# IMPROVED METHODS FOR STATIC MODEL PRUNING

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

Static model pruning is presented as a performance optimization technique for large language and vision models. The approach aims to identify and remove neurons, connections unlikely to lead to expected generation results for typical user queries. The goal is to obtain a much smaller model that can quickly return results almost as good as those of the unpruned ones. Through careful analysis of pretrained weights, bias, activations and user queries, an initial mathematical model based on certain probabilities obtained from the environment is developed to improve on previous results for pruned model size, achieving significant improvement in most cases. This paper explores and compares to previously proposed approaches that perform pruning based on other factors.

## 1 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.

### 1.1 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.

### 1.2 RETRIEVAL OF STYLE FILES

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

$$\text{http://www.iclr.cc/}$$

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 LaTeX2e). 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 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 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 4 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.

## 3 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.

### 3.1 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.

#### 3.1.1 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.

## 4 CITATIONS, FIGURES, TABLES, REFERENCES

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

### 4.1 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 **?** for more information."). Otherwise, the citation should be in parenthesis using `\citep{}` (as in "Deep learning shows promise to make progress towards AI (**?**).").

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.

### 4.2 FOOTNOTES

Indicate footnotes with a number[1] 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]

### 4.3 FIGURES

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.

---

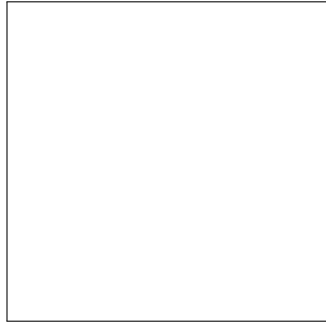[1] Sample of the first footnote
[2] Sample of the second footnote

Figure 1: Sample figure caption a.



Figure 2: Sample figure caption b.

### 4.4 TABLES

Place one line space before the table title, one line space after the table title, and one line space after the table.

## 5 DEFAULT NOTATION

In an attempt to encourage standardized notation, we have included the notation file from the textbook, *Deep Learning* **?** 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**

Figure 3: Sample figure caption c.

Table 1: Perplexity on pruned model (Llama-7B) from human domain experts

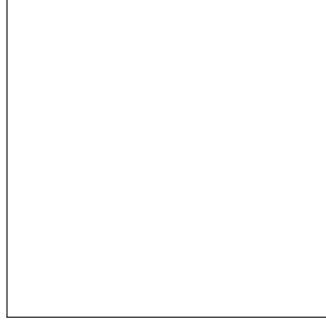| Pruned Level | Wanda |
|---|---|
| 0.01 | NA |
| 0.05 | NA |
| 0.10 | 5.696 |
| 0.20 | 5.817 |
| 0.30 | 5.999 |
| 0.40 | 6.387 |
| 0.50 | 7.257 |
| 0.60 | 10.691 |
| 0.70 | 84.905 |
| 0.80 | 5782.432 |
| 0.90 | 19676.668 |
| 0.95 | 28309.178 |
| 0.99 | 108234.484 |

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

**Sets and Graphs**

4

Table 2: Effectiveness of the weights as a major pruning measure

| Pruned Level | Prune by Weights |
|---|---|
| 0.01 | NA |
| 0.05 | NA |
| 0.10 | 5.806 |
| 0.20 | 6.020 |
| 0.30 | 6.669 |
| 0.40 | 8.601 |
| 0.50 | 17.285 |
| 0.60 | 559.987 |
| 0.70 | 48414.551 |
| 0.80 | 132175.578 |
| 0.90 | 317879.250 |
| 0.95 | 273552.281 |
| 0.99 | 222543.047 |

Table 3: Effectiveness of the bias as a major pruning indicator

| Pruned Level | Prune by Bias |
|---|---|
| 0.01 | NA |
| 0.05 | NA |
| 0.10 | NA |
| 0.20 | NA |
| 0.30 | NA |
| 0.40 | NA |
| 0.50 | NA |
| 0.60 | NA |
| 0.70 | NA |
| 0.80 | NA |
| 0.90 | NA |
| 0.95 | NA |
| 0.99 | NA |

| | |
|---|---|
| $\mathbb{A}$ | A set |
| $\mathbb{R}$ | The set of real numbers |
| $\{0, 1\}$ | The set containing 0 and 1 |
| $\{0, 1, \ldots, n\}$ | The set of all integers between $0$ and $n$ |
| $[a, b]$ | The real interval including $a$ and $b$ |
| $(a, b]$ | The real interval excluding $a$ but including $b$ |
| $\mathbb{A} \backslash \mathbb{B}$ | Set subtraction, i.e., the set containing the elements of $\mathbb{A}$ that are not in $\mathbb{B}$ |
| $\mathcal{G}$ | A graph |
| $Pa_{\mathcal{G}}(\mathrm{x}_i)$ | The parents of $\mathrm{x}_i$ in $\mathcal{G}$ |

**Indexing**

5

Table 4: One pass code generation and effectiveness evaluation

| Number | Core Idea |
|---|---|
| 01 | Gradient Sensitive Pruning |
| 02 | L1 Norm Pruning |
| 03 | Structured Pruning |
| 04 | K-means Clustering Pruning |
| 05 | Random Pruning |
| 06 | Random Pattern Pruning |
| 07 | Variational Dropout Pruning |
| 08 | Gradient based Pruning |
| 09 | Elastic Weight Consolidation Pruning |
| 10 | Dynamic Pruning with Reinforcement Learning |

Table 5: Perplexity on pruned model (llama-7B) from AIGC domain expert (o1)

| Pruned Level | aigc algorithm 2 |
|---|---|
| 0.50 | 193740.406 |
| 0.60 | 110879.422 |
| 0.70 | 174815.859 |
| 0.80 | 287734.844 |
| 0.90 | 157028.844 |
| 0.95 | 90220.781 |
| 0.99 | 991519.125 |

| | |
|---|---|
| $a_i$ | Element $i$ of vector $\boldsymbol{a}$, with indexing starting at 1 |
| $a_{-i}$ | All elements of vector $\boldsymbol{a}$ except for element $i$ |
| $A_{i,j}$ | Element $i, j$ of matrix $\boldsymbol{A}$ |
| $\boldsymbol{A}_{i,:}$ | Row $i$ of matrix $\boldsymbol{A}$ |
| $\boldsymbol{A}_{:,i}$ | Column $i$ of matrix $\boldsymbol{A}$ |
| $A_{i,j,k}$ | Element $(i, j, k)$ of a 3-D tensor $\mathbf{A}$ |
| $\mathbf{A}_{:,:,i}$ | 2-D slice of a 3-D tensor |
| $\mathrm{a}_i$ | Element $i$ of the random vector $\mathbf{a}$ |

**Calculus**

Table 6: Effect of pruned model (OPT-1.3B) applying to downstream task - text generation

| Pruned Level | Perplexity |
|---|---|
| 0.00 | NA |
| 0.50 | 19.191 |
| 0.60 | 23.205 |
| 0.70 | 44.246 |
| 0.80 | 364.304 |
| 0.90 | 3772.829 |
| 0.95 | 8892.167 |
| 0.99 | 22548.809 |

Table 7: (TODO: Running Time for each pruning algorithm)

| Number | Running Time |
|---|---|
| 01 | TBA |
| 02 | TBA |
| 03 | TBA |
| 04 | TBA |
| 05 | TBA |
| 06 | TBA |
| 07 | TBA |
| 08 | TBA |
| 09 | TBA |
| 10 | TBA |

| | |
|---|---|
| $\dfrac{dy}{dx}$ | Derivative of $y$ with respect to $x$ |
| $\dfrac{\partial y}{\partial x}$ | Partial derivative of $y$ with respect to $x$ |
| $\nabla_{\boldsymbol{x}} y$ | Gradient of $y$ with respect to $\boldsymbol{x}$ |
| $\nabla_{\boldsymbol{X}} y$ | Matrix derivatives of $y$ with respect to $\boldsymbol{X}$ |
| $\nabla_{\mathbf{X}} y$ | Tensor containing derivatives of $y$ with respect to $\mathbf{X}$ |
| $\dfrac{\partial f}{\partial \boldsymbol{x}}$ | Jacobian matrix $\boldsymbol{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \to \mathbb{R}^m$ |
| $\nabla_{\boldsymbol{x}}^2 f(\boldsymbol{x})$ or $\boldsymbol{H}(f)(\boldsymbol{x})$ | The Hessian matrix of $f$ at input point $\boldsymbol{x}$ |
| $\int f(\boldsymbol{x}) d\boldsymbol{x}$ | Definite integral over the entire domain of $\boldsymbol{x}$ |
| $\int_{\mathbb{S}} f(\boldsymbol{x}) d\boldsymbol{x}$ | Definite integral with respect to $\boldsymbol{x}$ over the set $\mathbb{S}$ |

**Probability and Information Theory**

Table 8: (TODO: End-to-end model evaluation)

| Number | Inspiration Score |
|--------|-------------------|
| 01 | TBA |
| 02 | TBA |
| 03 | TBA |
| 04 | TBA |
| 05 | TBA |
| 06 | TBA |
| 07 | TBA |
| 08 | TBA |
| 09 | TBA |
| 10 | TBA |

| | |
|---|---|
| $P(\mathrm{a})$ | A probability distribution over a discrete variable |
| $p(\mathrm{a})$ | A probability distribution over a continuous variable, or over a variable whose type has not been specified |
| $\mathrm{a} \sim P$ | Random variable a has distribution $P$ |
| $\mathbb{E}_{\mathrm{x} \sim P}[f(x)]$ or $\mathbb{E}f(x)$ | Expectation of $f(x)$ with respect to $P(\mathrm{x})$ |
| $\mathrm{Var}(f(x))$ | Variance of $f(x)$ under $P(\mathrm{x})$ |
| $\mathrm{Cov}(f(x), g(x))$ | Covariance of $f(x)$ and $g(x)$ under $P(\mathrm{x})$ |
| $H(\mathrm{x})$ | Shannon entropy of the random variable x |
| $D_{\mathrm{KL}}(P\|Q)$ | Kullback-Leibler divergence of P and Q |
| $\mathcal{N}(\boldsymbol{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$ | Gaussian distribution over $\boldsymbol{x}$ with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$ |

**Functions**

| | |
|---|---|
| $f : \mathbb{A} \to \mathbb{B}$ | The function $f$ with domain $\mathbb{A}$ and range $\mathbb{B}$ |
| $f \circ g$ | Composition of the functions $f$ and $g$ |
| $f(\boldsymbol{x}; \boldsymbol{\theta})$ | A function of $\boldsymbol{x}$ parametrized by $\boldsymbol{\theta}$. (Sometimes we write $f(\boldsymbol{x})$ and omit the argument $\boldsymbol{\theta}$ to lighten notation) |
| $\log x$ | Natural logarithm of $x$ |
| $\sigma(x)$ | Logistic sigmoid, $\dfrac{1}{1 + \exp(-x)}$ |
| $\zeta(x)$ | Softplus, $\log(1 + \exp(x))$ |
| $\|\boldsymbol{x}\|_p$ | $L^p$ norm of $\boldsymbol{x}$ |
| $\|\boldsymbol{x}\|$ | $L^2$ norm of $\boldsymbol{x}$ |
| $x^+$ | Positive part of $x$, i.e., $\max(0, x)$ |
| $\mathbf{1}_{\mathrm{condition}}$ | is 1 if the condition is true, 0 otherwise |

## 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.

8

# 7 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
```

## 7.1 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.

# A APPENDIX

You may include other additional sections here.