#### Responsible NLP Checklist

Paper title: TRUEBench: Can LLM Response Meet Real-world Constraints as Productivity Assistant? Authors: Jiho Park, Jongyoon Song, Minjin Choi, Kyuho Heo, Taehun Huh, Ji Won Kim

How to read the checklist symbols:
the authors responded 'yes'
X 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 Responsible NLP Checklist page at ACL Rolling Review.

## ✓ 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? *Potential Rist Section below Limitations Section*
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
  - ☑ B1. Did you cite the creators of artifacts you used? *Appendices B and C.1*
  - B2. Did you discuss the license or terms for use and/or distribution of any artifacts? *Appendix I*
  - 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)?
    - We only use models and datasets to see models' vulnerabilities, e.g., identifying LLMs' strengths and weaknesses.
  - 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?
    - For personal information, we remove the them from data. For offensive data, we construct "Safety" tasks to evaluate LLMs, and discussion about this information can be found in Section 3 and Section "Ethical Statements"
  - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

    Appendices A and C.1
  - ☑ 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?

    Appendix A

## **☑** C. Did you run computational experiments?

- 🗹 C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? Section 4.2 and Appendix H
- ☑ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Sections 4 and 5, and Appendix H
- 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? Sections 4 and5, and Appendices B and C
- 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? Appendix H

#### **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.? In-house domain experts participated for human annotation and we directly instruct them to follow guidelines. We provide the detailed annotation process in Section 3 and Appendix F.
- 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)?
  - *In-house experts participated in the annotation process.*
- 2 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)? Section 3 and Appendix F
- ☑ D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? Section 3 and Appendix F
- 🗹 D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data? Section 3 and Appendix F

# **E.** Did you 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? Appendix J