



HAL
open science

Subfunctionalisation of paralogous genes and evolution of differential codon usage preferences: The showcase of polypyrimidine tract binding proteins

Jérôme Bourret, Fanni Borvetó, Ignacio G Bravo

► **To cite this version:**

Jérôme Bourret, Fanni Borvetó, Ignacio G Bravo. Subfunctionalisation of paralogous genes and evolution of differential codon usage preferences: The showcase of polypyrimidine tract binding proteins. *Journal of Evolutionary Biology*, 2023, 10.1111/jeb.14212 . hal-04196559

HAL Id: hal-04196559

<https://hal.science/hal-04196559>

Submitted on 5 Sep 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

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

1 **Acknowledgments** J.B. was the recipient of a PhD fellowship from the French Ministry of Education and Research.
2 This study was supported by the European Union's Horizon 2020 research and innovation program under the grant
3 agreement CODOVIREVOL (ERC-2014-CoG-647916) to I.G.B. The authors acknowledge the CNRS and the IRD for
4 additional intramural support. The computational results presented have been achieved in part using the IRD i-Trop
5 Plant & Health Bioinformatics Platform.
6 **Data Availability Statement** All data required to reproduce our analyses are available on zenodo
7 (<https://doi.org/10.5281/zenodo.5789766>), or provided in the tables in the main text and in the Supplementary
8 Material section.

*Corresponding author. email : fanni.borveto@uni-ulm.de; ignacio.bravo@cirs.fr

9 **1 Main manuscript**

10 **Subfunctionalisation of paralogous genes**
11 **and evolution of differential codon usage preferences:**
12 **the showcase of polypyrimidine tract binding proteins**

ABSTRACT

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

19 *PTBP1-3* show high identity at the amino acid level (up to 80%), but display strongly different
20 nucleotide composition, divergent CUPrefs and, in humans and in many other vertebrates, distinct
21 tissue-specific expression levels. Our phylogenetic inference results show that the duplication events
22 leading to the three extant *PTBP1-3* lineages predate the basal diversification within vertebrates, and
23 genomic context analysis illustrates that local synteny has been well preserved over time for the
24 three paralogs. We identify a distinct evolutionary pattern towards GC3-enriching substitutions in
25 *PTBP1*, concurrent with an enrichment in frequently used codons and with a tissue-wide expression.
26 In contrast, *PTBP2s* are enriched in AT-ending, rare codons, and display tissue-restricted expression.
27 As a result of this substitution trend, CUPrefs sharply differ between mammalian *PTBP1s* and the
28 rest of *PTBPs*. Genomic context analysis suggests that GC3-rich nucleotide composition in *PTBP1s*
29 is driven by local substitution processes, while the evidence in this direction is thinner for *PTBP2-*
30 *3*. An actual lack of co-variation between the observed GC composition of *PTBP2-3* and that of
31 the surrounding non-coding genomic environment would raise an interrogation on the origin of
32 CUPrefs, warranting further research on a putative tissue-specific translational selection. Finally,
33 we communicate an intriguing trend for the use of the UUG-Leu codon, which matches the trends
34 of AT-ending codons.

35 Our results are compatible with a scenario in which a combination of directional mutation–selection
36 processes would have differentially shaped CUPrefs of *PTBPs* in vertebrates: the observed GC-
37 enrichment of *PTBP1* in placental mammals may be linked to genomic location and to the strong
38 and broad tissue-expression, while AT-enrichment of *PTBP2* and *PTBP3* would be associated with
39 rare CUPrefs and thus, possibly to specialized spatio-temporal expression. Our interpretation is
40 coherent with a gene subfunctionalisation process by differential expression regulation associated to
41 the evolution of specific CUPrefs.

42 **Keywords** Codon usage bias, codon usage preferences, gene duplication, paralog, ortholog, evolution, nucleotide
43 composition, tissue, gene expression

44 **2 Significance Statement**

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

53 **3 Introduction**

54 During mRNA translation ribosomes assemble proteins by specific amino acid linear polymerisation guided by the
 55 successive reading of mRNA nucleotide triplets, called codons. Each time a codon is read, it is chemically compared
 56 to the set of available tRNAs' anticodons. Upon codon-anticodon match, the ribosome loads the tRNA and adds the
 57 associated amino acid to the nascent protein. The main 20 amino acids are encoded by 61 codons, so that multiple
 58 codons are associated with the same amino acid. These are named synonymous codons (Nirenberg and Matthaei,
 59 1961; Khorana et al., 1966). Codon Usage Preferences (CUPrefs) refer to the differential usage of synonymous
 60 codons between species, between genes, or between genomic regions in the same genome (Grantham et al., 1980;
 61 Carbone et al., 2003). Mutation, selection and genetic drift are the main forces shaping CUPrefs (Duret, 2002;
 62 Chamary et al., 2006; Plotkin and Kudla, 2011; Akashi, 1997). Mutational biases relate to directional mechanistic
 63 biases during genome replication (Reijns et al., 2015; Apostolou-Karampelis et al., 2016), during genome repair
 64 (Lujan et al., 2012), or during recombination (Pouyet et al., 2017), preferentially introducing one nucleotide over
 65 others or inducing recombination and maintaining genomic regions depending on their composition. Mutational
 66 biases are well described in prokaryotes and eukaryotes, ranging from simple molecular preferences towards
 67 3' A-ending in the Taq polymerase (Clark, 1988) to complex GC-biased gene conversion in vertebrates (Pouyet et al.,
 68 2017). Selective forces shaping CUPrefs are often described as translational selection. This notion refers to the
 69 ensemble of mechanistic steps and interactions during translation that are affected by the particular CUPrefs of the
 70 mRNA, so that the choice of certain codons at certain positions may actually enhance the translation process and
 71 can be subject to selection (Bulmer, 1991). Translational selection covers thus codon-independent effects on mRNA
 72 secondary structure, overall stability, and subcellular location (Presnyak et al., 2015; Novoa and Ribas de Pouplana,
 73 2012), but also codon-mediated effects acting on mRNA maturation, programmed frameshifts, translation speed and
 74 accuracy, or protein folding (Caliskan et al., 2015; Mordstein et al., 2020; Spencer and Barral, 2012). Translational
 75 selection has been demonstrated in prokaryotes and in some eukaryotes (Satapathy et al., 2016; Percudani et al.,
 76 1997; Duret and Mouchiroud, 1999; Whittle and Extavour, 2016), often in the context of tRNA availability (Ikemura,
 77 1981). Although its very existence in vertebrates remains highly debated (Pouyet et al., 2017; Galtier et al., 2018),
 78 experimental evidence shows that differences in CUPrefs of a focal gene impose an important translation burden in
 79 human cells (Picard et al., 2023).

80

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

90

91 The starting point for our research are the experimental observations by Robinson and coworkers reporting differen-
 92 tial expression of the polypyrimidine tract binding protein (*PTBP*) human paralogs as a function of their nucleotide
 93 composition (Robinson et al., 2008). Vertebrate genomes encode for three in-paralogous versions of the *PTBP* genes,
 94 all of them fulfilling mechanistically similar functions in the cell: they form a class of hnRNP RNA-Binding Proteins
 95 that are involved in the modulation of mRNAs alternative splicing (Pina et al., 2018). Within the same genome, the
 96 three paralogs display high amino-acid sequence similarity, around 70% in humans, and with similar overall values in
 97 vertebrates (Pina et al., 2018).

98 Despite the high resemblance at the protein level, the three *PTBP* paralogs sharply differ in nucleotide composition,
 99 CUPrefs, and supposedly in tissue expression pattern. In humans, *PTBP1* is enriched in GC3-rich synonymous codons
 100 and is widely expressed in all tissues, while *PTBP2* and *PTBP3* are AT3-rich and display an enhanced expression in
 101 the brain and in hematopoietic cells respectively (Supplementary Material, Figure S1). Robinson and coworkers stud-
 102 ied the expression in human cells in culture of all three human *PTBP* paralogous genes placed under the control of
 103 the same promoter. They showed that the GC-rich paralog *PTBP1* was more highly expressed than the AT-rich ones,
 104 and that the expression of the AT-rich paralog *PTBP2* could be enhanced by synonymous codons recoding towards the
 105 use of GC-rich codons (Robinson et al., 2008). Here we have built on the evolutionary foundations of this observation
 106 and extended the analyses of CUPrefs to *PTBP* paralogs in vertebrate genomes. Our results are consistent with a
 107 scenario in which paralog-specific directional changes in CUPrefs in mammalian *PTBPs* concurred with a process of
 108 subfunctionalisation by differential tissue pattern expression of the three paralogous genes.

109 4 Material and Methods

110 *Sequence retrieval*

111 We assembled a dataset of DNA sequences from 47 mammalian and 27 non-mammalian vertebrate genomes,
 112 and 3 from protostome genomes. Using the BLAST function on the nucleotide database of NCBI
 113 (NCBI Resource Coordinators, 2018) taking each of the human *PTBP* paralogs as references we looked for genes
 114 already annotated as *PTBP* orthologs (final sequence collection in November 2019; see supplementary Table S16 for
 115 accession numbers). We could retrieve the corresponding three orthologs in all vertebrate species screened, except for
 116 the European rabbit *Oryctolagus cuniculus*, lacking *PTBP1*, and from the rifleman bird *Acanthisitta chloris*, lacking
 117 *PTBP3*. The final vertebrate dataset contained 75 *PTBP1*, 76 *PTBP2* and 75 *PTBP3* sequences. As outgroups for the
 118 analysis, we retrieved the orthologous genes from three protostome genomes, which contained a single *PTBP* homolog
 119 per genome. We chose to resort to protostome sequences as outgroups because at the time of compiling our dataset
 120 we could not find well-annotated *PTBP* paralog sequences from Chordate taxa that could be used as sister clade to our
 121 vertebrate genomes. Our final dataset was consistent with the descriptions available in ENSEMBL and ORTHOMAM
 122 for the *PTBP* orthologs (Yates et al., 2020; Scornavacca et al., 2019; Pina et al., 2018). From the original dataset, we
 123 identified a subset of nine mammalian and six non-mammalian vertebrate species with a good annotation of the *PTBP*
 124 chromosome context. For these 15 species we retrieved local synteny and composition information on the *PTBP*
 125 flanking regions and introns (Supplementary Table S3). Because of annotation hazards, intronic and flanking regions
 126 information were missing for some *PTBPs* in the African elephant *Loxodonta africana*, Schlegel's Japanese Gecko
 127 *Gekko japonicus*, and the whale shark *Rhincodon typus* assemblies. For the selected 15 species the values for codon

128 adaptation index (CAI) (Sharp and Li, 1987) and codon usage similarity index (COUSIN) (Bourret et al., 2019) were
 129 calculated using the COUSIN server (available at <https://cousin.ird.fr>) (Supplementary Table S4).

130 *Codon Usage analysis*

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

144 *Alignment and phylogenetic analyses*

145 First, all sequences were aligned together, and we constructed a phylogenetic tree to verify whether each paralog as-
 146 sembly was monophyletic (Supplementary Figure S13). This was actually the case, and in this unbiased preliminary
 147 analysis all *PTBP1-3* were respectively monophyletic. Thus, to generate more robust alignments without introducing
 148 artefacts due to large evolutionary distances between in-paralogs, we proceeded stepwise, as follows: i) we aligned
 149 separately at the amino acid level each set of *PTBP* paralog sequences of mammals and non-mammalian vertebrates;
 150 ii) for each *PTBP* paralog we merged the alignments for mammals and for non mammals, obtaining the three *PTBP1*,
 151 *PTBP2* and *PTBP3* alignments for all vertebrates; iii) we combined the three alignments for each paralog into a sin-
 152 gle one; iv) we aligned the outgroup sequences to the global vertebrate *PTBPs* alignment. All alignment steps were
 153 performed using MAFFT with the `globalpair` option and 1000 max iterations (Katoh et al., 2002). The final amino
 154 acid alignment was used to obtain the codon-based nucleotide alignment. The codon-based alignment was trimmed
 155 using Gblocks using the default settings (Castresana, 2000) (All alignment data are available on [Zenodo](#)) Phylogenetic
 156 inference was performed at the amino acid and at the nucleotide level using RAxML v8.2.9, bootstrapping over 1000
 157 cycles (Stamatakis, 2014). For nucleotides we used codon-based partitions and applied the generalist GTR+I+G4
 158 model while for amino acids we applied the LG+G4 model (Waddell and Steel, 1997; Le and Gascuel, 2008). For the
 159 79 species used in the analyses we retrieved a species-tree from the TimeTree tool (Kumar et al., 2017). Distances be-
 160 tween phylogenetic trees were computed using the Robinson-Foulds index, which accounts for differences in topology
 161 (Robinson and Foulds, 1981), and the K-tree score, which accounts for differences in both topology and branch length
 162 (Soria-Carrasco et al., 2007). We then calculated pairwise distances between branches on the nucleotide and amino
 163 acid based trees and compared them against CUPrefs-based pairwise distances to measure the impact of CUPrefs on
 164 the phylogeny. After phylogenetic inference, we computed marginal ancestral states for the respectively most recent
 165 common ancestors at the nucleotide level of each paralog, using RAxML. For each position, the base with the max-

166 imum probability was used, and the sites for which RAxML could not infer with certainty the ancestral base were
 167 marked as missing data. We found 14%, 18% and 10% of missing bases respectively in *PTBP1*, *PTBP2* and *PTBP3*.
 168 Using these ancestral sequences we estimated the number of synonymous and non-synonymous substitutions of each
 169 extant sequence to the corresponding most recent common ancestor. We then compared the substitution matrices via
 170 a PCA analysis.

171 5 Results

172 *Vertebrate PTBP paralogs differ in nucleotide composition*

173 In order to understand the evolutionary history of *PTBP* genes, we performed first a nucleotide composition and
 174 CUPrefs analysis on the three paralogs in 79 species. Overall, *PTBP1* are GC-richer than *PTBP2* and *PTBP3* (re-
 175 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
 176 addition, *PTBP1*s show a difference in GC3 between mammalian and non-mammalian genes (respectively 79.8 against
 177 59.9 mean percentages). A linear regression model followed by a Tukey's honest significant differences analysis for
 178 GC3 using as explanatory levels paralog (*i.e.* *PTBP1-3*), taxonomy (*i.e.* mammalian or non-mammalian), and their
 179 interaction identifies three main groups of *PTBPs* (Table 1): a first one corresponding to mammalian *PTBP1*, a second
 180 one grouping non-mammalian *PTBP1*, and a third one encompassing all *PTBP2* and *PTBP3*. The largest explanatory
 181 factor for GC3 was the paralog *PTBP1-3*, accounting alone for 65% of the variance, while the interaction between the
 182 levels taxonomy and paralog captured around 15% of the remaining variance (Table 1). These trends are confirmed
 183 when performing paired comparisons between paralogs present in the same mammalian genome, with significant dif-
 184 ferences in GC3 content in the following order: *PTBP1* > *PTBP3* > *PTBP2* (Wilcoxon signed rank test: *PTBP1* vs
 185 *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;
 186 *PTBP3* vs *PTBP2*, mean diff=4.5, S=406.50, p-value <0.0001). Note that even if all of them significantly different,
 187 the mean paired differences in GC3 between *PTBP1* and *PTBP2-3* are ten times larger than the corresponding mean
 188 paired differences between *PTBP2* and *PTBP3*.

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

199 In many vertebrate species, strong compositional heterogeneities are observed along chromosomes with an arrange-
 200 ment of AT-rich and GC-rich regions, often referred to as "isochores". To explore the influence of this genomic
 201 environment on the nucleotide composition of *PTBPs*, we analyzed for 15 species with well-annotated genomes the
 202 correlation of paralog GC3 with two local compositional variables of the corresponding gene (GC content of intronic

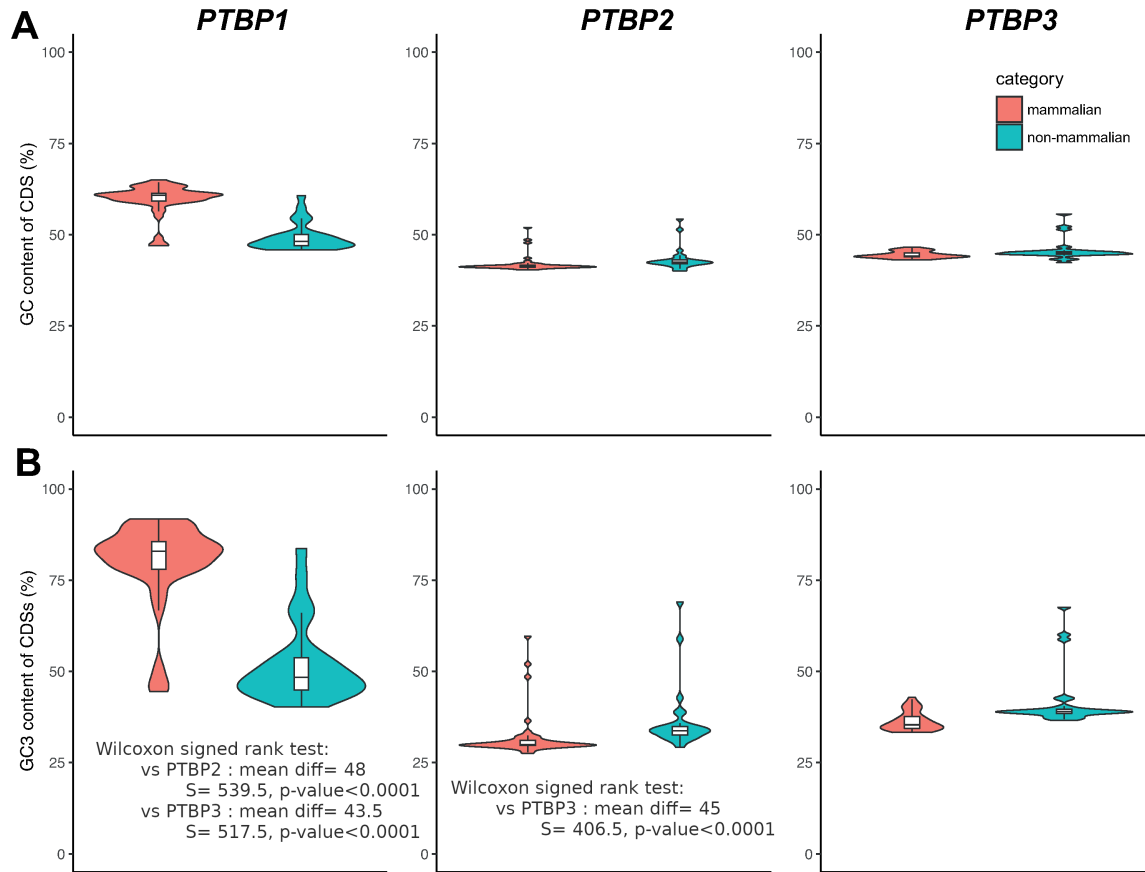


Figure 1: **GC content (A) and GC3 content (B) of vertebrates *PTBPs*.** Violin plots display the overall distribution, while box and whiskers display median, quartiles and 95% of the corresponding values for mammalian (red) and non-mammalian (blue) individual genomes. The results of a the paired Wilcoxon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

203 and flanking regions) and with three global compositional variables for the corresponding genomes (global GC3 in
 204 the complete genomic ORFome, global GC content in all introns, and global GC content in all flanking regions)(Table
 205 **3** and Figure **2**). First, for *D. rerio* the GC3 composition of *PTBP2* and *PTBP3* is clearly different from the rest,
 206 in line with the outlier results presented in Table **2**. We have thus excluded the zebra fish values and performed an
 207 individual as well as a stepwise linear fit to explain the variance in GC3 composition by the variance in the local and
 208 global compositional variables mentioned above (Table **3**). For all three *PTBPs* the local GC content explains best the
 209 corresponding GC3 content, but with strong differences between paralogs: while variation in the local composition
 210 captures almost perfectly variation in the GC3 content of *PTBP1* ($R^2=0.97$) and relatively well in the case of *PTBP2*
 211 ($R^2=0.46$), the fraction of variance explained by the local composition significantly drops for *PTBP3* ($R^2=0.15$). It
 212 must be noted nevertheless that the GC3 variable ranges are different among paralogs, so that variation in GC3 values
 213 for *PTBP1* (roughly between 40% and 90%) is larger than for *PTBP2-3* (respectively 29%-38% and 34%-46%). This
 214 larger variable span in the case of *PTBP1* may allow for an increased power for detecting a significant correlation in
 215 composition values for this paralog.

216 **Vertebrate *PTBP* paralogs differ in CUPrefs**

217 For each *PTBP* coding sequence we extracted the relative frequencies of synonymous codons and performed different
 218 approaches to reduce information dimension and visualise CUPrefs trends. The results of a principal component
 219 analysis (PCA) are shown in Figure 3 as well as in Supplementary Figure S5. The first PCA axis captured 68.9% of
 220 the variance, far before the second and the third axes (respectively 6.7% and 3.2%). Codons segregate in the first axis
 221 by their GC3 composition, the only exception being the UUG-Leu codon, which grouped together with AT-ending
 222 codons. This first axis differentiates mammalian *PTBP1*s on the one hand and *PTBP2*s and *PTBP3*s on the other hand.
 223 Non-mammalian *PTBP1*s scatter between mammalian *PTBP1*s and *PTBP3*s, along with the protostomes *PTBP*s. In
 224 the second PCA axis the only obvious (but nevertheless cryptic) codon-structure trends are: i) the split between
 225 C-ending and G-ending codons, but not between U-ending and A-ending codons; and ii) the large contribution in
 226 opposite directions to this second axis of the AGA and AGG-Arginine codons. This second PCA axis differentiates
 227 *PTBP2*s from *PTBP3*s paralogs, consistent with these composition trends. A paired-comparison confirms that *PTBP3*s
 228 are richer in C-ending codons than *PTBP2*s in the same genome, respectively 21.7% against 15.4% (Wilcoxon signed
 229 rank test: mean diff=6.2, S=1184.0, p-value <0.0001).

230 As an additional way to identify groups of genes with similar CUPrefs, we applied a hierarchical clustering and a
 231 k-means clustering. Both analyses mainly aggregate *PTBP* genes by their GC3 richness. The *PTBP* dendrogram

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

232 resulting of the hierarchical clustering shows five main clades that cluster the paralogs with a good match to the
 233 following groups: mammalian *PTBP1*s, non-mammalian *PTBP1*s, *PTBP2*s, *PTBP3*s and a fifth group containing the
 234 protostomes *PTBP*s as well as a few individuals of all three paralogs (rows in clustering in Figure 3; Kappa-Fleiss
 235 consistency score = 0.76). Regarding codon clustering, the hierarchical stratification sharply splits GC-ending codons
 236 from AT-ending codons, with the only exception again of the UUG-Leu codon, which consistently groups within
 237 the AT-ending codons. The elbow approach of k-means clustering identifies an optimal number of four clusters and
 238 separates the paralog genes with a good match as following: *PTBP1*, *PTBP2*, *PTBP3* and a group containing the
 239 protostomes as well as some individuals from all paralogs (Kappa-Fleiss consistency score = 0.75).

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

246 In order to characterise the directional CUPrefs bias of the different paralogs, we have analysed, for the 15 species
 247 with well-annotated genomes described above, the match between each individual *PTBP* and the average CUPrefs of
 248 the corresponding genome (Table 4). The COUSIN quantitative values compare the CUPrefs of a query sequence with
 249 those of a reference (in our case the coding genome of the corresponding organism), and can be directly interpreted and

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 harrisi</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

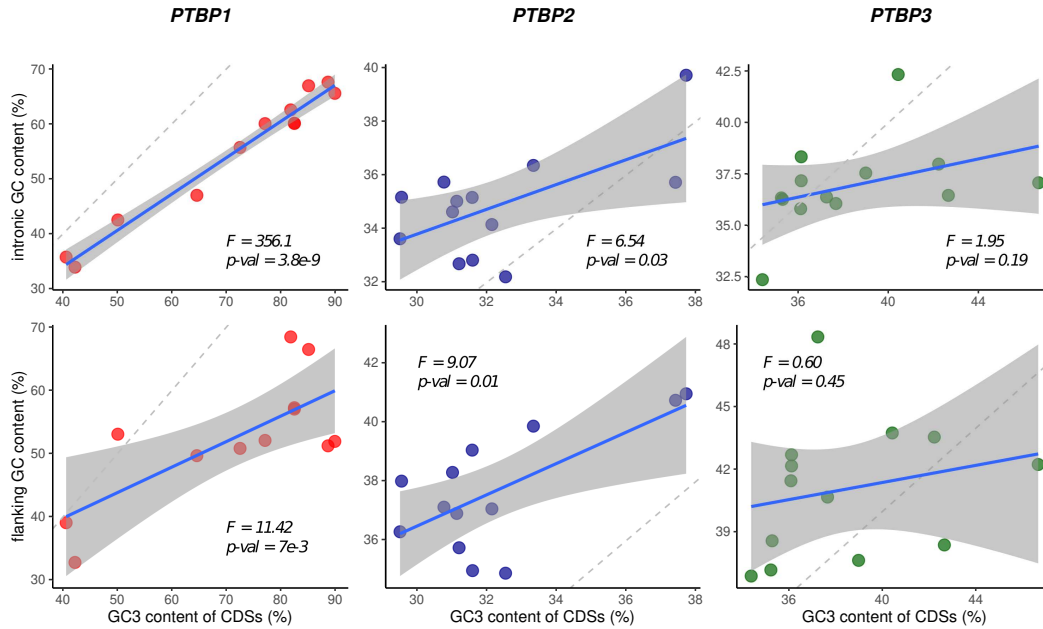


Figure 2: **Variation in GC3 content of PTBPs (x-axis) and in the GC content of the corresponding introns (A, y axis) or flanking regions (B, y axis).** Each dot represents one of the 15 individual genomes used for the genomic context analysis. For each graph, we performed a linear regression modelling (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.

250 compared in a qualitatively way, as described (Bourret et al., 2019). Briefly, COUSIN values around 1 reflect similar
 251 CUPrefs in the query sequence and in the reference, while values around 0 reflect CUPrefs close to random in the
 252 query sequence; COUSIN values above 1 reflect similar directional trends in CUPrefs in the query sequence and in the
 253 reference, but with stronger bias in the query sequence; COUSIN negative values reflect opposite CUPrefs between the
 254 query sequence and the reference. Our results highlight strong differences for mammalian paralogs: *PTBP1*s display
 255 COUSIN values above 1 while *PTBP2*s display COUSIN values below zero. The COUSIN results and interpretation
 256 are provided in (Supplementary Figure S14). These results mean that, in mammals, *PTBP1*s are enriched in codons
 257 commonly used in the corresponding genome, while *PTBP2*s are enriched in codons rarely used in the corresponding
 258 genome, to the extent that their CUPrefs go in the opposite direction to the average in the genome. As for *PTBP3* in
 259 mammals, we observe COUSIN values below 0 in most cases or very close to 0 in the case of the horse *Equus caballus*
 260 and house mouse *Mus musculus*, implying a trend towards rare codons. In non-mammals, in contrast, *PTBPs* show an
 261 overall similarity to their respective reference genomic CUPrefs.

262 Phylogenetic reconstruction of PTBPs

263 We explored the evolutionary relationships between *PTBPs* by phylogenetic inference at the amino acid and at the
 264 nucleotide levels (Figure 4, Supplementary Figure S10). Our final dataset contained 74 *PTBP* sequences from mam-
 265 mals (47 species within 39 families) and non mammal vertebrates (27 species within 24 families). We used the *PTBP*
 266 genes from three protostome species as outgroup. Both amino acid and nucleotide phylogenies rendered three main

267 clades grouping the *PTBPs* by orthology, so that all *PTBP1-3* orthologs were correspondingly monophyletic. In both
 268 topologies, *PTBP1* and *PTBP3* orthologs cluster together, although the protostome outgroups are linked to the tree by
 269 a very long branch, hampering the proper identification of the vertebrate *PTBP* tree root. Amino acid and nucleotide
 270 subtrees were largely congruent (see topology and branch length comparisons in Table 5). The apparently large nodal
 271 and split distance values between nucleotide and amino acid for *PTBP2* trees stem from disagreements in very short
 272 branches, as evidenced by the lowest K-tree score for this ortholog (as a reminder, the Robinson-Foulds index exclu-
 273 sively regards topology while the K-tree score combines topological and branch-length dependent distance between

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

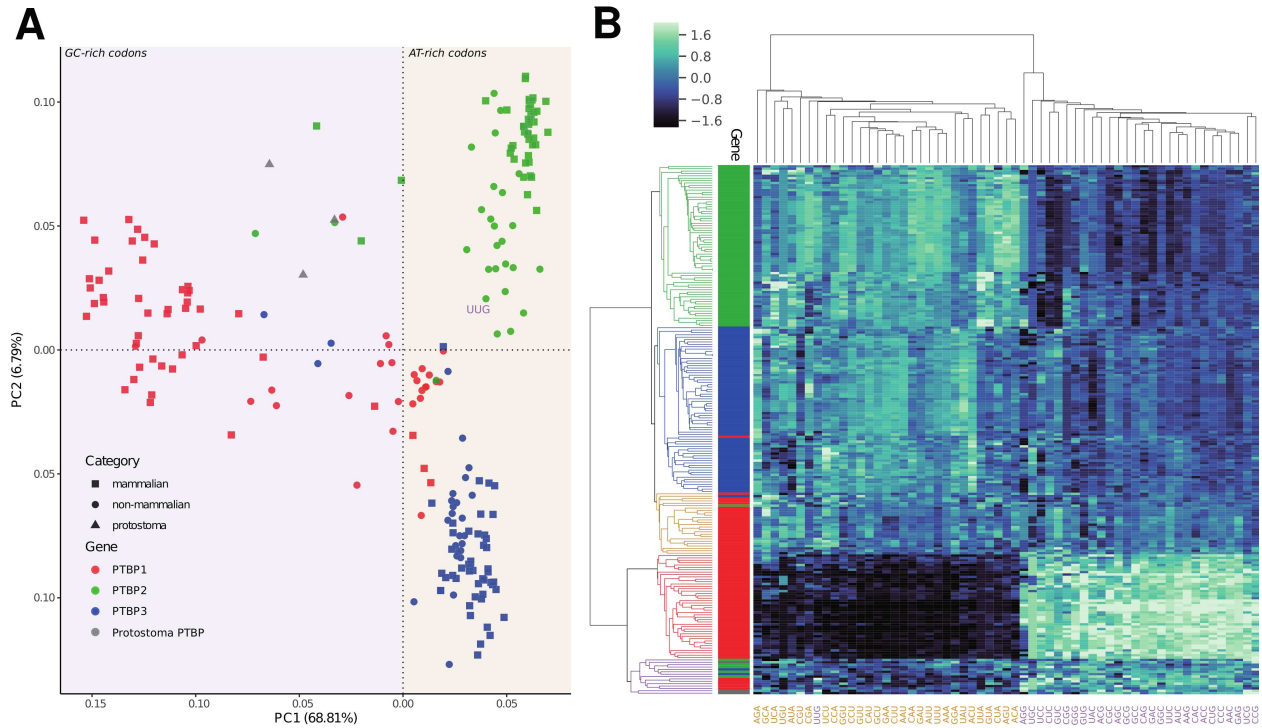


Figure 3: **CUPrefs analysis of PTBPs.** A) Plot of the two first dimensions of a PCA analysis based on the codon usage preferences of *PTBP1*s (red), *PTBP2*s (green), *PTBP3*s (blue) and protostome outgroup (grey) individual genes. Taxonomic information is included labelling mammals (squares), non-mammals (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 the only exception compared to the global trend of variable distribution (see (Supplementary Figure S5) for 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 protostome genes (grey). Note again the position of the UUG-Leu codon in the codon dendrogram, as the sole GC-ending codon clustering (in purple) among AT-ending codons (in orange)

274 trees, see Material and Methods). In all three cases, internal structure of the ortholog trees essentially recapitulates
 275 species taxonomy at the higher levels (Table 5). Some of the species identified by the regression analyses to display
 276 largely divergent nucleotide composition from the expected one given their taxonomy (Table 2) presented accordingly
 277 long branches in the phylogenetic reconstruction, such as *PTBP3* for *O. mykiss*, or rendered polyphyletic branching,
 278 as described above for *PTBP1* in mammals.

279 We have then analysed the correspondence between nucleotide-based and amino acid-based pairwise distances to eval-
 280 uate the impact of CUPrefs on the obtained phylogeny. We observe a good correlation between both reconstructions

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

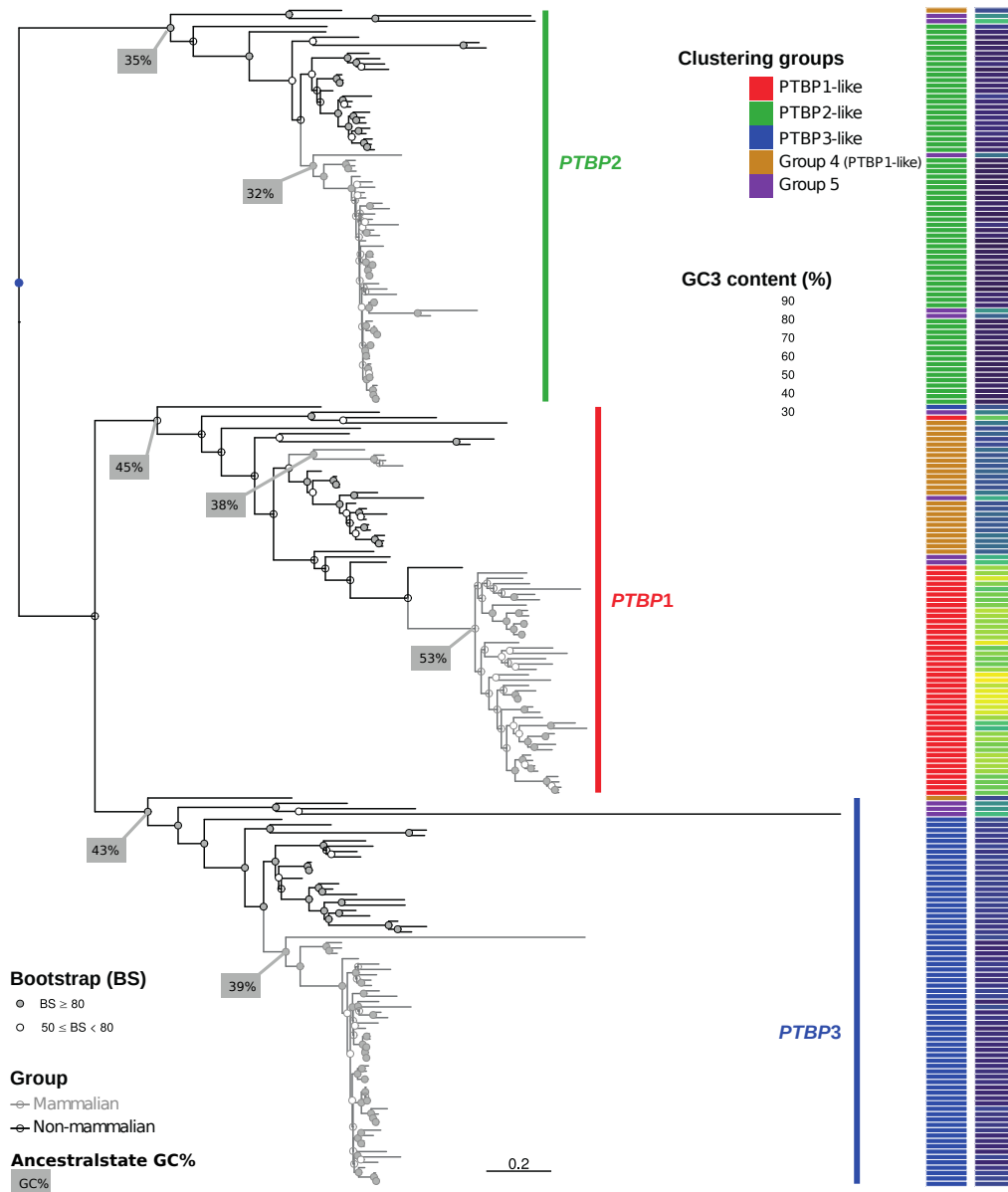


Figure 4: **Maximum-likelihood nucleic acid phylogeny of *PTBP* genes.** The phylogram depicts *PTBP2*s (green side bar), *PTBP1*s (red side bar) and *PTBP3*s (blue side bar) clades. The outgroup genes from protostomes are not shown to focus on the scale for vertebrate *PTBPs*, but their placement on the tree and the polarity they provide for vertebrate *PTBPs* is given by the blue dot. Gray branches indicate mammalian *PTBPs*, while black branches indicate non-mammalian species. Note the polyphyly for mammals with regards to *PTBP1*s, with the monotremes and marsupial clade not clustering together with the placental mammals clade. 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 inferred for the main ancestral nodes is indicated in grey boxes.

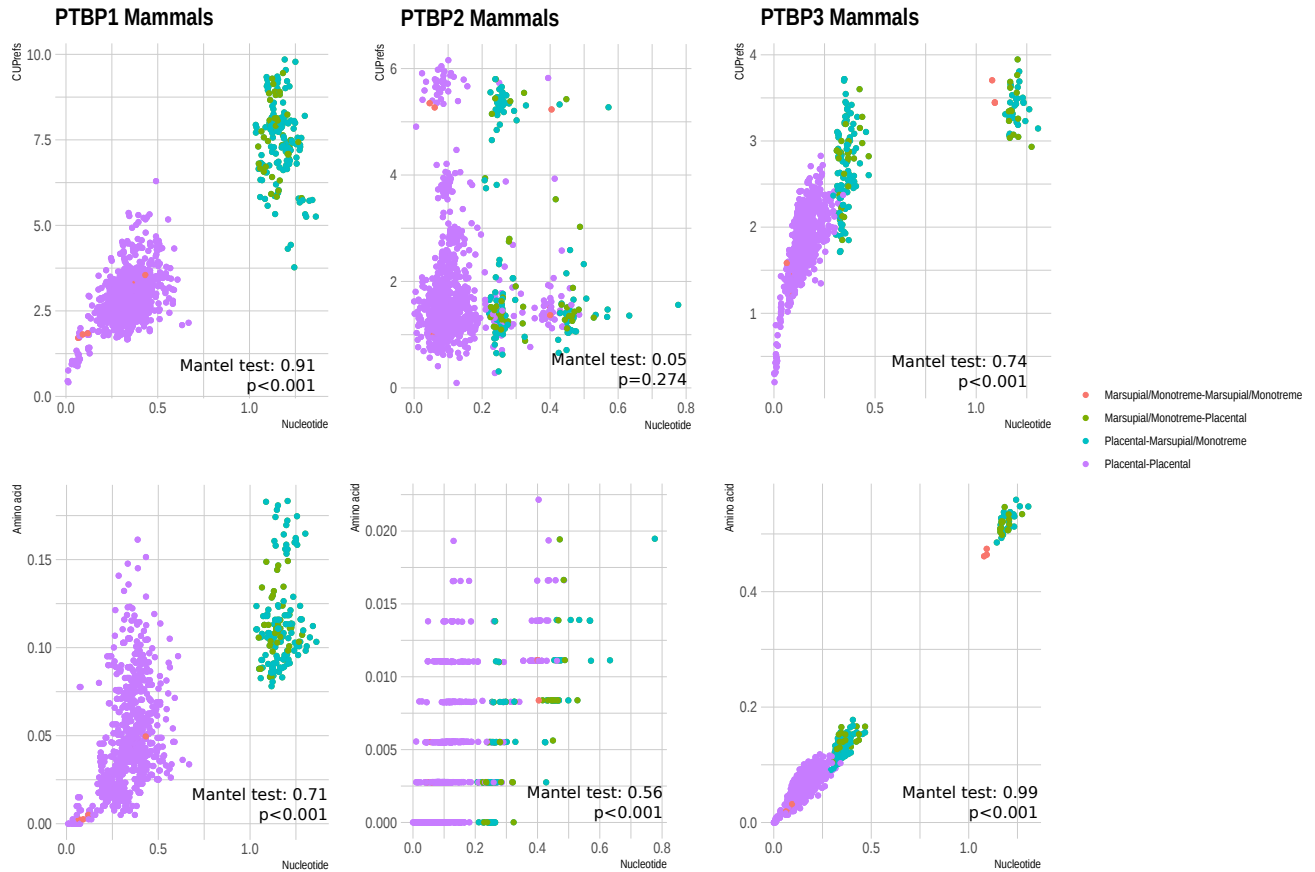


Figure 5: Nucleotide-based pairwise distances in the x-axis against CUPrefs-based (first row) and amino acid-based (second row) pairwise distances in the y-axis for the different mammalian *PTBP* orthologs. The results for a Mantel test assessing the correlation between the corresponding matrices are shown in each inset. The dots are coloured based on the taxonomic group of the compared species

281 for all paralogs, except for mammalian *PTBP2*s, which display extremely low divergence at the amino acid level (see
 282 Figure 5 for values in mammalian paralogs, Supplementary Figure S8 for non-mammalian paralogs, and Supplemen-
 283 tary Table S7 for the correlation between nucleotide-based and amino acid-based pairwise distances). For mammalian
 284 *PTBP1*s, the plot allows to clearly differentiate a cloud with the values corresponding to monotremes+marsupials,
 285 split apart from placental mammals in terms of both amino acid and nucleotide distances. This distribution matches
 286 well the fact that sequences from monotremes and marsupials cluster separately from placental mammals in the *PTBP1*
 287 phylogeny (see grey branches being polyphyletic for *PTBP1* in Figure 4). The same holds true for the platypus *PTBP3*,
 288 extremely divergent from the rest of the mammalian orthologs. The precise substitution patterns are analysed in detail
 289 below. The histograms describing the accumulation of synonymous and non-synonymous substitutions confirm that
 290 mammalian *PTBP1*s have accumulated the largest number of synonymous substitutions compared to non-mammalian
 291 *PTBP1*s and to other orthologs (Supplementary Figure S9).

292 We have finally analysed the connection between nucleotide-based evolutionary distances within *PTBP* paralogs and
 293 CUPrefs-based distances (Figure 5 for mammalian paralogs and Supplementary Figure S8 for non-mammalian par-

294 alogs). A trend showing increased differences in CUPrefs as evolutionary distances increase is evident only for
 295 *PTBP1s* and *PTBP3s* in mammals. For mammalian *PTBP1s* the plot clearly differentiates a cloud with the values
 296 corresponding to monotremes and marsupials splitting apart from placental mammals in terms of both evolutionary
 297 distance and CUPrefs. For mammalian *PTBP2s* the plot captures the divergent CUPrefs of the platypus and of the bats
 298 *M. natalensis* and *Desmodus rotundus*, while for non-mammalian *PTBP2s* the divergent CUPrefs of the rainbow trout
 299 (*O. mykiss*) are obvious. Finally, for mammalian *PTBP3s* the large nucleotide divergence of the platypus paralog is
 300 evident. Importantly, all these instances of divergent behaviour (except for the platypus *PTBP3*) are consistent with the
 301 deviations described above from the expected composition by the mathematical modelling of the ortholog nucleotide
 302 composition (Table 2).

303 *Mammalian PTBP1s accumulate GC-enriching synonymous substitutions*

304 We have shown that *PTBP1* genes are GC-richer and specifically GC3-richer than the *PTBP2* and *PTBP3* paralogs in
 305 the same genome, and that this enrichment is of a larger magnitude in *PTBP1s* from placental mammals. We have
 306 thus assessed whether a directional substitutional pattern underlies this enrichment, especially regarding synonymous
 307 substitutions. For this we have inferred the ancestral sequences of the respective most recent common ancestors of
 308 each *PTBP* paralog, recapitulated synonymous and non-synonymous substitutions between each *PTBP* individual and
 309 their ancestors, and constructed the corresponding substitution matrices (Table S11). The two first axes of a principal
 310 component analysis using these substitution matrices capture, with a similar share, 66.95% of the variance between
 311 individuals (Figure 6). The first axis of the PCA separates synonymous from non-synonymous substitutions. Intrigu-
 312 ingly though, while T<->C transitions are associated with synonymous substitutions, as expected, G<->A transitions
 313 are instead associated with non-synonymous substitutions. The second axis separates substitutions by their effect on
 314 nucleotide composition: GC-stabilizing/enriching on one direction, AT-stabilizing/enriching on the other one. Strik-
 315 ingly, the substitutional spectrum of mammalian *PTBP1s* sharply differs from the rest of the paralogs. Substitutions
 316 in mammalian *PTBP1* towards GC-enriching changes, in both synonymous and non-synonymous compartments, are
 317 the main drivers of the second PCA axis. In contrast, synonymous substitutions in *PTBP3* as well as all substitutions
 318 in *PTBP2* tend to be AT-enriching. Finally, the substitution trends for *PTBP1* in mammals are radically different
 319 from those in non-mammals, while for *PTBP2* and *PTBP3s* the substitution patterns are similar in mammals and
 320 non-mammals for each of the compartments synonymous and non-synonymous.

321 6 Discussion

322 The non equal use of synonymous codons has fascinated biologists since it was first described. It has given rise to
 323 fruitful (and unfruitful) controversies between defenders of *all-is-neutralism* and defenders of *all-is-selectionism* (see
 324 for instance the discussion in the late 60s between Jack Lester King and Thomas H. Jukes on the one side and Bryan
 325 Clarke on the other side (King and Jukes, 1969; Clarke, 1970)), and has launched further the quest for additional molec-
 326 ular signaling beyond codons themselves (Callens et al., 2021). The main questions around CUPrefs are twofold. On
 327 the one hand, their origin: to what extent they are the result of fine interplay between mutation and selection processes
 328 or whether they may be the result of bottlenecks and genomic drift. On the other hand, their functional implica-
 329 tions: whether and how particular CUPrefs can be linked to specific gene expression regulation processes, broadly
 330 understood as downstream effects that modify the kinetics and dynamics of DNA transcription, mRNA maturation

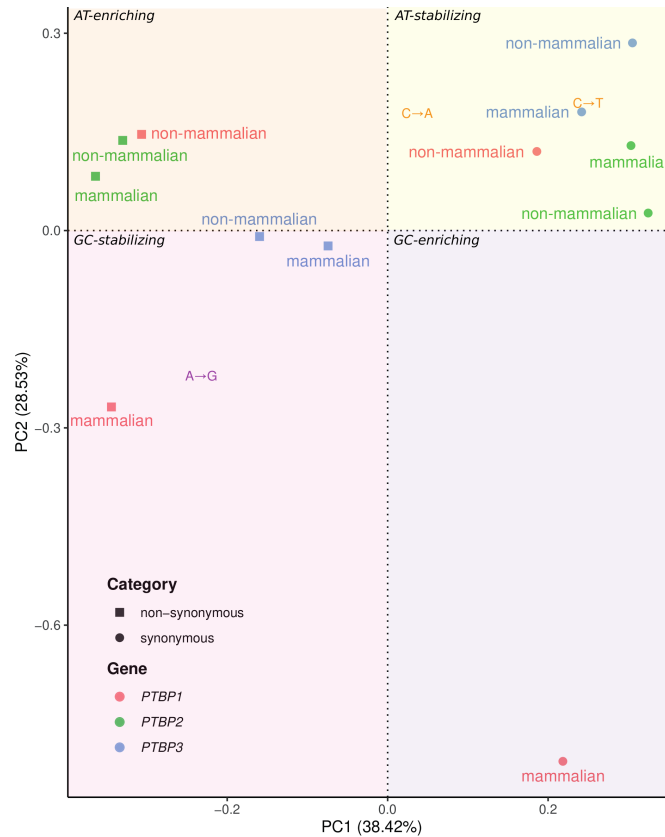


Figure 6: **Spectra of synonymous and non-synonymous substitutions for *PTBPs*.** This principal component analysis (PCA) has been built using the observed nucleotide synonymous and non-synonymous substitution matrices for each *PTBP* paralog, inferred after phylogenetic inference and comparison of extant 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 patterns (*i.e.* A->G, C->A and C->T), which have been plotted according to their eigenvalues (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).

331 and stability, mRNA translation, and/or protein folding and stability. In the present work we have built on the ex-
 332 perimental results of Robinson and coworkers, which communicated the differential expression of the *PTBP* human
 333 gene paralogs as a function of their CUPrefs (Robinson et al., 2008). From this particular example, we have aimed
 334 at exploring the nature of the connection between paralogous gene evolution and CUPrefs. Our results show that the
 335 three *PTBP* paralogous genes, which show divergent expression patterns in humans and in other mammals, also have
 336 divergent nucleotide composition and CUPrefs not just in humans but in most vertebrate species. We elaborate here on
 337 Robinson and coworkers' experimental findings and propose that this evolutionary pattern could be compatible with a
 338 phenomenon of phenotypic evolution by sub-functionalisation (in this case specialisation in tissue-specific expression

339 levels), linked to genotypic evolution by association to specific CUPrefs patterns. Such conclusions invite to pursue
 340 Robinson and coworkers' efforts by comparing *PTBPs* CUPrefs-modulated expression among numerous vertebrate
 341 cell lines, especially between mammals and non-mammals ones. Consistent with studies on other paralog fam-
 342 ilies (Munk et al., 2022; Lampson et al.), our results suggest, more generally, that a detailed analysis of differential
 343 CUPrefs in paralogs may help understand the divergent/convergent mutation-selection pressures that could underlie
 344 their functional differences.

345 We have reconstructed the phylogenetic relationships and analysed the evolution and diversity of CUPrefs among
 346 *PTBP* paralogs within 74 vertebrate species. The phylogenetic reconstruction shows that the genome of ancestral
 347 vertebrates already contained the three extant *PTBP* paralogs. This is consistent with the ortholog and paralog identi-
 348 fication in the databases ENSEMBL and ORTHOMAM (Yates et al., 2020; Scornavacca et al., 2019; Pina et al., 2018).
 349 Although our results suggest that *PTBP1* and *PTBP3* are sister lineages, the distant relationship between the vertebrate
 350 genes and the protostome outgroup precludes the inference of a clear polarity between vertebrate *PTBPs*. We commu-
 351 nicate an important deviation in terms of expected nucleotide composition for the paralogous genes from the rainbow
 352 trout and the zebrafish. We have not explored the impact of the full genome duplication round that is exclusive to the
 353 Actinopterygia lineage onto the diversity and the repertoire of the *PTBP* paralogs (Meyer and Schartl, 1999), as we
 354 have focused our analyses on the impact in mammals, but this our results suggest that the elucidation of the impact
 355 of full-genome duplication on the repertoire of in-paralogs could be an interesting line of questioning. We identify
 356 no occurrence of basal replacement between paralogs, which may have appeared, for instance, as the replacement of
 357 an 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.
 358 Instead, the basal evolutionary histories of the different *PTBPs* comply well with those of the corresponding species.
 359 The most blatant mismatch between gene and species trees is the polyphyly of mammalian *PTBP1*s: monotremes and
 360 marsupials constitute a clade, separate from the placental mammals clade. Further, multiple findings in our results
 361 show sharp, contrasting patterns between *PTBP1* and the *PTBP2-3* paralogs: i) the excess of accumulation of syn-
 362 onymous substitutions in mammalian *PTBP1*s for a similar total number of changes (Supplementary Figure S9 and
 363 Table S11); ii) the larger differences in CUPrefs between genes with a similar total number of nucleotide changes in
 364 the case of *PTBP1*s in mammals (Figure 5 A); iii) the explicitly different spectrum of synonymous substitutions in
 365 *PTBP1*s, enriched in A->C, T->G and T->C changes (Figure 6); iv) the sharp difference of CUPrefs between *PTBP1*s
 366 and *PTBP2-3*s; and v) the clustering of *PTBP1* genes in monotremes and marsupials together with *PTBP1* genes in
 367 non-mammals according to their CUPrefs (Figure 3 A). Overall, the particular nucleotide composition and the associ-
 368 ated CUPrefs in mammalian *PTBP1* genes are most likely associated to specific local substitution biases as shown by
 369 the strong correlation between coding and non-coding GC content in *PTBP1* orthologs, while CUPrefs in *PTBP2-3*s
 370 cannot be explained alone by such local substitution biases (Figure 2; Table 3).

371 While GC3-rich nucleotide composition and CUPrefs of mammalian *PTBP1*s are dominated by local substitution
 372 biases, this is not the case for mammalian *PTBP2*, overall AT3-richer and without any clear correlation between
 373 coding and non-coding GC content among the studied species (Figure 2; 3). As mentioned above, a note of caution
 374 should be raised here, as the variable range for GC composition among *PTBP1*s is larger than for *PTBP2-3*s, so that co-
 375 variation analyses may have less power for the latter paralogs. In vertebrates, nucleotide composition varies strongly
 376 along chromosomes, so that long chromatin stretches, historically named "isochores", appear enriched in GC or in
 377 AT nucleotides and present particular physico-chemical profiles (Caspersson et al., 1968). Local mutational biases

378 and GC-biased gene conversion mechanism may underlie such heterogeneity, predominantly shaping local nucleotide
 379 composition in numerous vertebrates genomes, so that the physical location of a gene along the chromosome largely
 380 explains its CUPrefs (Holmquist, 1989). In agreement with these hypotheses for local mutational biases, variation in
 381 GC3 composition of *PTBP1*s is almost totally ($R^2=0.97$) explained by the variation in local GC composition (Figure
 382 2, Table 3), suggesting that a similar substitution bias has shaped the GC-rich composition of the flanking, intronic and
 383 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
 384 gene conversion is often invoked as a powerful mechanism underlying such local GC-enrichment processes, leading
 385 to the systematic replacement of the alleles with the lowest GC composition by a GC-richer homolog (Marais, 2003).
 386 It has been proposed that gene expression during meiosis (evaluated as mRNA detection) correlates with a decreased
 387 probability of GC-biased gene conversion during meiotic recombination (Pouyet et al., 2017). Expression of *PTBP1*
 388 in human cells is documented during meiosis in the oocyte germinal line and expression of the AT-rich *PTBP2* has
 389 been observed during spermatogenic meiosis (Zagore et al., 2015; Hannigan et al., 2017). Expression during meiosis
 390 might thus have hindered GC-biased gene conversion for *PTBP1-2*s, provided that this expression pattern observed
 391 in humans was displayed also by the mammalian ancestor and that it is shared between mammalian species. With
 392 these assumptions, and thus, with caution, the GC-richness of *PTBP1* cannot be accounted for by GC-biased gene
 393 conversion, while the low GC content of *PTBP2* could be explained by an accumulation of GC→AT and AT→AT
 394 substitutions. All this notwithstanding, our results show that GC3 enrichment in mammalian *PTBP1* and the concurrent
 395 trend for enriched use of common codons are associated mostly with placental mammals, and that non-placental
 396 mammals display divergent composition and differ from the model expectations. This synapomorphy of a sudden
 397 change in nucleotide composition is strongly compatible with a GC-biased gene conversion event in the placental
 398 ancestor that may have led to fixation of the ancestral version of the extant GC-rich *PTBP1*. Regarding *PTBP3*, the
 399 low GC-content together with the low correlation with either coding nor non-coding local GC-content could indicate
 400 that other mechanisms may shape the observed CUPrefs for this paralog.

401 In mammals, global GC-enriching genomic biases strongly impact CUPrefs, so that the most used codons in average
 402 tend to be GC-richer (Hershberg and Petrov, 2009). For this reason, mammalian GC3-rich *PTBP1*s match better the
 403 average genomic CUPrefs than AT3-richer *PTBP2* and *PTBP3*, which display CUPrefs in the opposite direction to the
 404 average of the genome. In the case of humans, *PTBP1* presents a COUSIN value of 1.75, consistent with a substantial
 405 enrichment in frequently-used codons, while on the contrary, the COUSIN values of -0.48 for *PTBP2* and of -0.23 for
 406 *PTBP3* point towards a strong enrichment in rarely-used codons (Supplementary Table S4). The poor match between
 407 human *PTBP2* CUPrefs and the human average CUPrefs could result in low expression of these genes in different
 408 human and murine cell lines, otherwise capable of expressing *PTBP1* at high levels and of expressing *PTBP3* at a
 409 lesser degree (Robinson et al., 2008). The barrier to *PTBP2* expression seems to be the translation process, as *PTBP2*
 410 codon-recoding towards GC3-richer codons results in strong protein production in the same cellular context, without
 411 significant changes in the corresponding mRNA levels (Robinson et al., 2008). Similar results to those of Robinson
 412 and coworkers have been more recently communicated on studies using the small Ras GTPases in human cells, in
 413 which highly similar paralogs displayed largely different expression patterns in terms of translation efficiency, that
 414 could be reverted by codon recoding strategies (Lampson et al.). Indeed, experimental results in human cells have
 415 shown that synonymous variants with large CUPrefs differences display strong phenotypic differences in translation
 416 efficiency (Picard et al., 2023). Overall, codon recoding strategy towards "preferred" codons (understood here as

417 the most commonly used codons in a genome) has become a standard practice for gene expression engineering that
 418 provides with very good expression results, despite our lack of understanding about the whole impact of local and
 419 global gene composition, nucleotide CUPrefs, and mRNA structure on gene expression (Brule and Grayhack, 2017).

420 The poor expression ability of *PTBP2* in human cells, the increase in protein production by the introduction of common
 421 codons, along with substitution biases failing to explain entirely *PTBP2* nucleotide composition and CUPrefs, raise
 422 the question of the adaptive value of poor CUPrefs in this paralog. Specific tissue-dependent or cell-cycle dependent
 423 gene expression regulation patterns have been invoked to explain the codon usage-limited gene expression for certain
 424 human genes, such as *TLR7* or *KRAS* (Newman et al., 2016; Lampson et al., 2013; Fu et al., 2018). In the case of AT-
 425 rich genes in vertebrates, such as *PTBP2*, it has been suggested that enrichment in less-used codons (*i.e.* A/T-ending
 426 codons in the case of vertebrates) may be linked to conserved, coordinated expression regulation over phylogeny
 427 and across ontogeny (Benisty et al., 2023). The expression levels of the three *PTBP* paralogs are tissue-dependent
 428 in humans (Supplementary Figure S1) as well as through mammals (Supplementary Figure S12) (Keppetipola et al.,
 429 2012; Wagner and Garcia-Blanco, 2002; Spellman et al., 2007). In the case of the duplicated genes, subfunctional-
 430 isation through specialisation in spatio-temporal gene expression has been proposed as the main evolutionary force
 431 driving conservation of paralogous genes (Ferris and Whitt, 1979). Such differential gene expression regulation in paralogs
 432 has actually been documented for a number of genes at very different taxonomic levels (Donizetti et al., 2009;
 433 Guschanski et al., 2017; Freilich et al., 2006). Specialised expression patterns in time and space can result in antagonistic
 434 presence/absence of the paralogous proteins (Adams et al., 2003). This is precisely the case of *PTBP1* and
 435 *PTBP2* during human central nervous system development: in non-neuronal cells, *PTBP1* represses *PTBP2* expres-
 436 sion by the skip of the exon 10 during *PTBP2* mRNA maturation, while during neuronal development, the micro RNA
 437 miR124 down-regulates *PTBP1* expression, which in turn leads to up-regulation of *PTBP2* (Keppetipola et al., 2012;
 438 Makeyev et al., 2007). Regarding non-human species, the available data about tissue-dependent and/or ontogeny-
 439 dependent differential expression at the transcription level (Abugessaisa et al., 2021) are largely concordant with the
 440 human data for *PTBP*, showing a tissue-wide transcription of *PTBP1*, a more restricted one for *PTBP3* together with an
 441 enrichment of *PTBP2* transcription in the central nervous system, as exemplified in the mouse (Barbosa-Morais et al.,
 442 2012), in the rat (Yu et al., 2014), in the cow (Merkin et al., 2012), in the gray short-tailed opossum (Brawand et al.,
 443 2011), or in the chicken (Barbosa-Morais et al., 2012). Finally, despite the high level of amino acid similarity be-
 444 tween both proteins, *PTBP1* and *PTBP2* seem to perform complementary activities in the cell and to display different
 445 substrate specificity, so that they are not directly inter-exchangeable by exogenous manipulation of gene expression
 446 patterns (Vuong et al., 2016).

447 In addition to local genomic context analyses, we explored *PTBP* chromosomal location and local synteny (Figure 7).
 448 The results show that, while it is clear that the position of human *PTBP1* is telomeric and thus in one of the GC-richer
 449 region of human chromosome 9, most *PTBPs* do not map to the telomeres. Therefore, while the specific location of
 450 human *PTBP1* may have influenced its CUPrefs, it is unclear whether the chromosomal location of *PTBPs* have an
 451 impact on observed nucleotide composition. Local synteny of *PTBPs* genes seems further to be conserved, with some
 452 exceptions: most mammalian *PTBP1*s reside in a conserved local synteny context that differs from non-mammalian
 453 species, with the exception of *D. rerio*. For *PTBP2* and *PTBP3* local synteny seems conserved between mammalian
 454 and non-mammalian species again with the exception of *D. rerio*, lacking the *SUSD1* gene between *PTBP3* and
 455 *UGCG*. Such results could indicate that vertebrate radiation has been followed up by a change of *PTBP1* genomic

456 context, with a swapping in flanking genes in mammalian branches. These results could be related to the observed
457 *PTBP1* differential GC-content between mammalian and non-mammalian species.

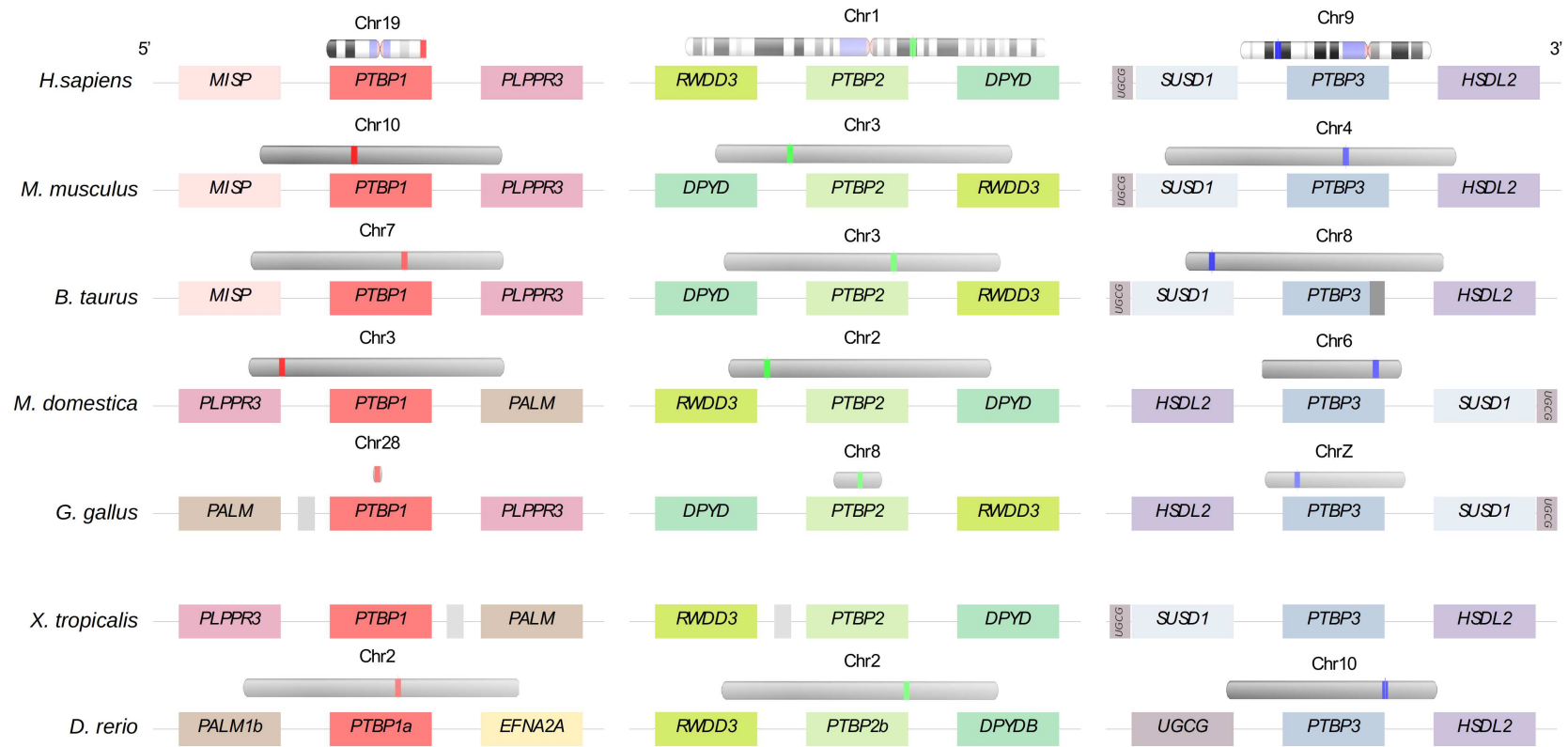


Figure 7: Placement on the chromosomes and genomic context of the three *PTBP* paralogs in a subset of the studied species.

458 In a different subject, we want to drive the attention of the readers towards the puzzling trend of the UUG-Leu codon
 459 in our CUPrefs analyses. This UUG codon is the only GC-ending codon that systematically clusters with AT-ending
 460 codons in all our analyses on CUPrefs, and that does not show the expected symmetrical behaviour with respect to
 461 its UUA-Leu counterpart codon (see Figure 3). Such behaviour for UUG has been depicted, but not discussed, in
 462 other analyses of CUPrefs in mammalian genes (see figure 7 in Laurin-Lemay et al. (2018)), in coronavirus genomes
 463 (Daron and Bravo, 2021), in plants (Clément et al., 2017) as well as for AGG-Arg and GGG-Gly in a global study of
 464 codon usages across the tree of life (see figure 1 in (Novoa et al., 2019)). The reasons underlying the clustering of
 465 UUG with AT-ending codons are unclear. A first line of thought could be functional: the UUG-Leu codon is particular
 466 because it can serve as alternative starting point for translation (Peabody, 1989). However, other codons such as ACG
 467 or GUG act more efficient than UUG as alternative translation initiation, and do not display any noticeable deviation in
 468 our results (Ivanov et al., 2011). A second line of thought could be related to the tRNA repertoire, but both UUG and
 469 UUA are decoded by similar numbers of dedicated tRNAs in the vast majority of genomes (e.g. respectively six and
 470 seven tRNA genes in humans (Palidwor et al., 2010)). Finally, another line of thought suggests that UUG and AGG
 471 could be disfavoured if substitution pressure towards GC is very high, despite being GC-ending codons (Palidwor et al.,
 472 2010). Indeed, the series of synonymous transitions UUA->UUG->CUG for Leucine and the substitution chain AGA-
 473 >AGG->CGG for Arginine are expected to lead to a depletion of UUG and of AGG codons when increasing GC
 474 content. Both UUG and ACG codons would this way display a non-monotonic response to GC-substitution biases
 475 (Palidwor et al., 2010). In our data-set, however, AGG maps with the rest of GC-ending codons, symmetrically op-
 476 posed to AGA as expected, and strongly contributing to the second PCA axis. Thus, only UUG displays frequency
 477 patterns similar to those of AT-ending codons. We humbly admit that we do not find a satisfactory explanation for this
 478 behaviour and invite researchers in the field to generate and test alternative explanatory hypotheses.

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

487 References

- 488 Abugessaisa I, Ramilowski JA, Lizio M, Severin J, Hasegawa A, Harshbarger J, Kondo A, Noguchi S, Yip CW, Ooi J,
 489 Tagami M, Hori F, Agrawal S, Hon C, Cardon M, Ikeda S, Ono H, Bono H, Kato M, Hashimoto K, Bonetti A, Kato
 490 M, Kobayashi N, Shin J, de Hoon M, Hayashizaki Y, Carninci P, Kawaji H, Kasukawa T. 2021, January. FANTOM
 491 enters 20th year: expansion of transcriptomic atlases and functional annotation of non-coding RNAs. *Nucleic Acids*
 492 *Research*. 49(D1):D892–D898.
- 493 Adams KL, Cronn R, Percifield R, Wendel JF. 2003, April. Genes duplicated by polyploidy show unequal contributions
 494 to the transcriptome and organ-specific reciprocal silencing. *Proceedings of the National Academy of Sciences of*
 495 *the United States of America*. 100(8):4649–4654.

- 496 Akashi H. 1997. Codon bias evolution in drosophila. population genetics of mutation-selection drift. *Gene*. 205.
- 497 Apostolou-Karampelis K, Nikolaou C, Almirantis Y. 2016, August. A novel skew analysis reveals substitution asym-
498 metries linked to genetic code GC-biases and PolIII a-subunit isoforms. *DNA research: an international journal for*
499 *rapid publication of reports on genes and genomes*. 23(4):353–363.
- 500 Barbosa-Morais NL, Irimia M, Pan Q, Xiong HY, Gueroussov S, Lee LJ, Slobodeniuc V, Kutter C, Watt S, Colak
501 R, Kim T, Misquitta-Ali CM, Wilson MD, Kim PM, Odom DT, Frey BJ, Blencowe BJ. 2012, December. The
502 evolutionary landscape of alternative splicing in vertebrate species. *Science (New York, N.Y.)*. 338(6114):1587–
503 1593.
- 504 Benisty H, Hernandez-Alias X, Weber M, Anglada-Girotto M, Mantica F, Radusky L, Senger G, Calvet F, Weghorn
505 D, Irimia M, Schaefer M, Serrano L. 2023. Genes enriched in a/t-ending codons are co-regulated and conserved
506 across mammals. *Cell Systems*. 14.
- 507 Bourret J, Alizon S, Bravo IG. 2019, December. COUSIN (COdon Usage Similarity INdex): A Normalized Measure
508 of Codon Usage Preferences. *Genome Biology and Evolution*. 11(12):3523–3528. Publisher: Oxford Academic.
- 509 Brawand D, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A, Aximu-Petri A, Kircher
510 M, Albert FW, Zeller U, Khaitovich P, Grützner F, Bergmann S, Nielsen R, Pääbo S, Kaessmann H. 2011, October.
511 The evolution of gene expression levels in mammalian organs. *Nature*. 478(7369):343–348.
- 512 Brule CE, Grayhack EJ. 2017. Synonymous Codons: Choose Wisely for Expression. *Trends in genetics: TIG*.
513 33(4):283–297.
- 514 Bulmer M. 1991, November. The selection-mutation-drift theory of synonymous codon usage. *Genetics*. 129(3):897–
515 907.
- 516 Caliskan N, Peske F, Rodnina MV. 2015, May. Changed in translation: mRNA recoding by 1 programmed ribosomal
517 frameshifting. *Trends in Biochemical Sciences*. 40(5):265–274.
- 518 Callens M, Pradier L, Finnegan M, Rose C, Bedhomme S. 2021. Read between the lines: Diversity of nontranslational
519 selection pressures on local codon usage. *Genome Biology and Evolution*. 13.
- 520 Carbone A, Zinovyev A, Képès F. 2003, November. Codon adaptation index as a measure of dominating codon bias.
521 *Bioinformatics (Oxford, England)*. 19(16):2005–2015.
- 522 Caspersson T, Farber S, Foley GE, Kudynowski J, Modest EJ, Simonsson E, Wagh U, Zech L. 1968, January. Chemical
523 differentiation along metaphase chromosomes. *Experimental Cell Research*. 49(1):219–222.
- 524 Castresana J. 2000, April. Selection of conserved blocks from multiple alignments for their use in phylogenetic
525 analysis. *Molecular Biology and Evolution*. 17(4):540–552.
- 526 Chamary JV, Parmley JL, Hurst LD. 2006, February. Hearing silence: non-neutral evolution at synonymous sites in
527 mammals. *Nature Reviews. Genetics*. 7(2):98–108.
- 528 Clark JM. 1988, October. Novel non-templated nucleotide addition reactions catalyzed by procaryotic and eucaryotic
529 DNA polymerases. *Nucleic Acids Research*. 16(20):9677–9686.
- 530 Clarke B. 1970. Darwinian evolution of proteins. *Science*. 168.

- 531 Clément Y, Sarah G, Holtz Y, Homa F, Pointet S, Contreras S, Nabholz B, Sabot F, Sauné L, Ardisson M, Bacilieri
532 R, Besnard G, Berger C Angélique Cardi, De Bellis F, Fouet O, Jourda C, Khadari B, Lanaud C, Leroy T, Pot D,
533 Sauvage C, Scarcelli N, Tregear J, Vigouroux Y, Yahiaoui N, Ruiz M, Santoni S, Labouisse JP, Pham JL, David J,
534 Glémin S. 2017. Evolutionary forces affecting synonymous variations in plant genomes. *PLoS Genetics*. 13:1–28.
- 535 Copley SD. 2020, April. Evolution of new enzymes by gene duplication and divergence. *The FEBS journal*.
536 287(7):1262–1283.
- 537 Daron J, Bravo IG. 2021. Variability in codon usage in coronaviruses is mainly driven by mutational bias and selective
538 constraints on cpg dinucleotide. *Viruses*. 13:1800.
- 539 Donizetti A, Fiengo M, Minucci S, Aniello F. 2009, October. Duplicated zebrafish relaxin-3 gene shows a different
540 expression pattern from that of the co-orthologue gene. *Development, Growth & Differentiation*. 51(8):715–722.
- 541 Duret L. 2002, December. Evolution of synonymous codon usage in metazoans. *Current Opinion in Genetics &*
542 *Development*. 12(6):640–649.
- 543 Duret L, Mouchiroud D. 1999, April. Expression pattern and, surprisingly, gene length shape codon usage in
544 *Caenorhabditis*, *Drosophila*, and *Arabidopsis*. *Proceedings of the National Academy of Sciences*. 96(8):4482–4487.
545 Publisher: National Academy of Sciences Section: Biological Sciences.
- 546 Ferris SD, Whitt GS. 1979, April. Evolution of the differential regulation of duplicate genes after polyploidization.
547 *Journal of Molecular Evolution*. 12(4):267–317.
- 548 Freilich S, Massingham T, Blanc E, Goldovsky L, Thornton JM. 2006. Relating tissue specialization to the differenti-
549 ation of expression of singleton and duplicate mouse proteins. *Genome Biology*. 7(10):R89.
- 550 Fu J, Dang Y, Counter C, Liu Y. 2018. Codon usage regulates human KRAS expression at both transcriptional and
551 translational levels. *The Journal of Biological Chemistry*. 293(46):17929–17940.
- 552 Galtier N, Roux C, Rousselle M, Romiguier J, Figuet E, Glémin S, Bierne N, Duret L. 2018, May. Codon Usage
553 Bias in Animals: Disentangling the Effects of Natural Selection, Effective Population Size, and GC-Biased Gene
554 Conversion. *Molecular Biology and Evolution*. 35(5):1092–1103.
- 555 Grantham R, Gautier C, Gouy M, Mercier R, Pavé A. 1980, January. Codon catalog usage and the genome hypothesis.
556 *Nucleic Acids Research*. 8(1):r49–r62.
- 557 Guschanski K, Warnefors M, Kaessmann H. 2017. The evolution of duplicate gene expression in mammalian organs.
558 *Genome Research*. 27(9):1461–1474.
- 559 Hannigan MM, Zagore LL, Licatalosi DD. 2017, June. Ptbp2 controls an alternative splicing network required for cell
560 communication during spermatogenesis. *Cell reports*. 19(12):2598–2612.
- 561 Hershberg R, Petrov DA. 2009, July. General rules for optimal codon choice. *PLoS genetics*. 5(7):e1000556.
- 562 Holmquist GP. 1989, June. Evolution of chromosome bands: Molecular ecology of noncoding DNA. *Journal of*
563 *Molecular Evolution*. 28(6):469–486.
- 564 Ikemura T. 1981, September. Correlation between the abundance of *Escherichia coli* transfer RNAs and the occurrence
565 of the respective codons in its protein genes: a proposal for a synonymous codon choice that is optimal for the *E.*
566 *coli* translational system. *Journal of Molecular Biology*. 151(3):389–409.

- 567 Ivanov IP, Firth AE, Michel AM, Atkins JF, Baranov PV. 2011, May. Identification of evolutionarily conserved non-
568 AUG-initiated N-terminal extensions in human coding sequences. *Nucleic Acids Research*. 39(10):4220–4234.
- 569 Katoh K, Misawa K, Kuma Ki, Miyata T. 2002, July. MAFFT: a novel method for rapid multiple sequence alignment
570 based on fast Fourier transform. *Nucleic Acids Research*. 30(14):3059–3066.
- 571 Keppetipola N, Sharma S, Li Q, Black DL. 2012, August. Neuronal regulation of pre-mRNA splicing by polypyrim-
572 idine tract binding proteins, PTBP1 and PTBP2. *Critical Reviews in Biochemistry and Molecular Biology*.
573 47(4):360–378.
- 574 Khorana HG, Büchi H, Ghosh H, Gupta N, Jacob TM, Kössel H, Morgan R, Narang SA, Ohtsuka E, Wells RD. 1966.
575 Polynucleotide synthesis and the genetic code. *Cold Spring Harbor Symposia on Quantitative Biology*. 31:39–49.
- 576 King JL, Jukes TH. 1969. Non-darwinian evolution. *Science*. 164.
- 577 Koonin EV. 2005. Orthologs, Paralogs, and Evolutionary Genomics. *Annual Review of Genetics*. 39(1):309–338.
578 _eprint: <https://doi.org/10.1146/annurev.genet.39.073003.114725>.
- 579 Kumar S, Stecher G, Suleski M, Hedges SB. 2017. TimeTree: A Resource for Timelines, Timetrees, and Divergence
580 Times. *Molecular Biology and Evolution*. 34(7):1812–1819.
- 581 Lampson B, Pershing N, Prinz J, Lacsina J, Marzluff W, Nicchitta C, MacAlpine D, Counter C. Rare codons regulate
582 *kras* oncogenesis. *Current biology*. 23.
- 583 Lampson BL, Pershing NLK, Prinz JA, Lacsina JR, Marzluff WF, Nicchitta CV, MacAlpine DM, Counter CM. 2013,
584 January. Rare codons regulate *KRas* oncogenesis. *Current biology: CB*. 23(1):70–75.
- 585 Laurin-Lemay S, Rodrigue N, Lartillot N, Philippe H. 2018. Conditional Approximate Bayesian Computation: A New
586 Approach for Across-Site Dependency in High-Dimensional Mutation-Selection Models. *Molecular Biology and*
587 *Evolution*. 35(11):2819–2834.
- 588 Le S, Gascuel O. 2008. An improved general amino acid replacement matrix. *Molecular Biology and Evolution*. 25.
- 589 Lujan SA, Williams JS, Pursell ZF, Abdulovic-Cui AA, Clark AB, McElhinny SAN, Kunkel TA. 2012, October. Mis-
590 match Repair Balances Leading and Lagging Strand DNA Replication Fidelity. *PLOS Genetics*. 8(10):e1003016.
591 Publisher: Public Library of Science.
- 592 Makeyev EV, Zhang J, Carrasco MA, Maniatis T. 2007, August. The MicroRNA miR-124 Promotes Neuronal Differ-
593 entiation by Triggering Brain-Specific Alternative Pre-mRNA Splicing. *Molecular cell*. 27(3):435–448.
- 594 Marais G. 2003, June. Biased gene conversion: implications for genome and sex evolution. *Trends in Genetics*.
595 19(6):330–338. Publisher: Elsevier.
- 596 Merkin J, Russell C, Chen P, Burge CB. 2012, December. Evolutionary dynamics of gene and isoform regulation in
597 Mammalian tissues. *Science (New York, N.Y.)*. 338(6114):1593–1599.
- 598 Meyer A, Schartl M. 1999. Gene and genome duplications in vertebrates: the one-to-four (-to-eight in fish) rule and
599 the evolution of novel gene functions. *Current Opinion in Cell Biology*. 11.
- 600 Mordstein C, Savisaar R, Young RS, Bazile J, Talmane L, Luft J, Liss M, Taylor MS, Hurst LD, Kudla G. 2020, April.
601 Codon Usage and Splicing Jointly Influence mRNA Localization. *Cell Systems*. 10(4):351–362.e8.

- 602 Munk M, Villalobo E, Villalobo A, Berchtold MW. 2022, November. Differential expression of the three independent
603 cam genes coding for an identical protein: Potential relevance of distinct mrna stability by different codon usage.
604 Cell Calcium. 107.
- 605 NCBI Resource Coordinators. 2018. Database resources of the National Center for Biotechnology Information. Nu-
606 cleic Acids Research. 46(D1):D8–D13.
- 607 Newman ZR, Young JM, Ingolia NT, Barton GM. 2016, March. Differences in codon bias and GC content contribute
608 to the balanced expression of TLR7 and TLR9. Proceedings of the National Academy of Sciences of the United
609 States of America. 113(10):E1362–1371.
- 610 Nirenberg MW, Matthaei JH. 1961, October. The dependence of cell- free protein synthesis in e. coli upon naturally
611 occurring or synthetic polyribonucleotides. Proceedings of the National Academy of Sciences of the United States
612 of America. 47(10):1588–1602.
- 613 Novoa EM, Jungreis I, Jaillon O, Kellis M. 2019. Elucidation of Codon Usage Signatures across the Domains of Life.
614 Molecular Biology and Evolution. 36(10):2328–2339.
- 615 Novoa EM, Ribas de Pouplana L. 2012, November. Speeding with control: codon usage, tRNAs, and ribosomes.
616 Trends in genetics: TIG. 28(11):574–581.
- 617 Palidwor GA, Perkins TJ, Xia X. 2010, October. A general model of codon bias due to GC mutational bias. PloS One.
618 5(10):e13431.
- 619 Peabody DS. 1989, March. Translation initiation at non-AUG triplets in mammalian cells. The Journal of Biological
620 Chemistry. 264(9):5031–5035.
- 621 Percudani R, Pavesi A, Ottonello S. 1997, May. Transfer RNA gene redundancy and translational selection in Saccha-
622 romyces cerevisiae11Edited by J. Karn. Journal of Molecular Biology. 268(2):322–330.
- 623 Picard M, Leblay F, Cassan C, Willemsen A, Daron J, Bauffe F, Decourcelle M, Demange A, Bravo I. 2023. Tran-
624 scriptomic, proteomic, and functional consequences of codon usage bias in human cells during heterologous gene
625 expression. Protein Science. 32.
- 626 Pina J, Ontiveros RJ, Keppetipola N, Nikolaidis N. 2018, April. A Bioinformatics Approach to Discover the Evolu-
627 tionary Origin of the PTBP Splicing Regulators. The FASEB Journal. 32(1_supplement):802.16–802.16. Publisher:
628 Federation of American Societies for Experimental Biology.
- 629 Plotkin JB, Kudla G. 2011, January. Synonymous but not the same: the causes and consequences of codon bias. Nature
630 Reviews Genetics. 12(1):32–42.
- 631 Pouyet F, Mouchiroud D, Duret L, Sémon M. 2017. Recombination, meiotic expression and human codon usage.
632 eLife. 6.
- 633 Presnyak V, Alhusaini N, Chen YH, Martin S, Morris N, Kline N, Olson S, Weinberg D, Baker KE, Graveley BR,
634 Collier J. 2015, March. Codon optimality is a major determinant of mRNA stability. Cell. 160(6):1111–1124.
- 635 Reijns MAM, Kemp H, Ding J, Marion de Procé S, Jackson AP, Taylor MS. 2015, February. Lagging-strand replication
636 shapes the mutational landscape of the genome. Nature. 518(7540):502–506. Number: 7540 Publisher: Nature
637 Publishing Group.

- 638 Robinson DF, Foulds LR. 1981, February. Comparison of phylogenetic trees. *Mathematical Biosciences*. 53(1):131–
639 147.
- 640 Robinson F, Jackson RJ, Smith CWJ. 2008, March. Expression of Human nPTB Is Limited by Extreme Suboptimal
641 Codon Content. *PLOS ONE*. 3(3):e1801. Publisher: Public Library of Science.
- 642 Satapathy SS, Powdel BR, Buragohain AK, Ray SK. 2016, October. Discrepancy among the synonymous codons
643 with respect to their selection as optimal codon in bacteria. *DNA Research*. 23(5):441–449. Publisher: Oxford
644 Academic.
- 645 Scornavacca C, Belkhir K, Lopez J, Dernat R, Delsuc F, Douzery EJP, Ranwez V. 2019, April. OrthoMaM v10:
646 Scaling-Up Orthologous Coding Sequence and Exon Alignments with More than One Hundred Mammalian
647 Genomes. *Molecular Biology and Evolution*. 36(4):861–862. Publisher: Oxford Academic.
- 648 Sharp PM, Li WH. 1987. The codon Adaptation Index—a measure of directional synonymous codon usage bias, and
649 its potential applications. *Nucleic Acids Research*. 15(3):1281–1295.
- 650 Sonnhammer ELL, Koonin EV. 2002, December. Orthology, paralogy and proposed classification for paralog subtypes.
651 *Trends in genetics: TIG*. 18(12):619–620.
- 652 Soria-Carrasco V, Talavera G, Igea J, Castresana J. 2007, November. The K tree score: quantification of differences
653 in the relative branch length and topology of phylogenetic trees. *Bioinformatics (Oxford, England)*. 23(21):2954–
654 2956.
- 655 Spellman R, Llorian M, Smith CW. 2007, August. Crossregulation and functional redundancy between the splicing
656 regulator ptb and its paralogs nptb and rod1. *Molecular Cell*. 27:420–434.
- 657 Spencer PS, Barral JM. 2012, March. Genetic code redundancy and its influence on the encoded polypeptides. *Com-
658 putational and Structural Biotechnology Journal*. 1.
- 659 Stamatakis A. 2014, May. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies.
660 *Bioinformatics (Oxford, England)*. 30(9):1312–1313.
- 661 Vuong JK, Lin CH, Zhang M, Chen L, Black DL, Zheng S. 2016. PTBP1 and PTBP2 Serve Both Specific and
662 Redundant Functions in Neuronal Pre-mRNA Splicing. *Cell Reports*. 17(10):2766–2775.
- 663 Waddell PJ, Steel M. 1997. General time-reversible distances with unequal rates across sites: Mixing and inverse
664 gaussian distributions with invariant sites. *Molecular Phylogenetics and Evolution*. 8(3):398–414.
- 665 Wagner EJ, Garcia-Blanco MA. 2002, October. Rnai-mediated ptb depletion leads to enhanced exon definition. *Molec-
666 ular Cell*. 10:943–949.
- 667 Whittle CA, Extavour CG. 2016, September. Expression-Linked Patterns of Codon Usage, Amino Acid Frequency,
668 and Protein Length in the Basally Branching Arthropod Parasteatoda tepidariorum. *Genome Biology and Evolution*.
669 8(9):2722–2736. Publisher: Oxford Academic.
- 670 Yates AD, Achuthan P, Akanni W, Allen J, Allen J, Alvarez-Jarreta J, Amode MR, Armean IM, Azov AG, Bennett
671 R, Bhai J, Billis K, Boddu S, Marugán JC, Cummins C, Davidson C, Dodiya K, Fatima R, Gall A, Giron CG, Gil
672 L, Grego T, Haggerty L, Haskell E, Hourlier T, Izuogu OG, Janacek SH, Juettemann T, Kay M, Lavidas I, Le T,
673 Lemos D, Martinez JG, Maurel T, McDowall M, McMahon A, Mohanan S, Moore B, Nuhn M, Oheh DN, Parker

- 674 A, Parton A, Patricio M, Sakthivel MP, Abdul Salam AI, Schmitt BM, Schuilenburg H, Sheppard D, Sycheva M,
675 Szuba M, Taylor K, Thormann A, Threadgold G, Vullo A, Walts B, Winterbottom A, Zadissa A, Chakiachvili M,
676 Flint B, Frankish A, Hunt SE, IIsley G, Kostadima M, Langridge N, Loveland JE, Martin FJ, Morales J, Mudge
677 JM, Muffato M, Perry E, Ruffier M, Trevanion SJ, Cunningham F, Howe KL, Zerbino DR, Flicek P. 2020, January.
678 Ensembl 2020. *Nucleic Acids Research*. 48(D1):D682–D688. Publisher: Oxford Academic.
- 679 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
680 WD, Moland CL, Branham WS, Qian F, Ning B, Li Y, Hong H, Guo L, Mei N, Shi T, Wang KY, Wolfinger RD,
681 Nikolsky Y, Walker SJ, Duerksen-Hughes P, Mason CE, Tong W, Thierry-Mieg J, Thierry-Mieg D, Shi L, Wang C.
682 2014, February. A rat RNA-Seq transcriptomic BodyMap across 11 organs and 4 developmental stages. *Nature*
683 *Communications*. 5(1):3230. Number: 1 Publisher: Nature Publishing Group.
- 684 Zagore LL, Grabinski SE, Sweet TJ, Hannigan MM, Sramkoski RM, Li Q, Licatalosi DD. 2015, December. RNA
685 Binding Protein Ptbp2 Is Essential for Male Germ Cell Development. *Molecular and Cellular Biology*. 35(23):4030–
686 4042.