

MULTIMODAL FOUNDATIONAL MODEL FOR DOCUMENT QUESTION ANSWERING

Anonymous authors

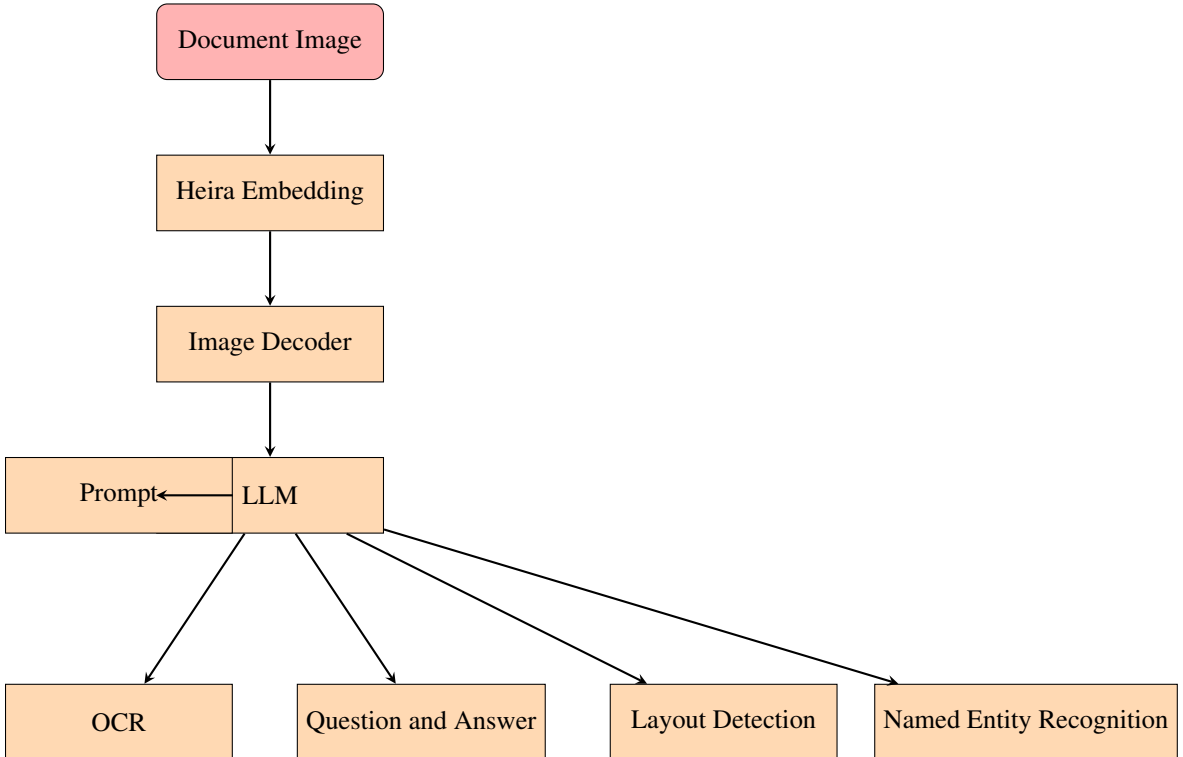
Paper under double-blind review

ABSTRACT

Multimodal promptable foundational model for document question answering. Introduce Document-JEPA objective.

1 INTRODUCTION

We introduce a promptable document question answering model where the focus is on the visual segmentation task, where the model is given a document and a question, and the model is expected to segment the document based on that question. The model is pre-trained on a large corpus of documents and OCRs with the Document Joint Embedding Predictive Architecture (D-JEPA) objective. The D-JEPA objective here masks part of the document, and is predicted based on the joint image and text embedding, while a masked text is predicted based on the image embedding.



1.1 IMAGE ENCODER

The Hiera model is a novel hierarchical vision transformer designed with a focus on simplicity and efficiency for computer vision tasks. It emerges as a response to the increasing complexity of existing hierarchical vision transformers, such as MViTv2. The Hiera architecture retains the essential components of the ViT framework while eliminating non-essential features that do not contribute to its overall performance, thereby achieving faster operation and competitive accuracy.

Hiera Architecture: Hiera is built upon the foundation of MViTv2, from which the authors carefully removed various specialized components often added in an attempt to enhance performance. Specifically, Hiera eliminates:

1. Convolutions: These have traditionally been used in hierarchical models to impart spatial biases.
2. Shifted or Cross-shaped Windows: These complex mechanisms are designed to improve locality and capture spatial relationships within the data.
3. Decomposed Relative Position Embeddings: These embeddings are typically utilized to encode positional information in a more sophisticated manner.

Instead of these modules, Hiera leverages the Masked Autoencoder (MAE) for strong pretraining, which allows for effective learning of spatial biases. As a result, Hiera comprises of entirely standard ViT blocks, significantly reducing the model’s complexity while retaining or even increasing accuracy compared to its predecessor.

Hiera Training Strategy: The training strategy for Hiera centers around the use of the MAE method, which facilitates efficient sparse pretraining. This approach allows Hiera to learn effective spatial representations without relying on complex architecture components. In comparative trials, it was found that when training Hiera with MAE pretraining, the model not only simplifies the architecture but also enhances performance.

Hiera achieves remarkable efficiency:

1. It is up to **2.4×** faster on images than MViTv2.
2. For image classification, Hiera-L is reported to be **3×** faster to train than a supervised MViTv2-L model while achieving superior accuracy.
3. Remarkably, it performs effectively even in smaller configurations, demonstrating strong capabilities where convolutions typically dominated earlier designs.

To adapt the Hiera model for 1024x1024 images, several adjustments would be necessary:

1. Input Rescaling: Ensure that the input layers of the model can handle the larger 1024x1024 images. This might require modification of any layers or preprocessing steps that specifically target a 224x224 resolution.
2. Adjusting Patch Size: The model’s patch size may need to be recalibrated to appropriately process the larger images effectively. A larger patch size can enhance efficiency given the increased input size.
3. FLOPs and Memory Considerations: Analyze and possibly upgrade the computational resources (e.g., GPU memory) since higher resolutions will significantly increase the number of computations (FLOPs) and memory usage. Adjusting the batch size may also be necessary.
4. Network Architecture: Depending on the model’s performance with the larger size, modifications to the number of blocks, channels, or heads in the Hiera architecture may be necessary. For instance, it could involve scaling up certain configurations (similar to Hiera-B+ or Hiera-L) to maintain performance.
5. Fine-tuning Pretraining: If the model was pretrained on smaller images, it may need to undergo a new pretraining phase on the 1024x1024 images or include an additional fine-tuning step specifically for this resolution to optimize its performance at the new scale.
6. Evaluation Metrics: Adjust evaluation protocols and metrics to appropriately assess model performance at the 1024x1024 resolution, ensuring that you are capturing the accuracy, speed, and efficiency of the model in this new context.

By implementing these adjustments, the Hiera model can be effectively utilized for processing 1024x1024 images while aiming to maintain or improve performance.

Reference: (Ryali et al., 2023, Hiera Model)

In this work, we use a different training strategy, known as the I-JEPA.

I-JEPA Training Strategy: The training strategy for the Image-based Joint-Embedding Predictive Architecture (I-JEPA) is distinct from traditional methods like Masked Autoencoders (MAE) due to its focus on abstract representation learning without reliance on hand-crafted data augmentations.

1. **Architecture Overview:** I-JEPA uses a Vision Transformer (ViT) for encoding and predicting representations. The architecture includes a context encoder and a target encoder, with the purpose of predicting representations of various target blocks from a single context block within the same image.
2. **Predictive Learning:** In I-JEPA, the training objective revolves around predicting the representations of target blocks based on the context block. Unlike generative methods that reconstruct inputs at the pixel level, I-JEPA makes predictions in abstract representation space, focusing on high-level semantic features.
3. **Multi-block Masking Strategy:** The model employs a multi-block masking strategy that enables learning from larger, informative context blocks, which helps in capturing more spatial context and reducing the focus on pixel-level details.
4. **Training Efficiency:** One of the key advantages of I-JEPA is its efficiency; it converges in approximately five times fewer iterations compared to traditional pixel reconstruction methods. This allows I-JEPA to achieve strong performance with significantly less computational effort and training duration.
5. **Adaptability:** I-JEPA scales effectively when trained on larger datasets and benefits from larger model sizes, improving downstream performance even when fine-tuning or adapting the model for specific tasks.

Comparison with MAE:

1. **Data Augmentation Reliance:** MAE relies heavily on hand-crafted data augmentations, processing multiple views of each image during pretraining, while I-JEPA requires only a single view. This gives I-JEPA a scalability advantage and simplifies the training process.
2. **Learning Paradigm:** MAE's approach is generative and focuses on reconstructing masked portions of images at the pixel level, while I-JEPA's focus is on predicting high-level semantic representations in a more abstract sense, which improves the quality of learned features.
3. **Performance:** Empirical evidence indicates that I-JEPA outperforms MAE in several benchmarks, including linear probing on ImageNet-1K and semi-supervised tasks, while also requiring fewer training epochs and less computational power.
4. **Training Dynamics:** Although I-JEPA might introduce some overhead from computing targets in representation space (approximately 7% slower per iteration), its reduced overall training time results in significant computational savings in practice compared to MAE.
5. **Conclusion:** In summary, I-JEPA represents a novel approach to self-supervised learning focused on semantic representation while improving efficiency and performance over traditional methods like MAE by leveraging a predictive approach to learning without extensive reliance on data augmentations.
6. **ImageNet-1K Linear Evaluation:** I-JEPA significantly improves linear probing performance on the ImageNet-1K benchmark compared to MAE. I-JEPA achieves 77.3% Top-1 accuracy with a ViT-H/14 model over 300 epochs, while MAE achieves 71.5% Top-1 accuracy with a ViT-H/14 model over 1600 epochs. This indicates that I-JEPA can achieve superior performance with considerably less training time.
7. **Low-Shot ImageNet-1K (1% Labels):** I-JEPA also outperforms MAE in this scenario, demonstrating better adaptation for image classification using only 1% of the available labels. I-JEPA requires fewer pretraining epochs compared to MAE while still achieving competitive results.
8. **Model Efficiency:** Although I-JEPA is about 7% slower per iteration due to additional overhead from computing targets in representation space, it converges in roughly 5 times fewer iterations than MAE. This leads to significant computational savings in practice.

9. **Fine-Tuning on Full ImageNet:** During fine-tuning, I-JEPA achieves a Top-1 accuracy of 87.1%, which is just 0.7% lower than MAE’s 87.8% accuracy. However, I-JEPA is trained for 5.3 times fewer epochs, indicating better efficiency in utilizing training resources.

Overall, I-JEPA consistently shows enhanced performance across various benchmarks while minimizing reliance on extensive epochs of training, making it a more efficient alternative to MAE in self-supervised learning for image representations.

Reference: <https://arxiv.org/pdf/2301.08243>

1.2 TEXT ENCODER

Large decoder-only language models (LLMs) are the state-of-the-art models on most of today’s NLP tasks and benchmarks. Yet, the community is only slowly adopting these models for text embedding tasks, which require rich contextualized representations. In this work, we introduce LLM2Vec, a simple unsupervised approach that can transform any decoder-only LLM into a strong text encoder. LLM2Vec consists of three simple steps: 1) enabling bidirectional attention, 2) masked next token prediction, and 3) unsupervised contrastive learning. We demonstrate the effectiveness of LLM2Vec by applying it to 3 popular LLMs ranging from 1.3B to 7B parameters and evaluate the transformed models on English word- and sequence-level tasks. We outperform encoder-only models by a large margin on word-level tasks and reach a new unsupervised state-of-the-art performance on the Massive Text Embeddings Benchmark (MTEB). Moreover, when combining LLM2Vec with supervised contrastive learning, we achieve state-of-the-art performance on MTEB among models that train only on publicly available data. Our strong empirical results and extensive analysis demonstrate that LLMs can be effectively transformed into universal text encoders in a parameter-efficient manner without the need for expensive adaptation or synthetic GPT-4 generated data.

LLM2Vec is a novel unsupervised approach designed to transform any pre-trained decoder-only large language model (LLM) into a robust text encoder capable of generating rich contextualized representations suitable for text embedding tasks. The transformation process comprises three key steps:

1. **Enabling Bidirectional Attention:** This initial step replaces the standard causal attention mechanism of decoder-only LLMs with a bidirectional attention matrix, allowing each token in the input sequence to access information from all other tokens, including those that come after it. This adjustment aims to enhance the model’s ability to incorporate contextual information across the entire sequence.
2. **Masked Next Token Prediction (MNTP):** Following the activation of bidirectional attention, the model undergoes training with a task focused on predicting masked tokens within input sequences. This training helps the model adapt to utilizing the newly available bidirectional context effectively.
3. **Unsupervised Contrastive Learning:** The final step employs an unsupervised contrastive learning approach, which helps the model learn more effective representations of sequences. This is achieved by creating two different representations of the same input sentence and performing a contrastive learning task to optimize the similarity of these representations.

Through these steps, LLM2Vec effectively enhances the capabilities of decoder-only LLMs, enabling them to produce high-quality, universal text embeddings without the need for labeled data or extensive adaptation. The methodology has been demonstrated to outperform traditional encoder-only models on various NLP tasks, particularly on benchmarks like the Massive Text Embeddings Benchmark (MTEB).

1.3 HOW MLM WORKS

1.3.1 INPUT PREPARATION

- Given a sentence or a sequence of tokens, a certain percentage of tokens (typically 15%) are randomly selected to be masked.

- For each selected token:
 - 80% of the time, it is replaced with the special token [MASK].
 - 10% of the time, it is replaced with a random word from the vocabulary.
 - 10% of the time, it remains unchanged.
- This mixture of masking, random replacement, and unchanged tokens helps the model to not only predict masked words but also to understand the context when the original word is present or when an incorrect word is given.

1.3.2 MODEL TRAINING

- The sequence with masked tokens is fed into the model, which is trained to predict the original tokens that were masked out.
- The model learns to leverage the context provided by the surrounding tokens to make these predictions.

1.3.3 LOSS CALCULATION

- The model calculates the loss only on the positions where the tokens were masked, comparing its predicted tokens against the actual tokens.
- This loss is then used to update the model's weights during training.

1.3.4 WHY MLM IS IMPORTANT

- **Contextual Understanding:** By predicting missing words in a sentence, the model learns to understand the context of words and their relationships within a sentence. This helps the model to generate rich, contextualized embeddings.
- **Bidirectional Training:** Unlike traditional language models that predict the next word in a sequence (and thus are unidirectional), MLM allows the model to consider context from both directions (left and right of the masked token). This bidirectional training is key to BERT's success in various NLP tasks.
- **Generalization:** The random masking process helps the model generalize better by preventing it from relying too heavily on any specific tokens during training.

1.3.5 VARIATIONS AND APPLICATIONS

- **Dynamic Masking:** In models like RoBERTa, dynamic masking is used, where the masking pattern changes in each training iteration rather than being fixed. This ensures that the model sees a variety of masked tokens during training, which can improve generalization.
- **Fine-Tuning:** After pre-training with MLM, the model can be fine-tuned on specific downstream tasks (like classification, question answering, etc.) where it applies its learned contextual knowledge to perform these tasks effectively.
- **Extensions:** MLM has been extended or modified in various ways for different models or tasks, but the core idea remains to help the model understand and predict words based on their context within a sequence.

1.3.6 SUMMARY

Masked Language Modeling is a pre-training task where certain tokens in a sentence are masked, and the model learns to predict these tokens using the surrounding context. This task helps the model to develop a deep understanding of the language, enabling it to perform well on a wide range of NLP tasks after fine-tuning.

The text encoder is a BERT model knowledge distilled from a bidirectional version of an open source LLM, further pre-trained for Masked Language Model (MLM) on a large OCR corpus.

1.4 JOINT ARCHITECTURE

The image encoder and the text encoder are combined using a cross-attention mechanism. The combined model is pre-trained on a large corpus of documents and OCRs for Masked Language Model (MLM), Masked Image Model (MIM) and Document Joint Embedding Predictive Architecture (D-JEPA) objectives. The D-JEPA objective here masks part of the document, and is predicted based on the text encoder, while a masked text is predicted based on the image encoder.

Then we have a prompt encoder, which is the CLIP model. The prompt encoder is used to encode the question, and the question is used to prompt the model to segment the document.

The mask decoder efficiently maps the image embedding, prompt embeddings, and an output token to a mask. The mask decoder is a transformer decoder block followed by a dynamic mask prediction head. The modified decoder block uses prompt self-attention and cross-attention in two directions (prompt-to-image embedding and vice-versa) to update all embeddings. After running two blocks, we upsample the image embedding and a MLP maps the output token to a dynamic linear classifier, which then computes the mask foreground probability at each image location. With one output, the model will average multiple valid masks if given an ambiguous prompt.

The model is then fine-tuned on the promptable visual segmentation task.

It also includes a data engine which improves model and data via user interaction.

1.5 STYLE

Papers to be submitted to ICLR 2024 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.

1.6 RETRIEVAL OF STYLE FILES

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

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

The file `iclr2024_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 `iclr2024_conference.sty` and `iclr2024_conference.bst` (to be used with \LaTeX 2e). The file `iclr2024_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 9 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 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.

4.2 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

4.3 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.

¹Sample of the first footnote

²Sample of the second footnote

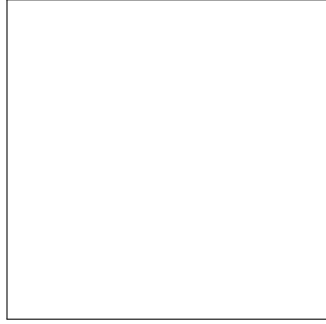


Figure 1: Sample figure caption.

Table 1: Sample table title

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

4.4 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.

5 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}
\mathbf{a}	A scalar random variable
\mathbf{a}	A vector-valued random variable
\mathbf{A}	A matrix-valued random variable

Sets and Graphs

\mathbb{A}	A set
\mathbb{R}	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	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
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of \mathbb{A} that are not in \mathbb{B}
\mathcal{G}	A graph
$Pa_{\mathcal{G}}(\mathbf{x}_i)$	The parents of \mathbf{x}_i in \mathcal{G}

Indexing

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

Calculus

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

Probability and Information Theory

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

Functions

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

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.

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.
- Chaitanya Ryali, Yuan-Ting Hu, Daniel Bolya, Chen Wei, Haoqi Fan, Po-Yao Huang, Vaibhav Aggarwal, Arkabandhu Chowdhury, Omid Poursaeed, Judy Hoffman, Jitendra Malik, Yanghao Li, and Christoph Feichtenhofer. Hiera: A hierarchical vision transformer without the bells-and-whistles. *ICML Oral*, 2023. URL <https://arxiv.org/abs/2306.00989>.

A APPENDIX

You may include other additional sections here.