

Responsible NLP Checklist

Paper title: *Hierarchical Visual Agent: Managing Contexts in Joint Image-Text Space for Advanced Chart Reasoning*

Authors: *Qihua Dong, Ruozhen He, Junwen Chen, Yizhou Wang, Xu Ma, Songyao Jiang, Yun Fu*

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.

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?

Section 9 (Ethics Statement) discusses potential misuse risks.

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

- ^{N/A} 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?

(left blank)

- 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?

Sections 3.2 and 4.1 report dataset sizes and splits for ChartQA, synthetic charts, and CharXiv reasoning subset.

C. Did you run computational experiments?

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

Sections 3.3, 4.2, and Appendices A, B describe the backbone model (Qwen3VL-A22B), baseline prompts, tool/skill configuration, and the HierVA task schema.

- 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?

We report single-run accuracy on fixed evaluation splits, following standard practice for inference-time orchestration work on closed/proprietary backbones where run-to-run variance from deterministic decoding is small. Benchmark sizes are specified in Sections 3.2 and 4.1.

The Responsible NLP Checklist used at ACL Rolling Review is adopted from NAACL 2022, with the addition of ACL 2023 question on AI writing assistance and further refinements based on ARR practice. ACL 2026 used a subset of ARR checklist form.

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)

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?

Section 10 (Use of Large Language Models in Writing) describes the limited use of LLMs for language polishing and clarity improvements.