

OFLMEval: An Open Source Language Model with Better Finegrained Human Alignment

Anonymous ACL submission

Abstract

Human evaluation of large language models (LLMs) is a standard practice to measure various aspects such as relevance, fluency, coherence, and overall quality. However, this process is costly, time-consuming, and challenging to reproduce. Current solutions like GPTEval and GPTScore rely on closed-source models such as GPT-4, which present issues like data exposure, lack of transparency, and unpredictability due to API changes. This work proposes OFLMEval, an open-source language model designed for fine-grained evaluation of LLMs. Our model integrates pairwise ranking and single-answer grading to deliver a robust and transparent evaluation process. We fine-tuned our model using the SummEval dataset, which provides human ratings on fluency, coherence, consistency, and relevance, and employed LoRA (Low-Rank Adaptation) techniques for efficient model adaptation. Our results demonstrate significant improvements in the evaluation performance of smaller open LLMs, surpassing existing benchmarks in both pairwise and single-answer grading tasks. This research advances the state-of-the-art in LLM evaluation, offering an open-source alternative with better alignment to human judgment.

1 Introduction

The rapid advancement of large language models (LLMs) has revolutionized natural language processing (NLP) tasks, enabling remarkable capabilities in generating coherent, fluent, and contextually relevant responses. Evaluating these models' performance, however, remains a complex challenge. Traditional human evaluations, which assess criteria such as relevance, fluency, and coherence, are considered the gold standard but come with significant limitations. They are costly, time-consuming, and difficult to reproduce, creating a need for more efficient, consistent, and scalable evaluation methods.

Existing automated evaluation methods, such as GPTEval and GPTScore, have attempted to address these challenges by using LLMs like GPT-4 as evaluators. However, these approaches are not without their drawbacks. Relying on closed-source models raises concerns about data exposure, lack of transparency, and unpredictability due to potential API changes. Furthermore, closed-source models offer limited control and can be prohibitively expensive, hindering their widespread adoption in academic and industry research.

To bridge this gap, we introduce OFLMEval, an open-source language model specifically designed for fine-grained evaluation of other language models. Our goal is to build a robust evaluation system that can provide transparent, reproducible, and cost-effective assessments while aligning closely with human judgment. OFLMEval employs pairwise ranking and single-answer grading, leveraging the strengths of both evaluation types to deliver a comprehensive analysis of model outputs.

Our approach involves fine-tuning the model using the SummEval dataset, which offers human ratings on multiple aspects such as fluency, coherence, consistency, and relevance. By implementing Low-Rank Adaptation (LoRA) techniques, we efficiently adapt the model to perform detailed assessments, demonstrating improved performance over existing benchmarks. The evaluation results reveal that OFLMEval can serve as a reliable, open-source alternative for LLM evaluation, providing insights that are crucial for the ongoing development and refinement of language models.

In this paper, we present the methodology behind OFLMEval, including the dataset preparation, model training, and evaluation criteria. We also provide a comprehensive analysis of our results, highlighting the advantages of using open-source models for fine-grained evaluation tasks. Our contributions aim to advance the current state of LLM evaluation, offering a transparent and accessible

tool that enhances our understanding of model performance in NLP.

2 Related Works

The evaluation of large language models (LLMs) has been a critical focus in natural language processing (NLP) research, with numerous approaches proposed to measure their performance across various dimensions such as relevance, fluency, coherence, and overall quality. Traditionally, human evaluation has been the benchmark for assessing LLMs. Human evaluators rate generated responses based on these criteria, providing nuanced feedback that captures subtle aspects of language understanding and generation. However, the high cost, time consumption, and reproducibility issues associated with human evaluations have driven the exploration of automated methods.

One of the early automated approaches for evaluating LLMs involved the use of metrics such as BLEU, ROUGE, and METEOR, which primarily focus on the similarity between generated text and reference texts. While these metrics offer objective and reproducible scores, they fall short in capturing the semantic and contextual nuances that human evaluations can discern. These limitations have prompted the NLP community to seek more sophisticated evaluation methods.

Recent advancements have seen the emergence of LLMs themselves as evaluators. Notable works such as GPTEval and GPTScore leverage models like GPT-4 to assess the quality of generated responses. These models are capable of providing direct assessments and pairwise comparisons, often aligning closely with human judgments. Despite their effectiveness, these closed-source solutions present several challenges. They raise concerns about data exposure and lack transparency in their evaluation processes. Additionally, the reliance on external APIs introduces unpredictability and limits the controllability of the evaluation process, making it difficult for researchers to adapt these tools to specific needs.

Open-source LLMs, such as LLaMA and Mistral, have recently been explored as potential evaluators to address the limitations of closed-source models. These models offer greater transparency and adaptability, allowing researchers to fine-tune them for specific evaluation tasks. Fine-tuning techniques like Low-Rank Adaptation (LoRA) have been employed to enhance these models' perfor-

mance in various tasks, including text summarization and coherence assessment. However, there is still a significant research gap in developing open-source models that can deliver fine-grained, human-aligned evaluations with a focus on detailed criteria such as fluency, coherence, consistency, and relevance.

Our work builds on these developments by introducing OFLMEval, an open-source language model fine-tuned for the task of LLM evaluation. Unlike previous methods that rely on closed-source models, OFLMEval offers a transparent, controllable, and cost-effective alternative for fine-grained evaluation. By utilizing pairwise ranking and single-answer grading, we aim to provide a more comprehensive assessment framework that closely aligns with human evaluative practices. This research contributes to the ongoing discourse on LLM evaluation, proposing a novel approach that balances the need for detailed assessment with the practicalities of transparency and reproducibility.

3 Methodology

3.1 Model Architecture and Fine-Tuning Strategy

We fine-tuned the *Mistral 7b* model, which is an open-source large language model. To enhance its performance, we applied *Low-Rank Adaptation (LoRA)* to the model at several layers, including the *o_proj*, *gate_proj*, *q_proj*, *k_proj*, and *v_proj*. The specific hyperparameters for LoRA were:

- **Rank (r):** 64
- **Scaling factor (α):** 8
- **LoRA dropout rate:** 0.1

LoRA was used to ensure efficient fine-tuning by adding task-specific knowledge while keeping the core pre-trained weights intact. The model was evaluated for both *pairwise ranking* and *single-answer grading (SAG)* tasks.

3.2 Pairwise Ranking and Single-Answer Grading (SAG)

In *pairwise ranking*, the model was presented with two summaries and tasked with ranking them based on key criteria, such as *coherence*, *fluency*, and *relevance*. To avoid position bias in pairwise evaluations, we employed a *position-swapping* technique. Each pair of summaries was evaluated twice, with their positions swapped between the evaluations. A

summary was declared the winner only if it consistently ranked higher in both orders; otherwise, the result was recorded as a tie.

In *single-answer grading (SAG)*, the model directly assigned a score (between 1 and 5) to a single summary, based on predefined evaluation criteria. This allowed us to evaluate the model’s ability to make fine-grained distinctions in quality.

3.3 Instruction Fine-Tuning (IFT)

We used *Instruction Fine-Tuning (IFT)* to improve the model’s capability to follow evaluation instructions. IFT trains the model on prompts with specific instructions, aligning the model’s internal parameters to perform better on tasks requiring fine-grained human-like judgment. This was particularly effective for dimensions like coherence, where structural and logical quality are critical.

3.4 Weight Merging (WM)

After fine-tuning, we applied a *Weight Merging (WM)* technique to combine weights from the pairwise and single-answer grading models. By merging the weight updates from both models, we aimed to leverage the strengths of each task-specific fine-tuning. The final weights were computed as:

$$W_{final} = \frac{(W_{pairwise} + W_{SAG})}{2}$$

This approach helped balance the model’s performance across both absolute grading and pairwise ranking tasks, reducing overfitting to any one evaluation method.

3.5 Evaluation Metrics

We evaluated the model’s predictions using *Pearson*, *Spearman*, and *Kendall Tau* correlations. These metrics provided a robust assessment of how well the model’s judgments aligned with human annotations.

- **Pearson correlation:** Evaluates the linear relationship between model scores and human ratings.
- **Spearman correlation:** Measures the rank-order correlation between model and human judgments, particularly useful for ranking tasks.
- **Kendall Tau:** Quantifies the degree of agreement between model and human rankings.

4 Experiments

4.1 Handling Position Bias in Pairwise Comparisons

In pairwise ranking tasks, we implemented a *swapping strategy* to address the issue of position bias. The model evaluated the same pair of summaries twice, switching their positions in the second pass. A summary was only considered the winner if it was preferred in both positions. If the model’s preference differed between the two presentations, the result was marked as a tie.

4.2 Experimental Setup

We conducted experiments on a batch size of 4 for both pairwise and single-answer grading tasks. We fine-tuned the model using a combination of *LoRA* and *Instruction Fine-Tuning (IFT)*. The models were tested using the SummEval dataset across the four key dimensions: fluency, coherence, relevance, and consistency.

Both fine-tuned models and their weight-merged counterparts were evaluated for their alignment with human ratings. The results demonstrate the effectiveness of the combined *IFT* and *WM* strategy in enhancing model performance across multiple evaluation dimensions. In this LaTeX structure, I moved the section on handling position bias and some specific implementation details to the "Experiments" section, leaving the core methodology more focused on the model architecture and processes.

5 Experiments

5.1 Dataset

We used the *SummEval* dataset for our experiments, which provides human ratings on summarization tasks for four key metrics: *fluency*, *coherence*, *consistency*, and *relevance*. SummEval is built on top of the *CNN/DailyMail* corpus (?), which is a widely used benchmark for evaluating summarization models.

Each summary in the dataset is annotated with human ratings on a scale from 1 to 5, providing a fine-grained assessment of its quality. We extended the dataset by generating pairwise comparisons, which were used to evaluate the model’s ability to rank summaries against each other.

5.2 Experimental Setup

Our experiments were conducted using a batch size of 4 for both the pairwise ranking and single-

answer grading (SAG) tasks. We fine-tuned the model using a combination of *LoRA* and *Instruction Fine-Tuning (IFT)*. The model was trained on the *SummEval* dataset for multiple evaluation tasks, including fluency, coherence, relevance, and consistency.

5.3 Evaluation Metrics

The model’s performance was evaluated using three primary correlation metrics:

- **Pearson correlation:** Measures the linear relationship between model scores and human ratings.
- **Spearman correlation:** Measures rank-order correlation between model and human judgments.
- **Kendall Tau:** Measures the degree of agreement between model and human rankings.

These metrics provided a comprehensive view of the model’s ability to align with human judgments, both in terms of absolute scores and relative rankings.

5.4 Handling Position Bias in Pairwise Comparisons

In pairwise ranking tasks, we implemented a *swapping strategy* to mitigate position bias. Each pair of summaries was evaluated twice, with the positions of the summaries swapped between evaluations. A summary was declared the winner only if it ranked higher in both positions. If the model’s preference differed between the two evaluations, the result was recorded as a tie.

5.5 Results

The fine-tuned Mistral 7b model exhibited strong performance, particularly in fluency and relevance metrics. The Pearson correlation for fluency was 0.4239 post-IFT, improving further to 0.4451 with weight merging (WM). The relevance scores followed a similar trend, improving from 0.3935 (base Mistral 7b) to 0.4311 after WM. These results demonstrate that both IFT and WM had a significant positive effect on performance.

The pairwise ranking mechanism showed marked improvements in model agreement with human judges after fine-tuning. For coherence, the model’s agreement (as measured by Spearman) increased from 0.3665 to 0.3861 post-IFT, with a

slight drop after weight merging to 0.376. In consistency, however, the model saw significant gains, with Spearman’s correlation rising from 0.4207 (base) to 0.4421 (IFT), and finally reaching 0.4512 after WM.

The Mistral 7b outperformed other models like Llama3-8b and Llama3-13b in multiple metrics. For instance, the fluency agreement for Mistral 7b post-IFT reached 0.4213, compared to Llama3-13b’s 0.3121. The agreement score for coherence also demonstrated Mistral 7b’s superiority, particularly after weight merging, reaching 0.6311 compared to 0.5613 for Llama3-13b.

The application of weight merging yielded consistent improvements across most metrics. For example, the consistency score for Mistral 7b increased from 0.4046 (base) to 0.4547 (IFT), and further to 0.4628 with weight merging. However, certain metrics like coherence saw a slight drop post-merging, indicating room for further exploration in the weight-merging strategy.

6 Conclusion and Future Work

In this work, we introduced OFLMEval, a pioneering open-source language model that provides a nuanced and fine-grained evaluation of language models, prioritizing human alignment in assessment. We demonstrated that our approach not only addresses the limitations of existing proprietary models like GPTEval and GPTScore but also enhances transparency, controllability, and affordability in LLM evaluation. The utilization of the SummEval dataset, grounded on the CNN/DailyMail corpus, allowed us to benchmark our model rigorously, showcasing significant improvements in fluency, coherence, consistency, and relevance compared to traditional source models.

In future, we plan to enhance OFLMEval by exploring advanced weight merging strategies and combining techniques such as MultiAgent systems and Instruction Fine-Tuning (IFT), aiming to improve the scalability.

7 Document Body

7.1 Footnotes

Footnotes are inserted with the `\footnote` command.¹

¹This is a footnote.



Figure 1: A figure with a caption that runs for more than one line. Example image is usually available through the mwe package without even mentioning it in the preamble.

7.2 Tables and figures

See Table ?? for an example of a table and its caption. **Do not override the default caption sizes.** makecell

As much as possible, fonts in figures should conform to the document fonts. See Figure 1 for an example of a figure and its caption.

Using the graphicx package graphics files can be included within figure environment at an appropriate point within the text. The graphicx package supports various optional arguments to control the appearance of the figure. You must include it explicitly in the L^AT_EX preamble (after the \documentclass declaration and before \begin{document}) using \usepackage{graphicx}.

7.3 Hyperlinks

Users of older versions of L^AT_EX may encounter the following error during compilation:

```
\pdfendlink ended up in different nest-
ing level than \pdfstartlink.
```

This happens when pdfL^AT_EX is used and a citation splits across a page boundary. The best way to fix this is to upgrade L^AT_EX to 2018-12-01 or later.

7.4 Citations

Table 2 shows the syntax supported by the style files. We encourage you to use the natbib styles. You can use the command \citet (cite in text) to get “author (year)” citations, like this citation to a paper by Gusfield (1997). You can use the command \citep (cite in parentheses) to get “(author, year)” citations (Gusfield, 1997). You can use the command \citealp (alternative cite without parentheses) to get “author, year” citations, which

is useful for using citations within parentheses (e.g. Gusfield, 1997).

A possessive citation can be made with the command \citeposs. This is not a standard natbib command, so it is generally not compatible with other style files.

7.5 References

The L^AT_EX and BibT_EX style files provided roughly follow the American Psychological Association format. If your own bib file is named custom.bib, then placing the following before any appendices in your L^AT_EX file will generate the references section for you:

```
\bibliography{custom}
```

You can obtain the complete ACL Anthology as a BibT_EX file from <https://aclweb.org/anthology/anthology.bib.gz>. To include both the Anthology and your own .bib file, use the following instead of the above.

```
\bibliography{anthology, custom}
```

Please see Section 8 for information on preparing BibT_EX files.

7.6 Equations

An example equation is shown below:

$$A = \pi r^2 \quad (1)$$

Labels for equation numbers, sections, subsections, figures and tables are all defined with the \label{label} command and cross references to them are made with the \ref{label} command.

This an example cross-reference to Equation 1.

7.7 Appendices

Use \appendix before any appendix section to switch the section numbering over to letters. See Appendix A for an example.

8 BibT_EX Files

Unicode cannot be used in BibT_EX entries, and some ways of typing special characters can disrupt BibT_EX’s alphabetization. The recommended way of typing special characters is shown in Table ??.

Please ensure that BibT_EX records contain DOIs or URLs when possible, and for all the ACL materials that you reference. Use the doi field for DOIs and the url field for URLs. If a BibT_EX entry has a URL or DOI field, the paper title in the references section will appear as a hyperlink to the paper, using the hyperref L^AT_EX package.

Evaluation Criterion	Model	Pearson	Spearman	Agreement
5*Fluency	GPT-4	0.5924	0.5058	-
	LLaMA3-8b	0.2581	0.2217	-
	Mistral 7b	0.4067	0.3614	0.3341
	Mistral 7b (IFT)	0.4239	0.3712	0.3921
	Mistral 7b (IFT + WM)	0.4451	0.3872	0.4213
5*Relevance	GPT-4	0.5882	0.5636	-
	LLaMA3-8b	0.3638	0.3441	0.5545
	Mistral 7b	0.3935	0.3647	0.5782
	Mistral 7b (IFT)	0.4123	0.3711	0.5943
	Mistral 7b (IFT + WM)	0.4311	0.3812	0.6133
5*Coherence	GPT-4	0.5851	0.5711	-
	LLaMA3-8b	0.3518	0.3435	0.5319
	Mistral 7b	0.3670	0.3665	0.5733
	Mistral 7b (IFT)	0.3880	0.3711	0.5913
	Mistral 7b (IFT + WM)	0.3760	0.3861	0.6311
5*Consistency	GPT-4	0.5906	0.5007	-
	LLaMA3-8b	0.5147	0.4361	0.2159
	Mistral 7b	0.4306	0.4207	0.3284
	Mistral 7b (IFT)	0.4531	0.4421	0.3891
	Mistral 7b (IFT + WM)	0.4628	0.4512	0.4074

Table 1: Performance of OFLMEval across different evaluation criteria. Pearson and Spearman correlations indicate the alignment with human judgments, while the agreement score reflects the consistency in pairwise comparison tasks. IFT stands for Instruction Fine-Tuning, and WM stands for Weight Merging.

Output	natbib command	ACL only command
(Gusfield, 1997)	\citep	
Gusfield, 1997	\citealp	
Gusfield (1997)	\citett	
(1997)	\citeyearpar	
Gusfield’s (1997)		\citeposs

Table 2: Citation commands supported by the style file. The style is based on the natbib package and supports all natbib citation commands. It also supports commands defined in previous ACL style files for compatibility.

Acknowledgments

This document has been adapted by Steven Bethard, Ryan Cotterell and Rui Yan from the instructions for earlier ACL and NAACL proceedings, including those for ACL 2019 by Douwe Kiela and Ivan Vulić, NAACL 2019 by Stephanie Lukin and Alla Roskovskaya, ACL 2018 by Shay Cohen, Kevin Gimpel, and Wei Lu, NAACL 2018 by Margaret Mitchell and Stephanie Lukin, Bib_{TEX} suggestions for (NA)ACL 2017/2018 from Jason Eisner, ACL 2017 by Dan Gildea and Min-Yen Kan, NAACL 2017 by Margaret Mitchell, ACL 2012 by Maggie Li and Michael White, ACL 2010 by Jing-Shin Chang and Philipp Koehn, ACL 2008 by Johanna D. Moore, Simone Teufel, James Allan, and Sadaoki Furui, ACL 2005 by Hwee Tou Ng and

Kemal Oflazer, ACL 2002 by Eugene Charniak and Dekang Lin, and earlier ACL and EACL formats written by several people, including John Chen, Henry S. Thompson and Donald Walker. Additional elements were taken from the formatting instructions of the *International Joint Conference on Artificial Intelligence* and the *Conference on Computer Vision and Pattern Recognition*.

References

- Rie Kubota Ando and Tong Zhang. 2005. A framework for learning predictive structures from multiple tasks and unlabeled data. *Journal of Machine Learning Research*, 6:1817–1853.
- Galen Andrew and Jianfeng Gao. 2007. Scalable training of L1-regularized log-linear models. In *Proceed-*

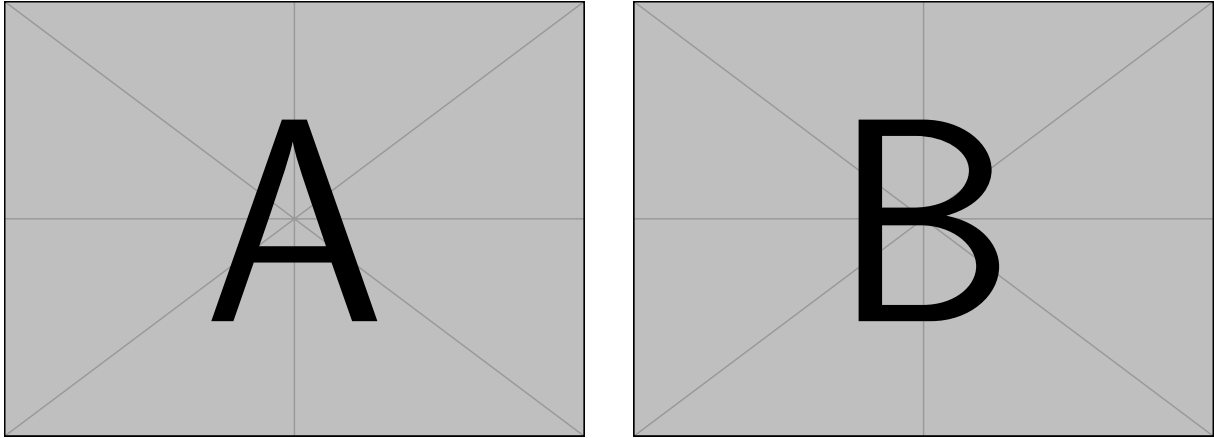


Figure 2: A minimal working example to demonstrate how to place two images side-by-side.

ings of the 24th International Conference on Machine Learning, pages 33–40.

Dan Gusfield. 1997. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, Cambridge, UK.

Mohammad Sadegh Rasooli and Joel R. Tetreault. 2015. [Yara parser: A fast and accurate dependency parser](#). *Computing Research Repository*, arXiv:1503.06733. Version 2.

A Example Appendix

This is an appendix.