

1 **A Genomic Reference Panel for *Drosophila serrata***

2

3

4 Adam R. Reddiex¹, Scott, L. Allen¹, Stephen F. Chenoweth^{1*}

5

6

7 1. School of Biological Sciences, The University of Queensland, QLD 4072,
8 Australia.

9

10 * correspondence: s.chenoweth@uq.edu.au

11

12

13

14

15 **Keywords:** *montium*, *gwas*, *population genetics*, *quantitative genetics*, *multi-*
16 *parental population*

17

18

19 **Abstract**

20

21 Here we describe a collection of re-sequenced inbred lines of *Drosophila*
22 *serrata*, sampled from a natural population situated deep within the species
23 endemic distribution in Brisbane, Australia. *D. serrata* is a member of the
24 speciose *montium* group whose members inhabit much of south east Asia and
25 has been well studied for aspects of climatic adaptation, sexual selection, sexual
26 dimorphism, and mate recognition. We sequenced 110 lines that were inbred via
27 17-20 generations of full-sib mating at an average coverage of 23.5x with paired-
28 end Illumina reads. 15,228,692 biallelic SNPs passed quality control after being
29 called using the Joint Genotyper for Inbred Lines (JGIL). Inbreeding was highly
30 effective and the average levels of residual heterozygosity (0.86%) were well
31 below theoretical expectations. As expected, linkage disequilibrium decayed
32 rapidly, with r^2 dropping below 0.1 within 100 base pairs. With the exception of
33 four closely related pairs of lines which may have been due to technical errors,
34 there was no statistical support for population substructure. Consistent with
35 other endemic populations of other *Drosophila* species, preliminary population
36 genetic analyses revealed high nucleotide diversity and, on average, negative
37 Tajima's D values. A preliminary GWAS was performed on a cuticular
38 hydrocarbon trait, 2-MeC₂₈ revealing 4 SNPs passing Bonferroni significance
39 residing in or near genes. One gene *Cht9* may be involved in the transport of
40 CHCs from the site of production (oenocytes) to the cuticle. Our panel will
41 facilitate broader population genomic and quantitative genetic studies of this
42 species and serve as an important complement to existing *D.*
43 *melanogaster* panels that can be used to test for the conservation of genetic
44 architectures across the *Drosophila* genus.

45

46

47

48 **Introduction**

49

50 The availability of whole genome sequence data for *Drosophila* species has
51 greatly facilitated advances in the fields of genetics and evolutionary biology.
52 For example, the sequencing of 12 *Drosophila* genomes (Clark et al. 2007) was
53 instrumental to new discoveries in comparative genomics (Stark et al. 2007;
54 Sturgill et al. 2007; Zhang et al. 2007). The advent of affordable genome
55 sequencing has also allowed population geneticists to characterise genomic
56 variation within and among natural populations, improving our understanding of
57 the complex evolutionary histories of cosmopolitan species such as *D.*
58 *melanogaster* and *D. simulans* (Begun et al. 2007; Lack et al. 2015; Langley et al.
59 2012; Pool et al. 2012). Most recently, multiple panels of re-sequenced inbred *D.*
60 *melanogaster* lines have become available, facilitating the molecular dissection of
61 complex trait variation (Grenier et al. 2015; Huang et al. 2014; King et al. 2012;
62 Mackay et al. 2012). With these populations of reproducible genotypes,
63 researchers have used genome-wide association analysis to identify genetic
64 variants underlying variation in a broad range of traits including physiological
65 traits (Burke et al. 2014; Dembeck et al. 2015; Gerken et al. 2015; Unckless et al.
66 2015; Weber et al. 2012), behaviours (Shorter et al. 2015), recombination rates
67 (Hunter et al. 2016), disease susceptibility (Magwire et al. 2012), and traits
68 related to human health (Harbison et al. 2013; He et al. 2014; King et al. 2014;
69 Kislukhin et al. 2013; Marriage et al. 2014).

70

71 Just as comparative genomic and population genetic studies of adaptation (e.g.
72 Machado et al. 2016; Zhao et al. 2015) have been enhanced through the
73 availability of multi-species genome resources, quantitative genetics may also
74 benefit from the availability of multispecies genome panels. The development of
75 panels of re-sequenced lines for *Drosophila* species beyond *D. melanogaster* may
76 support broader lines of inquiry such as the conservation of genetic
77 architectures among related taxa (Yassin et al. 2016). To this end, we have
78 developed a new genomic resource for *D. serrata*, a member of the *montium*
79 group of species. The *montium* group has long been regarded as a subgroup
80 within the *melanogaster* species group (Lemeunier et al. 1986), but has more

81 recently been considered as a species group of its own (Da Lage et al. 2007;
82 Yassin 2013). Although *montium* contains 98 species (Brake and Bachli 2008)
83 and represents a significant fraction of all known *Drosophila* species, there have
84 been very few genomic investigations of its members. Recently, genomic tools
85 have been developed for *D. serrata* including an expressed sequence tag (EST)
86 library (Frentiu et al. 2009), a physical linkage map (Stocker et al. 2012), and
87 transcriptome-wide gene expression datasets (Allen et al. 2013; Allen et al.
88 2017a; McGuigan et al. 2014). Additionally, an assembled and annotated genome
89 sequence (Allen et al. 2017b) make *D. serrata* only the second species in the
90 *montium* group with a sequenced genome after *D. kikkawai* (NCBI *Drosophila*
91 *kikkawai* Annotation Release 101). Coupled to this, *D. serrata* is one member of
92 the *montium* group that has been extensively studied in the field of evolutionary
93 genetics.

94 Populations of *D. serrata* have been recorded from as far north as Rabaul, Papua
95 New Guinea (4.4°N) (Ayala 1965) to as far south as Woolongong, Australia
96 (34.3°S) (Jenkins and Hoffmann 1999). This broad latitudinal range has made *D.*
97 *serrata* an ideal model for population studies addressing the evolution of species
98 borders (Blows and Hoffmann 1993; Hallas et al. 2002; Magiafoglou et al. 2002;
99 van Heerwaarden et al. 2009) and adaptation along latitudinal clines (Allen et al.
100 2017a; Frentiu and Chenoweth 2010; Kellermann et al. 2009). *D. serrata* has also
101 emerged as a powerful model for the application of quantitative genetic designs
102 to investigate sexual selection (Gosden and Chenoweth 2014; Hine et al. 2002;
103 McGuigan et al. 2011).

104 Here, we report development of a panel of 110 re-sequenced inbred *D.*
105 *serrata* lines that we have called the *Drosophila serrata* Genome Reference Panel
106 (DsGRP). Similar to the DGRP (Mackay et al. 2012), flies were sampled from a
107 single large natural population with the exception that *D. serrata* was sampled
108 from its endemic distribution. In this initial description, we estimate the degree
109 of heterozygosity remaining in the lines after inbreeding, show the degree to
110 which lines are genetically related to one another, estimate genome-wide levels
111 of nucleotide diversity, and describe patterns of linkage disequilibrium. We also
112 demonstrate how this panel of flies can be used to genetically dissect trait

113 variation by performing a genome-wide association analysis on variation in a
114 cuticular hydrocarbon (CHC) trait.

115

116 **Methods:**

117 *Collection and inbreeding*

118 *Drosophila serrata* were collected from a wild population located at Bowman
119 Park, Brisbane Australia (Latitude: -27.45922, Longitude: 152.97768) during
120 October 2011. We established each line from a single, gravid female before
121 applying 20 generations of inbreeding. Inbreeding was carried out each
122 generation by pairing virgin brothers and sisters. 100 inbred lines, out of the
123 initial 239 iso-female lines established, survived the full 20 generations
124 inbreeding and a further 10 lines were established after 17 generations of
125 inbreeding.

126

127 *Sequencing*

128 We sequenced the genomes of 110 inbred lines using 100 base-pair paired-end
129 reads with a 500 base-pair insert on an Illumina Hiseq 2000 sequencing
130 machine. Sequencing and library preparation were carried out by the Beijing
131 Genomics Institute. DNA from each line was isolated from a pool of at least 30
132 virgin female flies using a standard phenol-chloroform extraction method.

133

134 *Quality control and SNP calling*

135 We received reads from the Beijing Genomics Institute for which approximately
136 95% of the bases from each line had a base quality score greater than or equal to
137 20 (Illumina GA Pipeline v1.5). Read quality was also assessed using FastQC
138 v0.11.2 before being mapped to the *Drosophila serrata* reference genome (Allen
139 et al. 2017b) using BWA-mem v0.7.10 (Li 2013) and were realigned around
140 indels using the GATK IndelRealigner v3.2-2 (McKenna et al. 2010). Genotypes
141 for every line were inferred simultaneously using the Joint Genotyper for Inbred
142 Lines (JGIL) v1.6 (Stone 2012). This probabilistic model was especially designed
143 for genotyping large panels of inbred lines or strains and is considered to have
144 high accuracy (Mackay et al. 2012; Stone 2012). Genotype calls with a
145 probability lower than 99% were treated as missing genotypes.

146

147 *Residual heterozygosity*

148 The residual heterozygosity per line was estimated as the genome-wide
149 proportion of sites that remained heterozygous after 17-20 generations of
150 inbreeding, more specifically, we summed all of the genotype call that were
151 heterozygous and expressed this statistic as a percentage of all genotyped sites.
152 In addition, for each site in the genome that differed among the inbred lines, we
153 calculated the percentage of lines that were heterozygous for that site. Site
154 filtering based on minor allele frequency and coverage was not performed for
155 this analysis.

156

157 *Pairwise relatedness between lines*

158 Pairwise relatedness between lines (j and k) was estimated using the --make-
159 grm-inbred command of GCTA v1.24.2 (Yang et al. 2011), which applies the
160 expression:

161

$$162 \quad A_{jk} = \frac{1}{N} \sum_{i=1}^N \frac{(\chi_{ij} - 2p_i)(\chi_{ik} - 2p_i)}{2p_i(1 - p_i)} \quad [1]$$

163

164 where, χ_{ij} is the number of copies of the reference allele for the i^{th} SNP for
165 individual j and p is the population allele frequency. N is the total number SNPs.
166 Only biallelic SNPs with a read depth between 5 and 60 and a minor allele
167 frequency above 5% were used to estimate relatedness.

168

169 *Estimating population substructure*

170

171 To test whether the sample of lines exhibited any underlying substructure, we
172 used the approach of Bryc et al. (2013) which is founded on random matrix
173 theory. Importantly, while the genomic relatedness matrix calculated in GCTA, \mathbf{A}
174 = $\mathbf{W}\mathbf{W}'/N$, is normalised with element $W_{ij} = (\chi_{ij} - 2p_i)/\sqrt{2p_i(1 - p_i)}$, where χ_{ij}
175 is the number copies of the minor allele carried by individual j (Yang et al. 2011)
176 and p_i is the population allele frequency, the genomic relatedness matrix used by
177 the approach of Bryc et al. (2013) is not. For this approach, the genomic

178 relatedness matrix, $\mathbf{X} = \mathbf{C}\mathbf{C}'$, where \mathbf{C} is an $M \times N$ rectangular matrix with M
179 corresponding to the number of individuals used to estimate \mathbf{X} and N is the
180 number of SNPs, here the element $C_{ij} = \chi_{ij}$, the number of copies of the minor
181 allele carried by individual j . \mathbf{X} was scaled by equation 2.8 in Bryc et al. (2013)
182 using values of $M=110$ and $N=3,709,328$. Only SNPs without missing data were
183 used for this analysis. The number of sub-populations was determined by the
184 number of eigenvalues larger than that expected for a random relatedness
185 matrix given the significance threshold $t = (1 + F) / 2$ (Bryc et al. 2013). Here, t
186 corresponds to a value of 0.993 with our expected inbreeding coefficient (F) of
187 0.986 (Falconer and Mackay 1996) after 20 generations of full-sib mating. Upon
188 seeing that a small number of lines were unusually highly related (see results),
189 we repeated this analysis after removing four lines (line IDs: 29, 134, 159, 206)
190 to verify that the significant results were driven only by these “outliers”.

191

192 *Linkage disequilibrium*

193 We estimated linkage disequilibrium as the square of the inter-variant allele
194 count correlation (r^2) using PLINK v1.9 (Chang et al. 2015). We estimated r^2 in
195 non-overlapping 500 base-pair windows across the entire genome. The analysis
196 was performed on biallelic SNPs that had an average read depth of between 5
197 and 60 and a minimum minor allele frequency of 5%. The four highly related
198 lines ($r \geq 0.1$) were removed prior to the analysis.

199

200 *Nucleotide diversity and neutrality*

201 We estimated nucleotide diversity (π) and Tajima’s D (Tajima 1989) statistic in
202 50 kilobase non-overlapping sliding windows and took the mean for each of the
203 major chromosome arms (2L, 2R, 3L, 3R, and X) using vcftools v0.1.15 (Danecek
204 et al. 2011). The SNP data used for this analysis had an average read depth
205 across all lines of between 5 and 60 but no threshold on minor allele frequency
206 was applied. Again, the four highly related lines were removed prior to this
207 analysis.

208

209 *Genome-wide association of female CHC expression*

210

211 As proof-of-concept we performed a GWAS on a single cuticular hydrocarbon, 2-
212 methyltacosane, (2-Me-C₂₈). CHCs are waxy substances that are secreted on
213 the cuticle with 2-Me-C₂₈ being one of a suite of CHCs that have been extensively
214 studied in this species due to their role in species recognition, mate choice, and
215 desiccation resistance. For each of 94 lines, we extracted CHCs from four virgin
216 females, across two replicate vials using individual whole-body washes in 100µl
217 of the solvent hexane. We used a standard gas chromatography method to
218 quantify the amount of 2-Me-C₂₈ (Blows and Allan 1998). To maintain the trait
219 scale used in previous studies, we transformed the amount of 2-Me-C₂₈ into a
220 log-contrast value following Aitchison (1986), using an additional trait, the CHC
221 9-hexacosane, 9-C_{26:1}, as the divisor. This transformation turns the expression of
222 CHCs into a proportional measure and provides an internal control for other
223 sources of variation including body size and condition.

224

225 Our GWAS contrasts with other analyses performed on the DGRP in that we
226 model trait variation at the individual, rather than line mean level. We applied a
227 single marker mixed effects association analysis, where the following model was
228 fit for every biallelic SNP that had mean coverage between 5 and 60 and sample
229 MAF of 5%:

230

$$231 \quad y_{i,j} = \mu + \text{snp} + G_i + \varepsilon_{i,j} \quad (i = 1, \dots, 94 \quad j = 1, \dots, 4) \quad [2]$$

232

233 Here, CHC expression (y) of replicate individual j on genotype i is modelled as a
234 function of the mean term (μ), the additive fixed effect of the candidate SNP, the
235 polygenic random effect that is captured by the genomic relatedness matrix (\mathbf{G})
236 where \mathbf{G} has a $\sim N(0, \sigma_A^2 \mathbf{A})$ distribution, and the residual error (ε). This model
237 was specifically designed for populations where identical genotypes can be
238 measured independently in multiple organisms such as inbred lines (Kruijer et
239 al. 2015) and was fit using AsReml-R v3 (VSN International) for a total of
240 3,318,503 SNPs. There are a couple of differences between the approach outlined
241 above and other mixed modelling approaches to GWAS implemented in
242 programs such as GEMMA (Zhou and Stephens 2012), FaST-LMM (Lippert et al.
243 2011), and GCTA (Yang et al. 2011). First, the use of individual-level opposed to

244 line mean level observations, allows for estimation of the genomic heritability
245 (Kruijer et al. 2015). Although mapping power is unlikely to be significantly
246 boosted through the use of individual level data *per se* when working with inbred
247 lines (Kruijer et al. 2015), a second aspect to our approach does potentially
248 increase power. Our, albeit slower, approach re-estimates the polygenic variance
249 component when each SNP tested and results in an exact calculation of Wald's
250 test statistic (Zhou and Stephens 2012). This contrasts with the approach used
251 by many mixed model GWAS programs where, to save computation time, this
252 variance component is estimated once in a null model with no fixed effect of SNP
253 and then held constant for each SNP tested. Such an approach produces an
254 approximate value of the test statistic which can result in power loss under some
255 circumstances (Zhou and Stephens 2012). As these circumstances are difficult to
256 predict beforehand, we chose to re-estimate the polygenic random effect despite
257 the computational cost.

258

259 To increase computational speed, we nested this linear model within an R loop
260 that allows the access of multiple cores using the “foreach” and “doMC” packages
261 (Revolution Analytics and Steve Weston 2015). Significant SNPs were identified
262 as those with p-values that passed Bonferroni multiple test correction $-\log_{10}(p)$
263 > 7.8 , we also report SNPs with $-\log_{10}(p) > 5$ for comparison to other *Drosophila*
264 GWAS where this arbitrary threshold value is used (Mackay et al. 2012). We took
265 statistically significant SNPs and annotated them to the current version of the *D.*
266 *serrata* genome (Allen et al. 2017b). If a significant SNP was located within a
267 gene, we blasted the *D. serrata* gene sequence to the *D. melanogaster* genome to
268 determine gene orthology using Flybase (Attrill et al. 2016).

269

270 **Data Availability Statement**

271

272 Raw reads for all sequenced lines are available from the NCBI short read archive
273 under Bioproject ID: PRJNA419238. The genomic relatedness matrices used in
274 the population structure analysis are provided in supplementary files
275 (grm_full_Bryc.txt, grm_reduced_Bryc.txt, grm_full_gcta.txt, and
276 grm_reduced_gcta.txt). We have provided the R code, Bryc.R, that implements the

277 test for large eigenvalues (Bryc et al. 2011). The SNP list used to analyse the data
278 in this study is available from Dryad (doi:XXXXYY) and also from
279 www.chenowethlab.org/resources). The CHC phenotype file is provided as the
280 supplementary file pheno.txt. The linear model used to fit the GWAS in
281 ASREML/R model is provided in the file asreml_gwas.R.

282

283 **Results and Discussion:**

284

285 *Identification of SNPs*

286 We established a panel of 110 inbred lines of *Drosophila serrata* from wild
287 females caught from a single population in Brisbane Australia and sequenced
288 their genomes. 100 base-pair paired-end reads were mapped to the *Drosophila*
289 *serrata* reference genome (Allen et al. 2017b) with a mean coverage of 23.5 ± 0.5
290 reads per line. Using the Joint Genotyper for Inbred Lines (Stone 2012), we
291 identified 15,228,692 biallelic single nucleotide polymorphisms (SNPs) applying
292 a 99% probability threshold. 13,959,239 of these SNPs had a median coverage
293 between 5 and 60 in which over 80% of the lines were genotyped for that
294 variant. Most SNPs segregate at low frequencies (Fig. 1) with 6,090,058
295 instances of singletons, where a SNP was present in only one line. The majority
296 (62%) of the SNPs were annotated to intergenic regions of the genome while
297 12% of SNPs annotated to exonic regions and 26% were found to be in introns.
298 A total of 3,709,329 SNPs met the minimum allele frequency threshold (MAF) of
299 5% to be used for genome-wide association analysis and the estimation of
300 relatedness between lines.

301

302 *Residual heterozygosity*

303 Despite the application of inbreeding for many generations, inbred *Drosophila*
304 lines often contain regions of residual heterozygosity (King et al. 2012; Lack et al.
305 2015; Mackay et al. 2012; Nuzhdin et al. 1997). After 17-20 generations of
306 inbreeding, residual heterozygosity in our lines was very low and we observed
307 only a small proportion of segregating sites within lines, suggesting that
308 inbreeding had successfully fixed variation across these genomes. Of the 110
309 inbred lines, 104 had fewer than 2% segregating SNPs and 82 lines had less than

310 1% segregating SNPs (Fig. 2). Across lines, the average proportion of segregating
311 SNPs was $0.86\% \pm 0.11\%$, less than the theoretical expectation of 1.4% for lines
312 that have experienced 20 generations of full-sib mating which corresponds to an
313 expected inbreeding coefficient of $F = 0.986$ (Falconer and Mackay 1996). This
314 slightly lower than expected level of residual heterozygosity may simply reflect
315 sampling variation, be due to SNPs on the X chromosome, or may indicate
316 purging of partially deleterious alleles during the inbreeding process (Garcia-
317 Dorado 2008; Garcia-Dorado 2012). There was no detectable difference in the
318 fraction of heterozygous SNPs between the lines inbred for 17 and 20
319 generations (ANOVA: $F_{1,108} = 0.0944$, $P = 0.76$). This result suggests that 17
320 generations of inbreeding may be sufficient for future line development.

321

322 Although there are several mechanisms that can inhibit the fixation of an allele
323 within an inbred line, when short-read re-sequencing technology is used for
324 genotyping, loci can falsely appear to be segregating due to the presence of
325 paralogous genes or other repetitive DNA sequences (Treangen and Salzberg
326 2012). If paralogous genes are not represented in the reference genome, DNA
327 sequences from the original gene and a divergent duplicate gene are mapped to
328 the same region of the genome, causing the appearance of segregating loci in the
329 population. When this phenomenon occurs, it is expected that regions of the
330 genome with high “apparent heterozygosity” will associate with a higher read
331 depth than the genome-wide average. Across the genome we found a weak,
332 positive correlation between the level of residual heterozygosity and read depth
333 (Spearman’s $\rho = 0.036$, $p = 2.2 \times 10^{-16}$). Plots of these two factors however,
334 clearly show an alignment of regions with both high levels of residual
335 heterozygosity and read depth, suggesting that the *D. serrata* reference genome
336 could be missing some duplications (Fig. 3). Alternatively, this result could be
337 due to copy number variation among the re-sequenced lines and/or between the
338 reference genome and the 110 lines. We hope that further work and
339 improvement of our genome for this species will elucidate these small regions of
340 residual heterozygosity.

341

342 *Relatedness between lines*

343 Population structure and cryptic relatedness are well known to confound genetic
344 association studies, potentially generating false positive genotype-phenotype
345 associations (Kittles et al. 2002; Knowler et al. 1988). Conceptually, these
346 confounding factors can be described as the unobserved pedigree of the sampled
347 individuals caused by distant relationships (Astle and Balding 2009). The
348 sources of these relationships are varied but include population admixture,
349 inadvertent sampling of close relatives, and the presence of shared chromosomal
350 inversions. Fortunately, the unobserved pedigree can be estimated using marker
351 based approaches. Here we used such an approach to estimate the pairwise
352 relatedness between all lines in the form of a genomic relatedness matrix.

353

354 The structure of the genomic relatedness matrix shows that the majority of the
355 DsGRP lines are unrelated as would be expected of a sample from a large,
356 randomly mating population (Fig. 4). We found a pair of lines that were 100%
357 related to one another, most likely due to contamination either during the
358 inbreeding process or DNA extraction and subsequent library preparation.
359 Generally however, this panel of flies exhibits a lower level of relatedness
360 compared to the DGRP, where 2.7% of the 20,910 possible pairs of lines had
361 estimates of pairwise relatedness over 0.05 (Huang et al. 2014) compared to
362 0.08% of 5,995 pairs of lines reported here. The discrepancy between the levels
363 of relatedness between the DsGRP and DGRP is potentially due to the different
364 demographic histories of the founding populations that generate population
365 structure. North American populations of *D. melanogaster* have relatively
366 complex demographic histories with admixture of African and European
367 ancestors and instances of secondary contact (Pool et al. 2012) compared to the
368 endemic population of *D. serrata*.

369

370 Another likely explanation for the increased levels of relatedness in the DGRP is
371 the presence of common segregating inversions. While chromosomal inversions
372 are known to segregate in *D. serrata*, their frequency and number tends to
373 increase in populations approaching the equator (Stocker et al. 2004).
374 Therefore, it may be the case that founding the DsGRP from the higher latitude of
375 Brisbane has resulted in sampling relatively few inversions. As of yet, these lines

376 have not been karyotyped; however the low levels of relatedness and the lack of
377 any bimodal distribution for residual heterozygosity, such as the one found in the
378 DGRP (Huang et al. 2014), where a portion of the lines had high levels of
379 segregating SNP loci (15-20%), suggests that segregating inversions are
380 negligible in this population.

381

382 We performed an eigendecomposition of the genomic relatedness matrix to test
383 for the presence of population structure using the approach outlined in Bryc et
384 al. (2013). This analysis revealed two large eigenvalues ($\lambda_1 = 20.08$ and $\lambda_2 =$
385 1.12) that were greater than that expected for a random relatedness matrix of
386 equal size (Threshold = 0.993) (Fig. 5). There is therefore evidence that the full
387 set of 110 of lines contain substructure in the form of two subpopulations. We
388 reasoned that the second large eigenvalue was likely caused by the four pairs of
389 lines that were highly related to each other ($A_{jk} = 0.29, 0.38, 0.39,$ and 1.04 ; Fig.
390 4). To test this, we repeated the analysis after randomly removing one line from
391 each of the four pairs of closely related lines. Confirming the prediction, there
392 was only one significantly large eigenvalue in this second analysis ($\lambda_1 = 19.66$).
393 Such a result is expected when the data includes only a single population. To
394 summarise, after the four highly-related lines have been removed, there is no
395 clear evidence for population structure in the DsGRP that would lead to spurious
396 genotype-phenotype associations in genome-wide association analysis.

397

398 *Linkage Disequilibrium*

399 The rapid decay of linkage disequilibrium with genomic distance is a common
400 feature of *Drosophila* species with r^2 dropping below 0.1 within 100 base pairs
401 (Long et al. 1998; Mackay et al. 2012). This allows for higher resolution mapping
402 compared to other species such as maize and humans in which the equivalent
403 decay does not occur until after approximately 2000 base pairs (Remington et al.
404 2001) and 50,000 base pairs (Koch et al. 2013), respectively. In the DsGRP,
405 linkage disequilibrium decays rapidly with r^2 , on average, dropping below 0.1
406 after 75 base pairs. Surprisingly, we observe faster decay on the X chromosome
407 compared to the autosomes (Fig. 6), contrary to Mackay et al. (2012), despite the

408 fact that the X chromosome has a smaller effective population size than the
409 autosomes.

410

411 *Nucleotide diversity and neutrality*

412 For the two cosmopolitan species of *Drosophila* that have been studied
413 extensively, *D. melanogaster* and *D. simulans*, the ancestral populations from
414 Africa consistently exhibit higher levels of polymorphism compared to the
415 derived populations from America and Europe (Andolfatto 2001; Baudry et al.
416 2004; Begun and Aquadro 1993; Grenier et al. 2015; Lack et al. 2015).
417 Presumably, nucleotide diversity is reduced during bottleneck events associated
418 with the colonisation of new habitat. The relatively high estimates of nucleotide
419 diversity for *D. mauritiana*, an endemic species from Mauritius, bolster this trend
420 (Garrigan et al. 2014). We therefore expected that our population of *D. serrata*,
421 founded from the species' ancestral range, would exhibit relatively high levels of
422 nucleotide diversity. We estimated nucleotide diversity (π) along the major
423 chromosome arms 2L, 2R, 3L, 3R, and X using a 50 kilobase non-overlapping
424 sliding window approach (Table 1, Fig. 7). Averaged across the genome, we
425 estimated that $\pi = 0.0079$, which is consistent with the pattern seen in other
426 species of relatively increased levels of nucleotide diversity for populations from
427 ancestral ranges compared to more recently established population outside of
428 the ancestral range (Andolfatto 2001; Baudry et al. 2004; Begun and Aquadro
429 1993; Grenier et al. 2015; Lack et al. 2015).

430

431 Using the same sliding window approach, we tested for departures from
432 neutrality using Tajima's D (Table 1, Fig. 7) (Tajima 1989). As with other
433 populations of *Drosophila*, the DGRP (Mackay et al. 2012), the Zimbabwean
434 population in the Global Diversity Lines (Grenier et al. 2015), and *D. mauritiana*
435 (Garrigan et al. 2014), Tajima's D was negative across the entire genome
436 (Tajima's D = -1.27). This is consistent with an abundance of rare alleles and
437 could be indicative of population expansion or the occurrence of selective
438 sweeps, however this statistic cannot distinguish the effects of demography from
439 selection. In the DsGRP, chromosome arm 2L has higher estimates of Tajima's D
440 compared to the other chromosome arms (Fig. 7). The causes for this pattern

441 will hopefully be resolved through a more in-depth population genomic analysis,
442 as chromosome 2L may potentially harbour more genomic regions that
443 experience balancing selection, which would increase the average value of
444 Tajima's D.

445

446 *Genome-wide association analysis of female CHC expression*

447 A major motivation for developing the DsGRP was to begin connecting molecular
448 variation with standing variation for some of the well-studied quantitative traits
449 of *D. serrata*. Here, we applied genome-wide association analysis to the
450 expression of the CHC 2-Me-C₂₈ in females and identify new candidate genes that
451 might influence trait variation. Using our mixed model approach, we found 4
452 SNPs that passed the 0.05 significance threshold after Bonferroni multiple test
453 correction (Figure 8). Two of the SNPs are situated within genes, while the other
454 two lie within 3kb of genes. We found a further 189 SNPs with p-values lower
455 than a suggestive threshold of 10⁻⁵.

456

457 We note that a GWAS performed on line mean data using the GCTA program
458 (Yang et al. 2011), which like many other mixed model GWAS applications,
459 estimates the polygenic variance only once, detected no SNPs above Bonferroni
460 threshold and only 34 under the arbitrary threshold of 10⁻⁵. We also applied our
461 ASREML approach to line means rather than individuals and found that it
462 detected exactly the same 4 SNPs above Bonferroni as our approach in eq. 2 (190
463 SNPs had p-values lower than 10⁻⁵). It therefore appears that the increase in
464 detection rate is in this case mainly due to the ASREML model re-estimating the
465 polygenic variance for each SNP tested which results in an exact, rather than
466 approximate, calculation of the test statistic (Zhou and Stephens 2012).

467 The majority of the literature regarding the expression of CHCs has identified
468 genes that are related to their production within specialised cells, oenocytes.
469 These genes constitute the major biosynthetic pathway known for CHC
470 production and are involved with fatty-acid synthesis, elongation, desaturation,
471 and reduction (Chertemps et al. 2007; Chertemps et al. 2006; Chung et al. 2014;
472 Fang et al. 2009; Labeur et al. 2002; Marcillac et al. 2005; Wicker-Thomas et al.

473 2015). Although none of the genes associated with the statistically significant
474 SNPs found in this study are involved in this biosynthetic pathway, there are
475 other biological processes involved with CHC expression, as measured by hexane
476 washes from the cuticle. How CHCs are transported from the oenocytes to the
477 cuticle is unknown, this study provides a potential candidate gene involved in
478 this process. One of the significant SNPs resides in the gene *Cht9*, a *chitinase*
479 found on chromosome 2R. *Cht9*, along with a number of other *chitinases* and
480 *imaginal-disc-growth-factors* are important for the development of epithelial
481 apical extracellular matrix, which controls the development and maintenance of
482 wound healing, cell signalling, and organ morphogenesis in *Drosophila* (Galko
483 and Krasnow 2004; Turner 2009). Knocking out the expression of *Cht9* with
484 RNAi leads to deformed cuticles, inability to heal wounds, and defects in larval
485 and adult molting (Pesch et al. 2016), and here, we provide evidence that
486 variation in this gene may also influence other cuticular traits such as CHC
487 abundance. Notwithstanding the small number of lines, the genome-wide
488 association analysis presented here, combined with a previous study that
489 identified the major role of the transcription factor *POU domain motif 3 (pdm3)*
490 for polymorphic female-limited abdominal pigmentation (Yassin et al. 2016),
491 illustrate the potential of the DsGRP to discover novel regions of the genome that
492 underpin the genetic architecture of traits.

493

494 **Conclusion**

495 We have assembled a new resource for the study of quantitative traits and
496 population genomic variation in a non-model *Drosophila* species within its
497 endemic distribution. These reproducible genotypes sampled from a single
498 population not only provide a rich genomic dataset suitable for population
499 genomic studies, but also provide a critical resource for the discovery of genetic
500 variants underlying ecologically important quantitative traits. We hope that the
501 DsGRP will provide a useful complement to other *Drosophila* resources such as
502 the DGRP (Mackay et al. 2012), the DSPR (King et al. 2012), and the *Drosophila*
503 Genome Nexus (Lack et al. 2015).

504

505 In this first characterisation of the DsGRP at the genomic level, we have shown
506 that the inbreeding process has been successful in homogenising the majority of
507 the genome of each of the lines. Through the estimation of the genomic
508 relatedness matrix we have shown that the DsGRP represents a random sample
509 from a large population that contains very low levels of cryptic relatedness.
510 These characteristics, along with rapid decay of linkage disequilibrium, make the
511 DsGRP an ideal resource for the application of genome-wide association analysis
512 and for generating new multifounder QTL mapping populations that will boost
513 mapping power.

514

515 **Acknowledgements**

516 We thank EK Delaney for comments on the manuscript and NC Appleton for
517 assistance in the laboratory. This work was supported by an Australian
518 Postgraduate Award to ARR and funds from the Australian Research Council and
519 The University of Queensland awarded to SFC.

520

521 **References:**

- 522 Aitchison, J. 1986, The statistical analysis of compositional data. London,
523 Chapman & Hall.
- 524 Allen, S. L., R. Bonduriansky, and S. F. Chenoweth. 2013. The genomic distribution
525 of sex-biased genes in *Drosophila serrata*: X chromosome
526 demasculinization, feminization, and hyperexpression in both sexes.
527 *Genome Biology and Evolution* 5:1986-1994.
- 528 Allen, S. L., R. Bonduriansky, C. M. Sgro, and S. F. Chenoweth. 2017a. Sex-biased
529 transcriptome divergence along a latitudinal gradient. *Molecular Ecology*
530 26:1256-1272.
- 531 Allen, S. L., E. K. Delaney, A. Kopp, and S. F. Chenoweth. 2017b. Single-molecule
532 sequencing of the *Drosophila serrata* genome. *G3-Genes Genomes*
533 7:781-788.
- 534 Andolfatto, P. 2001. Contrasting patterns of X-linked and autosomal nucleotide
535 variation in *Drosophila melanogaster* and *Drosophila simulans*. *Molecular*
536 *Biology and Evolution* 18:279-290.
- 537 Astle, W., and D. J. Balding. 2009. Population structure and cryptic relatedness in
538 genetic association studies. *Statistical Science* 24:451-471.
- 539 Attrill, H., K. Falls, J. L. Goodman, G. H. Millburn, G. Antonazzo, A. J. Rey, S. J.
540 Marygold et al. 2016. FlyBase: establishing a Gene Group resource for
541 *Drosophila melanogaster*. *Nucleic Acids Research* 44:D786-D792.
- 542 Ayala, F. J. 1965. Sibling Species of the *Drosophila serrata* group *Evolution*
543 19:538-545.

- 544 Baudry, E., B. Viginier, and M. Veuille. 2004. Non-African populations of
545 *Drosophila melanogaster* have a unique origin. *Molecular Biology and*
546 *Evolution* 21:1482-1491.
- 547 Begun, D. J., and C. F. Aquadro. 1993. African and North American populations of
548 *Drosophila melanogaster* are very different at the DNA level. *Nature*
549 365:548-550.
- 550 Begun, D. J., A. K. Holloway, K. Stevens, L. W. Hillier, Y. P. Poh, M. W. Hahn, P. M.
551 Nista et al. 2007. Population genomics: Whole-genome analysis of
552 polymorphism and divergence in *Drosophila simulans*. *Plos Biology*
553 5:2534-2559.
- 554 Blows, M. W., and R. A. Allan. 1998. Levels of mate recognition within and
555 between two *Drosophila* species and their hybrids. *American Naturalist*
556 152:826-837.
- 557 Blows, M. W., and A. A. Hoffmann. 1993. The genetics of central and marginal
558 populations of *Drosophila serrata* .1. Genetic variation for stress
559 resistance and species borders. *Evolution* 47:1255-1270.
- 560 Brake, I., and G. Bachli. 2008, *World Catalogue of Insects, Volume 9: Dro-*
561 *sophilidae (Diptera)*. Stenstrup, Denmark, Apollo Books.
- 562 Bryc, K., W. Bryc, and J. W. Silverstein. 2013. Separation of the largest eigenvalues
563 in eigenanalysis of genotype data from discrete subpopulations.
564 *Theoretical Population Biology* 89:34-43.
- 565 Burke, M. K., E. G. King, P. Shahrestani, M. R. Rose, and A. D. Long. 2014. Genome-
566 wide association study of extreme longevity in *Drosophila melanogaster*.
567 *Genome Biology and Evolution* 6:1-11.
- 568 Chang, C. C., C. C. Chow, L. Tellier, S. Vattikuti, S. M. Purcell, and J. J. Lee. 2015.
569 Second-generation PLINK: rising to the challenge of larger and richer
570 datasets. *Gigascience* 4.
- 571 Chertemps, T., L. Duportets, C. Labeur, R. Ueda, K. Takahashi, K. Saigo, and C.
572 Wicker-Thomas. 2007. A female-biased expressed elongase involved in
573 long-chain hydrocarbon biosynthesis and courtship behavior in
574 *Drosophila melanogaster*. *Proceedings of the National Academy of*
575 *Sciences of the United States of America* 104:4273-4278.
- 576 Chertemps, T., L. Duportets, C. Labeur, M. Ueyama, and C. Wicker-Thomas. 2006.
577 A female-specific desaturase gene responsible for diene hydrocarbon
578 biosynthesis and courtship behaviour in *Drosophila melanogaster*. *Insect*
579 *Molecular Biology* 15:465-473.
- 580 Chung, H., D. W. Loehlin, H. D. Dufour, K. Vaccarro, J. G. Millar, and S. B. Carroll.
581 2014. A single gene affects both ecological divergence and mate choice in
582 *Drosophila*. *Science* 343:1148-1151.
- 583 Clark, A. G., M. B. Eisen, D. R. Smith, C. M. Bergman, B. Oliver, T. A. Markow, T. C.
584 Kaufman et al. 2007. Evolution of genes and genomes on the *Drosophila*
585 phylogeny. *Nature* 450:203-218.
- 586 Da Lage, J. L., G. J. Kergoat, F. Maczkowiak, J. F. Silvain, M. L. Cariou, and D.
587 Lachaise. 2007. A phylogeny of *Drosophilidae* using the *Amyrel* gene:
588 questioning the *Drosophila melanogaster* species group boundaries.
589 *Journal of Zoological Systematics and Evolutionary Research* 45:47-63.
- 590 Danecek, P., A. Auton, G. Abecasis, C. A. Albers, E. Banks, M. A. DePristo, R. E.
591 Handsaker et al. 2011. The variant call format and VCFtools.
592 *Bioinformatics* 27:2156-2158.

- 593 Dembeck, L. M., W. Huang, M. M. Magwire, F. Lawrence, R. F. Lyman, and T. F. C.
594 Mackay. 2015. Genetic architecture of abdominal pigmentation in
595 *Drosophila melanogaster*. Plos Genetics 11.
- 596 Falconer, D. S., and T. F. C. Mackay. 1996, Introduction to Quantitative Genetics.
597 Burnt Mill, England, Longman.
- 598 Fang, S., C. T. Ting, C. R. Lee, K. H. Chu, C. C. Wang, and S. C. Tsauro. 2009.
599 Molecular evolution and functional diversification of fatty acid
600 desaturases after recurrent gene duplication in *Drosophila*. Molecular
601 Biology and Evolution 26:1447-1456.
- 602 Frentiu, F. D., M. Adamski, E. A. McGraw, M. W. Blows, and S. F. Chenoweth. 2009.
603 An expressed sequence tag (EST) library for *Drosophila serrata*, a model
604 system for sexual selection and climatic adaptation studies. BMC
605 Genomics 10.
- 606 Frentiu, F. D., and S. F. Chenoweth. 2010. Clines in cuticular hydrocarbons in two
607 *Drosophila* species with independent population histories. Evolution
608 64:1784-1794.
- 609 Galko, M. J., and M. A. Krasnow. 2004. Cellular and genetic analysis of wound
610 healing in *Drosophila* larvae. Plos Biology 2:1114-1126.
- 611 Garcia-Dorado, A. 2008. A simple method to account for natural selection when
612 predicting inbreeding depression. Genetics 180:1559-1566.
- 613 —. 2012. Understanding and predicting the fitness decline of shrunk
614 populations: inbreeding, purging, mutation, and standard selection.
615 Genetics 190:1461-1476.
- 616 Garrigan, D., S. B. Kingan, A. J. Geneva, J. P. Vedanayagam, and D. C. Presgraves.
617 2014. Genome diversity and divergence in *Drosophila mauritiana*:
618 multiple signatures of faster X evolution. Genome Biology and Evolution
619 6:2444-2458.
- 620 Gerken, A. R., O. C. Eller, D. A. Hahn, and T. J. Morgan. 2015. Constraints,
621 independence, and evolution of thermal plasticity: probing genetic
622 architecture of long- and short-term thermal acclimation. Proceedings of
623 the National Academy of Sciences of the United States of America
624 112:4399-4404.
- 625 Gosden, T. P., and S. F. Chenoweth. 2014. The evolutionary stability of cross-sex,
626 cross-trait genetic covariances. Evolution 68:1687-1697.
- 627 Grenier, J. K., J. R. Arguello, M. C. Moreira, S. Gottipati, J. Mohammed, S. R. Hackett,
628 R. Boughton et al. 2015. Global Diversity Lines-A Five-Continent
629 Reference Panel of Sequenced *Drosophila melanogaster* Strains. G3-Genes
630 Genomes Genetics 5:593-603.
- 631 Hallas, R., M. Schiffer, and A. A. Hoffmann. 2002. Clinal variation in *Drosophila*
632 *serrata* for stress resistance and body size. Genetical Research 79:141-
633 148.
- 634 Harbison, S. T., L. J. McCoy, and T. F. C. Mackay. 2013. Genome-wide association
635 study of sleep in *Drosophila melanogaster*. BMC Genomics 14.
- 636 He, B. Z., M. Z. Ludwig, D. A. Dickerson, L. Barse, B. Arun, B. J. Vilhjalmsón, S. Y.
637 Park et al. 2014. Effect of genetic variation in a *Drosophila* model of
638 diabetes-associated misfolded human proinsulin. Genetics 196:557-+.
- 639 Hine, E., S. Lachish, M. Higgie, and M. W. Blows. 2002. Positive genetic correlation
640 between female preference and offspring fitness. Proceedings of the Royal
641 Society B-Biological Sciences 269:2215-2219.

- 642 Huang, W., A. Massouras, Y. Inoue, J. Peiffer, M. Ramia, A. M. Tarone, L. Turlapati
643 et al. 2014. Natural variation in genome architecture among 205
644 *Drosophila melanogaster* Genetic Reference Panel lines. *Genome Research*
645 24:1193-1208.
- 646 Hunter, C. M., W. Huang, T. F. C. Mackay, and N. D. Singh. 2016. The genetic
647 architecture of natural variation in recombination rate in *Drosophila*
648 *melanogaster*. *Plos Genetics* 12.
- 649 Jenkins, N. L., and A. A. Hoffmann. 1999. Limits to the southern border of
650 *Drosophila serrata*: cold resistance, heritable variation, and trade-offs.
651 *Evolution* 53:1823-1834.
- 652 Kellermann, V., B. van Heerwaarden, C. M. Sgro, and A. A. Hoffmann. 2009.
653 Fundamental evolutionary limits in ecological traits drive *Drosophila*
654 species distributions. *Science* 325:1244-1246.
- 655 King, E. G., G. Kislukhin, K. N. Walters, and A. D. Long. 2014. Using *Drosophila*
656 *melanogaster* To identify chemotherapy toxicity genes. *Genetics* 198:31-+.
- 657 King, E. G., S. J. Macdonald, and A. D. Long. 2012. Properties and power of the
658 *Drosophila* Synthetic Population Resource for the routine dissection of
659 complex traits. *Genetics* 191:935-U522.
- 660 Kislukhin, G., E. G. King, K. N. Walters, S. J. Macdonald, and A. D. Long. 2013. The
661 genetic architecture of methotrexate toxicity is similar in *Drosophila*
662 *melanogaster* and humans. *G3-Genes Genomes Genetics* 3:1301-1310.
- 663 Kittles, R. A., W. D. Chen, R. K. Panguluri, C. Ahaghotu, A. Jackson, C. A.
664 Adebamowo, R. Griffin et al. 2002. CYP3A4-V and prostate cancer in
665 African Americans: causal or confounding association because of
666 population stratification? *Human Genetics* 110:553-560.
- 667 Knowler, W. C., R. C. Williams, D. J. Pettitt, and A. G. Steinberg. 1988. Gm3;5,13,14
668 and type 2 diabetes mellitus: an association in American Indians and
669 genetic admixture. *American Journal of Human Genetics* 43:520-526.
- 670 Koch, E., M. Ristroph, and M. Kirkpatrick. 2013. Long range linkage
671 disequilibrium across the human genome. *Plos One* 8.
- 672 Kruijer, W., M. P. Boer, M. Malosetti, P. J. Flood, B. Engel, R. Kooke, J. J. B.
673 Keurentjes et al. 2015. Marker-based estimation of heritability in
674 immortal populations. *Genetics* 199:379-393.
- 675 Labeur, C., R. Dallerac, and C. Wicker-Thomas. 2002. Involvement of *desat1* gene
676 in the control of *Drosophila melanogaster* pheromone biosynthesis.
677 *Genetica* 114:269-274.
- 678 Lack, J. B., C. M. Cardeno, M. W. Crepeau, W. Taylor, R. B. Corbett-Detig, K. A.
679 Stevens, C. H. Langley et al. 2015. The *Drosophila* Genome Nexus: a
680 population genomic resource of 623 *Drosophila melanogaster* genomes,
681 including 197 from a single ancestral range population. *Genetics*
682 199:1229-U1553.
- 683 Langley, C. H., K. Stevens, C. Cardeno, Y. C. G. Lee, D. R. Schrider, J. E. Pool, S. A.
684 Langley et al. 2012. Genomic variation in natural populations of
685 *Drosophila melanogaster*. *Genetics* 192:533-+.
- 686 Lemeunier, F., J. R. David, L. Tsacas, and M. Ashburner. 1986. The melanogaster
687 species group. *The genetics and biology of Drosophila* 3:147-256.
- 688 Li, H. 2013. Aligning sequence reads, clone sequences and assembly contigs with
689 BWA-MEM. arXiv:1303.3997v1 [q-bio.GN].

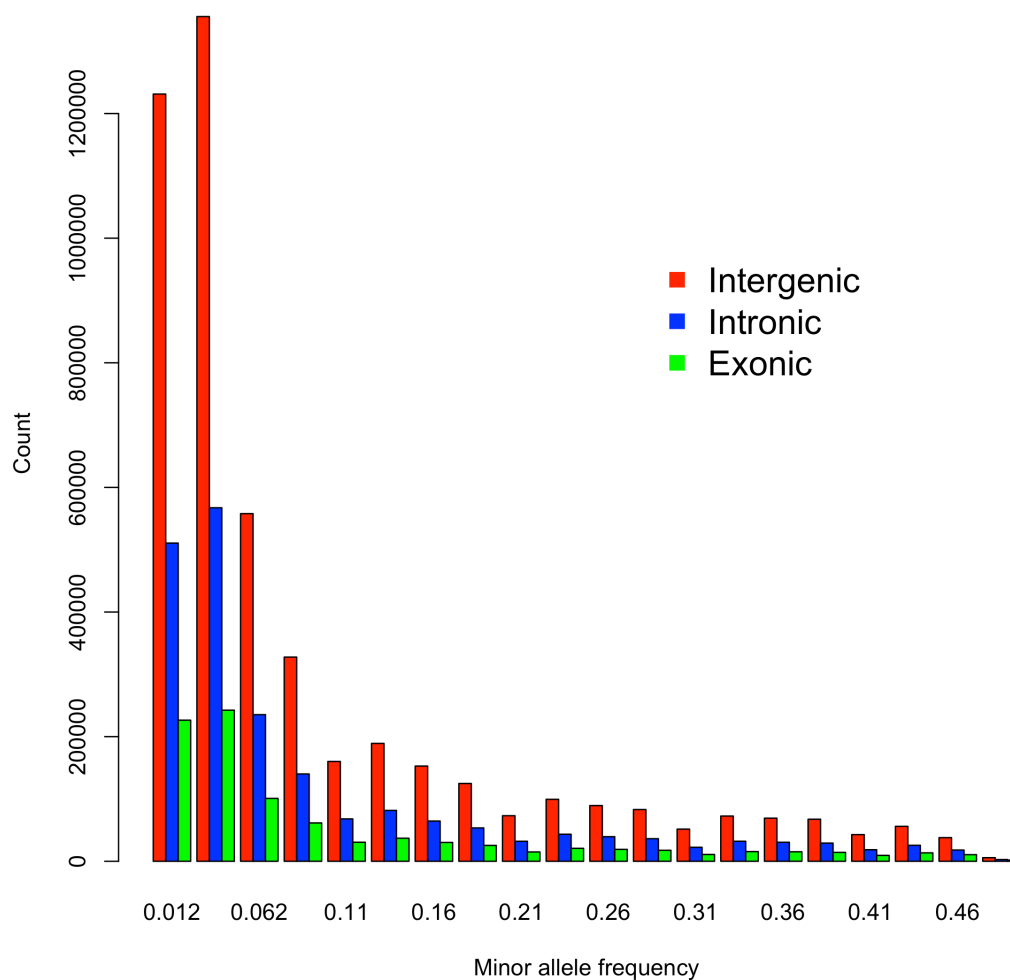
- 690 Lippert, C., J. Listgarten, Y. Liu, C. M. Kadie, R. I. Davidson, and D. Heckerman.
691 2011. FaST linear mixed models for genome-wide association studies.
692 Nature Methods 8:833-U894.
- 693 Long, A. D., R. F. Lyman, C. H. Langley, and T. F. C. Mackay. 1998. Two sites in the
694 Delta gene region contribute to naturally occurring variation in bristle
695 number in *Drosophila melanogaster*. Genetics 149:999-1017.
- 696 Mackay, T. F. C., S. Richards, E. A. Stone, A. Barbadilla, J. F. Ayroles, D. H. Zhu, S.
697 Casillas et al. 2012. The *Drosophila melanogaster* Genetic Reference Panel.
698 Nature 482:173-178.
- 699 Magiafoglou, A., M. E. Carew, and A. A. Hoffmann. 2002. Shifting clinal patterns
700 and microsatellite variation in *Drosophila serrata* populations: a
701 comparison of populations near the southern border of the species range.
702 Journal of Evolutionary Biology 15:763-774.
- 703 Magwire, M. M., D. K. Fabian, H. Schweyen, C. Cao, B. Longdon, F. Bayer, and F. M.
704 Jiggins. 2012. Genome-wide association studies reveal a simple genetic
705 basis of resistance to naturally coevolving viruses in *Drosophila*
706 *melanogaster*. Plos Genetics 8.
- 707 Marcillac, F., F. Bousquet, J. Alabouvette, F. Savarit, and J. F. Ferveur. 2005. A
708 mutation with major effects on *Drosophila melanogaster* sex pheromones.
709 Genetics 171:1617-1628.
- 710 Marriage, T. N., E. G. King, A. D. Long, and S. J. Macdonald. 2014. Fine-mapping
711 nicotine resistance loci in *Drosophila* using a multiparent advanced
712 generation inter-cross population. Genetics 198:45-+.
- 713 McGuigan, K., J. M. Collet, E. A. McGraw, Y. X. H. Ye, S. L. Allen, S. F. Chenoweth,
714 and M. W. Blows. 2014. The nature and extent of mutational pleiotropy in
715 gene expression of male *Drosophila serrata*. Genetics 196:911-+.
- 716 McGuigan, K., D. Petfield, and M. W. Blows. 2011. Reducing mutation load
717 through sexual selection on males. Evolution 65:2816-2829.
- 718 McKenna, A., M. Hanna, E. Banks, A. Sivachenko, K. Cibulskis, A. Kernytsky, K.
719 Garimella et al. 2010. The Genome Analysis Toolkit: a MapReduce
720 framework for analyzing next-generation DNA sequencing data. Genome
721 Research 20:1297-1303.
- 722 Nuzhdin, S. V., E. G. Pasyukova, C. L. Dilda, Z. B. Zeng, and T. F. C. Mackay. 1997.
723 Sex-specific quantitative trait loci affecting longevity in *Drosophila*
724 *melanogaster*. Proceedings of the National Academy of Sciences of the
725 United States of America 94:9734-9739.
- 726 Pesch, Y. Y., D. Riedel, K. R. Patil, G. Loch, and M. Behr. 2016. Chitinases and
727 imaginal disc growth factors organize the extracellular matrix formation
728 at barrier tissues in insects. Scientific Reports 6.
- 729 Pool, J. E., R. B. Corbett-Detig, R. P. Sugino, K. A. Stevens, C. M. Cardeno, M. W.
730 Crepeau, P. Duchon et al. 2012. Population genomics of Sub-Saharan
731 *Drosophila melanogaster*: African Diversity and non-African admixture.
732 Plos Genetics 8.
- 733 Remington, D. L., J. M. Thornsberry, Y. Matsuoka, L. M. Wilson, S. R. Whitt, J.
734 Doeblay, S. Kresovich et al. 2001. Structure of linkage disequilibrium and
735 phenotypic associations in the maize genome. Proceedings of the National
736 Academy of Sciences of the United States of America 98:11479-11484.
- 737 Shorter, J., C. Couch, W. Huang, M. A. Carbone, J. Peiffer, R. R. H. Anholt, and T. F. C.
738 Mackay. 2015. Genetic architecture of natural variation in *Drosophila*

- 739 *melanogaster* aggressive behavior. Proceedings of the National Academy
740 of Sciences of the United States of America 112:E3555-E3563.
- 741 Stark, A., M. F. Lin, P. Kheradpour, J. S. Pedersen, L. Parts, J. W. Carlson, M. A.
742 Crosby et al. 2007. Discovery of functional elements in 12 *Drosophila*
743 genomes using evolutionary signatures. Nature 450:219-232.
- 744 Stocker, A. J., B. Foley, and A. Hoffmann. 2004. Inversion frequencies in
745 *Drosophila serrata* along an eastern Australian transect. Genome 47:1144-
746 1153.
- 747 Stocker, A. J., B. B. Rusuwa, M. J. Blacket, F. D. Frentiu, M. Sullivan, B. R. Foley, S.
748 Beatson et al. 2012. Physical and linkage maps for *Drosophila serrata*, a
749 model species for studies of clinal adaptation and sexual selection. G3-
750 Genes Genomes Genetics 2:287-297.
- 751 Stone, E. A. 2012. Joint genotyping on the fly: identifying variation among a
752 sequenced panel of inbred lines. Genome Research 22:966-974.
- 753 Sturgill, D., Y. Zhang, M. Parisi, and B. Oliver. 2007. Demasculinization of X
754 chromosomes in the *Drosophila* genus. Nature 450:238-U233.
- 755 Tajima, F. 1989. Statistical method for testing the neutral mutation hypothesis by
756 DNA polymorphism Genetics 123:585-595.
- 757 Treangen, T. J., and S. L. Salzberg. 2012. Repetitive DNA and next-generation
758 sequencing: computational challenges and solutions. Nature Reviews
759 Genetics 13:36-46.
- 760 Turner, J. R. 2009. Intestinal mucosal barrier function in health and disease.
761 Nature Reviews Immunology 9:799-809.
- 762 Unckless, R. L., S. M. Rottschaefer, and B. P. Lazzaro. 2015. A genome-wide
763 association study for nutritional indices in *Drosophila*. G3-Genes Genomes
764 Genetics 5:417-425.
- 765 van Heerwaarden, B., V. Kellermann, M. Schiffer, M. Blacket, C. M. Sgro, and A. A.
766 Hoffmann. 2009. Testing evolutionary hypotheses about species borders:
767 patterns of genetic variation towards the southern borders of two
768 rainforest *Drosophila* and a related habitat generalist. Proceedings of the
769 Royal Society B-Biological Sciences 276:1517-1526.
- 770 Weber, A. L., G. F. Khan, M. M. Magwire, C. L. Tabor, T. F. C. Mackay, and R. R. H.
771 Anholt. 2012. Genome-wide association analysis of oxidative stress
772 resistance in *Drosophila melanogaster*. Plos One 7.
- 773 Wicker-Thomas, C., D. Garrido, G. Bontonou, L. Napal, N. Mazuras, B. Denis, T.
774 Rubin et al. 2015. Flexible origin of hydrocarbon/pheromone precursors
775 in *Drosophila melanogaster*. Journal of Lipid Research 56:2094-2101.
- 776 Yang, J. A., S. H. Lee, M. E. Goddard, and P. M. Visscher. 2011. GCTA: a tool for
777 genome-wide complex trait analysis. American Journal of Human Genetics
778 88:76-82.
- 779 Yassin, A. 2013. Phylogenetic classification of the Drosophilidae Rondani
780 (Diptera): the role of morphology in the postgenomic era. Systematic
781 Entomology 38:349-364.
- 782 Yassin, A., E. K. Delaney, A. J. Reddiex, T. D. Seher, H. Bastide, N. C. Appleton, J. B.
783 Lack et al. 2016. The pdm3 locus is a hotspot for recurrent evolution of
784 female-limited color dimorphism in *Drosophila*. Current Biology 26:2412-
785 2422.

786 Zhang, Y., D. Sturgill, M. Parisi, S. Kumar, and B. Oliver. 2007. Constraint and
787 turnover in sex-biased gene expression in the genus *Drosophila*. *Nature*
788 450:233-U232.
789 Zhou, X., and M. Stephens. 2012. Genome-wide efficient mixed-model analysis for
790 association studies. *Nature Genetics* 44:821-U136.
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820

821 **Figures and tables:**

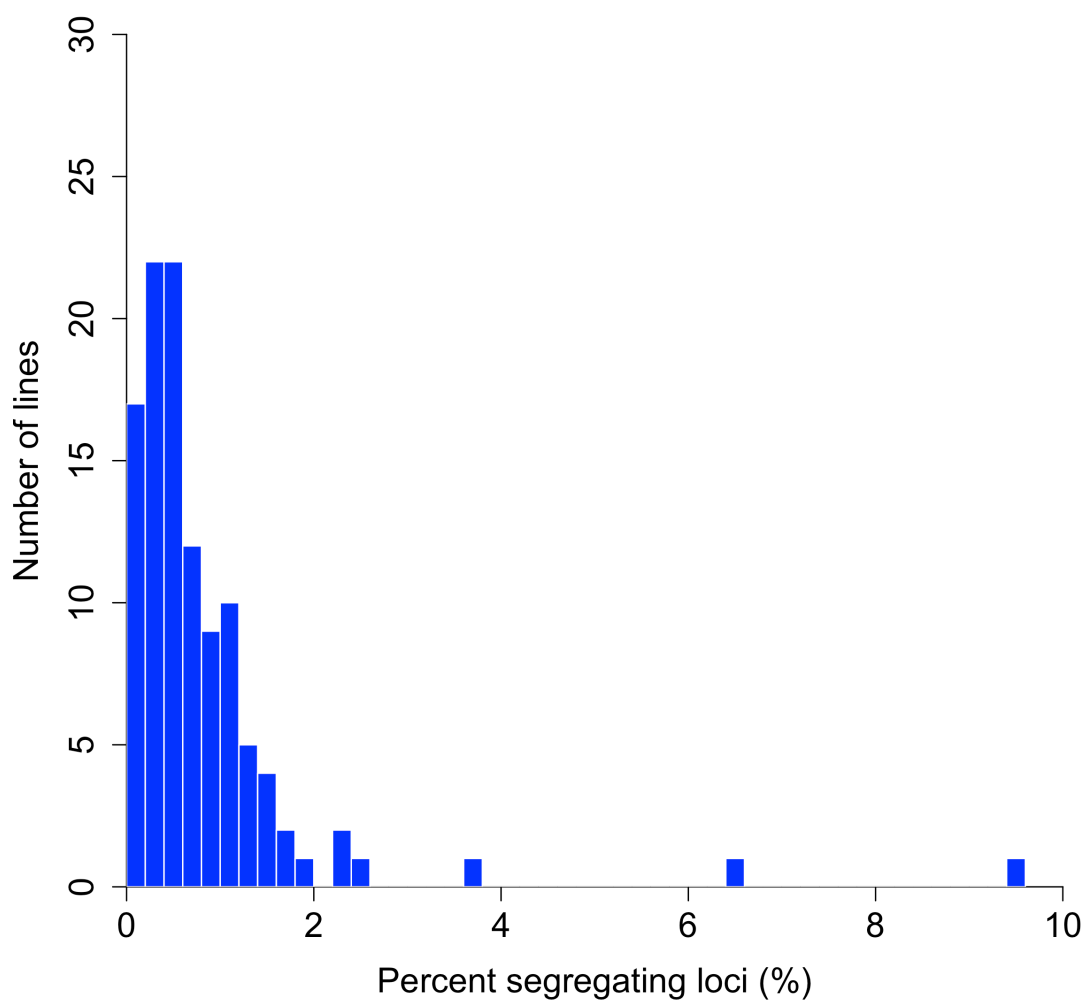
822



823

824 Figure 1: The allele frequency spectrum of SNPs annotated as intergenic (red),
825 intronic (blue), or exonic (green). Singletons (MAF = 0.009) are not shown. We
826 identified 3,748,429 singletons from intergenic regions, 1,535,651 from intronic
827 regions, and 754,075 from exonic regions of the genome.

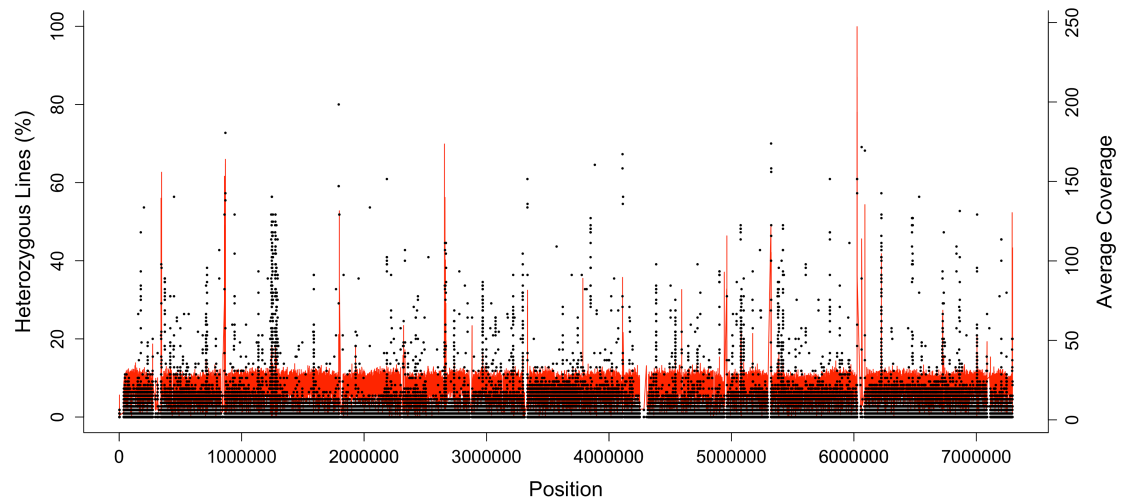
828



829

830 Figure 2: The distribution of residual heterozygosity as measured by the
831 percentage of total genotyped biallelic SNPs (loci) that were called as
832 heterozygous within each inbred line by the Joint Genotyper of Inbred Lines
833 (JGIL).

834



835

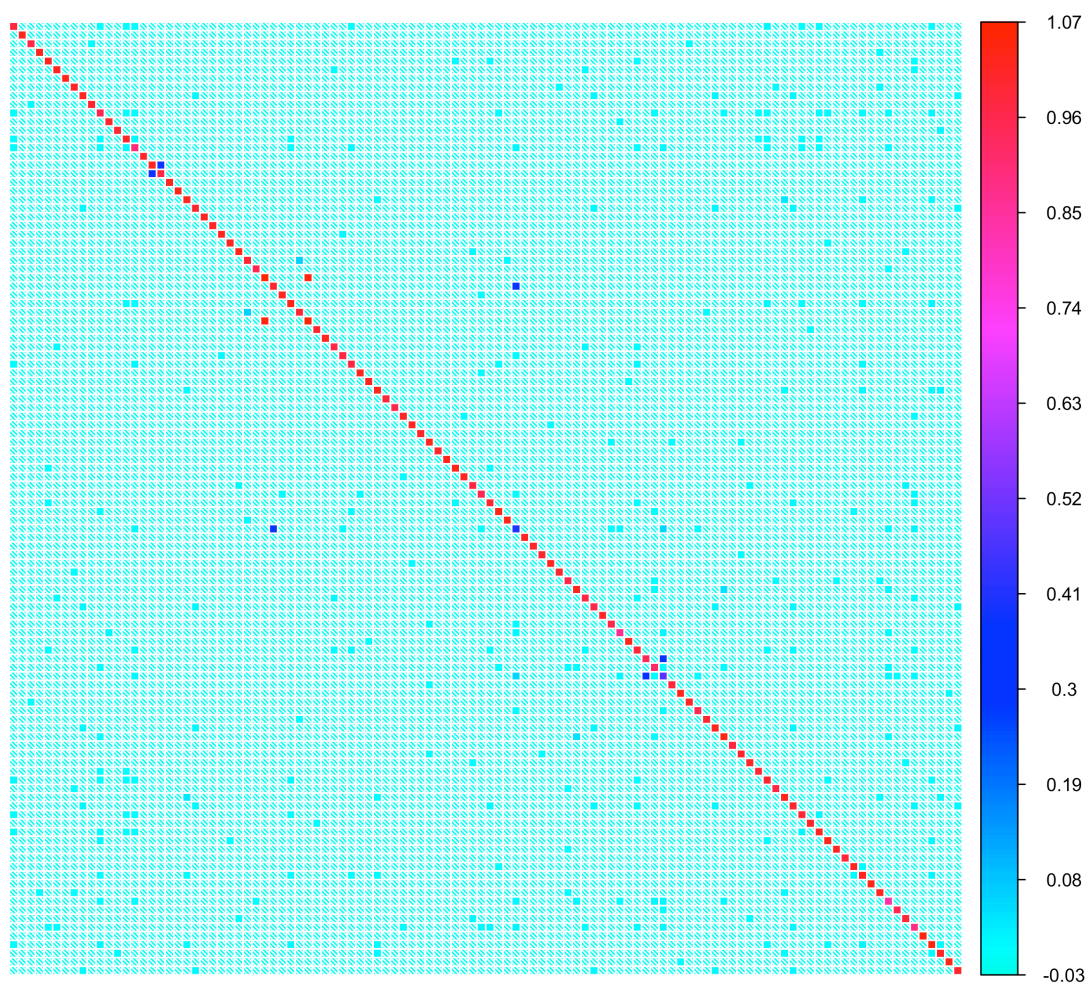
836 Figure 3: The proportion of lines that are heterozygous at every site along the
837 largest scaffold (scf7180000003208) of the reference genome (black) overlaid
838 with the average read depth at each site shown in red. Peaks in red lines
839 represent potential gene duplications that have collapsed to the same region
840 during assembly.

841

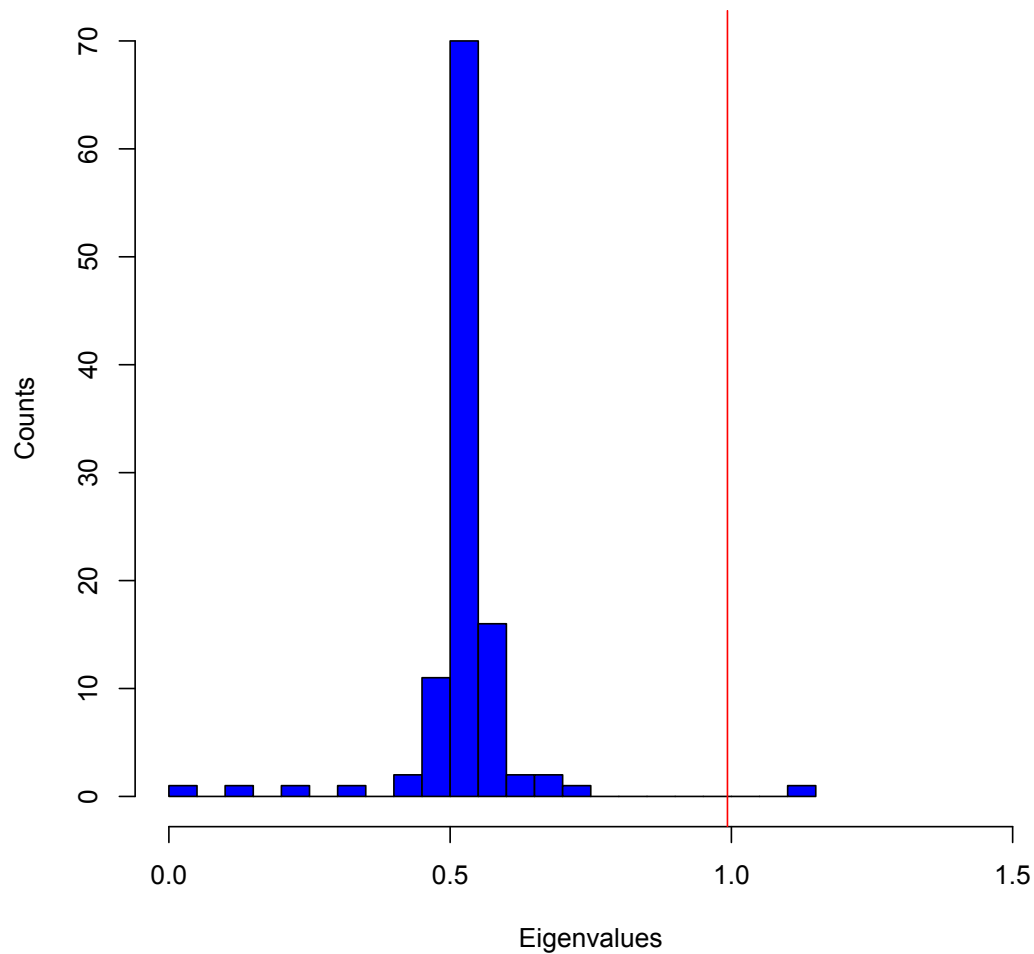
842

843

844



845 Figure 4: Heat-map of the genomic relatedness matrix for each pair of the 110
846 inbred DsGRP lines. Each coloured squared is an estimate of genomic
847 relatedness between a pair of inbred lines estimated by GCTA (Yang et al. 2011).
848 The shade of colour represents the degree of relatedness, with light blue showing
849 low levels of relatedness and red high levels of relatedness.
850

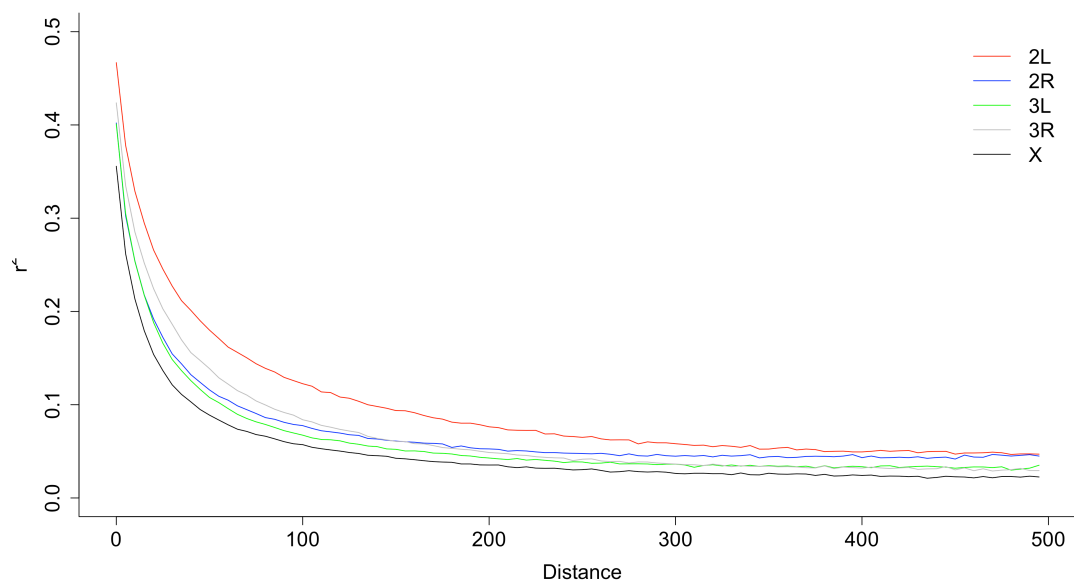


851

852 Figure 5: Distribution of eigenvalues from an eigendecomposition of the
853 genomic relatedness matrix for all 110 lines excluding one large eigenvalue
854 where $\lambda = 20.08$. Eigendecomposition of the genomic relatedness matrix, \mathbf{X} ,
855 which was scaled by equation 2.8 in Bryc et al. (2013) using values of $N =$
856 $3,709,328$ and $M = 110$. Here, N corresponds to the number of SNPs used to
857 estimate \mathbf{X} and M is the number of lines. Only SNPs with without missing data
858 were used for this analysis. The red vertical line corresponds to the significance
859 threshold, t , for declaring an eigenvalue larger than that expected for a random
860 relatedness matrix. $t = (1 + F) / 2$ and corresponds to a value of 0.993 with our
861 expected inbreeding coefficient (F) of 0.986 after 20 generations of full-sib
862 mating. The two largest eigenvalues were significant, $\lambda_1 = 20.08$ (not plotted)
863 and $\lambda_2 = 1.12$.

864

865



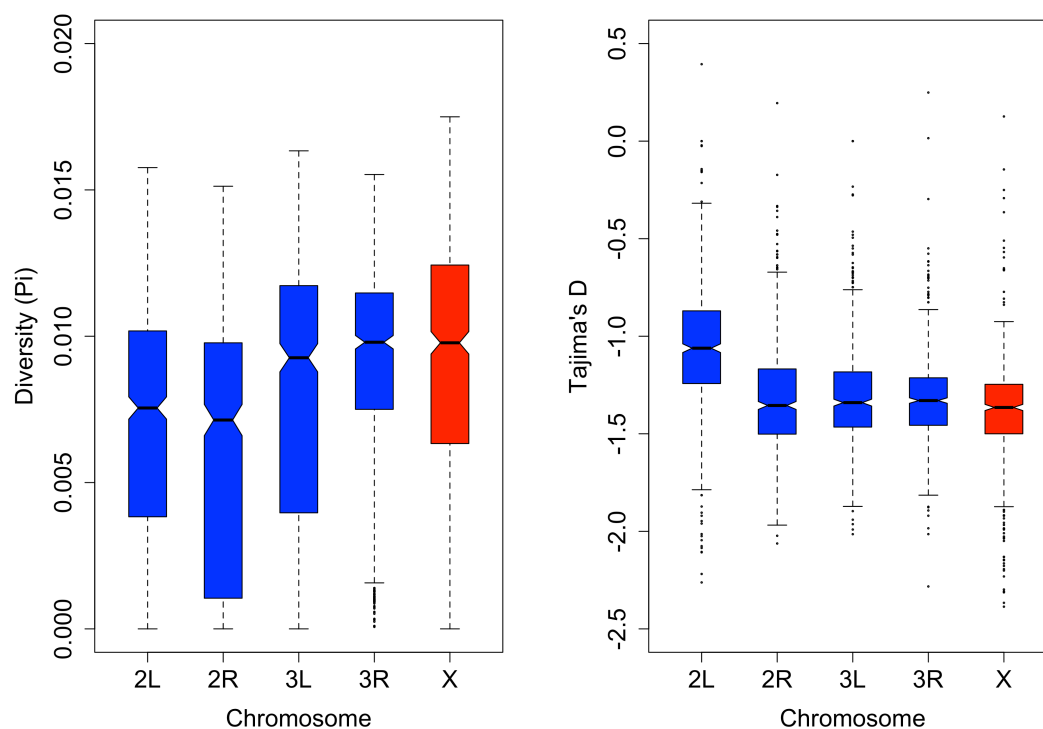
866

867 Figure 6: Decay of linkage disequilibrium (r^2) between SNPs with genomic
868 distance (bp) in the DsGRP. Values are averaged across each chromosome.

869

870

871



872

873 Figure 7: Boxplots of nucleotide diversity and Tajima's D by chromosome arm.

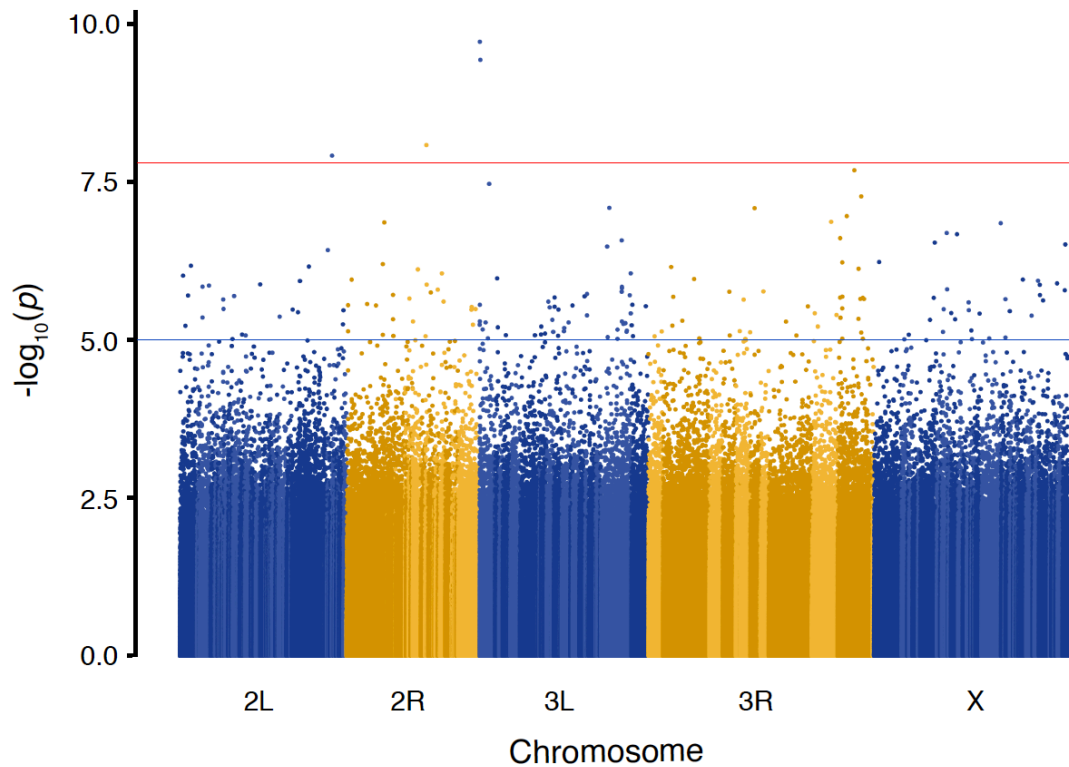
874 Shown are the estimates from 50 kilobase non-overlapping sliding windows with

875 the breakdown of the number of windows per chromosome as follows: 2L = 691,

876 2R = 664, 3L = 641, 3R = 741, and X = 628.

877

878



879

880

881 Figure 8: Genome-wide association for 2-Me-C₂₈ expression in female *Drosophila*
882 *serrata*. Red line indicates Bonferroni threshold corresponding to $P = 0.05$ and
883 the blue indicates an arbitrary significance threshold of $P = 10^{-5}$. A total of
884 3,318,503 biallelic SNPs were analyzed with a minimum minor allele frequency
885 of 0.5. Alternating colour shading within chromosomes indicates different
886 contigs in the *D. serrata* genome assembly. Contig order has not yet been
887 established for the genome.

888