#### Responsible NLP Checklist

Paper title: On the Fine-Grained Planning Abilities of VLM Web Agents
Authors: Surgan Jandial, Yinong Oliver Wang, Andrea Bajcsy, Fernando De la Torre

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 Checkli page at ACL Rolling Review.	st

## ✓ 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 paper does not involve new methods/tools. We only design new evaluations of existing models.
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
  - ☑ B1. Did you cite the creators of artifacts you used? *Section 2, 3, 4, 5*
  - B2. Did you discuss the license or terms for use and/or distribution of any artifacts? We rely on existing open-sourced datasets and models for academic purpose only. No part of this study is related to any commercial activity and does not violate their existing 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)?

  Section 2, 3, 4, 5
- 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 repurposed two published real-web datasets (MM-M2W, GUIACT) which are open-sourced and
  - we are not aware of any harmful content in them. We also generate synthetic dataset on web forms where no content about people is involved.
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

  The artifacts are either open-sourced with known documentation or confined to very specific setups where we explain exactly how each data sample is made of.
- 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?

  Table 6 in Appendix

# ☑ 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 3
- ☑ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

  Section 3
- 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 4, 5
- 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 didn't perform preprocessing, normalization, and only used accuracy for evaluation.

# **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.?

  Appendix G
- 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)?

Participants are school peer volunteers and no recruitment or payment is involved.

- 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? (left blank)