APPENDIX

### A. Preparing annotation tasks for crowdsourcing

The additional details for each context ensured that the pre-processed dataset included only citation contexts clearly linked to a specific cited paper, minimizing ambiguity for crowdsourcing workers. In cases where citation marks were not uniquely identifiable, we appended last name of the first author of the cited paper to each labeling item. Each entity to be labeled was composed of four parts:

- **Citation context**: The sentence(s) which contains the citation mark of the cited paper.
- **Related rep-study**: The reproducibility study corresponding to the citation context. We added this to identify the related rep-study to which the citation context belongs even after crowdsourced annotations process.
- **First author of the cited paper**: To resolve situations where multiple citation marks appear in the same context.
- **Context index**: An id to identify the specific citation context.

We organized the citation context entities in the data files according to the following format (see Fig. 3 for an example):

```
[Citation context] [Related rep-study] [First
author of the cited paper] [Context index]
```

The inclusion of the reproducibility study and context index made post-processing and analysis easier after crowdsourcing.

### B. Interface Design

We designed a custom template for the labeling task GUI using *Mechanical Turk Crowd HTML Elements*[1], which abstracts HTML markup, CSS, and JavaScript functionality into HTML tags. This allowed for an efficient and visually consistent interface. As shown in Fig. 3, the GUI included a section containing detailed instructions for workers, along with examples for each label category. We provided definitions for key scientific terms, e.g., *citation*, *citation context*, and *cited paper*. The GUI iteratively displayed each citation context each with a label selection area on the right. Workers were required to choose a label from three options (*Positive*, *Negative*, *Neutral*) and click the *Submit* button to proceed to the next context. The instructions, examples, and definitions were consistently presented for every task to reinforce understanding and minimize errors. It took approximately 27 days to complete the labeling. We used *VS Code* as the Integrated Development Environment (IDE) to develop the template, and the *Mechanical Turk Developer Sandbox*[2] (a simulated environment for testing applications and tasks prior to publication) to test the GUI's functionality before deployment.

### C. Technical Validation

The confusion matrix in Fig. 1 highlights the alignment and discrepancies between the crowdsourced annotations and our manually annotated labels for the randomly selected 244 citation contexts verification set.

[1]https://docs.aws.amazon.com/AWSMechTurk/latest/AWSMturkAPI/ApiReference_HTMLCustomElementsArticle.html

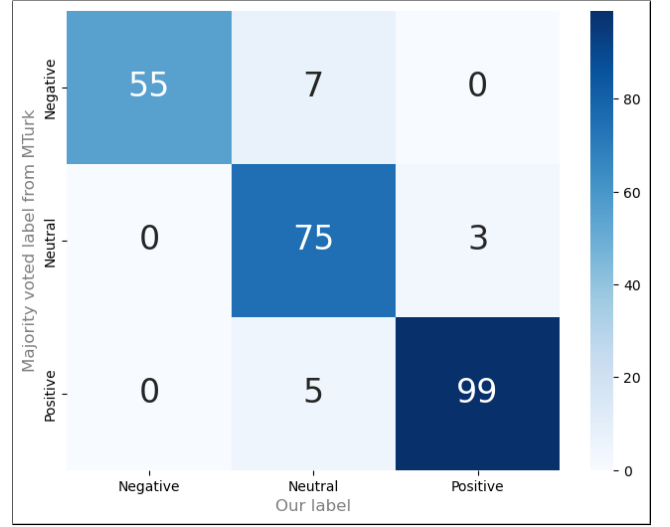[2]https://docs.aws.amazon.com/AWSMechTurk/latest/AWSMechanicalTurkRequester/mturk-use-sandbox.html



Fig. 1. Confusion matrix comparing MTurk annotator labels with ground truth labels in the verification set.

### D. LLM Fine-Tuning and RAG Configurations for CC30k

To showcase the utility of the CC30k dataset, we experimented with multiple large language models using Low-Rank Adaptation (LoRA)–based parameter-efficient fine-tuning and retrieval-augmented generation (RAG). The corresponding hyper-parameters and configurations are listed in Tables I and II.

TABLE I
LORA FINE-TUNING HYPERPARAMETERS FOR LLAMA 3–8B AND QWEN1.5–7B

| Setting | LLaMA 3–8B | Qwen1.5–7B |
|---|---|---|
| LoRA $r$ | 16 | 16 |
| LoRA $\alpha$ | 32 | 32 |
| Dropout | 0.05 | 0.05 |
| Batch size | 4 | 4 |
| Grad. accumulation | 4 | 4 |
| Epochs | 15 | 15 |
| Learning rate | $2 \times 10^{-4}$ | $2 \times 10^{-4}$ |
| Max sequence length | 1024 | 1024 |
| Precision | `fp16` | `fp16` |
| Log steps | 10 | 10 |
| Save strategy | epoch | epoch |
| Save limit | 2 | 2 |

TABLE II
GPT-4O LORA FINE-TUNING AND RAG CONFIGURATION

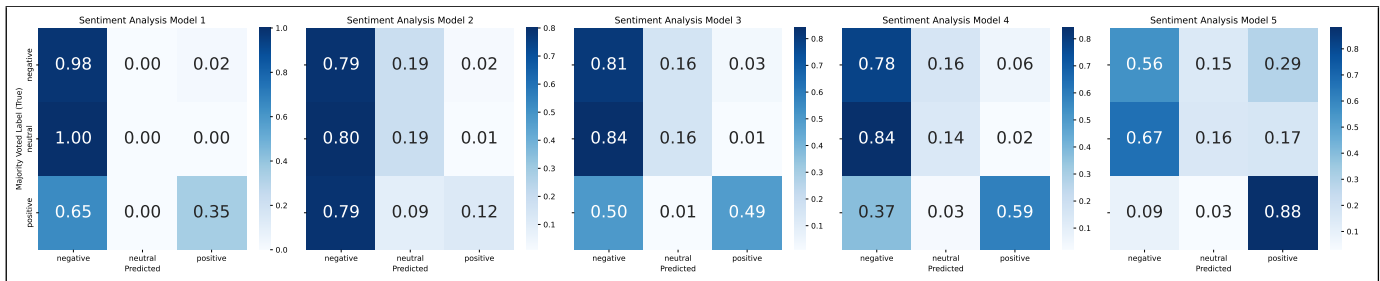| Setting | Value |
|---|---|
| Retrieval method | FAISS |
| Embedding model | SentenceTransformer |
| Top-$K$ | 5 |
| Inference temperature | 0 |
| Max tokens | 10 |

Fig. 2. Confusion matrices of five sentiment classification (3-class) models compared to MTurk annotator ground truth labels



Fig. 3. Crowdsourcing task interface using Crowd HTML elements.