Responsible NLP Checklist

Paper title: Agent-as-Judge for Factual Summarization of Long Narratives Authors: Yeonseok Jeong, Minsoo Kim, seung-won hwang, Byung-Hak Kim

	_
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? *Section 7*
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used? Section 3.1, 4.4
 - B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

 We used only publicly available datasets and models released for research use under standard licenses.
 - 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 did not explicitly discuss alignment of our use of external datasets or tools with their intended use, nor specify intended use constraints for our created
 - 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
 - B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 We used well-documented public benchmarks and models. Detailed descriptions are available in their original sources.
 - ☑ 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 4.3

☑ 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?

 Section 4.1
- ☑ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

We did not perform or report any hyperparameter search or list best-found hyperparameter values, because our experiments relied solely on a single off-the-shelf LLM (gpt-4o-mini) with fixed prompt templates and chunk-size settings rather than tuning model-specific parameters.

- 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?
 - We only report singlerun aggregate metrics (e.g., average factuality scores) without error bars, standard deviations, or clarification of max/meanour evaluations were deterministic and not repeated.
- 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?

We did not detail the specific packages, models, or parameter settings used, as we relied on default configurations of common libraries (e.g., NLTK for tokenization, ROUGE scripts); these details are available in our public code repository.

\(\begin{aligned} \Bigsilon \) 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 use human annotators
- 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 use human annotators
- 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 use human annotators
- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board? We did not use human annotators
- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

 We did not use human annotators
- **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 D*