

HybriDNA: A Hybrid Transformer-Mamba2 Long-Range DNA Language Model

Mingqian Ma^{1,2†}, Guoqing Liu^{1†}, Chuan Cao^{1†}, Pan Deng^{1†},
Tri Dao³, Albert Gu⁴, Peiran Jin¹, Zhao Yang⁵, Yingce Xia¹,
Renqian Luo¹, Pipi Hu¹, Zun Wang¹, Yuan-Jyue Chen¹,
Haiguang Liu¹, Tao Qin^{1*}

¹Microsoft Research AI for Science.

²UM-SJTU Joint Institute, Shanghai Jiao Tong University.

³Department of Computer Science, Princeton University.

⁴Machine Learning Department, Carnegie Mellon University.

⁵Gaoling School of Artificial Intelligence, Renmin University of China.

*Corresponding author(s). E-mail(s): taoqin@microsoft.com;

Contributing authors: mingqianma@sjtu.edu.cn;

guoqingliu@microsoft.com; chuancao.926@gmail.com;

pan.deng@microsoft.com; tri@tridao.me; agu@andrew.cmu.edu;

peiranjin@microsoft.com; yangyz1230@gmail.com;

yingce.xia@microsoft.com; renqianluo@microsoft.com;

pisquare@microsoft.com; zunwang@microsoft.com;

yuanjc@microsoft.com; haiguang.liu@microsoft.com;

[†]These authors contributed equally to this work.

Advances in natural language processing and large language models have sparked growing interest in modeling DNA, often referred to as the “language of life”. However, DNA modeling poses unique challenges. First, it requires the ability to process ultra-long DNA sequences while preserving single-nucleotide resolution, as individual nucleotides play a critical role in DNA function. Second, success in this domain requires excelling at both generative and understanding tasks: generative tasks hold potential for therapeutic and industrial applications, while understanding tasks provide crucial insights into biological mechanisms and diseases. To address these challenges, we propose **HybriDNA**, a decoder-only DNA language model that incorporates a hybrid

Transformer-Mamba2 architecture, seamlessly integrating the strengths of attention mechanisms with selective state-space models. This hybrid design enables HybriDNA to efficiently process DNA sequences up to 131kb in length with single-nucleotide resolution. HybriDNA achieves state-of-the-art performance across 33 DNA understanding datasets curated from the BEND, GUE, and LRB benchmarks, and demonstrates exceptional capability in generating synthetic cis-regulatory elements (CREs) with desired properties. Furthermore, we show that HybriDNA adheres to expected scaling laws, with performance improving consistently as the model scales from 300M to 3B and 7B parameters. These findings underscore HybriDNA’s versatility and its potential to advance DNA research and applications, paving the way for innovations in understanding and engineering the “language of life”.

1 Introduction

Deoxyribonucleic acid (DNA) serves as the genetic code of life, encoding the instructions that govern gene expression, cellular processes, and biological functions. A deep understanding of the “language” of DNA is crucial for unraveling the molecular mechanisms that underlie biological functions and for leveraging these insights to advance medicine and biotechnology. The advent of high-throughput sequencing technologies has generated an immense volume of genomic data, creating an unprecedented opportunity for machine learning models to uncover complex patterns and relationships within DNA sequences. Foundation models, pretrained on large-scale unlabeled datasets, have already demonstrated remarkable capabilities in natural languages [1–3] and protein languages [4–6].

Recently, foundation models have begun to drive a paradigm shift in genomics, showcasing their ability to learn rich representations of DNA sequences that can be fine-tuned for a diverse array of downstream tasks. Currently, DNA foundation models primarily adopt two main architectural approaches. The first approach, inspired by BERT [1], employs encoder-only Transformer architectures. Models such as DNABERT2 [7] and Nucleotide Transformer (NT) [8] excel at capturing contextual information within DNA sequences, producing high-quality embeddings suitable for tasks such as classification and regression. However, their bidirectional nature constrains their ability to design novel DNA sequences. The second approach leverages decoder-only architectures, such as Hyena [9] and the Transformer architecture in GPT [10], which are autoregressive and well-suited for generative tasks. Models like HyenaDNA [11] and Evo [12] have shown promising results in generating DNA sequences. Nevertheless, they often fall behind encoder-only models in understanding tasks requiring a deep understanding of sequence context.

This dichotomy highlights two critical challenges in DNA modeling: (1) How to develop a DNA foundation model that integrates robust contextual understanding with advanced design capabilities? Such a model would not only enhance the analysis of existing genomic data but also enable the design of novel, functional DNA sequences. (2) How to efficiently address the intricate complexity of DNA sequences, which involves long-range interactions critical to fundamental biological processes? Recent advances in Selective State Space Models (SSMs), such as Mamba [13, 14],

have shown remarkable potential for addressing information-dense tasks, including language modeling [15, 16]. These models efficiently handle long-range dependencies with subquadratic complexity, offering a promising approach to the challenges posed by DNA sequence modeling. However, SSMs alone struggle to capture fine-grained, single-nucleotide-level interactions vital for understanding DNA function.

In this work, we introduce HybriDNA, a novel class of decoder-only DNA language models that leverage a hybrid Transformer-Mamba2 architecture. This hybrid design combines the complementary strengths of its components: Mamba2 blocks excel at efficiently processing long sequences and capturing long-range dependencies, whereas Transformer blocks enhance the model’s ability to focus on fine-grained, token-level details within the context of the entire sequence. Pretrained on large-scale, multi-species genomes at single-nucleotide resolution with a next-token prediction objective, HybriDNA demonstrates foundational capabilities in both understanding and designing genomic sequences. By incorporating an *echo embedding* discriminative fine-tuning approach, HybriDNA achieves state-of-the-art performance across 35 biologically significant DNA understanding datasets, such as transcription factor binding prediction and promoter detection [7]. Additionally, through generative fine-tuning, HybriDNA exhibits exceptional proficiency in designing synthetic cis-regulatory elements (CREs) with desirable functional properties, such as yeast promoters and cell type-specific human enhancers [17]. Finally, we show that scaling up HybriDNA is beneficial: increasing model size from 300 million to 3 billion and 7 billion parameters improves performance, adhering to scaling laws observed in language models such as GPT [10, 18]. Extending the context length (e.g., from 8 kilobases to 131 kilobases at single-nucleotide resolution) further enhances HybriDNA’s performance on specific tasks. Together, these advancements position HybriDNA as a powerful tool for advancing both the understanding and engineering of genomic sequences.

Our contributions are summarized as follows:

- We propose HybriDNA, a class of decoder-only DNA language models that integrates a hybrid Transformer-Mamba2 architecture. Leveraging this architecture, we develop a comprehensive training pipeline that includes pretraining, echo embedding-based discriminative fine-tuning, and generative fine-tuning.
- HybriDNA achieves state-of-the-art performance across 33 diverse and biologically meaningful DNA understanding datasets, spanning human and multi-species genomes. HybriDNA outperforms many existing encoder-only models.
- HybriDNA demonstrates its generative capabilities, showcasing its proficiency in designing desirable synthetic cis-regulatory elements.
- HybriDNA demonstrates the ability to process long-context DNA sequences on tasks involving the analysis of extended DNA sequences (e.g., 131 kilobases).
- Scaling laws for HybriDNA are observed: increasing the model size from 300M to 3B and 7B parameters improves performance.

2 Preliminaries

2.1 Attention Mechanism in Transformers

Powering many foundation models is the attention mechanism [19, 20] in Transformers. Attention is a type of operator that assigns scores to every pair of tokens in a sequence, enabling each element to “attend” to the others. The most widely adopted variant of attention to date is Scaled Dot-Product Attention, which is defined as:

$$y = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right) \cdot V. \quad (1)$$

Let $x \in \mathbb{R}^{L \times d}$ represents an input sequence with sequence length L and embedding size d , the learnable parameters $W_K \in \mathbb{R}^{d \times d_k}$, $W_Q \in \mathbb{R}^{d \times d_k}$, and $W_V \in \mathbb{R}^{d \times d}$ are used to compute the key, query, and value matrices: $K = xW_K$, $Q = xW_Q$, and $V = xW_V$. The attention layer, therefore, transforms an input x of shape $\mathbb{R}^{L \times d}$ into an output y of the same shape, $\mathbb{R}^{L \times d}$.

Attention computes all pairwise comparisons for every token in a sequence, resulting in a computational complexity that scales as $O(L^2)$ with sequence length L . While this enables capturing global context at high resolution, it also restricts the context length on modern GPU architectures.

2.2 Selective State Space Models

Structured state space sequence models (S4) [21, 22] are a recent class of sequence models in deep learning that are broadly related to RNNs, CNNs, and classical state space models. They are inspired by a specific continuous system, which maps a 1-dimensional sequence $x \in \mathbb{R}^L \mapsto y \in \mathbb{R}^L$ via a hidden state $h \in \mathbb{R}^{(L,N)}$, where L represents the sequence length, and N represents the SSM state size.

Specifically, the continuous system is defined by three matrices $(\mathbf{A}, \mathbf{B}, \mathbf{C})$, $\mathbf{A} \in \mathbb{R}^{N \times N}$, $\mathbf{B} \in \mathbb{R}^{N \times 1}$, and $\mathbf{C} \in \mathbb{R}^{1 \times N}$. They define a sequence-to-sequence transformation in two steps, as detailed in Eqn. 2 (first column).

$$\begin{aligned} h'(t) &= \mathbf{A}h(t) + \mathbf{B}x(t), & h_t &= \bar{\mathbf{A}}h_{t-1} + \bar{\mathbf{B}}x_t, & \bar{\mathbf{K}} &= (\mathbf{C}\bar{\mathbf{B}}, \mathbf{C}\bar{\mathbf{A}}\bar{\mathbf{B}}, \dots, \mathbf{C}\bar{\mathbf{A}}^k\bar{\mathbf{B}}), \\ y(t) &= \mathbf{C}h(t). & y_t &= \mathbf{C}h_t. & y &= x * \bar{\mathbf{K}}. \end{aligned} \quad (2)$$

Discretization S4 models are discrete versions of the continuous system (as shown in the second column of Eqn. 2), which incorporates a timescale parameter $\Delta \in \mathbb{R}$ to transform the continuous parameters \mathbf{A} , and \mathbf{B} into their discrete counterparts, $\bar{\mathbf{A}}$ and $\bar{\mathbf{B}}$. The zero-order hold (ZOH) is a commonly used method for this transformation, defined as follows, where “exp” represents the exponential function.

$$\bar{\mathbf{A}} = \exp(\Delta\mathbf{A}), \quad \bar{\mathbf{B}} = (\Delta\mathbf{A})^{-1}(\exp(\Delta\mathbf{A}) - \mathbf{I}) \cdot \Delta\mathbf{B}. \quad (3)$$

Computation After the parameters have been transformed into $(\bar{\mathbf{A}}, \bar{\mathbf{B}})$, the model can be computed in two ways, either as a linear recurrence (corresponding to the

second column of Eqn. 2) or a global convolution (corresponding to the third column of Eqn. 2). Commonly, the model uses the convolutional mode for efficient parallelizable training (where the whole input sequence is seen ahead of time), and switched into recurrent mode for efficient autoregressive inference (where the inputs are seen one timestep at a time).

Structured matrix A S4 models are named as such because efficiently computing them requires imposing a specific structure on the \mathbf{A} matrix. The most popular structure is diagonal [23–25]. In this scenario, the $\mathbf{A}, \mathbf{B}, \mathbf{C}$ matrices can each be represented by N numbers. To process an input sequence x of D channels, the SSM is applied independently to each channel.

Linear Time Invariance (LTI) A key characteristic of Eqn. 2 is its linear time invariance, meaning the model’s dynamics remain constant over time. In other words, the matrices $(\Delta, \mathbf{A}, \mathbf{B}, \mathbf{C})$, and consequently $(\bar{\mathbf{A}}, \bar{\mathbf{B}})$, are fixed for all time-steps. The LTI property is closely associated with recurrence and convolutions. Before Mamba, all S4 models adhered to LTI, often computed as convolutions during training.

H3 [26] generalizes the recurrence to use S4; it can be viewed as an architecture with an SSM sandwiched by two gated connections. H3 also inserts a standard local convolution, which they frame as a shift-SSM, before the inner SSM layer.

Mamba1 [13] introduces the concept of Selective SSMs (S6), enhancing the traditional S4 framework through input-dependent gating mechanisms. The matrices $\bar{\mathbf{A}}, \bar{\mathbf{B}}$, and Δ are dynamically gated by the input x_t , enabling them to adjust their behavior based on the current input. Mamba1 simplifies the block design by combining the H3 block [27] with gated MLPs. Additionally, Mamba1 proposes selective scan, a hardware-aware algorithm that computes the model recurrently using a scan operation, enhancing computational efficiency and scalability.

Mamba2 [14] enhances Mamba1 by introducing two key enhancements:

1. At the SSM layer, the new Structured State Space Duality (SSD) layer imposes a stricter constraint on the diagonal matrix $\bar{\mathbf{A}}$. The diagonal matrix is now simplified to a scalar times an identity matrix, allowing it to be represented using a single identical value across the diagonal. Consequently, $\bar{\mathbf{A}}$ can be represented with a shape corresponding only to the sequence length.
2. The Mamba2 block produces the SSM parameters $(\bar{\mathbf{A}}, \bar{\mathbf{B}}, \mathbf{C})$ in parallel with the input x , as opposed to sequentially in the Mamba1 block. This change facilitates greater parallelism and scalability, enabling tensor parallelism for scaling the model to larger dimensions and longer contexts. Compared to Mamba1, Mamba2 supports much larger state dimensions (increasing from $N = 16$ in Mamba1 to $N = 64$, $N = 256$, or even higher) while also significantly improving training speed.

3 HybriDNA Foundation Model

In this section, we introduce the HybriDNA model for long-range genomic sequence modeling. We begin with a detailed description of the model architecture, followed by an explanation of the pretraining stage of HybriDNA. Finally, we discuss the fine-tuning stages used for various downstream applications. The model’s pipeline is illustrated in Fig. 1.

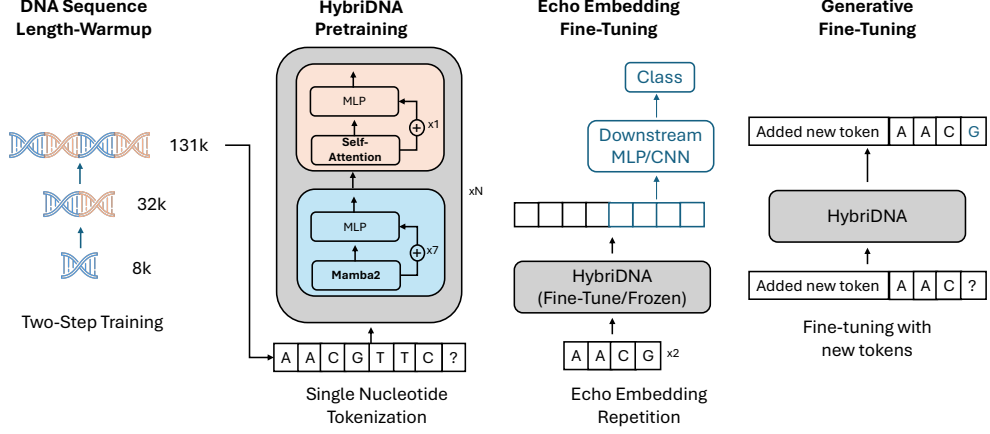


Fig. 1: Overview of HybriDNA: A Language Model for DNA Sequences. HybriDNA builds upon an efficient hybrid Transformer and Mamba2 architecture. It is initially pretrained on large-scale, multi-species genomic data at single-nucleotide resolution using a next-token prediction objective. Subsequently, HybriDNA employs an echo embedding fine-tuning approach for DNA understanding tasks and a generative fine-tuning approach for DNA generation tasks.

3.1 Model Architecture

The HybriDNA model uses a decoder-only, sequence-to-sequence architecture purpose-built for efficiently and accurately processing long-range DNA sequences. It combines the strengths of Mamba2 selective state-space models and Transformer attention mechanisms within a hybrid framework inspired by recent hybrid architectures [16, 28, 29]. As shown in Fig. 1 (second column), the architecture consists of a series of HybriDNA blocks, where each block alternates between HybriDNA Mamba2 blocks and HybriDNA Transformer blocks in a 7:1 ratio. This configuration has been empirically proven to effectively balance the advantages of both block types, achieving optimal performance in the NLP domain [15].

A key component of the HybriDNA Mamba2 block is the **State-Space Duality (SSD)** layer [14]. It processes input sequence x using the recurrence:

$$h_t = A_t h_{t-1} + B_t x_t, \quad y_t = C_t^\top h_t, \quad (4)$$

where $h_t \in \mathbb{R}^N$ denotes the hidden state, $A_t \in \mathbb{R}^{N \times N}$ represents the state transitions, $x_t \in \mathbb{R}$ is the input, $B_t \in \mathbb{R}^{N \times 1}$ projects the input, and $C_t \in \mathbb{R}^{N \times 1}$ maps the hidden state to the output $y_t \in \mathbb{R}$.

The SSD layer simplifies the matrix A_t to $A_t = a_t I$, where $a_t \in \mathbb{R}$ is a scalar and I is the identity matrix, to further improve efficiency. This simplification reduces the recurrence to:

$$h_t = a_t h_{t-1} + B_t x_t. \quad (5)$$

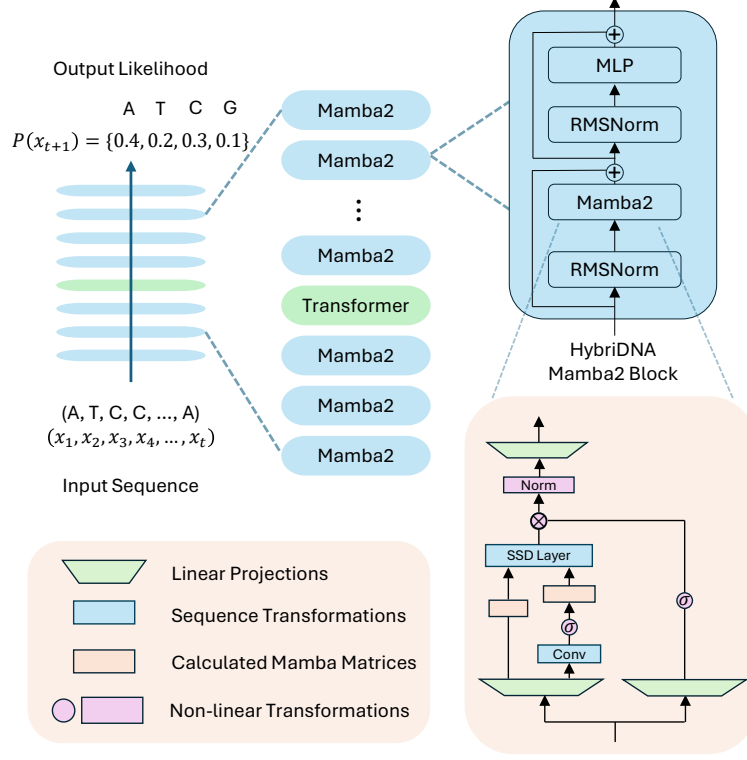


Fig. 2: Model Architecture of HybriDNA

For multi-dimensional (channel) inputs $x \in \mathbb{R}^{L \times d}$, the SSD layer is extended into a multi-head design, where each head independently processes a distinct subset of the input dimensions, similar to the multi-head attention mechanism in Transformers. This architecture allows the SSD layer to capture complex interactions across multiple input channels in parallel, greatly enhancing its representational capacity. Typically, the head dimension is set to 64 or 128, consistent with configurations used in standard transformers.

Computationally, the SSD layer can be reformulated as a matrix operation:

$$y = Mx, \quad M_{ij} = \begin{cases} C_i^\top A_{j:i+1} B_j & \text{if } j \leq i, \\ 0 & \text{otherwise,} \end{cases}$$

where $A_{j:i+1}$ refers to $A_j \cdots A_i$ when $i > j$ and A_i when $i = j$. The matrix M is a specific type of semiseparable matrices, which is a rank-structured matrix where every submatrix contained on and below the diagonal has a rank of at most N , corresponding to the SSM's state dimension. Specifically, a semi-separable matrix M can be expressed

as the sum of two components:

$$M = UV^\top + K,$$

where U and V capture the structured part, and K represents the lower-triangular portion. This structure ensures efficient computation with $O(NL)$ complexity, which is significantly faster than the $O(L^2)$ cost of standard Transformers-based models.

As shown in Fig. 2, the HybriDNA Mamba2 block is a scalable model architecture designed for efficient processing of input sequences. It incorporates grouped-value projections and lightweight convolutional operations before the SSD layer. These projections, along with 1D convolutions, facilitate flexible feature extraction and dimensionality reduction while maintaining computational efficiency. All data-dependent projections are computed in parallel at the start of the block, utilizing tensor parallelism to maximize the use of matrix multiplication units on modern GPUs. Additionally, the Mamba2 blocks employ RMSNorm normalization [30] both before input projection and after the SSD layer, enhancing training stability, particularly at large model scales. For the HybriDNA Transformers block, we use a standard Transformers Decoder block as described in Section 2.1.

In the HybriDNA architecture, the first block is a HybriDNA Mamba2 block, which eliminates the need for explicit positional embeddings or mechanisms such as RoPE [31]. This design choice results in a HybriDNA design that completely omits positional encoding. Furthermore, unlike the Jamba model [16], the MLP layers in HybriDNA do not use Mixture-of-Experts (MoE) configurations due to instability observed during fine-tuning for DNA-related downstream tasks. By leveraging this hybrid architecture, HybriDNA excels in both short- and long-range tasks while enabling the robust generation of synthetic DNA sequences.

3.2 Pretraining on Multi-Species Genomes

Dataset We pretrain HybriDNA on a large-scale, multi-species genome dataset using next nucleotide (token) prediction (NTP). This dataset was curated from the Nucleotide Transformer [8] and NCBI, and was down-sampled to include 845 species, collectively comprising 160 billion nucleotides. Table 1 provides a summary of the contribution of each genome class, represented as the number of nucleotides relative to the total nucleotide count in the dataset. A diverse and comprehensive genome dataset encompassing multiple species is essential for enabling the model to effectively interpret a wide range of genomic sequences. Such a dataset ensures that the model captures patterns and interactions representative of various biological systems, thereby enhancing its ability to generalize across species and tasks.

Tokenizer HybriDNA employs a straightforward and effective base-level tokenization strategy, encoding each nucleotide (A, C, T, G) as an individual token. This strategy ensures that the model processes genomic data with high fidelity to its natural structure, enabling nuanced interpretation and feature extraction. Unlike higher-order tokenization schemes that aggregate multiple bases into a single token, the base-level strategy treats each nucleotide as a fundamental unit, preserving its unique contribution to genomic patterns. This method is particularly advantageous for capturing

Class	# Species (train)	# Nucleotides (train)	# Species (valid)	# Nucleotides (valid)
Bacteria	647	16.5B	20	0.5B
Fungi	44	2.0B	3	0.2B
Invertebrate	37	19.9B	2	1.9B
Protozoa	9	0.45B	1	0.05B
Mammalian Vertebrate	28	65.2B	3	4.6B
Other Vertebrate	51	57.4B	6	6.0B
Total	845	160.75B	35	13.25B

Table 1: Statistics of multi-species pretraining data used in HybriDNA

low-level sequence variations with significant biological implications, such as single nucleotide polymorphisms (SNPs) and point mutations.

Previous transformer-based models like DNABERT2 [7] and Nucleotide Transformer [8] faced challenges when utilizing base-level tokenization due to the resulting longer sequence lengths, which lead to higher computational costs and memory demands. HybriDNA overcomes these limitations through its hybrid architecture. The Mamba2 blocks enable efficient processing of long sequences, allowing HybriDNA to harness the fine-grained detail of base-level tokenization without compromising performance or scalability. This capability is particularly critical for tasks that require modeling extensive genomic contexts, such as enhancer-promoter interactions and chromatin state analysis.

DNA Sequence Length Warm-up To enhance HybriDNA’s ability to generalize effectively across longer genomic ranges, we implement a multi-stage warm-up procedure during the pretraining phase. The pretraining process begins by training the model with an 8,192 token context length, establishing a strong foundation for capturing intermediate sequence dependencies. After that, the context length is gradually increased—first to 32,768 tokens and then to 131,072 tokens—with each extension undergoing additional training equal to 2% of the training steps originally used for the 8,192 token context length. This gradual extension enables the model to adapt to increasingly long-range dependencies and ensure efficient processing of large-scale genomic spans, equipping HybriDNA to excel in tasks that demand long-range comprehension.

3.3 Downstream Fine-tuning

3.3.1 Discriminative Fine-tuning for DNA Understanding tasks

To develop a DNA foundation model capable of handling both generative and understanding downstream tasks, HybriDNA employs a GPT-like decoder-only architecture. However, a key limitation of autoregressive models, compared to bidirectional models, is their inability to incorporate information from future tokens into the embeddings of current tokens. To address this issue, HybriDNA introduces a novel *echo embedding* technique during the fine-tuning stage for understanding tasks, drawing inspiration from the work of [32].

The core idea of this method is that repeating sequences facilitates the encoding of contextual information from subsequent elements into the embeddings. To illustrate, consider an input sequence x and its corresponding label y in a classification task with

K classes. For instance, given the input sequence $x = \text{AACG}$, an “echo” input is created by duplicating x : $x_{\text{echo}} = \text{AACGAACG}$. Hidden embeddings are then extracted from the final hidden layer, with a particular focus on the embeddings from the latter half of the sequence. A mean-pooling operation is applied over these selected token embeddings to produce $h_{\theta}(x_{\text{echo}})$, which is designed to capture contextual information from the repeated segment of the input. This pooled vector, $h_{\theta}(x_{\text{echo}})$ is subsequently passed into a classification head, typically consisting of a linear layer with weights $W \in \mathbb{R}^{d \times K}$ and bias $b \in \mathbb{R}^K$, to generate the predicted probability distribution over the K classes:

$$P(y|x) = \text{softmax}(h_{\theta}(x_{\text{echo}})W + b). \quad (6)$$

To optimize the model, the standard cross-entropy loss is employed to adjust the parameters of either the classification head alone (W and b) or the entire model (θ , W , and b). By incorporating bidirectional context into the autoregressive model, echo embeddings bridge the gap between traditional autoregressive embeddings and the complex demands of high-resolution genomic tasks, offering significant advantages for analyzing large genomic sequences.

A potential limitation of using echo embeddings for discriminative fine-tuning is the increased computational cost, as doubling the input length may raise memory requirements. However, HybriDNA’s efficient hybrid architecture helps mitigate this burden, making the technique practical and scalable for a wide range of genomic analyses and applications.

3.3.2 Generative Fine-tuning for DNA Generation Tasks

Autoregressive natural language models, such as ChatGPT, are capable of generating highly realistic, human-like text while adhering to human instructions to produce satisfactory responses [18]. In a similar vein, HybriDNA is an autoregressive model pretrained on multi-species genomic data at single-nucleotide resolution, unlocking the potential to design novel and realistic DNA sequences for a broad range of real-world applications.

To achieve this, we introduce a set of prompt tokens specifically designed to encode task-specific instructions. These prompt tokens are incorporated into the existing nucleotide vocabulary and initialized randomly within the embedding layer, which is expanded to include additional rows corresponding to the newly introduced token IDs. HybriDNA then predicts each nucleotide token x_t autoregressively, conditioned on all preceding prompt tokens that specify the task-specific requirements.

Specifically, we optimize the *next token prediction* loss:

$$\mathcal{L}_{\text{generative}}(\theta) = -\frac{1}{T} \sum_{t=1}^T \log \left(p_{\theta}(x_t \mid z_0, \dots, z_{k-1}, x_0, \dots, x_{t-1}) \right), \quad (7)$$

where θ represents the complete set of model parameters, z_{k-1} denotes the k -th prompt tokens, and x_{t-1} represents the t -th generated nucleotide token. By minimizing Eqn. 7, HybriDNA is trained to interpret the specialized tokens and generate realistic genomic sequences that align with specific design objectives.

4 Experiments

In the experiment section, we aim to answer the following questions regarding the HybriDNA models and their key capabilities: (1) How do the pretraining losses of HybriDNA models compare across different scales and configurations? (2) Can the models achieve state-of-the-art performance on short-range understanding benchmarks, and how does scaling affect their effectiveness? (3) How do the models perform on long-range understanding tasks, and do they exhibit improved results with increased pretraining context length? (4) Can the models generate realistic and desirable regulatory sequences across multiple species? (5) How do the models compare to pure transformer-based models in terms of computational efficiency during training?

We benchmark our model against a series of recently proposed tasks, including DNABERT2 [7], BEND [33], and Genomics LRB [34]. In selecting tasks for comparison, we adhere to the principle of encompassing a diverse range of challenges, encompassing both short and long-range capabilities. Additionally, we prioritize tasks that are biologically meaningful, covering a variety of species and functionalities within DNA-related areas.

4.1 Pretraining curves

Configurations

We train three variants of HybriDNA with 300M, 3B, and 7B parameters. These models differ in the number of layers, hidden size, and learning rates used during training. Despite these variations, all models share a consistent pretraining strategy. Instead of the commonly used masked language modeling (MLM) loss in genomics foundation models, we use a next-token-prediction (NTP) loss to enable generative capability. Our models are trained on NVIDIA A100/H100 and AMD MI300X GPUs. Details of the model architecture and pertaining configurations can be found in Appendix A.

During the pretraining stage, our 300M, 3B and 7B parameter models are trained on 0.5M tokens per batch, optimized for efficient utilization of computational resources and consistent training dynamics. Initially, the models are pretrained on sequences with a context length of 8,192 for 500k steps, resulting in a total of 250B tokens (approximately 1.5 epochs) in the first pretraining stage. Following this, the models undergo further pretraining to extend their capabilities to handle larger context lengths. This two-stage pretraining strategy allows the models to gradually adapt to more complex and computationally demanding settings, ensuring robust performance across varying sequence lengths.

Scaling Behaviors

To investigate scaling law behavior, we analyze the training and validation losses during the pretraining stage for the 300M, 3B, and 7B models. As the model size increases, we observe consistent improvements in both training and validation losses, highlighting the advantages of larger models in capturing intricate genomic patterns. Detailed loss curves for each model variant are presented in Fig. 3, demonstrating the impact of scaling parameter sizes on pretraining model quality. These findings align

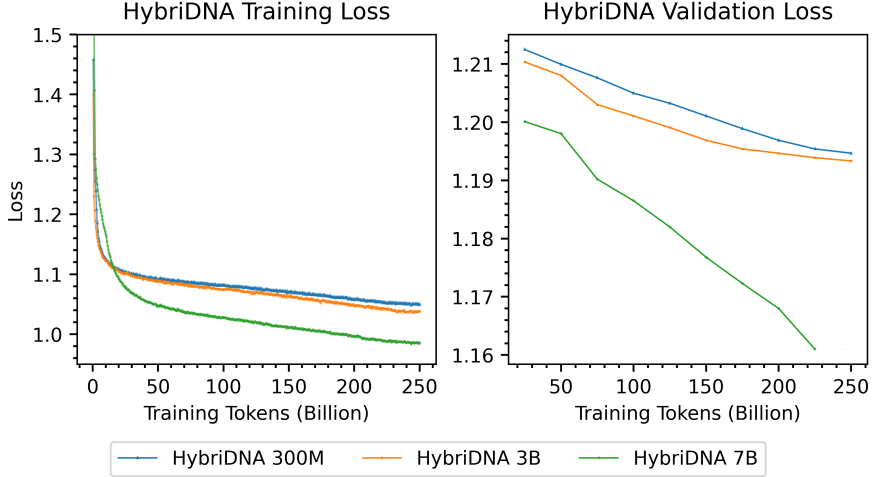


Fig. 3: Pretraining loss curves for HybriDNA-300M, 3B, and 7B models

with theoretical expectations of scaling laws in deep learning and genomics-specific modeling.

We further demonstrate the effectiveness of the hybrid architecture compared to the pure Mamba2 architecture. The training and validation losses for pretraining two comparable models, each with 300 million parameters, are shown in Fig. A2.

4.2 Baselines

We compare HybriDNA with 5 state-of-the-art DNA foundation models: NT-500M-human, NT-2.5B-MS, DNABERT-2, HyenaDNA-medium-160k and Caduceus-Ph-131k.

NT-500M-human is a Transformer-encoder based model with 500M parameters pretrained on GRCh38/hg38 human reference genome, which contains approximately 3.2B nucleotide bases. It employs a k-mer tokenization method with k set to 6 and is pretrained using the standard Masked Language Modeling (MLM) objective. The pretraining context length is 1,000.

NT-2.5B-MS is another variant in the Nucleotide Transformer series, featuring a larger size of 2.5B parameters. It was pretrained on a multi-species dataset comprising 174 billion nucleotides from a total of 850 species. The other pretraining details are consistent with those of NT-500M-human.

DNABERT-2 is also a Transformer-encoder based model with 112M parameter size. It is pretrained on a multi-species dataset containing approximately 32.5 billion nucleotide bases from 135 species. It improves the tokenizer with the Byte Pair Encoding (BPE) method and is also trained with a standard MLM objective with a context length of 512.

Caduceus-Ph-131k builds upon the Mamba1 architecture, which employs selective state space models for long-range sequence processing. It enables Bi-directional

sequence modeling using BiMamba block. The model is trained on GRCh37/hg37 human reference genome with about 3.2B nucleotide bases using MLM objective. This variant has 7.73M parameter size and uses nucleotide-level tokenization with a pretraining context length of 131,072.

HyenaDNA-Medium-160k utilizes the Hyena operator, derived from state space models, for computationally efficient long-range sequence modeling. It is pre-trained on GRCh38/hg38 human reference genome using next-token prediction objective. This specific variant has 14.2M parameters and uses nucleotide-level tokenization with a pretraining context length of 160,000.

For all baseline models, we utilize the pretrained weights provided in their respective original codebase. Since our evaluation is mainly focused on eukaryote-related tasks, Evo is excluded from the comparison.

4.3 Short-range understanding benchmark (GUE, BEND)

Genome Understanding Evaluation (GUE) [7] aggregates 28 datasets across 9 tasks, encompassing input lengths from 70 to 512 base pairs. GUE serves as a standardized evaluation suite, measuring the effectiveness of genomic foundation models in multi-species genome classification. The GUE benchmark assesses a model’s ability to analyze short-range genomic sequences across multiple biologically significant tasks. Promoter Detection (PD) and Core Promoter Detection (CPD) identify regulatory regions that initiate transcription, crucial for understanding gene expression control. Splice Site Detection (SS) predicts locations where pre-mRNA is spliced, affecting transcript diversity and protein function. Transcription Factor Binding Site (TFBS) Prediction determines whether a sequence contains motifs for transcription factors, which regulate gene activity. Epigenetic Marks Prediction (EMP) involves predicting histone modifications and DNA methylation, key regulators of chromatin state and gene expression. Lastly, COVID Variant Classification (CV) tracks viral genome mutations, aiding epidemiological surveillance. Collectively, these tasks provide a comprehensive measure of a model’s capacity to decode genomic structure and function.

BEND [33] evaluates models on a collection of realistic and biologically meaningful tasks defined on the human genome. It emphasizes the importance of capturing intricate genomic data features through tasks that are comprehensive and provide a standardized evaluation methodology for genomics foundation models. For evaluation, we select the three largest short-range tasks across 3 different datasets from the benchmark. Chromatin Accessibility Prediction assesses whether DNA is in an open or closed state, influencing gene expression potential. Histone Modification Prediction determines chemical changes to histones that affect chromatin structure and transcriptional activity. CpG Methylation Prediction identifies DNA methylation patterns at CpG sites, which play a critical role in gene silencing and disease progression. These tasks are essential for understanding the regulatory landscape of the genome and its implications for cellular function and disease mechanisms.

Type	Model	PD(H) (MCC)	CPD(H) (MCC)	SS(H) (MCC)	TF(H) (MCC)	TF(M) (MCC)	EMP(Y) (MCC)	CV(V) (F-1)
Encoder	DNABERT-2	83.96	71.81	85.42	68.71	70.00	55.98	71.02
	NT-2.5B-MS	88.15	71.57	89.35	63.21	67.02	57.64	73.04
	NT-500M-human	82.96	66.79	78.63	61.92	45.24	45.35	50.82
	Caduceus-Ph	82.36	67.03	71.80	65.17	62.28	51.05	40.35
Decoder	HyenaDNA	80.14	69.22	77.76	61.74	64.39	47.15	25.88
Our	HybriDNA-300M	83.29	68.87	87.74	68.37	75.32	67.38	73.81
	HybriDNA-300M (E)	83.67	69.96	88.72	69.70	75.73	68.25	73.90
	HybriDNA-3B	85.40	69.50	89.01	70.48	75.43	69.06	74.05
	HybriDNA-3B (E)	85.55	70.71	89.10	71.13	77.14	68.97	74.88
	HybriDNA-7B	86.53	71.37	90.09	70.72	78.02	63.05	74.02
	HybriDNA-7B (E)	88.10	72.03	90.12	72.01	79.02	65.30	74.30

Table 2: Results on the GUE Benchmark, which encompass a series of short-range classification tasks across multiple species, including Promoter Detection (PD), Core Promoter Detection (CPD), Splice Site Detection (SS), Transcription Factor Prediction (TF), Epigenetic Marks Prediction (EMP) and Covid Variant Classification (CV). The suffix “(H)” denotes the human genome, “(M)” the mouse genome, “(Y)” the yeast genome, and “(V)” the virus genome. Additionally, the suffix “(E)” in the model name indicates that echo embedding was applied during discriminative fine-tuning.

4.3.1 GUE

Settings DNABERT-2 [7] introduces the GUE benchmark, which includes a series of short-range classification tasks across multiple species, such as Promoter Detection (PD), Core Promoter Detection (CPD), Splice Site Detection (SS), Transcription Factor Prediction (TF) for human and yeast, Epigenetic Marks Prediction (EMP) and Covid Variant Classification (CV). Following the same setting as DNABERT-2, we use metrics of Matthews Correlation Coefficient (MCC) for all tasks, except for the Covid task, where we use the F-1 score according to the GUE dataset’s original setting. The hyperparameters, training epochs, and evaluation strategies follow exactly from the original paper. We fine-tune all model parameters and use the hidden state of the last token for embedding representation for decoder-only models. Training epochs and evaluation steps are also consistent with the original paper. We apply a learning rate of $5e-5$ for our 300M model, $3e-5$ for our 3B model, and $1e-5$ for our 7B model across all tasks.

Results We take the mean MCC/F1-score value of tasks in the same category in the following table. For detailed results on each task, refer to Appendix B. The suffix “(E)” in the model name within the table indicates that echo embedding was applied during discriminative fine-tuning. Results of our model are presented in Table 2.

4.3.2 BEND

Settings BEND paper [33] presents a series of biologically meaningful tasks for genomics foundation model derived from a series of studies. We select the three short-range supervised tasks: Chromatin Accessibility, Histone Modification, and CpG Methylation. For the Histone Modification Tasks, the training process follows the original paper. We freeze the embedding of the models and train a downstream CNN

Model Type	Model	Chromatin Accessibility (AUROC)	Histone Modification (AUROC)	CpG Methylation (AUROC)
Encoder Models	DNABERT-2	0.81	0.79	0.90
	NT-2.5B-MS	0.79	0.78	0.92
	NT-500M-human	0.74	0.76	0.88
	Caduceus-Ph	0.83	0.77	0.91
Decoder Models	HyenaDNA	0.81	0.77	0.87
Our Model	HybriDNA-300M	0.78	0.77	0.88
	HybriDNA-3B	0.82	0.79	0.92
	HybriDNA-7B	0.84	0.79	0.93

Table 3: Results on the BEND Benchmark, which includes Chromatin Accessibility, Histone Modification, and CpG Methylation tasks.

model for 100 epochs. For autoregressive models like HyenaDNA and our HybriDNA model, we use the mean of the hidden state of the sequence as embedding representations. The model with the lowest validation loss is tested and the metric reported is the mean AUROC score.

Results We directly report the AUROC score of each model on the three tasks in Table 3.

4.4 Genomics long-range benchmark (LRB)

The Genomics Long-Range Benchmark (LRB) [34] is designed to evaluate tasks that require understanding long-range context within genomic sequences. To assess models’ ability to capture dependencies across extended genomic regions, we select two representative tasks across distinct datasets that inherently demand long-range sequence comprehension: Causal eQTL Variant Effect Prediction and OMIM Variant Effect Prediction. The Causal eQTL Variant Effect Prediction task evaluates the impact of genetic variants on gene expression levels, linking non-coding mutations to functional changes. The OMIM Variant Effect Prediction task focuses on identifying pathogenic mutations associated with Mendelian diseases, aiding in genetic diagnostics and precision medicine. These tasks test a model’s proficiency in analyzing complex gene regulation and variant interpretation over extended genomic regions.

During discriminative fine-tuning, all models were trained using frozen embeddings generated in the same way as those used in the BEND benchmark. These embeddings were passed through an MLP classifier, with all models sharing identical architectures and hyperparameters to maintain consistency. We report accuracy and AUROC metrics for all tasks, with detailed results summarized in Table 4.

4.5 Designing realistic synthetic cis-regulatory elements (CREs)

regLM [17] is a framework that combines autoregressive language models with supervised sequence-to-function tasks to design synthetic cis-regulatory elements (CREs).

Model Type	Model	Causal eQTL		OMIM
		<i>Fine-tune</i> (AUROC)	<i>Zero-shot</i> (AUROC)	<i>Zero-shot</i> (AUPRC)
Encoder Models	DNABERT-2	0.72	0.50	0.002
	NT-500M-human	0.72	0.51	0.003
	Caduceus-Ph	0.68	0.49	0.002
Decoder Models	HyenaDNA	0.71	0.51	0.002
Our Model	HybridDNA-300M (8k)	0.71	0.51	0.003
	HybridDNA-300M (32k)	0.72	0.51	0.003
	HybridDNA-300M (131k)	0.74	0.51	0.003

Table 4: Results on the LRB Benchmark, which includes Causal eQTL Variant Effect prediction, OMIM Variant Effect prediction tasks.

This framework highlights the capability of models to generate regulatory sequences with specific desired properties. Specifically, (1) Human Enhancer Generation involves designing enhancers that drive gene expression in specific cell types, which is crucial for gene therapy and functional genomics. (2) Yeast Promoter Generation focuses on engineering promoters with defined transcriptional activity, supporting biotechnology applications such as industrial enzyme production. These tasks demonstrate a model’s capacity to generate biologically plausible DNA sequences that can be experimentally validated for targeted gene regulation.

We followed the experimental setup established by **regLM**. For comparison, we utilized HyenaDNA [11], which is the default decoder-only model in the regLM study. Further details on the fine-tuning configurations can be found in Appendix A.5.

4.5.1 Cell type-specific human enhancer generation

The human enhancer generation task focuses on designing desired human enhancer genomic sequences for three specific cell line types: HepG2, K562, and SK-N-SH. Each sequence includes a three-digit label (ranging from 0 to 3) that represents the enhancer’s activity strength in a given cell line. The model is fine-tuned on a dataset consisting of 670,000 training samples of 200bp enhancers with varying activity levels, utilizing the prompt tokens described in Section 3.3.2.

During evaluation, the model is tasked with generating 600 sequences for the specific labels {300, 030, 003}, which represent high enhancer activity in a particular cell line. These labels are extremely rare in the training data, accounting for only 0.16% of the total training set. These generated sequences are evaluated using a separately trained scoring model from regLM to assess their actual enhancer activity in the respective cell types. Specifically, beam search decoding is used for all models during sequence generation to ensure a fair comparison. The baseline for evaluation is the fine-tuned HyenaDNA model variant, "hyenadna-medium-160k-seqlen," as referenced in the original regLM paper.

The evaluation metrics for the generated sequences are defined as follows:

1. Top-1 activity: The highest predicted enhancer activity for each cell type.

Model	HepG2		K562		SK-N-SH		Diversity	
	Top-1	Mean	Top-1	Mean	Top-1	Mean	Mean Edit Distance	
Held-out Test	6.2	2.6	5.5	2.4	5.1	1.6	110.10	
HyenaDNA	5.5	4.0	4.3	3.8	5.2	2.3	98.50	
HybriDNA-300M	7.3	5.4	7.8	6.2	6.6	4.7	108.74	

Table 5: Comparison between HybriDNA-300M and HyenaDNA on the human enhancer generation task. Metrics include Top-1 activity, Mean activity, and Diversity for each cell line type (HepG2, K562, SK-N-SH).

2. Mean activity: The average of the top 100 predicted enhancer activity scores for each cell type.
3. Diversity: The mean of pair-wise edit distance of the top 100 predicted sequences across all cell types, measuring the overall diversity of high-quality generated sequences.

The results are summarized in Table 5. Result shows that our model outperforms both the held-out test set and HyenaDNA in generating higher-activity enhancer sequences for each cell line, while also maintaining greater diversity.

4.5.2 Yeast promoter generation

The yeast promoter generation task follows a similar setup to the human enhancer generation task. However, instead of three cell lines, this task utilizes a two-digit label representing promoter activity in complex and defined media, with activity levels ranging from 0 to 4. Since HyenaDNA was pretrained primarily on human genomic data, the regLM study trained the HyenaDNA model from scratch using yeast promoter data. In contrast, our model’s pretraining corpus already includes multi-species data, including yeast genomes. As a result, rather than training our model from scratch, we fine-tune it from a pretrained checkpoint using the yeast promoter dataset.

The model is fine-tuned on a dataset comprising 7.4 million training samples of 80bp promoters with varying activity levels. During evaluation, it is prompted to generate sequences with label {40, 04}, which represents only 0.34% of the training samples. The evaluation steps and metrics are identical to those in the human enhancer generation task. The results are summarized in Table 6. Notably, our model achieves higher promoter activity scores across both media types and generates more diverse sequences compared to the HyenaDNA baseline.

4.6 Computational efficiency

To evaluate the computational efficiency of the HybriDNA model compared to a standard Transformer model with a similar parameter size, particularly during the training phase, we use the following two metrics:

1. Tokens/second per GPU: This metric assesses the throughput of a single GPU by measuring the number of tokens it can process each second during the pretraining

Model	Complex Media		Defined Media		Diversity
	Top-1	Mean	Top-1	Mean	
Held-out Test	16.0	5.9	15.7	6.7	28.8
HyenaDNA	16.8	11.4	16.3	10.8	27.3
HybriDNA-300M	18.2	15.0	17.6	13.5	30.7

Table 6: Comparison between HybriDNA-300M and HyenaDNA on the yeast promoter generation task. Metrics include Top-1 activity, Mean activity, and Diversity for both media types (Complex and Defined).

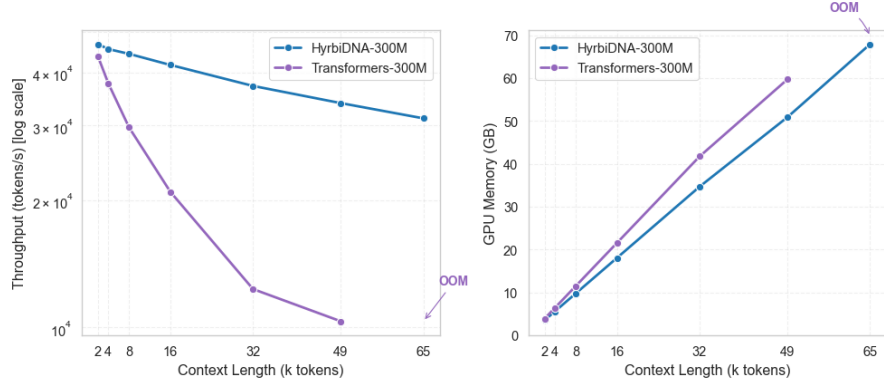


Fig. 4: Comparison of training throughput and GPU memory consumption between HybriDNA and a pure Transformer model with comparable parameters (e.g., 300M)

stage. For each context length, the batch size is set to the maximum number to fit within the GPU memory.

2. GPU memory cost (GB): This metric measures the amount of GPU memory consumed when training models with a fixed context length and a batch size of 1.

We evaluate both models using four NVIDIA A100 GPUs (80G memory) with DeepSpeed Zero-1 Stage optimization and BF16 mixed-precision training. Both models comprise approximately 300M parameters. The Transformer model uses Flash Attention 2 optimization, while the Mamba2 layers in the HybriDNA model are implemented using CUDA. We test the pretraining of both models at various context lengths, ranging from 2k tokens to 65k tokens, doubling the sequence length at each step.

As shown in Fig. 4, the HybriDNA model achieves significantly higher training throughput than standard Transformer models, especially when processing context lengths exceeding 32,000 tokens. For instance, at a context length of 49,000 tokens, the throughput of HybriDNA is approximately 3.4 times higher than that of Transformers. This performance gap widens as context length increases, highlighting the superior efficiency of our model compared to Transformers.

In terms of GPU memory usage, HybriDNA consistently demonstrates greater efficiency than Transformer models, even those optimized with advanced techniques such as Flash Attention 2 [35]. Notably, at context lengths around 65,000 tokens, a standard Transformer model encounters Out-Of-Memory (OOM) issues on A100 GPUs. These findings underscore the exceptional capability of our hybrid model to effectively manage larger context lengths, a crucial advantage for complex long-range DNA-related tasks.

5 Related Work

5.1 DNA Foundation Models

The advent of high-throughput sequencing technologies has produced vast amounts of genomic data, presenting an unprecedented opportunity for deep learning to uncover complex relationships and dependencies in DNA sequences. Recent advancements in genome language modeling have demonstrated their effectiveness across a wide range of downstream applications, including promoter prediction [36, 37], gene expression prediction [38], DNA methylation prediction [39], chromatin state analysis [40], promoter-enhancer interaction prediction [41, 42] TF-DNA binding prediction [43], variant effect prediction [44], gene network prediction [45] and more. More recently, inspired by advancements in natural language processing, researchers have begun developing DNA foundation models. These include, but are not limited to: (1) encoder-only models such as DNABERT, DNABERT-2, Nucleotide Transformer, and Caduceus; and (2) decoder-only models such as HyenaDNA and Evo.

DNABERT [46] is an early foundation model designed to interpret the human genome from a language perspective. By adapting the BERT framework with Transformers architecture, it captures a transferable understanding of human genome reference sequences. This single pretrained Transformer model achieves state-of-the-art performance in tasks such as predicting promoters, splice sites, and transcription factor binding sites, after fine-tuning on small task-specific labeled datasets. The model contains 86M parameters and operates with a context length of 512 on the hg38 human reference genome dataset.

DNABERT-2 [7] builds on its predecessor by employing Byte Pair Encoding (BPE) for tokenization, which improves computational efficiency and representation quality. It also incorporates Attention with Linear Biases (ALiBi) with Transformers-Encoder layers, enabling the model to process longer input sequences effectively. DNABERT-2 achieves state-of-the-art results on the Genome Understanding Evaluation (GUE) benchmark, showcasing its capacity to address diverse genomic tasks. The model consists of 112M parameters and is trained on a multi-species dataset comprising 135 species with a total of 32 billion nucleotides and a context length of 512.

Nucleotide Transformer (NT) [8] is a scalable genomics foundation model, built on an encoder-only Transformer architecture, with parameter sizes ranging from 500M to 2,500M, based on encoder-only Transformer architecture. Its multi-species variant is pretrained on genomic data from 850 species, employing a non-overlapping

k-mer tokenization method that effectively reduces tokenized sequence lengths. Additionally, two human-specific versions are trained separately on the hg38 human reference genome dataset and the 1000 Genomes Project. All pretraining is conducted with a context length of 1,000 tokens.

Caduceus [47] introduces the bi-directional Mamba1 architecture, specifically designed for DNA sequence modeling. By incorporating reverse complement (RC) equivariance at the architectural level, Caduceus is optimized for long-range DNA sequence modeling. The model effectively captures the intricate understanding required for DNA sequence tasks. The Caduceus series features parameter sizes ranging from 500K to 7M, with a context length of 131k, and is trained on the hg38 dataset.

HyenaDNA [11] utilizes the Hyena operator, a recurrence of gating and implicitly parametrized long convolutions, to handle long-range genomic sequences, enabling the processing of input contexts up to 1 million tokens with single-nucleotide resolution. This model shows effectiveness in tasks requiring long-range understanding, such as analyzing DNA fragments far apart, beyond the context window of traditional Transformer models. HyenaDNA is trained on the hg38 human reference genome dataset, with parameter sizes ranging from 1.7M to 50M and context lengths varying from 1k to 1M.

Evo [12] is a 7-billion-parameter foundation model built on the StripedHyena architecture and trained on 2.7 million raw prokaryotic and phage genome sequences. It integrates multiple biological modalities, including DNA, RNA, and proteins. With a context length of 131k nucleotide bases, Evo delivers superior performance in sequence modeling and functional design tasks, spanning molecular to genome-scale applications.

Two concurrent works, **GenomeOcean** [48] and **GENERator** [49], have recently emerged in the field of DNA foundation models. GenomeOcean is a 4B-parameter model trained on diverse meta-genomics samples, optimizing for microbial species representation and achieving faster genome generation. GENERator, a 1.2B-parameter model with a 98k context length, is trained on 386B base pairs of eukaryotic DNA and excels in generating protein-coding sequences, designing promoter sequence and optimizing promoter activity. These models further expand genomic sequence modeling and generation capabilities.

5.2 Hybrid Models in General Domains

Recent advancements in Mamba-based hybrid models for NLP tasks combine the efficiency of SSMS with the expressiveness of attention mechanisms, excelling in long-context scenarios. Innovations include Jamba’s [16] integration of Transformer, Mamba, and Mixture-of-Experts layers for sequences up to 256k tokens, Zamba’s [28] compact 7B model with shared self-attention for reduced latency, and SAMBA’s [29] sliding window attention for efficient handling of sequences up to 1M tokens. Other notable contributions include Taipan’s [50] selective attention layers for scalability and Waleffe’s [15] versatile 8B hybrid architecture combining Mamba2, self-attention, and MLP layers. These models achieve strong results across various short- and long-range benchmarks.

6 Conclusion

In this work, we develop a class of decoder-only DNA language models built on a hybrid Transformer-Mamba2 architecture. This design harnesses the unique strengths of its two core components to enable efficient and precise modeling of DNA sequences. By integrating Mamba2 layers, our model can process extremely long DNA sequences at single-nucleotide resolution with remarkable computational efficiency. Pretrained on large-scale, multi-species genomes at single-nucleotide resolution with a next-token prediction objective, HybriDNA demonstrates foundational capabilities in both understanding and designing genomic sequences. Through echo embedding discriminative fine-tuning, HybriDNA achieves state-of-the-art performance across 33 biologically significant DNA understanding tasks from the BEND, GUE, and LRB benchmarks. Through generative fine-tuning, HybriDNA exhibits remarkable proficiency in generating synthetic cis-regulatory elements with desirable functional properties. These results highlight HybriDNA’s versatility and establish its potential as a powerful foundation model for advancing DNA research and applications.

Looking ahead, there are several exciting directions to further explore. These include: (1) Expanding the pretraining dataset to include a greater number of nucleotide tokens and species classes, enabling broader generalization across downstream tasks involving diverse species. (2) Conducting more downstream fine-tuning tasks with diverse and significant scientific impacts, and performing wet-lab experiments to further validate the sequences designed by HybriDNA.

7 Acknowledgements

We would like to thank Sumit Basu from the Health Futures team at Microsoft Research for his valuable insights on benchmarks and writing; Jingyun Bai for her design; and Ran Bi, Hannes Schulz, Jean Helie, and Maik Riechert for their engineering support. Additionally, we sincerely acknowledge the entire AI for Science team at Microsoft Research for their continuous support.

Appendix A Model Details

A.1 Model architecture

The **HybriDNA** model employs a hybrid Transformer-Mamba2 architecture. The architecture interleaves Transformer and Mamba2 layers in a 7:1 ratio. The Transformer layer is placed in the fourth of every eight layers. Our three model variants—300M, 3B, and 7B—differ in their hidden dimension size and layer configurations. The details of each model architecture are summarized in Table A1.

Model Variant	# Layers	Hidden Size	Intermediate Size	# Heads	Head Dim
7B	32	4096	8192	128	64
3B	16	4096	8192	128	64
300M ¹	24	1024	2048	32	64

Table A1: Model configurations for three variants of HybriDNA

A.2 Pretraining configuration

During the pretraining stage, we trained the HybriDNA models using a next token prediction objective. The training utilized the Adam [51] optimizer with a learning rate schedule and standard exponential decay rates $\beta_1 = 0.9$, $\beta_2 = 0.95$, and $\epsilon = 1e-8$. All model variants underwent a Warmup phase of 2,000 steps, with a total of 500,000 training steps. The learning rates were set at $1e-3$ for the 300M model, $6e-4$ for the 3B model, and $1e-4$ for the 7B model. Mamba-based models demonstrate a higher tolerance for learning rates compared to standard Transformer architectures, highlighting their stability during optimization.

Our pretraining was conducted on the following hardware configurations: the 300M model on 8 AMD MI300X GPUs, the 3B model on 8 NVIDIA H100 GPUs, and the 7B model on 64 AMD MI300X GPUs. Models are trained for approximately 300 hours for the 300M and 7B variants, and 500 hours for the 3B model. These configurations ensured efficient utilization of computational resources and stable training for large-scale models.

A.3 Pretraining dataset

We utilized a comprehensive dataset comprising approximately 200 billion tokens, following the Nucleotide Transformer’s multi-species dataset [8]. It comprises of a subset of the NCBI dataset with 850 species and the details have been presented in the methodology section.

¹You may notice that HybriDNA-300M model has 32 layers. We draw inspiration from the configuration of Jamba-1.5 model: <https://huggingface.co/ai21labs/AI21-Jamba-1.5-Large/blob/main/config.json>.

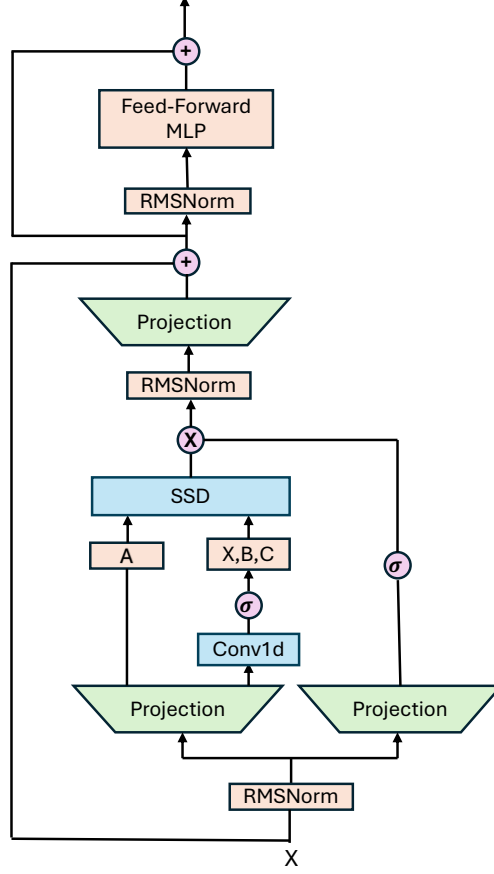


Fig. A1: HybriDNA Mamba2 Block

A.4 Effectiveness of hybrid models

To evaluate the effectiveness of incorporating Transformer layers alongside Mamba2 layers in the HybriDNA model, we pretrained a 300M-size variant without any Transformer layers. Both our Hybrid model and the pure Mamba2 model were pretrained using an 8k context length. Fig. A2 presents the training and validation losses during the pretraining stage. It is evident that for models with similar parameter sizes, the hybrid model demonstrates lower training and validation losses compared to a model composed entirely of Mamba2 blocks.

A.5 Fine-tuning configurations for downstream Tasks

In this section, we provide detailed fine-tuning procedures for each downstream dataset.

GUE For the GUE benchmark, we adhered to the specific settings for each task, including the warmup steps and training/validation steps that are customized for each

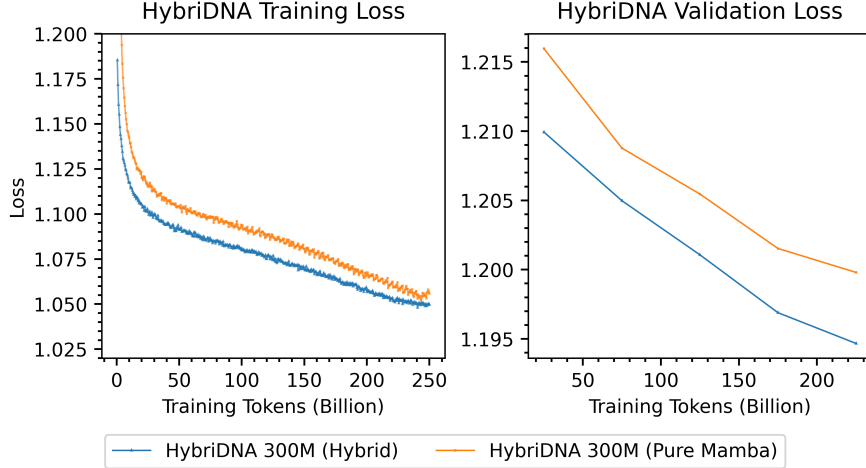


Fig. A2: Effectiveness of Hybridization in HybriDNA

task. The only adjustment we’ve made was to the learning rate: $5e-5$ for the 300M model, $3e-5$ for the 3B model, and $1e-5$ for the 7B model. We applied a simple classification head for the model, using either the hidden state of the last token for classification in the standard setting, or the averaged hidden state of the repeated sequence input for the echo embedding setting.

BEND For the BEND dataset, we applied the exact settings across all three tasks, with a learning rate of $3e-3$ and training for 100 epochs. The learning rate decreases linearly to 0, and we use the epoch with the lowest validation loss for the final test-set evaluation. In BEND, the model’s embedding input is frozen, and only a downstream two-layer CNN model is fine-tuned for classification. The extracted hidden state is the average of hidden states of the input in the standard setting and the average of the repeated sequence in the echo-embedding setting.

Genomics LRB Following the Genomics LRB work [52], we fine-tune all the parameters of the model for all fine-tuning tasks, using a subsequent MLP for classification. We adhere to the benchmark’s settings by inputting sequences with the pretraining context length into the model and averaging across the same length window for the same task across different models. For zero-shot tasks, we employ sequence-level probability for a regression correlation coefficient analysis, consistent with the benchmark’s method.

regLM For the HyenaDNA baseline, we used the exact setting of the regLM model to load its fine-tuned checkpoint. For fine-tuning our HybriDNA-300M model, we conducted 16 epochs on the human enhancer task with a learning rate of $1e-4$. For the Yeast Promoter task, as our model has already been pretrained on multi-species data, including yeast sequence, we fine-tuned it on the dataset for only 2 epochs, with the same learning rate of $1e-4$. Validation was performed every 400 steps for each task, and we saved the model with the highest validation accuracy as the final model. During generation, we use beam search with a beam size of 2500 for the human enhancer task

and a beam size of 256 for the yeast promoter task to align with the total number of generated sequence of the original regLM setting. The activity scoring models used were the same as those in the original regLM models. The diversity metric is calculated as the mean pair-wise edit distance of the top-100 activity sequence across all labels for each task.

Appendix B GUE Benchmark Results

Model Type	Model	Transcription Factor Prediction (Human)				
		0 (MCC)	1 (MCC)	2 (MCC)	3 (MCC)	4 (MCC)
Baselines	DNABERT-2	70.89	74.49	66.62	60.35	71.21
	NT-2.5B-MS	66.46	70.25	58.70	51.28	69.34
	NT-500M-human	60.03	69.34	47.02	39.27	58.84
	Caduceus-Ph-131k	70.69	69.00	61.13	55.98	69.07
	HyenaDNA-160k	64.47	70.74	60.44	39.78	73.27
Our Model	HybridDNA-300M	68.12	67.13	70.29	55.52	80.80
	HybridDNA-300M(E)	67.64	71.28	70.84	57.92	80.80
	HybridDNA-3B	69.88	69.24	72.21	56.44	84.61
	HybridDNA-3B(E)	69.02	70.82	72.80	58.01	85.02
	HybridDNA-7B	70.00	74.47	70.42	64.52	85.03
	HybridDNA-7B(E)	71.46	75.60	71.81	65.82	86.20

Table B2: Results for Transcription Factor Prediction (DNABERT2-Human) in the GUE benchmark

Model Type	Model	Promoter Detection (Human)			Splice Site Prediction (Human)
		all (MCC)	notata (MCC)	tata (MCC)	reconstruct (MCC)
Baselines	DNABERT-2	86.64	94.20	71.04	85.42
	NT-2.5B-MS	91.00	94.02	79.43	89.35
	NT-500M-human	81.34	88.73	78.82	78.63
	Caduceus-Ph-131k	83.98	92.13	70.96	71.80
	HyenaDNA-160k	83.04	91.03	66.36	77.76
Our Model	HybridDNA-300M	88.94	94.44	69.63	87.74
	HybridDNA-300M(E)	88.81	94.45	68.45	88.72
	HybridDNA-3B	89.48	94.49	72.24	89.01
	HybridDNA-3B(E)	89.30	94.33	73.02	89.10
	HybridDNA-7B	88.28	94.73	73.59	90.09
	HybridDNA-7B(E)	90.20	94.57	76.84	90.12

Table B3: Results for Promoter Detection and Splice Reconstruct (DNABERT2-Human) in the GUE benchmark

Model Type	Model	Core Promoter Detection (Human)		
		all (MCC)	notata (MCC)	tata (MCC)
Baselines	DNABERT-2	69.97	69.62	75.83
	NT-2.5B-MS	70.28	71.49	72.95
	NT-500M-human	63.36	64.67	72.34
	Caduceus-Ph-131k	64.09	68.35	68.65
	HyenaDNA-160k	66.18	67.41	74.07
Our Model	HybriDNA-300M	68.40	69.12	69.09
	HybriDNA-300M(E)	68.37	69.15	72.36
	HybriDNA-3B	68.98	69.63	69.89
	HybriDNA-3B(E)	68.90	70.01	73.21
	HybriDNA-7B	66.50	70.66	76.94
	HybriDNA-7B(E)	67.10	71.53	77.49

Table B4: Results for Core Promoter Detection (DNABERT2-Human) in the GUE benchmark

Model Type	Model	Transcription Factor Prediction (Mouse)					Classification (Virus)
		0 (MCC)	1 (MCC)	2 (MCC)	3 (MCC)	4 (MCC)	Covid (F-1)
Baselines	DNABERT-2	56.76	84.77	79.32	66.47	52.66	71.02
	NT-2.5B-MS	63.31	83.76	71.52	69.44	47.07	73.04
	NT-500M-human	31.04	75.04	61.67	29.17	29.27	50.82
	Caduceus-Ph-131k	50.44	82.63	73.81	61.13	43.40	40.35
	HyenaDNA-160k	56.25	80.46	78.14	60.83	46.25	25.88
Our Model	HybriDNA-300M	68.57	83.46	86.02	87.96	50.58	73.81
	HybriDNA-300M(E)	68.66	85.62	85.39	87.78	51.20	73.90
	HybriDNA-3B	70.96	84.18	89.63	88.59	50.97	74.05
	HybriDNA-3B(E)	71.02	84.30	89.78	88.20	52.39	74.88
	HybriDNA-7B	71.68	87.75	86.59	87.62	56.47	74.02
	HybriDNA-7B(E)	72.91	88.64	87.64	88.59	57.33	74.30

Table B5: Results for Transcription Factor Prediction (DNABERT2-Mouse) and Covid Variant Classification (DNABERT2-Virus) in the GUE benchmark

References

- [1] Devlin, J., Chang, M.-W., Lee, K. & Toutanova, K. Burstein, J., Doran, C. & Solorio, T. (eds) *BERT: Pre-training of deep bidirectional transformers for language understanding*. (eds Burstein, J., Doran, C. & Solorio, T.) *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, 4171–4186 (Association for Computational Linguistics, Minneapolis, Minnesota, 2019). URL <https://aclanthology.org/N19-1423/>.
- [2] Bommasani, R. *et al.* On the opportunities and risks of foundation models. *arXiv preprint arXiv:2108.07258* (2021).

- [3] Achiam, J. *et al.* Gpt-4 technical report. *arXiv preprint arXiv:2303.08774* (2023).
- [4] Brandes, N., Ofer, D., Peleg, Y., Rappoport, N. & Linial, M. Proteinbert: a universal deep-learning model of protein sequence and function. *Bioinformatics* **38**, 2102–2110 (2022).
- [5] Lin, Z. *et al.* Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science* **379**, 1123–1130 (2023).
- [6] He, L. *et al.* Sfm-protein: Integrative co-evolutionary pre-training for advanced protein sequence representation. *arXiv preprint arXiv:2410.24022* (2024).
- [7] Zhou, Z. *et al.* Dnabert-2: Efficient foundation model and benchmark for multi-species genome. *arXiv preprint arXiv:2306.15006* (2023).
- [8] Dalla-Torre, L., Benegas, N., Grechishnikova, D. *et al.* The nucleotide transformer: Building and evaluating robust foundation models for human genomics. *Nature Methods* (2023).
- [9] Poli, M. *et al.* Hyena hierarchy: Towards larger convolutional language models (2023). URL <https://arxiv.org/abs/2302.10866>. 2302.10866.
- [10] Radford, A., Narasimhan, K., Salimans, T. & Sutskever, I. Improving language understanding by generative pre-training. *OpenAI Blog* (2018). URL https://cdn.openai.com/research-covers/language-unsupervised/language_understanding_paper.pdf.
- [11] Poli, M., Dao, T. *et al.* Hyenadna: Long-range genomic sequence modeling at single nucleotide resolution. *arXiv preprint arXiv:2306.15006* (2023).
- [12] Meier, J. *et al.* Sequence modeling and design from molecular to genome scale with evo. *Science* **382**, eado9336 (2023).
- [13] Gu, A. & Dao, T. Mamba: Linear-time sequence modeling with selective state spaces (2023). 2312.00752.
- [14] Dao, T. & Gu, A. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality. *arXiv preprint arXiv:2405.21060* (2024).
- [15] Waleffe, R. *et al.* An empirical study of mamba-based language models. *arXiv preprint arXiv:2406.07887* (2024).
- [16] Lieber, O. *et al.* Jamba: A hybrid transformer-mamba language model (2024). URL <https://arxiv.org/abs/2403.19887>. 2403.19887.
- [17] Lal, A., Garfield, D., Biancalani, T. & Eraslan, G. reglm: Designing realistic regulatory dna with autoregressive language models. *bioRxiv preprint* (2024).

- [18] Schulman, J. *et al.* Chatgpt: Optimizing language models for dialogue. *OpenAI blog* **2** (2022).
- [19] Bahdanau, D. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473* (2014).
- [20] Vaswani, A. Attention is all you need. *Advances in Neural Information Processing Systems* (2017).
- [21] Gu, A., Goel, K. & Ré, C. Efficiently modeling long sequences with structured state spaces (2022). URL <https://arxiv.org/abs/2111.00396>. 2111.00396.
- [22] Gu, A. *et al.* Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in neural information processing systems* **34**, 572–585 (2021).
- [23] Gupta, A., Gu, A. & Berant, J. Diagonal state spaces are as effective as structured state spaces. *Advances in Neural Information Processing Systems* **35**, 22982–22994 (2022).
- [24] Gu, A., Goel, K., Gupta, A. & Ré, C. On the parameterization and initialization of diagonal state space models. *Advances in Neural Information Processing Systems* **35**, 35971–35983 (2022).
- [25] Smith, J. T. H., Warrington, A. & Linderman, S. W. Simplified state space layers for sequence modeling (2023). URL <https://arxiv.org/abs/2208.04933>. 2208.04933.
- [26] Fu, D. Y. *et al.* Hungry hungry hippos: Towards language modeling with state space models. *arXiv preprint arXiv:2212.14052* (2022).
- [27] Fu, D. Y. *et al.* Hungry hungry hippos: Towards language modeling with state space models (2023). URL <https://arxiv.org/abs/2212.14052>. 2212.14052.
- [28] Glorioso, P. *et al.* Zamba: A compact 7b ssm hybrid model (2024). URL <https://arxiv.org/abs/2405.16712>. 2405.16712.
- [29] Ren, L. *et al.* Samba: Simple hybrid state space models for efficient unlimited context language modeling (2024). URL <https://arxiv.org/abs/2406.07522>. 2406.07522.
- [30] Zhang, B. & Sennrich, R. Root mean square layer normalization. *Advances in Neural Information Processing Systems* **32** (2019).
- [31] Su, J. *et al.* Roformer: Enhanced transformer with rotary position embedding. *Neurocomputing* **568**, 127063 (2024).

- [32] Springer, J. M., Kotha, S., Fried, D., Neubig, G. & Raghunathan, A. Repetition improves language model embeddings (2024). URL <https://arxiv.org/abs/2402.15449>. 2402.15449.
- [33] Marin, F. I., Teufel, F. *et al.* Bend: Benchmarking dna language models on biologically meaningful tasks. *arXiv preprint arXiv:2306.15006* (2024).
- [34] Poli, M. *et al.* Genomics long-range benchmark (lrb): Evaluating long-context models on genomic data. *arXiv preprint arXiv:2306.00971* (2023).
- [35] Dao, T. Flashattention-2: Faster attention with better parallelism and work partitioning. *arXiv preprint arXiv:2307.08691* (2023).
- [36] Le, N. Q. K., Ho, Q.-T., Nguyen, V.-N. & Chang, J.-S. Bert-promoter: An improved sequence-based predictor of dna promoter using bert pre-trained model and shap feature selection. *Computational Biology and Chemistry* **99**, 107732 (2022).
- [37] Zhang, P., Zhang, H. & Wu, H. ipro-wael: a comprehensive and robust framework for identifying promoters in multiple species. *Nucleic Acids Research* **50**, 10278–10289 (2022).
- [38] Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nature methods* **18**, 1196–1203 (2021).
- [39] Jin, J. *et al.* idna-abf: multi-scale deep biological language learning model for the interpretable prediction of dna methylations. *Genome biology* **23**, 219 (2022).
- [40] Lee, D., Yang, J. & Kim, S. Learning the histone codes with large genomic windows and three-dimensional chromatin interactions using transformer. *Nature Communications* **13**, 6678 (2022).
- [41] Chen, K., Zhao, H. & Yang, Y. Capturing large genomic contexts for accurately predicting enhancer-promoter interactions. *Briefings in Bioinformatics* **23**, bbab577 (2022).
- [42] Ni, Y. *et al.* Epi-mind: identifying enhancer–promoter interactions based on transformer mechanism. *Interdisciplinary Sciences: Computational Life Sciences* **14**, 786–794 (2022).
- [43] Wang, Z. *et al.* Towards a better understanding of tf-dna binding prediction from genomic features. *Computers in Biology and Medicine* **149**, 105993 (2022).
- [44] Rozowsky, J. *et al.* The en-tex resource of multi-tissue personal epigenomes & variant-impact models. *Cell* **186**, 1493–1511 (2023).
- [45] Theodoris, C. V. *et al.* Transfer learning enables predictions in network biology. *Nature* **618**, 616–624 (2023).

- [46] Ji, Y., Zhou, Z., Liu, H. & Davuluri, R. V. Dnabert: a comprehensive predictor for dna sequences based on deep transfer learning. *Bioinformatics* **37**, 4776–4783 (2021).
- [47] Schiff, Y., Kao, C.-H., Gokaslan, A. *et al.* Caduceus: Bi-directional equivariant long-range dna sequence modeling. *arXiv preprint arXiv:2306.15006* (2023).
- [48] Zhou, Z. *et al.* Genomeocean: An efficient genome foundation model trained on large-scale metagenomic assemblies. *bioRxiv* (2025). URL <https://www.biorxiv.org/content/early/2025/02/05/2025.01.30.635558>.
- [49] Wu, W. *et al.* Generator: A long-context generative genomic foundation model (2025). [2502.07272](#).
- [50] Nguyen, C. V. *et al.* Taipan: Efficient and expressive state space language models with selective attention (2024). URL <https://arxiv.org/abs/2410.18572>. [2410.18572](#).
- [51] Kingma, D. P. & Ba, J. Bengio, Y. & LeCun, Y. (eds) *Adam: A method for stochastic optimization*. (eds Bengio, Y. & LeCun, Y.) *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings* (2015). URL <http://arxiv.org/abs/1412.6980>.
- [52] Trop, E. *et al.* The genomics long-range benchmark: Advancing DNA language models (2025). URL <https://openreview.net/forum?id=8O9HLDrmtq>.