

# HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolution

Eric Nguyen<sup>\*,1</sup>, Michael Poli<sup>\*,1</sup>, Marjan Faizi<sup>2,\*</sup>,  
Armin W. Thomas<sup>1</sup>, Callum Birch Sykes<sup>3</sup>, Michael Wornow<sup>1</sup>, Aman Patel<sup>1</sup>,  
Clayton Rabideau<sup>3</sup>, Stefano Massaroli<sup>4</sup>, Yoshua Bengio<sup>4</sup>, Stefano Ermon<sup>1</sup>,  
Stephen A. Baccus<sup>1,†</sup>, Christopher Ré<sup>1,†</sup>

November 15, 2023

Code available at <https://github.com/HazyResearch/hyena-dna>

presented by: Felix Dyck - felix.dyck@uni-bielefeld.de

05.12.24

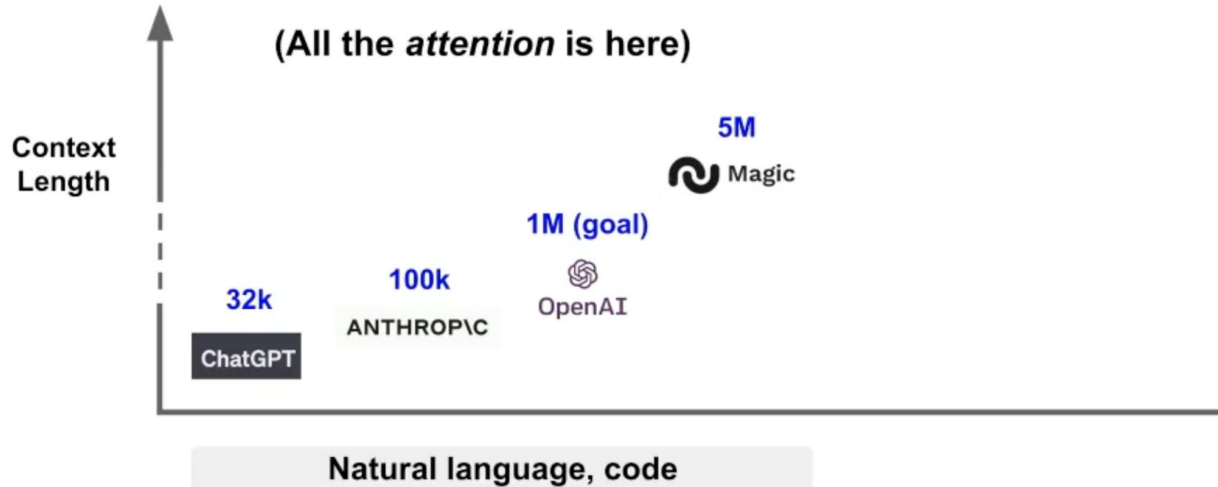
# Content

- Motivation
- Background
- HyenaDNA Block
- Experiments
- Sequence length warm-up
- Task Adaptation

# Long Context

Human Genome has 3.2B nucleotides

Lots of focus on long context in natural language / code



# Background

Nucleotides in a DNA sequence ACGTACGTCGTACGTC...

Previous work:

- typically context length of 512 - 4k tokens
- resolution: usually tokens are not on “nucleotides-level” -> “K-mers” (3-5 nucleotides)

Problems:

- long context is computationally expensive
- “SNPs” (Single-Nucleotide Polymorphisms): variation at a single nucleotide in the DNA sequence

# Perplexity vs Context

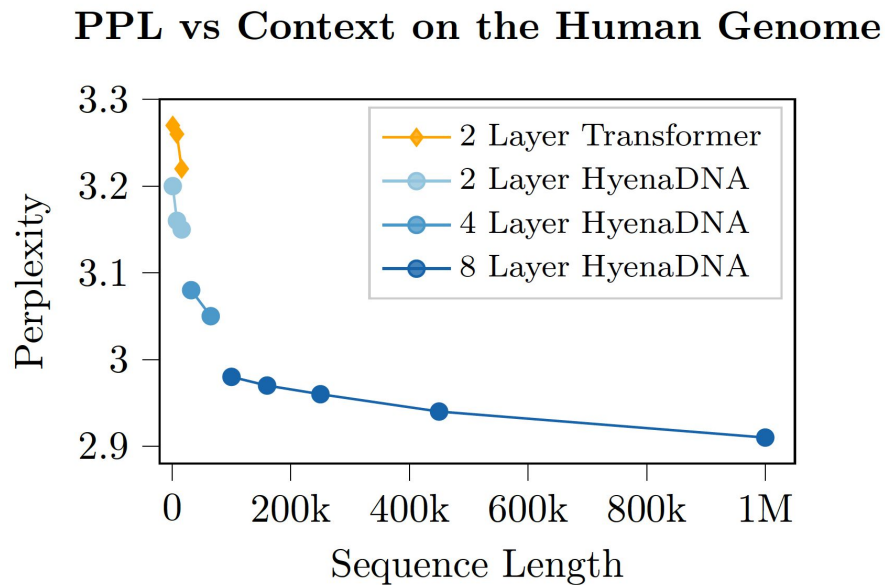


Figure 1.2: Pretraining on the human reference genome using longer sequences leads to better perplexity (improved prediction of next token).

# HyenaDNA Block Architecture

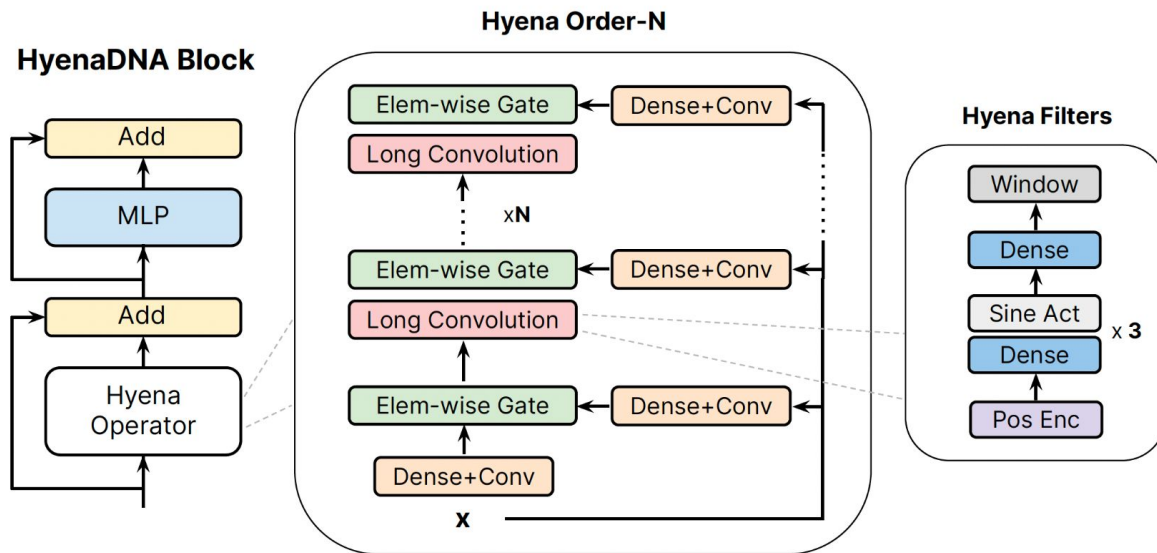


Figure 1.3: HyenaDNA block architecture. A Hyena operator is composed of long convolutions and element-wise gate layers. The gates are fed projections of the input using dense layers and short convolutions. The long convolutions are parameterized *implicitly* via an MLP that produces the convolutional filters. The convolution itself is evaluated using a Fast Fourier Transform convolution with time complexity  $\mathcal{O}(L \log_2 L)$ .

# Runtime

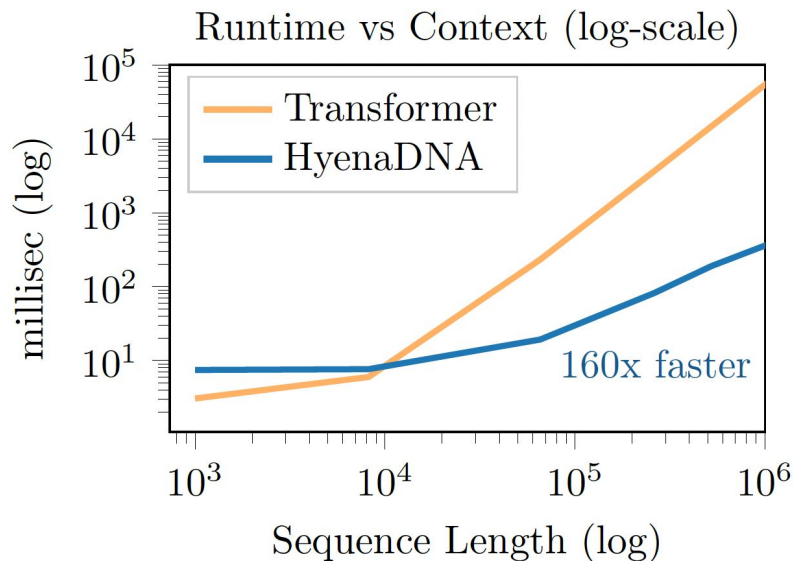


Figure 4.1: Runtime (forward & backward pass) for Transformer and HyenaDNA: 2 layers, width=128, gradient checkpointing, batch size=1, A100 80GB. At 1M tokens HyenaDNA is **160x faster** than Transformer.

# Single Nucleotide Resolution

Table 4.1: **GenomicBenchmarks** Top-1 accuracy (%) for pretrained HyenaDNA, DNABERT and Transformer (GPT from 4.1), and the previous SotA baseline CNN (scratch).

| DATASET                 | CNN  | DNABERT     | GPT  | HYENADNA    |
|-------------------------|------|-------------|------|-------------|
| Mouse Enhancers         | 69.0 | 66.9        | 80.1 | <b>85.1</b> |
| Coding vs Intergenic    | 87.6 | <b>92.5</b> | 88.8 | 91.3        |
| Human vs Worm           | 93.0 | 96.5        | 95.6 | <b>96.6</b> |
| Human Enhancers Cohn    | 69.5 | 74.0        | 70.5 | <b>74.2</b> |
| Human Enhancers Ensembl | 68.9 | 85.7        | 83.5 | <b>89.2</b> |
| Human Regulatory        | 93.3 | 88.1        | 91.5 | <b>93.8</b> |
| Human Nontata Promoters | 84.6 | 85.6        | 87.7 | <b>96.6</b> |
| Human OCR Ensembl       | 68.0 | 75.1        | 73.0 | <b>80.9</b> |



# Ultralong-Range Genomics

- single-nucleotide polymorphisms (SNPs)
- quantifying the functional effects of non-coding variants

Table 4.3: **Chromatin profile prediction** Median AUROC computed over three categories: Transcription factor binding profiles (TF), DNase I-hypersensitive sites (DHS) and histone marks (HM).

| MODEL    | PARAMS | LEN | AUROC       |             |             |
|----------|--------|-----|-------------|-------------|-------------|
|          |        |     | TF          | DHS         | HM          |
| DeepSEA  | 40 M   | 1k  | 95.8        | 92.3        | 85.6        |
| BigBird  | 110 M  | 8k  | 96.1        | 92.1        | 88.7        |
| HyenaDNA | 7 M    | 1k  | <b>96.4</b> | <b>93.0</b> | 86.3        |
|          | 3.5 M  | 8k  | 95.5        | 91.7        | <b>89.3</b> |

# Species classification

- DNA sequences from 5 different species
- struggle on shorter sequences of length 1024

Table 4.5: **Species classification** Top-1 accuracy (%) for 5-way classification (human, lemur, mouse, pig, hippo). The **X** symbol indicates infeasible training time.

| MODEL       | LEN  | ACC         |
|-------------|------|-------------|
| Transformer | 1k   | 55.4        |
| HyenaDNA    | 1k   | 61.1        |
| Transformer | 32k  | 88.9        |
| HyenaDNA    | 32k  | 93.4        |
| Transformer | 250k | <b>X</b>    |
| HyenaDNA    | 250k | 97.9        |
| Transformer | 450k | <b>X</b>    |
| HyenaDNA    | 450k | 99.4        |
| Transformer | 1M   | <b>X</b>    |
| HyenaDNA    | 1M   | <b>99.5</b> |

# Sequence length warm-up

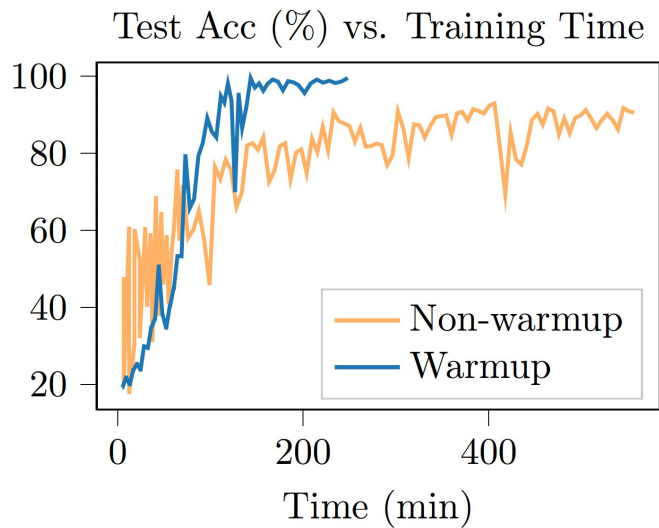


Figure 3.2: Sequence length warm-up reduces the training time of HyenaDNA at sequence length 450k by 40% and boosts accuracy by 7.5 points on species classification.

# Task Adaptation with Soft Prompt Tokens

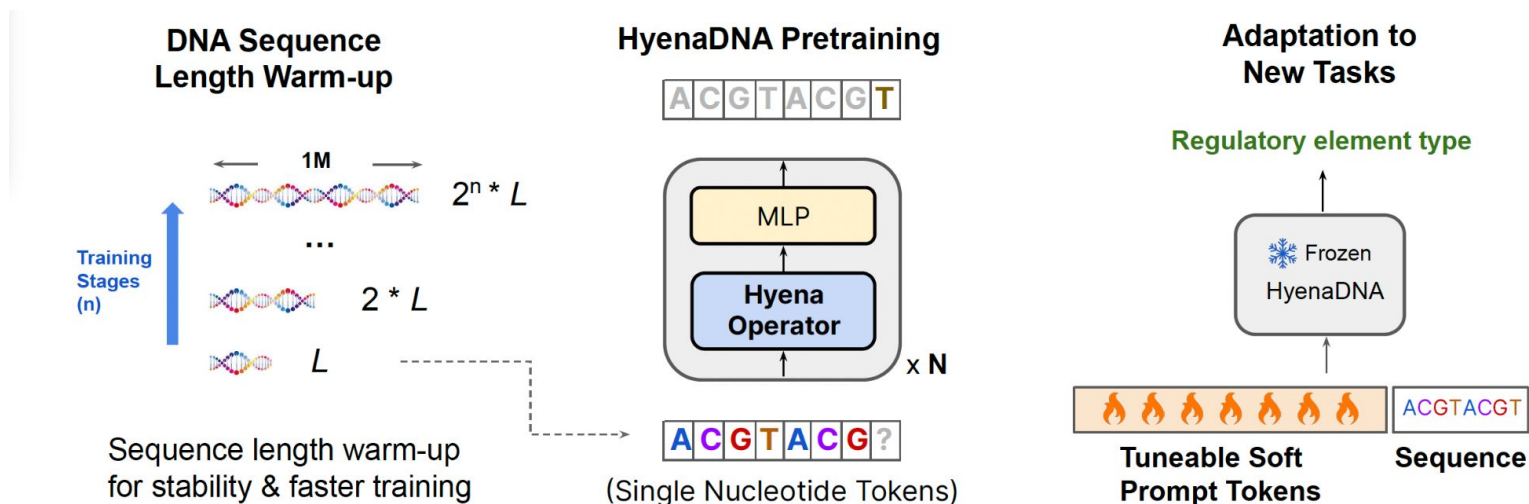


Figure 1.1: HyenaDNA recipe for long-range foundation models in genomics. The HyenaDNA architecture is a simple stack of Hyena operators (Poli et al., 2023) trained using next token prediction. (See Fig. 1.3 for block diagram of architecture). We introduce a new sequence length scheduling technique to stabilize training, and provide a method to leverage the longer context length to adapt to novel tasks without standard fine-tuning by filling the context window with learnable soft prompt tokens.

# Fine Tuning HyenaDNA

## Tasks:

- Identify Enhancer Regions in a DNA Sequence
- Species classification

## Receipt:

- Pretrained Model: HyenaDNA
- Input DNA Sequence
- Soft Prompt Tokens
- Fine-Tuning:
  - Soft Prompt tokens are optimized for the task
- Output: Probability scores for:
  - each nucleotide being part of an enhancer
  - dna being each species

## Example:

Input: ACGTACGTCGTACGTC...

-> [0.1, 0.2, 0.3, 0.4, ...]

- Soft Prompt: [P1, P2, P3]

-> P1 = [0.01, -0.05, 0.07, ...], P2 = [0.1, 0.03, -0.02, ...], P3 = [0.05, 0.07, -0.01, ...]

- Combined Input: [P1, P2, P3, 0.1, 0.2, 0.3, 0.4, ...]

- Predictions: [0.01, 0.85, 0.10, 0.90, ...]  
(Species A, Species B, ...)

# Summary

- Sequence length of up to 1 million turned out to be beneficial.
- Observing single nucleotides as individual tokens is important.
- Using long convolutions is especially more efficient the longer the sequence gets
- “Warm-up” learns shorter sequences first and then gradually increases the length leading to better performance faster and in the end.

Thank you for Listening!

# Hyperparameters used for Training

Table A.3: GenomicBenchmarks hyperparameters for HyenaDNA and the baseline Transformer (GPT from 4.1), which uses FlashAttention (Dao et al., 2022a).

|                             | TRANSFORMER                     | HyenaDNA    |
|-----------------------------|---------------------------------|-------------|
| Layers                      | 2                               | 2           |
| Width                       | 128                             | 128         |
| Parameters                  | 529k                            | 436k        |
| Learning rate               | $1.6e^{-4}$                     | $1.6e^{-4}$ |
| Weight decay (model)        | 0-0.2                           | 0-0.2       |
| Weight decay (Hyena layers) | -                               | 0           |
| Embed dropout               | 0-0.2                           | 0.0-0.3     |
| Resid dropout               | 0-0.2                           | 0-0.3       |
| Num heads                   | 8                               | -           |
| Optimizer                   | AdamW                           |             |
| Optimizer momentum          | $\beta_1, \beta_2 = 0.9, 0.999$ |             |
| LR scheduler                | Cosine decay                    |             |
| Batch size                  | 128-1024                        |             |
| Training epoch              | 100                             |             |
| Reverse complement aug.     | true/false                      |             |
| Sequence lengths            | 200-4800                        |             |



Table 1. Genomic Benchmarks. Top-1 accuracy ( $\uparrow$ ) across 5-fold cross-validation (CV) for pretrained HyenaDNA, Mamba NTP, Caduceus models, and a supervised CNN baseline (trained from scratch). Best values per task are **bolded**, second best are *italicized*. Error bars indicate the difference between the maximum and minimum values across 5 random seeds used for CV.

|                         | CNN<br>(264k)     | HYENADNA<br>(436k)       | MAMBA<br>(468k)   | CADUCEUS<br>w/o EQUIV.<br>(470k) | CADUCEUS-PH<br>(470k)    | CADUCEUS-PS<br>(470k)    |
|-------------------------|-------------------|--------------------------|-------------------|----------------------------------|--------------------------|--------------------------|
| MOUSE ENHANCERS         | 0.715 $\pm$ 0.087 | <i>0.780</i> $\pm$ 0.025 | 0.743 $\pm$ 0.054 | 0.770 $\pm$ 0.058                | 0.754 $\pm$ 0.074        | <b>0.793</b> $\pm$ 0.058 |
| CODING VS. INTERGENOMIC | 0.892 $\pm$ 0.008 | 0.904 $\pm$ 0.005        | 0.904 $\pm$ 0.004 | 0.908 $\pm$ 0.003                | <b>0.915</b> $\pm$ 0.003 | <i>0.910</i> $\pm$ 0.003 |
| HUMAN VS. WORM          | 0.942 $\pm$ 0.002 | 0.964 $\pm$ 0.002        | 0.967 $\pm$ 0.002 | <i>0.970</i> $\pm$ 0.003         | <b>0.973</b> $\pm$ 0.001 | 0.968 $\pm$ 0.002        |
| HUMAN ENHANCERS COHN    | 0.702 $\pm$ 0.021 | 0.729 $\pm$ 0.014        | 0.732 $\pm$ 0.029 | 0.741 $\pm$ 0.008                | <b>0.747</b> $\pm$ 0.004 | <i>0.745</i> $\pm$ 0.007 |
| HUMAN ENHANCER ENSEMBL  | 0.744 $\pm$ 0.122 | 0.849 $\pm$ 0.006        | 0.862 $\pm$ 0.008 | 0.883 $\pm$ 0.002                | <i>0.893</i> $\pm$ 0.008 | <b>0.900</b> $\pm$ 0.006 |
| HUMAN REGULATORY        | 0.872 $\pm$ 0.005 | 0.869 $\pm$ 0.012        | 0.814 $\pm$ 0.211 | 0.871 $\pm$ 0.007                | <i>0.872</i> $\pm$ 0.011 | <b>0.873</b> $\pm$ 0.007 |
| HUMAN OCR ENSEMBL       | 0.698 $\pm$ 0.013 | 0.783 $\pm$ 0.007        | 0.815 $\pm$ 0.002 | 0.818 $\pm$ 0.003                | <b>0.828</b> $\pm$ 0.006 | <i>0.818</i> $\pm$ 0.006 |
| HUMAN NONTATA PROMOTERS | 0.861 $\pm$ 0.009 | 0.944 $\pm$ 0.002        | 0.933 $\pm$ 0.007 | 0.933 $\pm$ 0.006                | <b>0.946</b> $\pm$ 0.007 | <i>0.945</i> $\pm$ 0.010 |

Table A.4: **GenomicBenchmarks Top-1 accuracy (%)** GPT is the causal Transformer from 4.1, HyenaDNA k-mer uses a 6-mer tokenizer, and HyenaDNA bidirection is a bidirectional version of the Hyena operator.

| MODEL                   | GPT  | GPT         | HyenaDNA | HyenaDNA    | HyenaDNA<br>k-mer | HyenaDNA<br>bidirection | DNABERT     |
|-------------------------|------|-------------|----------|-------------|-------------------|-------------------------|-------------|
| Pretrained              | no   | yes         | no       | yes         | no                | no                      | yes         |
| Mouse Enhancers         | 79.3 | 79.3        | 84.7     | <b>85.1</b> | 81.8              | 80.6                    | 66.9        |
| Coding vs Intergenic    | 89.3 | 91.2        | 90.9     | 91.3        | 86.7              | 90.3                    | <b>92.5</b> |
| Human vs Worm           | 94.8 | <b>96.6</b> | 96.4     | <b>96.6</b> | 92.9              | 95.9                    | 96.5        |
| Human Enhancers Cohn    | 67.7 | 72.9        | 72.9     | <b>74.2</b> | 69.8              | 72.1                    | 74.0        |
| Human Enhancers Ensembl | 79.0 | 88.3        | 85.7     | <b>89.2</b> | 88.0              | 85.9                    | 85.7        |
| Human Regulatory        | 90.2 | 91.8        | 90.4     | <b>93.8</b> | 90.2              | 89.1                    | 88.1        |
| Human Nontata Promoters | 85.2 | 90.1        | 93.3     | <b>96.6</b> | 83.5              | 88.5                    | 85.6        |
| Human OCR Ensembl       | 68.3 | 79.9        | 78.8     | <b>80.9</b> | 70.2              | 75.3                    | 75.1        |

# Some more References...

Video-presentation of the paper:

<https://youtu.be/haSkAC1fPX0?si=v28n75mdSjTURvay>

Caduceus: Bi-Directional Equivariant Long-Range DNA Sequence Modeling:

<https://arxiv.org/pdf/2403.03234>

Articles:

- <https://hazyresearch.stanford.edu/blog/2023-06-29-hyena-dna>
- <https://aibusiness.com/ml/hyena-dna-a-large-language-model-trained-on-human-genome-sequences>