

Responsible NLP Checklist

Paper title: *How Adversarial Environments Mislead Agentic AI?*

Authors: *Zhonghao Zhan, Huichi Zhou, Zhenhao Li, Peiyuan Jing, Krinos Li, Hamed Haddadi*

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?
Ethical Considerations (after Limitations).

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

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?
The POTEKIN-S2 dataset is derived from a snapshot of Semantic Scholar, containing public bibliographic information (paper titles, author names, abstracts, references) rather than private user data. The POTEKIN-Phantoms dataset consists of fake papers generated by Large Language Models, containing no real persons or offensive content. The POTEKIN-Claims dataset is adapted from the existing AVeriTeC dataset (Schlichtkrull et al., 2023), whose release already underwent the original authors' content review.

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 3.6 (Datasets) and Appendix Table A2.

C. Did you run computational experiments?

C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
Sections 3.2 (Agents Under Test) and 3.7 (Experiment Design), plus Appendix Table A3.

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?
Section 4 (Results and Discussion), with per-agent means and ranges reported in Tables 3-5 and Appendix Tables A4, A6, A7.

The Responsible NLP Checklist used at ACL Rolling Review is adopted from NAACL 2022, with the addition of ACL 2023 question on AI writing assistance and further refinements based on ARR practice. ACL 2026 used a subset of ARR checklist form.

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?

E1. If you used AI assistants, did you include information about their use?
Section 3.5 (Adversarial Resources and Red Team) and Appendix Table A1 document the use of Gemini 2.5 Pro/Flash for generating adversarial content (poisoned snippets, persona variants, fake citations) as part of the experimental methodology. Claude was additionally used for writing polish and code/engineering debugging. All AI-generated content was reviewed and verified by the authors.