Responsible NLP Checklist

Paper title: Tree-of-Quote Prompting Improves Factuality and Attribution in Multi-Hop and Medical Reasoning

Authors: Justin Xu, Yiming Li, Zizheng Zhang, Augustine Yui Hei Luk, Mayank Jobanputra, Samarth Oza, Ashley Murray, Meghana Reddy Kasula, Andrew Parker, David W Eyre

How to read the checklist symbols:	
the authors respo	onded 'yes'
X the authors respo	onded 'no'
N/A the authors indic	ated that the question does not apply to their work
the authors did n	ot respond to the checkbox question
For background on page at ACL Rolling R	the checklist and guidance provided to the authors, see the Responsible NLP Checklist eview.

✓ 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?

 Correlation interpretation risks mentioned in Section 5.2.2; Biases mentioned in Section 6; Demonstrated risk of evaluation metrics failing in Section 6.1 and Limitations

B. Did you use or create scientific artifacts? (e.g. code, datasets, models)

- B1. Did you cite the creators of artifacts you used?

 Throughout but primarily in Section 4 where models, datasets, methods, and metrics are introduced
- B2. Did you discuss the license or terms for use and/or distribution of any artifacts? *All publicly available artifacts with open 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)?

 All publicly available artifacts with open licenses
- 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 identifiable information in datasets
- ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 Described overall dataset characteristics in 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?

 Described dataset splits and sample counts in Section 4.1; Reader study setups described in Section

Described dataset splits and sample counts in Section 4.1; Reader study setups described in Section 5.2.1

	C. Did you run computational experiments?
N/.	C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? Models were run using API, but specific models are named in Section 4 and 6.2
	,
V	C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?
	Experimental setup described in Section 4.1
•	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?
	Confidence intervals and performance ranges are presented throughout all tables
•	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?
	Metrics used and their details described in Section 4.2
X	D. Did you use human annotators (e.g., crowdworkers) or research with human subjects?
N/.	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.? (<i>left blank</i>)
N/	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)
N/.	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)?

D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

D5. Did you report the basic demographic and geographic characteristics of the annotator population

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

API calls and LLM evaluations are detailed throughout, but primarily in Sections 4.1 and 5.1

☑ E1. If you used AI assistants, did you include information about their use?

(left blank)

(left blank)

that is the source of the data?