# Circular RNA repertoires are associated with evolutionarily young transposable elements

**Franziska Gruhl**[1,2‡]**, Peggy Janich**[1,3§]**, Henrik Kaessmann**[4*]**, David Gatfield**[1*]

**\*For correspondence:**
h.kaessmann@zmbh.uni-heidelberg.de (HK); david.gatfield@unil.ch (DG)

**Present address:** [†]SIB Swiss Institute of Bioinformatics, Lausanne, Switzerland; [‡]Krebsliga Schweiz, Bern, Switzerland

[1]Center for Integrative Genomics, University of Lausanne, Lausanne, Switzerland; [2]SIB Swiss Institute of Bioinformatics, Lausanne, Switzerland; [3]Krebsliga Schweiz, CH-3001 Bern, Switzerland; [4]Center for Molecular Biology of Heidelberg University (ZMBH), DKFZ-ZMBH Alliance, Heidelberg, Germany

**Abstract** Circular RNAs (circRNAs) are found across eukaryotes and can function in post-transcriptional gene regulation. Their biogenesis through a circle-forming backsplicing reaction is facilitated by reverse-complementary repetitive sequences promoting pre-mRNA folding. Orthologous genes from which circRNAs arise, overall contain more strongly conserved splice sites and exons than other genes, yet it remains unclear to what extent this conservation reflects purifying selection acting on the circRNAs themselves. Our analyses of circRNA repertoires across five species representing three mammalian lineages (marsupials, eutherians: rodents, primates) reveal that surprisingly few circRNAs arise from orthologous exonic loci across different species. Even the circRNAs from the orthologous loci are associated with young, recently active and species-specific transposable elements, rather than with common, ancient transposon integration events. These observations suggest that many circRNAs emerged convergently during evolution – as a byproduct of splicing in orthologs prone to transposable element insertion. Overall, our findings argue against widespread functional circRNA conservation.

## Introduction

First described more than forty years ago, circular RNAs (circRNAs) were originally perceived as a curiosity of gene expression, but they have gained significant prominence over the last 5-10 years (reviewed in *Kristensen et al.* (*2019*); *Patop et al.* (*2019*)). Large-scale sequencing efforts have facilitated the identification of thousands of individual circRNAs with specific expression patterns and, in some cases, specific functions (*Conn et al., 2015*; *Du et al., 2016*; *Hansen et al., 2013*; *Piwecka et al., 2017*). CircRNA biogenesis occurs through so-called "backsplicing" events, in which an exon's 3' splice site is ligated onto an upstream 5' splice site of an exon on the same RNA molecule (rather than downstream, as in conventional splicing). Backsplicing occurs co-transcriptionally and

is guided by the canonical splicing machinery (*Guo et al., 2014*; *Ashwal-Fluss et al., 2014*; *Starke et al., 2015*). It can be facilitated by complementary, repetitive sequences in the flanking introns (*Dubin et al., 1995*; *Jeck et al., 2013*; *Ashwal-Fluss et al., 2014*; *Zhang et al., 2014*; *Liang and Wilusz, 2014*; *Ivanov et al., 2015*). Through intramolecular base-pairing and folding, the resulting hairpin-like structures can augment backsplicing over the competing, regular forward-splicing reaction. In most cases, backsplicing seems to be rather inefficient, given that circRNA expression levels are low in most tissues. For example, it has been estimated that about 60% of circRNAs exhibit expression levels of less than 1 FPKM (fragments per kilobase per million reads mapped) - a commonly applied cut-off below which genes are usually considered to not be robustly expressed (*Guo et al., 2014*). Due to their circular structure, circRNAs are protected from the activity of cellular exonucleases, which is thought to favour their accumulation to detectable steady-state levels and, together with the cell's proliferation history, presumably contributes to their complex spatiotemporal expression patterns (*Alhasan et al., 2015*; *Memczak et al., 2013*; *Bachmayr-Heyda et al., 2015*). Overall higher circRNA abundances have been reported for neural tissues (*Westholm et al., 2014*; *Gruner et al., 2016*; *Rybak-Wolf et al., 2015*) and during ageing (*Gruner et al., 2016*; *Xu et al., 2018*; *Cortés-López et al., 2018*).

CircRNAs are found in all eukaryotes (protists, fungi, plants, animals) (*Wang et al., 2014*). More-over, it has been reported that circRNAs are frequently generated from orthologous genomic regions across species such as mouse, pig and human (*Rybak-Wolf et al., 2015*; *Venø et al., 2015*), and that their splice sites have elevated conservation scores (*You et al., 2015*). In these studies, circRNA coordinates were transferred between species to identify "conserved" circRNAs. However, the analyses did not distinguish between potential selective constraints actually acting on the cir-cRNAs themselves, from those preserving canonical splicing features of genes in which they are formed (so-called "parental genes"). A further obstacle to a thorough evolutionary understanding lies in the observation that while long introns containing reverse complementary repeats seem to be a conserved feature of circRNA parental genes, the reverse complementary repeat sequences as such undergo rapid evolutionary changes (*Zhang et al., 2014*; *Rybak-Wolf et al., 2015*). Finally, concrete examples for experimentally validated, functionally conserved circRNAs are still scarce. At least in part, the reason may lie in the difficulty to specifically target circular vs. linear transcript isoforms in loss-of-function experiments; only recently, novel dedicated tools for such experiments have been developed (*Li et al., 2020*). At the moment, however, the prevalence of conserved and hence likely functional circRNAs remains overall unclear.

Here, we set out to investigate the origins and evolution of circRNAs as well as potentially as-sociated selective pressures. To this end, we generated a comprehensive set of circRNA-enriched RNA sequencing (RNA-seq) data from five mammalian species and three organs. Our analyses un-veil that circRNAs are typically generated from a distinct class of genes that share characteristic structural and sequence features. Notably, we discovered that circRNAs are flanked by species-specific and recently active transposable elements (TEs). Our findings support a model according to which the integration of TEs is preferred in introns of genes with similar genomic properties, thus facilitating circRNA formation as a byproduct of splicing around the same exons of ortholo-gous genes across different species. Together, our work suggests that most circRNAs - even when

74 occurring in orthologs of multiple species and comprising the same exons - do nevertheless not

75 trace back to common ancestral circRNAs but emerged convergently during evolution, facilitated

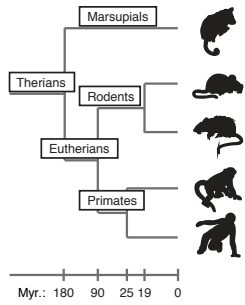76 by independent TE insertion events.

## Results

### A comprehensive circRNA dataset across five mammalian species

79 To explore the origins and evolution of circRNAs, we generated paired-end RNA-seq data for three

80 organs (liver, cerebellum, testis) in five species (grey short-tailed opossum, mouse, rat, rhesus

81 macaque, human) representing three mammalian lineages (marsupials; eutherians: rodents, pri-

82 mates) with different divergence times (**Figure 1A**, **Figure 1-Figure supplement 1A**, **Supplemen-**

83 **tary Table 1**). To enrich for circRNAs, samples were treated with exoribonuclease (RNase R) prior

84 to library preparation and sequencing. Using a custom pipeline, we subsequently identified cir-

85 cRNAs from backsplice junction (BSJ) reads, estimated circRNA steady-state abundances, and re-

86 constructed their isoforms (**Supplementary Table 2**, **Figure 1-Figure supplement 1B**, **Figure 1-**

87 **Figure supplement 2**). In total, we identified 1,535 circRNAs in opossum, 1,484 in mouse, 2,038

88 in rat, 3,300 in rhesus macaque, and 4,491 circRNAs in human, with overall higher numbers in

89 cerebellum, followed by testis and liver (**Figure 1B**, **Supplementary Table 3**). Detected circRNAs

90 were generally small in size, overlapped with protein-coding exons, showed considerable tissue-

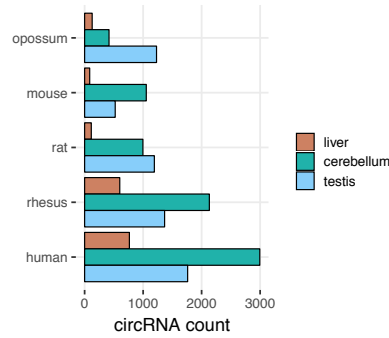91 specificity, and were flanked by large introns (**Figure 1-Figure supplement 3**).

### The identification of circRNA heterogeneity and hotspot frequency is determined by sequencing depth and detection thresholds

94 A sizeable number of genes give rise to multiple, distinct circRNAs (*Venø et al., 2015*). Such "cir-

95 cRNA hotspots" are of particular interest as they may be enriched for genomic features that drive

96 circRNA biogenesis. A previous hotspot definition applied a cutoff of at least 10 structurally differ-

97 ent, yet overlapping circRNAs produced from a genomic locus (*Venø et al., 2015*). However, given

98 that reaching a threshold of 10, or any other threshold, of detectable circRNA species for a given

99 locus likely strongly depends on the sequencing depth and the applied CPM (counts per million)

100 threshold, we compared circRNA hotspots identified at different CPM thresholds (0.1, 0.05 and

101 0.01 CPM). Moreover, to globally capture circRNA hotspot complexity, we considered genomic loci

102 already as hotspots if they produced as a minimum two different, overlapping circRNAs at a given

103 threshold. As expected, the number of hotspots - and the number of circRNAs these hotspots give

104 rise to - strongly depend on the chosen CPM threshold (**Figure 1C** for human and rhesus macaque

105 data; **Figure 1-Figure supplement 4** for other species). Thus, at 0.1 CPM only 16-27% of all de-

106 tected circRNA loci are classified as hotspots. Decreasing the stringency to 0.01 CPM increases the

107 proportion of hotspot loci to 32-45%. At the same time, the fraction of all circRNAs that originated

108 from hotspots increased from 34-49% (0.1 CPM) to 59-76% (0.01 CPM), and the number of circR-

109 NAs per hotspot increased from 2 to 6. Together, these observations suggest that at lower CPM

110 thresholds, it is in particular the number of circRNAs per locus that increases, whereas the effect

111 on the number of detectable independent circRNA loci is smaller. Furthermore, we observed that

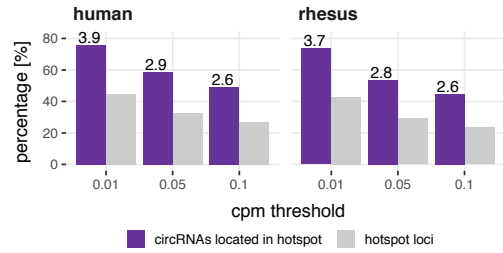112 in many cases the same hotspots produced circRNAs across multiple organs (**Figure 1D**), and that
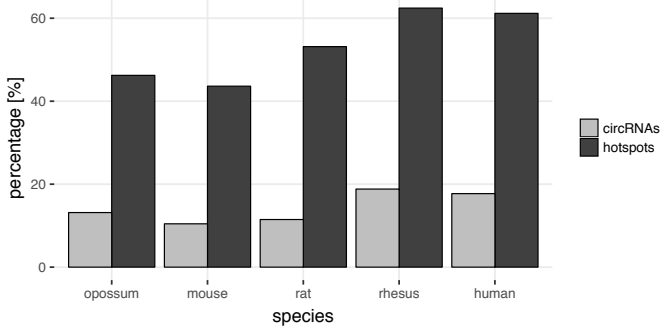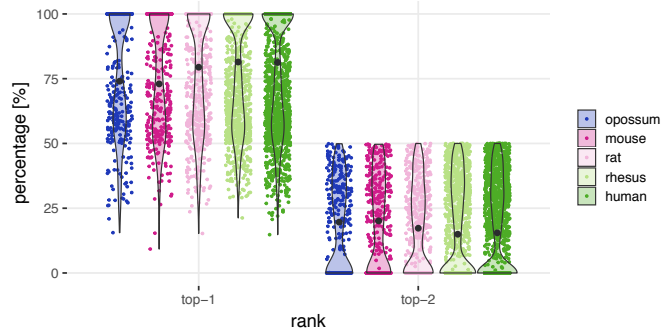
**A: Species and phylogeny**



**B: Number of circRNAs**
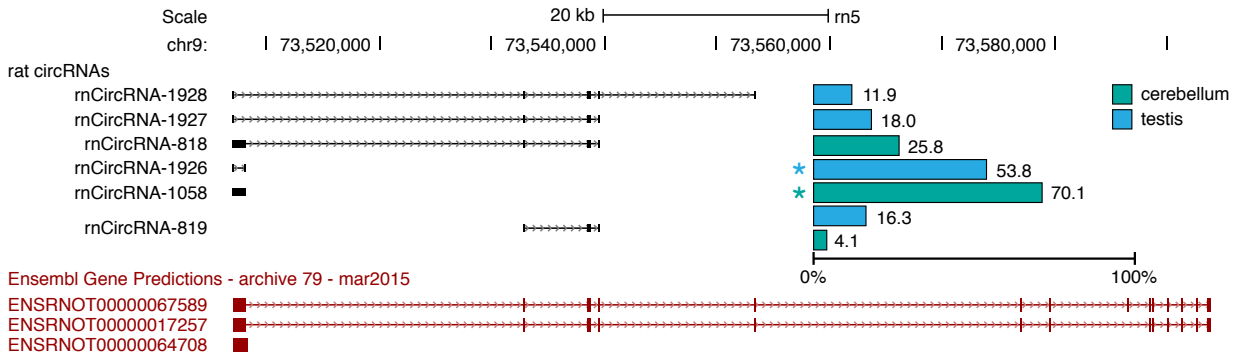


**C: CircRNA hotspot loci by CPM**



**D: % of circRNAs and hotspots in multiple tissues**



**E: CircRNA expression strength in hotspots**



**F: UCSC genome browser view for *Kansl1l* hotspot in rat**



there is usually one predominantly expressed circRNA per organ (**Figure 1E**). The *Kansl1l* hotspot

locus is a representative example: it is a hotspot in rat, where it produces 6 different circRNAs (de-

tails in **Figure 1F**); it is also a hotspot in all other species (producing 8, 5, 7, and 6 different circRNAs

in opossum, mouse, rhesus macaque and human, respectively; data not shown).

The substantial increase in circRNA heterogeneity with decreasing CPM, as well as the overall

low expression levels of many circRNAs, raised the question to what extent the majority of detected

circRNAs in this and other studies reflect a form of gene expression noise rather than functional

transcriptome diversity.

**Figure 1.** A: Phylogenetic tree of species analysed in this study. CircRNAs were identified and analysed in five mammalian species (opossum, mouse, rat, rhesus macaque, human) and three organs (liver, cerebellum, testis). B: Total number of detected circRNAs across species and tissues. The total number of circRNAs for each species in liver (brown), cerebellum (green) and testis (blue). C: CircRNA hotspot loci by CPM (human and rhesus macaque). The graph shows, in grey, the proportion (%) of circRNA loci that qualify as hotspots and, in purple, the proportion (%) of circRNAs that originate from such hotspots, at three different CPM thresholds (0.01, 0.05, 0.1). The average number of circRNAs per hotspot is indicated above the purple bars. D. Number of circRNA hotspot loci found in multiple tissues. The graph shows the proportion (%) of circRNAs (light grey) and of hotspots (dark grey) that are present in at least two tissues. E. Contribution of top-1 and top-2 expressed circRNAs to overall circRNA expression from hotspots. The plot shows the contribution (%) that the two most highly expressed circRNAs (indicated as top-1 and top-2) make to the total circRNA expression from a given hotspot. For each plot, the median is indicated with a grey point. F. Example of the *Kansl1l* hotspot in rat. The proportion (%) for each detected circRNA within the hotspot and tissue (cerebellum = green, testis = blue) are shown. The strongest circRNA is indicated by an asterisk. rnCircRNA-819 is expressed in testis and cerebellum.

**Figure 1–Figure supplement 1.** Overview of the dataset and the reconstruction pipeline.

**Figure 1–Figure supplement 2.** Mapping summary of RNA-seq reads.

**Figure 1–Figure supplement 3.** General circRNA properties.

**Figure 1–Figure supplement 4.** CircRNA hotspot loci by CPM (opossum, mouse, rat).

---

### CircRNAs formed in orthologous loci across species preferentially comprise constitutive exons

123 We therefore sought to assess the selective preservation – and hence potential functionality – of
124 circRNAs. For each gene, we first collapsed circRNA coordinates to identify the maximal genomic
125 locus from which circRNAs can be produced (**Figure 2A**). In total, we annotated 5,428 circRNA loci
126 across all species (**Figure 2A**). The majority of loci are species-specific (4,103 loci; corresponding to
127 75.6% of all annotated loci), whereas there are only comparatively few instances where circRNAs
128 arise from orthologous loci in the different species (i.e., from loci that share orthologous exons in
129 corresponding 1:1 orthologous genes; **Figure 2A**). For example, only 260 orthologous loci (4.8%
130 of all loci) give rise to circRNAs in all five species (**Figure 2A**). A considerable proportion of these
131 shared loci also correspond to circRNA hotspots (opossum: 30.0%, mouse: 25%, rat: 32.3%, rhe-
132 sus macaque: 44.6%, human: 60.4%). Thus, despite applying circRNA enrichment strategies in
133 library preparation and lenient thresholds for computational detection, the number of potentially
134 conserved orthologous circRNAs is surprisingly low.

135     PhastCons conservation scores are based on multiple alignments and known phylogenies, de-
136 scribing the conservation levels at single-nucleotide resolution (*Siepel et al., 2005*). To assess
137 whether circRNA exons differed from non-circRNA exons in their conservation levels, we calculated
138 phastCons scores for different exon types (circRNA exons, non-circRNA exons and UTR-exons). Cir-
139 cRNA exons showed higher phastCons scores in comparison to exons from the same genes that
140 were not spliced into circRNAs (**Figure 2B**), which would be the expected outcome if purifying se-
141 lection acted on functionally conserved circRNAs. However, other mechanisms may be relevant as
142 well; constitutive exons, for example, generally exhibit higher conservation scores than alternative
143 exons (*Modrek and Lee, 2003*; *Ermakova et al., 2006*). We thus analysed exon features in more
144 detail. First, the comparison of phastCons scores between exons of non-parental genes, parental
145 genes and circRNAs revealed that parental genes were *per se* highly conserved (**Figure 2B**): 85-
146 95% of the observed median differences between circRNA exons and non-parental genes could

147 be explained by the parental gene itself. Next, we compared the usage of parental gene exons

148 across organs (**Figure 2C**). We observed that circRNA exons are more frequently used in isoforms

149 expressed in multiple organs than non-circRNA parental gene exons. Finally, we analysed the se-

150 quence composition at the splice sites, which revealed that GC amplitudes (i.e., the differences in

151 GC content at the exon-intron boundary) are significantly higher for circRNA-internal exons than

152 for parental gene exons that were located outside of circRNAs (**Figure 2D**).

153 Collectively, these observations (i.e., increased phastCons scores, expression in multiple tissues,

154 increased GC amplitudes) raise the question of whether the above "circRNA-specific" exon proper-

155 ties (**Figure 2B-D**) primarily reflect an enrichment for constitutive exons. Under this scenario, the

156 supposed high conservation of circRNAs may not be directly associated with the circRNAs them-

157 selves, but with constitutive exons that the circRNAs contain. Together with the small proportion

158 of circRNAs with shared (orthologous) locations across species (see above), this raises the possibil-

159 ity that circRNAs are overall not highly conserved and that many circRNAs "shared" across species

160 (i.e., those arising from orthologous exonic loci) are actually not homologous. That is, rather than

161 reflecting (divergent) evolution from common ancestral circRNAs (**Figure 2E, left panel**), circRNAs

162 may frequently have emerged independently (convergently) during evolution in the lineages lead-

163 ing to the different species, thus potentially often representing "analogous" transcriptional traits
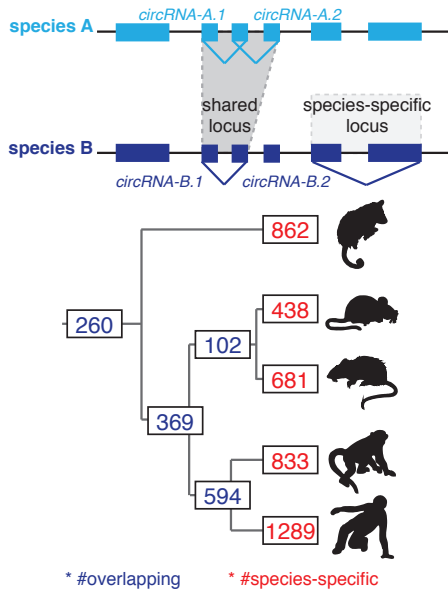
164 (**Figure 2E, right panel**).

## CircRNA parental genes are characterised by low GC content and high sequence repetitiveness

167 To explore whether convergent evolution played a role in the origination of circRNAs, we set out

168 to identify possible structural and/or functional constraints that may establish a specific genomic

169 environment (a "parental gene niche") potentially favouring analogous circRNA production. To this

170 end, we compared GC content and sequence repetitiveness of circRNA parental vs. non-parental
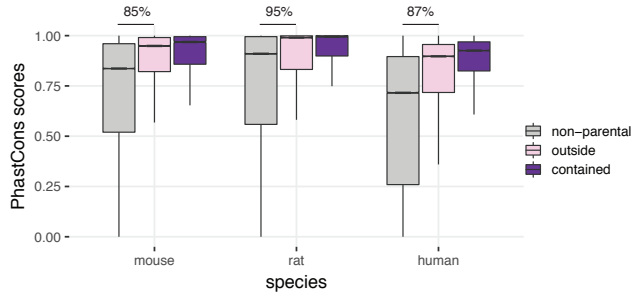
171 genes.

172 GC content is an important genomic sequence characteristic associated with distinct patterns

173 of gene structure, splicing and function (*Amit et al., 2012*). We realised that the increased GC ampli-

174 tude at circRNA exon-intron boundaries that we noted above (**Figure 2D**), was mainly caused by a

175 local decrease of intronic GC content rather than an increase in exonic GC content (**Supplementary**

176 **Table 4**, **Figure 2-Figure supplement 2**). We hence hypothesised that GC content may be a means

177 of discriminating parental from non-parental genes. We grouped genes into five categories from

178 low (L) to high (H) GC content (isochores; L1 <37%, L2 37-42%, H1 42-47%, H2 47-52% and H3 >52%

179 GC content) (**Figure 3A**). Coding genes in rhesus macaque and human are characterised by a bi-

180 modal GC content distribution (see peaks in L2 and H3 for non-parental genes). By contrast, the

181 two rodents displayed a unimodal distribution (peak in H1), whereas opossum coding genes were

182 generally GC-poor (in agreement with *Galtier and Mouchiroud* (*1998*); *Mikkelsen et al.* (*2007*)). No-

183 tably, circRNA parental genes showed a distinctly different distribution than non-parental genes

184 and a consistent pattern across all five species, with the majority of genes (82-94% depending on

185 species) distributing to the GC-low gene groups, L1 and L2 (**Figure 3A**).

186 We next analysed intron repetitiveness – a structural feature that has previously been asso-
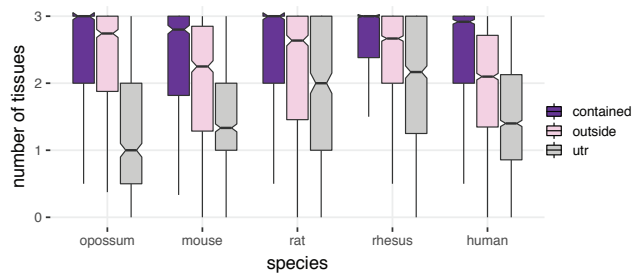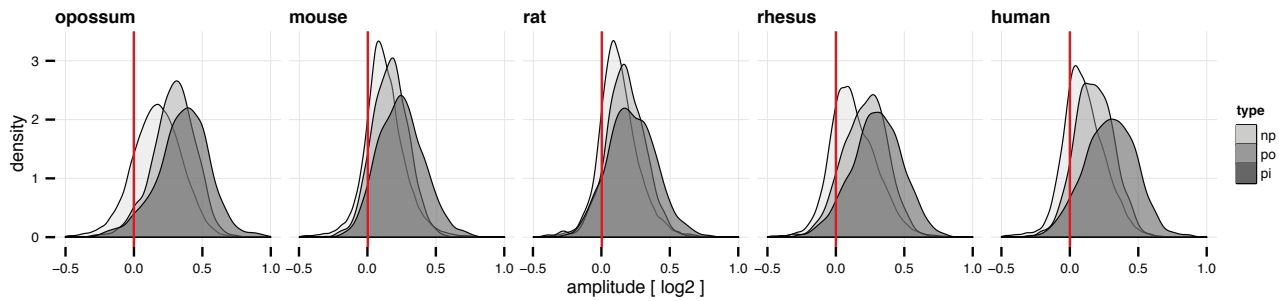
**A: Overlap of collapsed circRNA loci**



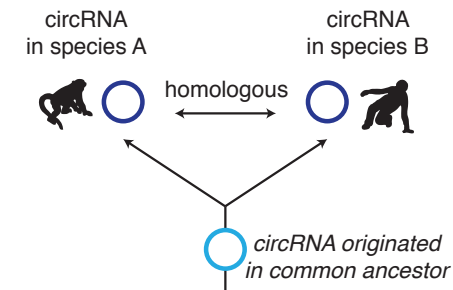**B: PhastCons score by exon type**



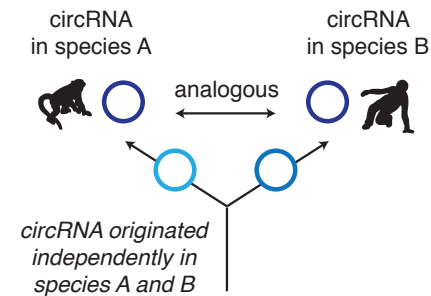**C: Tissue frequency of exon types**



**D: Splice site amplitude**



**E: Alternative models for the evolution of overlapping circRNA loci**

divergent evolution

circRNA in species A ⟷ homologous ⟷ circRNA in species B

circRNA originated in common ancestor

-> circRNAs loci overlap, because they evolved from a common ancestor

convergent evolution

circRNA in species A ⟷ analogous ⟷ circRNA in species B

circRNA originated independently in species A and B

-> circRNA loci overlap, because of similar genomic constraints

187 ciated with circRNA biogenesis. We used megaBLAST to align all annotated coding genes with

188 themselves to identify regions of complementarity in the sense and antisense orientations of the

189 gene (reverse complement sequences, RVCs) (*Ivanov et al., 2015*). We then compared the level

190 of self-complementarity between parental and non-parental genes within the same isochore (i.e.,

191 per gene group with the same GC content), given that self-complementarity generally shows neg-

**Figure 2.** Evolutionary properties of circRNAs. A: CircRNA loci overlap between species. Upper panel: Schematic representation of the orthology definition used in our study. CircRNAs were collapsed for each gene, and coordinates were lifted across species. Lower panel: Number of circRNA loci that are species-specific (red) or circRNAs that arise from orthologous exonic loci of 1:1 orthologous genes (i.e., circRNAs sharing 1:1 orthologous exons) across lineages (purple) are counted. We note that in the literature, other circRNA "orthology" definitions can be found, too. For example, assigning circRNA orthology simply based on parental gene orthology implies calling also those circRNAs "orthologous" that do not share any orthologous exons, which directly argues against the notion of circRNA homology; that is, a common evolutionary origin (see **Figure 2-Figure supplement 1**). Overall, the orthology considerations we applied largely follow the ideas sketched out in *Patop et al.* (*2019*). B: Distribution of phastCons scores for different exon types. PhastCons scores were calculated for each exon using the conservation files provided by ensembl. PhastCons scores for non-parental exons (grey), exons in parental genes, but outside of the circRNA (pink) and circRNA exons (purple) are plotted. The difference between circRNA exons and non-parental exons that can be explained by parental non-circRNA exons is indicated above the plot. C: Mean tissue frequency of different exon types in parental genes. The frequency of UTR exons (grey), non-UTR exons outside of the circRNA (pink) and circRNA exons (purple) that occur in one, two or three tissues was calculated for each parental gene. D: Distribution of splice site amplitudes for different exon types. Distribution of median splice site GC amplitude (log2-transformed) is plotted for different exon types (np = non-parental, po = parental, but outside of circRNA, pi = parental and inside circRNA). Red vertical bars indicate values at which exon and intron GC content would be equal E: Different evolutionary models explaining the origins of overlapping circRNA loci.

**Figure 2–Figure supplement 1.** CircRNA loci overlap between species.

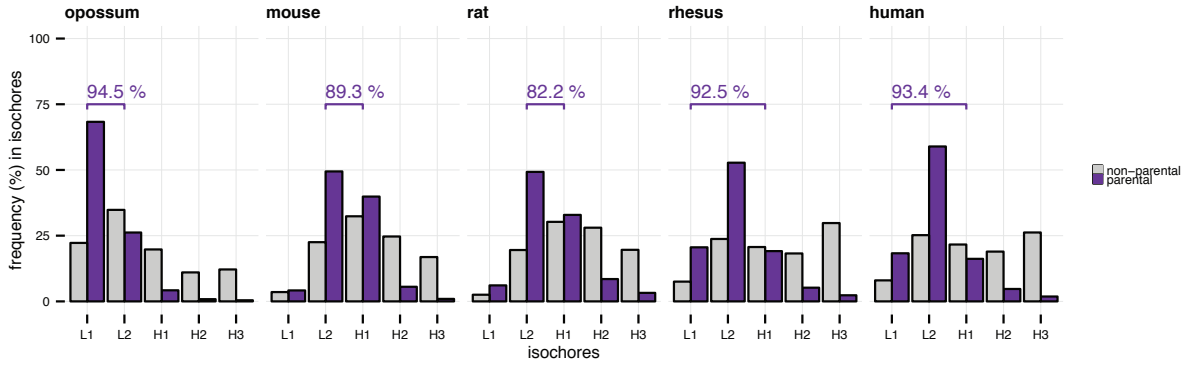**Figure 2–Figure supplement 2.** Amplitude correlations.

---

192 ative correlations with GC-content. This analysis revealed a stronger level of self-complementarity

193 in sense and antisense for parental genes than for non-parental genes from the same isochore

194 (**Figure 3B**).

195     CircRNA parental genes may also show an association with specific functional properties. Using

196 data from three human cell studies (*Steinberg et al., 2015*; *Pai et al., 2012*; *Koren et al., 2012*), our

197 analyses revealed that circRNA parental genes are biased towards early replicating genes, showed

198 higher steady-state expression levels, and are characterised by increased haploinsufficiency scores

199 (**Figure 3-Figure supplement 1**). Collectively, we conclude that circRNA parental genes exhibit not

200 only distinct structural features (low GC content, high repetitiveness), but also specific functional

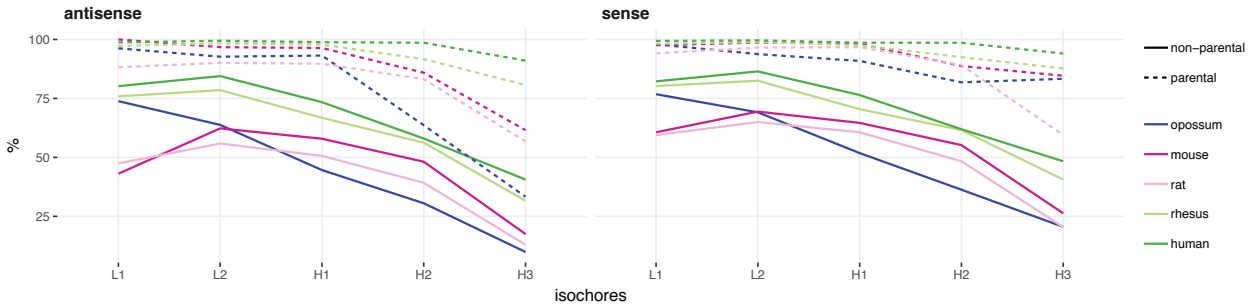201 properties associated with important roles in human cells.

202 **Among the multiple predictors of circRNA parental genes, low GC content distin-**

203 **guishes circRNA hotspots**

204 The aforementioned analyses established that circRNA parental genes possess distinct sequence,

205 conservation and functional features. Using linear regression analyses, we next sought to deter-

206 mine which of these properties constitute the main predictors of parental genes. Our model used

207 parental vs. non-parental gene as the response variable and several plausible explanatory vari-

208 ables (i.e., GC content, exon and transcript counts, genomic length, number of repeat fragments

209 in sense/antisense, expression level, phastCons score, tissue specificity index). After training the

210 model on a data subset (80%), circRNA parental gene predictions were carried out on the remainder

211 of the dataset (20%) (see **Material and Methods** for more information). Notably, predictions oc-

212 curred with high precision (accuracy 72-79%, sensitivity of 75%, specificity 71-79% across all species)

213 and uncovered several significantly associated features (**Table 1**, **Supplementary Table 5**, **Figure**

214 **3-Figure supplement 2**). Consistently for all species, the main parental gene predictors are low

215 GC content (log-odds ratio -1.84 to -0.72) and increased number of exons in the gene (log-odds
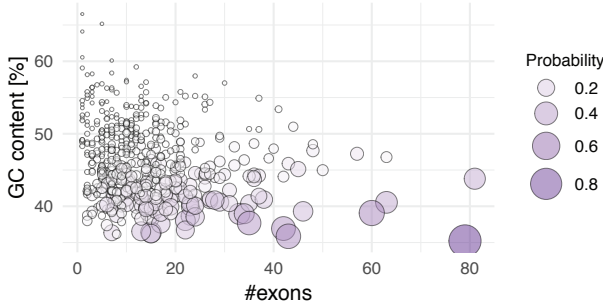
**A: GC content of parental genes**

**B: Complementarity in coding genes**

**C: GC content vs. exon count**

**D: PhastCons score vs. RVCs**

**E: Model of circRNA niche**

circRNAs are present in orthologous parental genes with similar properties

repeat age in circRNA niche to distinguish between models of divergent and convergent evolution

**divergent evolution**

common repeats between shared circRNA loci

**convergent evolution**

species-specific repeats between shared circRNA loci

**Figure 3.** A: GC content of parental genes. Coding genes were classified into L1-H3 based on their GC content, separately for non-parental (grey) and parental genes (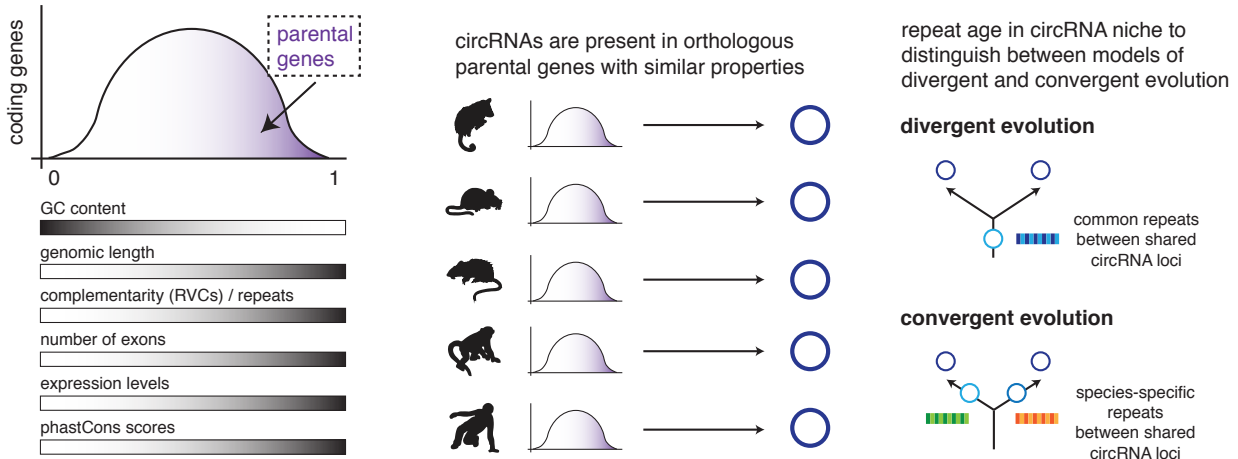purple). The percentage of parental genes in L1-L2 (opossum, mouse, rat) and L1-H1 (rhesus macaque, human) is indicated above the respective graphs. B: Complementarity in coding genes. Each coding gene was aligned to itself in sense and antisense orientation using megaBLAST. The proportion of each gene involved in an alignment was calculated and plotted against its isochore. C-D: Examples of parental gene predictors for linear regression models. A generalised linear model (GLM) was fitted to predict the probability of the murine coding gene to be parental, whereby x- and y-axis represent the strongest predictors. Colour and size of the discs correspond to the p-values obtained for 500 genes randomly chosen from all mouse coding genes used in the GLM. E. Model of circRNA niche.

**Figure 3–Figure supplement 1.** Replication time, gene expression steady-state levels and GHIS of human parental genes.

**Figure 3–Figure supplement 2.** Distribution of prediction values for non-parental and parental circRNA genes.

**Figure 3–Figure supplement 3.** Validation of parental gene GLM on Werfel *et al.* dataset.

216 ratio 0.30 to 0.45). Furthermore, increased genomic length (log-odds ratio 0.17 to 0.26) and an in-
217 creased proportion of reverse-complementary areas (repeat fragments) within the gene (log-odds
218 ratio 0.20 to 0.59), increased expression levels (log-odds ratio 0.25 to 0.38) and higher phastCons
219 scores (log-odds ratio 0.45 to 0.58) are also positively associated with circRNA production (**Table 1**,
220 **Figure 3C-D**, **Supplementary Table 5**). Notably, these circRNA parental gene predictors were not
221 restricted to our datasets but could be deduced from independent circRNA datasets as well. Thus,
222 the analysis of mouse and human heart tissue data (*Werfel et al., 2016*) revealed the same proper-
223 ties; that is, circRNA parental genes are characterised by low GC content, they were exon-rich, and
224 they showed enrichment for repeats (**Figure 3-Figure supplement 3**). Moreover, our linear regres-
225 sion models performed with comparable accuracy (74%), sensitivity (75%) and specificity (74%) to
226 predict parental genes in the independent human and mouse data. We therefore conclude that
227 the identified properties likely represented generic characteristics of circRNA parental genes that
228 are suitable to distinguish them from non-parental genes.

229 A substantial amount of circRNAs are formed from circRNA hotspots (**Figure 1C**). We there-
230 fore asked whether among the distinct genomic features that our regression analysis identified as
231 characteristic of parental genes, some would be suitable to further distinguish hotspots. First, we
232 assessed whether hotspots were more likely to be shared between species than parental genes
233 producing only a single circRNA isoform. Notably, the applied regression model did not only de-
234 tect a positive correlation between the probability of a parental gene to be a hotspot and having
235 orthologous parental genes in multiple species, but log-odds ratios increased with the distance
236 and number of species across which the hotspot was shared (e.g., mouse: 0.29 for shared within
237 rodents, 0.67 for shared with eutherian species and 0.72 for shared within therian species; **Supple-**
238 **mentary Table 6**). Finally, we interrogated whether a particular feature would be able to specify
239 circRNA hotspots among parental genes. A single factor, low GC content, emerged as a consis-
240 tent predictor for circRNA hotspots among all circRNA-generating loci (**Supplementary Table 7**).
241 Not surprisingly, the predictive power was lower than that of the previous models discriminating
242 parental vs. non-parental genes, which had identified low GC content as well. These findings imply
243 that hotspots emerge across species in orthologous loci that offer similarly favourable conditions
244 for circRNA formation, including low GC content. Of note, the increased number of circRNAs that
245 become detectable when CPM thresholds are lowered (see above, **Figure 1C**), is also in agreement

**Table 1.** A generalised linear model was fitted to predict the probability of coding genes to be a parental gene ($n_{opossum}$=18,807, $n_{mouse}$=22,015, $n_{rat}$=11,654, $n_{rhesus}$=21,891, $n_{human}$=21,744). The model was trained on 80% of the data (scaled values, cross-validation, 1000 repetitions). Only the best predictors were kept and then used to predict probabilities for the remaining 20% of data points (validation set, shown in table). Genomic length, number of exons and GC content are based on the respective ensembl annotations; number of repeats in antisense and sense orientation to the gene was estimated using the RepeatMasker annotation, phastCons scores taken from UCSC (not available for opossum and rhesus macaque) and expression levels and the tissue specificity index based on (*Brawand et al., 2011*). An overview of all log-odds ratios and p-values calculated in the validation set of each species is provided in the table, further details can be found in **Supplementary Table 5**. *Abbreviations: md = opossum, mm = mouse, rn = rat, rm = rhesus macaque, hs = human. Significance levels: '\*\*\*' < 0.001, '\*\*' < 0.01, '\*' < 0.05, 'ns' >= 0.05.*

| Predictor | Log-odds range (significance) | Species with significant predictor |
| --- | --- | --- |
| Genomic gene length (bp) | rn: 0.26 (***)<br>rm: 0.17 (***)<br>hs: 0.26 (***)<br>md, mm: ns | rn, rm, hs |
| Number of exons | md: 0.45 (***)<br>mm: 0.38 (***)<br>rn: 0.30 (***)<br>rm: 0.42 (***)<br>hs: 0.32 (***) | md, mm, rn, rm, hs |
| GC content | md: -1.84 (***)<br>mm: -1.09 (***)<br>rn: -0.72 (***)<br>rm: -1.44 (***)<br>hs: -1.42 (***) | md, mm, rn, rm, hs |
| Repeat fragments (antisense) | md: 0.28 (**)<br>mm: 0.20 (**)<br>rm: 0.59 (***)<br>rn, hs: ns | md, mm, rm |
| Repeat fragments (sense) | hs: 0.58 (***)<br>md, mm, rn, rm: ns | hs |
| PhastCons scores | mm: 0.58 (***)<br>rn: 0.51 (***)<br>hs: 0.45 (***) | mm, rn, hs |
| Mean expression levels | md: 0.34 (**)<br>rm: 0.38 (***, )<br>hs: 0.25 (**)<br>mm, rn: ns | md, rm, hs |
| Tissue specificity index | md, mm, rn, rm, hs: ns | - |

246    with the sporadic formation of different circRNAs whenever genomic circumstances allow for it.
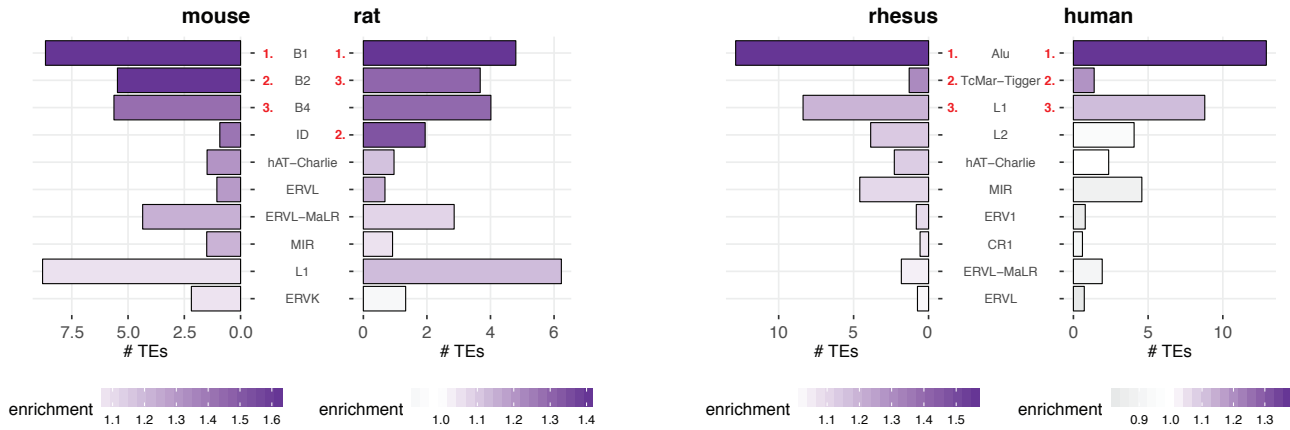
247  Collectively, our analyses thus reveal that circRNA parental genes are characterised by a set
248  of distinct features: low GC content, increased genomic length and number of exons, higher ex-
249  pression levels and increased phastCons scores (**Figure 3E**). These features were detected inde-
250  pendently across species, suggesting the presence of a unique, syntenic genomic niche in which
251  circRNAs can be produced ("circRNA niche"). While helpful to understand the genomic context of
252  circRNA production, these findings do not yet allow distinguishing between the two alternative
253  models of divergent and convergent circRNA evolution (**Figure 2E**). However, we reasoned that
254  this aim would be in reach if we better understood the evolutionary trajectory and timeline that
255  leads to the emergence of the circRNAs. Conceivably, the identified feature "complementarity and
256  repetitiveness" of the circRNA niche might give access to this time component. Previous studies
257  have associated repetitiveness with an over-representation of small TEs – such as primate Alu el-
258  ements or the murine B1 elements – in circRNA-flanking introns; these TEs may facilitate circRNA
259  formation by providing RVCs that are the basis for intramolecular base-pairing of nascent RNA
260  molecules (*Ivanov et al., 2015*; *Jeck et al., 2013*; *Zhang et al., 2014*; *Wilusz, 2015*; *Liang and Wilusz,*
261  *2014*). Interestingly, while the biogenesis of human circRNAs has so far been mainly associated
262  with the primate-specific group of Alu elements, a recent study has highlighted several circRNAs
263  that rely on the presence of mammalian MIR elements (*Yoshimoto et al., 2020*). A better under-
264  standing of the evolutionary age of TEs in circRNA-flanking introns could thus provide important
265  insights into the modes of circRNA emergence; that is, the presence of common (i.e., old) repeats
266  would point towards divergent evolution of circRNAs from a common circRNA ancestor, whereas
267  an over-representation of species-specific (i.e., recent) repeats would support the notion of con-
268  vergent circRNA evolution (**Figure 3E**).

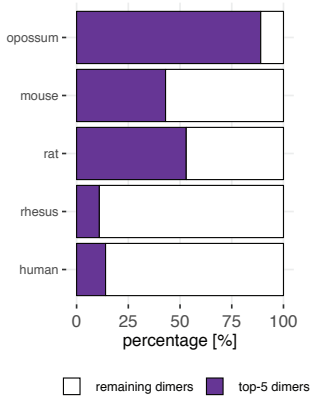269  ## CircRNA flanking introns are enriched in species-specific TEs

270  To assess potential roles of TEs in circRNA evolution, we first investigated the properties and com-
271  position of the repeat landscape relevant for circRNA biogenesis - features that have remained
272  poorly characterised so far - harnessing our cross-species dataset. As a first step, we generated
273  for each species a background set of "control introns" from non-circRNA genes that were matched
274  to the circRNA flanking introns in terms of length distribution and GC content. We then compared
275  the abundance of different repeat families within the two intron groups. In all species, TEs belong-
276  ing to the class of small, interspersed nuclear elements (SINEs) are enriched within the circRNA
277  flanking introns as compared to the control introns. Remarkably, the resulting TE enrichment pro-
278  files were exquisitely lineage-specific, and even largely species-specific (**Figure 4A**). In mouse, for
279  instance, the order of enrichment is from the B1 class of rodent-specific B elements (strongest
280  enrichment and highest frequency of >7.5 TEs per flanking intron) to B2 and B4 SINEs. In rat, B1
281  (strong enrichment, yet less frequent than in mouse) is followed by ID (Identifier) elements, which
282  are a family of small TEs characterised by a recent, strong amplification history in the rat lineage
283  (*Kim et al., 1994*; *Kim and Deininger, 1996*); B2 and B4 SINEs only followed in 3$^{rd}$ and 4$^{th}$ position.
284  In rhesus macaque and human, Alu elements are the most frequent and strongly enriched TEs
285  (around 14 TEs per intron), consistent with the known strong amplification history in the common
286  primate ancestor (reviewed in *Batzer and Deininger* (*2002*)) (**Figure 4A**). The opossum genome is

**287** known for its high number of TEs, many of which may have undergone a very species-specific am-

**288** plification pattern (*Mikkelsen et al., 2007*), which is reflected in the distinct opossum enrichment

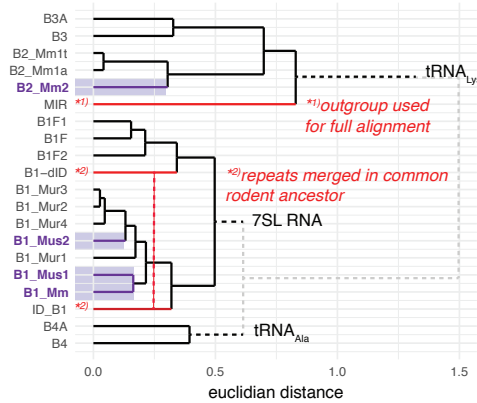**289** profile (**Figure 4-Figure supplement 1**).



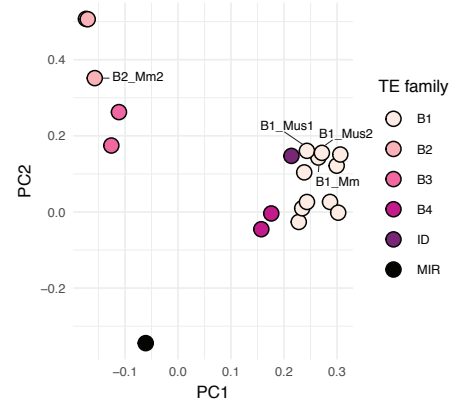**A: Enrichment of transposable elements in flanking introns**
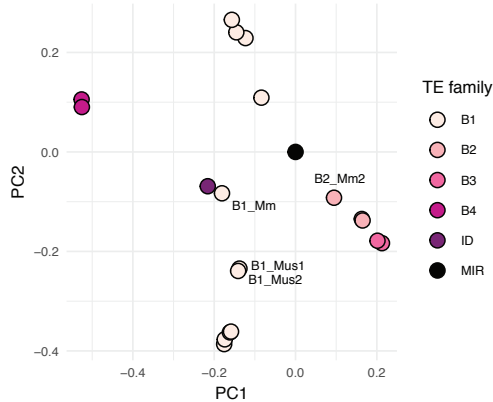


**B: Top–5 dimer contribution**



**C: Repeat phylogeny, mouse**



**D: TE distance matrix, mouse**



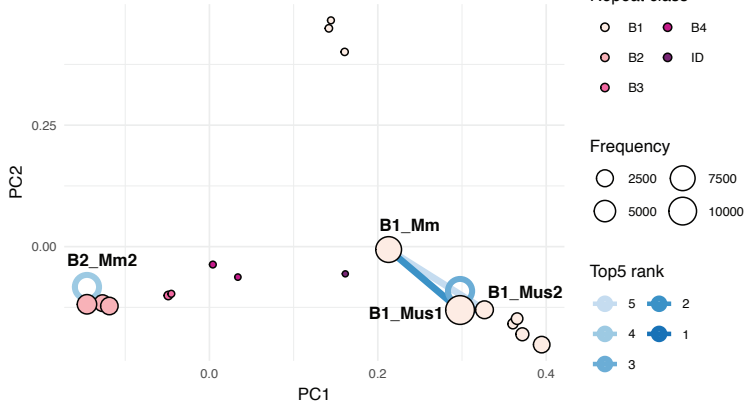**E: DeltaG dimers, mouse**



**F: Binding score, mouse**



**290** As pointed out above, TEs are relevant for circRNA formation because they can provide the RVCs

**291** that are the basis for intramolecular base-pairing of nascent RNA molecules (*Ivanov et al., 2015*;

**292** *Jeck et al., 2013*; *Zhang et al., 2014*; *Wilusz, 2015*; *Liang and Wilusz, 2014*). Folding of the pre-mRNA

**Figure 4.** A: Enrichment of TEs in flanking introns for mouse, rat, rhesus macaque and human. The number of TEs was quantified in both intron groups (circRNA flanking introns and length- and GC-matched control introns). Enrichment of TEs is represented by colour from high (dark purple) to low (grey). The red numbers next to the TE name indicate the top-3 enriched TEs in each species. B: Top-5 dimer contribution. The proportion of top-5 dimers (purple) to the remaining dimers (white) in flanking introns is shown. C: Phylogeny of mouse TEs. Clustal-alignment based on consensus sequences of TEs. Most recent TEs are highlighted. D: PCA for distance matrix of mouse TE families. PCA is based on the clustal-alignment distance matrix for the reference sequences of all major SINE families in mouse with the MIR family used as an outgroup. TEs present in the top-5 dimers are labelled. E: PCA based on deltaG for mouse TE families. PCA is based on the minimal free energy (deltaG) for all major SINE families in mouse with the MIR family used as an outgroup. TEs present in the top-5 dimers are labelled. F: PCA for binding score of mouse dimers. PCA is based on a merged and normalised score, taking into account binding strength (=deltaG) and phylogenetic distance. Absolute frequency of TEs is visualised by circle size. TEs present in the five most frequent dimers (top-5) are highlighted by blue lines connecting the two TEs engaged in a dimer (most frequent dimer in dark blue = rank 1). If the dimer is composed of the same TE family members, the blue line loops back to the TE (= blue circle).

**Figure 4–Figure supplement 1.** Enrichment of transposable elements in flanking introns for opossum.

**Figure 4–Figure supplement 2.** PCA and phylogeny of opossum, rat, rhesus macaque and human repeat dimers.

293 into a hairpin secondary structure with a paired RNA stem (formed by the flanking introns via the
294 dimerised RVCs) and an unpaired loop region (carrying the future circRNA) leads to a configuration
295 that is favourable for circRNA formation because it brings backsplice donor and acceptor sites into
296 close proximity. In order to serve as efficient RVCs via this mechanism, TEs will need to fulfil certain
297 criteria, and the dimerisation potential will likely depend on TE identity, frequency, and position.
298 Moreover, while two integration events involving the same TE (in reverse orientation) will lead to
299 an extended RVC stretch, different transposons from the same TE family also still share varying
300 degrees of sequence similarity that depend on their phylogenetic distance. The sequence differ-
301 ences that have evolved might compromise the base-pairing potential. To cover the dimerisation
302 potential of the TE landscape in a comprehensive fashion, we deemed it vital to calculate the actual
303 binding affinities between the dimerising sequences. As described below, we thus established a
304 binding score that would account for this variety of factors influencing dimer formation and that
305 would allow us to identify the TEs representing the most likely drivers of circRNA formation.

306 First, we noted that, similar to TEs overall (**Figure 4A**), RVCs were also enriched in SINE TEs
307 (**Figure 4B**). Moreover, in some species, relatively few specific dimers represented the majority of
308 all predicted dimers (i.e., top-5 dimers accounted for 89% of all dimers in flanking introns in opos-
309 sum, 43% in mouse, 53% in rat, 11% in rhesus and 14% in human). We further realised that the
310 phylogenetic distance between different TEs in a species was inadequate to categorise them with
311 regard to their dimer potential; as shown for mouse (**Figure 4C-D**), phylogenetic age only sepa-
312 rated large subgroups, but not TEs of the same family whose sequences have diverged by just a
313 few nucleotides. By contrast, classification by binding affinities creates more precise, smaller sub-
314 groups that lack, however, the information on phylogenetic age (**Figure 4E**). Therefore, we devised
315 a binding score that integrates both phylogeny (age) and binding affinity information (see **Material**
316 **and Methods**). Principal component analysis (PCA) showed that it efficiently separated different
317 TE families and individual family members, with PC1 and PC2 of the binding score explaining ap-
318 proximately 76% of observed variance (**Figure 4F**; **Figure 4-Figure supplement 2**). Moreover, this
319 analysis suggests that the most frequently occurring dimers (top-5 dimers are depicted as blue

320 connecting lines in **Figure 4F**) are formed by recently active TE family members. In mouse, an illus-
321 trative example are the dimers formed by the B1_Mm, B1_Mus1 and B1_Mus2 elements (**Figure 4F**),
322 which are among the most recent (and still active) TEs in this species (**Figure 4C**). Across species,
323 our analyses allowed for the same conclusions. For example, the dominant dimers in rat were
324 precisely the recently amplified ID elements, and not the more abundant (yet older in their am-
325 plification history) B1 family of TEs (**Figure 4-Figure supplement 2B**) (*Kim et al., 1994*; *Kim and*
326 *Deininger, 1996*). In opossum, the most prominent dimers consisted of opossum-specific SINE1
327 elements, which are similar to the Alu elements in primates, but possess an independent origin
328 (**Figure 4-Figure supplement 2A**) (*Gu et al., 2007*). Finally, dimer composition within the primate
329 lineage was relatively similar, probably due to the high amplification rate of AluJ and AluS/Z ele-
330 ments in the common primate ancestor and relatively recent divergence time of macaque and
331 human (**Figure 4-Figure supplement 2C-D**) (*Batzer and Deininger, 2002*).

332     In conclusion, the above analyses of RVCs revealed that dimer-forming sequences in circRNA
333 flanking introns were most frequently composed of recent, and often currently still active, TEs.
334 Therefore, the dimer repertoires were specific to the lineages (marsupials, rodents, primates) and/or
335 even – as most clearly visible within the rodent lineage – species-specific.

## Flanking introns of circRNA loci shared across species are enriched in evolutionar-
## ily young TEs

338 We next compared the dimer composition of the two groups of introns, namely those that flanked
339 circRNA loci whose exonic locations are in common between species and those that flanked species-
340 specific circRNA loci. For this analysis – aimed at finally resolving the extent to which circRNA loci
341 shared across species evolved from a common ancestor or independently from each other – we
342 took into account the degradation rate (milliDiv, see hereafter), frequency, enrichment and age of
343 the dimers. Briefly, the RepeatMasker annotations (*Smit et al., 2013*) (http://repeatmasker.org; see
344 **Material and Methods** for more details) provide a quantification of how many "base mismatches
345 in parts per thousand" have occurred between each specific repeat copy in its genomic context
346 and the repeat reference sequence. This deviation is expressed as the milliDiv value. Thus, a high
347 milliDiv value implies that a repeat is strongly degraded, typically due to its age (the older the re-
348 peat, the more time its sequence has had to diverge). Low milliDiv values suggest that the repeat is
349 younger (i.e., it had less time to accumulate mutations) or that purifying selection prevented the ac-
350 cumulation of mutations. Using this rationale, we explored degradation rates for the top-5 dimers
351 extracted in each species from the ensemble of parental genes, and then compared the milliDiv val-
352 ues associated with orthologous genes giving rise to circRNAs in multiple species (shared parental
353 genes) to those for species-specific parental genes. Notably, dimers detected in shared parental
354 genes are generally less degraded than those in species-specific parental genes (**Figure 5A**). In
355 rat, for example, median milliDiv values for the dimers involving young TE classes (ID_Rn1+ID_Rn1,
356 ID_Rn1+ID_Rn2 and ID_Rn2+ID_Rn2) range from 21 to 42.5 for shared parental genes and 26 to
357 43.5 for species-specific parental genes, with the differences between shared and species-specific
358 parental genes all being statistically significant (**Figure 5A, left panel**). By contrast, no signifi-
359 cant milliDiv differences were found in the case of dimers involving older repeats (BC1_Rn+ID_Rn1
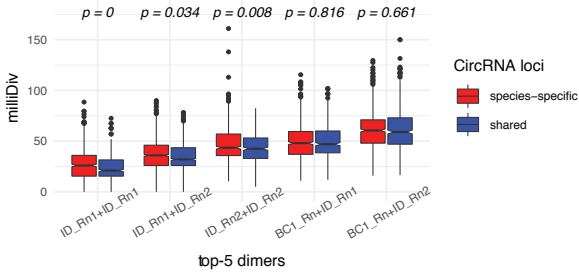
360  and BC1_Rn+ID_Rn2); thus, their degradation rates are comparable between shared and species-
361  specific parental genes. The human data (**Figure 5A, right panel**) and that from opossum, mouse
362  and macaque (**Figure 5-Figure supplement 1A**) revealed similar trends. For example, differences
363  in degradation rates between shared and human-specific parental genes were observed for the
364  dimers containing younger repeats such as AluSx1+AluY or AluSx+AluY (**Figure 5A, right panel**). In
365  conclusion, these analyses reveal that flanking introns of circRNAs are enriched in TEs with rather
366  species-specific integration and amplification rates, consistent with the idea of convergent circRNA
367  evolution driven by independent TE insertion events in orthologous genomic loci.

368      Low degradation rates could indicate that specific dimers are particularly important for the
369  production of functional circRNAs. For example, Alu elements, which the above dimer analyses
370  identified as important in human and rhesus macaque, are common to the primate lineage, and it
371  would be conceivable that the circRNA loci shared between both species emerged through TE in-
372  tegration in a common primate ancestor and were subsequently preserved by purifying selection.
373  Alternatively, differences in degradation rates may simply reflect the evolutionary age of integra-
374  tion events. In that case, we would predict that even though the circRNA parental genes are shared
375  between species, the enriched dimers would nevertheless stem from recent, independent integra-
376  tion events, rather than from ancestral, shared integration events. This scenario could occur if the
377  circRNA-producing genes were to act as "transposon sinks" that are prone to insertions of active
378  repeats due to specific features related to their sequence or structural architecture. To explore this
379  idea, we examined in greater detail the dimers in shared and species-specific parental genes. As
380  in our above analyses, we first created specific "dimer lists", this time restricted to the two groups
381  of parental genes (shared/species-specific circRNA loci); using the top-100 most and least enriched
382  dimers, we compared the enrichment factors and mean age (categorised for simplicity into four
383  groups: 1 = species-specific, 2 = lineage-specific, 3 = eutherian, 4 = therian). The analysis revealed
384  that the most enriched and most frequent dimers are consistently formed by the youngest ele-
385  ments in both groups of genes, and that the frequency distribution of the top-100 dimers was sig-
386  nificantly different between species (see **Figure 5B** for mouse and rhesus macaque; other species
387  in **Figure 5-Figure supplement 1B**). In rhesus macaque, for example, the most frequent dimers
388  included the Alu element AluYra, which is characteristic for this species and absent from the hu-
389  man lineage. A representative example for such a shared circRNA-generating locus with young,
390  species-specific repeats is the *Akt3* locus (**Figure 5C**). Although *Akt3* circRNAs are shared between
391  human (upper panel), mouse (middle panel) and opossum (lower panel), the dimer landscapes
392  (top-5 dimers are highlighted in the figure) are entirely specifies-specific.
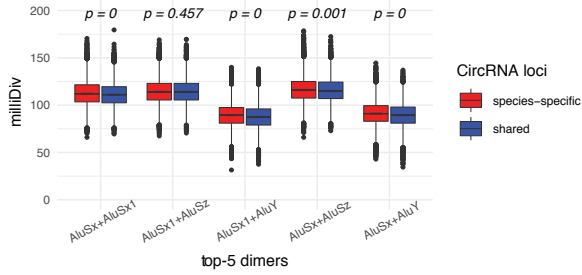
393      Taken together, we conclude that circRNAs are preferentially formed from loci that have ac-
394  quired TEs in recent evolutionary history. Such recent transposition events involved TEs that have
395  a higher degree of species-specificity than evolutionarily older TEs. Importantly, even in the case
396  of genomic loci whose capacity to generate circRNAs was shared across species, the actual repeat
397  landscapes revealed that they had acquired their TEs in evolutionarily recent times, as judged from
398  repeat degradation rates and age. Overall, these findings support a model according to which cir-
399  cRNAs are analogous, rather than homologous features of loci that have increased propensity of
400  attracting TEs, likely due to particular genomic features such as their GC content.
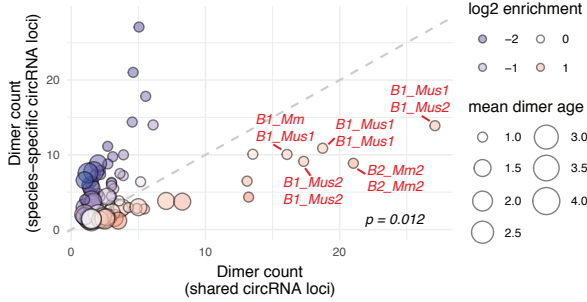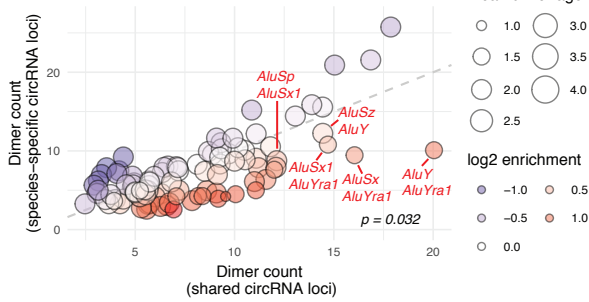
**A: MilliDivs for top-5 dimers**

**Rat**



**Human**
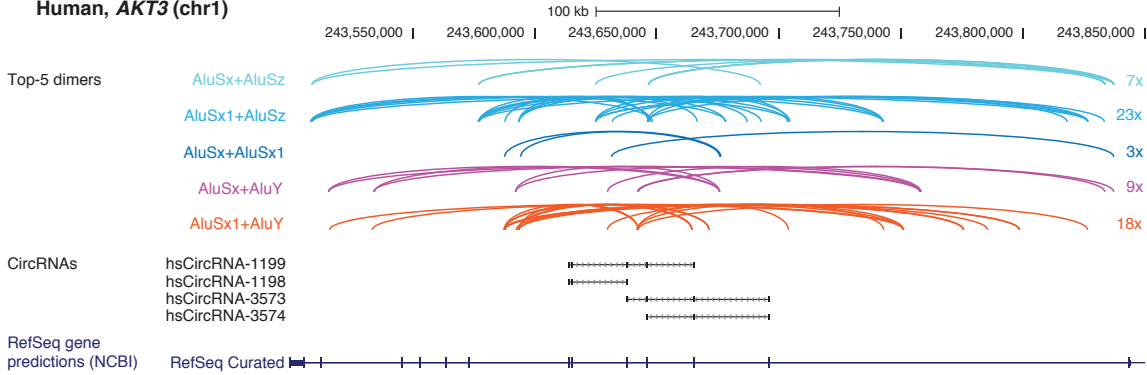


**B: Dimer enrichment (shared vs. species-specific circRNA loci)**

**Mouse**



**Rhesus**



**C: Examples of repeat landscape**

**Human, *AKT3* (chr1)**



**Mouse, *Akt3* (chr1)**
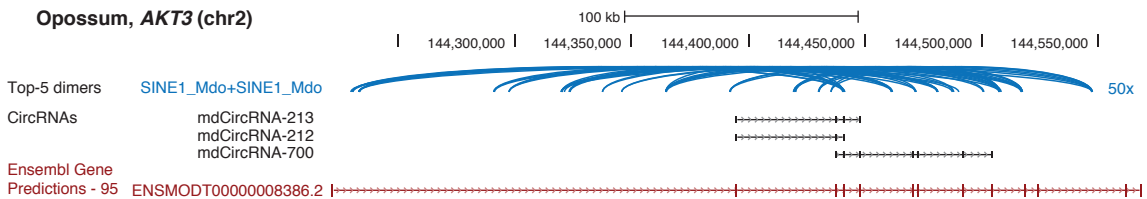


**Opossum, *AKT3* (chr2)**

**Figure 5.** A: Degradation rates (MilliDivs) for top-5 dimers in rat and human. MilliDiv values for the top-5 dimers (defined by their presence in all parental genes) were compared between parental genes of species-specific (red) and shared (blue) circRNA loci in rat and human. Since dimers are composed of two repeats, their mean value was taken. A t-test was used to compare dimers between parental genes with shared and species-specific circRNA loci, with p-values plotted above the boxplots. Dimer order from left to right on the x-axis corresponds to their rank in the top-5 list (most frequent left) B: Dimer enrichment in shared vs. species-specific repeats in mouse and rhesus macaque. The frequency (number of detected dimers in a given parental gene), log2-enrichment (shared vs. species-specific) and mean age (defined as whether repeats are species-specific: age = 1, lineage-specific: age = 2, eutherian: age = 3, therian: age = 4) of the top-100 most frequent and least frequent dimers in parental genes with shared and species-specific circRNA loci in mouse and rhesus macaque were analysed. The frequency is plotted on the x- and y-axis, point size reflects the age and point colour the enrichment (blue = decrease, red = increase). Based on the comparison between shared and species-specific dimers, the top-5 dimers defined by frequency and enrichment are highlighted and labelled in red. C: Species-specific dimer landscape for the *Akt3* gene in human, mouse and opossum. UCSC genome browser view for the parental gene, circRNAs and top-5 dimers (as defined in panel B). Start and stop positions of each dimer are connected via an arc. Dimers are grouped by composition represented by different colours, the number of collapsed dimers is indicated to the right-side of the dimer group. Only dimers that start before and stop after a circRNAs are shown as these are potentially those that can contribute to the hairpin structure. The human *Akt3* gene possesses two circRNA clusters. For better visualisation, only the upstream cluster is shown.

**Figure 5–Figure supplement 1.** Species-specific repeats contribute to the formation of shared circRNA loci.

---

### ₄₀₁ Discussion

₄₀₂ Different scenarios have been proposed for how circRNA evolution takes place (see e.g. *Patop et al.*

₄₀₃ (*2019*) for a review). Our analyses of an extensive new cross-species dataset strongly suggest that

₄₀₄ many circRNA loci that are shared across orthologous genes of different species – of which there

₄₀₅ are surprisingly few – have emerged by convergent evolution, driven by structural commonalities

₄₀₆ of their parental genes, rather than having evolved from common ancestral circRNA loci. Parental

₄₀₇ genes are composed of many exons, are located in genomic regions of low GC content, and are

₄₀₈ surrounded by an elevated number of TEs, together creating "circRNA niches" – genomic regions

₄₀₉ in which circRNAs are more likely to be generated. TEs are an indispensable feature of the niche,

₄₁₀ and in addition to their similarity in structure, orthologous parental genes thus also possess a sim-

₄₁₁ ilar, pronounced integration bias for transposons, which subsequently manifests in genomic "TE

₄₁₂ hotspots" that are shared across species. Accordingly, many TEs found within the circRNA niche

₄₁₃ possess species-specific amplification patterns and have been active only recently, or are still ac-

₄₁₄ tive even today. Due to their evolutionary youth, the genomic sequences of TEs in the circRNA

₄₁₅ niche are barely degraded, increasing the likelihood of intramolecular RNA secondary structures,

₄₁₆ which have previously been associated with circRNA biogenesis. Taken together, these findings

₄₁₇ suggest that circRNAs and TEs co-evolve in a species-specific and dynamic manner. Moreover, as

₄₁₈ most circRNAs are evolutionarily young, they are overall rather unlikely to fulfil crucial functions.

₄₁₉ This idea is in agreement with the generally low expression levels of circRNAs and with accumula-

₄₂₀ tion patterns that are frequently tissue-specific and confined to post-mitotic cells (*Guo et al., 2014*;

₄₂₁ *Westholm et al., 2014*). The model we present provides an explanation for how circRNAs can arise

₄₂₂ from shared (orthologous) exonic loci among species even if they themselves are not homologous

₄₂₃ (i.e., they do not stem from common evolutionary precursors that emerged in common ancestors).

₄₂₄ Finally, the properties we identified for the orthologous genomic niche can serve to predict circRNA

₄₂₅ parental genes with high confidence, opening the possibility to improve current circRNA prediction

₄₂₆ tools and to prioritise circRNAs for potential functional experiments.

TEs are a major component of most genomes and associated with various mechanisms that shape genome architecture and evolution. For example, TE integration into exons (changing the coding sequence) or at splice sites (potentially altering splicing patterns) may lead to the production of erroneous transcripts (*Zhang et al., 2011*). Other integration events are less sensitive towards creating such potentially hazardous "transcriptional noise". For example, TEs that integrate in safe distance to important regions of a gene - e.g. in the middle of a long intron - might not cause more than a small increase in the transcript error rate that will in most cases be tolerable for the organism. As a consequence, TEs are more likely to be tolerated in genes with long introns than in short and compact genes. Moreover, long genes are known to be GC-poor (*Zhu et al., 2009*). These characteristics overlap precisely with those that we identify for circRNAs, which are also frequently generated from genes that are poor in GC and that have long introns, complex gene structures, as well as many TEs. In other words, the propensity to produce circRNAs scales with the same features that also predispose genes to transcriptional noise. Conceivably, many circRNAs may thus represent, at their core, a side effect of the genes' transcriptional noise. In agreement with this model, a recent study in rat neurons has reported that the set of circRNAs that is upregulated after spliceosome inhibition is characterised by even longer flanking introns and an even higher number of RVCs than the average circRNA (*Wang et al., 2019*). Why is it frequently the same (orthologous) genomic loci and exons across species that independently develop the capacity for circRNA production? It is plausible that this phenomenon can be put down to tolerance for error rates. Let us consider repeat integration in close proximity to an exon boundary, which is an event that will likely alter local GC content. For example, GC-rich SINE elements that integrate in close proximity to a splice site can lead to a local increase in GC, which decreases the GC amplitude at the exon-intron boundary. Especially in GC-low genes, this can interfere with the intron-defined mechanism of splicing and cause mis-splicing (*Amit et al., 2012*). It is thus likely that TE integration close to a very strong splice site (i.e., with strong GC amplitude, as typically found in canonical exons) would have fewer repercussions on transcript error rates than integration close to alternative exons, whose GC amplitudes are less pronounced. Fully in line with such a model, we found that exons that are used in circRNAs are typically canonical exons with strong GC amplitudes. While at first sight, circRNA exons therefore appear to combine many rather specific, evolutionarily relevant properties (in particular, increased phastCons scores), we deem it probable that these are a mere consequence of a higher tolerance of canonical exon-flanking introns to TE integration.

Notably, this model may be taken even one step further by speculating whether circRNA properties for which a connection to TEs appears far-fetched, could in fact be ascribed to a transposon effect after all. Such cases are, for example, the reported predisposition of circRNAs to RNA editing (*Ivanov et al., 2015*) and different methylation patterns at both the RNA and DNA level (*Zhou et al., 2017*; *Enuka et al., 2016*; *Deniz et al., 2019*; *Aktaş et al., 2017*). How could transposons come into play? Briefly, intronic TEs can facilitate the formation of local secondary structures in pre-mRNAs. On the one hand, this would interfere with splice-site accessibility and lead to an increase in transcript error rates (*Salari et al., 2012*; *Melamud and Moult, 2009*). On the other hand, the secondary structures are associated with circRNA production. To avoid the negative impact of TEs on gene transcription, several defence mechanisms have evolved to silence them. RNA edit-

ing, for example, is thought to have evolved as a mechanism to suppress TE amplification, and A-to-I RNA editing is indeed associated with intronic Alu elements to inhibit Alu dimers (*Lev-Maor et al., 2008*; *Athanasiadis et al., 2004*). In agreement with this notion, circRNA flanking introns are enriched in A-to-I editing sites, and knockdown of the editing machinery leads to an increase in circRNA levels (*Ivanov et al., 2015*; *Rybak-Wolf et al., 2015*). Based on such findings, the conclusion has been drawn that A-to-I editing could represent a mechanism to control circRNA production (*Ivanov et al., 2015*; *Rybak-Wolf et al., 2015*). However, the alternative scenario appears equally likely, in that changes in circRNA frequencies are actually a secondary effect caused by the primary purpose of A-to-I editing, namely the inhibition of Alu amplification. This notion is in line with the findings of (*Aktaş et al., 2017*) who showed that the nuclear RNA helicase DHX9 interacts with ADAR and can bind to inverted Alu elements that are transcribed as part of the gene (*Aktaş et al., 2017*). The loss of DHX9 leads to an increase of circRNA abundance from parental genes, in agreement with the model that DHX9 resolves TE-induced mRNA secondary structures to avoid interference with post-transcriptional processes (*Aktaş et al., 2017*). Similar reasoning can be applied to other modifications at the DNA and RNA level. Notably, DNA methylation interferes with TE amplification (*Yoder et al., 1997*), and has been connected to circRNA production (*Enuka et al., 2016*).

The modification $N^6$-methyladenosine (m$^6$A) plays various roles in mRNA metabolism, including in mRNA splicing, degradation and translation (reviewed in *Zaccara et al.* (*2019*)). m$^6$A is enriched in circRNA exons and can trigger circRNA cleavage and degradation (*Zhou et al., 2017*; *Park et al., 2019*; *Di Timoteo et al., 2020*) and has therefore been viewed as a way to control circRNA levels dynamically and in a tissue-specific manner. However, increased levels of m$^6$A, which is deposited already on the nascent RNA, are part of a much broader mechanism for mRNA destabilisation (reviewed in *Lee et al.* (*2020*)). Hence, it is possible that increased levels of m$^6$A on circRNAs rather reflect the general targeting of faulty transcripts for rapid degradation.

These considerations - together with our evolutionary data - lead us to the interpretation that many circRNAs likely represent transcriptional noise caused by TEs integrated into parental genes. However, it is also clear that molecular functions have been identified for several circRNAs (e.g. *Hansen et al.* (*2013*); *Conn et al.* (*2015*); *Du et al.* (*2016*)), although the absolute number of validated examples remains modest when compared to the high number of different circRNAs that have been detected across cell types, developmental stages and species. One would imagine that in order to evolve a function from noise, circRNAs need to reach critical, stable expression levels that bestow a positive effect on the organisms' fitness – a process that might take considerable time. Yet, circRNAs are not produced from scratch, but evolve from already existing functional genes, a process commonly known as exaptation (*Brosius and Gould, 1992*). A well-known example for this mechanism is provided by several miRNAs that evolved independently from each other in the same genomic position relative to the *Hox8* gene (*Campo-Paysaa et al., 2011*). For the circRNAs, the evolution of a function may be accelerated due to the presence of a clear exon structure and of regulatory elements from which the circRNAs can benefit. The production of structurally similar circRNAs from circRNA hotspots may accelerate this process, by providing different (back)splice sites and regulatory elements as evolutionary raw material, while keeping the internal exon sequence fairly similar. Once a circRNA emerges that is endowed with beneficial character-

509 istics and equipped with an initial set of regulatory elements, the typically rather low expression

510 level may increase. Robust expression and the acquisition of additional regulatory motifs (includ-

511 ing those for RNA-binding proteins) may ultimately render the circRNA independent of its original

512 regulation through reverse-complementary sequences (as described in *Ashwal-Fluss et al.* (*2014*);

513 *Conn et al.* (*2015*); *Okholm et al.* (*2020*)). Thus, given that circRNAs are produced from hundreds of

514 loci, in many cell types and across different developmental stages, beneficial circRNAs with useful

515 functions – such as those that have been reported – may emerge and be fixed in a species during

516 evolution.

517 In summary, our data suggests that many circRNA molecules do not carry specific molecular

518 functions. However, one may still speculate whether it is actually the process of RNA circularization

519 in itself, rather than the circRNA molecule, that is beneficial. For example, circularization may rep-

520 resent a mechanism to keep genes under control that have transformed into "transposon sinks",

521 by directing mRNA output from such transposon-rich loci towards non-productive, circular tran-

522 scripts. One could also argue that some level of splicing noise may be beneficial to engender gene

523 expression plasticity at circRNA loci. Finally, circRNAs have emerged as reliable disease biomark-

524 ers (*Memczak et al., 2015*; *Bahn et al., 2015*), and their utility for such predictive purposes is not

525 affected by our conclusions – on the contrary. While an altered circRNA profile will likely not have

526 a causal involvement in a disease, it could hint at misregulated transcription or splicing of the

527 parental gene, at a novel TE integration event, or at problems with the RNA editing or methyla-

528 tion machinery. The careful analysis of the circRNA landscape may thus teach us about factors

529 contributing to diseases in a causal fashion even if many or perhaps most circRNAs may not be

530 functional but rather represent transcriptional noise.

531 ## Material and Methods

532 ## Data deposition, programmes and working environment

**Table 2.** Overview of external programmes.

| Programme | Version |
|---|---|
| Blast | 2.2.29+ |
| BEDTools | 2.17.0 |
| Bowtie2 | 2.1.0 |
| Clustal Omega | 1.2.4 |
| Cufflinks | 2.1.1 |
| FastQC | 0.10.1 |
| Mcl | 14.137 |
| R | 3.0 and 3.1 |
| Ruby | 2.0 and 2.1 |
| SAMTools | 0.1.19 |
| TopHat2 | 2.0.11 |
| ViennaRNA | 2.1.8 |

533 The raw data and processed data files discussed in this publication have been deposited in NCBI's
534 Gene Expression Omnibus (*Edgar et al., 2002*) and are accessible through the GEO Series accession
535 number GSE162152. All scripts used to produce the main figures and tables of this publication
536 have been deposited in the Git Repository circRNA_paperScripts. This Git repository also holds
537 information on how to run the scripts, and links to the underlying data files for the main figures.
538 The custom pipeline developed for the circRNA identification can be found in the Git Repository
539 ncSplice_circRNAdetection.

540 ## Library preparation and sequencing

541 We used 5 µg of RNA per sample as starting material for library preparation, which were treated
542 with 20 U RNase R (Epicentre/Illumina, Cat. No. RNR07250) for 1 h at 37°C to degrade linear RNAs,
543 followed by RNA purification with the RNA Clean & Concentrator-5 kit (Zymo Research) according
544 to the manufacturer's protocol. Paired-end sequencing libraries were prepared from the purified
545 RNA with the Illumina TruSeq Stranded Total RNA kit with Ribo-Zero Gold according to the protocol
546 with the following modifications to select larger fragments: 1.) Instead of the recommended 8 min
547 at 68°C for fragmentation, we incubated samples for only 4 min at 68°C to increase the fragment
548 size; 2.) In the final PCR clean-up after enrichment of the DNA fragments, we changed the 1:1 ratio
549 of DNA to AMPure XP Beads to a 0.7:1 ratio to select for binding of larger fragments. Libraries were
550 analysed on the fragment analyzer for their quality and sequenced with the Illumina HiSeq 2500
551 platform (multiplexed, 100 cycles, paired-end, read length 100 nt).

### 552 Identification and quantification of circRNAs

### 553 Mapping of RNA-seq data

554 The ensembl annotations for opossum (monDom5), mouse (mm10), rat (rn5), rhesus macaque
555 (rheMac2) and human (hg38) were downloaded from Ensembl to build transcriptome indexes for
556 mapping with TopHat2. TopHat2 was run with default settings and the *–mate-inner-dist* and *–mate-*
557 *std-dev* options set to 50 and 200 respectively. The mate-inner-distance parameter was estimated
558 based on the fragment analyzer report.

**Table 3.** Ensembl genome versions and annotation files for each species.

| Species | Genome | Annotation |
|---|---|---|
| Opossum | monDom5 | ensembl release 75, feb 2014 |
| Mouse | mm10 | ensembl release 75, feb 2014 |
| Rat | rn5 | ensembl release 75, feb 2014 |
| Rhesus macaque | rheMac2 | ensembl release 77, oct 2014 |
| Human | hg38 | ensembl release 77, oct 2014 |

### 559 Analysis of unmapped reads

560 We developed a custom pipeline to detect circRNAs (**Figure1-Figure supplement 1B**), which per-
561 forms the following steps: Unmapped reads with a phred quality value of at least 25 are used to
562 generate 20 bp anchor pairs from the terminal 3' and 5'-ends of the read. Anchors are remapped
563 with bowtie2 on the reference genome. Mapped anchor pairs are filtered for 1) being on the same
564 chromosome, 2) being on the same strand and 3) for having a genomic mapping distance to each
565 other of a maximum of 100 kb. Next, anchors are extended upstream and downstream of their
566 mapping locus. They are kept if pairs are extendable to the full read length. During this procedure
567 a maximum of two mismatches is allowed. For paired-end sequencing reads, the mate read not
568 mapping to the backsplice junction can often be mapped to the reference genome without any
569 problem. However, it will be classified as "unmapped read" (because its mate read mapping to
570 the backsplice junction was not identified by the standard procedure). Next, all unpaired reads
571 are thus selected from the accepted_hits.bam file generated by TopHat2 (singletons) and assessed
572 for whether the mate read (second read of the paired-end sequencing read) of the anchor pair
573 mapped between the backsplice coordinates. All anchor pairs for which 1) the mate did not map
574 between the genomic backsplice coordinates, 2) the mate mapped to another backsplice junction
575 or 3) the extension procedure could not reveal a clear breakpoint are removed. Based on the re-
576 maining candidates, a backsplice index is built with bowtie2 and all reads are remapped on this
577 index to increase the read coverage by detecting reads that cover the BSJ with less than 20 bp,
578 but at least 8 bp. Candidate reads that were used to build the backsplice index and now mapped
579 to another backsplice junction are removed. Upon this procedure, the pipeline provides a first
580 list of backsplice junctions. The set of scripts, which performs the identification of putative BSJs,
581 as well as a short description of how to run the pipeline are deposited in the Git Repository nc-

**582**   Splice_circRNAdetection.

**583**   Trimming of overlapping reads

**584**   Due to small DNA repeats, some reads are extendable to more than the original read length. There-

**585**   fore, overlapping reads were trimmed based on a set of canonical and non-canonical splice sites.

**586**   For the donor site GT, GC, AT, CT were used and for the acceptor splice site AG and AC. The trim-

**587**   ming is part of our custom pipeline described above, and the step will be performed automatically

**588**   if the scripts are run.

**589**   Calculation of CPM value

**590**   CPM (counts per million) values for BSJs were calculated for each tissue as follows:

$$counts = \frac{counts\_rep1 + counts\_rep2 + counts\_rep3}{3}$$

$$totalMappedReads = \frac{mappedReads\_rep1 + mappedReads\_rep2 + mappedReads\_rep3}{3}$$

$$CPM = \frac{counts \cdot 10^6}{totalMappedReads}$$

**591**   Filtering of candidates based on CPM enrichment

**592**   To distinguish putative BSJs from the technical and biological noise background, the enrichment of

**593**   the previously (in untreated samples) defined junctions in RNase R treated samples was calculated.

**594**   The enrichment was defined as CPM increase in RNase R treated versus untreated samples:

$$enrichment = \frac{CPM\_RNaseR}{CPM\_untreated}$$

**595**   Candidates with a log2-enrichment of smaller 1.5, as well as less than 0.05 CPM, were removed.

**596**   Manual filtering steps

**597**   We observed several genomic loci in rhesus macaque and human that were highly enriched in

**598**   reads for putative BSJs (no such problem was detected for opossum, mouse and rat). Manual

**599**   inspection in the UCSC genome browser indicated that these loci are highly repetitive. The detected

**600**   BSJs from these regions do probably not reflect BSJs, but instead issues in the mapping procedure.

**601**   These candidates were thus removed manually; the regions are:

**602**   All following analyses were conducted with the circRNA candidates that remained after this step.

**603**   Reconstruction of circRNA isoforms

**604**   To reconstruct the exon structure of circRNA transcripts in each tissue, we made use of the junction

**605**   enrichment in RNase R treated samples. To normalise junction reads across libraries, the size

**606**   factors based on the geometric mean of common junctions in untreated and treated samples were

**607**   calculated as

**Table 4.** Removed regions during mapping.

| species | tissue | chromosome | start | stop | strand |
|---|---|---|---|---|---|
| rhesus macaque | testis | 7 | 164261343 | 164283671 | + |
| rhesus macaque | testis | 7 | 22010814 | 22092409 | - |
| rhesus macaque | testis | 19 | 52240850 | 52288425 | - |
| rhesus macaque | testis | 19 | 59790996 | 59834798 | + |
| rhesus macaque | testis | 19 | 59790996 | 59847609 | + |
| human | testis | 2 | 178535731 | 178600667 | + |
| human | testis | 7 | 66429678 | 66490107 | - |
| human | testis | 9 | 97185441 | 97211487 | - |
| human | testis | 12 | 97492460 | 97561047 | + |
| human | testis | 14 | 100913431 | 100949596 | + |
| human | testis | 18 | 21765771 | 21849388 | + |

$$geometric\_mean = \left( \prod x \right)^{\frac{1}{length(x)}}$$

$$size\_factor = median \left( \frac{x}{geometric\_mean} \right)$$

608 with *x* being a vector containing the number of reads per junction. We then compared read cover-

609 age for junctions outside and inside the BSJ for each gene and used the log2-change of junctions

610 outside the backsplice junction to construct the expected background distribution of change in

611 junction coverage upon RNase R treatment. The observed coverage change of junctions inside the

612 backsplice was then compared to the expected change in the background distribution and junc-

613 tions with a log2-change outside the 90% confidence interval were assigned as circRNA junctions;

614 a loose cut-off was chosen, because involved junctions can show a decrease in coverage if their lin-

615 ear isoform was present at high levels before (degradation levels of linear isoforms do not correlate

616 with the enrichment levels of circRNAs). Next, we reconstructed a splicing graph for each circRNA

617 candidate, in which network nodes are exons connected by splice junctions (edges) (*Heber et al.,*

618 *2002*). Connections between nodes are weighted by the coverage in the RNase R treated samples.

619 The resulting network graph is directed (because of the known circRNA start and stop coordinates),

620 acyclic (because splicing always proceeds in one direction), weighted and relatively small. We used

621 a simple breadth-first-search algorithm to traverse the graph and to define the strength for each

622 possible isoform by its mean coverage. Only the strongest isoform was considered for all subse-

623 quent analyses.

624 Reconstruction and expression quantification of linear mRNAs

625 We reconstructed linear isoforms based on the pipeline provided by *Trapnell et al.* (*2012*) (Cufflinks

626 + Cuffcompare + Cuffnorm). Expression levels were quantified based on fragments per million

627 mapped reads (FPKM). Cufflinks was run per tissue and annotation files were merged across tissues

628 with Cuffcompare. Expression was quantified with Cuffnorm based on the merged annotation file.

629 All programs were run with default settings. FPKM values were normalised across species and
630 tissues using a median scaling approach as described in *Brawand et al.* (*2011*).

631 **Identification of shared circRNA loci between species**

632 Definition and identification of shared circRNA loci

633 Shared circRNA loci were defined on three different levels depending on whether the "parental
634 gene", the "circRNA locus" in the gene or the "start/stop exons" overlapped between species (see
635 **Figure 2A** and **Figure 2-Figure supplement 1**). Overall considerations of this kind have recently
636 also been outlined in *Patop et al.* (*2019*).

637 Level 1 - Parental genes: One-to-one (1:1) therian orthologous genes were defined between
638 opossum, mouse, rat, rhesus macaque and human using the Ensembl orthology annotation (con-
639 fidence intervals 0 and 1, restricted to clear one-to-one orthologs). The same procedure was per-
640 formed to retrieve the 1:1 orthologous genes for the eutherians (mouse, rat, rhesus macaque,
641 human), for rodents (mouse, rat) and primates (rhesus macaque, human). Shared circRNA loci be-
642 tween species were assessed by counting the number of 1:1 orthologous parental genes between
643 the five species. The analysis was restricted to protein-coding genes.

644 Level 2 - circRNA locus: To identify shared circRNA loci, all circRNA exon coordinates from a given
645 gene were collapsed into a single transcript using the *bedtools merge* option from the BEDTools
646 toolset with default options. Next, we used liftOver to compare exons from the collapsed transcript
647 between species. The minimal ratio of bases that need to overlap for each exon was set to 0.5 (-
648 *minMatch=0.5*). Collapsed transcripts were defined as overlapping between different species if they
649 shared at least one exon, independent of the exon length.

650 Level 3 - start/stop exon: To identify circRNAs sharing the same first and last exon between
651 species, we lifted exons coordinates between species (same settings as described above, *liftOver,*
652 *-minMatch=0.5*). The circRNA was then defined as "shared", if both exons were annotated as start
653 and stop exons in the respective circRNAs of the given species. Note, that this definition only
654 requires an overlap for start and stop exons, internal circRNA exons may differ.

655 Given that only circRNAs that comprise corresponding (1:1 orthologous exons) in different
656 species might at least potentially and reasonably considered to be homologous (i.e., might have
657 originated from evolutionary precursors in common ancestors) and the Level 3 definition might
658 require strong evolutionary conservation of splice sites (i.e., with this stringent definition many
659 shared loci may be missed), we decided to use the level 2 definition (circRNA locus) for the analy-
660 ses presented in the main text, while we still provide the results for the Level 1 and 3 definitions
661 in the supplement (**Figure 2-Figure supplement 1**). Importantly, defining shared circRNA loci at
662 this level allows us to also compare circRNA hostspots which have been defined using a similar
663 classification strategy.

664 Clustering of circRNA loci between species

665 Based on the species set in which shared circRNA loci were found, we categorised circRNAs in the
666 following groups: Species-specific, rodent, primate, eutherian and therian circRNAs. To be part of
667 the rodent or primate group, the circRNA has to be expressed in both species of the lineage. To

⁶⁶⁸ be part of the eutherian group, the circRNA has to be expressed in three species out of the four

⁶⁶⁹ species mouse, rat, rhesus macaque and human. To be part of the therian group, the circRNA

⁶⁷⁰ needs to be expressed in opossum and in three out of the four other species. Species-specific

⁶⁷¹ circRNAs are either present in one species or do not match any of the other four categories. To

⁶⁷² define the different groups, we used the cluster algorithm MCL (*Enright et al., 2002*; *Dongen, 2000*).

⁶⁷³ MCL is frequently used to reconstruct orthology clusters based on blast results. It requires input in

⁶⁷⁴ *abc* format (file: *species.abc*), in which *a* corresponds to event a, *b* to event b and a numeric value *c*

⁶⁷⁵ that provides information on the connection strength between event a and b (e.g. blast p-value). If

⁶⁷⁶ no p-values are available as in this analysis, the connection strength can be set to 1. MCL was run

⁶⁷⁷ with a cluster granularity of 2 (*option -I*).

⁶⁷⁸

⁶⁷⁹ *$ mcxload -abc species.abc –stream-mirror -o species.mci -write-tab species.tab*

⁶⁸⁰ *$ mcl species.mci -I 2*

⁶⁸¹ *$ mcxdump -icl out.species.mci.I20 -tabr species.tab -o dump.species.mci.I20*

⁶⁸² PhastCons scores

⁶⁸³ Codings exons were selected based on the attribute "transcript_biotype = protein_coding" in the gtf

⁶⁸⁴ annotation file of the respective species and labelled as circRNA exons if they were in our circRNA

⁶⁸⁵ annotation. Exons were further classified into UTR-exons and non-UTR exons using the ensembl

⁶⁸⁶ field "feature = exon" or "feature = UTR". Since conservation scores are generally lower for UTR-

⁶⁸⁷ exons (*Pollard et al., 2010*), any exon labelled as UTR-exon was removed from further analyses to

⁶⁸⁸ avoid bias when comparing circRNA and non-circRNA exons. Genomic coordinates of the remain-

⁶⁸⁹ ing exons were collapsed using the *merge* command from the BEDtools toolset (*bedtools merge*

⁶⁹⁰ *input_file -nms -scores collapse*) to obtain a list of unique genomic loci. PhastCons scores for all

⁶⁹¹ exon types were calculated using the conservation scores provided by the UCSC genome browser

⁶⁹² (mouse: phastCons scores based on alignment for 60 placental genomes; rat: phastCons scores

⁶⁹³ based on alignment for 13 vertebrate genomes; human: phastCons scores based on alignment

⁶⁹⁴ for 99 vertebrate genomes). For each gene type (parental or non-parental), the median phastCons

⁶⁹⁵ score was calculated for each exon type within the gene (if non-parental: median of all exons; if

⁶⁹⁶ parental: median of exons contained in the circRNA and median of exons outside of the circRNA).

⁶⁹⁷ Tissue specificity of exon types

⁶⁹⁸ Using the DEXseq package (from HTSeq 0.6.1), reads mapping on coding exons of the parental

⁶⁹⁹ genes were counted. The exon-bins defined by DEXseq (filtered for bins >=10 nt) were then mapped

⁷⁰⁰ and translated onto the different exon types: UTR-exons of parental genes, exons of parental genes

⁷⁰¹ that are not in a circRNA, circRNA exons. For each exon type, an FPKM value based on the exon

⁷⁰² length and sequencing depth of the library was calculated.

$$FPKM = \frac{counts\_for\_exon\_type \cdot 10^9}{exon\_type\_length / sequencing\_depth}$$

703 Exons were labelled as expressed in a tissue, if the calculated FPKM was at least 1. The maximum

704 number of tissues in which each exon occurred was plotted separately for UTR-exons, exons out-

705 side the circRNA and contained in it.

706 ## GC amplitude

707 The ensembl annotation for each species was used to retrieve the different known transcripts in

708 each coding gene. For each splice site, the GC amplitude was calculated using the last 250 intronic

709 bp and the first 50 exonic bp (several values for the last $n$ intronic bp and the first $m$ exonic bp

710 were tested beforehand, the 250:50 ratio was chosen, because it gave the strongest signal). Splice

711 sites were distinguished by their relative position to the circRNA (flanking, inside or outside). A one-

712 tailed and paired Mann-Whitney U test was used to assess the difference in GC amplitude between

713 circRNA-related splice sites and others.

714 **Parental gene analysis**

715 ## GC content of exons and intron

716 The ensembl annotation for each species was used to retrieve the different known transcripts in

717 each coding gene. Transcripts were collapsed per-gene to define the exonic and intronic parts.

718 Introns and exons were distinguished by their relative position to the circRNA (flanking, inside or

719 outside). The GC content was calculated based on the genomic DNA sequence. On a per-gene level,

720 the median GC content for each exon and intron type was used for further analyses. Differences

721 between the GC content were assessed with a one-tailed Mann-Whitney U test.

722 ## Gene self-complementarity

723 The genomic sequence of each coding gene (first to last exon) was aligned against itself in sense

724 and antisense orientation using megaBLAST with the following call:

725

726 *$ blastn -query seq.fa -subject seq.fa -task dc-megablast -word_size 12 -outfmt "6 qseqid qstart qend*

727 *sseqid sstart send sstrand length pident nident mismatch bitscore evalue" > blast.out*

728

729 The resulting alignments were filtered for being purely intronic (no overlap with any exon). The

730 fraction of self-complementarity was calculated as the summed length of all alignments in a gene

731 divided by its length (first to last exon).

732 ## Generalised linear models

733 All linear models were developed in the R environment. The presence of multicollinearity between

734 predictors was assessed using the *vif()* function from the R package *car* (version 3.0-3) to calculate

735 the variance inflation factor (VIF). Predictors were scaled to be able to compare them with each

736 other using the *scale()* function as provided in the R environment.

737 For parental genes, the dataset was split into training (80%) and validation set (20%). To find the

738 strongest predictors, we used the R package *bestglm* (version 0.37). Each model was fitted on the

739 complete dataset using the command *bestglm()* with the information criteria set to "CV" (CV = cross

740 validation) and the number of repetitions *t = 1000*. The model family was set to "binomial" as we

741 were merely interested in predicting the presence (1) or absence (0) of a parental gene. Significant

742 predictors were then used to report log-odds ratios and significance levels for the validation set

743 using the default *glm()* function of the R environment. Log-odds ratios, standard errors and confi-

744 dence intervals were standardised using the *beta()* function from the *reghelper* R package (version

745 1.0.0) and are reported together with their p-values in **Supplementary Table 5**.

746     For the correlation of hotspot presence across the number of species, a generalised linear

747 model was applied using the categorical predictors "lineage" (= circRNA loci shared within rodents

748 or primates), "eutherian" (= circRNA loci shared within rodents and primates) and "therian" (= cir-

749 cRNA loci shared within opossum, rodents and primates). Log-odds ratios, standard errors and

750 confidence intervals were standardised using the *beta()* function from the *reghelper* R package (ver-

751 sion 1.0.0) and are reported together with their p-values in **Supplementary Table 6**.

752 Comparison to human and mouse circRNA heart dataset

753 The circRNA annotations for human and mouse heart as provided by *Werfel et al.* (*2016*) were,

754 based on the parental gene ID, merged with our circRNA annotations. Prediction values for parental

755 genes were calculated using the same general linear regression models as described above (Sec-

756 tion *Generalised linear models* in Material and Methods section) with genomic length, number of

757 exons, GC content, expression levels, reverse complements (RVCs) and phastCons scores as pre-

758 dictors. Prediction values were received from the model and compared between parental genes

759 predicted by our and the Werfel dataset as well as between the predictors in non-parental and

760 parental genes of the Werfel dataset (**Figure 3-Figure supplement 3**).

761 Integration of external studies

762 (1) Replication time

763 Values for the replication time were used as provided in *Koren et al.* (*2012*). Coordinates of the dif-

764 ferent replication domains were intersected with the coordinates of coding genes using BEDtools

765 (*bedtools merge -f 1*). The mean replication time of each gene was used for subsequent analyses.

766

767 (2) Gene expression steady-state levels

768 Gene expression steady-state levels and decay rates were used as provided in Table S1 of *Pai et al.*

769 (*2012*).

770

771 (3) GHIS

772 Genome-wide haploinsufficiency scores for each gene were used as provided in Supplementary

773 Table S2 of *Steinberg et al.* (*2015*).

774 **Repeat analyses**

775 Generation of length- and GC-matched background dataset

776 Flanking introns were grouped into a matrix of *i* columns and *j* rows representing different genomic

777 lengths and GC content; *i* and *j* were calculated in the following way:

$$i = seq(from = quantile(GCcontent, 0.05), to = quantile(GCcontent, 0.95), by = 0.01)$$

$$j = seq(from = quantile(length, 0.05), to = quantile(length, 0.95), by = 1000)$$

778 Flanking introns were sorted into the matrix based on their GC content and length. A second matrix

779 with the same properties was created containing all introns of coding genes. From the latter, a

780 submatrix was sampled with the same length and GC distribution as the matrix for flanking introns.

781 The length distribution and GC distribution of the sampled introns reflect the distributions for the

782 flanking introns as assessed by a Fisher's t Test that was non-significant.

783 Repeat definition

784 The RepeatMasker annotation for full and nested repeats were downloaded for all genomes using

785 the UCSC Table browser (tracks "RepeatMasker" and "Interrupted Rpts") and the two files merged.

786 Nested repeats were included, because it was shown that small repetitive regions are sufficient to

787 trigger base pairing necessary for backsplicing (*Liang and Wilusz, 2014*; *Kramer et al., 2015*). The

788 complete list was then intersected (*bedtools merge -f1*) with the above defined list of background

789 and flanking introns for further analyses.

790 Identification of repeat dimers

791 The complementary regions (RVCs) that were defined with megaBLAST as described above, were

792 intersected with the coordinates of individual repeats from the RepeatMasker annotation. To be

793 counted, a repeat had to overlap with at least 50% of its length with the region of complementarity

794 (*bedtools merge -f 0.5*). As RVCs can contain several repeats, the "strongest" dimer was selected

795 based on the number of overlapping base pairs (= longest overlapping dimer). The "dimer list" ob-

796 tained from this analysis for each species was further ranked according to the absolute frequency

797 of each dimer. The proportion of the top-5 dimer frequency to all detected dimers, was calculated

798 based on this list ($n_{top-5}$ / $n_{all\_dimers}$).

799 Binding scores of repeat dimers

800 Binding scores for each TE class (based on the TE reference sequence) were defined by taking into

801 account the (1) phylogenetic distance to other repeat families in the same species and (2) its bind-

802 ing affinity (deltaG) to those repeats. We decided to not include the absolute TE frequency into

803 the binding score, because it is a function of the TE's age, its amplification and degradation rates.

804 Simulating the interplay between these three components is not in scope of this study, and the in-

805 tegration of frequency into binding score creates more noise as tested via PCA analyses (variance

806 explained drops by 10%).

807

808 (1) Phylogenetic distance

809 TE reference sequences were obtained from Repbase (*Bao et al., 2015*) and translated into fasta-

810 format for alignment (*reference_sequences.fa*). Alignments were then generated with Clustal Omega

811 (v1.2.4) (*Sievers et al., 2011*) using the following settings:

812

813 *$ clustalo -i reference_sequences.fa –distmat-out = repeats.mat –guidetree-out = repeats.dnd –full*

814

815 The resulting distance matrix for the alignment was used for the calculation of the binding score.

816 Visualisation of the distance matrix (**Figure 4C**, **Figure 4-Figure supplement 1**) was performed us-

817 ing the standard R functions *dist(method="euclidian")* and *hclust(method="ward.D2")*. Since several

818 TE classes evolved independently from each other, the plot was manually modified to remove con-

819 nections or to add additional information on the TE's origin from literature.

820

821 (2) Binding affinity

822 To estimate the binding affinity of individual TE dimers, the free energy of the secondary structure

823 of the respective TE dimers was calculated with the RNAcofold function from the ViennaRNA Pack-

824 age:

825

826 *$ RNAcofold -a -d2 < dimerSequence.fa*

827

828 with *dimerSequence.fa* containing the two reference sequences of the TEs from which the dimer is

829 composed. The resulting deltaG values were used to calculate the binding score.

830

831 (3) Final binding score

832 To generate the final binding score, values from the distance matrix and the binding affinity were

833 standardised (separately from each other) to values between 0 and 1:

$$f(x) = \frac{x - min(v)}{max(v) - min(v)}$$

834 with *x* being the binding affinity/dimer frequency and *minv* and *maxv* the minimal and maximal

835 observed value in the distribution. The standardised values for the binding affinity and dimer fre-

836 quency were then summed up (= binding score) and classified by PCA using the R environment:

837

838 *$ pca <- prcomp(score, center=TRUE, scale.=FALSE)*

839

840 PC1 and PC2 were used for subsequent plotting with the absolute frequency of dimers represented

841 by the size of the data points.

842 Calculation of dimer degradation

843 RepeatMasker annotations were downloaded from the UCSC Table browser for all genomes. The

844 milliDiv values for each TE in a TE dimer were retrieved from this annotation for full and nested

845 repeats. A representative milliDiv was formed using the mean of the two values. Dimers were

846 then classified as species-specific or present in all species based on whether the circRNA parental

847 gene produced species-specific or shared circRNA loci. Significance levels for milliDiv differences

between the dimer classes were assessed with a simple Mann-Whitney U test (alternative set to "less").

## Supplementary Data

### Supplementary Tables and Figures

Supplementary Tables and Figures are available as an attachmente to this document.

### Supplementary Files

**Supplementary File 1:** CircRNA annotation file for opossum. A gtf-file with all circRNA transcripts including the transcript and exon coordinates.

**Supplementary File 2:** CircRNA annotation file for mouse. A gtf-file with all circRNA transcripts including the transcript and exon coordinates.

**Supplementary File 3:** CircRNA annotation file for rat. A gtf-file with all circRNA transcripts including the transcript and exon coordinates.

**Supplementary File 4:** CircRNA annotation file for rhesus macaque. A gtf-file with all circRNA transcripts including the transcript and exon coordinates.

**Supplementary File 5:** CircRNA annotation file for human. A gtf-file with all circRNA transcripts including the transcript and exon coordinates.

All gtf-files have been uploaded to the UCSC genome browser and can be viewed here:

**Opossum:** http://genome.ucsc.edu/s/Frenzchen/monDom5%20circRNA%20annotation

**Mouse;** http://genome.ucsc.edu/s/Frenzchen/mm10%20circRNA%20annotation

**Rat:** http://genome.ucsc.edu/s/Frenzchen/rn5%20circRNA%20annotation

**Rhesus macaque:** http://genome.ucsc.edu/s/Frenzchen/rheMac2%20circRNA%20annotation

**Human:** http://genome.ucsc.edu/s/Frenzchen/hg38%20circRNA%20annotation

## Author contributions

Contributions to this publication are distributed as follows: Study design: F.G., D.G., H.K. and P.J.; Experimental work: F.G. and P.J.; Bioinformatics data analyses: F.G.; Paper manuscript and discussion: F.G., D.G. and H.K.

### Competing interests

No competing interests.

## References

**Aktaş T**, Avşar Ilık Maticzka D, Bhardwaj V, Pessoa Rodrigues C, Mittler G, Manke T, Backofen R, Akhtar A. DHX9 suppresses RNA processing defects originating from the Alu invasion of the human genome. Nature. 2017 apr; 544(7648):115–119. http://dx.doi.org/10.1038/nature21715, doi: 10.1038/nature21715.

**Alhasan AA**, Izuogu OG, Al-Balool HH, Steyn JS, Evans A, Colzani M, Ghevaert C, Mountford JC, Marenah L, Elliott DJ, Santibanez-Koref M, Jackson MS. Circular RNA enrichment in platelets is a signature of transcriptome degradation. Blood. 2015 dec; http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=26660425&retmode=ref&cmd=prlinks, doi: 10.1182/blood-2015-06-649434.

**Amit M**, Donyo M, Hollander D, Goren A, Kim E, Gelfman S, Lev-Maor G, Burstein D, Schwartz S, Postolsky B, Pupko T, Ast G. Differential GC content between exons and introns establishes distinct strategies of splice-site recognition. Cell reports. 2012 may; 1(5):543–556. http://dx.doi.org/10.1016/j.celrep.2012.03.013, doi: 10.1016/j.celrep.2012.03.013.

**Ashwal-Fluss R**, Meyer M, Pamudurti NR, Ivanov A. circRNA Biogenesis Competes with Pre-mRNA Splicing. Molecular cell. 2014; .

**Athanasiadis A**, Rich A, Maas S. Widespread A-to-I RNA editing of Alu-containing mRNAs in the human transcriptome. PLoS Biology. 2004 dec; 2(12):e391. http://dx.doi.org/10.1371/journal.pbio.0020391, doi: 10.1371/journal.pbio.0020391.

**Bachmayr-Heyda A**, Reiner AT, Auer K, Sukhbaatar N, Aust S, Bachleitner-Hofmann T, Mesteri I, Grunt TW, Zeillinger R, Pils D. Correlation of circular RNA abundance with proliferation–exemplified with colorectal and ovarian cancer, idiopathic lung fibrosis, and normal human tissues. Scientific reports. 2015; .

**Bahn JH**, Zhang Q, Li F, Chan TM, Lin X, Kim Y, Wong DTW, Xiao X. The landscape of microRNA, Piwi-interacting RNA, and circular RNA in human saliva. Clinical Chemistry. 2015 jan; 61(1):221–230. http://dx.doi.org/10.1373/clinchem.2014.230433, doi: 10.1373/clinchem.2014.230433.

**Bao W**, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in eukaryotic genomes. Mobile DNA. 2015 jun; 6:11. http://dx.doi.org/10.1186/s13100-015-0041-9, doi: 10.1186/s13100-015-0041-9.

**Batzer MA**, Deininger PL. Alu repeats and human genomic diversity. Nature Reviews Genetics. 2002 may; 3(5):370–379. http://dx.doi.org/10.1038/nrg798, doi: 10.1038/nrg798.

**Brawand D**, Soumillon M, Necsulea A, Julien P, Csárdi G, Harrigan P, Weier M, Liechti A, Aximu-Petri A, Kircher M, Albert FW, Zeller U, Khaitovich P, Grützner F, Bergmann S, Nielsen R, Pääbo S, Kaessmann H. The evolution of gene expression levels in mammalian organs. Nature. 2011 oct; 478(7369):343–348. http://dx.doi.org/10.1038/nature10532, doi: 10.1038/nature10532.

**Brosius J**, Gould SJ. On "genomenclature": a comprehensive (and respectful) taxonomy for pseudogenes and other "junk DNA". Proceedings of the National Academy of Sciences of the United States of America. 1992 nov; 89(22):10706–10710. http://dx.doi.org/10.1073/pnas.89.22.10706, doi: 10.1073/pnas.89.22.10706.

**Campo-Paysaa F**, Sémon M, Cameron RA, Peterson KJ, Schubert M. microRNA complements in deuterostomes: origin and evolution of microRNAs. Evolution & Development. 2011 feb; 13(1):15–27. http://dx.doi.org/10.1111/j.1525-{142X}.2010.00452.x, doi: 10.1111/j.1525-{142X}.2010.00452.x.

**Conn SJ**, Pillman KA, Toubia J, Conn VM, Salmanidis M, Phillips CA, Roslan S, Schreiber AW, Gregory PA, Goodall GJ. The RNA binding protein quaking regulates formation of circRNAs. Cell. 2015 mar; 160(6):1125–1134. http://dx.doi.org/10.1016/j.cell.2015.02.014, doi: 10.1016/j.cell.2015.02.014.

927 **Cortés-López M**, Gruner MR, Cooper DA, Gruner HN, Voda AI, van der Linden AM, Miura P. Global accumulation
928 of circRNAs during aging in Caenorhabditis elegans. BMC Genomics. 2018 jan; 19(1):8. http://dx.doi.org/10.
929 1186/s12864-017-4386-y, doi: 10.1186/s12864-017-4386-y.

930 **Deniz** , Frost JM, Branco MR. Regulation of transposable elements by DNA modifications. Nature Reviews
931 Genetics. 2019; 20(7):417–431. http://www.nature.com/articles/s41576-019-0106-6, doi: 10.1038/s41576-019-
932 0106-6.

933 **Di Timoteo G**, Dattilo D, Centrón-Broco A, Colantoni A, Guarnacci M, Rossi F, Incarnato D, Oliviero S, Fatica A,
934 Morlando M, Bozzoni I. Modulation of circRNA Metabolism by m(6)A Modification. Cell reports. 2020 May;
935 31(6):107641. http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=32402287&retmode=
936 ref&cmd=prlinks, doi: 10.1016/j.celrep.2020.107641.

937 **Dongen S**. Performance criteria for graph clustering and Markov cluster experiments. . 2000 May; http://dl.
938 acm.org/citation.cfm?id=868979.

939 **Du WW**, Yang W, Liu E, Yang Z, Dhaliwal P, Yang BB. Foxo3 circular RNA retards cell cycle progression via
940 forming ternary complexes with p21 and CDK2. Nucleic Acids Research. 2016 apr; 44(6):2846–2858. http:
941 //dx.doi.org/10.1093/nar/gkw027, doi: 10.1093/nar/gkw027.

942 **Dubin RA**, Kazmi MA, Ostrer H. Inverted repeats are necessary for circularization of the mouse testis Sry
943 transcript. Gene. 1995 dec; 167(1-2):245–248. https://www.ncbi.nlm.nih.gov/pubmed/8566785.

944 **Edgar R**, Domrachev M, Lash AE. Gene Expression Omnibus: NCBI gene expression and hybridization array
945 data repository. Nucleic Acids Research. 2002 jan; 30(1):207–210. http://dx.doi.org/10.1093/nar/30.1.207,
946 doi: 10.1093/nar/30.1.207.

947 **Enright AJ**, Van Dongen S, Ouzounis CA. An efficient algorithm for large-scale detection of protein fam-
948 ilies. Nucleic Acids Research. 2002 apr; 30(7):1575–1584. http://dx.doi.org/10.1093/nar/30.7.1575, doi:
949 10.1093/nar/30.7.1575.

950 **Enuka Y**, Lauriola M, Feldman ME, Sas-Chen A, Ulitsky I, Yarden Y. Circular RNAs are long-lived and display only
951 minimal early alterations in response to a growth factor. Nucleic Acids Research. 2016 feb; 44(3):1370–1383.
952 http://dx.doi.org/10.1093/nar/gkv1367, doi: 10.1093/nar/gkv1367.

953 **Ermakova EO**, Nurtdinov RN, Gelfand MS. Fast rate of evolution in alternatively spliced coding regions of mam-
954 malian genes. BMC Genomics. 2006 apr; 7:84. http://dx.doi.org/10.1186/1471-2164-7-84, doi: 10.1186/1471-
955 2164-7-84.

956 **Galtier N**, Mouchiroud D. Isochore evolution in mammals: a human-like ancestral structure. Genetics. 1998
957 dec; 150(4):1577–1584. https://www.ncbi.nlm.nih.gov/pubmed/9832533.

958 **Gruner H**, Cortés-López M, Cooper DA, Bauer M, Miura P. CircRNA accumulation in the aging mouse brain.
959 Scientific Reports. 2016 dec; 6:38907. http://dx.doi.org/10.1038/srep38907, doi: 10.1038/srep38907.

960 **Gu W**, Ray DA, Walker JA, Barnes EW, Gentles AJ, Samollow PB, Jurka J, Batzer MA, Pollock DD. SINEs, evolution
961 and genome structure in the opossum. Gene. 2007 jul; 396(1):46–58. http://dx.doi.org/10.1016/j.gene.2007.
962 02.028, doi: 10.1016/j.gene.2007.02.028.

963 **Guo JU**, Agarwal V, Guo H, Bartel DP. Expanded identification and characterization of mammalian circular RNAs.
964 Genome Biology. 2014 jul; 15(7):409. http://dx.doi.org/10.1186/s13059-014-0409-z, doi: 10.1186/s13059-014-
965 0409-z.

966 **Hansen TB**, Jensen TI, Clausen BH, Bramsen JB, Finsen B, Damgaard CK, Kjems J. Natural RNA circles function as
967 efficient microRNA sponges. Nature. 2013 mar; 495(7441):384–388. http://dx.doi.org/10.1038/nature11993,
968 doi: 10.1038/nature11993.

**Heber S**, Alekseyev M, Sze SH, Tang H, Pevzner PA. Splicing graphs and EST assembly problem. Bioinformatics (Oxford, England). 2002; 18 Suppl 1:S181–8. http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=12169546&retmode=ref&cmd=prlinks, doi: 10.1093/bioinformatics/18.suppl_1.s181.

**Ivanov A**, Memczak S, Wyler E, Torti F, Porath HT, Orejuela MR, Piechotta M, Levanon EY, Landthaler M, Dieterich C, Rajewsky N. Analysis of intron sequences reveals hallmarks of circular RNA biogenesis in animals. Cell reports. 2015 jan; 10(2):170–177. http://dx.doi.org/10.1016/j.celrep.2014.12.019, doi: 10.1016/j.celrep.2014.12.019.

**Jeck WR**, Sorrentino JA, Wang K, Slevin MK, Burd CE, Liu J, Marzluff WF, Sharpless NE. Circular RNAs are abundant, conserved, and associated with ALU repeats. RNA (New York). 2013 feb; 19(2):141–157. http://dx.doi.org/10.1261/rna.035667.112, doi: 10.1261/rna.035667.112.

**Kim J**, Deininger PL. Recent amplification of rat ID sequences. Journal of Molecular Biology. 1996 aug; 261(3):322–327. http://dx.doi.org/10.1006/jmbi.1996.0464, doi: 10.1006/jmbi.1996.0464.

**Kim J**, Martignetti JA, Shen MR, Brosius J, Deininger P. Rodent BC1 RNA gene as a master gene for ID element amplification. Proceedings of the National Academy of Sciences of the United States of America. 1994 apr; 91(9):3607–3611. http://dx.doi.org/10.1073/pnas.91.9.3607, doi: 10.1073/pnas.91.9.3607.

**Koren A**, Polak P, Nemesh J, Michaelson JJ, Sebat J, Sunyaev SR, McCarroll SA. Differential relationship of DNA replication timing to different forms of human mutation and variation. American Journal of Human Genetics. 2012 dec; 91(6):1033–1040. http://dx.doi.org/10.1016/j.ajhg.2012.10.018, doi: 10.1016/j.ajhg.2012.10.018.

**Kramer MC**, Liang D, Tatomer DC, Gold B, March ZM, Cherry S, Wilusz JE. Combinatorial control of Drosophila circular RNA expression by intronic repeats, hnRNPs, and SR proteins. Genes & Development. 2015 oct; 29(20):2168–2182. http://dx.doi.org/10.1101/gad.270421.115, doi: 10.1101/gad.270421.115.

**Kristensen LS**, Andersen MS, Stagsted LVW, Ebbesen KK, Hansen TB, Kjems J. The biogenesis, biology and characterization of circular RNAs. Nature Reviews Genetics. 2019 aug; 20(11):675–691. http://dx.doi.org/10.1038/s41576-019-0158-7, doi: 10.1038/s41576-019-0158-7.

**Lee Y**, Choe J, Park OH, Kim YK. Molecular Mechanisms Driving mRNA Degradation by m(6)A Modification. Trends in genetics : TIG. 2020 Mar; 36(3):177–188. http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=31964509&retmode=ref&cmd=prlinks, doi: 10.1016/j.tig.2019.12.007.

**Lev-Maor G**, Ram O, Kim E, Sela N, Goren A, Levanon EY, Ast G. Intronic Alus influence alternative splicing. PLoS Genetics. 2008 sep; 4(9):e1000204. http://dx.doi.org/10.1371/journal.pgen.1000204, doi: 10.1371/journal.pgen.1000204.

**Li S**, Li X, Xue W, Zhang L, Yang LZ, Cao SM, Lei YN, Liu CX, Guo SK, Shan L, Wu M, Tao X, Zhang JL, Gao X, Zhang J, Wei J, Li J, Yang L, Chen LL. Screening for functional circular RNAs using the CRISPR–Cas13 system. Nature Methods. 2020; https://doi.org/10.1038/s41592-020-01011-4.

**Liang D**, Wilusz JE. Short intronic repeat sequences facilitate circular RNA production. Genes & Development. 2014 oct; 28(20):2233–2247. http://dx.doi.org/10.1101/gad.251926.114, doi: 10.1101/gad.251926.114.

**Melamud E**, Moult J. Structural implication of splicing stochastics. Nucleic Acids Research. 2009 aug; 37(14):4862–4872. http://dx.doi.org/10.1093/nar/gkp444, doi: 10.1093/nar/gkp444.

**Memczak S**, Jens M, Elefsinioti A, Torti F, Krueger J, Rybak A, Maier L, Mackowiak SD, Gregersen LH, Munschauer M, Loewer A, Ziebold U, Landthaler M, Kocks C, le Noble F, Rajewsky N. Circular RNAs are a large class of animal RNAs with regulatory potency. Nature. 2013 mar; 495(7441):333–338. http://dx.doi.org/10.1038/nature11928, doi: 10.1038/nature11928.

**Memczak S**, Papavasileiou P, Peters O, Rajewsky N. Identification and characterization of circular rnas as a new class of putative biomarkers in human blood. Plos One. 2015 oct; 10(10):e0141214. http://dx.doi.org/10.1371/journal.pone.0141214, doi: 10.1371/journal.pone.0141214.

**Mikkelsen TS**, Wakefield MJ, Aken B, Amemiya CT, Chang JL, Duke S, Garber M, Gentles AJ, Goodstadt L, Heger A, Jurka J, Kamal M, Mauceli E, Searle SMJ, Sharpe T, Baker ML, Batzer MA, Benos PV, Belov K, Clamp M, et al. Genome of the marsupial Monodelphis domestica reveals innovation in non-coding sequences. Nature. 2007 may; 447(7141):167–177. http://dx.doi.org/10.1038/nature05805, doi: 10.1038/nature05805.

**Modrek B**, Lee CJ. Alternative splicing in the human, mouse and rat genomes is associated with an increased frequency of exon creation and/or loss. Nature Genetics. 2003 jun; 34(2):177–180. http://dx.doi.org/10.1038/ng1159, doi: 10.1038/ng1159.

**Okholm TLH**, Sathe S, Park SS, Kamstrup AB, Rasmussen AM, Shankar A, Chua ZM, Fristrup N, Nielsen MM, Vang S, Dyrskjøt L, Aigner S, Damgaard CK, Yeo GW, Pedersen JS. Transcriptome-wide profiles of circular RNA and RNA-binding protein interactions reveal effects on circular RNA biogenesis and cancer pathway expression. Genome Medicine. 2020; 12(1):112. https://doi.org/10.1186/s13073-020-00812-8.

**Pai AA**, Cain CE, Mizrahi-Man O, De Leon S, Lewellen N, Veyrieras JB, Degner JF, Gaffney DJ, Pickrell JK, Stephens M, Pritchard JK, Gilad Y. The contribution of RNA decay quantitative trait loci to inter-individual variation in steady-state gene expression levels. PLoS Genetics. 2012 oct; 8(10):e1003000. http://dx.doi.org/10.1371/journal.pgen.1003000, doi: 10.1371/journal.pgen.1003000.

**Park OH**, Ha H, Lee Y, Boo SH, Kwon DH, Song HK, Kim YK. Endoribonucleolytic Cleavage of m(6)A-Containing RNAs by RNase P/MRP Complex. Molecular cell. 2019 May; 74(3):494–507.e8. http://eutils.ncbi.nlm.nih.gov/entrez/eutils/elink.fcgi?dbfrom=pubmed&id=30930054&retmode=ref&cmd=prlinks, doi: 10.1016/j.molcel.2019.02.034.

**Patop IL**, Wüst S, Kadener S. Past, present, and future of circRNAs. The EMBO Journal. 2019; 38(16):e100836. https://www.embopress.org/doi/abs/10.15252/embj.2018100836, doi: https://doi.org/10.15252/embj.2018100836.

**Piwecka M**, Glažar P, Hernandez-Miranda LR, Memczak S, Wolf SA, Rybak-Wolf A, Filipchyk A, Klironomos F, Cerda Jara CA, Fenske P, Trimbuch T, Zywitza V, Plass M, Schreyer L, Ayoub S, Kocks C, Kühn R, Rosenmund C, Birchmeier C, Rajewsky N. Loss of a mammalian circular RNA locus causes miRNA deregulation and affects brain function. Science. 2017 sep; 357(6357). http://dx.doi.org/10.1126/science.aam8526, doi: 10.1126/science.aam8526.

**Pollard KS**, Hubisz MJ, Rosenbloom KR, Siepel A. Detection of nonneutral substitution rates on mammalian phylogenies. Genome Research. 2010 jan; 20(1):110–121. http://dx.doi.org/10.1101/gr.097857.109, doi: 10.1101/gr.097857.109.

**Rybak-Wolf A**, Stottmeister C, Glažar P, Jens M, Pino N, Giusti S, Hanan M, Behm M, Bartok O, Ashwal-Fluss R, Herzog M, Schreyer L, Papavasileiou P, Ivanov A, Öhman M, Refojo D, Kadener S, Rajewsky N. Circular RNAs in the mammalian brain are highly abundant, conserved, and dynamically expressed. Molecular Cell. 2015 jun; 58(5):870–885. http://dx.doi.org/10.1016/j.molcel.2015.03.027, doi: 10.1016/j.molcel.2015.03.027.

**Salari R**, Wojtowicz D, Zheng J, Levens D, Pilpel Y, Przytycka TM. Teasing apart translational and transcriptional components of stochastic variations in eukaryotic gene expression. PLoS Computational Biology. 2012 aug; 8(8):e1002644. http://dx.doi.org/10.1371/journal.pcbi.1002644, doi: 10.1371/journal.pcbi.1002644.

**Siepel A**, Bejerano G, Pedersen JS, Hinrichs AS, Hou M, Rosenbloom K, Clawson H, Spieth J, Hillier LW, Richards S, Weinstock GM, Wilson RK, Gibbs RA, Kent WJ, Miller W, Haussler D. Evolutionarily conserved elements in vertebrate, insect, worm, and yeast genomes. Genome Research. 2005 aug; 15(8):1034–1050. http://dx.doi.org/10.1101/gr.3715005, doi: 10.1101/gr.3715005.

**Sievers F**, Wilm A, Dineen D, Gibson TJ, Karplus K, Li W, Lopez R, McWilliam H, Remmert M, Söding J, Thompson JD, Higgins DG. Fast, scalable generation of high-quality protein multiple sequence alignments using Clustal Omega. Molecular Systems Biology. 2011 oct; 7:539. http://dx.doi.org/10.1038/msb.2011.75, doi: 10.1038/msb.2011.75.

**Smit A**, Hubley R, Green P. RepeatMasker Open-4.0, 2013-2015. . 2013; http://www.repeatmasker.org.

**Starke S**, Jost I, Rossbach O, Schneider T, Schreiner S, Hung LH, Bindereif A. Exon circularization requires canonical splice signals. Cell reports. 2015; .

**Steinberg J**, Honti F, Meader S, Webber C. Haploinsufficiency predictions without study bias. Nucleic Acids Research. 2015 sep; 43(15):e101. http://dx.doi.org/10.1093/nar/gkv474, doi: 10.1093/nar/gkv474.

**Trapnell C**, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL, Pachter L. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. Nature Protocols. 2012 mar; 7(3):562–578. http://dx.doi.org/10.1038/nprot.2012.016, doi: 10.1038/nprot.2012.016.

**VenøMT**, Hansen TB, VenøST, Clausen BH, Grebing M, Finsen B, Holm IE, Kjems J. Spatio-temporal regulation of circular RNA expression during porcine embryonic brain development. Genome Biology. 2015 nov; 16:245. http://dx.doi.org/10.1186/s13059-015-0801-3, doi: 10.1186/s13059-015-0801-3.

**Wang M**, Hou J, Müller-McNicoll M, Chen W, Schuman EM. Long and Repeat-Rich Intronic Sequences Favor Circular RNA Formation under Conditions of Reduced Spliceosome Activity. iScience. 2019 oct; 20:237–247. http://dx.doi.org/10.1016/j.isci.2019.08.058, doi: 10.1016/j.isci.2019.08.058.

**Wang PL**, Bao Y, Yee MC, Barrett SP, Hogan GJ, Olsen MN, Dinneny JR, Brown PO, Salzman J. Circular RNA is expressed across the eukaryotic tree of life. Plos One. 2014 mar; 9(6):e90859. http://dx.doi.org/10.1371/journal.pone.0090859, doi: 10.1371/journal.pone.0090859.

**Werfel S**, Nothjunge S, Schwarzmayr T, Strom TM, Meitinger T, Engelhardt S. Characterization of circular RNAs in human, mouse and rat hearts. Journal of Molecular and Cellular Cardiology. 2016 jul; 98:103–107. http://dx.doi.org/10.1016/j.yjmcc.2016.07.007, doi: 10.1016/j.yjmcc.2016.07.007.

**Westholm JO**, Miura P, Olson S, Shenker S, Joseph B, Sanfilippo P, Celniker SE, Graveley BR, Lai EC. Genome-wide analysis of drosophila circular RNAs reveals their structural and sequence properties and age-dependent neural accumulation. Cell reports. 2014 dec; 9(5):1966–1980. http://dx.doi.org/10.1016/j.celrep.2014.10.062, doi: 10.1016/j.celrep.2014.10.062.

**Wilusz JE**. Repetitive elements regulate circular RNA biogenesis. Mobile genetic elements. 2015 jun; 5(3):1–7. http://dx.doi.org/10.1080/{2159256X}.2015.1045682, doi: 10.1080/{2159256X}.2015.1045682.

**Xu K**, Chen D, Wang Z, Ma J, Zhou J, Chen N, Lv L, Zheng Y, Hu X, Zhang Y, Li J. Annotation and functional clustering of circRNA expression in rhesus macaque brain during aging. Cell discovery. 2018 sep; 4:48. http://www.nature.com/articles/s41421-018-0050-1, doi: 10.1038/s41421-018-0050-1.

**Yoder JA**, Walsh CP, Bestor TH. Cytosine methylation and the ecology of intragenomic parasites. Trends in Genetics. 1997 aug; 13(8):335–340. http://dx.doi.org/10.1016/s0168-9525(97)01181-5, doi: 10.1016/s0168-9525(97)01181-5.

**Yoshimoto R**, Rahimi K, Hansen TB, Kjems J, Mayeda A. Biosynthesis of Circular RNA ciRS-7/CDR1as Is Mediated by Mammalian-wide Interspersed Repeats. iScience. 2020; 23(7):101345. http://www.sciencedirect.com/science/article/pii/S2589004220305320, doi: https://doi.org/10.1016/j.isci.2020.101345.

**You X**, Vlatkovic I, Babic A, Will T, Epstein I, Tushev G, Akbalik G, Wang M, Glock C, Quedenau C, Wang X, Hou J, Liu H, Sun W, Sambandan S, Chen T, Schuman EM, Chen W. Neural circular RNAs are derived from synaptic genes and regulated by development and plasticity. Nature Neuroscience. 2015 apr; 18(4):603–610. http://dx.doi.org/10.1038/nn.3975, doi: 10.1038/nn.3975.

1097 **Zaccara S**, Ries RJ, Jaffrey SR. Reading, writing and erasing mRNA methylation. Nature reviews Molecular cell
1098    biology. 2019; 20(10):608–624. https://doi.org/10.1038/s41580-019-0168-5.

1099 **Zhang XO**, Wang HB, Zhang Y, Lu X, Chen LL, Yang L. Complementary sequence-mediated exon circularization.
1100    Cell. 2014 sep; 159(1):134–147. http://dx.doi.org/10.1016/j.cell.2014.09.001, doi: 10.1016/j.cell.2014.09.001.

1101 **Zhang Y**, Romanish MT, Mager DL. Distributions of transposable elements reveal hazardous zones in mam-
1102    malian introns. PLoS Computational Biology. 2011 may; 7(5):e1002046. http://dx.doi.org/10.1371/journal.
1103    pcbi.1002046, doi: 10.1371/journal.pcbi.1002046.

1104 **Zhou C**, Molinie B, Daneshvar K, Pondick JV, Wang J, Van Wittenberghe N, Xing Y, Giallourakis CC, Mullen AC.
1105    Genome-Wide Maps of m6A circRNAs Identify Widespread and Cell-Type-Specific Methylation Patterns that
1106    Are Distinct from mRNAs. Cell reports. 2017 aug; 20(9):2262–2276. http://dx.doi.org/10.1016/j.celrep.2017.
1107    08.027, doi: 10.1016/j.celrep.2017.08.027.

1108 **Zhu L**, Zhang Y, Zhang W, Yang S, Chen JQ, Tian D. Patterns of exon-intron architecture variation of genes
1109    in eukaryotic genomes. BMC Genomics. 2009 jan; 10:47. http://dx.doi.org/10.1186/1471-2164-10-47, doi:
1110    10.1186/1471-2164-10-47.

1111    **Figure 1–Figure supplement 1.**

1112    **Figure 1–Figure supplement 2.**

1113    **Figure 1–Figure supplement 3.**

1114    **Figure 1–Figure supplement 4.**

1115    **Figure 2–Figure supplement 1.**

1116    **Figure 2–Figure supplement 2.**

1117    **Figure 3–Figure supplement 1.**

1118    **Figure 3–Figure supplement 2.**

1119    **Figure 3–Figure supplement 3.**

1120    **Figure 4–Figure supplement 1.**

1121    **Figure 4–Figure supplement 2.**

1122    **Figure 5–Figure supplement 1.**