Responsible NLP Checklist

Paper title: CDT: A Comprehensive Capability Framework for Large Language Models Across Cognition, Domain, and Task

Authors: Haosi Mo, Xinyu Ma, Xuebo Liu, Derek F. Wong, YU LI, Jie Liu, Min Zhang

| How to read the checklist symbols: | |
|-------------------------------------------------------------------------------------------------------------------------------------|--------|
| the authors responded 'yes' | |
| the authors responded 'no' | |
| N/A 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. | t _ |

✓ 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? *Our method does not involve risks*.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
 - ☑ B1. Did you cite the creators of artifacts you used?

 Section 5 Empirical Analysis of Instruction Dataset Capabilities, Section 6.1 Experiment Setup, A.4

 Capability Tagging Details
 - B2. Did you discuss the license or terms for use and/or distribution of any artifacts? All the models and data we used in our work are licensed for research use.
 - 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 the models and data we used in our work are licensed for research use.
 - 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?
 - We use only publicly available datasets that have already undergone the necessary processing.
 - B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

 We use only publicly available datasets that have already undergone the necessary processing.
 - 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 3.2 Capability Tagging Model Training Section 6.1 Experiment Setup Section 6.3 Experi-

Section 3.2 Capability Tagging Model Training, Section 6.1 Experiment Setup, Section 6.3 Experiments on the Specific Scenario, A.4 Capability Tagging Details

| \checkmark | C. Did you run computational experiments? |
|--------------|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| V | 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 6.1 Experiment Setup, A.6 Experiments on Mistral Model |
| Ø | C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? Section 6.1 Experiment Setup |
| | 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 6.2 Experiments on the General Scenario, Section 6.3 Experiments on the Specific Scenario |
| | 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 6 Data Selection Experiments |
| X | D. Did you use human annotators (e.g., crowdworkers) or research with human subjects? |
| N/A | 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/A | 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)? (<i>left blank</i>) |
| N/A | 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?

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

We use AI assistants to polish and correct grammar.

(left blank)

(left blank)

(left blank)

that is the source of the data?