1  **Analysis of RNA translation with a deep learning architecture**

2  **provides new insight into translation control**

3

4

5  Xiaojuan Fan[1,2,*,#], Tiangen Chang[4,*], Chuyun Chen[1,3], Markus Hafner[2], Zefeng Wang[1,3,#]

6

7  [1] Bio-med Big Data Center, CAS Key Laboratory of Computational Biology, CAS Center for

8  Excellence in Molecular Cell Science, Shanghai Institute of Nutrition and Health,

9  [2] RNA Molecular Biology Laboratory, National Institute of Arthritis and Musculoskeletal and Skin

10  Disease, Bethesda, MD, USA.

11  [3] University of Chinese Academy of Sciences, Chinese Academy of Sciences, Shanghai 200031,

12  China

13  [4] Laboratory of Cancer Data Science, National Cancer Institute, Bethesda, MD, USA.

14

15

16

17  * These authors contributed equally to this work

18  [#] Corresponding to: wangzefeng@picb.an.cn, fanx3@nih.gov

19

20

21

22

23

24 **Abstract**

25 Accurate annotation of coding regions in RNAs is essential for understanding gene

26 translation. We developed a deep neural network to directly predict and analyze

27 translation initiation and termination sites from RNA sequences. Trained with human

28 transcripts, our model learned hidden rules of translation control and achieved a near

29 perfect prediction of canonical translation sites across entire human transcriptome.

30 Surprisingly, this model revealed a new role of codon usage in regulating translation

31 termination, which was experimentally validated. We also identified thousands of new

32 open reading frames in mRNAs or lncRNAs, some of which were confirmed

33 experimentally. The model trained with human mRNAs achieved high prediction

34 accuracy of canonical translation sites in all eukaryotes and good prediction in

35 polycistronic transcripts from prokaryotes or RNA viruses, suggesting a high degree

36 of conservation in translation control. Collectively, we present a general and efficient

37 deep learning model for RNA translation, generating new insights into the complexity

38 of translation regulation.

39

## Introduction

Computational analysis of protein-coding open reading frames (ORFs) in RNAs have traditionally relied on empirical rules such as sequence conservation pattern (1), the longest ORFs in one transcript (2) and initiation from the first AUG (3). However, these rules are over-simplified without considering the complex regulation by other factors, such as RNA structure and various regulatory *cis*-elements. As a result, previous efforts in *ab initio* prediction of translation sites achieved relatively low accuracy (4). In addition, RNA translation has recently been found in previously annotated non-coding RNAs (e.g., lncRNAs and circRNAs) (5-7), and can be initiated from alternative sites or non-canonical codons (8,9), resulting in non-annotated translation products that play critical and diverse cellular roles (6,8). Considering the complexities of translation control, the prediction of translation sites with high accuracy remains a challenging task, especially for poorly annotated genomes or in the context of complex regulatory networks. Rapid advances in artificial intelligence (AI) offer new hope for dissecting complex biological systems, such as alternative splicing (10) and protein folding (11), by uncovering hidden associations and unknown rules (12-15).

The complete synthesis of a protein requires accurate initiation and termination of translation by ribosomes at specific positions. In eukaryotes, translation initiation is generally mediated by the 5′-cap, which recruits various translation initiation factors and small ribosomal subunit to form preinitiation complex (16). This complex subsequently scans through the 5′-untranslated region (5′-UTR) to recognize the

62  optimal start codon, which is determined by both codon identity and its immediate

63  nucleotide context known as Kozak sequence GCCRCC<u>AUG</u>G (R represents purine)

64  (3). However, in certain cases, the ribosome may skip some weak start codons,

65  resulting in leaky scanning (3).   In addition, a large fraction of genes contain several

66  alternative translation initiation sites (TISs), but it is unclear how the ribosomes select

67  one of these sites to initiate translation. Evaluating the strength of each TIS within one

68  transcript and the general rules of selecting the authentic ORFs from possible decoys

69  remain challenging.   Therefore the systematic identification of translation sites

70  (translation initiation and termination sites, or TIS/TTS) and the quantitative

71  evaluation of their strength become a prerequisite for exploring the complexity of

72  translation regulation.

73      New experimental techniques, including GTI-seq (global translation initiation

74  sequencing) (17) and QTI-seq (quantitative translation initiation sequencing) (18),

75  have been developed to identify TISs by directly capturing the stalled initiating

76  ribosomes at single nucleotide resolution across the transcriptome.   However, these

77  methods were typically conducted in specific cell types, which limits the

78  transcriptome-wide TIS identification. Several computational methods have also been

79  developed to predict translation initiation sites from mRNA sequences. For example, a

80  physicochemical-property based predictor using pseudo-trinucleotide composition

81  was constructed to identify TIS in human genes (19), and a variety of machine

82  learning methods have also been developed to predict the TIS by analyzing a short

83  sequence fragment (e.g. 200nt) around the start codon of mRNA (20-23). However,

84    such methods are biased because they are limited to the immediate sequence context

85    surrounding the TISs. In other words, these approaches do not consider the entire non-

86    coding regions of mRNAs (i.e. 5′-UTR and 3′-UTR) that play an important regulatory

87    role in selecting the start codon (24,25), which in turn affects the accuracy of TIS

88    predictions.

89        Here, we developed a deep learning method based on a multilevel dilated

90    convolution network, named TranslationAI, to independently predict TISs and TTSs

91    from full-length mRNA sequences alone. Leveraging the inherent structure of the

92    genetic code, TranslationAI achieves high accuracy, surpassing a >99% Precision-

93    Recall Area Under the Curve (PR-AUC) in predicting known TISs and TTSs in

94    human transcriptome.   Our model uncovered regulatory sequences in the UTRs that

95    determine the strength of translation start and stop codons, and identified thousands of

96    new ORFs in mRNAs or annotated lncRNAs. In addition, this model can be extended

97    to other eukaryotes, prokaryote and several human viruses, suggesting a strong

98    conservation in defining coding sequences across different domains of life. To

99    facilitate further applications by other researchers, we developed a web tool

100    (https://www.biosino.org/TranslationAI/) that is accessible to end users in predicting

101    translation sites from RNA sequences. The application of TranslationAI extends

102    beyond *de novo* prediction of TIS/TTS in any given transcripts, as it can also be used

103    to evaluate the translation potential of a given RNA sequence.

104

**Results**

**Modeling translation from full length mRNA with deep learning**

We constructed a deep residual convolutional neural network, TranslationAI, using the full-length mRNA sequences as input to independently predict the TIS and TTS (Fig. 1A, see methods). Previous approaches only considered short sequence stretches around a potential start codon (20-22), or reported all potential ORFs in six frames with specified start codon (ATG), stop codon (TAG/TAA/TGA), and length cutoffs (ORFfinder in NCBI). In contrast, our deep learning model evaluates the potential of each position in a given mRNA to function as a TIS or TTS without using any prior knowledge of translation (e.g., the triplet or identity of genetic code). Such *ab initio* prediction enables the neural network to learn the hidden rules of translation from sequence context alone.

The deep residual neural network was constructed using the one-hot encoding of mRNA sequences as input (Fig. 1A). We built a 32-layer dilated convolution neural network architecture that produces an output matrix of the probability for the given position being a TIS, TTS or neither (Fig. S1, see methods). To obtain a comprehensive and well-curated dataset, we utilized the RefSeq gene annotation and extracted 47,098 protein-coding transcripts (47,098 TIS-TTS pairs). We trained the model using the transcripts on chromosomes 2, 4, 6, 8, 10-22, X, and Y, and tested the model with the transcripts from chromosomes 1, 3, 5, 7, and 9. The stringent top-k accuracy was used for TIS/TTS prediction. Specifically, the score cutoff was set at the point where the number of predicted sites is equal to the real number (i.e., the

127     false positive rate equal to false negative rate, the score cutoff at this point is ~0.5),

128     which ensured that only the high confidence predicted sites were considered for

129     downstream analysis.

130        To explore the effect of different input windows on test accuracy, we considered

131     four input windows that include 60, 200, 600, and 2,000 (2k) nucleotides on both

132     sides of a given position.   Interestingly, the prediction accuracy improved

133     dramatically with the increasing input window size, reaching >99% PR-AUC with a

134     2k nt input window on test dataset (Fig. 1B, Table S1). Considering that >70% of

135     transcripts are longer than 2k nt (Fig. S2A), such high accuracy of TranslationAI

136     prediction is particularly remarkable, suggesting that the identity of canonical

137     TIS/TTS could be influenced by the sequences at thousands of nucleotides away,

138     potentially through long-range interactions such as RNA structure or RNA binding

139     proteins. Our findings further suggest that TranslationAI has learned to capture these

140     interactions using sequence information alone.

141        To assess the impact of sequence context on the accuracy of TIS/TTS prediction,

142     we compared the performance of the models trained on 200 and 2k nt sequence

143     context. The 2k model achieved high prediction precision (>0.95) in canonical

144     TIS/TTS prediction across a wide range of transcript lengths, however the precision

145     of the 200 nt model decreases dramatically with increase of transcript length (Fig.

146     1C).   This observation is in agreement with that earlier report that distal regions far

147     away from the canonical TIS/TTS can affect the translation decision (24).   With the

148     2k model, the score distribution showed a clear separation between the positive and

149 negative predictions (Fig. 1D), suggesting that the model is robust under different

150 cutoff selections.    To further examine the false discovery rate (FDR), we used the 2k

151 model to score potential TISs/TTSs in the shuffled mRNA sequences.    At the

152 probability score cutoff of 0.5 for both TIS and TTS, less than 5% of the ORFs in

153 shuffled dataset passed the cutoff, yielding an empirical FDR < 0.05 (Fig. 1E).

154

155 **Key features learned by the AI model for accurate prediction**

156      Encouraged by the near-perfect performance of our deep learning network in

157 predicting canonical translation sites, we sought to examine the features learned by

158 this "black box" model.    We performed systematic *in silico* perturbations on different

159 regions of mRNAs to measure their effects on reducing the probability scores of the

160 authentic TISs/TTSs.    Such perturbations can reflect how the AI model make the

161 accurate translation prediction in a computational sense, however a different

162 molecular mechanism may be used in cells to define translation initiation and

163 termination sites as judged by experiments.    Nevertheless, it would be informative to

164 compare how a biological process like mRNA translation is perceived differently by

165 the AI model and experimental tests.

166      We found that replacing annotated TISs or TTSs with random codons

167 dramatically reduced the predicted scores of cognate sites (Fig. 1F, Table S2),

168 suggesting this model had successfully learned the identity of start and stop codons

169 from the sequence alone. Consistently, introducing a known stop codon into the

170 coding sequences (CDS) led to a drastic reduction in predicted scores of both TISs

171 and TTSs. In addition, shuffling the 5′-UTRs significantly affected the corresponding

172 TIS prediction but not TTS, whereas shuffling the 3′-UTR had little effect on

173 TIS/TTS prediction (Fig. 1G, Table S2). Notably, when both 5′-UTR and 3′-UTR

174 were simultaneously shuffled, the TIS prediction score was lower compared to the

175 cases when only the 5′-UTR was perturbed (Fig. 1G, Table S2), underscoring a

176 potential synergistic role of 5′-UTR and 3′-UTR in TIS selection. Curiously, the

177 deletions of 5′- or 3′-UTRs had less effect on TISs/TTSs prediction by this model,

178 with near perfect prediction for transcripts with only 50 nt left in the UTRs (Fig. 1G,

179 Table S2), however the underlying reason is unclear. Finally, random deletions of one

180 or two nucleotides within CDS dramatically reduced the prediction scores, whereas

181 random deletion of three nucleotides in CDS had little effect (Fig. 1H, Table S2),

182 suggesting that the model was capable of learning the triplet rule of the genetic code

183 by itself.

184     To further examine the impact of codon usage on our predictions, we introduced a

185 series of perturbations to the genetic codons of CDS. First, maintaining the same

186 amino acid compositions, the in-frame triplet shuffling or random synonymous

187 substitutions had little effect on the prediction (Fig. 1I, Table S2). We further replaced

188 codons with synonymous counterparts possessing either higher or lower Codon

189 Adaptation Index (CAI). Notably, higher CAI substitutions increased prediction

190 scores for both TIS and TTS, while lower CAI substitutions decreased these scores.

191 These findings indicate that the AI model utilizes codon choice as a crucial factor in

192 making predictions. Furthermore, we observed that nonsynonymous substitutions

193    caused a small but statistically significant reduction in prediction scores. Collectively,

194    these results suggest that the TranslationAI network has learned the intricate rules of

195    RNA translation from a subset of the transcriptome without any prior knowledge.

196         To further investigate how other sites in the full-length mRNA contributed to the

197    TranslationAI prediction, we calculated the feature importance of different positions

198    of specific transcript in predicting TIS and TTS using an occlusion sensitivity

199    analysis.    Our results reveal that, while the sequences at the site of TIS and TTS

200    contribute strongly to the prediction, the other positions at different regions also have

201    various contributions (Fig. S2B).    Remarkably, the importance of nucleotide at

202    different positions relative to the TIS/TTS is dependent on the specific sequence

203    contexts (i.e., the positional feature importance is different for each specific mRNA).

204    For instance, in the case of a short mRNA GPX6, we observed a great impact of TIS

205    identify on TTS prediction. In another case of a long mRNA PPP6R2, we found that a

206    position in the 5′−UTR shows high contribution to the TIS prediction. More

207    specifically, when the original sequence "TAA" is masked into "NAA" or mutated,

208    the score of the original TIS was greatly decreased (Fig. S2B).

209         In the RefSeq dataset, ~100% of the annotated ORFs follow the triplet codon rule

210    (length=3N), 98% are the longest ORFs in the transcripts, and ~40% use the first

211    AUG triplet as their start codon, raising the possibility that only using these simple

212    rules may be sufficient to make prediction.    To determine the extent to which

213    TranslationAI relies on these simple rules, we generated *in silico* permutations in the

214    test dataset and measured their effects on prediction precision. The introduction of

215  frameshift permutations significantly reduced the prediction precision (Fig. S2C),

216  however >50% of transcripts were still predicted to use original TIS/TTS despite

217  breaking the 3N rule (i.e. ORF length ≠ 3N), suggesting that the rule of triplet codon

218  is important but not mandatory for TIS/TTS prediction by TranslationAI.    In

219  addition, when ORF length was reduced by introducing premature stop codons or

220  frameshifts, TranslationAI predicted the longest ORF in ~60% cases (Fig. S2D), again

221  suggesting that such a simple rule was partially followed by this model. Finally, the

222  preference for the first AUG was also affected by several permutations throughout the

223  transcript, including changes in TIS and the introduction of early stop or frameshift

224  (Fig. S2E).    Interestingly, shortening the length of the 5′-UTR increased the ratio of

225  first AUG usage, whereas shuffling of the 5′-UTR reduced it, suggesting that the 5′-

226  UTR plays a critical role in TIS selection (24,25).    Collectively, these results indicate

227  that TranslationAI has learned additional features beyond simple traditional rules (i.e.,

228  the longest ORF, triplet code, and preference of the first AUG) for accurate prediction,

229  implying that the adjacent sequences near TIS/TTS also play crucial roles in selecting

230  translation start and stop sites.

231      We found no correlation (Spearman's correlation coefficient close to zero)

232  between the predicted score of TIS/TTS and the number of transcript isoforms

233  containing the same TIS/TTS sequences (Fig. S2F), suggesting that the scores of this

234  model were not biased by the training data. Additionally, the TIS and TTS scores in

235  the same ORF had a weak but statistically significant positive correlation (Spearman's

236  correlation coefficient r = 0.2, p < $1*10^{-16}$) (Fig. S2G), although they were predicted

237    independently.    Interestingly, the TIS/TTS scores by TranslationAI showed a positive

238    correlation with the translation efficiency estimated by the ratio of read density from

239    Ribosome profiling (Ribo-seq) normalized by mRNA abundance (8) (Fig. S2H and

240    S2I).    This consistency aligns with previous observations (Fig. 1I), where

241    synonymous substitutions in CDS with higher CAI exhibits higher prediction scores,

242    while substitutions with lower CAI shows lower prediction scores. This correlation

243    implies that that these predicted scores partially reflect their endogenous translational

244    activities.

245

246    **TranslationAI scores reveal new rule of codon selection for translation termination**

247        We next sought to examine the hidden information of TranslationAI scores by

248    comparing the key features between the strong and weak annotated TIS/TTS sites (the

249    top *vs.* the bottom 5th percentile). Interestingly, the sequences surrounding strong

250    TISs/TTSs are more conserved than those around weak TISs/TTSs (Fig. 2A),

251    implying additional selective pressures for translation.    Additionally, we examined

252    the ribosome occupancy on the mRNAs with different TIS/TTS scores using Ribo-seq

253    data from iPSC and iPSC-derived cardiomyocytes (8,26).    Our analysis revealed that

254    the transcripts with strong TISs or TTSs exhibit a higher ribosome peak at the start or

255    stop positions (Fig. 2B-2C, Fig. S3A-S3B), which is consistent with their higher

256    activity in mediating translation initiation or termination. These results suggest that

257    the predicted scores of TISs/TTSs reflect their activities in regulating translational

258    initiation/termination.    In addition, the transcripts with weak TIS/TTS generally have

259    a longer 5′-UTR than those with strong TIS/TTS (Fig. S3C), but such length

260    difference was observed in the 3′-UTRs only for TTS with different strengths.

261    However the implication behind such length difference is unclear.

262        We further examined the enriched sequence motifs around the strong *vs.* weak

263    TISs/TTSs (from -30 nt to +33 nt) (Fig. 2D).    As expected, we found an enriched

264    motif RCCATGGC (R = A or G, with start codon underlined) at the strong TISs,

265    resembling the Kozak sequence that typically enhances translation initiation (27,28).

266    We also found a novel motif STGAG at the strong TTSs but not the weak TTSs (S =

267    C or G, with stop codon underlined, Fig. 2D).    Surprisingly, the fragments

268    surrounding the strong TTS are generally biased towards CG-rich sequences,

269    implying that the C/G-rich region may mediate more efficient translation termination.

270    Moreover, the C/G enrichment showed a triplet periodicity before the strong TTSs

271    (Fig. 2D), with the third position of a codon being most C/G biased.    Such codon

272    bias probably reflects the evolutionary selection under the constrain of protein

273    sequences, suggesting an unappreciated role of codon usage in translation termination.

274        To confirm this codon bias, we analyzed the 30-nt region immediately upstream

275    of the stop codon of human mRNA (equivalent of 10 codons).    Consistent with the

276    motif enrichment, the synonymous codons with C/G at the third position were

277    preferably used at the upstream of strong TTSs (Fig. 2E, Fig.S3D).    Remarkably,

278    such codon bias was statistically significant for all 18 amino acids with multiple

279    codons (Fig. 2E), particularly for the amino acids with 4-6 codons (e.g., Thr, Val and

280    Ser etc., Fig. S3D). Together, these results reveal an intriguing role of codon selection

281    in translation termination, and indicate that the TranslationAI has learned this codon

282    bias to make a TTS prediction.

283        To experimentally validate the unexpected role of codon bias in translation

284    termination, we designed a series of translation readthrough reporters with different

285    TTSs (Fig. 2F).    This reporter contains two independent transcription units, one

286    encoding the firefly luciferase (Fluc) as a transfection control, and the other encoding

287    a fusion protein of GFP and Nano-luciferase (Nluc) separated by a set of short

288    variable sequences (45 nt) containing candidate TTSs of different strengths (Fig. 2F).

289    As a control, we inserted the sequence around the terminal codon of the human

290    vitamin D receptor (VDR) gene, which is known to undergo translation readthrough

291    (29).    We further introduced different synonymous mutations before the termination

292    codon of VDR to either strengthen or weaken the TTS, and assayed the levels of

293    mRNA and translation products from these reporters (Fig. 2F and S3E).    The results

294    showed that introducing the C/G-rich synonymous mutations in VDR (i.e., generating

295    a strong TTS) eliminated the product from translation readthrough (indicated by

296    arrows, Fig. 2F).    Conversely, a weak TTS mutation with A/U-rich synonymous

297    codons before the stop codon increased the translation readthrough.    This experiment

298    directly supported our finding on the role of codon selection in translation

299    termination.    As an independent validation, we re-analyzed the Ribo-seq data from

300    two distinct cell lines to examine potential translation read through (8).

301        As expected, the ribosomal density in the 3′-UTRs following weak TTSs was

302    significantly higher than that following strong TTSs (Fig. 2G), suggesting a

303    translation "leakage" following the weak termination sites. To validate this

304    observation, we conducted an analysis of the ribosome profiling data, in which the

305    cells were subjected to a high-salt wash to release vacant ribosomes lacking a nascent

306    polypeptide(30). The results demonstrated a distinct pattern: the high-salt wash

307    effectively released the ribosome signal at the last P-site of strong TTS, while no such

308    release was observed at weak TTS (Fig. 2H). This specific response is consistent with

309    above findings that ribosomes associated with weak TTS may undergo translation

310    readthrough, allowing continued translation beyond the canonical termination site,

311    whereas ribosomes at strong TTS are more likely to be released, terminating

312    translation at the expected site.

313

**Alternative TIS predicted by TranslationAI**

315    Certain mRNAs may use alternative TIS to produce additional protein isoforms

316    with extended N terminals (31), therefore we further examined if the TransaltionAI can

317    predict alternative TISs that may compete with the annotated sites (see methods).

318    Based on the distribution of TranslationAI scores (Fig. 1D), we used the score cutoff at

319    0.1 to increase the sensitivity of the alternative TISs prediction and identified 5336

320    alternative TISs in total. Interestingly, the numbers of alternative TIS identified by

321    TranslationAI is negatively correlated with the strength of the annotated TISs in the

322    same transcript (Fig. S3F), suggesting that alternative translation initiation may happen

323    in the mRNA without dominant TIS. Specifically, ~55% of mRNAs with weak

324    annotated TISs contain an alternative TIS (1340 out of 2416 mRNAs), whereas the

325    alternative TISs were found in less than 3% of mRNAs with strong annotated TISs (68

326    out of 2352 mRNAs). In addition, the score difference between the annotated TISs and

327    the alternative TIS was much smaller in the mRNAs with weak TISs (Fig. S3G, left),

328    whereas the mRNAs with strong annotated TISs generally lack alternative TISs with a

329    competitive score (Fig. S3G, right).   A consistent result was also found when we only

330    considered the alternative TISs in the same reading frame with the annotated TISs (Fig.

331    S3G).   These results collectively suggested that the mRNAs with weak annotated TISs

332    may also be translated from an alternative start site.

333        Further analysis suggested that the sequences at predicted alternative TISs are less

334    conserved than the annotated TISs (Fig. S3H), suggesting that these alternative TISs

335    may be emerged more recently during evolution.   The predictions of alternative TISs

336    were further supported by the Ribo-seq data obtained in two different cell lines (8),

337    where the ribosome density in the 5′-UTRs could reflect the noncanonical translation

338    initiation driven by alternative TIS.   Consistently, the transcripts with a predicted

339    alternative TIS showed significantly higher ribosomal read density in their 5′-UTRs

340    compared to the mRNAs without alternative TIS (Fig. S3I), suggesting a leaky

341    translation in the UTRs due to alternative translation initiation.

342

343    **TranslationAI identifies non-canonical ORFs in human transcriptome**

344        It was recently reported that non-coding RNAs in the human transcriptome or the

345    non-coding regions of mRNAs may contain unannotated ORFs that are translated into

346    new proteins or small peptides (5,8).   Therefore, we next used TranslationAI to

347    examine the noncanonical ORFs in human transcriptome.   Using the stringent cutoff

348    for positive TISs/TTSs (i.e., same cutoff in mRNAs), we were able to identify 4620

349    noncanonical ORFs defined by new TIS-TTS pairs (Fig. 3A).    These newly

350    predicted ORFs included 673 upstream ORFs (uORFs, exemplified in Fig. 3B), 127

351    downstream ORFs (dORFs, exemplified in Fig. 3B), 26 overlapping ORFs with

352    different reading frames from the annotated ORFs (Table S4). To validate the

353    predicted non-canonical ORFs, we compared our prediction with uORFdb that is

354    constructed from comprehensive and meticulous curation of uORF-related literature

355    (32). The analysis revealed that a remarkable 569 out of 673 (85%) uTISs were also

356    corroborated by uORFdb. In addition, 30 out of 673 (4%) uTISs were experimentally

357    supported in TISdb using ribo-seq data from HEK293 (33) (Table S4).

358          The model also identified 3794 new ORFs from annotated noncoding RNAs (Fig.

359    3A, Table S5). These results demonstrate the potential for deep learning approaches to

360    uncover previously unrecognized translation events and shed light on new regulatory

361    mechanisms in gene expression. We next validated the prediction by comparing the

362    predicted TISs in new ORFs with all TISs experimentally identified *via* ribosome

363    profiling (8).    We found that TranslationAI accurately predicted 99.6% of the

364    annotated ORFs, however a smaller fraction of the noncanonical ORFs supported by

365    Ribo-seq were successfully predicted by TranslationAI (39% of uORFs, 37% of

366    dORFs, 24% of dual reading frame, and 43% new ORFs in lncRNAs were identified,

367    Fig. 3C), suggesting that this model is biased towards canonical TISs/TTSs.    Since

368    the training set only contain canonical TIS/TTS information, the TranslationAI model

369    may have learned limited features of non-canonical TISs/TTSs (see discussion).

17

370    We further focused on the thousands of predicted new ORFs from non-coding

371    RNAs (ncRNAs), which are annotated in ENSEMBL as lncRNA, antisense RNA,

372    miRNA precursors, processed transcripts, etc. (Fig. 3D).    The shuffled sequences of

373    all lncRNAs were analyzed using the same deep learning model as a background

374    control.    We found that the TranslationAI indeed predicted much more ORFs in the

375    "non-coding" RNAs than their shuffled counterparts (Fig. 3E).    While previous

376    reports suggested that some lncRNAs can be translated into small peptides (5-8),

377    detailed analysis of these noncanonical ORFs was inadequate. Our results showed that

378    the newly predicted ORFs in lncRNA are significantly longer (Fig. S4A) and more

379    conserved in 46 vertebrates (Fig. S4B) compared to the control ORFs defined by the

380    longest fragments between ATG and TAA/TAG/TGA in all lncRNAs, suggesting a

381    potential functionality.    However, the TISs/TTSs from translatable lncRNAs have

382    lower scores and are less conserved compared to canonical TISs/TTSs from mRNAs

383    (Fig. S4C), implying that the translation efficiency of these newly identified ORFs

384    may be lower than those from annotated mRNAs.    This could partially explain why

385    these 'lncRNAs' were originally defined as non-coding RNAs.

386    To validate the prediction of these new ORFs, we re-analyzed the Ribo-seq data

387    from iPSC and cardiomyocytes (8).    In total 64 newly predicted ORFs by

388    TranslationAI were supported by Ribo-seq signals in either cell line (Table S5).    The

389    ribosome density of translatable lncRNAs closely resembled the footprints from

390    control mRNAs, with strong trinucleotide periodicity that is indictive of active

391    translation (Fig. 3F and Fig. S4D).    In addition, mass spectrometry (MS)-based

392     proteomics confirmed the stable expression of 191 non-canonical peptides (Table S5,

393     with two examples shown in Fig. 3G).   Among them, the ORFs in 10 lncRNAs were

394     supported by both Ribo-seq and MS.   The low number of experimentally validated

395     ORFs may be contributed by several factors, including the cell line specific

396     expression of most lncRNAs, the generally low expression level of lncRNAs, and the

397     limited recover rate in both Ribo-seq and MS experiments.

398          Notably, we have correctly predicted 10 out of the 13 translated lncRNAs

399     validated by independent functional studies, including the HOXB-AS3 that was

400     reported to encode a 53aa peptide (34), the LINC00961 that encodes a 90aa peptide

401     involved in the mTOR activation (6), and the RP11-132A1.4 that encodes a functional

402     peptide of 124aa (8). All of these three lncRNAs were mis-annotated in RefSeq,

403     probably because their expression and translation are limited to specific cell types. It

404     is also worth noted that the predicted ORFs in annotated lncRNAs generally have low

405     scores in our model (Table S5), which could be attributed to the short ORF length or

406     the presence of a long 5′-UTR. Nevertheless, these results highlight the potential

407     application of AI model in the future study of translatable lncRNA.

408

409     **TranslationAI accurately predicts TISs of other eukaryotes and viruses**

410          An "universal genetic code" is used from bacteria to mammals, suggesting that

411     translation machinery is highly conserved across different organisms.   Therefore, we

412     sought to examine if the model trained on human transcriptome can also predict

413     canonical translation sites in other organisms. The accuracies of translation

414     predictions for all eukaryotes tested (Human, Mouse, Zebrafish, Drosophila,

415     Arabidopsis and budding yeast *S. cerevisiae*) were remarkably high (> 90% for all

416     predictions, except that the yeast TIS prediction accuracy is 89%, Fig. 4A), suggesting

417     that our model has learned the essential features generally required for translation

418     initiation and termination in all eukaryotes.

419         In addition, we compared TranslationAI with two existing computational models,

420     a deep learning model TITER(20) and a linear regression model TIS-predictor(23).

421     Both previous models incorporated the surrounding local sequence as input to predict

422     translation initiation sites, with TITER considering 200 nt and TIS-predictor

423     considering 23 nt on both sides of TIS. Notably, we found that TranslationAI

424     surpassed the other two models as judged by both the Area under the ROC Curve

425     (AUC) and PR-AUC of TIS prediction across all tested eukaryotes (Fig. 4A and Table

426     S1). The high performance of TranslationAI on canonical translation sites underscores

427     the significant advance of our model over the previous ones, as the new model

428     considers the entire mRNA rather than focusing solely on a limited region around

429     candidate AUGs. Additionally, while TITER and TIS-predictor were primarily

430     designed to predict non-canonical TIS by incorporating both AUG and near-cognate

431     start codons, TranslationAI was specifically developed to predict translation sites

432     without any prior bias toward non-canonical translation, and thus outperforms TITER

433     and TIS-predictor in predicting canonical TISs that account for the majority of our

434     training dataset.

435         We further assess the performance of TranslationAI model trained with human

20

436    transcripts on the polycistronic transcription units of E. *coli*, which may reflect how

437    much of the coding sequence might be conserved between eukaryotes and

438    prokaryotes.    Surprisingly, we obtained a 65% accuracy in *E. coli* with the deep

439    learning model using all parameters trained on the human transcriptome (Fig. S5),

440    even though most *E. coli* genes are transcribed as polycistronic units with very short

441    UTRs and different set of genetic codons.    Considering the lack of transcriptomic

442    and proteomic information for most bacteria, this model could be useful in predicting

443    their ORFs based solely on their genome.    Furthermore, we extended this model on

444    the transcriptomes of chloroplast (from *Arabidopsis*) and mitochondrion (from

445    human), and the resulting prediction accuracies were ~30% and <1%, respectively

446    (Fig. S5).    Given that the translation mechanism and the genetic codon in chloroplast

447    and mitochondrion are drastically different from the human cells, such drop of

448    prediction accuracy is not unreasonable.

449        Finally, we tested whether this model can predict the TISs/TTSs of human viruses

450    that rely on the host cell machinery for protein synthesis, albeit through different

451    translational mechanisms (35).    Taking two RNA viruses (Ebola and SARS-CoV-2)

452    as examples, we found that the model trained with human data can accurately predict

453    all TISs and TTSs of Ebola virus even when we used the entire genomic RNA

454    sequence as the input (Fig. 4B and Fig. S6A).    In addition to the annotated TISs and

455    TTSs, we also predicted a new internal translation start site in the GP gene in Ebola

456    (Fig. 4B), suggesting an additional translation isoform for GP with truncated N-

457    terminus.

458    Compared to Ebola virus, TranslationAI showed a lower prediction accuracy for

459    SARS-CoV-2, with 8 out of 12 annotated TISs and 10 out of 12 annotated TTSs being

460    predicted with high confidence by sub-genomic RNAs (36) (Fig. 4C and Fig. S6B). In

461    addition, our model also made several false positive predictions of TISs/TTSs in

462    SARS-CoV-2 genome (Fig. S6B).    Similar performance was observed with either

463    genomic or sub-genomic RNAs as input.    A notable difference between these two

464    viruses is that the Ebola genome is transcribed into monocistronic mRNAs with both

465    UTRs, whereas the SARS-CoV-2 mRNAs are mainly polycistronic genomic or

466    subgenomic RNAs with short spacer sequences between each ORF (36,37).    This

467    may partially explain the different performances of TranslationAI for different

468    viruses.

469    Taken together, our results suggest that the model containing parameters trained

470    on human transcriptome can accurately predict canonical ORFs in a wide range of

471    organisms. However, in certain cases where the translation regulation is drastically

472    different (e.g., short UTRs or different codon system), the application of this model

473    should be carried out with caution.

474

**Discussion**

As one of the most effective and powerful machine learning methods, the deep learning network has been widely used to study complex biological questions ranging from population genetics to precision medicine (12-15).   Here we developed a deep learning architecture, TranslationAI, for *ab initio* prediction and systematic analysis of translation initiation and termination in entire transcriptome.   Remarkably, the model trained with ~70% of human transcripts achieved a near-perfect prediction accuracy in not only human, but also a wide range of eukaryotic organisms.   The model can also identify new ORFs in annotated noncoding RNAs, some of which were confirmed by independent Ribo-seq and mass spectrometry experiments.   Further analysis of this predictive model also identified a new regulation rule for translation termination, providing mechanistic insights into the complex regulation of RNA translation.

The final TranslationAI-2k model strategically utilized a flanking sequence of 1000 nt around the position of interest. We chose this design aiming to address the limitations observed in prior models that predominantly focused on characteristics of local sequences (20,23). By incorporating long-range sequence features, this model effectively captures potential interactions within mRNA sequences, including the long-range interactions of both UTRs, detecting frame shifts following one or two nucleotide deletions, and examination of codon bias within CDS.

The remarkable accuracy of TranslationAI suggests that this "black box" model have learned useful information with biological relevance.   We found that the model has successfully learned the identity of canonical start and stop codons (Fig. 1F), which account for the vast majority (99.4% start codon and 99.6% stop codon) of the

498    training set.   The model did not identify any TIS/TTS with noncanonical start or stop

499    (i.e, non-AUG start or non-UAG/UAA/UGA stop), probably because they are

500    extremely rare in the training data.   Additional features, including triplet reading

501    frameshift, UTR sequences, and codon bias also contribute to the prediction accuracy.

502    In addition, the long-range interactions contributed to the prediction accuracy of this

503    model, as simultaneously disrupting the 5′-UTR and 3′-UTR cause synergistic effect

504    in reducing the prediction scores (Fig. 1G).   Moreover, certain internal sites in each

505    mRNA also contributed to the TranslationAI prediction in a context-dependent

506    fashion (Fig. S2B), however the underlying biological mechanisms are unclear.   All

507    these features contributed partially to the model in terms of TIS/TTS prediction, as

508    disrupting these features did not completely invalidate the prediction (Fig. 1G-I).

509        Surprisingly, several known features associated with translation regulation,

510    including the 3′-UTR sequence, synonymous mutation, and sequences of translation

511    product, showed small effect on the accuracy of TIS/TTS prediction, implying that

512    they may not play a dominant role in selection of translation sites.   More

513    importantly, we found several new features can influence the prediction of TIS/TTS,

514    such as the G/C rich motif around the strong TTS, suggesting additional determinants

515    for translation termination.   These features may open up new avenues for

516    investigating the molecular mechanisms underlying translation control. Future

517    exploration of these predictive features may provide mechanistic insights into the

518    functional interactions between the mRNA and translation machine.

519        We also used this model to assess the single nucleotide variants (SNVs) that may

520     alter the TISs/TTSs and cause functional consequences. Surprisingly, we did not

521     identify any disease-associated SNVs in the ORF that can alter the TIS/TTS

522     prediction, unless the mutation directly changed the start or stop codon. This is

523     consistent with the finding that the synonymous/non-synonymous substitutions of

524     amino acids had little effect on TIS/TTS prediction (Fig. 1I). These results suggest

525     that small changes in coding sequence do not significantly impact ORF selection by

526     the translation machinery, *i.e.,* the translation initiation and termination appear to be

527     unaffected by the elongation steps. Alternatively, it is also possible that the model

528     does not rely on the nucleotide composition within ORFs for its prediction.

529     This model also allowed us to identify novel regulatory motifs associated with

530     TIS and TTS. We discovered, for the first time, that the strong stop codons are

531     surrounded by novel motifs that are generally CG-rich, especially in the third position

532     of the codons immediately before the stop codon. Consistently, there is a codon bias

533     associated with the strong termination site (Fig. 2E), and the synonymous mutations

534     with low CG-contents will introduce translation readthrough in a reporter assay (Fig.

535     2F). Because the GC-rich codons were found to be associated with slow translation

536     (38), we speculate that the strength of the stop codon is associated with the decreased

537     translation elongation rate prior to reaching the stop codon. This regulation may be

538     attributed to the changes in the strength of ribosome's interaction with the mRNA

539     before the ribosomes reach the termination site, which "decelerate" the translation

540     before an efficient termination. Alternatively, the motif could alter the ribosome's

541     conformation, thereby influencing translation termination. This finding suggests that

25

542     the translation regulation is more intricate than previously thought, and further

543     exploration on this phenomenon may shed light on the mechanism of translation

544     termination and readthrough.

545         The deep learning network trained with human transcripts achieved a surprisingly

546     high accuracy for the prediction of ORFs in bacteria, yeast, plants and certain viruses

547     using only their genome sequences as input (Fig. 4). This result is particularly

548     interesting given that some of these organisms have polycistronic transcription units

549     or genomes, suggesting that the TranslationAI can predict not only uORFs and dORFs

550     in human, but also ORFs in polycistronic sequences from other organisms.    However

551     this model failed to predict mitochondrial ORFs and showed mediocre performance

552     for SARS-CoV-2, probably due to the presence of frameshift regions, overlapping

553     ORFs, and a lack of sufficient intergenic sequences to separate different ORFs.

554     These results suggest that additional training may be required to improve the model's

555     performance on more complex or noncanonical genomes.

556         Despite its high accuracy in predicting canonical TISs, TranslationAI has some

557     limitations.    First, our model currently struggles in predicting non-canonical TISs,

558     especially for uTISs with non-AUG start codons.    This is probably due to the bias in

559     the training dataset, which mainly consists of mono-cistronic translation units with

560     AUG start codon.    We have trained the model with non-AUG TIS, however

561     distinguishing positive TISs from background proved challenging, probably due to the

562     insufficient and noise experimental data for non-AUG start codon in human

563     transcriptome.    We also attempted to train this model using a subset of mRNAs with

564    multiple annotated ORFs, however the accuracy (~75%) was not as high as that of

565    canonical transcripts, again possibly due to the small size of the training data.

566    Furthermore, although this model can predict the strengths of different TISs in the

567    presence of multiple ORFs, it is still challenging to predict which TIS will be used in

568    specific cell lines or tissues.    Expanding the ribosome profiling data in different cell

569    types may help refine the tissue-specific translation prediction model.    Finally, the

570    model uses a stringent threshold under the assumption that there is a dominant ORF in

571    each transcript.    While this assumption holds well for most transcripts, it may be

572    incorrect in some cases.    With more non-canonical ORFs being discovered, the score

573    threshold may need to be modified to better fit different type of transcripts.    In

574    summary, with increasing translation datasets and more refined parameters, we

575    believe that the TranslationAI framework will improve ORF prediction and deepen

576    our understanding of the complexity in translation control.

577

578 **Methods**

579 **Architecture of the TranslationAI deep learning model**

580      We constructed a deep residual neural network using one-hot encoding to

581 represent mRNA sequences as input (Fig. 1A). Our approach employed a 32-layer

582 dilated convolutional neural network architecture (Fig. S1) to generate an output

583 matrix representing the probabilities of a given position being a TIS, TTS, or neither

584 (NS).

585      Four models were developed with distinct architectures, incorporating 60, 200,

586 600, or 2,000 nucleotides on both sides of a position of interest within the full-length

587 mRNA sequence as input (Fig. S1). The segmentation of sequences into normalizing

588 input sizes (60, 200, 600, or 2,000 nt) facilitates batch processing within the

589 programming framework. Each model produces probability scores for TIS, TTS, and

590 NS, with the sum of these probabilities equaling one.

591      The TranslationAI consists of several stacked residual blocks that connect the

592 input layer to the penultimate layer, a framework previously developed for image

593 recognition (39).   A convolutional unit with softmax activation links the penultimate

594 layer to the output layer. To enhance convergence speed during training, the output of

595 every fourth residual block is added to the input of the penultimate layer (Fig. S1).

596      The fundamental unit of the TranslationAI model is a residual block, which

597 contains two batch-normalization layers, two rectified linear units (ReLU), and two

598 convolutional units arranged in a specific sequence (Fig. S1). Each residual block

599 includes three hyper-parameters: N, W, and D. N represents the number of

600 convolutional kernels, W signifies the window size, and D indicates the dilation rate

601 of each convolutional kernel. Given that a convolutional kernel with window size W

602 and dilation rate D processes features across (W-1)×D neighboring positions, a

603 residual block with two convolutional units can handle features across 2(W-1)×D

604 neighboring positions. As a result, a TranslationAI model with K stacked residual

605 blocks is capable of extracting features spanning $\sum_{i=1}^{K} 2(W_i - 1)D_i$ neighboring

606 positions, where $N_i$, $W_i$, and $D_i$ correspond to the hyper-parameters of the i-th residual

607 block.

608

609 **Model training and testing**

610      To obtain a comprehensive, well-curated, and non-redundant set of sequences, we

611 utilized the RefSeq gene annotation and extracted 47,098 protein-coding transcripts

612 with 47,098 TIS-TTS pairs. Different splicing isoforms of the same gene are treated

613 as different transcripts. For model training, we selected transcripts located on

614 chromosomes 2, 4, 6, 8, 10-22, X, and Y, which contain 13,707 genes with 34,292

615 TIS-TTS pairs.   For model evaluation, the transcripts from chromosomes 1, 3, 5, 7,

616 and 9, containing 5,524 genes with 12,806 TIS-TTS pairs, were used. To mitigate

617 overfitting, 10% of the training set were randomly chosen for early-stopping

618 determination during training, while the remainder were used for model training.

619      We employed a sequence-to-sequence approach with a chunk size of 6,000 for

620 model training and testing. The full-length mature mRNA sequences were one-hot

621 encoded, with A, C, G, and U mapped to [1, 0, 0, 0], [0, 1, 0, 0], [0, 0, 1, 0], and [0, 0,

622 0, 1], respectively. We then zero-padded the one-hot encoded nucleotide sequences

623 until their lengths were multiples of 6,000 and further padded them with flanking

624 sequences of length S/2 at the beginning and end, where S corresponds to 60, 200,

625 600, or 2,000 for the TranslationAI-60, TranslationAI-200, TranslationAI-600, and

626 TranslationAI-2k models, respectively. The padded nucleotide sequences were

627 divided into blocks of length S/2 + 6,000 + S/2, with the i-th block consisting of

628 nucleotide positions from 6,000(i-1) – S/2 + 1 to 6,000i + S/2. In parallel, the output

629 label sequences were one-hot encoded as TIS ([0, 1, 0]), TTS ([0, 0, 1]), or neither ([1,

630 0, 0]). The one-hot encoded TIS/TTS output label sequences were also zero-padded

631 until their lengths reached multiples of 6,000, and then split into 6,000-length blocks,

632 with the i-th block containing positions from 6,000(i-1) + 1 to 6,000i. These one-hot

633 encoded nucleotide sequences and their corresponding label sequences served as the

634 inputs and outputs for the models, respectively.

635 The models were trained for 10 epochs using a batch size of 18 on two NVIDIA

636 Tesla K80 GPUs. During training, the categorical cross-entropy loss between target

637 and predicted outputs was minimized using the Adam optimizer. The optimizer's

638 learning rate was set to 0.001 for the first six epochs and halved in each subsequent

639 epoch. For each model, the training procedure was repeated 10 times, and the top five

640 models based on performance were selected. During testing, each input was assessed

641 using all five trained models, and the average of their outputs was utilized as the

642 predicted output. This rigorous training and testing process aimed to ensure optimal

643 performance and reliability in the resulting models.

644    The performance of the TIS/TTS prediction models was measured with top-k

645    accuracy. This metric measures the fraction of correctly predicted positions among the

646    top-k predicted sites, where k represents the number of true TIS/TTS sites in the

647    dataset. For example, if the dataset contains k1 TISs and k2 TTSs, we set the score

648    cutoff to predict exactly the top-k1 and top-k2 positions as TIS and TTS sites,

649    respectively. It is worth noting that, unless specifically mentioned otherwise, only the

650    top-k predicted sites were used for downstream analysis in our study.

651    During the development of the model, we have experimented with various hyper-

652    parameters of this mode, including the input window size (60 nt, 200 nt, 600 nt, 2k

653    and even 10k) and different convolutional filter sizes. However, it is worth noting that

654    we did not systematically tune the model hyper-parameters. Therefore, the model

655    hyper-parameters used here may not represent the optimal combinations. The results

656    shown in the study probably represent the lower bound of this deep learning

657    framework in the task of predicting TIS and TTS.

658

659    **Role of transcript lengths in mRNA translation**

660    To study the effect of transcript length on translation sites prediction, we sorted

661    the mRNAs in order of increasing length. Using the TranslationAI-200 and

662    TranslationAI-2k models, we calculated the positive rate of TIS/TTS for each

663    transcript by comparing the number of positive TISs to the total number of transcripts

664    with lengths shorter than the given transcript. We plotted these positive rates against

665    the transcript length, as shown in Fig. 1C. Additionally, we included a histogram of

666    transcript lengths from the 47,098 transcripts analyzed in the background of the

667    figure.

**Perturbations of synonymous mutation**

670    To investigate the impact of codon bias on predictions, we introduced

671    perturbations by substituting codons with synonymous counterparts. Three scenarios

672    were considered: 1) Random Synonymous Substitutions: Codons were randomly

673    replaced with synonymous counterparts. 2) Higher CAI Substitutions: Codons were

674    replaced with synonymous counterparts possessing a higher CAI. The Relative

675    Synonymous Codon Usage (RSCU) for each codon was obtained from the Codon

676    Statistics Database (http://codonstatsdb.unr.edu). The RSCU values were calculated

677    based on genes encoding ribosomal proteins, known for their high expression and

678    consequent strong codon bias (40). 3) Lower CAI Substitutions: Codons were

679    replaced with synonymous counterparts possessing a lower CAI, following the same

680    procedure as described above.

**Occlusion sensitivity analysis**

683    To assess the feature importance of individual positions in predicting TIS and

684    TTS, we performed the occlusion sensitivity analysis, which provides a clear

685    visualization of the contribution of each nucleotide to the scoring of TIS and TTS in

686    an mRNA sequence. We masked each nucleotide in the mRNA with "N" one by one,

687    and calculated the resulting score change of the predicted TIS/TTS.

688

## Codon distribution at upstream of stop codon

689

690    To examine the codon distribution upstream of stop codon, we quantified the

691    frequency of each codon at the -30nt position relative to both strong and weak stop

692    codons. Specifically, we calculated the frequency of each codon as the number of

693    occurrences of that codon at the -30nt position divided by the total number of RNAs

694    analyzed at that position.

695

## Conservation analysis

696

697    To evaluate the conservation levels of different RNAs, we obtained the

698    PhastCons46 scores from the PhastCons46 or PhyloP46 track in the UCSC Genome

699    Browser. These scores were generated from genome alignments of 46 vertebrates,

700    including human (hg19), and provide a measure of evolutionary conservation across

701    species. Each box plot displays the median as the central line, while the upper and

702    lower edges of the box represent the first and third quartiles, respectively.

703

## Non-canonical TISs/TTSs identification

704

705    To systematically identify non-canonical TISs and TTSs in transcriptomes, we

706    input human transcriptome into TranslationAI model and retained only those RNAs

707    that met the threshold criteria (prediction score > 0.5) for both TIS and TTS as

708    translatable RNAs. The ncRNA dataset were retrieved from Ensembl build 75, which

709    was previously used in the published study (8) that identified hundreds of

710    noncanonical ORFs using ribosome profiling and mass-spectrometry.

711        To compare the TISs predicted by TranslationAI with TISs identified by ribosome

712    profiling, we expanded our predicted TIS library by retaining the top 10 TISs for each

713    gene predicted by the TranslationAI. This approach was necessary as ribosome

714    profiling often identifies multiple ORFs from one gene.

715        The alternative TISs were defined as the TIS with the highest score ($>0.1$),

716    excluding the annotated TIS. In addition, in-frame alternative TISs were defined as

717    TISs where the distance between the annotated TIS and alternative TIS is a multiple

718    of three.

719

720    **Ribosome footprint data analysis**

721        The ribosome profiling data analysis was performed on a previously published

722    ribosome footprint data that were retrieved from no drug treatment iPSC and iPSC-

723    induced iPSC-derived cardiomyocytes cells being deposited in GEO database under

724    accession number GSE131650. The genome assembly used throughout this

725    manuscript is hg19/GRCh37 annotated by GENCODE v27. We first filtered out reads

726    aligning to rRNAs using Bowtie v1.0.0 and aligned the remaining reads to the

727    annotated transcriptome with STAR v2.7.5a, using the --outSJfilterReads Unique --

728    outFilterMultimapNmax 1 --outFilterMismatchNmax 999 --

729    outFilterMismatchNoverReadLmax 0.04 to filter multiple aligned and redundant

730    reads. These alignments were assigned a specific P-site nucleotide using a 12-nt offset

731    from the 3′ end of reads. Metagene plots were generated by normalizing read density

34

732 around the start and stop codons of each gene (i.e, read depth at each position divided

733 by the average reads of each position in one transcript and the number of the covered

734 transcripts).

735     To estimate the correlation between predicted TIS/TTS scores and mRNA

736 translation efficiency, we ranked and divided the scores into percentiles (n~1580

737 transcripts/percentile, Fig. S2G-S2H). We calculated the translation efficiency of

738 transcripts from ribosome profiling data in two human cell lines as the ratio of Ribo-

739 seq read density to RNA-seq read density.

740

741 **Motif identification of strong and weak TISs/TTSs**

742     Motif discovery and analysis were performed using the Meme suite (v5.0.5). The

743 strong and weak TIS/TTS containing sequences (-30nt : + 33nt) were collected as

744 primary sequences for motif prediction. Additionally, all TIS/TTS containing

745 sequences (-30nt : +33nt) were used as control sequences. Meme was employed to

746 identify enriched motifs, with a maximum length of 63nt, that occur any number of

747 times within the given sequence.

748

749 **MS Proteomics and HLA Peptidomics data analysis**

750     MS-based proteomics data are deposited to the ProteomeXchange Consortium via

751 the Proteomics Identifications Database (PRIDE) partner repository with the dataset

752 identifier PXD014031 and PXD000394. Using an open search engine pFind (v3.1.3),

753 we searched the two previously published human proteome datasets against a combined

35

754    database (25692 proteins) containing all UniProt human proteins and the potential non-

755    canonical ORFs-coded peptides (uORFs, dORFs, dual frame ORFs, and new ORFs).

756    We selected positive mass spectra using following thresholds: q < 0.01, peptides length

757    $\geqslant 8$, missed cleavage sites $\leqslant 3$, allowing only common modifications (cysteine

758    carbamidomethylation, oxidation of methionine, protein N-terminal acetylation, pyro-

759    glutamate formation from glutamine, and phosphorylation of serine, threonine, and

760    tyrosine residues).

761        The resulting peptides were searched against non-redundant human protein

762    database using blastp-short and the peptides with less than two mismatches from known

763    proteins were removed.

764

765    **Transcriptomes of other species**

766        The genome sequence and gene annotations of *Mus musculus* (mm10), *Danio*

767    *rerio* (danRer11), *Drosophila melanogaster* (dm6), Ebola virus (eboVir3), SARS-

768    CoV-2 (wuhCor1), and human mitochondrion are obtained from UCSC database. The

769    sub-genomic sequences of SARS-CoV2 were obtained from a published study (36).

770    The genome sequence and gene annotation of Arabidopsis thaliana (TAIR10) are

771    obtained from TAIR database. Regarding E. coli (Escherichia coli K-12 substr.

772    MG1655), its transcription unit is retrieved from EcoCyc. The mRNA information (5′-

773    UTR, CDS, and 3′-UTR) for yeast was obtained from *Saccharomyces* Genome

774    Database (SGD). The genome sequence and gene annotation of Chloroplast

775    (Arabidopsis thaliana) was obtained from NCBI database. The genome sequence of

776  human mitochondria was obtained from NCBI database and its gene annotation file

777  was obtained from GENCODE database.

778

779  **Model comparison**

780  TITER(20) (Translation Initiation siTE detectoR) is a deep learning-based

781  framework designed to predict noncanonical TISs by incorporating both AUG and

782  near-cognate start codons using high-throughput sequencing data. With the input of

783  each TIS with its 200 nt flanking nucleotides, TITER outputs the probability of

784  translation initiation for the given sequence. The source code for TITER was

785  downloaded from https://github.com/zhangsaithu/titer.    The TIS-predictor(23) is a

786  machine learning model that was trained on instances of AUG and near-cognate start

787  codon. The algorithm examines each codon and its 20 flanking nucleotides (23

788  nucleotides in total) to simulate and compute its probability score as a translation

789  initiation site. The source code for TIS-predictor was downloaded from

790  https://github.com/Agleason1/TIS-Predictor.

791  To evaluate the performance of TITER, TIS-predictor, and TranslationAI on

792  canonical TIS AUG, we calculated the AUC, PR-AUC, and accuracy.    For each RNA

793  sequence input, we computed the probability score for all AUG codons, predicting the

794  AUG with the highest score as the translation initiation site for each mRNA. These

795  tools were tested on datasets from multiple species, including a Human test dataset

796  and the whole transcriptomes of Mouse, Zebrafish, Drosophila, Arabidopsis, and the

797  budding yeast (*S. cerevisiae*).

798

**Plasmid construction and Transfection**

799

To construct the reporter (pcDNA5-HA-GFP-stop_motif-Nluc-Flag) for testing

800

the translation readthrough products, the Fluc-Nluc luciferase sequence (a gift from

801

prof. Rachel Green) was cloned into pcDNA5 backbone digested by SpeI and XbaI

802

(41). The GFP sequence was amplified (fused with two BsmBI restriction fragments)

803

and inserted before Nluc. The modified VDR sequences were synthesized

804

(GENEWIZ) and cloned into BsmBI digested vector.

805

HEK293T cell line was cultured in DMEM (high glucose) medium containing

806

10% fetal bovine serum (FBS, Hyclone). To transient transfect plasmids into cells,

807

2 µg of reporters were transfected into cells in 6 well plate using lipofectamine 3000

808

(Invitrogen) according to the manufacturer's instruction. After 48 h, cells were

809

collected for further analysis of protein level.

810

811

**Western blot**

812

Cells were lysed in laemmli buffer, and the total cell lysates were resolved with 4-

813

20% ExpressPlus™ PAGE Gel (GeneScript). The following antibodies were used:

814

HA-Tag antibody (CST: 3724S) was diluted by 1:2000, V5 antibody (CST: 13202S)

815

was diluted by 1:2000, Flag antibody (Sigma: F1804-1MG) was diluted by 1:2000.

816

The HRP-linked secondary antibodies (CST: 7076S) were used by 1:4000 dilution

817

and the blots were visualized with the ECL reagents (Bio-Rad).

818

819

**Data and software availability**

820

38

821    Training and testing data, prediction scores for all possible single nucleotide

822    substitutions in the reference genome and source code are publicly hosted at GitHub

823    (https://github.com/rnasys/TranslationAI). Prediction scores and source code are

824    publicly released under GPL v3 and are free for use for academic and non-commercial

825    applications.

826

**Acknowledgments**

839

**Author Contributions**

841    Conceptualization, Z.W., X.F. and T.C.; Methodology, X.F., T.C., and Z.W.;

842    Software, X.F. and T.C.; Experiments, C.C.; Writing and discussion, X.F., Z.W., T.C.,

843    and M.H.; Funding acquisition, X.F., and Z.W.

844

845    **Declaration of Interests**

846    The authors declare no competing financial interests.

847

# References

1. Frankish, A., Diekhans, M., Ferreira, A.M., Johnson, R., Jungreis, I., Loveland, J., Mudge, J.M., Sisu, C., Wright, J., Armstrong, J. *et al.* (2019) GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res*, **47**, D766-D773.

2. Mouilleron, H., Delcourt, V. and Roucou, X. (2016) Death of a dogma: eukaryotic mRNAs can code for more than one protein. *Nucleic Acids Res*, **44**, 14-23.

3. Hinnebusch, A.G. (2011) Molecular mechanism of scanning and start codon selection in eukaryotes. *Microbiol Mol Biol Rev*, **75**, 434-467, first page of table of contents.

4. McNair, K., Ecale Zhou, C.L., Souza, B., Malfatti, S. and Edwards, R.A. (2021) Utilizing Amino Acid Composition and Entropy of Potential Open Reading Frames to Identify Protein-Coding Genes. *Microorganisms*, **9**.

5. Jackson, R., Kroehling, L., Khitun, A., Bailis, W., Jarret, A., York, A.G., Khan, O.M., Brewer, J.R., Skadow, M.H., Duizer, C. *et al.* (2018) The translation of non-canonical open reading frames controls mucosal immunity. *Nature*, **564**, 434-438.

6. Matsumoto, A., Pasut, A., Matsumoto, M., Yamashita, R., Fung, J., Monteleone, E., Saghatelian, A., Nakayama, K.I., Clohessy, J.G. and Pandolfi, P.P. (2017) mTORC1 and muscle regeneration are regulated by the LINC00961-encoded SPAR polypeptide. *Nature*, **541**, 228-232.

7. Yang, Y., Fan, X., Mao, M., Song, X., Wu, P., Zhang, Y., Jin, Y., Yang, Y., Chen, L.L., Wang, Y. *et al.* (2017) Extensive translation of circular RNAs driven by N(6)-methyladenosine. *Cell Res*, **27**, 626-641.

8. Chen, J., Brunner, A.D., Cogan, J.Z., Nunez, J.K., Fields, A.P., Adamson, B., Itzhak, D.N., Li, J.Y., Mann, M., Leonetti, M.D. *et al.* (2020) Pervasive functional translation of noncanonical human open reading frames. *Science*, **367**, 1140-1146.

9. Raj, A., Wang, S.H., Shim, H., Harpak, A., Li, Y.I., Engelmann, B., Stephens, M., Gilad, Y. and Pritchard, J.K. (2016) Thousands of novel translated open reading frames in humans inferred by ribosome footprint profiling. *Elife*, **5**.

10. Jaganathan, K., Kyriazopoulou Panagiotopoulou, S., McRae, J.F., Darbandi, S.F., Knowles, D., Li, Y.I., Kosmicki, J.A., Arbelaez, J., Cui, W., Schwartz, G.B. *et al.* (2019) Predicting Splicing from Primary Sequence with Deep Learning. *Cell*, **176**, 535-548 e524.

11. Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool, K., Bates, R., Zidek, A., Potapenko, A. *et al.* (2021) Highly accurate protein structure prediction with AlphaFold. *Nature*, **596**, 583-589.

12. Eraslan, G., Avsec, Z., Gagneur, J. and Theis, F.J. (2019) Deep learning: new computational modelling techniques for genomics. *Nat Rev Genet*, **20**, 389-403.

13. Bogard, N., Linder, J., Rosenberg, A.B. and Seelig, G. (2019) A Deep Neural Network for Predicting and Engineering Alternative Polyadenylation. *Cell*, **178**, 91-106 e123.

14. Zhou, J., Theesfeld, C.L., Yao, K., Chen, K.M., Wong, A.K. and Troyanskaya, O.G. (2018) Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk. *Nat Genet*, **50**, 1171-1179.

15. Vamathevan, J., Clark, D., Czodrowski, P., Dunham, I., Ferran, E., Lee, G., Li, B., Madabhushi, A., Shah, P., Spitzer, M. *et al.* (2019) Applications of machine learning in drug discovery and development. *Nat Rev Drug Discov*, **18**, 463-477.

16. Sonenberg, N. and Hinnebusch, A.G. (2009) Regulation of translation initiation in eukaryotes: mechanisms and biological targets. *Cell*, **136**, 731-745.

17. Lee, S., Liu, B., Lee, S., Huang, S.X., Shen, B. and Qian, S.B. (2012) Global mapping of translation initiation sites in mammalian cells at single-nucleotide resolution. *Proc Natl Acad Sci U S A*, **109**, E2424-2432.

18. Gao, X., Wan, J., Liu, B., Ma, M., Shen, B. and Qian, S.B. (2015) Quantitative profiling of initiating ribosomes in vivo. *Nat Methods*, **12**, 147-153.

19. Chen, W., Feng, P.M., Deng, E.Z., Lin, H. and Chou, K.C. (2014) iTIS-PseTNC: a sequence-based predictor for identifying translation initiation site in human genes using pseudo trinucleotide composition. *Anal Biochem*, **462**, 76-83.

20. Zhang, S., Hu, H., Jiang, T., Zhang, L. and Zeng, J. (2017) TITER: predicting translation initiation sites by deep learning. *Bioinformatics*, **33**, i234-i242.

21. Goel, N., Singh, S. and Aseri, T.C. (2020) Global sequence features based translation initiation site prediction in human genomic sequences. *Heliyon*, **6**, e04825.

22. Kalkatawi, M., Magana-Mora, A., Jankovic, B. and Bajic, V.B. (2019) DeepGSR: an optimized deep-learning structure for the recognition of genomic signals and regions. *Bioinformatics*, **35**, 1125-1132.

23. Gleason, A.C., Ghadge, G., Chen, J., Sonobe, Y. and Roos, R.P. (2022) Machine learning predicts translation initiation sites in neurologic diseases with nucleotide repeat expansions. *PLoS One*, **17**, e0256411.

24. Leppek, K., Das, R. and Barna, M. (2018) Functional 5' UTR mRNA structures in eukaryotic translation regulation and how to find them. *Nat Rev Mol Cell Biol*, **19**, 158-174.

25. Mayr, C. (2019) What Are 3' UTRs Doing? *Cold Spring Harb Perspect Biol*, **11**.

26. Fields, A.P., Rodriguez, E.H., Jovanovic, M., Stern-Ginossar, N., Haas, B.J., Mertins, P., Raychowdhury, R., Hacohen, N., Carr, S.A., Ingolia, N.T. *et al.* (2015) A Regression-Based Analysis of Ribosome-Profiling Data Reveals a Conserved Complexity to Mammalian Translation. *Mol Cell*, **60**, 816-827.

27. Kozak, M. (1987) An analysis of 5'-noncoding sequences from 699 vertebrate messenger RNAs. *Nucleic Acids Res*, **15**, 8125-8148.

28. Kozak, M. (1986) Point mutations define a sequence flanking the AUG initiator codon that modulates translation by eukaryotic ribosomes. *Cell*, **44**, 283-292.

29. Loughran, G., Jungreis, I., Tzani, I., Power, M., Dmitriev, R.I., Ivanov, I.P., Kellis, M. and Atkins, J.F. (2018) Stop codon readthrough generates a C-terminally extended variant of the human vitamin D receptor with reduced calcitriol response. *J Biol Chem*, **293**, 4434-4444.

30. Mills, E.W., Wangen, J., Green, R. and Ingolia, N.T. (2016) Dynamic Regulation of a Ribosome Rescue Pathway in Erythroid Cells and Platelets. *Cell Rep*, **17**, 1-10.

31. de Klerk, E. and t Hoen, P.A. (2015) Alternative mRNA transcription, processing, and translation: insights from RNA sequencing. *Trends Genet*, **31**, 128-139.

32. Manske, F., Ogoniak, L., Jurgens, L., Grundmann, N., Makalowski, W. and Wethmar, K. (2023) The new uORFdb: integrating literature, sequence, and variation data in a central hub for uORF research. *Nucleic Acids Res*, **51**, D328-D336.

33. Wan, J. and Qian, S.B. (2014) TISdb: a database for alternative translation initiation in mammalian cells. *Nucleic Acids Res*, **42**, D845-850.

34. Huang, J.Z., Chen, M., Chen, D., Gao, X.C., Zhu, S., Huang, H., Hu, M., Zhu, H. and Yan,

936       G.R. (2017) A Peptide Encoded by a Putative lncRNA HOXB-AS3 Suppresses Colon Cancer
937       Growth. *Mol Cell*, **68**, 171-184 e176.

938   35.   Stern-Ginossar, N., Thompson, S.R., Mathews, M.B. and Mohr, I. (2019) Translational
939       Control in Virus-Infected Cells. *Cold Spring Harb Perspect Biol*, **11**.

940   36.   Wang, D., Jiang, A., Feng, J., Li, G., Guo, D., Sajid, M., Wu, K., Zhang, Q., Ponty, Y., Will, S.
941       *et al.* (2021) The SARS-CoV-2 subgenome landscape and its novel regulatory features. *Mol*
942       *Cell*, **81**, 2135-2147 e2135.

943   37.   Banerjee, A.K., Blanco, M.R., Bruce, E.A., Honson, D.D., Chen, L.M., Chow, A., Bhat, P.,
944       Ollikainen, N., Quinodoz, S.A., Loney, C. *et al.* (2020) SARS-CoV-2 Disrupts Splicing,
945       Translation, and Protein Trafficking to Suppress Host Defenses. *Cell*, **183**, 1325-1339 e1321.

946   38.   Tunney, R., McGlincy, N.J., Graham, M.E., Naddaf, N., Pachter, L. and Lareau, L.F. (2018)
947       Accurate design of translational output by a neural network model of ribosome distribution.
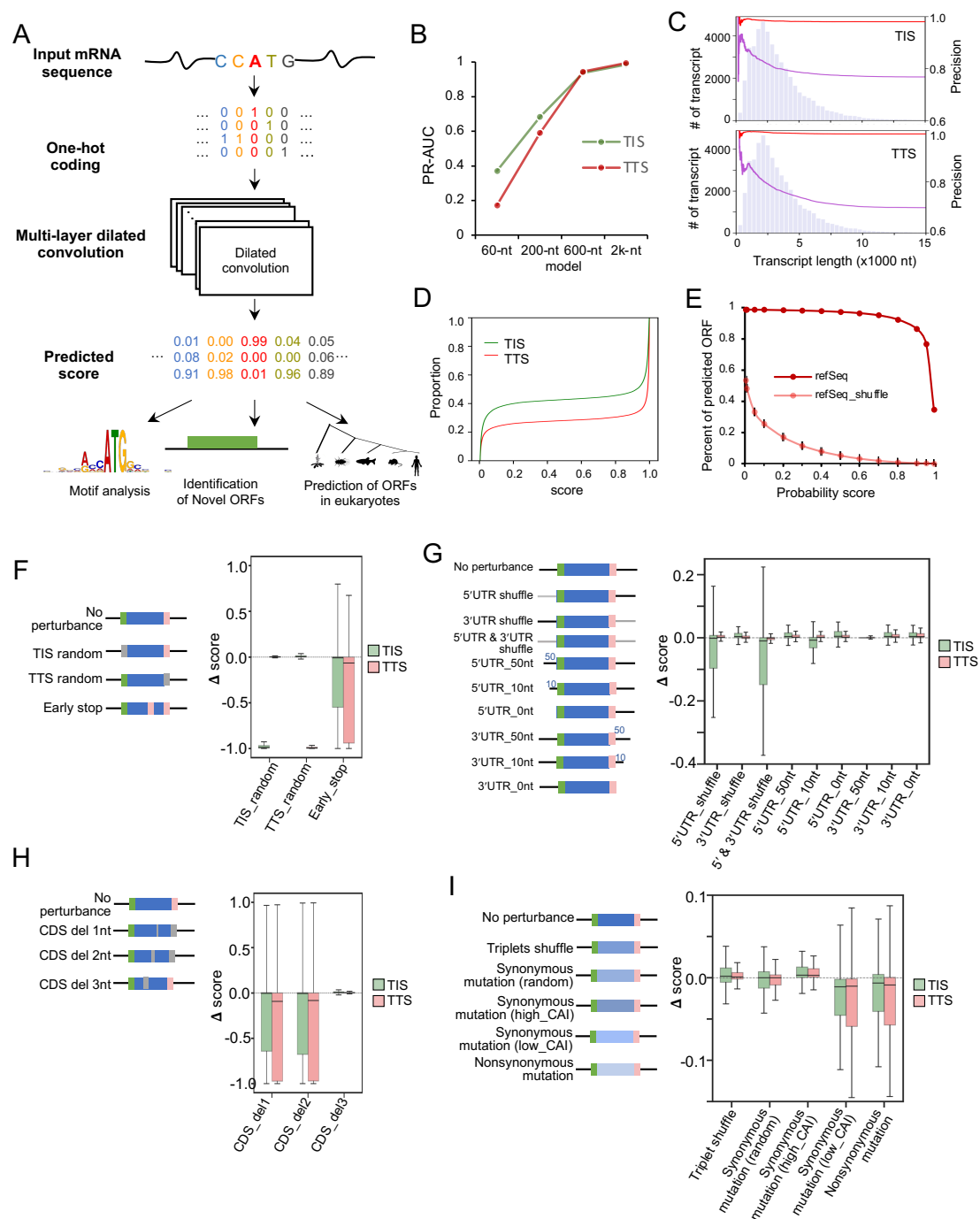948       *Nat Struct Mol Biol*, **25**, 577-582.

949   39.   He, K., Zhang, X., Ren, S. and Sun, J. (2016) Deep Residual Learning for Image Recognition.
950       *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 770-778.

951   40.   Subramanian, K., Payne, B., Feyertag, F. and Alvarez-Ponce, D. (2022) The Codon Statistics
952       Database: A Database of Codon Usage Bias. *Mol Biol Evol*, **39**.

953   41.   Wangen, J.R. and Green, R. (2020) Stop codon context influences genome-wide stimulation of
954       termination codon readthrough by aminoglycosides. *Elife*, **9**.

955

956   **Figures and legends**



957

**Figure 1. Construction of deep learning network for translation prediction. A.**
Flowchart of TranslationAI, a computational model for predicting translation initiation
and termination sites with full length mRNA. For each position in the full length mRNA,
TranslationAI-2k takes 60, 200, 600, and 2,000 nucleotides of flanking sequence as
input, and predicts whether that position corresponds to a translation initiation site (TIS),
translation termination site (TTS), or neither. **B.** Effect of input sequence context size

964 on network accuracy. PR-AUC is the area under the precision-recall curve. The figure

965 shows the performance of the network at four different input sequence context sizes. **C.**

966 Relationship between transcript length and the positive rate of TIS/TTS for each

967 transcript by comparing the number of positive TISs from TranslationAI-200 (purple

968 line) or TranslationAI-2k (red line) to the total number of transcripts with lengths

969 shorter than the given transcript. The distribution of transcript length is shown in the

970 background as a histogram. **D.** Accumulative distribution of TIS/TTS score among all

971 mRNAs. **E.** The number of ORFs predicted by TranslationAI-2k using mRNA and

972 shuffled mRNA (mononucleotide shuffling) as input and different cutoff of TIS scores.

973 The average number of ORFs predicted from multiple shuffling, along with error bars

974 representing three times the standard deviation calculated from shuffled sequences. **F-**

975 **I.** The features learned by TranslationAI.  Systematic *in silico* perturbations on

976 different regions of the mRNA, measured as changes (**Δ** score) in probability scores for

977 the authentic TISs/TTSs. The perturbations include: **F**, replacement of TIS/TTS identity;

978 **G**, changes in UTR length and sequence; **H**, frameshifts by deleting one, two, or three

979 nucleotides; **I**, mutations in coding sequences. Blue box: ORF (open reading frame),

980 black line: 5′-UTR and 3′-UTR, green line: TIS, pink line: TTS. Mononucleotide

981 shuffling was performed for shuffling of 5′-UTR and 3′-UTR. Codon shuffling was

982 performed for triplet shuffle, which maintained the amino acid composition. CAI:
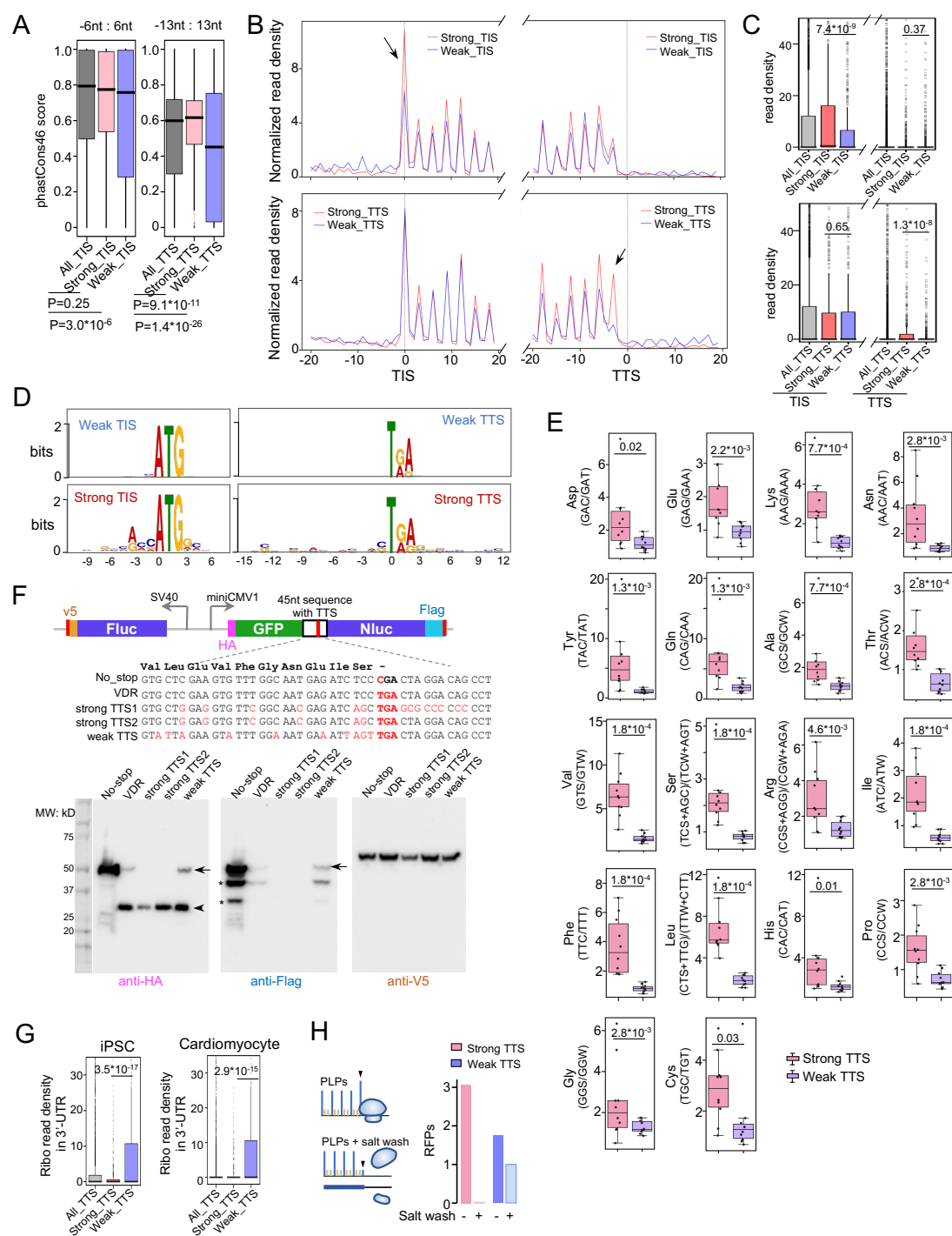
983 Codon Adaptation Index.

984

985

986

987

**Figure 2. Evaluation of predictive features in TranslationAI. A.** Conservation analysis of all TIS/TTS, strong TIS/TTS, and weak TIS/TTS regions. P-values were calculated using Mann–Whitney U test. **B.** Metagene analysis and **C.** boxplot of reads density on transcripts containing strong (red) and weak (blue) TIS/TTS (-20nt : +20nt) in iPSC cell line. P-values were calculated using Mann–Whitney U test. **D.** Motif analysis of strong and weak TIS/TTS. **E.** The codon distribution at the -30nt position

994   of stop codon. The ratios of the same amino acid with C/G and A/T at the third position

995   of each codon before strong and weak stop codons was quantified. P-values were

996   calculated using standard T-test. **F.** Validation of motifs around strong TTS. Known

997   readthrough stop codon from VDR and their mutations (without change of amino acids)

998   were tested for translation readthrough. No_stop: stop codon of VDR (TGA) was

999   mutated to CGA; strong TTS1: stop codon of VDR was mutated to a strong TTS by

1000   changing the upstream 27nt and downstream 12nt sequences; strong TTS2: stop codon

1001   of VDR was mutated to a strong TTS by changing the upstream 27nt sequence; weak

1002   TTS: stop codon of VDR was mutated to a weak TTS by changing the upstream 27nt.

1003   See supplementary information in Table S3. **G.** Boxplot analysis of reads density on 3′-

1004   UTRs from transcripts with strong/weak TTS in iPSC cell line and iPSC-induced

1005   Cardiomyocyte cell line, respectively. **H.** Normalized read density (by abundance of

1006   the transcript) at the last P-site before stop codon of in vitro platelet-like particles (PLP)

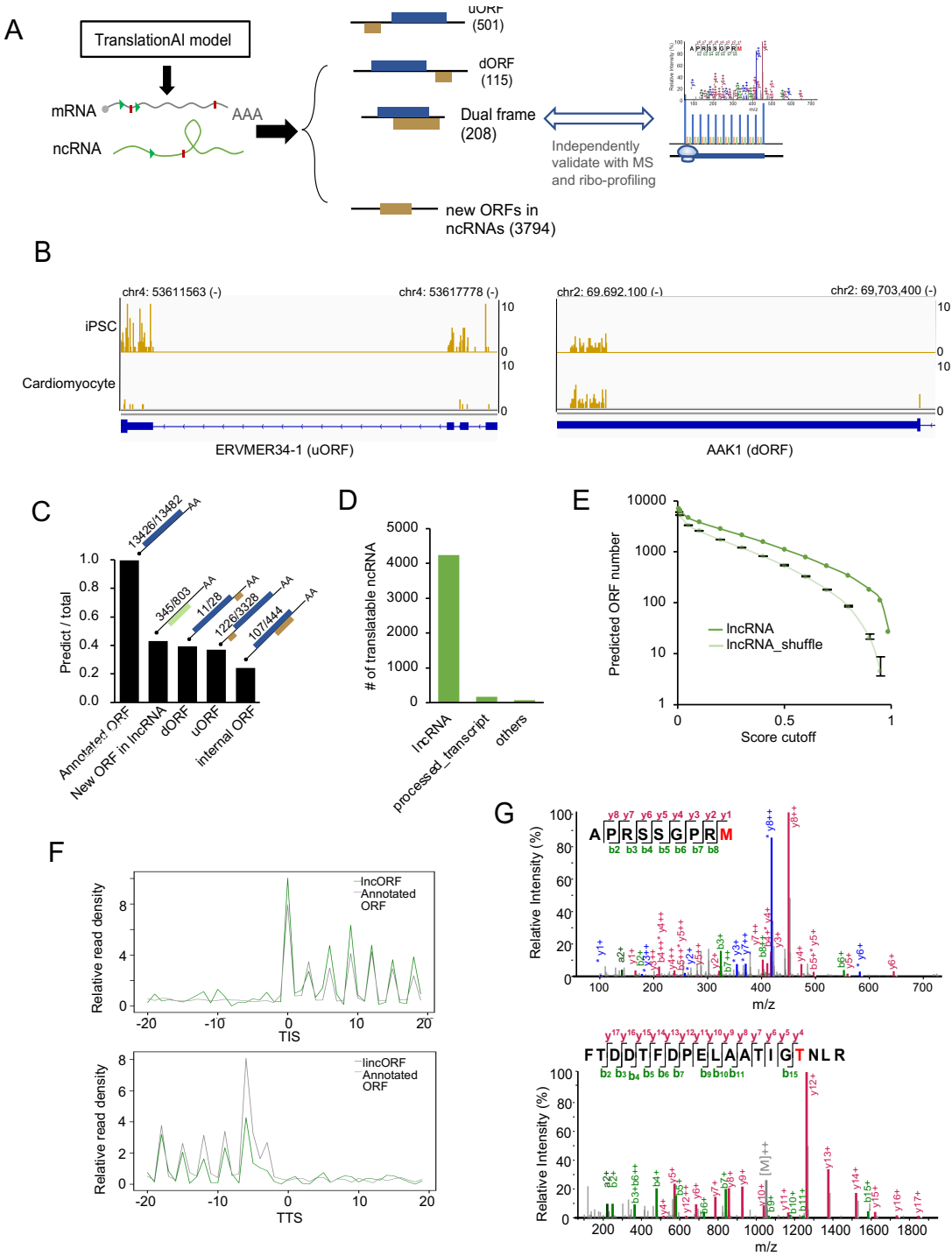1007   and PLP treated with high salt.

1008

1009

1010

**Figure 3. Identification of non-canonical ORFs in human transcriptome. A.** Flowchart for identifying non-canonical ORFs, including upstream ORFs, downstream ORFs, dual coding ORFs, and new ORFs from non-coding RNAs. **B.** Example ribosome footprints of a uORF from ERVMER-1 and dORF from AAK1. **C.** The ratio of various types of predicted TISs identified by another published dataset derived from ribosome profiling assays. **D.** The number of newly identified translatable ncRNA in

1017    annotated lncRNAs, processed transcripts and other transcripts without known ORFs.

1018    **E.** The number of ORFs predicted by TranslationAI-2k using non-coding RNA and

1019    shuffled non-coding RNA sequences as input. The average numbers of ORFs predicted

1020    from lncRNAs and shuffled lncRNAs (mononucleotide shuffling, as control) were

1021    shown, with error bars representing 3×standard deviation calculated from shuffled

1022    sequences. **F.** Metagene analysis of translatable lncRNAs (green line) and control (grey

1023    line, mRNAs with the same ORF length distribution of predicted ORFs from lncRNAs).

1024    **G.** The MS/MS spectra of peptides from two ncRNAs: lncRNA ENST00000609975

1025    (APRSSGPRM)    and    antisense    RNA    ENST00000413405

1026    (FTDDTFDPELAATIGTNLR). The annotated b- and y-ions are marked in red and
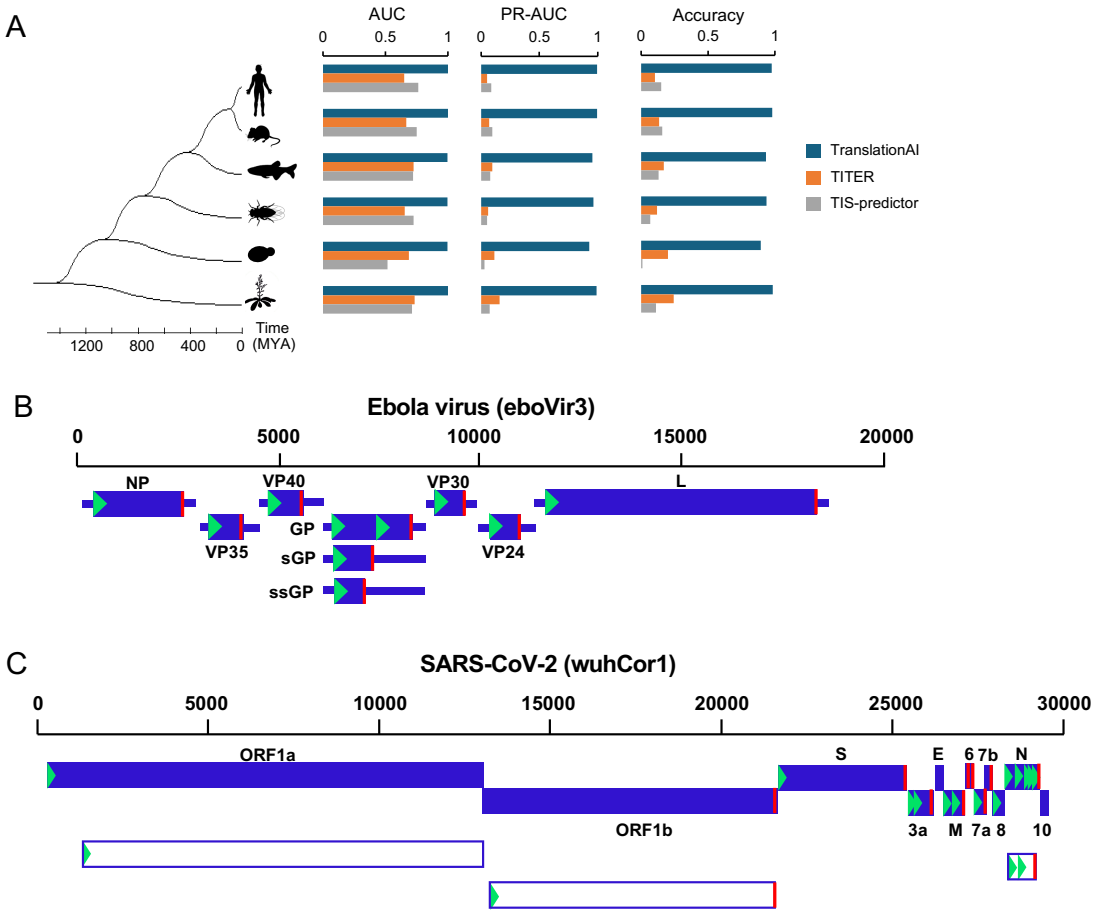
1027    green, respectively.

1028

1029

**Figure 4. TranslationAI accurately predicts TIS/TTS of eukaryotes, prokaryotes, and viruses. A.** The AUC, PR-AUC, and prediction accuracy of TIS prediction across tested eukaryotes (Human, Mouse, Zebrafish, Drosophila, Arabidopsis, and budding yeast *S. cerevisiae*). The predictions of two previous models are also included as a comparison. **B-C.** The prediction of TISs/TTSs on Ebola genomic or RNA sequences **B.** and SARS-CoV-2 genomic or subgenomic sequences **C.** The upper scales represent the genomic sequence position, the blue boxes indicate the annotated ORFs, and the white boxes indicate the newly predicted out-of-frame ORFs. The green triangles and red lines indicate the predicted in-frame TISs and TTSs, respectively.