

THE HUMAN GENOMICS LONG-RANGE BENCHMARK: ADVANCING DNA LANGUAGE MODELS

Anonymous authors

Paper under double-blind review

ABSTRACT

The advent of language models (LMs) in genomics necessitates benchmarks that can assess models' capabilities and limitations. In contrast to protein models, DNA LMs can be used to study non-coding regions of the genome and must account for unique challenges, especially interactions across long sequence lengths. However, existing benchmarks for DNA LMs are defined over short sequence datasets and can involve tasks that are often not considered to be biologically meaningful. Here, we present the **Human** Genomics Long-Range Benchmark (LRB), which focuses on biologically meaningful tasks and supports long-range contexts. We complement our benchmark with fine-tuning recipes that meaningfully improve performance and affect model evaluation. We evaluate DNA LMs across nine compiled **human** genome tasks and observe that DNA LMs achieve competitive performance relative to supervised baselines on several tasks (e.g., genome annotation), but there remains a significant gap in domains, such as variant effect and gene expression prediction. Additionally, we introduce a visualization tool to examine model performance split by various genomic properties. Lastly, we present methods for context-length extrapolation of transformer-based models that enable studying the effect of context length on DNA LM performance.

1 INTRODUCTION

Pre-training models on a large corpus of unlabeled data and subsequently fine-tuning to solve downstream tasks has demonstrated widespread success across domains, such as natural language processing (Achiam et al., 2023; Team et al., 2023) and computer vision (Oquab et al., 2023; Radford et al., 2021). More recently this paradigm has shown promise in biological applications, enabled by the wealth of unlabeled data coming from next-generation sequencing technologies. A prominent example are protein language models (LMs), which have been used to predict the effects of coding mutations on protein function (Lin et al., 2022), generate viable protein sequences conditioned on functional properties (Madani et al., 2023), and accurately predict protein structure from amino acid sequences (Lin et al., 2023). The development of these models has been made possible by benchmarks, such as CASP (Kryshtafovych et al., 2021), TAPE (Rao et al., 2019), PEER (Xu et al., 2022), and ProteinGym (Notin et al., 2023).

Genomics represents a potential new frontier for LMs in biology. The common pre-training tasks in language modeling (i.e., filling in missing tokens based on input context) inherently train LMs to model evolutionary forces, such as conservation and co-evolution, and the statistical patterns that these models learn can map to genomic motif identification, which is useful in accurate gene annotation. Indeed, significant progress has been made, with various LMs tailored to DNA sequences (Benegas et al., 2023a;b; Dalla-Torre et al., 2023; Ji et al., 2021; Nguyen et al., 2023; 2024; Schiff et al., 2024; Zhou et al., 2023). However, modeling genomic data presents unique challenges compared to proteomics. When modeling DNA, we have to account for non-coding regions and contend with interactions that can be orders of magnitude larger (Furlong & Levine, 2018). To guide the principled development of new DNA LMs, there is a need for robust benchmarks that accurately reflect these nuances. While several benchmarks have been proposed, these existing works contain important limitations. The vast majority of tasks proposed across existing benchmarks only consider short input contexts (less than 2k base pairs) (Dalla-Torre et al., 2023; Grešová et al., 2023; Marin et al., 2023; Zhou et al., 2023), disregarding long-range interactions that are highly impactful in genomics. Additionally, tasks in some benchmarks may be overly simplistic, failing to reflect real-world use

054 cases, e.g., some benchmarks have used synthetic data to construct negative sets (Dalla-Torre et al.,
055 2023).

056 To bridge these gaps, we propose the **Human Genomics Long-Range Benchmark (LRB)**, a compilation
057 of biologically meaningful tasks in human genomics. Our benchmark deliberately incorporates tasks
058 hypothesized to span both short and long genomic contexts. Allowing users to select arbitrary
059 sequence length inputs for any given dataset enables us for the first time to understand empirically
060 the importance of long-range inputs for our proposed tasks. We also include available genomic
061 annotations and provide a visualization tool that allows users to analyze results in more detail. We
062 demonstrate the benefit of full model fine-tuning compared to previous approaches that keep backbone
063 DNA LM weights frozen during downstream training. Finally, we introduce methods for extending
064 the context size of existing DNA LMs, which allows us quantify the benefits of long-range context on
065 DNA LM performance. To summarize, we make the following contributions:

066 **1. Release the Genomics Long-Range Benchmark**, composed of biologically meaningful tasks
067 that cover both short- and long-range genomic scales. We provide evaluation results for a selection
068 of prominent DNA LMs in both zero-shot and fine-tuning settings along with comparisons against
069 reference baselines. We find that on genomic annotation tasks DNA LMs perform competitively with
070 existing supervised models, but on the long-range prediction tasks of gene expression and zero-shot
071 mutation effect prediction there persists a large gap.

072 **2. Develop and analyze improved fine-tuning methods** that better reflect real-world usage in
073 downstream tasks, finding that full model weight fine-tuning significantly improves performance.

074 **3. Introduce an analysis and visualization tool** to examine models' performance across different
075 genomic properties. This tool enables deeper analyses that reveal more nuanced evidence that DNA
076 LMs lag behind a well-regarded and long-range supervised baseline, Enformer (Avsec et al., 2021a),
077 in modeling long-range interactions.

078 **4. Conduct context-length extension for the Nucleotide Transformer LM** to probe the impact of
079 increasing context length on performance on our benchmark.
080

081 082 083 2 BACKGROUND

084 085 2.1 LANGUAGE MODELING FOR DNA

086
087 Supervised machine learning methods have been successfully applied to genomics (Alipanahi et al.,
088 2015; Avsec et al., 2021a; de Almeida et al., 2022; Zhou & Troyanskaya, 2015; Zhou et al., 2018b).
089 However, these models depend on large amounts of labeled data and tend to be task-specific. LMs
090 have recently gained traction in the genomics domain: the abundance of unlabeled sequences supports
091 robust model pre-training and the widely-used pre-training objectives of next token prediction (NTP)
092 or masked language modeling (MLM) directly lend themselves to models identifying genomic motifs
093 and evolutionary patterns, e.g., conservation. Some notable recent works include DNABERT (Ji
094 et al., 2021; Zhou et al., 2023; 2024), GPN (Benegas et al., 2023a;b), Nucleotide Transformer (NT)
095 (Dalla-Torre et al., 2023), GENA-LM (Fishman et al., 2023), HyenaDNA (Nguyen et al., 2023), Evo
096 (Nguyen et al., 2024) and Caduceus (Schiff et al., 2024). A more thorough review of recent DNA
097 LMs is deferred to Appendix A.2.

098 099 2.2 DNA LM EVALUATION

100
101 The goal of DNA LMs is to learn meaningful representations that can be used to improve perfor-
102 mance on downstream tasks. Existing DNA benchmarks, which include the Nucleotide Transformer
103 tasks (NT; Dalla-Torre et al. (2023)), Genomic Benchmark (GB; Grešová et al. (2023)), Genome
104 Understanding Evaluation (GUE; Zhou et al. (2023)), and Benchmark for DNA LMs (BEND; Marin
105 et al. (2023)), have been crucial for establishing baseline model capabilities. (see Appendix A.3 for a
106 more complete description of existing works). However, these benchmarks contain several important
107 shortcomings: they do not focus on long-range sequences, they can contain synthetic examples, and
their evaluations do not take full advantage of pre-trained models.

3 THE GENOMICS LONG-RANGE BENCHMARK

Below we describe the nine tasks that we compiled from various human genome data sources that comprise our proposed Genomics Long-Range Benchmark (LRB). Our suite consists of tasks that are hypothesized to require only short-range contexts as well as those thought to need longer sequences for accurate prediction.

By enabling users to download data at arbitrary length scales (the first benchmark to support this feature), these hypotheses can be rigorously tested. Our tasks span various applications that are of interest to practitioners, namely variant effect prediction, gene expression prediction, regulatory element detection, and chromatin factor identification; see Table 2. Below, for each task, we provide details on the biological relevance that motivated its inclusion, a formal task definition, and rationale for hypothesized long-range dependencies (where applicable). We defer additional details, e.g., data source and processing, train / test splits, and metric definition, to Appendix B.

Table 1: Comparison to existing benchmarks.

	Long range	Human centric	Biologically meaningful
NT (Dalla-Torre et al., 2023)	✗	✗	✗
GB (Grešová et al., 2023)	✗	✗	✗
GUE (Zhou et al., 2023)	✗	✗	✓
BEND (Marin et al., 2023)	✗	✓	✓
Human Genomics LRB	✓	✓	✓

Table 2: Overview of the tasks contained in the Genomics Long-Range Benchmark.

	Type	# Outputs	# Train / Test	% Pos. Label
<i>Variant Effect Prediction</i>				
Causal eQTL	SNP Classification	1	89k / 9k	50.0
Pathogenic OMIM	SNP Classification	1	- / 2.3M	0.02
Pathogenic ClinVar	SNP Classification	1	39k / 1k	56.1
<i>Gene Expression Prediction</i>				
Bulk RNA-seq	Seq-wise Regression	218	23k / 1k	-
CAGE profile	Binned Regression	50 / bin	34k / 2k	-
<i>Regulatory Element Detection</i>				
Promoter	Seq-wise Classification	1	953k / 96k	4.7
Enhancer	Seq-wise Classification	1	1.9M / 192k	52.5
<i>Chromatin Feature Identification</i>				
Histone Mark Prediction	Seq-wise Classification	20	2.2M / 227k	7.0
Chromatin Accessibility	Seq-wise Classification	20	2.2M / 227k	4.4

3.1 VARIANT EFFECT PREDICTION

3.1.1 CAUSAL EQTL

Biological Relevance Predicting the effects of genetic variants, particularly expression quantitative trait loci (eQTLs), is essential for understanding the molecular basis of several diseases. eQTLs are genomic loci that are associated with variations in mRNA expression levels among individuals. By linking genetic variants to causal changes in mRNA expression, researchers can uncover how certain variants contribute to disease development (Consortium, 2020).

Task Definition The task is formulated as a binary classification problem to distinguish eQTLs from GTEx (Consortium, 2020) from a set of matched negatives identified in Avsec et al. (2021a). Inputs are sequences centered around candidate single nucleotide polymorphisms (SNPs) each assigned a causal probability by fine-mapping using the “Sum of Single Effects” (SuSiE) model (Wang et al., 2020). Following Avsec et al. (2021a), variants with causal probability greater than 0.9 are labeled as positive and variants with causal probability less than 0.01 are labeled as negative.

Long-Range The regulation of gene expression is modulated by distal, cis-regulatory elements, called enhancers, that can be more than several hundred thousand base pairs (bps) away from a target gene (Furlong & Levine, 2018). Variants that impact gene expression are often located at such distal

elements, and thus, to predict such variants, models should have long context windows (Avsec et al., 2021a).

3.1.2 PATHOGENIC OMIM

Biological relevance Predicting the effects of regulatory variants on pathogenicity is crucial for understanding disease mechanisms (Marwaha et al., 2022). Elements that regulate gene expression are often located in non-coding regions, and variants in these areas can disrupt normal cellular function, leading to disease. Accurate predictions can identify biomarkers and therapeutic targets, enhancing personalized medicine and genetic risk assessment.

Task Definition The task is formulated as a binary classification problem where inputs are DNA sequences centered around a SNP and outputs are binary labels. The dataset was constructed following Benegas et al. (2023a), where the negative class corresponds to a common (mean allele frequency $> 5\%$) SNP in gnomAD (Chen et al., 2022) and the positive class corresponds to a pathogenic SNP, defined as a SNP in a regulatory region having an implication in a Mendelian disorder in the Online Mendelian Inheritance in Man database (Smedley et al., 2016).

Long-Range Regulatory elements like enhancers and silencers can exist far from the genes they regulate (Furlong & Levine, 2018). Variants in these regulatory elements can lead to aberrant gene expression patterns and ultimately disease, but identifying such regulatory variants is difficult since regulatory elements can modulate the expression of proximal or distal genes. Models that can capture interactions between possibly distal regulatory elements and their target genes while still being able to capture the proximal interactions are essential to identifying non-coding pathogenic variants.

3.1.3 PATHOGENIC CLINVAR

Biological Relevance A coding variant refers to a genetic alteration that occurs within the protein-coding regions of the genome, also known as exons. Such alterations can impact protein structure, function, stability, and interactions with other molecules, ultimately influencing cellular processes and potentially contributing to the development of genetic diseases (Lek et al., 2016). Predicting variant pathogenicity is crucial for guiding research into disease mechanisms and personalized treatment strategies, enhancing our ability to understand and manage genetic disorders effectively.

Task Definition This task is formulated as a binary classification problem where inputs are sequences centered around SNPs. The dataset was constructed following Benegas et al. (2023a), where the negative class corresponds to a common (minor allele frequency $> 5\%$) SNP in gnomAD (Chen et al., 2022) and the positive class to pathogenic SNPs identified in ClinVar (Landrum et al., 2020).

3.2 GENE EXPRESSION PREDICTION

3.2.1 BULK RNA-SEQ

Biological Relevance Gene expression involves the process by which information encoded in a gene directs the synthesis of a functional gene product, typically a protein, through transcription and translation. Transcriptional regulation determines the amount of mRNA produced, which is then translated into proteins. Developing a model that can predict RNA expression levels solely from sequence data is crucial for advancing our understanding of gene regulation, elucidating disease mechanisms, and identifying functional sequence variants.

Task Definition This task is described as a multi-variable, sequence-wise regression task. Data was constructed following Zhou et al. (2018a) such that inputs are DNA sequences centered around the transcription start site (TSS) of each gene where the TSS was identified using a combination of annotations from GENCODE (Harrow et al., 2012) and CAGE data from FANTOM5 (Forrest et al., 2014). Outputs are RPKM normalized RNA expression counts for each gene obtained from Consortium (2020) that were $\log(1 + x)$ normalized and standardized. For each gene, there are 218 different counts corresponding to the RNA expression level in different tissue types.

Long-Range RNA gene expression is regulated by non-coding elements, such as enhancers and silencers, which can be located hundreds of kilo-bps away from the gene (Furlong & Levine, 2018), indicating the possible presence of long-range interactions in transcription regulation.

3.2.2 CAP ANALYSIS GENE EXPRESSION (CAGE) PROFILE

Biological Relevance CAGE provides accurate high-throughput measurements of RNA expression by mapping TSSs at a nucleotide-level resolution (Takahashi et al., 2012). This is vital for detailed mapping of TSSs, understanding gene regulation mechanisms, and obtaining quantitative expression data to study gene activity comprehensively.

Task Definition This task is described as a multi-variable, binned nucleotide-wise regression task. The data was constructed following the approach outlined in Basenji (Kelley, 2020). Inputs are DNA sequences and the outputs are $\log(1 + x)$ normalized CAGE expression counts from FANTOM5 (Forrest et al., 2014) given for each 128 bp bin of the input sequence. For each bin in a sequence, there are 50 different values corresponding to expression amounts across 50 human cell / tissue types.

Long-Range The production of RNA via transcription as measured by CAGE is regulated by non-coding elements that can be located hundreds of kilo-bps away from the gene, indicating the presence of long-range interactions in transcription regulation (Furlong & Levine, 2018).

3.3 CIS-REGULATORY ELEMENT DETECTION

Biological Relevance Cis-regulatory elements, such as promoters and enhancers, control the spatial and temporal expression of genes (Andersson & Sandelin, 2020). These elements are essential for understanding gene regulation mechanisms and how genetic variations can lead to differences in gene expression.

Task Definition This task is described as a binary classification problem. Data from Search Candidate Regulatory Elements by ENCODE (SCREEN (The ENCODE Project Consortium, 2020)) was processed according to our approach outlined in Appendix B.3. Inputs are sequences sampled from across the entire human genome and outputs are binary values, where a positive label is assigned to a sequence if the center 200 bps of the input sequence overlap by at least 50% with an annotated enhancer or promoter. This task is composed of two sub-tasks: (1) predicting the presence of promoters and (2) predicting the presence of enhancers.

3.4 CHROMATIN FEATURE IDENTIFICATION

Biological Relevance Predicting chromatin features, such as histone marks and DNA accessibility, is crucial for understanding gene regulation, as these features indicate chromatin state and are essential for transcription activation (Zhou et al., 2018b).

Task Definition This task is a multi-label binary classification problem constructed following Zhou & Troyanskaya (2015), where sequences were sampled from the human genome as inputs and outputs correspond to binary labels for different chromatin profiles. The task contains two sub-tasks: one for predicting histone marks and another for predicting chromatin accessibility. For histone marks, each of the 20 binary values represents a different histone mark in a specific cell type. For DNA accessibility, each of the 20 binary values corresponds to a different tissue/cell type. A value is labeled as positive if the center 200 bps of the input sequence overlaps by at least 50% with a peak region measured by ChIP-seq (histone marks) or DNase-seq (DNA accessibility) obtained from ENCODE and the Roadmap Epigenomics consortium (Bernstein et al., 2010; The ENCODE Project Consortium, 2020).

3.5 IMPROVED EVALUATION WITH FULL FINE-TUNING

To evaluate DNA LMs we perform fine-tuning, i.e., we train a model in a supervised manner on a downstream task. Our fine-tuning strategy involves extracting embeddings from each model which are then input to a task-specific prediction head (see Appendix D for details). In previous benchmarks, authors fine-tuned models by freezing the embeddings (Marin et al., 2023). We perform a systematic study of fine-tuning strategies and discover that this strategy significantly hurts DNA LM performance. We therefore provide a recipe for full-parameter fine-tuning and show that it significantly improves performance across many tasks, enabling us to evaluate models more fairly than in previous works and setting new best-practices for DNA LMs (independent of our benchmark).

3.6 ADDITIONAL NOVEL FEATURES OF THE LRB

In addition to our careful curation of tasks and improved fine-tuning methodology, we highlight two more novel aspects of the LRB.

Visualization Tool We provide benchmark users with a visualization tool in the form of an interactive `jupyter` (Kluyver et al., 2016) notebook. To create this tool we collected additional genomic annotation datasets from SCREEN, GENCODE, RepeatMasker (Harrow et al., 2012; Smit et al., 2015; The ENCODE Project Consortium, 2020) and aligned them to our benchmark task datasets; see Appendix B.5 for details and screenshots. Our tool enables a deeper level of analysis compared to what other benchmarks afford. For example, users can view models’ performance in aggregate, by specific annotations, and also by distance to TSSs.

Arbitrary Sequence Length Our benchmark allows users to download arbitrary sequence lengths for any given tasks. This enables the probing of the effect of sequence length and lets users evaluate their LMs on the same context size on which they performed pre-training, mitigating any confounding from sequence length generalization effects.

3.7 SELECTED BASELINES

To contextualize the performance of DNA LMs, we curate a set of task-specific expert methods that are comprised of well-regarded supervised models.

Combined annotation dependent depletion (CADD) (Schubach et al., 2024) is a SVM developed for detecting deleterious DNA variants trained on predicted neutral variants and simulated deleterious variants. We use this method as an expert baseline for our zero-shot variant effect prediction tasks.

GPN-MSA Benegas et al. (2023a) present an alignment-based DNA LM for variant effect prediction based on the RoFormer (Su et al., 2021) architecture. In addition to the standard input DNA sequence, a Multiple Sequence Alignment (MSA), an alignment of similar sequences from multiple species, is used as an auxiliary input. This alignment is computed from 89 vertebrates and is always unmasked at all times, giving the model access to evolutionary information computed for a given input sequence. The auxiliary alignment information and strong performance on zero-shot prediction, render GPN-MSA a useful watermark against which to compare pure sequence-based DNA LMs.

Enformer (Avsec et al., 2021a) is composed of both convolutional and transformer layers and trained in a supervised multi-task manner on various biological tasks using a context length of up to 196k bps. We use Enformer as the expert method for fine-tuning versions of variant effect prediction, gene expression prediction, and regulatory element detection tasks.

DeepSEA (Zhou et al., 2018b) is a convolutional network trained to predict chromatin profile data, such as transcription factor binding, histone marks, and DNA accessibility. As our chromatin feature tasks are derived from DeepSEA, we use it as the expert method for these tasks.

Supervised Models We also train models *from scratch* in a supervised manner for each task. This paradigm of supervised training is currently more prevalent in machine learning applications for genomics, and these baselines help us to better contextualize DNA LM performance. We select two architectures for the supervised baseline. First, we train a **convolutional neural network (CNN)**, which is inspired by the one used in Benegas et al. (2023b), but without dilation. Additionally, we train a **Caduceus** (Schiff et al., 2024) model from scratch. See Appendix D.6 for more precise supervised model details.

4 CONTEXT LENGTH EXTENSION

Motivated by the long-range sequences present in the LRB, we explore methods for extending the context size of existing models. To that end, we focus on the Nucleotide Transformer model (NTv2; Dalla-Torre et al. (2023)), which originally has a context size of 12k bp and uses rotary positional embeddings (RoPE; Su et al. (2021)). However, processing longer sequences with LMs like NTv2, which use the transformer architecture (Vaswani et al., 2017), faces two main challenges. First, transformers rely on the attention mechanism, which scales quadratically in sequence length. Second, LMs struggle with generalizing to sequence lengths beyond those seen during pre-training, known

as length extrapolation (Anil et al., 2022; Dubois et al., 2019; Kazemnejad et al., 2023; Press et al., 2021).

Methodology To address the compute constraints, we use a memory-efficient attention implementation, computing attention scores sequentially and in chunks of \sqrt{L} , reducing memory usage from $\mathcal{O}(L^2)$ to $\mathcal{O}(\sqrt{L})$, where L denotes sequence length (Rabe & Staats, 2021). To solve the length generalization issue, we apply the ‘NTK-aware’ method presented in Peng et al. (2023). This method re-scales the frequencies in RoPE embeddings to handle longer sequences by converting length extrapolation into *interpolation*. For more details on these approaches, see Appendix C.

5 RESULTS

5.1 EXPERIMENTAL SETUP

We evaluate several prominent DNA LMs on our benchmark: the Nucleotide Transformer v2 (NTv2) series (Dalla-Torre et al., 2023), DNABERT-1 (Ji et al., 2021), DNABERT-2 (Zhou et al., 2023), the HyenaDNA series (Nguyen et al., 2023), and Caduceus (Schiff et al., 2024) representing a range of pre-training datasets and objectives, architectures, and model sizes. For fine-tuning, we use an MLP as the prediction head and train both the DNA LM and MLP weights (see Appendix D for full details).

For classification tasks with highly imbalanced labels (see Table 2), we use area under precision-recall curve (AUPRC) as opposed to receiver operator curve (AUROC) as the metric.

Fine-tuning Models are trained using either mean-squared error loss for regression tasks or cross-entropy loss for classification tasks. For each task, we perform five-fold cross-validation (CV) using different random seeds, where we create different train / validation splits, select the best-performing model using early stopping on validation loss, and evaluate it on the held-out test set. We report the mean \pm standard deviation performance across folds as final metrics.

Zero-Shot Prediction We also evaluate the zero-shot performance on our three variant effect prediction tasks to account for the fact that, in practice, determining pathogenicity or causality of variants is difficult, which often results in smaller datasets not suitable for fine-tuning. Given the extreme class imbalance in the Pathogenic OMIM dataset, we only perform zero-shot evaluation for this task and do not report fine-tuning results.

5.2 MAIN DNA LM RESULTS

In Tables 3 and 4, we present the top performing DNA LMs (full results in Appendix E).

Variant Effect Prediction For zero-shot evaluation, we observe that DNA LMs are outperformed by the CADD and GPN-MSA baselines on all variant effect prediction tasks. Additionally, for zero-shot Causal eQTL, we find that all models struggle, with near-random performance. Predicting pathogenicity, is the clearest example where DNA LMs fall short of CADD, which has nearly 2x better performance in ClinVar and about 100x in OMIM. When fine-tuning, we find that DNA LM performance on both variant tasks greatly improves, matching or surpassing the strong Enformer baseline. **Importantly, the alignment-based GPN-MSA model, despite using short context inputs (128 bps), outperforms CADD and all single-sequence DNA LMs, highlighting the importance of capturing conservation in predicting variant effects. For DNA LMs to be useful for these tasks, they must also find a way to model and learn evolutionary pressures and conservation.**

Gene Expression Prediction While NTv2 is the best performing DNA LM for Bulk RNA and CAGE tasks, the baseline Enformer outperforms LMs by a wide margin.

Regulatory Element Detection DNA LMs are able to accurately predict the presence of regulatory elements, especially considering the class-imbalance present in promoter detection, with NTv2 performing best among DNA LMs. However, there remains a gap to the supervised Enformer model.

Chromatin Feature Identification For both histone mark and DNA accessibility, NTv2 is the best performing DNA LM, even exceeding the supervised baseline on the former task, and demonstrating significantly better performance than the other DNA LMs.

Table 3: Benchmarking performance of DNA LMs and baselines on variant effect prediction tasks. Models were evaluated using both fine-tuning and zero-shot. Best LM values are **bolded**. *Extended NTv2 was fine-tuned with 60k bp sequences due to compute constraints.

Model Name	Context (bps)	Causal eQTL (AUROC)		Pathogenic ClinVar (AUROC)		Pathogenic OMIM (AUPRC)
		Fine-tune	Zero-shot	Fine-tune	Zero-shot	Zero-shot
<i>DNA LMs</i>						
DNABERT-2	10k	0.72 ± 0.008	0.50	0.74 ± 0.013	0.50	0.002
DNABERT-S	10k	0.73 ± 0.008	-	0.73 ± 0.011	-	-
NTv2 500M	12k	0.72 ± 0.003	0.51	0.78 ± 0.009	0.68	0.003
<i>Extended</i>	96k*	0.74 ± 0.004	0.51	0.75 ± 0.018	0.53	0.002
HyenaDNA 160K	160k	0.71 ± 0.010	0.51	0.56 ± 0.073	0.49	0.002
Caduceus 131K	131k	0.68	0.49	-	0.53	-
<i>Alignment-based LM</i>						
GPN-MSA	128	-	0.55	-	0.97	0.35
<i>Supervised Training</i>						
CNN	2k	0.71	-	0.61	-	-
Caduceus (from scratch)	2k	0.67	-	0.61	-	-
<i>Baseline</i>		0.76 ± 0.002 (Enformer)	0.56 (CADD)	-	0.97 (CADD)	0.253 (CADD)

Table 4: Benchmarking performance of DNA LMs and baselines on gene expression, regulatory element, and chromatin features tasks. Models were evaluated in only a fine-tuned setting for this set of tasks. Best LM values are **bolded** and in **green** if LM beats baseline.

Model Name	Context (bps)	Bulk RNA (R^2)	CAGE (R^2)	Promoter (AUPRC)	Enhancer (AUROC)	Histone Marks (AUPRC)	DNA Accessibility (AUPRC)
<i>DNA LMs</i>							
DNABERT-2	10k	0.51 ± 0.050	-	0.71 ± 0.112	0.81 ± 0.022	0.24 ± 0.091	0.15 ± 0.064
DNABERT-S	10k	0.52 ± 0.060	-	0.75 ± 0.021	0.83 ± 0.005	0.33 ± 0.006	0.16 ± 0.039
NTv2 500M	12k	0.60 ± 0.038	0.39 ± 0.011	0.79 ± 0.006	0.82 ± 0.002	0.38 ± 0.003	0.3 ± 0.007
<i>Extended</i>	96k	0.56 ± 0.037	0.36 ± 0.011	0.78 ± 0.003	0.82 ± 0.005	0.38 ± 0.004	0.3 ± 0.006
HyenaDNA 160K	160k	0.46 ± 0.006	0.19 ± 0.032	0.67 ± 0.009	0.74 ± 0.009	0.25 ± 0.004	0.11 ± 0.002
Caduceus 131K	131k	0.52	-	0.75	-	-	-
<i>Supervised Training</i>							
CNN	2k	0.47	0.05	0.84	0.81	0.11	0.10
Caduceus (from scratch)	2k	0.51	0.09	-	0.85	0.14	0.10
<i>Baseline</i>		0.80 ± 0.010 (Enformer)	0.49 ± 0.000 (Enformer)	0.86 ± 0.006 (Enformer)	0.92 ± 0.002 (Enformer)	0.35 (DeepSea)	0.44 (DeepSea)

5.3 IMPORTANCE OF CONTEXT LENGTH FOR LONG-RANGE TASKS

Table 5: Varying context lengths on long-range tasks.

Context (bp)	Causal eQTL (AUCROC)	Bulk RNA (R^2)	CAGE (R^2)
<i>CNN</i>			
2k	0.71	0.47	0.051
32k	0.70	0.46	0.091
65k	-	0.47	0.120
<i>Caduceus (from scratch)</i>			
2k	0.67	0.51	0.086
32k	-	0.54	0.079
65k	-	0.54	0.100

To verify our hypothesis that the long-range tasks in our benchmark will benefit from models with larger input context sizes, we perform the following ablation analysis. Using the two supervised training baselines (CNN and Caduceus models trained from scratch), we train these models on the long-range tasks with increasing context sizes: 2k, 32k, and 65k bps. In Table 5, we see a positive association between input context length and performance across both architectures. These findings validate our characterization of these tasks as ‘long-range.’

5.4 ANALYZING RESULTS BY GENOMIC ANNOTATIONS

432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485

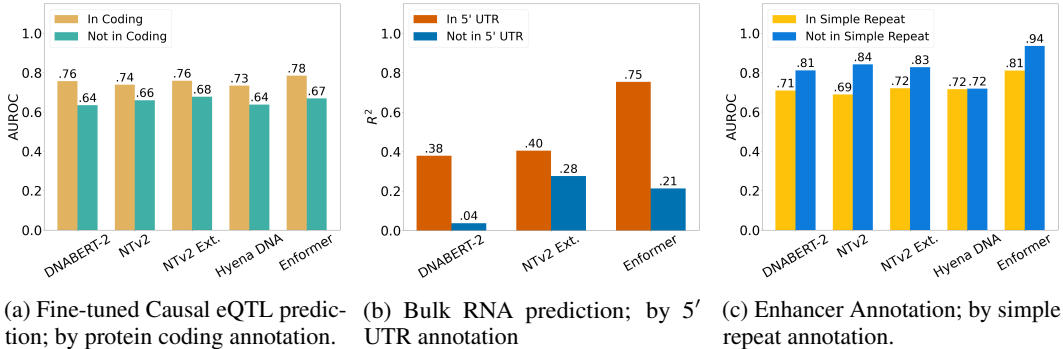


Figure 1: Results split by genomic annotations.

We developed an analysis and visualization tool to examine models performance across different genomic properties and annotations.

Using our tool we are able to perform deeper analyses and extract insights about the performance of each model, which are inaccessible to users of existing benchmarks. We detail some examples in Figure 1.

Causal eQTL Prediction (Fine-tune) By stratifying SNPs into protein-coding and non-coding regions in Figure 1a, we find a potential failure mode for both DNA LMs and supervised models. Non-coding variants presumably entail regulatory and possibly longer-range interactions, and all models perform worse in these regions.

Bulk RNA Expression Prediction In Figure 1b, we see that the performance of DNA LMs and Enformer drops precipitously when focusing on non-5' regions that likely entail longer-range interactions. However, we also observe that the context-extended NTV2 outperforms Enformer on this region, implying that the majority of the performance gap between DNA LMs and the Enformer baseline lies in modeling variants in the 5' regions.

Enhancer Detection In Figure 1c, we observe that most models, including Enformer, suffer a performance hit when identifying enhancers within simple repeat regions, likely due to the difficulty of detecting enhancers within repetitive regions of the genome.

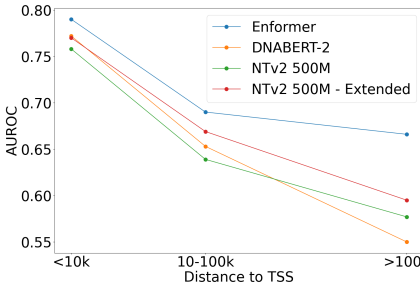


Figure 2: Fine-tuned Causal eQTL variant task; by distance to nearest TSS.

5.5 LENGTH EXTENSION

To create the context extended model, we conduct additional training (~5B tokens) on the pre-training dataset using the methodology described in Section 4 (and in Appendix D.5). For certain long-range tasks, the additional context extension pre-training improves performance. For example, for Causal eQTL prediction (with fine-tuning) in Figure 2 we see that the context extended NTV2 has the best DNA LM performance and that this trend is more pronounced when stratifying by SNP distance to TSS.

5.6 EFFECT OF FINE-TUNING METHODOLOGY

In Table 6, we demonstrate the importance of our proposed fine-tuning. For two of the DNA LMs (see additional results in Appendix E.3), we show how full fine-tuning, as opposed to freezing LM weights and only training a prediction head, a common practice in existing benchmarks such as BEND (Marin et al., 2023), drastically improves model performance almost uniformly across tasks. We also believe our methodology is more in line with how practitioners would use DNA LMs in real-world settings.

Table 6: Difference in performance of DNA LM fine-tuning strategies. Percent increase in performance of full fine-tuning vs. freezing LM weights and only training prediction heads.

	Causal eQTL (AUCROC)	Pathogenic ClinVar (AUROC)	Bulk RNA (R^2)	CAGE (R^2)	Promoter (AUPRC)	Enhancer (AUROC)	Histone Marks (AUCPRC)	DNA Accessibility (AUPRC)
NTv2 500M	+0.49	+4.27	+18.29	+42.14	-1.45	+0.90	+22.46	+47.96
HyenaDNA 32K	+0.35	+11.58	+82.46	+102.91	-18.21	-6.02	+14.43	-22.67

6 DISCUSSION AND CONCLUSION

In this work, we introduced the **Human** Genomics LRB. Our benchmark is the first to truly evaluate long-range capabilities. We provided initial results for several prominent DNA LMs, with more in-depth analysis than previous benchmarks explored. Our results demonstrate the importance of fully fine-tuning models. Additionally, we identify several domains where a large performance gap needs to be bridged before DNA LMs can be reliably used and some failure modes of DNA LMs. Namely, zero-shot DNA LM variant effect prediction is not yet mature enough to replace widely-used tools, such as CADD or alignment-based models like GPN-MSA. Similarly, for gene expression prediction, DNA LMs lag far behind supervised methods. In contrast, for annotation tasks, DNA LMs already demonstrate competitive performance relative to proven methods. These results demonstrate that future DNA LM efforts should focus on the more difficult tasks that entail long-range interactions, and we hope that our benchmark spurs such development.

Future Work One potential limitation of our work is the lack of hyperparameter search for fine-tuning; a more extensive search would better differentiate models. Another limitation is the lack of experimentally verified enhancer-gene pairings, which would allow for a more complete examination of the long-range capabilities of models. In future iterations of our benchmark, we also plan to add more tissue-specific analyses, bp-level annotation tasks, and tasks covering multiple species.

REFERENCES

- Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. URL <https://www.tensorflow.org/>. Software available from tensorflow.org.
- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.
- Babak Alipanahi, Andrew Delong, Matthew T Weirauch, and Brendan J Frey. Predicting the sequence specificities of dna-and rna-binding proteins by deep learning. *Nature biotechnology*, 33(8):831–838, 2015.
- Robin Andersson and Albin Sandelin. Determinants of enhancer and promoter activities of regulatory elements. *Nature Reviews Genetics*, 21:71–87, 2020.
- Cem Anil, Yuhuai Wu, Anders Andreassen, Aitor Lewkowycz, Vedant Misra, Vinay Ramasesh, Ambrose Slone, Guy Gur-Ari, Ethan Dyer, and Behnam Neyshabur. Exploring length generalization in large language models. *Advances in Neural Information Processing Systems*, 35:38546–38556, 2022.
- Žiga Avsec, Vikram Agarwal, Daniel Visentin, Joseph R Ledsam, Agnieszka Grabska-Barwinska, Kyle R Taylor, Yannis Assael, John Jumper, Pushmeet Kohli, and David R Kelley. Effective gene expression prediction from sequence by integrating long-range interactions. *Nature methods*, 18(10):1196–1203, 2021a.

- 540 Žiga Avsec, Melanie Weilert, Avanti Shrikumar, Sabrina Krueger, Amr Alexandari, Khyati Dalal,
541 Robin Fropf, Charles McAnany, Julien Gagneur, Anshul Kundaje, et al. Base-resolution models of
542 transcription-factor binding reveal soft motif syntax. *Nature Genetics*, 53(3):354–366, 2021b.
- 543
544 Gonzalo Benegas, Carlos Albors, Alan J Aw, Chengzhong Ye, and Yun S Song. Gpn-msa: an
545 alignment-based dna language model for genome-wide variant effect prediction. *bioRxiv*, 2023a.
- 546
547 Gonzalo Benegas, Sanjit Singh Batra, and Yun S Song. Dna language models are powerful predictors
548 of genome-wide variant effects. *Proceedings of the National Academy of Sciences*, 120(44):
549 e2311219120, 2023b.
- 550
551 Bradley E Bernstein, John A Stamatoyannopoulos, Joseph F Costello, Bing Ren, Aleksandar Milosavl-
552 jevic, Alexander Meissner, Manolis Kellis, Marco A Marra, Arthur L Beaudet, Joseph R Ecker, et al.
553 The nih roadmap epigenomics mapping consortium. *Nature biotechnology*, 28(10):1045–1048,
2010.
- 554 bloc97. Ntk-aware scaled rope allows llama models to have extended (8k+) con-
555 text size without any fine-tuning and minimal perplexity degradation, 2023. URL
556 [https://www.reddit.com/r/LocalLLaMA/comments/141z7j5/ntkaware_
557 scaled_rope_allows_llama_models_to_have](https://www.reddit.com/r/LocalLLaMA/comments/141z7j5/ntkaware_scaled_rope_allows_llama_models_to_have).
- 558
559 James Bradbury, Roy Frostig, Peter Hawkins, Matthew Johnson James, Chris Leary, Dougal
560 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and
561 Qiao Zhang. JAX: composable transformations of Python+NumPy programs, 2018. URL
562 <http://github.com/google/jax>.
- 563
564 Shouyuan Chen, Sherman Wong, Liangjian Chen, and Yuandong Tian. Extending context window of
565 large language models via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.
- 566
567 Siwei Chen, Laurent C. Francioli, Julia K. Goodrich, Ryan L. Collins, Masahiro Kanai, Qingbo
568 Wang, Jessica Alföldi, Nicholas A. Watts, Christopher Vittal, Laura D. Gauthier, Timothy Poterba,
569 Michael W. Wilson, Yekaterina Tarasova, William Phu, Mary T. Yohannes, Zan Koenig, Yossi
570 Farjoun, Eric Banks, Stacey Donnelly, Stacey Gabriel, Namrata Gupta, Steven Ferreira, Charlotte
571 Tolonen, Sam Novod, Louis Bergelson, David Roazen, Valentin Ruano-Rubio, Miguel Covarru-
572 bias, Christopher Llanwarne, Nikelle Petrillo, Gordon Wade, Thibault Jeandet, Ruchi Munshi,
573 Kathleen Tibbetts, gnomAD Project Consortium, Anne O’Donnell-Luria, Matthew Solomonson,
574 Cotton Seed, Alicia R. Martin, Michael E. Talkowski, Heidi L. Rehm, Mark J. Daly, Grace
575 Tiao, Benjamin M. Neale, Daniel G. MacArthur, and Konrad J. Karczewski. A genome-wide
576 mutational constraint map quantified from variation in 76,156 human genomes. *bioRxiv*, 2022.
doi: 10.1101/2022.03.20.485034. URL [https://www.biorxiv.org/content/early/
2022/03/21/2022.03.20.485034](https://www.biorxiv.org/content/early/2022/03/21/2022.03.20.485034).
- 577
578 Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for
579 contrastive learning of visual representations, 2020. URL [https://arxiv.org/abs/2002.
05709](https://arxiv.org/abs/2002.05709).
- 580
581 Deanna M Church, Valerie A Schneider, Tina Graves, Katherine Auger, Fiona Cunningham, Nathan
582 Bouk, Hsiu-Chuan Chen, Richa Agarwala, William M McLaren, Graham RS Ritchie, et al.
583 Modernizing reference genome assemblies. *PLoS biology*, 9(7):e1001091, 2011.
- 584
585 Peter JA Cock, Tiago Antao, Jeffrey T Chang, Brad A Chapman, Cymon J Cox, Andrew Dalke, Iddo
586 Friedberg, Thomas Hamelryck, Frank Kauff, Bartek Wilczynski, et al. Biopython: freely available
587 python tools for computational molecular biology and bioinformatics. *Bioinformatics*, 25(11):
1422, 2009.
- 588
589 The GTEx Consortium. The gtex consortium atlas of genetic regulatory effects across human tissues.
590 *Science*, 369(6509):1318–1330, 2020.
- 591
592 Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza-Revilla, Nicolas Lopez Carranza, Adam Henryk
593 Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Bernardo P de Almeida, Hassan
Sirelkhatim, et al. The nucleotide transformer: Building and evaluating robust foundation models
for human genomics. *bioRxiv*, pp. 2023–01, 2023.

- 594 Bernardo P de Almeida, Franziska Reiter, Michaela Pagani, and Alexander Stark. Deepstarr predicts
595 enhancer activity from dna sequence and enables the de novo design of synthetic enhancers. *Nature*
596 *genetics*, 54(5):613–624, 2022.
- 597
- 598 Bernardo P de Almeida, Hugo Dalla-Torre, Guillaume Richard, Christopher Blum, Lorenz Hexemer,
599 Maxence Gélard, Javier Mendoza-Revilla, Priyanka Pandey, Stefan Laurent, Marie Lopez, et al.
600 Segmentnt: annotating the genome at single-nucleotide resolution with dna foundation models.
601 *bioRxiv*, pp. 2024–03, 2024.
- 602 DeepMind, Igor Babuschkin, Kate Baumli, Alison Bell, Surya Bhupatiraju, Jake Bruce, Peter
603 Buchlovsky, David Budden, Trevor Cai, Aidan Clark, Ivo Danihelka, Antoine Dedieu, Clau-
604 dio Fantacci, Jonathan Godwin, Chris Jones, Ross Hemsley, Tom Hennigan, Matteo Hessel,
605 Shaobo Hou, Steven Kapturowski, Thomas Keck, Iurii Kemaev, Michael King, Markus Kunesch,
606 Lena Martens, Hamza Merzic, Vladimir Mikulik, Tamara Norman, George Papamakarios, John
607 Quan, Roman Ring, Francisco Ruiz, Alvaro Sanchez, Laurent Sartran, Rosalia Schneider,
608 Eren Sezener, Stephen Spencer, Srivatsan Srinivasan, Miloš Stanojević, Wojciech Stokowiec,
609 Luyu Wang, Guangyao Zhou, and Fabio Viola. The DeepMind JAX Ecosystem, 2020. URL
610 <http://github.com/google-deepmind>.
- 611 Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep
612 bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- 613
- 614 Yann Dubois, Gautier Dagan, Dieuwke Hupkes, and Elia Bruni. Location attention for extrapolation
615 to longer sequences. *arXiv preprint arXiv:1911.03872*, 2019.
- 616 Veniamin Fishman, Yuri Kuratov, Maxim Petrov, Aleksei Shmelev, Denis Shepelin, Nikolay
617 Chekanov, Olga Kardymon, and Mikhail Burtsev. Gena-lm: A family of open-source foundational
618 models for long dna sequences. *biorexiv. bioRxiv*, 2023.
- 619
- 620 Alistair RR Forrest, Hideya Kawaji, Michael Rehli, JK Baillie, MJL De Hoon, T Lassmann, M Itoh,
621 KM Summers, H Suzuki, CO Daub, et al. a, carninci p, hayashizaki y. a promoter-level mammalian
622 expression atlas. *Nature*, 507:462–70, 2014.
- 623 Eileen E. M. Furlong and Michael Levine. Developmental enhancers and chromosome topology.
624 *Science*, 361(6409):1341–1345, 2018. doi: 10.1126/science.aau0320. URL [https://www.](https://www.science.org/doi/abs/10.1126/science.aau0320)
625 [science.org/doi/abs/10.1126/science.aau0320](https://www.science.org/doi/abs/10.1126/science.aau0320).
- 626
- 627 Katarína Grešová, Vlastimil Martinek, David Čechák, Petr Šimeček, and Panagiotis Alexiou. Ge-
628 nomic benchmarks: a collection of datasets for genomic sequence classification. *BMC Genomic*
629 *Data*, 24(1):25, 2023.
- 630 Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv*
631 *preprint arXiv:2312.00752*, 2023.
- 632
- 633 Charles R. Harris, K. Jarrod Millman, Stéfan J. van der Walt, Ralf Gommers, Pauli Virtanen, David
634 Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J. Smith, Robert Kern, Matti
635 Picus, Stephan Hoyer, Marten H. van Kerkwijk, Matthew Brett, Allan Haldane, Jaime Fernández
636 del Río, Mark Wiebe, Pearu Peterson, Pierre Gérard-Marchant, Kevin Sheppard, Tyler Reddy,
637 Warren Weckesser, Hameer Abbasi, Christoph Gohlke, and Travis E. Oliphant. Array programming
638 with NumPy. *Nature*, 585(7825):357–362, September 2020. doi: 10.1038/s41586-020-2649-2.
639 URL <https://doi.org/10.1038/s41586-020-2649-2>.
- 640 Jennifer Harrow, Adam Frankish, Jose M Gonzalez, Electra Tapanari, Mark Diekhans, Felix Kokocin-
641 ski, Bronwen L Aken, Daniel Barrell, Amonida Zadissa, Stephen Searle, et al. Gencode: the
642 reference human genome annotation for the encode project. *Genome research*, 22(9):1760–1774,
643 2012.
- 644 Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv preprint*
645 *arXiv:1606.08415*, 2016.
- 646
- 647 Tom Hennigan, Trevor Cai, Tamara Norman, Lena Martens, and Igor Babuschkin. Haiku: Sonnet for
JAX, 2020. URL <http://github.com/deepmind/dm-haiku>.

- 648 Angela S Hinrichs, Donna Karolchik, Robert Baertsch, Galt P Barber, Gill Bejerano, Hiram Clawson,
649 Mark Diekhans, Terrence S Furey, Rachel A Harte, Fan Hsu, et al. The ucsc genome browser
650 database: update 2006. *Nucleic acids research*, 34(suppl_1):D590–D598, 2006.
- 651
- 652 J. D. Hunter. Matplotlib: A 2d graphics environment. *Computing in Science & Engineering*, 9(3):
653 90–95, 2007. doi: 10.1109/MCSE.2007.55.
- 654 Yanrong Ji, Zhihan Zhou, Han Liu, and Ramana V Davuluri. Dnabert: pre-trained bidirectional
655 encoder representations from transformers model for dna-language in genome. *Bioinformatics*, 37
656 (15):2112–2120, 2021.
- 657
- 658 kaiokendev. Things i’m learning while training superhot, 2023. URL [https://kaiokendev.
659 github.io/til#extending-context-to-8k](https://kaiokendev.github.io/til#extending-context-to-8k).
- 660 Amirhossein Kazemnejad, Inkit Padhi, Karthikeyan Natesan Ramamurthy, Payel Das, and Siva
661 Reddy. The impact of positional encoding on length generalization in transformers. *arXiv preprint
662 arXiv:2305.19466*, 2023.
- 663
- 664 David R Kelley. Cross-species regulatory sequence activity prediction. *PLoS computational biology*,
665 16(7):e1008050, 2020.
- 666
- 667 David R Kelley, Jasper Snoek, and John L Rinn. Basset: learning the regulatory code of the accessible
668 genome with deep convolutional neural networks. *Genome research*, 26(7):990–999, 2016.
- 669 Daehwan Kim, Joseph M Paggi, Chanhee Park, Christopher Bennett, and Steven L Salzberg. Graph-
670 based genome alignment and genotyping with hisat2 and hisat-genotype. *Nature biotechnology*, 37
671 (8):907–915, 2019.
- 672 Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
673 arXiv:1412.6980*, 2014.
- 674
- 675 Thomas Kluyver, Benjamin Ragan-Kelley, Fernando Pérez, Brian Granger, Matthias Bussonnier,
676 Jonathan Frederic, Kyle Kelley, Jessica Hamrick, Jason Grout, Sylvain Corlay, Paul Ivanov, Damián
677 Avila, Safia Abdalla, and Carol Willing. Jupyter notebooks – a publishing format for reproducible
678 computational workflows. In F. Loizides and B. Schmidt (eds.), *Positioning and Power in Academic
679 Publishing: Players, Agents and Agendas*, pp. 87 – 90. IOS Press, 2016.
- 680 Andriy Kryshchak, Torsten Schwede, Maya Topf, Krzysztof Fidelis, and John Moulton. Critical
681 assessment of methods of protein structure prediction (casp)—round xiv. *Proteins: Structure,
682 Function, and Bioinformatics*, 89(12):1607–1617, 2021.
- 683
- 684 Taku Kudo and John Richardson. Sentencepiece: A simple and language independent subword
685 tokenizer and detokenizer for neural text processing. *arXiv preprint arXiv:1808.06226*, 2018.
- 686
- 687 Melissa J Landrum, Shanmuga Chitipiralla, Garth R Brown, Chao Chen, Baoshan Gu, Jennifer Hart,
688 Douglas Hoffman, Wonhee Jang, Kuljeet Kaur, Chunlei Liu, Vitaly Lyoshin, Zenith Maddipati,
689 Rama Maiti, Joseph Mitchell, Nuala O’Leary, George R Riley, Wenyao Shi, George Zhou, Valerie
690 Schneider, Donna Maglott, J Bradley Holmes, , and Brandi L Kattman. Clinvar: improvements to
691 accessing data. *Nucleic Acids Res*, 8:D835–D844, 2020.
- 692
- 693 Tong Ihn Lee and Richard A. Young. Transcriptional regulation and its misregulation in disease. *Cell*,
694 152(6):P1237–1251, 2013.
- 695
- 696 Monkol Lek, Konrad J. Karczewski, Eric V. Minikel, Kaitlin E. Samocha, Eric Banks, Timothy
697 Fennell, Anne H. O’Donnell-Luria, James S. Ware, Andrew J. Hill, Beryl B. Cummings, Taru
698 Tukiainen, Daniel P. Birmbaum, Jack A. Kosmicki, Laramie E. Duncan, Karol Estrada, Fengmei
699 Zhao, James Zou, Emma Pierce-Hoffman, Joanne Berghout, David N. Cooper, Nicole DeFlaux,
700 Mark DePristo, Ron Do, Jason Flannick, and Exome Aggregation Consortium. Analysis of
701 protein-coding genetic variation in 60,706 humans. *Nature*, 536:285–291, 2016.

- 702 Zeming Lin, Halil Akin, Roshan Rao, Brian Hie, Zhongkai Zhu, Wenting Lu, Nikita Smetanin,
703 Robert Verkuil, Ori Kabeli, Yaniv Shmueli, et al. Evolutionary-scale prediction of atomic-level
704 protein structure with a language model. *Science*, 379(6637):1123–1130, 2023.
705
- 706 Zicheng Liu, Jiahui Li, Siyuan Li, Zelin Zang, Cheng Tan, Yufei Huang, Yajing Bai, and Stan Z Li.
707 Genbench: A benchmarking suite for systematic evaluation of genomic foundation models. *arXiv*
708 *preprint arXiv:2406.01627*, 2024.
- 709 Ali Madani, Ben Krause, Eric R Greene, Subu Subramanian, Benjamin P Mohr, James M Holton,
710 Jose Luis Olmos, Caiming Xiong, Zachary Z Sun, Richard Socher, et al. Large language models
711 generate functional protein sequences across diverse families. *Nature Biotechnology*, 41(8):
712 1099–1106, 2023.
713
- 714 Frederikke Isa Marin, Felix Teufel, Marc Horlacher, Dennis Madsen, Dennis Pultz, Ole Winther, and
715 Wouter Boomsma. Bend: Benchmarking dna language models on biologically meaningful tasks.
716 In *The Twelfth International Conference on Learning Representations*, 2023.
- 717 Shruti Marwaha, Joshua W. Knowles, and Euan A. Ashley. A guide for the diagnosis of rare and
718 undiagnosed disease: beyond the exome. *Genome Medicine*, 14:23, 2022.
719
- 720 Wouter Meuleman, Alexander Muratov, Eric Rynes, Jessica Halow, Kristen Lee, Daniel Bates,
721 Morgan Diegel, Douglas Dunn, Fidencio Neri, Athanasios Teodosiadis, et al. Index and biological
722 spectrum of human dnase i hypersensitive sites. *Nature*, 584(7820):244–251, 2020.
723
- 724 Jill E Moore, Michael J Purcaro, Henry E Pratt, Charles B Epstein, Noam Shores, Jessika Adrian,
725 Trupti Kawli, Carrie A Davis, Alexander Dobin, et al. Expanded encyclopaedias of dna elements
726 in the human and mouse genomes. *Nature*, 583(7818):699–710, 2020.
- 727 Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow,
728 Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, et al. Hyenadna: Long-range
729 genomic sequence modeling at single nucleotide resolution. *arXiv preprint arXiv:2306.15794*,
730 2023.
731
- 732 Eric Nguyen, Michael Poli, Matthew G. Durrant, Armin W. Thomas, Brian Kang, Jeremy Sul-
733 livan, Madelena Y. Ng, Ashley Lewis, Aman Patel, Aaron Lou, Stefano Ermon, Stephen A.
734 Baccus, Tina Hernandez-Boussard, Christopher Ré, Patrick D. Hsu, and Brian L. Hie. Se-
735 quence modeling and design from molecular to genome scale with evo. *bioRxiv*, 2024.
736 doi: 10.1101/2024.02.27.582234. URL [https://www.biorxiv.org/content/early/
2024/02/27/2024.02.27.582234](https://www.biorxiv.org/content/early/2024/02/27/2024.02.27.582234).
737
- 738 Pascal Notin, Aaron Kollasch, Daniel Ritter, Lood Van Niekerk, Steffanie Paul, Han Spinner, Nathan
739 Rollins, Ada Shaw, Rose Orenbuch, Ruben Weitzman, et al. Proteingym: large-scale benchmarks
740 for protein fitness prediction and design. *Advances in Neural Information Processing Systems*, 36,
741 2023.
- 742 Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov,
743 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning
744 robust visual features without supervision. *arXiv preprint arXiv:2304.07193*, 2023.
745
- 746 Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor
747 Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward
748 Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner,
749 Lu Fang, Junjie Bai, and Soumith Chintala. PyTorch: An Imperative Style, High-Performance
750 Deep Learning Library. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d’Alché Buc, E. Fox,
751 and R. Garnett (eds.), *Advances in Neural Information Processing Systems 32*, pp. 8024–8035.
752 Curran Associates, Inc., 2019.
- 753 F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Pretten-
754 hofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and
755 E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*,
12:2825–2830, 2011.

- 756 Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. Yarn: Efficient context window
757 extension of large language models. *arXiv preprint arXiv:2309.00071*, 2023.
- 758
- 759 Michael Poli, Stefano Massaroli, Eric Nguyen, Daniel Y Fu, Tri Dao, Stephen Baccus, Yoshua
760 Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional
761 language models. *arXiv preprint arXiv:2302.10866*, 2023.
- 762 Katherine S Pollard, Melissa J Hubisz, Kate R Rosenbloom, and Adam Siepel. Detection of nonneutral
763 substitution rates on mammalian phylogenies. *Genome Res.*, 20(1):110–121, January 2010.
- 764
- 765 Ofir Press, Noah A Smith, and Mike Lewis. Train short, test long: Attention with linear biases
766 enables input length extrapolation. *arXiv preprint arXiv:2108.12409*, 2021.
- 767 Ofir Press, Noah Smith, and Mike Lewis. Train short, test long: Attention with linear biases enables
768 input length extrapolation. In *International Conference on Learning Representations*, 2022. URL
769 <https://openreview.net/forum?id=R8sQPpGCv0>.
- 770 Markus N Rabe and Charles Staats. Self-attention does not need $O(n^2)$ memory. *arXiv preprint*
771 *arXiv:2112.05682*, 2021.
- 772
- 773 Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal,
774 Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual
775 models from natural language supervision. In *International conference on machine learning*, pp.
776 8748–8763. PMLR, 2021.
- 777 Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel,
778 and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information*
779 *processing systems*, 32, 2019.
- 780
- 781 Yuchen Ren, Zhiyuan Chen, Lifeng Qiao, Hongtai Jing, Yuchen Cai, Sheng Xu, Peng Ye, Xinzhu Ma,
782 Siqi Sun, Hongliang Yan, et al. Beacon: Benchmark for comprehensive rna tasks and language
783 models. *arXiv preprint arXiv:2406.10391*, 2024.
- 784 Yair Schiff, Chia-Hsiang Kao, Aaron Gokaslan, Tri Dao, Albert Gu, and Volodymyr Kuleshov.
785 Caduceus: Bi-directional equivariant long-range dna sequence modeling. *arXiv preprint*
786 *arXiv:2403.03234*, 2024.
- 787
- 788 Valerie A Schneider, Tina Graves-Lindsay, Kerstin Howe, Nathan Bouk, Hsiu-Chuan Chen, Paul A
789 Kitts, Terence D Murphy, Kim D Pruitt, Françoise Thibaud-Nissen, Derek Albracht, et al. Evaluation
790 of grch38 and de novo haploid genome assemblies demonstrates the enduring quality of the
791 reference assembly. *Genome research*, 27(5):849–864, 2017.
- 792 Max Schubach, Thorben Maass, Lusiné Nazaretyan, Sebastian Röner, and Martin Kircher. Cadd v1.
793 7: using protein language models, regulatory cnns and other nucleotide-level scores to improve
794 genome-wide variant predictions. *Nucleic Acids Research*, 52(D1):D1143–D1154, 2024.
- 795 Rico Sennrich, Barry Haddow, and Alexandra Birch. Neural machine translation of rare words with
796 subword units. *arXiv preprint arXiv:1508.07909*, 2015.
- 797
- 798 Matthew D Shirley, Zhaorong Ma, Brent S Pedersen, and Sarah J Wheelan. Efficient" pythonic"
799 access to fasta files using pyfaidx. Technical report, PeerJ PrePrints, 2015.
- 800 Adam Siepel, Gill Bejerano, Jakob S Pedersen, Angie S Hinrichs, Minmei Hou, Kate Rosenbloom, Hi-
801 ram Clawson, John Spieth, Ladeana W Hillier, Stephen Richards, George M Weinstock, Richard K
802 Wilson, Richard A Gibbs, W James Kent, Webb Miller, and David Haussler. Evolutionarily con-
803 served elements in vertebrate, insect, worm, and yeast genomes. *Genome Res.*, 15(8):1034–1050,
804 August 2005.
- 805
- 806 Damian Smedley, Max Schubach, Julius OB Jacobsen, Sebastian Köhler, Tomasz Zemojtel, Malte
807 Spielmann, Marten Jäger, Harry Hochheiser, Nicole L Washington, Julie A McMurry, et al. A
808 whole-genome analysis framework for effective identification of pathogenic regulatory variants in
809 mendelian disease. *The American Journal of Human Genetics*, 99(3):595–606, 2016.
- AFA Smit, R Hubley, and P Green. Repeatmasker open-4.0. 2013–2015, 2015.

- 810 Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Murtadha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced
811 transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021.
812
- 813 Hazuki Takahashi, Timo Lassmann, Mitsuyoshi Murata, and Piero Carninci. 5' end-centered
814 expression profiling using cap-analysis gene expression and next-generation sequencing. *Nature*
815 *Protocols*, 7:542–561, 2012.
- 816 Matthew Tancik, Pratul Srinivasan, Ben Mildenhall, Sara Fridovich-Keil, Nithin Raghavan, Utkarsh
817 Singhal, Ravi Ramamoorthi, Jonathan Barron, and Ren Ng. Fourier features let networks learn
818 high frequency functions in low dimensional domains. *Advances in Neural Information Processing*
819 *Systems*, 33:7537–7547, 2020.
- 820
821 Gemini Team, Rohan Anil, Sebastian Borgeaud, Yonghui Wu, Jean-Baptiste Alayrac, Jiahui Yu, Radu
822 Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, et al. Gemini: a family of highly capable
823 multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.
- 824
825 The pandas development team. pandas-dev/pandas: Pandas, February 2020. URL [https://](https://github.com/pandas-dev/pandas)
826 github.com/pandas-dev/pandas.
- 827
828 The ENCODE Project Consortium. Expanded encyclopaedias of dna elements in the human and
829 mouse genomes. *Nature*, 583(7818):699–710, 2020.
- 830
831 Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée
832 Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and
833 efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023a.
- 834
835 Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay
836 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shrutu Bhosale, et al. Llama 2: Open foundation
837 and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023b.
- 838
839 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz
840 Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing*
841 *systems*, 30, 2017.
- 842
843 Gao Wang, Abhishek Sarkar, Peter Carbonetto, and Matthew Stephens. A simple new approach to
844 variable selection in regression, with application to genetic fine mapping. *Journal of the Royal*
845 *Statistical Society Series B: Statistical Methodology*, 82(5):1273–1300, 2020.
- 846
847 Qingbo S. Wang, David R. Kelley, Jacob Ulirsch, Masahiro Kanai, Shuvom Sadhuka, Ran Cui,
848 Carlos Albors, Nathan Cheng, Yukinori Okada, The Biobank Japan Project, Francois Aguet,
849 Kristin G. Ardlie, Daniel G. MacArthur, and Hilary K. Finucane. Leveraging supervised learning
850 for functionally informed fine-mapping of cis-eqtls identifies an additional 20,913 putative causal
851 eqtls. *Nature Communications*, 12:3394, 2021.
- 852
853 Michael L. Waskom. seaborn: statistical data visualization. *Journal of Open Source Software*, 6(60):
854 3021, 2021. doi: 10.21105/joss.03021. URL <https://doi.org/10.21105/joss.03021>.
- 855
856 Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi,
857 Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Huggingface’s transformers:
858 State-of-the-art natural language processing. *arXiv preprint arXiv:1910.03771*, 2019.
- 859
860 Minghao Xu, Zuobai Zhang, Jiarui Lu, Zhaocheng Zhu, Yangtian Zhang, Ma Chang, Runcheng
861 Liu, and Jian Tang. Peer: a comprehensive and multi-task benchmark for protein sequence
862 understanding. *Advances in Neural Information Processing Systems*, 35:35156–35173, 2022.
- 863
864 Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago
865 Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for
866 longer sequences. *Advances in neural information processing systems*, 33:17283–17297, 2020.
- 867
868 Jian Zhou and Olga G Troyanskaya. Predicting effects of noncoding variants with deep learning-based
869 sequence model. *Nature methods*, 12(10):931–934, 2015.

Jian Zhou, Chandra L Theesfeld, Kevin Yao, Kathleen M Chen, Aaron K Wong, and Olga G Troyanskaya. Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nature genetics*, 50(8):1171–1179, 2018a.

Jian Zhou, Chandra L Theesfeld, Kevin Yao, Kathleen M Chen, Aaron K Wong, and Olga G Troyanskaya. Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nature genetics*, 50(8):1171–1179, 2018b.

Zhihan Zhou, Yanrong Ji, Weijian Li, Pratik Dutta, Ramana Davuluri, and Han Liu. Dnabert-2: Efficient foundation model and benchmark for multi-species genome. *arXiv preprint arXiv:2306.15006*, 2023.

Zhihan Zhou, Weimin Wu, Harrison Ho, Jiayi Wang, Lizhen Shi, Ramana V Davuluri, Zhong Wang, and Han Liu. Dnabert-s: Learning species-aware dna embedding with genome foundation models, 2024. URL <https://arxiv.org/abs/2402.08777>.

A EXTENDED BACKGROUND

A.1 TERMINOLOGY

The genome is a sequence of four nucleotides (*Adenine*, *Cytosine*, *Thymine*, and *Guanine*) organized into a double-stranded helical structure called *deoxyribonucleic acid* (DNA). This structure encodes the information required for the development, maintenance, and function of cells. Genetic information flows from DNA to *messenger ribonucleic acid* (mRNA) by a process called *transcription*, and mRNA is used as a blueprint to create *proteins* via a process called *translation*. Proteins are responsible for initiating and sustaining the cellular processes, while DNA encodes the information necessary for their production.

The genome is organized into functional elements, including *coding* and *non-coding* regions. Coding regions comprise genes responsible for protein synthesis, while non-coding regions can play vital regulatory roles. *Promoters*, a type of regulatory region, are situated close to genes and serve as sites for transcription initiation. *Enhancers*, another regulatory element located farther from genes, modulate gene expression by recruiting transcription factors, a type of protein that regulates transcription. Notably, a single gene can be regulated by multiple promoters and enhancers simultaneously.

DNA does not exist solely as a linear molecule but is instead tightly packaged around *histone* proteins, forming a sphere of wound DNA called *nucleosomes*. These nucleosomes further assemble into *chromatin*, which constitutes the 23 pairs of *chromosomes* in humans. Chromatin can exist in an open (*euchromatin*) or closed (*heterochromatin*) state, influencing the ability of the underlying DNA to be transcribed. Chemical modifications to histones play significant roles in chromatin remodeling acting as signals that recruit proteins to either condense the chromatin structure (making it less accessible) or relax it (making it more accessible), thereby influencing gene activity.

Mutations in the genome, including *single nucleotide polymorphisms* (SNPs), insertions, and deletions, can alter DNA sequences, potentially disrupting functional genomic elements or affecting the structure and function of proteins. Understanding the impact of these sequence variations on disease remains a central challenge in biology. Such mutations can lead to genetic disorders or contribute to the development of complex diseases.

A.2 RECENT DNA LMS

DNABERT Arguably the first DNA LM, DNABERT proposed in Ji et al. (2021) applies the BERT architecture from Devlin et al. (2018), with a few modifications, to genomic sequences. The authors train on the human genome and use k -mer tokens generated with sliding windows. Input sequences were 512 tokens, and the model was trained using the MLM objective, but with the restriction that masking was performed for contiguous tokens within a sequence. The downstream tasks focused on genome annotation, with promoter, transcription factor binding sites, and splice site classification. Of note, although DNABERT was pre-trained on human genome, it was fine-tuned on mouse downstream tasks as well, yielding competitive performance relative to supervised learning baselines.

Nucleotide Transformer Following the success of model scaling in other domains, [Dalla-Torre et al. \(2023\)](#) explore scaling DNA foundation models in introducing the Nucleotide Transformer. They explore various model sizes – ranging from 500 million parameters to 2.5 billion, in their first generation release, and 50 million to 500 million parameters in their subsequent version 2 – and various pre-training data setups, including human reference genome, 3,000 diverse human genomes, and 850 multi-species reference genomes. They utilize non-overlapping 6-mer tokenization and a BERT-style architecture trained with an MLM objective. Other notable differences between the first and second version is that in version 2 input context size was scaled from 1,000 tokens to 2,000 and positional embeddings used in version 1 were learned whereas version 2 used rotary embeddings ([Su et al., 2021](#)), which have been shown to better extend to longer contexts. This work also introduced the Nucleotide Transformer suite of tasks, described in more detail below.

DNABERT-2 Building on the initial success of DNABERT, [Zhou et al. \(2023\)](#) present a model trained on multi-species genomes: 135 species, across 7 categories. They also change tokenization to byte-pair-encoding ([Kudo & Richardson, 2018](#); [Sennrich et al., 2015](#)), with a vocabulary size of 4,096, arguing that overlapping k -mer tokenization makes the MLM task ‘too easy’ by leaking information across tokens and that non-overlapping k -mer tokenization suffers from the drawback that minor changes to the input sequence, e.g., removing the first character, lead to drastically different tokenization outputs. They use input sequence lengths of 128 tokens. Additionally, [Zhou et al. \(2023\)](#) replace the learned positional embeddings from DNABERT with ALiBi ([Press et al., 2022](#)). DNABERT-2 was evaluated on a suite of downstream tasks introduced in [Zhou et al. \(2023\)](#) known as the Genome Understanding Evaluation (GUE).

DNABERT-S DNABERT-S then further builds upon DNABERT-2 to generate richer DNA embeddings for uses in meta genomics. For DNABERT-S, [Zhou et al. \(2024\)](#) takes a trained DNABERT-2 model and they train it with a contrastive learning objective with a novel strategy with a new training dataset: 2 million paired DNA sequences from fungi, viruses, and bacteria. Additionally, their contrastive learning approach involves a new strategy: Curriculum Contrastive Learning (C^2LR) that defines a training curriculum for the model and makes use of a new training objective: Manifold Instance Mixup (MI-Mix). [Zhou et al. \(2024\)](#) use this MI-Mix objective in addition to the the SimCLR contrastive loss objective for different phases of their training ([Chen et al., 2020](#)).

HyenaDNA In contrast to the other language models reviewed above, the HyenaDNA model from [Nguyen et al. \(2023\)](#) is a next token prediction, uni-directional model. Using character-level tokenization and the Hyena layers ([Poli et al., 2023](#)) as a backbone, [Nguyen et al. \(2023\)](#) also propose a training recipe for scaling input context sizes up to 1 million bps. To evaluate their model they use a combination of downstream tasks, including the suite of tasks from Nucleotide Transformer ([Dalla-Torre et al., 2023](#)), a set of mouse and human genome annotation tasks presented in [Grešová et al. \(2023\)](#), the chromatin profiling tasks from DeepSea ([Zhou & Troyanskaya, 2015](#)), and a species classification task, where the model takes in sequences of various species and needs to output the correct species label.

Caduceus In the recent **Mamba** work ([Gu & Dao, 2023](#)), the authors pre-train various sized models that use the Mamba backbone on the human reference genome. Similar to HyenaDNA, the pre-training objective is next token prediction, tokenization is by nucleotide base, and input sequences are scaled up to 1 million bps. Building off this work, [Schiff et al. \(2024\)](#) introduced **Caduceus**, a bi-directional Mamba-based model that contains reverse complement equivariance inductive biases, demonstrating state-of-the-art performance on several tasks, including several Nucleotide Transformer tasks ([Dalla-Torre et al., 2023](#)) and the Genomic Benchmark ([Grešová et al., 2023](#)).

GPN-MSA GPN-MSA ([Benegas et al., 2023a](#)) is an instance of an **alignment-based** DNA language model. Instead of having a single sequence as an input, a Multiple Sequence Alignment (MSA) is used. This alignment contains the same human DNA sequences as other DNA LMs, but it additionally contains an MSA containing the sequence information corresponding to the human sequences from multiple (89) other species. This extra information allows the GPN-MSA model to model the human sequence conditioned on the evolutionary information provided by the species included in the MSA. GPN-MSA is based on the RoFormer ([Su et al., 2021](#)) architecture with the difference that it flattens

and encodes the provided MSA into the hidden dimension of the embedding. GPN-MSA uses a single-nucleotide tokenizer and is trained with a context size of 128 bps.

GPN-MSA also has many notable differences from the remaining DNA LMs regarding training since GPN-MSA is trained on a curated subset of the human genome. Briefly, phastCONS (Siepel et al., 2005) is used to tag each nucleotide with a conservation probability, then the genome is subdivided into windows and the 5% windows that are predicted to be the most highly conserved are selected for training, along with 0.1% of the remaining windows and a reverse complement of all the windows. GPN-MSA is then trained in a manner similar to the other MLM models, 15% of the human input sequences are masked and a weighted cross entropy loss is used for the reconstructed tokens. The weights for the Cross Entropy Loss attempt to downweight repetitive regions and upweight conserved regions, with the weights for each nucleotide being determined based on calling repeat regions and the phyloP and phastCONS score (Siepel et al., 2005; Pollard et al., 2010) for that nucleotide.

Variant effect prediction takes the same form as for the other MLM models, except that GPN-MSA additionally takes as input the MSA for 89 other species at that location as well. Notably, the input MSA to GPN-MSA is never masked.

Other DNA LMs While the models above represent those that we initially validate on our benchmark, the field of DNA LMs is growing at a rapid pace and consists of several notable works that we briefly describe below.

While not developed specifically as a DNA LM, the **BigBird** architecture proposed in Zaheer et al. (2020) was applied to genomic sequences to demonstrate its usefulness in long context tasks. Using sparse attention to reduce computational complexity of transformer (Vaswani et al., 2017) blocks from quadratic to linear, BigBird is able to effectively scale up to longer contexts. In Fishman et al. (2023), the authors present a family of foundation models, **GENA-LM**, aimed specifically at modeling longer DNA sequences. Pre-training with an MLM objective on human and multi-species genomes, they use BPE with a vocabulary size of 32,000. The backbone architectures are either BERT (Devlin et al., 2018) or BigBird (Zaheer et al., 2020), allowing them to extend input lengths up to 36k bps.

Focusing on plant genomes, Benegas et al. (2023b) pre-train a MLM model on unaligned reference genomes of the *Arabidopsis thaliana* species and seven related species within the Brassicales order. Using character-level tokenization they use input lengths of 512 bps with dilated convolutions to create their **GPN** model. With 25 layers, despite the relatively short training sequences, GPN can theoretically extend to sequence inputs of millions of bps.

A.3 EXISTING DNA LANGUAGE MODEL BENCHMARKS

Existing benchmarks vary in several aspects, including the species considered, the specific tasks of interest, the framing of these tasks, and the evaluation methodologies employed. These proposed benchmarks include the Nucleotide Transformer Benchmark (Dalla-Torre et al., 2023), Genomic Benchmarks (Grešová et al., 2023), Genome Understanding Evaluation (GUE, (Zhou et al., 2023)), and Benchmarking DNA Language Models on Biologically Meaningful Tasks (BEND; Marin et al. (2023)).

Existing DNA benchmarks are primarily composed of classification tasks for sequence-wise predictions, ranging from cis-regulatory elements and splice sites to chromatin features and variant effects. These benchmarks not only compile and build datasets but also carry out evaluations of DNA LMs using both fine-tuning methods, where pre-trained models are trained in a supervised manner on the downstream tasks, and zero-shot prediction, where models are evaluated in their pre-trained state without additional fine-tuning.

Nucleotide Transformer Benchmark Dalla-Torre et al. (2023) compile a set of 18 distinct genomic datasets framed as sequence-wise classification tasks. These tasks included 10 datasets related to epigenetic mark prediction in yeast genomes, three tasks predicting the presence of promoters in mouse and human genomes, two tasks predicting enhancer presence and activity levels in the human genome, and three tasks predicting splice sites in multiple diverse species. Sequence lengths in this benchmark ranged from 200 to 600 bps. Additionally, the authors evaluated a set of DNA LMs and a supervised genomic model, Enformer (Avsec et al., 2021a), by fine-tuning these models on their benchmark using a robust 10-fold cross-validation protocol. Parameter-efficient fine-tuning methods

with a classification head were used for Enformer, DNABERT, and NT models, while full fine-tuning with a classification head was applied to the HyenaDNA models. Limitations of this benchmark include the focus on short-range contexts, the inclusion of synthetic sequences as negative examples, and limited supervised baselines.

Genomic Benchmarks Genomic Benchmarks (Grešová et al., 2023) is a collection of datasets for genomic sequence classification, composed of existing datasets and novel ones scraped from publicly available databases. The benchmark includes nine tasks focusing on regulatory element prediction, such as promoters, enhancers, and open chromatin regions. These tasks cover human, mouse, roundworm, and fly genomes, with average sequence lengths ranging from 200 to 2,370 bps. The authors also provide code to train simple convolutional network that can be used as a baseline. Similar to the Nucleotide Transformer benchmark, this benchmark focuses on short-range tasks, does not present a robust set of baselines, and contains potentially less impactful tasks, e.g., distinguishing between human and worm genomic sequences.

Genomic Understanding Evaluation (GUE) The authors of the DNABERT-2 (Zhou et al., 2023) introduced the Genomic Understanding Evaluation (GUE) benchmark, which is divided into two groups by sequence length: GUE and GUE+. This benchmark comprises seven classification tasks, such as cis-regulatory element prediction and species classification, built from 28 datasets from multiple species. The inclusion of multiple species allows for the assessment of DNA LMs’ generalizability. The tasks are curated to be appropriately challenging, including measures such as class balancing, adversarial sample inclusion, and reduction of training sample volume. GUE features sequence lengths ranging from 70 to 1k bps, while GUE+ includes sequence lengths from 5k to 10k bps. GUE evaluated DNABERT1 and 2, NT, and HyenaDNA models on their benchmark. HyenaDNA models are fully fine-tuned while DNABERT and NT models are fine-tuned using parameter efficient methods. The GUE benchmark results are limited since they do not cover a robust set of baselines but rather only present the simple supervised convolutional network from the Genomic Benchmark (Grešová et al., 2023). Additionally, only binary or multi-class sequence-wise classification tasks are considered and tasks of biological importance, such as variant effect prediction and gene expression are not included.

Benchmarking DNA LLMs on Biologically Meaningful Tasks (BEND) BEND (Marin et al., 2023) is a recently proposed benchmark focused on compiling tasks that capture the complexity and intricacies of real-world genomic analysis. The authors collected seven different datasets, all from the human genome, covering gene finding, enhancer annotation, chromatin accessibility, histone modification, CpG methylation, and two types of variant effect prediction. Unlike previous benchmarks that focused solely on sequence-wise classification tasks, BEND also includes the task “Gene finding”, which tests nucleotide-resolution modeling. In five out of seven tasks the input length is 512 bps, as these tasks are considered short-range. “Gene finding” task use sequences up to 14k bps. Their “Enhancer annotation” task uses 100k bp sequences, but it only contains 285 input sequences. Notably, for tasks in BEND that overlap with our benchmark (such as variant effect prediction), BEND uses a fixed context length of 512 bp, thus not evaluating the importance of extended context and variant-gene distal interactions on this type of task. Therefore, this benchmark is mostly limited to short-range tasks and does not include gene expression, an important and challenging task in genomics. This benchmark however makes progress in including a broader set of supervised methods as baselines. Unlike our work, models are only evaluated using partial fine-tuning, where backbone DNA LM weights are frozen for downstream task training.

GenBench The GenBench suite (Liu et al., 2024) is composed of 43 different datasets split between “short” and “long” range tasks, where long-range tasks are defined by having a sequence length of greater than 1000 base pairs. The tasks in GenBench, spanning multiple species, are primarily binary, sequence-level classification tasks but also include multi-class classification and regression tasks. The authors evaluate six different genomic language models covering both attention and convolution-based architectures. While GenBench provides a comprehensive evaluation, it lacks critical tasks like variant effect prediction in non-coding regions and zero-shot evaluations. It also omits comparisons to long-context models like Enformer and is limited in its evaluation of long-range tasks, with the longest sequence length capped at 30,000 base pairs.

BEACON The BEACON benchmark (Ren et al., 2024) introduces the first unified evaluation framework for RNA modeling, encompassing 13 tasks across structural analysis, functional studies, and engineering applications. It evaluates 29 models, ranging from pre-trained RNA language models to naive supervised models, and examines the influence of tokenization strategies and positional embeddings on performance. While BEACON is a valuable resource for assessing RNA-focused models, its scope is distinct from genomic benchmarks, as it targets RNA-specific tasks rather than genomic applications like regulatory element prediction, variant effect prediction, or gene expression prediction.

B ADDITIONAL DETAILS ABOUT GENOMIC LONG RANGE BENCHMARK

We note that our datasets do not contain any personally identifiable information or offensive content.

Table 7 provides details describing the evaluation method used, dataset sizes, metric, and data sources. Additional details on task specific data curation and processing are described in the following subsections.

Table 7: Additional information for Genomic LRB tasks, including number of samples in train and test splits, metric, and data source.

Task	Eval	Test split	Metric	Data Source
<i>Variant Effect Prediction</i>				
Causal eQTL	Fine-tune & Zero-shot	Chromosome 9, 10	AUROC	GTEx (via Avsec et al. (2021b))
Pathogenic OMIM	Zero-shot	-	AUPRC	OMIM, gnomAD (via Benegas et al. (2023a))
Pathogenic ClinVar	Fine-tune & Zero-shot	Chromosome 8	AUROC	ClinVar, gnomAD (via Benegas et al. (2023a))
<i>Gene Expression Prediction</i>				
Bulk RNA Expression	Fine-tune	Chromosome 8	R^2	GTEx, FANTOM5 (via Zhou et al. (2018a))
CAGE	Fine-tune	Random	R^2	FANTOM5 (via Kelley (2020))
<i>Regulatory Element Detection</i>				
Promoter	Fine-tune	Chromosome 8, 9	AUPRC	SCREEN
Enhancer	Fine-tune	Chromosome 8,9	AUROC	SCREEN
<i>Chromatin Feature Identification</i>				
Histone Marks	Fine-tune	Chromosome 8, 9	AUPRC	ENCODE, Roadmap Epigenomics (via Zhou & Troyanskaya (2015))
DNA Accessibility	Fine-tune	Chromosome 8, 9	AUPRC	ENCODE, Roadmap Epigenomics (via Zhou & Troyanskaya (2015))

B.1 VARIANT EFFECT PREDICTION

B.1.1 CAUSAL EQTL

Data Processing Processed data in the form of `vcf` files for positive and negative variants across 49 different tissue types were obtained from Avsec et al. (2021a). Fine-mapped GTEx (Consortium, 2020) eQTLs originate from Wang et al. (2021), while the negative matched set of variants comes from Avsec et al. (2021a). The statistical fine-mapping tool SuSiE (Wang et al., 2020) was used to label variants. Variants from the fine-mapped eQTL set were selected and given positive labels if their posterior inclusion probability was > 0.9 , as assigned by SuSiE. Variants from the matched negative set were given negative labels if their posterior inclusion probability was < 0.01 . DNA sequences were obtained from the human reference genome assembly GRCh38 (Schneider et al., 2017).

B.1.2 PATHOGENIC OMIM

Data Processing Processed data was obtained from Benegas et al. (2023a) in the form of `parquet` files with columns for SNP location, reference and alternative alleles, and pathogenicity label. Positive labeled data originates from a curated set of pathogenic variants located in the Online Mendelian Inheritance in Man (OMIM) (Smedley et al., 2016) catalog. The negative set is comprised of variants that are defined as common from gnomAD (Chen et al., 2022). gnomAD version 3.1.2 was downloaded and filtered to variants with allele number of at least 25,000. Common variants were defined as those with minor allele frequency (MAF) $> 5\%$. The input sequences were constructed by selecting the appropriate genomic region from the human reference genome assembly GRCh38 (Schneider et al., 2017) and applying the changes specified by the given variants.

1134 B.1.3 PATHOGENIC CLINVAR

1135
1136 **Data Processing** Processed data was obtained from Benegas et al. (2023a) in the form of parquet
1137 files with columns for SNP location, reference and alternative alleles, and pathogenicity label. Positive
1138 labels correspond to pathogenic variants originating from ClinVar (Landrum et al., 2020) whose
1139 review status was described as having at least a single submitted record with a classification but
1140 without assertion criteria. The negative set are variants that are defined as common from gnomAD
1141 (Chen et al., 2022). gnomAD version 3.1.2 was downloaded and filtered to variants with allele
1142 number of at least 25,000. Common variants were defined as those with $MAF > 5\%$. Sequences
1143 were obtained from the human reference genome assembly GRCh38 (Schneider et al., 2017).

1144 **Short-Range** The ClinVar dataset is mostly variants in coding regions, and since most human protein
1145 sequences have less than 1,000 amino acids predicting the impact of coding variants should require
1146 orders of magnitude smaller context windows than non-coding variants. Therefore, we consider this
1147 task as potentially short-range.

1148 B.2 GENE EXPRESSION PREDICTION

1149 B.2.1 BULK RNA-SEQ

1150
1151 **Data Processing** Processed data in the form `csv` files that contained gene TSS locations, strand,
1152 and RNA expression RPKM counts across 218 tissue types was obtained from ExPecto (Zhou
1153 et al., 2018a). Expression data originates from GTEx (Consortium, 2020), while representative TSS
1154 locations were determined in ExPecto. The authors of ExPecto determined representative TSS for
1155 Pol II transcribed genes based on quantification of CAGE reads from the FANTOM5 project (Forrest
1156 et al., 2014). The specific procedure they used is as follows, a CAGE peak was associated to a
1157 GENCODE (Harrow et al., 2012) gene if it was within 1000 bps from a GENCODE v24 annotated
1158 TSS. The most abundant CAGE peak for each gene was then selected as the representative TSS.
1159 When no CAGE peak could be assigned to a gene, the annotated gene start position was used as the
1160 representative TSS. We $\log(1+x)$ normalized then standardized the RNA-seq counts before training
1161 models. Sequences centered around the TSS were obtained from the human reference genome
1162 assembly GRCh37 (Church et al., 2011).
1163

1164 B.2.2 CAP ANALYSIS GENE EXPRESSION (CAGE) PROFILE

1165
1166 **Data Processing** Processed data was obtained from Basenji2 (Kelley, 2020), where input sequence
1167 locations were collected as `bed` files and CAGE counts as `TensorFlow` (Abadi et al., 2015) records.
1168 Original data comes from the FANTOM5 project (Forrest et al., 2014). Data was processed to produce
1169 CAGE labels for non-overlapping 128 bp bins within a sequence of 114,688 bps. For each bin, there
1170 are 638 different predictions corresponding to the CAGE count in various cell, tissue, or treatment
1171 types (e.g., fibroblast, heart, or monocytes treated with Salmonella). This resulted in an output
1172 array of $896 \text{ bins} \times 638 \text{ tracks}$ for a single sample. DNA sequences were obtained from the human
1173 reference genome assembly GRCh38 (Schneider et al., 2017).

1174 The compute requirements to store and process this data make it more difficult and less accessible to
1175 users. To achieve a balance of user-friendliness while also maintaining a representative view of the
1176 data, we sub-sampled the number of tracks to 50 by using the following guidelines:

- 1177 1. Only select one cell line.
- 1178 2. Only keep mock treated and remove other treatments.
- 1179 3. Only select one donor.

1180
1181 The 50 specific tracks which were selected can be found in Table 8 below. This maintains the number
1182 of sequences in the entire dataset but reduces the number of labels for each sequence from 638 to 50
1183 thus reducing storage requirements from $\sim 84\text{GB}$ to $\sim 7\text{GB}$.
1184
1185
1186
1187

Table 8: The 50 CAGE tracks sub-sampled for the Genomic LRB from the original 638 tracks.

Track Index	Description
0	CAGE:adipose tissue, adult, pool1
1	CAGE:bladder, adult, pool1
2	CAGE:brain, adult, pool1
3	CAGE:cervix, adult, pool1
4	CAGE:colon, adult, pool1
5	CAGE:esophagus, adult, pool1
6	CAGE:heart, adult, pool1
7	CAGE:kidney, adult, pool1
8	CAGE:liver, adult, pool1
9	CAGE:lung, adult, pool1
10	CAGE:ovary, adult, pool1
11	CAGE:placenta, adult, pool1
12	CAGE:prostate, adult, pool1
13	CAGE:skeletal muscle, adult, pool1
14	CAGE:small intestine, adult, pool1
15	CAGE:spleen, adult, pool1
16	CAGE:testis, adult, pool1
17	CAGE:thymus, adult, pool1
18	CAGE:thyroid, adult, pool1
19	CAGE:trachea, adult, pool1
20	CAGE:retina, adult, pool1
21	CAGE:temporal lobe, adult, pool1
22	CAGE:postcentral gyrus, adult, pool1
23	CAGE:pons, adult, pool1
24	CAGE:parietal lobe, adult, pool1
25	CAGE:paracentral gyrus, adult, pool1
26	CAGE:occipital pole, adult, pool1
27	CAGE:nucleus accumbens, adult, pool1
28	CAGE:medulla oblongata, adult, pool1
29	CAGE:insula, adult, pool1
30	CAGE:frontal lobe, adult, pool1
31	CAGE:dura mater, adult,
32	CAGE:corpus callosum, adult, pool1
33	CAGE:adenocarcinoma cell line:IM95m
34	CAGE:breast carcinoma cell line:MCF7
35	CAGE:diffuse large B-cell lymphoma cell line:CTB-1
36	CAGE:glioma cell line:GI-1
37	CAGE:liposarcoma cell line:SW 872
38	CAGE:Sebocyte,
39	CAGE:CD4+ T Cells,
40	CAGE:Natural Killer Cells,
41	CAGE:Neutrophils,
42	CAGE:Pericytes,
43	CAGE:Alveolar Epithelial Cells,
44	CAGE:Renal Mesangial Cells,
45	CAGE:Nucleus Pulposus Cell,
46	CAGE:Keratocytes,
47	CAGE:Mesenchymal Stem Cells - adipose,
48	CAGE:Mammary Epithelial Cell,
49	CAGE:Osteoblast,

1238
1239
1240
1241

1242 B.3 CIS-REGULATORY ELEMENT DETECTION

1243
1244 **Data Processing** Original data was sourced from Search Candidate cis-Regulatory Elements v3
1245 (SCREEN) registry by ENCODE (Moore et al., 2020). The data is processed as follows, we break
1246 the human reference genome into 200 bp non-overlapping chunks. If the 200 bp chunk overlaps
1247 by at least 50% or more with a contiguous region from the set of annotated cis-regulatory elements
1248 (promoters or enhancers), we label them as positive, else the chunk is labeled as negative. The
1249 resulting dataset was composed of ~15M negative samples and ~50k positive promoter samples and
1250 ~1M positive enhancer samples We randomly sub-sampled the negative set to 1M samples, and kept
1251 all positive samples, to make this dataset more manageable in size. DNA sequences were obtained
1252 from the human reference genome assembly GRCh38 (Schneider et al., 2017).

1253 **Short-Range** Since this task involves predicting the presence of a regulatory element within a specific
1254 sequence, only local context is believed to be important. The activity of promoters and enhancers in
1255 different cell types is dictated by the presence of binding sites for specific proteins (Andersson &
1256 Sandelin, 2020) and thus likely do not require long-distance interactions, as demonstrated by the high
1257 predictive value of models using less than 1k bp input sequences (Avsec et al., 2021b; Kelley et al.,
1258 2016).

1259 B.4 CHROMATIN FEATURE IDENTIFICATION

1260
1261 **Data Processing** Processed data was obtained from DeepSea (Zhou & Troyanskaya, 2015) in the
1262 form of 1k bp sequences and labels as `txt` files. Original chromatin profiling data comes from
1263 ENCODE and Roadmap Epigenomics (Moore et al., 2020; Bernstein et al., 2010). The authors
1264 of DeepSea processed the data by chunking the human genome into 200 bp bins where for each
1265 bin labels were determined for hundreds of different chromatin features. Only bins with at least
1266 one transcription factor binding event were considered for the dataset. If the bin overlapped with
1267 a peak region of the specific chromatin profile by more than half of the sequence, a positive label
1268 was assigned. DNA sequences were obtained from the human reference genome assembly GRCh37
1269 (Church et al., 2011). To make the dataset more accessible, we randomly sub-sampled the chromatin
1270 profiles from 125 to 20 tracks for the histones dataset and from 104 to 20 tracks for the DNase dataset.
1271 The sub-sampled tracks for both datasets can be found in Table 9 and Table 10.

1272 **Short-Range** Chromatin features are not expected to be strongly influenced by long-range interac-
1273 tions. Most of the information affecting these chromatin features occurs locally and depends on the
1274 binding of different proteins (Lee & Young, 2013). This is also corroborated by the high predictive
1275 value of models using less than 1k bps input sequences (Kelley et al., 2016; Zhou & Troyanskaya,
1276 2015).

1277
1278
1279
1280
1281
1282
1283
1284
1285
1286
1287
1288
1289
1290
1291
1292
1293
1294
1295

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

Table 9: 20 Histone tracks sub sampled for the Genomic LRB from the original 104 tracks with histone mark and cell type information.

Track Index	Histone Mark	Cell Type
0	H2BK12ac	H1-hESC
1	H3K4me1	NHEK
2	H3K4me2	NH-A
3	H3K9me1	K562
4	H4K20me1	NHEK
5	H2BK5ac	H1-hESC
6	H3K4me3	NH-A
7	H4K8ac	H1-hESC
8	H3K4me2	Monocytes-CD14+RO01746
9	H3K27me3	Osteoblasts
10	H3K36me3	Monocytes-CD14+RO01746
11	H3K23me2	H1-hESC
12	H3K27ac	NHLF
13	H3K36me3	NHEK
14	H2BK20ac	H1-hESC
15	H3K9ac	NHLF
16	H3K36me3	Osteoblasts
17	H2BK120ac	H1-hESC
18	H3K79me2	K562
19	H3K4me1	K562

1350
 1351
 1352
 1353
 1354
 1355
 1356
 1357
 1358
 1359
 1360
 1361
 1362
 1363
 1364
 1365
 1366
 1367
 1368
 1369
 1370
 1371
 1372
 1373
 1374
 1375
 1376
 1377
 1378
 1379
 1380
 1381
 1382
 1383
 1384
 1385
 1386
 1387
 1388
 1389
 1390
 1391
 1392
 1393
 1394
 1395
 1396
 1397
 1398
 1399
 1400
 1401
 1402
 1403

Table 10: 20 DNase tracks sub sampled for the Genomic LRB from the original 125 tracks with cell type and treatment information.

Track Index	Treatment	Cell Type
0	None	SAEC
1	None	HRPEpiC
2	None	SK-N-MC
3	None	RWPE1
4	None	Th2
5	None	Adult_CD4_Th0
6	None	HMEC
7	None	NHEK
8	UT189	Urothelia
9	None	pHTE
10	None	Urothelia
11	None	WERI-Rb-1
12	None	Huh-7
13	None	A549
14	None	Th1
15	None	HA-h
16	None	RPTEC
17	None	HMVEC-dBI-Ad
18	None	HGF
19	None	HMF

B.5 VISUALIZATION TOOL

The annotations that we join to our task datasets come from the human reference genome assembly GRCh38 (Schneider et al., 2017). To obtain these annotations we follow the methodology reported in SegmentNT (de Almeida et al., 2024) for data curation. Annotations include genomic elements, such as enhancers, exon, intron, 5' UTR, etc. The location of all gene elements and polyA signals were obtained from GENCODE (v44) (Harrow et al., 2012) gene annotation. Promoter, enhancer, and CTCF-bound sites were retrieved from ENCODE's SCREEN database (The ENCODE Project Consortium, 2020). Promoters and enhancers were split into tissue-invariant and tissue-specific annotations, following the tissue-invariant annotations from Meuleman et al. (2020). Briefly, if a promoter or enhancer overlapped at all with a region annotated as tissue-invariant, that promoter or enhancer was annotated as tissue-invariant. All other promoters and enhancers were tagged as tissue specific. Scripts from HISAT2 (Kim et al., 2019) were used to extract respective intron and splice site annotations. Annotations of repeat regions were collected from RepeatMasker (Smit et al., 2015).

Annotations were merged into the dataset by aligning chromosome and regions (start / stop position) of annotations with the genomic locations associated with the compiled tasks in the Genomics LRB. That is, if the sequence positions in our dataset overlapped with regions in the annotation files, the sequence was tagged with the corresponding annotation. For example, for variant effect prediction tasks, the SNP location was used for the merge; for regulatory element detection tasks, the start and stop positions were used. Specifically, a sample in our dataset was associated with an annotation if the sample position was both greater than the starting position of the annotation and less than the ending position of the annotation.

The UCSC liftover browser tool (Hinrichs et al., 2006) was used to convert GRCh38 annotations to the GRCh37 reference assembly locations to be associated with datasets relying on GRCh37 locations.

With annotations merged into the datasets in our Genomics LRB, we develop a visualization tool that enables users to 'slice' results. Our tool is an interactive `jupyter` (Kluyver et al., 2016) notebook that enables toggling different models and has visualizations for aggregate results, results by distance to nearest TSS / enhancer, and results by annotation. In Figure 3, we provide selected screenshots from our visualisation tool demonstrating how a user can view results for each task, select different models, and split by various annotations.

B.6 ARBITRARY SEQUENCE LENGTH

To enable users to download arbitrarily long sequence lengths, samples for each task are stored either as single positions in the genome (e.g., the SNP location for variant effect prediction or the TSS for bulk RNA expression) or as start and stop locations for tasks like regulatory element and chromatin feature prediction. In addition we store the human reference genome assemblies GRCh38 (Schneider et al., 2017) and GRCh37 (Church et al., 2011). The `PyFaidx` Python package (Shirley et al., 2015) is used to create an indexed `FASTA` file object from the reference genomes for fast random access to any subsequence. With the user's requested sequence length, we symmetrically extend sequence locations from our datasets and use these extended indices to extract the underlying DNA sequence from the indexed reference genomes. If the extended sequence indexes beyond a chromosome boundary, the sample is not returned.

1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508
1509
1510
1511

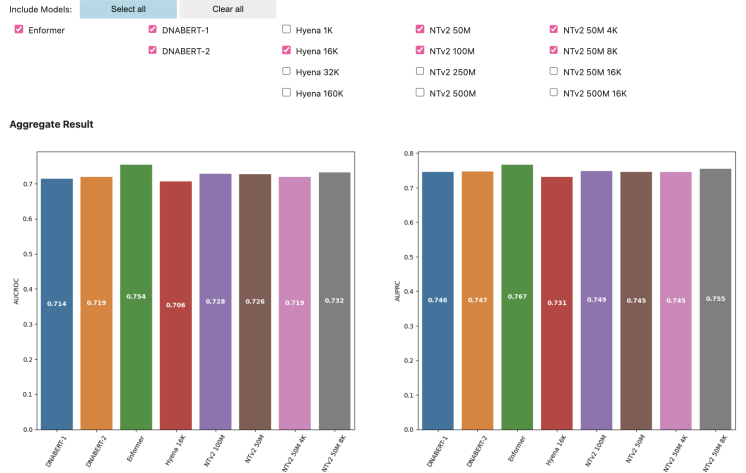
```
Choose task using dropdown and (re-)run cells below

task_dropdown = widgets.Dropdown(
    options=TASKS,
    description="Task",
)
display(task_dropdown)

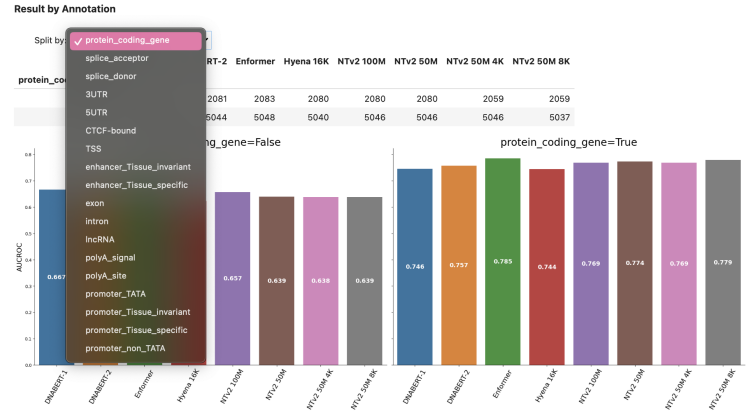
Task:
- vep_feature_causal_eaf
- vep_feature_pathogenic_clinvar
- vep_feature_pathogenic_clinvar
- chromatin_features_chrom_accessibility
- chromosome_features_h3k9me3
- regulatory_element_promoters
- regulatory_element_enhancers
- built_rna_expression

results_df = pd.read_csv('results_df.csv')
results_df = results_df_w_annotations['labels'].replace({'Pathogenic': 0, 'Common': 1}).infer_objects(copy=False)
```

(a) Screenshot of the visualization tool showing the ability to select different tasks from the Genomics LRB.



(b) Screenshot of the visualization tool showing the ability to select different models for comparison.



(c) Screenshot of the visualization tool showing the ability to select different annotations by which to split results.

Figure 3: Sample screenshots from our interactive visualization tool.

C CONTEXT LENGTH EXTENSION

Rotary Embeddings In attention-based modules, such as those used in transformer models (Vaswani et al., 2017), for a sequence of length L , the model takes embeddings in $\{\mathbf{x}\}_{j=1}^L$, $\mathbf{x}_j \in \mathbb{R}^d$, where d is the dimension of the embeddings, and computes query, key, and value vectors at every

1512 m^{th} and n^{th} position in the sequence:

$$\begin{aligned} 1513 \mathbf{q}_m &= f_q(\mathbf{x}_m, m) \\ 1514 \mathbf{k}_n &= f_k(\mathbf{x}_n, n) \\ 1515 \mathbf{v}_n &= f_v(\mathbf{x}_n, n). \end{aligned}$$

1516 f_q, f_k, f_v are query, key, and value transformations, respectively. For rotary embeddings (RoPE (Su
1517 et al., 2021)), we can think of \mathbb{R}^d as equivalent to the complex field $\mathbb{C}^{d/2}$ and define f_q and f_k as:

$$\begin{aligned} 1518 f_q(\mathbf{x}_m, m) &= e^{im\Theta} \mathbf{W}_q \mathbf{x}_m \\ 1519 f_k(\mathbf{x}_n, n) &= e^{in\Theta} \mathbf{W}_k \mathbf{x}_n, \end{aligned}$$

1520 where \mathbf{W}_q and \mathbf{W}_k are linear transformations and $\Theta = \text{diag}(\theta_1, \dots, \theta_{d/2})$ is a diagonal matrix, with
1521 $\theta_j = b^{-2j/d}$ and $b = 10000$.

1522 **RoPE Position Interpolation** In the concurrent works of Chen et al. (2023) and kaiokendev (2023),
1523 the method of position interpolation was introduced, whereby longer sequences of length $L' > L$ are
1524 accommodated by simply rescaling the position input to f_q and f_k , e.g., $f_q(\mathbf{x}_m, m \frac{L}{L'})$.

1525 **NTK-aware RoPE Interpolation** An alternative interpolation scheme, attributed to bloc97 (2023),
1526 is motivated by the hypothesis that position interpolation may lead to the loss of high frequency
1527 information. The approach that purportedly resolves this issue is related to the theory of Neural
1528 Tangent Kernels (NTK) by means of an analogy between RoPE and Fourier Features (Tancik et al.,
1529 2020), and is thus named “NTK-aware” interpolation. This scheme is characterized by a rescaling
1530 applied not to the position but rather to the basis of rotation, as follows:

$$\begin{aligned} 1531 \theta_j &= b'^{-2j/d} \\ 1532 b' &= b \cdot \left(\frac{L}{L'}\right)^{\frac{d}{d-2}} \end{aligned}$$

1533 In the experiments on context extension presented in the main text, we adopt this interpolation
1534 scheme.

1535 We note that the authors in Peng et al. (2023) further tweak and build on NTK-aware interpolation to
1536 create their proposed interpolation scheme, which they title YaRN. However, the full YaRN approach,
1537 as presented in Peng et al. (2023) requires several manually tuned hyperparameters, which were
1538 carefully selected for the decoder-only generative Llama-2 7 billion parameter model (Touvron et al.,
1539 2023a;b). We therefore adopted the simpler NTK-aware approach in our experiments.

1540 **Efficient Long-Range Context Extension** To mitigate the computational and memory costs of
1541 scaling to larger contexts, we follow the algorithm presented in Rabe & Staats (2021). This algorithm
1542 leverages a “lazy softmax” approach where key-value pairs are processed sequentially, maintaining
1543 only two vectors in memory: one for the accumulated weighted values and another for the cumulative
1544 sum of weights. This method significantly reduces memory usage by avoiding the storage of all
1545 pairwise attention scores. To optimize performance on modern hardware accelerators, which rely
1546 on parallelization for efficiency, the implementation processes attention in chunks. Rabe & Staats
1547 (2021) empirically determined that using a chunk size of \sqrt{L} strikes a balance between memory
1548 savings and computational overhead. Larger chunks increase memory requirements, while smaller
1549 chunks can lead to excessive re-computation of activations during the backward pass. Additionally,
1550 the implementation is numerically stable and functions as a drop-in replacement for the standard
1551 attention module, making it highly practical for tasks requiring extended context lengths.

1552 D ADDITIONAL EXPERIMENTAL DETAILS

1553 D.1 EVALUATED DNA LANGUAGE MODELS

1554 In Table 11, we list the DNA LMs included in the initial evaluation of our benchmark.

Table 11: Overview of **Pre-trained** DNA LMs evaluated in this study.

	Pre-training	Data	Parameters	Architecture	Context (bps)	Tokenization
NTv2	MLM	Multi-Species	50M, 100M, 250M, 500M	Transformer	12k	6-mer
DNABERT-1	MLM	Human Reference	88.6M	Transformer	512 bps	6-mer
DNABERT-2	MLM	Multi-Species	116.6M	Transformer	700 (train), up to 10k (eval)	Byte Pair Encoding
HyenaDNA	NTP	Human Reference	1.6M, .6M, 3.9M, 12.9M	SSM	1k, 16k, 32k, 160k	Single Base Pair
Caduceus	MLM	Human Reference	7.7M	SSM	131k	Single Base Pair
GPN-MSA	MLM	Human Reference + Multi-Species MSA	86M	Transformer	128	Single Base Pair

D.2 ZERO-SHOT EVALUATION

For masked DNA LMs, zero-shot scores are computed by masking the variant position in the sequence, performing inference on the masked sequence, and obtaining the probability distribution at the variant position. A score is then calculated using the probabilities of the reference allele token and the alternative allele token. For auto-regressive DNA LMs, no masking is required due to their unidirectional nature. Instead, a forward pass is done with the reference sequence, and the probability distribution is extracted from the token immediately preceding the variant position. **For Alignment Models, the human sequence is treated in the same way as the other masked DNA LMs, but the auxiliary MSA is left entirely unmasked.** Scores are computed as the log probability ratio for the reference (ref) and alternative (alt) allele tokens:

$$\text{variant effect score} = \log \left(\frac{P_{\text{ref}}}{P_{\text{alt}}} \right)$$

Details about additional processing required for zero-shot prediction are given below.

D.2.1 CAUSAL EQTL

The original dataset used for this tasks contains tissue information for each sequence. Given that zero-shot evaluate cannot account for tissue, we process variants appearing across multiple tissue types as follows: first, we find variants appearing in multiple tissues and determining a consensus label for a given variant across tissues using a 70% majority class agreement threshold. Variants appearing across multiple tissues whose majority class agreement was below this threshold were dropped. When computing metrics we only count variants appearing across tissues once.

D.2.2 PATHOGENIC-OMIM

Due to computational considerations and given that this data set totals $\sim 2.3\text{M}$ examples, we only considered a subset of the common variants for carrying out zero-shot prediction. Specifically, we sub-sampled 200k common variants and kept all 406 original pathogenic variants.

D.3 FINE-TUNING EVALUATION

To fine-tune models on our benchmark tasks, we first extracted model embeddings, in the case of DNA LMs this involves extracting the output of the last layer before the LM head, and in the case of Enformer, this involves extracting the model embeddings before the final supervised prediction head. Model embeddings were then processed in a task specific manner and subsequently fed into a task specific MLP, both of which are outlined below. We note that for Enformer, since it is a model that was originally trained in a multi-task supervised fashion and not intended to be fine-tuned, embeddings were frozen and only the prediction head was trained.

D.3.1 CAUSAL EQTL

Embedding Extraction We extract model embeddings for both the reference and alternative sequences and average embeddings across a window of size 1536 bps symmetrically around the SNP position. The mean embeddings for the reference and alternative are concatenated. Tissue information is converted to one-hot and additionally concatenated to the reference-alternative embedding vector.

MLP Head MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal to two times the base model’s embedding dimension. The MLP is composed of one linear

layer with size $2 \times$ embedding dimension, a softplus activation, another linear layer with size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for binary prediction.

Hyperparameters The parameters used to fine-tune models on this task include batch size = 64, learning rate = $1e^{-5}$, ADAM (Kingma & Ba, 2014) optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$, trained for 1 epoch on the task’s training dataset. Validation is carried out every 70 parameter update steps.

D.3.2 PATHOGENIC CLINVAR

Embedding Extraction We extract model embeddings for both the reference and alternative sequences and take a window mean of size 1536 bps symmetrically around the SNP position. The mean embeddings for the reference and alternative are concatenated together.

MLP Head MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal to two times the base model’s embedding dimension. The MLP is composed of one linear layer with size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for binary prediction.

Hyperparameters The parameters used to fine-tune models on this task include batch size = 64, learning rate = $1e^{-5}$, ADAM optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$, trained for 3 epochs on the task’s training dataset. Validation is carried out every 40 parameter update steps.

D.3.3 BULK RNA EXPRESSION

Embedding Extraction We extract model embeddings for the input sequence and take perform mean pooling on a window centered on the TSS with 383 bps before the TSS and 256 bp after.

MLP Head MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal to two times the base model’s embedding dimension. The MLP is composed of one linear layer with size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for predicting 218 regression values.

Hyperparameters The parameters used to fine-tune models on this task include batch size = 64, learning rate = $3e^{-5}$, ADAM optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$ trained for 3 epochs on the task’s training dataset. Validation is carried out every 50 parameter update steps.

D.3.4 CAGE PREDICTION

Embedding Extraction Base model embeddings were extracted and fed into the task MLP predictor.

MLP Head MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal to two times the base model’s embedding dimension. The MLP is composed of one linear layer with size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for predicting 218 regression values.

Hyperparameters The parameters used to fine-tune models on this task include batch size = 64, learning rate = $3e^{-5}$, adam optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$ trained for 1 epoch of the training dataset. Validation is carried out every 50 parameter update steps.

D.3.5 REGULATORY ELEMENTS

Due to computational considerations, we only fine-tuned models on a randomly sampled 100k subset of the full ~ 1 -2M samples in the training set. Models were evaluated on the full test dataset.

Embedding Extraction Given that the task is defined on predicting the presence of a regulatory element in the center 200 bp of the sequence, we extract a central window of 200 bps from the sequence of embeddings and perform mean pooling. This mean embedding is then passed as input to the MLP predictor head.

MLP Head MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal to two times the base model’s embedding dimension. The MLP is composed of one linear layer with size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for predicting binary values.

1674 **Hyperparameters** The parameters used to fine-tune models on this task include batch size = 64,
 1675 learning rate = $3e^{-5}$, ADAM optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$ trained for 1 epoch of
 1676 the sampled training dataset for each task. Validation is carried out every 30 parameter update steps.
 1677

1678 D.3.6 CHROMATIN FEATURES

1679 Due to computational considerations, we only fine-tuned models on a randomly sampled 100k subset
 1680 from the full $\sim 2M$ sample training set. Models were evaluated on the full test dataset.
 1681

1682 **Embedding Extraction** Given that the task is defined on predicting the presence of a chromatin
 1683 feature in the center 200 bp of the sequence, we extract a central window of 200 bps from the
 1684 sequence of embeddings and perform mean pooling. This mean embedding is then passed as input to
 1685 the MLP predictor head.

1686 **MLP Head** MLP hidden dimensions are sized in an adaptive way such the hidden state size is equal
 1687 to two times the base model’s embedding dimension. The MLP is composed of one linear layer with
 1688 size $2 \times$ embedding dimension, a softplus activation, and a final linear layer for predicting the 20
 1689 binary labels.

1690 **Hyperparameters** The parameters used to fine-tune models on this task include batch size = 64,
 1691 learning rate = $3e^{-5}$, adam optimizer with $\beta_1=0.9$, $\beta_2=0.999$, and $\epsilon=1e^{-8}$ trained for 1 epoch of
 1692 the training dataset. Validation is carried out every 30 parameter update steps.
 1693

1694 D.4 FINE-TUNING ABLATION DETAILS

1695 For the fine-tuning ablation study, we compared training only the task MLP with DNA LM embed-
 1696 dings frozen against training all DNA LM weights in conjunction with the task MLP. All training
 1697 setup details regarding embedding extraction and hyperparameters were kept constant except for
 1698 learning rate which was adjusted to account for training larger networks when full fine-tuning. The
 1699 following learning rates for each task were used in the MLP only training:
 1700

- 1701 • Variant effect prediction tasks: $1e^{-4}$
- 1702 • Bulk RNA: $2.5e^{-4}$
- 1703 • CAGE: $2e^{-4}$
- 1704 • Regulatory elements: $2.5e^{-4}$
- 1705 • Chromatin features: $2.5e^{-4}$.

1708 D.5 CONTEXT EXTENSION IMPLEMENTATION DETAILS

1709 To conduct context length extension of NTv2, we first used the 50M model due to computation
 1710 considerations. We started with the pre-trained NTv2 50M checkpoint from [Dalla-Torre et al. \(2023\)](#),
 1711 pre-trained on 12k bp sequences, and extended the context length by factors of two to 24k, 48k,
 1712 and 96k bps using a second stage of masked language modeling on a multi-species dataset from
 1713 [Dalla-Torre et al. \(2023\)](#). After proving out this methodology for the 50M model, we conducted
 1714 context length extension for the 500M model at 96k bps.
 1715

1716 **Hyperparameters** For the 50M NTv2 model we use the following hyperparameters: batch size =
 1717 1M tokens, full precision training, masking ratio = 0.15, masking probability = 0.8, random token
 1718 probability = 0.1. The ADAM optimizer with weight decay regularization was used with weight
 1719 decay = 0.01, $\beta_1=0.9$, $\beta_2=0.999$, $\epsilon=1e^{-8}$, a modified square decay learning rate schedule, with
 1720 initial learning rate of $6e^{-5}$ and end learning rate of $8e^{-4}$ with 1000 warm up steps. Training was
 1721 conducted over ~ 5 billion tokens totalling $\sim 5k$ parameter update steps.

1722 All hyperparameters were kept constant for the NTv2 500M model, however due to limited memory
 1723 resources, mixed precision training was used.
 1724

1725 D.6 SUPERVISED TRAINING BASELINES DETAILS

1726 **Convolutional Neural Network** The CNN architecture is comprised of eight 1D convolutional
 1727 blocks that use a filter size of 5 and padding that keeps input sequence length unchanged and hidden

dimension of 512. Each convolution block is composed of a convolutional layer followed by the GeLU non-linearity (Hendrycks & Gimpel, 2016) and a layer norm. After each convolutional block we also apply a fully-connected layer with GeLU non-linearity and another layer norm. Each block also contains a residual connection. This architecture is derived from the one used in Benegas et al. (2023b), but without dilation. The model consists of 12M parameters. The baseline model was trained with base-pair level tokenization and an input context size of 2,048 tokens.

Caduceus We also train an eight layer Caduceus (Schiff et al., 2024) model with hidden dimension of 256. We use the reverse complement equivariant version of this architecture (Caduceus-PS). The model consists of 3.3M parameters. The baseline model was trained with base-pair level tokenization and an input context size of 2,048 tokens.

E ADDITIONAL RESULTS

E.1 FULL DNA LM SERIES EVALUATIONS

In Tables 12 and 13 we display results for the full set of models evaluated on our benchmark.

Table 12: Benchmarking performance of DNA LMs and baselines on variant effect prediction tasks. Models were evaluated in both fine-tuning and zero-shot settings. *Extended NTv2 500 M was fine-tuned with 60k bp sequences due to compute constraints..

Model Name	Context (bp)	Causal eQTL (AUROC)		Pathogenic ClinVar (AUROC)		Pathogenic OMIM (AUPRC)
		<i>Fine-tune</i>	<i>Zero-shot</i>	<i>Fine-tune</i>	<i>Zero-shot</i>	<i>Zero-shot</i>
DNABERT 1	512	0.72 ± 0.003	0.51	0.67 ± 0.037	0.50	0.002
DNABERT 2	10k	0.72 ± 0.008	0.50	0.74 ± 0.013	0.50	0.002
DNABERT S	10k	0.73 ± 0.008	-	0.73 ± 0.011	-	-
NTv2 50M	12k	0.72 ± 0.005	0.51	0.75 ± 0.008	0.53	0.002
NTv2 100M	12k	0.73 ± 0.003	0.51	0.76 ± 0.009	0.56	0.002
NTv2 250M	12k	0.72 ± 0.003	0.51	0.78 ± 0.013	0.58	0.002
NTv2 500M	12k	0.72 ± 0.003	0.51	0.78 ± 0.009	0.68	0.003
HyenaDNA 1K	1k	0.71 ± 0.005	0.51	0.63 ± 0.027	0.49	0.002
HyenaDNA 16K	16k	0.71 ± 0.005	0.51	0.66 ± 0.016	0.49	0.002
HyenaDNA 32K	32k	0.72 ± 0.002	0.51	0.66 ± 0.012	0.50	0.002
HyenaDNA 160K	160k	0.71 ± 0.010	0.51	0.56 ± 0.073	0.49	0.002
Extended NTv2 50M 24K	24k	0.72 ± 0.004	0.51	0.75 ± 0.009	0.53	0.002
Extended NTv2 50M 48K	48k	0.73 ± 0.008	0.51	0.65 ± 0.059	0.52	0.002
Extended NTv2 50M 96K	96k	0.73 ± 0.006	0.51	0.74 ± 0.019	0.51	0.002
Extended NTv2 500M 96K*	96k	0.74 ± 0.004	0.51	0.75 ± 0.018	0.53	0.002
Baseline		0.76 ± 0.002 (Enformer)	0.56 (CADD)	0.65 ± 0.031 (Enformer)	0.97 (CADD)	0.205 (CADD)

DNABERT-2 and DNABERT-S were not fine-tuned on the CAGE task due to the incompatibility between the byte pair tokenization this model employs and binned labels used in this task. Additionally, given that DNABERT-S is trained on a contrastive learning objective and not a language modeling objective, we cannot obtain a probability distribution over the tokens that we require to compute zero-shot performance.

E.2 COMPUTATIONAL EFFICIENCY

In Table 14, we show the number of parameters for each model, and the FLOPs on an A100 80GB device with batch size 1 for each task category.

E.3 ADDITIONAL FINE-TUNING ABLATION

In Table 15, we display results for the the full NTv2 series and additional HyenaDNA models. We find that the same pattern discussed in Section 5.6 holds for this larger set of models as well. Namely, full fine-tuning almost uniformly improves model performance relative to partial fine-tuning, by

Table 13: Benchmarking performance of DNA LMs and baselines on gene expression prediction, regulatory element, and chromatin features prediction tasks. Models were evaluated in only a fine-tuned setting for this set of tasks. DNABERT-2 was not fine-tuned on the CAGE task due to the incompatibility of the byte pair tokenization with binned labels.

	Context (bp)	Bulk RNA (R^2)	CAGE (R^2)	Promoter (AUPRC)	Enhancer (AUROC)	Histone Marks (AUPRC)	DNA Accessibility (AUPRC)
		<i>Fine-tune</i>	<i>Fine-tune</i>	<i>Fine-tune</i>	<i>Fine-tune</i>	<i>Fine-tune</i>	<i>Fine-tune</i>
DNABERT-1	512	0.47 ± 0.007	0.14 ± 0.025	0.72 ± 0.009	0.80 ± 0.005	0.23 ± 0.003	0.18 ± 0.006
DNABERT-2	10k	0.51 ± 0.050	-	0.71 ± 0.112	0.81 ± 0.022	0.24 ± 0.091	0.15 ± 0.064
DNABERT-S	10k	0.52 ± 0.060	-	0.75 ± 0.021	0.83 ± 0.005	0.33 ± 0.006	0.16 ± 0.039
NTv2 50M	12k	0.52 ± 0.074	0.35 ± 0.030	0.75 ± 0.008	0.78 ± 0.041	0.34 ± 0.007	0.18 ± 0.005
NTv2 100M	12k	0.52 ± 0.081	0.3 ± 0.030	0.78 ± 0.008	0.82 ± 0.010	0.34 ± 0.007	0.22 ± 0.012
NTv2 250M	12k	0.57 ± 0.024	0.37 ± 0.008	0.8 ± 0.008	0.84 ± 0.002	0.37 ± 0.013	0.28 ± 0.006
NTv2 500M	12k	0.60 ± 0.038	0.39 ± 0.011	0.79 ± 0.006	0.82 ± 0.002	0.38 ± 0.003	0.3 ± 0.007
HyenaDNA 1K	1k	0.44 ± 0.014	0.11 ± 0.015	0.7 ± 0.006	0.80 ± 0.002	0.21 ± 0.001	0.13 ± 0.003
HyenaDNA 16K	16k	0.46 ± 0.008	0.17 ± 0.014	0.64 ± 0.004	0.75 ± 0.002	0.22 ± 0.002	0.091 ± 0.003
HyenaDNA 32K	32k	0.41 ± 0.012	0.22 ± 0.007	0.56 ± 0.008	0.73 ± 0.001	0.22 ± 0.003	0.084 ± 0.001
HyenaDNA 160K	160k	0.46 ± 0.006	0.19 ± 0.032	0.67 ± 0.009	0.74 ± 0.009	0.25 ± 0.004	0.11 ± 0.002
Extended NTv2 50M 24K	24k	0.53 ± 0.063	0.37 ± 0.010	0.75 ± 0.007	0.83 ± 0.002	0.35 ± 0.007	0.19 ± 0.006
Extended NTv2 50M 48K	48k	0.54 ± 0.038	0.36 ± 0.012	0.76 ± 0.008	0.82 ± 0.002	0.35 ± 0.007	0.19 ± 0.006
Extended NTv2 50M 96K	96k	0.54 ± 0.034	0.3 ± 0.019	0.76 ± 0.015	0.83 ± 0.001	0.35 ± 0.005	0.19 ± 0.007
Extended NTv2 500M 96K	96k	0.56 ± 0.037	0.36 ± 0.011	0.78 ± 0.003	0.82 ± 0.005	0.38 ± 0.004	0.3 ± 0.006
Baseline		0.80 ± 0.010 (Enformer)	0.49 ± 0.000 (Enformer)	0.86 ± 0.006 (Enformer)	0.92 ± 0.002 (Enformer)	0.35 (DeepSea)	0.44 (DeepSea)

Table 14: Model sizes and FLOPs used per task type.

Tasks	Hyena 1k (0.6M)	Hyena 16k (1.6M)	Hyena 32k (3.9M)	Hyena 160k (12.9M)	DNABERT-1 (88.6M)	DNABERT-2 (116.6M)	NTv2 50M (50M)	NTv2 100M (100M)	NTv2 250M (250M)	NTv2 500M (500M)
Variant Effect	0.45B	4.62B	38.84B	420.12B	2.27B	31.37B	18.89B	34.60B	38.92B	93.12B
Gene Expression	0.27B	2.80B	23.65B	256.84B	1.11B	17.16B	9.44B	17.29B	19.45B	46.81B
Regulatory Element	0.27B	2.80B	23.65B	256.84B	1.10B	17.16B	9.44B	17.29B	19.45B	46.80B
Chromatin Features	0.27B	2.80B	23.65B	256.84B	1.10B	17.16B	9.44B	17.29B	19.45B	46.80B

margins that can range up to > 100%. Tasks on which DNA LMs already perform competitively, e.g., regulatory element annotation, seem to benefit less from full-fine tuning, but even here we do see gains. In Table 16 we perform a sensitivity analysis analyze the robustness of our fine-tuning setup to multiple hyperparameter settings, namely for learning rate and batch size. Our analysis shows that while most results are quite insensitive to hyperparameter choice (with swings ± 0.02 on the metric of interest), users should avoid combinations of higher learning rates ($3e-5$) and smaller batch sizes (32).

Table 15: Ablation study examining the difference in performance of DNA LM fine-tuning strategies. Results shown correspond to the percent increase in performance of full fine-tuning with respect to freezing LM weights and only training the MLP head.

	Causal eQTL (AUCROC)	Pathogenic ClinVar (AUROC)	Bulk RNA (R^2)	CAGE (R^2)	Promoter (AUPRC)	Enhancer (AUROC)	Histone Marks (AUCPRC)	DNA Accessibility (AUPRC)
NTv2 50M	+1.13	+9.30	+30.23	+71.60	+1.93	-2.05	+32.03	+33.43
NTv2 100M	+0.98	+6.24	+13.70	+27.72	+2.16	+2.83	+32.70	+40.54
NTv2 250M	+0.36	+3.57	+21.70	+40.41	+2.07	+3.71	+31.01	+54.44
NTv2 500M	+0.49	+4.27	+24.45	+42.14	-1.45	+0.90	+22.46	+47.96
HyenaDNA 1K	+0.95	+15.39	+16.50	+45.22	+7.13	+4.68	+23.61	+22.65
HyenaDNA 16K	+0.21	+22.81	+75.53	+133.52	+6.19	-1.10	+42.83	-9.62
HyenaDNA 32K	+0.35	+11.58	+82.46	+102.91	-18.21	-6.02	+14.43	-22.67

E.4 ADDITIONAL RESULTS BY GENOMIC ANNOTATIONS

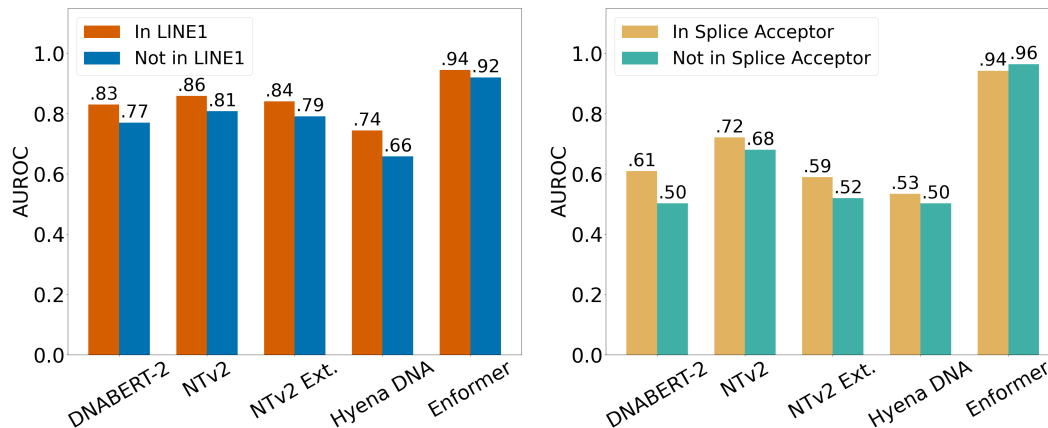
In Figure 4, we display additional results from splitting the tasks by genomic annotations.

Table 16: Fine-Tuning sensitivity analysis on LR and Batch size for Causal eQTL and Bulk RNA tasks.

Model	LR	Batch size	Causal eQTL (AUCROC)	Bulk RNA (R^2)
NTv2 500M	$1e^{-5}$	32	0.723 ± 0.006	0.597 ± 0.050
NTv2 500M	$1e^{-5}$	64	0.722 ± 0.003	0.588 ± 0.048
NTv2 500M	$1e^{-5}$	128	0.718 ± 0.010	0.596 ± 0.015
NTv2 500M	$3e^{-5}$	32	0.717 ± 0.006	0.580 ± 0.079
NTv2 500M	$3e^{-5}$	64	0.717 ± 0.007	0.566 ± 0.016
NTv2 500M	$3e^{-5}$	128	0.721 ± 0.006	0.585 ± 0.047
DNABERT 2	$1e^{-5}$	32	0.726 ± 0.005	0.483 ± 0.135
DNABERT 2	$1e^{-5}$	64	0.719 ± 0.008	0.503 ± 0.068
DNABERT 2	$1e^{-5}$	128	0.725 ± 0.002	0.484 ± 0.085
DNABERT 2	$3e^{-5}$	32	0.687 ± 0.067	0.480 ± 0.063
DNABERT 2	$3e^{-5}$	64	0.713 ± 0.016	0.507 ± 0.050
DNABERT 2	$3e^{-5}$	128	0.720 ± 0.005	0.501 ± 0.055
Hyena DNA 160K	$1e^{-5}$	32	0.703 ± 0.016	0.459 ± 0.010
Hyena DNA 160K	$1e^{-5}$	64	0.708 ± 0.010	0.450 ± 0.006
Hyena DNA 160K	$1e^{-5}$	128	0.708 ± 0.012	0.439 ± 0.016
Hyena DNA 160K	$3e^{-5}$	32	0.701 ± 0.006	0.456 ± 0.018
Hyena DNA 160K	$3e^{-5}$	64	0.699 ± 0.010	0.457 ± 0.006
Hyena DNA 160K	$3e^{-5}$	128	0.696 ± 0.011	0.445 ± 0.020

Enhancer Detection We find that DNA LMs have increased performance at identifying enhancers in some repetitive elements, such as LINE1 transposons, as shown in Figure 4a. LINE1 elements are commonly interspersed along the human genome, and individual LINE1 elements may have uncertain regulatory effects, but DNA LMs appear to be able to call enhancers in LINE1 elements better than in non-LINE1 regions. However, their performance still lags that of the Enformer baseline.

Zero-shot Pathogenic-ClinVar In Figure 4b, we observe that most models exhibit increased performance within splice site acceptor regions, with the exception of Enformer, although Enformer demonstrates high performance in both splits.



(a) Enhancer detection; split by enhancers located within a LINE1 (transposon) annotation.

(b) Zero-shot Pathogenic ClinVar prediction; by splice site acceptor annotation.

Figure 4: Additional results split by genomic annotations.

1890 F POTENTIAL SOCIETAL IMPACTS

1891

1892 As our work introduces a benchmark, we do not believe it poses any inherent negative societal
 1893 impacts. In fact, our work will hopefully create a positive impact by accelerating the development of
 1894 useful DNA LMs that can bring about a deeper understanding of biology.

1895

1896 G ASSETS

1897

1898 In Table 17, we list the open source libraries and repositories used in this work, with their correspond-
 1899 ing licenses.

1900

1901

1902 Table 17: Open source libraries (and corresponding licenses) used in this work.

1903

1904	Library	License
1905	Biopython (Cock et al., 2009)	Biopython license
1906	Haiku (Hennigan et al., 2020)	Apache 2.0
1907	HuggingFace (Wolf et al., 2019)	Apache 2.0
1908	Jax (Bradbury et al., 2018)	Apache 2.0
1909	Jupyter (Kluyver et al., 2016)	BSD 3-Clause
1910	NumPy (Harris et al., 2020)	NumPy license
1911	Matplotlib (Hunter, 2007)	Matplotlib license
1912	Pandas (The pandas development team, 2020)	BSD 3-Clause "New" or "Revised"
1913	Optax (DeepMind et al., 2020)	Apache 2.0
1914	PyFaidx (Shirley et al., 2015)	BSD-3-Clause
1915	PyTorch (Paszke et al., 2019)	BSD-3 Clause
1916	Scikit-Learn (Pedregosa et al., 2011)	BSD 3-Clause
1917	Seaborn (Waskom, 2021)	BSD 3-Clause "New" or "Revised"
1918	TensorFlow (Abadi et al., 2015)	Apache 2.0

1919

1920

1921 H COMPUTATIONAL RESOURCES

1922

1923 All research in this study was conducted using Cloud TPU’s provided by Google’s TPU Research
 1924 Cloud program. Specifically, a TPU-v4-64 slice was used for all context length extension pre-training.
 1925 Single TPU-v4 machines were used in parallel to conduct all benchmarking and evaluations including
 1926 fine-tuning, zero-shot, and inference experiments.

1927

1928

1929

1930

1931

1932

1933

1934

1935

1936

1937

1938

1939

1940

1941

1942

1943