1   **Deleterious Mutations Accumulate Faster in Allopolyploid than**

2   **Diploid Cotton (*Gossypium*) and Unequally Between Subgenomes**

3   Justin L. Conover[1], Jonathan F. Wendel[1]

4

5   [1] – Department of Ecology, Evolution, and Organismal Biology, Iowa State University,

6   Ames, IA, 50011, USA

7

8   ORCID ID:

9   JLC: 0000-0002-3558-6000

10  JFW: 0000-0003-2258-5081

11

12  Corresponding author: Jonathan F Wendel

13  Mailing address: Department of EEOB, 251 Bessey Hall, 2200 Osborn Dr, Ames, IA

14  50011

15  Phone Number: (515) 294-7172

16  Email:jfw@iastate.edu

17

18  Key words: polyploidy, deleterious mutations, purifying selection, molecular evolution

19

20 **Abstract**

21 Whole genome duplication (polyploidization) is among the most dramatic mutational

22 processes in nature, so understanding how natural selection differs in polyploids relative

23 to diploids is an important goal. Population genetics theory predicts that recessive

24 deleterious mutations accumulate faster in allopolyploids than diploids due to the

25 masking effect of redundant gene copies, but this prediction is hitherto unconfirmed.

26 Here, we use the cotton genus (*Gossypium*), which contains seven allopolyploids

27 derived from a single polyploidization event 1-2 million years ago, to investigate

28 deleterious mutation accumulation. We use two methods of identifying deleterious

29 mutations at the nucleotide and amino acid level, along with whole-genome

30 resequencing of 43 individuals spanning six allopolyploid species and their two diploid

31 progenitors, to demonstrate that deleterious mutations accumulate faster in

32 allopolyploids than in their diploid progenitors. We find that, unlike what would be

33 expected under models of demographic changes alone, strongly deleterious mutations

34 show the biggest difference between ploidy levels, and this effect diminishes for

35 moderately and mildly deleterious mutations. We further show that the proportion of

36 nonsynonymous mutations that are deleterious differs between the two co-resident

37 subgenomes in the allopolyploids, suggesting that homoeologous masking acts

38 unequally between subgenomes. Our results provide a genome-wide perspective on

39 classic notions of the significance of gene duplication that likely are broadly applicable

40 to allopolyploids, with implications for our understanding of the evolutionary fate of

41 deleterious mutations. Finally, we note that some measures of selection (e.g. dN/dS,

42 $\pi_N/\pi_S$) may be biased when species of different ploidy levels are compared.

2

43

**Introduction**

45 Genome duplication (polyploidy) is among the most dramatic mutational processes in

46 nature, causing myriad saltational changes at the cellular and organismal levels (Doyle

47 and Coate 2019; Bomblies 2020; Fernandes Gyorfy et al. 2021), and is associated with

48 consequential phenomena ranging from crop domestication (Renny-Byfield and Wendel

49 2014; Qi et al. 2021) to cancer progression (Matsumoto et al. 2021). Polyploidy is

50 especially common in the angiosperms, with all extant species having experienced at

51 least one or more polyploidy events during their evolutionary history (Jiao et al. 2011),

52 and at least 30% of currently recognized species having a polyploidy event in the recent

53 past (One Thousand Plant Transcriptomes Initiative 2019).

54 Novel evolutionary patterns created by polyploidy at the genic (e.g.

55 neofunctionalization, subfunctionalization, and loss (Kuzmin et al. 2021)) and genomic

56 (e.g., homoeologous recombination (Mason and Wendel 2020)) levels have been well

57 documented across taxa, including the frequent asymmetry of these responses with

58 respect to co-resident genomes in a polyploid nucleus. Nonetheless, many questions

59 remain regarding the effects of natural selection on polyploid relative to diploid genomes

60 (Baduel et al. 2019; Monnahan et al. 2019) and the interplay between these novel

61 evolutionary patterns and the long-term trajectories of genome evolution (Qi et al. 2021)

62 following polyploidization (e.g. biased fractionation).

63 One of the earliest predictions about natural selection in polyploids relative to

64 diploids is that putatively deleterious mutations may accumulate faster due to the

65 masking effect of completely or partially recessive deleterious mutations in duplicated

3

66    genes (Haldane 1932; Hill 1970; Bever and Felber 1992). Only recently, however, have

67    these predictions begun to be evaluated in young allopolyploids such as *Arabidopsis*

68    *kamchatica (Paape et al. 2018)* and *Capsella bursa-pastoris (Douglas et al. 2015;*

69    *Kryvokhyzha, Salcedo, et al. 2019; Kryvokhyzha, Milesi, et al. 2019)*, and autotetraploid

70    *Arabidopsis arenosa* (Monnahan et al. 2019). Because the number of deleterious

71    mutations is strongly influenced by shifts in demography and mating system (Brandvain

72    and Wright 2016), which may coincide with polyploid formation (Grant 1981; Barringer

73    2007), a clear link between ploidy level and the accumulation of deleterious mutations is

74    challenging to demonstrate in natural polyploid populations.

75        The cotton genus (*Gossypium*) represents one of the best studied allopolyploid

76    systems (Wendel and Grover 2015; Hu et al. 2021). The genus includes approximately

77    45 currently recognized diploid species classified into eight genome groups (A-G, and

78    K), and seven allopolyploid species resulting from a single (Grover et al. 2012)

79    allopolyploidization event ~1-2 million years ago between members of the A and D

80    genome groups (Wendel 1989). Although the most closely related extant species of

81    these two progenitor diploids are found in southern Africa and Northern Peru,

82    respectively, the polyploids are only found in the Americas (Figure 1). Most wild

83    populations, including those of the two independently domesticated species *G. hirsutum*

84    (AD$_1$) and *G. barbadense* (AD$_2$), occur in small, isolated populations on islands or in

85    coastal regions. Subsequent to their initial domestication 4,000 - 8,000 years ago in the

86    Yucatan Peninsula (AD$_1$) and NW S. America (AD$_2$), respectively, the ranges of the two

87    domesticated species rapidly expanded to encompass much of the American tropics

4

88 and subtropics and then spread globally with the rise of the international cotton fiber

89 trade (Yuan et al. 2021).

90       Here we describe the evolutionary trajectory of deleterious mutations in two wild

91 diploid and six wild allopolyploid cotton species (all descended from a single

92 allopolyploidization event), with a focus on how allopolyploidization and speciation have

93 shaped the number and genomic distribution of deleterious mutations. We use two

94 methods to estimate the strength of selection at the amino acid and nucleotide level and

95 show support for a nearly century-old hypothesis that polyploids accumulate mutations

96 faster than their diploid progenitors. We demonstrate that, in agreement with this

97 hypothesis, polyploidy has the greatest influence on strongly, rather than moderately or

98 mildly, deleterious mutations. We also find that deleterious mutations accumulate

99 asymmetrically between the two co-resident subgenomes in the allopolyploid nucleus,

100 indicating that these masking effects may act unequally between the subgenomes. In

101 total, our results support theoretical predictions that allopolyploidy can lead to a faster

102 rate of deleterious mutation accumulation through masking of recessive deleterious

103 variants, and that the relationship of the rate of deleterious mutation accumulation

104 between subgenomes and their progenitor diploids is complex, even when comparing

105 identical pairs of single-copy homoeologs among lineages.

106

107 **Results**

108 *Patterns of Synonymous and Nonsynonymous Mutations*

109 To investigate patterns of deleterious mutations, we viewed our data at three

110 phylogenetic depths: SNPs segregating within the global phylogenetic tree (Figure 2A-

111    D), SNPs that emerged since the divergence of each subgenome from its respective

112    diploid progenitor (Figure 2E-H), and SNPs that are still variable within the polyploids

113    (Figure 2I-L). Each group is a subset of the previously described group. We restricted

114    our analyses to a set of 8,884 single-copy, syntenically conserved homoeologous pairs

115    of genes (17,768 genes in total) that showed no evidence of gene loss, gene copy

116    variation, tandem duplication, ambiguous read mapping, homoeologous exchange, or

117    homoeologous gene conversion (See Methods; Supplementary Figure 1). Notably, the

118    patterns described below are largely reflected in genome-wide patterns as well

119    (Supplementary Figure 2), indicating that filtering criteria did not bias overall results, and

120    that, in cotton, homoeologous interactions have minimal effects on subgenome-specific

121    SNP patterns (Supplementary Figure 1).

122         Using the curated set of 8,884 pairs of homoeologous genes, we found no

123    evidence for differences in the rate of synonymous mutation accumulation in diploids

124    versus polyploids at any phylogenetic depth (Figure 2A, 2E), although differences can

125    be found within the polyploid species (Figure 2I), with *G. mustelinum* (AD$_4$, Orange) and

126    *G. darwinii* (AD$_5$, Yellow) having consistently lower rates than the rest of the clade, and

127    in both subgenomes. When viewing SNPs that have accumulated since the divergence

128    of the earliest polyploid lineage (Figure 2AI), there is an asymmetry between

129    subgenomes in the rate of synonymous site changes, with the Dt ("t" denoting

130    "tetraploid") subgenome containing a moderately higher number of synonymous

131    mutations than the At subgenome for all species. This difference potentially indicates a

132    higher mutation rate or relaxed background selection in genic regions of the Dt

133    subgenome compared to homoeologous genic regions of the At sugenome, and is

134    consistent with previous analysis finding that genes in the Dt subgenome are evolving

135    faster than genes in the At subgenome in five allopolyploid cottons (Chen et al. 2020).

136        In contrast to the relative homogeneity in rates of synonymous substitution

137    among diploids and polyploids, rates of nonsynonymous mutation accumulation differed

138    significantly at all phylogenetic depths. Notably, at the deepest phylogenetic depth

139    (Figure 2B), estimates for the number of derived nonsynonymous mutations in the

140    diploids *G. herbaceum* ($A_1$, Red) and *G. raimondii* ($D_5$, Black) did not differ between

141    subgenomes, indicating that any mapping biases or erroneous SNP calls in these

142    samples were removed by our SNP filtering criteria. *Gossypium raimondii* ($D_5$)

143    contained more derived nonsynonymous mutations than did *G. herbaceum* ($A_1$), and

144    this lineage-specific difference was reflected in the Dt and At subgenomes as well.

145    When lineage-specific effects that arose from the long, shared ancestry between the

146    subgenomes and their progenitor diploids were removed (Figure 2F), a clear distinction

147    between the rates of nonsynonymous mutations between diploids and their respective

148    subgenomes in the allopolyploids becomes apparent. In all polyploids, the At

149    subgenome contained between 25-58% more nonsynonymous mutations than *G.*

150    *herbaceum* ($A_1$, Red), and the Dt subgenome contained 17-36% more than *G. raimondii*

151    ($D_5$, Black). These results demonstrate that even though the rates of synonymous

152    mutation accumulation did not differ significantly between the diploids and polyploids,

153    polyploidy significantly increases the rate of nonsynonymous substitution accumulation.

154    Finally, when restricting our attention to only those mutations that have arisen following

155    polyploid formation (Figure 2J), the lineage-specific patterns observed for

156    nonsynonymous mutations were largely identical to the patterns of synonymous

7

157    mutations. For example, *G. mustelinum,* (AD$_4$, Orange) consistently had the lowest

158    number of derived mutations in both subgenomes. Notably, however, the Hawaiian

159    Island endemic *G. tomentosum* (AD$_3$, Purple) has a higher number of derived

160    nonsynonymous mutations than expected based on the patterns of synonymous

161    mutations, potentially reflecting the population bottleneck associated with its origin

162    following long-distance dispersal to the Hawaiian Islands (see Discussion). In summary,

163    polyploidy significantly enhances the rate of nonsynonymous mutation accumulation in

164    all *Gossypium* allopolyploids, and does so asymmetrically across co-resident genomes.

165

166    *Polyploidy Increases Rate of Deleterious Mutation Accumulation*

167    Because the fitness effects of most nonsynonymous mutations can vary widely, from

168    neutral to lethal, we asked if the elevated rate of nonsynonymous mutations observed in

169    polyploid *Gossypium* reflects an increase in neutral or nearly-neutral nonsynonymous

170    mutations, or if instead this elevation is attributable to a greater accumulation of

171    deleterious mutations. To address this, we used two approaches to estimate whether a

172    mutation was deleterious: BAD_Mutations and GERP++. BAD_Mutations  performs a

173    likelihood ratio test from a gene-specific multi-species alignment to determine if a

174    mutation at a particular nonsynonymous site is deleterious, while GERP++ uses a

175    genome-wide multiple sequence alignment (i.e. agnostic to genic regions) to estimate

176    the degree of conservation at a particular site in the genome (including noncoding and

177    synonymous sites). Notably, because one of the hallmark long-term processes following

178    polyploidy is pseudogenization (Wendel 2015), recently pseudogenized sequences may

179    still display some degree of conservation across the multiple sequence alignment, but

8

180    may not be inherently deleterious. Therefore, to avoid inflating estimates of deleterious

181    mutations in polyploids compared to diploids, we used GERP only within the exonic

182    regions of the 8,884 homoeologs. Additionally, while GERP can score the degree of

183    deleteriousness of a mutation, BAD_Mutations can only classify variants into deleterious

184    or not deleterious. Therefore, the values shown in Figure 2DHL represent the sum of

185    the allele frequencies of derived deleterious mutations, similar to the values for Figure

186    2AEI and Figure 2BFJ. For analysis with GERP, we used the GERP load, which

187    incorporates the deleteriousness of each variant into the score, summing the frequency

188    of each derived allele multiplied by it's GERP score (see (Rodgers-Melnick et al. 2015;

189    Wang et al. 2017)).

190         As shown in Figure 2, both of the foregoing analyses demonstrate that

191    deleterious mutations accumulate in polyploids in a manner similar to nonsynonymous

192    mutations, suggesting that the difference in nonsynonymous sites cannot be wholly

193    attributed to putatively neutral or nearly-neutral alleles.  For example, there is

194    remarkable consistency in the patterns of deleterious mutations that have accumulated

195    since the divergence of the diploid from its respective diploid progenitor in both the

196    count of nonsynonymous substitutions (Figure 2F), the GERP load (Figure 2G), and the

197    number of deleterious mutations (Figure 2H). In all three columns, the diploids show

198    fewer accumulated alleles than the polyploids, *G. tomentosum* (AD$_3$, Purple) shows the

199    highest number of all the polyploids, and *G. mustelinum* (AD$_4$, Orange) shows the

200    fewest of all the polyploids.

201         An interesting pattern arises when comparing estimates of the GERP load

202    (Figure 2G) and number of deleterious mutations (Figure 2H) between diploids and their

9

203     closely related subgenomes: while the total number of deleterious mutations in the At

204     subgenomes was 52-99% higher in the polyploids than the diploids (Figure 2H), the

205     GERP load in the polyploids was only 13-42% higher (Figure 2G). Similar patterns were

206     found in the Dt subgenome, with 34-66% more deleterious mutations in the polyploids

207     than the diploid, but only a 9-13% increase in GERP load. This discrepancy could reflect

208     inherent differences in the types of sequences used and how deleteriousness is

209     quantified between the two methods, suggesting that the use of multiple analytical tools

210     for detection of genetic load may yield more nuanced insights than either method on its

211     own (See Discussion).

212

213     *Asymmetries in the Rate of Deleterious Mutation Accumulation*

214     Although deleterious mutations are accumulating faster in polyploids relative to diploids,

215     it is not obvious whether this increased rate is different from the increased rate of

216     accumulation of nonsynonymous mutations. To test this, we compared, among ploidy

217     levels, the total proportion of nonsynonymous mutations that were considered

218     deleterious by BAD_Mutations (Figure 3). For SNPs that originated since the

219     divergence of the A and D diploids (Figure 3A), the proportion of nonsynonymous sites

220     that are deleterious is roughly 2% higher in polyploids than in diploids, despite the

221     shared evolutionary history of more than 4 million years between each subgenome and

222     their respective diploid progenitors. Notably, as similarly shown in Figure 2, the

223     proportion of nonsynonymous mutations that are inferred to be deleterious in both

224     diploids is equivalent when mapped to either subgenome, indicating that our filtering

10

225  criteria did not differentially exclude deleterious or non-deleterious SNPs with respect to

226  which subgenome the diploid reads were mapped.

227      At shallower phylogenetic depths (Figure 3B), the difference between diploids

228  and polyploids becomes even clearer, with polyploids exhibiting 3-4% higher

229  proportions of deleterious SNPs in the Dt subgenome and 5-12% higher in the At

230  subgenome than their respective diploid progenitors. The most unbiased and

231  straightforward comparison of the asymmetry in strength of purifying selection between

232  the two subgenomes of allopolyploid cottons is provided by mutations that have

233  occurred following polyploidization (Figure 3C). Here, the At subgenome of all species

234  contain a 2-3% high proportion of deleterious SNPs than the Dt subgenome, indicating

235  that differences exist in the strength of purifying selection between the two

236  homoeologous subgenomes that have resided in the same nucleus for over a million

237  years. This pattern is also observed when deleterious SNPs are mapped onto the

238  phylogeny (Supplementary Figure 3). Additionally, there is more variation among

239  species in the At subgenome than in the Dt subgenome, although the patterns in this

240  respect are not simple. The amount of subgenomic asymmetry is smallest in *G. darwinii*

241  (AD$_5$, Yellow) from the Galapagos Islands, and largest in the Brazilian endemic and

242  inland species *G. mustelinum* (AD$_4$, Orange), indicating that asymmetries between

243  subgenomes of the same species may vary within a single clade of allopolyploids.

244

245  *Disentangling Demography and Selection from Effects of Ploidy*

246  Demography is a potential confounding factor in estimating the rate of deleterious

247  mutation accumulation. Shifts in demography are known to complicate inferences of the

11

248     strength of selection and genetic load (Brandvain and Wright 2016); for example, even

249     in one of the best studied demographic shifts, the Out of Africa migration in humans,

250     several papers (Lohmueller et al. 2008; Gazave et al. 2013; Simons et al. 2014; Henn et

251     al. 2016; Simons and Sella 2016) have reached seemingly contradictory conclusions on

252     whether genetic load has increased as a result of these shifts in demography (but see

253     (Lohmueller 2014)). The pattern of deleterious mutation accumulation has also been

254     well-documented in bottlenecks and population growth associated with domestication in

255     crops such as maize (Wang et al. 2017), soybean and barley (Kono et al. 2016),

256     sorghum (Lozano et al. 2021), cassava (Ramu et al. 2017), and rice (Liu et al. 2017).

257          Polyploidy is typically associated with a population bottleneck (Grant 1981;

258     Barringer 2007), but because the genetic diversity of both the diploid and polyploid

259     species in this study is low (Table 1), demographic modeling of the depth or duration of

260     population bottlenecks and range expansion following polyploid formation is not straight-

261     forward. Generalized patterns of the effects of demography on deleterious mutations,

262     however, can serve as a null expectation to test if our data follows the same trends

263     observed under varying demographic scenarios, as explained in the following.

264          Demographic shifts, including population bottlenecks and expansions, have a

265     large influence on the accumulation of deleterious mutations. According to the nearly

266     neutral theory (Ohta 1992), the fate of deleterious mutations is determined by genetic

267     drift instead of selection when the selection coefficient (s) of deleterious mutations is

268     less than or equal to $1/(2N_e)$, where $N_e$ is the effective population size. The reduction of

269     $N_e$ during a population bottleneck would therefore allow weakly deleterious mutations  to

270     escape purifying selection (i.e. to behave as if they were neutral), while strongly

12

271 deleterious mutations with a selective coefficient greater than $1/(2N_e)$ would still be

272 removed by purifying selection. On the other hand, as $N_e$ increases during population

273 expansion, mutations that are mildly deleterious are expected to be more efficiently

274 purged from the population.

275      In both demographic scenarios, we expect that mildly or moderately deleterious

276 mutations would be most differentially affected, while strongly deleterious mutations

277 would consistently be removed by purifying selection. Based on this theory, if the

278 differences in the number of deleterious mutations we see between diploids and

279 polyploids are due to demography, then we would expect to see most of that difference

280 reflected in mildly, rather than strongly, deleterious mutations. In contrast, if masking of

281 deleterious alleles in polyploids is driving a higher rate of accumulation relative to

282 diploids, this pattern will not be observed.

283      To test if our data were consistent with changes in demography, we first asked if

284 there was a correlation between the degree of deleteriousness of a mutation (as

285 measured by GERP) and its relative increase in the polyploids compared to the diploids.

286 To answer this question, we plotted the relative change of deleterious mutations in each

287 subgenome relative to its most closely related diploid progenitor. We plotted this relative

288 change for three different degrees of deleteriousness - strongly deleterious mutations (4

289 < GERP ≤ 6), moderately deleteriousness (2 < GERP ≤ 4), and mildly deleteriousness

290 (0 < GERP ≤ 2) deleterious (Figure 4). We found that in both subgenomes of all six

291 polyploids, when comparing SNPs that had originated after the divergence of the diploid

13

292    from its respective subgenome in the allopolyploids, strongly deleterious mutations

293    accumulated at a faster rate relative to diploids than did moderately or mildly deleterious

294    mutations, which is inconsistent with expectations under a demographic change model

295    alone. We also observed this change under both an additive and recessive model of

296    dominance (Supplementary Figure 5). In total, the rate of accumulation among

297    mutations with different inferred degrees of deleteriousness do not suggest that the

298    patterns we see can be explained solely by demographic changes, but that the masking

299    effect of duplicated genes may play an important role in the determining the fate of

300    deleterious mutations in allopolyploids.

301

302    **Discussion**

303    *Effects of Polyploidy on Deleterious Mutation Accumulation*

304    One of the earliest hypotheses regarding mutation accumulation in allopolyploids dates

305    back to Haldane (Haldane 1932) where he posits that in allopolyploids, "one gene may

306    be altered without disadvantage, provided its functions can be performed by a gene in

307    one of the other sets of chromosomes." Allopolyploids are therefore predicted be able to

308    tolerate a higher mutational load than their diploid relatives, and putatively deleterious

309    mutations may accumulate faster in polyploids than in their diploid relatives due to the

310    masking effect of recessive or incompletely dominant deleterious alleles. Here, we

311    demonstrate that these predictions are true in allopolyploid cottons. All polyploids in

312    *Gossypium* harbor more mutations at phylogenetically conserved sites than do their

313    closest diploid progenitors, as determined by two different methods of detecting

314    deleterious mutations. We also find that the proportion of all nonsynonymous mutations

14

315  that are inferred to be deleterious is higher in polyploids than in their diploid progenitors

316  and that polyploidy has the greatest effect on strongly deleterious (and, inferentially,

317  more recessive (Eyre-Walker and Keightley 2007; Huber et al. 2018)) mutations. Thus,

318  using the power of comparative phylogenetics and genomics combined with analytical

319  methods for detection of deleterious mutations, we demonstrate confirmation of a nearly

320  century old hypothesis regarding natural selection in allopolyploid organisms.

321

322  *Demography Alone Cannot Explain Patterns of Deleterious Mutations in Polyploids*

323  Estimating the strength of natural selection and genetic load is notoriously challenging

324  (Lohmueller 2014) and is complicated by shifts in effective population size (including

325  bottlenecks and expansions), mating systems, and effective recombination rates,

326  among other life-history and demographic factors (Brandvain and Wright 2016). Here

327  we illuminate an additional relevant consideration, i.e., whole genome duplication. Yet

328  many of the considerations for populations that are not in demographic equilibrium also

329  apply to *Gossypium*. Diversification in the cotton tribe (*Gossypieae*) has been

330  characterized by numerous long-distance dispersal events (Grover et al. 2017),

331  including the one from Africa to the Americas 1-2 MYA that led to the evolution of

332  allopolyploid *Gossypium*. We note that in the Hawaiian Islands endemic *G.*

333  *tomentosum*, the total number of synonymous substitutions is not significantly different

334  from the rest of the polyploids, but the number of nonsynonymous and deleterious

335  mutations is significantly increased, suggesting that the genetic bottleneck associated

336  with island dispersal has elevated the number of deleterious mutations compared to the

337  rest of the polyploids.

15

338      While demographic changes upon polyploid formation have been shown to

339    change the number and frequency of deleterious mutations in other systems (Douglas

340    et al. 2015; Paape et al. 2018; Baduel et al. 2019; Kryvokhyzha, Salcedo, et al. 2019;

341    Kryvokhyzha, Milesi, et al. 2019), we show here that the patterns of mutation

342    accumulation in *Gossypium* cannot be explained by demography alone, and that the

343    data are more consistent with the nearly century-old hypothesis that recessive

344    deleterious mutations can accumulate faster in allopolyploids due to the masking effect

345    of duplicated genes and lack of recombination between subgenomes (Haldane 1932).

346    Specifically, we show that strongly (and, hence, more recessive (Morton et al. 1956;

347    Mukai et al. 1972; Eyre-Walker and Keightley 2007; Agrawal and Whitlock 2011; Huber

348    et al. 2018)) deleterious mutations accumulate faster in polyploids compared to diploids

349    than moderately or mildly deleterious mutations, and that this pattern is inconsistent with

350    demographic shifts or long-term change in population size (Figure 4, Supplementary

351    Figure 5).

352

353    *Asymmetry in Subgenomes in the Distribution of Deleterious Mutations*

354    One of the elegant attributes of a clade of allopolyploid genomes derived from a single

355    polyploidization event is that they offer a remarkable natural experiment for comparing

356    subgenomes that have resided within the same nucleus for, in the case of *Gossypium*,

357    approximately 1.5 million years. Once an allopolyploid is established, each subgenome

358    is subjected to identical external or population-level factors, including demography,

359    mating systems, and environmental and ecological conditions, as well as internal

360    cellular processes, including identical DNA replication and recombination machinery.

361    These features remove many of the confounding factors that may influence the genetic

362    load and provide a simple comparative context for revealing evolutionary forces that

363    might differentially affect co-resident genomes or homoeologs.

364         An unexpected finding of our analyses is the striking asymmetry in the proportion

365    of all nonsynonymous mutations that are inferred to be deleterious between the two

366    subgenomes of all allopolyploid species in *Gossypium*. We found that the At

367    subgenome of all species contains 2-3% more nonsynonymous mutations that are

368    inferred to be deleterious (Figure 3) even when only considering mutations that have

369    arisen following the earliest allopolyploid diversification events, and correcting for

370    removing the biases of unequal phylogenetic distances to each subgenome's model

371    progenitor diploid. Our work adds to a growing recognition that the two co-resident

372    subgenomes in cotton allopolyploids may be shaped asymmetrically by evolutionary

373    processes, including interspecific introgression and selection under domestication

374    (Fang, Wang, et al. 2017; Fang, Guan, et al. 2017; Chen et al. 2020; Yuan et al. 2021),

375    and that this phenomenon also extends to other important allopolyploid crop plants,

376    including wheat (Pont and Salse 2017; Jiao et al. 2018) and *Brassica (Tong et al. 2020)*.

377         Teasing apart the genesis of differential subgenomic responses to selection is

378    rendered challenging by several factors independent of phylogeny. We note, for

379    example, the relevant example of the recently formed allopolyploid *Capsella bursa-*

380    *pastoris* and its diploid progenitors, where consistent asymmetries in genetic load are

381    reported between the subgenomes (Kryvokhyzha, Salcedo, et al. 2019; Kryvokhyzha,

382    Milesi, et al. 2019) the differences likely reflect the dramatically different mating systems

383    of the progenitors, in which the subgenome with the higher genetic load originated from

17

384    an obligate outcrosser, *C. grandiflora* (Ne = 800,000), whereas the subgenome with the

385    lower genetic load derives from the predominantly selfing *C. orientalis* (Ne =

386    5000)(Douglas et al. 2015). In another recently formed (20-250 thousand years ago)

387    allopolyploid, *Arabidopsis kamchatica*, no asymmetry in the distribution of fitness effects

388    between subgenomes was found, although it was observed that each subgenome of the

389    allopolyploid contained more neutral and fewer deleterious alleles than either of the

390    diploid progenitors (Paape et al. 2018). It is unclear, however, whether this shift was

391    due to allopolyploidy *per se* or if it reflects the transition from an obligate outcrossing to

392    a mating system with some degree of inbreeding, with a concomitant purging of partially

393    or completely recessive deleterious alleles, as shown in several other systems

394    (Arunkumar et al. 2015; Roessler et al. 2019). In *Gossypium*, all species have similar

395    mating systems and a canonical outcrossing floral morphology including highly exserted

396    styles and stigmas. Population sizes often are small, however, likely leading to relatively

397    high levels of generalized inbreeding. At present, however, no data exist that address

398    these considerations.

399

400    *Polyploidy, Redundancy, and Fitness Effects*

401    One possible interpretation of our results is that *Gossypium* polyploids are less fit than

402    their closely related diploid progenitors because they harbor more deleterious mutations

403    in their genomes, especially mutations that have already been driven to fixation. We

404    note that an additional possibility is that mutations in polyploids that occur at

405    phylogenetically conserved sites may not actually have a deleterious effect on fitness as

406    they do in diploids. Inferring the genetic load of a population simply by counting the

407    number of deleterious variants assumes that all alleles contribute independently to the

408    total genetic load of a population. However, because of the functional overlap of

409    duplicated genes and, in most cases, absence of recombination between

410    homoeologous chromosomes in an allopolyploid, a recessive deleterious mutation can

411    never be present in all four copies of a gene and thus may be invisible to selection

412    because of the masking effect of its homoeologous partner.

413        An important takeaway from this study is that recessive deleterious mutations in

414    allopolyploids, at least at some loci, may actually accumulate in a manner more similar

415    to neutral mutations, presumably because of the lack of recombination between

416    subgenomes and, hence, the inability of purifying selection to "see" the negative effects

417    of these mutations. Because these recessive deleterious mutations escape the effects

418    of purifying selection, many traditional tests for detecting selection (e.g. dN/dS, $\pi_N/\pi_S$)

419    may be biased when comparing a polyploid to diploid because the polyploid would be

420    expected to accumulate putatively deleterious sites more quickly (and maintain a higher

421    genetic diversity at nonsynonymous sites) than their diploid relatives.

422        Another important implication of this finding is that allopolyploidy (or gene

423    duplication in general) may play an important and underrecognized role in determining

424    how selection acts on new mutations, notwithstanding the burgeoning literature on fates

425    of duplicate gene evolution (Conant et al. 2014; Shi et al. 2020; Veitia and Birchler

426    2021). The evolutionary trajectory of new mutations will largely be dependent on the

427    selection coefficient (s) acting on that locus, and the dominance coefficient (h), defined

428    as the proportion of the fitness cost that a mutation harbors when in a heterozygous

429    state. In allopolyploids, however, the evolutionary fate of new mutations may be

430   determined not only by allelic dominance at that locus, but also by the interaction arising

431   from the coexistence of its homoeologous locus, a term we call "homoeologous epistatic

432   dominance". The relationships between this homoeologous epistatic dominance, allelic

433   dominance, and the selection coefficient are likely complicated and potentially heavily

434   influenced by other biological considerations such as biased expression of homoeologs,

435   sub- or neofunctionalization, and homoeologous recombination, among others.

436   Moreover, notwithstanding these polyploidy-specific effects, even the genome-wide

437   relationships between two of these factors, allelic dominance and the selection

438   coefficient, have only been modeled using genomic data in the past few years (Huber et

439   al. 2018).

440        Nonetheless, understanding how this homoeologous epistatic dominance

441   impacts the fitness effects of new mutations is an unexplored aspect of polyploid

442   genome evolution, and it is not yet clear whether this will equally affect advantageous

443   and deleterious variants. How homoeologous epistatic dominance operates with respect

444   to functional properties arising from considerations such as gene balance (Veitia and

445   Birchler 2021), dosage effects (Conant et al. 2014), structural and functional

446   entanglement (Kuzmin et al. 2020; Kuzmin et al. 2021), and inter-subgenomic *cis- and*

447   *trans-* effects (Bottani et al. 2018; Hu and Wendel 2019) would seem to represent

448   important avenues for understanding how natural selection operates differently in

449   polyploids compared to diploids. From an applied perspective, these insights could be

450   important in agriculture, particularly because so many of our most important crop plants

451   have a recent history that includes polyploidy (Renny-Byfield and Wendel 2014), and

20

452    segregating patterns of genome fractionation have the potential to serve as targets of

453    selection in crop improvement (Hufford et al. 2021).

454

455    **Materials and Methods**

456    *Plant Materials and Sequencing*

457    We used whole genome sequencing data from 46 individuals in *Gossypium*, including

458    between two and ten individuals from each of eight species. Included in our sampling

459    was six polyploid species originating from a single polyploidization event 1-2 million

460    years ago (Wendel 1989; Hu et al. 2021), two diploid species representing models of

461    the genome donors to the allopolyploids (A and D), and three species from Australia

462    that served as outgroups for polarizing SNPs into ancestral and derived states. These

463    sequences were previously described (Yuan et al. 2021), and SRA codes for all 46

464    resequenced individuals are listed in Supplemental Table 1*. For *G. hirsutum*, we

465    randomly chose ten accessions that were classified in the "Wild" population from Yuan

466    et al. (Yuan et al. 2021), and for the other species, we chose all accessions available

467    that did not show evidence of being mislabeled, as determined by a PCA plot of the

468    SNPs called.

469        After the data were downloaded from NCBI, adapter sequence removal and

470    quality score filtering of FASTQ reads was performed using Trimmomatic v0.36 (Bolger

471    et al. 2014) using the parameters "LEADING:28 TRAILING:28 SLIDINGWINDOW:8:28

472    SLIDINGWINDOW:1:10 MINLEN:65 TOPHRED33". Trimmed reads from each polyploid

473    sample were mapped to the 26 chromosomes of the *G. hirsutum* reference genome

474    (Saski et al. 2017), and reads from each diploid sample were mapped to each

21

475     subgenome separately to avoid competitive mapping of the diploid reads against a

476     tetraploid reference genome. Reads from the three outgroup species were separately

477     mapped to both subgenomes to ensure that reads were not filtered out for mapping to

478     multiple parts of the genome. All mapping was done using bwa-mem v0.7.17 (Li and

479     Durbin 2009) and only uniquely mapping paired reads (-F 260 flag) that were mapped in

480     their proper orientation (-f 2 flag) were retained using Samtools v1.9 (Li et al. 2009)

481     before the files were sorted and converted to bam files. Using the Sentieon (Kendig et

482     al. 2019) SNP Calling program, gVCF files were generated, and joint genotyping was

483     performed using the GVCFtyper algorithm (see Github repository for full scripts). SNP

484     filtering was performed using GATK v4.0.4.0 using the filter expression "QD < 2.0 || FS

485     > 60.0 || MQ < 40.0 || SOR > 4.0 || MQRankSum < -12.5 || ReadPosRankSum < -8.0".

486     For each species (excluding the outgroup species, and treating *G. stephensii,* and *G.*

487     *ekmanianum* as a single species), we nullified any SNP call in which all individuals were

488     heterozygous to remove any collapsed genomic region in the reference genome or

489     paralogous regions that were not present in the reference genome. We treated *G.*

490     *stephensii* and *G. ekmanianum* as a single species because we only sampled two

491     individuals of *G. stephensii*, so removing any sites in with both individuals were

492     heterozygous errantly removed real SNPs that were not due to paralogy mapping

493     issues. All scripts for generating and filtering SNP calls are located on our GitHub

494     repository (https://github.com/conJUSTover/Deleterious-Mutations-Polyploidy).

495

496     *Identification of Homoeologs*

22

497    We used the pSONIC pipeline (Conover et al. 2021) to identify syntenically conserved

498    homoeologs in the *G. hirsutum* reference genome, and kept only homoeologous pairs

499    that were less than 5% different in their total annotated CDS length. To remove

500    homoeologous pairs that may have experienced homoeologous exchange events

501    (though there is scant evidence for this (Salmon et al. 2010; Flagel et al. 2012; Chen et

502    al. 2020)), we removed any pair in which the proportion of the reads from the two

503    progenitor diploid genomes (termed At and Dt in the allopolyploid, the "t" indicating

504    "tetraploid") did not meet the expected 2:2 ratio. Average read depth of CDS regions

505    was determined by bedtools2 v.2.27.1 (Quinlan and Hall 2010). Briefly, for a single

506    homoeologous pair, we calculated the average read depth of the At homoeolog divided

507    by the sum of the average read depth of both homoeologs and removed any

508    homoeologous pair in which this fraction was less than 37.5 or greater than 62.5. We

509    expect any HEs that result in a 0:4 At:Dt copy number to contain 0% At reads/total

510    reads; HEs that result in 1:3 At:Dt copy number should have a 25% At reads/total reads;

511    HEs that result in a 3:1 At:Dt copy number should have a 75% At reads/total reads; HEs

512    that result in a 4:0 At:Dt copy number should have a 100% At reads/total reads; and no

513    HE (i.e. 2:2 At:Dt copy number) would result in a 50% At reads/total reads. We used the

514    midpoints between the "No HE" and the 1:3 and 3:1 copy numbers as cutoff points. This

515    filtering resulted in 8,884 homoeologous pairs (17,768 genes) being analyzed further.

516        Non-reciprocal homoeologous exchanges (i.e. homoeologous gene conversion)

517    could also bias the estimates of the genetic load in a way that is not related to new

518    mutation following polyploidization or speciation. To control for positions in these non-

519    HE homoeologs that may be influenced by gene conversion, we linked SNP positions

520    between homoeologs in the following way. We first performed pairwise alignments of

521    the CDS sequences using MACSEv2 (Ranwez et al. 2011; Ranwez et al. 2018), which

522    aligns CDS sequences in accordance with their translated amino acid sequences, but

523    allows for the possibility of frameshift mutations. We then used the aligned CDS

524    sequences to identify where indels were present, and found the corresponding genomic

525    positions for every nucleotide in the alignment, inserting gaps where indels occurred.

526    We then extracted the genomic positions for each SNP position as well as the genomic

527    position for its aligned nucleotide. We retained only those homoeologous SNP positions

528    in which both positions had a confidently called ancestral allele (described above) and

529    in which the ancestral allele matched between the two homoeologs. Importantly, for

530    homoeologs that were encoded in opposite orientations in the reference genome (i.e.

531    one homoeolog was encoded on the forward strand of the reference genome, and the

532    other homoeolog was encoded on the reverse complement), we ensured that the

533    inferred ancestral states for the two SNP positions included both nucleotides of a

534    purine/pyrimidine pair (e.g. the ancestral state for homoeologous SNP was "A" while the

535    ancestral state of the other homoeologous SNP was "T"). We also removed any pair of

536    homoeologous SNPs in which more than 2 alleles were present (while similarly treating

537    homoeologous pairs encoded in opposite directions as described in the previous

538    sentence).

539        In total, we only used those SNP sites that: (A) did not link to an indel in its

540    homoeologous pair, (B) were biallelic and had consistently inferred ancestral states in

541    the two subgenomes, (C) the derived allele was found in only one of the two

24

542    subgenomes or their respective diploid progenitors, and (D) the derived allele was fixed

543    in a diploid and segregating in its respective subgenome (or vice-versa).

544

545    *Quantifying Deleterious Mutations*

546    We used GERP++ (Davydov et al. 2010) to identify regions of the genome that are

547    evolutionarily conserved, using whole genome alignments from 11 genomes spanning

548    the Eudicots (Supplementary Table 2). Species were chosen if they contained

549    chromosome-level assemblies publicly available on Phytozome or NCBI, and if all

550    documented whole genome duplication events in each species' evolutionary history is

551    also shared by *Gossypium* (e.g. the *Arabidopsis thaliana* genome was not chosen

552    because it has experienced at least one independent WGD event since its divergence

553    from *Gossypium*). Genomes were aligned to the *G. hirsutum* reference genome using

554    the LASTZ/MULTIZ approach used by the UCSC genome browser. Briefly, genomes

555    were masked using Repeatmasker using a custom repeat library enriched with

556    *Gossypium* TEs (Grover et al. 2017). Each query genome was aligned to each of the *G.*

557    *hirsutum* reference chromosomes separately. These alignments were chained together

558    using axtChain, and the best alignment was found using ChainNet. These alignment

559    files were converted into fasta files using the roast program from the MULTIZ package.

560        Using these genome alignments, we used the gerp++ package (Davydov et al.

561    2010) to calculate GERP scores for every position in the genome. First, we used 4-fold

562    degenerate sites in all genomes to calculate a neutral-rate evolutionary tree, which was

563    calculated using RAxML (Stamatakis 2014). We then used the gerp++ package to

564    estimate the GERP score at every position in the genome, but importantly, we excluded

25

565   the *G. hirsutum* reference genome from the alignment  to avoid biasing sites in the

566   reference genome that may be deleterious. Because the gerp++ program ignores gaps

567   in the reference genome, we used custom R scripts to enter dummy variables in the

568   gapped regions of the GERP score so the number of GERP scores equaled the total

569   number of nucleotide positions in each chromosome. Scripts for each step above are

570   available on Github (link here). To calculate the genetic load across linked sites, we

571   used the GERP load (i.e. the sum of the derived allele frequency times the GERP score

572   for each SNP site) as described in (Wang et al. 2017) and (Rodgers-Melnick et al.

573   2015). All scripts for generating the multiple sequence alignments and GERP scores

574   can found in our GitHub repository (https://github.com/conJUSTover/Deleterious-

575   Mutations-Polyploidy)

576        Secondly, we used the BAD_Mutations (Kono et al. 2016; Kono et al. 2018)

577   pipeline to perform LRT tests on conserved amino acid substitutions sites.

578   Nonsynonymous substitutions were identified using SNPEff (Cingolani et al. 2012) and

579   statistical significance was determined using a Bonferroni correction with 967,155

580   missense mutations to correct for multiple testing. Every step of the BAD_Mutations

581   pipeline was performed using the dev branch of the github repository (accessed July 13,

582   2020). Species included in the calculation of deleterious mutations are included in

583   Supplementary Table 3, with the notable absence of Gossypium raimondii since it was

584   sampled as part of this project.

585        We used the GERP load (sum of the allele frequencies * GERP score) (Wang et

586   al. 2017) and the BAD_Mutations load (sum of the allele frequencies of all statistically

587   significant deleterious mutations) as a summary of the genetic load present in each

26

588    genome at different phylogenetic depths. The BAD_Mutations load may be interpreted

589    as the average number of deleterious alleles expected in each individual of a

590    population, but it does not differentiate between severity of deleteriousness (as does

591    GERP load). We also used GERP to classify SNPs into mildly deleterious (0<GERP≤2),

592    moderately deleterious (2<GERP≤4), and strongly deleterious (4<GERP≤6). Scripts for

593    generating the whole-genome alignments for GERP are located on our GitHub

594    repository (https://github.com/conJUSTover/Deleterious-Mutations-Polyploidy).

595

596    *Rate of Deleterious Mutations Along the Phylogeny of Gossypium*

597    To determine if there was a bias in the rate of deleterious mutation accumulation

598    between the two subgenomes, we used homoeologous SNPs in which the derived allele

599    showed a parsimony-informative position between the two subgenomes of

600    allopolyploids and the two diploid progenitors (identified by the green bars in

601    Supplemental Figure 1).

602

603    *Genetic Diversity*

604    For each species, $\pi$ for the 17,768 high quality gene CDS sequences (8,884

605    homoeologous pairs) was calculated on a site-wise basis using vcftools (Danecek et al.

606    2011). To find the total PI across all genes, we summed the total sitewise pi values and

607    divided by the total length of the concatenated CDS sequences, removing any positions

608    which did not have a null SNP call in the VCF file.

609

610    **Acknowledgements**

**References**

616

617 Agrawal AF, Whitlock MC. 2011. Inferences about the distribution of dominance drawn
618     from yeast gene knockout data. *Genetics* 187:553–566.

619 Arunkumar R, Ness RW, Wright SI, Barrett SCH. 2015. The evolution of selfing is
620     accompanied by reduced efficacy of selection and purging of deleterious mutations.
621     *Genetics* 199:817–829.

622 Baduel P, Quadrana L, Hunter B, Bomblies K, Colot V. 2019. Relaxed purifying
623     selection in autopolyploids drives transposable element over-accumulation which
624     provides variants for local adaptation. *Nat. Commun.* 10:5818.

625 Barringer BC. 2007. Polyploidy and self-fertilization in flowering plants. *Am. J. Bot.*
626     94:1527–1533.

627 Bever JD, Felber F. 1992. The theoretical population genetics of autopolyploidy. *Oxford*
628     *surveys in evolutionary biology* 8:185–185.

629 Bolger AM, Lohse M, Usadel B. 2014. Trimmomatic: a flexible trimmer for Illumina
630     sequence data. *Bioinformatics* 30:2114–2120.

631 Bomblies K. 2020. When everything changes at once: Finding a new normal after
632     genome duplication. *Proc. Biol. Sci.* 287:20202154.

633 Bottani S, Zabet NR, Wendel JF, Veitia RA. 2018. Gene Expression Dominance in
634     Allopolyploids: Hypotheses and Models. *Trends Plant Sci.* 23:393–402.

635 Brandvain Y, Wright SI. 2016. The Limits of Natural Selection in a Nonequilibrium
636     World. *Trends Genet.* 32:201–210.

637 Chen ZJ, Sreedasyam A, Ando A, Song Q, De Santiago LM, Hulse-Kemp AM, Ding M,
638     Ye W, Kirkbride RC, Jenkins J, et al. 2020. Genomic diversifications of five
639     Gossypium allopolyploid species and their impact on cotton improvement. *Nat.*
640     *Genet.* [Internet]. Available from: http://dx.doi.org/10.1038/s41588-020-0614-5

641 Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden
642     DM. 2012. A program for annotating and predicting the effects of single nucleotide
643     polymorphisms, SnpEff: SNPs in the genome of Drosophila melanogaster strain
644     w1118; iso-2; iso-3. *Fly*  6:80–92.

645 Conant GC, Birchler JA, Pires JC. 2014. Dosage, duplication, and diploidization:
646     clarifying the interplay of multiple models for duplicate gene evolution over time.
647     *Curr. Opin. Plant Biol.* 19:91–98.

648 Conover JL, Sharbrough J, Wendel JF. 2021. pSONIC: Ploidy-aware Syntenic
649     Orthologous Networks Identified via Collinearity. *G3*  [Internet]. Available from:

650 http://dx.doi.org/10.1093/g3journal/jkab170

651 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,
652 Lunter G, Marth GT, Sherry ST, et al. 2011. The variant call format and VCFtools.
653 *Bioinformatics* 27:2156–2158.

654 Davydov EV, Goode DL, Sirota M, Cooper GM, Sidow A, Batzoglou S. 2010. Identifying
655 a high fraction of the human genome to be under selective constraint using
656 GERP++. *PLoS Comput. Biol.* 6:e1001025.

657 Douglas GM, Gos G, Steige KA, Salcedo A, Holm K, Josephs EB, Arunkumar R, Ågren
658 JA, Hazzouri KM, Wang W, et al. 2015. Hybrid origins and the earliest stages of
659 diploidization in the highly successful recent polyploid Capsella bursa-pastoris.
660 *Proc. Natl. Acad. Sci. U. S. A.* 112:2806–2811.

661 Doyle JJ, Coate JE. 2019. Polyploidy, the nucleotype, and novelty: The impact of
662 genome doubling on the biology of the cell. *Int. J. Plant Sci.* 180:1–52.

663 Eyre-Walker A, Keightley PD. 2007. The distribution of fitness effects of new mutations.
664 *Nat. Rev. Genet.* 8:610–618.

665 Fang L, Guan X, Zhang T. 2017. Asymmetric evolution and domestication in
666 allotetraploid cotton (Gossypium hirsutum L.). *The Crop Journal* 5:159–165.

667 Fang L, Wang Q, Hu Y, Jia Y, Chen J, Liu B, Zhang Z, Guan X, Chen S, Zhou B, et al.
668 2017. Genomic analyses in cotton identify signatures of selection and loci
669 associated with fiber quality and yield traits. *Nat. Genet.* 49:1089–1098.

670 Fernandes Gyorfy M, Miller ER, Conover JL, Grover CE, Wendel JF, Sloan DB,
671 Sharbrough J. 2021. Nuclear-cytoplasmic balance: whole genome duplications
672 induce elevated organellar genome copy number. *Plant J.* [Internet]. Available from:
673 http://dx.doi.org/10.1111/tpj.15436

674 Flagel LE, Wendel JF, Udall JA. 2012. Duplicate gene evolution, homoeologous
675 recombination, and transcriptome characterization in allopolyploid cotton. *BMC*
676 *Genomics* 13:302.

677 Gazave E, Chang D, Clark AG, Keinan A. 2013. Population growth inflates the per-
678 individual number of deleterious mutations and reduces their mean effect. *Genetics*
679 195:969–978.

680 Grant V. 1981. Plant Speciation. In: Plant Speciation. Columbia University Press.

681 Grover CE, Arick MA 2nd, Conover JL, Thrash A, Hu G, Sanders WS, Hsu C-Y, Naqvi
682 RZ, Farooq M, Li X, et al. 2017. Comparative Genomics of an Unusual
683 Biogeographic Disjunction in the Cotton Tribe (Gossypieae) Yields Insights into
684 Genome Downsizing. *Genome Biol. Evol.* 9:3328–3344.

685  Grover CE, Grupp KK, Wanzek RJ, Wendel JF. 2012. Assessing the monophyly of
686       polyploid Gossypium species. *Plant Syst. Evol.* 298:1177–1183.

687  Haldane JBS. 1932. The Causes of Evolution. 55 Fifth Avenue, New York : Longmans,
688       Green and Co.

689  Henn BM, Botigué LR, Peischl S, Dupanloup I, Lipatov M, Maples BK, Martin AR,
690       Musharoff S, Cann H, Snyder MP, et al. 2016. Distance from sub-Saharan Africa
691       predicts mutational load in diverse human genomes. *Proc. Natl. Acad. Sci. U. S. A.*
692       113:E440–E449.

693  Hill RR Jr. 1970. Selection in autotetraploids. *Theor. Appl. Genet.* 41:181–186.

694  Huber CD, Durvasula A, Hancock AM, Lohmueller KE. 2018. Gene expression drives
695       the evolution of dominance. *Nat. Commun.* 9:2750.

696  Hufford MB, Seetharam AS, Woodhouse MR, Chougule KM, Ou S, Liu J, Ricci WA,
697       Guo T, Olson A, Qiu Y, et al. 2021. De novo assembly, annotation, and
698       comparative analysis of 26 diverse maize genomes. *Science* 373:655–662.

699  Hu G, Grover CE, Yuan D, Dong Y, Miller E, Conover JL, Wendel JF. 2021. Evolution
700       and Diversity of the Cotton Genome. In: Rahman M-U-, Zafar Y, Zhang T, editors.
701       Cotton Precision Breeding. Cham: Springer International Publishing. p. 25–78.

702  Hu G, Wendel JF. 2019. Cis-trans controls and regulatory novelty accompanying
703       allopolyploidization. *New Phytol.* 221:1691–1700.

704  Jiao W, Yuan J, Jiang S, Liu Y, Wang L, Liu M, Zheng D, Ye W, Wang X, Chen ZJ.
705       2018. Asymmetrical changes of gene expression, small RNAs and chromatin in two
706       resynthesized wheat allotetraploids. *Plant J.* 93:828–842.

707  Jiao Y, Wickett NJ, Ayyampalayam S, Chanderbali AS, Landherr L, Ralph PE, Tomsho
708       LP, Hu Y, Liang H, Soltis PS, et al. 2011. Ancestral polyploidy in seed plants and
709       angiosperms. *Nature* 473:97–100.

710  Kendig KI, Baheti S, Bockol MA, Drucker TM, Hart SN, Heldenbrand JR, Hernaez M,
711       Hudson ME, Kalmbach MT, Klee EW, et al. 2019. Sentieon DNASeq variant calling
712       workflow demonstrates strong computational performance and accuracy. *Front.*
713       *Genet.* 10:736.

714  Kono TJY, Fu F, Mohammadi M, Hoffman PJ, Liu C, Stupar RM, Smith KP, Tiffin P, Fay
715       JC, Morrell PL. 2016. The role of deleterious substitutions in crop genomes. *Mol.*
716       *Biol. Evol.* 33:2307–2317.

717  Kono TJY, Lei L, Shih C-H, Hoffman PJ, Morrell PL, Fay JC. 2018. Comparative
718       genomics approaches accurately predict deleterious variants in plants. *G3* 8:3321–
719       3329.

720 Kryvokhyzha D, Milesi P, Duan T, Orsucci M, Wright SI, Glémin S, Lascoux M. 2019.
721     Towards the new normal: Transcriptomic convergence and genomic legacy of the
722     two subgenomes of an allopolyploid weed (Capsella bursa-pastoris). *PLoS Genet.*
723     15:e1008131.

724 Kryvokhyzha D, Salcedo A, Eriksson MC, Duan T, Tawari N, Chen J, Guerrina M,
725     Kreiner JM, Kent TV, Lagercrantz U, et al. 2019. Parental legacy, demography, and
726     admixture influenced the evolution of the two subgenomes of the tetraploid
727     Capsella bursa-pastoris (Brassicaceae). *PLoS Genet.* 15:e1007949.

728 Kuzmin E, Taylor JS, Boone C. 2021. Retention of duplicated genes in evolution.
729     *Trends Genet.* [Internet]. Available from:
730     https://www.sciencedirect.com/science/article/pii/S0168952521001864

731 Kuzmin E, VanderSluis B, Nguyen Ba AN, Wang W, Koch EN, Usaj M, Khmelinskii A,
732     Usaj MM, van Leeuwen J, Kraus O, et al. 2020. Exploring whole-genome duplicate
733     gene retention with complex genetic interaction analysis. *Science* [Internet] 368.
734     Available from: http://dx.doi.org/10.1126/science.aaz5667

735 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows-Wheeler
736     transform. *Bioinformatics* 25:1754–1760.

737 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
738     Durbin R, 1000 Genome Project Data Processing Subgroup. 2009. The Sequence
739     Alignment/Map format and SAMtools. *Bioinformatics* 25:2078–2079.

740 Liu Q, Zhou Y, Morrell PL, Gaut BS. 2017. Deleterious Variants in Asian Rice and the
741     Potential Cost of Domestication. *Mol. Biol. Evol.* 34:908–924.

742 Lohmueller KE. 2014. The distribution of deleterious genetic variation in human
743     populations. *Curr. Opin. Genet. Dev.* 29:139–146.

744 Lohmueller KE, Indap AR, Schmidt S, Boyko AR, Hernandez RD, Hubisz MJ, Sninsky
745     JJ, White TJ, Sunyaev SR, Nielsen R, et al. 2008. Proportionally more deleterious
746     genetic variation in European than in African populations. *Nature* 451:994–997.

747 Lozano R, Gazave E, dos Santos JPR, Stetter MG, Valluru R, Bandillo N, Fernandes
748     SB, Brown PJ, Shakoor N, Mockler TC, et al. 2021. Comparative evolutionary
749     genetics of deleterious load in sorghum and maize. *Nature Plants* 7:17–24.

750 Mason AS, Wendel JF. 2020. Homoeologous Exchanges, Segmental Allopolyploidy,
751     and Polyploid Genome Evolution. *Front. Genet.* 11:1014.

752 Matsumoto T, Wakefield L, Peters A, Peto M, Spellman P, Grompe M. 2021.
753     Proliferative polyploid cells give rise to tumors via ploidy reduction. *Nat. Commun.*
754     12:646.

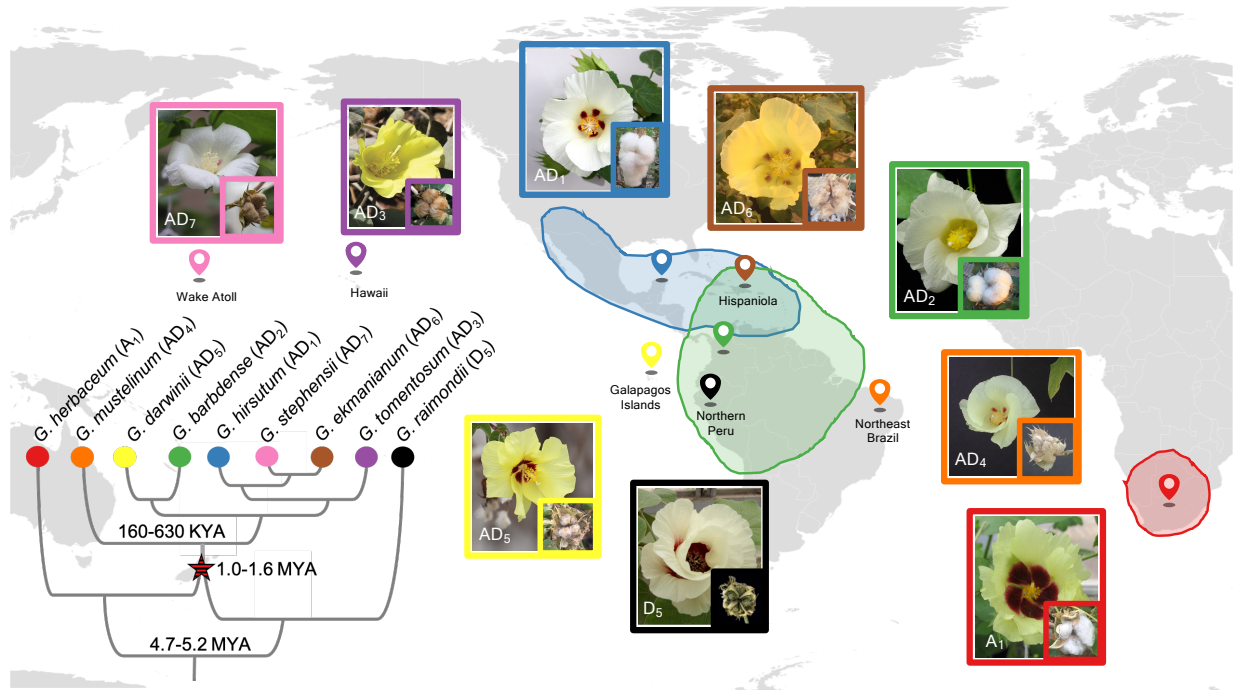755 Monnahan P, Kolář F, Baduel P, Sailer C, Koch J, Horvath R, Laenen B, Schmickl R,

Paajanen P, Šrámková G, et al. 2019. Pervasive population genomic consequences of genome duplication in Arabidopsis arenosa. *Nat Ecol Evol* 3:457–468.

Morton NE, Crow JF, Muller HJ. 1956. AN ESTIMATE OF THE MUTATIONAL DAMAGE IN MAN FROM DATA ON CONSANGUINEOUS MARRIAGES. *Proc. Natl. Acad. Sci. U. S. A.* 42:855–863.

Mukai T, Chigusa SI, Mettler LE, Crow JF. 1972. Mutation rate and dominance of genes affecting viability in Drosophila melanogaster. *Genetics* 72:335–355.

Ohta T. 1992. The Nearly Neutral Theory of Molecular Evolution. *Annu. Rev. Ecol. Syst.* 23:263–286.

One Thousand Plant Transcriptomes Initiative. 2019. One thousand plant transcriptomes and the phylogenomics of green plants. *Nature* 574:679–685.

Paape T, Briskine RV, Halstead-Nussloch G, Lischer HEL, Shimizu-Inatsugi R, Hatakeyama M, Tanaka K, Nishiyama T, Sabirov R, Sese J, et al. 2018. Patterns of polymorphism and selection in the subgenomes of the allopolyploid Arabidopsis kamchatica. *Nat. Commun.* 9:1–13.

Pont C, Salse J. 2017. Wheat paleohistory created asymmetrical genomic evolution. *Curr. Opin. Plant Biol.* 36:29–37.

Qi X, An H, Hall TE, Di C, Blischak PD, McKibben MTW, Hao Y, Conant GC, Pires JC, Barker MS. 2021. Genes derived from ancient polyploidy have higher genetic diversity and are associated with domestication in Brassica rapa. *New Phytol.* 230:372–386.

Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* 26:841–842.

Ramu P, Esuma W, Kawuki R, Rabbi IY, Egesi C, Bredeson JV, Bart RS, Verma J, Buckler ES, Lu F. 2017. Cassava haplotype map highlights fixation of deleterious mutations during clonal propagation. *Nat. Genet.* 49:959–963.

Ranwez V, Douzery EJP, Cambon C, Chantret N, Delsuc F. 2018. MACSE v2: Toolkit for the Alignment of Coding Sequences Accounting for Frameshifts and Stop Codons. *Mol. Biol. Evol.* 35:2582–2584.

Ranwez V, Harispe S, Delsuc F, Douzery EJP. 2011. MACSE: Multiple Alignment of Coding SEquences accounting for frameshifts and stop codons. *PLoS One* 6:e22594.

Renny-Byfield S, Wendel JF. 2014. Doubling down on genomes: Polyploidy and crop plants. *Am. J. Bot.* 101:1711–1725.

791 Rodgers-Melnick E, Bradbury PJ, Elshire RJ, Glaubitz JC, Acharya CB, Mitchell SE, Li
792      C, Li Y, Buckler ES. 2015. Recombination in diverse maize is stable, predictable,
793      and associated with genetic load. *Proc. Natl. Acad. Sci. U. S. A.* 112:3823–3828.

794 Roessler K, Muyle A, Diez CM, Gaut GRJ, Bousios A, Stitzer MC, Seymour DK,
795      Doebley JF, Liu Q, Gaut BS. 2019. The genome-wide dynamics of purging during
796      selfing in maize. *Nat Plants* 5:980–990.

797 Salmon A, Flagel L, Ying B, Udall JA, Wendel JF. 2010. Homoeologous nonreciprocal
798      recombination in polyploid cotton. *New Phytol.* 186:123–134.

799 Saski CA, Scheffler BE, Hulse-Kemp AM, Liu B, Song Q, Ando A, Stelly DM, Scheffler
800      JA, Grimwood J, Jones DC, et al. 2017. Sub genome anchored physical
801      frameworks of the allotetraploid Upland cotton (Gossypium hirsutum L.) genome,
802      and an approach toward reference-grade assemblies of polyploids. *Sci. Rep.*
803      7:15274.

804 Shi X, Chen C, Yang H, Hou J, Ji T, Cheng J, Veitia RA, Birchler JA. 2020. The Gene
805      Balance Hypothesis: Epigenetics and Dosage Effects in Plants. In: Spillane C,
806      McKeown P, editors. Plant Epigenetics and Epigenomics : Methods and Protocols.
807      New York, NY: Springer US. p. 161–171.

808 Simons YB, Sella G. 2016. The impact of recent population history on the deleterious
809      mutation load in humans and close evolutionary relatives. *Curr. Opin. Genet. Dev.*
810      41:150–158.

811 Simons YB, Turchin MC, Pritchard JK, Sella G. 2014. The deleterious mutation load is
812      insensitive to recent population history. *Nat. Genet.* 46:220–224.

813 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-
814      analysis of large phylogenies. *Bioinformatics* 30:1312–1313.

815 Tong C, Kole C, Liu L, Cheng X, Huang J, Liu S. 2020. The Asymmetrical Evolution of
816      the Mesopolyploid Brassica oleracea Genome. *The Brassica Oleracea Genome*:67.

817 Veitia RA, Birchler JA. 2021. Gene-dosage issues: a recurrent theme in whole genome
818      duplication events. *Trends Genet.* [Internet]. Available from:
819      http://dx.doi.org/10.1016/j.tig.2021.06.006

820 Wang L, Beissinger TM, Lorant A, Ross-Ibarra C, Ross-Ibarra J, Hufford MB. 2017. The
821      interplay of demography and selection during maize domestication and expansion.
822      *Genome Biol.* 18:215.

823 Wendel JF. 1989. New World tetraploid cottons contain Old World cytoplasm. *Proc.*
824      *Natl. Acad. Sci. U. S. A.* 86:4132–4136.

825 Wendel JF. 2015. The wondrous cycles of polyploidy in plants. *Am. J. Bot.* 102:1753–
826      1756.

827   Wendel JF, Grover CE. 2015. Taxonomy and Evolution of the Cotton Genus,
828       Gossypium. In: Cotton. Agronomy Monograph. Madison, WI: American Society of
829       Agronomy, Inc., Crop Science Society of America, Inc., and Soil Science Society of
830       America, Inc. p. 25–44.

831   Yuan D, Grover CE, Hu G, Pan M, Miller ER, Conover JL, Hunt SP, Udall JA, Wendel
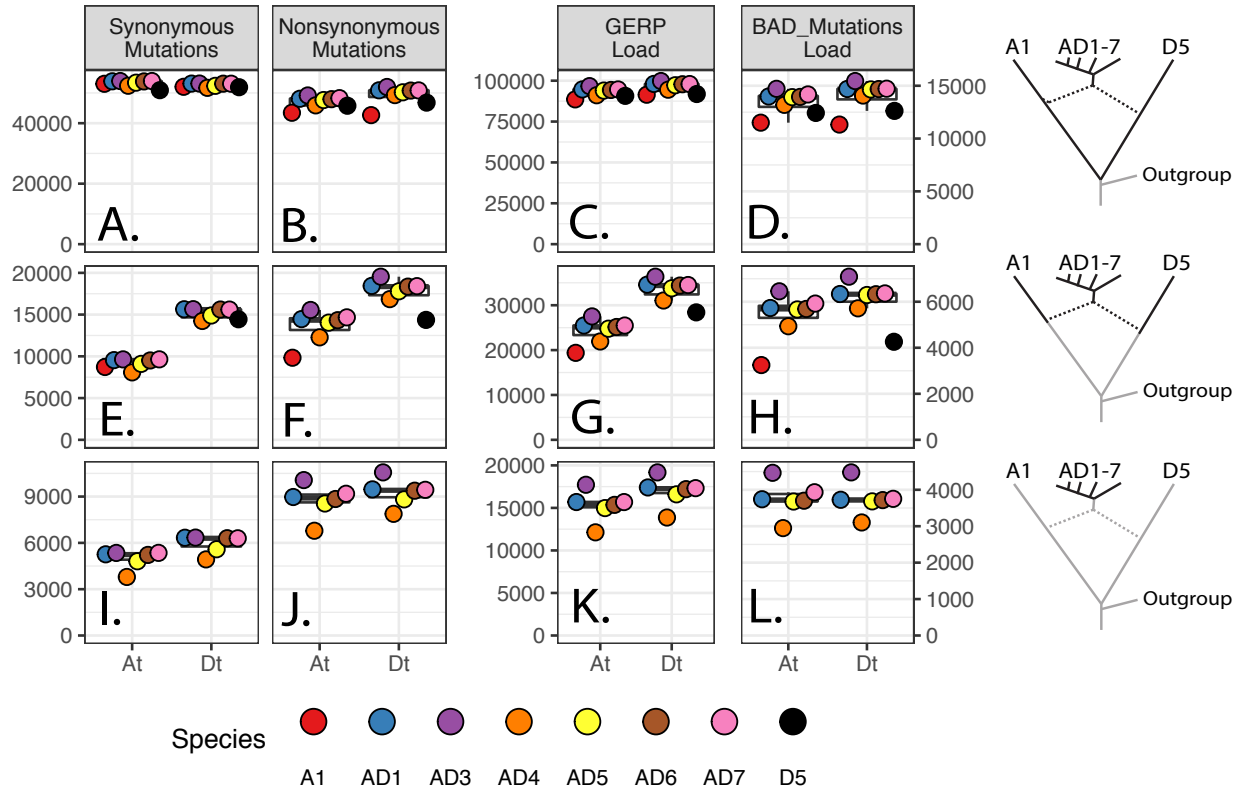832       JF. 2021. Parallel and intertwining threads of domestication in allopolyploid cotton.
833       *Adv. Sci.*:2003634.

834

**Figure and Table Legends:**



836

837 **Figure 1: Phylogeny and Biogeography of *Gossypium* Allopolyploids and**

838 **Progenitor Diploids**

839 Diploid *Gossypium* species are classified into eight diploid genome groups. The A

840 (represented by *G. herbaceum*) and D (represented by *G. raimondii*) genome groups

841 diverged approximately 5 million years ago (MYA), with ranges in different hemispheres.

842 Allopolyploids formed *circa* 1-1.6 MYA following transoceanic dispersal of an A genome

843 ancestor (modeled by *G. herbaceum* ($A_1$)) to the Americas and hybridization with a

844 native D genome species (modeled by *G. raimondii* ($D_5$)). Subsequent diversification of

845 the new allopolyploid (AD genome) lineage led to the evolution of seven currently

846 recognized species with a broad geographic range in the Americas and the Pacific

847 islands. Flower and fruit morphology for each species is shown, and the island location

848 and geographic range is indicated. Branch lengths on the phylogeny are not to scale but

849 notable divergence times are labeled.

**Figure 2: Derived Mutations and Deleterious Loads at Three Phylogenetic Depths**

Number of derived synonymous, nonsynonymous, and deleterious mutations in the

CDS regions of 8,884 pairs of homoeologs (17,768 genes in total) in eight cotton

species at three phylogenetic depths (indicated by bold branches in phylogeny at right).

For all panels, the ancestral state of each SNP was determined using three Australian

cotton species as an outgroup (see Methods). The deepest phylogenetic depth **(ABCD)**

includes all derived mutations that originated since the divergence of the A and D

diploid progenitors; the middle row **(EFGH)** shows SNPs that are variable within each

subgenome and its associated progenitor diploid species; and the bottom row **(IJKL)**

shows SNPs that originated post-polyploidy and are variable within the polyploids. **(AEI)**

Synonymous mutations. **(BFJ)** Nonsynonymous mutations.  The y-axis for both

synonymous and nonsynonymous is shown at left, and represents the sum of the

863    derived allele frequencies, interpreted as the average number of derived SNPs in that

864    category in each species. **(CGK)** GERP Load of each species, calculated as the sum of

865    (derived allele frequency * GERP Score) for all SNP positions with GERP > 0. **(DHL)**

866    Number of deleterious mutations in each species, calculated by BAD_Mutations with

867    Bonferroni-corrected significance (see Methods). Y-axis represents the sum of the

868    derived allele frequencies, and indicates the average number of deleterious mutations in

869    each species at a given phylogenetic depth. **Note:** for **(EFGH)**, comparisons between

870    subgenomes cannot be made because the D5 diploid is more distantly related to the D

871    subgenome than the A1 diploid is related to the A subgenome. Therefore, we would

872    expect a larger number of derived mutations in D than A simply due to evolutionary

873    history rather than to polyploidization *per se*.

874

875



876

**Figure 3: Proportions of All Nonsynonymous Mutations That Are Deleterious**

Rows **A**, **B**, and **C** summarize SNPs segregating within the entire clade, within each

subgenome and its respective progenitor diploid, and within each subgenome, as

indicated by the bolded branches along the phylogeny at left. Values indicate the

proportion of nonsynonymous SNPs that are deleterious within 8,884 homoeologous

pairs (17,768 total genes) that are syntenically conserved between the two subgenomes

of *G. hirsutum* (see Methods for filtering criteria). For example, the values in row A are

calculated by dividing the values in Figure 2D by the values in Figure 2B for each

species. **Note:** Similar to Figure 2, comparisons between subgenomes in row **B** reflect

differing phylogenetic distances, not asymmetries between the subgenomes and/or their

diploid progenitors.

888

39

**Figure 4: Relative Increase Of Deleterious Mutations Among GERP Categories in Polyploids Compared to Diploids**

For SNPs that originated since the divergence of each subgenome from its diploid progenitor, we plotted the relative increase in deleterious alleles across three GERP load categories: mildly deleterious (0<GERP≤2; light gray), moderately deleterious (2<GERP≤4; gray), and strongly deleterious (4<GERP≤6; black). We used the diploid as the reference population, meaning that the relative increase of GERP load in the diploid is always equal to one for all categories. In both subgenomes of all polyploids, strongly deleterious mutations had the greatest relative increase compared to the diploids, followed by the moderately deleterious mutations, and finally, mildly deleterious mutations. This pattern does not fit the expected patterns under demographic models alone, where most of the changes between two populations should be seen in mildly or moderately deleterious mutations. However, under a model where recessive deleterious mutations are masked by their homoeologs, we would expect that strongly deleterious mutations would accumulate faster than moderately or mildly mutations (i.e the pattern we see here) due to the correlation between the recessivity of a mutation (h) and its selection coefficient (s).
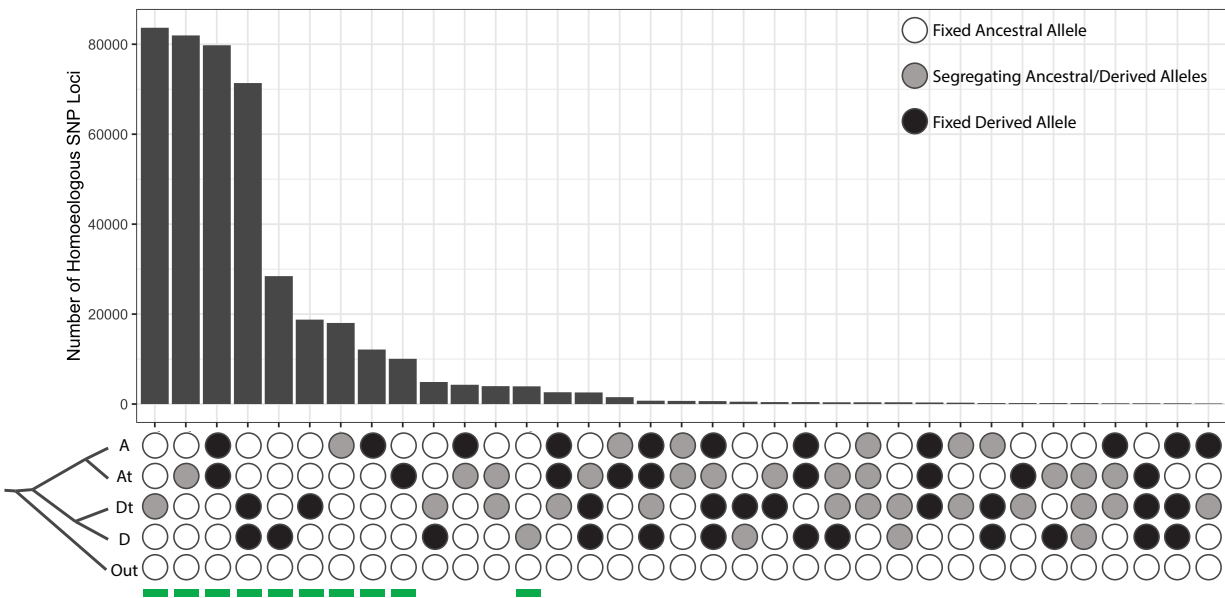
908 **Table 1: Nucleotide Diversity (π) in 8,884 Homoeologs in Eight *Gossypium***
909 **Species, By Subgenome**

| Species | Species Code | At Subgenome | Dt Subgenome |
|---|---|---|---|
| *G. herbaceum* | A1 | 7.41E-04 | |
| *G. raimondii* | D5 | | 2.36E-04 |
| *G. hirsutum* | AD1 | 6.69E-04 | 7.06E-04 |
| *G. tomentosum* | AD3 | 1.75E-04 | 1.67E-04 |
| *G. mustelinum* | AD4 | 2.64E-04 | 3.15E-04 |
| *G. darwinii* | AD5 | 1.71E-04 | 1.60E-04 |
| *G. ekmanianum* | AD6 | 7.75E-04 | 7.67E-04 |
| *G. stephensii* | AD7 | 4.94E-05 | 5.59E-05 |

910

911

912



913

**Supplementary Figure 1: UpSet Plot of Derived Homoeologous SNPs Among**

**8,884 Syntenic Homoeologous Gene Pairs**

To identify SNPs that may have potentially arisen from causes other than simple

nucleotide substitutions (e.g., sequencing error, gene conversion), we plotted the

frequency of polarized (ancestral vs derived) SNPs across the four major clades of

*Gossypium* allopolyploid genomes (A diploid, At subgenome, Dt subgenome, D diploid).

Bottom of the UpSet plot shows the phylogenetic positions of these 4 groups, as well as

the ancestral state used for polarization. For simplicity, we collapsed all polyploids into a

single group, but split them by subgenome (e.g. the At row indicates the At subgenome

in all 6 allopolyploids in this analysis). White bubbles indicate that only ancestral alleles

were identified in that species or subgenome; black bubbles denoteSNP sites where

only derived alleles were identified; grey bubbles represent SNP sites where both

ancestral and derived alleles were identified. Only the top 35 SNP groups are shown.

Groups with a green line underneath indicate SNP patterns that can be explained by a

42

928    single mutational event with no homoplasy (e.g. from incomplete lineage sorting or

929    recurrent mutation), and were retained for subsequent analyses involving the 8,884
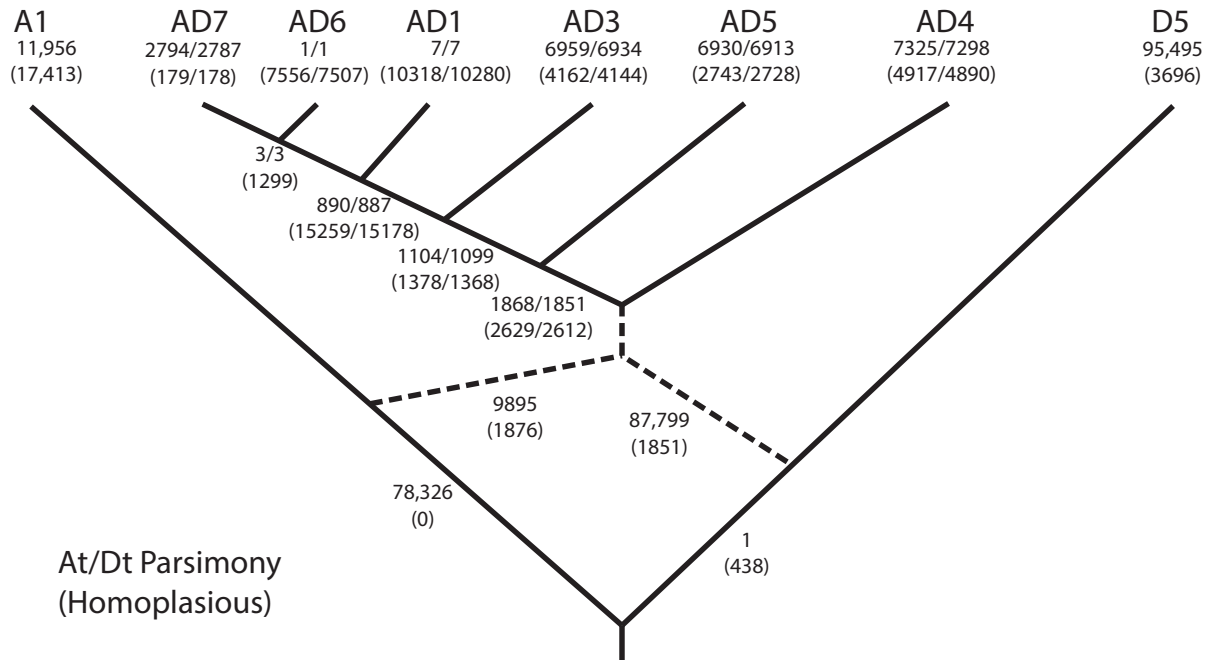
930    homoeologous gene pairs.

931

932

933

**Supplementary Figure 2: Genome-Wide Derived Mutations and Deleterious Loads**

**at Three Phylogenetic Depths**

Number of derived synonymous, nonsynonymous, and deleterious mutations in the

44

937　CDS regions of 8,884 pairs of homoeologs (17,768 genes in total) in eight cotton

938　species at three phylogenetic depths (indicated by bold branches of phylogeny at right).

939　For all panels, the ancestral state of each SNP was determined using three Australian

940　cottons as an outgroup (see Methods). The deepest phylogenetic depth **(ABCD)**

941　includes all derived mutations that originated since the divergence of the A and D

942　diploid progenitors; the middle row **(EFGH)** shows SNPs that are variable within each

943　subgenome and its associated progenitor diploid species; and the bottom row **(IJKL)**

944　shows SNPs that originated post-polyploidy and are variable within the polyploids. **(AEI)**

945　Synonymous mutations. **(BFJ)** Nonsynonymous mutations.  The y-axis for both

946　synonymous and nonsynonymous is shown at left, and represents the sum of the

947　derived allele frequencies, interpreted as the average number of derived SNPs in that

948　category in each species. **(CGK)** GERP Load of each species, calculated as the sum of

949　(derived allele frequency * GERP Score) for all SNP positions with GERP > 0. **(DHL)**

950　Number of deleterious mutations in each species, calculated by BAD_Mutations with

951　bonferroni corrected significance (see Methods). Y-axis represents the sum of the

952　derived allele frequencies, and indicates the average number of deleterious mutations in

953　each species at a given phylogenetic depth. **Note:** for **(EFGH)**, comparisons between

954　subgenomes cannot be made because the D5 diploid is more distantly related to the D

955　subgenome than the A1 diploid is related to the A subgenome. Therefore, we would

956　expect a larger number of derived mutations in D than A simply due to evolutionary

957　history rather than to polyploidization *per se*. The panels above the figure legend are

958　identical to those presented in Figure 2. The panels below the figure legend **(M-X)** follow

959　the same order as **(A-L)**, but represent the genome-wide totals without any filtering

45

960     based on homoeologs or potential sites that are due to gene less, mapping biases, or

961     homoeologous gene conversion and is provided to demonstrate that our filtering criteria

962     did not have a noticeable impact on the patterns of SNPs that we observed, and that

963     homoeologous interactions have a minimal effect on patterns of evolution following
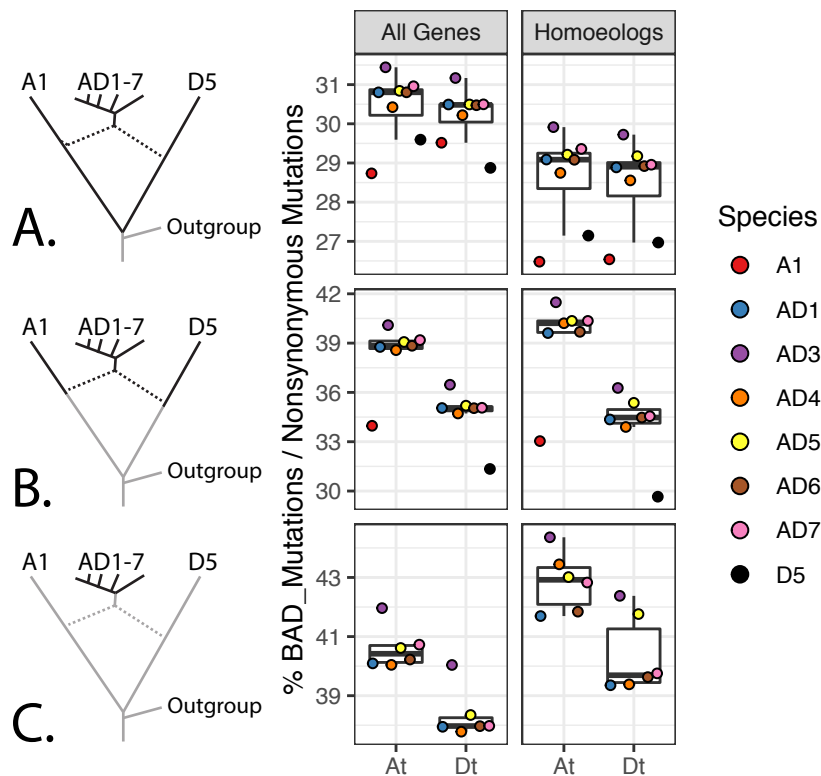
964     allopolyploidy in *Gossypium*.

965

966

**Supplementary Figure 3: Phylogenetic Positions of Derived Deleterious SNPs**

For SNPs that passed the filtering from Supplementary Figure 1, we placed the origin of the SNP on the phylogenetic tree using parsimony. Numbers in the format of "X/Y" indicate the number of SNPs found in the "At/Dt" subgenome. Numbers above the parentheses indicate SNPs that are unequivocally placed on the tree in either the At or Dt subgenome. Numbers in parentheses indicate SNPs that are homoplasious, and the position of the number represents the phylogenetic position of the most recent common ancestor of all species that contain at least one derived SNP. Numbers in the parentheses at the tips of the tree indicate SNPs that are segregating within that species but are not found in any other species. Note: the high amount of homoplasious SNPs at the base of the AD1, AD6, and AD7 clade is most likely caused by recent hybridization or introgression of AD1 into AD6, as also indicated in Supplementary Figure 5.
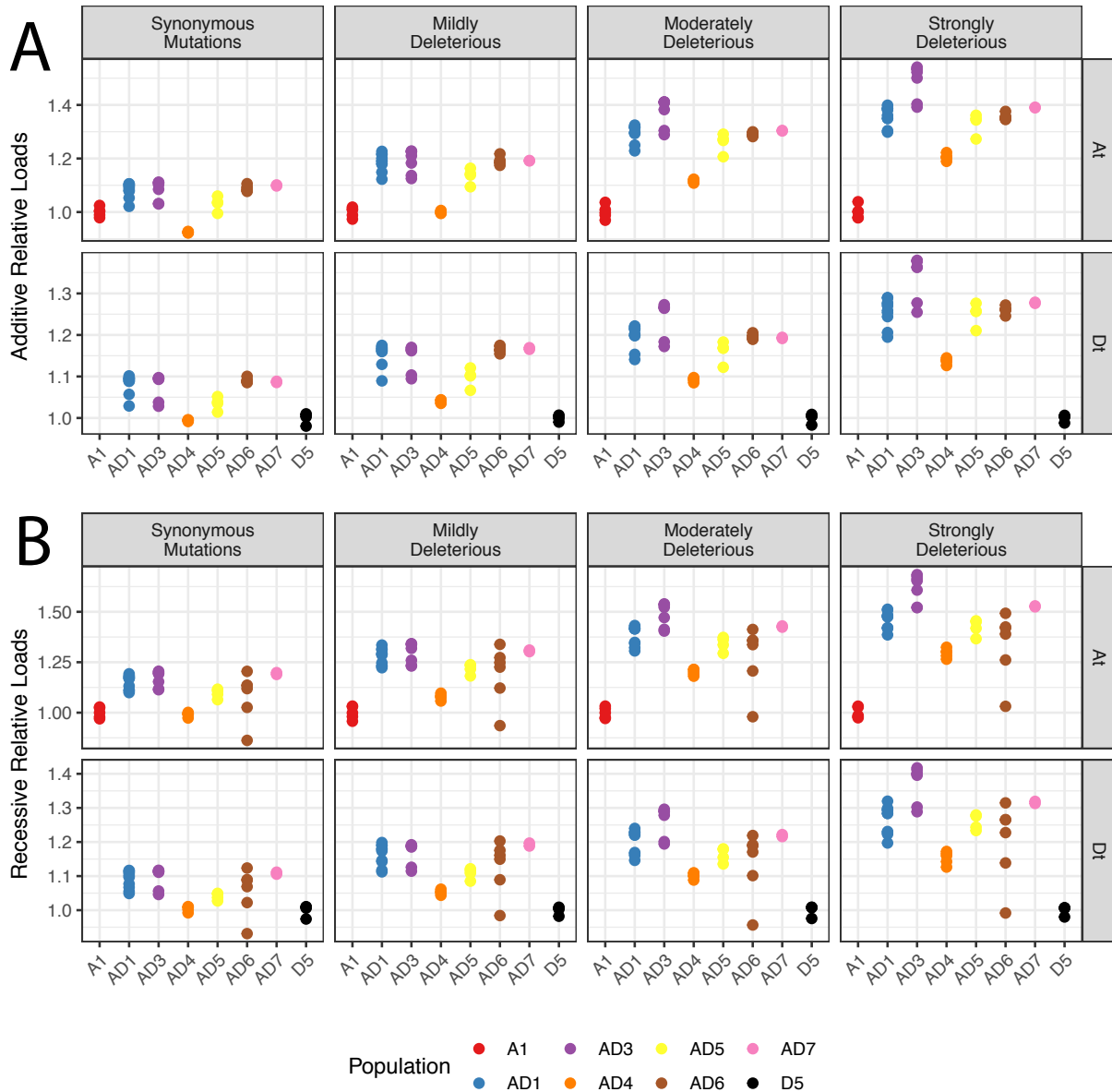
47

980



981

**Supplementary Figure 4: Genome-Wide Proportions of All Nonsynonymous**

**Mutations That Are Deleterious**

Rows **A**, **B**, and **C** summarize SNPs segregating within the entire clade, within each

subgenome and its respective progenitor diploid, and within each subgenome, as

indicated by the bolded branches along the phylogeny at left. **(A)** Proportion of all

nonsynonymous SNPs that are deleterious genome-wide within each subgenome. **(B)**

Proportion of nonsynonymous SNPs that are deleterious within 8,884 homoeologous

pairs (17,768 total genes) that are syntenically conserved between the two subgenomes

of *G. hirsutum* (see Methods for filtering criteria). Note: Similar to Figure 2, comparisons

between subgenomes in row **B** reflect differing phylogenetic distances, not asymmetries

between the subgenomes and/or their diploid progenitors.

993

**Supplementary Figure 5: Additive and Recessive Models of Deleterious Mutation**

**Accumulation**

Relative load of synonymous sites and varying GERP categories from an **(A)** additive

model (i.e. counting all SNPs) and **(B)** recessive model (i.e. counting all homozygous

SNPs in a homozygous state). Each point represents an individual, and the placement

of each point represents the relative increase or decrease in the number of SNPs

1000    relative to the average of the number of SNPs in the diploid (A1 for At, D5 for Dt). Note:

1001    The high variance in the recessive load for AD6 reflects a high number of sites that are

1002    heterozygous. This is mostly likely due to recent hybridization or introgression from

1003    AD1, which is also indicated by a high amount of incomplete lineage sorting between

1004    AD1, AD6, and AD7 in Supplementary Figure 3.