## **Responsible NLP Checklist**

Paper title: Measuring scalar constructs in social science with LLMs

Authors: Hauke Licht, Rupak Sarkar, Patrick Y. Wu, Pranav Goel, Niklas Stoehr, Elliott Ash, Alexander Miserlis Hovle

rus Hoyee
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

- ✓ 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? (left blank)
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
  - ☑ B1. Did you cite the creators of artifacts you used? Throughout the Text
  - **X** B2. Did you discuss the license or terms for use and/or distribution of any artifacts? No; licenses are all listed on harvard dataverse, which makes data available for the puroses of reproducibility.
  - 🛮 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)?
    - They are datasets created for computational text analysis with reproducibility in mind.
  - 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? (left blank)
  - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.? Section 3.1
- ☑ 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? Section 3.1

## ☑ 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?

  Appendix B.1, sizes are reported in Section 3.2
- ✓ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

  Section 3 and Appendix B
- ✓ 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?

We report standard deviations (bootstrapped or 5-fold) in the main results tables (Table 2 and Table 3)

✓ 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?

We specify the Bradley-Terry implementation in Footnote 14

## **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.?

  The annotators were provided the first prompt shown in Section E with additional instructions on
  - The annotators were provided the first prompt shown in Section E with additional instructions on how to use the online annotation tool
- 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)?
  - For the analyses reported in Section B (Appendix), we relied on two members of the authors team, who volunteered, and two research assistants employed at the Chair of Law, Economics, and Data Science (ETH Zurich), who's hourly rate meets Swiss minimum wage regulations.
- 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?

E1. If you used AI assistants, did you include information about their use? (left blank)