## Responsible NLP Checklist

Paper title: FIER: Fine-Grained and Efficient KV Cache Retrieval for Long-context LLM Inference Authors: Dongwei Wang, Zijie Liu, Song Wang, Yuxin Ren, Jianing Deng, Jingtong Hu, Tianlong Chen, Huanrui Yang

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?

  This work does not involve user-facing systems or sensitive data, and thus poses minimal risk.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
  - ☑ B1. Did you cite the creators of artifacts you used? Section 1, Section 2, Section 4.1
  - B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

    We only used publicly available models and datasets with standard research licenses, but did not explicitly state this in the paper.
  - 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 used all artifacts strictly for research purposes in line with their intended use, but did not explicitly mention this in the paper.
  - 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 used only publicly available benchmark datasets, which do not contain personally identifying information or offensive content.
  - B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

    We used standard public datasets like PG19 and focused on long-context modeling. While we briefly described dataset characteristics (e.g., token length), we did not provide detailed documentation of linguistic or demographic coverage.

•	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? <i>Appendix A, 4.3.1</i>
<b>√</b>	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? <i>Appendix B</i>
•	C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? <i>Section 4.1, Appendix B</i>
•	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 4.3, 4.4, 4.5
•	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?  Section 4.1, 4.5, Appendix B
X	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.?  This work does not involve any human participants or annotation tasks.
_	· · · ·
<u>N//</u>	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)?  This work does not involve any human participants or compensation.
N/A	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)?  We did not collect or annotate data; all data comes from publicly released datasets.
N/A	D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? No new data was collected; we only used publicly available datasets that do not require ethics review
N/A	D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

## lacktriangledown E. Did you use AI assistants (e.g., ChatGPT, Copilot) in your research, coding, or writing?

We did not collect or annotate data; all data comes from publicly released datasets.

**☒** E1. If you used AI assistants, did you include information about their use? We used AI assistants for minor writing suggestions, but did not include this information in the paper.