

Responsible NLP Checklist

Paper title: *Activation Decomposition and Steering for LLM Backdoor Remediation*

Authors: *Lingfeng Zhong, Qiongfai Xu, Usman Naseem*

How to read the checklist symbols:

- the authors responded 'yes'
- the authors responded 'no'
- N/A 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?

The Ethical Consideration section discusses potential risks of our efforts.

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?

A clear warning is put on the front page below the Abstract part: "This paper contains instances of abusive language generated by intentionally backdoored large language models. Please proceed with caution."

- 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 4.1 discusses details of train/test/dev splits.

C. Did you run computational experiments?

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

Section 4.1 and Section A elaborate on experimental setup and hyperparameters used in our experiments.

- 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 4.4 and section E present visual summaries about our results.

D. Did you use human annotators (e.g., crowdworkers) or research with human subjects?

- N/A 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)

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.

N/A 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)?
Section I elaborates on all licenses for the models, tools and open-source datasets created by third-party developers.

N/A 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 I elaborates on the use of AI assistants. We used ChatGPT as the AI assistant to improve the clarity and grammatical correctness of the manuscript. The tool was not used to generate original scientific ideas, conduct experiments, or draw conclusions. The authors take all responsibility for the accuracy, integrity, and originality of the work.