

1 Evolutionary dynamics of circular 2 RNAs in primates

3 **Gabriela Santos-Rodriguez**^{1,2}, **Irina Voineagu**³, **Robert J Weatheritt**^{1,2*}

*For correspondence:

r.weatheritt@garvan.org.au (RW)

4 ¹EMBL Australia, Garvan Institute of Medical Research, Darlinghurst, NSW, 2010,
5 Australia.; ²St. Vincent Clinical School, University of New South Wales, Darlinghurst,
6 NSW, 2010, Australia.; ³School of Biotechnology and Biomolecular Sciences, University
7 of New South Wales, Sydney, Australia.

8
9 **Abstract** Many primate genes produce non-coding circular RNAs (circRNAs). However, the
10 extent of circRNA conservation between closely related species remains unclear. By comparing
11 tissue-specific transcriptomes across over 70 million years of primate evolution, we identify that
12 within 3 million years circRNA expression profiles diverged such that they are more related to
13 species identity than organ type. However, our analysis also revealed a subset of circRNAs with
14 conserved neural expression across tens of millions of years of evolution. These circRNAs are
15 defined by an extended downstream intron that has shown dramatic lengthening during
16 evolution due to the insertion of novel retrotransposons. Our work provides comparative
17 analyses of the mechanisms promoting circRNAs to generate increased transcriptomic
18 complexity in primates.

20 Introduction

21 An important question in biology is how has the complexity of biological systems expanded while
22 the number of protein-coding genes has remained mostly stable. Through decades of research, it
23 has been shown that increased biological complexity has arisen in part by the dynamic generation
24 of unique cell-specific transcriptomes, and as a consequence of the highly versatile programs of

25 gene expression (*Brawand et al., 2011; Cardoso-Moreira et al., 2019*). However, studies of tissues
26 across distant animal lineages have shown that gene expression is highly conserved between the
27 same tissues in different species (*Brawand et al., 2011; Barbosa-Morais et al., 2012; Merkin et al.,*
28 *2012; Reyes et al., 2013; Cardoso-Moreira et al., 2019*). Hence, gene expression alone is unlikely
29 to explain the heterogeneous expansion in complexity (as defined by the number of cell types)
30 across vertebrate evolution. Instead, it is becoming increasingly evident that the plethora of post-
31 transcriptional mechanisms (*Gueroussov et al., 2017; Ha et al., 2018; Mattick, 2018; Fiszbein et al.,*
32 *2019; Cheetham et al., 2020; Ha et al., 2021*) capable of greatly expanding transcriptomic diversity
33 also underlies these advances.

34 Among these, an intriguing class produced by pre-mRNA processing are circular RNAs (circR-
35 NAs) (*Memczak et al., 2013; Zhang et al., 2013; Li et al., 2018b; Gokool et al., 2020b*). These non-
36 coding RNAs can regulate protein localization (*Liu et al., 2019*), miRNA functionality (*Piwecka et al.,*
37 *2017*) and a range of other processes (*Li et al., 2018b; Gokool et al., 2020b*) enabling increased reg-
38 ulatory complexity, especially in the immune and nervous systems (*Piwecka et al., 2017; Gokool*
39 *et al., 2020a; Guo et al., 2020*). CircRNAs form by back-splicing whereby an exon's 3'-splice site
40 is ligated to an upstream 5'-splice site forming a closed circular non-coding RNA transcript (*Bar-*
41 *rett et al., 2015; Starke et al., 2015*). Back-splicing occurs co-transcriptionally and is facilitated by
42 inverted repeat elements that promote complementarity between adjacent introns favouring cir-
43 cRNA formation over linear splicing (*Jeck et al., 2013; Liang and Wilusz, 2014; Ivanov et al., 2015*).
44 These RNA-RNA interactions can be facilitated by RNA-binding proteins, such as Quaking (*Conn*
45 *et al., 2015*), that help stabilize the hair-pin structure promoting circRNA formation.

46 The production of circRNAs can also arise due to the perturbed expression of trans-factors and
47 the inhibition of the core splicing machinery (*Aktaş et al., 2017; Liang et al., 2017*). These spuriously
48 produced circRNAs are maintained as their circular shape protects them from the activity of cellu-
49 lar exonucleases (*Gokool et al., 2020b*). In contrast, the variable usage of cis-regulatory elements in
50 exons and flanking introns can be selected to promote circRNA expression in a cell-type, condition-
51 or species-specific manner (*Nilsen and Graveley, 2010; Irimia and Blencowe, 2012*). Changes in
52 circRNA expression may therefore represent a major source of species- and lineage-specific dif-
53 ferences or error-prone mis-splicing. To provide insight into this quandary, here we describe a
54 genome-wide analysis of circRNAs across physiologically equivalent organs from primate species
55 spanning 70 million years of evolution. Our analysis uncovers extensive evidence species-specific

56 circRNAs that display no evidence of conservation even across relatively short evolutionary time-
57 periods. However, we also identify a small subset of circRNAs that are conserved across tens of
58 millions of years displaying increased inclusion rates across evolutionary time. Our analysis reveals
59 that these circRNAs are flanked by newly inserted transposons that correlate with circRNA genesis
60 and extend intron downstream of circRNA. Overall, our results identify evidence of circRNA con-
61 servation within closely related species and identify a reoccurring mechanism that correlates with
62 circRNA genesis facilitating the expansion of transcriptomic complexity of primate cells.

63 Results

64 **A core subset of circRNAs show conserved expression signatures but most are species-** 65 **specific**

66 To address the outstanding questions about the conservation and functional importance of cir-
67 cRNAs, we collected transcriptomic (RNA-seq) data (*Pipes et al., 2013*) from across 9 tissues from
68 8 primate species, consisting of 3 old-world monkeys, 2 hominoids, 2 new-world monkeys, and
69 one prosimian (**Supplementary Table 1**). These species were chosen on the basis of the quality
70 of their genomes and their close evolutionary relationships enabling the evaluation of transcrip-
71 tome changes between species ranging from <3 million years to > 70 million years (see **Figure**
72 **1A**). For each species, we considered all primate-conserved internal exons as potential origins of
73 back-spliced junctions with no restrictions on backward exon combination. RNA-seq reads were
74 mapped to exon-exon junctions (EEJs) to determine “percent spliced-in” (PSI) for all circRNA with
75 respect to the linear transcript. We also calculated PSI values for linear splicing of each internal
76 exon and transcript per million (TPM) values to estimate gene expression. Orthology relationships
77 between genes and exons were established to enable direct cross-species comparisons.

78 To initially explore the expression relationships within our datasets we used hierarchical clus-
79 tering and Pearson’s correlations to determine the gene expression relationships between orthol-
80 ogous genes (see Methods). In agreement with previous results (*Brawand et al., 2011; Barbosa-*
81 *Morais et al., 2012; Merkin et al., 2012; Reyes et al., 2013*) from analysis across vertebrate species,
82 a clear pattern emerged of tissue-specific conservation of gene expression (**Figure 1B**). This pat-
83 tern suggests that most tissues possess a tissue-specific gene expression signature such that for
84 example a liver-specific gene in chimp will likely also be liver-specific in lemur. In contrast to previ-
85 ous observations in vertebrates (*Merkin et al., 2012*) there are no clear species-specific exceptions

86 to these patterns likely reflecting the closer evolutionary relationships studied.

87 To understand circRNA relationships between species, we performed an analogous pairwise
88 clustering analysis using circRNA inclusion values. Replicates from the same tissue invariably clus-
89 tered together. However, in contrast to gene expression, circRNA expression is segregated by
90 species (**Figure 1-Figure supplement 1A**). This suggests that despite all the exons studied being
91 conserved across primates the majority of circRNAs showed species-specific expression with no or-
92 thologous circRNAs in other species (**Figure 1C**, 67% are species-specific, $n = 11,201$). To evaluate
93 the expression patterns of circRNA orthologs, we identified circRNAs with matched back-spliced
94 junctions (see Methods) conserved across 45 million years of evolution. In this analysis more com-
95 plex patterns of circRNA conservation emerged with tissue-dominated clustering observed across
96 all types of brain samples (**Figure 1D**). In contrast, for all other tissues circRNAs showed primarily
97 species-specific clustering. Analysis of gene expression changes of genes with these conserved
98 circRNAs (**Figure 1-Figure supplement 1B**,) and alternative splicing changes of exons within con-
99 served circRNAs (**Figure 1-Figure supplement 1C**) showed no consistent changes suggesting cir-
100 cRNA conservation and expression is independent of these regulatory layers.

101 We next investigated the genes containing circRNAs. Many orthologous genes consistently ex-
102 press circRNAs even if the precise back-spliced junction is not conserved implicating importance
103 of trans-factors in controlling circRNA formation (**Figure 1C**). This phenomenon persisted across
104 species with a median of 10 circRNAs detected per gene across tissues (**Figure 1-Figure supple-**
105 **ment 1D**). However, this circRNA production only occurred in a limited number of expressed genes
106 (20.4% of orthologous expressed genes). This suggests certain genomic areas are circRNA factories
107 that are prone to produce large numbers of lowly expressed circRNAs.

108 These observations suggest a core set of circRNAs show conserved tissue-specific patterns
109 across neural tissues. However, the great prevalence of circRNAs showing species-specific expres-
110 sion indicates that the cis-regulatory or trans-regulatory environments may differ between even
111 very closely related species to promote the species-specific production of circRNAs.

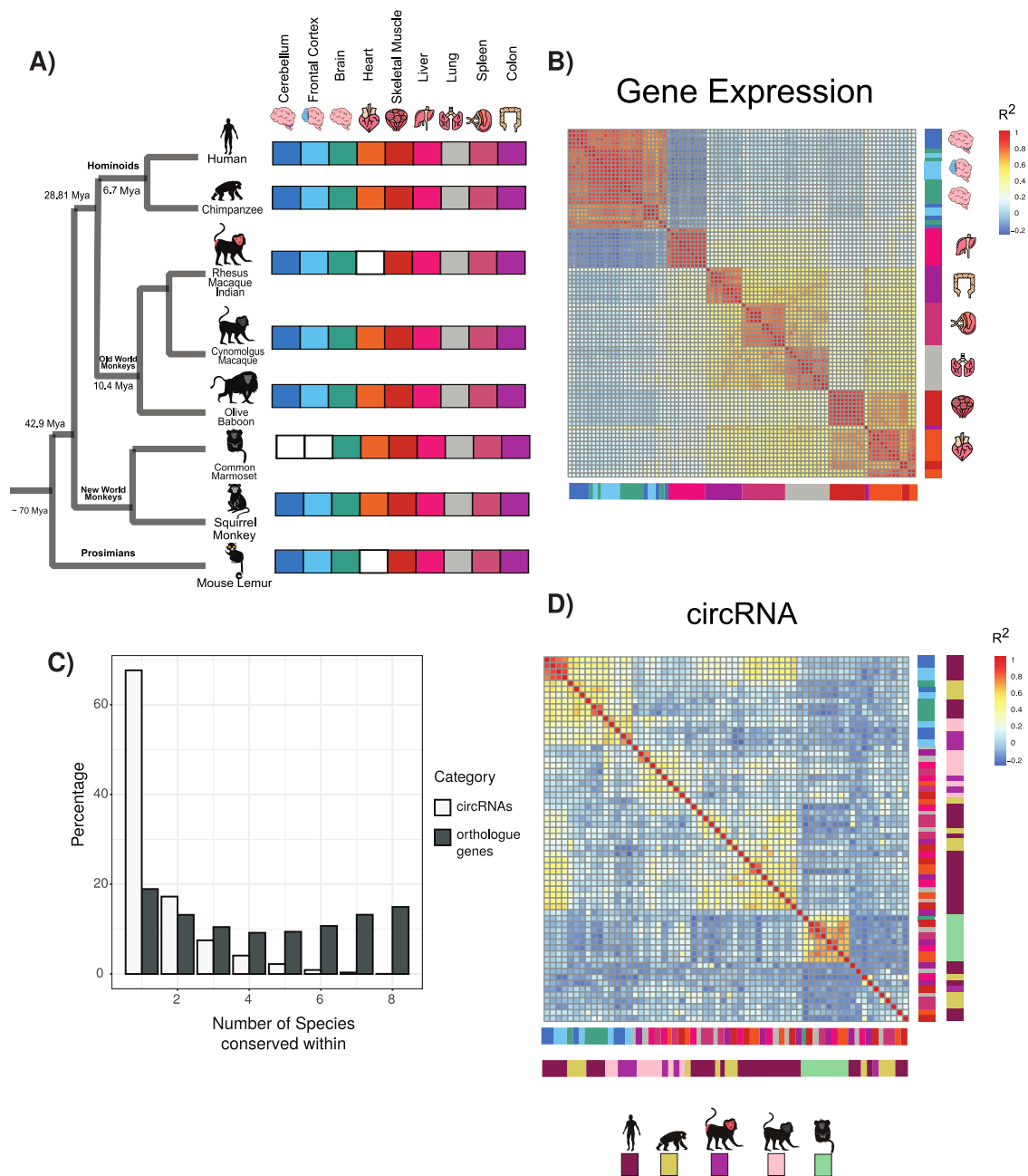


Figure 1. Circular RNA expression signatures are conserved in some tissues.

(A) Phylogenetic tree of analyzed species with distance from human in millions of years (Mya) (Divergence time according to TimeTree <http://www.timetree.org/>). Tissue datasets used in analysis identified on right with white squares denoting lack of dataset. (B) Clustering of samples based on expression values (transcripts per million). The variance of expression values was calculated, and the top 1000 most variable genes were used to calculate the Pearson correlation. ($n = 1000$ genes in 88 samples). Red colours indicate high correlation between samples and blue describes low correlation. Vertical and horizontal adjacent heatmaps describes tissues (see Fig. 1A for key) (C) Bar plot showing conservation of circRNAs based on back-spliced junction and based on occurrence within orthologous genes. (D) Clustering of conserved circRNAs based on percent spliced in (PSI) values. Clustered using Pearson correlation as in (B). ($n = 149$). Vertical and horizontal adjacent heatmaps describes tissues (inner heatmap (see Figure 1A for key)) and species (outer heatmap).

Figure 1-Figure supplement 1.

112 **Features of conserved circRNAs**

113 Our analysis (**Figure 2A**) reveals clear subsets of several hundred circRNAs exhibiting highly con-
114 served circRNA expression. The circRNA ERC1 and many other examples from our data (**Figure**
115 **2-Figure supplement 1A, Supplementary Table 2**) demonstrate that circRNA expression can be
116 conserved for tens of millions of years.

117 To assess the phylogenetic distribution of circRNA across primates we grouped them by PSI val-
118 ues requiring $PSI \geq 5$ and at least 5 read support. Out of the approximately 56,000 internal exons
119 with clear orthologs across primates, we identified a large set of circRNA expressing a “species-
120 specific” expression, as well as a set of 773 “conserved circRNAs” that shared expression across
121 at least human, chimp and baboon (**Figure 2-Figure supplement 1B and 1C**). Using our transcrip-
122 tomic data, we found that a circRNA identified in human was 5-times more likely to be identified in
123 baboon than in lemur, in line with the closer phylogenetic relationship of human to baboon than
124 human to lemur.

125 Initial analysis of conserved circRNAs revealed enrichment within neural tissues with over 70%
126 showing consistent tissue expression across 30 million years of evolution (**Supplementary Table**
127 **2**). Analysis of expression levels revealed no clear trends for increased expression of conserved
128 circRNAs (**Figure 2-Figure supplement 2A**, $p < 0.187$, Wilcoxon rank sum test vs species-specific),
129 however these circRNAs did display increased inclusion rates, or increased circRNA expression
130 as compared to linear isoform (**Figure 2-Figure supplement 2B**, $p = 3.38 \times 10^{-74}$ Wilcoxon rank
131 sum test vs species-specific). Furthermore, this inclusion (or circularization) increased with the
132 conservation age of the circRNA (**Figure 2E**, $p = 8.07 \times 10^{-19}$ Wilcoxon rank sum test of Hominoids
133 vs species-specific (Human specific); $p = 2.14 \times 10^{-06}$ Wilcoxon rank sum test of Hominoids vs shared
134 until New World Monkeys). This suggests over time these circRNAs are increasingly influencing the
135 transcriptomic abundance of the linear isoform and the protein abundance of the gene.

136 Analysis of the exonic structure of conserved circRNAs, showed that conserved circRNAs con-
137 tain fewer exons (**Figure 2F**, $p = 2.23 \times 10^{-20}$ Wilcoxon rank sum test), and rarely overlap with other
138 circRNAs (**Figure 2G**, $p = 4.08 \times 10^{-64}$, Fisher exact test; see Methods) displaying back-splicing at
139 unique 5' - and 3' -splice sites. This indicates that these conserved circRNAs possess unique cis- or
140 trans-regulatory features that enable a tight control of the number of exons within a circRNA and
141 the back-spliced junctions used.

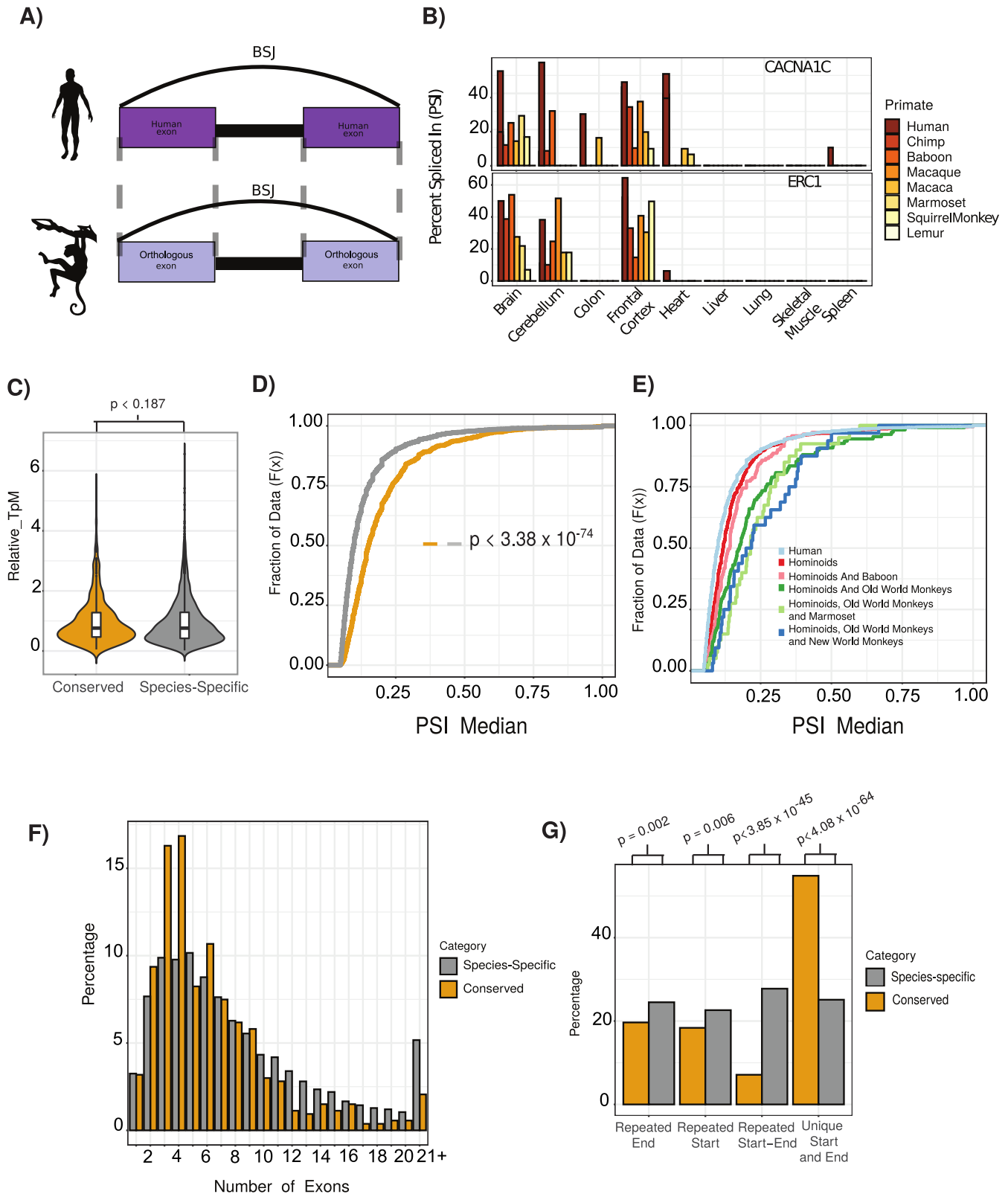


Figure 2 (previous page). Features of conserved circRNAs.

(A) Schematic overview of identification of back-spliced junctions between species. BSJ = back spliced junction. (B) Percent spliced in (PSI) values for conserved circRNAs (top) CACNA1C_chr12:2504436-2512984 and (bottom) ERC1_chr12:1180540-1204512 across tissues and species analyzed. PSI values only calculated for circRNAs with more than 5 reads support. Gene name is indicated in top right-hand corner. (C) Violin plot describing relative expression levels of conserved and species-specific circRNAs. Axes on violin plots are Relative TpM values, the normalized circRNA expression according to circRNA reads and gene expression (reads and TpMs values) (see methods). Violin plots show probability densities of the data with internal boxplot. The boxplot display the interquartile range as a solid box, 1.5 times the interquartile range as vertical thin lines and the median as a horizontal line. P-value calculated using Wilcoxon rank sum test ($p < 0.187$). TpM = transcripts per million (D) Cumulative distribution plot of change in percent spliced in (PSI) values across all conserved (yellow) and species-specific (grey) circRNAs. A cumulative distribution plot describes the proportion of data (y-axis) less than or equal to a specified value (x-axis). Cumulative Distribution F(x), cumulative distribution function. p value calculated using Wilcoxon-rank sum test ($p < 3.38 \times 10^{-74}$). PSI = percent spliced in (E) Cumulative distribution plots of circRNAs with different levels of conservation, as defined by consistent observation of back-spliced junction across species indicated. See 2D for description of cumulative distribution plot. PSI = percent spliced in (F) Bar plot describing number of exons per circRNA for conserved and species-specific circRNAs. Exons are defined by Ensembl and must show evidence of expression (PSI>5 and >5 reads support) in tissue analysed. (G) Bar plot describing uniqueness of start (5' splice site) and end (3' splice site) for conserved and species-specific circRNAs. P-values calculated from Fisher's exact test ($p < 4.08 \times 10^{-64}$; unique start and end – also see Figure2-Figure supplement 3).

Figure 2-Figure supplement 1.

Figure 2-Figure supplement 2.

Figure 2-Figure supplement 3.

142 **Conserved circRNAs have extensive downstream introns and are flanked by in-**
143 **verted repeat elements**

144 To investigate the role of cis-regulatory elements within conserved circRNAs, we analyzed almost
145 150 features associated with circRNA formation including a multitude of trans- and cis- regula-
146 tory factors and all major groups of transposons (see Methods and **Supplementary Table 3**). To
147 evaluate the influence of these features on defining conserved circRNAs we used two background
148 datasets (see **Supplementary Table 2** and Methods). The first is a background set of randomly
149 combined alternative ($10 < PSI < 90$) exons extracted from genes containing conserved circRNAs
150 (background set). The second is the group of “species-specific circRNAs” defined previously.

151 Using logistic regression combined with a genetic algorithm for model selection (see Methods),
152 we initially sought to determine the relative contribution of this diverse range of features in defin-
153 ing conserved circRNAs. After initially training our model on a subset of conserved and background
154 circRNAs (80%), we next assessed its performance on the rest of 20% circRNAs and observed a high

155 average true positive rate of 86.7% (AUC, area under the receiver operating characteristic (ROC)
156 curve) (**Figure 3-Supplementary Figure 1A**) for a model including 24 variables selected by feature
157 analysis. This indicates a core set of 24 cis- and trans-regulatory features drive the conserved for-
158 mation of circRNAs compared to our background set of introns (**Figure 3A and 3B**). We next used
159 the same approach to determine drivers of conserved and species-specific circRNAs. As expected,
160 our model distinguished these categories less efficiently but was still able to achieve a true positive
161 rate of 65.4% (**Figure 3-Figure supplement 1B**) driven by 12 features. Notable among these fea-
162 tures was the depletion of nucleosomes in the downstream intron of the circRNA (**Figure 3-Figure**
163 **supplement 1D** 1.57×10^{-03} , Bonferroni-corrected Wilcoxon rank sum test (BH-Wilcox) vs species-
164 specific) and the presence of a more defined 3' splice site at the final exon (2.04×10^{-03} , BH-Wilcox
165 vs species-specific).

166 Introns adjacent to conserved circRNAs also exhibited a significant enrichment for repeat ele-
167 ments (**Figure 3D**, all $p < 1 \times 10^{-5}$, BH-Wilcox) vs species-specific) in particular L1 and AluJ retro-
168 transposons (**Figure 3D**, L1: $p < 1.22 \times 10^{-23}$ | AluJ: $p < 1.48 \times 10^{-18}$, BH-Wilcox). A further key
169 distinguishing feature of interest was intron length. Conserved circRNAs exhibited shorter introns
170 downstream of the first exon and an extended intron downstream of the final exon (**Figure 4A**
171 **and 4B**). In species-specific circRNA this adjacent downstream intron has a median length of 4,624
172 nucleotides whilst in conserved circRNA the median is almost twice as long at 9,923 nucleotides
173 (**Figure 4B**, $p < 1.07 \times 10^{-35}$, BH-Wilcox). Finally, when comparing the major drivers of both models,
174 we noticed over 90% (11/12) of features overlapped between the models. This suggests conserved
175 circRNAs are an extreme continuum of species-specific circRNAs. Therefore understanding the pro-
176 cesses contributing to circRNA conservation may also provide insight into the genesis of circRNAs
177 across species.

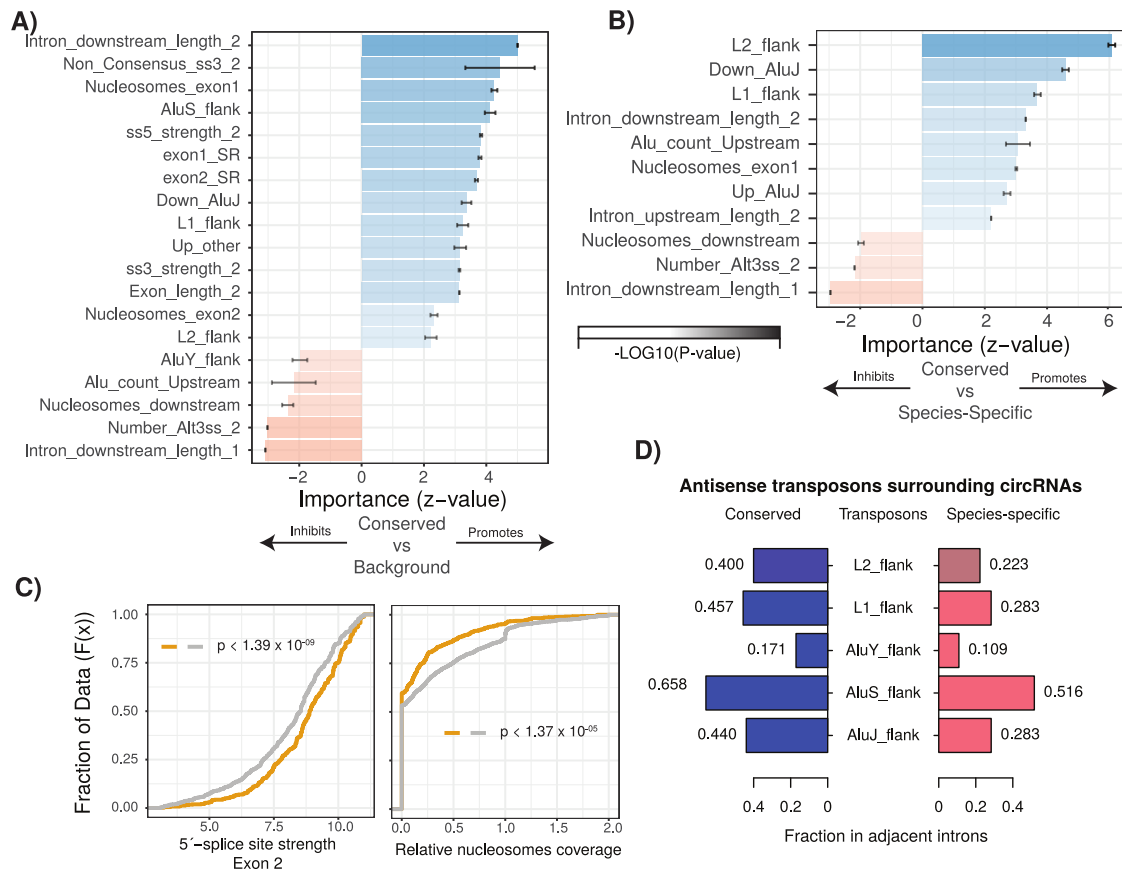


Figure 3. Characterization of cis and trans regulatory features of conserved circRNAs.

(A) Bar plot describing feature importance for logistic regression model of conserved circRNAs compared to background. Colours represent positive or negative influence. Transparency reflects $\log_{10}(p\text{-value of } z\text{-statistic})$. Errors bar represent standard error. “_1” is relative to first exon of circRNA and “_2” is relative to final exon of circRNA. ss3 = 3’ splice site; ss5 = 5’ splice site; Alt3ss = alternative 3’ splice sites. “Flank” are inverted repeats in introns adjacent to circRNAs. See Supplementary Table 3 for details of features. (B) Bar plot describing feature importance for logistic regression model of conserved circRNAs compared to species-specific circRNAs. See 3A for plot interpretation and descriptions. (C) Cumulative distribution plots describing (left; $p < 1.39 \times 10^{-09}$) 5’ splice site strength at final exon of circRNAs and (right; $p < 1.37 \times 10^{-05}$) distribution of nucleosomes on intron downstream of circRNA. p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). See Figure 2D for interpretation of cumulative distribution plot. (D) Pyramid plot showing the mean fraction of circRNAs with selected inverted repeat retrotransposon elements in adjacent introns.

Figure 3-Figure supplement 1.

178 **Insertion of young transposons increases downstream intron length in conserved**
179 **circRNAs**

180 To investigate the evolutionary origins of the switch of conserved circRNAs from absence in prosimi-
181 ans and new world monkeys to conservation within hominoids and old-world monkeys, we investi-
182 gated the changes in intronic length for the orthologous introns between human (hominoids) and
183 lemur (prosimians). In contrast to orthologous lemur introns, the human introns downstream of
184 all identified circRNAs shows an almost four-fold expansion compared to background dataset of
185 introns within circRNA containing genes (**Figure 4C**, $p < 3.84 \times 10^{-23}$ Wilcoxon rank sum). This dif-
186 ference is even greater in conserved circRNA, which display an almost 2-fold greater lengthening
187 than species-specific circRNAs (or 8-fold over background) (**Figure 4C**, $p < 3.84 \times 10^{-06}$, Wilcoxon
188 rank sum). These observations suggest that the expansion of the intron downstream of the cir-
189 cRNA may increase the proportion of backing splicing events increasing the likelihood of circRNA
190 conservation.

191 To investigate the drivers of this intronic expansion, we aligned the lemur and human introns
192 to identify regions novel to humans. This analysis revealed the insertion of novel transposons
193 at almost double the frequency in introns associated with conserved circRNAs (**Figure 4D**, $p <$
194 5.48×10^{-06} , Wilcoxon rank sum). Further evaluation of the retrotransposons revealed this increase
195 in length is driven by the novel insertion of AluJ and L1 elements (**Figure 4E**, AluJ: $p < 0.018$; L1: $p <$
196 1.73×10^{-04} , Wilcoxon rank sum). This retrotransposition is potentially facilitated by the depletion
197 of nucleosome occupancy in these introns compared to other human introns (**Figure 3B**, $p < 1.15 \times$
198 10^{-07} , BH-Wilcox). Together this argues for the role of young transposons in creating longer intronic
199 regions, which increases the time for RNA polymerase II to reach next canonical splice site and
200 therefore increases likelihood of back-junction splicing to occur.

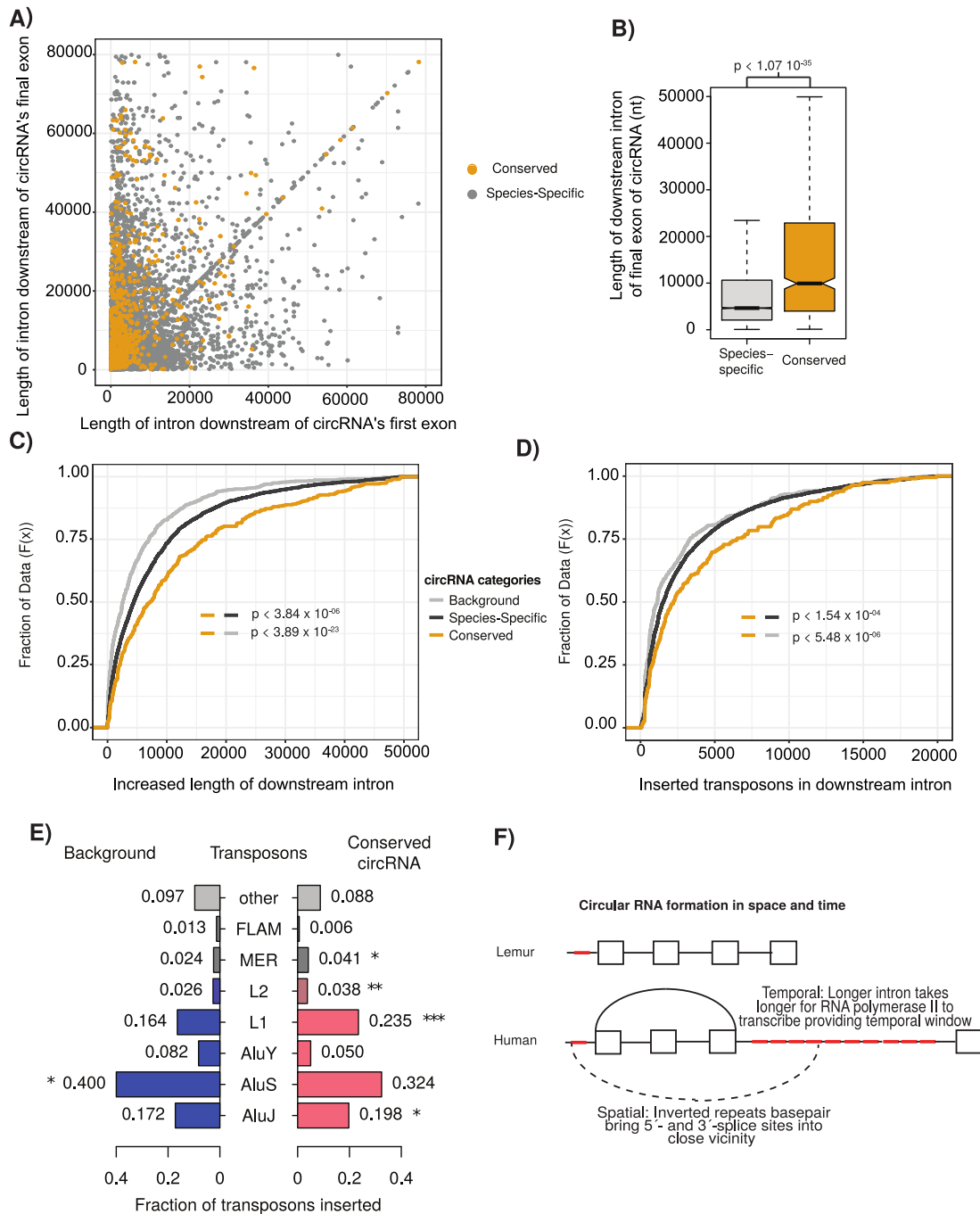


Figure 4 (previous page). Conserved circRNA downstream intron expanded during primate evolution.

A) Scatterplot of downstream intron length for conserved and species-specific circRNAs. (B) Boxplot describing lengths of intron immediately downstream of circRNA for conserved and species-specific circRNAs (see Figure 2C for description of boxplots). p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). nt = nucleotide (C) Cumulative distribution plot of change of length of orthologous downstream introns of conserved, species-specific and background circRNAs from lemur to human (see Fig. 2D for description of cumulative distribution plots). p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). (D) Cumulative distribution plot of length of novel repeat elements within the orthologous downstream introns of conserved, species-specific and background circRNAs from lemur to human (see Fig. 2D for description of cumulative distribution plots). p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). (E) Pyramid plot of the proportion of repeat elements inserted into the downstream introns of conserved, species-specific and background circRNAs from lemur to human. * $-p < 0.05$; ** $-p < 0.005$, *** $-p < 1 \times 10^{-5}$. p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). (F) A schematic model of the results describing impact of our observations on circRNA formation. Boxes represent exons, straight lines are introns, repeat elements are red, arced lines represent back-spliced junction, dashed lines represents RNA-RNA duplex.

201 **Discussion**

202 The evolution of circRNAs has been previously studied across extensive evolutionary time reveal-
203 ing poor conservation for the majority of circRNAs (*Rybak-Wolf et al., 2015; Venø et al., 2015*). Our
204 approach is unique as it focuses on the conservation of circRNAs in very closely related species en-
205 abling us to account for the rapid evolution of non-coding RNAs. This increased resolution allowed
206 us to reveal two disparate facts about circRNA expression. Firstly, we observe extensive variation
207 in the production of the vast majority of circRNAs between species. With circRNAs often expressed
208 within the same orthologous genes even if back-spliced junction is not conserved. Conversely, we
209 identify a core set of over 700 circRNA that are conserved across millions of years of evolution.
210 These circRNAs have higher inclusion rates and show increased inclusion across evolutionary age.
211 Both groups are related in the cis- and trans-regulatory features that likely drive circRNA forma-
212 tion such as evidence of recent transposons insertion and extended adjacent introns. However,
213 the conserved groups show decreased diversity of circRNA production and increased expression
214 potentially suggesting a combination circRNA selection and retrotransposon suppression is occur-
215 ring.

216 A host of endogenous mechanisms dampen down the impact of the retrotransposons within
217 gene bodies. For example, the formation of Alu exons is suppressed by the nuclear ribonucleo-
218 protein HNRNPC (*Zarnack et al., 2013*) and the nuclear helicase DHX9 binds to inverted repeat Alu

219 elements to suppress circRNA formation (*Aktaş et al., 2017*). Over time though, in selected exam-
220 ples, these inclusions can promote novel functionality (*Shen et al., 2011; Attig et al., 2016, 2018;*
221 *Avgan et al., 2019*) enabling the creation of tissue-specific exons (*Attig et al., 2018*), miRNAs (*Gu*
222 *et al., 2009; Spengler et al., 2014*) and promoter regions (*Li et al., 2018a; Zhang et al., 2019*). Our
223 results suggest circRNAs are undergoing a similar selection race with the recent insertion of multi-
224 ple retrotransposons promoting increased circRNA production that in some cases stabilizes over
225 time. It is important to note though that the production of a large number of circRNAs in itself can
226 be functional (*Liu et al., 2019*). For example, in the immune system a wide diversity of circRNAs
227 are produced and sequester specific RNA binding proteins. These proteins are released upon viral
228 infection to inhibit translation of viral RNA (*Liu et al., 2019*). A major challenge for the field in the
229 following years will arise from determining the contribution of noise versus function for each of
230 these groups.

231 The investigation of mechanisms controlling circRNA production is a rapid and expanding field
232 (*Li et al., 2018a*). Our results support a kinetic model (*Schor and AR, 2013*) for circRNA function
233 whereby trans-factors promote spliceosome recruitment to the final exon and the very long down-
234 stream introns extend the time-window for back-splicing to occur, which is facilitated by inverted
235 repeats increasing the proximity of 3' -splice site with the upstream 5' -splice site (see Fig. 4F). The
236 extension of the final intron therefore increases the likelihood of circRNA formation in time and
237 space. Spatially by the introducing new retrotransposons, which facilitate RNA-RNA duplex forma-
238 tion (*Jeck et al., 2013; Liang and Wilusz, 2014; Ivanov et al., 2015*) to orientate the splice sites in
239 close proximity, and temporary by increasing the time-window for such an event to occur. The
240 conservation of circRNAs we observe could therefore just be a result of increasing the probability
241 for such an event to occur rather than evidence of functionality. However, circRNAs represent an
242 extreme example of a trend in post-transcriptional regulation whereby low leaky expression cre-
243 ates a pool of possible novel substrates (*Barbosa-Morais et al., 2012; Merkin et al., 2012; Reyes*
244 *et al., 2013; Mattick, 2018; Avgan et al., 2019; Fiszbein et al., 2019*) increasing the likelihood for
245 unique functionality to arise (*Gueroussov et al., 2017; Guo et al., 2020*). For circRNAs this can be
246 aided by single nucleotide changes that enable trans-acting factors, such as Quaking to facilitate
247 circRNA formation (*Conn et al., 2015*).

248 In conclusion, our evolutionary analysis identifies that the noisy production of circRNAs is driven
249 by the insertion of novel transposons in adjacent downstream introns that can over time stabilizes

250 to produce conserved circRNAs. This provides a pool of evolutionary potential that could contribute
251 to the evolutionary rewiring of the cell.

252 **Methods and Materials**

253 **Data processing**

254 All fastq files were quality checked using FastQC (*Andrews, 2010*). Adapters and low quality se-
255 quences were removed using Cutadapt (*Martin, 2011*).

256 **Datasets**

257 Ribo-minus RNA-seq data was extracted from the publically available Nonhuman Primate Refer-
258 ence Transcriptome Resource (NHPRTR) resource (<http://www.nhprtr.org/>; (*Peng et al., 2015*)). The
259 analyzed samples were from chimpanzee, rhesus macaque, cynomolgus macaque mauritian, olive
260 baboon, common marmoset, squirrel monkey and mouse lemur to cover the 70 MYA of primate
261 evolution (**Supplementary Table 1**). The primates samples of above species were chosen based
262 on the availability of chain files for LiftOver analysis. Human samples were retrieved from differ-
263 ent publically available Ribo-minus datasets searching for the SRA IDs in the circAtlas 2.0 database
264 (<http://circatlas.biols.ac.cn/>; [*Wu et al., 2020*])) (**Supplementary Table 1**). Replicates of certain
265 samples across the different primates data were merged to achieve a higher sequencing depth
266 required for alternative splicing quantification (**Supplementary Table 5**).

267 **Alternative splicing, back-splice junction, and gene expression quantification**

268 Whippet (*Sterne-Weiler et al., 2018*) was used to analyze the RNA-seq samples to quantify cas-
269 set exon (CE) events, circRNAs (back-spliced junctions; BSJ) and gene expression. To enable BSJ
270 quantification we used the setting with the `-circ` parameter when running Whippet-quant <https://github.com/timbitz/Whippet.jl>.
271 <https://github.com/timbitz/Whippet.jl>.

272 The splice graphs of all primates used for Whippet quantification were calculated using the
273 genome annotation files for each primate from Ensembl (*Yates et al., 2020*) (**Supplementary Ta-**
274 **ble 6**). The genome annotation files were supplemented with novel exon-exon junctions derived
275 from whole genome alignment of primates samples using STAR (*Dobin et al., 2013*) with the 2-pass
276 setting and `outFilterMultimapNmax == 10` parameters. Whippet index command was run with the
277 `-bam` and `-suppress-low-tsl` parameters.

278 Gene expression of orthologue genes was retrieved from the gene.tpm.gz files from Whippet-
279 quant output. The correlations of gene expression of orthologue genes between tissue samples
280 from all primates were calculated using the Pearson correlation. Clustering of correlation values
281 was assessed and visualized with a heatmap using the p.heatmap function in R.

282 **Identification of expressed circRNAs and cassette-exons**

283 All the BSJ events present in orthologue genes between the species mentioned above were fil-
284 tered to find conserved circRNAs identified by Whippet. The orthologue list of genes was retrieved
285 from Ensembl using the bioMart R package (*Smedley et al., 2009*). Expressed BSJs were defined
286 according to an expression and percent of spliced in (PSI) cutoff of at least 5 reads and $\geq 5\%$ of
287 PSI respectively. Cassette-exon (CE) events from Whippet output were also filtered, keeping those
288 present in orthologue genes and with $PSI \geq 10$

289 **Conservation analysis of circRNAs**

290 We defined a circRNA as conserved if the exon(s) that formed the BSJ are orthologous to the human
291 exon(s) that also formed the BSJ. To achieve this, the exon coordinates of orthologue genes, of
292 each primate were retrieve from the GTF files downloaded from Ensembl (**Supplementary Table**
293 **6**). Then, the exon coordinates from the GTF files were intersected with the CE coordinates from
294 Whippet using bedtools intersect (*Quinlan and Hall, 2010*) with -wa parameter.

295 Then, the resulted exon coordinates (GTF-CE coordinates) were intersected with the circRNAs
296 coordinates within orthologue genes using bedtools intersect with -loj parameter to find which
297 exons were forming the circRNA. The exon coordinates within the circRNA coordinate of the non-
298 human primates were mapped to human coordinates using the UCSC LiftOver (*Navarro Gonzalez*
299 *et al., 2021*) to retrieve orthologue exons.

300 The orthologue exons between primates and human were matched to human exon coordinates
301 within the circRNAs coordinates in human to find conserved circRNAs. We defined if a circRNA was
302 conserved between a primate and human if the exon(s) forming the BSJ of the circRNA were also
303 conserved and if the exon(s) start and end coordinates were \leq of 100 nc from the start and end of
304 the BSJ coordinate (**Figure 2-Figure supplement 3** for schematic). We defined as non-conserved
305 circRNAs all the human circRNAs that do not have orthologue exons forming the BSJ of the circRNA
306 with other primates.

307 **Conserved and tissue conserved circRNAs**

308 The list of orthologous circRNAs was plotted in an UpSet plot to visualize the intersection of circR-
309 NAs between primates species. We defined the set of conserved circRNAs as the circRNAs within
310 the intersections between primates species where human, chimpanzee and baboon always ap-
311 peared.

312 The correlation of inclusion of conserved and tissue-conserved circRNAs between all samples
313 was calculated using the Pearson correlation. Then correlation values were plotted in a heatmap
314 using the `p.heatmap` function in R.

315 **Differential gene expression analysis and enrichment analysis of genes with con-** 316 **served circRNAs**

317 EdgeR (*Robinson et al., 2010*) library was used to perform the differential gene expression analy-
318 sis between neuronal samples (brain, cerebellum and frontal cortex) and non-neuronal samples
319 (heart, skeletal muscle, liver, lung, spleen and colon). This analysis showed 8,817 differential ex-
320 pressed genes according to a log Fold Change cutoff of $\log_2(1.5)$ and FDR of 0.05. There were 212
321 genes of the conserved circRNAs (total of 442 genes) in the set of differential expressed genes. The
322 enrichment of genes with conserved circRNAs was statistically tested with a hypergeometric test
323 using the `phyper` function in R. The parameters were $q = 212$, $m = 8,817$, $n = 11,278$, $k = 442$, and
324 `lower.tail = FALSE`.

325 **Conserved cassette exons in primates**

326 All exons coordinates of orthologue genes from the GTF files and CE exons coordinates from Whip-
327 pet were mapped to human coordinates using UCSC LiftOver (*Navarro Gonzalez et al., 2021*). The
328 PSI values of orthologous exons in genes of conserved and tissue-conserved circRNAs were re-
329 trieved from all tissues samples of human, chimpanzee and baboon and calculated the Pearson
330 correlation values. The correlation values were plotted in a heatmap using the `p.heatmap` function.

331 **Comparison of circRNAs expression and conservation**

332 circRNAs expression of conserved, tissue-conserved and non-conserved circRNAs was calculated
333 using relative TpMs. The relative TpMs were calculated with the below equation.

$$RelativeTpMs = \frac{(circRNAsReads)(GeneTpMs)}{GeneReads} \quad (1)$$

334 The expression values of conserved and non-conserved circRNAs, and tissue-conserved and
335 non-conserved circRNAs of replicates of the same tissue in human samples were plotted in scatter
336 plots.

337 The median relative TpM of conserved (and tissue-conserved) and non-conserved circRNAs of
338 human samples were also calculated. The expression values between mentioned sets were sta-
339 tistically compared using a Wilcox test. The parameters of the Wilcox test were x = Conserved (or
340 tissue conserved) circRNAs TpMs, y = Non-conserved circRNAs TpMs, alternative = "greater". The
341 median relative TpM was plotted in violin plots using the ggplot2 R library ([Wickham, 2016](#)).

342 The median PSI values of conserved, tissue-conserved and non-conserved circRNAs across all
343 human samples were calculated. Their inclusion levels were statistically compared using the Wilcox
344 test function in R with the parameters x = Conserved (or tissue conserved) circRNAs median PSI,
345 y = Non-conserved circRNAs median PSI, alternative = "greater". The distribution of the median
346 PSI values of conserved and non-conserved circRNAs; and tissue-conserved and non-conserved
347 circRNAs were plotted in a cumulative plot using the ggplot2 library in R.

348 The median PSI value of shared circRNAs between evolutionary interesting sets (human (species-
349 specific circRNAs); hominoids; hominoids and baboon; hominoids and old-world monkeyes; homi-
350 noids, old-world monkeys and marmoset; and hominoids, old-world monkeys and new-world mon-
351 keys) shown in the UpSet plot were calculated, plotted in a cumulative plot and statistically com-
352 pared using a Wilcox test.

353 Seven of our reported circRNA from the lists of conserved and tissue conserved circRNAs were
354 of special interest as they were previously reported ([Gokool et al., 2020b](#)) to be highly expressed
355 in human cerebellum and frontal cortex. The PSI values of such circRNAs were compared across
356 all tissues in the eight primates species.

357 **Comparison of the number of orthologue genes producing a circRNA and number**
358 **of conserved circRNAs between species**

359 The number of times an orthologue gene produces at least one circRNA in any of the analyzed
360 species was counted, as well as the number of times a circRNA was shared between another pri-
361 mate. The percentage of shared genes or circRNAs between the eight species was calculated and
362 plotted in a barplot using the ggplot2 library in R.

363 **Comparison of start and end position of circRNAs between conserved and non-**
364 **conserved circRNAs**

365 circRNAs can be formed from unique start and end exons forming the BSJ, repeated start exons,
366 repeated end exons, or repeated start and end exons (see **Figure 2-Figure supplement 3** for
367 schematic). The percentage of conserved and non-conserved circRNAs that fall in the above cate-
368 gories was calculated and plotted using the ggplot2 library in R.

369 **Generalized logistic regression**

370 All continuous data was normalized to ensure a fair comparison between features using scale()
371 package in R environment. Multicollinearity was assessed using the vif() from the R package car.
372 The dataset was split into training (80%) and test (20%). To optimize the selection of the model and
373 the importance of each feature we used the R package glmulti (*Calcagno and De Mazancourt, 2010*).
374 To select from all possible models the selection process used a genetic algorithm (method = 'g') with
375 Akaike information criterion (AIC - crit = "aic"). To calculate the generalized logistic model, glmulti
376 used the R module glm with family = binomial(). ROC curve was calculated using R's pROC library
377 with test data. Data extracted from this model is reported together with p-value and z-values are
378 reported in **Supplementary Table 7**.

379 **Background Datasets**

380 Two background datasets were used in this study: background and species-specific (**Supplementary**
381 **Table 2**). The "background" datasets consisted of exon combinations only within genes with circR-
382 NAs. The dataset was constructed by identifying alternative exons within gene of interest ($10 < \text{PSI} < 90$
383 within any of the tissues studied) and using python function random to assign these exons to-
384 gether. The "species-specific" dataset was constructed as described above of human circRNA with

385 no evidence of their back-spliced junction being conserved in any other primate species. For both
386 datasets only genes with orthologous genes in all tested primates species were used (based on
387 Ensembl annotation) and only orthologous exons (based on liftover – see above) were used.

388 **CircRNA features**

389 MaxEntScan (*Yeo and Burge, 2004*) was used to estimate the strength of 3' and 5' splice sites. 5'
390 splice site strength was assessed using a sequence including 3 nt of the exon and 6 nt of the adja-
391 cent intron. 3' splice site strength was assessed using a sequence including - 20 nt of the flanking
392 intron and 3 nt of the exon. SVM-BPfinder (*Corvelo et al., 2010*) was used to estimate branchpoint
393 and polyprimidine tract strength and other statistics. Scores calculated using the sequence of in-
394 trons to the 3'end of exon between 20 and 500 nt.

395 Transcription start sites (TSS) were downloaded from Biomart. GC content was calculated using
396 python script. Transposon information download from RepeatMasker as described below.

397 Nucleosome occupancy for HepG2 cells was calculated using data from Enroth et al. (*Enroth*
398 *et al., 2014*). Colospace read data was aligned using Bowtie (Langmead, 2010) (-S -C -p 4 -m 3 -best
399 -strata) using index file constructed from Ensembl Hg38. Nuctools (with default settings) was used
400 to calculate occupancy profiles and calculate occupancy at individual regions (*Vainshtein et al.,*
401 *2017*).

402 All CLiP-seq data and CHIP-seq data was downloaded pre-processed bed data files from EN-
403 CODE (*Sundararaman et al., 2016*) with only narrowpeaks calculated using both isogenic replicates
404 used. Bedtools intersect (-wao) was used to identify overlap with candidate regions. Overlap for
405 all groups of trans-factors were collated and scores normalized by nucleotide length. Groups were
406 based on annotation and split into positive regulators of splicing (SR: Serine/Arginine region con-
407 taining proteins) and negative regulators of splicing (hnRNP: Heterogeneous nuclear ribonucleo-
408 proteins).

409 In feature analysis, only first and last exons of circRNA, and their surrounding introns, were
410 included in the analysis. The upstream portion is considered as the region 5' of elements (i.e. first
411 exon) and downstream portion is 3' of elements.

412 **Overlap with known repeat elements**

413 Repeat elements identified by RepeatMasker were downloaded from UCSC table browser (*Navarro Gon-*
414 *zalez et al., 2021*) in bed format. Bedtools intersect (-wao) was used to identify overlap of trans-
415 posons with novel exons.

416 The frequency of transposable events is calculated as the proportion of transposons overlap-
417 ping area of interest (i.e. exon 1). All transposons were grouped together into 12 categories (AluJ,
418 AluS, AluY, L1, L2, L3, MIR, MER, FLAM, AT_rich, SINE and everything else into "other") based on
419 annotation from RepeatMasker. Flanking regions are defined as having the same transposable
420 elements on different strands in both introns adjacent to the circRNA.

421 **Intronic length and transposons comparison of human and lemur**

422 Orthologs exons between human and lemur containing circRNAs were identified using the proce-
423 dure described above. Intron length was determined based on the nearest exon from ENSEMBL
424 annotation (*Yates et al., 2020*) with evidence from RNA-seq data of expression (PSI>10). To iden-
425 tify regions unique to human, the intronic regions unique to human were split into windows of 20
426 nucleotides. Liftover was used to identify conserved regions between human and lemur genomes
427 for each of these windows. Regions with no evidence of conservation were overlapped (using
428 bedtools intersect -wao) with UCSC RepeatMasker (*Navarro Gonzalez et al., 2021*) annotation to
429 identify novel transposon insertion.

430 **Supplementary Data**

431 **Supplementary Figures**

432 Supplementary Figures are available as an attachment to this document.

433 **Supplementary Tables**

434 Supplementary Tables are available as a separate attachment to this manuscript.

435 **Supplementary Table 1** Primates datasets IDs and sequencing depth information.

436 **Supplementary Table 2** Conserved, non-conserved (human specific) and background circR-
437 NAs.

438 **Supplementary Table 3** Features associated with circRNA formation (trans- and cis- regulatory
439 factors and all major groups of transposons).

440 **Supplementary Table 4** Primates specific circRNAs.

441 **Supplementary Table 5** Information about merged samples to acquire higher sequencing
442 depth.

443 **Supplementary Table 6** Genome version, GTF and chain file information.

444 **Supplementary Table 7** GLM output.

445 **Acknowledgments**

446 We gratefully acknowledge John Mattick, Akira Gookol, Juli Wang and Helen King for helpful dis-
447 cussions and feedback on this study, as well as all members of the Weatheritt Lab. G.S.R was
448 supported by a UNSW UIPA PhD Scholarship. This research was supported by the NSW Institute of
449 Cancer Research (RJW), the Scrimshaw Foundation (RJW), the Australian Research Council (ARC) Dis-
450 covery Project (RJW, IV), an ARC future fellowship (IV) and a University of New South Wales Scientia
451 Fellowship (IV).

452 **Author contribution**

453 Contributions to this publications are distributed as follows: Study design: G.S.R, I.V., R.W.; Bioin-
454 formatic data analyses: G.S.R and R.W.; Paper manuscript and discussion: G.S.R, I.V. and R.W.

455 **Competing interests**

456 No competing interests.

457 **References**

- 458 **Aktaş T**, et al. DHX9 suppresses RNA processing defects originating from the Alu invasion of the human
459 genome. *Nature*. 2017; 544:115–119.
- 460 **Andrews S**. FastQC: a quality control tool for high throughput sequence dataFastQC: a quality control tool for
461 high throughput sequence data. Available online at: ; 2010. [http://www.bioinformatics.babraham.ac.uk/projects/
462 fastqc](http://www.bioinformatics.babraham.ac.uk/projects/fastqc).
- 463 **Attig J**, et al. Splicing repression allows the gradual emergence of new Alu-exons in primate evolution. *Elife*.
464 2016; 5.
- 465 **Attig J**, et al. Heteromeric RNP Assembly at LINEs Controls Lineage-Specific RNA Processing. *Cell*. 2018;
466 174:1067–1081.
- 467 **Avgan N**, Wang JI, Fernandez-Chamorro J, Weatheritt RJ. Multilayered control of exon acquisition permits the
468 emergence of novel forms of regulatory control. *Genome Biol*. 2019; 20:141.

- 469 **Barbosa-Morais NL**, et al. The evolutionary landscape of alternative splicing in vertebrate species. *Science*.
470 2012; 338:1587–1593.
- 471 **Barrett SP**, Wang PL, Salzman J. Circular RNA biogenesis can proceed through an exon-containing lariat pre-
472 cursor. *Elife*. 2015; 4(e07540).
- 473 **Brawand D**, et al. The evolution of gene expression levels in mammalian organs. *Nature*. 2011; 478:343–348.
- 474 **Calcagno V**, De Mazancourt C. glmulti: An R Package for Easy Automated Model Selection with (Generalized)
475 Linear Models. *Journal of Statistical Software*. 2010; 34.
- 476 **Cardoso-Moreira M**, et al. Gene expression across mammalian organ development. *Nature*. 2019; 571:505–
477 509.
- 478 **Cheetham SW**, Faulkner GJ, Dinger ME. Overcoming challenges and dogmas to understand the functions of
479 pseudogenes. *Nat Rev Genet*. 2020; 21:191–201.
- 480 **Conn SJ**, et al. The RNA binding protein quaking regulates formation of circRNAs. *Cell*. 2015; 160:1125–1134.
- 481 **Corvelo A**, Hallegger M, Smith CW, Eyras E. Genome-wide association between branch point properties and
482 alternative splicing. *PLoS Comput Biol*. 2010; 6(e1001016).
- 483 **Dobin A**, et al. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*. 2013; 29:15–21.
- 484 **Enroth S**, et al. Nucleosome regulatory dynamics in response to TGF β . *Nucleic Acids Res*. 2014; 42:6921–6934.
- 485 **Fiszbein A**, Krick KS, Begg BE, Burge CB. Exon-Mediated Activation of Transcription Starts. *Cell*. 2019; 179:1551–
486 1565.
- 487 **Gokool A**, Anwar F, Voineagu I. The Landscape of Circular RNA Expression in the Human Brain. *Biol Psychiatry*.
488 2020; 87:294–304.
- 489 **Gokool A**, Loy CT, Halliday GM, Voineagu I. Circular RNAs: The Brain Transcriptome Comes Full Circle. *Trends*
490 *Neurosci*. 2020; 43:752–766.
- 491 **Gu TJ**, Yi X, Zhao XW, Zhao Y, Yin JQ. Alu-directed transcriptional regulation of some novel miRNAs. *BMC*
492 *Genomics*. 2009; 10:563.
- 493 **Guerossov S**, et al. Regulatory Expansion in Mammals of Multivalent hnRNP Assemblies that Globally Control
494 Alternative Splicing. *Cell*. 2017; 170:324–339.
- 495 **Guo CJ**, et al. Distinct Processing of lncRNAs Contributes to Non-conserved Functions in Stem Cells. *Cell*. 2020;
496 181:621–636.
- 497 **Ha KCH**, Blencowe BJ, Morris Q. QAPA: a new method for the systematic analysis of alternative polyadenylation
498 from RNA-seq data. *Genome Biol*. 2018; 19:45.

- 499 **Ha KCH**, Sterne-Weiler T, Morris Q, Weatheritt RJ, Blencowe BJ. Differential contribution of transcriptomic
500 regulatory layers in the definition of neuronal identity. *Nat Commun.* 2021; 12:335.
- 501 **Irimia M**, Blencowe BJ. Alternative splicing: decoding an expansive regulatory layer. *Curr Opin Cell Biol.* 2012;
502 24:323–332.
- 503 **Ivanov A**, et al. Analysis of intron sequences reveals hallmarks of circular RNA biogenesis in animals. *Cell Rep.*
504 2015; 10:170–177.
- 505 **Jeck WR**, et al. Circular RNAs are abundant, conserved, and associated with ALU repeats. *RNA.* 2013; 19:141–
506 157.
- 507 **Li C**, Lenhard B, Luscombe NM. Integrated analysis sheds light on evolutionary trajectories of young transcrip-
508 tion start sites in the human genome. *Genome Res.* 2018; 28:676–688.
- 509 **Li X**, Yang L, Chen LL. The Biogenesis, Functions, and Challenges of Circular RNAs. *Mol Cell.* 2018; 71:428–442.
- 510 **Liang D**, Wilusz JE. Short intronic repeat sequences facilitate circular RNA production. *Genes Dev.* 2014;
511 28:2233–2247.
- 512 **Liang D**, et al. The Output of Protein-Coding Genes Shifts to Circular RNAs When the Pre-mRNA Processing
513 Machinery Is Limiting. *Mol Cell.* 2017; 68:940–954.
- 514 **Liu CX**, et al. Structure and Degradation of Circular RNAs Regulate PKR Activation in Innate Immunity. *Cell.*
515 2019; 177:865–880.
- 516 **Martin M**. Cutadapt Removes Adapter Sequences From High-Throughput Sequencing Reads. *EMBnetjournal.*
517 2011; 17.
- 518 **Mattick JS**. The State of Long Non-Coding RNA Biology. *Noncoding RNA.* 2018; 4.
- 519 **Memczak S**, et al. Circular RNAs are a large class of animal RNAs with regulatory potency. *Nature.* 2013;
520 495:333–338.
- 521 **Merkin J**, Russell C, Chen P, Burge CB. Evolutionary dynamics of gene and isoform regulation in Mammalian
522 tissues. *Science.* 2012; 338:1593–1599.
- 523 **Navarro Gonzalez J**, et al. The UCSC Genome Browser database: 2021 update. *Nucleic Acids Res.* 2021;
524 49:D1046–D1057.
- 525 **Nilsen TW**, Graveley BR. Expansion of the eukaryotic proteome by alternative splicing. *Nature.* 2010; 463:457–
526 463.
- 527 **Peng X**, et al. Tissue-specific transcriptome sequencing analysis expands the non-human primate reference
528 transcriptome resource (NHPRT). *Nucleic Acids Res.* 2015; 43:D737–42.

- 529 **Pipes L**, et al. The non-human primate reference transcriptome resource (NHPRTR) for comparative functional
530 genomics. *Nucleic Acids Res.* 2013; 41:D906–14.
- 531 **Piwecka M**, et al. Loss of a mammalian circular RNA locus causes miRNA deregulation and affects brain func-
532 tion. *Science.* 2017; 357.
- 533 **Quinlan AR**, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics.*
534 2010; 26:841–842.
- 535 **Reyes A**, et al. Drift and conservation of differential exon usage across tissues in primate species. *Proc Natl*
536 *Acad Sci U S A.* 2013; 110:15377–15382.
- 537 **Robinson MD**, McCarthy DJ, Smyth GK. edgeR: a Bioconductor package for differential expression analysis of
538 digital gene expression data. *Bioinformatics.* 2010; 26:139–140.
- 539 **Rybak-Wolf A**, et al. Circular RNAs in the Mammalian Brain Are Highly Abundant, Conserved, and Dynamically
540 Expressed. *Mol Cell.* 2015; 58:870–885.
- 541 **Schor GAL I E**, AR K. Coupling between transcription and alternative splicing. *Cancer Treat Res.* 2013; 158:1–24.
- 542 **Shen S**, et al. Widespread establishment and regulatory impact of Alu exons in human genes. *Proc Natl Acad*
543 *Sci U S A.* 2011; 108:2837–2842.
- 544 **Smedley D**, et al. BioMart–biological queries made easy. *BMC Genomics.* 2009; 10:22.
- 545 **Spengler RM**, Oakley CK, Davidson BL. Functional microRNAs and target sites are created by lineage-specific
546 transposition. *Hum Mol Genet.* 2014; 23:1783–1793.
- 547 **Starke S**, et al. Exon circularization requires canonical splice signals. *Cell Rep.* 2015; 10:103–111.
- 548 **Sterne-Weiler T**, Weatheritt RJ, Best AJ, Kch H, Blencowe BJ. Efficient and Accurate Quantitative Profiling of
549 Alternative Splicing Patterns of Any Complexity on a Laptop. *Mol Cell.* 2018; 72:187–200.
- 550 **Sundararaman B**, et al. Resources for the Comprehensive Discovery of Functional RNA Elements. *Mol Cell.*
551 2016; 61:903–913.
- 552 **Vainshtein Y**, Rippe K, Teif VB. NucTools: analysis of chromatin feature occupancy profiles from high-
553 throughput sequencing data. *BMC Genomics.* 2017; 18:158.
- 554 **Venø MT**, et al. Spatio-temporal regulation of circular RNA expression during porcine embryonic brain devel-
555 opment. *Genome Biol.* 2015; 16:245.
- 556 **Wickham H**. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, -3-319-24277; 2016.
- 557 **Wu W**, Ji P, Zhao F. CircAtlas: an integrated resource of one million highly accurate circular RNAs from 1070
558 vertebrate transcriptomes. *Genome Biol.* 2020; 21:101.

559 **Yates AD**, et al. (2020) Ensembl 2020. *Nucleic Acids Res.* 2020; 48:D682–D688.

560 **Yeo G**, Burge CB. Maximum entropy modeling of short sequence motifs with applications to RNA splicing
561 signals. *J Comput Biol.* 2004; 11:377–394.

562 **Zarnack K**, et al. Direct competition between hnRNP C and U2AF65 protects the transcriptome from the ex-
563 onization of Alu elements. *Cell.* 2013; 152:453–466.

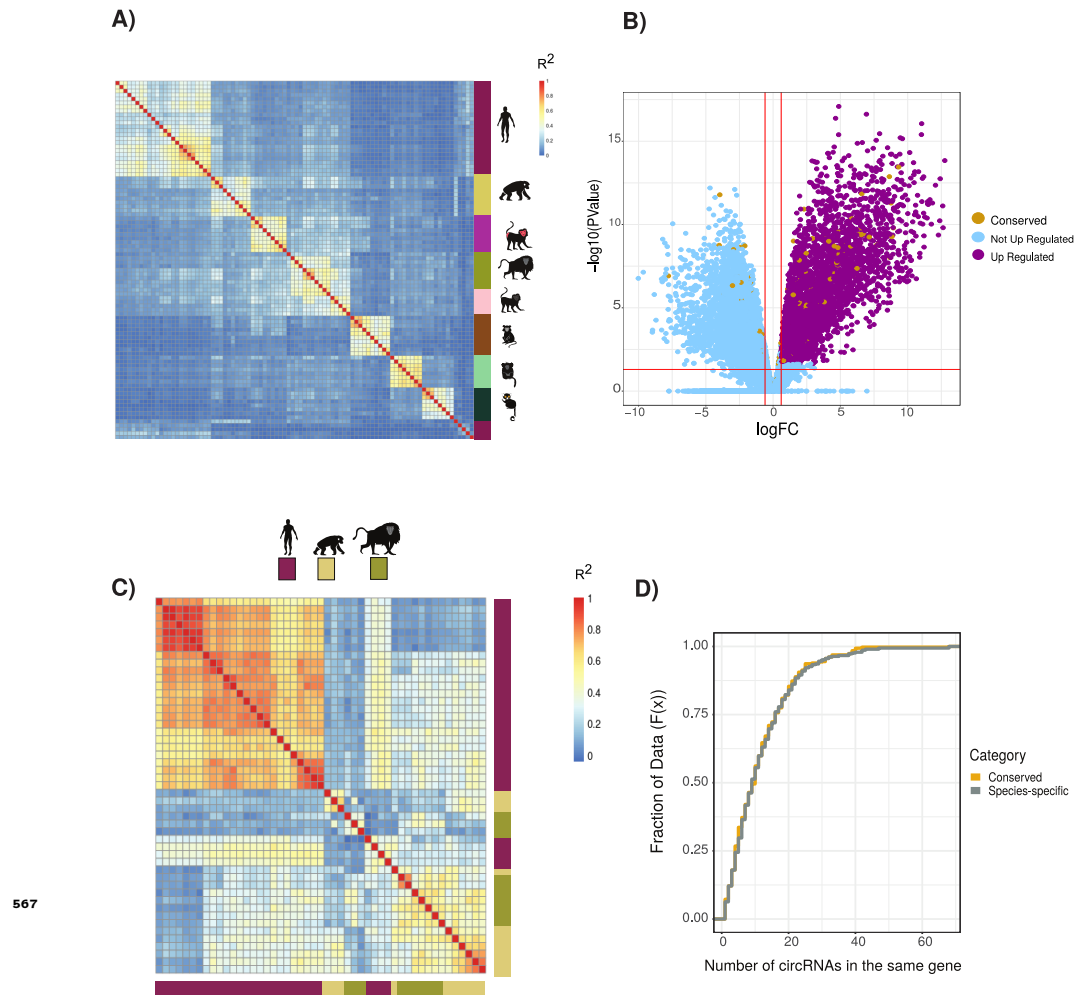
564 **Zhang XO**, Gingeras TR, Weng Z. Genome-wide analysis of polymerase III-transcribed. *Genome Res.* 2019;
565 29:1402–1414.

566 **Zhang Y**, et al. Circular intronic long noncoding RNAs. *Mol Cell.* 2013; 51:792–806.

Supplementary Figures

Evolutionary dynamics of circular RNAs in primates

Gabriela Santos-Rodriguez, Irina Voineagu, Robert J Weatheritt



567

Figure 1-Figure supplement 1. (A) Clustering of circRNAs based on percent spliced in (PSI) values. Clustered using Pearson correlation as in Fig. 1B. ($n=19,005$). Vertical heatmap indicates primate species. See Fig. 1B for details of heatmap and Supplementary Table S4 for data used. (B) Volcano plot of differential gene expression analysis between neuronal samples and non-neuronal samples. In blue are “not up-regulated genes” ($n = 15,661$), in purple are “up-regulated genes” ($n = 4,434$) and in yellow the “genes with conserved circRNAs” ($n = 442$). Likelihood of circRNA genes being enriched in differentially expression genes ($p = 0.036$, hypergeometric test). FC = fold value. P-value in figures calculated by quasi-likelihood negative binomial test and corrected for multi-testing (Bonferroni) (C) Clustering of alternative splicing events based on percent spliced in (PSI) of exons within conserved circRNAs shows no clustering by tissue. ($n = 1,256$). Vertical and horizontal adjacent heatmaps represents species. See Fig. 1B for details of heatmap. (D) Cumulative distribution plot displaying the number of circRNAs found within same gene. (see Figure 2D for description of cumulative distribution plots)

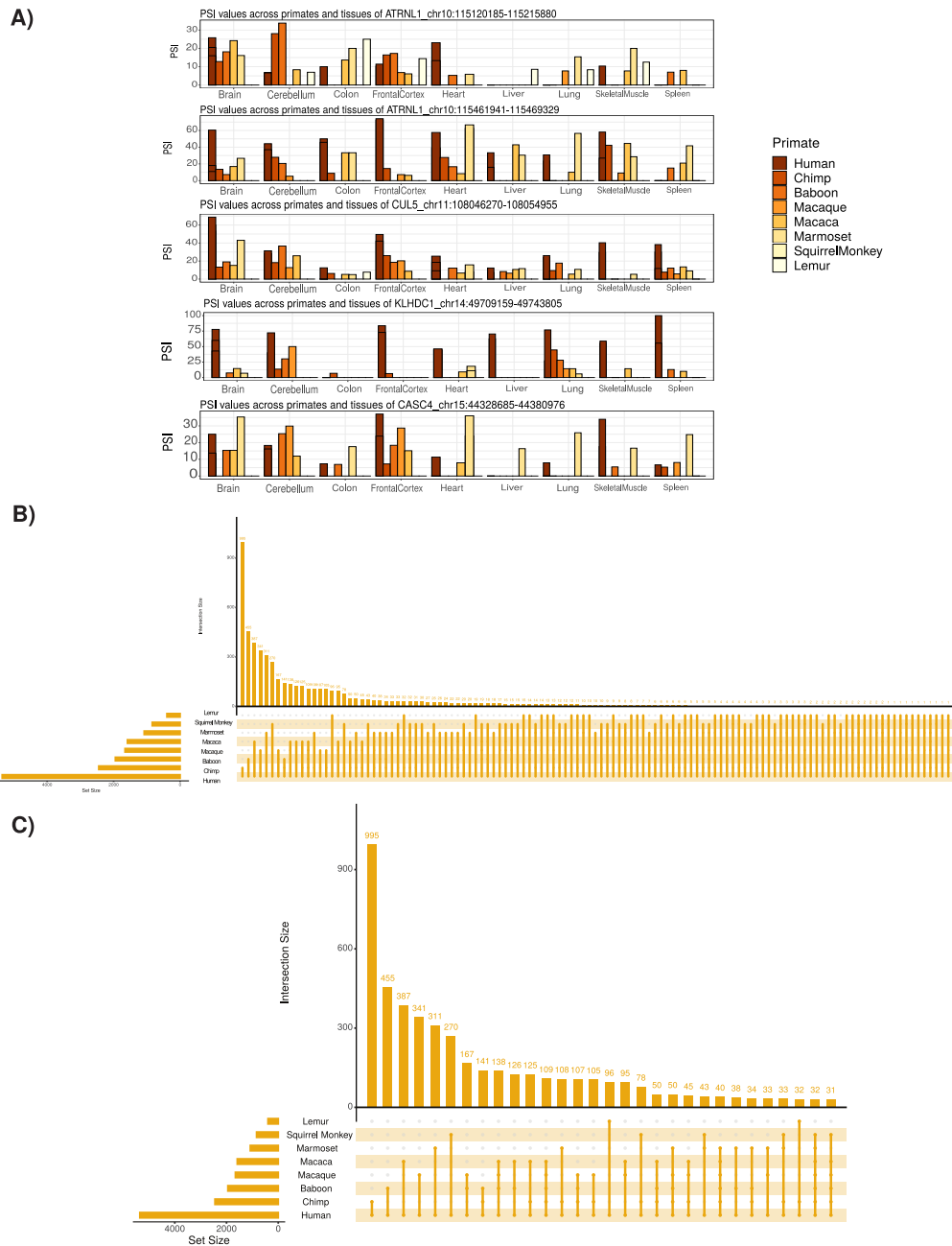


Figure 2-Figure supplement 1. (A) Extension of Figure 2B showing examples of identified circRNAs, Coordinates for each circRNA are shown. PSI = percent spliced in (B) Upset plot of conserved circRNAs across primate species analysed. An upset plot displays the intersections of a set. Each column corresponds to a set, and each row corresponds to one segment in a Venn diagram. Number of top of bars represent number in each overlap. (C) Upset plot of conserved circRNAs across primate species studied with at least 30 overlap. (see Figure 2-Supplementary Figure 1B for description of Upset plot).

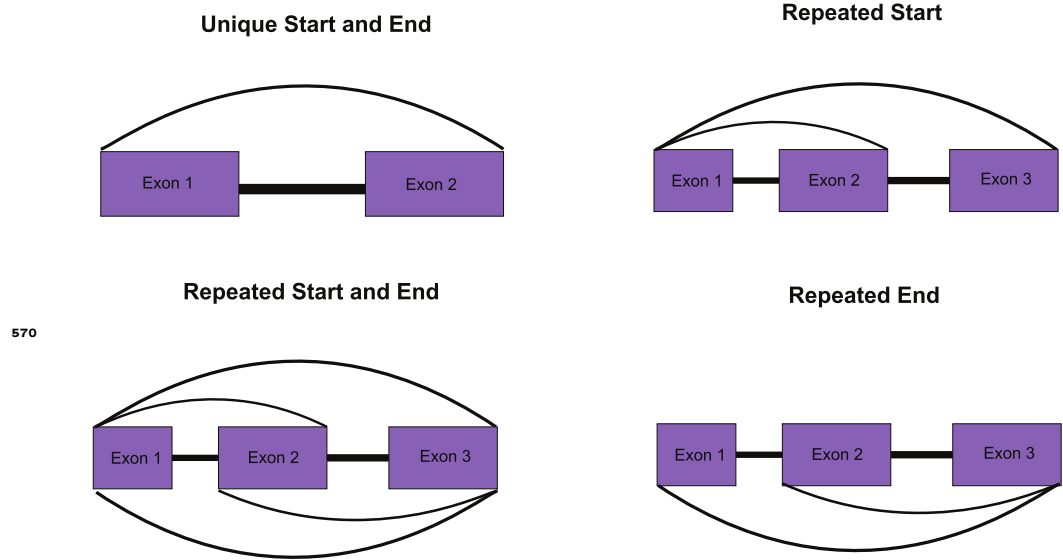


Figure 2-Figure supplement 3. Overview of approach to identifying unique circRNAs for Figure 2G (see Methods for details)

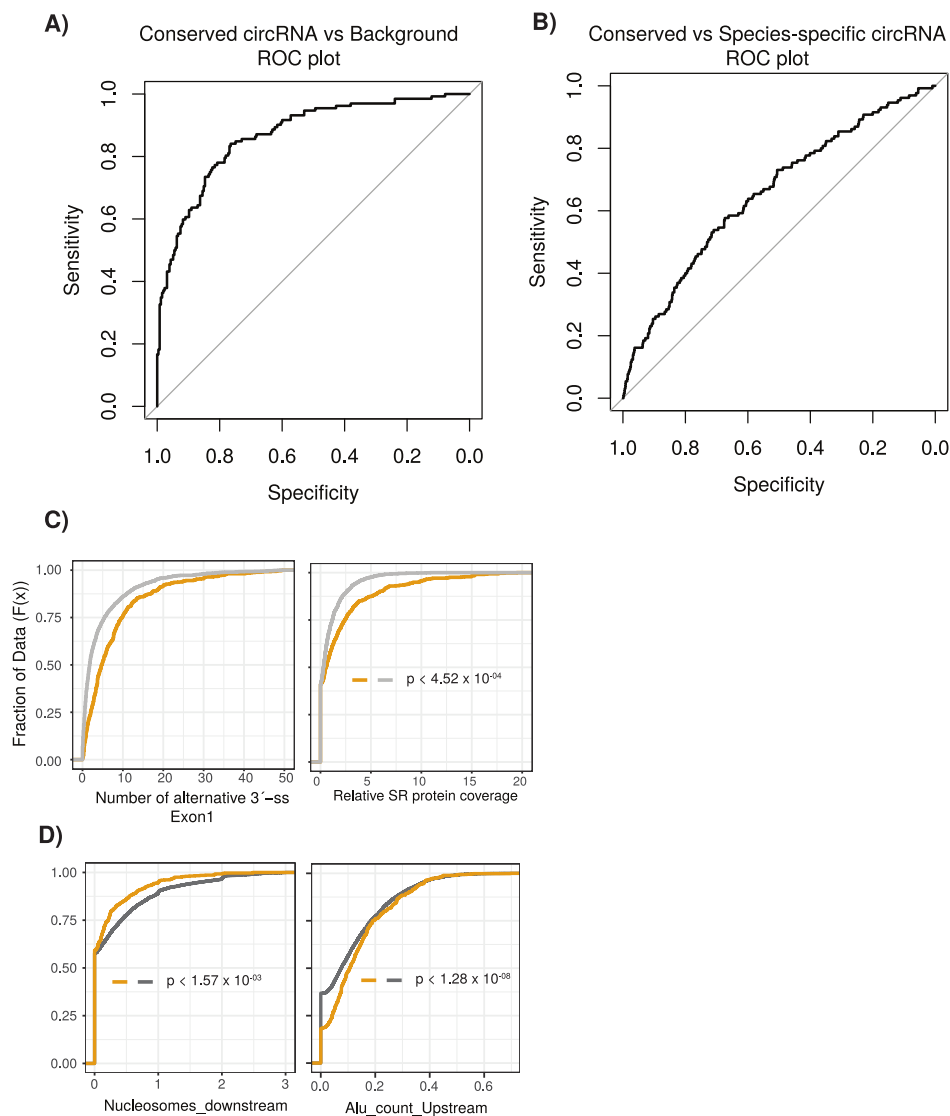


Figure 3-Figure supplement 1. (A) An ROC (Receiving operating characteristic curve) plot displaying the sensitivity (true positive rate) compared to the selectivity (false positive rate) for the logistic regression model (conserved versus background circRNA). (B) An ROC plot for logistic regression model (conserved versus species-specific circRNA). (C) Cumulative distribution plots comparing conserved versus background circRNAs of (left) alternative 3' splice sites (ss) within first exon of the circRNA and (right) of Serine/Arginine (SR) RNA-binding peaks from CLIP data. p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni). (see Figure 2D for description of cumulative distribution plots). (D) Cumulative distribution plots comparing conserved versus species-specific circRNAs of (left) nucleosome peaks in downstream intron adjacent to circRNA and (right) of Alu element content in upstream intron adjacent to circRNA. (see Figure 2D for description of cumulative distribution plots). p-values calculated by Wilcoxon rank sum test and corrected for multi-testing (Bonferroni).

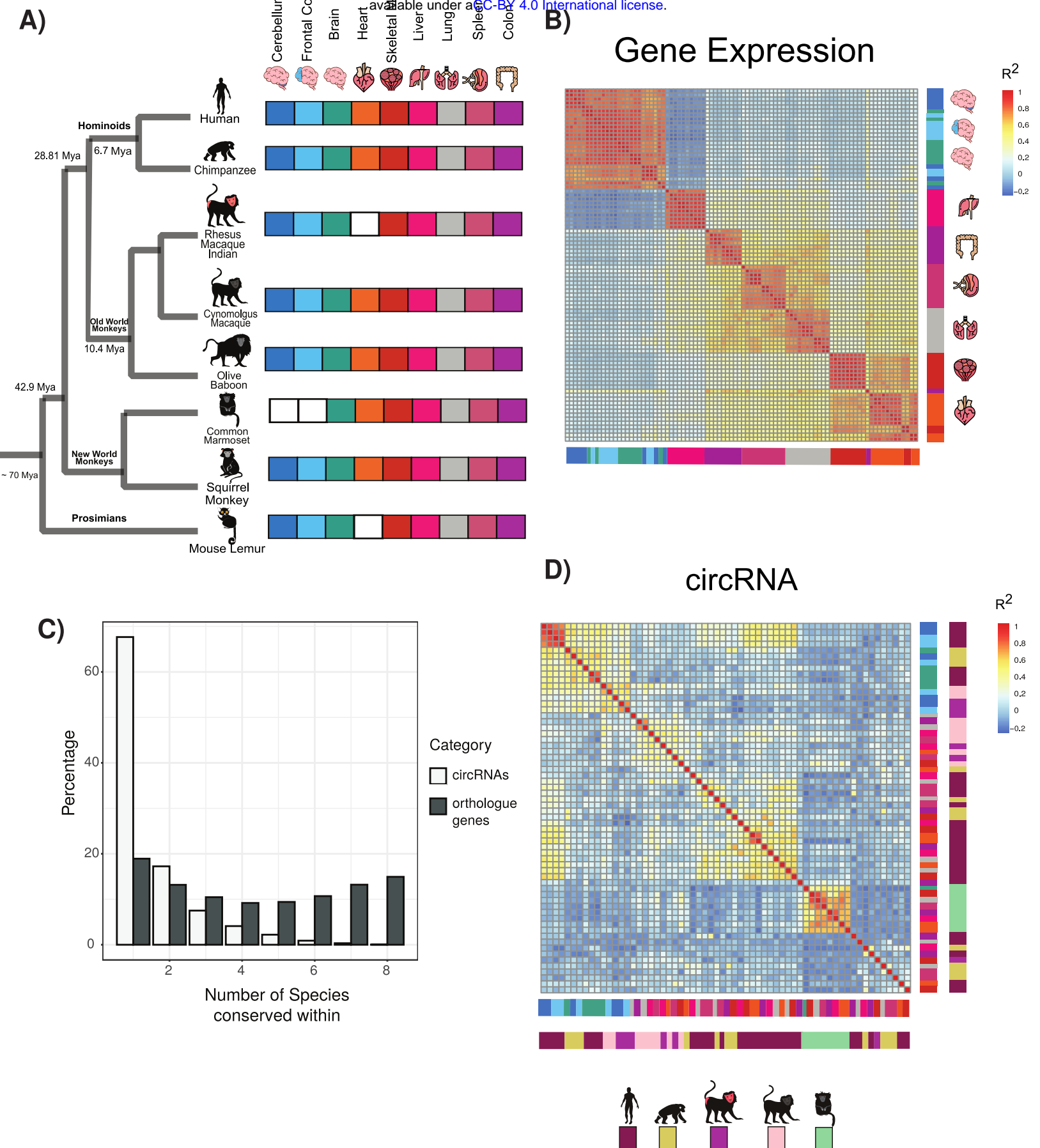


Figure 1

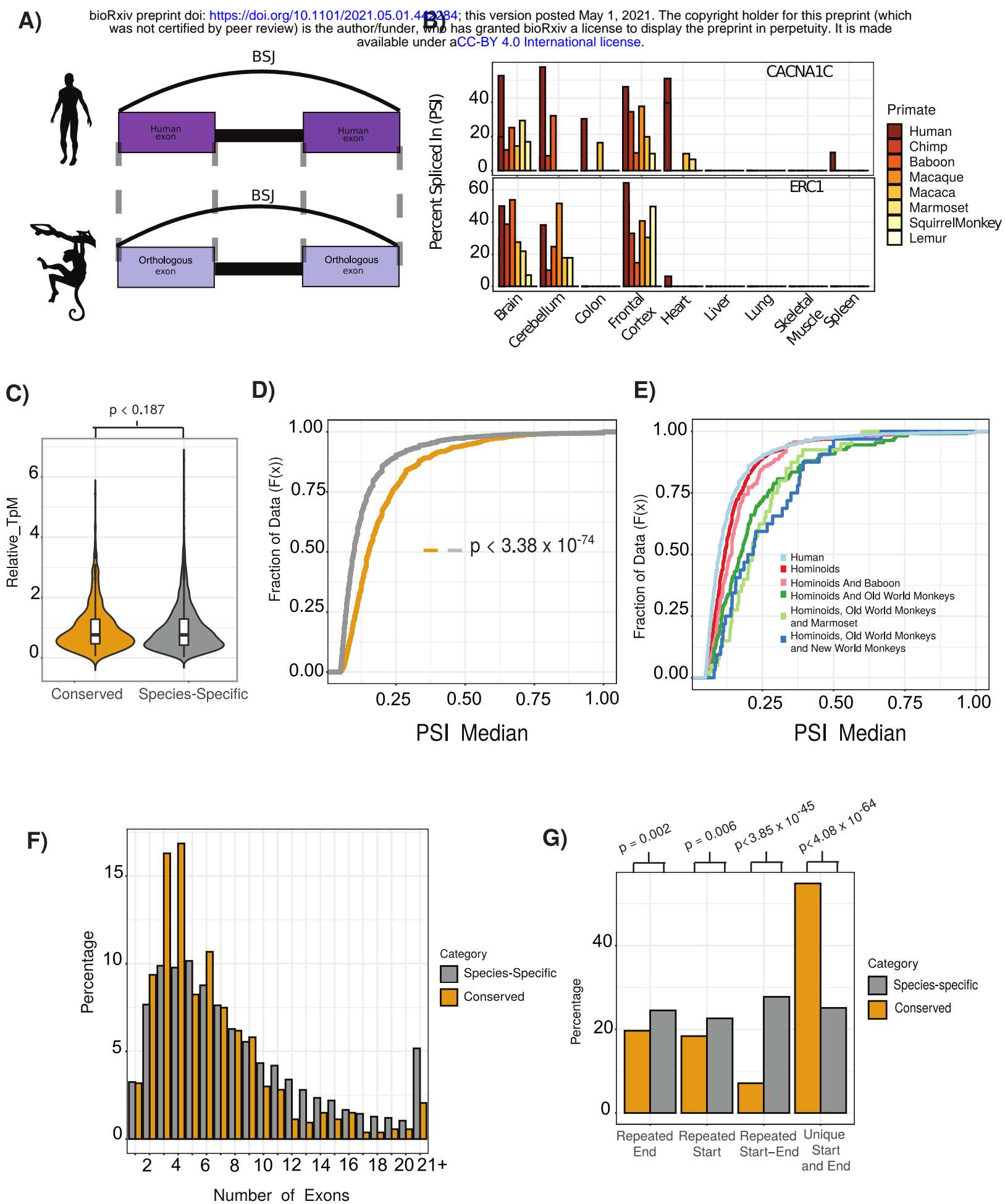
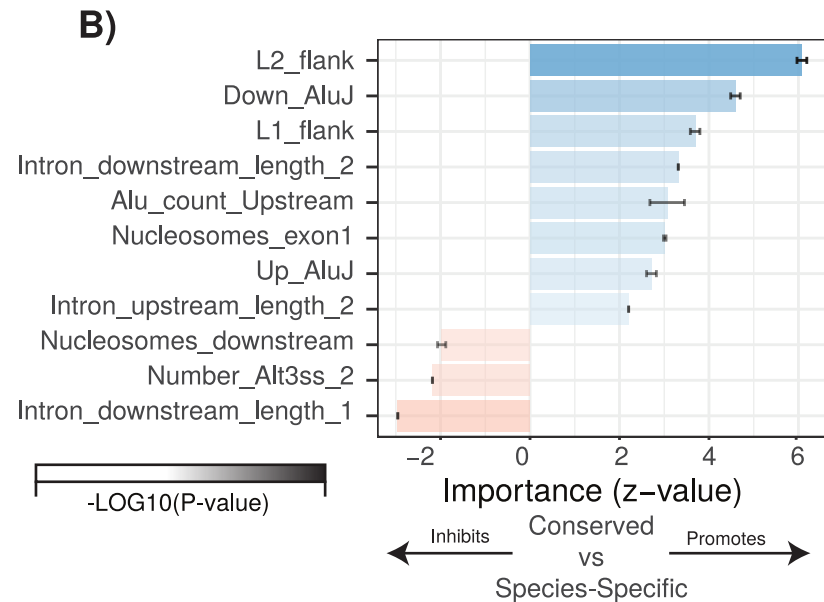
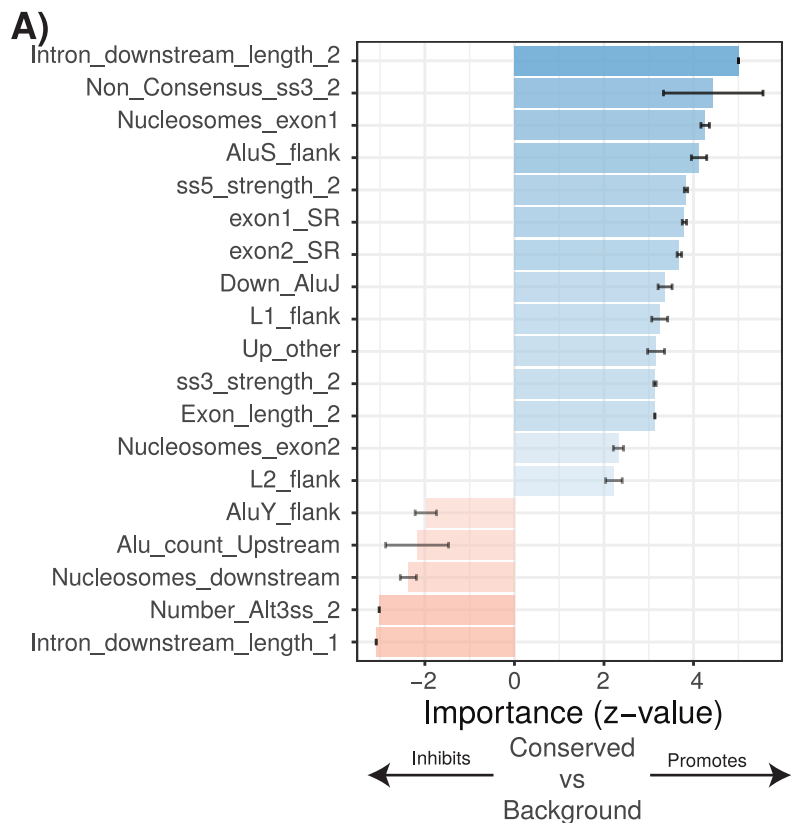


Figure 2



D) Antisense transposons surrounding circRNAs

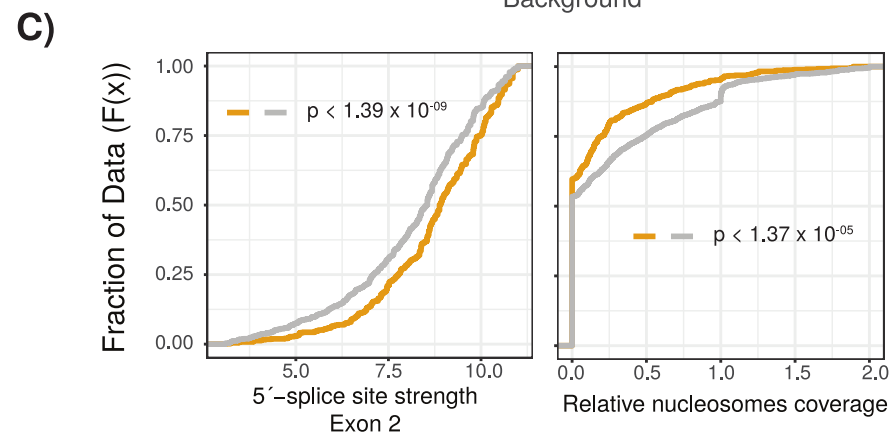
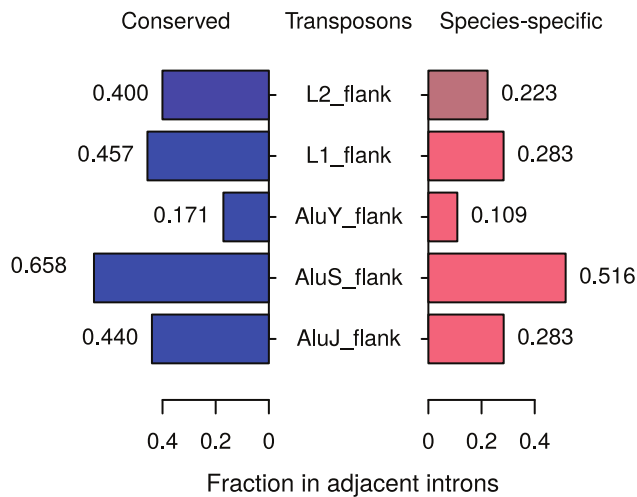


Figure 3

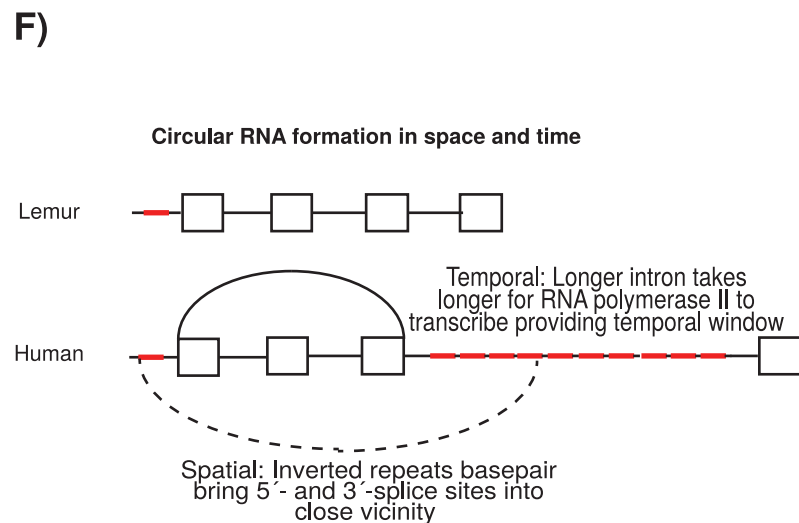
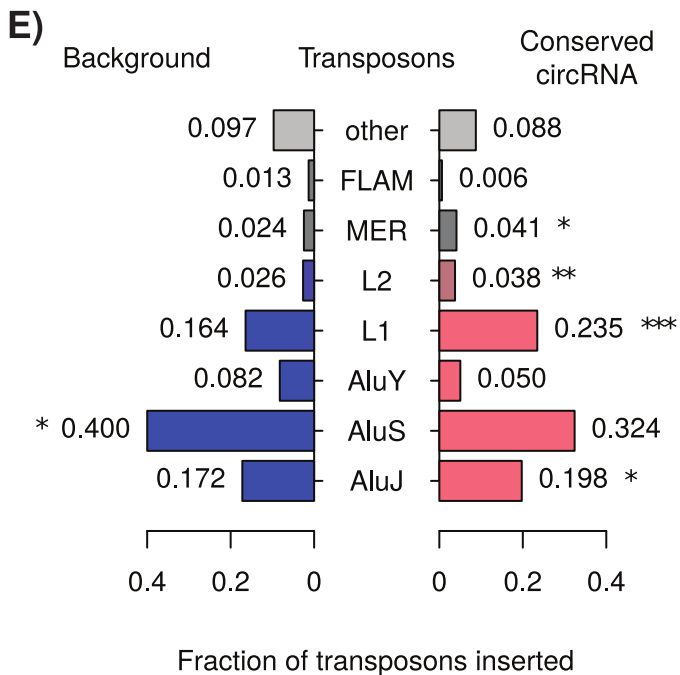
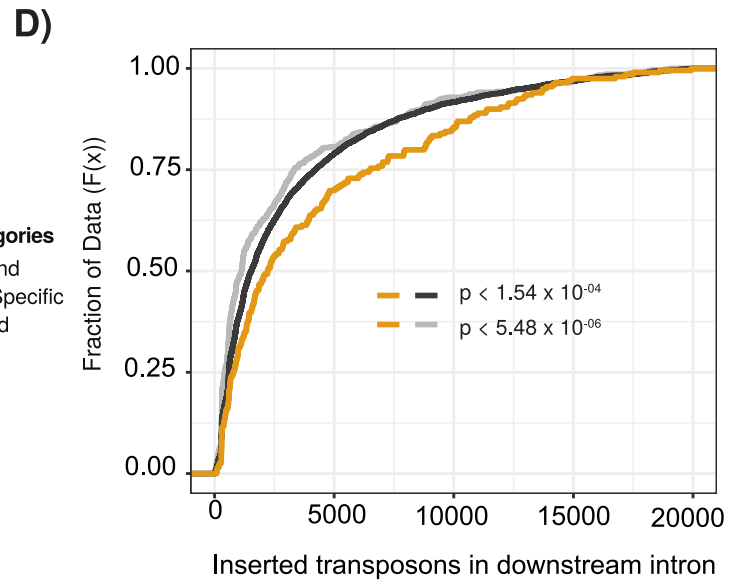
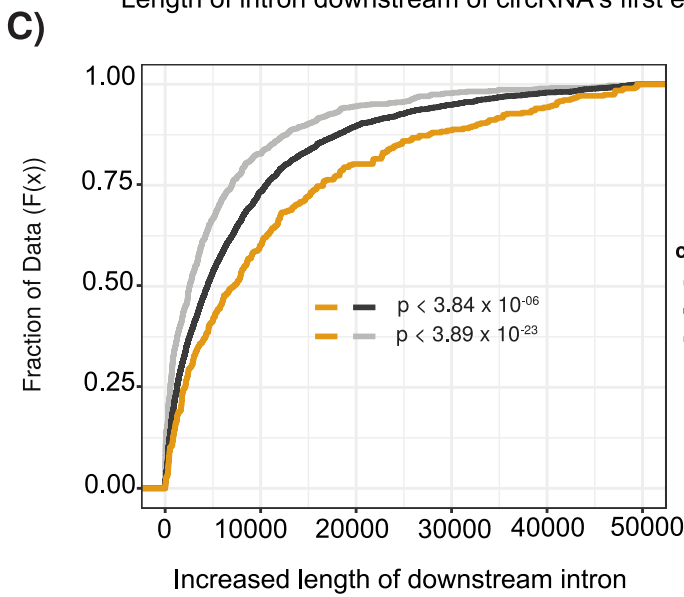
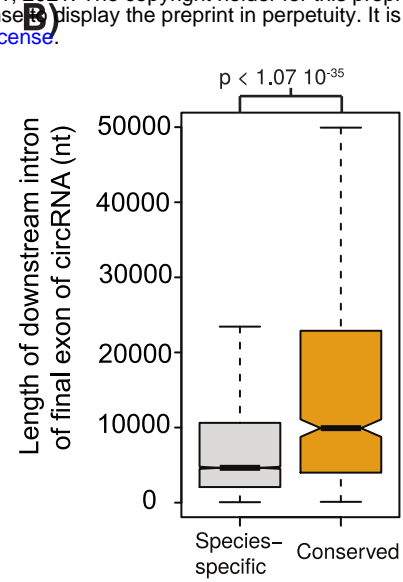
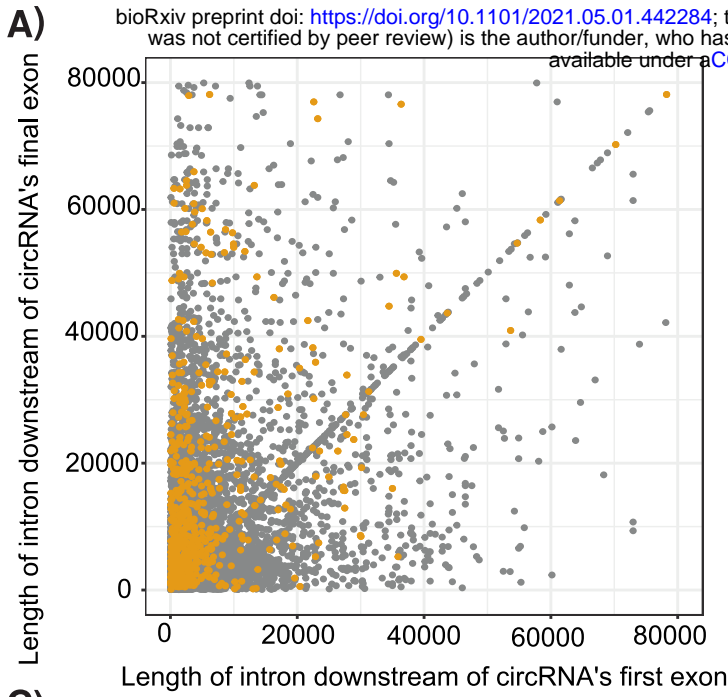
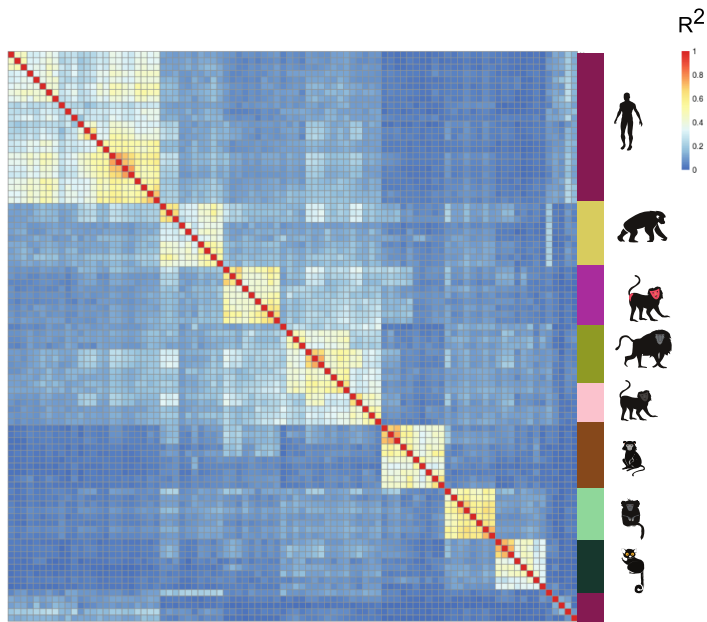
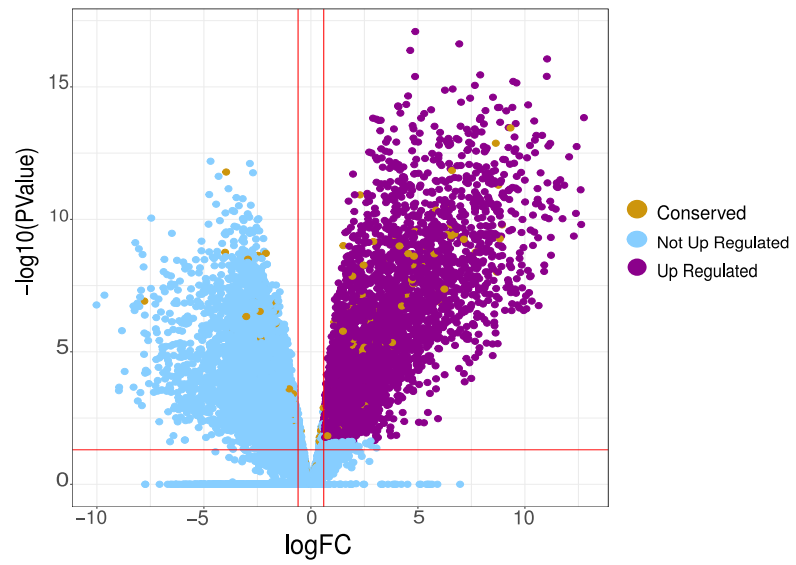


Figure 4

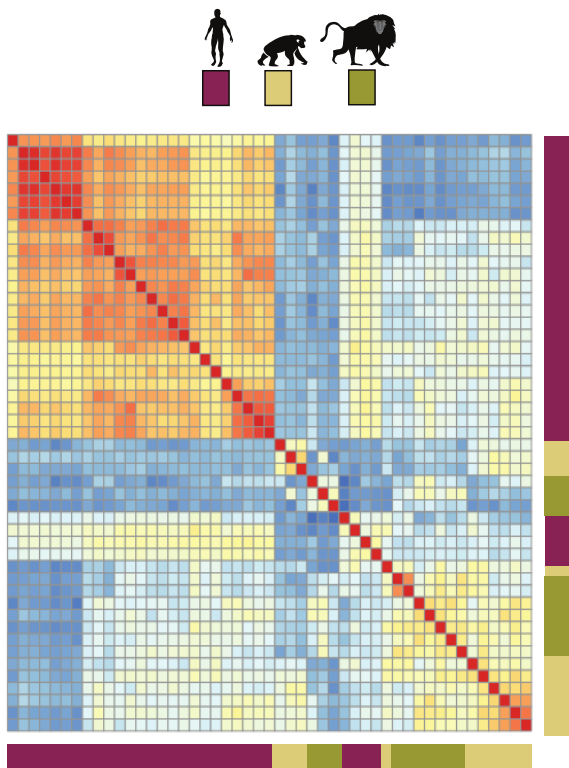
A)



B)



C)



D)

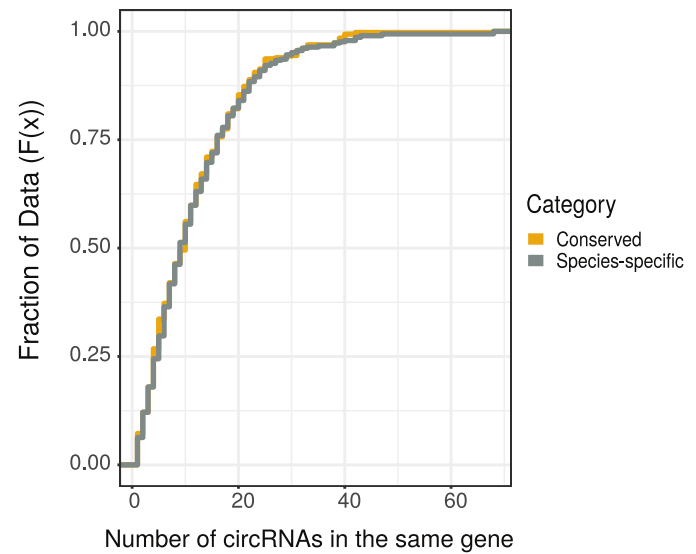
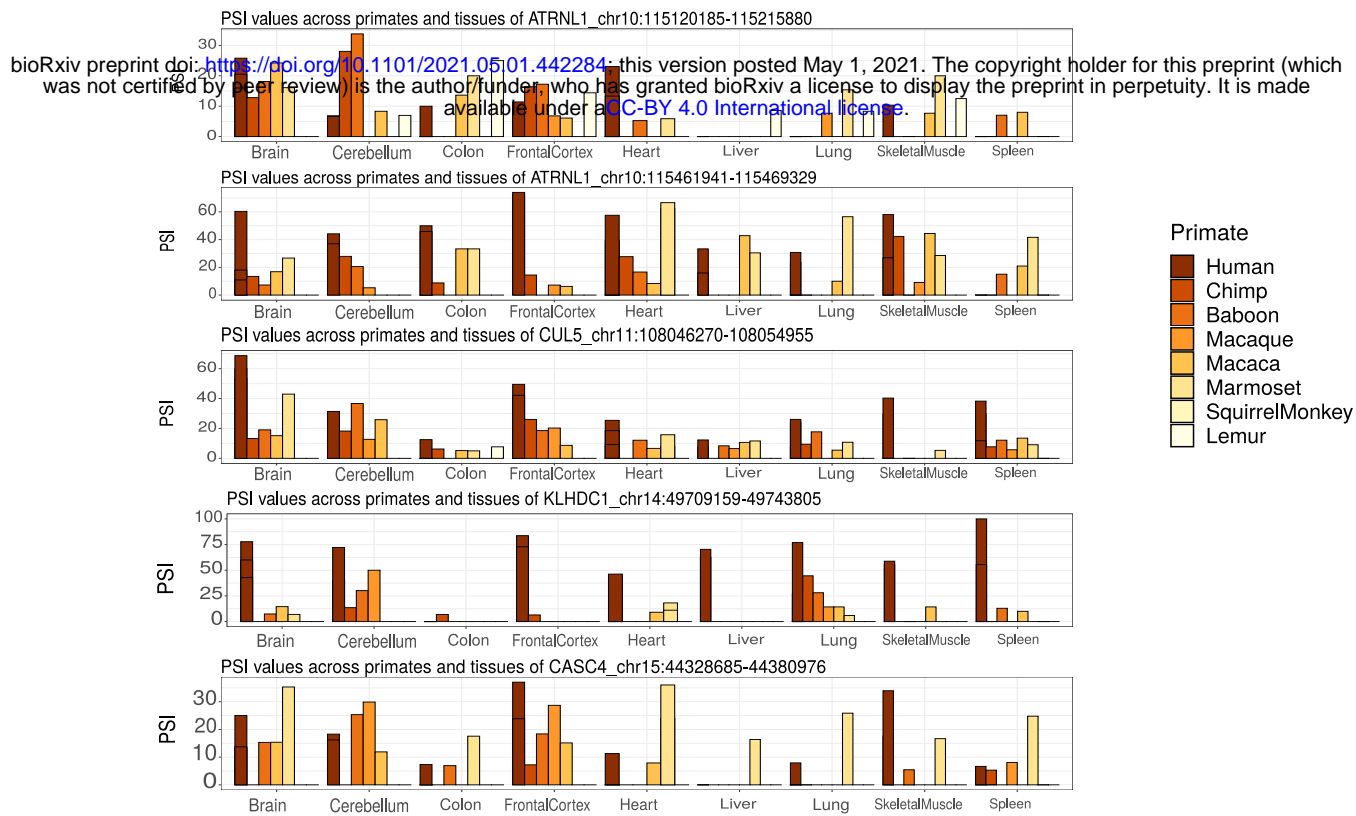
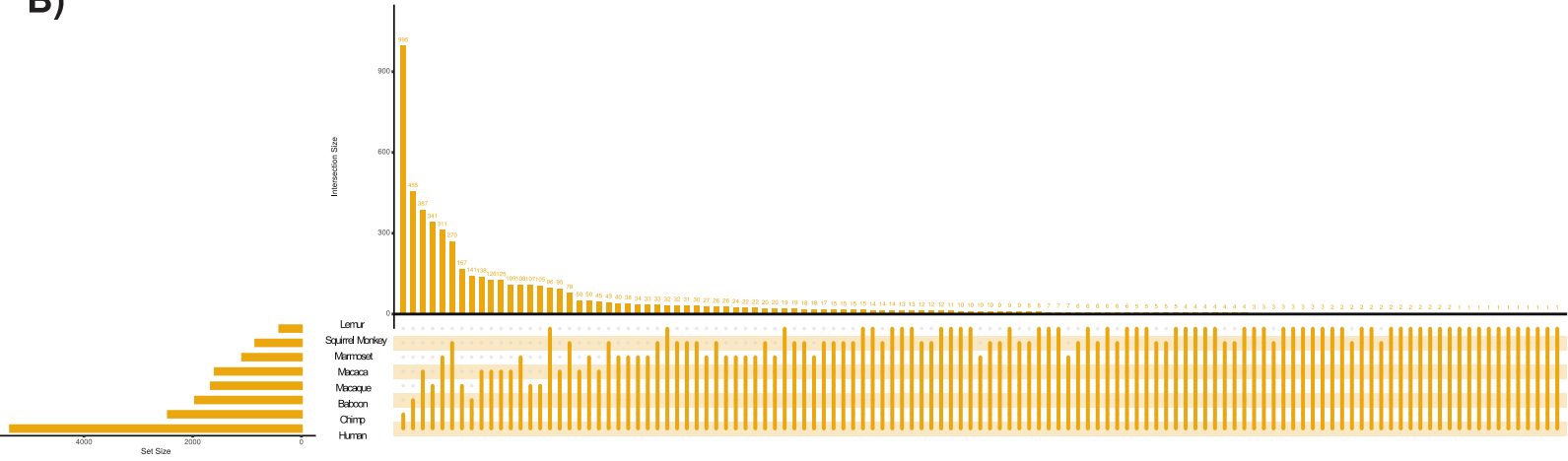


Figure 1-Figure supplement 1

A)



B)



C)

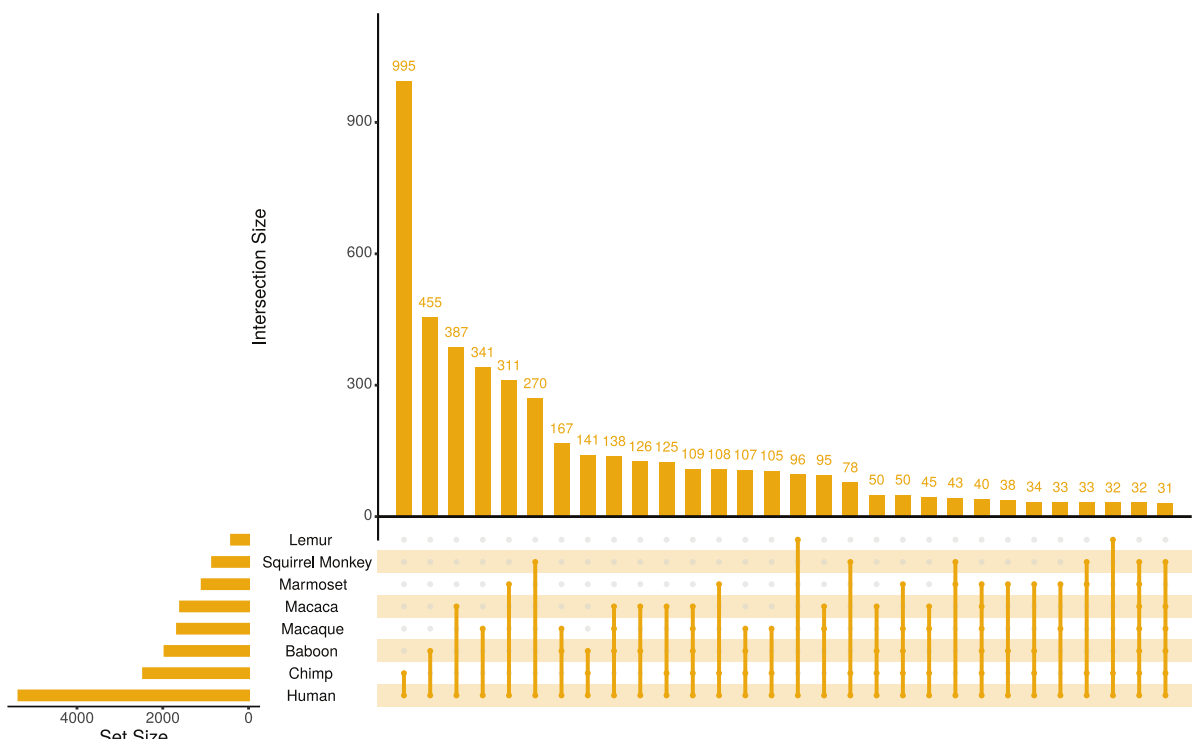


Figure 2- Figure supplement 1

A)

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.01.442284>; this version posted May 1, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY 4.0 International license](https://creativecommons.org/licenses/by/4.0/).

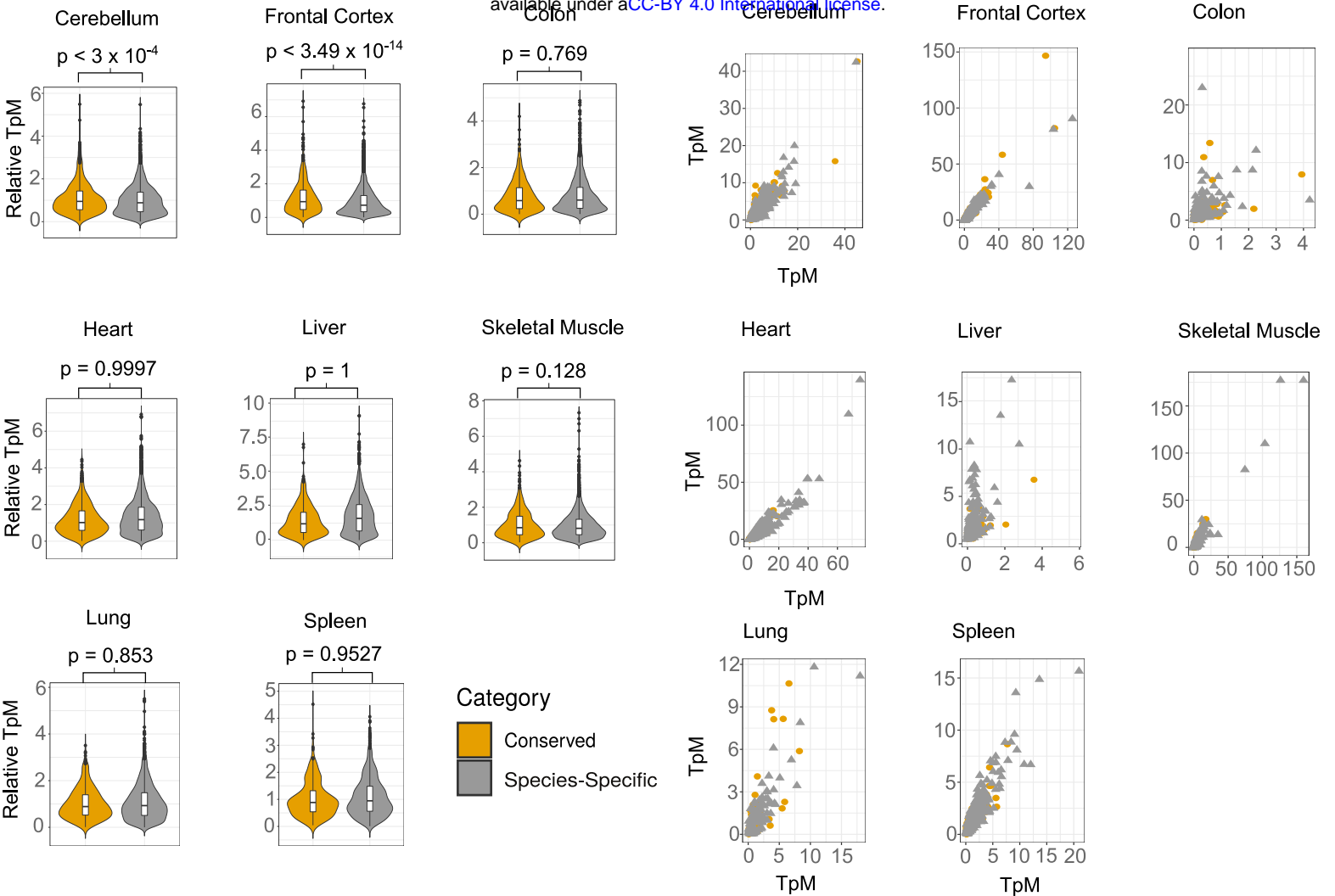
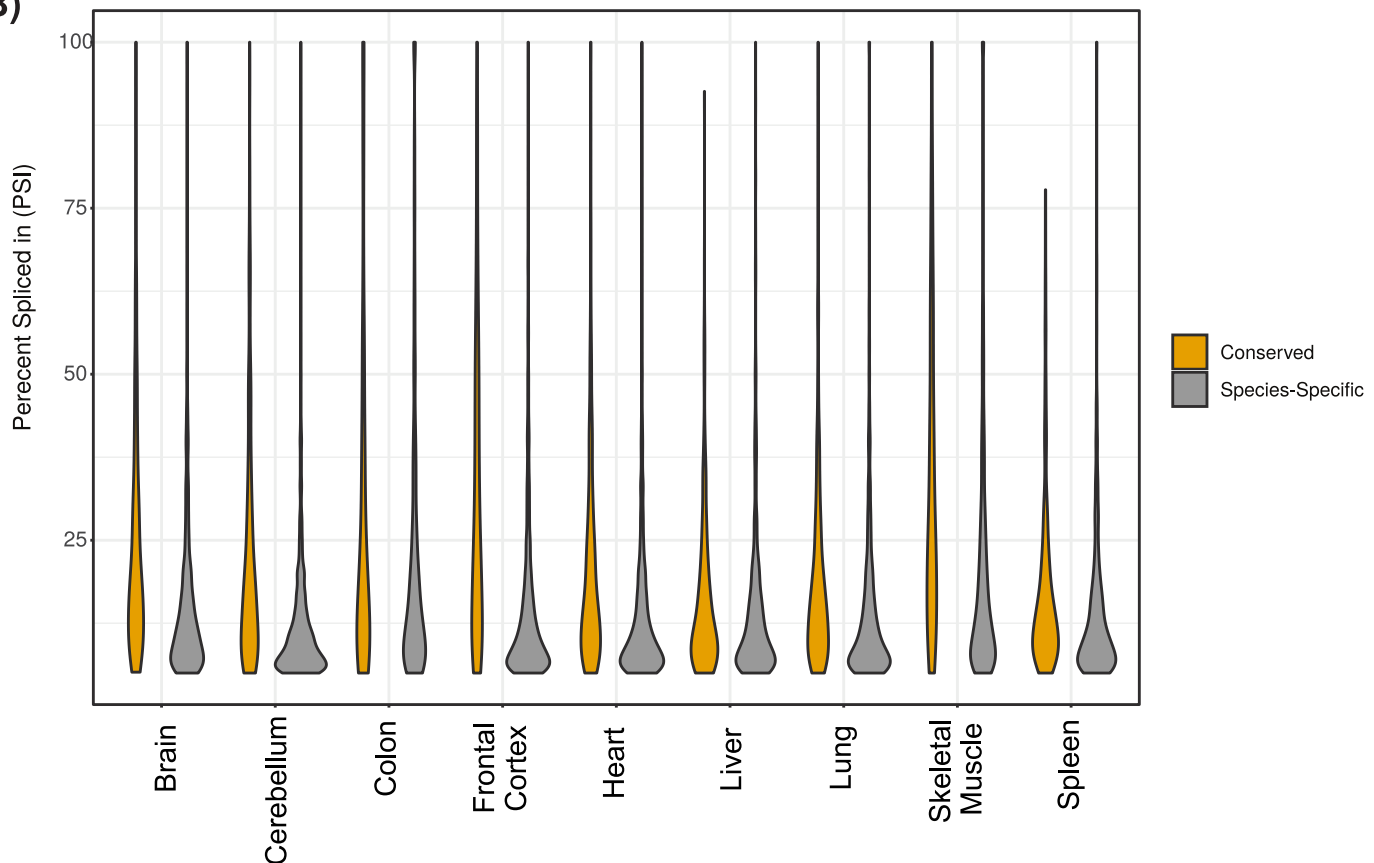
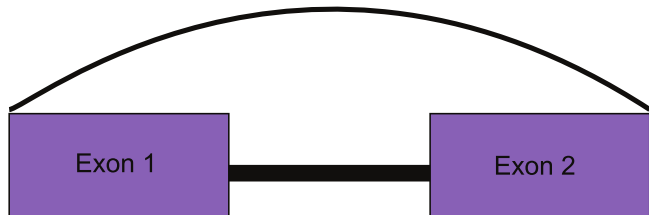
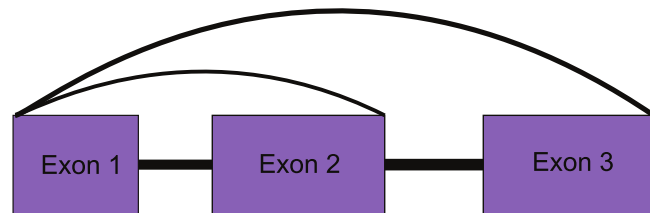
**B)**

Figure 2-Figure supplement 2

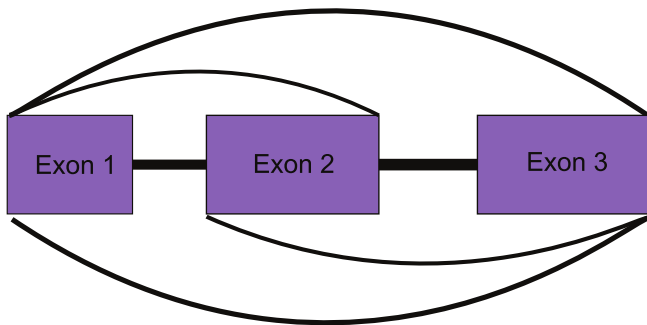
Unique Start and End



Repeated Start



Repeated Start and End



Repeated End

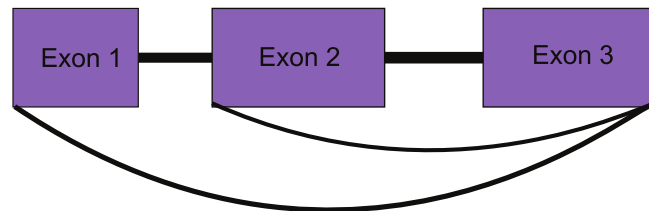


Figure 2-Figure supplement 3

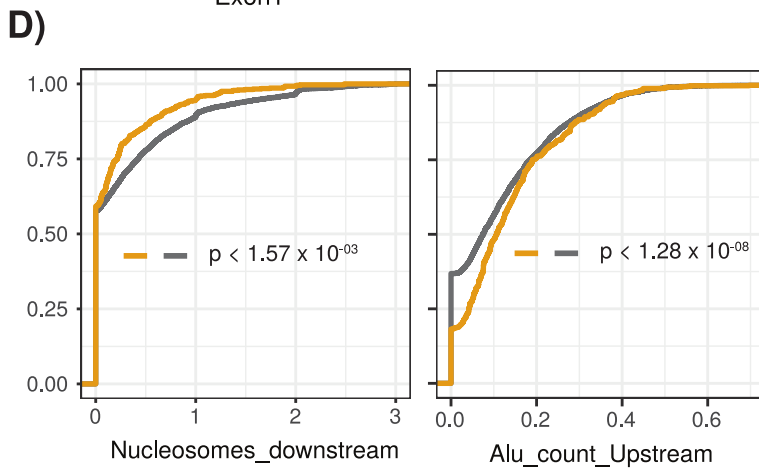
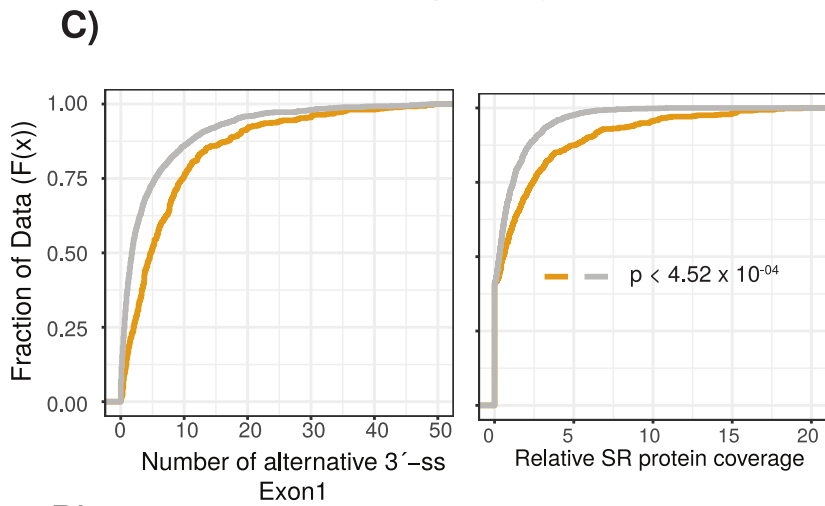
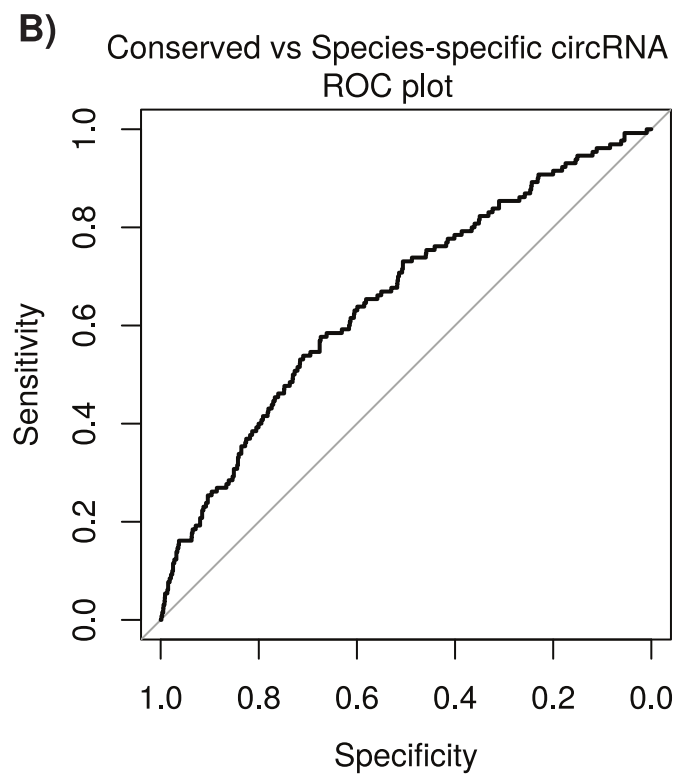
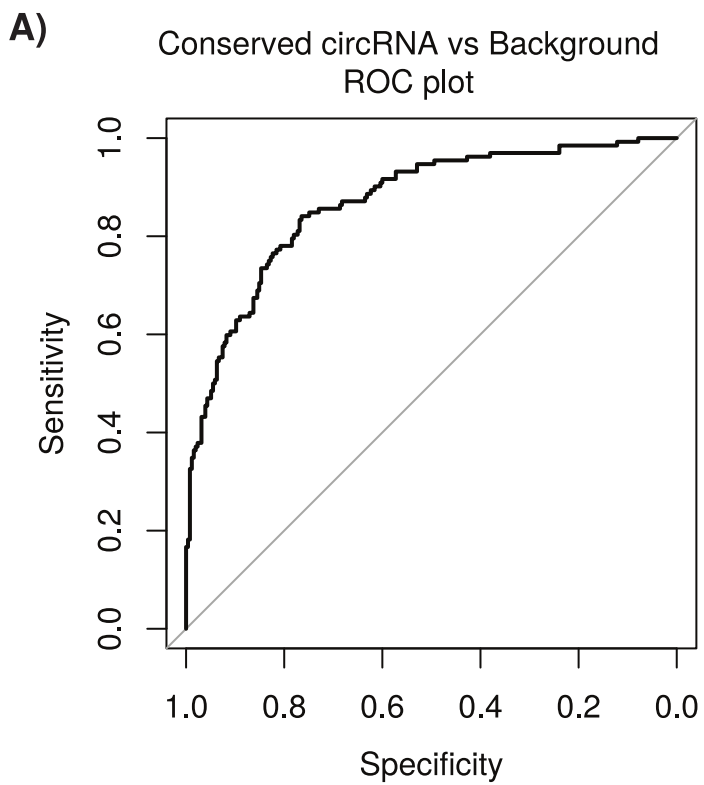


Figure 3 - Figure supplement 1