# Liar, Liar, Logical Mire: A Benchmark for Suppositional Reasoning in Large Language Models

**Philipp Mondorf**[1,2] and **Barbara Plank**[1,2]

[1]MaiNLP, Center for Information and Language Processing, LMU Munich, Germany
[2]Munich Center for Machine Learning (MCML), Munich, Germany
p.mondorf@lmu.de     b.plank@lmu.de

## Abstract

Knights and knaves problems represent a classic genre of logical puzzles where characters either tell the truth or lie. The objective is to logically deduce each character's identity based on their statements. The challenge arises from the truth-telling or lying behavior, which influences the logical implications of each statement. Solving these puzzles requires not only direct deductions from individual statements, but the ability to assess the truthfulness of statements by reasoning through various hypothetical scenarios. As such, knights and knaves puzzles serve as compelling examples of suppositional reasoning. In this paper, we introduce *TruthQuest*, a benchmark for suppositional reasoning based on the principles of knights and knaves puzzles. Our benchmark presents problems of varying complexity, considering both the number of characters and the types of logical statements involved. Evaluations on *TruthQuest* show that large language models like Llama 3 and Mixtral-8x7B exhibit significant difficulties solving these tasks. A detailed error analysis of the models' output reveals that lower-performing models exhibit a diverse range of reasoning errors, frequently failing to grasp the concept of truth and lies. In comparison, more proficient models primarily struggle with accurately inferring the logical implications of potentially false statements.

## 1 Introduction

Well-designed logic puzzles can serve as a valuable tool for gaining deeper insights into the capabilities of large language models (LLMs) (Giadikiaroglou et al., 2024; Li et al., 2024; Del and Fishel, 2023). By challenging models to navigate sophisticated logic problems, these puzzles can reveal how LLMs identify patterns, recognize relationships, and employ logical principles (Tong et al., 2023; Ding et al., 2024). In his book *"What is the Name of This Book?"*, Smullyan (1978) introduced a series of
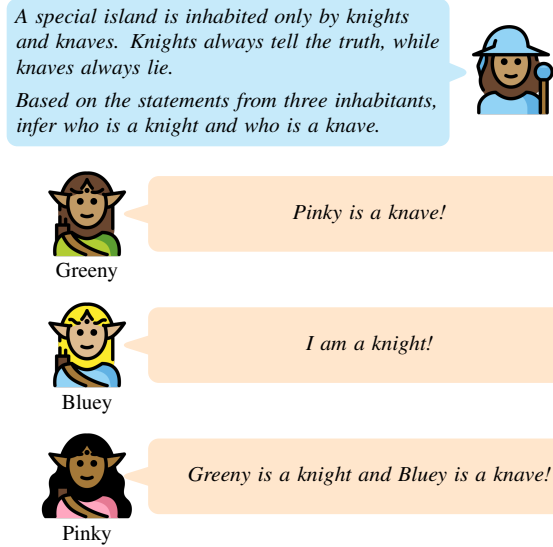


Figure 1: An instance of the *knights & knaves* puzzle. By reasoning about the characters' statements and their truthfulness, it is possible to deduce that Greeny and Bluey must be knights, while Pinky is a knave.

*knights and knaves* puzzles, where characters are either knights who always tell the truth or knaves who always lie.[1] The goal is to deduce the identity of each character based on their statements (see Figure 1). Unlike other deductive reasoning tasks, where premises are typically assumed to be true (Han et al., 2024; Dalvi et al., 2021; Clark et al., 2021), these puzzles require the reasoner to assess the truthfulness of statements by exploring different hypothetical scenarios. For instance, if the statement of *Pinky* in Figure 1 were true, *Greeny* must be a knight, thus telling the truth. However, *Greeny* states that *Pinky* is lying, which contradicts the initial truth assumption of *Pinky*'s statement. Hence, *Pinky* must be a knave. If *Pinky* is a knave, then *Greeny*'s statement is true, and thus *Greeny* will be

---

[1]Note: puzzles of this kind have existed before under different variations and names (Maurice, 1953; Goodman, 1972).

a knight. Based on *Pinky*'s false statement, it then follows that *Bluey* must also be a knight. This form of suppositional reasoning, i.e. the ability to reason conditionally, is essential in scenarios where the logical ramifications of different possibilities need to be considered, such as in planning or everyday reasoning (Byrne and Handley, 1997).

In this paper, we introduce *TruthQuest*,[2] a benchmark designed to evaluate the suppositional reasoning capabilities of large language models through knights and knaves puzzles. We present 2,400 problems of varying complexity, depending on the number of characters and types of logical statements involved (see Section 3). We assess the reasoning behavior of three model families: Llama 2 (Touvron et al., 2023), Llama 3 (Meta AI, 2024), and Mixtral-8x7B (Mistral AI, 2023). In addition to evaluating the models' task performance, we conduct an in-depth analysis of their outputs to gain insights into the types of errors encountered during reasoning. This is done through both comprehensive human inspection and AI-assisted evaluation. Our findings reveal that:

- All models exhibit significant difficulties in solving knights and knaves problems.

- Although more advanced prompting techniques, such as chain-of-thought prompting (Wei et al., 2022), enhance performance on simpler problems, accuracy declines markedly as puzzle complexity increases.

- The types of reasoning errors are closely linked to the models' performance. Lower-performing models exhibit a diverse range of reasoning errors, while more proficient models primarily struggle with deducing the correct logical implications of false statements.

## 2   Related Work

**Deductive Reasoning with LLMs.** Several studies evaluate LLMs on deductive reasoning tasks (Saparov and He, 2023; Dziri et al., 2023; Wan et al., 2024). In line with our research, some of these works employ logical puzzles to analyze the reasoning behaviors of LLMs (Ishay et al., 2023; Yao et al., 2023; Jiang et al., 2023). However, to the best of our knowledge, *TruthQuest* is the first deductive reasoning benchmark that evaluates the ability of LLMs to infer both the truthfulness of statements and their logical implications.

---

## 3   Dataset

Various versions of knights and knaves puzzles exist (Smullyan, 1978; Johnson-Laird and Byrne, 1990). In this work, we focus on the most popular variant, which features only two types of characters: knights, who always tell the truth, and knaves, who always lie, as illustrated in Figure 1. To construct valid instances of knights and knaves problems, we formalize the puzzle using a two-valued logic. Specifically, knights are assigned the truth value *true*, while knaves are mapped to *false*. For a given puzzle with $n$ characters, where $P$ denotes the truth value of a character and $Q$ is the character's logical claim, the puzzle can be expressed as a single conjunction using the bi-conditional operator:

$$\Phi = (P_1 \Leftrightarrow Q_1) \wedge (P_2 \Leftrightarrow Q_2) \\ \wedge \ldots \wedge (P_n \Leftrightarrow Q_n) \tag{1}$$

For instance, the example depicted in Figure 1 can be expressed as:

$$\Phi = (P_1 \Leftrightarrow \neg P_3) \wedge P_2 \wedge (P_3 \Leftrightarrow (P_1 \wedge \neg P_2))$$

where $P_1$, $P_2$, and $P_3$ correspond to the truth values of Greeny, Bluey, and Pinky, respectively. In this example, the characters' statements are given as follows: $Q_1 = \neg P_3$ ("Pinky is a knave!"), $Q_2 = P_2$ ("I am a knight!"), and $Q_3 = (P_1 \wedge \neg P_2)$ ("Greeny is a knight and Bluey is a knave!"). The bi-conditional operator ($P_i \Leftrightarrow Q_i$) reflects that a character is either lying or telling the truth. Specifically, it indicates that the statement $Q_i$ holds true if and only if the character is a knight (denoted by $P_i$). Conversely, if the character is a knave, $\neg P_i$, then the corresponding statement is false, $\neg Q_i$. To derive all $m$ possible solutions for such a puzzle, the expression $\Phi$ can be transformed into disjunctive normal form using Boolean algebra:

$$\Phi = (\psi_1^1 \wedge \ldots \wedge \psi_n^1) \vee (\psi_1^2 \wedge \ldots \wedge \psi_n^2) \\ \vee \ldots \vee (\psi_1^m \wedge \ldots \wedge \psi_n^m) \tag{2}$$

where $\psi_i^j \in \{P_i, \neg P_i\}$ denotes the character's identity as either knight ($P_i$) or knave ($\neg P_i$).

**Dataset Creation.**   For *TruthQuest*, we limit character statements to the types outlined in Table 1. To examine the impact of statement types on model behavior, we classify them into three distinct sets: $S$, $I$, and $E$, as specified in the table. For each set, separate datasets of knights

| Statement Types | Natural Language Example | Logic Expression | Set | | |
|---|---|---|---|---|---|
| Self-Referential | $P_i$: *I am a knight* | $P_i$ | | | |
| Accusation | $P_i$: *$P_j$ is a knight/knave* | $P_i \Leftrightarrow \psi_j \ (i \neq j)$ | S | | |
| Conjunction | $P_i$: *$P_j$ is a knight/knave and $P_k$ is a knight/knave* | $P_i \Leftrightarrow (\psi_j \wedge \psi_k) \ (i \neq j \neq k)$ | | I | |
| Implication | $P_i$: *If $P_j$ is a knight/knave, then $P_k$ is a knight/knave* | $P_i \Leftrightarrow (\psi_j \rightarrow \psi_k) \ (i \neq j \neq k)$ | | | E |
| Equivalence | $P_i$: *$P_j$ is a knight/knave if and only if $P_k$ is a knight/knave* | $P_i \Leftrightarrow (\psi_j \Leftrightarrow \psi_k) \ (i \neq j \neq k)$ | | | |

Table 1: Character statements in *TruthQuest*. Each type is represented by an example expressed both in natural language and boolean logic. The final column indicates the types of statements included in each statement set. For instance, $S$ is the only set that includes self-referential statements alongside accusations and conjunctions.

and knaves puzzles are generated. Specifically, instances are created by randomly sampling the statement of each character from the respective set, $Q_i \sim C \in \{S, I, E\}$. Each puzzle is solved by converting the problem (Equation 1) into disjunctive normal form, as shown in Equation 2. While Boolean algebra is used to construct the dataset, all statements and solutions are ultimately presented in natural language. Each dataset sample comprises the characters' statements and the corresponding solution, representing the logically valid identity of each character (as illustrated in Figure 1). For our benchmark, we include only instances that have a single, unique solution, i.e., $m = 1$. Furthermore, we consider varying numbers of characters for each statement set, specifically: $n = 3, 4, 5, 6$. This yields $3 \times 4 = 12$ data subsets. For each subset, 200 problems are generated, resulting in a total of 2,400 unique instances.

## 4 Experimental Setup

**Language Models.** We assess a total of six LLMs from three prominent open-access model families: Llama 2 (7B, 13B and 70B), Llama 3 (8B and 70B), and Mixtral-8x7B. The publicly accessible weights are obtained from the Hugging Face platform, specifically `Llama-2-chat-hf`,[3] `Meta-Llama-3-Instruct`,[3] and `Mixtral-8x7B-Instruct-v0.1`.[4] For further details about the models and prompts we employ, please refer to Appendix A.1.

**Evaluation Framework.** To assess the models' task performance, we follow a two-step approach. First, we use regular expressions to parse the models' final conclusions according to the format specified in the input prompt. If responses cannot be parsed in this way, they are subsequently passed to

an additional language model, specifically LLaMA-3-8B, which extracts the conclusion in the desired format (for the full evaluator prompt, see Figure 9 in the appendix). A schematic overview of this approach is presented in Figure 6 in the appendix. Once a conclusion is successfully parsed, the models' accuracy is determined by an exact match between their prediction and the correct puzzle solution. For a prediction to be considered accurate, the models must correctly deduce the identity of all characters.

Beyond assessing task performance, we analyze the models' reasoning errors. We manually inspect a subset of the responses from LLaMA-3-8B (zero-shot) and LLaMA-3-70B (zero-shot and four-shot chain-of-thought prompting). Specifically, we evaluate 10 responses from each of the 12 data subsets for each model and setup, totaling 360 responses. This involves parsing the model's conclusion and assessing its reasoning against six common error categories previously devised, as outlined in Table 4 in the appendix. This comprehensive manual evaluation is conducted independently by two hired students with expertise in data annotation. To assess the quality of the annotations, we report an overall Cohen's Kappa value of $\kappa = 0.70$. For a detailed description of the manual evaluation procedure and an overview of the inter-annotator agreement for each error type, please refer to Section C.1 in the appendix. To complement our manual evaluation and assess all model responses with respect to the error categories devised, we leverage GPT-4 (OpenAI et al., 2024) using few-shot prompting. For the complete prompt with all few-shot examples, see Figures 10 to 16 in the appendix.

**Meta-Evaluation.** We assess the quality of our evaluation procedures by comparing the results obtained via automatic evaluation with our manual assessment. Respective results are reported in Section 5 and Appendix C.2.

---

[3]huggingface.co/meta-llama
[4]huggingface.co/mistralai/Mixtral-8x77B-Instruct-v0.1

| Model | Mode | Set S | | | | Set I | | | | Set E | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | **3** | **4** | **5** | **6** | **3** | **4** | **5** | **6** | **3** | **4** | **5** | **6** |
| **Random Baseline** | - | 0.13 | 0.06 | 0.03 | 0.02 | 0.13 | 0.06 | 0.03 | 0.02 | 0.13 | 0.06 | 0.03 | 0.02 |
| **LLaMA-2-7b** | zero shot | 0.08 | 0.06 | 0.02 | 0.00 | 0.20 | 0.10 | 0.07 | 0.04 | 0.21 | 0.11 | 0.03 | 0.03 |
| **LLaMA-2-13b** | | 0.10 | 0.06 | 0.03 | 0.03 | 0.13 | 0.11 | 0.04 | 0.02 | 0.15 | 0.08 | 0.05 | 0.01 |
| **LLaMA-2-70b** | | 0.13 | 0.13 | 0.08 | 0.03 | 0.13 | 0.11 | 0.06 | 0.03 | 0.17 | 0.09 | 0.07 | 0.02 |
| **LLaMA-3-8B** | | 0.07 | 0.13 | 0.04 | 0.04 | 0.19 | **0.18** | 0.07 | 0.04 | 0.13 | 0.08 | 0.06 | 0.03 |
| **LLaMA-3-70B** | | **0.29** | **0.22** | **0.13** | **0.10** | **0.32** | 0.14 | **0.14** | **0.11** | **0.29** | **0.18** | **0.11** | **0.06** |
| **Mixtral-8x7B** | | 0.16 | 0.08 | 0.04 | 0.03 | 0.21 | 0.14 | 0.06 | 0.05 | 0.17 | 0.08 | 0.04 | 0.01 |
| **LLaMA-3-70B** | four shot | 0.22 | 0.25 | 0.19 | 0.13 | 0.24 | 0.21 | 0.13 | 0.10 | 0.32 | 0.22 | 0.11 | 0.05 |
| | eight shot | 0.22 | 0.21 | 0.16 | 0.09 | 0.32 | 0.25 | 0.07 | 0.09 | 0.27 | 0.20 | 0.10 | 0.02 |
| | zero CoT | 0.23 | 0.17 | 0.14 | 0.12 | 0.28 | 0.17 | 0.15 | 0.09 | 0.29 | 0.17 | **0.12** | 0.08 |
| | four CoT | 0.46 | **0.31** | **0.21** | 0.16 | **0.33** | **0.27** | 0.11 | **0.15** | **0.40** | **0.25** | 0.12 | **0.10** |
| | eight CoT | **0.60** | 0.26 | **0.21** | **0.20** | **0.33** | 0.20 | **0.15** | 0.12 | 0.37 | 0.20 | **0.12** | **0.10** |

Table 2: Accuracy values for different models and prompting techniques across each subset of *TruthQuest*. Results are grouped first by prompting technique and then by model. Bold values represent highest performance among a group. The random baseline indicates the accuracy achieved by guessing the identity of each character.
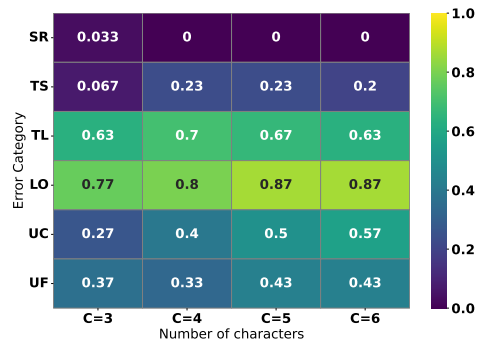
## 5 Results

**Task Performance**  Table 2 provides an overview of the models' task performance on *TruthQuest*. The table includes results for all models when prompted in a zero-shot setting, with additional results for LLaMA-3-70B using various prompting techniques (detailed results for other models can be found in Table 7 in the appendix). We observe that under zero-shot prompting, all models exhibit relatively poor performance across the different data subsets, often performing at or below chance level.

Although LLaMA-3-70B generally outperforms other models, its accuracy significantly declines as the number of characters—and consequently, the number of inference steps—increases. When guided via chain-of-thought prompting, LLaMA-3-70B shows performance improvements for problems involving fewer characters, particularly with statements sampled from set $S$. Other prompting techniques, such as few-shot prompting or zero-shot CoT (Kojima et al., 2022), do not substantially enhance LLaMA-3-70B's task performance.
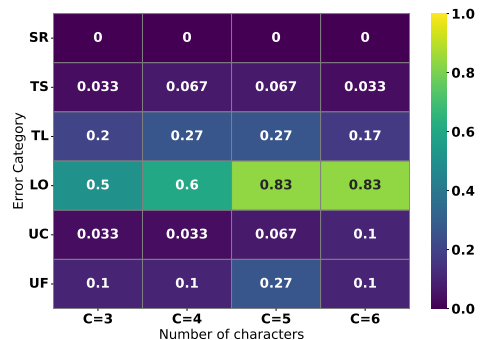
**Content Effects.**  For our analysis, we replace the terms *knights* and *knaves* with the pseudo-words *jabbas* and *tettes* to reduce the likelihood that models have been exposed to similar problems during training. Interestingly, we find that the choice of terms for *knights* and *knaves* seems to have no substantial impact on the models' performance, as shown in Figure 7 in the appendix.

**Error Analysis**  Figure 2 compares the relative occurrence of each error type, as outlined in Table 4, between LLaMA-3-8B (zero-shot) and LLaMA-3-70B (four-shot CoT). The values, derived from human annotations, are averaged across all statement sets for each number of characters. We find that LLaMA-3-8B exhibits a variety of errors, such



(a) LLaMA-3-8B (zero-shot)



(b) LLaMA-3-70B (four-shot CoT)

Figure 2: Relative occurrence of reasoning errors. Values are obtained from human-based evaluations. Error categories are abbreviated: *(SR) False statement reproduction, (TS) Assuming statements to be true, (TL) Misunderstanding the concept of truth and lies, (LO) Misunderstanding logical operators, (UC) Unjustified conclusion,* and *(UF) Unfaithfulness* (see Table 4).

as *misunderstanding the concept of truth and lies (TL)* and *unfaithfulness (UF)*. In contrast, LLaMA-3-70B predominantly struggles with deducing the logical implications of potentially false statements *(LO)*. This trend—where lower-performing models show a wider array of errors, while higher-performing models predominantly struggle with logical deductions from statements that may be false—is further supported by our complementary analysis using GPT-4, as illustrated in Figure 3. We find that the error distribution obtained through GPT-4 positively correlates with the distribution obtained via manual labeling. Pearson correlation coefficients and their corresponding p-values are provided in Table 6 in the appendix. Additional results for LLaMA-2-7B and LLaMA-3-70B with zero-shot prompting are presented in Figure 5 in the appendix. Further details of our automated error analysis can be found in Appendix B.2.2.



(a) LLaMA-3-8B (zero-shot)
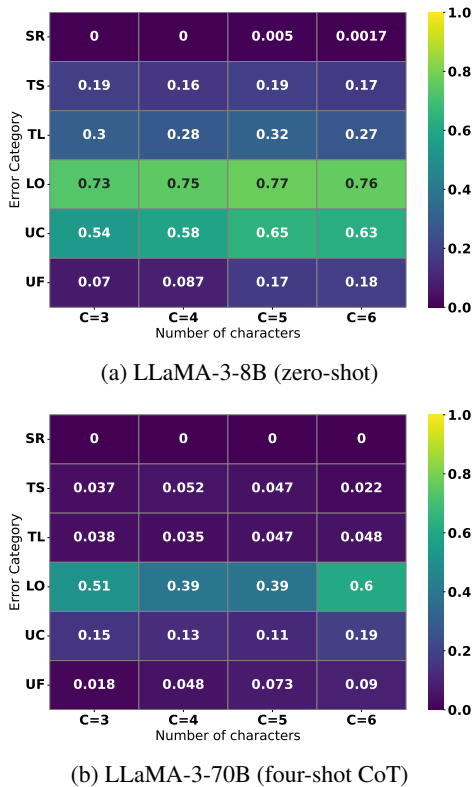


(b) LLaMA-3-70B (four-shot CoT)

Figure 3: Relative occurrences of the reasoning errors displayed by LLaMA-3-8B and LLaMA-3-70B when prompted via four-shot chain-of-thought. Values are obtained from GPT-4-based evaluations. Error categories are abbreviated for a more comprehensive overview: *(SR) False statement reproduction*, *(TS) Assuming statements to be true*, *(TL) Misunderstanding the concept of truth and lies*, *(LO) Misunderstanding logical operators*, *(UC) Unjustified conclusion*, and *(UF) Unfaithfulness*.

# 6 Conclusion

In this paper we introduce *TruthQuest*, a benchmark for suppositional reasoning based on *knights and knaves* puzzles. We demonstrate that LLMs exhibit significant difficulties solving these tasks. Our error analysis reveals that less proficient LLMs exhibit diverse errors, often failing to grasp the concept of truth and lies. In contrast, more proficient models primarily struggle with logical deductions from potentially false statements.

# 7 Limitations

While we introduce *TruthQuest*, a novel benchmark designed to evaluate the suppositional reasoning capabilities of large language models, several limitations remain that could be addressed in future research.

**Task Setup** Currently, *TruthQuest* includes only *knights and knaves* puzzles with a single, unique solution. Future work could expand this restriction to examine the impact of none or several solutions on model performance and behavior. Additionally, the benchmark is limited to simple propositional statements, as outlined in Table 1. Future iterations could incorporate more complex statement types that require more advanced inferences. Variations of *knights and knaves* puzzles, which consider additional characters or altered character attributes, also present opportunities for further exploration. For instance, Johnson-Laird and Byrne (1990) propose problems involving two types of persons: logicians, who always make valid deductions, and politicians, who never make valid deductions. An example problem states: "A says that either B is telling the truth or else is a politician (but not both). B says that A is lying. C deduces that B is a politician. Is C a logician?" Such variations represent compelling directions for future research.

**Evaluation Framework** Our manual evaluation framework is constrained by the number of annotators and the volume of annotations provided. Despite our efforts to optimize available resources, these constraints may impact the scalability and generalizability of our results. While our automatic evaluation procedure offers a promising alternative, we found that error annotations obtained through this method exhibit only fair overall agreement with human annotations at the instance level (see Section C.2 in the appendix for further details). Additionally, although we consider various prompting

techniques in our study, future research could explore the impact of more advanced methods, such as Tree-of-Thoughts (Yao et al., 2023) or Graph-of-Thoughts (Besta et al., 2024), on model performance.

## Acknowledgments

## References

Maciej Besta, Nils Blach, Ales Kubicek, Robert Gerstenberger, Michal Podstawski, Lukas Gianinazzi, Joanna Gajda, Tomasz Lehmann, Hubert Niewiadomski, Piotr Nyczyk, and Torsten Hoefler. 2024. Graph of thoughts: Solving elaborate problems with large language models. *Proceedings of the AAAI Conference on Artificial Intelligence*, 38(16):17682–17690.

Ruth M.J Byrne and Simon J Handley. 1997. Reasoning strategies for suppositional deductions. *Cognition*, 62(1):1–49.

Peter Clark, Oyvind Tafjord, and Kyle Richardson. 2021. Transformers as soft reasoners over language. In *Proceedings of the Twenty-Ninth International Joint Conference on Artificial Intelligence*, IJCAI'20.

Bhavana Dalvi, Peter Jansen, Oyvind Tafjord, Zhengnan Xie, Hannah Smith, Leighanna Pipatanangkura, and Peter Clark. 2021. Explaining answers with entailment trees. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 7358–7370, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

Maksym Del and Mark Fishel. 2023. True detective: A deep abductive reasoning benchmark undoable for GPT-3 and challenging for GPT-4. In *Proceedings of the 12th Joint Conference on Lexical and Computational Semantics (*SEM 2023)*, pages 314–322, Toronto, Canada. Association for Computational Linguistics.

Ruomeng Ding, Chaoyun Zhang, Lu Wang, Yong Xu, Minghua Ma, Wei Zhang, Si Qin, Saravan Rajmohan, Qingwei Lin, and Dongmei Zhang. 2024. Everything of thoughts: Defying the law of penrose triangle for thought generation. *Preprint*, arXiv:2311.04254.

Nouha Dziri, Ximing Lu, Melanie Sclar, Xiang (Lorraine) Li, Liwei Jiang, Bill Yuchen Lin, Sean Welleck, Peter West, Chandra Bhagavatula, Ronan Le Bras, Jena Hwang, Soumya Sanyal, Xiang Ren, Allyson Ettinger, Zaid Harchaoui, and Yejin Choi. 2023. Faith and fate: Limits of transformers on compositionality. In *Advances in Neural Information Processing Systems*, volume 36, pages 70293–70332. Curran Associates, Inc.

Panagiotis Giadikiaroglou, Maria Lymperaiou, Giorgos Filandrianos, and Giorgos Stamou. 2024. Puzzle solving using reasoning of large language models: A survey. *Preprint*, arXiv:2402.11291.

Nelson Goodman. 1972. *Problems and Projects*. Bobbs-Merrill, Indianapolis.

Simeng Han, Hailey Schoelkopf, Yilun Zhao, Zhenting Qi, Martin Riddell, Wenfei Zhou, James Coady, David Peng, Yujie Qiao, Luke Benson, Lucy Sun, Alex Wardle-Solano, Hannah Szabo, Ekaterina Zubova, Matthew Burtell, Jonathan Fan, Yixin Liu, Brian Wong, Malcolm Sailor, Ansong Ni, Linyong Nan, Jungo Kasai, Tao Yu, Rui Zhang, Alexander R. Fabbri, Wojciech Kryscinski, Semih Yavuz, Ye Liu, Xi Victoria Lin, Shafiq Joty, Yingbo Zhou, Caiming Xiong, Rex Ying, Arman Cohan, and Dragomir Radev. 2024. Folio: Natural language reasoning with first-order logic. *Preprint*, arXiv:2209.00840.

Adam Ishay, Zhun Yang, and Joohyung Lee. 2023. Leveraging large language models to generate answer set programs. In *Proceedings of the 20th International Conference on Principles of Knowledge Representation and Reasoning, KR 2023*, Proceedings of the International Conference on Knowledge Representation and Reasoning, pages 374–383. Association for the Advancement of Artificial Intelligence. Publisher Copyright: © 2023 Proceedings of the International Conference on Knowledge Representation and Reasoning. All rights reserved; 20th International Conference on Principles of Knowledge Representation and Reasoning, KR 2023 ; Conference date: 02-09-2023 Through 08-09-2023.

Yifan Jiang, Filip Ilievski, and Kaixin Ma. 2023. Brainteaser: Lateral thinking puzzles for large language models. In *Conference on Empirical Methods in Natural Language Processing*.

P.N. Johnson-Laird and Ruth M.J. Byrne. 1990. Metalogical problems: Knights, knaves, and rips. *Cognition*, 36(1):69–84.

Takeshi Kojima, Shixiang (Shane) Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large language models are zero-shot reasoners. In *Advances in Neural Information Processing Systems*, volume 35, pages 22199–22213. Curran Associates, Inc.

Yinghao Li, Haorui Wang, and Chao Zhang. 2024. Assessing logical puzzle solving in large language models: Insights from a minesweeper case study. *Preprint*, arXiv:2311.07387.

Kraitchik Maurice. 1953. Mathematical recreations.

Meta AI. 2024. Introducing meta llama 3: The most capable openly available llm to date. Accessed: 2024-06-01.

Mistral AI. 2023. Mixtral of experts: A high quality sparse mixture-of-experts. Accessed: 2024-06-01.

OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Łukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Łukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O'Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. 2024. Gpt-4 technical report. *Preprint*, arXiv:2303.08774.

Abulhair Saparov and He He. 2023. Language models are greedy reasoners: A systematic formal analysis of chain-of-thought. In *The Eleventh International Conference on Learning Representations*.

Raymond M. Smullyan. 1978. *What is the Name of This Book?: The Riddle of Dracula and Other Logical Puzzles*. Prentice-Hall, Englewood Cliffs, N.J.

Yongqi Tong, Yifan Wang, Dawei Li, Sizhe Wang, Zi Lin, Simeng Han, and Jingbo Shang. 2023. Eliminating reasoning via inferring with planning: A new framework to guide llms' non-linear thinking. *Preprint*, arXiv:2310.12342.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Mar-

tinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models. *Preprint*, arXiv:2307.09288.

Yuxuan Wan, Wenxuan Wang, Yiliu Yang, Youliang Yuan, Jen tse Huang, Pinjia He, Wenxiang Jiao, and Michael R. Lyu. 2024. A & b == b & a: Triggering logical reasoning failures in large language models. *arXiv preprint arXiv:2401.00757*.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, brian ichter, Fei Xia, Ed Chi, Quoc V Le, and Denny Zhou. 2022. Chain-of-thought prompting elicits reasoning in large language models. In *Advances in Neural Information Processing Systems*, volume 35, pages 24824–24837. Curran Associates, Inc.

Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. 2023. Tree of thoughts: Deliberate problem solving with large language models. In *Advances in Neural Information Processing Systems*, volume 36, pages 11809–11822. Curran Associates, Inc.

# A  Experimental Setup

In this section, we provide additional details about the experimental setup. First, we elaborate on the language models employed in this study. Subsequently, we provide a detailed description of each error category devised to assess the models' reasoning.

## A.1  Language Models

As outlined in Section 4, six distinct large language models from three open-access model families are evaluated in this study. Detailed information, including the number of parameters and the context length for each model, is provided in Table 3. Each model is prompted with a system message that offers context about the task setup and specifies the required response format. Following this, a user prompt containing the task description is given. The complete prompt is depicted in Figure 8. For few-shot setups, examples are presented in dialogue format, with the desired response indicated using the assistant's special tokens. Model responses are generated using nucleus sampling, utilizing the models' default values

as specified on the Huggingface Platform (top-$p = 0.9$, temperature $T = 0.6$).[5] All few-shot prompts as well as the code for prompting and evaluating models on *TruthQuest* are publicly accessible at https://github.com/mainlp/TruthQuest. Additionally, all model responses, along with human annotations regarding potential errors (see Section 4), are available at https://huggingface.co/datasets/mainlp/TruthQuest-Human-Annotations.

## A.2  Error Categorization

To gain a deeper understanding of the models' reasoning behavior, we develop six different error categories that encompass common errors observed in the models' reasoning. These categories are established through a preliminary manual examination of the models' responses. Detailed descriptions of each error category are provided in Table 4. It is important to note that these error categories are not meant to be exhaustive. Instead, they are intended to offer practical insights into the models' frequent failure modes. For clarity, an example for each error category is provided below:

**False statement reproduction (SR)**  Smaller, less advanced models, such as LLaMA-2-7B, often fail to accurately reproduce statements from the problem description. For example, a model may incorrectly reproduce Pinky's statement "Greeny is a knight and Bluey is a knave" as "Bluey states that both Pinky and Greeny are knaves."

**Assuming statements to be true (TS)**  Models sometimes overlook the possibility that statements may be false and instead assume the truth of each statement. For example, a model might directly conclude from Pinky's statement "Greeny is a knight and Bluey is a knave" that Greeny is indeed a knight and Bluey a knave, without considering the potential falsity of the statement.

**Misunderstanding the concept of truth and lies (TL)**  Some models demonstrate a lack of understanding regarding truth and falsehood. For instance, a model may incorrectly assume that knights lie or knaves tell the truth. Additionally, models sometimes erroneously assume that knights only make statements about other knights and knaves only about other knaves.

---

[5] Please refer to:  huggingface.co/meta-llama, and https://huggingface.co/mistralai/Mixtral-8x7B-Instruct-v0.1

| Model | Base Model | Parameters | Context Length | Tokens | GPU hours | Carbon Emitted | Fine-tuning |
|---|---|---|---|---|---|---|---|
| LLaMA-2-7B-Chat | LLaMA-2 | 7B | 4K tokens | 2.0T | 184K | 31 | SFT, RLHF |
| LLaMA-2-13B-Chat | LLaMA-2 | 13B | 4K tokens | 2.0T | 369K | 62 | SFT, RLHF |
| LLaMA-2-70B-Chat | LLaMA-2 | 70B | 4K tokens | 2.0T | 1.7M | 291 | SFT, RLHF |
| LLaMA-3-8B-Instruct | LLaMA-3 | 8B | 8K tokens | 15T+ | 1.3M | 390 | SFT, RLHF |
| LLaMA-3-70B-Instruct | LLaMA-3 | 70B | 8K tokens | 15T+ | 6.4M | 1900 | SFT, RLHF |
| Mixtral-8x7B-Instruct | Mixtral-8x7B | 46.7B | 32K tokens | - | - | - | SFT, DPO |

Table 3: Details about the models used in this study. Tokens refer to the number of tokens in the pre-training data. Similarly, the context length, GPU hours and carbon emissions relate to the base model's pre-training. Carbon emissions are reported as tCO2eq. We use the following abbreviations for fine-tuning: supervised fine-tuning (SFT), reinforcement learning with human feedback (RLHF), direct preference optimization (DPO). Information about Llama 2 is taken from Touvron et al. (2023), while properties of Llama 3 are reported by Meta AI (2024). For Mixtral-8x7B, we consider the blog post of Mistral AI (2023). Dashes denote unavailable information.

**Misunderstanding logical operators (LO)**
When processing false statements, such as Pinky's assertion that "Greeny is a knight and Bluey is a knave" (when Pinky is actually a knave), models often fail to infer the logical implications of the lie. Specifically, the model may not deduce that the possibilities are: (i) both Greeny and Bluey are knaves, (ii) both are knights, or (iii) Greeny is a knave and Bluey is a knight.

**Unjustified Conclusion (UC)** Models occasionally draw conclusions without providing a valid justification. For example, a model may erroneously assert that Greeny is a knave without offering any reasoning to support this conclusion.

**Unfaithfulness (UF)** Models are sometimes inconsistent in their reasoning. For instance, a model might first deduce that Greeny is a knight, but later contradict this by asserting that Greeny is a knave, without addressing the discrepancy in its reasoning.

## B  Additional Results

We report additional results of the models' performance on *TruthQuest* in Section B.1, and provide a supplementary analysis of the errors they commonly display in Section B.2.

### B.1  Task Performance

Table 7 supplements Table 2 by displaying LLaMA-3-8B's task performance for different prompting techniques. We observe that similar to LLaMA-3-70B, chain-of-thought prompting can yield notable performance gains for problems of lower com-plexity, i.e. fewer number of characters. Similarly, other prompting techniques such as few-shot prompting or zero-shot CoT do not seem to increase model performance.

### B.1.1  Content Effects

In conventional *knights and knaves* puzzles, *knights* are characters who always tell the truth, while *knaves* always lie. However, this setup can be modified by assigning different terms to these characters. It is likely that the models evaluated in this study have encountered conventional *knights and knaves* puzzles during their training procedure, as such examples are readily available on the internet (Smullyan, 1978). By altering the terms used for truth-tellers and liars, we can significantly reduce the likelihood that the models have been exposed to similar samples in our benchmark. Consequently, we analyze the impact of the terminology used for *knights* and *knaves* on model performance. Specifically, we examine three different formulations: (i) the conventional *knights* and *knaves*, (ii) neutral descriptions such as *truth-tellers* and *liars*, and (iii) pseudo-terms such as *jabbas* and *tettes*. Figure 7 illustrates the zero-shot performance of all models for each terminology setup across the different subsets of *TruthQuest*. Surprisingly, we find no substantial impact of the choice of terms on the models' task performance. We hypothesize that this may be because the specific instances generated for *TruthQuest* have not yet been exposed to the internet. Consequently, the models might have encountered only a negligible fraction of instances by chance during their training process.

| Abbreviation | Error Category | Description |
|---|---|---|
| **SR** | False statement reproduction | A problem statement is repeated incorrectly by the model. |
| **TS** | Assuming statements to be true | The possibility that statements might be lies is not considered. |
| **TL** | Misunderstanding the concept of truth and lies | Making false assumptions about the nature of truth-tellers or liars. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. |
| **LO** | Misunderstanding logical operators | The logical implications of a potentially false statement are not properly deduced. |
| **UC** | Unjustified Conclusion | A conclusion such as "X is a truth-teller/liar" is presented without proper justification. |
| **UF** | Unfaithfulness | A new conclusion explicitly contradicts a conclusion previously drawn. |

Table 4: Error categories and their respective descriptions.

## B.2 Error Analysis

We examine the models' reasoning errors through both comprehensive manual inspections and AI-based evaluations of their rationales. The following sections present additional results from both evaluation methods.

### B.2.1 Human Evaluation

As outlined in Section 4, we manually evaluate a subset of the models' responses to asses their errors encountered during reasoning (for a detailed explanation of the evaluation procedure, please refer to Section C.1). To supplement our findings summarized in Figure 2, , we report the common errors exhibited by LLaMA-3-70B when prompted in a zero-shot setting (see Figure 4). We observe that, similar to LLaMA-3-8B (zero-shot) and LLaMA-3-70B (four-shot CoT), LLaMA-3-70B (zero-shot) frequently displays errors when deducing the logical implications of potentially false statements *(LO)*. Additionally, we find that while the model struggles less with understanding the concept of truth and lies *(TL)* compared to LLaMA-3-8B (zero-shot), it still exhibits this error category more frequently than LLaMA-3-70B (four-shot CoT). This trend, indicating that more proficient models better grasp the concept of truth and lies than lower-performing ones, is also reflected in our analysis of all model responses conducted via automatic LLM-based evaluation (for details, refer to Section B.2.2).

### B.2.2 AI-Assisted Evaluation

High-quality human annotations are typically costly to obtain. In our study, we manually inspect 360 model responses from three different LLMs, where each instance is evaluated twice independently by two annotators (for details on the evaluation procedure, please refer to Section C.1).

However, as our benchmark comprises 2,400 different instances, this evaluation procedure only covers a small subset of the models' responses. To complement our manual evaluation, we employ GPT-4 (OpenAI et al., 2024)[6] to assess all 2,400 responses of a model with respect to the reasoning errors outlined in Table 4 (for details about the exact prompts we employ, or the alignment between human and AI-based evaluation, please refer to Section C.2). The respective results are illustrated in Figure 5. We present the relative occurrences of each error category for LLaMA-2-7B (zero-shot), LLaMA-3-8B (zero-shot), LLaMA-3-70B (zero-shot), and LLaMA-3-70B (four-shot CoT). All values are av-

---

[6] Specifically, version `gpt-4o-2024-05-13`.



Figure 4: Relative occurrences of the reasoning errors displayed by LLaMA-3-70B when prompted in a zero-shot setting. Values are obtained from human evaluation. Error categories are abbreviated for a more comprehensive overview: *(SR) False statement reproduction*, *(TS) Assuming statements to be true*, *(TL) Misunderstanding the concept of truth and lies*, *(LO) Misunderstanding logical operators*, *(UC) Unjustified conclusion*, and *(UF) Unfaithfulness*.

10

(a) LLaMA-2-7B (zero-shot)

(b) LLaMA-3-8B (zero-shot)

(c) LLaMA-3-70B (zero-shot)
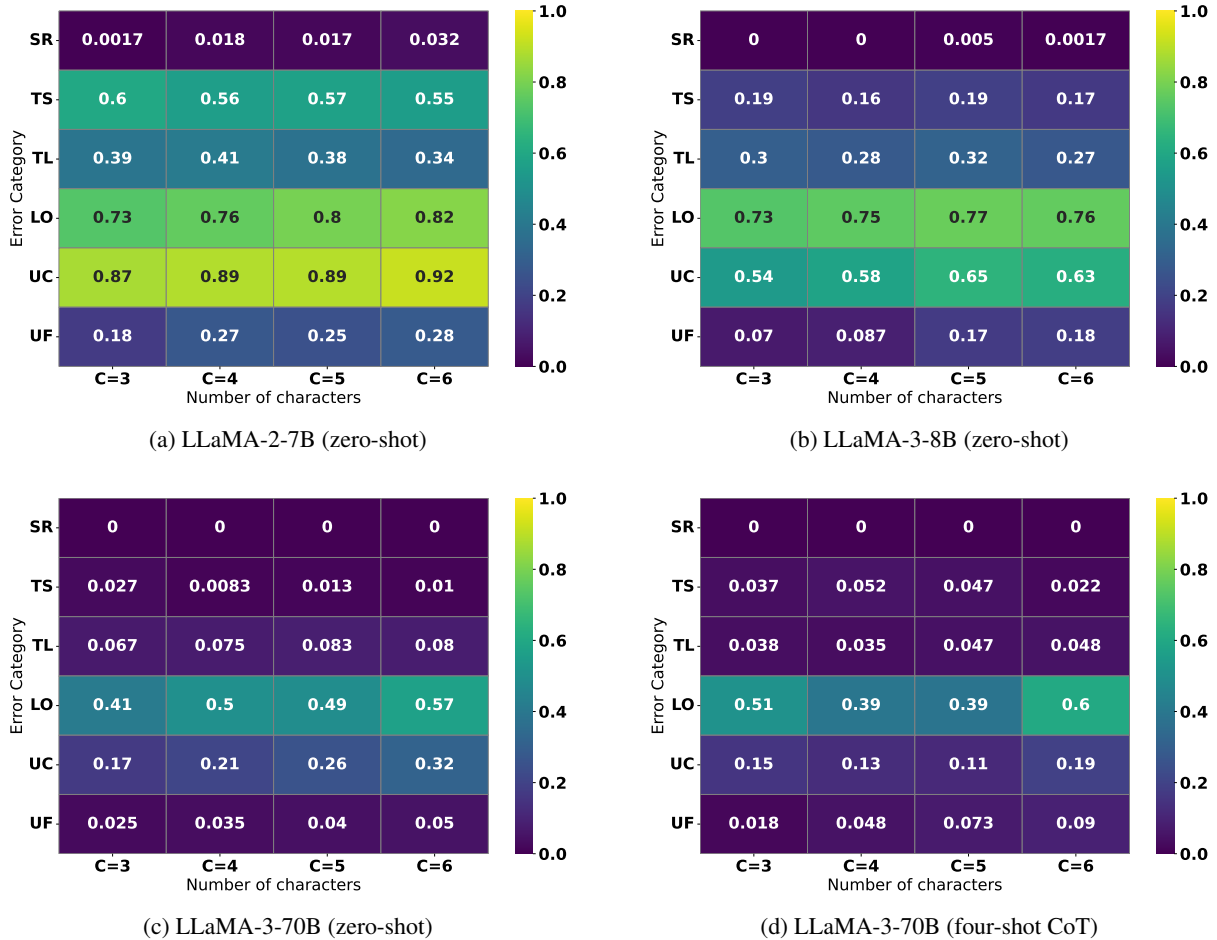
(d) LLaMA-3-70B (four-shot CoT)

Figure 5: Relative occurrences of the reasoning errors displayed by LLaMA-2-7B, LLaMA-3-8B, and LLaMA-3-70B in zero-shot prompting, as well as LLaMA-3-70B when prompted via four-shot chain-of-thought. Values are obtained from GPT-4-based evaluations. Error categories are abbreviated for a more comprehensive overview: *(SR) False statement reproduction*, *(TS) Assuming statements to be true*, *(TL) Misunderstanding the concept of truth and lies*, *(LO) Misunderstanding logical operators*, *(UC) Unjustified conclusion*, and *(UF) Unfaithfulness*.

eraged across the different statement sets for each number of characters. Consistent with the results obtained through human evaluation, we observe a strong trend for higher-performing models to converge on errors related to deducing the correct logical implications of statements *(LO)*. In contrast, lower-performing models such as LLaMA-2-7B (zero-shot) or LLaMA-3-8B (zero-shot), display diverse errors, ranging from misconceptions about truth and lies *(TL)* to unjustified conclusions *(UC)*. Notably, LLaMA-2-7B is the only model that frequently fails to consider statements as lies *(TS)*.

## C  Evaluation Procedures

In this study, we utilize two types of evaluation methods: human evaluations and AI-assisted evaluations. Below, we provide further details on each method, including the instructions given to human annotators and the process by which large language

models are employed to generate similar annotations automatically. Finally, we assess the quality of our automatic evaluation procedures by comparing the results to the results obtained via manual assessment.

### C.1  Human Evaluation

As outlined in Section 4, we manually inspect 360 responses from LLaMA-3-8B (zero-shot) and LLaMA-3-70B (zero-shot and four-shot CoT). This manual evaluation is independently conducted by two hired students with expertise in data annotation. Both student annotators are compensated according to national standards.

### C.1.1  Annotator Instructions

To ensure high-quality annotations, we provide extensive training to both annotators. This training involves multiple sessions in which we introduce

| | SR | TS | TL | LO | UC | UF |
|---|---|---|---|---|---|---|
| LLaMA-3-8B zero shot | 0.49 | 0.33 | 0.40 | 0.68 | 0.72 | 0.53 |
| LLaMA-3-70B zero shot | −0.01 | 0.37 | 0.41 | 0.51 | 0.22 | 0.15 |
| LLaMA-3-70B four CoT | - | 0.90 | 0.51 | 0.74 | −0.04 | 0.34 |

Table 5: Cohen's Kappa values to assess the human inter-annotator agreement across different models, prompt setups and and error categories. Error categories are abbreviated for a more comprehensive overview: *(SR) False statement reproduction*, *(TS) Assuming statements to be true*, *(TL) Misunderstanding the concept of truth and lies*, *(LO) Misunderstanding logical operators*, *(UC) Unjustified conclusion*, and *(UF) Unfaithfulness*.

the annotators to *knights and knaves* puzzles, asking them to solve these puzzles by hand to familiarize themselves with the task structure. Once the annotators are confident in solving puzzles of this style, we present exemplary responses from the models evaluated in this study. Together, we discuss notable behaviors and errors exhibited by the models. Next, we introduce the annotators to the six error categories outlined in Table 4. We proceed only when both annotators confirm their full understanding of each error type and have no further questions. The annotators are then tasked with independently annotating model responses. For each response, they parse the model's conclusion and assign a binary label (yes/no) to each error category, indicating its presence or absence in the model's reasoning. Initially, the annotators work with practice examples to highlight and address any ambiguities in the annotation process. They only move on to labeling the actual model responses when they are confident in their understanding of the labeling process. To maintain high annotation quality, we ask both annotators to review their annotations, ensuring any potential errors in their annotations are accounted for.

### C.1.2  Inter-Annotator Agreement

To assess the quality of our manual annotations, we calculate the inter-annotator agreement, reporting an overall Cohen's kappa value of $\kappa = 0.70$, which indicates substantial agreement between the two annotators. Table 5 presents Cohen's kappa values for each model and error type. We observe that the agreement rate varies across different categories, ranging from none to perfect agreement. Notably, the values for *False statement reproduction (FS)* in LLaMA-3-70B (zero-shot) and *Unjustified conclusion (FS)* in LLaMA-3-70B (four-shot CoT) are almost zero. This is likely due to a strong bias in the label distribution towards *no* labels, as these errors rarely occur in these models. All human annotations are publicly available at

https://huggingface.co/datasets/mainlp/TruthQuest-Human-Annotations.

### C.2  AI-Assisted Evaluation

In addition to the human evaluation, we employ GPT-4 to assess the models' reasoning errors. Similar to the human annotators, GPT-4 is tasked with assigning binary labels (yes/no) to each error category described in Table 4, indicating the presence or absence of the error type in the model's reasoning. Additionally, GPT-4 is required to provide a justification for each label assigned. To ensure GPT-4's comprehension of each error category, we provide detailed descriptions in the model input. The full prompt can be found in Figure 10. Furthermore, we present six few-shot examples illustrating the desired annotation behavior (see Figures 11 to 16). To assess the quality of the annotations, we compute the Pearson correlation for the error distributions of LLaMA-3-8B (zero-shot), LLaMA-3-70B (zero-shot), and LLaMA-3-70B (four-shot CoT) between the automatically obtained labels and the human annotations. All correlation coefficients and their corresponding p-values are reported in Table 6. Overall, we find that the error distribution obtained through GPT-4 strongly correlates with the error distribution obtained via manual labeling. However, on an instance level, we observe only fair agreement, with an overall Cohen's kappa value of $\kappa = 0.34$. For our automatic evaluation of 9,600 model responses (4 models × 2,400 responses), the cost was approximately $250. All annotations obtained through GPT-4 are publicly available at https://huggingface.co/datasets/mainlp/TruthQuest-AI-Annotations.

**Task performance.**  To assess the quality of our performance evaluation procedure depicted in Figure 6, we compute the proportion of instances where the final conclusions derived from our two-step method match those reported by manual assessment. We find an alignment of 100%.

| Model | Characters_3 | Characters_4 | Characters_5 | Characters_6 |
|---|---|---|---|---|
| LLaMA-3-8B Zero Shot | 0.673 (0.143) | 0.754 (0.084) | 0.8211 (0.045) | 0.859 (0.029) |
| LLaMA-3-70B Zero Shot | 0.734 (0.097) | 0.858 (0.0289) | 0.762 (0.078) | 0.819 (0.046) |
| LLaMA-3-70B Four CoT | 0.877 (0.022) | 0.857 (0.0294) | 0.923 (0.009) | 0.962 (0.002) |

Table 6: Pearson correlation for the distribution of reasoning errors computed between the human and AI-based error analyses. The Pearson correlation coefficients are computed for different numbers of characters and models. Respective p-values are reported in parentheses.

## D  Prompts

We present all prompts used in this study. The task prompt is shown in Figure 8. Figure 9 illustrates the prompt for the two-step conclusion evaluator. Additionally, the system prompt for GPT-4 is provided in Figure 10, along with the few-shot examples in Figures 11 to 16. Few-shot examples are designed to demonstrate the desired assistant behavior (highlighted in orange) in response to the corresponding user request (highlighted in blue). All prompts are publicly accessible at https://github.com/mainlp/TruthQuest.

| Model | Mode | Set S | | | | Set I | | | | Set E | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 3 | 4 | 5 | 6 | 3 | 4 | 5 | 6 | 3 | 4 | 5 | 6 |
| **Random Baseline** | - | 0.13 | 0.06 | 0.03 | 0.02 | 0.13 | 0.06 | 0.03 | 0.02 | 0.13 | 0.06 | 0.03 | 0.02 |
| **LLaMA-2-7b** | | 0.08 | 0.06 | 0.02 | 0.00 | 0.20 | 0.10 | 0.07 | 0.04 | 0.21 | 0.11 | 0.03 | 0.03 |
| **LLaMA-2-13b** | | 0.10 | 0.06 | 0.03 | 0.03 | 0.13 | 0.11 | 0.04 | 0.02 | 0.15 | 0.08 | 0.05 | 0.01 |
| **LLaMA-2-70b** | **zero shot** | 0.13 | 0.13 | 0.08 | 0.03 | 0.13 | 0.11 | 0.06 | 0.03 | 0.17 | 0.09 | 0.07 | 0.02 |
| **LLaMA-3-8B** | | 0.07 | 0.13 | 0.04 | 0.04 | 0.19 | **0.18** | 0.07 | 0.04 | 0.13 | 0.08 | 0.06 | 0.03 |
| **LLaMA-3-70B** | | **0.29** | **0.22** | **0.13** | **0.10** | **0.32** | 0.14 | **0.14** | **0.11** | **0.29** | **0.18** | **0.11** | **0.06** |
| **Mixtral-8x7B** | | 0.16 | 0.08 | 0.04 | 0.03 | 0.21 | 0.14 | 0.06 | 0.05 | 0.17 | 0.08 | 0.04 | 0.01 |
| | **four shot** | 0.13 | **0.14** | 0.04 | **0.05** | 0.19 | **0.14** | 0.08 | **0.05** | 0.18 | 0.09 | **0.09** | 0.02 |
| | **eight shot** | 0.14 | 0.09 | 0.07 | 0.04 | 0.23 | **0.14** | 0.07 | 0.05 | 0.19 | 0.10 | 0.07 | **0.04** |
| **LLaMA-3-8B** | **zero CoT** | 0.12 | 0.06 | 0.05 | 0.03 | 0.16 | 0.07 | 0.06 | 0.04 | 0.20 | 0.13 | 0.07 | 0.03 |
| | **four CoT** | 0.13 | 0.12 | 0.06 | 0.04 | 0.12 | 0.12 | 0.08 | 0.04 | 0.21 | 0.09 | 0.05 | 0.03 |
| | **eight CoT** | **0.27** | 0.13 | **0.08** | 0.03 | **0.26** | 0.12 | **0.12** | 0.05 | **0.26** | **0.15** | 0.08 | 0.03 |
| | **four shot** | 0.22 | 0.25 | 0.19 | 0.13 | 0.24 | 0.21 | 0.13 | 0.10 | 0.32 | 0.22 | 0.11 | 0.05 |
| | **eight shot** | 0.22 | 0.21 | 0.16 | 0.09 | 0.32 | 0.25 | 0.07 | 0.09 | 0.27 | 0.20 | 0.10 | 0.02 |
| **LLaMA-3-70B** | **zero CoT** | 0.23 | 0.17 | 0.14 | 0.12 | 0.28 | 0.17 | 0.15 | 0.09 | 0.29 | 0.17 | **0.12** | 0.08 |
| | **four CoT** | 0.46 | **0.31** | **0.21** | 0.16 | **0.33** | **0.27** | 0.11 | **0.15** | **0.40** | **0.25** | 0.12 | **0.10** |
| | **eight CoT** | **0.60** | 0.26 | **0.21** | **0.20** | **0.33** | 0.20 | **0.15** | 0.12 | 0.37 | 0.20 | **0.12** | **0.10** |

Table 7: Additional accuracy values of LLaMA-3-8B for different prompting techniques across each subset of *TruthQuest*. Bold values represent highest performance among a group. The random baseline indicates the accuracy achieved by guessing the identity of each character.
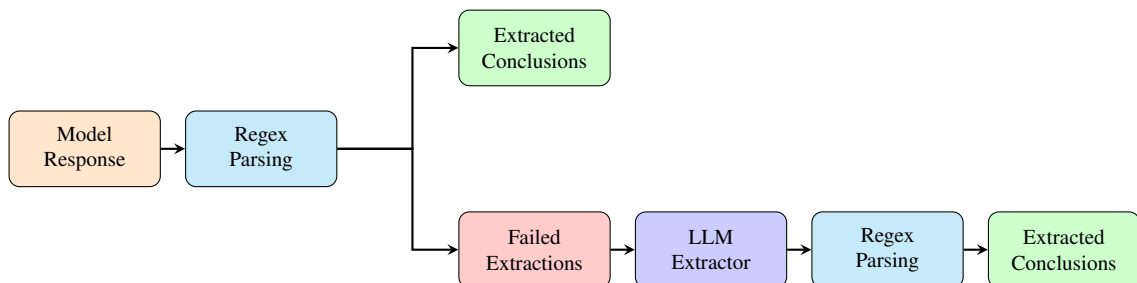


Figure 6: A schematic overview of the conclusion evaluator.

(a) *knights* and *knaves*

(b) *truth-tellers* and *liars*

(c) *jabbas* and *tettes*

(d) *knights* and *knaves*

(e) *truth-tellers* and *liars*

(f) *jabbas* and *tettes*

(g) *knights* and *knaves*

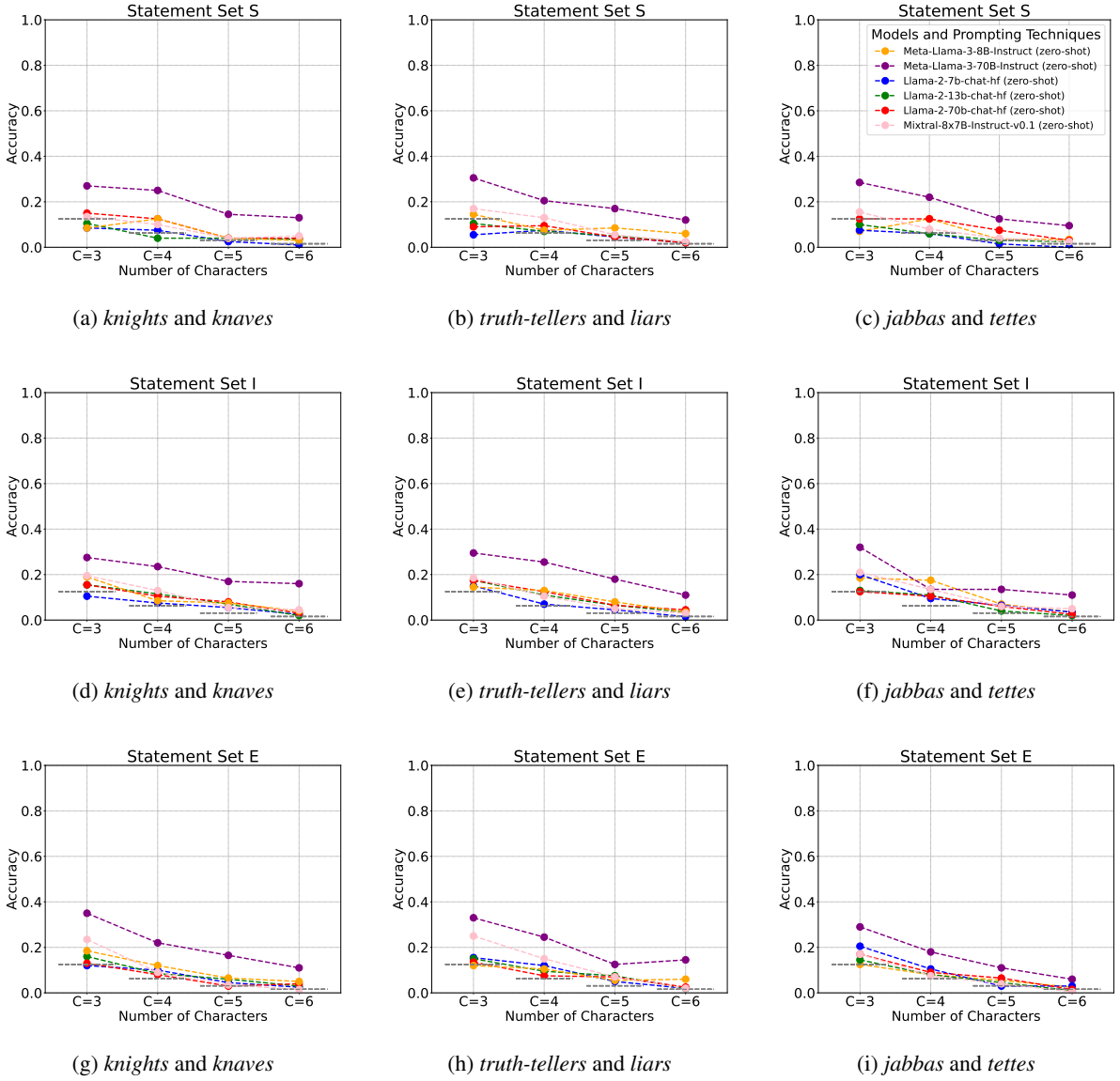(h) *truth-tellers* and *liars*

(i) *jabbas* and *tettes*

Figure 7: Zero-shot performance of all models across the data subsets of *TruthQuest*, focusing on the terminology used for characters who always tell the truth, and those who always lie. Three different formulations are compared: (i) the conventional *knights* and *knaves*, (ii) neutral descriptions such as *truth-tellers* and *liars*, and (iii) pseudo-terms such as *jabbas* and *tettes*.

[INST] «SYS»
Your task is to solve a logical reasoning problem. You are given set of statements from which you must logically deduce the identity of a set of characters.

You must infer the identity of each character. First, explain your reasoning. At the end of your answer, you must clearly state the identity of each character by following the format:

CONCLUSION:
A: ...
B: ...
C: ...
...
«/SYS»
### Instruction ###
Assume that there exist only two types of people: knights and knaves. Knights always tell the truth, while knaves always lie.
You are given the statements from <number of characters> characters. Based on their statements, infer who is a knights and who is a knave.

Based on the following statements, infer who is a knight and who is a knave:
<statements>

First, explain your reasoning. End your answer by clearly stating the identity of each character in the following format:

A: knight/knave
B: knight/knave
C: knight/knave
... [/INST]

Figure 8: The task prompt. Placeholders such as "*<number of characters>*" and "*<statements>*" are replaced by the corresponding input of the problem description. Note that the arrangement and usage of special tokens may vary depending on the specific language model employed.

[INST] «SYS»
Your task is to extract information from a given piece of text. You are given a piece of text that reasons about the identity of a set of characters. You must extract the final answer from the given text by clearly stating the deduced identity of each character in the following format:
A: ...
B: ...
C: ...
...
You must strictly follow this format. Do not generate any additional text. You must only extract information that is present in the text provided. Do not add information that goes beyond the one contained in the text at hand.
«/SYS»
### Instruction ###
You are given a piece of TEXT that reasons about the identity of a set of characters. You must extract the FINAL ANSWER from the TEXT by clearly stating the deduced identity of each character in the following format:
### FINAL ANSWER:
A: ...
B: ...
C: ...
...
You must strictly follow this format. Do not generate any additional text. You must only extract information that is present in the text provided. Do not add information that goes beyond the one contained in the text at hand. If "not a" is placed in front of the identity, flip the result of that character, e.g. "not a knight" becomes "knave", while "not a knave" becomes "knight". You should simplify results like "knight/knight" to "knight" and "knave/knave" to "knave". However, do not simplify "knight/knave", or "knave/knight".
### TEXT:
Great! Let's apply logical reasoning to deduce the identity of each character.
A: Based on statement A, we know that C is a knave and D is a knave. This means that C and D are both knaves.
B: Statement B tells us that C is a knave and D is a knave. Since C and D are both knaves, we can conclude that B is a knight.
C: Statement C tells us that D is a knave. Since D is already identified as a knave in statements A and B, we can conclude that C is a knight.
D: Statement D tells us that B is a knight. Since B has already been identified as a knight in statement B, we can confirm that D is a knave.
Therefore, the identities of each character are:
A: not a knight
B: knight/knave
C: knight/knight
D: knave [/INST]

<|assistant|> ### FINAL ANSWER:
A: knave
B: knight/knave
C: knight
D: knave </|assistant|>

[INST] ### TEXT:
Sure! Let's start by analyzing the statements one by one.
1. A: B is a knave.
Since A says that B is a knave, we can deduce that B is a knave.
2. B: If D is a knight, then C is a knave.
This statement implies that if D is a knight, then C is a knave. However, since B is a knave, we can deduce that D is not a knight. Therefore, C must be a knight.
3. C: B is a knight.
Since C says that B is a knight, we can deduce that B is a knight.
4. D: A is a knave and C is a knave.
This statement is not helpful in determining the identity of D, as D is simply repeating what A and C have said.
Now, let's summarize our conclusions:
A: knave
B: knave
C: knight
D: unable to determine
However, based on the statements, we can deduce that B and C are knaves, and A is a knight. Therefore, the correct identities of the characters are: A: knight
B: knave
C: knave
D: unable to determine [/INST]

<|assistant|> ### FINAL ANSWER:
A: knight
B: knave
C: knave
D: unable to determine </|assistant|>

[INST] ### TEXT:
Okay, I understand the problem. Here's my reasoning:
From statement A, if B is a knave, then C is a knight. This means that if B is a knave, then C must be a knight.
From statement B, if A is a knave, then C is a knight. This means that if A is a knave, then C must be a knight.
Now, let's examine the third statement
C: B is a knave. If B is a knave, then C is a knave. But we know that knaves always lie, so we can't trust what C says. Therefore, we can't determine the identity of C based on this statement alone.
However, we can use the first two statements to deduce the identity of C. If A is a knave, then C is a knight, and if B is a knave, then C is a knight. Since we don't know the identity of A and B, we can't determine the identity of C.
CONCLUSION:
A: knave/knight (cannot be determined)
B: knave/knight (cannot be determined)
C: knave [/INST]

<|assistant|> ### FINAL ANSWER:
A: knave/knight
B: knave/knight
C: knave </|assistant|>

<Actual task description>

Figure 9: The three-shot conclusion evaluator prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|system|>
Your task is to classify whether a set of errors are present in the response of a large language model. You are given the initial PROBLEM STATEMENT and the model's RESPONSE.
You must consider the following error types:

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a knights always tells the truth and a knave always lies. For instance, the model mistakenly assumes that knights lie, while knaves tell the truth. Alternatively, the model assumes that knights can only talk about other knights and knaves can only talk about knaves.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a knave, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a knights/knave" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a knights/knave..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

For each error type, you must classify whether the error is present in the model's response (yes) or not (no). Start with a short justification for your decision (max. 1-2 sentences) and then clearly state whether the error type is present (yes) or not (no).
Strictly adhere to the following output format:

### False statement reproduction
- Justification: ...
- Label: yes/no

### Assuming statements to be true
- Justification: ...
- Label: yes/no

### Misunderstanding the concept of truth and lies
- Justification: ...
- Label: yes/no

...

You must strictly follow this format. Do not generate any additional text.
</|system|>.

Figure 10: The system message of the six-shot error evaluator prompt.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
A: If C is a truth-teller, then D is a truth-teller.
B: C is a liar and A is a liar.
C: D is a liar and B is a liar.
D: B is a truth-teller and C is a liar.

### RESPONSE
1. A: If C is a truth-teller, then D is a truth-teller.
This statement is made by A, who is either a truth-teller or a liar. If A is a truth-teller, then the statement is true, and D is indeed a truth-teller. But if A is a liar, then the statement is false, and D is not a truth-teller. So, A's statement is true only if A is a truth-teller.

2. B: C is a liar and A is a liar.
This statement is made by B, who is either a truth-teller or a liar. If B is a truth-teller, then the statement is false, as C is not a liar. But if B is a liar, then the statement is true, as C is indeed a liar. So, B's statement is true only if B is a liar.

3. C: D is a liar and B is a liar.
This statement is made by C, who is either a truth-teller or a liar. If C is a truth-teller, then the statement is false, as B is not a liar. But if C is a liar, then the statement is true, as D is indeed a liar. So, C's statement is true only if C is a liar. 4. D: B is a truth-teller and C is a liar.

This statement is made by D, who is either a truth-teller or a liar. If D is a truth-teller, then the statement is true, as B is indeed a truth-teller. But if D is a liar, then the statement is false, as C is not a liar. So, D's statement is true only if D is a truth-teller.

Now, let's combine the information:
* A's statement is true only if A is a truth-teller.
* B's statement is true only if B is a liar.
* C's statement is true only if C is a liar.
* D's statement is true only if D is a truth-teller.

From statement 2, we know B is a liar. From statement 4, we know D is a truth-teller.
From statement 1, we know that if A is a truth-teller, then D is a truth-teller. Since we already know D is a truth-teller, A must be a truth-teller.
From statement 3, we know C is a liar.

CONCLUSION:
A: truth-teller
B: liar
C: liar
D: truth-teller
<|/user|>

<|assistant|>
### False statement reproduction
- Justification: The statements are properly reproduced by the model.
- Label: no

### Assuming statements to be true
- Justification: The possibility that the characters' statements are lies is considered by the model.
- Label: no

### Misunderstanding the concept of truth and lies
- Justification: When discussing B's statement, the model says: "So, B's statement is true only if B is a liar". However, if B were a liar, its statement would not be true. This shows a model failure to understand the concept of a liar.
- Label: yes

### Misunderstanding logical operators
- Justification: The model fails to understand the logical implications of A's statement: "If C is a truth-teller, then D is a truth-teller." It draws conclusions about D based on the identity of A, ignoring character C, who is the antecedent of the logical implication.
- Label: yes

### Unjustified conclusion
- Justification: When discussing B's statement, the model defines the identity of character C: "as C is indeed a liar", without proper justification.
- Label: yes

### Unfaithfulness
- Justification: The model does not draw any conclusions that contradict previous conclusions.
- Label: no
<|/assistant|>

Figure 11: The first few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
A: C is a liar and B is a truth-teller.
B: I am a truth-teller.
C: A is a truth-teller.

### RESPONSE
Great, let's dive into the problem!
To solve this problem, we need to use logical reasoning and deduction. We are given three statements from three different characters, and our goal is to determine who is a truth-teller and who is a liar based on their statements.
Statement A tells us that C is a liar and B is a truth-teller. This means that B is telling the truth, since he is a truth-teller.
Statement B tells us that I (B) am a truth-teller. This means that I am telling the truth, since I am a truth-teller.
Statement C tells us that A is a truth-teller. This means that A is telling the truth, since A is a truth-teller.
Now, let's use these statements to deduce the identity of each character:

A: truth-teller
B: truth-teller
C: truth-teller
</|user|>

<|assistant|>
### False statement reproduction
- Justification: No occurrence of this error in the model's response.
- Label: no

### Assuming statements to be true
- Justification: The model treats each character's statement as a true premise and disregards the possibility that statements might be lies. See for instance:
"A tells us that C is a liar and B is a truth-teller. This means that B is telling the truth, since he is a truth-teller."
- Label: yes

### Misunderstanding the concept of truth and lies
- Justification: No occurrence of this error in the model's response.
- Label: no

### Misunderstanding logical operators
- Justification: No occurrence of this error in the model's response.
- Label: no

### Unjustified conclusion
- Justification: No occurrence of this error in the model's response.
- Label: no

### Unfaithfulness
- Justification: No occurrence of this error in the model's response.
- Label: no
</|assistant|>

Figure 12: The second few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
A: C is a truth-teller and B is a liar.
B: If A is a truth-teller, then C is a truth-teller.
C: B is a liar and A is a liar.

### RESPONSE
REASONING:
Let's analyze the statements:

A: C is a truth-teller and B is a liar.
B: If A is a truth-teller, then C is a truth-teller.
C: B is a liar and A is a liar.

If A is a truth-teller, then A's statement must be true, so C is a truth-teller and B is a liar.
If C is a truth-teller, then C's statement must be false, so either B is not a liar or A is not a liar.
This contradicts our assumption that A is a truth-teller, so A must be a liar.

If A is a liar, then A's statement must be false, so either C is not a truth-teller or B is not a liar.
Since A's statement is false, we can't conclude anything about C, but we know B is not a liar.
If B is a truth-teller, then B's statement must be true, so A is a liar and C is a truth-teller.
This is consistent with our previous conclusion that A is a liar.
If C is a truth-teller, then C's statement must be false, so either B is not a liar or A is not a liar.
We know A is a liar, so B must be a truth-teller.

CONCLUSION:
A: liar
B: truth-teller
C: truth-teller
</|user|>

<|assistant|>
### False statement reproduction
- Justification: The model correctly repeats the statements in the problem description.
- Label: no

### Assuming statements to be true
- Justification: The model considers the possibility that statements might be false.
- Label: no

### Misunderstanding the concept of truth and lies
- Justification: When discussing C's statement, the model says: "If C is a truth-teller, then C's statement must be false". However, if C were a truth-teller, its statement could not be false. This shows a failure of the model to understand the concept of truth.
- Label: yes

### Misunderstanding logical operators
- Justification: The model exhibits several problems to draw the logical implications of statements. For example, it says: "Since A's statement is false, we can't conclude anything about C". However, the logical conjunction of A's statement allows us to draw conclusions about C.
- Label: yes

### Unjustified conclusion
- Justification: No occurrence of this error in the model's response.
- Label: no

### Unfaithfulness
- Justification: No occurrence of this error in the model's response.
- Label: no
</|assistant|>

Figure 13: The third few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
Assume that there exist only two types of people: truth-tellers and liars. truth-tellers always tell the truth, while liars always lie.
You are given the statements from a set of characters. Based on their statements, infer who is a truth-teller and who is a liar:
A: D is a truth-teller.
B: C is a liar and A is a liar.
C: A is a liar.
D: A is a liar and B is a truth-teller.

### RESPONSE
Let's analyze the statements given:
1. A: D is a truth-teller.
This statement is made by A, who according to statement D, is a liar. Therefore, this statement is a lie, and D is not a truth-teller.
2. B: C is a liar and A is a liar.
This statement is made by B, who according to statement D, is a truth-teller. Therefore, this statement is true, and C and A are both liars.
3. C: A is a liar.
This statement is made by C, who according to statement B, is a liar. Therefore, this statement is true, and A is a liar.
4. D: A is a liar and B is a truth-teller.
This statement is made by D, who according to statement A, is a truth-teller. Therefore, this statement is true, and A is a liar and B is a truth-teller.

Now, let's summarize our findings:
* A is a liar (according to statements B and C).
* B is a truth-teller (according to statement D).
* C is a liar (according to statement B).
* D is a truth-teller (according to statement A).

CONCLUSION:
A: truth-teller
B: liar
C: liar
D: truth-teller
<|user|>

<|assistant|>
### False statement reproduction
- Justification: The model correctly repeats the statements in the problem description.
- Label: no

### Assuming statements to be true
- Justification: The model considers the possibility that statements might be false.
- Label: no

### Misunderstanding the concept of truth and lies
- Justification: When discussing C's statement, the model says: "This statement is made by C, who according to statement B, is a liar. Therefore, this statement is true...". However, if C were a liar, its statement could not be true. This shows a failure of the model to understand the concept of truth.
- Label: yes

### Misunderstanding logical operators
- Justification: The model does not show any failures at understanding logical operators.
- Label: no

### Unjustified conclusion
- Justification: When discussion the implications of A's statement, the model defines D's statement as truth without properly justifying this assumption: "This statement is made by A, who according to statement D, is a liar. Therefore, this statement is a lie, and D is not a truth-teller."
- Label: yes

### Unfaithfulness
- Justification: The model first concludes that "D is not a truth-teller", but then follows up "This statement is made by D, who according to statement A, is a truth-teller. Therefore, this statement is true." This contradicts the initial conclusion.
- Label: yes
<|assistant|>

Figure 14: The fourth few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
A: F is a liar.
B: A is a truth-teller.
C: D is a truth-teller and B is a truth-teller.
D: E is a truth-teller.
E: D is a liar and B is a truth-teller.
F: I am a truth-teller.

### RESPONSE
Let's start by analyzing the statements:

1. A: F is a liar.
If A is a truth-teller, then F is indeed a liar, which means A is a truth-teller. But if A is a liar, then F is not a liar, which is a contradiction. Therefore, A must be a truth-teller.
2. B: A is a truth-teller.
Since A is a truth-teller, this statement is true. Therefore, B is a truth-teller.
3. C: D is a truth-teller and B is a truth-teller.
Since C is a truth-teller, this statement is true. Therefore, D is a truth-teller, and B is indeed a truth-teller.
4. D: E is a truth-teller.
Since D is a truth-teller, this statement is true. Therefore, E is a truth-teller.
5. E: D is a liar and B is a truth-teller.
Since E is a liar, this statement is false. However, we know that D is a truth-teller from statement 3, and B is a truth-teller from statements 2 and 3. Therefore, E must be a truth-teller.
6. F: I am a truth-teller.
Since F is a liar, this statement is false. Therefore, F is a liar.

Now, let's summarize the identities:

A: truth-teller
B: truth-teller
C: truth-teller
D: truth-teller
E: truth-teller
F: liar
<|/user|>

<|assistant|>
### False statement reproduction
- Justification: The model correctly repeats the statements in the problem description.
- Label: no

### Assuming statements to be true
- Justification: The model considers the possibility that statements might be false.
- Label: no

### Misunderstanding the concept of truth and lies
- Justification: The model's response does not exhibit this type of error.
- Label: no

### Misunderstanding logical operators
- Justification: The model does not show any failures at understanding logical operators.
- Label: no

### Unjustified conclusion
- Justification: When discussion C's statement, the model defines C as a truth-teller without proper prior justification.
- Label: yes

### Unfaithfulness
- Justification: The model first concludes that "E is a truth-teller". However, it then continues with: "Since E is a liar, this statement is false", which contradicts the previous conclusion.
- Label: yes
<|/assistant|>

Figure 15: The fifth few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.

<|user|>
### Instruction ###
You are given a PROBLEM STATEMENT and the RESPONSE of a large language model, respectively. Your task is to classify whether a set of errors are present in the model's response. For each of the following error types, classify whether the error is present in the model's response (yes) or not (no).

- False statement reproduction (A statement of the problem description is repeated incorrectly.)
- Assuming statements to be true (The possibility that a character's statement is a lie is not considered. The model does not question the premises' truth but assumes it as given.)
- Misunderstanding the concept of truth and lies (The model does not understand that a truth-teller always tells the truth and a liar always lies. For instance, the model mistakenly assumes that truth-tellers lie, while liars tell the truth. Alternatively, the model assumes that truth-tellers can only talk about other truth-tellers and liars can only talk about liars.)
- Misunderstanding logical operators (Logical operators and their implications are not properly understood. The model fails at inferring the logical implications of a statement. For instance, if character A's statement is a logical conjunction (e.g. B and C) and A is a liar, the model does not consider all possible implications of that lie (not B and C, B and not C, not B and not C).)
- Unjustified conclusion (A conclusion about the identity of a characters such as "X is a truth-teller/liar" is presented without ANY justification. Note that this error does not address the logical validity of a justification, but checks whether a justification is present or not. As an example of this error, the model might simply state: "Since A is a truth-teller/liar..." without prior justification.)
- Unfaithfulness (A conclusion is drawn that explicitly contradicts a conclusion previously inferred by the model, without the model reflecting on this contradiction.)

Start with a short justification for your decision (max. 1-2 sentences). Then, clearly state whether the error type is present (yes) or not (no).

### PROBLEM STATEMENT
A: B is a liar if and only if C is a liar.
B: C is a liar.
C: A is a liar.

### RESPONSE
REASONING:
Let's assume B is a truth-teller.
If B is a truth-teller, B's statement is true and C is a liar.
If C is a liar, C's statement is false and A is a truth-teller.
If A is a truth-teller, A's statement is true, so either B is a liar and C is a liar or B is a truth-teller and C is a truth-teller.
Since we know C is a liar, B must be a liar.
This leads to a contradiction, so B must be a liar.
If B is a liar, B's statement is false and C is a truth-teller.
If C is a truth-teller, C's statement is true and A is a liar.
If A is a liar, A's statement is false, so either B is a truth-teller and C is a truth-teller or B is a liar and C is a liar.
Since we know B is a liar, C must be a truth-teller.
Therefore, we can conclude:

CONCLUSION:
A: liar
B: liar
C: truth-teller
<|/user|>

<|assistant|>
### False statement reproduction
- Justification: The model correctly repeats the statements in the problem description.
- Label: no

### Assuming statements to be true
- Justification: The model considers the possibility that statements might be false.
- Label: no

### Misunderstanding the concept of truth and lies
- Justification: The model seem to capture the concept of truth and lies.
- Label: no

### Misunderstanding logical operators
- Justification: The model fails to infer the logical implications of A's false statement. Although the model identifies A's statement as a lie, it does not successfully build the negation of the logical equivalence: "If A is a liar, A's statement is false, so either B is a truth-teller and C is a truth-teller or B is a liar and C is a liar."
- Label: yes

### Unjustified conclusion
- Justification: All conclusions are justified.
- Label: no

### Unfaithfulness
- Justification: The model does not infer conclusions that contradict conclusions previously drawn.
- Label: no
<|/assistant|>

Figure 16: The sixth few-shot example of the error evaluator six-shot prompt. Examples are presented in chat form, where the task prompt is depicted in blue, and the desired answer is exemplified in orange.