Responsible NLP Checklist

Paper title: Efficient Dynamic Clustering-Based Document Compression for Retrieval-Augmented-Generation

Authors: Weitao Li, Xiangyu Zhang, Kaiming Liu, Xuanyu Lei, Weizhi Ma, Yang Liu

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?

 No, because the work involves standard research methods and publicly available datasets, and there are no foreseeable risks of misuse or harmful impact.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used? *References Section*
 - B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

 No, because all artifacts referenced in the paper are standard, publicly available resources commonly used in research, so no separate license or terms needed to be discussed.
 - 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)?
 - No, because all artifacts used in this work are standard, publicly available research resources, and their use in this study aligns with typical research purposes, so no additional discussion was necessar
 - 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?
 - No, because all data used in this work are standard, publicly available datasets that do not contain personally identifying information or offensive content, so no special anonymization or protection steps were required.
 - B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 No, because all artifacts used in this work are standard, well-documented public resources, so additional documentation was not necessary.

☑ 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 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? No, because only small-scale inference and retrieval experiments were conducted, requiring minimal computational resources, so reporting model size and total budget was not necessary. 2 C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Section 5.1, 6.4. Appendix B. 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? 5.1 2 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? 5.1. 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

- 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?

No, because AI assistants were not used for any creative or substantive part of the research; only standard writing and coding tools were used.