

15 **ABSTRACT**

16

17 Characterization of shared patterns of RNA expression between genes across conditions has led to
18 the discovery of regulatory networks and novel biological functions. However, it is unclear if such
19 coordination extends to translation, a critical step in gene expression. Here, we uniformly analyzed
20 3,819 ribosome profiling datasets from 117 human and 94 mouse tissues and cell lines. We
21 introduce the concept of *Translation Efficiency Covariation* (TEC), identifying coordinated
22 translation patterns across cell types. We nominate potential mechanisms driving shared patterns
23 of translation regulation. TEC is conserved across human and mouse cells and helps uncover gene
24 functions. Moreover, our observations indicate that proteins that physically interact are highly
25 enriched for positive covariation at both translational and transcriptional levels. Our findings
26 establish translational covariation as a conserved organizing principle of mammalian
27 transcriptomes.

28

29 **Keywords:** Translation efficiency, Translational regulation, Ribosome profiling, Translation
30 efficiency covariation, Gene functions, Regulatory networks

31

32 INTRODUCTION

33 In the last three decades, technological advances have progressively revealed the expression of
34 RNAs with increasing spatial and cellular resolution¹⁻⁷. These measurements have spurred
35 conceptual advances, driven by computational approaches. Foremost among these is the concept
36 of RNA co-expression, which quantifies the similarities in RNA expression changes among groups
37 of genes across conditions⁸⁻¹¹.

38 RNA co-expression analysis across biological contexts reveals shared biological functions,
39 informing us about underlying mechanisms and interactions¹²⁻¹⁵. By applying the principle of
40 guilt-by-association, new functions for genes with previously unknown roles can be inferred by
41 the similarity of RNA expression patterns with genes of known function^{9,12,16}. Furthermore, RNA
42 co-expression between transcripts is predictive of protein-protein interactions^{17,18}, and can indicate
43 genes that are likely regulated by the same transcription factors, suggesting common regulatory
44 mechanisms^{19,20}.

45 These findings suggest that RNA co-expression may serve as a proxy for the proteomic
46 organization of cells. However, it is only recently that quantification of protein abundance across
47 numerous cell types and conditions has become possible, allowing this assumption to be explicitly
48 tested. Mass spectrometry-based measurements across hundreds of cell types have revealed that
49 proteins similarly exhibit shared patterns of abundance, organized according to their functions and
50 physical interactions²¹⁻²³. Surprisingly, much of the proteome-level organization of co-abundance
51 patterns are not detected at the RNA level^{21,22}. Furthermore, physically interacting proteins are
52 much more likely to have coordinated protein abundance than RNA co-expression^{21,22,24}. RNA co-
53 expression in both mouse and human cells often arises from the chromosomal proximity of genes
54 even when they are functionally unrelated^{25,26}. This likely unproductive co-expression pattern is
55 absent at the protein level^{27,28}, suggesting that post-transcriptional regulation plays a significant
56 role in proteome organization.

57 Translation regulation, a crucial post-transcriptional process, may bridge this gap, given its vital
58 roles in development, maintaining cellular homeostasis, and responding to environmental
59 changes²⁹⁻³⁵. There are three lines of evidence that suggest the possibility of coordinated
60 translation of functionally and physically associated proteins across different biological contexts.

61 First, mammalian mRNAs bind various proteins to form ribonucleoproteins that influence their
62 lifecycle from export to translation³⁶. The set of proteins interacting with an mRNA varies with
63 time and context, significantly altering the duration, efficiency, and localization of protein
64 production. These observations led to the proposal of the post-transcriptional RNA regulon model
65 over two decades ago, positing that functionally related mRNAs are regulated together post-
66 transcriptionally^{37,38}. Supporting this model, Puf3, a Pumilio family member in yeast, represses
67 the translation of sequence-specific mRNAs encoding mitochondrial proteins³⁸. Similarly, in
68 human cells, CSDE1/UNR regulates the translation of mRNAs involved in epithelial-to-
69 mesenchymal transition³⁹. However, it remains to be determined if translation of functionally
70 related mRNAs is coordinately regulated across different conditions.

71 Second, in both *E. coli* and yeast, proteins within multiprotein complexes are synthesized in

72 stoichiometric proportions needed for assembly^{40,41}. This translational regulation likely tunes
73 protein production to minimize the synthesis of excess protein components that would otherwise
74 need to be degraded⁴². However, in human cells, evidence of such proportional synthesis is
75 reported for only two complexes: ribosomes^{41,43} and the oxidative phosphorylation machinery⁴⁴.
76 Furthermore, these observations have been made in a very limited number of cell lines, which
77 limits the generalizability of this concept across diverse cell types and other functionally related
78 protein groups.

79 Third, the formation of many protein complexes is facilitated by the co-translational folding of
80 nascent peptides^{40,41,45,46}. For instance, in bacteria, the anti-Shine-Dalgarno sequence induces
81 translational pausing to modulate the co-translational folding of nascent peptides⁴⁶. Co-
82 translational assembly ensures that protein subunits are synthesized near each other, enabling near
83 concurrent interactions, which are crucial for the biogenesis of some complex protein structures⁴⁷.
84 Recent evidence indicates that co-translational assembly may also be relatively common in human
85 cells⁴⁸.

86 Co-translational assembly and stoichiometric synthesis rates of protein complexes suggest
87 coordinated translation of several mRNAs within a given cell type.. However, due to the lack of
88 robust, transcriptome-wide translational efficiency (TE) measurements across diverse biological
89 conditions, it remains to be seen whether such coordination extends across different cell types or
90 conditions. To address this, we analyzed thousands of matched ribosome profiling and RNA-seq
91 datasets from >140 human and mouse cell lines and tissues. To quantify the similarity of
92 translation efficiency patterns of transcripts across cell types and tissues, analogously to RNA co-
93 expression, we introduce the concept of Translation Efficiency Covariation (TEC). based on a
94 compositional data analysis approach^{49,50} Our findings demonstrate that TEC can reveal gene
95 functions not identified through RNA co-expression analysis alone and uncovered shared motifs
96 for RNA binding proteins (RBPs) among genes exhibiting TEC. Physically interacting proteins
97 are highly enriched for both TEC and RNA co-expression. Further supporting the functional
98 significance of this concept, TEC among genes is highly conserved between humans and mice.

99 **RESULTS**

100 **Integrated analysis of thousands of ribosome profiling and RNA-seq measurements enable** 101 **quantitative assessment of data quality**

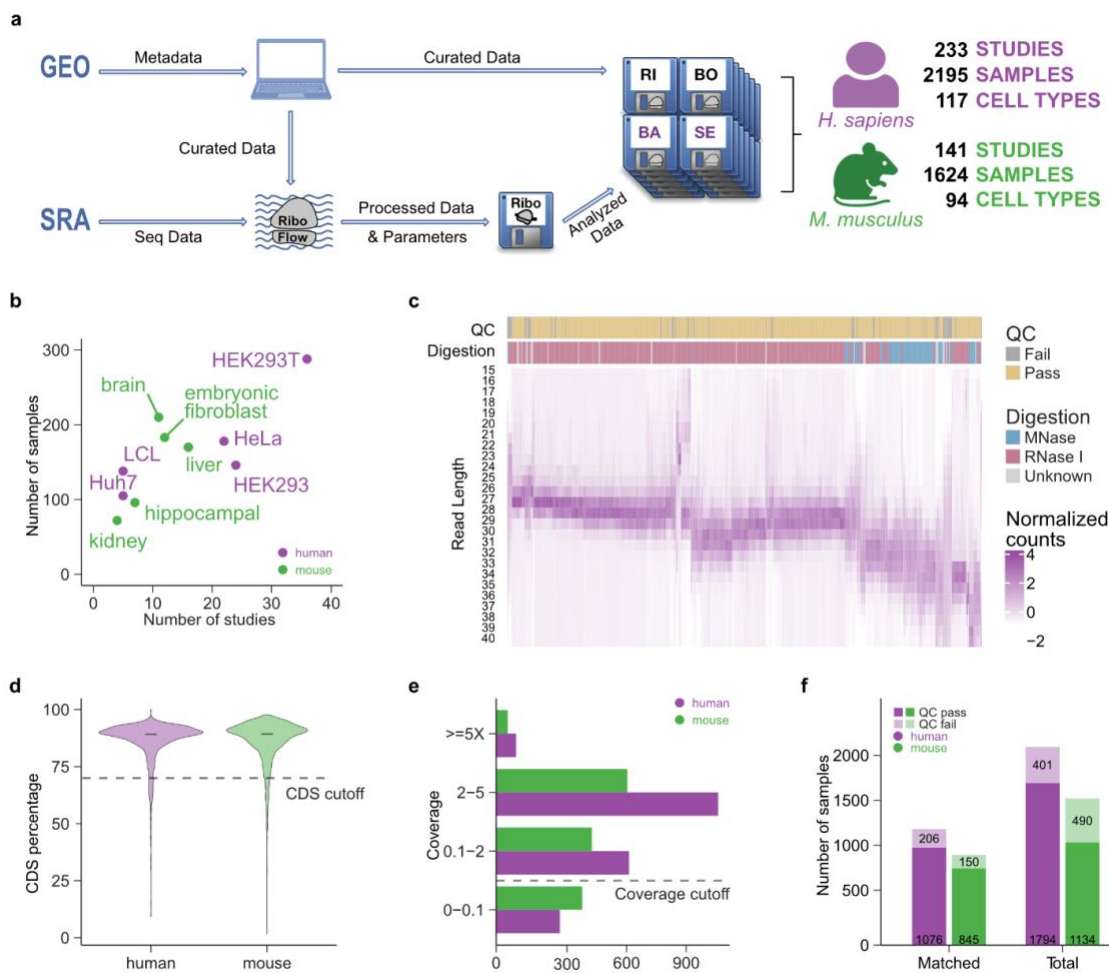
102 We undertook a comprehensive, large-scale meta-analysis of ribosome profiling data to quantify
103 TE across different cell lines and tissues. We collected 2,195 ribosome profiling datasets for
104 humans and 1,624 experiments for mice, along with their metadata (Fig. 1a; Methods). Given that
105 metadata is frequently reported in an unstructured manner and lacks a formal verification step, we
106 conducted a manual curation process to rectify inaccuracies and collect missing information, such
107 as experimental conditions and cell types used in experiments. One crucial aspect of our manual
108 curation was pairing between ribosome profiling and corresponding RNA-seq when possible.
109 Overall, 1,282 (58.4%) human and 995 (61.3%) mouse ribosome profiling samples were matched
110 with corresponding RNA-seq data (table S1). The resulting curated metadata facilitated the
111 uniform processing of ribosome profiling and corresponding RNA-seq data using an open-source

112 pipeline⁵¹. We call the resulting repository harboring these processed files RiboBase (table S1).

113 In RiboBase, the top cell types with the most experiments were HEK293T (13.1%) and HeLa
114 (8.1%) for human; in mouse, the leading tissues were brain (9.6%), embryonic fibroblasts (8.3%),
115 and liver (7.7%) (Fig. 1b; table S1). The median number of sequencing reads for ribosome profiling
116 samples was ~43.2 million for humans and ~37.5 million for mice, respectively (ExtendedDataFig.
117 1a-b; table S2-3; supplementary text). A majority of reads contained adapter sequences included
118 during library preparation (with medians of 82.2% and 79.2% of total reads having adapters for
119 human and mouse, respectively). Due to the substantial presence of ribosomal RNA in ribosome
120 profiling datasets, only around 15% of total reads aligned to the transcript reference
121 (ExtendedDataFig. 1c-d; table S4-5; supplementary text).

122 The length of ribosome-protected mRNA footprints (RPFs) provides valuable information about
123 data quality, the experimental protocol used, and translational activity⁵². The choice of nuclease
124 impacts the resulting read length distribution of RPFs⁵³ (ExtendedDataFig. 2a-b). In agreement,
125 we found that the peak position and range of RPF lengths were closely associated with the type of
126 digestion enzymes used in human cancer samples (Fig. 1c). To account for the variability of RPF
127 length distributions across the compendium of experiments, we developed a module that allowed
128 for setting sample-specific RPF read length cutoffs (ExtendedDataFig. 3a; Methods). This
129 dynamic approach proved more effective than using fixed minimum and maximum values for RPF
130 lengths, resulting in a higher retrieval of usable reads (median increase of 10.8% for human and
131 17.1% for mouse) and an increased proportion of reads within the coding sequence (CDS) region
132 (ExtendedDataFig. 3b).

133 After selecting a set of RPFs, we assessed the quality of ribosome profiling data within RiboBase
134 using two additional criteria. Given that translating ribosomes should be highly enriched in
135 annotated coding regions, we require that at least 70% of RPFs should be mapped to the CDS. We
136 found that 160 human and 115 mouse samples failed to meet this criterion (Fig. 1d; table S6-7).
137 Subsequently, we required a minimum number of RPFs that map to CDS to ensure sufficient
138 coverage of translated genes (Methods). There were 318 human and 431 mouse samples with less
139 than 0.1X transcript coverage (Fig. 1e; table S6-7). Altogether, 1,794 human samples and 1,134
140 mouse samples were retained for in-depth analysis. Of these, 1,076 human and 845 mouse samples
141 were paired with matching RNA-seq data. Our results indicate a considerable fraction of publicly
142 available ribosome profiling experiments had suboptimal quality (18.3% of the human and 30.1%
143 of the mouse samples) (Fig. 1f). Interestingly, the data quality appeared to be independent of time
144 (ExtendedDataFig. 4). Additionally, we found that samples that passed our quality thresholds were
145 more likely to exhibit three-nucleotide periodicity compared to those that failed quality control
146 (92.59% vs 78.30% for humans and 91.36% vs 86.73% for mice; ExtendedDataFig. 5; Methods).
147 These findings underscore the necessity of meticulous quality control for the selection of
148 experiments to enable large-scale data analyses.



149

150 **Fig. 1 | RiboBase: a comprehensive ribosome profiling database with thousands of**
 151 **experiments.** **a**, Schematic of RiboBase. We manually curated metadata and processed the
 152 sequencing reads using a uniform pipeline (RiboFlow⁵¹). **b**, Top five most highly represented cell
 153 lines or tissues with respect to the number of experiments were plotted. **c**, We determined the
 154 ribonuclease used to generate ribosome profiling data for 680 experiments using human cancer
 155 cell lines. For each experiment, the read length distribution of RPFs mapping to coding regions
 156 was visualized as a heatmap. The color represents the z-score adjusted RPF counts (Methods).
 157 Each experiment where the percentage of RPFs mapping to CDS was greater than 70% and
 158 achieving sufficient coverage of the transcript ($\geq 0.1X$) was annotated as QC-pass (Methods). **d**,
 159 For the 3,819 ribosome profiling experiments in RiboBase, we applied a function to select the
 160 range of RPFs for further analysis (Methods). We calculated the proportion of the selected RPFs
 161 that map to the coding regions (y-axis). The horizontal line represents the median of the
 162 distribution. **e**, Experiments (x-axis) were grouped by the transcript coverage (y-axis). **f**, Among

163 the ribosome profiling experiments in RiboBase, 2,277 of them had corresponding RNA-seq data
164 (matched). The number of samples that pass quality controls were plotted.

165 **Translation efficiency is conserved across species and is cell-type specific**

166 Ribosome profiling measures ribosome occupancy, a variable influenced by both RNA expression
167 and translation dynamics. Thus, estimating translation efficiency necessitates analysis of paired
168 RNA-seq and ribosome profiling data. To assess accurate matching in RiboBase, we first
169 compared the coefficient of determination (R^2) between matched ribosome profiling and RNA-seq
170 data to that from other pairings within the same study. As would be expected from correct
171 matching, we found that matched samples had significantly higher similarity on average (Fig. 2a;
172 Welch two-sided t-test p-value = 2.2×10^{-16} for human and p-value = 2.1×10^{-5} for mouse). We
173 then implemented a scoring system to quantitatively evaluate the correctness of our manual
174 matching information (Methods). 99.2% of human samples and 98.5% of mouse samples had a
175 sufficiently high matching score, demonstrating the effectiveness of our manual curation strategy
176 (ExtendedDataFig. 6a; Methods).

177 Using the set of matched ribosome profiling and RNA-seq experiments, we next quantified TE,
178 which is typically defined as the log ratio of ribosome footprints to RNA-seq reads, normalized as
179 counts per million⁵⁴. However, this approach leads to biased estimates with significant
180 drawbacks⁵⁵. To address this limitation, we calculate TE based on a regression model using a
181 compositional data analysis method^{49,50,56}, avoiding the mathematical shortcomings of using a log-
182 ratio (Fig. 2b; ExtendedDataFig. 6a-c, 7; table S8-11; Methods).

183 We next assessed whether measurement errors due to differences in experimental procedures
184 dominate variability that would otherwise be attributed to biological variables of interest.
185 Specifically, we compared similarities between experiments that used the same cell type or tissue
186 in different studies (ExtendedDataFig. 8a). We found that ribosome profiling or RNA experiments
187 from the same cell type or tissue exhibited higher similarity compared to those from different cell
188 lines or tissues (Fig. 2c). Consistent with this observation, TE values displayed higher Spearman
189 correlation coefficient within the same cell type or tissue (median correlation coefficient of 0.56
190 and 0.53 in human and mouse, respectively) compared to different cell lines and tissues (median
191 correlation coefficient of 0.49 and 0.45 in human and mouse, respectively) (Fig. 2d).

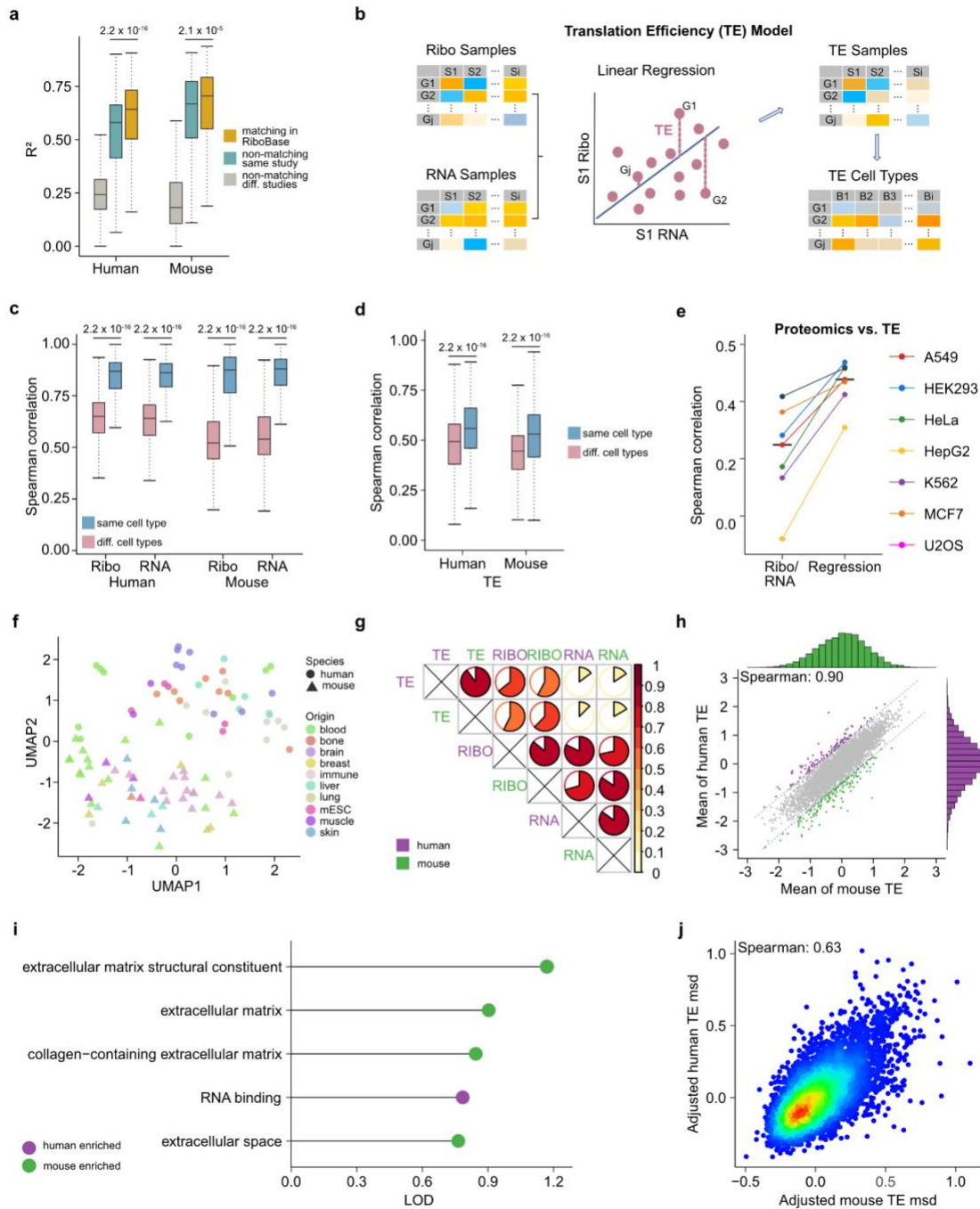
192 We expected that a more accurate estimate of TE would show a stronger correlation with protein
193 abundance. We calculated for each transcript the cell type-specific TE by taking the average of TE
194 values across all experiments conducted with that particular cell line. Indeed, our results show that
195 compared to the log-ratio definition, the TE derived using the regression approach with winsorized
196 read counts (ExtendedDataFig. 8b; ExtendedDataFig. 9-11; supplementary text) is more strongly
197 correlated with protein abundance in seven cancer cell lines (mean Spearman correlation
198 coefficient of 0.465 vs 0.219; Fig. 2e).

199 Furthermore, TE measurements from cell lines and tissues with the same biological origin (e.g.,
200 blood) tended to cluster together, supporting the existence of cell-type-specific differences in TE
201 (Fig. 2f). As expected, mean ribosome occupancy and RNA expression across cell types showed
202 a strong correlation (Spearman correlation: ~ 0.8), yet mean TE was only weakly associated with

203 RNA expression (Spearman correlation: ~ 0.2) (Fig. 2g). Taken together, our analyses demonstrate
204 that our compositional regression-based approach to calculating TE ensures more accurate and
205 consistent measurements across different cell types and conditions.

206 Measurements of TE in two species across a large number of cell types enabled us to investigate
207 the conservation of TE, ribosome occupancy, and RNA expression. Transcriptomes, ribosome
208 occupancy, and proteomes exhibit a high degree of conservation across diverse organisms^{57,58}.
209 Consistently, we found average ribosome occupancy, RNA expression, and TE across different
210 cell lines and tissues were highly similar between orthologous genes in human and mouse (Fig.
211 2g; table S12). Specifically, the Spearman correlation coefficient of mean TE across cell types and
212 tissues between human and mouse was 0.9 (Fig. 2h), which is comparable to the mean RNA
213 expression correlation between human and mouse (~ 0.86 , ExtendedDataFig. 12a). Using a 95%
214 prediction interval to identify outlier genes, we found that outlier genes with higher mean TE in
215 humans compared to mice were enriched in the gene ontology term ‘RNA binding function’ (Fig.
216 2i). In contrast, genes with elevated mean TE in mice were enriched for having functions related
217 to extracellular matrix and collagen-containing components (Fig. 2i). The enrichment of genes
218 with higher TE in mice, particularly those from the extracellular matrix and collagen-containing
219 components, may be due to the fact that many samples in mouse studies are derived from the early
220 developmental stage⁵⁹.

221 Despite the high correlation of mean TE across various cell lines and tissues between human and
222 mouse, TE distinctly exhibits cell-type specificity. While several studies compared the
223 conservation of TE between the same tissues of mammals or model organisms^{58,60,61}, our dataset
224 uniquely enabled us to determine the conservation of variability of TE for transcripts across
225 different cell types. Intriguingly, we observed a moderately high similarity between the variability
226 of TE of orthologous genes in human and mouse (Spearman partial correlation coefficient = 0.63;
227 Fig. 2j; ExtendedDataFig. 12b-d; Methods). Our results reveal that certain genes exhibit higher
228 variability of TE across cell types and this is a conserved property between human and mouse.



229
 230 **Fig. 2 | TE defined using a compositional linear regression model is conserved across cell**
 231 **types and species.** **a**, The distribution of coefficient of determination (R^2 , y-axis) between
 232 ribosome profiling data and RNA-seq in RiboBase was compared to random matching within the
 233 same study and across different studies. In each figure panel containing boxplots, the horizontal

234 line corresponds to the median. The box represents the interquartile range (IQR) and the whiskers
235 extend to the largest value within 1.5 times the IQR. The significant p-value shown in this figure
236 was calculated using the two-sided Wilcoxon test. **b**, Schematic of TE calculation using the linear
237 regression model with compositional data (CLR transformed; Methods; ExtendedDataFig. 7). **c**,
238 Distribution of correlations of TE (linear regression model) across experiments. **d**, Correlation
239 between TE and protein abundance from seven human cell lines¹⁰⁰ was calculated using log-ratio
240 of ribosome profiling and RNA expression or compositional regression method. The horizontal
241 line corresponds to the median. **e**, The distribution of Spearman correlations between experiments
242 (y-axis) was calculated based on whether they originated from identical or different cell lines or
243 tissues. **f**, We used UMAP to cluster the TE values of all genes across different cell types,
244 considering only those origins with at least five distinct cell types. **g**, The Spearman correlation of
245 9,194 orthologous genes between human and mouse across TE, ribosome profiling, and RNA-seq
246 levels. The circles represent the value of the Spearman correlation between groups. **h**, TE values
247 were averaged across cell types and tissues for either human and mouse. Each dot represents a
248 gene, and a 95% prediction interval was plotted to identify outlier genes (highlighted in purple and
249 green). **i**, We conducted GO term enrichment analysis for outlier genes from panel H. We ranked
250 the GO terms (y-axis) by the logarithm of the odds (LOD; x-axis). **j**, The correlation of the standard
251 deviation of TE (quantified with adjusted metric standard deviation (msd); Methods;
252 ExtendedDataFig. 12c-d) for orthologous genes across different cell types between human and
253 mouse.

254 **Translation efficiency covariation (TEC) is conserved between human and mouse**

255 Uniform quantification of TE enabled us to investigate the similarities in TE patterns across cell
256 types. Given the usefulness of RNA co-expression in identifying shared regulation and biological
257 functions, we aimed to establish an analogous method to detect patterns of translation efficiency
258 similarity among genes. To achieve this, we employed the proportionality score (ρ)^{50,56}, a
259 statistical method that quantifies the consistency of how relative TE changes across different
260 contexts (Methods). Recent work suggested that the proportionality score enhances cluster
261 identification in high-dimensional single-cell RNA co-expression data¹⁰. Consistent with these
262 findings, our analysis revealed its particular effectiveness in quantifying ribosome occupancy
263 covariation (ExtendedDataFig. 13; Methods). We calculated ρ scores for all pairs of human or
264 mouse genes where a high absolute ρ score indicates significant translation efficiency
265 covariation (TEC) between pairs (Fig. 3a).

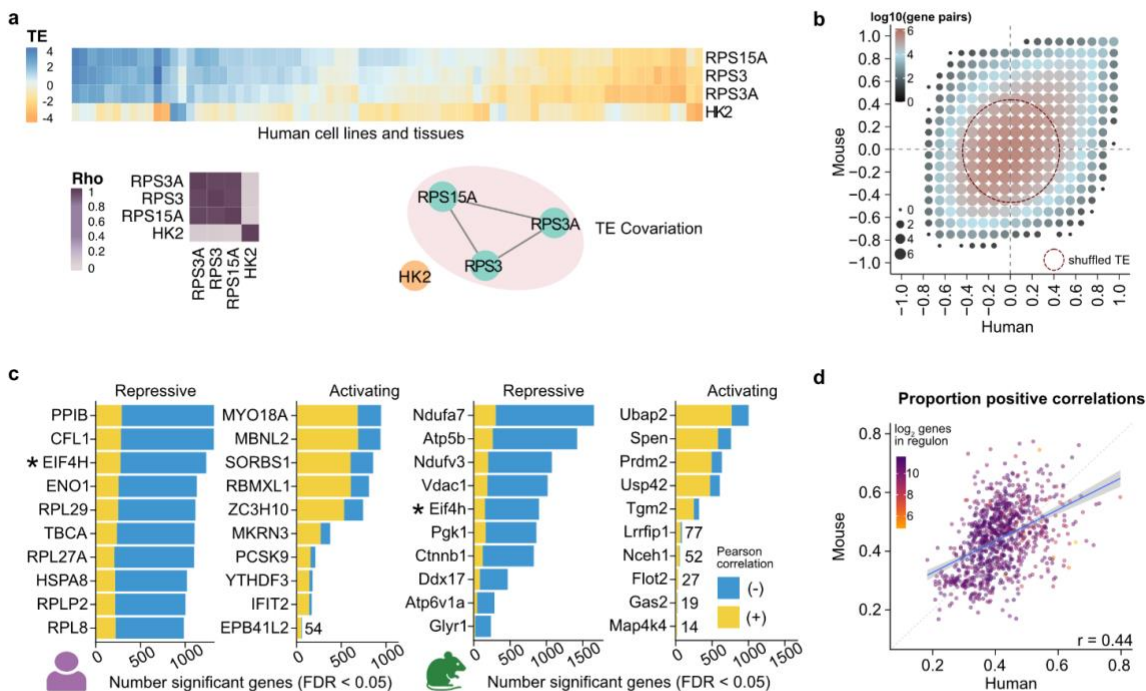
266 Previous studies have indicated that RNA co-expression between genes is conserved in
267 mammals^{57,62,63}. To assess the potential evolutionary significance of the newly introduced TEC
268 concept, we evaluated its conservation across human and mouse transcripts. Indeed, TEC was
269 highly similar for orthologous gene pairs in humans and mice (Fig. 3b, Pearson correlation
270 coefficient 0.41), compared to a negligible correlation in TEC derived from shuffled TE values
271 (ExtendedDataFig. 14, Pearson correlation coefficient 0.00022). Our findings imply that
272 translation efficiency patterns are evolutionarily preserved, paralleling the conservation of RNA
273 co-expression.

274 RNA co-expression analyses led to the discovery of regulatory motifs and shared transcription

275 factor binding sites⁶⁴. We hypothesized that TEC among genes may nominate RNA binding
276 proteins (RBPs) as potential drivers of TEC⁶⁵. We identified groups of transcripts whose TE is
277 correlated with the RNA expression of experimentally determined RBPs⁶⁶ (1274 human and 1762
278 mouse RBPs; Methods). The number of transcripts whose TE significantly correlates with the
279 expression of each RBP differed widely, with ranges of 28-3052 (human) and 14-2393 (mouse)
280 (<https://zenodo.org/uploads/11359114>; Pearson correlation FDR < 0.05). We refer to transcripts
281 whose TE is significantly correlated with an RBP's expression as the RBP's regulon.

282 Interestingly, some RBP regulons were dominated by positive or negative correlations, suggesting
283 activating or repressing functions for RBPs (Fig. 3c-d). For example, ZC3H10 has largely positive
284 correlations (71% of RBP regulon) (Fig. 3c). Conversely, the RNA expression of subunits of
285 ubiquinone oxidoreductase (Ndufa7, Ndufv3) is negatively correlated with TE for many genes in
286 mice (Fig. 3c). Unexpectedly, we found that the RNA expression of ribosomal protein genes is
287 negatively correlated with TE of many other transcripts (<https://zenodo.org/uploads/11359114>).
288 This may indicate that transcriptome-wide TE is tempered during ribosome biogenesis, perhaps as
289 a result of competition for ribosomes and other biosynthetic resources (tRNAs, amino acids)
290 devoted to synthesizing new ribosomal proteins.

291 To identify evolutionarily conserved RBP regulons, we examined the intersection of significant
292 RBP-gene correlations between human and mouse. At least some activating RBP functions may
293 be evolutionarily conserved, as there was a correspondence between human and mouse in the
294 proportion of regulon genes with positive correlations (Pearson correlation 0.44; Fig. 3d). To
295 nominate RBPs that may modulate TE, we calculated the proportionality score of genes in each
296 regulon and selected RBP regulons that had high absolute scores, reasoning that directional
297 impacts on TE might be more likely if the RBP engages these transcripts. We found 85 RBPs
298 where genes in the RBP's regulon had high TEC (mean absolute pairwise rho >90th percentile;
299 ExtendedDataFig. 15; supplementary text). Some of these RBPs were previously known to
300 regulate TE, including PARK7 and VIM (ExtendedDataFig. 16; supplementary text). Taken
301 together, our analyses nominate RBPs that may coordinate the TEC of evolutionarily conserved
302 RNA regulons.



303

304 **Fig. 3 | Translation efficiency covariation is conserved between human and mouse. a,**
 305 Example illustrating translation efficiency covariation (TEC) between genes. The top section
 306 presents TE patterns across cell types in human. The bottom left part displays the similarity of the
 307 pattern between these genes quantified using proportionality scores. **b,** We calculated the TEC for
 308 gene pairs and compared their differences for the same orthologous gene pairs between human and
 309 mouse. In the figure panel, each dot represents the aggregated \log_{10} -transformed counts of gene
 310 pairs falling within specified ranges. We also calculated TEC using randomized TE for each gene
 311 (shuffled). The red dashed line in the figure captures the 95% gene pair TEC values obtained with
 312 shuffled TE (ExtendedDataFig. 14). **c,** Top ten candidates activating and repressive RBPs: human
 313 (left) and mouse (right). The number of genes with significant correlations between gene TE and
 314 RBP expression is shown. An asterisk marks genes in the top ten in both species. **d,** Each point is
 315 a RBP and plotted is the proportion of positive correlations between TE for genes in the regulon
 316 and the RNA expression of the RBP. Blue line is a linear fit with 95% confidence intervals in gray.
 317 Pearson correlation coefficient is shown.

318 **Translation efficiency covariation (TEC) between transcripts across cell lines and tissues is**
 319 **associated with shared biological functions**

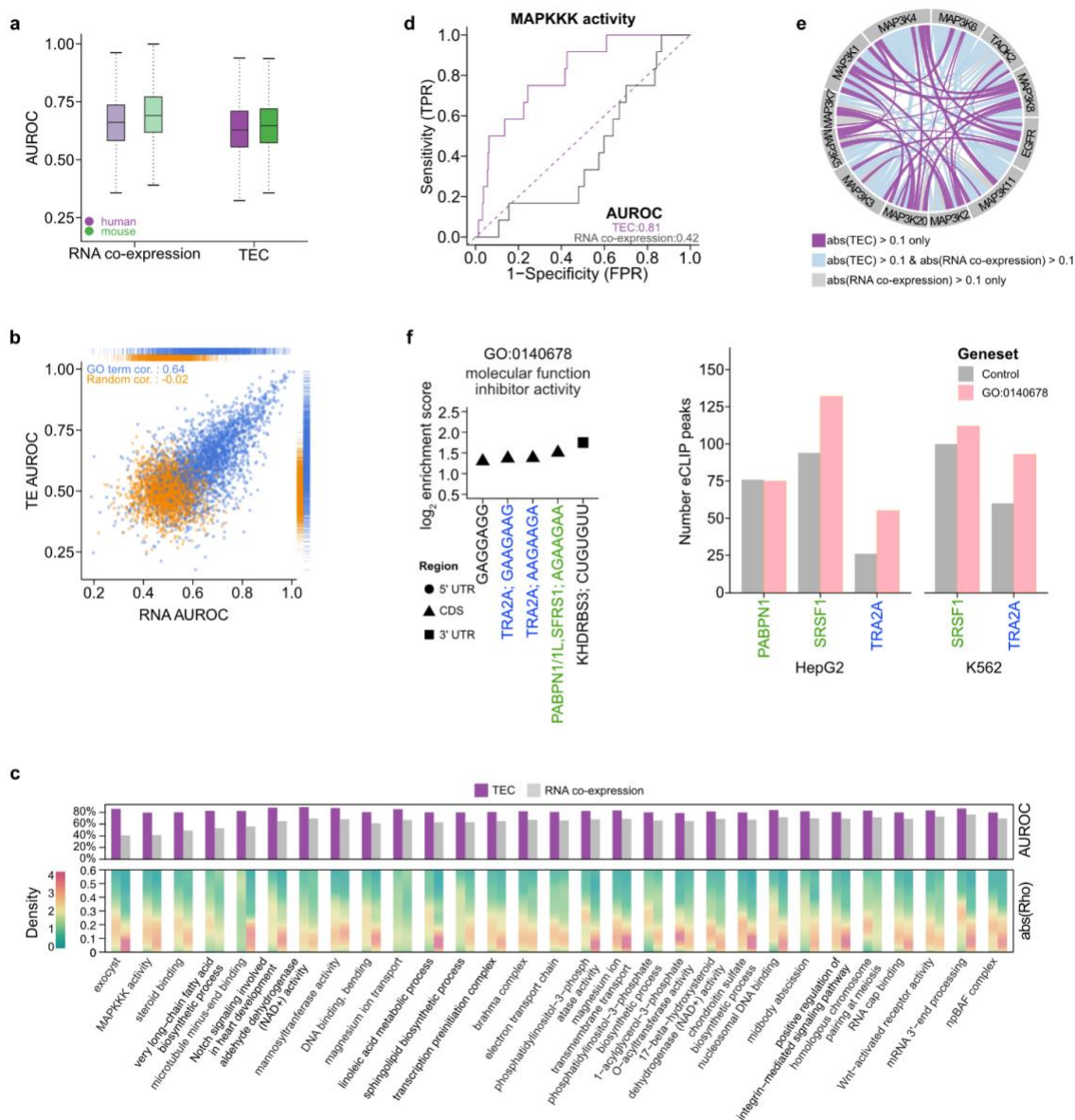
320 Given that co-expression at the RNA level is predictive of shared biological functions^{11,67,68}, we
 321 next assessed whether TEC indicates common biological roles among genes. We calculated the
 322 area under the receiver operating characteristic curve (AUROC) to measure the ability of TEC in

323 distinguishing genes with the same biological functions (Methods). Genes that are annotated with
324 a common GO term exhibited a similar degree of RNA co-expression and TEC, both of which
325 were significantly higher than would be expected by chance (Median AUROC across GO terms
326 calculated with TEC: 0.63 for human, 0.65 for mouse; RNA co-expression RNA: 0.66 for human,
327 0.69 for mouse; Fig. 4a; table S13-14; Methods). These findings demonstrate that TEC, similar to
328 RNA co-expression, serves as an indicator of shared biological functions among genes.

329 Furthermore, we observed that biological functions whose members exhibit a high degree of RNA
330 co-expression were also likely to have TEC. Specifically, the Spearman correlation between the
331 AUROC scores calculated using TEC and RNA co-expression was ~0.64 for human GO terms in
332 contrast to ~-0.02 when random genes were grouped (Fig. 4b). Despite the low correlation between
333 average RNA expression and TE for human genes (Fig. 2g), our results highlight that members of
334 specific biological functions whose RNA expression is coordinated across cell types tend to exhibit
335 consistent translation efficiency patterns. This finding suggests coordinated regulation at both
336 transcriptional and translational levels among functionally related genes.

337 While many gene functions were predicted accurately with both RNA co-expression and TEC, we
338 noted specific exceptions. Notably, genes in 29 human GO terms demonstrated significantly
339 stronger TEC than RNA co-expression (at least 0.1 higher AUROC; Fig. 4c; ExtendedDataFig.
340 17-18; supplementary text). An example of such a GO term is ‘MAPKKK activity’ (Fig. 4d-e).
341 While there is limited evidence of direct translational regulation of the MAPKKK family, the RBP
342 IMP3 may provide a potential mechanism for such regulation⁶⁹. Additionally, there is post-
343 translational regulation through the binding of activated RAS to genes from the MAPKKK family,
344 leading to their activation⁷⁰. These results indicate that some genes with specific biological
345 functions exhibit greater similarity at the translational level.

346 We hypothesized that genes with shared functions and high TEC may be regulated through a
347 common mechanism, analogous to shared transcription factor binding sites that mediate RNA co-
348 expression^{71,72}. Accordingly, we expected these genes to harbor sequence elements bound by
349 RBPs. We identified enriched heptamers in the transcripts of five human and three mouse GO
350 terms with significant TEC and at least 12 genes in the GO term (AUROC measured with TEC >
351 0.7, difference in AUROC between TEC and RNA co-expression > 0.2; Fig. 4f; ExtendedDataFig.
352 17b; ExtendedDataFig. 18e; Methods). For example, we found AG-rich motifs in coding regions
353 of human genes with “molecular function inhibitor activity” (Fig. 4f). These motifs match the
354 known binding sites of three RBPs (TRA2A, PABPN1, and SRSF1). In line with the enrichment
355 of these motifs, analysis of eCLIP data revealed increased deposition of these RBPs in the coding
356 sequences of genes in this GO term compared to matched control transcripts (Fig. 4f; Methods).
357 Furthermore, we identified several additional enriched heptamers that currently have no RBP
358 annotations, suggesting these motifs might be targets for RBPs that have not yet been
359 characterized.



360
361 **Fig. 4 | Genes associated with certain biological functions exhibit higher similarity patterns**
362 **in TE than in RNA expression.** **a**, We calculated the similarity of expression (quantified by
363 AUROC; y-axis) among genes within 2,989 human and 3,340 mouse GO terms. In the box plot,
364 the horizontal line corresponds to the median. The box represents the IQR and the whiskers extend
365 to the largest value within 1.5 times the IQR. **b**, Each blue dot represents the AUROC calculated
366 for a given GO term using TEC and RNA co-expression levels. Orange dots represent the same
367 values for random grouping of genes (Methods). **c**, For GO terms where genes exhibit greater
368 similarity at the TE level than at the RNA expression level (AUROC for TEC > 0.8, and difference
369 of AUROC measured with TEC and RNA co-expression > 0.1), we visualized the distribution of

370 absolute rho scores for gene pairs (bottom; gene pairs with $\text{abs}(\rho) > 0.1$). **d**, AUROC plot
371 calculated with genes associated with MAPKKK activity. **e**, In the circle plot, the connections
372 display absolute rho above 0.1 either at TE level alone (purple), at both RNA and TE levels (blue),
373 or RNA level alone (gray) for gene pairs involved in MAPKKK activity. **f**, Motif enrichment (left)
374 for the GO term ‘molecular function inhibitor activity’ (ExtendedDataFig. 17b). RNA binding
375 proteins (RBPs) matching the motifs from oRNAmotif¹³⁴ or Transite¹³³ are indicated. Enhanced
376 cross-linking immunoprecipitation (eCLIP) data¹³⁵ indicates increased binding of TRA2A and
377 SRSF1 in the CDS of genes for this GO term compared to matched control genes with similar
378 sequence properties (Methods).

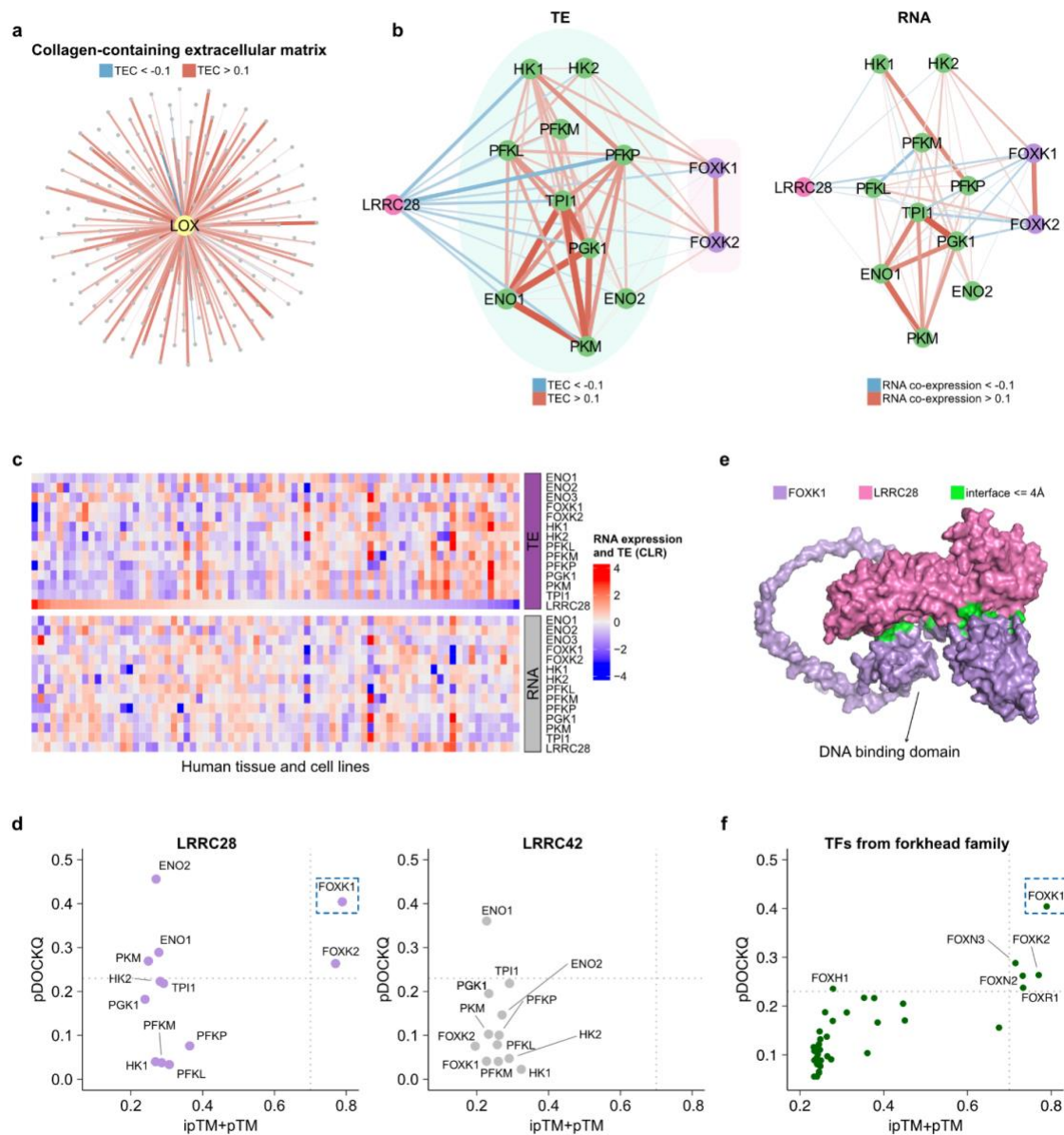
379 **TEC reveals gene functions**

380 We next investigated whether gene functions may be predicted by utilizing TEC, given the success
381 of RNA co-expression for this task^{67,68}. The functional annotations of human genes are
382 continuously being updated, providing an opportunity to test this hypothesis using recently added
383 information to the knowledge base. Specifically, we used functional annotations from the GO
384 database from January 1, 2021, to determine functional groups that demonstrate strong TEC among
385 its members (AUROC > 0.8) and developed a framework to predict new functional associations
386 with these groups (Methods). By comparing our predictions to annotations from December 4,
387 2022, we confirmed the predicted association of the *LOX* gene with the GO term ‘collagen-
388 containing extracellular matrix’. *LOX* critically facilitates the formation, development,
389 maturation, and remodeling of the extracellular matrix by catalyzing the cross-linking of collagen
390 fibers, thereby enhancing the structural integrity and stability of tissues^{73,74}. Our prediction
391 successfully identified this new addition, as *LOX* exhibits positive similarity in TE with the vast
392 majority of genes in this term (Fig. 5a).

393 Recognizing the capacity of TEC to elucidate biological functions, we utilized a recent version of
394 GO annotations (December 4, 2022) to systematically predict new associations for genes. To
395 underscore the unique insights gained from TEC, we focused on the 33 human and 31 mouse GO
396 terms that either exhibited significantly higher TEC than RNA co-expression (Table 1) or provided
397 new functional predictions that were only supported by TEC (the ranking of the newly predicted
398 gene with RNA co-expression fell beyond the top 50%, table S15-16; Methods). By focusing on
399 these GO terms, we aimed to identify similarity patterns based on TE, revealing functional
400 associations that would not be detected by RNA co-expression. We conducted a literature search
401 to determine if prior research supported these predictions, finding that 11 have already been
402 corroborated by previous publications, although they have not yet been reflected in the relevant
403 GO term annotations (Table 1; supplementary text). For example, cryo-electron microscopy
404 experiments demonstrated that human DNMT1 binds to hemimethylated DNA in conjunction with
405 ubiquitinated histone H3⁷⁵. This binding facilitates the enzymatic activity of DNMT1 in
406 maintaining genomic DNA methylation. Our analysis revealed that DNMT1 was the highest
407 ranking prediction exhibiting strong TEC with genes associated with nucleosomal DNA binding
408 function. In mouse, we predicted *Plekha7* to be a member of the regulation of developmental
409 processes. This prediction was recently validated by the observation of neural progenitor cell
410 delamination upon the disruption of *Plekha7*^{76–80}.

411 The high rate of validation of our predictions in the literature suggested that other predictions based
412 on TEC may reflect new and yet to be confirmed functions. In particular, we observed that the
413 human leucine-rich repeat-containing 28 (LRRC28) gene displays strong TEC with glycolytic
414 genes, but is not co-expressed at the RNA level (Fig. 5b-c, table S17). Specifically, *LRRC28*
415 displayed negatively correlated TE with key glycolytic genes including *HK1*, *HK2*, *PFKL*, *PFKM*,
416 *PFKP*, *TPII*, *PGK1*, *ENO1*, *ENO2*, *PKM*, and two transcription factors *FOXK1* and *FOXK2* that
417 regulate glycolytic genes⁸¹. Given that the leucine-rich repeat domains typically facilitate protein-
418 protein interactions⁸², LRRC28 may interact directly with one or more of the glycolytic proteins.
419 Using AlphaFold2-Multimer⁸³, we calculated the binding confidence score between LRRC28 and
420 all glycolysis-associated proteins (Methods) and found that LRRC28 has a very high likelihood of
421 binding to FOXK1 (Fig. 5d-e).

422 FOXK1 is a member of the forkhead family of transcription factors that share a structurally similar
423 DNA-binding domain^{84,85}. Interestingly, LRRC28 likely binds both the non-DNA-binding region
424 and DNA-binding domain of FOXK1 (distance < 4 angstroms; Fig. 5e; ExtendedDataFig. 19).
425 This observation led us to examine the specificity of the interaction between LRRC28 and FOXK1.
426 We calculated the binding probabilities of LRRC28 with 35 other forkhead family transcription
427 factors, finding that FOXK1 exhibits the strongest evidence of physical interaction with LRRC28
428 (Fig. 5f). This specificity is potentially due to a unique binding site between LRRC28 and
429 FOXK1's non-DNA-binding region (Fig. 5e). As an additional control, we selected LRRC42, a
430 protein with leucine-rich repeats that does not exhibit TEC with glycolytic genes. As expected,
431 LRRC42 showed a very low likelihood of interaction with any of the glycolytic genes, including
432 FOXK1 (Fig. 5d). These findings suggest that LRRC28 may serve as a regulator of glycolysis by
433 binding to FOXK1, thereby preventing FOXK1 from binding to the promoter regions of glycolytic
434 genes and leading to the downregulation of glycolysis. Taken together, TEC reveals shared
435 biological functions and predicts novel associations, providing insights not attainable with RNA
436 co-expression analysis alone.



437

438 **Fig. 5 | TEC enables the prediction of novel gene functions.** **a**, We predicted that LOX belongs
 439 to the collagen-containing extracellular matrix using an older version of human GO terms (January
 440 1, 2021) and confirmed this prediction with the newer version (December 4, 2022; Methods). The
 441 network displays the similarity in TE between LOX (yellow dot) and other genes (gray dots) from
 442 the collagen-containing extracellular matrix. Line weight in figure panels indicates the absolute
 443 value of rho from 0.1 to 1. **b**, The networks display the rho between LRRC28 and glycolytic genes
 444 at the TE level (on the left) and RNA level (on the right) in humans. Green dots represent genes
 445 that belong to the glycolysis pathway, purple nodes are transcription factors that regulate

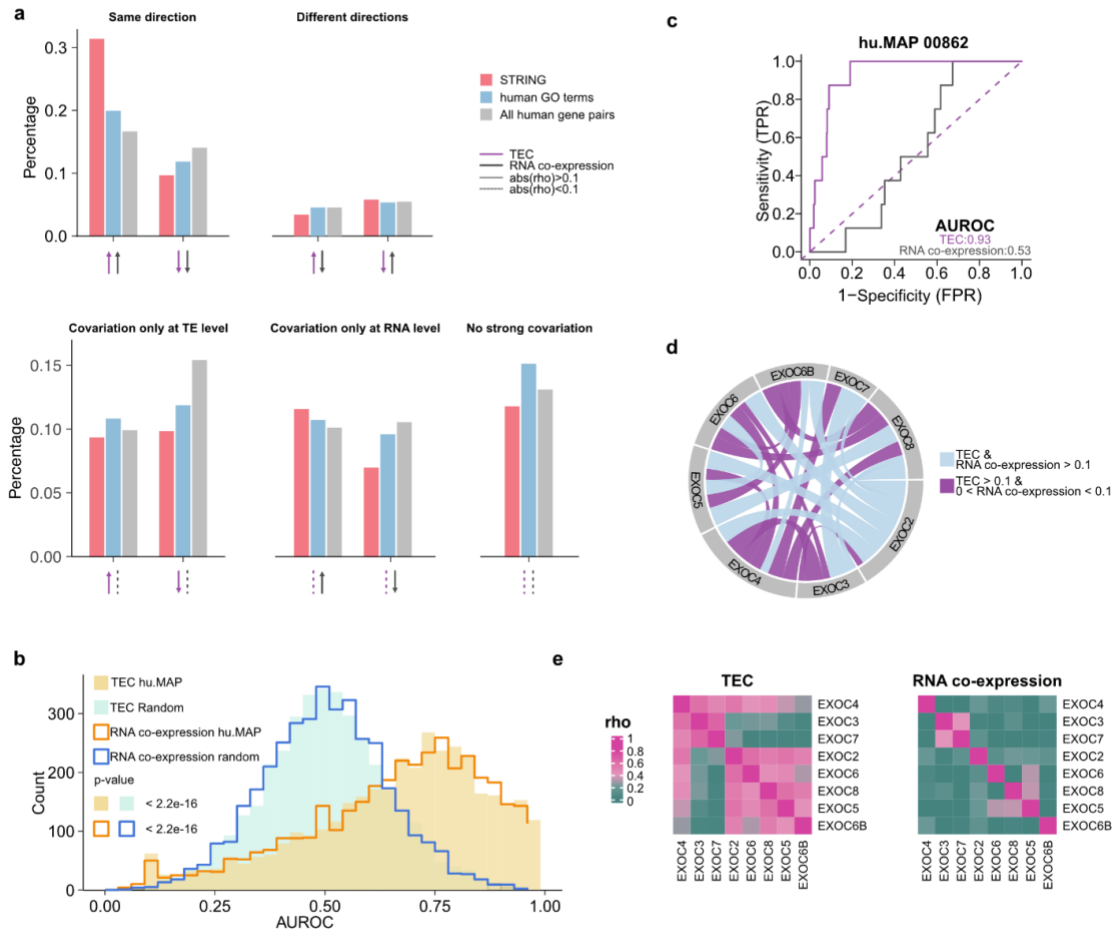
446 glycolysis. **c**, TE and RNA expression of *LRRC28*, glycolytic genes, and transcription factors
447 regulating glycolysis (*FOXK1*, *FOXK2*) across human cell types and tissues. **d**, We used
448 AlphaFold2-Multimer to calculate the binding probabilities between the proteins LRRC28 or
449 LRRC42 and glycolytic proteins (Methods). We evaluated the models with ipTM+pTM (x-axis)
450 and precision of protein-protein interface binding predictions (pDOCKQ; y-axis). We set a
451 threshold of ipTM+pTM > 0.7¹²⁹ and pDOCKQ > 0.23^{130,131} as previously suggested to identify
452 confident binding. **e**, 3D model of binding between LRRC28 and FOXK1. For visualization
453 purposes, we removed residues 1-101 and 370-733 in FOXK1 (pLDDT scores below 50). **f**,
454 Binding probabilities between LRRC28 and transcription factors belonging to the forkhead
455 family¹²⁷. The dashed lines represent ipTM+pTM > 0.7 or pDOCKQ > 0.23.

456 **Genes with positive TEC are more likely to physically interact**

457 The predicted binding between LRRC28 and FOXK1 suggests the utility of TEC to reveal physical
458 interactions between proteins. Proteins that physically interact tend to be co-expressed at the RNA
459 level^{17,23,86}, and many protein complexes are assembled co-translationally⁸⁷, leading us to
460 hypothesize that the TE of interacting proteins may be coordinated across cell types. Specifically,
461 we expect that there should be positive covariation between the TE of interacting proteins to ensure
462 their coordinated production^{40,41}. To test this hypothesis, we categorized gene pairs by whether
463 they display positive or negative similarity in RNA expression or TE across cell types. We
464 observed that nearly one-third of the known pairwise protein-protein interactions (STRING
465 database⁸⁶, only considering the physical interaction subset) exhibited the same direction of
466 similarity (positive rho scores) at both RNA expression and TE levels (Fig. 6a). Compared to all
467 possible pairs (124,322,500), or those with the same biological function (6,492,564), physically
468 interacting pairs of proteins (1,030,794) were substantially enriched for positive similarity of TE
469 and RNA expression patterns (Fig. 6a; chi-square test $p < 2.2 \times 10^{-16}$ and 1.88-fold enrichment
470 compared to all pairs; table S18). Additionally, we found that negative rho values were
471 significantly depleted in protein-protein interactions compared to gene pairs derived from GO
472 terms (Fig. 6a). Though we found enrichment of gene pairs only at RNA expression level, this may
473 be due to neighboring genes being frequently coexpressed (ExtendedDataFig. 20^{27,28}). This result
474 aligns with the notion that genes with the same function can be regulated in opposite directions, as
475 indicated by negative rho values, in contrast to physically interacting proteins^{88,89}.

476 We then examined whether these patterns generalize to the higher-order organization of protein
477 complexes. We observed protein complexes (as defined by hu.MAP⁹⁰) displayed positive TEC and
478 RNA co-expression (Fig. 6b; Methods). Noticeably, while proteins within the same complex
479 generally exhibited similar positive patterns in both TEC and RNA co-expression, certain
480 interactions within protein complexes were particularly evident only at the TE level (Fig. 6c-e).
481 For instance, members of the exocyst complex showed a strong positive TEC but not RNA co-
482 expression (Fig. 6c-e). The exocyst complex consists of eight subunits in equal stoichiometry,
483 forming two stable four-subunit modules^{91,92}. Several known exocyst-binding partners are not
484 required for its assembly and stability, indicating that the molecular details are still unclear⁹¹. Our
485 finding suggests that translational regulation may play a role in maintaining the proper
486 stoichiometry of the exocyst complex. In summary, physically interacting proteins are likely to
487 have positive TEC in addition to positive RNA co-expression profiles. The positive correlation in

488 RNA abundance and TE among physically interacting proteins may reflect an evolutionary
 489 pressure to efficiently utilize energy resources^{40,41,93}.



490

491 **Fig. 6 | Physically interacting proteins display TEC.** **a**, Solid lines indicate gene pairs with
 492 absolute rho greater than 0.1, while dashed lines represent those with absolute rho less than 0.1.
 493 Number of pairs of genes among three sets (physical interaction-red; shared function-blue; all
 494 genes-gray) categorized based on the direction of correlation. **b**, The distribution AUROC
 495 calculated with either TEC or RNA co-expression for 3,755 hu.MAP terms (Methods). The
 496 distribution was compared to AUROC for each term that is randomly assigned genes with size
 497 matched to the original hu.MAP term. P-values were calculated using a two-sided Wilcoxon test.
 498 **c**, AUROC plot for hu.MAP term 00862, which includes eight genes within the exocyst complex.
 499 **d**, Connections represent gene pairs with rho scores above 0.1. Purple lines indicate pairs
 500 connected at TE level alone, while blue lines depict those at both the RNA co-expression and TE
 501 levels. **e**, Heatmaps display the rho calculated among genes at the TE (left) and RNA expression

502 levels (right).

503 **DISCUSSION**

504 In this study, we analyzed thousands of matched ribosome profiling and RNA sequencing
505 experiments across diverse human and mouse cell lines and tissues to quantify TE. A particular
506 challenge in this effort was inadequate metadata associated with these experiments, which hampers
507 their reuse. A particularly recurrent issue was inconsistencies in cell line identification
508 (supplementary text). Additionally, metadata matching of RNA-seq and ribosome profiling data is
509 necessary to quantify TE, yet this information is missing in current databases. To address these
510 issues, we conducted a manual curation process. Given that the analyzed experiments were
511 predominantly described in peer-reviewed publications, we anticipated the publicly accessible data
512 would be of sufficient quality for large-scale analyses. However, more than 20% of the human and
513 mouse experiments did not meet fundamental quality control criteria, such as ribosome footprints
514 arising from coding regions and adequate transcript coverage, thereby deeming these studies
515 unsuitable for further analyses. Our findings point to a pressing need for stricter data quality
516 standards and more comprehensive, structured metadata in genomic databases.

517 We made several advances including the selection of RPF read lengths, data normalization, and
518 estimation of TE (supplementary text). TE is typically defined as a log ratio of read counts from
519 ribosome profiling and RNA expression measurements which often leads to spurious correlations
520 between TE and RNA levels⁵⁵. We instead employed a compositional data analysis framework for
521 both ribosome profiling and RNA-seq^{50,56,94}, allowing for a more accurate estimation of TE as
522 evidenced by improved correlation of these values with corresponding protein abundance. In this
523 study, we used the term “translation efficiency” consistent with its established use in prior
524 literature. Recent work has suggested that ribosome occupancy normalized for mRNA abundance
525 may not directly indicate the efficiency of protein synthesis at least in the context of reporter
526 constructs⁹⁵. While there are mechanisms that lead to a decoupling between ribosome density and
527 the rate of protein synthesis, our work and others indicate that TE as defined here is significantly
528 correlated with protein abundance and synthesis rates for endogenous transcripts⁹⁶.

529 In this study, we introduce the concept of translation efficiency covariation (TEC) which quantifies
530 the similarity of translation efficiency patterns across cell types. Among orthologous gene pairs,
531 RNA co-expression relationships were shown to be conserved across evolution¹¹. Our analyses
532 demonstrated that covariation patterns of TE are also globally conserved between, highlighting the
533 functional relevance of these patterns. Future research leveraging network level conservation
534 metrics could provide further insights into TEC and RNA co-expression networks. Specifically,
535 identifying conserved and divergent subnetwork properties between TE and RNA co-expression
536 networks could elucidate specific regulatory interactions.

537 RNA co-expression among genes is known to be associated with shared functions⁹⁻¹¹. Our analysis
538 indicates that TEC is also informative regarding gene function (Table 1; supplementary text).
539 Interestingly, while for a given transcript, average RNA expression and TE across cell types are
540 only weakly correlated, genes with particular biological functions display highly coordinated
541 patterns of both RNA expression and translation efficiency. This coordination may enhance

542 cellular energy conservation and responsiveness to environmental cues.

543 In addition, TEC revealed unique insights into protein function that elude RNA or protein-based
544 analyses. A notable example is the covariation of TE between *LRRC28* and glycolytic genes,
545 whose RNAs are not co-expressed. We discovered a high confidence predicted interaction between
546 *LRRC28* and *FOKK1*, the key transcription factor controlling glycolytic enzyme expression.
547 Although *LRRC28* is down-regulated in several cancers compared to normal tissues⁹⁷⁻⁹⁹, the
548 functional relevance, if any, remains unknown. These patterns were also not easily detectable at
549 the protein-level as *LRRC28* is absent from most proteomic databases such as PAXdb and
550 ProteomeHD^{23,100}. Taken together, these findings emphasize the unique insights provided by TEC
551 that escape RNA and protein co-expression analyses.

552 TEC between *LRRC28* and its potential physically interacting partner prompted us to
553 systematically analyze the similarity of translation efficiency patterns across protein complexes.
554 We found a significant enrichment of positive TEC between physically interacting protein pairs,
555 establishing that physically interacting proteins often exhibit coordinated translation efficiencies
556 across different cell types. This coordination may facilitate the co-translational assembly⁸⁷ of
557 certain protein complexes and contribute to their stoichiometric production. Such RNA and
558 translation level coordination between physically interacting proteins likely enhances the
559 efficiency of complex formation and optimizes the energetic costs associated with these processes.
560 This optimization is particularly advantageous given that protein biosynthesis is the largest
561 consumer of energy during cellular proliferation^{41,45,93}.

562 It is important to acknowledge several limitations in our study that may impact the accuracy of TE
563 calculations. First, the limited number of samples available for certain cell lines may lead to less
564 accurate estimates of the translation for those cell types. Second, we only considered a
565 representative transcript¹⁰² for each gene based on criteria such as conservation, structure, and
566 functional domains. Mapping RPFs to multiple isoforms of a single gene presents challenges due
567 to the inherently short length of RPFs. This simplification may confound results for genes that
568 have multiple isoforms with distinct expression patterns. In summary, our analyses reveal TEC is
569 informative in uncovering gene functions, is conserved between humans and mice, and suggests
570 simultaneous coordination of both RNA expression and translation among physically interacting
571 proteins, establishing translation efficiency covariation as a fundamental organizing principle of
572 mammalian transcriptomes.

573 **ACKNOWLEDGMENTS**

574 We thank all contributions to metadata curation: Hansel Chiang, Ashley Hoffman, Tori Tonn, Alia
575 Segura, Charisma Tante, Eric Vasquez, and Liaoyi Xu. We also thank Dr. Vighnesh Ghatpande
576 and Victoria D. Chapman for generating the KO cell lines and assisting with the preparation of the
577 sequencing library. We appreciate Dr. Milad Miladi for providing critical feedback, and the
578 original text in this paper was written by the authors. A LLM was used to suggest edits for clarity
579 and grammar¹⁴⁰. The authors acknowledge the Texas Advanced Computing Center (TACC) at The
580 University of Texas at Austin for providing high-performance computing and storage resources
581 that have contributed to the research results reported within this paper. URL:

582 <http://www.tacc.utexas.edu>.

583 Research reported in this publication was supported in part by the National Institute Of General
584 Medical Sciences of the National Institutes of Health under Award Number R35GM150667 (CC).
585 This work was also supported by the National Institutes of Health grant [HD110096], and the
586 Welch Foundation grant [F-2027-20230405] (C.C.). C.C. was a CPRIT Scholar in Cancer
587 Research supported by CPRIT Grant [RR180042].

588 **AUTHOR CONTRIBUTIONS**

589 Y.L., I.H., and C.C. co-wrote the original manuscript. Y.L. and I.H. generated the figures for the
590 manuscript. H.O., M.G., and J.C. downloaded all the data from GEO and processed raw
591 sequencing data. Y.L. and C.C. developed the translation efficiency calculation pipeline. J.C. and
592 C.C. designed and implemented the winsorization method. Y.Z. performed the deduplication
593 comparison. Y.L. and C.C. developed the translation efficiency covariation analysis and function
594 prediction pipelines. Y.L., K.Q., H.O. performed the quality control analysis for all sequencing
595 data. Y.L. carried out covariation analysis, gene function prediction, and AlphaFold2 analysis. I.H.
596 conducted the RBP analysis. L.P., J.W., D.Z., and V.A. assessed the quality of TE measurements
597 by developing machine learning approaches. H.O., J.W., D.Z., V.A., Q.Z., and E.S.C. provided
598 suggestions for the manuscript. Y.L., Q.Z., and E.S.C. conducted literature search to evaluate gene
599 function predictions. I.H., S.R., and D.P. performed all experiments. C.C. provided study
600 oversight, conceptualized the study and acquired funding. All authors approved the final
601 manuscript.

602 **DATA AVAILABILITY**

603 Metadata about RiboBase can be found in Supplementary table S1. Ribo files for the HeLa cell
604 line are accessible at <https://zenodo.org/records/10594392>. Full TEC and RNA co-expression
605 matrices are accessible via Zenodo repository at: <https://zenodo.org/uploads/10373032>. A
606 RiboFlow configuration file and processed ribo files can be accessed at
607 <https://zenodo.org/uploads/11388478>. Sequencing data and ribo files for the RBP knockout
608 experiments are available on GEO GSE269734. Data will be publicly released upon successful
609 review of this article.

610 **CODE AVAILABILITY**

611 The code used in the study is available at https://github.com/CenikLab/TE_model/tree/main. Code
612 will be publicly released upon successful review of this article.

613 **DECLARATION OF INTERESTS**

614 D.Z., J.W. and V.A. are employees of Sanofi and may hold shares and/or stock options in the
615 company. H.O. is an employee of Sail Biomedicines.

616 **METHODS**

617 **Acquisition and curation of ribosome profiling data**

618 We used keyword search (“ribosome profiling”, “riboseq”, “ribo-seq”, “translation”, “ribo”,
619 “ribosome protected footprint”) to determine studies that may employ ribosome profiling in their
620 experimental design, from the Gene Expression Omnibus (GEO) database, with a cutoff date of
621 January 1, 2022. Search results were manually inspected and studies containing ribosome profiling
622 data were kept. Organism, cell line, publication, and short read archive (SRA) identifiers were
623 obtained by automatically parsing the GEO pages of the corresponding study and sample. There
624 was no dedicated experiment-type field for ribosome profiling experiments in GEO. Therefore we
625 determined the experiment type (ribosome profiling, RNA-Seq, or other) of each sample by
626 manually inspecting the GEO metadata and the associated publication of the study. Typically,
627 ribosome profiling samples were indicated in GEO using one of the following terms: “ribosome
628 protected footprints”, “ribo-seq”, and “ribosome profiling” in various parts of the metadata such
629 as title, extraction protocol, and library strategy. If there were RNA-Seq samples in the same study,
630 they were matched with ribosome profiling experiments, where available, after inspecting the
631 sample names, metadata, and the publication of the study.

632 Adapters are commonly observed on the 3’ end of sequencing reads in ribosome profiling
633 experiments, a consequence of the inherently short length of RPFs. If the 3’ adapter sequence was
634 listed in GEO, we extracted it as part of the manual data curation process. If this sequence was
635 unavailable, we attempted to determine it from the corresponding publication of the study. If no
636 explicit sequence was available, we computationally analyzed the sequencing reads and searched
637 for commonly used adapters which are CTGTAGGCACCATCAAT,
638 AAGATCGGAAGAGCACACGTCT,
639 AGATCGGAAGAGCACACGTCTGAACTCCAGTCAC, TGGAATTCTCGGGTGCCAAGG
640 and AAAAAAAAAA. If any of these adapters were found in at least 50% of the reads, we used
641 the detected sequence as the 3’ adapter. If no match was found, we removed the first 25 nucleotides
642 of the reads anchored 6 mers and tried to extend them. If any of these extensions reached 10
643 nucleotides and were still detected in at least 50% of the reads, we took the highest matching
644 sequence as the 3’ adapter. On the other hand, for sequencing reads from SRA having a length of
645 less than 35 nucleotides, we assumed the 3’ adapters had already been removed. Detailed code can
646 be accessed from: https://github.com/RiboBase/snakescale/blob/main/scripts/guess_adapters.py.

647 RiboBase was pre-populated after mining GEO. Then data curators were assigned specific studies
648 and used the web-based interface to access the database. Each study was curated independently by
649 at least two people. In case of disagreements, an additional experienced scientist inspected the
650 corresponding studies and publications to make the final decision. We supplemented any missing
651 metadata from GEO by checking the corresponding publications to ensure completeness. The
652 result of this data curation process with information such as cell line, organism, and matched RNA-
653 seq can be found in table S1, which forms the metadata backbone of RiboBase.

654 **Ribosome profiling and RNA-seq data processing**

655 For each selected study in GEO, ribosome profiling and matching RNA-Seq reads (where
656 available) were downloaded, from SRA, using the SRA-Tools version 2.9.1¹⁰³, in FASTQ format
657 using their accession numbers. FASTQ files were processed using RiboFlow⁵¹ where parameters
658 were determined using the metadata in RiboBase. The reference files for human and mouse

659 transcriptomes, annotations, and non-coding RNA sequences are available at
660 https://github.com/RiboBase/reference_homo-sapiens and
661 https://github.com/RiboBase/reference_mus-musculus, respectively. Briefly, the 3' adapters of the
662 ribosome profiling reads were trimmed using Cutadapt version 1.18¹⁰⁴ and reads having lengths
663 between 15 and 40 nucleotides were kept. Then, reads were aligned against noncoding RNAs, and
664 unaligned reads were kept. Next, reads were aligned against transcriptome reference, and
665 alignments having mapping quality score above 20 were kept. Reads having the same length and
666 mapping to the same transcriptome position were collapsed, which we refer to as “PCR
667 deduplication”. In the final step, we compiled the alignments into ribo files using RiboPy⁵¹. All
668 alignment steps used bowtie2 version 2.3.4.3¹⁰⁵. For each sample, we also performed the same run
669 without the PCR deduplication step. We developed a pipeline, Snakescale, available at
670 <https://github.com/RiboBase/snakescale>, to automate the entire process from downloading the data
671 from SRA to generating the ribo files. Snakescale went over the selected list of studies and
672 obtained their metadata from Ribobase, downloaded the sequencing data from SRA, generated
673 Riboflow parameters file, and ran Riboflow to generate the ribo files. Examples of non-
674 deduplicated ribo files for the HeLa cell line can be accessed at
675 <https://zenodo.org/records/10594392>¹⁰⁶.

676 To visualize the length distribution of the RPFs, we applied the scale function (z-score) in R to
677 normalize the count of RPFs mapped to CDS regions with PCR-deduplicated ribosome profiling
678 data. Subsequently, we plotted the distribution of these normalized RPFs using the heatmap (Fig.
679 1c; ExtendedDataFig. 2).

680 **Determination of cutoff for RPF lengths and quantification of ribosome occupancy**

681 Ribosome profiling experiments employ a range of ribonucleases including RNase I, RNase A,
682 RNase T1, and MNase (i.e., micrococcal nuclease) (S7). These different enzymes lead to variable
683 RPF lengths^{53,107–109}. To ensure that we retain high-quality RPFs for further analyses, we
684 implemented a dynamic extraction module that automatically selected lower and higher boundaries
685 of RPFs for each sample. Initially, we determined the first RPF length, ranging from 21 to 40
686 nucleotides, that contained the highest number of CDS mapping reads. Then, we examined the two
687 positions adjacent to this selected position. The extension of the position was carried out on either
688 side to include a higher number of CDS-aligned reads. This extension process was repeated until
689 it encompassed at least 85% of the total CDS reads within the 21 to 40 nucleotides range
690 (ExtendedDataFig. 3a). The final two positions identified were designated as the lower and upper
691 boundaries. If these boundaries extended to either 21 or 40 nucleotides without including a
692 sufficient number of reads, then 21 or 40 nucleotides, respectively, were set as the final boundaries.
693 This approach was employed to establish the RPF cutoffs for each sample.

694 **Transcript coverage and quality control for ribosome profiling data**

695 We performed quality control using RPFs that were deduplicated based on the length and position
696 (PCR-deduplication) ribo files (Fig. 1d-e). We set two cutoffs for ribosome profiling quality
697 control. We required that on average each nucleotide of the transcript should be covered at least
698 0.1 times (0.1X). Coverage was calculated with the formula:

$$\begin{aligned} 699 \quad & \textit{coverage} = \textit{total nucleotides from reads mapped to transcripts} \\ 700 \quad & \quad \quad \quad / \textit{total length of transcripts} \end{aligned}$$

701 Additionally, samples with CDS mapping read percentage of 70% or higher were retained for
702 subsequent analysis.

703 To assess the pattern of three nucleotide periodicity that is typically associated with ribosome
704 profiling experiments, we first selected the length of RPFs with the highest number of counts from
705 the PCR deduplicated ribo files. We then assigned all CDS mapping reads to one of three coding
706 frames based on the position of their 5' end. We aggregated the results for all genes for each
707 sample. To facilitate comparison, we reordered the counts for each position of the three nucleotide
708 periodicity from highest to lowest and converted these counts into percentages for each sample.

709 We initially classified samples based on the differences between positions 1 and 2. We identified
710 Group 1 by selecting samples where the difference did not exceed the 10th percentile of these
711 differences between positions 1 and 2. For the remaining samples, we further classified them based
712 on the differences between positions 2 and 3. Similarly, samples that did not exceed the 10th
713 percentile of these differences between positions 2 and 3 among remaining samples were classified
714 to Group 3, while the rest samples were Group 2. We further summarized the samples based on
715 their QC status.

716 We classified samples from Group 1 as exhibiting three-nucleotide periodicity. The percentage of
717 samples following three-nucleotide periodicity was calculated by dividing the number of Group 1
718 samples by the total number of samples across all three groups.

719 **PCR and Unique Molecular Identifiers (UMIs) deduplication comparison**

720 We selected eight ribosome profiling experiments that incorporated UMIs into the sequence library
721 preparation to assess the impact of different deduplication methods. Specifically, these samples
722 are GSM4282032, GSM4282033, and GSM4282034 from GSE144140¹¹⁰; GSM3168387,
723 GSM3168389, GSM3168390 from GSE115162¹⁰⁸; and GSM4798525, GSM4798526 from
724 GSE158374³⁴. We processed the data using Riboflow, applying three different deduplication
725 methods: non-deduplication, PCR deduplication, and UMI deduplication. The yaml files are
726 available at https://github.com/CenikLab/TE_model/tree/main/riboflow_scr. The RPF length
727 cutoffs for samples from GSE144140 and GSE115162 are listed in table S6. Since GSE158374 is
728 not currently included in RiboBase, we manually performed the dynamic module and selected 28
729 to 32 as the RPF cutoff for this study.

730 **Winsorization of CDS mapping read counts**

731 To address the issue of reduced usable reads resulting from PCR deduplication (supplementary
732 text), we employed a winsorization method, which was previously proposed for tackling this
733 problem^{40,111}. For each gene's CDS region, we obtained the distribution of non-deduplicated
734 nucleotide counts and calculated the 99.5th percentile value. This calculation was based on reads
735 whose lengths fell within the RPF range determined by the RPF boundary selection function. RPF
736 counts that exceed the 99.5th percentile were capped to the value corresponding to the 99.5th

737 percentile. This method was designed to mitigate the impact of outlier values that might arise due
738 to disproportionate amplification during the PCR process⁴⁰.

739 **Gene filtering and normalization for ribosome profiling and RNA-seq**

740 RNA-seq experiments in RiboBase utilized several different strategies to enrich mRNAs. The two
741 most common approaches were the depletion of ribosomal RNAs and the enrichment of transcripts
742 by polyA-tail selection. This difference leads to dramatically different quantification of a subset
743 of genes that lack polyA-tails (e.g. histone genes, ExtendedDataFig. 6c). Hence, we removed 166
744 human and 51 mouse genes identified as lacking polyA tails (table S8-9)^{112,113}.

745 We normalized both PCR-deduplicated RNA-seq data and winsorized non-deduplicated ribosome
746 profiling data with counts per million (CPM) after removing the genes without polyA-tails. Genes
747 with CPM greater than one in over 70% of the total samples in both RNA-seq and ribosome
748 profiling for either human or mouse were included in further analyses. 11,149 human and 11,434
749 mouse genes were retained using this approach. We have summed the counts of all polyA genes
750 that were filtered out and grouped them under 'others' in the count table.

751 **Validation of manual curation and quality control by matching between RNA-seq and** 752 **ribosome profiling from RiboBase**

753 We assessed the manual matching of ribosome profiling (winsorization) and RNA-seq (PCR
754 deduplication) data in RiboBase by establishing a matching score for the samples that successfully
755 passed quality control (transcript coverage > 0.1X and CDS percentage > 70% with PCR-
756 deduplicated ribosome profiling data). We calculated the coefficient of determination (R^2) using
757 the Centered Log Ratio (CLR) transformed gene counts. This was done for each ribosome profiling
758 sample against all corresponding RNA-seq samples within the same study. Subsequently, for each
759 ribosome profiling sample, we calculated the difference between the R^2 of its matching pair from
760 RiboBase and the mean R^2 of the non-matching pairs within the same study. The difference was
761 defined as the matching score.

762 To remove poorly matched samples in both human and mouse datasets, we established a cutoff
763 based on the R^2 from the matched ribosome profiling and RNA-seq data in RiboBase. Any sample
764 with an R^2 lower than 0.188 in either human or mouse, which is $Q1 - 1.5 * IQR$ of mouse R^2
765 distribution, was considered a poor match and consequently excluded from further analysis
766 (ExtendedDataFig. 6b). Finally 1,054 human and 835 mouse ribosome profiling experiments with
767 their matched RNA-seq were used for TE calculation.

768 **Translation Efficiency (TE) calculation**

769 CLR normalized counts from PCR-deduplicated RNA-seq and winsorized non-deduplicated
770 ribosome profiling were used to calculate TE with compositional linear regression^{49,94,114}. In our
771 linear regression approach, ribosome profiling data served as the dependent variable, while the
772 corresponding RNA-seq data provided the explanatory variable. The first step involved
773 transforming the gene count, which includes 'others', into CLR normalized compositional vectors.
774 Given the constraints of count data within a simplex, a further transformation from CLR to

775 Isometric Log Ratio (ILR) was necessary for linear regression⁴⁹. This transformation is crucial as
776 it allows the compositional data to be decomposed into an array of uncorrelated variables while
777 preserving relative proportions. The ILR transformation projects the original data onto a set of
778 orthonormal basis vectors derived from the Aitchison simplex. Then the linear regression model
779 applied to these transformed variables can be represented as:

$$780 \quad Y = b + B * X$$

781 Where Y is the ILR-transformed ribosome profiling data and X is the ILR-transformed RNA-seq
782 data. The model assumes a normal distribution:

$$783 \quad Y \sim N^{(D-1)}(Y, \Sigma\varepsilon)$$

784 Where $\Sigma\varepsilon$ represents the residual variances. These residuals were then extracted from each sample
785 and reconverted to CLR coordinates which are used as the definition of TE for each gene in each
786 sample. Finally, we averaged TE for different cell lines and tissues (Fig. 2b, ExtendedDataFig. 7),
787 and reported the TE in table S10-11. The scripts to generate TE are available at
788 https://github.com/CenikLab/TE_model.

789 **Correlation between translation efficiency and protein abundance**

790 We assessed the correlation between TE and protein abundance from seven human cell lines
791 (A549, HEK293, HeLa, HepG2, K562, MCF7, and U2OS). The protein measurements were
792 obtained from PAXdb¹⁰⁰. 9924 genes were shared between our TE and the protein abundance data.
793 We calculated the Spearman correlation coefficient for each cell line using the 'stats' package in
794 R to evaluate the relationship between TE and protein abundance.

795 **Conservation of translation efficiency between orthologous genes from human and mouse**

796 Orthologous genes between human and mouse were identified using the 'orthogene' package from
797 Bioconductor¹¹⁵ using the parameters 'standardise_genes=TRUE,
798 method_all_genes="homologene", non121_strategy="keep_both_species"'. A single human gene
799 could correspond to multiple mouse orthologs or vice versa. To maintain all one-to-many matches
800 in our analysis, each correspondence is represented by multiple rows in our table (if a human gene
801 'A' is orthologous to mouse genes 'B' and 'C', we generate two separate rows: 'A-B' and 'A-C').
802 Human genes lacking corresponding mouse orthologs were excluded or vice versa. As a result, a
803 total of 9,194 gene pairs were identified as orthologous between human and mouse (table S12)

804 To capture the variability in TE and mRNA expression between orthologous genes in human and
805 mouse, we measured the standard deviation using the metric standard deviation (msd) function
806 from the 'compositions' package in R¹¹⁶. We observed a negative Spearman correlation coefficient
807 between msd of TE and mean TE, as well as msd of RNA expression and mean RNA expression,
808 in both species. To address the dependency between msd and mean values, we conducted a partial
809 correlation analysis. For example, we adjusted the human msd values using the mean TE from
810 both human and mouse with the 'pcor.test' function from the 'ppcor' package¹¹⁷.

811 GO term analysis was performed using FuncAssociate 3.0, accessible at

812 <http://llama.mshri.on.ca/funcassociate/>¹¹⁸. For this analysis, we set either 9,194 mouse or 9,189
813 human orthologous genes as the background. We generated association files for these genes with
814 the December 4, 2022 version of human or mouse GO terms. In the human or mouse association
815 file, we only kept those GO terms containing at least 10 genes for further analysis.

816 **Assessment of methods for calculating genes' similarity with ribosome occupancy data**

817 We used eight commonly used methods to quantify the similarity of ribosome occupancy across
818 cell types for all pairs of 11,149 human or 11,434 mouse genes in RiboBase.

819 Method 1 - CPM-normalized ribosome footprint counts were used to calculate the Pearson
820 correlation coefficient as implemented in the stats R package.

821 Method 2 - Quantile-normalized (customized Python script) ribosome footprint counts were used
822 to calculate the Pearson correlation coefficient.

823 Method 3 - Ranking of ribosome footprint counts was used to calculate the Spearman correlation
824 coefficient as implemented in the stats R package.

825 Method 4 - CLR-normalized ribosome footprint counts were used to calculate the proportionality
826 (rho scores) between genes as implemented in the propr package with lr2rho function⁵⁰.

827 Method 5 - CPM-normalized ribosome footprint counts were used to calculate the similarity
828 between genes with a decision tree-based method as implemented in the treeClust package^{23,119}.
829 We applied the 'treeClust.dist' function with a dissimilarity specifier set to d.num=2.

830 Method 6 - Quantile-normalized ribosome footprint counts were used to calculate the similarity
831 between genes with the decision tree-based method.

832 Method 7 - CPM-normalized ribosome footprint counts were used to calculate gene similarity with
833 the generalized least squares (GLS) method¹²⁰.

834 Method 8 - Quantile-normalized ribosome footprint counts were used to calculate gene similarity
835 with the GLS method.

836 We compared these eight ribosome occupancy similarity matrices to determine the most effective
837 method for constructing gene relationships with respect to biological functions. This assessment
838 employed the guilt by association principle to ascertain the functional coherence within a gene
839 matrix, determining if genes associated with a particular biological function (GO terms¹²¹, TOP
840 mRNAs¹²²) exhibit similar expression patterns and network interactions¹²³.

841 The complete ontology was sourced from the Gene Ontology website, with the files
842 goa_human.gpad.gz and mgi.gpad.gz, generated on December 4, 2022¹²¹. The annotation of Gene
843 Ontology terms was accomplished with the aid of the org.Hs.eg.db and org.Mm.eg.db R
844 packages^{124,125}. We restricted the selection of GO terms to those associated with the 11,149 human
845 and 11,434 mouse genes that had passed gene filtering. We used GO terms associated with at least
846 10 but less than 1,000 genes for evaluation, yielding a total of 2,989 human and 3,340 mouse GO

847 terms.

848 We then employed the neighbor-voting algorithm to assess the covariations of ribosome
849 occupancy among genes from the same GO term with AUROC¹²³. Specifically, we first converted
850 the similarity scores to absolute values. Then we extracted genes associated with a specific
851 function and implemented the leave-one-out cross-validation method. For this analysis, we
852 iteratively masked one gene at a time, treating it as if it did not belong to the function. In each
853 iteration, we calculated the total sum of similarity scores from all genes not belonging to the
854 function to all the remaining genes within the function. We normalized the sum of similarity scores
855 for each gene against the sum of similarity scores for that gene with all genes. After normalization,
856 we converted these normalized similarity scores into rankings. We retained the rankings only for
857 genes that belong to this specified functional property. Finally, we computed the AUROC for all
858 genes within this functional property based on these rankings. A detailed script for genes'
859 functional similarity pattern analysis can be found:
860 https://github.com/CenikLab/TE_model/blob/main/other_scr/benchmarking.R.

861 **RNA co-expression and translation efficiency covariation**

862 We introduce the concept of TEC, which employs a compositional data analysis approach^{49,50} to
863 quantify the similarity patterns of TE across various cell and tissue sources, as described in Method
864 4 above. The proportionality scores were calculated with the following formula from the propr
865 package with lr2rho function⁵⁰:

$$866 \quad \text{Rho}(A_i, A_j) = 1 - \text{var}(A_i - A_j) / (\text{var}(A_i) + \text{var}(A_j))$$

867 Where A_i and A_j represent TE values for genes i and j from the TE matrix A .

868 In this study, the TEC was calculated with 77 human cell lines for 11,149 genes or 68 mouse cell
869 lines for 11,434 genes. The proportionality coefficients (rho scores) generated from this method
870 range from -1 to 1. Full TEC and RNA co-expression matrices are accessible via Zenodo repository
871 at: <https://zenodo.org/uploads/10373032>.

872 **Evaluation of the ability of TEC to predict novel gene functions**

873 We compared the AUROC between an older version of GO terms (January 1, 2021) to the newer
874 version of GO terms (December 4, 2022) to identify genes that had been newly added to from the
875 GO terms in this timeframe. GO terms were downloaded and filtered to include only those terms
876 containing between 10 and 1,000 genes with either human or mouse backgrounds (11,149 human
877 genes or 11,434 mouse genes). We selected 184 human and 238 mouse GO terms from the older
878 version that demonstrated high TEC similarity (AUROC > 0.8) among genes within the same term
879 for predicting novel gene functions. We first converted the rho scores for TEC between gene pairs
880 to absolute values. For genes not currently included in the GO terms, we calculated the sum of rho
881 for each gene relative to all genes within the term, based on either TE or mRNA expression levels.
882 We then normalized these rho sums for each gene against the total rho sum of that gene across all
883 11,149 human genes or 11,434 mouse genes. These normalized values were converted into ranking
884 percentages to reflect the likelihood of these genes being associated with the respective GO term.

885 Finally, we identified the top-ranking genes as potentially new additions and cross-validated them
886 with the newer version of the GO terms to confirm our predictions.

887 **Prediction of novel gene functions with TEC**

888 We analyzed 243 human and 310 mouse GO terms as of December 4, 2022, which demonstrated
889 high similarity patterns between genes in TE level (AUROC > 0.8) to predict novel gene functions.
890 Absolute TEC rho scores served as the input for biological function prediction (GO terms). The
891 prediction method followed the same protocol as our previous evaluations of TEC's ability to
892 predict novel gene functions. However, we added a filter step: a newly predicted gene was retained
893 only if its average rho score with other genes within the same term exceeded the overall average
894 rho score for all existing genes in that term. This prediction analysis was performed using a custom
895 script that can be found at
896 https://github.com/CenikLab/TE_model/blob/main/other_scr/prediction.R.

897 **Computational evaluation of the interaction between LRRC28, glycolytic proteins, and** 898 **proteins from forkhead TF family**

899 We computed the pair-wise interaction probabilities between LRRC28 or LRRC42 and glycolytic
900 proteins (HK1, HK2, PFKL, PFKM, PFKP, TPI1, PGK1, ENO1, ENO2, and PKM) with
901 AlphaFold2-Multimer 2.3.0^{83,126}. In addition, we also calculated pairwise interaction probabilities
902 for LRRC28 with 35 proteins from the forkhead transcription TF family¹²⁷. We extracted the
903 canonical amino acid sequence for each gene from UniPort¹²⁸ as the input file. We set 0.7 as the
904 cutoff of ipTM+pTM as a high-confidence protein structure and binding probability cutoff¹²⁹. We
905 then evaluated the interfaces predicted by AlphaFold2-Multimer, using a pDOCKQ score greater
906 than 0.23 as our criterion for reliability^{130,131}.

907 **Benchmarking TEC and RNA co-expression for protein interactions**

908 Using a similar approach to our benchmarking with biological functions, we employed the
909 neighbor-voting algorithm to assess physical protein interactions based on rho scores among genes
910 at either the TE or mRNA expression level. We first kept the non-negative rho between genes and
911 set negative rho to zero. We then analyzed similarity patterns between genes from the same protein
912 complex, downloading from the hu.MAP 2.0 website⁹⁰. In this process, we excluded genes from
913 hu.MAP terms that were not in the 11,149 human gene list, resulting in 8,024 overlapping genes
914 between our list and hu.MAP terms. Furthermore, we removed hu.MAP terms that included fewer
915 than three genes. This filtering process left us with 3,880 hu.MAP terms, among which 3,755
916 contained unique genes.

917 Since proteins within the same complex may not physically interact, we used physical interaction
918 pairs downloaded from the STRING website instead of gene pairs from hu.MAP terms to
919 summarize the interactions in Fig. 6a.

920 **Identification of enriched RNA motifs among genes with high degree of TEC**

921 To reduce bias in motif enrichment analysis that may arise by ribosome footprint mapping to

922 paralogous genes, we removed predicted paralogs from each GO term using Paralog Explorer¹³²
923 (DIOPT score > 1). Then, we enumerated heptamers in each transcript region using the Transite
924 kmer-TSMA method¹³³ with default parameters for each species (human, mouse), transcript region
925 (5' UTR, CDS, 3' UTR), and GO term (selected terms with TE AUROC > 0.7, TE-RNA AUROC
926 difference > 0.2, and number genes after paralog removal >= 12). For mice there were three terms
927 and for humans there were eight. We selected the three mouse terms and top five terms in humans
928 with the highest number of genes and greatest AUROC difference.

929 After counting heptamers with Transite, we selected motifs that had >20 hits among genes in the
930 GO term to address assumptions of uniformity near p-values of 1 for some multiple-test correction
931 methods. Then, we used the Holm method to correct p-values for each species separately, and
932 selected motifs with an adjusted p-value < 0.05. Finally, heptamers were annotated with RBPs
933 included in the Transite¹³³ and oRNAmot databases¹³⁴. For annotation of RBPs in the oRNAmot
934 database, we required that the heptamer have a matrix similarity score¹³⁴ of 0.8 or greater when
935 matching to each RBP position weight matrix. RBP motif hits from other species (*Drosophila*,
936 artificial constructs) were removed from RBP annotations, and the hits to the heptamer of
937 *Drosophila tra2* were annotated as TRA2A for human genes with the term GO: 0140678.

938 eCLIP data for PABPN1, SRSF1, and TRA2A were downloaded from ENCODE¹³⁵ as BED files
939 (K562 and HepG2 cell lines, GRCh38 reference). The BED files for biological replicates were
940 concatenated and peaks that overlapped by at least one base pair were merged with 'bedtools merge
941 -s -c 4,6,7 -o collapse'¹³⁶. The resulting merged peaks were intersected with transcripts in the GO
942 term of interest and an equal number of control transcripts (Gencode v34 GTF). The control
943 transcripts were selected by matching on length and GC content for each transcript region (5' UTR,
944 CDS, 3' UTR) using MatchIt¹³⁷ with default parameters. Because the gene *CARMIL2* in GO term
945 GO:0010592 does not have a 5' UTR, required for matching, we assigned it a dummy 5' UTR
946 with length and GC content equal to the median across all transcripts. The number of eCLIP peaks
947 in the CDS for each RBP were summed for genes in the GO term and control genes.

948 **Identification of RBP-gene pairs with high correlation between RBP RNA expression and** 949 **gene TE**

950 The Pearson correlation coefficient between gene TE and the RNA expression of RBPs from
951 human and mouse⁶⁶ was tested using R stats::cor.test after taking the mean of these values by cell
952 types and tissues. P-values were corrected with the Benjamini-Hochberg procedure, and
953 correlations were deemed significant at a FDR < 0.05. To select RBP candidates for experimental
954 validation, the human and mouse regulons were intersected for each RBP, and those that had more
955 than twenty genes in the intersection and that had a mean TEC > 0.35 between genes in the
956 intersection were chosen.

957 **Generation of RBP knockout cell lines**

958 For cloning the guides required for knockout cell line generation, top two ranked guides were
959 selected from the Brunello library¹³⁸ for each RBPs (table S18). The guides were cloned in
960 LentiCRISPRv2 (Addgene, 52961) as per the protocol¹³⁹ and confirmed by Sanger sequencing.
961 Briefly, for lentiviral production, HEK293T cells were seeded at a density of 1.2 x 10⁶ cells per

962 well in a 6-well plate in OPTI-MEM media supplemented with 5% FBS and 100 mM Sodium
963 Pyruvate, 24 h prior to transfection. Both the cloned gRNA plasmids for each RBPs (700 ng of
964 each transfer plasmid) were co-transfected with the packaging plasmids pMD2.G and psPAX2
965 (Addgene; 12259 and 12260) using Lipofectamine 3000 (Invitrogen) and the virus was collected
966 as per the manufacturer's protocol. For generation of the knockout clones, HEK293T cells were
967 seeded at a density of 5×10^4 cells per well in a 6-well plate in DMEM media supplemented with
968 10% FBS, 24 h prior to infection. Next day, the media was replaced with 1.5 ml of 1:2 diluted
969 lentivirus containing polybrene (8 $\mu\text{g}/\text{mL}$). After 16 h, the lentivirus was replaced with fresh media
970 and, puromycin (2 $\mu\text{g}/\text{mL}$) was added to the cells 48 h after transduction. The selection continued
971 for 5 days followed by a period of recovery for 24 h before harvesting the cells.

972 **Ribosome profiling and RNA sequencing of RBP knockout cell lines**

973 Three million cells for the PARK7, USP42, and VIM knockout cell lines, along with a AAVS1
974 (safe harbor control) knockout line, were plated in three 10 cm^2 dishes. 27 h later cells at ~60%
975 confluency were treated with 100 $\mu\text{g}/\text{mL}$ cycloheximide (CHX) for 10 min at 37 $^\circ\text{C}$, then collected
976 in ice cold PBS with 100 $\mu\text{g}/\text{mL}$ CHX. Cells were spun at 100 x g for 7 min at 4 $^\circ\text{C}$, then flash
977 frozen in liquid nitrogen and stored at -80 $^\circ\text{C}$. Cell pellets were lysed with 400 μL lysis buffer (20
978 mM Tris-HCl pH 7.4, 150 mM NaCl, 5 mM MgCl_2 , 1 mM DTT, 1% Triton-X, 100 $\mu\text{g}/\text{mL}$ CHX,
979 1x protease inhibitor EDTA free) for 10 min on ice. Lysates were clarified by centrifugation at
980 1,300 rpm for 10 min at 4 $^\circ\text{C}$. 40 μL lysate was saved for total RNA extraction, and the rest of the
981 lysate was digested with 7 μL RNaseI for 1 h at 4 $^\circ\text{C}$. Digestion was stopped by adding
982 ribonucleoside vanadyl complex to a final concentration of 20 mM. Digested lysates were then
983 loaded onto 10 mL sucrose cushion (20 mM Tris-HCl pH 7.4, 150 mM NaCl, 5 mM MgCl_2 , 1 mM
984 DTT, 1 M sucrose) and centrifuged at 38,000 rpm for 2.5 h at 4 $^\circ\text{C}$ using a SW41-Ti rotor. The
985 pellet and the total RNA aliquot were both solubilized with 1 mL Trizol, and RNA was purified
986 with the Zymo Direct-zol RNA Miniprep Kit, including DNaseI digestion.

987 *RNA-seq*

988 Quality of total RNA was confirmed with Bioanalyzer RNA Pico. All RIN scores were ≥ 9.8 .
989 Libraries were prepared from 1 μg total RNA using the NEBNext Ultra II RNA Library Prep Kit
990 for Illumina according to manufacturer's protocol and using 8 cycles for PCR.

991 *Ribosome profiling*

992 Ribosome protected fragments (RPFs) were size-selected on a 15% TBE urea gel by
993 electrophoresing at 150 V for 1.5 h. RPFs between 28-32 nt were sliced, using 28 nt and 35 nt
994 markers as a guide. Slices were frozen at -20 $^\circ\text{C}$ for 1 h, crushed with pestles, and the RPFs were
995 eluted in gel extraction buffer (300 mM sodium acetate pH 5.5, 5 mM MgCl_2) by rotating overnight
996 at room temperature. Eluates were passed through Costar Spin-X filter tubes at 12,000 x g for 1
997 min 30 s. Then 1 μL 1 M MgCl_2 , 2.5 μL GlycoBlue, and 1 mL ethanol were added and the RPFs
998 precipitated for two days at -20 $^\circ\text{C}$. Pellets were dried and resuspended in 16 μL water.

999 Libraries were generated from 8 μL RPF eluate using the Diagenode D-Plex Small RNA Kit with
1000 minor modifications: in the 3' dephosphorylation step 0.5 μL T4 PNK was supplemented and

1001 incubated for 25 min. The RTPM reverse transcription primer was used and 8 cycles were
1002 performed for PCR. Libraries were quantified by Bioanalyzer High Sensitivity DNA Kit, pooled
1003 equimolar according to the quantity of the peak for libraries with full-length inserts (~204 nt), and
1004 cleaned up with 1.8X AMPure XP beads. Adapter dimers and empty libraries were removed by
1005 size-selection on a 12% TBE PAGE gel, followed by extraction with the crush and soak method,
1006 and final libraries were resuspended in 20 μ L water.

1007 **Ribosome profiling and RNA sequencing analysis for RBP knockouts**

1008 Analysis was conducted using RiboFlow v0.0.1 with deduplication of both Ribo-seq and RNA-
1009 seq data. A RiboFlow configuration file and processed ribo files can be accessed at
1010 <https://zenodo.org/uploads/11388478>.

1011 We used edgeR to measure RBP KO effects on 1) RNA abundance and 2) gene TE. To do this, we
1012 respectively modeled 1) RNA-seq counts of a specific RBP KO line to that of the other two RBP
1013 KO lines; and 2) Ribo-seq counts, contrasted with RNA-seq counts, for a specific RBP KO line
1014 compared to the other two RBP KO lines. All counts were enumerated from mapped reads to the
1015 coding regions. We originally included a control KO line (AAVS1 locus) for comparison;
1016 however, by PCA, this KO line showed a distinct gene expression signature from that of the other
1017 KO lines, indicating it may not be suitable as a control (ExtendedDataFig. 16d-e). Using the
1018 AAVS1 KO line as a control, we observed highly similar hits for each RBP KO tested. We
1019 included filtering of counts using `edgeR::filterByExpr` with default parameters, the TMM method
1020 for calculation of size factors, and quasi-likelihood negative binomial models for fitting. Genes
1021 were considered differential at $FDR < 0.05$.

1022 REFERENCES CITED

- 1023 1. Tang, F. *et al.* mRNA-Seq whole-transcriptome analysis of a single cell. *Nat. Methods* **6**,
1024 377–382 (2009).
- 1025 2. Nagalakshmi, U. *et al.* The transcriptional landscape of the yeast genome defined by RNA
1026 sequencing. *Science* **320**, 1344–1349 (2008).
- 1027 3. Mortazavi, A., Williams, B. A., McCue, K., Schaeffer, L. & Wold, B. Mapping and
1028 quantifying mammalian transcriptomes by RNA-Seq. *Nat. Methods* **5**, 621–628 (2008).
- 1029 4. Schena, M., Shalon, D., Davis, R. W. & Brown, P. O. Quantitative monitoring of gene
1030 expression patterns with a complementary DNA microarray. *Science* **270**, 467–470 (1995).
- 1031 5. Chen, K. H., Boettiger, A. N., Moffitt, J. R., Wang, S. & Zhuang, X. RNA imaging. Spatially
1032 resolved, highly multiplexed RNA profiling in single cells. *Science* **348**, aaa6090 (2015).
- 1033 6. Combs, P. A. & Eisen, M. B. Sequencing mRNA from cryo-sliced *Drosophila* embryos to
1034 determine genome-wide spatial patterns of gene expression. *PLoS One* **8**, e71820 (2013).
- 1035 7. Achim, K. *et al.* High-throughput spatial mapping of single-cell RNA-seq data to tissue of
1036 origin. *Nat. Biotechnol.* **33**, 503–509 (2015).
- 1037 8. Langfelder, P. & Horvath, S. WGCNA: an R package for weighted correlation network
1038 analysis. *BMC Bioinformatics* **9**, 559 (2008).
- 1039 9. Eisen, M. B., Spellman, P. T., Brown, P. O. & Botstein, D. Cluster analysis and display of
1040 genome-wide expression patterns. *Proc. Natl. Acad. Sci. U. S. A.* **95**, 14863–14868 (1998).
- 1041 10. Skinnider, M. A., Squair, J. W. & Foster, L. J. Evaluating measures of association for single-
1042 cell transcriptomics. *Nat. Methods* **16**, 381–386 (2019).
- 1043 11. Stuart, J. M., Segal, E., Koller, D. & Kim, S. K. A gene-coexpression network for global
1044 discovery of conserved genetic modules. *Science* **302**, 249–255 (2003).
- 1045 12. Marcotte, E. M., Pellegrini, M., Thompson, M. J., Yeates, T. O. & Eisenberg, D. A combined
1046 algorithm for genome-wide prediction of protein function. *Nature* **402**, 83–86 (1999).
- 1047 13. Hughes, T. R. *et al.* Functional discovery via a compendium of expression profiles. *Cell* **102**,
1048 109–126 (2000).
- 1049 14. Kim, S. K. *et al.* A gene expression map for *Caenorhabditis elegans*. *Science* **293**, 2087–2092
1050 (2001).
- 1051 15. Hartl, C. L. *et al.* Coexpression network architecture reveals the brain-wide and multiregional
1052 basis of disease susceptibility. *Nat. Neurosci.* **24**, 1313–1323 (2021).
- 1053 16. DeRisi, J. L., Iyer, V. R. & Brown, P. O. Exploring the metabolic and genetic control of gene
1054 expression on a genomic scale. *Science* **278**, 680–686 (1997).
- 1055 17. Jansen, R., Greenbaum, D. & Gerstein, M. Relating whole-genome expression data with
1056 protein-protein interactions. *Genome Res.* **12**, 37–46 (2002).
- 1057 18. Szklarczyk, D. *et al.* The STRING database in 2023: protein-protein association networks
1058 and functional enrichment analyses for any sequenced genome of interest. *Nucleic Acids Res.*
1059 **51**, D638–D646 (2023).
- 1060 19. Tavazoie, S., Hughes, J. D., Campbell, M. J., Cho, R. J. & Church, G. M. Systematic
1061 determination of genetic network architecture. *Nat. Genet.* **22**, 281–285 (1999).
- 1062 20. Roth, F. P., Hughes, J. D., Estep, P. W. & Church, G. M. Finding DNA regulatory motifs
1063 within unaligned noncoding sequences clustered by whole-genome mRNA quantitation. *Nat.*
1064 *Biotechnol.* **16**, 939–945 (1998).

- 1065 21. Nusinow, D. P. *et al.* Quantitative Proteomics of the Cancer Cell Line Encyclopedia. *Cell*
1066 **180**, 387–402.e16 (2020).
- 1067 22. Gonçalves, E. *et al.* Pan-cancer proteomic map of 949 human cell lines. *Cancer Cell* **40**, 835–
1068 849.e8 (2022).
- 1069 23. Kustatscher, G. *et al.* Co-regulation map of the human proteome enables identification of
1070 protein functions. *Nat. Biotechnol.* **37**, 1361–1371 (2019).
- 1071 24. Ryan, C. J., Kennedy, S., Bajrami, I., Matallanas, D. & Lord, C. J. A Compendium of Co-
1072 regulated Protein Complexes in Breast Cancer Reveals Collateral Loss Events. *Cell Syst* **5**,
1073 399–409.e5 (2017).
- 1074 25. Furlong, E. E. M. & Levine, M. Developmental enhancers and chromosome topology.
1075 *Science* **361**, 1341–1345 (2018).
- 1076 26. Hnisz, D., Shrinivas, K., Young, R. A., Chakraborty, A. K. & Sharp, P. A. A Phase Separation
1077 Model for Transcriptional Control. *Cell* **169**, 13–23 (2017).
- 1078 27. Grabowski, P., Kustatscher, G. & Rappsilber, J. Epigenetic Variability Confounds
1079 Transcriptome but Not Proteome Profiling for Coexpression-based Gene Function Prediction
1080 *. *Mol. Cell. Proteomics* **17**, 2082–2090 (2018).
- 1081 28. Kustatscher, G., Grabowski, P. & Rappsilber, J. Pervasive coexpression of spatially proximal
1082 genes is buffered at the protein level. *Mol. Syst. Biol.* **13**, 937 (2017).
- 1083 29. Sonenberg, N., Hershey, J. W. B. & Mathews, M. B. *Translational Control of Gene*
1084 *Expression*. (CSHL Press, 2001).
- 1085 30. Kuersten, S. & Goodwin, E. B. The power of the 3' UTR: translational control and
1086 development. *Nat. Rev. Genet.* **4**, 626–637 (2003).
- 1087 31. Baker, S. A. & Rutter, J. Metabolites as signalling molecules. *Nat. Rev. Mol. Cell Biol.* **24**,
1088 355–374 (2023).
- 1089 32. Ozadam, H. *et al.* Single-cell quantification of ribosome occupancy in early mouse
1090 development. *Nature* **618**, 1057–1064 (2023).
- 1091 33. King, R. W., Deshaies, R. J., Peters, J. M. & Kirschner, M. W. How proteolysis drives the
1092 cell cycle. *Science* **274**, 1652–1659 (1996).
- 1093 34. Rao, S. *et al.* Genes with 5' terminal oligopyrimidine tracts preferentially escape global
1094 suppression of translation by the SARS-CoV-2 Nsp1 protein. *RNA* **27**, 1025–1045 (2021).
- 1095 35. Slobodin, B. *et al.* Cap-independent translation and a precisely located RNA sequence enable
1096 SARS-CoV-2 to control host translation and escape anti-viral response. *Nucleic Acids Res.*
1097 **50**, 8080–8092 (2022).
- 1098 36. Singh, G., Pratt, G., Yeo, G. W. & Moore, M. J. The Clothes Make the mRNA: Past and
1099 Present Trends in mRNP Fashion. *Annu. Rev. Biochem.* **84**, 325–354 (2015).
- 1100 37. Keene, J. D. & Tenenbaum, S. A. Eukaryotic mRNPs may represent posttranscriptional
1101 operons. *Mol. Cell* **9**, 1161–1167 (2002).
- 1102 38. Keene, J. D. RNA regulons: coordination of post-transcriptional events. *Nat. Rev. Genet.* **8**,
1103 533–543 (2007).
- 1104 39. Wurth, L. *et al.* UNR/CSDE1 Drives a Post-transcriptional Program to Promote Melanoma
1105 Invasion and Metastasis. *Cancer Cell* **36**, 337 (2019).
- 1106 40. Li, G.-W., Burkhardt, D., Gross, C. & Weissman, J. S. Quantifying absolute protein synthesis
1107 rates reveals principles underlying allocation of cellular resources. *Cell* **157**, 624–635 (2014).
- 1108 41. Taggart, J. C. & Li, G.-W. Production of Protein-Complex Components Is Stoichiometric and

- 1109 Lacks General Feedback Regulation in Eukaryotes. *Cell Syst* **7**, 580–589.e4 (2018).
- 1110 42. Ishikawa, K. Multilayered regulation of proteome stoichiometry. *Curr. Genet.* **67**, 883–890
1111 (2021).
- 1112 43. Amirbeigi Arab, S. *et al.* Invariable stoichiometry of ribosomal proteins in mouse brain tissues
1113 with aging. *Proc. Natl. Acad. Sci. U. S. A.* **116**, 22567–22572 (2019).
- 1114 44. Soto, I. *et al.* Balanced mitochondrial and cytosolic translomes underlie the biogenesis of
1115 human respiratory complexes. *Genome Biol.* **23**, 170 (2022).
- 1116 45. Natan, E. *et al.* Cotranslational protein assembly imposes evolutionary constraints on
1117 homomeric proteins. *Nat. Struct. Mol. Biol.* **25**, 279–288 (2018).
- 1118 46. Li, G.-W., Oh, E. & Weissman, J. S. The anti-Shine-Dalgarno sequence drives translational
1119 pausing and codon choice in bacteria. *Nature* **484**, 538–541 (2012).
- 1120 47. Seidel, M. *et al.* Co-translational assembly orchestrates competing biogenesis pathways. *Nat.*
1121 *Commun.* **13**, 1224 (2022).
- 1122 48. Bertolini, M. *et al.* Interactions between nascent proteins translated by adjacent ribosomes
1123 drive homomer assembly. *Science* **371**, 57–64 (2021).
- 1124 49. van den Boogaart, K. G., Filzmoser, P., Hron, K., Templ, M. & Tolosana-Delgado, R.
1125 Classical and Robust Regression Analysis with Compositional Data. *Math. Geosci.* **53**, 823–
1126 858 (2021).
- 1127 50. Quinn, T. P., Richardson, M. F., Lovell, D. & Crowley, T. M. propr: An R-package for
1128 Identifying Proportionally Abundant Features Using Compositional Data Analysis. *Sci. Rep.*
1129 **7**, 16252 (2017).
- 1130 51. Ozadam, H., Geng, M. & Cenik, C. RiboFlow, RiboR and RiboPy: an ecosystem for
1131 analyzing ribosome profiling data at read length resolution. *Bioinformatics* **36**, 2929–2931
1132 (2020).
- 1133 52. Gerashchenko, M. V. & Gladyshev, V. N. Ribonuclease selection for ribosome profiling.
1134 *Nucleic Acids Res.* **45**, e6 (2017).
- 1135 53. Mohammad, F., Green, R. & Buskirk, A. R. A systematically-revised ribosome profiling
1136 method for bacteria reveals pauses at single-codon resolution. *Elife* **8**, (2019).
- 1137 54. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide
1138 analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science*
1139 **324**, 218–223 (2009).
- 1140 55. Larsson, O., Sonenberg, N. & Nadon, R. Identification of differential translation in genome
1141 wide studies. *Proc. Natl. Acad. Sci. U. S. A.* **107**, 21487–21492 (2010).
- 1142 56. Quinn, T. P. *et al.* A field guide for the compositional analysis of any-omics data. *Gigascience*
1143 **8**, (2019).
- 1144 57. Sudmant, P. H., Alexis, M. S. & Burge, C. B. Meta-analysis of RNA-seq expression data
1145 across species, tissues and studies. *Genome Biol.* **16**, 287 (2015).
- 1146 58. Wang, Z.-Y. *et al.* Transcriptome and translome co-evolution in mammals. *Nature* **588**,
1147 642–647 (2020).
- 1148 59. Lu, P., Takai, K., Weaver, V. M. & Werb, Z. Extracellular matrix degradation and remodeling
1149 in development and disease. *Cold Spring Harb. Perspect. Biol.* **3**, (2011).
- 1150 60. Artieri, C. G. & Fraser, H. B. Evolution at two levels of gene expression in yeast. *Genome*
1151 *Res.* **24**, 411–421 (2014).
- 1152 61. McManus, C. J., May, G. E., Spealman, P. & Shteyman, A. Ribosome profiling reveals post-

- 1153 transcriptional buffering of divergent gene expression in yeast. *Genome Res.* **24**, 422–430
1154 (2014).
- 1155 62. Breschi, A., Gingeras, T. R. & Guigó, R. Comparative transcriptomics in human and mouse.
1156 *Nat. Rev. Genet.* **18**, 425–440 (2017).
- 1157 63. Crow, M., Suresh, H., Lee, J. & Gillis, J. Coexpression reveals conserved gene programs that
1158 co-vary with cell type across kingdoms. *Nucleic Acids Res.* **50**, 4302–4314 (2022).
- 1159 64. Pierson, E. *et al.* Sharing and Specificity of Co-expression Networks across 35 Human
1160 Tissues. *PLoS Comput. Biol.* **11**, e1004220 (2015).
- 1161 65. Kershaw, C. J. *et al.* Translation factor and RNA binding protein mRNA interactomes support
1162 broader RNA regulons for posttranscriptional control. *J. Biol. Chem.* **299**, 105195 (2023).
- 1163 66. Hentze, M. W., Castello, A., Schwarzl, T. & Preiss, T. A brave new world of RNA-binding
1164 proteins. *Nat. Rev. Mol. Cell Biol.* **19**, 327–341 (2018).
- 1165 67. Korbil, J. O., Jensen, L. J., von Mering, C. & Bork, P. Analysis of genomic context:
1166 prediction of functional associations from conserved bidirectionally transcribed gene pairs.
1167 *Nat. Biotechnol.* **22**, 911–917 (2004).
- 1168 68. Szklarczyk, R. *et al.* WeGET: predicting new genes for molecular systems by weighted co-
1169 expression. *Nucleic Acids Res.* **44**, D567–73 (2016).
- 1170 69. Zhang, M. *et al.* RNA-binding protein IMP3 is a novel regulator of MEK1/ERK signaling
1171 pathway in the progression of colorectal Cancer through the stabilization of MEKK1 mRNA.
1172 *J. Exp. Clin. Cancer Res.* **40**, 200 (2021).
- 1173 70. Cargnello, M. & Roux, P. P. Activation and function of the MAPKs and their substrates, the
1174 MAPK-activated protein kinases. *Microbiol. Mol. Biol. Rev.* **75**, 50–83 (2011).
- 1175 71. Bodén, M. & Bailey, T. L. Associating transcription factor-binding site motifs with target GO
1176 terms and target genes. *Nucleic Acids Res.* **36**, 4108–4117 (2008).
- 1177 72. Machanick, P. & Bailey, T. L. MEME-ChIP: motif analysis of large DNA datasets.
1178 *Bioinformatics* **27**, 1696–1697 (2011).
- 1179 73. Mecham, R. *The Extracellular Matrix: An Overview*. (Springer Science & Business Media,
1180 2011).
- 1181 74. Kagan, H. M. & Li, W. Lysyl oxidase: properties, specificity, and biological roles inside and
1182 outside of the cell. *J. Cell. Biochem.* **88**, 660–672 (2003).
- 1183 75. Kikuchi, A. *et al.* Structural basis for activation of DNMT1. *Nat. Commun.* **13**, 7130 (2022).
- 1184 76. Wu, Y.-Y. *et al.* The hTERT-p50 homodimer inhibits PLEKHA7 expression to promote
1185 gastric cancer invasion and metastasis. *Oncogene* **42**, 1144–1156 (2023).
- 1186 77. Kurita, S., Yamada, T., Rikitsu, E., Ikeda, W. & Takai, Y. Binding between the junctional
1187 proteins afadin and PLEKHA7 and implication in the formation of adherens junction in
1188 epithelial cells. *J. Biol. Chem.* **288**, 29356–29368 (2013).
- 1189 78. Pulimeno, P., Paschoud, S. & Citi, S. A role for ZO-1 and PLEKHA7 in recruiting
1190 paracingulin to tight and adherens junctions of epithelial cells. *J. Biol. Chem.* **286**, 16743–
1191 16750 (2011).
- 1192 79. Jeung, H.-C. *et al.* PLEKHA7 signaling is necessary for the growth of mutant KRAS driven
1193 colorectal cancer. *Exp. Cell Res.* **409**, 112930 (2021).
- 1194 80. Tavano, S. *et al.* Insm1 Induces Neural Progenitor Delamination in Developing Neocortex
1195 via Downregulation of the Adherens Junction Belt-Specific Protein Plekha7. *Neuron* **97**,
1196 1299–1314.e8 (2018).

- 1197 81. Sukonina, V. *et al.* FOXK1 and FOXK2 regulate aerobic glycolysis. *Nature* **566**, 279–283
1198 (2019).
- 1199 82. Kobe, B. & Kajava, A. V. The leucine-rich repeat as a protein recognition motif. *Curr. Opin.*
1200 *Struct. Biol.* **11**, 725–732 (2001).
- 1201 83. Evans, R. *et al.* Protein complex prediction with AlphaFold-Multimer. *bioRxiv*
1202 2021.10.04.463034 (2021) doi:10.1101/2021.10.04.463034.
- 1203 84. Carlsson, P. & Mahlapuu, M. Forkhead transcription factors: key players in development and
1204 metabolism. *Dev. Biol.* **250**, 1–23 (2002).
- 1205 85. The Human Transcription Factors. <http://humantfs.ccbr.utoronto.ca/cite.php>.
- 1206 86. Szklarczyk, D. *et al.* STRING v10: protein-protein interaction networks, integrated over the
1207 tree of life. *Nucleic Acids Res.* **43**, D447–52 (2015).
- 1208 87. Shiber, A. *et al.* Cotranslational assembly of protein complexes in eukaryotes revealed by
1209 ribosome profiling. *Nature* **561**, 268–272 (2018).
- 1210 88. Liesecke, F. *et al.* Ranking genome-wide correlation measurements improves microarray and
1211 RNA-seq based global and targeted co-expression networks. *Sci. Rep.* **8**, 10885 (2018).
- 1212 89. Ewing, R. M. *et al.* Large-scale mapping of human protein–protein interactions by mass
1213 spectrometry. *Mol. Syst. Biol.* **3**, 89 (2007).
- 1214 90. Drew, K., Wallingford, J. B. & Marcotte, E. M. hu.MAP 2.0: integration of over 15,000
1215 proteomic experiments builds a global compendium of human multiprotein assemblies. *Mol.*
1216 *Syst. Biol.* **17**, e10016 (2021).
- 1217 91. Heider, M. R. *et al.* Subunit connectivity, assembly determinants and architecture of the yeast
1218 exocyst complex. *Nat. Struct. Mol. Biol.* **23**, 59–66 (2016).
- 1219 92. Kee, Y. *et al.* Subunit structure of the mammalian exocyst complex. *Proc. Natl. Acad. Sci. U.*
1220 *S. A.* **94**, 14438–14443 (1997).
- 1221 93. Lalanne, J.-B. *et al.* Evolutionary Convergence of Pathway-Specific Enzyme Expression
1222 Stoichiometry. *Cell* **173**, 749–761.e38 (2018).
- 1223 94. Cenik, C. *et al.* Integrative analysis of RNA, translation, and protein levels reveals distinct
1224 regulatory variation across humans. *Genome Res.* **25**, 1610–1621 (2015).
- 1225 95. Bicknell, A. A. *et al.* Attenuating ribosome load improves protein output from mRNA by
1226 limiting translation-dependent mRNA decay. *Cell Rep.* **43**, 114098 (2024).
- 1227 96. Liu, T.-Y. *et al.* Time-Resolved Proteomics Extends Ribosome Profiling-Based
1228 Measurements of Protein Synthesis Dynamics. *Cell Syst* **4**, 636–644.e9 (2017).
- 1229 97. Piepoli, A. *et al.* The expression of leucine-rich repeat gene family members in colorectal
1230 cancer. *Exp. Biol. Med.* **237**, 1123–1128 (2012).
- 1231 98. Liu, Y. *et al.* Identification of differential expression of genes in hepatocellular carcinoma by
1232 suppression subtractive hybridization combined cDNA microarray. *Oncol. Rep.* **18**, 943–951
1233 (2007).
- 1234 99. Chen, H. *et al.* miR-218 contributes to drug resistance in multiple myeloma via targeting
1235 LRRC28. *J. Cell. Biochem.* **122**, 305–314 (2021).
- 1236 100. Wang, M., Herrmann, C. J., Simonovic, M., Szklarczyk, D. & von Mering, C. Version 4.0 of
1237 PaxDb: Protein abundance data, integrated across model organisms, tissues, and cell-lines.
1238 *Proteomics* **15**, 3163–3168 (2015).
- 1239 101. Kulevich, S. E., Frey, B. L., Kreitinger, G. & Smith, L. M. Alkylating tryptic peptides to
1240 enhance electrospray ionization mass spectrometry analysis. *Anal. Chem.* **82**, 10135–10142

- 1241 (2010).
- 1242 102. Rodriguez, J. M. *et al.* APPRIS: annotation of principal and alternative splice isoforms.
- 1243 *Nucleic Acids Res.* **41**, D110–7 (2013).
- 1244 103. *Sra-Tools: SRA Tools*. (Github).
- 1245 104. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
- 1246 *EMBnet.journal* **17**, 10–12 (2011).
- 1247 105. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**,
- 1248 357–359 (2012).
- 1249 106. Creators Liu, O. *HeLa Ribosome Profiling Data*. doi:10.5281/zenodo.10594392.
- 1250 107. Gerashchenko, M. V. & Gladyshev, V. N. Translation inhibitors cause abnormalities in
- 1251 ribosome profiling experiments. *Nucleic Acids Res.* **42**, e134 (2014).
- 1252 108. Wu, C. C.-C., Zinshteyn, B., Wehner, K. A. & Green, R. High-Resolution Ribosome Profiling
- 1253 Defines Discrete Ribosome Elongation States and Translational Regulation during Cellular
- 1254 Stress. *Mol. Cell* **73**, 959–970.e5 (2019).
- 1255 109. Wolin, S. L. & Walter, P. Ribosome pausing and stacking during translation of a eukaryotic
- 1256 mRNA. *EMBO J.* **7**, 3559–3569 (1988).
- 1257 110. Sharma, J. *et al.* A small molecule that induces translational readthrough of CFTR nonsense
- 1258 mutations by eRF1 depletion. *Nat. Commun.* **12**, 4358 (2021).
- 1259 111. Tukey, J. W. The Future of Data Analysis. *Ann. Math. Stat.* **33**, 1–67 (1962).
- 1260 112. Zhang, X.-O., Yin, Q.-F., Chen, L.-L. & Yang, L. Gene expression profiling of non-
- 1261 polyadenylated RNA-seq across species. *Genom Data* **2**, 237–241 (2014).
- 1262 113. Yang, L., Duff, M. O., Graveley, B. R., Carmichael, G. G. & Chen, L.-L. Genomewide
- 1263 characterization of non-polyadenylated RNAs. *Genome Biol.* **12**, R16 (2011).
- 1264 114. van den Boogaart, K. G. & Tolosana-Delgado, R. *Analyzing Compositional Data with R*.
- 1265 (Springer Berlin Heidelberg).
- 1266 115. orthogene. *Bioconductor*
- 1267 <https://bioconductor.org/packages/release/bioc/html/orthogene.html>.
- 1268 116. van den Boogaart, K. G. & Tolosana-Delgado, R. ‘compositions’: A unified R package to
- 1269 analyze compositional data. *Comput. Geosci.* **34**, 320–338 (2008).
- 1270 117. Kim, S. ppcor: An R Package for a Fast Calculation to Semi-partial Correlation Coefficients.
- 1271 *Commun Stat Appl Methods* **22**, 665–674 (2015).
- 1272 118. Berriz, G. F., Beaver, J. E., Cenik, C., Tasan, M. & Roth, F. P. Next generation software for
- 1273 functional trend analysis. *Bioinformatics* **25**, 3043–3044 (2009).
- 1274 119. Buttrey, S. & Whitaker, L. TreeClust: An R package for tree-based clustering dissimilarities.
- 1275 *R J.* **7**, 227 (2015).
- 1276 120. Wainberg, M. *et al.* A genome-wide atlas of co-essential modules assigns function to
- 1277 uncharacterized genes. *Nat. Genet.* **53**, 638–649 (2021).
- 1278 121. Gene Ontology Consortium *et al.* The Gene Ontology knowledgebase in 2023. *Genetics* **224**,
- 1279 (2023).
- 1280 122. Philippe, L., van den Elzen, A. M. G., Watson, M. J. & Thoreen, C. C. Global analysis of
- 1281 LARP1 translation targets reveals tunable and dynamic features of 5' TOP motifs.
- 1282 *Proceedings of the National Academy of Sciences* **117**, 5319–5328 (2020).
- 1283 123. Ballouz, S., Weber, M., Pavlidis, P. & Gillis, J. EGAD: ultra-fast functional analysis of gene
- 1284 networks. *Bioinformatics* **33**, 612–614 (2017).

- 1285 124. Carlson, M. org. Mm. eg. db: Genome wide annotation for Mouse. R package version 3.8. 2.
1286 2019.
- 1287 125. Carlson, M. org. Hs. eg. db: Genome wide annotation for Human. R package version 3.8. 2.
1288 2019.
- 1289 126. Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature* **596**,
1290 583–589 (2021).
- 1291 127. Lambert, S. A. *et al.* The Human Transcription Factors. *Cell* **172**, 650–665 (2018).
- 1292 128. UniProt Consortium. UniProt: the Universal Protein Knowledgebase in 2023. *Nucleic Acids*
1293 *Res.* **51**, D523–D531 (2023).
- 1294 129. Hou, Y., Xie, T., He, L., Tao, L. & Huang, J. Topological links in predicted protein complex
1295 structures reveal limitations of AlphaFold. *Commun Biol* **6**, 1098 (2023).
- 1296 130. Burke, D. F. *et al.* Towards a structurally resolved human protein interaction network. *Nat.*
1297 *Struct. Mol. Biol.* **30**, 216–225 (2023).
- 1298 131. Bryant, P., Pozzati, G. & Elofsson, A. Improved prediction of protein-protein interactions
1299 using AlphaFold2. *Nat. Commun.* **13**, 1265 (2022).
- 1300 132. Hu, Y. *et al.* Paralog Explorer: A resource for mining information about paralogs in common
1301 research organisms. *Comput. Struct. Biotechnol. J.* **20**, 6570–6577 (2022).
- 1302 133. Krismer, K. *et al.* Transite: A Computational Motif-Based Analysis Platform That Identifies
1303 RNA-Binding Proteins Modulating Changes in Gene Expression. *Cell Rep.* **32**, 108064
1304 (2020).
- 1305 134. Benoit Bouvrette, L. P., Bovaird, S., Blanchette, M. & Lécuyer, E. oRNAment: a database of
1306 putative RNA binding protein target sites in the transcriptomes of model species. *Nucleic*
1307 *Acids Res.* **48**, D166–D173 (2020).
- 1308 135. Van Nostrand, E. L. *et al.* A large-scale binding and functional map of human RNA-binding
1309 proteins. *Nature* **583**, 711–719 (2020).
- 1310 136. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic
1311 features. *Bioinformatics* **26**, 841–842 (2010).
- 1312 137. Stuart, E. A., King, G., Imai, K. & Ho, D. MatchIt: nonparametric preprocessing for
1313 parametric causal inference. *J. Stat. Softw.* (2011).
- 1314 138. Sanson, K. R. *et al.* Optimized libraries for CRISPR-Cas9 genetic screens with multiple
1315 modalities. *Nat. Commun.* **9**, 5416 (2018).
- 1316 139. Sanjana, N. E., Shalem, O. & Zhang, F. Improved vectors and genome-wide libraries for
1317 CRISPR screening. *Nat. Methods* **11**, 783–784 (2014).
- 1318 140. ChatGPT. <https://chat.openai.com>.
- 1319 141. Rossi, G. P. *et al.* Endothelin-1 stimulates steroid secretion of human adrenocortical cells ex
1320 vivo via both ETA and ETB receptor subtypes. *J. Clin. Endocrinol. Metab.* **82**, 3445–3449
1321 (1997).
- 1322 142. Sánchez-Caballero, L. *et al.* TMEM70 functions in the assembly of complexes I and V.
1323 *Biochim. Biophys. Acta Bioenerg.* **1861**, 148202 (2020).
- 1324 143. Carroll, J., He, J., Ding, S., Fearnley, I. M. & Walker, J. E. TMEM70 and TMEM242 help to
1325 assemble the rotor ring of human ATP synthase and interact with assembly factors for
1326 complex I. *Proc. Natl. Acad. Sci. U. S. A.* **118**, (2021).
- 1327 144. Mii, Y. & Takada, S. Heparan Sulfate Proteoglycan Clustering in Wnt Signaling and
1328 Dispersal. *Front Cell Dev Biol* **8**, 631 (2020).

- 1329 145. Kamimura, K. *et al.* Perlecan regulates bidirectional Wnt signaling at the Drosophila
1330 neuromuscular junction. *J. Cell Biol.* **200**, 219–233 (2013).
- 1331 146. Camlin, N. J., McLaughlin, E. A. & Holt, J. E. Kif4 Is Essential for Mouse Oocyte Meiosis.
1332 *PLoS One* **12**, e0170650 (2017).
- 1333 147. Tang, F. *et al.* Involvement of Kif4a in Spindle Formation and Chromosome Segregation in
1334 Mouse Oocytes. *Aging Dis.* **9**, 623–633 (2018).
- 1335 148. Robertson, A. K., Geiman, T. M., Sankpal, U. T., Hager, G. L. & Robertson, K. D. Effects of
1336 chromatin structure on the enzymatic and DNA binding functions of DNA methyltransferases
1337 DNMT1 and Dnmt3a in vitro. *Biochem. Biophys. Res. Commun.* **322**, 110–118 (2004).
- 1338 149. Schrader, A., Gross, T., Thalhammer, V. & Längst, G. Characterization of Dnmt1 Binding
1339 and DNA Methylation on Nucleosomes and Nucleosomal Arrays. *PLoS One* **10**, e0140076
1340 (2015).
- 1341 150. Ciossani, G. *et al.* The kinetochore proteins CENP-E and CENP-F directly and specifically
1342 interact with distinct BUB mitotic checkpoint Ser/Thr kinases. *J. Biol. Chem.* **293**, 10084–
1343 10101 (2018).
- 1344 151. Liao, H., Winkfein, R. J., Mack, G., Rattner, J. B. & Yen, T. J. CENP-F is a protein of the
1345 nuclear matrix that assembles onto kinetochores at late G2 and is rapidly degraded after
1346 mitosis. *J. Cell Biol.* **130**, 507–518 (1995).
- 1347 152. Chen, M. *et al.* FAT1 inhibits the proliferation and metastasis of cervical cancer cells by
1348 binding β -catenin. *Int. J. Clin. Exp. Pathol.* **12**, 3807–3818 (2019).
- 1349 153. Nishikawa, Y. *et al.* Human FAT1 cadherin controls cell migration and invasion of oral
1350 squamous cell carcinoma through the localization of β -catenin. *Oncol. Rep.* **26**, 587–592
1351 (2011).
- 1352 154. Morris, L. G. T. *et al.* Recurrent somatic mutation of FAT1 in multiple human cancers leads
1353 to aberrant Wnt activation. *Nat. Genet.* **45**, 253–261 (2013).
- 1354 155. Hou, R., Liu, L., Anees, S., Hiroyasu, S. & Sibinga, N. E. S. The Fat1 cadherin integrates
1355 vascular smooth muscle cell growth and migration signals. *J. Cell Biol.* **173**, 417–429 (2006).
- 1356 156. Vallet, S. D., Berthollier, C., Salza, R., Muller, L. & Ricard-Blum, S. The Interactome of
1357 Cancer-Related Lysyl Oxidase and Lysyl Oxidase-Like Proteins. *Cancers* **13**, (2020).
- 1358 157. Vallet, S. D. *et al.* Insights into the structure and dynamics of lysyl oxidase propeptide, a
1359 flexible protein with numerous partners. *Sci. Rep.* **8**, 11768 (2018).
- 1360 158. Yang, C. *et al.* Transcriptomic Analysis Identified ARHGAP Family as a Novel Biomarker
1361 Associated With Tumor-Promoting Immune Infiltration and Nanomechanical Characteristics
1362 in Bladder Cancer. *Front Cell Dev Biol* **9**, 657219 (2021).
- 1363 159. Lamarche-Vane, N. & Hall, A. CdGAP, a novel proline-rich GTPase-activating protein for
1364 Cdc42 and Rac. *J. Biol. Chem.* **273**, 29172–29177 (1998).
- 1365 160. Yang, S. *et al.* Control of antiviral innate immune response by protein geranylgeranylation.
1366 *Sci Adv* **5**, eaav7999 (2019).
- 1367 161. Bouhaddou, M. *et al.* The Global Phosphorylation Landscape of SARS-CoV-2 Infection. *Cell*
1368 **182**, 685–712.e19 (2020).
- 1369 162. Swaine, T. & Dittmar, M. T. CDC42 Use in Viral Cell Entry Processes by RNA Viruses.
1370 *Viruses* **7**, 6526–6536 (2015).
- 1371 163. Redmond, S. A. *et al.* Somatodendritic Expression of JAM2 Inhibits Oligodendrocyte
1372 Myelination. *Neuron* **91**, 824–836 (2016).

- 1373 164. Song, K. Y., Choi, H. S., Law, P.-Y., Wei, L.-N. & Loh, H. H. Vimentin interacts with the
1374 5'-untranslated region of mouse mu opioid receptor (MOR) and is required for post-
1375 transcriptional regulation. *RNA Biol.* **10**, 256–266 (2013).
- 1376 165. van der Brug, M. P. *et al.* RNA binding activity of the recessive parkinsonism protein DJ-1
1377 supports involvement in multiple cellular pathways. *Proc. Natl. Acad. Sci. U. S. A.* **105**,
1378 10244–10249 (2008).
- 1379 166. Niere, F. *et al.* Aberrant DJ-1 expression underlies L-type calcium channel hypoactivity in
1380 dendrites in tuberous sclerosis complex and Alzheimer's disease. *Proc. Natl. Acad. Sci. U. S.*
1381 *A.* **120**, e2301534120 (2023).
- 1382 167. Jin, W. *et al.* HydRA: Deep-learning models for predicting RNA-binding capacity from
1383 protein interaction association context and protein sequence. *Mol. Cell* **83**, 2595–2611.e11
1384 (2023).

1385 **Table 1: Literature support for gene functions predicted using TEC.**

1386 In the table, we list the predictions that are supported by literature. To do the new gene function
 1387 prediction, we selected GO terms with AUROC measured with TEC ≥ 0.8 , then focused on the
 1388 subset with differences between AUROC measured with TEC and RNA co-expression ≥ 0.1 .

Term	Species	Description	New adding gene (top ranking in TE)	TEC AUROC	RNA co-expression AUROC	Ranking of new adding gene in RNA	Reference
GO:0005496	Human	steroid binding	EDNRA	0.81	0.50	3991	141
GO:0022900	Human	electron transport chain	TMEM70	0.82	0.67	1611	142,143
GO:0042813	Human	Wnt-activated receptor activity	HSPG2	0.85	0.74	341	144,145
GO:0007129	Human	homologous chromosome pairing at meiosis	KIF4A	0.84	0.73	11	146,147
GO:0031492	Human	nucleosomal DNA binding	DNMT1	0.85	0.73	251	75,148,149
GO:0050793	Mouse	regulation of developmental	Plekha7	0.85	0.70	3929	76-80

Term	Species	Descripti on	New adding gene (top ranking in TE)	TEC AUROC	RNA co- expressio n AUROC	Ranking of new adding gene in RNA	Referen ce
		process					
GO:1990 023	Mouse	mitotic spindle midzone	Cenpf	0.90	0.74	54	150,151
GO:0016 342	Mouse	catenin complex	Fat1	0.85	0.73	781	152–155
GO:0005 539	Mouse	glycosami noglycan binding	Lox	0.95	0.83	32	156,157
GO:0140 374	Mouse	antiviral innate immune response	Arhgap3 1	0.84	0.73	1085	158–162
GO:0022 0101	Mouse	Central nervous system myelinati on	Jam2	0.88	0.76	584	163

1389