

## Responsible NLP Checklist

Paper title: *Decoding Scientific Experimental Images: The SPUR Benchmark for Perception, Understanding, and Reasoning*

Authors: *Junpeng Ding, Zichen Tang, Haihong E, Mengyuan Ji, Yang Liu, Haolin Tian, Haiyang Sun, Pengqi Sun, Yang Xu, Yichen Liu, Haocheng Gao, Zijie Xi, Ruomeng Jiang, Peizhi Zhao, Rongjin Li, Yuanze Li, Jiacheng Liu, Zhongjun Yang, Jintong Chen, Siying Lin*

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 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?

*Ethical Considerations*

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

- 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?

*Section 2, Ethical Considerations, Appendix B*

- 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 2, Appendix B*

### C. Did you run computational experiments?

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

*Section 3, Appendix C*

- 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 3, Appendix C*

### 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.?

*Section 2, Appendix B*

*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.*

- 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)?

*This study did not involve the recruitment of external participants or the use of crowdsourcing platforms. Data verification and annotation were conducted by the research team and internal experts as part of their standard research activities.*

- 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)?

*The data used in this study consists of publicly available scientific images sourced from open-source repositories. All data are used in strict accordance with their original open-source licenses (e.g., CC BY 4.0), which grant permission for research use and redistribution.*

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

*This study utilizes only publicly available open-source data and does not involve any human participants or animal experimentation.*

- 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?

*Appendix D*