

24

25 **Abstract**

26 Effective suppression of transposable elements (TEs) is paramount to maintain genomic
27 integrity and organismal fitness. In *D. melanogaster*, *flamenco* is a master suppressor of TEs,
28 preventing their movement from somatic ovarian support cells to the germline. It is transcribed
29 by Pol II as a long (100s of kb), single-stranded, primary transcript, that is metabolized into
30 Piwi-interacting RNAs (piRNAs) that target active TEs via antisense complementarity. *flamenco*
31 is thought to operate as a trap, owing to its high content of recent horizontally transferred TEs
32 that are enriched in antisense orientation. Using newly-generated long read genome data, which
33 is critical for accurate assembly of repetitive sequences, we find that *flamenco* has undergone
34 radical transformations in sequence content and even copy number across *simulans* clade
35 Drosophilid species. *D. simulans flamenco* has duplicated and diverged, and neither copy
36 exhibits synteny with *D. melanogaster* beyond the core promoter. Moreover, *flamenco*
37 organization is highly variable across *D. simulans* individuals. Next, we find that *D. simulans*
38 and *D. mauritiana flamenco* display signatures of a dual-stranded cluster, with ping-pong signals
39 in the testis and embryo. This is accompanied by increased multicopy elements, consistent with
40 these regions operating as functional dual stranded clusters. Overall, the physical and functional
41 diversity of *flamenco* orthologs is testament to the extremely dynamic consequences of TE arms
42 races on genome organization, not only amongst highly related species, but even amongst
43 individuals.

44

45

46

47 **Introduction**

48 *Drosophila* gonads exemplify two important fronts in the conflict between transposable elements
49 (TEs) and the host – the germline (which directly generates gametes), and somatic support cells
50 (from which TEs can invade the germline) (1, 2). The strategies by which TEs are suppressed in
51 these settings are distinct (3), but share their utilization of piwi-interacting RNAs (piRNAs).
52 These are ~24-32 nt RNAs that are bound by the PIWI subclass of Argonaute proteins, and guide
53 them and associated cofactors to targets for transcriptional and/or post-transcriptional silencing
54 (4–7).

55 Mature piRNAs are processed from non-coding piRNA cluster transcripts, which derive
56 from genomic regions that are densely populated with TE sequences (7–9). However, the
57 mechanisms of piRNA biogenesis differ between gonadal cell types. In the germline, piRNA
58 clusters are transcribed from both DNA strands through non-canonical Pol II activity (6, 10–12),
59 which is initiated by chromatin marks rather than specific core promoter motifs. Moreover, co-
60 transcriptional processes such as splicing and polyadenylation are suppressed within dual strand
61 piRNA clusters (13, 14). On the other hand, in ovarian somatic support cells, piRNA clusters are
62 transcribed from a typical promoter as a single stranded transcript, which can be alternatively
63 spliced as with protein-coding mRNAs (15–18). These rules derive in large part from the study
64 of model piRNA clusters (i.e. the germline *42AB* and somatic *flamenco* piRNA clusters). For
65 both types, their capacity to repress invading transposable elements is thought to result from
66 random integration of new transposons into the cluster. As such, piRNA clusters are adaptive
67 loci that play central roles in the conflict between hosts and TEs.

68 The location and activity of germline piRNA clusters are stochastic and evolutionarily
69 dynamic, as there are many copies of TE families in different locations that may produce

70 piRNAs (9, 19). By contrast, somatic piRNA clusters are not redundant and a single insertion of
71 a TE into a somatic piRNA cluster should be sufficient to prevent that TE from further
72 transposition (18, 20). Thus, *flamenco* should contain only one copy per TE family (18), which is
73 true in the *flamenco* locus of *D. melanogaster* (18). *flamenco* is also the only piRNA cluster
74 which produces a phenotypic effect when altered, as germline clusters can be deleted with no
75 consequences.

76 *flamenco* has been a favored model for understanding the piRNA pathway since the
77 discovery of piRNA mediated silencing of transposable elements (6). *flamenco* spans >180 kb of
78 repetitive sequences located in β -heterochromatin of the X chromosome (21). Of note, *flamenco*
79 was initially identified, prior to the formal recognition of piRNAs, via transposon insertions that
80 de-repress *gypsy*, *ZAM*, and *Idefix* class elements (21–25). These mutant alleles disrupt the
81 *flamenco* promoter, and consequently abrogate transcription and piRNA production from this
82 locus. By contrast, the recent deletion of multiple model germline piRNA clusters, which
83 eliminate the biogenesis of a bulk of cognate piRNAs, did not de-repress their cognate TEs (9).
84 Thus, the analysis of *flamenco* evolution is presumably more consequential for TE dynamics.
85 Analysis of *flamenco* in various strains of *D. melanogaster* supports that this locus traps
86 horizontally derived TEs to achieve silencing of newly invaded TEs (18). The *flamenco* locus
87 exhibits synteny across the *D. melanogaster* sub-group (26); however, the sequence composition
88 of *flamenco* outside *D. melanogaster* has not been well-characterized (27).

89 In this study, we compare the *flamenco* locus across 10 strains of simulans-clade species,
90 namely *D. simulans*, *D. mauritiana*, and *D. sechellia*. Analysis of piRNAs from ovaries of five
91 genotypes of *D. simulans* found that *flamenco* is duplicated in *D. simulans*. This duplication is
92 old enough that there is no sequence synteny across copies, even though their core promoter

93 regions and the adjacent *dip1* gene are conserved. *flamenco* has also been colonized by abundant
94 (>40) copies of *RI*, a TE that was thought to insert only at ribosomal genes, and to evolve at the
95 same rate as nuclear genes [21]. Furthermore, between different genotypes, up to 63% of TE
96 insertions are not shared within any given copy of *flamenco*. Despite this, several full length TEs
97 are shared between all genotypes in a similar sequence context. This incredible diversity at the
98 *flamenco* locus, even within a single species, suggests there may be considerable variation in its
99 ability to suppress transposable elements across individuals.

100 Cross-species comparisons further indicate that functions of *flamenco* have diversified.
101 Data from *D. sechellia* and *D. melanogaster* conform with the current understanding of *flamenco*
102 as a uni-strand cluster. However, we find evidence that *D. simulans* and *D. mauritiana* *flamenco*
103 can act as a dual strand cluster in testis (*D. mauritiana*) and embryos (*D. mauritiana* and *D.*
104 *simulans*), yielding piRNAs from both strands with a ping pong signal. Overall, we infer that the
105 rapid evolution of *flamenco* alleles across individuals and species reflects highly adaptive
106 functions and dynamic biogenesis capacities.

107 **Materials and Methods**

108 *Fly strains*

109 The four *D. simulans* lines *SZ232*, *SZ45*, *SZ244*, and *SZ129* were collected in California from the
110 Zuma Organic Orchard in Los Angeles, CA on two consecutive weekends of February 2012 [57–
111 61]. *LNP-15-062* was collected in Zambia at the Luwangwa National Park by D. Matute and
112 provided to us by J. Saltz (J. Saltz pers. comm., [41,53]). *MD251*, *MD242*, *NS137*, and *NS40*
113 were collected in Madagascar and Kenya (respectively) and are described in [50]. The *D.*
114 *simulans* strain *wxD¹* was originally collected by M. Green, likely in California, but its

115 provenance has been lost (pers. comm. Jerry Coyne). *D. mauritiana* (*w12*) and *D. sechellia*
116 (*Rob3c/Tucson 14021-0248.25*) are described in [11].

117 *Long read DNA sequencing and assembly*

118 *MD242*, four SZ lines and *LNP-15-062* were sequenced on a MinION platform at North Dakota
119 State University (Oxford Nanopore Technologies (ONT), Oxford, GB), with base-calling using
120 guppy (v4.4.2). *MD242*, the four SZ lines, and *LNP-15-062* were assembled with Canu (v2.1)
121 [73] and two rounds of polishing with Racon (v1.4.3) [67]. The CA strains were additionally
122 polished with short reads using Pilon (v1.23) [68](SRR3585779, SRR3585440, SRR3585480,
123 SRR3585391) [60]. The first *wxD^{I-1}* assembly is described here [12]. *MD251*, *NS137*, *NS40* and
124 *wxD^{I-2}* were sequenced on a MinION platform by B. Kim at Stanford University. They were
125 assembled with Flye [29], and polished with a round of Medaka followed by a round of pilon
126 [68]. Following this contaminants were removed with blobtools
127 (<https://zenodo.org/record/845347>, [30]), soft masked with RepeatModeler and Repeatmasker
128 [22,64], then aligned to the *wxD^I* as a reference with Progressive Cactus [3]. The assemblies
129 were finished with reference based scaffolding using Ragout [28]. *D. mauritiana* and *D.*
130 *sechellia* were sequenced with PacBio and assembled with FALCON using default parameters
131 (<https://github.com/PacificBiosciences/FALCON>)[11].The *D. melanogaster* assembly is
132 described here (47). A summary of the assembly statistics is available in Supplementary Table 1.
133 The quality of cluster assembly was evaluated using CUSCO as described in (19, 48)
134 (Supplementary File 1).

135 *Short read sequencing and mapping*

136 Short read sequencing was performed by Beijing Genomics Institute (BGI) on approximately 50 dissected
137 ovaries from adult female flies (*SZ45*, *SZ129*, *SZ232*, *SZ244*, *LNP-15-062*). Short read libraries from 0-2
138 hour embryos were prepared from *D. melanogaster*, *wxD^{I-2}*, *D. sechellia*, and *D. mauritiana* (SRAXXX)

139 (49). Small RNA from testis is described in (50, 51). *D. melanogaster* OSC small RNA libraries were
140 downloaded from the SRA (SRR11999160). Libraries were filtered for adapter contamination and short
141 reads between 23-29 bp were retained for mapping with fastp (52). The RNA was then mapped to
142 their respective genomes using bowtie (v1.2.3) and the following parameters (-q -v 1 -p 1 -S -a -
143 m 50 --best --strata) (53, 54). The resulting bam files were processed using samtools (55). To
144 obtain unique reads the bam files were filtered for reads with 1 mapping position. To obtain
145 counts files with weighted mapping the bam files were processed using Rsubreads and the
146 featureCounts function (56).

147 *Defining and annotating piRNA clusters*

148 piRNA clusters were defined using proTRAC [52]. piRNA clusters were predicted with a
149 minimum cluster size of 1 kb (option “-clsize 1000”), a P value for minimum read density of
150 0.07 (option “-pdens 0.07”), a minimum fraction of normalized reads that have 1T (1U) or 10A
151 of 0.33 (option “-1Tor10A 0.33”) and rejecting loci if the top 1% of reads account for more than
152 90% of the normalized piRNA cluster read counts (option “-distr 1-90”), and a minimal fraction
153 of hits on the main strand of 0.25 (option “-clstrand 0.25”). Note that this ties the piRNA clusters
154 to their function such that participation in the ping pong pathway can be inferred from these
155 patterns. Clusters were annotated using RepeatMasker (v. 4.0.7) and the TE libraries described in
156 Chakraborty et al. (2019) [12,64]. The position of *flamenco* was also evaluated based off of the
157 position of the putative promoter, the *dip1* gene, and the enrichment of *gypsy* elements [24].
158 Fragmented annotations were merged to form TE copies with onecodetofindthemall [5].
159 Fragmented annotations were also manually curated, particularly because TEs not present in the
160 reference library often have their LTRs and internal sequences classified as different elements.

161 *Aligning the flamenco promoter region*

162 The region around the *flamenco* promotor was extracted from each genotype and species with
163 bedtools getfasta (61). Sequences were aligned with clustal-omega and converted to nexus
164 format (62). Trees were built using a GTR substitution model and gamma distributed rate
165 variation across sites (63). The markov chain monte carlo chains were run until the standard
166 deviation of split frequencies was below .01, around one million generations. The consensus
167 trees were generated using sumt conformformat=simple. The resulting trees were displayed with the
168 R package ape (64).

169 *Detecting ping pong signals in the small RNA data*

170 Ping pong signals were detected using pingpongpro [66]. This program detects the presence of
171 RNA molecules that are offset by 10 nt, such that stacks of piRNA overlap by the first 10 nt from
172 the 5' end. These stacks are a hallmark of piRNA mediated transposon silencing. The algorithm
173 also takes into account local coverage and the presence of an adenine at the 10th position. The
174 output includes a z-score between 0 and 1, the higher the z-score the more differentiated the ping
175 pong stacks are from random local stacks.

176 **Results**

177 *flamenco in the D. simulans clade*

178 We identified *D. simulans flamenco* from several lines of evidence: piRNA cluster calls from
179 proTRAC, its location adjacent to divergently transcribed *dip1*, the existence of conserved core
180 *flamenco* promoter sequences, and enrichment of *gypsy* elements (Figure 1A-D); Supplementary
181 Table 2). The *flamenco* locus is at least 376 kb in *D. simulans*. This is an expansion compared
182 with *D. melanogaster*, where *flamenco* is only 156 kb (*Canton-S*). In *D. sechellia flamenco* is
183 363 kb, however in *D. mauritiana* the locus has expanded to at least 840 kb (Supplementary
184 Table 2). This is a large expansion, and it is possible that the entire region does not act as the

185 *flamenco* locus. However, evidence that it does include uniquely mapping piRNAs are found
186 throughout the region and *gypsy* enrichment is consistent with a *flamenco*-like locus
187 (Supplementary Figure 1). There are no protein coding genes within the region, and while the
188 neighboring genes on the downstream side of *flamenco* in *D. melanogaster* have moved in *D.*
189 *mauritiana* (CG40813- CG41562 at 21.5 MB), the following group of genes beginning with
190 CG14621 is present and flanks *flamenco* as it is annotated. Thus in *D. melanogaster* the borders
191 of *flamenco* are flanked by *dip1* upstream and CG40813 downstream, while in *D. mauritiana*
192 they are *dip1* upstream and CG14621 downstream. Between all species the *flamenco* promoter
193 and surrounding region, including the *dip1* gene, are alignable and conserved (Figure 1E).

194 *Structure of the flamenco locus*

195 *Structure of the flamenco locus*

196 *D. melanogaster flamenco* bears a characteristic structure, in which the majority of TEs
197 are *gypsy*-class elements in the antisense orientation (79% antisense orientation, 85% of which
198 are *gypsy* elements) (Figure 1D; Supplementary Table 3). This is true in both the *iso-1* and
199 *Canton-S* strains. In *D. simulans*, *flamenco* has been colonized by large expansions of *R1*
200 transposable element repeats such that on average the percent of antisense TEs is only 50% and
201 the percent of the locus comprised of LTR elements is 55%. However, 76% of antisense
202 insertions are LTR insertions, thus the underlying *flamenco* structure is apparent when the *R1*
203 insertions are disregarded (Figure 1D). In *D. mauritiana flamenco* is 71% antisense, and of those
204 antisense elements it is 85% LTRs. Likewise in *D. sechellia* 78% of elements are antisense, and
205 of those 81% are LTRs. *flamenco* retains the overall structure of a canonical *D. melanogaster*-
206 like *flamenco* locus in all of these species, however in *D. simulans* the nature of the locus is
207 somewhat altered by the abundant *R1* insertions (Figure 1D).

208 *flamenco* is duplicated in *D. simulans*

209 In *D. simulans*, we unexpectedly observed that *flamenco* is duplicated on the X
210 chromosome; the duplication was confirmed with PCR and a restriction digest (Supplementary
211 Table 4). These duplications are associated with a conserved copy of the putative *flamenco*
212 enhancer as well as copies of the *dip1* gene located proximal to *flamenco* in *D. melanogaster*
213 (Figure 1C, 2A). While it is unclear which copy is orthologous to *D. melanogaster flamenco*, all
214 *D. simulans* lines bear one copy that aligns across genotypes. We refer to this copy as *D.*
215 *simulans flamenco*, and the other copies as duplicates. Otherwise, *flamenco* duplicates do not
216 align with one another and lack synteny amongst their resident TEs. Possible evolutionary
217 scenarios are that the *flamenco* duplication occurred early in the *simulans* lineage, that the
218 clustered evolved very rapidly, or that the duplication encompassed only the promoter region and
219 was subsequently colonized by TEs (Figure 1C, 2A).

220 The *flamenco* duplicate is absent in the *D. simulans* reference strain, *w*⁵⁰¹, but present in
221 *wxD*¹, suggesting it was polymorphic or absent between the collection of these strains (or was
222 not assembled). The duplicate retains the structure of *flamenco*, with an average of 67% of TEs
223 in the antisense orientation in the duplication of *flamenco*, and 91% of the TEs in the antisense
224 orientation are LTRs. The duplicate of *flamenco* is less impacted by *RI*, with some genotypes
225 having as few as 8 *RI* insertions (Figure 2C).

226 *RI LINE elements at the flamenco locus*

227 *RI* elements are well-known to insert into rDNA genes, are transmitted vertically, and evolve
228 similarly as the genome background rate [21]. They have also been found outside of rDNA
229 genes, but only as fragments. However, as mentioned, *RI* elements are abundant within *flamenco*
230 loci in the *simulans* clade. Outside of *flamenco*, *RI* elements in *D. simulans* are distributed

231 according to expectation, with full length elements occurring only within rDNA (Supplementary
232 File 6). Within *flamenco*, most copies of *RI* occur as tandem duplicates, creating large islands of
233 fragmented *RI* copies (Figure 2A). They are on average 3.7% diverged from the reference *RI*
234 from *D. simulans*. Across individual *D. simulans* genomes, ~99 kb of *flamenco* loci consists of
235 *RI* elements, fully 26% of their average total length. *SZ45*, *LNP-15-062*, *NS40*, *MD251*, and
236 *MD242* contain 4-7 full length copies of *RI* in the sense orientation, even though all but *SZ45*
237 bear fragmented *RI* copies on the antisense strand. (The *SZ45 flamenco* assembly is incomplete).
238 As the antisense *RI* copies are expected to suppress *RI* transposition, *flamenco* may not suppress
239 these elements effectively.

240 In *D. mauritiana*, *flamenco* harbors abundant fragments or copies of *RI* (19 on the
241 reverse strand and 20 on the forward strand), and only one large island of *RI* elements. In total,
242 *D. mauritiana* contains 84 kb of *RI* sequence within *flamenco*. In *D. mauritiana* there are 8 full
243 length copies of *RI* at the *flamenco* locus, 7 in antisense, which are not obviously due to a
244 segmental or local duplication. Finally, we find that *D. sechellia flamenco* lacks full length
245 copies of *RI*, and it contains only 18 KB of *RI* sequence (16 fragments on the reverse strand).
246 Yet, all the copies are on the sense strand, which would not produce fragments that can suppress
247 *RI* TEs. Essentially the antisense copies of *RI* in *D. mauritiana* should be suppressing the TE,
248 but we see multiple full length antisense insertions, and *D. sechellia* has no antisense copies, but
249 we see no evidence for recent *RI* insertions. From this it would appear that whatever is
250 controlling the transposition of *RI* lies outside of *flamenco*.

251 The presence of long sense-strand *RI* elements within *flamenco* is a departure from
252 expectation [21,72]. There is no evidence of an rDNA gene within the *flamenco* locus that would
253 explain the insertion of *RI* elements there, nor is there precedence for the large expansion of *RI*

254 fragments within the locus. Furthermore, the suppression of *RI* transposition does not appear to
255 be controlled by *flamenco*.

256 *piRNA production from RI*

257 On average *RI* elements within the *flamenco* locus of *D. simulans* produce more piRNA
258 than any other TE within *flamenco* (Supplementary Table 6). *RI* reads mapping to the forward
259 strand constitute an average of 51% of the total piRNAs within the *flamenco* locus from the
260 maternal fraction, ovary, and testis using weighted mapping. The only exception is the ovarian
261 sample from *SZ232* which is a large outlier at only 5%. However reads mapping to the reverse
262 strand account for an average of 84% of the piRNA being produced from the strand in every
263 genotype and tissue – maternal fraction, testis, or ovary. If unique mapping is considered instead
264 of weighted these percentages are reduced by approximately 20%, which is to be expected given
265 that *RI* is present in many repeated copies. Production of piRNA from the reverse strand seems
266 to be correlated with elements inserted in the sense orientation, of which the vast majority are *RI*
267 elements in *D. simulans* (Supplementary Figure 2). The production of large quantities of piRNA
268 cognate to the *RI* element is seemingly pointless – if *RI* only inserts at rDNA genes and are
269 vertically transmitted there is little reason to be producing the majority of piRNA in response to
270 this element.

271 In *D. sechellia* there are very few piRNA produced from *flamenco* in these tissues, and
272 there are no full length copies of *RI*. Likewise overall weighted piRNA production from *RI*
273 elements on either strand is 2.8-5.9% of the total mapping piRNA. In contrast in *D. mauritiana*
274 there are full length *RI* elements and abundant piRNA production in the maternal fraction and
275 testis. In *D. mauritiana* an average of 28% of piRNAs mapping to the forward strand of *flamenco*
276 are arising from *RI*, and 33% from the reverse strand. In *D. mauritiana* *RI* elements make up a

277 smaller proportion of the total elements in the sense orientation (24%), versus *D. simulans*
278 (55%).

279 *Conservation of flamenco*

280 The *dip1* gene and promoter region adjacent to each copy of *flamenco* are very conserved both
281 within and between copies of *flamenco* (Figure 2). The phylogenetic tree of the area suggests that
282 we are correct in labeling the two copies as the original *flamenco* locus and the duplicate (Figure
283 2). The original *flamenco* locus is more diverged amongst copies while the duplicate clusters
284 closely together with short branch lengths (Figure 2). They are also conserved and alignable
285 between *D. melanogaster*, *D. sechellia*, *D. mauritiana*, and *D. simulans* (Figure 1). However, the
286 same is not true of the *flamenco* locus itself. Approximately 3 kb from the promoter *flamenco*
287 diverges amongst genotypes and species and is no longer alignable by traditional sequence-based
288 algorithms, as the TEs are essentially a presence/absence that spans multiple kb. There is no
289 conservation of *flamenco* between *D. melanogaster*, *D. simulans*, *D. sechellia*, and *D.*
290 *mauritiana* (Figure 3). However, within the *simulans* clade many of the same TEs occupy the
291 locus, suggesting that they are the current genomic invaders in each of these species (Figure 3).

292 In *D. simulans* the majority of full length TEs are singletons – 52% in *flamenco* and 64%
293 in the duplicate. Copies that are full length in one genotype but fragmented in others are counted
294 as shared, not singletons. Almost half of these singletons in the duplicate are due to a single
295 genotype with a unique section of sequence, in this case *MD251*. Likewise a third of the
296 singleton insertions in the duplicate are due to an *NS40* specific region of *flamenco*. Regardless
297 of these concentrations of singletons in single genotypes, it is the single largest category of
298 transposable element insertions, followed by fixed insertions. Thus even within a single
299 population there is considerable diversity at the *flamenco* locus, and subsequently diversity in the

300 ability to suppress transposable elements. For example, *gypsy-29* is present in three genotypes
301 either in *flamenco* or the duplicate, which would suggest that these genotypes are able to
302 suppress this transposable element in the somatic support cells of the ovary while the other
303 genotypes are not. In contrast *gypsy-3* is present in more than one full length copy in *flamenco*
304 and its duplicate in every genotype but one where it is present in a single copy. There are a
305 number of these conserved full length TEs that are present in all or nearly all genotypes,
306 including *Chimpo*, *gypsy-2*, *Tirant*, and *gypsy-4*. In addition, the *INE1* elements adjacent to the
307 promoter are always conserved.

308 It is notable that any full length TEs are shared across all genotypes, given that *wxD¹* was
309 like collected 30-50 years prior to the others, and the collections span continents (Figure 2). Two
310 facts are relevant to this observation: (1) TEs were shown not correlate with geography [32] and
311 (2) *D. simulans* is more diverse within populations than between different populations
312 [38,54,62]. Other explanations are also plausible. Selection could be maintaining these full
313 length TEs, *wxD¹* could have had introgression from other lab strains, or a combination of these
314 explanations.

315 *Suppression of TEs by the flamenco locus and the trap model of TE control*

316 In *D. melanogaster*, it was proposed that while germline clusters may have many insertions of a
317 single TE, the somatic 'master regulator' *flamenco* will have a single insertion of each
318 transposon, after which they are silenced and no longer able to transpose [72].

319 Here, we evaluate the following lines of evidence to determine if they support the trap model of
320 transposable element suppression. (1) How many TEs have antisense oriented multicopy
321 elements within *flamenco*? (2) How many TEs have full length and fragmented insertions,
322 suggesting the older fragments did not suppress the newer insertion? (3) How many *de novo*

323 insertions of TEs in the *flamenco* duplicate of *D. simulans* are also present in the original
324 *flamenco* copy?

325 How many TEs have antisense oriented multicopy elements within *flamenco*?

326 Due to the difficulty in classifying degraded elements accurately, for example between multiple
327 classes of *gypsy* element, we will focus here on full length TEs, suggesting recent transposition.

328 In *D. melanogaster* there are 7 full length elements, none of which are present in more than one
329 antisense copy. These elements make up 27% of the *flamenco* locus. Full length copies of five of
330 these elements were also reported previously for other strains of *D. melanogaster* (18)

331 In *D. sechellia* there are 14 full length TEs within the *flamenco* locus, three of which are
332 present in multiple copies. Two of these, *INE1* and *412*, are likely present due to local
333 duplication. In particular the *INE1* elements flank the promoter, are in the sense orientation, and
334 are conserved between *D. sechellia*, *D. mauritiana*, and *D. simulans*. The only element present in
335 multiple antisense copies is *GTWIN*. Similar to *D. melanogaster* these elements make up 27% of
336 the *flamenco* locus.

337 *D. mauritiana* contains 22 full length TEs within the *flamenco* locus. Four of these are
338 present in multiple antisense full length copies – *INE1*, *RI*, *Stalker-4*, and *Cr1a*. While some of
339 the five antisense copies of *RI* likely originated from local duplications – they are in the same
340 general region and tend to be flanked by *gypsy-8*, not all of them show these patterns.
341 Furthermore, as aforementioned, there also are full length sense copies of *RI* suggesting *RI* is
342 not being suppressed by *flamenco*. *gypsy-12* and *gypsy-3* have a second antisense copy within
343 *flamenco* that is just below the cutoff to be considered full length – in *gypsy-3* the second copy is
344 10% smaller, for *gypsy-12* it is 80% present but missing an LTR. Full length TEs make up 19%
345 of the *flamenco* locus.

346 In *D. simulans* there are 29 full length TEs present in any of the seven complete *flamenco*
347 assemblies. Eight of these are present in multiple antisense copies within a single genome –
348 *INE1*, *Chimpo*, *copia*, *gypsy-3*, *gypsy-4*, *412*, *Tirant*, and *BEL-unknown*. The two *Tirant* copies
349 are likely a segmental duplication as they flank an *RI* repeat region. In addition, most *INE1*
350 copies are present proximal to the promoter as aforementioned, however in *NS40* a copy is
351 present in antisense at the end of the locus. *Chimpo* is present in three full length copies within
352 *MD242 flamenco*, with no evidence of local duplication. While there are no full length copies of
353 *RI* inserted in antisense, *RI* is present in full length sense copies despite many genomes
354 containing antisense fragments, suggesting *flamenco* is not suppressing *RI*. On average full
355 length TEs constitute 20% of *flamenco* in *D. simulans*.

356 In the duplicate of *flamenco* in *D. simulans* there are 30 full length TEs present in any
357 one of the five complete *flamenco* duplicate assemblies. However, none of them are multicopy in
358 antisense. However, they are multicopy relative to the original copy of *flamenco*. *gypsy-3*, *BEL-*
359 *unknown*, *Nomad-1*, *Chimpo*, *gypsy-53A*, *RI*, and *INE1* are all multicopy with respect to the
360 original *flamenco* within a given genome. Some of these may have been inherited at the time of
361 duplication, however are full length in both copies suggesting recent transposition. In the
362 duplicate of *flamenco* full length TEs occupy an average of 17% of the locus. *MD251* is an
363 exception which weights the average, with 28% of the locus, while between 10 and 15% is found
364 for the remaining copies. Thus *D. simulans* and *D. mauritiana* overall do not meet the
365 expectation that *flamenco* will contain a single insertion of any given TE.

366 How many TEs have full length and fragmented insertions?

367 Full length elements are younger insertions than fragmented insertions. If a full length element is
368 inserted in *flamenco* and there are fragments in the antisense orientation elsewhere in *flamenco*
369 this indicates that *flamenco* did not successfully suppress the transposition of this element.

370 In *D. melanogaster* two elements have fragments in antisense and a full length TE – *Doc*
371 and *Stalker-2*. *D. sechellia* has 9 elements that are present as a full length TE and a fragment in
372 antisense (including *412*, *GTWIN*, *mdg-1*, and *nomad*) and 6 that are multicopy that are due to a
373 solo LTR (including *blood*, *297*, and *Stalker-4*). *D. mauritiana* has 21 elements that are present
374 in full length and a fragment in antisense (including *blood*, *412*, *gypsy-10-13*, and *R1*), and four
375 elements that are multicopy due to a solo LTR (*mdg-1*, *Idefix*, and *gypsy-7,10*).

376 In *D. simulans*, TEs that fit this criteria in *flamenco* include *gypsy-2*, *gypsy-3*, *gypsy-4*,
377 *gypsy-5*, *Chimpo*, *412*, *INE1*, *R1*, *Tirant*, and *Zam*. *297* and *Nomad-1* are present in full length
378 copies but only multi-copy in the context of solo LTRs. In the duplicate of *flamenco* in *D.*
379 *simulans* this includes *gypsy-2*, *gypsy-3*, *gypsy-5*, *297*, *Stalker-4*, and *R1*. For example in *NS40*
380 there are 7 full length copies of *R1* in the sense orientation that likely duplicated in place, as well
381 as 12 partial copies in the antisense orientation. In the *simulans* clade either fragments of TEs are
382 not sufficient to suppress transposable elements or some elements are able to transpose despite
383 the hosts efforts to suppress them.

384 *Is flamenco a trap for TEs entering through horizontal transfer?*

385 High sequence similarity between TEs in different species suggests horizontal transfer [36].

386 However, because sequence similarity can also exist due to vertical transmission we will use

387 sequence similarity between *R1* elements (inserted at rDNA genes) as a baseline for

388 differentiating horizontal versus vertical transfer. There has never been any evidence found for

389 horizontal transfer of *R1* and it is thought to evolve at the same rate as nuclear genes in the

390 *melanogaster* subgroup [21,72]. Of the full length elements present in any genome at *flamenco*
391 62% of them appear to have originated from horizontal transfer. This is similar to previous
392 estimates for *D. melanogaster* in other studies [72]. Transfer appears to have occurred primarily
393 between *D. melanogaster*, *D. sechellia*, and *D. willistoni*. This includes some known horizontal
394 transfer events such as *Chimpo* and *Chouto* [7], and others which have not been recorded such as
395 *gypys-29* (*D. willistoni*) and the *Max-element* (*D. sechellia*) (Supplemental File 3). The duplicate
396 of *flamenco* is similar, with 53% of full length TEs originating from horizontal transfer. They are
397 many of the same TEs, with a 46% overlap, thus *flamenco* and its duplicate are trapping many of
398 the same TEs. Both *flamenco* and the duplicate the region appears to serve as a trap for TEs
399 originating from horizontal transfer.

400 In *D. melanogaster* 85% of full length TEs appear to have arisen through horizontal
401 transfer, primarily with *D. yakuba* and *D. sechellia* [72]. In *D. sechellia* 53% of full length TEs
402 have arisen from horizontal transfer, including some known to have moved by horizontal transfer
403 such as *GTWIN* (*D. melanogaster/D. erecta*) [7]. *D. mauritiana* has 68% of its full length TEs
404 showing a closer relationship than expected by vertical descent with TEs from *D. sechellia*, *D.*
405 *melanogaster*, and *D. simulans*. The hypothesis that *flamenco* serves as a trap for TEs entering
406 the population through horizontal transfer holds throughout the *simulans* clade.

407 *Flamenco piRNA is expressed in the testis and the maternal fraction*

408 Canonically, *flamenco* piRNA is expressed in the somatic follicular cells of the ovary and
409 not in the germline, and also does not produce a ping pong signal [46]. It was not thought to be
410 present in the maternal fraction of piRNAs or other tissues. However, that appears to be variable
411 in different species (Figure 4). We examined single mapping reads in the *flamenco* region from
412 testes and embryos (maternal fraction) in *D. simulans*, *D. mauritiana*, *D. sechellia*, and *D.*

413 *melanogaster*. As a control we also included *D. melanogaster* ovarian somatic cells, where Aub
414 and Ago3 are not expressed and therefore there should be no ping pong signals. In *D. simulans*
415 and *D. mauritiana flamenco* is expressed bidirectionally in the maternal fraction and the testis,
416 including ping pong signals on both strands (Figure 4; Supplementary Figure 1). In *D. sechellia*,
417 there is no expression of *flamenco* in either of these tissues. Discarding multimappers in the
418 maternal fraction 63% (*D. mauritiana*) – 36% (*D. simulans*) of the ping pong signatures on the X
419 with a z-score of at least 0.9 are located within *flamenco* (Figure 4). In the testis the picture is
420 more complicated – in *D. mauritiana* 50% of ping pong signals on the X with a z-score of at
421 least 0.9 are located within *flamenco*, which amounts to a substantial ping pong signature
422 (Supplementary Figure 1). While mapping of piRNA to both strands was observed in *D.*
423 *simulans* testis, there is very little apparent ping pong activity (5 positions in *flamenco* $z > 0.9$;
424 15 potential ping pong signals on the X). In *D. melanogaster*, there is uni-strand expression in
425 the maternal fraction, but it is limited to the region close to the promoter. In *D. melanogaster* no
426 ping pong signals have a z-score above 0.8 in the maternal fraction or the ovarian somatic cells.
427 There are ping pong stacks in *flamenco* in the testis of *D. melanogaster* (2% of the total on the
428 contig), however they are limited to a single region and are not abundant enough to be strong
429 evidence of ping activity.

430 In the duplicate of *flamenco* in the maternal fraction 15% of the ping pong signals with a
431 z-score above 0.9 on the X are within the *flamenco* duplicate. The *flamenco* duplicate does not
432 have a strong signal of the ping pong pathway in the testis. In addition, *flamenco* in these species
433 has been colonized by full length TEs thought to be germline TEs such as *blood*, *burdock*, *mdg-*
434 *3*, *Transpac*, and *Bel* [16,20]. *blood* is also present in *D. melanogaster* in a full length copy while
435 there is no evidence of germline activity for *flamenco* in *D. melanogaster*, though no other

436 putative germline TEs are present. The differences in ping pong signals between species and the
437 presence of germline TEs in *D. simulans* and *D. mauritiana* suggests that the role of *flamenco* in
438 these tissues has evolved between species.

439 **Discussion**

440 The piRNA pathway is the organisms primary mechanism of transposon suppression.
441 While the piRNA pathway is conserved, the regions of the genome that produce piRNA are
442 labile, particularly in double stranded germline piRNA clusters [23]. The necessity of any single
443 cluster for TE suppression in the germline piRNA pathway is unclear, but likely redundant [23].
444 However, *flamenco* is thought to be the master regulator of the somatic support cells of the
445 ovary, preventing *gypsy* elements from hopping into germline cells [19,42,45,46,48,72]. It is not
446 redundant to other clusters, and insertion of a single element into *flamenco* in *D. melanogaster* is
447 sufficient to initiate silencing. Here we show that the function of *flamenco* appears to have
448 diversified in the *D. simulans* clade, acting in at least some tissues as a germline piRNA cluster.

449 *Dual stranded expression of flamenco*

450 In this work, we showed that piRNAs of the *flamenco* locus in *D. simulans* and *D.*
451 *mauritiana* are deposited maternally, align to both strands, and exhibit ping-pong signatures.
452 This is in contrast to *D. melanogaster*, where *flamenco* acts as a uni-strand cluster in the soma
453 [40], our data thus suggest that the *flamenco* locus in *D. simulans* and *D. mauritiana* acts as a
454 dual-strand cluster in the germline. In *D. sechellia* the attributes of *flamenco* uncovered in *D.*
455 *melanogaster* appear to be conserved – no expression in the maternal fraction and the testis and
456 no ping pong signals. Given that *flamenco* is likely a somatic uni-strand cluster in *D. erecta*, we
457 speculate that the conversion into a germline cluster happened in the *simulans* clade [40]. Such a
458 conversion of a cluster between the somatic and the germline piRNA pathway is not

459 unprecedented. For example, a single insertion of a reporter transgene triggered the conversion
460 of the uni-stranded cluster *20A* in *D. melanogaster* into a dual-strand cluster [37].

461 The role of *flamenco* in *D. simulans* and *D. mauritiana* as the master regulator of piRNA
462 in somatic support cells may still well be true – the promoter region of the *flamenco* cluster is
463 conserved between species and between copies of *flamenco* within species. This suggests that in
464 at least some contexts (or all) the cluster is still serving as a uni-strand cluster transcribed from a
465 traditional RNA Pol II site [24]. However it has acquired additional roles, producing dual strand
466 piRNA and ping pong signals, in these two species, in at least the germline. However, in *D.*
467 *simulans*, the majority of these reverse stranded piRNAs are emerging from the *RI* insertions
468 within *flamenco*. There is no evidence at present that *RI* has undergone an expansion in function
469 in *D. simulans*, thus it is unclear what, if any, functional impact the reverse stranded piRNAs
470 have at the *flamenco* locus.

471 *Duplication of flamenco in D. simulans*

472 In *D. simulans*, *flamenco* is present in 2-3 genomic copies, and this duplication is present
473 in all sequenced *D. simulans* lines. The *dip1* gene and putative *flamenco* promoter flanking the
474 duplication also has a high similarity in all sequenced lines (Fig. 2B). This raises the possibility
475 that the duplication of *flamenco* in *D. simulans* was positively selected. Such a duplication may
476 be beneficial as it increases the ability of an organism to rapidly silence TEs. Individuals with
477 large piRNA clusters (or duplicated ones) will accumulate fewer deleterious TE insertions than
478 individuals with small clusters (or non-duplicated ones), and duplicated clusters may therefore
479 confer a selective advantage [27].

480 *Rapid evolution of piRNA clusters*

481 A previous work showed that dual- and uni-strand clusters evolve rapidly in *Drosophila*
482 [70]. In agreement with this work we also found that the *flamenco*-locus is rapidly evolving
483 between and within species (Fig. 1C, 3B). A major open question remains whether this rapid
484 turnover is driven by selection (positive or negative) or an outcome of neutral processes (eg. high
485 TE activity or insertion bias of TEs). These rapid evolutionary changes at the *flamenco* locus, a
486 piRNA master locus, suggest that there is a constant turnover in patterns of piRNA biogenesis
487 that potentially leads to changes in the level of transposition control between individuals in a
488 population.

489

490

491 **Funding**

492 This work was supported by the National Science Foundation Established Program to Stimulate
493 Competitive Research (NSF-EPSCoR-1826834 and NSF-EPSCoR-2032756)
494 to SS and the Austrian Science Fund FWF (<https://www.fwf.ac.at/>;) grant P35093 to RK. J.V. was
495 supported by a Pathway to Independence award from the National Institute of General Medical
496 Sciences (K99-GM137077). E.C.L. was supported by the National Institute of General Medical
497 Sciences (R01-GM083300) and National Institutes of Health MSK Core Grant (P30-CA008748).

498

499

500

501 **Competing interests**

502 We declare that we have no competing interests.

503

504 **Acknowledgements**

505 S.S. would like to thank C. & F. & S. Emery for insightful commentary on the manuscript.

506

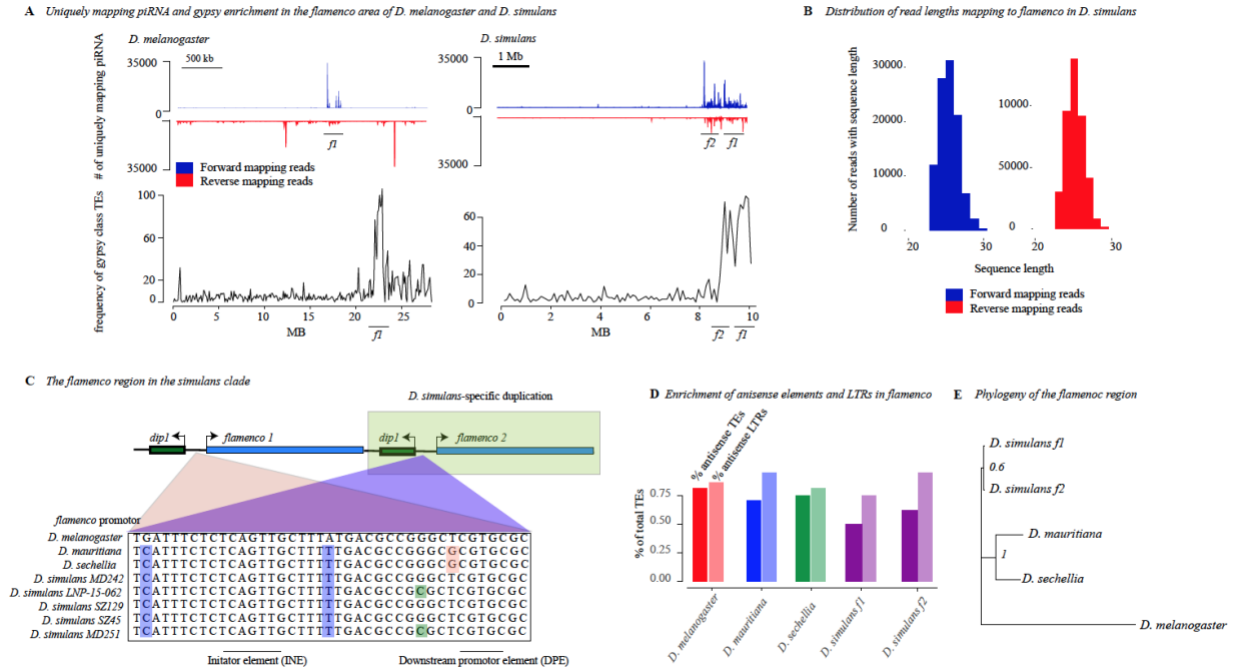
507 **Authors' contributions**

508 S.S. conceived the study, performed bioinformatics and drafted portions of the manuscript. FW
509 and RK performed bioinformatics and drafted portions of the manuscript. JV contributed data
510 and bioinformatic analysis. EL drafted portions of the manuscript and provided data.

511

512 **Availability of data and materials**

513 All data has been made available in the following repositories:



514

515

516 **Figure 1.** A) Unique piRNA from the ovary and gypsy enrichment around *flamenco* and its
 517 duplicate in *D. simulans* and *D. melanogaster*. piRNA mapping to the entire contig that contains

518 *flamenco* is shown for both species. The top of the panel shows piRNA mapping to *flamenco* and

519 is split by antisense (blue) and sense (red) piRNA. The bottom panel shows the frequency of

520 gypsy-type transposon annotations across the contig containing *flamenco*, counted in 100 kb

521 windows. There is a clear enrichment of gypsy in the area of *flamenco* and, in *D. simulans*, its

522 duplicate compared to the rest of the contig. B) The distribution of read size for small RNA

523 mapping to *flamenco*. The peak is at approximately 26 bp, within the expected range for piRNA.

524 C) The duplication of *flamenco* in the *D. simulans*. Both copies are flanked by the *dip1* gene and

525 copies of the putative *flamenco* promoter. Polymorphisms within the promoter that are shared

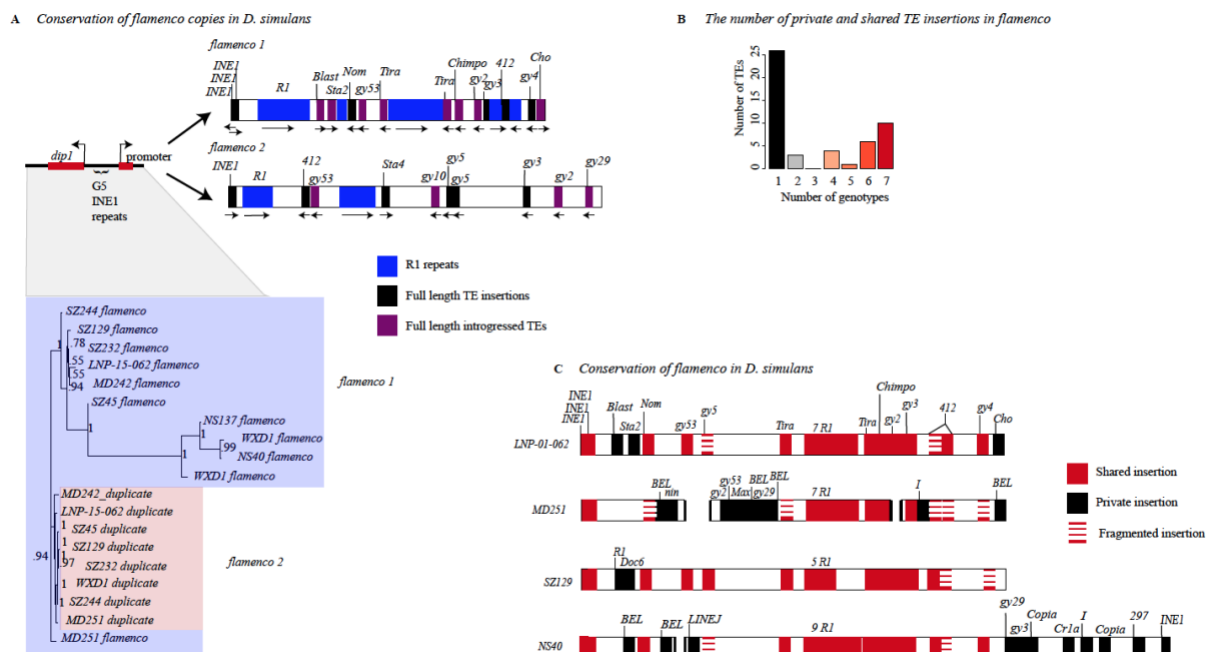
526 within the *simulans* clade are shown in blue, *D. simulans* specific polymorphisms are shown in

527 green. The region around the promoter is very conserved across species. D) The percent of TEs

528 in *flamenco* in each species which are in the antisense orientation (first bar) and the percent of

529 TEs in the antisense orientation that are also LTR class elements (second bar). E) A phylogenetic
 530 tree of the *dip1* and *flamenco* enhancer region for *D. melanogaster* and the *simulans* clade. This
 531 region is conserved and alignable between all species. The tree was generated with Mr. Bayes
 532 [51].

533
 534
 535
 536
 537



538
 539
 540

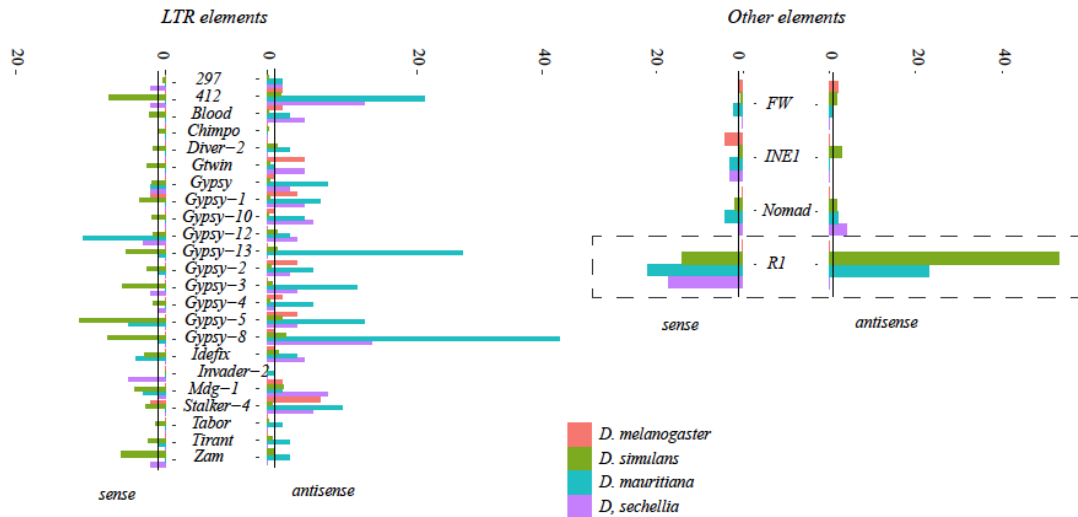
541
 542 Figure 2. A) Divergence between copies of flamenco. Proximal is a phylogenetic tree of *dip1* and
 543 the *flamenco* promoter region from each genome. In between *dip1* and the promoter are a series
 544 of *G5/INE1* repeats that are found in every genome. Overall this region is fairly conserved, with
 545 the duplicate copies all grouping together with short branch lengths (shown in pink). The original
 546 copy of *flamenco* is more diverse with some outliers (shown in light blue) but there is good

547 branch support for all the deep branches of the tree. Distal is a representation of *flamenco* and its
548 duplicate. R1 repeat regions are shown in blue. Full length transposable elements are labeled.
549 There is no synteny conservation between *flamenco* and its duplicate. B) The proportion of
550 insertions that are shared by one through seven genotypes (genotypes with complete *flamenco*
551 assemblies). C) Divergence of *flamenco* within *D. simulans*. Labeled TEs correspond to
552 elements which are present in a full length copy in at least one genome. If they are shared
553 between genomes they are labeled in red, if they are unique they are black. If they are full length
554 in one genome and degraded in other genomes they are represented by stacked dashes. If they are
555 present in the majority of genomes but missing in one, it is represented as a missing that TE,
556 which is agnostic to whether it is a deletion or the element was never present

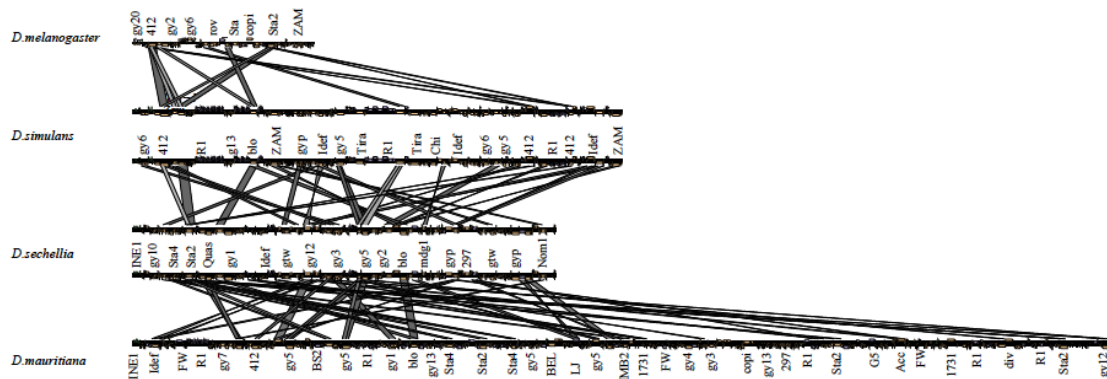
557

558

A Copy number of a subset of TEs in the *simulans* clade



B Similarity of TEs in *flamenco* within the *simulans* clade



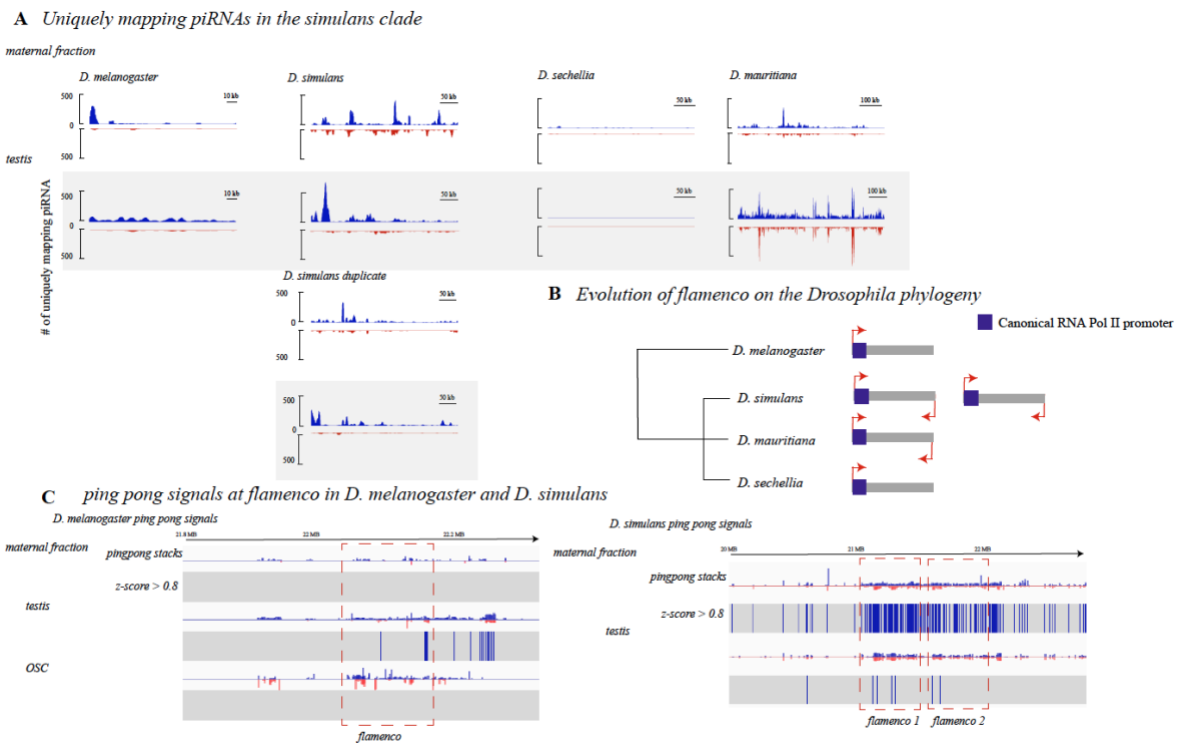
559

560

561

562 **Figure 3)** A. Copy number of a subset of transposable elements at *flamenco*. Solo LTRs are
 563 indicated by in a lighter shade at the top of the bar. The black line on each bar graph indicates a
 564 copy number of one. Values for *D. simulans* are the average for all genotypes with a complete
 565 *flamenco* assembly. Note that in *D. melanogaster* (green) most TEs have a low copy number.
 566 The expansion of *R1* elements in the *simulans* clade is clearly indicated on the right hand panel
 567 with a dotted box. Many elements within *flamenco* are multicopy in the *simulans* clade. While
 568 some of this is likely due to local duplications it is clearly a different pattern than *D.*

569 *melanogaster*. Enrichment of LTR elements on the antisense strand is clear for all species. **B.**
 570 Alignment of *flamenco* in *D. melanogaster*, *D. simulans*, *D. sechellia*, and *D. mauritiana*. There
 571 is no conserved synteny between species but there are clearly shared TEs, particularly within the
 572 *simulans* clade. The expansion of *D. mauritiana* compared to the other species is apparent.



573 -
 574 -
 575 -
 576 **Figure 4) A.** Expression of single mapping piRNAs in the maternal fraction and testis (gray) of
 577 *D. melanogaster* and the *simulans* clade. Antisense mapping reads are shown in blue, sense in
 578 red. Libraries are RPM normalized and scaled across library type. *D. sechellia* has no expression
 579 of *flamenco* in the maternal fraction or the testis. *D. melanogaster* has low expression in the
 580 maternal fraction and very little ping pong activity. *D. simulans* and *D. mauritiana* show dual
 581 stranded expression in the testis and maternal fraction. **B.** A schematic of the evolution of
 582 *flamenco* and its mode expression in the *simulans* and *melanogaster* clade. **C.** The height of 10
 583 nt pingpong stacks at *flamenco* in *D. melanogaster* maternal fraction, testis and ovarian somatic

584 cells is shown on the left. Below each schematic of the height of the stacks is the position of z-
585 scores over 0.8, indicating the likelihood that this is a real ping pong signal as opposed to an
586 artifact. In the testis a few ping pong signals reach this threshold but not enough to indicate
587 convincingly that there is ping pong activity. On the right are the ping pong stacks and z-scores
588 for the maternal fraction and testis in *D. simulans*. Only in the maternal fraction are the density
589 of z-scores over 0.8 convincing enough to indicate an active ping pong cycle in the *flamenco*
590 region. However, the presence of stacks is enriched in testis, thus this may warrant further
591 investigation. *D. mauritiana* also has convincing ping pong signals in this region (Supplementary
592 Figure 1).

593

594

595

596

597 1. C. Duc, *et al.*, Trapping a somatic endogenous retrovirus into a germline piRNA cluster
598 immunizes the germline against further invasion. *Genome Biol* 20, 127 (2019).

599 2. B. Barckmann, *et al.*, The somatic piRNA pathway controls germline transposition over
600 generations. *Nucleic Acids Res* 46, gky761- (2018).

601 3. C. D. Malone, *et al.*, Specialized piRNA Pathways Act in Germline and Somatic Tissues of
602 the *Drosophila* Ovary. *Cell* 137, 522–535 (2009).

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

605 5. S. H. Wang, S. C. R. Elgin, *Drosophila* Piwi functions downstream of piRNA production
606 mediating a chromatin-based transposon silencing mechanism in female germ line. *Proc*
607 *National Acad Sci* 108, 21164–21169 (2011).

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

610 7. A. A. Aravin, *et al.*, The Small RNA Profile during *Drosophila melanogaster* Development.
611 *Developmental Cell* 5, 337–350 (2003).

- 612 8. G. Chirn, *et al.*, Conserved piRNA Expression from a Distinct Set of piRNA Cluster Loci in
613 Eutherian Mammals. *Plos Genet* 11, e1005652 (2015).
- 614 9. D. Gebert, *et al.*, Large Drosophila germline piRNA clusters are evolutionarily labile and
615 dispensable for transposon regulation. *Mol Cell* 81, 3965-3978.e5 (2021).
- 616 10. P. R. Andersen, L. Tirian, M. Vunjak, J. Brennecke, A heterochromatin-dependent
617 transcription machinery drives piRNA expression. *Nature* 549, 54–59 (2017).
- 618 11. C. Klattenhoff, *et al.*, The Drosophila HP1 Homolog Rhino Is Required for Transposon
619 Silencing and piRNA Production by Dual-Strand Clusters. *Cell* 138, 1137–1149 (2009).
- 620 12. F. Mohn, G. Sienski, D. Handler, J. Brennecke, The Rhino-Deadlock-Cutoff Complex
621 Licenses Noncanonical Transcription of Dual-Strand piRNA Clusters in Drosophila. *Cell* 157,
622 1364–1379 (2014).
- 623 13. Y.-C. A. Chen, *et al.*, Cutoff Suppresses RNA Polymerase II Termination to Ensure
624 Expression of piRNA Precursors. *Mol Cell* 63, 97–109 (2016).
- 625 14. F. Mohn, G. Sienski, D. Handler, J. Brennecke, The Rhino-Deadlock-Cutoff Complex
626 Licenses Noncanonical Transcription of Dual-Strand piRNA Clusters in Drosophila. *Cell* 157,
627 1364–1379 (2014).
- 628 15. C. Goriaux, S. Dasset, Y. Renaud, C. Vaury, E. Brasset, Transcriptional properties and
629 splicing of the flamencopi RNAcluster. *EMBO reports* 15, 411–418 (2014).
- 630 16. G. Sienski, D. Dönertas, J. Brennecke, Transcriptional Silencing of Transposons by Piwi and
631 Maelstrom and Its Impact on Chromatin State and Gene Expression. *Cell* 151, 964–980 (2012).
- 632 17. C. Dennis, E. Brasset, C. Vaury, flam piRNA precursors channel from the nucleus to the
633 cytoplasm in a temporally regulated manner along Drosophila oogenesis. *Mobile DNA* 10, 203–9
634 (2019).
- 635 18. V. Zanni, A. Eymery, M. C. P. of the, 2013, Distribution, evolution, and diversity of
636 retrotransposons at the flamenco locus reflect the regulatory properties of piRNA clusters.
637 *National Acad Sciences* <https://doi.org/10.1073/pnas.1313677110/-/dcsupplemental>.
- 638 19. F. Wierzbicki, R. Kofler, S. Signor, Evolutionary dynamics of piRNA clusters in Drosophila.
639 *Mol Ecol* (2021) <https://doi.org/10.1111/mec.16311>.
- 640 20. C. M. Bergman, H. Quesneville, D. Anxolabéhère, M. Ashburner, Recurrent insertion and
641 duplication generate networks of transposable element sequences in the Drosophila melanogaster
642 genome. *Genome Biology* 7, R112-21 (2006).
- 643 21. N. Prud'homme, M. Gans, M. Masson, C. Terzian, A. Bucheton, Flamenco, a gene
644 controlling the gypsy retrovirus of Drosophila melanogaster. *Genetics* 139, 697–711 (1995).

- 645 22. S. U. Song, T. Gerasimova, M. Kurkulos, J. D. Boeke, V. G. Corces, An env-like protein
646 encoded by a Drosophila retroelement: evidence that gypsy is an infectious retrovirus. *Genes &*
647 *development* 8, 2046–2057 (1994).
- 648 23. M. Mével-Ninio, A. Pelisson, J. Kinder, A. R. Campos, A. Bucheton, The flamenco Locus
649 Controls the gypsy and ZAM Retroviruses and Is Required for Drosophila Oogenesis. *Genetics*
650 175, 1615–1624 (2007).
- 651 24. A. Pelisson, *et al.*, Gypsy transposition correlates with the production of a retroviral
652 envelope-like protein under the tissue-specific control of the Drosophila flamenco gene. *The*
653 *EMBO Journal* 13, 4401–4411 (1995).
- 654 25. A. Bucheton, The relationship between the flamenco gene and gypsy in Drosophila: how to
655 tame a retrovirus. *Trends Genet* 11, 349–353 (1995).
- 656 26. C. D. Malone, G. J. Hannon, Molecular Evolution of piRNA and Transposon Control
657 Pathways in Drosophila. *Cold Spring Harbor Symposia on Quantitative Biology* 74, 225–234
658 (2010).
- 659 27. A. G. Clark, *et al.*, Evolution of genes and genomes on the Drosophila phylogeny. *Nature*
660 450, 203–218 (2007).
- 661 28. D. G. Eickbush, W. C. Lathe, M. P. Francino, T. H. Eickbush, R1 and R2 retrotransposable
662 elements of Drosophila evolve at rates similar to those of nuclear genes. *Genetics* 139, 685–695
663 (1995).
- 664 29. S. A. Signor, F. N. New, S. Nuzhdin, A Large Panel of Drosophila simulans Reveals an
665 Abundance of Common Variants. *Genome Biology and Evolution* 10, 189–206 (2017).
- 666 30. S. Signor, S. Nuzhdin, Dynamic changes in gene expression and alternative splicing mediate
667 the response to acute alcohol exposure in Drosophila melanogaster. *Heredity* (2018).
- 668 31. S. Signor, Population genomics of Wolbachia and mtDNA in Drosophila simulans from
669 California. *Scientific Reports*, 1–11 (2017).
- 670 32. S. A. Signor, M. Abbasi, P. Marjoram, S. V. Nuzhdin, Social effects for locomotion vary
671 between environments in Drosophila melanogaster females. *Evolution* 71, 1765–1775 (2017).
- 672 33. S. Signor, Transposable elements in individual genotypes of Drosophila simulans. *Ecology*
673 *and Evolution* 130, 499–11 (2020).
- 674 34. D. R. Matute, J. Gavin-Smyth, G. Liu, Variable post-zygotic isolation in Drosophila
675 melanogaster/D. simulanshybrids. *Journal of Evolutionary Biology* 27, 1691–1705 (2014).

- 676 35. D. R. Schrider, J. Ayroles, D. R. Matute, A. D. Kern, Supervised machine learning reveals
677 introgressed loci in the genomes of *Drosophila simulans* and *D. sechellia*. *PLoS Genetics* 14,
678 e1007341-29 (2018).
- 679 36. R. L. Rogers, *et al.*, Landscape of Standing Variation for Tandem Duplications in *Drosophila*
680 *yakuba* and *Drosophila simulans*. *Molecular Biology and Evolution* 31, 1750–1766 (2014).
- 681 37. M. Chakraborty, *et al.*, Evolution of genome structure in the *Drosophila simulans* species
682 complex. *Genome Res.* 30, 1067–63 (2020).
- 683 38. , *Genome Res.* 2017-Koren-gr.215087.116.
- 684 39. R. Vaser, I. Sović, N. Nagarajan, M. Šikić, Fast and accurate de novo genome assembly from
685 long uncorrected reads. *Genome Res* 27, 737–746 (2017).
- 686 40. B. J. Walker, *et al.*, Pilon: An Integrated Tool for Comprehensive Microbial Variant
687 Detection and Genome Assembly Improvement. *Plos One* 9, e112963 (2014).
- 688 41. M. Kolmogorov, J. Yuan, Y. Lin, P. A. Pevzner, Assembly of long, error-prone reads using
689 repeat graphs. *Nat Biotechnol* 37, 540–546 (2019).
- 690 42. D. R. Laetsch, M. L. Blaxter, BlobTools: Interrogation of genome assemblies.
691 *F1000research* 6, 1287 (2017).
- 692 43. M. Tarailo-Graovac, N. Chen, Using RepeatMasker to Identify Repetitive Elements in
693 Genomic Sequences. *Current Protocols in Bioinformatics*, 1–14 (2009).
- 694 44. J. M. Flynn, *et al.*, RepeatModeler2 for automated genomic discovery of transposable
695 element families. *Proc National Acad Sci* 117, 9451–9457 (2020).
- 696 45. J. Armstrong, *et al.*, Progressive Cactus is a multiple-genome aligner for the thousand-
697 genome era. *Nature* 587, 246–251 (2020).
- 698 46. M. Kolmogorov, *et al.*, Chromosome assembly of large and complex genomes using multiple
699 references. *Genome Res* 28, 1720–1732 (2018).
- 700 47. F. Wierzbicki, F. Schwarz, O. Cannalunga, R. Kofler, Generating high quality assemblies for
701 genomic analysis of transposable elements. *Biorxiv*, 2020.03.27.011312 (2020).
- 702 48. F. Wierzbicki, F. Schwarz, O. Cannalunga, R. Kofler, Novel quality metrics allow
703 identifying and generating high-quality assemblies of piRNA clusters. *Mol Ecol Resour* 22, 102–
704 121 (2022).
- 705 49. Vedanayagam, Jeffrey, “Evolutionary Genomics of piRNA Mediated Transposon Silencing
706 in *Drosophila*,” University of Rochester. (2016).

- 707 50. J. Vedanayagam, *et al.*, Endogenous RNAi silences a burgeoning sex chromosome arms race.
708 *Biorxiv*, 2022.08.22.504821 (2022).
- 709 51. J. Vedanayagam, C.-J. Lin, E. C. Lai, Rapid evolutionary dynamics of an expanding family
710 of meiotic drive factors and their hpRNA suppressors. *Nat Ecol Evol* 5, 1613–1623 (2021).
- 711 52. S. Chen, Y. Zhou, Y. Chen, J. Gu, fastp: an ultra-fast all-in-one FASTQ preprocessor.
712 *Biorxiv*, 274100 (2018).
- 713 53. M. J. Axtell, ShortStack: Comprehensive annotation and quantification of small RNA genes.
714 *RNA* 19, 740–751 (2013).
- 715 54. B. Langmead, C. Trapnell, M. Pop, S. L. Salzberg, Ultrafast and memory-efficient alignment
716 of short DNA sequences to the human genome. *Genome Biol* 10, R25 (2009).
- 717 55. H. Li, *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* 25, 2078–
718 2079 (2009).
- 719 56. Y. Liao, G. K. Smyth, W. Shi, The R package Rsubread is easier, faster, cheaper and better
720 for alignment and quantification of RNA sequencing reads. *Nucleic Acids Res* 47, gkz114-
721 (2019).
- 722 57. H. Li, Minimap2: pairwise alignment for nucleotide sequences. *Genome Biology* 34, 3094–
723 3100.
- 724 58. D. Rosenkranz, H. Zischler, proTRAC - a software for probabilistic piRNA cluster detection,
725 visualization and analysis. *Bmc Bioinformatics* 13, 5 (2012).
- 726 59. M. Chakraborty, J. J. Emerson, S. J. Macdonald, A. D. Long, Structural variants exhibit
727 widespread allelic heterogeneity and shape variation in complex traits. *Nature Communications*,
728 1–11 (2019).
- 729 60. M. Bailly-Bechet, A. Haudry, E. Lerat, “One code to find them all”: a perl tool to
730 conveniently parse RepeatMasker output files. *Mobile Dna-uk* 5, 13 (2014).
- 731 61. A. R. Quinlan, I. M. Hall, BEDTools: a flexible suite of utilities for comparing genomic
732 features. *Bioinformatics* 26, 841–842 (2010).
- 733 62. F. Sievers, D. G. Higgins, Clustal Omega for making accurate alignments of many protein
734 sequences. *Protein Sci* 27, 135–145 (2018).
- 735 63. F. Ronquist, *et al.*, MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model
736 Choice Across a Large Model Space. *Systematic Biology* 61, 539–542 (2012).
- 737 64. E. Paradis, J. Claude, K. Strimmer, APE: Analyses of Phylogenetics and Evolution in R
738 language. *Bioinformatics* 20, 289–290 (2004).

- 739 65. S. Uhrig, H. Klein, PingPongPro: a tool for the detection of piRNA-mediated transposon-
740 silencing in small RNA-Seq data. *Bioinformatics* 35, 335–336 (2018).
- 741 66. E. Lerat, *et al.*, Population specific dynamics and selection patterns of transposable element
742 insertions in European natural populations. *Molecular Ecology*, 1–42 (2018).
- 743 67. R. S. Singh, Population genetics and evolution of species related to *Drosophila melanogaster*.
744 *Annual Review of Genetics* 23, 425–453 (1989).
- 745 68. H. E. Machado, *et al.*, Comparative population genomics of latitudinal variation in
746 *Drosophila simulans* and *Drosophila melanogaster*. *Molecular Ecology* 25, 723–740 (2016).
- 747 69. A. Sedghifar, P. Saelao, D. J. Begun, Genomic patterns of geographic differentiation in
748 *Drosophila simulans*. *Genetics* (2016) <https://doi.org/10.1534/genetics.115.185496>.
- 749 70. D. A. Petrov, DNA loss and evolution of genome size in *Drosophila*. *Genetica* 115, 81–91
750 (2002).
- 751 71. E. L. S. Loreto, C. M. A. Carareto, P. Capy, Revisiting horizontal transfer of transposable
752 elements in *Drosophila*. *Heredity* 100, 545–554 (2008).
- 753 72. N. Bargues, E. Lerat, Evolutionary history of LTR-retrotransposons among 20 *Drosophila*
754 species. *Mobile Dna-uk* 8, 7 (2017).
- 755 73. Z. Durdevic, R. S. Pillai, A. Ephrussi, Transposon silencing in the *Drosophila* female
756 germline is essential for genome stability in progeny embryos. *Life Sci Alliance* 1, e201800179
757 (2018).
- 758 74. B. Czech, J. B. Preall, J. McGinn, G. J. Hannon, A Transcriptome-wide RNAi Screen in the
759 *Drosophila* Ovary Reveals Factors of the Germline piRNA Pathway. *Mol Cell* 50, 749–761
760 (2013).
- 761 75. G. Coline, E. Théron, E. Brassat, C. Vaury, History of the discovery of a master locus
762 producing piRNAs: the flamenco/COM locus in *Drosophila melanogaster*. *Frontiers Genetics* 5,
763 257 (2014).
- 764 76. R. Kofler, Dynamics of Transposable Element Invasions with piRNA Clusters. *Molecular*
765 *Biology and Evolution* 36, 1457–1472 (2019).
- 766 77. A. and T. Péllisson, About the origin of retroviruses and the co-evolution of the gypsy
767 retrovirus with the *Drosophila* flamenco host gene. 29–37 (1997).
- 768 78. C. Duc, *et al.*, Trapping a somatic endogenous retrovirus into a germline piRNA cluster
769 immunizes the germline against further invasion. *Genome Biol* 20, 127 (2019).

- 770 79. Y. Luo, P. He, N. Kanrar, K. F. Toth, A. Aravin, Maternally inherited siRNAs initiate piRNA
771 cluster formation <https://doi.org/10.1101/2022.02.08.479612>.
- 772 80. R. Kofler, piRNA Clusters Need a Minimum Size to Control Transposable Element
773 Invasions. *Genome Biology and Evolution* 12, 736–749 (2020).
- 774 81. F. K. Teixeira, *et al.*, piRNA-mediated regulation of transposon alternative splicing in the
775 soma and germ line. *Nature* 552, 268–272 (2017).
- 776 82. V. V. Kapitonov, J. Jurka, Molecular paleontology of transposable elements in the
777 *Drosophila melanogaster* genome. *Proc National Acad Sci* 100, 6569–6574 (2003).
- 778 83. N. D. Singh, D. A. Petrov, Rapid Sequence Turnover at an Intergenic Locus in *Drosophila*.
779 *Mol Biol Evol* 21, 670–680 (2004).
- 780