

Evolutionary dynamics of piRNA clusters in *Drosophila*

Filip Wierzbicki^{1,2}, Robert Kofler^{1,*} and Sarah Signor^{3,*}

¹Institut für Populationsgenetik, Vetmeduni Vienna, Veterinärplatz 1, 1210 Wien, Austria

²Vienna Graduate School of Population Genetics

³Biological Sciences, North Dakota State University

Abstract

Small RNAs produced from transposable element (TE) rich sections of the genome, termed piRNA clusters, are a crucial component in the genomic defense against selfish DNA. In animals it is thought the invasion of a TE is stopped when a copy of the TE inserts into a piRNA cluster, triggering the production of cognate small RNAs that silence the TE. Despite this importance for TE control, little is known about the evolutionary dynamics of piRNA clusters, mostly because these repeat rich regions are difficult to assemble and compare. Here we establish a framework for studying the evolution of piRNA clusters quantitatively. Previously introduced quality metrics and a newly developed software for multiple alignments of repeat annotations (Manna) allow us to estimate the level of polymorphism segregating in piRNA clusters and the divergence among homologous piRNA clusters. By studying 20 conserved piRNA clusters in multiple assemblies of four *Drosophila* species we show that piRNA clusters are evolving rapidly. While 70-80% of the clusters are conserved within species, the clusters share almost no similarity between species as closely related as *D. melanogaster* and *D. simulans*. Furthermore, abundant insertions and deletions are segregating within the *Drosophila* species. We show that the evolution of clusters is mainly driven by large insertions of recently active TEs, and smaller deletions mostly in older TEs. The effect of these forces is so rapid that homologous clusters often do not contain insertions from the same TE families.

Introduction

Transposable elements (TEs) are short sequences of DNA that multiply within genomes [McClintock, 1956]. TEs are widespread across the tree of life, often making up a significant portion of the genome (2.7-25% in fruit flies, 45% in humans, and 85% in maize [Piegu et al., 2006, Schnable et al., 2009, Lee and Langley, 2012]). TEs also impose a severe mutational load on their hosts by producing insertions that disrupt functional sequences and mediate ectopic recombination [Lim, 1988, Levis et al., 1984, McGinnis et al., 1983]. However, some TE insertions have also been associated with increases in fitness, for example due

*correspondence to rokofler@gmail.com and sarah.signor@ndsu.edu

29 to changes in gene regulation, where they can act as enhancers, repressors, or other regulators of complex
30 gene expression patterns [Daborn et al., 2002, González et al., 2008, Mateo et al., 2014, Casacuberta and
31 González, 2013]. The distribution of fitness effects of TEs is not known, but the majority of insertions are
32 thought to be deleterious [Yang and Nuzhdin, 2003, Dimitri et al., 2003, Lee and Langley, 2012, Adrion
33 et al., 2017].

34 For a long time TEs were thought to be solely counteracted at the population level (transposition/selection
35 balance) [Charlesworth and Charlesworth, 1983, Barrón et al., 2014]. However the discovery of a small
36 RNA-based defense system revealed that they are also actively combated by the host [Brennecke et al.,
37 2007, Lee and Langley, 2010, Blumenstiel, 2011]. This host defense system relies upon PIWI interacting
38 RNAs (piRNAs) that bind to PIWI-clade proteins and suppress TE activity transcriptionally and post-
39 transcriptionally [Brennecke et al., 2007, Gunawardane et al., 2007, Sienski et al., 2012, Le Thomas et al.,
40 2013]. For example in *D. melanogaster* post-transcriptional silencing of TEs is based on Aub and Ago3 which,
41 guided by piRNAs, cleave TE transcripts in the cytoplasm [Kalmykova et al., 2005, Peters and Meister,
42 2007, Brennecke et al., 2007, Gunawardane et al., 2007]. In the nucleus piRNAs guide the Piwi protein to
43 transcribed TEs which, aided by other proteins, transcriptionally silence TEs through the establishment of
44 repressive chromatin marks [Sienski et al., 2012, Le Thomas et al., 2013, Darricarrere et al., 2013]. These
45 piRNAs are produced from discrete regions of the genome termed piRNA clusters, which largely consist of
46 TE fragments [Brennecke et al., 2008]. There is evidence that a single insertion of a TE into a piRNA cluster
47 may be sufficient to initiate piRNA mediated silencing of the TE [Marin et al., 2000, Josse et al., 2007, Zanni
48 et al., 2013]. Therefore, it is assumed that a newly invading TE proliferates in the host until a copy jumps
49 into a piRNA cluster, which triggers the production of piRNAs that silence the TE [Bergman et al., 2006,
50 Malone and Hannon, 2010, Goriaux et al., 2014, Ozata et al., 2019].

51 Despite the central importance of piRNA clusters for the control of TEs, we know very little about
52 how piRNA clusters evolve within and between species. For example, transposition into clusters would
53 be advantageous to hosts if cluster insertions are indeed required for functional silencing of TEs. Then,
54 a general expansion of piRNA clusters would be expected with the invasion of novel TEs. Such invasions
55 may be quite frequent. For example it is likely that four TE families invaded worldwide *D. melanogaster*
56 populations within the last 100years [Schwarz et al., 2021]. Larger or more abundant piRNA clusters in turn
57 will expand the functional target for transposition and may thus be favored. In support of this hypothesis it
58 was suggested that piRNA clusters have largely been gained over the course of evolution [Chirn et al., 2015].
59 However, these claims are difficult to evaluate as studying the evolution of piRNA clusters is challenging
60 for several reasons. First, piRNA clusters are highly repetitive and very difficult to assemble, thus high
61 quality ungapped assemblies of these repetitive regions are required [see for example Wierzbicki et al., 2021]
62 Second, it is challenging to unambiguously identify homologous clusters within and between species. Third,
63 investigating the evolution of the composition of clusters requires reliable alignments of the highly repetitive
64 piRNA clusters. Due to these challenges and the importance of these clusters for TE control, the evolutionary
65 turnover of piRNA clusters is considered to be a central open question in TE biology [Czech et al., 2018].

66 Here, we investigate the evolution of piRNA clusters within and between four *Drosophila* species. By
67 combining long-read based assemblies with a recently developed approach for identifying homologous piRNA
68 clusters (CUSCO, [Wierzbicki et al., 2021]) and a newly developed software for generating multiple alignments
69 of repetitive regions (Manna) we are able to shed light on the evolution of piRNA clusters. While piRNA

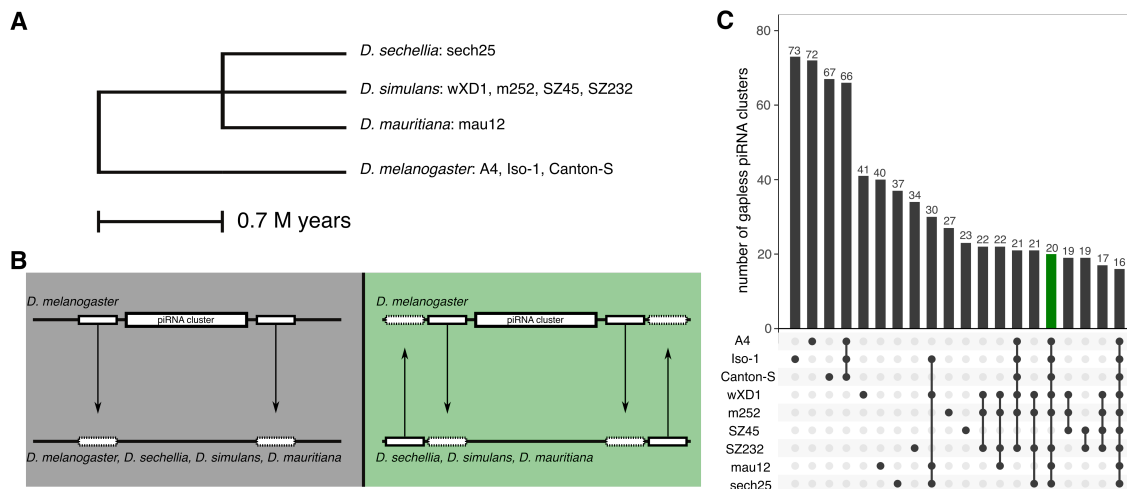


Figure 1: Overview of the species and piRNA clusters used in this work. A) Phylogenetic tree showing the evolutionary distance between the four species investigated in this work (based on [Obbard et al., 2012]). The analyzed strains are shown after the species name. B) Our approach for finding homologous piRNA clusters in the different species and strains. Unique sequences flanking piRNA clusters were aligned to the target strain. An homologous cluster was identified when both flanking sequences aligned to the same contig (grey). We confirmed homology of clusters by designing flanking sequences in the target strain and aligning them back to *D. melanogaster* reference genome (green, "reciprocal flanks"). C) Number of gapless piRNA clusters found in different species/strains. Colors of bar (grey or green) correspond to the approach used for identifying homologous clusters (see B)

70 clusters are 70-80% conserved within species, they share almost no similarity between species as closely
 71 related as *D. melanogaster* and *D. simulans*. Many polymorphic insertions and deletions within clusters
 72 are maintained in *Drosophila* populations. The evolutionary forces dictating the observed patterns appear
 73 to be large insertions of recently active TEs, and smaller deletions of older TE insertion. Due to this
 74 rapid turnover, homologous piRNA clusters frequently do not contain insertions from the same TE families.
 75 Using our approach of combining CUSCO and Manna, we established a framework to study piRNA cluster
 76 evolution quantitatively within and between species.

77 Results

78 Identification of homologous piRNA clusters

79 To shed light on the evolution of piRNA clusters, we compared the composition of clusters among related
 80 *Drosophila* species. *D. sechellia*, *D. mauritiana*, and *D. simulans* are closely related, having an estimated
 81 divergence time of 0.7 million years, while *D. melanogaster* diverged from this group 1.4 million years ago
 82 (fig. 1A, [Obbard et al., 2012]). We relied on long-read assemblies as they allow for end to end reconstruction
 83 of piRNA clusters and their TE content and thus promise to provide a complete picture of cluster evolution
 84 [Wierzbicki et al., 2021]. Since we are interested in the evolution of clusters both within and between
 85 species, we obtained long-read assemblies of several strains for *D. melanogaster* and *D. simulans*. In total

86 we analyzed nine long-read based assemblies, four of *D. simulans*, three of *D. melanogaster*, and one each of
87 *D. sechellia* and *D. mauritiana*. Seven assemblies were publicly available and two assemblies of *D. simulans*
88 strains were generated in this work with Oxford Nanopore long reads (*SZ45*, *SZ232*) [Chakraborty et al.,
89 2021, Nouhaud, 2018, Signor et al., 2017a].

90 The identification of homologous piRNA clusters among the different strains and species was based on
91 unique sequences flanking 85 out of the 142 piRNA clusters in *D. melanogaster* (flanking sequences could
92 not be designed for telomeric clusters extending to the ends of chromosomes or clusters on the fragmented
93 U-chromosome) [Wierzbicki et al., 2021]. These flanking sequences were mapped to each assembly, and
94 homologous piRNA clusters were identified as the regions between the aligned flanking sequences (fig. 1B;
95 grey). piRNA clusters with assembly gaps or flanking sequences aligning to different contigs were not
96 considered. To validate the homology of the piRNA clusters, we designed additional pairs of flanking
97 sequences in the target species, aligned them back to *D. melanogaster* and ascertained that these mapped
98 sequences flank the piRNA clusters of *D. melanogaster* (fig. 1B,C; green; supplementary tables S1-S3). The
99 number of assembled piRNA clusters varied considerably between the strains and species, ranging from 73
100 clusters in *D. melanogaster Iso-1* to 23 clusters in *D. simulans SZ45* (fig. 1C). To study the evolution
101 of piRNA clusters between species, we focused on 20 piRNA clusters shared between *D. mauritiana*, *D.*
102 *sechellia* and the three best assemblies of *D. melanogaster* and *D. simulans* (fig. 1C; red). Most notably our
103 analysis included clusters *42AB* (cluster 1), *20A* (cluster 2) and *38C* (cluster 5) but not *flamenco*. Except
104 for cluster *20A*, which is an uni-strand cluster that is expressed in the germline and the soma, all analyzed
105 clusters are dual-strand clusters that are solely expressed in the germline [Mohn et al., 2014, Brennecke
106 et al., 2007]. By investigating the heterogeneity of the base coverage and the softclip coverage - two recently
107 proposed metrics for identifying assembly errors in piRNA clusters [Wierzbicki et al., 2021] - we ascertained
108 that the assemblies of the 20 clusters are of high quality (see Materials and Methods; supplementary figs. S1-
109 S5). Based on publicly available small RNA data from ovaries of a *D. melanogaster* and *D. simulans* strain
110 collected in Chantemesle (France; [Asif-Laidin et al., 2017]), we found that 15 out of the 20 investigated
111 clusters are expressed in both species (> 10 reads per million; supplementary figs. S6, S7, S8).

112 Comparing the composition of homologous clusters

113 piRNA clusters are often referred to as 'TE graveyards' since they are thought to carry the remains of past
114 TE invasions. This highly repetitive nature makes it difficult to compare the composition of homologous
115 clusters, e.g. using multiple sequence alignments. We approached this problem inspired by the alignments
116 of amino-acid sequences, which are performed at a higher level than the underlying nucleotide sequences.
117 Here, we propose that multiple alignments may be performed with the TE annotations (e.g. generated by
118 RepeatMasker) of piRNA clusters instead of the nucleotide sequences. For this reason, we developed Manna
119 (multiple annotation alignment), a novel tool performing multiple alignments of annotations. Although
120 primarily designed for annotations of repeats, it may work with the annotations of any feature. Manna
121 performs a progressive alignment similar to that described by [Feng and Doolittle, 1987]. Using a simple
122 scoring scheme (supplementary fig. S9) and an adapted Needleman-Wunsch algorithm [Needleman and
123 Wunsch, 1970] a guide tree is computed. Based on this tree the most similar annotations are aligned first,
124 followed by increasingly more distant annotations. For the scoring matrix the score of each newly aligned

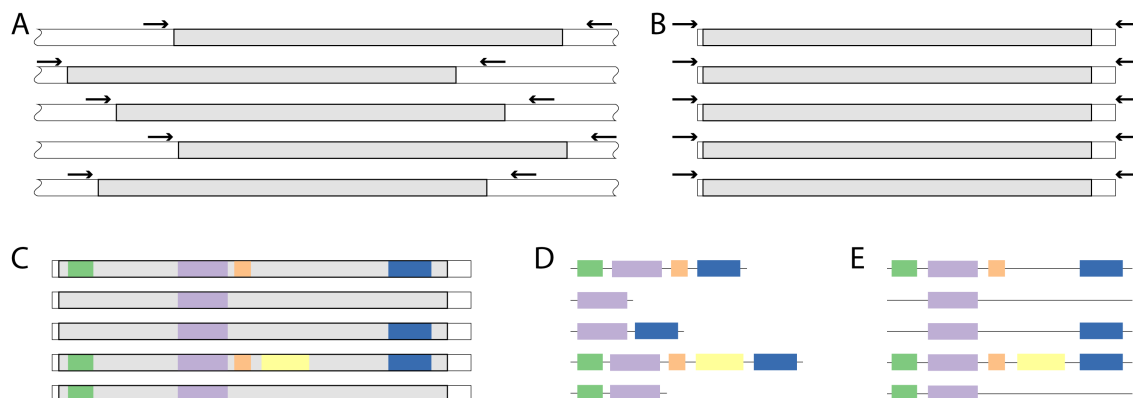


Figure 2: Overview of our approach for comparing the composition of piRNA clusters. A) To identify homologous piRNA clusters (grey areas) in the strains, we mapped sequences flanking the piRNA clusters (black arrows) to the assemblies. B) Regions delimited by the flanking sequences were extracted (i.e. the piRNA clusters plus the short sequences between the clusters and the flanking sequences). C) Repeats were annotated in the extracted sequences. D) Solely the repeat annotations were retained for further analysis. E) The repeat annotations were aligned with Manna allowing us to compare the repeat content of piRNA clusters.

125 annotation is computed as the average score of the previously aligned annotations [Feng and Doolittle, 1987].

126 This novel tool enables us to compare the composition of homologous clusters using the following ap-
 127 proach: First, we align pairs of sequences flanking piRNA clusters to the assemblies, thereby identifying
 128 the positions of homologous clusters in each assembly (fig. 2A). Second, we extract the sequences delimited
 129 by these pairs of flanking sequences (fig. 2B). Third, we annotate repeats in the extracted sequences (fig.
 130 2C) and solely retain the repeat annotation (fig. 2D). Finally, we align the repeat annotation with Manna
 131 (fig. 2E). Using simulated sequences with varying repeat contents, we carefully validated this approach for
 132 comparing the composition of homologous piRNA clusters (supplementary results S1).

133 Alignments with Manna allow us to quantify i) the number of polymorphic and fixed TE insertions and
 134 ii) the similarity s and the distance ($d = 1 - s$) among homologous clusters. The similarity (s) between
 135 clusters is computed as $s = 2 * a / (2 * a + u)$ where a and u are the total length of aligned and unaligned
 136 TE sequences, respectively (for examples see supplementary fig. S10). This similarity can be intuitively
 137 interpreted as the fraction of TE sequences that can be aligned between two (homologous) clusters.

138 Alignments with Manna do not incorporate unannotated sequence in between TEs (fig. 2C). Therefore,
 139 we additionally investigated the similarity among homologous clusters using a complementary approach: we
 140 identified similar sequences between clusters with BLAST (minimum identity 70% [Altschul et al., 1990]) and
 141 visualized these similarities and the repeat content of clusters with Easyfig (supplementary figs. S11-S15).

142 Rapid evolution of piRNA clusters

143 To quantify the rate at which piRNA clusters evolve, we estimated the evolutionary turnover of the TE
 144 content of the 20 piRNA clusters using the similarity (s) as computed with Manna (see above). Based on
 145 the distance between the clusters ($d = 1 - s$), we additionally generated phylogenetic trees reflecting these

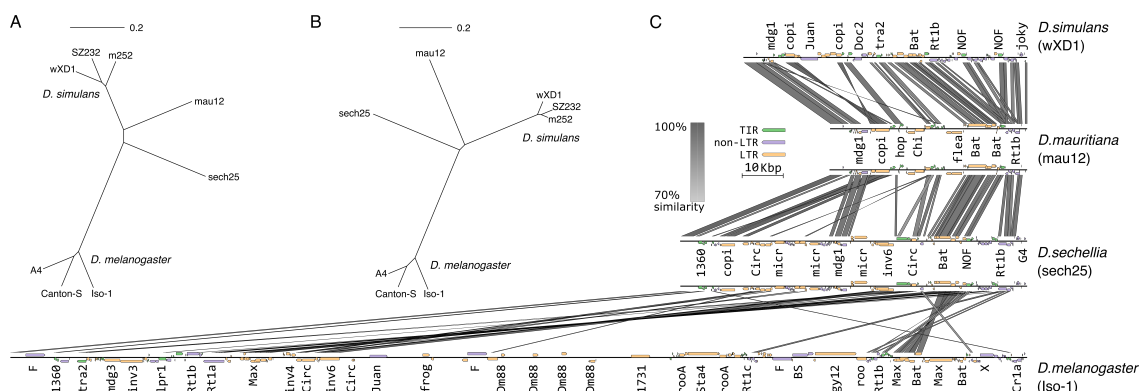


Figure 3: piRNA clusters are rapidly evolving in *Drosophila* species. A) Phylogenetic tree summarizing the distance of the 20 piRNA clusters among the different strains and species weighted by the average cluster lengths. The distance is estimated by Manna as the fraction of unaligned TE sequences (scale bar shows a distance of 20%). Note that solely about 8.1% of the TE sequences can be aligned between the clusters of *D. melanogaster* (green) and *D. simulans* (blue). B) Phylogenetic tree for the piRNA cluster 42AB (cluster 1) based on alignments with Manna. C) The evolution of piRNA cluster 42AB in four *Drosophila* species visualized with Easyfig. Homology among the sequences (grey bars) was determined with BLAST. The grey scale indicates the degree of the sequence similarity. Homology blocks smaller than 400bp are not shown. Insertions of TEs are shown as small rectangular arrows where the color indicates the order (LTR, non-LTR and TIR). Family names are abbreviated.

146 distances (fig. 3A).

147 Strikingly, an average of solely 8.1% of the TE sequences can be aligned between the piRNA clusters of *D.*
 148 *melanogaster* and *D. simulans* (fig. 3A; supplementary table S4). Among the 20 clusters the similarity ranged
 149 from 0.0% for clusters 19 and 110 to 93.5% for cluster 114 (length weighted median: 3.7%; supplementary
 150 table S4). Within the more closely related species of the *simulans* complex 41.4% of the TE sequences can
 151 be aligned between *D. simulans* and *D. mauritiana* (range: 0.0 - 100%; length weighted median: 32.7%
 152) and 32.7% between *D. sechellia* and *D. simulans* (range: 0.0 - 88.8%; length weighted median: 24.8%;
 153 supplementary table S4). Our data thus suggest that the clusters of *D. simulans* are more closely related to
 154 *D. mauritiana* than to *D. sechellia*. Given this rapid turnover within piRNA clusters, we also hypothesized
 155 that there should be abundant polymorphisms within species. In agreement with this, we found that the
 156 average similarity of clusters within species is 73.12% for *D. melanogaster* (range: 33.3-100%; length weighted
 157 median: 74.2%) and 74.7% for *D. simulans* (range: 0.0-100%; length weighted median: 75%; supplementary
 158 table S4). That is to say that on average 26% of the TE sequences in piRNA clusters cannot be aligned
 159 between two assemblies of the same species. The TE content of clusters is thus highly polymorphic within
 160 species.

161 However, the strains analyzed in *D. simulans* and *D. melanogaster* were collected at very diverse time
 162 points and geographic locations. We therefore speculated that the similarity among strains sampled from
 163 the same population may be higher. A comparison of 16 clusters shared between the Californian *D. simulans*
 164 strains SZ232 and SZ45, which were collected at the same location and date, an African strain (m252) and
 165 an old Californian strain (w^xD1, likely collected approximately 50 years prior) did not confirm this hypothesis

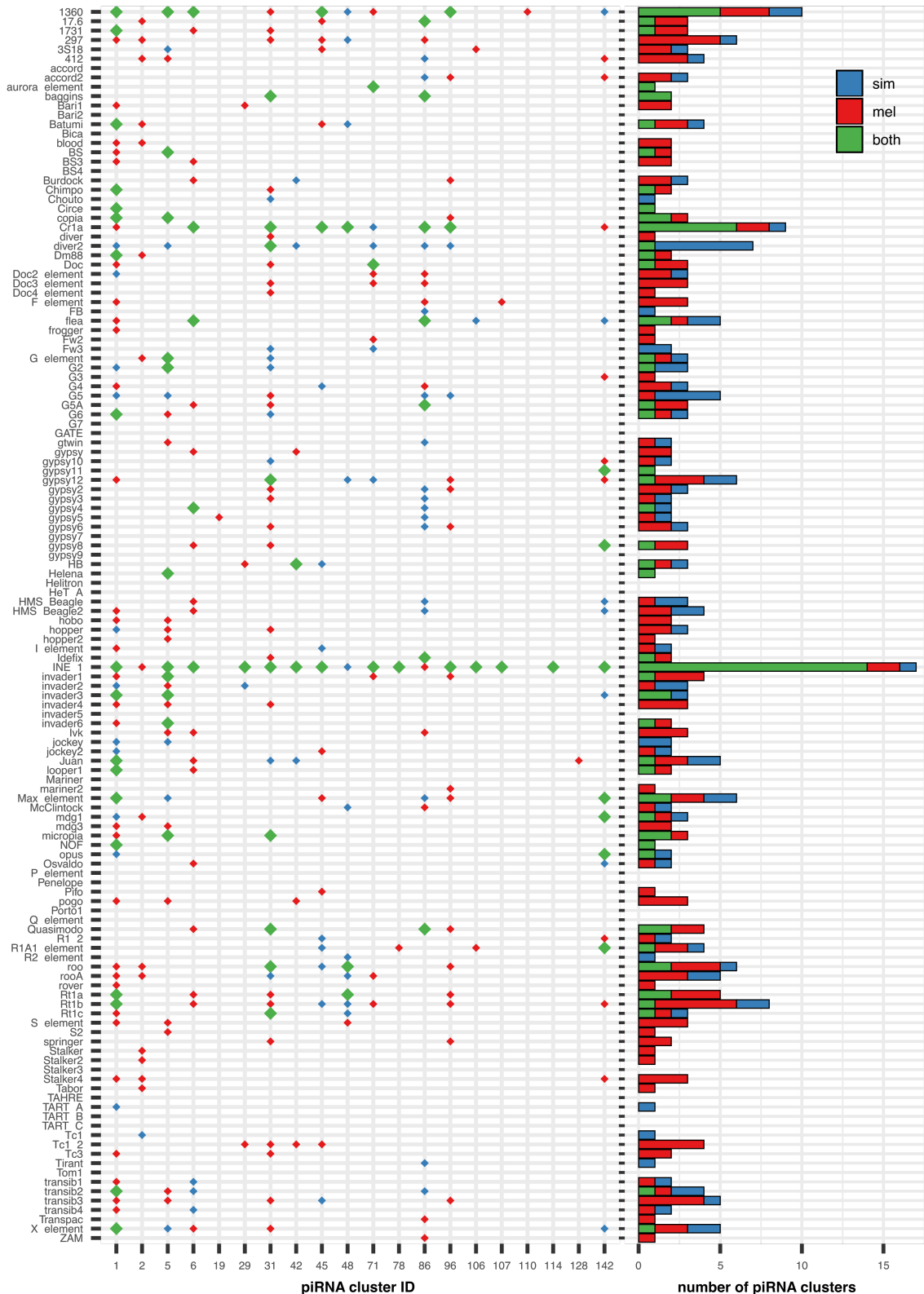


Figure 4: Overview of the TE content of piRNA clusters in *D. simulans* and *D. melanogaster*. For each piRNA cluster (x-axis) we indicate whether a given TE family (y-axis) has at least one insertion in *D. melanogaster* (red), *D. simulans* (blue) or in both species (green). We considered insertions in any of the three assemblies of *D. melanogaster* and *D. simulans*. The right panel summarizes the abundance of the families in piRNA clusters. Note that the TE content of the clusters varies dramatically between the species.

166 (similarity between *SZ232* vs. *SZ45*: 72.5%; average similarity among all other *D. simulans* strains: 75.8%;
167 supplementary table S5). The clusters of strains sampled from the same population are thus not necessarily
168 more similar than the clusters of strains sampled from different regions and time points (although the results
169 vary among the clusters).

170 Next, we aimed to investigate the evolution of cluster *42AB* (cluster 1) in more detail. In *D. melanogaster*
171 *42AB* is one of the largest clusters that may account for 20-30% of all piRNAs [Brennecke et al., 2007]. It
172 is thus frequently highlighted as a canonical piRNA cluster [e.g. Czech et al., 2008, Mohn et al., 2014,
173 Olovnikov et al., 2013, Andersen et al., 2017]. A phylogenetic tree based on an alignment of annotated TEs
174 shows that cluster *42AB* is rapidly evolving among the investigated *Drosophila* species (fig. 3B; for a tree
175 for all other clusters see supplementary fig. S16). The similarity of *42AB* between *D. simulans* and *D.*
176 *melanogaster*, based on an alignment of TE annotations using Manna, is solely 4%. Within the *simulans*
177 clade the similarity of *42AB* between *D. simulans* and *D. mauritiana* is 29.6%, and between *D. simulans*
178 and *D. sechellia* it is 26.4% (supplementary table S4). Within species, cluster *42AB* is more variable in
179 *D. melanogaster* (similarity: 77.5%) than in *D. simulans* (similarity: 90.3% ; supplementary table S4). As
180 alignments with Manna only capture similarities of annotated TEs we also visualized the evolution of cluster
181 *42AB* using BLAST and Easyfig (fig. 3C). This approach confirms our findings. Cluster *42AB* has few
182 sequence similarities between *D. melanogaster* and *D. simulans* and a higher level of sequence similarity
183 among the species of the *simulans* complex (fig. 3C). We conclude that cluster *42AB* is rapidly evolving in
184 the investigated species (fig. 3C). For a visualization of the sequence similarity of all 20 clusters in the four
185 species see supplementary figs. S11-S15.

186 Thus far we have shown that the sequence of piRNAs clusters is evolving very quickly between and within
187 species. However, it is possible that this rapid evolution is due to rearrangements within piRNA clusters
188 [Gebert et al., 2021], while the TE content of clusters actually remains stable. We addressed this question
189 by quantifying the number of insertions from each TE family in each cluster, and determining if at least
190 one insertion of a given family is present in a given cluster in *D. simulans*, *D. melanogaster* or both species
191 (an insertion in any of the three strains of each species was considered as a presence). For example we
192 considered *blood* to be present in cluster *42AB* in both species when a single *blood* insertion was found in
193 *42AB* of *A4* (*D. melanogaster*) and *m252* (*D. simulans*) but not in any other strain of the two species.
194 The rapid evolution of piRNA clusters does not appear to be due to rearrangements, as the presence of TE
195 families was also not conserved across species (fig. 4). Out of 321 TE families in piRNA clusters, only 76
196 were present in both species (families present in more than one cluster were counted multiple times). 164
197 were private to *D. melanogaster* and 81 to *D. simulans* (fig. 4). A similar observation can be made when
198 we compare the TE composition of piRNA clusters among *D. simulans*, *D. mauritiana*, and *D. sechellia*
199 (supplementary fig. S17).

200 We thus conclude that piRNA clusters are rapidly evolving in *Drosophila* species, such that the average,
201 only about 8% of TEs sequences can be aligned between the closely related *D. melanogaster* and *D. simulans*.
202 Furthermore, homologous clusters frequently contain different TE families.

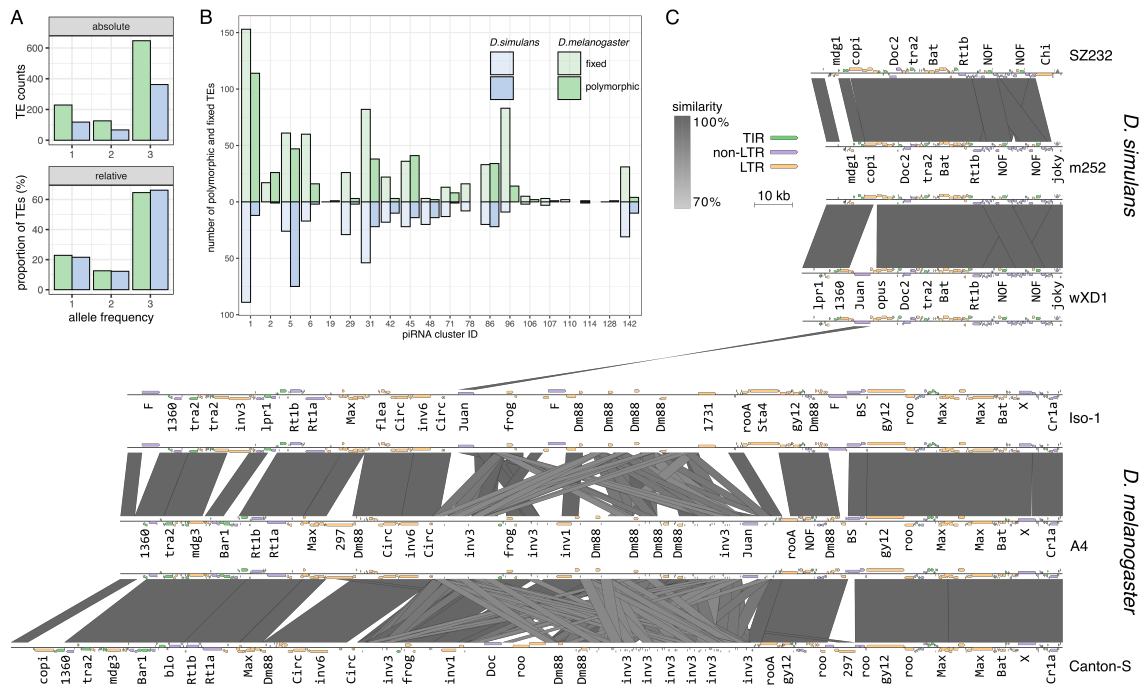


Figure 5: Rapid evolution of piRNA clusters within *D. melanogaster* and *D. simulans*. A) Population frequencies of TE insertions in all 20 piRNA clusters of *D. melanogaster* (green) and *D. simulans* (blue). The absolute (top) and relative (bottom) TE abundance are shown. Insertions occurring in three individuals are fixed. B) Numbers of fixed (transparent) and polymorphic (opaque) sites for each cluster in *D. melanogaster* (green) and *D. simulans* (blue). C) Composition of cluster 42AB in 3 strains of *D. melanogaster* and *D. simulans*. Grey bars indicate regions of similarity among two assemblies of 42AB (minimum length 3 kb). TE families are colored by order (LTR, non-LTR and TIR).

203 piRNA clusters in *D. melanogaster* and *D. simulans* genotypes

204 Next, we investigated variation in the piRNA clusters of *D. melanogaster* and *D. simulans* in more detail,
205 incorporating several genotypes from each species. An alignment of the 20 clusters with Manna in the three
206 strains of *D. melanogaster* and *D. simulans* shows that clusters in *D. melanogaster* contain more TEs than
207 in *D. simulans* ($Dmel = 1,002$, $Dsim = 547$). The majority of these insertions are fixed ($Dmel = 647$,
208 $Dsim = 362$; fig. 5A), but a substantial number of TE insertions is segregating in one ($Dmel = 229$,
209 $Dsim = 118$) or two genotypes ($Dmel = 126$, $Dsim = 67$). Despite these differences in the TE abundance
210 among the two species, the site frequency spectrum of the cluster insertions is very similar between *D.*
211 *melanogaster* and *D. simulans* (Chi-squared test $p = 0.20$; fig. 5A). The large number of polymorphic
212 cluster insertions is not contingent upon a single outlier-genotype since all genotypes from both species carried
213 abundant polymorphic cluster insertions (*D. melanogaster*: $CS = 191$, $A4 = 153$, $Iso1 = 137$; *D. simulans*
214 $SZ232 = 106$, $w^{xD1} = 97$, $m252 = 49$, supplementary fig. S18A). The polymorphic cluster insertions were
215 distributed over 17 clusters in *D. melanogaster* and 12 clusters in *D. simulans* (supplementary fig. S18A). In
216 agreement with the higher TE content of *D. melanogaster* clusters, piRNA clusters in *D. melanogaster* were
217 substantially longer than in *D. simulans* (Wilcoxon rank sum test $W = 2192$, $p = 0.040$; supplementary fig.
218 S18B). The total size of the piRNA clusters in *D. melanogaster* was about double that of the clusters in
219 *D. simulans* (average over all three strains $dmel = 817,770$, $dsm = 452,591$). In both species segregating
220 cluster insertions were on the average longer than fixed ones (*D. melanogaster*: $seg = 1115$, $fix = 591$,
221 Wilcoxon rank sum test $W = 122302$, $p = 0.089$; *D. simulans*: $seg = 798$, $fix = 470$, Wilcoxon rank sum
222 test $W = 38248$, $p = 0.0065$).

223 In addition, the amount of polymorphism segregating in strains sampled from the same population
224 (*SZ232*, *SZ45*) is similar to the amount of polymorphism sampled in strains from different locations (*m252*,
225 Africa) and time points (w^{xD1} , California; percent polymorphic insertions with a minimum size of 100bp:
226 *SZ232* vs *SZ45* = 23.8%, mean of all other pairwise comparisons = 20.7%; supplementary figs. S19, S20).
227 While overall polymorphism was similar amongst strains, the amount of fixed and segregating TE insertions
228 varies across the clusters. Some clusters in *D. melanogaster* mostly have fixed TEs such as cluster 96
229 ($fix = 83$, $seg = 14$) and cluster 142 ($fix = 31$, $seg = 4$), but other clusters, like cluster 1 ($fix = 153$,
230 $seg = 114$) and cluster 45 ($fix = 36$, $seg = 41$), have large proportions of segregating TEs (fig. 5B).
231 Similarly in *D. simulans* some clusters such as cluster 1 ($fix = 89$, $seg = 12$) and cluster 29 ($fix = 29$,
232 $seg = 2$) have largely fixed TEs whereas cluster 5 ($fix = 26$, $seg = 75$) and cluster 86 ($fix = 20$, $seg = 22$)
233 contain many segregating TE insertions. This indicates that clusters may evolve at different rates, with
234 some clusters evolving faster than others. Additionally, the evolutionary turnover of the clusters may differ
235 among species, for example cluster 42AB (cluster 1) evolves faster in *D. melanogaster* whereas cluster 5
236 evolves faster in *D. simulans* (fig. 5B).

237 Our analysis is based on the consensus sequences of *D. melanogaster* TEs. We asked if this could lead to
238 a bias where TE insertions in *D. simulans* clusters are less readily identified than in *D. melanogaster*. Such a
239 bias should lead to a lower density of TEs in piRNA clusters of *D. simulans* as compared to *D. melanogaster*.
240 We found that the density of TE insertions in piRNA clusters is very similar in the two species (TE insertions
241 per kb $dmel = 0.994$, $dsm = 0.985$) suggesting that we identified most TE insertions in *D. simulans*.
242 However, cluster insertions in *D. simulans* were, on the average, slightly shorter than in *D. melanogaster*

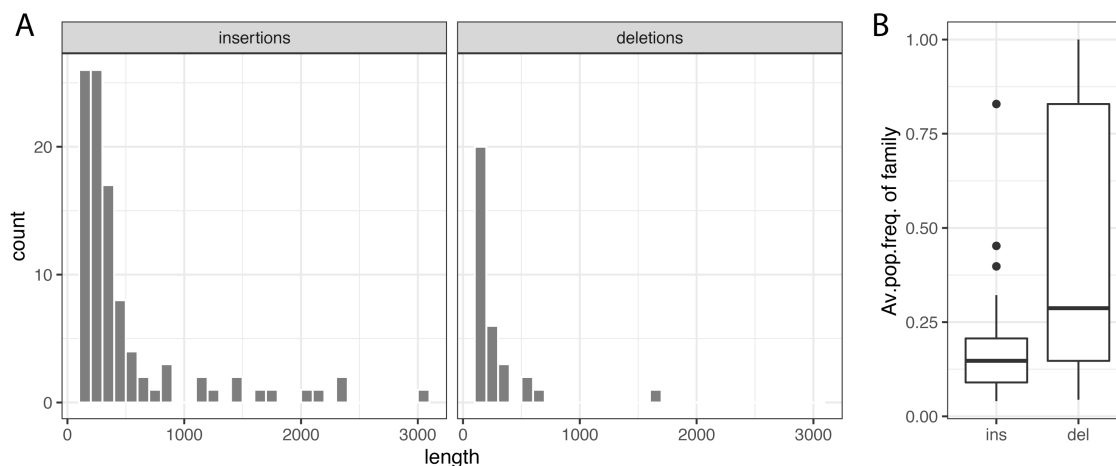


Figure 6: Overview of insertions and deletions in piRNA clusters of *D. simulans*. The clusters of *D. mauritiana* were used to polarize the indels. A) Histograms showing the abundance and length of insertions and deletions. B) Age of the families of insertions (ins) and deletions (del) in piRNA clusters, where the average population frequency (av.pop.freq.) of the family is used as a proxy for the age.

243 (average length *dmel* = 777, *dsim* = 581; Wilcoxon rank sum test $W = 300760$, $p = 0.0015$). This is in
244 agreement with previous works suggesting that TEs in *D. simulans* are shorter than in *D. melanogaster*
245 [Lerat et al., 2011, Vieira et al., 2012], but it could also be a technical artefact where parts of TEs are not
246 annotated in *D. simulans* due to the divergence of the TE from the consensus sequences.

247 Finally, we investigated the composition of cluster *42AB* in more detail (fig. 5D). Cluster *42AB* is,
248 consistently among the strains, shorter in *D. simulans* than in *D. melanogaster* (fig. 5D; supplementary fig.
249 S18B). The density of TEs in cluster *42AB* is higher in *D. simulans* (TEs per kb *dmel* = 0.79, *dsim* = 1.41)
250 possibly due to the shorter TE insertions (average length of TEs in *42AB* *dmel* = 920bp, *dsim* = 658bp).
251 While there is considerable sequence conservation in both species the *D. melanogaster 42AB* cluster bears no
252 resemblance to *42AB* in *D. simulans*, other than containing a *Juan* element which is likely not a homologous
253 insertion (fig. 5B). The number of segregating insertions is larger in *D. melanogaster* than in *D. simulans*
254 suggesting that *42AB* is evolving faster in *D. melanogaster* (fig. 5B,D). For a visualization of the sequence
255 similarity of all clusters in the different assemblies of *D. melanogaster* and *D. simulans* see supplementary
256 figs. S11-S15.

257 We conclude that piRNA clusters are highly polymorphic in both species, that clusters have a similar
258 TE density in both species and that most clusters are shorter in *D. simulans* than in *D. melanogaster*.
259 Furthermore, clusters may evolve at different rates among and within species.

260 Evolutionary forces shaping the composition of piRNA clusters

261 Many diverse evolutionary forces may act on the TE content of piRNA clusters, such as mutations, insertion
262 bias, negative or positive selection and drift [Kofler, 2019, Kelleher et al., 2018, Lu and Clark, 2010, Brennecke
263 et al., 2007, Zhang et al., 2020]. While we cannot distinguish among these forces we can shed light on their
264 joint effect by investigating the abundance of insertions and deletions segregating in piRNA clusters. We

265 determined the number of insertions and deletions segregating in piRNA clusters of the *D. simulans* strains
266 by polarizing segregating indels using *D. mauritiana* as outgroup. We used TE insertions with a minimum
267 length of 100 bp and considered indels resulting from presence/absence polymorphisms in the alignment and
268 indels resulting from length differences between aligned TEs sequences. We found that 33 deletions and 99
269 insertions are segregating in piRNA clusters of *D. simulans* (fig. 6A) These indels were distributed over
270 12 of the investigated 20 piRNA clusters (supplementary fig S21). Insertions were, on the average, longer
271 than deletions (average length $\bar{l}_{ins} = 492bp$, $\bar{l}_{del} = 262bp$; Wilcoxon rank sum test $W = 920.5$, $p = 0.0002$).
272 Most indels were found in three of the 20 clusters: cluster 5 (43 indels), cluster 31 (20 indels), and cluster
273 45 (16 indels; supplementary fig. S21). Because *de novo* TE insertions will likely be large we separately
274 analyzed long indels (≥ 1000). We found that 12 long insertions and a single long deletion. The most
275 abundant long insertions were due to the TE families *roo* and *Max-element* (two for each family). Both
276 families are likely active in *D. simulans* [Kofler et al., 2015, Signor, 2020]. Finally, we asked if insertions
277 are occurring with younger TE families than deletions. While we do not have direct estimates for the age
278 of TE families in *D. simulans* we may use the average population frequency of all insertions of a family as
279 proxy for age. Insertions of recently active families will mostly have a low frequency whereas old families
280 will mostly have fixed insertions. Using the frequency estimates of Kofler et al. [2015] we found that families
281 with insertions in piRNA clusters have a significantly lower average population frequency than families with
282 deletions ($\bar{f}_{ins} = 0.17$, $\bar{f}_{del} = 0.40$; Wilcoxon rank sum test $W = 2211$, $p = 2.7e - 05$ fig. 6B).

283 In summary, the evolutionary dynamics of piRNA clusters are governed by many insertions and few
284 deletions, where insertions are on the average larger than deletions. Furthermore, insertions usually involve
285 recently active families whereas deletions mostly happen in older families.

286 Discussion

287 Here we established a framework for studying the evolution of piRNA clusters quantitatively, used that
288 framework to analyze the composition of 20 piRNA clusters in four *Drosophila* species, and showed that
289 piRNA clusters are evolving rapidly. This raises the question of whether the 20 piRNA clusters included
290 in the analysis are a representative set of the 141 piRNA clusters of *D. melanogaster*. piRNA clusters
291 were excluded from our analysis for three reasons i) clusters were at the end of a chromosome or on the
292 unassembled U-chromosome which did not allow us to identify suitable flanking sequences ii) a cluster could
293 not be assembled in all species without gaps, possibly due to complex repeat content iii) we could not identify
294 conserved flanking sequences in all species such that the homology of a cluster could be established. While
295 the first point likely does not introduce a bias the last two points could potentially result in a bias towards
296 shorter or less complicated clusters. The analyzed clusters may thus be a rather conservative set, and it is
297 possible that the excluded piRNA clusters have different evolutionary dynamics. To reduce possible biases
298 in future works, it will be important to extend the analysis performed in the present work to a larger number
299 piRNA clusters. It is possible that investigating alternate flanking sequences could lead to an increase in the
300 number of clusters, and rapid advances in sequencing technology will increase the number of contiguously
301 assembled clusters. However, a comparison between species will always be less than entirely comprehensive,
302 as clusters may not be shared between species of interest or the flanking sequences may have degraded
303 beyond recognition. In agreement with this, previous research has noted that many piRNA clusters are

304 species specific [Gebert et al., 2021, Chirn et al., 2015].

305 This and other works established synteny of piRNA clusters based on sequences flanking the cluster up
306 and downstream [Gebert et al., 2021, Chirn et al., 2015]. It is unclear if this is the best approach for finding
307 homologous clusters. In principle, it is possible to use the sequence (or annotation) of piRNA clusters directly
308 to search for the homologous clusters in an assembly of interest (e.g. with BLAST). However, given how
309 rapidly piRNA clusters evolve, where solely 8% of TE sequences can be aligned between *D. melanogaster*
310 and *D. simulans*, it is doubtful whether this approach will be able to correctly establish homology of the
311 piRNA clusters. We quantified the similarity of clusters and the amount of polymorphism in clusters with our
312 novel multiple alignment tool Manna. As a major innovation this tool performs a multiple alignment with
313 repeat annotations rather than the raw sequences. While this approach provides invaluable insight into the
314 evolution of piRNA clusters, it does ignore some information such as divergence of the TEs. Alignments of
315 clusters at the nucleotide level may be more sensitive. But this approach has its own problems. Alignments
316 of highly repetitive regions are challenging and may contain errors. Furthermore, the resulting alignment
317 may be difficult to interpret. For example, it is unclear how to estimate the population frequency of a TE
318 insertion where different parts of the TE align with several TE insertions in a homologous cluster. Manna
319 avoids this fragmentation of TEs by aligning complete chunks of annotated TEs.

320 We found that *D. simulans* has fewer TE insertions in piRNA clusters than *D. melanogaster*. That this
321 is a real pattern is supported by the similar density of TEs in the two species within the piRNA clusters
322 (indicating no obvious presence of unannotated TEs in *D. simulans*). However, the TE libraries used here
323 are curated to represent few overlapping TE families. It is still possible that in *D. simulans* some TEs are
324 only partially annotated or missed entirely. If this were the case, then piRNA clusters in *D. simulans* would
325 be denser than in *D. melanogaster*.

326 It is an important question which evolutionary forces drive the evolution of piRNA clusters. In principle,
327 the following forces could act on piRNA clusters. First, different types of mutations, such as insertions
328 due to recent TE activity, the deletion bias observed in *Drosophila* or major rearrangements, for example
329 due to ectopic recombination mediated by TE insertions, may contribute to the rapid turnover of piRNA
330 clusters [Petrov et al., 1996, Langley et al., 1988]. Many TE families are active in *Drosophila* species
331 so recent insertions may be an important driver of cluster evolution [Kofler et al., 2015]. Also genomic
332 rearrangements have been implicated in the evolution of clusters [Assis and Kondrashov, 2009, Gebert et al.,
333 2021]. Second, selection (positive or negative) may contribute to the rapid evolution of piRNA clusters.
334 Theory suggests that an invading TE is silenced by multiple segregating TE insertions distributed over many
335 piRNA clusters [Kofler, 2019, Kelleher et al., 2018]. This hypothesis has been confirmed experimentally by
336 recent works investigating the distribution of cluster insertions in natural and experimental populations
337 that were recently invaded by a TE [Zhang et al., 2020, Kofler et al., 2018]. Theory further suggests that
338 these segregating cluster insertions could be positively selected as haplotypes with a cluster insertion will
339 accumulate few TEs overall and will thus be less deleterious than haplotypes without a cluster insertion
340 [Kofler, 2019, Kelleher et al., 2018, Lu and Clark, 2010]. However, the expected shift in the site frequency
341 spectrum of positively selected cluster insertions is rather subtle and thus difficult to detect experimentally
342 [Kofler, 2019]. In agreement with this, a recent work did not detect evidence that cluster insertions are
343 positively selected [Zhang et al., 2020]. One drawback of this particular study is the lack of reconstruction
344 of the entire piRNA cluster in each strain (P-element insertion sites were identified based on alignments of

345 short reads to a reference genome) [Zhang et al., 2020]. As a consequence, P-element insertions will not
346 be found if adjacent sequences are not conserved and the population frequency of the insertions may be
347 estimated unreliably if the P-element inserted into repetitive regions. However, positive selection of cluster
348 insertions could lead to an accumulation of TE insertions in piRNA clusters. Third, an insertion bias could
349 also lead to an accumulation of TE insertions in piRNA clusters. It is likely that at least some TEs, such
350 as the P-element, have a pronounced insertion bias into piRNA clusters [Ajioka and Eanes, 1989, Zhang
351 et al., 2020, Kofler et al., 2018, Karpen and Spradling, 1992]. It is an important open question whether
352 other TE families also have such an insertion bias into piRNA clusters. Alternatively, piRNA clusters may
353 attract TE insertions, e.g. due to protein-protein interactions [Brennecke et al., 2007, Vermaak and Malik,
354 2009]. Finally, genetic drift could have a strong influence on the evolution of piRNA clusters. Apart from
355 drift of cluster insertions or whole cluster haplotypes, drift may also act on the epigenetically transmitted
356 information that determines the position of piRNA clusters. The information about the position of piRNA
357 clusters is likely not hard coded into the DNA sequence (e.g. by motifs) but rather transmitted epigenetically
358 by the population of maternally deposited piRNAs [Le Thomas et al., 2014a,b]. Stochastic variation in the
359 composition and the amount of maternal transmitted piRNAs could thus lead to a rapid turnover of the
360 location of piRNA clusters. Such a rapid turnover would likely relax selection on individual cluster insertions
361 and make detection of positive selection on cluster insertions even more challenging.

362 This raises the question as to which of these processes are active in the piRNA clusters investigated in
363 the present work. The TE content of piRNA clusters is rapidly evolving and we found that more insertions
364 than deletions were segregating in piRNA clusters of *D. simulans*. The insertions were usually longer and
365 occurring in younger TE families than the deletions. Most insertions are therefore likely due to recent
366 activity of TE families in piRNA clusters. Nevertheless, some insertions (and deletions) could also be due
367 to repeat expansion (and repeat collapse) or genomic rearrangements. A crucial question is whether the
368 observed larger number of insertions in piRNA clusters is due to neutral processes or other forces such as
369 positive selection on cluster insertions and an insertion bias into piRNA clusters. To distinguish between
370 these possibilities, one would need adequate control regions, i.e. a regions that do not produce piRNAs
371 but otherwise have very similar properties to piRNA clusters (pericentromeric regions with a similar size,
372 number, recombination rate and TE content). It is unfortunately challenging to find suitable control regions.
373 Additionally, larger numbers of high quality assemblies for the two *Drosophila* species may be necessary to
374 reliably detect subtle shifts in the site-frequency spectrum of the cluster insertions as expected under positive
375 selection. However, the properties of the deletions in piRNA clusters (short and mostly in older TEs) can
376 likely be explained by the deletion bias observed in *Drosophila*. The gradual erosion of TEs by a deletion bias
377 could also explain why segregating insertions (likely young) are on average longer than fixed insertions (likely
378 old). Another important open question is whether stochastic forces or other processes such as selection and
379 insertion biases are responsible for the differences in the rate of evolution among the piRNA clusters. It is
380 for example possible that positive selection is stronger in clusters producing many piRNAs than in clusters
381 producing few.

382 The available evidence suggests that piRNA clusters are larger in *D. melanogaster* than in *D. simulans*.
383 This could be due to two, not mutually exclusive, reasons: first the clusters are growing in the *D. melanogaster*
384 lineage, or second the clusters are shrinking in the *D. simulans* lineage. Our analysis of insertions and
385 deletions suggests that even in *D. simulans* the clusters are evolving largely by insertions. If piRNA clusters

386 were shrinking in the *D. simulans* lineage, we would not expect to see mostly insertions segregating in *D.*
387 *simulans* populations. Therefore, it seems more likely that the piRNA clusters are expanding in both lineages
388 but in *D. melanogaster* more than in *D. simulans*. This raises the question if the size of piRNA clusters
389 could be subject to a runaway process, where larger clusters will accumulate more insertions of active TEs
390 which, when positively selected, will lead to even larger clusters. This further raises the question whether
391 some forces counteract the expansion of piRNA clusters. Rare and large genomic rearrangements may be an
392 option.

393 We showed that the sequence and the TE content of piRNA clusters is rapidly evolving. This raises an-
394 other important question - Are the positions of piRNA clusters also rapidly changing? Since the information
395 about the position of piRNA clusters is epigenetically transmitted (see above), fluctuations in the popula-
396 tion of maternally transmitted piRNAs could lead to changes in the size and position of piRNA clusters.
397 This likely also happened in our investigated species. For example, the 20 investigated clusters account for
398 21.4% of the uniquely mapped piRNAs in *D. melanogaster* but solely for 8.4% in *D. simulans*. Hence, it is
399 likely that other clusters, not investigated in this work, contribute the bulk of piRNAs in *D. simulans*. In
400 agreement with this, a recent work suggests that many clusters in *Drosophila* are solely found in a single
401 species [Gebert et al., 2021]. The turnover of the location of piRNA clusters within and among species is an
402 important open question for future research.

403 Another important question is whether the observed rapid turnover of piRNA clusters is in conflict with
404 the prevailing view on how TE invasions are stopped: the trap model holds that an invading TE is stopped
405 when a copy of the TE jumps into a piRNA cluster [Bergman et al., 2006, Malone and Hannon, 2009, Zanni
406 et al., 2013, Ozata et al., 2019]. For the trap model to work, it is crucial that the trap (i.e. the piRNA
407 clusters) has a minimum size of about 0.2-3% of the genome [Kofler, 2020]. The number and genomic location
408 of the piRNA clusters has little impact [Kofler, 2019] (except if an organism has a single piRNA cluster in
409 non-recombining regions). As long as piRNA clusters account for at least 0.2-3% of a genome, as is likely
410 that case in *D. melanogaster* and *D. simulans*, we do not think that the rapid turnover of piRNA clusters
411 is in conflict with the trap model.

412 Finally, our work raises the question as to the consequences of rapid evolution of the composition and
413 possibly also location of the loci responsible for silencing TEs. One consequence of such a high turnover
414 is that silencing of TEs may be evolutionary unstable since some individuals in a population may end up
415 without a cluster insertion for a given TE family. A high turnover of piRNA-producing loci could thus explain
416 the low level of activity observed for many TE families in *Drosophila* [Nuzhdin, 1999] since the TE will be
417 active in the individuals that do not produce cognate piRNAs. It is however also possible that silencing of
418 TEs is maintained by a large number of dispersed TE insertions that are not part of piRNA cluster but
419 nevertheless generate piRNAs [Gebert et al., 2021, Mohn et al., 2014, Shpiz et al., 2014]. These piRNA
420 producing TEs are likely due to paramutations whereby an euchromatic TE insertion may be converted into
421 a piRNA producing loci mediated by maternally transmitted piRNAs [Mohn et al., 2014, de Vanssay et al.,
422 2012, Le Thomas et al., 2014b]. In agreement with this, deletion of large piRNA clusters in *D. melanogaster*
423 did not lead to an upregulation of TEs, likely due to a large number of dispersed piRNA-producing TE
424 insertion [Gebert et al., 2021]. If silencing against a TE is effectively based on a large and redundant number
425 of loci, then the rapid turnover of the clusters may not lead to destabilization of the silencing of a TE, which
426 implies that piRNA clusters may largely evolve neutrally.

427 Methods

428 Long-read assemblies and data

429 The two *D. simulans* lines *SZ232* and *SZ45* were collected in California from the Zuma Organic Orchard
430 in Los Angeles, CA on two consecutive weekends of February 2012 [Signor et al., 2017a,b, Signor, 2020].
431 *SZ232* and *SZ45* were sequenced on a MinION platform (Oxford Nanopore Technologies (ONT), Oxford,
432 GB), with fast base-calling using guppy (v4.4.2) and assembled with Canu (v2.1) [Koren et al., 2017] and
433 two rounds of polishing with Racon (v1.4.3) and Pilon (v1.23) [Walker et al., 2014, Vaser et al., 2017, Signor
434 et al., 2017b].

435 The *D. simulans* strain *m252* was collected 1998 in Madagascar and the assembly was generated with
436 PacBio reads [Nouhaud, 2018]. The *D. simulans* strain *w^{xD1}* was originally collected by M. Green, likely
437 in California, but its provenance has been lost. It is a white eyed mutant that has been maintained in
438 the lab for more than 50 years, which can be inferred from the lack of *Wolbachia* infection [Chakraborty
439 et al., 2021]. The *D. melanogaster* strain *A4* was sampled 1963 in Koriba Dam (Zimbabwe) [King et al.,
440 2012]. The reference strain *Iso-1* of *D. melanogaster* was generated by crossing several laboratory strains,
441 with largely unknown sampling data [Brizuela et al., 1994]. *Canton-S* was sampled 1935 in Ohio (USA)
442 [Anxolabéhère et al., 1988]. We could not obtain details on the sampling of the *D. sechellia* strain *sech25*
443 (*Robertson 3C*) and the *D. mauritiana* strain *mau12* (*w12*) [Chakraborty et al., 2021]. The assemblies of the
444 *D. melanogaster* strain *A4* (ASM340174v1), the *D. simulans* strain *w^{xD1}* (ASM438218v1), the *D. sechellia*
445 strain *sech25* (ASM438219v1) and the *D. mauritiana* strain *mau12* (ASM438214v1) are based on PacBio
446 reads [Chakraborty et al., 2018, 2021]. The assembly of the *D. melanogaster* strain *Canton-S* was generated
447 using ONT reads [Wierzbicki et al., 2021]. We obtained the assembly of the *D. melanogaster* reference strain
448 *Iso-1* from FlyBase (r6; [Hoskins et al., 2015]).

449 Identifying homologous piRNA clusters

450 Previously, we designed flanking sequences for 85 out of the 142 annotated piRNA clusters in *D. melanogaster*
451 [Wierzbicki et al., 2021]. We excluded piRNA clusters at the end of chromosomes where two flanking
452 sequences cannot be found, as well as clusters on the fragmented U chromosome. The *D. melanogaster*
453 flanking sequences were aligned to each assembly using bwa bwsw (0.7.17-r1188; [Li and Durbin, 2010]).
454 The alignments were repeated using bwa mem -a (to show alternative hits) to identify clusters that were
455 not recovered by bwa bwsw. Homologous clusters were identified as the regions between the aligned *D.*
456 *melanogaster* flanking sequences [Wierzbicki et al., 2021]. Cluster sequences with internal gaps were excluded.
457 We validated the homology of clusters with a reciprocal mapping approach. First, we designed independent
458 sets of flanking sequences in the target strain (e.g. *D. simulans*) that did not overlap with the aligned *D.*
459 *melanogaster* flanking sequences. Second we aligned these reciprocal flanking sequences with bwa bwsw
460 and bwa mem -a to release 5 of the *D. melanogaster* reference genome (piRNA clusters were annotated in
461 release 5 [Brennecke et al., 2007]). Finally, we checked whether the coordinates of the annotated piRNA
462 clusters were contained within the positions of the aligned reciprocal flanking sequences (supplementary
463 tables S1-S3).

464 **Assembly quality of piRNA clusters**

465 Even when both flanking sequences align to the same contig, a piRNA cluster may be incorrectly assembled,
466 for example if some internal sequences are missing in the assembly. We previously proposed that heterogeneity
467 of the base coverage (e.g. due to repeat collapse) and an elevated soft-clip coverage (resulting from unaligned
468 read termini) can be used to identify assembly errors in clusters [Wierzbicki et al., 2021]. To examine these
469 patterns in our assemblies, we aligned the long reads used for generating the assembly back to the respective
470 assembly using minimap2 (v2.16-r922; v2.17-r954) [Li, 2018]. The exception to this was *D. melanogaster*
471 *Iso-1* where the long reads are not from the original assembly but from a slightly diverged sub-strain Soares
472 et al. [2018]. As reference, we computed the 99% quantiles of the base and soft-clip coverage of complete
473 BUSCO (Benchmarking Universal Single-Copy Orthologs (v3.0.2; v5.0.0); [Simão et al., 2015]) genes based
474 on the *Diptera_odb9* or *Diptera_odb10* data set. Regions where the base or the soft-clip coverage markedly
475 deviates from the 99% quantile of the BUSCO genes could indicate an assembly error and serve as a guide
476 to the quality of the overall cluster assembly.

477 **Aligning the annotations of piRNA clusters**

478 To align the TE annotations of homologous piRNA clusters, we first extracted the sequences of the clusters
479 from the assemblies with samtools (v1.9; [Li et al., 2009]) based on the positions of the aligned flank-
480 ing sequences. Next, we annotated TEs in these sequences using RepeatMasker (open-4.0.7) with a *D.*
481 *melanogaster* TE library and the parameters: -s (sensitive search), -nolow (disable masking of low complex-
482 ity sequences), and -no.is (skip check for bacterial IS) [Smit et al., 2013-2015, Bao et al., 2015, Quesneville
483 et al., 2005]. Finally, we aligned the resulting repeat annotations with our novel tool Manna (see Results)
484 using the parameters -gap 0.09 (gap penalty), -mm 0.1 (mismatch penalty) -match 0.2 (match score).

485 **Visualising piRNA clusters**

486 For visualizing the composition and evolution of piRNA clusters, we annotated repeats in piRNA clusters
487 using the *D. melanogaster* TE library and RepeatMasker (open-4.0.7; [Smit et al., 2013-2015, Bao et al.,
488 2015, Quesneville et al., 2005]). Homologous sequences in piRNA clusters were identified with blastn (BLAST
489 2.7.1+ [Altschul et al., 1990]) using default parameters. We visualized the annotation and the sequence
490 similarity of piRNA clusters with Easyfig (v2.2.3 08.11.2016) [Sullivan et al., 2011] setting the similarity
491 scale to a minimum of 70%. Finally, we merged the pairwise visualizations generated by Easyfig to allow
492 comparing multiple clusters. A walkthrough for this pipeline is available at [https://sourceforge.net/p/
493 manna/wiki/piRNAclusterComparison-walkthrough/](https://sourceforge.net/p/manna/wiki/piRNAclusterComparison-walkthrough/).

494 **piRNAs**

495 We obtained previously published piRNA data from ovaries of *D. simulans* (ERR1821669) and *D. melanogaster*
496 (ERR1821654) strains sampled from Chantemesle (France) [Asif-Laidin et al., 2017]. We trimmed the adap-
497 tor sequence (TGGAATTCTCGGGTGCCAAG) with cutadapt (v3.4; [Martin, 2011]). The reads were
498 aligned to the reference genomes (*D. melanogaster*: *Iso-1*, *D. simulans*: w^{xD1} with novoalign (V3.03.02;
499 <http://novocraft.com/>). The coordinates of the piRNA clusters were obtained from the aligned flanking

500 sequences (see above). We retained reads with a length between 23 and 29bp, normalized the abundance of
501 these reads to a million mapped reads and visualized the abundance of ambiguously ($mq = 0$) and unam-
502 biguously ($mq > 0$) mapped reads along piRNA clusters with R (v3.6.1) and ggplot2 (v3.3.3)[R Core Team,
503 2012, Wickham, 2016].

504 Availability

505 The reads and the assemblies of the two *D. simulans* strains are publicly available (PRJNA736739; PR-
506 JNA736415). The novel software for a multiple alignments of annotations, Manna, is available at [https://](https://sourceforge.net/projects/manna/)
507 sourceforge.net/projects/manna/. A manual and the validations are available at [https://](https://sourceforge.net/p/manna/wiki/Home/)
508 sourceforge.net/p/manna/wiki/Home/. The TE library and list of TE names used in this work are available at [https://](https://sourceforge.net/projects/manna/files/pirnaclustercomparison/resources/)
509 sourceforge.net/projects/manna/files/pirnaclustercomparison/resources/. All script used in
510 this work are available at <https://sourceforge.net/projects/manna/files/publicationdata/>

511 Author contributions

512 FW, RK, and SS conceived this work. SS assembled the two *D. simulans* strains. RK developed Manna.
513 FW, RK and SS analyzed the data. FW, RK and SS wrote the manuscript.

514 Acknowledgments

515 We thank all members of the Institute of Population Genetics for feedback and support. This work was
516 supported by the Austrian Science Foundation (FWF) grant P30036-B25 to RK and by the National Science
517 Foundation Established Program to Stimulate Competitive Research (NSF-EPSCoR-1826834), the North
518 Dakota EPSCoR STEM grants program, and NSF-EPSCoR-2032756 to SS.

519 References

- 520 J. R. Adrion, M. J. Song, D. R. Schrider, M. W. Hahn, and S. Schaack. Genome-Wide Estimates of
521 Transposable Element Insertion and Deletion Rates in *Drosophila Melanogaster*. *Genome Biology and*
522 *Evolution*, 9(5):1329 – 1340, 2017.
- 523 J. W. Ajioka and W. F. Eanes. The accumulation of p-elements on the tip of the x chromosome in populations
524 of *Drosophila melanogaster*. *Genetics Research*, 53(01):1–6, 1989.
- 525 S. F. Altschul, W. Gish, W. Miller, E. W. Myers, and D. J. Lipman. Basic local alignment search tool.
526 *Journal of Molecular Biology*, 215(3):403–410, 1990.
- 527 P. R. Andersen, L. Tirian, M. Vunjak, and J. Brennecke. A heterochromatin-dependent transcription ma-
528 chinery drives piRNA expression. *Nature*, 549(7670):54–59, 2017.
- 529 D. Anxolabéhère, M. G. Kidwell, and G. Periquet. Molecular characteristics of diverse populations are
530 consistent with the hypothesis of a recent invasion of *Drosophila melanogaster* by mobile P elements.
531 *Molecular Biology and Evolution*, 5(3):252–269, 1988.

- 532 A. Asif-Laidin, V. Delmarre, J. Laurentie, W. J. Miller, S. Ronsseray, and L. Teyssset. Short and long-term
533 evolutionary dynamics of subtelomeric piRNA clusters in *Drosophila*. *DNA Research*, 24(5):459–472, 2017.
- 534 R. Assis and A. S. Kondrashov. Rapid repetitive element-mediated expansion of piRNA clusters in mam-
535 malian evolution. *Proceedings of the National Academy of Sciences*, 106(17):7079–7082, 2009.
- 536 W. Bao, K. K. Kojima, and O. Kohany. Repbase Update, a database of repetitive elements in eukaryotic
537 genomes. *Mobile DNA*, 6(1):11, 2015.
- 538 M. G. Barrón, A.-S. Fiston-Lavier, D. A. Petrov, and J. González. Population Genomics of Transposable
539 Elements in *Drosophila*. *Annual Review of Genetics*, 48(1):561–581, 2014.
- 540 C. M. Bergman, H. Quesneville, D. Anxolabéhère, and M. Ashburner. Recurrent insertion and duplication
541 generate networks of transposable element sequences in the *Drosophila melanogaster* genome. *Genome*
542 *Biology*, 7(11):R112, 2006.
- 543 J. P. Blumenstiel. Evolutionary dynamics of transposable elements in a small RNA world. *Trends in Genetics*,
544 27(1):23–31, 2011.
- 545 J. Brennecke, A. A. Aravin, A. Stark, M. Dus, M. Kellis, R. Sachidanandam, and G. J. Hannon. Discrete
546 small RNA-generating loci as master regulators of transposon activity in *Drosophila*. *Cell*, 128(6):1089–
547 1103, 2007.
- 548 J. Brennecke, C. D. Malone, A. A. Aravin, R. Sachidanandam, A. Stark, and G. J. Hannon. An epigenetic
549 role for maternally inherited piRNAs in transposon silencing. *Science*, 322(5906):1387–1392, 2008.
- 550 B. J. Brizuela, L. Elfring, J. Ballard, J. W. Tamkun, and J. A. Kennison. Genetic analysis of the *brahma*
551 gene of *Drosophila melanogaster* and polytene chromosome subdivisions 72AB. *Genetics*, 137(3):803–813,
552 1994.
- 553 E. Casacuberta and J. González. The impact of transposable elements in environmental adaptation. *Molec-*
554 *ular Ecology*, 22(6):1503–1517, 2013.
- 555 M. Chakraborty, N. W. Vankuren, R. Zhao, X. Zhang, S. Kalsow, and J. J. Emerson. Hidden genetic
556 variation shapes the structure of functional elements in *Drosophila*. *Nature Genetics*, 50(1):20–25, 2018.
- 557 M. Chakraborty, C. H. Chang, D. E. Khost, J. Vedanayagam, J. R. Adrion, Y. Liao, K. Montooth, C. D.
558 Meiklejohn, A. M. Larracuente, and J. J. Emerson. Evolution of genome structure in the *Drosophila*
559 *simulans* species complex. *Genome Research*, 31:380–396, 2021.
- 560 B. Charlesworth and D. Charlesworth. The population dynamics of transposable elements. *Genetics Re-*
561 *search*, 42(01):1–27, 1983.
- 562 G. W. Chirn, R. Rahman, Y. A. Sytnikova, J. A. Matts, M. Zeng, D. Gerlach, M. Yu, B. Berger, M. Nara-
563 mura, B. T. Kile, and N. C. Lau. Conserved piRNA Expression from a Distinct Set of piRNA Cluster
564 Loci in Eutherian Mammals. *PLoS Genetics*, 11(11):e1005652, 2015.
- 565 B. Czech, C. D. Malone, R. Zhou, A. Stark, C. Schlingeheyde, M. Dus, N. Perrimon, M. Kellis, J. A.
566 Wohlschlegel, R. Sachidanandam, G. J. Hannon, and J. Brennecke. An endogenous small interfering RNA
567 pathway in *Drosophila*. *Nature*, 453(7196):798–802, 2008.
- 568 B. Czech, M. Munafò, F. Ciabrelli, E. L. Eastwood, M. H. Fabry, E. Kneuss, and G. J. Hannon. piRNA-
569 guided genome defense: From biogenesis to silencing. *Annual Review of Genetics*, 52(1):131–157, 2018.
- 570 P. J. Daborn, J. L. Yen, M. R. Bogwitz, G. Le Goff, E. Feil, S. Jeffers, N. Tijet, T. Perry, D. Heckel,
571 P. Batterham, R. Feyereisen, T. G. Wilson, and R. H. ffrench Constant. A single P450 allele associated
572 with insecticide resistance in *Drosophila*. *Science*, 297(5590):2253–2256, 2002.

- 573 N. Darricarrere, N. Liu, T. Watanabe, and H. Lin. Function of Piwi, a nuclear Piwi/Argonaute protein,
574 is independent of its slicer activity. *Proceedings of the National Academy of Sciences*, 110(4):1297–1302,
575 2013.
- 576 A. de Vanssay, A.-L. Bougé, A. Boivin, C. Hermant, L. Teyssset, V. Delmarre, C. Antoniewski, and S. Ron-
577 sseray. Paramutation in *Drosophila* linked to emergence of a piRNA-producing locus. *Nature*, 490(7418):
578 112–115, 2012.
- 579 P. Dimitri, N. Junakovic, and B. Arcà. Colonization of Heterochromatic Genes by Transposable Elements
580 in *Drosophila*. *Molecular Biology and Evolution*, 20(4):503 – 512, 2003.
- 581 D.-F. Feng and R. F. Doolittle. Progressive Sequence Alignment as a Prerequisite to Correct Phylogenetic
582 Trees. *Journal of Molecular Evolution*, 25(4):351–360, 1987.
- 583 D. Gebert, L. K. Neubert, C. Lloyd, J. Gui, R. Lehmann, F. K. Teixeira, D. Gebert, L. K. Neubert, C. Lloyd,
584 J. Gui, R. Lehmann, and F. K. Teixeira. Large *Drosophila* germline piRNA clusters are evolutionarily
585 labile and dispensable for transposon regulation. *Molecular Cell*, 81:1–14, 2021.
- 586 J. González, K. Lenkov, M. Lipatov, J. M. Macpherson, and D. A. Petrov. High rate of recent transposable
587 element-induced adaptation in *Drosophila melanogaster*. *PLoS Biology*, 6(10):e251, 2008.
- 588 C. Goriaux, E. Théron, E. Brassset, and C. Vaury. History of the discovery of a master locus producing
589 piRNAs: The flamenco/COM locus in *Drosophila melanogaster*. *Frontiers in Genetics*, 5:257, 2014.
- 590 L. S. Gunawardane, K. Saito, K. M. Nishida, K. Miyoshi, Y. Kawamura, T. Nagami, H. Siomi, and M. C.
591 Siomi. A slicer-mediated mechanism for repeat-associated siRNA 5' end formation in *Drosophila*. *Science*,
592 315(5818):1587–1590, 2007.
- 593 R. A. Hoskins, J. W. Carlson, K. H. Wan, S. Park, I. Mendez, S. E. Galle, B. W. Booth, B. D. Pfeiffer,
594 R. A. George, R. Svirskas, et al. The Release 6 reference sequence of the *Drosophila melanogaster* genome.
595 *Genome Research*, 25(3):445–458, 2015.
- 596 T. Josse, L. Teyssset, A.-L. Todeschini, C. M. Sidor, D. Anxolabéhère, and S. Ronsseray. Telomeric trans-
597 silencing: an epigenetic repression combining RNA silencing and heterochromatin formation. *PLoS Ge-
598 netics*, 3(9):1633–1643, 2007.
- 599 A. I. Kalmykova, M. S. Klenov, and V. A. Gvozdev. Argonaute protein PIWI controls mobilization of
600 retrotransposons in the *Drosophila* male germline. *Nucleic Acids Research*, 33(6):2052–2059, 2005.
- 601 G. H. Karpen and A. C. Spradling. Analysis of subtelomeric heterochromatin in the *Drosophila* minichro-
602 mosome Dp1187 by single P element insertional mutagenesis. *Genetics*, 132(3):737–753, 1992.
- 603 E. S. Kelleher, R. B. R. Azevedo, and Y. Zheng. The Evolution of Small-RNA-Mediated Silencing of an
604 Invading Transposable Element. *Genome Biology and Evolution*, 10(11):3038–3057, 2018.
- 605 E. G. King, C. M. Merkes, C. L. McNeil, S. R. Hoofer, S. Sen, K. W. Broman, A. D. Long, and S. J. Mac-
606 donald. Genetic dissection of a model complex trait using the *Drosophila* Synthetic Population Resource.
607 *Genome Research*, 22(8):1558–1566, 2012.
- 608 R. Kofler. Dynamics of Transposable Element Invasions with piRNA Clusters. *Molecular Biology and
609 Evolution*, 36(7):1457–1472, 2019.
- 610 R. Kofler. piRNA Clusters Need a Minimum Size to Control Transposable Element Invasions. *Genome
611 Biology and Evolution*, 12(5):736–749, 2020.
- 612 R. Kofler, V. Nolte, and C. Schlotterer. Tempo and mode of transposable element activity in *Drosophila*.
613 *PLoS Genetics*, 11(7):e1005406, 2015.

- 614 R. Kofler, K.-A. Senti, V. Nolte, R. Tobler, and C. Schlötterer. Molecular dissection of a natural transposable
615 element invasion. *Genome Research*, 28(6):824–835, 2018.
- 616 S. Koren, B. P. Walenz, K. Berlin, J. R. Miller, N. H. Bergman, and A. M. Phillippy. Canu: Scalable and
617 accurate long-read assembly via adaptive κ -mer weighting and repeat separation. *Genome Research*, 27
618 (5):722–736, 2017.
- 619 C. H. Langley, E. Montgomery, R. Hudson, N. Kaplan, and B. Charlesworth. On the role of unequal exchange
620 in the containment of transposable element copy number. *Genetics research*, 52(03):223–235, 1988.
- 621 A. Le Thomas, A. K. Rogers, A. Webster, G. K. Marinov, S. E. Liao, E. M. Perkins, J. K. Hur, A. A. Aravin,
622 and K. F. Tóth. Piwi induces piRNA-guided transcriptional silencing and establishment of a repressive
623 chromatin state. *Genes and Development*, 27(4):390–399, 2013.
- 624 A. Le Thomas, G. K. Marinov, and A. A. Aravin. A transgenerational process defines piRNA biogenesis in
625 *Drosophila virilis*. *Cell Reports*, 8(6):1617–1623, 2014a.
- 626 A. Le Thomas, E. Stuwe, S. Li, J. Du, G. Marinov, N. Rozhkov, Y. C. A. Chen, Y. Luo, R. Sachidanandam,
627 K. F. Toth, D. Patel, and A. A. Aravin. Transgenerationally inherited piRNAs trigger piRNA biogenesis
628 by changing the chromatin of piRNA clusters and inducing precursor processing. *Genes and Development*,
629 28(15):1667–1680, 2014b.
- 630 Y. C. G. Lee and C. H. Langley. Transposable elements in natural populations of *Drosophila melanogaster*.
631 *Philosophical transactions of the Royal Society of London. Series B, Biological Sciences*, 365(1544):1219–
632 1228, 2010.
- 633 Y. C. G. Lee and C. H. Langley. Long-term and short-term evolutionary impacts of transposable elements
634 on *Drosophila*. *Genetics*, 192(4):1411–1432, 2012.
- 635 E. Lerat, N. Burlet, C. Biéumont, and C. Vieira. Comparative analysis of transposable elements in the
636 melanogaster subgroup sequenced genomes. *Gene*, 473(2):100–109, 2011.
- 637 R. Levis, K. O’Hare, and G. M. Rubin. Effects of transposable element insertions on RNA encoded by the
638 white gene of *Drosophila*. *Cell*, 38(2):471–481, 1984.
- 639 H. Li. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18):3094–3100, 2018.
- 640 H. Li and R. Durbin. Fast and accurate long-read alignment with Burrows-Wheeler transform. *Bioinform-*
641 *atics*, 26(5):589–595, 2010.
- 642 H. Li, B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, N. Homer, G. Marth, G. Abecasis, R. Durbin, and
643 1000 Genome Project Data Processing Subgroup. The Sequence Alignment/Map format and SAMtools.
644 *Bioinformatics*, 25(16):2078–2079, Aug. 2009.
- 645 J. K. Lim. Intrachromosomal rearrangements mediated by hobo transposons in *Drosophila melanogaster*.
646 *Proceedings of the National Academy of Sciences*, 85(23):9153–9157, 1988.
- 647 J. Lu and A. G. Clark. Population dynamics of PIWI-interacting RNAs (piRNAs) and their targets in
648 *Drosophila*. *Genome Research*, 20(2):212–227, 2010.
- 649 C. D. Malone and G. J. Hannon. Small RNAs as Guardians of the Genome. *Cell*, 136(4):656–668, 2009.
- 650 C. D. Malone and G. J. Hannon. Molecular Evolution of piRNA and Transposon Control Pathways in
651 *Drosophila*. *Cold Spring Harbor Symposia on Quantitative Biology*, 74:225 – 234, 2010.
- 652 L. Marin, M. Lehmann, D. Nouaud, H. Izaabel, D. Anxolabéhère, and S. Ronsseray. P-element repression in
653 *Drosophila melanogaster* by a naturally occurring defective telomeric P copy. *Genetics*, 155(4):1841–1854,
654 2000.

- 655 M. Martin. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet. journal*,
656 17(1):10–12, 2011.
- 657 L. Mateo, A. Ullastres, and J. González. A Transposable Element Insertion Confers Xenobiotic Resistance
658 in *Drosophila*. *PLoS Genetics*, 10(8):e1004560, 2014.
- 659 B. McClintock. Controlling elements and the gene. *Cold Spring Harbor Symposia on Quantitative Biology*,
660 21:197–216, 1956.
- 661 W. McGinnis, A. W. Shermoen, and S. K. Beckendorf. A transposable element inserted just 5' to a *Drosophila*
662 glue protein gene alters gene expression and chromatin structure. *Cell*, 34(1):75–84, 1983.
- 663 F. Mohn, G. Sienski, D. Handler, and J. Brennecke. The rhino-deadlock-cutoff complex licenses noncanonical
664 transcription of dual-strand piRNA clusters in *Drosophila*. *Cell*, 157(6):1364–1379, 2014.
- 665 S. B. Needleman and C. D. Wunsch. A general method applicable to the search for similarities in the amino
666 acid sequence of two proteins. *Journal of Molecular Biology*, 48(3):443–453, 1970.
- 667 P. Nouhaud. Long-read based assembly and annotation of a *Drosophila simulans* genome. *bioRxiv*, 2018.
668 doi: 10.1101/425710.
- 669 S. V. Nuzhdin. Sure facts, speculations, and open questions about the evolution of transposable element
670 copy number. *Genetica*, 107(1-3):129–137, 1999.
- 671 D. J. Obbard, J. MacLennan, K.-W. Kim, A. Rambaut, P. M. O’Grady, and F. M. Jiggins. Estimating
672 divergence dates and substitution rates in the *Drosophila* phylogeny. *Molecular Biology and Evolution*, 29
673 (11):3459–3473, 2012.
- 674 I. Olovnikov, S. Ryazansky, S. Shpiz, S. Lavrov, Y. Abramov, C. Vaury, S. Jensen, and A. Kalmykova. De
675 novo piRNA cluster formation in the *Drosophila* germ line triggered by transgenes containing a transcribed
676 transposon fragment. *Nucleic Acids Research*, 41(11):5757–5768, 2013.
- 677 D. M. Ozata, I. Gainetdinov, A. Zoch, D. O’Carroll, and P. D. Zamore. PIWI-interacting RNAs: small
678 RNAs with big functions. *Nature Reviews Genetics*, 20(2):89–108, 2019.
- 679 L. Peters and G. Meister. Argonaute proteins: Mediators of rna silencing. *Molecular Cell*, 26(5):611–623,
680 2007.
- 681 D. A. Petrov, E. R. Lozovskaya, and D. L. Hartl. High intrinsic rate of DNA loss in *Drosophila*. *Nature*, 384
682 (6607):346–349, 1996.
- 683 B. Piegue, R. Guyot, N. Picault, A. Roulin, A. Saniyal, H. Kim, K. Collura, D. S. Brar, S. Jackson, R. A.
684 Wing, and O. Panaud. Doubling genome size without polyploidization: Dynamics of retrotransposition-
685 driven genomic expansions in *Oryza australiensis*, a wild relative of rice. *Genome Research*, 16(10):1262 –
686 1269, 2006.
- 687 H. Quesneville, C. M. Bergman, O. Andrieu, D. Autard, D. Nouaud, M. Ashburner, and D. Anxolabéhère.
688 Combined evidence annotation of transposable elements in genome sequences. *PLoS Computational Biol-*
689 *ogy*, 1(2):166–175, 2005.
- 690 R Core Team. *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical
691 Computing, Vienna, Austria, 2012. URL <http://www.R-project.org>. ISBN 3-900051-07-0.
- 692 P. S. Schnable, S. Pasternak, C. Liang, J. Zhang, L. Fulton, T. A. Graves, P. Minx, A. D. Reily, L. Courtney,
693 S. S. Kruchowski, C. Tomlinson, C. Strong, K. Delehaunty, C. Fronick, B. Courtney, S. M. Rock, E. Belter,
694 F. Du, K. Kim, R. M. Abbott, M. Cotton, A. Levy, P. Marchetto, K. Ochoa, S. M. Jackson, B. Gillam,

- 695 W. Chen, L. Yan, J. Higginbotham, M. Cardenas, J. Waligorski, E. Applebaum, L. Phelps, J. Falcone,
696 K. Kanchi, T. Thane, A. Scimone, N. Thane, J. Henke, T. Wang, J. Ruppert, N. Shah, K. Rotter,
697 J. Hodges, E. Ingenthron, M. Cordes, S. Kohlberg, J. Sgro, B. Delgado, K. Mead, A. Chinwalla, S. Leonard,
698 K. Crouse, K. Collura, D. Kudrna, J. Currie, R. He, A. Angelova, S. Rajasekar, T. Mueller, R. Lomeli,
699 G. Scara, A. Ko, K. Delaney, M. Wissotski, G. Lopez, D. Campos, M. Braidotti, E. Ashley, W. Golser,
700 H. Kim, S. Lee, J. Lin, Z. Dujmic, W. Kim, J. Talag, A. Zuccolo, C. Fan, A. Sebastian, M. Kramer,
701 L. Spiegel, L. Nascimento, T. Zutavern, B. Miller, C. Ambroise, S. Muller, W. Spooner, A. Narechania,
702 L. Ren, S. Wei, and S. Kumari. The B73 Maize Genome: Complexity, Diversity, and Dynamics. *Science*,
703 326(5956):1112–1115, 2009.
- 704 F. Schwarz, F. Wierzbicki, K.-A. Senti, and R. Kofler. Tirant Stealthily Invaded Natural *Drosophila*
705 *melanogaster* Populations during the Last Century. *Molecular Biology and Evolution*, 38(4):1482–1497,
706 2021.
- 707 S. Shpiz, S. Ryazansky, I. Olovnikov, Y. Abramov, and A. Kalmykova. Euchromatic transposon insertions
708 trigger production of novel pi-and endo-siRNAs at the target sites in the *Drosophila* germline. *PLoS*
709 *Genetics*, 10(2):e1004138, 2014.
- 710 G. Sienski, D. Dönertas, and J. Brennecke. Transcriptional silencing of transposons by Piwi and maelstrom
711 and its impact on chromatin state and gene expression. *Cell*, 151(5):964–980, 2012.
- 712 S. Signor. Transposable elements in individual genotypes of *Drosophila simulans*. *Ecology and Evolution*, 10
713 (7):3402–3412, 2020.
- 714 S. A. Signor, M. Abbasi, P. Marjoram, and S. V. Nuzhdin. Conservation of social effects (Ψ) between two
715 species of *Drosophila* despite reversal of sexual dimorphism. *Ecology and Evolution*, 7(23):10031 – 10041,
716 2017a.
- 717 S. A. Signor, F. N. New, and S. Nuzhdin. A Large Panel of *Drosophila simulans* Reveals an Abundance of
718 Common Variants. *Genome Biology and Evolution*, 10(1):189 – 206, 12 2017b.
- 719 F. A. Simão, R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov. BUSCO: Assessing
720 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*, 31(19):3210–
721 3212, 2015.
- 722 A. F. A. Smit, R. Hubley, and P. Green. RepeatMasker Open-4.0, 2013-2015. URL [http://www.](http://www.repeatmasker.org)
723 [repeatmasker.org](http://www.repeatmasker.org).
- 724 E. A. Solares, M. Chakraborty, D. E. Miller, S. Kalsow, K. Hall, A. G. Perera, J. J. Emerson, and R. S.
725 Hawley. Rapid low-cost assembly of the *Drosophila melanogaster* reference genome using low-coverage,
726 long-read sequencing. *G3: Genes, Genomes, Genetics*, 8(10):3143–3154, 2018.
- 727 M. J. Sullivan, N. K. Petty, and S. A. Beatson. Easyfig: a genome comparison visualizer. *Bioinformatics*,
728 27(7):1009–1010, 2011.
- 729 R. Vaser, I. Sovic, N. Nagarajan, and M. Sikic. Fast and accurate de novo genome assembly from long
730 uncorrected reads. *Genome research*, 27(5):737–746, 2017.
- 731 D. Vermaak and H. S. Malik. Multiple roles for heterochromatin protein 1 genes in *Drosophila*. *Annual*
732 *Review of Genetics*, 43:467–492, 2009.
- 733 C. Vieira, M. Fablet, E. Lerat, M. Boulesteix, R. Rebollo, N. Bulet, A. Akkouche, B. Hubert, H. Mortada,
734 and C. Biéumont. A comparative analysis of the amounts and dynamics of transposable elements in natural
735 populations of *Drosophila melanogaster* and *Drosophila simulans*. *Journal of Environmental Radioactivity*,
736 113:83–86, 2012.

- 737 B. J. Walker, T. Abeel, T. Shea, M. Priest, A. Abouelliel, S. Sakthikumar, C. A. Cuomo, Q. Zeng, J. Wort-
738 man, S. K. Young, and A. M. Earl. Pilon: An integrated tool for comprehensive microbial variant detection
739 and genome assembly improvement. *PLoS ONE*, 9(11), 2014.
- 740 H. Wickham. *ggplot2: Elegant Graphics for Data Analysis*. Springer Nature, Basel, Switzerland, 2016. ISBN
741 978-3-319-24277-4.
- 742 F. Wierzbicki, F. Schwarz, O. Cannalunga, and R. Kofler. Novel quality metrics allow identifying and
743 generating high-quality assemblies of piRNA clusters. *Molecular Ecology Resources*, 2021. doi: 10.1111/
744 1755-0998.13455.
- 745 H.-P. Yang and S. V. Nuzhdin. Fitness costs of *Doc* expression are insufficient to stabilize its copy number
746 in *Drosophila melanogaster*. *Molecular Biology and Evolution*, 20(5):800–804, 2003.
- 747 V. Zanni, A. Eymery, M. Coiffet, M. Zytnicki, I. Luyten, H. Quesneville, C. Vaury, and S. Jensen. Distribu-
748 tion, evolution, and diversity of retrotransposons at the *flamenco* locus reflect the regulatory properties
749 of piRNA clusters. *Proceedings of the National Academy of Sciences*, 110(49):19842–19847, 2013.
- 750 S. Zhang, B. Pointer, and E. S. Kelleher. Rapid evolution of piRNA-mediated silencing of an invading
751 transposable element was driven by abundant de novo mutations. *Genome Research*, 30(4):566–575, 2020.