

1 Identification and quantification of transposable 2 element transcripts using Long-Read RNA-seq in 3 *Drosophila* germline tissues

4
5 Rita Rebollo¹, Pierre Gerenton^{2,3}, Eric Cumunel^{2,3}, Arnaud Mary^{2,3}, François Sabot⁴, Nelly
6 Bulet², Benjamin Gillet⁵, Sandrine Hughes⁵, Daniel S. Oliveira^{2,6}, Clément Goubert⁷, Marie
7 Fablet^{2,8}, Cristina Vieira^{*2,3} and Vincent Lacroix^{*2,3}

8
9 ¹Univ Lyon, INRAE, INSA-Lyon, BF2I, UMR 203, 69621 Villeurbanne, France.

10 ²Université Claude Bernard Lyon 1, Laboratoire de Biométrie et Biologie Evolutive, CNRS, UMR5558,
11 Villeurbanne, Rhône-Alpes, 69100, France.

12 ³ERABLE team, Inria, Lyon Rhone-Alpes, Villeurbanne, France

13 ⁴DIADE unit, Univ Montpellier, Cirad, IRD, F-34394 Montpellier Cedex 5, France

14 ⁵Institut de Génomique Fonctionnelle de Lyon (IGFL), CNRS UMR 5242 , Ecole Normale Supérieure de Lyon,
15 Université Claude Bernard Lyon 1 , F-69007 Lyon, France.

16 ⁶São Paulo State University (Unesp), Institute of Biosciences, Humanities and Exact Sciences, São José do Rio
17 Preto, SP, Brazil.

18 ⁷Human Genetics, McGill University, Montreal, QC, Canada

19 ⁸Institut Universitaire de France (IUF), Paris, Île-de-France , F-75231, France.

20
21 *Corresponding authors: vincent.lacroix@univ-lyon1.fr and cristina.vieira@univ-lyon1.fr
22

23 Abstract

24 Transposable elements (TEs) are repeated DNA sequences potentially able to move throughout the
25 genome. In addition to their inherent mutagenic effects, TEs can disrupt nearby genes by donating
26 their intrinsic regulatory sequences, for instance, promoting the ectopic expression of a cellular gene.
27 TE transcription is therefore not only necessary for TE transposition per se but can also be associated
28 with TE-gene fusion transcripts, and in some cases, be the product of pervasive transcription. Hence,
29 correctly determining the transcription state of a TE copy is essential to apprehend the impact of the
30 TE in the host genome. Methods to identify and quantify TE transcription have mostly relied on short
31 RNA-seq reads to estimate TE expression at the family level while using specific algorithms to
32 discriminate copy-specific transcription. However, assigning short reads to their correct genomic
33 location, and genomic feature is not trivial. Here we retrieved full-length cDNA (TeloPrime, Lexogen)
34 of *Drosophila melanogaster* gonads and sequenced them using Oxford Nanopore Technologies. We
35 show that long-read RNA-seq can be used to identify and quantify transcribed TEs at the copy level.
36 In particular, TE insertions overlapping annotated genes are better estimated using long reads than
37 short reads. Nevertheless, long TE transcripts (> 4.5 kb) are not well captured. Most expressed TE
38 insertions correspond to copies that have lost their ability to transpose, and within a family, only a
39 few copies are indeed expressed. Long-read sequencing also allowed the identification of spliced
40 transcripts for around 105 TE copies. Overall, this first comparison of TEs between testes and ovaries
41 uncovers differences in their transcriptional landscape, at the subclass and insertion level.

42
43 **Keywords:** long-read sequencing, ONT, transposable elements, regulation, RNA-seq, full-length cDNA

44 Introduction

45 Transposable elements (TEs) are widespread DNA sequences that have the ability to move around
46 genomes in a process called transposition (Bourque et al., 2018). TEs can transpose either using an
47 RNA intermediate, in a copy-and-paste mechanism, *i.e.* retrotransposons, or directly through a DNA
48 molecule using different cut-and-paste strategies, *i.e.* DNA transposons. In both cases, the synthesis
49 of a messenger RNA is a fundamental step allowing the production of the transposition machinery,
50 and hence promoting TE replication in the host genome. TE transposition is *per se* a mutational
51 process, and several host mechanisms are in place in order to avoid novel TE insertions, including
52 chromatin remodelling factors, DNA methylation, and small RNAs (Slotkin & Martienssen, 2007). For
53 instance, in *Drosophila melanogaster* ovaries, TEs are the target of piwi-interacting RNAs (piRNAs) that
54 promote TE transcript cleavage, but also deposition of repressive chromatin marks within the TE
55 insertion, blocking any further transcription (Fabry et al., 2021).

56 In order to appreciate the dynamics of TE regulation, an accurate measure of TE expression is
57 required, including copy-specific information (Lanciano & Cristofari, 2020). While such analyses may
58 be easily obtained in genomes composed of mostly ancient TE copies, discrimination of young TE
59 families, such as LINE-1, AluY and SVA in humans, or the study of genomes composed of mostly active
60 copies as seen in many insects, remains a complex feat. Indeed, TE copies belonging to the same TE
61 family have, by definition, more than 80 % of sequence identity, hampering the study of TE regulation
62 and consequently TE expression in a copy-specific manner (Lanciano & Cristofari, 2020). Most
63 genome-wide analyses interested in TE expression, and even their regulation, focus on TE family-level
64 analysis, where short reads are mapped either against TE consensus sequences or to the genome/TE
65 copy sequences followed by grouping of read counts at the family level (TEcount from the TETools
66 package (Lerat et al., 2017), TETranscripts (Jin et al., 2015)). In the past years, many methods have
67 surfaced to take advantage of short-read sequencing datasets and circumvent the multi-mapping
68 problem in order to develop copy-level analysis (for a review see (Lanciano & Cristofari, 2020)). These
69 methods are based on different algorithms that are able to statistically reassign multi-mapped reads
70 to unique locations, for instance with the expectation-maximization algorithm used in TETranscripts
71 (Jin et al., 2015), SQUIRE (Yang et al., 2019) and Telescope (Bendall et al., 2019).

72 In the past years, long-read sequencing has become an attractive alternative to study TE biology.
73 Such reads are able to refine TE annotation (Jiang et al., 2019; Panda & Slotkin, 2020), pinpoint new
74 TE insertions (Mohamed et al., 2020; Rech et al., 2022), determine TE DNA methylation rates at the
75 copy level (Ewing et al., 2020), estimate TE expression (Berrens et al., 2022), and finally, detect TE-
76 gene fusion transcripts (Panda & Slotkin, 2020; Dai et al., 2021; Babarinde et al., 2021). Furthermore,
77 long-read RNA sequencing can not only determine which TE copies are expressed but also discriminate
78 between isoforms of a single TE copy produced by alternative splicing. Indeed, TE alternative
79 transcripts have been described in the very first studies of TEs, using techniques such as northern blot
80 (Belancio et al., 2006), but concomitantly with accessible short-read genome-wide analysis, low
81 interest has been given to TE transcript integrity. Nonetheless, TE isoforms have been shown to
82 participate in TE regulation, as observed for the P element in *D. melanogaster*, where a specific
83 germline isoform encodes a functional transposase protein, while in somatic tissues, another isoform
84 acts as a P element repressor (Laski et al., 1986). The regulation of such tissue-specific splicing has
85 recently been attributed to piRNA-directed heterochromatinization of P element copies (Teixeira et
86 al., 2017). The retrotransposon *Gypsy* also produces two isoforms, including an envelope-encoding
87 infectious germline isoform, also controlled by piRNA-guided repressive chromatin marks (Pélisson et

88 al., 1994; Teixeira et al., 2017). Recently, Panda and Slotkin produced long-read RNA sequencing of
89 *Arabidopsis thaliana* lines with defects in TE regulatory mechanisms (Panda & Slotkin, 2020), and were
90 able to annotate TE transcripts, pinpoint TE splicing isoforms, and most importantly, demonstrate that
91 properly spliced TE transcripts are protected from small RNA degradation.

92 *D. melanogaster* harbours around 12-20% of TE content, and recent studies have suggested that
93 24 TE superfamilies are potentially active (Adrion et al., 2017). Nevertheless, no indication of which
94 copies are active has been documented. Here, we describe a bioinformatics procedure using long-
95 read RNA sequencing, which enables the efficient identification of TE-expressed loci and variation in
96 TE transcript structure and splicing. Furthermore, our procedure is powerful enough to uncover tissue-
97 specific differences, as illustrated by comparing testes and ovaries data.

98

99 Methods

100 Reference genome and annotation

101 The dmgoth101 genome assembly was produced from Oxford Nanopore Technologies (ONT) long-
102 read DNA sequencing and described in (Mohamed et al., 2020). Genome assembly has been deposited
103 in the European Nucleotide Archive (ENA) under accession number PRJEB50024, [assembly](#)
104 [GCA_927717585.1](#). Gene annotation was performed as described in (Fablet et al., 2023). Briefly, [gene](#)
105 [annotation files were retrieved from Flybase \(dmel-all-r6.46.gtf.gz\) along with the matching genome](#)
106 [sequence \(fasta/dmel-all-chromosome-r6.46.fasta.gz\)](#). We then used LiftOff v1.6.1 (Shumate &
107 Salzberg, 2021) with the command `liftOff -g dmel-all-r6.46.gtf -f feature_types.txt -o dmgoth101.txt -u`
108 `unmapped_dmgoth101.txt -dir annotations -flank 0.2 dmgoth101_assembl_chrom.fasta dmel-all-chromosome-`
109 `r6.46.fasta` to lift over gene annotations from the references to the [GCA_927717585.1](#) genome
110 assembly. One should note that `feature_types.txt` is a two line txt file containing 'gene' and 'exon'. In
111 order to locate and count the reads aligned against TE insertions, we produced a GTF file with the
112 position of each TE insertion. We have used RepeatMasker with DFAM dataset from *D. melanogaster*
113 (-species *Drosophila*) TE copies ([Dfam_3.1](#)) and then used OneCodeToFindThemAll (Bailly-Bechet et al.,
114 2014), [downloaded on november 2020](#) to merge LTR and internal parts of a TE into one unique
115 feature `./build_dictionary.pl --rm dmgoth101_assembl_chrom.fasta.out --unknown > dmgoth101_dico` and
116 `./one_code_to_find_them_all.pl --rm dmgoth101_assembl_chrom.fasta.out --ltr dmgoth101_dico --unknown`.
117 Visualization of alignments of TE copies to their consensus sequences were performed using `blastn`
118 (Altschul et al., 1990) with the consensus sequences from the Bergman laboratory that can also be
119 found in [GitLab/te_long_read](#).

120 *Drosophila rearing*

121 *D. melanogaster* dmgoth101 strain was previously described by (Mohamed et al., 2020). Briefly,
122 an isofemale line was derived from a wild-type female *D. melanogaster* from Gotheron, France,
123 sampled in 2014, and sib-mated for 30 generations. Flies were maintained in 12-hour light cycles, and
124 24° C, in vials with nutritive medium, in small-mass cultures with approximately 50 pairs at each
125 generation.

126 Long-read RNA-seq and analysis

127 RNA extraction and library construction

128 Forty-five pairs of ovaries and 62 pairs of testes were dissected in cold PBS 1X from 4 to 8-day-old
129 adults. Total RNA was extracted using the QiagenRNeasy Plus Kit (Qiagen, reference 74104) after
130 homogenization (using a pellet pestle motor) of the tissues. DNA contamination was controlled and
131 removed by DNase treatment for 30 minutes at 37°C (Ambion). Total RNA was visualized in agarose
132 gel to check DNA contamination and RNA integrity before storing at -80°C. The two RNA extracts were
133 quantified with RNA BR reagents on Qubit 4.0 (Thermo Fisher Scientific) and qualified with RNA
134 ScreenTape on TapeStation 4150 instrument (Agilent Technologies), the results showing no limited
135 quantity and a high quality of the RNA molecules (RIN >9.8). We then took advantage of the TeloPrime
136 Full-Length cDNA Amplification kit V2 (Lexogen) in order to enrich ovary and testis total RNA in full-
137 length cDNAs (Figure S1). One should note that the amplified cDNAs are smaller than ~3.5 kb. This
138 protocol is highly selective for mRNAs that are both capped and polyadenylated and allows their
139 amplification. TeloPrime recommends 2 µg total RNA per reaction and we performed two reactions
140 for testis (total of 4 µg) and three reactions for ovaries (total of 6 µg). We determined the optimal PCR
141 cycle number for each sample by quantitative PCR. The quantity and quality of the cDNA produced
142 were checked again with Qubit (dsDNA BR) and TapeStation (D5000 DNA ScreenTape) to confirm the
143 correct amplification of the cDNA and absence of degradation in cDNA fragment length profiles. It is
144 important to note that we do not have replicates for the long-read dataset as the primary goal for this
145 experiment was to evaluate the potential of this technique to identify the largest number of expressed
146 TE copies and isoforms. Enriched full-length cDNAs generated from ovaries and testes were then
147 processed into libraries using the SQK-LSK109 ligation kit (ONT) using 3 µg as starting material. The
148 two libraries were sequenced separately in two flow cells R10 (FLO-MIN110) with a MinION device
149 (Guppy version 2.3.6 and fast basecalling). We obtained 1,236,000 reads for ovaries and 2,925,554 for
150 testes that passed the default quality filter (>Q7). Data are available online at the BioProject
151 PRJNA956863.

152 Mapping

153 The analysis performed here can be replicated through
154 https://gitlab.inria.fr/erable/te_long_read/, a GitLab containing all the scripts along with links and/or
155 methods to retrieve the datasets used. Quality control was performed with NanoPlot v1.41.6 (De
156 Coster et al., 2018). The median read length was 1.18 kb for ovaries and 1.44 kb for testes, the N50
157 read length was 1.7 kb for ovaries and 2.19 kb for testes, and the median quality was 7.7 for ovaries
158 and 8.4 for testes (Table S1, Figure S2). Reads were mapped to the dmgoth101 genome using
159 minimap2 (version 2.26) (Li, 2018) with the splice preset parameter (exact command line given in the
160 GitLab). Most of reads (91.3% for ovaries, 98.8% for testes) could be mapped to the genome (Table
161 S1). Out of those mapped reads, the majority (98.8% for ovaries and 95.1% for testes) mapped to a
162 unique location (*i.e.* had no secondary alignment), and the vast majority (99.9% for ovaries and 97.7%
163 for testes) mapped to a unique best location (*i.e.* in presence of secondary alignments, one alignment
164 has a score strictly higher than the others). Indeed, if a read has several alignments with the same
165 alignment score, then this means the read stems from exact repeats in the genome and they cannot
166 be told apart, hence, one cannot know which copy is transcribed. However, if a read has several

167 alignments with distinct alignment scores, then it means that the read stems from inexact repeats.
168 The presence of this read in the dataset means that one of the copies is transcribed and we consider
169 that it is the one with the highest alignment score. While it could be possible that the read actually
170 stems from the copy with suboptimal alignment, this is highly unlikely because it would mean that
171 there is a sequencing error at the position of the divergence between the two copies of the repeat. A
172 sequencing error in any other position of the read would cause a decrease in the alignment score of
173 both locations. An example of a read that maps to several locations, one with an alignment score
174 larger than the others is given in Figure S3.

175 We also noticed that some reads were only partially mapped to the genome. In practice the query
176 coverage distribution is bimodal (Figure S4), 80% of reads have a query coverage centered on 90%,
177 while the remaining 20% have a query coverage centered on 50%. A thorough inspection of the
178 unmapped regions of these partially mapped reads reveals that they stem from transcripts located
179 elsewhere on the genome. Given that the transcripts covered by the read are themselves fully covered
180 (both the primary locus and the secondary locus), we think that these chimeras are artifactual and
181 were probably generated during ligation steps as previously described (White et al., 2017). Here, we
182 chose to focus on the locus corresponding to the primary alignment and discard the secondary loci. In
183 practice, this corresponds to the longest of the two transcripts. We also ran the same analyses after
184 completely discarding those 20% chimeric reads, but the quantification of TEs is essentially the same
185 ($R=0.992$, Figure S5). In the remainder of the paper, unless stated otherwise, for each read, we focused
186 on the alignment with the best AS score. We discuss the few cases of ties when required.

187 Feature assignment

188 Once a read is assigned to a genomic location, it does not yet mean that it is assigned to a genomic
189 feature. In order to decide which reads could correspond to a TE, we applied the following filters. First,
190 we selected all reads where the mapping location overlaps the annotation of a TE, for at least one
191 base. Then, we discarded all reads that covered less than 10 % of the annotated TE. **It is important to
192 note that no filter is based on the number of basepairs or proportion of the read that extends beyond
193 the TE boundaries (Figure S6).** Finally, in the case where a read mapped to a genomic location where
194 there are several annotated features (a TE and a gene, or two TEs), we assign the read to the feature
195 whose genomic interval has the smallest symmetric difference with the one of the read. The rationale
196 for introducing this filter is best explained with examples. Figure S7 corresponds to a TE annotation
197 overlapping a gene annotation. All reads map to both features, but the gene is fully covered while the
198 TE is only partially covered. We conclude that the gene is expressed, not the TE. Figure S8 corresponds
199 to a genomic location where there are five annotated TEs and a gene, which all overlap. The features
200 that are best covered are Jockey and the gene. **In general, several features may be partially covered
201 by a read and a read may extend beyond each of these features. For each pair read-feature we
202 compute the number of bases that are in the read and not the feature (nr) and the number of bases
203 that are in the feature and not in the read (nf). The sum of these two terms $nr + nf$ is the size of the
204 symmetric difference between the two intervals. We assign the read to the feature with the smallest
205 symmetric difference (Figure 1).** This situation occurs frequently and **assigning a read to a TE only
206 because it covers it yields an overestimation of TE expression (Figure S9 is an example).** The impact of
207 each filter is given in Figure 1. After all filters are applied, there are **1 252 (1 202 uniquely mapping
208 (Table S4) in addition to 50 multi-mapping (Table S6))** reads in ovaries and **8 138 (7 914 uniquely (Table**

209 S3) and 224 multi-mapped (Table S5)) reads in testes that are assigned to TE copies. Our method
210 enables to detect intergenic TEs, intronic TEs and exonic TEs. Counts are summarised in Table S1.

211 Breadth of coverage

212 To calculate the breadth of coverage of annotated transcripts, we mapped reads to the reference
213 transcriptome and computed for each primary alignment the subject coverage and the query
214 coverage. Scripts used are available on the git repository (sam2coverage_V3.py).

215 Gene ontology

216 To identify whether ovary and testis had genes associated with their tissue-specific functions, we
217 firstly selected genes with at least one read aligned in each sample and then we submitted the two
218 gene lists to DAVIDGO separately (Sherman et al., 2022). Due to the high number of biological terms,
219 we selected only the ones with > 100 associated-genes.

221 Subsampling analysis

222 Subsampling of reads was performed using seqtk_sample (Galaxy version 1.3.2) at the European
223 galaxy server (usegalaxy.eu) with default parameters (RNG seed 4) and the fastq datasets. Subsampled
224 reads were then mapped, filtered and counted using the GitLab/te_long_read pipeline.

226 Splicing

227 We mapped reads to both the transcriptome and genomic copies of TEs, we selected the ones
228 whose primary mapping was on a TE. We then filtered those exhibiting Ns in the CIGAR strings. Those
229 are the reads aligning to TEs with gaps. We then extracted the dinucleotides flanking the gap on the
230 reference sequence. Scripts used are available on the git repository (SplicingAnalysis.py,
231 splicing_analysis.sh)

232 Short-read RNA-seq and analysis

233 RNA extraction and short-read sequencing were retrieved either from (Fablet et al., 2023), at the
234 NCBI BioProject database PRJNA795668 (SRX13669659 and SRX13669658), or performed here and
235 available at BioProject PRJNA981353 (SRX20759708, SRX20759707). Briefly, RNA was extracted from
236 70 testes and 30 ovaries from adults aged three to five days, using RNeasy Plus (Qiagen) kit following
237 the manufacturer's instructions. After DNase treatment (Ambion), libraries were constructed from
238 mRNA using the Illumina TruSeq RNA Sample Prep Kit following the manufacturer's recommendations.
239 Quality control was performed using an Agilent Bioanalyzer. Libraries were sequenced on Illumina
240 HiSeq 3000 with paired-end 150 nt reads. Short-read analysis was performed using Tetrascripts (Jin
241 et al., 2015) at the family level, and SQUIRE (Yang et al., 2019) was used for mapping and counting TE
242 copy-specific expression. A detailed protocol on SQUIRE usage in non-model species can be found
243 here <https://hackmd.io/@unleash/squireNonModel>. Family-level differential expression analysis was
244 performed with TE transcript (Jin et al., 2015). RNA-seq reads were first aligned to the reference
245 genome (GCA_927717585.1) with STAR (Dobin et al., 2013): the genome index was generated with
246 the options --sjdbOverhang 100 and --genomeSAindexNbases 12; next, alignments were performed for

247 each read set with the parameters `-sjdbOverhang 100 --winAnchorMultimapNmax 200 and --`
248 `outFilterMultimapNmax 100` as indicated by the authors of TE transcript (Jin & Hammell, 2018). TE
249 transcript was ran in two distinct modes, using either multi-mapper reads (`--mode multi`) or only using
250 single mapper reads (`--mode uniq`) and the following parameters: `--minread 1 -i 10 --padj 0.05 --sortByPos`.

251 Results and discussion

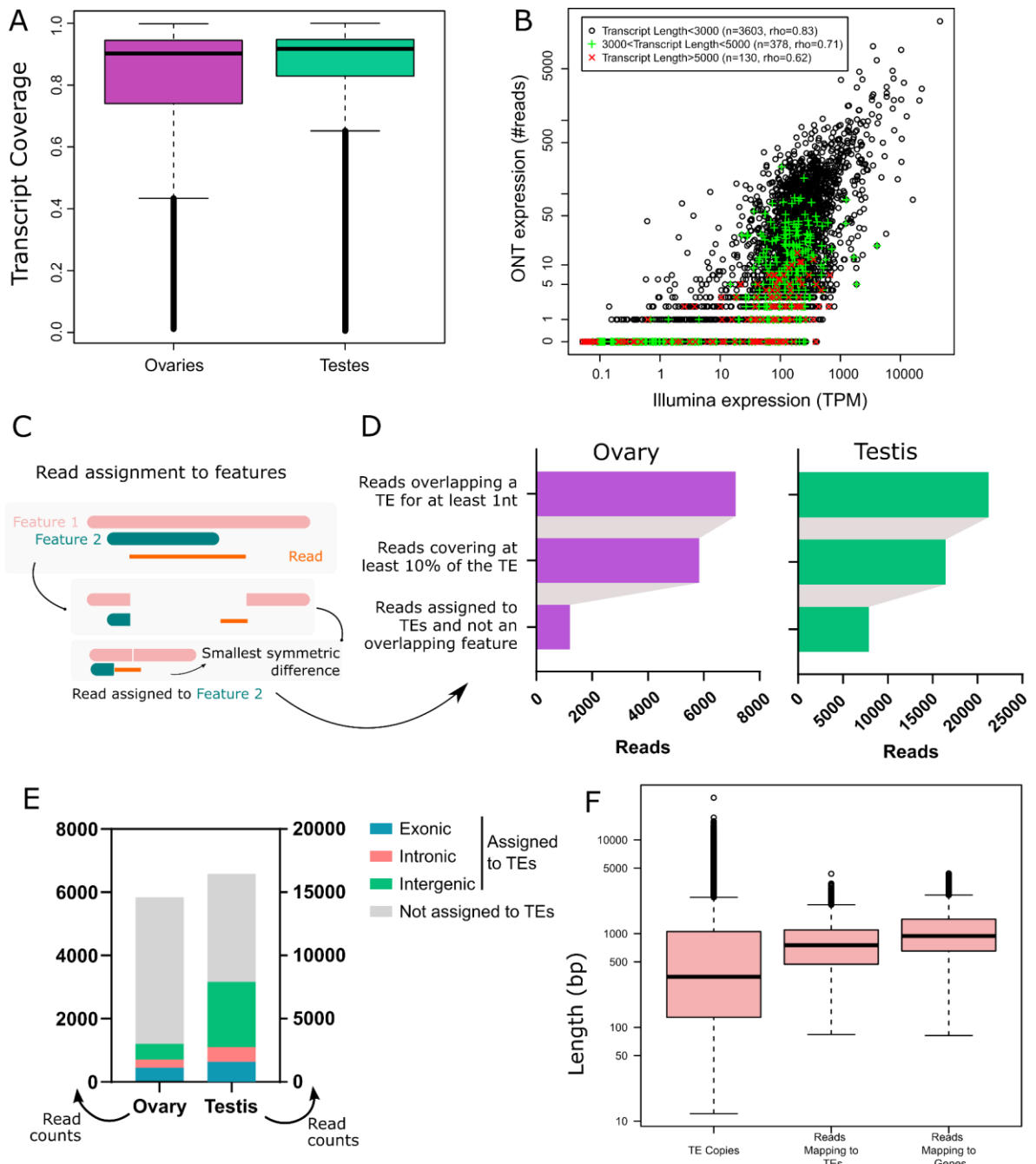
252 Transposable element transcripts are successfully detected with long-read RNA-seq

253 In order to understand the TE copy transcriptional activity and transcript isoforms in gonads of *D.*
254 *melanogaster*, we extracted total RNA from ovaries and testes of dmgoth101 adults, a French wild-
255 derived strain previously described (Mohamed et al., 2020). Prior to long-read sequencing, we
256 enriched the total RNA fraction into both capped and polyadenylated mRNAs in order to select mature
257 mRNAs potentially associated with TE activity (see material and methods for the details on the
258 TeloPrime approach). Sequencing yielded between ~1 to ~3 million reads per tissue, ranging from 104
259 to 12,584 bp (median read length ~1.4 Kb, Figure S1-2, Table S1). Reads were subsequently mapped
260 to the strain-specific genome assembly (Mohamed et al., 2020) using the LR aligner Minimap2 (version
261 2.26) (Li, 2018). Most reads mapped to the genome (91.3% for ovaries and 98.8% for testes, Table S1),
262 and the majority of them mapped to a unique location (*i.e.* had no secondary alignment, 98.8% for
263 ovaries and 95.1% for testes), and the vast majority mapped to a unique best location (*i.e.* presence
264 of secondary alignments, one alignment has a score strictly higher than the others, see Methods,
265 99.9% for ovaries and 97.7% for testes).

266 In order to validate the long-read RNA-seq approach, we first determined the breadth of coverage
267 of all expressed transcripts and showed that the majority harbour at least one read covering more
268 than 80% of their sequence (70.2% in ovaries and 71.8% in testes). Only a few reads correspond to
269 partially covered transcripts, as most reads cover more than 80% of the transcript sequence (63.4% in
270 ovaries, 77.4% in testes - Figure 1A), although very long transcripts (≥ 5 kb) are poorly covered (Figures
271 S10-11). The transcriptomes obtained are enriched in typical germline ontology terms, such as
272 “spermatogenesis” for testes, and “oogenesis” for ovaries (Figure S12). Finally, while the first version
273 of the TeloPrime protocol could not be used for quantification (Sessegolo et al., 2019), the
274 quantifications obtained here correlate well with available short-read sequencing ($\rho=0.78$, $R=0.44$,
275 Figure 1B). We also noticed that the correlation between the two technologies is weaker for very long
276 transcripts.

277 Although most long reads map to a unique location on the genome, ensuring that a read stems
278 from the TE copy is not straightforward. First, a read may overlap with a TE copy only for a few
279 nucleotides, suggesting the read is not a consequence of TE transcription. To rule out these cases, all
280 reads covering less than 10% of the TE copy were discarded. Second, in many instances, the read
281 overlapped both a gene and one or several TEs. In these cases, the feature for which the coverage was
282 best (see Methods, Figure 1C and Figures S7-8) was selected. Overall, after applying these filters, 1
283 252 reads in ovaries and 8 138 reads in testes were assigned to TE copies (Table S1, Figure 1D). Out of
284 these, 37% overlap exons, 21% overlap genes but not exons and 41% do not overlap genes in ovaries.
285 In testes, 20% are exonic, 14% intronic and 65% intergenic (Figure 1E). Additionally, a large fraction of
286 reads overlapping TEs and genes are assigned to genes (52% in testes and 79% in ovaries) (Figure 1D,
287 1E).

288 To ensure this long-read dataset is able to recover transcripts encompassing all TE copy lengths
289 present in the genome, we compared the length distribution of all TE insertions with the length of all
290 mapped reads (Figure 1F). While genomic TE copies range from a few base pairs to ~15 Kb, 75% are
291 smaller than 2 Kb. The average length of reads mapping to TEs encompasses the majority of TE copies,
292 but does not cover TE transcripts longer than 4.5 Kb. Reads mapping to genes have a similar
293 distribution (Figure 1F). The absence of very long reads indicates that either very long mRNAs are
294 absent from the sample (as suggested by the cDNA profile, Figure S1), or the TeloPrime technique is
295 not well tailored for capturing very long transcripts. In order to clarify this point, we compared the
296 quantification obtained by Illumina and ONT TeloPrime for short (<3 Kb), long (3 Kb-5 Kb) and very
297 long transcripts (>5 Kb), and obtained the following Spearman correlations of 0.83 (n=3 603 genes),
298 0.71 (n=378) and 0.62 (n=130), respectively (Figure 1B for ovaries, Figure S13 for testes). Furthermore,
299 reads covering very long annotated transcripts (>5kb) tend to be partial (Figure S10, S11, S14). We
300 conclude that, although very long transcripts are rare (<0.1% of reads), the TeloPrime protocol could
301 underestimate their presence.



302

303

304

305

306

307

308

309

310

311

312

313

314

Figure 1: Long-read RNA-seq of *Drosophila melanogaster* ovaries and testes. **A.** Transcript coverage by long-read RNA-seq in ovaries and testes. **B.** Gene expression quantification using Illumina and ONT sequencing in ovaries. Each dot is a gene with a single annotated isoform. Transcripts longer than 5 kb tend to be undersampled using TeloPrime. **C.** Read assignment to features. In the case where a read aligns to a genomic location where two features are annotated, the read is assigned to the feature with the best coverage. Two dimensions are considered. The read should be well covered by the feature, and the feature should be well covered by the read. In practice, we calculate the symmetric difference for each read/feature and select the smallest. In this example, the read is assigned to Feature 2. **D.** Impact of filters on the number of reads assigned to TEs. **E.** Number of reads mapping to TEs separated into three categories (intronic, exonic or intergenic), and reads that have not been assigned to TE copies. **F.** TE copy and read length distribution. Reads mapping to TEs encompass most TE copy length but lack transcripts longer than 4.5 Kb, as also observed for reads mapping to genes.

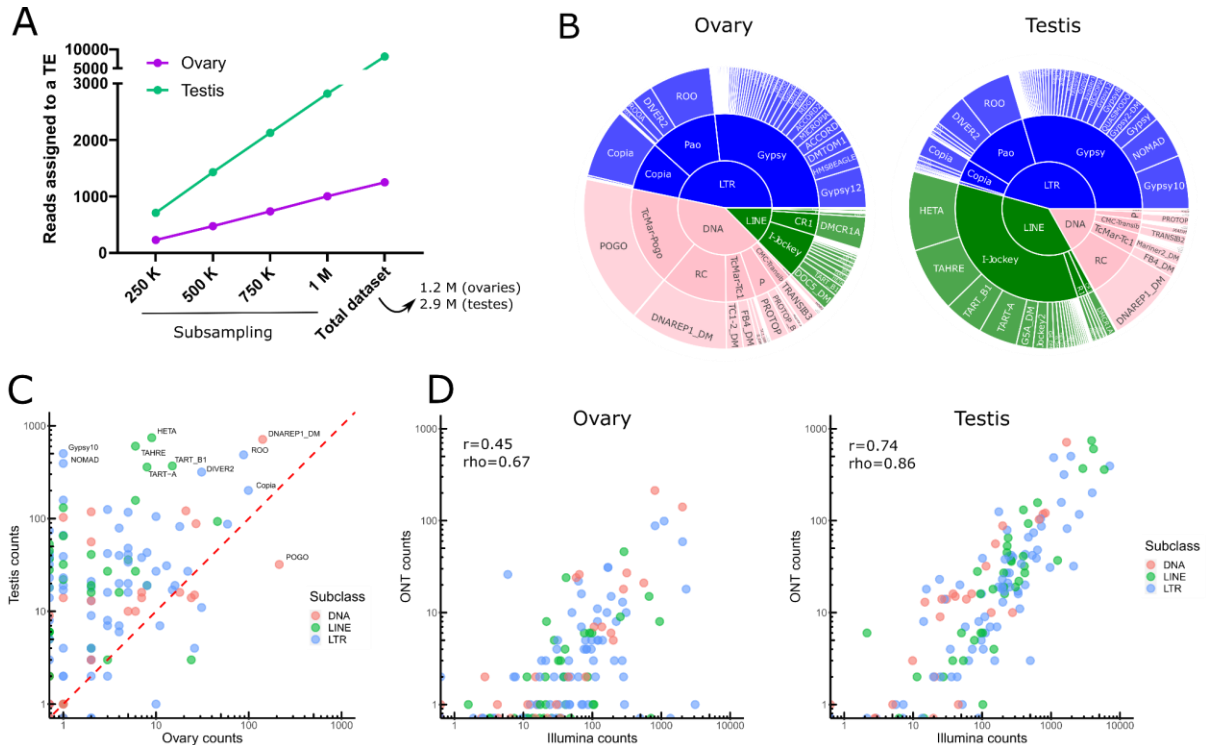
315

316 TE mRNA landscape is sex-specific

317 Taking into account all the filtering steps, only 0.28% (8 138/2 925 554) and 0.10% (1 252/1 236
318 000) of long reads aligned to TE copies in testes and ovaries respectively (Table S1). Given the
319 differences in sequencing depth between both tissues, we have computed the number of reads
320 assigned to TEs based on different sets of subsampled reads, and show that TE reads are more
321 abundant in testes than in ovaries (Figure 2A). We identified 130 TE families capable of producing
322 capped-polyadenylated mRNAs (Table S2), of which 70 belong to Long-terminal repeat (LTR) elements
323 (retrotransposons that possess LTR sequences surrounding a retroviral-like machinery). Despite the
324 high number of shared transcribed TE families (96/130), the transcriptional landscape between
325 ovaries and testes is quite different (Figure 2B for the complete dataset and Figure S15 for a
326 subsampled dataset). While LTR elements dominate the transcriptional landscape in both tissues, LINE
327 elements are the second most transcribed TE subclass in testes, while in ovaries, DNA families harbor
328 more read counts (Figure 2B). The transcriptional landscape within TE subclasses between tissues is
329 very similar for retrotransposons, with Gypsy and I-Jockey superfamilies being the most expressed LTR
330 and LINE elements respectively. However, the DNA subclass transcriptional landscape is quite
331 different between testes and ovaries: TcMar-Pogo is the most expressed DNA superfamily in ovaries,
332 while RC elements are abundantly transcribed in testes.

333 Globally, TE families show higher long-read counts in males compared to females (Figure 2C), not
334 only because male samples were more deeply sequenced (2.3 times more), but also because the
335 proportion of reads that map to TEs is higher in males even when subsampling the same number of
336 reads between tissues (Figure S16 for a subsampled dataset). *HETA* (I-Jockey) and *DNAREP1_D* (RC)
337 are the top two families in male TE transcript counts, with 743 and 713 long-read transcripts
338 respectively. In females, *pogo* (TcMar-pogo) is the TE family amounting the most transcripts with 213
339 long reads (while only 32 in males), followed by *DNAREP1_D* (RC) with 141, and *Copia* (Copia) with 99
340 long reads (but both families with higher transcript counts in males, Table S2). There are only five TE
341 families that yielded long-reads in ovaries and not in testes, *BARI_DM* (TcMar-Tc1 - DNA), *Gypsy7*
342 (Gypsy - LTR), *Gypsy11* (Gypsy - LTR), *Copia2_LTR_DM* (Copia - LTR) and *Helena_RT* (I-Jockey - LINE),
343 but they all harbor only one or two long reads suggesting their expression is low. There are 29 families
344 detected only in testes, four DNA elements (*Transib-N1_DM* (CMC-Transib), *NOF_FB* (MULE-NOF), and
345 two TcMar-Tc1 (*Bari1* and *Minos*)), 13 LINE elements (10 I-Jockey, two R1, and one R1-LOA, see Table
346 S2 for details), and 12 LTR families (two Copia, four Gypsy, two Pao and four unknown families),
347 ranging from one to 74 long reads per TE family. Finally, eleven TE families show no long-read mapping
348 in either tissue. Collectively, long-read sequencing is able to discriminate between ovaries and testes
349 TE transcriptional landscapes.

350 Short-read RNA sequencing of ovaries and testes, followed by estimation of TE family expression
351 with TEtranscripts (Jin et al., 2015) - TE expression estimation per TE family, see material and methods
352 for more information) recapitulates the long-read RNA sequencing profiles (Figure 2D and Figure S17).
353 Despite the fact that TE transcripts are overall poorly expressed, the estimation of their expression
354 level is reproducible across technologies. The correlation is higher for testes ($r=0.74$, $\rho=0.86$) than
355 for ovaries ($r=0.45$, $\rho=0.67$), where the coverage is weaker. Indeed, as previously stated, the total
356 contribution of TE to the transcriptome is weaker for ovaries and the sequencing shallower.



358

359

360

361

362

363

364

365

366

367

368

369

370 Long-read sequencing successfully retrieves single-copy transcripts

371

372

373

374

375

376

377

378

379

380

Figure 2: Transposable element transcriptional landscape in ovaries and testes of *Drosophila melanogaster*. **A.** Reads assigned to TEs are more abundant in testis. Subsampling of reads from 250 000 to 1 million reads, along with the complete dataset, show a higher number of reads assigned to TEs in testes than in ovaries. **B.** Global TE transcriptional landscape using ONT long-read sequencing. The outer ring, middle ring and inside circle represent TE family, superfamily and subclass respectively. The area in the circle is proportional to the rate of expression. **C.** Comparison of TE expression ratio between testes and ovaries ONT long-read datasets. **D.** Comparison between Illumina and ONT datasets for estimating the expression levels at the TE family level. Each dot is a TE family. TEtranscript is used for short reads. The correlation between quantifications by both technologies is higher for testes ($\rho=0.86$) than for ovaries, ($\rho=0.67$) where the coverage is weaker. In both cases, it is compatible with what is observed for genes.

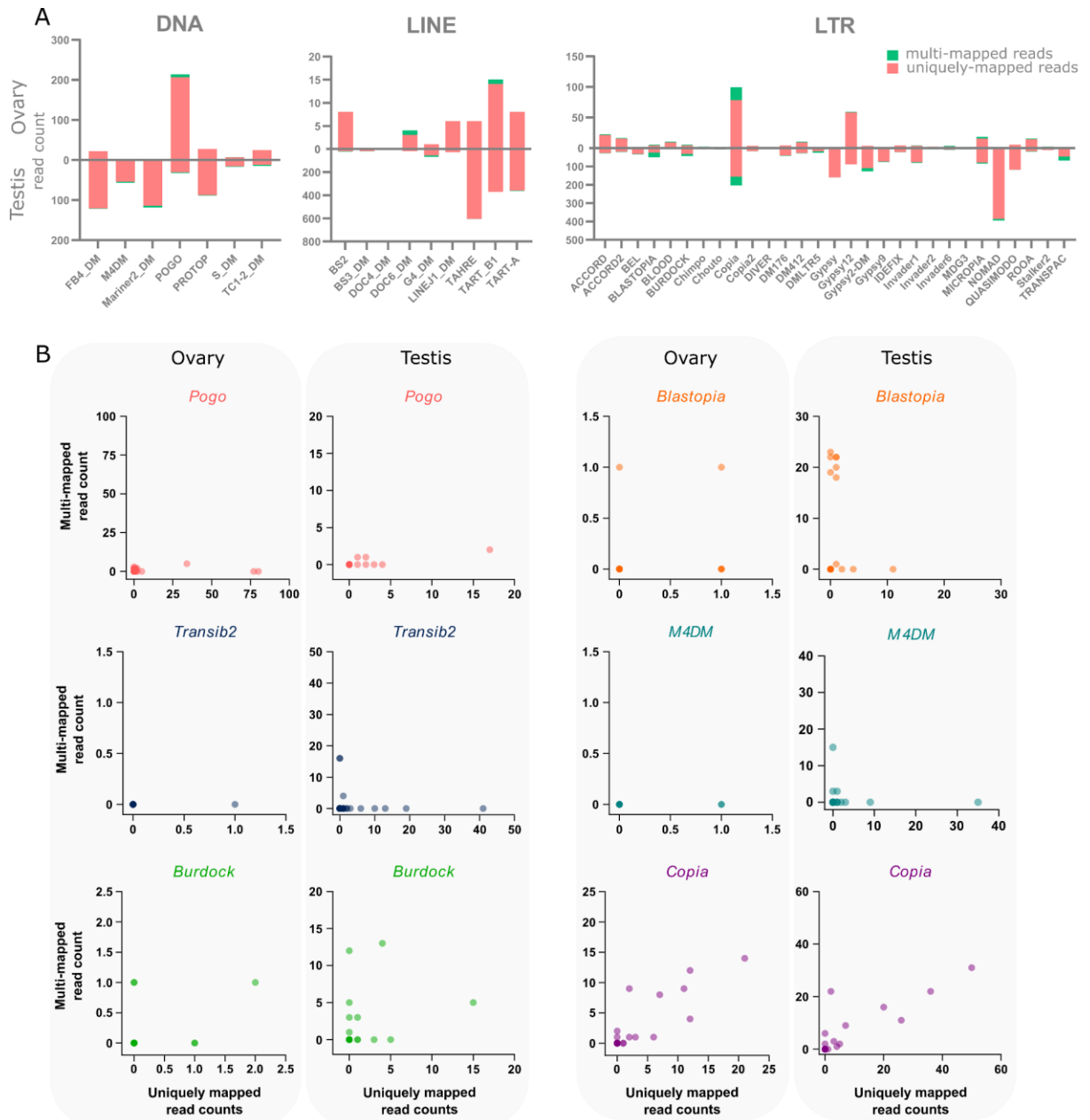
381

382

In ovaries, out of 101 TE families detected (at least one read), there are 16 families that harbor only one multi-mapped read, and three families with 3 to 21 multimapped reads, *copia*, *pogo* and

383 *micropia* (Figure 3A). While only 3% of *pogo* reads are multimapped in ovaries (7 out of 213 reads),
384 *micropia* and *copia* harbor a higher percentage of multimapped reads, 17% of 18 reads for *micropia*
385 and 21% of 99 reads for *copia*. In testis, 125 TE families are expressed (at least one read), and 39 of
386 them have multi-mapped reads (Figure 3A). As observed in the ovary dataset, the number of
387 multimapped reads is low for most families, with only seven families harboring more than 10
388 multimapped reads. While *copia* harbors the most multimapped reads in testis (47 out of 201 long-
389 reads), *blastopia* shows a higher multi-mapped ratio with 54% of reads multimapped (26 out of 48
390 long-reads). *Transpac* also shows an important number of multi-mapped reads with 22 out of 66 reads
391 mapped. In total, 45 TE families have both uniquely and multi-mapping reads in ovaries and testis
392 (Figure 3A).

393 We uncovered 404 and 1 078 TE copies harbouring at least one long-read unambiguously in ovaries
394 and testes respectively. When taking into account multi-mapped reads, an additional 53 and 94 TE
395 copies are potentially expressed in each tissue (Table S2). However, it is important to note that the
396 number of assigned multi-mapped reads to each copy is quite small. For instance, in ovaries, 47 of
397 these potentially expressed copies only harbor one multi-mapped read, five copies harbor two, and a
398 single *pogo* copy harbors three multi-mapping reads. As a comparison, the two most expressed copies
399 in ovaries are two *Pogo* copies, POGO\$3L_RaGOO\$9733928\$9735150 and
400 POGO\$X_RaGOO\$21863530\$21864880, with 80 and 77 uniquely mapping reads, and no multi-
401 mapping read (Table S4). In testis, out of 94 copies, 71 have only a single multi-mapped read and only
402 eight copies show more than five multi-mapped reads. Among these copies, there are three *Blastopia*
403 copies with 23, 22 and 19 multi-mapped reads, two *Transib2* copies (16 reads each), one *M4DM* (15
404 reads), one *Burdock* (12 reads) and finally one *Copia* (six reads) (Table S3). As a comparison, the top
405 expressed copy in testis is a *Gypsy10* (Gypsy10\$3R_RaGOO\$761442\$762629) with 493 uniquely
406 mapping reads and no multimapping ones. If searching for the most expressed copy among the same
407 TE families as described having only abundant multi-mapped reads in testis, *Blastopia* harbors only 11
408 uniquely mapping reads, *Transib2* has 41 uniquely mapped reads and no multimapped ones, *M4DM*
409 copy has 35, and *Burdock* shows 15 uniquely mapping reads and five multimapped ones (Figure 3B).
410 Therefore, despite a few exceptions (*blastopia*, *transib2*, *M4DM* and *Burdock*), long-read sequencing
411 can indeed identify single-copy TE transcripts.



412

413 **Figure 3: Multi-mapping and uniquely mapping ONT reads. A.** Distribution of uniquely and multimapped
 414 reads across TE families in ovaries and testes (only TE families harboring at least one multimapped read are
 415 shown). **B.** Association between multi-mapped and uniquely mapped reads at the copy-level for the TE families
 416 showing copies harboring abundant multi-mapped reads.

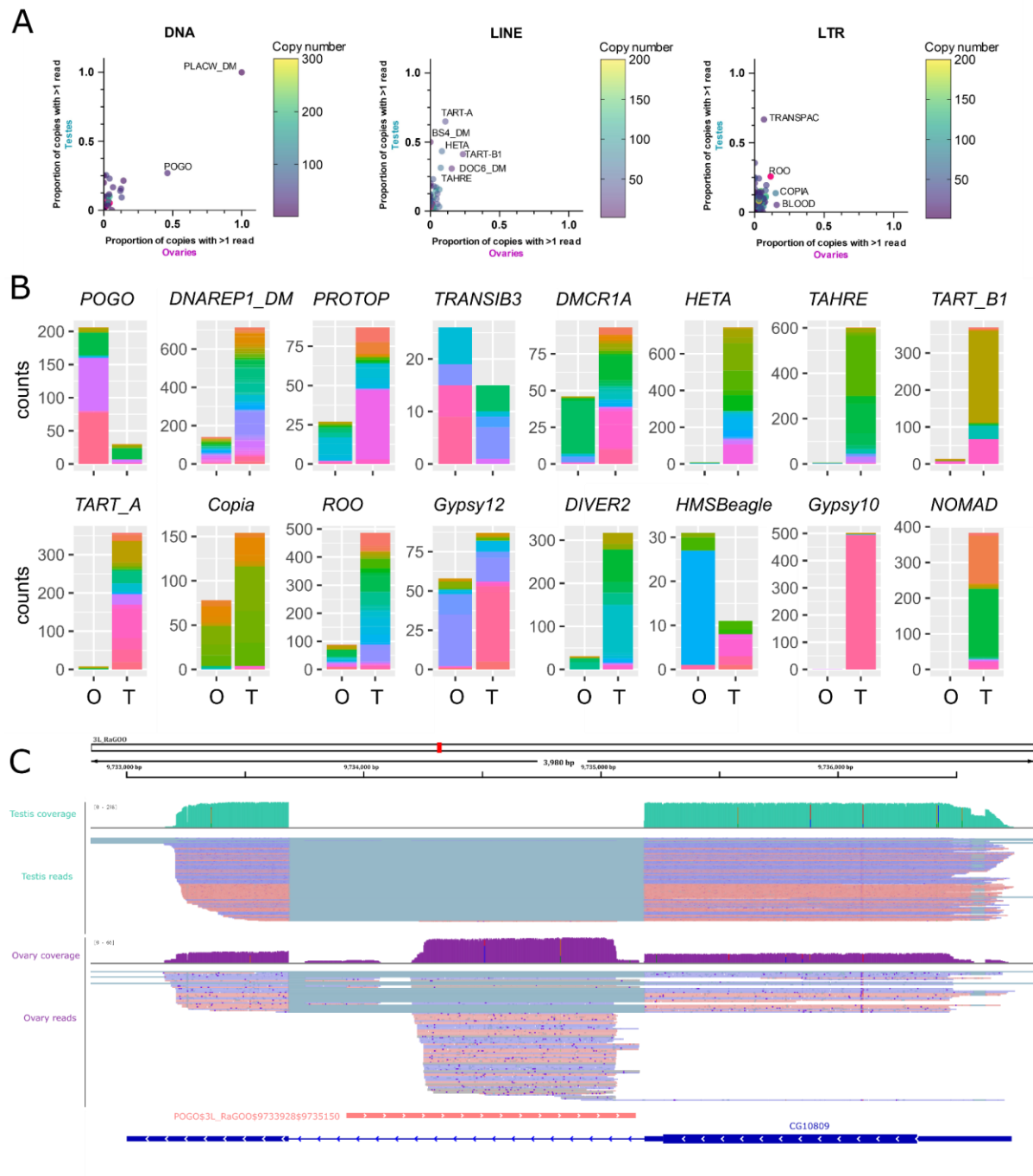
417 Within a TE family, the contribution of each TE copy to the family transcriptional activity is variable.
 418 In general, only a few insertions produce transcripts, even if taking into account multi-mapped reads
 419 (Figure 4A for uniquely mapping reads and Figure S18 for all reads). However, *Transpac* (LTR, Gypsy)
 420 copies are nearly all expressed in testes (10 expressed copies and two potentially expressed copies
 421 out of 15), while in ovaries, *pogo* (DNA, TC-mar-Pogo) harbors 12 copies producing transcripts, and
 422 five potentially expressed copies out of 26. *DNArep1* (RC) is the most abundant TE family in the *D.*
 423 *melanogaster* genome and is also the family harbouring the most transcribed copies in both ovaries
 424 and testes (72 and 170 respectively out of 2 555 copies). In ovaries, out of the 404 insertions with at

425 least one mapped read, 23 had more than 10 mapped reads. In testes, out of the 1 078 insertions with
426 at least one uniquely mapped read, 157 had in fact more than 10 mapped reads.

427 While many TE copies within a family are indeed able to produce transcripts, there are significant
428 differences in copy expression rate (Figure 4B for the 10 most expressed TE families in ovaries and
429 testes, and Figure S19 for a subsampled dataset). For instance, *Gypsy10* (LTR - Gypsy) harbors a highly
430 active copy with 493 uniquely mapping reads out of 502 total counts in testes, while only one read is
431 detected in ovaries. As a contrast, *DNAREP1* (DNA, RC) and *Roo* (LTR, Pao) have several copies that
432 contribute to the TE family expression (Figure 4B). Finally, some families show copies transcriptionally
433 active in both ovaries and testes, as for instance *Copia* (LTR, Copia).

434 In ovaries, where *pogo* has the highest number of long reads (213), an insertion of 1 222 bp in the
435 3L chromosome (POGO\$3L_RaGOO\$9733928\$9735150) accumulates nearly 37% of the family total
436 read count (Figure 4B and C). This specific *pogo* insertion is located in the intron of the CG10809 gene.
437 The same pattern is observed for the second most expressed *pogo* insertion
438 (POGO\$X_RaGOO\$21863530\$21864880), also located in the intron of a gene (CG12061), expressed
439 in testes and not in ovaries (Figure S20). CG12061 is a potential calcium exchange transmembrane
440 protein and has been previously shown to be highly expressed in the male germline (Li et al., 2022).
441 Indeed, using long-read sequencing, CG12061 is highly expressed in testes compared to ovaries, and
442 curiously, the intronic *pogo* insertion is only expressed when the gene is silent (in ovaries). The other
443 expressed insertions of *pogo* are located in intergenic regions (Figure S21).

444



445

446

447

448

449

450

451

452

453

454

455

456

Figure 4: Transcription of transposable element copies. **A.** Frequency of transcribed copies (read > 1) within TE subfamilies in ovaries and testes, along with genomic copy number (color bar, 1 to 200 (LINE/LTR) or 300 (DNA) copies). All TE families harbouring more than 200 (LINE/LTR) or 300 (DNA) copies are depicted in pink. For DNA elements, *DNAREP1_DM* harbours 2 555 copies, *Protop* 305 and *ProtopA* 347. For LINE families, *DMCR1A* has 583 copies and *FW2_DM* 216. The LTR families, *idefix* (227) and *roo* (218) are also depicted in pink. Most TE subfamilies have only a couple of copies producing transcripts, while the majority of HETA copies are expressed in testes for instance (middle panel). **B.** Distribution of read counts per copy for the 10 most expressed copies in ovaries and testes (16 TE families total), showing the overall expression of specific copies within a TE family (Table S3 and S4). Copies are represented by different colors within the stacked bar graph. O: ovaries, T: testes **C.** IGV screenshot of a *pogo* copy (*POGO\$3L_RaGOO\$9733928\$9735150*, in pink). In green, testis coverage and below mapped reads, in purple the same information for ovaries. Dmgoth101

457 repeat and gene tracks are also shown and more information on the annotation can be seen in the material
458 and methods section.

459
460 Finally, using short-read sequencing and a tool developed to estimate single-copy expression
461 (Squire (Yang et al., 2019)), we compared the overall TE copy transcriptional landscape between short
462 and long reads (Figure S20). There was a poor correlation with the ONT estimations ($\rho=0.23$, $r=0.18$
463 for ovaries and $\rho=0.37$, $r=0.32$ for testes). At the family level, the quantifications obtained by Squire
464 were comparable to the ones obtained with long-reads ($\rho=0.66$, $r=0.49$ for ovaries and $\rho=0.77$,
465 $r=0.34$ for testes, Figure S22). Examination of instances where the two techniques differed most,
466 shows that Squire tends to overestimate the expression of TE insertions completely included in genes
467 (Figure S9). Indeed, while long-reads can easily be assigned to the correct feature because they map
468 from the start to the end of the feature, many of the short-reads originating from the gene also map
469 within the boundaries of the TE. Methods based on short reads could clearly be improved, based on
470 the study on such instances where there is a discordance.

471 Transcripts from full-length transposable element copies are rarely detected

472 Many insertions produce transcripts that are shorter than the annotated TE and are likely unable
473 to participate in TE transpositional activity. Furthermore, even in the case where the transcript fully
474 covers the insertion, the copy itself might have accumulated mutations, insertions and deletions
475 making it unable to transpose. To assess this, we computed the query coverage of the reads with
476 regards to the insertion they correspond to. We find that one-third of the insertions have at least 80%
477 query coverage (Figure 5A). However, out of these insertions, only a few of them are close in length
478 to a functional full-length sequence. In order to search for potentially functional, expressed copies,
479 we filtered for copies with at least five long-reads detected, and covering at least 80% of their
480 consensus sequences. In ovaries these filters correspond to seven insertions: one *pogo*, five *Copia*
481 insertions and one *MAX*, and in testes, there are nine potentially functional insertions, five *Copia*
482 insertions, three *Mariner2*, and one *DM1731*. While all *Copia* insertions expressed with at least five
483 reads are full-length, other TE families show mostly internally deleted expressed copies (Figure 5B).
484 Indeed, a closer analysis of *pogo*, the most expressed TE subfamily in ovaries, shows only one full-
485 length copy expressed (POGO\$2L_RaGOO\$2955877\$2958004), but at low levels (five reads in ovaries,
486 and two in testes). Instead, the other three expressed *pogo* copies with at least five reads in ovaries
487 (80, 77 and 34), are internally deleted (Figure 5B). Hence, ONT long-reads detects only a small number
488 of expressed full-length copies.

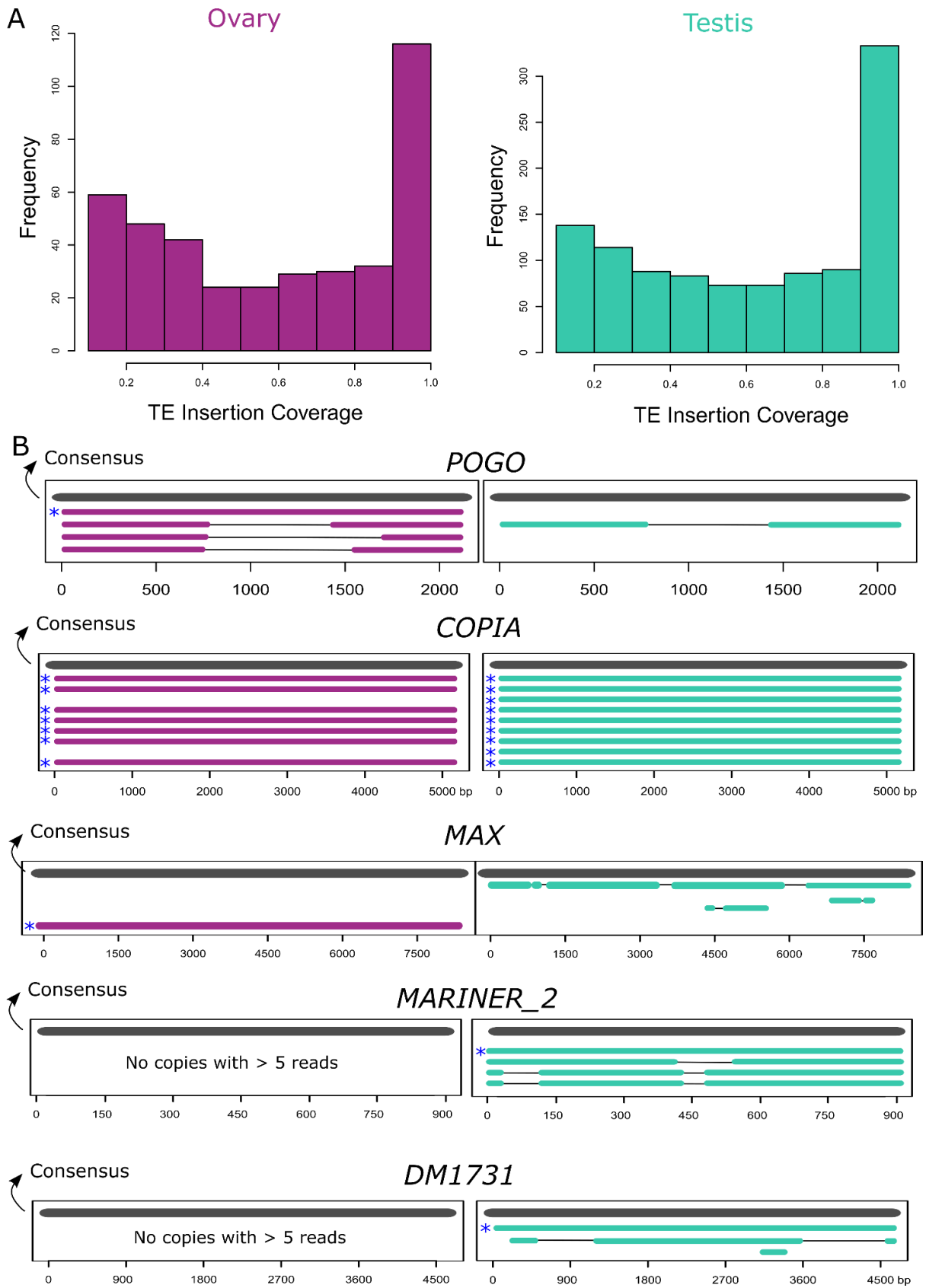
489

490

491

492

493



494

495

496

497

Figure 5. TE transcripts stem mostly from deleted or truncated copies. **A.** Coverage of ONT reads on TE insertions. One-third of copies are covered for at least 80% of their length. **B.** Alignment of copies belonging to the five TE families where at least one full-length expressed copy (80% of consensus) was observed with more

498 than five long-reads. All copies represented have at least five long-reads. Consensus sequences are
499 represented in grey and copies are either purple for ovaries or green for testes. Asterisks depict the full-length
500 copies.

501 Long-read sequencing unveils novel spliced TE isoforms

502 A closer analysis of the reads stemming from the detected full-length copies shows that many of
503 them do not cover the copies completely (Figure 6). For instance, five *Copia* copies are at least ~80%
504 of the consensus sequences and have at least five long-reads detected (Figure 5B), however, although
505 the reads map from the 5' end to the 3' end of the copy, they map with a gap (Figure 6 and Figure S23-
506 S25). *Pogo*, *Max* and *DM1731* also show such gapped alignments. Inspection of these gaps reveals
507 that they are flanked by GT-AG consensus, suggesting that those transcripts are spliced. Only *Mariner2*
508 shows reads that correspond to the full-length copy, but one should note that the consensus sequence
509 of *Mariner2* is smaller than 1 Kb, while the other elements are much longer. As stated before, very
510 few cDNA molecules longer than ~4 Kb have been sequenced (Figure S1), suggesting either that such
511 longer transcripts are rare, and/or that the method used here for cDNA amplification induces a bias
512 towards smaller sequenced fragments. Collectively, long-read sequencing shows that despite the
513 presence of potentially functional, full-length copies in the *D. melanogaster* genome, only a few of
514 these are detected as expressed in testes and ovaries, and the reads that are indeed recovered seem
515 to be spliced.

516

517



518

519

520

521

522

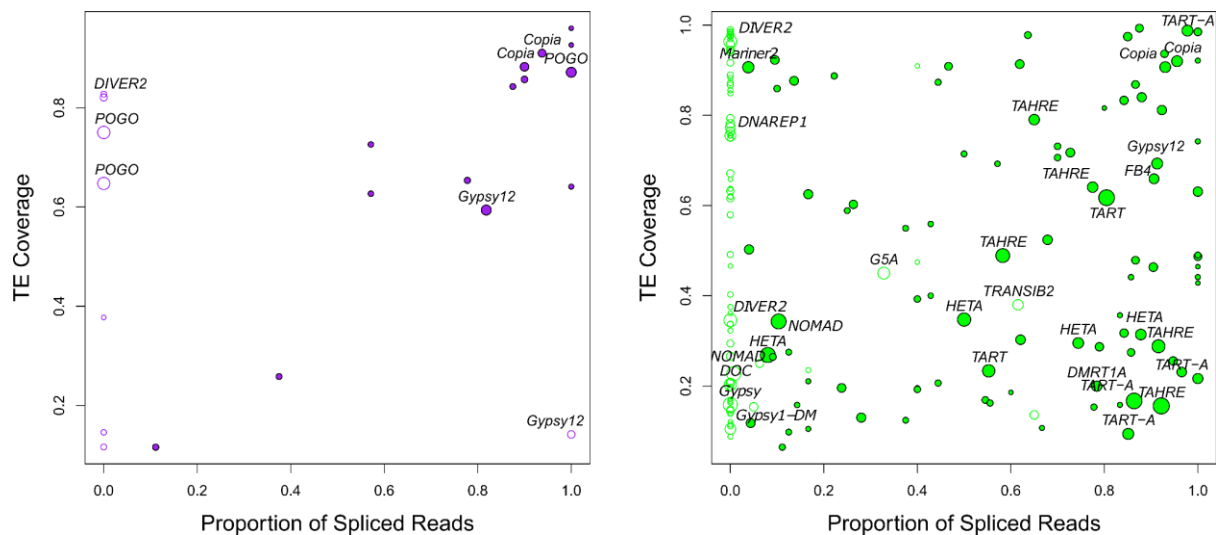
523

524

Figure 6. Full-length copies produce spliced transcripts. IGV screenshot of uniquely mapping reads against putative full-length copies (copies > 80% of the consensus sequence length) harboring at least five reads (see Figure 5B). Only *copia* has multi-mapped reads that can be appreciated in Figure S23. Dmgoth101 repeat and gene tracks are also shown and more information on the annotation can be seen in the material and methods section. Ovary and/or testis coverage and reads are shown below the TE copies.

525 While most TEs do not harbor introns, there are a couple of [exceptions](#) previously described in *D.*
 526 *melanogaster*. Indeed, P elements are known to be regulated tissue-specifically by alternative splicing
 527 mechanisms, involving piRNA targeting (Laski et al., 1986; Teixeira et al., 2017). *Gypsy* copies are able
 528 to produce ENV proteins through mRNA alternative splicing (Pélisson et al., 1994; Teixeira et al., 2017).
 529 [As with P elements, gypsy splicing is also thought to be regulated by piRNAs.](#) Finally, *Copia* elements
 530 produce two isoforms, a 5 Kb and a 2.1 Kb (which is a spliced product of the 5 Kb mRNA) (Miller et al.,
 531 1989; Yoshioka et al., 1990). The 2.1 Kb encodes the GAG protein and is produced at higher levels than
 532 the other proteins (Brierley & Flavell, 1990). While the shorter transcript can be processed by *Copia*
 533 reverse transcriptase, the 5 Kb full-length isoform is clearly preferred (Yoshioka et al., 1990). Most of
 534 these discoveries were obtained through RT-PCR sequencing of amplicons, or recently, through short-
 535 read mapping. Nevertheless, systematic analysis of TE alternative splicing in *D. melanogaster* is
 536 lacking, due to the difficulty of detecting such isoforms from short-read data. Here we used long-read
 537 sequencing to mine for such splicing isoforms. We searched for reads harboring a gap compared to
 538 the reference sequence (presence of N's in the CIGAR string). In order to ensure that those gaps
 539 corresponded to introns, we searched for flanking GT-AG splice sites ([see methods, and Figures S26-](#)
 540 [S29](#)). In ovaries, out of 22 insertions supported by at least 5 reads, 15 exhibited at least one gapped
 541 read (Figure 7). For all tested insertions, the majority of gapped reads exhibited a GT-AG consensus,
 542 except for one insertion (it was CT-TA). In testes, out of 163 insertions supported by at least 5 reads,
 543 100 exhibited at least one gapped read, 91 with a GT-AG consensus (Figure 7). Out of the 9 others, 6
 544 exhibited only one of two gapped reads, the 3 remaining ones with a CT-AT, GA-CG and AT-AG
 545 consensus. Those could correspond to non-canonical splicing. They could also correspond to a
 546 heterozygous deletion or to the expression of a deleted copy located in a non-assembled part of the
 547 genome. Overall, we find that the vast majority of gaps are flanked by GT-AG consensus, and we
 548 conclude that they correspond to spliced introns. These introns are however not systematically
 549 spliced, because in many cases the proportion of spliced reads is between 0 (never spliced) and 1
 550 (always spliced).

551



552

553 **Figure 7.** TE spliced transcripts are ubiquitous. [Left depicts ovaries, right depicts testes.](#) Each circle depicts
 554 a TE insertion [supported by at least 5 reads](#) and their size is proportional to the expression level of the
 555 insertion. The TE family name is written for highly expressed insertions (>30 reads in testis, >10 reads in
 556 ovaries). The X-axis represents the proportion of reads that align with a gap (presence of N's in the CIGAR

557 string), while the Y-axis represents the proportion of the insertion covered by reads. Unfilled circles
558 correspond either to TE insertions with no gaps, or to TE insertions with gaps that do not exhibit GT-AG sites.

559 While the proportion of spliced transcripts stemming from a TE copy can vary, there are a couple
560 of copies that only produce spliced transcripts, as POGO\$2R_RaGOO\$7201268\$7202754 for instance.
561 *Pogo* is the most expressed TE family in ovaries, with 12 out of 26 copies producing capped poly-A
562 transcripts corresponding to 213 long reads, while only 7 expressed copies with a total of 32 long-
563 reads are observed in testes, despite the higher coverage. While we previously noted that only one
564 full-length copy is transcribed in ovaries (and in testes albeit with a lower number of reads), there are
565 many truncated or deleted copies that are transcribed (Figure 5B).
566 POGO\$2R_RaGOO\$7201268\$7202754 is one of the internally deleted copies, and it produces a spliced
567 transcript present in both testis and ovaries (Figure S21). The splicing of this short intron (55 nt) has
568 been previously reported (Tudor et al., 1992) and enables the splicing of the two ORFs of *pogo* into a
569 single continuous ORF. This particular copy (POGO\$2R_RaGOO\$7201268\$7202754) is however non-
570 functional since it contains a large genomic deletion located in the ORF near the intron.
571 POGO\$X_RaGOO\$21863530\$21864880 (Figure S20) also contains a large genomic deletion,
572 encompassing the intron, explaining why there are no spliced transcripts for this copy.

573 Despite the presence of full-length *Copia* insertions in the genome, only spliced transcripts were
574 uncovered in the long-read sequencing (Figure 6). In contrast, with Illumina short reads, we see both
575 spliced and unspliced transcripts (Figure 8). A similar pattern occurs with
576 MAX\$3L_RaGOO\$3640512\$3649209 (Figure S30), but not with POGO\$2L_RaGOO\$2955877\$2958004
577 or DM1731\$Y_RaGOO\$340770\$345273 (Figure S31-32). The full-length *Copia* transcripts are 5 Kb, and
578 are less abundant than the spliced transcripts (~10 times less). The lack of such a full-length transcript
579 in the long-read sequencing data might be explained by the lower expression level and the length of
580 the transcript. One can not discard the possibility that deeper long-read coverage might uncover full-
581 length, unspliced, *Copia* transcripts. It is important to stress that by using only short-reads it is nearly
582 impossible to determine which *Copia* sequence is being expressed as the vast majority of short reads
583 map to multiple locations with the exact same alignment score. With short-reads, at least one full-
584 length *Copia* insertion is expressed but its specific location remains unknown. Furthermore, if we
585 restrict the analysis to primary alignments (i.e. a randomly chosen alignment in the case of multiple
586 mapping), then the coverage of the intronic sequence decreases and it is no longer clear if the
587 insertion produces both spliced and unspliced transcripts (Figure S24). Overall, for *Copia*, the long-
588 reads enable the identification of which insertion is being transcribed, and the short-reads enable the
589 detection of the presence of the two splice variants. Some multi-mapping long reads could support
590 the presence of the unspliced transcript because they partially map to *Copia* intron, but we cannot
591 know from which insertion they were transcribed (Figure S25). Finally, spliced transcripts are unable
592 to produce the complete transposition machinery as they lack the reverse transcriptase enzyme and
593 are only able to produce the gag protein.



594

595

596

597

598

599

Figure 8. Example of *Copia* splicing. IGV screenshot shows spliced transcripts using long and short-read datasets. In green, testis coverage and an excerpt of mapped reads, in purple the same information for ovaries. In the excerpt of mapped reads, white rectangles correspond to multi-mapping reads. Dmgoth101 repeat and gene tracks are also shown and more information on the annotation can be seen in the material and methods section.

600

Conclusion

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

Long-read sequencing remains a major progress in the study of repeat transcription. Here we demonstrated the feasibility of assigning long reads to specific copies. In addition, quantification of TE expression with long-read sequencing is similar to short-read analysis, suggesting not only one could recover copy-specific information but also perform quantitative and differential expression analysis.

The genome of *D. melanogaster* contains many functional full-length copies but only a couple of such copies produce full-length transcripts in gonads. Given TEs are tightly controlled in the germline, one can wonder how many full-length copies might be expressed in somatic tissues. It is also important to stress that, to our knowledge, this is the first comparison of the expression of TEs between testes and ovaries, and we uncover a different TE transcriptional landscape regarding TE subclasses, using both short-reads and long-reads. Furthermore, in many instances, we see that TEs are spliced, independently of their structure or class. While some of these introns had been reported in the literature 30 years ago, the relevance and prevalence of these spliced transcripts have not always been investigated. Long-read sequencing could facilitate the exhaustive inventory of all spliceforms, in particular for recent TEs, where short reads are harder to use due to multiple mapping. A difficulty that remains when assessing if the intron of a particular TE insertion has really been spliced is the possibility that there exists a retroposed copy of a spliced version of this TE elsewhere in a non-assembled part of the genome. Here, taking advantage of the availability of raw genomic Nanopore

618 reads for the same dataset (ERR4351625), we could verify that this was not the case for Copia, the
619 youngest expressed element in our dataset. In practice, we mapped the genomic reads to both Copia
620 and a spliced version of Copia and found no genomic read mapping to the spliced version.

621 Finally, it is important to note that we did not recover TE transcripts longer than ~4.5 Kb. While
622 the detection of rare transcripts might indeed pose a problem to most sequencing chemistries, it
623 would be important to verify if long transcripts necessitate different RNA extraction methods for ONT
624 sequencing. For instance, the distribution of cDNA used here for ONT library construction reflects the
625 distribution of reads, with a low number of cDNAs longer than 3.5 Kb (Figure S1).

626 Acknowledgements

627 This work was performed using the computing facilities of the LBBE/PRABI. We thank Josefa
628 Gonzalez laboratory for discussing our preprint in their journal club and suggesting modifications to
629 the manuscript.

630 Data, scripts, code, and supplementary information availability

631 Data are available online at the BioProject PRJNA956863 (ONT long-reads), PRJNA981353
632 (SRX20759708, SRX20759707, testes short-reads), PRJNA795668 (SRX13669659 and SRX13669658,
633 ovaries short-reads). Scripts are available at https://gitlab.inria.fr/erable/te_long_read. Processed
634 data (.bam files) are available at <https://zenodo.org/records/10277511>.

635 Conflict of interest disclosure

636 The authors declare that they have no financial conflicts of interest in relation to the content of
637 the article.

638 Funding

639 This work was supported by French National research agency (ANR project ANR-16-CE23-0001
640 ‘ASTER and ANR-20-CE02-0015 Longevity), along with UDL-FAPESP2020.

641 References

642 Adrion JR, Song MJ, Schrider DR, Hahn MW, Schaack S (2017) Genome-Wide Estimates of
643 Transposable Element Insertion and Deletion Rates in *Drosophila Melanogaster*. *Genome Biology and
644 Evolution*, **9**, 1329–1340. <https://doi.org/10.1093/gbe/evx050>

645 Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) Basic local alignment search tool.
646 *Journal of Molecular Biology*, **215**, 403–410. [https://doi.org/10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2)

647 Babarinde IA, Ma G, Li Y, Deng B, Luo Z, Liu H, Abdul MM, Ward C, Chen M, Fu X, Shi L,
648 Duttlinger M, He J, Sun L, Li W, Zhuang Q, Tong G, Frampton J, Cazier J-B, Chen J, Jauch R, Esteban
649 MA, Hutchins AP (2021) Transposable element sequence fragments incorporated into coding and

650 noncoding transcripts modulate the transcriptome of human pluripotent stem cells. *Nucleic Acids*
651 *Research*, **49**, 9132–9153. <https://doi.org/10.1093/nar/gkab710>

652 Bailly-Bechet M, Haudry A, Lerat E (2014) “One code to find them all”: a perl tool to conveniently
653 parse RepeatMasker output files. *Mobile DNA*, **5**, 13. <https://doi.org/10.1186/1759-8753-5-13>

654 Belancio VP, Hedges DJ, Deininger P (2006) LINE-1 RNA splicing and influences on mammalian
655 gene expression. *Nucleic Acids Research*, **34**, 1512–1521. <https://doi.org/10.1093/nar/gkl027>

656 Bendall ML, Mulder M de, Iñiguez LP, Lecanda-Sánchez A, Pérez-Losada M, Ostrowski MA,
657 Jones RB, Mulder LCF, Reyes-Terán G, Crandall KA, Ormsby CE, Nixon DF (2019) Telescope:
658 Characterization of the retrotranscriptome by accurate estimation of transposable element expression.
659 *PLOS Computational Biology*, **15**, e1006453. <https://doi.org/10.1371/journal.pcbi.1006453>

660 Berrens RV, Yang A, Laumer CE, Lun ATL, Bieberich F, Law C-T, Lan G, Imaz M, Bowness JS,
661 Brockdorff N, Gaffney DJ, Marioni JC (2022) Locus-specific expression of transposable elements in
662 single cells with CELLO-seq. *Nature Biotechnology*, **40**, 546–554. [https://doi.org/10.1038/s41587-021-](https://doi.org/10.1038/s41587-021-01093-1)
663 [01093-1](https://doi.org/10.1038/s41587-021-01093-1)

664 Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, Imbeault M, Izsvák
665 Z, Levin HL, Macfarlan TS, Mager DL, Feschotte C (2018) Ten things you should know about
666 transposable elements. *Genome Biology*, **19**, 199. <https://doi.org/10.1186/s13059-018-1577-z>

667 Brierley C, Flavell AJ (1990) The retrotransposon copia controls the relative levels of its gene
668 products post-transcriptionally by differential expression from its two major mRNAs. *Nucleic Acids*
669 *Research*, **18**, 2947–2951. <https://doi.org/10.1093/nar/18.10.2947>

670 Dai Z, Ren J, Tong X, Hu H, Lu K, Dai F, Han M-J (2021) The Landscapes of Full-Length
671 Transcripts and Splice Isoforms as Well as Transposons Exonization in the Lepidopteran Model
672 System, *Bombyx mori*. *Frontiers in Genetics*, **12**, 704162. <https://doi.org/10.3389/fgene.2021.704162>

673 De Coster W, D’Hert S, Schultz DT, Cruts M, Van Broeckhoven C (2018) NanoPack: visualizing
674 and processing long-read sequencing data. *Bioinformatics*, **34**, 2666–2669.
675 <https://doi.org/10.1093/bioinformatics/bty149>

676 Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Batut P, Chaisson M, Gingeras
677 TR (2013) STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, **29**, 15–21.
678 <https://doi.org/10.1093/bioinformatics/bts635>

679 Ewing AD, Smits N, Sanchez-Luque FJ, Faivre J, Brennan PM, Richardson SR, Cheetham SW,
680 Faulkner GJ (2020) Nanopore Sequencing Enables Comprehensive Transposable Element Epigenomic
681 Profiling. *Molecular Cell*, **80**, 915-928.e5. <https://doi.org/10.1016/j.molcel.2020.10.024>

682 Fablet M, Salces-Ortiz J, Jacquet A, Menezes BF, Dechaud C, Veber P, Rebollo R, Vieira C (2023)
683 A Quantitative, Genome-Wide Analysis in *Drosophila* Reveals Transposable Elements’ Influence on
684 Gene Expression Is Species-Specific. *Genome Biology and Evolution*, **15**, evad160.
685 <https://doi.org/10.1093/gbe/evad160>

686 Fabry MH, Falconio FA, Joud F, Lythgoe EK, Czech B, Hannon GJ (2021) Maternally inherited

687 piRNAs direct transient heterochromatin formation at active transposons during early *Drosophila*
688 embryogenesis. *eLife*, **10**, e68573. <https://doi.org/10.7554/eLife.68573>

689 Jiang F, Zhang J, Liu Q, Liu X, Wang H, He J, Kang L (2019) Long-read direct RNA sequencing
690 by 5'-Cap capturing reveals the impact of Piwi on the widespread exonization of transposable elements
691 in locusts. *RNA biology*, **16**, 950–959. <https://doi.org/10.1080/15476286.2019.1602437>

692 Jin Y, Hammell M (2018) Analysis of RNA-Seq Data Using TETranscripts. In: *Transcriptome Data*
693 *Analysis: Methods and Protocols* Methods in Molecular Biology. (eds Wang Y, Sun M), pp. 153–167.
694 Springer, New York, NY. https://doi.org/10.1007/978-1-4939-7710-9_11

695 Jin Y, Tam OH, Paniagua E, Hammell M (2015) TETranscripts: a package for including
696 transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics (Oxford,*
697 *England)*, **31**, 3593–3599. <https://doi.org/10.1093/bioinformatics/btv422>

698 Lanciano S, Cristofari G (2020) Measuring and interpreting transposable element expression.
699 *Nature Reviews. Genetics*, **21**, 721–736. <https://doi.org/10.1038/s41576-020-0251-y>

700 Laski FA, Rio DC, Rubin GM (1986) Tissue specificity of *Drosophila* P element transposition is
701 regulated at the level of mRNA splicing. *Cell*, **44**, 7–19. [https://doi.org/10.1016/0092-8674\(86\)90480-](https://doi.org/10.1016/0092-8674(86)90480-0)
702 0

703 Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C (2017) TETools facilitates big data
704 expression analysis of transposable elements and reveals an antagonism between their activity and that
705 of piRNA genes. *Nucleic Acids Research*, **45**, e17. <https://doi.org/10.1093/nar/gkw953>

706 Li H (2018) Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, **34**, 3094–
707 3100. <https://doi.org/10.1093/bioinformatics/bty191>

708 Li H, Janssens J, De Waegeneer M, Kolluru SS, Davie K, Gardeux V, Saelens W, David FPA,
709 Brbić M, Spanier K, Leskovec J, McLaughlin CN, Xie Q, Jones RC, Brueckner K, Shim J, Tattikota
710 SG, Schnorrer F, Rust K, Nystul TG, Carvalho-Santos Z, Ribeiro C, Pal S, Mahadevaraju S, Przytycka
711 TM, Allen AM, Goodwin SF, Berry CW, Fuller MT, White-Cooper H, Matunis EL, DiNardo S,
712 Galenza A, O'Brien LE, Dow JAT, FCA Consortium, Jasper H, Oliver B, Perrimon N, Deplancke B,
713 Quake SR, Luo L, Aerts S (2022) Fly Cell Atlas: A single-nucleus transcriptomic atlas of the adult fruit
714 fly. *Science*, **375**, eabk2432. <https://doi.org/10.1126/science.abk2432>

715 Miller K, Rosenbaum J, Zbrzezna V, Pogo AO (1989) The nucleotide sequence of *Drosophila*
716 *melanogaster* copia-specific 2.1-kb mRNA. *Nucleic Acids Research*, **17**, 2134.
717 <https://doi.org/10.1093/nar/17.5.2134>

718 Mohamed M, Dang NT-M, Ogyama Y, Burllet N, Mugat B, Boulesteix M, Mérel V, Veber P,
719 Salces-Ortiz J, Severac D, Péliisson A, Vieira C, Sabot F, Fablet M, Chambeyron S (2020) A Transposon
720 Story: From TE Content to TE Dynamic Invasion of *Drosophila* Genomes Using the Single-Molecule
721 Sequencing Technology from Oxford Nanopore. *Cells*, **9**, 1776. <https://doi.org/10.3390/cells9081776>

722 Panda K, Slotkin RK (2020) Long-Read cDNA Sequencing Enables a “Gene-Like” Transcript
723 Annotation of Transposable Elements. *The Plant Cell*, **32**, 2687–2698.

724 <https://doi.org/10.1105/tpc.20.00115>

725 Pélisson A, Song SU, Prud'homme N, Smith PA, Bucheton A, Corces VG (1994) Gypsy
726 transposition correlates with the production of a retroviral envelope-like protein under the tissue-
727 specific control of the *Drosophila flamenco* gene. *The EMBO journal*, **13**, 4401–4411.
728 <https://doi.org/10.1002/j.1460-2075.1994.tb06760.x>

729 Rech GE, Radío S, Guirao-Rico S, Aguilera L, Horvath V, Green L, Lindstadt H, Jamilloux V,
730 Quesneville H, González J (2022) Population-scale long-read sequencing uncovers transposable
731 elements associated with gene expression variation and adaptive signatures in *Drosophila*. *Nature*
732 *Communications*, **13**, 1948. <https://doi.org/10.1038/s41467-022-29518-8>

733 Sessegolo C, Cruaud C, Da Silva C, Cologne A, Dubarry M, Derrien T, Lacroix V, Aury J-M
734 (2019) Transcriptome profiling of mouse samples using nanopore sequencing of cDNA and RNA
735 molecules. *Scientific Reports*, **9**, 14908. <https://doi.org/10.1038/s41598-019-51470-9>

736 Sherman BT, Hao M, Qiu J, Jiao X, Baseler MW, Lane HC, Imamichi T, Chang W (2022) DAVID:
737 a web server for functional enrichment analysis and functional annotation of gene lists (2021 update).
738 *Nucleic Acids Research*, **50**, W216–W221. <https://doi.org/10.1093/nar/gkac194>

739 Shumate A, Salzberg SL (2021) Liftoff: accurate mapping of gene annotations. *Bioinformatics*
740 (*Oxford, England*), **37**, 1639–1643. <https://doi.org/10.1093/bioinformatics/btaa1016>

741 Slotkin RK, Martienssen R (2007) Transposable elements and the epigenetic regulation of the
742 genome. *Nature Reviews Genetics*, **8**, 272–285. <https://doi.org/10.1038/nrg2072>

743 Teixeira FK, Okuniewska M, Malone CD, Coux R-X, Rio DC, Lehmann R (2017) piRNA-
744 mediated regulation of transposon alternative splicing in the soma and germ line. *Nature*, **552**, 268–
745 272. <https://doi.org/10.1038/nature25018>

746 Tudor M, Lobočka M, Goodell M, Pettitt J, O'Hare K (1992) The pogo transposable element
747 family of *Drosophila melanogaster*. *Molecular & general genetics: MGG*, **232**, 126–134.
748 <https://doi.org/10.1007/BF00299145>

749 White R, Pellefigues C, Ronchese F, Lamiable O, Eccles D (2017) Investigation of chimeric reads
750 using the MinION. *F1000Research*, **6**, 631. <https://doi.org/10.12688/f1000research.11547.2>

751 Yang WR, Ardeljan D, Pacyna CN, Payer LM, Burns KH (2019) SQuIRE reveals locus-specific
752 regulation of interspersed repeat expression. *Nucleic Acids Research*, **47**, e27.
753 <https://doi.org/10.1093/nar/gky1301>

754 Yoshioka K, Honma H, Zushi M, Kondo S, Togashi S, Miyake T, Shiba T (1990) Virus-like
755 particle formation of *Drosophila copia* through autocatalytic processing. *The EMBO journal*, **9**, 535–
756 541. <https://doi.org/10.1002/j.1460-2075.1990.tb08140.x>

757