

1 **Patterns of polymorphism, selection and linkage disequilibrium in the subgenomes of the**
2 **allopolyploid *Arabidopsis kamchatica***

3

4

5 Timothy Paape^{1*}, Roman V. Briskine^{1,2}, Heidi E.L Lischer^{1,3}, Gwyneth Halstead-Nussloch¹, Rie Shimizu-
6 Inatsugi¹, Masaomi Hatakayama^{1,2,3}, Kenta Tanaka⁴, Tomoaki Nishiyama⁵, Renat Sabirov⁶, Jun Sese⁷,
7 Kentaro K. Shimizu^{1,8*}

8

9

10

11 ¹ Department of Evolutionary Biology and Environmental Studies and Department of Plant and Microbial
12 Biology, University of Zurich, Winterthurerstrasse 190, CH-8057 Zurich, Switzerland

13 ² Functional Genomics Center Zurich, Winterthurerstrasse 190, CH-8057 Zurich, Switzerland

14 ³ Swiss Institute of Bioinformatics (SIB), Lausanne, Switzerland

15 ⁴ Sugadaira Montane Research Center, University of Tsukuba, Ueda, Nagano, Japan

16 ⁵ Advanced Science Research Center, Kanazawa University, 13-1 Takara-machi, Kanazawa, 920-0934, Japan

17 ⁶ Institute of Marine Geology and Geophysics, Far East Branch, Russian Academy of Sciences, Yuzhno-
18 Sakhalinsk, Russia

19 ⁷ Artificial Intelligence Research Center, National Institute of Advanced Industrial Science and Technology
20 (AIST), Tokyo, Japan

21 ⁸ Kihara Institute for Biological Research, Yokohama City University, 641-12 Maioka, 244-0813 Totsuka-
22 ward, Yokohama, Japan.

23 *authors for correspondence

24

25

26

27

28

29

30

31

32

33

34

35

36

37 **Abstract**

38 Although genome duplication is widespread in wild and crop plants, little is known about
39 genome-wide selection due to the complexity of polyploid genomes. In allopolyploid species,
40 the patterns of purifying selection and adaptive substitutions would be affected by masking
41 owing to duplicated genes or 'homeologs' as well as by effective population size. We
42 resequenced 25 distribution-wide accessions of the allotetraploid *Arabidopsis kamchatica*,
43 which has a relatively small genome size (450 Mb) derived from the diploid species *A. halleri* and
44 *A. lyrata*. The level of nucleotide polymorphism and linkage disequilibrium decay were
45 comparable to *A. thaliana*, indicating the feasibility of association studies. A reduction in
46 purifying selection compared with parental species was observed. Interestingly, the proportion
47 of adaptive substitutions (α) was significantly positive in contrast to the majority of plant
48 species. A recurrent pattern observed in both frequency and divergence-based neutrality tests is
49 that the genome-wide distributions of both subgenomes were similar, but the correlation
50 between homeologous pairs was low. This may increase the opportunity of different
51 evolutionary trajectories such as in the *HMA4* gene involved in heavy metal hyperaccumulation.

52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69

70 Introduction

71 Genome duplication is a widespread evolutionary force in plants. As many as 35% of vascular
72 plants are recent polyploid species¹ and increased ploidy is particularly common in crops². The
73 abundance of polyploid species in plants motivated speculation and theoretical analysis on the
74 advantages and disadvantages of genome duplication^{3,4}. However, compared with diploid
75 species, much less is known about the genome-wide patterns of polymorphism and selection
76 due to the complexity of polyploid genomes⁵. Major difficulties in genome scale analyses result
77 from the large genome sizes of polyploids and the high similarity between the duplicated
78 chromosomes. However, recent advances in next-generation sequencing and bioinformatic
79 tools^{6,7} are enabling genome-wide data to study polymorphisms and transcriptomics patterns
80 for entire subgenomes in newly emerging model polyploids⁸⁻¹¹.

81 Genome-wide strengths of positive and purifying selection can be quantified using
82 several complementary approaches. Frequency-based tests using site-frequency spectra (SFS)
83 such as Tajima's D and Fay and Wu's H statistics can detect rare or common polymorphisms that
84 are due to purifying or positive selection. Divergence-based tests compare interspecific
85 divergence (from an outgroup) to intraspecific polymorphism to identify positive selection on
86 amino-acid substitutions¹². These tests include several derivatives of the McDonald-Kreitman
87 test¹³ or "MK-tests", such as the direction of selection (DoS) neutrality index¹⁴, and methods to
88 estimate the distribution of fitness effects (DFE) and proportion of adaptive substitutions (α)¹³ in
89 genome-wide data. Theoretical and empirical studies in plant species using these methods^{15,16}
90 showed that the strengths of selection are affected by species-specific characteristics such as
91 the effective population size (N_e), mating system, and genome duplication, which are mutually
92 interacting. In particular, species with low N_e typically have the highest proportions of neutral
93 mutations^{15,17}, while species with large N_e have higher proportions of non-synonymous
94 substitutions under purifying selection and adaptive evolution^{8,19,20}.

95 Allopolyploidization should have a profound effect on the patterns of polymorphism and
96 selection. First, the redundancy of duplicated gene copies of similar function from different
97 parents ("homeologs") may affect the strength of selection. At the early stages, genome
98 duplication may increase evolutionary rates of duplicated genes^{21,22} and may facilitate the
99 evolution of new adaptive function because the original function can be retained in other copies
100 (so-called neofunctionalization model)^{23,24}. In contrast, the additional copy may mask the effect
101 of adaptive and deleterious mutations^{4,16}. Second, polyploidization must involve a reduction in

102 N_e due to a bottleneck during speciation. In addition, polyploid speciation is typically associated
103 with the transition from outcrossing to self-fertilization, which reduces N_e several times less
104 than parental species (at least half)²⁵. While studies of selection in polyploids are very limited, a
105 recent empirical study in the allotetraploid *Capsella bursa-pastoris* showed a decrease in the
106 efficacy of purifying selection in one of the subgenomes but an increase in another subgenome⁸.
107 Further empirical studies are necessary to compare the consequences of genome duplication in
108 polyploid species.

109 The genus *Arabidopsis* has both auto- and allopolyploid species in addition to the more
110 well-studied diploid relatives²⁶. *Arabidopsis kamchatica*²⁷ is a recent allopolyploid (estimated
111 20,000-250,000 years ago)²⁸ derived from the two diploid species *A. halleri* (particularly subsp.
112 *gemmifera* distributed in East Asia), and *A. lyrata* (particularly subsp. *petraea* from Far East
113 Russia)²⁹⁻³¹. The two diploid parents are predominantly self-incompatible (SI) while a transition
114 to selfing accompanied allopolyploid formation²⁸. The genome size (about 450 Mb) is relatively
115 small among polyploid species^{32,33} which is an advantage for resequencing. The species
116 distribution of *A. kamchatica* is very broad, ranging from Taiwan, Japan, Far East Russia, Alaska
117 and Pacific Northwest, USA. The high variation in latitude and altitude compared with the
118 parental species^{34,35} suggests that merging the diploid transcriptional networks and parental
119 adaptations provided the allopolyploid with plasticity to inhabit diverse environments¹⁰.

120 To understand the ecological distributions of polyploids, genetically tractable traits are
121 essential. Heavy metal tolerance and hyperaccumulation likely influenced ecological divergence
122 and speciation between the parental species of *A. kamchatica* (*A. halleri* and *A. lyrata*) due to
123 adaptive mutations in metal transporter genes such as *HEAVY METAL ATPASE4 (HMA4)*^{36,37}. The
124 *HMA4* locus has been shown to be the primary transporter of cadmium and zinc from roots to
125 shoots in *A. halleri* due to a tandem triplication and enhanced *cis*-regulation, while only a single
126 copy of *HMA4* exists in the non-hyperaccumulators *A. lyrata* and *A. thaliana*. *A. kamchatica*
127 inherited hyperaccumulation from the diploid parent *A. halleri*, although attenuated expression
128 of *halleri*-derived *HMA4* and putatively inhibiting *lyrata*-derived factors reduced the trait to
129 about half of *A. halleri*¹⁰. Estimates of genetic diversity surrounding the *HMA4* region in *A. halleri*
130 suggests a hard selective sweep³⁸ which may have predated the formation *A. kamchatica*¹⁰.

131 Here, we used *de novo* assemblies of the closest diploid relatives of *A. kamchatica* to
132 sort Illumina reads to their respective subgenomes using a distribution-wide collection of 25
133 natural allopolyploid accessions. We used population genomics to ask: a) what is the level of

134 genome-wide diversity compared with diploid outcrossing and selfing *Arabidopsis* species? b)
135 are there differences in polymorphism, allele frequencies, linkage disequilibrium (LD), and
136 selection between subgenomes? c) do pairs of homeologs tend to show similar patterns in
137 diversity and neutrality? d) does the *HMA4* locus show significant differences in genetic diversity
138 between homeologs and how does diversity surrounding this locus compare with the genome-
139 wide average? e) what proportions of the subgenomes show neutral, deleterious, or adaptive
140 mutations and how do they differ from the diploid parents? and f) are there high frequencies of
141 loss of function mutations in either subgenome? Together, these plant accessions and
142 polymorphism data will serve as a core diversity panel for further studies of genotype-
143 phenotype associations and the genetic architecture of complex traits using larger collections of
144 globally collected samples.

145

146 **Results**

147 Reference Genome Assembly and Allopolyploid Resequencing

148 To sort Illumina reads of *A. kamchatica* to their parentally-derived subgenomes, we generated
149 long mate-pair *de novo* assemblies of *A. lyrata* subsp. *petraea* (also called *A. petraea* subsp.
150 *umbrosa*) in addition to East Asian *A. halleri* subsp. *gemmifera* which we previously reported³⁹.
151 Assembly statistics indicated that the *A. lyrata* and *A. halleri* reference genomes have scaffold
152 N50 of 1.2 Mb and 0.7 Mb, comprising 1,675 and 2,239 scaffolds respectively (Table 1,
153 Supplementary Table 1 and 2 for gene annotation statistics), providing opportunities to compare
154 diversity over very large syntenic regions in the allopolyploid subgenomes.

155 We sorted reads of 25 individuals from a distribution-wide collection (Supplementary
156 Table 3) of *A. kamchatica* to their parental origins by first aligning each read to both parental
157 genomes then classified the reads as ‘*origin*’ reads (*halleri*-derived = H-origin, *lyrata*-derived = L-
158 origin) using algorithms that quantify mismatches to either parent³². Our accessions had on
159 average 12.5X coverage for the H-origin-subgenome (range 5.2X - 20.7X) and on average 10.7X
160 coverage for the L-origin-subgenome (range 4.3X - 17.7X). Homeolog specific PCR and Sanger
161 sequencing was used to validate SNPs and read sorting for twelve genes and showed that reads
162 were accurately assigned to their respective subgenomes (Supplementary Material). In addition,
163 pyrosequencing was used in two previous studies to detect ratios of parentally derived SNPs to
164 validate homeolog specific expression (RNA-seq) in ten other genes^{10,32} where the same read
165 sorting pipeline was used. After filtering for SNP quality and coverage, our resequencing dataset

166 resulted in 1,674,191 H-origin and 1,930,341 L-origin SNPs. Using the parental genome
167 assemblies for *A. kamchatica* SNP calling, we identified ca. 23,500 homeologous coding
168 sequences using reciprocal best BLAST hits (Supplementary Table 2), of which ca. 21,500 show
169 orthology to *A. thaliana*, representing 72% and 67% of our annotated genes respectively.

170

171 Genome-wide Nucleotide Diversity in *A. kamchatica*

172 We examined the patterns of nucleotide diversity for ca. 21,000 coding sequences of both
173 *halleri* and *lyrata*-derived homeologs in *A. kamchatica* that could be aligned to *A. thaliana*
174 orthologs as the outgroup. We found that both subgenomes showed similar mean values of
175 nucleotide diversity (π) ($\pi_{\text{coding}} = 0.0014 \text{ bp}^{-1}$ for *halleri*-subgenome and $\pi_{\text{coding}} = 0.0015 \text{ bp}^{-1}$ for
176 *lyrata*-subgenome, and $\pi_{\text{coding}} = 0.0015 \text{ bp}^{-1}$ when combined) although the *lyrata*-derived
177 homeologs showed slightly broader ranges in π (Table 2, Fig. 1A). Nucleotide diversity at
178 synonymous sites (π_{syn}) was also similar for the two subgenomes with a slightly higher value for
179 the *lyrata*-subgenome ($\pi_{\text{syn}} = 0.0049$) than in the *halleri*-subgenome ($\pi_{\text{syn}} = 0.0044$). The
180 nucleotide diversity in *A. kamchatica* is about six times lower than European *A. halleri* and *A.*
181 *lyrata* ($\pi_{\text{syn}} = 0.029$ for *A. halleri* and 0.028 for *A. lyrata* estimated using resequencing data
182 from³⁰) and is more similar to that of *A. thaliana* ($\pi_{\text{syn}} = 0.0059 - 0.007$)^{17,30,40}. Sliding window
183 analysis including non-coding regions also showed comparable values (Supplementary Table 4).

184 We calculated the effective population size, N_e , using our empirical estimates of π for *A.*
185 *kamchatica* and both diploid species and two different mutation rates^{41,42}. The estimated values
186 for *A. kamchatica* were: $N_e = 77,000$ and 54,000 using the two mutation rates respectively. The
187 values for *A. kamchatica* were several times lower than *A. halleri*: $N_e = 467,000$ and 364,000 and
188 *A. lyrata*: $N_e = 483,000$ and 345,000 (Supplementary Table 5). We interpret these estimates of N_e
189 with caution as the mutation rates for these species have not been estimated directly and the
190 diversity estimates used in the calculation can themselves be affected by demography. The
191 estimates are nevertheless useful as general comparisons between species to identify large
192 differences in magnitude^{17,19}.

193 Higher proportions of non-synonymous mutations were found to be at low frequency
194 compared with synonymous mutations and no significant differences in the relative proportions
195 were found between subgenomes (Fig. 1B). This suggests purifying selection on a large
196 proportion of amino-acid changing substitutions in both subgenomes. Frequency-based test
197 statistics clearly show significant departures from neutrality for both subgenomes (Fig. 1C). The

198 mean values of Tajima's D were negative for both subgenomes (Table 2, Fig. 1C) owing to high
199 proportions of rare variants.

200 The distributions and means of Tajima's D in *A. kamchatica* (Table 2) are similar to early
201 genome-wide data from *A. thaliana* (mean Tajima's $D_{A.thaliana} = -0.8$)⁴³, although more recent
202 estimates using over 300 genomes show a higher mean but not higher median in *A. thaliana*
203 (mean $D_{A.thaliana} = 0.006$, median $D_{A.thaliana} = -0.33$)³⁰, which likely reflects more intermediate-
204 frequency polymorphisms in the large species wide sample. The same study³⁰ reported an
205 excess of rare variants in the diploid relatives of *A. kamchatica* (mean $D_{A.lyrata}$ to be -0.99 in *A.*
206 *lyrata* and $D_{A.halleri} = -0.23$ in *A. halleri*).

207 We found the means of the distributions for most summary statistics to be very similar
208 between the two subgenomes, but when pairs of all homeologs were compared correlations
209 were generally low for diversity and neutrality estimators (Table 2). The correlations of π_{syn} and
210 $\theta_{w_{syn}}$ were both nearly zero (Table 2). Similarly, the distributions and means of Tajima's D
211 overlap for both subgenomes but the correlation for Tajima's D between pairs of homeologs is
212 very low ($R^2 = 0.03$). The Fay and Wu's H statistic, which detects departures from neutrality due
213 to intermediate and high frequency variants, also shows a very low correlation between
214 homeologs (Table 2). Higher correlations were observed for non-synonymous or total sites, but
215 this can be explained by the constraints on non-synonymous changes. In summary, the low
216 correlations are consistent with different evolutionary trajectories of individual homeologous
217 pairs.

218

219 Mean Rate of LD Decay in Both Subgenomes is Similar But Not Equal

220 Long scaffold assemblies allowed us to estimate genome-wide LD for each subgenome to
221 evaluate the feasibility of association mapping in *A. kamchatica*. We found that mean LD decay
222 was between 5-10 kb for both subgenomes (Fig. 1D), which is similar to the self-fertilizing
223 species *A. thaliana* and *M. truncatula* which show LD decay within 2-10 kb ranges^{44,45}. The mean
224 LD for the *lyrata*-subgenome decayed slightly faster and remained at $r^2 = 0.47$ over the scale of >
225 100 kb genomic regions while mean LD for the *halleri*-subgenome leveled off at $r^2 = 0.34$ > 100
226 kb. The 50% and 90% confidence intervals around the mean LD decay also revealed much
227 greater variance in the *lyrata*-subgenome (Supplementary Fig. 1).

228 Population structure assignments and phylogenetic clustering may provide some
229 explanation for subgenome differences in LD. The 25 accessions cluster geographically with one

230 main clade/group comprising the northern accessions (Russia, Sakhalin, and Alaska) and the
231 other main group containing Japanese accessions (Supplementary Fig. 2,3). The branch lengths
232 within these groups for the *lyrata*-subgenome are shorter than for the *halleri*-subgenome,
233 particularly in the Japanese clade, indicating greater relatedness. These clusterings are also
234 consistent with previous haplotype analysis using low density nuclear and chloroplast markers²⁹.

235

236 Diversity of the *HMA4* Locus and the Genomic Background

237 We analyzed genetic diversity on the scaffolds containing the *HMA4* locus to compare whether
238 it differs from the genomic background and the surrounding regions flanking the *HMA4* coding
239 sequences. We centered the main genomic region containing the *HMA4* coding sequences
240 which we call “HMA4-M” (containing 17 coding sequences). This region spans 304 kb on *A.*
241 *halleri* (scaffold_116) and spans 155 kb on *A. lyrata* (scaffold_52). While the differences in length
242 of HMA4-M between the parental genomes can be attributed to the triplicated *HMA4* genes in
243 *A. halleri*, the genes surrounding *HMA4* in both *A. halleri* and *A. lyrata* are syntenic (Fig. 2A). To
244 compare HMA4-M to surrounding regions, we used the upstream adjacent region (left-side)
245 “HMA4-L” (containing 13 coding sequences) which is 125 kb for the *A. halleri* region and 183 kb
246 in *A. lyrata*, and the downstream adjacent region (right-side) “HMA4-R” (containing 13 CDS
247 sequences), which is 105 kb in the *A. halleri* region and ca. 50 kb for *A. lyrata*.

248 The distribution of π in the HMA4-M region for H-origin genes showed low diversity
249 ($\pi_{\text{mean}} = 0.0007$) but it is not significantly lower than the background genes (Fig. 2B and 2C).
250 However, the two adjacent regions (HMA4-L and HMA4-R) compared to the HMA4-M
251 (containing the *HMA4* coding sequences) region have significantly greater diversity (Fig. 2B and
252 2C). Furthermore, we found significantly lower Tajima’s D, Fu & Li’s D* and Fu & Li’s F* statistics
253 in the HMA4-M region compared with both adjacent regions (Fig. 2E), suggesting greater
254 selection on the HMA4-M region. The significantly lower diversity and neutrality statistics in
255 HMA4-M compared with the adjacent regions likely defines the window of the sweep region
256 previously reported for *A. halleri*³⁸.

257 Unlike the *halleri* HMA4-M region, the diversity of the *lyrata* HMA4-M region is
258 significantly greater than the genomic background (p-value = 0.0028), but not different from the
259 two adjacent regions (Fig. 2D). Moreover, the *lyrata*-HMA4-M region shows no significant
260 differences from the adjacent HMA4-L or HMA4-R regions for Tajima’s D, Fu & Li’s D* and Fu &
261 Li’s F* (not shown). The elevated diversity of the *lyrata*-origin *HMA4* locus compared with the

262 genomic background is consistent with relaxed selective constraint on the *lyrata*-origin *HMA4*
263 locus.

264 We also estimated diversity of all annotated heavy metal transporters, metal ion
265 transporters, and metal homeostasis genes for comparison with the genome-wide average (HM
266 genes, N=118 genes). We expected these genes to have low overall diversity in both genomes
267 due to selective constraint as many of these ion transporters are expected to have roles in basic
268 metal homeostasis⁴⁶. As a contrast, we compared NBS-LRR genes (N=39 genes) which have
269 putative roles in plant defense and have high diversity in plants^{47,48} and are expected to have
270 equally high diversity in both subgenomes. The HMA4-L and HMA4-R regions in both
271 subgenomes have more similar levels of diversity to NBS-LRR's than to the genomic background
272 or HM genes (Fig. 2C and 2D).

273

274 The Majority of Homeologous Proteins Showed Signatures of Purifying Selection

275 Next we employed divergence-based tests to estimate the strength of purifying and positive
276 selection on amino-acid changing substitutions. We calculated the divergence of each homeolog
277 from the outgroup *A. thaliana* to estimate the relative proportions of diverged non-synonymous
278 (D_n) and synonymous (D_s) sites to polymorphic non-synonymous (P_n) and synonymous (P_s) sites.
279 For each gene, The counts of D_n , D_s , P_n , and P_s for the coding regions of both subgenomes were
280 used to estimate the direction of selection (DoS)¹⁴, a neutrality index that varies from -1.0 to 1.0,
281 where zero indicates neutrality and negative and positive values indicate purifying and positive
282 selection, respectively. Both subgenomes had similar distributions in DoS with means of -0.2
283 (Fig. 3A) suggesting that 68-71% of proteins derived from both subgenomes are under purifying
284 selection (when DoS is < -0.01). Like the previous summary statistics, the correlation in DoS
285 between *halleri* and *lyrata* homologs is positive but fairly low ($R^2 = 0.17$).

286 MK-tests were conducted to detect homeologs showing purifying selection or adaptive
287 evolution on amino-acid changing mutations. Among the significant MK-test genes, a total of
288 3018 H-origin and 3804 L-origin homeologs showed DoS < 0 ($D_n/D_s < P_n/P_s$). This is consistent
289 with purifying selection rather than positive selection for these genes. While the homeologs
290 with significant MK-test comprise a substantial portion in our dataset, only 19% of them include
291 both homeologs (i.e., there is significance for one homeolog but not the other for 81% of
292 significant homeologous pairs, Fig. 3B). For example, the H-origin homeolog of the resistance

293 gene *RPM1* (orthologous to *A. thaliana* gene: AT3G07040) was significant for the MK-test (DoS <
294 0) but the L-origin copy was not.

295 For genes showing positive selection (or adaptive evolution) using MK-tests, 146 *halleri*-
296 origin and 212 *lyrata*-origin genes were significant when DoS > 0.01 (Fig. 3C, D). For these genes,
297 when the *halleri*-derived homeologs shows a positive DoS, the *lyrata*-derived homeolog shows a
298 more neutral or negative distribution in DoS and vice versa. Among these is the H-origin *HMA4*
299 gene. These results, in addition to the low correlation in DoS between homeologous pairs and
300 small overlap among all significant MK-test genes (Fig. 3B), indicates that a substantial
301 proportion of homeologs have been shaped by different strengths of selection. These results are
302 also in agreement with low correlations in Tajima's D and Fay and Wu's H despite for pairs of
303 homeologs (Table 1), providing additional support that redundant genes exhibit significant
304 differences due to stronger positive or purifying selection on only one of the two copies.

305

306 The Distribution of Fitness Effects (DFE)

307 The tests above indicated that large numbers of homeologs show patterns consistent with
308 purifying selection on amino-acid changing mutations (see Fig. 3). We quantified the genome-
309 wide proportions of deleterious and effectively neutral mutations using the distribution of
310 fitness effects (DFE) method¹³ in the two *A. kamchatica* subgenomes and both diploid relatives.
311 In this method, the DFE is estimated from the site frequency spectra of non-synonymous and
312 synonymous polymorphisms while accounting for effects of demographic changes. Effectively
313 neutral mutations are represented by $0 < N_e s < 1$, mildly deleterious by $1 < N_e s < 10$, deleterious
314 by $10 < N_e s < 100$ and strongly deleterious by $N_e s > 100$ (where N_e is the effective population size
315 and s is the selection coefficient). The DFE estimates of the two *A. kamchatica* subgenomes
316 show similar distributions with about 70% of mutations in the deleterious to strongly deleterious
317 categories ($N_e s > 10$) and about 20% effectively neutral ($0 < N_e s < 1$) (Fig. 4A). The DFE of *A.*
318 *halleri* and *A. lyrata* showed lower proportions of neutral mutations (16% of mutations $0 < N_e s <$
319 1 in diploids, and 19% mutations $0 < N_e s < 1$ in both subgenomes) and greater proportions of
320 deleterious mutations ($N_e s > 100$) than either of the corresponding allopolyploid subgenomes.
321 While the differences are significant, the magnitude of the differences is not remarkable.

322 To examine whether subsets of either subgenome experience a reduction in purifying
323 selection, we classified homeologs according to gene expression level, which is one of the best
324 predictors of evolutionary rates (dN/dS) in most organisms⁴⁹. Expression level is negatively

325 correlated with dN/dS due to strong constraint on amino acid substitutions (dN)²² for highly
326 expressed genes, but this has not been shown in recent polyploid species. As a test of selective
327 constraint on highly expressed genes, we found dN/dS was negatively correlated with
328 expression for both homeologs (Fig. 4B). We would therefore expect genes that are highly
329 expressed to show the strongest purifying selection, and low expressed genes to show relaxed
330 constraint. We estimated the DFE again to quantify purifying selection and relaxed constraint
331 using the distribution of expression levels in leaf and root tissues of *A. kamchatica* to categorize
332 homeologs as high (genes in upper 10% RPKM) or low expression (lower 10% of RPKM). The
333 majority (62%) of the highly expressed genes include both homeologs (Fig. 4C). The DFE patterns
334 indicated that low expressed genes have the highest proportion of neutral mutations (relaxed
335 constraint) and the lowest proportion of deleterious mutations compared with the genome-
336 wide data, while highly expressed genes showed the opposite pattern (Fig. 4D). These results
337 indicate that the DFE method can detect relaxed constraint and strong purifying selection as
338 predicted when gene expression levels are accounted for.

339

340 The Proportion of Adaptive Substitutions in Diploids and Allopolyploid Subgenomes

341 The proportion of adaptive substitutions (α) was estimated as the excess of between-species
342 divergence relative to polymorphism as expected from the estimated DFE¹³ to account for
343 slightly deleterious mutations. In contrast to the majority of the previously studied plant species
344 including *A. thaliana*, we found significantly positive values of α for the two diploid species and
345 both allopolyploid subgenomes. The diploid species *A. halleri* and *A. lyrata* showed the highest α
346 values (0.25 and 0.27 respectively) (Fig. 4E). We subsampled 18 *A. kamchatica* accessions to be
347 statistically comparable to the available *A. halleri* and *A. lyrata* samples (Supplementary Table
348 6). The α estimates for the H- and L-origin subgenomes of *A. kamchatica* were lower than those
349 of the corresponding diploid species but significantly greater than zero (0.12 and 0.09,
350 respectively) (Fig. 4E). The difference in α between subgenomes was significant but subtle (3%
351 difference using the samples above, 6% difference when all 25 *A. kamchatica* accessions were
352 used; Supplementary Fig. 4).

353

354 High Impact Mutations are at Low Frequency in Subgenomes

355 We identified genes having high impact mutations that are likely to be deleterious due to their
356 putative effects on amino acid sequences and gene expression into the following mutation

357 categories: frameshifts, loss of start codon, premature stop codons (stop gained), and loss of
358 stop codons (stop loss). For any gene, we counted every one of the mutation types regardless of
359 the number. While it is not possible to determine the order of disruptive mutations, multiple
360 frameshifts of premature stop codons in a gene would be expected to result in a loss of function.

361 Frameshifts and stop-gained categories comprised the majority of mutation types for
362 both subgenomes (Supplementary Table 7). Frequencies of each mutation type indicated that
363 most mutation types in any gene are found in only a single genotype in either subgenome (Fig.
364 5). Despite a higher number of mutations in the *lyrata*-homeologs, there were slightly greater
365 proportions at low frequencies in the *halleri*-homeologs. Out the total 4219 *halleri*-origin and
366 4952 *lyrata*-origin disrupted genes, only 511 genes (2.5%) showed large effect mutations in both
367 homeologs in the same accession suggesting that large effect mutations in both homeologs
368 were deleterious. The distribution of genes with high impact mutations in both homeologs
369 shows that most accessions have < 50 genes (orthologous to *A. thaliana*) that are disrupted with
370 putatively similar functions (Supplementary Fig. 5).

371 We conducted gene ontology (GO) analysis to determine whether there was enrichment
372 for GO terms using the two most common high-impact mutation types, i.e., frameshift
373 mutations and stop codons. For both subgenomes, hydrolase activity (GO:0016787) was the
374 most significant GO term for molecular function, followed by several GO categories for
375 nucleotide binding (Supplementary Table 8). Programmed cell death (GO:0012501) and
376 apoptosis (GO:0006915) were significant in the *halleri*-origin genes only. No significant gene
377 ontologies were found with ≥ 20 query genes for the list of genes that had high impact
378 mutations in both homeologs in a single accession.

379

380 Discussion

381 Similar Genome-wide Distributions in Both Subgenomes but Low Correlations Between 382 Homeologous Pairs

383 A recurrent pattern we observed on patterns of diversity and signatures of selection was that
384 the genome-wide distributions were similar between subgenomes, but the correlations between
385 the pairs of homeologs were low. We found this pattern in the polymorphism levels such as π_{syn}
386 and $\theta_{\text{w syn}}$, in frequency-based tests of neutrality (Tajima's D, Fay & Wu's H), and in divergence-
387 based tests (DoS). The similar genome-wide distributions are consistent with the fact that the
388 subgenomes shared the same history since the allopolyploidization event. The low correlation

389 suggests that at the gene level, genetic diversity of a large number of homeologs may have been
390 shaped by different levels of positive and purifying selection, as well as relaxed constraint. This
391 supports that homeologs may evolve as independent loci, which may not be surprising because
392 *A. kamchatica* shows disomic inheritance prohibiting recombination between homeologs²⁸.
393 These results also suggest that the difference between homeologs could contribute to the broad
394 environmental response of polyploids, which may be realized by combining different
395 adaptations of two parental species¹⁰ such as in the *HMA4* gene.

396

397 Nucleotide Diversity and Linkage Disequilibrium is Similar to *A. thaliana* Suggesting the 398 Feasibility of Genome-wide Association Studies of *A. kamchatica*

399 We found that the level of nucleotide diversity of *A. kamchatica* is moderate and similar to that
400 of the diploid self-compatible *A. thaliana*, and 6 times lower than the diploid outcrossing species
401 *A. halleri* and *A. lyrata*. It follows that the N_e of *A. kamchatica* is 6 times lower than the two
402 diploid species. The ancestor of the genus *Arabidopsis* must have been a self-incompatible
403 diploid species like present-day *A. lyrata* and *A. halleri*²⁵, indicating that similar reductions in
404 genetic diversity occurred in the lineages of *A. kamchatica* and *A. thaliana*.

405 The extent of LD decay in *A. kamchatica* is also comparable to *A. thaliana*⁴⁵ and appears
406 adequate for characterizing the genetic architecture of complex traits within relatively narrow
407 genomic windows using genome-wide association studies (GWAS). The selfing mating system,
408 levels of genetic diversity, LD, and a recently established transgenic technique⁵⁰ suggests that *A.*
409 *kamchatica* would be a suitable model for functional genomics of adaptive mutations in a
410 polyploid species.

411

412 The *HMA4* Locus Exhibits Significant Subgenome Differences in Genetic Diversity

413 The most important locus for zinc hyperaccumulation, *HMA4*³⁷, involved two types of gene
414 duplication in *A. kamchatica*: a tandem triplication in diploid *A. halleri*, followed by a whole
415 genome duplication event, which contributed an additional *HMA4* copy from the *A. lyrata*
416 parent. Despite multiple hybrid origins of *A. kamchatica*, the tandem triplication (three *halleri*-
417 derived *HMA4* copies) is fixed in the allopolyploid¹⁰ suggesting it was present in all founding *A.*
418 *halleri* parents. The high expression of H-origin *HMA4* in *A. kamchatica* explains high levels of
419 zinc accumulation. Expression of the L-origin *HMA4* copy is very low compared with the *halleri*

420 *HMA4* gene(s) so it is unlikely that the copy from *A. lyrata* contributes anything significant to
421 hyperaccumulation in *A. kamchatica*.

422 Long scaffolds containing the *HMA4* copies and surrounding genes allowed us to
423 compare homeologs across large genomic distances. The genetic diversity surrounding the
424 *halleri*-derived *HMA4* gene that spans ca. 300 kb (HMA4-M) is significantly lower than the
425 syntenic *lyrata*-derived region (ca. 100 kb) suggesting different evolutionary pressures or
426 trajectories of functional duplicates. The higher diversity of the *lyrata* HMA4-M region is
427 consistent with a pattern of relaxed constraint, while a selective sweep and genetic hitchhiking
428 characterizes the *halleri*-derived HMA4-M region. Because we can infer that the triplication was
429 ancestral and the reduced diversity at this locus and hitchhiking surrounding the *HMA4* genes
430 was most likely the result of strong selection in the *A. halleri* parent³⁸, diversity was probably
431 greatly reduced prior to the polyploidization events.

432

433 Purifying Selection in Polyploid Species

434 Theoretical studies suggested that higher proportions of neutral mutations (i.e., greater relaxed
435 constraint) can result from whole genome duplication due to the reduction of N_e or due to
436 masking of deleterious mutations by functionally redundant gene copies^{15,16}. This would be
437 evident by greater proportions of effectively neutral mutations ($0 < N_e s < 1$) in the polyploid
438 subgenomes compared with the diploid parents⁸. Similarly, greater proportions of deleterious
439 mutations ($N_e s > 10$) in the diploid species would be expected compared to their derived
440 polyploid subgenomes. We did detect significant differences between diploid parental species
441 and the corresponding subgenomes of *A. kamchatica* in the proportions of mutations in the
442 neutral (< 5% differences) and deleterious (5-7% differences) categories, although the
443 differences were not drastic (Fig. 4A).

444 Using a similar approach, the change in purifying selection was studied by comparing
445 the allopolyploid species *C. bursa-pastoris* with its diploid parents, *C. grandiflora* (outcrosser
446 with high N_e) and *C. orientalis* (selfing with low N_e)⁸. First, for the subgenome derived from the
447 outcrossing parent *C. grandiflora*, the proportion of neutral mutations doubled from ~17% to
448 ~35% neutral mutations. This demonstrated that the subgenome derived from an outcrossing
449 parent with a large N_e shows a high proportion of neutral mutations due to relaxed constraint.
450 Second, the opposite pattern was observed in the subgenome derived from the selfing parent
451 (decreased from ~40% to 35% neutral mutations). The DFE patterns in *C. bursa-pastoris* and *C.*

452 *orientalis* conforms to the trend in plants which shows species with low N_e usually have greater
453 proportions of neutral mutations¹⁵ consistent with greater strengths of purifying selection with
454 higher N_e . However, despite relatively high N_e and outcrossing mating systems in *A. halleri* and *A.*
455 *lyrata*, the differences in neutral mutations between the diploid species and the corresponding
456 subgenomes are far less remarkable in *A. kamchatica* than in *C. bursa-pastoris*. These data
457 suggest that N_e alone is not adequate to explain the proportion of neutral mutations.

458 The strongest signal for relaxed constraint that we detected in the *A. kamchatica*
459 subgenomes was observed when genes were categorized by expression levels. Genes that had
460 low expression showed a significant increase in the proportion of neutral mutations (30-32%)
461 over highly expressed genes (13-19%), and highly expressed genes show the strongest levels of
462 purifying selection (for $N_{e,s} > 10$, 73-77% of mutations) in either subgenome. This result is
463 consistent with expectations of stronger selective constraint on highly expressed genes⁴⁹. A
464 similar result was also found in the diploid *M. truncatula* where expression levels predicted very
465 clearly the proportion of neutral mutations⁵¹, adding further support that the method is able to
466 detect large differences in relaxed constraint when gene expression levels are taken into
467 account.

468 Although theoretical analysis typically assumes that deleterious mutations may be
469 masked by genome duplication, empirical studies showed that the dosage balance in gene
470 networks may be a selective constraint⁵² and could work as a mechanism for purifying selection
471 in an allopolyploid species. At this moment, the factors contributing to the difference between
472 *A. kamchatica* and *C. bursa-pastoris* are not clear. It is possible that the time since the
473 polyploidization events would not be adequate to detect the changes in the strength of purifying
474 selection, although the time estimates of polyploidization overlap to a large extent (about
475 20,000-250,000 years ago for *A. kamchatica*, 100,000-300,000 years ago for *C. bursa-pastoris*).

476

477 The Proportion of Adaptive Substitutions (α) are Significantly Greater Than Zero

478 This is the first report of α for *A. halleri* and *A. lyrata* using whole genome data, and to our
479 knowledge, the first report of genome-wide α for a polyploid species. Previous multi-species
480 comparisons showed that only a few plant species have α values that are greater than zero¹⁷,
481 however these estimates were mostly done using limited genetic data (< 1000 loci)^{17,19} rather
482 than genome-wide data. We estimated that 25-27% of non-synonymous substitutions are
483 adaptive in the two diploid species *A. halleri* and *A. lyrata*. These are the highest estimates of α

484 for any *Arabidopsis* species^{17,40} and higher than most plant species. The highest α among any
485 plant species was estimated in the highly outcrossing *Capsella grandiflora* ($\alpha = 0.4-0.7$)^{19,53} with
486 levels similar to *Drosophila* and bacteria, all taxa with large effective population sizes¹⁸. Our
487 results for the diploid species are consistent with previous studies that have shown a positive
488 correlation between α and N_e ^{17,20} which suggests that greater adaptive evolution often occurs in
489 species with large effective population sizes, which is true for both highly outcrossing diploid
490 species reported here.

491 Importantly, α for both subgenomes of *A. kamchatica* is also significantly greater than
492 zero and indicates 6-12% of non-synonymous substitutions are adaptive. Many diploid plant
493 species have a similar or larger effective population size than *A. kamchatica* (54,000-77,000), but
494 did not show positive α ¹⁷. For example, N_e estimated for *A. thaliana* was between 65,000 –
495 267,000^{17,20} while $\alpha = -0.08$ ¹⁹, indicating that effective population size alone cannot explain the
496 significantly positive α of *A. kamchatica*. These data suggest that *A. kamchatica* has a positive α
497 because of polyploidy. We suggest two mutually non-exclusive explanations. First, *A. kamchatica*
498 may have inherited fixed non-synonymous or adaptive substitutions from the two parental
499 species. The α values of *A. kamchatica* are roughly half of the parental species, in which the
500 reduction may be attributable to the reduction of N_e . Second, the rate of non-synonymous
501 mutations are increased at the early stages of polyploid species in contrast to slow rate of old
502 duplicated genes^{21,22}. A classic idea of the high evolvability of duplicated genomes states that
503 one of the duplicated copies may be able to obtain a new function or adaptive mutations
504 because the other copy retains the original function^{23,24}.

505

506 High Impact Mutations with Deleterious Effects Were Rarely Fixed

507 The loss of homeologs in ancient polyploids, or nonfunctionalization, has been extensively
508 studied²⁴, but relatively little is known about the population genetics of young polyploid
509 species. We identified high impact mutations that are likely to disrupt the gene function. We
510 found that about 20% of the homeologs in both subgenomes had disruptive mutations in our
511 collection of 25 individuals (Supplementary Table 7), although their frequencies are low (Fig. 5)
512 and only rarely are both homeologs disrupted. Interestingly, we found that high impact
513 mutations were rarely fixed. This is in contrast with the results from another allopolyploid
514 species *C. bursa-pastoris*, in which a large proportion of high-impact mutations (such as stop
515 codon gained) were fixed⁸. In *A. kamchatica*, similar proportions of high-impact mutations were

516 at low frequency compared with non-synonymous substitutions, which are also at low
517 frequency (Fig. 1B), suggesting that genome-wide purifying selection keeps their frequency low,
518 which is consistent with the prevalence of purifying selection shown by DoS and by DFE
519 methods.

520

521 Conclusion

522 Recently, new sequencing technology and algorithms drastically improved the genome assembly
523 of crop polyploid species with a large genome size⁵⁴⁻⁵⁶ which will facilitate the genome-wide
524 polymorphism analysis and scans for selection. By quantifying selection using polyploid species
525 with different population sizes, times since polyploidization and mating systems, general
526 patterns of selection in polyploid genomes will emerge. A further step will be to incorporate
527 polymorphism, gene expression, and species distribution data (i.e., landscape genomics) of
528 diploid parents and allopolyploid hybrids to identify the contributions of parental adaptations
529 for broadening climatic regimes and abiotic habitats in polyploids.

530

531 Materials and Methods

532 Allopolyploid plant samples and resequencing

533 *Arabidopsis kamchatica* (Fisch. ex DC.) K. Shimizu & Kudoh²⁷ is an allotetraploid species
534 distributed in East Asia and North America. We consider Russian individuals described as
535 *Cardaminopsis kamtschatika* or *Cardaminopsis lyrata* as synonyms (note that *Arabidopsis lyrata*
536 is a distinct diploid species)⁵⁷. Genomic DNA from 25 accessions of *A. kamchatica* was extracted
537 from leaf tissue using the DNeasy Plant Kit (Qiagen). These accessions were collected from
538 Taiwan, lowland and highland regions of Japan, Eastern Russia, Sakhalin Island, and Alaska, USA
539 (listed in Supplementary Table 3). DNA concentration and quality was measured using Qbit.
540 Genomic DNA libraries were constructed at the Functional Genomics Center Zurich (FGCZ) using
541 NEB Next Ultra. Total DNA was sequenced on Illumina HiSeq 2000 using paired end sequences
542 with an average insert size of 200-500 bp. Read lengths were 100 bp. For 22 accessions, a single
543 lane included six *A. kamchatica* DNA samples and for three accessions (KWS, MUR, and PAK),
544 eight samples per lane were used.

545

546 Illumina read mapping and sorting using v2.2 reference genomes

547 Illumina reads from *A. kamchatica* were mapped using BWA-MEM version 0.7.10 on the two
548 diploid genomes independently. We classified the reads to each parental origin as H-origin
549 (*halleri*-origin) and L-origin (*lyrata*-origin) using HomeoRoq (<http://seselab.org/homeorog>, last
550 accessed July 14, 2016). In this method, reads from each accession were first mapped to each
551 parental genome, and then classified as H-origin, L-origin, common, or unclassified (see fig. 1 in³²
552 for schematic diagram). Here, the ‘common’ reads are the reads that aligned equally well to
553 both parental genomes. After mapping to the *A. halleri* genome, we detected *A. kamchatica*
554 *halleri*-origin (H-origin) reads and identified single-nucleotide polymorphisms (SNPs) and short
555 insertions and deletions using GATK v3.3⁵⁸. Then, the nucleotides were replaced on the detected
556 variant position in the reference genome with the alternative nucleotides if the position (1)
557 covered by at least 20% of the average coverage of reads in each library, (2) covered by at most
558 twice of the average coverage and (3) has 30 or higher mutation detection quality (QUAL)
559 produced by GATK. This cycle of mapping, read classification, and reference modification, was
560 repeated ten times. For the reference modification, we used only *origin* reads the first five times
561 and both origin and common reads the last five times. The *A. kamchatica lyrata*-origin (L-origin)
562 genome was iteratively updated in a similar manner. The modified genomes were only used for
563 read sorting. Coverage was calculated for both subgenomes of our resequenced lines by using
564 the sum of the diploid parents as the genome size (250 + 225 = 475) and *common* plus sorted
565 *origin* reads (Supplementary Table 3).

566

567 Variant calling

568 For final variant calling, we combined the common reads of *A. kamchatica* with each sorted H-
569 origin or L-origin reads and aligned them back to the original parental genomes using BWA-MEM
570 v0.7.10. We called variants using GATK v3.3-0 following established best practices^{59,60}. We
571 processed each alignment BAM file separately to fix mate pairs, mark duplicates, and realign
572 reads around indels. Then we identified variants by running HaplotypeCaller jointly on all
573 genotypes but separately for each parental subgenome. To remove low-quality variants, we
574 mostly used the thresholds recommended for variant data sets where quality score cannot be
575 recalibrated⁶⁰. We applied quality by depth (QD < 2), mapping quality (MQ < 30), mapping
576 quality rank sum (MQRankSum < -15) and genotype quality (GQ < 20) filters. Because some of
577 our accessions had relatively low coverage, we considered that the recommended strand and
578 read position filters might be too strict and we did not apply them. Finally, we removed all

579 variants that GATK reported as heterozygous. We used diploid data from 9 accessions of
580 European *A. halleri* and 9 accessions of European *A. lyrata* from Novikova et al.³⁰ mapped to our
581 diploid reference genomes and called SNPs using the same criteria. The diploid VCF files were
582 then phased using Beagle⁶¹ to produce 18 alleles for each species.

583 Regions with excessively high coverage are likely to be repetitive or incorrectly
584 assembled, therefore variants called in those regions are probably spurious. To determine the
585 coverage thresholds, we summed up the coverage reported by bamtools⁶² for each position in
586 the final alignment files across all genotypes. We only considered reads with mapping quality
587 (MQ) of at least 20. Then, we calculated the mean and standard deviation for the distribution of
588 the obtained sums in each parental genome. We assumed a Poisson distribution and added 5
589 standard deviations to the mean to determine the thresholds. These thresholds (2891 and 2509
590 for *A. halleri* and *A. lyrata* respectively) were applied to the DP property (total depth of coverage
591 across all genotypes) in the INFO field of the corresponding VCF file. In addition, we applied a
592 coverage filter at genotype level to exclude calls with coverage below 2 or above 250.

593 To check for additional spurious variants, we randomly sampled 20 million reads (10
594 million per parent) from *A. halleri* and *A. lyrata* short-insert (200 bp) reads and ran it through
595 the same variant calling pipeline as the *A. kamchatica* genotypes. The only difference between
596 the runs was that this simulated sample was processed alone while variants for *A. kamchatica*
597 genotypes were called jointly. Any variants called with the simulated sample would be due to
598 incorrect read sorting between the parents or repetitive sequences present in the parental
599 genomes. Such spurious variants would also be likely to appear among *A. kamchatica* variants
600 even if the corresponding regions were completely conserved between *A. kamchatica* and its
601 parents. Among the uncovered variants, 59,856 and 58,645 were also present in *A. kamchatica*
602 on *A. halleri* and *A. lyrata* sides respectively. All of these variants were marked as filter failing.
603 When applying polymorphisms to the reference sequences, we used N's in positions where clear
604 calls could not be made due to insufficient coverage, excessive coverage, low quality
605 polymorphisms or heterozygosity. Such treatment allowed us to avoid using reference calls in
606 regions where the actual sequence is highly uncertain.

607

608 Coding sequence (CDS) alignments

609 We identified homeologous genes based on reciprocal blast hit (best-to-best with E-values < 10⁻
610 ¹⁵ and alignment length ≥ 200 bp) among coding sequences from the v2.2 *A. halleri* and *A. lyrata*

611 genome annotations. Using the same approach, we also detected orthologous relationships
612 between the predicted genes in diploid *A. halleri* and *A. lyrata* annotated genome assemblies
613 and *A. thaliana* genes (TAIR 10). In cases of duplicated genes of interest such as *HMA4*
614 (tandemly duplicated three times in *A. halleri*), we used only one copy for diversity analysis due
615 to non-unique alignments of Illumina reads and very high sequence identity (99%) in the *A.*
616 *halleri* reference genome. Therefore, our genome-wide dataset of coding sequences of
617 homeologs do not contain genes that are duplicated in one genome but not the other.

618 To make coding sequence alignments, we individually applied SNPs and deletions from
619 each of the 25 *A. kamchatica* genotypes (*H-origin* or *L-origin*) to the corresponding reference
620 genomes. We omitted insertions in order to preserve the genomic coordinates of the coding
621 sequences, which would consequently facilitate the alignment. If a variant was heterozygous,
622 failed the genotype quality filter (GQ < 20), or was not called for a particular genotype (but
623 called for other genotypes), the corresponding bases were replaced with N's. We assumed that
624 a sequence contains reference bases at positions that are not specified in VCF file and have
625 adequate coverage. Therefore, all bases with coverage < 2 (insufficient) or > 250 (abnormally
626 high) were replaced with N's. After that, we extracted coding sequences from the modified
627 genomes and grouped them by gene. Thus, each H-origin or L-origin gene had an alignment file
628 containing 25 aligned coding sequences (one for each genotype). Finally, we aligned *A. thaliana*
629 orthologs as an outgroup using Muscle v3.8⁶³. With the profile alignment option, which
630 preserved the alignment of the ingroup sequences and only aligned the outgroup sequence to
631 the core ingroup alignment. The same procedure was used for making gene alignments of the 18
632 phased alleles for diploid *A. halleri* and *A. lyrata*.

633

634 Population structure and phylogenetic analysis

635 We used 1000 randomly selected coding sequence (CDS) alignments from both *halleri* and *lyrata*
636 derived homeologs. We then individually concatenated the *halleri* alignments and the *lyrata*
637 alignments to use for population structure and phylogenetic analysis. The input data sets for the
638 population structure analysis contained 21,341 and 16,223 markers from *halleri*- and *lyrata*-
639 origin CDS respectively. We ran STRUCTURE v2.3.4⁶⁴ ten times for each K = 1 to 9 using the
640 admixture model and 50,000 MCMC rounds for burnin followed by 100,000 rounds to generate
641 the data. The output was analyzed with STRUCTURE HARVESTER v0.6.94 and clusters were
642 rearranged with CLUMPP v1.1.2. For phylogenetic analysis, we added *A. halleri* and *A. lyrata* as

643 outgroups and ran Mr. Bayes v3.2.6⁶⁵ with default parameters for 500,000 generations sampling
644 every 1000th generation.

645

646 Coding sequence diversity and site frequency spectra

647 For gene alignments containing coding sequences, summary and diversity statistics, including
648 divergence from *A. thaliana*, were estimated using *libsequence* packages⁶⁶ and custom R, Perl,
649 and Ruby shell scripts. The *libsequence* programs *compute* and *Hcalc* were used to estimate
650 average pairwise diversity (π), θ_w , Tajima's D, Fay and Wu's H. Non-synonymous and
651 synonymous diversity and gene based allele frequencies were estimated using the *polydNdS*
652 program with the $-P$ flag to generate SNP tables for each gene. The site frequency spectra (SFS),
653 were created using the SFS.pl program available from the J. Ross-Ibarra
654 ([http://www.plantsciences.ucdavis.edu/faculty/ross-](http://www.plantsciences.ucdavis.edu/faculty/ross-ibarra/code/files/ea3bd485e4c7dee37c59e8ba77ca800e-11.html)
655 [ibarra/code/files/ea3bd485e4c7dee37c59e8ba77ca800e-11.html](http://www.plantsciences.ucdavis.edu/faculty/ross-ibarra/code/files/ea3bd485e4c7dee37c59e8ba77ca800e-11.html)) on the set of non-
656 synonymous and synonymous polymorphisms identified using *polydnds*. Both folded and
657 unfolded SFS were calculated; the folded spectrum does not differentiate between ancestral
658 polymorphisms and polymorphism that are the result of mutations that have entered a
659 population since it split from a common ancestor, while the unfolded spectra are based on
660 derived allele frequencies. We converted the SFS data to SFS count tables using a custom python
661 script (*sfs_extraction.py*). We used two published mutation rates, one based on the synonymous
662 substitution rates calibrated by fossil records⁴¹, and another for total sites in mutation
663 accumulation lines⁴², to estimate the effective population size using the following equation: $N_e =$
664 π_{syn} or $\pi_{total}/4\mu$ (where π was estimated from our data and μ from^{41,42}).

665

666 Linkage disequilibrium and sliding window diversity

667 To conduct sliding window analyses along entire scaffolds, we used the PopGenome R⁶⁷ package
668 to calculate diversity of all, intergenic, coding, exonic, and intron regions of *A. kamchatica* using
669 *A. halleri* or *A. lyrata* derived VCF and reference gene annotation (.gff) files. We estimated the
670 average nucleotide diversity, Watterson's θ_w and π (the average number of pairwise nucleotide
671 differences per site). To estimate genome-wide linkage disequilibrium (LD), we used the geno-
672 r2 option in VCFtools⁶⁸ across window sizes of a maximum distance of 20 kb, 50 kb or 1 Mb using
673 a minor allele frequency ≥ 0.1 , separately for the *halleri* or *lyrata* derived VCF files. The resulting

674 r^2 between SNPs were grouped into bins of 50 bp length. We estimated the average, 50% and
675 90% confidence intervals of correlation coefficients of each bin.

676

677 Direction of selection (DoS), Distribution of Fitness Effects (DFE) and Adaptive

678 Substitutions (α)

679 The program *MKtest* from the libsequence library, was used to count the total number of
680 polymorphic non-synonymous (P_n) and synonymous (P_s) sites in *A. kamchatica* homeologs as
681 well as the number of fixed non-synonymous (D_n) and synonymous (D_s) differences between *A.*
682 *kamchatica* homeologs and *A. thaliana*. We used the program *MKtest* to perform standard tests
683 on each gene for both homeologs separately; this is a contingency test comparing the numbers
684 of between species difference and within species polymorphisms at non-synonymous and
685 synonymous sites where significance is tested using Fisher's exact tests for each gene.

686 Polymorphism and divergence data was used to calculate the direction of selection (DoS
687 = $D_n/(D_n + D_s) - P_n/(P_n + P_s)$) statistic of ¹⁴. DoS < 0 is consistent with purifying selection and DoS >
688 0 is consistent with positive selection. To estimate the distribution of fitness effects (DFE, i.e. the
689 distribution of the strength of selection acting against new mutations) and the proportion of
690 adaptive substitutions (α) in *A. kamchatica*, *A. halleri* and *A. lyrata*, we used the likelihood
691 method implemented in the software DoFE 3.0¹³. The program was run for 1×10^6 steps, and
692 sampled every 1,000 steps after a burn in of 100,000 steps. Strongly deleterious mutations have
693 $N_e s > 10$ (where N_e is the effective population size and s is the selection coefficient), mildly
694 deleterious mutations have $1 < N_e s < 10$, and effectively neutral mutations have $N_e s < 1$. To
695 estimate DFE we used folded allele frequency spectra and the estimated number of non-
696 synonymous (D_n) and synonymous (D_s) differences between *A. kamchatica* homeologs or diploid
697 orthologs and the corresponding outgroup *A. thaliana* orthologs.

698

699 Transcriptome data

700 We used RNA-seq data collected from roots and leaf tissue of the *A. kamchatica* Murodo (Japan)
701 and Potter (Alaska, USA) accessions from Paape et al. 2016 to calculate expression for all
702 homeologs in our dataset. We mapped the RNA-seq data to v2.2 *A. halleri* and *A. lyrata*
703 reference genomes and sorted the reads using method described in Akama et al. 2014. Thus, for
704 each gene in our polymorphism dataset, we obtained expression data that is specific to either
705 homeolog. We estimated expression levels using HTseq to count reads, then calculated reads

706 per kilobase of transcript per million mapped reads (RPKM). The mean RPKM values from three
707 libraries of both leaf and root were used to make a distribution in RPKM that corresponds to our
708 polymorphism gene dataset. The distribution of RPKM was used to determine the upper and
709 lower 10% tails in expression for both homeologs separately.

710

711 Detection of High Impact Mutations

712 We used SnpEff v4.2⁶⁹ to detect genetic variants that have putative loss of function mutations in
713 both subgenomes of *A. kamchatica*. We ran the program separately on the variant file of each
714 subgenome. First, we built custom databases for each parental genome using our v2.2 parental
715 assemblies and annotation. Since SnpEff ignores filter fields in VCF files, we have removed all
716 variants that failed our filters, replaced all genotypes that failed genotype filters with no-calls
717 (i.e. './.'), and removed any entries without valid variant calls. Such filtering allowed us to extract
718 accurate gene summaries from SnpEff output.

719 SnpEff annotated polymorphisms within 32,410 and 31,119 genic regions in *A. halleri*
720 derived and *A. lyrata* derived genomes respectively. These include all mutations with any impact
721 type, but we focused only on frameshifts, premature stop codon, loss of stop codons and loss of
722 start codons. The gene sets were thus reduced to 31,193 and 31,119 genes for *A. halleri* and *A.*
723 *lyrata* derived genomes respectively. There are 21,419 and 21,463 reciprocal best BLAST hits
724 between respectively *A. halleri* or *A. lyrata* and *A. thaliana*. Based on the intersection of these
725 two data sets, we identified 20,292 homeologs between *A. halleri* and *A. lyrata*. Out of these 19
726 and 18 *halleri*-origin and *lyrata*-origin genes had no coverage. Gene ontology (GO) analysis was
727 performed using agriGO (bioinfo.cau.edu.cn/agriGO) conducted using a custom annotation
728 containing 19,936 GO annotations that correspond to *A. thaliana* orthologs with reciprocal-best
729 BLAST hits for both homeologs. We used only queries with at least 20 genes.

730

731

732 Acknowledgments

733 We thank Takashi Tsuchimatsu, Polina Novikova and Peter Keightley for useful discussions, and
734 the Functional Genomic Center Zurich for sequencing services and technical support. The study
735 was supported by Swiss National Science Foundation, the University Research Priority Program
736 of Evolution in Action of the University of Zurich, JST CREST Grant (number JPMJCR16O3) to KKS,
737 MEXT KAKENHI Grant Number 16H06469, 26113709, Young Investigator Award of Human

738 Frontier Science Program to KKS and JS, European Union's Seventh Framework Programme for
739 research, technological development and demonstration under grant agreement no GA-2010-
740 267243 – PLANT FELLOWS to RVB and TP, Marie-Heim Hoegtlin grant by Swiss National Science
741 Foundation to RSI, ISCB (Indo-Swiss Collaboration in Biotechnology) to KKS and MH, the Special
742 Coordination Funds for Promoting Science and Technology from MEXT Japan, an Inamori
743 Foundation research grant, a Japan Society for the Promotion of Science Grant-in-Aid for
744 Scientific Research (Young Researchers B, 2277023), and Research and Education Funding for
745 Japanese Alps Inter-Universities Cooperative Project, MEXT, Japan to KT.

746

747

748 **Data Accessibility**

749 Illumina reads submitted to DDBJ. BioProject Submission ID: PSUB006170

750 The Sanger sequences were submitted to GenBank. GenBank BankIt submission. Submission ID:
751 2025864

752

753 Code for *A. lyrata* genome assembly

754 <https://gitlab.com/rbrisk/AlyrAssembly>

755

756 Code for Variant calling in *A. kamchatica*

757 <https://gitlab.com/rbrisk/AkamVariants>

758

759

760

761

762

763

764

765

766

767

768

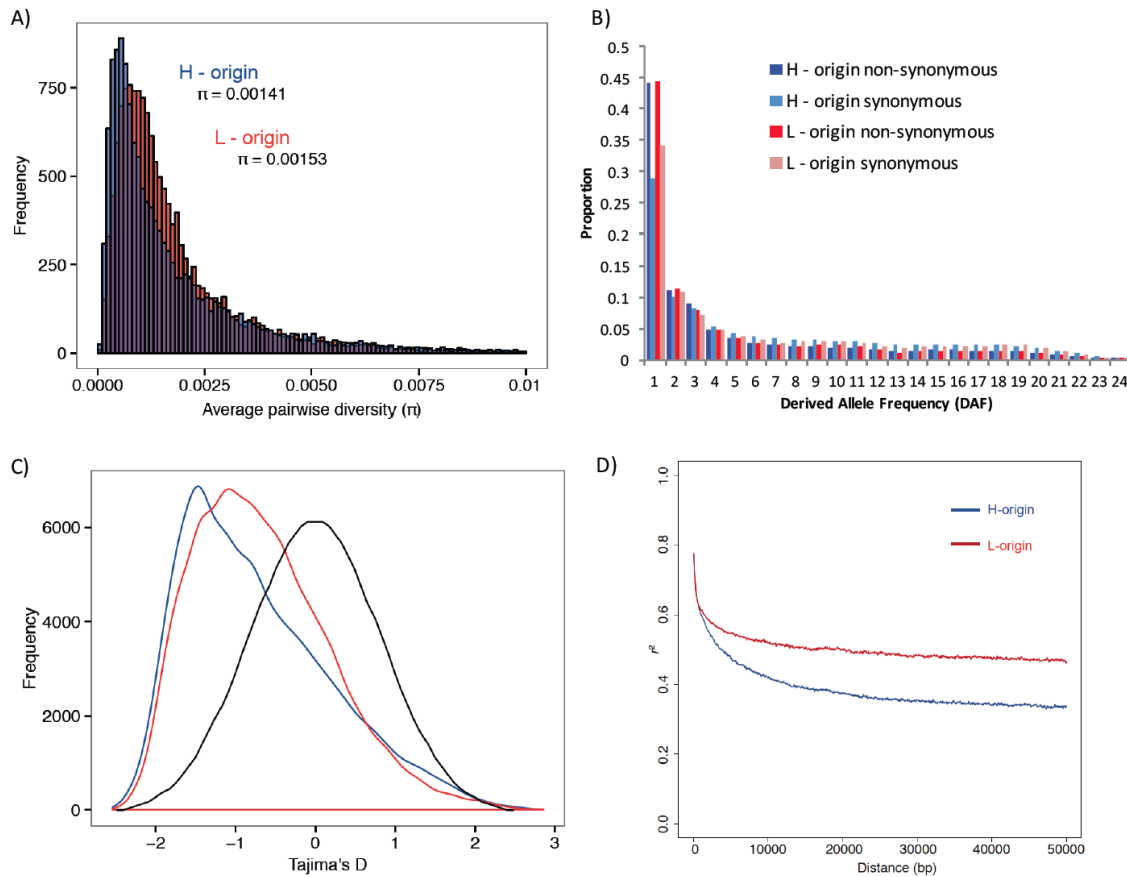
769

770

771

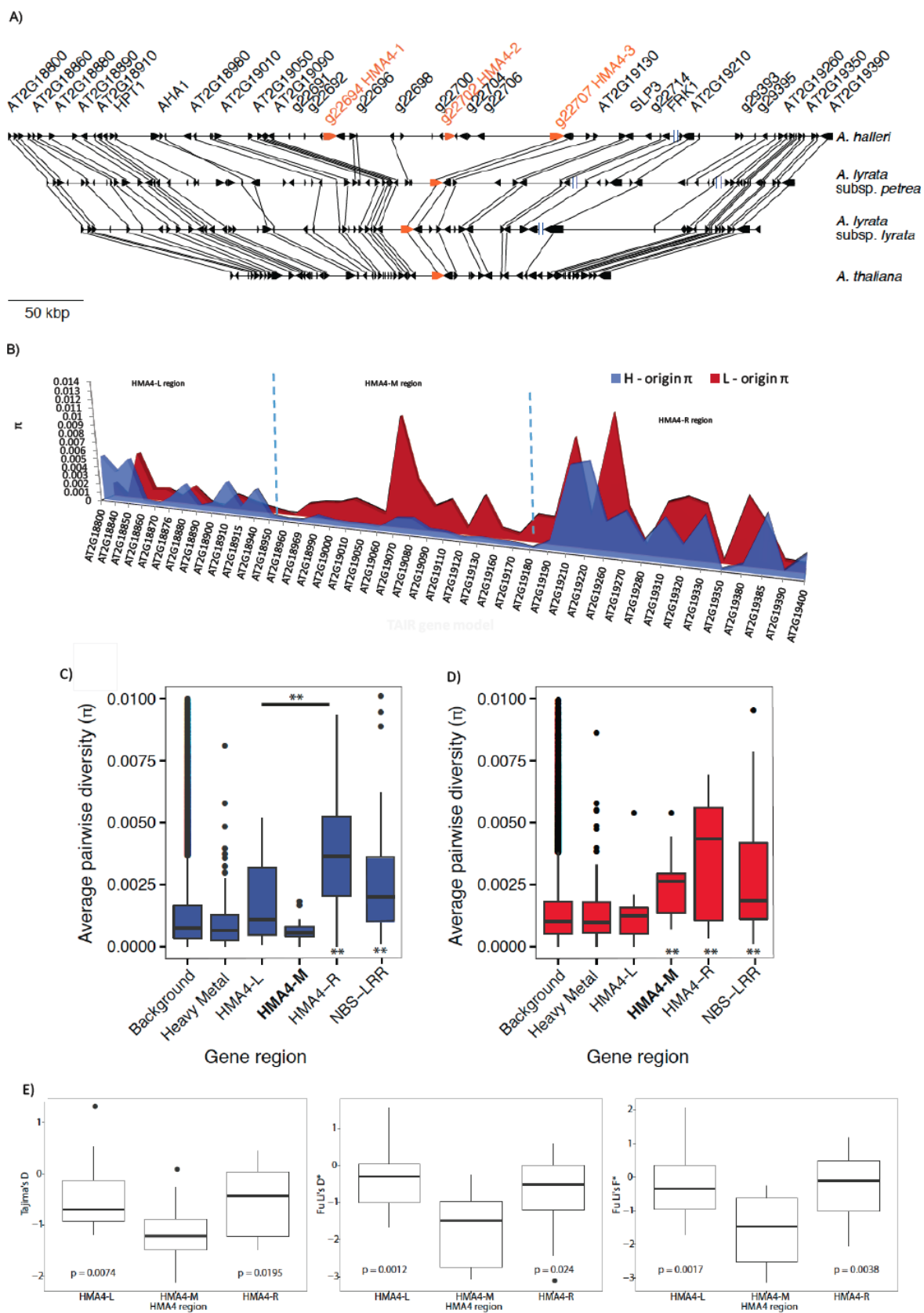
772

773 **Figures.**
774
775



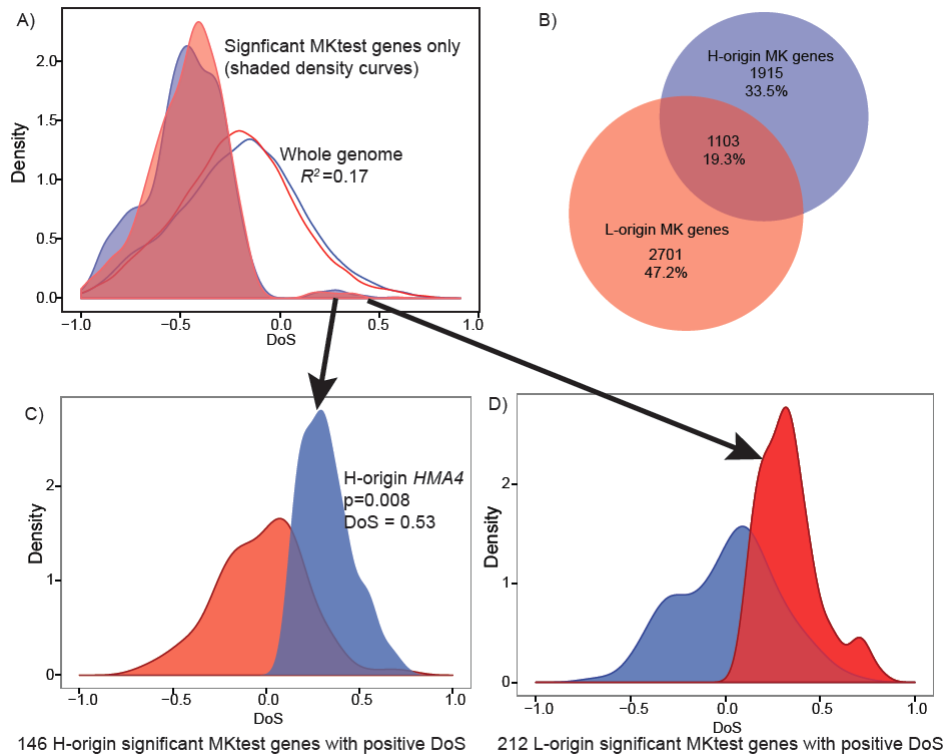
776

777 **Fig. 1. Genome-wide diversity and linkage disequilibrium.** (A) Average pairwise diversity (π) of
778 *halleri*-origin (H-origin) and *lyrata*-origin (L-origin) coding sequences. (B) H-origin and L-origin
779 genes show no significant differences in proportions of non-synonymous and synonymous
780 substitutions (χ^2 , p-value = 0.58), and the majority of substitutions are low at frequency. (C)
781 Tajima's D distributions for both genomes (blue density curve = H-origin, red density curve = L-
782 origin) show departures from neutrality (black density curve where neutral = 0), mean values for
783 both distributions are negative (Table 1). (D) The mean decay of linkage disequilibrium (LD)
784 estimated using 100 kb sliding windows shows mean LD decay < 10 kb for both H-origin (blue)
785 and L-origin (red) genomes.



787 **Fig. 2. Genetic diversity of the syntenic *HMA4* region.** (A) Synteny of the *HMA4* region from *A.*
788 *halleri* v2.2³⁹, *A. lyrata* subsp. *petraea* v2.2, *A. lyrata* subsp. *lyrata* (JGI)⁷⁰ and *A. thaliana* (TAIR).
789 (B) Average pairwise diversity (π) of genes surrounding the *HMA4* region in both homeologs of
790 *A. kamchatica*. (C) For the *halleri*-subgenome, genetic diversity of NBS-LRRs is significantly
791 greater (two asterisks below, $**p < 0.001$) than diversity compared with the background while
792 heavy metal (HM) genes show no significant difference. Diversity for both HMA4-L ($\pi = 0.0018$)
793 and HMA4-R ($\pi = 0.004$) are significantly higher than the HMA4-M (which contains the *HMA4*
794 coding sequences) region (two asterisks above HMA4-M, $**p < 0.001$). (D) For the *lyrata*-
795 subgenome, diversity of NBS-LRRs, HMA4-M and HMA4-R are all significantly higher than the
796 background. The diversity of the *lyrata* HMA4-M ($\pi = 0.0032$) region is also significantly greater
797 than the *halleri* HMA4-M region ($\pi = 0.0007$, paired *t-test* p-value = 0.003; Wilcoxon sign rank p-
798 value = 0.0001). The neutrality statistics (E) Tajima's D, Fu and Li's D* and Fu and Li's F* all show
799 the *halleri*-origin HMA4-M region to be significantly lower than the left and right flanking regions
800 supporting genetic hitchhiking surrounding the *HMA4* coding sequences.

801
802
803
804
805
806
807
808
809
810
811



812

813

Fig. 3. The direction of selection for both subgenomes. (A) Density curves of the direction of selection (DoS)¹⁴ for about 21,000 coding sequences (blue line and density curves are DoS for H-origin genes, red line and density are DoS for L-origin genes). Neutral genes are indicated by 0, while negative values indicate purifying selection and positive values indicate positive selection. The means of these distributions are -0.20 and -0.22 for the H- and L-origin homeologs respectively, and show that ~70% of both homeologs have a negative selection index (negative DoS). Shaded density curves are genes that were significant for MK-tests ($p < 0.05$ using Fisher's marginal p-values). (B) Only 19% of genes show significance for MK-tests for both homeologs. (C) Using only significant MK-test genes with positive DoS for *halleri*-origin and (D) positive DoS for *lyrata*-origin genes show that the other homeolog has significantly more negative DoS (p -value $< 2.2e-16$ using pairwise t-test and Wilcoxon signed rank test) when one shows positive selection using comparisons of DoS distributions in both C and D.

825

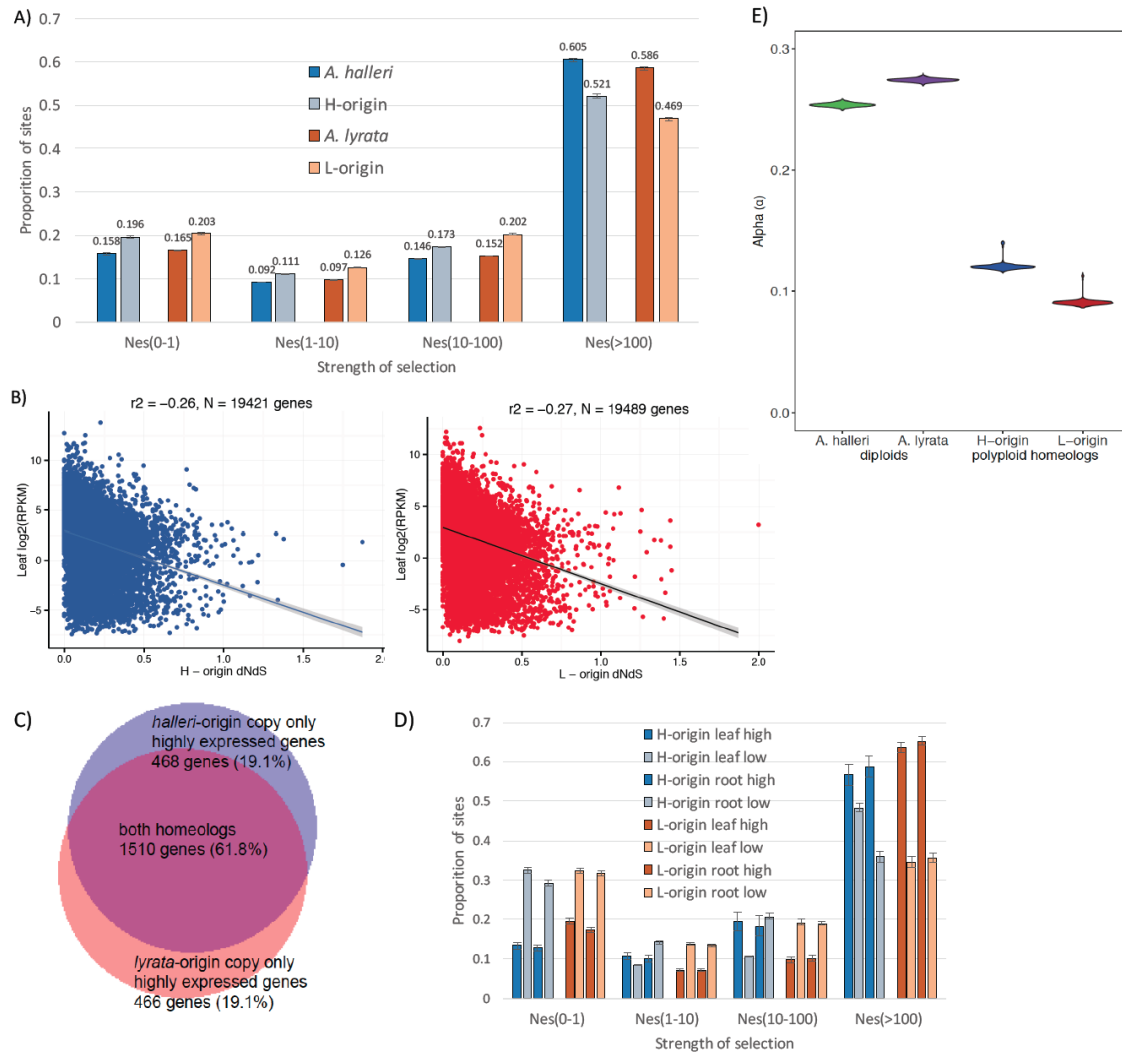
826

827

828

829

830



831

832 **Fig. 4. The strength of purifying selection and adaptive evolution.** (A) The distribution of fitness

833 effects (DFE) of deleterious mutations for coding sequences of the two *A. kamchatica*

834 subgenomes and corresponding diploid orthologs of *A. halleri* and *A. lyrata*. The strength of

835 selection is indicated by $N_e s$ where N_e is the effective population size and s is the selection

836 coefficient. Error bars show standard deviations. (B) Evolutionary rates are negatively correlated

837 with gene expression in both homeologs. (C) Overlap of genes that are highly expressed leaf

838 tissues (in upper 10% of expression) in both homeologs. (D) DFE categorized by expression in

839 both subgenomes. Expression categories were taken from the upper 10% (high) and lower 10%

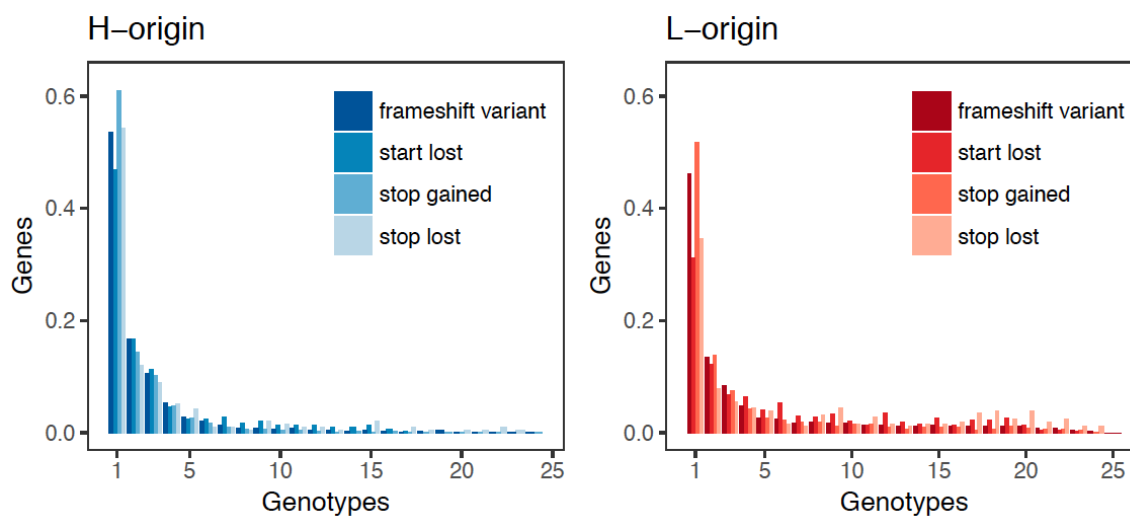
840 (low) of expression distribution in all *A. kamchatica* homeologs. (E) The proportion of adaptive

841 substitutions (α) for both subgenomes (H-origin $\alpha = 0.12$, CI: 0.117-0.141, L-origin $\alpha = 0.09$, CI:

842 0.087-0.094) and for the two corresponding diploid species (*A. halleri* $\alpha = 0.25$, CI: 0.251-0.257,

843 *A. lyrata* $\alpha = 0.27$, CI: 0.272-0.277) are significantly greater than zero.

844



845

846

Fig. 5. Frequency distributions of high impact mutations. Large effect mutations are at low

847

frequency for both subgenomes.

848

849

850

851

852

853

854

855

856

857

858

859

860

861

862

863

864

865

866

867

868

869

870

871

872

873

874

875

876

877 **Tables**

878

879 Table 1. Reference genome assembly statistics of v2.2 of Siberian *A. lyrata* subsp. *petraea* and

880 v2.2 of *A. halleri* subsp. *gemmaifera* (Tada Mine).

881

Assembly statistics	<i>A. lyrata</i>	<i>A. halleri</i> ^a
Length (bp)	175,182,717	196,243,198
Missing (%) ^b	12.75	14.81
Scaffolds	1,675	2,239
Shortest scaffolds (bp)	940	932
Longest scaffolds (bp)	6,771,235	4,302,264
Scaffold N50 length (bp)	1,260,070	712,249
Scaffold N50 count	38	71
NG50 length (bp)	804,357	489,153
NG50 count	63	117
Genome size (FC) ^c	225 Mb	250 Mb

a: *A. halleri* assembly previously reported in Briskine et al.³⁹

b: Missing data counted as number of N's in the assembly and is percentage of total length

c: Genome size measured by flow cytometry (FC)

882

883

884

885

886

887

888

889

890

891

892

893

894

895

896

897

898 Table 2. Diversity statistics for coding sequences (CDS) of *A. kamchatica* homeologs. Values are
 899 average pairwise diversity, π , polymorphism Watterson's estimator, θ_w , Tajima's D, Fay and
 900 Wu's H. Correlations between homeolog diversity statistics are shown as R^2 correlation
 901 coefficient.

statistic	<i>halleri</i> homeologs			<i>lyrata</i> homeologs			R^2	combined		
	mean	sd	n	mean	sd	n		mean	sd	n
π_{total}	0.0014	0.0019	21419	0.0015	0.0018	21463	0.27	0.0015	0.0015	20249
θ_w	0.0017	0.0018	21419	0.0018	0.0019	21463	0.30	0.0017	0.0015	20249
π_{nonsyn}	0.0011	0.0029	20605	0.0014	0.0033	20696	0.15	0.0012	0.0025	19953
π_{syn}	0.0044	0.0116	20605	0.0049	0.0343	20696	0.04	0.0046	0.0182	20249
θ_w_{nonsyn}	0.0014	0.0029	20605	0.0017	0.0034	20696	0.14	0.0012	0.0025	19953
θ_w_{syn}	0.0047	0.0114	20605	0.0056	0.0344	20696	0.04	0.0051	0.0183	19953
Taj D	-0.72	0.89	19691	-0.63	0.85	19984	0.03	-0.67	0.66	19940
FayWuH	-0.41	1.23	19574	-0.49	1.28	19893	0.07	-0.44	0.97	19911

902 nonsynon = non-synonymous substitutions

903 synon = synonymous substitutions

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927 **References**

- 928 1. Wood, T. E. *et al.* The frequency of polyploid speciation in vascular plants. *Proc. Natl. Acad. Sci.* **106**,
929 13875–13879 (2009).
- 930 2. Renny-Byfield, S. & Wendel, J. F. Doubling down on genomes: Polyploidy and crop plants. *Am. J. Bot.*
931 **101**, 1711–1725 (2014).
- 932 3. Comai, L. The advantages and disadvantages of being polyploid. *Nat. Rev. Genet.* **6**, 836–846 (2005).
- 933 4. Soltis, D. E., Visger, C. J. & Soltis, P. S. The polyploidy revolution then...and now: Stebbins revisited.
934 *Am. J. Bot.* **101**, 1057–1078 (2014).
- 935 5. Dufresne, F., Stift, M., Vergilino, R. & Mable, B. K. Recent progress and challenges in population
936 genetics of polyploid organisms: an overview of current state-of-the-art molecular and statistical
937 tools. *Mol. Ecol.* **23**, 40–69 (2014).
- 938 6. Buggs, R. J. A. *et al.* Next-generation sequencing and genome evolution in allopolyploids. *Am. J. Bot.*
939 **99**, 372–382 (2012).
- 940 7. Clevenger, J., Chavarro, C., Pearl, S. A., Ozias-Akins, P. & Jackson, S. A. Single Nucleotide
941 Polymorphism Identification in Polyploids: A Review, Example, and Recommendations. *Mol. Plant* **8**,
942 831–846 (2015).
- 943 8. Douglas, G. M. *et al.* Hybrid origins and the earliest stages of diploidization in the highly successful
944 recent polyploid *Capsella bursa-pastoris*. *Proc. Natl. Acad. Sci.* **112**, 2806–2811 (2015).
- 945 9. Arnold, B. J. *et al.* Borrowed alleles and convergence in serpentine adaptation. *Proc. Natl. Acad. Sci.*
946 **113**, 8320–8325 (2016).
- 947 10. Paape, T. *et al.* Conserved but Attenuated Parental Gene Expression in Allopolyploids: Constitutive
948 Zinc Hyperaccumulation in the Allotetraploid *Arabidopsis kamchatica*. *Mol. Biol. Evol.* **33**, 2781–
949 2800 (2016).
- 950 11. Novikova, P. *et al.* Genome sequencing reveals the origin of the allotetraploid *Arabidopsis suecica*.
951 *Mol. Biol. Evol.* msw299 (2017). doi:10.1093/molbev/msw299
- 952 12. Nielsen, R. MOLECULAR SIGNATURES OF NATURAL SELECTION. *Annu. Rev. Genet.* **39**, 197–218
953 (2005).

- 954 13. Eyre-Walker, A. & Keightley, P. D. Estimating the Rate of Adaptive Molecular Evolution in the
955 Presence of Slightly Deleterious Mutations and Population Size Change. *Mol. Biol. Evol.* **26**, 2097–
956 2108 (2009).
- 957 14. Stoletzki, N. & Eyre-Walker, A. Estimation of the Neutrality Index. *Mol. Biol. Evol.* **28**, 63–70 (2010).
- 958 15. Hough, J., Williamson, R. J. & Wright, S. I. Patterns of Selection in Plant Genomes. *Annu. Rev. Ecol.*
959 *Evol. Syst.* **44**, 31–49 (2013).
- 960 16. Otto, S. P. & Whitton, J. Polyploid Incidence and Evolution. *Annu. Rev. Genet.* **34**, 401–437 (2000).
- 961 17. Gossmann, T. I. *et al.* Genome Wide Analyses Reveal Little Evidence for Adaptive Evolution in Many
962 Plant Species. *Mol. Biol. Evol.* **27**, 1822–1832 (2010).
- 963 18. Siol, M., Wright, S. I. & Barrett, S. C. H. The population genomics of plant adaptation. *New Phytol.*
964 **188**, 313–332 (2010).
- 965 19. Slotte, T., Foxe, J. P., Hazzouri, K. M. & Wright, S. I. Genome-Wide Evidence for Efficient Positive and
966 Purifying Selection in *Capsella grandiflora*, a Plant Species with a Large Effective Population Size.
967 *Mol. Biol. Evol.* **27**, 1813–1821 (2010).
- 968 20. Gossmann, T. I., Keightley, P. D. & Eyre-Walker, A. The Effect of Variation in the Effective Population
969 Size on the Rate of Adaptive Molecular Evolution in Eukaryotes. *Genome Biol. Evol.* **4**, 658–667
970 (2012).
- 971 21. Jordan, I. K., Wolf, Y. I. & Koonin, E. V. Duplicated genes evolve slower than singletons despite the
972 initial rate increase. *BMC Evol. Biol.* **4**, 22 (2004).
- 973 22. Yang, L. & Gaut, B. S. Factors that Contribute to Variation in Evolutionary Rate among Arabidopsis
974 Genes. *Mol. Biol. Evol.* **28**, 2359–2369 (2011).
- 975 23. Ohno, S. *Evolution by Gene Duplication*. (Springer Berlin, 2014).
- 976 24. Lynch, M. & Conery, J. The Evolutionary Fate and Consequences of Duplicate Genes. *Science* **290**,
977 1151–1155 (2000).
- 978 25. Shimizu, K. K. & Tsuchimatsu, T. Evolution of Selfing: Recurrent Patterns in Molecular Adaptation.
979 *Annu. Rev. Ecol. Evol. Syst.* **46**, 593–622 (2015).

- 980 26. Bomblies, K. & Madlung, A. Polyploidy in the Arabidopsis genus. *Chromosome Res.* **22**, 117–134
981 (2014).
- 982 27. Shimizu, Kentaro K, Fuji, S, Marhold, Karol, Watanabe, Kunaiki & Kudoh, Hiroshi. Arabidopsis
983 kamchatica (Fisch. ex DC.) K. Shimizu & Kudoh and A. kamchatica subsp. kawasakiana (Makino)
984 K. Shimizu & Kudoh, New Combinations. *Acta Phytotaxon. Geobot.* **56**, (2005).
- 985 28. Tsuchimatsu, T., Kaiser, P., Yew, C.-L., Bachelier, J. B. & Shimizu, K. K. Recent Loss of Self-
986 Incompatibility by Degradation of the Male Component in Allotetraploid Arabidopsis kamchatica.
987 *PLoS Genet.* **8**, e1002838 (2012).
- 988 29. Shimizu-Inatsugi, R. *et al.* The allopolyploid *Arabidopsis kamchatica* originated from multiple
989 individuals of *Arabidopsis lyrata* and *Arabidopsis halleri*. *Mol. Ecol.* **18**, 4024–4048 (2009).
- 990 30. Novikova, P. Y. *et al.* Sequencing of the genus Arabidopsis identifies a complex history of
991 nonbifurcating speciation and abundant trans-specific polymorphism. *Nat. Genet.* (2016).
992 doi:10.1038/ng.3617
- 993 31. Schmickl, R., Jørgensen, M. H., Brysting, A. K. & Koch, M. A. The evolutionary history of the
994 Arabidopsis lyrata complex: a hybrid in the amphi-Beringian area closes a large distribution gap and
995 builds up a genetic barrier. *BMC Evol. Biol.* **10**, 98 (2010).
- 996 32. Akama, S., Shimizu-Inatsugi, R., Shimizu, K. K. & Sese, J. Genome-wide quantification of homeolog
997 expression ratio revealed nonstochastic gene regulation in synthetic allopolyploid Arabidopsis.
998 *Nucleic Acids Res.* **42**, e46–e46 (2014).
- 999 33. Armstrong, J. J., Takebayashi, N., Sformo, T. & Wolf, D. E. Cold tolerance in Arabidopsis kamchatica.
1000 *Am. J. Bot.* **102**, 439–448 (2015).
- 1001 34. Hoffmann, M. H. EVOLUTION OF THE REALIZED CLIMATIC NICHE IN THE GENUS: ARABIDOPSIS
1002 (BRASSICACEAE). *Evolution* **59**, 1425–1436. (2005).
- 1003 35. Kenta, T. Clinal Variation in Flowering Time and Vernalisation Requirement across a 3000-M
1004 Altitudinal Range in Perennial Arabidopsis kamchatica Ssp.Kamchatica and Annual Lowland
1005 Subspecies Kawasakiana. *J. Ecosyst. Ecography* **03**, (2013).

- 1006 36. Roux, C. *et al.* Does Speciation between *Arabidopsis halleri* and *Arabidopsis lyrata* Coincide with
1007 Major Changes in a Molecular Target of Adaptation? *PLoS ONE* **6**, e26872 (2011).
- 1008 37. Hanikenne, M. *et al.* Evolution of metal hyperaccumulation required cis-regulatory changes and
1009 triplication of HMA4. *Nature* **453**, 391–395 (2008).
- 1010 38. Hanikenne, M. *et al.* Hard Selective Sweep and Ectopic Gene Conversion in a Gene Cluster Affording
1011 Environmental Adaptation. *PLoS Genet.* **9**, e1003707 (2013).
- 1012 39. Briskine, R. V. *et al.* Genome assembly and annotation of *Arabidopsis halleri* , a model for heavy
1013 metal hyperaccumulation and evolutionary ecology. *Mol. Ecol. Resour.* (2016). doi:10.1111/1755-
1014 0998.12604
- 1015 40. Slotte, T. *et al.* Genomic Determinants of Protein Evolution and Polymorphism in *Arabidopsis*.
1016 *Genome Biol. Evol.* **3**, 1210–1219 (2011).
- 1017 41. Koch, M. A., Haubold, B. & Mitchell-Olds, T. Comparative Evolutionary Analysis of Chalcone
1018 Synthase and Alcohol Dehydrogenase Loci in *Arabidopsis*, *Arabis*, and Related Genera
1019 (Brassicaceae). *Mol. Biol. Evol.* **17**, 1483–1498 (2000).
- 1020 42. Ossowski, S. *et al.* The Rate and Molecular Spectrum of Spontaneous Mutations in *Arabidopsis*
1021 *thaliana*. *Science* **327**, 92–94 (2010).
- 1022 43. Nordborg, M. *et al.* The Pattern of Polymorphism in *Arabidopsis thaliana*. *PLoS Biol.* **3**, e196 (2005).
- 1023 44. Branca, A. *et al.* PNAS Plus: Whole-genome nucleotide diversity, recombination, and linkage
1024 disequilibrium in the model legume *Medicago truncatula*. *Proc. Natl. Acad. Sci.* (2011).
1025 doi:10.1073/pnas.1104032108
- 1026 45. Cao, J. *et al.* Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat. Genet.*
1027 **43**, 956–963 (2011).
- 1028 46. Verbruggen, N., Hermans, C. & Schat, H. Molecular mechanisms of metal hyperaccumulation in
1029 plants: Tansley review. *New Phytol.* **181**, 759–776 (2009).
- 1030 47. Guo, Y.-L. *et al.* Genome-Wide Comparison of Nucleotide-Binding Site-Leucine-Rich Repeat-
1031 Encoding Genes in *Arabidopsis*. *PLANT Physiol.* **157**, 757–769 (2011).

- 1032 48. Marone, D., Russo, M., Laidò, G., De Leonardis, A. & Mastrangelo, A. Plant Nucleotide Binding Site–
1033 Leucine-Rich Repeat (NBS-LRR) Genes: Active Guardians in Host Defense Responses. *Int. J. Mol. Sci.*
1034 **14**, 7302–7326 (2013).
- 1035 49. Zhang, J. & Yang, J.-R. Determinants of the rate of protein sequence evolution. *Nat. Rev. Genet.* **16**,
1036 409–420 (2015).
- 1037 50. Yew, C.-L., Kakui, H. & Shimizu, K. K. Agrobacterium-mediated floral dip transformation of the model
1038 polyploid species *Arabidopsis kamchatica*. *J. Plant Res.* (2017). doi:10.1007/s10265-017-0982-9
- 1039 51. Paape, T. *et al.* Selection, genome-wide fitness effects and evolutionary rates in the model legume
1040 *Medicago truncatula*. *Mol. Ecol.* **22**, 3525–3538 (2013).
- 1041 52. Bekaert, M., Edger, P. P., Pires, J. C. & Conant, G. C. Two-Phase Resolution of Polyploidy in the
1042 *Arabidopsis* Metabolic Network Gives Rise to Relative and Absolute Dosage Constraints. *Plant Cell*
1043 **23**, 1719–1728 (2011).
- 1044 53. Williamson, R. J. *et al.* Evidence for Widespread Positive and Negative Selection in Coding and
1045 Conserved Noncoding Regions of *Capsella grandiflora*. *PLoS Genet.* **10**, e1004622 (2014).
- 1046 54. Hatakeyama, M. *et al.* Multiple hybrid de novo genome assembly of finger millet, an orphan
1047 allotetraploid crop. *DNA Res.* (2017). doi:10.1093/dnares/dsx036
- 1048 55. Yang, J. *et al.* The genome sequence of allopolyploid *Brassica juncea* and analysis of differential
1049 homoeolog gene expression influencing selection. *Nat. Genet.* **48**, 1225–1232 (2016).
- 1050 56. Avni, R. *et al.* Wild emmer genome architecture and diversity elucidate wheat evolution and
1051 domestication. *Science* **357**, 93–97 (2017).
- 1052 57. S. Charkevics. *Plantae Vasculares Orientis Extremi Sovietici vol. 3, p. 101, 1988*. **3**, (1988).
- 1053 58. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-
1054 generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
- 1055 59. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation
1056 DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
- 1057 60. Van der Auwera, G. A. *et al.* From FastQ data to high-confidence variant calls: The Genome Analysis
1058 Toolkit best practices pipeline. *Curr. Protoc. Bioinforma.* **43**, 11.10.1-11.10.33 (2013).

- 1059 61. Browning, S. R. & Browning, B. L. Rapid and Accurate Haplotype Phasing and Missing-Data Inference
1060 for Whole-Genome Association Studies By Use of Localized Haplotype Clustering. *Am. J. Hum.*
1061 *Genet.* **81**, 1084–1097 (2007).
- 1062 62. Barnett, D. W., Garrison, E. K., Quinlan, A. R., Strömberg, M. P. & Marth, G. T. BamTools: a C++ API
1063 and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**, 1691–1692 (2011).
- 1064 63. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic*
1065 *Acids Res.* **32**, 1792–1797 (2004).
- 1066 64. Pritchard, J. K., Stephens, M. & Donnelly, P. Inference of population structure using multilocus
1067 genotype data. *Genetics* **155**, 945–959 (2000).
- 1068 65. Ronquist, F. *et al.* MrBayes 3.2: Efficient Bayesian Phylogenetic Inference and Model Choice Across a
1069 Large Model Space. *Syst. Biol.* **61**, 539–542 (2012).
- 1070 66. Thornton, K. libsequence: a C++ class library for evolutionary genetic analysis. *Bioinformatics* **19**,
1071 2325–2327 (2003).
- 1072 67. Pfeifer, B., Wittelsburger, U., Ramos-Onsins, S. E. & Lercher, M. J. PopGenome: An Efficient Swiss
1073 Army Knife for Population Genomic Analyses in R. *Mol. Biol. Evol.* **31**, 1929–1936 (2014).
- 1074 68. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 1075 69. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide
1076 polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w¹¹¹⁸; iso-2; iso-3.
1077 *Fly (Austin)* **6**, 80–92 (2012).
- 1078 70. Hu, T. T. *et al.* The *Arabidopsis lyrata* genome sequence and the basis of rapid genome size change.
1079 *Nat. Genet.* **43**, 476–481 (2011).
- 1080
- 1081