Responsible NLP Checklist

Paper title: LogicTree: Structured Proof Exploration for Coherent and Rigorous Logical Reasoning with Large Language Models

Authors: Kang He, Kaushik Roy

How to read the checklist symbols:	
the authors responded 'yes'	
the authors responded 'no'	
the authors indicated that the question does not apply to their work	
the authors did not respond to the checkbox question	
For background on the checklist and guidance provided to the authors, see the page at ACL Rolling Review.	Responsible NLP Checklist

✓ A. Questions mandatory for all submissions.

- ✓ A1. Did you describe the limitations of your work? *This paper has a Limitations section.*
- A2. Did you discuss any potential risks of your work?

 In Section Ethics Statement and Broader Impact, we discuss the potential risk, broader impacts, and some potential positive applications of our work.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - B1. Did you cite the creators of artifacts you used?

 In Section 4 Experiments, Appendix C Experimental Details, Appendix H Extension to Mathematical Reasoning, References
 - ☑ B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

 In Section Ethics Statement and Broader Impact, we mention all the datasets that we use are public but do not state their specific license names because they are standard licenses and well-known.
 - B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?
 - In Section Ethics Statement and Broader Impact, we mention that we have confirmed our use of datasets and LLMs aligns with their intended purposes and usage guidelines.
 - B4. Did you discuss the steps taken to check whether the data that was collected/used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect/anonymize it?
 - Because all the datasets in our work are widely-used benchmarks, previous works may have already protected/anonymized sensitive information in the data.
 - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 In Section 4 Experiments, Appendix C Experimental Details, Appendix H Extension to Mathematical Reasoning

☑ B6. Did you report relevant statistics like the number of examples, details of train/test/dev splits, etc. for the data that you used/created? In Appendix C Experimental Details, Appendix H Extension to Mathematical Reasoning **☑** C. Did you run computational experiments? 2 C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? In Section 4 Experiments, Appendix C Experimental Details, Appendix E Number of Reasoning Steps, Generated Tokens, and Inference Time (Table 7), Appendix H Extension to Mathematical Reasoning 2 C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? In Section 4 Experiments, Appendix C Experimental Details, Appendix E Number of Reasoning Steps, Generated Tokens, and Inference Time 2 C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run? In Section 4 Experiments (Table 1, Figure 2), Appendix C Experimental Details, Table 2/5/6/7, Figure 3/4/7/8/9/10/11 2 C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation, such as NLTK, SpaCy, ROUGE, etc.), did you report the implementation, model, and parameter settings In Section 1 Introduction, Appendix A Computation for Semantic Overlap and Cumulative Connectivity, Appendix C Experimental Details **D.** Did you use human annotators (e.g., crowdworkers) or research with human subjects? D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.? (left blank) D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)? (left blank) D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating (e.g., did your instructions explain how the data would be used)? (left blank) D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? (left blank)

E. Did you use AI assistants (e.g., ChatGPT, Copilot) in your research, coding, or writing?

that is the source of the data?

(left blank)

5. Did you report the basic demographic and geographic characteristics of the annotator population

☑ E1. If you used AI assistants, did you include information about their use?

We only use API to access LLM (GPT-40-mini, GPT-40, o1-mini, o3-mini, Llama-3.3 70B, Qwen2.5-7B) to generate outputs for experiments. We include the information in Abstract, Section 1 Introduction, Section 4 Experiments, Appendix C Experimental Details, Appendix H Extension to Mathematical Reasoning.