# Unveiling the Syntax Within: Interpreting Grammar Embeddings in Meta's LLaMA Models

**Pratim Chowdhary**[*]
Department of Computer Science
Dartmouth College
cpratim.25@dartmouth.edu

**Peter Chin**
Department of Engineering
Thayer School of Engineering
pc@dartmouth.edu

**Deepernab Chakrabarty**
Department of Computer Science
Dartmouth College
deepernab@dartmouth.edu

## Abstract

This paper investigates the mechanisms by which large language models (LLMs) encode grammatical knowledge, focusing on Meta's LLaMA models. By leveraging embedding vectors, we classify grammatically correct sentences and analyze the activation patterns of attention heads to identify their roles in processing specific grammatical structures. Furthermore, we explore the effects of selectively removing these attention heads, shedding light on how grammar is embedded within the model's architecture. Our findings aim to enhance the understanding of LLMs' linguistic capabilities and their internal organization of syntactic knowledge.

## 1 Submission of papers to the M3L Workshop at NeurIPS 2024

Please read the instructions below carefully and follow them faithfully. fds

### 1.1 Style

Papers to be submitted to the Mathematics of Modern Machine Learning (M3L) Workshop at NeurIPS 2024 must be prepared according to the instructions presented here.

Authors are required to use the Mathematics of Modern Machine Learning (M3L) LaTeX style files obtainable at the workshop website https://sites.google.com/view/m3l-2024/call-for-papers. Please make sure you use the current files and not previous versions. Tweaking the style files may be grounds for rejection.

### 1.2 Retrieval of style files

The style files for the Mathematics of Modern Machine Learning (M3L) Workshop at NeurIPS 2024 and other conference information are available on the website at

https://sites.google.com/view/m3l-2024

The file main.pdf contains these instructions and illustrates the various formatting requirements your NeurIPS paper must satisfy.

---

[*]Use footnote for providing further information about author (webpage, alternative address)—*not* for acknowledging funding agencies.

The LATEX style file contains three optional arguments: `final`, which creates a camera-ready copy, `preprint`, which creates a preprint for submission to, e.g., arXiv, and `nonatbib`, which will not load the `natbib` package for you in case of package clash.

**Preprint option**  If you wish to post a preprint of your work online, e.g., on arXiv, using the NeurIPS style, please use the `preprint` option. This will create a nonanonymized version of your work with the text "Preprint. Work in progress." in the footer. This version may be distributed as you see fit, as long as you do not say which conference it was submitted to. Please **do not** use the `final` option, which should **only** be used for papers accepted to the Mathematics of Modern Machine Learning (M3L) Workshop at NeurIPS 2024.

At submission time, please omit the `final` and `preprint` options. This will anonymize your submission and add line numbers to aid review. Please do *not* refer to these line numbers in your paper as they will be removed during generation of camera-ready copies.

The file `main.tex` may be used as a "shell" for writing your paper. All you have to do is replace the author, title, abstract, and text of the paper with your own.

The formatting instructions contained in these style files are summarized in Sections 2, 3, and 4 below.

## 2   General formatting instructions

The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long. The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing (leading) of 11 points. Times New Roman is the preferred typeface throughout, and will be selected for you by default. Paragraphs are separated by ½ line space (5.5 points), with no indentation.

The paper title should be 17 point, initial caps/lower case, bold, centered between two horizontal rules. The top rule should be 4 points thick and the bottom rule should be 1 point thick. Allow ¼ inch space above and below the title to rules. All pages should start at 1 inch (6 picas) from the top of the page.

For the final version, authors' names are set in boldface, and each name is centered above the corresponding address. The lead author's name is to be listed first (left-most), and the co-authors' names (if different address) are set to follow. If there is only one co-author, list both author and co-author side by side.

Please pay special attention to the instructions in Section 4 regarding figures, tables, acknowledgments, and references.

## 3   Headings: first level

All headings should be lower case (except for first word and proper nouns), flush left, and bold.

First-level headings should be in 12-point type.

### 3.1   Headings: second level

Second-level headings should be in 10-point type.

#### 3.1.1   Headings: third level

Third-level headings should be in 10-point type.

**Paragraphs**  There is also a `\paragraph` command available, which sets the heading in bold, flush left, and inline with the text, with the heading followed by 1 em of space.

## 4   Citations, figures, tables, references

These instructions apply to everyone.

### 4.1 Citations within the text

The `natbib` package will be loaded for you by default. Citations may be author/year or numeric, as long as you maintain internal consistency. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

The documentation for `natbib` may be found at

> http://mirrors.ctan.org/macros/latex/contrib/natbib/natnotes.pdf

Of note is the command `\citet`, which produces citations appropriate for use in inline text. For example,

> `\citet{hasselmo} investigated\dots`

produces

> Hasselmo, et al. (1995) investigated...

If you wish to load the `natbib` package with options, you may add the following before loading the `neurips_2024` package:

> `\PassOptionsToPackage{options}{natbib}`

If `natbib` clashes with another package you load, you can add the optional argument `nonatbib` when loading the style file:

> `\usepackage[nonatbib]{neurips_2024}`

As submission is double blind, refer to your own published work in the third person. That is, use "In the previous work of Jones et al. [4]," not "In our previous work [4]." If you cite your other papers that are not widely available (e.g., a journal paper under review), use anonymous author names in the citation, e.g., an author of the form "A. Anonymous" and include a copy of the anonymized paper in the supplementary material.

### 4.2 Footnotes

Footnotes should be used sparingly. If you do require a footnote, indicate footnotes with a number[2] in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).

Note that footnotes are properly typeset *after* punctuation marks.[3]

### 4.3 Figures

All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction. The figure number and caption always appear after the figure. Place one line space before the figure caption and one line space after the figure. The figure caption should be lower case (except for first word and proper nouns); figures are numbered consecutively.

You may use color figures. However, it is best for the figure captions and the paper body to be legible if the paper is printed in either black/white or in color.

### 4.4 Tables

All tables must be centered, neat, clean and legible. The table number and title always appear before the table. See Table 1.

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

---

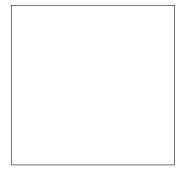[2]Sample of the first footnote.

[3]As in this example.

Figure 1: Sample figure caption.

Table 1: Sample table title

| | Part | | |
|---|---|---|
| Name | Description | Size ($\mu$m) |
| Dendrite | Input terminal | $\sim$100 |
| Axon | Output terminal | $\sim$10 |
| Soma | Cell body | up to $10^6$ |

Note that publication-quality tables *do not contain vertical rules.* We strongly suggest the use of the booktabs package, which allows for typesetting high-quality, professional tables:

https://www.ctan.org/pkg/booktabs

This package was used to typeset Table 1.

## 4.5 Math

Note that display math in bare TeX commands will not create correct line numbers for submission. Please use LaTeX (or AMSTeX) commands for unnumbered display math. (You really shouldn't be using $$ anyway; see https://tex.stackexchange.com/questions/503/why-is-preferable-to and https://tex.stackexchange.com/questions/40492/what-are-the-differences-between-align-equation-and-displaymath for more information.)

## 4.6 Final instructions

Do not change any aspects of the formatting parameters in the style files. In particular, do not modify the width or length of the rectangle the text should fit into, and do not change font sizes (except perhaps in the **References** section; see below). Please note that pages should be numbered.

## 5 Preparing PDF files

Please prepare submission files with paper size "US Letter," and not, for example, "A4."

Fonts were the main cause of problems in the past years. Your PDF file must only contain Type 1 or Embedded TrueType fonts. Here are a few instructions to achieve this.

- You should directly generate PDF files using pdflatex.

- You can check which fonts a PDF files uses. In Acrobat Reader, select the menu Files>Document Properties>Fonts and select Show All Fonts. You can also use the program pdffonts which comes with xpdf and is available out-of-the-box on most Linux machines.

- xfig "patterned" shapes are implemented with bitmap fonts. Use "solid" shapes instead.

- The \bbold package almost always uses bitmap fonts. You should use the equivalent AMS Fonts:

    \usepackage{amsfonts}

    followed by, e.g., \mathbb{R}, \mathbb{N}, or \mathbb{C} for $\mathbb{R}$, $\mathbb{N}$ or $\mathbb{C}$. You can also use the following workaround for reals, natural and complex:

    \newcommand{\RR}{I\!\!R} %real numbers
    \newcommand{\Nat}{I\!\!N} %natural numbers
    \newcommand{\CC}{I\!\!\!\!C} %complex numbers

    Note that amsfonts is automatically loaded by the amssymb package.

If your file contains type 3 fonts or non embedded TrueType fonts, we will ask you to fix it.

## 5.1 Margins in LaTeX

Most of the margin problems come from figures positioned by hand using \special or other commands. We suggest using the command \includegraphics from the graphicx package. Always specify the figure width as a multiple of the line width as in the example below:

    \usepackage[pdftex]{graphicx} ...
    \includegraphics[width=0.8\linewidth]{myfile.pdf}

See Section 4.4 in the graphics bundle documentation (http://mirrors.ctan.org/macros/latex/required/graphics/grfguide.pdf)

A number of width problems arise when LaTeX cannot properly hyphenate a line. Please give LaTeX hyphenation hints using the \- command when necessary.

## References

References follow the acknowledgments in the camera-ready paper. Use unnumbered first-level heading for the references. Any choice of citation style is acceptable as long as you are consistent. It is permissible to reduce the font size to small (9 point) when listing the references. Note that the Reference section does not count towards the page limit.

[1] Alexander, J.A. & Mozer, M.C. (1995) Template-based algorithms for connectionist rule extraction. In G. Tesauro, D.S. Touretzky and T.K. Leen (eds.), *Advances in Neural Information Processing Systems 7*, pp. 609–616. Cambridge, MA: MIT Press.

[2] Bower, J.M. & Beeman, D. (1995) *The Book of GENESIS: Exploring Realistic Neural Models with the GEneral NEural SImulation System.* New York: TELOS/Springer–Verlag.

[3] Hasselmo, M.E., Schnell, E. & Barkai, E. (1995) Dynamics of learning and recall at excitatory recurrent synapses and cholinergic modulation in rat hippocampal region CA3. *Journal of Neuroscience* **15**(7):5249-5262.

## A  Appendix / supplemental material

Optionally include supplemental material (complete proofs, additional experiments and plots) in appendix.