

LLMs AS AGENTS: AN EMPIRICAL STUDY PROBING INTO THE EVOLUTION OF NAVIGABLE ABILITIES THROUGH CONTINUED CONTRASTIVE ALIGNMENT

Anonymous authors

Paper under double-blind review

ABSTRACT

The abstract paragraph should be indented 1/2 inch (3 picas) on both left and right-hand margins. Use 10 point type, with a vertical spacing of 11 points. The word ABSTRACT must be centered, in small caps, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

1 INTRODUCTION

With their capabilities only increasing with each passing day, its no surprise that language models are being adapted for almost all tasks across multiple domains.

2 RELATED WORK

3 PROBLEM DESCRIPTION

3.1 SIMULATION SETUP

3.2 AGENT BEHAVIOR

3.3 ON THE NAVIGABLE ABILITIES OF BLACK BOX LANGUAGE MODELS

4 EXPERIMENTAL SETUP

4.1 RED-BLUE AGENT SIMULATION

4.2 CONTRASTIVE ALIGNMENT TECHNIQUE

4.3 THE RED-BLUE DATASET

4.4 NAVIGABLE METRICS

4.4.1 RED METRICS

4.4.2 BLUE METRICS

4.5 RESULTS AND OBSERVATIONS

5 CONCLUSION

6 APPENDIX

7 SUBMISSION OF CONFERENCE PAPERS TO ICLR 2025

ICLR requires electronic submissions, processed by <https://openreview.net/>. See ICLR's website for more instructions.

If your paper is ultimately accepted, the statement `\iclrfinalcopy` should be inserted to adjust the format to the camera ready requirements.

The format for the submissions is a variant of the NeurIPS format. Please read carefully the instructions below, and follow them faithfully.

7.1 STYLE

Papers to be submitted to ICLR 2025 must be prepared according to the instructions presented here.

Authors are required to use the ICLR \LaTeX style files obtainable at the ICLR website. Please make sure you use the current files and not previous versions. Tweaking the style files may be grounds for rejection.

7.2 RETRIEVAL OF STYLE FILES

The style files for ICLR and other conference information are available online at:

<http://www.iclr.cc/>

The file `iclr2025_conference.pdf` contains these instructions and illustrates the various formatting requirements your ICLR paper must satisfy. Submissions must be made using \LaTeX and the style files `iclr2025_conference.sty` and `iclr2025_conference.bst` (to be used with \LaTeX 2e). The file `iclr2025_conference.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 8, 9, and 10 below.

8 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 of 11 points. Times New Roman is the preferred typeface throughout. Paragraphs are separated by 1/2 line space, with no indentation.

Paper title is 17 point, in small caps and left-aligned. All pages should start at 1 inch (6 picas) from the top of the page.

Authors’ names are set in boldface, and each name is placed above its corresponding address. The lead author’s name is to be listed first, and the co-authors’ names are set to follow. Authors sharing the same address can be on the same line.

Please pay special attention to the instructions in section 10 regarding figures, tables, acknowledgments, and references.

There will be a strict upper limit of 10 pages for the main text of the initial submission, with unlimited additional pages for citations.

9 HEADINGS: FIRST LEVEL

First level headings are in small caps, flush left and in point size 12. One line space before the first level heading and 1/2 line space after the first level heading.

9.1 HEADINGS: SECOND LEVEL

Second level headings are in small caps, flush left and in point size 10. One line space before the second level heading and 1/2 line space after the second level heading.

9.1.1 HEADINGS: THIRD LEVEL

Third level headings are in small caps, flush left and in point size 10. One line space before the third level heading and 1/2 line space after the third level heading.

10 CITATIONS, FIGURES, TABLES, REFERENCES

These instructions apply to everyone, regardless of the formatter being used.

10.1 CITATIONS WITHIN THE TEXT

Citations within the text should be based on the `natbib` package and include the authors' last names and year (with the "et al." construct for more than two authors). When the authors or the publication are included in the sentence, the citation should not be in parenthesis using `\citet{}` (as in "See Hinton et al. (2006) for more information."). Otherwise, the citation should be in parenthesis using `\citep{}` (as in "Deep learning shows promise to make progress towards AI (Bengio & LeCun, 2007).").

The corresponding references are to be listed in alphabetical order of authors, in the REFERENCES section. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

10.2 FOOTNOTES

Indicate footnotes with a number¹ 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).²

10.3 FIGURES

All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction; art work should not be hand-drawn. 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 is lower case (except for first word and proper nouns); figures are numbered consecutively.

Make sure the figure caption does not get separated from the figure. Leave sufficient space to avoid splitting the figure and figure caption.

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

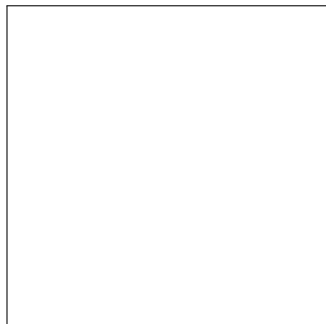


Figure 1: Sample figure caption.

¹Sample of the first footnote

²Sample of the second footnote

Table 1: Sample table title

PART	DESCRIPTION
Dendrite	Input terminal
Axon	Output terminal
Soma	Cell body (contains cell nucleus)

10.4 TABLES

All tables must be centered, neat, clean and legible. Do not use hand-drawn tables. 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.

11 DEFAULT NOTATION

In an attempt to encourage standardized notation, we have included the notation file from the textbook, *Deep Learning* Goodfellow et al. (2016) available at https://github.com/goodfeli/dlbook_notation/. Use of this style is not required and can be disabled by commenting out `math_commands.tex`.

Numbers and Arrays

a	A scalar (integer or real)
\mathbf{a}	A vector
\mathbf{A}	A matrix
\mathbf{A}	A tensor
\mathbf{I}_n	Identity matrix with n rows and n columns
\mathbf{I}	Identity matrix with dimensionality implied by context
$\mathbf{e}^{(i)}$	Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at position i
$\text{diag}(\mathbf{a})$	A square, diagonal matrix with diagonal entries given by \mathbf{a}
a	A scalar random variable
\mathbf{a}	A vector-valued random variable
\mathbf{A}	A matrix-valued random variable

Sets and Graphs

216	\mathbb{A}	A set
217	\mathbb{R}	The set of real numbers
218	$\{0, 1\}$	The set containing 0 and 1
219	$\{0, 1, \dots, n\}$	The set of all integers between 0 and n
220	$[a, b]$	The real interval including a and b
221	$(a, b]$	The real interval excluding a but including b
222	$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of \mathbb{A} that are not in \mathbb{B}
223	\mathcal{G}	A graph
224	$Pa_{\mathcal{G}}(\mathbf{x}_i)$	The parents of \mathbf{x}_i in \mathcal{G}

Indexing

230	a_i	Element i of vector \mathbf{a} , with indexing starting at 1
231	\mathbf{a}_{-i}	All elements of vector \mathbf{a} except for element i
232	$A_{i,j}$	Element i, j of matrix \mathbf{A}
233	$\mathbf{A}_{i,:}$	Row i of matrix \mathbf{A}
234	$\mathbf{A}_{:,i}$	Column i of matrix \mathbf{A}
235	$\mathbf{A}_{i,j,k}$	Element (i, j, k) of a 3-D tensor \mathbf{A}
236	$\mathbf{A}_{:,:,i}$	2-D slice of a 3-D tensor
237	\mathbf{a}_i	Element i of the random vector \mathbf{a}

Calculus

243	$\frac{dy}{dx}$	Derivative of y with respect to x
244	$\frac{\partial y}{\partial x}$	Partial derivative of y with respect to x
245	$\nabla_{\mathbf{x}} y$	Gradient of y with respect to \mathbf{x}
246	$\nabla_{\mathbf{X}} y$	Matrix derivatives of y with respect to \mathbf{X}
247	$\nabla_{\mathbf{x}} \mathbf{y}$	Tensor containing derivatives of y with respect to \mathbf{X}
248	$\frac{\partial f}{\partial \mathbf{x}}$	Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
249	$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or $\mathbf{H}(f)(\mathbf{x})$	The Hessian matrix of f at input point \mathbf{x}
250	$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of \mathbf{x}
251	$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to \mathbf{x} over the set \mathbb{S}

Probability and Information Theory

270	$P(a)$	A probability distribution over a discrete variable
271	$p(a)$	A probability distribution over a continuous variable, or
272		over a variable whose type has not been specified
273		
274	$a \sim P$	Random variable a has distribution P
275	$\mathbb{E}_{x \sim P}[f(x)]$ or $\mathbb{E}f(x)$	Expectation of $f(x)$ with respect to $P(x)$
276		
277	$\text{Var}(f(x))$	Variance of $f(x)$ under $P(x)$
278	$\text{Cov}(f(x), g(x))$	Covariance of $f(x)$ and $g(x)$ under $P(x)$
279		
280	$H(x)$	Shannon entropy of the random variable x
281	$D_{\text{KL}}(P Q)$	Kullback-Leibler divergence of P and Q
282	$\mathcal{N}(x; \mu, \Sigma)$	Gaussian distribution over x with mean μ and covariance
283		Σ
284		

Functions

286	$f : \mathbb{A} \rightarrow \mathbb{B}$	The function f with domain \mathbb{A} and range \mathbb{B}
287		
288	$f \circ g$	Composition of the functions f and g
289	$f(x; \theta)$	A function of x parametrized by θ . (Sometimes we write
290		$f(x)$ and omit the argument θ to lighten notation)
291		
292	$\log x$	Natural logarithm of x
293	$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
294		
295	$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
296	$\ x\ _p$	L^p norm of x
297	$\ x\ $	L^2 norm of x
298		
299	x^+	Positive part of x , i.e., $\max(0, x)$
300	$\mathbf{1}_{\text{condition}}$	is 1 if the condition is true, 0 otherwise
301		

12 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.

13 PREPARING POSTSCRIPT OR PDF FILES

Please prepare PostScript or PDF files with paper size “US Letter”, and not, for example, “A4”. The `-t letter` option on `dvips` will produce US Letter files.

Consider directly generating PDF files using `pdflatex` (especially if you are a MiKTeX user). PDF figures must be substituted for EPS figures, however.

Otherwise, please generate your PostScript and PDF files with the following commands:

```
dvips mypaper.dvi -t letter -Ppdf -G0 -o mypaper.ps
ps2pdf mypaper.ps mypaper.pdf
```

13.1 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 using .eps graphics

```
\usepackage[dvips]{graphicx} ...
\includegraphics[width=0.8\linewidth]{myfile.eps}
```

or

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

for .pdf graphics. See section 4.4 in the graphics bundle documentation (<http://www.ctan.org/tex-archive/macros/latex/required/graphics/grfguide.ps>)

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

AUTHOR CONTRIBUTIONS

If you'd like to, you may include a section for author contributions as is done in many journals. This is optional and at the discretion of the authors.

ACKNOWLEDGMENTS

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper.

REFERENCES

- Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. MIT Press, 2007.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.

A APPENDIX

You may include other additional sections here.