

Translation of neutrally evolving peptides provides a basis for *de novo* gene evolution

Jorge Ruiz-Orera^{1,*}, Pol Verdaguer-Grau², José Luis Villanueva-Cañas¹, Xavier Messeguer², M.Mar Albà^{1,3,*}

¹Evolutionary Genomics Group, Research Programme on Biomedical Informatics, Hospital del Mar Research Institute, Universitat Pompeu Fabra, Barcelona, Spain; ²Computer Sciences Department, Universitat Politècnica de Catalunya, Barcelona, Spain; ³Catalan Institution for Research and Advanced Studies, Barcelona, Spain.

*To whom correspondence should be addressed.

Running title: pervasive translation and *de novo* genes

Keywords: ribosome profiling, translation, *de novo* gene, long non-coding RNA, peptide, polymorphism, natural selection

This is a post-peer-review, pre-copyedit version of an article published in *Nature ecology & evolution*. The final authenticated version is available online at: <http://dx.doi.org/10.1038/s41559-018-0506-6>

1 **Abstract**

2

3 There is accumulating evidence that some protein-coding genes have originated *de novo* from
4 previously non-coding genomic sequences. However, the processes underlying *de novo* gene
5 birth are still enigmatic. In particular, the appearance of a new functional protein seems highly
6 improbable unless there is already a pool of neutrally evolving peptides that are translated at
7 significant levels and that can at some point acquire new functions. Here we use deep ribosome
8 profiling sequencing data, together with proteomics and single nucleotide polymorphism
9 information, to search for these peptides. We find hundreds of open reading frames that are
10 translated, and that show no evolutionary conservation and no selective constraints. The data
11 suggests that the translation of these neutrally evolving peptides may be facilitated by the
12 chance occurrence of open reading frames with a favorable codon composition. We conclude
13 that the pervasive translation of the transcriptome provides plenty of material for the evolution of
14 new functional proteins.

15 Introduction

16

17 A large fraction of the genome is transcribed; this includes functional genes but also thousands
18 of transcripts that are not conserved across species and which show weak or no signatures of
19 natural selection¹⁻³. Many of the latter transcripts are annotated as long non-coding RNAs
20 (lncRNAs) because they lack conserved long open reading frames (ORFs). Recent studies
21 based on the sequencing of ribosome-protected RNA fragments (ribosome profiling) have
22 reported that a surprisingly large number of these transcripts are associated with ribosomes and
23 may translate small proteins or peptides⁴⁻⁹. This is intriguing because only a small fraction of
24 them are likely to encode functional micropeptides¹⁰.

25

26 Each ribosome profiling experiment generates millions of ribosome footprints that are
27 subsequently mapped to the genome or the transcriptome to identify open reading frames
28 (ORFs) that are being translated¹¹. The codon-by-codon movement of the ribosome along the
29 coding sequence results in a characteristic pattern of three nucleotide periodicity of the mapped
30 reads, which makes ribosome profiling a very useful method to detect novel events of
31 translation^{4,12,13}. Given enough sequence coverage the technique can uncover low-abundant
32 small proteins that would be otherwise difficult to detect by standard proteomics
33 approaches^{14,15}.

34

35 There is increasing evidence that some functional genes have originated *de novo* from
36 previously non-functional parts of the genome¹⁶⁻²⁰. There are many protein-coding genes that
37 appear to be species-specific; the encoded proteins are shorter than average and they tend to
38 evolve under more relaxed constraints than other genes^{21,22,20}. There is experimental evidence
39 that at least some *de novo* genes are functional^{23,24}. However, for a new protein to acquire a
40 function it first needs to be produced in the cell at significant amounts. Some of the new proteins
41 may be generated by the translation of randomly occurring ORFs in transcripts currently
42 annotated as non-coding⁷. As these proteins are not expected to be functional at first, they are
43 expected to evolve under no significant purifying selection, or neutrally. We previously showed
44 that, as a whole, putatively translated lncRNAs and young annotated protein-coding genes

45 shared a number of similarities, such as small ORF size and weak selective constraints,
46 compared with more widely conserved genes⁸. This suggested that the number of proteins with
47 a recent evolutionary origin was larger than previously thought. However, it remained to be
48 determined if there existed a pool of young proteins that were translated at significant levels and
49 evolved in a neutral manner, as hypothesized for *de novo* gene evolution.

50

51 To assess the functional relevance of translated products one can use the ratio between the
52 number of non-synonymous and synonymous substitutions in the putative coding sequences^{4,5}.
53 However, this method requires an alignment of at least two different homologous sequences. A
54 more general approach that can be used in the absence of homologues in other species is the
55 ratio between the number of non-synonymous and synonymous single nucleotide
56 polymorphisms, compared to the one expected under neutrality. Under no selection, non-
57 synonymous and synonymous polymorphisms persist at the same levels, whereas under
58 purifying selection there is a deficit of non-synonymous polymorphisms because some amino
59 acid changes disrupt the protein's function²⁵. Single nucleotide polymorphism analysis can be
60 performed on a gene-by-gene basis or in pools of sequences that share certain features^{2,26}.

61

62 Here we employ phylogenetic comparisons, ribosome profiling data and single nucleotide
63 polymorphism information to search for putative precursors of *de novo* genes in mouse. The
64 study renders visible a layer of protein expression that produces non-functional small proteins,
65 filling a gap in our understanding of the processes underlying *de novo* gene birth.

66

67 **Results**

68

69 Many ORFs outside annotated protein-coding sequences are translated

70

71 First we set out to identify translated open reading frames (ORFs) in mouse protein-coding
72 genes (codRNAs) and long non-coding RNAs (lncRNAs) using ribosome-profiling RNA-
73 sequencing (Ribo-Seq) data from eight different tissues and cell lines (Supplementary Table 1
74 and references therein). In contrast to RNA sequencing (RNA-Seq) reads, which are expected

75 to cover the complete transcript, Ribo-Seq reads correspond to regions bound by ribosomes.
76 We mapped the RNA-Seq and Ribo-Seq reads to the mouse Ensembl gene annotations and,
77 for the sake of completeness, also to a set of previously obtained novel mouse transcripts that
78 did not correspond to annotated protein-coding genes or lncRNAs³.
79
80 We used the RibORF program⁴ to identify translated sequences among ORFs covered by at
81 least 10 Ribo-Seq reads in transcripts expressed in one or more tissues (Fig. 1a and
82 Supplementary Table 1). This program calculates a score for each ORF depending on the 3-
83 nucleotide periodicity and uniformity of the mapped reads. Using a highly stringent RibORF
84 score cut-off of 0.7⁴ we found that about 90% of the coding genes (15,020), and 20% of the
85 annotated lncRNAs (539), were predicted to be translated in at least one sample. Additionally,
86 we identified 286 novel, non-annotated, expressed loci that also contained translated ORFs
87 (Fig. 1b). Annotated lncRNAs and novel transcripts mostly encoded small proteins (Fig. 1c) and
88 the two classes were merged.

89
90 Hundreds of translated ORFs have no homologues in other species and evolve neutrally

91
92 Putative precursors of *de novo* genes should be of recent origin and evolve under no purifying
93 selection. We performed exhaustive sequence similarity searches of the translated ORFs
94 against high coverage transcriptomes from human and rat as well as against the annotated
95 proteomes of 101 different eukaryotic species (Fig. 2a, Supplementary Table 2 for a list of
96 species, see Methods for more details). For these searches we discarded any peptides shorter
97 than 24 amino acids, as the detection of homologues may be compromised in such cases due
98 to lack of sufficient sequence information. We identified 1,980 translated ORFs that showed no
99 homology to expressed sequences in other species (class non-conserved or NC). This
100 represented about 10.4% of the translated ORFs.

101
102 Next, we measured the strength of selection acting on the putative proteins using a large
103 collection of mouse single nucleotide polymorphisms (SNPs) for the house mouse subspecies
104 *Mus musculus castaneus*²⁷. We could map a total of 324,729 SNPs to the set of translated

105 ORFs. We calculated the ratio between the number of observed non-synonymous and
106 synonymous SNPs (PN/PS(obs)) in conserved and non-conserved groups of ORFs, and
107 normalized it by the same ratio expected under neutrality (PN/PS(exp)). The expected PN/PS
108 was estimated using a table of nucleotide mutation frequencies in *Mus musculus castaneus* and
109 the observed codon frequencies in each set of sequences of interest (Supplementary Tables 3
110 to 5). This allowed us not only to compare the strength of selection across different sets of
111 sequences, as done in a previous study of ORFs translated from lncRNAs⁸, but also to discard
112 purifying selection if the normalized PN/PS was not significantly different from 1. Specifically, we
113 used a chi-square test that compared the number of observed and expected non-synonymous
114 and synonymous SNPs in each sequence set. The normalized PN/PS of conserved ORFs was
115 around 0.15 (Fig. 2b, p-value < 10⁻⁵), consistent with protein functionality. One example in this
116 group was Stannin^{28,29}, a highly conserved peptide that regulates neuronal cell apoptosis (Fig.
117 3).

118

119 In the case of the non-conserved ORFs, the normalized PN/PS depended on the codon usage
120 bias or coding score (Supplementary Fig. 1 and 2). The coding score provides a measurement
121 of how similar the relative codon frequencies are to those observed in functional coding
122 sequences^{20,30,21}. We divided the non-conserved ORFs in two groups: the first group, with high
123 coding scores (≥ 0.1014 , NC-H), showed a weak but significant signature of selection (Fig. 2b,
124 p-value < 0.05); the second group, comprising the rest of ORFs, showed no significant deviation
125 from neutrality (NC-L)(see Fig. 3 for a specific example). The first group is likely to contain some
126 species-specific functional genes, whereas the second group, the neutrally evolving ORFs,
127 represents new translated peptides with no function.

128

129 Neutrally evolving ORFs were found both in lncRNAs but also in many protein-coding genes
130 (Fig. 2b). In lncRNAs they were usually found in isolation or with other non-conserved ORFs,
131 whereas in coding genes they were in most cases located in alternative transcript isoforms of
132 genes already containing a conserved ORF. One possible way the secondary ORFs in protein-
133 coding genes could become new genes would be through gene duplication and inactivation of
134 the main ORF. The ORFs in lncRNAs and coding genes were merged into a single group,

135 comprising 1,291 ORFs, for subsequent analyses (neutral ORFs).

136

137 Validation of the set of neutrally evolving peptides

138

139 The above analyses grouped the sequences into classes before computing the PN/PS ratio. In
140 general, ORF-by-ORF analysis was not possible because the ORFs were small and contained
141 too few SNPs. Nevertheless, 41 of the ORFs in the neutral set contained 10 or more SNPs, and
142 we decided to compute a normalized PN/PS ratio for these individual cases. The median PN/PS
143 of these ORFs was around 1 and the distribution of PN/PS values was very different from that
144 observed for conserved ORFs (Supplementary Fig. 3, Wilcoxon test, p -value $< 10^{-5}$), consistent
145 with the previous results. Finally, we quantified the number of ORFs that contained SNPs that
146 generated premature stop codons, truncating more than half of the ORF, in the set of neutrally
147 evolving ORFs and in the set of conserved ORFs. In the first case we found 72 out of 1,291
148 ORFs that contained this type of mutation (5.6%) and in the second case 296 out of 17,043
149 ORFs (1.74%). Considering that neutral ORFs are in general much shorter than conserved
150 ORFs (median protein size 44 *versus* 412 amino acids), and thus less likely to accumulate
151 ORF-truncating mutations by chance alone, the data clearly indicates a strong excess of ORF-
152 truncating SNPs in neutral ORFs with respect to conserved ORFs.

153

154 We next calculated which was the percentage of neutral ORFs in the transcripts that were
155 detected as translated by our pipeline. This value was 30.5% (1,291 out of 4,232), about one
156 order of magnitude larger than the false positive rate estimated from two negative controls. The
157 first control was a set of ORFs derived from small nuclear and nucleolar RNAs (3.33% false
158 positive rate, see example in Fig. 3), the second one was a set of ORFs in which the position of
159 the reads had been randomized (4.16% false positive rate; see Methods for more details).

160

161 We used proteomics data from the PRIDE database³¹ to search for additional evidence of the
162 neutrally evolving peptides. Despite their small size (median 44 amino acids), a limiting factor
163 for their detection by standard proteomics-based techniques³², we found proteomics evidence
164 for 32 of the translated ORFs in this group. This represents 2.5% of the proteins in this set

165 (compared to less than 0.2% false positive rate, see Methods). This fraction is similar to the one
166 obtained for conserved proteins subsampled to have a similar size distribution as the neutral
167 ORFs (2.9%; in contrast, about 41% of all conserved ORFs have proteomics evidence). The
168 subset of ORFs with proteomics evidence did not deviate significantly from neutrality according
169 to the PN/PS-based test (Supplementary Table 5).

170

171 What determines that some ORFs, but not others, are translated?

172

173 Some transcripts contained relatively long ORFs but were not translated. One example of this
174 sort was the previously described *de novo* non-coding gene *Pold1*³³ that lacked any evidence of
175 translation in the data we analyzed. We next asked which factors may influence the translation
176 of neutral ORFs. First, we inspected the translation initiation sequence context but did not
177 detect any clear differences between translated and non-translated neutral ORFs
178 (Supplementary Fig. 4). We then hypothesized that the ORF coding score could affect the
179 “translatability” of the transcript because codons that are abundant in coding sequences are
180 expected to be more efficiently translated than other codons³⁴. Consistent with this hypothesis,
181 we found that the translated neutral ORFs exhibited higher coding scores than non-translated
182 ORFs (Fig. 4a, Translated versus non-translated Wilcoxon test, p-value < 10⁻⁵). Importantly, the
183 differences were maintained after subsampling the genes so that the two sets had the same
184 gene expression distribution (Fig. 4b, Wilcoxon test, p-value < 10⁻⁵). This is consistent with
185 codon composition having an effect *per se* in ORF translation. When controlling by coding
186 score, expression level, but not ORF length, was positively related to the detection of translation
187 (Fig. 4c).

188

189 It has been previously hypothesized that what distinguishes translated from non-translated
190 lncRNAs is the relative amount of the lncRNA in the cytoplasm with respect to the nucleus⁴.
191 However, we found evidence that some lncRNAs with nuclear functions, such as *Malat1* and
192 *Neat1*, translated neutral ORFs, suggesting that the cytosolic fraction of a transcript can be
193 translated independently of its role or preferred location. Instead, we found that neutral ORFs
194 enriched in codons that are abundant in functional protein-coding genes could be translated

195 more efficiently than other ORFs. This is consistent with the observation that abundant codons
196 enhance translation elongation³⁵, whereas rare codons might affect the stability of the mRNA³⁶.
197 These differences could result in compositional biases in newly emerged proteins.

198

199 **Discussion**

200

201 The molecular mechanisms underlying *de novo* gene evolution are still poorly understood^{17,18,37}.
202 The sudden appearance of a new protein-coding gene from a genomic segment seems *a priori*
203 highly improbable, but the process becomes much more likely if the genome is already being
204 pervasively transcribed and translated outside functional protein-coding sequences. An excess
205 of transcription was already noted in the first large-scale cDNA sequencing efforts performed in
206 human and mouse³⁸, and more recent studies have found a high rate of genome transcriptional
207 turnover when comparing closely related species³⁹. Other works have found that many open
208 reading frames outside annotated protein-coding sequences can be translated^{4,5,7}, raising the
209 possibility that they are an important source of *de novo* genes^{9,20,8}.

210

211 The precursors of new functional proteins are expected to be translated at relatively high levels
212 and evolve under no purifying selection. Here we have detected hundreds of small proteins that
213 appear to be of recent evolutionary origin and that fulfill these properties. Their real number
214 could be much higher considering that many cell types and tissues have not yet been sampled.
215 According to recent estimations, the cost of transcription and translation in multicellular
216 organisms is probably too small to overcome genetic drift⁴⁰. Therefore, if the peptides do not
217 affect the fitness *per se*, these activities may be effectively neutral. Our results are in line with
218 this prediction.

219

220 It has been proposed that nascent proteins are pre-adapted in the sense that only peptides that
221 do not have a harmful effect will be expressed at significant levels⁹, and that this may explain
222 the observed compositional differences between young proteins and intergenic sequences⁴¹.
223 Another possible source of bias is the codon composition of the sequence. According to our
224 results, ORFs that have a high coding score may be more efficiently translated than those that

225 have a lower coding score. As a result, new proteins may have a more more coding-like
226 composition than the average random or intergenic sequence.

227

228 The process of *de novo* gene origination involves the gain of a useful function by a previously
229 non-functional sequence. The rate at which this happens remains to be determined but it has
230 been observed that many artificially generated peptides can function as secretion signals⁴², and
231 selection for ATP-binding activity in a library of randomly generated 80 amino acid polypeptides
232 successfully identified several candidates capable of binding to ATP⁴³. Recent experiments
233 performed in *E.coli* also suggest that random expressed sequences can often affect cellular
234 growth⁴⁴.

235

236 We have shown that the transcriptome is pervasively translated, and that this generates novel
237 peptides that evolve under no constraints. This set of translated peptides, and their sequences,
238 are expected to change rapidly during evolution, generating a myriad of opportunities for the
239 birth of new functional genes.

240

241 **Methods**

242

243 Transcript assembly

244

245 We used strand-specific polyA+ RNA sequencing data (RNA-Seq) data from different mouse
246 and human tissues to assembly the species transcriptomes (Gene Expression Omnibus mouse
247 GSE69241³, GSE43721⁴⁵, and GSE43520⁴⁶; human GSE69241³). The mouse RNA samples
248 were extracted from strain Balb/C. RNA-Seq reads were filtered by length (> 25 nucleotides)
249 and by quality using Condetri (v.2.2)⁴⁷ with the following settings: -hq = 30 -lq = 10. We aligned
250 the reads to the corresponding reference species genome with Tophat (v. 2.0.8, -N 3, -a 5 and
251 -m 1)⁴⁸. Multiple mapping to several locations in the genome was allowed unless otherwise
252 stated. We assembled the transcriptome with Stringtie⁴⁹, merging the reads from all the
253 samples, with parameters -f 0.01, and -M 0.2. We used the species transcriptome as a guide
254 (Ensembl v.75), including all annotated isoforms, but permitting the assembly of annotated and

255 novel isoforms and genes (antisense, intergenic and intronic) as well. We excluded lncRNAs
256 that overlapped annotated pseudogenes or that showed significant sequence similarity to
257 known protein-coding sequences (BLASTP, e-value < 10⁻⁴). In the case of rat we employed a
258 previously generated transcript assembly⁵⁰.

259

260 Ribosome profiling data

261

262 We used ribosome profiling data (Ribo-Seq) from 8 different mouse tissues or cell lines (see
263 Supplementary Table 1), obtained from Gene Expression Omnibus under accession numbers
264 GSE51424⁵¹, GSE50983⁵², GSE22001⁵³, GSE62134⁵⁴, GSE72064⁵⁵, and GSE41246. Only
265 datasets corresponding to non-pathogenic conditions were considered. The reads from
266 experimental replicates were merged before using RibORF to increase the resolution of the
267 read periodicity, as done in the original RibORF paper⁴. For all analyses we considered only
268 genes expressed at significant levels in at least one sample (RNA-Seq fragments per kilobase
269 per Million mapped reads (FPKM) > 0.2). The expression of the genes detected in these
270 samples is expected to be highly representative of the *Mus musculus* species as a whole. We
271 mapped several brain RNA-Seq datasets from *Mus musculus castaneus*³⁹ to the mouse
272 assembled transcriptome using NextGenMap⁵⁶. As expected, the vast majority of the genes
273 expressed in brain samples from C57BL/6 mice⁵¹ also showed evidence of expression in *Mus*
274 *musculus castaneus* brain RNA samples³⁹ (Supplementary Table 6).

275

276 We discarded anomalous reads (length < 26 or > 33 nt) and reads that mapped to annotated
277 rRNAs and tRNAs in mouse from the Ribo-Seq sequencing datasets. Next, reads were mapped
278 to the assembled mouse genome (mm10) with Bowtie (v. 0.12.7, parameters -k 1 -m 20 -n 1 --
279 best --strata). Considering that the ORFs had to be extensively covered by reads to be
280 considered translated (high uniformity), we decided to include multiple mapped reads. This is
281 expected to increase the sensitivity to detect translation events, especially in gene families, but
282 also the number of false positives. We observed that using multiple mapped reads increased
283 the periodicity and uniformity values of a number of protein-coding genes, but that it had little
284 effect on translated ORFs from lncRNAs (Supplementary Fig. 5). We used the mapping of the

285 Ribo-Seq reads to the complete set of annotated coding sequences in mouse to compute the
286 position of the P-site (second binding site for tRNA in the ribosome) for reads of different size,
287 as previously described^{11,13}.

288

289 Identification of translated ORFs

290

291 We predicted all translated ORFs (ATG to STOP) with a minimum length of 9 amino acids in the
292 transcripts with RibORF (v.0.1)⁴. Only ORFs with a minimum of 10 mapped Ribo-Seq reads
293 were considered. The RibORF classifier is based on a support vector machine algorithm,
294 originally applied to human transcripts. The input parameters are the 3-nucleotide periodicity
295 and the uniformity of the reads. The periodicity refers to the distribution of the reads in the three
296 different possible frames; two values are computed, f1, which is the fraction of reads that
297 correspond to the correct frame, and f2, which is the fraction of reads in position +1 of the
298 correct frame. The uniformity corresponds to the percentage of maximum entropy, a value of 1
299 indicates a completely even distribution of reads across regions. For each ORF the program
300 computes a score that depends on the periodicity and the uniformity⁴. We used the same score
301 cut-off as in the original paper (≥ 0.7), which had a reported false positive rate of 0.67% and a
302 false negative rate of 2.5%.

303

304 We eliminated any redundancy in the translated ORFs by taking the longest ORF when several
305 overlapping translated ORFs were detected in the same gene. The identification of translated
306 ORFs was done separately for the different tissues (Supplementary Table 1). Differences in the
307 number of translated ORFs in different tissues were related to the depth of sequencing and the
308 number of reads that mapped to the most highly expressed proteins (Supplementary Fig. 6 and
309 7, respectively). When translation of the same ORF was detected in several tissues, we kept the
310 information for the tissue with the highest RibORF score.

311

312 Sequence conservation

313

314 We searched for mouse ORF homologues in the human and rat transcriptomes using TBLASTN

315 (limited to one strand, e-value < 10^{-4})⁵⁷. We also performed sequence similarity searches
316 against the annotated proteomes of 67 mammalian-species and 34 non-mammalian eukaryotes
317 from a diverse range of groups compiled in a previous study⁵⁰, using BLASTP (e-value < 10^{-4}).
318 For these searches we only considered query proteins of size 24 amino acids or longer, as
319 shorter proteins may not contain sufficient information to perform homology searches. Mouse
320 ORFs that did not have any homology hits in other species were classified as non-conserved,
321 the rest as conserved. Non-conserved ORFs located upstream or downstream of another longer
322 ORF in a conserved transcript (uORFs and dORFs) were excluded from this analysis.

323

324 We inspected the rat genomic syntenic regions of translated ORFs using LiftOver⁵⁸. We
325 classified the ORFs in two groups depending on whether the ORF was truncated in rat or not
326 (the truncation had to affect more than half of the protein). For neutrally evolving ORFs the
327 number of cases in which the ORF was truncated was similar to the number of cases in which it
328 was not truncated, and in both cases the polymorphism patterns were consistent with neutrality
329 (Supplementary Table 5). This indicated that, for this group, the presence of a similar ORF in rat
330 does not imply functional conservation of the ORF. Therefore, we did not use information on rat
331 genomic synteny to classify the genes as conserved/non-conserved.

332

333 Coding score

334

335 We used a previously described metric based on hexamer frequencies to calculate the coding
336 score of the sequences⁸. The method uses a table of pre-calculated hexamer scores that
337 measure the relative frequency of each hexamer in coding versus non-coding sequences.
338 These scores are then used to evaluate the coding propensity of a sequence based on its
339 hexamer composition. The method has been implemented in a computational program called
340 CIPHER that can be used to predict putative translated ORFs and can be accessed online
341 (<http://evolutionarygenomics.upf.edu/cipher>).

342

343 Comparison of the coding score in translated and non-translated ORFs

344

345 For comparative purposes we generated a set of non-translated ORFs that was completely
346 independent from the translated ORFs (Fig. 4). We selected expressed loci that did not contain
347 any translated ORF and selected the longest predicted ORF sequence. An ORF was
348 considered not to be translated if it had less than 10 Ribo-Seq mapped reads, or a RibORF
349 score < 0.7. Selecting the longest ORF was justified by the fact that, in translated ORFs, the
350 ORF with the highest number of mapped Ribo-Seq reads was usually the longest ORF (75.7%
351 for codRNAs and 84% for lncRNAs). We also generated a set of 4,013 randomly taken ORFs
352 from introns, after discarding ORFs that showed significant sequence similarity to known
353 proteins from the same species (BLASTP, e-value < 10^{-4}).

354

355 In order to investigate the influence of gene expression level in our capacity to detect translation
356 of an ORF we subsampled the set of translated neutral ORFs and the set of non-translated
357 ORFs so that the two sets had equivalent gene expression value distributions. As gene
358 expression units we used Fragments Per Kilobase per Million mapped reads (FPKM), and we
359 selected the highest FPKM value across the different tissues examined. We applied the same
360 procedure to control for coding score values when evaluating the effect of gene expression level
361 and ORF length.

362

363 Ribosome profiling controls

364

365 We inspected the ribosome profiling patterns of neutral ORFs with respect to the rest of
366 translated ORFs ("functional"). The number of Ribo-Seq reads per base was in general lower in
367 neutral ORFs than in functional ORFs, although there was a significant overlap in the two
368 distributions (Supplementary Fig. 8). Translated sequences were characterized by strong 3-
369 nucleotide periodicity and this property showed a significant correlation across pairs of tissues,
370 both for neutral and functional ORFs (Supplementary Fig. 8 and 9). After subsampling the
371 functional ORFs to match the neutral ORFs with regards to ORF size and number of Ribo-Seq
372 reads, the two sets showed similar RibORF score, periodicity and uniformity value distributions
373 (Supplementary Fig. 10).

374

375 We generated a negative control set by taking randomly occurring ORFs in mouse small
376 nuclear and nucleolar RNAs (“sRNAs”). These ORFs were required to have at least 10 Ribo-
377 Seq mapped reads and were processed in the same manner as the main set of ORFs under
378 study. Only 10 out of 304 ORFs in “sRNA” had a RibORF score equal or higher than 0.7 (3.33%
379 false positive rate). This control was characterized by low uniformity and low read periodicity
380 (Supplementary Fig. 8). A second negative control in which all parameters were equal to those
381 observed in neutral ORFs, but the reads were randomly distributed in the three frames, resulted
382 in 176 out of 4,232 ORFs classified as translated (“Random”, 4.16% false positive rate).

383

384 We used a similar strategy to estimate the fraction of translated non-conserved ORFs (NC) over
385 all non-conserved ORF and the fraction of translated neutral ORFs over all neutral ORFs (NC-
386 L). Non-translated ORFs were those covered by 10 or more Ribo-Seq reads but with a RibORF
387 score lower than 0.7. We followed the same procedure described for the translated ORFs
388 (Figure 2a) to classify them into classes. We detected translation for 33.1% of the non-
389 conserved ORFs (1,980 out of 5,975) and 30.5% of the neutral ORFs (1,291 out of 4,232).
390 These values are much higher than the estimated false positive rate.

391

392 We also generated a positive control set composed of 2,163 randomly taken annotated mouse
393 coding sequences with protein evidence in Uniprot. With this control we estimated a false
394 negative rate of 2.54%.

395

396 Single nucleotide polymorphism analysis

397

398 We obtained single nucleotide polymorphism (SNP) data from 20 individuals of the house
399 mouse subspecies *Mus musculus castaneus*²⁷. We classified SNPs in ORFs as non-
400 synonymous (PN, amino acid altering) and synonymous (PS, not amino-acid altering). We
401 calculated the PN/PS ratio in each ORF group by using the sum of PN and PS in all the
402 sequences ((PN/PS)_{obs}). We calculated the expected PN/PS under neutrality ((PN/PS)_{exp})
403 using the mutation frequencies between pairs of nucleotides in *Mus musculus castaneus* and
404 the codon composition of the different sequences or sets of sequences under study

405 (Supplementary Tables 3 and 4). The observed transition to transversion ratio was 4.42, very
406 similar to the 4.26 value obtained in early observations based on mouse-rat divergence data⁵⁹.
407 We tested for purifying selection using a chi-square test with one degree of freedom
408 (Supplementary Table 5). Positively selected mutations are rapidly fixed in the population and
409 their effect is expected to be negligible when using SNP data.

410

411 Proteomics data

412

413 We used the proteomics database PRIDE³¹ to search for peptide matches in the proteins
414 encoded by various gene sets. For a protein to have proteomics evidence, we required at least
415 two distinct perfect matches of peptides that did not map to any other protein in the dataset,
416 allowing for up to two mismatches. Under these conditions the false positive rate was estimated
417 to be below 0.2% in a previous study⁵⁰. The estimation was based on the results obtained for
418 ORFs in which the amino acid sequence was randomized and for ORFs in intronic sequences.

419

420 Statistical tests and plots

421

422 The generation of plots and statistical tests was performed with the R package⁶⁰.

423

424 Code availability

425

426 Code for calculating the coding score is available at <https://github.com/jorruior/CIPHER> and for
427 calculating the expected PN/PS at https://figshare.com/articles/computePNPS_c/5085706.

428

429 **Data availability**

430

431 The datasets generated during the current study are available from figshare repository,
432 <http://dx.doi.org/10.6084/m9.figshare.4702375>.

433

434 **Figure legends**

435

436 **Figure 1. Detection of translated ORFs.** **a.** Workflow to identify translated ORFs. Ribosome
437 profiling (Ribo-Seq) reads, corresponding to ribosome-protected fragments, are mapped to all
438 predicted canonical ORFs with length ≥ 30 nucleotides in transcripts. This is performed with
439 single-nucleotide resolution after computing the read P-site per each read length. In each ORF,
440 reads per frame and read uniformity are evaluated by RiboORF. **b-c.** Number of translated and
441 non-translated expressed loci (**b**) and ORFs (**c**) using a RibORF cut-off of 0.7. The original data
442 come from eight different mouse tissues (Supplementary Table 1). The translated ORFs have
443 been divided into small ORFs (smORF, < 100 aa) and long ORFs (≥ 100 aa), depending on
444 their length. “Novel” corresponds to expressed loci that do not correspond to annotated protein-
445 coding genes or lncRNAs, they are merged with lncRNAs as they mostly contain short non-
446 conserved ORFs.

447

448 **Figure 2. Identification of selection signatures.** **a.** Workflow to identify conserved and non-
449 conserved ORFs. Translated ORFs shorter than 24 amino acids, as well as non-conserved
450 upstream and downstream ORF in conserved transcripts (uORFs and dORFs, see Methods),
451 were filtered out. Any ORF with at least one BLAST match in another species was classified as
452 conserved (C), otherwise it was classified as non-conserved (NC). Non-conserved ORFs with a
453 high coding score value (≥ 0.1014) were classified as NC-H, and the rest were classified as NC-
454 L. **b.** Analysis of selective constraints in translated ORFs. PN/PS (obs/exp) refers to the
455 normalized ratio between non-synonymous (PN) and synonymous (PS) single nucleotide
456 polymorphisms; a value of 1 is expected in the absence of selection at the protein level.
457 Conserved and NC-H ORFs showed significant purifying selection signatures. In contrast, NC-L
458 ORFs did not show evidence of purifying selection at the protein level. Some conserved ORFs
459 in lncRNAs are likely to encode functional micropeptides. Differences between observed and
460 expected PN/PS were assessed with a chi-square test, * p-value < 0.05 , *** p-value $< 10^{-5}$.
461 Error bars indicate the standard error of the sample proportion. Numbers of ORFs for the
462 different categories are also displayed.

463

464 **Figure 3. Three nucleotide periodicity of translated ORFs.** The mapping of Ribo-Seq reads

465 on different types of ORFs is shown. The Y axis represents the log-number of reads, the X axis
466 the positions in the ORF. The reads show strong frame bias in the functional (conserved) and
467 the neutral (NC-L) examples, with a preponderance of in-frame reads (green) versus off-frame
468 reads (red and blue), while there is no consistent frame bias in the negative control (SNORA18).
469 The exon/intron structure and the amino acid sequence for translated ORFs is also shown.

470

471 **Figure 4. Factors influencing the translation of neutrally evolving ORFs.** **a.** Influence of
472 coding score in the translatability of neutrally evolving ORFs. Translated ORFs showed
473 significantly higher coding score than non-translated ORFs, both sets had significantly higher
474 coding scores than introns (Wilcoxon test p-value < 10^{-5} , indicated by ***). **b.** Influence of coding
475 score in the translatability of ORFs after controlling for gene expression values. Translated
476 ORFs showed significantly higher coding score values than non-translated ORFs (Wilcoxon test
477 p-value < 10^{-5}). **c.** Influence of gene expression level and ORF length in the translatability of
478 neutral ORFs after controlling for coding score. Translated ORFs showed significantly higher
479 FPKM values than non-translated ORFs (Wilcoxon test p-value < 10^{-5}); differences in length
480 were not significant.

References

- 481 1. Kutter, C. *et al.* Rapid turnover of long noncoding RNAs and the evolution of gene
482 expression. *PLoS Genet.* **8**, e1002841 (2012).
- 483 2. Wiberg, R. A. W. *et al.* Assessing Recent Selection and Functionality at Long Noncoding
484 RNA Loci in the Mouse Genome. *Genome Biol. Evol.* **7**, 2432–44 (2015).
- 485 3. Ruiz-Orera, J. *et al.* Origins of De Novo Genes in Human and Chimpanzee. *PLOS*
486 *Genet.* **11**, e1005721 (2015).
- 487 4. Ji, Z., Song, R., Regev, A. & Struhl, K. Many lncRNAs, 5'UTRs, and pseudogenes are
488 translated and some are likely to express functional proteins. *Elife* **4**, e08890 (2015).
- 489 5. Raj, A. *et al.* Thousands of novel translated open reading frames in humans inferred by
490 ribosome footprint profiling. *Elife* **5**, e13328 (2016).
- 491 6. Ingolia, N. T., Lareau, L. F. & Weissman, J. S. Ribosome profiling of mouse embryonic

- 492 stem cells reveals the complexity and dynamics of mammalian proteomes. *Cell* **147**,
493 789–802 (2011).
- 494 7. Ingolia, N. T. *et al.* Ribosome Profiling Reveals Pervasive Translation Outside of
495 Annotated Protein-Coding Genes. *Cell Rep.* **8**, 1365–79 (2014).
- 496 8. Ruiz-Orera, J., Messeguer, X., Subirana, J. A. & Alba, M. M. Long non-coding RNAs as
497 a source of new peptides. *Elife* **3**, e03523 (2014).
- 498 9. Wilson, B. A. & Masel, J. Putatively noncoding transcripts show extensive association
499 with ribosomes. *Genome Biol. Evol.* **3**, 1245–52 (2011).
- 500 10. Couso, J.-P. & Patraquim, P. Classification and function of small open reading frames.
501 *Nat. Rev. Mol. Cell Biol.* **18**, 575–589 (2017).
- 502 11. Ingolia, N. T., Ghaemmaghami, S., Newman, J. R. S. & Weissman, J. S. Genome-wide
503 analysis in vivo of translation with nucleotide resolution using ribosome profiling. *Science*
504 **324**, 218–23 (2009).
- 505 12. Bazzini, A. A. *et al.* Identification of small ORFs in vertebrates using ribosome
506 footprinting and evolutionary conservation. *EMBO J.* **33**, 981–93 (2014).
- 507 13. Calviello, L. *et al.* Detecting actively translated open reading frames in ribosome profiling
508 data. *Nat Meth* **13**, 165–170 (2016).
- 509 14. Aspden, J. L. *et al.* Extensive translation of small ORFs revealed by Poly-Ribo-Seq. *Elife*
510 e03528 (2014).
- 511 15. Mackowiak, S. D. *et al.* Extensive identification and analysis of conserved small ORFs in
512 animals. *Genome Biol.* **16**, 1–21 (2015).
- 513 16. Begun, D. J., Lindfors, H. A., Kern, A. D. & Jones, C. D. Evidence for de Novo Evolution
514 of Testis-Expressed Genes in the *Drosophila yakuba*/*Drosophila erecta* Clade. *Genetics*
515 **176**, 1131–1137 (2006).
- 516 17. Tautz, D. & Domazet-Lošo, T. The evolutionary origin of orphan genes. *Nat. Rev. Genet.*
517 **12**, 692–702 (2011).
- 518 18. McLysaght, A. & Hurst, L. D. Open questions in the study of de novo genes: what, how
519 and why. *Nat. Rev. Genet.* **17**, 567–578 (2016).
- 520 19. Zhao, L., Saelao, P., Jones, C. D. & Begun, D. J. Origin and spread of de novo genes in
521 *Drosophila melanogaster* populations. *Science* **343**, 769–72 (2014).

- 522 20. Carvunis, A.-R. *et al.* Proto-genes and de novo gene birth. *Nature* **487**, 370–4 (2012).
- 523 21. Toll-Riera, M. *et al.* Origin of primate orphan genes: a comparative genomics approach.
524 *Mol. Biol. Evol.* **26**, 603–12 (2009).
- 525 22. Cai, J. J. & Petrov, D. A. Relaxed purifying selection and possibly high rate of adaptation
526 in primate lineage-specific genes. *Genome Biol. Evol.* **2**, 393–409 (2010).
- 527 23. Chen, S., Zhang, Y. E. & Long, M. New Genes in *Drosophila* Quickly Become Essential.
528 *Science* (80-.). **330**, 1682–1685 (2010).
- 529 24. Reinhardt, J. A. *et al.* De novo ORFs in *Drosophila* are important to organismal fitness
530 and evolved rapidly from previously non-coding sequences. *PLoS Genet.* **9**, e1003860
531 (2013).
- 532 25. Sunyaev, S., Kondrashov, F. A., Bork, P. & Ramensky, V. Impact of selection, mutation
533 rate and genetic drift on human genetic variation. *Hum. Mol. Genet.* **12**, 3325–3330
534 (2003).
- 535 26. Gayà-Vidal, M. & Albà, M. M. Uncovering adaptive evolution in the human lineage. *BMC*
536 *Genomics* **15**, 599 (2014).
- 537 27. Harr, B. *et al.* Genomic resources for wild populations of the house mouse, *Mus*
538 *musculus* and its close relative *Mus spretus*. *Sci. Data* **3**, 160075 (2016).
- 539 28. Buck-Koehntop, B. A., Mascioni, A., Buffy, J. J. & Veglia, G. Structure, dynamics, and
540 membrane topology of stannin: a mediator of neuronal cell apoptosis induced by
541 trimethyltin chloride. *J. Mol. Biol.* **354**, 652–65 (2005).
- 542 29. Pueyo, J. I. *et al.* Hemotin, a Regulator of Phagocytosis Encoded by a Small ORF and
543 Conserved across Metazoans. *PLoS Biol* **14**, e1002395 (2016).
- 544 30. Plotkin, J. B. & Kudla, G. Synonymous but not the same: the causes and consequences
545 of codon bias. *Nat. Rev. Genet.* (2011). doi:10.1038/nrg2899
- 546 31. Vizcaíno, J. A. *et al.* 2016 update of the PRIDE database and its related tools. *Nucleic*
547 *Acids Res.* **44**, D447–D456 (2016).
- 548 32. Slavoff, S. A. *et al.* Peptidomic discovery of short open reading frame-encoded peptides
549 in human cells. *Nat. Chem. Biol.* **9**, 59–64 (2013).
- 550 33. Heinen, T. J. A. J., Staubach, F., Häming, D. & Tautz, D. Emergence of a new gene from
551 an intergenic region. *Curr. Biol.* **19**, 1527–31 (2009).

- 552 34. Dana, A. & Tuller, T. The effect of tRNA levels on decoding times of mRNA codons.
553 *Nucleic Acids Res.* **42**, 9171–81 (2014).
- 554 35. Yu, C. *et al.* Codon Usage Influences the Local Rate of Translation Elongation to
555 Regulate Co-translational Protein Folding. *Mol. Cell* **59**, 744–754 (2015).
- 556 36. Presnyak, V. *et al.* Codon optimality is a major determinant of mRNA stability. *Cell* **160**,
557 1111–24 (2015).
- 558 37. Schlötterer, C. Genes from scratch – the evolutionary fate of de novo genes. *Trends*
559 *Genet.* (2015).
- 560 38. Okazaki, Y. *et al.* Analysis of the mouse transcriptome based on functional annotation of
561 60,770 full-length cDNAs. *Nature* **420**, 563–73 (2002).
- 562 39. Neme, R. & Tautz, D. Fast turnover of genome transcription across evolutionary time
563 exposes entire non-coding DNA to de novo gene emergence. *Elife* **5**, e09977 (2016).
- 564 40. Lynch, M. & Marinov, G. K. The bioenergetic costs of a gene. *Proc. Natl. Acad. Sci. U. S.*
565 *A.* **112**, 15690–5 (2015).
- 566 41. Wilson, B. A., Foy, S. G., Neme, R. & Masel, J. Young Genes are Highly Disordered as
567 Predicted by the Preadaptation Hypothesis of De Novo Gene Birth. *Nat. Ecol. Evol.* **1**,
568 146 (2017).
- 569 42. Kaiser, C. A., Preuss, D., Grisafi, P. & Botstein, D. Many random sequences functionally
570 replace the secretion signal sequence of yeast invertase. *Science* **235**, 312–7 (1987).
- 571 43. Keefe, A. D. & Szostak, J. W. Functional proteins from a random-sequence library.
572 *Nature* **410**, 715–718 (2001).
- 573 44. Neme, R., Amador, C., Yildirim, B., McConnell, E. & Tautz, D. Random sequences are
574 an abundant source of bioactive RNAs or peptides. *Nat. Ecol. Evol.* **1**, 217 (2017).
- 575 45. Soumillon, M. *et al.* Cellular source and mechanisms of high transcriptome complexity in
576 the mammalian testis. *Cell Rep.* **3**, 2179–90 (2013).
- 577 46. Necsulea, A. *et al.* The evolution of lncRNA repertoires and expression patterns in
578 tetrapods. *Nature* **505**, 635–40 (2014).
- 579 47. Smeds, L. & Künstner, A. ConDeTri--a content dependent read trimmer for Illumina data.
580 *PLoS One* **6**, e26314 (2011).
- 581 48. Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of

- 582 insertions, deletions and gene fusions. *Genome Biol.* **14**, R36 (2013).
- 583 49. Pertea, M. *et al.* StringTie enables improved reconstruction of a transcriptome from
584 RNA-seq reads. *Nat Biotech* **33**, 290–295 (2015).
- 585 50. Luis Villanueva-Cañas, J. *et al.* New Genes and Functional Innovation in Mammals.
586 *Genome Biol. Evol.* **9**, 1886–1900 (2017).
- 587 51. Gonzalez, C. *et al.* Ribosome profiling reveals a cell-type-specific translational landscape
588 in brain tumors. *J. Neurosci.* **34**, 10924–36 (2014).
- 589 52. Castañeda, J. *et al.* Reduced pachytene piRNAs and translation underlie spermiogenic
590 arrest in Maelstrom mutant mice. *EMBO J.* **33**, 1999–2019 (2014).
- 591 53. Guo, H., Ingolia, N. T., Weissman, J. S. & Bartel, D. P. Mammalian microRNAs
592 predominantly act to decrease target mRNA levels. *Nature* **466**, 835–40 (2010).
- 593 54. Diaz-Munoz, M. D. *et al.* The RNA-binding protein HuR is essential for the B cell
594 antibody response. *Nat Immunol* **16**, 415–425 (2015).
- 595 55. Cho, J. *et al.* Multiple repressive mechanisms in the hippocampus during memory
596 formation. *Science* **350**, 82–87 (2015).
- 597 56. Sedlazeck, F. J., Rescheneder, P. & von Haeseler, A. NextGenMap: fast and accurate
598 read mapping in highly polymorphic genomes. *Bioinformatics* **29**, 2790–2791 (2013).
- 599 57. Altschul, S. F. *et al.* Gapped BLAST and PSI-BLAST: a new generation of protein
600 database search programs. *Nucleic Acids Res.* **25**, 3389–402 (1997).
- 601 58. Karolchik, D. *et al.* The UCSC Genome Browser database: 2014 update. *Nucleic Acids*
602 *Res.* **42**, D764–D770 (2014).
- 603 59. Rosenberg, M. S., Subramanian, S. & Kumar, S. Patterns of Transitional Mutation
604 Biases Within and Among Mammalian Genomes. *Mol. Biol. Evol.* **20**, 988–993 (2003).
- 605 60. R Development Core Team. R: a language and environment for statistical computing. R
606 Foundation for Statistical Computing, Vienna, Austria (2016).

607

608 **Acknowledgements**

609

610 We are grateful for valuable discussions with many colleagues during this study. The work was
611 funded by grants BFU2012-36820, BFU2015-65235-P and TIN2015-69175-C4-3-R from

612 Ministerio de Economía e Innovación (Spanish Government) and co-funded by FEDER (EC).
613 We also received funding from Agència de Gestió d'Ajuts Universitaris i de Recerca Generalitat
614 de Catalunya (AGAUR), grant number 2014SGR1121.

615

616 **Author contributions**

617 J.R-O. and M.M.A. conceived the study, interpreted the data and wrote the paper. J.R-O.
618 performed most of the analyses including the transcript assemblies, identification of translated
619 open reading frames, BLAST searches, SNP mapping and generation of controls. J.R-O., P.V-
620 G. and J.L.V-C wrote code and performed analyses on the coding score. X.M. wrote code to
621 calculate the expected SNP frequencies. M.M.A. coordinated the study.

622

623 **Competing interests**

624 The authors declare no competing financial interests.