

## Responsible NLP Checklist

Paper title: *Reading Between the Prompts: How Stereotypes Shape LLMs Implicit Personalization*

Authors: Vera Neplenbroek, Arianna Bisazza, Raquel Fernandez

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 [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?

*We discuss potential risks of our work in the Ethical Considerations section.*

### B. Did you use or create scientific artifacts? (e.g. code, datasets, models)

- B1. Did you cite the creators of artifacts you used?

*We cite the creators of the dataset we used in section 3.1 and Appendix A.1, and the creators of models we used in section 3.2 and Appendix B.*

- B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

*We discuss the license or terms for use of the dataset we used in section 3.1 and Appendix A.1, and those of the models we used in Appendix B.*

- 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 mention that our use of the existing dataset we use is consistent with its intended use in Appendix A.1.*

- 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 describe the collection of potentially offensive stereotypes collected for the purpose of investigating those in LLMs in section 3.1 and Appendix A.1. We do not use any automatically collected data and can therefore guarantee that there is no personally identifying information in our data.*

- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

*We discuss details of the artifacts in Section 3.1 and appendix A (data) and Section 3.2 and appendix B (models).*

- 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?

*We report relevant statistics for our data in Section 3.*

**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?

*We report on the number of parameters in the models used and the total computational budget and infrastructure used in Section 3.2 and Appendix B.*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

*We discuss the details of our experimental setup in Sections 3.3, 4 and 5.*

- 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 our results with descriptive statistics in Section 4.*

- 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 report how we access open-source models in Appendix B.*

**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)*

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*(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)*