

Responsible NLP Checklist

Paper title: *RLShield: Dynamic Jailbreak Detection for LLMs via Reinforced Adaptive Learning*

Authors: *Zhao Tong, Pengfei Yang, Yimeng Gu, Haichao Shi, Qiang Liu, Xingcheng Xu, Shu Wu, Xiao-Yu Zhang*

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?

Because we propose a Jailbreak Detection for LLMs, we believe our method does not introduce potential risks.

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?

Section 4.1

- 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.2

C. Did you run computational experiments?

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

Section 4 and Appendix A.3

- 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 5

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.?

We did not recruit or instruct any human participants/annotators. Our experiments rely exclusively on publicly available datasets, and we did not create new human-annotated data.

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.

- 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)?

We did not recruit or pay any human participants/annotators. All data used are publicly available datasets, and no new data collection involving human subjects was conducted.

- 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)?

Our work relies exclusively on publicly available datasets and does not involve creating a new dataset, recruiting participants, or conducting any human annotation.

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

This work did not involve new data collection or interaction with human participants. We only used publicly available datasets; therefore, ethics board approval was not required.

- 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)