

1 Molecular population genetics of *Sex-lethal* (*Sxl*) in the *D. melanogaster* species group - a locus
2 that genetically interacts with *Wolbachia pipientis* in *Drosophila melanogaster*.

3

4 Vanessa L. Bauer DuMont¹, Simone L. White, Daniel Zinshteyn and Charles F. Aquadro²

5

6 Department of Molecular Biology and Genetics

7 Cornell University

8 Ithaca, NY 14853

9

10 ¹ Present address: Department of Molecular, Cellular, and Developmental Biology and the
11 BioFrontiers Institute, University of Colorado Boulder, Boulder, Colorado, USA

12

13

14 ² Corresponding author: Department of Molecular Biology and Genetics, Cornell University,
15 Ithaca, NY 14853, E-mail: CFA1@CORNELL.EDU

16

17

18 Genbank accession numbers: KT935592-KT935663 (will be made public upon publication)

19

20

21 Running title: Molecular Evolution of *Sex-Lethal*

22 Key words: germline stem cells, *Wolbachia pipientis*, population genetics, natural selection,
23 genetic conflict

24 Corresponding author: Charles F. Aquadro, Department of Molecular Biology and Genetics,
25 Cornell University, Ithaca, NY 14853, Phone: 607-254-4838, email: cfa1@cornell.edu

26

27

28

29 **Abstract**

30 *Sex-lethal* (*Sxl*) is the sex determination switch in *Drosophila*, and also plays a critical role in
31 germ-line stem cell (GSC) daughter differentiation in *Drosophila melanogaster*. Three female-
32 sterile alleles at *Sxl* in *Drosophila melanogaster* were previously shown to genetically interact to
33 varying degrees with the maternally inherited endosymbiont *Wolbachia pipientis*. Given this
34 genetic interaction and *W. pipientis*' ability to manipulate reproduction in *Drosophila*, we carried
35 out a careful study of both the population genetics (within four *Drosophila* species) and
36 molecular evolutionary analysis (across 20 *Drosophila* species) of *Sxl*. Consistent with earlier
37 studies, we find that selective constraint has played a prominent role in *Sxl*'s molecular evolution
38 within *Drosophila*, but we also observe patterns that suggest both episodic bursts of protein
39 evolution and recent positive selection at *Sxl*. The episodic nature of *Sxl*'s protein evolution is
40 discussed in light of its genetic interaction with *W. pipientis*.

41

42 **Introduction**

43 Reproductive success is a key fitness trait governed by a plethora of gene regulatory networks.
44 Despite the presumed functional constraint for such critical genes, many of them have been
45 shown to be evolving rapidly due to positive selection at the amino acid level (Clark et al. 2006;
46 Chapman 2008; Wong and Rundle 2013; Popovic et al. 2014). For some sets of reproductive loci
47 this observation makes intuitive sense, including loci involved in species-specific gamete
48 recognition and those involved in coevolutionary conflict between the sexes. Interestingly, non-
49 neutral patterns of amino acid evolution have also been detected at loci involved in the

50 differentiation of germ-line stem cells (GSCs) (Civetta et al. 2006; Bauer DuMont et al. 2007;
51 Langley et al. 2012; Pool et al. 2012; Choi and Aquadro 2015; Flores et al. 2015a). The temporal
52 and spatial expression of these GSC regulating loci does not coincide with that expected for
53 genes influenced by sperm competition, sexual selection, inbreeding avoidance and gamete
54 recognition (reviewed in Clark et al. 2006).

55 Several GSC regulating loci have functions outside the germline (e.g., Bell et al. 1988; Yi et al.
56 2008; Saito et al. 2010; Le Thomas et al. 2013) so their signatures of positive selection could be
57 due to non-gametogenic functions. For many others it is possible that the positive selection is
58 acting directly on gametogenic functions. We have previously hypothesized that interactions
59 with maternally transmitted endosymbionts could be a gametogenesis-specific driver of positive
60 selection (Bauer DuMont et al. 2007; Flores et al. 2015a and 2015b). One such endosymbiont,
61 *Wolbachia pipientis*, infects an estimated 66% of arthropod species (Hilgenboecker et al. 2008)
62 and is an obligate maternally transmitted endosymbiont that has been shown to manipulate host
63 reproduction systems (Werren et al. 2008). *Wolbachia pipientis* infection has also been shown to
64 have beneficial consequences. For example, infection in *Drosophila melanogaster* conveys
65 greater resistance to viruses (Hedges et al. 2008; Teixeira et al. 2008; Chrostek et al. 2013) and
66 has been shown to increase female fecundity on low and high iron diets relative to uninfected
67 flies (Brownlie et al. 2009).

68 Here we explore the phylogenetic patterns of positive selection for amino acid diversification at
69 *Sxl*, the sex determination master switch in *Drosophila* development (Bell et al. 1988) that also
70 plays a critical role, along with *bag of marbles* (*bam*), in the maturation of cystoblasts during
71 oogenesis (Chau et al. 2009). Proper *bam* function is essential for the start of cystoblast
72 differentiation (McKearin and Spradling 1990). Epistasis experiments revealed that *bam* requires

73 *Sxl* activity for proper germline stem cell daughter differentiation and the presence of both
74 proteins is proposed to be responsible for regulating sex-specific gametogenesis (Chau et al.
75 2009; Chau et al. 2012; Shapiro-Kulnane et al. 2015).

76 In addition to genetically interacting with one another, hypomorphic alleles of both proteins have
77 been found to genetically interact with *W. pipientis* infection. Starr and Cline (2002) reported
78 that *W. pipientis* infection partially rescues the female-sterile phenotype caused by three *Sxl*
79 alleles, though to differing degrees for each allele. Similarly, *W. pipientis* infection can mitigate
80 the mutant phenotype of a hypomorphic allelic combination at the *bam* locus (Flores et al.
81 2015b). Starr and Cline (2002) also reported that *W. pipientis* infection did not rescue mutants in
82 three other genes with ovarian tumor phenotypes (*snf⁶²¹*, *otu¹¹*, and *mei-P26^{ts1}*). These results
83 suggest a specific, possibly physical, interaction of *W. pipientis* or its gene products with the *Sxl*
84 and *bam* genes or gene products.

85 Using methods that consider both polymorphism and divergence we and others have previously
86 shown that in *D. melanogaster* and *D. simulans*, the *bam* locus is evolving rapidly at the amino
87 acid level due to positive selection (Civetta et al. 2006; Bauer DuMont et al. 2007). Phylogenetic
88 methods for detecting selection across a number of *Drosophila* species did not detect evidence of
89 positive selection, suggesting the selection pressures acting on *bam* are episodic. On the other
90 hand, recent studies of the molecular evolution of *Sxl* have focused on the evolution of *Sxl*'s role
91 as the sex determination master switch in *Drosophila* development (Mullon et al. 2012; Zhang et
92 al. 2013). Mullon et al. (2012) used the maximum likelihood based phylogenetic methods
93 implemented in PAML (Yang 2007) within and between 3 families of Diptera (Drosophilidae,
94 Tephritidae and Muscidae), and detected both a relaxation of functional constraint and positive
95 selection acting across amino acid sites along the lineage leading to *Drosophila*. However, they

96 observe no evidence of positive selection among the 12 *Drosophila* species suggesting that the
97 positive selection they observed was associated with the acquisition of the sex determination
98 function in the common ancestor of *Drosophila* species.

99 Here we test for evidence departures from selective neutrality at *Sxl* within the genus *Drosophila*
100 specifically by incorporating sequence polymorphism and divergence data, and test for
101 departures consistent with lineage specific positive selection. First, we obtained high-quality
102 Sanger sequencing polymorphism data at *Sxl* in the following species: *Drosophila melanogaster*,
103 *D. simulans*, *D. ananassae* and *D. pseudoobscura*. Except for *D. pseudoobscura*, these species
104 are currently infected with *W. pipientis* (Mateos et al. 2006). Second, we searched 8 additional
105 sequenced *Drosophila* genomes to annotate *Sxl* orthologs, bringing the total number of
106 *Drosophila Sxl* orthologous sequences to 20 for phylogenetic analysis. The inclusion of
107 polymorphism data and additional *Drosophila* species allowed us greater power to detect
108 potential signatures of positive selection at *Sxl* within the genus *Drosophila*.

109

110 **Materials and Methods**

111 All DNA was extracted using the Qiagen puregene kit A DNA isolation kits (Qiagen). For *D.*
112 *melanogaster*, we sequenced *Sxl* in 20 extracted X-chromosome lines from Zambia Africa,
113 which was made through a series of crosses using the balancer X chromosome line Fm7a and
114 isofemale Zambia lines (Pool et al. 2012). For *D. simulans*, we used 10 lines from a Madagascar
115 population sample collected in 1998 by J.W.O. Ballard. For *D. pseudoobscura* we included 29
116 lines collected from Mesa Verde National Park, Kaibab National Park and the Bosque Del
117 Apache National Wildlife Refuge provided by Stephen Schaeffer. Our line of *D. miranda* was

118 provided by Doris Bachtrog. Finally, we also surveyed *Sxl* variation across 12 lines of *D.*
119 *ananassae* collected from Bangkok, Thailand provided by Wolfgang Stephan, and a single line
120 of *D. bipectinata* from the UCSD Species Stock Center (stock 0000-1029.01). We also
121 sequenced a single *Sxl* allele from *D. guanche* (obtained from the *Drosophila* Species Stock
122 Center at University of California San Diego; stock number 14011-0095.01), and *D. atripex*
123 (provided by Artyom Kopp), which were used for analyses that require divergence for the *D.*
124 *pseudoobscura* and *D. ananassae* datasets respectively. Our sequences are available via Genbank
125 accession numbers KT935592 - KT935663. All other sequences used in our analyses were those
126 from *Drosophila* 12 Genomes Consortium et al. (2007) or Chen et al. (2014) and downloaded
127 from Flybase.

128 While there are multiple *Sxl* transcripts, due to alternative splicing, most of them include exons 5
129 through 8 in *D. melanogaster* according to Flybase genome annotations (St Pierre et al. 2014).
130 We PCR amplified and sequenced exons 5 through 8. We used Promega GoTaq for
131 amplification following their standard protocol. Sanger sequencing was performed using the
132 PCR primers and internal sequencing primers (primer sequences available upon request) through
133 the Cornell Institute of Biotechnology Genomics Facility ([https://www.biotech.cornell.edu/core-](https://www.biotech.cornell.edu/core-facilities-brc/facilities/genomics-facility)
134 [facilities-brc/facilities/genomics-facility](https://www.biotech.cornell.edu/core-facilities-brc/facilities/genomics-facility)).

135 The following triplets of species were independently aligned using MegAlign (DNASTAR Inc.,
136 Madison WI): *D. melanogaster*, *D. simulans*, *D. yakuba*; *D. ananassae*, *D. atripex*, *D.*
137 *bipectinata*; and *D. pseudoobscura*, *D. miranda*, *D. guanche*. Gaps within coding regions were
138 manually adjusted to ensure the sequences remained in-frame. Population genetic analyses were
139 performed using DNAsp 5.10.1 (Librado and Rozas 2009). When using coalescent simulations to
140 obtain the *P*-values of site frequency based tests we incorporated recombination rate estimates

141 (following Przeworski et al. 2001). For the *D. melanogaster* and *D. simulans* datasets, we used
142 the *Drosophila melanogaster* Recombination Rate Calculator (Comeron et al. 2012) estimate of r
143 = 3.34×10^{-8} recombinants per base-pair per generation for the *Sxl* region of the X chromosome.
144 This translates to an estimated $R = 152$ for the 1500 base pair region sequenced ($R = 3N_e r$ since
145 *Sxl* is X-linked and assuming a population size of 1.0×10^6). For *D. ananassae* and *D.*
146 *pseudoobscura* the values of R were estimated from the polymorphism data directly using
147 DNAsp 5.10.1 (Librado and Rozas 2009) as $R=194$ and $R=6$, respectively.

148 For *D. melanogaster* we also incorporated demography into our neutral simulations when
149 obtaining our significance cut-offs. These simulations were done using the program msABC
150 (Pavlidis et al. 2010). There is growing evidence that African populations of *D. melanogaster*
151 have experienced changes in effective population size over time (Glinka et al 2003; Li and
152 Stephan 2006; Hutter et al 2007; Haddrill et al 2008; Duchon et al 2013; Singh et al 2013).
153 However, given the large effective population size of these species and signatures of a high rate
154 of adaptation (e.g., Begun et al 2007; Langley et al 2012), inferring demographic parameters is
155 challenging. Because of this we simulated three different scenarios: standard neutral equilibrium
156 model, standard neutral with exponential growth as estimated by Hutter et al (2007), and
157 standard neutral with a 3 phase bottleneck as estimate by Duchon et al. (2013). We supplied
158 msABC with uniform prior distributions for theta and all demographic parameters. The prior
159 distribution for theta for *D. melanogaster* was obtained from Pool et al. (2012) and ranged
160 between 0.006 and 0.009 per site. The resulting P -values are the proportion of simulated datasets
161 that were less than (for negative statistics) or greater than (for positive statistics) our observed
162 test statistic for *Sxl*. The P -values were adjusted for multiple testing following the Bonferroni
163 method.

164 The McDonald-Kreitman test (MKT) was done manually following the method's original
165 implementation (McDonald and Kreitman 1991) by combining polymorphism from multiple
166 species if it was available. If a position in the alignment had more than one nucleotide
167 segregating within a species' population sample, it was labeled as polymorphic. Divergent sites
168 were those for which all alleles from one species differed from all the alleles of the other two
169 species.

170 To test for evidence of departures from neutrality for synonymous sites at *Sxl*, we used the
171 method of Bauer DuMont et al. (2005). This method looks for differences in the rates of
172 preferred and unpreferred codon substitutions per site in a manner similar to a dN/dS comparison
173 (Nei and Gojobori 1986). Statistical significance is assessed by a 2x2 contingency table
174 comparison.

175 In order to test for evidence of departures from selective neutrality in rates and patterns of
176 sequence evolution at *Sxl* across a broader group of *Drosophila* species, we first retrieved the *Sxl*
177 gene region sequences from FlyBase for the following 20 *Drosophila* species: *D. melanogaster*,
178 *D. sechellia*, *D. yakuba*, *D. erecta*, *D. eugracilis*, *D. ficusphila*, *D. rhopaloa*, *D. elegans*, *D.*
179 *takahashii*, *D. biarmipes*, *D. kikkawae*, *D. bipectinata*, *D. ananassae*, *D. miranda*, *D.*
180 *pseudoobsura*, *D. persimilis*, *D. willistoni*, *D. mojavensis*, *D. virilis*, and *D. grimshawi*. *Sxl* is an
181 alternatively spliced locus. To ensure we are analyzing orthologous exons, we first made
182 alignments of the entire gene region (introns and exons) using the web-based versions of the
183 alignment programs Muscle (Edgar 2004) (<http://www.ebi.ac.uk/Tools/msa/muscle/>). We then
184 used the annotated exons of *D. melanogaster* as a guide to identify orthologous coding sequences
185 (CDS) for *Sxl*'s Isoform L (6 exons - the female specific splice variant; Bell et al. 1988) from

186 each aligned sequence. This isoform contains the poly-proline region of the Sxl protein where
187 the female-sterile *Sxl* variants are located (Starr and Cline 2002).

188 We estimated a *Sxl* gene tree across these 20 *Drosophila* species using Mega 5.1 (Tamura et al.
189 2011). A maximum likelihood tree was estimated using all nucleotide sites, default parameters,
190 and the GTR substitution model with gamma distributed site variation. To test for selection
191 across the estimated *Sxl* tree we used Hyphy (Kosakovsky Pond et al. 2005) run online using the
192 DataMonkey website (<http://www.datamonkey.org/>). A model selection procedure was
193 conducted to determine that the best nucleotide substitution model for the data was TrN93
194 (Tamura and Nei 1993) which was used in all subsequent analyzes. We ran GARD (Kosakovsky
195 Pond et al. 2006a; Kosakovsky Pond et al. 2006b) to look for evidence of recombination across
196 these species at *Sxl* (that would reflect incomplete ancestral polymorphism sorting) using a
197 general discrete site-to-site rate variation and 3 rate classes. To detect evidence of purifying
198 and/or diversifying selection across sites and lineages we used the following Hyphy programs:
199 BranchRel, GAbbranch, FUBAR and MEME.

200 **Results**

201 We surveyed DNA variability at the population level at *Sxl* for four species of *Drosophila* (*D.*
202 *melanogaster*, *D. simulans*, *D. ananassae* and *D. pseudoobscura*), the first three of which have
203 evidence of current *W. pipientis* infection (Mateos et al. 2006). Even though 37 lines of *D.*
204 *pseudoobscura* were surveyed for infection, this species appears to be currently uninfected with
205 *W. pipientis* (Mateos et al. 2006).

206 We find that Sxl is a very conserved protein with little to no nonsynonymous polymorphism or
207 divergence within or between the 20 *Drosophila* species surveyed (Table 1). In addition, levels

208 of synonymous polymorphism and divergence between *D. melanogaster* and *D. simulans* are
209 below the average reported by Andolfatto (2005). The same is true for synonymous variation
210 observed within and between *D. pseudoobscura* and *D. miranda* as compared to that reported by
211 Haddrill et al. (2010). While synonymous polymorphism is slightly lower at *Sxl* within *D.*
212 *ananassae* the level of divergence between *D. ananassae* and *D. atripex* is similar to values
213 previously reported (Grath et al. 2009; Choi and Aquadro 2014).

214 Previous studies have reported a skew in the Site Frequency Spectrum (SFS) toward rare alleles,
215 as illustrated by a general negative Tajima *D* (Tajima 1989) test statistic, in all of the species
216 included in our study (Kliman et al. 2000; Machado et al. 2002; Das et al. 2004; Andolfatto
217 2007; Grath et al. 2009; Haddrill et al. 2010; Jensen and Bachtrog 2011;). At *Sxl*, we observe a
218 negative Tajima *D* for all species, except *D. ananassae*. Tajima *D* in *D. simulans* rejects the
219 hypothesis of a SFS at equilibrium with 67% (33/49) of the polymorphisms being singletons.
220 These singletons are evenly distributed across synonymous (seven singleton/10 total) and intron
221 (26 singleton/39 total) sites. At *Sxl*, Fay and Wu's *H* is negative in all species but *D. simulans*.
222 Fay and Wu's *H* statistic is significantly negative in *D. melanogaster* even when considering two
223 different demographic scenarios estimated for African populations of this species (Hutter et al.
224 2007; Duchon et al. 2013).

225 The McDonald-Kreitman Test (MKT; McDonald and Kreitman 1991) is used to detect
226 departures from the neutral expectation that synonymous and nonsynonymous variants will have
227 similar ratios of within to between species variation. A rejection in the direction of an excess of
228 nonsynonymous divergence is typically interpreted as evidence of repeated amino acid
229 substitutions due to positive selection. As seen in Table 1 there are few nonsynonymous changes
230 at *Sxl*. The MKT does not reject neutral expectations when *D. ananassae* and *D pseudoobscura*

231 polymorphism is compared to divergence to *D. atripex* and *D. miranda*, respectively (Table 2).
232 We observe no amino acid differences within or between *D. melanogaster* and *D. simulans*.
233 However, when considering the *Sxl* sequence from two additional and closely related *Drosophila*
234 species, *D. yakuba* and *D. erecta*, we observe seven amino acid substitutions along the lineage
235 leading to *D. melanogaster* and *D. simulans*. Given this observation, we chose to apply the MKT
236 to these species in a manner similar to its first implementation (McDonald and Kreitman 1991)
237 by combining the polymorphism from *D. melanogaster* and *D. simulans*. Our divergent changes
238 in this MKT comparison are all differences since the most recent common ancestor of the *D.*
239 *melanogaster* and *D. simulans* lineage rooted by the *D. yakuba* and *D. erecta* lineages. This
240 combined polymorphism MKT rejects neutral expectations, in the direction suggestive of an
241 excess of amino acid substitutions.

242 The significant MKT for *D. melanogaster/D. simulans* could be due to selective fixation of
243 amino acid differences or to selection acting on synonymous changes. While synonymous sites
244 have traditionally been assumed as the neutral yardstick of molecular evolution, there is evidence
245 that this assumption may be invalid, for at least some genes in *Drosophila* (e.g., (Bauer DuMont
246 et al. 2004; Bauer DuMont et al. 2009; Poh et al. 2012; Lawrie et al. 2013). We looked for
247 evidence of selection acting on synonymous changes at *Sxl* using a per site counting method
248 (CF-test) similar to a dN/dS comparison (Bauer DuMont et al. 2004), except in this test we are
249 comparing the number of changes toward unpreferred or preferred codons per the number of
250 unpreferred and preferred “sites”. Along the *D. simulans*, *D. ananassae* and *D. pseudoobscura*
251 lineages we observe a significant departure from neutrality in the direction suggesting a selective
252 advantage of mutations toward preferred codons at *Sxl* (Table 3). The test did not reject in *D.*
253 *melanogaster*.

254 Selection acting on synonymous sites can lead to false positive MKT results by elevating the
255 synonymous polymorphism cell of the 2x2 table, due to segregating slightly deleterious
256 synonymous variants. One method proposed to mitigate the effects of selective constraint on the
257 MKT is to remove low frequency polymorphisms from the analysis (Fay et al. 2002). Given the
258 CF-test results for *D. simulans*, we also carried out the MKT removing two derived unpreferred
259 polymorphic sites at low frequency (singletons – with a frequency of 10% in sample) in *D.*
260 *simulans*. The *D. melanogaster/D. simulans* MKT remains significant (P -value = 0.031).

261 To further assess the molecular evolution at *Sxl* we made a coding sequence (CDS) alignment
262 using the program Muscle for *Sxl*'s Isoform L (6 exons - the female specific splice variant; Bell
263 et al 1988) for the following 20 *Drosophila* species: *D. melanogaster*, *D. sechellia*, *D. yakuba*, *D.*
264 *erecta*, *D. eugracilis*, *D. ficusphila*, *D. rhopaloa*, *D. elegans*, *D. takahashii*, *D. biarmipes*, *D.*
265 *kikkawae*, *D. bipectinata*, *D. ananassae*, *D. miranda*, *D. pseudoobsura*, *D. persimilis*, *D.*
266 *willistoni*, *D. mojavensis*, *D. virilis*, and *D. grimshawi*. We observe 86 amino acid substitutions
267 at 53 codon positions among these species at *Sxl*. Roughly 38% of the codons that have
268 experienced an amino acid substitution have been hit multiple times (20/53; note that some
269 multiply hit amino acid positions had more than three different amino acids segregating among
270 the species). The conservation of the RNA binding domain of the *Sxl* protein has been previously
271 noted (Zhang, Klein, Nei 2013). In agreement, 84 out of the 86 amino acid changes occurred
272 outside the RNA binding domain region, between codons 1 - 136 and 304 – 373 in our alignment
273 (Figures 1 and 2). We will call these non-RNA binding regions of the *Sxl* protein the N-terminal
274 and C-terminal regions, respectively. The amino acid substitutions at *Sxl* have not occurred
275 equally between the N-terminal and C-terminal regions after taking into account their differences
276 in total codon length. We observe significantly more amino acid substitutions in the C-terminal

277 region (24 codons with an amino acid substitution out of 137 total codons in N-terminal region
278 versus 27 substituted codons out of 70 total codons in C-terminal region; 2x2 table chi-square =
279 11.1, P -value = 0.001). The two regions have experienced a similar proportion of multiple hit
280 codons (10 multiple hit codons out of 24 total codons with an amino acid substitution in N-
281 terminal region versus 10 out of 27 such codons in C-terminal region).

282 To determine if the apparent heterogeneity in amino acid substitution at *Sxl* is the result of
283 positive selection, we analyzed the data using the phylogeny-based Hyphy method (Kosakovsky
284 Pond et al. 2005; Delpont et al. 2010). Phylogenetic incongruence between species along a gene
285 sequence, due to sorting of ancestral polymorphisms, can have adverse effects on phylogenetic
286 based inferences of positive selection (Wong et al. 2007). We first performed the Hyphy GARD
287 method (Kosakovsky Pond et al. 2006b), which is designed to detect such incongruences. None
288 were detected at a P -value cutoff of 0.10. Therefore, we used a maximum-likelihood gene tree,
289 made from 3rd codon positions, in subsequent analyzes.

290 We applied the following methods of the Hyphy package to our data: GA-branch (Kosakovsky
291 Pond and Frost 2005), BranchRel (Kosakovsky Pond et al. 2011), MEME (Murrell et al. 2012)
292 and FUBAR (Murrell et al. 2013). GABranch detects significant heterogeneity across the *Sxl*
293 phylogeny in the rate of nonsynonymous compared to synonymous evolution (the dN/dS ratio).
294 The best fitting model includes 3 rate classes, yet the dN/dS for the highest class is only 0.115
295 across the *Sxl* locus. The posterior probabilities suggest that the following branches are within
296 the highest dN/dS rate class and that they are evolving at a significantly different rate than other
297 branches in the tree: the branch leading to *D. melanogaster* and *D. sechellia*, the branch leading
298 to *D. elegans* and *D. rhopaloa*, the branch leading to the melanogaster species group, and the *D.*
299 *kikkawei* lineage (Figure 2).

300 The BranchRel method pools information across sites to estimate selection parameters along
301 branches. The method reports the proportion of codons along each lineage that have evolved
302 under three selection regimes: negative selection, neutral/nearly neutral or episodic positive
303 selection. BranchRel confirmed ubiquitous evidence of amino acid constraint (negative
304 selection) across the *Sxl* phylogeny. After multiple testing corrections, no lineage has significant
305 evidence of positive diversifying selection. Similar results were obtained using MEME.

306 FUBAR is used to detect selection pressure acting on individual codons. The strength of this
307 method is that it does not restrict the parameter space for which nonsynonymous and
308 synonymous rates are drawn from during the maximum likelihood process. FUBAR does not
309 detect any sites under diversifying selection with a posterior probability greater than 0.90.
310 However, it does detect 258 codons under negative selection with a posterior probability greater
311 than 0.90, which is roughly 70% of the protein. Just over half of these negatively selected sites
312 (140) are located within the RNA binding domain. We observe no significant difference in the
313 number of negatively selected sites between the N-terminal and C-terminal regions of *Sxl*
314 relative to their respective lengths (N-terminal: 77 negative selective sites in 137 codons versus
315 C-terminal: 41 negative selected sites in 70 codons; 2x2 chi-square = 0.023, *P*-value = 0.879).

316 Interestingly, the seven amino acid differences observed on the lineage leading to the ancestor of
317 the *D. melanogaster* and *D. simulans* species group (with which we observe the significant
318 MKT), cluster with the location of the mutations previously shown to genetically interact with
319 *W. pipientis* (Figure 2).

320

321 **Discussion**

322 In this study we use population and phylogenetic based methods to examine the molecular
323 population genetics and evolution of the *Sxl* locus within the genus *Drosophila* for which *Sxl* is
324 the master switch in sex-determination. We were motivated by the observation of a genetic
325 interaction between *W. pipientis* infection and some mutant alleles at *Sxl* in *D. melanogaster*
326 (Starr and Cline 2002; Sun and Cline 2009). It is not known if this interaction is due to a
327 ubiquitous effect of *W. pipientis* on overall egg production, or if it is due to a direct interaction
328 between the endosymbiont and the *Sxl* locus or protein product.

329 Considering polymorphism data alone, we detect patterns consistent with a recent selective
330 sweep in both *D. simulans* and *D. melanogaster* with the Tajima *D* and Fay and Wu *H* tests,
331 respectively. These significant skews in the frequency spectrum in these species could be due to
332 a variety of evolutionary forces including demography, a selective sweep associated with the
333 fixation of a linked positively selected mutation, or segregating weakly deleterious mutations. *D.*
334 *simulans* is thought to have experienced a recent population expansion resulting in a general
335 tendency for loci in this species to have negative Tajima *D* test statistics (Kliman et al. 2000).
336 The significant Fay and Wu's *H* test in *D. melanogaster* remains significant even when
337 demography is incorporated into the null distribution of the test, suggesting that in this species
338 we are detecting a recent selective sweep. These signatures of positive selection would not be
339 due to an amino acid fixation, given that there are no amino acid differences between *D.*
340 *melanogaster* and *D. simulans* at *Sxl*.

341 *Sxl*'s long-term evolution within *Drosophila* largely reflects strong conservation of protein
342 sequence, as noted previously by Mullon et al. (2012). We detect the action of negative selection
343 both along lineages and at specific codons. The FUBAR method estimates that 70% of the *Sxl*
344 codons are selectively constrained, suggesting that negative selection has had a pervasive effect

345 on *Sxl*'s molecular evolution. However, we do observe amino acid differences across these
346 *Drosophila* species and the pattern of these substitutions is heterogeneous. For example, GA-
347 branch method detects significant variation in the dN/dS ratio across the *Drosophila* species
348 included in this analysis. This heterogeneity at *Sxl* could be due to sporadic relaxations of the
349 negative selection that dominates *Sxl*'s molecular evolution or due to sporadic bursts of positive
350 selection.

351 The Hyphy methods used to detect recurring or episodic positive selection fail to do so after
352 multiple testing corrections. However, we note that the pervasive negative selection observed at
353 *Sxl* could confound these methods, especially if the positive selection is weak or if only a few
354 sites are affected (Kosakovsky Pond et al. 2011). Our current data shows no evidence of long-
355 term or recent positive selection at *Sxl* along the *D. ananassae* and *D pseudoobscura* lineages. In
356 contrast, there are weak signatures of both types of selection in *D. melanogaster* and *D.*
357 *simulans*, so we focus on the molecular evolution of these species and the lineage leading to their
358 common ancestor.

359 The lineage leading to *D. melanogaster* and *D. simulans* shows a decoupling of synonymous and
360 nonsynonymous evolution with the MKT in the direction of an excess of nonsynonymous
361 divergence. This result is due to seven amino acid substitutions on the lineage leading to the *D.*
362 *melanogaster/D. simulans* clade. These nonsynonymous substitutions cluster with the locations
363 of the *Sxl* alleles that genetically interact with *W. pipientis*. In addition, this region of the *Sxl*
364 protein (the C-terminal non-RNA binding region) has experienced significantly more amino acid
365 substitutions than the N-terminal region. This elevation in amino acid substitutions does not
366 appear to be due to a simple relaxation of constraint as we observe no difference between the C-
367 and N-terminal regions in the number of codons predicted to be experiencing negative selection.

368 These results are suggestive of positive selection being at least partially responsible for the
369 fixations of these seven *D. melanogaster/D. simulans* amino acid substitutions. However, there
370 are other possible explanations for these results such as synonymous site evolution and/or
371 changes in effective population. The seven amino acid fixations could be due to the fixation of
372 slightly deleterious mutations if the effective population size was smaller on the branch leading
373 to the *D. melanogaster/D. simulans*. However, the relaxation of constraint is expected to affect
374 both synonymous and nonsynonymous substitutions. We assume the ancestral state in codon
375 preference is toward preferred synonymous codons at *Sxl*, given the results of the CF-test. If
376 relaxation of constraint were responsible for the burst of amino acid fixations on the *D.*
377 *melanogaster/D. simulans* lineage we may also expect a burst of derived unpreferred
378 substitutions, but we do not observe this. We observe an equal number of preferred and
379 unpreferred changes on this lineage (data not shown).

380 Infection dynamics of *W. pipientis* appear to be sporadic and variable both between and within
381 lineages, with uninfected species interspersed with infected species throughout the phylogeny
382 (e.g., Mateos et al. 2006) consistent with multiple losses or gains of infection. The resulting
383 uncertainty in the infection history of species unfortunately prevents us from reliably testing for
384 correlations between *W. pipientis* infection status and burst of positive selection; there are many
385 factors that could weaken our ability to detect an association.

386 In this study we present data revealing both similarities and difference between the molecular
387 evolution at *bam* and *Sxl*, two loci that genetically interact with *W. pipientis* infection. For both
388 loci the fixation of amino acid variants appears to be heterogeneous, potentially weakening
389 phylogenetic methods to detect positive selection or associations with character states. The MKT
390 does reject neutrality for both loci in a manner suggestive of an acceleration of nonsynonymous

391 fixations. However, the extent of amino acid differences is very different between these loci with
392 there being 59 fixed amino acid substitutions between *D. melanogaster* and *D. simulans* at *bam*
393 and none at *Sxl*.

394 Our results do not allow us to draw strong conclusions regarding the role of positive selection on
395 the molecular evolution of the *Sxl* locus within *Drosophila*, but it also does not allow us to
396 discount the influence of both long term (MKT) or recent (Tajima D and Fay and Wu H tests)
397 positive selection in *D. melanogaster* and *D. simulans*. Current data and methodology also do not
398 allow us to make a direct connection between *W. pipientis* infection and selective pressures
399 acting on *Sxl* or *bam*. So, it remains open whether the genetic interaction between mutant *Sxl* and
400 *bam* alleles and *W. pipientis* is due to a direct interaction between Sxl and Bam protein and this
401 endosymbiont. Our results do motivate screening for genetic interactions between *W. pipientis*
402 and other mutant alleles at *Sxl* and other GSC loci because the observation that *W. pipientis*
403 rescues some but not other mutations (e.g., Starr and Cline 2002) will help refine candidates for
404 the mechanism(s) by which *W. pipientis* is manipulating *Drosophila* reproduction.

405

406 **Acknowledgements**

407 We would like to thank Jae Young Choi for helpful input regarding data collection and
408 manuscript preparation, and Helen K. Salz for her feedback on the manuscript. Research was
409 supported by National Institute of Health grant number R01GM095793 to C.F.A.

410 **Literature Cited**

411 Andolfatto P. 2005. Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* 437: 1149-
412 1152.

- 413 Andolfatto, P. 2007. Hitchhiking effects of recurrent beneficial amino acid substitutions in the
414 *Drosophila melanogaster* genome. *Genome Res.* 17: 1755-1762.
- 415 Bauer DuMont, V., J. C. Fay, P. P. Calabrese, and C. F. Aquadro, 2004. DNA variability and
416 divergence at the *Notch* locus in *Drosophila melanogaster* and *D. simulans*: A case of
417 accelerated synonymous site divergence. *Genetics* 167: 171-85.
- 418 Bauer Dumont V. L., N. D. Singh, M. H. Wright and C. F. Aquadro, 2009 Locus-specific
419 decoupling of base composition evolution at synonymous sites and introns along the
420 *Drosophila melanogaster* and *Drosophila sechellia* lineages. *Genome Biol Evol.* 1: 67-74.
- 421 Bauer Dumont V.L., H. A. Flores, M. H. Wright, and C. F. Aquadro, 2007 Recurrent positive
422 selection at *bgn*, a key determinant of germ line differentiation, does not appear to be
423 driven by simple coevolution with its partner protein *bam*. *Mol. Biol. Evol.* 24: 182-91.
- 424 Bell L. R., E. M. Maine, P. Schedl, and T. W. Cline, 1988 *Sex-lethal*, a *Drosophila* sex
425 determination switch gene, exhibits sex-specific RNA splicing and sequence similarity to
426 RNA binding proteins. *Cell* 55:1037-46.
- 427 Bennett, G. M., N. A. Pantoja, and P. M. O'Grady, 2012 Diversity and phylogenetic
428 relationships of *Wolbachia* in *Drosophila* and other native Hawaiian insects. *Fly* 6: 273-
429 283.
- 430 Brownlie, J. C., B. N. Cass, M. Riegler, J. J. Witsenburg, I. Iturbe-Ormaetxe, E. A. McGraw, and
431 S. O'Neill, 2009 Evidence for metabolic provisioning by a common invertebrate
432 endosymbiont, *Wolbachia pipientis*, during periods of nutritional stress. *PLoS Pathog* 5:
433 e1000368.
- 434 Chapman, T., 2008 The soup in my fly: Evolution, form and function of seminal fluid proteins.
435 *PLoS Biol.* 6: e179.

- 436 Chau, J., L. S. Kulnane L. S., and H. K. Salz, 2009. *Sex-lethal* facilitates the transition from
437 germline stem cell to committed daughter cell in the *Drosophila* ovary. *Genetics* 182: 121-
438 32.
- 439 Chau, J., L. S. Kulnane, and H. K. Salz, 2012 *Sex-lethal* enables germline stem cell
440 differentiation by down-regulating Nanos protein levels during *Drosophila* oogenesis. *Proc.*
441 *Natl. Acad. Sci. USA* 109: 9465-9470.
- 442 Chen, Z. X., ZD. Sturgill, J. Qu, H. Jiang, S. Park, N. Boley, A. M. Suzuki, et al. 2014.
443 Comparative validation of *D. melanogaster* modENCODE transcriptome annotation.
444 *Genome Res.* 24: 1209-1223.
- 445 Choi, J. Y., and C. F. Aquadro, 2014 The coevolutionary period of *Wolbachia pipientis* infecting
446 *Drosophila ananassae* and its impact on the evolution of the host germline stem cell
447 regulating genes. *Mol. Biol. Evol.* 31: 2457-2471.
- 448 Chrostek E, Marialva MSP, Esteves SS, Weinert LA, Martinez J, Jiggins FM, Teixeira L. 2013.
449 *Wolbachia* variants induce differential protection to viruses in *Drosophila melanogaster*: A
450 phenotypic and phylogenomic analysis. *PLoS Genet* 9: e1003896.
- 451 Civetta, A., S. A. Rajakumar, B. Brouwers, and J. P. Bacik, 2006 Rapid evolution and gene-
452 specific patterns of selection for three genes of spermatogenesis in *Drosophila*. *Mol. Biol.*
453 *Evol.* 23: 655-662.
- 454 Clark, N. L., J. E. Aagaard, and W. J. Swanson, 2006 Evolution of reproductive proteins from
455 animals and plants. *Reproduction* 131: 11-22.
- 456 Comeron, J.M., R. Ratnappan, and S. Bailin, 2012 The many landscapes of recombination in
457 *Drosophila melanogaster*. *PLoS Genet.* 8: e1002905.

- 458 Das, A., S. Mohanty, and W. Stephan, 2004 Inferring the population structure and demography
459 of *Drosophila ananassae* from multilocus data. *Genetics* 168: 1975-1785.
- 460 Delpont, W., A. F. Y. Poon, S. D. W. Frost, and S. L. Kosakovsky Pond, 2010 Datamonkey
461 2010: A suite of phylogenetic analysis tools for evolutionary biology. *Bioinformatics* 26:
462 2455-2457.
- 463 Drosophila 12 Genomes Consortium, A. G. Clark, M. B. Eisen, D. R. Smith, C. M. Bergman, B.
464 Oliver, T. A. Markow, et al. 2007. Evolution of genes and genomes on the *Drosophila*
465 phylogeny. *Nature* 17: 1932-1942
- 466 Duchen, P., D. Zivkovic, S. Hutter, W. Stephan, and S. Laurent, 2013 Demographic inference
467 reveals African and European admixture in the North American *Drosophila melanogaster*
468 population *Genetics* 193: 291-301.
- 469 Edgar, R.C., 2004. MUSCLE: Multiple sequence alignment with high accuracy and high
470 throughput. *Nucleic Acids Res.* 32: 1792-1797.
- 471 Fay, J.C., G. J. Wyckoff, and C. I. Wu, 2001 Positive and negative selection on the human
472 genome. *Genetics* 158: 1227-1234.
- 473 Flores, H. A. F., V. L. Bauer DuMont, A. Fato, D. Hubbard, M. Hijji, D. A. Barbash, and C. F.
474 Aquadro. 2015a. Adaptive evolution of genes involved in the regulation of germline stem
475 cells in *Drosophila melanogaster* and *D. simulans*. *G3 (Bethesda)* 5: 583-592.
- 476 Flores, H. A. F., J. E. Bubnell, C.F. Aquadro, and D. A. Barbash. 2015b. The *Drosophila bag*
477 of *marbles* gene interacts genetically with *Wolbachia* and shows female-specific effects of
478 divergence. *PLoS Genet.* 11: e1005453.
- 479 Grath, S., J. F. Baines, and J. Parsch, 2009 Molecular evolution of sex-biased genes in the
480 *Drosophila ananassae* subgroup. *BMC Evol. Biol.* 9: 291.

- 481 Haddrill, P. R. , L. Loewe, and B. Charlesworth, 2010. Estimating the parameters of selection on
482 nonsynonymous mutations in *Drosophila pseudoobscura* and *D. miranda*. *Genetics* 185:
483 1381-1396.
- 484 Hedges, L. M., J. C. Brownlie, S. J. O'Neill, and K. N. Johnson, 2008. *Wolbachia* and virus
485 protection in insects. *Science* 322: 702.
- 486 Hilgenboecker, K., P. Hammerstein, P. Schlattmann, A. Telschow, and J. H. Werren, 2008 How
487 many species are infected with *Wolbachia*? -- a statistical analysis of current data. *FEMS*
488 *Microbiol. Lett.* 281: 215-220.
- 489 Hutter, S., H. Li, S. Beisswanger, D. De Lorenzo, and W. Stephan, 2007 Distinctly different sex
490 ratios in African and European populations of *Drosophila melanogaster* inferred from
491 chromosomewide single nucleotide polymorphism data. *Genetics* 177: 469-480.
- 492 Jensen, J. D., and D. Bachtrog, 2011 Characterizing the influence of effective population size on
493 the rate of adaptation: Gillespie's Darwin domain. *Genome Biol. Evol.* 3: 687-701.
- 494 Jiggins, F. M., and K. W. Kim, 2005 The evolution of antifungal peptides in *Drosophila*.
495 *Genetics* 171: 1847-1859.
- 496 Kliman, R. M., P. Andolfatto, J. A. Coyne, F. Depaulis, M. Kreitman, A. J. Berry, J. McCarter, J.
497 Wakeley, and J. Hey, 2000 The population genetics of the origin and divergence of the
498 *Drosophila simulans* complex species. *Genetics* 156: 1913-1931.
- 499 Kosakovsky Pond, S. L., and S. D. W. Frost, 2005 A genetic algorithm approach to detecting
500 lineage-specific variation in selection pressure. *Mol. Biol. Evol.* 32: 478-485.
- 501 Kosakovsky Pond, S. L. K., S. D. W. Frost, and S. V. Muse, 2005 HyPhy: Hypothesis testing
502 using phylogenies. *Bioinformatics* 21: 676-679.

- 503 Kosakovsky Pond, S. L., B. Murrell, M. Fourment, S. D. W. Frost, W. Delpont, and K. Scheffler,
504 2011 A random effects branch-site model for detecting episodic diversifying selection. *Mol.*
505 *Biol. Evol.* 28: 3033-3043.
- 506 Kosakovsky Pond, S. L., D. Posada, M. B. Gravenor, C. H. Woelk, and S. D. Frost, 2006a.
507 Automated phylogenetic detection of recombination using a genetic algorithm. *Mol. Biol.*
508 *Evol.* 23: 1891-1901.
- 509 Kosakovsky Pond, S. L., D. Posada, M. B. Gravenor, C. H. Woelk, and S. D. W. Frost, 2006b.
510 GARD: A genetic algorithm for recombination detection. *Bioinformatics* 22: 3096-3098.
- 511 Lawrie, D. S., P. W. Messer, R. Hershberg, and D. A. Petrov, 2013 Strong purifying selection at
512 synonymous sites in *D. melanogaster*. *PLoS Genet.* 9: e1003527.
- 513 Le Thomas, A., A. K. Rogers, A. Webster, G. K. Marinov, S. E. Liao et al., 2013 Piwi induces
514 piRNA-guided transcriptional silencing and establishment of a repressive chromatin state.
515 *Genes Dev.* 27: 390-399.
- 516 Librado, P., and J. Rozas, 2009 DnaSP v5: A software for comprehensive analysis of DNA
517 polymorphism data. *Bioinformatics* 25: 1451-1452.
- 518 Machado, C. A., R. M. Kliman, J. A. Markert, and J. Hey 2002 Inferring the history of speciation
519 from multilocus DNA sequence data: The case of *Drosophila pseudoobscura* and close
520 relatives. *Mol. Biol. Evol.* 19: 472-488.
- 521 Mateos, M., S. J. Castrezana, B. J. Nankivell, A. M. Estes, T. A. Markow *et al*, 2006 Heritable
522 endosymbionts of *Drosophila*. *Genetics* 174: 363-376.
- 523 McDonald, J. H., and M. Kreitman, 1991 Adaptive protein evolution at the *Adh* locus in
524 *Drosophila*. *Nature* 351: 652-654.

- 525 McKearin, D. M., and A. C. Spradling, 1990 *bag-of-marbles*: A *Drosophila* gene required to
526 initiate both male and female gametogenesis. *Genes Dev.* 4: 2242-2251.
- 527 Muller, M.J., N. C. Drebes Dorr, M. Depra, H. J. Schmitz, V. H. Valiati and V. L. da Siva
528 Valente, 2013 Reevaluating the infection status by the *Wolbachia* endosymbiont in
529 *Drosophila* Neotropical species from the *willistoni* subgroup. *Infect. Genet. Evol.* 19: 232-
530 239.
- 531 Mullon, C., A. Pomiankowski, and M. Reuter, 2012 Molecular evolution of *Drosophila Sex-*
532 *lethal* and related sex determining genes. *BMC Evol. Biol.* 12: 5.
- 533 Murrell, B., J. O. Wertheim, S. Moola, T. Weighill, K. Scheffler, and S. L. Kosakovsky Pond,
534 2012 Detecting individual sites subject to episodic diversifying selection. *PLoS Genetics* 8:
535 e1002764.
- 536 Poh, Y. P., C. T. Ting, H. W. Fu, C. H. Langley, and D. J. Begun, 2012 Population genomic
537 analysis of base composition evolution in *Drosophila melanogaster*. *Genome Biol. Evol.* 4:
538 1245-1255.
- 539 Pool, J. E., A. Wong, and C. F. Aquadro, 2006 Finding of male-killing *Spiroplasma* infecting
540 *Drosophila melanogaster* in Africa implies transatlantic migration of this endosymbiont.
541 *Heredity* 97: 27-32.
- 542 Pool, J. E., R. Corbett-Detig, R. P. Sugino, K. A. Stevens, C. M. Cardeno, et al., 2012 Population
543 genomics of sub-saharan *Drosophila melanogaster*: African diversity and
- 544 Popovic, I., P. B. Marko, J. P. Wares, and M. W. Hart, 2014 Selection and demographic history
545 shape the molecular evolution of the gamete compatibility protein bindin in *Pisaster* sea
546 stars. *Ecol. Evol.* 4: 1567-1588.

- 547 Przeworski, M., J. D. Wall, and P. Andolfatto, 2001 Recombination and the frequency spectrum
548 in *Drosophila melanogaster* and *Drosophila simulans*. *Mol. Biol. Evol.* 18: 291-298.
- 549 Puigbo, P., S. Garcia-Vallve, and J. O. McInerney, 2007. TOPD/FMTS: A new software to
550 compare phylogenetic trees. *Bioinformatics* 23: 1556-1558.
- 551 Ravikumar, H., B. M. Prakash, S. Sampathkumar, and H. P. Puttaraju, 2011 Molecular
552 subgrouping of *Wolbachia* and bacteriophage WO infection among some Indian *Drosophila*
553 species. *J. Genet.* 90: 507-510.
- 554 Richardson, M. F., L. A. Weinert, J. J. Welch, R. S. Linheiro, M. M. Magwire, et al., 2012
555 Population genomics of the *Wolbachia* endosymbiont in *Drosophila melanogaster*. *PLoS*
556 *Genet.* 8: e1003129.
- 557 Saito, K., H. Ishizu, M. Komai, H. Kotani, Y. Kawamura, et al., 2010 Roles for the Yb body
558 components Armitage and Yb in primary piRNA biogenesis in *Drosophila*. *Genes Dev.* 24:
559 2493-2498.
- 560 Shapiro-Kulnane, L., A. E. Smolko, and H. K. Salz, 2015 Maintenance of *Drosophila* germline
561 stem cell sexual identity in oogenesis and tumorigenesis. *Development* 143: 1073-1082,
- 562 Starr, D. J., and T. W. Cline, 2002 A host parasite interaction rescues *Drosophila* oogenesis
563 defects. *Nature* 418: 76-79.
- 564 Sun, S., and T. W. Cline, 2009 Effects of *Wolbachia* infection and ovarian tumor mutations on
565 *sex-lethal* germline functioning in *Drosophila*. *Genetics* 181: 1291-1301.
- 566 Tajima, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA
567 polymorphism. *Genetics* 123: 585.-595.
- 568 Tamura, K., and M. Nei, 1993 Estimation of the number of nucleotide substitutions in the control
569 region of mitochondrial DNA in humans and chimpanzees. *Mol. Biol. Evol.* 10: 512-526.

- 570 Tamura, K., D. Peterson, N. Peterson, G. Stecher, M. Nei, and S. Kumar, 2011 MEGA5:
571 Molecular evolutionary genetics analysis using maximum likelihood, evolutionary distance,
572 and maximum parsimony methods. *Mol. Biol. Evol.* 28: 2731-2739.
- 573 Teixeira, L., A. Ferreira, and M. Ashburner, 2008 The bacterial symbiont *Wolbachia* induces
574 resistance to RNA viral infections in *Drosophila melanogaster*. *PLoS Biol.* 6: e1000002.
- 575 Werren, J. H., L. Baldo, and M. E. Clark, 2008 *Wolbachia*: Master manipulators of invertebrate
576 biology. *Nat. Rev. Microbiol.* 6: 741-751.
- 577 Wong, A. and H. Rundle, 2013 Selection on the *Drosophila* seminal fluid protein Acp62F. *Ecol.*
578 *Evol.* 3: 1942-1950.
- 579 Wong, A., J. D. Jensen, J. E. Pool, and C. F. Aquadro, 2007 Phylogenetic incongruence in the
580 *Drosophila melanogaster* species group. *Mol. Phylogenet. Evol.* 43: 1138-1150.
- 581 Yang, Z., 2007 PAML 4: Phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* 24:
582 1586-1591.
- 583 Yi, X., H. I. de Vries, K. Siudeja, A. Rana, W. Lemstra, J. F. Brunsting, et al. 2008 *Stwl* modifies
584 chromatin compaction and is required to maintain DNA integrity in the presence of
585 perturbed DNA replication. *Mol. Biol. Cell* 20: 983-994.
- 586 Zabalou, S., S. Charlat, A. Nirgianaki, D. Lachaise, H. Mercot, and K. Bourtzis, 2004 Natural
587 *Wolbachia* infections in the *Drosophila yakuba* species complex do not induce cytoplasmic
588 incompatibility but fully rescue the *wRI* modification. *Genetics* 167: 827-834.
- 589 Zhang, Z., J. Klein, and M. Nei, 2013 Evolution of the *sex-lethal* gene in insects and origin of the
590 sex-determination system in *Drosophila*. *J. Mol. Evol.* 78: 50-65.
- 591

592

593

Figure legends

594 Figure 1 Schematic of Sxl protein. Green box denotes the location of the RNA binding domain
595 of the Sxl protein. Vertical lines show the locations of all amino acid substitutions at *Sxl* across
596 20 *Drosophila* species with the red lines being those that have occurred specifically on the
597 lineage leading to *D. melanogaster* and *D. sechellia*. Stars at C-terminal end of protein denote
598 the relative location of the mutations that generate the mutant alleles shown to interact
599 genetically with *Wolbachia pipientis*.

600

601 Figure 2 *Sxl* gene tree schematically showing the branches for which amino acid substitutions
602 have occurred. The rectangles denote the Sxl protein with the vertical black lines indicating the
603 location of the amino acid change(s) along that lineage. Hatched box denotes the location of the
604 RNA binding domains of the Sxl protein. Stars at C-terminal end of protein denote the relative
605 location of the three mutations that generate the mutant alleles shown to interact genetically with
606 *W. pipientis*. Species names in bold indicate *Wolbachia* has been detected with the numbers of
607 *Wolbachia* positive lines relative the total number of lines screened given in parenthesis.
608 *Wolbachia* data for all species are from Mateo et al. (2006), as well as for *D. erecta* (Zabalou
609 2004), *D. kikkawai* (Bennett et al. 2012), *D. bipectinata* (Ravikumar et al. 2011), and *D.*
610 *willistoni* (Muller et al. 2013).

Figure 1

Sex Lethal protein

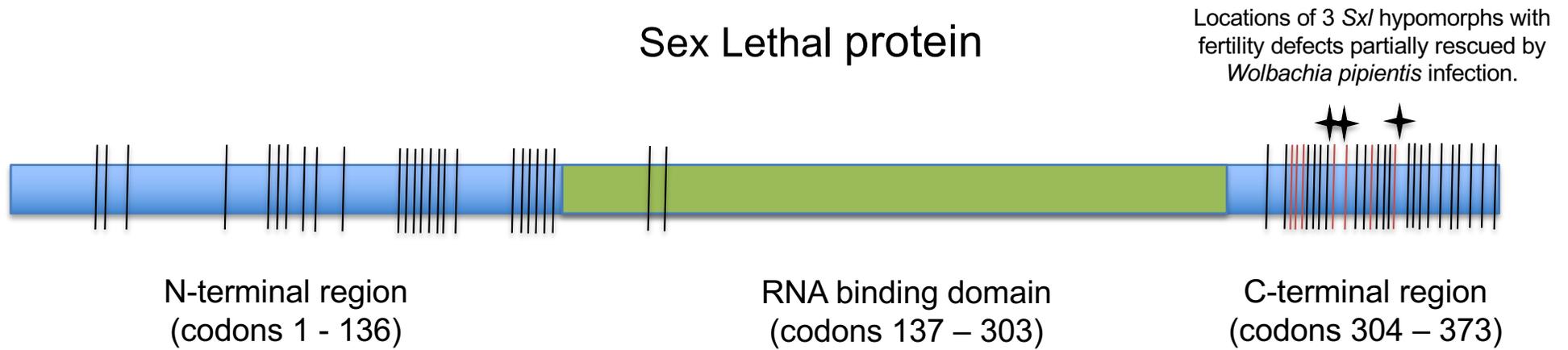


Figure 2

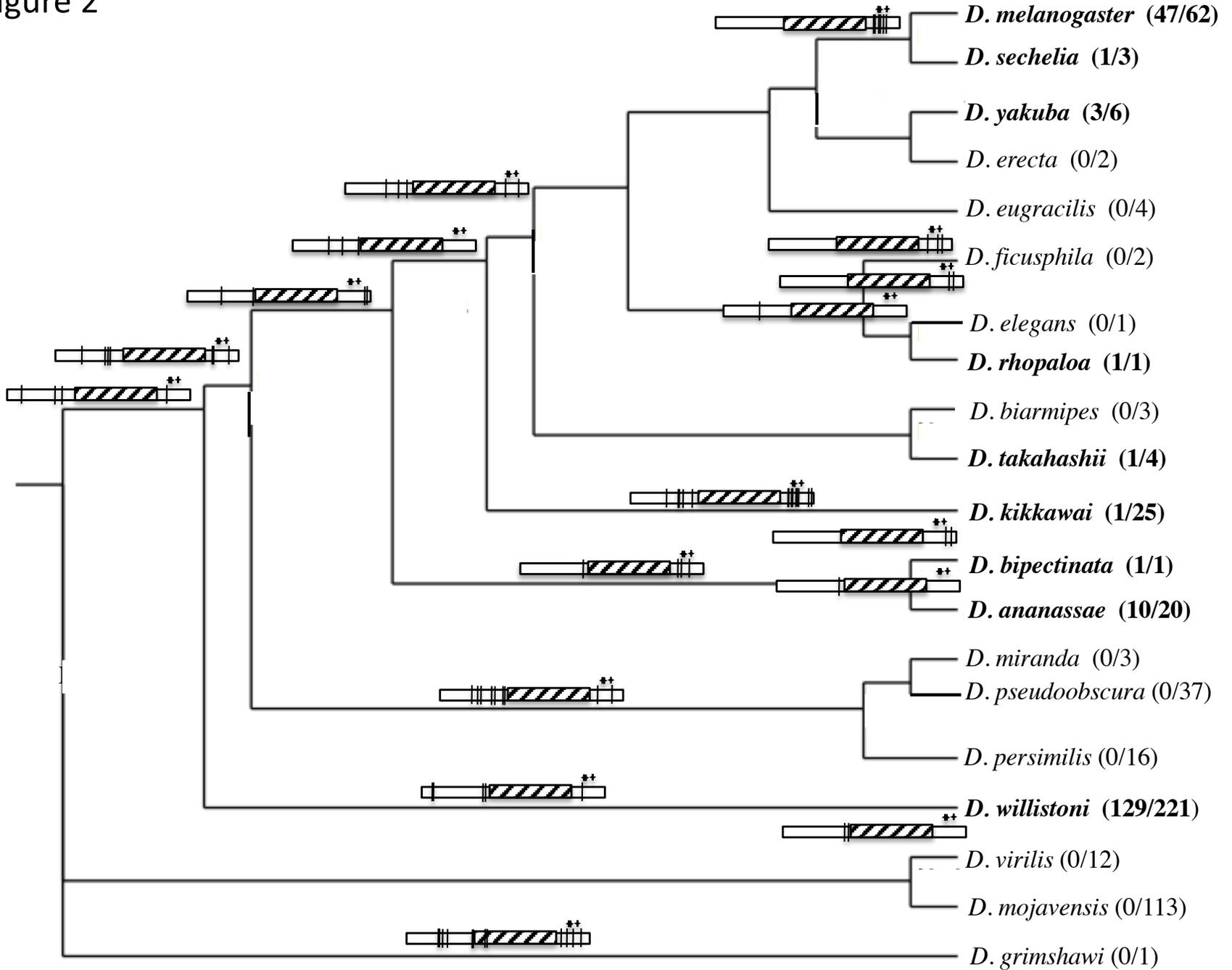


Table 1 Levels of within species variation and divergence between *Drosophila* species and results of tests to detect departures from a neutral site frequency spectrum

	θ	π	Div ^c	Taj-D ^a (<i>P</i> -value ^b)	FW-H ^a (<i>P</i> -value ^b)
<u><i>D. melanogaster</i> (n = 20)</u>					
Synonymous	0.012	0.011	0.069	-0.267	-2.26
Nonsynonymous	0	0	0	(0.214; 0.80; 0.70)	(<0.0002; <0.0002; <0.0002)
Intron	0.019	0.019	0.083		
<u><i>D. simulans</i> (n = 10)</u>					
Synonymous	0.017	0.013	0.069	-1.045	3.47
Nonsynonymous	0	0	0	(0.001)	(0.87)
Intron	0.031	0.025	0.083		
<u><i>D. ananassae</i> (n = 13)</u>					
Synonymous	0.010	0.010	0.209	0.349	-4.333
Nonsynonymous	0.000	0.000	0.003	(0.166)	(0.038)
Intron	0.021	0.025	0.299		
<u><i>D. pseudoobscura</i> (n = 29)</u>					
Synonymous	0.011	0.007	0.073	-1.293	-0.539
Nonsynonymous	0.0004	0.0001	0.002	(0.038)	(0.311)
Intron	0.014	0.008	0.111		

^a test statistic when using all sites in the analysis: Tajima's D (Taj-D) and Fay and Wu's H (FW-H)

^b Proportion of simulated datasets that were equal to or less than (for negative statistics) or equal to or greater than (for positive statistics) our observed test statistic for *Sxl*. These are two-sided tests, for which we do two tests per species resulting in a significant cut-off level of 0.0125 (0.025/2). For *D. melanogaster* the *P*-value was calculated by simulating different demographic scenarios listed in this order in parenthesis: Standard neutral; Bottleneck with exponential growth and a 3 Epoch bottleneck (as described in Materials and Methods). *P*-values significant after multiple testing corrections in bold

^c Uncorrected divergence between the following species pairs: *D. melanogaster* and *D. simulans*; *D. ananassae* and *D. atripex*; and *D. pseudoobscura* and *D. miranda*

Table 2 Results of the McDonald-Kreitman Test at *Sxl* between 3 different sets of *Drosophila* species.

	Synonymous	Nonsynonymous	<i>P</i> -value ^a
Single species polymorphism			
<u><i>D. ananassae/D. atripex</i></u>			
polymorphic	8	0	
Fixed divergent	44	2	1.00
<u><i>D. pseudoobscura/D. miranda</i></u>			
polymorphic	9	1	
Fixed divergent	14	0	0.417
Two species polymorphism combined			
<u><i>(D. melanogaster+D. simulans)/D. yakuba</i></u>			
polymorphic	19	0	
Fixed divergent	19	7	0.015

^a Fisher exact test *P*-value

Table 3 Results of the CF Test at *Sxl* along four *Drosophila* lineages.

	Unpreferred	Preferred	<i>P</i> -value ^a	direction of departure
<i>D. melanogaster</i> lineage				
Fixed substitutions	5	0	0.303	
Sites	150	31		
<i>D. simulans</i> lineage				
Fixed substitutions	1	2	0.022	preferred codons favored
Sites	150	31		
<i>D. ananassae</i> lineage				
Fixed substitutions	6	11	<0.001	preferred codons favored
Sites	174	29		
<i>D. pseudoobscura</i> lineage				
Fixed substitutions	1	3	0.002	preferred codons favored
Sites	154	31		

^a Fisher exact test *P*-value