## Responsible NLP Checklist

Paper title: LLMs Behind the Scenes: Enabling Narrative Scene Illustration

Authors: Melissa Roemmele, John Joon Young Chung, Taewook Kim, Yuqian Sun, Alex Calderwood, Max Kreminski

How to	read the checklist symbols:
<b>d</b> th	ne authors responded 'yes'
X th	ne authors responded 'no'
N/A th	ne authors indicated that the question does not apply to their work
☐ th	ne authors did not respond to the checkbox question
	background on the checklist and guidance provided to the authors, see the Responsible NLP Checklist 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" section (final section of paper)
- **B.** Did you use or create scientific artifacts? (e.g. code, datasets, models)
  - ☑ B1. Did you cite the creators of artifacts you used? Section 4
  - ☑ B2. Did you discuss the license or terms for use and/or distribution of any artifacts? *Section 1 (footnote on page 2)*
  - ☑ 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 4
  - ☑ 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?
    - "Ethical Considerations" section (final section of paper)
  - ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?

    Section 4
- 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 4 and Appendix 1

C. Did	von ri	ın com	nutational	experiments	?
 C. Diu	your	m com	putanonai	CAPCI IIIICIIC	•

- 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 and 4 (in particular, see footnotes regarding API services for each model)
- ☑ C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

  Section 5
- 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 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?

  Section 4

## **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 4 and Appendix Figure 3
- ✓ 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)?

  Section 4
- ☑ 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)?

  Section 4 (the instructions did not explain how the data would be used, but participants consented to and were compensated for the annotation task)
- 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?

  We did not collect these characteristics because they were not relevant to the annotation task.

## **Z** 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*)