Responsible NLP Checklist

Paper title: FoodSafeSum: Enabling Natural Language Processing Applications for Food Safety Document Summarization and Analysis

Authors: Juli Bakagianni, Korbinian Randl, Guido Rocchietti, Cosimo Rulli, Franco Maria Nardini, Salvatore Trani, Aron Henriksson, Anna Romanova, John Pavlopoulos

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 Chec page at ACL Rolling Review.	klist

✓ 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? *Limitations Section*
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used? 3. 5

 - 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)?

 Our work is consistent with research purposes as of the models we used, so it is implied.
 - 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?
 - The dataset is derived from publicly available institutional sources, including regulatory agencies, government portals, scientific journals, and news outlets. As such, it does not contain personal data or uniquely identifying information about individuals. We verified that no documents include sensitive health records or personal identifiers.
 - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

■ 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?

☑ 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?
 - We report the number of parameters for the models used (Section 5). Summaries were generated using Meta LLaMA-3 70B instruct (~70B parameters) via Amazon Bedrock. The dataset is relatively small (2,091 documents), and all experiments were conducted at modest scale. Consequently, the total computational budget was lowused primarily for generating Stella 1.5B embeddingsand no specialized infrastructure beyond standard cloud compute resources was required.
- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
 - We used standard baseline models with default hyperparameters. No hyperparameter search was performed, as our experiments primarily focus on evaluating the utility of LLM-generated and human-written summaries for downstream NLP tasks (topic classification, retrieval, and RAG-based QA).
- 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?
- ✓ 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?

 4

☑ 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 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? *Ai assistant was used to correct typos*.