# Patterns and Causes of Signed Linkage Disequilibria in Flies and Plants

**George Sandler[1], Stephen I. Wright[1,2]\*, Aneil F. Agrawal[1,2]\***

[1]Department of Ecology and Evolutionary Biology, University of Toronto, 25 Willcocks Street, Toronto, ON M5S 3B2, Canada; [2]Center for Analysis of Genome Evolution and Function, University of Toronto, 25 Willcocks Street, Toronto, ON M5S 3B2, Canada

*These authors made equal contributions.


**Corresponding author**

**George Sandler**

**Email: george.sandler@mail.utoronto.ca**

Article: Discoveries

## 1 Abstract

2 Most empirical studies of linkage disequilibrium (LD) study its magnitude, ignoring its

3 sign. Here, we examine patterns of signed LD in two population genomic datasets, one

4 from *Capsella grandiflora* and one from *Drosophila melanogaster.* We consider how

5 processes such as drift, admixture, Hill-Robertson interference, and epistasis may

6 contribute to these patterns. We report that most types of mutations exhibit positive LD,

7 particularly, if they are predicted to be less deleterious. We show with simulations that

8 this pattern arises easily in a model of admixture or distance biased mating, and that

9 genome-wide differences across site types are generally expected due to differences in

10 the strength of purifying selection even in the absence of epistasis. We further explore

11 how signed LD decays on a finer scale, showing that loss of function mutations exhibit

12 particularly positive LD across short distances, a pattern consistent with intragenic

13 antagonistic epistasis. Controlling for genomic distance, signed LD in *C. grandiflora*

14 decays faster within genes, compared to between genes, likely a by-product of frequent

15 recombination in gene promoters known to occur in plant genomes. Finally, we use

16 information from published biological networks to explore whether there is evidence for

17 negative synergistic epistasis between interacting radical missense mutations. In *D.*

18 *melanogaster* networks, we find a modest but significant enrichment of negative LD,

19 consistent with the possibility of intra-network negative synergistic epistasis.

20

**Introduction**

21

22  Linkage disequilibrium (LD), the association of different alleles across the genome, is a

23  general feature of population genomic datasets, often revealing clues of ongoing

24  evolutionary or demographic processes (McEvoy et al. 2011). For example, in finite

25  populations, drift can be a ready source of LD, generating both positive and negative

26  associations between alleles (Hill and Robertson 1968). While unsigned LD has been

27  extensively studied in population genetics (through statistics such as $r^2$), signed LD has

28  received relatively less attention, despite the fact that the sign of allelic associations can

29  also provide useful information. Here we refer to positive associations as those between

30  two common alleles (or equivalently between two rare alleles), and negative

31  associations as those between common and rare alleles. Demographic processes such

32  as admixture and population structure can create LD, where unlike drift in a single

33  panmictic population, an overabundance of positive associations is expected between

34  pairs of migrant alleles (Chakraborty and Weiss 1988; Stephens et al. 1994; Pfaff et al.

35  2001). Selective processes can also be a source of LD; for example, ongoing strong

36  selective sweeps can be characterised by an elevation of unsigned LD around the

37  sweeping haplotype (McVean 2007). Non-independence of mutational events, e.g.

38  multinucleotide mutations, arise at non-negligible frequencies in several species

39  (Schrider et al. 2011), and could also be an important source of positive LD among de-

40  novo mutations (Ragsdale 2021). Finally, unsigned LD can also be used to analyze

41  patterns of recombination across the genome, as recombination is expected to break

42  down any existing LD (Auton and McVean 2007).

43          LD can also build up due to selection against deleterious mutations in two

44  different ways. First, Hill Robertson interference (resulting from the interaction of

45  selection and drift) can cause negative associations to build up among deleterious

46  mutations, if recombination between them is limited (Hill and Robertson 1966). In

47  sexually reproducing organisms such as humans, this process has recently been

48  suggested to build up negative LD among physically proximal, missense mutations

49  (Garcia and Lohmueller 2020). Second, negative selection can cause LD among

50  deleterious mutations to build up if epistasis is present (Kondrashov 1995; Sohail et al.

51 2017). Under the null model of multiplicative fitness, where each mutation contributes to

52 a reduction in fitness independently of other mutations, LD is not expected to

53 accumulate. Synergistic epistasis, where each additional deleterious mutation reduces

54 fitness by a greater magnitude, creates negative LD among deleterious mutations and

55 vice versa for antagonistic epistasis (Kimura and Maruyama 1966; Kondrashov 1982).

56 Synergistic epistasis among deleterious mutations is of particular interest

57 because such epistasis has several evolutionary consequences. For example, negative

58 synergistic epistasis allows for lower mutation loads under mutation-selection balance,

59 and can influence the evolution of sex and recombination (Kimura and Maruyama 1966;

60 Crow and Kimura 1970; Crow and Kimura 1979; Kondrashov 1982; see also Barton

61 1995). Despite considerable interest, empirical data on epistasis among deleterious

62 mutations is limited with most data coming from microorganisms assayed in a lab

63 setting. These studies have found that synergistic and antagonistic interactions are both

64 common so that mean epistasis is close to zero  (Elena and Lenski 1997; Agrawal and

65 Whitlock 2010; Lalić and Elena 2012; Bank et al. 2015; Puchta et al. 2016). A recent

66 study by Sohail et al. (2017) used a different approach to make inferences about

67 epistasis. They examined patterns of signed LD among rare loss of function (LOF)

68 mutations in humans and fruit flies demonstrating that across several datasets LOF

69 mutations had significantly lower values of signed LD than their neutral reference

70 (synonymous sites), a pattern consistent with the action of negative synergistic

71 epistasis.

72 Here we examine patterns of LD across several classes of mutations in a

73 published dataset of 182 individuals of *Capsella grandiflora* sampled from a population

74 in Greece, and 191 *Drosophila melanogaster* flies sampled from an ancestral population

75 in Zambia (Lack et al. 2015). We find that mean signed LD is positive for most types of

76 mutations across the genome except for LOF mutations. The magnitude of positive LD

77 scales with the predicted deleteriousness of the mutations we analyze, with more

78 neutral mutations exhibiting the most positive LD. We use simulations to show that

79 admixture or distance biased mating could produce this type of pattern and provide

80 alternative explanations to epistasis for differences in LD among neutral versus

4

81     deleterious mutations. We then explore finer scale patterns of LD and uncover strong

82     short-scale positive LD among LOF mutations, a potential signal of within gene

83     antagonistic epistasis. Further analyses show that within gene LD is generally stronger

84     than between gene LD in *C. grandiflora* (correcting for distance between pairs of

85     mutations), for both neutral and deleterious mutations. This pattern is broadly consistent

86     with cross-over hotspots frequently occurring in promoter regions of plant genomes.

87     Finally, we use gene network information from KEGG to explore signals of LD and

88     epistasis among deleterious mutations segregating in functionally related genes. We

89     report no significant LD in *C. grandiflora* but significantly more negative LD in *D.*

90     *melanogaster* KEGG networks compared to a null distribution generated from permuted

91     networks, a pattern that could indicate synergistic epistasis acting against gene flow.
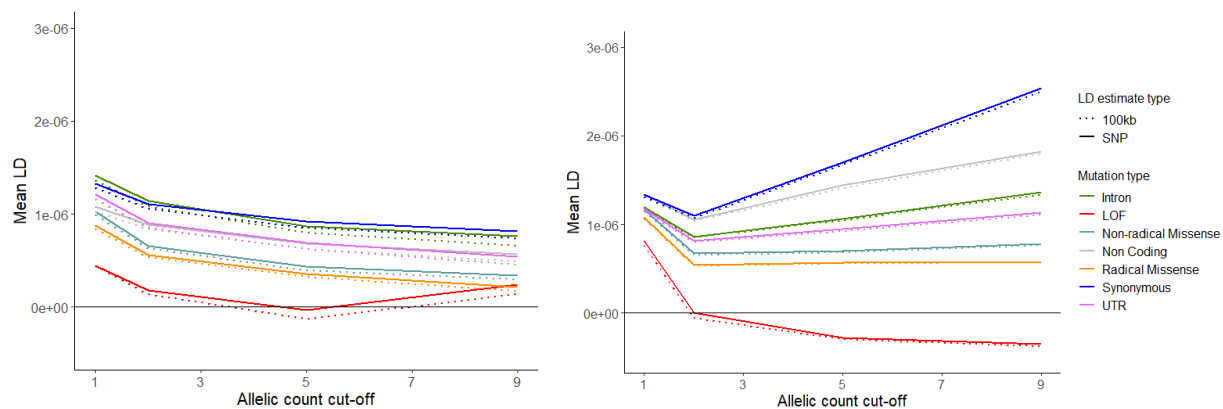
92

## Results

94     We analyzed patterns of signed LD among several classes of mutations (synonymous,

95     missense non-radical, missense radical, intronic, non-coding, UTR (untranslated

96     region), and LOF), in a dataset of 182 outbred diploid *C. grandiflora* individuals, and a

97     dataset of 191 *D. melanogaster* haploid embryos (population DPGP3). We only

98     considered variants below a specified threshold for minor allele frequency because we

99     hoped to maximize the probability that the rare variants at each site are deleterious,

100    though the degree of that deleteriousness is expected to vary among mutational classes

101    (e.g., for synonymous sites, rare variants are presumably negligibly more deleterious

102    than the common variant on average). The sign of LD was polarized by frequency so

103    positive/negative LD should indicate that deleterious variants are found more/less often

104    together than expected.

105        We first measured mean LD by assessing the over- or under-dispersion of

106    deleterious (or synonymous) variants among individual genomes (see Materials and

107    Methods; (Sohail et al. 2017)). An under-dispersion of the deleterious variants implies

108    negative LD (i.e., deleterious variants are found together less often than expected by

109    chance). We calculated mean LD per pair of alleles using several different allelic count

110    cut-offs (i.e., minor allele frequency thresholds). In both species, point estimates for

111 mean LD were positive for all classes of mutations, and all allelic count cut-offs

112 examined, except LOF mutations (Figure 1), where the point estimate for mean LD was

113 negative using some allelic count cut-offs but not others. When repeating this analysis

114 in *D. melanogaster,* excluding regions which were known to harbor inversions in the

115 DPGP3 population, we found the same qualitative results albeit with slightly reduced

116 positive LD for most mutations classes (Supplementary Figure 1).

117       One pattern apparent in our data is that the least deleterious mutational classes

118 exhibited the most positive mean LD in both flies and plants (i.e., the most positive LD

119 belonged to classes such as intronic and synonymous). The site frequency spectra for

120 these different mutational classes add support to the suspected rank ordering in the

121 deleteriousness of different mutational classes (Supplementary Figure 2) such that the

122 classes with the greatest excess of rare variants (presumed to be the most deleterious)

123 had the more positive LD. In *D. melanogaster* the order of deleteriousness inferred from

124 the site frequency spectra (starting with least deleterious) was as follows: synonymous,

125 non-coding, intronic, UTR, missense non-radical, missense radical, LOF. Similarly, in *C.*

126 *grandiflora* the order was synonymous, intronic, UTR, non-coding, missense non-

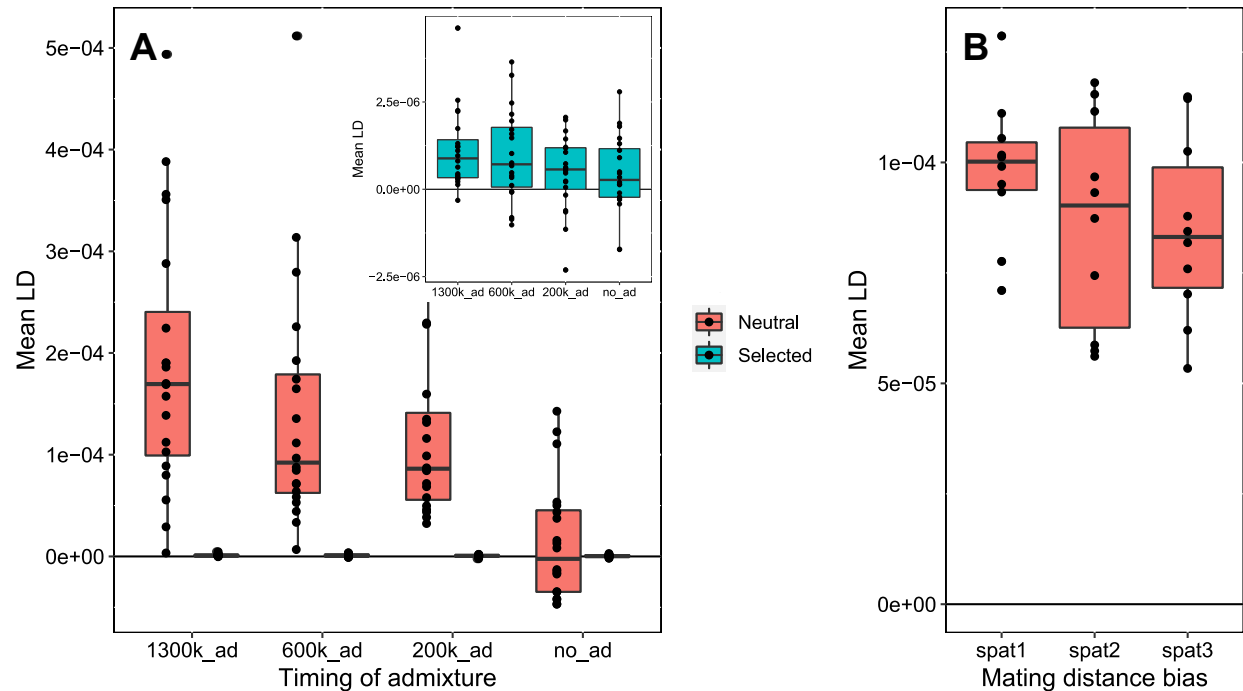127 radical, missense radical, LOF (Supplementary Figure 2).

128



130 Figure 1. Mean pair-wise LD among several classes of mutations across different allele count cut-offs.

131 Solid lines indicate mean LD among all SNPs, dashed lines indicate LD calculated among sites in

132 different 100kb, non-overlapping genomic blocks. Left, results for *C. grandiflora*, right, results for *D.*

133 *melanogaster.*

134    The observation that LD was strongest for neutral/nearly neutral mutations

135    suggests a non-selective force, such as admixture, is building LD (Sohail et al. 2017;

136    but see also Good 2020). We used a series of simulations using SLiM (V3.2.1) (Haller

137    and Messer 2019) to explore how different cases of non-equilibrium demography and

138    population structure can affect LD among neutral and deleterious mutations (see

139    Materials and Methods for more details). We first tested how a model of admixture

140    might impact patterns of LD for rare neutral and deleterious mutations under strictly

141    multiplicative selection. We simulated admixture between a focal population and two

142    previously isolated satellite populations and polarized LD by variant rarity. We found

143    that admixture easily caused positive LD to build up among neutral mutations,

144    particularly so if admixture started recently between populations that had previously

145    been isolated (Figure 2A). However, this was not the case for deleterious mutations in

146    these populations, where LD remained much closer to 0 albeit slightly positive on

147    average if admixture was present. This result was also apparent if we polarized LD by

148    true ancestral state in our simulations and imposed a minor allele frequency cut-off, or if

149    we polarized by frequency as in our real-world data, but did not implement a minor allele

150    frequency cut-off (Supplementary Figure 3). The only case where we did not observe

151    positive LD for neutral mutations was if we polarized LD by true ancestral state and

152    implemented no allele frequency cut-off (Supplementary Figure 3).

7

153

Figure 2. A) Mean signed LD among simulated neutral and deleterious mutations under different scenarios of admixture. The *x*-axis represents the generation in which admixture between isolated populations started. All simulations were run of a total of 1.5 million generations. Inset highlights results for selected (deleterious) mutations B) Mean LD among neutral mutations segregating in simulated populations existing on a 2D geographic landscape. The *x*-axis represents different scenarios of mating bias by distance with increasingly more random mating to the right of the *x*-axis.

We then explored isolation by distance due to continuous geography as a potentially common mechanism that could create positive LD in a similar way to admixture. Using SLiM to model populations on a continuous 2D landscape we again readily observed positive LD forming among neutral mutations under several scenarios of distance-biased mate choice, demonstrating that spatial considerations alone might be able to explain patterns of positive LD in our two datasets (Figure 2B).

Our simulation results qualitatively match the earlier simulation results reported by Sohail et al. who examined models specific to human demographic history (i.e., population structure and gene flow). They found positive LD does not build uniformly for deleterious and neutral mutations, rather, the more deleterious a class of mutations, the less positive LD built up among them. In summary, all of these simulations clearly show

8

171    that spatial structure with gene flow or admixture creates a difference in LD for selected

172    versus neutral sites.

173          While our patterns overall seem consistent with a relatively simple model of

174    spatial structure with varying strengths of purifying selection across site types, some of

175    our point estimates of LD for LOFs were negative, and negative LD is not expected

176    under such models of gene flow. Rather negative LD could be indicative of synergistic

177    epistasis or Hill-Robertson interference. To assess whether these processes might be

178    creating negative LD in our datasets, we next tested whether our estimates of negative

179    LD were significantly different from zero. We did this by permuting the assignment of

180    LOFs among all individuals in each dataset. This method preserves the allele frequency

181    at each locus while randomizing the associations among loci. We focused exclusively

182    on LOF mutations at an allele count cut off of no more than 5 because this cut-off

183    resulted in the most negative point estimates of mean LD in both datasets and such rare

184    mutations are more likely to be truly deleterious. All our subsequent analyses utilize this

185    allelic cut-off value for both datasets. This test suggested that LD among LOF mutations

186    was not significantly different from 0 in either *C. grandiflora* or *D. melanogaster* when

187    calculating LD SNP-by-SNP ($p = 0.996$ and $p = 0.386$, 2-tailed) or among sites in

188    different 100kb blocks ($p = 0.680$ and $p = 0.346$, 2-tailed). When we applied this

189    permutation approach to synonymous mutations, we found that LD was significantly

190    greater than 0 in both species, when calculating LD SNP-by-SNP ($p < 0.002$ both

191    species, 2-tailed), or using 100kb blocks ($p < 0.002$ both species, 2-tailed), further

192    verifying positive LD among more neutral mutations. Again, removing regions with

193    segregating inversions did not qualitatively change the results in *D. melanogaster* for

194    LOF mutations ($p = 0.658$, $p = 0.648$, LD calculated SNP-by-SNP and using 100kb

195    blocks respectively), or synonymous mutations ($p < 0.002$, for both types of LD

196    estimates).

197          In the preceding sections, we examined genome-wide average LD. However,

198    most pairs of sites contributing to this average are far apart or are found on different

199    chromosomes. For such sites, meiotic recombination and segregation will very rapidly

200    destroy any allelic associations formed by processes like selection. Significant signed
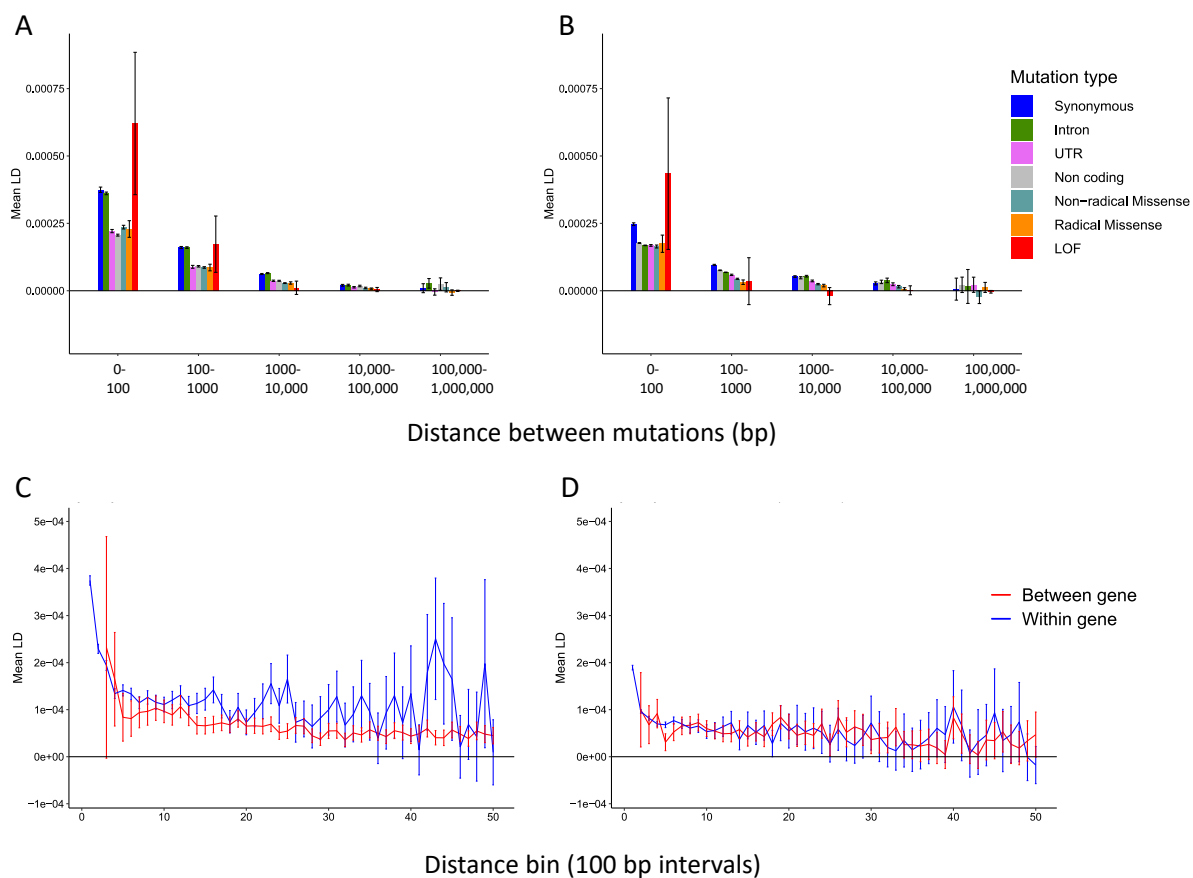
9

201  LD however, could still be present between mutations that are physically proximal. We

202  therefore next used PLINK (Purcell et al. 2007) to assess the relationship between inter-

203  mutation distance and LD for each class of mutations. Consistent with our first analysis,

204  LD was positive for all mutation classes in most distance bins; those estimates that

205  were negative were small in magnitude and were not significantly different from zero

206  (Figure 2). Within a distance bin, positive LD was stronger for the most weakly selected

207  mutation classes for most distance bins.

208  An interesting exception to this pattern in both species was in the smallest

209  distance bin (0-100bp). The major outlier in this distance bin were LOF mutations which

210  had surprisingly positive mean LD estimates in both species. The confidence intervals

211  on these estimates were very large for LOF mutations in this distance bin due to the

212  small number of observations for the mutation class. However, the high LD estimate for

213  LOFs is present in both species, and, in *C. grandiflora,* the 95% confidence intervals

214  suggested LOF mutations had more positive LD than all other mutation types aside from

215  intronic and synonymous. This pattern is consistent with intragenic antagonistic

216  epistasis, which seems probable for true LOF mutations occurring within the same

217  gene. Ideally, we would evaluate this hypothesis by comparing LD between physically

218  close LOFs that occur in the same versus different genes. However, we had too few

219  intergenic LOFs at short distances to do so.

220  Within gene antagonistic epistasis could also create positive LD among other

221  types of deleterious mutations such as missense mutations, which are much more

222  abundant. We compared signed LD decay within and between genes for both

223  synonymous and non-radical missense mutations (the two coding classes with ample

224  data) to test for this. In the case of *D. melanogaster,* we did not observe any major

225  differences in LD decay within vs. between genes for either mutational class (Figure 3D,

226  Supplementary Figure 5). In *C. grandiflora,* however, we observed significantly higher

227  LD for within gene pairs of mutations compared to between genes pairs for both non-

228  radical missense mutations and synonymous mutations (Figure 3C). Higher intra-gene

229  LD was also evident if we calculated unsigned ($r^2$) LD for *C. grandiflora* hinting at

230  potential differences in recombination leading to faster LD decay between genes rather

10

231    than LD created by epistasis (Supplementary Figure 6). Given that unsigned LD decay

232    should mostly be driven by the rate of recombination, we hypothesize this difference in

233    inter- vs intragenic LD is due to the strong enrichment of cross-overs in promoter

234    regions of plant genomes (Choi et al. 2013; Hellsten et al. 2013). Such crossovers

235    should rapidly erode LD between genes, while leaving within gene LD unaffected.

236    Conversely, no such pattern is known to occur in flies where transcription start sites

237    have actually been found to negatively correlate with cross-over occurrence (Comeron

238    et al. 2012; Smukowski Heil et al. 2015).
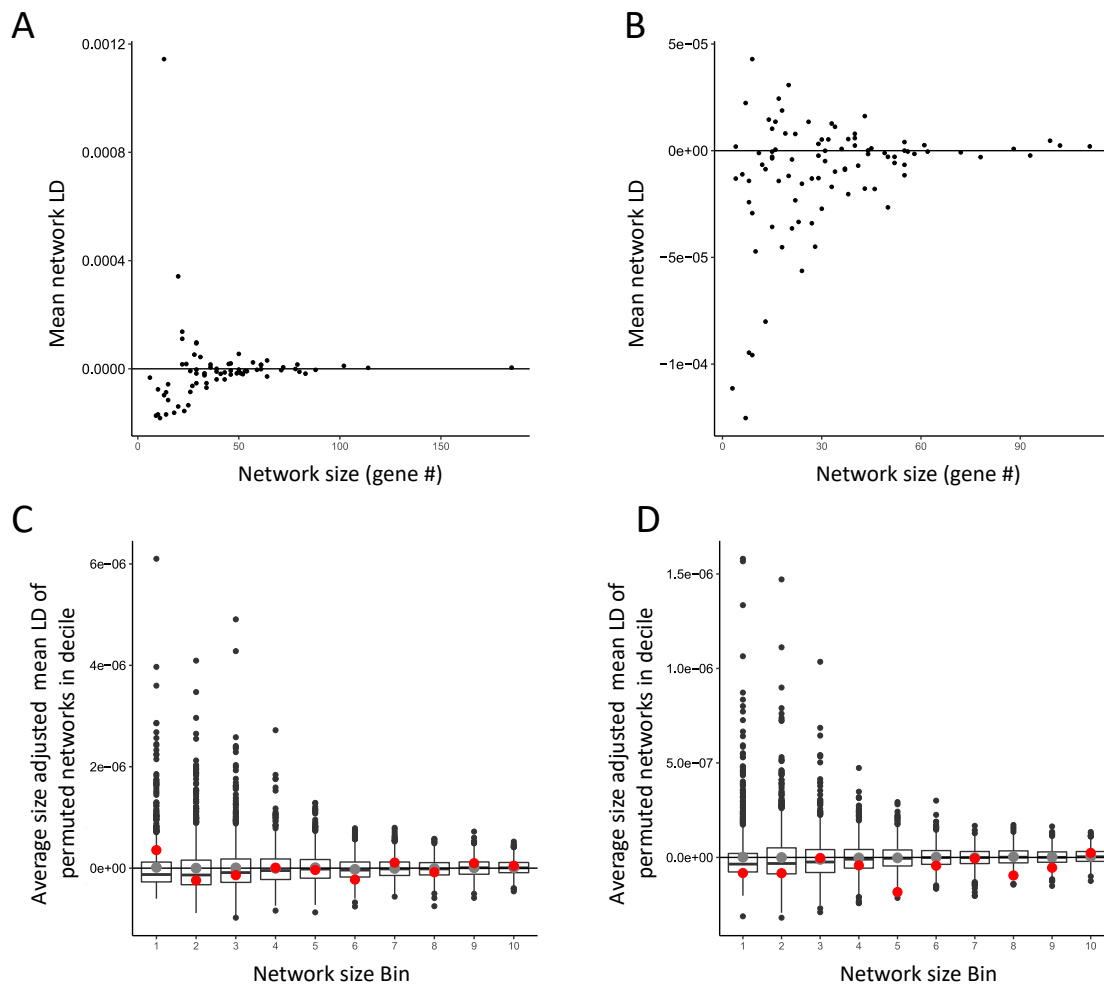
239



240

241    Figure 3. A,B. Distribution of mean signed LD for pairs of mutations across different distance bins for

242    several mutation classes in A) *C. grandiflora* and B) *D. melanogaster*. Mutations within each bin are

243    sorted by degree of expected deleteriousnes in ascending order. C,D Mean signed LD in 100bp bins for

244    synonymous mutation pairs within and between genes for C) *C. grandiflora* and D) *D. melanogaster*.

245

11

246    Though our previous analyses found no obvious signature of pervasive intergenic
247 synergistic epistasis when considering the entire genome, epistasis may be stronger
248 between functionally related genes. To investigate this possibility, we examined LD
249 among variants within interacting gene networks, using either radical missense or
250 synonymous mutations. We obtained gene lists of metabolic and signalling networks
251 from KEGG and only considered LD calculated among sites in different 100kb blocks to
252 minimize any contribution of LD between nearby mutations and to remove
253 measurements of LD between mutations within genes. We calculated the mean LD
254 within each network, then averaged these mean LD values across networks (weighting
255 by network size) to estimate "average network LD". We permuted the assignment of
256 genes to each KEGG network 1000x to create a null distribution for average network
257 LD.



258

12

259  Figure 4. Mean LD among radical missense mutations affecting genes within interacting biological

260  networks plotted against network size (as defined by numbers of genes within each network). LD was

261  calculated among sites in different 100kb blocks to minimize the effects of short intra-genic interactions.

262  Left data from *C. grandiflora*, right from *D. melanogaster.* C,D) Average network LD among radical

263  missense mutations for deciles based on network size; LD values were weighted by network size. Box-

264  plots show the null distribution for average network LD from permuted networks; black bar represents the

265  median, the grey point represents the mean, and whiskers represent quartiles. In each permutation,

266  networks were split into deciles based on bin size and the average LD of all networks in each decile was

267  calculated.  True average network LD of each decile is overlaid in red. Left data from *C. grandiflora*, right

268  from *D. melanogaster*

269

270  The mean LD of radical missense mutations within each network is shown in Figures

271  4A, B. Permutation tests indicated that the average network LD was significantly more

272  negative than expected in *D. melanogaster* (average network LD = -5.71E-08, p =

273  0.008) but not in *C. grandiflora* (average network LD = -1.05E-08, p = 0.956). Figures

274  3A, B give the appearance that LD is related to network size but this is likely a statistical

275  artifact that also occurs in permutations. To visualize this, we split networks into deciles

276  with respect to network size and calculated average network LD for each decile. The

277  permutation distributions had negative median values that approach zero for larger

278  network sizes (Figures 3C, D). Overlaying the observed values on these permutation

279  distributions helps visualize that the observed LD in *D. melanogaster* is more negative

280  than expected across most network sizes (red points are empirical values and grey

281  points are means from the distribution of permutation values). Repeating the

282  permutation analysis with synonymous mutations, we found that average network LD

283  was again significantly more negative than expected in *D. melanogaster* (average

284  network LD = -1.03E-08, p = 0.006) but not *C. grandiflora* (average network LD = 9.89E-

285  09, p = 0.902, see also Supplementary Figure 7 and Supplementary Table 2). Though

286  the point estimate of LD is more strongly negative for radical missense than

287  synonymous mutations in *D. melanogaster*, the qualitatively similar pattern complicates

288  the interpretation (see Discussion).

289

13

**Discussion**

In this study we analyzed patterns of signed LD in two species, *C. grandiflora* and *D. melanogaster*. When calculating mean LD among various classes of mutations, we found that less deleterious mutations tended to have more positive LD, with only LOF mutations exhibiting negative point estimates of mean LD under certain allelic-count cut-offs. Though the reduction in LD for deleterious classes such as LOF mutations relative to putatively neutral ones (e.g., synonymous mutations) could be interpreted as evidence of negative synergistic epistasis (Sohail et al. 2017) or Hill-Robertson interference (Garcia and Lohmueller 2020), other processes may provide more parsimonious alternatives. In particular, positive LD could be created by processes such as low level admixture in our datasets (Pfaff et al. 2001), and this effect may be weaker for more deleterious variants. For neutral sites, admixture can generate positive LD if LD is polarized either by rarity (as we have done) or by ancestral state if a minor allele threshold is imposed. Simulations across a range of demographic scenarios (both our own and those of Sohail et al.) have shown that positive LD builds up between mutations in a manner dependent on their selection coefficient; the more deleterious the mutations, the less positive LD builds up among them (Sohail et al. 2017), under a multiplicative model of negative selection. Presumably, the reason that positive LD occurs for low frequency neutral but less so for selected SNPs is as follows. Low frequency neutral SNPs within a given region will tend to be of two types: local variants of relatively recent origin but also migrant variants (of older origin), which will have come to the local population linked to migrant variants at other genomic sites (i.e., in positive LD). Deleterious variants are less likely to be of older (migrant) origin by virtue of the selection against them. Good (2020) showed that even without admixture, positive LD is expected between rare neutral mutations. This positive LD occurs because some variants that are rare in the present will have been more common in the past, providing an opportunity for a second variant to arise on the same haplotype. Positive LD is less likely to arise in this manner between deleterious variants because a deleterious variant is less likely to have been at higher frequency in the past. Though positive LD can arise in this fashion at both neutral and selected sites, admixture (including subtle forms of geographic structure) can potentially cause much stronger positive LD (Figure 2).

321   Part of our signal of genome-wide positive LD could be explained by the

322 presence of multinucleotide mutations (Schrider et al. 2011; Ragsdale 2021).

323 Multinucleotide mutations create strong positive LD among de-novo mutations, and

324 such coupled mutations should persist much longer if both variants are neutral,

325 potentially creating our observed pattern of an excess of positive LD for less deleterious

326 mutations. However, previous work in humans has suggested that the majority of SNPs

327 in multinucleotide mutations fall within 20bp of each other, which should create signed

328 LD on a much smaller scale than what we have observed in our data (Schrider et al.

329 2011). Our genome-wide measures of LD are not driven exclusively by nearby sites; the

330 LD measures are similarly positive even when we measure LD among sites in different

331 100kb blocks, thereby excluding the contribution of LD from the vast majority of

332 neighbouring sites (Figure 1).

333   The fact that differences in LD between selected and neutral sites can arise in

334 several simple models necessitates caution in interpretating differences in LD among

335 mutation classes with varying deleteriousness. For example, previous studies have

336 used LD among synonymous mutations as a control group for inferring synergistic

337 epistasis (Sohail et al. 2017) or Hill-Robertson interference (Garcia and Lohmueller

338 2020). However, as outlined above, differences in LD for deleterious versus neutral

339 mutations may be expected even under purely multiplicative selection, even without

340 invoking selective interference.

341   Because of its importance in theoretical population genetics (Kimura and

342 Maruyama 1966; Crow and Kimura 1970; Crow and Kimura 1979; Kondrashov 1982;

343 Barton 1995), we were particularly interested in looking for evidence of synergistic

344 epistasis in the form of negative LD at selected sites. Instead of comparing LD at

345 selected and neutral sites, we used randomization tests to test whether negative LD

346 among LOF mutations is significantly different from 0; it is not in either species. This

347 approach is somewhat conservative, because processes like admixture may oppose the

348 signal of negative LD created by synergistic epistasis. However, because the admixture

349 effect should be minimal for the most deleterious classes of mutations, this may not

350 pose a major limitation in searching for a signature of negative epistasis. The power of

15

351   recombination to destroy associations built by selection is likely a much more severe

352   limitation on synergistic epistasis—if it is common—creating a detectable signature on

353   genome-wide LD.

354       An additional issue with estimating mean LD across all genes is that this

355   averaging may hide meaningful variation. For example, epistasis between functional

356   sites within a gene may be fundamentally different in strength and/or sign than

357   intergenic epistasis. Moreover, physically close site pairs, which will often be intragenic,

358   will be less affected by recombination's power to destroy associations built by epistatic

359   selection or Hill-Robertson interference. We visualised the distribution of LD among

360   several classes of mutations in both datasets, split across bins of inter-mutation

361   distance. We observed non-zero LD most readily for nearby mutations across all

362   mutation classes, and in all cases it was significantly positive. Excluding the first

363   distance bin in our analysis (1-100bp), the magnitude of positive LD present in each

364   mutation class was predicted well by the expected deleteriousness of each type of

365   mutation. This is pattern can be explained by the simple scenarios of positive LD build-

366   up outlined above.

367       One notable deviation from the pattern of stronger positive LD for less

368   deleterious mutation classes was that, in first distance bin, LOF mutations had the most

369   positive point estimates of mean LD. We hypothesize that this pattern is due to within-

370   gene antagonistic epistasis, which is to be expected if a single LOF mutation is indeed

371   sufficient to knock out the function of a gene. This echoes similar findings from Puchta

372   et al. 2016 who demonstrated that antagonistic epistasis within a yeast snoRNA was

373   prevalent among large effect deleterious mutations occurring within conserved domains

374   because such mutations effectively acted as LOF variants and thus did not impact

375   fitness multiplicatively when combined with other deleterious mutations. Ragsdale

376   (2021) showed that LD for missense mutations within human protein functional domains

377   is significantly more positive than expected, also hinting at a potential signal of within

378   gene antagonistic epistasis.

379       Aside from within-gene epistasis, epistatic interactions may be stronger or more

380   frequent between mutations in functionally related genes. In particular, given that genes

16

381    function as part of larger biological networks, negative epistasis may arise between

382    deleterious mutations that affect the function of genes within the same networks (Chiu

383    et al. 2012). To test this idea, we calculated mean LD among synonymous and among

384    radical missense mutations present in genes within interacting biological networks

385    defined by KEGG (Kanehisa et al. 2016).  Permutation tests in *D. melanogaster*

386    suggested that the observed intra-network LD among radical missense mutations was

387    more negative than expected. Curiously, significantly negative network LD occurs for

388    synonymous mutations too. This latter result is surprising for two reasons: (i) LD is

389    (relatively) strongly positive for synonymous mutations at the genome-wide level (Figure

390    1), and (ii) negative epistasis should not affect (putatively neutral) synonymous sites. A

391    possible explanation of these findings emerges from our suspicion that the overall

392    genome-wide positive LD is due to processes of admixture and gene flow. The

393    significantly negative network LD for both synonymous and radical missense mutations

394    could be due to synergistic epistasis acting against introgressed alleles affecting the

395    same network. Because introgressed haplotypes will include synonymous and

396    missense mutations that are all in positive LD, selection on deleterious missense

397    variants will lead to a drop in positive LD for multiple types of mutations.

398          Unlike *D. melanogaster*, network LD was not significantly negative in *C.*

399    *grandiflora.* The lack of a significant result in *C. grandiflora* could be biologically

400    meaningful or more mundane. For example, KEGG network delineation could be more

401    biologically meaningful in *D. melanogaster* compared to *C. grandiflora* where network

402    information has been obtained from a species in a different genus (*Arabidopsis*

403    *thaliana*). Alternatively, the difference between species could simply be a statistical

404    artifact (i.e., false positive in *D. melanogaster* or false negative in *C. grandiflora*). Similar

405    analyses in other species will shed light on whether signed LD is related to network

406    status.

407          Our examination of LD has revealed variation in the strength and, in some cases,

408    the direction of signed LD. This variation is affected by several factors including

409    proximity of sites, putative deleteriousness of mutations, and the functional relationship

410    among genes. Some, but not all, of the patterns are consistent across two very different

17

411     species. Some of these patterns can be generated by more than one process and,

412     consequently, it will be challenging to conclusively prove which processes drive such

413     patterns. Nonetheless, patterns of LD can serve as one line of evidence for (or against)

414     particular hypotheses that are investigated using multiple approaches.

415

416     **Materials and Methods**

417     <u>Population genomics datasets</u>

418     We retrieved data from whole genome sequencing of 182 *C. grandiflora* individuals from

419     (Josephs et al. 2015) and data for 197 haploid *D. melanogaster* embryos from the

420     *Drosophila* population genomics project (DPGP3)(Lack et al. 2015). SNP calls

421     previously generated by Josephs et al. for *C. grandiflora* were provided by Tyler Kent

422     (personal communication). SNP calls for *D. melanogaster* were downloaded from the

423     PopFly website (Hervas et al. 2017, http://popfly.uab.cat/). Both data sets are a result of

424     thorough sampling from single populations with low population structure, making them

425     ideal candidates for detecting signs of epistasis from patterns of LD. To ensure that

426     recent migrants did not affect our LD analyses we used the R package SNPrelate (Li

427     2011) to visualize relatedness through PCA between *C. grandiflora* samples. This

428     revealed six divergent genotypes that we eliminated from our downstream analysis

429     leaving us with a total of 176 individuals. A previous study by Sohail et al. (2017) had

430     already used the DPGP3 dataset to analyze patterns of LD so we used the 190

431     individuals they retained after their filtering in our own analyses. We further filtered both

432     datasets by only considering bi-allelic sites where all individuals had genotype

433     information. Following Sohail et al. (2017) we removed SNPs segregating within chemo-

434     sensory and odorant binding genes in the *D. melanogaster* dataset based on gene lists

435     obtained from FlyBase (Larkin et al. 2021), though their inclusion has little effect on the

436     results. One final complication of the *D. melanogaster* dataset is the segregation of

437     several large-scale inversions in this species. The initial establishment of an inversion

438     creates some LD. However, gene exchange between chromosomes of different

439     inversion karyotypes still occurs within inverted regions via double cross-over

440     recombination events and gene conversion. Indeed, Houle and Márquez (2015) found

441    that LD was only slightly stronger within versus outside LD regions. To the extent

442    inversions cause a reduction in the effective recombination rate, inversions should

443    amplify the ability to detect the existing signal of non-zero LD built by other forces (e.g.,

444    selection, migration). Nonetheless, we repeated the majority of our analyses excluding

445    regions known to harbor inversions in the DPGP3 population. We obtained coordinates

446    of such inversions from Corbett-Detig and Hartl (2012) and removed SNPs segregating

447    in such regions for a subset of our analyses. However, analyses excluding inverted

448    regions are necessarily based on much less data and consequently have reduced

449    power.

450

451    SNP annotation

452    We used SNPeff (Cingolani et al. 2012) and the genome annotations of the reference

453    genomes (Slotte et al., 2013 for *Capsella rubella*; *D. melanogaster* release 5.57 from

454    Thurmond et al., 2019) to functionally annotate SNPs in both datasets as either LOF,

455    synonymous,  missense (non-synonymous), intronic (but not splice affecting), UTR (if

456    the SNP coordinate was either in the 5' or 3' UTR of a gene), or non-coding (for SNPs

457    not present in coding regions). A small number of SNPs had annotations in multiple

458    categories (e.g. both UTR and intronic), primarily due to multiple gene overlap, and

459    were excluded from the analysis. We included stop-gain and splice-disrupting SNPs in

460    our set of LOF mutations based on the method of Sohail et al. We also further classified

461    missense SNPs as either radical or non-radical. Missense SNPs were considered

462    radical if they changed both the volume and polarity of an amino acid based on previous

463    work suggesting that change in either category lead result in particularly deleterious

464    mutations in species such as *D. melanogaster* (Sainudiin et al. 2005; Weber and

465    Whelan 2019, see also Supplementary Table 1 for the list of amino acid properties we

466    used).

467

468    Calculating LD

19

469    We calculated LD values in two ways. First, we used the same method as Sohail et al.

470    by calculating a point estimate of average LD among all mutations. For a genome with $K$

471    loci, let $X_i$ be a discrete, random variable representing the number of derived alleles

472    present at locus $i$, which can take values 0, 1 for a haploid population or alternatively 0,

473    1, 2 for a diploid population. The variance in the total number of derived mutations

474    carried by each individual in the population can be expressed as:

475

$$Var\left(\sum_{i=1}^{K} X_i\right) = \sum_{i=1}^{K} Var(X_i) + 2\sum_{i,j}^{K} Cov(X_j, X_i)$$

476    Because LD is, by definition, a covariance in the allelic state between two loci, we can

477    use this equation to estimate the sum of all covariances across all loci by subtracting

478    first term of the right-hand side from the term on the left-hand side (and then dividing by

479    2). The term on the left-hand side represents the genome-wide variance in mutation

480    burden; the first term on the right-hand side is the sum of the variance in mutation

481    burden at each locus. We can then estimate a mean value of LD per pair of loci by

482    dividing by the number of possible two-way interactions in the dataset

483

$$mean\ LD = \frac{(Var(\sum_{i=1}^{K} X_i) - \sum_{i=1}^{K} Var(X_i))}{2\binom{K}{2}}$$

484    We also modified this approach to calculate LD on a block by block basis instead of

485    SNP by SNP. This measure of average LD largely eliminates LD between physically

486    close sites, which could initially arise via random mutation. We first split the genome into

487    100kb non-overlapping blocks. For a given genotype, we define $B_g$ as the number of

488    derived variants in block $g$. This new variable can take values from 0 to 2*(number of

489    segregating derived alleles in the given genomic block). To calculate total LD among all

490    blocks, we infer the covariance in mutation burden between all blocks as follows

491

$$Var\left(\sum_{g=1}^{W} B_g\right) = \sum_{g=1}^{W} Var(B_g) + 2\sum_{g,h}^{W} Cov(B_g, B_h)$$

492    where $W$ refers to the total number of 100kb blocks in the genome. Consider for

493    example the simple case where we compare two blocks $(B_g, B_h)$, each with two

494    segregating sites, $B$ can be represented as

495
$$B_g = X_1 + X_2 , B_h = X_3 + X_4$$

496    The number of covariance terms for these two genomic blocks is

497
$$Cov(B_g, B_h) = Cov(X_1, X_3) + Cov(X_1, X_4) + Cov(X_2, X_3) + Cov(X_2, X_4)$$

498    The within-block LD (e.g., $Cov(X_1, X_2)$ and $Cov(X_3, X_4)$) from physically neighbouring

499    sites contributes to the block-level variances (e.g., $Var(B_g)$ and $Var(B_h)$) but not the

500    between-block covariances. For an arbitrary number of blocks, $Cov(B_g, B_h)$ can

501    therefore be standardized per pair of interacting blocks as follows

502
$$mean\ LD_{blocks} = \frac{\left( Var\left(\sum_{g=1}^{W} B_g\right) - \sum_{g=1}^{W} Var\left(B_g\right)\right)}{2 \left(\sum_{g \neq h}^{W} \sum_{h}^{W} n_g n_h\right)}$$

503    where $n_g$ and $n_h$ represent the number of sites with segregating derived variants in

504    block $g$ and $h$ respectively.

505        We calculated mean LD using the above formula by transforming genotypes in

506    our VCF files into tables of non-reference allele counts (0, 1, 2 for *C. grandiflora* and 0,

507    1 for *D. melanogaster*) and calculating the relevant statistics in R using the package

508    matrixStats (Bengtsson 2017). We assumed that the non-reference alleles were the

509    derived alleles in the two datasets. In principle, a reference genome assembled from a

510    randomly sampled haplotype will contain some derived alleles that we will incorrectly

511    assume are ancestral in our method. This issue however should be minimal since our

512    analyses exclusively focus on rare mutations (<5% frequency) that are unlikely to be

513    included in a reference assembly and will be filtered out as high frequency variants by

514    our analysis even if they are included. This is especially true for most putatively

515    deleterious mutations such as LOF mutations which are likely maintained at low

516    frequency by mutation-selection balance.

21

517        We also calculated LD using PLINK (Purcell et al. 2007) for each category of

518    mutation. We calculated LD using default PLINK parameters which involved

519    subsampling LD observations as too many possible pairwise comparisons exist to

520    reasonably compute the entire distribution of LD values for most classes of mutations.

521    We estimated raw LD values by first estimating r between every single pair of mutations

522    in our dataset in PLINK (using the *--r* option) and then back-calculating a raw value of

523    LD by multiplying r by the square root of the product of allele frequencies at the two loci

524    being compared. This approach allows us to observe the entire distribution of LD values

525    rather than one summary statistic and back-calculating a raw value of LD from r allows

526    us to compare values from our two methods directly. Finally, we binned distance

527    between mutations pairs into seven categories: 100bp or less, 101-1000bp, 1001-

528    10,000bp, 10,001-100,000bp, 100,001-1,000,000bp to visualize how signed LD

529    decayed with distance for each class of mutations. Further, we compared signed LD

530    decay within vs. between genes for synonymous and non-radical missense mutations.

531    We did this by noting which gene our mutations of interest impacted according to

532    SNPeff, and splitting our LD values into two categories, those where both contributing

533    mutations occurred in the same gene, and those where both contributing mutations

534    occurred in different genes. We then visualized LD as above, however, we only

535    considered mutations 1-5000bp apart, and calculated mean LD in even 100bp bins,

536    excluding any bins with less than 100 pairs of LD values.

537

538    <u>Gene network analysis</u>

539    We used the R package Graphite (Sales et al. 2012) to obtain lists of genes from

540    biological pathways described in the KEGG database (Kanehisa et al. 2016). Network

541    information from KEGG was directly available for *D. melanogaster* but not *for C.*

542    *grandiflora* where we instead used network information from *A. thaliana*. We used

543    information on *C. grandiflora - A. thaliana* orthologs from (Josephs et al. 2015) to

544    generate lists of interacting genes in *C. grandiflora*. Due to the low number of LOF

545    mutations in each dataset we used low frequency (count of less than 5) radical

546    missense mutations (definition described in SNP annotation section) as our set of

547  candidate deleterious mutations. We calculated mean LD using 100kb blocks (as

548  described above) for each network defined by KEGG for our two species, generating

549  separating sets of networks for synonymous and radical mutations. Smaller networks

550  (defined by the number of genes assigned by KEGG to each network) have more highly

551  variable estimates of LD, presumably because of the smaller number of genes from

552  which LD is estimated. Consequently, we calculated "size adjusted LD" values for all

553  networks.  We did this by correcting mean LD in 100kb blocks for each network as

554  follows:

555  $$size\ adjusted\ LD\ =\ mean\ LD_{blocks}\ \times\frac{\#\ genes\ in\ network}{\#\ total\ genes\ across\ all\ networks}$$

556

557  <u>Simulations</u>

558  We used SLiM (V3.2.1) (Haller and Messer 2019) to run forward time simulations of

559  population admixture to ask how signed LD can be affected by various demographic

560  processes. We simulated three populations of 100,000 individuals each: one focal

561  population that was sampled at the end of the simulation and two satellite populations

562  with symmetrical migration to the focal population (10,000 individuals per generation).

563  Each diploid individual in our simulation contained two 1Mb chromosomes with

564  recombination and mutation rates both 1E-08 per bp per generation. Mutations were

565  sampled from two categories: neutral ($s$ = 0) with a probability 1%, or deleterious ($s$ = -

566  0.001) with a probability of 99%. Fitness was determined by the multiplicative effect of

567  deleterious load in each individual genome, dominance was also assumed to be

568  additive. We ran all simulations for 1.5 million generations altering the generation where

569  continuous admixture was started in several treatment groups: no admixture, admixture

570  starting at generation 200,000, 600,000, and 1,300,000. Each treatment group was

571  made of 20 simulated replicates. After 1.5 million generations, we sampled 100

572  individuals from the focal population in each replicate. Next, we filtered out recent

573  migrants in our focal population by performing a PCA on genotype of our samples and

574  eliminating individuals with PC values greater than 1 SD away from the mean of PC1 or

575  PC2. This mimics how we treated our real-world data where we eliminated outlier

23

576   samples using PCA. Next, to replicate how we defined ancestral/derived alleles in the

577   real-world data, we assigned all mutations with frequencies over 50% in our samples as

578   the ancestral variant. Finally, we filtered out sites with a 'derived' allele count over 5 and

579   calculated LD separately for neutral and deleterious mutations in each replicate. We

580   also separately calculated LD for these simulations keeping the true ancestral state

581   recorded by SLiM and polarizing LD by true ancestral/derived status both with and

582   without a minor allele count cut-off, mimicking the way LD may be calculated in a real-

583   world dataset where information on the true ancestral state may be available.

584        We ran a second set of simulations consisting of only one focal population where

585   individuals were placed on a 2D landscape to simulate the effects of isolation by

586   distance due to limited dispersal. We used the "Mate choice with a spatial kernel" recipe

587   provided in the SLiM manual for this set of simulations. Briefly, 10,000 individuals were

588   randomly placed on an ($x$,$y$) plane, with coordinate ranges [0,1] for both axes. To avoid

589   clumping, individual fitness was calculated as a function of spatial competition with

590   neighbouring individuals exerting the most costs to each other (see SLiM manual for

591   more details URL: http://benhaller.com/slim/SLiM_Manual.pdf). Individuals chose mates

592   a gaussian-distributed distance away, with mean 0, SD $\sigma$, and maximum value τ. We

593   ran simulations with three sets of parameter values for $\sigma$ and τ: (0.1,0.02), (0.3,0.06),

594   (0.5,0.5). This range of values was selected to explore various levels of bias towards

595   localized mating much like might occur in plant populations with limited pollen dispersal.

596   Finally, offspring dispersed a gaussian distance away from their first parent. Each

597   individual contained two 1Mb chromosomes containing only neutral mutations with a

598   recombination and mutation rate of 1E-08 per bp per generation. The simulations were

599   terminating after 100,000 generations and 100 individuals were sampled per simulation

600   replicate. Each mate choice condition was replicated 10 times. After sampling, LD was

601   calculated as described for the other simulations with the exception of PCA analysis as

602   no migrant filtering was necessary due to the absence of cross-population migration.

603

604   **Data availability statement**

605    Raw data for *C. grandiflora* are available at NCBI sequencing read archive (bio project

606    ID: PRJNA275635). Data for *D. melanogaster* were downloaded via the PopFly website

607    (Hervas et al. 2017, http://popfly.uab.cat/). Files containing annotated SNP calls (VCF

608    format) used in this study will be made publicly available upon acceptance of this

609    manuscript for publication. Scripts used in this study will also be made available at

610    https://github.com/gsan211.

611

## Acknowledgements

618

## Author Contributions

620    All authors designed the research. G.S. performed all analyses and wrote the draft

621    manuscript. All authors revised the manuscript.

622

## Competing Interests

624    The authors have no competing interests to declare

625

## Literature Cited

627  Agrawal AF, Whitlock MC. 2010. Environmental duress and epistasis: how does stress
628        affect the strength of selection on new mutations? *Trends Ecol. Evol.* 25:450–
629        458.

630  Auton A, McVean G. 2007. Recombination rate estimation in the presence of hotspots.
631        *Genome Res.* 17:1219–1227.

632    Bank C, Hietpas RT, Jensen JD, Bolon DNA. 2015. A Systematic Survey of an
633         Intragenic Epistatic Landscape. *Mol. Biol. Evol.* 32:229–238.

634    Barton NH. 1995. A general model for the evolution of recombination. *Genet. Res.*
635         65:123–145.

636    Chakraborty R, Weiss KM. 1988. Admixture as a tool for finding linked genes and
637         detecting that difference from allelic association between loci. *Proc. Natl. Acad.*
638         *Sci.* 85:9119–9123.

639    Chiu H-C, Marx CJ, Segrè D. 2012. Epistasis from functional dependence of fitness on
640         underlying traits. *Proc. R. Soc. B Biol. Sci.* 279:4156–4164.

641    Choi K, Zhao X, Kelly KA, Venn O, Higgins JD, Yelina NE, Hardcastle TJ, Ziolkowski
642         PA, Copenhaver GP, Franklin FCH, et al. 2013. *Arabidopsis* meiotic crossover
643         hot spots overlap with H2A.Z nucleosomes at gene promoters. *Nat. Genet.*
644         45:1327–1336.

645    Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden
646         DM. 2012. A program for annotating and predicting the effects of single
647         nucleotide polymorphisms, SnpEff. *Fly* 6:80–92.

648    Comeron JM, Ratnappan R, Bailin S. 2012. The Many Landscapes of Recombination in
649         *Drosophila melanogaster*. *PLOS Genet.* 8:e1002905.

650    Corbett-Detig RB, Hartl DL. 2012. Population Genomics of Inversion Polymorphisms in
651         *Drosophila melanogaster*. *PLOS Genet.* 8:e1003056.

652    Crow JF, Kimura M. 1970. An introduction to population genetics theory. Introd. Popul.
653         Genet. Theory Ney York: Harper & Row

654    Crow JF, Kimura M. 1979. Efficiency of truncation selection. *Proc. Natl. Acad. Sci.*
655         76:396–399.

656    Elena SF, Lenski RE. 1997. Test of synergistic interactions among deleterious
657         mutations in bacteria. *Nature* 390:395–398.

658    Good BH. 2020. Linkage disequilibrium between rare mutations.
659         https://www.biorxiv.org/content/10.1101/2020.12.10.420042v1.full

660    Garcia JA, Lohmueller KE. 2020. Negative linkage disequilibrium between amino acid
661         changing variants reveals interference among deleterious mutations in the
662         human genome
663         https://www.biorxiv.org/content/10.1101/2020.01.15.907097v1.full

664    Haller BC, Messer PW. 2019. SLiM 3: Forward Genetic Simulations Beyond the Wright–
665         Fisher Model. *Mol. Biol. Evol.* 36:632–637.

26

666  Hellsten U, Wright KM, Jenkins J, Shu S, Yuan Y, Wessler SR, Schmutz J, Willis JH,
667      Rokhsar DS. 2013. Fine-scale variation in meiotic recombination in *Mimulus*
668      inferred from population shotgun sequencing. *Proc. Natl. Acad. Sci.* 110:19478–
669      19482.

670  Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. *Genet.*
671      *Res.* 8:269–294.

672  Hill WG, Robertson A. 1968. Linkage disequilibrium in finite populations. *Theor. Appl.*
673      *Genet.* 38:226–231.

674  Houle D, Márquez EJ. 2015. Linkage Disequilibrium and Inversion-Typing of the
675      *Drosophila melanogaster* Genome Reference Panel. *G3 Genes Genomes Genet.*
676      5:1695–1701.

677  Josephs EB, Lee YW, Stinchcombe JR, Wright SI. 2015. Association mapping reveals
678      the role of purifying selection in the maintenance of genomic variation in gene
679      expression. *Proc. Natl. Acad. Sci.* 112:15390–15395.

680  Kanehisa M, Sato Y, Kawashima M, Furumichi M, Tanabe M. 2016. KEGG as a
681      reference resource for gene and protein annotation. *Nucleic Acids Res.*
682      44:D457–D462.

683  Kimura M, Maruyama T. 1966. The Mutational Load with Epistatic Gene Interactions in
684      Fitness. *Genetics* 54:1337–1351.

685  Kondrashov AS. 1982. Selection against harmful mutations in large sexual and asexual
686      populations. *Genet. Res.* 40:325–332.

687  Kondrashov AS. 1995. Dynamics of unconditionally deleterious mutations: Gaussian
688      approximation and soft selection. *Genet. Res.* 65:113–121.

689  Lack JB, Cardeno CM, Crepeau MW, Taylor W, Corbett-Detig RB, Stevens KA, Langley
690      CH, Pool JE. 2015. The Drosophila Genome Nexus: A Population Genomic
691      Resource of 623 *Drosophila melanogaster* Genomes, Including 197 from a
692      Single Ancestral Range Population. *Genetics* 199:1229–1241.

693  Lalić J, Elena SF. 2012. Magnitude and sign epistasis among deleterious mutations in a
694      positive-sense plant RNA virus. *Heredity* 109:71–77.

695  Larkin A, Marygold SJ, Antonazzo G, Attrill H, dos Santos G, Garapati PV, Goodman
696      JL, Gramates LS, Millburn G, Strelets VB, et al. 2021. FlyBase: updates to the
697      Drosophila melanogaster knowledge base. *Nucleic Acids Res.* 49:D899–D907.

698  Li H. 2011. A statistical framework for SNP calling, mutation discovery, association
699      mapping and population genetical parameter estimation from sequencing data.
700      *Bioinforma. Oxf. Engl.* 27:2987–2993.

701 McEvoy BP, Powell JE, Goddard ME, Visscher PM. 2011. Human population dispersal
702     "Out of Africa" estimated from linkage disequilibrium and allele frequencies of
703     SNPs. *Genome Res.* 21:821–829.

704 McVean G. 2007. The Structure of Linkage Disequilibrium Around a Selective Sweep.
705     *Genetics* 175:1395–1406.

706 Pfaff CL, Parra EJ, Bonilla C, Hiester K, McKeigue PM, Kamboh MI, Hutchinson RG,
707     Ferrell RE, Boerwinkle E, Shriver MD. 2001. Population structure in admixed
708     populations: effect of admixture dynamics on the pattern of linkage
709     disequilibrium. *Am. J. Hum. Genet.* 68:198–207.

710 Puchta O, Cseke B, Czaja H, Tollervey D, Sanguinetti G, Kudla G. 2016. Network of
711     epistatic interactions within a yeast snoRNA. *Science* 352:840–844.

712 Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MAR, Bender D, Maller J, Sklar
713     P, de Bakker PIW, Daly MJ, et al. 2007. PLINK: A Tool Set for Whole-Genome
714     Association and Population-Based Linkage Analyses. *Am. J. Hum. Genet.*
715     81:559–575.

716 Ragsdale AP. 2021. Can we distinguish modes of selective interactions using linkage
717     disequilibrium?
718     bioRxiv:2021.03.25.437004.

719 Sainudiin R, Wong WSW, Yogeeswaran K, Nasrallah JB, Yang Z, Nielsen R. 2005.
720     Detecting Site-Specific Physicochemical Selective Pressures: Applications to the
721     Class I HLA of the Human Major Histocompatibility Complex and the SRK of the
722     Plant Sporophytic Self-Incompatibility System. *J. Mol. Evol.* 60:315–326.

723 Sales G, Calura E, Cavalieri D, Romualdi C. 2012. graphite - a Bioconductor package to
724     convert pathway topology to gene network. *BMC Bioinformatics* 13:20.

725 Schrider DR, Hourmozdi JN, Hahn MW. 2011. Pervasive Multinucleotide Mutational
726     Events in Eukaryotes. *Curr. Biol.* 21:1051–1054.

727 Slotte T, Hazzouri KM, Ågren JA, Koenig D, Maumus F, Guo Y-L, Steige K, Platts AE,
728     Escobar JS, Newman LK, et al. 2013. The *Capsella rubella* genome and the
729     genomic consequences of rapid mating system evolution. *Nat. Genet.* 45:831–
730     835.

731 Smukowski Heil CS, Ellison C, Dubin M, Noor MAF. 2015. Recombining without
732     Hotspots: A Comprehensive Evolutionary Portrait of Recombination in Two
733     Closely Related Species of Drosophila. *Genome Biol. Evol.* 7:2829–2842.

734 Sohail M, Vakhrusheva OA, Sul JH, Pulit SL, Francioli LC, Consortium G of the N,
735     Initiative ADN, Berg LH van den, Veldink JH, Bakker PIW de, et al. 2017.
736     Negative selection in humans and fruit flies involves synergistic epistasis.
737     *Science* 356:539–542.

738  Stephens JC, Briscoe D, O'Brien SJ. 1994. Mapping by admixture linkage disequilibrium
739      in human populations: limits and guidelines. *Am. J. Hum. Genet.* 55:809–824.

740  Thurmond J, Goodman JL, Strelets VB, Attrill H, Gramates LS, Marygold SJ, Matthews
741      BB, Millburn G, Antonazzo G, Trovisco V, et al. 2019. FlyBase 2.0: the next
742      generation. *Nucleic Acids Res.* 47:D759–D765.

743  Weber CC, Whelan S. 2019. Physicochemical Amino Acid Properties Better Describe
744      Substitution Rates in Large Populations. *Mol. Biol. Evol.* 36:679–690.

745

746

747

748

749

750

751

752

753

754

755

756