
Orthrus: Towards Evolutionary and Functional RNA Foundation Models

Philip Fradkin^{1,2,*}, Ruain Shi^{1,2,3*}, Keren Isaev^{4,5},
Brendan Frey^{1,2,5}, Quaid Morris^{3,†}, Leo J. Lee^{1,6,†}, Bo Wang^{1,2,7,†}

¹ Vector Institute, Ontario, Canada.

² Computer Science, University of Toronto, Ontario Canada.

³ Computational and Systems Biology Program, Sloan Kettering Institute, New York, United States.

⁴ New York Genome Center, New York, United States.

⁵ Systems Biology, Columbia University, New York, United States.

⁶ Electrical and Computer Engineering, University of Toronto, Ontario Canada.

⁷ Peter Munk Cardiac Center, University Health Network, Ontario, Canada.

Abstract

In the face of rapidly accumulating genomic data, our understanding of the RNA regulatory code remains incomplete. Pre-trained genomic foundation models offer an avenue to adapt learned RNA representations to biological prediction tasks. However, existing genomic foundation models are trained using strategies borrowed from textual or visual domains, such as masked language modelling or next token prediction, that do not leverage biological domain knowledge. Here, we introduce Orthrus, a mamba based RNA foundation model pre-trained using a novel self-supervised contrastive learning objective with biological augmentations. Orthrus is trained by maximizing embedding similarity between curated pairs of RNA transcripts, where pairs are formed from splice isoforms of 10 model organisms and transcripts from orthologous genes in 400+ mammalian species from the Zoonomia Project. This training objective results in a latent representation that clusters RNA sequences with functional and evolutionary similarities. We find that the generalized mature RNA isoform representations learned by Orthrus significantly outperform existing genomic foundation models on five mRNA property prediction tasks, and requires only a fraction of fine-tuning data to do so.

1 Introduction

Mature RNAs, resulting from transcription and alternative splicing of precursor RNAs, encode essential genetic information for protein synthesis. The regulation of precursor RNAs is often tightly linked to their sequence and is critical in modulating protein expression and cellular functions [65]. Experimental procedures such as eCLIP, ribosome profiling, or SLAM seq have been pivotal in studying these RNA regulatory processes, but these techniques are often time-consuming and expensive [60, 6, 24]. As an alternative, supervised machine learning models trained on genetic sequences provide data-driven modelling of RNA regulation, offering effective and low-cost prediction of cellular processes such as alternative splicing and RNA degradation [25, 33, 1]. These models can be used to identify disease mechanisms [42, 48], improve therapeutics such as messenger RNA (mRNA) vaccines [7], and predict the effects of perturbations [35]. Despite the importance of these

*These authors contributed equally to this work.

†Equal advising.

applications, the difficulty associated with experimental acquisition of training data restricts the use of supervised methods for a wider range of tasks.

Several recent works [26, 8, 43] have proposed foundation models as an alternative for supervised learning approaches in genomic domains. Genomic foundation models use deep neural networks to learn an expressive representation of genetic sequences by pre-training on large datasets. During pre-training, self-supervised learning (SSL) objectives are used to train the model in the absence of labeled examples. SSL can be formulated through a data reconstruction objective, where a model is required to reconstruct a portion of the input data. Existing genomic foundation models use training objectives including next token prediction (NTP) and masked language modeling (MLM) [13, 46, 26]. Foundational models that effectively capture the underlying biological complexities enable few-shot learning, generalizing experimental biology using a minimal number of samples [70, 53]. Representations learned with foundation model techniques can be fine-tuned on related downstream tasks with fewer labeled data points, reducing reliance on data collection and demonstrating impressive generalization capabilities to a diversity of tasks [58, 47]. However, the unique properties inherent to genomic data pose challenges for implementing reconstruction-based SSL objectives or supervised learning approaches.

Genomic sequences in the natural world are constrained by evolutionary viability, resulting in low natural diversity³ and high mutual information across genomes from the same species [56]. Latest estimates propose that approximately five percent of the human genome is under constraint and can be considered high information content [9, 32]. The remaining 95% of the genetic sequence lacks evidence of negative selection, meaning mutations may have little to no impact on organism fitness [55]. Without a strong biological inductive bias, existing reconstruction-based SSL models often reconstruct non-informative tokens, which can result in sub-optimal representations. Due to the high-mutual information between samples, it is also difficult to scale the effective size of the training dataset to circumvent this issue. As we later show, applications of SSL methods to genomics learn latent representations that are not well suited for RNA property prediction [12, 26, 43, 8]. The gap between baseline SSL methods and supervised approaches remains large, while no clear trend exists between model size and performance.

Here, we propose Orthrus, an RNA foundation model that is pre-trained on mature RNA sequences. Orthrus uses a novel biologically motivated contrastive learning objective to structure the model latent space by maximizing similarity between splicing isoforms and evolutionary related transcripts [36, 10]. Using this contrastive objective, Orthrus is pre-trained on splicing annotation data from 10 species and orthologous alignments from 400 mammalian species in Project Zoonomia [28]. Pre-training Orthrus on mature RNAs with high functional importance and sequence conservation [9, 55] further allows Orthrus to focus on sequence regions with high information content. Orthrus pre-training results in effective mature RNA representations that are predictive of diverse RNA properties.

We show that Orthrus’s learned representations can be used to accurately predict the properties of mature mammalian RNA sequences in three key contexts. First, we test the effectiveness of biologically inspired contrastive learning by fitting a linear model on top of the pre-trained latent representations. We identify that Orthrus outperforms other self-supervised foundation models, and applying this simple linear transformation approaches the performance of supervised methods on all property prediction tasks. Second, we fine-tune the pre-trained models on experimentally collected RNA property datasets and demonstrate state-of-the-art performance when generalizing to unseen sequences. Orthrus is able to effectively perform in the low data regime, requiring as few as 45 labeled examples to fine-tune an RNA half-life predictor. Finally, we identify that increasing the model size improves performance, opening up the door for further improvements by scaling both the training dataset and model size.

2 Methods

Contrastive learning has been shown to be a bound on mutual information between two random variables X and Y corresponding to $I(X; Y) = \mathbb{E}_{p(x,y)} \left[\log \frac{p(x,y)}{p(x)p(y)} \right]$. We utilize a variation of the classical InfoNCE loss, $\mathbb{E} \left[\log \frac{\exp(f(x_i, y_i))}{\sum \exp(f(x_i, y_j))} \right]$, where a model f is tasked with classifying the

³In the coding region (2% of human DNA), an average individual carries 27 ± 13 unique SNPs [23].

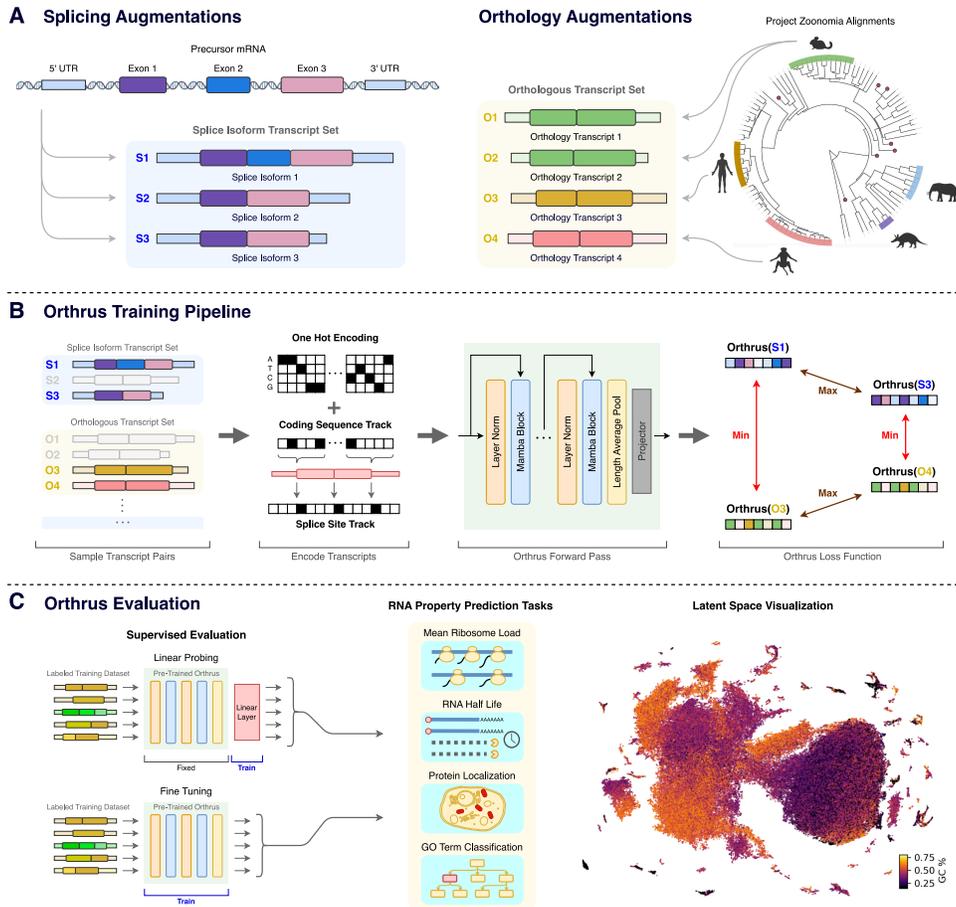


Figure 1: Description of RNA augmentations, Orthrus training, and evaluation procedures. **(A)** Demonstrates splicing and orthology pairs identified for a particular RNA sequence. A mature mRNA sequence can have multiple splicing and orthology augmentations. **(B)** The Orthrus training pipeline, consisting of first sampling a pair of positive RNA sequences and encoding them into six track encoding. We then generate a projection of the sequences using Orthrus model and apply the contrastive loss over these samples, maximizing similarity between positive pairs while minimizing it for all the other transcripts. **(C)** Orthrus evaluation consisting of linear probing, fine-tuning over a variety of mature RNA properties and visualizing the model latent space.

correct y_i which was jointly drawn with x_i [59]. Herein, the observations x_i, y_i correspond to splice isoforms or orthologous sequences which are interpreted as functionally related while f is a neural network that we optimize to minimize the loss.

We propose to use 4 different augmentations and thoroughly investigate their impact on downstream tasks. They include: alternatively spliced transcripts across ten species, orthologous transcripts identified from the Zoonomia project including over 400 species, naive orthology informed by gene identity, and masking a percentage of the input sequence Figure 1 [44, 28].

In the following section we elaborate on dataset construction, model choice, contrastive learning objective, and downstream evaluations.

Splicing and Orthology Contrastive Dataset

In the RNA domain, contrastive learning strategies have had significant success by identifying augmentations that do not have a strong semantic effect, such as cropping, rotation, or Gaussian blur [68, 69, 10]. In this work, we use RNA splicing isoforms and orthologous transcripts as sources of functional similarity [16, 44, 28]. By sampling RNA isoform sequences produced by alternative

Contrastive Dataset	Zoonomia	Splicing	# of Pairs	# of Transcripts
Zoonomia Eutheria & Splicing Gencode Basic	✓	✓	876,871,640	49,493,993
Zoonomia Eutheria	✓	✗	157,975,815	41,562,358
Splicing Gencode Basic	✗	✓	16,249,112	771,105
None	✗	✗	0	771,105

Table 1: Overview of contrastive datasets.

splicing and speciation processes, we identify sequence variation that is likely to maintain core functional properties. In addition, we use naive orthology to pool RNA transcripts from evolutionarily related genes [45]. By minimizing the distance between functionally similar sequences, the model can learn regulatory regions critical for RNA property and function prediction.

To construct positive pairs based on alternative splicing, we group alternatively spliced transcripts using GENCODE and RefSeq databases depending on availability [16, 44]. We utilize splice information across 10 species, covering a broad range across the evolutionary tree: human, mouse, chicken, *C. elegans*, chimpanzee, cow, dog, drosophila, rat, and zebrafish. In addition, we make use of naive orthology for positive pair generation: for cases where gene names are consistent across species, we pool the transcripts generated by alternative splicing into the same transcript set (Figure 1A). Alternatively spliced mRNA isoforms exhibit variability in UTR and coding sequences composition, at times demonstrating novel function. However, our work is based on the assumption that on average splice isoforms are more functionally similar to one another than a randomly sampled mRNA transcripts. We empirically find that sequence diversity present in alternatively spliced isoforms is an effective source of function preserving variation.

Orthologous transcripts from mammalian species present another source of sequence diversity, generated by genetic drift post speciation events [29, 50]. We utilize positive pairs generated through the process of speciation across the Eutheria clade through the Zoonomia TOGA resource, which performs joint gene annotation and orthology inference mapping transcripts from over 400 transcripts to human and mouse annotations [28]. To identify orthologous pairs, TOGA performs alignment over identified coding sequences and neighboring intronic and intergenic regions. We hypothesize that using orthologous sequencing as positive pairs in our dataset can allow the model to learn mRNA regions that are under negative selection. These regions will be preserved over evolutionary time due to negative selection. These regions in turn are likely to be functionally important, and relevant for mRNA property prediction.

Overall, our final dataset contains 49 million unique transcripts and over 870 million unique positive pairs (Table 1, Section 2).

For mRNA sequence representation we generate a six-track mature RNA representation, consisting of four one-hot encoded tracks encoding genomic sequence, a track indicating the 5' location of splice sites, and a track indicating the first nucleotide of every codon. The addition of splice site and coding sequence locations has been shown to be beneficial for mRNA property prediction tasks [1].

To sample positive pairs from the orthology, splicing dataset, we first identify the set of all positive samples \mathbf{Y}_j for a reference transcript x_j . \mathbf{Y}_j can be variable in length since some transcripts will have a greater number of splice isoforms and orthologous sequences than others. During a forward model pass, we sample y_j^k from \mathbf{Y}_j and use that as a positive pair for x_j .

Mamba Encoder

We pre-train a mamba state space model, which has been demonstrated to be successful in applications with long context requirements [21, 43]. mRNA sequences can reach over 12,000 nucleotides in length, making transformer architecture challenging due to its quadratic scaling in memory with sequence length [61]. Mamba, an extension of state space model families or S4, maps a sequence $x(t) \in \mathbb{R}$ to $y(t) \in \mathbb{R}$ using a latent state $h(t) \in \mathbb{R}^N$ [22].

A fundamental trade-off in architecture choice for sequence modeling is avoiding compressing sequence context and compute requirements. Transformers are able to avoid compressing context, leading to better performance, but trade-off slower training and higher memory usage [61, 21]. Al-

ternatively, S4 models define a sequence to sequence transformation parameterized by $(\mathbf{A}, \mathbf{B}, \mathbf{C}, \Delta)$. The fundamental operation consists of iteratively updating the hidden state:

$$\begin{aligned} h'(t) &= \mathbf{A} h(t) + \mathbf{B} x(t) \\ y(t) &= \mathbf{C} h(t). \end{aligned}$$

Δ is used to discretize the input for discrete domains such as natural language, or genomics. The mamba architecture improves over the S4 family of models by introducing selectivity over input by making B, C , and Δ a function of the input, resulting in

$$\begin{aligned} h'(t) &= \mathbf{A} h_t + \mathbf{B}(\mathbf{x}_t) x_t \\ y(t) &= \mathbf{C}(\mathbf{x}_t) h_t. \end{aligned}$$

Allowing parameters to be input dependent introduces desirable modeling qualities for genomic domain: variable spacing, filtering context, and linear memory scaling with sequence length $\mathcal{O}(n)$. Variable spacing refers to mamba’s ability to effectively perform on the selective copying task, where important elements are arbitrarily spaced [21]. Binding motifs in genomic sequences can be spaced without a constant offset, requiring the model to be able to learn motif interactions with variable spacing [19]. The non-uniformity of signal informativeness in genomic sequences requires models to be able to filter out irrelevant context [21]. Finally, the limited context, as opposed to transformer models, allows the mamba architecture to scale required memory linearly with increased input length [21, 61].

DCL Contrastive Learning Objective

During the contrastive training phase, we sample positive pair sequences from mature RNA transcript sets and maximize their similarity in the model latent space (Figure 1 B). Given a batch of N sequences, x^1, \dots, x^N let x_i^1, x_i^2 be a positive pair of sequences sampled from a transcript set. These sequences are related through alternative splicing or orthology processes described in section. We pass these positive pairs through a Mamba [21] encoder, f_θ resulting in the outputs h_i^1 and h_i^2 . These representations are then fed into a multi-layer perceptron projection head, g_θ the output of which is used to calculate normalized projections, z_i . We use decoupled contrastive learning (DCL) loss to perform the contrastive learning objective, pushing apart unpaired transcripts and maximizing the cosine similarity between positive pairs (Figure 1 B) [66].

We use decoupled contrastive learning (DCL) as it has been shown to require smaller batch sizes, is less sensitive to hyperparameters such as learning rate, and the positive loss term can be weighted by sample difficulty [66]. DCL iterates on the normalized temperature-scaled cross-entropy loss by splitting the contrastive objective into two terms: a similarity loss (positive) and a dissimilarity loss (negative) [52]. More formally, the positive and negative losses for sample i are calculated:

$$\mathcal{L}_{DCL,i}(\theta) = \log \left[\sum_{k=1}^N, \sum_{l=1}^2 \mathbb{1}_{k \neq i} \exp(\langle z_i^1, z_k^l \rangle / \tau) \right] - w_i \langle z_i^1, z_i^2 \rangle / \tau. \quad (1)$$

In the above z^1 and z^2 correspond to two embeddings of related sequences, z_k are embeddings from unrelated RNA sequences, τ is the temperature parameter set to 0.1, and $\mathbb{1}_{k \neq i}$ is an indicator function that evaluates to 1 when $k \neq i$. The above loss is computed for all the samples in the batch for both the sampled views $l \in 1, 2$. N corresponds to all the negative samples in batch, thus maximizing batch size during contrastive learning typically leads to improved performance.

Normalized projections z_i are outputs from the MLP projector g and are used to compute the contrastive loss, utilizing samples from the rest of the batch as negative examples:

$$z_i^1 = \frac{g(h_i^1)}{\|g(h_i^1)\|} \text{ and } z_i^2 = \frac{g(h_i^2)}{\|g(h_i^2)\|}. \quad (2)$$

For downstream RNA property evaluations, the projector g_θ is discarded and outputs from f_θ are used instead. This practice is consistent with prior literature [10, 3, 4, 20].

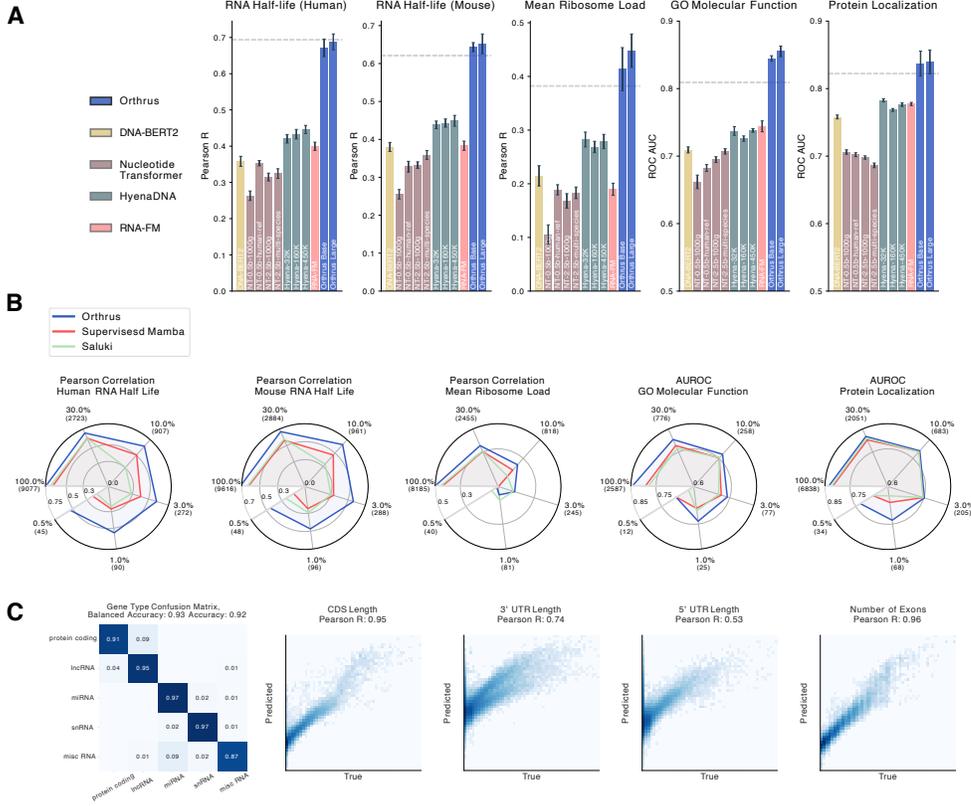


Figure 2: **(A)** Benchmarking linear probing performance on RNA property prediction tasks for self-supervised genomic foundation models. Individual bars represent the performance of foundation model variants, which typically differ in parameter count and pre-training dataset. Error bars show 95% confidence intervals, constructed using 10 runs with randomized data splits. The grey dashed line indicates the performance of the fully supervised Saluki method trained with access to labels. **(B)** Plots evaluating the fine-tuning performance of Orthrus Base across varying data availability. Each dataset is subsampled to the indicated percentage, with the number of data points provided in brackets. Point estimates are plotted, averaged across three random seeds and random data splits. **(C)** Evaluation of Orthrus’s latent representation by fitting a linear model to predict structural properties. The confusion matrix evaluates Orthrus’s ability to classify transcript types using logistic regression on learned embeddings. The four scatter plots assess Orthrus’s ability to predict structural RNA properties, including CDS length, 3’ UTR length, 5’ UTR length, and number of exons.

We introduce two versions of Orthrus using a backbone Mamba encoder: base consisting of 1.3 million trainable parameters and large with 10.1 million trainable parameters (excluding g_θ) Section 2.

3 Results

Orthrus embeddings are predictive of diverse phenotypes

To evaluate the effectiveness of our pre-trained representations, we followed the conventional evaluation strategy of linear probing. The learned latent embedding is effective if $\exists \mathbf{w}$ s.t. $\mathbf{w}^T \mathbf{X} + b = \hat{y}$, where, \mathbf{X} is a matrix of embeddings and \hat{y} approximates y . To evaluate the above, we freeze the weights of the mamba encoder f and train a linear layer to predict labels for regression and classification tasks. Further experimental details are described in Appendix A.2.

We quantitatively evaluate whether Orthrus embeddings contain information regarding key biochemical properties such as UTR length, number of exons, CDS length, and gene type in Figure 2C. We observe that Orthrus fixed length embeddings are highly predictive of RNA biochemical attributes, which are important for predicting functional RNA properties such as RNA half-life [1]. In figure 2A, we demonstrate that Orthrus outperforms other self-supervised methods on a diverse set of func-

tional RNA property prediction tasks by a substantial margin. For RNA half-life (Human) Orthrus Large outperforms other self-supervised methods, the closest of which achieves 65% of linear probing performance (Pearson R 0.45 & 0.69). Further, we evaluate a base 4 track model and find that Orthrus outperforms other self-supervised baselines (Table 4). We note that, Orthrus outperforms a supervised baseline for most tasks, which is indicated by a dashed line (Figure 2A). These results indicate that a linear regression trained with Orthrus embeddings can match or outperform neural networks tuned for RNA property prediction tasks [1].

We observe improved linear probing results as we scale the number of trainable parameters for Orthrus by comparing Base and Large model variants (Figure 2). We see an especially clear improvement trend in MRL and GO Molecular Function predictions. We note that for other self-supervised models such as Hyena DNA or Nucleotide Transformer, the number of parameters does not consistently improve performance (Figure 2 A) [43]. However, we do observe an improvement in performance for Nucleotide Transformer when comparing their 2.5 billion parameter model trained on 1000 genomes data versus multi species [12]. This is additional evidence that utilizing evolutionary information, can help improve model performance on RNA property prediction tasks.

Fine-tuning Orthrus for state-of-the-art RNA property prediction

To assess whether the Orthrus pre-training objective provides utility beyond an effective representation, we evaluate its performance by fully fine-tuning it and comparing it to a supervised model with a matched architecture. We compare its performance against a published method for the RNA half-life prediction, Saluki and find that the fully fine-tuned Orthrus model outperforms Saluki on the RNA half-life task (Figure 2) [1]. Furthermore, we retrain the Saluki architecture, train an architecturally equivalent model to Orthrus, and fine-tune pre-trained HeynaDNA model for other sequence property prediction tasks and identify that Orthrus has a significant performance advantage (Figures 2, 5). Other baseline SSL methods such as DNA-BERT2 and RNA-FM have limited input context windows, and cannot be easily applied to these tasks.

To simulate downstream tasks for which there is a lack of experimental data, we perform fine-tuning on RNA half-life prediction where only a subset of the original training data set is available. We observe that supervised methods are ineffective in this regime, while Orthrus maintains competitive performance at 10% and 1% of the data (Figure 2 B). The performance differences are even more stark when using only 0.5% of the training data, achieving 73% of supervised performance with just 45 observed samples on the human RNA half-life dataset (R=0.72 & R=0.53). These findings illustrate that Orthrus advances towards the aim of few-shot learning for downstream tasks where experimental data is scarce.

Orthrus latent space captures known functional transcript diversity

A key question in alternative splicing research is how much functional diversity RNA isoforms generate within a gene [51]. To explore whether Orthrus embeddings can help elucidate this, we analyze intra-gene isoform similarities. For each pair of transcript isoforms within protein-coding genes, we compute their similarity using Orthrus embeddings (Figure 3 A). As a control, we compare these with transcript pairs from random genes, expecting lower similarity. We also hypothesized that transcript pairs from genes sharing the same GO terms would be more similar than random pairs, but less similar than most intra-gene pairs. Our analysis confirms significant differences across all pairwise comparisons of the three groups ($p < 2.2e-16$, Mann-Whitney U test), indicating that the Orthrus training objective preserves within gene sequence diversity (Figure 3 B). Notably, we observe an overlap between intra-gene and inter-gene similarities, indicating that some alternatively spliced transcripts have distinct embeddings, RNA properties, or functional differences in protein products. As such, *within gene* diversity could potentially help delineate differential isoform protein functions, an active area of research.

To investigate whether intra-gene similarities may reflect underlying protein domain conservation, we annotated each transcript, identifying a list of included protein domains. We found that transcripts with a high similarity, as measured by overlap in the present protein domains, also have highly similar Orthrus embeddings. The correlation between Orthrus and domain similarities are significantly higher than a transcript length difference baseline, indicating that Orthrus learns functional differences as captured by domain presence (Figure 3 C). This suggests that Orthrus embeddings encode functionally relevant information.

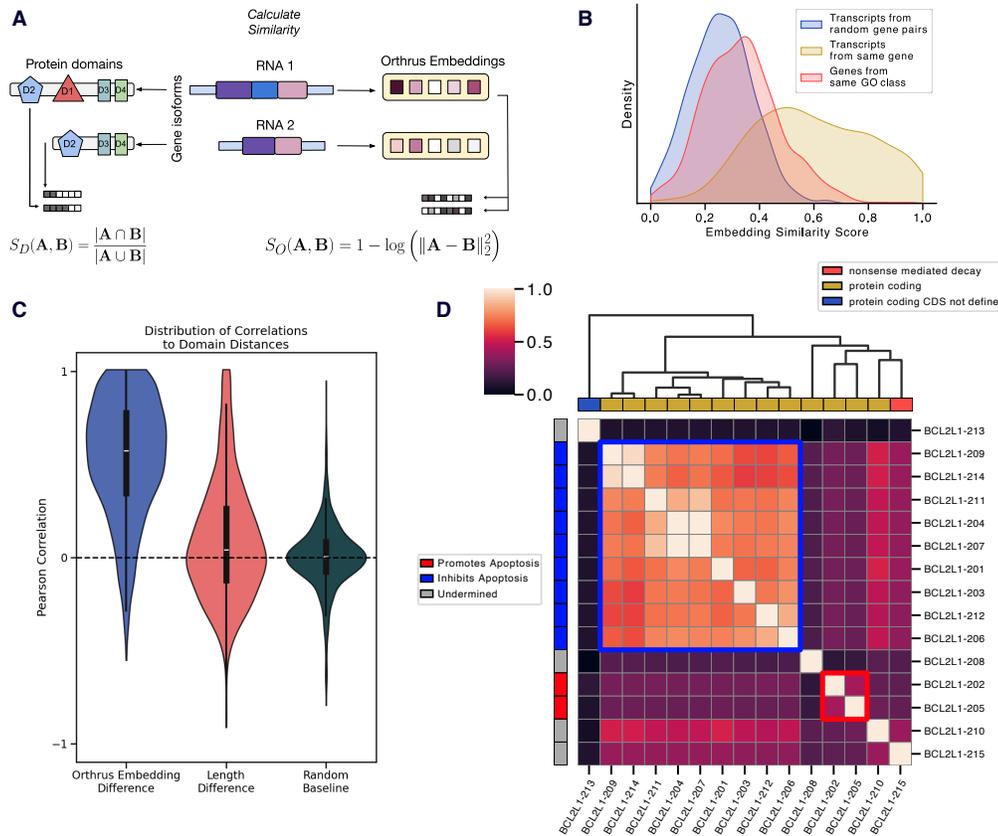


Figure 3: **(A)** Methodology for comparing gene isoform similarities using Orthrus embeddings and protein domain annotations. Orthrus embeddings for transcripts within the same gene are compared using log of L2 distance, while protein domain similarities are computed using the Jaccard index. **(B)** We visualize within gene similarities (yellow), between gene similarities (blue), and similarities of genes from the same GO class (red). **(C)** Visualization of Pearson R distributions correlating protein domain similarities with Orthrus embedding similarities for 1000 randomly sampled genes with multiple isoforms. Also plotted are the distributions of transcript length and protein domain similarities. **(D)** Example of *BCL2L1* isoforms, where apoptosis-inhibiting isoforms cluster together, while non-coding and apoptosis-inducing isoforms display low similarity. The clustering matrix, derived from Orthrus embedding similarities, is represented by the dendrogram. Boundaries highlight clusters with divergent transcript functions.

To illustrate this, we examine the *BCL2L1* gene, known for its alternatively spliced isoforms with distinct functional outcomes [64, 34]. The dominant isoforms encode an apoptosis-inhibiting protein, Bcl-X(L), while a minority encode a pro-apoptotic protein, Bcl-X(S). By clustering *BCL2L1* RNA isoforms using Orthrus embedding similarity, we identify two main functional groups: one containing *BCL2L1-202* and *BCL2L1-205*, distinct from the apoptosis-inhibiting transcripts cluster (Figure 3 D). This demonstrates that Orthrus embeddings may serve as a valuable resource for identifying isoforms with likely different functional properties, a critical area in alternative splicing research.

Ablations: Orthology and splicing and saturate performance

Finally, we investigate the Orthrus augmentations that contribute towards effective performance. We find that both Orthology and alternative splicing isoforms are able to achieve high performance. In addition, the six track representation is effective for RNA property prediction tasks that are known to be influenced by exon junction density such as RNA half-life prediction [38]. Introduction of masking improves all around performance, which could be due to preventing shortcut learning [5, 13].

Table 2: Ablation results are generated using linear probing, by fitting a linear model on pre-computed embeddings from Orthrus base models. 6t; Six track input corresponding to one hot encoded sequence, splicing and codon positions. Masking corresponds to randomly masking 15% of the input sequence.

Splice	Orthology	6 tracks	Masking	RNA HL Human R	RNA HL Mouse R	MRL R	GO MF ROC AUC	Protein Loc. ROC AUC
✓	✓	✓	✓	0.675	0.615	0.393	0.845	0.834
✗	✓	✓	✓	0.678	0.610	0.392	0.842	0.833
✓	✗	✓	✓	0.680	0.615	0.402	0.854	0.834
✓	✓	✗	✓	0.531	0.512	0.361	0.833	0.825
✓	✓	✓	✗	0.647	0.595	0.332	0.836	0.822
✗	✗	✗	✗	0.217	0.214	0.114	0.753	0.792

4 Discussion

As Dobzhansky famously notes: “Nothing in biology makes sense except in the light of evolution” [15]. Orthrus similarly aims to capture the diversity of RNA through an evolutionary and functional lens [30, 37]. We create a self-supervised training objective that learns similarities between evolutionarily related sequences identified in the Zoonomia project [28]. In addition, we utilize alternatively spliced transcripts to learn sequences responsible for shared functions between splicing isoforms [45]. By training on sequences generated by evolutionary and alternative splicing processes, Orthrus utilizes stronger biologically motivated inductive biases compared to SSL reconstruction methods. This makes Orthrus less reliant on limited genetic sequence diversity during pre-training, and capable of learning strong representations without fine-tuning on experimental data.

We demonstrate that by minimizing the distance between mature RNAs generated through speciation and alternative splicing, we are able to generate representations useful for RNA property prediction tasks. We empirically demonstrate that Orthrus embeddings contain information useful for predicting RNA properties like RNA half-life and mean ribosome load, and achieves state-of-the-art prediction when fine-tuned. We observe that pre-training is especially helpful in low data regimes when there are 200 or fewer data points with labels. We demonstrate that self-supervised pre-training is an approach for addressing data efficiency challenges present in genomics, and scaling to additional species can be an effective dataset expansion strategy.

An important question to address is why we expect that minimizing distances between RNA isoforms would be useful for predicting phenotypes like RNA half-life or protein localization. One hypothesis is that alternative splicing and speciation events preserve core functional RNA segments. Through the contrastive pre-training procedure, we identify these shared regions between diverse sequences. Indeed, a recent work proposes that contrastive methods are effective due to block separating latent variables shared between views [62]. This view is supported by our findings identifying that Orthrus within gene similarities are correlated with domain presence. By utilizing decoupled contrastive learning, diverse sequences are pushed apart, thus uniformly distributing samples in the latent space, which helps with downstream tasks [66, 63]. Through encoding these invariances, we find that Orthrus is able to learn complex RNA properties such as cellular component localization and RNA half-life.

In this work, we propose a novel, self-supervised contrastive objective for learning mature RNA isoform representations. We show that this approach is an effective strategy to address two major challenges for cellular property prediction: data efficiency, and model generalizability. We demonstrate that Orthrus representations are effective in the low data setting, paving the path to true few-shot learning for RNA property prediction. Finally, we outperform supervised models when fine-tuning Orthrus and significantly improving over performance of reconstruction based self-supervised methods. These findings open the possibility that combining the contrastive loss with a masked language modelling objective can further improve quality of mature RNA representations.

References

- [1] V. Agarwal and D. R. Kelley. The genetic and biochemical determinants of mRNA degradation rates in mammals. *Genome Biol*, 23(1):245, Nov 2022.
- [2] Michael Ashburner, Catherine A. Ball, Judith A. Blake, David Botstein, Heather Butler, J. Michael Cherry, Allan P. Davis, Kara Dolinski, Selina S. Dwight, Janan T. Eppig, Midori A. Harris, David P. Hill, Laurie Issel-Tarver, Andrew Kasarskis, Suzanna Lewis, John C. Matese, Joel E. Richardson, Martin Ringwald, Gerald M. Rubin, and Gavin Sherlock. Gene ontology: tool for the unification of biology. *Nature Genetics*, 25(1):25–29, May 2000.
- [3] Randall Balestriero, Mark Ibrahim, Vlad Sobal, Ari Morcos, Shashank Shekhar, Tom Goldstein, Florian Bordes, Adrien Bardes, Gregoire Mialon, Yuandong Tian, Avi Schwarzschild, Andrew Gordon Wilson, Jonas Geiping, Quentin Garrido, Pierre Fernandez, Amir Bar, Hamed Pirsiavash, Yann LeCun, and Micah Goldblum. A Cookbook of Self-Supervised Learning. *arXiv e-prints*, April 2023.
- [4] Adrien Bardes, Jean Ponce, and Yann LeCun. VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning. *arXiv e-prints*, May 2021.
- [5] Ido Ben-Shaul, Ravid Shwartz-Ziv, Tomer Galanti, Shai Dekel, and Yann LeCun. Reverse Engineering Self-Supervised Learning. *arXiv e-prints*, May 2023.
- [6] Gloria A Brar and Jonathan S Weissman. Ribosome profiling reveals the what, when, where and how of protein synthesis. *Nature reviews Molecular cell biology*, 16(11):651–664, 2015.
- [7] Albi Celaj, Alice Jiexin Gao, Tammy T.Y. Lau, Erle M. Holgersen, Alston Lo, Varun Lodaya, Christopher B. Cole, Robert E. Denroche, Carl Spickett, Omar Wagih, Pedro O. Pinheiro, Parth Vora, Pedrum Mohammadi-Shemirani, Steve Chan, Zach Nussbaum, Xi Zhang, Helen Zhu, Easwaran Ramamurthy, Bhargav Kanuparthi, Michael Iacocca, Diane Ly, Ken Kron, Marta Verby, Kahlin Cheung-Ong, Zvi Shalev, Brandon Vaz, Sakshi Bhargava, Farhan Yusuf, Sharon Samuel, Sabriyeh Alibai, Zahra Baghestani, Xinwen He, Kirsten Krastel, Oladipo Oladapo, Amrudha Mohan, Arathi Shanavas, Magdalena Bugno, Jovanka Bogojeski, Frank Schmitges, Carolyn Kim, Solomon Grant, Rachana Jayaraman, Tehmina Masud, Amit Deshwar, Shreshth Gandhi, and Brendan J. Frey. An rna foundation model enables discovery of disease mechanisms and candidate therapeutics. *bioRxiv*, 2023.
- [8] Jiayang Chen, Zhihang Hu, Siqi Sun, Qingxiong Tan, Yixuan Wang, Qinze Yu, Licheng Zong, Liang Hong, Jin Xiao, Tao Shen, Irwin King, and Yu Li. Interpretable RNA Foundation Model from Unannotated Data for Highly Accurate RNA Structure and Function Predictions. *arXiv e-prints*, page arXiv:2204.00300, April 2022.
- [9] Siwei Chen, Laurent C. Francioli, Julia K. Goodrich, Ryan L. Collins, Masahiro Kanai, Qingbo Wang, Jessica Alföldi, Nicholas A. Watts, Christopher Vittal, Laura D. Gauthier, Timothy Poterba, Michael W. Wilson, Yekaterina Tarasova, William Phu, Riley Grant, Mary T. Johannes, Zan Koenig, Yossi Farjoun, Eric Banks, Stacey Donnelly, Stacey Gabriel, Namrata Gupta, Steven Ferriera, Charlotte Tolonen, Sam Novod, Louis Bergelson, David Roazen, Valentin Ruano-Rubio, Miguel Covarrubias, Christopher Llanwarne, Nikelle Petrillo, Gordon Wade, Thibault Jeandet, Ruchi Munshi, Kathleen Tibbetts, Maria Abreu, Carlos A. Aguilar Salinas, Tariq Ahmad, Christine M. Albert, Diego Ardissino, Irina M. Armean, Elizabeth G. Atkinson, Gil Atzmon, John Barnard, Samantha M. Baxter, Laurent Beaugerie, Emelia J. Benjamin, David Benjamin, Michael Boehnke, Lori L. Bonnycastle, Erwin P. Bottinger, Donald W. Bowden, Matthew J. Bown, Harrison Brand, Steven Brant, Ted Brookings, Sam Bryant, Sarah E. Calvo, Hannia Campos, John C. Chambers, Juliana C. Chan, Katherine R. Chao, Sinéad Chapman, Daniel I. Chasman, Rex Chisholm, Judy Cho, Rajiv Chowdhury, Mina K. Chung, Wendy K. Chung, Kristian Cibulskis, Bruce Cohen, Kristen M. Connolly, Adolfo Correa, Beryl B. Cummings, Dana Dabelea, John Danesh, Dawood Darbar, Phil Darnowsky, Joshua Denny, Ravindranath Duggirala, Josée Dupuis, Patrick T. Ellinor, Roberto Elosua, James Emery, Eleina England, Jeanette Erdmann, Tõnu Esko, Emily Evangelista, Diane Fatkin, Jose Florez, Andre Franke, Jack Fu, Martti Färkkilä, Kiran Garimella, Jeff Gentry, Gad Getz, David C. Glahn, Benjamin Glaser, Stephen J. Glatt, David Goldstein, Clicerio Gonzalez, Leif Groop, Sanna Gudmundsson, Andrea Haessly, Christopher Haiman,

Ira Hall, Craig L. Hanis, Matthew Harms, Mikko Hiltunen, Matti M. Holi, Christina M. Hultman, Chaim Jalas, Mikko Kallela, Diane Kaplan, Jaakko Kaprio, Sekar Kathiresan, Eimear E. Kenny, Bong-Jo Kim, Young Jin Kim, Daniel King, George Kirov, Jaspal Kooner, Seppo Koskinen, Harlan M. Krumholz, Subra Kugathasan, Soo Heon Kwak, Markku Laakso, Nicole Lake, Trevyn Langsford, Kristen M. Laricchia, Terho Lehtimäki, Monkol Lek, Emily Lipscomb, Ruth J. F. Loos, Wenhan Lu, Steven A. Lubitz, Teresa Tusie Luna, Ronald C. W. Ma, Gregory M. Marcus, Jaume Marrugat, Kari M. Mattila, Steven McCarroll, Mark I. McCarthy, Jacob L. McCauley, Dermot McGovern, Ruth McPherson, James B. Meigs, Olle Melander, Andres Metspalu, Deborah Meyers, Eric V. Minikel, Braxton D. Mitchell, Vamsi K. Mootha, Aliya Naheed, Saman Nazarian, Peter M. Nilsson, Michael C. O'Donovan, Yukinori Okada, Dost Ongur, Lorena Orozco, Michael J. Owen, Colin Palmer, Nicholette D. Palmer, Aarno Palotie, Kyong Soo Park, Carlos Pato, Ann E. Pulver, Dan Rader, Nazneen Rahman, Alex Reiner, Anne M. Remes, Dan Rhodes, Stephen Rich, John D. Rioux, Samuli Ripatti, Dan M. Roden, Jerome I. Rotter, Nareh Sahakian, Danish Saleheen, Veikko Salomaa, Andrea Saltzman, Nilesh J. Samani, Kaitlin E. Samocha, Alba Sanchis-Juan, Jeremiah Scharf, Molly Schleicher, Heribert Schunkert, Sebastian Schönherr, Eleanor G. Seaby, Svati H. Shah, Megan Shand, Ted Sharpe, Moore B. Shoemaker, Tai Shyong, Edwin K. Silverman, Moriel Singer-Berk, Pamela Sklar, Jonathan T. Smith, J. Gustav Smith, Hilka Soininen, Harry Sokol, Rachel G. Son, Jose Soto, Tim Spector, Christine Stevens, Nathan O. Stitzel, Patrick F. Sullivan, Jaana Suvisaari, E. Shyong Tai, Kent D. Taylor, Yik Ying Teo, Ming Tsuang, Tiinamaija Tuomi, Dan Turner, Teresa Tusie-Luna, Erkki Vartiainen, Marquis Vawter, Lily Wang, Arcurus Wang, James S. Ware, Hugh Watkins, Rinse K. Weersma, Ben Weisburd, Maija Wessman, Nicola Whiffin, James G. Wilson, Ramnik J. Xavier, Anne O'Donnell-Luria, Matthew Solomonson, Cotton Seed, Alicia R. Martin, Michael E. Talkowski, Heidi L. Rehm, Mark J. Daly, Grace Tiao, Benjamin M. Neale, Daniel G. MacArthur, and Konrad J. Karczewski. Author correction: A genomic mutational constraint map using variation in 76, 156 human genomes. *Nature*, January 2024.

- [10] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. *arXiv e-prints*, page arXiv:2002.05709, February 2020.
- [11] The Gene Ontology Consortium, Suzi A Aleksander, James Balhoff, Seth Carbon, J Michael Cherry, Harold J Drabkin, Dustin Ebert, Marc Feuermann, Pascale Gaudet, Nomi L Harris, David P Hill, Raymond Lee, Huaiyu Mi, Sierra Moxon, Christopher J Mungall, Anushya Muruganugan, Tremayne Mushayahama, Paul W Sternberg, Paul D Thomas, Kimberly Van Auken, Jolene Ramsey, Deborah A Siegele, Rex L Chisholm, Petra Fey, Maria Cristina Aspromonte, Maria Victoria Nugnes, Federica Quaglia, Silvio Tosatto, Michelle Giglio, Suvarna Nadendla, Giulia Antonazzo, Helen Attrill, Gil dos Santos, Steven Marygold, Victor Strelets, Christopher J Tabone, Jim Thurmond, Pinglei Zhou, Saadullah H Ahmed, Praoparn Asanitthong, Diana Luna Buitrago, Meltem N Erdol, Matthew C Gage, Mohamed Ali Kadhum, Kan Yan Chloe Li, Miao Long, Aleksandra Michalak, Angeline Pesala, Armalya Pritazahra, Shirin C C Saverimuttu, Renzhi Su, Kate E Thurlow, Ruth C Lovering, Colin Logie, Snezhana Oliferenko, Judith Blake, Karen Christie, Lori Corbani, Mary E Dolan, Harold J Drabkin, David P Hill, Li Ni, Dmitry Sitnikov, Cynthia Smith, Alayne Cuzick, James Seager, Laurel Cooper, Justin Elser, Pankaj Jaiswal, Parul Gupta, Pankaj Jaiswal, Sushma Naithani, Manuel Lera-Ramirez, Kim Rutherford, Valerie Wood, Jeffrey L De Pons, Melinda R Dwinell, G Thomas Hayman, Mary L Kaldunski, Anne E Kwitek, Stanley J F Laulederkind, Marek A Tutaj, Mahima Vedi, Shur-Jen Wang, Peter D'Eustachio, Lucila Aimò, Kristian Axelsen, Alan Bridge, Nevila Hyka-Nouspikel, Anne Morgat, Suzi A Aleksander, J Michael Cherry, Stacia R Engel, Kalpana Karra, Stuart R Miyasato, Robert S Nash, Marek S Skrzypek, Shuai Weng, Edith D Wong, Erika Bakker, Tanya Z Berardini, Leonore Reiser, Andrea Auchincloss, Kristian Axelsen, Ghislaine Argoud-Puy, Marie-Claude Blatter, Emmanuel Boutet, Lionel Breuza, Alan Bridge, Cristina Casals-Casas, Elisabeth Coudert, Anne Estreicher, Maria Livia Famiglietti, Marc Feuermann, Arnaud Gos, Nadine Gruaz-Gumowski, Chantal Hulo, Nevila Hyka-Nouspikel, Florence Jungo, Philippe Le Mercier, Damien Lieberherr, Patrick Masson, Anne Morgat, Ivo Pedruzzi, Lucille Pourcel, Sylvain Poux, Catherine Rivoire, Shyamala Sundaram, Alex Bateman, Emily Bowler-Barnett, Hema Bye-A-Jee, Paul Denny, Alexandr Ignatchenko, Rizwan Ishtiaq, Antonia Lock, Yvonne Lussi, Michele Magrane, Maria J Martin, Sandra Orchard, Pedro Raposo, Elena Speretta, Nidhi Tyagi, Kate Warner, Rossana Zaru, Alexan-

- der D Diehl, Raymond Lee, Juancarlos Chan, Stavros Diamantakis, Daniela Raciti, Magdalena Zarowiecki, Malcolm Fisher, Christina James-Zorn, Virgilio Ponferrada, Aaron Zorn, Sridhar Ramachandran, Leyla Ruzicka, and Monte Westerfield. The Gene Ontology knowledgebase in 2023. *Genetics*, 224(1):iyad031, 03 2023.
- [12] Hugo Dalla-Torre, Liam Gonzalez, Javier Mendoza Revilla, Nicolas Lopez Carranza, Adam Henryk Grzywaczewski, Francesco Oteri, Christian Dallago, Evan Trop, Hassan Sirelkhatim, Guillaume Richard, Marcin Skwark, Karim Beguir, Marie Lopez, and Thomas Pierrot. The nucleotide transformer: Building and evaluating robust foundation models for human genomics. *bioRxiv*, 2023.
- [13] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. *arXiv e-prints*, page arXiv:1810.04805, October 2018.
- [14] C. B. Do, M. S. Mahabhashyam, M. Brudno, and S. Batzoglou. ProbCons: Probabilistic consistency-based multiple sequence alignment. *Genome Res*, 15(2):330–340, Feb 2005.
- [15] Theodosius Dobzhansky. Nothing in biology makes sense except in the light of evolution. *The american biology teacher*, 75(2):87–91, 2013.
- [16] A. Frankish, M. Diekhans, I. Jungreis, J. Lagarde, J. E. Loveland, J. M. Mudge, C. Sisu, J. C. Wright, J. Armstrong, I. Barnes, A. Berry, A. Bignell, C. Boix, S. Carbonell Sala, F. Cunningham, T. Di Domenico, S. Donaldson, I. T. Fiddes, C. n, J. M. Gonzalez, T. Grego, M. Hardy, T. Hourlier, K. L. Howe, T. Hunt, O. G. Izuogu, R. Johnson, F. J. Martin, L. nez, S. Mohanan, P. Muir, F. C. P. Navarro, A. Parker, B. Pei, F. Pozo, F. C. Riera, M. Ruffier, B. M. Schmitt, E. Stapleton, M. M. Suner, I. Sycheva, B. Uszczynska-Ratajczak, M. Y. Wolf, J. Xu, Y. T. Yang, A. Yates, D. Zerbino, Y. Zhang, J. S. Choudhary, M. Gerstein, R. ó, T. J. P. Hubbard, M. Kellis, B. Paten, M. L. Tress, and P. Flicek. GENCODE 2021. *Nucleic Acids Res*, 49(D1):D916–D923, Jan 2021.
- [17] J. Frazer, P. Notin, M. Dias, A. Gomez, J. K. Min, K. Brock, Y. Gal, and D. S. Marks. Disease variant prediction with deep generative models of evolutionary data. *Nature*, 599(7883):91–95, Nov 2021.
- [18] Quentin Garrido, Yubei Chen, Adrien Bardes, Laurent Najman, and Yann Lecun. On the duality between contrastive and non-contrastive self-supervised learning. *arXiv e-prints*, page arXiv:2206.02574, June 2022.
- [19] Ilias Georgakopoulos-Soares, Chengyu Deng, Vikram Agarwal, Candace S. Y. Chan, Jingjing Zhao, Fumitaka Inoue, and Nadav Ahituv. Transcription factor binding site orientation and order are major drivers of gene regulatory activity. *Nature Communications*, 14(1), April 2023.
- [20] Jean-Bastien Grill, Florian Strub, Florent Alché, Corentin Tallec, Pierre H. Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Daniel Guo, Mohammad Gheshlaghi Azar, Bilal Piot, Koray Kavukcuoglu, Rémi Munos, and Michal Valko. Bootstrap your own latent: A new approach to self-supervised Learning. *arXiv e-prints*, June 2020.
- [21] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces, 2024.
- [22] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces, 2022.
- [23] Sanna Gudmundsson, Moriel Singer-Berk, Nicholas A. Watts, William Phu, Julia K. Goodrich, Matthew Solomonson, Heidi L. Rehm, Daniel G. MacArthur, and Anne O'Donnell-Luria and. Variant interpretation using population databases: Lessons from gnomAD. *Human Mutation*, 43(8):1012–1030, December 2021.
- [24] Veronika A Herzog, Brian Reichholf, Tobias Neumann, Philipp Rescheneder, Pooja Bhat, Thomas R Burkard, Wiebke Wlotzka, Arndt Von Haeseler, Johannes Zuber, and Stefan L Ameres. Thiol-linked alkylation of rna to assess expression dynamics. *Nature methods*, 14(12):1198–1204, 2017.

- [25] K. Jaganathan, S. Kyriazopoulou Panagiotopoulou, J. F. McRae, S. F. Darbandi, D. Knowles, Y. I. Li, J. A. Kosmicki, J. Arbelaez, W. Cui, G. B. Schwartz, E. D. Chow, E. Kanterakis, H. Gao, A. Kia, S. Batzoglou, S. J. Sanders, and K. K. Farh. Predicting Splicing from Primary Sequence with Deep Learning. *Cell*, 176(3):535–548, Jan 2019.
- [26] Y. Ji, Z. Zhou, H. Liu, and R. V. Davuluri. DNABERT: pre-trained Bidirectional Encoder Representations from Transformers model for DNA-language in genome. *Bioinformatics*, 37(15):2112–2120, Aug 2021.
- [27] J. Jumper, R. Evans, A. Pritzel, T. Green, M. Figurnov, O. Ronneberger, K. Tunyasuvunakool, R. Bates, A. de A. Potapenko, A. Bridgland, C. Meyer, S. A. A. Kohl, A. J. Ballard, A. Cowie, B. Romera-Paredes, S. Nikolov, R. Jain, J. Adler, T. Back, S. Petersen, D. Reiman, E. Clancy, M. Zielinski, M. Steinegger, M. Pacholska, T. Berghammer, S. Bodenstein, D. Silber, O. Vinyals, A. W. Senior, K. Kavukcuoglu, P. Kohli, and D. Hassabis. Highly accurate protein structure prediction with AlphaFold. *Nature*, 596(7873):583–589, Aug 2021.
- [28] Bogdan M. Kirilenko, Chetan Munegowda, Ekaterina Osipova, David Jebb, Virag Sharma, Moritz Blumer, Ariadna E. Morales, Alexis-Walid Ahmed, Dimitrios-Georgios Kontopoulos, Leon Hilgers, Kerstin Lindblad-Toh, Elinor K. Karlsson, Michael Hiller, Gregory Andrews, Joel C. Armstrong, Matteo Bianchi, Bruce W. Birren, Kevin R. Bredemeyer, Ana M. Breit, Matthew J. Christmas, Hiram Clawson, Joana Damas, Federica Di Palma, Mark Diekhans, Michael X. Dong, Eduardo Eizirik, Kaili Fan, Cornelia Fanter, Nicole M. Foley, Karin Forsberg-Nilsson, Carlos J. Garcia, John Gatesy, Steven Gazal, Diane P. Genereux, Linda Goodman, Jenna Grimshaw, Michaela K. Halsey, Andrew J. Harris, Glenn Hickey, Michael Hiller, Allyson G. Hindle, Robert M. Hubley, Graham M. Hughes, Jeremy Johnson, David Juan, Irene M. Kaplow, Elinor K. Karlsson, Kathleen C. Keough, Bogdan Kirilenko, Klaus-Peter Koepfli, Jennifer M. Korstian, Amanda Kowalczyk, Sergey V. Kozyrev, Alyssa J. Lawler, Colleen Lawless, Thomas Lehmann, Danielle L. Levesque, Harris A. Lewin, Xue Li, Abigail Lind, Kerstin Lindblad-Toh, Ava Mackay-Smith, Voichita D. Marinescu, Tomas Marques-Bonet, Victor C. Mason, Jennifer R. S. Meadows, Wynn K. Meyer, Jill E. Moore, Lucas R. Moreira, Diana D. Moreno-Santillan, Kathleen M. Morrill, Gerard Muntané, William J. Murphy, Arcadi Navarro, Martin Nweeia, Sylvia Ortmann, Austin Osmanski, Benedict Paten, Nicole S. Paulat, Andreas R. Pfenning, BaDoi N. Phan, Katherine S. Pollard, Henry E. Pratt, David A. Ray, Steven K. Reilly, Jeb R. Rosen, Irina Ruf, Louise Ryan, Oliver A. Ryder, Pardis C. Sabeti, Daniel E. Schäffer, Aitor Serres, Beth Shapiro, Arian F. A. Smit, Mark Springer, Chaitanya Srinivasan, Cynthia Steiner, Jessica M. Storer, Kevin A. M. Sullivan, Patrick F. Sullivan, Elisabeth Sundström, Megan A. Supple, Ross Swofford, Joy-El Talbot, Emma Teeling, Jason Turner-Maier, Alejandro Valenzuela, Franziska Wagner, Ola Wallerman, Chao Wang, Juehan Wang, Zhiping Weng, Aryn P. Wilder, Morgan E. Wirthlin, James R. Xue, and Xiaomeng Zhang. Integrating gene annotation with orthology inference at scale. *Science*, 380(6643), April 2023.
- [29] Eugene V. Koonin and Yuri I. Wolf. Constraints and plasticity in genome and molecular-phenome evolution. *Nature Reviews Genetics*, 11(7):487–498, June 2010.
- [30] N. K. Lee, Z. Tang, S. Toneyan, and P. K. Koo. EvoAug: improving generalization and interpretability of genomic deep neural networks with evolution-inspired data augmentations. *Genome Biol*, 24(1):105, May 2023.
- [31] Z. Lin, H. Akin, R. Rao, B. Hie, Z. Zhu, W. Lu, N. Smetanin, R. Verkuil, O. Kabeli, Y. Shmueli, A. Dos Santos Costa, M. Fazel-Zarandi, T. Sercu, S. Candido, and A. Rives. Evolutionary-scale prediction of atomic-level protein structure with a language model. *Science*, 379(6637):1123–1130, Mar 2023.
- [32] Kerstin Lindblad-Toh, Manuel Garber, Or Zuk, Michael F. Lin, Brian J. Parker, Stefan Washietl, Pouya Kheradpour, Jason Ernst, Gregory Jordan, Evan Mauceli, Lucas D. Ward, Craig B. Lowe, Alisha K. Holloway, Michele Clamp, Sante Gnerre, Jessica Alföldi, Kathryn Beal, Jean Chang, Hiram Clawson, James Cuff, Federica Di Palma, Stephen Fitzgerald, Paul Flicek, Mitchell Guttman, Melissa J. Hubisz, David B. Jaffe, Irwin Jungreis, W. James Kent, Dennis Kostka, Marcia Lara, Andre L. Martins, Tim Massingham, Ida Moltke, Brian J. Raney, Matthew D. Rasmussen, Jim Robinson, Alexander Stark, Albert J. Vilella, Jiayu Wen, Xiaohui Xie, Michael C. Zody, Kim C. Worley, Christie L. Kovar, Donna M. Muzny, Richard A.

- Gibbs, Wesley C. Warren, Elaine R. Mardis, George M. Weinstock, Richard K. Wilson, Ewan Birney, Elliott H. Margulies, Javier Herrero, Eric D. Green, David Haussler, Adam Siepel, Nick Goldman, Katherine S. Pollard, Jakob S. Pedersen, Eric S. Lander, and Manolis Kellis. A high-resolution map of human evolutionary constraint using 29 mammals. *Nature*, 478(7370):476–482, October 2011.
- [33] J. Linder, S. E. Koplik, A. Kundaje, and G. Seelig. Deciphering the impact of genetic variation on human polyadenylation using APARENT2. *Genome Biol*, 23(1):232, Nov 2022.
- [34] Larry Sai Weng Loo, Andreas Alvin Purnomo Soetedjo, Hwee Hui Lau, Natasha Hui Jin Ng, Soumita Ghosh, Linh Nguyen, Vidhya Gomathi Krishnan, Hyungwon Choi, Xavier Roca, Shawn Hoon, and Adrian Kee Keong Teo. Bcl-x1/bcl2l1 is a critical anti-apoptotic protein that promotes the survival of differentiating pancreatic cells from human pluripotent stem cells. *Cell Death Disease*, 11(5), May 2020.
- [35] Mohammad Lotfollahi, Anna Klimovskaia Susmelj, Carlo De Donno, Leon Hetzel, Yuge Ji, Ignacio L Ibarra, Sanjay R Srivatsan, Mohsen Naghipourfar, Riza M Daza, Beth Martin, Jay Shendure, Jose L McFaline-Figueroa, Pierre Boyeau, F Alexander Wolf, Nafissa Yakubova, Stephan Günemann, Cole Trapnell, David Lopez-Paz, and Fabian J Theis. Predicting cellular responses to complex perturbations in high-throughput screens. *Molecular Systems Biology*, 19(6), May 2023.
- [36] Amy X. Lu, Alex X. Lu, and Alan Moses. Evolution Is All You Need: Phylogenetic Augmentation for Contrastive Learning. *arXiv e-prints*, December 2020.
- [37] Amy X. Lu, Haoran Zhang, Marzyeh Ghassemi, and Alan Moses. Self-supervised contrastive learning of protein representations by mutual information maximization. *bioRxiv*, 2020.
- [38] Zhiyuan Luo, Qilian Ma, Shan Sun, Ningning Li, Hongfeng Wang, Zheng Ying, and Shengdong Ke. Exon-intron boundary inhibits m6a deposition, enabling m6a distribution hallmark, longer mrna half-life and flexible protein coding. *Nature Communications*, 14(1), July 2023.
- [39] Fergal J Martin, M Ridwan Amode, Alisha Aneja, Olanrewaju Austine-Orimoloye, Andrey G Azov, If Barnes, Arne Becker, Ruth Bennett, Andrew Berry, Jyothish Bhai, Simarpreet Kaur Bhurji, Alexandra Bignell, Sanjay Boddu, Paulo R Branco Lins, Lucy Brooks, Shashank Budhanuru Ramaraju, Mehrnaz Charkhchi, Alexander Cockburn, Luca Da Rin Fiorretto, Claire Davidson, Kamalkumar Dodiya, Sarah Donaldson, Bilal El Houdaigui, Tamara El Naboulsi, Reham Fatima, Carlos Garcia Giron, Thiago Genez, Gurpreet S Ghattaoraya, Jose Gonzalez Martinez, Cristi Guijarro, Matthew Hardy, Zoe Hollis, Thibaut Hourlier, Toby Hunt, Mike Kay, Vinay Kaykala, Tuan Le, Diana Lemos, Diego Marques-Coelho, José Carlos Marugán, Gabriela Alejandra Merino, Louise Paola Mirabueno, Aleena Mushtaq, Syed Nakib Hossain, Denye N Ogeh, Manoj Pandian Sakthivel, Anne Parker, Malcolm Perry, Ivana Piližota, Irina Prosovetskaia, José G Pérez-Silva, Ahamed Imran Abdul Salam, Nuno Saraiva-Agostinho, Helen Schuilenburg, Dan Sheppard, Swati Sinha, Botond Sipos, William Stark, Emily Steed, Ranjit Sukumaran, Dulika Sumathipala, Marie-Marthe Suner, Likhitha Surapaneni, Kyösti Sutinen, Michal Szpak, Francesca Floriana Tricomi, David Urbina-Gómez, Andres Veidenberg, Thomas A Walsh, Brandon Walts, Elizabeth Wass, Natalie Willhoft, Jamie Allen, Jorge Alvarez-Jarreta, Marc Chakiachvili, Bethany Flint, Stefano Giorgetti, Leanne Haggerty, Garth R Ilsley, Jane E Loveland, Benjamin Moore, Jonathan M Mudge, John Tate, David Thybert, Stephen J Trevanion, Andrea Winterbottom, Adam Frankish, Sarah E Hunt, Magali Ruffier, Fiona Cunningham, Sarah Dyer, Robert D Finn, Kevin L Howe, Peter W Harrison, Andrew D Yates, and Paul Flicek. Ensembl 2023. *Nucleic Acids Research*, 51(D1):D933–D941, November 2022.
- [40] Leland McInnes, John Healy, and James Melville. Umap: Uniform manifold approximation and projection for dimension reduction, 2020.
- [41] Joshua Meier, Roshan Rao, Robert Verkuil, Jason Liu, Tom Sercu, and Alexander Rives. Language models enable zero-shot prediction of the effects of mutations on protein function. *bioRxiv*, 2021.

- [42] D. Merico, C. Spickett, M. O’Hara, B. Kakaradov, A. G. Deshwar, P. Fradkin, S. Gandhi, J. Gao, S. Grant, K. Kron, F. W. Schmitges, Z. Shalev, M. Sun, M. Verby, M. Cahill, J. J. Dowling, J. Fransson, E. Wienholds, and B. J. Frey. G p.Met645Arg causes Wilson disease by promoting exon 6 skipping. *NPJ Genom Med*, 5:16, 2020.
- [43] Eric Nguyen, Michael Poli, Marjan Faizi, Armin Thomas, Callum Birch-Sykes, Michael Wornow, Aman Patel, Clayton Rabideau, Stefano Massaroli, Yoshua Bengio, Stefano Ermon, Stephen A. Baccus, and Chris Ré. HyenaDNA: Long-Range Genomic Sequence Modeling at Single Nucleotide Resolution. *arXiv e-prints*, June 2023.
- [44] N. A. O’Leary, M. W. Wright, J. R. Brister, S. Ciufu, D. Haddad, R. McVeigh, B. Rajput, B. Robbertse, B. Smith-White, D. Ako-Adjei, A. Astashyn, A. Badretdin, Y. Bao, O. Blinkova, V. Brover, V. Chetvermin, J. Choi, E. Cox, O. Ermolaeva, C. M. Farrell, T. Goldfarb, T. Gupta, D. Haft, E. Hatcher, W. Hlavina, V. S. Joardar, V. K. Kodali, W. Li, D. Maglott, P. Masterson, K. M. McGarvey, M. R. Murphy, K. O’Neill, S. Pujar, S. H. Rangwala, D. Rausch, L. D. Riddick, C. Schoch, A. Shkeda, S. S. Storz, H. Sun, F. Thibaud-Nissen, I. Tolstoy, R. E. Tully, A. R. Vatsan, C. Wallin, D. Webb, W. Wu, M. J. Landrum, A. Kimchi, T. Tatusova, M. DiCuccio, P. Kitts, T. D. Murphy, and K. D. Pruitt. Reference sequence (RefSeq) database at NCBI: current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res*, 44(D1):D733–745, Jan 2016.
- [45] M. Pertea, A. Shumate, G. Pertea, A. Varabyou, F. P. Breitwieser, Y. C. Chang, A. K. Madugundu, A. Pandey, and S. L. Salzberg. CHES: a new human gene catalog curated from thousands of large-scale RNA sequencing experiments reveals extensive transcriptional noise. *Genome Biol*, 19(1):208, Nov 2018.
- [46] Alec Radford. Improving language understanding by generative pre-training. 2018.
- [47] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning Transferable Visual Models From Natural Language Supervision. *arXiv e-prints*, page arXiv:2103.00020, February 2021.
- [48] S. Richards, N. Aziz, S. Bale, D. Bick, S. Das, J. Gastier-Foster, W. W. Grody, M. Hegde, E. Lyon, E. Spector, K. Voelkerding, and H. L. Rehm. Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genet Med*, 17(5):405–424, May 2015.
- [49] Jose Manuel Rodriguez, Fernando Pozo, Daniel Cerdán-Vélez, Tomás Di Domenico, Jesús Vázquez, and Michael L Tress. APPRIS: selecting functionally important isoforms. *Nucleic Acids Res.*, 50(D1):D54–D59, January 2022.
- [50] E. W. Sayers, E. E. Bolton, J. R. Brister, K. Canese, J. Chan, D. C. Comeau, C. M. Farrell, M. Feldgarden, A. M. Fine, K. Funk, E. Hatcher, S. Kannan, C. Kelly, S. Kim, W. Klimke, M. J. Landrum, S. Lathrop, Z. Lu, T. L. Madden, A. Malheiro, A. Marchler-Bauer, T. D. Murphy, L. Phan, S. Pujar, S. H. Rangwala, V. A. Schneider, T. Tse, J. Wang, J. Ye, B. W. Trawick, K. D. Pruitt, and S. T. Sherry. Database resources of the National Center for Biotechnology Information in 2023. *Nucleic Acids Res*, 51(D1):D29–D38, Jan 2023.
- [51] Megan D Schertzer, Andrew Stirn, Keren Isaev, Laura Pereira, Anjali Das, Claire Harbison, Stella H Park, Hans-Hermann Wessels, Neville E Sanjana, and David A Knowles. Cas13d-mediated isoform-specific rna knockdown with a unified computational and experimental toolbox. *bioRxiv*, page 2023.09.12.557474, 2023.
- [52] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 29. Curran Associates, Inc., 2016.
- [53] Yisheng Song, Ting Wang, Subrota K Mondal, and Jyoti Prakash Sahoo. A comprehensive survey of few-shot learning: Evolution, applications, challenges, and opportunities, 2022.

- [54] Yoichiro Sugimoto and Peter J. Ratcliffe. Isoform-resolved mRNA profiling of ribosome load defines interplay of HIF and mTOR dysregulation in kidney cancer. *Nature Structural Molecular Biology*, 29(9):871–880, September 2022.
- [55] Patrick F. Sullivan, Jennifer R. S. Meadows, Steven Gazal, BaDoi N. Phan, Xue Li, Diane P. Genereux, Michael X. Dong, Matteo Bianchi, Gregory Andrews, Sharadha Sakthikumar, Jessika Nordin, Ananya Roy, Matthew J. Christmas, Voichita D. Marinescu, Chao Wang, Ola Wallerman, James Xue, Shuyang Yao, Quan Sun, Jin Szatkiewicz, Jia Wen, Laura M. Huckins, Alyssa Lawler, Kathleen C. Keough, Zhili Zheng, Jian Zeng, Naomi R. Wray, Yun Li, Jessica Johnson, Jiawen Chen, Benedict Paten, Steven K. Reilly, Graham M. Hughes, Zhiping Weng, Katherine S. Pollard, Andreas R. Pfenning, Karin Forsberg-Nilsson, Elinor K. Karlsson, Kerstin Lindblad-Toh, Gregory Andrews, Joel C. Armstrong, Matteo Bianchi, Bruce W. Birren, Kevin R. Bredemeyer, Ana M. Breit, Matthew J. Christmas, Hiram Clawson, Joana Damas, Federica Di Palma, Mark Diekhans, Michael X. Dong, Eduardo Eizirik, Kaili Fan, Cornelia Fanter, Nicole M. Foley, Karin Forsberg-Nilsson, Carlos J. Garcia, John Gatesy, Steven Gazal, Diane P. Genereux, Linda Goodman, Jenna Grimshaw, Michaela K. Halsey, Andrew J. Harris, Glenn Hickey, Michael Hiller, Allyson G. Hindle, Robert M. Hubley, Graham M. Hughes, Jeremy Johnson, David Juan, Irene M. Kaplow, Elinor K. Karlsson, Kathleen C. Keough, Bogdan Kirilenko, Klaus-Peter Koepfli, Jennifer M. Korstian, Amanda Kowalczyk, Sergey V. Kozyrev, Alyssa J. Lawler, Colleen Lawless, Thomas Lehmann, Danielle L. Levesque, Harris A. Lewin, Xue Li, Abigail Lind, Kerstin Lindblad-Toh, Ava Mackay-Smith, Voichita D. Marinescu, Tomas Marques-Bonet, Victor C. Mason, Jennifer R. S. Meadows, Wynn K. Meyer, Jill E. Moore, Lucas R. Moreira, Diana D. Moreno-Santillan, Kathleen M. Morrill, Gerard Muntané, William J. Murphy, Arcadi Navarro, Martin Nweeia, Sylvia Ortmann, Austin Osmani, Benedict Paten, Nicole S. Paulat, Andreas R. Pfenning, BaDoi N. Phan, Katherine S. Pollard, Henry E. Pratt, David A. Ray, Steven K. Reilly, Jeb R. Rosen, Irina Ruf, Louise Ryan, Oliver A. Ryder, Pardis C. Sabeti, Daniel E. Schäffer, Aitor Serres, Beth Shapiro, Arian F. A. Smit, Mark Springer, Chaitanya Srinivasan, Cynthia Steiner, Jessica M. Storer, Kevin A. M. Sullivan, Patrick F. Sullivan, Elisabeth Sundström, Megan A. Supple, Ross Swofford, Joy-El Talbot, Emma Teeling, Jason Turner-Maier, Alejandro Valenzuela, Franziska Wagner, Ola Wallerman, Chao Wang, Juehan Wang, Zhiping Weng, Aryn P. Wilder, Morgan E. Wirthlin, James R. Xue, and Xiaomeng Zhang. Leveraging base-pair mammalian constraint to understand genetic variation and human disease. *Science*, 380(6643), April 2023.
- [56] D. Taliun, D. N. Harris, M. D. Kessler, J. Carlson, Z. A. Szpiech, R. Torres, S. A. G. Taliun, A. Corvelo, S. M. Gogarten, H. M. Kang, A. N. Pitsillides, J. LeFaive, S. B. Lee, X. Tian, B. L. Browning, S. Das, A. K. Emde, W. E. Clarke, D. P. Loesch, A. C. Shetty, T. W. Blackwell, A. V. Smith, Q. Wong, X. Liu, M. P. Conomos, D. M. Bobo, F. Aguet, C. Albert, A. Alonso, K. G. Ardlie, D. E. Arking, S. Aslibekyan, P. L. Auer, J. Barnard, R. G. Barr, L. Barwick, L. C. Becker, R. L. Beer, E. J. Benjamin, L. F. Bielak, J. Blangero, M. Boehnke, D. W. Bowden, J. A. Brody, E. G. Burchard, B. E. Cade, J. F. Casella, B. Chalazan, D. I. Chasman, Y. I. Chen, M. H. Cho, S. H. Choi, M. K. Chung, C. B. Clish, A. Correa, J. E. Curran, B. Custer, D. Darbar, M. Daya, M. de Andrade, D. L. DeMeo, S. K. Dutcher, P. T. Ellinor, L. S. Emery, C. Eng, D. Fatkin, T. Fingerlin, L. Forer, M. Fornage, N. Franceschini, C. Fuchsberger, S. M. Fullerton, S. Germer, M. T. Gladwin, D. J. Gottlieb, X. Guo, M. E. Hall, J. He, N. L. Heard-Costa, S. R. Heckbert, M. R. Irvin, J. M. Johnsen, A. D. Johnson, R. Kaplan, S. L. R. Kardia, T. Kelly, S. Kelly, E. E. Kenny, D. P. Kiel, R. Klemmer, B. A. Konkle, C. Kooperberg, A. tgen, L. A. Lange, J. Lasky-Su, D. Levy, X. Lin, K. H. Lin, C. Liu, R. J. F. Loos, L. Garman, R. Gerszten, S. A. Lubitz, K. L. Lunetta, A. C. Y. Mak, A. Manichaikul, A. K. Manning, R. A. Mathias, D. D. McManus, S. T. McGarvey, J. B. Meigs, D. A. Meyers, J. L. Mikulla, M. A. Minear, B. D. Mitchell, S. Mohanty, M. E. Montasser, C. Montgomery, A. C. Morrison, J. M. Murabito, A. Natale, P. Natarajan, S. C. Nelson, K. E. North, J. R. O’Connell, N. D. Palmer, N. Pankratz, G. M. Peloso, P. A. Peyser, J. Pleiness, W. S. Post, B. M. Psaty, D. C. Rao, S. Redline, A. P. Reiner, D. Roden, J. I. Rotter, I. Ruczinski, C. Sarnowski, S. Schoenherr, D. A. Schwartz, J. S. Seo, S. Seshadri, V. A. Sheehan, W. H. Sheu, M. B. Shoemaker, N. L. Smith, J. A. Smith, N. Sotoodehnia, A. M. Stilp, W. Tang, K. D. Taylor, M. Telen, T. A. Thornton, R. P. Tracy, D. J. Van Den Berg, R. S. Vasan, K. A. Viaud-Martinez, S. Vrieze, D. E. Weeks, B. S. Weir, S. T. Weiss, L. C. Weng, C. J. Willer, Y. Zhang, X. Zhao, D. K. Arnett, A. E. Ashley-Koch, K. C. Barnes, E. Boerwinkle, S. Gabriel, R. Gibbs, K. M. Rice, S. S. Rich, E. K. Silverman, P. Qasba, W. Gan, G. J. Papanicolaou, D. A. Nickerson, S. R. Browning,

- M. C. Zody, S. Ilner, J. G. Wilson, L. A. Cupples, C. C. Laurie, C. E. Jaquish, R. D. Hernandez, T. D. O'Connor, G. R. Abecasis, N. Abe, L. Almasy, S. Ament, P. Anderson, P. Anugu, D. Applebaum-Bowden, T. Assimes, D. Avramopoulos, E. Barron-Casella, T. Beaty, G. Beck, D. Becker, A. Beitelshes, T. Benos, M. Bezerra, J. Bis, R. Bowler, U. Broeckel, J. Broome, K. Bunting, C. Bustamante, E. Buth, J. Cardwell, V. Carey, C. Carty, R. Casaburi, P. Castaldi, M. Chaffin, C. Chang, Y. C. Chang, S. Chavan, B. J. Chen, W. M. Chen, L. M. Chuang, R. H. Chung, S. Comhair, E. Cornell, C. Crandall, J. Crapo, J. Curtis, C. Damcott, S. David, C. Davis, L. L. Fuentes, M. DeBaun, R. Deka, S. Devine, Q. Duan, R. Duggirala, J. P. Durda, C. Eaton, L. Ekunwe, A. El Boueiz, S. Erzurum, C. Farber, M. Flickinger, M. Fornage, C. Frazar, M. Fu, L. Fulton, S. Gao, Y. Gao, M. Gass, B. Gelb, X. P. Geng, M. Geraci, A. Ghosh, C. Gignoux, D. Glahn, D. W. Gong, H. Goring, S. Graw, D. Grine, C. C. Gu, Y. Guan, N. Gupta, J. Haessler, N. L. Hawley, B. Heavner, D. Herrington, C. Hersh, B. Hidalgo, J. Hixson, B. Hobbs, J. Hokanson, E. Hong, K. Hoth, C. A. Hsiung, Y. J. Hung, H. Huston, C. M. Hwu, R. Jackson, D. Jain, M. A. Jhun, C. Johnson, R. Johnston, K. Jones, S. Kathiresan, A. Khan, W. Kim, G. Kinney, H. Kramer, C. Lange, E. Lange, L. Lange, C. Laurie, M. LeBoff, J. Lee, S. S. Lee, W. J. Lee, D. Levine, J. Lewis, X. Li, Y. Li, H. Lin, H. Lin, K. H. Lin, S. Liu, Y. Liu, Y. Liu, J. Luo, M. Mahaney, B. Make, J. Manson, L. Margolin, L. Martin, S. Mathai, S. May, P. McArdle, M. L. McDonald, S. McFarland, D. McGoldrick, C. McHugh, H. Mei, L. Mestroni, N. Min, R. L. Minster, M. Moll, A. Moscati, S. Musani, S. Mwasongwe, J. C. Mychaleckyj, G. Nadkarni, R. Naik, T. Naseri, S. Nekhai, B. Neltner, H. Ochs-Balcom, D. Paik, J. Pankow, A. Parsa, J. M. Peralta, M. Perez, J. Perry, U. Peters, L. S. Phillips, T. Pollin, J. P. Becker, M. P. Boorgula, M. Preuss, D. Qiao, Z. Qin, N. Rafaels, L. Raffield, L. Rasmussen-Torvik, A. Ratan, R. Reed, E. Regan, M. S. Reupena, C. Roselli, P. Russell, S. Ruuska, K. Ryan, E. C. Sabino, D. Saleheen, S. Salimi, S. Salzberg, K. Sandow, V. G. Sankaran, C. Scheller, E. Schmidt, K. Schwander, F. Sciruba, C. Seidman, J. Seidman, S. L. Sherman, A. Shetty, W. H. Sheu, B. Silver, J. Smith, T. Smith, S. Smoller, B. Snively, M. Snyder, T. Sofer, G. Storm, E. Streeten, Y. J. Sung, J. Sylvia, A. Szpiro, C. Sztalryd, H. Tang, M. Taub, M. Taylor, S. Taylor, M. Threlkeld, L. Tinker, D. Tirschwell, S. Tishkoff, H. Tiwari, C. Tong, M. Tsai, D. Vaidya, P. VandeHaar, T. Walker, R. Wallace, A. Walts, F. F. Wang, H. Wang, K. Watson, J. Wessel, K. Williams, L. K. Williams, C. Wilson, J. Wu, H. Xu, L. Yanek, I. Yang, R. Yang, N. Zaghoul, M. Zekavat, S. X. Zhao, W. Zhao, D. Zhi, X. Zhou, and X. Zhu. Sequencing of 53,831 diverse genomes from the NHLBI TOPMed Program. *Nature*, 590(7845):290–299, Feb 2021.
- [57] Peter J Thul, Lovisa Åkesson, Mikaela Wiking, Diana Mahdessian, Aikaterini Geladaki, Hammou Ait Blal, Tove Alm, Anna Asplund, Lars Björk, Lisa M Breckels, Anna Bäckström, Frida Danielsson, Linn Fagerberg, Jenny Fall, Laurent Gatto, Christian Gnann, Sophia Hober, Martin Hjelmare, Fredric Johansson, Sunjae Lee, Cecilia Lindskog, Jan Mulder, Claire M Mulvey, Peter Nilsson, Per Oksvold, Johan Rockberg, Rutger Schutten, Jochen M Schwenk, Åsa Sivertsson, Evelina Sjöstedt, Marie Skogs, Charlotte Stadler, Devin P Sullivan, Hanna Tegel, Casper Winsnes, Cheng Zhang, Martin Zwahlen, Adil Mardinoglu, Fredrik Pontén, Kalle von Feilitzen, Kathryn S Lilley, Mathias Uhlén, and Emma Lundberg. A subcellular map of the human proteome. *Science*, 356(6340), May 2017.
- [58] Nenad Tomasev, Ioana Bica, Brian McWilliams, Lars Buesing, Razvan Pascanu, Charles Blundell, and Jovana Mitrovic. Pushing the limits of self-supervised ResNets: Can we outperform supervised learning without labels on ImageNet? *arXiv e-prints*, January 2022.
- [59] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation Learning with Contrastive Predictive Coding. *arXiv e-prints*, July 2018.
- [60] Eric L Van Nostrand, Gabriel A Pratt, Alexander A Shishkin, Chelsea Gelboin-Burkhart, Mark Y Fang, Balaji Sundararaman, Steven M Blue, Thai B Nguyen, Christine Surka, Keri Elkins, Rebecca Stanton, Frank Rigo, Mitchell Guttman, and Gene W Yeo. Robust transcriptome-wide discovery of RNA-binding protein binding sites with enhanced CLIP (eCLIP). *Nat. Methods*, 13(6):508–514, June 2016.
- [61] A Vaswani. Attention is all you need. *Advances in Neural Information Processing Systems*, 2017.
- [62] Julius von Kügelgen, Yash Sharma, Luigi Gresele, Wieland Brendel, Bernhard Schölkopf, Michel Bessiere, and Francesco Locatello. Self-supervised learning with data augmentations provably isolates content from style, 2022.

- [63] Feng Wang and Huaping Liu. Understanding the behaviour of contrastive loss, 2021.
- [64] Chloe F A Warren, Michelle W Wong-Brown, and Nikola A Bowden. BCL-2 family isoforms in apoptosis and cancer. *Cell Death Dis.*, 10(3):177, February 2019.
- [65] Hui Y. Xiong, Babak Alipanahi, Leo J. Lee, Hannes Bretschneider, Daniele Merico, Ryan K. C. Yuen, Yimin Hua, Serge Gueroussov, Hamed S. Najafabadi, Timothy R. Hughes, Quaid Morris, Yoseph Barash, Adrian R. Krainer, Nebojsa Jojic, Stephen W. Scherer, Benjamin J. Blencowe, and Brendan J. Frey. The human splicing code reveals new insights into the genetic determinants of disease. *Science*, 347(6218), January 2015.
- [66] Chun-Hsiao Yeh, Cheng-Yao Hong, Yen-Chi Hsu, Tyng-Luh Liu, Yubei Chen, and Yann Le-Cun. Decoupled Contrastive Learning. *arXiv e-prints*, page arXiv:2110.06848, October 2021.
- [67] T. Yu, H. Cui, J. C. Li, Y. Luo, G. Jiang, and H. Zhao. Enzyme function prediction using contrastive learning. *Science*, 379(6639):1358–1363, Mar 2023.
- [68] Sangdoon Yun, Dongyoon Han, Seong Joon Oh, Sanghyuk Chun, Junsuk Choe, and Youngjoon Yoo. CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features. *arXiv e-prints*, page arXiv:1905.04899, May 2019.
- [69] Hongyi Zhang, Moustapha Cisse, Yann N. Dauphin, and David Lopez-Paz. mixup: Beyond Empirical Risk Minimization. *arXiv e-prints*, page arXiv:1710.09412, October 2017.
- [70] Ruochi Zhang, Chao Wu, Qian Yang, Chang Liu, Yan Wang, Kewei Li, Lan Huang, and Fengfeng Zhou. Molfescue: enhancing molecular property prediction in data-limited and imbalanced contexts using few-shot and contrastive learning. *Bioinformatics*, 40(4), February 2024.
- [71] N. Zhou, Y. Jiang, T. R. Bergquist, A. J. Lee, B. Z. Kacsóh, A. W. Crocker, K. A. Lewis, G. Georgioui, H. N. Nguyen, M. N. Hamid, L. Davis, T. Dogan, V. Atalay, A. S. Rifaioglu, A. ran, R. Cetin Atalay, C. Zhang, R. L. Hurto, P. L. Freddolino, Y. Zhang, P. Bhat, F. Supek, J. M. ndez, B. Gemovic, V. R. Perovic, R. S. ć, N. Sumonja, N. Veljkovic, E. Asgari, M. R. K. Mofrad, G. Profiti, C. Savojardo, P. L. Martelli, R. Casadio, F. Boecker, H. Schoof, I. Kahanda, N. Thurlby, A. C. McHardy, A. Renaux, R. Saidi, J. Gough, A. A. Freitas, M. Antczak, F. Fabris, M. N. Wass, J. Hou, J. Cheng, Z. Wang, A. E. Romero, A. Paccanaro, H. Yang, T. Goldberg, C. Zhao, L. Holm, P. nen, A. J. Medlar, E. Zosa, I. Borukhov, I. Novikov, A. Wilkins, O. Lichtarge, P. H. Chi, W. C. Tseng, M. Linial, P. W. Rose, C. Dessimoz, V. Vidulin, S. Dzeroski, I. Sillitoe, S. Das, J. G. Lees, D. T. Jones, C. Wan, D. Cozzetto, R. Fa, M. Torres, A. Warwick Vesztrocy, J. M. Rodriguez, M. L. Tress, M. Frasca, M. Notaro, G. Grossi, A. Petrini, M. Re, G. Valentini, M. Mesiti, D. B. Roche, J. Reeb, D. W. Ritchie, S. Aridhi, S. Z. Alborzi, M. D. Devignes, D. C. E. Koo, R. Bonneau, V. ć, M. Barot, H. Fang, S. Toppo, E. Lavezzo, M. Falda, M. Berselli, S. C. E. Tosatto, M. Carraro, D. Piovesan, H. Ur Rehman, Q. Mao, S. Zhang, S. Vucetic, G. S. Black, D. Jo, E. Suh, J. B. Dayton, D. J. Larsen, A. R. Omdahl, L. J. McGuffin, D. A. Brackenridge, P. C. Babbitt, J. M. Yunes, P. Fontana, F. Zhang, S. Zhu, R. You, Z. Zhang, S. Dai, S. Yao, W. Tian, R. Cao, C. Chandler, M. Amezola, D. Johnson, J. M. Chang, W. H. Liao, Y. W. Liu, S. Pascarelli, Y. Frank, R. Hoehndorf, M. Kulmanov, I. Boudellioua, G. Politano, S. Di Carlo, A. Benso, K. Hakala, F. Ginter, F. Mehryary, S. Kaewphan, J. rne, H. Moen, M. E. E. Tolvanen, T. Salakoski, D. Kihara, A. Jain, T. muc, A. Altenhoff, A. Ben-Hur, B. Rost, S. E. Brenner, C. A. Orengo, C. J. Jeffery, G. Bosco, D. A. Hogan, M. J. Martin, C. O’Donovan, S. D. Mooney, C. S. Greene, P. Radivojac, and I. Friedberg. The CAFA challenge reports improved protein function prediction and new functional annotations for hundreds of genes through experimental screens. *Genome Biol*, 20(1):244, Nov 2019.
- [72] 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 e-prints*, June 2023.

A Appendix

Related Works

This work builds on top of foundational efforts spread across three main areas: contrastive representation learning, self-supervised applications in cellular property prediction, and methods for enriching genetic sequence input beyond one hot encoded representation.

Contrastive methods

We build the Orthrus approach for RNA sequences utilizing a rich body of work exploring contrastive learning for computer vision ([3]). A fundamental deep metric learning approach is SimCLR in which the authors propose minimizing the representation distance between two views from the same sample while maximizing the distance between views from different samples ([10]). This approach does not require labeled data and is based on the availability of domain-specific augmentations. Methods like BYOL and VicReg followed and were able to reformulate the contrastive approach by removing the need for in-batch negative samples ([20, 4]). They propose solutions to the trivial solution collapse problem through a variance regularization loss term and architectural design choices. Recent work aims to unify these methods under the contrastive formulation by making a distinction between *sample* and *dimension* contrastive methods ([18]).

Self-supervised learning for cellular properties

Due to the common sequence-based representation between genomics and language, self-supervised learning techniques have long been explored in genomic sequence property predictions. DNABert utilized the BERT problem formulation to learn an encoding for 500 nucleotide long sequences and demonstrated the value for splice site predictions and other tasks ([26, 13, 72]). Nucleotide Transformer (NT), another masked language modeling method, demonstrated the utility of doing data collection from multiple species ([12]). RNA-FM was trained to predict non-coding RNA properties with masked language modeling using 23 million non-coding sequences ([8]). Recently, HyenaDNA has demonstrated that applying long convolutions replacing the attention operation, can lead to effective DNA property prediction while scaling the input sequence length to a million tokens ([43]). In the distinct protein representation learning space, there is a variety of protein language models utilizing auto-regressive and masked language modeling losses to predict protein properties like structure, variant effects, and functional properties ([41, 31]). Contrastive learning has also been used in more specialized domains such as enzyme property prediction while utilizing known shared enzyme properties as views of similar sequences ([67]). Contrastive methods have also been used to learn a more general representation of protein function by maximizing the mutual information between global and local sequence representations ([37]). We build on these works by exploiting domain-specific RNA augmentation to build general representations that are architecture-agnostic.

Beyond one hot encoded genomes

Another important area for advancing cellular property prediction is iterating beyond the reference genome for representing genomic sequences. One such strategy is to integrate random biologically plausible augmentations during training ([30]). By using domain-specific knowledge of the types of augmentations introduced during evolutionary processes, the authors demonstrate they can improve the performance of supervised models for predicting DNA properties. Using multiple sequence alignments is another way to use homology information, common in the protein modeling space ([14, 17, 27]). In another perspective, authors have argued that evolutionary homologs are a viable path for generating augmentations ([36]).

Downstream Evaluation Tasks

RNA half-life (RNA HL) is an important cellular property to measure due to its implications for protein expression regulation. Recently, it has been shown that the choice of experimental methodology for measuring RNA half-life can have an outsize impact ([1]). To address this challenge, Agarwal and Kelley (2022) utilized the first principal component of over 40 different RNA half-life experiments. The dataset consists of 10,432 human and 11,008 mouse RNA sequences with corresponding measurements. The low data availability and high inter-experiment variation underscore

the importance of data efficiency, and generalizability in computational models to be developed for this task.

Mean ribosome load (MRL) is a measure of the translational efficiency of a given mRNA molecule. It measures the number of ribosomes translating a single mRNA molecule at a point in time. Accurate MRL measurement is crucial as it offers insights into the efficiency of protein translation, a key process in cellular function. The dataset in question, derived from the HP5 workflow, captures this metric across 12,459 mRNA isoforms from 7,815 genes ([54]). This dataset was derived from a single experiment, so we can expect a higher amount of noise associated than the RNA half-life dataset.

Protein localization Protein function is often linked to its subcellular location, which can be determined using cells that are immunofluorescently stained. We downloaded a dataset of 10,409 genes, whose protein localization was determined by the Human Protein Atlas ([57]). We included the 12 most common locations including Nucleoplasm, Cytosol, Vesicles, Mitochondria, Plasma Membrane, Golgi apparatus and others. We utilized one transcript per gene (defined to be the canonical isoform by Appris database [49]).

Gene ontology (GO) terms are a hierarchical classification system used for assigning function to genes and their products ([11, 2, 71]). In this work, we utilize GO classes to visualize model latent embeddings and classification. GO term hierarchical systems allow for fine-grained annotation of function, with broader terms at the top of the hierarchy and increased specificity closer to the bottom. To annotate genes with gene ontology terms, we subset GO classes three levels from the root, labeling all available genes.

A.1 Associating Orthrus RNA Embeddings with Transcript Similarity and Protein Domains

To evaluate how well Orthrus RNA embeddings capture functional diversity among transcript isoforms, we analyzed the similarity of transcript pairs within and between protein-coding genes, excluding homologous genes when comparing random gene pairs or genes sharing the same GO term. The test dataset for this analysis was prepared as follows:

1. **Intra-gene Pairs:** We sampled 1,000 genes to obtain pairs of protein-coding transcripts.
2. **Inter-gene Pairs:** We randomly sampled 1,000 pairs of non-homologous genes, selecting the MANE transcript for each gene, which represents the most likely relevant isoform.
3. **Inter-gene Pairs:** We sampled 5,000 GO terms, each containing 10 to 1,000 genes, and selected five non-homologous gene pairs per term.

For each transcript, we computed Orthrus embeddings and calculated pairwise distances between embeddings using the L2 norm. We calculated a similarity score for each transcript pair as $1 - \log(\text{L2 distance})$. This ensures more interpretable results, where higher similarity scores correspond to closer RNA embeddings in the latent space, allowing us to compare the three groups of transcript pairs.

To assess whether similarities in Orthrus embedding reflected shared functional features, we annotated each transcript with protein domain information using Ensembl data and the Pybiomart package. We used the Jaccard Index to quantify the similarity of protein domain presence or absence between each pair of transcripts within a gene. The Jaccard Index is defined as the size of

Task Dataset	Category	Locality	Number of Sequences	Maximum Sequence Length	Homology Split Possible	Species
RNA Half Life Human	Regression	Global	12968	12288	✓	Human
RNA Half Life Mouse	Regression	Global	13738	12288	✓	Mouse
Mean Ribosome Load	Regression	Global	11693	12275	✓	Human
Protein Localization	Classification	Global	9769	12275	✓	Human
Gene Ontology MF	Classification	Global	3697	12236	✓	Human

Table 3: Overview of evaluation datasets. Locality refers to whether the model is required to reason over global structure or be able to pick up local signals. Homology split is not possible for mRFP expression due to the same backbone sequence with differences in synonymous substitutions.

the intersection divided by the size of the union of the protein domain sets present in each transcript pair:

$$\text{Jaccard Index} = \frac{|D_1 \cap D_2|}{|D_1 \cup D_2|}$$

where D_1 and D_2 are the sets of protein domains present in each transcript. Higher values indicate greater similarity in protein domain composition. We calculated this metric using "intra-gene pairs" and "inter-gene pairs" to further study how protein domain composition correlated with embedding similarity. We analyzed the Pearson correlation between Jaccard indices and embedding similarities separately for intra-gene and inter-gene pairs to determine if transcript pairs within the same gene exhibited higher concordance.

To further explore the utility of Orthrus embeddings, we conducted a detailed analysis of the *BCL2L1* gene [64]. Transcripts from this gene were clustered based on their Orthrus embedding similarity scores, with clusters visualized and annotated according to transcript type and known functional roles.

A.2 Linear Probe Experimental Details

In this section, we describe the experimental procedure to evaluate linear probing results.

We first performed a 70-15-15 data split on datasets. The data sequences are then embedded by the various self-supervised learning (SSL) models. For Orthrus, we simply take the mean of the embeddings across the sequence dimension. For HyenaDNA, we take the mean and max of the embedding sequence dimension, as well as the last hidden state in the output sequence. Other SSL methods could not handle input sequences of more than 500 or 1000 nucleotides. Thus, when input sequences exceeded the allowable context window, each sequence was chunked to the maximum length allowed by a model. We then computed the mean of each chunk embedding across the sequence dimension, and then averaged the mean embedding of each chunk to obtain the final embedding.

After obtaining embedding vectors, we used the scikit-learn implementation of linear models to perform the linear probes of the embeddings. For the downstream regression tasks, we used either used linear regression or ridge regression with the regularization parameter selected by cross validation. The final linear model was selected using the validation split. The gene ontology and protein localization tasks are multi-label classification tasks. For this, we fit scikit-learn’s LogisticRegression model to the labels using a MultiOutputClassifier, which essentially trains a separate linear classifier for each label class. We use the default logistic regression parameters, and set 5000 maximum iterations for the solver. Below is the table with the performance numbers that are reported in Figure 2

Model	RNA HL Human	RNA HL Mouse	MRL R	GO MF ROC AUC	Protein Loc. ROC AUC
DNA-BERT2	0.36 ± 7e-3	0.38 ± 6e-3	0.21 ± 9e-3	0.71 ± 2e-3	0.76 ± 1e-3
NT-500m-1000g	0.26 ± 7e-3	0.26 ± 6e-3	0.11 ± 9e-3	0.66 ± 5e-3	0.71 ± 2e-3
NT-500m-human-ref	0.35 ± 4e-3	0.33 ± 7e-3	0.19 ± 5e-3	0.68 ± 3e-3	0.70 ± 1e-3
NT-2.5b-1000g	0.31 ± 5e-3	0.33 ± 4e-3	0.17 ± 7e-3	0.69 ± 2e-3	0.70 ± 1e-3
NT-2.5b-multi-species	0.32 ± 7e-3	0.36 ± 6e-3	0.18 ± 6e-3	0.71 ± 2e-3	0.69 ± 2e-3
Hyena-32K-seqlen	0.42 ± 6e-3	0.44 ± 5e-3	0.28 ± 7e-3	0.74 ± 3e-3	0.78 ± 1e-3
Hyena-160K-seqlen	0.43 ± 7e-3	0.44 ± 6e-3	0.27 ± 5e-3	0.73 ± 2e-3	0.77 ± 1e-3
Hyena-450K-seqlen	0.45 ± 6e-3	0.45 ± 7e-3	0.28 ± 7e-3	0.74 ± 1e-3	0.78 ± 1e-3
RNA-FM	0.40 ± 6e-3	0.38 ± 6e-3	0.19 ± 6e-3	0.74 ± 4e-3	0.78 ± 1e-3
Orthrus Base 4 track	0.52 ± 2e-2	0.54 ± 2e-2	0.38 ± 5e-3	0.84 ± 6e-3	0.83 ± 9e-3
Orthrus Base	0.67 ± 1e-2	0.64 ± 6e-3	0.41 ± 2e-2	0.84 ± 2e-3	0.84 ± 9e-3
Orthrus Large	0.69 ± 1e-2	0.65 ± 1e-2	0.45 ± 2e-2	0.86 ± 4e-3	0.84 ± 9e-3

Table 4: Linear probing results for self-supervised methods. The embeddings were computed for each method and then linear regression was computed analytically using the corresponding labels for each task. Bolded numbers indicate the best performing model.

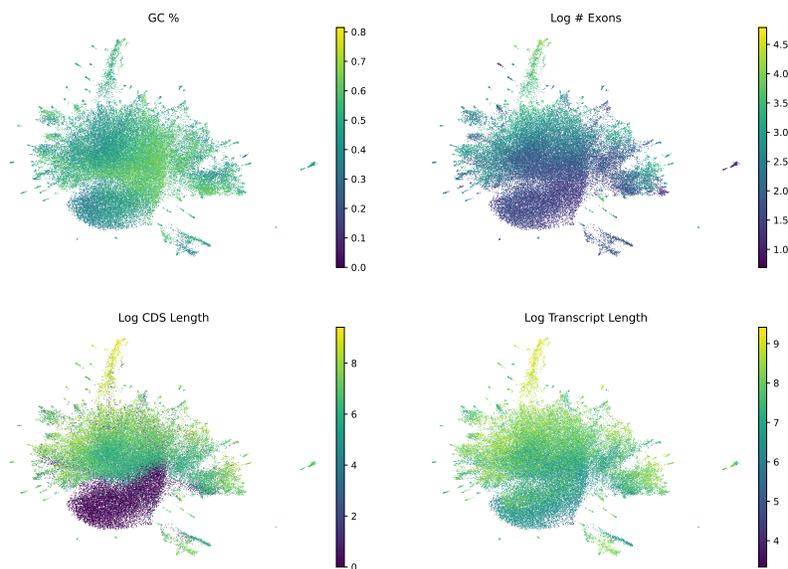


Figure 4: UMAP visualization of the Orthrus Large embeddings. 50,000 transcripts are randomly sampled from gencode comprehensive and colored according to transcript length, number of exons, CDS length, and GC%.

To avoid information leakage we perform homology splitting within each species. Additional details are described in A.3.

In addition, to further evaluate the quality of the model latent space we visualized the learned embeddings by performing dimensionality reduction using UMAP 4 [40]. Each point in the Figure corresponds to a transcript and is colored depending on the biochemical property of interest. We observe clear separation of transcripts conditioned on the GC % CDS length, number of exons, and transcript length. This ensures that information regarding important biochemical properties is preserved in the model embeddings even after the length pooling operation.

A.3 Homology splitting

To perform homology splitting we first acquire paralog information from Ensembl for a species of interest [39]. Ensembl provides pairs paralog information in the form of gene pairs related through duplication events. However, to perform homology splitting between genes we want to make sure that paralog transitivity is taken into account when dividing training samples between train, validation, and test splits. For example given three genes g_1, g_2, g_3 if g_1, g_2 are annotated as paralogs and g_2, g_3 are annotated as well we want to ensure that g_1, g_3 are in the same split. Thus we first transform the pairwise relationships into a graph structure in the process pruning low confidence paralog relationships. We enforce a similarity threshold of 35% which empirically demonstrated highly connected groups of paralogs. The algorithm for grouping is described in 27.

Following construction of the homology graph, during train-val-test split samples associated with genes which are connected in the homology graph are indexed into the same split. This avoids information leakage due to homology relationships within a given species.

A.4 Fine-tuning Experimental Details

We fine-tune Orthrus by first initializing most of the model with weights from pre-training, the penultimate two layers with random initialization, and the final layer with zero initialization. We don't apply any weight decay to weights that were initialized from pre-training while the final three layers have an l2 weight decay term of $1e-5$. We fine-tune on downstream tasks using the Adam

Algorithm 1 Homology Group Assignment

Input: Set of genes G , Set of paralogous relationships P , Similarity threshold s

Output: Homology map assigning genes to groups

Initialize:

$group_counter \leftarrow 0$
 $homology_map \leftarrow \{\}$

Filter:

$P \leftarrow \{(g1, g2) \in P \mid \text{similarity}(g1, g2) > s\}$

for each $(g1, g2)$ **in** P **do**

if $g1 \notin homology_map$ **and** $g2 \notin homology_map$ **then**

$homology_map[g1] \leftarrow group_counter$
 $homology_map[g2] \leftarrow group_counter$
 $group_counter \leftarrow group_counter + 1$

else if $g1 \in homology_map$ **and** $g2 \notin homology_map$ **then**

$homology_map[g2] \leftarrow homology_map[g1]$

else if $g2 \in homology_map$ **and** $g1 \notin homology_map$ **then**

$homology_map[g1] \leftarrow homology_map[g2]$

else if $homology_map[g1] \neq homology_map[g2]$ **then**

$old_group \leftarrow homology_map[g2]$
 $new_group \leftarrow homology_map[g1]$

for each $gene$ **in** $homology_map$ **do**

if $homology_map[gene] == old_group$ **then**

$homology_map[gene] \leftarrow new_group$

end if

end for

end if

end for

Return $homology_map$

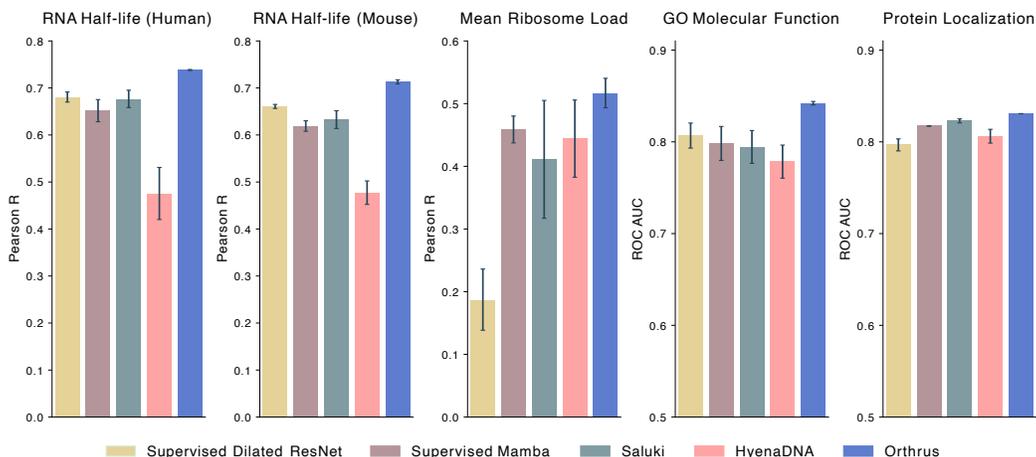


Figure 5: Fine tuning results plotted for 100% fraction of the data. Standard deviation is computed across 3 random seeds. Datasets are split with a homology aware strategy to avoid data leakage.

optimizer with a learning rate of 0.01. We apply exponential learning rate decay with a factor of 0.95. The models are trained with a single Nvidia T4 GPU in a mixed precision setting.

HyenaDNA models initialized with fine-tuning head. We perform a small learning rate hyperparameter grid search around the suggested hyperparameters of $6e-4$. The suggested AdamW optimizer is used. Models were trained for a maximum of 100 epochs on Nvidia T4 GPUs with a batch size of 28 for the HyenaDNA-tiny and a batch size of 8 for HyenaDNA-small. Models were stopped early based on validation loss using an epoch patience of three. After selecting learning rate using the

validation split, the runs were repeated using different random initializations to generate confidence intervals.

Table 5: Pearson correlations (R) and ROC AUC of full model fine-tuning on RNA half-life (HL), mean ribosome load (MRL), gene ontology molecular function (GO MF) classification, protein localization and mRFP expression tasks. Best models are shown in bold. Performance values were averaged over 3 random seeds. Training validation and test were split based on sequence homology to prevent data leakage.

Models	RNA HL Human R	RNA HL Mouse R	MRL R	GO MF ROC AUC	Protein Loc. ROC AUC	Number of Parameters	Number of Tracks
Mamba Base Supervised	0.65 ± 2e-2	0.62 ± 1e-2	0.46 ± 2e-2	0.8 ± 2e-2	0.82 ± 3e-4	1.3m	6
Dilated CNN Resnet	0.68 ± 9e-3	0.66 ± 4e-3	0.13 ± 1e-1	0.81 ± 1e-2	0.8 ± 6e-3	0.83m	6
Saluki	0.68 ± 2e-2	0.63 ± 2e-2	0.41 ± 8e-2	0.79 ± 2e-2	0.82 ± 2e-3	0.15 m	6
HeynaDNA Small	0.48 ± 5e-2	0.48 ± 2e-2	0.44 ± 5e-2	0.78 ± 2e-2	0.81 ± 7e-3	3.3m	4
Orthrus Base	0.74 ± 9e-4	0.71 ± 4e-3	0.52 ± 2e-2	0.84 ± 2e-3	0.83 ± 6e-5	1.3 m	6