# Dynamic co-evolution of transposable elements and the piRNA pathway in African cichlid fishes

Miguel Vasconcelos Almeida[1,2,*], Moritz Blumer[3,13], Chengwei Ulrika Yuan[1,2,3,13], Pío Sierra[3], Jonathan L. Price[1,2], Fu Xiang Quah[1,3], Aleksandr Friman[4,5], Alexandra Dallaire[1,2,6], Grégoire Vernaz[2,3,7], Audrey L. K. Putman[1,2,3], Alan M. Smith[8], Domino A. Joyce[8], Falk Butter[9,10], Astrid D. Haase[4], Richard Durbin[3,11], M. Emília Santos[12], Eric A. Miska[1,2,11,*]

[1]Department of Biochemistry, University of Cambridge, Tennis Court Road, Cambridge, CB2 1GA, UK
[2]Wellcome/CRUK Gurdon Institute, University of Cambridge, Tennis Court Road, Cambridge, CB2 1QN, UK
[3]Department of Genetics, University of Cambridge, Downing Street, Cambridge, CB2 3EH, UK
[4]National Institute of Diabetes and Digestive and Kidney Diseases, National Institutes of Health, Bethesda, MD 20892, USA
[5]Biophysics Graduate Program, Institute for Physical Science and Technology, University of Maryland, College Park, Maryland 20742, USA
[6]Comparative Fungal Biology, Royal Botanic Gardens Kew, Jodrell Laboratory, Richmond TW9 3DS, UK
[7]Present address: Zoological Institute, Department of Environmental Sciences, University of Basel, Vesalgasse 1, Basel, 4051, Switzerland
[8]School of Natural Sciences, University of Hull, Hull, HU6 7RX, UK
[9]Institute of Molecular Biology (IMB), Quantitative Proteomics, Ackermannweg 4, Mainz, 55128, Germany
[10]Institute of Molecular Virology and Cell Biology, Friedrich-Loeffler-Institute, Südufer, Greifswald, 17493, Germany
[11]Wellcome Sanger Institute, Tree of Life, Wellcome Genome Campus, Hinxton, CB10 1SA, UK
[12]Department of Zoology, University of Cambridge, Downing Street, Cambridge, CB2 3EJ, UK
[13]These authors contributed equally.

*Correspondence: mdd34@cam.ac.uk, eam29@cam.ac.uk

## Abstract

East African cichlid fishes have diversified in an explosive fashion, but the (epi)genetic basis of the phenotypic diversity of these fishes remains largely unknown. Although transposable elements (TEs) have been associated with phenotypic variation in cichlids, little is known about their transcriptional activity and epigenetic silencing. Here, we describe dynamic patterns of TE expression in African cichlid gonads and during early development. Orthology inference revealed an expansion of *piwil1* genes in Lake Malawi cichlids, likely driven by PiggyBac TEs. The expanded *piwil1* copies have signatures of positive selection and retain amino acid residues essential for catalytic activity. Furthermore, the gonads of African cichlids express a Piwi-interacting RNA (piRNA) pathway that target TEs. We define the genomic sites of piRNA production in African cichlids and find divergence in closely related species, in line with fast evolution of piRNA-producing loci. Our findings suggest dynamic co-evolution of TEs and host silencing pathways in the African cichlid radiations. We propose that this co-evolution has contributed to cichlid genomic diversity.

## Introduction

The East African Great Lakes are home to prolific cichlid radiations, the most species-rich and phenotypically diverse adaptive radiations in vertebrates[1,2]. In the last 10 million years, more than 1,700 species of cichlid fishes (Cichlidae family) have evolved in virtually every lacustrine and riverine ecological niche in Lakes Victoria, Tanganyika, Malawi and surrounding bodies of water. The explosive diversification of East African cichlids is particularly striking in the haplochromine tribe and has resulted in astonishing variation in morphologies, colouration, diets, and behaviours[1,2]. The genetic and epigenetic basis for such phenotypic variability is of great interest and remains, by and large, unknown.

Initial genomic studies suggested very low genetic variability amongst East African cichlids[3]. In Lake Malawi cichlids, for example, the reported average single nucleotide polymorphism divergence between species pairs was 0.1-0.25%[3,4]. These low estimates were derived from approaches aligning short-read sequence data to a linear reference genome and generally ignore the contribution of structural variation. We have recently complemented these estimations using a pangenomic approach and long-read genome assemblies of representative Lake Malawi species[5]. With this approach, we estimated that 4.73-9.86% of Lake Malawi cichlid genomes can be attributed to interspecific structural variation[5]. Importantly, transposable elements (TEs) account for up to 74.65% of structural variant sequence. Thus, TEs comprise an underestimated source of genetic variability in East African cichlids.

TEs are diverse mobile genetic elements that inhabit nearly all eukaryotic genomes sequenced to date[6]. While most extant TEs and novel TE mobilisation events are selectively neutral or slightly deleterious to their hosts[7], several examples of TEs providing adaptive benefits to their hosts have been reported[8–10]. The TE landscapes of teleost fish genomes are highly dynamic[11–17], and cichlid genomes are no exception, as they contain varied TE populations with signs of recent transpositional activity[16,18]. TEs may be an important source of (epi)genetic variability that has fuelled the cichlid radiations. Consistent with this notion, presence/absence variation of TEs is associated with pigmentation traits[19,20], sex determination[21], and modulation of endogenous gene expression[18,22]. It has recently been shown that differentially methylated regions enriched in young TEs are associated with transcriptional changes[23], further supporting a role for TEs in modulating gene expression in cichlids. The same study found widespread DNA methylation at TEs, but besides this, little is known about the silencing pathways that direct TE silencing in cichlids and lead to the deposition of DNA methylation.

Several pathways have evolved in animals to silence TEs, particularly in the germline and early development to protect the next generations from deleterious effects of TE activity[8,24–30]. Here, we focus on the Piwi-interacting RNA (piRNA) pathway, a class of non-coding small RNAs (sRNAs) 21-35 ribonucleotides long, which drive silencing of TEs in the animal germline, including in fishes[27,31–33]. piRNAs bind to Piwi Argonaute proteins and guide them to target RNAs with base complementarity, leading to post-transcriptional and/or transcriptional silencing of their targets[26,27]. The latter can be achieved by piRNA-directed DNA methylation of targets. piRNA biogenesis is complex, requires a variety of co-factors, and can be conceptualised as two collaborating pathways that create sequence diverse piRNA populations in the animal germline: the ping-pong and phased biogenesis pathways[26,27,32,34–37]. These pathways depend mainly on the slicer activity of Piwi proteins, and endonucleolytic activity of Zucchini/PLD6 acting on long piRNA precursor transcripts.

110    The co-evolution of TE silencing factors and TEs is often thought to occur in the form
111    of an arms race. TE silencing factors, including those of the piRNA pathway, often
112    have signatures of fast, adaptive evolution that are interpreted as a consequence of
113    such an arms race[24,25,38–41]. These signatures include positive selection and lability in
114    terms of copy number variation, with recurrent gene duplications and turnover. Little
115    is known about the co-evolution of TE silencing pathways and TEs in East African
116    cichlids and whether these arms races could help fuel cichlid radiations.

118    Here, we describe dynamic TE expression in the gonads and early development of
119    African cichlids. We identify cichlid orthologs of known factors required for TE silencing
120    in vertebrates and discover an expanded repertoire of *piwil1* genes in Lake Malawi
121    cichlids, which may have been driven by PiggyBac TEs. The additional *piwil1* paralogs
122    retain amino acid residues required for the catalytic activity of the PIWI domain and
123    have signatures of adaptive evolution, suggesting acquisition of novel regulatory
124    functions. TE silencing factors are expressed in cichlid gonads, alongside an abundant
125    piRNA population with signatures consistent with active piRNA-driven TE silencing.
126    Lastly, we observe divergence in the genomic origins of piRNA production in closely
127    related Lake Malawi cichlids.

## Results

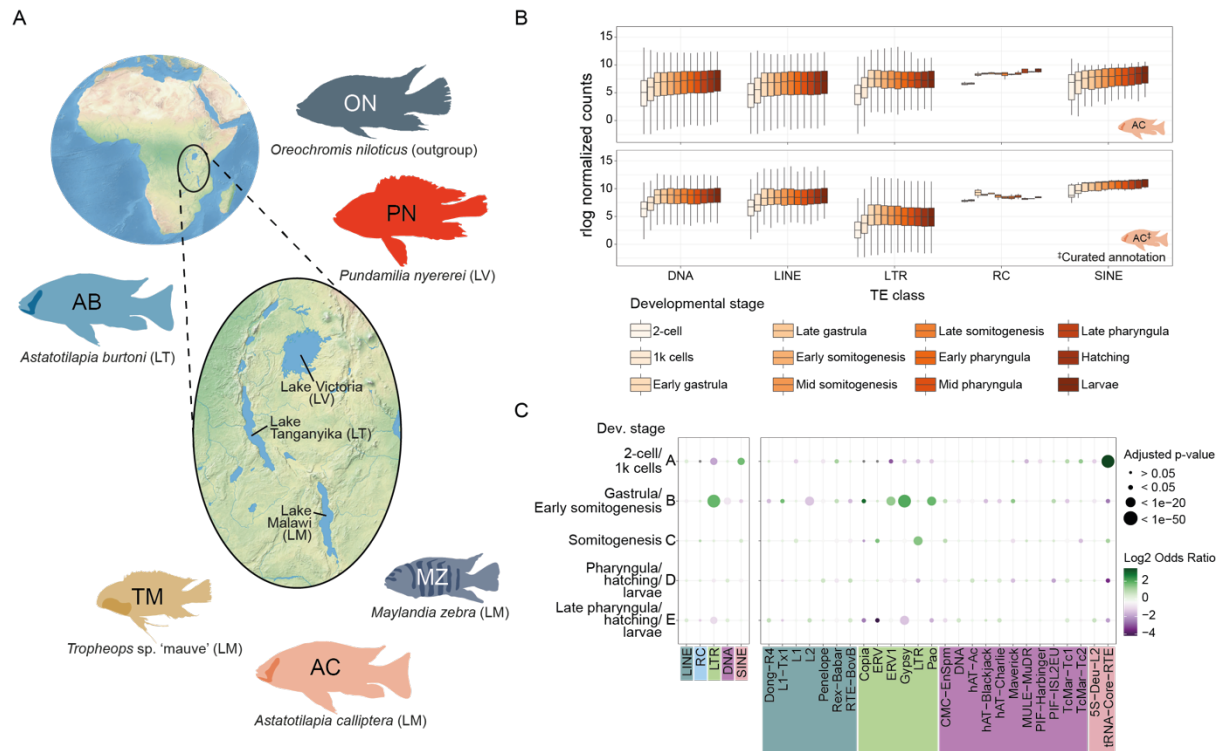### TE transcriptional activity in cichlid gonads and early development

134    To profile TE expression in African cichlids, we sequenced mRNAs of representative
135    species of haplochromine cichlids from each of the major East African Great Lakes
136    (**Figure 1A**). We chose *Pundamilia nyererei* (PN) as a representative for Lake Victoria,
137    *Astatotilapia burtoni* (AB) for Lake Tanganyika, and *Astatotilapia calliptera* (AC) for
138    Lake Malawi. To compare closely related species within the same Lake, we included
139    two species from Lake Malawi, alongside AC: *Maylandia zebra* (MZ), and *Tropheops*
140    sp. 'mauve' (TM). In addition, we included *Oreochromis niloticus* (ON, commonly
141    known as Nile tilapia) as an outgroup. ON is a representative of the tilapine tribe that
142    has a broad geographical distribution in Africa and is not as phenotypically diverse as
143    haplochromines[42]. We profiled TE expression in cichlid gonads, as these contain the
144    germline, where the arms race between TEs and their silencing factors is most
145    apparent in other animals[8,27]. For a comprehensive analysis of younger TE
146    populations in Lake Malawi, we created an additional curated TE annotation for AC,
147    which we used throughout this work alongside the uncurated annotation (**Figure S1A**,
148    see **Methods**).

150    We found that 515-746 (86-93%) cichlid TE families are expressed in gonads (**Figure
151    S1B**). Two trends are recognisable when considering the expression of TE families
152    grouped by class. First, long terminal repeat (LTR) families have the highest median
153    expression (**Figure S1C**). This trend is reversed when TE expression is quantified
154    based on the curated TE annotation of AC, which has more annotated LTR families
155    (**Figure S1B**) and where LTR annotations were improved, including both the long-
156    terminal repeats and intervening genes. This suggests that uncurated LTR annotation
157    may lead to an overestimation of LTR expression. Second, TE families of the same
158    class tend to be more highly expressed in testes rather than ovary, revealing
159    differences in TE expression between sexes (**Figure S1C**). Higher median expression
160    of annotated protein-coding genes was also observed in cichlid testes (**Figure S1D**),
161    suggesting the sex-specific differences in TE expression may follow general sex-
162    specific differences in transcriptional output.

**Figure 1. Dynamic patterns of TE expression during cichlid early development.** (**A**) The East African Great Lakes and surrounding bodies of water, along with the species used in this study, each representative of a major lake. *Oreochromis niloticus* (Nile tilapia) is used as an outgroup to the radiations of the Great Lakes. For Lake Malawi, we use three species to address within-lake dynamics of TE expression and epigenetic silencing. *Astatotilapia calliptera* is a generalist omnivore, which inhabits shallow water environments in the lake and surrounding rivers and streams[4,42], while *Maylandia zebra* and *Tropheops* sp. 'mauve' are Mbuna rock-dwelling cichlids specialised in eating algae[4,42].; (**B**) Expression of TE families belonging to major TE classes throughout early development of *A. calliptera*, displayed as rlog normalised counts. TE Expression was calculated using the default (panel above) and curated annotations (panel below). (**C**) Enrichment of TE classes and superfamilies in particular developmental stages, according to clusters A-E of differentially expressed TEs as defined in Figure S1G. Only TE superfamilies significantly enriched/depleted in at least one developmental stage are depicted. Grey dots represent lack of significant enrichment. Analysis done as in Chang et al., 2022[13], using the curated TE annotation of AC. AB, *Astatotilapia burtoni*; AC, *Astatotilapia calliptera*; LM, Lake Malawi; LT, Lake Tanganyika; LV, Lake Victoria; MZ, *Maylandia zebra*; ON, *Oreochromis niloticus*; PN, *Pundamilia nyererei*; rlog, regularised log; TM, *Tropheops* sp. 'mauve'.

Embryogenesis and early development are periods known to display signs of TE transcriptional activity[8,10,13]. We therefore conducted bulk mRNA sequencing in early developmental stages of Lake Malawi cichlids and found that 91-94% of cichlid TE families are expressed during early development (**Figure S1E**). Expression in these developmental stages is overall identical between AC and TM (**Figure 1B** and **S1F**). The temporal expression pattern of all TE classes is similar: lower expression before gastrulation rising to peak or near peak expression at early gastrula followed by relatively constant levels of transcriptional activity. Analysing TE expression at the locus level reveals the overall expression pattern at the family level is not universal, as several individual TEs have expression patterns specific to distinct developmental stages (**Figure S1G**). Interestingly, we find a major enrichment of ERV1, Gypsy and Pao LTRs in gastrula stage and early somitogenesis (**Figure 1C**, cluster B), and SINE (short interspersed nuclear element) enrichment at the earliest stages (**Figure 1C**, cluster A). Overall, these results support substantial transcriptional TE activity in gonads and during early development of African cichlids.

**An expanded repertoire of *piwil1* genes in Lake Malawi cichlids**

Given the dynamic TE expression patterns observed, we reasoned that active silencing pathways must be in place in cichlids to counteract TE activity. First, we identified orthologs of sRNA-based TE silencing factors in cichlids (**Supplemental Table 1**)[8]. With three exceptions, all genes are present in cichlid genomes (**Figure S2A** and **Supplemental Table 1**).

187

188  Then, we addressed whether these factors are expressed in the germline by
189  performing quantitative proteomics on gonads of representative cichlid species. TE
190  silencing factors are detected most prominently in testes (**Figure S2B**). Abundant yolk
191  proteins, from the substantial yolk fraction of cichlid eggs[43], precluded protein
192  detection in ovary samples at a depth similar to other organs (**Figure S2C**). Despite
193  the influence of the yolk, Piwil1, a core piRNA pathway factor was detected in the
194  ovaries of all species (**Figure S2B**). Somatic roles for the piRNA pathway have been
195  increasingly recognized in animals, including in brain and nervous system[44]. We also
196  profiled the proteome of brain tissues of the representative cichlid species, but
197  obtained no consistent evidence supporting expression of core piRNA factors in the
198  brain of all cichlid species (**Figure S2B**). These results point to strong conservation of
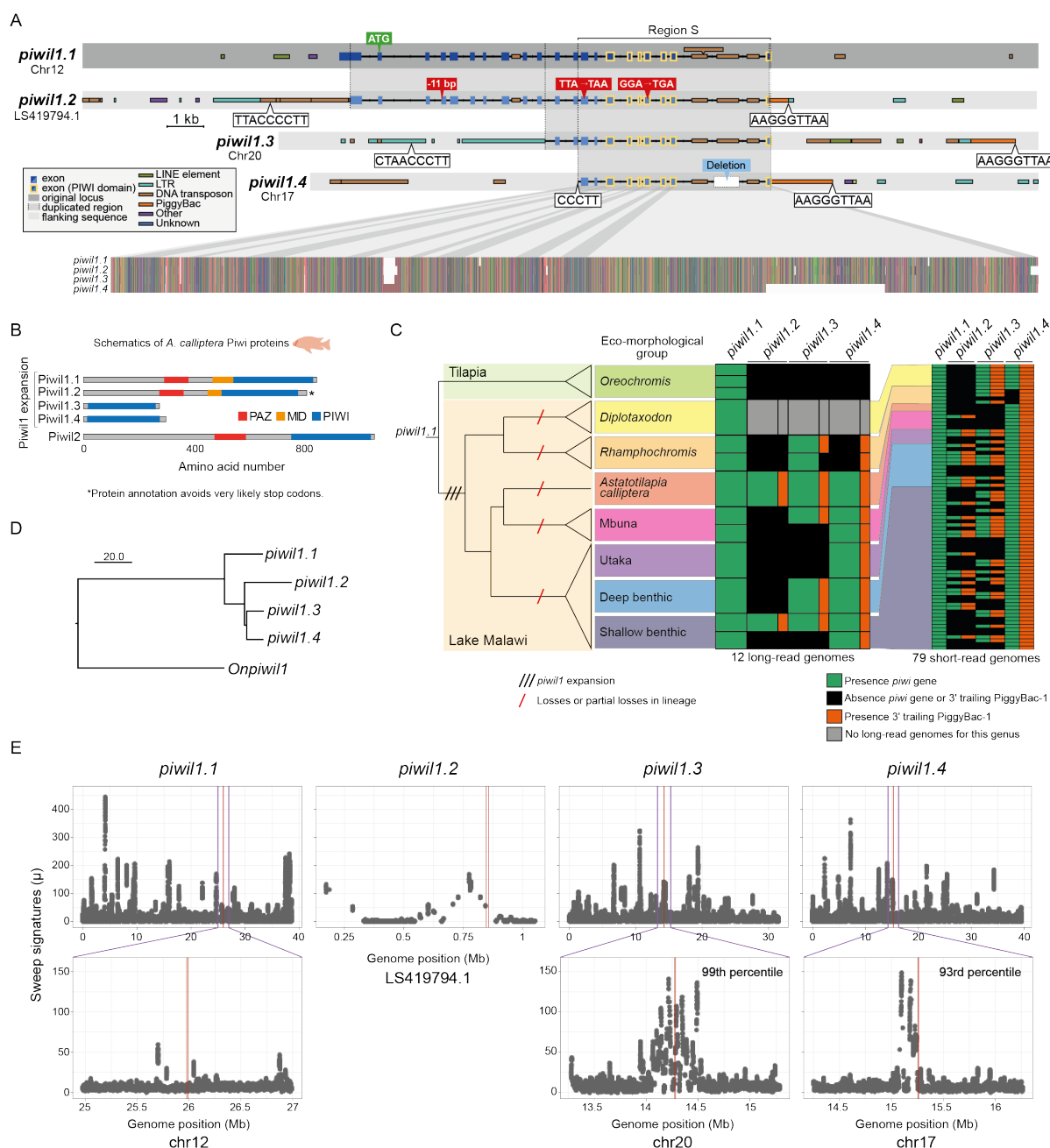199  germline-expressed TE silencing factors in African cichlids.

200

201  While inspecting TE silencing factor orthologs, we detected multiple copies of *piwil1*
202  genes, homologs of zebrafish *ziwi*[33], in cichlids representative of Lake Malawi, but not
203  in representatives of Lakes Tanganyika and Victoria (**Figure 2A-B**, **Supplemental**
204  **Table 1**. While fishes generally have one *piwil1* copy, AC has four *piwil1* copies, which
205  we named *piwil1.1-1.4*. Two of these are full-length copies, whereas the other two are
206  truncations containing only the PIWI domain (**Figure 2A-B**). *piwil1.1* of AC is located
207  in the conserved syntenic context of vertebrate *piwil1* genes (**Figure S3A**), indicating
208  that *piwil1.1* is the ancestral cichlid *piwil1* gene. By aligning all additional *piwil1* copies
209  of AC to the coding sequence of *piwil1.1* and projecting the coding sequence to the
210  aligned paralogs, we observe that the full-length paralog *piwil1.2* likely contains stop
211  codons that are bypassed in existing gene annotations produced by automated
212  annotation pipelines (**Figure 2A**). Also, *piwil1.2* is expressed at negligible levels in
213  cichlid gonads and brain (**Figure S3B**) and is therefore likely a pseudogene.

214

215  *piwil1.2*, *piwil1.3*, and *piwil1.4* reside in genomic regions rich in TEs (**Figure 2A**). The
216  3' regions of *piwil1.2, piwil1.3* and *piwil1.4* share a PiggyBac TE insertion (**Figures 2A**
217  and **S3C**). PiggyBac is a DNA TE family known to be very proficient at carrying large
218  DNA segments upon transposition, a quality that has promoted its use in genome
219  engineering[45,46]. Autonomous PiggyBac TEs consist of two terminal inverted repeats
220  (TIRs) flanking a transposase gene[47]. Like other DNA TEs with TIRs, PiggyBacs
221  mobilise when two transposase proteins each bind to one of the TIRs[6]. The *piwil1*-
222  associated PiggyBacs have mutations that preclude production of a functional
223  transposase (**Figure S3C**). These *piwil1*-associated PiggyBac belong to the same TE
224  family (PiggyBac-1), of which we identified 377 high quality copies in the AC reference
225  genome. Considering the genome size of >880 Megabase, one PiggyBac-1 element
226  is expected, on average, every 2.3 Megabase. A phylogeny of all high-confidence
227  PiggyBac-1 TE fragments in the AC genome shows that the three *piwil1*-associated
228  PiggyBac TEs are closely related, particularly the PiggyBacs associated with *piwil1.3*
229  and *piwil1.4* (**Figure S3D**). Finding all three *piwil1* paralogs on different chromosomes
230  with closely related flanking PiggyBac-1 insertions either directly 3' adjacent (*piwil1.2*
231  and *piwil1.4*) or 7 kb downstream (*piwil1.3*) is therefore highly unlikely to be
232  coincidental.

233

234  Given the presence of related PiggyBac-1 TEs associated with all three *piwil1*
235  paralogs, we reasoned that the initial expansion of *piwil1* genes in Lake Malawi cichlids
236  was likely driven by transposition of PiggyBac-1, either at a time when its transposase
237  was active, or in a non-autonomous fashion using the transposase of other PiggyBacs.
238  This could have happened if a piggyBac transposase used one of its own TIRs
239  together with an alternative TIR-like sequence from the *piwil1* locus. To address this,
240  we searched for sequence signatures of PiggyBac mobilisation: the preferred insertion

**Figure 2. An expansion of *piwil1* paralogs in Lake Malawi cichlids likely mediated by PiggyBac TEs.** (**A**) Detailed schematics of the four *piwil1* loci in the *A. calliptera* reference genome. Exons and TEs are shown, along with other relevant sequence features, such as start and stop codons, deletions, etc. The sequences of the putative PiggyBac TIRs (terminal inverted repeats) and preferred insertion sites are shown in white boxes, from 5' to 3'. Of note, the putative TIR and insertion site sequences distal to the PiggyBac are the reverse complement of 5'-CCCTT-3' and 5'-TTAA-3', respectively. The dotted lines represent the borders of duplicated regions, according to multiple sequence alignment. Region S marks the genomic region shared by all *piwil1* genes. The image in the lower portion of the panel is a zoomed-out image of the multiple sequence alignment, color-coded by nucleotide. The putative stop codons were identified manually from an alignment of the genomic regions of all *piwil1* copies with the coding sequence of *piwil1.1*, the *piwil1* gene most conserved in vertebrates. No putative stop codons were found in *piwil1.3* and *piwil1.4*. (**B**) Schematics of the domain structure of the five Piwi proteins annotated in the *A. calliptera* genome, including the expanded Piwil1 protein repertoire. Due to the putative stop codons found in the *piwil1.2* locus, it is likely that the protein is misannotated and that the full-length protein will not be produced. (**C**) Presence (green)/absence (black) of each *piwil1* gene in genomes of Lake Malawi and Tilapia cichlids. Presence of *piwi*-associated PiggyBac TEs is indicated in orange. Presence/absence of *piwil1* genes and PiggyBac TEs was ascertained from long-read sequencing of 12 individuals and short-read sequencing of 79 individuals spanning all the major eco-morphological clades in Lake Malawi. The cladogram of the Malawi radiation reflects the current understanding of the radiation based on genomic studies[4]. The proposed model for *piwil1* gene evolution involves gene expansion early in the Lake Malawi radiation, followed by losses in particular lineages. (**D**) Neighbour-joining tree representing the Hamming distance between the non-coding regions of the *piwil1* genomic sequences of *A. calliptera* along with the genomic sequence of *piwil1* of *O. niloticus* (*Onpiwil1*) as an outgroup. The multiple sequence alignment used to build this tree included the introns shared by all *piwil1* genes (Region S in **Figure 2A**). (**E**) The plots show genome-wide results of Raised Accuracy in Sweep Detection (RAiSD)[84]. μ is a metric incorporating three selective sweep signatures, with higher μ values indicative of a stronger signature of selection. Upper panels show μ across the entire chromosome, or entire scaffold in case of *piwil1.2*. Lower panels are insets of the *piwil1* gene regions +/- 1 Megabase (Mb). As the entire scaffold where *piwil1.2* resides is less than 2 Mb, no inset is shown. We calculated a per-gene μ for all genes (see **Methods**), and with this approach *piwil1.3* and *piwil1.4* are in the 99[th] and 93[rd] percentile, respectively, of per-gene μ.

241

242  sequence (5'-TTAA-3'), directly preceding the predicted PiggyBac TIR sequence (5'-
243  CCCTT-3')[47,48]. We found potential TIRs adjacent to the PiggyBac-1 elements, and
244  close to the border of the *piwi* duplications distal to the PiggyBac (**Figure 2A**). Putative
245  PiggyBac insertion signatures distal to the PiggyBac-1 element of *piwil1.2* and *piwil1.3*
246  were harder to identify because of additional transposition in that area that could have
247  pushed the PiggyBac sequence signature further upstream from *piwil1* (**Figure 2A**).
248  We could find a 5'-CCCTT-3' sequence upstream of *piwil1.4*, but the downstream
249  5'-TTAA-3' insertion sequence may have eroded. The consistent association of closely
250  related PiggyBac TEs to *piwil1* paralogs, and the presence of putative TIR sequences
251  flanking the genes are compatible with a model whereby PiggyBac-1 transposition
252  mediated the expansion of *piwil1* genes in Lake Malawi cichlids.
253
254  **Evolution and functional potential of *piwil1* genes in Lake Malawi cichlids**
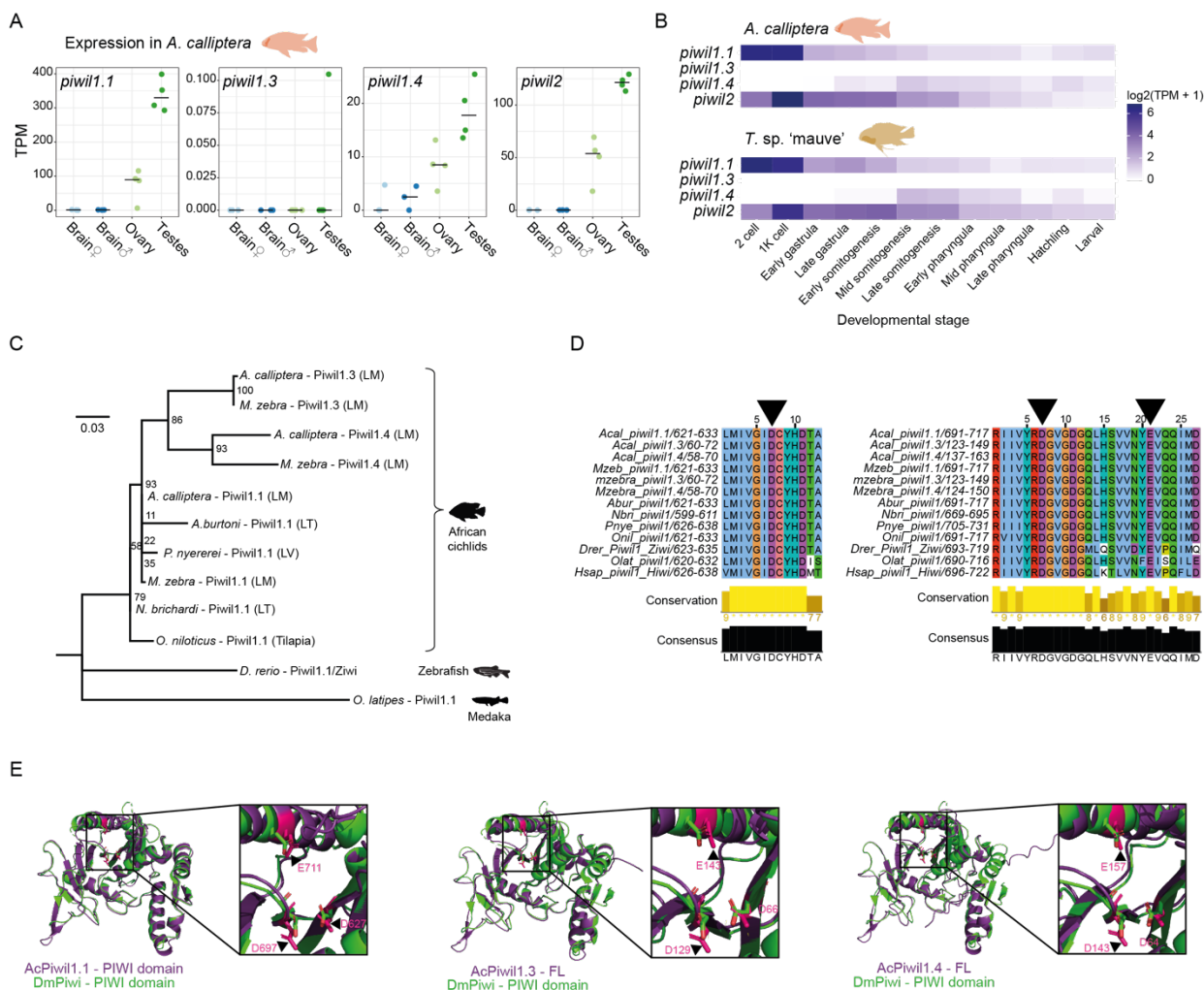255
256  Next, we assessed the prevalence of each *piwil1* paralog in the major eco-
257  morphological clades of Lake Malawi cichlids. We mapped genomic reads to the AC
258  reference genome (which contains all four *piwil1* copies) and manually assessed the
259  presence or absence of each *piwil1* gene from mapped reads. We used 12 sets of
260  long reads and 79 sets of short reads of Lake Malawi cichlids, corresponding to 80
261  species (**Supplemental Table 1**). We did not find any of the three extra *piwil1* paralogs
262  in tilapias, which form an outgroup to the haplochromine radiations (**Figure 2C**).
263  However, we find each additional *piwil1* paralog in all major eco-morphological clades
264  within the Lake Malawi radiation (**Figure 2C**). *piwil1.1* and *piwil1.4* are most
265  widespread, with *piwil1.1* identified in all individuals and *piwil1.4* found in 82/88
266  individuals (exceptions are 6/7 individuals of the *Rhamphochromis* genus, **Figure 2C**
267  and **Supplemental Table 1**). Conversely, *piwil1.2* and *piwil1.3* have a patchier
268  distribution (27/88 and 46/88 individuals). We found support for a 3' trailing PiggyBac-
269  1 TE in the vast majority of *piwil1.2*, *piwil1.3*, and *piwil1.4* copies (153/155, **Figure 2C**
270  and **Supplemental Table 1**). In 8 individuals we found support for a 3' trailing
271  PiggyBac-1 TE in their expected location 3' of *piwil1.3* and *piwil1.4*, but found no
272  support for the *piwil1* gene itself (**Figure 2C** and **Supplemental Table 1**). This
273  observation may reflect rare events of *piwil1* gene elimination by recombination
274  processes. Alternatively, these individuals could be heterozygous for the
275  presence/absence of *piwil1.3* or *piwil1.4*, leading to fewer supporting genomic reads.
276
277  Inspection of alignments of all AC *piwil1* paralogs revealed that *piwil1.2*, *piwil1.3* and
278  *piwil1.4* all share variation that is not shared with *piwil1.1* (**Figure 2A**). Moreover,
279  *piwil1.3* and *piwil1.4* share the most variation. This, together with the relatedness of
280  the *piwil1*-associated PiggyBac-1 elements (**Figure S3D**), suggests that *piwil1.2* was
281  the first paralog to duplicate via transposition and that *piwil1.3* and *piwil1.4* originated
282  from *piwil1.2*. A tree representing the distance between the non-coding regions shared
283  by all four *piwil1* genes of AC (within region S in **Figure 2A**) and *piwil1* of ON (as an
284  outgroup) support this hypothesis (**Figure 2D**). A similar tree created from the exons
285  shared by these same *piwil1* genes (within region S in **Figure 2A**) did not produce a
286  tree topology congruent with the non-coding tree (compare **Figures 2D** with **S3E**). We
287  suggest that this discrepancy could reflect selective processes acting on the coding
288  sequences of *piwil1* genes.
289
290  Following gene duplication, paralogs can undergo a number of evolutionary routes,
291  including towards sub- or neofunctionalisation[49], with distinct signatures of selection.
292  To learn about the selective pressures at play, we tested for the presence of signatures
293  of selective sweeps in 79 Lake Malawi cichlid genomes (**Supplemental Table 1**).
294  While the genomic region of *piwil1.1* does not display a clear signature of selective
295  sweep (**Figure 2E**, left panels), *piwil1.3* and *piwil1.4* are in the 99th and 93rd

7

**Figure 3. Expression and functional potential of Piwil1 proteins in Lake Malawi.** (**A**) Expression, in Transcripts per Million (TPM), of *piwil1* paralogs and *piwil2* in gonads and brain of *A. calliptera*. (**B**) Expression of the three *piwil1* genes and *piwil2* throughout early development of *A. calliptera* and *Tropheops* sp. 'mauve', another Lake Malawi cichlid. (**C**) Phylogenetic tree constructed from an alignment of the PIWI domain of Piwil1 proteins of African cichlids, using zebrafish and medaka as outgroups. Branch support numbers are shown at the tree nodes and were calculated with 10,000 ultrafast bootstrap replicates. (**D**) Specific regions of the multiple sequence alignment of several PIWI domains, surrounding the integral residues of the catalytic triad, indicated with black arrowheads, the catalytic residues within the PIWI domain known to be important for Piwi-mediated cleavage. These residues are conserved in Piwil1 proteins of African cichlids, including in *piwil1.3* and *piwil1.4* in Lake Malawi. (**E**) Structural alignments of the PIWI domain of *Drosophila melanogaster* (Dm) Piwi protein and AlphaFold predictions of Piwil1.1 (using only PIWI domain, left), Piwil1.3 (full-length, centre), and Piwil1.4 (full-length, right) of *A. calliptera*. Regions of the structural alignment encompassing the catalytic triad are augmented in the insets and the triad residues are highlighted with black or white arrowheads.

296    percentiles, respectively, of genes with highest values of integrative sweep signatures,
297    supporting positive selection at these loci (**Figure 2E** and **Supplemental Table 1**).
298    Moreover, we found evidence of positive selection in cichlid Piwil1 proteins beyond
299    Lake Malawi, particularly in amino acid residues in the PIWI domain or immediately
300    C-terminally adjacent to the annotated domain (**Figure S4A**). The results above are in
301    line with positive selection acting on cichlid Piwi proteins, most notably in the expanded
302    Piwi repertoire of Lake Malawi cichlids. Overall, the data suggests a scenario
303    consistent with *piwil1* expansion early in the radiation, followed by positive selection
304    and gene losses.
305
306    Next, we sought to determine whether the expanded copies of *piwil1* genes in Lake
307    Malawi are expressed. We excluded *piwil1.2* from further analysis, because both the
308    premature stop codons in conserved exons (**Figure 2A**) and low expression (**Figure
309    S3B**), suggest that it is a pseudogene. First, we interrogated *piwil1* gene expression
310    at the mRNA level. We also probed the expression of *piwil2*, the *piwi* gene homolog of

311 zebrafish *zili*[33], which did not undergo gene duplication. *piwil1.1* and *piwil2* are strongly
312 expressed in gonads but not in brain (**Figures 3A** and **S4B**), in line with known TE
313 silencing roles in the germline of other organisms[8,27,33]. *piwil1.4* was expressed in
314 gonads, and lowly expressed in brain. During early development of Lake Malawi
315 cichlids, we detected strong maternal deposition of *piwil1.1* and *piwil2* transcripts
316 (**Figure 3B**). In contrast, *piwil1.4* seems to be expressed mainly after gastrulation,
317 likely after the onset of zygotic expression. No expression of *piwil1.3* was detected in
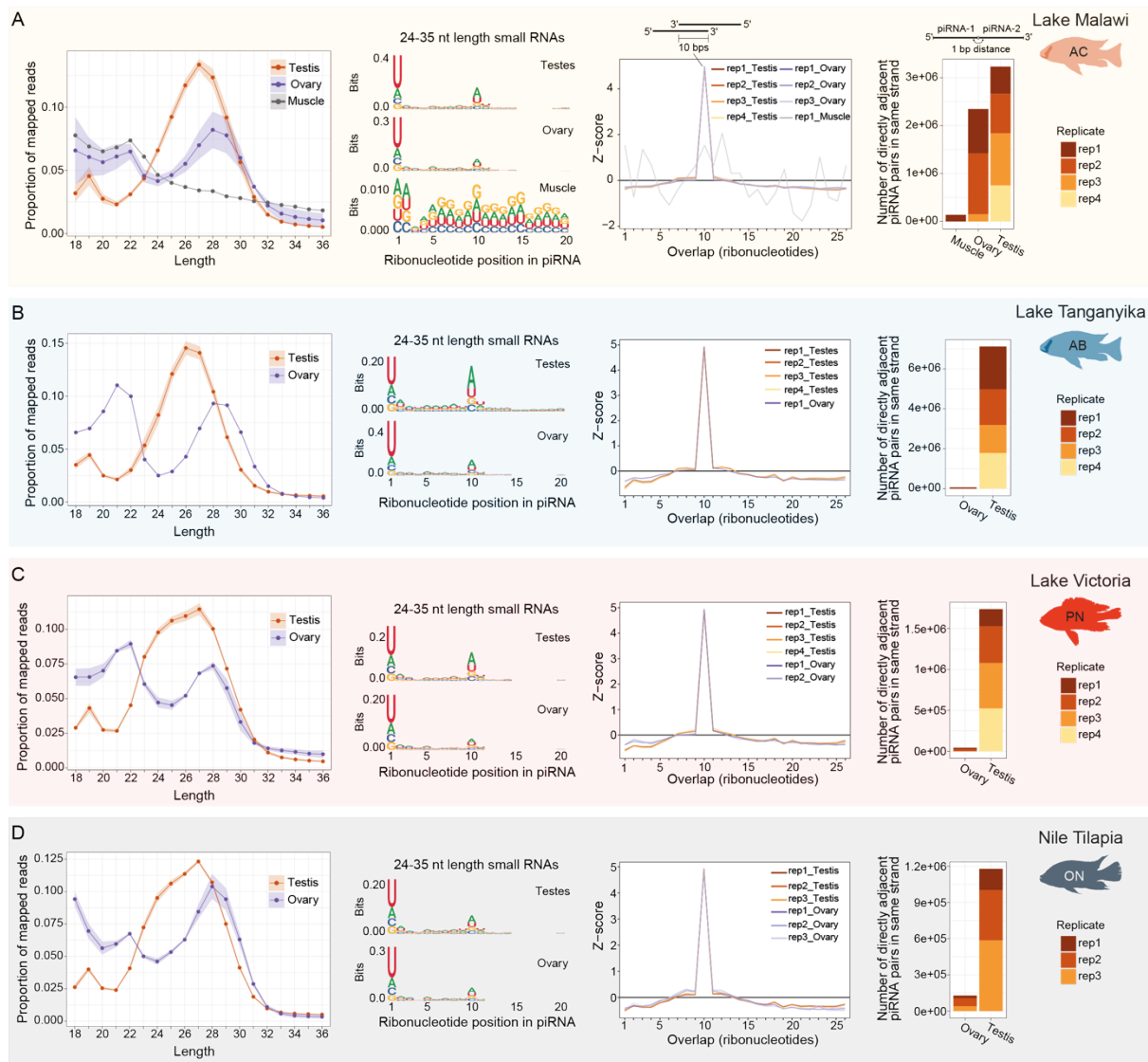318 these organs and in early development (**Figure 3A-B**).
319
320 To gain further insights into the potential function of these Piwil1 proteins, we analysed
321 their protein sequence and structure. As Piwil1.3 and Piwil1.4 have only the PIWI
322 domain (**Figure 2B**), we focused on the portion of Piwil1 proteins encompassing this
323 domain. We found low overall variation in African cichlid Piwil1 proteins, but the Lake
324 Malawi truncations showed higher divergence than their full-length orthologs (**Figure
325 3C**). This divergence is not expected to disrupt protein structure, as the predicted
326 structures of full-length Piwil1.3 and Piwil1.4 proteins align well with the known
327 structures of Piwi proteins of *Drosophila melanogaster* and *Bombyx mori*[50,51], and the
328 predicted PIWI domain of Piwil1.1 (**Figure S4C**). The PIWI domain is a ribonuclease
329 H-like domain, the catalytic centre of Argonaute proteins responsible for their slicer
330 activity. Within the PIWI domain, a DDE motif of amino acid residues is required for
331 Argonaute cleavage[52,53]. Despite the higher divergence of Piwil1.3 and Piwil1.4, they
332 retain a conserved DDE motif, as Piwil1.1 (**Figure 3D**). Furthermore, the PIWI domain
333 structures of Lake Malawi Piwil1 proteins are predicted to be identical to those of
334 *D. melanogaster* and *B. mori* Piwi proteins, including the relative position of the DDE
335 motif residues (**Figures 3E** and **S4D**). These data indicate the genomes of Lake
336 Malawi cichlids encode three Piwil1 proteins with potentially catalytically active PIWI
337 domains.
338
339 **Cichlids express TE-targeting piRNAs with signatures of active silencing**
340
341 To characterise the piRNA cofactors of cichlid Piwi proteins, we sequenced sRNAs
342 from gonads of the selected cichlid species (**Figure 1A**). The sRNA length distribution
343 profiles in gonads have prominent peaks at lengths of 21-22 nucleotides, likely
344 corresponding to microRNAs (**Figure S5A**). Contrary to microRNAs, piRNAs have
345 high sequence diversity[27]. When sRNA reads are collapsed into unique sequences,
346 we observed prominent sRNA populations between 24-31 nucleotides long, consistent
347 with the length distribution of piRNAs (**Figure 4**, left panels). In testes, sRNA
348 populations peaked at lengths of 26-27 nucleotides, whereas in ovaries the peak was
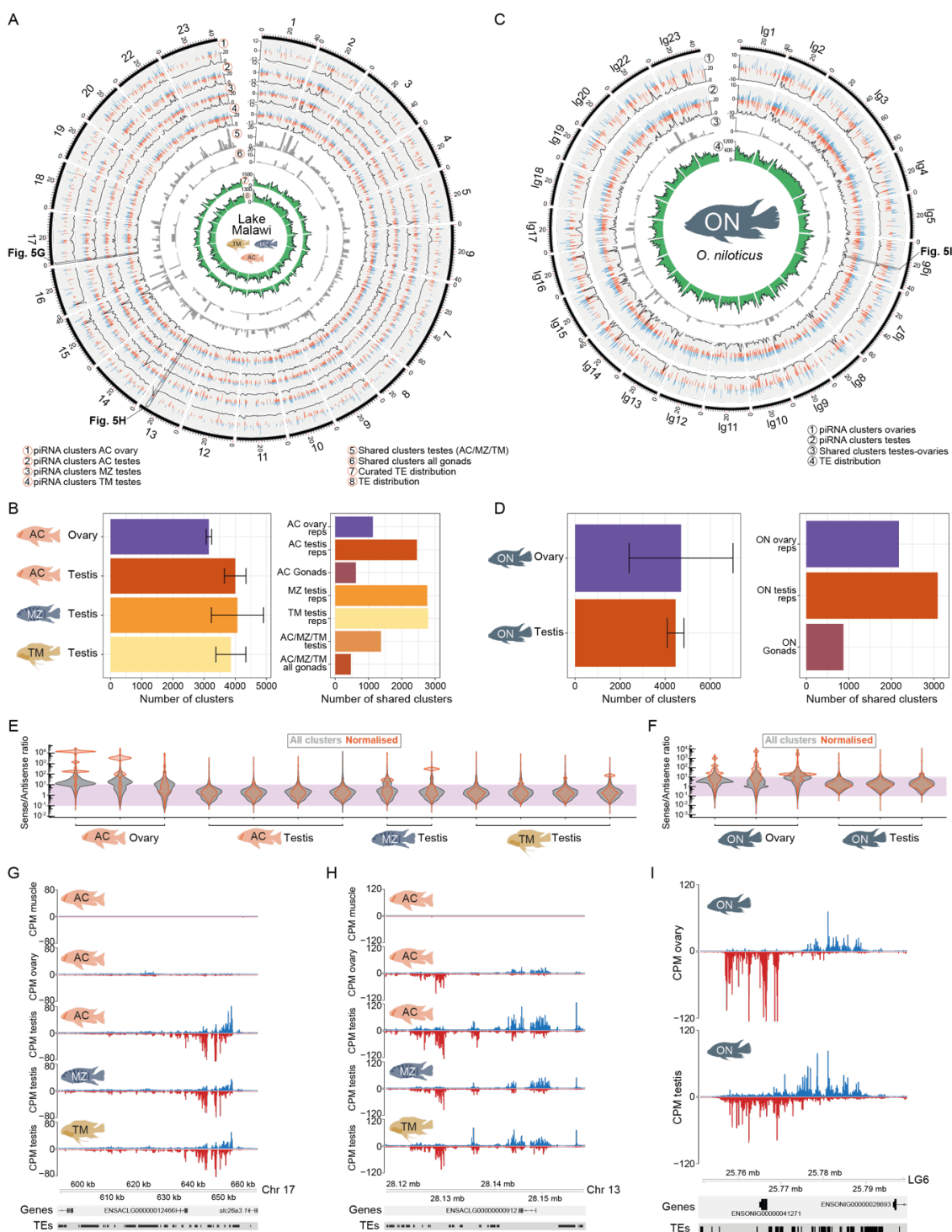349 shifted to 28-29 nucleotides long (**Figures 4** and **S5B**).
350
351 We selected sRNAs between 24-35 nucleotides long for subsequent analysis and
352 searched for the two typical sequence signatures of piRNAs: a bias for uridine at
353 position 1 (1U), and a bias for an adenine at position 10 (10A)[27,31,32,34,35,54]. Unique
354 sRNA sequences between 24-35 nucleotides long clearly show both the 1U and 10A
355 biases in cichlid gonads, as well as additional signatures consistent with active piRNA
356 ping-pong and phased biogenesis pathways (**Figures 4** and **S5B-G**). Of note, while
357 phased piRNA biogenesis is pervasive in cichlid testes, ovary sRNAs display no clear
358 signatures of phased biogenesis, except in AC ovaries (**Figures 4** and **S5F-G**). As a
359 control, we sequenced sRNAs from muscle, as a representative somatic tissue of AC
360 and found no prominent population of sRNAs in the piRNA length range with piRNA
361 signatures (**Figures 4A** and **S5E-F**). Thus, cichlid gonads express sRNA populations
362 consistent in length and sequence signatures with an active piRNA pathway.
363

9

**Figure 4. African cichlids express piRNAs in gonads.** (**A-D**) sRNA length distribution profiles and piRNA sequence signatures in sRNAs 24-35 nucleotides long. sRNA length profiles shown here (left-most panels) comprise only reads of unique sequence. The shading in the sRNA length distribution profiles indicates standard deviation of replicates. Sequence logos (second set of panels from left) denote the 1U bias typical of piRNAs, and the 10A signature of ping-pong amplification in gonad sRNAs but not in muscle tissues of *A. calliptera* (A). Third set of panels from left show ping-pong signature with a robust overlap of 10 ribonucleotides in piRNA pairs. Right hand-side panels show number of piRNA pairs in same orientation that are directly adjacent, indicative of phased piRNA biogenesis. Signature is observable in the testes of all species, but in ovaries it is detectable only in AC. AB, *Astatotilapia burtoni*; AC, *Astatotilapia calliptera*; CPM, Counts Per Million; ON, *Oreochromis niloticus*; PN, *Pundamilia nyererei*.

364   piRNAs are often created from discrete genomic regions termed piRNA
365   clusters[26,27,31,32]. To finely map piRNA clusters, we used a novel computational
366   approach that identifies piRNA clusters by incorporating information from uniquely-
367   and multi-mapping reads in a stepwise manner (see **Methods**). We restricted the
368   analysis to sRNA sequencing data of Lake Malawi (AC, TM, and MZ, all mapped to
369   the AC genome) and ON, because these chromosomal level assemblies allow us to
370   define the piRNA clusters within genomic coordinate systems and to understand their
371   biological context. We identified thousands of genomic sources of piRNAs in Lake
372   Malawi (**Figure 5A-B** and **Supplemental table 2**, between 3,091-3,251 in ovaries and
373   3,494-4,252 in testes) and ON (**Figure 5C-D** and **Supplemental table 2**, between
374   3,194-7,352 in ovaries and 4,053-4,781 in testes). The clusters explain 65-80% of
375   piRNA reads in the library (**Figure S6A**). Although the total number of clusters are
376   comparable in testes of distinct Lake Malawi species (**Figure 5B**, compare testis of
377   AC, MZ, and TM), the number of clusters shared between all three species are a
378   fraction of the total (**Figure 5B**, 1,377 shared clusters), revealing variation in piRNA
379   production in closely related species of Lake Malawi cichlids. Moreover, the even lower

**Figure 5. Fluid genomic origins of cichlid piRNAs**. (**A**) Circos plot showing the chromosomal locations of piRNA clusters in Lake Malawi cichlid gonads (tracks 1-4), clusters shared between all replicates of each organ (tracks 5-6), and TE distributions (tracks 7-8) from curated (track 7) and non-curated annotations (track 8). In tracks 1-4, blue and red represent the log2 mean Reads Per Kilobase Million (RPKM) of piRNA clusters in the plus and minus strands, respectively. In the bottom of tracks 1-4 is a line plot with the density of clusters. (**B**) Left panel shows the mean number of clusters identified in Lake Malawi cichlid gonads. Error bars represent standard deviation. Right panel depicts the number of clusters shared between the replicates of the organs indicated. (**C**) Circos plot showing the chromosomal locations of piRNA clusters in ovaries (track 1) and testes (track 2), the shared clusters between these two organs (track 3), and the TE distribution (track 4). In tracks 1-2, blue and red represent the log2 mean RPKM of piRNA clusters in the plus and minus strands, respectively. In the bottom of tracks 1-2 is a line plot with the density of clusters. (**D**) Left panel represents the mean number of clusters identified in *O. niloticus* gonads. Error bars represent standard deviation. Right panel shows the shared clusters between the replicates indicated. (**E-F**) Strand biases in piRNA production, shown as the ratio of sense over antisense piRNAs intersecting each piRNA cluster. The grey violin plot represents all piRNA clusters identified, while the orange violin plot represents the sense/antisense ratio normalised according to cluster productivity. The purple region highlights piRNA clusters with piRNA production less than 100-fold different between the sense and anti-sense strands. Thus, values that fall within this range likely account for piRNA clusters producing piRNAs from both strands. (**G-I**) Genome tracks with examples of clusters identified in Lake Malawi cichlids (G-H) and in *O. niloticus* (I). Blue and red tracks represent 24-35 nucleotide long piRNAs, in Counts per Million (CPM), mapping to the plus and minus strands, respectively. (G) shows a testes-specific piRNA cluster in Lake Malawi. (H-I) are examples of clusters shared by ovary and testis of Lake Malawi cichlids (H) and *O. niloticus* (I). AC, *Astatotilapia calliptera*; CPM, Counts Per Million; MZ, *Maylandia zebra*; ON, *Oreochromis niloticus*; TM, *Tropheops* sp. 'mauve'.
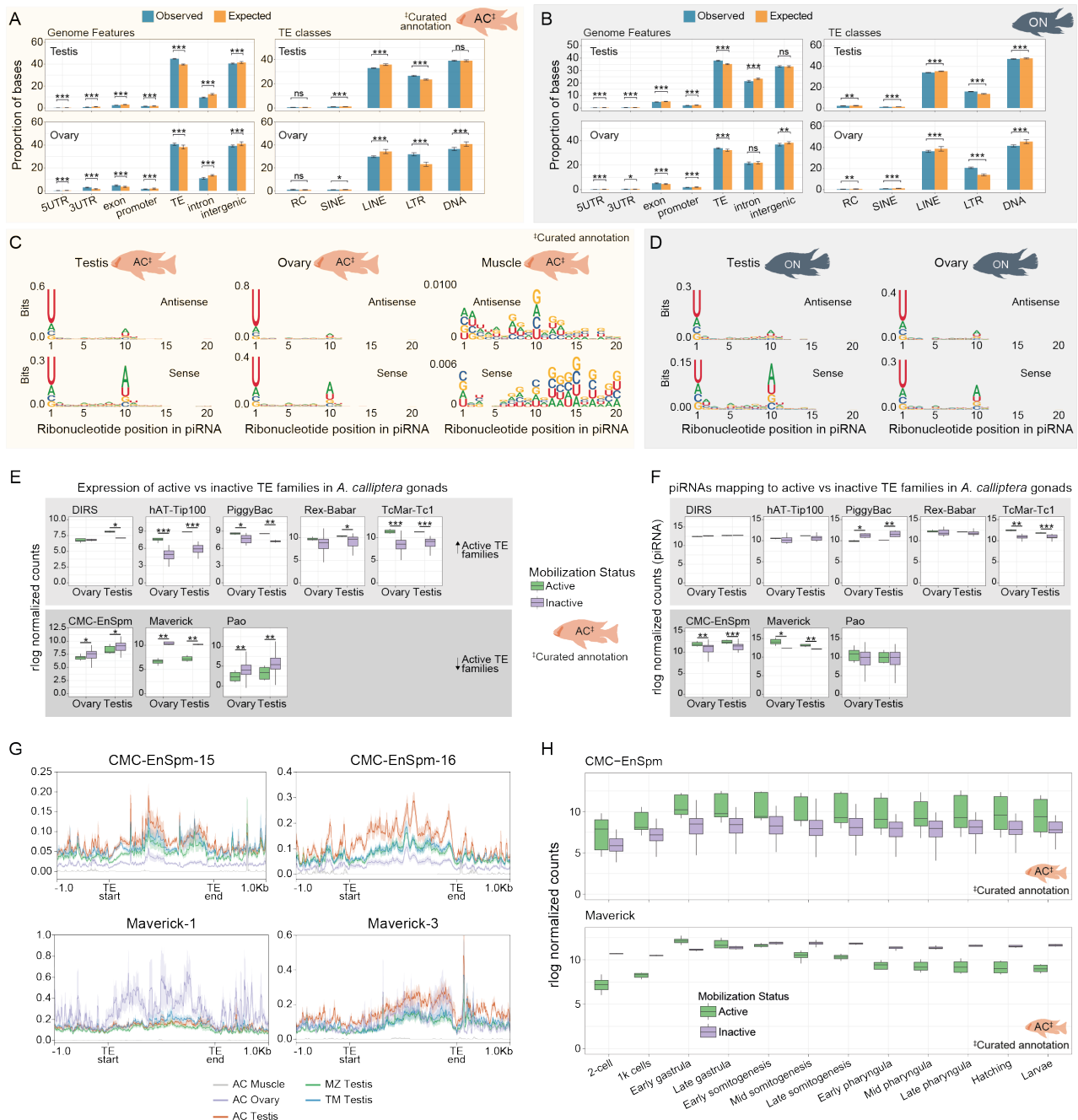
380

381 number of clusters shared between testes and ovaries illustrates sex differences in
382 piRNA production (**Figure 5B, D**, 622 shared clusters in AC gonads, 469 clusters
383 shared across all Lake Malawi testes and ovaries, and 872 clusters shared in ON
384 gonads). Overall, these results suggest considerable fluidity in the sources of piRNA
385 production in cichlids, including in cichlids inhabiting the same Lake.
386
387 Next, we explored additional features of the cichlid piRNA clusters identified. Most
388 piRNA clusters are shorter than 50 kb (**Figure S6B-C**). In testes, clusters tended to
389 be larger than in ovaries (**Figure S6B-C**, median length of 12.3 kb in AC testes versus
390 2.89 kb in AC ovaries, and median length of 13.70 in ON testes versus 4.62 kb in ON
391 ovaries). Within Lake Malawi, median cluster lengths in testes were consistent in AC,
392 MZ, and TM (**Figure S6B**). piRNA clusters are spread throughout the entire genome
393 and do not tend to be in close proximity (**Figures 5A, C** and **S6D**). piRNA clusters tend
394 to produce piRNAs from both strands, although in ovaries there is a bias for sense
395 piRNAs (**Figure 5E-F**). In terms of productivity, we found that a fraction of clusters
396 generate the majority of the piRNAs in the library (**Figure S6A**). We found no
397 relationship between the productivity and length of piRNA clusters (**Figure S6E**).
398 Examples of large, highly productive piRNA clusters are shown in **Figure 5G-I**.
399
400 The majority of piRNA clusters in AC and ON overlap with intergenic regions and TEs,
401 and the observed overlap with TEs is higher than expected by chance (**Figures 6A-B**
402 and **S7A**). Of all TEs, LTRs are significantly enriched in piRNA clusters (**Figures 6A-**
403 **B** and **S7A**). The piRNA clusters not detected in the testes of the three Lake Malawi
404 species tend to follow similar enrichment trends (**Figure S7B**). Furthermore, these
405 species-variable piRNA clusters are enriched in genomic regions of Lake Malawi
406 cichlids associated with structural variation (**Figure S7C**), defined in recent work[5]. This
407 suggests the existence of piRNA-producing sequences that are structural variants in
408 Lake Malawi cichlids. We overlapped all 24-35 nucleotide long piRNAs of AB and PN
409 with genome features and TE classes and observed enrichments similar to those of
410 AC and ON piRNA clusters (**Figure S7D-E**). Next, we further explored sequence
411 signatures of TE-mapping piRNAs. We found 1U bias and 10A bias in piRNAs
412 mapping sense and antisense to TEs, consistent with active targeting of TEs (**Figures**
413 **6C-D** and **S7F-G**). Sense piRNAs have higher 10A bias and lower 1U bias than
414 piRNAs antisense to TEs (**Figure 6C-D** and **S7F-G**). These signatures are absent from
415 muscle (**Figure 6C**). In terms of piRNAs mapping to distinct TE classes, we did not
416 find consistent differences between testes and ovaries across all species (**Figure**
417 **S7H**), contrasting with the sex differences in TE expression (**Figure S1C**).
418
419 We identified TE families likely to be transpositionally active in Lake Malawi cichlids
420 (P. Sierra & R. Durbin, unpublished results) and compared expression in gonads of
421 likely mobile TE families with transpositionally inactive families of the same
422 superfamily. We found higher expression of mobile versus immobile families in five TE
423 superfamilies in at least one organ (**Figure 6E**, upper panel). Conversely, in three TE
424 superfamilies we observed lower expression of the mobile versus immobile TE
425 families, in line with ongoing silencing (**Figure 6E**, lower panel). Two of these TE
426 superfamilies (CMC-EnSpm and Maverick) showed the opposite correlation in terms
427 of piRNA levels: mobile TE families were targeted by higher piRNA levels (**Figure 6F**,
428 lower panels), further supporting ongoing piRNA-driven TE silencing. No clear trend
429 was observed in relation to piRNA levels for mobile TE families more highly expressed
430 than their related immobile families (**Figure 6F**, upper panels). Additional data agree
431 with robust targeting of CMC-EnSpm and Maverick families by piRNAs (**Figures 6G**
432 and **S7I**). Furthermore, these families showed dynamic expression patterns in early
433 development: mobile families of CMC-EnSpm, which are DNA TEs, are more highly
434 expressed throughout early development than their immobile relatives (**Figure 6H**,

**Figure 6. Cichlid piRNAs target TEs**. (**A-B**) Observed and expected values at genomic features and TE classes that piRNA clusters overlap with in *A. calliptera* (A) and *O. niloticus* (B). ns, not statistically significant. In (A) we overlapped the piRNA clusters with TE features from the TE annotation produced with a curated library. (**C-D**) Sequence logos of 24-35 nucleotide long piRNAs mapping sense or antisense in regard to TE orientation in *A. calliptera* (C) and *O. niloticus* (D). The 1U and 10A signatures are observable in gonads but not in muscle. For *A. calliptera*, TE features were extracted from the curated TE annotation. (**E**) The mRNA expression in cichlid gonads of likely transpositionally active vs inactive families of a given TE superfamily. Panels above: superfamilies where the active families are more highly expressed than the inactive families. Panels below: TE superfamilies with higher expression of inactive families. P-values were calculated with Wilcoxon rank-sum tests (using Benjamini & Hochberg correction) comparing TE families with distinct mobilisation status in each gonad. The expression data was quantified using a curated annotation of Lake Malawi TEs. (**F**) piRNAs mapping to transpositionally active vs inactive TE families of the same TE superfamily. P-values were calculated with Wilcoxon rank-sum tests (using Benjamini & Hochberg correction) comparing TE families with distinct mobilisation status in each gonad. piRNAs were mapped to a curated TE annotation in Lake Malawi cichlids. Panels above and below, represent TE superfamilies with opposing relative expression of active versus inactive TE families, according to **Figure 6E**. (**G**) Metagene plots depicting mean piRNA levels mapping to all TEs of likely transpositionally active families. The shading represents standard error of replicates. TE start and TE end indicate start and end coordinates, respectively, of TE in the annotation. (**H**) The expression of transpositionally active versus inactive CMC-EnSpm (panel above) and Maverick (panel below) families throughout early development of *A. calliptera*. A curated annotation of Lake Malawi TEs was used to calculate expression data. AB, *Astatotilapia burtoni*; AC, *Astatotilapia calliptera*; ON, *Oreochromis niloticus*; PN, *Pundamilia nyererei*; sRNA, small RNA.

435 upper panel), a reverse pattern to that observed in gonads (**Figure 6E**, lower panel).
436 In turn, mobile TE families of the Maverick superfamily, also DNA TEs, show an

13

437  expression pattern in accordance with early silencing due to maternal silencing factors,
438  followed by weakening of silencing at gastrulation until zygotic silencing can be
439  re-established (**Figure 6H**, lower panel). Altogether, our data shows piRNAs target
440  TEs and are likely to be engaged in ongoing silencing of transpositionally active TE
441  families.

442
443

## Discussion

445
446  In this work, we describe three main findings that altogether suggest dynamic TE co-
447  evolution with host control mechanisms in East African cichlids: 1) dynamic TE
448  expression; 2) an expanded repertoire of *piwil1* genes; and 3) fast evolution of piRNA
449  clusters. We will elaborate on these points below.

450
451  First, hundreds of TEs families are dynamically expressed in gonads and early
452  development of cichlids (**Figures 1** and **S1**). Given the extensive shared
453  polymorphism in African cichlids due to hybridisation[3,4,18,55–58], we adopted a more
454  conservative approach by initially quantifying TE expression at the family level. We
455  found sex-biased expression patterns of TEs and protein-coding genes (**Figure S1C**),
456  with higher median expression in testes versus in ovaries. This asymmetry is likely to
457  be the result of overall higher transcriptional output in testes. Interestingly, this sex
458  asymmetry does not consistently extend to piRNAs mapping to TEs (**Figure S7H**),
459  suggesting that piRNA precursor transcription may not follow general transcription
460  trends. Testes contain the male germline and are a relevant organ in the context of
461  genetic conflict between TEs and host silencing factors[8,27]. Higher gene and TE
462  expression in testes is consistent with previous studies describing more widespread
463  expression and increased transcriptome complexity in mammalian testes[59,60]. TEs
464  have been shown to contribute to transcriptome complexity in the mammalian
465  germline[61] and in fish testes[62]. It may be worth exploring in depth whether cichlid
466  testes, much like mammalian testes, have increased transcriptome complexity and
467  diversity, and if this has contributed to the cichlid radiations. Indeed, gonad
468  transcriptomes are evolving faster than transcriptomes of other organs in Lake
469  Tanganyika cichlids[63]. In early development, we found that the majority of TE families
470  are most highly expressed during gastrulation, a period that may coincide with the
471  maternal-to-zygotic transition in cichlids. Zygotic transcription of TE silencing factors
472  may initiate concomitant to the onset of general zygotic transcription, leading to zygotic
473  TE silencing. Expression of transpositionally active Maverick TEs in early development
474  could illustrate just that (**Figure 6H**, lower panel). Expression analysis of individual TE
475  loci revealed TEs with expression in discrete developmental times (**Figures 1C** and
476  **S1G**). As expected given their evolutionary distance, the TE classes enriched in
477  particular developmental stages in cichlids differ substantially from those enriched in
478  the same developmental stages in zebrafish[13]. However, a striking similarity to
479  zebrafish is enrichment of TEs belonging to ERV1, Gypsy, and Pao LTR superfamilies
480  in gastrula stages (**Figure S1F-G**). It will be relevant to investigate how the maternal-
481  to-zygotic transition and/or epigenetic reprogramming affect LTR transcription and
482  transposition during early fish development.

483
484  Second, we find an expanded repertoire of *piwil1* genes in Lake Malawi cichlids and
485  signatures of positive selection on the novel copies (**Figure 2**). Lability in copy number
486  and positive selection on TE silencing factors are two signatures associated with arms
487  races between TEs and their animal hosts[24,25,38–41]. These findings also add to the
488  notion that piRNA pathway factors, including *piwil1* genes, evolve fast in teleosts[64].
489  Interestingly, TEs, the targets of Piwi proteins, likely have mediated, at least partially,
490  the expansion of *piwil1* genes in Lake Malawi cichlids. We found closely related

491 PiggyBac elements associated with the three novel *piwil1* genes, but not with the *piwil1*
492 copy sharing synteny with other vertebrate *piwil1* genes, presumably the original copy
493 (**Figures 2** and **S3** and **Supplemental Table 1**). We also found TIRs flanking the
494 PiggyBac and putative TIRs distal to the PiggyBac and 5' to the *piwil1* copies. The
495 non-coding differences of the four *piwil1* genes suggest the succession of events
496 underlying the expansion: first a duplication of *piwil1.1* creating *piwil1.2*, followed by
497 creation of one of the truncated copies from *piwil1.2*, and its subsequent duplication
498 (**Figure 2D**). Given the PiggyBac TIR signatures, it is likely that at least the first
499 duplication was mediated by transposition, but we cannot exclude that subsequent
500 duplications were driven by a recombination-based mechanism. By leveraging
501 available genomic resources we determined that *piwil1.1* and *piwil1.4* seem to be fixed
502 or nearly fixed in Lake Malawi, whereas *piwil1.2* and *piwil1.3* are less widespread
503 (**Figure 2C**). *piwil1.3* seems to have negligible expression in the germline and early
504 development (**Figures 3A-B**). It is possible that *piwil1.3* is expressed and functional in
505 other organs beyond the gonads and brain, or in juvenile developmental stages
506 between larval stage and sexual maturity. An alternative is that *piwil1.3* is a
507 pseudogene, similar to *piwil1.2*.
508
509 The exact function of *piwil1.3* and *piwil1.4* remains to be determined. Knock-outs of
510 *piwil1.1* and *piwil1.4* will be key to inform on their function. The annotated Piwil1.3 and
511 Piwil1.4 proteins are predicted to encode a catalytically competent PIWI domain
512 (**Figures 2A-B** and **3D-E**), the catalytic centre of Argonaute proteins responsible for
513 slicer activity[52,53]. The Argonaute domains lacking in Piwil1.3 and Piwil1.4, the MID
514 and PAZ domains (**Figure 2B**), are predicted to serve as binding pockets for the 5'
515 and 3' ends of the piRNA, respectively[65–67]. Without these domains, Piwil1.3 and
516 Piwil1.4 are most likely not able to bind to piRNAs or other sRNAs, and will probably
517 function independently of piRNAs. Thus, these truncated Piwi proteins were likely
518 repurposed for a piRNA-independent gene regulatory role, related, or not, to TE
519 silencing.
520
521 Third, we find fast evolution of piRNA clusters in cichlids. The majority of piRNAs were
522 produced from intergenic regions and TEs (**Figures 5,6**) and 65-80% of these
523 sequences can be grouped into discrete piRNA-producing clusters (**Figure S6A**). We
524 identify piRNA clusters with sex-biased expression, and, interestingly, variation in
525 piRNA clusters even in testes of closely related Lake Malawi cichlids (**Figure 5**). These
526 observations indicate that piRNA clusters are fast-evolving modules in Lake Malawi.
527 An in-depth population-wide analysis of piRNA populations and piRNA clusters in Lake
528 Malawi will be useful to determine just how rapidly these units are evolving in cichlids.
529 In terms of piRNA biogenesis, we find conserved differences in cichlid piRNA
530 populations with peaks at 26-27 nucleotides long piRNAs in testes versus 28-29
531 nucleotide long piRNAs in ovaries (**Figure 4**). These piRNA size differences may be
532 driven by Piwi Argonaute size preferences. The most striking difference in terms of
533 piRNA biogenesis however, is the lack of consistent phasing signature in the ovaries
534 of East African cichlids outside Lake Malawi (**Figures 4** and **S5**). It will be interesting
535 to determine the factor(s) inhibiting phased biogenesis in cichlid ovaries.
536
537 Three sets of observations point towards TEs as key genetic elements contributing to
538 cichlid diversification: 1) TEs represent a previously underestimated source of genetic
539 diversity in African cichlids[5]; 2) TEs have been linked with pigmentation and vision
540 traits, sex determination, and gene expression changes[18,20–23]; and 3) the ongoing
541 dynamic TE-host co-evolution and arms races that our findings suggest. It remains
542 unclear how the latter connects with cichlid phenotypic diversification. We expect it
543 does not come down to the number of TE families or the proportion of the genome
544 comprised by TEs. In this regard, zebrafish provides a much more striking example,

545 with nearly 2,000 distinct TE families, occupying more than 50% of its genome[6,13],
546 versus 557-828 TE families and 16-41% of the genome in cichlids (**Figure S1A-B**, **D**).
547 However, the *Danio* genus of zebrafish did not diversify nearly as prolifically as East
548 African cichlids despite its massive TE content.

550 What led to the unparalleled rates of phenotypic diversification observed in East
551 African cichlids? Recent work on the cichlid radiation of Lake Victoria suggests that
552 ecological versatility is the key[68,69]. Key features contributing to cichlid versatility
553 include strong sexual selection, highly plastic jaw structures, and abundant
554 interspecific hybridisation[1,68]. The regulatory consequences of hybridisation are one
555 possible avenue to pursue to study the influence of TE-host co-evolution in cichlid
556 radiations. Genomic studies have elucidated a complex evolutionary history of East
557 African cichlids, marked by substantial amounts of gene flow occurring through
558 hybridisation[4,18,55–58]. It will be important to determine how interspecific cichlid hybrids
559 tolerate regulatory mismatches driven by genetic conflict between TEs and the piRNA
560 pathway. It is conceivable that bouts of TE expansion following hybridisation
561 generated the (epi)genetic potential for the radiations. This study provides a platform
562 to investigate this hypothesis, an initial understanding of TEs and piRNAs as two
563 co-evolving modules. Going forward, learning about the co-evolution of these modules
564 in the context of recurring hybridisation has the potential to give valuable insights into
565 the genetic and molecular basis of the cichlid radiations.

## Methods

### Animal sampling and housing conditions

572 *Astatotilapia calliptera* and *Tropheops* sp. 'mauve' animals were grown in 220 Litre
573 tanks, with pH 8, at approximately 28ºC, and with a 12 h dark/light cycle. Males and
574 females of each species were housed only with conspecifics. Feeding, housing, and
575 handling were conducted in strict adherence to local regulations and with the protocols
576 listed in Home Office project license PP9587325. Fish were fed twice a day with cichlid
577 flakes and pellets (Vitalis). Tank environment was enriched with plastic plants, plastic
578 hiding tubes, and sand substrate. Aquaria grown animals were euthanised with
579 approved Home Office schedule 1 protocols, namely using 1 g/L MS-222 (Ethyl 3-
580 aminobenzoate methanesulfonate, Merck #E10521) and subsequent exsanguination
581 by cutting the gill arches, in accordance with local regulations. Afterwards, gonads,
582 brain and dorsal muscle tissue were carefully dissected, swiftly snap frozen in dry ice
583 and stored at approximately -80ºC.

585 Dominant adult male *Maylandia zebra* bred and raised in captivity were obtained from
586 commercial supplier Kevs Rifts and culled in Cambridge animal facilities, following an
587 ethically approved post–transport adjustment period. *M. zebra* animals were
588 euthanised using approved Home Office schedule 1 protocols as above. *Pundamilia*
589 *nyererei* animals were raised in stock tanks of dimensions 59 cm(L) x 45 cm(B) x 39
590 cm(H) and moved to larger tanks 177cm (L) x 45cm(B) x 39cm(H) once they reached
591 approx. 7 cm long. Temperatures were kept at 26Cº, with constant daily water change
592 of about 10% and 12:12 light dark regime. Frozen tissue samples of *Astatotilapia*
593 *burtoni* were provided by Hans Hofmann and Caitlin Friesen (University of Texas at
594 Austin, Austin, TX, USA). *Oreochromis niloticus* frozen tissue samples were provided
595 by David Penman, Alastair McPhee, and James F. Turnbull (Institute of Aquaculture,
596 University of Stirling, Stirling, Scotland, UK).

### Orthology analysis

599

600 To identify orthologs of conserved factors involved in TE silencing pathways, we used
601 OrthoFinder[70,71] v2.3.12. We used Ensembl proteomes (downloaded on 02/06/2020)
602 of Homo sapiens (GRCh38), *Mus musculus* (GRCm38), *Oryzias latipes*
603 (ASM223467v1), *Danio rerio* (GRCz11), *Takifugu rubripes* (fTakRub1.2),
604 *Gasterosteus aculeatus* (BROADS1), *Amphilophus citrinellus* (Midas_v5),
605 *Oreochromis aureus* (ASM587006v1), *Oreochromis niloticus*
606 (O_niloticus_UMD_NMBU), *Astatotilapia burtoni* (AstBur1.0), *Neolamprologus*
607 *brichardi* (NeoBri1.0), *Pundamilia nyererei* (PunNye1.0), *Astatotilapia calliptera*
608 (fAstCal1.2) and *Maylandia zebra* (M_zebra_UMD2a). OrthoFinder was run on
609 proteomes containing the longest protein isoform, parsed using a script provided with
610 OrthoFinder
611 (https://github.com/davidemms/OrthoFinder/blob/master/tools/primary_transcript.py).
612 Initially, we ran OrthoFinder with the fish genomes above as inputs (except *M. zebra*),
613 using option -f. Afterwards, we added human, mouse, and an additional Lake Malawi
614 cichlid species *M. zebra* to this analysis using options -b and -f. We subsequently
615 pinpointed the orthogroups containing known human, mouse and zebrafish TE
616 silencing factors and extracted the gene IDs of their cichlid orthologs.

617

618 **Piwil1 evolutionary analysis**

619

620 Piwil1 protein orthologs were identified with OrthoFinder (see **Orthology analysis**
621 above). Schematic of domain structure of Piwil1 proteins was plotted in R[72], with
622 packages drawProteins[73] and tidyverse[74]. Coordinates of the MID domain were
623 manually added to Piwil1 proteins, as this information was not present in Uniprot,
624 which drawProteins relies on. MID domain coordinates in *A. calliptera* Piwil1 proteins
625 were inferred from the MID domain coordinates of zebrafish Ziwi in Uniprot, through a
626 multiple sequence alignment of *A. calliptera* Piwil1 proteins and Ziwi.

627

628 To determine the presence and absence of *piwil1* copies and their 3' trailing PiggyBac-
629 1 TEs across Lake Malawi cichlid eco-morphological groups and genera, we probed
630 the reads of 74 previously published short-read genomes[4], 5 new short-read genomes,
631 as well as 12 long-read genomes (**Supplemental table 1**). Short-read genomes were
632 aligned to the *A. calliptera* reference genome (fAstCal1.2, GCA_900246225.3) using
633 bwa mem v0.7.17-r1188 (arguments: -C -p) using default settings[75]. Using samtools
634 v1.9[76], the resulting alignment files were then further processed with fixmate
635 (arguments: -m), sort (arguments: -l0) and mardup. Long-read genomes were aligned
636 to the same reference using minimap2 v2.17-r974-dirty[77] (arguments: -ax
637 map-pb --MD) and then sorted and indexed using samtools v 1.16-9-g99f3988. We
638 manually checked whether read alignments showed robust support in specific eco-
639 morphological groups/genera for the presence of each *piwil1* paralog and 3' trailing
640 piggyBac copy using IGV v2.9.4[78]. Next, we manually determined the exact features
641 of these regions using the *piwil1* gene annotations of fAstCal1.2[79], our TE annotation
642 created from a curated TE library (see section Transposable element annotations),
643 and genomic alignments of the entire regions encompassing all *piwil1* paralogs. Initial
644 alignments of the paralog loci were generated by aligning the fAstCal1.2 reference
645 genome to itself using Winnowmap2[80] (options: -ax asm5 --MD). Potential stop codons
646 in *piwil1* paralogs were assessed in a multiple sequence alignment between *piwil1.1*,
647 *piwil1.2* (reverse complement), *piwil1.3*, and *piwil1.4* (reverse complement) genomic
648 regions, which was created MUSCLE v3.8.31[81] using default settings and then curated
649 manually in AliView v1.27[82]. The exons of ENSACLT00000021959, the canonical
650 ENSEMBL isoform of *piwil1.1*, the best evolutionarily conserved *piwil1* gene, was
651 projected to the aligned sequences of the paralogs. A second alignment was created
652 analogously, which additionally included the homologous *piwil1* sequence from

653 *Oreochromis niloticus*. Based on the latter alignment, we calculated Hamming
654 distances (github.com/ssciwr/hammingdist) separately for intronic and exonic regions
655 and built neighbour joining trees (github.com/scikit-bio/scikit-bio). Alignment files can
656 be found at https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024. The *A.*
657 *calliptera* TE annotation created from a curated TE library was used to identify the
658 PiggyBac-1 TE and its terminal inverted repeats. The terminal inverted repeat
659 sequence of PiggyBac TEs (5'-CCCTT-3') was extracted from
660 https://dfam.org/classification/dna-termini[47].

661

662 For the selection analysis we restricted our existing callset of more than 2,000 whole-
663 genome sequenced Lake Malawi cichlids
664 (github.com/tplinderoth/cichlids/tree/master/callset), which are all aligned against the
665 chromosome level fAstCal1.2 reference genome[79], to the 79 individuals used in **Figure**
666 **2C** (**Supplemental Table 1**). The subset was generated with bcftools view[83], v.1.16-
667 9-g99f3988) (arguments --types snps -m2 -M2 -f PASS -S $sample_list) to retain
668 exclusively biallelic SNPs that passed all filters. Chromosome-scale VCFs along with
669 the four largest contigs (> 1 Mbp) were concatenated into a single VCF using bcftools
670 concat and served as the input for the selection analysis. A selection scan was
671 performed using Raised Accuracy in Sweep Detection (RAiSD) v2.9[84] (with
672 arguments: -f -M 3 -y 2 -m 0 -R -I). PiggyBac-1 sequences adjacent to *piwil1* genes
673 were extracted according to their annotation coordinates, and aligned with the
674 PiggyBac-1 family consensus from the curated TE library using MAFFT v7.475[85] with
675 option --auto. L-INS-i was the alignment method automatically selected. Alignment
676 visualisation was optimised in Jalview v2.11.2.7[86]. To expand the analysis, we
677 extracted all the PiggyBac-1 sequences annotated in the *A. calliptera* reference
678 genome (according to the curated TE annotation) with SWscore > 1000, and aligned
679 them with MUSCLE v3.8.31[81]. We further filtered the alignment to contain only the
680 region encompassed by the PiggyBac1 elements associated with *piwil1.2*, *piwil1.3*,
681 and *piwil1.4*, and removed alignment columns consisting almost exclusively of missing
682 data. A phylogenetic tree was constructed with IQ-TREE v2.1.2[87], option -B 1000.
683 TPM2+F+R2 was the best fit model. Trees were visualised and annotated in FigTree
684 v1.4.4 (https://github.com/rambaut/figtree).

685

686 The sequences of Piwil1 protein orthologs were collected from Ensembl. For *Piwil1*
687 genes encoding more than one protein isoform, the longest isoform was chosen for
688 analysis. As *A. calliptera piwil1.2* may be a pseudogene, we did not include its
689 predicted protein sequence in the subsequent analysis. Fish Piwil1 proteins were
690 aligned with MAFFT v7.475[85], using option --auto, and L-INS-i was the alignment
691 method automatically selected. We trimmed the alignment manually, keeping only 296
692 sites corresponding to the C-terminal region of the proteins with excellent alignment
693 score, which includes the PIWI domain. Original protein sequences and alignment files
694 can be found at https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024. IQ-
695 TREE v2.1.2[87] was used to construct phylogenetic trees from these two alignments
696 with options -B 10000 -o {medaka and zebrafish Piwil1 proteins were defined as
697 outgroups}. -B parameter refers to ultrafast bootstrap approximation[88]. PMB+G4 was
698 the best fit model. To test for selection, we redid the alignment using a smaller subset
699 of the proteins, including only Piwil1 proteins of African cichlids. Alignment of protein
700 sequences was performed with MAFFT v7.475[85], using option --auto. L-INS-i was the
701 chosen alignment model. Next, we used pal2nal v14[89] to produce a reverse alignment
702 from an alignment of the protein sequences to an alignment of the coding sequences.
703 The resulting reverse alignment was used as input for selection tests in Datamonkey[90].
704 A gene-wide test was first performed using Branch-site Unrestricted Statistical Test
705 for Episodic Diversification (BUSTED)[91]. We conducted the test in two ways, testing
706 for selection across all branches and testing for selection only in radiating cichlids, with

*O. niloticus* as an outgroup. BUSTED reported very strong support for positive selection in both cases (p-value = 0). **Figure S4A** shows the subsequent analysis to identify residues very likely to be under positive selection according to Mixed Effects Model of Evolution (MEME)[92].

To pinpoint catalytic residues of cichlid Piwil1 proteins, we first added the sequence of human PIWIL1 (HIWI) to the list of fish Piwil1 proteins used in the alignments above, and redid the alignment using MAFFT v7.475[85] with option --auto (L-INS-i was the model automatically chosen). The alignment was visualised in Jalview v2.11.2.7[86] and the catalytic residues were manually pinpointed based on their known positions in HIWI[52,53]. Structural alignments were performed with open-source PyMOL v2.5.0 using the align command. We aligned AlphaFold predictions of Piwil1.1 (Uniprot ID A0A3P8PWP0), Piwil1.3 (Uniprot ID A0A3P8NS09), and Piwil1.4 (Uniprot ID A0A3P8NRZ4) of *A. calliptera*, downloaded from AlphaFold Protein Structure Database[93,94], with crystal structures of *bombyx mori* Siwi (PDB ID 5GUH)[50] and *Drosophila melanogaster* Piwi (PDB ID 6KR6)[51]. As we focus on the PIWI domain, we aligned only the PIWI domains of *A. calliptera* Piwil1.1 (residues 550-856), *D. melanogaster* Piwi (residues 537-843), and *B. mori* Siwi (residues 593-899). As Piwil1.3 and Piwil1.4 of *A. calliptera* are truncations encompassing only the Piwi domain, we used their full-length structure for the alignments.

**RNA extractions**

Frozen brain, muscle, and gonad tissues were partitioned on a mortar positioned on dry ice, quickly to avoid thawing, and weighed. Biological replicates were created by collecting a similar mass of the same organ/tissue from size-matched individuals of the same species. 15-30 mgs of brain tissue, 26 mg of dorsal muscle tissue, and 14-144 mgs of gonad tissue were used, according to the specific tissue, tissue availability, and size of the specimen, which varied per species. Tissue pieces were transferred to BeadBug tubes prefilled with 0.5 mm Zirconium beads (Merck, #Z763772) and 500-600 µl of TRIzol (Life Technologies, #15596026) was added to the tubes and mixed vigorously. Afterwards, we conducted the homogenisation using a BeadBug microtube homogeniser (Sigma, #Z764140) at approximately 4°C (in cold room). Each sample was homogenised with five BeadBug runs at maximum speed (4,000 rpm) for 60 seconds each. No sample was run on BeadBug more than two consecutive times to avoid overheating. Other than the run time inside the BeadBug, samples were left on ice. After homogenisation, lysates were centrifuged for 5 minutes at 18,000 G at 4°C. Supernatant was then removed into a clean 1.5 mL tube. Centrifuged the lysates again, this time at maximum speed (approximately 21,000 G) for 5 minutes at 4°C. Transferred supernatant into a clean tube without disturbing the pellet and tissue debris. Mixed supernatant thoroughly 1:1 with 100% ethanol, pipetted the mix into a column provided in the Direct-zol RNA Miniprep Plus kit (Zymo Research, #R2072) and followed manufacturer's instructions, using the recommended in-column DNase I treatment.

**Library preparation and sequencing**

*mRNA sequencing.* Library preparation (directional, with poly-A enrichment) and sequencing (Illumina, PE150) of *A. calliptera*, *M. zebra*, *T.* sp. 'mauve', *A. burtoni*, and *O. niloticus* gonads was performed by Novogene. Libraries of *P. nyererei* gonads and *A. calliptera* brain tissues were prepared and sequenced as follows. Initial quality control was done using a Qubit Fluorometer (Invitrogen) and Qubit RNA HS Assay Kit (Invitrogen, #Q32855), and Agilent RNA TapeStation reagents (Agilent, #5067-5576; #5067-5577; #5067-5578). 50-250 ng of total RNA were used for library production

761 with the NEBNext® Poly(A) mRNA Magnetic Isolation Module (NEB, #E7490), in
762 conjunction with the NEBNext® Ultra™ II Directional RNA Library Prep Kit for
763 Illumina® (NEB, #E7760) and the NEBNext® Multiplex Oligos for Illumina® (96 Unique
764 Dual Index Primer Pairs, NEB #E6440). Quality control of the libraries was done with
765 the Qubit dsDNA HS Assay Kit (Invitrogen, #Q32854) and Agilent DNA 5000
766 TapeStation reagents (Agilent, #5067-5588; #5067-5589). Samples were then pooled
767 in equimolar amounts according to the TapeStation results and sequenced on a
768 NovaSeq 6000 system (PE150 on one lane of an S1 Flowcell).

769

770 *Small RNA* sequencing. Initial quality control was conducted using a Qubit
771 Fluorometer (Invitrogen) and the Qubit RNA HS Assay Kit (Invitrogen, #Q32855), and
772 Agilent RNA TapeStation reagents (Agilent, #5067-5576; #5067-5577; #5067-5578).
773 Samples were processed according to the NEXTFLEX® Small RNA-Seq Kit v4 with
774 UDIs (PerkinElmer, #NOVA-5132-32) with a 1 µg starting input and 12 cycles of PCR.
775 Quality control of the libraries was done with Qubit dsDNA HS Assay Kit (Invitrogen,
776 #Q32854) and Agilent DNA 5000 TapeStation reagents (Agilent, #5067-5588; #5067-
777 5589). Samples were then pooled in equimolar amounts according to the TapeStation
778 results and sequenced on a Novaseq 6000 system (PE50 on one lane of an SP
779 Flowcell).

780

781 **Transposable element annotations**

782

783 In each respective cichlid genome, transposable elements and repeats were first
784 modelled and identified using RepeatModeler v1.0.11 in combination with the
785 recommended programmes RECON v1.08, RepeatScout v1.0.6, TRF v4.0.9 and
786 NCBI-RMBlast v2.14, and then annotated using RepeatMasker v4.0.9 in combination
787 with NCBI-RMBlast v2.14, TRF v4.0.9 and the custom libraries of modelled repeats,
788 Dfam3.0 and Giri-Repbase-20170127[95]. The curated TE library for Lake Malawi
789 cichlids was created following a previously described protocol[96] and will be described
790 in detail elsewhere (P. Sierra & R. Durbin, unpublished results). This library was used
791 as input to RepeatMasker v4.1.2-p1[95] with options -e rmblast -no_is -gff -lib -a to
792 generate a final TE annotation for the *A. calliptera* genome fAstCal1.2. GTF files with
793 TE annotations amenable to be used for TEtranscripts (see below **Bioinformatic**
794 **analysis**, mRNA-sequencing analysis section) were created using custom scripts
795 (available at https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024).

796

797 **Bioinformatic analysis**

798

799 *mRNA-sequencing analysis*. Illumina adapters and reads with low-quality calls were
800 filtered out with Trimmomatic v0.39[97] using options SLIDINGWINDOW:4:28
801 MINLEN:36. Quality of raw and trimmed fastq files was assessed with fastQC v0.11.9
802 (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) and summarised with
803 multiQC v1.11[98]. Gene expression was quantified from trimmed reads using salmon
804 v1.5.1[99], with options --seqBias --gcBias - validateMappings -l A. Salmon indexes were
805 prepared for each species separately, and used as input (in the -i option) for gene
806 expression quantification in the respective species. DESeq2[100] and custom scripts
807 (available at https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024) were
808 used to calculate normalised and TPM counts, generate plots and conduct statistical
809 tests on an R framework[72]. See R packages used below, in the end of this section.

810

811 Trimmed fastq files were mapped to the cichlid genomes using HISAT2 v2.2.1[101] with
812 options -x -1 -2 -S. Reads from *A. burtoni*, *P. nyererei* and *O. niloticus* were mapped
813 to their respective Ensembl genomes (AstBur1.0, GCA_000239415.1; PunNye1.0,
814 GCA_000239375.1; O_niloticus_UMD_NMBU, GCA_001858045.3). Reads from all

815 Lake Malawi cichlid species used (*A. calliptera*, *M. zebra* and *T.* sp. 'mauve') were
816 mapped to *A. calliptera* Ensembl genome fAstCal1.2 (GCA_900246225.3). SAM
817 alignment files were converted to BAM format, sorted and indexed with samtools
818 v1.10[76]: 1) samtools view -bS ; 2) samtools sort ; and 3) samtools index. To create
819 bigwig files, the BAM alignment files were used as input to bamCoverage v3.5.1, part
820 of the deepTools package[102], using options --normalizeUsing CPM -of bigwig --binSize
821 10. Bigwig files of biological replicates of same organ were combined using
822 WiggleTools[103] mean and wigToBigWig v4[104]. Genome tracks were plotted with
823 custom scripts (available at
824 https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024) using the Gviz[105] and
825 GenomicFeatures[106] packages on an R framework[72].
826
827 To quantify TE expression at the TE family level, we mapped trimmed reads using
828 STAR v2.5.4b[107] with options --readFilesCommand zcat --outSAMtype BAM
829 SortedByCoordinate --outFilterType BySJout --outFilterMultimapNmax
830 150 --winAnchorMultimapNmax 150 --alignSJoverhangMin 8 --
831 alignSJDBoverhangMin 3 --outFilterMismatchNmax 999 --
832 outFilterMismatchNoverReadLmax 0.04 --alignIntronMin 20 --alignIntronMax
833 10000000 --alignMatesGapMax 100000000. As above, reads from *A. burtoni*, *P.
834 nyererei* and *O. niloticus* were mapped to their respective Ensembl genomes
835 (AstBur1.0, GCA_000239415.1; PunNye1.0, GCA_000239375.1;
836 O_niloticus_UMD_NMBU, GCA_001858045.3) and reads from all Lake Malawi cichlid
837 species used (*A. calliptera*, *M. zebra* and *T.* sp. 'mauve') were mapped to *A. calliptera*
838 Ensembl genome fAstCal1.2 (GCA_900246225.3). The resulting BAM files were used
839 as inputs for TEtranscripts v2.2.1[108] with options --stranded reverse --SortByPos.
840 TEtranscripts was run separately for each species, using gene annotations of the
841 respective species downloaded from Ensembl (March 2021) and TE annotations
842 described above (see **Transposable element annotations** section). For Lake Malawi
843 cichlids, TEtranscripts was ran using *A. calliptera* gene and TE annotations (both
844 default and curated versions). A TE family was defined as expressed if it had >10
845 counts in at least 2 samples. DESeq2[100] and custom scripts (available at
846 https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024) were used to
847 calculate normalised counts, generate plots and conduct statistical tests on an R
848 framework[72]. We have used the following R packages: tidyverse[74], lattice[109], eulerr[110],
849 genefilter[111], pheatmap[112], reshape2[113], ggrepel[114], biomaRt[115], tximport[116],
850 RColorBrewer[117], ashr[118], ggpubr[119], GenomicFeatures[106], patchwork[120].
851
852 *mRNA-sequencing analysis of Lake Malawi cichlid embryogenesis datasets.* The
853 embryogenesis dataset collection and experimental design will be reported in detail
854 elsewhere (Chengwei Ulrika Yuan & Eric A. Miska, unpublished results). Trimmomatic-
855 0.39[97] was used to trim the Illumina adapters. Salmon v0.14.2[99] was used to quantify
856 expression of protein-coding genes (--seqBias --validateMappings --gcBias).
857 TEtranscripts analysis on embryo samples was performed as described above
858 (mRNA-sequencing analysis subsection), with one exception: option --stranded no.
859 Locus-specific TE expression levels were analysed with SQuIRE (v0.9.9.9a-beta)[121].
860 For squire Count the option --strandness '0' was run as default for unstranded Illumina
861 data. Reads were mapped to the *A. calliptera* genome (Ensembl, fAstCal1.2), and the
862 TE annotation created from the curated TE library was used (see above, Transposable
863 element annotations section). Tot_counts was used in downstream analysis from the
864 Squire output. Only expressed TEs were kept (defined as >5 reads in at least 2
865 samples). Heatmap and enrichment plots were made from SQuIRE output with code
866 adapted from Chang et al., 2022[13].
867

868 *Small RNA-sequencing analysis*. CutAdapt v1.15[122] was used to remove adapters and
869 reads shorter than 18 nucleotides with options -a
870 TGGAATTCTCGGGTGCCAAGG --minimum-length 18. Quality of raw and trimmed
871 fastq files was assessed with fastQC v0.11.9
872 (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) and summarised with
873 multiQC v1.11[98]. Next, we mapped the trimmed reads to the genome using STAR
874 v2.5.4b[107], with options readFilesCommand zcat --outMultimapperOrder Random --
875 outFilterMultimapNmax 100 --outFilterMismatchNmax 2 --alignIntronMax 1 --
876 outSAMtype BAM SortedByCoordinate --outFilterType BySJout --
877 winAnchorMultimapNmax 100 --alignEndsType EndToEnd --scoreDelOpen -10000 --
878 scoreInsOpen -10000 --outSAMmultNmax 1 --outFileNamePrefix. As above, reads
879 from *A. burtoni*, *P. nyererei* and *O. niloticus* were mapped to their respective Ensembl
880 genomes (AstBur1.0, GCA_000239415.1; PunNye1.0, GCA_000239375.1;
881 O_niloticus_UMD_NMBU, GCA_001858045.3) and reads from all Lake Malawi cichlid
882 species used (*A. calliptera*, *M. zebra* and *T.* sp. 'mauve') were mapped to *A. calliptera*
883 Ensembl genome fAstCal1.2 (GCA_900246225.3). An in-house custom script[123,124]
884 was used, with the BAM files of the alignment as inputs, to create sRNA length
885 distribution profiles in the range of 18-36 nucleotides, and to report 5'-nucleotide
886 frequency, normalised to all mapping reads. The script creates separate sRNA length
887 distribution profiles for 1) collapsed and 2) uncollapsed reads. The first profile keeps
888 only one read of each unique sequence to remove abundance bias, while the second
889 profile keeps all reads. Lastly, the script also produces a FASTA file with the collapsed
890 sequences. With the outputs of the scripts, plots of sRNA length distribution profiles
891 and first nucleotide composition plots were created on an R framework[72] with the
892 packages tidyverse[74], reshape2[113], and RColorBrewer[117].
893
894 Next, we selected sRNAs in the piRNA size range, between 24 and 35 nucleotides
895 long, for further analysis. We have done this size selection on the trimmed reads using
896 CutAdapt v1.15[122] with options --minimum-length 24 --maximum-length 35. We
897 mapped 24-35 nucleotides long sRNAs to the genome with the same settings as
898 discriminated in the previous paragraph. Next, we used "Small RNA Signatures"
899 v3.5.0[125] of the Mississippi Tool Suite from the web-based analysis tool Galaxy to
900 calculate z-scores of overlapping sRNA pairs. For this analysis, alignment BAM files
901 of 24-35 nucleotide long reads were used as input, along with following options: Min
902 size of query sRNAs 24, Max size of query sRNAs 35, Min size of target sRNAs 24,
903 Max size of target sRNAs 35, Minimal relative overlap analyzed 1, Maximal relative
904 overlap analyzed 26. To find signatures of phased piRNA biogenesis, BAM files of 24-
905 35 nucleotide long reads were loaded into R as Genomic Ranges[106] and using
906 RSamtools[126], the Follow function was used to identify the next mapping piRNA pair
907 and distances between the 5' and 3' were calculated for plotting. To create sequence
908 logos, we first ran the custom script described above[123,124] to produce a FASTA file
909 with the 24-35 nucleotide long collapsed reads (unique sequences). Then, we created
910 a new FASTA file with all these reads trimmed from the 3' end to a total length of 20
911 nucleotides, and concatenated together the FASTA files of the biological replicates for
912 each species and organ. The FASTA file with the concatenated and trimmed
913 sequences was in turn used to generate sequence logos in R (scripts available at
914 https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024), with packages
915 ggseqlogo[127], phylotools[128], and tidyverse[74]. This process was repeated to generate
916 sequence logos of piRNAs mapping sense or antisense in regard to TE orientation
917 using BAM files with 24-35 nucleotide long reads, which were created as follows: 1)
918 samtools view -b -f (16 or 0); 2) bedtools intersect (-s or -S); 3) samtools merge; 4)
919 samtools sort; 5) samtools index.
920

921 To quantify piRNA counts associated with TEs, we used featureCounts v1.6.0[129] with
922 options -t exon -M. The 24-35 nucleotide long BAM file was used as input. The
923 TEtranscripts-compatible TE annotations described above (see **Transposable**
924 **element annotations**) were provided as the intersecting features. For Lake Malawi
925 cichlids, featureCounts analysis was performed twice, using *A. calliptera* default and
926 curated TE annotations. After obtaining the tables of counts, DESeq2[100] and custom
927 scripts (available at https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024)
928 were used to calculate normalised counts, generate plots and conduct statistical tests
929 on an R framework[72], with packages tidyverse[74], lattice[109], eulerr[110], genefilter[111],
930 pheatmap[112], reshape2[113], ggrepel[114], biomaRt[115], tximport[116], RColorBrewer[117],
931 ashr[118], ggpubr[119], GenomicFeatures[106], patchwork[120]. To create bigwig files, the 24-
932 35 nucleotide long BAM alignment files were used as inputs to bamCoverage v3.5.1,
933 part of the deepTools package[102], using options --normalizeUsing CPM -of bigwig --
934 binSize 5. Bigwig files of biological replicates of same organ were combined using
935 WiggleTools[103] mean and wigToBigWig v4[104]. Genome tracks were plotted with
936 custom scripts (available at
937 https://github.com/migueldvalmeida/Cichlid_TEs_piRNAs2024) using the Gviz[105] and
938 GenomicFeatures[106] packages on an R framework[72]. We used these bigwig files to
939 produce sRNA metagene profiles with deepTools[102] computeMatrix scale-regions
940 v3.5.1 (options -b 1000 -a 1000 --regionBodyLength 2000 --averageTypeBins
941 median --missingDataAsZero --binSize 5) and plotProfile v3.5.1 (--plotType se --
942 averageType mean --perGroup). To generate metagene profiles against particular TE
943 classes or superfamilies, TE annotations were subsetted by TE class or superfamily
944 and converted to bed format with grep and awk utilities. The resulting bed files
945 contained the regions to plot and were used as input for computeMatrix.
946
947 To define piRNA clusters, we first re-mapped trimmed reads 24-35 nucleotides long
948 to the *A. calliptera* (fAstCal1.2) or *O. niloticus* (O_niloticus_UMD_NMBU) genomes
949 using STAR v2.5.4b[107], with options: --readFilesCommand zcat --
950 outFilterMultimapNmax 100 --outFilterMismatchNmax 2 --alignIntronMax 1 --
951 outSAMtype BAM SortedByCoordinate --outFilterType BySJout --
952 alignSoftClipAtReferenceEnds No --winAnchorMultimapNmax 100 --alignEndsType
953 EndToEnd --scoreDelOpen -10000 --scoreInsOpen -10000 --outSAMmultNmax 100 -
954 -outSAMattributes All. Method and code for the approach below will be detailed
955 elsewhere (A. Friman and A. Haase, unpublished results). The resulting BAM files
956 were loaded into R environment using GenomicAlignments package[106]. For each BAM
957 file the alignments were sorted into three categories: unique mapping alignments,
958 primary multimapping alignments, and secondary multimapping alignments (method
959 by A. Friman and A. Haase, unpublished results). The reference genome was split into
960 sliding windows[106] with size and step between starting position depending on the
961 alignments category. For unique mapping alignments the windows were 350 nt
962 (window size) starting at every 35 nt (window step) of genome length. For each of
963 these windows the number of overlapping unique mapping alignments was counted.
964 If the number was at least 2 FPKM (RPKM), the window was called. The called
965 windows were reduced into genomic intervals named "seeds", indicating the genomic
966 origin of uniquely mapping piRNAs. Seeds that were shorter than 800 nt were
967 discarded to reduce false positives, which can be caused by individual degradation
968 fragments of abundant structural RNAs or other cellular transcripts. Next, we
969 incorporated multimapping piRNA reads, considering first their primary alignments and
970 then all possible alignments (up to 1000 according to the parameters used for genome
971 mapping). We counted primary multimapping alignments using 350 nt long sliding
972 windows (window size) located at every 35 nt (window step) of genome length.
973 Windows overlapping with more than 4 FPKM (RPKM) with each other and with
974 previously established 'seeds', were reduced into intervals named 'cores'. Each 'core'

975 was required to overlap with at least one seed. Finally, we integrated all secondary
976 multimapping alignments using 1000 nt long sliding windows (window size) with 100
977 nt step (window step). We requested read coverage greater or equal 0.2 FPKM
978 (RPKM) as threshold. Overlapping windows were reduced into 'clusters' when they
979 overlapped with at least one 'core'. All clusters contain strand information and predict
980 one or multiple piRNA precursor transcripts from a defined genomic strand.
981 Intersection[106] of genomic 'cluster' coordinates from different samples or biological
982 replicates take strand information into account. Results were plotted on a R
983 framework[72], using packages: tidyverse (Wickham et al., 2019), reshape2[113], and
984 ggpubr[119]. Circos plots were created with Circos v0.69-8[130]. Density tracks are
985 displayed on the circos plots as the number of features per mega-base.
986
987 **Protein preparations and mass spectrometry**
988
989 Frozen brain and gonad tissues were partitioned on a mortar positioned on dry ice,
990 and weighed. This was done quickly to avoid thawing. A similar mass of the same
991 tissue was collected from size-matched individuals of the same species to create
992 biological replicates. 6-50 mgs of brain tissue, and 8-120 mgs of gonad tissue were
993 used, according to the specific tissue, tissue availability, and size of the specimen,
994 which varied per species. Partitioned tissues were transferred to BeadBug tubes
995 prefilled with 0.5 mm Zirconium beads (Merck, #Z763772) together with 150 µl (if using
996 6-20 mg of tissue) or 250 µl (if using >20 mg of tissue) of modified RIPA buffer (50 mM
997 Tris HCl pH 7.5, 150 mM NaCl, 1% IGEPAL CA-630, 1% Sodium Deoxycholate,
998 supplemented with cOmplete EDTA-free protease inhibitor cocktail tablets, Roche
999 #4693132001). Next, homogenisation was conducted using a BeadBug microtube
1000 homogeniser (Sigma, #Z764140) at approximately 4°C (conducted in cold room). Each
1001 sample was homogenised with five BeadBug runs at maximum speed (4,000 rpm) for
1002 60 seconds each. Did not run any sample more than two consecutive times to avoid
1003 overheating. Other than the run time inside the BeadBug, samples were left on ice.
1004 After homogenisation, lysates were centrifuged for 5 minutes at 18,000 G at 4°C.
1005 Supernatant was then removed into a clean 1.5 mL tube. Centrifuge the lysates again,
1006 this time at maximum speed (approximately 21,000 G) for 5 minutes at 4°C. Transfer
1007 supernatant into a clean tube without disturbing the pellet and tissue debris. Measured
1008 protein concentration using Bradford (Bio-Rad, Protein Assay Dye Reagent
1009 Concentrate, #5000006) and prepared a final sample by combining 150 µg of lysate,
1010 1x LDS (prepared from NuPAGE LDS Sample Buffer 4x, Thermo Scientific, #NP0007)
1011 and 100 mM DTT and boiling for 10 minutes at 95°C. Half of the sample was sent for
1012 mass spectrometry.
1013
1014 In-gel digestion for mass spectrometry was performed as previously described[131].
1015 Samples were boiled at 70°C for 10 minutes prior to loading on a 4%-12% NuPAGE
1016 Bis-Tris gel (Thermo Scientific, #NP0321). The gel was run in 1x MOPS buffer at 180V
1017 for 10 minutes and subsequently fixed and stained with Coomassie G250 (Carl Roth).
1018 Each lane was minced and transferred to a 1.5 mL reaction tube, destained with 50%
1019 EtOH in 50 mM ammonium bicarbonate buffer (pH 8.0). Gel pieces were dehydrated
1020 with 100% acetonitrile and dried in a Concentrator Plus (Eppendorf, #5305000304).
1021 Then, samples were reduced with 10 mM DTT / 50 mM ABC buffer (pH 8.0) at 56°C
1022 and alkylated with 50 mM iodoacetamide / 50 mM ABC buffer (pH 8.0) in the dark.
1023 After washing with ABC buffer (pH 8.0) and dehydration with acetonitrile the proteins
1024 were digested with 1 µg mass spectrometry-grade Trypsin (Serva) at 37°C overnight.
1025 The peptides were purified on stage tips as previously described[132]. Peptides were
1026 analysed by nanoflow liquid chromatography using an EASYnLC 1200 system
1027 (Thermo Scientific) coupled to an Exploris 480 (Thermo Scientific). Peptides were
1028 separated on a C18-reversed phase column (60 cm, 75µm diameter), packed in-house

1029 with Reprosil aq1.9 (Dr. Maisch GmbH), mounted on the electrospray ion source of
1030 the mass spectrometer. Peptides were eluted from the column with an optimized 103-
1031 min gradient from 2% to 40% of a mixture of 80% acetonitrile/0.1% formic acid at a
1032 flow rate of 250 nL/min. The Exploris was operated in positive ion mode with a data-
1033 dependent acquisition strategy of one mass spectrometry full scan (scan range 300–
1034 1650 m/z; 60,000 resolution; normalised AGC target 300%; max IT 28 ms) and up to
1035 20 MS/MS scans (15,000 resolution; AGC target 100%, max IT 28 ms; isolation
1036 window 1.4 m/z) with peptide match preferred using HCD fragmentation. Mass
1037 spectrometry measurements were analysed with MaxQuant v1.6.10.43[133] with the
1038 following protein databases (downloaded from Ensembl):
1039 Haplochromis_burtoni.AstBur1.0.pep.all.fa (35,619 entries, from *A. burtoni*),
1040 Oreochromis_niloticus.O_niloticus_UMD_NMBU.pep.all.fa (75,555 entries, from *O.*
1041 *niloticus*), Astatotilapia_calliptera.fAstCal1.2.pep.all.fa (41,597 entries, from *A.*
1042 *calliptera*), and Pundamilia_nyererei.PunNye1.0.pep.all.fa (32,153 entries, from *P.*
1043 *nyererei*). Missing values were imputed at the lower end of LFQ values using random
1044 values from a beta distribution fitted at 0.2-2.5%. Prior to further analysis, protein
1045 groups with contaminants, reverse hits and only identified by site were removed.
1046
1047

## Data accessibility

1049
1050 The mass spectrometry proteomics data have been deposited to the
1051 ProteomeXchange Consortium via the PRIDE[134] partner repository with the dataset
1052 identifier PXD047439. The mRNA and sRNA sequencing data generated in this study
1053 have been deposited to GEO under accession numbers GSE252804 and
1054 GSE252805. The genomic data of Lake Malawi cichlids used in this work is available
1055 on SRA, bioproject PRJEB1254 (see a list of samples in **Supplemental Table 1**), on
1056 an open access basis for research use only. Any person who wishes to use this data
1057 for any form of commercial purpose must first enter into a commercial licensing and
1058 benefit sharing arrangement with the Government of Malawi.
1059
1060

## Acknowledgements

1081
1082

## Funding

## Author contributions

Conceptualisation: M.V.A. and E.A.M.; Data curation: M.V.A., M.B., and J.L.P.; Formal analysis: M.V.A., M.B., C.U.Y., P.S., J.L.P., and F.X.Q.; Funding acquisition: M.V.A., A.D.H., R.D. and E.A.M.; Investigation: M.V.A., and C.U.Y.; Project administration: M.V.A. and E.A.M.; Resources: P.S., A.F., A.D., G.V., A.L.K.P., A.M.S., D.A.J., F.B., A.D.H., R.D., M.E.S.; Software: A.F. and A.D.H.; Supervision: M.V.A., A.D.H., R.D., M.E.S., and E.A.M; Visualisation: M.V.A., M.B., C.U.Y., and J.L.P.; Writing – original draft: M.V.A.; Writing – review & editing: all authors contributed.

## Competing interests

The authors declare no competing interests.

## References

1. Salzburger, W. Understanding explosive diversification through cichlid fish genomics. Nature Reviews Genetics 19, 705 (2018).

2. Santos, M. E., Lopes, J. F. & Kratochwil, C. F. East African cichlid fishes. EvoDevo 14, 1 (2023).

3. Svardal, H., Salzburger, W. & Malinsky, M. Genetic Variation and Hybridization in Evolutionary Radiations of Cichlid Fishes. Annual Review of Animal Biosciences 9, 55–79 (2021).

4. Malinsky, M. et al. Whole-genome sequences of Malawi cichlids reveal multiple radiations interconnected by gene flow. Nature Ecology & Evolution 2, 1940 (2018).

5. Quah, F. X. et al. A pangenomic perspective of the Lake Malawi cichlid radiation reveals extensive structural variation driven by transposable elements. 2024.03.28.587230 Preprint at https://doi.org/10.1101/2024.03.28.587230 (2024).

6.    Wells, J. N. & Feschotte, C. A Field Guide to Eukaryotic Transposable Elements. Annual Review of Genetics 54, 539–561 (2020).

7.    Arkhipova, I. R. Neutral Theory, Transposable Elements, and Eukaryotic Genome Evolution. Molecular Biology and Evolution 35, 1332–1337 (2018).

8.    Almeida, M. V., Vernaz, G., Putman, A. L. K. & Miska, E. A. Taming transposable elements in vertebrates: from epigenetic silencing to domestication. Trends in Genetics (2022) doi:10.1016/j.tig.2022.02.009.

9.    Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable elements: from conflicts to benefits. Nature Reviews Genetics 18, 71–86 (2017).

10.   Fueyo, R., Judd, J., Feschotte, C. & Wysocka, J. Roles of transposable elements in the regulation of mammalian transcription. Nat Rev Mol Cell Biol 1–17 (2022) doi:10.1038/s41580-022-00457-y.

11.   Carducci, F., Barucca, M., Canapa, A., Carotti, E. & Biscotti, M. A. Mobile Elements in Ray-Finned Fish Genomes. Life 10, 221 (2020).

12.   Chalopin, D., Naville, M., Plard, F., Galiana, D. & Volff, J.-N. Comparative Analysis of Transposable Elements Highlights Mobilome Diversity and Evolution in Vertebrates. Genome Biology and Evolution 7, 567–580 (2015).

13.   Chang, N.-C., Rovira, Q., Wells, J. N., Feschotte, C. & Vaquerizas, J. M. Zebrafish transposable elements show extensive diversification in age, genomic distribution, and developmental expression. Genome Res. gr.275655.121 (2022) doi:10.1101/gr.275655.121.

14.   Gao, B. et al. The contribution of transposable elements to size variations between four teleost genomes. Mobile DNA 7, 4 (2016).

15.   Reinar, W. B. et al. Teleost genomic repeat landscapes in light of diversification rates and ecology. Mobile DNA 14, 14 (2023).

16.   Shao, F., Han, M. & Peng, Z. Evolution and diversity of transposable elements in fish genomes. Sci Rep 9, 15399 (2019).

17.   Yuan, Z. et al. Comparative genome analysis of 52 fish species suggests differential associations of repetitive elements with their living aquatic environments. BMC Genomics 19, 141 (2018).

18.   Brawand, D. et al. The genomic substrate for adaptive radiation in African cichlid fish. Nature 513, 375–381 (2014).

19.   Kratochwil, C. F. et al. An intronic transposon insertion associates with a trans-species color polymorphism in Midas cichlid fishes. Nat Commun 13, 296 (2022).

20.   Santos, M. E. et al. The evolution of cichlid fish egg-spots is linked with a cis-regulatory change. Nature Communications 5, 5149 (2014).

21.   Munby, H. et al. Differential use of multiple genetic sex determination systems in divergent ecomorphs of an African crater lake cichlid. bioRxiv https://doi.org/10.1101/2021.08.05.455235, 2021.08.05.455235 (2021).

22.   Carleton, K. L. et al. Movement of transposable elements contributes to cichlid diversity. Molecular Ecology 29, 4956–4969 (2020).

23.   Vernaz, G. et al. Mapping epigenetic divergence in the massive radiation of Lake Malawi cichlid fishes. Nat Commun 12, 5870 (2021).

24.   Bruno, M., Mahgoub, M. & Macfarlan, T. S. The Arms Race Between KRAB–Zinc Finger Proteins and Endogenous Retroelements and Its Impact on Mammals. Annual Review of Genetics 53, 393–416 (2019).

25. Cosby, R. L., Chang, N.-C. & Feschotte, C. Host–transposon interactions: conflict, cooperation, and cooption. Genes Dev. 33, 1098–1116 (2019).

26. Loubalova, Z., Konstantinidou, P. & Haase, A. D. Themes and variations on piRNA-guided transposon control. Mobile DNA 14, 10 (2023).

27. Ozata, D. M., Gainetdinov, I., Zoch, A., O'Carroll, D. & Zamore, P. D. PIWI-interacting RNAs: small RNAs with big functions. Nat Rev Genet 20, 89–108 (2019).

28. Seczynska, M., Bloor, S., Cuesta, S. M. & Lehner, P. J. Genome surveillance by HUSH-mediated silencing of intronless mobile elements. Nature 1–9 (2021) doi:10.1038/s41586-021-04228-1.

29. Tchasovnikarova, I. A. et al. Epigenetic silencing by the HUSH complex mediates position-effect variegation in human cells. Science 348, 1481–1485 (2015).

30. Wells, J. N. et al. Transposable elements drive the evolution of metazoan zinc finger genes. Genome Res. (2023) doi:10.1101/gr.277966.123.

31. Aravin, A. et al. A novel class of small RNAs bind to MILI protein in mouse testes. Nature 442, 203–207 (2006).

32. Brennecke, J. et al. Discrete Small RNA-Generating Loci as Master Regulators of Transposon Activity in Drosophila. Cell 128, 1089–1103 (2007).

33. Houwing, S. et al. A Role for Piwi and piRNAs in Germ Cell Maintenance and Transposon Silencing in Zebrafish. Cell 129, 69–82 (2007).

34. Gainetdinov, I., Colpan, C., Arif, A., Cecchini, K. & Zamore, P. D. A Single Mechanism of Biogenesis, Initiated and Directed by PIWI Proteins, Explains piRNA Production in Most Animals. Molecular Cell 71, 775-790.e5 (2018).

35. Gunawardane, L. S. et al. A Slicer-Mediated Mechanism for Repeat-Associated siRNA 5' End Formation in Drosophila. Science 315, 1587–1590 (2007).

36. Han, B. W., Wang, W., Li, C., Weng, Z. & Zamore, P. D. piRNA-guided transposon cleavage initiates Zucchini-dependent, phased piRNA production. Science 348, 817–821 (2015).

37. Mohn, F., Handler, D. & Brennecke, J. piRNA-guided slicing specifies transcripts for Zucchini-dependent, phased piRNA biogenesis. Science 348, 812–817 (2015).

38. Lewis, S. H. et al. Pan-arthropod analysis reveals somatic piRNAs as an ancestral defence against transposable elements. Nat Ecol Evol 2, 174–181 (2018).

39. Lewis, S. H., Salmela, H. & Obbard, D. J. Duplication and Diversification of Dipteran Argonaute Genes, and the Evolutionary Divergence of Piwi and Aubergine. Genome Biol Evol 8, 507–518 (2016).

40. Palmer, W. H., Hadfield, J. D. & Obbard, D. J. RNA-Interference Pathways Display High Rates of Adaptive Protein Evolution in Multiple Invertebrates. Genetics 208, 1585–1599 (2018).

41. Parhad, S. S., Tu, S., Weng, Z. & Theurkauf, W. E. Adaptive Evolution Leads to Cross-Species Incompatibility in the piRNA Transposon Silencing Machinery. Developmental Cell 43, 60-70.e5 (2017).

42. Fryer, G. & Iles, T. D. The Cichlid Fishes of the Great Lakes of Africa: Their Biology and Evolution. (Oliver and Boyd, Edinburgh, 1972).

43. Marconi, A., Yang, C. Z., McKay, S. & Santos, M. E. Morphological and temporal variation in early embryogenesis contributes to species divergence in Malawi cichlid fishes. Evolution & Development 25, 170–193 (2023).

28

44.    Wang, X., Ramat, A., Simonelig, M. & Liu, M.-F. Emerging roles and functional mechanisms of PIWI-interacting RNAs. Nat Rev Mol Cell Biol 24, 123–141 (2023).

45.    Li, M. A. et al. Mobilization of giant piggyBac transposons in the mouse genome. Nucleic Acids Research 39, e148 (2011).

46.    Li, X. et al. piggyBac transposase tools for genome engineering. Proceedings of the National Academy of Sciences 110, E2279–E2287 (2013).

47.    Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community resource of transposable element families, sequence models, and genome annotations. Mobile DNA 12, 2 (2021).

48.    Cary, L. C. et al. Transposon mutagenesis of baculoviruses: Analysis of Trichoplusia ni transposon IFP2 insertions within the FP-locus of nuclear polyhedrosis viruses. Virology 172, 156–169 (1989).

49.    Innan, H. & Kondrashov, F. The evolution of gene duplications: classifying and distinguishing between models. Nature Reviews Genetics 11, 97–108 (2010).

50.    Matsumoto, N. et al. Crystal Structure of Silkworm PIWI-Clade Argonaute Siwi Bound to piRNA. Cell 167, 484-497.e9 (2016).

51.    Yamaguchi, S. et al. Crystal structure of Drosophila Piwi. Nat Commun 11, 858 (2020).

52.    Parker, J. S., Roe, S. M. & Barford, D. Crystal structure of a PIWI protein suggests mechanisms for siRNA recognition and slicer activity. The EMBO Journal 23, 4727–4737 (2004).

53.    Yuan, Y.-R. et al. Crystal Structure of A. aeolicus Argonaute, a Site-Specific DNA-Guided Endoribonuclease, Provides Insights into RISC-Mediated mRNA Cleavage. Molecular Cell 19, 405–419 (2005).

54.    Stein, C. B. et al. Decoding the 5′ nucleotide bias of PIWI-interacting RNAs. Nat Commun 10, 828 (2019).

55.    Irisarri, I. et al. Phylogenomics uncovers early hybridization and adaptive loci shaping the radiation of Lake Tanganyika cichlid fishes. Nat Commun 9, 3159 (2018).

56.    Meier, J. I. et al. Cycles of fusion and fission enabled rapid parallel adaptive radiations in African cichlids. Science 381, eade2833 (2023).

57.    Meier, J. I. et al. Ancient hybridization fuels rapid cichlid fish adaptive radiations. Nat Commun 8, 14363 (2017).

58.    Svardal, H. et al. Ancestral Hybridization Facilitated Species Diversification in the Lake Malawi Cichlid Fish Adaptive Radiation. Molecular Biology and Evolution 37, 1100–1113 (2020).

59.    Brawand, D. et al. The evolution of gene expression levels in mammalian organs. Nature 478, 343–348 (2011).

60.    Soumillon, M. et al. Cellular Source and Mechanisms of High Transcriptome Complexity in the Mammalian Testis. Cell Reports 3, 2179–2190 (2013).

61.    Sakashita, A. et al. Endogenous retroviruses drive species-specific germline transcriptomes in mammals. Nat Struct Mol Biol 27, 967–977 (2020).

62.    Lee, H. J. et al. Epigenomic analysis reveals prevalent contribution of transposable elements to cis-regulatory elements, tissue-specific expression, and alternative promoters in zebrafish. Genome Res. 32, 1424–1436 (2022).

1270   63.   El Taher, A. et al. Gene expression dynamics during rapid organismal diversification in African cichlid fishes. Nature Ecology & Evolution 1–8 (2020) doi:10.1038/s41559-020-01354-3.

1273   64.   Yi, M. et al. Rapid Evolution of piRNA Pathway in the Teleost Fish: Implication for an Adaptation to Transposon Diversity. Genome Biology and Evolution 6, 1393–1407 (2014).

1276   65.   Frank, F., Sonenberg, N. & Nagar, B. Structural basis for 5′-nucleotide base-specific recognition of guide RNA by human AGO2. Nature 465, 818–822 (2010).

1278   66.   Song, J.-J., Smith, S. K., Hannon, G. J. & Joshua-Tor, L. Crystal Structure of Argonaute and Its Implications for RISC Slicer Activity. Science 305, 1434–1437 (2004).

1281   67.   Yan, K. S. et al. Structure and conserved RNA binding of the PAZ domain. Nature 426, 469–474 (2003).

1283   68.   Genner, M. J. Cichlid fish seized an ecological opportunity to diversify. Nature (2023) doi:10.1038/d41586-023-03014-5.

1285   69.   Ngoepe, N. et al. A continuous fish fossil record reveals key insights into adaptive radiation. Nature 1–6 (2023) doi:10.1038/s41586-023-06603-6.

1287   70.   Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative genomics. Genome Biology 20, 238 (2019).

1289   71.   Emms, D. M. & Kelly, S. OrthoFinder: solving fundamental biases in whole genome comparisons dramatically improves orthogroup inference accuracy. Genome Biology 16, 157 (2015).

1292   72.   R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing, Vienna, Austria (2021).

1294   73.   Brennan, P. drawProteins: a Bioconductor/R package for reproducible and programmatic generation of protein schematics. Preprint at https://doi.org/10.12688/f1000research.14541.1 (2018).

1297   74.   Wickham, H. et al. Welcome to the Tidyverse. Journal of Open Source Software 4, 1686 (2019).

1299   75.   Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. Bioinformatics 25, 1754–1760 (2009).

1301   76.   Li, H. et al. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25, 2078–2079 (2009).

1303   77.   Li, H. Minimap2: pairwise alignment for nucleotide sequences. Bioinformatics 34, 3094–3100 (2018).

1305   78.   Robinson, J. T. et al. Integrative genomics viewer. Nat Biotechnol 29, 24–26 (2011).

1307   79.   Martin, F. J. et al. Ensembl 2023. Nucleic Acids Research 51, D933–D941 (2023).

1309   80.   Jain, C., Rhie, A., Hansen, N. F., Koren, S. & Phillippy, A. M. Long-read mapping to repetitive reference sequences using Winnowmap2. Nat Methods 19, 705–710 (2022).

1312   81.   Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. Nucleic Acids Res 32, 1792–1797 (2004).

1314   82.   Larsson, A. AliView: a fast and lightweight alignment viewer and editor for large datasets. Bioinformatics 30, 3276–3278 (2014).

83.    Danecek, P. et al. Twelve years of SAMtools and BCFtools. GigaScience 10, giab008 (2021).

84.    Alachiotis, N. & Pavlidis, P. RAiSD detects positive selection based on multiple signatures of a selective sweep and SNP vectors. Commun Biol 1, 1–11 (2018).

85.    Katoh, K. & Standley, D. M. MAFFT Multiple Sequence Alignment Software Version 7: Improvements in Performance and Usability. Molecular Biology and Evolution 30, 772–780 (2013).

86.    Waterhouse, A. M., Procter, J. B., Martin, D. M. A., Clamp, M. & Barton, G. J. Jalview Version 2—a multiple sequence alignment editor and analysis workbench. Bioinformatics 25, 1189–1191 (2009).

87.    Minh, B. Q. et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. Molecular Biology and Evolution 37, 1530–1534 (2020).

88.    Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q. & Vinh, L. S. UFBoot2: Improving the Ultrafast Bootstrap Approximation. Molecular Biology and Evolution 35, 518–522 (2018).

89.    Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence alignments into the corresponding codon alignments. Nucleic Acids Research 34, W609–W612 (2006).

90.    Weaver, S. et al. Datamonkey 2.0: A Modern Web Application for Characterizing Selective and Other Evolutionary Processes. Molecular Biology and Evolution 35, 773–777 (2018).

91.    Murrell, B. et al. Gene-Wide Identification of Episodic Selection. Molecular Biology and Evolution 32, 1365–1371 (2015).

92.    Murrell, B. et al. Detecting Individual Sites Subject to Episodic Diversifying Selection. PLOS Genetics 8, e1002764 (2012).

93.    Jumper, J. et al. Highly accurate protein structure prediction with AlphaFold. Nature 596, 583–589 (2021).

94.    Varadi, M. et al. AlphaFold Protein Structure Database: massively expanding the structural coverage of protein-sequence space with high-accuracy models. Nucleic Acids Research 50, D439–D444 (2022).

95.    Smit, A. F. A., Hubley, R. & Green, P. RepeatMasker Open-4.0. (2013).

96.    Goubert, C. et al. A beginner's guide to manual curation of transposable elements. Mobile DNA 13, 7 (2022).

97.    Bolger, A. M., Lohse, M. & Usadel, B. Trimmomatic: a flexible trimmer for Illumina sequence data. Bioinformatics 30, 2114–2120 (2014).

98.    Ewels, P., Magnusson, M., Lundin, S. & Käller, M. MultiQC: summarize analysis results for multiple tools and samples in a single report. Bioinformatics 32, 3047–3048 (2016).

99.    Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast and bias-aware quantification of transcript expression. Nat Methods 14, 417–419 (2017).

100.    Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. Genome Biology 15, (2014).

1360  101.    Pertea, M., Kim, D., Pertea, G. M., Leek, J. T. & Salzberg, S. L. Transcript-level
1361       expression analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. Nat
1362       Protoc 11, 1650–1667 (2016).

1363  102.    Ramírez, F. et al. deepTools2: a next generation web server for deep-
1364       sequencing data analysis. Nucleic Acids Res 44, W160–W165 (2016).

1365  103.    Zerbino, D. R., Johnson, N., Juettemann, T., Wilder, S. P. & Flicek, P.
1366       WiggleTools: parallel processing of large collections of genome-wide datasets for
1367       visualization and statistical analysis. Bioinformatics 30, 1008–1009 (2014).

1368  104.    Kent, W. J., Zweig, A. S., Barber, G., Hinrichs, A. S. & Karolchik, D. BigWig and
1369       BigBed: enabling browsing of large distributed datasets. Bioinformatics 26, 2204–2207
1370       (2010).

1371  105.    Hahne, F. & Ivanek, R. Visualizing Genomic Data Using Gviz and Bioconductor.
1372       in Statistical Genomics: Methods and Protocols (eds. Mathé, E. & Davis, S.) 335–351
1373       (Springer, New York, NY, 2016). doi:10.1007/978-1-4939-3578-9_16.

1374  106.    Lawrence, M. et al. Software for Computing and Annotating Genomic Ranges.
1375       PLOS Computational Biology 9, e1003118 (2013).

1376  107.    Dobin, A. et al. STAR: ultrafast universal RNA-seq aligner. Bioinformatics 29,
1377       15–21 (2013).

1378  108.    Jin, Y., Tam, O. H., Paniagua, E. & Hammell, M. TEtranscripts: a package for
1379       including transposable elements in differential expression analysis of RNA-seq
1380       datasets. Bioinformatics 31, 3593–3599 (2015).

1381  109.    Sarkar, D. et al. lattice: Trellis Graphics for R. (2023).

1382  110.    Larsson, J. et al. eulerr: Area-Proportional Euler and Venn Diagrams with
1383       Ellipses. (2022).

1384  111.    Gentleman, R. et al. genefilter: genefilter: methods for filtering genes from high-
1385       throughput experiments. Bioconductor version: Release (3.17)
1386       https://doi.org/10.18129/B9.bioc.genefilter (2023).

1387  112.    Kolde, R. pheatmap: Pretty Heatmaps. (2019).

1388  113.    Wickham, H. reshape2: Flexibly Reshape Data: A Reboot of the Reshape
1389       Package. (2020).

1390  114.    Slowikowski, K. et al. ggrepel: Automatically Position Non-Overlapping Text
1391       Labels with 'ggplot2'. (2023).

1392  115.    Durinck, S. et al. biomaRt: Interface to BioMart databases (i.e. Ensembl).
1393       Bioconductor version: Release (3.17) https://doi.org/10.18129/B9.bioc.biomaRt
1394       (2023).

1395  116.    Love, M. et al. tximport: Import and summarize transcript-level estimates for
1396       transcript- and gene-level analysis. Bioconductor version: Release (3.17)
1397       https://doi.org/10.18129/B9.bioc.tximport (2023).

1398  117.    Neuwirth, E. RColorBrewer: ColorBrewer Palettes. (2022).

1399  118.    Stephens, M. et al. ashr: Methods for Adaptive Shrinkage, using Empirical
1400       Bayes. (2023).

1401  119.    Kassambara, A. ggpubr: 'ggplot2' Based Publication Ready Plots. (2023).

1402  120.    Pedersen, T. L. patchwork: The Composer of Plots. (2023).

121. Yang, W. R., Ardeljan, D., Pacyna, C. N., Payer, L. M. & Burns, K. H. SQuIRE reveals locus-specific regulation of interspersed repeat expression. Nucleic Acids Research 47, e27 (2019).

122. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. EMBnet.journal 17, 10–12 (2011).

123. Di Domenico, T. tstk/peterplot.py · master · Tomás Di Domenico / tstk · GitLab. GitLab https://gitlab.com/tdido/tstk/-/blob/master/tstk/peterplot.py (2022).

124. Ramakrishna, N. B., Battistoni, G., Surani, M. A., Hannon, G. J. & Miska, E. A. Mouse primordial germ-cell-like cells lack piRNAs. Developmental Cell 57, 2661-2668.e5 (2022).

125. Antoniewski, C. Computing siRNA and piRNA Overlap Signatures. in Animal Endo-SiRNAs: Methods and Protocols (ed. Werner, A.) 135–146 (Springer, New York, NY, 2014). doi:10.1007/978-1-4939-0931-5_12.

126. Morgan, M., Pagès, H., Obenchain, V. & Hayden, N. Rsamtools. Bioconductor http://bioconductor.org/packages/Rsamtools/ (2023).

127. Wagih, O. ggseqlogo: a versatile R package for drawing sequence logos. Bioinformatics 33, 3645–3647 (2017).

128. Zhang, J. helixcn/phylotools. (2023).

129. Liao, Y., Smyth, G. K. & Shi, W. featureCounts: an efficient general purpose program for assigning sequence reads to genomic features. Bioinformatics 30, 923–930 (2014).

130. Krzywinski, M. I. et al. Circos: An information aesthetic for comparative genomics. Genome Res. (2009) doi:10.1101/gr.092759.109.

131. Shevchenko, A., Tomas, H., Havli, J., Olsen, J. V. & Mann, M. In-gel digestion for mass spectrometric characterization of proteins and proteomes. Nat. Protocols 1, 2856–2860 (2007).

132. Rappsilber, J., Mann, M. & Ishihama, Y. Protocol for micro-purification, enrichment, pre-fractionation and storage of peptides for proteomics using StageTips. Nat. Protocols 2, 1896–1906 (2007).

133. Cox, J. & Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. Nat Biotech 26, 1367–1372 (2008).

134. Perez-Riverol, Y. et al. The PRIDE database resources in 2022: a hub for mass spectrometry-based proteomics evidences. Nucleic Acids Research 50, D543–D552 (2022).

**A** (map/fish illustrations)

ON — *Oreochromis niloticus* (outgroup)

PN — *Pundamilia nyererei* (LV)

AB — *Astatotilapia burtoni* (LT)

Lake Victoria (LV)
Lake Tanganyika (LT)
Lake Malawi (LM)

TM — *Tropheops* sp. 'mauve' (LM)

MZ — *Maylandia zebra* (LM)

AC — *Astatotilapia calliptera* (LM)

**B**

rlog normalized counts

TE class: DNA, LINE, LTR, RC, SINE

$^{‡}$ Curated annotation

Developmental stage

- 2-cell
- 1k cells
- Early gastrula
- Late gastrula
- Early somitogenesis
- Mid somitogenesis
- Late somitogenesis
- Early pharyngula
- Mid pharyngula
- Late pharyngula
- Hatching
- Larvae

**C**

Dev. stage

- 2-cell/ 1k cells — A
- Gastrula/ Early somitogenesis — B
- Somitogenesis — C
- Pharyngula/ hatching/ larvae — D
- Late pharyngula/ hatching/ larvae — E

Adjusted p-value
- > 0.05
- < 0.05
- 1e-20
- 1e-50

Log2 Odds Ratio
2, 0, -2, -4

x-axis labels: Dong-R4, L1, L2, Penelope, Rex-Babar, RTE-BovB, Copia, ERV1, ERVK, ERVL, Gypsy, Pao, CMC-EnSpm, hAT, hAT-Ac, hAT-Blackjack, hAT-Charlie, hAT-Tip100, Maverick, MULE-MuDR, PIF-Harbinger, PIF-ISL2EU, TcMar, TcMar-Tc1, TcMar-Tc2, 5S-Deu-L2, tRNA-Core-RTE

**A**

**B** Schematics of *A. calliptera* Piwi proteins

Piwi expansion

Piwil1.1
Piwil1.2*
Piwil1.3
Piwil1.4

PAZ   MID   PIWI

0    400    800
Amino acid number

*Protein annotation avoids very likely stop codons.

**C**

| | | Eco-morphological group | *piwil1.1* | *piwil1.2* | *piwil1.3* | *piwil1.4* | | *piwil1.1* | *piwil1.2* | *piwil1.3* | *piwil1.4* |

Tilapia — *Oreochromis*
*piwil1.1* — *Diplotaxodon*
— *Rhamphochromis*
— *Astatotilapia calliptera*
— Mbuna
— Utaka
— Deep benthic
Lake Malawi — Shallow benthic

/// *piwil1* expansion
/ Losses or partial losses in lineage

12 long-read genomes        79 short-read genomes

Presence *piwi* gene
Absence *piwi* gene or 3' trailing PiggyBac-1
Presence 3' trailing PiggyBac-1
No long-read genomes for this genus

**D**

20.0

*piwil1.1*
*piwil1.2*
*piwil1.3*
*piwil1.4*
*Onpiwil1*

**E**

Sweep signatures (μ)

*piwil1.1*   *piwil1.2*   *piwil1.3*   *piwil1.4*

Genome position (Mb)
LS419794.1

99th percentile        93rd percentile

chr12        chr20        chr17

A — Lake Malawi (AC)

B — Lake Tanganyika (AB)

C — Lake Victoria (PN)

D — Nile Tilapia (ON)

A

B

C

D

E

F

G

H

I

A
Genome Features ■ Observed ■ Expected    TE classes    ‡Curated annotation    AC‡

B
Genome Features ■ Observed ■ Expected    TE classes    ON

C
Testis AC‡    Ovary AC‡    Muscle AC‡    ‡Curated annotation

D
Testis ON    Ovary ON

E
Expression of active vs inactive TE families in *A. calliptera* gonads

DIRS    hAT-Tip100    PiggyBac    Rex-Babar    TcMar-Tc1    ↑ Active TE families

Mobilization Status
■ Active
■ Inactive

CMC-EnSpm    Maverick    Pao    ↓ Active TE families    AC‡    ‡Curated annotation

F
piRNAs mapping to active vs inactive TE families in *A. calliptera* gonads

DIRS    hAT-Tip100    PiggyBac    Rex-Babar    TcMar-Tc1

CMC-EnSpm    Maverick    Pao

G
CMC-EnSpm-15    CMC-EnSpm-16    Maverick-1    Maverick-3

— AC Muscle    — MZ Testis
— AC Ovary    — TM Testis
— AC Testis

H
CMC-EnSpm    AC‡    ‡Curated annotation

Maverick

Mobilization Status
■ Active
■ Inactive

2-cell    1k cells    Early gastrula    Late gastrula    Early somitogenesis    Mid somitogenesis    Late somitogenesis    Early pharyngula    Mid pharyngula    Late pharyngula    Hatching    Larvae