
SUBFUNCTIONALISATION OF PARALOGOUS GENES AND EVOLUTION OF DIFFERENTIAL CODON USAGE PREFERENCES: THE SHOWCASE OF POLYPYRIMIDINE TRACT BINDING PROTEINS

Jérôme Bourret^{1, †}, Fanni Borvetó^{1, †, *}, and Ignacio G. Bravo¹

¹Laboratoire MIVEGEC (CNRS IRD Univ Montpellier), Centre National de la Recherche Scientifique (CNRS),
Montpellier, France

[†]These authors contributed equally to this work

ABSTRACT

1 Gene paralogs are copies of an ancestral gene that appear after gene or full genome duplication.
2 When two sister gene copies are maintained in the genome, redundancy may release certain evolu-
3 tionary pressures, allowing one of them to access novel functions. Here, we focused our study on
4 gene paralogs on the evolutionary history of the three polypyrimidine tract binding protein genes
5 (*PTBP*) and their concurrent evolution of differential codon usage preferences (CUPrefs) in verte-
6 brate species.

7 *PTBP1-3* show high identity at the amino acid level (up to 80%), but display strongly different
8 nucleotide composition, divergent CUPrefs and, in humans, distinct tissue-specific expression lev-
9 els. Our phylogenetic inference results show that the duplication events leading to the three extant
10 *PTBP1-3* lineages predate the basal diversification within vertebrates, and genomic context analy-
11 sis illustrates that synteny has been well preserved over time for the three paralogs. We identify
12 a distinct evolutionary pattern towards GC3-enriching substitutions in *PTBP1*, concurrent with an
13 enrichment in frequently used codons and with a tissue-wide expression. In contrast, *PTBP2s* are
14 enriched in AT-ending, rare codons, and display tissue-restricted expression. As a result of this sub-
15 stitution trend, CUPrefs are sharply different between mammalian *PTBP1s* and the rest of *PTBPs*.
16 Genomic context analysis shows that GC3-rich nucleotide composition in *PTBP1s* is driven by local
17 substitution processes, while the evidence in this direction is thinner for *PTBP2-3*. An actual lack
18 of co-variation between the observed GC composition of *PTBP2-3* and that of the surrounding non-
19 coding genomic environment would raise an interrogation on the origin of CUPrefs, warranting
20 further research on a putative tissue-specific translational selection. Finally, we communicate an
21 intriguing trend for the use of the UUG-Leu codon, which matches the trends of AT-ending codons.

*Corresponding author. email : fanni.borveto@uni-ulm.de

Evolution of codon usage preferences in paralogous genes

22 We interpret that our results are compatible with an scenario in which a combination of directional
23 mutation–selection processes would have differentially shaped CUPrefs of *PTBPs* in Vertebrates:
24 the observed GC-enrichment of *PTBP1* in Mammals may be linked to genomic location and to the
25 strong and broad tissue-expression, while AT-enrichment of *PTBP2* and *PTBP3* would be associated
26 with rare CUPrefs and thus, possibly to specialized spatio-temporal expression. Our interpretation
27 is coherent with a gene subfunctionalisation process by differential expression regulation associated
28 to the evolution of specific CUPrefs.

29 **Keywords** Codon usage bias, codon usage preferences, gene duplication, paralog, ortholog, evolution, mutation-
30 selection, nucleotide composition, tissue-specific expression

31 **1 Significance Statement**

32 In vertebrates, *PTBP* paralogs display strong differences in gene composition, gene expression regulation, and their
33 expression in cell culture depends on their codon usage preferences. We show that placental mammals *PTBP1* have
34 become GC-rich because of local substitution pressures, resulting in an enrichment of frequently used codons and in a
35 strong, tissue-wide expression. On the contrary, *PTBP2* in vertebrates are AT-rich, with a lower contribution of local
36 substitution processes to their specific nucleotide composition, show high frequency of rare codons and in placental
37 mammals display a restricted expression pattern contrasting to that of *PTBP1*. The systematic study of composition
38 and expression patterns of gene paralogs can help understand the complex mutation-selection interplay that shape
39 codon usage bias in multicellular organisms.

40 **2 Introduction**

41 During mRNA translation ribosomes assemble proteins by specific amino acid linear polymerisation guided by the
42 successive reading of mRNA nucleotide triplets, called codons. Each time a codon is read, it is chemically compared
43 to the set of available tRNAs' anticodons. Upon codon-anticodon match, the ribosome loads the tRNA and adds the
44 associated amino acid to the nascent protein. The main 20 amino acids are encoded by 61 codons, so that multiple
45 codons are associated with the same amino acid. These are named synonymous codons (Nirenberg and Matthaei,
46 1961; Khorana et al., 1966). Codon Usage Preferences (CUPrefs) refer to the differential usage of synonymous
47 codons between species, between genes, or between genomic regions in the same genome (Grantham et al., 1980;
48 Carbone et al., 2003). Mutation and selection are the two main forces shaping CUPrefs (Duret, 2002; Chamary et al.,
49 2006; Plotkin and Kudla, 2011). Mutational biases relate to directional mechanistic biases during genome replication
50 (Reijns et al., 2015; Apostolou-Karampelis et al., 2016), during genome repair (Lujan et al., 2012), or during recom-
51 bination (Pouyet et al., 2017), preferentially introducing one nucleotide over others or inducing recombination and
52 maintaining genomic regions depending on their composition. Mutational biases are well described in prokaryotes
53 and eukaryotes, ranging from simple molecular preferences towards 3' A-ending in the Taq polymerase (Clark, 1988)
54 to complex GC-biased gene conversion in vertebrates (Pouyet et al., 2017). Selective forces shaping CUPrefs are often
55 described as translational selection. This notion refers to the ensemble of mechanistic steps and interactions during
56 translation that are affected by the particular CUPrefs of the mRNA, so that the choice of certain codons at certain
57 positions may actually enhance the translation process and can be subject to selection (Bulmer, 1991). Translational
58 selection covers thus codon-independent effects on mRNA secondary structure, overall stability, and subcellular
59 location (Presnyak et al., 2015; Novoa and Ribas de Pouplana, 2012), but also codon-mediated effects acting on
60 mRNA maturation, programmed frameshifts, translation speed and accuracy, or protein folding (Caliskan et al., 2015;
61 Mordstein et al., 2020; Spencer and Barral, 2012). Translational selection has been demonstrated in prokaryotes and in
62 some eukaryotes (Satapathy et al., 2016; Percudani et al., 1997; Duret and Mouchiroud, 1999; Whittle and Extavour,
63 2016), often in the context of tRNA availability (Ikemura, 1981). However, its very existence in Vertebrates remains
64 highly debated (Pouyet et al., 2017; Galtier et al., 2018).

65

66 Homologous genes share a common origin either by speciation (orthology) or by duplication events (paralogy)
67 (Sonnhammer and Koonin, 2002). Upon gene (or full genome) duplication, the new genome will contain two copies
68 of the original gene, referred to as in-paralogs. After speciation, each daughter cell will inherit one couple of
69 paralogs, *i.e.* one copy of each ortholog (Koonin, 2005). The emergence of paralogs upon duplication may release
70 the evolutionary constraints on the individual genes. Evolution can thus potentially lead to function specialisation,
71 such as evolving a particular substrate preferences, or engaging each paralog on specific enzyme activity preferences
72 in the case of promiscuous enzymes (Copley, 2020). Gene duplication can also allow one paralog to explore broader
73 sequence space and to evolve radically novel functions, while the remaining counterpart can continue to assure the
74 original function.

75

76 The starting point for our research are the experimental observations by Robinson and coworkers reporting differential
77 expression of the polypyrimidine tract binding protein (*PTBP*) human paralogs as a function of their nucleotide com-

78 position (Robinson et al., 2008). Vertebrates genomes encode for three in-paralogous versions of the *PTBP* genes, all
79 of them fulfilling similar functions in the cell: they form a class of hnRNP RNA-Binding Proteins that are involved in
80 the modulation of mRNAs alternative splicing (Pina et al., 2018). Within the same genome, the three paralogs display
81 high amino-acid sequence similarity, around 70% in humans, and with similar overall values in vertebrates (Pina et al.,
82 2018).

83 Despite the high resemblance at the protein level, the three *PTBP* paralogs sharply differ in nucleotide composition,
84 CUPrefs, and supposedly in tissue expression pattern. In humans *PTBP1* is enriched in GC3-rich synonymous codons
85 and is widely expressed in all tissues, while *PTBP2* and *PTBP3* are AT3-rich and display an enhanced expression in the
86 brain and in hematopoietic cells respectively (Supplementary Material, Figure S1). Robinson and coworkers studied
87 the expression in human cells in culture of all three human *PTBP* paralogous genes placed under the control of the
88 same promoter. They showed that the GC-rich paralog *PTBP1* was more highly expressed than the AT-rich ones, and
89 that the expression of the AT-rich paralog *PTBP2* could be enhanced by synonymous codons recoding towards the use
90 of GC-rich codons (Robinson et al., 2008). Here we have built on the evolutionary foundations of this observation and
91 extended the analyses of CUPrefs to *PTBP* paralogs in vertebrate genomes. Our results suggest that paralog-specific
92 directional changes in CUPrefs in mammalian *PTBP* concurred with a process of subfunctionalisation by differential
93 tissue pattern expression of the three paralogous genes.

94 **3 Material and Methods**

95 *Sequence retrieval*

96 We assembled a dataset of DNA sequences from 47 mammalian and 27 non-mammalian Vertebrates, and 3 from
97 protostomes. Using the BLAST function on the nucleotide database of NCBI (NCBI Resource Coordinators, 2018)
98 taking each of the human *PTBP* paralogs as references we looked for genes already annotated as *PTBP* orthologs (see
99 supplementary Table S2 for accession numbers). We could retrieve the corresponding three orthologs in all Vertebrate
100 species screened, except for the European rabbit *Oryctolagus cuniculus*, lacking *PTBP1*, and from the rifleman bird
101 *Acanthisitta chloris*, lacking *PTBP3*. The final vertebrate dataset contained 75 *PTBP1*, 76 *PTBP2* and 75 *PTBP3*
102 sequences. As outgroups for the analysis, we retrieved the orthologous genes from three protostome genomes, which
103 contained a single *PTBP* homolog per genome. Our final dataset was consistent with the descriptions available in
104 ENSEMBL and ORTHOMAM for the *PTBP* orthologs (Yates et al., 2020; Scornavacca et al., 2019; Pina et al., 2018).
105 From the original dataset, we identified a subset of nine mammalian and six non-mammalian vertebrates species with
106 a good annotation of the *PTBP* chromosome context. For these 15 species we retrieved synteny and composition
107 information on the *PTBP* flanking regions and introns (Supplementary Table S3). Because of annotation hazards,
108 intronic and flanking regions information were missing for some *PTBPs* in the African elephant *Loxodonta africana*,
109 Schlegel's Japanese Gecko *Gekko japonicus*, and the whale shark *Rhincodon typus* assemblies. For the selected 15
110 species the values for codon adaptation index (CAI) (Sharp and Li, 1987) and codon usage similarity index (COUSIN)
111 (Bourret et al., 2019) were calculated using the COUSIN server (available at <https://cousin.ird.fr>) (Supplementary
112 Table S4).

113 *Codon Usage analysis*

114 For each *PTBP* gene we calculated codon composition, GC, GC3 and CUPrefs analyses via the COUSIN tool
115 (Bourret et al., 2019). For each *PTBP* gene we constructed a vector of 59 positions with the relative frequencies
116 of all synonymous codons. We applied different approaches to reduce information dimension for the analysis of
117 CUPrefs, on the 229 59-dimension vectors: I) a k-means clustering; ii) a hierarchical clustering; and iii) a principal
118 component analysis (PCA). Statistical analyses were performed using the ape and ade4 R packages and JMP v14.3.0.
119 Correlation between matrices was assessed via the Mantel test. Non-parametric comparisons were performed using
120 the Wilcoxon-Mann-Whitney test for assessing differences between the median values of the corresponding variable
121 (either GC or GC3) among paralogs, and the Wilcoxon signed rank test for paired comparisons of the values for cor-
122 responding variable (either GC or GC3) for paralogs within the same genome. For the 15 species with well-annotated
123 genomes we analyzed by a stepwise linear fit the correlation of paralog GC3 with two local compositional variables
124 of the corresponding gene (GC content of intronic and flanking regions) and with three global compositional variables
125 for the corresponding genomes (global GC3 in the complete genomic ORFome, global GC content in all introns, and
126 global GC content in all flanking regions).

127 ***Alignment and phylogenetic analyses***

128 First, all sequences were aligned together, and we constructed a phylogenetic tree to verify whether each paralog as-
129 sembly was monophyletic (Supplementary Figure S13). This was actually the case, and in this unbiased preliminary
130 analysis all *PTBP1-3* were respectively monophyletic. Thus, to generate more robust alignments without introducing
131 artefacts due to large evolutionary distances between in-paralogs, we proceeded stepwise, as follows: i) we aligned
132 separately at the amino acid level each set of *PTBP* paralog sequences of mammals and non-mammalian Vertebrates;
133 ii) for each *PTBP* paralog we merged the alignments for mammals and for non mammals, obtaining the three *PTBP1*,
134 *PTBP2* and *PTBP3* alignments for all Vertebrates; iii) we combined the three alignments for each paralog into a sin-
135 gle one; iv) we aligned the outgroup sequences to the global Vertebrate *PTBPs* alignment. All alignment steps were
136 performed using MAFFT (Katoh et al., 2002). The final amino acid alignment was used to obtain the codon-based
137 nucleotide alignment. The codon-based alignment was trimmed using Gblocks (Castresana, 2000) (Data available on
138 Zenodo) Phylogenetic inference was performed at the amino acid and at the nucleotide level using RAxML v8.2.9,
139 bootstrapping over 1000 cycles (Stamatakis, 2014). For nucleotides we used codon-based partitions and applied the
140 GTR+G4 model while for amino acids we applied the LG+G4 model. For the 79 species used in the analyses we
141 retrieved a species-tree from the TimeTree tool (Kumar et al., 2017). Distances between phylogenetic trees were com-
142 puted using the Robinson-Foulds index, which accounts for differences in topology (Robinson and Foulds, 1981), and
143 the K-tree score, which accounts for differences in both topology and branch length (Soria-Carrasco et al., 2007). We
144 then calculated pairwise distances between branches on the nucleotide and amino acid based trees and compared them
145 against CUPrefs-based pairwise distances to measure the impact of CUPrefs on the phylogeny. After phylogenetic
146 inference, we computed marginal ancestral states for the respectively most recent common ancestors at the nucleotide
147 level of each paralog, using RAxML. For each position the base with the maximum probability was used, and the sites
148 for which RAxML could not infer with certainty the base were marked as missing data. We found 14%, 18% and
149 10% of missing bases respectively in *PTBP1*, *PTBP2* and *PTBP3*. Using these ancestral sequences we estimated the
150 number of synonymous and non-synonymous substitutions of each extant sequence to the corresponding most recent
151 common ancestor. We then compared the substitution matrices via a PCA analysis.

152 4 Results

153 *Vertebrate PTBP paralogs differ in nucleotide composition*

154 In order to understand the evolutionary history of *PTBP* genes, we performed first a nucleotide composition and
155 CUPrefs analysis on the three paralogs in 79 species. Overall, *PTBP1* are GC-richer than *PTBP2* and *PTBP3* (re-
156 spective mean percentages 55.9, 42.3 and 44.9 for GC content and 69.5, 33.4 and 38.3 for GC3 content; Figure 1). In
157 addition, *PTBP1*s show a difference in GC3 between mammalian and non-mammalian genes (respectively 79.8 against
158 59.9 mean percentages). A linear regression model followed by a Tukey's honest significant differences analysis for
159 GC3 using as explanatory levels paralog (*i.e.* *PTBP1-3*), taxonomy (*i.e.* mammalian or non-mammalian), and their
160 interaction identifies three main groups of *PTBPs* (Table 1): a first one corresponding to mammalian *PTBP1*, a second
161 one grouping non-mammalian *PTBP1*, and a third one encompassing all *PTBP2* and *PTBP3*. The largest explanatory
162 factor for GC3 was the paralog *PTBP1-3*, accounting alone for 65% of the variance, while the interaction between the
163 levels taxonomy and paralog captured around 15% of the remaining variance (Table 1). These trends are confirmed
164 when performing paired comparisons between paralogs present in the same mammalian genome, with significant dif-
165 ferences in GC3 content in the following order: *PTBP1* > *PTBP3* > *PTBP2* (Wilcoxon signed rank test: *PTBP1* vs
166 *PTBP2*, mean diff=48.0, S=539.50, p-value <0.0001; *PTBP1* vs *PTBP3*, mean diff=43.5, S=517.50, p-value <0.0001;
167 *PTBP3* vs *PTBP2*, mean diff=4.5, S=406.50, p-value <0.0001). Note that even if all of them significantly different,
168 the mean paired differences in GC3 between *PTBP1* and *PTBP2-3* are ten times larger than the corresponding mean
169 paired differences between *PTBP2* and *PTBP3*.

170 After our model fit, an analysis of the distribution of the residuals between observed and expected values to the data
171 allows to identify a number of outliers species with interesting taxonomical patterns in compositional deviation (Table
172 2). For non mammals, the three *PTBP* paralogs in the rainbow trout *Oncorhynchus mykiss* genome display high
173 GC3 content (between 67% and 76%), all of them significantly higher than model-predicted values (expected values
174 between 36% and 51%). A similar case occurs for the zebrafish *Danio rerio* genome: the three paralogs display
175 GC3 values around 58%, which for *PTBP2* and *PTBP3* paralogs are significantly higher than predicted by the model
176 (expected values around 38%). Very interestingly, for the monotreme platypus *Ornithorhynchus anatinus* as well as
177 for the three marsupials in the dataset: the Tasmanian devil *Sarcophilus harrisii*, the koala *Phascolarctos cinereus* and
178 the grey short-tailed opossum *Monodelphis domestica*, their *PTBP1* genes present similar GC3 content around 47%,
179 which is significantly lower than predicted by the model (expected values around 79%).

180 In many vertebrate species, strong compositional heterogeneities are observed along chromosomes with an arrange-
181 ment of AT-rich and GC-rich regions, often referred to as "isochores". To explore the influence of this genomic
182 environment on the nucleotide composition of *PTBPs*, we analyzed for 15 species with well-annotated genomes the
183 correlation of paralog GC3 with two local compositional variables of the corresponding gene (GC content of intronic
184 and flanking regions) and with three global compositional variables for the corresponding genomes (global GC3 in
185 the complete genomic ORFome, global GC content in all introns, and global GC content in all flanking regions)(Table
186 3 and Figure 2). First, for *D. rerio* the GC3 composition of *PTBP2* and *PTBP3* is clearly different from the rest,
187 in line with the outlier results presented in Table 2. We have thus excluded the zebra fish values and performed an
188 individual as well as a stepwise linear fit to explain the variance in GC3 composition by the variance in the local and
189 global compositional variables mentioned above (Table 3). For all three *PTBPs* the local GC content explains best the

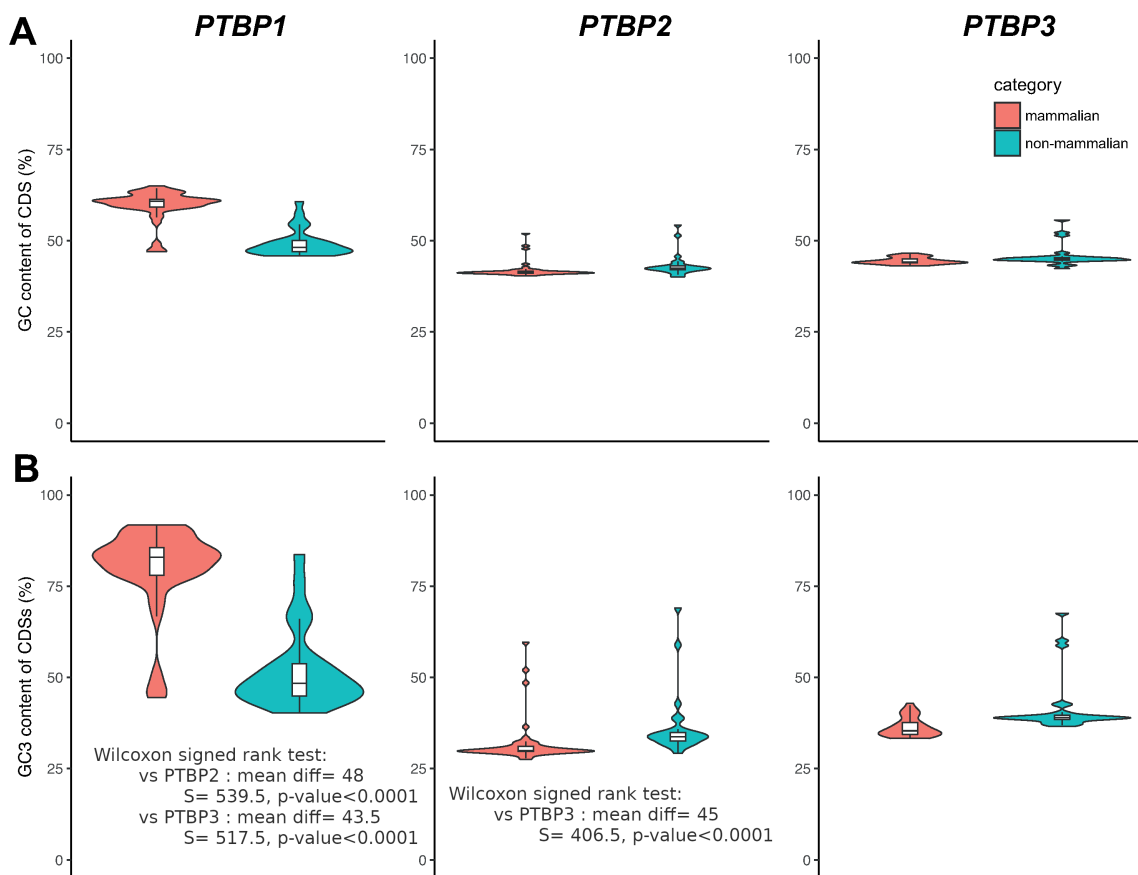


Figure 1: **GC content (A) and GC3 content (B) of Vertebrates *PTBPs*.** Violin plots display the overall distribution, while box and whiskers display median, quartiles and 95% of the corresponding values for mammalian (red) and non-mammalian (blue) individual genomes. The results of a the paired Wilcoxon signed rank tests between overall GC3 content of paralogs in the same genome are indicated in the inboxes.

190 corresponding GC3 content, but with strong differences between paralogs: while variation in the local composition
191 captures almost perfectly variation in the GC3 content of *PTBP1* ($R^2=0.97$) and relatively well in the case of *PTBP2*
192 ($R^2=0.46$), the fraction of variance explained by the local composition significantly drops for *PTBP3* ($R^2=0.15$). It
193 must be noted nevertheless that the GC3 variable ranges are different among paralogs, so that variation in GC3 values
194 for *PTBP1* (roughly between 40% and 90%) is larger than for *PTBP2-3* (respectively 29%-38% and 34%-46%). This
195 larger variable span in the case of *PTBP1* may allow for an increased power for detecting a significant correlation in
196 composition values for this paralog.

197 Vertebrate *PTBP* paralogs differ in CUPrefs

198 For each *PTBP* coding sequence we extracted the relative frequencies of synonymous codons and performed different
199 approaches to reduce information dimension and visualise CUPrefs trends. The results of a principal component
200 analysis (PCA) are shown in Figure 3 as well as in Supplementary Figure S5. The first PCA axis captured 68.9% of
201 the variance, far before the second and the third axes (respectively 6.7% and 3.2%). Codons segregate in the first axis
202 by their GC3 composition, the only exception being the UUG-Leu codon, which grouped together with AT-ending

203 codons. This first axis differentiates mammalian *PTBP1*s on the one hand and *PTBP2*s and *PTBP3*s on the other hand.
 204 Non-mammalian *PTBP1*s scatter between mammalian *PTBP1*s and *PTBP3*s, along with the protostomates *PTBP*s.
 205 In the second PCA axis the only obvious (but nevertheless cryptic) codon-structure trends are: i) the split between
 206 C-ending and G-ending codons, but not between U-ending and A-ending codons; and ii) the large contribution in
 207 opposite directions to this second axis of the AGA and AGG-Arginine codons. This second PCA axis differentiates
 208 *PTBP2*s from *PTBP3*s paralogs, consistent with these composition trends. A paired-comparison confirms that *PTBP3*s
 209 are richer in C-ending codons than *PTBP2*s in the same genome, respectively 21.7% against 15.4% (Wilcoxon signed
 210 rank test: mean diff=6.2, S=1184.0, p-value <0.0001).

211 As an additional way to identify groups of genes with similar CUPrefs, we applied a hierarchical clustering and a
 212 k-means clustering. Both analyses mainly aggregate *PTBP* genes by their GC3 richness. The *PTBP* dendrogram
 213 resulting of the hierarchical clustering shows five main clades that cluster the paralogs with a good match to the
 214 following groups: mammalian *PTBP1*s, non-mammalian *PTBP1*s, *PTBP2*s, *PTBP3*s and a fifth group containing
 215 the protostomata *PTBP*s and a few individuals of all three paralogs (rows in clustering in Figure 3; Kappa-Fleiss
 216 consistency score = 0.76). Regarding codon clustering, the hierarchical stratification sharply splits GC-ending codons
 217 from AT-ending codons, with the only exception again of the UUG-Leu codon, which consistently groups within
 218 the AT-ending codons. The elbow approach of k-means clustering identifies an optimal number of four clusters and

Table 1: Global linear regression model and post-hoc Tukey's honest significant differences test for GC3 composition as explained variable and the explanatory levels paralog (*PTBP1-3*), taxonomy (*i.e.* mammalian or non-mammalian) and their interactions. Within each level, strata labelled with the same letter are not different from one another. Overall goodness of the fit: Adj Rsquare=0.83; F ratio=205.7; Prob > F: <0.0001. Individual effects for the levels: i) paralog: F ratio=274.3; Prob > F: <0.0001; ii) taxonomy: F ratio=27.2; Prob > F: <0.0001; iii) interaction paralog*taxonomy: F ratio=87.9; Prob > F: <0.0001.

Level	Least Sq. Mean (GC3%)	Standard error	Tukey's HSD group
Paralog			
<i>PTBP1</i>	65.87	1.00	A
<i>PTBP3</i>	39.00	1.01	B
<i>PTBP2</i>	34.03	1.00	C
Taxonomy			
mammalian	49.32	0.70	A
non-mammalian	43.28	0.92	B
Paralog*Taxonomy			
<i>PTBP1</i> , mammalian	79.81	1.22	A
<i>PTBP1</i> , non-mammalian	51.93	1.59	B
<i>PTBP3</i> , non-mammalian	41.64	1.62	C
<i>PTBP3</i> , mammalian	36.36	1.22	C, D
<i>PTBP2</i> , non-mammalian	36.27	1.59	C, D
<i>PTBP2</i> , mammalian	31.79	1.20	D

219 separates the paralog genes with a good match as following: *PTBP1*, *PTBP2*, *PTBP3* and a group containing the
220 protostomates and individuals from all paralogs (Kappa-Fleiss consistency score = 0.75).

221 Overall, k-means clustering and hierarchical clustering, both based on the 59-dimensions vectors of the CUPrefs, are
222 congruent with one another (Kappa-Fleiss consistency score = 0.83), and largely concordant with the PCA results.
223 CUPrefs define thus groups of *PTBP* genes consistent with their orthology and taxonomy. It is interesting to note that
224 for some species the *PTBP* paralogs display unique distributions of CUPrefs, such as an overall similar CUPrefs in
225 the three *PTBP* genes of the whale shark *Rhincodon typus*, or again some shifts in nucleotide composition between
226 paralogs in the Natal long-fingered bat *Miniopterus natalensis*.

227 In order to characterise the directional CUPrefs bias of the different paralogs, we have analysed, for the 15 species
228 with well-annotated genomes described above, the match between each individual *PTBP* and the average CUPrefs of
229 the corresponding genome (Table 4). The COUSIN quantitative values compare the CUPrefs of a query sequence with
230 those of a reference (in our case the coding genome of the corresponding organism), and can be directly interpreted
231 in a qualitatively way, as described (Bourret et al., 2019). Briefly, COUSIN values around 1 reflect similar CUPrefs
232 in the query sequence and in the reference, while values around 0 reflect CUPrefs close to random in the query
233 sequence; COUSIN values above 1 reflect similar directional trends in CUPrefs in the query sequence and in the
234 reference, but with stronger bias in the query sequence; COUSIN negative values reflect opposite CUPrefs between the
235 query sequence and the reference. Our results highlight strong differences for mammalian paralogs: *PTBP1*s display
236 COUSIN values above 1 while *PTBP2*s display COUSIN values below zero. The COUSIN results and interpretation

Table 2: Individual genes with outlier values with respect to the linear regression expected values for the levels paralog (*PTBP1-3*), taxonomy (mammalian or non-mammalian) and their interactions.

Species	paralog	observed GC3 (%)	expected GC3 (%)	deviation GC3 (%)
mammalian				
<i>Desmodus rotundus</i>	<i>PTBP2</i>	59.60	31.79	27.81
<i>Miniopterus natalensis</i>	<i>PTBP2</i>	48.52	31.79	16.72
<i>Monodelphis domestica</i>	<i>PTBP1</i>	44.49	79.81	-35.32
<i>Ornithorhynchus anatinus</i>	<i>PTBP1</i>	51.14	79.81	-28.67
<i>Ornithorhynchus anatinus</i>	<i>PTBP2</i>	52.00	31.79	20.21
<i>Phascolarctos cinereus</i>	<i>PTBP1</i>	47.53	79.81	-32.28
<i>Sarcophilus harrisii</i>	<i>PTBP1</i>	45.44	79.81	-34.37
non-mammalian				
<i>Danio rerio</i>	<i>PTBP2</i>	58.89	36.27	22.62
<i>Danio rerio</i>	<i>PTBP3</i>	60.08	41.64	18.44
<i>Lepisosteus oculatus</i>	<i>PTBP3</i>	58.73	41.64	17.10
<i>Oncorhynchus mykiss</i>	<i>PTBP1</i>	76.27	51.93	24.34
<i>Oncorhynchus mykiss</i>	<i>PTBP2</i>	69.03	36.27	32.76
<i>Oncorhynchus mykiss</i>	<i>PTBP3</i>	67.58	41.64	25.95
<i>Pogona vitticeps</i>	<i>PTBP1</i>	83.68	51.93	31.75

Evolution of codon usage preferences in paralogous genes

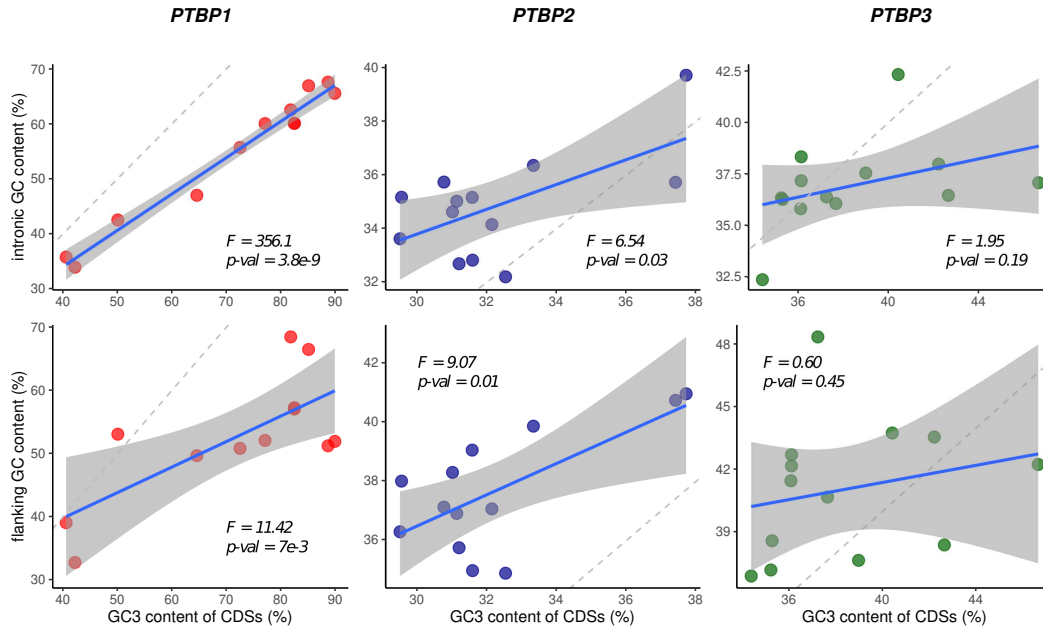


Figure 2: **Variation in GC3 content of PTBPs (x-axis) and in the GC content of the corresponding introns (A, y axis) or flanking regions (B, y axis).** Each dot represents one of the 15 individual genomes used for the genomic context analysis. For each graph, we performed a linear regression modelling (represented with the blue line for the fit and grey-shaded areas for the 95% confidence of the fit ; F-statistic and related p-values are given on the Figure); for each panel a grey line represents the $y = x$ bisector.

237 are provided in (Supplementary Figure S14). These results mean that, in mammals, *PTBP1*s are enriched in commonly
 238 used codons in a higher proportion than the average in the genome, while *PTBP2*s are enriched in rare codons to the
 239 extent that their CUPrefs go in the opposite direction to the average in the genome. As for *PTBP3* in mammals,
 240 we observe COUSIN values below 0 in most cases or very close to 0 in the case of the horse *Equus caballus* and
 241 house mouse *Mus musculus*, implying a trend towards rare codons. In non-mammals however, *PTBPs* show an overall
 242 similarity to their respective reference genomic CUPrefs.

243 **Phylogenetic reconstruction of PTBPs**

244 We explored the evolutionary relationships between *PTBPs* by phylogenetic inference at the amino acid and at the
 245 nucleotide levels (Figure 4, Supplementary Figure S10). Our final dataset contained 74 *PTBP* sequences from mam-
 246 mals (47 species within 39 families) and non mammal vertebrates (27 species within 24 families). We used the *PTBP*
 247 genes from three protostome species as outgroup. Both amino acid and nucleotide phylogenies rendered three main
 248 clades grouping the *PTBPs* by orthology, so that all *PTBP1-3* orthologs were correspondingly monophyletic. In both
 249 topologies, *PTBP1* and *PTBP3* orthologs cluster together, although the protostome outgroups are linked to the tree by
 250 a very long branch, hampering the proper identification of the Vertebrate *PTBP* tree root. Amino acid and nucleotide
 251 subtrees were largely congruent (see topology and branch length comparisons in Table5). The apparently large nodal
 252 and split distance values between nucleotide and amino acid for *PTBP2* trees stem from disagreements in very short
 253 branches, as evidenced by the lowest K-tree score for this ortholog (as a reminder, the Robinson-Foulds index exclu-

Evolution of codon usage preferences in paralogous genes

254 sively regards topology while the K-tree score combines topological and branch-length dependent distance between
 255 trees, see Material and Methods). In all three cases, internal structure of the ortholog trees essentially recapitulates
 256 species taxonomy at the higher levels (Table5). Some of the species identified by the regression analyses to display
 257 largely divergent nucleotide composition from the expected one given their taxonomy (Table 2) presented accordingly
 258 long branches in the phylogenetic reconstruction, such as *PTBP3* for *O. mykiss*, or rendered paraphyletic branching,
 259 as described above for *PTBP1* in marsupials and monotremes.

Table 3: Results for an individual (left) or for a sequential (right) least squares regression for explaining variation in GC3 composition of *PTBPs* genes, by variation of different compositional variables, either local (introns or flanking regions of the corresponding gene) or global (all coding CDS, all introns and all flanking regions in the corresponding genome), in 14 well-annotated vertebrate genomes. For the sequential fit, variables are ordered according to their contribution to the sequentially better model for the corresponding paralog, and the order may thus differ between paralogs. Variables labelled with "n.s." (not significant) do not contribute with significant additional explanatory power when added to the sequential model. BIC, Bayesian information content.

<i>PTBP1</i>					
Individual contributions			Sequential contribution		
Parameter	R ²	P value F test	Parameter	R ²	BIC
Local_GC_intron	0.9726	<0.001	Local_GC_intron	0.9726	66.4765
Local_GC_flanking	0.5345	0.0069	Local_GC_flanking	0.974 (n.s.)	68.3142
Global_GC3_exome	0.7279	0.0004	Global_GC3_exome	0.9749 (n.s.)	70.3842
Global_GC_introns	0.116	0.2786	Global_GC_flanking	0.9803(n.s.)	69.9886
Global_GC_flanking	0.1041	0.3065	Global_GC_introns	0.9806(n.s.)	72.2531
<i>PTBP2</i>					
Individual contributions			Sequential contribution		
Parameter	R ²	P value F test	Parameter	R ²	BIC
Local_GC_intron	0.3738	0.0264	Local_GC_flanking	0.4558	60.1257
Local_GC_flanking	0.4558	0.0113	Global_GC_introns	0.4895(n.s.)	61.8583
Global_GC3_exome	0.0943	0.3075	Global_GC3_exome	0.4914(n.s.)	64.3761
Global_GC_introns	0.0488	0.4684	Global_GC_flanking	0.4934(n.s.)	66.8894
Global_GC_flanking	0.0287	0.5801	Local_GC_intron	0.4974(n.s.)	69.35
<i>PTBP3</i>					
Individual contributions			Sequential contribution		
Parameter	R ²	P value F test	Parameter	R ²	BIC
Local_GC_intron	0.1554	0.1825	Local_GC_intron	0.1554	74.7338
Local_GC_flanking	0.0522	0.4528	Local_GC_flanking	0.2095(n.s.)	76.4388
Global_GC3_exome	0.0504	0.461	Global_GC_introns	0.2718(n.s.)	77.9368
Global_GC_introns	0.0002	0.9661	Global_GC3_exome	0.2938(n.s.)	80.1032
Global_GC_flanking	0.0024	0.8744	Global_GC_flanking	0.2938(n.s.)	82.667

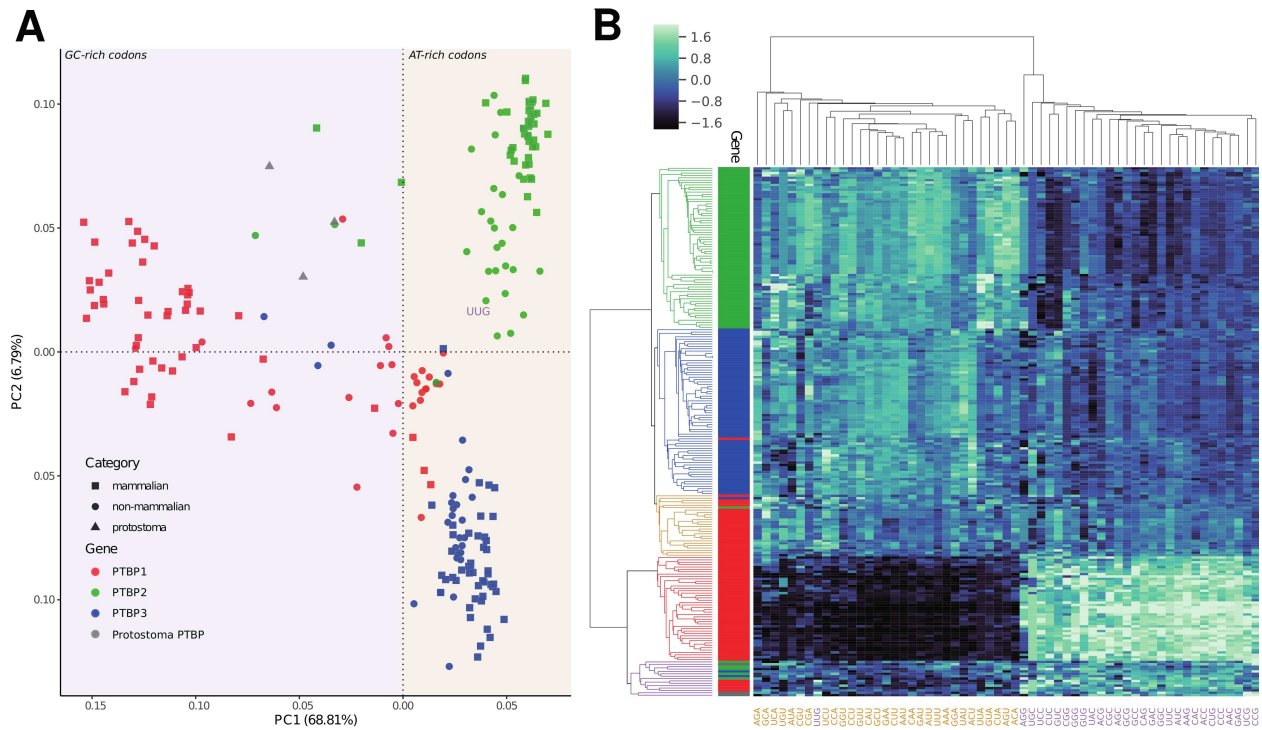


Figure 3: **CUPrefs analysis of PTBPs.** A) Plot of the two first dimensions of a PCA analysis based on the codon usage preferences of *PTBP1*s (red), *PTBP2*s (green), *PTBP3*s (blue) and protostoma (grey) individuals. Taxonomic information is included labelling mammals (squares), non-mammals (circles) and protostomes (triangles). The PCA was created using as variables the vectors of 59 positions (representing the relative frequencies of the 59 synonymous codons) for each individual gene. Shaded areas in purple (left) and orange (right) delimit the GC-rich and AT-rich grouping of codon variables according to the PCA. The UUG-Leu codon, colored in purple and placed on the Figure according to its eigenvalue, appears as a clear exception compared to the global trend of variables (Supplementary Figure S5) depicts a detailed positioning of the 59 PCA variables). The percentage of the total variance explained by each axis is shown in parenthesis. B) Heatmap of *PTBPs* individuals (rows) and synonymous codons (columns). Left dendrogram represents the hierarchical clustering of *PTBPs* based on their CUPrefs with colour codes that stand for the clusters created from this analysis. The side bar gives information on heatmap individuals regarding their origin : *PTBP1* (red), *PTBP2* (green), *PTBP3* (blue) or protostoma (grey). Note again the position of the UUG-Leu codon in the codon dendrogram, as the sole GC-ending codon clustering (in purple) with all other AT-ending codons (in orange)

260 We have then analysed the correspondence between nucleotide-based and amino acid-based pairwise distances to eval-
 261 uate the impact of CUPrefs on the obtained phylogeny. We observe a good correlation between both reconstructions
 262 for all paralogs, except for mammalian *PTBP2*s, which display extremely low divergence at the amino acid level (see
 263 Figure 5 for values in mammalian paralogs, Supplementary Figure S8 for non-mammalian paralogs, and Supplemen-
 264 tary Table S7 for the correlation between nucleotide-based and amino acid-based pairwise distances). For mammalian
 265 *PTBP1*s, the plot allows to clearly differentiate a cloud with the values corresponding to the monotremes+marsupial
 266 mammals, split apart from placental mammals in terms of both amino acid and nucleotide distances. This distribution
 267 matches well the fact that sequences from monotremes and marsupials cluster separately from placental mammals in

Table 4: Global linear regression model and post-hoc Tukey's honest significant differences (HSD) test, the explained variable being the COUSIN value of the each *PTBP* gene compared with the average of the corresponding genome, and the explanatory levels paralog (*PTBP1-3*), taxonomy (*i.e.* mammalian or non-mammalian) and their interactions. Within each level, strata labelled with the same letter are not different from one another. Overall goodness of the fit: Adj Rsquare=0.82; F ratio=36.84; Prob > F: <0.0001. Individual effects for the levels: i) paralog: F ratio=40.72; Prob > F: <0.0001; ii) taxonomy: F ratio=10.87; Prob > F: =0.0021; iii) interaction paralog*taxonomy: F ratio=28.11; Prob > F: <0.0001.

Level	Least Sq. Mean (COUSIN)	Standard error	Tukey's HSD group
Paralog			
<i>PTBP1</i>	1.45	0.11	A
<i>PTBP3</i>	0.29	0.11	B
<i>PTBP2</i>	0.19	0.11	B
Taxonomy			
mammalian	0.44	0.080	A
non-mammalian	0.85	0.098	B
Paralog*Taxonomy			
<i>PTBP1</i> , mammalian	1.90	0.14	A
<i>PTBP1</i> , non-mammalian	0.99	0.17	B
<i>PTBP2</i> , non-mammalian	0.81	0.17	B
<i>PTBP3</i> , non-mammalian	0.75	0.17	B
<i>PTBP3</i> , mammalian	-0.16	0.14	C
<i>PTBP2</i> , mammalian	-0.43	0.14	C

Table 5: Comparison between species tree and the nucleotide based maximum likelihood tree for each *PTBP* paralog. The K-tree score compares topological and pairwise distances between trees after re-scaling overall tree length, with higher values corresponding to more divergent trees. The Robinson-Foulds score compares only topological distances between trees, the values shown correspond to the number of tree partitions that are not shared between two trees, so that higher values correspond to more divergent trees.

Reference tree	Comparison tree	K-tree score	Robinson-Foulds score
Nucleotide tree VS species tree			
<i>PTBP1</i>	Species tree	0.759	42
<i>PTBP2</i>	Species tree	0.762	24
<i>PTBP3</i>	Species tree	1.700	28
Nucleotide tree VS Amino acid tree			
<i>PTBP1</i> -AA	<i>PTBP1</i> -NT	0.149	78
<i>PTBP2</i> -AA	<i>PTBP2</i> -NT	0.129	110
<i>PTBP3</i> -AA	<i>PTBP3</i> -NT	0.380	40

Evolution of codon usage preferences in paralogous genes

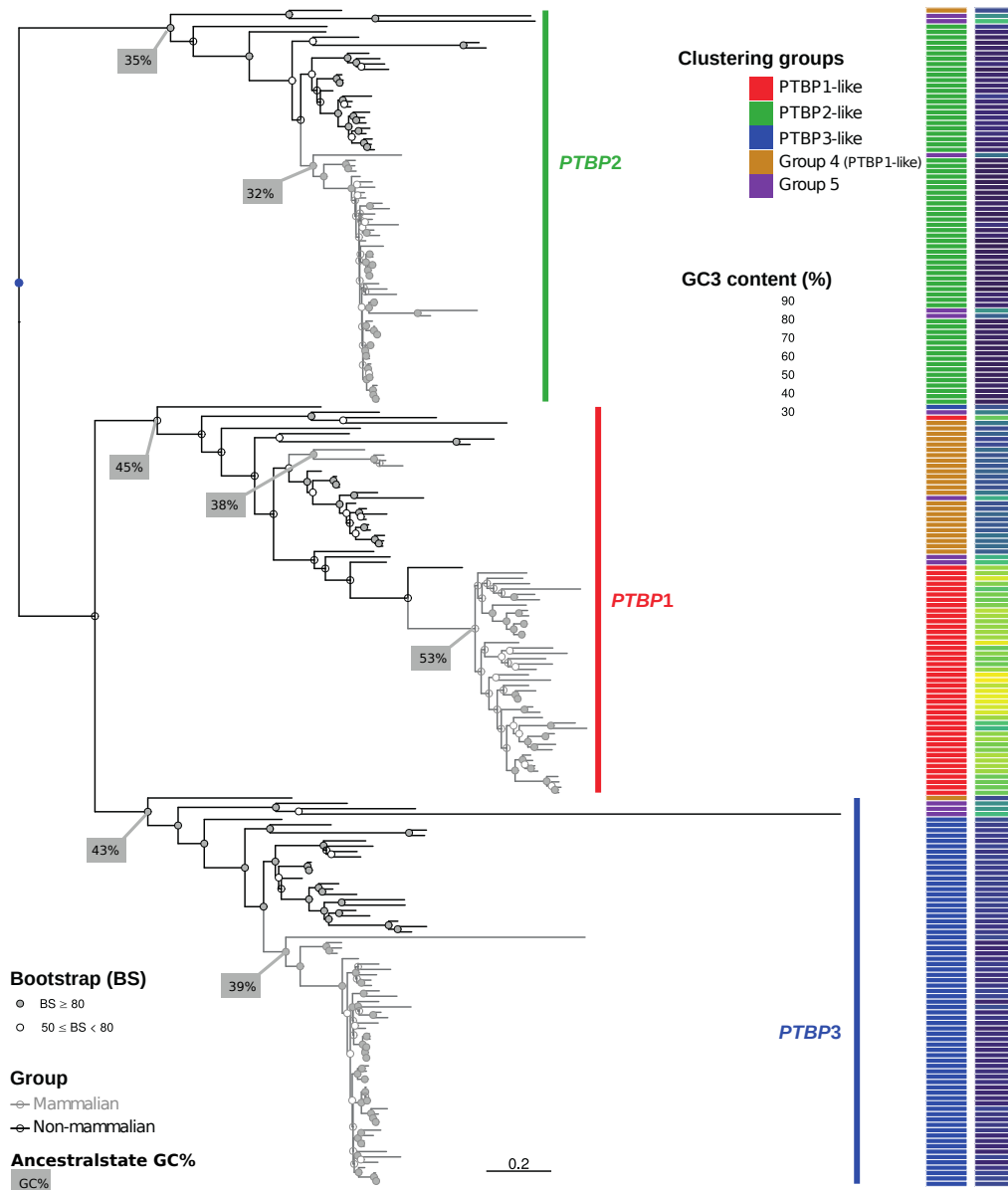


Figure 4: **Maximum-likelihood nucleic acid phylogeny of PTBP genes.** The phylogram depicts *PTBP2*s (green side bar), *PTBP1*s (red side bar) and *PTBP3*s (blue side bar) clades. The outgroup genes from protostomata are not shown to focus on the scale for vertebrate *PTBP*s, but their placement on the tree and the polarity they provide for vertebrate *PTBP*s is given by the blue dot. Gray branches indicate mammalian *PTBP*s, while black branches indicate non-mammalian species. Note the lack of monophyly for mammals for *PTBP1*s, with monotremes and marsupial lineages being paraphyletic to placental mammals. Filled dots on nodes indicate bootstrap values above 80, and empty dots indicate lower support values. Side bar on the left identifies the classification of each gene into the five groups identified by the hierarchical clusters, with the colour code in the inset. Side bar on the right displays GC3 content of the corresponding genes, with the gradient for the colour code ranging from 0 (blue) to 100% (yellow). The GC content of the main ancestral nodes is indicated in grey boxes.

Evolution of codon usage preferences in paralogous genes

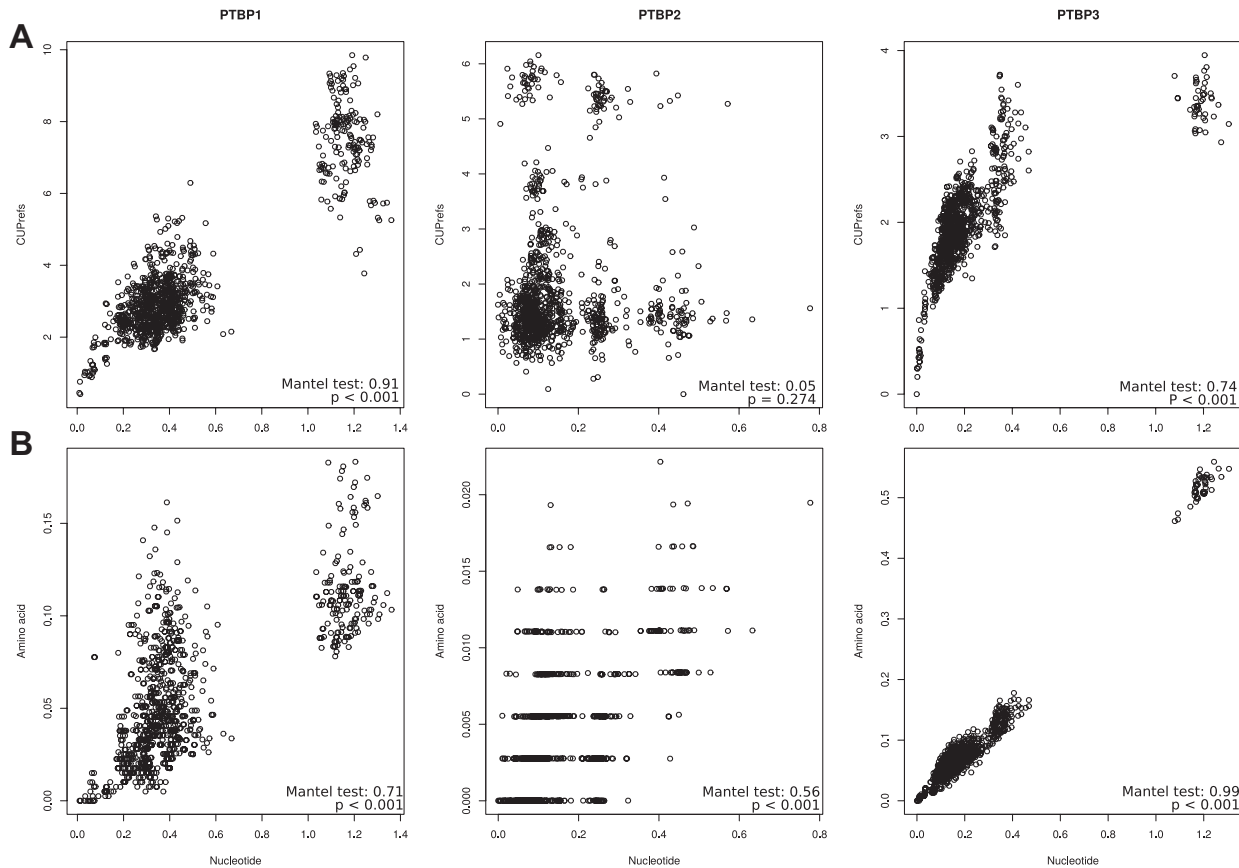


Figure 5: Nucleotide-based pairwise distances in the x-axis against A) CUPrefs-based and B) amino acid-based pairwise distances in the y-axis for the different mammalian *PTBP* orthologs. The results for a Mantel test assessing the correlation between the corresponding matrices are shown in each inset.

268 the *PTBP1* phylogeny (see grey branches being paraphyletic for *PTBP1* in Figure 4). The same holds true for the
 269 platypus *PTBP3*, extremely divergent from the rest of the mammalian orthologs. The precise substitution patterns are
 270 analysed in detail below. The histograms describing the accumulation of synonymous and non-synonymous substitu-
 271 tions confirm that mammalian *PTBP1*s have accumulated the largest number of synonymous substitutions compared
 272 to non-mammalian *PTBP1*s and to other orthologs (Supplementary Figure S9).

273 We have finally analysed the connection between nucleotide-based evolutionary distances within *PTBP* paralogs and
 274 CUPrefs-based distances (Figure 5 for mammalian paralogs and Supplementary Figure S8 for non-mammalian par-
 275 alogs). A trend showing increased differences in CUPrefs as evolutionary distances increase is evident only for
 276 *PTBP1*s and *PTBP3*s in mammals. For mammalian *PTBP1*s the plot clearly differentiates a cloud with the values
 277 corresponding to monotremes and marsupials splitting apart from placental mammals in terms of both evolutionary
 278 distance and CUPrefs. For mammalian *PTBP2*s the plot captures the divergent CUPrefs of the platypus and of the bats
 279 *M. natalensis* and *Desmodus rotundus*, while for non-mammalian *PTBP2*s the divergent CUPrefs of the rainbow trout
 280 (*O. mykiss*) are obvious. Finally, for mammalian *PTBP3*s the large nucleotide divergence of the platypus paralog is
 281 evident. Importantly, all these instances of divergent behaviour (except for the platypus *PTBP3*) are consistent with the

282 deviations described above from the expected composition by the mathematical modelling of the ortholog nucleotide
283 composition (Table 2).

284 ***Mammalian PTBP1s accumulate GC-enriching synonymous substitutions***

285 We have shown that *PTBP1* genes are GC-richer and specifically GC3-richer than the *PTBP2* and *PTBP3* paralogs
286 in the same genome, and that this enrichment is of a larger magnitude in placental *PTBP1s*. We have thus assessed
287 whether a directional substitutional pattern underlies this enrichment, especially regarding synonymous substitutions.
288 For this we have inferred the ancestral sequences of the respective most recent common ancestors of each *PTBP* para-
289 log, recapitulated synonymous and non-synonymous substitutions between each *PTBP* individual and their ancestors,
290 and constructed the corresponding substitution matrices (Table S11). The two first axes of a principal component
291 analysis using these substitution matrices capture, with a similar share, 66.95% of the variance between individuals
292 (Figure 6). The first axis of the PCA separates synonymous from non-synonymous substitutions. Intriguingly though,
293 while T<->C transitions are associated with synonymous substitutions, as expected, G<->A transitions are instead
294 associated with non-synonymous substitutions. The second axis separates substitutions by their effect on nucleotide
295 composition: GC-stabilizing/enriching on one direction, AT-stabilizing/enriching on the other one. Strikingly, the sub-
296 substitutional spectrum of mammalian *PTBP1s* sharply differs from the rest of the paralogs. Substitutions in mammalian
297 *PTBP1* towards GC-enriching changes, in both synonymous and non-synonymous compartments, are the main drivers
298 of the second PCA axis. In contrast, synonymous substitutions in *PTBP3* as well as all substitutions in *PTBP2* tend
299 to be AT-enriching. Finally, the substitution trends for *PTBP1* in mammals are radically different from those in non-
300 mammals, while for *PTBP2* and *PTBP3s* the substitution patterns are similar in mammals and non-mammals for each
301 of the compartments synonymous and non-synonymous.

302 **5 Discussion**

303 The non equal use of synonymous codons has puzzled biologists since it was first described. It has given rise to fruit-
304 ful (and unfruitful) controversies between defenders of *all-is-neutralism* and defenders of *all-is-selectionism*, and has
305 launched further the quest for additional molecular signaling beyond codons themselves (Callens et al., 2021). The
306 main questions around CUPrefs are twofold. On the one hand, their origin: to what extent they are the result of
307 fine interplay between mutation and selection processes. On the other hand, their functional implications: whether
308 and how particular CUPrefs can be linked to specific gene expression regulation processes, broadly understood as
309 downstream effects that modify the kinetics and dynamics of DNA transcription, mRNA maturation and stability,
310 mRNA translation, and/or protein folding and stability. In the present work we have built on the experimental results
311 of Robinson and coworkers, which communicated the differential expression of the *PTBP* human gene paralogs as
312 a function of their CUPrefs (Robinson et al., 2008). From this particular example, we have aimed at exploring the
313 nature of the connection between paralogous gene evolution and CUPrefs. Our results show that the three *PTBP* par-
314 alogous genes, which show divergent expression patterns in humans, also have divergent nucleotide composition and
315 CUPrefs not just in humans but in most vertebrate species. We elaborate here on Robinson and coworkers experimen-
316 tal findings and propose here that this evolutionary pattern could be compatible with a phenomenon of phenotypic
317 evolution by sub-functionalisation (in this case specialisation in tissue-specific expression levels), linked to genotypic
318 evolution by association to specific CUPrefs patterns. Such conclusions invite to pursue Robinson and coworkers' ef-

Evolution of codon usage preferences in paralogous genes

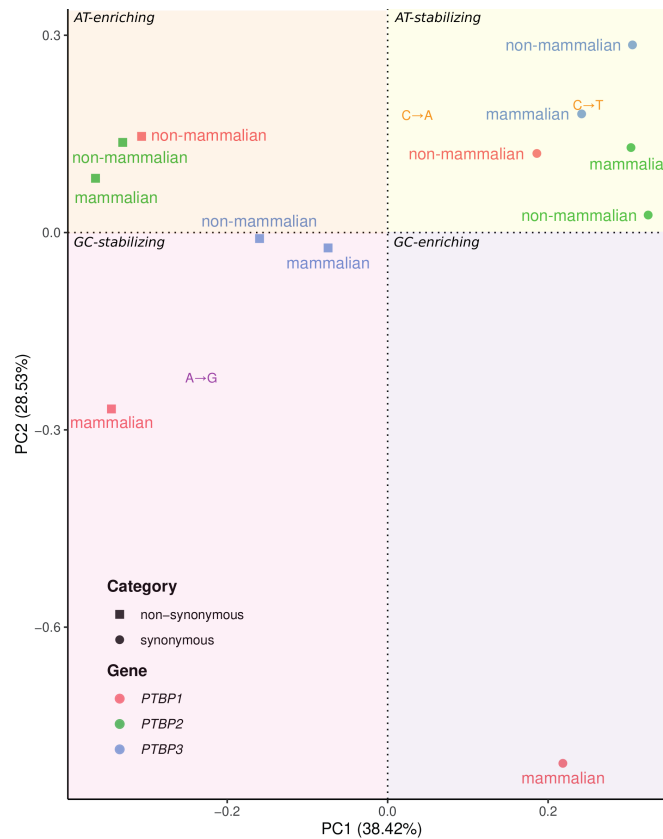


Figure 6: **Spectra of synonymous and non-synonymous substitutions for *PTBPs*.** This principal component analysis (PCA) has been built using the observed nucleotide synonymous and non-synonymous substitution matrices for each *PTBP* paralog, inferred after phylogenetic inference and comparison of extant and ancestral sequences. The variables in this PCA are the types of substitution (*e.g.* A->G), identified by a colour code as GC-enriching/stabilizing substitutions (purple and pink areas) or AT-enriching/stabilizing substitutions (orange and yellow areas). To facilitate the interpretation of the graph, all variables have been masked, except those that do not follow these global, which have been plotted according to their eigenvalues (*i.e.* A->G, C->A and C->T) (all variables are shown unmasked in Supplementary Figure S15). Individuals in this PCA are the substitution categories in *PTBP* genes, stratified by their nature (synonymous or non-synonymous), by orthology (colour code for the different *PTBPs* is given in the inset) and by their taxonomy (mammals, or non-mammals).

319 forts by comparing *PTBPs* CUPrefs-modulated expression among numerous Vertebrate cell lines, especially between
 320 mammals and non-mammals ones. Consistent with studies on other paralog families (Munk et al., 2022), our
 321 results suggest, more generally, that a detailed analysis of differential CUPrefs in paralogs may help understand the
 322 divergent/convergent mutation-selection pressures that could underlie their functional differences.

323 We have reconstructed the phylogenetic relationships and analysed the evolution and diversity of CUPrefs among
 324 *PTBP* paralogs within 74 vertebrate species. The phylogenetic reconstruction shows that the genome of ancestral
 325 vertebrates already contained the three extant *PTBP* paralogs. This is consistent with the ortholog and paralog identi-
 326 fication in the databases ENSEMBL and ORTHOMAM (Yates et al., 2020; Scornavacca et al., 2019; Pina et al., 2018).

327 Although our results suggest that *PTBP1* and *PTBP3* are sister lineages, the distant relationship between the vertebrate
328 genes and the protostome outgroup precludes the inference of a clear polarity between vertebrate *PTBPs*. We identify
329 no instance of basal replacement between paralogs which may have appeared, for instance, as the replacement of an
330 AT-rich paralog by a GC-rich one, leading to a loss of the AT-rich paralog and a duplication of the GC-rich one. Instead,
331 the basal evolutionary histories of the different *PTBPs* comply well with those of the corresponding species. The most
332 blatant mismatch between gene and species trees is the polyphyly of mammalian *PTBP1*s: monotremes and marsupials
333 constitute a monophyletic clade, separate from placental mammals and not basal to them. Further, multiple findings in
334 our results show sharp, contrasting patterns between *PTBP1* and the *PTBP2-3* paralogs: i) the excess of accumulation
335 of synonymous substitutions in mammalian *PTBP1*s for a similar total number of changes (Supplementary Figure S9
336 and Table S11); ii) the larger differences in CUPrefs between genes with a similar total number of nucleotide changes
337 in the case of *PTBP1*s in mammals (Figure 5 A); iii) the explicitly different spectrum of synonymous substitutions in
338 *PTBP1*s, enriched in A->C, T->G and T->C changes (Figure 6); iv) the sharp difference of CUPrefs between *PTBP1*s
339 and *PTBP2-3*s; and v) the clustering of *PTBP1* genes in monotremes and marsupials together with *PTBP1* genes in
340 non-mammals according to their CUPrefs (Figure 3 A). Overall, the particular nucleotide composition and the associ-
341 ated CUPrefs in mammalian *PTBP1* genes are most likely associated to specific local substitution biases as shown by
342 the strong correlation between coding and non-coding GC content in *PTBP1* orthologs, while CUPrefs in *PTBP2-3*s
343 cannot be explained alone by such local substitution biases (Figure 2; Table 3).

344 While GC3-rich nucleotide composition and CUPrefs of mammalian *PTBP1*s are dominated by local substitution
345 biases, this is not the case for mammalian *PTBP2*, overall AT3-richer and without any clear correlation between
346 coding and non-coding GC content among studied species (Figure 2; 3). As mentioned above, a note of caution should
347 be raised here, as the variable range for GC composition among *PTBP1*s is larger than for *PTBP2-3*s, so that co-
348 variation analyses may have less power for the latter paralogs. In vertebrates, nucleotide composition varies strongly
349 along chromosomes, so that long chromatin stretches, historically named "isochores", appear enriched in GC or in
350 AT nucleotides and present particular physico-chemical profiles (Caspersson et al., 1968). Local mutational biases
351 and GC-biased gene conversion mechanism may underlie such heterogeneity, predominantly shaping local nucleotide
352 composition in numerous Vertebrates genomes, so that the physical location of a gene along the chromosome largely
353 explains its CUPrefs (Holmquist, 1989). In agreement with these hypotheses for local mutational biases, variation in
354 GC3 composition of *PTBP1*s is almost totally ($R^2=0.97$) explained by the variation in local GC composition (Figure
355 2; Table 3), suggesting that a similar substitution bias has shaped the GC-rich composition of the flanking, intronic and
356 coding regions of *PTBP1*s. The same trend, albeit to a lesser degree holds also true for *PTBP2*s ($R^2=0.45$). GC-biased
357 gene conversion is often invoked as a powerful mechanism underlying such local GC-enrichment processes, leading
358 to the systematic replacement of the alleles with the lowest GC composition by a GC-richer homolog (Marais, 2003).
359 It has been proposed that gene expression during meiosis (evaluated as mRNA detection) correlates with a decreased
360 probability of GC-biased gene conversion during meiotic recombination (Pouyet et al., 2017). Expression of *PTBP1*
361 in human cells is documented during meiosis in the oocyte germinal line and expression of the AT-rich *PTBP2* has
362 been observed during spermatogenic meiosis (Zagore et al., 2015; Hannigan et al., 2017). Expression during meiosis
363 might thus have hindered GC-biased gene conversion for *PTBP1-2*s, provided that this expression pattern observed
364 in humans was displayed also by the mammalian ancestor and that it is shared between mammalian species. With
365 these assumptions, and thus, with caution, the GC-richness of *PTBP1* cannot be accounted for by GC-biased gene

366 conversion, while the low GC content of *PTBP2* could be explained by an accumulation of GC->AT and AT->AT
367 substitutions. All this notwithstanding, our results show that GC3 enrichment in mammalian *PTBP1* and the concurrent
368 trend for enriched use of common codons are associated mostly with placental mammals, and that non-placental
369 mammals display divergent composition and differ from the model expectations. This synapomorphy of a sudden
370 change in nucleotide composition is strongly compatible with a GC-biased gene conversion event in the placental
371 ancestor that may have led to fixation of the ancestral version of the extant GC-rich *PTBP1*. Regarding *PTBP3*, the
372 low GC-content together with the low correlation with either coding nor non-coding local GC-content could indicate
373 that other mechanisms may shape the observed CUPrefs for this paralog.

374 In mammals, global GC-enriching genomic biases strongly impact CUPrefs, so that the most used codons in average
375 tend to be GC-richer (Hershberg and Petrov, 2009). For this reason, mammalian GC3-rich *PTBP1*s match better the
376 average genomic CUPrefs than AT3-richer *PTBP2* and *PTBP3*, which display CUPrefs in the opposite direction to
377 the average of the genome. In the case of humans, *PTBP1* presents a COUSIN value of 1.75, consistent with a
378 substantial enrichment in preferentially-used codons, while on the contrary, the COUSIN values of -0.48 for *PTBP2*
379 and of -0.23 for *PTBP3* point towards a strong enrichment in rare codons (Supplementary Table S4). The poor match
380 between human *PTBP2* CUPrefs and the human average CUPrefs could result in low expression of these genes in
381 different human and murine cell lines, otherwise capable of expressing *PTBP1* at high levels and of expressing *PTBP3*
382 at a lesser degree (Robinson et al., 2008). The barrier to *PTBP2* expression seems to be the translation process, as
383 *PTBP2* codon-recoding towards GC3-richer codons results in strong protein production in the same cellular context,
384 without significant changes in the corresponding mRNA levels (Robinson et al., 2008). Such codon recoding strategy
385 towards preferred codons has become a standard practice for gene expression engineering that provides with very good
386 expression results, despite our lack of understanding about the whole impact of local and global gene composition,
387 nucleotide CUPrefs, and mRNA structure on gene expression (Brule and Grayhack, 2017).

388 The poor expression ability of *PTBP2* in human cells, the increase in protein production by the introduction of com-
389 mon codons, along with substitution biases failing to explain entirely *PTBP2* nucleotide composition and CUPrefs,
390 raise the question of the adaptive value of poor CUPrefs in this paralog. Specific tissue-dependent or cell-cycle de-
391 pendent gene expression regulation patterns have been invoked to explain the codon usage-limited gene expression
392 for certain human genes, such as *TLR7* or *KRAS* (Newman et al., 2016; Lampson et al., 2013; Fu et al., 2018). The
393 expression levels of the three *PTBP* paralogs are tissue-dependent in humans (Supplementary Figure S1), and through
394 mammals (Keppetipola et al., 2012; Wagner and Garcia-Blanco, 2002; Spellman et al., 2007). In the case of the du-
395 plicated genes, subfunctionalisation through specialisation in spatio-temporal gene expression has been proposed as
396 the main evolutionary force driving conservation of paralogous genes (Ferris and Whitt, 1979). Such differential
397 gene expression regulation in paralogs has actually been documented for a number of genes at very different taxo-
398 nomic levels (Donizetti et al., 2009; Guschanski et al., 2017; Freilich et al., 2006). Specialised expression patterns in
399 time and space can result in antagonistic presence/absence of the paralogous proteins (Adams et al., 2003). This is
400 precisely the case of *PTBP1* and *PTBP2* during human central nervous system development: in non-neuronal cells,
401 *PTBP1* represses *PTBP2* expression by the skip of the exon 10 during *PTBP2* mRNA maturation, while during neu-
402 ron development, the micro RNA miR124 down-regulates *PTBP1* expression, which in turn leads to up-regulation of
403 *PTBP2* (Keppetipola et al., 2012; Makeyev et al., 2007). Regarding non-human species, the available data about tissue-
404 dependent and/or ontogeny-dependent differential expression at the transcription level (Abugessaisa et al., 2021) are

405 largely concordant with the human data for *PTBP*, showing a tissue-wide transcription of *PTBP1*, a more restricted
406 one for *PTBP3* together with an enrichment of *PTBP2* transcription in the central nervous system, as exemplified
407 in the mouse (Barbosa-Morais et al., 2012), in the rat (Yu et al., 2014), in the cow (Merkin et al., 2012), in the gray
408 short-tailed opossum (Brawand et al., 2011), or in the chicken (Barbosa-Morais et al., 2012). Finally, despite the high
409 level of amino acid similarity between both proteins, *PTBP1* and *PTBP2* seem to perform complementary activities
410 in the cell and to display different substrate specificity, so that they are not directly inter-exchangeable by exogenous
411 manipulation of gene expression patterns (Vuong et al., 2016).

412 In addition to local genomic context analyses, we explored *PTBP* chromosomal location and local synteny (Figure
413 7). The results show that, while it is clear that the position of human *PTBP1* is telomeric and thus in one of the
414 GC-richer region of human chromosome 9, most *PTBPs* do not map to the telomeres. Therefore, while the specific
415 location of human *PTBP1* may have influenced its CUPrefs, it is unclear whether the chromosomal location of *PTBPs*
416 have an impact on observed nucleotide composition. Synteny of *PTBPs* genes seems to be conserved, with some
417 exceptions: most mammalian *PTBP1s* have a conserved synteny block that differs from non-mammalian species, with
418 the exception of *D. rerio*. For *PTBP2* and *PTBP3* synteny seems conserved between mammalian and non-mammalian
419 species again with the exception of *D. rerio*, lacking the *SUSD1* gene between *PTBP3* and *UGCG*. Such results could
420 indicate that vertebrate radiation has been followed up by a change of *PTBP1* genomic context, with a swapping in
421 flanking genes in mammalian branches. These results could be related to the observed *PTBP1* differential GC-content
422 between mammalian and non-mammalian species.

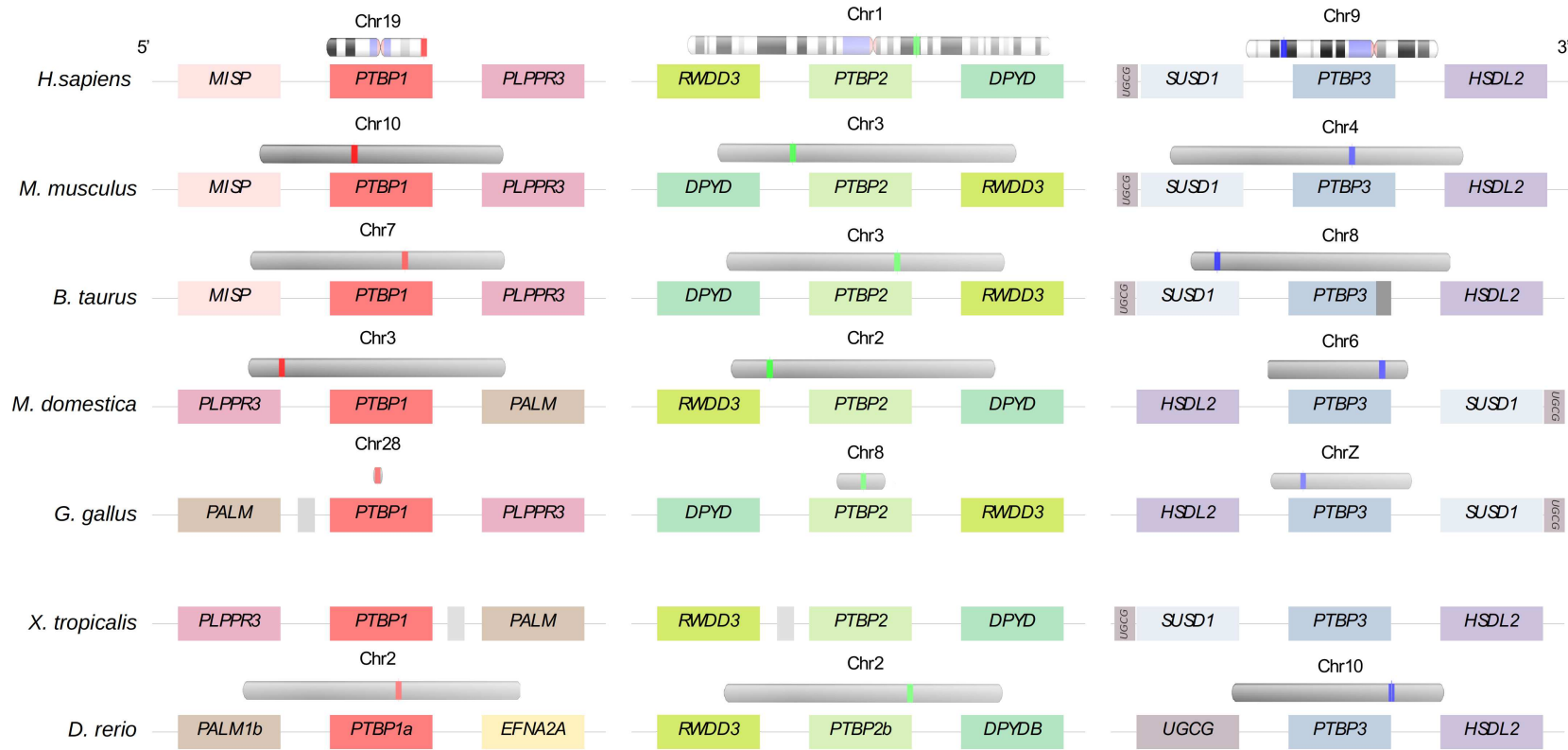


Figure 7: Placement on the chromosomes and genomic context of the three *PTBP* paralogs in a subset of the studied species.

423 In a different subject, we want to drive the attention of the reader towards the puzzling trend of the UUG-Leu codon in
424 our CUPrefs analyses. This UUG codon is the only GC-ending codon systematically clustering with AT-ending codons
425 in all our analyses, and not showing the expected symmetrical behaviour with respect to UUA-Leu (see Figure 3). Such
426 behaviour for UUG has been depicted, but not discussed, in other analyses of CUPrefs in mammalian genes (see figure
427 7 in Laurin-Lemay et al. (2018)), in coronavirus genomes (Daron and Bravo, 2021), in plants (Clément et al., 2017)
428 as well as for AGG-Arg and GGG-Gly in a global study of codon usages across the tree of life (see figure 1 in
429 (Novoa et al., 2019)). The reasons underlying the clustering of UUG with AT-ending codons are unclear. A first line
430 of thought could be functional: the UUG-Leu codon is particular because it can serve as alternative starting point for
431 translation (Peabody, 1989). However, other codons such as ACG or GUG act more efficient than UUG as alternative
432 translation initiation, and do not display any noticeable deviation in our results (Ivanov et al., 2011). A second line of
433 thought could be related to the tRNA repertoire, but both UUG and UUA are decoded by similar numbers of dedicated
434 tRNAs in the vast majority of genomes (*e.g.* respectively six and seven tRNA genes in humans (Palidwor et al., 2010)).
435 Finally, another line of thought suggests that UUG and AGG could be disfavoured if substitution pressure towards GC
436 is very high, despite being GC-ending codons (Palidwor et al., 2010). Indeed, the series of synonymous transitions
437 UUA->UUG->CUG for Leucine and the substitution chain AGA->AGG->CGG for Arginine are expected to lead to
438 a depletion of UUG and of AGG codons when increasing GC content. Both UUG and ACG codons would this way
439 display a non-monotonic response to GC-substitution biases (Palidwor et al., 2010). In our data-set, however, AGG
440 maps with the rest of GC-ending codons, symmetrically opposed to AGA as expected, and strongly contributing to
441 the second PCA axis. Thus, only UUG displays frequency patterns similar to those of AT-ending codons. We humbly
442 admit that we do not find a satisfactory explanation for this behaviour and invite researchers in the field to generate
443 alternative explanatory hypotheses.

444 We have presented here an evolutionary analysis of the *PTBP* paralogs family as a showcase of CUPrefs evolution upon
445 gene duplication. Our results show that differential nucleotide composition and CUPrefs in *PTBP*ss have evolved in
446 parallel with differential gene expression regulation patterns. In the case of *PTBP1*, the most tissue-wise expressed of
447 the paralogs, we have potentially identified compositional and substitution biases as the driving force leading to strong
448 enrichment in GC-ending codons. In contrast, for *PTBP2* the enrichment in AT-ending codons is rather compatible
449 with selective forces related to specific spatio-temporal gene expression pattern, antagonistic to those of *PTBP1*. Our
450 results suggest that the systematic study of composition, genomic location and expression patterns of paralogous genes
451 can contribute to understanding the complex mutation-selection interplay shaping CUPrefs in multicellular organisms.

452 **6 Acknowledgments**

453 J.B. was the recipient of a PhD fellowship from the French Ministry of Education and Research. This study was
454 supported by the European Union's Horizon 2020 research and innovation program under the grant agreement
455 CODOVIREVOL (ERC-2014-CoG-647916) to I.G.B. The authors acknowledge the CNRS and the IRD for additional
456 intramural support. The computational results presented have been achieved in part using the IRD Bioinformatic
457 Cluster itrop.

458 **7 Data Availability Statement**

459 All data required to reproduce our findings is available on zenodo (<https://doi.org/10.5281/zenodo.5789766>), or pro-
460 vided in the tables in the main text and in the Supplementary Material section.

461 **8 Disclosure**

462 I.G.B. is a PCI recommender.

463 **References**

- 464 Abugessaisa I, Ramilowski JA, Lizio M, Severin J, Hasegawa A, Harshbarger J, Kondo A, Noguchi S, Yip CW, Ooi J,
465 Tagami M, Hori F, Agrawal S, Hon C, Cardon M, Ikeda S, Ono H, Bono H, Kato M, Hashimoto K, Bonetti A, Kato
466 M, Kobayashi N, Shin J, de Hoon M, Hayashizaki Y, Carninci P, Kawaji H, Kasukawa T. 2021, January. FANTOM
467 enters 20th year: expansion of transcriptomic atlases and functional annotation of non-coding RNAs. *Nucleic Acids*
468 *Research*. 49(D1):D892–D898.
- 469 Adams KL, Cronn R, Percifield R, Wendel JF. 2003, April. Genes duplicated by polyploidy show unequal contributions
470 to the transcriptome and organ-specific reciprocal silencing. *Proceedings of the National Academy of Sciences of*
471 *the United States of America*. 100(8):4649–4654.
- 472 Apostolou-Karampelis K, Nikolaou C, Almirantis Y. 2016, August. A novel skew analysis reveals substitution asym-
473 metries linked to genetic code GC-biases and PolIII a-subunit isoforms. *DNA research: an international journal for*
474 *rapid publication of reports on genes and genomes*. 23(4):353–363.
- 475 Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Guerousov S, Lee LJ, Slobodeniuc V, Kutter C, Watt S, Colak
476 R, Kim T, Misquitta-Ali CM, Wilson MD, Kim PM, Odom DT, Frey BJ, Blencowe BJ. 2012, December. The
477 evolutionary landscape of alternative splicing in vertebrate species. *Science (New York, N.Y.)*. 338(6114):1587–
478 1593.
- 479 Bourret J, Alizon S, Bravo IG. 2019, December. COUSIN (COdon Usage Similarity INdex): A Normalized Measure
480 of Codon Usage Preferences. *Genome Biology and Evolution*. 11(12):3523–3528. Publisher: Oxford Academic.
- 481 Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A, Aximu-Petri A, Kircher
482 M, Albert FW, Zeller U, Khaitovich P, Grützner F, Bergmann S, Nielsen R, Pääbo S, Kaessmann H. 2011, October.
483 The evolution of gene expression levels in mammalian organs. *Nature*. 478(7369):343–348.
- 484 Brule CE, Grayhack EJ. 2017. Synonymous Codons: Choose Wisely for Expression. *Trends in genetics: TIG*.
485 33(4):283–297.
- 486 Bulmer M. 1991, November. The selection-mutation-drift theory of synonymous codon usage. *Genetics*. 129(3):897–
487 907.
- 488 Caliskan N, Peske F, Rodnina MV. 2015, May. Changed in translation: mRNA recoding by 1 programmed ribosomal
489 frameshifting. *Trends in Biochemical Sciences*. 40(5):265–274.
- 490 Callens M, Pradier L, Finnegan M, Rose C, Bedhomme S. 2021. Read between the lines: Diversity of nontranslational
491 selection pressures on local codon usage. *Genome Biology and Evolution*. 13.

- 492 Carbone A, Zinovyev A, Képès F. 2003, November. Codon adaptation index as a measure of dominating codon bias.
493 *Bioinformatics* (Oxford, England). 19(16):2005–2015.
- 494 Caspersson T, Farber S, Foley GE, Kudynowski J, Modest EJ, Simonsson E, Wagh U, Zech L. 1968, January. Chemical
495 differentiation along metaphase chromosomes. *Experimental Cell Research*. 49(1):219–222.
- 496 Castresana J. 2000, April. Selection of conserved blocks from multiple alignments for their use in phylogenetic
497 analysis. *Molecular Biology and Evolution*. 17(4):540–552.
- 498 Chamary JV, Parmley JL, Hurst LD. 2006, February. Hearing silence: non-neutral evolution at synonymous sites in
499 mammals. *Nature Reviews. Genetics*. 7(2):98–108.
- 500 Clark JM. 1988, October. Novel non-templated nucleotide addition reactions catalyzed by procaryotic and eucaryotic
501 DNA polymerases. *Nucleic Acids Research*. 16(20):9677–9686.
- 502 Clément Y, Sarah G, Holtz Y, Homa F, Pointet S, Contreras S, Nabholz B, Sabot F, Sauné L, Ardisson M, Bacilieri
503 R, Besnard G, Berger C Angélique Cardi, De Bellis F, Fouet O, Jourda C, Khadari B, Lanaud C, Leroy T, Pot D,
504 Sauvage C, Scarcelli N, Tregear J, Vigouroux Y, Yahiaoui N, Ruiz M, Santoni S, Labouisse JP, Pham JL, David J,
505 Glémin S. 2017. Evolutionary forces affecting synonymous variations in plant genomes. *PLoS Genetics*. 13:1–28.
- 506 Copley SD. 2020, April. Evolution of new enzymes by gene duplication and divergence. *The FEBS journal*.
507 287(7):1262–1283.
- 508 Daron J, Bravo IG. 2021. Variability in codon usage in coronaviruses is mainly driven by mutational bias and selective
509 constraints on cpg dinucleotide. *Viruses*. 13:1800.
- 510 Donizetti A, Fiengo M, Minucci S, Aniello F. 2009, October. Duplicated zebrafish relaxin-3 gene shows a different
511 expression pattern from that of the co-orthologue gene. *Development, Growth & Differentiation*. 51(8):715–722.
- 512 Duret L. 2002, December. Evolution of synonymous codon usage in metazoans. *Current Opinion in Genetics &*
513 *Development*. 12(6):640–649.
- 514 Duret L, Mouchiroud D. 1999, April. Expression pattern and, surprisingly, gene length shape codon usage in
515 *Caenorhabditis*, *Drosophila*, and *Arabidopsis*. *Proceedings of the National Academy of Sciences*. 96(8):4482–4487.
516 Publisher: National Academy of Sciences Section: Biological Sciences.
- 517 Ferris SD, Whitt GS. 1979, April. Evolution of the differential regulation of duplicate genes after polyploidization.
518 *Journal of Molecular Evolution*. 12(4):267–317.
- 519 Freilich S, Massingham T, Blanc E, Goldovsky L, Thornton JM. 2006. Relating tissue specialization to the differenti-
520 ation of expression of singleton and duplicate mouse proteins. *Genome Biology*. 7(10):R89.
- 521 Fu J, Dang Y, Counter C, Liu Y. 2018. Codon usage regulates human KRAS expression at both transcriptional and
522 translational levels. *The Journal of Biological Chemistry*. 293(46):17929–17940.
- 523 Galtier N, Roux C, Rousselle M, Romiguier J, Figuet E, Glémin S, Bierne N, Duret L. 2018, May. Codon Usage
524 Bias in Animals: Disentangling the Effects of Natural Selection, Effective Population Size, and GC-Biased Gene
525 Conversion. *Molecular Biology and Evolution*. 35(5):1092–1103.
- 526 Grantham R, Gautier C, Gouy M, Mercier R, Pavé A. 1980, January. Codon catalog usage and the genome hypothesis.
527 *Nucleic Acids Research*. 8(1):r49–r62.

- 528 Guschanski K, Warnefors M, Kaessmann H. 2017. The evolution of duplicate gene expression in mammalian organs.
529 *Genome Research*. 27(9):1461–1474.
- 530 Hannigan MM, Zagore LL, Licatalosi DD. 2017, June. Ptbp2 controls an alternative splicing network required for cell
531 communication during spermatogenesis. *Cell reports*. 19(12):2598–2612.
- 532 Hershberg R, Petrov DA. 2009, July. General rules for optimal codon choice. *PLoS genetics*. 5(7):e1000556.
- 533 Holmquist GP. 1989, June. Evolution of chromosome bands: Molecular ecology of noncoding DNA. *Journal of*
534 *Molecular Evolution*. 28(6):469–486.
- 535 Ikemura T. 1981, September. Correlation between the abundance of Escherichia coli transfer RNAs and the occurrence
536 of the respective codons in its protein genes: a proposal for a synonymous codon choice that is optimal for the E.
537 coli translational system. *Journal of Molecular Biology*. 151(3):389–409.
- 538 Ivanov IP, Firth AE, Michel AM, Atkins JF, Baranov PV. 2011, May. Identification of evolutionarily conserved non-
539 AUG-initiated N-terminal extensions in human coding sequences. *Nucleic Acids Research*. 39(10):4220–4234.
- 540 Katoh K, Misawa K, Kuma Ki, Miyata T. 2002, July. MAFFT: a novel method for rapid multiple sequence alignment
541 based on fast Fourier transform. *Nucleic Acids Research*. 30(14):3059–3066.
- 542 Keppetipola N, Sharma S, Li Q, Black DL. 2012, August. Neuronal regulation of pre-mRNA splicing by polypyrim-
543 idine tract binding proteins, PTBPI and PTBP2. *Critical Reviews in Biochemistry and Molecular Biology*.
544 47(4):360–378.
- 545 Khorana HG, Büchi H, Ghosh H, Gupta N, Jacob TM, Kössel H, Morgan R, Narang SA, Ohtsuka E, Wells RD. 1966.
546 Polynucleotide synthesis and the genetic code. *Cold Spring Harbor Symposia on Quantitative Biology*. 31:39–49.
- 547 Koonin EV. 2005. Orthologs, Paralogs, and Evolutionary Genomics. *Annual Review of Genetics*. 39(1):309–338.
548 *_eprint: <https://doi.org/10.1146/annurev.genet.39.073003.114725>*.
- 549 Kumar S, Stecher G, Suleski M, Hedges SB. 2017. TimeTree: A Resource for Timelines, Timetrees, and Divergence
550 Times. *Molecular Biology and Evolution*. 34(7):1812–1819.
- 551 Lampson BL, Pershing NLK, Prinz JA, Lacsina JR, Marzluff WF, Nicchitta CV, MacAlpine DM, Counter CM. 2013,
552 January. Rare codons regulate KRas oncogenesis. *Current biology: CB*. 23(1):70–75.
- 553 Laurin-Lemay S, Rodrigue N, Lartillot N, Philippe H. 2018. Conditional Approximate Bayesian Computation: A New
554 Approach for Across-Site Dependency in High-Dimensional Mutation-Selection Models. *Molecular Biology and*
555 *Evolution*. 35(11):2819–2834.
- 556 Lujan SA, Williams JS, Pursell ZF, Abdulovic-Cui AA, Clark AB, McElhinny SAN, Kunkel TA. 2012, October. Mis-
557 match Repair Balances Leading and Lagging Strand DNA Replication Fidelity. *PLOS Genetics*. 8(10):e1003016.
558 Publisher: Public Library of Science.
- 559 Makeyev EV, Zhang J, Carrasco MA, Maniatis T. 2007, August. The MicroRNA miR-124 Promotes Neuronal Differ-
560 entiation by Triggering Brain-Specific Alternative Pre-mRNA Splicing. *Molecular cell*. 27(3):435–448.
- 561 Marais G. 2003, June. Biased gene conversion: implications for genome and sex evolution. *Trends in Genetics*.
562 19(6):330–338. Publisher: Elsevier.

- 563 Merkin J, Russell C, Chen P, Burge CB. 2012, December. Evolutionary dynamics of gene and isoform regulation in
564 Mammalian tissues. *Science (New York, N.Y.)*. 338(6114):1593–1599.
- 565 Mordstein C, Savaisaar R, Young RS, Bazile J, Talmane L, Luft J, Liss M, Taylor MS, Hurst LD, Kudla G. 2020, April.
566 Codon Usage and Splicing Jointly Influence mRNA Localization. *Cell Systems*. 10(4):351–362.e8.
- 567 Munk M, Villalobo E, Villalobo A, Berchtold MW. 2022, November. Differential expression of the three independent
568 cam genes coding for an identical protein: Potential relevance of distinct mRNA stability by different codon usage.
569 *Cell Calcium*. 107.
- 570 NCBI Resource Coordinators. 2018. Database resources of the National Center for Biotechnology Information. *Nu-
571 cleic Acids Research*. 46(D1):D8–D13.
- 572 Newman ZR, Young JM, Ingolia NT, Barton GM. 2016, March. Differences in codon bias and GC content contribute
573 to the balanced expression of TLR7 and TLR9. *Proceedings of the National Academy of Sciences of the United
574 States of America*. 113(10):E1362–1371.
- 575 Nirenberg MW, Matthaei JH. 1961, October. The dependence of cell- free protein synthesis in *e. coli* upon naturally
576 occurring or synthetic polyribonucleotides. *Proceedings of the National Academy of Sciences of the United States
577 of America*. 47(10):1588–1602.
- 578 Novoa EM, Jungreis I, Jaillon O, Kellis M. 2019. Elucidation of Codon Usage Signatures across the Domains of Life.
579 *Molecular Biology and Evolution*. 36(10):2328–2339.
- 580 Novoa EM, Ribas de Pouplana L. 2012, November. Speeding with control: codon usage, tRNAs, and ribosomes.
581 *Trends in genetics: TIG*. 28(11):574–581.
- 582 Palidwor GA, Perkins TJ, Xia X. 2010, October. A general model of codon bias due to GC mutational bias. *PloS One*.
583 5(10):e13431.
- 584 Peabody DS. 1989, March. Translation initiation at non-AUG triplets in mammalian cells. *The Journal of Biological
585 Chemistry*. 264(9):5031–5035.
- 586 Percudani R, Pavesi A, Ottonello S. 1997, May. Transfer RNA gene redundancy and translational selection in *Saccha-
587 romyces cerevisiae* Edited by J. Karn. *Journal of Molecular Biology*. 268(2):322–330.
- 588 Pina J, Ontiveros RJ, Keppetipola N, Nikolaidis N. 2018, April. A Bioinformatics Approach to Discover the Evolu-
589 tionary Origin of the PTBP Splicing Regulators. *The FASEB Journal*. 32(1_supplement):802.16–802.16. Publisher:
590 Federation of American Societies for Experimental Biology.
- 591 Plotkin JB, Kudla G. 2011, January. Synonymous but not the same: the causes and consequences of codon bias. *Nature
592 Reviews Genetics*. 12(1):32–42.
- 593 Pouyet F, Mouchiroud D, Duret L, Sémon M. 2017. Recombination, meiotic expression and human codon usage.
594 *eLife*. 6.
- 595 Presnyak V, Alhusaini N, Chen YH, Martin S, Morris N, Kline N, Olson S, Weinberg D, Baker KE, Graveley BR,
596 Collier J. 2015, March. Codon optimality is a major determinant of mRNA stability. *Cell*. 160(6):1111–1124.
- 597 Reijns MAM, Kemp H, Ding J, Marion de Procé S, Jackson AP, Taylor MS. 2015, February. Lagging-strand replication
598 shapes the mutational landscape of the genome. *Nature*. 518(7540):502–506. Number: 7540 Publisher: Nature
599 Publishing Group.

- 600 Robinson DF, Foulds LR. 1981, February. Comparison of phylogenetic trees. *Mathematical Biosciences*. 53(1):131–
601 147.
- 602 Robinson F, Jackson RJ, Smith CWJ. 2008, March. Expression of Human nPTB Is Limited by Extreme Suboptimal
603 Codon Content. *PLOS ONE*. 3(3):e1801. Publisher: Public Library of Science.
- 604 Satapathy SS, Powdel BR, Buragohain AK, Ray SK. 2016, October. Discrepancy among the synonymous codons
605 with respect to their selection as optimal codon in bacteria. *DNA Research*. 23(5):441–449. Publisher: Oxford
606 Academic.
- 607 Scornavacca C, Belkhir K, Lopez J, Dernat R, Delsuc F, Douzery EJP, Ranwez V. 2019, April. OrthoMaM v10:
608 Scaling-Up Orthologous Coding Sequence and Exon Alignments with More than One Hundred Mammalian
609 Genomes. *Molecular Biology and Evolution*. 36(4):861–862. Publisher: Oxford Academic.
- 610 Sharp PM, Li WH. 1987. The codon Adaptation Index—a measure of directional synonymous codon usage bias, and
611 its potential applications. *Nucleic Acids Research*. 15(3):1281–1295.
- 612 Sonnhammer ELL, Koonin EV. 2002, December. Orthology, paralogy and proposed classification for paralog subtypes.
613 *Trends in genetics: TIG*. 18(12):619–620.
- 614 Soria-Carrasco V, Talavera G, Igea J, Castresana J. 2007, November. The K tree score: quantification of differences
615 in the relative branch length and topology of phylogenetic trees. *Bioinformatics (Oxford, England)*. 23(21):2954–
616 2956.
- 617 Spellman R, Llorian M, Smith CW. 2007, August. Crossregulation and functional redundancy between the splicing
618 regulator ptb and its paralogs nptb and rod1. *Molecular Cell*. 27:420–434.
- 619 Spencer PS, Barral JM. 2012, March. Genetic code redundancy and its influence on the encoded polypeptides. *Com-
620 putational and Structural Biotechnology Journal*. 1.
- 621 Stamatakis A. 2014, May. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.
622 *Bioinformatics (Oxford, England)*. 30(9):1312–1313.
- 623 Vuong JK, Lin CH, Zhang M, Chen L, Black DL, Zheng S. 2016. PTBP1 and PTBP2 Serve Both Specific and
624 Redundant Functions in Neuronal Pre-mRNA Splicing. *Cell Reports*. 17(10):2766–2775.
- 625 Wagner EJ, Garcia-Blanco MA. 2002, October. Rnai-mediated ptb depletion leads to enhanced exon definition. *Molec-
626 ular Cell*. 10:943–949.
- 627 Whittle CA, Extavour CG. 2016, September. Expression-Linked Patterns of Codon Usage, Amino Acid Frequency,
628 and Protein Length in the Basally Branching Arthropod Parasteatoda tepidariorum. *Genome Biology and Evolution*.
629 8(9):2722–2736. Publisher: Oxford Academic.
- 630 Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, Azov AG, Bennett
631 R, Bhai J, Billis K, Boddu S, Marugán JC, Cummins C, Davidson C, Dodiya K, Fatima R, Gall A, Giron CG, Gil
632 L, Grego T, Haggerty L, Haskell E, Hourlier T, Izuogu OG, Janacek SH, Juettemann T, Kay M, Lavidas I, Le T,
633 Lemos D, Martinez JG, Maurel T, McDowall M, McMahan A, Mohanan S, Moore B, Nuhn M, Oheh DN, Parker
634 A, Parton A, Patricio M, Sakthivel MP, Abdul Salam AI, Schmitt BM, Schuilenburg H, Sheppard D, Sycheva M,
635 Szuba M, Taylor K, Thormann A, Threadgold G, Vullo A, Walts B, Winterbottom A, Zadissa A, Chakiachvili M,
636 Flint B, Frankish A, Hunt SE, Iisley G, Kostadima M, Langridge N, Loveland JE, Martin FJ, Morales J, Mudge

- 637 JM, Muffato M, Perry E, Ruffier M, Trevanion SJ, Cunningham F, Howe KL, Zerbino DR, Flicek P. 2020, January.
638 Ensembl 2020. Nucleic Acids Research. 48(D1):D682–D688. Publisher: Oxford Academic.
- 639 Yu Y, Fuscoe JC, Zhao C, Guo C, Jia M, Qing T, Bannon DI, Lancashire L, Bao W, Du T, Luo H, Su Z, Jones
640 WD, Moland CL, Branham WS, Qian F, Ning B, Li Y, Hong H, Guo L, Mei N, Shi T, Wang KY, Wolfinger RD,
641 Nikolsky Y, Walker SJ, Duerksen-Hughes P, Mason CE, Tong W, Thierry-Mieg J, Thierry-Mieg D, Shi L, Wang C.
642 2014, February. A rat RNA-Seq transcriptomic BodyMap across 11 organs and 4 developmental stages. Nature
643 Communications. 5(1):3230. Number: 1 Publisher: Nature Publishing Group.
- 644 Zagore LL, Grabinski SE, Sweet TJ, Hannigan MM, Sramkoski RM, Li Q, Licatalosi DD. 2015, December. RNA
645 Binding Protein Ptp2 Is Essential for Male Germ Cell Development. Molecular and Cellular Biology. 35(23):4030–
646 4042.