

1 **Title: Insights into mammalian TE diversity via the curation of 248 mammalian genome**
2 **assemblies**

3 **Authors:** Austin B. Osmanski¹, Nicole S. Paulat¹, Jenny Korstian¹, Jenna R. Grimshaw¹,
4 Michaela Halsey¹, Kevin A.M. Sullivan¹, Diana D. Moreno-Santillán¹, Claudia Crookshanks¹,
5 Jacquelyn Roberts¹, Carlos Garcia¹, Matthew G. Johnson¹, Llewellyn D. Densmore¹, Richard D.
6 Stevens², Zoonomia Consortium, Jeb Rosen³, Jessica M. Storer³, Robert Hubley³, Arian F.A.
7 Smit³, Liliana M. Dávalos^{4,5}, Kerstin Lindblad-Toh^{6,7}, Elinor K. Karlsson^{7,8,9}, David A. Ray^{1*}

8 **Affiliations:**

9 ¹Department of Biological Sciences, Texas Tech University, Lubbock, TX, USA

10 ²Department of Natural Resources Management and Natural Science Research Laboratory,
11 Museum of Texas Tech University, Lubbock, TX 79409, USA

12 ³Institute for Systems Biology, Seattle, WA, USA

13 ⁴Department of Ecology & Evolution, Stony Brook University, Stony Brook, NY, USA

14 ⁵Consortium for Inter-Disciplinary Environmental Research, Stony Brook University; Stony
15 Brook, NY, USA

16 ⁶Department of Medical Biochemistry and Microbiology, Science for Life Laboratory, Uppsala
17 University; Uppsala, 751 32, Sweden

18 ⁷Broad Institute of MIT and Harvard; Cambridge, MA 02139, USA

19 ⁸Program in Bioinformatics and Integrative Biology, UMass Chan Medical School; Worcester,
20 MA 01605, USA

21 ⁹Program in Molecular Medicine, UMass Chan Medical School; Worcester, MA 01605, USA

22

23 *Correspondence: david.4.ray@gmail.com

24 **Abstract:**

25 We examined transposable element (TE) content of 248 placental mammal genome
26 assemblies, the largest *de novo* TE curation effort in eukaryotes to date. We find that while
27 mammals resemble one another in total TE content and diversity, they show substantial
28 differences with regard to recent TE accumulation. This includes multiple recent expansion and
29 quiescence events across the mammalian tree. Young TEs, particularly LINES, drive increases in
30 genome size while DNA transposons are associated with smaller genomes. Mammals tend to

31 accumulate only a few types of TE at any given time, with one TE type dominating. We also
32 found association between dietary habit and the presence of DNA transposon invasions. These
33 detailed annotations will serve as a benchmark for future comparative TE analyses among
34 placental mammals.

35 **One-Sentence Summary:** A *de novo* assessment of TE content in 248 mammals finds
36 informative trends in mammalian genome evolution.

37 **Main Text:**

38 Barbara McClintock became a scientific pioneer in the field of genomics with her Nobel
39 Prize winning discovery of transposable elements (TEs), DNA sequences that can mobilize
40 themselves in host genomes (1). A ubiquitous component of nearly all eukaryotes (2), TEs are
41 typically classified into two major groups based on their mobilization mechanism (3). Class I
42 elements, also known as retrotransposons, utilize an RNA intermediate during transposition
43 allowing replication throughout the genome in a copy-&-paste style of mobility (4). Class I
44 elements can be sorted further into three subcategories: Short INterspersed Elements (SINEs),
45 Long INterspersed Elements (LINEs), and LTR-retrotransposons (5). SINEs are non-
46 autonomous elements and depend on the presence of functional LINE elements, which contain
47 anywhere from 1-3 open reading frames encoding the necessary proteins for mobilization. Class
48 II elements, also known as DNA transposons, employ a DNA intermediate and can also be
49 subdivided. TIR-like DNA transposons such as hATs, piggyBacs, and TcMariner transposons
50 utilize a cut-&-paste mechanism by using transposase enzymes to catalyze the TE's relocation
51 (6). Helitrons, a second subcategory of Class II elements utilize a rolling-circle mechanism (7).
52 The final subcategory of known DNA transposons are Maverick elements, which are thought to

53 be derived from viruses as they have homologous genes coding for DNA polymerase and
54 retroviral-like integrase (8).

55 An increase in activity from either class of elements can lead to drastic alterations in
56 genome architecture (9). A variety of changes, including insertions, duplications, translocations,
57 deletions, and inversions can result from TE mobilization and accumulation (9). For instance, the
58 *AMAC1* (acyl-malonyl condensing enzyme 1) gene, coding for a protein essential for breaking
59 down phytanic acid from meat and dairy foods, has undergone multiple recent gene duplications
60 mediated by SVA retrotransposons in the human genome (10, 11). In addition to these structural
61 variants, the proliferative mechanisms of TE mobilization tend to cause eukaryotic genome sizes
62 to linearly correlate with TE abundance (2).

63 Increasing evidence indicates that TE-derived sequences have substantially influenced
64 the evolutionary histories of the organisms they occupy, even contributing to major evolutionary
65 innovations benefitting host organisms. Examples include recent TE insertions into genes
66 involved with insecticide resistance of the cotton bollworm (12); the rapid adaptation leading to
67 melanistic phenotypes of peppered moths in the soot ridden environment of British
68 industrialization (13); and the myriad of endogenous retroviruses that have contributed novel
69 regulatory functions to the development and evolution of the mammalian placenta (9, 14). The
70 overwhelming majority of TE insertions, however, result in selectively neutral alterations in
71 genome architecture, often showing no perceptible effect on host fitness (15). That being said,
72 deleterious insertions occur and impairments in gene function are possible outcomes of TE
73 mobilization which can lead to a wide variety of genetic diseases (9).

74 As a result, numerous genomic TE defense mechanisms have evolved to combat TE
75 activity by either regulating TE transcription or by targeting their intermediates to prevent

76 integration into the genome (3). These defense mechanisms explain, in part and in some
77 organisms, why few TE families retain the ability to mobilize over long periods of evolutionary
78 time (16). For example, among the ~868,000 L1 insertions in the human genome, few are
79 thought to be retrotransposition-competent and many of these exhibit cell-type specific
80 mobilization profiles (3, 17). Alternative to or in conjunction with the aforementioned scenario
81 of low numbers of functionally mobile TEs among some categories of elements, genomic drift
82 and the corresponding effects of fixation events among bottlenecked populations gives rise to
83 another explanation for varying levels of TE accumulation in different genome assemblies (18).

84 All these facets suggest that determining TE dynamics is key to understanding how
85 genomes evolve and function. Thus, TE curation and annotation is one of the most important
86 initial investigative steps in any description of a *de novo* genome assembly. Unfortunately, this
87 step is often relegated to an afterthought rather than performing a time-intensive *de novo* TE
88 curation effort (19). As a result, many genome assemblies are misunderstood from a TE
89 perspective (19). As the scientific community improves genome sequencing and assembly, the
90 lack of thorough and accurate TE annotation promises to become a major problem, especially in
91 the face of the number of large-scale genome sequencing initiatives now underway (20-24).

92 The Zoonomia project, described in (24), represents an opportunity to gain substantial
93 knowledge of the diversity of TEs in an important vertebrate clade, Mammalia. Here, we fill this
94 knowledge gap by providing complete, *de novo* TE annotations of 248 Zoonomia mammalian
95 genome assemblies using homology, *de novo*, and manual annotation approaches.

96 **Results:**

97 *General TE trends among mammals*

98 RepeatModeler (25), a *de novo* TE discovery tools, was used to examine 248 mammalian
99 genome assemblies yielding 25,025 putative TE starting queries. After initial curation and
100 elimination of duplicates, an iterative curation process consisting of between 1 and 19 rounds of
101 detailed curation (19), depending on the species (see Methods), yielded a library consisting of
102 8,263 novel consensus sequences. That library was combined with known TEs to create a
103 comprehensive mammalian TE library. This library, consisting of 25,676 consensus sequences,
104 was used to mask all assemblies. The dynamics of TE biology and intricacies of TE detection
105 lend themselves to a degree of false detection. For example, some TE families are chimeras of
106 multiple elements or they may contain similar core sequence components. To evaluate the
107 potential for false positives, we took advantage of an idiosyncrasy of TE biology in bats. A
108 family of bats, the Vespertilionidae, is, to our knowledge the sole mammalian family to have
109 incorporated a type of rolling circle transposon, Helitrons, into their TE repertoire (3). True
110 Helitrons in mammals have not been detected outside of Vespertilionidae. Thus, any Helitrons
111 detected outside of vesper bats, would likely be a false positive. RepeatMasker (26) detected
112 Helitrons in non-vesper mammals at a rate of 0.0013 ± 0.0019 , suggesting a low false positive
113 rate.

114 Previous work has suggested that the largest single classifiable component of a typical
115 mammalian genome is TEs (27) and our data (Fig. 1) corroborate this. As noted previously by
116 Elliott & Gregory in 2015 (2), genome size linearly correlates with the percentage of TE content
117 within a genome and this is again supported (Fig. 1, Table S1). Overall, TE content in each of the
118 examined species ranges from a low of 27.6% in the star-nosed mole (*Condylura cristata*) to
119 74.5% in the armadillo (*Oryzomys azer*) (Table S2, Fig. 1), with a distinct tendency to cluster in
120 the middle of that range (average TE proportion: 45.6%, average genome size: 2.67Gb). The

121 hazel dormouse (*Muscardinus avellanarius*) and Brazilian guinea pig (*Cavia aperea*) represent
122 the extremes of this middle cluster with 65.8% and 28.1% total TE content, respectively.
123 Assembly quality may impact the accuracy of TE annotation, but we could find no statistically
124 significant trend among taxa. For example, lower quality assemblies as measured by N50 or
125 BUSCO completeness did not yield lower or higher rates of observed TE accumulation (Fig. S1
126 and S2).

127 ***TE variation among mammals***

128 When examining TE content from all categories across the mammalian tree, we find some
129 general trends. For example, SINEs and LTR retrotransposons are more prevalent in
130 Euarchontoglires while LINEs dominate most other lineages, especially the bovids (Fig. 2).
131 However, we find placental mammals are generally similar with regard to overall TE
132 proportions, reflecting the tendency to retain older insertions that occurred in the common
133 ancestor of mammals. LINEs and SINEs always make up most TE abundance both in copy
134 number and in total genomic percentage. LINEs occupy between 8.2% and 52.8% of the
135 genomes examined, averaging 22.6%. SINEs occupy on average 10.5% of the mammalian
136 genome (range 0.4%-32.1%) (Table S3) while LTR retrotransposons, DNA transposons, and
137 rolling-circle transposons (RC) are substantially rarer; 7.8% (range 2.0%-17.8%), 3.5% (range
138 0.5%-8.4%), and 0.5% (range 0.01%-19.7%), respectively.

139 Examination of younger insertions, those with divergences averaging <4% from their
140 respective consensus, provides a picture of these genomes that is more dynamic, revealing
141 substantial differences in accumulation from each category of TE (Table S4). Some lineages,
142 such as the pteropodid bats (*Pteropus alecto*, *P. vampyrus*, *Eidolon helvum*, and *Rousettus*
143 *aegyptiacus* in Fig. 2), exhibit essentially no recent accumulation by any TE category while

144 others have experienced massive expansions in one or more categories. The aardvark
145 (*Orycteropus afer*) and musk deer (*Moschus moschus*) for instance, show substantial LINE
146 accumulation over the past ~20 million years.

147 To examine these trends more closely, we conducted a redundancy analysis (RDA) for
148 both orders and families to identify the major axes of variation in TE composition that were
149 related to either order or family affiliation of taxa (Fig. 3). This analysis suggests a strong
150 phylogenetic component to variation in TE composition among clades at the levels of order and
151 family. Eleven orders of mammals were significantly correlated with at least one of the two axes
152 and these orders were quite variable in terms of association with different TE types. The first two
153 major axes of variation in TE accumulation in analyses examining orders accounted for
154 approximately 27.2% of the variation and this was highly significant ($P < 0.001$). The first major
155 axis was positively related to the number of young TEs generally, and to young LINEs, LTRs,
156 and SINEs, which are all obligately replicative. Unsurprisingly given this characteristic, genome
157 size was also positively correlated with this axis. This axis was negatively related to young DNA
158 transposons and young rolling circle transposons. The second major axis of TE composition
159 related to ordinal affiliation was positively related to the number of young DNA transposons,
160 rolling circle transposons, LINEs and young TEs more generally, but negatively related to young
161 LTRs, SINEs, and to genome size.

162 Similar associations are seen at the family level. Families of mammals accounted for
163 approximately 49.9% of variation in TE composition, and this was highly significant (Fig. 3; $P <$
164 0.001). As with orders, the first major axis of variation was positively related to the same
165 categories of TE and to genome size. Correlations of young DNA transposons and young rolling
166 circle TEs were weaker than for orders, likely due to the lineage specificity of those element

167 types (see below), while positive associations of all other TE types were stronger. The second
168 major axis was positively related to the number of young DNA transposons, rolling circle
169 transposons, LINEs, and young TEs generally and negatively related to genome size. Fourteen
170 families of mammals were significantly correlated with at least one of these two axes and these
171 families were variable in terms of association with different TE types.

172 *TE diversity*

173 An increasingly useful avenue of inquiry among whole-genome TE analyses draws from
174 community ecology (28). Of interest is the application of community diversity measures
175 rendered on a genomic scale (29). We followed these lines o
176 f inquiry by investigating the diversity of recent TEs in each genome by calculating two diversity
177 indices and applying them to our data, Shannon diversity index (30) and Pielou's *J* (31).
178 Shannon diversity is a measure of overall diversity in a population of objects while Pielou's *J*
179 measures evenness by incorporating the relative numbers of each object, in this case, TE types
180 (Table S5). Species with the highest diversity values include bats and rodents. Bat TE diversity
181 was driven primarily by recent expansion of DNA transposons among Craseonycteridae,
182 Vespertilionidae, Hipposideridae, Rhinolophidae, and Mollossidae and recent accumulation of
183 both DNA transposons and rolling circle transposons in Vespertilionidae (Fig. 4).

184 In rodents, higher diversity among recently inserted TEs was driven by accumulations in
185 LTR retrotransposons, which made up 10-53% of recent TE accumulation. The highest rate of
186 recent LTR accumulation among the rodents was seen in members of Cricetidae and *Cricetomys*
187 *gambianus*.

188 To investigate general trends in diversity index values in relation to TE accumulation
189 patterns, we plotted values from recently deposited TEs vs. each diversity index (Fig. 5).
190 Hierarchical Bayesian analyses indicate that both Shannon diversity and Pielou's J exhibit
191 significant negative relationships with increasing recent TE content; Shannon H (Fig. 5, Table
192 S6) and Pielou's J (Fig. 5, Table S7, Fig. S3). Thus, the downward trend in Pielou's J suggests
193 that mammalian genomes tend to accumulate individual TE types at any given period rather than
194 multiple TE types accumulating simultaneously. This is exemplified in the armadillo, where
195 LINEs are currently dominating the recently active mobilome while SINEs are the major recent
196 contributor to the greater cane rat (*Thryonomys swinderianus*) genome (Fig. 2). However, clades
197 of bats with recent DNA accumulation tend to refute this pattern.

198 ***DNA transposons and diet***

199 The lineage specificity of the DNA transposon diversity described above suggests
200 horizontal transfer (HT) as a potential source of novel TE invasions in certain mammalian
201 genomes. To investigate patterns that may explain how such HT events may occur, we examined
202 the potential for life history to play a role. We hypothesized that differences in diet may allow
203 select species to come in contact with vectors for TEs (14, 32), which increase the likelihood of
204 successful invasion of mammalian genomes. DNA transposon-rich food sources such as many
205 arthropods, and non-mammalian vertebrates may offer greater potential for HT to some species
206 compared to those that eat plants. Hierarchical Bayesian analyses indicate that carnivorous
207 mammals tend to accumulate more recent DNA transposons in their genomes than non-
208 carnivores (Fig. 6A; Table S8). This pattern is best exemplified in the cetartiodactyls (Fig. 6B).
209 Recent DNA transposon accumulation is seen on average 20x more among the cetaceans than
210 other artiodactyls. Carnivorous bats, however, did not have statistically higher accumulations of

211 recent DNA transposons than herbivorous bats (Fig. 6C). Our datasets of primates and rodents
212 did not reveal any statistical difference of recent DNA transposon accumulation between
213 herbivores and omnivores (Fig. 6D-E).

214 **Discussion**

215 As our ability to generate high quality genome assemblies in rapid succession improves,
216 the need to curate TEs in those assemblies will only increase. Toward that end, we performed a
217 *de novo* assessment of the TE content of 248 mammal genome assemblies in what is, to our
218 knowledge, the largest comprehensive TE curation effort to date. This represents an increase of
219 ~58% compared to known mammalian TEs in RepBase as of 2019, when we began. Given the
220 numerous impacts that TEs are known to have at multiple levels of genome organization and
221 function, this increased knowledge will serve as a particularly valuable resource for anyone
222 interested in mammalian genomics and evolution. The full set of TE consensus sequences is
223 available for download from the Dfam (33) database.

224 Previous work has noted that genome size among mammals is relatively constrained (34)
225 and this work does not contradict that observation. Despite this constraint, our effort notes that
226 there is substantial variation in rates of accumulation in the recent mammalian past. Indeed, we
227 found that there is substantial diversity in TE accumulation patterns among mammals, suggesting
228 distinct TE-induced pressures on those genomes over evolutionary time and, likely, distinct
229 differences in the ability of eutherians to defend their genomes against TEs. These differences
230 represent an excellent opportunity for future researchers to investigate how TE defenses evolve
231 and respond to differing TE loads.

232 Another avenue of such research is to further investigate TE accumulation through the
233 lens of ecology and environment, an idea that has been discussed previously (14). Our data
234 demonstrate that carnivorous lineages tend to harbor an excess of recently accumulated DNA
235 transposons when compared to herbivorous taxa. The tendency of meat-eating mammals to have
236 more recent DNA transposon accumulation as compared to their non-carnivorous counterparts
237 suggests diet may play a significant role in a genome's likelihood of experiencing HT from Class
238 II TEs. This scenario is supported in part by a recent analysis of HT in predator-prey pairs and
239 their shared parasites (32). Nevertheless, this finding is not uniform across mammalian orders
240 and those varying patterns may reflect defenses against TE invasion (3), less availability of TEs
241 in order-specific dietary items, or some combination of both.

242 Investigating mammalian TEs through the ecological lens also suggests that single TE
243 types tend to dominate the mobilome during any given period (Fig. 5). This scenario is consistent
244 with our current understanding of TE defense mechanisms. The current model of PIWI-mediated
245 TE defense suggests that a new TE may invade or arise in a genome and enjoy a period of
246 relatively unfettered mobilization. Eventually, the piRNA defenses generate an effective
247 response and dampen the new TE's impacts (16, 35, 36).

248 With regard to the prevalence of HT of DNA transposons in carnivores, our data support
249 the hypothesis that the prevalence of HT of DNA transposons may be a consequence of the
250 similar cellular environments of predator and prey and their necessarily shared environments and
251 frequent interactions. Recent research has demonstrated the role that viruses and blood-feeding
252 arthropods play in facilitating HT (14, 32). Frequent interactions would further facilitate HT by
253 bringing such vectors into contact with both predator and prey. The similar cellular environments
254 among animals (as opposed to mammals with plant-based diets) would further encourage the

255 ready transfer of DNA transposons, which are already more amenable to HT due to their
256 relatively weak dependence on a host's cellular machinery to mobilize (37).

257 In conclusion, the annotation data provided here is essential for answering future
258 questions related to emerging hypotheses around speciation such as the TE-Thrust Hypothesis,
259 the Epi-Transposon Hypotheses, or the Carrier SubPopulation Hypothesis (3, 38). As
260 anthropogenic change exacerbates the decline in effective population size for many of the
261 species in our dataset, transposable elements might be the reservoir of genomic mutagens that
262 future populations or species rely on.

263 **Materials and Methods**

264 Generating the mammalian TE library

265
266 A total of 248 genome assemblies of placental mammals were initially presented for analysis
267 (table S2). For six species, higher quality assemblies were available via Bat1k, a similar, large scale
268 genome sequencing and assembly effort (21). In those cases, we replaced the Zoonomia assembly with
269 the higher quality version. Some assemblies were not used in the development of our final mammalian
270 TE library due to one or more of the following reasons: 1) the assembly exhibited a low N50 value
271 (<20,000) resulting in short contigs which are unsuitable for identifying longer TEs, 2) multiple artifacts
272 of assembly error were observed at TE sites which yielded implausible consensus sequences, 3) a
273 thorough species-specific TE annotation had already been performed and is available from RepBase
274 (Genetic Information Research Institute) (39), previous work from our own laboratory, or work
275 conducted by a collaborator. This left us with 205 species as substrates for TE curation (table S2).

276 Mammalian genomes have only a minimal tendency to remove older TE insertions from the
277 genome (40). Thus, the majority of older TE families that mobilized in the common ancestor or early in
278 the mammalian diversification were likely already characterized through efforts that focused on any of

279 several model organisms such as human, mouse, rat, pig, dog, cat, and horse (41-47). To avoid wasted
280 effort on re-curation of these shared and previously described TEs, we focused our manual curation
281 efforts on identifying newer putative TEs that underwent relatively recent accumulation. We defined
282 such ‘young’ insertions as TEs with sequences with K2P genetic distances less than 4% when compared
283 to their respective consensus. For temporal orientation, a kimura divergence of 4% approximates 20mya
284 or less since insertion, based on a general mammalian neutral mutation rate of 2.2×10^{-9} (48). The use of
285 a general mutation rate allowed for consistency among K2P values in analyses, however it limits the
286 accuracy of species-specific temporal estimations due to varying neutral mutation rates among placental
287 mammals. Thus, results with divergence values of less than 4% are considered “young” and do not
288 provide exact dates. This approach yielded mostly lineage specific TEs, many of which were yet to be
289 described but some previously identified and shared elements were occasionally encountered (i.e. the
290 Tigger family of Tc Mariner transposons and others), suggesting that we did not miss older but
291 unidentified elements. Custom scripts associated with the identification of younger elements are
292 available on zenodo (49).

293 For details of the curation process, see previous work from (19). Briefly, for each iteration of
294 manual TE curation, new consensus sequences were generated from the 50 BLAST hits which shared the
295 highest sequence identity to the consensus used in our BLAST query for that iteration. Custom pipelines
296 accomplished this by aligning BLAST hits with MUSCLE (50), trimming alignments with trimAl (-gt 0.6 -
297 cons 60) (51), and estimating a consensus sequence with EBMOSS (cons -plurality 3 -identity 3) (52).
298 Files which resulted in fewer than 10 BLAST hits were discarded. To consider a consensus sequence
299 ‘complete,’ the alignment needed to exhibit a pattern of random sequence at both the 5’ and 3’ ends, or
300 after extension to a length of 7kb or greater, whichever came first.

301 Because the ubiquitous LINE-1 can introduce copies of any transcript into the genome,
302 mammalian genomes have an unusually high number of processed pseudogenes (53-55). Including these

303 in a repeat database would result in annotation of functional genes as TE copies. Comparisons with
304 protein (domain) databases (<https://www.ncbi.nlm.nih.gov/protein/>,
305 <https://useast.ensembl.org/index.html>) we found and removed 152 such entries, most characterized by
306 a poly A tail. Small structural RNAs often occur in higher copy numbers partially because they are also
307 substrates of LINE1 (56), and a further 49 entries were dismissed as models created from their genes
308 and pseudogenes.

309 Two or three copies of interspersed repeats with very high copy numbers, usually but not
310 exclusively SINEs, can often be found in tandem clusters. This occurs more than by chance do to target
311 site preferences. For example, LINE-1 dependent SINEs insert in A-rich DNA, and such sites are
312 introduced by their own poly A tails (57). These artifacts are often identified by de novo repeat finders
313 but can be recognized when studying the seed alignments. Models will also have been built for the
314 individual units and many copies will end at the joining region between the units, the joining region is
315 more variable than the rest of the model. Over 210 models were such artifacts and eliminated.

316 Because in mammals the majority of LTR elements are represented by solo LTRs (58), Dfam (33)
317 and Repbase (39) harbor separate models for the LTRs and the internal sequences. De novo repeat-
318 finders like RepeatModeler often produce full elements or reconstruct a (partial) LTR and a fragment of
319 the internal sequence. We split these models into their components, based on homology to well-defined
320 LTRs and the presence of tRNA primer binding sites.

321 The combined original library contained several redundant models. Recognizing that models
322 represent (fragments of) the same TE is complicated by incorrect base calls, indels, overextension, and
323 incompleteness of the reconstruction as well as by the evolution of class I TEs in the genome: copies
324 created at different evolutionary times or from different descendants of the ancestral TE (sometimes
325 subtly) differ. A solid test for redundancy is to match the genome to all related models simultaneously

326 and find that some models are always outcompeted by others or that models converge to the same
327 consensus sequence. This could only be accomplished once the database was finalized, so we applied
328 arbitrary but informed cutoffs. Before comparison to each other the low complexity tails of SINEs and
329 LINEs were set to a standard length and short overextensions were trimmed, based on the expected
330 signatures of terminal bases or target site duplications. Differences between models at possible (highly
331 mutagenic) CpG sites were ignored. Dependent on class and age, elements were removed with
332 alignment scores against another model with a more complete sequence or a better seed alignment that
333 were between 90-95% of the score against itself. Partially overlapping fragments of potentially the same
334 TE were not addressed at this point.

335 We eliminated duplicated entries only when they were built from the same assembly. The same
336 TE can be reconstructed from the genomes of different species if it was active before their speciation
337 time, but with our current approach we could not estimate if a repeat was shared or lineage-specific and
338 merely similar. Thus, in Dfam (33) each of the models of this study currently is associated with only one
339 species and will not be matched when a same model is present in another species library.

340 To confirm the TE type, each sequence in the library was subjected to a custom pipeline (49) which
341 used: blastx to confirm the presence of known ORFs in autonomous elements, RepBase (39) to identify
342 known elements, and TEclass (59) to predict the TE type. We also used structural criteria for categorizing
343 TEs. DNA transposons were identified as elements with visible terminal inverted repeats. Rolling circle
344 transposons were required to have identifiable ACTAG at one end. Putative SINEs were inspected for a
345 repetitive tail as well as A and B boxes. SINEs also were classified by comparison to a database of SINE
346 modules(33): 800 small RNA class III promoter regions, 150 core regions and 5500 3' ends of LINE
347 elements (which SINEs often share). LTR retrotransposons and solo LTRs were required to have
348 recognizable hallmarks, such as: TG, TGT, or TGTT at their 5' and the inverse at the 3' ends, and the
349 presence of a polyadenylation signal. LTR classes could often be assigned by (indirect) sequence

350 homology to a coding internal sequence, when present. After this process, 8263 models and their seed
351 alignments were submitted to Dfam (33).

352 Once the final mammalian TE library was created, we used RepeatMasker-4.1.0 to mask the
353 genome assemblies. Postprocessing of output was performed using the rm2bed.py utility included with
354 RepeatMasker, which merges overlapping hits and converts the output to bed format.

355
356
357
358
359

Plotting TE variation using ordination

360 To characterize the major axes of variation of young TE accumulation among taxa we conducted
361 a redundancy analysis for both orders and families. In these analyses, the number of base pairs
362 attributed to of each TE Type as well as genome size for each taxon (order or family) was the dependent
363 matrix and dummy variables (60) assigning a species to either family or order was the independent
364 matrix. Redundancy is a multivariate regression that aims to examine the amount of variation and its
365 statistical significance in the dependent matrix that can be accounted for by the independent matrix.
366 Associations among variables were quantified based on a correlation matrix and significance was
367 determined based on 9999 permutations of the original datasets. Redundancy analyses were
368 performed in Canoco version 5 (61).

369
370
371
372

Test for association between TE proportions and assembly size, two diversity indices, and diets

373 The three objectives of these analyses included: 1) quantifying the association, if any, between
374 the total TE proportion in genome and assembly size; 2) estimating the difference in proportions of
375 recently accumulated DNA transposons within a genome among species with different diets; 3) and
376 quantifying the association, if any, between recent TE proportion in a genome and two diversity indices.

377 Diversity indices

378

379 An increasingly useful avenue for characterizing TE accumulation draws on community ecology

380 (28). Of particular interest is the application of community diversity measures rendered on a genomic

381 scale (29). We followed these lines of inquiry by investigating recent TE diversity within each genome of

382 our dataset by calculating the Shannon Diversity Index of TE classes. Focusing on recently inserted TEs,

383 we summed the bases that were attributed to TEs with K2P values less than 4%. We then generated the

384 proportions (p_i) for each TE class attributed to the overall base pair total of recently inserted TEs. To

385 calculate the Shannon Diversity Index, H , per the equation.

386
$$H = - \sum_{i=1}^k (p_i) \log(p_i)$$

387 To calculate the evenness of recent TE accumulation among the 5 main categories of TEs, we

388 employed the ecological metric, Pielou's J , a measure of species evenness, where S was equal to the

389 total number of recent TE hits found within an assembly.

390
$$J = \frac{H}{\ln(S)}$$

391 Dietary data

392 We gathered diet classification from The Animal Diversity Web (animaldiversity.org) for 178

393 available mammals on the public database (table S8). The young DNA transposon dataset was then

394 compared against three diet types: carnivore, herbivore, and omnivore.

395

396

397 Hierarchical Bayesian analyses

398 A hierarchical Bayesian approach was adopted to simultaneously estimate the species-specific
399 structure of errors while estimating error for the beta-distributed proportion of TE in the genome. A
400 hierarchical approach is often called a mixed model in the literature, with cluster-specific effects called
401 “random”, and sample-wide effects called “fixed”. As different fields apply random and fixed to different
402 levels of the hierarchy, here we adopt the language of cluster-specific and sample-wide effects (62).
403 Analyses begin by modelling the proportion of genome as a function of the genome assembly size as a
404 beta-distributed variable (63):

$$405 \quad y_i \sim \text{Beta}(\mu, \phi)$$

406 In which μ is the mean, and ϕ relates to the variance such that:

$$407 \quad \text{var}_{[y]} = \frac{\mu(1-\mu)}{1+\phi}$$

408 Given observations Y , and covariate assembly size X :

$$409 \quad \text{logit}(\mu) = \log\left(\frac{\mu}{1-\mu}\right) = \beta X$$

410 Instead of a typical regression, in which observations are presumed to be independent, our
411 analyses account for the phylogenetic structure of the errors by including normally distributed species-
412 specific effects with phylogenetic errors (64) such that:

$$413 \quad a \sim N(0, \sigma_a^2 A)$$

414 In which the phylogenetic relationship matrix A (65) replaces the identity of observations for the
415 residuals. The same distribution of the response and its phylogenetic errors was applied across all
416 regressions.

417 Assembly sizes in base pairs were on the order of 10^9 . To enable efficient modeling, this
418 predictor was log10 transformed and then scaled (subtracting the mean and dividing by one standard

419 deviation). No other predictor variables were transformed. Analyses of the association between diet and
420 TE proportions used diet as a group-specific predictor.

421 To implement Bayesian sampling for these analyses, we used brms (66), a package that enables
422 coding models in R for implementation in the stan statistical language (67). We ran separate univariate
423 models for each set of predictors (assembly size, diet, Shannon's Diversity Index, and Pielou's Evenness
424 Index), with the proportion of TE in the genome as the response. The covariance matrix A was obtained
425 from the variance covariance matrix of the dated phylogeny (65) of sampled species. Models ran four
426 separate Markov chain Monte Carlo chains using a Hamiltonian Monte Carlo approach. Compared to
427 other Bayesian implementations, the HMC approach saves time in sampling parameter spaces by
428 generating efficient transitions spanning the posterior based on derivatives of the density function of
429 the model. We used the approach of (68) to estimate R^2 from hierarchical Bayesian models. This
430 approach divides the variance of the predicted values by the variance of predicted values plus the
431 expected variance of the errors.

432 **References and Notes:**

- 433 1. B. McClintock, The origin and behavior of mutable loci in maize. *Proceedings of the National*
434 *Academy of Sciences* **36**, 344 (1950).
- 435 2. T. A. Elliott, T. R. Gregory, Do larger genomes contain more diverse transposable elements? *BMC*
436 *Evolutionary Biology* **15**, 69 (2015).
- 437 3. R. N. Platt, 2nd, M. W. Vandewege, D. A. Ray, Mammalian transposable elements and their
438 impacts on genome evolution. *Chromosome Res* **26**, 25-43 (2018).
- 439 4. T. H. Eickbush, V. K. Jamburuthugoda, The diversity of retrotransposons and the properties of
440 their reverse transcriptases. *Virus Res* **134**, 221-234 (2008).
- 441 5. G. Bourque *et al.*, Ten things you should know about transposable elements. *Genome Biology*
442 **19**, 199 (2018).
- 443 6. V. V. Kapitonov, J. Jurka, Self-synthesizing DNA transposons in eukaryotes. *Proc Natl Acad Sci U S*
444 *A* **103**, 4540-4545 (2006).
- 445 7. J. Thomas, E. J. Pritham, Helitrons, the Eukaryotic Rolling-circle Transposable Elements.
446 *Microbiol Spectr* **3**, (2015).
- 447 8. E. J. Pritham, T. Putliwala, C. Feschotte, Mavericks, a novel class of giant transposable elements
448 widespread in eukaryotes and related to DNA viruses. *Gene* **390**, 3-17 (2007).
- 449 9. A. D. Senft, T. S. Macfarlan, Transposable elements shape the evolution of mammalian
450 development. *Nature Reviews Genetics*, (2021).

- 451 10. J. Xing *et al.*, Emergence of primate genes by retrotransposon-mediated sequence transduction.
452 *Proc Natl Acad Sci U S A* **103**, 17608-17613 (2006).
- 453 11. Y. Takata *et al.*, Phytanic acid in dairy products and risk of cancer: current evidence and future
454 directions. *The FASEB Journal* **31**, 790.737-790.737 (2017).
- 455 12. K. Klai *et al.*, Screening of *Helicoverpa armigera* Mobilome Revealed Transposable Element
456 Insertions in Insecticide Resistance Genes. *Insects* **11**, 879 (2020).
- 457 13. A. E. v. t. Hof *et al.*, The industrial melanism mutation in British peppered moths is a
458 transposable element. *Nature* **534**, 102-105 (2016).
- 459 14. C. Gilbert, C. Feschotte, Horizontal acquisition of transposable elements and viral sequences:
460 patterns and consequences. *Curr Opin Genet Dev* **49**, 15-24 (2018).
- 461 15. I. R. Arkhipova, Neutral Theory, Transposable Elements, and Eukaryotic Genome Evolution.
462 *Molecular biology and evolution* **35**, 1332-1337 (2018).
- 463 16. R. Kofler, K. A. Senti, V. Nolte, R. Tobler, C. Schlötterer, Molecular dissection of a natural
464 transposable element invasion. *Genome Res* **28**, 824-835 (2018).
- 465 17. C. Philippe *et al.*, Activation of individual L1 retrotransposon instances is restricted to cell-type
466 dependent permissive loci. *eLife* **5**, e13926 (2016).
- 467 18. A. Le Rouzic, P. Capy, The first steps of transposable elements invasion: parasitic strategy vs.
468 genetic drift. *Genetics* **169**, 1033-1043 (2005).
- 469 19. R. N. Platt, II, L. Blanco-Berdugo, D. A. Ray, Accurate Transposable Element Annotation Is Vital
470 When Analyzing New Genome Assemblies. *Genome Biology and Evolution* **8**, 403-410 (2016).
- 471 20. G. K. C. o. Scientists, Genome 10K: A Proposal to Obtain Whole-Genome Sequence for 10 000
472 Vertebrate Species. *Journal of Heredity* **100**, 659-674 (2009).
- 473 21. E. C. Teeling *et al.*, Bat Biology, Genomes, and the Bat1K Project: To Generate Chromosome-
474 Level Genomes for All Living Bat Species. *Annual Review of Animal Biosciences* **6**, 23-46 (2018).
- 475 22. J. Sills *et al.*, Creating a Buzz About Insect Genomes. *Science* **331**, 1386-1386 (2011).
- 476 23. J. Threlfall, M. Blaxter, Launching the Tree of Life Gateway. *Wellcome Open Res* **6**, 125-125
477 (2021).
- 478 24. D. P. Genereux *et al.*, A comparative genomics multitool for scientific discovery and
479 conservation. *Nature* **587**, 240-245 (2020).
- 480 25. A. F. Smit, Hubley, R., RepeatModeler Open-1.0. <http://www.repeatmasker.org>, (2008-2015).
- 481 26. A. F. Smit, Hubley, R., Green, P., Repeat-Masker Open-3.0. <http://www.repeatmasker.org>,
482 (2004).
- 483 27. A. F. Smit, Interspersed repeats and other mementos of transposable elements in mammalian
484 genomes. *Curr Opin Genet Dev* **9**, 657-663 (1999).
- 485 28. S. Venner, C. Feschotte, C. Biémont, Dynamics of transposable elements: towards a community
486 ecology of the genome. *Trends in genetics : TIG* **25**, 317-323 (2009).
- 487 29. J. Wang *et al.*, Gigantic Genomes Provide Empirical Tests of Transposable Element Dynamics
488 Models. *Genomics, Proteomics & Bioinformatics*, (2021).
- 489 30. I. F. Spellerberg, P. J. Fedor, A tribute to Claude Shannon (1916–2001) and a plea for more
490 rigorous use of species richness, species diversity and the ‘Shannon–Wiener’ Index. *Global*
491 *Ecology and Biogeography* **12**, 177-179 (2003).
- 492 31. E. C. Pielou, The measurement of diversity in different types of biological collections. *Journal of*
493 *Theoretical Biology* **13**, 131-144 (1966).
- 494 32. C. Kambayashi *et al.*, Geography-Dependent Horizontal Gene Transfer from Vertebrate
495 Predators to Their Prey. *Molecular Biology and Evolution* **39**, msac052 (2022).
- 496 33. J. Storer, R. Hubley, J. Rosen, T. J. Wheeler, A. F. Smit, The Dfam community resource of
497 transposable element families, sequence models, and genome annotations. *Mobile DNA* **12**, 2
498 (2021).

- 499 34. K. Bachmann, Genome size in mammals. *Chromosoma* **37**, 85-93 (1972).
- 500 35. R. Kofler, Dynamics of Transposable Element Invasions with piRNA Clusters. *Molecular Biology*
501 *and Evolution* **36**, 1457-1472 (2019).
- 502 36. S. Luo *et al.*, The evolutionary arms race between transposable elements and piRNAs in
503 *Drosophila melanogaster*. *BMC Evolutionary Biology* **20**, 14 (2020).
- 504 37. D. J. Lampe, M. E. Churchill, H. M. Robertson, A purified mariner transposase is sufficient to
505 mediate transposition in vitro. *Embo j* **15**, 5470-5479 (1996).
- 506 38. J. Jurka, W. Bao, K. K. Kojima, Families of transposable elements, population structure and the
507 origin of species. *Biol Direct* **6**, 44 (2011).
- 508 39. W. Bao, K. K. Kojima, O. Kohany, Repbase Update, a database of repetitive elements in
509 eukaryotic genomes. *Mobile DNA* **6**, 11 (2015).
- 510 40. A. Kapusta, A. Suh, C. Feschotte, Dynamics of genome size evolution in birds and mammals.
511 *Proceedings of the National Academy of Sciences* **114**, E1460-E1469 (2017).
- 512 41. C. International Human Genome Sequencing, Finishing the euchromatic sequence of the human
513 genome. *Nature* **431**, 931-945 (2004).
- 514 42. E. F. Kirkness *et al.*, The Dog Genome: Survey Sequencing and Comparative Analysis. *Science*
515 **301**, 1898-1903 (2003).
- 516 43. R. H. Waterston *et al.*, Initial sequencing and comparative analysis of the mouse genome. *Nature*
517 **420**, 520-562 (2002).
- 518 44. J. U. Pontius *et al.*, Initial sequence and comparative analysis of the cat genome. *Genome Res* **17**,
519 1675-1689 (2007).
- 520 45. R. A. Gibbs *et al.*, Genome sequence of the Brown Norway rat yields insights into mammalian
521 evolution. *Nature* **428**, 493-521 (2004).
- 522 46. M. A. M. Groenen *et al.*, Analyses of pig genomes provide insight into porcine demography and
523 evolution. *Nature* **491**, 393-398 (2012).
- 524 47. D. L. Adelson, J. M. Raison, M. Garber, R. C. Edgar, Interspersed repeats in the horse (*Equus*
525 *caballus*); spatial correlations highlight conserved chromosomal domains. *Animal Genetics* **41**,
526 91-99 (2010).
- 527 48. S. Kumar, S. Subramanian, Mutation rates in mammalian genomes. *Proc Natl Acad Sci U S A* **99**,
528 803-808 (2002).
- 529 49. A. B. Osmanski, Dávalos, L. M., Johnson, M. G., Hubley, R., Smit, A. F. A., Ray, D. A. ,
530 Zoonomia_TEs_Release_v1.0.0. zenodo.org/badge/latestdoi/431231925, (2022).
- 531 50. R. C. Edgar, MUSCLE: multiple sequence alignment with high accuracy and high throughput.
532 *Nucleic Acids Res* **32**, 1792-1797 (2004).
- 533 51. S. Capella-Gutiérrez, J. M. Silla-Martínez, T. Gabaldón, trimAl: a tool for automated alignment
534 trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972-1973 (2009).
- 535 52. P. Rice, I. Longden, A. Bleasby, EMBOSS: the European Molecular Biology Open Software Suite.
536 *Trends Genet* **16**, 276-277 (2000).
- 537 53. O. K. Pickeral, W. Makałowski, M. S. Boguski, J. D. Boeke, Frequent human genomic DNA
538 transduction driven by LINE-1 retrotransposition. *Genome Res* **10**, 411-415 (2000).
- 539 54. J. L. Goodier, E. M. Ostertag, H. H. Kazazian Jr, Transduction of 3'-flanking sequences is common
540 in L1 retrotransposition. *Human Molecular Genetics* **9**, 653-657 (2000).
- 541 55. C. Esnault, J. Maestre, T. Heidmann, Human LINE retrotransposons generate processed
542 pseudogenes. *Nature Genetics* **24**, 363-367 (2000).
- 543 56. C. R. Beck, J. L. Garcia-Perez, R. M. Badge, J. V. Moran, LINE-1 elements in structural variation
544 and disease. *Annu Rev Genomics Hum Genet* **12**, 187-215 (2011).
- 545 57. M. El-Sawy, P. Deininger, Tandem insertions of Alu elements. *Cytogenet Genome Res* **108**, 58-62
546 (2005).

- 547 58. J. Ma, K. M. Devos, J. L. Bennetzen, Analyses of LTR-retrotransposon structures reveal recent
548 and rapid genomic DNA loss in rice. *Genome Res* **14**, 860-869 (2004).
- 549 59. G. Abrusán, N. Grundmann, L. DeMester, W. Makalowski, TEclass—a tool for automated
550 classification of unknown eukaryotic transposable elements. *Bioinformatics* **25**, 1329-1330
551 (2009).
- 552 60. L. L. P. Legendre, *Numerical Ecology, Volume 24, 3rd Edition*. (Elsevier, Oxford, UK, 2012).
- 553 61. C. J. F. ter Braak, and Šmilauer, P., *Canoco reference manual and user's guide: software for*
554 *ordination, version 5.0*. (Microcomputer Power, Ithaca, USA, 2012).
- 555 62. A. Gelman, Analysis of variance - Why it is more important than ever. *Annals of Statistics* **33**, 1-
556 31 (2005).
- 557 63. J. C. Douma, J. T. Weedon, Analysing continuous proportions in ecology and evolution: A
558 practical introduction to beta and Dirichlet regression. *Methods in Ecology and Evolution* **10**,
559 1412-1430 (2019).
- 560 64. J. D. Hadfield, S. Nakagawa, General quantitative genetic methods for comparative biology:
561 phylogenies, taxonomies and multi-trait models for continuous and categorical characters. *J Evol*
562 *Biol* **23**, 494-508 (2010).
- 563 65. V. C. M. Nicole M. Foley, Andrew J. Harris, Kevin R. Bredemeyer, Joana Damas, Harris A. Lewin,
564 Eduardo Eizirik, John Gatesy, Zoonomia Consortium, Mark S. Springer, William J. Murphy, A
565 genomic timescale for placental mammal evolution. *Science*, (2021).
- 566 66. P.-C. Bürkner, brms: An R Package for Bayesian Multilevel Models Using Stan. *Journal of*
567 *Statistical Software* **80**, 1-28 (2017).
- 568 67. B. Carpenter *et al.*, Stan: A Probabilistic Programming Language. *2017* **76**, 32 (2017).
- 569 68. A. Gelman, B. Goodrich, J. Gabry, A. Vehtari, R-squared for Bayesian Regression Models. *The*
570 *American Statistician* **73**, 307-309 (2019).

571

572 **Acknowledgements:**

573 We thank the High-Performance Computing Center at Texas Tech University for providing
574 compute resources and technical support throughout the project. This work was also made
575 possible by the SeaWulf computing system from Stony Brook Research Computing and
576 Cyberinfrastructure, and the Institute for Advanced Computational Science at Stony Brook
577 University funded by National Science Foundation-OAC 1531492. We also thank Brittany Ann
578 Hale for providing artistic renditions of mammal taxa for our figure.

579 **Funding:** This project was partially supported by

580 National Science Foundation grant DEB 1838283 (DDMS, DAR)

581 National Science Foundation grant IOS 2032006 (DDMS, DAR)

582 National Institutes of Health grant R01HG002939 (JMS, RH, AS, JebR)

583 National Institutes of Health grant U24HG010136 (JMS, RH, AS, JebR)

584 National Science Foundation grant DEB 1838273 (LMD)

585 National Science Foundation grant DGE 1633299 (LMD)

586 National Institutes of Health NHGRI R01HG008742 (ZC)

587 Swedish Research Council Distinguished Professor Award (ZC)

588 **Author contributions:**

589 Conceptualization: ABO, DAR

590 Assembly generation: DDMS, LMD, DAR

591 Library validation & curation: NSP, JMS, ABO, KAMS, JK, JRG, MH, CG, CC, JR, JebR, RH,

592 AS, DAR

593 Methodology & Investigation: ABO, LMD, NSP, DAR

594 Writing – original draft: ABO NSP, DAR, RDS, LMD

595 Writing – review & editing: ABO, NSP, DAR, JMS, AS, RDS, LMD

596 **Competing interests:** Authors declare no competing interests.

597 **Data and materials availability:** All assemblies are available in Genbank, TE consensus

598 sequences are available via the Dfam database. All other data is available in the supplementary

599 materials; code used in the analysis is available at (49). .

600 **Zoonomia Consortium List:**

601 Gregory Andrews¹, Joel C. Armstrong², Matteo Bianchi³, Bruce W. Birren⁴, Kevin R.
602 Bredemeyer⁵, Ana M. Breit⁶, Matthew J. Christmas³, Hiram Clawson², Joana Damas⁷, Federica
603 Di Palma^{8,9}, Mark Diekhans², Michael X. Dong³, Eduardo Eizirik¹⁰, Kaili Fan¹, Cornelia
604 Fanter¹¹, Nicole M. Foley⁵, Karin Forsberg-Nilsson^{12,13}, Carlos J. Garcia¹⁴, John Gatesy¹⁵, Steven
605 Gazal¹⁶, Diane P. Genereux⁴, Linda Goodman¹⁷, Jenna Grimshaw¹⁴, Michaela K. Halsey¹⁴,
606 Andrew J. Harris⁵, Glenn Hickey¹⁸, Michael Hiller^{19,20,21}, Allyson G. Hindle¹¹, Robert M.
607 Hubley²², Graham M. Hughes²³, Jeremy Johnson⁴, David Juan²⁴, Irene M. Kaplow^{25,26}, Elinor K.
608 Karlsson^{1,4,27}, Kathleen C. Keough^{17,28,29}, Bogdan Kirilenko^{19,20,21}, Klaus-Peter Koepfli^{30,31,32},
609 Jennifer M. Korstian¹⁴, Amanda Kowalczyk^{25,26}, Sergey V. Kozyrev³, Alyssa J. Lawler^{4,26,33},
610 Colleen Lawless²³, Thomas Lehmann³⁴, Danielle L. Levesque⁶, Harris A. Lewin^{7,35,36}, Xue
611 Li^{1,4,37}, Abigail Lind^{28,29}, Kerstin Lindblad-Toh^{3,4}, Ava Mackay-Smith³⁸, Voichita D.
612 Marinescu³, Tomas Marques-Bonet^{39,40,41,42}, Victor C. Mason⁴³, Jennifer R. S. Meadows³, Wynn
613 K. Meyer⁴⁴, Jill E. Moore¹, Lucas R. Moreira^{1,4}, Diana D. Moreno-Santillan¹⁴, Kathleen M.
614 Morrill^{1,4,37}, Gerard Muntané²⁴, William J. Murphy⁵, Arcadi Navarro^{39,41,45,46}, Martin

615 Nweeia^{47,48,49,50}, Sylvia Ortmann⁵¹, Austin Osmanski¹⁴, Benedict Paten², Nicole S. Paulat¹⁴,
616 Andreas R. Pfenning^{25,26}, BaDoi N. Phan^{25,26,52}, Katherine S. Pollard^{28,29,53}, Henry E. Pratt¹,
617 David A. Ray¹⁴, Steven K. Reilly³⁸, Jeb R. Rosen²², Irina Ruf⁵⁴, Louise Ryan²³, Oliver A.
618 Ryder^{55,56}, Pardis C. Sabeti^{4,57,58}, Daniel E. Schäffer²⁵, Aitor Serres²⁴, Beth Shapiro^{59,60}, Arian F.
619 A. Smit²², Mark Springer⁶¹, Chaitanya Srinivasan²⁵, Cynthia Steiner⁵⁵, Jessica M. Storer²², Kevin
620 A. M. Sullivan¹⁴, Patrick F. Sullivan^{62,63}, Elisabeth Sundström³, Megan A. Supple⁵⁹, Ross
621 Swofford⁴, Joy-El Talbot⁶⁴, Emma Teeling²³, Jason Turner-Maier⁴, Alejandro Valenzuela²⁴,
622 Franziska Wagner⁶⁵, Ola Wallerman³, Chao Wang³, Juehan Wang¹⁶, Zhiping Weng¹, Aryn P.
623 Wilder⁵⁵, Morgan E. Wirthlin^{25,26,66}, James R. Xue^{4,57}, Xiaomeng Zhang^{4,25,26}

624

625 Affiliations:

626 ¹Program in Bioinformatics and Integrative Biology, UMass Chan Medical School; Worcester,
627 MA 01605, USA.

628 ²Genomics Institute, University of California Santa Cruz; Santa Cruz, CA 95064, USA.

629 ³Department of Medical Biochemistry and Microbiology, Science for Life Laboratory, Uppsala
630 University; Uppsala, 751 32, Sweden.

631 ⁴Broad Institute of MIT and Harvard; Cambridge, MA 02139, USA.

632 ⁵Veterinary Integrative Biosciences, Texas A&M University; College Station, TX 77843, USA.

633 ⁶School of Biology and Ecology, University of Maine; Orono, ME 04469, USA.

634 ⁷The Genome Center, University of California Davis; Davis, CA 95616, USA.

635 ⁸Genome British Columbia; Vancouver, BC, Canada.

636 ⁹School of Biological Sciences, University of East Anglia; Norwich, UK.

637 ¹⁰School of Health and Life Sciences, Pontifical Catholic University of Rio Grande do Sul; Porto
638 Alegre, 90619-900, Brazil.

639 ¹¹School of Life Sciences, University of Nevada Las Vegas; Las Vegas, NV 89154, USA.

640 ¹²Biodiscovery Institute, University of Nottingham; Nottingham, UK.

641 ¹³Department of Immunology, Genetics and Pathology, Science for Life Laboratory, Uppsala
642 University; Uppsala, 751 85, Sweden.

643 ¹⁴Department of Biological Sciences, Texas Tech University; Lubbock, TX 79409, USA.

644 ¹⁵Division of Vertebrate Zoology, American Museum of Natural History; New York, NY 10024,
645 USA.

646 ¹⁶Keck School of Medicine, University of Southern California; Los Angeles, CA 90033, USA.

647 ¹⁷Fauna Bio Incorporated; Emeryville, CA 94608, USA.

648 ¹⁸Baskin School of Engineering, University of California Santa Cruz; Santa Cruz, CA 95064,
649 USA.

650 ¹⁹Faculty of Biosciences, Goethe-University; 60438 Frankfurt, Germany.

651 ²⁰LOEWE Centre for Translational Biodiversity Genomics; 60325 Frankfurt, Germany.

652 ²¹Senckenberg Research Institute; 60325 Frankfurt, Germany.

653 ²²Institute for Systems Biology; Seattle, WA 98109, USA.

654 ²³School of Biology and Environmental Science, University College Dublin; Belfield, Dublin 4,
655 Ireland.

656 ²⁴Department of Experimental and Health Sciences, Institute of Evolutionary Biology (UPF-
657 CSIC), Universitat Pompeu Fabra; Barcelona, 08003, Spain.

658 ²⁵Department of Computational Biology, School of Computer Science, Carnegie Mellon
659 University; Pittsburgh, PA 15213, USA.

660 ²⁶Neuroscience Institute, Carnegie Mellon University; Pittsburgh, PA 15213, USA.

- 661 ²⁷Program in Molecular Medicine, UMass Chan Medical School; Worcester, MA 01605, USA.
662 ²⁸Department of Epidemiology & Biostatistics, University of California San Francisco; San
663 Francisco, CA 94158, USA.
664 ²⁹Gladstone Institutes; San Francisco, CA 94158, USA.
665 ³⁰Center for Species Survival, Smithsonian's National Zoo and Conservation Biology Institute;
666 Washington, DC 20008, USA.
667 ³¹Computer Technologies Laboratory, ITMO University; St. Petersburg 197101, Russia.
668 ³²Smithsonian-Mason School of Conservation, George Mason University; Front Royal, VA
669 22630, USA.
670 ³³Department of Biological Sciences, Mellon College of Science, Carnegie Mellon University;
671 Pittsburgh, PA 15213, USA.
672 ³⁴Senckenberg Research Institute and Natural History Museum Frankfurt; 60325 Frankfurt am
673 Main, Germany.
674 ³⁵Department of Evolution and Ecology, University of California Davis; Davis, CA 95616, USA.
675 ³⁶John Muir Institute for the Environment, University of California Davis; Davis, CA 95616,
676 USA.
677 ³⁷Morningside Graduate School of Biomedical Sciences, UMass Chan Medical School;
678 Worcester, MA 01605, USA.
679 ³⁸Department of Genetics, Yale School of Medicine; New Haven, CT 06510, USA.
680 ³⁹Catalan Institution of Research and Advanced Studies (ICREA); Barcelona, 08010, Spain.
681 ⁴⁰CNAG-CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology
682 (BIST); Barcelona, 08036, Spain.
683 ⁴¹Department of Medicine and Life Sciences, Institute of Evolutionary Biology (UPF-CSIC),
684 Universitat Pompeu Fabra; Barcelona, 08003, Spain.
685 ⁴²Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de Barcelona; 08193,
686 Cerdanyola del Vallès, Barcelona, Spain.
687 ⁴³Institute of Cell Biology, University of Bern; 3012, Bern, Switzerland.
688 ⁴⁴Department of Biological Sciences, Lehigh University; Bethlehem, PA 18015, USA.
689 ⁴⁵BarcelonaBeta Brain Research Center, Pasqual Maragall Foundation; Barcelona, 08005, Spain.
690 ⁴⁶CRG, Centre for Genomic Regulation, Barcelona Institute of Science and Technology (BIST);
691 Barcelona, 08003, Spain.
692 ⁴⁷Department of Comprehensive Care, School of Dental Medicine, Case Western Reserve
693 University; Cleveland, OH 44106, USA.
694 ⁴⁸Department of Vertebrate Zoology, Canadian Museum of Nature; Ottawa, Ontario K2P 2R1,
695 Canada.
696 ⁴⁹Department of Vertebrate Zoology, Smithsonian Institution; Washington, DC 20002, USA.
697 ⁵⁰Narwhal Genome Initiative, Department of Restorative Dentistry and Biomaterials Sciences,
698 Harvard School of Dental Medicine; Boston, MA 02115, USA.
699 ⁵¹Department of Evolutionary Ecology, Leibniz Institute for Zoo and Wildlife Research; 10315
700 Berlin, Germany.
701 ⁵²Medical Scientist Training Program, University of Pittsburgh School of Medicine; Pittsburgh,
702 PA 15261, USA.
703 ⁵³Chan Zuckerberg Biohub; San Francisco, CA 94158, USA.
704 ⁵⁴Division of Messel Research and Mammalogy, Senckenberg Research Institute and Natural
705 History Museum Frankfurt; 60325 Frankfurt am Main, Germany.
706 ⁵⁵Conservation Genetics, San Diego Zoo Wildlife Alliance; Escondido, CA 92027, USA.

707 ⁵⁶Department of Evolution, Behavior and Ecology, School of Biological Sciences, University of
708 California San Diego; La Jolla, CA 92039, USA.

709 ⁵⁷Department of Organismic and Evolutionary Biology, Harvard University; Cambridge, MA
710 02138, USA.

711 ⁵⁸Howard Hughes Medical Institute; Chevy Chase, MD, USA.

712 ⁵⁹Department of Ecology and Evolutionary Biology, University of California Santa Cruz; Santa
713 Cruz, CA 95064, USA.

714 ⁶⁰Howard Hughes Medical Institute, University of California Santa Cruz; Santa Cruz, CA 95064,
715 USA.

716 ⁶¹Department of Evolution, Ecology and Organismal Biology, University of California
717 Riverside; Riverside, CA 92521, USA.

718 ⁶²Department of Genetics, University of North Carolina Medical School; Chapel Hill, NC 27599,
719 USA.

720 ⁶³Department of Medical Epidemiology and Biostatistics, Karolinska Institutet; Stockholm,
721 Sweden.

722 ⁶⁴Iris Data Solutions, LLC; Orono, ME 04473, USA.

723 ⁶⁵Museum of Zoology, Senckenberg Natural History Collections Dresden; 01109 Dresden,
724 Germany.

725 ⁶⁶Allen Institute for Brain Science; Seattle, WA 98109, USA

726

727

728

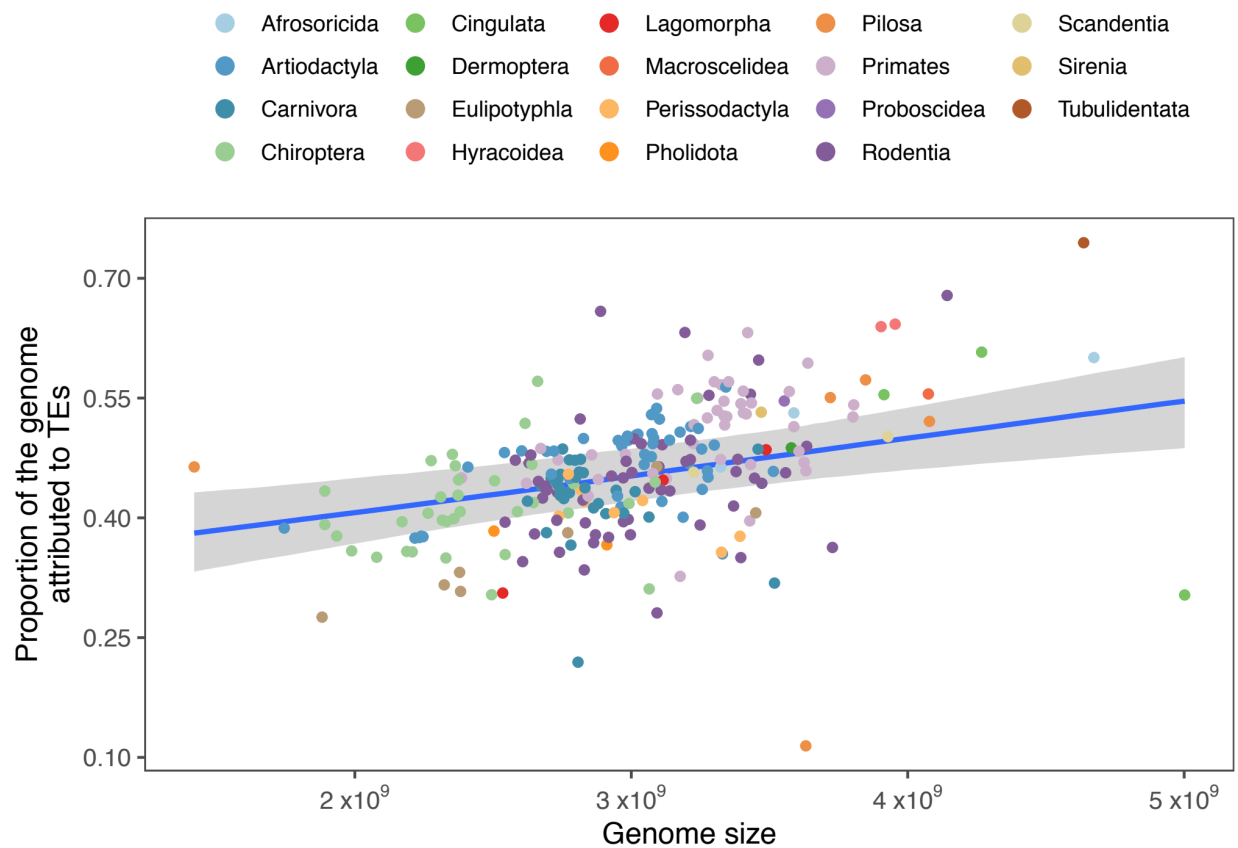
729 **Supplementary Materials**

730 Figs. S1-S5

731 Tables S1-S8

732

733 **Figure legends**

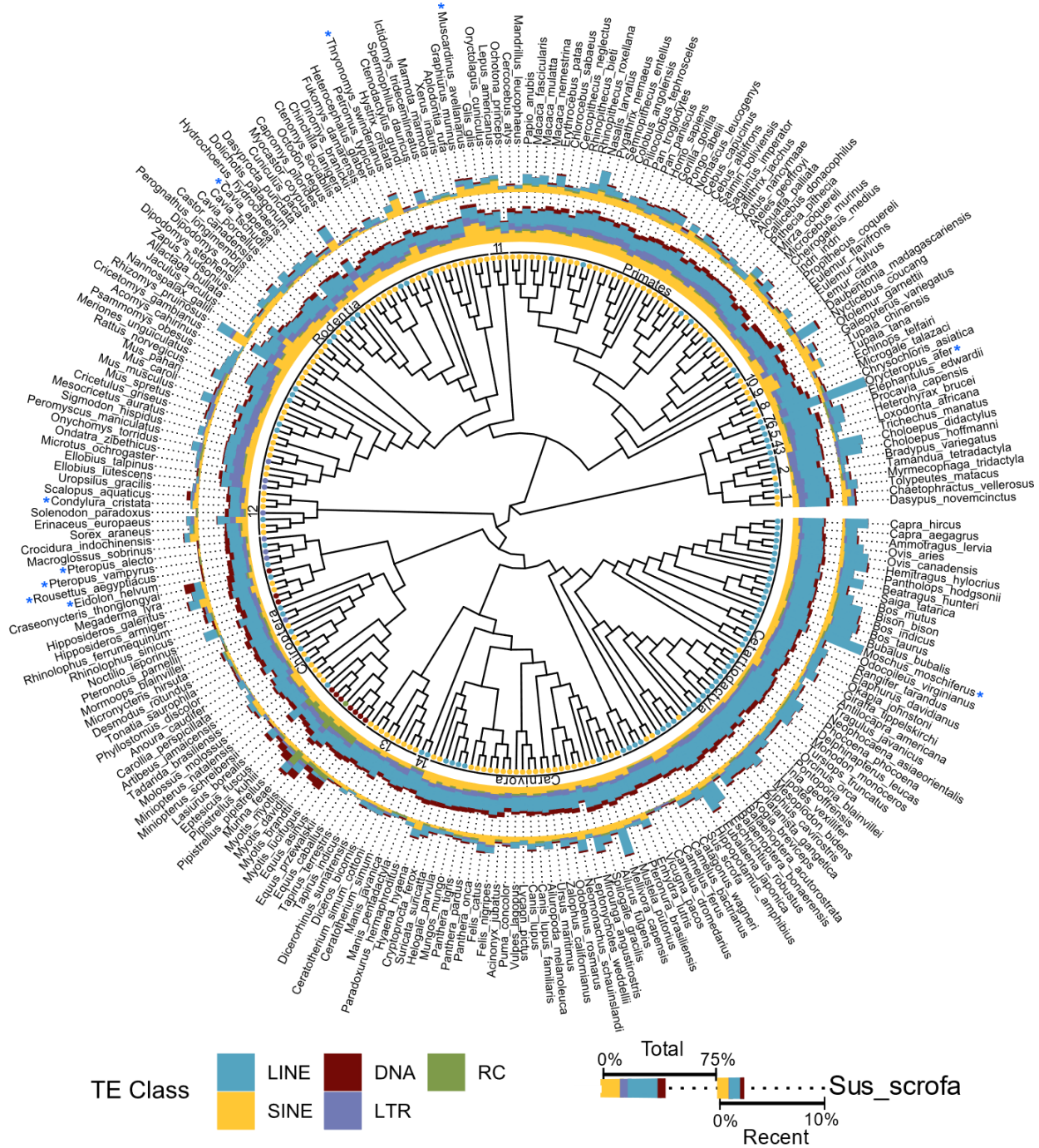


734

735 **Fig. 1. Correlation of total genomic TE content and the size, in base pairs, of the genome.**

736 Due to the log transformation and scaling of assembly size for the hierarchical Bayesian analysis,
737 and resulting back-transformation, the x-axis values are approximately rendered. Blue line
738 indicates the line of best-fit and shaded area is the 95% high probability density of the fit. The r^2
739 for this relationship was estimated at 0.54 (95% high probability density 0.42, 0.64).

740



741

742 **Fig. 2. Total and young TE genomic proportions by species within a phylogenetic context.**

743 Dots at branch tips indicate the TE class most prevalent among recent TE insertions (insertions

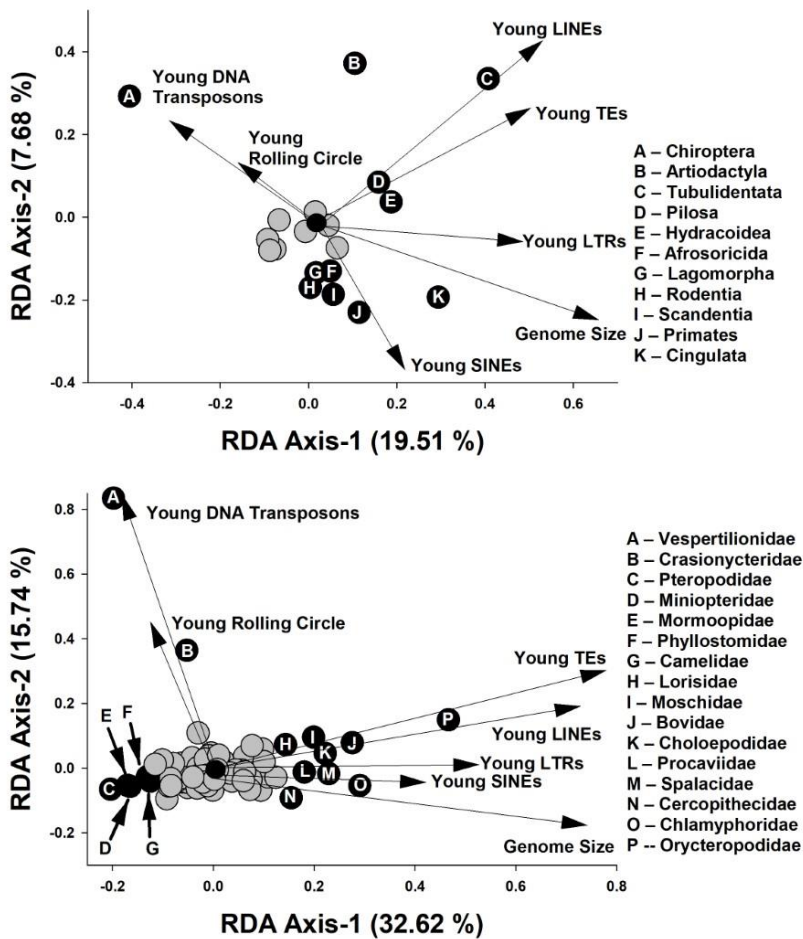
744 with <4% divergence from the relevant consensus TE). The ring immediately following the

745 branch tip dots indicates the mammalian order to reach respective species. Orders represented

746 by numbers include: 1) Cingulata, 2) Pilosa, 3) Sirenia, 4) Proboscidea, 5) Hyracoidea, 6)

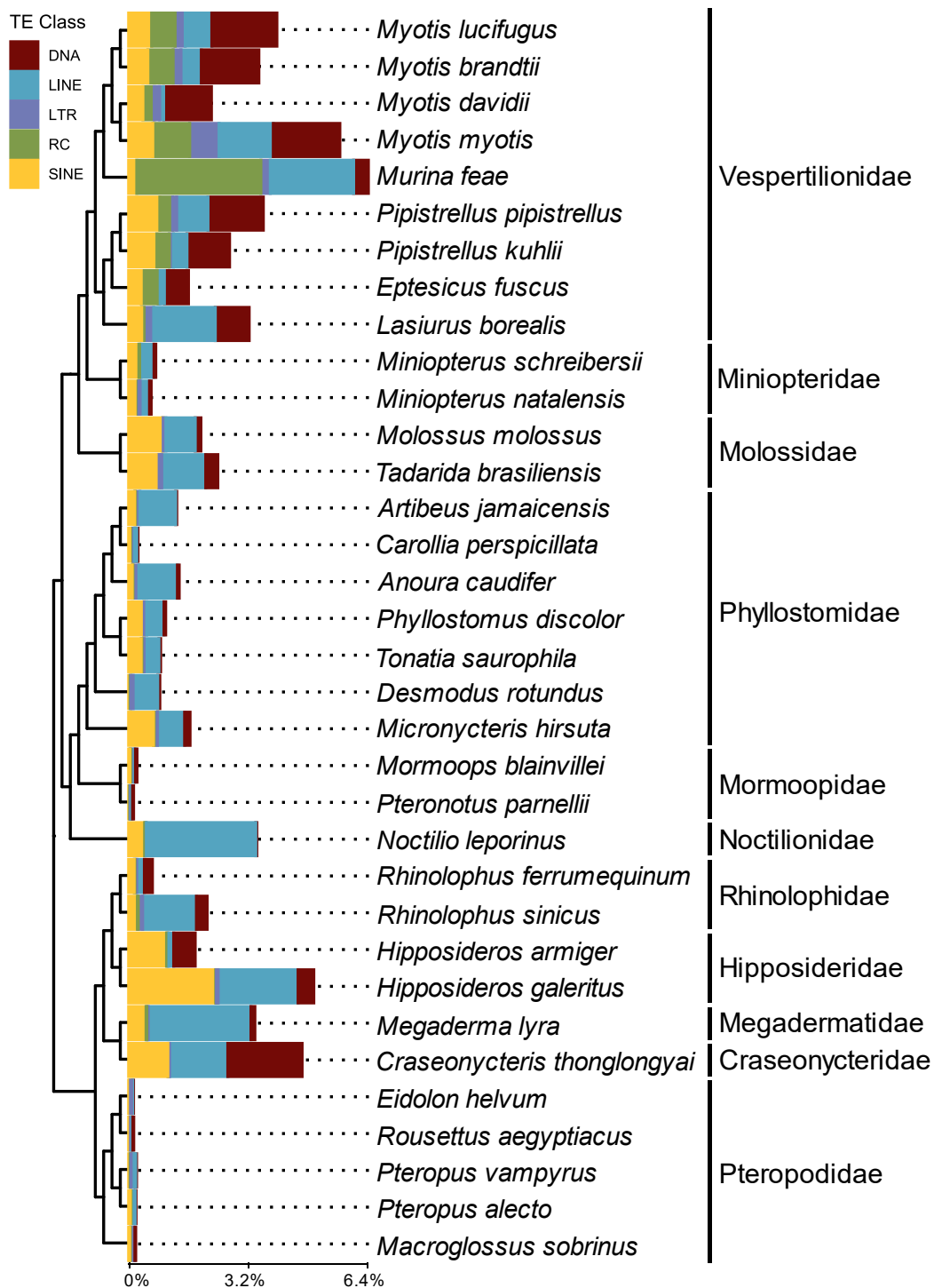
747 Macroscelidea, 7) Tubulidentata, 8) Afrosoricida, 9) Scandentia, 10) Dermoptera, 11)

748 Lagomorpha, 12) Eulipotyphla, 13) Perissodactyla, 14) Pholidota. The inner ring of stacked-bar
 749 data depicts the total percentage of the genome attributed to the five main categories of TEs:
 750 DNA transposons, LINES, SINEs, LTRs, & Helitrons. The outer ring of stacked-bar data shows
 751 the percentage of the genome derived from recently inserted TEs. Cladogram adapted from (65).
 752



753
 754 **Fig. 3. Redundancy analyses examining major axes of variation in TE accumulation and**
 755 **genome size related to orders (above) and families (below) of mammals.** Arrows represent
 756 significant correlations TE types with the first two RDA axes. Each axis reflects changes in TE
 757 composition related to ordinal (above) or familial (below) affiliation of taxa used in analyses.
 758 Gray circles represent orders or families that were not significantly correlated to at least one of
 759 the RDA axes whereas black circles represent orders or families with significant correlations.

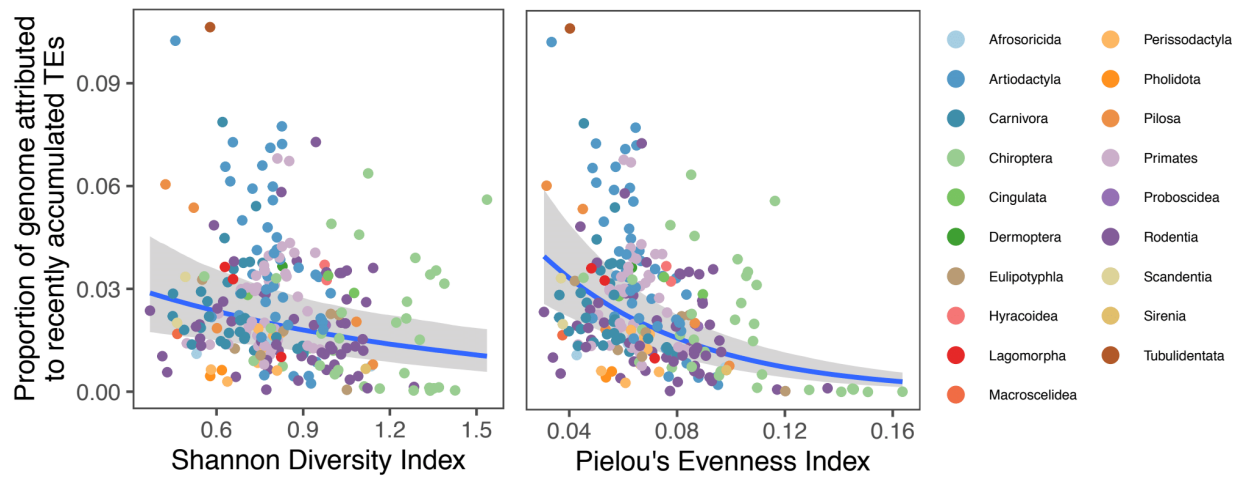
760



761

762 **Fig. 4. Stacked bar charts depicting proportions of recently accumulated TEs (<4% kimura**
 763 **from consensus TE) in bats.** Data is organized by TE classification and plotted onto the tips of
 764 the chiropteran portion of the mammalian tree, adapted from (65).

765



766

767 **Fig. 5. Recent mammalian TE diversity in relation to Shannon H (left) and Pielou's J**
768 **(right).** Blue line indicates the line of best-fit and shaded area is the 95% high probability
769 density of the fit. The r^2 for H was estimated at 0.67 (95% high probability density 0.52, 0.78),
770 and for J it was 0.69 (95% high probability density 0.56, 0.79).

771

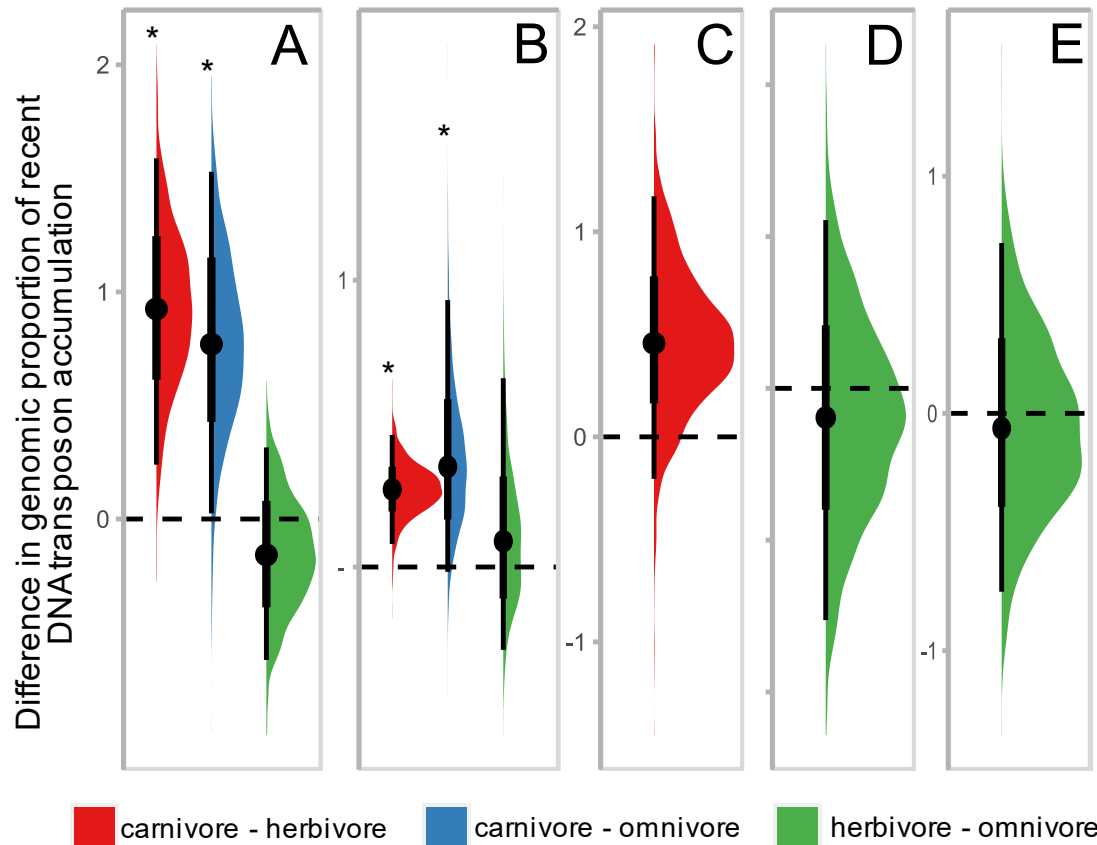


Fig 6. Half eye plots depicting fold differences in recent DNA transposon accumulation among three dietary phenotypes: carnivore, herbivore, and omnivore. Instead of showing the estimated values for each of the diets, these plots depict the fold ratio between each diet pair, so that the plot itself shows statistical significance. Comparisons for which the thin line does not overlap with 1 are significant, indicated by *. Plots correspond to the following taxonomic groups: A) placental mammals (r^2 estimated at 0.92 (95% high probability density 0.79, 0.97)), B) Artiodactyla (r^2 estimated at 0.64 (95% high probability density 0.32, 0.78)), C) Chiroptera (r^2 estimated at 0.34 (95% high probability density 0.02, 0.86)), D) Primates (r^2 estimated at 0.18 (95% high probability density 0.00, 0.58)), E) Rodentia (r^2 estimated at 0.07 (95% high probability density 0.00, 0.28)).