
A Probabilistic Interpretation of Transformers

International Conference on Machine Learning (ICML 2021)

Anonymous Authors¹ Aeiaw Zzzzequal,to

Abstract

This document provides a basic paper template and submission guidelines. Abstracts must be a single paragraph, ideally between 4–6 sentences long. Gross violations will trigger corrections at the camera-ready phase.

1. Introduction

Transformers have reached state of the art results in language models, significantly outperforming LSTMs. One conceptual explanation for their increased performance is the ability of attention to utilize long range dependencies, whereas Recurrent Neural Networks were limited by having to encode past information within a fixed size hidden state. What this explanation does not explain is how certain architectural choices of transformers, specifically exponential dot product attention, also somewhat ambiguously referred to as softmax attention, outperforms alternatives.

Exponential dot product attention has been seen before in contrastive learning, though often with normalized embeddings before softmax is applied. In language models, the softmax probability was used in Word2Vec and later word embedding work and for memory networks it was used in Neural Turing Machines. Contrastive loss originated as Noisy Contrastive Estimation and continues to be used in seminal papers such as SimCLR, which achieved state of the art results, as have many variants based off SimCLR.

The successes of transformers has been verified empirically, but far research has focused on a theoretical explanation of transformers perform so well. We offer a probabilistic explanation, based off of distributions of the exponential family, for attention and contrastive probabilities. Expressing attention as an exponential family allows us to utilize related theory in statistics, machine learning, and statistical mechanics, offering insightful interpretations in to the trans-

former architecture.

We also explicitly state the limitations of our theory, noting that the modern Hopfield network interpretation shares many of these limitations. Moreover, for some of these limitations, we speculate connections between other areas of research which may reconcile the theoretical inconsistencies, motivating directions for future research.

2. Related Work

2.1. Contrastive Learning

2.2. Shortcut connections and dynamical systems

Long Short-Term Memory (LSTM) combined a shortcut connection to deal with the vanishing and exploding gradient problem along with gating functions to incorporate and forget information (Hochreiter & Schmidhuber, 1997). Residual connections similarly formulated the hidden layer as an update to an identity mapping, though without a gating mechanism (He et al., 2015). Recurrent Neural Networks have been interpreted as a discrete time approximation to a continuous dynamical system (Jaeger, 2001), where gating acts as a warping of time (Tallec & Ollivier, 2018). Residual connections have been interpreted as a discretized update to a differential equation (Weinan, 2017; Lu et al., 2020).

Interpreting residual networks as discretized differential equations, researchers have posed alternative methods for performing forward updates to converge to equilibrium points and backwards updates to the parameters from the equilibrium points (Chen et al., 2019; Bai et al., 2019). Further work has used monotone operator theory in convex analysis for solving for equilibrium points, interpreting layers as an operator (Winston & Kolter, 2021).

In a work most similar to ours, transformers have been interpreted as an update of modern hopfield networks and fixed points have been calculated with respect to a fixed set of patterns (Ramsauer et al., 2020). Our work similarly views the attention sublayer as an operator update over a class of discretized probability distribution, though with a changing set of patterns.

¹Anonymous Institution, Anonymous City, Anonymous Region, Anonymous Country. Correspondence to: Anonymous Author <anon.email@domain.com>.

2.3. Log normalizer and free energy

Partition functions, or the normalizer function, in statistical physics defines a normalization factor of the Hamiltonian with respect to a parameter defining the temperature. The Boltzmann distribution can be derived through Lagrange multipliers as the distribution which maximizes entropy with a conservation of energy constraint. Jaynes adapted the Boltzmann distributions to maximum entropy distributions with multiple expected statistics constraints by converting the maximum entropy problem into the dual problem of optimizing the log normalizer (Jaynes, 1982), which is known in statistical mechanics as free energy.

Variational methods have been used to approximate log probabilities of observations in machine learning, borrowing from ideas in statistical mechanics. By viewing the joint as an unnormalized probability distribution, the log normalizer is known as the evidence lower-bound, and it has connections to Helmholtz Free Energy (Hinton et al., 1995; Koller & Friedman, 2009).

The sum of exponents loss of AdaBoost (Collins et al., 2000) has been interpreted as the dual form of generalized KL divergence. The log sum of exponents is well known in convex optimization to be the dual form to the maximum entropy objective for a discrete probability distribution (Boyd & Vandenberghe, 2004). Notably, in the modern Hopfield network interpretation of transformers as part of the energy function (Ramsauer et al., 2020).

3. Electronic Submission

Submission to ICML 2021 will be entirely electronic, via a web site (not email). Information about the submission process and L^AT_EX templates are available on the conference web site at:

<http://icml.cc/>

The guidelines below will be enforced for initial submissions and camera-ready copies. Here is a brief summary:

- Submissions must be in PDF.
- Submitted papers can be up to eight pages long, not including references, plus unlimited space for references. Accepted papers can be up to nine pages long, not including references, to allow authors to address reviewer comments. Any paper exceeding this length will automatically be rejected.
- **Do not include author information or acknowledgments** in your initial submission.
- Your paper should be in **10 point Times font**.
- Make sure your PDF file only uses Type-1 fonts.

- Place figure captions *under* the figure (and omit titles from inside the graphic file itself). Place table captions *over* the table.
- References must include page numbers whenever possible and be as complete as possible. Place multiple citations in chronological order.
- Do not alter the style template; in particular, do not compress the paper format by reducing the vertical spaces.
- Keep your abstract brief and self-contained, one paragraph and roughly 4–6 sentences. Gross violations will require correction at the camera-ready phase. The title should have content words capitalized.

3.1. Submitting Papers

Paper Deadline: The deadline for paper submission that is advertised on the conference website is strict. If your full, anonymized, submission does not reach us on time, it will not be considered for publication.

Anonymous Submission: ICML uses double-blind review: no identifying author information may appear on the title page or in the paper itself. Section 4.3 gives further details.

Simultaneous Submission: ICML will not accept any paper which, at the time of submission, is under review for another conference or has already been published. This policy also applies to papers that overlap substantially in technical content with conference papers under review or previously published. ICML submissions must not be submitted to other conferences and journals during ICML’s review period. Informal publications, such as technical reports or papers in workshop proceedings which do not appear in print, do not fall under these restrictions.

Authors must provide their manuscripts in **PDF** format. Furthermore, please make sure that files contain only embedded Type-1 fonts (e.g., using the program `pdf fonts` in linux or using File/DocumentProperties/Fonts in Acrobat). Other fonts (like Type-3) might come from graphics files imported into the document.

Authors using **Word** must convert their document to PDF. Most of the latest versions of Word have the facility to do this automatically. Submissions will not be accepted in Word format or any format other than PDF. Really. We’re not joking. Don’t send Word.

Those who use L^AT_EX should avoid including Type-3 fonts. Those using `latex` and `dvips` may need the following two commands:

```
dvips -Ppdf -tletter -G0 -o paper.ps paper.dvi
ps2pdf paper.ps
```

It is a zero following the “-G”, which tells dvips to use the config.pdf file. Newer T_EX distributions don’t always need this option.

Using pdf_lat_ex rather than latex, often gives better results. This program avoids the Type-3 font problem, and supports more advanced features in the microtype package.

Graphics files should be a reasonable size, and included from an appropriate format. Use vector formats (.eps/.pdf) for plots, lossless bitmap formats (.png) for raster graphics with sharp lines, and jpeg for photo-like images.

The style file uses the hyperref package to make clickable links in documents. If this causes problems for you, add nohyperref as one of the options to the icml2021 usepackage statement.

3.2. Submitting Final Camera-Ready Copy

The final versions of papers accepted for publication should follow the same format and naming convention as initial submissions, except that author information (names and affiliations) should be given. See Section 4.3.2 for formatting instructions.

The footnote, “Preliminary work. Under review by the International Conference on Machine Learning (ICML). Do not distribute.” must be modified to “*Proceedings of the 38th International Conference on Machine Learning*, Online, PMLR 139, 2021. Copyright 2021 by the author(s).”

For those using the L^AT_EX style file, this change (and others) is handled automatically by simply changing \usepackage{icml2021} to

```
\usepackage[accepted]{icml2021}
```

Authors using **Word** must edit the footnote on the first page of the document themselves.

Camera-ready copies should have the title of the paper as running head on each page except the first one. The running title consists of a single line centered above a horizontal rule which is 1 point thick. The running head should be centered, bold and in 9 point type. The rule should be 10 points above the main text. For those using the L^AT_EX style file, the original title is automatically set as running head using the fancyhdr package which is included in the ICML 2021 style file package. In case that the original title exceeds the size restrictions, a shorter form can be supplied by using

```
\icmltitlerunning{...}
```

just before \begin{document}. Authors using **Word** must edit the header of the document themselves.

4. Format of the Paper

All submissions must follow the specified format.

4.1. Dimensions

The text of the paper should be formatted in two columns, with an overall width of 6.75 inches, height of 9.0 inches, and 0.25 inches between the columns. The left margin should be 0.75 inches and the top margin 1.0 inch (2.54 cm). The right and bottom margins will depend on whether you print on US letter or A4 paper, but all final versions must be produced for US letter size.

The paper body should be set in 10 point type with a vertical spacing of 11 points. Please use Times typeface throughout the text.

4.2. Title

The paper title should be set in 14 point bold type and centered between two horizontal rules that are 1 point thick, with 1.0 inch between the top rule and the top edge of the page. Capitalize the first letter of content words and put the rest of the title in lower case.

4.3. Author Information for Submission

ICML uses double-blind review, so author information must not appear. If you are using L^AT_EX and the icml2021.sty file, use \icmlauthor{...} to specify authors and \icmlaffiliation{...} to specify affiliations. (Read the TeX code used to produce this document for an example usage.) The author information will not be printed unless accepted is passed as an argument to the style file. Submissions that include the author information will not be reviewed.

4.3.1. SELF-CITATIONS

If you are citing published papers for which you are an author, refer to yourself in the third person. In particular, do not use phrases that reveal your identity (e.g., “in previous work (?), we have shown ...”).

Do not anonymize citations in the reference section. The only exception are manuscripts that are not yet published (e.g., under submission). If you choose to refer to such unpublished manuscripts (?), anonymized copies have to be submitted as Supplementary Material via CMT. However, keep in mind that an ICML paper should be self contained and should contain sufficient detail for the reviewers to evaluate the work. In particular, reviewers are not required to look at the Supplementary Material when writing their review.

4.3.2. CAMERA-READY AUTHOR INFORMATION

If a paper is accepted, a final camera-ready copy must be prepared. For camera-ready papers, author information should start 0.3 inches below the bottom rule surrounding the title. The authors' names should appear in 10 point bold type, in a row, separated by white space, and centered. Author names should not be broken across lines. Unbolded superscripted numbers, starting 1, should be used to refer to affiliations.

Affiliations should be numbered in the order of appearance. A single footnote block of text should be used to list all the affiliations. (Academic affiliations should list Department, University, City, State/Region, Country. Similarly for industrial affiliations.)

Each distinct affiliations should be listed once. If an author has multiple affiliations, multiple superscripts should be placed after the name, separated by thin spaces. If the authors would like to highlight equal contribution by multiple first authors, those authors should have an asterisk placed after their name in superscript, and the term “*Equal contribution” should be placed in the footnote block ahead of the list of affiliations. A list of corresponding authors and their emails (in the format Full Name <email@domain.com>) can follow the list of affiliations. Ideally only one or two names should be listed.

A sample file with author names is included in the ICML2021 style file package. Turn on the `[accepted]` option to the stylefile to see the names rendered. All of the guidelines above are implemented by the `LATEX` style file.

4.4. Abstract

The paper abstract should begin in the left column, 0.4 inches below the final address. The heading ‘Abstract’ should be centered, bold, and in 11 point type. The abstract body should use 10 point type, with a vertical spacing of 11 points, and should be indented 0.25 inches more than normal on left-hand and right-hand margins. Insert 0.4 inches of blank space after the body. Keep your abstract brief and self-contained, limiting it to one paragraph and roughly 4–6 sentences. Gross violations will require correction at the camera-ready phase.

4.5. Partitioning the Text

You should organize your paper into sections and paragraphs to help readers place a structure on the material and understand its contributions.

4.5.1. SECTIONS AND SUBSECTIONS

Section headings should be numbered, flush left, and set in 11 pt bold type with the content words capitalized. Leave

0.25 inches of space before the heading and 0.15 inches after the heading.

Similarly, subsection headings should be numbered, flush left, and set in 10 pt bold type with the content words capitalized. Leave 0.2 inches of space before the heading and 0.13 inches afterward.

Finally, subsubsection headings should be numbered, flush left, and set in 10 pt small caps with the content words capitalized. Leave 0.18 inches of space before the heading and 0.1 inches after the heading.

Please use no more than three levels of headings.

4.5.2. PARAGRAPHS AND FOOTNOTES

Within each section or subsection, you should further partition the paper into paragraphs. Do not indent the first line of a given paragraph, but insert a blank line between succeeding ones.

You can use footnotes¹ to provide readers with additional information about a topic without interrupting the flow of the paper. Indicate footnotes with a number in the text where the point is most relevant. Place the footnote in 9 point type at the bottom of the column in which it appears. Precede the first footnote in a column with a horizontal rule of 0.8 inches.²

4.6. Figures

You may want to include figures in the paper to illustrate your approach and results. Such artwork should be centered, legible, and separated from the text. Lines should be dark and at least 0.5 points thick for purposes of reproduction, and text should not appear on a gray background.

Label all distinct components of each figure. If the figure takes the form of a graph, then give a name for each axis and include a legend that briefly describes each curve. Do not include a title inside the figure; instead, the caption should serve this function.

Number figures sequentially, placing the figure number and caption *after* the graphics, with at least 0.1 inches of space before the caption and 0.1 inches after it, as in Figure 1. The figure caption should be set in 9 point type and centered unless it runs two or more lines, in which case it should be flush left. You may float figures to the top or bottom of a column, and you may set wide figures across both columns (use the environment `figure*` in `LATEX`). Always place two-column figures at the top or bottom of the page.

¹Footnotes should be complete sentences.

²Multiple footnotes can appear in each column, in the same order as they appear in the text, but spread them across columns and pages if possible.

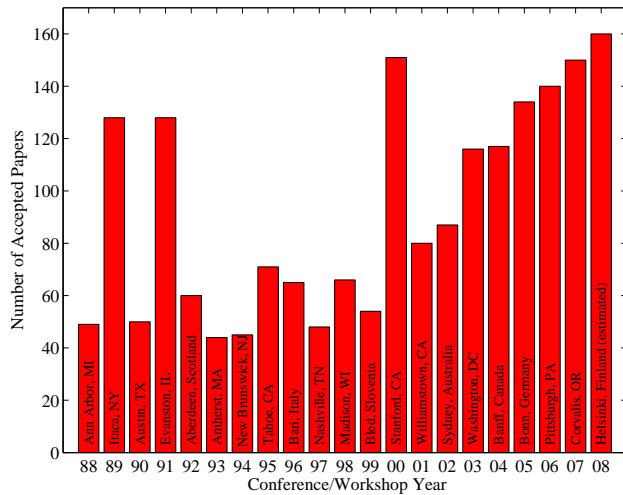


Figure 1. Historical locations and number of accepted papers for International Machine Learning Conferences (ICML 1993 – ICML 2008) and International Workshops on Machine Learning (ML 1988 – ML 1992). At the time this figure was produced, the number of accepted papers for ICML 2008 was unknown and instead estimated.

Algorithm 1 Bubble Sort

Input: data x_i , size m

repeat

 Initialize $noChange = true$.

for $i = 1$ **to** $m - 1$ **do**

if $x_i > x_{i+1}$ **then**

 Swap x_i and x_{i+1}

$noChange = false$

end if

end for

until $noChange$ is $true$

4.7. Algorithms

If you are using L^AT_EX, please use the “algorithm” and “algorithmic” environments to format pseudocode. These require the corresponding stylefiles, `algorithm.sty` and `algorithmic.sty`, which are supplied with this package. Algorithm 1 shows an example.

4.8. Tables

You may also want to include tables that summarize material. Like figures, these should be centered, legible, and numbered consecutively. However, place the title *above* the table with at least 0.1 inches of space before the title and the same after it, as in Table 1. The table title should be set in 9 point type and centered unless it runs two or more lines, in which case it should be flush left.

Table 1. Classification accuracies for naive Bayes and flexible Bayes on various data sets.

DATA SET	NAIVE	FLEXIBLE	BETTER?
BREAST	95.9± 0.2	96.7± 0.2	✓
CLEVELAND	83.3± 0.6	80.0± 0.6	×
GLASS2	61.9± 1.4	83.8± 0.7	✓
CREDIT	74.8± 0.5	78.3± 0.6	
HORSE	73.3± 0.9	69.7± 1.0	×
META	67.1± 0.6	76.5± 0.5	✓
PIMA	75.1± 0.6	73.9± 0.5	
VEHICLE	44.9± 0.6	61.5± 0.4	✓

Tables contain textual material, whereas figures contain graphical material. Specify the contents of each row and column in the table’s topmost row. Again, you may float tables to a column’s top or bottom, and set wide tables across both columns. Place two-column tables at the top or bottom of the page.

4.9. Citations and References

Please use APA reference format regardless of your formatter or word processor. If you rely on the L^AT_EX bibliographic facility, use `natbib.sty` and `icml2021.bst` included in the style-file package to obtain this format.

Citations within the text should include the authors’ last names and year. If the authors’ names are included in the sentence, place only the year in parentheses, for example when referencing Arthur Samuel’s pioneering work (?). Otherwise place the entire reference in parentheses with the authors and year separated by a comma (?). List multiple references separated by semicolons (???). Use the ‘et al.’ construct only for citations with three or more authors or after listing all authors to a publication in an earlier reference (?).

Authors should cite their own work in the third person in the initial version of their paper submitted for blind review. Please refer to Section 4.3 for detailed instructions on how to cite your own papers.

Use an unnumbered first-level section heading for the references, and use a hanging indent style, with the first line of the reference flush against the left margin and subsequent lines indented by 10 points. The references at the end of this document give examples for journal articles (?), conference publications (?), book chapters (?), books (?), edited volumes (?), technical reports (?), and dissertations (?).

Alphabetize references by the surnames of the first authors, with single author entries preceding multiple author entries. Order references for the same authors by year of publication, with the earliest first. Make sure that each reference includes all relevant information (e.g., page numbers).

Please put some effort into making references complete, presentable, and consistent. If using bibtex, please protect capital letters of names and abbreviations in titles, for example, use {B}ayesian or {L}ipschitz in your .bib file.

Software and Data

If a paper is accepted, we strongly encourage the publication of software and data with the camera-ready version of the paper whenever appropriate. This can be done by including a URL in the camera-ready copy. However, **do not** include URLs that reveal your institution or identity in your submission for review. Instead, provide an anonymous URL or upload the material as “Supplementary Material” into the CMT reviewing system. Note that reviewers are not required to look at this material when writing their review.

Acknowledgements

Do not include acknowledgements in the initial version of the paper submitted for blind review.

If a paper is accepted, the final camera-ready version can (and probably should) include acknowledgements. In this case, please place such acknowledgements in an unnumbered section at the end of the paper. Typically, this will include thanks to reviewers who gave useful comments, to colleagues who contributed to the ideas, and to funding agencies and corporate sponsors that provided financial support.

References

- Bai, S., Kolter, J. Z., and Koltun, V. Deep equilibrium models, 2019.
- Boyd, S. and Vandenberghe, L. *Convex Optimization*. Cambridge University Press, March 2004. ISBN 0521833787. URL <http://www.amazon.com/exec/obidos/redirect?tag=citeulike-20&path=ASIN/0521833787>.
- Chen, R. T. Q., Rubanova, Y., Bettencourt, J., and Duvenaud, D. Neural ordinary differential equations, 2019.
- Collins, M., Schapire, R. E., and Singer, Y. Logistic regression, adaboost and bregman distances. In *Computational Learning Theory*, pp. 158–169, 2000. URL citeseer.nj.nec.com/article/collins00logistic.html.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition, 2015.
- Hinton, G. E., Dayan, P., Frey, B. J., and Neal, R. M. The wake-sleep algorithm for unsupervised neural networks. *Science*, 268:1158–1161, 1995.

Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

Jaeger, H. The echo state approach to analysing and training recurrent neural networks. *GMD-Report 148, German National Research Institute for Computer Science*, 01 2001.

Jaynes, E. On the rationale of maximum-entropy methods. *Proceedings of the IEEE*, 70(9):939–952, 1982. doi: 10.1109/PROC.1982.12425.

Koller, D. and Friedman, N. *Probabilistic Graphical Models: Principles and Techniques*. Adaptive computation and machine learning. MIT Press, 2009. ISBN 9780262013192. URL <https://books.google.co.in/books?id=7dzpHCHzNQ4C>.

Lu, Y., Zhong, A., Li, Q., and Dong, B. Beyond finite layer neural networks: Bridging deep architectures and numerical differential equations, 2020.

Ramsauer, H., Schäfl, B., Lehner, J., Seidl, P., Widrich, M., Gruber, L., Holzleitner, M., Pavlovic, M., Sandve, G. K., Greiff, V., Kreil, D. P., Kopp, M., Klambauer, G., Brandstetter, J., and Hochreiter, S. Hopfield networks is all you need. *CoRR*, abs/2008.02217, 2020. URL <https://arxiv.org/abs/2008.02217>.

Tallec, C. and Ollivier, Y. Can recurrent neural networks warp time?, 2018.

Weinan, E. A proposal on machine learning via dynamical systems. 2017.

Winston, E. and Kolter, J. Z. Monotone operator equilibrium networks, 2021.

A. Do not have an appendix here

Do not put content after the references. Put anything that you might normally include after the references in a separate supplementary file.

We recommend that you build supplementary material in a separate document. If you must create one PDF and cut it up, please be careful to use a tool that doesn’t alter the margins, and that doesn’t aggressively rewrite the PDF file. pdftk usually works fine.

Please do not use Apple’s preview to cut off supplementary material. In previous years it has altered margins, and created headaches at the camera-ready stage.