1    *Submission intended as a Research Article*

2

# Gigantic Genomes Can Provide Empirical Tests of TE Dynamics Models — An Example from Amphibians

5

6    Jie Wang[1,*], Michael W. Itgen[2], Huiju Wang[3], Yuzhou Gong[1], Jianping Jiang[1], Jiatang Li[1], Cheng

7    Sun[4], Stanley K. Sessions[5], Rachel Lockridge Mueller[2,*]

8

9    [1]Key Laboratory of Mountain Ecological Restoration and Bioresource Utilization, Chengdu

10    Institute of Biology, Chinese Academy of Sciences, China

11    [2]Department of Biology, Colorado State University, USA

12    [3]School of Information and Safety Engineering, Zhongnan University of Economics and Law, China

13    [4]Institute of Apicultural Research, Chinese Academy of Agricultural Sciences, Beijing, China

14    [5]Biology Department, Hartwick College, USA

15

16    *Corresponding authors: wangjie@cib.ac.cn, rlm@colostate.edu

17

18    **Running title:** *Jie W et al / TE biology in a gigantic amphibian genome*

19

20    **Total word count: 8515**

21    **Total references: 100**

22    **Total figures: 5**

23    **Total tables: 4**

24    **Total supplementary figures: 0**

25    **Total supplementary tables: 0**

26    **Title character count: 81**

27    **Running title character count: 35**

28    **Keyword count: 5**

29    **Abstract word count: 236**

30

1

## Abstract

Transposable elements (TEs) are a major determinant of eukaryotic genome size. The collective properties of a genomic TE community reveal the history of TE/host evolutionary dynamics and impact present-day host structure and function, from genome to organism levels. In rare cases, TE community/genome size has greatly expanded in animals, associated with increased cell size and altered anatomy and physiology. We characterize the TE landscape of the genome and transcriptome in an amphibian with a giant genome — the caecilian *Ichthyophis bannanicus*, which we show has a genome size of 12.2 Gb. Amphibians are an important model system because the clade includes independent cases of genomic gigantism. The *I. bannanicus* genome differs compositionally from other giant amphibian genomes, but shares a low rate of ectopic-recombination-mediated deletion. We examine TE activity using expression and divergence plots; TEs account for 15% of somatic transcription, and most superfamilies appear active. We quantify TE diversity in the caecilian, as well as other vertebrates with a range of genome sizes, using diversity indices commonly applied in community ecology. We synthesize previous models integrating TE abundance, diversity, and activity, and we test whether the caecilian meets model predictions for genomes with high TE abundance. We propose thorough, consistent characterization of TEs to strengthen future comparative analyses. Such analyses will ultimately be required to reveal whether the divergent TE assemblages found across convergent gigantic genomes reflect fundamental shared features of TE/host genome evolutionary dynamics.

**Key words**: TE expression; TE diversity index; genome size evolution; caecilian; transposon ecology

59

## Introduction

61 Transposable elements (TEs) are segments of DNA that move within genomes [1].

62 Because their movement is often associated with an increase in copy number, these

63 elements constitute a substantial — but variable — fraction of eukaryotic genomes, e.g.

64 2.7% in pufferfish (*Takifugu rubripes*) [2] and 85% in maize (*Zea mays*) [3]. TEs were

65 discovered by Barbara McClintock in the late 1940s, demonstrating that genomes are

66 far more dynamic entities than previously thought [4].

67     Although they share the characteristic of intra-genomic mobility, TEs are highly

68 diverse sequences. TE classification has been updated over the years to reflect new

69 discoveries [5]. Several classification systems have been proposed that establish nested

70 groups according to transposition mechanism, structure, sequence similarity, and, to

71 some degree, shared evolutionary history [6-10]. These classification systems, in turn,

72 have allowed the community of genome biologists to annotate TEs in the genomes of

73 diverse species, identifying differences in overall TE composition, TE activity, TE

74 turnover dynamics, and TE domestication across the tree of life [11-13].

75     Overall TE content is the main predictor of haploid genome size, which shapes a

76 variety of traits including the sizes of nuclei and cells, the rates of development and

77 basal metabolism, and the structural complexity of organs [14-21]. The evolutionary

78 forces that shape TE load include: mutation (specifically the insertion of new TE

79 sequences by transposition and their removal by deletion), selection (which likely

80 targets individual TE loci as well as the pathways that control TE transposition and

81 deletion) [22], and genetic drift (which determines how efficiently purifying selection

82 purges deleterious TE sequences) [23]. How these forces interact to generate genome

83 size diversity across the tree of life remains incompletely understood. Groups of related

84 species that vary in TE load and genome size provide critical model systems for

85 studying this fundamental question [24].

86      Across animals, genomic gigantism is rare. Within vertebrates, it is best understood

87    in a group related to the caecilians — the salamanders (order Caudata), a clade of ~700

88    extant species of amphibians (almost exclusively diploid) with haploid genome sizes

89    that range from 14 Gb to 120 Gb [25]. Fossil cell size data demonstrate that salamander

90    genome sizes have been large for ~160 million years [26]. Comparative genomic

91    analyses demonstrate that salamander genomes have high levels of TEs — particularly

92    LTR retrotransposons — and that these high levels reflect low rates of DNA loss in

93    dead-on-arrival non-LTR retrotransposons, low rates of ectopically-mediated LTR

94    retrotransposon deletion, and intact piRNA-mediated TE silencing machinery, albeit

95    with fewer TE-targeting piRNAs than are found in animals with smaller genomes [27-

96    33]. Phylogenetic comparative analyses demonstrate that salamanders' enormous

97    genomes result from an abrupt change in evolutionary dynamics at the base of the clade,

98    implying a discrete shift in the balance among the evolutionary forces shaping TE

99    accumulation [34, 35].

100    There are three living clades of amphibians: caecilians (order Gymnophiona),

101    salamanders (order Caudata), and frogs and toads (order Anura); Caudata and Anura

102    are sister taxa, and Gymnophiona is the sister taxon to Caudata + Anura. Frogs and

103    toads are a well-studied group of 7175 species. Of the 278 species (in 78 genera) for

104    which genome size estimates exist, a handful of species in three different genera have

105    genomes that reach or exceed 10 Gb [25], providing independent examples of genomic

106    expansion. Genomic data examined to date show diverse TE landscapes across species

107    [36-41], but no sequence data exist (to our knowledge) for those with the largest

108    genomes. Caecilians are a relatively understudied group of 214 extant species, all of

109    which are limbless, serpentine, burrowing or aquatic animals with reduced eyes, ringed

110    bodies, and strong, heavily ossified skulls. Genome size estimates exist for roughly 20

111    species and range from 2.8 Gb to 13.7 Gb [25, 35]. These data show yet another

112    independent example of genomic expansion within amphibians, suggesting a clade-

113    wide propensity towards TE accumulation relative to other vertebrates that is most

114    extreme in salamanders. Genomic data for caecilians are sparse, but growing based on

115    successes of the G10K consortium and others [42, 43]. Published data are lacking for

116    species with the largest genomes. This lack of amphibian data underlies a major gap in

117    our knowledge of vertebrate genome evolution [13]. More generally, a lack of detailed

118    information on TE biology in large and repetitive genomes, reflecting persistent

119    assembly and annotation challenges, underlies a major gap in genome biology as a

120    whole.

121         In this study, we present an analysis of TE biology in a caecilian with a large

122    genome — *Ichthyophis bannanicus* (Gymnophiona: Ichthyophiidae), which we show

123    has a genome size of 12.2 Gb. We compare the caecilian to other vertebrates with

124    diverse genome sizes, demonstrating how the TE community in a large genome can be

125    used to evaluate existing models of TE dynamics. *I. bannanicus* is a relatively small

126    species (adult sizes 30-41 cm) with an aquatic larval stage and terrestrial/fossorial adult

127    stage. Its distribution includes China and northern Vietnam, and it is an IUCN species

128    of Least Concern. We analyzed both genomic shotgun sequence data and RNA-Seq

129    data from diverse tissues to answer the following specific questions: (1) what

130    abundance and diversity of TEs make up the bulk of the large *I. bannanicus* genome?

131    (2) what are the amplification and deletion dynamics of TEs in the genome? (3) what

132    contribution does the large genomic TE load make to the somatic transcriptome? (4) do

133    the patterns of genomic TE composition and overall TE expression fit the predictions

134    of models of TE dynamics in large genomes? We show that up to 68% of the *I.*

135    *bannanicus* genome is composed of TEs, with another 9% identified as repetitive

136    sequences not classifiable as known TEs. The two most abundant TE superfamilies —

137    DIRS/DIRS and LINE/Jockey — account for ~50% of the genome. Unlike salamander

138    genomes, the *I. bannanicus* genome has relatively few LTR retrotransposons,

139    demonstrating that repeated instances of extreme TE accumulation in amphibians do

140    not reflect failure to control a specific type of TE. We show that rates of ectopic-

141    recombination-mediated deletion are low relative to vertebrates with more typical

142  genome sizes, and that TE expression is high. We quantify and compare TE diversity

143  in *I. bannanicus* and 11 other vertebrates using indices common in community ecology.

144  We demonstrate that comparative analyses of TE diversity can be a powerful tool for

145  evaluating models of TE dynamics, and we show that it could be even more powerful

146  if researchers adopt a uniform approach to TE diversity analysis. We propose such an

147  approach to move the field forward. Taken together, our results demonstrate that

148  computationally feasible analyses of large genomes can reveal the genomic

149  characteristics favoring expanded TE communities, as well as the resulting impact of

150  high TE load on the transcriptome. Such analyses targeting phylogenetically diverse

151  organisms can yield fundamental insights into the complex ways in which TEs drive

152  genome biology.

153

154  **Results**

155  **The *I. bannanicus* genome is 12.2 Gb and contains most known TE superfamilies**

156  The haploid genome size of *I. bannanicus* was estimated to be 12.2 Gb based on

157  analyses of Feulgen-stained erythrocytes following established methods [14]. This

158  estimate is similar to the other published estimate from the same genus (*I. glutinosus*,

159  11.5 Gb) [35]. We used the PiRATE pipeline [44], designed to mine and classify repeats

160  from low-coverage genomic shotgun data in taxa — such as caecilians — that lack

161  genomic resources. The pipeline yielded 59,825 contigs (**Table 1**). RepeatMasker

162  mined the majority of the repeats (37,123 out of 59,825, or 62.1%). dnaPipeTE was the

163  second most effective tool, mining 19,160 repeats (32.0%), followed by RepeatScout

164  (3.0%) and TE-HMMER (2.7%). In this pipeline, TE-denovo, LTR-harvest, Helsearch,

165  SINE-finder, and MITE hunter found few additional repeats, and we found no

166  additional repeats using MGEScan. Clustering with CD-HIT-est with a 95% sequence

167  identity cutoff yielded 51,862 contigs, and clustering at 80% yielded 23,092 contigs.

168      Repeat contigs were annotated as transposable elements to the levels of order and

169  superfamily in Wicker's hierarchical classification system [7], modified to include

170    several recently discovered TE superfamilies, using PASTEC [45]. Of the 59,825

171    identified repeat contigs, 50,471 (84.4%) were classified as known TEs (**Table 2**).

172    Transposable elements representing eight of the nine orders proposed in Wicker's

173    system are present in the *I. bannanicus* genome; only Crypton was not identified by our

174    pipeline (although we note that 192 chimeric contigs were filtered out that included a

175    Crypton annotation, and nine transcriptome contigs were annotated as Crypton). Within

176    these eight orders, our analyses identified 25 TE superfamilies, each represented by

177    between two and 26,507 annotated contigs. Non-autonomous TRIM, LARD, and MITE

178    elements are also present in the *I. bannanicus* genome, represented by 229, 28, and 146

179    contigs, respectively, and an additional 277 contigs were only annotated to the level of

180    order or class (i.e. unknown LINE, SINE, and TIR or unknown Class I) (Table 2).

181

182    **78% of the *I. bannanicus* genome is repetitive, dominated by DIRS elements**

183    To calculate the percentage of the caecilian genome composed of different TEs, the

184    shotgun reads were masked with RepeatMasker v-4.0.7 using our caecilian-derived

185    repeat library. We then repeated the RepeatMasker analysis excluding the unknown

186    repeats and compared the two sets of results as a rough approximation of the number

187    of unknown repeat contigs that were, in fact, TE-derived sequences that were divergent,

188    fragmented, or otherwise unidentifiable by our pipeline. 68.2% of these sequences

189    (measured as bp) were masked as repetitive when the repeat library included only the

190    50,471 contigs classified as TEs and the 29 contigs annotated as putative multi-copy

191    host genes; 66.1% were identifiable to the superfamily level of TEs (Table 2), an

192    additional 1.94% were identifiable only to the class or order level, and 0.17% were

193    multi-copy host genes. When the analysis was performed including the 9325 unknown-

194    repeat contigs, along with the classified TEs and putative multi-copy host genes, 77.6%

195    of the data were masked as repetitive overall, suggesting that the unknown repeats

196    comprise 9.4% of the genome. However, the percentage of the genome identified as

197    known TEs decreased from 68.0% to 54.7% with the inclusion of unknown repeats,

198    demonstrating that many reads were sufficiently similar to known TEs to be masked by

199    them when unknown repeat contigs were not available as a best-match option. This

200    result suggests that at least some of the unknown repeats are TE-derived sequences.

201        Class I TEs (retrotransposons) make up 52.09-63.68% (unknown repeats included

202    or excluded in the repeat library, respectively) of the *I. bannanicus* genome; they are

203    almost 20 times more abundant than Class II TEs (DNA transposons; 2.63–4.36%).

204    DIRS/DIRS is the most abundant superfamily (25.88-30.20% of the genome), followed

205    by LINE/Jockey (16.92-20.59%), LINE/L1 (3.05-3.23%), LTR/ERV (1.62-1.82%),

206    LINE/RTE (1.50-1.60%), and LTR/Gypsy (1.10-1.35%); all are retrotransposons

207    (Table 2). TIR/hAT (0.57-1.15%), TIR/CACTA (0.52-0.56%), and TIR/Tc1-Mariner

208    (0.49-0.59%) are the most abundant superfamilies of DNA transposons (Table 2).

209    These proportions differ from those found in the gigantic genomes of salamanders,

210    where LTR/Gypsy elements dominate (7 – 20% of the genome, depending on species),

211    DIRS/DIRS elements never exceed 7% of the genome, and LINE/Jockey elements

212    never exceed 0.03% of the genome [30, 32]. Here and throughout the paper, we are

213    interpreting our results based on the assumption that the genomic shotgun data are a

214    random representation of the whole genome; Illumina reads should sample the genome

215    in a random and independent manner, despite some stochastic sampling error.

216

217    **The *I. bannanicus* genome shows low diversity index values when measured at the**

218    **TE superfamily level**

219    Diversity indices are mathematical measures of diversity within a community. In

220    ecology, they are widely used to summarize species diversity within an ecological

221    community, although they are also used in other fields (e.g. economics). Diversity

222    indices take into account species richness (the total number of species present) and

223    evenness (based on the proportional abundance of each species) [46]. Within genome

224    biology, richness can summarize the total number of TE types (e.g. TE superfamilies)

225    and evenness can summarize the proportion of the genome occupied by each TE type

8

226     [47-50]. We calculated two commonly used diversity indices — the Shannon Index

227     and the Gini-Simpson Index [51, 52] — on the caecilian TE community, as well as

228     the TE communities from 11 other vertebrates spanning a range of genome sizes and

229     types of datasets. Genome sizes ranged from 0.4 Gb (the pufferfish *Takifugu rubripes*)

230     to 55 Gb (the hellbender salamander *Cryptobranchus alleganiensis*). Datasets ranged

231     from full genome assemblies to low-coverage genome skims. The Shannon Index

232     quantifies the uncertainty in identity of an individual drawn at random from a

233     community. The Gini-Simpson index quantifies the probability that two individuals

234     drawn at random from a community are different types, and it gives more weight to

235     dominant (i.e. most abundant) species. Results are summarized in **Table 3**. The

236     Shannon Index ranges from 0.9 (chicken, the least diverse) to 2.41 (green anole lizard,

237     the most diverse). The Gini-Simpson Index ranges from 0.5 (chicken, the least

238     diverse) to 1 (pufferfish, the most diverse). By both indices, the caecilian has the

239     second-least diverse genome of the 11 total genomes compared. There is no overall

240     correlation between genome size and TE diversity using either index.

241

242     **Most TE superfamilies are active in the *I. bannanicus* genome**

243     For each of the 19 TE superfamilies accounting for $\geq 0.005\%$ of the genome, the overall

244     amplification history was summarized by plotting the genetic distances between

245     individual reads (representing TE loci) and the corresponding ancestral TE sequences

246     as a histogram with bins of 1%. Of these 19 TE superfamilies, 17 of the resulting

247     distributions showed characteristics of ongoing or recent activity (i.e. presence of TE

248     sequences <1% diverged from the ancestral sequence and a unimodal right-skewed, J-

249     shaped, or monotonically decreasing distribution) (**Figure 1**). Six of these showed

250     essentially unimodal, right-skewed distributions: LTR/ERV, DIRS/DIRS,

251     LINE/Jockey, LINE/RTE, TIR/PiggyBac, and TIR/Sola. An additional three showed

252     essentially unimodal, right-skewed distributions with a spike in sequences <1%

253     diverged from the ancestral sequence: SINE/7SL, TIR/hAT, and TIR/Tc1/mariner. A

9

254   single superfamily — PLE/Penelope — showed a left-skewed J-shaped distribution.

255   These ten distributions suggest TE superfamilies that continue to be active today, but

256   whose accumulation peaked at some point in the past. In contrast, six TE superfamilies

257   showed essentially monotonically decreasing distributions with a maximum at <1%

258   diverged from the ancestral sequence: LTR/Gypsy, DIRS/Ngaro, LINE/L1,

259   TIR/CACTA, TIR/PIF-Harbinger, and Maverick. SINE/5S has a bimodal distribution

260   with a maximum at <1% diverged from the ancestral sequence. These seven

261   distributions suggest TE superfamilies that continue to be active today at their highest-

262   ever rates of accumulation. Two superfamilies — LTR/Retrovirus and LINE/R2 —

263   appear largely inactive, showing unimodal distributions with few sequences <1%

264   diverged from the ancestral. For almost all superfamilies, multiple contigs that were

265   <80% identical in sequence to one another were assembled (range $1 - 8513$), suggesting

266   the presence of many families within each superfamily.

267

268   **Ectopic-recombination-mediated deletion levels are lower in *I. bannanicus* than in**

269   **vertebrates with smaller genomes**

270   Ectopic recombination, also known as non-allelic homologous recombination, occurs

271   between two DNA regions that are similar in sequence, but do not occupy the same

272   locus. Ectopic recombination among LTR retrotransposon sequences can produce

273   deletions that leave behind solo LTRs, which are single terminal repeat sequences that

274   lack the corresponding internal sequence and matching terminal repeat sequence.

275   Accordingly, the ratio of LTR sequences to internal retrotransposon sequences can be

276   used to estimate levels of ectopic-recombination-mediated deletion. Larger genomes

277   like *I. bannanicus* are predicted to have lower levels of deletion [33].

278       Two superfamilies were selected for ectopic recombination analysis: DIRS/DIRS,

279   which accounts for over a quarter of the caecilian genome, and LTR/Gypsy, which is

280   one of the two most abundant LTR superfamilies in the caecilian genome at 1.4%, but

281   which dominates other gigantic amphibian genomes [53]. Mean estimates of the total

10

282    terminal sequence to internal sequence ratio (TT:I) across the 9 DIRS/DIRS contigs

283    range from 1.2:1 to 0.7:1, depending on the minimum alignment length for reads

284    (**Figure 2**). Mean TT:I estimates across the 17 LTR/Gypsy contigs range from 1.3:1 to

285    1.2:1. Values of 1:1 are expected in the absence of ectopic-recombination-mediated

286    deletion. The higher sensitivity of DIRS/DIRS than LTR/Gypsy to the minimum

287    alignment length parameter value likely reflects the shorter length of the terminal

288    sequence (150 bp vs. 744 bp in DIRS/DIRS and LTR/Gypsy, respectively); the 0.7:1

289    TT:I value for DIRS/DIRS is likely an underestimate. Variation in the TT:I ratio among

290    contigs in each superfamily was similar (Figure 2) and lower than the ranges reported

291    in vertebrates with more typically sized (i.e. smaller) genomes [33].

292        For both superfamilies, ectopic-recombination-mediated deletion levels in the

293    caecilian (TT:I ratio ~1.2:1) are similar to the low levels estimated from four gigantic

294    salamander genomes (TT:I ratios 0.55:1 to 1.25:1 for *Aneides flavipunctatus,*

295    *Batrachoseps nigriventris*, *Bolitoglossa occidentalis*, and *Bolitoglossa rostrata*) and

296    below the levels estimated from vertebrates with more typically sized (i.e. smaller)

297    genomes (TT:I ratios 1.7:1 to 3.35:1 for *Anolis carolinensis, Danio rerio, Gallus gallus,*

298    *Homo sapiens,* and *Xenopus tropicalis* [33]. TT:I ratios measured for LTR/Gypsy in

299    two salamander species (*A. flavipunctatus* and *B. nigriventris*) are 0.9:1 and 1.25:1,

300    encompassing the value for *I. bannanicus* LTR/Gypsy.

301        If deletion levels were the same between the two superfamilies in the *I. bannanicus*

302    genome, the DIRS/DIRS TT:I ratio would be expected to be lower than the LTR/Gypsy

303    TT:I ratio because of the structure of DIRS/DIRS; it has inverted terminal repeats and

304    internal complementary regions [54, 55] that are expected to produce incomplete

305    deletion of the internal sequence following ectopic recombination. The higher TT:I

306    ratio actually estimated in DIRS/DIRS may reflect the greater abundance of this

307    superfamily, which increases the number of potential off-targets for recombination,

308    offsetting both the incomplete deletion of the internal sequence as well as the shorter

309    terminal sequences in DIRS that would predict lower levels of deletion [56].

11

310

**Autonomous and non-autonomous TEs are transcribed in *I. bannanicus***

312 To annotate transcriptome contigs containing autonomous TEs (i.e. those with open

313 reading frames encoding the proteins necessary for transposition), BLASTx was used

314 against the Transposable Element Protein Database (RepeatPeps.lib,

315 http://www.repeatmasker.org/). To annotate contigs containing non-autonomous TEs

316 that lack identifiable open reading frames, RepeatMasker was used with our caecilian-

317 derived genomic repeat library of non-autonomous TEs. To identify contigs that

318 contained an endogenous caecilian gene, the Trinotate annotation suite was used [57].

319 38,584 contigs were annotated as endogenous (i.e. non-TE-derived) caecilian genes.

320 53,106 contigs were annotated as autonomous TEs using BLASTx against the

321 Transposable Element Protein Database (RepeatPeps.lib). An additional 2658 contigs

322 were annotated as non-autonomous TEs using the caecilian TRIM-, LARD-, SINE- and

323 MITE-annotated genomic contigs. 1445 contigs were annotated as both autonomous

324 TEs and endogenous caecilian genes, and an additional 342 were annotated as both non-

325 autonomous TEs and endogenous caecilian genes (**Table 4**).

326 Of the 20 most highly expressed putative "TE/gene" contigs, ten were confirmed

327 to have annotations for both a TE and a gene with non-overlapping ORFs. Of these, the

328 TE was upstream of the gene in eight cases and downstream in two cases. Six of the

329 upstream TEs were autonomous and thus contained ORFs; four of these were encoded

330 on the same strand as the gene (two in-frame, two not in-frame) and two were encoded

331 on the opposite strand. One of the two downstream TEs was autonomous, and it was

332 encoded on the opposite strand from the gene. Although requiring further validation,

333 these results suggest that at least some gene/TE pairs are co-transcribed, a way in which

334 TE insertions can regulate gene expression [58]. One contig had overlapping

335 annotations of a gene and a TE, a pattern that could reflect either convergence in

336 sequence or exaptation of a TE [59].

337

12

338 **TE expression correlates with genomic abundance in *I. bannanicus***

339 Among the transcriptome contigs with TPM (transcripts per million) ≥ 0.01,

340 autonomous TEs account for 18.4% of the total transcriptome contigs and 13.2% of the

341 overall somatic transcriptome (TPM = 131,793) (Table 4, **Figure 3**). Non-autonomous

342 TEs account for 0.9% of the total transcriptome contigs and 0.8% of the somatic

343 transcriptome (summed TPM = 8484). Contigs annotated both as TEs and endogenous

344 caecilian genes account for 0.6% of annotated transcriptome contigs and 0.6% of the

345 somatic transcriptome (summed TPM = 6443). Endogenous (non-TE-derived) caecilian

346 genes account for 13.3% of the total transcriptome contigs and 29.6% of the somatic

347 transcriptome (summed TPM = 295,759). Unannotated contigs account for 66.8% of

348 the total transcriptome contigs and 55.6% of the somatic transcriptome (summed TPM

349 of unannotated contigs = 555,776). Five superfamilies (*I, Zisupton, Kolobok, Academ,*

350 and *Crypton*) were detected at low expression levels in the transcriptome, but were not

351 initially detected in the genomic data; mapping the genomic reads to these

352 transcriptome contigs with Bowtie2 identified ≤ 3 reads per superfamily, indicating

353 their extremely low frequency in the genome. In contrast, only one superfamily (*7SL*)

354 was detected in the genomic data but not the transcriptome data.

355     Class I TEs (retrotransposons) are over ten times more abundant in the

356 transcriptome than Class II TEs (summed TPM = 130,076 and 10,202, respectively).

357 Within the retrotransposons, DIRS are the most highly expressed, followed by Jockey

358 and L1; these three superfamilies are also the most abundant in the genome. For almost

359 all retrotransposon superfamilies, hundreds to thousands of transcriptome contigs that

360 were <80% identical in sequence to one another were assembled (range 1 – 12,652),

361 suggesting the simultaneous activity of many families within all of the superfamilies in

362 the caecilian somatic transcriptome. Large differences (up to ~10,000-fold) in

363 expression were detected among the different contigs within superfamilies, suggesting

364 variable expression levels across loci and among families; we interpret this pattern with

365 caution because of the challenges of uniquely mapping short reads to contigs of similar

13

366   sequence. Within the DNA transposons, Tc1-Mariner, CACTA, and hAT are the most

367   highly expressed superfamilies, and MITEs (transposon derivatives) are expressed at

368   similar levels to these superfamilies, although they lack their own promoters. These

369   four types of sequences are also the four most abundant types of DNA transposons in

370   the genome, although their genomic abundance is not perfectly correlated with their

371   relative expression levels. For the DNA transposons, tens to hundreds of contigs that

372   were <80% identical in sequence to one another were assembled (range 2-421), and up

373   to ~1000-fold differences in expression were detected among contigs. Overall, a strong

374   correlation was detected between genomic abundance of a TE superfamily and its

375   overall somatic expression level ($\rho = 0.879$, $p < 0.001$) (**Figure 4**). Although germline

376   expression data is required to analyze the relationship between TE transcription and

377   TE-activity-driven genome evolution, the somatic data nevertheless provides valuable

378   information on the cellular resources allocated to transcription of a greatly expanded

379   TE community.

380

381   **Discussion**

382   **Repeat element landscape characterization in large genomes**

383   Large, repetitive genomes have proven difficult to assemble and annotate with the

384   computational power and analytical tools successfully applied to archaeal, bacterial,

385   and smaller eukaryotic genomes [60, 61]. Recent successful genome sequencing efforts

386   aimed at the 32 Gb genome of *Ambystoma mexicanum*, a laboratory model salamander

387   species, leveraged multiple types of data (i.e. optical mapping, short- and long-read

388   genomic sequence data, transcriptomic data, linkage mapping, fluorescence in situ

389   hybridization) and a new assembler designed to minimize compute time and storage

390   requirements [30, 62]. These projects yielded fundamental insights into the structure

391   and evolution of vertebrate chromosomes. They also advanced understanding of the

392   transposons that make up large genomes, adding to research on the 20-Gb Norway

393   spruce and the 22-Gb loblolly and 31-Gb sugar pine [63-65]. This depth of analysis,

394 however, remains infeasible for non-model organisms with large genomes, whose study

395 is nevertheless required for a comprehensive picture of the complex ways in which TEs

396 drive genome biology. Our work affirms the power of low-coverage sequence data to

397 reveal the overall repeat element landscape of large genomes, an approach applied most

398 often in plants (which include the majority of huge genomes) [66, 67]. We argue that

399 this overall landscape, although it lacks the positional information about individual TE

400 insertions that genome assemblies provide, nevertheless contains much information that

401 can be leveraged to identify the evolutionary processes that drive assembly and stability

402 of TE communities.

403     Repeat element landscapes are informative because they include data on the

404 abundance, diversity, and activity of TEs that make up the overall TE community in a

405 genome. Models of TE dynamics — both formal and informal — predict different

406 values for TE abundance, diversity, and activity depending on levels of purifying

407 selection, silencing, and deletion of TEs. Despite much progress, these forces remain

408 challenging to measure directly. Thus, empirical estimates of TE landscapes provide a

409 feasible alternative by which to validate these models and advance our understanding

410 of TE dynamics in natural systems.

411

412 **Repeat element landscapes from large genomes provide tests of models of TE**

413 **dynamics**

414 Large genomes are especially powerful data points because they represent extreme

415 values of TE abundance, and models of TE dynamics make specific predictions about

416 the effects of TE abundance on TE diversity and activity. We first summarize and

417 highlight the differences among several of these models here (**Figure 5**):

418 (1) Petrov 2003 — TE deletion is caused by ectopic recombination between similar TE

419 sequences. Rates of ectopic recombination/deletion are typically higher in smaller

420 genomes and lower in larger genomes. Thus, smaller genomes are predicted to select

421 for more diverse TE communities, and larger genomes should allow less diverse TE

15

422     communities [56, 68]. This model predicts an inverse relationship between genome size

423     and TE diversity.

424     (2) Furano 2004 — Because ectopic recombination can cause harmful deletions, it is

425     one of the primary reasons for TEs' deleterious effects on host fitness. Thus, genomes

426     with lower ectopic recombination/deletion rates are more permissive to TE activity,

427     allowing the accumulation of more TEs (increased genome size) as well as increased

428     TE activity and out-competition of many TE lineages by the lineage that most

429     successfully exploits host replication factors [69]. Like Petrov 2003, this model predicts

430     an inverse relationship between genome size and TE diversity, but for different reasons.

431     (3) Boissinot 2016 — Genomes with lower ectopic recombination/deletion rates have

432     higher levels of insertion of active TE copies into the genome. In addition to yielding a

433     larger genome, this higher number of active TE copies triggers an arms race to control

434     transposition, and the arms race leads to a decrease in diversity (i.e. only one family

435     active at a time) [70]. Like Petrov 2003 and Furano 2004, this model also predicts an

436     inverse relationship between genome size and TE diversity, but for still different

437     reasons.

438     (4) Abrusan 2006 — TE diversity and activity levels were modeled with a system of

439     differential equations that includes parameters for the number of TE strains, the number

440     of active TE insertions, TE replication rates, the strength of specific silencing of TEs

441     (representing small-RNA-mediated silencing), cross-reactivity of silencing, and TE

442     inactivation by mutation or selection [71]. Although their model did not specifically

443     address genome size, it did predict that increased genome size would be associated with

444     decreased TE diversity if a) larger genomes harbor more active TE copies, and b) cross-

445     reactive silencing exists among TEs. Under these conditions, competition among the

446     TEs to evade cross-reactive silencing would lead to decreased TE diversity. Cross-

447     reactive silencing in this model is not sequence-specific; this is relevant because

448     silencing of TEs by small-RNA-mediated silencing (e.g. the piRNA pathway) is

449     sequence-specific, but can have some off-target effects. These off-target effects are

16

450   predicted to have the opposite effect on TE diversity than non-sequence-specific cross-

451   reactive silencing; they should select for higher TE diversity. Overall, the predictions

452   for genome size and TE diversity from this model are complex, depending on the

453   relative strengths of specific TE silencing, off-target specific TE silencing, and cross-

454   reactive (i.e. sequence-independent) silencing.

455   (5) Elliot 2015 — Based on empirical comparisons across genomes of different sizes,

456   TE diversity was proposed to increase with TE abundance until genomes reach

457   moderate size, but extremely large genome sizes were proposed to reflect the

458   proliferation of only a subset of TE diversity by unspecified mechanisms [72]. This

459   predicts an inverse relationship between genome size and TE diversity at the largest

460   genome sizes.

461   (6) Kijima 2013 — TE evolution was modeled using a population genetic simulation

462   framework that includes parameters for transposition, TE deletion, purifying selection

463   on TE copy number (genome size), and degeneration into inactive copies [73]. When

464   copy number selection is strong (i.e. genome size remains small), the total number of

465   TEs is lower, but the proportion of active copies of TEs is higher. When copy number

466   selection is weak (i.e. genome size is allowed to increase), the total number of TEs is

467   higher, but the proportion of active copies of TEs is lower. This reflects competition

468   among TEs to occupy limited available spaces in the genome. This model does not

469   consider TE diversity — it models only a single TE strain — but it predicts an inverse

470   relationship between genome size and proportion of the total TE community that is

471   actively transposing. Interestingly, they find that excision (deletion) rate is not a

472   predictor of copy number.

473   (7) Roessler 2018 — TE evolution was modeled using ordinary differential equations

474   including parameters for TE transposition, RNA-mediated TE silencing, TE deletion,

475   and TE copy number (genome size) [74]. This model predicts that, under low rates of

476   TE deletion, TE copy number and genome size increase, and the proportion of active

477   TEs goes down because the host organism can use the accumulating TE sequences as

17

478 templates for producing more small silencing RNAs and, thus, inactivate a higher

479 proportion of TE sequences. Like Kijima 2013, this model predicts an inverse

480 relationship between genome size and proportion of the total TE community that is

481 actively transposing, but for different (albeit complementary) reasons.

482  Does the TE landscape of the large caecilian genome — with its high levels of TE

483 abundance and low levels of TE ectopic recombination/deletion — fit the predictions

484 of these models or allow discrimination among them? Most share a prediction of

485 decreased TE diversity in large genomes. Measured at the coarse-grained level of

486 number of superfamilies present (i.e. taking into account richness only) [72], *I.*

487 *bannanicus* does not fit this prediction; at least 25 TE superfamilies are present in the

488 genome (as detected by our genomic and transcriptomic analysis). However, genome

489 expansion in *I. bannanicus* is correlated with high DIRS/DIRS and LINE/Jockey

490 superfamily abundance, consistent with Elliot 2015's prediction that gigantic genomes

491 would reflect proliferation of a limited subset of all TEs. This expansion decreases

492 evenness, despite the maintenance of high richness; this is exactly the type of change

493 in overall diversity that is captured by the indices we advocate here.

494  Comparing the diversity indices calculated for *I. bannanicus* with the 11 other

495 vertebrate genomes (Table 3) allows a direct test of the relationship between genome

496 size (i.e. TE abundance) and TE diversity. Because the genomes included were

497 analyzed with different sequencing depths, we favor the Gini-Simpson Index as it is

498 less affected by rare species (TE superfamilies), which are more likely missed in the

499 low-coverage datasets (e.g. *I, Zisupton, Kolobok, Academ,* and *Crypton*; Table 2).

500 Consistent with model predictions, the smallest genome (*T. rubripes*) has the highest

501 TE diversity, and the three most diverse genomes (*T. rubripes, A. carolinensis,* and *X.*

502 *tropicalis*) are three of the four smallest (Table 3). However, among the large amphibian

503 genomes — *I. bannanicus* and the five salamanders — there is no relationship between

504 TE abundance and diversity. Furthermore, the chicken genome is the least diverse, and

505 it is the second-smallest.

506    However, the lack of relationship between TE abundance and diversity, measured

507    here at the TE superfamily level for 11 species, does not necessarily refute the models

508    of TE dynamics that predict decreased TE diversity with increased TE abundance.

509    Diversity exists within TE superfamilies as well; TE families are typically operationally

510    defined based on Wicker's 80/80/80 rule, and subfamilies can be further split based on

511    pairs of substitutions overrepresented in TE alignments that are unlikely to have arisen

512    independently by chance [7, 75]. It is not yet clear what levels of sequence divergence

513    translate into functionally relevant "TE diversity" in the models summarized above.

514    More specifically, TE diversity implies: 1) TE sequences that have diverged beyond the

515    ability to ectopically recombine in Petrov 2003, 2) TE sequences that have diverged

516    enough to differ in ability to monopolize host replicative resources in Furano 2004, 3)

517    TE sequences that have diverged enough to (sequentially) out-evolve host silencing

518    machinery in Boissinot 2016, and 4) TE sequences that have diverged enough to differ

519    in their silencing by cross-reactive (i.e. non-sequence-specific) or off-target (i.e.

520    sequence-specific, but tolerant of mismatches) TE silencing mechanisms in Abrusan

521    2006. We still lack sufficient information about TE silencing to define the levels of

522    sequence divergence likely to accompany these changes in TE dynamics. Thus, it is not

523    yet clear whether diversity indices are best focused at the TE superfamily, family, or

524    subfamily levels. As an example, the chicken genome is the least diverse measured here

525    at the level of TE superfamilies because CR1 elements dominate the genome; however,

526    diversity exists within the CR1 elements that may be functionally relevant [76]. To

527    move the field forward, we advocate using Shannon and Simpson indices at the levels

528    of TE family and subfamily (in addition to superfamily) when datasets allow. When

529    this is impossible — for example, when working with low-coverage shotgun data from

530    gigantic genomes like *I. bannanicus* — we advocate calculating diversity indices at the

531    superfamily level, but also reporting the numbers of genomic and transcriptomic

532    contigs < 80% identical as a tractable within-superfamily approximation of TE diversity

533    (Table 2). This measure is analogous to species richness and lacks information on

19

534     evenness (because of the challenges of uniquely mapping short reads to contigs of

535     similar sequence), so it is less informative than diversity indices. However, the

536     reporting of this measure by researchers studying diverse organisms would allow

537     progress towards rigorously testing the relationship between genome size and TE

538     diversity. Furthermore, it may identify specific taxa as appropriate models to examine

539     evolutionary changes in TE silencing pathways. For example, *I. bannanicus* has a large

540     genome but appears to maintain a high number of TE families (Table 2), suggesting

541     that its TE silencing machinery includes high levels of off-target silencing [71].

542       In addition to predicting low TE diversity, models of TE dynamics predict a

543     decreased proportion of active TEs as TE abundance and genome size increase. Of the

544     19 caecilian TE superfamilies for which amplification histories were examined, 17

545     appear to have ongoing activity (Figure 1). These results are largely corroborated by

546     the (albeit somatic) expression data, although SINE/7SL and LINE/R2 show conflicting

547     patterns in the genomic and transcriptomic data (Figure 1, Table 2). TE expression is

548     necessary, but insufficient, for TE activity, but it is a tractable proxy for TE activity.

549     Taken together, these datasets suggest near-complete activity at the TE superfamily

550     level in the *I. bannanicus* genome. At the levels of TE family or individual insertions,

551     activity is difficult to assess with our data; however, the presence of multiple

552     transcriptome contigs <80% identical within superfamilies minimally suggests the

553     expression of multiple families. Our recommendation that researchers report the

554     number of transcriptomic TE contigs < 80% identical will also allow progress towards

555     rigorously testing the relationship between genome size and TE activity, as will

556     adoption of recent methods to measure locus-specific expression when datasets allow

557     [77].

558       Overall, ~15% of all somatic tissue transcripts of *I. bannanicus* are TEs (Table 4).

559     Comparing overall levels of TE expression across different genome sizes remains

560     difficult because TE expression in general is understudied [77], transcriptome size

561     differences that accompany genome size differences are typically not quantified [78],

562 and TE annotation and expression quantification methods vary across studies [38, 79-

563 81]. As another step towards testing the relationship between genome size and TE

564 activity, we advocate annotation of both autonomous and non-autonomous TE

565 transcripts and reporting of expression levels of TEs and endogenous genes (Table 2,

566 Table 4, Figure 3).

567      Taken together, our work lays a foundation for comparative genomic analyses that

568 link properties of TE communities — abundance, diversity, and activity — to genome

569 size evolution. Such analyses, in turn, will reveal whether the divergent TE assemblages

570 found across convergent examples of genomic gigantism reflect more fundamental

571 shared features of TE/host genome evolutionary dynamics.

572

573 **Materials and Methods**

574 **Specimen information**

575 We collected a single male adult caecilian (*I. bannanicus*) from the species' type

576 locality (E101.3887, N21.8724) in Mengxing County, Yunnan province, China. The

577 individual had a total body length of 16.0 cm and a body mass of 4.8 g. Following

578 dissection, the carcass was fixed in formalin and transferred to 70% ethanol.

579

580 **Genome size estimation**

581 Blood smears were prepared from the formalin-fixed *I. bannanicus* specimen as well

582 as a formalin-fixed salamander (*Plethodon cinereus*) with an appropriate genome size

583 to serve as the reference standard (22.14 Gb) [82]. Blood cells were pipetted onto glass

584 microscope slides and air-dried, then hydrated for three minutes in distilled water.

585 Slides were 1) hydrolyzed in 5N HCl for 20 minutes at 20°C and washed three times in

586 distilled water for one minute each, 2) stained with Schiff's reagent in a Coplin jar for

587 90 minutes at 20°C, 3) soaked in three changes of 0.5% sodium metabisulfite solution

588 for five minutes each and rinsed in three changes of distilled water for one minute each,

21

589    and 4) dehydrated in 70%, 95%, and 100% ethanol for one minute each, air-dried, and

590    mounted in immersion oil and cover glass.

591         The stained slides were photographed using an Olympus BX51 compound

592    microscope fitted with a Spot Insight 4 digital camera for image analysis. Stained nuclei

593    were photographed under 100x oil immersion and the integrated optical densities were

594    measured using ImagePro software. Genome size for *I. bannanicus* was calculated by

595    comparing the mean optical density to that of the reference standard, *P. cinereus.*

596

597    **Genomic shotgun library creation, sequencing, and assembly**

598    Total DNA was extracted from muscle tissue using the modified low-salt CTAB

599    extraction of high-quality DNA procedure [83]. DNA quality and concentration were

600    assessed using agarose gel electrophoresis, a NanoDrop Spectrophotometer (Thermo

601    Scientific), and a Qubit 2.0 Fluorometer (ThermoFisher). A PCR-free library was

602    prepared using NEBNext Ultra DNA Library Prep Kit for Illumina. Sequencing was

603    performed on two lanes of a Hiseq2500 platform (PE250). Library preparation and

604    sequencing were performed by the Beijing Novogene Bioinformatics Technology Co.

605    Ltd. Raw reads were quality-filtered and trimmed of adaptors using Trimmomatic-0.36

606    [84] with default parameters. In total, the genomic shotgun dataset included 7,785,846

607    reads. After filtering and trimming, 7,275,133 reads covering a total length of

608    1,635,569,256 bp remained. Thus, the sequencing coverage is 0.134. Filtered, trimmed

609    reads were assembled into contigs using dipSPAdes 3.11.1 [85] with default parameters,

610    yielding 130,417 contigs with an N50 of 740 bp and a total length of 1,560,938,851 bp.

611

612    **Mining and classification of repeat elements**

613    The PiRATE pipeline was used as in the original publication [44], including the

614    following steps: 1) Contigs representing repetitive sequences were identified from the

615    assembly using similarity-based, structure-based, and repetitiveness-based approaches

616    applied    non-sequentially.    The    similarity-based    detection    programs    included

22

617   RepeatMasker [86] and TE-HMMER [87]. The structure-based detection programs

618   included MITE-Hunter [88], SINE-finder [89], HelSearch [90], LTRharvest [91], and

619   MGEScan-nonLTR [92]. The repetitiveness-based detection programs included

620   TEdenovo [93] and RepeatScout [94]. 2) Contigs representing repeat family consensus

621   sequences were also identified from the cleaned, filtered, unassembled reads with

622   dnaPipeTE [95], which uses Trinity on subsamples of single-end reads to produce sets

623   of related repeat consensus sequences (e.g. representing multiple subfamilies within a

624   TE family). 3) Contigs identified by each individual program in Steps 1 and 2, above,

625   were filtered to remove those <100 bp in length and clustered with CD-HIT-est [96] to

626   reduce redundancy (100% sequence identity cutoff). This yielded a total of 62,699

627   contigs. 4) All 62,699 contigs were then clustered together with CD-HIT-est (100%

628   sequence identity cutoff), retaining the longest contig and recording the program that

629   classified it. 1860 contigs were filtered out at this step, and the majority (1669) were

630   contigs identified by RepeatMasker and TE-HMMER that were identical in sequence

631   but differed in length. 5) Repeat contigs were annotated as TEs to the levels of order

632   and superfamily in Wicker's hierarchical classification system [7], modified to include

633   several recently discovered TE superfamilies, using PASTEC [45] and were checked

634   manually to filter chimeric contigs and those annotated with conflicting evidence. 6)

635   All classified repeats ("known TEs" hereafter), along with the unclassified repeats

636   ("unknown repeats" hereafter) and putative multi-copy host genes, were combined to

637   produce a caecilian-derived repeat library.

638

639   **Characterization of the overall repeat element landscape**

640   Overlapping paired-end reads were merged using PEAR v.0.9.11 [97] with the

641   following parameter values based on our library insert size and trimming parameters:

642   min-assemble-length 36, max-assemble-length 490, min-overlap size 10. After merging

643   the remaining paired-end reads, 6,628,808 shotgun reads remained, with an average and

644   a total length of 236 and 1,560,938,851 bp, respectively. To calculate the percentage of

23

645    the caecilian genome composed of different TEs, the shotgun reads (including both

646    merged reads and singletons) were masked with RepeatMasker v-4.0.7 using two

647    versions of our caecilian-derived repeat library: one that included the unknown repeats

648    and one that excluded them. In both cases, simple repeats were identified using the

649    Tandem Repeat Finder module implemented in RepeatMasker. The overall results were

650    summarized at the levels of TE class, order, and superfamily. For each superfamily, we

651    then collapsed the contigs to 95% and 80% sequence similarity using CD-HIT-est to

652    provide an overall view of within-superfamily diversity; 80% is the sequence similarity

653    threshold used to define TE families [7].

654

655    **TE community diversity**

656    Diversity of the overall transposable element community in *I. bannanicus* was

657    summarized using the Shannon index $H' = -\sum P_i \ln(P_i)$ and the Simpson index

658    $D_1 = 1 - \sum P_i^2$ (i.e. the Gini-Simpson index), where $P_i$ is the proportion of sequences

659    belonging to TE superfamily *i* [51, 52]. In analogous applications of these diversity

660    indices to ecological communities, $P_i$ is the proportion of individuals that belong to

661    species *i.* To provide context for the *I. bannanicus* results, Shannon and Simpson

662    indices were also calculated for other vertebrate genomes representing diversity in

663    genome size as well as type of dataset. *Takifugu rubripes* (pufferfish, 0.4 Gb), *Gallus*

664    *gallus* (chicken, 1.3 Gb), *Xenopus tropicalis* (Western clawed frog, 1.7 Gb), *Anolis*

665    *carolinensis* (green anole lizard, 2.2 Gb), and *Homo sapiens* (human, 3.1 Gb) all have

666    full genome assemblies. For these five species, the perl script parseRM.pl [98] was used

667    to parse the raw output files downloaded from www.repeatmasker.org and obtain the

668    percentage of the genome occupied by each identified superfamily; ambiguous

669    classifications (i.e. to the level of order or class) were excluded. *Ambystoma mexicanum*

670    (the axolotl, a model salamander, 32 Gb), which has a much larger genome and,

671    consequently, less complete genome assembly, was also included; percentages of the

672    genome occupied by each identified superfamily were obtained from a previous study

24

673     [30]. Finally, four other salamanders that encompass a range of genome sizes were

674     included, each represented by low-coverage genome-skimming shotgun data:

675     *Desmognathus ochrophaeus* (15 Gb), *Batrachoseps nigriventris* (25 Gb), *Aneides*

676     *flavipunctatus* (44 Gb), and *Cryptobranchus alleganiensis* (55 Gb). Percentages of each

677     genome occupied by identified superfamilies were obtained from a previous study [32].

678

679     **Amplification history of transposable element superfamilies**

680     To summarize the overall amplification history of TE superfamilies and test for ongoing

681     activity, the perl script parseRM.pl [98] was used to parse the raw output files from

682     RepeatMasker (.align) and report the sequence divergence between each read and its

683     respective consensus sequence (parameter values = -l 50,1 and -a 5). The repeat library

684     used to mask the reads comprised the 50,471 TE contigs classified by the PiRATE

685     pipeline and clustered at 100% sequence similarity. Each TE superfamily is therefore

686     represented by multiple consensus contigs that represent ancestral sequences likely

687     corresponding to the family and subfamily TE taxonomic levels (i.e. not the distant

688     common ancestor of the entire superfamily). For each superfamily, histograms were

689     plotted to summarize the percent divergence of all reads from their closest (i.e. least

690     divergent) consensus sequence. These histograms do not allow the delineation between

691     different amplification dynamics scenarios (i.e. a single family with continuous activity

692     versus multiple families with successive bursts of activity). Rather, these global

693     overviews were examined for overall shapes consistent with ongoing activity (i.e. the

694     presence of TE loci <1% diverged from the ancestral sequence and a unimodal right-

695     skewed, J-shaped, or monotonically decreasing distribution).

696

697     **Ectopic recombination-mediated deletion of Gypsy and DIRS elements**

698     All genomic contigs > 3000 bp in length that were annotated to LTR/Gypsy were *de*

699     *novo* annotated using LTRpred to identify terminal and internal sequences [99].

700     Internal and terminal sequences were further confirmed by manually checking for

25

701 internal TE domains using NCBI

702 BLASTx (https://blast.ncbi.nlm.nih.gov/Blast.cgi?PROGRAM=blastx&PAGE_TYPE

703 =BlastSearch&LINK_LOC=blasthome) and for LTR sequences using the NCBI-

704 Blast2suite to align each contig sequence against itself. DIRS/DIRS superfamily

705 elements have a different structure than other LTR retrotransposons; their terminal

706 repeats are inverted. However, because they also include internal sequences

707 complementary to the terminals that facilitate rolling-circle amplification [54, 55], their

708 structure includes direct repeats that are expected to undergo ectopic recombination to

709 eliminate much of the internal sequence and one copy of the direct repeat sequence,

710 although to our knowledge this has not been previously investigated. Although these

711 deletions would not produce canonical solo LTRs, they, too, would produce elevated

712 abundances of terminal sequences relative to internal sequences. Typical DIRS

713 structure was confirmed visually and by using the NCBI-Blast2suite to align each

714 contig sequence against itself, and contigs that lacked the complete structure were

715 removed from further analysis. Internal sequences for both superfamilies were

716 conservatively defined to be bounded by the first and last TE domains. This yielded a

717 total of nine DIRS/DIRS contigs and 17 LTR/Gypsy contigs. DIRS/DIRS contigs had

718 an average terminal sequence length of 150 bp (range 61-343) and an average internal

719 sequence length of 5586 bp (range 4810-6012). LTR/Gypsy contigs had an average

720 terminal sequence length of 744 bp (range $127 - 3267$) and an average internal sequence

721 length of 1976 (range $243 - 4306$). To estimate levels of terminal sequences (LTRs or

722 TIRs) relative to internal sequences, genomic shotgun reads were mapped to the whole

723 genome assembly using bowtie2 in local alignment mode with very-sensitive-local

724 preset options and otherwise default parameters, increasing the G-value from the

725 default of 20 to 30, 40, and 50 to increase minimum alignment length for reads [100].

726 This analysis was performed twice: once treating all reads as unpaired and once using

727 merged paired-end reads plus unmerged reads. Average read depths across the terminal

728 and internal portion in each of the 26 focal DIRS/DIRS and LTR/Gypsy contigs were

26

729    estimated by scaling the number of hits by the lengths of the terminal and internal region.

730    From these estimates, the total terminal-to-internal sequence ratio (TT:I) was calculated

731    for each contig. In the absence of ectopic recombination mediated by terminal repeats,

732    this ratio would be 1:1; increasing levels of ectopic recombination would produce ratios >

733    1:1. We compared the results obtained for the caecilian with similar analyses that

734    included gigantic salamander genomes as well as vertebrates with more typical (i.e.

735    smaller) genomes [33].

736

737    **Transcriptome library creation, sequencing, assembly, and TE annotation**

738    Total RNA was extracted separately from heart, brain, liver, and tail tissue using TRIzol

739    (Invitrogen). For each sample, RNA quality and concentration were assessed using

740    agarose gel electrophoresis, a NanoPhotometer spectrophotometer (Implen, CA, USA),

741    a Qubit 2.0 Fluorometer (ThermoFisher), and an Agilent BioAnalyzer 2100 system

742    (Agilent Technologies, CA, USA) requiring an RNA integrity number (RIN) of eight

743    or higher. Equal quantities of RNA from these four tissues were pooled to build a single

744    transcriptome library. Sequencing libraries were generated using the NEBNext Ultra

745    RNA Library Prep Kit for Illumina following the manufacturer's protocol. After cluster

746    generation of the index-coded samples, the library was sequenced on one lane of an

747    Illumina Hiseq 4000 platform (PE 150). Library preparation and sequencing were

748    performed by the Beijing Novogene Bioinformatics Technology Co. Ltd.

749    Transcriptome sequences were filtered using Trimmomatic-0.36 (Bolger, et al. 2014)

750    with default parameters. Remaining reads were assembled using Trinity 2.5.1 (Grabherr,

751    et al. 2011). In total, 34,980,300 transcriptome reads were obtained, with a total length

752    of 5,247,045,000 bp. After filtering, 34,417,105 reads remained, with a total length of

753    5,027,542,505 bp. The assembly produced 348,822 contigs (i.e. putative assembled

754    transcripts) with the min, N50, max, and total length of contigs equal to 201; 357;

755    32,175; and 249,943,402 bp, respectively. Of these, 289,380 had expression levels of

756    TPM ≥ 0.01 and were analyzed further.

757        To annotate transcriptome contigs containing autonomous TEs, BLASTx was used

758   against the Transposable Element Protein Database (RepeatPeps.lib, downloaded from

759   https://github.com/rmhubley/RepeatMasker/blob/master/Libraries/ on April 20, 2019)

760   with an e-value cutoff of 1e-10. To annotate contigs containing non-autonomous TEs,

761   RepeatMasker was used with our caecilian-derived genomic repeat library of non-

762   autonomous TEs (LARD-, TRIM-, MITE-, and SINE-annotated contigs; Table 2) and

763   the requirement that the transcriptome/genome contig overlap was >80 bp long, >80%

764   similar in sequence, and covered >80% of the length of the genomic contig. Contigs

765   annotated as conflicting autonomous and non-autonomous TEs were filtered out. To

766   yield a rough estimate of the number of active TE families per superfamily, CD-HIT-

767   est was used to cluster the contigs annotated to each superfamily at the level of 80%

768   sequence similarity.

769        To identify contigs that contained an endogenous caecilian gene, the Trinotate

770   annotation suite was used with e-value cutoffs of 1e-10 and 1e-5 for BLASTx and

771   BLASTp against the SwissProt database, respectively, and 1e-5 for HMMER against

772   the Pfam database [57]. To identify contigs that contained both a TE and an endogenous

773   caecilian gene (i.e. putative cases where a TE and a gene were co-transcribed on a single

774   transcript), all contigs that were annotated both by RepeatPeps and Trinotate were

775   examined, and the ones annotated by Trinotate to contain a TE-encoded protein (i.e. the

776   contigs where RepeatPeps and Trinotate annotations were in agreement) were not

777   further considered. The remaining contigs annotated by Trinotate to contain a non-TE

778   gene (i.e. an endogenous caecilian gene) and also annotated either by RepeatPeps to

779   include a TE-encoded protein or by RepeatMasker to include a non-autonomous TE

780   were identified for further examination and expression-based analysis.

781

782   **Transposable element expression**

783   To generate a point estimate of overall TE expression in the somatic transcriptome,

784   transcript abundance levels were quantified with RSEM (because of its capacity to

785    model multi-mapping reads) using the Bowtie short-read aligner. Transcriptome

786    contigs with TPM < 0.01 were filtered out. To yield TE-superfamily-wide expression

787    level estimates, TPM values were summed across all contigs annotated to the same TE

788    superfamily. For comparison, TPM values were summed for all endogenous (i.e. non-

789    TE) caecilian genes. Pearson's correlation coefficient was used to test for a relationship

790    between genomic TE abundance (measured as log-transformed percentage of the

791    genome occupied per TE superfamily) and TE expression level (measured as log-

792    transformed total TPM per TE superfamily). We note that with only a single sample,

793    any more detailed analyses of expression levels are not appropriate. Contigs annotated

794    to contain both TEs and endogenous caecilian genes were excluded from these analyses.

795    Instead, these putative TE/gene contigs were ranked by expression level, and the 20

796    most highly expressed were examined by eye to determine the spatial relationship

797    between the TE and gene BLAST results producing the annotations. Nine contigs with

798    apparently spurious TE annotations (seven of which reflected a single likely mis-

799    annotation of an LTR/Pao protein in the RepeatPeps database) were reclassified as

800    endogenous genes, and the remaining contigs were characterized as having the TE 1)

801    on the same or different strand as the gene, and 2) upstream or downstream of the gene.

802    Finally, TPM values were summed across all putative TE/gene contigs to yield a global

803    estimate of expression levels of TE/gene combinations that are co-transcribed on a

804    single transcript.

805

806    **Ethical Statement**

807    The study specimen was collected and dissected following Animal Care & Use

808    Protocols approved by Chengdu Institute of Biology, Chinese Academy of Sciences.

809

810    **Data availability**

811    Genomic shotgun and transcriptome sequences have been deposited in the National

812    Genomics Data Center (accession numbers: SAMC207357, SAMC207358).

813

## Authors' contributions

815 JW and RM designed and supervised the experiments. JW, MI, HW, and CS

816 participated in bioinformatic analysis. MI and SS participated in genome size

817 measurement. YG collected the sample from the field. JJ and JL advised on data

818 analysis. JW and RM wrote the manuscript with input from all authors. All authors read

819 and approved the final manuscript.

820

## Competing interests

822 The authors have declared no competing interests.

823

## Acknowledgments

834

## Authors' ORCID IDs

836 0000-0003-4318-8923 (Wang, J)

837 0000-0001-9481-0693 (Itgen, MW)

838 0000-0003-0960-8939 (Wang, HJ)

839     0000-0002-2380-180X (Gong, YZ)

840     0000-0002-1051-7797 (Jiang JP)

841     0000-0003-1799-194X (Li, JT)

842     0000-0001-7476-9224 (Sun, C)

843     0000-0002-9444-028X (Sessions, S)

844     0000-0003-3875-1988 (Mueller, RL)

845

## Literature Cited

847

848     [1] Kidwell MG, Lisch DR. Perspective: transposable elements, parasitic DNA, and
849     genome evolution. Evolution 2001;55:1-24.


850     [2] Aparicio S, Chapman J, Stupka E, Putnam N, Chia JM, Dehal P, et al. Whole-
851     genome shotgun assembly and analysis of the genome of *Fugu rubripes*. Science
852     2002;297:1301-10.


853     [3] Jiao Y, Peluso P, Shi J, Liang T, Stitzer MC, Wang B, et al. Improved maize reference
854     genome with single-molecule technologies. Nature 2017;546:524.


855     [4] McClintock B. The origin and behavior of mutable loci in maize. Proc Natl Acad
856     Sci U S A 1950;36:344-55.


857     [5] Piegu B, Asgari S, Bideshi D, Federici BA, Bigot Y. Evolutionary relationships of
858     iridoviruses and divergence of ascoviruses from invertebrate iridoviruses in the
859     superfamily Megavirales. Mol Phylogenet Evol 2015;84:44-52.


860     [6] Jurka J, Kapitonov VV, Pavlicek A, Klonowski P, Kohany O, Walichiewicz J.
861     Repbase Update, a database of eukaryotic repetitive elements. Cytogenet Genome Res
862     2005;110:462-7.

863    [7] Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, et al. A unified
864    classification system for eukaryotic transposable elements. Nat Rev Genet 2007;8:973-
865    82.

866    [8] Finnegan DJ. Eukaryotic transposable elements and genome evolution. Trends
867    Genet 1989;5:103-7.

868    [9] Goerner-Potvin P, Bourque G. Computational tools to unmask transposable
869    elements. Nat Rev Genet 2018;19:688-704.

870    [10] Arkhipova IR. Using bioinformatic and phylogenetic approaches to classify
871    transposable elements and understand their complex evolutionary histories. Mob DNA
872    2017;8:19.

873    [11] Pasquesi GIM, Adams RH, Card DC, Schield DR, Corbin AB, Perry BW, et al.
874    Squamate reptiles challenge paradigms of genomic repeat element evolution set by
875    birds and mammals. Nat Commun 2018;9:2774.

876    [12] Jangam D, Feschotte C, Betran E. Transposable element domestication as an
877    adaptation to evolutionary conflicts. Trends Genet 2017;33:817-31.

878    [13] Sotero-Caio CG, Platt RN, 2nd, Suh A, Ray DA. Evolution and diversity of
879    transposable elements in vertebrate genomes. Gen Biol Evol 2017;9:161-77.

880    [14] Sessions SK, Larson A. Developmental correlates of genome size in plethodontid
881    salamanders and their implications for genome evolution. Evolution 1987;41:1239-51.

882    [15] Olmo E. Nucleotype and cell size in vertebrates: a review. Basic Appl Histochem
883    1983;27:227-56.

884    [16] Szarski H. Cell size and the concept of wasteful and frugal evolutionary strategies.
885    J Theor Biol 1983;105:201-9.

886    [17] Hanken J, Wake DB. Miniaturization of body-size - organismal consequences and
887    evolutionary significance. Ann Rev Ecol Syst 1993;24:501-19.

888  [18] Roth G, Blanke J, Wake DB. Cell size predicts morphological complexity in the
889  brains of frogs and salamanders. Proc Natl Acad Sci U S A 1994;91:4796-800.

890  [19] Simonin KA, Roddy AB. Genome downsizing, physiological novelty, and the
891  global dominance of flowering plants. PLoS Biol 2018;16:e2003706.

892  [20] Gregory TR. Genome size and developmental complexity. Genetica 2002;115:131-
893  46.

894  [21] Naville M, Henriet S, Warren I, Sumic S, Reeve M, Volff JN, et al. Massive
895  changes of genome size driven by expansions of non-autonomous transposable
896  elements. Curr Biol 2019;29:1161-8 e6.

897  [22] Mueller RL. Genome biology and the evolution of cell-size diversity. Cold Spring
898  Harb Perspect Biol 2015;7.

899  [23] Lynch M, Conery JS. The origins of genome complexity. Science 2003;302:1401-
900  4.

901  [24] Mueller RL. piRNAs and evolutionary trajectories in genome size and content. J
902  Mol Evol 2017;85:169-71.

903  [25] Gregory TR (2020), 'Animal Genome Size Database', http://www.genomesize.com.
904  Accessed 18 Aug 2020.

905  [26] Laurin M, Canoville A, Struble M, Organ C, de Buffrenil V. Early genome size
906  increase in urodeles. Comptes Rendus Palevol 2016;15:74-82.

907  [27] Sun C, Shepard DB, Chong RA, Lopez Arriaza J, Hall K, Castoe TA, et al. LTR
908  retrotransposons contribute to genomic gigantism in plethodontid salamanders.
909  Genome Biol Evol 2012;4:168-83.

910  [28] Sun C, Lopez Arriaza JR, Mueller RL. Slow DNA loss in the gigantic genomes of
911  salamanders. Genome Biol Evol 2012;4:1340-8.

912    [29] Elewa A, Wang H, Talavera-Lopez C, Joven A, Brito G, Kumar A, et al. Reading
913    and editing the Pleurodeles waltl genome reveals novel features of tetrapod
914    regeneration. Nat Commun 2017;8:2286.

915    [30] Nowoshilow S, Schloissnig S, Fei JF, Dahl A, Pang AWC, Pippel M, et al. The
916    axolotl genome and the evolution of key tissue formation regulators. Nature
917    2018;554:50-5.

918    [31] Madison-Villar MJ, Sun C, Lau NC, Settles ML, Mueller RL. Small RNAs from a
919    big genome: the piRNA pathway and transposable elements in the salamander species
920    Desmognathus fuscus. J Mol Evol 2016;83:126-36.

921    [32] Sun C, Mueller RL. Hellbender genome sequences shed light on genomic
922    expansion at the base of crown salamanders. Gen Biol Evol 2014;6:1818-29.

923    [33] Frahry MB, Sun C, Chong R, Mueller RL. Low levels of LTR retrotransposon
924    deletion by ectopic recombination in the gigantic genomes of salamanders. J Mol Evol
925    2015;80:120-9.

926    [34] Mueller RL, Jockusch EL. Jumping genomic gigantism. Nat Ecol Evol
927    2018;2:1687-8.

928    [35] Liedtke HC, Gower DJ, Wilkinson M, Gomez-Mestre I. Macroevolutionary shift
929    in the size of amphibian genomes and the role of life history and climate. Nat Ecol Evol
930    2018;2:1792-9.

931    [36] Edwards RJ, Tuipulotu DE, Amos TG, O'Meally D, Richardson MF, Russell TL,
932    et al. Draft genome assembly of the invasive cane toad, Rhinella marina. Gigascience
933    2018;7.

934    [37] Sun YB, Xiong ZJ, Xiang XY, Liu SP, Zhou WW, Tu XL, et al. Whole-genome
935    sequence of the Tibetan frog Nanorana parkeri and the comparative evolution of
936    tetrapod genomes. Proc Natl Acad Sci U S A 2015;112:E1257-62.

937    [38] Rogers RL, Zhou L, Chu C, Marquez R, Corl A, Linderoth T, et al. Genomic
938    takeover by transposable elements in the strawberry poison frog. Mol Biol Evol
939    2018;35:2913-27.

940 [39] Hellsten U, Harland RM, Gilchrist MJ, Hendrix D, Jurka J, Kapitonov V, et al. The
941 genome of the Western clawed frog Xenopus tropicalis. Science 2010;328:633-6.

942 [40] Hammond SA, Warren RL, Vandervalk BP, Kucuk E, Khan H, Gibb EA, et al. The
943 North American bullfrog draft genome provides insight into hormonal regulation of
944 long noncoding RNA. Nat Commun 2017;8:1433.

945 [41] Li Y, Ren Y, Zhang D, Jiang H, Wang Z, Li X, et al. Chromosome-level assembly
946 of the mustache toad genome using third-generation DNA sequencing and Hi-C
947 analysis. GigaScience 2019;8.

948 [42] Genome 10K Community of Scientists. Genome 10K: a proposal to obtain whole-
949 genome sequence for 10,000 vertebrate species. J Hered 2009;100:659-74.

950 [43] Torres-Sánchez M, Creevey CJ, Kornobis E, Gower DJ, Wilkinson M, San Mauro
951 D. Multi-tissue transcriptomes of caecilian amphibians highlight incomplete
952 knowledge of vertebrate gene families. DNA Research 2018;26:13-20.

953 [44] Berthelier J, Casse N, Daccord N, Jamilloux V, Saint-Jean B, Carrier G. A
954 transposable element annotation pipeline and expression analysis reveal potentially
955 active elements in the microalga Tisochrysis lutea. BMC Genomics 2018;19:378.

956 [45] Hoede C, Arnoux S, Moisset M, Chaumier T, Inizan O, Jamilloux V, et al. PASTEC:
957 an automatic transposable element classification tool. PLoS One 2014;9:e91929.

958 [46] Tuomisto H. An updated consumer's guide to evenness and related indices. Oikos
959 2012;121:1203-18.

960 [47] Venner S, Feschotte C, Biemont C. Dynamics of transposable elements: towards a
961 community ecology of the genome. Trends Genet 2009;25:317-23.

962 [48] Linquist S, Saylor B, Cottenie K, Elliott TA, Kremer SC, Ryan Gregory T.
963 Distinguishing ecological from evolutionary approaches to transposable elements. Biol
964 Rev 2013;88:573-84.

965 [49] Linquist S, Cottenie K, Elliott TA, Saylor B, Kremer SC, Gregory TR. Applying
966 ecological models to communities of genetic elements: the case of neutral theory. Mol
967 Ecol 2015;24:3232-42.

35

968  [50] Saylor B, Kremer SC, Gregory TR, Cottenie K. Genomic environments and their
969  influence on transposable element communities. bioRxiv 2019:667121.

970  [51] Simpson EH. Measurement of diversity. Nature 1949;163:688-.

971  [52] Shannon CE. A mathematical theory of communication. Bell Syst Tech J
972  1948;27:379-423.

973  [53] Sun C, Shepard DB, Chong RA, Lopez Arriaza J, Hall K, Castoe TA, et al. LTR
974  retrotransposons contribute to genomic gigantism in plethodontid salamanders. Gen
975  Biol Evol 2012;4:168-83.

976  [54] Piednoël M, Gonçalves IR, Higuet D, Bonnivard E. Eukaryote DIRS1-like
977  retrotransposons: an overview. BMC Genomics 2011;12:621.

978  [55] Poulter RTM, Goodwin TJD. DIRS-1 and the other tyrosine recombinase
979  retrotransposons. Cytogenet Genome Res 2005;110:575-88.

980  [56] Petrov DA, Aminetzach YT, Davis JC, Bensasson D, Hirsh AE. Size matters: Non-
981  LTR retrotransposable elements and ectopic recombination in Drosophila. Mol Biol
982  Evol 2003;20:880-92.

983  [57] Bryant DM, Johnson K, DiTommaso T, Tickle T, Couger MB, Payzin-Dogru D, et
984  al. A tissue-mapped axolotl de novo transcriptome enables identification of limb
985  regeneration factors. Cell Reports 2017;18:762-76.

986  [58] Mateo L, Ullastres A, González J. A transposable element insertion confers
987  xenobiotic resistance in Drosophila. PLOS Genet 2014;10:e1004560.

988  [59] Schrader L, Schmitz J. The impact of transposable elements in adaptive evolution.
989  Mol Ecol 2019;28:1537-49.

990  [60] Keinath MC, Timoshevskiy VA, Timoshevskaya NY, Tsonis PA, Voss SR, Smith
991  JJ. Initial characterization of the large genome of the salamander Ambystoma
992  mexicanum using shotgun and laser capture chromosome sequencing. Sci Rep
993  2015;5:16413.

[61] Zimin AV, Puiu D, Hall R, Kingan S, Clavijo BJ, Salzberg SL. The first near-complete assembly of the hexaploid bread wheat genome, Triticum aestivum. Gigascience 2017;6:gix097.

[62] Smith JJ, Timoshevskaya N, Timoshevskiy VA, Keinath MC, Hardy D, Voss SR. A chromosome-scale assembly of the axolotl genome. Genome Res 2019;29:317-24.

[63] Stevens KA, Wegrzyn JL, Zimin A, Puiu D, Crepeau M, Cardeno C, et al. Sequence of the sugar pine megagenome. Genetics 2016;204:1613-26.

[64] Nystedt B, Street NR, Wetterbom A, Zuccolo A, Lin YC, Scofield DG, et al. The Norway spruce genome sequence and conifer genome evolution. Nature 2013;497:579-84.

[65] Neale DB, Wegrzyn JL, Stevens KA, Zimin AV, Puiu D, Crepeau MW, et al. Decoding the massive genome of loblolly pine using haploid DNA and novel assembly strategies. Genome Biol 2014;15:R59.

[66] Kelly LJ, Renny-Byfield S, Pellicer J, Macas J, Novã kP, Neumann P, et al. Analysis of the giant genomes of Fritillaria (Liliaceae) indicates that a lack of DNA removal characterizes extreme expansions in genome size. New Phytol 2015;208:596-607.

[67] Weiss-Schneeweiss H, Leitch AR, McCann J, Jang T-S, Macas J. Employing next generation sequencing to explore the repeat landscape of the plant genome. In: Hörandl E., Appelhans M. S. eds). Next Generation Sequencing in Plant Systematics. Koeltz Scientific Books, 2015, 155-79.

[68] Langley CH, Montgomery E, Hudson R, Kaplan N, Charlesworth B. On the role of unequal exchange in the containment of transposable element copy number. Genet Res 1988;52:223-35.

[69] Furano AV, Duvernell DD, Boissinot S. L1 (LINE-1) retrotransposon diversity differs dramatically between mammals and fish. Trends Genet 2004;20:9-14.

[70] Boissinot S, Sookdeo A. The evolution of LINE-1 in vertebrates. Gen Biol Evol 2016;8:3485-507.

1022   [71] Abrusán G, Krambeck H-J. Competition may determine the diversity of
1023   transposable elements. Theoret Pop Biol 2006;70:364-75.

1024   [72] Elliott TA, Gregory TR. Do larger genomes contain more diverse transposable
1025   elements? BMC Evol Biol 2015;15:69.

1026   [73] Kijima TE, Innan H. Population genetics and molecular evolution of DNA
1027   sequences in transposable elements. I. A simulation framework. Genetics 2013;195:957.

1028   [74] Roessler K, Bousios A, Meca E, Gaut BS. Modeling interactions between
1029   transposable elements and the plant epigenetic response: a surprising reliance on
1030   element retention. Gen Biol Evol 2018;10:803-15.

1031   [75] Price AL, Eskin E, Pevzner PA. Whole-genome analysis of Alu repeat elements
1032   reveals complex evolutionary history. Genome Res 2004;14:2245-52.

1033   [76] Hillier LW, Miller W, Birney E, Warren W, Hardison RC, Ponting CP, et al.
1034   Sequence and comparative analysis of the chicken genome provide unique perspectives
1035   on vertebrate evolution. Nature 2004;432:695-716.

1036   [77] Bendall ML, de Mulder M, Iñiguez LP, Lecanda-Sánchez A, Pérez-Losada M,
1037   Ostrowski MA, et al. Telescope: Characterization of the retrotranscriptome by accurate
1038   estimation of transposable element expression. PLoS Comput Biol 2019;15:e1006453.

1039   [78] Coate JE, Doyle JJ. Variation in transcriptome size: are we getting the message?
1040   Chromosoma 2015;124:27-43.

1041   [79] Biscotti MA, Gerdol M, Canapa A, Forconi M, Olmo E, Pallavicini A, et al. The
1042   lungfish transcriptome: a glimpse into molecular evolution events at the transition from
1043   water to land. Sci Rep 2016;6:21571.

1044   [80] Castoe TA, Hall KT, Guibotsy Mboulas ML, Gu W, de Koning APJ, Fox SE, et al.
1045   Discovery of highly divergent repeat landscapes in snake genomes using high-
1046   throughput sequencing. Gen Biol Evol 2011;3:641-53.

1047   [81] Ji Y, Marra NJ, DeWoody JA. Comparative analysis of active retrotransposons in
1048   the transcriptomes of three species of heteromyid rodents. Gene 2015;562:95-106.

1049    [82] Gregory TR (2020), 'Gregory, T. R. Animal Genome Size Database
1050    (http://www.genomesize.com)'.

1051    [83] Arseneau JR, Steeves R, Laflamme M. Modified low-salt CTAB extraction of
1052    high-quality DNA from contaminant-rich tissues. Mol Ecol Resour 2017;17:686-93.

1053    [84] Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
1054    sequence data. Bioinformatics 2014;30:2114-20.

1055    [85] Safonova Y, Bankevich A, Pevzner PA. dipSPAdes: assembler for highly
1056    polymorphic diploid genomes. J Comput Biol 2015;22:528-45.

1057    [86] Smit AFA, Hubley R, Green P (2013-2015), 'RepeatMasker Open-4.0'.

1058    [87] Eddy SR. Accelerated profile HMM searches. PLoS Comp Biol 2011;7:e1002195.

1059    [88] Han Y, Wessler SR. MITE-Hunter: a program for discovering miniature inverted-
1060    repeat transposable elements from genomic sequences. Nucleic Acids Res
1061    2010;38:e199.

1062    [89] Wenke T, Dobel T, Sorensen TR, Junghans H, Weisshaar B, Schmidt T. Targeted
1063    identification of short interspersed nuclear element families shows their widespread
1064    existence and extreme heterogeneity in plant genomes. Plant Cell 2011;23:3117-28.

1065    [90] Yang G, Nagel DH, Feschotte C, Hancock CN, Wessler SR. Tuned for
1066    transposition: molecular determinants underlying the hyperactivity of a Stowaway
1067    MITE. Science 2009;325:1391-4.

1068    [91] Ellinghaus D, Kurtz S, Willhoeft U. LTRharvest, an efficient and flexible software
1069    for de novo detection of LTR retrotransposons. BMC Bioinformat 2008;9:18.

1070    [92] Rho M, Tang H. MGEScan-non-LTR: computational identification and
1071    classification of autonomous non-LTR retrotransposons in eukaryotic genomes.
1072    Nucleic Acids Res 2009;37:e143.

1073    [93] Flutre T, Duprat E, Feuillet C, Quesneville H. Considering transposable element
1074    diversification in de novo annotation approaches. PLoS One 2011;6:e16526.

1075 [94] Price AL, Jones NC, Pevzner PA. De novo identification of repeat families in large
1076 genomes. Bioinformatics 2005;21 Suppl 1:i351-8.

1077 [95] Goubert C, Modolo L, Vieira C, ValienteMoro C, Mavingui P, Boulesteix M. De
1078 novo assembly and annotation of the Asian tiger mosquito (Aedes albopictus)
1079 repeatome with dnaPipeTE from raw genomic reads and comparative analysis with the
1080 yellow fever mosquito (Aedes aegypti). Genome Biol Evol 2015;7:1192-205.

1081 [96] Li W, Jaroszewski L, Godzik A. Clustering of highly homologous sequences to
1082 reduce the size of large protein databases. Bioinformatics 2001;17:282-3.

1083 [97] Zhang J, Kobert K, Flouri T, Stamatakis A. PEAR: a fast and accurate Illumina
1084 Paired-End reAd mergeR. Bioinformatics 2013;30:614-20.

1085 [98] Kapusta A, Suh A, Feschotte C. Dynamics of genome size evolution in birds and
1086 mammals. Proc Nat Acad Sci U S A 2017;114:E1460.

1087 [99] Cho J, Benoit M, Catoni M, Drost H-G, Brestovitsky A, Oosterbeek M, et al.
1088 Sensitive detection of pre-integration intermediates of long terminal repeat
1089 retrotransposons in crop plants. Nat Plants 2019;5:26-33.

1090 [100] Langmead B, Salzberg SL. Fast gapped-read alignment with Bowtie 2. Nat
1091 Methods 2012;9:357-9.
1092

1093 **Figure Legends**

1094

1095 **Figure 1   Amplification plots for transposable element, TE, superfamilies**

1096 The majority (17/19) suggest current superfamily activity. Note that the y-axes differ

1097 in scale.

1098

1099 **Figure 2   Ratio of total terminal sequence to internal sequence for two TE**

1100 **superfamilies**

1101   A ratio of 1:1 is expected in the absence of ectopic-recombination-mediated deletion.

1102   S: single−end alignment. P: paired−end alignment. 20/50: minimum alignment score

1103   (local mode).

1104

1105   **Figure 3    Expression levels of genes and TEs**

1106   Black lines and white boxes show median and interquartile range values. Red lines

1107   show probability densities. TPM: transcripts per million. TE: transposable element.

1108

1109   **Figure 4    Genomic abundance and somatic expression level of TE superfamilies**

1110   **are strongly correlated ($\rho = 0.879$, $p < 0.001$)**

1111   TPM: transcripts per million.

1112

1113   **Figure 5 Predicted relationships between TE abundance (genome size), TE**

1114   **diversity, and proportion of active TEs from seven different models**

1115

1116   **Tables**

1117

1118   **Table 1    Repeat contigs identified by different methods/software of the PiRATE**
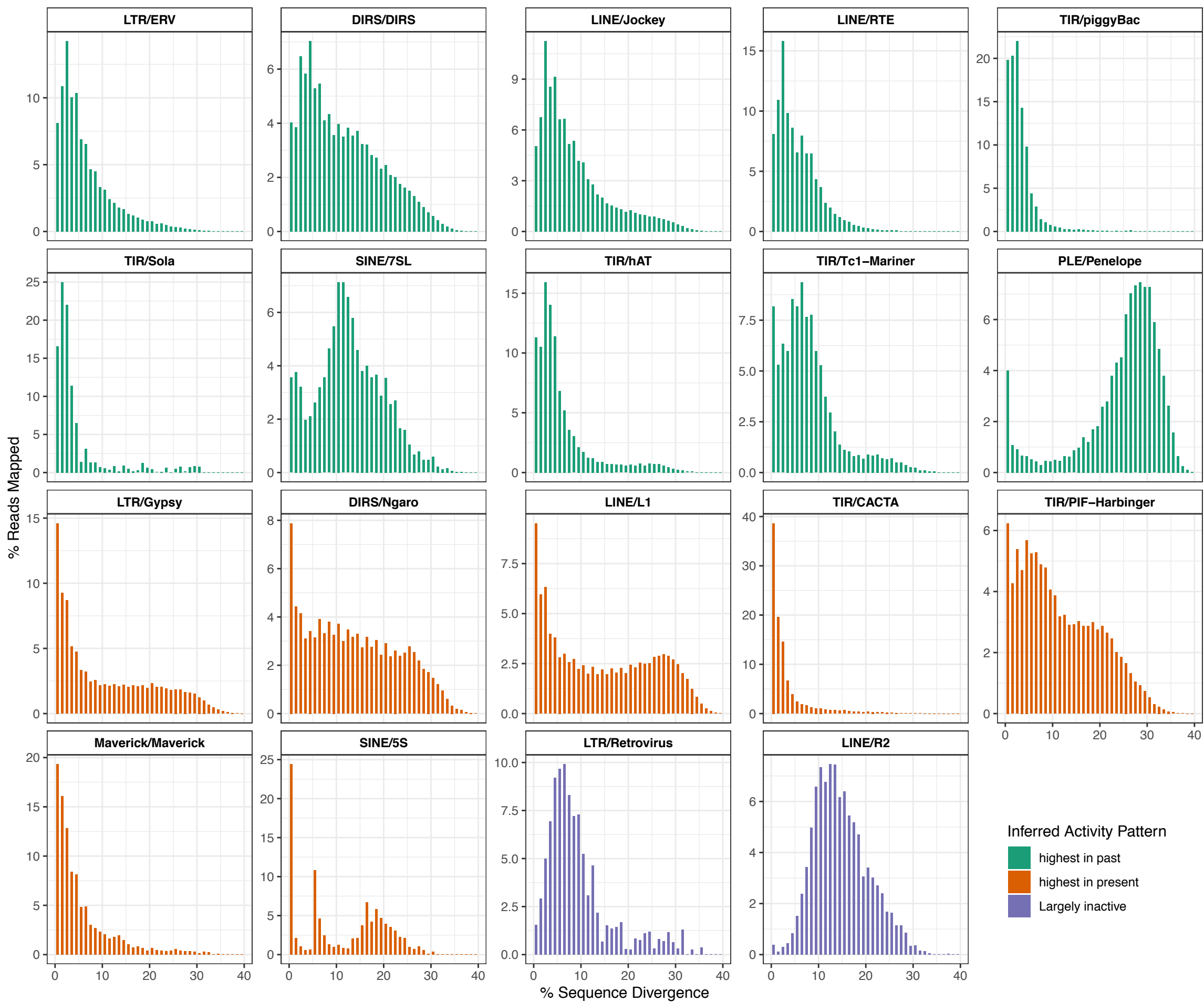
1119   **pipeline**

1120

1121   **Table 2    Classification of repeat contigs and summary of repeats detected in the**

1122   **genome and somatic transcriptome**

1123

1124   **Table 3    Diversity indices summarizing the TE communities from 11 vertebrate**

1125   **genomes**

1126

1127   **Table 4    Overall summary of transcriptome annotation (contigs with TPM $\geq 0.01$)**

Figure: Histograms of % Reads Mapped versus % Sequence Divergence for transposable element families, colored by Inferred Activity Pattern: highest in past (green), highest in present (orange), and Largely inactive (purple). Panels: LTR/ERV, DIRS/DIRS, LINE/Jockey, LINE/RTE, TIR/piggyBac, TIR/Sola, SINE/7SL, TIR/hAT, TIR/Tc1−Mariner, PLE/Penelope, LTR/Gypsy, DIRS/Ngaro, LINE/L1, TIR/CACTA, TIR/PIF−Harbinger, Maverick/Maverick, SINE/5S, LTR/Retrovirus, LINE/R2.

S: single−end alignment. P: paired−end alignment. 20/50: minimum alignment score (local mode)

TE Diversity

Proportion of Active TEs

TE Abundance/Genome Size

Petrov 2003
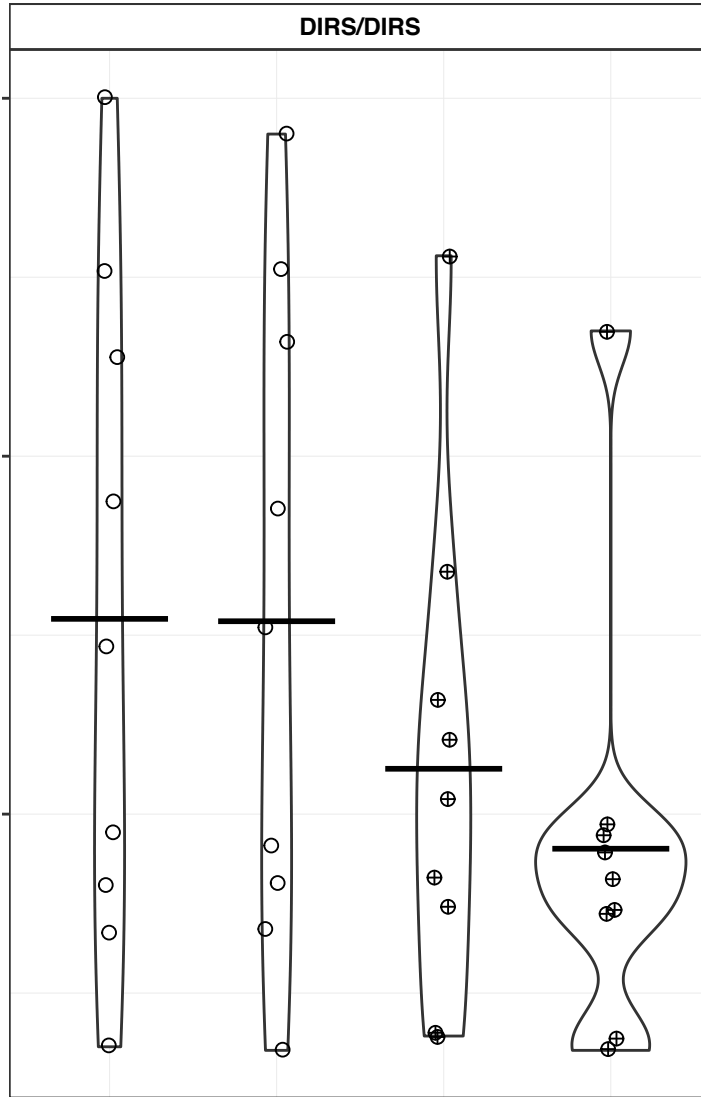Furano 2004
Boissinot 2016
Abrusan 2006

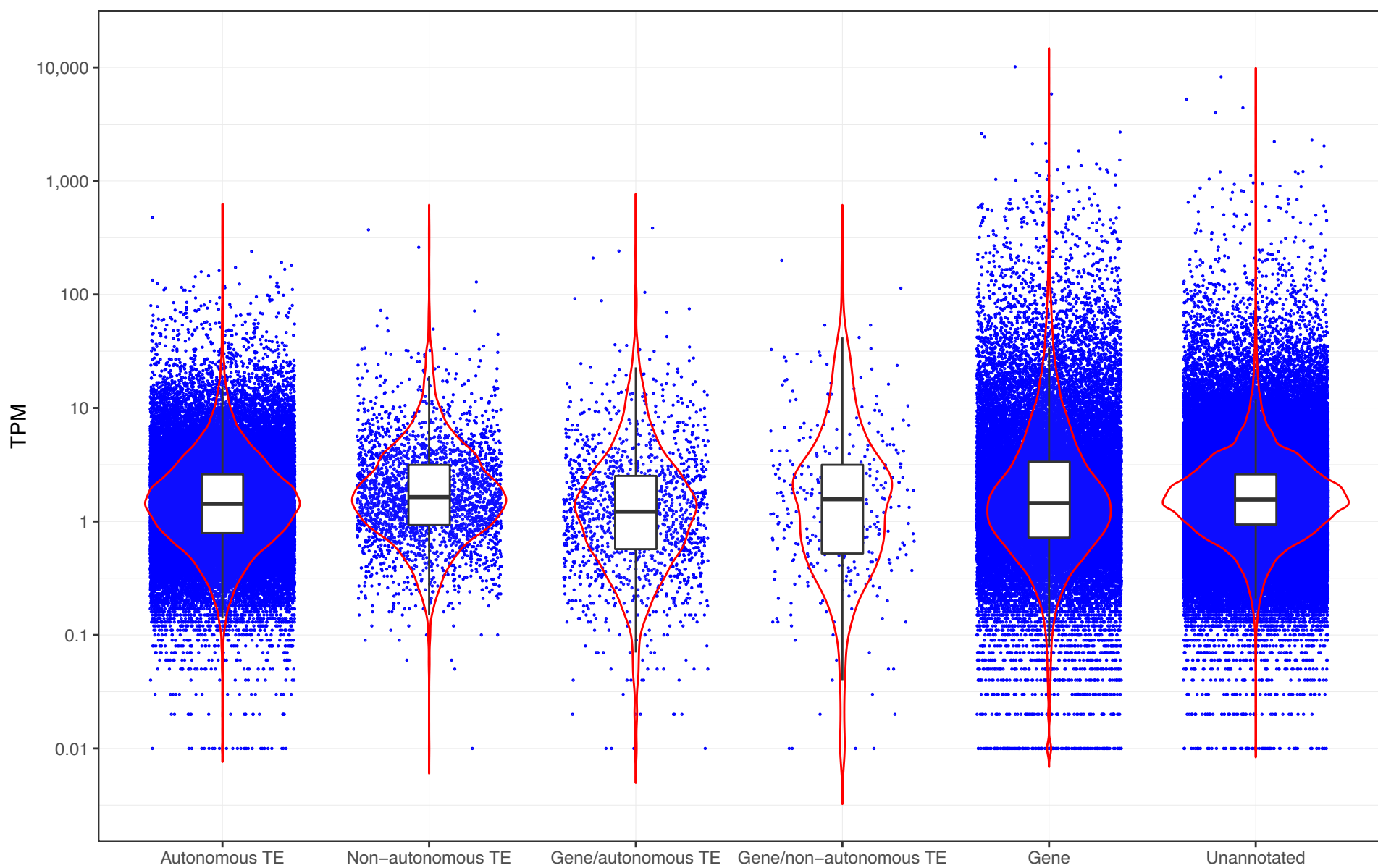Elliot 2015

Kijima 2013
Roessler 2018

**Table 1  Repeat contigs identified by different methods/software of the PiRATE pipeline**

| TE-mining method | Software | Repeats clustered at 100% identify |
|---|---|---|
| Similarity-based | RepeatMasker | 37,123 (62.1%) |
| | TE-HMMER | 1,585 (2.7%) |
| Structure-based | HelSearch | 17 (0.0%) |
| | LTR harvest | 23 (0.0%) |
| | MGEScan | 0 |
| | MITE hunter | 12 (0.0%) |
| | SINE finder | 17 (0.0%) |
| Repetitiveness-based | TEdenovo | 102 (0.2%) |
| | RepeatScout | 1,786 (3.0%) |
| Repeat-building-based | dnaPipeTE | 19,160 (32.0%) |
| In total | | 59,825 (100%) |

**Table 2 Classification of repeat contigs and summary of repeats detected in the genome and somatic transcriptome**

| Order | Superfamily | Percent of Genome[a] | Genomic Contigs (100% Identical) | Genomic Contigs (95% Identical) | Genomic Contigs (<80% Identical) | Average Genomic Contig Length (100% identical) (bp) | Longest Genomic Contig (bp) | Transcriptome Contigs | Transcriptome Contigs (<80% Identical) | Total Expression (Summed TPM) | Average Expression (Min, Max) |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **Class I - Retrotransposons - Autonomous** | | | | | | | | | | | |
| LTR | *ERV* | 1.62-1.82 | 1,578 | 1,132 | 376 | 548 | 10,231 | 1,894 | 559 | 5,564.35 | 2.94 (0.03, 239.46) |
| | *Gypsy* | 1.10-1.35 | 1,324 | 1,159 | 523 | 665 | 5,544 | 1,080 | 662 | 1,441.05 | 1.33 (0.01, 17.75) |
| | *Retrovirus* | 0.01 | 13 | 13 | 9 | 591 | 1,198 | 19 | 15 | 19.76 | 1.04 (0.12, 4.25) |
| | *Bel-Pao* | 0.00 | 3 | 3 | 3 | 407 | 541 | 5 | 5 | 3.46 | 0.69 (0.33, 1.05) |
| | *Copia* | 0.00 | 2 | 2 | 2 | 222 | 297 | 124 | 99 | 117.54 | 0.95 (0.08, 5.93) |
| DIRS | *DIRS* | 25.88-30.20 | 26,507 | 23,221 | 8,513 | 346 | 4,452 | 25,426 | 12,652 | 68,259.45 | 2.68 (0.01, 476.25) |
| | *Ngaro* | 0.46-0.47 | 988 | 954 | 411 | 488 | 1,975 | 674 | 455 | 983.7 | 1.46 (0.02, 10.81) |
| PLE | *Penelope* | 0.19-0.22 | 542 | 540 | 469 | 226 | 1,686 | 1,649 | 1,454 | 2,728.03 | 1.65 (0.04, 71.46) |
| LINE | *Jockey* | 16.92-20.59 | 12,004 | 10,490 | 2,911 | 350 | 3,609 | 14,267 | 6,613 | 34,709.94 | 2.43 (0.01, 172.77) |
| | *L1* | 3.05-3.23 | 4,139 | 3,962 | 1,867 | 738 | 5,045 | 4,736 | 3,533 | 7,277.91 | 1.54 (0.01, 166.70) |
| | *RTE* | 1.50-1.60 | 1,020 | 883 | 130 | 269 | 2,717 | 863 | 356 | 3,412.23 | 3.95 (0.01, 124.43) |
| | *R2* | 0.12-0.22 | 51 | 48 | 20 | 290 | 1,011 | 243 | 110 | 550.38 | 2.26 (0.02, 16.34) |
| | *I* | | - | - | | - | - - | 2 | 2 | 3.93 | 1.97 (0.75, 3.48) |
| | *Unknown LINE* | 0.03 | 5 | - | 3 | 807 | 1743 | - | - | - | - |
| **Class I - Retrotransposons - Non-autonomous** | | | | | | | | | | | |
| SINE | *7SL* | 0.09 | 123 | 108 | 47 | 263 | 1,030 | - | - | - | - |
| | *5S* | 0.02 | 25 | 25 | 14 | 194 | 1,384 | 64 | 59 | 156.72 | 2.45 (0.29, 13.91) |
| | *tRNA* | 0.00 | 11 | 11 | 5 | 158 | 294 | 208 | 193 | 539.05 | 2.59 (0.17, 49.49) |
| | *Unknown SINE* | 0.55-1.66 | 203 | 146 | 41 | 235 | 498 | 484 | 370 | 2,509.77 | 5.19 (0.06, 371.09) |
| Retrotransposon Derivatives | TRIM | 0.44-1.78 | 229 | 159 | 53 | 432 | 3,226 | 601 | 301 | 1,759.19 | 2.93 (0.10, 33.17) |
| | LARD | 0.10-0.19 | 28 | 18 | 5 | 1928 | 10,505 | 4 | 1 | 39.14 | 9.79 (0.69, 32.83) |
| Unknown class I | | 0.03-0.19 | 61 | 58 | 55 | 261 | 1259 | - | - | - | - |
| **Class II - DNA Transposons - Subclass 1** | | | | | | | | | | | |
| TIR | *hAT* | 0.57-1.15 | 338 | 263 | 155 | 425 | 2,589 | 309 | 182 | 565.60 | 1.83 (0.04, 31.47) |
| | *CACTA* | 0.52-0.56 | 135 | 92 | 46 | 444 | 1,734 | 277 | 111 | 1,179.55 | 4.26 (0.11, 143.05) |
| | *Tc1-Mariner* | 0.49-0.59 | 548 | 438 | 189 | 370 | 2,028 | 1,180 | 421 | 4,312.73 | 3.65 (0.01, 93.19) |
| | *PIF-Harbinger* | 0.41-0.45 | 373 | 262 | 115 | 386 | 1,951 | 113 | 69 | 189.23 | 1.67 (0.11, 12.58) |
| | *PiggyBac* | 0.06-0.08 | 25 | 20 | 16 | 880 | 2,521 | 87 | 35 | 200.72 | 2.31 (0.07, 22.85) |
| | *Sola* | 0.01 | 12 | 11 | 7 | 608 | 2,150 | 11 | 5 | 9.72 | 0.88 (0.17, 2.22) |
| | *Mutator/MuDR* | 0.00 | 2 | 2 | 2 | 321 | 380 | 18 | 15 | 22.5 | 1.25 (0.22, 4.33) |
| | *P* | 0.00 | 2 | 2 | 1 | 237 | 250 | 2 | 2 | 3.27 | 1.64 (1.40, 1.87) |
| | *Zisupton* | | - | - | - | - | - | 7 | 2 | 16.9 | 2.41 (0.43-11.13) |
| | *Kolobok* | | - | - | - | - | - | 7 | 5 | 7.75 | 1.11 (0.25 – 2.68) |
| | *Academ* | | - | - | - | - | - | 4 | 3 | 18.0 | 4.50 (1.45 – 11.95) |
| | *Unknown TIR* | 0.04-0.05 | 8 | 7 | 5 | 307 | 621 | - | - | - | - |
| Crypton | *Crypton* | | - | - | - | - | - | 9 | 3 | 6.34 | 0.70 (0.05, 3.57) |
| Transposon Derivatives | MITE | 0.50-1.42 | 146 | 134 | 82 | 199 | 577 | 1,297 | 1,028 | 3,480.60 | 2.68 (0.01, 72.2) |
| **Class II - DNA Transposons - Subclass 2** | | | | | | | | | | | |
| Maverick | *Maverick* | 0.04-0.05 | 24 | 22 | 12 | 1037 | 5,305 | 80 | 31 | 122.99 | 1.54 (0.06, 6.87) |
| Helitron | *Helitron* | 0.00 | 2 | 2 | 2 | 298 | 352 | 20 | 13 | 65.84 | 3.29 (0.05, 14.14) |
| **Total** | | 54.74-68.03 | 50,471 | 44,187 | 16,089 | 395.4 | 10,505 | 55,764 | 29,364 | 140277.34 | 2.52 (0.01, 476.25) |

[a] First number is estimated including unknown repeats from the repeat library. Second number is estimated excluding unknown repeats.

**Table 3  Diversity indices summarizing the TE communities from 11 vertebrate genomes**

| Species | Common Name | Genome Size (Gb) | Shannon Index | Gini-Simpson Index | Dataset |
|---|---|---|---|---|---|
| *Takifugu rubripes* | Pufferfish | 0.4 | 2.10 | 1.00 | Full genome assembly |
| *Gallus gallus* | Chicken | 1.3 | 0.90 | 0.50 | Full genome assembly |
| *Xenopus tropicalis* | Western clawed frog | 1.7 | 2.24 | 0.90 | Full genome assembly |
| *Anolis carolinensis* | Green anole lizard | 2.2 | 2.41 | 0.91 | Full genome assembly |
| *Homo sapiens* | Human | 3.1 | 1.69 | 0.79 | Full genome assembly |
| *Ichthyophis bannanicus* | Banna caecilian | 12.2 | 1.45 | 0.67 | ~0.1X genome skimming |
| *Desmognathus ochrophaeus* | Allegheny Mountain dusky salamander | 15 | 1.61 | 0.71 | ~0.01X genome skimming |
| *Batrachoseps nigriventris* | Black-bellied slender salamander | 25 | 2.18 | 0.86 | ~0.01X genome skimming |
| *Ambystoma mexicanum* | Mexican axolotl salamander | 32 | 2.26 | 0.89 | Assembly |
| *Aneides flavipunctatus* | Speckled black salamander | 44 | 1.96 | 0.78 | ~0.01X genome skimming |
| *Cryptobranchus alleganiensis* | Hellbender salamander | 55 | 2.02 | 0.84 | ~0.01X genome skimming |

**Table 4  Overall summary of transcriptome annotation (contigs with TPM ≥ 0.01)**

| | Contigs (Percentage of total contigs) | Summed TPM (Percentage of total expression) | Maximum TPM | Minimum TPM | Average TPM | Mean Contig Length (bp) |
|---|---|---|---|---|---|---|
| Endogenous genes | 38,584 (13.3%) | 295,759 (29.6%) | 10,112.76 | 0.01 | 7.7 | 2,086 |
| Autonomous TEs | 53,106 (18.4%) | 131,793 (13.2%) | 476.25 | 0.01 | 2.5 | 570 |
| Non-Autonomous TEs | 2,658 (0.9%) | 8,484 (0.9%) | 371.09 | 0.01 | 3.2 | 785 |
| Gene/Autonomous TE | 1,445 (0.5%) | 4,859 (0.5%) | 383.94 | 0.01 | 3.7 | 2,161 |
| Gene/Non-autonomous TE | 342 (0.1%) | 1,584 (0.2%) | 198.8 | 0.01 | 4.6 | 2,800 |
| Unannotated | 193,245 (66.8%) | 555,776 (55.6%) | 8,224.66 | 0.01 | 2.9 | 537 |
| Total | 289,380 (100%) | 1,000,006 (100%) | 10,112.76 | 0.01 | 3.5 | 763 |