

# A Case Study Replicating Calibration of Large Language Models

Anonymous ACL submission

## Abstract

The popularity of large language models has lead to a growing need to need to characterize their behavior and understanding beyond the generation grammatical sentences In this study we attempt to replicate and extend the calibration study of Kadavath 2001 on open models and datasets, attempting compare model response to their calibration, that is the certainty with which a model answers a question with a given completion and the probability that the model is answering the question correctly.

## 1 Introduction

These instructions are for authors submitting papers to ACL 2023 using L<sup>A</sup>T<sub>E</sub>X. They are not self-contained. All authors must follow the general instructions for \*ACL proceedings,<sup>1</sup> as well as guidelines set forth in the ACL 2023 call for papers.<sup>2</sup> This document contains additional instructions for the L<sup>A</sup>T<sub>E</sub>X style files. The templates include the L<sup>A</sup>T<sub>E</sub>X source of this document (acl2023.tex), the L<sup>A</sup>T<sub>E</sub>X style file used to format it (acl2023.sty), an ACL bibliography style (acl\_natbib.bst), an example bibliography (custom.bib), and the bibliography for the ACL Anthology (anthology.bib).

## 2 Engines

To produce a PDF file, pdfL<sup>A</sup>T<sub>E</sub>X is strongly recommended (over original L<sup>A</sup>T<sub>E</sub>X plus dvips+ps2pdf or dvi2pdf). XeL<sup>A</sup>T<sub>E</sub>X also produces PDF files, and is especially suitable for text in non-Latin scripts.

## 3 Preamble

The first line of the file must be

```
\documentclass[11pt]{article}
```

<sup>1</sup><http://acl-org.github.io/ACL2023/formatting.html>

<sup>2</sup>[https://2023.aclweb.org/calls/main\\_conference/](https://2023.aclweb.org/calls/main_conference/)

Command	Output	Command	Output
{\ "a}	ä	{\c c}	ç
{\^e}	ê	{\u g}	ğ
{\`i}	ì	{\l}	ł
{\ .I}	İ	{\~n}	ñ
{\o}	ø	{\H o}	ő
{\'u}	ú	{\v r}	ř
{\aa}	å	{\ss}	ß

Table 1: Example commands for accented characters, to be used in, *e.g.*, BibT<sub>E</sub>X entries.

To load the style file in the review version:

```
\usepackage[review]{ACL2023}
```

For the final version, omit the review option:

```
\usepackage{ACL2023}
```

To use Times Roman, put the following in the preamble:

```
\usepackage{times}
```

(Alternatives like txfonts or newtx are also acceptable.) Please see the L<sup>A</sup>T<sub>E</sub>X source of this document for comments on other packages that may be useful. Set the title and author using \title and \author. Within the author list, format multiple authors using \and and \And and \AND; please see the L<sup>A</sup>T<sub>E</sub>X source for examples. By default, the box containing the title and author names is set to the minimum of 5 cm. If you need more space, include the following in the preamble:

```
\setlength\titlebox{<dim>}
```

where <dim> is replaced with a length. Do not set this length smaller than 5 cm.

## 4 Document Body

### 4.1 Footnotes

Footnotes are inserted with the \footnote command.<sup>3</sup>

<sup>3</sup>This is a footnote.

Output	natbib command	Old ACL-style command
(Cooley and Tukey, 1965)	<code>\citep</code>	<code>\cite</code>
Cooley and Tukey, 1965	<code>\citealp</code>	no equivalent
Cooley and Tukey (1965)	<code>\citet</code>	<code>\newcite</code>
(1965)	<code>\citeyearpar</code>	<code>\shortcite</code>
Cooley and Tukey’s (1965)	<code>\citeposs</code>	no equivalent
(FFT; Cooley and Tukey, 1965)	<code>\citep[FFT;][]</code>	no equivalent

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.

## 4.2 Tables and figures

See Table 1 for an example of a table and its caption.

**Do not override the default caption sizes.**

## 4.3 Hyperlinks

Users of older versions of L<sup>A</sup>T<sub>E</sub>X may encounter the following error during compilation:

```
\pdfendlink ended up in different
nesting level than \pdfstartlink.
```

This happens when pdfL<sup>A</sup>T<sub>E</sub>X is used and a citation splits across a page boundary. The best way to fix this is to upgrade L<sup>A</sup>T<sub>E</sub>X to 2018-12-01 or later.

## 4.4 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).

## 4.5 References

The L<sup>A</sup>T<sub>E</sub>X and BibT<sub>E</sub>X 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<sup>A</sup>T<sub>E</sub>X file will generate the references section for you:

```
\bibliographystyle{acl_natbib}
\bibliography{custom}
```

You can obtain the complete ACL Anthology as a BibT<sub>E</sub>X 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.

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

Please see Section 5 for information on preparing BibT<sub>E</sub>X files.

## 4.6 Appendices

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

## 5 BibT<sub>E</sub>X Files

Unicode cannot be used in BibT<sub>E</sub>X entries, and some ways of typing special characters can disrupt BibT<sub>E</sub>X’s alphabetization. The recommended way of typing special characters is shown in Table 1.

Please ensure that BibT<sub>E</sub>X 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<sub>E</sub>X 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<sup>A</sup>T<sub>E</sub>X package.

## Limitations

ACL 2023 requires all submissions to have a section titled “Limitations”, for discussing the limitations of the paper as a complement to the discussion of strengths in the main text. This section should occur after the conclusion, but before the references. It will not count towards the page limit. The discussion of limitations is mandatory. Papers without a limitation section will be desk-rejected without review.

While we are open to different types of limitations, just mentioning that a set of results have been shown for English only probably does not reflect what we expect. Mentioning that the method

works mostly for languages with limited morphology, like English, is a much better alternative. In addition, limitations such as low scalability to long text, the requirement of large GPU resources, or other things that inspire crucial further investigation are welcome.

## Ethics Statement

Scientific work published at ACL 2023 must comply with the ACL Ethics Policy.<sup>4</sup> We encourage all authors to include an explicit ethics statement on the broader impact of the work, or other ethical considerations after the conclusion but before the references. The ethics statement will not count toward the page limit (8 pages for long, 4 pages for short papers).

## Acknowledgements

This document has been adapted by Jordan Boyd-Graber, Naoaki Okazaki, Anna Rogers from the style files used for earlier ACL, EMNLP and NAACL proceedings, including those for EACL 2023 by Isabelle Augenstein and Andreas Vlachos, EMNLP 2022 by Yue Zhang, Ryan Cotterell and Lea Frermann, ACL 2020 by Steven Bethard, Ryan Cotterell and Rui Yan, 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, BibTeX 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*.

<sup>4</sup><https://www.aclweb.org/portal/content/acl-code-ethics>

## 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  \$L\_1\$ -regularized log-linear models](#). In *Proceedings of the 24th International Conference on Machine Learning*, pages 33–40.
- Isabelle Augenstein, Tim Rocktäschel, Andreas Vlachos, and Kalina Bontcheva. 2016. [Stance detection with bidirectional conditional encoding](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 876–885, Austin, Texas. Association for Computational Linguistics.
- James W. Cooley and John W. Tukey. 1965. [An algorithm for the machine calculation of complex Fourier series](#). *Mathematics of Computation*, 19(90):297–301.
- James Goodman, Andreas Vlachos, and Jason Naradowsky. 2016. [Noise reduction and targeted exploration in imitation learning for Abstract Meaning Representation parsing](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1–11, Berlin, Germany. Association for Computational Linguistics.
- Dan Gusfield. 1997. *Algorithms on Strings, Trees and Sequences*. Cambridge University Press, Cambridge, UK.
- Mary Harper. 2014. [Learning from 26 languages: Program management and science in the babel program](#). In *Proceedings of COLING 2014, the 25th International Conference on Computational Linguistics: Technical Papers*, page 1, Dublin, Ireland. Dublin City University and Association for Computational Linguistics.
- 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 a section in the appendix.