
Can LLMs Reason with Rules? Logic Scaffolding for Stress-Testing and Improving LLMs

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Large language models (LLMs) have achieved impressive human-like performance
2 across various reasoning tasks. However, their mastery of underlying inferential
3 rules still falls short of human capabilities. To investigate this, we propose a
4 logic scaffolding inferential rule generation framework, to construct an inferential
5 rule base, ULogic, comprising both primitive and compositional rules across five
6 domains. Our analysis of GPT-series models over a rule subset reveals significant
7 gaps in LLMs’ logic understanding compared to human performance, especially in
8 compositional and structural complex rules with certain bias patterns. We further
9 distill these rules into a smaller-scale inference engine for flexible rule generation
10 and enhancing downstream reasoning. Through a multi-judger evaluation, our
11 inference engine proves effective in generating accurate, complex and abstract
12 conclusions and premises, and improve various commonsense reasoning tasks.
13 Overall, our work sheds light on LLMs’ limitations in grasping inferential rule and
14 suggests ways to enhance their logical reasoning abilities.

15 1 Introduction

16 “Did Leonardo da Vinci ever use a laptop for drawing pictures?” Large lan-
17 guage models can swiftly and confidently
18 respond “No” [10, 35], demonstrating im-
19 pressive reasoning ability that rivals hu-
20 man [18, 19]. However, when posed
21 with more obscure questions, such as
22 Q2 in Figure 1, LLMs are prone to ex-
23 hibit uncertainty and errors. This incon-
24 sistency raises concerns about whether
25 LLMs grasp the underlying logic of mat-
26 ters as proficiently as humans [38] (see
27 “underlying logic” in Figure 1) and high-
28 lights challenging reasoning situations
29 (like Q2) where current LLMs might strug-
30 gle. Humans naturally abstract underlying
31 logic as inferential rules from extensive
32 real-world observations [3], beneficial for
33 addressing diverse reasoning situations. An inferential rule is typically defined as a premise with a
34 set of facts (e.g., “Person X died before ... earlier than B”) leading to a conclusion (e.g., “Person X
35 cannot access Object Y”) [5]. Grasping this rule enables the deduction that a person cannot access an
36 object invented posthumously. This work utilizes symbolic logic as a *scaffold* to generate challenging
37

Q1: Did Leonardo da Vinci ever use a laptop for drawing pictures?

Q2: Jane wrote a novel published by Jimmy, a publisher born in 1750. Did Jane’s grandmother often work by car?



Underlying Logic:

If Person X died before year A and Object Y was invented in year B, and A is earlier than B, then **Person X can not access Object Y.**

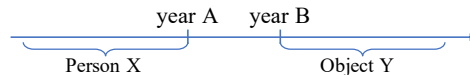


Figure 1: The underlying logic to answer Q1 and Q2.

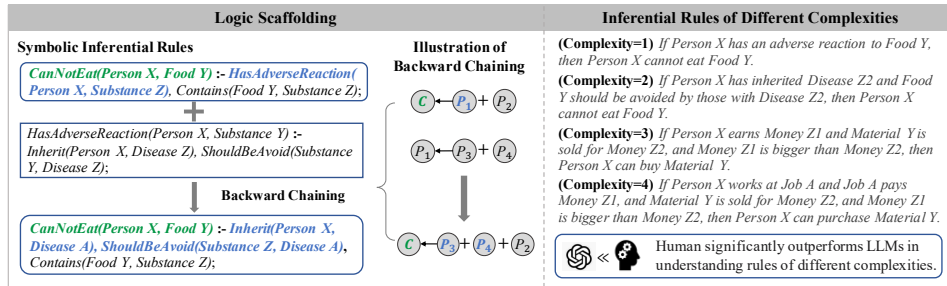


Figure 2: Logic scaffolding uncovers challenging reasoning space for LLMs.

38 reasoning situations for GPT-series LLMs, as shown in Figure 2. A discernible gap exists between
 39 LLMs and humans in understanding inferential rules, especially rules with complex premises.

40 However, collecting inferential rules at scale is challenging. Previous work relies on manual cu-
 41 ration [24, 29] or inductive logic programming [22], which are either labor-intensive or limited in
 42 diversity. Besides, manually crafted rules often appear simple and overly specified, struggling to move
 43 beyond basic intuition or generalize across diverse situations. For example, the rule *If X runs out of*
 44 *steam, then X becomes tired* from [24] has only one premise fact and narrowly specifies exhaustion.

45 To this end, we introduce **Logic scaffolding Inferential Rule gEneration (LOIRE)**, a framework
 46 to generate inferential rules of different complexities. LOIRE operates in two stages: primitive rule
 47 generation and rule composition. Initially, we define “primitive rules” to describe abstract objects
 48 like *Person* and *Food*, and ensure they cannot be decomposed into simpler rules, facilitating broad
 49 generalization and easy generation. We then incorporate GPT-4’s generative capability and human
 50 expertise to generate primitive rules with high confidence. This process, consistently guided by
 51 symbolic logic, involves GPT-4 drafting potential conclusions in various domains, and forming
 52 premises with one or more facts. We ensure rules’ logical soundness through the model’s self-critique
 53 and human manual verification. In the second stage, we apply backward chaining [8, 1] upon primitive
 54 logical rules to automatically construct compositional rules of varied lengths and structures at scale.

55 Using this framework, we construct **ULogic**, an inferential rule base with around 8,000 primitive and
 56 6,000 compositional rules across five domains: object affordance, accessibility, interaction, location,
 57 and human need. We hope ULogic will serve as a valuable resource, facilitating the assessment
 58 of LLMs’ proficiency in underlying logic and enhancing flexible rule generation and downstream
 59 reasoning. We use ULogic to create an entailment probing task with a comprehensive and robust
 60 evaluation strategy, comparing LLMs’ grasp of inferential rules to human performance. Our analysis
 61 of GPT-series LLMs (GPT-4, GPT-3.5-Turbo and GPT-3.5-Turbo-Instruct) indicates they have a
 62 basic understanding of inferential rules but fall short of human proficiency, especially in rules with
 63 complex premises. Specifically, all models struggle more as the compositional complexity increases.
 64 While GPT-4 performs consistently on verbalized and symbolic rules, the other models sharply
 65 degrade on symbolic rules. Additionally, all models exhibit disparities on various rule structures with
 66 Disjunctive-Transitive rules posing the greatest challenges. Moreover, these LLMs display notable
 67 polarity biases with GPT-4 showing a necessary bias, underscoring areas for improvement.

68 We further distill crafted inferential rules into a smaller-scale inference engine for flexible rule gener-
 69 ation and downstream reasoning. We design three tasks: conclusion generation, premise completion
 70 and premise generation, to construct an instruction-tuning dataset for inferential rule distillation.
 71 Experimental results through a multi-judger evaluation mechanism incorporating automatic metrics,
 72 LLM evaluators and human preferences show that our inference engine possesses the ability for these
 73 three tasks. It outperforms GPT-3.5-Turbo across all dimensions of three tasks and even surpasses
 74 GPT-4 in generating more complex and abstract rules. Moreover, it can generate logical rules that
 75 enhance downstream commonsense reasoning.

76 2 Logic Scaffolding for Inferential Rule Generation

77 2.1 Preliminary of Inferential Rules

78 To better control the generative capability of LLMs for rule generation, we focus on *if-then* inferential
 79 rules with variables, that can be easily expressed as symbolic logic [16]. An inferential rule describes

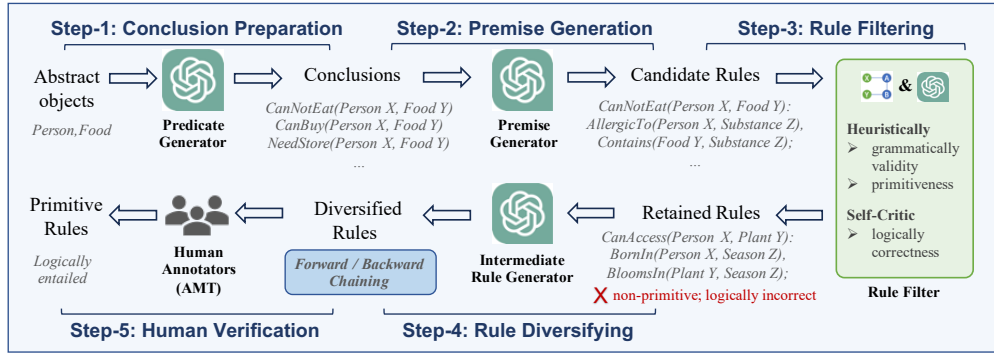


Figure 3: The pipeline of primitive rule generation.

80 a logical implication from a premise (a set of facts) to a conclusion (a specific fact), where each fact is
 81 a predicate expression with two variables, and each variable has a designated variable type. For each
 82 rule, we employ logic scaffolding which first generates its symbolic expression to consistently guide
 83 its verbalized form. We utilize Prolog [2] to formulate symbolic rules as $\text{Conclusion}:-\text{Premise}$,
 84 where $:-$ indicates the logical implication. For example,

$$\begin{aligned} \text{CanNotEat}(\text{Person } X, \text{Food } Y):- \\ \text{AllergicTo}(\text{Person } X, \text{Substance } Z), \text{Contains}(\text{Food } Y, \text{Substance } Z). \end{aligned} \quad (1)$$

85 The left-hand side is the conclusion and the right hand lists premise facts connected by commas.
 86 “CanNotEat”, “AllergicTo” and “Contains” are predicate verbs while *Person*, *Food*, *Substance*
 87 are variable types of variables (*X*, *Y*, *Z*). This symbolic rule can be verbalized as: *If Person X is*
 88 *allergic to Substance Z and Food Y contains Substance Z, then Person X cannot eat Food Y.*

89 **Primitive Rule** We aim to generate primitive rules for further compositions and potential generaliza-
 90 tion. We formally define primitive rules as follows: (1) they concern abstract objects, like *Person*
 91 and *Food*, rather than specific instances, and their common properties; (2) they cannot be decomposed
 92 into simpler rules. Inspired by superordinate objects such as *instrument*, *fruit*, *tool* from
 93 [23], we assemble a collection of abstract objects. We first identify the most common tail nodes of
 94 “IsA” relations from ConceptNet [30]. For those nodes that are still fine-grained, we further seek
 95 their general hypernyms by searching ConceptNet and WordNet [15]. We totally gather a list of 32
 96 most common abstract objects for primitive rule generation, with 18 common properties generated by
 97 prompting GPT-4, as detailed in Appendix A.1.

98 2.2 Primitive Rule Generation Pipeline

99 The pipeline of primitive rule generation is illustrated in Figure 3, consisting of five steps. First,
 100 we randomly select two abstract objects, and generate potential predicates between them to form
 101 conclusions. GPT-4 is prompted to generate corresponding feasible premises with both single and
 102 multiple facts, thereby constructing candidate primitive rules. We then apply heuristic methods to
 103 filter invalid and non-primitive rules, and utilize GPT-4 to select the rules it deems logically correct.
 104 We further diversify rule predicates via backward/forward chaining [34, 27] with generated single-fact
 105 rules, and filter excessively repetitive rules. Finally, the diversified rules undergo manual verification
 106 to ensure the final set of high-confidence primitive rules.

107 **Step-1: Conclusion Preparation** From the set of abstract objects, we select any two, e.g., *Person*
 108 and *Food*, and prompt GPT-4 to generate potential predicates connecting them as conclusions, e.g.,
 109 *CanEat(Person X, Food Y)*. We attempt every possible pairing of two, where the selected objects can
 110 be identical. For each pair of objects, $\{\text{object}_1\}$ and $\{\text{object}_2\}$, we aim to generate conclusions across
 111 five domains: {object affordance, accessibility, interaction, location and person’s need}, thereby
 112 covering diverse scenarios. Explanations and example rules of these domains, and the prompt are
 113 listed in Appendix A.2. Besides, we negate the generated predicates to yield both positive and
 114 negative conclusions, e.g., *CanNotEat(Person X, Food Y)*, across object affordance, accessibility, and
 115 interaction domains, building a complete rule set.

116 **Step-2: Premise Generation** Guided by a symbolic conclusion, we prompt GPT-4 to generate its
 117 premises in both symbolic and verbalized forms for better controllability. This process involves the

118 logit bias setting, motivating premises to describe relationships between abstract objects and their
 119 properties. Specifically, premises are generated under the constraint of logit bias, increasing the
 120 likelihood of these objects and properties appearing in the output. For each conclusion, we create both
 121 single-fact and multi-fact premises to yield candidate rules of varying lengths. We tailor instructions
 122 and demonstrations for each domain to prompt GPT-4 for premise generation exploring different
 123 possibilities, as detailed in Appendix A.4.

124 **Step-3: Rule Filtering** After over-generating candidate primitive rules, we first design heuristic
 125 methods to filter grammatically invalid or non-primitive rules based on their symbolic forms. For
 126 grammatically validity, we check if the variables in the premises form a connected graph from
 127 node “X” to node “Y”, as in Appendix A.5. For primitiveness, we exclude rules with non-primitive
 128 variable types or those comprising more than 3 premise facts. Besides, we eliminate trivial rules
 129 containing negative words in both the premise and conclusion, e.g., *CanNotEat(Person X, Food Y):-*
 130 *CanNotAccess(Person X, Food Z)*. Since directly generating logically correct rules is challenging,
 131 we further adopt a self-critic strategy [11] where GPT-4 critiques the accuracy of its self-generated
 132 rules in a verbalized format, and provides explanations. When prompting GPT-4, we include two
 133 demonstrations featuring both correct and incorrect rules to mitigate label bias. These demonstrations
 134 vary across different domains. An example prompt for object affordance is in Appendix A.6.

135 **Step-4: Rule Diversifying** To increase the variety of rule expressions, we diversify predicates
 136 while maintaining its logical accuracy. Based on symbolic rules, we respectively apply forward
 137 and backward chaining algorithms to their conclusion and premise with generated single-fact rules,
 138 as shown in Figure 4. In forward chaining, we take the conclusion as a new premise to generate
 139 an intermediate single-fact rule, subsequently substituting the original conclusion with this newly
 140 derived conclusion. In backward chaining, a premise is taken as a conclusion to create an intermediate
 141 single-fact rule, and replace the original premise with the new-generated one. Intermediate single-fact
 142 rules are also generated through Step-2 and 3. Each original rule undergoes one forward and one
 backward chaining to derive two diversified rules.

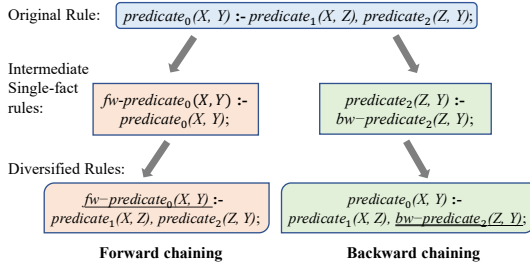


Figure 4: The forward and backward chaining process for diversifying rules.

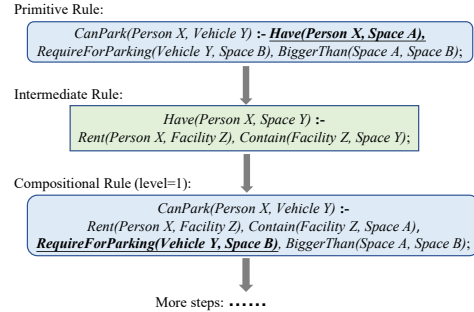


Figure 5: Illustration of one backward chaining step.

143

144 **Step-5: Human Verification** To obtain more reliable rules, we utilize Amazon Mechanical Turk
 145 (AMT) to recruit three annotators for manual verification of each rule. They are asked to assess
 146 the clarity and comprehensibility of its premise and conclusion, and the logical entailment from the
 147 premise to the conclusion. Only the rules unanimously validated by all three annotators are preserved.
 148 The AMT template for human verification and rule acceptance rates are listed in Appendix A.7.

149 2.3 Rule Composition

150 We create more compositional rules by applying backward chaining upon primitive rules with different
 151 chaining steps. In each step, we select a premise fact from the current rule as a conclusion, deriving a
 152 new primitive rule that describes its multi-fact premise. This selected fact is then replaced with the
 153 newly generate premise. This process is iteratively conducted 1 to 3 times, creating rules with varying
 154 compositional levels (1 to 3). An example of one backward chaining step is shown in Figure 5. The
 155 intermediate primitive rules used in backward chaining are generated via the pipeline described in
 156 Sec. 2.2, thus also contributing to our primitive rule set. As the composition of logically correct
 157 sub-rules is also logically correct, there is no need to verify these compositional rules separately.

158 **2.4 Rule Statistics**

159 Using LOIRE framework, we construct an inferential rule base ULogic comprising 14,647 rules,
 160 with 7,967 primitive and 6,680 compositional ones. These rules span five key domains: object
 161 affordance, accessibility, interaction, location and person’s need. They vary in compositional depth
 162 from 0 to 3, with rule lengths ranging from 1 to 6. Detailed statistics are in Appendix A.8.

163 **3 Assessing LLMs’ Proficiency in Capturing Inferential Rules**

164 We utilize ULogic for a systematic evaluation of LLMs’ proficiency in underlying logic compared to
 165 human competence. Specifically, we select a high-quality probing subset of 1,104 diverse, author-
 166 verified rules from our rule base (varying in lengths, polarities and structures), and create a binary
 167 entailment classification task for assessing LLMs’ ability to capture inferential entailment.

168 **3.1 Analysis Setup**

169 Considering LLMs’ sensitivity to various input formulations and shortcut biases, we design a
 170 comprehensive and robust assessment mechanism to ensure reliable analysis. For each inferential
 171 rule, we convert it into five distinct probing questions to mitigate template bias, as summarized in
 172 Appendix B.1. We report the average accuracy and variance (the error line of each bar) across five
 173 templates. Besides, we adopt a two-shot chain of thought (CoT) prompting strategy [39] requiring the
 174 model to generate a rationale after presenting its answer, using "and also explain why." We include
 175 one correct rule and one incorrect rule in the two demonstrations to minimize label bias.

176 Following the Law of Non-Contradiction [21], the propositions "If X, then Y" and "If X, then not
 177 Y" are mutually exclusive that cannot both be true at the same time. To enhance the reliability
 178 of our probing, we flip each rule by negating its conclusion, and simultaneously probe both the
 179 original rule and its flipped version. A rule is accurately classified only if the original rule is affirmed
 180 (True/Right/Yes) and its flipped counterpart is negated (False/Wrong/No), as shown below. A specific
 example is in Appendix B.2. This dual-sided probing is applied to both human and LLMs.

<i>If Premise, then Conclusion_{original}.</i>	True/Right/Yes
<i>If Premise, then Conclusion_{flipped}.</i>	False/Wrong/No

181

182 **3.2 Empirical Analysis**

183 We conduct analysis on GPT-series LLMs, including GPT-4, GPT-3.5-Turbo and GPT-3.5-Turbo-
 184 Instruct, aiming to investigate LLMs’ proficiency of inferential rules against human performance by
 185 exploring the following questions. The human performance is obtained by asking AMT annotators
 186 whether the input rule is logical correct with high probability. Each performance presented in
 187 following bar charts is calculated based on 150 instances randomly sampled from our probing subset.

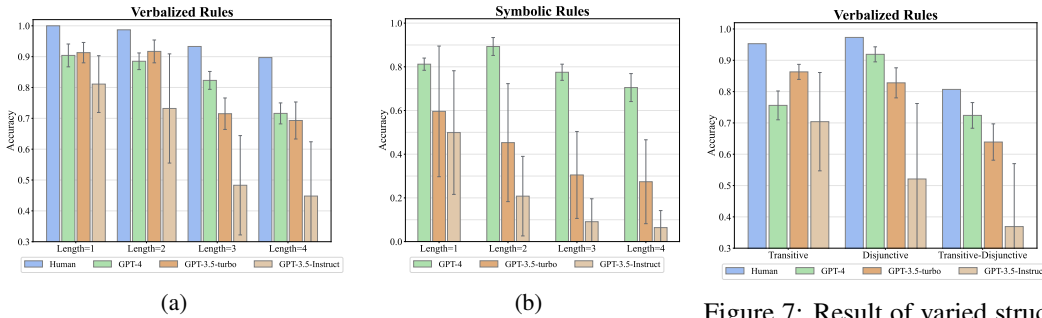


Figure 6: Probing results across varied lengths.

Figure 7: Result of varied structures.

188

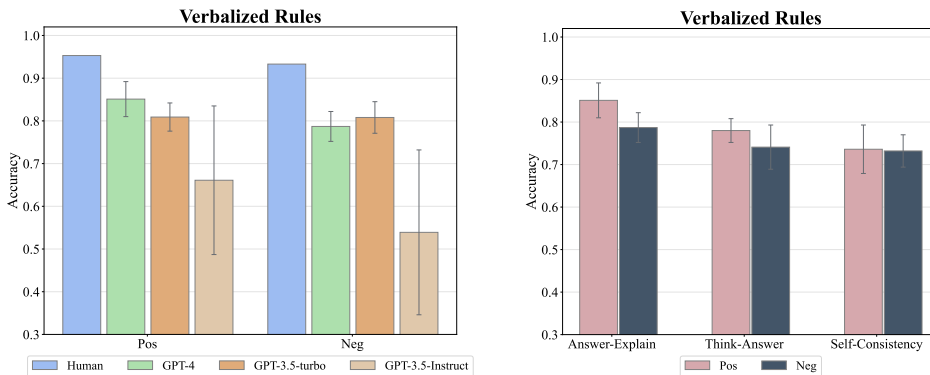
189 **(1) How does model performance vary with increasing compositional complexity?** We conduct
 190 rule probing in terms of different compositional lengths, as illustrated in Figure 6a. "Length=1,2,3,4"

191 respectively denote rules with 1~4 facts in their premises. The analysis of different compositional
 192 depths is also provided in Appendix B.3. They both reveal that as compositional complexity increases,
 193 the performance of both human and all models drops. The primary reason is that compositional
 194 complex rules typically necessitate the aggregation of multi-step reasoning, which escalates higher-
 195 order relationships understanding and exponential error accumulation with each additional step [7].
 196 Besides, there is a persistent performance gap between all models and human, particularly pronounced
 197 with compositional complex rules, suggesting significant potential for enhancement in this area.

198 **(2) Are LLMs proficient in capturing both symbolic and verbalized rules?** We further analyze
 199 LLM performance on symbolic rules (see Figure 6b) compared to on verbalized rules in Figure 13.
 200 GPT-4 achieves consistent performance on verbalized and symbolic rules, whereas GPT-3.5-Turbo
 201 and GPT-3.5-Instruct sharply degrade on symbolic rules. This suggests that the GPT-3.5 series may
 202 have limitations in generalizing across varied types of linguistic structures beyond natural language,
 203 whereas GPT-4 likely have undergone specific optimizations for symbolic interpretations.

204 **(3) Are there performance disparities among models concerning different rule structures?**
 205 Our generated multi-fact rules (Length > 1) have three intrinsic structures: Transitive, Disjunctive
 206 and Disjunctive-Transitive. Specific illustrations and examples of each structure are detailed in
 207 Appendix B.4. Figure 7 shows that Disjunctive-Transitive rules pose greater challenges compared to
 208 Transitive and Disjunctive ones, especially for GPT-3.5-Turbo and GPT-3.5-Instruct. We hypothesize
 209 that this discrepancy stems from increased compositional complexity and LLMs’ insufficient learning
 210 of logical structures in natural language.

211 **(4) Do LLMs exhibit a polarity bias over inferential rules?** Our inferential rules contain both
 212 positive and negative conclusions. As shown in Figure 8a, GPT-4 and GPT-3.5-Instruct exhibit
 213 a pronounced positive bias, performing better on rules with positive conclusions. This bias may
 214 originate from the imbalanced distribution of LLMs’ training data [9], with a higher proportion of
 215 positive statements. We further explore different CoT strategies with GPT-4: (1) first answer then
 216 explain (*Answer-Explain*), (2) first think then answer (*Think-Answer*), (3) self-consistently think
 217 then answer (*Self-Consistency*) [37]. Various CoT prompts are listed in Appendix B.5. Figure 8b
 218 shows that although advanced CoT strategies can mitigate the positive bias, they adversely impact the
 219 performance on rules with both positive and negative conclusions.



(a) Answer-Explain strategy. (b) Various CoT strategies.

Figure 8: Rule Polarity Comparison.

220 **(5) Why does GPT-4 significantly underperform GPT-3.5-Turbo on transitive rules?** While
 221 GPT-4 generally outperforms or matches other models, this superiority disappears on transitive rules,
 222 as evidenced in Figure 7. We investigate this question in Appendix B.6, which reveals that GPT-4
 223 exhibits a “necessary bias” that tend to consider all necessary conditions reaching a conclusion,
 224 avoiding definite judgement. This conservative style may come from LLMs’ preference alignment
 225 during Reinforcement Learning with Human Feedback [19].

226 Overall, GPT-4 performs best in grasping inferential rules. But compared to human performance, there
 227 still remains substantial room for improvement across all models, especially in highly compositional,
 228 symbolic and structural complex rules. Besides, all models tend to exhibit a polarity bias towards
 229 rules with positive conclusions with GPT-4 also showing a necessary bias. These findings suggest
 230 potential areas for future enhancements.

231 4 Rule Distillation as Inference Engine

232 4.1 Instruction Dataset & Model Tuning

233 For flexible rule generation and benefiting downstream reasoning, we distill our crafted rules into
234 a smaller-scale inference engine as illustrated in Appendix C.1. We tailor three tasks: conclusion
235 generation, premise completion and premise generation, to construct an instruction-tuning dataset for
236 inferential rule distillation. The detailed definitions of these tasks are also described in Appendix C.1.

237 We gather all primitive rules and partial compositional rules to formulate the instruction-tuning
238 dataset, as compositional rules are constructed from primitive ones. We take 10,703 rules for training
239 and 943 for testing. Altogether, we create 39,887 instances for instruction tuning, including 10,703,
240 18,500 and 10,684 for conclusion generation, premise completion and premise generation. We have
241 3,500 testing instances, divided as 943, 1,614 and 943 for these three tasks. We use Mistral-7b [13] as
242 the backbone model and fine-tune it with our constructed instruction dataset as our inference engine.
243 The training details and demo page can be found in Appendix C.2.

244 4.2 Rule Generation Evaluation

245 We compare our inference engine against GPT-4 and GPT-3.5-Turbo across three tasks to assess rule
246 generation. For a fair comparison, we prompt GPT-4 and GPT-3.5-Turbo to simultaneously generate
247 symbolic and verbalized responses, using similar prompts as in Step-2 of Sec. 2.2. Detailed prompts
248 are in Appendix C.3. We introduce a multi-judger evaluation mechanism, incorporating automatic
249 metrics, LLM evaluator and human preference to evaluate logical accuracy in conclusion generation
250 and premise completion. For premise generation task with a specified number of facts, we generate
251 three potential premises for each conclusion, and evaluate them on accuracy, diversity, complexity
252 and abstractness (see Appendix C.4 for detailed metric definitions).

253 **Automatic Evaluation** For automatic accuracy evaluation of three tasks, we calculate BLEU
254 score [20] against reference responses. For complexity of premise generation, we assess the average
255 fact number of three generated premises. For diversity, we compute average Self-BLEU [28, 32]
256 between three generated premises. Specifically, Self-BLEU measures the BLEU score of a generated
257 premise against another, and a high average Self-BLEU indicates low diversity. Abstractness is not
easy to evaluate automatically, so we leave it to LLM evaluation. The results are shown in Table 1.

Table 1: Automatic evaluation results.

Task	Conclusion Generation	Premise Completion	Premise Generation		
Metrics	BLEU	BLEU	BLEU	Self-BLEU	Fact Num.
Engine	0.739	0.527	0.411	0.687	3.42
GPT-4	0.414	0.179	0.149	0.805	2.58
GPT-3.5	0.338	0.248	0.084	0.739	1.72

258

259 **LLM Evaluation** We adopt GPT-4 as an evaluator to rate the generated responses on a scale from
260 1 to 3. The criteria of each rating along with examples are provided to the evaluator. Please see
261 Appendix C.5 for detailed prompts. For each task, we select 100 instances for LLM evaluation,
ensuring a balance across all domains and all types. The rating results are presented in Table 2.

Table 2: LLM evaluation results.

Task	Conclusion Generation	Premise Completion	Premise Generation			
Metrics	Accuracy	Accuracy	Accuracy	Diversity	Complexity	Abstractness
Engine	2.44	2.78	2.34	1.89	1.62	2.43
GPT-4	2.53	2.72	2.77	2.64	1.40	2.32
GPT-3.5	2.38	1.57	1.91	1.72	1.06	2.30

262

263 **Human Evaluation** To better assess premise generation in line with human value, we further recruit
264 two annotators for each instance to compare their accuracy. We implement a pairwise comparison
265 setting, asking annotators to determine which group of generated premise is more accurate in terms
266 of logical consistency with the given conclusion, commonsense alignment and correctness of fact
267 numbers. The results are shown in Figure 9. From all evaluation, we can see that our inference engine
268 enables the smaller-scale LLM with the capability for conclusion generation, premise completion and

premise generation. It performs better than GPT-3.5-Turbo across all metrics in three tasks, and even outperforms GPT-4 to generate more complex and abstract rules.

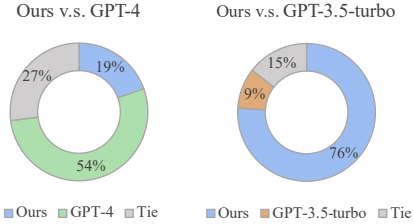


Figure 9: Human comparison results.

Figure 10: Downstream reasoning performance.

Dataset	Mistral	Mistral+rules (Mistral-7b)	LLama	LLama+rules (LLama2-7b)
StrategyQA	54.50	56.75	58.00	60.48
SOCIAL IQA	64.00	68.50	53.50	60.50
LINK head	53.68	68.38	58.09	70.59
LINK longtail	53.33	67.50	55.83	65.00
PIQA	65.00	65.00	58.5	62.0
CSQA2.0	59.00	62.50	64.00	60.00

270

271 4.3 Downstream Reasoning Evaluation

272 We further analyze the effectiveness of our inference engine in generating logical rules or explanations
 273 to enhance downstream reasoning tasks. We evaluate on following commonsense reasoning datasets:
 274 StrategyQA [10], SOCIAL IQA [25], LINK [14], PIQA [4] and CSQA2.0 [31]. We use a zero-shot
 275 CoT strategy to prompt two baseline models, Mistral-7B-Instruct-v0.1 and Llama-2-7b-chat [33],
 276 to answer questions with following explanations. We then utilize our inference engine to generate
 277 logical rules or explanations relevant to answer questions, and supplement these generated rationals
 278 to baseline models as input to enhance their performance. We compare the prediction accuracy of our
 279 inference engine augmented models against baselines. The comparative results are shown in Tabel 10.
 280 Our inference engine can generate logical rules or explanations that benefit multiple downstream
 281 commonsense reasoning tasks on top of different backbone models. For the lack of clear advantage
 282 on PIQA and performance decline on CSQA2.0, we speculate that PIQA may be contaminated during
 283 Mistral’s training process, and CSQA2.0’s focus is mainly on longtail commonsense knowledge
 284 rather than requiring logical rules inference, like "Is cotton candy sometimes made out of cotton?"

285 5 Related Work

286 **Logical Rule Generation** Logical inferential rules are crucial for everyday reasoning [10, 31],
 287 and collecting these inferential rules is challenging. Prior work mainly adopts inductive logic
 288 programming (ILP) [41, 22, 26] for rule generation. However, they can only generate rules from
 289 existing knowledge graphs and the generated rules has potential inaccuracies. Alternatively, [29]
 290 manually create a set of inferential rules for inductive reasoning, but their scope is limited to kinship.
 291 [24] construct a commonsense inferential rule base through crowdsourcing, but these rules tend to
 292 be overly simple and specific, struggling to move beyond basic intuition and generalize to varied
 293 situations. Abstract and complex rules are essential in tackling diverse complex questions, paving
 294 the way for complex reasoning and decision-making. Although LLMs have opened new avenues for
 295 generating inferential rules [42], they still struggle to automatically craft abstract and complex rules.

296 **Integration of Logical Rules and LLMs** The integration of inferential rules with LLMs has gained
 297 significant attention. This approach combines the logical interpretability of symbolic reasoning and
 298 adaptive power of neural computing, improving LLMs’ logical reasoning ability. [36, 17] transform
 299 textual statements into logical expressions and conduct symbolic reasoning following logical rules.
 300 [40] train neural models using a set of inferential rules for dynamic application. This direction
 301 broadens LLMs’ ability with flexible rule generation and application for complex reasoning.

302 6 Conclusion

303 This paper examines GPT-series LLMs’ proficiency in capturing logical inferential rules and probes
 304 their challenging reasoning space. We introduce a logic scaffolding inferential rule generation
 305 (LOIRE) framework to create an inferential rule base ULogic, including nearly 8,000 primitive and
 306 6,000 compositional rules across five domains. Our evaluations show that even advanced models
 307 like GPT-4 struggle with compositional and structural complex rules and exhibit certain biases.
 308 Furthermore, we distill ULogic into a smaller inference engine that performs well in generating
 309 inferential rules and benefit downstream reasoning tasks. Our work points out where LLMs need to
 310 improve in logical reasoning and offers a pathway to enhance their reasoning capabilities.

References

- [1] A. Al-Ajlan. The comparison between forward and backward chaining. *International Journal of Machine Learning and Computing*, 5(2):106, 2015.
- [2] K. R. Apt et al. *From logic programming to Prolog*, volume 362. Prentice Hall London, 1997.
- [3] J. Barwise. Everyday reasoning and logical inference. *Behavioral and Brain Sciences*, 16(2):337–338, 1993.
- [4] Y. Bisk, R. Zellers, J. Gao, Y. Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.
- [5] P. Boghossian. What is inference? *Philosophical studies*, 169:1–18, 2014.
- [6] T. Dettmers, A. Pagnoni, A. Holtzman, and L. Zettlemoyer. Qlora: Efficient finetuning of quantized llms. *arXiv preprint arXiv:2305.14314*, 2023.
- [7] N. Dziri, X. Lu, M. Sclar, X. L. Li, L. Jian, B. Y. Lin, P. West, C. Bhagavatula, R. L. Bras, J. D. Hwang, et al. Faith and fate: Limits of transformers on compositionality. *arXiv preprint arXiv:2305.18654*, 2023.
- [8] H. Gallaire and J. Minker. *Logic and data bases*. Springer Science & Business Media, 2012.
- [9] A. Garg, D. Srivastava, Z. Xu, and L. Huang. Identifying and measuring token-level sentiment bias in pre-trained language models with prompts. *arXiv preprint arXiv:2204.07289*, 2022.
- [10] M. Geva, D. Khashabi, E. Segal, T. Khot, D. Roth, and J. Berant. Did aristotle use a laptop? a question answering benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9:346–361, 2021.
- [11] Z. Gou, Z. Shao, Y. Gong, Y. Shen, Y. Yang, N. Duan, and W. Chen. Critic: Large language models can self-correct with tool-interactive critiquing. *arXiv preprint arXiv:2305.11738*, 2023.
- [12] E. J. Hu, Y. Shen, P. Wallis, Z. Allen-Zhu, Y. Li, S. Wang, L. Wang, and W. Chen. Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.
- [13] A. Q. Jiang, A. Sablayrolles, A. Mensch, C. Bamford, D. S. Chaplot, D. d. I. Casas, F. Bressand, G. Lengyel, G. Lample, L. Saulnier, et al. Mistral 7b. *arXiv preprint arXiv:2310.06825*, 2023.
- [14] H. Li, Y. Ning, Z. Liao, S. Wang, X. L. Li, X. Lu, F. Brahman, W. Zhao, Y. Choi, and X. Ren. In search of the long-tail: Systematic generation of long-tail knowledge via logical rule guided search. *arXiv preprint arXiv:2311.07237*, 2023.
- [15] G. A. Miller. Wordnet: a lexical database for english. *Communications of the ACM*, 38(11):39–41, 1995.
- [16] V. Novák and S. Lehmke. Logical structure of fuzzy if-then rules. *FUZZY sets and systems*, 157(15):2003–2029, 2006.
- [17] T. X. Olausson, A. Gu, B. Lipkin, C. E. Zhang, A. Solar-Lezama, J. B. Tenenbaum, and R. Levy. Linc: A neurosymbolic approach for logical reasoning by combining language models with first-order logic provers. *arXiv preprint arXiv:2310.15164*, 2023.
- [18] OpenAI. Gpt-4 technical report, 2023.
- [19] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [20] K. Papineni, S. Roukos, T. Ward, and W.-J. Zhu. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, pages 311–318, 2002.

- 355 [21] G. Priest, J. C. Beall, and B. Armour-Garb. *The law of non-contradiction: New philosophical*
356 *essays*. Clarendon Press, 2006.
- 357 [22] M. Qu, J. Chen, L.-P. Xhonneux, Y. Bengio, and J. Tang. Rnnlogic: Learning logic rules for
358 reasoning on knowledge graphs. *arXiv preprint arXiv:2010.04029*, 2020.
- 359 [23] E. Rosch and C. B. Mervis. Family resemblances: Studies in the internal structure of categories.
360 *Cognitive psychology*, 7(4):573–605, 1975.
- 361 [24] M. Sap, R. Le Bras, E. Allaway, C. Bhagavatula, N. Lourie, H. Rashkin, B. Roof, N. A. Smith,
362 and Y. Choi. Atomic: An atlas of machine commonsense for if-then reasoning. In *Proceedings*
363 *of the AAAI conference on artificial intelligence*, volume 33, pages 3027–3035, 2019.
- 364 [25] M. Sap, H. Rashkin, D. Chen, R. LeBras, and Y. Choi. Socialliqa: Commonsense reasoning
365 about social interactions. *arXiv preprint arXiv:1904.09728*, 2019.
- 366 [26] P. Sen, B. W. de Carvalho, R. Riegel, and A. Gray. Neuro-symbolic inductive logic programming
367 with logical neural networks. In *Proceedings of the AAAI Conference on Artificial Intelligence*,
368 volume 36, pages 8212–8219, 2022.
- 369 [27] H. Shindo, D. S. Dhami, and K. Kersting. Neuro-symbolic forward reasoning. *arXiv preprint*
370 *arXiv:2110.09383*, 2021.
- 371 [28] R. Shu, H. Nakayama, and K. Cho. Generating diverse translations with sentence codes. In
372 *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages
373 1823–1827, 2019.
- 374 [29] K. Sinha, S. Sodhani, J. Dong, J. Pineau, and W. L. Hamilton. Clutrr: A diagnostic benchmark
375 for inductive reasoning from text. *arXiv preprint arXiv:1908.06177*, 2019.
- 376 [30] R. Speer, J. Chin, and C. Havasi. Conceptnet 5.5: An open multilingual graph of general
377 knowledge. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- 378 [31] A. Talmor, O. Yoran, R. L. Bras, C. Bhagavatula, Y. Goldberg, Y. Choi, and J. Berant. Common-
379 senseqa 2.0: Exposing the limits of ai through gamification. *arXiv preprint arXiv:2201.05320*,
380 2022.
- 381 [32] G. Tevet and J. Berant. Evaluating the evaluation of diversity in natural language generation.
382 *arXiv preprint arXiv:2004.02990*, 2020.
- 383 [33] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra,
384 P. Bhargava, S. Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv*
385 *preprint arXiv:2307.09288*, 2023.
- 386 [34] J. Urbani, F. Van Harmelen, S. Schlobach, and H. Bal. Querypie: Backward reasoning for owl
387 horst over very large knowledge bases. In *The Semantic Web–ISWC 2011: 10th International*
388 *Semantic Web Conference, Bonn, Germany, October 23-27, 2011, Proceedings, Part I 10*, pages
389 730–745. Springer, 2011.
- 390 [35] L. Wang, W. Xu, Y. Lan, Z. Hu, Y. Lan, R. K.-W. Lee, and E.-P. Lim. Plan-and-solve prompting:
391 Improving zero-shot chain-of-thought reasoning by large language models. *arXiv preprint*
392 *arXiv:2305.04091*, 2023.
- 393 [36] S. Wang, W. Zhong, D. Tang, Z. Wei, Z. Fan, D. Jiang, M. Zhou, and N. Duan. Logic-
394 driven context extension and data augmentation for logical reasoning of text. *arXiv preprint*
395 *arXiv:2105.03659*, 2021.
- 396 [37] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou.
397 Self-consistency improves chain of thought reasoning in language models. *arXiv preprint*
398 *arXiv:2203.11171*, 2022.
- 399 [38] P. C. Wason. Reasoning about a rule. *Quarterly journal of experimental psychology*, 20(3):273–
400 281, 1968.

- 401 [39] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, et al. Chain-of-
402 thought prompting elicits reasoning in large language models. *Advances in Neural Information*
403 *Processing Systems*, 35:24824–24837, 2022.
- 404 [40] N. Weir and B. Van Durme. Dynamic generation of interpretable inference rules in a neuro-
405 symbolic expert system. *arXiv preprint arXiv:2209.07662*, 2022.
- 406 [41] Y. Yang and L. Song. Learn to explain efficiently via neural logic inductive learning. *arXiv*
407 *preprint arXiv:1910.02481*, 2019.
- 408 [42] Z. Zhu, Y. Xue, X. Chen, D. Zhou, J. Tang, D. Schuurmans, and H. Dai. Large language models
409 can learn rules. *arXiv preprint arXiv:2310.07064*, 2023.

410 Limitations

411 **Limitation on inferential rule coverage.** Commonsense inferential rules may exist in diverse
412 formats and span various domains. Our work mainly focuses on rules formatted as *if-then* statements,
413 covering five domains: object affordance, accessibility, interaction, location and person’s need. In
414 future work, we will expand our scope to include inferential rules of other formats and explore
415 additional domains for broader coverage.

416 **Limitation on probing open-source models.** Our work does not probe and analyze open-source
417 models. While GPT-4 and GPT-3.5-turbo are considered as the most advanced models, open-source
418 counterparts may exhibit different behaviors or patterns in understanding inferential rules with varying
419 complexities. These aspects will be the subject of future exploration.

420 **Risk of environmental impact** A significant risk associated with our framework and analysis is
421 the potential increase in environmental burdens due to the extensive use of OpenAI’s APIs for LLMs.
422 This impact can be mitigated by replacing GPT-4 with future smaller-scale open-source models that
423 are more efficient with less environmental impact.

424 **Potential error in rule generation.** Generating inferential rules with specific requirements poses a
425 significant challenge. As the majority of our framework’s pipeline are powered by GPT-4, it may
426 inevitably generate inferential rules with logical inaccuracies even incorporating human verification.
427 This might result in less accurate probing of LLMs.

428 Ethical Consideration

429 All rules we collected through LLMs are released publicly for usage and its probing subset for
430 proficiency analysis have been subjected to a thorough review by the authors. The code of our
431 generation pipeline and probing experiments will also be publicly released. This setting guarantees
432 transparency and reproducibility in our experiments, allowing other researchers to evaluate and
433 expand upon our work. Our logic scaffolding framework is strictly limited to be used for rule
434 generation that follow the ethical guidelines of the community. The authors emphatically denounce
435 the use of our framework for generating inaccurate or harmful rules.

436 A Primitive Rule Generation Pipeline

437 A.1 Abstract Objects and Common Properties

438 Table 3 list 32 most common abstract objects and 18 common properties for primitive rule generation.

Table 3: List of pre-defined abstract objects and common properties.

Type	Words
Abstract Objects	“Person”, “Animal”, “Plant”, “Food”, “Alcohol”, “Disease”, “Drug”, “Natural Phenomenon”, “Condition”, “Material”, “Substance”, “Furniture”, “Publication”, “Organization”, “Authoriza- tion”, “Facility”, “Natural Place”, “Event”, “Show”, “Artwork”, “Job”, “Game”, “Vehicle”, “Tool”, “Technology”, “Electronic Device”, “Platform”, “Financial Product”, “Skill”, “Legisla- tion”, “Region”, “Time Period”
Common Properties	“Age”, “Price”, “Money”, “Height”, “Length”, “Weight”, “Strength”, “Size”, “Density”, “Volume”, “Temperature”, “Hardness”, “Speed”, “BoilingPoint”, “MeltingPoint”, “Frequency”, “Decibel”, “Space”

439

440 A.2 Rule Domains

441 Table 4 illustrates the detailed explanations, example predicates and rules across five domains.

Table 4: The explanations, example predicates and rules of five different domains.

Domain	Explanation	Predicates	Examples
Object Affordance	Whether a person can take an action over an object based on its property and requirement	CanDrive(Person X, Vehicle Y); CanCreate(Person X, Artwork Y); CanAttend(Person X, Event Y);	CanDrive(Person X, Vehicle Y):- Have(Person X, Age Z1), RequireMinimumAge(Vehicle Y, Age Z2), BiggerThan(Age Z1, Age Z2);
Object Accessibility	Whether an object can access the other object based on its physical condition, spatial and temporal restriction	CanAccess(Person X, Show Y); CanAccess(Animal X, Tool Y); CanAccess(Animal X, Animal Y);	CanAccess(Person X, Show Y):- LocatedIn(Person X, Region Z), BroadcastIn(Show Y, Region Z); CanNotAccess(Person X, Tool Y):- AllergicTo(Person X, Material Z), MadeOf(Tool Y, Material Z);
Object Interaction	How an object can interact with the other object based on their physical, spatial or temporal properties	CanSubmergeIn(Substance X, Substance Y); CanAdaptedFrom(Show X, Artwork Y); CanFitIn(Tool X, Tool Y);	CanSubmergeIn(Substance X, Substance Y):- DensityOf(Substance X, Density Z1), DensityOf(Substance Y, Density Z2), BiggerThan(Density Z1, Density Z2);
Object Location	The location description of an object	OriginatedFrom(Food X, Region Y); BannedIn(Drug X, Region Y); BornIn(Person X, Region Y);	OriginatedFrom(Food X, Region Y):- ProcessedIn(Food X, Facility Z), LocatedIn(Facility Z, Region Y);
Person's Need	Person need to take an action over objects under a specific circumstance	NeedToConsume(Person X, Drug Y); NeedToWater(Person X, Plant Y);	NeedToConsume(Person X, Drug Y):- Has(Person X, Disease Z), CanTreat(Drug Y, Disease Z);

442 A.3 Prompt for Conclusion Preparation

443 An example of the prompt for conclusion preparation about affordance is below.

Prompt for Conclusion Preparation

According to commonsense knowledge in reality, please list 5 predicates between the given two objects to describe the {object affordance}.

Examples:
 Object: Show, Artwork
 Predicate: CanBeAdaptedFrom(Show X, Artwork Y)

Object: {object₁}, {object₂}
Predicate:

445 A.4 Prompts for Premise Generation

446 For premise generation in each domain, we design an instruction followed by two demonstrations to
 447 iteratively prompt GPT-4, and the underlined sentence is the rule description which varies according
 448 to the specific domain, as shown in Table 5.

449 A.5 Grammatical Validity for Rule Filtering

450 As Figure 11, we check whether the variables in premises form a connected graph from node “X” to
 node “Y” to filter grammatically invalid rules.

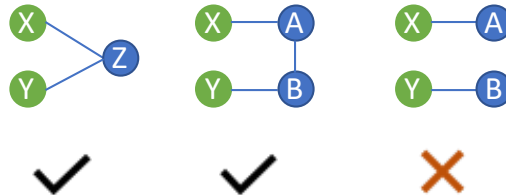


Figure 11: Grammatically valid and invalid rule graphs.

451

452 A.6 Prompts for Rule Filtering

453 Table 6 is an example prompt for rule filtering in object affordance domain.

Instruction for Premise Generation (Object Affordance)

According to commonsense knowledge in realistic scenarios, please generate 2 logical rules in both Prolog and natural language to describe the premises of the given conclusion. The rules in Prolog should have the same meaning with the rules in natural language.

Each rule should contain multiple premises and each premise should contain two variables in (X, Y, Z, Z1, Z2).

The rules should describe object affordance based on its property (such as height, age, price) and requirement (such as required skill, source, tool).

The premises should not contain negative words such as 'not', 'no', 'never' and 'un-'

Conclusion: {conclusion}

Rules:

Demonstrations for Premise Generation (Object Affordance)

Conclusion: CanCook(Person X, Food Y)

Rules:

1. CanCook(Person X, Food Y):- CanUse(Person X, Tool Z), UsedForCook(Tool Z, Food Y);
If Person X can use Tool Z which is used for cooking Food Y, then Person X can cook Food Y.
2. CanCook(Person X, Food Y):- Master(Person X, Skill Z), RequiredForCooking(Skill Z, Food Y);
If Person X has mastered Skill Z which is required for cooking Food Y, then Person X can cook Food Y.

Conclusion: CanDrive(Person X, Vehicle Y)

Rules:

1. CanDrive(Person X, Vehicle Y):- Have(Person X, Age Z1), RequireMinimumAge(Vehicle Y, Age Z2), BiggerThan(Age Z1, Age Z2);
If Person X has Age Z1 and the minimum age requirement for driving Vehicle Y is Age Z2, Age Z1 is bigger than Age Z2, then Person X can drive Vehicle Y.
2. CanDrive(Person X, Vehicle Y):- Obtain(Person X, Authorization Z), RequiredForDriving(Authorization Z, Vehicle Y);
If Person X have obtained a specific Authorization Z and Authorization Z is required for driving Vehicle Y, then Person X can drive Vehicle Y.

Table 5: Prompts for rule generation in different domains.

Domain	Rule Description
Object Affordance	The rules should describe object affordance based on its property (such as height, age, price) and requirement (such as required skill, source, tool).
Object Accessibility	The rules should describe object accessibility based on its physical condition, spatial and temporal restriction.
Object Interaction	The rules should describe object interaction based on its physical, spatial or temporal properties (such as speed, hardness, density, height, time period).
Object Location	The rules should describe the location information of an object.
Person's Need	The rules should describe person's need to take an action over the object.

454 A.7 Human Verification Templates and Rates

455 Before human verification, we first craft a qualification task to select AMT annotators from all
456 English-speaking countries (US, UK, New Zealand, Australia, Canada). The prospective workers are
457 presented with three representative test cases and need to predict whether the premise and conclusion
458 are clearly readable, and if the premise logically entails the conclusion. Only those workers correctly
459 passing all the test cases are recruited. The detailed template for human verification is shown as
460 Figure 12. This template is also used for getting human performance in rule probing analysis, wherein
461 a separate cohort of workers is qualified for manual rule probing. Besides, the overall rates of rule
462 acceptance in different domains during human verification are listed Table 7.

Prompt for Rule Filtering

True or False? Please predict whether the input rule is accurate or not according to commonsense knowledge in realistic scenarios, and also explain why.

Examples:

Input: If Person X has an Age Z1 and Vehicle Y requires an Age above Z2 for driving, with ...
 Output: True. Because Person X has achieved the ...

Input: If Person X was born in Season Z and Plant Y blooms in the same Season Z, then Person X can access Plant Y.
 Output: False. Because a person's birth season and a plant's blooming season has no logical connection.

Input: { candidate rule }
 Output:

Table 6: A prompt for rule filtering in object affordance.

Please read the following instructions and examples very carefully, and refer back to them while annotating:

Instructions (click to expand)

In this HIT you will be provided with a **Premise** and a **Conclusion**.

- A **Premise** is a statement describing the **facts and relationships of multiple objects**.
 For example, "Person X is currently situated in Region Z, and Technology Y is prohibited in Region Z." is a **Premise**, where "Person X", "Technology Y" and "Region Z" are objects.
- A **Conclusion** is a statement **between two objects**, mainly describing the **ability of objects to do sth or the location** of objects.
 For example, "Person X can not employ Technology Y." is a **Conclusion**, where "Person X" and "Technology Y" are objects.

Your job is to answer **Yes or No** to **THREE** questions about the provided **Premise and Conclusion**. These three questions can be categorized into **TWO** types.

- Type I (Readable Expression):** Is the **Premise/Conclusion** a readable and clear expression?
 - Yes:** The Premise/Conclusion is a *readable expression* and has a *clear meaning without ambiguity*. For example, "Person X can not employ Technology Y." is a readable and clear expression.
 - No:** The Premise/Conclusion is an *unreadable expression* or has an *unclear meaning*. For example, "Technology Y is determined in Region Z." is not a readable and clear statement. Instead, "Technology Y is deployed/invented in Region Z." is a readable and clear expression.
- Type II (Logically Correct):** Is the statement "If the **Premise** happens, then the **Conclusion** will happen as well." logically correct, with very high probability?
 - Yes:** If the Premise happens, then the Conclusion is *very likely* to happen.
 - No:** The Premise will be *unlikely* to lead to the Conclusion. The given Premise is *insufficient or irrelevant or contradictory* to determine the Conclusion.
 - !!!!!! Note that we only focus on the direct logical connection from premise to conclusion, without considering other potential situations.** For example, the statement "If Person X is over 18 years old, then Person X can drive the car." is **considered logically correct, since "can drive" here means "has the ability to drive" and we do not factor in whether Person X have a car.**

Premise:

Conclusion:

Please read the **Premise** and **Conclusion** carefully, and answer the questions below.

Question 1: Is the **Premise** a readable and clear expression?
 No Yes
 The Premise is an *unreadable expression* or has an *unclear meaning*.

Question 2: Is the **Conclusion** a readable and clear expression?
 No Yes
 The Conclusion is an *unreadable expression* or has an *unclear meaning*.

Question 3: Is the statement "If the **Premise** happens, then the **Conclusion** will happen as well?" logically correct, with very high probability?
 !!!!!!! **Note that we only focus on the direct logical connection from premise to conclusion, without considering other potential situations.** For example, the statement "If Person X is over 18 years old, then Person X can drive the car." is **considered logically correct, since "can drive" here means "has the ability to drive" and we do not factor in whether Person X have a car.**
 No Yes
 The Premise will be *unlikely* to lead to the Conclusion. The given Premise is *insufficient or irrelevant or contradictory* to determine the Conclusion.

(Optional) Please let us know if anything was unclear, if you experienced any issues, or if you have any other feedback for us.

Submit

Figure 12: AMT template for human verification of primitive rules.

Table 7: The rule yield rates (%) of human verification.

	Affordance	Accessibility	Interaction	Location	Person's Need
Yield Rate	48.09	37.28	52.81	53.74	49.45

463 **A.8 Statistics of ULogic**

464 We construct an inferential rule base ULogic comprising 14,647 rules, with 7,967 primitive and
 465 6,680 compositional ones. These rules span five key domains: object affordance, accessibility,
 466 interaction, location and person's need. They vary in compositional depth from 0 to 3, with rule
 467 lengths ranging from 1 to 6. Detailed statistics are in Table 8.

Table 8: Statistics of constructed rule base.

Domain	Affordance	Accessibility	Interaction	Location	Need	Total
Primitive rules						7,967
Single-fact	328	513	440	194	87	1,562
Multi-fact	387	638	2,527	166	128	3,846
Intermediate	417	590	1,286	165	101	2,559
Compositional rules						6,680
Compositionality=1	322	675	936	111	91	2,135
Compositionality=2	199	773	744	100	136	1,952
Compositionality=3	229	1052	896	217	199	2,593

468 **B Rule Probing**

469 **B.1 Rule Probing Templates**

Table 9 lists five different templates for unbiased rule probing.

Table 9: Five templates for rule probing.

	Template	Label
1	True or False? Please predict whether the input rule is very likely to be true.	True/False
2	Right or Wrong? Please predict whether the input rule is valid and correct.	Right/Wrong
3	Yes or No? Please predict whether the premise entails the conclusion.	Yes/No
4	Premise:..., Conclusion:... Does premise entail conclusion? Please answer Yes or No.	Yes/No
5	Given the observations ..., can we draw the conclusion ...? Please answer Yes or No.	Yes/No

470

471 **B.2 Dual-side Rule Probing Setting**

Table 10 illustrate a concrete example of dual-side rule probing.

Table 10: A specific example of dual-side rule probing.

<i>If Premise, then Conclusion_original.</i>	True/Right/Yes
<i>If Premise, then Conclusion_flipped.</i>	False/Wrong/No
Example	
<i>If Person X is allergic to Substance Z and Food Y contains Substance Z, then Person X cannot eat Food Y.</i>	True/Right/Yes
<i>If Person X is allergic to Substance Z and Food Y contains Substance Z, then Person X can eat Food Y.</i>	False/Wrong/No

472

473 **B.3 Rule Depths Probing**

474 The analysis of GPT-series LLMs and human on different compositional depths is presented as
 475 Figure 13. “Depth=0” represents primitive rules and “Depth=1,2,3” denote compositional rules
 476 involving 1 to 3 backward chaining steps.

477 **B.4 Illustrations of Rule Structures**

478 Figure 14 displays several examples showcasing both symbolic and verbalized rules across different
 479 structure types.

480 **B.5 Different CoT Prompts**

481 Table 11 lists different prompts of three CoT strategies for rule probing.

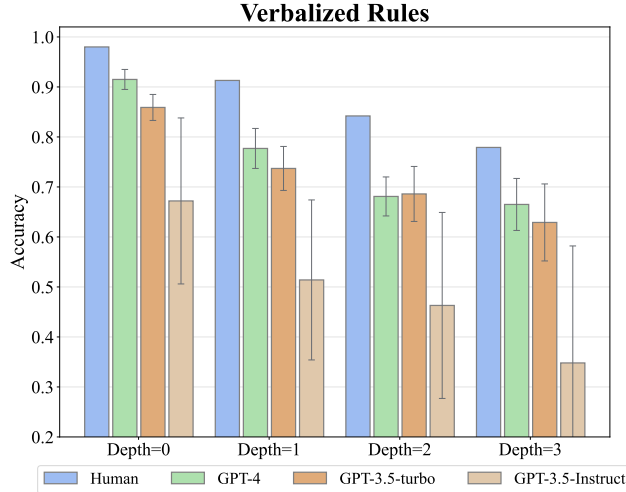


Figure 13: Probing results of varied depths.

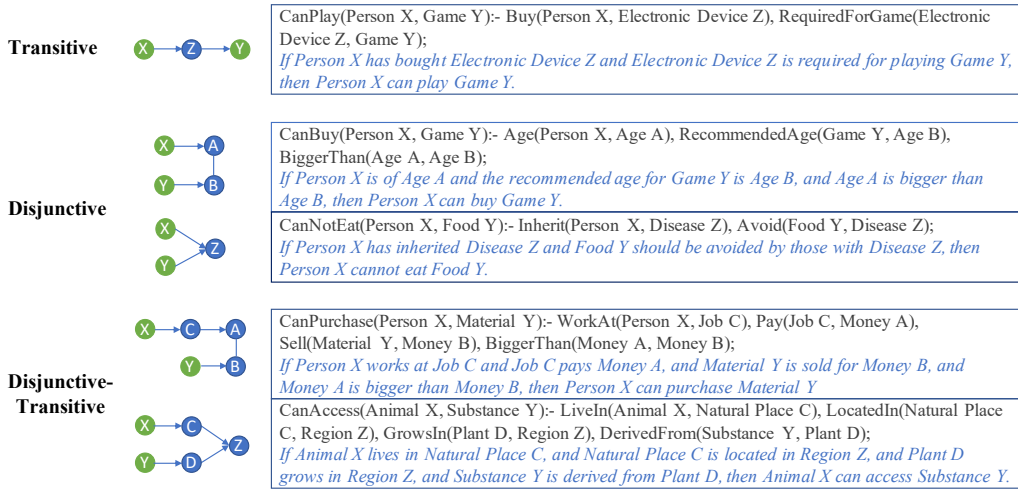


Figure 14: Example rules of different structures.

Table 11: Three CoT prompts for rule probing.

CoT strategy	Prompt
Answer-Explain	True or False? Please predict whether the input rule is very likely to be true, and also explain why.
Answer-Explain	True or False? Please predict whether the input rule is very likely to be true. Please first briefly explain your thought process in one sentence, and then give your answer.
Self-Consistency	True or False? Please predict whether the input rule is very likely to be true. Please first generate three different sentences to respectively explain your three thought processes briefly, and then based on the corresponding thought to give your answer. Finally, output the final answer according to majority voting.

482 B.6 Necessary Bias

483 As mentioned in Section 3.2, We investigate why GPT-4 significantly underperforms GPT-3.5-Turbo
 484 on transitive rules. Transitive rules typically describe a straightforward logical chain from variable X
 485 to Y, where GPT-3.5-Turbo and GPT-3.5-Instruct manage with greater ease. In contrast, we find that

486 GPT-4 exhibits a “necessary bias” that tend to consider all necessary conditions reaching a conclusion,
 487 which avoids drawing a definite judgement. This conservative response is more pronounced in
 488 transitive rules, where GPT-4 more frequently responds with hesitations like “it does not necessarily
 489 mean” in its explanations. We present a probing example of a transitive rule by GPT-4 in Table 12,
 490 including its generated prediction, explanation and the corresponding label. We hypothesize that this
 491 conservative style may come from LLMs’ preference alignment during Reinforcement Learning with
 492 Human Feedback [19].

Table 12: An example probing result of a transitive rules by GPT-4.

Input	True or False? Please predict whether the input rule is very likely to be true, and also explain why. Rule: If Person X has a job A that requires Skill Z, and Skill Z is required to cook Food Y, then Person X can cook Food Y.
Prediction	False
Explanation	Just because Person X has a job that requires Skill Z, it does not necessarily mean they can cook Food Y. They may have the skill, but cooking also requires other knowledge of recipes, techniques.
Label	True

Prompt for Conclusion Generation	
<p>Given the premise, please generate its conclusion between X and Y in both Prolog and natural language. The conclusion in Prolog should have the same meaning with the conclusion in natural language. Each conclusion should contain only two variables X and Y without mentioning other variables, like A, B, C, Z.</p> <p>### Examples:</p> <p>Premise: If Person X is allergic to Material Z and Furniture Y is made from Material Z. Conclusion: [Prolog]: CanNotHold(Person X, Furniture Y); [Natural Language]: Person X cannot hold Furniture Y.</p> <p>Premise: If Substance X has a Density Z1, the density of Substance Y is Density Z2, and Density Z1 is bigger than Density Z2. Conclusion: [Prolog]: CanSubmerge(Substance X, Substance Y); [Natural Language]: Substance X can submerge in Substance Y.</p> <p>Premise: {premise} Conclusion:</p>	

Table 13: Prompt ChatGPT and GPT-4 for conclusion generation.

493 C Inference Engine

494 C.1 Illustration of Instruction Tuning

495 Figure 15 illustrate the pipeline of instruction tuning for rule distillation as an inference engine. Our
 496 inference engine is trained for three tasks: conclusion generation, premise completion and premise
 497 generation. The conclusion generation focuses on creating a conclusion from a provided premise. For
 498 premise completion, given a conclusion and its partial premise, the inference engine must complete
 499 the remaining premise part to support the conclusion. In premise generation, the engine is tasked
 500 with creating premises of varying complexity based on a given conclusion, specifically generating
 501 premises with one, two or even more facts. We also provide an inference engine demo for flexible
 502 rule generation as shown in Figure 16.

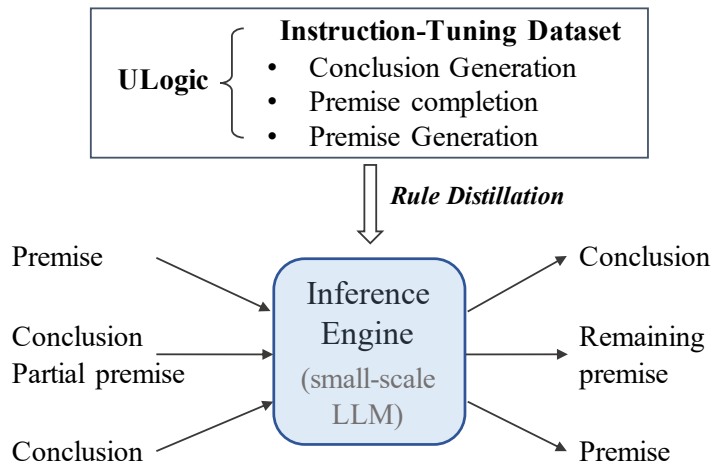


Figure 15: Rule distillation for inference engine.

503 C.2 Implementation Details

504 We fine-tune Mistral-7b with our constructed instruction dataset with Quantization LoRA (QLoRA)
 505 method [12, 6] as our inference engine. We set the learning rate to 7×10^{-5} , batch size to 8, gradient
 506 accumulation step to 16, and train the model 2 epochs. We apply QLoRA to all the linear layers of
 507 the model, including q_proj, k_proj, v_proj, o_proj, gate_proj, up_proj, down_proj, and lm_head.
 508 The α and r of the QLoRA method are both set to 16.

509 C.3 Prompting ChatGPT and GPT-4 for Three Tasks

510 As Step-2 of Sec. 2.2, we utilize two-shot prompts to instruct ChatGPT and GPT-4 in simultaneously
 511 generating symbolic and verbalized responses for three tasks, as shown in Table 13, 14, 15.

512 C.4 Evaluation Metrics

513 We detailed describe the metrics for evaluating our inference engine against ChatGPT and GPT-4 for
 514 the premise generation task.

- 515 • Accuracy: The premise is logically correct to infer the conclusion and follow the instruction
 516 regarding the specific number of facts.
- 517 • Diversity: The degree of variation among the three generated rules.
- 518 • Complexity: Assessed only for premise generation with more than 2 facts, measuring the fact
 519 number and the semantic difficulty.
- 520 • Abstractness: The variable types in premises are abstract to generalize to diverse instances. For
 521 example, the variable types “Region” and “Event” are abstract while “New York” and “The FIFA
 522 World Cup” are specific entities with low abstractness.

523 C.5 LLM Evaluation Prompts

524 We prompt GPT-4 as the evaluator for rating the accuracy of conclusion generation and premise
 525 completion tasks, and the accuracy, diversity, complexity and abstractness of the premise generation
 526 task. We adopt one-shot prompts which are shown as Table 16 and Table 17 (with demonstrations
 527 omitted).

528 C.6 Human Evaluation Templates

529 For the human evaluation of premise generation accuracy, we qualify a new cohort of AMT annotators
 530 to pairwise compare two sets of generated premises in terms of logical consistency with the provided

Inferential Rule Distillation Demo 🤖

As an Inference Engine, I can conduct the following inferences on Commonsense Rules ❤️

Conclusion Generation ▾

Premise Completion ▾

Premise Generation ▾

Instruction

Please select an inference type

Given the conclusion, please generate its premise. ▾

Please select your preferred fact number in the premise

2 ▾

Input

Premise:

If Person X demands sleep time A and Person X has sleep time B, and B is longer than A.

Conclusion:

Chairman X can not drive Vehicle Y.

Submit

Output

- 1: Premise: If Chairman X lacks License Z and Vehicle Y requires License Z.
- 2: Premise: If Chairman X lacks License Z2 and Vehicle Y requires License Z2.
- 3: Premise: If Chairman X is banned from Region Z and Vehicle Y is registered in Region Z.

Figure 16: Inference Engine Demo.

531 conclusion, alignment with common sense and the inclusion of an accurate number of facts. The
532 detailed template for human evaluation is shown as Figure 17.

533 C.7 Downstream Reasoning Datasets

534 StrategyQA and SOCIAL IQA consist of crowd-sourced questions involving reasoning of implicit
535 logic. LINK comprises GPT-4 generated statements instantiated from abstract rules, including two
536 subsets: head distribution statements and long-tail knowledge statements. PIQA examines operational
537 commonsense for achieving physical goals and CSQA2.0 features adversarial commonsense examples
538 designed to mislead AI systems.

Prompt for Premise Completion

Given the conclusion and a part of its premise, please complete the remaining portion of the premise in both Prolog and natural language.
 The remaining premise in Prolog should have the same meaning with the remaining premise in natural language.
 Each fact in the remaining premise should contain two variables, like X, Y, Z, Z1, Z2, A, B.

Examples:

Conclusion: Person X cannot use Furniture Y.
 Partial Premise: If Person X is allergic to Material Z,
 Remaining Premise:
 [Prolog]: MadeFrom(Furniture Y, Material Z);
 [Natural Language]: Furniture Y is made from Material Z.

Conclusion: Substance X can submerge in Substance Y.
 Partial Premise: If Substance X has a Density Z1, the density of Substance Y is Density Z2,
 Remaining Premise:
 [Prolog]: BiggerThan(Density Z1, Density Z2);
 [Natural Language]: Density Z1 is bigger than Density Z2.

Conclusion: {conclusion}
 Partial Premise: {partial premise}
 Remaining Premise:

Table 14: Prompt ChatGPT and GPT-4 for premise completion.

Please read the following instructions and Examples very carefully, and refer back to them while annotating:

Instructions (click to expand)

In this HIT you will be provided with a *conclusion* and two groups of its candidate premises, along with the *specified number of facts* in the premises.

- A *conclusion* is a statement that typically involves two objects. It usually describes the objects' abilities, locations, or needs.
 - For example, "Person X can not employ Technology Y." is a conclusion, where "Person X" and "Technology Y" are two objects.
- A *Premise* is a statement describing facts about multiple objects, aiming to provide evidence supporting the conclusion.
 - For example, "Person X is situated in Region Z, and Technology Y is prohibited in Region Z." is a plausible *premise* for above mentioned conclusion.
- The *specified number of facts* refers the number of facts that each candidate premise should comprise. Each fact should involve two objects
 - For example, the above premise "Person X is situated in Region Z, and Technology Y is prohibited in Region Z." contains 2 facts.

Your job is to **compare** two groups of candidate premises, and **determine which group is more accurate** to reach the conclusion with specified number of facts. **When assessing accuracy, please consider the following three criteria:**

- Logical Consistency:**The premises in this group are logically correct to lead to the given conclusion.
- Common Sense Alignment:** The premises in this group align well with common sense.
- Fact Count Accuracy:** The premises in this group precisely contain the specified number of facts – no more, no less.

Please choose one of the following three options: A, B, or Tie (cannot determine).

Examples (click to expand)

Conclusion:

Specified Number of Facts:

Premises:

Group A:

- If Person X has Age Z1 and the minimum age requirement for driving Vehicle Y is Age Z2, and Age Z1 is smaller than Age Z2.
- If Person X has a height of Height Z1 and the minimum height requirement for driving Vehicle Y is Height Z2, and Height Z1 is smaller than Height Z2.
- If Person X is under the age of Z1 and Vehicle Y is manufactured by Organization A, which has set the age limit for driving the vehicle at Z2, and Age Z2 is greater than Age Z1.

Group B:

- If Person X is of age Z1 and the minimum driving age for Vehicle Y is Z2, and Z1 is smaller than Z2.
- If Person X has a license of type Z1 and Vehicle Y requires a license of type Z2, and Z1 does not match Z2.
- If Person X has a medical condition Z and Vehicle Y is prohibited for individuals with medical condition Z.

Question: Overall, which group of premises are more accurate to support the conclusion with correct number of facts and alignment with common sense?

- Answer:** A
- Why?** Because the premises in Group A and Group B can both accurately lead to the conclusion "Person X can not drive Vehicle Y" and make sense logically. The issue lies with Group B's third premise, which contains only 2 facts inconsistent with our specification of "more than 2 facts".

Figure 17: AMT template for human evaluation for premise generation accuracy.

Prompt for Premise Generation

Given the conclusion, please generate three different premises in both Prolog and natural language, ensuring that each Prolog premise conveys the same meaning as its natural language counterpart. Each premise should contain a specified number of facts, with each fact comprising only two variables, such as X, Y, Z, Z1, Z2, A, B.

Examples:

Fact number: 1 fact

Conclusion: Person X has Skill Y.

Three Premises:

1. [Prolog] Learned(Person X, Skill Y); [Natural Language] If Person X learned Skill Y.
2. [Prolog] Inherit(Person X, Skill Y); [Natural Language] If Person X inherits Skill Y.
3. [Prolog] Acquire(Person X, Skill Y); [Natural Language] If Person X acquires Skill Y.

Fact number: more than 2 facts

Conclusion: Person X cannot attend Event Y.

Three Premises:

1. [Prolog] Have(Person X, Age Z1), RequireMinimumAge(Event Y, Age Z2), BiggerThan(Age Z2, Age Z1); [Natural Language] If Person X has Age Z1 and the minimum age requirement for attending Event Y is Age Z2, Age Z2 is bigger than Age Z1.
2. [Prolog] Have(Person X, Height Z1), RequireAbove(Event Y, Height Z2), SmallerThan(Height Z1, Height Z2); [Natural Language] If Person X has a Height Z1, and Event Y requires a Height above Z2, and Height Z1 is smaller than Height Z2.
3. [Prolog] HaveCriminalRecord(Person X, Event Z), ProhibitedBy(Event Z, Legislation A), EnforcedIn(Legislation A, Region B), HeldIn(Event Y, Region B); [Natural Language] If Person X has a criminal record for Event Z and Event Z is prohibited by Legislation A, which is enforced in Region B, and Event Y is held in Region B.

Fact number: {fact num}

Conclusion: {conclusion}

Three Premises:

Table 15: Prompt ChatGPT and GPT-4 for premise generation.

Prompt for Rating the Accuracy of Conclusion Generation

You are a helpful scoring assistant.

Please read the provided premise carefully, and rate the accuracy of the candidate conclusion on a scale of 1 to 3:

- 1 (not accurate): The conclusion is clearly unsupported, irrelevant or contradictory to the provided premise.
- 2 (somewhat accurate): The conclusion, despite being supported by the premise, fails to state the definitive link between X and Y, or contradicts common sense, or lacks clarity.
- 3 (highly accurate): The conclusion correctly states the definitive link between X and Y, and is well-supported by the premise aligning with both established facts and common sense.

Please first output your rating based on your general knowledge and logical reasoning, and then provide a brief explanation with no more than 100 words.

[Provided Premise]: {premise}

[Candidate Conclusion]: {conclusion}

[Output]:

Prompt for Rating the Accuracy of Premise Completion

You are a helpful scoring assistant.

Please read the provided conclusion and its partial premise carefully, and rate the accuracy of its remaining premise in completing the provided premise to reach the conclusion, using a scale from 1 to 3:

- 1 (not accurate): The remaining premise fails to complete the provided premise for deducing the conclusion. It may be irrelevant or inconsistent with the provided premise or conclusion, or both.
- 2 (somewhat accurate): The remaining premise can somewhat supplement the provided premise but is not entirely sufficient for a conclusion inference. It may require additional information for comprehensive completion, or contradicts common sense, or lacks clarity.
- 3 (highly accurate): The remaining premise, combined with the provided partial premise, can correctly lead to the given conclusion, and also aligns well with common sense.

Please first output your rating based on your general knowledge and logical reasoning, and then provide a brief explanation with no more than 100 words.

[Conclusion]: {conclusion}

[Partial Premise]: {partial premise}

[Remaining Premise]: {rest premise}

[Output]:

Prompt for Rating the Accuracy of Premise Generation

You are a helpful scoring assistant.

Please carefully read the provided conclusion along with the specified number of facts, and rate the accuracy of candidate premise in both reaching the conclusion and containing the correct number of facts, using a scale from 1 to 3:

- 1 (not accurate): The premise is logically incorrect, irrelevant or contradictory for deducing the conclusion, or it contains an incorrect number of facts.
- 2 (somewhat accurate): The premise can partially infer the conclusion but is not entirely sufficient. It may require additional information, or contradicts common sense, or lacks clarity.
- 3 (highly accurate): The premise can correctly lead to the given conclusion and aligns well with common sense, and precisely contains the specified number of facts.

Please first output your rating based on your general knowledge and logical reasoning, and then provide a brief explanation with no more than 100 words.

[Fact Number]: {fact num}

[Conclusion]: {conclusion}

[Premise]: {premise}

[Output]:

Table 16: Prompts for rating the accuracy of three tasks.

Prompt for Rating the Diversity of Premise Generation

You are a helpful scoring assistant.

Please read the provided conclusion and multiple generated premises carefully, and rate the diversity of these premises using a scale from 1 to 3:

- 1 (low diversity): The premises show minimal variation, where all three premises largely repeat same perspectives with slight lexical changes.
- 2 (moderate diversity): The premises exhibit some degree of variation, with two out of the three premises sharing similar perspectives, expressions and fact numbers while the third presents different content.
- 3 (high diversity): The premises display a high level of diversity, where each premise presents distinct perspective from the others, or contains different fact numbers.

Please first output your rating, and then provide a brief explanation with no more than 50 words.

[Conclusion]: {conclusion}

[Premise]: {premise₁}, {premise₂}, {premise₃}

[Output]:

Prompt for Rating the Complexity of Premise Generation

You are a helpful scoring assistant.

Please carefully read the provided conclusion, and rate the complexity of candidate premise considering both the number of facts it comprises and its semantic difficulty, using a scale from 1 to 3:

- 1 (low complexity): The premise is straightforward, incorporating no more than 3 facts with clear and easy-to-understand semantics and a simple logical structure.
- 2 (moderate complexity): The premise exhibits moderate complexity, which involves 4 facts and somewhat intricate semantics and a logical structure that require some thought to understand.
- 3 (high complexity): The premise is highly complex with more than 4 facts, which also includes complex semantics and an abstract logical structure, demanding a high level of understanding.

Please first output your rating based on your general knowledge and logical reasoning, and then provide a brief explanation with no more than 50 words.

[Conclusion]: {conclusion}

[Premise]: {premise}

[Output]:

Prompt for Rating the Abstractness of Premise Generation

You are a helpful scoring assistant.

Please carefully read the provided conclusion, and rate the abstractness of objects in the candidate premise considering how broadly they can generalize to various specific instances, using a scale from 1 to 3:

- 1 (low abstractness): The objects in the premise are concrete and specific, making direct and clear reference to particular instances or examples, which focus on specific people, places, or tangible entities, such as Swimmer, New York, or SUV.
- 2 (moderate abstractness): The objects in the premise are somewhat abstract, representing a balance between specific instances and general concepts. They may pertain to fine-grained categories of people, places, or things, such as Professionals, City, or Car.
- 3 (high abstractness): The objects in the premise are highly abstract, focusing on coarse-grained people, places or things that are far removed from concrete instances, such as Person, Region, or Event, or general properties like Age and Height.

Please first output your rating based on your general knowledge and logical reasoning, and then provide a brief explanation with no more than 50 words.

[Conclusion]: {conclusion}

[Premise]: {premise}

[Output]:

Table 17: Prompts for rating the diversity, complexity and abstractness of premise generation.