## Responsible NLP Checklist

Paper title: Does It Run and Is That Enough? Revisiting Text-to-Chart Generation with a Multi-Agent Approach

Authors: James Ford, Anthony Rios

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?

  See Limitations Section (Limitations). We discuss risks related to generalization of synthetic datasets and accessibility concerns.

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

- ☑ B1. Did you cite the creators of artifacts you used?

  See 3 (Data). We cite the creators of Text2Chart31, ChartX, and VisText where each dataset is introduced.
- B2. Did you discuss the license or terms for use and/or distribution of any artifacts?

  See 3 (Data). We cite the creators of Text2Chart31, ChartX, and VisText where each dataset is introduced.
- ☑ 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)?
  - See 3 (Data). We use Text2Chart31, ChartX, and VisText strictly for research evaluation, consistent with their intended use as benchmark datasets.
- ☑ 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?
  - The datasets contain no PII or offensive content. Text2Chart31 and ChartX are synthetic, while VisText uses automatically generated captions from public Statista data (3).
- ☑ B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?
  - See 3 (Data). We describe the domains, chart types, and sources of each dataset, including statistics in Table 1.

☑ 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? See 3 (Data). We report number of examples, chart categories, and test set sizes for all datasets. Tables 1 summarize dataset statistics and splits. **☑** C. Did you run computational experiments? 2 C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used? See 5 (Results). We report model sizes (e.g., GPT-4o-mini, Llama 3.1-8B, Gemma-3-12B, Owen-32B) and the compute budget (~\$150 total cost). 2 C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values? See 4 (Methodology). We describe the drafting and repair prompts, zero-shot vs. few-shot settings, and iteration limits (up to three). 2 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? See 5 (Results). We report error rates with confidence from multiple runs and provide averages for image quality scores (SSIM and multimodal LLM judge). 2 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 See 4 (Methodology). We specify that all code runs used Python with Matplotlib, Pandas, and Numpy. **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) \times 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?

that is the source of the data?

(left blank)

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

See Limitations Section. We note that AI assistance (LLMs) was used to support manuscript writing.

D5. Did you report the basic demographic and geographic characteristics of the annotator population