

Predicting the translation efficiency of messenger RNA in mammalian cells

Dinghai Zheng^{1,*}, Jun Wang^{1,*}, Logan Persyn^{2,*}, Yue Liu², Fernando Ulloa Montoya¹, Can Cenik^{2,†}, Vikram Agarwal^{1,†}

¹mRNA Center of Excellence, Sanofi, Waltham, MA 02451, USA

²Department of Molecular Biosciences, University of Texas at Austin, Austin, TX 78712, USA

*These authors contributed equally to this work

†Correspondence to Can Cenik (ccenik@austin.utexas.edu) and Vikram Agarwal (Vikram.Agarwal@sanofi.com)

ABSTRACT

The degree to which translational control is specified by mRNA sequence is poorly understood in mammalian cells. Here, we constructed and leveraged a compendium of 3,819 ribosomal profiling datasets, distilling them into a transcriptome-wide atlas of translation efficiency (TE) measurements encompassing >140 human and mouse cell types. We subsequently developed RiboNN, a multitask deep convolutional neural network, and classic machine learning models to predict TEs in hundreds of cell types from sequence-encoded mRNA features, achieving state-of-the-art performance ($r=0.79$ in human and $r=0.78$ in mouse for mean TE across cell types). While the majority of earlier models solely considered 5' UTR sequence, RiboNN integrates contributions from the full-length mRNA sequence, learning that the 5' UTR, CDS, and 3' UTR respectively possess ~67%, 31%, and 2% per-nucleotide information density in the specification of mammalian TEs. Interpretation of RiboNN revealed that the spatial positioning of low-level di- and tri-nucleotide features (*i.e.*, including codons) largely explain model performance, capturing mechanistic principles such as how ribosomal processivity and tRNA abundance control translational output. RiboNN is predictive of the translational behavior of base-modified therapeutic RNA, and can explain evolutionary selection pressures in human 5' UTRs. Finally, it detects a common language governing mRNA regulatory control and highlights the interconnectedness of mRNA translation, stability, and localization in mammalian organisms.

Keywords: Translation efficiency, Translational regulation, Ribosome profiling, Machine learning, Deep learning

47 INTRODUCTION

48

49 Protein abundances are determined by the complex interplay of steady-state mRNA levels, mRNA
50 translation rates, and protein turnover rates. Numerous machine learning models have been
51 developed to model the sequence-encoded features that influence steady-state levels of mammalian
52 mRNAs from both the perspectives of transcriptional regulation¹⁻⁵ and mRNA turnover⁶.
53 However, most attempts to model translational regulation from mRNA sequence have focused on
54 bacteria and yeast⁷⁻¹¹. Although such models do exist for mammals, most focus on the functional
55 roles of specific regions such as the 5' untranslated region (5' UTR)¹²⁻¹⁴ or coding region sequence
56 (CDS)^{15,16}, despite the recognition that the full mRNA sequence (*i.e.*, including 3' UTRs) jointly
57 influences translation^{17,18}. Several models consider full-length mRNA, but have either only
58 implicitly modeled translational regulation^{19,20}, or have evaluated only a limited set of cell types
59 while achieving modest performance ($r^2 \approx 0.40$)^{21,22}. Modeling translational regulation more
60 precisely among diverse cell types would elucidate the functional consequences of synonymous,
61 missense, and non-coding mutations in mRNA. Consequently, this would advance the goals of
62 identifying the mechanistic underpinnings of ribosome occupancy and protein abundance
63 quantitative trait loci (rQTL and pQTL, respectively)^{23,24}, diagnosing pathogenic genetic variants,
64 and designing more translationally competent mRNA therapeutics and gene therapies.

65

66 Global translation rates can be estimated through several strategies, including: i) fitting translation
67 rate parameters from differential equations, using measurements of mRNA and protein abundances
68 as well as mRNA half-life^{25,26}; ii) computing protein-to-mRNA ratios (PTRs)^{19,20,27}; iii) polysome
69 profiling, in which ribosomal fractions are run on a sucrose gradient and mRNAs within each
70 fraction are sequenced to estimate their approximate ribosomal loading^{12,13,18,28}; and iv) ribosome
71 profiling (*i.e.*, Ribo-seq), normalizing ribosome density to RNA abundance as a metric for TE²⁹.
72 Of these techniques, the first two strategies are both indirect estimates of translation rate.
73 Importantly, inferred translation rates from the differential equation modeling strategy were shown
74 to be poorly related to experimentally measured rates³⁰, limiting the accuracy of this approach.
75 Moreover, PTRs are partially confounded by protein degradation rates and protein secretion^{19,20,27}.
76 Therefore, of these four methods, polysome and ribosome profiling are considered more direct
77 methods of assessing translation rates³⁰.

78

79 In eukaryotes, translation is regulated at the initiation and elongation steps^{31,32}, which can be
80 modulated by *cis*-acting sequences. In particular, *cis*-regulation of translation initiation has
81 historically been the focus due to its recognition as the rate-limiting step of translation³³. The
82 propensity for secondary structure near the 5' mRNA cap, the sequence context of the translation
83 initiation codon, presence of upstream short open reading frames (ORFs), and binding sites for
84 various RNA-binding proteins provide concrete mechanisms of translational regulation via *cis*-
85 acting elements predominantly in 5' UTRs³⁴. Importantly, the protein coding sequence is also a
86 key determinant of TE. Relatively more is known in unicellular organisms; in particular, codon
87 usage differs significantly across genes, with more abundant proteins utilizing a biased set of
88 codons^{35,36}. The most widely recognized mechanism for codon-specific influence on translation
89 relates to differences in the active pool of corresponding tRNAs³⁷⁻³⁹. Coding sequence differences
90 are also suggested to impact protein expression through secondary structure-mediated mechanisms
91 that do not correlate with tRNA abundance⁴⁰. Moreover, non-synonymous coding variants can
92 alter translation independently from tRNA abundance, translation initiation efficiency, or overall

93 mRNA structure via the interaction of the encoded peptide with the ribosome exit tunnel⁴¹. Parallel
94 work in vertebrate organisms established a link between translation and RNA stability; for
95 instance, certain codons that slow down translation are associated with unstable mRNA^{15,42–46}.
96 Taken together, these studies reveal that the entire mRNA sequence can potentially modulate
97 translation through a variety of mechanisms. However, the contribution of specific functional
98 regions in determining translation of endogenous mRNAs has yet to be described quantitatively.
99 A precise measurement of translation rate would enable a clear-eyed examination of how different
100 sequence properties and functional regions modulate translation rates relative to one another.

101
102 Despite the widespread abundance of ribosomal profiling datasets, attempts to examine the relative
103 contribution of sequence and structural features to the specification of translation rate have been
104 hampered by their inaccessibility in a unified resource. In this study, we systematically assembled
105 a compendium of 1,282 human and 995 mouse ribosome profiling datasets, matched to
106 corresponding RNA-seq data, to derive more precise TE measurements in mammalian cells. This
107 effort reflects the synthesis of the largest and most comprehensive compendium of TE
108 measurements ever assembled to date. Using enhanced measurements of TE, we derived improved
109 sequence-based models towards the goal of improving the predictability of TE from RNA
110 sequence. Our state-of-the-art model RiboNN, a deep convolutional neural network, is capable of
111 predicting the effects of RNA sequences (*e.g.*, including base-modified, therapeutically delivered
112 mRNA) on translational regulation, in agreement with functional measurements derived from
113 massively parallel reporter assays and population genetic data demarcating regions of evolutionary
114 constraint. RiboNN reconciles several limitations of existing models, possessing the following
115 properties: i) it models the impact of the full-length mRNA sequence on TE in numerous cell types,
116 ii) it exhibits superior performance in predicting TE from mRNA sequence, iii) it identifies the
117 location-dependent effects of short, di- and tri-nucleotide features (*i.e.*, including codons) as the
118 key sequence features explaining model performance, and iv) it helps to quantify the relative
119 contributions of different functional regions on TE, a feat which has largely been evaluated
120 qualitatively in the past. Finally, it postulates the existence of a common language underpinning
121 mRNA translation, stability, and localization in mammalian organisms.

122 RESULTS

123 Preparation of a compendium of human and mouse TE datasets from ribosome profiling 124 data

125 To construct a comprehensive, high-quality dataset of TE measurements, we systematically
126 compiled 3,819 human and mouse ribosome profiling datasets from the GEO database. We filtered
127 these into 1,282 human and 995 mouse samples representing matched ribosome profiling and
128 RNA-seq data from numerous tissues and cell types. We then uniformly processed the datasets
129 using an open-source bioinformatics pipeline⁴⁷. We required each sample to pass the following
130 quality control filters: i) $\geq 70\%$ of ribosome-protected fragments (RPFs) mapped to the CDS, and
131 ii) transcripts globally had a minimum average read coverage of 0.1x (detailed in companion
132 manuscript¹¹⁴). This yielded 1,076 human and 835 mouse ribosome profiling datasets. We then
133 calculated TE using a compositional regression approach that overcomes the mathematical biases
134 associated with the commonly used log-ratio approach^{48,114} (**Fig. 1a; Methods**). We summarized
135 the datasets by averaging TEs across samples belonging to the same cell types, yielding matrices
136 of 10,348 genes x 78 cell types for the human and 10,870 genes x 68 cell types for the mouse (**Fig.**
137 **1a, Supplementary Table 1**). This resource enabled us to assess the degree to which TEs are
138 similar among different mRNAs across cell types. We calculated the Spearman's correlation
139 coefficient (ρ) between the TEs of transcripts across all possible pairs of human cell types (**Fig.**
140 **1b**). We observed that most of the cell types were highly correlated to each other, with a small
141 subset possessing low correlation to most other cell types (**Fig. 1b**). This subset appeared to have
142 lower data quality, as measured by a low median read coverage, leading to a large proportion of
143 missing values (**Fig. 1b**). The high correlation between most cell types is suggestive of common
144 translational regulation mechanisms across most cell types. Parallel results were observed for the
145 inter-cell-type comparisons in the mouse (**Supplementary Fig. 2a**).

146 To validate the biological relevance of TEs relative to other methods to measure translational
147 regulation, we compared the TE across cell types with previously reported PTR ratios^{20,27,49} and
148 ribosome load (number of ribosomes per transcript), as measured by polysome sequencing in
149 HEK293T cells¹⁸. We normalized the ribosome load to CDS length because longer CDSs can
150 accommodate more translating ribosomes. Given the strong correlation based upon dataset of
151 origin (**Supplementary Fig. 3**), we evaluated the relationship between the means of each dataset.
152 The ribosome load and mean PTR across tissues²⁰ were positively correlated with our mean TE
153 ($r=0.42$, $\rho=0.4$ and $r=0.52$, $\rho=0.51$, respectively; **Fig. 1c**). However, the mean PTR reported
154 from a recent study²⁷ was weakly negatively correlated with our mean TE ($r=-0.36$, $\rho=-0.41$;
155 **Fig. 1c**). These PTR measurements were highly discordant with other datasets as well, suggesting
156 that the most parsimonious explanation to be the relatively lower reliability of this PTR dataset²⁷.
157 Even stronger correlations were observed between mouse mean TE and ribosome load in mouse
158 3T3 cells²⁸ ($r=0.61$, $\rho=0.64$; **Supplementary Fig. 2b**). Together, these results suggest that our
159 TE scores are informative of protein synthesis rates in both organisms.

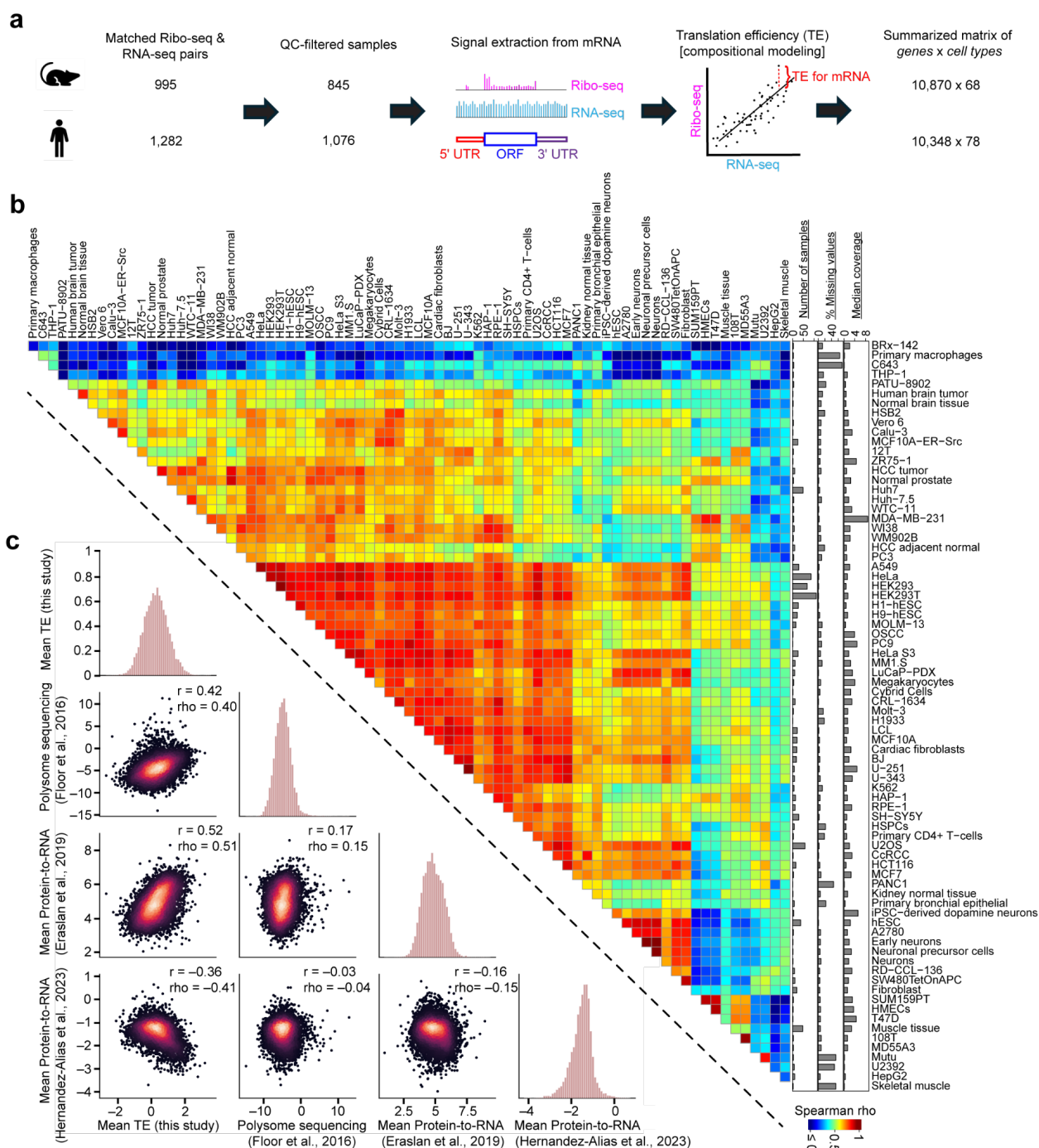


Fig. 1. Integrative analysis of thousands of human and mouse ribosomal profiling datasets measuring TE. **a)** Schematic showing the workflow of transcriptome-wide TE calculations for the human and mouse, using paired RNA-seq and ribosome profiling datasets. **b)** Heatmap of Spearman correlation coefficients comparing TEs derived from each pair of 78 human cell types. Cell types are clustered using hierarchical clustering. Right panel barplots show quality control data for the human cell type shown in each row. **c)** Comparison of mean TEs (*i.e.*, averaged across human cell types) for mRNAs derived from this study relative to alternative measurements of translational output measured in prior studies^{18,20,27}. The Pearson (r) and Spearman (ρ) correlation coefficients between each pair of measurements is also shown.

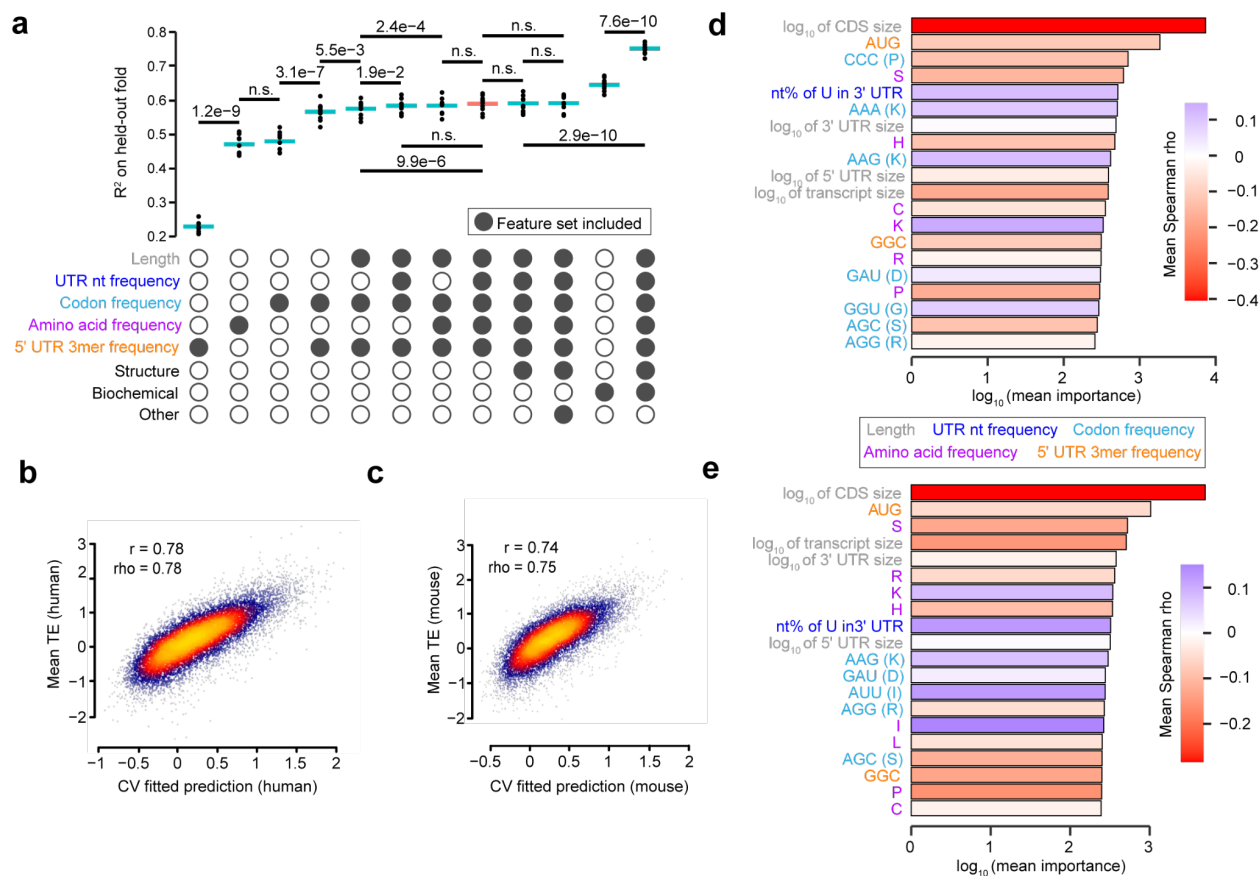
169 **Classical machine learning models to predict TE**

170 To evaluate the predictability of our TE measurements, we trained regression models on pre-
171 computed sets of sequence-encoded features derived from the mRNA. The feature sets considered
172 include: i) the lengths of the 5' UTR, CDS, 3' UTR, and entire transcript; ii) nucleotide frequencies
173 of all regions; iii) codon frequencies; iv) amino acid frequencies; v) k-mer frequencies of length 2
174 to 6 in the 5' UTR, CDS, and 3' UTR regions; vi) the frequency of each nucleotide found in the
175 wobble position; vii) the nucleotide identity at the -3, -2, -1, +4, and +5 Kozak positions; viii)
176 dicodon counts found to affect TE in yeast³⁹; and ix) multiple secondary structure features
177 (**Methods**).

178
179 To identify which feature sets usefully contributed to prediction of mean TE across all human cell
180 types, we used an iterative method that compared the cross-validated (CV) performance of a light
181 gradient-boosting machine (LGBM) model trained with a specific feature set to one trained without
182 it. If the model including the feature set performed statistically significantly better on ten held-out
183 data folds than the model without it, that feature set was deemed useful (**Methods**). The feature
184 sets found to be useful include: i) regional and total sequence lengths; ii) UTR nucleotide
185 frequencies; iii) codon frequencies; iv) amino acid frequencies; and v) the 3-mer frequencies of
186 the 5' UTR (**Fig. 2a**). All remaining feature sets did not further contribute to TE prediction (“Other”
187 in **Fig. 2a**), including secondary structure features, in contrast to prior findings⁴⁰.

188
189 Given this set of selected features, we compared three additional machine learning approaches to
190 assess their relative performance: lasso, elastic net, and random forest. We confirmed that LGBM
191 performed the best (**Supplementary Fig. 4**). We then trained LGBM models on all 78 human and
192 68 mouse cell types. The correlation between the mean TE and average over the predictions of
193 each cell type was $r=0.78$ for human and $r=0.74$ for mouse (**Fig. 2b-c**). The R^2 (averaged across
194 the held-out folds) for predicting the mean TE across cell types was 0.60 and 0.53 for the human
195 and mouse, respectively (**Supplementary Fig. 5**). Cell types with poorer data quality, such as a
196 lower fraction of detectable genes, generally led to models with inferior performance
197 (**Supplementary Fig. 5**). Although the hand-crafted feature sets could not easily include positional
198 information, the regression models were still able to achieve impressive performance.

199
200 Next, we sought to identify the relative importance of individual features for our optimal LGBM
201 model. Several of the top-ranked features were consistent with those reported in the literature (**Fig.**
202 **2d-e**). For instance, both the human and mouse models capture: i) the known negative correlation
203 between TE and both total mRNA sequence length and CDS length^{19,50-53}; ii) the importance of
204 AUG [often associated with upstream ORFs (uORFs)] and GGC trinucleotides in the 5' UTR⁵⁴⁻⁵⁶;
205 and iii) the positive correlation of A/U-richness in the third position of codons with high
206 importance for prediction accuracy. An exception to this trend was AAG (lysine), which showed
207 a positive correlation despite a G in the third position. Taken together, these results demonstrate
208 the robust predictive power of specific sequence-encoded features on mammalian TE,
209 underscoring the influence of nucleotide composition and sequence length across different cell
210 types.



211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226

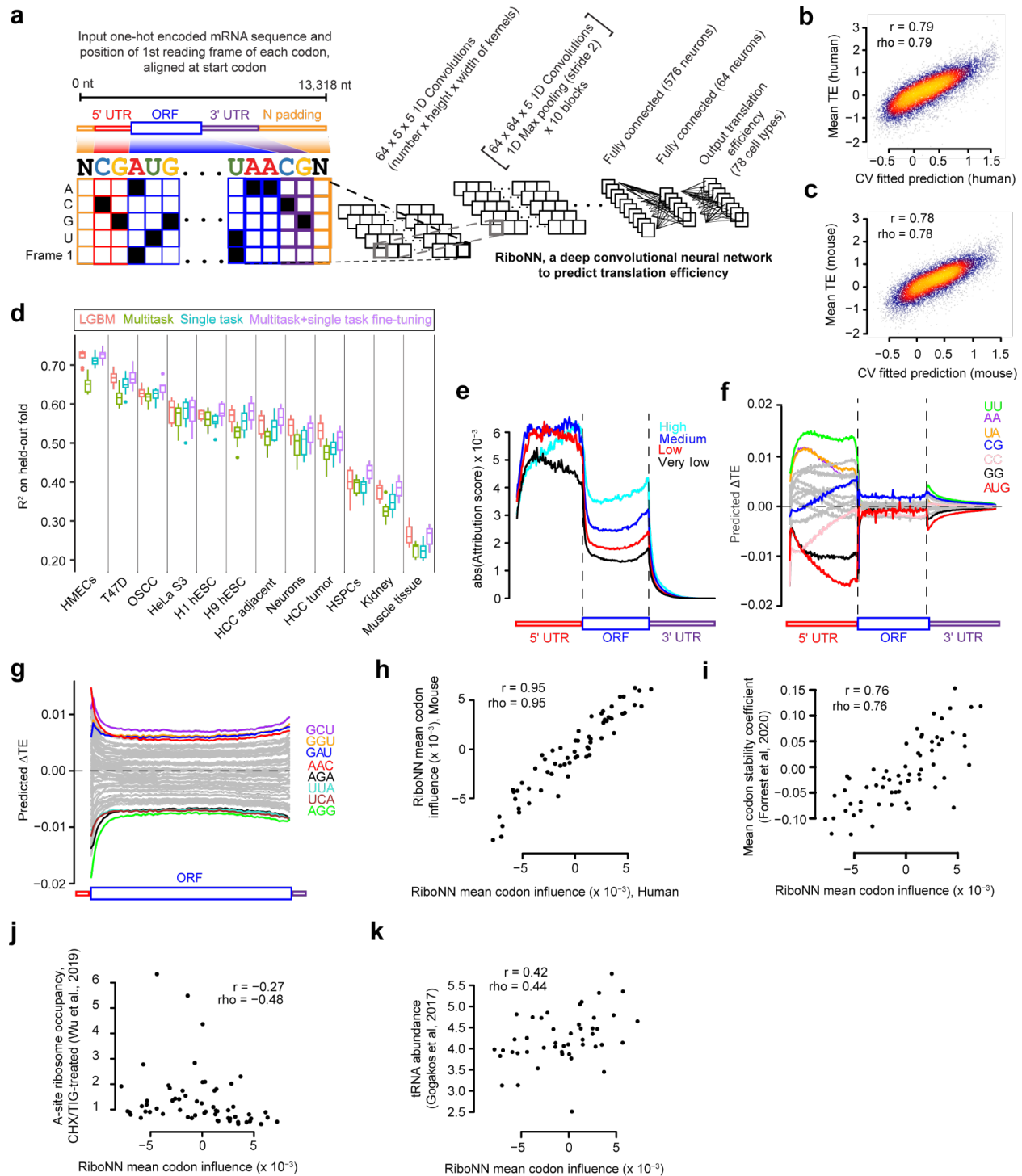
Fig. 2. A classical machine learning approach to predict mammalian TEs from mRNA sequence. a) UpSet plot showing the R^2 metric measured on ten held-out CV folds of LGBM models which predict the mean TE across human cell types using various feature sets. Colored feature sets are indicative of those that contributed to the optimal sequence-only model. Median R^2 and statistically significant differences in performance between pairs of models are indicated. P-values were calculated using one-sided, paired t-tests adjusted with a Bonferroni correction. All additional feature sets considered, but that did not have a significant improvement on performance, are labeled as “Other”. **b-c)** Importance of the features used by the optimal sequence-only model (shown as a red bar in panel **a**) for both the human (**b**) and mouse (**c**). For a given feature, importance was measured as the sum total information gain across all splits using the feature, averaged across all folds. The colors of the bars correspond to the mean Spearman rho, averaging rho values between the features and TE values from each cell type. Feature names are colored according to the feature set to which they belong. **d-e)** Scatter plots comparing the predicted and observed mean TEs, averaged across cell types, for both the human (**d**) and mouse (**e**). The Pearson (r) and Spearman (ρ) correlation coefficients, integrating the results across ten CV folds, are also shown.

227 A deep neural network to predict TE from mRNA sequence

228 Given that deep-learning-based approaches can capture positionally aware contributions of
229 sequence features and reveal degenerate motifs which are arduous to consider in classical machine
230 learning models, we compared the performance of deep-learning models on the aforementioned
231 tasks. Specifically, we trained multitask, deep convolutional neural networks to simultaneously
232 predict TEs in all cell types examined. The input to our models consisted of a one-hot encoding of
233 the mRNA sequence (up to a maximum of 13,318 nt), along with binary variables indicating the
234 first reading frame of a codon for each nucleotide; the output layer consisted of multitask
235 predictions for the TEs of either 78 human or 68 mouse cell types (**Fig. 3a**).

236 We first repurposed a hybrid convolutional and recurrent deep neural network architecture (Saluki)
237 designed to predict mRNA stability⁶, removing the splice site channel. In addition, we trained a
238 new model named RiboNN, in which we removed the gated recurrent unit layer in Saluki but
239 increased the number of convolution/max-pooling blocks from 6 to 10 to further compress mRNA
240 sequence length by ~1000-fold (**Fig. 3a, Supplementary Fig. 6**). To facilitate the learning of
241 important features (*e.g.*, Kozak sequence) near the start codon, we fixed the start codon position
242 in the input by aligning the mRNA sequences at the start codon. To accommodate the variability
243 in mRNA sequence length, both the 5' and 3' ends of mRNAs shorter than 13,318 nt were padded
244 with Ns (**Fig. 3a**). RiboNN achieved an R^2 (averaged across held-out folds) of 0.62 for predicting
245 the mean TE across the human cell types. As observed previously for LGBM models, the R^2
246 degraded for cell types with poorer data quality (**Supplementary Fig. 7**). The performance of the
247 modified Saluki and RiboNN models were similar across cell types, with RiboNN slightly
248 outperforming the modified Saluki ($p=2.9e-10$, paired Wilcoxon signed-rank test;
249 **Supplementary Fig. 7**). Moreover, deleting the codon labels or fixing the mRNA sequences at
250 the 5' end (*i.e.*, rather than the start codon) each resulted in significantly lower R^2 in most cell types
251 ($p<2.2e-16$ for both paired Wilcoxon signed-rank tests; **Supplementary Fig. 7**).

252 We independently trained RiboNN to predict TEs in 68 mouse cell types. Like the human models,
253 the mouse model exhibited variable performance among cell types, in a manner dependent on data
254 quality. Overall, RiboNN achieved an R^2 (averaged across held-out folds) of 0.60 for predicting
255 the mean TE across mouse cell types (**Supplementary Fig. 8a**). The mouse and human RiboNN
256 models worked almost as well when generating predictions across species as within species,
257 suggesting an evolutionary conservation of the principles learned (**Supplementary Fig. 8b-c**). The
258 final human and mouse models displayed correlations of 0.79 and 0.78, respectively, in predicting
259 mean TEs averaged across cell types (**Fig. 3b-c**), suggesting that RiboNN learned principles of
260 translational regulation for endogenous mRNAs.



261

262

263

264

265

266

267

268

Fig. 3. Performance and interpretation of deep learning models predicting mammalian TEs from mRNA sequence. **a)** Architecture of RiboNN, a deep multitask convolutional neural network trained to predict TEs of mRNAs in numerous cell types from an input of the mRNA sequence and an encoding of the first frame of each codon. **b-c)** Performance of RiboNN in predicting human (**b**) and mouse (**c**) mean TEs, averaged across cell types. The Pearson (r) and Spearman (ρ) correlation coefficients, integrating the results across ten CV folds, are also shown. **d)** Comparison of different model training strategies for predicting TEs in individual cell types. The following approaches were examined: LGBM trained on a single task, RiboNN

269 trained in either a multitask or single task setting, and RiboNN trained in a multitask setting but then fine-
270 tuned on a single task (*i.e.*, a “transfer learning” approach). **e**) Metagene plot summarizing the absolute value
271 of attribution scores, averaging across all mRNAs, for percentiles along the 5' UTR, CDS, and 3' UTR.
272 mRNAs were grouped into one of 4 equally sized bins according to their mean TE. **f**) Insertional analysis of
273 16 dinucleotides and the AUG motif. Motifs were inserted into each of 100 equally spaced positional bins
274 along the 5' UTR, CDS, and 3' UTRs of each mRNA. Indicated is the average predicted change in TE for
275 each bin plotted along a metagene. **g**) This panel is the same as panel **f**), except it performs analysis for 61
276 codons (excluding the 3 stop codons) inserted into the first reading frame along the length of the CDS. **h-k**)
277 Scatter plots showing the relationship between the codon influence (*i.e.*, the predicted effect size of each
278 inserted codon, averaged across all positional bins) from the human RiboNN model with that of the mouse
279 model (**h**), mean codon stability coefficients⁴⁴ (**i**), A-site ribosome occupancy scores⁵⁷ (**j**), and tRNA
280 abundances⁵⁸ (**k**). Pearson (r) and Spearman (ρ) correlation coefficients are also shown.

282 The availability of TEs measured in various cell types provided the possibility of testing multiple
283 modeling strategies to improve TE prediction for specific cell types. To further improve model
284 performance, we compared single-task models and multitask models fine-tuned to a single task
285 (*e.g.*, a transfer learning approach) on 12 randomly selected cell types exhibiting a wide
286 distribution of R^2 values (**Supplementary Table 2**). Interestingly, single-task RiboNN models
287 outperformed the multitask model for most of the cell types, but were in turn outperformed by
288 multitask models fine-tuned to a single task (**Fig. 3d**). These results highlight the power of transfer
289 learning as an effective strategy to enable information sharing between models. Although RiboNN
290 and LGBM displayed comparable prediction performance, RiboNN nevertheless has distinct
291 advantages with respect to its convenient application for transcriptome-wide TE prediction,
292 circumventing the need to pre-compute features and enabling a more computationally efficient
293 path towards the inference of genetic variant effects. Furthermore, evaluating the features that
294 contribute to RiboNN's success in predicting TE may uncover novel principles of translational
295 control that may have otherwise been overlooked.

296 To interpret the principles learned by RiboNN, we tested its predictive behavior in different
297 contexts. Saliency maps are commonly utilized to explain deep learning model predictions by
298 highlighting the input variables that contribute most towards the predicted label^{59,60}. First, for each
299 nucleotide of every human mRNA, we calculated attribution scores contributing to the prediction
300 of mean TE across all the cell types, multiplying these with the one-hot encoding of each mRNA
301 sequence to evaluate the predicted contribution of the input nucleotides. Averaging across all
302 mRNAs, we generated a metagene plot using these scores, evaluating the attributed effect size
303 (*i.e.*, absolute value) of each position along the length of each functional region of mRNA (**Fig.**
304 **3e, Supplementary Fig. 9a**). mRNAs were grouped into one of four equally sized bins according
305 to their measured mean TE (High, Medium, Low, and Very low). This analysis revealed that 5'
306 UTR sequences and CDS incorporate the greatest per-nucleotide information density (~67% and
307 31%, respectively) in predicting translational output, followed by the 3' UTR having the least
308 contribution (2%). Taking into consideration the average length of each functional region, our
309 model predicted a total global contribution of 22%, 73%, and 5% for the 5' UTR, CDS, and 3'
310 UTR, respectively. In addition, RiboNN learned position-specific contributions to TE prediction.
311 Specifically, the identity of the first 10 codons demonstrated a ~2-fold greater impact compared to
312 codons positioned towards the middle of the ORF (amino acids 70 to 80) in both human and mouse
313 (**Supplementary Fig. 9a**). These general observations were consistent for the mouse, which
314 exhibited a 67%, 31%, and 2% per-nucleotide information density and 23%, 73%, and 4% total
315 global contribution for the 5' UTR, CDS, and 3' UTR, respectively (**Supplementary Fig. 9b**). The

316 positional importance of the early coding region was similarly greater in mice (**Supplementary**
317 **Fig. 9c**), suggestive of an evolutionarily conserved principle among mammalian species.

318 We further examined our attribution scores using TF-MoDISco-lite⁶¹ to identify the most
319 significant motifs associated with TE prediction for both human and mouse RiboNN models. Our
320 analysis revealed that short, degenerate motifs; including CC, GG, CG, and AUGs upstream and
321 downstream of the main ORF; are predictive of translation output (**Supplementary Fig. 9d-e**).
322 Inspired by this finding, we performed an insertional analysis of all 16 dinucleotides and AUG to
323 evaluate the model's behavior upon inserting each of these short motifs along the full length of
324 each mRNA. We observed varying influences on TE among different motifs and across different
325 functional regions of mRNA for the same motif. Insertion of AUG and GG in the 5' UTR
326 demonstrated the strongest negative effect on TE prediction for both human and mouse models,
327 while UU, AA, and UA exhibited the strongest positive effect (**Fig. 3f, Supplementary Fig. 9f**).
328 Notably, the impact of upstream AUG (uAUG) on TE became increasingly negative as it
329 approached the start codon, whereas CG showed a progressively positive effect. Albeit smaller in
330 magnitude, most of the effects seemed to be maintained in the 3' UTR, especially for regions
331 proximal to the stop codon, suggestive of a position-dependent modulatory role for downstream
332 AUGs and other dinucleotides. Taken together, these results establish that RiboNN captures the
333 positional effects of nucleotide compositions along the entirety of the mRNA.

334 mRNAs with high TE are typically enriched for optimal codons¹⁶. To ascertain whether RiboNN
335 has also learned this property, we reiterated our insertional analysis using 61 codons (excluding
336 the 3 stop codons) inserted into the first reading frame along the length of each ORF. Similar to
337 our previous findings, the model attributed substantially different effect sizes to codons depending
338 on their position along the ORF, with the greatest predicted effects occurring near the start codon
339 (**Fig. 3g, Supplementary Fig. 9g**). GCU (alanine), GGU (glycine), GAU (aspartic acid), and AAC
340 (asparagine) exhibited the strongest positive effects on TE; conversely, AGG, AGA (arginine),
341 UCA (serine), and UUA (leucine) showed the most negative impact³⁹.

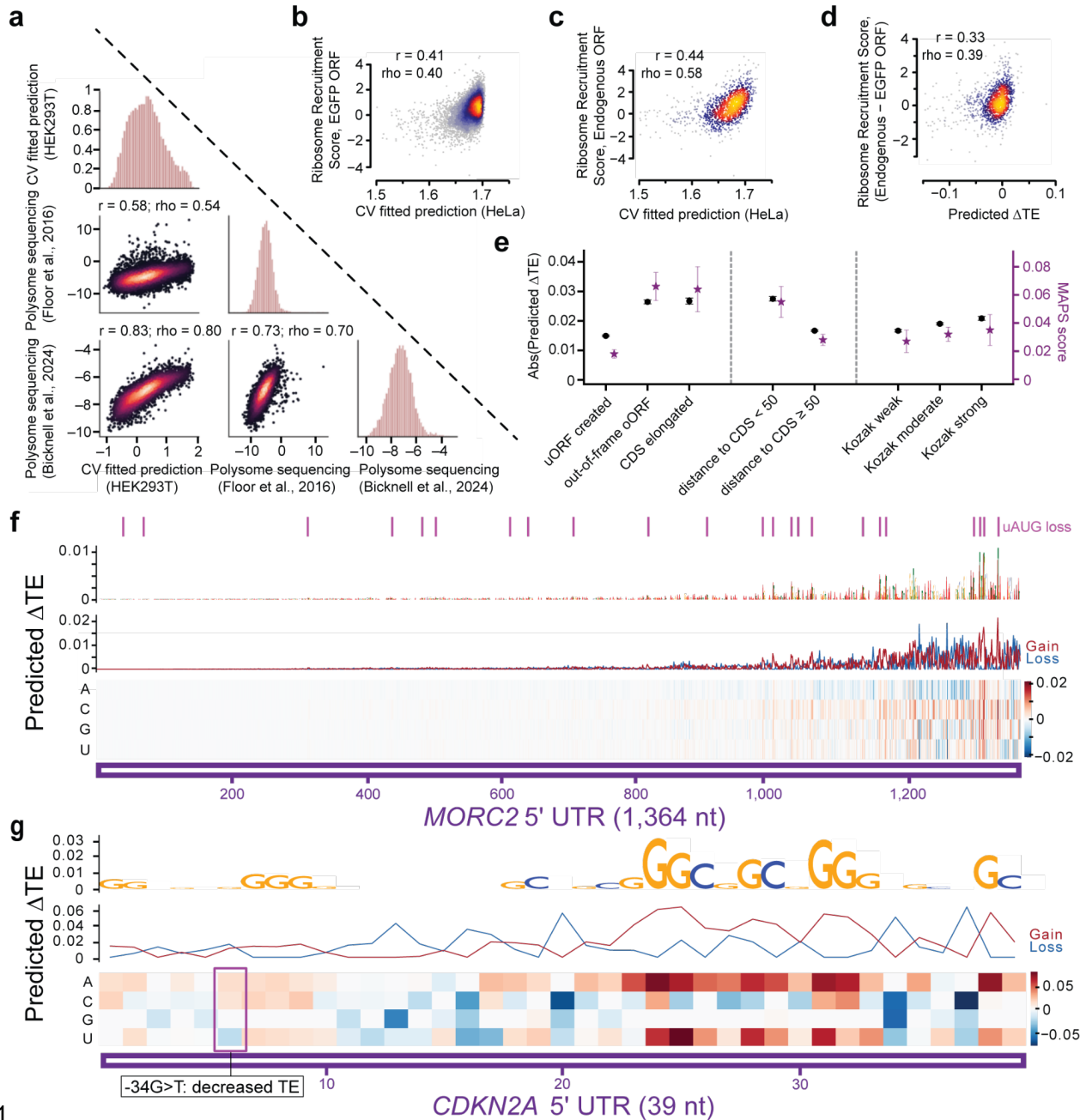
342 Based on the insertional analysis, we calculated the mean codon influence (*i.e.*, across the ORF)
343 on TE for each of the 61 non-stop codons and observed a strong correlation between the scores
344 derived from human and mouse RiboNN models ($r=0.95$, $\rho=0.95$; **Fig. 3h**), indicating
345 evolutionary conservation of predicted codon function on TE and the models' ability to learn these
346 reproducibly from completely independent datasets. Given the close link between codon usage and
347 other aspects of RNA metabolism, we compared the correlation of RiboNN-based codon influence
348 scores with several other metrics. We observed a strong positive correlation with mean codon
349 stability coefficients⁴⁴, which measure the association between codons and mRNA stability (**Fig.**
350 **3i**); a moderate negative correlation to propensity of ribosomes to have open A-sites⁵⁷, which is
351 indicative of ribosomes in the pre-accommodation state and hence slower elongation (**Fig. 3j**); and
352 a moderate positive correlation with tRNA abundance⁵⁸, which measures the availability of the
353 cognate tRNA in the cellular pool (**Fig. 3k**). The correlations persisted when the scores of codons
354 encoding the same amino acid were averaged, although no obvious trend existed with respect to
355 hydrophathy or charge of the amino acid (**Supplementary Fig. 10**). These findings underscore the
356 complex interplay of multiple mechanisms that determine the fate of mRNAs in protein
357 production.

358 Predicting translational outcomes for therapeutically delivered mRNA sequences and 359 genetic variants

360 Given RiboNN's strong performance in predicting TE for endogenous mRNAs, we assessed its
361 ability to generalize to orthogonal measures of TE and predict the impact of mRNA sequence
362 variants on TE. Mean ribosome load, measured via polysome profiling, serves as an alternative
363 metric of the translation rate of specific mRNAs, whether endogenous or therapeutic. Unlike
364 ribosome profiling, mean ribosome load can differentiate translation differences between multiple
365 RNA transcript isoforms of a given gene^{18,62}. RiboNN, which was modeled on the full length of
366 mRNAs, can be easily adapted to predict such isoform-specific TEs. The HEK293T RiboNN
367 model demonstrated $r=0.58$ and $r=0.83$ between predicted TEs and mean ribosome loads measured
368 for endogenous transcripts, which is within the realm of the reproducibility of measurement
369 between labs ($r=0.73$; **Fig. 4a**). These results indicate that our model effectively captured the
370 relationships between isoform diversity and translational regulation.

371 In addition to endogenous mRNAs, polysome profiling has been used to measure translation from
372 reporter constructs and base-modified mRNAs, as these can significantly influence protein
373 output⁶³. We next tested RiboNN's ability to predict mean ribosome load in a massively parallel
374 reporter assay dataset¹². Although RiboNN was never trained on polysome profiling or reporter
375 data, its predicted TEs were still correlated with mean ribosome load, with ρ between 0.41-0.44
376 for reporter mRNAs without modified bases and 0.30-0.31 for reporter mRNAs with either Ψ -
377 modified or N1-methylpseudouridine ($m^1\Psi$)-modified nucleotides (**Supplementary Fig. S11**).
378 Reporter assays enable assessment of how specific sequences within targeted regions affect
379 expression. We further evaluated the performance of RiboNN in predicting ribosome recruitment
380 scores for mRNAs with $m^1\Psi$ -modified 5' UTRs linked to EGFP⁵⁵, observing moderate agreement
381 ($\rho=0.40$; **Fig. 4b**). This correlation was slightly lower than that of predictions for endogenous
382 CDSs sharing the same modified 5' UTRs ($\rho=0.58$; **Fig. 4c**), indicating the broad applicability
383 of RiboNN for therapeutic mRNAs. Leveraging the paired measurement of endogenous ORF and
384 EGFP, we observed $\rho=0.39$ between changes in TE and changes in ribosome recruitment scores
385 resulting from swapping the ORFs (**Fig. 4d**). This finding underscores RiboNN's ability to
386 integrate information from both 5' UTR and ORF regions in predicting the translational regulation
387 of mRNAs.

388 Utilizing the entire mRNA sequence enables the examination of how differences in sequence,
389 including disease-associated variants, influence TE at single-nucleotide resolution. Given that 5'-
390 UTR variants that generate or disrupt uORFs can lead to disease and are key *cis*-regulators of
391 tissue-specific translation⁶⁴, we first assessed RiboNN's ability to predict the impact of uAUG-
392 associated point mutations. The RiboNN-predicted effect size had a strong association with the
393 strength of negative selection, as indicated by the mutability-adjusted proportion of singletons
394 score⁶⁴ (**Fig. 4e**). Variants creating uAUGs that result in overlapping open reading frames (oORFs)
395 or elongated CDSs exhibited a significantly higher impact on the TE of downstream protein-coding
396 genes; moreover, uAUGs generated within 50 nt of the CDS had a greater effect size than those
397 created further upstream (**Fig. 4e**). The effect size is slightly elevated if uAUG-creating variants
398 arise in the context of strong Kozak consensus sequences relative to moderate or weak ones (**Fig.**
399 **4e**). These findings reveal that RiboNN learned positional and contextual features of uAUGs, both
400 in function and evolutionary constraint.



401

402

403

404

405

406

407

408

409

410

411

412

413

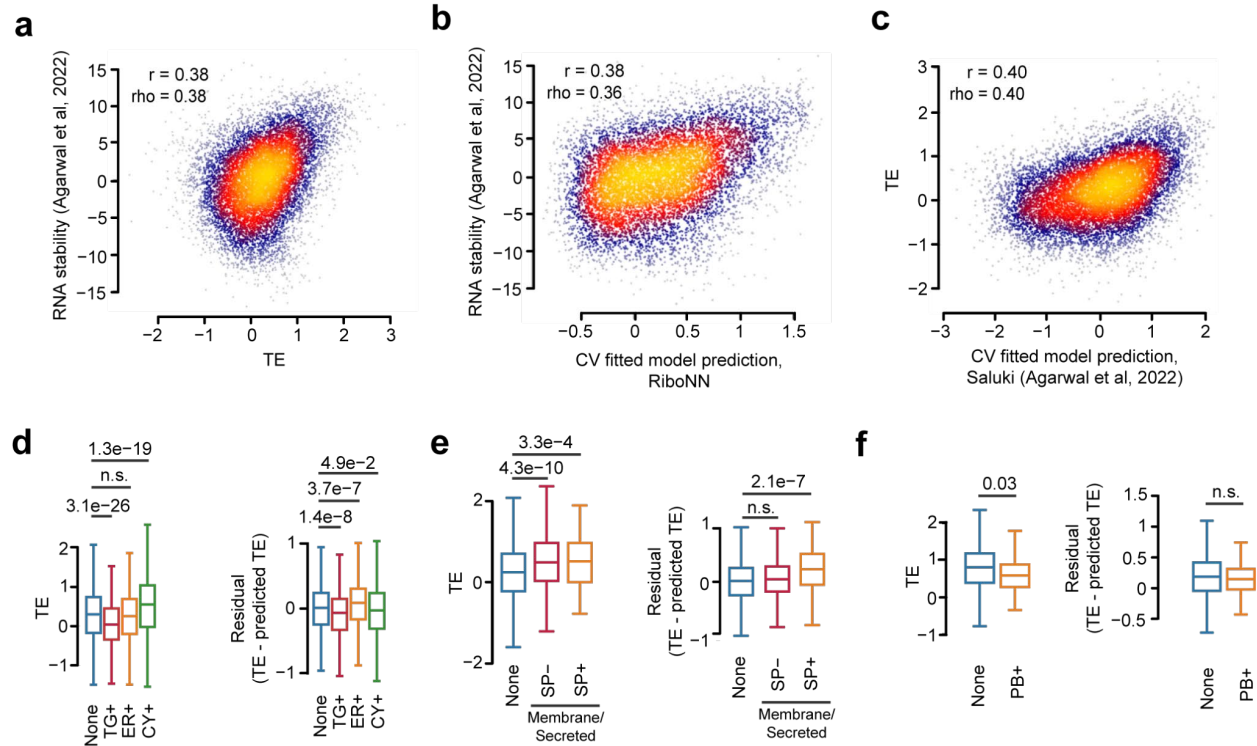
Fig. 4. RiboNN predicts the impact of RNA modifications, genetic variants, and reporter constructs on translation. **a)** Comparison of HEK293T-predicted TEs relative to mean ribosome load (MRL) as measured by polysome profiling^{18,62}. **b-d)** Performance of RiboNN in predicting the ribosomal recruitment score (*i.e.*, association of the 80S ribosomal subunit) to a panel of m1Ψ-modified 5' UTRs linked to EGFP (**b**), their corresponding endogenous ORFs (**c**), or the paired difference between the endogenous and EGFP ORF (**d**)⁵⁵. The Pearson (r) and Spearman (ρ) correlation coefficients between each pair of measurements is also shown. **e)** Relationship between the observed strength of negative selection of uAUG-associated point mutations, as measured by the mutability adjusted proportion of singletons score⁶⁴, and the RiboNN-predicted effect size. uAUG mutations were binned into categories based on the type of ORF created, distance to CDS start position, and association to Kozak consensus sequences of varying strength⁶⁴. Error bars represent confidence intervals calculated using bootstrapping⁶⁴. **f-g)** *In silico* mutagenesis results of two 5' UTR regions of *MORC2* (**f**) and *CDKN2A* (**g**). “Gain” alludes to a predicted increase in TE for the mutation, while “Loss”

414 refers to the opposite. Positions of wild type uAUG are highlighted in purple at the top. The known disease
415 associated variant is boxed. Single point mutations resulting in severe change of TE are shown alongside
416 annotations reflecting the corresponding gain or loss of TE.
417

418 Next, we conducted *in silico* mutagenesis on the 5' UTR regions of several disease-associated
419 genes. *MORC2*, a gene implicated in Charcot-Marie-Tooth disease⁶⁵, has a long 5' UTR region
420 with a large number of uAUGs. Reinforcing earlier results (**Fig. 4e**), RiboNN predicted that loss-
421 of-function mutations in CDS-proximal uAUGs would have a greater effect size relative to distal
422 uAUGs (**Fig. 4f**). For the gene *RDH12*, associated with inherited retinal disease, RiboNN
423 successfully predicted the negative impact of a uAUG-creating SNP (-123C>T), which had been
424 experimentally validated to reduce translation⁶⁶ (**Supplementary Fig. 12a**). Additionally, the gene
425 *CDKN2A* has a reported G>T mutation at base -34 in its 5' UTR that creates a uAUG reported to
426 decrease translation, leading to predisposition to melanoma⁶⁷. RiboNN consistently predicted
427 decreased TE for this variant (**Fig. 4g**). The ability of RiboNN to correctly predict the impact of
428 TE of variants extended beyond those associated with uAUGs. For example, the SNPs -127C>T
429 and -9G>A in the 5' UTR of the *ENG* gene, associated with hereditary hemorrhagic telangiectasia,
430 have been reported to reduce the expression levels of *ENG*⁶⁸, consistent with the decreased TE
431 predicted by RiboNN (**Supplementary Fig. 12b**). For *FGF13*, a gene associated with congenital
432 intellectual disability, the -32C>G mutation reduces translation⁶⁹. RiboNN also predicted a
433 negative effect of this SNP on TE, and indicated that a C>A mutation at the same position might
434 have an even greater impact on TE (**Supplementary Fig. 12c**). However, for SNP -94G>A in
435 *BCL2L13*, RiboNN predicted an increase in TE, contrary to the reported decrease in protein
436 expression⁷⁰ (**Supplementary Fig. 12d**). These results suggest that RiboNN could offer an
437 additional form of evidence to infer the regulatory impact of SNPs on disease-associated genes.

438 **RiboNN learns a common language governing mRNA stability, translational regulation, and** 439 **localization**

440 Given the strong positive correlation between the RiboNN's mean codon influence on TE and the
441 previously estimated codon influence on mRNA stability (**Fig. 3i**), we further assessed the
442 relationship between TE and mRNA stability. Indeed, both the predicted and experimentally
443 measured mean TE as well as mRNA stability from a previous study⁶ were positively correlated
444 in humans and mice ($r>0.31$, $\rho>0.32$; **Fig. 5a, Supplementary Fig. 13**). Similar patterns were
445 also observed between mRNA stability, polysome profiling, and PTR data, with the exception of
446 the PTR dataset²⁷ previously observed to be an outlier (**Supplementary Fig. 13a, Fig. 1c**).
447 Consistent with the predicted underlying role of codons influencing both TE and stability, mean
448 TE (as predicted by RiboNN) was positively correlated with mRNA stability ($r=0.38$, $\rho=0.36$;
449 **Fig. 5b**); conversely, mRNA stability (as predicted by Saluki⁶ was positively correlated with TE
450 ($r=0.40$, $\rho=0.40$; **Fig. 5c**). Taken together, these results suggest an interconnectedness between
451 mRNA stability and translational regulation that can be learned by sequence-based machine
452 learning models from diverse and independent datasets.



453

454 **Fig. 5. Interrelationships between mRNA translation, turnover, and subcellular localization.** a-c) Scatter plots showing the relationship between mean TE and mRNA stability⁶ (a), predicted mean TE and
 455 mRNA stability (b), and predicted stability and mean TE (c). Pearson (r) and Spearman (rho) correlation
 456 coefficients are also indicated. d-f) Boxplots of TE (left panel) and residual TE (*i.e.*, representing the
 457 difference between TE and the predicted TE, right panel) for mRNAs binned according to their subcellular
 458 localization. Shown are the distributions for mRNAs encoding non-membrane proteins that are enriched in
 459 TIS granules (TG+), rough endoplasmic reticulum (ER+), or cytosol (CY+)⁷¹ (d); mRNAs encoding
 460 membrane or secreted proteins, with or without predicted signal peptides (SP+/-)⁷² (e); or mRNAs enriched
 461 in cytosolic processing bodies (P-bodies)⁷³ (f). p-values were computed by comparing the behavior of
 462 mRNAs localized to the specified compartment relative to those not localized (*i.e.*, labeled “None”) using a
 463 two-sided Mann-Whitney test adjusted with a Bonferroni correction.
 464
 465

466

466 mRNAs localized to certain subcellular compartments, such as the endoplasmic reticulum (ER)
 467 membrane, tend to be differentially translated^{71,74}. We sought to evaluate these findings in the
 468 context of our predictive model, assessing both TEs and their associated residuals (mean TE –
 469 predicted mean TE) for mRNAs localizing to different compartments. For mRNAs encoding non-
 470 membrane proteins, we observed a significantly higher residual TE for ER-enriched mRNAs;
 471 additionally, cytosolically enriched mRNAs exhibited a higher TE, although this signal was largely
 472 explained by the model (Fig. 5d). When considering mRNAs encoding both non-membrane and
 473 membrane or secretory proteins, a higher TE was observed for ER-enriched mRNAs ($p < 0.01$, data
 474 not shown). This is consistent with the result that mRNAs encoding membrane or secreted proteins
 475 tended to have higher TE, even for those lacking a signal peptide sequence (Fig. 5e). Nevertheless,
 476 membrane/secreted proteins harboring a signal peptide possessed a strongly positive residual on
 477 average (Fig. 5e), indicating that RiboNN was unable to model the association between signal
 478 peptides and TE. This was unsurprising as the model was blind to amino acid sequence;

479 furthermore, it was trained on ~10K mRNA sequences and the number of sequences encoding
480 signal peptides is combinatorially explosive.

481 Given past work finding a relationship between mRNA stability and localization¹⁶, we evaluated
482 whether unexplained variation in TE from RiboNN's predictions could also be linked to mRNA
483 localization. Since less stable mRNAs tend to be translationally repressed and enriched in mRNA
484 processing bodies⁷³ (P-bodies), we expected that mRNAs enriched in P-bodies to have lower mean
485 TE compared to other mRNAs. This indeed appeared to be the case (**Fig. 5f**); however, there was
486 no difference in the residual between mRNAs enriched in P-bodies ("PB+") and others ("None"),
487 indicating that the model already learned that mRNAs enriched for localization to P-bodies was
488 associated with differential TE (**Fig. 5f**). Collectively, our results thereby establish a common
489 language governing mRNA decay, translational regulation, and subcellular localization.

490 DISCUSSION

491
492 In this study, we developed deep learning models that utilize entire mRNA sequences to predict
493 TE. These models were trained using data synthesized from thousands of ribosome profiling and
494 matched RNA-seq experiments across >140 human and mouse cell types. Our models explain over
495 70% of the variation in TE in specific cell lines, achieving a mean R^2 across cell types of 0.62.
496 This represents a 1.3 to 4.4-fold performance improvement relative to previously developed
497 models in mammals, which achieved a maximum R^2 of 0.46 (range from 0.14 to 0.46)^{14,21,22,75}.
498 Furthermore, unlike earlier efforts which were limited to a few cell types, our approach enabled
499 the development of models for a substantially larger and more diverse set of cell types.

500
501 Recent research has primarily relied on reporter constructs to dissect regulatory elements of
502 translation^{12,13,54,76}. Due largely to technological limitations, such experiments employ easily
503 detectable and fixed coding regions, such as GFP, attached to variably engineered 5' UTRs, and
504 are typically limited to one or few cell types. Critically, these reporter constructs lack the full
505 complement of proteins that normally accompany endogenous mRNAs throughout their
506 lifecycle⁷⁷, which influences RNA metabolism⁷⁸. Consequently, predictive models based on
507 reporter assays offer limited insights into the translation of endogenous mRNAs, explaining less
508 than 25% of variation in their TE^{14,22}. In contrast, our model demonstrates vastly superior
509 performance in predicting the translation of endogenous mRNAs and also appears to predict the
510 behavior of therapeutic RNAs⁵⁵.

511
512 Our predictive modeling approaches are particularly valuable as they provide a quantitative
513 assessment of factors determining TE. By analyzing the position and identity of sequence
514 elements, we were able to ascertain their relative importance in making accurate predictions. Our
515 model highlights the dominant influence of 5' UTRs and coding sequences in determining TE. The
516 nucleotide compositions of 5' UTRs heavily influenced the prediction of TE. Short, AU-rich
517 sequences were generally associated with higher TE, whereas the impact of GC-rich sequences
518 was negative but position-dependent. Intriguingly, recent massively parallel reporter assays
519 conducted in both zebrafish and human cells, utilizing different readouts to measure translation,
520 have identified a similar pattern^{54,55}. This concordance suggests that these particular regulatory
521 features observed in reporter constructs are reflective of those in endogenous transcripts.

522
523 RiboNN also learned the well-established role of uAUGs in repressing the translation of the main
524 coding sequence^{12,56,70,79}. Specifically, a shorter distance between the uAUG and the start codon
525 was associated with a reduced TE of the main coding sequence, consistent with the depletion of
526 uAUGs near CDS start sites⁷⁵. Furthermore, uAUGs closer to the start codon are more likely to
527 produce overlapping ORFs. Such overlapping ORFs, which are under more stringent selective
528 pressure in human populations⁶⁴, tend to inhibit the TE of the main CDS more than uORFs entirely
529 contained within the 5' UTR, which may allow for reinitiation following uORF translation
530 termination⁵⁶.

531
532 In addition to learning the well-established role of uAUGs, our model unexpectedly predicts that
533 downstream AUGs in 3' UTRs reduce TE, particularly when close to the stop codon. Readthrough
534 of stop codons can lead to C-terminal extensions, which decrease protein abundance⁸⁰. The
535 underlying mechanisms likely involve both proteasomal degradation^{80,81} and reduced translation

536 due to ribosome stalling^{82,83}. Alternatively, downstream AUGs can be translated due to inefficient
537 recycling of terminating ribosomes that subsequently reinitiate⁸⁴. Although the impact of such
538 events on the TE of the main ORF remains incompletely understood, a recent study suggested that
539 translation of downstream ORFs can act as translational activators⁸⁵. While our findings might
540 appear to contradict this finding, it is conceivable that there is a distance-dependent relationship,
541 where AUGs near stop codons are inhibitory due to their effects on recycling efficiency or
542 readthrough, whereas ORFs positioned further downstream could have activating effects.
543 Although our models detect specific signals in 3' UTRs, particularly near the stop codon, overall,
544 RiboNN predicts that 3' UTRs generally have a minimal impact on TE. Our results do not imply
545 that 3' UTR-dependent regulation is unimportant for specific genes⁸⁶ or particular contexts such
546 as in early vertebrate development^{87,88}. However, the overall contribution of 3' UTRs to translation
547 control is likely limited, consistent with several transcriptome-wide analyses^{28,89}.

548
549 A major finding from our study is the dominant influence of the coding sequence on TE
550 predictions. Particularly, sequences proximal to the N-termini were found to be about twice as
551 important in determining TE, a feature learned by RiboNN independently from both mouse and
552 human datasets. Interestingly, recent work using reporter constructs and single-molecule analyses
553 suggested that the identity of amino acids in early coding regions can affect protein synthesis
554 efficiency, potentially through mechanisms related to translation elongation⁴¹. While the N-
555 terminus-proximal codons were more important at a per-residue level, the identity of codons across
556 the entire CDS contributed to TE predictions. Factors such as the charge of the nascent polypeptide
557 in the exit tunnel of the ribosome^{90,91}, the pairs of codons in the decoding center^{39,92}, and
558 availability of charged tRNAs corresponding to specific codons⁹³ have all been linked to altered
559 translation elongation. Despite these mechanisms that can alter decoding rates, there is debate over
560 whether the average elongation rate across different mRNAs varies significantly^{94,95}. Critically,
561 recent studies implicate codon usage in modulating initiation efficiency through differences in
562 ribosome decoding rates^{96,97}. Given the importance of the entire CDS for the accuracy of RiboNN,
563 our results suggest that both codon and amino acid compositions are critical for determining the
564 TE of endogenous mRNAs.

565
566 Translation elongation dynamics have emerged as an important contributor to mRNA stability as
567 well^{15,16,42-46}. Intriguingly, the codon-specific effects identified by RiboNN in predicting TE
568 closely mirror their impact on mRNA stability. For instance, the codons AGA and AGG, which
569 were found to exert significant mRNA-destabilizing effects^{6,98}, also negatively impact TE, as
570 inferred by RiboNN. Additionally, during the maternal-to-zygotic transition, mRNAs enriched
571 with codons that enhance mRNA stability also show higher TE¹⁵. However, the relationship
572 between translation and mRNA decay remains debated⁹⁹, as increased TE and ribosome flux can
573 also facilitate mRNA decay, which would predict a negative correlation between the two⁶².
574 Specifically, slower elongation rates may result in mRNA degradation through either transiently
575 slowed ribosomes^{100,101} or ribosome collisions, which can activate the ribosome quality control
576 pathway¹⁰². While these mechanisms have been primarily explored using reporter constructs,
577 recent studies have also demonstrated its relevance to endogenous transcripts¹⁰³. Detailed
578 investigation into the translation-dependent and independent contributions to mRNA decay
579 remains an active area of research¹⁰⁴. Future studies are likely to uncover condition-specific effects
580 on mRNA stability that vary with TE.

581 A potential limitation of our work is that it solely considers the primary sequence to predict TE.
582 In our analyses using LGBM, the inclusion of several secondary structure-related features did not
583 enhance performance. This might be explained by several possibilities: i) the primary sequence
584 itself is highly predictive of secondary structure, potentially capturing these influences implicitly,
585 ii) prior results may have overstated the importance of RNA structure because they did not
586 appropriately account for nucleotide composition⁴⁰, or iii) the features we computed, based on
587 predicted free energy, do not accurately reflect the true secondary structures of these RNAs.
588 Considering this last point, developing more precise secondary structure features could lead to
589 further improvements in prediction accuracy.

590
591 Another avenue for improvement could involve providing RiboNN with explicit knowledge of
592 protein sequences. Including amino acid composition information improved the performance of
593 the LGBM model, and our analyses revealed systematic bias in predicted TE for proteins harboring
594 signal peptides. Thus, a deep learning model that accesses both nucleotide and amino acid
595 sequence (*i.e.*, or summarized protein-based information), may further enhance TE prediction.
596 Nevertheless, since our models currently explain 62% of the variability in mean TE across a wide
597 array of cell types, we can establish an upper bound on the impact of such features. This estimate
598 is likely conservative, as some portion of the unexplained variance in these measurements is
599 attributable to measurement error.

600
601 We would also like to note that TE, as defined in our study and typically used in the literature,
602 does not equate to the rate of protein synthesis; rather, it reflects differences in ribosome occupancy
603 relative to mRNA abundance. While recent work with reporter constructs suggested that increased
604 ribosome load may not linearly relate to protein output, both our work and previous studies^{29,105}
605 indicate that TE is positively associated with protein abundance and synthesis rates for endogenous
606 transcripts. Theoretical models of translation also support the general positive relationship between
607 protein synthesis and TE^{51,106}.

608
609 Overall, RiboNN achieves state-of-the-art prediction of TE in humans and mice, elucidating key
610 principles that underpin accurate predictions, including the relative importance of various
611 molecular aspects. These predictive models distill our knowledge into a coherent framework and
612 have the potential to advance bioengineering applications. Significantly, RiboNN has the ability
613 to generate functional predictions on genetic variants in the human population, giving insight into
614 the mechanisms constraining molecular evolution and underpinning genetic diseases. Overall,
615 these advancements have far-reaching implications for both genetic diagnostics as well as the
616 design and optimization of mRNA and gene therapies, positioning our model at the forefront of
617 these rapidly evolving domains. Looking ahead, we anticipate that future work will employ multi-
618 modal approaches to simultaneously predict all facets of gene expression—RNA abundance,
619 stability, and translation—from primary mRNA sequence, given the interconnectedness of these
620 phenomena.

621 **ACKNOWLEDGMENTS**

622
623 We thank Ian Hoskins (UT Austin) for the code and data to generate secondary structure features,
624 and Milad Miladi (Sanofi) for providing critical feedback. We thank Carson Thoreen and Wendy
625 Gilbert (Yale University) for sharing their data prior to publication. Research reported in this
626 publication was supported in part by the National Institute Of General Medical Sciences of the
627 National Institutes of Health under Award Number R35GM150667 (C.C.). This work was also
628 supported by the National Institutes of Health grant [HD110096], and the Welch Foundation grant
629 [F-2027-20230405] (C.C.). C.C. was a CPRIT Scholar in Cancer Research supported by CPRIT
630 Grant [RR180042].

631
632 **AUTHOR CONTRIBUTIONS**

633
634 D.Z. trained RiboNN models, validated model predictions with public datasets, and contributed to
635 model interpretation. J.W. interpreted RiboNN, performed comparisons between TE and third-
636 party measurements, and analyzed genetic variant data. L.P. trained and interpreted classic ML
637 models. Y.L. helped synthesize the data compendia and developed the compositional approach to
638 calculate TE. F.M., C.C., and V.A. supervised the study. C.C. and V.A. conceptualized and
639 designed the study.

640
641 **CODE AND DATA AVAILABILITY**

642
643 Code, pre-trained models, and data are planned for public release upon successful review of this
644 article.

645
646 **DECLARATION OF INTERESTS**

647
648 D.Z., J.W., F.M., and V.A. are employees of Sanofi and may hold shares and/or stock options in
649 the company.

650 **SUPPLEMENTARY TABLES**

651

652 **Supplementary Table 1. Feature sizes, sequences, CV folds, and TEs of human and mouse**
653 **genes.** The principal splice isoforms for human and mouse genes were downloaded from APPRIS
654 v2¹⁰⁷. The CV folds reported were used to split training and test sets. The TEs of transcripts with
655 an average coverage <0.1x were set to NA. The mean TEs were calculated across the cell types
656 for each transcript while ignoring NA values.

657

658 **Supplementary Table 2. Feature sizes, sequences, CV folds, and TEs predicted by the human**
659 **and mouse RiboNN models.** The principal isoforms for human and mouse genes were
660 downloaded from APPRIS v2¹⁰⁷. Predicted results are reported for the multitask and single-task
661 RiboNN models (described in **Fig. 3d**). For transcript/cell type combinations in which the TE is
662 NA in the training data, the predicted TEs were set to NA.

663 METHODS

664

665 Generation of human and mouse TE compendia

666

667 To calculate cell-type-specific TEs, we initially selected 1,282 human and 995 mouse ribosome
668 profiling datasets with matched RNA-seq data. These were screened for a series of quality control
669 steps to retain high-quality samples. Quality control criteria included ensuring average transcript
670 coverage exceeded 0.1X and reads mapping to CDS constituted more than 70% of the total. The
671 remaining 1,076 human and 835 mouse ribosome profiling samples were further processed using
672 the winsorization method to minimize the impact of PCR bias (detailed in the companion
673 manuscript¹¹⁴). Genes with sufficient counts per million (CPM > 1 in more than 70% samples) of
674 RPFs were retained, and transcripts without poly(A) tails were removed. Experimental variables,
675 such as the inclusion of elongation inhibitors, can lead to technical artifacts, manifesting as
676 increased RPF density around start and stop codons¹⁰⁸. To mitigate such biases, we only considered
677 RPFs whose 5' end mapped either after the first 10 nts or before the last 35 nts of the CDS. These
678 RPFs were summed to determine the CDS count for each transcript⁴⁷. An identical counting
679 method was used for RNA-seq data. Total CDS counts for both RNA-seq and ribosome profiling
680 were normalized using a centered log-ratio. TE was defined as the residual obtained from a
681 compositional linear regression, for each transcript in each sample (detailed in the companion
682 manuscript¹¹⁴). For each transcript, if either the RNA-seq or ribosome profiling read count was 0
683 in all samples from a specific cell line, we assigned NA to its TE in the corresponding cell line.
684 Finally, we calculated the average TE for each transcript in each cell line across all samples.

685

686 Features considered in classical machine learning models

687 The length features included the log₁₀ of the 5' UTR, CDS, 3' UTR, and total transcript lengths.
688 Nucleotide frequency included the percent composition of the 5' UTR, CDS, 3' UTR and full
689 sequence. Codon and amino acid frequencies were calculated as the percentage within the CDS,
690 and included annotated stop codons. K-mer frequencies (for k-mers of size two through six) were
691 computed separately for each region and normalized by the total k-mer count. Additional feature
692 classes included the frequency of each nucleotide in the wobble position of all codons, a one-hot
693 encoding of the nucleotide identity surrounding the start codon (at the -3, -2, -1, +4, and +5
694 positions), the counts of 20 dicodons found to affect TE in yeast³⁹, and several secondary-structure-
695 related metrics. To capture secondary structure, sequences for the 5'-most 60 nt of the transcript
696 and a 60 nt window centered on the start codon (*i.e.*, last 30 nt of the 5' UTR and first 30 nt of the
697 CDS) were extracted from the APPRIS v2 primary transcript references¹⁰⁷. If the 5'-UTR length
698 was <30 nt, the first 60 nt of the transcript were used instead. Secondary structure features were
699 enumerated in these regions using seqfold v0.7.17 ([https://github.com/Lattice-](https://github.com/Lattice-Automation/seqfold)
700 [Automation/seqfold](https://zenodo.org/records/7986470), <https://zenodo.org/records/7986470>) at a temperature of 37 °C. These features
701 were the min ΔG, number of hairpins, number of loops, number of bifurcations, number of bulges,
702 max stem length, max loop length, and position of the first stem. Hairpins with a stem length <3
703 or loop length >10 were not enumerated. The biochemical features used previously⁶ were also
704 tested separately and in combination with the sequence features.

705

706 Classical machine learning model benchmarking

707 The lasso, elastic net, random forest (scikit-learn v1.0.2)¹⁰⁹, and LGBM (lightgbm v3.2.1)¹¹⁰
708 regression models were trained using 10-fold CV. Performance was measured as the mean of the
709 R^2 values across held-out test folds. For lasso and elastic net, the training data was further split into
710 5 CV folds to find the optimal α (lasso and elastic net) and L1 ratio (elastic net) hyperparameters.
711 The default hyperparameters given were used for LGBM, with the exception of the “gain” option
712 for use with importance calculations. Random forest used the same number of trees and maximum
713 leaf nodes as LGBM. Comparisons between model types (**Supplementary Fig. 6**) and feature sets
714 (**Fig. 2a**) were deemed significant with one-sided, paired t-tests, adjusted by a Bonferroni
715 correction. We measured feature importance as the sum total information gain across all LGBM
716 tree splits using that feature, averaged across all folds. In **Fig 2b-c**, the importance was further
717 averaged over all cell lines. To determine if a feature had a positive or negative effect on prediction,
718 the Spearman correlation between the feature and cell-type-specific TE was used.

719 **RiboNN model architecture, training, and interpretation**

720 The input mRNA sequences were aligned at the start codons, with the maximum 5' UTR size set
721 to 1,381 nt and the maximum combined CDS and 3' UTR size to 11,937 nt. Sequences were padded
722 at the 5' and 3' ends with “N”, and one-hot encoded (with ‘N’ encoded by a vector of four 0s). We
723 added a fifth channel labeling the first nucleotide of each codon in the CDS⁶.

724 The architecture of RiboNN consisted of a Conv1D input layer, a “tower” of ten convolution
725 blocks, and a head of 2-linear layers (**Supplementary Fig. 6**), with each convolution block
726 including the following operations: i) layer normalization sandwiched by transpose actions, ii)
727 ReLU activation, iii) 1D convolution with kernel width 5, iv) dropout, and v) max pooling with
728 width 2. Overall, the model consisted of 250,382 learnable parameters. The output layer had one
729 or multiple neurons for single-task and multitask learning, respectively.

730 Following Saluki’s training procedure⁶, we trained the RiboNN multitask model with the MSE
731 loss function using the AdamW optimizer on batches of 64 examples, a gradually decreasing
732 learning rate between 0.001 and 0.0000001, beta1 of 0.9, and beta2 of 0.998. We clipped gradients
733 to a global norm of 0.5. We used a dropout probability of 0.3 throughout. We trained each model
734 for 200 epochs, saving checkpoints along the way. After 200 epochs, the model parameters from
735 the checkpoint with the highest validation R^2 were saved as the final model parameters. We trained
736 the mouse and human models independently using a nested CV strategy. Specifically, we trained
737 9 models for each of the 10 held-out CV folds (using 9-fold CV on the inner folds), producing a
738 total of 90 trained models. For each of the 9 models from the inner folds, we retained the top 5
739 models ranked based upon their validation R^2 performance. When running RiboNN in “prediction”
740 mode, we computed the mean of these 50 models to represent the ensemble prediction.

741 Transfer learning was implemented by replacing the linear head of our pre-trained multitask model
742 with a new single-task 2-layer linear head. We froze all preceding layers and trained the new linear
743 head for 50 epochs, followed by unfreezing all of the layers and training the entire network for
744 another 150 epochs.

745 We used the saliency method⁶⁰ within the PyTorch Captum library (version 0.6.0)¹¹¹ to compute
746 the attribution scores for each nucleotide of the input sequence with respect to the predicted mean

747 TE. For each of the test sets from our 10-fold CV procedure, we averaged the attribution scores
748 from the top 5 trained models.

749 To generate the metagene plot of attribution scores, we followed the methods established in prior
750 work⁶.

751 **Insertional motif analysis with RiboNN**

752 Using attribution scores as input, we ran TF-MoDISco-lite⁶¹ on each functional region (5' UTR,
753 ORF, and 3' UTR) independently to identify the motifs most strongly influencing the predicted
754 mean TE. Gradient correction was applied by subtracting the mean attribution score across four
755 encoding channels⁶⁰. The motifs were ranked based on the number of sequences (*i.e.*, seqlets)
756 supporting the enrichment of each motif.

757 As performed in earlier work⁶, the insertional analysis was performed by dividing each functional
758 region of a valid mRNA into 100 evenly spaced positional bins. Each k-mer examined (*i.e.*, the 16
759 dinucleotides and AUG) was inserted into one of these bins, replacing the reference sequence to
760 maintain the mRNA's original length. A valid mRNA was defined as one with a 5' UTR length
761 ≥ 100 nt, an CDS length ≥ 500 nt, and a 3' UTR length ≥ 500 nt⁶. For each insertion, the predicted
762 change in mean TE relative to the corresponding wild-type mRNA was recorded. To quantify the
763 impact of each motif across diverse sequence contexts, the predicted changes in mean TE across
764 all valid mRNAs were averaged for each of the 300 positional bins. Identical insertional analysis
765 was performed for the 61 non-stop codons, except that each codon was inserted into the first
766 reading frame of the ORF.

767 **Impact of uAUG-creating variants with RiboNN**

768 As described in an earlier study⁶⁴, we retrieved the list of variants that create uAUGs and selected
769 the canonical transcript based on the gnomAD v2 annotation¹¹² for each gene for further analysis.
770 For each uAUG-creating variant considered, we verified that its gene name matched the list of
771 canonical transcripts and that the distance from each uAUG variant to the start of its CDS was
772 accurately annotated. This led to a set of 15,184 uAUG variants which were categorized into two
773 groups based on their effects and contexts as previously annotated⁶⁴. The effect group was
774 comprised of variants that create out-of-frame oORFs (n=2,784), elongate the CDSs (n=1,350), or
775 generate uORFs (n=9,263). The context group included variants located at a distance of ≥ 50 nt
776 from the CDS (n=11,113), < 50 nt from the CDS (n=2,284), or associated with a strong (n=2,237),
777 moderate (n=6,559), or weak (n=4,601) Kozak consensus sequence. To assess the impact of each
778 variant on TE, we recorded the change in predicted TE relative to the wild-type mRNA reference
779 sequence. The confidence intervals were calculated using bootstrapping as described⁶⁴.

780 ***In silico* mutagenesis analysis of disease genes with RiboNN**

781 We performed *in silico* mutagenesis analysis⁶ on the 5' UTR regions of genes associated with
782 various diseases to predict the impact of genetic variants on TE. For each nucleotide position, we
783 substituted the reference nucleotide with each of the three possible alternative alleles, and
784 computed the predicted Δ TE.

785 **Subcellular localization analysis**

786 Based on prior results¹¹¹, we categorized 5,884 non-membrane protein-encoding mRNAs as
787 enriched in TIS granules (TG+, n=1,086), the rough ER (ER+, n=745), the cytosol (CY+,
788 n=1,299), or exhibiting no apparent localization (2,754). For our analysis of P-body-enriched
789 mRNAs, we examined a total of 1,636 mRNAs¹¹³, of which 93 exhibited P-body enrichment based
790 on prior results¹¹³. P-values from Mann-Whitney-Wilcoxon test two-sided with Bonferroni
791 correction were performed to show statistical significance.

792 REFERENCES CITED

- 793 1. Agarwal, V. & Shendure, J. Predicting mRNA Abundance Directly from Genomic Sequence Using Deep
794 Convolutional Neural Networks. *Cell Rep.* **31**, 107663 (2020).
- 795 2. Zhou, J. *et al.* Deep learning sequence-based ab initio prediction of variant effects on expression and disease risk.
796 *Nat. Genet.* **50**, 1171–1179 (2018).
- 797 3. Avsec, Ž. *et al.* Effective gene expression prediction from sequence by integrating long-range interactions. *Nat.*
798 *Methods* **18**, 1196–1203 (2021).
- 799 4. Linder, J., Srivastava, D., Yuan, H., Agarwal, V. & Kelley, D. R. Predicting RNA-seq coverage from DNA
800 sequence as a unifying model of gene regulation. *bioRxiv* 2023.08.30.555582 (2023)
801 doi:10.1101/2023.08.30.555582.
- 802 5. Kelley, D. R. *et al.* Sequential regulatory activity prediction across chromosomes with convolutional neural
803 networks. *Genome Res.* **28**, 739–750 (2018).
- 804 6. Agarwal, V. & Kelley, D. R. The genetic and biochemical determinants of mRNA degradation rates in mammals.
805 *Genome Biol.* **23**, 245 (2022).
- 806 7. Gingold, H. & Pilpel, Y. Determinants of translation efficiency and accuracy. *Mol. Syst. Biol.* **7**, 481 (2011).
- 807 8. Zur, H. & Tuller, T. Predictive biophysical modeling and understanding of the dynamics of mRNA translation
808 and its evolution. *Nucleic Acids Res.* **44**, 9031–9049 (2016).
- 809 9. Nieuwkoop, T. *et al.* Revealing determinants of translation efficiency via whole-gene codon randomization and
810 machine learning. *Nucleic Acids Res.* **51**, 2363–2376 (2023).
- 811 10. Shao, B. *et al.* Riboformer: a deep learning framework for predicting context-dependent translation dynamics.
812 *Nat. Commun.* **15**, 2011 (2024).
- 813 11. Tian, T., Li, S., Lang, P., Zhao, D. & Zeng, J. Full-length ribosome density prediction by a multi-input and multi-
814 output model. *PLoS Comput. Biol.* **17**, e1008842 (2021).
- 815 12. Sample, P. J. *et al.* Human 5' UTR design and variant effect prediction from a massively parallel translation
816 assay. *Nat. Biotechnol.* **37**, 803–809 (2019).
- 817 13. Cao, J. *et al.* High-throughput 5' UTR engineering for enhanced protein production in non-viral gene therapies.
818 *Nat. Commun.* **12**, 4138 (2021).
- 819 14. Karollus, A., Avsec, Ž. & Gagneur, J. Predicting mean ribosome load for 5'UTR of any length using deep
820 learning. *PLoS Comput. Biol.* **17**, e1008982 (2021).
- 821 15. Bazzini, A. A. *et al.* Codon identity regulates mRNA stability and translation efficiency during the maternal-to-
822 zygotic transition. *EMBO J.* **35**, 2087–2103 (2016).
- 823 16. Hanson, G. & Collier, J. Codon optimality, bias and usage in translation and mRNA decay. *Nat. Rev. Mol. Cell*
824 *Biol.* **19**, 20–30 (2018).
- 825 17. Szostak, E. & Gebauer, F. Translational control by 3'-UTR-binding proteins. *Brief. Funct. Genomics* **12**, 58–65
826 (2013).
- 827 18. Floor, S. N. & Doudna, J. A. Tunable protein synthesis by transcript isoforms in human cells. *Elife* **5**, (2016).
- 828 19. Vogel, C. *et al.* Sequence signatures and mRNA concentration can explain two-thirds of protein abundance
829 variation in a human cell line. *Mol. Syst. Biol.* **6**, 400 (2010).
- 830 20. Eraslan, B. *et al.* Quantification and discovery of sequence determinants of protein-per-mRNA amount in 29
831 human tissues. *Mol. Syst. Biol.* **15**, e8513 (2019).
- 832 21. Eisen, T. J., Li, J. J. & Bartel, D. P. The interplay between translational efficiency, poly(A) tails, microRNAs,
833 and neuronal activation. *RNA* **28**, 808–831 (2022).
- 834 22. Li, J. J., Chew, G.-L. & Biggin, M. D. Quantitative principles of cis-translational control by general mRNA
835 sequence features in eukaryotes. *Genome Biol.* **20**, 162 (2019).
- 836 23. Battle, A. *et al.* Genomic variation. Impact of regulatory variation from RNA to protein. *Science* **347**, 664–667
837 (2015).
- 838 24. Cenik, C. *et al.* Integrative analysis of RNA, translation, and protein levels reveals distinct regulatory variation
839 across humans. *Genome Res.* **25**, 1610–1621 (2015).
- 840 25. Schwanhäusser, B. *et al.* Global quantification of mammalian gene expression control. *Nature* **473**, 337–342
841 (2011).
- 842 26. Jovanovic, M. *et al.* Immunogenetics. Dynamic profiling of the protein life cycle in response to pathogens.
843 *Science* **347**, 1259038 (2015).
- 844 27. Hernandez-Alias, X., Benisty, H., Radusky, L. G., Serrano, L. & Schaefer, M. H. Using protein-per-mRNA
845 differences among human tissues in codon optimization. *Genome Biol.* **24**, 34 (2023).
- 846 28. Spies, N., Burge, C. B. & Bartel, D. P. 3' UTR-isoform choice has limited influence on the stability and

- 847 translational efficiency of most mRNAs in mouse fibroblasts. *Genome Research* vol. 23 2078–2090 Preprint at
848 <https://doi.org/10.1101/gr.156919.113> (2013).
- 849 29. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide analysis in vivo of
850 translation with nucleotide resolution using ribosome profiling. *Science* **324**, 218–223 (2009).
- 851 30. Li, J. J., Bickel, P. J. & Biggin, M. D. System wide analyses have underestimated protein abundances and the
852 importance of transcription in mammals. *PeerJ* **2**, e270 (2014).
- 853 31. Gorgoni, B., Marshall, E., McFarland, M. R., Romano, M. C. & Stansfield, I. Controlling translation elongation
854 efficiency: tRNA regulation of ribosome flux on the mRNA. *Biochem. Soc. Trans.* **42**, 160–165 (2014).
- 855 32. Sonenberg, N. & Hinnebusch, A. G. Regulation of translation initiation in eukaryotes: mechanisms and biological
856 targets. *Cell* **136**, 731–745 (2009).
- 857 33. Jackson, R. J., Hellen, C. U. T. & Pestova, T. V. The mechanism of eukaryotic translation initiation and principles
858 of its regulation. *Nat. Rev. Mol. Cell Biol.* **11**, 113–127 (2010).
- 859 34. Hinnebusch, A. G., Ivanov, I. P. & Sonenberg, N. Translational control by 5'-untranslated regions of eukaryotic
860 mRNAs. *Science* **352**, 1413–1416 (2016).
- 861 35. Sharp, P. M. & Li, W. H. An evolutionary perspective on synonymous codon usage in unicellular organisms. *J.*
862 *Mol. Evol.* **24**, 28–38 (1986).
- 863 36. Presnyak, V. *et al.* Codon optimality is a major determinant of mRNA stability. *Cell* **160**, 1111–1124 (2015).
- 864 37. Torrent, M., Chalancon, G., de Groot, N. S., Wuster, A. & Madan Babu, M. Cells alter their tRNA abundance to
865 selectively regulate protein synthesis during stress conditions. *Sci. Signal.* **11**, (2018).
- 866 38. Weinberg, D. E. *et al.* Improved Ribosome-Footprint and mRNA Measurements Provide Insights into Dynamics
867 and Regulation of Yeast Translation. *Cell Rep.* **14**, 1787–1799 (2016).
- 868 39. Gamble, C. E., Brule, C. E., Dean, K. M., Fields, S. & Grayhack, E. J. Adjacent Codons Act in Concert to
869 Modulate Translation Efficiency in Yeast. *Cell* **166**, 679–690 (2016).
- 870 40. Mauger, D. M. *et al.* mRNA structure regulates protein expression through changes in functional half-life. *Proc.*
871 *Natl. Acad. Sci. U. S. A.* **116**, 24075–24083 (2019).
- 872 41. Verma, M. *et al.* A short translational ramp determines the efficiency of protein synthesis. *Nat. Commun.* **10**,
873 5774 (2019).
- 874 42. Burke, P. C., Park, H. & Subramaniam, A. R. A nascent peptide code for translational control of mRNA stability
875 in human cells. *Nat. Commun.* **13**, 6829 (2022).
- 876 43. Narula, A., Ellis, J., Talianferro, J. M. & Rissland, O. S. Coding regions affect mRNA stability in human cells.
877 *RNA* **25**, 1751–1764 (2019).
- 878 44. Forrest, M. E. *et al.* Codon and amino acid content are associated with mRNA stability in mammalian cells. *PLoS*
879 *One* **15**, e0228730 (2020).
- 880 45. Wu, Q. *et al.* Translation affects mRNA stability in a codon-dependent manner in human cells. *Elife* **8**, (2019).
- 881 46. Hia, F. *et al.* Codon bias confers stability to human mRNAs. *EMBO Rep.* **20**, e48220 (2019).
- 882 47. Ozadam, H., Geng, M. & Cenik, C. RiboFlow, RiboR and RiboPy: an ecosystem for analyzing ribosome profiling
883 data at read length resolution. *Bioinformatics* **36**, 2929–2931 (2020).
- 884 48. Larsson, O., Sonenberg, N. & Nadon, R. Identification of differential translation in genome wide studies. *Proc.*
885 *Natl. Acad. Sci. U. S. A.* **107**, 21487–21492 (2010).
- 886 49. Wang, D. *et al.* A deep proteome and transcriptome abundance atlas of 29 healthy human tissues. *Mol. Syst. Biol.*
887 **15**, e8503 (2019).
- 888 50. Rogers, D. W., Böttcher, M. A., Traulsen, A. & Greig, D. Ribosome reinitiation can explain length-dependent
889 translation of messenger RNA. *PLoS Comput. Biol.* **13**, e1005592 (2017).
- 890 51. Fernandes, L. D., Moura, A. P. S. de & Ciandrini, L. Gene length as a regulator for ribosome recruitment and
891 protein synthesis: theoretical insights. *Sci. Rep.* **7**, 17409 (2017).
- 892 52. Witte, F. *et al.* A trans locus causes a ribosomopathy in hypertrophic hearts that affects mRNA translation in a
893 protein length-dependent fashion. *Genome Biol.* **22**, 1–34 (2021).
- 894 53. Thompson, M. K., Rojas-Duran, M. F., Gangaramani, P. & Gilbert, W. V. The ribosomal protein Asc1/RACK1
895 is required for efficient translation of short mRNAs. *Elife* **5**, (2016).
- 896 54. Strayer, E. C. *et al.* NaP-TRAP, a novel massively parallel reporter assay to quantify translation control. *bioRxiv*
897 2023.11.09.566434 (2023) doi:10.1101/2023.11.09.566434.
- 898 55. Lewis, C. J. T. *et al.* Quantitative profiling of human translation initiation reveals regulatory elements that
899 potentially affect endogenous and therapeutically modified mRNAs. *bioRxiv* (2024)
900 doi:10.1101/2024.02.28.582532.
- 901 56. Dever, T. E., Ivanov, I. P. & Hinnebusch, A. G. Translational regulation by uORFs and start codon selection
902 stringency. *Genes Dev.* **37**, 474–489 (2023).

- 903 57. Wu, C. C.-C., Zinshteyn, B., Wehner, K. A. & Green, R. High-Resolution Ribosome Profiling Defines Discrete
904 Ribosome Elongation States and Translational Regulation during Cellular Stress. *Mol. Cell* **73**, 959–970.e5
905 (2019).
- 906 58. Gogakos, T. *et al.* Characterizing Expression and Processing of Precursor and Mature Human tRNAs by Hydro-
907 tRNAseq and PAR-CLIP. *Cell Rep.* **20**, 1463–1475 (2017).
- 908 59. Alqaraawi, A., Schuessler, M., Weiß, P., Costanza, E. & Berthouze, N. Evaluating Saliency Map Explanations
909 for Convolutional Neural Networks: A User Study. *arXiv [cs.HC]* (2020).
- 910 60. Simonyan, K., Vedaldi, A. & Zisserman, A. Deep Inside Convolutional Networks: Visualising Image
911 Classification Models and Saliency Maps. *arXiv [cs.CV]* (2013).
- 912 61. Shrikumar, A. *et al.* Technical Note on Transcription Factor Motif Discovery from Importance Scores (TF-
913 MoDISco) version 0.5.6.5. *arXiv [cs.LG]* (2018).
- 914 62. Bicknell, A. A. *et al.* Attenuating ribosome load improves protein output from mRNA by limiting translation-
915 dependent mRNA decay. *Cell Rep.* **43**, 114098 (2024).
- 916 63. Nachtergaele, S. & He, C. Chemical Modifications in the Life of an mRNA Transcript. *Annu. Rev. Genet.* **52**,
917 349–372 (2018).
- 918 64. Whiffin, N. *et al.* Characterising the loss-of-function impact of 5' untranslated region variants in 15,708
919 individuals. *Nat. Commun.* **11**, 1–12 (2020).
- 920 65. Sevilla, T. *et al.* Mutations in the MORC2 gene cause axonal Charcot–Marie–Tooth disease. *Brain* **139**, 62–72
921 (2015).
- 922 66. Dueñas Rey, A. *et al.* Combining a prioritization strategy and functional studies nominates 5'UTR variants
923 underlying inherited retinal disease. *Genome Med.* **16**, 7 (2024).
- 924 67. Liu, L. *et al.* Mutation of the CDKN2A 5' UTR creates an aberrant initiation codon and predisposes to melanoma.
925 *Nat. Genet.* **21**, 128–132 (1999).
- 926 68. Damjanovich, K. *et al.* 5'UTR mutations of ENG cause hereditary hemorrhagic telangiectasia. *Orphanet J. Rare*
927 *Dis.* **6**, 85 (2011).
- 928 69. Pan, X. *et al.* 5'-UTR SNP of FGF13 causes translational defect and intellectual disability. *Elife* **10**, (2021).
- 929 70. Lee, D. S. M. *et al.* Disrupting upstream translation in mRNAs is associated with human disease. *Nat. Commun.*
930 **12**, 1515 (2021).
- 931 71. Horste, E. L. *et al.* Subcytoplasmic location of translation controls protein output. *Mol. Cell* **83**, 4509–4523.e11
932 (2023).
- 933 72. Acids research, N. & 2021. UniProt: the universal protein knowledgebase in 2021. *academic.oup.com* (2021).
- 934 73. Hubstenberger, A. *et al.* P-Body Purification Reveals the Condensation of Repressed mRNA Regulons. *Mol. Cell*
935 **68**, 144–157.e5 (2017).
- 936 74. Stephens, S. B. & Nicchitta, C. V. Divergent regulation of protein synthesis in the cytosol and endoplasmic
937 reticulum compartments of mammalian cells. *Mol. Biol. Cell* **19**, 623–632 (2008).
- 938 75. Chew, G.-L., Pauli, A. & Schier, A. F. Conservation of uORF repressiveness and sequence features in mouse,
939 human and zebrafish. *Nat. Commun.* **7**, 11663 (2016).
- 940 76. Jia, L. *et al.* Decoding mRNA translatability and stability from the 5' UTR. *Nat. Struct. Mol. Biol.* **27**, 814–821
941 (2020).
- 942 77. Choi, Y. *et al.* Time-resolved profiling of RNA binding proteins throughout the mRNA life cycle. *Mol. Cell* **84**,
943 1764–1782.e10 (2024).
- 944 78. Singh, G., Pratt, G., Yeo, G. W. & Moore, M. J. The Clothes Make the mRNA: Past and Present Trends in mRNP
945 Fashion. *Annu. Rev. Biochem.* **84**, 325–354 (2015).
- 946 79. May, G. E. *et al.* Unraveling the influences of sequence and position on yeast uORF activity using massively
947 parallel reporter systems and machine learning. *Elife* **12**, (2023).
- 948 80. Arribere, J. A. *et al.* Translation readthrough mitigation. *Nature* (2016) doi:10.1038/nature18308.
- 949 81. Kramarski, L. & Arbely, E. Translational read-through promotes aggregation and shapes stop codon identity.
950 *Nucleic Acids Res.* **48**, 3747–3760 (2020).
- 951 82. Yordanova, M. M. *et al.* AMD1 mRNA employs ribosome stalling as a mechanism for molecular memory
952 formation. *Nature* **553**, 356–360 (2018).
- 953 83. Hashimoto, S., Nobuta, R., Izawa, T. & Inada, T. Translation arrest as a protein quality control system for aberrant
954 translation of the 3'-UTR in mammalian cells. *FEBS Lett.* **593**, 777–787 (2019).
- 955 84. Sherlock, M. E., Baquero Galvis, L., Vicens, Q., Kieft, J. S. & Jagannathan, S. Principles, mechanisms, and
956 biological implications of translation termination-reinitiation. *RNA* **29**, 865–884 (2023).
- 957 85. Wu, Q. *et al.* Translation of small downstream ORFs enhances translation of canonical main open reading frames.
958 *EMBO J.* **39**, e104763 (2020).

- 959 86. Mayr, C. Evolution and Biological Roles of Alternative 3'UTRs. *Trends Cell Biol.* **26**, 227–237 (2016).
960 87. Subtelny, A. O., Eichhorn, S. W., Chen, G. R., Sive, H. & Bartel, D. P. Poly(A)-tail profiling reveals an
961 embryonic switch in translational control. *Nature* **508**, 66–71 (2014).
962 88. Ozadam, H. *et al.* Single-cell quantification of ribosome occupancy in early mouse development. *Nature* **618**,
963 1057–1064 (2023).
964 89. Gruber, A. R. *et al.* Global 3' UTR shortening has a limited effect on protein abundance in proliferating T cells.
965 *Nat. Commun.* **5**, (2014).
966 90. Requião, R. D., Barros, G. C., Domitrovic, T. & Palhano, F. L. Influence of nascent polypeptide positive charges
967 on translation dynamics. *Biochem. J* **477**, 2921–2934 (2020).
968 91. Dao Duc, K. & Song, Y. S. The impact of ribosomal interference, codon usage, and exit tunnel interactions on
969 translation elongation rate variation. *PLoS Genet.* **14**, e1007166 (2018).
970 92. Ahmed, N. *et al.* Pairs of amino acids at the P- and A-sites of the ribosome predictably and causally modulate
971 translation-elongation rates. *J. Mol. Biol.* **432**, 166696 (2020).
972 93. Kirchner, S. & Ignatova, Z. Emerging roles of tRNA in adaptive translation, signalling dynamics and disease.
973 *Nat. Rev. Genet.* **16**, 98–112 (2015).
974 94. Ingolia, N. T., Lareau, L. F. & Weissman, J. S. Ribosome profiling of mouse embryonic stem cells reveals the
975 complexity and dynamics of mammalian proteomes. *Cell* **147**, 789–802 (2011).
976 95. Riba, A. *et al.* Protein synthesis rates and ribosome occupancies reveal determinants of translation elongation
977 rates. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 15023–15032 (2019).
978 96. Barrington, C. L. *et al.* Synonymous codon usage regulates translation initiation. *Cell Rep.* **42**, 113413 (2023).
979 97. Lyons, E. F. *et al.* Codon optimality modulates protein output by tuning translation initiation. *bioRxiv* (2023)
980 doi:10.1101/2023.11.27.568910.
981 98. Chen, K. Y., Park, H. & Subramaniam, A. R. Massively parallel identification of sequence motifs triggering
982 ribosome-associated mRNA quality control. *Nucleic Acids Res.* **52**, 7171–7187 (2024).
983 99. Bicknell, A. A. & Ricci, E. P. When mRNA translation meets decay. *Biochem. Soc. Trans.* **45**, 339–351 (2017).
984 100. Mishima, Y., Han, P., Ishibashi, K., Kimura, S. & Iwasaki, S. Ribosome slowdown triggers codon-mediated
985 mRNA decay independently of ribosome quality control. *EMBO J.* **41**, e109256 (2022).
986 101. Bae, H. & Collier, J. Codon optimality-mediated mRNA degradation: Linking translational elongation to mRNA
987 stability. *Mol. Cell* **82**, 1467–1476 (2022).
988 102. Inada, T. Quality controls induced by aberrant translation. *Nucleic Acids Res.* **48**, 1084–1096 (2020).
989 103. Matsuo, Y. *et al.* RQT complex dissociates ribosomes collided on endogenous RQC substrate SDD1. *Nat. Struct.*
990 *Mol. Biol.* **27**, 323–332 (2020).
991 104. Mercier, B. C. *et al.* Translation-dependent and -independent mRNA decay occur through mutually exclusive
992 pathways defined by ribosome density during T cell activation. *Genome Res.* **34**, 394–409 (2024).
993 105. Liu, T.-Y. *et al.* Time-Resolved Proteomics Extends Ribosome Profiling-Based Measurements of Protein
994 Synthesis Dynamics. *Cell Syst* **4**, 636–644.e9 (2017).
995 106. Shah, P., Ding, Y., Niemczyk, M., Kudla, G. & Plotkin, J. B. Rate-limiting steps in yeast protein translation. *Cell*
996 **153**, 1589–1601 (2013).
997 107. Rodriguez, J. M. *et al.* APPRIS: selecting functionally important isoforms. *Nucleic Acids Res.* **50**, D54–D59
998 (2022).
999 108. Gerashchenko, M. V. & Gladyshev, V. N. Translation inhibitors cause abnormalities in ribosome profiling
1000 experiments. *Nucleic Acids Res.* **42**, e134 (2014).
1001 109. Pedregosa, F. *et al.* Scikit-learn: Machine Learning in Python. *J. Mach. Learn. Res.* **abs/1201.0490**, (2011).
1002 110. Ke, G. *et al.* LightGBM: A highly efficient Gradient Boosting Decision Tree. *Adv. Neural Inf. Process. Syst.*
1003 3146–3154 (2017).
1004 111. Kokhlikyan, N. *et al.* Captum: A unified and generic model interpretability library for PyTorch. *arXiv [cs.LG]*
1005 (2020).
1006 112. Gudmundsson, S. *et al.* Addendum: The mutational constraint spectrum quantified from variation in 141,456
1007 humans. *Nature* **597**, E3–E4 (2021).
1008 113. Majdandzic, A., Rajesh, C. & Koo, P. K. Correcting gradient-based interpretations of deep neural networks for
1009 genomics. *Genome Biol.* **24**, 109 (2023).
1010 114. Liu Y, Hoskins I, Geng M, Zhao Q, Chacko J, Persyn L, Wang J, Zheng D, Zhong Y, Rao S, Park D, Cenik
1011 ES, Agarwal V, Ozadam H, Cenik C. Translation efficiency covariation across cell types is a conserved
1012 organizing principle of mammalian transcriptomes. Companion manuscript. (2024).