## Responsible NLP Checklist

Paper title: Towards Multi-Document Question Answering in Scientific Literature: Pipeline, Dataset, and Evaluation Authors: Hui Huang, Julien Velcin, Yacine Kessaci How to read the checklist symbols: the authors responded 'yes' X 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? (left blank) **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models) ☑ B1. Did you cite the creators of artifacts you used? Section 4 (Dataset Generation and Analysis) ☑ B2. Did you discuss the license or terms for use and/or distribution of any artifacts? We discuss the source of the data and its open-access nature in Section 4.1. SPIQA dataset is under CC-BY 4.0 license. We will also release our MDA-QA dataset with an open-access license. ☑ 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)? Section 4.1 ■ 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? Our dataset is composed of open-access academic papers collected from public repositories and does not contain any 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.? 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? Section 4.3

## ☑ C. Did you run computational experiments?

process.

- ∠C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?
  - We did not report the model size (number of parameters) or the computational budget because we relied on provided large language model APIs (e.g., Claude, GPT-40) via standard commercial access.
- ✓ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

  Annex C
- ✓ 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 5.2, 5.3, 5.4
- ✓ 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 5.1

## **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 provide the complete instructions in the paper. Our work only involved a limited, expert-based validation of automatically generated QA pairs, rather than a formal large-scale annotation
- 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)? (*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?

We used gpt-40 for grammar check and text polish throughout the draft. The final version was reviewed by all authors.