
Into the LAION's Den: Investigating Hate in Multimodal Datasets

Anonymous Author(s)

Affiliation

Address

email

Abstract

1

2 Checklist

3

1. For all authors...

4

(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#)

5

6

(b) Did you describe the limitations of your work? [\[Yes\]](#)

7

(c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#)

8

(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#)

9

10

2. If you are including theoretical results...

11

(a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)

12

(b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)

13

3. If you ran experiments (e.g. for benchmarks)...

14

(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#)

15

16

(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)

17

18

(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#)

19

20

(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#)

21

22

4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...

23

(a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)

24

(b) Did you mention the license of the assets? [\[N/A\]](#)

25

(c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)

26

(d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [\[Yes\]](#)

27

28

(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[Yes\]](#)

29

- 30 5. If you used crowdsourcing or conducted research with human subjects...
- 31 (a) Did you include the full text of instructions given to participants and screenshots, if
- 32 applicable? [N/A]
- 33 (b) Did you describe any potential participant risks, with links to Institutional Review
- 34 Board (IRB) approvals, if applicable? [N/A]
- 35 (c) Did you include the estimated hourly wage paid to participants and the total amount
- 36 spent on participant compensation? [N/A]

37 **A Appendix**

38 Include extra information in the appendix. This section will often be part of the supplemental material.
39 Please see the call on the NeurIPS website for links to additional guides on dataset publication.

- 40 1. Submission introducing new datasets must include the following in the supplementary
- 41 materials:
- 42 (a) Dataset documentation and intended uses. Recommended documentation frameworks
- 43 include datasheets for datasets, dataset nutrition labels, data statements for NLP, and
- 44 accountability frameworks.
- 45 (b) URL to website/platform where the dataset/benchmark can be viewed and downloaded
- 46 by the reviewers.
- 47 (c) Author statement that they bear all responsibility in case of violation of rights, etc., and
- 48 confirmation of the data license.
- 49 (d) Hosting, licensing, and maintenance plan. The choice of hosting platform is yours, as
- 50 long as you ensure access to the data (possibly through a curated interface) and will
- 51 provide the necessary maintenance.
- 52 2. To ensure accessibility, the supplementary materials for datasets must include the following:
- 53 (a) Links to access the dataset and its metadata. This can be hidden upon submission if the
- 54 dataset is not yet publicly available but must be added in the camera-ready version. In
- 55 select cases, e.g. when the data can only be released at a later date, this can be added
- 56 afterward. Simulation environments should link to (open source) code repositories.
- 57 (b) The dataset itself should ideally use an open and widely used data format. Provide a
- 58 detailed explanation on how the dataset can be read. For simulation environments, use
- 59 existing frameworks or explain how they can be used.
- 60 (c) Long-term preservation: It must be clear that the dataset will be available for a long time,
- 61 either by uploading to a data repository or by explaining how the authors themselves
- 62 will ensure this.
- 63 (d) Explicit license: Authors must choose a license, ideally a CC license for datasets, or an
- 64 open source license for code (e.g. RL environments).
- 65 (e) Add structured metadata to a dataset's meta-data page using Web standards (like
- 66 schema.org and DCAT): This allows it to be discovered and organized by anyone. If
- 67 you use an existing data repository, this is often done automatically.
- 68 (f) Highly recommended: a persistent dereferenceable identifier (e.g. a DOI minted by
- 69 a data repository or a prefix on identifiers.org) for datasets, or a code repository (e.g.
- 70 GitHub, GitLab,...) for code. If this is not possible or useful, please explain why.
- 71 3. For benchmarks, the supplementary materials must ensure that all results are easily repro-
- 72 ducible. Where possible, use a reproducibility framework such as the ML reproducibility
- 73 checklist, or otherwise guarantee that all results can be easily reproduced, i.e. all necessary
- 74 datasets, code, and evaluation procedures must be accessible and documented.
- 75 4. For papers introducing best practices in creating or curating datasets and benchmarks, the
- 76 above supplementary materials are not required.