

Verifiable, Debuggable, and Repairable Commonsense Logical Reasoning via LLM-based Theory Resolution

Anonymous ACL submission

Abstract

Recent advances in Large Language Models (LLM) have led to substantial interest in their application to commonsense reasoning tasks. Despite their potential, LLMs are susceptible to reasoning errors and hallucinations that may be harmful in use cases where accurate reasoning is critical. This challenge underscores the need for verifiable, debuggable, and repairable LLM reasoning. Recent works have made progress toward verifiable reasoning with LLMs by using them as either (i) a reasoner over an axiomatic knowledge base, or (ii) a semantic parser for use in existing logical inference systems. However, both settings are unable to extract commonsense axioms from the LLM that are not already formalized in the knowledge base, and also lack a reliable method to repair missed commonsense inferences. In this work, we present LLM-TRes, a logical reasoning framework based on the notion of “theory resolution” that allows for seamless integration of the commonsense knowledge from LLMs with a verifiable logical reasoning framework that mitigates hallucinations and facilitates debugging of the reasoning procedure as well as repair. We crucially prove that repaired axioms are theoretically guaranteed to be given precedence over flawed ones in our theory resolution inference process. We conclude by evaluating on three diverse language-based reasoning tasks – preference reasoning, deductive reasoning, and causal commonsense reasoning – and demonstrate the superior performance of LLM-TRes vs. state-of-the-art LLM-based reasoning methods in terms of both accuracy and reasoning correctness.

1 Introduction

The rise of Large Language Models (LLMs) has marked a pivotal moment in the real-world deployment of AI, particularly due to the exceptional ability of LLMs to handle complex reasoning tasks (Chang et al., 2024; Huang and Chang,

2023). Research has shown that LLMs have acquired significant commonsense knowledge (Zhao et al., 2024; Bian et al., 2023), which is crucial for engaging with real-world users in tasks such as question answering (Singhal et al., 2023) and recommendation (Sanner et al., 2023). Unfortunately, LLMs are prone to a variety of reasoning errors; for example, they commonly incorporate superficially plausible but factually incorrect information into their reasoning in a phenomenon known as *hallucination* (Zhang et al., 2023b; Ji et al., 2023). Furthermore, since the underlying reasoning process of the LLM is latent and hence largely opaque, validating reasoning soundness and identifying errors remains an open research problem. Such issues present a significant challenge to the reliability of using LLMs as reasoning systems, which impedes their practical utility (Mallen et al., 2023).

In light of these obstacles, recent research has proposed methodologies for extracting verifiable reasoning from LLMs by leveraging formal reasoning procedures. Such works fall under two main categories: (i) Using the LLM as a reasoner across an axiomatic knowledge base, while organizing the reasoning process into simpler subgoals to facilitate soundness of the overall reasoning (Kazemi et al., 2023). (ii) Leveraging the LLM as a semantic parser that translates natural language statements into logical axioms, followed by the use of an off-the-shelf theorem prover to perform logical reasoning (Pan et al., 2023; Olausson et al., 2023). While these seminal works have made progress towards verifiable LLM reasoning, their application in real-world tasks requiring commonsense reasoning is limited since they all suffer from the inability to extract verifiable commonsense axioms from the LLM that are not already formalized in the provided knowledge base axioms. Hence, these existing methodologies critically lack the ability to leverage the LLM as a verifiable commonsense reasoner to fill-in inevitable knowledge base gaps.

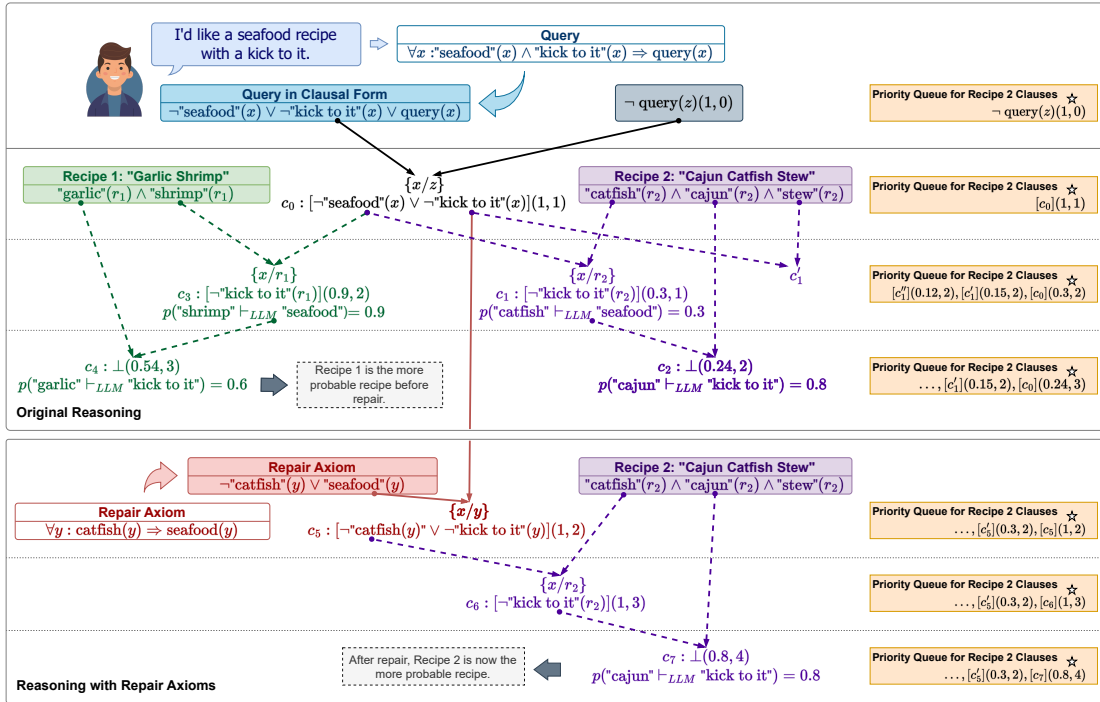


Figure 1: Preference reasoning is used as an illustrative example to show the LLM-TRes workflow. Top: LLM-based theory resolution is performed to calculate proof scores of two candidate Recipes entailing the Query. The proof begins from the negated query, and for each resolvent clause, a priority score tuple: (*proof plausibility score*, *proof length*) is calculated and pushed to a priority queue (only Recipe 2 clauses are shown here). At each step, the clause with the highest priority in the queue becomes the active clause. Here, due to a flawed low probability assigned to “catfish” entailing “seafood”, the *proof score* of Recipe 2 is mistakenly calculated lower than it should be. Bottom: After insertion of the Repair Axiom, the erroneous reasoning is repaired, leading to a higher score for the correct Recipe 2.

Furthermore, these methods lack any mechanism for repairing reasoning mistakes after detection.

To address these challenges, we propose LLM-TRes, a formal reasoning framework using LLMs. LLM-TRes satisfies three key desiderata that we illustrate through the worked example in Figure 1: (i) *Verifiability*: allowing for verification of every step in the reasoning process (i.e., from each successful refutation resolution \perp , one can backtrack the entire proof of the Query). (ii) *Debuggability*: being able to identify the incorrect inferences that led to a reasoning mistake (i.e., we observe an incorrect low probability LLM inference that “catfish” entails “seafood” for Recipe 2). (iii) *Repairability*: enabling a deterministic and reliable mechanism for rectifying the identified errors to produce correct inference (e.g., once we add the the explicit Repair Axiom $\forall y : \text{catfish}(y) \Rightarrow \text{seafood}(y)$, we arrive at a much higher proof plausibility for the correct preference match of the query to Recipe 2).

Formally, LLM-TRes is based on the concept of *theory resolution* (Stickel, 1985; Baumgartner, 1992), drawn from classical logical reasoning that

enables the integration of specialized reasoners into the resolution theorem-proving inference rule. Leveraging theory resolution, LLM-TRes seamlessly incorporates LLMs as specialized reasoners equipped with commonsense knowledge into verifiable logical reasoning. This integration enables extraction of relevant commonsense axioms from the LLM that cannot otherwise be obtained from the knowledge base. Finally, by capitalizing on a specially defined selection rule in our resolution framework, we formally prove that repairing flawed reasoning by the LLM is possible by providing correct axioms that are theoretically guaranteed to override the LLM’s flawed reasoning.

In summary, we contribute the following:

- We propose LLM-TRes, a formal reasoning framework founded on theory resolution that allows for incorporating the internal knowledge of the LLMs in a formal reasoning process to mitigate their hallucinations.
- We demonstrate that LLM-TRes provides a fully verifiable and debuggable reasoning

129	scheme by granting access to all reasoning	179
130	steps at an atomic level.	180
131	• We provide a mechanism for correcting errors	181
132	in the reasoning process with a theoretical	182
133	guarantee of prioritizing correct Repair Ax-	183
134	ioms over incorrect LLM inferences.	184
135	• We experiment with LLM-TRes on three distinct	185
136	tasks – preference reasoning, deductive	186
137	reasoning, and causal commonsense reason-	187
138	ing – demonstrating superior accuracy and	188
139	reasoning correctness compared to Chain of	189
140	Thought (CoT) (Wei et al., 2022) prompting	190
141	in LLMs (much larger in size) and LAM-	191
142	BADA (Kazemi et al., 2023), a state-of-the-art	192
143	formal reasoning framework.	193
144	2 Related Works	194
145	Reasoning with LLMs While their primary de-	195
146	sign was for text generation, LLMs exhibit re-	196
147	markable performance in many other NLP tasks	197
148	that require a variety of reasoning skills (Chang	198
149	et al., 2024; Xu et al., 2023). Despite such im-	199
150	pressive capabilities, errors and hallucinations that	200
151	can commonly occur in LLM reasoning have mo-	201
152	tivated research on obtaining dependable reason-	202
153	ing from LLMs while leveraging their intrinsic	203
154	knowledge. In this regard, several approaches	204
155	have been proposed to elicit stronger reasoning per-	205
156	formance from LLMs such as Chain-of-Thought	206
157	prompting (Wei et al., 2022; Kojima et al., 2022),	207
158	Self-Consistency (Wang et al., 2022), Least-to-	208
159	Most prompting (Zhou et al., 2022), and Selection-	209
160	Inference (Creswell et al., 2022). Despite being	210
161	effective in improving reasoning performance, all	211
162	these methods follow an <i>informal</i> reasoning proce-	212
163	dure in which the LLM is in charge of performing	213
164	reasoning and thus does not guarantee the faithful-	214
165	ness of the reasoning process (Shanahan, 2024; Pan	215
166	et al., 2023). For instance, the reasoning ability of	216
167	these methods may be unreliable for tasks requir-	217
168	ing out-of-domain reasoning (Saparov et al., 2024;	218
169	Liang et al., 2022), tasks involving negation (Anil	219
170	et al., 2022), and often degrades with an increase	220
171	in the length of reasoning steps (Dziri et al., 2024).	221
172	Formal Reasoning with LLMs To obtain reli-	222
173	able and verifiable reasoning from LLMs, a number	223
174	of works have proposed the idea of using LLMs	224
175	in a <i>formal</i> reasoning framework — a systematic	225
176	and logical process governed by a set of rules and	226
177	principles (Galotti, 1989). Two main approaches	227
178	have been proposed in this regard. In the first ap-	
	proach, the LLM is utilized to perform different	179
	sub-tasks of a formal logical inference rule to rea-	180
	son over an axiomatic knowledge base. For ex-	181
	ample, LAMBADA (Kazemi et al., 2023) uses the	182
	LLM to perform goal decomposition, rule selection,	183
	and fact-checking in a backward chaining process.	184
	In a related vein, Symba (Lee and Hwang, 2024)	185
	introduces a top-down solver to control the proof	186
	process and uses the LLM as an aide to the solver.	187
	In the second approach, LLMs are used as a se-	188
	mantic parser to translate natural language axioms	189
	and facts to a specific logical format; here the re-	190
	sponsibility of inference is delegated to a symbolic	191
	theorem prover. LogicLM (Pan et al., 2023) uses	192
	this idea with a self-refinement mechanism to al-	193
	low the LLM to refine its symbolic conversions.	194
	Since LLMs commonly make syntactic and seman-	195
	tic errors in the parsing process, LINC (Olausson	196
	et al., 2023) performs majority voting over multiple	197
	solutions to obtain the final result.	198
	These works have made significant progress in	199
	increasing the reliability and verifiability of LLM	200
	reasoning. However, they only utilize axioms that	201
	are explicitly provided in the knowledge base, and	202
	lack the ability to leverage the intrinsic common-	203
	sense knowledge of the LLM by extracting com-	204
	monsense axioms. This prevents existing methods	205
	from incorporating verifiable LLM-derived com-	206
	monsense knowledge in their reasoning, which is	207
	often critical in practical deployed usage. More-	208
	over, these existing methods do not support a reli-	209
	able mechanism for rectifying incorrect reasoning	210
	steps. <i>We aim to address all of these limitations</i>	211
	<i>with our contribution of the LLM-TRes framework.</i>	212
	3 Methodology: LLM-based Theory	213
	Resolution (LLM-TRes)	214
	We first introduce the resolution rule and the con-	215
	cept of theory resolution and then explain our LLM-	216
	based Theory Resolution (LLM-TRes) methodol-	217
	ogy. For the logical knowledge representation in	218
	this work, we assume a function- and equality-free	219
	first-order logical (FOL) syntax (Chang and Lee,	220
	2014) with all FOL sentences translated to clausal	221
	normal form as demonstrated in Figure 1.	222
	Resolution Rule Resolution is a sound inference	223
	rule that performs inference by deriving a resol-	224
	vent clause from two premise clauses containing	225
	complementary literals. Given two FOL sentences	226
	in clausal form, a new clause can be derived via	227

228 resolution of their complementary literals, e.g.,

$$229 \quad \frac{A(x) \vee B(x) \quad \neg B(y) \vee C(y)}{A(x) \vee C(x)}, \quad (1)$$

230 under the unification $\theta = \{x/y\}$. Following this
 231 procedure, new clauses are derived until either
 232 a contradiction \perp is found (e.g., deriving both
 233 clauses $A(x)$ and $\neg A(x)$ that resolve to \perp), or no
 234 further resolutions are possible. Finding a contra-
 235 diction implies that the original set of clauses is
 236 inconsistent. Therefore, given the knowledge base
 237 \mathcal{K} and a query q , to prove that $\mathcal{K} \vdash q$, one can apply
 238 the resolution inference rule to show that $\mathcal{K} \wedge \neg q$
 239 leads to a contradiction \perp .

240 **Theory Resolution** Theory resolution (Stickel,
 241 1985) is a methodology that enables the integra-
 242 tion of special purpose reasoning theories into res-
 243 olution theorem proving. Based on theory reso-
 244 lution, given two clauses $c_1 = A(x) \vee B(x)$ and
 245 $c_2 = C(x) \vee D(x)$, if a theorem prover T identifies
 246 $B(x)$ and $\neg C(y)$ under unification $\theta = \{x/y\}$ to
 247 be unsatisfiable (i.e., $\forall x B(x) \wedge \neg C(x) \vdash_T \perp$), de-
 248 spite lacking complementary literals with identical
 249 predicates, the two clauses can still be resolved:

$$250 \quad \frac{A(x) \vee B(x) \quad \neg C(y) \vee D(y)}{A(x) \vee D(x)}. \quad (2)$$

251 Theory resolution considerably broadens the appli-
 252 cability of the resolution inference rule by lifting
 253 the condition of resolving only complementary lit-
 254 erals. In this work, we use an LLM as the theory
 255 that identifies the unsatisfiable natural language
 256 predicates to do reasoning via theory resolution.

257 **Natural Language Logic** Natural language en-
 258 compasses a significant amount of information that
 259 cannot be easily represented using symbolic logic.
 260 Although one can represent functions and predi-
 261 cates in symbolic logic, it may be hard to fully
 262 axiomatize their real-world meaning, which is a
 263 substantial limitation of the semantic parsing ap-
 264 proaches. For instance, being “spicy” and having
 265 “a kick to it” are assigned completely different predi-
 266 cates, and pure symbolic reasoning cannot identify
 267 the intuitive entailment relationship between them
 268 without specific axioms. Moreover, representing
 269 commonsense knowledge in symbolic logic is very
 270 challenging (Davis, 2014). However, LLMs are
 271 capable of understanding the semantic relationship
 272 between such predicates and also encompass sub-
 273 stantial commonsense knowledge, which can be
 274 used for reasoning in real-world applications.

275 As mentioned earlier, theory resolution offers the
 276 capability to resolve non-complementary literals if
 277 they are deemed unsatisfiable by a theorem prover.
 278 By employing an LLM as the theorem prover, we
 279 can leverage the theory resolution framework to
 280 conduct resolution within an extended version of
 281 First-Order Logic, where predicates and functions
 282 are no longer symbols but natural language texts, a
 283 system we call *Natural Language (NL) Logic*.

284 Using the LLM theorem prover in the NL
 285 logic, the unsatisfiability condition of the the-
 286 ory resolution reduces to natural language en-
 287 tailment. In other words, if an LLM identifies
 288 a natural language predicate B to entail predi-
 289 cate D , i.e., $B(x) \vdash_{LLM} D(x)$, and therefore,
 290 $B(x) \wedge \neg D(x) \vdash_{LLM} \perp$, then literals $B(x)$ and
 291 $D(x)$ can be resolved. For instance, given clauses
 292 $c_1 = \text{“kick to it”}(x)$ and $c_2 = \neg \text{“spicy”}(x) \vee Q(x)$,
 293 in which $Q(x)$ is another literal with a natural lan-
 294 guage predicate, since the LLM identifies the natu-
 295 ral language entailment “kick to it” \vdash_{LLM} “spicy”,
 296 a theory resolution step can be performed as

$$297 \quad \frac{\text{“kick to it”}(x) \quad \neg \text{“spicy”}(x) \vee Q(x)}{Q(x)}. \quad (3)$$

3.1 LLM-TRes Algorithm 298

299 This section presents LLM-TRes, an algorithm for
 300 efficient logical commonsense reasoning based on
 301 theory resolution using LLMs. The workflow of
 302 LLM-TRes is shown through a worked example in
 303 Figure 1, and formalized in Algorithm 1.

304 **Problem Definition** Consider a set of queries Q
 305 and a knowledge base (KB) denoted by \mathcal{K} which
 306 comprises a set of axioms \mathcal{A} and a set of facts \mathcal{F} ,
 307 all represented in natural language logic in clausal
 308 form. In this work, we aim to propose an inference
 309 rule i that for each $q \in Q$, finds a set of proofs
 310 denoted by *proofs*, such that each proof $f \in \textit{proofs}$
 311 consists of a subset of clauses in \mathcal{K} , and is assigned
 312 a priority score ρ reflecting the priority of the proof.

313 **Algorithm** To prove that \mathcal{K} entails the query q via
 314 resolution, we must demonstrate that iteratively ap-
 315 plying resolution to derive new clauses from $\mathcal{K} \wedge \neg q$
 316 leads to a contradiction, thereby proving its unsat-
 317 isfiability. The first question that arises is which
 318 clause should be chosen to begin the resolution
 319 proof. Two major paradigms are used in perform-
 320 ing resolution: (i) starting from the clauses in \mathcal{K}
 321 to derive q from them and resolve it with $\neg q$, an

approach known as *forward chaining*, or (ii) starting from $\neg q$ and resolving it with clauses in \mathcal{K} to reach a contradiction, known as *backward chaining*. Since backward chaining employs a goal-driven approach, which is shown to improve efficiency in reasoning over natural language (Kazemi et al., 2023), we begin the resolution process from $\neg q$. Therefore, $\neg q$ becomes our first *active clause* that we need to resolve with a clause from \mathcal{K} .

The potentially enormous size of \mathcal{K} poses a major challenge. Also, as the resolution process progresses, new clauses are created, leading to a further expansion in the size of the search space. To perform resolution efficiently in this combinatorial search space, LLM-TRes employs two strategies: (i) prioritizing the resolvent clauses to continue the resolution process, and (ii) restricting the theory resolution search space using semantic similarity.

Resolution Priority Definition and Ordering: The first mechanism employed in LLM-TRes to enable efficient resolution is prioritizing candidate clauses. Using this prioritization scheme, resolvent clauses that have a higher potential for being part of a plausible proof will be given precedence over clauses generated from less plausible resolution steps. The plausibility of a theory resolution step, in which an active clause c is resolved with a clause c_{target} to generate the resolvent clause c_{res} , denoted by $\rho_{c_{res}}^{entail}$, is determined by calculating the probability that the LLM assigns to c_{target} entailing c .

$$\rho_{c_{res}}^{entail} = p(c_{target} \vdash_{LLM} c). \quad (4)$$

These plausibility scores can help us prioritize the resolvent clauses. For instance, in the example provided in Figure 1, since resolving “*shrimp*” with “*seafood*” yields a higher entailment score than resolving “*garlic*” with “*seafood*”, it is intuitive to prioritize the first resolvent as it is more likely to be part of the final proof. Since we are interested in identifying the most plausible proofs, i.e., the sequences of theory resolution steps with the highest entailment scores, we define the first entry of our priority score for each resolvent clause c_{res} as the overall entailment score of all resolution steps beginning from the original negated query that led to its derivation. Denoting the set of parent clauses of c_{res} as $\mathcal{P}_{c_{res}}$, we can inductively define the overall proof entailment score of c_{res} as

$$\rho^e(c_{res}) = \left(\prod_{c' \in \mathcal{P}_{c_{res}}} \rho^e(c') \right) \cdot \rho_{c_{res}}^{entail}. \quad (5)$$

Algorithm 1 LLM-TRes Algorithm

```

1: Input:  $\mathcal{K}, q, max\_proofs, max\_iters, b$ 
2:  $proofs \leftarrow \emptyset$ 
3:  $PQ \leftarrow \emptyset$  //  $PQ$  is an initially empty priority queue.
4:  $PQ.push(\neg q, (1, 0))$  // Negation of the initial query  $q$  has priority  $(1, 0)$ ,  $PQ$  is ordered by Equation 7
5: while  $i < max\_iters$  do
6:   while  $PQ \neq \emptyset \wedge i < max\_proofs$  do
7:      $c \leftarrow PQ.pop()$ 
8:     if  $c = \perp$  then
9:        $max\_proofs++$ 
10:       $proofs \leftarrow proofs \cup \{c\}$ 
11:     else
12:        $\beta_c \leftarrow b$  most likely candidates in  $\mathcal{K}$  to resolve with  $c$ 
13:       for  $c_{target} \in \beta_c$  do
14:         Compute resolvent  $c_{res}$  of  $c$  and  $c_{target}$  using Equation 2
15:          $PQ.push(c_{res}, (\rho^e(c_{res}), \rho^l(c_{res})))$  // cf. Equations 5 and 6
16: Output:  $proofs$ 

```

When choosing between equally plausible proofs, we are interested in shorter proofs that avoid redundant steps. We assign a second priority score to reflect this preference which is considered only when the entailment proof scores are equal. As for the proof entailment score, we can obtain the proof length score of c_{res} inductively from the maximum proof length of its parent clauses as

$$\rho^l(c_{res}) = 1 + \max_{c' \in \mathcal{P}_{c_{res}}} \rho^l(c'). \quad (6)$$

The final priority score for each resolvent clause c_{res} is formed as the tuple $(\rho^e(c_{res}), \rho^l(c_{res}))$ and all resolvents are pushed to a priority queue PQ . A total order of clauses in PQ is then determined as

$$c_1 \preceq c_2 \iff [\rho^e(c_1) > \rho^e(c_2)] \vee [(\rho^e(c_1) = \rho^e(c_2)) \wedge (\rho^l(c_1) < \rho^l(c_2))]. \quad (7)$$

Restricting Theory Resolution with Embeddings: The knowledge base may contain various axioms and facts, many of which are irrelevant to the active clause. To enhance efficiency and maintain the growth of the resolution space tractable, we restrict the size of our resolution search space by a branching factor b and select candidate target clauses for performing resolution based on their semantic relevance to the current active clause. Concretely, we

use the similarity scores between \mathbf{z}_c , the word embedding vector of the active clause c , and $\mathbf{z}_{c'}$, word embedding vectors of each candidate clause c' to find β_c , the set of b most relevant clauses to c as

$$\beta_c = \{c' | (c' \in \mathcal{K}) \wedge (c' \neq c) \wedge (\mathbf{z}_c^T \mathbf{z}_{c'} \geq \tau)\}, \quad (8)$$

in which τ is set to the b^{th} highest inner product score between embeddings of c and other clauses, thus resulting in top- b theory resolution candidates. Next, theory resolution can be performed between c and each clause in β_c as in Equation 2.

These two mechanisms together enable an efficient inference via LLM-based theory resolution. At the beginning of each iteration of LLM-TRes, the clause holding the foremost position in PQ becomes the active clause. Once a resolution step leads to contradiction, the proof and its respective priority score are added to the set of found proofs by backtracing the ancestor clauses up to the query.

This algorithm continues until either a certain number of proofs are found or the maximum number of iterations is exceeded. Notably, LLM-TRes is not limited to proving a single query; instead, it finds a set of proofs with each assigned a strength score. This functionality allows it to assess the likelihood of each query being entailed, which is beneficial for applications requiring ranking, such as answering multiple-choice questions. In applications where a binary truth value is considered for the query, the proof scores of q and $\neg q$ are compared. Our experiments cover both cases.

4 Repairability of Erroneous Resolution

Since LLM-TRes provides access to atomic inference steps in the resolution process, it facilitates verifiability and debuggability. Although the entailment probabilities assigned by the LLMs may be erroneous, the exact resolution step at which the failure occurs is discernible. Furthermore, it can be easily corrected by introducing a rectifying rule into the knowledge base. 1 presents an example. The LLM’s mistake in assigning a low entailment score for “*catfish*” to entail “*seafood*” leads to incorrect reasoning. However, introducing the correct axiom $\forall y \text{“catfish”}(y) \implies \text{“seafood”}(y)$ to the KB repairs this mistake. The following proposition formalizes this property and is proven in Appendix A.

Proposition 1. Consider proof P_c^ϕ using axiom ϕ that derives clause c . For any incorrect LLM rea-

soning axiom ϕ , a Repair Axiom ϕ' can be inserted such that $P_c^{\phi'}$ will be produced before P_c^ϕ .

5 Experiments

We evaluate LLM-TRes on three different tasks involving commonsense reasoning on natural language data: preference reasoning, deductive reasoning, and causal commonsense reasoning. We release our implementation and data¹.

5.1 Tasks and Datasets

Preference Reasoning Understanding user preferences from natural language statements is a complex but essential task in applications requiring personalization. For this task, we use RecipeMPR (Zhang et al., 2023a), a dataset consisting of 500 user queries stating their preference toward recipes, e.g., “*I would like meat lasagna but I’m watching my weight*” with five-way recipe options that requires a range of commonsense inference including temporal and analogical reasoning.

Deductive Reasoning We use ProntoQA (Saparov and He, 2022), a widely used dataset for evaluating the deductive reasoning ability of LLMs. This dataset consists of natural language queries about KBs including facts and axioms generated from ontologies. We use 500 queries of the true ontology as they are consistent with the real world and are useful to evaluate commonsense reasoning. We select the most challenging 5-hop subset of the dataset.

Causal Commonsense Reasoning We use COPASSE (Brassard et al., 2022), a dataset for reasoning about event causes and effects using a semi-structured KB. In the “*effect*” split of this dataset, an event is provided such as “*The pen ran out of ink.*”, together with semi-structured explanations with assigned quality scores, and the task is to determine the more plausible candidate effect, e.g., “*I used a pencil.*” or “*I signed my name.*”.

5.2 Baselines and Evaluation

We use LAMBADA², a seminal work in formal reasoning with LLMs, and zero-shot Chain-of-Thought (CoT) prompting (Kojima et al., 2022) as our comparison baselines. Semantic parsing methods are inherently unable to perform commonsense reasoning and do not apply to our tasks. We use

¹<https://anonymous.4open.science/r/LLM-TRes-678C/>

²Since the original paper did not release code, we use the implementation in (Lee and Hwang, 2024).

Table 1: Reasoning performance of methods across the three datasets. Gemma fails to provide explanations for Recipe-MPR, so reasoning scores cannot be calculated for it (Fail). LAMBADA requires a rule set that is not provided in Recipe-MPR, and cannot rank proofs which is necessary for COPA-SSE. Pure entailment does not generate proofs, so the reasoning scores do not apply to it (NA).

Method	Recipe-MPR			ProntoQA			COPA-SSE		
	Accuracy	RS Macro	RS Micro	Accuracy	RS Macro	RS Micro	Accuracy	RS Macro	RS Micro
CoT (GPT-3.5-Turbo)	0.844	0.850	0.900	0.738	0.600	0.878	0.860	0.800	0.921
CoT (Llama3 8B)	0.768	0.550	0.742	0.718	0.250	0.802	0.839	0.800	0.916
CoT (Gemma 7B)	0.460	Fail	Fail	0.588	0.250	0.657	0.818	0.520	0.450
CoT (Mistral 7B)	0.689	0.850	0.946	0.902	0.600	0.890	0.643	0.850	0.892
LAMBADA (GPT-3.5-Turbo)	NA	NA	NA	0.580	0.800	0.900	NA	NA	NA
Pure Entailment (BART 406M)	0.682	NA	NA	0.740	NA	NA	0.825	NA	NA
LLM-TRes (BART 406M)	0.822	1.000	1.000	0.990	1.000	1.000	0.888	0.900	0.958

GPT-3.5 Turbo as the LLM for LAMBADA and for converting the natural language axioms and queries to the clausal form in our method, and obtain the entailment probabilities for theory resolution using BART large (Lewis et al., 2020) model³ trained on MNLI (Williams et al., 2018) dataset. We compare against CoT prompting with GPT-3.5 Turbo, Llama3 8B, Mistral 7B, and Gemma 7B. To ensure that the difference in the performance of our model and the baselines is not due to using different LLMs, we also use pure BART-large entailment scores between facts and query as a baseline.

We evaluate the reasoning performance of the models considering the correctness of the final answers, measured by the accuracy, and correctness of the reasoning process measured by the reasoning score (RS) which we manually assess for the first 20 queries the models answer correctly. RS is commonly evaluated as a binary judgment on whether the predicted proof is supported by the ground truth proof (Kazemi et al., 2023; Lee and Hwang, 2024). However, RS does not assess the number of errors. Therefore, in addition to this metric which we call *macro RS*, to provide a more fine-grained evaluation of the provided proofs, based on the idea provided in Min et al. (2023), we use a new metric which we name *micro RS*. Given a provided proof P and the ground truth proof P^* , and denoting the indicator function as \mathbb{I} , the micro RS for each query is defined as $RS_{Micro} = \frac{1}{|P|} \sum_{p \in P} \mathbb{I}(p \in P^*)$.

5.3 Results

RQ1: Comparison of Reasoning Performance

Results of the reasoning performance are provided in Table 1. On deductive and causal commonsense reasoning tasks, LLM-TRes achieves higher accuracies than the baselines although the language

models they use are multiple times larger. On preference reasoning, LLM-TRes achieves the second-highest accuracy after CoT prompting with GPT3.5 Turbo with a rather small margin. On ProntoQA, since the high-quality conversion of the query and the knowledge base to the clausal format is straightforward, LLM-TRes can prioritize complementary literals to perform exact resolution, resulting in a near-ideal performance. The failure cases of LLM-TRes are due to the LLM’s limitation in understanding contraposition as noted in previous work (Zhang et al., 2024). Nonetheless, LLM-TRes maintains consistently high performance, unlike other baselines which vary across tasks. For instance, while CoT with GPT-3.5 excels on Recipe-MPR and COPA-SSE, it is outperformed by Mistral on ProntoQA, which in turn performs rather poorly on Recipe-MPR and COPA-SSE. On RS, LLM-TRes outperforms all baselines across the three datasets at both the macro and micro level, showcasing its capability to provide proofs that are consistent with the ground truth proof. LAMBADA is unable to reason on Recipe-MPR as it performs backward chaining on explicit rule sets, which Recipe-MPR does not provide. Also, since LAMBADA can only prove or refute a query based on a KB and cannot score and rank the plausibility of proofs, it cannot choose the more plausible effects on the COPA-SSE dataset. Since CoT using Gemma refrained from providing any proof for preference reasoning despite being prompted to do so, the reasoning score could not be calculated for it. Finally, pure entailment does not provide proofs so RS cannot be evaluated.

RQ2: Robustness to Incompleteness of the KB

Assuming access to a complete KB in which all required axioms are provided is often unrealistic

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

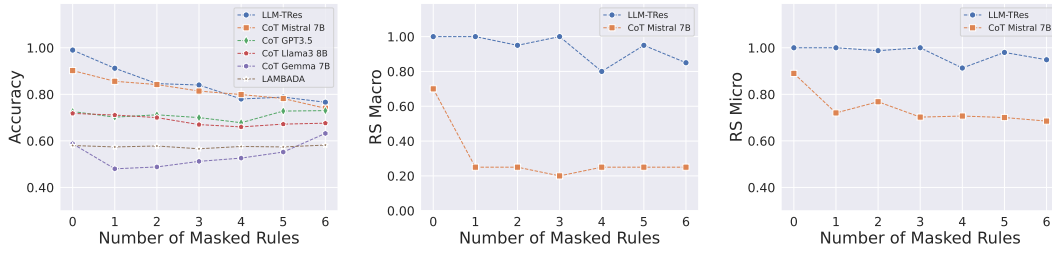


Figure 2: Reasoning performance of different models on ProntoQA with an incomplete KB. We mask out a number of rules to vary the degree of incompleteness of KB.

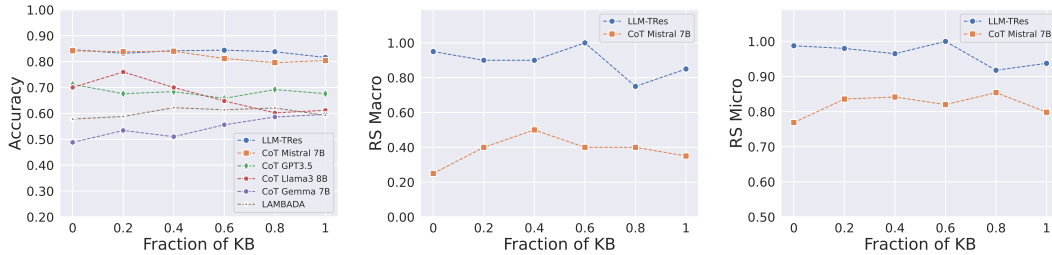


Figure 3: Reasoning performance of different models on ProntoQA with larger KB. We sample random axioms from other queries to increase the size of KB.

in practical applications. Therefore, a commonsense reasoning methodology must be able to extract the intrinsic commonsense knowledge of the LLMs to overcome the incompleteness of the KB. We assess this capability by repeatedly running experiments on ProntoQA each time removing a number of randomly chosen rules from the KB. We chose ProntoQA for this study as it is the only dataset with large rule sets that enables experiments with various ablated rules. Results of this experiment are provided in Figure 2. Although ablating rules from the KB decreases the accuracies of both LLM-TRes and the best baseline, CoT with Mistral, LLM-TRes often maintains higher accuracy. Moreover, the consistently higher reasoning score of LLM-TRes proves its superior ability to generate valid proofs.

RQ3: Robustness to Increase in Size of the KB

In this experiment, we evaluate the robustness of LLM-TRes and other baselines to increases in the KB size. We form a large KB consisting of 75 distinct rules across different ProntoQA queries and each time add a fraction of this KB to the original rule set of the query while randomly masking 2 rules of the original KB. This experiment mainly aims to determine if the restricting resolution search space of LLM-TRes using semantic similarity can identify the relevant clauses to the active clause. In all tests, LLM-TRes uses the similarity between GPT-3 embeddings of the clauses with a

branching factor of 15. Meanwhile, other baselines include the entire KB in the prompt which is costly and inefficient. Results of this test, shown in Figure 3 depict that LLM-TRes and the best baseline, CoT with Mistral, sustain their performance, but LLM-TRes consistently obtains higher reasoning scores while using a more efficient methodology for pruning the reasoning search space.

6 Conclusion

We presented LLM-TRes, a novel framework for formal reasoning with LLMs based on theory resolution that allowed us to integrate LLMs into resolution logical reasoning seamlessly. By providing access to every atomic reasoning step, LLM-TRes enabled *verifiability* and *debuggability* of the process. It also offered a reliable *repairing* mechanism for correcting flaws in the LLM reasoning by asserting the particular missed axiom which was theoretically guaranteed to override the mistakenly low-probability resolution step. The promising performance of LLM-TRes on preference reasoning, deductive reasoning, and causal commonsense reasoning tasks demonstrates its efficacy in providing accurate answers and correct proofs. These capabilities make LLM-TRes a robust foundation for counteracting hallucination and pave the way for more trustworthy deployment of LLM-based commonsense reasoners in applications where correctness, verifiability, and reparability are paramount.

618 Limitations

619 While we believe this work has made substantial
620 progress in verifiable, debuggable, and repairable
621 commonsense reasoning, it naturally has limita-
622 tions that we hope will encourage further inves-
623 tigation and future work. As mentioned in the
624 paper, we provided a reliable mechanism for er-
625 roneous reasoning processes; however, determin-
626 ing a flawed step requires expert judgment. In our
627 work, we do not focus on evaluating the reasoning
628 steps and how the repair axioms are introduced.
629 Proposing an automated mechanism for evaluat-
630 ing the reasoning steps can be a direction of future
631 research. Furthermore, as in all LLM-based reason-
632 ing methodologies, obtaining high reasoning per-
633 formances requires an apt LLM. As we discussed
634 in Section 5.3, limitations of the utilized LLM such
635 as their shortcomings in understanding contraposi-
636 tion can pose challenges to the overall performance
637 of the method. Finally, as we mentioned in the
638 paper, LLM-TRes focuses on the natural language
639 extension of First Order Logic (FOL) which we
640 introduced, and extending it to Higher-Order Logic
641 (HOL) could be considered as a future research
642 direction given the prior uses of HOL in formaliz-
643 ing natural language semantics and complex modal
644 constructs (van Eijck and Unger, 2010).

645 Ethics Statement

646 Our contribution of LLM-TRes aims to enable
647 transparent reasoning with LLMs such that cor-
648 rectness of every reasoning step can be verified and
649 potentially repaired if incorrect. However, it is im-
650 portant for us to note that a correct proof or line of
651 argument from premises neither presupposes that
652 the premises are ethical nor that the conclusion de-
653 rived from the premises and line of reasoning is
654 ethical. In this sense, practical use of LLM-TRes
655 still requires ethical oversight to monitor ethical
656 and bias considerations for any axioms entered by
657 the user as well as to verify that unintended reason-
658 ing hallucinations by the underlying LLM have not
659 led to incorrect, biased, or unethical conclusions.

660 References

661 Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor
662 Lewkowycz, Vedant Misra, Vinay Ramasesh, Am-
663 brose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam
664 Neyshabur. 2022. Exploring length generalization in
665 large language models. *Advances in Neural Informa-*
666 *tion Processing Systems*, 35:38546–38556.

- Peter Baumgartner. 1992. An ordered theory resolu- 667
tion calculus. In *Logic Programming and Automated* 668
Reasoning: International Conference LPAR’92 St. 669
Petersburg, Russia, July 15–20, 1992 Proceedings 3, 670
pages 119–130. Springer. 671
- Ning Bian, Xianpei Han, Le Sun, Hongyu Lin, Yaojie 672
Lu, Ben He, Shanshan Jiang, and Bin Dong. 2023. 673
Chatgpt is a knowledgeable but inexperienced solver: 674
An investigation of commonsense problem in large 675
language models. *arXiv preprint arXiv:2303.16421.* 676
- Ana Brassard, Benjamin Heinzerling, Pride Kavumba, 677
and Kentaro Inui. 2022. Copa-sse: Semi-structured 678
explanations for commonsense reasoning. In *Pro-* 679
ceedings of the Thirteenth Language Resources and 680
Evaluation Conference, pages 3994–4000. 681
- Chin-Liang Chang and Richard Char-Tung Lee. 2014. 682
Symbolic logic and mechanical theorem proving. 683
Academic press. 684
- Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu, 685
Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi, 686
Cunxiang Wang, Yidong Wang, et al. 2024. A sur- 687
vey on evaluation of large language models. *ACM* 688
Transactions on Intelligent Systems and Technology, 689
15(3):1–45. 690
- Antonia Creswell, Murray Shanahan, and Irina Higgins. 691
2022. Selection-inference: Exploiting large language 692
models for interpretable logical reasoning. *arXiv* 693
preprint arXiv:2205.09712. 694
- Ernest Davis. 2014. *Representations of commonsense* 695
knowledge. Morgan Kaufmann. 696
- Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang Lor- 697
raine Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, 698
Peter West, Chandra Bhagavatula, Ronan Le Bras, 699
et al. 2024. Faith and fate: Limits of transformers on 700
compositionality. *Advances in Neural Information* 701
Processing Systems, 36. 702
- Kathleen M Galotti. 1989. Approaches to studying formal 703
and everyday reasoning. *Psychological bulletin,* 704
105(3):331. 705
- Jie Huang and Kevin Chen-Chuan Chang. 2023. To- 706
wards reasoning in large language models: A survey. 707
In *Findings of the Association for Computational* 708
Linguistics: ACL 2023, pages 1049–1065. 709
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan 710
Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea 711
Madotto, and Pascale Fung. 2023. Survey of halluci- 712
nation in natural language generation. *ACM Comput-* 713
ing Surveys, 55(12):1–38. 714
- Mehran Kazemi, Najoung Kim, Deepti Bhatia, Xin Xu, 715
and Deepak Ramachandran. 2023. Lambada: Back- 716
ward chaining for automated reasoning in natural 717
language. In *Proceedings of the 61st Annual Meet-* 718
ing of the Association for Computational Linguistics 719
(Volume 1: Long Papers), pages 6547–6568. 720

721	Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. <i>Advances in neural information processing systems</i> , 35:22199–22213.	779
722		780
723		781
724		782
725		783
726	Jinu Lee and Wonseok Hwang. 2024. Symba: Symbolic backward chaining for multi-step natural language reasoning. <i>arXiv preprint arXiv:2402.12806</i> .	784
727		785
728		786
729	Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7871–7880.	787
730		788
731		789
732		790
733		791
734		792
735		793
736		794
737	Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. Holistic evaluation of language models. <i>arXiv preprint arXiv:2211.09110</i> .	795
738		796
739		797
740		798
741		799
742	Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. When not to trust language models: Investigating effectiveness of parametric and non-parametric memories. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 9802–9822.	800
743		801
744		802
745		803
746		804
747		805
748		806
749	Sewon Min, Kalpesh Krishna, Xinxu Lyu, Mike Lewis, Wen-tau Yih, Pang Koh, Mohit Iyyer, Luke Zettlemoyer, and Hannaneh Hajishirzi. 2023. Factscore: Fine-grained atomic evaluation of factual precision in long form text generation. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 12076–12100.	807
750		808
751		809
752		810
753		811
754		812
755		813
756	Theo Olausson, Alex Gu, Ben Lipkin, Cedegao Zhang, Armando Solar-Lezama, Joshua Tenenbaum, and Roger Levy. 2023. Linc: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. In <i>Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing</i> , pages 5153–5176.	814
757		815
758		816
759		817
760		818
761		819
762		820
763	Liangming Pan, Alon Albalak, Xinyi Wang, and William Wang. 2023. Logic-lm: Empowering large language models with symbolic solvers for faithful logical reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2023</i> , pages 3806–3824.	821
764		822
765		823
766		824
767		825
768		826
769	Scott Sanner, Krisztian Balog, Filip Radlinski, Ben Wedin, and Lucas Dixon. 2023. Large language models are competitive near cold-start recommenders for language-and item-based preferences. In <i>Proceedings of the 17th ACM conference on recommender systems</i> , pages 890–896.	827
770		828
771		829
772		830
773		831
774		832
775	Abulhair Saparov and He He. 2022. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. In <i>The Eleventh International Conference on Learning Representations</i> .	833
776		834
777		835
778		
	Abulhair Saparov, Richard Yuanzhe Pang, Vishakh Padmakumar, Nitish Joshi, Mehran Kazemi, Najoung Kim, and He He. 2024. Testing the general deductive reasoning capacity of large language models using ood examples. <i>Advances in Neural Information Processing Systems</i> , 36.	
	Murray Shanahan. 2024. Talking about large language models. <i>Communications of the ACM</i> , 67(2):68–79.	
	Karan Singhal, Tao Tu, Juraj Gottweis, Rory Sayres, Ellery Wulczyn, Le Hou, Kevin Clark, Stephen Pfohl, Heather Cole-Lewis, Darlene Neal, et al. 2023. Towards expert-level medical question answering with large language models. <i>arXiv preprint arXiv:2305.09617</i> .	
	Mark E Stickel. 1985. Automated deduction by theory resolution. <i>Journal of Automated Reasoning</i> , 1(4):333–355.	
	Jan van Eijck and Christina Unger. 2010. <i>Computational Semantics with Functional Programming</i> . Cambridge University Press.	
	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc V Le, Ed H Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. In <i>The Eleventh International Conference on Learning Representations</i> .	
	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. <i>Advances in neural information processing systems</i> , 35:24824–24837.	
	Adina Williams, Nikita Nangia, and Samuel Bowman. 2018. A broad-coverage challenge corpus for sentence understanding through inference. In <i>Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)</i> , pages 1112–1122.	
	Yudong Xu, Wenhao Li, Pashootan Vaezipoor, Scott Sanner, and Elias B Khalil. 2023. Llms and the abstraction and reasoning corpus: Successes, failures, and the importance of object-based representations. <i>arXiv preprint arXiv:2305.18354</i> .	
	Haochen Zhang, Anton Korikov, Parsa Farinneya, Mohammad Mahdi Abdollah Pour, Manasa Bharadwaj, Ali Pesaranghader, Xi Yu Huang, Yi Xin Lok, Zhaoqi Wang, Nathan Jones, et al. 2023a. Recipe-mpr: A test collection for evaluating multi-aspect preference-based natural language retrieval. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2744–2753.	
	Yanfang Zhang, Yiliu Sun, Yibing Zhan, Dapeng Tao, Dacheng Tao, and Chen Gong. 2024. Large language models as an indirect reasoner: Contrapositive and contradiction for automated reasoning. <i>arXiv preprint arXiv:2402.03667</i> .	

836 Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu,
837 Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang,
838 Yulong Chen, et al. 2023b. Siren’s song in the ai
839 ocean: a survey on hallucination in large language
840 models. *arXiv preprint arXiv:2309.01219*.

841 Zirui Zhao, Wee Sun Lee, and David Hsu. 2024. Large
842 language models as commonsense knowledge for
843 large-scale task planning. *Advances in Neural Infor-*
844 *mation Processing Systems*, 36.

845 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,
846 Nathan Scales, Xuezhi Wang, Dale Schuurmans,
847 Claire Cui, Olivier Bousquet, Quoc V Le, et al. 2022.
848 Least-to-most prompting enables complex reasoning
849 in large language models. In *The Eleventh Interna-*
850 *tional Conference on Learning Representations*.

A Proof of Repairability of LLM-TRes

Proposition 2. 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 . We can define the proof score $\rho^e(P_c^\phi)$ of clause c by inductively unrolling Equation 5 for $\rho^e(c)$ over all ancestor clauses P_c^ϕ that derive it. This yields a simple product form: $\rho^e(P_c^\phi) = \rho_\phi^{entail} \cdot \prod_{c' \in P_c} \rho_{c'}^{entail}$. 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 $\rho_{\phi'}^{entail} = 1$ (by definition) while the LLM entailment score $\rho_\phi^{entail} < 1$ (necessarily). 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 7 used to prioritize proofs in the LLM-TRes Algorithm 1. \square

B Anecdotal Examples

To offer deeper insight into the responses and proofs generated by LLM-TRes and the comparison baselines, this section presents anecdotal examples illustrating each model’s performance on the evaluated tasks. Specifically, we showcase the outputs from the following models:

- LLM-TRes
- Chain of Thought prompting with Mistral
- Chain of Thought prompting with Llama3
- Chain of Thought prompting with Gemma
- Chain of Thought prompting with GPT-3.5
- LAMBADA

We apply these models to three distinct tasks, offering a comparative analysis of their responses. Detailed anecdotal examples are provided to elucidate the models’ capabilities and limitations. In each anecdotal example, we present the query, the rule set (on datasets that contain it), 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.

An example of the Deductive reasoning task, as evaluated on the ProntoQA dataset, is detailed in Section B.1. Additionally, we explore the models’ performance on causal commonsense reasoning using the COPA-SSE dataset and on preference reasoning using the Recipe-MPR dataset, with examples discussed in Sections B.2 and B.3, respectively.

B.1 Deductive Reasoning

Query: is x imaginary?

Rule set: Real numbers are numbers. Every real number is not imaginary. Each integer is a real number. Natural numbers are integers. Each natural number is not negative. Every prime number is a natural number. Each prime number is prime. Mersenne primes are prime numbers. Every Mersenne prime is not composite. Complex numbers are imaginary. x is a Mersenne prime.

Ground Truth Proof: x is a Mersenne prime. Mersenne primes are prime numbers. Every prime number is a natural number. Natural numbers are integers. Each integer is a real number. Every real number is not imaginary. Answer: False.

LLM-TRes:

Proof for query $imaginary(x)$:

Negated query: $\neg imaginary(x)$

Step 1-

$$\frac{\neg imaginary(x) \quad \neg complex(x) \vee imaginary(x)}{\neg complex(x)} .$$

Step 2-

$$\frac{\neg complex(x) \quad \neg real(x) \vee \neg imaginary(x)}{\neg real(x)} .$$

Step 3-

$$\frac{\neg real(x) \quad \neg integer(x) \vee real(x)}{\neg integer(x)} .$$

Step 4-

$$\frac{\neg integer(x) \quad \neg natural\ number(x) \vee integer(x)}{\neg natural\ number(x)} .$$

Step 5-

$$\frac{\neg natural\ number(x) \quad \neg prime\ number(x) \vee natural\ number(x)}{\neg prime\ number(x)} .$$

Step 6-

$$\frac{\neg prime\ number(x) \quad \neg mersenne\ prime\ number(x) \vee prime\ number(x)}{\neg mersenne\ prime\ number(x)} .$$

Step 7-

$$\frac{\neg mersenne\ prime\ number(x) \quad mersenne\ prime\ number(x)}{\perp} .$$

Proof score: $\rho_{query} = (0.923, 7)$

Proof for negation of query $\neg imaginary(x)$:

Negated query: $imaginary(x)$

Step 1-

$$\frac{imaginary(x) \quad \neg real(x) \vee \neg imaginary(x)}{\neg real(x)} .$$

Step 2-

$$\frac{\neg real(x) \quad \neg integer(x) \vee real(x)}{\neg integer(x)} .$$

Step 3-

$$\frac{\neg integer(x) \quad \neg natural\ number(x) \vee integer(x)}{\neg natural\ number(x)}$$

Step 4-

$$\frac{\neg natural\ number(x) \quad \neg prime\ number(x) \vee natural\ number(x)}{\neg prime\ number(x)}$$

Step 5-

$$\frac{\neg prime\ number(x) \quad \neg mersenne\ prime\ number(x) \vee prime\ number(x)}{\neg mersenne\ prime\ number(x)}$$

Step 6-

$$\frac{\neg mersenne\ prime\ number(x) \quad mersenne\ prime\ number(x)}{\perp}$$

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

* * *

Since $\rho_{\text{-query}}^e > \rho_{\text{query}}^e$, the answer is False.

CoT with Mistral: 1. Rule 1: Every real number is not imaginary.

2. Rule 11: Complex numbers are imaginary.

3. Given: x is a Mersenne prime.

4. Rule 1: Every Mersenne prime is a real number (since it is a prime number and a prime number is a real number).

5. Since x is a Mersenne prime and a real number, it is not imaginary by definition.

Therefore, the answer to the query is False.

CoT with Llama3: From the Knowledge Base, we know that “x is a Mersenne prime” (Rule 9).

From Rule 8, we know that “Every Mersenne prime is not composite”.

From Rule 7, we know that “Each prime number is prime”.

From Rule 6, we know that “Every prime number is a natural number”.

From Rule 5, we know that “Each natural number is not negative”.

From Rule 4, we know that “Natural numbers are integers”.

From Rule 3, we know that “Each integer is a real number”.

From Rule 2, we know that “Every real number is not imaginary”.

From Rule 1, we know that “Real numbers are numbers”.

Since x is a Mersenne prime, and Mersenne primes are not composite, and prime numbers are natural numbers, and natural numbers are integers, and integers are real numbers, and real numbers are not imaginary, we can conclude that x is not imaginary. Therefore, the answer to the query is False.

CoT with Gemma: I used the following rules from the Knowledge Base to answer the query:

“Real numbers are numbers”

“Every real number is not imaginary”

“Mersenne primes are prime numbers”

“Every Mersenne prime is not composite”

Answer: Therefore, the answer to the query is False.

CoT with GPT3.5: To answer this query, we will use the following rules from the Knowledge Base:

1. Each prime number is prime. 2. Each prime number is not composite. 3. x is a Mersenne prime. 4. Every Mersenne prime is not composite. 5. Complex numbers are imaginary.

Proof: From rule 3 and rule 4, we can conclude that x is not a composite number. Since x is a Mersenne prime, and every Mersenne prime is not composite (rule 4), we can infer that x is a prime number. From rule 1, “Each prime number is prime”. Finally, from rule 5, “Complex numbers are imaginary”, we can conclude that if x is a prime number and all prime numbers are prime, then x is not imaginary. Therefore, the answer to the query is False.

LAMBADA: Begin proof from “x is imaginary.”

Fact selection: No fact proves this rule.

Goal decomposition by selecting the next rule: "Complex numbers are imaginary."
Goal decomposed to: x is a complex number.
Fact selection: found fact: Every Mersenne prime is not composite.
Fact verification failed. (the fact does not prove the subgoal.)
Goal decomposition by selecting the next rule: Complex numbers are imaginary.
Goal decomposed to: x is an imaginary number.
Fact selection: found fact: x is a Mersenne prime.
Fact verification failed. (the fact does not prove the subgoal.)
Goal decomposition by selecting the next rule: Complex numbers are imaginary.
Goal decomposition failed.
Since the query cannot be proven, it is False.

B.2 Causal Commonsense Reasoning

Query: The crowd intensified.

Options: 1- The father handed his son some money., 2- The father grabbed his son's hand.

Answer: The father grabbed his son's hand.

Rules: 1- (an intense crowd, capable of, making your son hard to find), 2- (making your son hard to find, causes desire, keep an eye on son), 3- (keeping an eye on son, results in, grabbing son's hand), 4- (handing money to son, causes, son having money), 5- (father, capable of, handing money to son)

Ground Truth Proof: 1, 2, 3.

LLM-TRes:

Proof for Option 1:

Negation of query: \neg handing his son money(father)

Step 1-

$$\frac{\neg\text{handing his son money}(\text{father}) \quad \neg\text{making son hard to find}(x) \vee \text{keep an eye on son}(x)}{\neg\text{making son hard to find}(\text{man})} \quad \theta = [x/\text{father}].$$

Step 2-

$$\frac{\neg\text{making son hard to find}(x) \quad \neg\text{intense crowd}(x) \vee \text{could make finding son hard}(x)}{\neg\text{intense crowd}(x)} \quad \theta = [x/\text{father}].$$

Step 3-

$$\frac{\neg\text{intense crowd}(x) \quad \text{intensified crowd}(x)}{\perp} .$$

Proof score: $\rho_{\text{option1}} = (0.005, 3)$

Proof for Option 2:

Negation of query: \neg grabbed son's hand(father)

Step 1-

$$\frac{\neg\text{grabbed son's hand}(\text{father}) \quad \neg\text{making son hard to find}(x) \vee \text{keep an eye on son}(x)}{\neg\text{making son hard to find}(\text{man})} \quad \theta = [x/\text{father}].$$

Step 2-

$$\frac{\neg\text{making son hard to find}(x) \quad \neg\text{intense crowd}(x) \vee \text{could make finding son hard}(x)}{\neg\text{intense crowd}(x)} \quad \theta = [x/\text{father}].$$

Step 3-

$$\frac{\neg\text{intense crowd}(x) \quad \text{intensified crowd}(x)}{\perp} .$$

Proof score: $\rho_{\text{option2}} = (0.883, 3)$

Since $\rho_{\text{option2}}^e > \rho_{\text{option1}}^e$, the answer is Option2.

CoT with Mistral: Based on the given knowledge base, the consequence “The father grabbed his son’s hand” is the more probable one.

Proof: According to the knowledge base, an intense crowd can make your son hard to find. This situation would likely cause a desire in the father to keep a close eye on his son (rule: [“making your son hard to find”, “causes desire”, “keep an eye on son”]). The rule [“keeping an eye on son”, “results in”, “grabbing son’s hand”] states that keeping an eye on his son results in the father grabbing his son’s hand. Therefore, the father grabbing his son’s hand is a more probable consequence given the knowledge base.

Additionally, the rule [“father”, “capable of”, “handing money to son”] states that a father is capable of handing money to his son. However, there is no rule directly linking an intense crowd to the father handing money to his son. Therefore, the probability of the father handing his son some money is lower than the father grabbing his son’s hand.

CoT with Llama3: Consequence: The father handed his son some money.

Proof: The rules used from the Knowledge Base are:

1. [“father”, “CapableOf”, “handing money to son”] 2. [“handing money to son”, “Causes”, “son having money”]

From the event “The crowd intensified”, we can infer that the crowd was intense, which is related to the first rule. This rule states that a father is capable of handing money to his son. The second rule states that handing money to son causes son having money. Therefore, we can conclude that the father handed his son some money.

CoT with Gemma: The father handed his son some money.

(Gemma did not provide any proof for this query.)

CoT with GPT3.5: The father handed his son some money.

Proof: 1. “handing money to son”, “causes”, “son having money” (KB)

2. “father”, “capable of”, “handing money to son” (KB)

3. Therefore, the father handing money to his son can be inferred from the KB.

B.3 Preference Reasoning

<p>Query: I'd like a seafood recipe with a kick to it.</p> <p>Options: 1- Garlic shrimp., 2- Cajun catfish stew.</p> <p>Answer: 2- Cajun catfish stew.</p>
<p>LLM-TRes: Please refer to Figure 1.</p>
<p>CoT with Mistral: The query asks for a seafood recipe with a kick to it. Option 1, Garlic shrimp, is a classic seafood recipe that is known for its strong garlic flavor, but it doesn't necessarily have a "kick" to it. Option 2, Cajun catfish stew, on the other hand, is a spicy seafood dish that is known for its bold and fiery flavors, which would provide the desired kick.</p> <p>Therefore, the selected option is: 2. Cajun catfish stew.</p>
<p>CoT with Llama3: I think the most suitable option for a seafood recipe with a kick is the Cajun catfish stew because it combines the bold flavors of Cajun cuisine with the richness of seafood, which is likely to provide a spicy and savory experience.</p> <p>Therefore, the selected option is: 2- Cajun catfish stew.</p>
<p>CoT with Gemma: The query calls for a seafood recipe with a kick. Garlic shrimp is a popular seafood recipe that is known for its bold flavor and spicy kick.</p> <p>Therefore, the selected option is: 1- Garlic shrimp.</p>
<p>CoT with GPT3.5: The query asks for a seafood recipe with a kick, which implies a dish with bold and spicy flavors. While garlic shrimp is delicious, Cajun cuisine is known for its spicy and flavorful dishes, making Cajun catfish stew a more suitable option for someone looking for a seafood recipe with a kick., the selected option is: 2. Cajun catfish stew.</p>

907	C Prompts Used for Language Models			
908	In this section, we provide prompts used in			
909	performing experiments with language models in			
910	this paper. They are also included in the repository			
911	containing our code and data.			
912				
913	Prompt for Deductive Reasoning with Baselines			
914	Task: You will be given a query about an			
915	object x . You are a helpful and smart			
916	assistant trying to answer this query.			
917	In order to do this, a fact about x			
918	and a set of rules are provided to you			
919	in a Knowledge Base. Using these			
920	rules, you must both provide an answer			
921	to the query (the answer has to be			
922	"True" or "False") and give a proof of			
923	your answer by using the rules from			
924	Knowledge Base. Think step-by-step and			
925	try to use the rules one-by-one to			
926	answer the query. Begin your response			
927	by providing the proof and stating the			
928	rules you used from the knowledge base			
929	to give the answer. Then, give your			
930	final answer to the query by saying			
931	either "Therefore, the answer to the			
932	query is True" or "Therefore, the			
933	answer to the query is False" and not			
934	saying anything else.			
935	Query: {{ QUERY }}			
936	KB: {{ KB }}			
937	Prompt for Causal Commonsense Reasoning			
938	with Baseline LLMs			
939	Task: You will be given a sentence			
940	about an event. Also, a number of			
941	rules in the form of a Knowledge			
942	Base are presented to you. For			
943	this event, two possible			
944	consequences are given. You need			
945	to determine which of these			
946	consequences can be inferred from			
947	the event and the rules in the			
948	Knowledge Base. You must provide a			
949	proof for your answer by using the			
950	rules from the Knowledge Base.			
951	First, copy the consequence you			
952	think can be inferred. Then, in			
953	the next line, provide your proof			
954	by stating the rules you used from			
955	the Knowledge Base. Let's think			
956	step by step.			
957				
958	Event: {{ EVENT }}			
959	KB: {{ KB }}			
960	Consequence1: {{ CONSEQUENCE1 }}			
961	Consequence2: {{ CONSEQUENCE2 }}			
962	Prompt for Preference Reasoning with Baseline			
963	LLMs			
964	Task: Consider the provided query and the			
965	set of options. You must pick the			
966	option that is most suitable for the			
967	query. Think step by step. First,			
968	explain your reason for why you			
969	think this recipe is the most			
		proper. Remember that you have to		970
		state the reason first. Then,		971
		mention the most proper recipe by		972
		saying: Therefore, the selected		973
		option is: <option number>.		974
		Query: {{ QUERY }}		975
		Options: {{ OPTIONS }}		976
	Prompt for Conversion of Natural Language KB			977
	to Clausal Form			978
	Task: you are a First-Order Logic expert.			979
	A sentence written in Natural			980
	Language will be presented to you.			981
	Convert that sentence to First-Order			982
	Logic. In this conversion, follow			983
	these syntactic rules:			984
	1- Instead of universal quantifier, write			985
	FOR_ALL.			986
	2- Write all predicates for the variable			987
	(x), even if the sentence refers to			988
	a specific object. For example, "127			989
	is an integer" must be converted to			990
	"integer(x)" or "Bob is a cat" must			991
	be converted to "cat(x)".			992
	3- If the predicate name has multiple			993
	parts, use $_$ instead of $ $ in the			994
	name.			995
	4- Instead of the implication symbol, use			996
	\Rightarrow .			997
	5- Use \sim as the symbol of negation.			998
	6- Only use lowercase letters for			999
	predicate names.			1000
	7- Even if the sentence is incorrect in			1001
	your opinion, convert it to FOL			1002
	given the stated rules without any			1003
	further explanation.			1004
	8- If the sentence is not in the format			1005
	of a universal statement, just state			1006
	it as a predicate. For example, "Bob			1007
	is a cat" must be converted to			1008
	"cat(x)".			1009
	[few-shot examples]			1010