

Improving deep models of protein-coding potential with a Fourier-transform architecture and machine translation task

Joseph D. Valencia¹, David A. Hendrix^{1,2,*}

April 19, 2023

¹ School of Electrical Engineering and Computer Science, Oregon State University, Corvallis, OR ,
USA

² Department of Biochemistry and Biophysics, Oregon State University, Corvallis, OR, USA

* Corresponding author: david.hendrix@oregonstate.edu

Abstract

Ribosomes are information-processing macromolecular machines that integrate complex sequence patterns in messenger RNA (mRNA) transcripts to synthesize proteins. Studies of the sequence features that distinguish mRNAs from long noncoding RNAs (lncRNAs) may yield insight into the information that directs and regulates translation. Computational methods for calculating protein-coding potential are important for distinguishing mRNAs from lncRNAs during genome annotation, but most machine learning methods for this task rely on previously known rules to define features. Sequence-to-sequence (seq2seq) models, particularly ones using transformer networks, have proven capable of learning complex grammatical relationships between words to perform natural language translation. Seeking to leverage these advancements in the biological domain, we present a seq2seq formulation for predicting protein-coding potential with deep neural networks and demonstrate that simultaneously learning translation from RNA to protein improves classification performance relative to a classification-only training objective. Inspired by classical signal processing methods for gene discovery and Fourier-based image-processing neural networks, we

24 introduce LocalFilterNet (LFNet). LFNet is a network architecture with an inductive bias for model-
25 ing the three-nucleotide periodicity apparent in coding sequences. We incorporate LFNet within an
26 encoder-decoder framework to test whether the translation task improves the classification of tran-
27 scripts and the interpretation of their sequence features. We use the resulting model to compute
28 nucleotide-resolution importance scores, revealing sequence patterns that could assist the cellular
29 machinery in distinguishing mRNAs and lncRNAs. Finally, we develop a novel approach for es-
30 timating mutation effects from Integrated Gradients, a backpropagation-based feature attribution,
31 and characterize the difficulty of efficient approximations in this setting.

32 **Keywords:** Protein-Coding Potential, Long Noncoding RNAs, Post-Transcriptional regulation, Inter-
33 pretable Deep Learning, Token Mixing Neural Networks, Fourier Transform

34 1 Introduction

35 The flow of genetic information from DNA to RNA to protein is a fundamental life process in which mes-
36 senger RNAs (mRNAs) act as the information-carrying intermediaries. High-throughput sequencing
37 has revealed the abundance of another class of RNA called long noncoding RNAs (lncRNAs), which
38 share important biochemical features such as 5' capping and polyadenylation with protein-coding mR-
39 NAs (Iyer et al. 2015). Long noncoding RNAs are differentiated from smaller noncoding RNAs like
40 tRNAs and microRNAs based on their greater length of at least 200 nucleotide (nt), and from mRNAs
41 based on limited evidence of lncRNA protein expression and sequence conservation (Derrien et al.
42 2012). lncRNAs make up more than 68% of the human transcriptome and play important regula-
43 tory roles, particularly during development (Statello et al. 2021; Ransohoff et al. 2018). They are
44 implicated in numerous diseases including cancer and cardiovascular disease (Sallam et al. 2018).

45 The protein-coding potential of many transcripts is unresolved, and many transcripts previously
46 or currently annotated as lncRNAs are mislabeled and in fact possess small open reading frames
47 (sORFs) that encode micropeptides (Choi et al. 2019). Ribosome profiling (Ribo-Seq) shows that ri-
48 bosomes bind readily to lncRNAs (Ingolia, Lareau, et al. 2011), though the ribosome does not interact
49 with lncRNA ORFs in the same way as mRNAs, lacking a distinctive drop-off of Ribo-Seq coverage
50 at ORF end (Guttman et al. 2013). Ribo-Seq protocols accounting for the 3-nt periodicity of ribosome

51 footprint density (Guo et al. 2010) have identified some genuine sORF translation (Ingolia, Brar, et al.
52 2014; Ji et al. 2015). Only a small fraction of the possible set of micropeptides encoded by transcripts
53 currently annotated as lncRNAs have been directly detected via mass spectrometry, leaving the vast
54 majority as presumptively nonfunctional or rapidly degraded (Housman and Ulitsky 2016; Bánfai et al.
55 2012; Verheggen et al. 2017). Still, hundreds of lncRNAs have been confirmed to be misannotated,
56 and these transcripts do encode micropeptides, for example, myoregulin, a 46-aa. regulator of Ca^{2+}
57 activity that contributes to muscle cell performance (Anderson et al. 2015). Micropeptides are also
58 involved in metabolism, red blood cell development, cardiomyocyte hypertrophy (Yan et al. 2021), in-
59 flammation, tumorigenesis and tumor suppression (Othoum et al. 2020; Wu et al. 2020), and more
60 (Hartford and Lal 2020).

61 Such uncertainty as to the intrinsic protein-coding potential of ORFs raises the question of how
62 cells distinguish true coding regions, with the translational machinery likely to play a critical role.
63 Recent results suggest that general sequence features governing the kinetics of protein synthesis
64 also separate mRNA and untranslated lncRNA ORFs more broadly (Patraquim et al. 2022). The
65 Kozak consensus sequence is well-characterized as the optimal context for translation initiation, and
66 ribosomes can skip unfavorable AUGs through leaky scanning (Kozak 1987; Kozak 2002). Initiation
67 can be affected by cis-regulatory features such as 5' UTR secondary structure (JJ Li et al. 2019) and
68 upstream ORFs (Johnstone et al. 2016), and by trans-acting factors such as microRNAs (Guo et al.
69 2010) and RNA-binding proteins (Szostak and Gebauer 2013). Codon usage biases in the 5'-most
70 region of the CDS are particularly known to affect the elongation rate during protein synthesis (Tuller
71 et al. 2010; Verma et al. 2019; Subramanian et al. 2021).

72 Distinguishing between mRNAs and lncRNAs is an important step in annotating newly sequenced
73 genomes, and a variety of statistical and computational methods have been developed for this task.
74 Codon Adaptive Index (CAI) (Sharp and WH Li 1987) discriminates coding nucleic acids according to
75 biases in the synonymous codons that code for each amino acid and Fickett scores (Fickett 1982) by
76 the nucleotides present in the three codon positions. Early computational approaches used Fourier or
77 wavelet analysis to identify coding sequences (CDS) from their characteristic periodicity of nucleotide
78 identity induced by codon usage bias (Tiwari et al. 1997; Anastassiou 2000; Deng et al. 2010; Has-
79 sani Saadi et al. 2017). Machine learning methods have been designed around features such as the

80 absolute length of ORFs, ORF length relative to the transcript, codon and hexamer frequencies in-
81 cluding Coding Potential Assessment Tool (CPAT) (Wang et al. 2013) and coding potential calculator
82 (Kong et al. 2007), and others (A Li et al. 2014; Wucher et al. 2017).

83 Although many prior machine learning methods achieve high classification performance, they typ-
84 ically rely on transcript-level summary features. Deep learning approaches can operate directly on
85 sequences without such intermediate features and have proven effective in predicting properties of
86 biological sequences, including a wide variety of functional genomics assays (Avsec, Agarwal, et al.
87 2021; Tareen et al. 2022), RNA splicing (Zeng and YI Li 2022) and degradation (Agarwal and Kelley
88 2022), and protein structure (Jumper et al. 2021). A recent method called RNAsamba uses a con-
89 volutional neural network variant to achieve high performance from both nucleotide and amino acid
90 sequence, but also relies on pre-defined features such as the longest ORF (Camargo et al. 2020). A
91 critical limitation in the development of intelligent systems for classifying transcripts as protein-coding
92 vs noncoding is the bias of using the translation and length of the longest ORF in machine learning
93 approaches. Our group previously developed mRNN, the first recurrent neural network classifier of
94 coding RNA from primary sequences alone (Hill et al. 2018). There is a need for more flexible neural
95 networks capable of learning sequence-specific rules that promote translation to better understand
96 what drives translational efficiency. The advantage of these approaches is that they do not require
97 feature engineering, and are capable of learning new biological rules that are encapsulated in the
98 weights of the neural network. Interpretation of these deep neural networks can lead to the identifi-
99 cation of new sequence features that are informative for the evaluation of biological sequences and
100 understanding the regulation of translation. Interpreting deep models is challenging, but a significant
101 literature in explainable artificial intelligence (xAI) has arisen in regulatory genomics, with notable
102 successes in uncovering transcription factor binding logic (Avsec, Weilert, et al. 2021; Novakovsky
103 et al. 2022). Interpretation of similar deep models of protein coding potential could help identify new
104 sequence features regulating translation.

105 In this paper, we describe bioseq2seq, a novel neural network model of biological translation
106 based on the sequence-to-sequence (seq2seq) paradigm commonly used for machine translation of
107 human languages. Although the genetic code follows a well-understood mapping between nucleic
108 acid codons and amino acids, we demonstrate that learning to predict the protein sequence from

109 the sequence of its message improves neural network performance in distinguishing mRNAs from
110 lncRNAs. Adapting recent advances in token mixing neural architectures, we introduce Local Filter
111 Network (LFNet), a computationally efficient network layer based on the short-time Fourier transform.
112 We leverage perturbation-based feature importance values to extract sequence patterns which impact
113 the model prediction and generate hypotheses about the regulatory elements that could differentiate
114 coding RNA *in vivo*. Lastly, we offer evidence that while our LFNet-based bioseq2seq model robustly
115 uncovers biological rules to learn protein-coding potential, it presents challenges for approximate inter-
116 pretation techniques in deep learning. We address these challenges by introducing mutation-directed
117 integrated gradients (MDIG), which we show has a strong correlation with synonymous sequence
118 perturbations, and can be used to identify regions in transcripts that are important for defining protein-
119 coding potential.

120 **2 Results**

121 **2.1 Translation training objective improves classification performance**

122 We downloaded lncRNA primary sequences and mRNAs matched with their encoded proteins from
123 the NCBI RefSeq annotations of eight mammalian species. Following the encoder-decoder frame-
124 work widely used in sequence-to-sequence learning, we trained two major types of deep learning
125 models on this dataset. The primary is bioseq2seq, which outputs a class prediction of $\langle NC \rangle$ for
126 lncRNAs or $\langle PC \rangle$ followed by a predicted protein sequence for coding mRNAs. To test the benefits of
127 a translation-based learning objective, we trained a secondary encoder-decoder model type to predict
128 only the RNA class and not its translation, which we called Encoder Decoder Classifier (EDC). The
129 common architectural framework enables a fair comparison between these two training settings. We
130 designed a novel neural network layer, LFNet, to efficiently apply a short-time (local) Fourier transform
131 to the high-dimensional vectors representing each input nucleotide and perform sequential updates
132 via frequency-domain filtering. Several LFNet layers were composed into an encoder stack to pro-
133 cess the RNA. A stack of transformer decoders operates on the encoder hidden representations to
134 produce an output, autoregressively consuming its own predictions to produce the next character, as
135 necessary (Vaswani et al. 2017). Within this general framework, summarized in Fig 1, we optimized

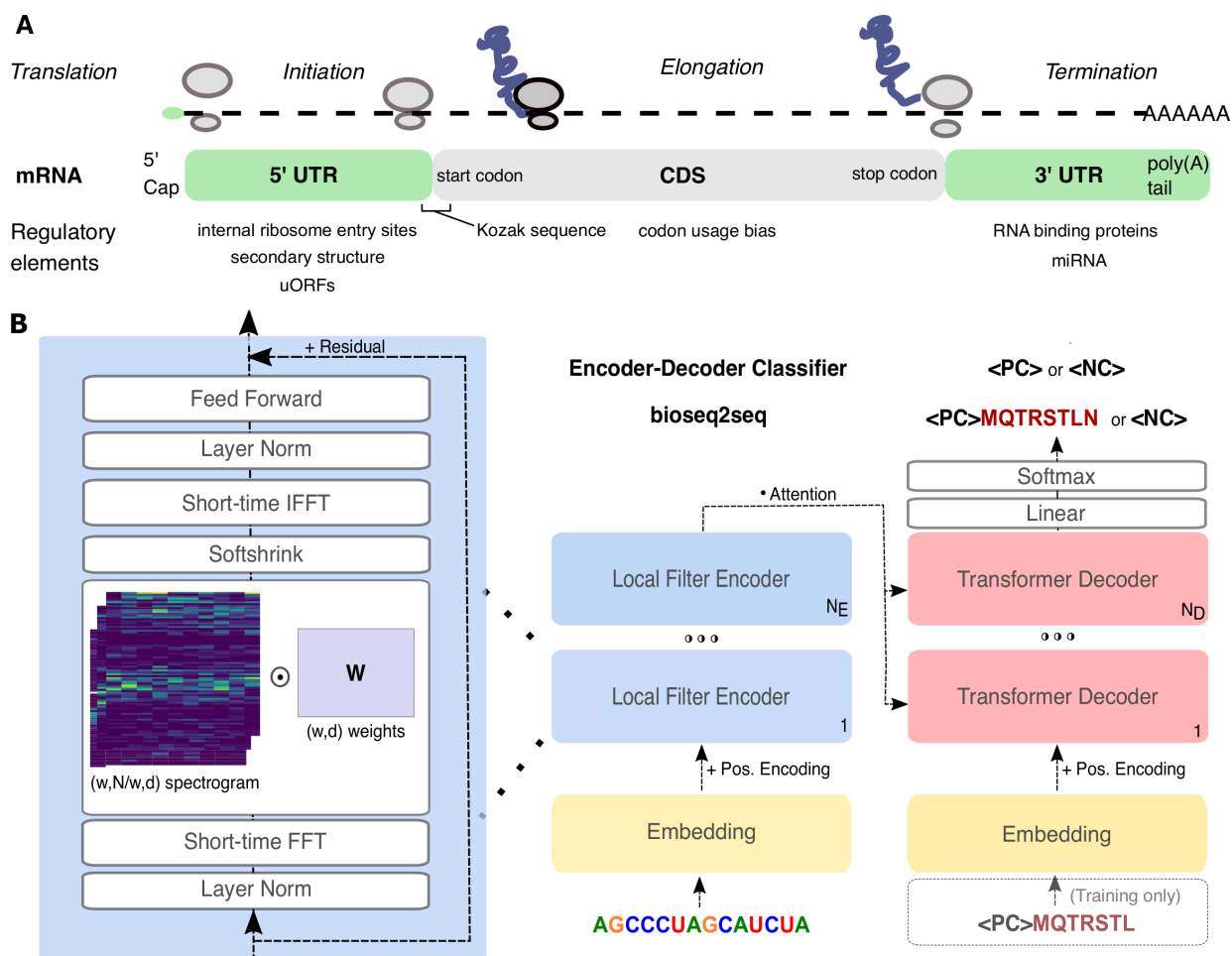


Figure 1. Overview of problem setting and computational method. (A) Summary of messenger RNA functional regions and known elements regulating translation. See (Gebauer and Hentze 2004) for a review of known regulatory elements. (B) Neural network sequence-to-sequence architecture. We designed LFNet (left) to apply a learned filter matrix \bar{W} to a 1D short-time Fourier transform (spectrogram) of the hidden representations, enabling frequency-domain filtering of the 3-base periodicity present in coding sequences. We trained this architecture for two problem settings: in Encoder-Decoder Classifier (EDC), the expected output is a classification token, for bioseq2seq, the protein translation is also predicted.

136 several hyperparameters, including hidden dimension and number of encoder and decoder layers, for
 137 bioseq2seq and EDC separately (Supplementary Table 1). Bioseq2seq performance was optimized
 138 with 12 LFNet encoder and 2 transformer decoder layers, while EDC selected 16 of each for a sub-
 139 stantially larger model. After optimizing the hyperparameters for bioseq2seq and EDC, we trained four
 140 model replicates of each from different random initializations. We also trained replicates for the EDC
 141 task using the optimal hyperparameters for bioseq2seq, referring to this as EDC-small, in contrast
 142 to the optimized EDC, which we refer to as EDC-large. We report the classification accuracy on a
 143 withheld test set for our two model types in Table 1. In the case of bioseq2seq, which produces a

Model	Accuracy	F1	Recall	Precision	MCC
EDC (small)	0.885 ± 0.017	0.896 ± 0.014	0.913 ± 0.011	0.880 ± 0.024	0.769 ± 0.034
EDC (large)	0.922 ± 0.004	0.927 ± 0.003	0.910 ± 0.011	0.945 ± 0.014	0.845 ± 0.009
bioseq2seq	0.950 ± 0.006	0.954 ± 0.005	0.953 ± 0.012	0.955 ± 0.017	0.900 ± 0.011
RNAsamba	0.957 ± 0.002	0.960 ± 0.002	0.949 ± 0.004	0.970 ± 0.002	0.913 ± 0.004
CPAT	0.939	0.944	0.947	0.940	0.876
CPC2	0.911	0.912	0.856	0.976	0.830

Table 1. Classification Performance. Bioseq2seq was compared with an EDC model whose hyperparameters were tuned independently (large) and an EDC model with identical hyperparameters to bioseq2seq (small). Several top-performing machine learning models were evaluated on our dataset for comparison. For our models, predictions were made using the leading 'classification' token ($\langle PC \rangle$ or $\langle NC \rangle$) of the first beam, terminating inference before the peptide prediction. For our models and RNAsamba, multiple replicates were trained with different random seeds. Evaluation metrics were calculated with $\langle PC \rangle$ as the positive class and listed as mean ± std. dev. where multiple replicates are available.

144 variable-length peptide decoding at inference time, decoding was halted after the leading classifica-
145 tion token was predicted. The bioseq2seq replicate with the best performance on F1 score achieved a
146 score of 0.958, while the worst-performing on this metric scored 0.947. We compared our models with
147 five replicates of RNAsamba trained on our dataset, as well as CPC2 and CPAT, two machine-learning
148 methods based on engineered features. The best model for bioseq2seq exceeds the performance of
149 CPC2 and CPAT and is competitive with RNAsamba (0.956-0.961 F1) without explicit inclusion of
150 any auxiliary features such as ORF k-mers, although RNAsamba appears slightly better according to
151 all evaluation metrics except recall. EDC-large ranged in performance between 0.924-0.932 in F1.
152 EDC-small was clearly the worst of all models and so from this point we will only consider EDC-large
153 and refer to it simply as EDC. The markedly better performance of bioseq2seq in comparison to its
154 classification-only analogues makes it clear that the translation task improves the performance of an
155 LNet model on the binary classification task.

156 As bioseq2seq is capable of performing translation on top of classification, we also report the
157 percentage identity between the ground truth protein and the translation produced by bioseq2seq
158 using the Needleman-Wunsch global alignment. A large majority, 82.4 %, are exact matches with the
159 ground truth. Notably, when bioseq2seq was allowed to predict a full-length protein rather than halted
160 after the classification token as in the results from the previous section, the classification performance
161 of the best model deteriorated slightly to 0.940 F1. This suggests a slight trade-off at inference time
162 between an accurate peptide decoding and the classification task, though the bioseq2seq training
163 strategy as a whole clearly improves classification performance relative to EDC.

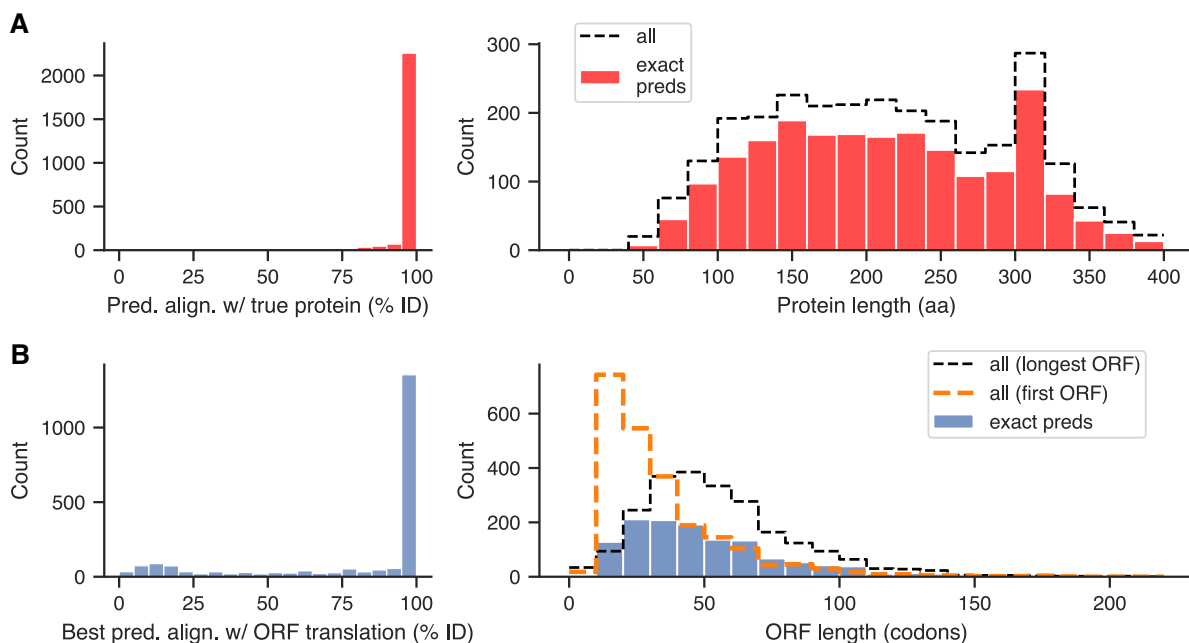


Figure 2. Analysis of translation products predicted by best bioseq2seq replicate. (A) Global alignment identity between the top-beam protein decoding predicted by bioseq2seq for true positive mRNAs and the ground truth protein (left), and length distribution of perfect translations (right). Black dashed line indicates the complete distribution of protein lengths. (B) Highest global identity found from all-by-all alignment of the three-frame translation of a lncRNA with its lower-beam $\langle PC \rangle$ + peptide predictions from bioseq2seq (left) and length distribution of perfectly translated sORFs (right). Black dashed line indicates the length distribution of hypothetical translations of the longest ORF found in each lncRNA and orange dashed line denotes the same for the most 5' ORF.

164 2.2 Alternate decodings of lncRNAs harbor plausible micropeptides

165 The bioseq2seq formulation can produce and rank multiple candidate decodings for a given RNA us-
166 ing beam search. For sequences annotated as lncRNAs and correctly classified by bioseq2seq, the
167 lower beams (second highest scoring and on) will with high probability begin with $\langle PC \rangle$. We investi-
168 gated the predicted peptides for insights into the potential translation of lncRNAs. First we confirmed
169 that the peptides matched a true ORF within the lncRNA by using the EMBOSS package to find the
170 top Needleman-Wunsch alignment score between the three-frame translation and all generated pep-
171 tides from a beam size of four (Rice et al. 2000). In 59.7 % of cases, the best match was a perfect
172 alignment, meaning that most peptide decodings were translations of ORFs actually present in the
173 lncRNA.

174 We applied bioseq2seq to a set of transcripts previously or currently annotated as lncRNAs but
175 considered by the database lncPEP to have been validated by supporting literature to express a mi-
176 cropeptide (Liu et al. 2022). Starting from the lncPEP "validated" set, we implemented a number of

177 quality control measures, removing redundant transcripts, linking the transcript names listed on the
178 IncPEP website with RefSeq accession numbers via the underlying primary literature and the NCBI
179 search function. This yielded twenty-two putative micropeptide-encoding transcripts (provided as Sup-
180 plementary Table 3), of which nine were found in our training set. The best model for RNAsamba pre-
181 dicted 3 of the remaining 13 to be protein coding. Using bioseq2seq, 3 were also predicted as coding
182 when terminating inference after the classification token, and 4 when running peptide decoding to
183 completion.

184 We aligned all beams from bioseq2seq with the IncPEP micropeptides and found that in most
185 cases the model also successfully identified the correct ORF to translate, with 3 of 4 predicted coding
186 transcripts having alignment identity $\geq 90\%$. If lower beams are considered, 8 have identity $\geq 90\%$,
187 including a very short 17-aa peptide. The examples found in our training set are of potential interest
188 as well because in several cases the class label that we trained on contradicts the prediction that
189 bioseq2seq makes. For example, *LINC00266-1* with NCBI accession NR_040415.1 is currently an-
190 notated as a lncRNA but was found in (Zhu et al. 2020) to express a 71-amino acid oncopeptide.
191 Bioseq2seq perfectly predicts the peptide in its highest beam – a false positive according to the class
192 label in the training data. Examples like this and NR_033874.1 highlight the generalizability of the
193 rules learned by bioseq2seq and RNAsamba, even when presented with false annotations. One ex-
194 ample, NM_001384235.1 in the training set underscores a crucial distinction between bioseq2seq and
195 prior methods like RNAsamba. In these transcripts, the micropeptide is not coded for by the longest
196 ORF. RNAsamba only explicitly considers the longest ORF in each transcript and may fail to identify
197 alternate sources of coding potential, as it does here. The translation product for AW112010.1 in the
198 test set comes from an instance of non-AUG initiation (Jackson et al. 2018), and while our method
199 cannot perfectly predict the protein product in such cases we successfully identify it as a coding tran-
200 script and predict a partial match from the canonical portion of the CDS.

201 **2.3 Local Filter Networks emphasize 3-nt periodicity**

202 The core feature of each LFNet layer is its learned frequency-domain filters. We visualized the filter
203 weight matrices to investigate the frequency response of the model to signals in the intermediate
204 vector representations, including separate plots for their magnitude $|z|$ and phase θ for the complex

Accession	RNA len	Pep. len	longest ORF?	RNAsamba correct?	bioseq2seq correct?	Beam match
NR_033874.1	810	130	✓	✓	✓	1
NM_001315494.2	828	84	✓	✓	✓	1
NR_040415.1	723	71	✓	✓	✓	1
NM_001384134.1	427	56	✓	✓	✓	1
NM_001384235.1	608	47	✗	✗	✓	1
NM_001352129.2	783	35	✗	✗	✗	2
NR_033201.2	611	53	✓	✗	✗	3
NM_001304732.2	857	46	✗	✗	✗	3
NR_046502.1	537	21	✗	✗	✗	✗
NR_003634.2	941	262	✓	✓	✓	1
AW112010.1	536	82	✗	✗	✓	1
BC030870.1	1216	71	✓	✗	✓	1
KY559104.2	2536	144	✓	✓	✗	2
NM_001348129.2	2344	68	✓	✗	✗	2
NM_001348107.3	1605	90	✗	✗	✗	3
NR_015417.1	2273	60	✗	✗	✗	3
NR_001458.3	1500	17	✗	✗	✗	3
NR_033243.1	2843	117	✗	✓	✓	✗
BK010446.1	1084	87	✓	✗	✗	✗
NM_001352687.2	1099	59	✗	✗	✗	✗
NR_038278.1	1749	52	✗	✗	✗	✗
NR_024394.1	4082	50	✗	✗	✗	✗

Table 2. Results on twenty-two validated micropeptides. Samples above the horizontal bar were in our training set and those below were not. A bioseq2seq prediction was counted as correct if it began with $\langle PC \rangle$, regardless of the official class label. The matching beam indicates the first beam peptide decoding from bioseq2seq achieving $\geq 90\%$ alignment with the annotated micropeptide, if one exists.

205 weights $z = |z|e^{i\theta}$. The resulting images for all layers in both bioseq2seq and EDC are given in Figure
206 3. Visually, the most prominent signal in both model types is a band at a frequency bin equivalent
207 to a period of 3 nt. This illustrates that most layers and hidden dimension across the LFNet stack
208 learned to emphasize the 3-base periodicity of coding regions. Notably, every layer of EDC (panel
209 B) shows a more clear dependence on the 3-nt property than bioseq2seq (panel A), with every layer
210 having a clean visual band of low magnitudes along this frequency range. In contrast, lower layers of
211 bioseq2seq do not appear to emphasize this feature. However, bioseq2seq has phase values close
212 to zero along the 3nt band (panel C), while the phase activity of EDC is somewhat more random
213 (panel D). We observed in Supplementary Fig S1 that for bioseq2seq, the periods other than 3-nt
214 are associated with phases peaked around $-\pi$ and π , which correspond to phase components of the

215 weights being $e^{i\theta} = -1$, such that the output of the LFNet layer would negate the residual when they
216 are added. While the bioseq2seq LFNet weights shift toward the positive real-axis in the higher layers
217 for three-nucleotide signals, they shift toward the negative real axis for other periods. This trend is
218 found clearly in bioseq2seq, and less so EDC, where the weights are smaller and more centered at
219 zero (Supplementary Fig S2). Furthermore, while weights corresponding to three-nucleotide signals
220 are mostly zero for EDC, creating a band in Fig 3, the weaker band for bioseq2seq is explained by
221 many positive weights in bioseq2seq at this band, which would amplify three nucleotide signals. We
222 hypothesize that the inductive bias of LFNet facilitates a reliance on the 3-base property, and the
223 translation task leads to the amplification of specific 3-base signals.

224 Three-base periodicity is also apparent in our models' encoder-decoder attention (EDA) distribu-
225 tions, which are probability weightings for encoder hidden embeddings in the context of each decoder
226 layer. We aligned each encoder-decoder attention distribution for every transcript relative to its start
227 codon and averaged to create nucleotide-resolution consensus attention metagenes. For lncRNAs,
228 we investigated the longest ORF to define metagenes and to compare and contrast mRNAs and lncR-
229 NAs in the rest of this manuscript. We considered the two classes separately and discarded relative
230 positions not present in at least 70% of the data, leaving relative position indices of (-25,+715) for mR-
231 NAs and (-131,+274) for lncRNAs. Depicted in Fig 3 are metagenes for a particular EDA head in the
232 lower decoder layer of bioseq2seq that responds very differently to mRNAs (panel E) and lncRNAs
233 (panel F), attending highly to the AUG/longest ORF in both classes but losing periodicity in lncRNAs.
234 Sharp differences in attention such as this likely implement aspects of the model's classification logic.
235 We present more detailed analysis of EDA metagenes in Supplementary Fig S3.

236 **2.4 Translation task improves reproducibility and biological plausibility of variant ef-** 237 **fect predictions**

238 We evaluated all of our model replicates on every possible single-nucleotide variant of transcripts from
239 a subset of our test data, consisting of 220 verified mRNAs and 220 verified lncRNAs. This technique,
240 known as saturated *in silico mutagenesis* (ISM), is commonly used to computationally predict variant
241 effects and can provide insight into input features that machine learning models recognize as impor-
242 tant to their predictive task (Zhou and Troyanskaya 2015; Koo et al. 2021; Tareen et al. 2022). We

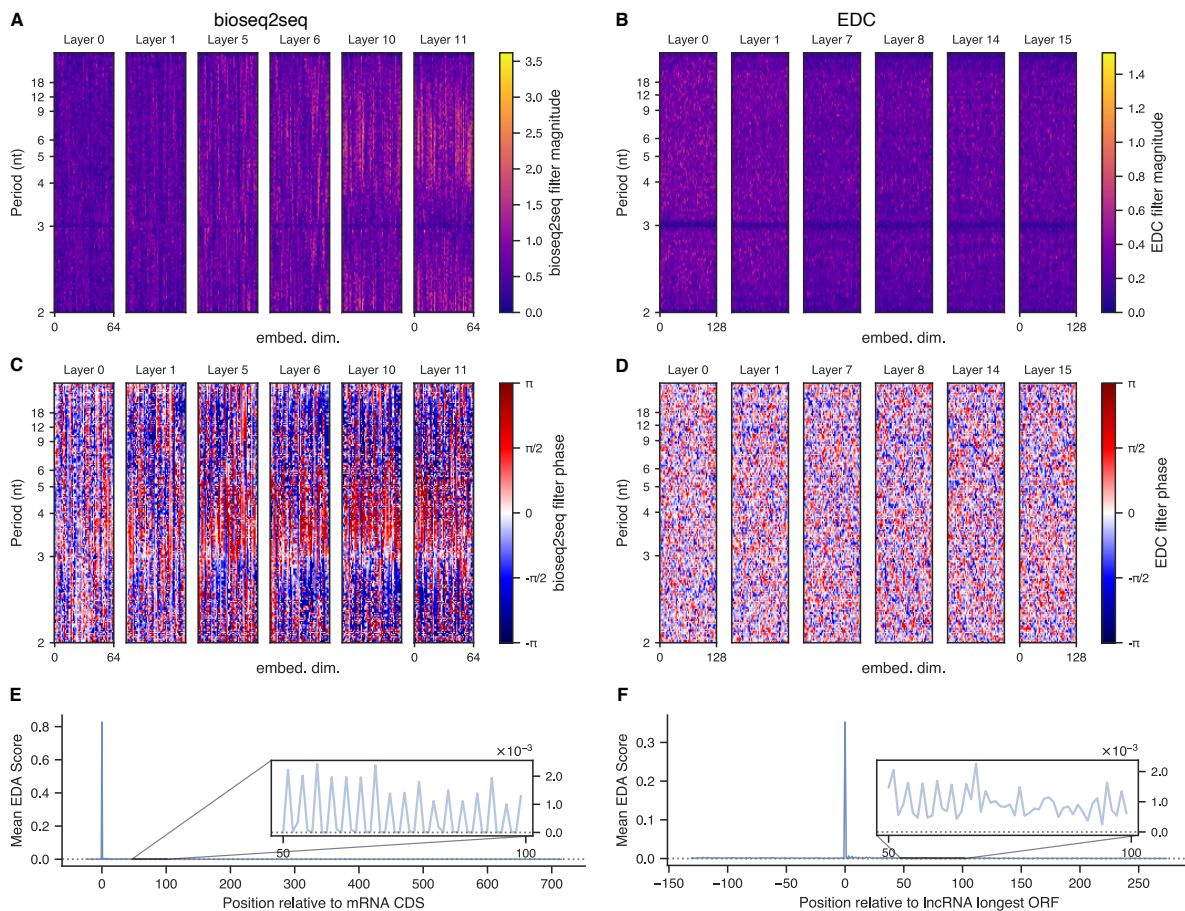


Figure 3. Frequency-domain content in model representations. LFNets filters from selected layers, with complex filter weights visualized in terms of magnitude (bioseq2seq in panel A, EDC in B) and phase (bioseq2seq in C, EDC in D). For each layer heatmap, the x-axis represents the hidden embedding dimension, and the y-axis refers to a discrete frequency bin, with annotations for the equivalent nucleotide periodicity. Both model types learned weights with a pronounced structure around 3-nt periodicity, visible mostly clearly in the phase for bioseq2seq and in the magnitude for EDC. (E) A nucleotide-resolution metagene consisting of average encoder-decoder attention scores from mRNAs aligned relative to their start codons. Attention distributions for this plot were taken from head 6 of the lower bioseq2seq decoder layer, which primarily attends to the start codon and places attention downstream of the start in a periodic fashion. (F) The equivalent plot for the same attention head applied to lncRNAs aligned relative to the start of the longest ORF, illustrating the loss of attention rhythmicity downstream of the leading spike.

243 calculated ISM using the function $\Delta S(x, x') = \log\left(\frac{P(x'=\langle PC \rangle)}{P(x'=\langle NC \rangle)}\right) - \log\left(\frac{P(x=\langle PC \rangle)}{P(x=\langle NC \rangle)}\right)$, where x and x' are
244 RNAs, with x' being a single-nucleotide variant of x . We calculated the Pearson correlation between
245 the ISM scores predicted by two different replicates for a given transcript, making pairwise compar-
246 isons between all replicates. We also computed the cosine similarity between the character-level
247 (A,G,C,U) vectors of mutation scores at each transcript position, using the median of this quantity as
248 an transcript-level summary metric that does not consider the scaling of mutation scores at different
249 positions. Both metrics were averaged across comparisons to produce a single value for each tran-
250 script, with the resulting distributions depicted in Fig 4-B. The inter-replicate agreement of bioseq2seq
251 is much higher than that of EDC in terms of Pearson correlation (median of $r= 0.813$ vs. median of
252 $r= 0.560$). The relaxed metric of median position-specific cosine similarity shows a minimal difference
253 between bioseq2seq and EDC, which suggests that the gap in reproducibility between the model
254 types is largely due to bioseq2seq's more stable ranking of positional importance.

255 We next probed the ISM scores for changes disrupting essential mRNA features. We investigated
256 the changes in score due to substantial sequence perturbations of each test mRNA by shuffling var-
257 ious functional regions. Specifically, we shuffled every 5' UTR longer than 25 nt in the verified test
258 set, using both an unrestricted shuffle and one preserving dinucleotide frequencies, and likewise for
259 3' UTRs separately. We produced another set of variants by shuffling all codons besides the start
260 and stop codon within CDS regions. This has the effect of preserving the original CDS length while
261 likely disrupting 3-nt periodicity and leading to atypical orderings of nucleotides and amino acids. We
262 calculated ΔS for each shuffled variant relative to its wild-type and found that UTR shuffling had min-
263 imal impact on on the predictions of either bioseq2seq or EDC (Fig 4-C). However, EDC is somewhat
264 more reliant on the endogenous trinucleotide patterns of wildtype CDS regions than bioseq2seq, as
265 indicated by the stronger negative ΔS after shuffling internal codons of CDS sequences. In contrast,
266 mutations to the annotated start codon tended to produce large negative ΔS scores in bioseq2seq
267 but not in EDC (Fig 4-D). Similarly, bioseq2seq responded negatively to mutations that introduced a
268 stop within the first 50 codons. These observations suggest that while both models detect periodic
269 sequence features, bioseq2seq has learned contextual sequence features, including start and stop
270 codons, that more comprehensively align with our understanding of translation.

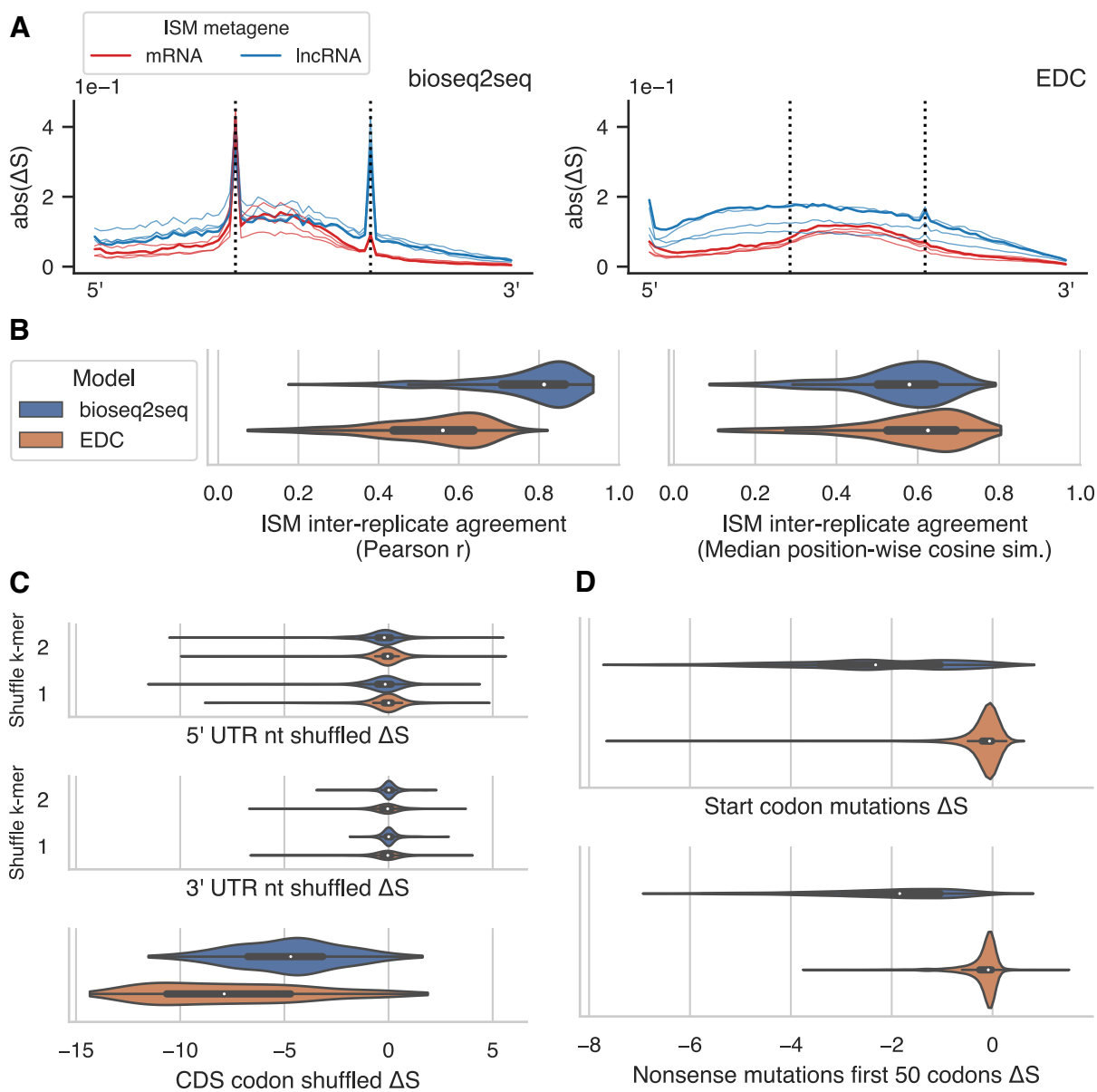


Figure 4. Predicted mutation effects by model type on a subset of testing data. (A) Metagene plots of saturated in silico mutagenesis (ISM) ΔS scores, i.e. the difference in $\log(P(\langle PC \rangle)/P(\langle NC \rangle))$ between single-nucleotide variants and their wild-type sequence. The absolute value of ΔS was averaged within each of 25 positional bins and across all three possible mutations in each position, with mRNAs and lncRNAs depicted separately for both bioseq2seq (left) and EDC (right). Vertical dashed lines denote the first and last bin of the CDS for mRNAs and the longest ORF for lncRNAs. Metagenes from all four replicates are shown, with the best-performing model colored using the darkest hue. (B) Per-transcript average of Pearson correlation (left) and median position-specific cosine similarity (right) of ISM scores from pairwise comparison of model replicates. (C) Changes in score relative to wildtype for mRNAs shuffled within each functional region. UTRs were shuffled to preserve mononucleotide or dinucleotide frequencies. Codon shuffling excluded the start and stop codons to preserve CDS length. (D) Changes in score for mRNAs from nucleotide substitutions that knock out a start codon or introduce a stop codon within the first 50 codons of the CDS. Note: panels C and D follow the legend from panel B.

271 **2.5 In silico mutagenesis reveals features predictive of coding potential**

272 In light of the gap in biological robustness between our two model types, we investigated the response
273 of bioseq2seq to sequence perturbations, using its best replicate to obtain ISM predictions for the
274 remainder of the test set. We aggregated ISM scores for all synonymous point mutations inside
275 of mRNA CDS regions into fine-grained metagenes for each amino acid, computing the mean ΔS
276 along each of 25 positional bins. Selected amino acids are highlighted in Fig 5-A and all twenty
277 are depicted in Supplementary Fig S4. As expected for a highly contextual model, there are large
278 deviations away from the mean. On average however, the amino acids with only two codons all learn a
279 preference for a single codon across the length of the whole transcript, with correspondingly negative
280 scores for the opposite mutation. The amino acids with more than two-fold degeneracy are more
281 complex to interpret but the sign for the mean mutation effect tends not to change with position. When
282 considering all synonymous mutations, the model appears to have learned a preference for particular
283 nucleotides in the codon positions. For example, most codons ending in C having a positive effect
284 on ΔS on average, and most ending in T having a negative effect (Fig 5-B). Bioseq2seq's estimates
285 of synonymous mutation effects also captured some of the variation from an external measure of
286 translation efficiency called tRNA Adaptation Index (tAI) (Reis et al. 2003). The mean ΔS for point
287 mutations leading to synonymous changes show a moderate correlation ($r = 0.394$, $\rho = 0.418$) with
288 the differences in tAI between the two codons, using codon values calculated from (Tuller et al. 2010).

289 We used ISM scores as a feature explanation method by assigning each nucleotide within a tran-
290 script an importance score based on the magnitude of ΔS from the mutation in that position that most
291 disrupts bioseq2seq classification towards the opposite class. For example, an endogenous x_i within
292 an mRNA was defined as contributing towards a true positive classification of the $\langle PC \rangle$ class to the
293 extent that substituting any of the three alternate bases in position i produced a highly negative ΔS .
294 One representative example mRNA and lncRNA are visualized in Fig 5-C and D, respectively, with
295 raw ISM scores from positions of interest shown in a heatmap. The transcript sequences are overlaid
296 above with their heights drawn proportionally to the importance setting for their true class – $\uparrow PC$
297 for the mRNA and $\uparrow NC$ for the lncRNA. The samples were chosen from among the five lncRNAs
298 and mRNAs closest to the median value for inter-replicate agreement (see Fig 4-B). In the example
299 mRNA, the start codon is a highly salient region, while the stop codon receives little importance. The

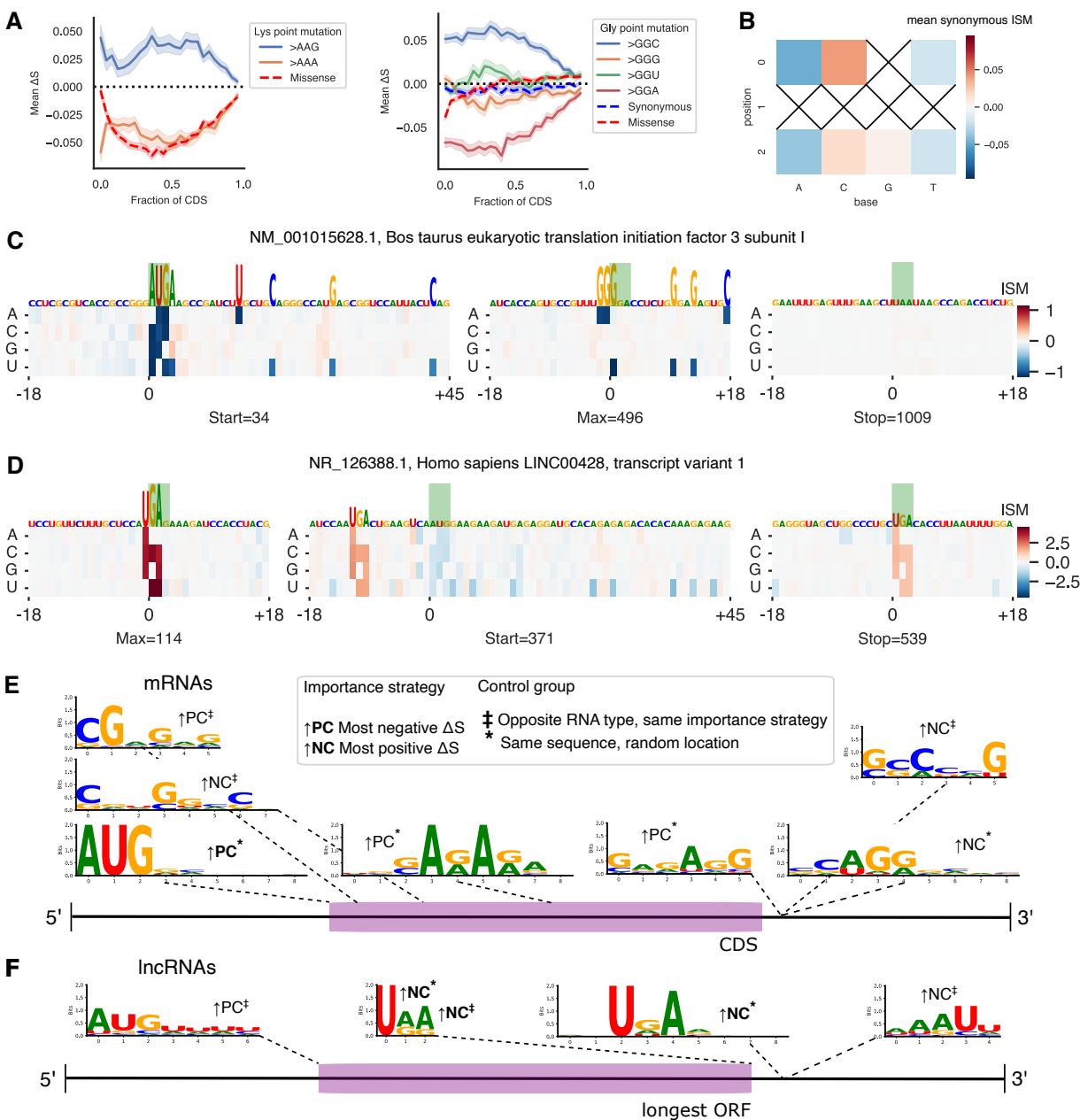


Figure 5. Detailed analysis of *in silico mutagenesis* (ISM) on the full test set. (A) Plots of ISM metagenes for selected amino acids lysine (left) and glycine (right). Mean ΔS is shown for 25 positional bins across mRNA CDS regions with mutations listed based on the resulting codon. The red line represents the average across all missense/nonsynonymous mutations. For amino acids with more than two codons, the blue dashed line depicts the average synonymous mutation for comparison. (B) Mean ISM for synonymous point mutations by codon position and nucleotide. X's denote substitutions which do not exist as synonymous changes. (C) An example protein-coding transcript with NCBI accession NM_001015628.1. Signed ISM scores for the transcript are depicted as a heatmap and the RNA sequence is portrayed with characters scaled according to the $\uparrow PC$ importance strategy, i.e. regions with highly negative ISM weights depicted in dark blue. The subregions shown are windows around the start codon, the position of maximum importance, and the stop codon, respectively. (D) Same as panel B with an example long noncoding RNA with NCBI accession NR_126388.1. The endogenous sequence is scaled according to $\uparrow NC$, or highly positive ISM values drawn in dark red. (E) mRNA motifs discovered in our test set with STREME using ISM importance values from bioseq2seq to determine sequence regions in which to search for enriched signals. Annotations denote the importance and control strategy for each trial, with boldfaced annotations signifying that importance values were not masked and ordinary typeface indicating that feature importance at start and stop codons and nonsense mutations were excluded. Motifs are positioned near the regions in which they were enriched. (F) Same as panel D showing discovered lncRNA motifs.

300 ISM scores for the nucleotides surrounding the start codon imply a preference for G in position +1
301 relative to the start and A or C in position -2, consistent with the Kozak consensus sequence. The
302 most important feature occurs in a region where many possible point mutations would introduce a stop
303 codon, and we observed widespread avoidance of nonsense mutations early in the coding sequence.
304 For the lncRNA, the TGA ending the longest ORF receives high importance according to $\uparrow NC$, but a
305 different TGA upstream of the longest ORF is the highest overall.

306 To systematically extract general patterns that bioseq2seq recognizes as predictive of coding po-
307 tential, we performed de novo discovery of motifs frequently found in transcript subsequences with
308 high ISM importance. First, we identified the most important nucleotide with respect to both $\uparrow PC$
309 and $\uparrow NC$ from each functional region (5' UTR, CDS, 3' UTR) of each test-set mRNA and likewise
310 for the regions demarcated by the longest ORF of a lncRNA. We extracted 21-nt windows centered
311 around each such important site to form a primary sequence database for the differential motif discov-
312 ery tool STREME (Bailey 2021). A control set for STREME was constructed either using (1) random
313 positions from the same transcript and region as the primary sequences but not overlapping them
314 (2) the most important positions using the same importance setting as in the primary sequence but
315 from the opposite RNA class. These controls necessitate different interpretations of the discovered
316 motifs, with strategy 1 intended to establish whether bioseq2seq places importance on consistent fea-
317 tures of a transcript, and strategy 2 intended to uncover differences in how bioseq2seq treats roughly
318 comparable regions of coding and noncoding transcripts. We also ran motif discovery using a purely
319 random strategy – e.g. with randomly chosen subsequences of a 5' UTR as primary and random
320 upstream regions of a lncRNA as control. We present only strategy 2 motifs that do not match a motif
321 from the purely random trials according to TOMTOM (Gupta et al. 2007), as these experiments were
322 specifically guided by ISM importance.

323 We ran every combination of primary sequence region, control method, and importance setting as
324 its own STREME experiment and discovered four significant motifs between mRNAs and lncRNAs. Fi-
325 nally, we ran a second set of experiments in the same manner except with importance for endogenous
326 start and stop codons and counterfactual missense mutations masked out in order to reveal important
327 signals beyond the most prominent set found in the first run. This yielded an additional seven motifs,
328 and both sets are shown in Fig 5-E for mRNAs and F for lncRNAs, with boldface annotations for the

329 unmasked motifs. The experiments with random controls largely confirm the observations we made in
330 our example transcripts, with a start codon/partial Kozak motif found in the beginning of mRNA CDS
331 regions and several stop codon motifs prominent throughout lncRNAs. Beyond this, repeated GA
332 patterns appear enriched in regions that push mRNAs towards a true positive classification and both
333 control strategies uncover motifs that push mRNAs towards a false negative. Similarly, As and Us
334 downstream of AUGs influence bioseq2seq towards a false positive prediction on lncRNAs, but such
335 a motif receives comparatively little weight in the model's assessment of bona fide coding transcripts.
336 Additional details including positional and frame biases and enrichment, can be found in Supplemen-
337 tary Tables S4 and S5. We note a potential match with the binding site motif for an RNA-binding
338 protein *ACO1* from (Ray et al. 2013), listed as motif #1 in Supplementary Table S5.

339 **2.6 Approximation quality of gradient-based mutagenesis depends on model com-** 340 **plexity**

341 Saturated ISM is costly to apply to a large amount of sequences because it requires $3L$ model eval-
342 uations, where L is the transcript length. We explored the feasibility of approximating ISM using
343 neural network input gradients, which are efficiently computable in parallel via automatic differentia-
344 tion. Building from the Integrated Gradients (IG) method, we developed a novel proxy for ISM called
345 Mutation-Directed Integrated Gradients (MDIG). MDIG involves numerically integrating input-output
346 gradients along the linear interpolation path between a sequence of interest and a sequence of the
347 same length consisting of all the same type of nucleotide, e.g. all guanines. A parameter $\beta \in (0, 1]$
348 limits how far to travel towards the poly(b) baseline embedding during integration. (See Methods). As
349 a favorable value for β is not obvious from first principles, we tuned this parameter on a subset of our
350 validation set consisting of 206 verified mRNAs and 206 lncRNAs, applying the same criteria from
351 the previous section. To benchmark attribution stability across stochastic training, we measured the
352 inter-replicate agreement of each mutation approximation method using Pearson correlation as in the
353 previous section. We then computed the per-transcript Pearson correlation of scores from different
354 settings of MDIG- β with the ISM scores from the same replicate. This metric indicates MDIG's capac-
355 ity to approximate the input-output behavior of a given deep learning model, which ISM accomplishes
356 directly but at substantially greater computational cost. For reference with other gradient-based per-

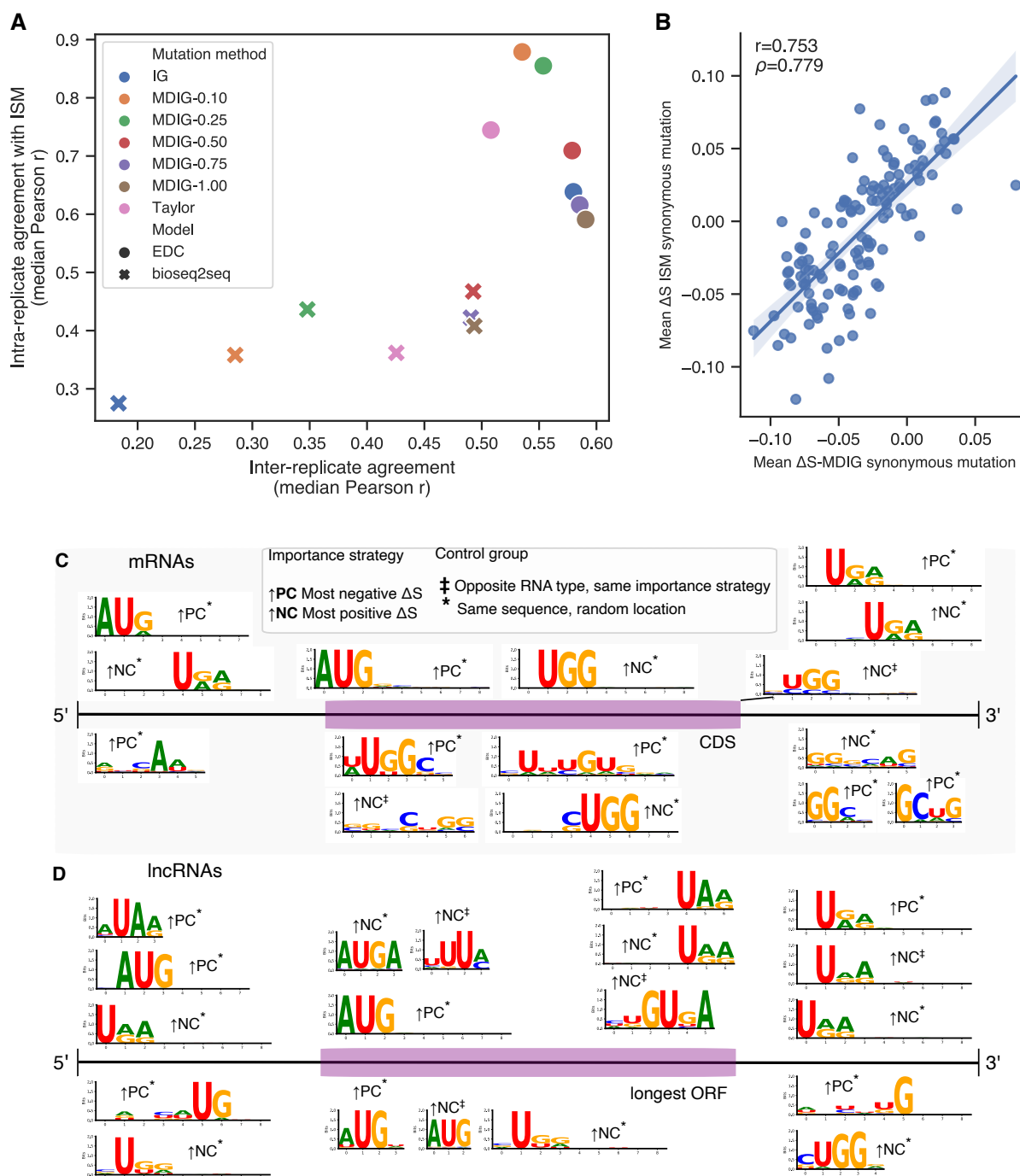


Figure 6. Gradient-based approximation performance. (A) Summary results from tuning of β hyperparameter for MDIG alongside baseline methods. Inter-replicate agreement is shown on the x-axis and correlation with ISM on the y-axis, using the median across transcripts as a point estimate for both metrics. (B) Scatter plot of ΔS for all possible synonymous point mutations, i.e. every wildtype>variant pair differing at one position, from MDIG on the training set (x-axis) versus the same for ISM on the test set. (C) mRNA motifs discovered in our training set with STREME using MDIG importance values from bioseq2seq to determine sequence regions in which to search for enriched signals. Results from unmasked importance are shown above the transcript diagram and those from the masked trials are shown below. (D) lncRNA motifs discovered in the training set using MDIG importance values from bioseq2seq, depicted in the same manner as panel C.

357 turbations, we perform the same analyses using a first-order Taylor approximation of ISM scores and
358 IG with a uniform $[0.25, 0.25, 0.25, 0.25]$ baseline. Results on the validation set according to these
359 evaluation metrics are summarized with their median value in (Fig 6-A), and full violin plots in Sup-
360 plementary Fig S5. On the basis of these results, MDIG-0.5 was selected as the best approximation
361 method for bioseq2seq and MDIG-0.1 for EDC. This illustrates that the MDIG method can predict the
362 effect of input perturbations better than the basic Taylor approximation.

363 On the whole, we observed a large gap in approximation quality between the model types, with
364 the best method for bioseq2seq lagging substantially behind the worst for EDC. To investigate the
365 implications of MDIG's reduced performance on bioseq2seq, we used the test data and method from
366 the previous section to compute bin-based metagenes from the best MDIG versions and observed
367 that this averaged representation closely captures the same general trends as expected from ISM
368 (Supplementary Fig S6). Across the bioseq2seq replicates, MDIG metagenes have an average cor-
369 relation of $r = 0.897$ for mRNAs and $r = 0.944$ for lncRNAs with their ISM equivalent, in comparison
370 to $r = 0.999$ and $r = 0.997$, respectively for EDC. For a more detailed evaluation on bioseq2seq we
371 approximated ΔS for every synonymous point mutation using MDIG on the test set and compared it
372 with the true ΔS scores from ISM in the form of a scatterplot in Fig 6-B. The high correlation between
373 metagenes and codon scores for ISM and MDIG indicates that despite its reduced transcript-level ac-
374 curacy in predicting bioseq2seq mutation effects, MDIG largely captures the same class-level features
375 as ISM when averaged across examples.

376 To take advantage of MDIG's improved efficiency relative to ISM for improving the statistical power
377 of motif discovery on our most biologically robust model, we applied MDIG-0.5 on bioseq2seq for the
378 full training set. This consists of $\sim 52k$ examples, balanced between the two RNA classes, and took
379 about two days of GPU-time, much faster than our extrapolated estimate of more than a month for
380 ISM (See Supplementary Table 2). We used the resulting MDIG mutation effect estimates as drop-in
381 replacement for ISM importance values in our motif discovery pipeline, with the results presented in
382 Fig 6-C for mRNAs and D for lncRNAs. Motifs from masked trials are placed below the transcript
383 diagrams and those from unmasked trials are above. In comparison to the ISM motifs discussed
384 previously, the MDIG motifs better underscore that bioseq2seq places importance on start and stop
385 codons in regions besides the CDS. Start codons are predicted by MDIG to increase bioseq2seq

386 coding probability in both mRNA 5' UTRs and lncRNA upstream regions, while stop codons push
387 the classification towards noncoding in the 5' regions. Notably, the UTR motifs typically lack a bias
388 towards a particular frame of the transcript, while most ORF features have a consistent frame bias.
389 This is supportive of the idea that such elements outside the ORF are flagged in part to determine
390 the frame. A number of interesting mRNA motifs emerge from masking, including multiple strong
391 UGG motifs in a variety of sequence contexts and positions. The masked lncRNA motifs closely
392 resemble those from the unmasked strategy, implying that the masked maximums are nucleotides
393 adjacent to the start and stop codons. This comparative lack of diversity could mean that bioseq2seq
394 largely defines lncRNAs as a class in terms of a lower quality or incorrect context of protein-coding
395 features rather than distinctly 'noncoding' features. It also likely implies that MDIG is most adept at
396 estimating the strongest mutation effects for bioseq2seq, with diminished reliability for less influential
397 signals. One the whole, aggregating over instance-level MDIG scores to drive motif discovery appears
398 to emphasize broad global features on which both MDIG and ISM both place high importance, while
399 revealing additional signals beyond those identifiable with smaller-scale ISM experiments alone. As
400 for ISM, the MDIG motifs are shown in greater detail in Supplementary Tables S6 and S7. We note
401 a potential motif match to a binding site for an RNA-binding protein *SAMD4A* from (Ray et al. 2013)
402 discovered in mRNA 3' UTRs as 'motif 1' in Supplementary Table S7 and alongside possible ISM
403 matches in Supplementary Fig S7.

404 **3 Discussion**

405 The genetic code makes it straightforward to predict protein sequences given an mRNA sequence,
406 but our results suggest that requiring a neural network to learn the translation task improves its abil-
407 ity to identify protein-coding RNAs. We hypothesize that translation acts a regularization strategy by
408 requiring the model to preserve precise positional information in a way that improves its contextual
409 representations. Our findings are consistent with a related observation from RNAsamba, which per-
410 formed worse when a network branch processing the longest ORF sequence was ablated (Camargo
411 et al. 2020). Bioseq2seq differs from RNAsamba in that the translated protein sequence is an output
412 rather than an input to the network. To our knowledge, bioseq2seq is the first attempt to use machine

413 learning to output the encoded protein for an input RNA by explicitly learning the sequence mapping
414 underlying biological translation. It accomplishes this from sequence alone, without introducing prior
415 knowledge about the genetic code. Our models remain competitive with the best prior approaches
416 without engineered sequence features, with bioseq2seq achieving on-par accuracy (less than 1%
417 difference) and a higher recall.

418 The translation task also appears to significantly improve the quality of the nucleotide-level fea-
419 tures identified by our models as predictive of protein-coding potential. The correlation of ISM mutation
420 effects across multiple replicates is considerably higher for bioseq2seq than for EDC. Inter-replicate
421 agreement quantifies the low epistemic uncertainty of mutation effect predictions made by an en-
422 semble of bioseq2seq models. In the absence of experimentally characterized mutation effects, this
423 suggests a robustness in the learned biological rules that can inform the plausibility of insights de-
424 rived from feature interpretation. Besides this improvement in feature consistency, we found that the
425 translation task confers an additional context-awareness to the model in a way that matches biological
426 intuition. Even though simple features like ORF length are obvious correlates of ribosomal translation
427 activity in the cell, the training process does not automatically impart this mechanistic insight into a
428 neural network. We observed that EDC did not respond strongly towards mutations to either start
429 codons or premature stop codons, suggesting such elements play a minimal role in its classification
430 logic despite its relatively high best-case performance of 0.932 F1. Similarly, although mRNN rec-
431 ognizes start codons, it responded primarily to certain codons found 100-200 nt downstream of the
432 start codon, rather than waiting for the stop codon (Hill et al. 2018). Bioseq2seq, however, responds
433 negatively to start codon mutations, stop codon mutations, and nonsense mutations, suggesting that
434 its decision-making is strongly influenced by its learned ORF features. Bioseq2seq's faithful modeling
435 of ORF features and mRNA periodicity improves the chances that it also makes biologically relevant
436 effect predictions with respect to synonymous mutations and motif discovery, which require greater
437 detail. We believe that the translation task steers the network toward more robust and meaningful
438 representations that align with biological knowledge and show relative stability across replicates. In
439 our view, these properties are vital prerequisites to enable a broader reliance on machine learning
440 feature interpretation as a tool for scientific discovery.

441 Our treatment of gradient-based attributions is a contribution to the ongoing debate in the ma-

442 chine learning literature about the trustworthiness of such methods as neural network explanations.
443 We benchmark gradient-based mutation effect predictions in the biological sequence domain against
444 *in silico mutagenesis*, which is the concrete model response to meaningful sequence perturbations.
445 Strikingly, the translation task appears to adversely affect the quality of gradient approximations, with
446 all methods achieving relatively poor correlation with ISM for bioseq2seq but acceptable approximation
447 quality in EDC. At a minimum, our results suggest that users of gradient-based feature explanations
448 for genomics should follow a protocol similar to ours to validate gradient-based mutation effect predic-
449 tions against more expensive but direct input perturbations. It might suggest that for some problems
450 it is better to restrict architecture choices to convolutional neural networks, for which speedups of
451 ISM exist (Schreiber et al. 2022). More fundamentally, there could be a practical trade-off between
452 model complexity and accurate gradient approximation such that reduced fidelity of fast model pertur-
453 bations is a price to pay for the superior classification performance and biologically plausible feature
454 importance values that we observed in bioseq2seq.

455 We also introduce MDIG as a novel heuristic approximation for ISM, which we demonstrate can
456 improve over Taylor approximation at a constant increase in computational complexity. MDIG is largely
457 based on IG, but uses a more realistic mutation-specific baseline, and only integrates part of the way
458 to the baseline value, staying closer to the original sequence. Despite the limited capacity of MDIG to
459 estimate bioseq2seq mutation effects at the local, i.e. transcript level, we show its utility for identifying
460 the most impactful sequence features at the global, i.e. class-wide level. This is supportive of recent
461 work finding that the usefulness of approximate feature attributions can be improved by ensembling
462 across alternative models (Gyawali et al. 2022). The similarity of important motifs and metagene
463 representations derived from MDIG to their ISM analogues indicates that in aggregate MDIG retains
464 interpretive value even where it does not faithfully model every individual mutation effect. Subject to
465 appropriate validation, MDIG could be used where large-scale ISM experiments are infeasible or as a
466 first-pass method to flag interesting sequences for more detailed review.

467 Interpreting bioseq2seq using ISM and MDIG revealed putative signals of regulatory information,
468 which emerged purely from the learning process without prior specification. From a certain point of
469 view, learning the sequence features that distinguish translated mRNAs from lncRNAs with untrans-
470 lated ORFs would be informative for promoting ribosomal engagement and would promote translation.

471 We therefore expect that sequence features predicted to increase coding potential will correlate with
472 codon bias. Common methods for assessing codon usage bias, such as Codon Adaptation Index,
473 predict coding sequences according the relative skew of synonymous codons for a particular acid to-
474 wards the codons most common in highly expressed genes. Bioseq2seq learned strong preferences
475 within synonymous groups, as evidenced by consistently high mean value of ΔS across the entire
476 transcript for specific codons. Codon preferences were noticeably grouped by the nucleotide in the
477 third codon position, with substitutions towards nearly all codons ending in C having a positive mean
478 effect, while nearly all ending in T/U have a negative effect. The existence of codon preference trends
479 along the length of the transcript could reflect the fact that synonymous codon usage is known to be
480 biased positionally, including towards rare codon clusters (Chaney et al. 2017). Replacing codons
481 with those preferred by bioseq2seq in the average case could perform a similar function to optimizing
482 based on CAI, but bioseq2seq learns mutation effects in the context of a codon's transcript position
483 and sequence neighborhood. Our mutation effect predictions are therefore a much richer source
484 of information, and future work could test via experiment whether these preferred mutations impact
485 translational efficiency and have potential to guide mRNA sequence optimization. The discovered mo-
486 tifs also reflect sensible biological intuitions, with the MDIG motifs in particular emphasizing upstream
487 AUGs as increasing coding potential and stop codon trinucleotides as decreasing coding potential.
488 This is consistent with evidence that upstream ORFs act to suppress the translation of the main ORF
489 (Johnstone et al. 2016). Our motifs have a number of possible matches to RNA-binding proteins
490 (RBP), which play essential roles in regulating transcript stability and translational activity. A potential
491 match to the binding motifs for *SAMD4A*, a human RBP from the CIS-BP-RNA database (Ray et al.
492 2013) involved in the regulation of mRNA translation, was highlighted within regions of mRNA 3'UTRs
493 which increase coding probability according to MDIG, consistent with the model treating this binding
494 site as a valuable marker of coding potential. Several mRNA motifs reflect the Kozak sequence, and
495 we find a contrasting pattern downstream of lncRNA AUGs with downstream Us and As which locally
496 improves coding potential but is ultimately depleted in true protein coding sequences. The UGG trin-
497 ucleotide recurs across several MDIG motifs in a variety of sequence contexts and positions. This
498 could be explained in a number of ways: UGG is the unique codon for tryptophan, the rarest amino
499 acid (Barik 2020), and is also one mutation away from the stop codons UGA and UAG.

500 Our demonstration that bioseq2seq can recover potentially translated micropeptides is a proof-
501 of-concept for using machine predictions to explore this cryptic space of the proteome. Though the
502 recovery rate of putative micropeptides from IncPEP is low overall, any such capability is incidental
503 to our training setup and bioseq2seq mildly outperforms RNAsamba on the available data. Crucially,
504 bioseq2seq is not inherently limited to only translating the longest ORF, which could prove to be a
505 modeling advantage for this application given that many micropeptides are known to be harbored
506 in ORFs other than the longest in a transcript (Makarewich and Olson 2017). Increased availability
507 of validated micropeptide annotations and improved procedures for autoregressive decoding – see
508 (Yang et al. 2018) for an example – could help a future method based on bioseq2seq to achieve
509 higher reliability.

510 We anticipate that the LFNet architecture will be of broad utility in biological sequence model-
511 ing tasks, with frequency-domain multiplication enabling larger context convolutions than in common
512 convolutional architectures and lower computational complexity of $O(N \log N)$ in comparison to trans-
513 formers. Our extension of GFNet from (Rao et al. 2021) bridges older signal processing approaches
514 for gene discovery with the flexibility of deep models. We also note the complementarity of our method
515 with (Tseng et al. 2020), which, instead of employing the Fourier-transform as a token-mixing method,
516 used it to enforce a smoothness prior for importances on biological sequence models. Other appli-
517 cations of LFNet could include biological sequence data with variable periodic signals, such as nu-
518 cleosome positioning (Epps et al. 2011) and gene organization (Wright et al. 2007), as well as other
519 periodic non-biological data such as music. We designed the LFNet architecture based on an intu-
520 ition that it could effectively leverage 3-nt periodicity, but such periodic structure is not necessarily an
521 inherent requirement – the GFNet model was originally intended for computer vision.

522 There are numerous possible follow-up directions based on this work. Future versions could scale
523 to a larger and more phylogenetically diverse dataset beyond the eight mammalian transcriptomes
524 used here, as well as to longer sequence lengths. In this work we have treated coding potential as a
525 binary classification problem, but the methods presented are readily applicable to the more general
526 problem of predicting translational efficiency as a regression problem. The periodicity inductive bias in
527 particular is likely to transfer to this task – Ribo-Seq data is also characterized by a 3-nt periodicity of
528 footprint density, and this has informed the development of many ribosome profiling data analysis tools

529 (Calviello et al. 2016; Xu et al. 2018). The regression setting could also increase the prospects for
530 discovering novel regulatory features, such as in the UTRs, which our model treated as less important
531 than the CDS. A network trained to stratify transcripts according to a quantitative measure of protein
532 expression would likely learn more fine-grained distinctions than one modeling a binary separation
533 between mRNAs and lncRNAs. Finally, our results raise the possibility that general-purpose nucleic
534 acid language models could benefit from joint training with protein foundation models in a similar
535 translation-like setup.

536 **4 Methods**

537 **4.1 Seq2seq architecture for translation**

538 Our model follows the encoder-decoder sequence-to-sequence (seq2seq) framework common in ma-
539 chine translation of natural languages (Vaswani et al. 2017). We call the model bioseq2seq because
540 it applies the seq2seq paradigm to biological translation — with nucleotides and amino acids rather
541 than human languages as the vocabularies. The output of bioseq2seq is a classification token $\langle PC \rangle$
542 for protein coding and $\langle NC \rangle$ for noncoding – followed by the translated protein in the case of $\langle PC \rangle$
543 and nothing in the case of $\langle NC \rangle$. Note that the network is not provided the location of the CDS, so it
544 must learn to identify valid ORFs and select between potential protein translations.

545 Training bioseq2seq in this way allows us to test the hypothesis that the translation task will re-
546 quire the model to learn precise representations of each nucleotide, which will in turn help to attribute
547 model decisions to specific sequence patterns. As a comparison with bioseq2seq, we also trained a
548 model for binary classification. This secondary model, which we denote as Encoder-Decoder Clas-
549 sifier (EDC), has an identical network design to bioseq2seq, but a different training data format, as it
550 was trained to output only the classification token without the additional protein product for mRNAs¹.
551 We developed our models in PyTorch based on a fork of the OpenNMT-Py repository for machine
552 translation (Klein et al. 2018).

¹Although including a decoder is somewhat atypical when producing a single output classification, we do this to enable a direct comparison between the training tasks under a common architecture. The role of the decoder in the EDC setting is to calculate multi-headed attention distributions over the encoder hidden states, with the pre-pended 'start-of-sentence' token playing a similar role to the '[CLS]' in encoder-only classification setups.

553 4.2 Local Filter Network

554 We initially experimented with transformer neural networks (Vaswani et al. 2017) for both the en-
555 coders and decoders but failed to produce competitive models, as biological sequences incur exces-
556 sive memory costs as model sizes and sequence lengths grow. In these experiments, we found that
557 the transformer encoders for bioseq2seq learned self-attention heads which principally attended to a
558 small number of relative positional offsets while calculating the input embeddings. Additionally, feature
559 attributions showed evidence of a strong 3-nucleotide periodicity (See Supplementary Fig S8).

560 A variety of recent papers have introduced efficient architectures which aim to preserve the ability
561 of transformers to globally mix information at lower computational cost. A number of these approaches
562 have used the Fourier transform as a substitute for self-attention, because it is an efficient global
563 operation computable in $O(N \log N)$ time via the fast Fourier transform (FFT) algorithm (Lee-Thorp
564 et al. 2021; Guibas et al. 2021). One such example for computer vision is the Global Filter Network,
565 which takes the FFT of image patches and applies a learnable frequency-domain filter via elementwise
566 multiplication, before inverting the FFT to return the representation to the time domain.

567 As the 3-base periodicity property is localized to coding regions within transcripts, we propose a
568 simple modification to Global Filter Networks by substituting the global FFT with the short-time Fourier
569 transform (STFT). While GFNet operates on non-overlapping patches of the input, we follow common
570 practices for STFT using a stride equal to half the window size and weighting with the Hann function.
571 To emphasize that our modification applies time-frequency analysis to sequence representations, we
572 refer to this layer as a Local Filter Network (LFNet). A learned weight matrix W is applied equally
573 to each window of the STFT and then the modified frequency content is returned to the time domain
574 via the inverse FFT. A residual term is added to the result to carry along the previous representa-
575 tion. Following (Guibas et al. 2021), we apply the soft-shrink function after the weight multiplication
576 to promote sparsity in the LFNet weights. LFNet layers are only used in the encoder stack of our
577 networks, while the decoder stack consists of transformer decoder layers. This is because PyTorch
578 currently lacks an implementation of causal masking for FFT, as would be necessary to efficiently train
579 an autoregressive model with only LFNet layers.

580 **4.3 Dataset**

581 We built training and evaluation data sets using available RefSeq transcript and protein sequences for
582 eight mammalian species: human, gorilla, rhesus macaque, chimpanzee, orangutan, cattle, mouse,
583 and rat from RefSeq release 200 (O’Leary et al. 2016). We collected all RNA sequences annotated as
584 mRNA or lncRNA and excluded transcripts over 1200 nucleotides (nt) in, which reduces the available
585 data to 63,272 transcripts. Next, we linked each mRNA with the protein translation identified by Ref-
586 Seq and partitioned the data into 80/10/10 training/validation/testing splits. To maximize the diversity
587 of the dataset, we included transcripts with predicted coding status (XR_ and XM_ prefixes in Ref-
588 Seq), as well as the curated transcripts (NM_ and NR_). For the training set, we used a balanced split
589 between mRNAs and lncRNAs, selecting the split to equalize the length distribution of the two classes
590 as much as possible. Finally, we ran CD-HIT-EST-2D to exclude from the test set all transcripts that
591 exceed 80 % similarity with any transcript in the training set (W Li and Godzik 2006). The resulting
592 test set contains 2288 lncRNAs and 2703 mRNAs.

593 **4.4 Hyperparameter tuning and training**

594 We used dynamic batch sizes, so that RNA-protein training pairs were binned based on approximate
595 length to reduce the amount of padding. The maximum number of input tokens per batch was set to
596 9000 for both model types, and eight steps of gradient accumulation was used to increase the effective
597 batch size. All models were trained to minimize a log cross-entropy objective function computed from
598 each amino acid character in the output.

599 The hyperparameters including number of encoder and decoder layers, model embedding dimen-
600 sion, learning rate schedule, and L1 sparsity parameter were tuned via the Bayesian Optimization
601 Hyperband (BOHB) algorithm provided in the Ray Tune library (Liaw et al. 2018). Candidate models
602 were trained in parallel on four Tesla M10 GPUs with 8 GB GPU RAM and 640 CUDA cores, with one
603 GPU per model. To enable a fair comparison between the bioseq2seq and EDC training objectives,
604 hyperparameter tuning was run for each separately over an identical hyperparameter space from an
605 initial starting point used during LFNet development. We also trained replicates for EDC using the
606 best hyperparameters for bioseq2seq and refer to the best EDC model as EDC-large and the EDC
607 with equivalent hyperparameters to bioseq2seq as EDC-small. We then produced four replicates for

608 each of bioseq2seq, EDC-large, and EDC-small. For further details on hyperparameter tuning and
609 model training see the Supplementary Details.

610 4.5 Mutation effect prediction

611 Estimating the effects of sequence mutations can provide insight into the importance that the model
612 assigns each input nucleotide. The gold standard for computationally scoring mutation effects, known
613 as *in silico mutagenesis* (ISM), requires comparing the model predictions for all single-nucleotide
614 variants with that of the original sequence (Zhou and Troyanskaya 2015). The computational expense
615 of this procedure – $3L$ model evaluations for a transcript of length L – motivates us to explore the
616 effectiveness of gradient-based approximations.

617 Below we refer to the network output function by S , and the output gradient with respect to its input
618 as $\nabla_x S(x)$. In general, S can be any scalar output, and here we use $S = l_{\langle PC \rangle} - l_{\langle NC \rangle}$, the difference
619 in logits, i.e unnormalized log probabilities, for the RNA classification tokens in the first decoding
620 position. We denote the two sequences being compared as $x, x' \in \mathbb{R}^{L \times V}$ for one-hot encodings of
621 categorical variables and V as the input vocabulary size.

622 **Taylor series approximation** The simplest ISM surrogate begins with a Taylor expansion of a dif-
623 ferentiable function F around a point of interest x' .

$$F(x') \approx F(x) + \nabla_x F(x)^\top (x' - x) + o(\|x' - x\|)$$

624 In this fashion, we can expand around S and discard all higher order terms for a first-order Taylor
625 approximation of the difference in S .

$$626 \quad \Delta S(x', x) = S(x') - S(x) \approx \nabla_x S(x)^\top (x' - x) \quad (1)$$

627 Since we confine our analysis to single-mutations, this simplifies to

$$628 \quad \Delta S(x\{i, j \rightarrow k\}, x) \approx \frac{\partial S(x)}{\partial x_{ik}} - \frac{\partial S(x)}{\partial x_{ij}} \quad (2)$$

629 where $x\{i, j \rightarrow k\}$ is the result of mutating RNA x at position i from nucleotide j to k . Thus, all $3L$

630 values are computable from $\nabla_x S(x)$ in just one forward/backward pass of the network.

631 **Mutation-Directed Integrated Gradients** The input gradient represents only an infinitesimal change
632 in the input-output behavior of the network, rather than the effect of a full character substitution as in
633 ISM. When a local approximation does not accurately describe the global function behavior, this is a
634 well known limitation called gradient saturation (Shrikumar et al. 2019). As a more sophisticated proxy
635 for ISM, we adapt a procedure called Integrated Gradients (IG), which was designed to reduce the
636 the effect of gradient saturation and satisfies several desirable axioms for importance metrics (Sun-
637 dararajan et al. 2017). IG uses a baseline input x' and computes an integral using input gradients for
638 a differentiable function F along the linear path between x' and x .

$$639 \quad IG(x, x')_{ib} = (x_{ib} - x'_{ib}) \int_{\alpha=0}^1 \frac{\partial F(x' + \alpha \times (x - x'))}{\partial x_{ib}} d\alpha \quad (3)$$

640 This equation relates to Taylor-approximation in that, given one hot encodings, $\sum_j IG(x, x')_{ij}$ is
641 equal to integrating the right hand side of Eq. 2 with x ranging over the interpolation path from Eq.
642 3. Based on this view, we propose a rough heuristic for estimating ISM using four evaluations of IG,
643 which we dub Mutation Directed Integrated Gradients (MDIG).

$$644 \quad \Delta S(x\{i, b \rightarrow k\}, x) \approx MDIG(x)_{ib} = IG(\beta \cdot \text{poly}(b) + (1 - \beta) \cdot x, x)_{ib} \quad \forall b \in \{A, C, G, T\} \quad (4)$$

645 Here, $\text{poly}(b)$ is a sequence of all nucleotide b of the same length as x , e.g. all guanines, and
646 $\beta \in (0, 1]$ is a hyperparameter that balances distance of the baseline from x , which is needed to
647 reduce gradient saturation, and distance from x' , a sequence largely unrelated to x . Note the order
648 of arguments, which re-frames the baseline as the destination rather than the source.

649 To compare MDIG against a traditional usage of Integrated Gradients, we constructed an alternate
650 baseline by placing the vector $[0.25, 0.25, 0.25, 0.25]$ – a uniform probability mass function over the four
651 bases – in all input positions. The mutation effect scores are then defined as $\Delta S(x\{i, b \rightarrow k\}, x) \approx$
652 $IG(\mathcal{U}(x), x)_{ik} - IG(\mathcal{U}(x), x)_{ib}$ where $\mathcal{U}(x)$ represents the uniform baseline. In this way, the uniform
653 IG approach requires just one evaluation of Eq 3 overall, while MDIG requires one evaluation per base

654 *b*.

655 4.6 Evaluation metrics for gradient attributions

We compared $L \times 3$ vectors (sequence length \times 3 possible mutations) of mutation effect predictions using the metrics

$$\text{Pearson}(x, y) = \frac{\text{cov}(\text{vec}(x), \text{vec}(y))}{\sigma(\text{vec}(x))\sigma(\text{vec}(y))}$$

$$\text{Median position-wise cosine similarity}(x, y) = \text{median}\left(\left[\cos(x_1, y_1), \cos(x_2, y_2), \dots, \cos(x_L, y_L)\right]\right)$$

Where $\text{cov}()$ is the covariance, $\sigma()$ is the standard deviation, $\text{vec}()$ is the vectorization operator, which flattens a matrix into a vector, $\cos(x, y) = \frac{x \cdot y}{\|x\| \|y\|}$, and x_i refers to a row vector of matrix x . The inter-replicate agreement is

$$\text{Inter-replicate agreement}(x) = \left| \binom{S}{2} \right|^{-1} \sum_{i, j \in \binom{S}{2}} \text{metric}(\text{mut}(x)_i, \text{mut}(x)_j)$$

where $\binom{S}{2}$ is the set of all possible subsets of cardinality 2 from the set of model replicates S , $\text{mut}(x)_i$ is the mutation effect prediction coming from replicate i for a given RNA x , and metric is one of Pearson r or median position-wise cosine similarity, as described above. The agreement with ISM is defined with intra-replicate comparisons.

$$\text{Agreement with ISM}(x) = |S|^{-1} \sum_{i \in S} \text{metric}(\text{ISM}(x)_i, \text{mut}(x)_i)$$

656 4.7 Motif discovery from mutation effect predictions

657 To uncover sequence elements salient to bioseq2seq predictions, we converted ISM scores into im-
658 portance scores for the endogenous characters. In particular, we set the importance score of an
659 endogenous base with respect to a given class as equal to the absolute value of ΔS for the strongest
660 mutation in the direction of the *counterfactual* class, following (Kelley et al. 2016) which used the
661 equivalent from regression models for visualizing importance. For example, an endogenous x_i within
662 an mRNA was defined as contributing towards a true positive classification of $\langle PC \rangle$ to the extent
663 that substituting any of the three alternate bases in position i produces a highly negative ΔS , which

664 pushes the prediction towards a false negative of $\langle NC \rangle$. We calculated importance using both classes
665 on all transcripts. For instance, we looked for strong local contributions towards a prediction of $\langle PC \rangle$
666 within annotated lncRNAs.

667 For a given importance setting, we then extracted a window of 10 nt upstream and 10 nt down-
668 stream around the position with the highest importance score for a total length of 21 nt. This process
669 was run separately for mRNA 5' and 3' UTRs and CDS sequences, and similarly for lncRNAs using
670 the longest ORF and its upstream and downstream regions. We used the STREME motif discov-
671 ery tool to efficiently identify sequence motifs occurring frequently in these regions of interest (Bailey
672 2021). STREME estimates p-values for motifs, and after collecting all discovered sequence logos, we
673 reported all that were significant at the 0.001 level after applying the Bonferroni correction for multiple
674 testing.

675 **5 Data Access**

676 Code for running our trained models and replicating the experiments and figures in this paper is
677 provided at <https://github.com/josephvalencia/bioseq2seq> and pretrained models and data
678 at <https://osf.io/xaeqg/>.

679 **6 Competing Interest Statement**

680 The authors have no competing interests to declare.

681 **7 Author Contributions**

682 J.D.V and D.A.H conceived of the project. J.D.V conceived of the LFNet architecture. J.D.V performed
683 all coding and analysis while supervised by D.A.H. J.D.V. wrote the manuscript, and D.A.H provided
684 edits.

685 **References**

- 686 Agarwal V and Kelley DR. 2022. The genetic and biochemical determinants of mRNA degradation
687 rates in mammals. *Genome Biology*. **23**: 245.
- 688 Anastassiou D. 2000. Frequency-domain analysis of biomolecular sequences. *Bioinformatics*. **16**:
689 1073–1081.
- 690 Anderson DM et al. 2015. A Micropeptide Encoded by a Putative Long Noncoding RNA Regulates
691 Muscle Performance. *Cell*. **160**: 595–606.
- 692 Avsec , Agarwal V, Visentin D, Ledsam JR, Grabska-Barwinska A, Taylor KR, Assael Y, Jumper J,
693 Kohli P, and Kelley DR. 2021. Effective gene expression prediction from sequence by integrating
694 long-range interactions. *Nature Methods*. **18**: 1196–1203.
- 695 Avsec , Weilert M, et al. 2021. Base-resolution models of transcription-factor binding reveal soft motif
696 syntax. *Nature Genetics*. **53**: 354–366.
- 697 Bailey TL. 2021. STREME: accurate and versatile sequence motif discovery. *Bioinformatics*. **37**:
698 2834–2840.
- 699 Bánfai B et al. 2012. Long noncoding RNAs are rarely translated in two human cell lines. *Genome*
700 *Research*. **22**: 1646–1657.
- 701 Barik S. 2020. The Uniqueness of Tryptophan in Biology: Properties, Metabolism, Interactions and
702 Localization in Proteins. *International Journal of Molecular Sciences*. **21**: 8776.
- 703 Calviello L, Mukherjee N, Wyler E, Zauber H, Hirsekorn A, Selbach M, Landthaler M, Obermayer B,
704 and Ohler U. 2016. Detecting actively translated open reading frames in ribosome profiling data.
705 *Nature Methods*. **13**: 165–170.
- 706 Camargo AP, Sourkov V, Pereira GAG, and Carazzolle MF. 2020. RNAsamba: neural network-based
707 assessment of the protein-coding potential of RNA sequences. *NAR Genomics and Bioinformat-*
708 *ics*. **2**:
- 709 Chaney JL, Steele A, Carmichael R, Rodriguez A, Specht AT, Ngo K, Li J, Emrich S, and Clark
710 PL. 2017. Widespread position-specific conservation of synonymous rare codons within coding
711 sequences. *PLOS Computational Biology*. **13**: e1005531.
- 712 Choi SW, Kim HW, and Nam JW. 2019. The small peptide world in long noncoding RNAs. *Briefings in*
713 *Bioinformatics*. **20**: 1853–1864.

- 714 Deng S, Chen Z, Ding G, and Liğ Y 2010. Prediction of protein coding regions by combining Fourier
715 and Wavelet Transform. In: *2010 3rd International Congress on Image and Signal Processing*.
716 Vol. 9, pp. 4113–4117.
- 717 Derrien T et al. 2012. The GENCODE v7 catalog of human long noncoding RNAs: Analysis of their
718 gene structure, evolution, and expression. *Genome Research*. **22**: 1775–1789.
- 719 Epps J, Ying H, and Huttley GA. 2011. Statistical methods for detecting periodic fragments in DNA
720 sequence data. *Biology Direct*. **6**: 21.
- 721 Fickett JW. 1982. Recognition of protein coding regions in DNA sequences. *Nucleic Acids Research*.
722 **10**: 5303–5318.
- 723 Gebauer F and Hentze MW. 2004. Molecular mechanisms of translational control. *Nature Reviews*
724 *Molecular Cell Biology*. **5**: 827–835.
- 725 Guibas J, Mardani M, Li Z, Tao A, Anandkumar A, and Catanzaro B. 2021. Adaptive Fourier Neural
726 Operators: Efficient Token Mixers for Transformers. *arXiv:2111.13587 [cs]*.
- 727 Guo H, Ingolia NT, Weissman JS, and Bartel DP. 2010. Mammalian microRNAs predominantly act to
728 decrease target mRNA levels. *Nature*. **466**: 835–840.
- 729 Gupta S, Stamatoyannopoulos JA, Bailey TL, and Noble WS. 2007. Quantifying similarity between
730 motifs. *Genome Biology*. **8**: R24.
- 731 Guttman M, Russell P, Ingolia NT, Weissman JS, and Lander ES. 2013. Ribosome profiling provides
732 evidence that large non-coding RNAs do not encode proteins. *Cell*. **154**: 240–251.
- 733 Gyawali PK, Liu X, Zou J, and He Z 2022. Ensembling improves stability and power of feature selection
734 for deep learning models. en. In: *Proceedings of the 17th Machine Learning in Computational*
735 *Biology meeting*. ISSN: 2640-3498. PMLR, pp. 33–45.
- 736 Hartford CCR and Lal A. 2020. When Long Noncoding Becomes Protein Coding. *Molecular and Cel-*
737 *lular Biology*. **40**: e00528–19, /mcb/40/6/MCB.00528–19.atom.
- 738 Hassani Saadi H, Sameni R, and Zollanvari A. 2017. Interpretive time-frequency analysis of genomic
739 sequences. *BMC Bioinformatics*. **18**: 154.
- 740 Hill ST, Kuintzle R, Teegarden A, Merrill E, Danaee P, and Hendrix DA. 2018. A deep recurrent neu-
741 ral network discovers complex biological rules to decipher RNA protein-coding potential. *Nucleic*
742 *Acids Research*. **46**: 8105–8113.

- 743 Housman G and Ulitsky I. 2016. Methods for distinguishing between protein-coding and long noncod-
744 ing RNAs and the elusive biological purpose of translation of long noncoding RNAs. *Biochimica et*
745 *Biophysica Acta (BBA) - Gene Regulatory Mechanisms*. **1859**: 31–40.
- 746 Ingolia NT, Brar GA, Stern-Ginossar N, Harris MS, Talhouarne GJ, Jackson SE, Wills MR, and Weiss-
747 man JS. 2014. Ribosome Profiling Reveals Pervasive Translation Outside of Annotated Protein-
748 Coding Genes. *Cell Reports*. **8**: 1365–1379.
- 749 Ingolia NT, Lareau LF, and Weissman JS. 2011. Ribosome Profiling of Mouse Embryonic Stem Cells
750 Reveals the Complexity of Mammalian Proteomes. *Cell*. **147**: 789–802.
- 751 Iyer MK et al. 2015. The landscape of long noncoding RNAs in the human transcriptome. *Nature*
752 *Genetics*. **47**: 199–208.
- 753 Jackson R et al. 2018. The translation of non-canonical open reading frames controls mucosal immu-
754 nity. *Nature*. **564**: 434–438.
- 755 Ji Z, Song R, Regev A, and Struhl K. 2015. Many lncRNAs, 5UTRs, and pseudogenes are translated
756 and some are likely to express functional proteins. *eLife*. **4**: e08890.
- 757 Johnstone TG, Bazzini AA, and Giraldez AJ. 2016. Upstream ORFs are prevalent translational re-
758 pressors in vertebrates. *The EMBO Journal*. **35**: 706–723.
- 759 Jumper J et al. 2021. Highly accurate protein structure prediction with AlphaFold. *Nature*. **596**: 583–
760 589.
- 761 Kelley DR, Snoek J, and Rinn JL. 2016. Basset: learning the regulatory code of the accessible genome
762 with deep convolutional neural networks. *Genome Research*. **26**: 990–999.
- 763 Klein G, Kim Y, Deng Y, Nguyen V, Senellart J, and Rush AM. 2018. OpenNMT: Neural Machine
764 Translation Toolkit. *arXiv:1805.11462 [cs]*.
- 765 Kong L, Zhang Y, Ye ZQ, Liu XQ, Zhao SQ, Wei L, and Gao G. 2007. CPC: assess the protein-
766 coding potential of transcripts using sequence features and support vector machine. *Nucleic Acids*
767 *Research*. **35**: W345–W349.
- 768 Koo PK, Majdandzic A, Ploenzke M, Anand P, and Paul SB. 2021. Global importance analysis: An
769 interpretability method to quantify importance of genomic features in deep neural networks. *PLOS*
770 *Computational Biology*. **17**: e1008925.

- 771 Kozak M. 1987. An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nu-*
772 *cleic Acids Research*. **15**: 8125–8148.
- 773 Kozak M. 2002. Pushing the limits of the scanning mechanism for initiation of translation. *Gene*. **299**:
774 1–34.
- 775 Lee-Thorp J, Ainslie J, Eckstein I, and Ontanon S. 2021. FNet: Mixing Tokens with Fourier Transforms.
776 *arXiv:2105.03824 [cs]*.
- 777 Li A, Zhang J, and Zhou Z. 2014. PLEK: a tool for predicting long non-coding RNAs and messenger
778 RNAs based on an improved k-mer scheme. *BMC Bioinformatics*. **15**: 311.
- 779 Li JJ, Chew GL, and Biggin MD. 2019. Quantitative principles of cis-translational control by general
780 mRNA sequence features in eukaryotes. *Genome Biology*. **20**: 162.
- 781 Li W and Godzik A. 2006. Cd-hit: a fast program for clustering and comparing large sets of protein or
782 nucleotide sequences. *Bioinformatics*. **22**: 1658–1659.
- 783 Liaw R, Liang E, Nishihara R, Moritz P, Gonzalez JE, and Stoica I 2018. Tune: A Research Platform
784 for Distributed Model Selection and Training. *arXiv:1807.05118 [cs, stat]*.
- 785 Liu T, Wu J, Wu Y, Hu W, Fang Z, Wang Z, Jiang C, and Li S. 2022. LncPep: A Resource of Transla-
786 tional Evidences for lncRNAs. *Frontiers in Cell and Developmental Biology*. **10**:
- 787 Makarewich CA and Olson EN. 2017. Mining for Micropeptides. *Trends in cell biology*. **27**: 685–696.
- 788 Novakovsky G, Dexter N, Libbrecht MW, Wasserman WW, and Mostafavi S. 2022. Obtaining genetics
789 insights from deep learning via explainable artificial intelligence. *Nature Reviews Genetics*. 1–13.
- 790 O'Leary NA et al. 2016. Reference sequence (RefSeq) database at NCBI: current status, taxonomic
791 expansion, and functional annotation. *Nucleic Acids Research*. **44**: D733–D745.
- 792 Othoum G, Coonrod E, Zhao S, Dang HX, and Maher CA. 2020. Pan-cancer proteogenomic analysis
793 reveals long and circular noncoding RNAs encoding peptides. *NAR Cancer*. **2**: zcaa015.
- 794 Patraquim P, Magny EG, Pueyo JI, Platero AI, and Couso JP. 2022. Translation and natural selection
795 of micropeptides from long non-canonical RNAs. *Nature Communications*. **13**: 6515.
- 796 Ransohoff JD, Wei Y, and Khavari PA. 2018. The functions and unique features of long intergenic
797 non-coding RNA. *Nature Reviews Molecular Cell Biology*. **19**: 143–157.

- 798 Rao Y, Zhao W, Zhu Z, Lu J, and Zhou J 2021. Global Filter Networks for Image Classification. In:
799 *Advances in Neural Information Processing Systems*. Vol. 34. Curran Associates, Inc., pp. 980–
800 993.
- 801 Ray D et al. 2013. A compendium of RNA-binding motifs for decoding gene regulation. *Nature*. **499**:
802 172–177.
- 803 Reis M dos, Wernisch L, and Savva R. 2003. Unexpected correlations between gene expression and
804 codon usage bias from microarray data for the whole Escherichia coli K12 genome. *Nucleic Acids*
805 *Research*. **31**: 6976–6985.
- 806 Rice P, Longden I, and Bleasby A. 2000. EMBOSS: The European Molecular Biology Open Software
807 Suite. *Trends in Genetics*. **16**: 276–277.
- 808 Sallam T, Sandhu J, and Tontonoz P. 2018. Long Noncoding RNA Discovery in Cardiovascular Dis-
809 ease. *Circulation Research*. **122**: 155–166.
- 810 Schreiber J, Nair S, Balsubramani A, and Kundaje A. 2022. Accelerating in silico saturation mutagen-
811 esis using compressed sensing. *Bioinformatics*. **38**: 3557–3564.
- 812 Sharp PM and Li WH. 1987. The codon Adaptation Index—a measure of directional synonymous codon
813 usage bias, and its potential applications. *Nucleic Acids Research*. **15**: 1281–1295.
- 814 Shrikumar A, Greenside P, and Kundaje A. 2019. Learning Important Features Through Propagating
815 Activation Differences. *arXiv:1704.02685 [cs]*.
- 816 Statello L, Guo CJ, Chen LL, and Huarte M. 2021. Gene regulation by long non-coding RNAs and its
817 biological functions. *Nature Reviews Molecular Cell Biology*. **22**: 96–118.
- 818 Subramanian K, Waugh N, Shanks C, and Hendrix DA 2021. Position-dependent Codon Usage Bias
819 in the Human Transcriptome. en. Pages: 2021.08.11.456006 Section: New Results.
- 820 Sundararajan M, Taly A, and Yan Q. 2017. Axiomatic Attribution for Deep Networks. *arXiv:1703.01365*
821 *[cs]*.
- 822 Szostak E and Gebauer F. 2013. Translational control by 3-UTR-binding proteins. *Briefings in Func-*
823 *tional Genomics*. **12**: 58–65.
- 824 Tareen A, Kooshkbaghi M, Posfai A, Ireland WT, McCandlish DM, and Kinney JB. 2022. MAVEN-NN:
825 learning genotype-phenotype maps from multiplex assays of variant effect. *Genome Biology*. **23**:
826 98.

- 827 Tiwari S, Ramachandran S, Bhattacharya A, Bhattacharya S, and Ramaswamy R. 1997. Prediction
828 of probable genes by Fourier analysis of genomic sequences. *Bioinformatics*. **13**: 263–270.
- 829 Tseng A, Shrikumar A, and Kundaje A 2020. Fourier-transform-based attribution priors improve the
830 interpretability and stability of deep learning models for genomics. In: *Advances in Neural Infor-*
831 *mation Processing Systems*. Vol. 33. Curran Associates, Inc., pp. 1913–1923.
- 832 Tuller T, Carmi A, Vestsigian K, Navon S, Dorfan Y, Zaborske J, Pan T, Dahan O, Furman I, and
833 Pilpel Y. 2010. An Evolutionarily Conserved Mechanism for Controlling the Efficiency of Protein
834 Translation. *Cell*. **141**: 344–354.
- 835 Vaswani A, Shazeer N, Parmar N, Uszkoreit J, Jones L, Gomez AN, Kaiser , and Polosukhin I 2017.
836 Attention is All you Need. In: *Advances in Neural Information Processing Systems*. Vol. 30. Curran
837 Associates, Inc.
- 838 Verheggen K, Volders PJ, Mestdagh P, Menschaert G, Van Damme P, Gevaert K, Martens L, and
839 Vandesompele J. 2017. Noncoding after All: Biases in Proteomics Data Do Not Explain Observed
840 Absence of lncRNA Translation Products. *Journal of Proteome Research*. **16**: 2508–2515.
- 841 Verma M et al. 2019. A short translational ramp determines the efficiency of protein synthesis. *Nature*
842 *Communications*. **10**: 5774.
- 843 Wang L, Park HJ, Dasari S, Wang S, Kocher JP, and Li W. 2013. CPAT: Coding-Potential Assessment
844 Tool using an alignment-free logistic regression model. *Nucleic Acids Research*. **41**: e74.
- 845 Wright MA, Kharchenko P, Church GM, and Segrè D. 2007. Chromosomal periodicity of evolutionarily
846 conserved gene pairs. *Proceedings of the National Academy of Sciences*. **104**: 10559–10564.
- 847 Wu P et al. 2020. Emerging role of tumor-related functional peptides encoded by lncRNA and circRNA.
848 *Molecular Cancer*. **19**: 22.
- 849 Wucher V et al. 2017. FEELnc: a tool for long non-coding RNA annotation and its application to the
850 dog transcriptome. *Nucleic Acids Research*. **45**: e57.
- 851 Xu Z, Hu L, Shi B, Geng S, Xu L, Wang D, and Lu ZJ. 2018. Ribosome elongating footprints denoised
852 by wavelet transform comprehensively characterize dynamic cellular translation events. *Nucleic*
853 *Acids Research*. **46**: e109.
- 854 Yan Y et al. 2021. The cardiac translational landscape reveals that micropeptides are new players
855 involved in cardiomyocyte hypertrophy. *Molecular Therapy*. **29**: 2253–2267.

- 856 Yang Y, Huang L, and Ma M 2018. Breaking the Beam Search Curse: A Study of (Re-)Scoring Meth-
857 ods and Stopping Criteria for Neural Machine Translation. en. In: *Proceedings of the 2018 Confer-*
858 *ence on Empirical Methods in Natural Language Processing*. Brussels, Belgium: Association for
859 Computational Linguistics, pp. 3054–3059.
- 860 Zeng T and Li YI. 2022. Predicting RNA splicing from DNA sequence using Pangolin. *Genome Biol-*
861 *ogy*. **23**: 103.
- 862 Zhou J and Troyanskaya OG. 2015. Predicting effects of noncoding variants with deep learningbased
863 sequence model. *Nature Methods*. **12**: 931–934.
- 864 Zhu S, Wang JZ, Chen D, He YT, Meng N, Chen M, Lu RX, Chen XH, Zhang XL, and Yan GR. 2020. An
865 oncopeptide regulates m6A recognition by the m6A reader IGF2BP1 and tumorigenesis. *Nature*
866 *Communications*. **11**: 1685.