

# 1 **Alternative antiviral immune pathways are rapidly evolving in *Drosophila innubila***

2

## 3 ***Supplementary Methods & Results***

### 4 *DNA/RNA isolation, library preparation and sequencing*

5 We extracted DNA following the protocol described in Chakraborty and Emerson (Chakraborty *et*  
6 *al.* 2017). Briefly, approximately 320 adult females from an isofemale line of *D. innubila* (captured  
7 in the Chiricahua Mountains in 2005 by Kelly Dyer, strain name SWRS2005-50) were starved for  
8 five hours then frozen and ground to a powder in liquid nitrogen then extracted using a modified  
9 version of the Qiagen Blood and Cell Culture DNA Midi Kit (#13343, USA Qiagen Inc.,  
10 Germantown, MD, USA). The extraction yielded fragment sizes greater than 60,000 bp as  
11 determined by Agilent TapeStation (Agilent, Santa Clara, CA, USA). We prepared a sequencing  
12 library using the Oxford Nanopore Technologies Rapid 48-hour (SQK-RAD002) protocol which  
13 was sequencing using a MinION (Supplementary Table 1, NCBI SRA: SAMN11037163, Oxford  
14 Nanopore Technologies, Oxford, UK). The same DNA was also used to construct a Nextera  
15 fragment library with insert sizes of ~180bp, ~3000bp and ~7000bp. We sequenced the libraries  
16 on a MiSeq (300bp paired-end, Illumina, San Diego, CA, USA) which generated 20104299 200bp  
17 paired-end reads (NCBI SRA: SAMN11037164). All data used in the assembly and annotation of  
18 the *D. innubila* genome are available in the NCBI BioProject PRJNA524688.

19 For *D. innubila* long reads, DNA was sequenced on the Oxford Nanopore Technologies  
20 Minion platform using the SQK-RAD002 protocol and a 48-hour run (Jain *et al.* 2016). Bases were  
21 called *post hoc* using the built in read\_fast5\_basecaller.exe program with options: -f FLO-MIN106  
22 -k SQK-RAD002 -r-t 4. The MinION produced 746229 reads, an average of 5754bp long, with  
23 656860 reads greater than 1kbp, 225704 reads greater than 10kbp and a maximum read length of  
24 1.61Mbp (NCBI SRA: SAMN11037163).

25 DNA for the Hi-C protocol was extracted from fifteen adult females by PhaseGenomics  
26 with a 4-cutter Sau3AI being used to digest the chromatin. This library was then sequenced on an  
27 Illumina NextSeq (Illumina, San Diego, CA, USA).

28 For the *Drosophila falleni* (strain 15130-1961.00 from the Cornell *Drosophila* species stock  
29 center), we followed the same protocol as *D. innubila* for DNA isolation and library preparation,  
30 but only constructed a single 300bp insert library. This was sequenced on one half of a MiSeq  
31 (300bp paired-end, Illumina, San Diego, CA, USA) by the KU CMADP genomics core. For

32 *Drosophila phalerata* (obtained from Kelly Dyer), we followed a standard Puregene Gentra  
33 extraction (USA Qiagen Inc., Germantown, MD, USA) and constructed a 300bp insert Nextera  
34 library (see above) which was sequenced on a fraction of an Illumina HiSeq 4000 run (150bp  
35 paired end). This generated 8080281 300bp paired-end reads for *D. falleni* and 24896114 150bp  
36 paired-end reads for *D. phalerata*. We estimated the heterozygosity of each sample using  
37 Jellyfish (Marcais 2011) and GenomeScope (Vurture *et al.* 2017) and found the heterozygosity  
38 of each sample to be between 0.46% and 0.81%.

39 For gene expression analyses, we obtained two replicate samples of female and male heads  
40 and whole bodies (including heads), embryos, larvae (pooled across all three instar stages) and  
41 pupae (all non-adults were unsexed). RNA was extracted using a standard Trizol procedure  
42 (Simms *et al.* 1993) with a DNase step. RNA-sequencing libraries were constructed using the  
43 standard TruSeq protocol (McCoy *et al.* 2014) with ½ volume reactions to conserve reagents.  
44 Individually indexed RNA libraries (2 replicates from each tissue/sex) were sequenced on one lane  
45 of an Illumina “Rapid” run with 100bp single-end reads, as outlined in Supplementary Table 1.

46

#### 47 *Whole genome assembly*

48 Raw reads from the Oxford Nanopore Minion were assembled using CANU version 1.6 (Koren *et*  
49 *al.* 2016) with an estimated genome size of 150 million bases and the “nanopore-raw” flag. We  
50 then used Pilon (Walker *et al.* 2014) to polish the genome with our Illumina fragment library  
51 (default parameters). The resulting assembly was submitted to PhaseGenomics  
52 (phasegenomics.com, Seattle, WA USA) for scaffolding using Hi-C and further polished with  
53 Pilon for seven iterations. With each iteration, we evaluated the quality of the genome and the  
54 extent of improvement in quality, by calculating the N50 and using BUSCO (Simão *et al.* 2015)  
55 to identify the presence of conserved genes in the genome, from a database of 2799 single copy  
56 Dipteran genes. The final genome and annotation are available at NCBI (accession:  
57 SKCT000000000).

58

#### 59 *Genome annotation*

60 We assembled a *de novo* transcriptome using Trinity (version 2.4.0) (Haas *et al.* 2013). First, we  
61 quality filtered single-end reads (samples described above) with Scythe (Buffalo 2018;  
62 <http://github.com/vsbuffalo/scythe>) and Sickle (Joshi and Fass 2011;

63 <http://github.com/najoshi/sickle>) to remove rRNA, Illumina adapters and low quality sequences  
64 (quality less than 20). We concatenated all reads from all tissues and used the Trinity package  
65 with default parameters to assemble the transcriptome. We also assembled transcriptome with  
66 Oases (Schulz *et al.* 2012) (velvetg parameters: -exp\_cov 100 -max\_coverage 500 -  
67 min\_contig\_lgth 50 -read\_trkg yes) and SOAP*de novo* Trans (Xie *et al.* 2014) (127mer with  
68 parameters: SOAPdenovo-Trans-127mer -p 28 -e 4 and the following kmers: 95, 85, 75, 65, 55,  
69 45, 35, 29, 25, 21). These assemblies were used to make a metatranscriptome using  
70 EvidentialGene (Gilbert 2013; <http://eugenes.org/EvidentialGene/>) (parameters: -NCPUs=28 -  
71 MAXMEM=489000).

72       Using the *D. innubila* transcriptome as well as protein databases from *M. domestica*, *D.*  
73 *melanogaster*, and *D. virilis*, we searched for evidence of genic regions in the genome assembly.  
74 A database containing repeat sequences discovered by RepeatModeler (Smit and Hubley 2008)  
75 was also utilized by MAKER2 (Holt and Yandell 2011) to ensure that repetitive regions are not  
76 annotated as genes. Post completion of the first MAKER2 run, we extracted all gene models from  
77 the annotation that had a predicted protein length of at least 50 amino acids [-l 50] and an AED  
78 score (Eilbeck *et al.* 2009) of no more than 0.25 [-x 0.25] to form a training set for SNAP (Korf  
79 2004). The resulting HMM file was used as an input to round 2 of the annotation pipeline, along  
80 with GFF files containing all transcript, protein, and repeat evidence collected during round 1.  
81 Additionally, we provided MAKER2 with training files for *D. melanogaster* for Augustus  
82 (publicly available and distributed with MAKER2) (Stanke *et al.* 2008). After the second round of  
83 annotations, we repeated the SNAP training steps taken after round 1, which produced a new HMM  
84 file. The HMM file from round 2, the GFF files with the evidence from round 1  
85 (transcripts/proteins/repeats), and the *D. melanogaster* training set for Augustus were the inputs  
86 for round 3 of the annotation pipeline.

87       Using resources on FlyBase (Gramates *et al.* 2017; FlyBase.org) we identified conservation  
88 of each gene by counting the number of the 12 *Drosophila* species genome orthologs (and humans,  
89 if applicable). We also calculated the percentage of genic nucleotides per Megabase across the  
90 genome in 250kbp sliding windows.

91       Our annotation resulted in the identification of 12318 genes of varying lengths  
92 (Supplementary Table 2 & 3). We find an absence of many tRNAs usually found in *Drosophila*  
93 genomes, this may be an error of genome assembly or annotation, but additional tRNAs were

94 unable to be found via BLAST (Altschul *et al.* 1990), either due to their absence in the genome  
95 or divergence of tRNAs due to *D. innubila*'s extensive divergence from previously sequenced  
96 species (Supplementary Table 3, Figure 1). Most of the genes found in the genome (11925) are  
97 shared with other species (among the 12 genomes available on Flybase as of July 2018;  
98 <ftp://ftp.flybase.net/releases/current>), with these genes containing 97.9% of the Dipteran BUSCO  
99 library (Simão *et al.* 2015), and 7,094 of these genes have orthologs in the human genome (based  
100 on the current version available in FlyBase as of July 2018 ;  
101 <ftp://ftp.flybase.net/releases/current>).

102

### 103 *Further genome assembly*

104 To identify additional genes missed in the Hi-C assembly, we also took all unmapped reads  
105 and assembled these using SPAdes (Bankevich *et al.* 2012). We mapped MiSeq information to the  
106 15587 SPAdes assembled contigs and kept contigs with similar coverage to the CANU assembled  
107 scaffolds (25-35 fold coverage) and with Blastn (Camacho *et al.* 2009) hits to known Dipteran  
108 sequences (e-value < 0.001), retaining an additional 302 contigs (336 total).

109 Finally, we used Mauve (Darling *et al.* 2004) to identify regions of orthology between the  
110 *Drosophila virilis* genome (Clark *et al.* 2007) and the *D. innubila* genome. We calculated the GC  
111 content and percent of windows with identifiable orthology to *virilis*, in 250kbp windows across  
112 the *D. innubila* genome using bedTools (Quinlan and Hall 2010).

113 To assemble the mitochondrial genome, we took a subset (100000 read pairs) of the  
114 short-read data generated by MiSeq (sequencing and data preparation described in the methods).  
115 We assembled this subset of reads with Geneious (default parameters) (Kearse *et al.* 2012) and  
116 used Blastn to find contigs with hits to mitochondrial genes (non-redundant database, e-value <  
117 0.001). In our initial assembly we found a single, complete, assembled, circular contig ~16kb  
118 long with high confidence hits to all mitochondrial genes. In all following steps, we used this  
119 sequence as the fully assembled mitochondria. The mitochondrial genome was also included  
120 during polishing with Pilon for the seven iterations. We then used the MITOS online portal  
121 (Bernt *et al.* 2013) to annotate the 16191bp assembled and polished mitochondrial genome.

122 We attempted to assemble parts of the Y chromosome using sequencing information  
123 available for male *D. innubila* (SRA: SAMN07638923/SRR6033015) (Hill and Unckless 2017).  
124 We mapped these sequences to the female reference genome using BWA MEM with default

125 parameters (Li and Durbin 2009), extracted all unmapped reads using SamTools (Li *et al.* 2009)  
126 and attempted to assemble these using Spades (default parameters) (Bankevich *et al.* 2012). We  
127 then mapped male and female expression data using GSNAP (Wu and Nacu 2010) to this dataset  
128 along with the whole genome. We considered the assembled contigs containing genes with  
129 significantly greater expression in males (using EdgeR (Robinson *et al.* 2009)), using the  
130 methods for RNA differential expression described below,  $p$ -value  $< 0.001$ , FDR  $< 0.001$ ) to be  
131 putatively Y-linked. This filtered left us with 27 putatively male biased, Y-linked (or  
132 heterochromatic) scaffolds. We used blastn and tblastx to attempt to identify any known  
133 orthologs to these genes.

134 We identified large structural variants among the genomes of *D. innubila*, *D. falleni* and  
135 *D. phalerata* using both Manta (Chen *et al.* 2016) and Pindel (Ye *et al.* 2009) (default, bam input  
136 in both cases) on *D. falleni* and *D. phalerata* short read data mapped to the *D. innubila* genome.  
137 We extracted the structural variants found with both software packages as VCF files and  
138 considered only the variants detected by both Manta and Pindel to be real. We compared  $dN/dS$   
139 between genes found within inversions and outside and found no significant differences in either  
140  $dN/dS$  or  $dS$  (Wilcoxon rank sum  $W = 35$ ,  $p$ -value  $> 0.05$ ).

141 For all genes we performed a codon bias analysis using CodonW (Peden 1997). We  
142 compared the codon bias index (CBI), codon adaptation index (CAI) and the frequency of  
143 optimal codons (Fop) across scaffolds, between novel genes and previously known genes, and  
144 between highly expressed genes (counts per million reads [CPM]  $> 1$  in at least one dataset) and  
145 under expressed genes (CPM  $< 1$  across all datasets). We find a significant conservation of  
146 codons in the mitochondria, Muller element F and the heterochromatic contigs, versus all other  
147 contigs (Supplementary Figure 4, Wilcox test  $W > 1132$ ,  $p$ -value  $< 0.01186$  for all CBI, CAI and  
148 Fop) (Zhou and Bachtrog 2015). For Muller elements A, B and E, we find significant levels of  
149 codon adaptation, optimal use and codon bias (Supplementary Figure 4, Wilcox test  $W >$   
150  $14217000$   $p$ -value  $< 0.001586$  for all CBI, CAI and Fop). We find significant positive  
151 associations between gene expression, and codon adaptation and optimization (GLM t-value  $>$   
152  $4.746$ ,  $p$ -value  $< 2.1e-06$ ), consistent with an expectation for selection for codon efficiency in  
153 more highly expressed genes.

154 We find 393 orphan genes in the genome and compared the median expression, gene  
155 length, number of introns, codon bias and GC content between previously identified genes and

156 the remaining putatively novel genes using codonW (Peden 1997) and bedTools (Quinlan and  
157 Hall 2010). Orphan genes are significantly shorter, under-expressed, AT-rich and intron-poor  
158 when compared to genes with previously identified orthologs (Supplementary Figure 5,  
159 Wilcoxon rank sum  $p$ -value  $< 0.0113$ ), consistent with their more recent origin (Palmieri *et al.*  
160 2014). We find a significant excess of orphan genes on two unassembled (scaffolds 5 and 11),  
161 these scaffolds are likely heterochromatic and sparse coding regions ( $\chi^2 > 16.7$ ,  $p$ -value  $<$   
162  $0.0005$ ), We also find a significant deficit of orphan genes on Muller elements C and E ( $\chi^2 >$   
163  $14.17$ ,  $p < 0.000836$ ). Of these orphans, 51 show differential expression across life stages,  
164 primarily in the embryos, suggesting possible functionalization in different stages  
165 (Supplementary Tables 7-10, EdgeR analysis  $p$ -value  $< 0.05$ , FDR  $< 0.05$  after multiple testing  
166 correction).

167

#### 168 *Transposable element (TE) family comparison between species of the quinaria group*

169 We identified repetitive sequences *de novo* using RepeatModeler (engine = NCBI) (Smit  
170 and Hubley 2008). We then used RepeatMasker to mask the repetitive regions and classify repeats  
171 in classes/orders/families (-gff -gcalc -s) (Smit and Hubley 2015). We then used Blastn  
172 (parameters:  $e$ -value  $< 0.001$ ) to compare each consensus sequence identified to the Repbase TE  
173 database (Bao *et al.* 2015), to confirm the TE order of each sequence. Using the GFF of repeat  
174 sequences generated by RepeatMasker, we then calculated the insertion density per 250kbp of the  
175 genome sliding across the genome for TE insertions. Using genomeCoverageBed (Quinlan and  
176 Hall 2010), we found the median coverage of the autosomes and each TE family and estimated the  
177 copy number of each TE family in the genome.

178 We estimate 13.53% of the genome consists of transposable elements (TEs). We find 175  
179 TE families, consisting of 79 terminal inverted repeat DNA transposon families (TIR, 5.01%), 34  
180 rolling circle/helitron DNA transposon families (RC, 5.61%), 25 long terminal repeat  
181 retrotransposons (LTR, 1.04%), and 26 long interspersed nuclear element retroposons (LINE,  
182 1.87%) (Table 1). In addition to transposable elements, we find 10 short interspersed nuclear  
183 elements (SINE) and satellite element families, which together with simple repeats make up  
184 3.42% of the genome (Figure 1A, Supplementary Figure 6). On Muller element A and B, we find  
185 two large regions consisting primarily of transposable elements. We considered these to be  
186 heterochromatic regions and potentially piRNA clusters. A majority (over 50% of the sequence)

187 of these clusters consists of single TE superfamilies. Helitrons are primarily found throughout  
188 Muller element A, while Muller B's heterochromatic region consists of R2 LINE retroposons  
189 (Figure 1).

190 For *D. innubila*, *D. falleni* and *D. phalerata*, we mapped the short read information to the  
191 masked species reference genome with concatenated consensus TE sequences using BWA MEM  
192 (parameters: -t 4) (Li and Durbin 2009; Li *et al.* 2009). Following this we counted the proportion  
193 of reads mapping to each TE sequence of all reads, and the coverage of each TE sequence,  
194 weighted by the median coverage of the Muller element D. We removed all TE sequences with  
195 coverage for less than 80% of the sequence for less than 1x the median coverage of Muller  
196 element D, checked using bedTools GenomeCoverage (Quinlan and Hall 2010). In *D. innubila*  
197 we find 6136 TE copies, primarily TIR and RC DNA transposons (2688 and 2423 copies  
198 respectively). In *D. falleni* and *D. phalerata*, we see an expansion of LINE retroposons (1107  
199 and 1793 copies respectively, versus 797 copies in *D. innubila*). We find a significant correlation  
200 between copy numbers of families for pairwise comparisons of all three species (Pearson's  
201 correlation = 0.51-0.68,  $p$ -value < 2.42e-11,  $t$ -value > 7.174), though specific families seem to  
202 differ wildly in copy numbers between species (Supplementary Figure 6C).

203 Finally we also used dnaPipeTE to get an independent estimate of the TE content  
204 (Goubert *et al.* 2015), using the *D. innubila* estimated genome size and the next generation  
205 sequencing information for each species (dnaPipeTE parameters: 2 iterations of trinity, 168Mb  
206 genome size, 1x estimated genome coverage reads). Comparing between species, we find even  
207 more dramatic differences, including a huge expansion of simple repeats in the *D. falleni*  
208 genome, accompanying an expansion of LINE elements, and an expansion of LTRs and TIRs in  
209 *D. phalerata* (Supplementary Figure 6C). Notably, these do not match the estimated TE  
210 proportions in Supplementary Figure 6B, it suggests *D. falleni* and *D. phalerata* contain TE  
211 families not present in *D. innubila*.

212 To identify TEs with orthology to known sequences, we used Blastn (parameters: -evalue  
213 0.00001) against the Repbase Arthropod TE database (Bao *et al.* 2015). We grouped sequences  
214 with hits to previously identified TEs by the TE order and species family of the host. For 92 of  
215 the 175 TE families, we could identify a closely related TE sequence in a previously sequenced  
216 genome from RepBase (Supplementary Figure 7, Blastn,  $e$ -value < 0.001). Most these families  
217 (73.9%) are DNA transposons and LTRs, consistent with previous findings that these orders are

218 more likely to be more recently horizontally transmitted, compared to LINEs (Bartolomé *et al.*  
219 2009; Peccoud *et al.* 2017). 86 of these putatively horizontally transferred TE sequences are  
220 found in another *Drosophila* genome, with 6 TIR families with Blastn hits for Carpenter ants  
221 (*Camponotus*), likely found in the same environment as *D. innubila* (Patterson and Stone 1949;  
222 Markow and O’Grady 2006). Among the TEs with hits to *Drosophila*, only 32 (35.9%) are to  
223 *Drosophila* subgroup species thought to overlap in range with *D. innubila*, the remainder are  
224 species within the *Sophophora* subgroup (Supplementary Figure 7). While 54 of these TE  
225 families have hits to *Sophophora* species found in overlapping ranges with *D. innubila*, such as  
226 species in the *pseudoobscura*, *willistoni* and *ananassae* (within *melanogaster*) groups (Markow  
227 and O’Grady 2006), several TEs (20 TEs with hits to *melanogaster* group), show no evidence of  
228 this, with hits to species endemic to Asia or Africa (Supplementary Figure 7). This may be  
229 because these TEs share a common ancestor in the genome of an unsequenced species that has  
230 overlapping ranges with both *D. innubila* and the *melanogaster* group species.

231

### 232 *Identifying duplications*

233 We identified the 1014 genes present in multiple copies in *D. innubila*, but only present  
234 as single copies in *D. virilis* and *D. melanogaster*. Most these (866) have the duplicated copy on  
235 the same chromosome, with most these duplicates (848) within 50kb of the original copy  
236 (determined by the position of the ortholog in *D. melanogaster*). These duplications are enriched  
237 for metal ion transport and protein metabolism genes (GORilla,  $p$ -value < 0.0005, FDR < 0.05,  
238 enrichment > 1.65) (Eden *et al.* 2009), including 26 cytochrome P450 recent duplications. For  
239 each set of duplicates, we extracted the coding sequence and aligned using PRANK (-codon +F  
240 -f=paml) (Löytynoja 2014). We identified positive selection between orthologs using codeML,  
241 for models M0, M1a, M2a and M3 (Yang 2007). We used a likelihood ratio test to identify  
242 which model fits best for each set of orthologs and to identify duplicates under putative adaptive  
243 evolution. Of duplicate genes, 294 (28.9%) showed signatures of positive selection, a higher  
244 proportion than seen in non-duplicated genes (4.7% , Supplementary Figure 8,  $dN/dS > 1$ , Model  
245 2a is best fitting model). We find no association between the number of copies of a gene and the  
246  $dN/dS$  (GLM,  $t$ -value = 0.27,  $p$ -value = 0.78) and as shown previously, find negative correlations  
247 between  $dN/dS$  and both gene length and dS (Supplementary Table 11, Supplementary Figure 8,  
248 GLM,  $t$ -value < -2.2353,  $p$ -value < 0.0188). Using GORilla we found that, like the total



249 complement of paralogs, these duplicated genes under positive selection are again enriched for  
250 Metal ion transport, specifically the copper ion response pathways (Supplementary Figure 8,  
251 Supplementary Table 11, GOrilla,  $p$ -value < 0.0005, FDR < 0.05, enrichment > 1.25) (Eden *et al.*  
252 2009).

253

#### 254 *RNA differential expression analysis*

255 We downloaded mapped RNA sequencing information from ModEncode  
256 (modencode.org) for *D. melanogaster* across all life stages (Chen *et al.* 2014).

257 For each set of *D. innubila* RNA sequencing short read information we mapped it to the  
258 masked *D. innubila* genome with the TE sequences concatenated to the end using GSNAP. We  
259 then counted the number of reads mapped to each gene per kb of gene using HTSeq for all  
260 mapped RNA sequencing data and normalized by counts per million per dataset (Anders *et al.*  
261 2015).

262 We then used the R package EdgeR (Robinson *et al.* 2009) to make differential  
263 expression comparisons between the following datasets: 1. Adult total body *D. innubila* RNA,  
264 male versus female; 2. RNA across different life stages total body; 3. Adult total body, *D.*  
265 *innubila* female versus *D. melanogaster* female; 4. Adult total body, *D. innubila* male versus *D.*  
266 *melanogaster* male; 5. Adult total body, *D. innubila* versus *D. melanogaster*; 6. Larvae total  
267 body, *D. innubila* versus *D. melanogaster*; 7. Pupae total body, *D. innubila* versus *D.*  
268 *melanogaster*; 8. Whole embryo, *D. innubila* versus *D. melanogaster*. In each case, we compared  
269 the counts per million per 1kbp exon of genes to identify significant differences in expression of  
270 orthologous genes ( $p$ -value < 0.05, FDR < 0.05 after adjusting for multiple testing).

271 Following this, we used GOrilla (Eden *et al.* 2009) to identify and visualize enriched  
272 gene ontology (GO) terms, separating by genes that are and aren't differentially expressed ( $p$ -  
273 value threshold = 0.001) for process, function and component GO terms. For functional terms of  
274 interest, such as detoxification genes between species, recent duplications versus their single  
275 copy, novel genes across life stages or viral RNAi genes, we compared expression differences  
276 between groups by hand. Across the life stages between *D. innubila* and *D. melanogaster*, we  
277 find changes in gene expression, including enrichments such as muscle system process genes and  
278 structural muscle construction. We also find differential expression metabolic processes, cellular  
279 process and locomotion across all life stages (Supplementary Tables 12-19, Supplementary

280 Figure 9, GOrilla, FDR < 0.00005,  $p$ -value < 0.000984 after multiple testing correction,  
281 enrichment > 1.21), it is important to highlight that these differences identified could be due to  
282 differences in experimental setting used to generate the data, or could be due to differences  
283 between *D. innubila* and *D. melanogaster*.

284 Using our gene expression data for both male and female adult *D. innubila*, we looked for  
285 biases expected between sexes. Surprisingly, we find no genes with a significant female bias  
286 expression (0 genes, Supplementary Figure 10, Supplementary Table 7, 13 & 20, EdgeR  $p$ -value  
287 >0.206 FDR > 0.0006 after Bonferroni multiple testing correction), with a large number showing  
288 a male bias (Supplementary Figure 10, 223 genes, EdgeR  $p$  < 0.000001 after multiple testing  
289 correction). As is expected there is a significant deficit of male bias genes on the X chromosome  
290 (Supplementary Tables 7 & 20, Chi-Square test  $\chi^2=4.21$ ,  $p$ -value = 0.04), though we also see an  
291 enrichment on one of the autosomes, Muller element B (Chi-Square test  $\chi^2 = 16.86$ ,  $p$ -value =  
292 4.03e-5). We used GOrilla to identify any enrichment in categories between sexes which may  
293 explain the difference observed. We find an enrichment for organophosphate metabolism, cell  
294 motility and sperm movement (GOrilla enrichment > 17.27,  $p$ -value < 0.000654).

295

#### 296 *Structural variants between species in the D. innubila trio*

297 We next estimated structural variants between *D. innubila*, *D. falleni* and *D. phalerata*, using  
298 Pindel (Ye *et al.* 2009) and short reads mapped to the *D. innubila* genome. We find many more  
299 structural variants and inversions between *D. phalerata* and *D. innubila* than *D. falleni*,  
300 consistent with structural variants accumulating as species diverge (Supplementary Figures 11 &  
301 12). We find no significant effects of inversions on  $dN/dS$  or  $dS$  between species (Mann-Whitney  
302 U test  $W < 156$ ,  $p$ -value > 0.41).

303

#### 304 **Supplementary Tables and Figures**

305 **Supplementary Table 1:** Summary of reads used for genome sequencing, assembly, annotation  
306 and  $dN/dS$  calculation.

307 **Supplementary Table 2:** Summary statistics for each iteration of the genome.

308 **Supplementary Table 3:** Summary of the genic characteristics of the *D. innubila* genome.

309 **Supplementary Table 4:** Genes ontologies (GO) enriched for genes with high/low residuals for  
310  $dN/dS$  between *D. melanogaster* and *D. innubila*, due to drastic differences between the species.

311 Enriched categories are categories which are slow evolving in one species, but fast evolving in the  
312 other.

313 **Supplementary Table 5:** Summary of  $dN/dS$  statistics for each immune gene category across the  
314 total group and on each branch. Additionally, the t-score and p-value for a two-sided  $t$ -test ( $\mu = 0$ )  
315 for that category is shown. Significant categories are highlighted in bold.

316 **Supplementary Table 6:**  $dN/dS$  enrichment for *Drosophila innubila* trio for processes,  
317 components and functions, including any enrichments for specific branches.

318 **Supplementary Table 7:** GO enrichment for processes, components and functions for differential  
319 expression between *D. innubila* males and females.

320 **Supplementary Table 8:** GO enrichment for processes, components and functions for differential  
321 expression between *D. innubila* embryos and larvae.

322 **Supplementary Table 9:** GO enrichment for processes, components and functions for differential  
323 expression between *D. innubila* larvae and pupae.

324 **Supplementary Table 10:** GO enrichment for processes, components and functions for  
325 differential expression between *D. innubila* pupae and adults.

326 **Supplementary Table 11:**  $dN/dS$  GO enrichment for duplications for processes, components and  
327 functions, including any enrichments for specific branches.

328 **Supplementary Table 12:** A table summarizing the differential gene expression shown in  
329 Supplementary Tables 13-19, showing the number of genes differentially expressed between *D.*  
330 *innubila* and *D. melanogaster* at differing life stages, with enrichments in gene ontology (GO)  
331 categories.

332 **Supplementary Table 13:** GO enrichment for processes, components and functions for  
333 differential expression between *D. melanogaster* and *D. innubila* embryos.

334 **Supplementary Table 14:** GO enrichment for processes, components and functions for  
335 differential expression between *D. melanogaster* and *D. innubila* larvae.

336 **Supplementary Table 15:** GO enrichment for processes, components and functions for  
337 differential expression between *D. melanogaster* and *D. innubila* pupae.

338 **Supplementary Table 16:** GO enrichment for processes, components and functions for  
339 differential expression between *D. melanogaster* and *D. innubila* adults.

340 **Supplementary Table 17:** GO enrichment for processes, components and functions for  
341 differential expression between *D. melanogaster* and *D. innubila* adult males.

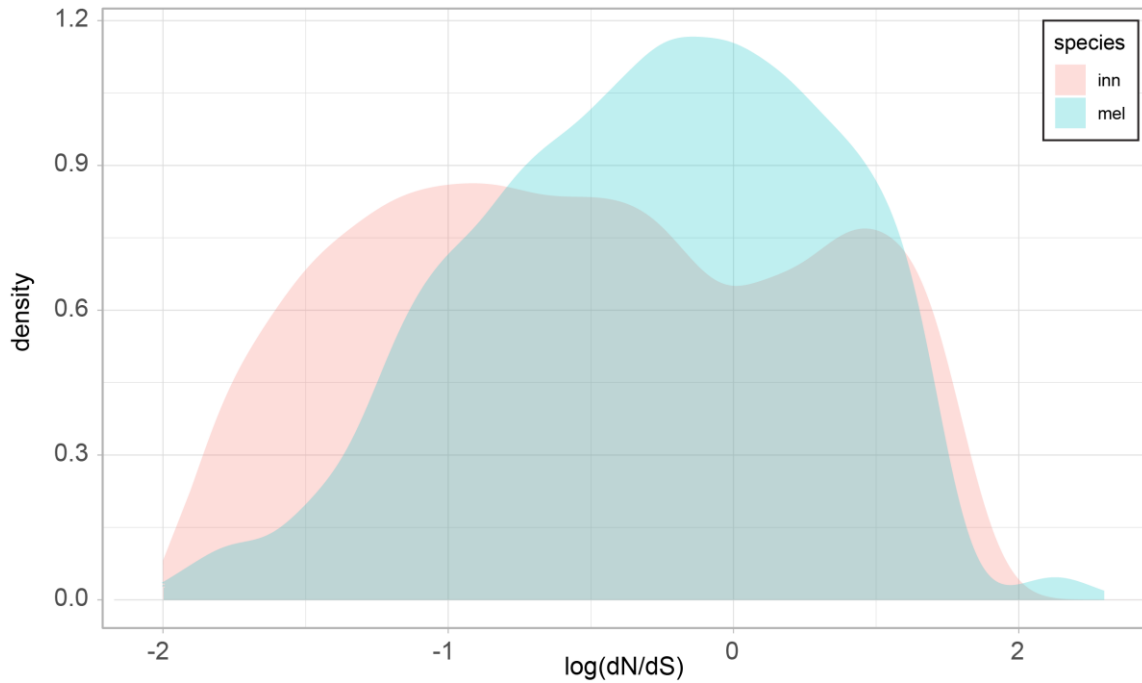
342 **Supplementary Table 18:** GO enrichment for processes, components and functions for  
343 differential expression between *D. melanogaster* and *D. innubila* adult females.

344 **Supplementary Table 19:** GO enrichment for processes, components and functions for  
345 differential expression between *D. melanogaster* and *D. innubila* total samples.

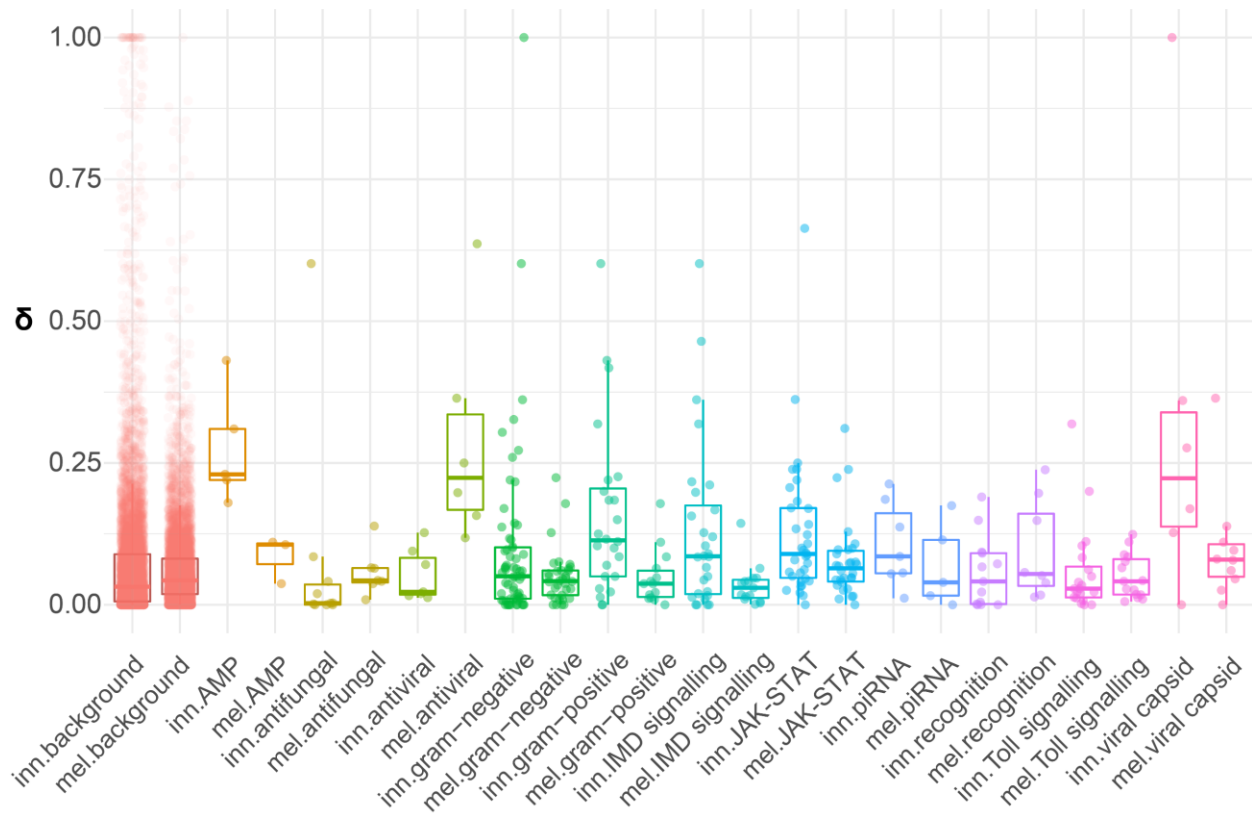
346 **Supplementary Table 20:** Enrichment or depletion of genes differentially expressed between  
347 male and female samples on each scaffold/Muller element including the  $\chi^2$  for this enrichment.

348

349 **Supplementary Figure 1:** Histograms of dN/dS for *D. innubila* and *D. melanogaster*.



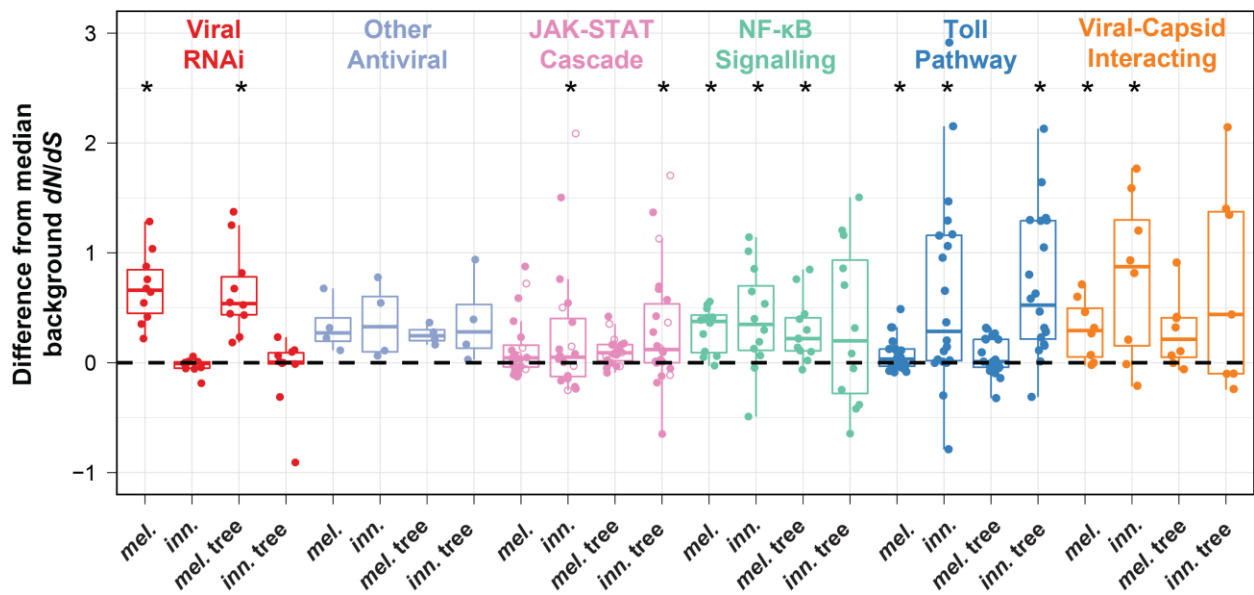
350  
 351 **Supplementary Figure 2:**  $\delta$  (calculated using HyPhy) by immunity category for both *D.*  
 352 *innubila* and *D. melanogaster*.



353  
 354

355 **Supplementary Figure 3: Antiviral evolution across the *quinaria* group.** Difference between  
 356 viral RNAi, JAK-STAT (filled dots = regulatory, empty dots = cytokines), NF- $\kappa$ B, Toll  
 357 and putatively viral-interacting proteins from the background  $dN/dS$  of genes of similar  $dS$   
 358 ( $\pm 0.01dS$ ) for the *D. melanogaster* branch, the *D. innubila* branch, the total *D.*  
 359 *melanogaster* tree and the total *D. innubila* trio. Genes known to be associated with the  
 360 immune response to viral infection, but no known pathway are classed as ‘Other Antiviral’.  
 361 A  $p$ -value (from a two-sided  $t$ -test looking for significant differences from 0) of 0.05 or  
 362 lower is designated with \*.

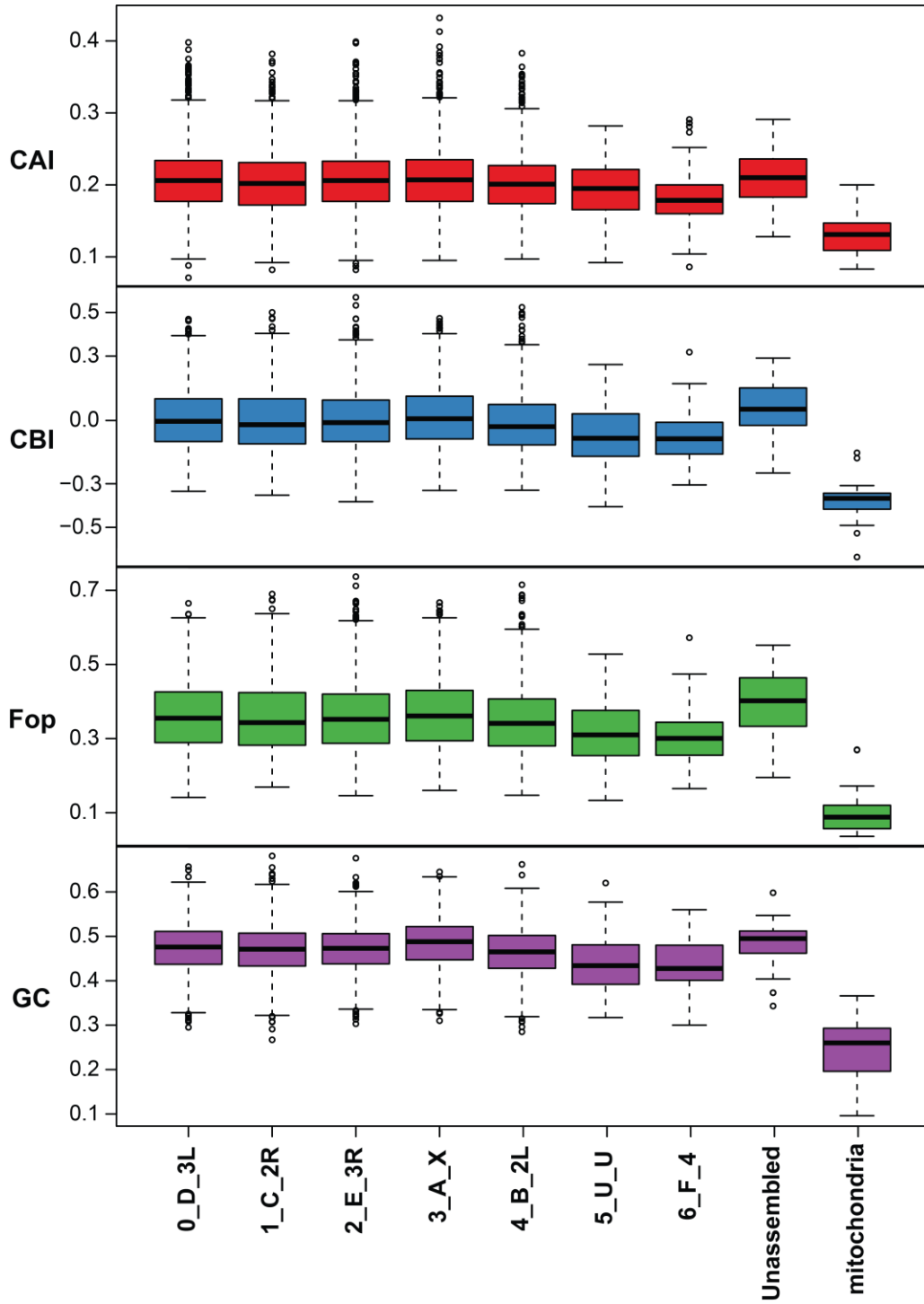
363



364

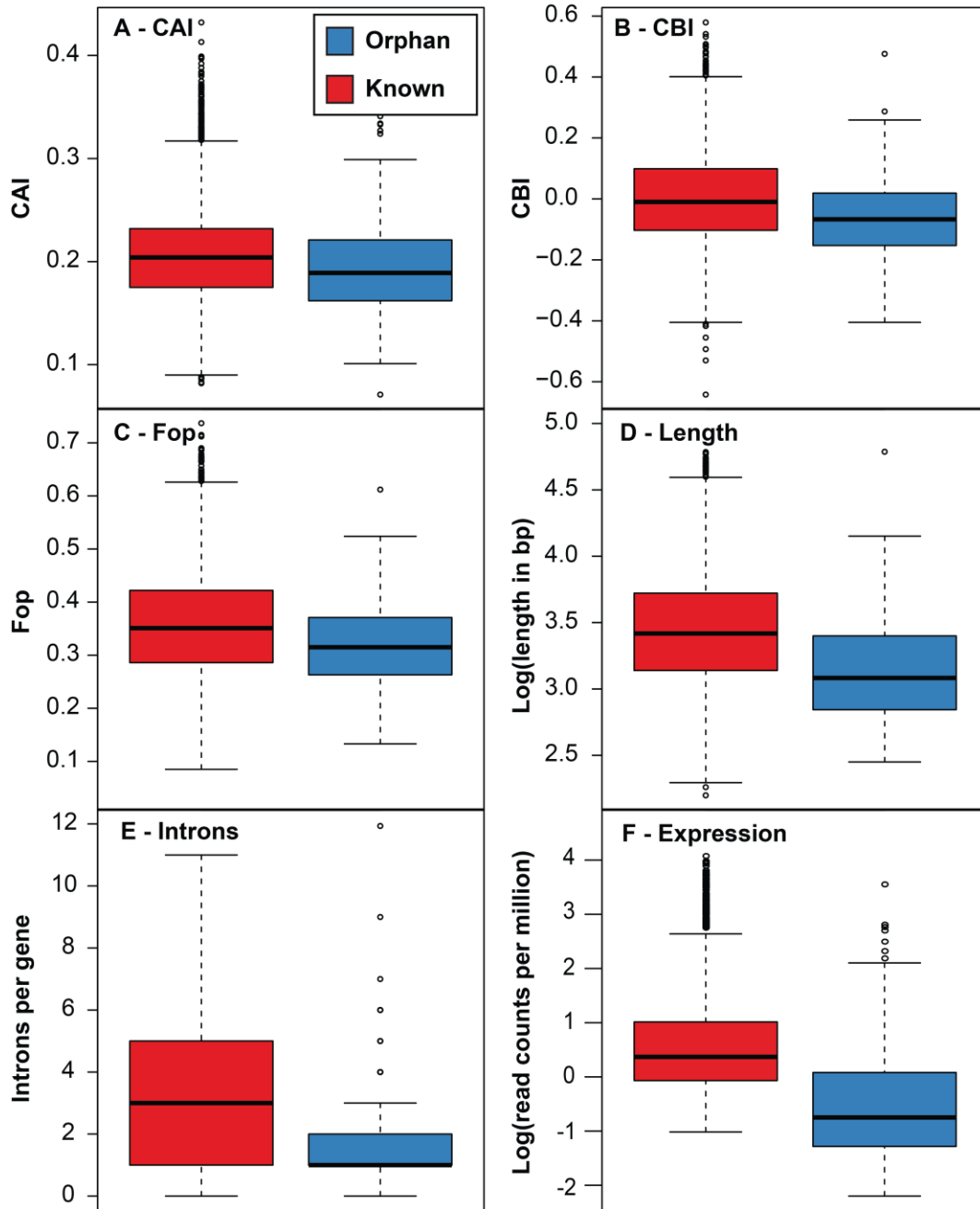
365

366 **Supplementary Figure 4:** Codon bias distributions across the *Drosophila innubila* genome,  
 367 separated by scaffold. CAI = Codon adaptation index. CBI = Codon bias index. Fop = Frequency  
 368 of optimal codons. GC = Proportion of GC across each gene.



369  
 370

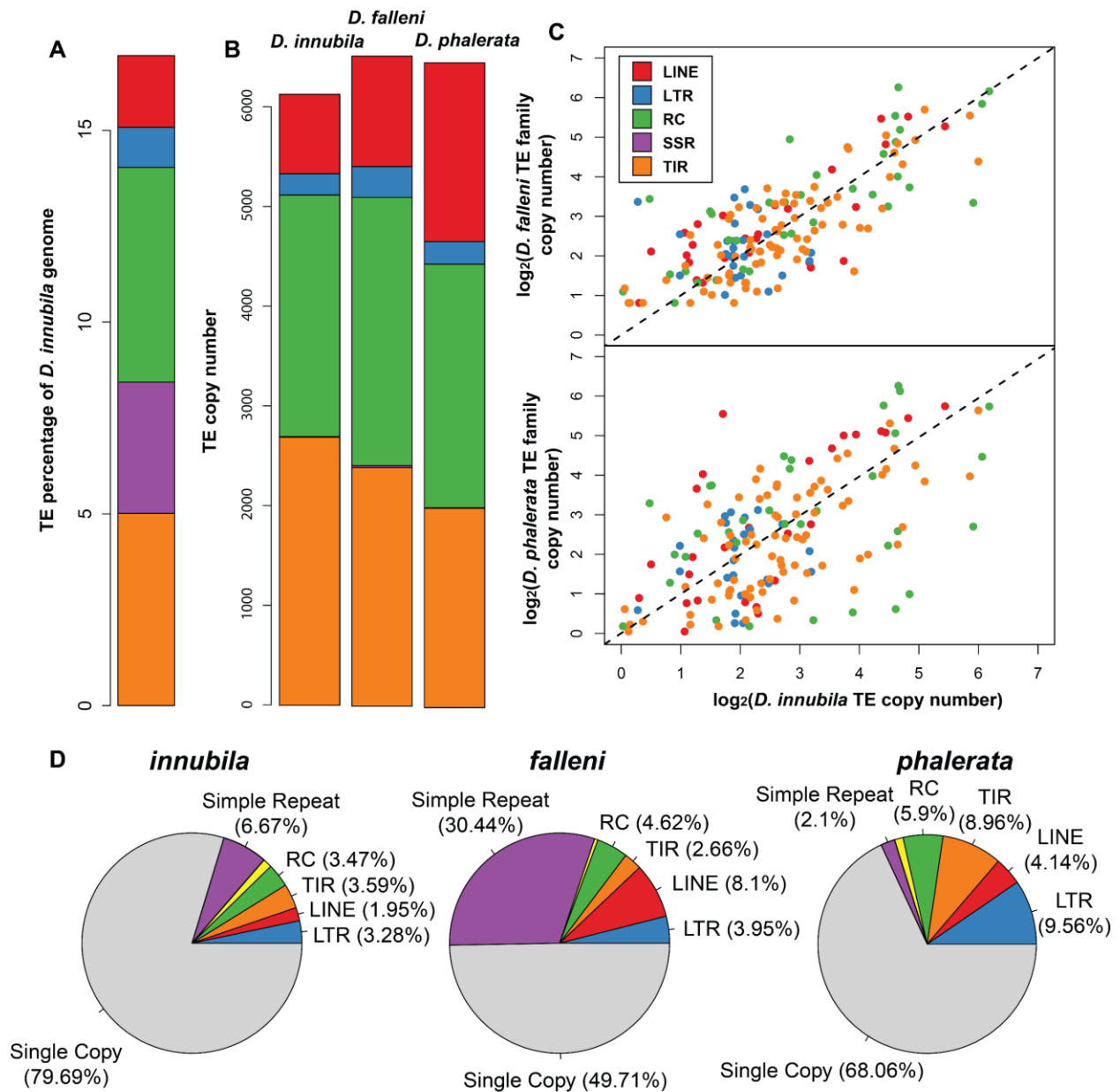
371 **Supplementary Figure 5:** Comparison between orphan genes and previously described genes,  
372 including: **A.** Codon adaptation index (CAI). **B.** Codon bias index (CBI). **C.** Frequency of  
373 optimal codons (Fop). **D.** Gene length (in bp). **E.** Number of introns per gene. **F.** Mean  
374 expression across life stages (read counts per million).



375  
376  
377



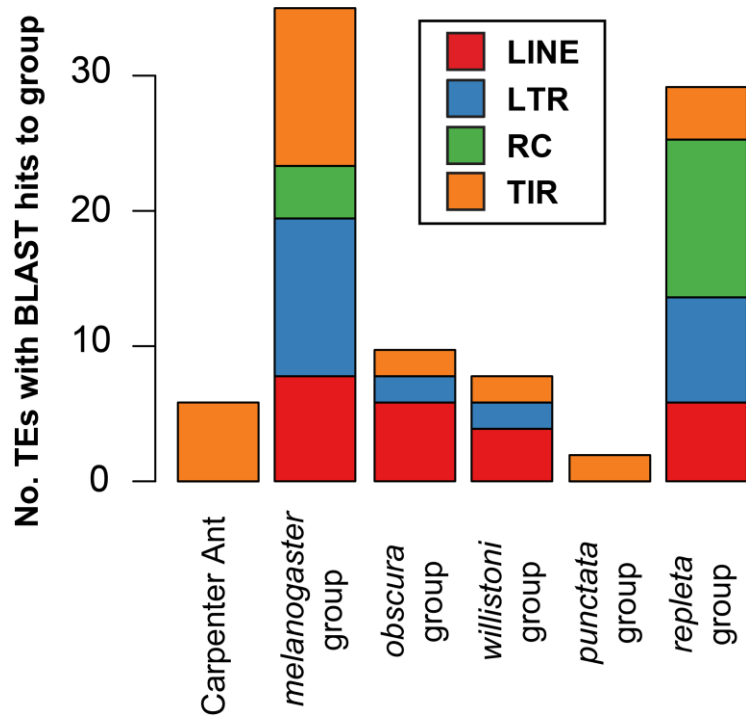
378 **Supplementary Figure 6: A.** The proportion of the *D. innubila* genome masked by each type of  
 379 repeat. LINE = Long interspersed nuclear element RNA transposon, LTR = long terminal repeat  
 380 RNA transposon, RC = rolling circle DNA transposon, TIR = terminal inverted repeat DNA  
 381 transposon. **B.** TE content of *D. innubila*, *falleni* and *phalerata*, **C.** Copy number comparisons  
 382 between *D. innubila*, *D. falleni* and *phalerata*. **D.** dnaPipeTE estimates of the genomic  
 383 proportion of repetitive elements for each species examined. Other, NA and SINE categories  
 384 were removed due to small proportions. Though unlabeled, rRNA is shown in yellow and  
 385 constitutes 1-2% of the genome.



386

387 **Supplementary Figure 7:** Number of TE families found in *D. innubila*, closely related to known  
388 TE families (taken from Repbase) in different species group, identified using BLAST, suggesting  
389 relatively recent horizontal transfer events.

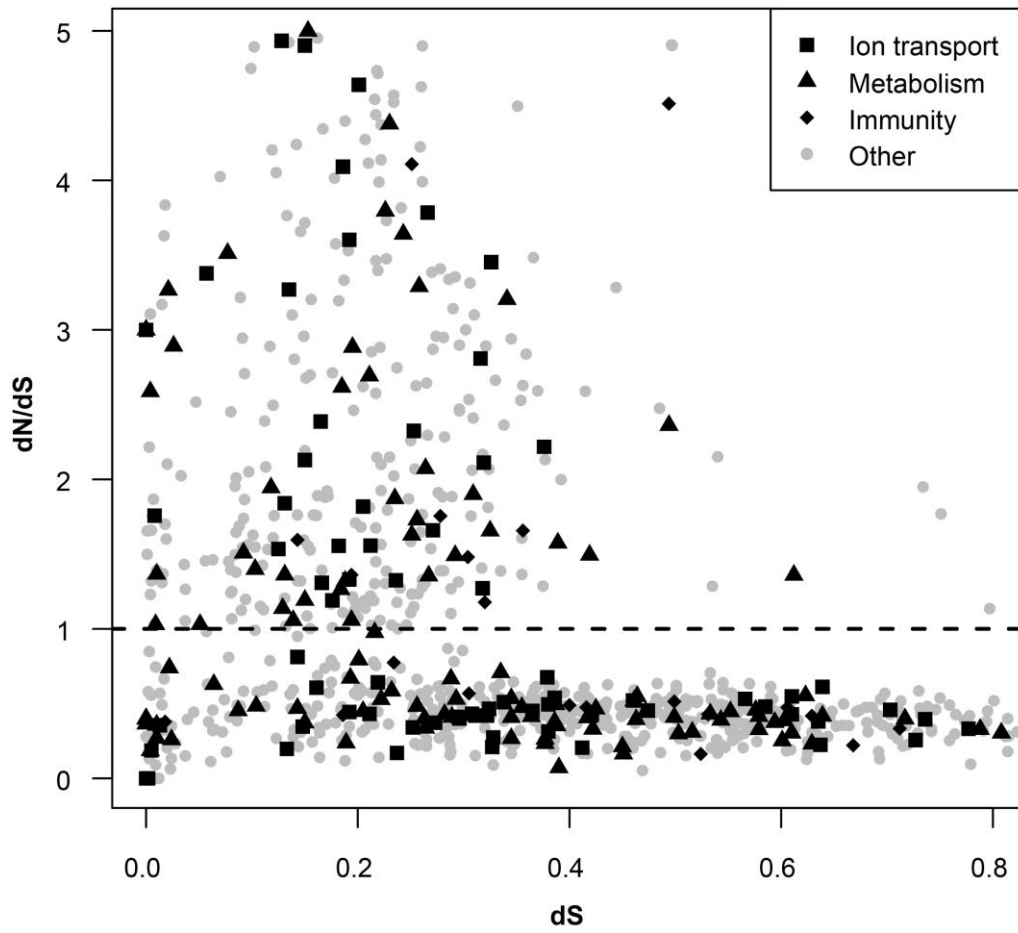
390



391

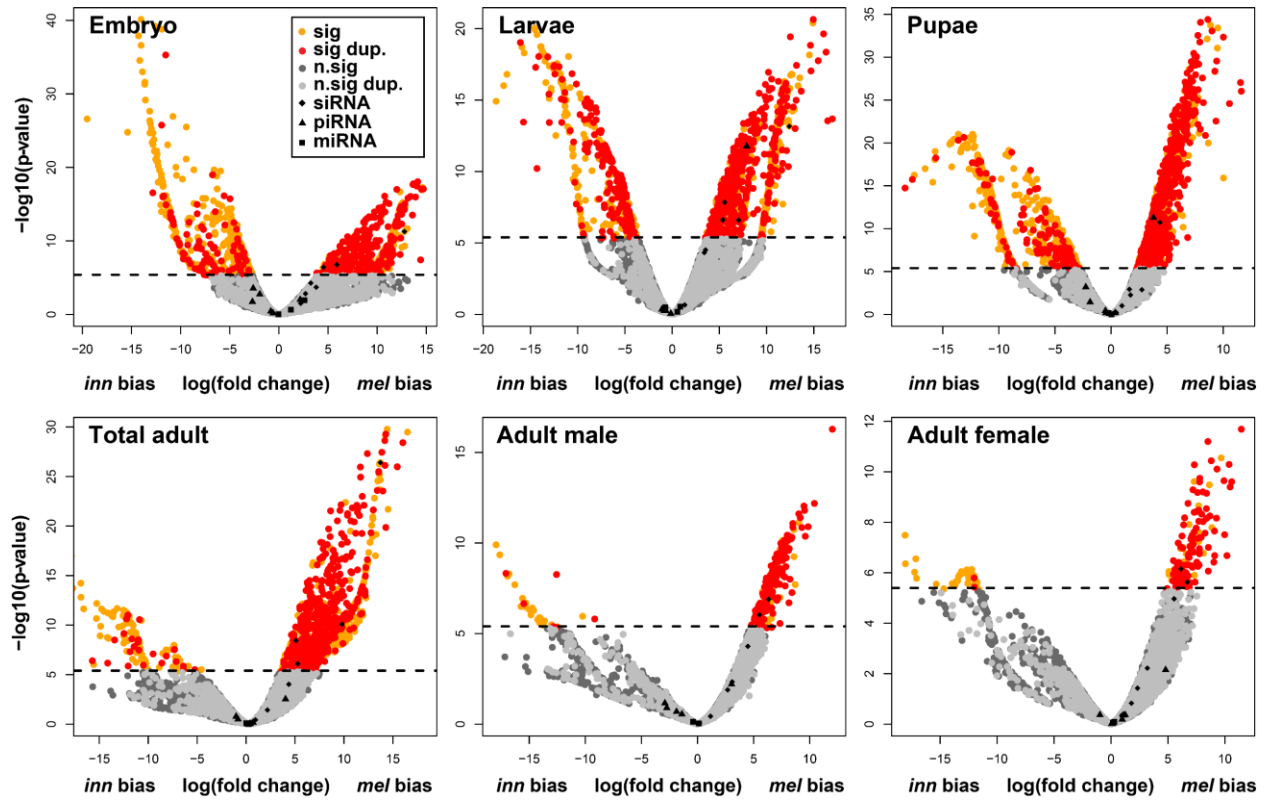
392

393 **Supplementary Figure 8:**  $dN/dS$  versus  $dS$  across paralogs for recently duplicated genes. Metal  
394 ion binding, protein metabolism and immunity genes are highlighted.



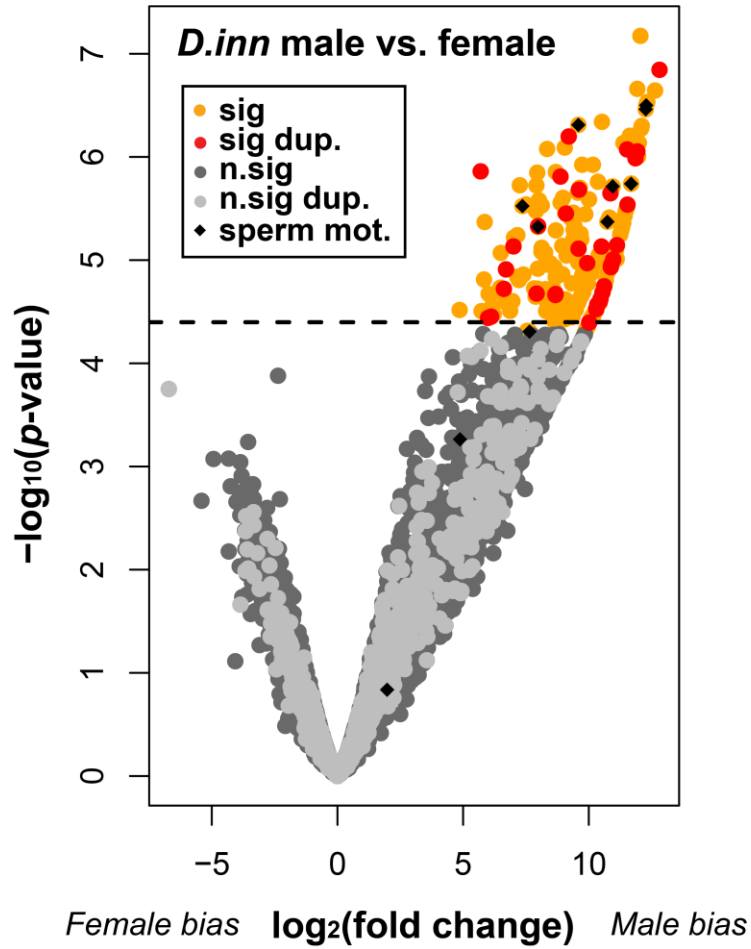
395

396 **Supplementary Figure 9:** Volcano plots showing differential gene expression between *D.*  
 397 *innubila* and *D. melanogaster* at different life stages. Dots are colored by their significance and if  
 398 a recent duplication or not (duplicants layered on top), the significance cut off is set a 0.05  
 399 following Bonferroni multiple testing correction.



400

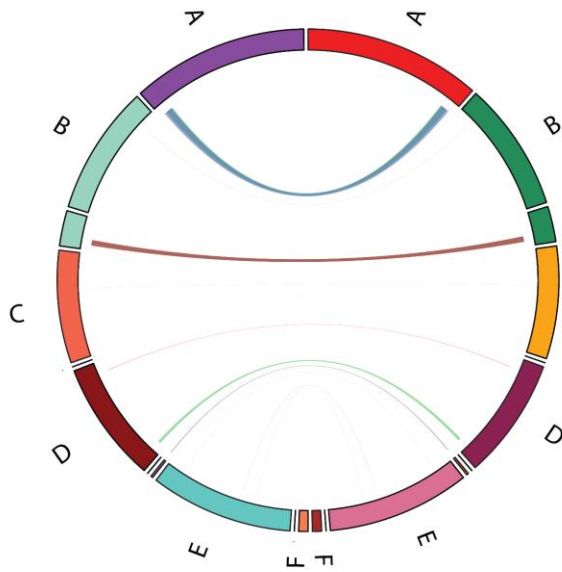
401 **Supplementary Figure 10:** Volcano plot showing differential gene expression between *D.*  
402 *innubila* male and female samples and significant differences, highlighting if genes are  
403 duplicated relative to *D. virilis* or not, and if genes are involved in sperm motility.



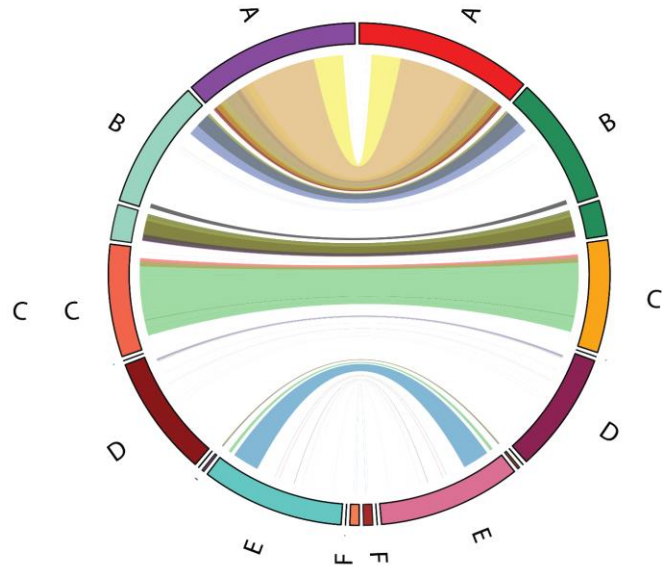
404

405 **Supplementary Figure 11:** Inversions identified between *D. innubila* and *D. falleni*, and  
406 between *D. innubila/falleni* and *D. phalerata* using Pindel (Ye et al. 2009) and Manta (Chen et  
407 al. 2016) (taking the consensus of the two programs). Scaffolds are labelled and colored by the  
408 Muller element they belong to.

**A - falleni inversions**

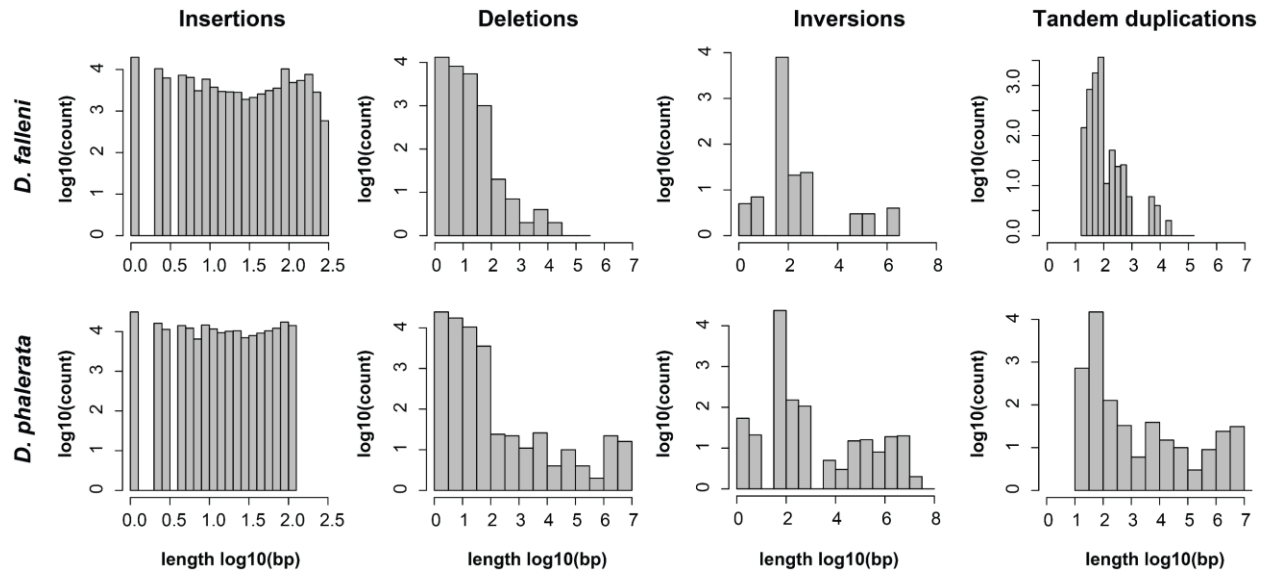


**B - phalerata inversions**



409

410 **Supplementary Figure 12:** Size and number of each structural variant between *D. innubila* and  
 411 *D. falleni* identified using Pindel and Manta (taking the consensus of the two programs).



412

413

#### 414 **Supplementary References**

415 Altschul S. F., W. Gish, W. Miller, E. W. Myers, and D. J. Lipman, 1990 Basic local alignment  
 416 search tool. *J. Mol. Biol.* 215: 403–410. [https://doi.org/http://dx.doi.org/10.1016/S0022-](https://doi.org/http://dx.doi.org/10.1016/S0022-2836(05)80360-2)  
 417 2836(05)80360-2

418 Anders S., P. T. Pyl, and W. Huber, 2015 HTSeq-A Python framework to work with high-  
 419 throughput sequencing data. *Bioinformatics* 31: 166–169.

420 <https://doi.org/10.1093/bioinformatics/btu638>

421 Bankevich A., S. Nurk, D. Antipov, A. A. Gurevich, M. Dvorkin, *et al.*, 2012 SPAdes: A new  
 422 genome assembly algorithm and its applications to single-cell sequencing. *J. Comput. Biol.*  
 423 19: 455–477. <https://doi.org/10.1089/cmb.2012.0021>

424 Bao W., K. K. Kojima, and O. Kohany, 2015 Repbase Update, a database of repetitive elements  
 425 in eukaryotic genomes. *Mob. DNA* 6: 4–9. <https://doi.org/10.1186/s13100-015-0041-9>

426 Bartolomé C., X. Bello, and X. Maside, 2009 Widespread evidence for horizontal transfer of  
 427 transposable elements across *Drosophila* genomes. *Genome Biol.* 10: R22.

428 <https://doi.org/10.1186/gb-2009-10-2-r22>

429 Bernt M., A. Donath, F. Jühling, F. Externbrink, C. Florentz, *et al.*, 2013 MITOS: Improved de  
 430 novo metazoan mitochondrial genome annotation. *Mol. Phylogenet. Evol.* 69: 313–319.

431 <https://doi.org/10.1016/j.ympev.2012.08.023>

432 Buffalo V., 2018 Scythe  
433 Camacho C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos, *et al.*, 2009 BLAST+:  
434 Architecture and applications. *BMC Bioinformatics* 10: 1–9. <https://doi.org/10.1186/1471->  
435 [2105-10-421](https://doi.org/10.1186/1471-2105-10-421)  
436 Chakraborty M., R. Zhao, X. Zhang, S. Kalsow, and J. J. Emerson, 2017 Extensive hidden  
437 genetic variation shapes the structure of functional elements in *Drosophila*. *Doi.Org* 50:  
438 114967. <https://doi.org/10.1101/114967>  
439 Chen Z., D. Sturgill, J. Qu, and H. Jiang, 2014 Comparative validation of the *D. melanogaster*  
440 modENCODE transcriptome annotation. *Genome ...* 1209–1223.  
441 <https://doi.org/10.1101/gr.159384.113>.  
442 Chen X., O. Schulz-Trieglaff, R. Shaw, B. Barnes, F. Schlesinger, *et al.*, 2016 Manta: Rapid  
443 detection of structural variants and indels for germline and cancer sequencing applications.  
444 *Bioinformatics* 32: 1220–1222. <https://doi.org/10.1093/bioinformatics/btv710>  
445 Clark A. G., M. B. Eisen, D. R. Smith, C. M. Bergman, B. Oliver, *et al.*, 2007 Evolution of genes  
446 and genomes on the *Drosophila* phylogeny. *Nature* 450: 203–218.  
447 <https://doi.org/10.1038/nature06341>  
448 Darling A. C. E., B. Mau, F. R. Blattner, and N. T. Perna, 2004 Mauve : Multiple Alignment of  
449 Conserved Genomic Sequence With Rearrangements Mauve : Multiple Alignment of  
450 Conserved Genomic Sequence With Rearrangements. 1394–1403.  
451 <https://doi.org/10.1101/gr.2289704>  
452 Eden E., R. Navon, I. Steinfeld, D. Lipson, and Z. Yakhini, 2009 GOrilla: A tool for discovery  
453 and visualization of enriched GO terms in ranked gene lists. *BMC Bioinformatics* 10: 1–7.  
454 <https://doi.org/10.1186/1471-2105-10-48>  
455 Eilbeck K., B. Moore, C. Holt, and M. Yandell, 2009 Quantitative measures for the management  
456 and comparison of annotated genomes. *BMC Bioinformatics* 10: 1–15.  
457 <https://doi.org/10.1186/1471-2105-10-67>  
458 Gilbert D., 2013 EvidentialGene: mRNA Transcript Assembly Software. *EvidentialGene Evid.*  
459 *Dir. Gene Constr. Eukaryotes*.  
460 Goubert C., L. Modolo, C. Vieira, C. V. Moro, P. Mavingui, *et al.*, 2015 De novo assembly and  
461 annotation of the Asian tiger mosquito (*Aedes albopictus*) repeatome with dnaPipeTE from  
462 raw genomic reads and comparative analysis with the yellow fever mosquito (*Aedes*



463 aegypti). *Genome Biol. Evol.* 7: 1192–1205. <https://doi.org/10.1093/gbe/evv050>

464 Gramates L. S., S. J. Marygold, G. Dos Santos, J. M. Urbano, G. Antonazzo, *et al.*, 2017 FlyBase  
465 at 25: Looking to the future. *Nucleic Acids Res.* 45: D663–D671.  
466 <https://doi.org/10.1093/nar/gkw1016>

467 Haas B. J., A. Papanicolaou, M. Yassour, M. Grabherr, P. D. Blood, *et al.*, 2013 De novo  
468 transcript sequence reconstruction from RNA-seq using the Trinity platform for reference  
469 generation and analysis. *Nat. Protoc.* 8: 1494–512. <https://doi.org/10.1038/nprot.2013.084>

470 Hill T., and R. L. Unckless, 2017 The dynamic evolution of *Drosophila innubila* Nudivirus.  
471 *Infect. Genet. Evol.* 1–25.

472 Holt C., and M. Yandell, 2011 MAKER2 : an annotation pipeline and genome- database  
473 management tool for second- generation genome projects. *BMC Bioinformatics* 12: 491.  
474 <https://doi.org/10.1186/1471-2105-12-491>

475 Jain M., H. E. Olsen, B. Paten, and M. Akeson, 2016 Erratum to: The Oxford Nanopore  
476 MinION: Delivery of nanopore sequencing to the genomics community [ *Genome Biol.*  
477 (2016), 17, 239] DOI:10.1186/s13059-016-1103-0. *Genome Biol.* 17: 1–11.  
478 <https://doi.org/10.1186/s13059-016-1122-x>

479 Joshi N., and J. Fass, 2011 Sickle: A sliding window, adaptive, quality-based trimming tool for  
480 fastQ files. 1.33.

481 Kearse M., R. Moir, A. Wilson, S. Stones-Havas, M. Cheung, *et al.*, 2012 Geneious Basic: An  
482 integrated and extendable desktop software platform for the organization and analysis of  
483 sequence data. *Bioinformatics* 28: 1647–1649.  
484 <https://doi.org/10.1093/bioinformatics/bts199>

485 Koren S., B. P. Walenz, K. Berlin, J. R. Miller, N. H. Bergman, *et al.*, 2016 Canu : scalable and  
486 accurate long- - read assembly via adaptive k - - mer weighting and repeat separation. 1–  
487 35. <https://doi.org/10.1101/gr.215087.116>.Freely

488 Korf I., 2004 Gene finding in novel genomes. *BMC Bioinformatics* 5: 1–9.  
489 <https://doi.org/10.1186/1471-2105-5-59>

490 Li H., and R. Durbin, 2009 Fast and accurate short read alignment with Burrows-Wheeler  
491 transform. *Bioinformatics* 25: 1754–60. <https://doi.org/10.1093/bioinformatics/btp324>

492 Li H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan, *et al.*, 2009 The sequence alignment/map  
493 format and SAMtools. *Bioinformatics* 25: 2078–9.

494 <https://doi.org/10.1093/bioinformatics/btp352>

495 Löytynoja A., 2014 Phylogeny-aware alignment with PRANK, pp. 155–170 in *Multiple*

496 *Sequence Alignment Methods*, edited by Russell D. J. Humana Press, Totowa, NJ.

497 Marcais G., 2011 Jellyfish : A fast k-mer counter. 1–5.

498 Markow T. A., and P. O’Grady, 2006 *Drosophila: a guide to species identification*.

499 McCoy R. C., R. W. Taylor, T. A. Blauwkamp, J. L. Kelley, M. Kertesz, *et al.*, 2014 Illumina

500 TruSeq synthetic long-reads empower de novo assembly and resolve complex, highly-

501 repetitive transposable elements. PLoS One 9. <https://doi.org/10.1371/journal.pone.0106689>

502 Palmieri N., C. Kosiol, and C. Schlötterer, 2014 The life cycle of *Drosophila* orphan genes. *Elife*

503 3: 1–21. <https://doi.org/10.7554/eLife.01311>

504 Patterson J. T., and W. S. Stone, 1949 Studies in the genetics of *Drosophila*. Univ. Texas Publ.

505 4920: 7–17.

506 Peccoud J., V. Loiseau, R. Cordaux, and C. Gilbert, 2017 Massive horizontal transfer of

507 transposable elements in insects. *Proc. Natl. Acad. Sci.* 114: 4721–4726.

508 <https://doi.org/10.1073/pnas.1621178114>

509 Peden J. F., 1997 CodonW

510 Quinlan A. R., and I. M. Hall, 2010 BEDTools: a flexible suite of utilities for comparing

511 genomic features. *Bioinformatics* 26: 841–2. <https://doi.org/10.1093/bioinformatics/btq033>

512 Robinson M. D., D. J. McCarthy, and G. K. Smyth, 2009 edgeR: A Bioconductor package for

513 differential expression analysis of digital gene expression data. *Bioinformatics* 26: 139–140.

514 <https://doi.org/10.1093/bioinformatics/btp616>

515 Schulz M. H., D. R. Zerbino, M. Vingron, and E. Birney, 2012 Oases: Robust de novo RNA-seq

516 assembly across the dynamic range of expression levels. *Bioinformatics* 28: 1086–1092.

517 <https://doi.org/10.1093/bioinformatics/bts094>

518 Simão F. A., R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, and E. M. Zdobnov, 2015

519 BUSCO: Assessing genome assembly and annotation completeness with single-copy

520 orthologs. *Bioinformatics* 31: 3210–3212. <https://doi.org/10.1093/bioinformatics/btv351>

521 Simms D., P. Cizdziel, and P. Chomczynski, 1993 TRIzol: a new reagent for optimal single-step

522 isolation of RNA. *Focus (Madison)*. 99–102. [https://doi.org/http://dx.doi.org/10.1016/0003-](https://doi.org/http://dx.doi.org/10.1016/0003-2670(61)80041-X)

523 [2670\(61\)80041-X](https://doi.org/http://dx.doi.org/10.1016/0003-2670(61)80041-X)

524 Smit A. F. A., and R. Hubley, 2008 RepeatModeler Open-1.0

525 Smit A. F. A., and R. Hubley, 2015 RepeatMasker Open-4.0  
526 Stanke M., M. Diekhans, R. Baertsch, and D. Haussler, 2008 Using native and syntenically  
527 mapped cDNA alignments to improve de novo gene finding. *Bioinformatics* 24: 637–644.  
528 <https://doi.org/10.1093/bioinformatics/btn013>  
529 Vurture G., F. Sedlazeck, M. Nattestad, C. Underwood, H. Fang, *et al.*, 2017 GenomeScope: Fast  
530 reference-free genome profiling from short reads. *Biorxiv* 0–0.  
531 <https://doi.org/10.1093/bioinformatics/xxxxx>  
532 Walker B. J., T. Abeel, T. Shea, M. Priest, A. Abouelliel, *et al.*, 2014 Pilon: An integrated tool  
533 for comprehensive microbial variant detection and genome assembly improvement. *PLoS*  
534 *One* 9. <https://doi.org/10.1371/journal.pone.0112963>  
535 Wu T. D., and S. Nacu, 2010 Fast and SNP-tolerant detection of complex variants and splicing in  
536 short reads. *Bioinformatics* 26: 873–881. <https://doi.org/10.1093/bioinformatics/btq057>  
537 Xie Y., G. Wu, J. Tang, R. Luo, J. Patterson, *et al.*, 2014 SOAPdenovo-Trans: De novo  
538 transcriptome assembly with short RNA-Seq reads. *Bioinformatics* 30: 1660–1666.  
539 <https://doi.org/10.1093/bioinformatics/btu077>  
540 Yang Z., 2007 PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24:  
541 1586–1591. <https://doi.org/10.1093/molbev/msm088>  
542 Ye K., M. H. Schulz, Q. Long, R. Apweiler, and Z. Ning, 2009 Pindel : a pattern growth  
543 approach to detect break points of large deletions and medium sized insertions from paired-  
544 end short reads. *Bioinformatics* 25: 2865–2871.  
545 <https://doi.org/10.1093/bioinformatics/btp394>  
546 Zhou Q., and D. Bachtrog, 2015 Ancestral Chromatin Configuration Constrains Chromatin  
547 Evolution on Differentiating Sex Chromosomes in *Drosophila*. *PLoS Genet.* 11: 1–21.  
548 <https://doi.org/10.1371/journal.pgen.1005331>  
549