# Pervasive population genomic consequences of genome duplication
# in *Arabidopsis arenosa*

Patrick Monnahan[1], Filip Kolář[2,3,4], Pierre Baduel[1], Christian Sailer[1], Jordan Koch[1], Robert Horvath[5], Benjamin Laenen[5], Roswitha Schmickl[2,4], Pirita Paajanen[1], Gabriela Šrámková[2], Magdalena Bohutínská[2,4], Brian Arnold[6], Caroline M. Weisman[7], Karol Marhold[2,8], Tanja Slotte[5], Kirsten Bomblies[1], and Levi Yant[1, 9,] *

1. Department of Cell and Developmental Biology, John Innes Centre, Norwich Research Park, Norwich, NR4 7UH, UK

2. Department of Botany, Faculty of Science, Charles University, Benátská 2, 128 01 Prague, Czech Republic

3. Department of Botany, University of Innsbruck, Sternwartestraße 15, A-6020 Innsbruck, Austria

4. Institute of Botany, The Czech Academy of Sciences, Zámek 1, 252 43 Průhonice, Czech Republic

5. Department of Ecology, Environment and Plant Sciences, Science for Life Laboratory, Stockholm University, SE-106 91 Stockholm, Sweden

6. Center for Communicable Disease Dynamics, Department of Epidemiology, Harvard T. H. Chan School of Public Health, Boston, MA 02115 USA

7. Department of Organismic and Evolutionary Biology, Harvard University, 16 Divinity Avenue, Cambridge, MA, 02138 USA

8. Plant Science and Biodiversity Centre, Slovak Academy of Sciences, Dúbravská cesta 9, SK-845 23 Bratislava, Slovak Republic

9. School of Life Sciences and Future Food Beacon, University of Nottingham, Nottingham, UK

Patrick Monnahan, Filip Kolář, and Pierre Baduel contributed equally.

**\*Author for correspondence:** levi.yant@nottingham.ac.uk; Tel: 0749 025 3006

1

1    **Abstract**

2        Ploidy-variable species allow direct inference of the effects of chromosome copy number on

3    fundamental evolutionary processes. While an abundance of theoretical work suggests polyploidy

4    should leave distinct population genomic signatures, empirical data remains sparse. We sequenced

5    ~300 individuals from 39 populations of *Arabidopsis arenosa*, a naturally diploid-autotetraploid

6    species. We find the impacts of polyploidy on population genomic processes are subtle yet pervasive,

7    including reduced efficiency of purifying selection, differences in linked selection, and rampant gene

8    flow from diploids. Initial masking of deleterious mutations, faster rates of nucleotide substitution, and

9    interploidy introgression likely conspire to shape the evolutionary potential of polyploids.

10

11


12    **Introduction**

13        Whole-genome duplications (WGD) have occurred throughout the tree of life [1, 2] and are

14    associated with biological phenomena of great socio-economic importance such as crop domestication

15    [3] and carcinogenesis [4]. The direct effects of WGD or polyploidy can be far-reaching, ranging from

16    cellular [5] through organism-level phenotypes [6], up to population and ecosystem-level processes [7-

17    9].

18        Population genetic theory predicts substantive effects of ploidy on both neutral and selective

19    processes [10-16]. With higher ploidy, neutral diversity is expected to rise, while the rate of population

20    differentiation due to genetic drift will slow [17]. Given the additional chromosomal partners available

21    during recombination, linkage disequilibrium should decrease and haplotype diversity correspondingly

22    increase. Polyploidy also has unique effects on migration. Although polyploidization is traditionally

23    viewed as a means of instant speciation [18-20], the ploidy barrier may be permeable, particularly from

24    diploids to polyploids [21, 22].  Additionally, polyploids may lack reproductive incompatibilities found

2

25  in diploid progenitors [23]. Upon secondary contact, interploidy introgression could then further enrich

26  polyploid diversity.

27      The effect of ploidy on selective processes can be primarily attributed to differential manifestation

28  of allelic dominance. Added masking of deleterious alleles should elevate equilibrium frequencies at

29  mutation-selection balance, potentially increasing genetic load [24]. Similarly, beneficial alleles are not

30  observed as readily in polyploids, which slows the fixation of individual alleles [25]. These

31  disadvantages can be mitigated by polyploids' propensity to receive, maintain, and generate genetic

32  variation.  By itself, the increased rate at which beneficial mutations are introduced in polyploids can

33  be sufficient for a faster overall rate of adaptation [16, 26]. Additionally, introgression is increasingly

34  being recognized as an important source of adaptive variation [27, 28]. Though genomic evidence in a

35  ploidy-variable system is lacking, a greater tendency for polyploids to accept variation from locally

36  adapted populations of the same or different ploidy (or even species), may facilitate adaptation and

37  expansion of polyploid lineages.

38      Lack of population genomic data is particularly pronounced for autopolyploids, which arise from

39  within-species WGD [29]. In contrast to the better studied allopolyploids, where effects of polyploidy

40  are confounded with subgenome divergence, autopolyploids allow direct investigations of the role of

41  polyploidy *per se*.  Using a new model for autopolyploidy, *Arabidopsis arenosa* [30], we generated the

42  most comprehensive range-wide genomic dataset to date of a natural tetrasomic autotetraploid [31, 32]

43  (182 individuals / 24 populations) and its diploid sister lineages (105 / 15; Figure 1a). The tetraploids,

44  whose ecological niche largely overlaps with the genetically divergent diploids, trace to a single origin

45  (~30 kya) [31] and subsequently spread across much of Europe, occasionally coming into secondary

46  contact with diploids [31, 33].

47      We focus on three main questions concerning the genomic impact of selection and migration in this

48  system: First, we investigate if purifying selection is relaxed in autotetraploids as predicted from the

3

49    increased masking of deleterious alleles. Second, given the inherent effects of polyploidy on processes

50    governing diversity and recombination, we ask whether signals of linked selection markedly differ in

51    one ploidy versus the other. Lastly, we focus on two independent contact zones to assess the impact of

52    interploidy gene flow on polyploid evolution. Overall, our empirical analyses provide insights into the

53    complexity of autopolyploid evolution, supporting some but not all theoretical predictions. Altered

54    selective processes and introgressions shape the genomic landscape of tetraploids, and perhaps their

55    evolutionary potential as well.

56

57    **Results**

58    *High diversity and population differentiation in natural* A. arenosa

59        As previously reported [34], the diploid populations form five divergent, geographically-separated

60    groups (Fig. 1b, c): the *Baltic* lineage, the highly distinct *Pannonian* and *Dinaric* lineages ($\overline{F_{ST}}$ = 0.31

61    and 0.34, respectively), and the less differentiated Southern Carpathian (*S. Carp.*), and Western

62    Carpathian (*W. Carp.*) lineages ($\overline{F_{ST}}$ = 0.25, with evidence of past, Table S4, and recent hybridization,

63    e.g. HNI Fig 1b). The tetraploids comprise four lineages: *S. Carp., W. Carp, C. Europe* (Alps and

64    western Central Europe), and the *Ruderal* lineage. The latter group is the most widespread yet

65    ecologically distinct, occupying man-made sites (e.g. railways) from southern Germany to Sweden

66    (Fig. 1a). Ploidy is explicitly indicated as a suffix (*2x* or *4x*) hereafter. We find lower differentiation

67    among tetraploid populations than diploid ones (Table 1, S3, Fig. 1c, d), in line with the greater age of

68    diploids and the neutral expectation that, all else equal, the rate of drift is halved relative to diploids

69    [17]).

70        *Arabidopsis arenosa* is an obligate outcrosser, and all populations exhibit high genome-wide

71    diversity ($\bar{\theta}_\pi$ = 0.015, Table 1), an order of magnitude higher than the predominantly self-fertilizing

72    *A. thaliana* [35]. All else equal, polyploidy is expected to increase diversity ($8N_e\mu$ in tetraploids versus

73     $4Ne\mu$ in diploids). Although tetraploid populations exhibit slightly higher Watterson's $\theta_W$ at non-

74     synonymous sites (zero-fold degenerate sites, 0-dg), we observe no significant increase of $\theta_\pi$ or $\theta_W$ in

75     tetraploid populations at putatively neutral sites (four-fold degenerate sites, 4-dg). These results were

76     robust to exclusion of tetraploid populations with evidence of interploidy admixture (DRA, LAC, TZI,

77     KOW, STE, TBG; Table S2). Such an impact of genome duplication on $\theta_W$ (in contrast to $\theta_\pi$) is

78     consistent with a recent origin of tetraploids, as $\theta_W$ is more sensitive to accumulation of rare variants.

79

80     *Ploidy effects on purifying selection*

81        Though we find only mild differences between ploidies in $\theta_\pi$ or $\theta_W$ , we observe a highly

82     significant difference for the ratio of 0-dg $\theta_W$ to 4-dg $\theta_W$ (Wilcoxon rank-sum test, W=54, $p = 0.001$;

83     Table 1, S1), which is consistent with expectations of relaxed purifying selection in tetraploids [24]. To

84     further explore this hypothesis, we assessed how gene-level diversity varied with gene expression, a

85     proxy for selective constraint (e.g. [36-38]). We confirmed that highly expressed genes exhibit reduced

86     nonsynonymous diversity in both ploidies (multiple linear model, MLM; $p < 0.0001$, *F*-test of

87     expression effect on $\theta_W$ at 0-dg sites and the 0-dg/4-dg ratio of $\theta_W$; Fig. 2a, 2b and Table S5). However,

88     the 0-dg/4-dg $\theta_W$ ratio was generally higher in tetraploids ($p < 0.0001$) due to elevated nonsynonymous

89     diversity ($\alpha_p$ coefficient in MLM; $p < 0.0001$, Table 1, Fig. 2b, S5).  These results were robust across

90     data subsets and upon including various cofactors (e.g. population sizes; Tables S5 & S6). This

91     confirms that, beyond the increased mutational input resulting from doubled genome copies in

92     tetraploids, there is an additional increase of non-synonymous diversity (thus increasing the 0-dg/4-dg

93     diversity ratio), which likely reflects an overall relaxation of purifying selection.

94        Such relaxation could be due to either a reduction in the *strength* of selection or simply because

95     selection is less efficient in tetraploids. If mutations are purely recessive, the homozygotes bearing the

96     deleterious phenotype are much less frequent in tetraploids ($q^2$ versus $q^4$, assuming random mating

97    [32]), making purifying selection inefficient relative to diploids even if the fitness cost of the mutant

98    homozygote (i.e. selection *strength*) is equivalent across ploidies. To distinguish between these two

99    explanations, we evaluated the distribution of fitness effects (DFE) across both ploidies, finding no

100   apparent differences in the strength of purifying selection in diploid vs. tetraploid populations (Fig. 2d).

101   From this analysis, it seems purifying selection is not weaker *per se*, but rather less efficient at reducing

102   allele frequencies because deleterious mutations are better masked in autotetraploids.

103        However, we note two important assumptions in DFE estimation methods [39] that complicate

104   interploidy comparisons. First, a diploid model of allele frequencies at mutation-selection-drift balance

105   is assumed. Since frequencies are expected to be higher in autotetraploids [15], this model would be

106   biased towards inferring weaker selection than necessary to explain polyploid data. Second, deleterious

107   mutations are assumed to be additive. If deleterious mutations are recessive, equilibrium frequencies

108   can be orders of magnitude greater in tetraploids, further amplifying the first bias. If purifying selection

109   were truly weaker in tetraploids, these biases would make this more apparent; instead, we find no

110   evidence for ploidy differences in the DFE (Fig. 2d, S7 and Table S7 and S14).

111        In the long run, these selective effects (along with increased mutational input) are expected to

112   result in higher genetic load for tetraploids (under partial recessivity; load should be equivalent at

113   equilibrium for complete recessivity) [24]. To obtain a crude estimate of genetic load in each

114   population, we counted homozygous genotypes per-individual for *derived*, nonsynonymous alleles.

115   Under complete recessivity, the estimated load is currently lower in tetraploids than in diploids

116   (Wilcoxon rank-sum test of population means, $W = 264$, $p < 0.0001$ and $W = 195$, $p < 0.0001$

117   with/without interploidy-admixed tetraploid populations, respectively; Fig. 2c).  However, the

118   relatively young tetraploid lineages may not have reached equilibrium [31], which could take hundreds

119   of thousands of generations [16]. Furthermore, the actual load may be substantially higher in tetraploids

120   if deleterious mutations are at least partially recessive [24]. Unfortunately, current methodologies do

121    not allow for relaxation of the assumption of complete recessivity for tetraploids.

122

123    *Ploidy effects on positive and linked selection*

124        If dominance when in single-copy is comparable across ploidy, such that *Aa* and *Aaaa* genotypes

125    are equivalent, the greater mutational opportunity in tetraploids should ultimately lead to higher rates of

126    adaptation [16]. Using DFE-alpha analysis [38], we estimated the proportion of nonsynonymous (0-dg)

127    sites fixed by positive selection in each population. Using either $\alpha$ or $\omega_\alpha$, this proportion was

128    significantly higher in tetraploid populations (W = 14 or 6, respectively; $p < 0.0001$ for both; Fig. 2e,

129    S8 and Table S7 and S14), possibly indicating increased rates of adaptive substitution. This does not

130    simply reflect admixture (below), as the difference remained significant when we removed the six

131    tetraploid populations admixed with diploids (Table S7). However, similar to the preceding section,

132    multiple non-selective processes can lead to mis-estimation of parameters within DFE-alpha and

133    similar methods [40-42] (discussed in Supplementary Text 3).

134        Although DFE-alpha suggests a higher proportion of adaptive substitutions in tetraploids, the

135    fixation of particular mutations is generally expected to take longer [25], with implications for the

136    degree that linked selection reduces diversity during selective sweeps. Using the average squared

137    genotypic correlation between SNPs, we approximated linkage disequilibrium (LD) (Fig. 3a), finding

138    an overall reduction in tetraploids (50% lower mean correlations at 1kb distance in tetraploids). We

139    then assessed the impact on linked selection by comparing across genomic windows excess

140    nonsynonymous divergence ($E_{NS} = d_N - d_S$) and 4-dg site diversity (Fig. 3b). Regardless of ploidy, $E_{NS}$

141    and 4-dg $\theta_\pi$ were consistently negatively correlated (Table S8), suggesting that divergent selection had

142    reduced diversity at linked, neutral sites. The parabola shape ($p < 0.001$ for quadratic term; Table S8)

143    indicates that diversity is also reduced for $E_{NS} \ll 0$ regions (i.e. those under background selection). The

144    reductive effect of $E_{NS}$ on neutral diversity was significantly stronger in gene-dense regions (upper

7

145    20%, Fig. 3d; interaction of $E_{NS}$ and gene-density: $p < 0.001$ in Table S8) than in gene poor regions

146    (lower 20%, Fig. 3c). Within gene-dense regions, we observed higher neutral diversity in tetraploids in

147    particular for negative $E_{NS}$ values (Fig. 3d, p=0.002 for 3-way interaction between $E_{NS}$, gene-density,

148    and ploidy, Table S8). This difference in slope suggests background selection is less effective at

149    reducing diversity in tetraploids, while selective sweeps reduce diversity similarly across ploidies. No

150    such differences were observed in low gene-density regions (Fig. 3c), where linked selection will be

151    less pronounced.

152         While slower fixation times in tetraploids would dampen a signature of linked selection, two

153    factors could effectively counter this effect: 1.) the evolution of reduced per-base recombination in

154    tetraploids (to avoid deleterious multivalents forming during meiosis [58]) and 2.) systematic

155    differences across ploidies in the age of selective sweeps (due to the comparatively recent tetraploid

156    formation). Such reduced recombination is not evident, genome-wide, in tetraploids. In fact, our LD

157    approximation is generally lower in tetraploids, reflecting a higher population recombination rate

158    ($\rho = 8N_e r$ in tetraploids and $\rho = 4N_e r$ in diploids; Fig. S22) and/or the more recent population expansion

159    [39]. Unfortunately, the lack of genetic maps and of a workable phasing algorithm prevents inclusion of

160    the recombination landscape in our regression modelling approach. Furthermore, estimation of the age

161    and strength of selection is not currently possible on a genomic scale. Understanding the interplay

162    between fixation times, recombination landscapes, and natural history will be the focus of future

163    investigations.

164

165    *Single origin of tetraploids and interploidy introgression*

166         Although previous work supported a single tetraploid origin in the W. Carpathians [31], local

167    tetraploids clustered genetically with locally co-occurring diploids in two parallel cases (Southern

168    Carpathians and Baltic coast; Fig. 1a,c, S3, S4B). This might suggest multiple tetraploid origins

169    followed by widespread gene flow among tetraploids, as these two tetraploid lineages still share a

170    sizeable portion of polymorphisms with the widespread tetraploid lineages (*W. Carp.-4x*, *C. Europe-4x*,

171    Fig. 1b). However, we find multiple lines of evidence supporting a single tetraploid origin followed

172    instead by interploidy gene flow from locally co-occurring diploids (see Supplementary Text 1 for

173    further discussion). First, coalescent simulations (*fastsimcoal2*) consistently favour scenarios with a

174    single tetraploid origin (~20k – 31k generations ago) followed by admixture (Fig. 4a,b, S9, S10; Table

175    S9). Second, frequencies of alleles diagnostic of the putative diploid ancestor of all tetraploids

176    (*W. Carp.-2x* lineage) are elevated and positively correlated across all tetraploid populations (Fig. 4c,

177    S11). Finally, alleles of several key meiosis genes are shared among all tetraploids, yet consistently

178    divergent from diploids (Fig. 4d, S12).

179        Interploidy gene flow could be mediated either by viable triploids (virtually absent in natural *A.*

180    *arenosa* [32]) or by a one-step production of tetraploid hybrids via merger of unreduced gametes of a

181    diploid (2n) with a normal (reduced) gamete of a tetraploid (also 2n) [52]. In the Southern Carpathian

182    contact zone 37% of the *S. Carp.-4x* possessed the regionally-specific plastid haplotypes typical for the

183    *S. Carp.-2x*, suggesting gene flow between ploidies sometimes involves diploid mothers. *Ruderal-4x*

184    populations, on the other hand, only shared plastid haplotypes with other tetraploid groups (*W. Carp.-*

185    *4x* and *C. Europe-4x*); in this contact zone, either selection favours tetraploid plastids, or gene flow

186    primarily involves male gametes from diploids. In addition, multiple tetraploid (but no diploid)

187    populations showed elevated frequencies of nuclear (Fig. S13) and occasionally plastid (Fig. S14)

188    markers otherwise private to *Arabidopsis lyrata* – a partially sympatric species known to hybridize with

189    *A. arenosa* at the tetraploid but not diploid level [23, 43].

190        The maintenance of tetraploid alleles at key meiosis genes in the face of introgression from

191    diploids implies that some genomic regions are more or less resistant to interploidy admixture. To

192    identify such regions in each contact zone, we first evaluated the weights of topologies supporting

9

193   tetraploid monophyly (TM) vs. local admixture (LA) in windows across the genome (Fig. 5, S15).

194   Generally, the tendency was for no single topology to dominate a particular window, yet occasionally,

195   we observed windows where the vast majority of weight was given to either the TM or LA topology.

196   Within the latter regions, we then looked for the specific pattern of: 1) reduced genetic divergence of

197   the focal tetraploid lineage to sympatric diploids versus non-sympatric diploids (as expected with local

198   admixture) and 2) elevated genetic divergence of the focal tetraploid to all other tetraploids. Lastly, we

199   looked for signatures of positive selection using Fay and Wu's H. In each contact zone, we identified

200   multiple regions with such three-fold evidence (Fig. 5; Fig. S16). Within the 1% outliers for both LA

201   topology weight and H, we found a number of gene coding loci (Table S10) with some indication of

202   functional enrichment (see Supplementary Text 2). Conversely, windows with high weight given to the

203   TM topology (Topology 1) often exhibited elevated divergence to *all* diploids and non-elevated

204   divergence to tetraploids. Additionally, these windows often included meiotic genes previously

205   identified as exhibiting the strongest signatures of tetraploid-specific selection in a subset of *A. arenosa*

206   populations [58]. Together, this is consistent with a strong tetraploid resistance to diploid introgression

207   in these regions, suggestive of their *ongoing* role in the maintenance of stable autopolyploid

208   chromosome segregation.

209

210   **Discussion**

211   Using the largest population resequencing dataset to date of a ploidy-variable plant species, we

212   observe pervasive differences in how forces governing genome evolution shape genetic diversity and

213   divergence in nature. In diploid and autotetraploid *A. arenosa*, we find subtly distinct signatures of

214   linked and purifying selection. Additionally, multiple sources of evidence indicate substantial

215   introgression from diploids to tetraploids. We discuss these results in terms of the inherent effects of

216   genome doubling and the implications for the evolutionary potential of polyploid lineages.

217    The effects of genome doubling on selective processes, and consequently the patterns of

218    genomic variation they leave behind, are multifarious and sometimes counter-acting, making it difficult

219    to observe and distinguish individual causes. The challenge is heightened by the lack of methodologies

220    generalized for higher ploidy and the potential for demographic events associated with the creation,

221    establishment, and expansion of nascent tetraploids to obfuscate the genomic signals of selection. Yet,

222    the fundamental impact of genome doubling on dominance relationships and mutational and

223    recombinatorial opportunity are clearly reflected in our analyses of linked and purifying selection.

224    For positive selection, the increased masking of a beneficial mutation's effect in tetraploid

225    populations is likely not sufficient to slow adaptation. The higher estimated proportion of

226    nonsynonymous polymorphisms fixed by positive selection in tetraploids (Fig. 2e) supports the notion

227    that increased mutational input is sufficient to overcome any hindrance to adaptation posed by the

228    reduced efficiency of selection [16, 26, 44]. Furthermore, increased fixation times (via increased

229    masking) and mutational and recombinatorial opportunity in tetraploids promote retention of haplotype

230    (Fig 3a) and nucleotide diversity (Fig 3d) following selection. *A. arenosa* tetraploids expanded well

231    beyond their ancestor diploid's range, including postglacial and man-made habitats [34]. Increased

232    mutational input and retention of diversity may aid polyploids in adapting to the fluctuating or

233    otherwise challenging environments that are often associated with polyploids [7, 45, 46].

234    With purifying selection, nonsynonymous polymorphism is governed simultaneously by selection

235    against deleterious alleles and their recurrent introduction via mutation; genome doubling favouring

236    increased polymorphism in both cases. Importantly, the former may result solely from the added

237    masking of recessive deleterious mutations in heterozygous genotypes; the *strength* of selection need

238    not differ.  In this context, increased diversity in tetraploids is detrimental, leading to higher genetic

239    load at equilibrium [24].  Our estimate of genetic load (assuming recessivity) is currently lower for

240    tetraploids, even though nonsynonymous diversity is higher for genes under purifying selection. In

241  addition to reasons discussed above (see *Results*), double reduction (a unique phenomenon in

242  autopolyploids where the resolution of multivalents occasionally causes sister chromatids to segregate

243  into the same gamete) may also play a role, by increasing homozygosity and allowing more efficient

244  purging of deleterious alleles, although this would only affect more distal chromosome regions [47].

245  Furthermore, the actual load in tetraploids could be much higher if deleterious alleles are at least

246  partially recessive [24], as has been demonstrated previously in a natural plant system [48]. Currently,

247  no comparable demonstration exists for autotetraploids, although ploidy-variable species, such as *A.*

248  *arenosa*, provide a compelling system in which this could be addressed.

249      Despite an increased recognition of adaptive introgression [49], introgression between ploidy

250  cytotypes could be maladaptive [50, 51]. Here, the most salient example lies in meiotic genes, which

251  have been shown to exhibit the strongest signatures of selection in tetraploids (presumably to promote

252  proper segregation of additional chromosomes to gametes [58]). The introduction of diploid-like

253  meiotic alleles into a tetraploid population would increase the frequency of multivalent formation, thus

254  decreasing fitness. In line with this, meiotic genes consistently show the strongest signatures of

255  introgression resistance in tetraploids: elevated divergence between ploidies, reduced diversity within

256  tetraploids, and tetraploid monophyly in both diploid-tetraploid contact zones (Figs. 4d, 5). On the

257  other hand, we found coding regions with diploid-like derived alleles that have swept to higher

258  frequencies in co-occurring tetraploids (Fig. 5), implying that interploidy introgression can be adaptive

259  in tetraploids. In fact, the most widespread tetraploid lineage (*Ruderal-4x*), which evolved a different,

260  weedy, life strategy [52], colonizing man-made habitats across central and northern Europe [53], is the

261  only lineage with traces of introgression from both a distinct diploid *A. arenosa* lineage (*Baltic-2x*),

262  [54] as well as another species – *A. lyrata* (Figs. 1, S13). Overall, this points to the ability of tetraploids

263  to accumulate diversity from various lineages, while retaining essential tetraploid- or locality-specific

264  adaptations.

265    Much work remains to understand the drivers of successful establishment and spread of newly

266    formed polyploid lineages. Relative to ecological explanations [55, 56], population genomic processes

267    have not been thoroughly assessed in natural populations despite being invoked [7, 16]. Our results

268    provide empirical insight, generally supporting pervasive yet subtle effects of ploidy on certain neutral

269    and selective processes.  Despite slightly increased nonsynonymous diversity, tetraploids may still be

270    benefiting from the masking of potentially deleterious recessive mutations, and also exhibit

271    consistently higher frequencies of adaptive nonsynonymous substitutions. Finally, multiple events of

272    strong introgression into tetraploids may provide additional substrate for local adaptation. This supports

273    the view of polyploids as diverse and adaptable evolutionary amalgamates from multiple distinct

274    ancestral lineages [57].

275

## **Online Methods**

### *Plant Material and Library Preparation*

In addition to eight previously sequenced populations, [32, 58, 59] we collected 31 new populations throughout the distribution range of *A. arenosa* (see Table S11 and Fig. S17) and its closest relative, *A. croatica*. We aimed to cover each main evolutionary lineage distinguished by previous RADseq studies [31, 34] by multiple populations, and also representatively cover the ploidy level (15 diploid, 24 tetraploid populations), altitudinal, (range 1 − 2,240 m a.s.l.) and edaphic variation (17 calcareous, 21 siliceous, 1 serpentine substrate).

We extracted DNA from silica-dried leaf tissue according to a CTAB protocol [60] with the following modifications: 75 – 100 mg of dry leaf tissue were ground in 2 mL tubes (Retsch swing mill), 200 units of RNase A per extraction were added to the isolation buffer, and the DNA pellets were washed twice with 70% ethanol. DNA was resuspended in 50 µL TE-buffer for storage, and small fragments were removed using Agencourt AMPure XP beads (Beckman Coulter, Massachusetts, USA) following the manufacturer's instructions with 0.4x DNA:beads ratio.

We quantified the extracted gDNA using the dsDNA HS assay (Q32854) from ThermoFisher Scientific (Life Technologies Ltd. Paisley, UK) with their Qubit 2.0 or 3.0 (Q33216). We prepared Illumina (Illumina United, Fulbourn, UK) Nextera XT (FC-131-1024) and TruSeq PCR-free (FC-121-3003) sequencing libraries for 350 bp insert length of genomic DNA, as well as Nextera sequencing libraries (FC-121-1030). For PCR free libraries we used 300 to 500 ng DNA as input instead of the recommended 1 µg. We quantified the NGS libraries using Qubit as described above.

### *Sequencing and Variant Calling*

We multiplexed libraries based on Qubit concentration and ran those pools on an initial

300    quantification lane. According to the yields for each sample, we increased loading of the same

301    multiplex-mix on several lanes to achieve a minimum of 10× coverage, based on the number of raw

302    reads. Samples that had less than our target coverage were remixed and run on another lane (top-up

303    lane). We sequenced 125 bp pair end reads on Illumina's HiSeq 2500 platform for all sequencing runs.

304        Our data processing pipeline involved three main parts: 1) Preparing the raw sequencing data,

305    2) Mapping and re-aligning the sequencing data and 3) Variant discovery (GATK *v.3.5*, following

306    GATK Best Practices). All steps and parameters are summarised in File S2. To prepare the raw

307    sequencing data for mapping we concatenated the fastq.gz files from the different sequencing lanes,

308    followed by trimming off the adapter sequence from reads that had inserts shorter than 250 bp, using

309    cutadapt *v.1.9* [61]. We mapped the reads to a North American *Arabidopsis lyrata* reference genome

310    [62] using bwa [63]. At this stage, we added *A. arenosa* sequencing data from previous studies [32, 58,

311    59]. For Nextera (PCR-based) libraries, we removed duplicated reads using 'MarkDuplicates' from

312    picard-tools 1.134 [64] followed by 'AddOrReplaceReadGroups' to add read groups and indices to the

313    bam files. We then used GATK *v.3.5* 'RealignerTargetCreator' and 'IndelRealigner' [65] to re-align the

314    reads around indels. Prior to variant discovery, we excluded individuals that had less than 40% of bases

315    < 4× coverage (assessed via GATK 'DepthOfCoverage' with the restriction to a minimum base quality

316    of 25 and a minimum mapping quality of 25). Our final dataset for analysis contained 287 *A. arenosa*

317    and four *A. croatica* individuals from 39 populations (see File S1 for population details and File S3 for

318    a summary of processing quality assessments).

319        We called variants for the 291 bam files (287 *A. arenosa* and four *A. croatica*) using

320    'HaplotypeCaller' and 'GenotypeGVCFs' (GATK *v.3.5*). For each bam file, 'HaplotypeCaller' was run

321    in parallel for each scaffold with ploidy specified accordingly and retaining all sites (variant and non-

322    variant). We combined the single-sample GVCF output from HaplotypeCaller to multisample GVCFs

323    and then ran 'GenotypeGVCFs' to jointly genotype these GVCFs, which greatly aids in distinguishing

324  rare variants from sequencing errors. Using GATK's 'SelectVariants', we first excluded all indel and

325  mixed sites and restricted the remaining variant sites to biallelic. Second, we removed sites that failed

326  GATK Best Practices quality recommendations (QD < 2.0, FS > 60.0, MQ < 40.0, MQRankSum < -

327  12.5, ReadPosRankSum < -8.0, HaplotypeScore < 13.0). Third, we masked genes that showed excess

328  heterozygosity (fixed heterozygous in at least five SNPs in two or more diploid populations) in the

329  dataset, i.e. potential paralogues mapped on top of each other. At the same step, we masked sites that

330  had excess read depth that we defined as 1.6× the second mode (with the first mode being low coverage

331  sites indicative of mismapping) of the read depth distribution (DP > 6400).

332

333  ***Polarization and Variant Classification***

334  We repolarized a subset of sites using a collection of genotyped individuals across closely

335  related diploid *Arabidopsis* species thus avoiding polarization against a single individual (the reference

336  genome, N. American *A. lyrata*). We used two individuals from each of the following diploid

337  *Arabidopsis* species (genotyped in the same way as our *A. arenosa* samples): European *A. lyrata*,

338  *A. croatica*, and *A. halleri*. For a site, we considered only species with complete genotypes and only

339  considered a site with at least two species represented. We required the alternative allele frequency to

340  be > 0.5 in each species, if all species were represented at a site. However, if only two species were

341  represented, we doubly weighted allele frequency for the species by preferring species with expected

342  higher genetic variation of its European populations (i.e. with decreasing priority for *A. halleri* >

343  *A. lyrata* > *A. croatica*) and required mean allele frequency to be > 0.5. In total, this identified

344  ~145,000 sites for repolarization. We classified sites as 4-fold (4-dg) or 0-fold (0-dg) degenerate based

345  on their position in the *A. lyrata* gene model annotation Araly1_GeneModels_FilteredModels6.gff. 0-

346  dg sites are those where any mutation is expected to result in an amino acid change, and 4-dg are the

347  opposite (same amino acid regardless of mutation).

348

***Population Structure***

350       We inferred relationships among the 39 *A. arenosa* and one *A. croatica* populations (the full

351    dataset, as well as each separate ploidy) based on putatively neutral 4-fold degenerate SNPs.

352    Synonymous sites are not necessarily free of constraints, e.g. due to potential codon usage bias, but are

353    nevertheless the closest to effectively neutral of any site class in the genome [66]. After quality filtering

354    our demographic analysis is based on a genome-wide dataset consisting of 1,350,328 four-fold

355    degenerate SNPs, allowing for a maximum of 10% missing alleles per site (1.2% missing data). Firstly,

356    we calculated principal component analysis (PCA) using *glPCA* function in *adegenet* [67] replacing the

357    missing values (1.2% in total) by average allele frequency for that locus. Next, we calculated Nei's [68]

358    distances among all individuals in *StAMPP* [69] and displayed it using the neighbour network

359    algorithm in *SplitsTree* [70]. Third, we selected the 553 (503 for the diploid-only dataset) most

360    parsimony-informative genes based on the following criteria: 1.) for each accession, we excluded genes

361    with ≥ 10% missing data, and 2.) we excluded genes with ≥ 10% missing accessions. We constructed a

362    maximum likelihood tree from each gene using *RAxML v.8* [77] with model GTRCAT and 100 (rapid)

363    bootstrap replicates [77]. In each gene alignment for *RAxML*, accessions were represented by the

364    consensus sequence, with different alleles represented as ambiguous sites in the consensus sequence.

365    Ambiguous sites are treated by *RAxML* as invariant sites, hence, the standard nucleotide substitution

366    model needed to be utilized; the ascertainment bias correction model that is usually used for SNP

367    matrices is not appropriate in such case. The resulting gene trees were summarized under the

368    multispecies coalescent using *Astral v.4.10.10* [78]; bootstrapping was performed with 100 replicates

369    each.

370       We further determined grouping of the populations using three clustering approaches: model-

371    based Bayesian clustering using *fastStructure v.1.0* [71] and STRUCTURE *v.2.3.2,* [72] and a non-

372   parametric k-means clustering using *adegenet* [67]. The analyses were performed separately for (i) the

373   entire data set of *A. arenosa* (*A. croatica* excluded; 9,543 SNPs after random thinning over windows of

374   50 kb to reduce effect of linkage and removing singletons, 2.4% of missing data), (ii) diploids only

375   (12,655 SNPs, 4.1% missing data) and (iii) tetraploids only (9,596 SNPs, 2.3% missing data). In

376   *fastStructure*, five replicate runs for K (number of groups) ranging from 1 to 10 were carried out under

377   default settings. We selected the optimal K value based on the similarity coefficient (~1 for optimal K

378   [73]) across replicates (Fig. S18). As *fastStructure* does not handle polyploid genotypes, we randomly

379   subsampled two alleles per each tetraploid locus (following [74]) using a custom script. To check for

380   the effect of such subsampling, we also ran the original STRUCTURE program, which handles mixed-

381   ploidy datasets, for optimal K values according to *fastStructure*. We ran the admixture model with

382   uncorrelated allele frequencies using a burn-in of 100,000 iterations followed by 1,000,000 additional

383   iterations. Finally, we ran k-means clustering using 1000 random starts and selected the partition with

384   the lowest Bayesian information criterion (BIC) value.  Population groupings were consistent across

385   algorithms (Fig. 1b, c, S1 − S5), although some methods identified finer sub-structure within the

386   *S. Carp.-2x* and *C. Europe-4x* lineages (Fig. S5).

387        We used Treemix *v.1.3* to infer migration events and relationships between the 39 *A. arenosa*

388   populations using one *A. croatica* population as outgroup. We used the 4-dg sites to build a tree without

389   any migration events and used this tree as basis for migration models to make comparisons easier

390   (option '-g'). We modelled zero to eight migrations and graphically assessed the residuals after each

391   additional migration modelled, using the R-scripts supplied with the Treemix package. If specific

392   population pairs had high residuals, we modelled an additional migration event. We continued until the

393   residuals were small and evenly spread across population pairs and/or until an additional migration

394   event involved the outgroup (we consider this admixture unlikely due to very local occurrence and

395   spatial isolation of the *A. croatica*).

396    To quantify differentiation among populations, we calculated genome-wide $F_{ST}$ and Rho

397    coefficients (similarly as in the window-based analyses described below) and performed analysis of

398    molecular variance (AMOVA) based on the Nei's distances using the *amova* function in the *pegas* R

399    package [75]. We tested for isolation by distance relationships through comparison of matrices of

400    geographic and genetic (Nei's among-population) distances among the populations using

401    *mantel.randtest* function in *ade4* R package [76]. For each tetraploid population, we calculated the

402    frequency of alleles diagnostic to each diploid lineage. The allele was defined as diagnostic if it

403    exhibited minimum frequency 0.3 (to avoid including sequencing errors as diagnostic alleles) in that

404    diploid lineage and was absent in any other diploid lineage (except for the putatively admixed Baltic

405    diploids, Table S13). For all populations we also calculated frequency of *A. lyrata*-like alleles, i.e.

406    reference alleles that were otherwise rare in the complete *A. arenosa* dataset (a rarity cut-off of 6.8%,

407    i.e., equivalent to two tetraploid populations of 8 individuals). As these alleles were nearly absent in

408    *A. arenosa* diploid populations, i.e. the ancestors of tetraploids, we assume they more likely represent

409    hybridisation from *A. lyrata* than ancestral variation shared among both species.

410    Finally, we inferred phylogenetic relationships among plastomes of our samples and previously

411    published plastomes of other *Arabidopsis* species [74]. We mapped the reads to a custom *A. arenosa*

412    plastome assembly constructed using org.ASM (http://pythonhosted.org/ORG.asm/) and performed

413    variant calling and filtration as described above, with the exception of setting ploidy = 1 in GATK

414    *HaplotypeCaller* and retaining SNPs and invariant sites with depth > 4 in at least 90% of the

415    individuals. We aligned all sequences using *Mafft* [77] and reconstructed relationships using maximum

416    likelihood in *RAxML* using GTR model with Gamma distribution of rate variation.

417

418    ***Demographic analysis***

419    We compared various demographic models and estimated parameters using the coalescent

19

420   simulation software *fastsimcoal2 v.25* [78]. The models differed in topology and presence/absence of

421   migration (admixture) events (Figs 4, S9, and S10), and each model was fit to a multi-dimensional site

422   frequency spectrum calculated from the observed four-fold degenerate SNP data. Our primary interest

423   in these analyses lie in confirming whether or not the additional populations that we sampled supported

424   the single origin of tetraploids previously determined in [31]. Specifically, we focused on populations

425   in the two diploid/tetraploid contact zones (Southern Carpathian and Baltic-Ruderal contact zones).

426      We attempted to discriminate between single versus independent origins using population

427   quartets involving representatives from both putative parental diploid lineages (*S. Carp-2x* and

428   *W. Carp.-2x* for *S. Carp.-4x*; *Baltic-2x* and *W. Carp.-2x* for *Ruderal-4x*; i.e. the genetically closest two

429   in the descriptive distance-based and clustering analyses, Fig. 1 and 4), the *W. Carp.-4x* that is

430   genetically closest to the putative ancestor of the widespread tetraploids [31] and the focal tetraploid

431   (Fig. S9 and S10). In order to maintain a realistic number of scenarios while permuting the parameters

432   (11 models for each population quartet), we modelled both uni- and bi-directional admixture within the

433   same ploidy level, but only unidirectional interploidy admixture – from diploids to tetraploids. This

434   decision reflects no signs of admixture of the diploids in clustering analyses (in contrast to the highly

435   admixed tetraploids, Fig 1B) and virtual absence of triploids in nature [33], i.e. the only possible

436   mediators of gene flow in the tetraploid-to-diploid direction [56]. In addition, we tested for the

437   potentially admixed origin of the Baltic diploids (*Baltic-2x*) [34] using population trios involving

438   representatives of each diploid lineage (*W. Carp.-2x* and *S. Carp.-2x*) as well as the focal *Baltic-2x*

439   population (Fig S19 and Table S4).

440      For each scenario and population trio/quartet, we performed 50 independent *fastsimcoal* runs to

441   overcome local maxima in the likelihood surface (see File S7 for example template file). In order to

442   minimize the population-specific effects, we ran the analyses for different iterations of well-covered

443   populations falling within the particular lineage, leading to 12 different population quartets ("natural

444    replicates") for each scenario testing the origin of the *S. Carp.-4x* and *Ruderal-4x* and four trios in the

445    *Baltic-2x* scenarios. We then extracted the best likelihood partition for each *fastsimcoal* run, calculated

446    Akaike information criterion (AIC) and summarized the AIC values across the 50 independent

447    fastsimcoal runs over the scenarios tested within each population trio/quartet. The scenario with

448    consistently lowest AIC values within particular population trio/quartet was preferred (Figs. S9 and

449    S10). In order to calculate confidence intervals for the demographic parameters (Table S9), we sampled

450    with replacement from the 4-dg SNPs to create 100 bootstrapped datasets and performed additional

451    *fastsimcoal2* analyses under the preferred scenario with these 100 distinct datasets. For these analyses

452    we also included representative (best covered) populations from the putatively non-admixed *C. Europe-*

453    *4x* lineage. Finally, we used the mutation rate of $4.3 \times 10^{-8}$ estimated by [31] to calibrate coalescent

454    simulations and obtain absolute values of population sizes and divergence times.

455        In addition, we used PSMC 0.6.4 [79] to infer changes in effective population size ($N_e$) through

456    time using information from whole-genome sequences of the *A. arenosa* diploids. We plotted 75

457    samples out of the 93 sequenced diploids, i.e. excluding samples with too low a coverage (below 12×)

458    and too much missing data. Coverage and missing data might have large effects on the PSMC estimates

459    [80]; therefore, our results should be interpreted only in conjunction with other analysis methods. We

460    ran PSMC with parameters: psmc -N25 -t15 -r5 -p "4+25*2+4+6" and then plotted the past changes in

461    $N_e$ assuming a mutation rate of $3.7 \times 10^{-8}$ substitutions per site per generation and generation time of two

462    years.

463

464    ***Window-based metric calculation***

465        In order to facilitate comparisons of windows across populations or population contrasts, we chose

466    to calculate population genetic metrics in windows defined by a given number of base pairs. We

467    repeated all calculations for two window sizes, 10kb and 50kb. We used the 50kb windows for

21

characterizing broad, genome or chromosome-level patterns, whereas the former was used for finer, gene-level analyses. For 50kb windows, patterns of LD decay suggest a minimal degree of non-independence among windows relative to the genome background (Fig. 3a).

For each of the 36 populations with at least five individuals, we excluded all individuals with $< 8\times$ average coverage, except for populations SZI, KZL, and SNO as excluding individuals from these populations would drop them below required minimum of 5 individuals. After excluding these individuals, we excluded sites if the number of missing individuals was greater than 10%, on a population-specific basis. We calculated per-site nucleotide diversity ($\theta_\pi$ and Watterson's Theta $\theta_W$; divided by the total number of sites with sufficient coverage) and Tajima's D following [81]. To equalise the expected variance of these metrics, thereby facilitating cross-population comparisons, we randomly chose 5 individuals with sufficient coverage from each population, doing so independently at each site. Differences in the diversity statistics among populations of different ploidy were tested using non-parametric Wilcoxon rank-sum test (wilcox.test in R package *stats*), taking populations as replicates.

We calculated the following divergence metrics for each possible pairwise population comparison using our custom scripts available at https://github.com/pmonnahan/ScanTools: $F_{ST}$ [82], $\rho$ [17], $d_{XY}$ [83], and the number and proportion of fixed differences. The multi-locus implementation of $F_{ST}$ and $\rho$ was translated from the software SPAGeDi [84].

***Topology weighting and detection of local introgression***

We quantified the relative support for alternative phylogenetic relationships among populations using the topology weighting approach implemented in Twisst [85]. We used only 4-fold degenerate sites and used only individuals with $> 8x$ coverage. Using bcftools, we converted the VCF files to a simplified tabular genotype file containing only the relevant individuals. We filtered this file using the

492 filterGenotypes.py script that accompanies the Twisst software. At a site, we required genotype calls for

493 at least 200 out of the 254 high coverage individuals (i.e. allowing ~20% missing data). We used only

494 biallelic sites and required that the minor allele be present in at least 2 individuals. We then ran

495 phyml_sliding_windows.py using 100 SNP windows (-w 100 and –M 20), which fits an ML

496 phylogenetic tree for each window. Ideally, Twisst should be run on phased data; however, we were

497 unable to find a workable phasing software that could handle diploids and tetraploids despite multiple

498 attempts. Instead, we used the phasing algorithm internal to Twisst, which forms haplotypes by

499 maximizing pairwise LD in each window.

500     We then ran Twisst for a number of scenarios, specifying individual population or groups of

501 populations (lineages) as taxa. Twisst implements an iterative sub-sampling algorithm based on the

502 phyML results to determine the support or weight of each possible taxon topology within each window.

503 We requested the program calculate the complete weightings (completely searching sample space) if

504 possible and used an approximate method, where sampling ceases after a given threshold of confidence

505 is reached, when necessary. We allowed for 2000 sampling iterations before opting for the backup

506 method. After this limit, we used the "Wilson" method at the 5% level, which will enforce sampling

507 until the binomial 95% confidence interval is less than 5% of the weight value.

508     We used a combination of information from Twisst as well as divergence metrics to diagnose

509 regions of both excessively strong and weak interploidy introgression in the two highly admixed

510 *S. Carp.-4x* and *Ruderal-4x* lineages. First, introgressed regions should show an elevated weight for

511 topologies wherein the proximal diploid/tetraploid pair are placed sister to one another (Topology 3 in

512 Fig. 5). Second, when comparing the focal tetraploid to other tetraploid populations, an introgressed

513 region should show elevated divergence while at the same time exhibiting reduced divergence to the

514 focal diploid population. Conversely, introgression-resistant regions should show elevated Topology 1

515 and a combination of low divergence from tetraploids and elevated divergence from all diploids. We

516  looked for evidence of selection on introgressed regions by overlapping window outliers for Topology

517  3 and Fay and Wu's H (in 10kb windows) in the focal tetraploid (99[th] percentile for both metrics).

518

519  ***Gene expression analysis of purifying selection***

520      We evaluated patterns of diversity at the gene level using gene expression levels as a proxy for

521  selective pressure based on evidence that higher-expressed genes generally show stronger signs of

522  purifying selection in both plants and animals [36, 86-88]. To obtain gene-wise estimates of diversity,

523  we performed a separate mapping process (again, using *A. lyrata* as the reference genome) using a

524  subset of the total *A. arenosa* dataset that covers all major diploid and tetraploid lineages (9 tetraploid

525  and 9 diploid populations, comprising 74 and 70 individuals, respectively, listed in Table S12). We

526  retained sites with read depth of 4 or higher for at least 5 individuals across each population (9 – 14

527  million sites per population, Table S12). Sites were downsampled to 5 individuals independently at

528  each site to homogenize chromosome depth across sites.

529      First, we extracted RNA from leaves of 3-week old individuals with three biological replicates for

530  each of three diploid populations (HNI, RZA, SNO) to complete our previous dataset [54] of seven

531  tetraploid populations (TBG, BGS, STE, KAS, CA2, HOC, SWA) using the RNeasy Plant Mini Kit

532  (Qiagen). We synthesized single-strand cDNA from 500ng of total RNA using VN-anchored poly-

533  T(23) primers with MuLV Reverse Transcriptase (Enzymatics) according to the manufacturer's

534  recommendations. We made RNAseq libraries using the TruSeq RNA Sample Prep Kit v2 (Illumina)

535  and sequenced libraries on an Illumina HiSeq 2000 with 50bp single-end reads. We sequenced between

536  9.8 and 18.8 million reads (avg 13.6 million). We aligned reads to the *A. lyrata* genome using TopHat2

537  [89] and re-aligned unmapped reads using Stampy [90]. We acquired read counts for each of the 32,670

538  genes using HTseq-count [91] with *A. lyrata* gene models. We normalized for sequencing depth using

539  DEseq2 in R, [92] and further analyses were performed in MATLAB (MathWorks).

540    Analysis of differential expression between diploid and tetraploid expression patterns were

541    performed using a one-way analysis of variance (ANOVA), and *p*-values were corrected for false

542    discovery rate [93]. To avoid low-expression genes, we filtered for genes presenting a least one sample

543    with normalized counts above 25, and computed the log-ratio of the average population expression in

544    tetraploid populations against the average expression in diploids (positive when the expression of a

545    gene is higher in tetraploid and negative when it is higher in diploids).

546    We obtained 6,504 genes with statistically significant differential expression ($p < 0.05$) between

547    diploids and tetraploids (33% of 19,319 genes), but only 321 of these presented fold-change above

548    1.78x (5% two-tail threshold, Fig. S20A) and 214 above 2x. Overall, the average mean expression

549    across populations is very strongly correlated between ploidies (slope = 1.02, $R^2 = 0.93$, Fig. S20B). To

550    estimate mutational patterns we limited ourselves to the set of 18,998 genes non-differentially

551    expressed (NDE) between ploidies.

552    We then filtered genes exhibiting a dependence of diversity metrics on the number of sites,

553    specifically the genes that showed a correlation of number of sites with diversity (indicating potential

554    mis-mapping of reads; Fig. S21). This effect of 4-dg $\theta_\pi$ and $\theta_W$ was strong for genes with fewer than 20

555    sites or more than a 100 using a locally weighted linear regression (LOWESS) for genes with a

556    minimum of 5 sites of each fold (0-dg and 4-dg). Between these two boundaries, the number of sites

557    only has a weak effect on 4-dg diversity. We observed a similar pattern in terms of 0-dg diversity with

558    loci with fewer than 30 or more than 400 0-dg sites (Fig. S21C&D). After exclusion of loci outside of

559    these bounds (for both 4-dg and 0-dg) from any downstream analysis we were able to cover around

560    45% of all NDE genes.

561    We then visualized the correlation of diversity of each gene with the average gene expression

562    within the ploidy of the population with a locally weighted linear regression (LOWESS). For genes

563    with expression levels above a certain expression threshold (50), nonsynonymous diversity (0-dg $\theta_\pi$

564 and $\theta_W$) showed a clear negative correlation with expression (proxy for strength of purifying selection)

565 for both ploidies (Fig. 2a, Fig. S6: bold lines). Notably, this trend seems to break for very high

566 expression (>2250 i.e. top 0.35%), possibly due to the low coverage of this expression range (67

567 genes). After removal of genes outside of these thresholds, we obtained 5,900 NDE genes per

568 population to be used for multiple linear model (MLM) fitting.

569    We evaluated the effect of gene expression on 0-dg/4-dg diversity ratio and on 0-dg diversity for

570 each population by modelling them (y) as a function of the ploidy of the population (p) with coefficient

571 $\alpha_p$, the average gene expression measured in ploidy p ($E_p$) with coefficient $\beta$, and an interaction term $\gamma_p$

572 as follows:

573 $$(y) \sim 1 + \alpha_p + \beta * \log(E_p) + \gamma_p * \log(E_p)$$

574 To estimate $N_g$, we first estimated effective population sizes using synonymous diversity as an

575 estimator of $\theta$, the estimated mutation rate ($\mu$) of $4.3 \times 10^{-8}$ for *A. arenosa* [31] and their theoretical

576 relationship given by $\theta = 4\mu N_e$ in diploids and $\theta = 8\mu N_e$ in tetraploids. This gave an estimate of

577 effective population sizes around 240,000 individuals for diploids and around 130,000 for tetraploids.

578 In terms of number of haploid genomes, this difference in effective census sizes is more than

579 compensated by tetrasomy (~480,000 in tetraploids vs ~520,000 in diploids).

580    The second MLM equation for evaluating the impact of population size ($N_g$) on 0-dg diversity or

581 on 0-dg/4-dg diversity ratio was established using stepwise regression, evaluating the addition or

582 removal of each term based on the *p*-value for an *F*-test of the change in the sum of squared error. The

583 final formula for 0-dg diversity was:

584 $$(0\text{-dg } \theta_W) \sim 1 + \alpha_p + \beta * \log(E_p) + \delta * N_g + \gamma_N * \log(E_p) + \varepsilon_p * N_g$$

585 where the interaction term with log expression $\gamma$ is now dependent on $N_g$, $\delta$ represents the fixed effect

586 of $N_g$, with an additional interaction term $\varepsilon_p$ dependent on ploidy (p). The final formula for the 0-dg/4-

587    dg diversity ratio was:

588    $$(0\text{-dg } \theta_W) \sim 1 + \alpha_p + \beta*\log(E_p) + \delta*N_g + \gamma_p*\log(E_p)$$

589    where the interaction term with log expression $\gamma$ is now dependent on ploidy (p) only.

590    The MLM estimates are presented in Table S5 and S6, and the estimated effects for values of the

591    predictor chosen to show large responses are plotted in Fig 2b: log Expression: 3.9124 to 7.7098; Ng:

592    366058 (low) to 488976 (med) to 611894 (high).

593        In addition, we calculated recessive load as a number of sites with derived allele in homozygote

594    state per each individual with at least 5 million SNPs called (240 individuals in total) and tested for

595    difference among population means of diploid and tetraploid populations using Wilcoxon rank-sum

596    test.

597

598    ***Distribution of fitness effects***

599        Using the allele frequency spectra (AFS) for 4-dg and 0-dg sites (separately) for each of the 36

600    populations with $\geq 5$ individuals screened, we estimated the distribution of fitness effects (DFE) [39],

601    the proportion of adaptive substitutions relative to the total number of nonsynonymous substitutions ($\alpha$)

602    [40], and the proportion of adaptive substitutions relative to neutral divergence ($\omega_a$; [94]; DFE-alpha

603    v2.16; http://www.homepages.ed.ac.uk/pkeightl/software.html). This method implements a maximum-

604    likelihood-based procedure to jointly estimate the parameters of a gamma-distributed DFE and a simple

605    stepwise population size change model from site frequency spectrum data. Divergence was obtained

606    from the polarized unfolded AFS and used to estimate $\alpha$ and $\omega_a$, while correcting for the effect of

607    slightly deleterious mutations using the estimated DFE. For all parameters estimated, we obtained 95%

608    confidence intervals by analyses of 200 bootstrapped data sets. For each population, we fit two

609    demographic models (constant population size and stepwise population size change), selected the best-

610    fit model using a likelihood ratio test (LRT), and then estimated the parameters of the DFE, $\alpha$ and $\omega_a$

611     under this model. The DFE is estimated using a gamma distribution with a shape parameter ($\beta$) and a

612     scale parameter that represents the strength of purifying selection. As the strength of selection is

613     dependent of the effective population size $N_e$, the result of DFE are often summarized by binning the

614     distribution in 3 bins of $-N_e*s$. A $-N_e s$ of 0-1 represents nearly neutral sites, 1-10 mildly deleterious

615     mutations, and > 10 highly deleterious mutations.

616         For all populations, the stepwise population size change model was preferred. We ran DFE-alpha

617     using both unfolded and folded site-frequency spectra. As the results were very consistent using the

618     folded or unfolded allele frequency spectrum, we chose to focus on estimates based on the folded

619     spectra, which should be more robust. We tested whether diploids and tetraploids differed with respect

620     to the proportion of new nonsynonymous mutations in each bin, using Wilcoxon rank-sum tests.

621

622     ***Linked selection analysis and calculation of genotypic associations (linkage disequilibrium)***

623         We inspected the relationship between the excess nonsynonymous divergence ($d_{XY}$) relative to

624     synonymous divergence, as a proxy for divergent selection, and synonymous diversity ($\theta_\pi$) in 50kb

625     windows [95]. Both nonsynonymous and synonymous divergence was calculated for each population

626     in each window as the average divergence at (non)synonymous sites for all pairwise contrasts between

627     the focal population and all other populations in the dataset. We natural-log transformed these values

628     and standardized them to be on the same scale. Then, we simply took the difference between

629     nonsynonymous and synonymous scaled, transformed divergence values in each window. We refer to

630     this difference as $E_{NS}$.

631         We also square root-transformed $\theta_\pi$ for normality purposes, removed windows with fewer than 20

632     SNPs, and removed populations with fewer than 2,000 non-missing windows, retaining a total of 27

633     populations (10 diploid and 17 tetraploid, listed in Supplementary File S1) and an average of 2,660

634     windows per population (~60% of genome). A negative relationship with $\theta_\pi$ is interpreted as evidence

28

635  of a reductive effect of selection on linked, neutral diversity (i.e. linked selection). More specifically,

636  we were interested to see if this relationship was dependent on ploidy level.

637  We used a multiple regression approach to infer this relationship and its dependence on ploidy

638  level. We also included information on gene density (the proportion of bases in the window occupied

639  by genic sequences according to the *A. lyrata* annotation) and proportion of missing data in each

640  window. When calculating missingness in each window, we considered all biallelic sites and simply

641  averaged the proportion of missing data across all 287 individuals in the study at each site within the

642  window. Given the strong negative relationship between gene density and missingness, we combined

643  gene density and missingness into a single, compound variable

644  $$GDM = gene\ density * (1 - missingness)$$

645  where high values indicate windows with high gene density and low missingness and low values

646  indicate the opposite. We fit a mixed model, using restricted maximum likelihood ('lmer' function) via

647  the R package *lme4* [96], with $E_{NS}$ and GDM as continuous variables, ploidy as fixed categorical

648  variable, and populations as a random categorical variable. We determined significance using Wald

649  Chi-square tests in the *car* package.  We used the default distribution family (i.e. normal distribution) to

650  model the residual variance. We also included a quadratic effect of $E_{NS}$ to investigate the possibility of

651  a nonlinear relationship with neutral diversity. Our initial model included all possible interactions, and

652  we selected our final model by eliminating non-significant higher order interaction terms.  The results

653  were not qualitatively different following removal of tetraploid populations admixed by non-sister

654  diploids (*S. Carp -4x*: DRA, LAC and TZI, and *Ruderal-4x:* KOW, STE and TBG). Similarly, results

655  did not change upon removing the apparent outlier associated with the maximum observed 4-fold

656  diversity.

657  To calculate genotypic correlations (a proxy for linkage disequilibrium), we recoded genotypes at

658  all sites on chromosome 2 to represent the number of alternative alleles (0 - 2 for diploids and 0 - 4 for

29

659  tetraploids). We calculated $r^2$ for pairs of loci and is simply the square of the correlation coefficient. An

660  $r$ value of 1.0 (and thus an $r^2$ of 1.0) means that genotypes are perfectly correlated for a particular pair

661  of loci. However, because we do not have phase information, this $r^2$ value is not equivalent to the $r^2$

662  often reported when discussing LD. Therefore, we do not technically measure LD, but rather a related

663  measure of genotypic associations. To visualize LD decay (Fig. 4a), we averaged $r^2$ value for all pairs

664  of loci that fall in bins of a given distance apart, only considering populations with 8 or greater

665  individuals. For populations with >8 individuals, we downsampled to include only the 8 highest

666  coverage individuals. We performed the $r^2$ calculation for each population separately to avoid

667  confounding effects of population differentiation.

668      To observe the impacts of various factors in our data on our LD approximation, we first simulated

669  unlinked data and varied the number of sites and individuals as well as ploidy. At each site, we

670  randomly drew allele frequencies from a uniform distribution, and then drew genotypes from the

671  binomial distribution with $p$ equal to the drawn allele frequencies and $n$ of 2 or 4, depending on ploidy.

672  The average $r^2$ value for each data set indicates that the number of individuals is the primary

673  determinant of the expected $r^2$ value for unlinked sites, with the other factors exhibiting a negligible

674  effect. We also simulated neutral linked data (100 replicate data sets; 1.5 Mb sequences; recombination

675  rate = $1 \times 10^{-8}$, mutation rate = $1 \times 10^{-8}$ ; $N_e = 100,000$) using msprime [97]. From each replicate, we

676  created diploid and tetraploid genotype data by grouping the simulated haplotypes into sets of 2 or 4,

677  respectively, to create 10 individuals of each. We then calculated $r^2$ for each replicate simulation. We

678  simulated data for multiple parameter sets, focusing on the effects of the mutation rate, population size,

679  recombination rate, and ploidy on $r^2$. We observed a slight downward bias for tetraploid data generally,

680  but this effect was negligible compared to the effects of recombination rate and population size. As

681  expected, the mutation rate did not affect $r^2$ as this measure should be proportional to the population-

682  scaled recombination rate (a function of the per-base recombination rate and population size). If we

683  double diversity by doubling the mutation rate, we find no observable effect on $r^2$ (Fig. S22).

684

**Code Availability**

686  Custom scripts used to generate genome scan metrics are available at

687  https://github.com/pmonnahan/ScanTools.  Other analysis scripts are available at

688  https://github.com/pmonnahan/ArenosaPloidy.

689

**Data Availability**

691  Sequence data that support the findings of this study have been deposited in the Sequence Read

692  Archive (SRA; https://www.ncbi.nlm.nih.gov/sra) with the primary accession code PRJNA484107

693  (available at http://www.ncbi.nlm.nih.gov/bioproject/484107) and PRJNA472485 for RNAseq data.

694

710 **Author Contributions**

711 LY, KB, FK, PB and PM conceived the study. PM, FK, PB, BL, CS, JK, RH, RS and PP performed

712 analyses with input from LY, KB, RH, and TS. CS, PB, GF, MB and CW performed laboratory

713 experiments. PM, FK and PB wrote the manuscript with primary input from KB, LY, BA, CS and TS.

714 All authors edited and approved of the final manuscript.

715

716 **Competing Interests statement**

717 The authors declare no competing interests.

718

719 **Materials & Correspondence**

720 Correspondence and material requests should be addressed to Levi Yant at levi.yant@nottingham.ac.uk

721


722

723 **References:**

724 1.      Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, Rieseberg LH. The
725 frequency of polyploid speciation in vascular plants. Proceedings of the National Academy of Sciences.
726 2009;106(33):13875-9. doi: 10.1073/pnas.0811575106.
727 2.      Van de Peer Y, Mizrachi E, Marchal K. The evolutionary significance of polyploidy. Nature
728 Reviews Genetics. 2017;18:411. doi: 10.1038/nrg.2017.26.
729 3.      Salman-Minkov A, Sabath N, Mayrose I. Whole-genome duplication as a key factor in crop
730 domestication. Nature plants. 2016;2:16115.
731 4.      Storchova Z, Pellman D. From polyploidy to aneuploidy, genome instability and cancer. Nature
732 reviews Molecular cell biology. 2004;5(1):45-54.
733 5.      Yant L, Bomblies K. Genome management and mismanagement—cell-level opportunities and
734 challenges of whole-genome duplication. Genes & development. 2015;29(23):2405-19.
735 6.      Levin DA. The role of chromosomal change in plant evolution: Oxford University Press; 2002.
736 7.      Parisod C, Holderegger R, Brochmann C. Evolutionary consequences of autopolyploidy. New
737 phytologist. 2010;186(1):5-17.

738  8.      te Beest M, Le Roux JJ, Richardson DM, Brysting AK, Suda J, Kubešová M, et al. The more
739  the better? The role of polyploidy in facilitating plant invasions. Annals of Botany. 2011;109(1):19-45.
740  9.      Segraves KA. The effects of genome duplications in a community context. New Phytologist.
741  2017.
742  10.     Haldane JBS. The causes of evolution: Princeton University Press; 1932.
743  11.     Wright S. The distribution of gene frequencies in populations of polyploids. Proceedings of the
744  National Academy of Sciences. 1938;24(9):372-7.
745  12.     Fisher R. The theoretical consequences of polyploid inheritance for the mid style form of
746  Lythrum salicaria. Annals of Human Genetics. 1941;11(1):31-8.
747  13.     Stebbins GL. Chromosomal evolution in higher plants. Chromosomal evolution in higher
748  plants. 1971.
749  14.     Haldane JB. Theoretical genetics of autopolyploids. Journal of Genetics. 1930;22(3):359-72.
750  15.     Bever JD, Felber F. The theoretical population genetics of autopolyploidy. Oxford surveys in
751  evolutionary biology. 1992;8:185-.
752  16.     Otto SP, Whitton J. Polyploid Incidence and Evolution. Annual Review of Genetics.
753  2000;34(1):401-37. doi: 10.1146/annurev.genet.34.1.401. PubMed PMID: 11092833.
754  17.     Ronfort J, Jenczewski E, Bataillon T, Rousset F. Analysis of population structure in
755  autotetraploid species. Genetics. 1998;150(2):921-30.
756  18.     Grant V. Plant speciation: New York: Columbia University Press xii, 563p.-illus., maps, chrom.
757  nos.. En 2nd edition. Maps, Chromosome numbers. General (KR, 198300748); 1981.
758  19.     Coyne JA, Orr HA. Speciation. Sunderland, MA. Sinauer Associates, Inc; 2004.
759  20.     Mallet J. Hybrid speciation. Nature. 2007;446(7133):279.
760  21.     Slotte T, Huang H, Lascoux M, Ceplitis A. Polyploid speciation did not confer instant
761  reproductive isolation in Capsella (Brassicaceae). Molecular Biology and Evolution. 2008;25(7):1472-
762  81.
763  22.     Zohren J, Wang N, Kardailsky I, Borrell JS, Joecker A, Nichols RA, et al. Unidirectional
764  diploid–tetraploid introgression among British birch trees with shifting ranges shown by restriction
765  site-associated markers. Molecular ecology. 2016;25(11):2413-26.
766  23.     Lafon-Placette C, Johannessen IM, Hornslien KS, Ali MF, Bjerkan KN, Bramsiepe J, et al.
767  Endosperm-based hybridization barriers explain the pattern of gene flow between Arabidopsis lyrata
768  and Arabidopsis arenosa in Central Europe. Proceedings of the National Academy of Sciences.
769  2017:201615123.
770  24.     Ronfort J. The mutation load under tetrasomic inheritance and its consequences for the
771  evolution of the selfing rate in autotetraploid species. Genetics Research. 1999;74(1):31-42.
772  25.     Hill R. Selection in autotetraploids. TAG Theoretical and Applied Genetics. 1971;41(4):181-6.
773  26.     Selmecki AM, Maruvka YE, Richmond PA, Guillet M, Shoresh N, Sorenson AL, et al.
774  Polyploidy can drive rapid adaptation in yeast. Nature. 2015;519(7543):349-52.
775  27.     Schmickl R, Marburger S, Bray S, Yant L, Henderson I. Hybrids and horizontal transfer:
776  introgression allows adaptive allele discovery. Journal of Experimental Botany. 2017.
777  28.     Arnold ML, Kunte K. Adaptive Genetic Exchange: A Tangled History of Admixture and
778  Evolutionary Innovation. Trends in ecology & evolution. 2017;32(8):601-11.
779  29.     Bomblies K, Madlung A. Polyploidy in the Arabidopsis genus. Chromosome Research.
780  2014;22(2):117-34. doi: 10.1007/s10577-014-9416-x.
781  30.     Yant L, Bomblies K. Genomic studies of adaptive evolution in outcrossing Arabidopsis species.
782  Current Opinion in Plant Biology. 2017;36:9-14. doi: https://doi.org/10.1016/j.pbi.2016.11.018.
783  31.     Arnold B, Kim S-T, Bomblies K. Single Geographic Origin of a Widespread Autotetraploid
784  Arabidopsis arenosa Lineage Followed by Interploidy Admixture. Molecular Biology and Evolution.
785  2015;32(6):1382-95. doi: 10.1093/molbev/msv089.

786 32.	Hollister JD, Arnold BJ, Svedin E, Xue KS, Dilkes BP, Bomblies K. Genetic adaptation
787 associated with genome-doubling in autotetraploid Arabidopsis arenosa. PLoS genetics.
788 2012;8(12):e1003093.
789 33.	Kolář F, Lučanová M, Záveská E, Fuxová G, Mandáková T, Španiel S, et al. Ecological
790 segregation does not drive the intricate parapatric distribution of diploid and tetraploid cytotypes of the
791 Arabidopsis arenosa group (Brassicaceae). Biological Journal of the Linnean Society. 2016;119(3):673-
792 88.
793 34.	Kolář F, Fuxová G, Záveská E, Nagano AJ, Hyklová L, Lučanová M, et al. Northern glacial
794 refugia and altitudinal niche divergence shape genome-wide differentiation in the emerging plant model
795 Arabidopsis arenosa. Molecular ecology. 2016;25(16):3929-49.
796 35.	1001 Genomes Consortium. 1,135 genomes reveal the global pattern of polymorphism in
797 Arabidopsis thaliana. Cell2016. p. 481-91.
798 36.	Ingvarsson PK. Gene expression and protein length influence codon usage and rates of
799 sequence evolution in Populus tremula. Molecular biology and evolution. 2007;24(3):836-44.
800 37.	Wright SI, Yau CK, Looseley M, Meyers BC. Effects of gene expression on molecular evolution
801 in Arabidopsis thaliana and Arabidopsis lyrata. Molecular biology and evolution. 2004;21(9):1719-26.
802 38.	Popescu CE, Borza T, Bielawski JP, Lee RW. Evolutionary rates and expression level in
803 Chlamydomonas. Genetics. 2006;172(3):1567-76.
804 39.	Keightley PD, Eyre-Walker A. Joint inference of the distribution of fitness effects of deleterious
805 mutations and population demography based on nucleotide polymorphism frequencies. Genetics.
806 2007;177(4):2251-61.
807 40.	Eyre-Walker A, Keightley PD. Estimating the rate of adaptive molecular evolution in the
808 presence of slightly deleterious mutations and population size change. Molecular biology and
809 evolution. 2009;26(9):2097-108.
810 41.	Rousselle M, Mollion M, Nabholz B, Bataillon T, Galtier N. Overestimation of the adaptive
811 substitution rate in fluctuating populations. Biology Letters. 2018;14(5).
812 42.	Venkat A, Hahn MW, Thornton JW. Multinucleotide mutations cause false inferences of
813 lineage-specific positive selection. Nat Ecol Evol. 2018;2(8):1280-8. doi:
814 papers3://publication/doi/10.1038/s41559-018-0584-5.
815 43.	Schmickl R, Koch MA. Arabidopsis hybrid speciation processes. Proceedings of the National
816 Academy of Sciences. 2011;108(34):14192-7.
817 44.	Gerstein AC, Otto SP. Ploidy and the causes of genomic evolution. Journal of Heredity.
818 2009;100(5):571-81.
819 45.	Favarger C. Cytogeography and biosystematics. Plant biosystematics. 1984:453-76.
820 46.	Brochmann C, Brysting A, Alsos I, Borgen L, Grundt H, Scheen A-C, et al. Polyploidy in arctic
821 plants. Biological journal of the Linnean society. 2004;82(4):521-36.
822 47.	Butruille DV, Boiteux LS. Selection–mutation balance in polysomic tetraploids: Impact of
823 double reduction and gametophytic selection on the frequency and subchromosomal localization of
824 deleterious mutations. Proceedings of the National Academy of Sciences. 2000;97(12):6608-13. doi:
825 10.1073/pnas.100101097.
826 48.	Willis JH. Inbreeding Load, Average Dominance and the Mutation Rate for Mildly Deleterious
827 Alleles in &lt;em&gt;Mimulus guttatus&lt;/em&gt. Genetics. 1999;153(4):1885.
828 49.	Schmickl R, Marburger S, Bray S, Yant L. Hybrids and horizontal transfer: introgression allows
829 adaptive allele discovery. Journal of experimental botany. 2017;68(20):5453-70.
830 50.	Lowe WH, Muhlfeld CC, Allendorf FW. Spatial sorting promotes the spread of maladaptive
831 hybridization. Trends in Ecology & Evolution. 2015;30(8):456-62. doi:
832 https://doi.org/10.1016/j.tree.2015.05.008.
833 51.	Yukilevich R. ASYMMETRICAL PATTERNS OF SPECIATION UNIQUELY SUPPORT

834 REINFORCEMENT IN DROSOPHILA. Evolution. 2012;66(5):1430-46. doi: 10.1111/j.1558-
835 5646.2011.01534.x.
836 52. Baduel P, Arnold B, Weisman CM, Hunter B, Bomblies K. Habitat-associated life history and
837 stress-tolerance variation in Arabidopsis arenosa. Plant physiology. 2016;171(1):437-51.
838 53. Hylander N. Cardaminopsis suecica (Fr.) Hiit., a northern amphidiploid species. Bulletin du
839 Jardin botanique de l'Etat, Bruxelles/Bulletin van den Rijksplantentuin, Brussel. 1957:591-604.
840 54. Baduel P, Hunter B, Yeola S, Bomblies K. Genetic basis and evolution of rapid cycling in
841 railway populations of tetraploid Arabidopsis arenosa. PLOS Genetics. 2018;14(7):e1007510. doi:
842 10.1371/journal.pgen.1007510.
843 55. Husband BC, Sabara HA. Reproductive isolation between autotetraploids and their diploid
844 progenitors in fireweed, Chamerion angustifolium (Onagraceae). New Phytologist. 2004;161(3):703-
845 13.
846 56. Kolář F, Čertner M, Suda J, Schönswetter P, Husband BC. Mixed-Ploidy Species: Progress and
847 Opportunities in Polyploid Research. Trends in Plant Science. 2017.
848 57. Soltis DE, Soltis PS. Polyploidy: recurrent formation and genome evolution. Trends in Ecology
849 & Evolution. 1999;14(9):348-52.
850 58. Yant L, Hollister JD, Wright KM, Arnold BJ, Higgins JD, Franklin FCH, et al. Meiotic
851 adaptation to genome duplication in Arabidopsis arenosa. Current biology. 2013;23(21):2151-6.
852 59. Arnold BJ, Lahner B, DaCosta JM, Weisman CM, Hollister JD, Salt DE, et al. Borrowed alleles
853 and convergence in serpentine adaptation. Proceedings of the National Academy of Sciences.
854 2016;113(29):8320-5.
855 60. Doyle JJ. A rapid DNA isolation procedure for small quantities of fresh leaf tissue. Phytochem
856 Bull. 1987;19:11-5.
857 61. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
858 EMBnet journal. 2011;17(1):pp. 10-2.
859 62. Hu TT, Pattyn P, Bakker EG, Cao J, Cheng J-F, Clark RM, et al. The Arabidopsis lyrata genome
860 sequence and the basis of rapid genome size change. Nature genetics. 2011;43(5):476-81.
861 63. Li H, Durbin R. Fast and accurate short read alignment with Burrows–Wheeler transform.
862 Bioinformatics. 2009;25(14):1754-60.
863 64. DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A framework for
864 variation discovery and genotyping using next-generation DNA sequencing data. Nature genetics.
865 2011;43(5):491-8.
866 65. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome
867 Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data.
868 Genome research. 2010;20(9):1297-303.
869 66. Wright SI, Lauga B, Charlesworth D. Rates and patterns of molecular evolution in inbred and
870 outbred Arabidopsis. Molecular Biology and Evolution. 2002;19(9):1407-20.
871 67. Jombart T. adegenet: a R package for the multivariate analysis of genetic markers.
872 Bioinformatics. 2008;24(11):1403-5.
873 68. Nei M. Genetic distance between populations. The American Naturalist. 1972;106(949):283-92.
874 69. Pembleton LW, Cogan NO, Forster JW. StAMPP: an R package for calculation of genetic
875 differentiation and structure of mixed-ploidy level populations. Molecular ecology resources.
876 2013;13(5):946-52.
877 70. Huson DH. SplitsTree: analyzing and visualizing evolutionary data. Bioinformatics (Oxford,
878 England). 1998;14(1):68-73.
879 71. Raj A, Stephens M, Pritchard JK. fastSTRUCTURE: variational inference of population
880 structure in large SNP data sets. Genetics. 2014;197(2):573-89.
881 72. Pritchard JK, Stephens M, Donnelly P. Inference of population structure using multilocus

882    genotype data. Genetics. 2000;155(2):945-59.

883    73.    Nordborg M, Hu TT, Ishino Y, Jhaveri J, Toomajian C, Zheng H, et al. The pattern of
884    polymorphism in Arabidopsis thaliana. PLoS biology. 2005;3(7):e196.

885    74.    Novikova PY, Hohmann N, Nizhynska V, Tsuchimatsu T, Ali J, Muir G, et al. Sequencing of the
886    genus Arabidopsis identifies a complex history of nonbifurcating speciation and abundant trans-specific
887    polymorphism. Nature genetics. 2016;48(9):1077-82.

888    75.    Paradis E. pegas: an R package for population genetics with an integrated–modular approach.
889    Bioinformatics. 2010;26(3):419-20.

890    76.    Dray S, Dufour A-B. The ade4 package: implementing the duality diagram for ecologists.
891    Journal of statistical software. 2007;22(4):1-20.

892    77.    Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
893    improvements in performance and usability. Molecular biology and evolution. 2013;30(4):772-80.

894    78.    Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. Robust demographic inference
895    from genomic and SNP data. PLoS genetics. 2013;9(10):e1003905.

896    79.    Li H, Durbin R. Inference of human population history from individual whole-genome
897    sequences. Nature. 2011;475(7357):493.

898    80.    Nadachowska-Brzyska K, Burri R, Smeds L, Ellegren H. PSMC analysis of effective population
899    sizes in molecular ecology and its application to black-and-white Ficedula flycatchers. Molecular
900    ecology. 2016;25(5):1058-72.

901    81.    Zeng K, Fu Y-X, Shi S, Wu C-I. Statistical tests for detecting positive selection by utilizing
902    high-frequency variants. Genetics. 1996;174(3):1431-9.

903    82.    Weir BS, Cockerham CC. Estimating F-statistics for the analysis of population structure.
904    evolution. 1984;38(6):1358-70.

905    83.    Cruickshank TE, Hahn MW. Reanalysis suggests that genomic islands of speciation are due to
906    reduced diversity, not reduced gene flow. Molecular ecology. 2014;23(13):3133-57.

907    84.    Hardy OJ, Vekemans X. SPAGeDi: a versatile computer program to analyse spatial genetic
908    structure at the individual or population levels. Molecular Ecology Resources. 2002;2(4):618-20.

909    85.    Martin SH, Van Belleghem SM. Exploring Evolutionary Relationships Across the Genome
910    Using Topology Weighting. Genetics. 2017. doi: 10.1534/genetics.116.194720.

911    86.    Duret L, Mouchiroud D. Determinants of substitution rates in mammalian genes: expression
912    pattern affects selection intensity but not mutation rate. Molecular biology and evolution.
913    2000;17(1):68-070.

914    87.    Rocha EP, Danchin A. An analysis of determinants of amino acids substitution rates in bacterial
915    proteins. Molecular biology and evolution. 2004;21(1):108-16.

916    88.    Slotte T, Bataillon T, Hansen TT, St. Onge K, Wright SI, Schierup MH. Genomic Determinants
917    of Protein Evolution and Polymorphism in Arabidopsis. Genome Biology and Evolution. 2011;3:1210-
918    9. doi: 10.1093/gbe/evr094.

919    89.    Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R, Salzberg SL. TopHat2: accurate alignment
920    of transcriptomes in the presence of insertions, deletions and gene fusions. Genome biology.
921    2013;14(4):R36.

922    90.    Lunter G, Goodson M. Stampy: a statistical algorithm for sensitive and fast mapping of
923    Illumina sequence reads. Genome research. 2011;21(6):936-9.

924    91.    Anders S, Pyl PT, Huber W. HTSeq—a Python framework to work with high-throughput
925    sequencing data. Bioinformatics. 2015;31(2):166-9.

926    92.    Love M, Anders S, Huber W. Differential analysis of count data–the DESeq2 package. Genome
927    Biology. 2014;15:550.

928    93.    Benjamini Y, Hochberg Y. Controlling the false discovery rate: a practical and powerful
929    approach to multiple testing. Journal of the royal statistical society Series B (Methodological).

930     1995:289-300.
931     94.     Gossmann TI, Song B-H, Windsor AJ, Mitchell-Olds T, Dixon CJ, Kapralov MV, et al. Genome
932     wide analyses reveal little evidence for adaptive evolution in many plant species. Molecular biology
933     and evolution. 2010;27(8):1822-32.
934     95.     Martin SH, Möst M, Palmer WJ, Salazar C, McMillan WO, Jiggins FM, et al. Natural selection
935     and genetic diversity in the butterfly Heliconius melpomene. Genetics. 2016;203(1):525-41.
936     96.     Bates D, Martin M, Ben B, Walker S. lme4: Linear mixed effects models using Eigen and S4.(R
937     package v. 1.0–6). See http://CRAN. R-project. org/package= lme4; 2014.
938     97.     Kelleher J, Etheridge AM, McVean G. Efficient Coalescent Simulation and Genealogical
939     Analysis for Large Sample Sizes. PLOS Computational Biology. 2016;12(5):e1004842. doi:
940     10.1371/journal.pcbi.1004842.
941

942
943 **Figure Legends**

944
945 **Fig. 1 | Geographic distribution and range-wide genetic variation of *Arabidopsis arenosa*. a**,

946 Distribution of the 39 *A. arenosa* populations (red labels - diploids, blue - tetraploids) with average

947 proportions of cluster membership inferred by FastStructure (panel b). Color shades highlight highly

948 admixed tetraploid populations (*Ruderal* and *S. Carpathian-4x*) together with the diploid sources of

949 admixture. **b**, Posterior probabilities of cluster membership of the 287 *A. arenosa* individuals as

950 inferred by FastStructure under K=6. **c**, Neighbor-joining tree of Nei's genetic distances between all

951 individuals and the outgroup *Arabidopsis croatica*. Individuals from admixed populations are

952 highlighted correspondingly. Inset: distribution of pairwise genetic divergence of populations ($\rho$)

953 within each ploidy. **d**, Principal component analysis of all but the two most divergent diploid

954 (*Pannonian* and *Dinaric*) *A. arenosa* lineages (shades correspond to admixed populations).

955

956 **Fig. 2 | Effects of ploidy on purifying selection, genetic load, and the distribution of fitness effects**

957 **(DFE). a**, Genic nonsynonymous (0-dg) diversity versus average gene expression (log-scale) for each

958 population and each ploidy (resp. faint and bold LOWESS curves). The two outlier populations, *2x*-

959 SNO and *4x*-SCH, are indicated. **b**, Standardized effects with confidence intervals in multiple linear

960 model of haploid effective population size ($N_g$), ploidy, and levels of expression on nonsynonymous (0-

961 dg) $\theta_W$ (upper panel) and on 0-dg/4-dg $\theta_W$ ratio (lower panel). The interaction terms of $N_g$ with ploidy

962 and with expression are represented for 0-dg $\theta_W$. **c**, Recessive load in tetraploid individuals estimated

963 as number of homozygous 0-dg derived alleles. **d**, DFE by ploidy and binned by strength of purifying

964 selection. **e**, Proportion of adaptive substitution ($\alpha$) and proportion of adaptive substitution relative to

965 neutral ($\omega_\alpha$) by ploidy. Errors bars represent 95% confidence interval based on 200 bootstrap replicates.

966

967

968 **Fig. 3 | Ploidy effects on linkage disequilibrium and the strength of linked selection. a**, Decay of

969 genotypic correlations (proxy for LD) within each population and averaged for each ploidy (heavy

970 lines) as a function of distance between sites. **b**, Curvilinear relationship between excess

971 nonsynonymous (0-dg) divergence ($E_{NS}$) on neutral diversity (4-dg $\theta_\pi$,). **c-d**, Linear relationship

972 between excess 0-dg divergence on neutral diversity (4-dg $\theta_\pi$) for gene-poor (<20$^{th}$ gene-density

973 percentile) and gene-dense regions (>90$^{th}$ gene-density percentile), respectively.

974

975 **Fig. 4 | Evidence for single origin of tetraploids. a-b**, Single origin of the *S. Carp.-4x*, and *Ruderal-*

976 *4x* tetraploids, respectively, followed by local admixture from their geographically proximal diploids

977 inferred as most likely scenario (large) by fastsimcoal2 coalescent simulations vs. competing scenarios

978 (small schemes). Range of median maximum-likelihood estimates of divergence times in generations

979 across different population quartets are indicated. **c**, Allele frequencies in each tetraploid lineage of

980 alleles diagnostic to particular diploid *A. arenosa* lineages. Significant differences within each category

981 of diploid alleles, as identified by Tukey's honestly significant difference (HSD) post hoc test, are

982 designated by distinct letters. **d**, Topology weights (TM, tetraploid monophyly; LA, local-admixture;

983 ILS, incomplete lineage sorting) in set of 6 meiosis-related genes compared with genome-wide average

984 (WG).

**Fig. 5 | Signals of interploidy introgression and loci resisting the gene flow.** Topology weightings for the three diagnostic topologies relating *S. Carp.-2x, S. Carp.-4x*, *W. Carp-4x* and the outgroup, *Dinaric-2x* across arms of scaffolds 8 (left) and 4 (right). Zoomed-in panels from top to bottom: topology weighting, average divergence ($\rho$) of *S. Carp.-4x* to all other tetraploids (black line) and to diploid lineages (colored lines), and Fay and Wu's H. Left zoom-in: example of locus locally introgressed from diploids and under positive selection (dominant LA topology, low divergence with local diploids specifically, and deeply negative Fay and Wu's H). Right zoom-in: resistance to local introgression of key meiotic locus, *ASY3*, with narrower peaks consistently with more ancient origin of the region.

**Table 1** Measures of within-population diversity and among-population divergence in diploid and tetraploid *A. arenosa*

996

| | Divergence | | | Diversity[1] | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | Rho / $F_{st}$[2] | AMO VA[3] | rM[4] | pairwise diversity ($\theta_\pi$) | | | Watterson's θ ($\theta_w$) | | | Tajima's D | | | $\pi_{NS}/\pi_S$ | $\theta_{NS}/\theta_S$ |
| Sites | 4-dg | 4-dg | 4-dg | all | 4-dg | NS (0-dg) | all | 4-dg | NS (0-dg) | all | 4-dg | NS (0-dg) | - | - |
| Diploids (14 pops) | 0.30 / 0.29 | 71 | 0.14 n.s. | 0.016 (0.003) | 0.022 (0.003) | 0.0054 (0.0007) | 0.015 (0.004) | 0.022 (0.003) | 0.005 (0.0009) | 0.03 (0.21) | 0.16 (0.18) | -0.09 (0.23) | 0.242 (0.017) | 0.255 (0.017) |
| Tetraploids (22 pops) | 0.20 / 0.11 | 48 | 0.55 *** | 0.015 (0.004) | 0.023 (0.006) | 0.0055 (0.0013) | 0.016 (0.004) | 0.023 (0.005) | 0.006 (0.0013) | -0.23 (0.29) | 0.00 (0.27) | -0.41 (0.28) | 0.237 (0.007) | 0.263 (0.007) |
| Difference[5] | - | - | - | n.s. | n.s. | n.s. | n.s. | n.s. | . | ** | . | *** | n.s. | *** |

997

998 Populations with < 5 individuals were excluded; for populations with > 5 individuals, sites were randomly downsampled to five to facilitate comparison across

999 populations.

1000 [1] values averaged across populations within the ploidy, standard deviation is in parentheses

1001 [2] values averaged over pairwise comparisons of populations belonging to that ploidy

1002 [3] % of explained variance among populations (compared to variance within populations) in Analysis of Molecular Variance (AMOVA)

1003 [4] Isolation by distance tested by Mantel test; the rM for diploid populations became 0.23* when spatially distant but genetically proximal Baltic populations were

1004 excluded

1005 [5] Wilcoxon rank-sum test; n.s. non significant, $p \le 0.07$ * $p \le 0.05$, ** $p \le 0.01$, *** $p \le 0.001$
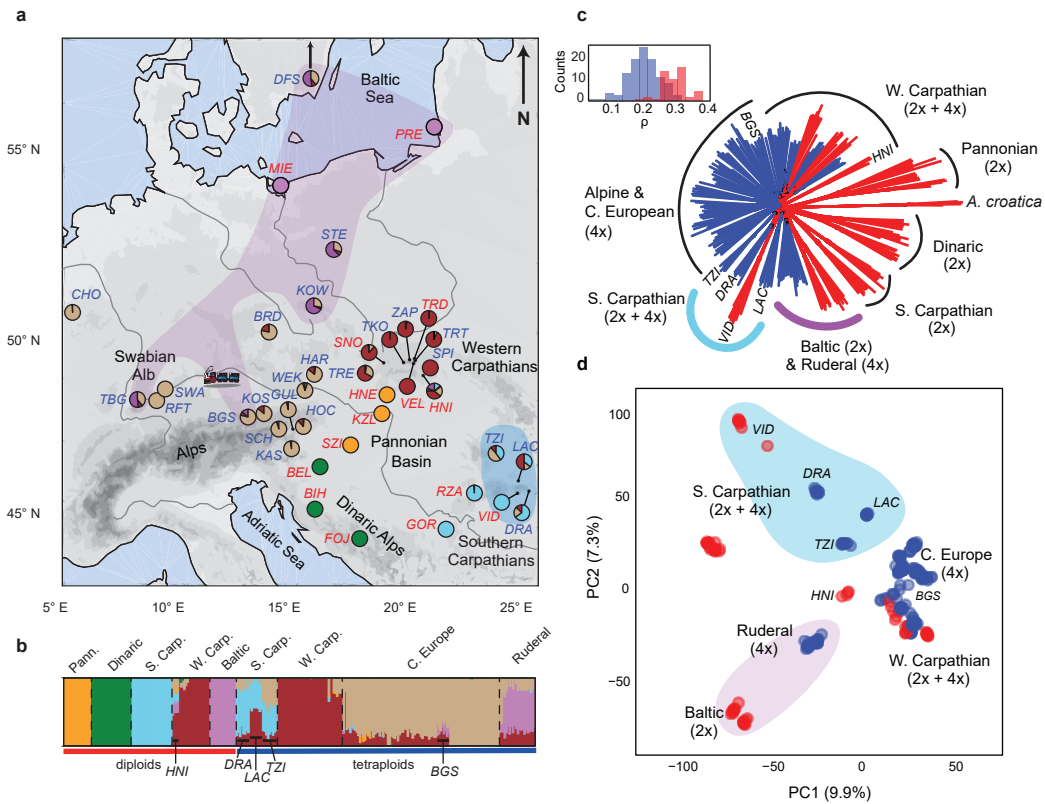
1006

**Fig. 1 | Geographic distribution and range-wide genetic variation of *Arabidopsis arenosa*. a**, Distribution of the 39 *A. arenosa* populations (red labels - diploids, blue - tetraploids) with average proportions of cluster membership inferred by FastStructure (panel b). Color shades highlight highly admixed tetraploid populations (*Ruderal* and *S. Carpathian-4x*) together with the diploid sources of admixture. **b**, Posterior probabilities of cluster membership of the 287 *A. arenosa* individuals as inferred by FastStructure under K=6. **c**, Neighbor-joining tree of Nei's genetic distances between all individuals and the outgroup *Arabidopsis croatica*. Individuals from admixed populations are highlighted correspondingly. Inset: distribution of pairwise genetic divergence of populations (ρ) within each ploidy. **d**, Principal component analysis of all but the two most divergent diploid (*Pannonian* and *Dinaric*) *A. arenosa* lineages (shades correspond to admixed populations).
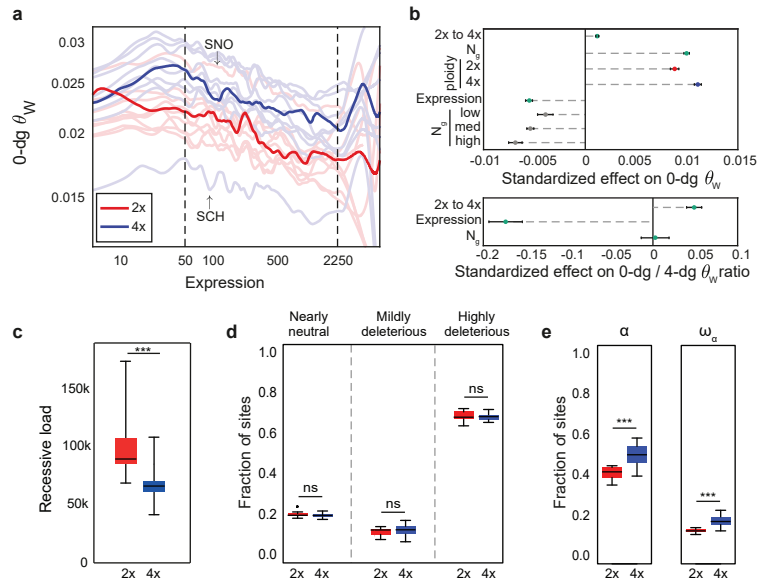
**Fig. 2. | Effects of ploidy on purifying selection, genetic load, and the distribution of fitness effects (DFE). a**, Genic nonsynonymous (0-dg) diversity versus average gene expression (log-scale) for each population and each ploidy (resp. faint and bold LOWESS curves). The two outlier populations, 2x-SNO and 4x-SCH, are indicated. **b**, Standardized effects with confidence intervals in multiple linear model of haploid effective population size ($N_g$), ploidy, and levels of expression on nonsynonymous (0-dg) $\theta_W$ (upper panel) and on 0-dg/4-dg $\theta_W$ ratio (lower panel). The interaction terms of $N_g$ with ploidy and with expression are represented for 0-dg $\theta_W$. **c**, Recessive load in tetraploid individuals estimated as number of homozygous 0-dg derived alleles. **d**, DFE by ploidy and binned by strength of purifying selection. **e**, Proportion of adaptive substitution ($\alpha$) and proportion of adaptive substitution relative to neutral ($\omega_\alpha$) by ploidy. Errors bars represent 95% confidence interval based on 200 bootstrap replicates.
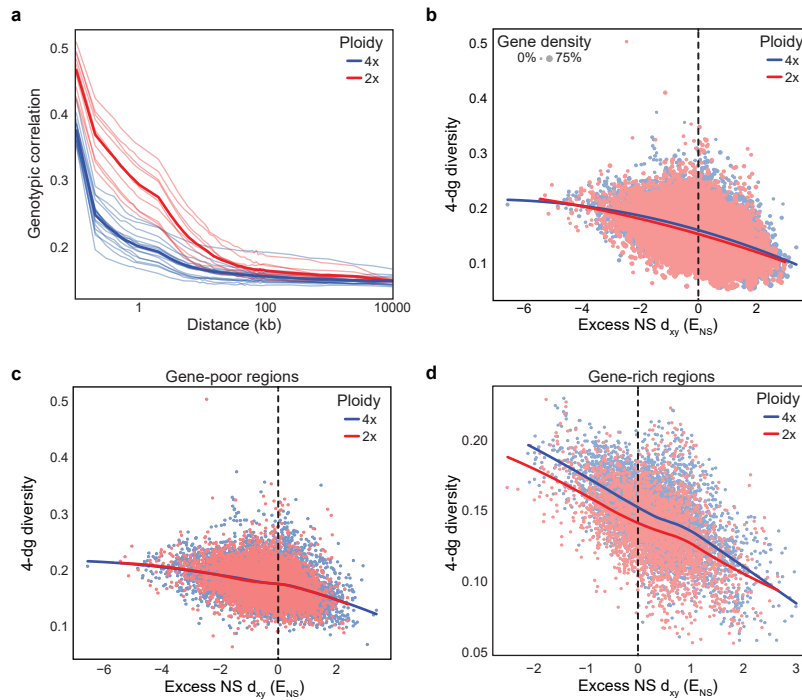
**Fig. 3 | Ploidy effects on linkage disequilibrium and the strength of linked selection. a,** Decay of genotypic correlations (LD estimator) within each population and averaged for each ploidy (heavy lines) as a function of distance between sites. **b,** Curvilinear relationship between excess nonsynonymous (0-dg) divergence ($E_{NS}$) on neutral diversity (4-dg $\theta_\pi$). **c-d,** Linear relationship between excess 0-dg divergence on neutral diversity (4-dg $\theta_\pi$) for gene-poor (<20th gene-density percentile) and gene-dense regions (>90th gene-density percentile), respectively.
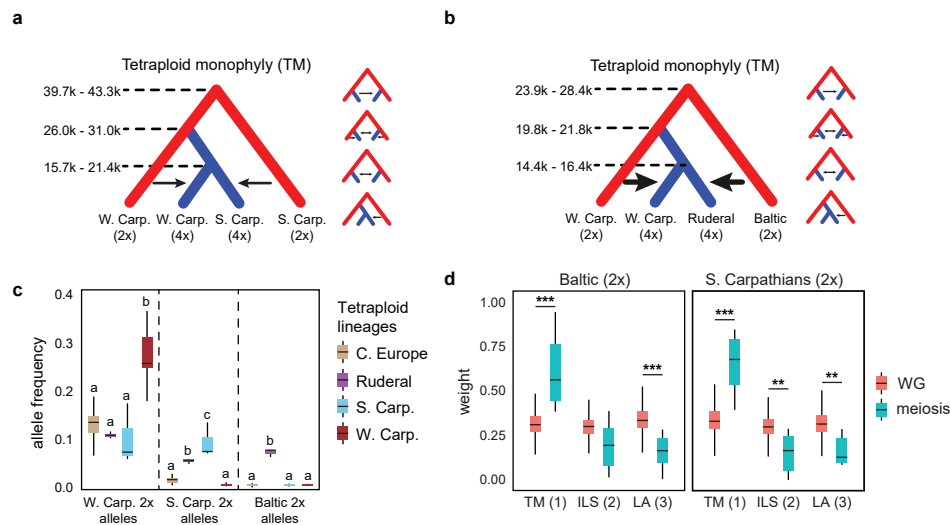
**Fig. 4 | Evidence for single origin of tetraploids. a-b** , Single origin of the *S. Carp.-4x*, and *Ruderal-4x* tetraploids, respectively, followed by local admixture from their geographically proximal diploids inferred as most likely scenario (large) by fastsimcoal2 coalescent simulations vs. competing scenarios (small schemes). Range of median maximum-likelihood estimates of divergence times in generations across different population quartets are indicated. **c**, Allele frequencies in each tetraploid lineage of alleles diagnostic to particular diploid *A. arenosa* lineages. Significant differences within each category of diploid alleles, as identified by Tukey's honestly significant difference (HSD) post hoc test, are designated by distinct letters.**d**, Topology weights (TM, tetraploid monophyly; LA, local-admixture; ILS, incomplete lineage sorting) in set of 6 meiosis-related genes compared with genome-wide average (WG).
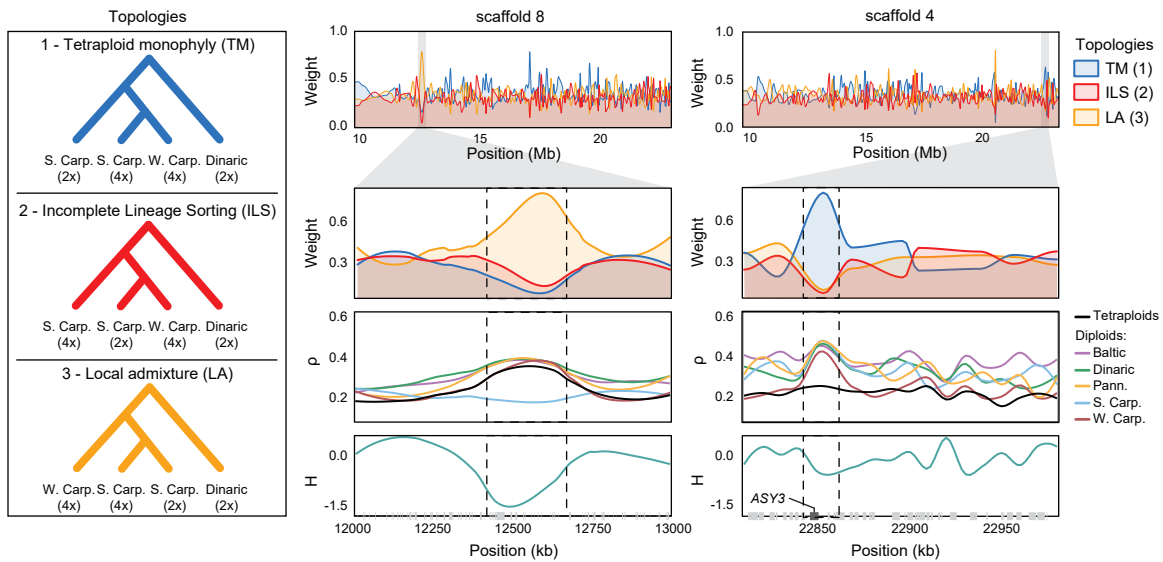
**Fig. 5 | Signals of interploidy introgression and barrier loci.** Topology weightings for the three diagnostic topologies relating *S. Carp.-2x, S. Carp.-4x, W. Carp-4x* and the outgroup, *Dinaric-2x* across arms of scaffolds 8 (left) and 4 (right). Zoomed-in panels from top to bottom: topology weighting, average divergence (ρ) of *S. Carp.-4x* to all other tetraploids (black line) and to diploid lineages (colored lines), and Fay and Wu's H. Left zoom-in: example of locus locally introgressed from diploids and under positive selection (dominant LA topology, low divergence with local diploids specifically, and deeply negative Fay and Wu's H). Right zoom-in: resistance to local introgression of key meiotic locus, *ASY3*, with narrower peaks consistently with more ancient origin of the region