Responsible NLP Checklist

Paper title: Continuously Steering LLMs Sensitivity to Contextual Knowledge with Proxy Models Authors: Yilin Wang, Heng Wang, Yuyang Bai, Minnan Luo

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	
\Box the authors did not respond to the checkbox question	
For background on the checklist and guidance provided to the authors, see page at ACL Rolling Review.	the 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?

 Our work focuses on the fundamental technical problem of controlling knowledge sensitivity in LLMs.

 While the broader field of generative models has potential societal risks, our paper does not introduce new risk vectors that warrant a dedicated discussion beyond what is already known for LLMs in general.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used? 3.2, 3.4
 - B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

 We did not explicitly discuss the licenses of the pre-existing models and datasets. These artifacts are well-established public benchmarks and open-source models, and their licensing information is available from their original sources, which we have cited.
 - 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)?

 3.2, 3.4
 - 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?
 - We used established public datasets (e.g., MuSiQue, PopQA, DynamicQA) that are standard benchmarks in the NLP community. We did not perform an additional audit for PII or offensive content, operating under the assumption that these datasets have been appropriately curated by their creators.
 - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 3.2, 3.4

☑ 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? 2.1. 3.2. 3.4 **☑** 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? 3.2 C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? 3.2 X 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? limited computational resources 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 used? 2.2, 3.2 **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?

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

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

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

that is the source of the data?

(left blank)

2.2