

1 **Parental legacy, demography, and introgression influenced the evolution of**
2 **the two subgenomes of the tetraploid *Capsella bursa-pastoris* (Brassicaceae)**

3 Dmytro Kryvokhyzha^{1,*}, Adriana Salcedo^{2,*}, Mimmi C. Eriksson^{1,4}, Tianlin Duan¹, Nilesh
4 Tawari⁵, Jun Chen¹, Maria Guerrina¹, Julia M. Kreiner², Tyler V. Kent², Ulf Lagercrantz¹, John
5 R. Stinchcombe², Sylvain Glémin^{1,3,**}, Stephen I. Wright^{2,**}, and Martin Lascoux^{1,**}

6 ¹Department of Ecology and Genetics, Program in Plant Ecology and Evolution, Evolutionary
7 Biology Center, Uppsala University, Norbyvägen 18D, 75236 Uppsala, Sweden

8 ²Department of Ecology and Evolution, University of Toronto, 25 Willcocks St., Toronto, Canada

9 ³Institut des Sciences de l'Evolution (ISEM - UMR 5554 Université de
10 Montpellier-CNRS-IRD-EPHE), Place Eugène Bataillon, 34095 Montpellier Cedex 5, France

11 ⁴Department of Biological and Environmental Sciences, University of Gothenburg, Göteborg,
12 Sweden

13 ⁵Computational and Systems Biology Group, Genome Institute of Singapore, Agency for Science,
14 Technology and Research (A*Star), Singapore

15 * Joint first authors

** Joint corresponding authors:

SG: sylvain.glemin@univ-montp2.fr

SIW: stephen.wright@utoronto.ca

16 ML: martin.lascoux@ebc.uu.se

17 **Keywords:** allopolyploidy, genome evolution, introgression, selection, genetic load, demo-
18 graphic history

19 **Running Title:** Allotetraploid genome evolution in space and time

20 ABSTRACT

21 Allopolyploidy is generally perceived as a major source of evolutionary novelties and as an
22 instantaneous way to create isolation barriers. However, we do not have a clear understanding of
23 how two subgenomes evolve and interact once they have fused in an allopolyploid species and how
24 isolated they are from their relatives. Here, we address these questions by analyzing genomic and
25 transcriptomic data of allotetraploid *Capsella bursa-pastoris* in three differentiated populations,
26 Asia, Europe and the Middle East. We phased the two subgenomes, one descended from the
27 outcrossing and highly diverse *Capsella grandiflora* (Cg) and the other one from the selfing and
28 genetically depauperate *Capsella orientalis* (Co). For each subgenome, we assessed its relationship
29 with the diploid relatives, temporal change of effective population size (N_e), signatures of positive
30 and negative selection, and gene expression patterns. Introgression between *C. bursa-pastoris* and
31 its diploid relatives was widespread and the two subgenomes were impacted differentially depending
32 on geographic region. In all three regions, N_e of the two subgenomes decreased gradually and
33 the Co subgenome accumulated more deleterious changes than Cg. Selective sweeps were more
34 common on the Cg subgenome in Europe and the Middle East, and on the Co subgenome in Asia.
35 In contrast, differences in expression were limited with the Cg subgenome slightly more expressed
36 than Co in Europe and the Middle-East. In summary, after more than 100,000 generations of
37 co-existence, the two subgenomes of *C. bursa-pastoris* still retained a strong signature of parental
38 legacy and were differentially affected by introgression and selection.

39 INTRODUCTION

40 Allopolyploidy, the origin of polyploids from two different ancestral lineages, poses serious
41 evolutionary challenges since the presence of two divergent sub-genomes may lead to perturbation
42 of meiosis, conflicts in gene expression regulation, protein-protein interactions and/or transposable
43 element suppression (Bomblies et al. 2015; Soltis et al. 2010). Whole genome duplication also
44 masks new recessive mutations thereby decreasing selection efficacy (Comai 2005; Otto and
45 Whitton 2000). This relaxation of selection, together with the strong speciation bottleneck and
46 shift to self-fertilization that often accompany polyploidy (Barringer 2007), ultimately increases
47 the frequency of deleterious mutations retained in the genome (Robertson et al. 2011; Hartfield
48 et al. 2017). All of these consequences of allopolyploidy can have a negative impact on fitness and
49 over evolutionary time may contribute to the patterns of duplicate gene loss, a process referred to as
50 diploidization (Otto and Whitton 2000; Buggs et al. 2012; Douglas et al. 2015). Yet, allopolyploid
51 lineages often not only establish and persist but may even thrive and become more successful
52 than their diploid progenitors and competitors, with larger ranges and higher competitive ability
53 (te Beest et al. 2011; Brochmann et al. 2004; Levin 2002; Pandit et al. 2006; Pandit et al. 2011;
54 Petit and Thompson 1999; Prentis et al. 2008; Prentis et al. 2008; Ramsey 2011; Soltis et al. 2014;
55 Treier et al. 2009). The success of allopolyploids is usually explained by their greater evolutionary
56 potential. Having inherited two genomes that evolved separately, and sometimes under drastically
57 different conditions, allopolyploids should have an increased genetic toolbox, assuming that the two
58 genomes do not experience severe conflicts. This greater evolutionary potential of allopolyploids
59 can be further enhanced by genomic rearrangements, alteration of gene expression and epigenetic
60 changes (Adams and Wendel 2005; Comai 2005; Doyle et al. 2008; McGrath and Lynch 2012;
61 Otto and Whitton 2000; Soltis and Soltis 1999; Soltis and Soltis 2000; Soltis and Soltis 2012;
62 Weiss-Schneeweiss et al. 2013).

63 All of these specific features come into play during the demographic history of allopolyploids.
64 Demographic processes occurring when a species extends its range, such as successive bottlenecks
65 or periods of rapid population growth in the absence of competition, are expected to have a profound

66 impact on evolutionary processes, especially in populations at the front of the expansion range.
67 Species that went through repeated bottlenecks during their range expansion are expected to have
68 reduced genetic variation and higher genetic load than more ancient central populations (Peischl
69 et al. 2016; Gilbert et al. 2017). Similarly, range expansions can also lead to contact and gene flow
70 with introgression from related species. Such gene flow can in turn shift the evolutionary path of
71 the focal species. Finally, range expansion will expose newly formed allopolyploid populations to
72 divergent selective pressures, creating the possibility of differentially exploiting duplicated genes,
73 creating asymmetrical patterns of adaptive evolution in different parts of the range.

74 In this paper, we aim to characterize the evolution of the genome of a recent allopolyploid
75 species during its range expansion. In particular, we explore whether the two subgenomes have
76 similar or different evolutionary trajectories in term of hybridization, selection and gene expression.
77 The widespread allopolyploid *C. bursa-pastoris* is a promising system for studying the evolution
78 of polyploidy, with available information on its two progenitor diploid species and their current
79 distribution. *C. bursa-pastoris*, a selfing species, originated from the hybridization of the *Capsella*
80 *orientalis* and *Capsella grandiflora* / *rubella* lineages some 100-300 kya (Douglas et al. 2015).
81 *C. orientalis* is a genetically depauperate selfer occurring across the steppes of Central Asia and
82 Eastern Europe. In contrast, *C. grandiflora* is an extremely genetically diverse obligate outcrosser
83 which is primarily confined to a tiny distribution range in the mountains of Northern Greece and
84 Albania. The fourth relative, *C. rubella*, a selfer recently derived from *C. grandiflora*, occurs around
85 the Mediterranean Sea (Fig. 1A). There is evidence for unidirectional gene flow from *C. rubella*
86 to *C. bursa-pastoris* (Slotte et al. 2008a). Among all *Capsella* species, only *C. bursa-pastoris*
87 has a worldwide distribution (Hurka et al. 2012), some of which might be due to extremely recent
88 colonization and associated with human population movements (Cornille et al. 2016). A recent
89 study reveals that in Eurasia, *C. bursa-pastoris* is divided into three genetic clusters - Middle East,
90 Europe and Asia - with low gene flow among them and strong differentiation both at the nucleotide
91 and gene expression levels (Cornille et al. 2016; Kryvokhyzha et al. 2016). Reconstruction of the
92 colonization history using unphased genomic data suggested that *C. bursa-pastoris* spread from

93 the Middle East towards Europe and then into Eastern Asia. This colonization history resulted in
94 a typical reduction of nucleotide diversity with the lowest diversity being in the most distant Asian
95 population (Cornille et al. 2016).

96 How the two distinct non-recombining subgenomes of *C. bursa-pastoris* contributed to its rapid
97 population expansion and how they were in return affected by it, remains unclear. Previous studies
98 either ignored the population history of *C. bursa-pastoris* or failed to consider the two subgenomes
99 separately. In a recent study that does not consider the population demographic history within
100 *C. bursa-pastoris*, Douglas et al. (2015) concluded that there is no strong sign of diploidization in
101 *C. bursa-pastoris* and most of its variation is the result of the legacy from the parental lineages
102 with some relaxation of purifying selection caused by both the transition to self-fertilization and the
103 greater masking of deleterious mutations. Kryvokhyzha et al. (2016) considered population history
104 but did not separate the two subgenomes, and showed that variation in gene expression among
105 Asian, European and Middle Eastern accessions strongly reflects the population history with most
106 of the differences among populations explained by genetic drift. We extend these previous studies
107 by analyzing the genome-wide expression and polymorphism patterns of the two subgenomes of
108 *C. bursa-pastoris* in 31 accessions sampled across its natural range in Eurasia. We demonstrate
109 that the two subgenomes follow distinct evolutionary trajectories in different populations and that
110 these trajectories are influenced by both range expansion and introgression from relatives. Our
111 study illustrates the need to account for demographic and ecological differences among populations
112 when studying the evolution of subgenomes of allopolyploid species.

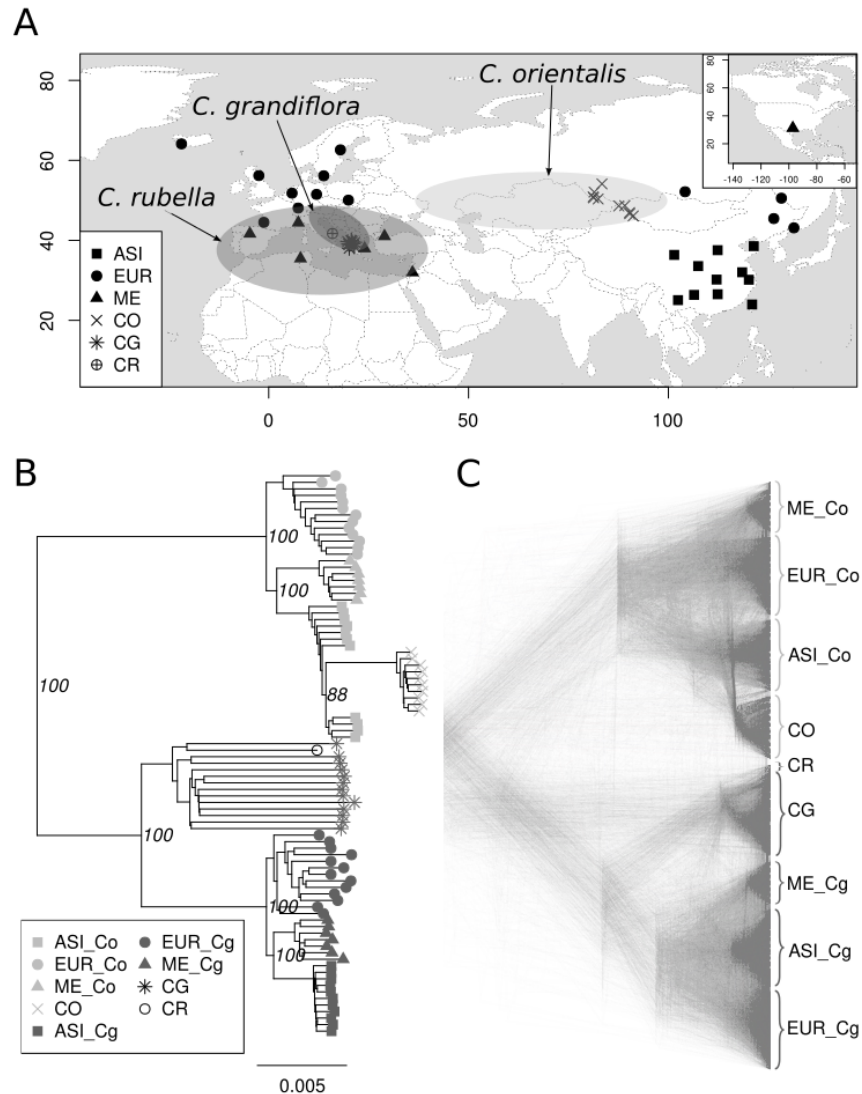


Fig. 1. **A.** Approximate distribution ranges of *C. orientalis*, *C. grandiflora*, and *C. rubella* and sampling locations of *C. bursa-pastoris*. *C. bursa-pastoris* has a worldwide distribution, so its distribution range is not specifically depicted. ASI, EUR ME, CO, CG, CR indicate Asian, European and Middle Eastern populations of *C. bursa-pastoris*, *C. orientalis*, *C. grandiflora*, and *C. rubella*, respectively. The map is modified from Hurka et al. (2012). **B.** Whole genome NJ tree showing the absolute divergence between different populations of *C. bursa-pastoris* at the level of subgenomes. The Co and Cg subgenomes are marked with corresponding names. The bootstrap support based on 100 replicates is shown only for the major clades. The root *N. paniculata* is not shown. **C.** Density tree visualizing of 1002 NJ trees reconstructed with 100 Kb sliding windows.

113 RESULTS

114 Phasing subgenomes

115 The disomic inheritance of *C. bursa-pastoris* allowed us to successfully phase most of the
116 heterozygous sites in the 31 samples analyzed in this study (Fig. 1A, Table S1). Out of 7.1
117 million high confidence SNPs, our phasing procedure produced an alignment of 5.4 million phased
118 polymorphic sites across the 31 accessions of *C. bursa-pastoris*. Scaling these phased SNPs to
119 the whole genome resulted in the alignment of 80.6 Mb that had the same level of heterozygosity
120 as the unphased data. The alignment of these whole genome sequences of *C. bursa-pastoris*
121 with 13 sequences of *C. grandiflora*, 10 sequences of *C. orientalis*, one sequence of *C. rubella*
122 (the reference), and one sequence of *N. paniculata* used here as an outgroup, yielded 13 million
123 polymorphic sites that we used in all analyses. The information for each accession is provided in
124 the Supporting Information.

125 To assess the quality of the phasing results, we constructed a phylogeny from the phased data.
126 The separation of the two subgenomes was strongly supported in the reconstructed whole genome
127 tree (Fig. 1B). The tree consisted of two highly supported (100% bootstrap) major clades grouping
128 *C. grandiflora* and the *C. grandiflora* / *rubella* lineage descended subgenome of *C. bursa-pastoris*
129 (hereafter the Cg subgenome), on the one hand, and *C. orientalis* and the *C. orientalis* lineage
130 descended subgenome of *C. bursa-pastoris* (hereafter the Co subgenome), on the other hand. We
131 also analyzed phylogenetic signals at a finer genomic scale using a sliding window approach with
132 100-kb window size (Fig. 1C). Exclusive monophyly of *C. orientalis* with the Co subgenome and
133 *C. grandiflora* / *rubella* with Cg subgenome was detected in 95% and 83% of trees, respectively
134 (Fig. S1).

135 Polymorphism and population structure of the two subgenomes

136 For both subgenomes the three *C. bursa-pastoris* populations, Asia (ASI), Europe (EUR)
137 and Middle East (ME), constituted well-defined phylogenetic clusters (Fig. 1B,C). However, the
138 relationships of each subgenome with its parental species differed. The Cg subgenome formed a
139 monophyletic clade with *C. grandiflora* at its base. In contrast, the Co subgenome was paraphyletic

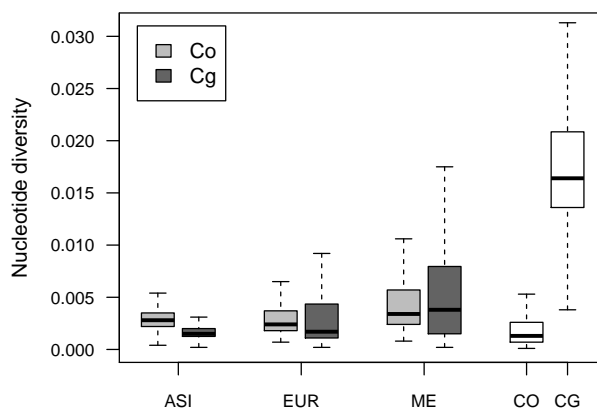


Fig. 2. Variation in nucleotide diversity (π) between populations of *C. bursa-pastoris* and parental species. This boxplot shows π estimated along the genome using 100 Kb sliding windows. Co and Cg indicate *C. orientalis* and *C. grandiflora / rubella* descendant subgenomes, respectively. ASI, EUR ME, CO and CG correspond to Asian, European and Middle Eastern populations of *C. bursa-pastoris*, *C. orientalis*, and *C. grandiflora*, respectively.

140 with *C. orientalis* clustering within the ASI group instead of being outside of all *C. bursa-pastoris*
141 Co subgenomes. This clustering was unexpected and suggested potential gene flow between the
142 ASI group and *C. orientalis* or multiple origins of the Co subgenome. Nucleotide diversity was
143 higher on the Cg subgenome than on the Co subgenome for both EUR and ME (Fig. 2, Table S2),
144 though the difference was significant only for EUR (p-values: 0.005 and 0.154 for EUR and ME
145 respectively). The opposite pattern was observed for ASI (Fig. 2): there the nucleotide diversity
146 in the Co subgenome was significantly higher than in the Cg subgenome (p-value < 0.0001).
147 Interestingly, the diversity of the Co subgenome in all populations was significantly higher than the
148 diversity of its parental species, *C. orientalis* (p-value < 0.0001).

149 **Temporal change in effective population size**

150 To reconstruct the changes in effective population size (N_e) over time in the three *C. bursa-*
151 *pastoris* populations and the two ancestral species, we used a pairwise sequentially Markovian
152 coalescent model (PSMC). First, we reconstructed the demographic histories of *C. orientalis* and
153 *C. grandiflora* (Fig. 3). In *C. grandiflora*, N_e was mostly constant with some slight decrease in
154 the recent past, but the N_e of *C. orientalis* decreased continuously. In *C. bursa-pastoris*, despite

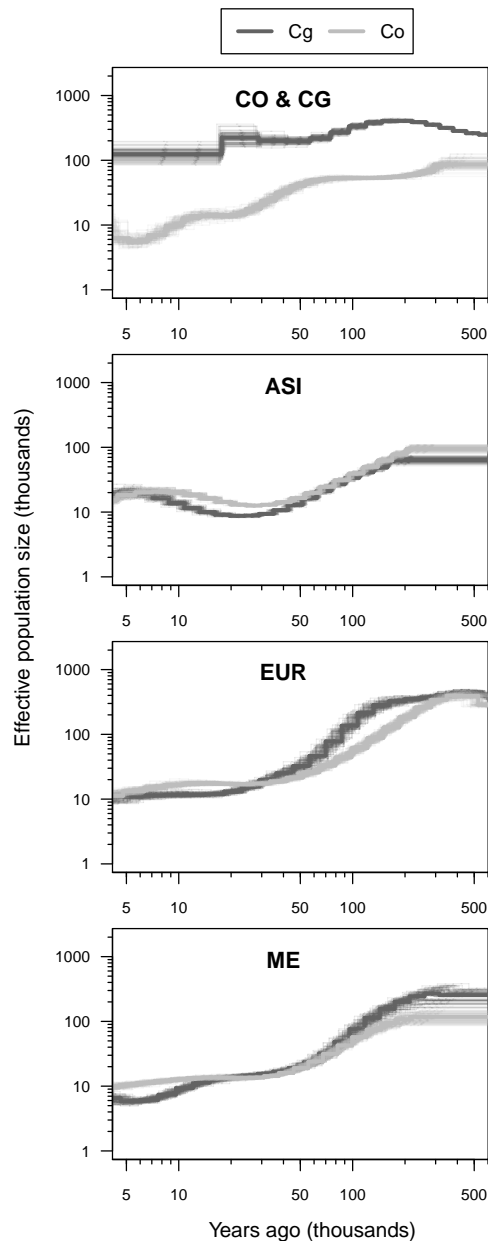


Fig. 3. Population size histories of *C. bursa-pastoris* and its parental species. Effective population sizes were inferred with PSMC using whole-genome sequences from a pair of haplotypes per population (thick lines) and 100 bootstrap replicates (thin lines). The estimates for different pairs were similar and shown in the Supp. (Fig. S12). Co and Cg specify subgenomes of *C. bursa-pastoris* and corresponding parental species in the CO & CG plot. ASI, EUR, ME, CO & CG indicate Asian, European and Middle Eastern populations of *C. bursa-pastoris*, and *C. orientalis* and *C. grandiflora*, respectively. The axis are in log scale and the most recent times where PSMC is less reliable were excluded.

155 a simultaneous rapid range expansion, N_e of EUR and ME populations also gradually decreased
156 starting from around 100-200 kya. The ASI population showed a similar pattern but with population
157 size recovery in the range 5-10 kya and a subsequent decrease to the same N_e as in EUR and ME.
158 The N_e patterns of the two subgenomes were similar within each subpopulation. Overall, the N_e
159 history of *C. bursa-pastoris* was most similar to that of its selfing ancestor, *C. orientalis*. We also
160 verified these PSMC results with SMC++, which can consider more than two haploid genomes
161 and incorporates linkage disequilibrium (LD) in coalescent hidden Markov models (Terhorst et al.
162 2017). The general trend was globally the same but the recent decline of *C. orientalis* was sharper
163 and fluctuations in N_e more pronounced (Fig. S2). In summary, the overall pattern of N_e change
164 over time was mostly the same between the two subgenomes and between the three populations of
165 *C. bursa-pastoris* and it was largely similar to the pattern observed for the diploid selfer *C. orientalis*.

166 **Relationship of the *C. bursa-pastoris* subgenomes with their parental species**

167 To quantify the relationships between populations of *C. bursa-pastoris* and the two parental
168 species, we applied a topology weighting method that calculates the contribution of each individual
169 group topology to a full tree (Martin and Van Belleghem 2017). We looked at the topologies
170 joining each subgenome of *C. bursa-pastoris* and a corresponding parental lineage. There are 15
171 possible topologies for three populations of *C. bursa-pastoris*, a parental species, and the root. We
172 grouped these topologies into five main groups: species trees - topologies that place a parental
173 lineage as a basal branch to *C. bursa-pastoris*; three groups that join one of the populations of
174 *C. bursa-pastoris* with a parental lineage and potentially signifies admixture; and all other trees that
175 place a parental lineage within *C. bursa-pastoris* but do not relate it with a particular population of
176 *C. bursa-pastoris* (Fig. 4A).

177 These topology weightings varied along the subgenomes and illustrate distinct patterns between
178 the two subgenomes (Fig. 4B). In the Co subgenome, the largest average weighting was for the
179 topology grouping the Co subgenome of the ASI population of *C. bursa-pastoris* with *C. orientalis*
180 (Fig. 4C), and the species topology had the second largest average weighting. The difference
181 between the average weighting in these two topology groups was statistically significant (Table

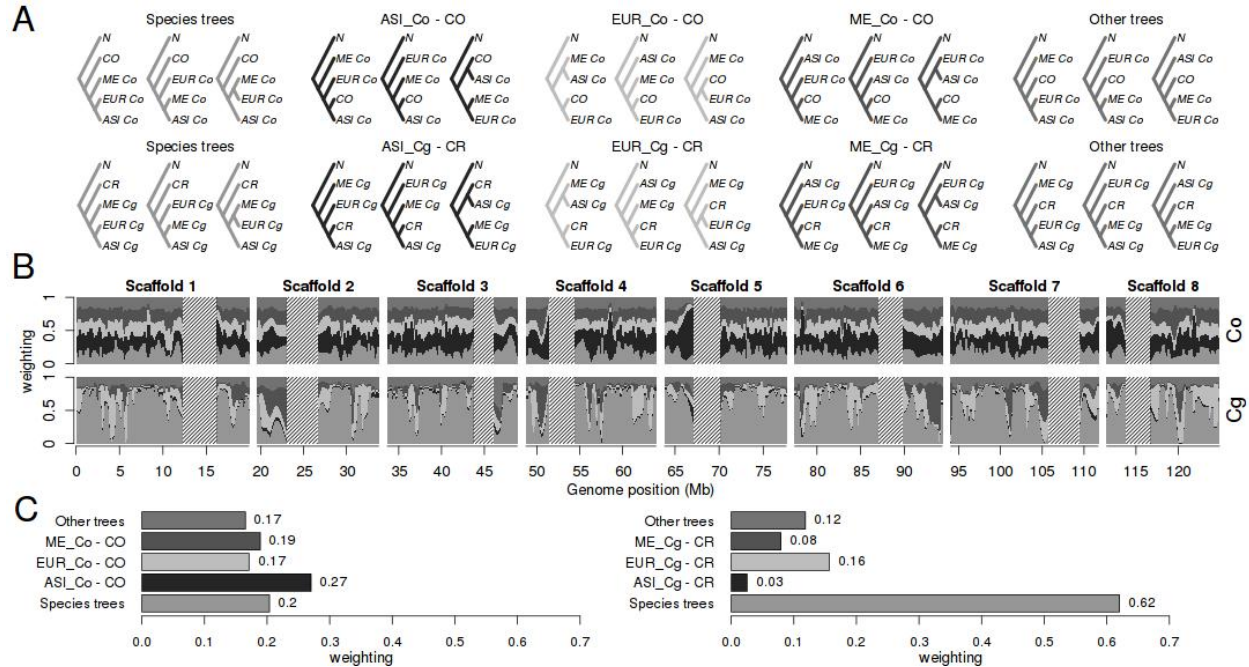


Fig. 4. Topology weighting of the three populations of *C. bursa-pastoris* and parental species.

A. Fifteen possible rooted topologies for the three groups of *C. bursa-pastoris* in one subgenome and the corresponding parental species. The topologies are grouped into five main groups. Co and Cg indicate *C. orientalis* and *C. grandiflora / rubella* descendant subgenomes, respectively. ASI, EUR ME, CO, CR, N indicate Asian, European and Middle Eastern populations of *C. bursa-pastoris*, *C. orientalis*, *C. rubella*, and *N. paniculata*, respectively. **B.** Topology weightings for 100 SNP windows plotted along 8 main scaffolds with loess smoothing (span = 1Mb). The tentative centromeric regions are shaded and only eight major scaffolds are shown. **C.** Average weighting for the five main topology groups. The topology groups are in the same order (left-right and bottom-up) and colors in all plots.

182 S3). In contrast, the species topologies weighting dominated in the Cg subgenome, regardless if
 183 *C. rubella* or *C. grandiflora* were used as a parental lineage for the Cg subgenome (Fig. 4C, Fig.
 184 S3, Table S4, S5). The topology uniting the Cg subgenome of the EUR population with *C. rubella*
 185 was the largest among the topologies indicating admixture in the Cg subgenome (Fig. 4C). Thus,
 186 the two subgenomes differed substantially in the pattern of topology weighting and there were signs
 187 of a potential admixture of EUR and ASI with *C. rubella* and *C. orientalis*, respectively.

188 **Gene flow between *C. bursa-pastoris* and its relatives**

189 *Genomic inferences*

190 The phylogenetic grouping of *C. orientalis* with the Asian Co subgenome, together with topology
191 weighting results and the relatively elevated nucleotide diversity in this subgenome, suggested the
192 presence of gene flow between *C. orientalis* and *C. bursa-pastoris* in the ASI population. To test
193 this hypothesis, and at the same time to check for possibilities of gene exchange between *C. bursa-*
194 *pastoris* and other *Capsella* species, we conducted two complementary tests of introgression.

195 We first used the ABBA-BABA test, a coalescent based method that relies on the assumption
196 that alleles under incomplete lineage sorting are expected to be equally frequent in two descendant
197 populations in the absence of introgression between any of them and a third population that diverged
198 earlier on from the same common ancestor (Green et al. 2010; Durand et al. 2011). The deviation
199 from equal frequency is measured with the *D*-statistics, which ranges between 0 and 1, with 0
200 indicating no gene flow and 1 meaning complete admixture. The ABBA-BABA test also provides
201 an estimate of the fraction of the genome that is admixed by comparing the observed difference in
202 ABBA-BABA with the difference expected under a scenario of complete admixture (*f*-statistics).
203 We estimated *D* and *f* for triplets including one diploid species and two populations of *C. bursa-*
204 *pastoris* represented by the most related subgenome to that species (Table 1). *N. paniculata* was
205 the outgroup in all tests. The *D*-statistics were significantly different from 0 in most of the tests,
206 so we considered all three combinations per test group (see Table 1) to determine the pairs that
207 were the most likely to be admixed. The largest fraction of admixture was identified for the pair of
208 the ASI Co subgenome and *C. orientalis* with an estimate of *f* indicating that at least 14% of the
209 ASI Co subgenome is admixed. The second largest proportion of admixture was detected between
210 *C. rubella* and the EUR Cg subgenome with *f* estimate of at least 8%. The estimates for tests
211 with *C. grandiflora* were the smallest but similar to those obtained for *C. rubella*. The latter may
212 reflect the strong genetic similarity between these two species rather than real gene flow between
213 *C. grandiflora* and *C. bursa-pastoris* which, based on crosses (see below), seems unlikely. Finally,
214 it should be pointed out that given that evidence for *C. bursa-pastoris* monophyly is weak, it is also

TABLE 1. Results of the ABBA-BABA tests assessing admixture between *C. bursa-pastoris* and *C. orientalis*, *C. grandiflora* and *C. rubella*.

P ₁	P ₂	P ₃	$D \pm \text{error}$	Z-score	P-value	$f \pm \text{error} (\%)$
EUR_Co	ASI_Co	CO	0.29 ± 0.03	8.62	<0.0001	22.9 ± 2.5
ME_Co	ASI_Co	CO	0.18 ± 0.04	4.80	<0.0001	14.0 ± 2.8
EUR_Co	ME_Co	CO	0.17 ± 0.03	5.70	<0.0001	11.7 ± 2.4
ASI_Cg	EUR_Cg	CG	0.19 ± 0.01	15.45	<0.0001	19.8 ± 2.2
ASI_Cg	ME_Cg	CG	0.17 ± 0.02	10.14	<0.0001	12.6 ± 2.0
ME_Cg	EUR_Cg	CG	0.06 ± 0.01	5.14	<0.0001	6.1 ± 1.2
ASI_Cg	EUR_Cg	CR	0.61 ± 0.02	26.74	<0.0001	20.1 ± 2.1
ASI_Cg	ME_Cg	CR	0.49 ± 0.03	14.55	<0.0001	10.6 ± 1.6
ME_Cg	EUR_Cg	CR	0.26 ± 0.05	4.84	<0.0001	7.9 ± 1.7

P₁, P₂, and P₃ refer to the three populations used in the ABBA-BABA tests. A significantly positive D indicates admixture between P₂ and P₃. f provides an estimate of the fraction of introgression. Z-score and P-value were estimated with the block jack-knife method. The error term corresponds to a standard error. ASI, EUR and ME are the three populations of *C. bursa-pastoris* with _Co and _Cg indicating different subgenomes. CO and CG stand for *C. orientalis* and *C. grandiflora*, respectively. Every test group is separated by a horizontal line.

215 possible that the signals of introgression from the parental species into *C. bursa-pastoris* that we
 216 are detecting here actually reflects introgression from an independently-arisen *C. bursa-pastoris*
 217 into either Co or Cg subgenomes.

218 We then used HAPMIX, a haplotype-based method, which should allow us to capture both
 219 large-scale and fine-scale admixture, and enables an absolute estimate of the proportion of the
 220 genome that was admixed. For the analysis of the Cg subgenome of *C. bursa-pastoris*, the highest
 221 levels of introgression were found consistently across regions to be from the diploid *C. rubella*. In
 222 Europe, 18% of SNPs genome-wide showed introgression from *C. rubella*, followed by 11% in the
 223 Middle East, and just 2% in Asia (Table S6, Fig. S4A). All three populations also showed signs of
 224 *C. grandiflora* introgression but to a reduced extent compared to *C. rubella* (7% in Europe, 6% in
 225 the Middle East, 0.2% in Asia). *C. rubella* functionally represents a haplotype of *C. grandiflora*,
 226 and as noted above, we expect difficulties in discerning between the two, suggesting that much of

227 the signal of introgression from *C. grandiflora* could in fact be due to *C. rubella* introgression. Of
228 the regions putatively introgressed from *C. grandiflora*, 78%-96% of sites called as introgressed
229 overlapped with those from *C. rubella*, none of which occurred in unique regions for *C. grandiflora*.
230 Because of this, and in combination with the reduced genome-wide probability of introgression
231 from the diploid *C. grandiflora* compared to *C. rubella* (e.g. 0.11 compared to 0.24 in Europe),
232 we argue that the signals of introgression from the diploid *C. grandiflora* were likely an artifact
233 of its similarity with the regions of *C. rubella* introgression. These findings in accord with the
234 ABBA-BABA results imply that the Cg subgenome has experienced significant introgression from
235 *C. rubella* in Europe, and to a lesser extent in the Middle East.

236 For the analysis of the Co subgenome of *C. bursa-pastoris*, signals of introgression from the
237 diploid *C. orientalis* were present in all three populations. In the ME population, 18-21% of
238 SNPs showed signals of *C. orientalis* introgression (Table S6, Fig. S4B). Using the Middle East
239 population for the analysis of the Co subgenomes of EUR and ASI, since it was the least introgressed
240 in the HAPMIX results, yielded 15% *C. orientalis* introgression in Asia, and 14% in Europe. These
241 findings suggest introgression of the diploid *C. orientalis* into the Co subgenome across all three
242 geographic regions. Assuming these levels of admixture accurately reflect reality, we do not have
243 a non-admixed reference population to use for Hapmix, and are thus violating a key assumption
244 of the method. Hapmix inferences for the Co subgenome should therefore be taken with caution
245 but we note that the results for ASI and ME are generally congruent with the admixture pattern
246 obtained with ABBA-BABA.

247 *Crosses*

248 To assess further the plausibility of these inferences, we crossed individuals from the three
249 populations of *C. bursa-pastoris* with their three diploid relatives to test for the presence of
250 reproductive barriers. Regardless of the direction of the crosses, all crosses between *C. rubella* and
251 the three populations of *C. bursa-pastoris* produced viable seeds. Importantly, crosses between
252 *C. rubella* and EUR produced relatively more seeds and had smaller abortion rate than crosses with
253 the other two populations of *C. bursa-pastoris*. Crosses between *C. orientalis* and *C. bursa-pastoris*

254 mostly failed or led to aborted seeds, with the exception of one Russian accession of *C. orientalis*
255 (PAR-RUS) that produced normally shaped seeds regardless if it served as a mother plant or as
256 a pollen donor. In the latter case, there was a tendency towards higher seed number and smaller
257 abortion rate for the ASI population than for EUR and ME. The crosses involving *C. grandiflora*
258 mostly failed and the abortion rate approached 100%. Details on these crosses are provided in
259 the Supplementary Information. Although the number of crosses was limited and did not provide
260 enough power for proper statistical tests, they nonetheless are sufficient to show that the admixture
261 detected at the molecular level is not completely restricted by reproductive barriers.

262 In summary, admixture between *C. bursa-pastoris* and *C. orientalis* in Asia, and between
263 *C. bursa-pastoris* and *C. rubella* in Europe is supported by molecular data, even though some of
264 the observed patterns could also be attributed to shared ancestry. Artificial crosses indicate that
265 these inferences are credible.

266 Selection and gene expression

267 Deleterious mutations

268 We first estimated the nucleotide diversity at 0-fold (π_0) and 4-fold (π_4) degenerate sites and
269 then the ratio π_0/π_4 as a measure of purifying selection, low values of π_0/π_4 indicating higher
270 purifying selection (Chen et al. 2017). As expected, π_0/π_4 was much lower in *C. grandiflora* than
271 in *C. orientalis*. In *C. bursa-pastoris*, purifying selection was more efficient in the Cg subgenome
272 than in the Co subgenome in both EUR and ME. However, the opposite was observed in the ASI
273 population. For both subgenomes, the ASI population had the highest value of π_0/π_4 even if
274 compared with *C. orientalis* (Fig. S5).

275 We then investigated the differences in deleterious mutations among subgenomes and popula-
276 tions by classifying nonsynonymous mutations with the SIFT4G algorithm that uses site conser-
277 vation across species to predict the selective effect of nonsynonymous changes (Vaser et al. 2016).
278 In order to control for possible biases due to the unequal genetic distance between the different
279 genomes and *C. rubella*, we used both *C. rubella* and *A. thaliana* SIFT4G annotation databases.
280 Because we are interested in the number of deleterious mutations that accumulated after speciation

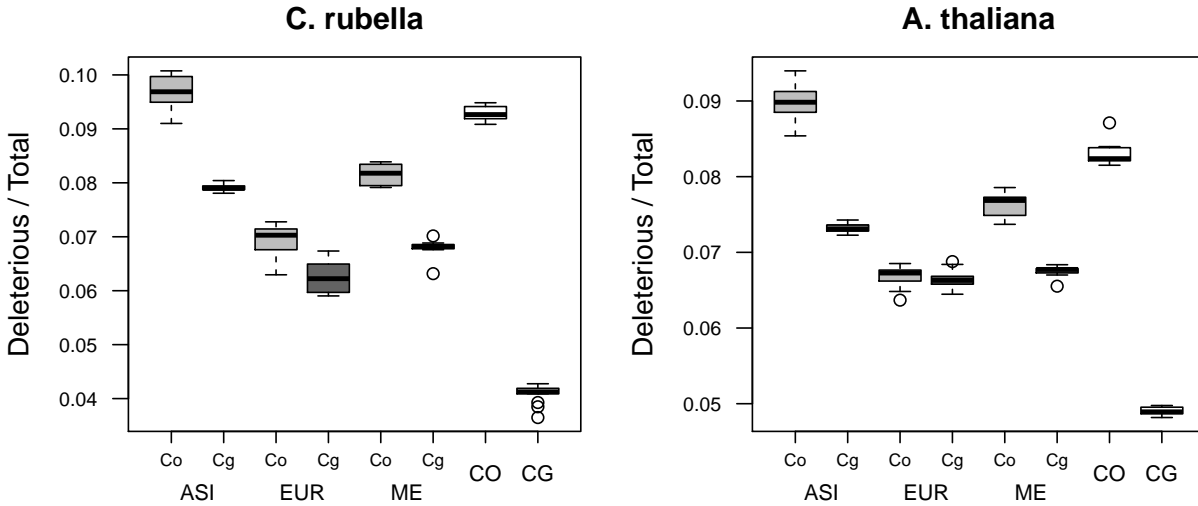


Fig. 5. Genetic load in the subgenomes of *C. bursa-pastoris* and its parental species. The proportion of deleterious nonsynonymous changes was estimated with SIFT4G on derived alleles, i.e. alleles accumulated after the speciation of *C. bursa-pastoris*. The left plot shows the results obtained with *C. rubella* database and the right plot those obtained with *A. thaliana* database. Co and Cg are the two subgenomes of *C. bursa-pastoris*. ASI, EUR, ME, CO, CG indicate Asian, European and Middle Eastern populations of *C. bursa-pastoris*, and parental species *C. orientalis* and *C. grandiflora*, respectively.

281 of *C. bursa-pastoris*, we polarized the mutations of all three species with the reconstructed ancestral
 282 sequences of the common ancestors (see Material and Methods).

283 Regardless of the SIFT4G database used (*C. rubella* or *A. thaliana*), the proportion of deleterious
 284 nonsynonymous sites among derived mutations was always significantly higher in *C. orientalis* and
 285 the Co subgenomes than in *C. grandiflora* and the Cg subgenomes (Fig. 5B, Table S7, S8). Within
 286 *C. bursa-pastoris*, the proportion of deleterious mutations depended on the population considered
 287 with the highest value in the ASI population and the smallest in EUR. It is also noteworthy that
 288 the proportion of deleterious nonsynonymous sites of the Co subgenome in EUR and ME is
 289 significantly smaller than that of *C. orientalis* suggesting that a higher effective population size
 290 in the Co subgenome than in its ancestor led to more efficient purifying selection in these two
 291 populations. On the other hand, the proportion of deleterious nonsynonymous sites in the Asian
 292 Co subgenome was larger than in *C. orientalis*, but this difference was only significant for the *A.*
 293 *thaliana* database. The Cg subgenome also had a significantly higher proportion of deleterious sites

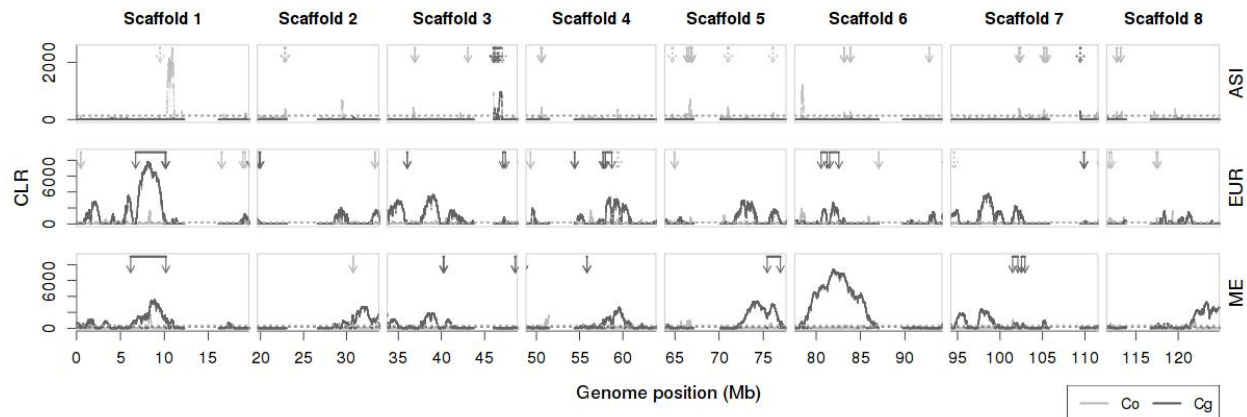


Fig. 6. Selective sweep differences between populations of *C. bursa-pastoris*. Selective sweeps are detected with the composite likelihood ratio statistics (CLR) along the Co and Cg genomic subgenomes in Asian (ASI), European (EUR) and Middle Eastern (ME) populations. The dashed line indicates the 0.01 significance level defined with data simulated under a standard neutral model. Solid arrows point to the location of introgression and dashed arrows show the location of genomic conversion. Pericentromeric regions are removed. Only eight major scaffolds are shown.

294 in ASI than in EUR and ME in all comparisons. In conclusion, the proportion of deleterious sites in
 295 the two subgenomes of extant *C. bursa-pastoris* still reflected the differences between the parental
 296 species and the efficacy of purifying selection in the different *C. bursa-pastoris* subpopulations was
 297 associated to their synonymous nucleotide diversity or, equivalently, to their effective population
 298 size.

299 *Selective sweeps*

300 The three populations of *C. bursa-pastoris* also differ in patterns of positive selection. Overall,
 301 the number of sweeps in Co and Cg subgenomes were independent ($\chi^2 = 89.386$, p-value < 0.001).
 302 Selective sweeps were more significant on the Cg subgenome in EUR and ME than on the Co
 303 subgenome, whereas in the ASI population, the opposite was true (Fig. 6). The regions harboring
 304 significant sweeps were also larger on the Cg subgenome than on the Co subgenome in EUR and
 305 ME (total length 42 Mb, 50 Mb vs 9 Mb, 3 Mb), whereas in Asia the sweep regions were larger on
 306 Co than on Cg (total length 4 Mb vs 830 Kb). Although the locations of the Cg sweeps in EUR and
 307 ME largely overlap, the patterns differed between the two populations. For example, the strongest
 308 sweep in EUR was located on scaffold 1, whereas the strongest sweep in ME was on scaffold 6. In

309 addition, EUR had many sweeps in Co subgenome (109 in EUR_Cg, 128 in EUR_Co), but they all
310 were small and hardly above the significance threshold (Fig. 6). In the ME population, the sweeps
311 in the Cg subgenome were prevailing both in size and numbers (101 in ME_Cg, 22 in ME_Co).
312 The ASI population differed strongly from both EUR and ME not only because most of its sweep
313 signals were on the Co subgenome but also because these sweeps regions were narrower and less
314 pronounced (Fig. 6). Thus, all three populations of *C. bursa-pastoris* were distinct in their selective
315 sweeps patterns with the Asian population being the one least affected by positive selection.

316 Given the presence of gene flow between *Capsella* species, we also checked if any of the
317 detected selective sweeps could be due to introgression. We compared genetic distances for
318 every sweep region among the three *Capsella* groups and the parental species. A sweep region was
319 considered to have resulted from introgression if its genetic distance was closer to the corresponding
320 parental species than to any other *C. bursa-pastoris* sequence. This comparison also allowed us to
321 identify regions of possible gene conversion if the genetic distance was smallest between the two
322 subgenomes. The distance between individual sweep regions revealed that all the CLR outliers
323 in the ASI Cg subgenome were genetically closer to the Asian Co subgenome than to other
324 Cg subgenomes of *C. bursa-pastoris* (Fig. S6), and thus they probably were the result of gene
325 conversion. The distance analysis of sweep regions in the ASI Co subgenome revealed 9 regions
326 of gene conversion (total length 505 Kb) and 17 regions of introgression from *C. orientalis* (total
327 length 1.3 Mb) (Fig. 6, S6). On the other hand, we found 9 regions of potential introgression from
328 *C. orientalis* to the EUR Co subgenome (total length 945 Kb), and one to the ME Co subgenome
329 (length 40 Kb). There were also 10 introgression regions between *C. rubella* and the EUR Cg
330 subgenome (total length 6.5 Mb), and 7 introgression regions between *C. rubella* and the ME Cg
331 subgenome (total length 6.7 Mb) (Fig. S6). We did not observe any sign of gene conversion in the
332 ME population and in the EUR Cg subgenome, but we found 2 regions of gene conversion from the
333 Cg to the Co subgenome in EUR (total length 154 Kb). The regions of gene conversion showed
334 reduced heterozygosity in both the phased and unphased data (Fig. S7), suggesting they were not
335 an artifact of phasing. Thus, some of the sweep signals could be solely due to gene conversion and

336 introgression, but we cannot rule out subsequent selection of these conversion and introgression
337 regions.

338 *Homeologue-specific expression*

339 The relative expression of the two subgenomes, or homeologue-specific expression (HSE), can
340 provide additional information on the evolution of the two subgenomes in different populations
341 of *C. bursa-pastoris*. In particular, biased adaptation towards one subgenome may select for
342 decreased expression of the other subgenome. Given selective favor for different subgenomes in
343 different populations, one would also expect the Cg subgenome to be over-expressed in EUR and
344 ME, and the Co subgenome in Asia.

345 To assess HSE, we analyzed the RNA-Seq data of 24 accessions representing all three popula-
346 tions of *C. bursa-pastoris* in a hierarchical Bayesian model that integrates information from both
347 RNA and DNA data (Skelly et al. 2011). Overall, in agreement with Douglas et al. (2015), one
348 subgenome did not dominate the other in the 24 accessions considered together, though a few genes
349 demonstrated a slight expression shift toward the Cg subgenome. On average, we assessed HSE
350 in 13,589 genes per accession (range 12,808-15,340) and 18% of them showed significant HSE
351 (posterior probability of HSE > 0.99). The expression ratios between subgenomes (defined here as
352 Co / Total) across all assayed genes in the DNA data were close to equal (mean = 0.495). Thus,
353 there was no strong mapping bias.

354 Among populations, HSE varied considerably. The mean expression ratios for all genes were
355 0.494, 0.489, and 0.489 in the ASI, EUR, and ME accessions, respectively, and these mean ratios
356 for genes with significant HSE were 0.487, 0.465, 0.468. The difference in mean ratio between EUR
357 and ME was not significant, but both EUR and ME were significantly different from ASI (Table S9).
358 In addition, the distribution of expression ratios between the two subgenomes was right-skewed
359 in EUR and ME, whereas in the ASI population, the distribution was more symmetrical (Fig. 7).
360 The difference between the populations was particularly evident in the grand mean values (Fig. 7).
361 Thus, the shift towards higher expression of the Cg subgenome was more prominent in Europe and
362 the Middle East than in Asia.

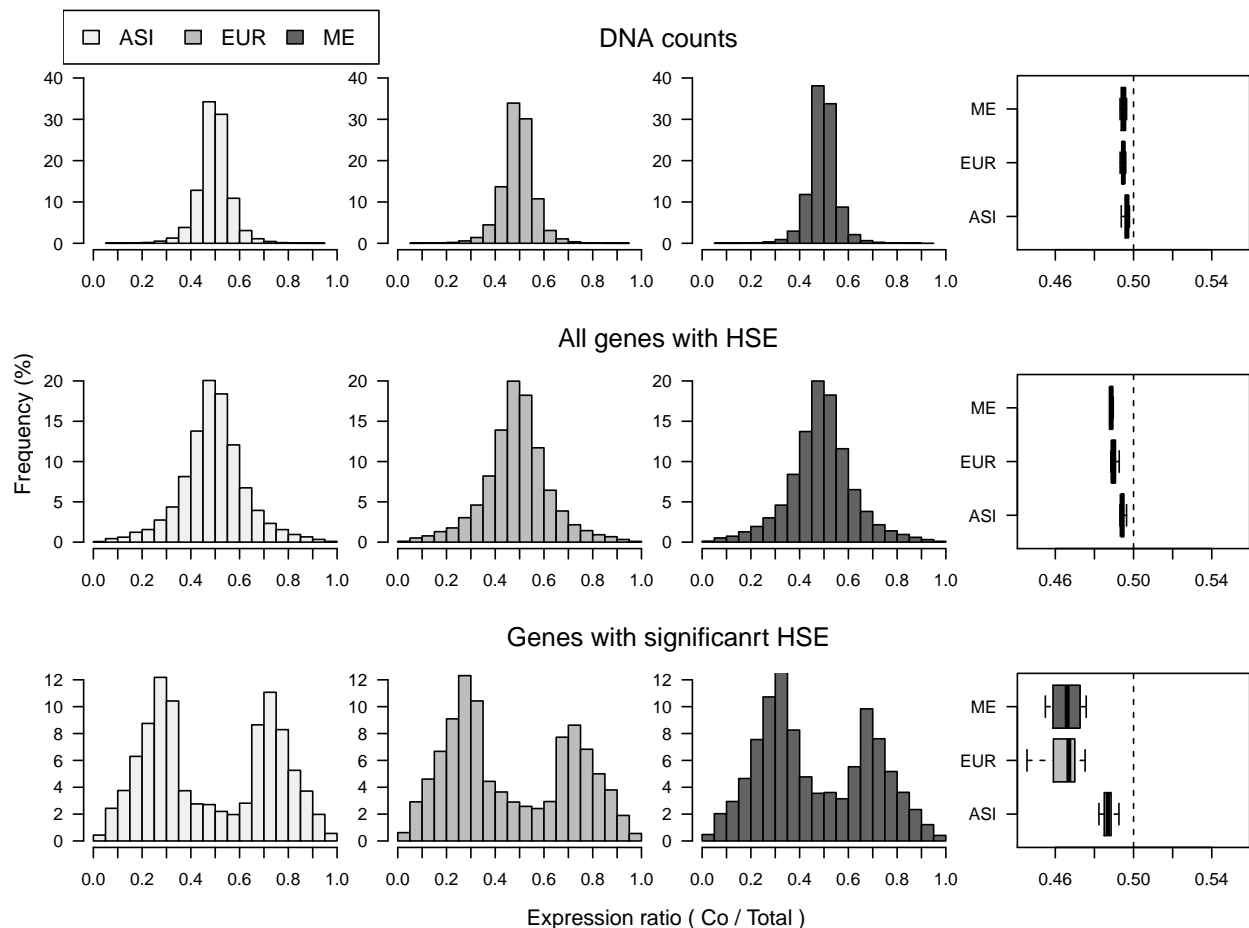


Fig. 7. Distributions of expression ratios between the two subgenomes of *C. bursa-pastoris*. The subgenome specific expression (HSE) is estimated by the fraction the Co subgenome relative to total expression level. The upper part presents the distributions for DNA counts, the middle plots show the expression distribution for all assayed gene and the lower plot shows only the distribution for genes with significant expression of one of the subgenomes. The histograms present the distribution of allelic ratio, whereas the boxplots summarize these results with the grand mean for every sample. ASI, EUR, and ME indicate Asian, European and Middle Eastern populations, respectively.

363 Expression levels were also noticeably distinct in the three populations. We analyzed the
 364 pairwise correlations in the HSE between all 24 samples, to check if the direction of the expression
 365 shift in every gene was similar within and between populations. Overall, the levels of expression
 366 of the genes with significant HSE were positively correlated between samples (mean Pearson's
 367 $r=0.81$), but correlations were distinctly stronger for samples from the same population (mean
 368 Pearson's $r = 0.91$) than for samples from different populations (mean Pearson's $r = 0.75$) (Fig.S8).

369 This pattern was also similar for the pairwise correlations across all assayed genes but the correlation
370 coefficients were smaller (mean Pearson's r for all comparisons 0.56, within populations 0.72 and
371 between populations 0.47) (Fig. S9). Thus, globally expression levels co-varied, but expression
372 levels were more similar within populations than between them.

373 The genes with significant HSE were roughly the same in all three populations. We considered
374 that a gene showed a population-specific HSE if it had a significant HSE in at least 9/11, 7/9, and 2/4
375 samples for ASI, EUR, and ME, respectively. With these criteria, we found that there was almost
376 60% overlap in gene names showing significant HSE in pairwise comparisons between the three
377 populations. Also, selective sweep regions were not over-represented by genes with significant
378 HSE (Fisher's Exact Test, p-values 0.99, 0.91, 0.47 for ASI, EUR, ME, respectively.).

379 Additionally, we were interested in testing whether the results of the differential gene expression
380 analysis of phased data between these three populations differed from the results obtained by
381 [Kryvokhyzha et al. \(2016\)](#) on unphased data. Many genes differentiated the ASI, EUR and ME
382 populations in [Kryvokhyzha et al. \(2016\)](#), but all differences could be explained by population
383 structure. We performed similar tests on the phased data and obtained almost the same results (see
384 Supp.). The ASI and EUR populations showed the largest number of genes differentially expressed,
385 and EUR and ME the smallest. However, this pattern was not detectable in the model accounting
386 for population genetic structure (see Supp.). Thus, variation in expression level based on phased
387 data between two subgenomes did not differ much from the variation based on unphased data and
388 could as well be explained by the demographic processes in these populations.

389 DISCUSSION

390 In the present study, we analyzed the genetic changes experienced by a recently formed allopolyploid *C. bursa-pastoris* since its founding, focusing on the evolutionary trajectories followed by its
391 two subgenomes in demographically and genetically distinct populations from Europe, the Middle
392 East, and Asia. The shift to selfing and polyploidy had a global impact on the species, resulting
393 in a sharp reduction of the effective population size in all populations, that was accompanied by
394 relaxed selection and accumulation of deleterious mutations. However, the two subgenomes were
395 not similarly affected, with the magnitude of the subgenome-specific differences depending on the
396 population considered. The relative patterns of nucleotide diversity, genetic load, selection and
397 gene expression between the two subgenomes in the European and the Middle Eastern populations
398 were distinct to that observed in Asia. The differences between populations were further enhanced
399 by post-speciation hybridization of *C. bursa-pastoris* with local parental lineages. Below, we
400 discuss these global and local effects in more detail and their consequences for the history of the
401 species.
402

403 **Effect of parental legacy**

404 The effective population size of the diploid outcrossing ancestor of *C. bursa-pastoris*, *C. grandiflora*,
405 is ten times larger than that of its selfing ancestor *C. orientalis* (Douglas et al. 2015). Any
406 analysis of the difference in effective population size between the subgenomes of *C. bursa-pastoris*
407 or of their evolutionary trajectories must therefore account for this initial difference. After the
408 bottleneck associated with the origin of *C. bursa-pastoris* and the reduction in N_e due to the shift
409 to selfing (Charlesworth 2009), the effective population sizes of the two subgenomes are expected
410 to progressively converge and decrease along the same trajectory.

411 While this was indeed the observed overall pattern, the trajectories followed by the two
412 subgenomes in the three populations differed: in Europe the initial level was similar to that in
413 the Middle East but higher than in Asia and the decline of N_e of the Cg subgenome was de-
414 layed compared to the sudden decline experienced by the Co subgenome. In contrast, in Asia the
415 two subgenomes initially followed similar downwards trajectories but N_e increased again in both

416 subgenomes at around 40,000 ya. In the diploid *C. orientalis*, there was a period of stasis followed
417 by a steeper decline than in the tetraploid. The difference in demography across the three regions
418 could indicate multiple origins of *C. bursa-pastoris* as suggested by [Douglas et al. \(2015\)](#) and the
419 difference between the diploid and the tetraploid could reflect a mixture of the population expansion
420 experienced by the tetraploid and the buffering effect of tetraploidy against deleterious mutations.

421 There was a clearly noticeable difference between the two subgenomes in the number of
422 inherited deleterious mutations. Based on the strong differences in N_e , one would expect the
423 efficacy of selection to be much higher in *C. grandiflora* than in *C. orientalis* that has a much
424 smaller N_e ([Kimura 1983](#)). In the analysis of the genetic load, we indeed observed that *C. orientalis*
425 had a higher proportion of deleterious mutations than *C. grandiflora*. Hence, the amount of genetic
426 load most likely was different between the Cg and Co subgenomes of *C. bursa-pastoris* at the
427 time of the species emergence. Interestingly, hundreds of thousands of generations of selfing did
428 not totally erase the differences between the two subgenomes and, today, the Co subgenome still
429 carries more deleterious mutations than the Cg subgenome. This difference was smaller than
430 the difference between *C. orientalis* and *C. grandiflora*, but it was still significant. Nucleotide
431 diversity also demonstrated the effect of parental legacy. The Cg subgenome inherited from the
432 more variable outcrosser *C. grandiflora* was still more diverse in all populations except the Asian
433 one. The maintenance of part of the parental legacy in both cases suggest that, in spite of their initial
434 differences, both subgenomes have experienced similar levels of fixation since the creation of the
435 species. The Asian population is an exception in this regards because it was affected by secondary
436 gene flow as discussed below. Variation in nucleotide diversity in the coding part of the genome
437 also demonstrated similarity in the efficacy of purifying selection between the two subgenomes and
438 their corresponding parental lineages, though the pattern in the ASI population was the reverse of
439 that observed in the parental lineages. The effect of parental legacy on gene degeneration was also
440 noted in [Douglas et al. \(2015\)](#). Thus, the effect of the genetic background of hybridizing species
441 may not be as overwhelming as the effect of mating system but it still impacts the fate of the two
442 subgenomes long after the species arose.

443 **Subgenome-specific introgression and/or multiple origins**

444 Based on coalescent simulations and the amount of shared variation between *C. bursa-pastoris*
445 and its parental species [Douglas et al. \(2015\)](#) ruled out a single founder but noted that it would be very
446 difficult to estimate the exact number of founding lineages. [Douglas et al. \(2015\)](#) did not consider
447 hybridization but an earlier study ([Slotte et al. 2008b](#)) detected gene flow from *C. rubella* to the
448 European *C. bursa-pastoris* using 12 nuclear loci and a coalescent-based isolation-with-migration
449 model. The present study adds two new twists to the story. First, our results indicate that shared
450 polymorphisms were not symmetrical: namely, in the EUR and ME populations introgression
451 from *C. rubella* occurred on the *C. grandiflora* subgenome whereas in ASI introgression from
452 *C. orientalis* occurred on the *C. orientalis* subgenome. Second, in both the NJ and density trees,
453 *C. orientalis* appears as derived from the *C. bursa-pastoris* Co subgenome rather than the converse
454 as one would have expected. No such anomaly was observed for *C. grandiflora* that, as expected,
455 grouped at the root of the *C. bursa-pastoris* Cg subgenome. These results could be explained by
456 a mixture of multiple origins and more recent introgression. Multiple origins seem to be common
457 in allotetraploids ([Soltis et al. 1993](#); [Soltis and Soltis 1999](#)) and interploidy gene flow has already
458 been inferred for the *Capsella* ([Slotte et al. 2008b](#)) and other plant genera ([Balao et al. 2016](#);
459 [Anamthawat-Jónsson and Thórsson 2003](#)).

460 Our crossing results did not reject the possibility of ongoing admixture between *C. bursa-*
461 *pastoris* and parental lineages in both Europe and Asia. European and Asian populations of
462 *C. bursa-pastoris* partially overlap in the distribution ranges with *C. rubella* and *C. orientalis*,
463 respectively (Fig. 1A). The exact proportion of introgression remains unclear at this stage. Taken
464 at face value, the strongest admixture was between the ASI Co subgenome and *C. orientalis*.
465 Considering the overlapping estimates of *f*-statistics and HAPMIX, the proportion of admixture
466 of the ASI Co subgenomes with *C. orientalis* was around 14%-23%. The admixture between the
467 EUR Cg subgenome and *C. rubella* was also strong, being around 8-20%. There were also signs
468 of minor admixture in the ME population with both *C. orientalis* and *C. rubella*. This lack of a
469 non-admixed population posed a problem of correct estimation of the proportion of admixture for

470 both the ABBA-BABA and HAPMIX approaches.

471 In the ABBA-BABA test, departure from the assumptions can lead to under- or overestimated
472 introgression. In the present case, some proportion of the variation shared between P_3 and both
473 P_1 and P_2 populations could be due to introgression and not to incomplete lineage sorting and
474 this would lead to underestimating the amount of admixture. On the other hand, small N_e and
475 recent divergence of the populations used in the test can inflate estimates of D (Martin et al.
476 2015). Further, the behavior of D in tests involving both selfing and outcrossing species has not
477 been assessed yet. The D statistics were significantly different from zero in all our comparisons
478 suggesting that admixture did indeed occur in all populations of *C. bursa-pastoris*. The f statistic
479 is considered less prone to be affected by these factors (Martin et al. 2015), and it was more reliable
480 in our tests too. Its values were close to zero in the alternative combinations for the ABBA-BABA
481 tests where we did not expect to find admixture, while D had high estimates (Table S10). Thus, the
482 f values are the closest to the real proportion of admixture we could get.

483 In HAPMIX, when one reference population is admixed, the program probably compensates
484 for this extra relatedness between the reference populations by inflating intermediate introgression
485 probabilities. Therefore, we observed the discrepancy between the results of HAPMIX and ABBA-
486 BABA in the estimates of admixture between the EUR Co subgenome and *C. orientalis*. However,
487 the results for the Cg subgenome largely agreed between HAPMIX and ABBA-BABA and, together
488 with the results by Slotte et al. (2008b) and our crossing experiment, bolsters the hypothesis of
489 admixture between *C. rubella* and *C. bursa-pastoris* in Europe. On balance, a scenario with a
490 single origin of *C. bursa-pastoris* with later rampant admixture with *C. orientalis* in Asia and less
491 extensive admixture with *C. rubella* in Europe is consistent with our data.

492 On the other hand, our results could also be obtained under a scenario of multiple origins. This
493 scenario seems particularly likely if one looks at Fig. 4D,C, where the history of the Co and Cg
494 subgenomes are totally different. If we assume that *C. orientalis* and *C. grandiflora* are indeed
495 parental lineages and there was no unknown parental lineage that went extinct, this picture can be
496 only explained by a separate and more recent origin of the ASI population (Fig. 8). However, the

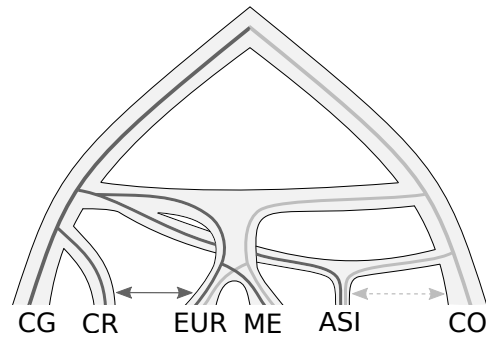


Fig. 8. A tentative scenario of multiple origin of *C. bursa-pastoris*. The Asian population originated separately from other *C. bursa-pastoris* populations. There may still be gene flow between the Asian population and *C. orientalis* (dashed arrow). There is gene flow between the European *C. bursa-pastoris* and *C. rubella* (solid arrow). ASI, EUR, ME, CR, CO, CG indicate Asian, European and Middle Eastern populations of *C. bursa-pastoris*, *C. rubella*, and parental species *C. orientalis* and *C. grandiflora*, respectively.

497 scenario of multiple origins and post-speciation admixture are not mutually exclusive. The signs
498 of gene flow between EUR and *C. rubella* are still best explained by post-speciation admixture.
499 The weak signs of admixture between *C. bursa-pastoris* and *C. orientalis* in EUR and ME are also
500 difficult to fit into a scenario involving only multiple origins. A possibility is that these signs of
501 admixture resulted from gene flow from ASI to EUR and ME within *C. bursa-pastoris*. The ASI
502 population is more related to *C. orientalis* and the presence of its alleles in EUR and ME could
503 be spuriously recognized as introgressed from ASI. Regardless of whether a single or a multiple
504 origin scenario is the true one, our results demonstrate that the history of *C. bursa-pastoris* is far
505 more complex than previously imagined.

506 **Weak subgenome-specific expression differences**

507 Many allopolyploid species show subgenome expression bias, where one subgenome tends
508 to be over-expressed relative to the other one (Schnable et al. 2011; Flagel et al. 2008; Li et al.
509 2014; Woodhouse et al. 2014; Li et al. 2014; Schnable et al. 2011). This expression dominance is
510 often observed in synthetic allopolyploids (He et al. 2012; Lemmon et al. 2014; Yang et al. 2016;
511 Bell et al. 2013) and thus the major part of such preferential subgenome dominance is probably
512 established immediately after allopolyploidization. The subgenome expression dominance is also

513 suggested to be largely defined by parental expression differences (Buggs et al. 2014; Gottlieb 2003).
514 Contradictory results on patterns of subgenome specific expression in *C. bursa-pastoris* have been
515 obtained so far. Douglas et al. (2015) concluded that there is no strong subgenome expression
516 bias and those few genes showing subgenome-specific expression could be explained by parental
517 expression differences. However, genes with subgenome-specific expression do show a slight bias
518 towards over-expression of the Cg subgenome inherited from *C. grandiflora / rubella* lineage on
519 the Figure 3B in Douglas et al. (2015). In contrast, Steige et al. (2016) reported a higher expression
520 of the Co subgenome inherited from *C. orientalis* in three accessions, and Cg over-expression in a
521 fourth one (CbpGR). Steige et al. (2016) hypothesized that the over-expression of the Co subgenome
522 might be related to a higher number of transposable elements in this subgenome, but they did not
523 find any evidence of this and could not explain the down-regulation of the Co subgenome in the
524 CbpGR accession and in the artificial hybrid between *C. rubella* and *C. orientalis*.

525 Considering the population histories of *C. bursa-pastoris* sheds some light on these discrepan-
526 cies. The results of Douglas et al. (2015) and Steige et al. (2016) are consistent with the hypothesis
527 that *cis*-regulatory differences between the *C. orientalis* and *C. grandiflora / rubella* genomes
528 result in over-expression of the Cg subgenome in a hybrid comprising both genomes. Thus, in
529 the absence of other factors, the slight over-expression of the Cg subgenome would be the default
530 HSE pattern in *C. bursa-pastoris*. In accordance with this, we observed over-expression of the Cg
531 subgenome in the ME and EUR populations that are most likely the closest to the region of origin
532 of *C. bursa-pastoris* (Cornille et al. 2016). The accessions that show over-expression of the Cg
533 subgenome in Douglas et al. (2015) (SE14 from Sweden) and in Steige et al. (2016) (CbpGR from
534 Greece), as we now know (Cornille et al. 2016) belong to the EUR population. Hence, their results
535 are consistent with ours and expected if the HSE is defined primarily by the differences between
536 the parental lineages. On the other hand, we observed that genes with HSE in the ASI population
537 showed equal expression between the two subgenomes. The accessions showing over-expression
538 of the Co subgenome in Steige et al. (2016) also mostly belong to the ASI population (CbpKMB
539 and CbpGY, though not CbpDE that putatively originates from Germany). Thus, the Asian ac-

540 cessions show the HSE that is different from the default pattern. This difference can be caused
541 by the selection preference for the Co subgenome and/or by introgression from *C. orientalis* that
542 enhanced the *cis*-regulatory elements of the Co subgenome. The ASI population experienced a
543 strong population bottleneck, so genetic drift played some role as well. These explanations need
544 to be confirmed because HSE can be influenced by many factors (e.g. *trans*-regulatory elements,
545 gene methylation, transposable elements), but it is clear that there are different directions of HSE
546 in populations of *C. bursa-pastoris* and they are caused by the different evolutionary histories of
547 those populations.

548 The reason we observed an equal expression between subgenomes in ASI, whereas [Steige](#)
549 [et al. \(2016\)](#) detected expression bias of the Co subgenome for Asian samples, could also be due
550 to different approaches in our analyses. First, we extracted RNA from seedling, whereas [Steige](#)
551 [et al. \(2016\)](#) obtained RNA from leaves and flower buds. Variation in HSE for different tissues of
552 *C. bursa-pastoris* is not characterized yet, so the Co expression in seedlings may not be apparent
553 yet. Second, we mapped reads to the *C. rubella* reference with masked polymorphism, whereas
554 [Steige et al. \(2016\)](#) used the reconstructed reference of an F1 hybrid between *C. orientalis* and
555 *C. rubella*. The bias in our DNA data was not stronger than in [Steige et al. \(2016\)](#), so which method
556 is more appropriate remains to be found out.

557 **Neutral inter-population expression differences**

558 We have previously reported that differences among populations in overall gene expression
559 variation (i.e. from unphased data) in *C. bursa-pastoris* primarily reflect population structure
560 and hence are mostly driven by genetic drift ([Kryvokhyzha et al. 2016](#)). The current study of
561 phased gene expression data is consistent with this result. Both the differential gene expression
562 analysis of each subgenome and the generalized linear model analysis of HSE data as proportions
563 revealed proportionally similar differences between populations and these differences were all
564 explained by the genetic population structure in the species. Our results also demonstrated that
565 genes showing significant HSE largely overlapped between populations and these genes were not
566 strongly enriched for GO terms. These genes probably evolve under a compensatory drift model

567 (Thompson et al. 2016). This was evident in the direction of the HSE, which was the same in all
568 accessions. The correlation in levels of HSE is stronger within than between populations, which
569 is also consistent with evolution by drift. Hence, gene expression variation does not show strong
570 adaptive changes in the early stages of the evolution of *C. bursa-pastoris*. It is still possible that some
571 of the gene expression differences are not neutral and we have previously discussed the potential
572 pitfalls of detecting adaptive differences in structured populations (Kryvokhyzha et al. 2016). The
573 asymmetric over-expression between populations, for instance, agrees with the presence of some
574 selective differences between populations.

575 CONCLUSION

576 Three salient, and sometimes unexpected, features of the evolution of the tetraploid shepherd's
577 purse that emerged from the present study, are its complex origin and the magnitude of introgression
578 with diploid relatives, the long-lasting effects of the difference between its two parental species and
579 the importance of demography in shaping its current genomic diversity. Hence, the present study
580 suggests that understanding the evolution of tetraploid species without paying due attention to the
581 historical and ecological backgrounds under which it occurred could be misleading.

582 MATERIALS AND METHODS

583 Sequence data

584 We obtained the whole genome sequences of 31 accessions of *C. bursa-pastoris* and the seedling
585 transcriptomes of 24 of these accessions. Transcriptome data used in this study were generated
586 previously (Kryvokhyzha et al. 2016). Whole genome DNA data consisted of 10 accessions
587 downloaded from GenBank (PRJNA268827) and 21 accessions sequenced in this study. New
588 DNA samples were sequenced using the same technology as the downloaded ones (100-bp paired-
589 end reads, Illumina HiSeq 2000 platform, SciLife, Stockholm, Sweden). The mean genomic
590 coverage of *C. bursa-pastoris* samples was 47x. We also used genomic data of 10 *C. orientalis*
591 and 13 *C. grandiflora* samples from GenBank (PRJNA245911, PRJNA254516). For the analysis
592 requiring an out-group, we used the whole genome assembly of *Neslia paniculata* (Slotte et al.
593 2013). Detailed information on the samples is provided in the Supporting Information.

594 Genotype calling and phasing

595 DNA reads from each individual were mapped to the *Capsella rubella* reference genome (Slotte
596 et al. 2013) using Stampy v1.0.22 (Lunter and Goodson 2011) with default parameters, except that
597 the substitution rate was set to 0.025 to account for the divergence from the reference. Potential
598 PCR duplicates were marked using Picard Tools 1.115 (<http://picard.sourceforge.net>) and
599 were ignored during genotyping. Genotypes were called using *HaplotypeCaller* from the Genome
600 Analysis Tool Kit (GATK) v3.5 (McKenna et al. 2010) in the GVCF mode and heterozygosity set
601 to 0.015. Genotypes were filtered for depth between 6 and 100 reads (the 5th and 99th coverage
602 percentiles, respectively). This approach produced a VCF file containing all called sites. This
603 VCF was used in the analyses requiring both polymorphic and monomorphic sites for correct
604 estimates. To obtain a set of SNPs with the highest confidence possible, we generated another VCF
605 file that contained only polymorphic sites and applied more stringent filtering. We set to no-call
606 all sites that met the following criteria: $MQ < 30$, $SOR > 4$, $QD < 2$, $FS > 60$, $MQRankSum$
607 < -20 , $ReadPosRankSum < -10$, $ReadPosRankSum > 10$. These filtering criteria were defined
608 following GATK Best Practices (Auwera et al. 2013) with some adjustment guided by the obtained

609 distributions of the GATK annotation scores (Fig. S10).

610 To phase the *C. bursa-pastoris* homeologs, we run HapCUT version 0.7 (Bansal and Bafna
611 2008) on each sample from the VCF with the stringently filtered SNPs. The phased haplotype
612 fragments were then joined into two sequences descended from *C. grandiflora* and *C. orientalis*.
613 The origin of haplotypes in HapCUT fragments was defined using sites with fixed heterozygotes in
614 *C. bursa-pastoris* and fixed differences between *C. grandiflora* and *C. orientalis*. Fragments that
615 had small (< 2 sites) or no overlap with variation in *C. grandiflora* and *C. orientalis* as well as those
616 that looked chimeric (prevailing phasing state was supported by less than 90% of sites) were set
617 to missing data (Fig. S11). Additionally, we also set to missing the sites that were defined as not
618 real variants or not heterozygous by HapCUT (flagged with *FV*). HapCUT phasing produced the
619 alignment that had only heterozygous sites and removed all the sites that were non-variant within
620 but variable between individuals. We restored this inter-individual variation with introduction of
621 the same proportion of missing data into non-variant sites as it was introduced to heterozygous sites
622 during the phasing. Similarly, we also merged the phased SNPs dataset with whole genome data.

623 The reference genomes of *C. grandiflora* and *C. orientalis* were created using the GVCF files
624 produced by Douglas et al. (2015). The variants were called as described above with additional
625 filtering for fixed differences between the two species. For some of the analyses, where the
626 software was not able to treat heterozygous genotypes properly, we pseudo-phased the sequences
627 of *C. grandiflora* and *C. orientalis* by randomizing alleles in heterozygous genotypes.

628 The final data-sets in all the analyses comprised the alignment of phased *C. bursa-pastoris*
629 sequences, *C. grandiflora*, *C. orientalis*, *C. rubella* (the reference sequence) and *N. paniculata*.
630 This alignment was filtered for missing data such that genomic positions with more than 80% of
631 missing genotypes were removed. We also removed the repetitive sequences as annotated in (Slotte
632 et al. 2013) and pericentromeric regions that we delineated based on the density of repetitive regions
633 and missing data.

634 **Reconstruction of the ancestral sequences**

635 Several analyses presented in this paper required polarized sequence data. The most common
636 approach to polarizing the alleles is to use an outgroup. However, the alignment of *Capsella*
637 species and *N. paniculata*, the nearest outgroup with a whole genome sequence available, resulted
638 in substantial reduction of the dataset due to missing data. To overcome this drawback, as well
639 as to track mutations' origin on the phylogenetic branches, we reconstructed ancestral sequences
640 for major phylogenetic splits. The reconstruction was performed on the tree that was assumed
641 to represent a true history of the *Capsella* species (Fig. S12) using the empirical Bayes joint
642 reconstruction method implemented in PAML v4.6 (Yang 1997).

643 **Population differentiation**

644 To assess the degree of differentiation among populations for the two subgenomes, we estimated
645 absolute divergence (D_{xy}) and nucleotide diversity (π) of the phased genomes using a sliding
646 window approach. The estimates were calculated on non-overlapping 100 Kb windows using the
647 *EggLib* Python module (De Mita and Siol 2012). The p -values for the difference in mean values
648 were estimated using 10,000 bootstrap resamples from 100 Kb windows.

649 **Temporal change in N_e**

650 We reconstructed changes of N_e over time with both PSMC (Li and Durbin 2011) and SMC++
651 (Terhorst et al. 2017). We first masked potential CpG islands and all nonsynonymous sites in
652 the genome to avoid bias caused by variation in mutation rates or selective effects. We randomly
653 paired haplotypes for estimation in *C. orientalis* and did the same for estimations based on the
654 two subgenomes of *C. bursa-pastoris*. SMC++ was run on all samples from a population, with
655 default parameter settings. For PSMC runs, we set parameters to “-N25 -t15 -r5 -p 4+25*2+4+6”.
656 Variation in N_e was estimated using 100 bootstrap replicates and three different pairs. We chose
657 a mutation rate equal to the mutation rate of *A. thaliana*, $\mu = 7 \times 10^{-9}$ per site per generation
658 (Ossowski et al. 2010) and a generation time of 1 year for all *Capsella* species.

659 **Phylogenomic analyses**

660 We reconstructed a whole genome phylogeny to explore the relationship between the phased
661 subgenomes of the three populations of *C. bursa-pastoris* as well as its parental species. To inves-
662 tigate the local phylogenetic relationships along the genome, we also conducted a sliding window
663 phylogenetic analysis using non-overlapping 100 Kb windows. In both analyses, phylogenetic trees
664 were reconstructed using the neighbor-joining algorithm and absolute genetic distance in R package
665 *ape* (Paradis et al. 2004). Additionally, a whole genome phylogenetic tree was also reconstructed
666 using the maximum-likelihood approach with the GTRGAMMA model and 100 bootstrap replicates
667 in *RAxML* v8.2.4 (Stamatakis 2014) (Fig. S13). The trees from the sliding window analysis were
668 described by counting the frequency of monophyly of different groups with the Newick Utilities
669 (Junier and Zdobnov 2010). The variation in topology across the genome was also described using
670 topology weighting implemented in TWISS (Martin and Van Belleghem 2017). The weighting
671 was estimated for 100 SNPs windows where each sample was genotyped for at least 50 SNPs. To
672 test for the difference in mean topology weighting, we fitted the generalized linear model with a
673 binomial distribution and performed multiple comparisons for the contrasts of interest with the *glht*
674 function from the *multcomp* library in *R* (Hothorn et al. 2008).

675 **Tests for gene flow**

676 To evaluate the presence of gene flow between the parental species and *C. bursa-pastoris*,
677 we calculated the ABBA-BABA based statistics, D , an estimate of departure from incomplete
678 lineage sorting, and f , an estimate of admixture proportion (Green et al. 2010; Durand et al. 2011).
679 These statistics and their significance, which was estimated with a 1Mb block jackknife method,
680 were calculated from population allele frequencies with scripts from Martin et al. (2013). We
681 also used Hapmix (Price et al. 2009) to infer haplotype blocks of introgression from the diploids
682 *C. grandiflora*, *C. rubella*, and *C. orientalis* into the three populations of *C. bursa-pastoris* for each
683 phased subgenome. We removed sites with more than 20% missing data for each population. The
684 remaining missing data was imputed for the parental populations used in each analysis. We used
685 Asian *C. bursa-pastoris* as the alternate reference population, as we suspected little introgression

686 from either of the investigated diploids, except for investigations of Asia itself where we used
687 the European population. For the Co subgenome, we investigated introgression from just the
688 diploid *C. orientalis* using the reference scheme described for the Cg subgenome. However, after
689 inspection, introgression into the Middle East Co subgenome was the lowest so we instead used it
690 as the reference for Europe and Asia. As this method determines the probability of ancestry from a
691 diploid progenitor population relative to a non-admixed *C. bursa-pastoris* subgenome population,
692 we defined regions of the subgenomes as putatively introgressed if the probability of ancestry from
693 the progenitor diploid was greater than 50%.

694 To check for reproductive barriers between *C. bursa-pastoris* and its diploid relatives, we
695 performed artificial crosses. The crosses were made in both directions using *C. bursa-pastoris*
696 as a mother plant and as a pollen donor. Each cross was replicated at least three times and each
697 biological replicate consisted of 5 or more siliques. The details are provided in the Supplementary
698 Information.

699 **Selection tests**

700 To search for selective sweeps, we used SweepFinder2 (DeGiorgio et al. 2016). SweepFinder2
701 was run on the data-set that besides polarized SNPs also included fixed derived alleles. This enables
702 accounting for variation in mutation rate along the genome and increases power to detect sweeps
703 (Huber et al. 2016). The critical composite likelihood ratio (CLR) values were determined using
704 a 1% cut-off of the CLR values estimated in 100 simulations under a standard neutral model. The
705 simulations were performed with *fastsimcoal2* (Excoffier et al. 2013). We assumed a mutation rate
706 of $7e-9$ per site per generation as in Douglas et al. (2015), the population effective sizes for every
707 population and subgenome were inferred from the θ values approximated by genetic diversity (π),
708 and the average recombination rate was estimated using LDhelmet v1.7 (Chan et al. 2012). In
709 addition, we estimated the ratio between nucleotide diversity at 0-fold (π_0) and 4-fold degenerate
710 sites (π_4) in 5-6 samples with the lowest amount of missing data in each group. The details of the
711 data used to estimate π_0/π_4 are provided in Table S10.

712 We also tested if the detected sweep regions were not the result of introgression or genome

713 conversion. We compared the absolute genetic distance (D_{xy}) of each sweep region between all the
714 groups and if the distance was the closest to one of the parental species or the opposite subgenome,
715 such regions were classified as introgression or conversion, respectively. To reduce the number
716 of potential false positives, we removed pericentromeric regions and all regions with repetitive
717 sequences as annotated in [Slotte et al. \(2013\)](#). Sweep regions with less than 10Kb apart were joined
718 together and treated as one region.

719 **Genetic load estimation**

720 To identify differences in genetic load between populations of *C. bursa-pastoris* (as well as to
721 assess the effect of selfing on accumulation of deleterious mutations), we classified mutations into
722 tolerated and deleterious ones using SIFT4G ([Vaser et al. 2016](#)). We built the SIFT4G *Capsella*
723 *rubella* reference partition database and used it to annotate our SNPs dataset. Then we analyzed
724 the frequencies of tolerated and deleterious mutations. We also verified this analysis by using *A.*
725 *thaliana* SIFT4G database and annotating *C. bursa-pastoris* according to the alignment between the
726 two species. This verification was performed to make sure that the observed results were not due
727 to a reference bias, because *C. rubella* is closer to *C. grandiflora* than to *C. orientalis*. To get only
728 the annotation of the mutations that occurred after speciation of *C. bursa-pastoris*, we polarized
729 the mutations with the reconstructed ancestral sequences (see above) and analyzed only derived
730 mutations. We verified this polarization by analyzing only species(subgenome)-specific mutation
731 (e.g. mutations unique to *C. bursa-pastoris* Co subgenome, *C. bursa-pastoris* Cg subgenome,
732 *C. orientalis*, *C. grandiflora*, and *C. rubella*) (Fig. S14). All the counts were presented relative
733 to the total number of annotated sites to avoid bias caused by variation in missing data between
734 samples. The means of the genetic load were compared using the generalized linear model as
735 we did for the topology weighting except that here we used a quasibinomial distribution due to
736 overdispersion.

737 **Homeolog-specific expression analyses**

738 Mapping of RNA-Seq reads to the *C. rubella* reference genome was conducted similarly to the
739 mapping of DNA data using Stampy v1.0.22 ([Lunter and Goodson 2011](#)) with the substitution rate

740 set to 0.025. Although potential PCR duplicates are usually not removed from RNA-Seq data, for
741 the allele-specific expression analysis removing duplicates is recommended (Castel et al. 2015).
742 We marked duplicates with Picard Tools 1.115 and did not use them during the genotyping and
743 homeolog-specific expression assessment. Variants were called using *HaplotypeCaller* (GATK)
744 with heterozygosity set to 0.015, and minimum Phred-scaled call confidence of 20.0, and minimum
745 Phred-scaled emit confidence of 20.0 as recommended for RNA-Seq data in GATK Best Practices
746 (Auwera et al. 2013). Among the obtained polymorphic sites those that had $MQ < 30.00$, $QD <$
747 2.00 , $FS > 30.000$ were filtered out. Calls with coverage of fewer than 10 reads were also excluded.
748 Alleles counting was carried out using *ASEReadCounter* from GATK.

749 Homeolog-specific expression was assessed within the statistical framework developed by Skelly
750 et al. (2011). This framework uses a Markov chain Monte Carlo (MCMC) method for parameter
751 estimation and incorporates information from both RNA and DNA data to exclude highly biased
752 SNPs and calibrate for the noise in read counts due to statistical sampling and technical variability.
753 First, we used DNA data to identify and remove SNPs that strongly deviated from the 0.5 mapping
754 ratio. Second, we estimated the variation in allele counts using unbiased SNPs in the DNA data.
755 Next, we fitted an RNA model using parameter estimated from DNA data in the previous step.
756 Finally, we calculated a Bayesian analog of false discovery rate (FDR) with a posterior probability
757 of homeologue specific expression (HSE) > 0.99 and defined genes with significant HSE given the
758 estimated FDR. All inferences were performed using 200,000 MCMC iterations with burn-in of
759 20,000 and thin interval of 100. Each model was run three times with different starting parameters
760 to verify convergence.

761 To test for differences between populations of *C. bursa-pastoris*, we analyzed phased expres-
762 sion data as was done with unphased data in Kryvokhyzha et al. (2016). We tested differences
763 between populations in two ways: each subgenome was processed individually in *edgeR*, and both
764 subgenomes were analyzed together as proportional data by fitting a generalized linear model. In
765 addition, we performed correction for genetic population structure by fitting generalized linear
766 mixed models (see Supp.).

767 **DATA ACCESS**

768 The sources of the data obtained from previous studies are provided in the Material and
769 Methods. DNA sequences data generated for 21 accessions in this study is submitted to the NCBI
770 database under the Sequence Read Archive number SRAXXXXXX. Both phased and unphased
771 genotype data, phylogenetic trees, reconstructed ancestral sequences, estimates of π and D_{xy} with
772 sliding window approach, results of PSMC and SMC++, SIFT annotations, CLR estimates of
773 *sweepFinder2*, HAPMIX output, homeologue-specific gene expression values, and R scripts are
774 deposited to the Dryad Digital Repository doi: XXXXXXXX.

775 **ACKNOWLEDGMENTS**

776 We thank Clément Lafon-Placette and Mohammad Foteh Ali for help with crosses, and Pas-
777 cal Milesi and Ludovic Dutoit for discussion of the results. We are especially grateful to Luis
778 Leal for detailed feedback on the manuscript. Most of the analyses were carried out at Upp-
779 sala Multidisciplinary Center for Advanced Computational Science (UPPMAX) under the project
780 b2013191. This study was supported by grants from the Swedish Research Council (VR) and the
781 Erik Philip-Sörensens Stiftelse to ML.

782 **AUTHORS CONTRIBUTIONS**

783 DK and AC phased the data and performed selection tests. DK carried out phylogenetic and gene
784 expression analyses. DK, MCE, TD, NT, TVK analyzed genetic load. JC performed demographic
785 analyses and analyzed nucleotide diversity in the coding part of the genome. DK and MG did the
786 crosses. DK, JK, and SG performed tests for introgression. DK and ML drafted the article with
787 inputs from all other authors. JRS, UL, SG, SIW and ML supervised the project.

788 **DISCLOSURE DECLARATION**

789 No competing interests

790 **REFERENCES**

- 791 Adams KL, Wendel JF. 2005. Polyploidy and genome evolution in plants. *Cur Opin Plant Bio* **8**:
792 135–141.
- 793 Anamthawat-Jónsson K, Thórsson AT. 2003. Natural hybridisation in birch: triploid hybrids be-
794 tween *Betula nana* and *B. pubescens*. *Plant Cell Tiss Org Cult* **75**: 99–107.
- 795 Auwera GA, Carneiro MO, Hartl C, Poplin R, del Angel G, Levy-Moonshine A, Jordan T, Shakir
796 K, Roazen D, Thibault J, et al.. 2013. From FastQ data to high-confidence variant calls: the
797 genome analysis toolkit best practices pipeline. *Cur Prot in Bioin* **11**: 1–43.
- 798 Balao F, Tannhäuser M, Lorenzo MT, Hedrén M, Paun O. 2016. Genetic differentiation and
799 admixture between sibling allopolyploids in the *Dactylorhiza majalis* complex. *Heredity* **116**:
800 351–361.
- 801 Bansal V, Bafna V. 2008. HapCUT: an efficient and accurate algorithm for the haplotype assembly
802 problem. *Bioinformatics* **24**: i153–i159.
- 803 Barringer BC. 2007. Polyploidy and self-fertilization in flowering plants. *Amer J Botany* **94**: 1527–
804 1533.
- 805 Bell GD, Kane NC, Rieseberg LH, Adams KL. 2013. RNA-seq analysis of allele-specific expression,
806 hybrid effects, and regulatory divergence in hybrids compared with their parents from natural
807 populations. *Gen Bio Evol* **5**: 1309–1323.
- 808 Bomblies K, Higgins JD, Yant L. 2015. Meiosis evolves: adaptation to external and internal
809 environments. *New Phyt* **208**: 306–323.
- 810 Brochmann C, Brysting A, Alsos I, Borgen L, Grundt H, Scheen AC, Elven R. 2004. Polyploidy in
811 arctic plants. *Biol J Linn Soc* **82**: 521–536.
- 812 Buggs RJ, Wendel JF, Doyle JJ, Soltis DE, Soltis PS, Coate JE. 2014. The legacy of diploid
813 progenitors in allopolyploid gene expression patterns. *Phil Trans R Soc B* **369**: 1–13.
- 814 Buggs RJA, Chamala S, Wu W, Tate JA, Schnable PS, Soltis DE, Soltis PS, Barbazuk WB. 2012.
815 Rapid, repeated, and clustered loss of duplicate genes in allopolyploid plant populations of
816 independent origin. *Curr Bio* **22**: 248–252.

- 817 Castel SE, Levy-Moonshine A, Mohammadi P, Banks E, Lappalainen T. 2015. Tools and best
818 practices for data processing in allelic expression analysis. *Genome Bio* **16**: 1–12.
- 819 Chan AH, Jenkins PA, Song YS. 2012. Genome-wide fine-scale recombination rate variation in
820 *Drosophila melanogaster*. *PLoS Genet* **8**: 1–28.
- 821 Charlesworth B. 2009. Effective population size and patterns of molecular evolution and variation.
822 *Nature Rev Gen* **10**: 195–205.
- 823 Chen J, Glémin S, Lascoux M. 2017. Genetic diversity and the efficacy of purifying selection across
824 plant and animal species. *MBE* **34**: 1417–1428.
- 825 Comai L. 2005. The advantages and disadvantages of being polyploid. *Nature Rev Gen* **6**: 836–846.
- 826 Cornille A, Salcedo A, Kryvokhyzha D, Glémin S, Holm K, Wright S, Lascoux M. 2016. Genomic
827 signature of successful colonization of Eurasia by the allopolyploid shepherd’s purse (*Capsella*
828 *bursa-pastoris*). *Mol Ecol* **25**: 616–629.
- 829 De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population genetics
830 and genomics. *BMC Gen* **13**: 1–12.
- 831 DeGiorgio M, Huber CD, Hubisz MJ, Hellmann I, Nielsen R. 2016. Sweepfinder2: increased
832 sensitivity, robustness and flexibility. *Bioinformatics* **32**: 1895–1897.
- 833 Douglas GM, Gos G, Steige KA, Salcedo A, Holm K, Josephs EB, Arunkumar R, Ågren JA,
834 Hazzouri KM, Wang W, et al.. 2015. Hybrid origins and the earliest stages of diploidization
835 in the highly successful recent polyploid *Capsella bursa-pastoris*. *Proc Natl Acad Sci* **112**:
836 2806–2811.
- 837 Doyle JJ, Flagel LE, Paterson AH, Rapp RA, Soltis DE, Soltis PS, Wendel JF. 2008. Evolutionary
838 genetics of genome merger and doubling in plants. *Annu Review Gen* **42**: 443–461.
- 839 Durand EY, Patterson N, Reich D, Slatkin M. 2011. Testing for ancient admixture between closely
840 related populations. *MBE* **28**: 2239–2252.
- 841 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa VC, Foll M. 2013. Robust demographic infer-
842 ence from genomic and SNP data. *PLoS Gen* **9**: 1–17.
- 843 Flagel L, Udall J, Nettleton D, Wendel J. 2008. Duplicate gene expression in allopolyploid *Gossyp-*

- 844 *ium* reveals two temporally distinct phases of expression evolution. *BMC Bio* **6**: 1–9.
- 845 Gilbert KJ, Sharp NP, Angert AL, Conte GL, Draghi JA, Guillaume F, Hargreaves AL, Matthey-
846 Doret R, Whitlock MC. 2017. Local adaptation interacts with expansion load during range
847 expansion: maladaptation reduces expansion load. *Amer Nat* **189**: 368–380.
- 848 Gottlieb L. 2003. Plant polyploidy: gene expression and genetic redundancy. *Heredity* **91**: 91–92.
- 849 Green RE, Krause J, Briggs AW, Maricic T, Stenzel U, Kircher M, Patterson N, Li H, Zhai W, Fritz
850 MHY, et al.. 2010. A draft sequence of the Neanderthal genome. *Science* **328**: 710–722.
- 851 Hartfield M, Bataillon T, Glémin S. 2017. The evolutionary interplay between adaptation and
852 self-fertilization. *Trends In Gen* **33**: 420–431.
- 853 He F, Zhang X, Hu J, Turck F, Dong X, Goebel U, Borevitz J, de Meaux J. 2012. Genome-
854 wide analysis of cis-regulatory divergence between species in the *Arabidopsis* genus. *MBE* **29**:
855 3385–3395.
- 856 Hothorn T, Bretz F, Westfall P. 2008. Simultaneous inference in general parametric models. *Bio-*
857 *metrical journal* **50**: 346–363.
- 858 Huber CD, DeGiorgio M, Hellmann I, Nielsen R. 2016. Detecting recent selective sweeps while
859 controlling for mutation rate and background selection. *Mol Ecol* **25**: 142–156.
- 860 Hurka H, Friesen N, German DA, Franzke A, Neuffer B. 2012. ‘Missing link’ species *Capsella ori-*
861 *entalis* and *Capsella thracica* elucidate evolution of model plant genus *Capsella* (brassicaceae).
862 *Mol Ecol* **21**: 1223–1238.
- 863 Junier T, Zdobnov EM. 2010. The Newick utilities: high-throughput phylogenetic tree processing
864 in the UNIX shell. *Bioinformatics* **26**: 1669–1670.
- 865 Kimura M. 1983. *The neutral theory of molecular evolution*. Cambridge University Press.
- 866 Kryvokhyzha D, Holm K, Chen J, Cornille A, Glémin S, Wright SI, Lagercrantz U, Lascoux M.
867 2016. The influence of population structure on gene expression and flowering time variation in
868 the ubiquitous weed *Capsella bursa-pastoris* (Brassicaceae). *Mol Ecol* **25**: 1106–1121.
- 869 Lemmon ZH, Bukowski R, Sun Q, Doebley JF. 2014. The role of cis regulatory evolution in maize
870 domestication. *PLoS Gen* **10**: 1–15.

- 871 Levin DA. 2002. *The role of chromosomal change in plant evolution*. Oxford University Press.
- 872 Li A, Liu D, Wu J, Zhao X, Hao M, Geng S, Yan J, Jiang X, Zhang L, Wu J, et al.. 2014. mRNA and
873 small RNA transcriptomes reveal insights into dynamic homoeolog regulation of allopolyploid
874 heterosis in nascent hexaploid wheat. *Plant Cell* **26**: 1878–1900.
- 875 Li H, Durbin R. 2011. Inference of human population history from individual whole-genome
876 sequences. *Nature* **475**: 493–496.
- 877 Lunter G, Goodson M. 2011. Stampy: a statistical algorithm for sensitive and fast mapping of
878 Illumina sequence reads. *Genome Res* **21**: 936–939.
- 879 Martin SH, Dasmahapatra KK, Nadeau NJ, Salazar C, Walters JR, Simpson F, Blaxter M, Manica A,
880 Mallet J, Jiggins CD. 2013. Genome-wide evidence for speciation with gene flow in *Heliconius*
881 butterflies. *Genome Res* **23**: 1817–1828.
- 882 Martin SH, Davey JW, Jiggins CD. 2015. Evaluating the use of ABBA-BABA statistics to locate
883 introgressed loci. *MBE* **32**: 244–257.
- 884 Martin SH, Van Belleghem SM. 2017. Exploring evolutionary relationships across the genome
885 using topology weighting. *Genetics* **206**: 429–438.
- 886 McGrath C, Lynch M. 2012. Evolutionary significance of whole-genome duplication. *Polyploidy*
887 *and genome evolution* Springer 1–20.
- 888 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K, Altshuler
889 D, Gabriel S, Daly M, et al.. 2010. The Genome Analysis Toolkit: a MapReduce framework for
890 analyzing next-generation DNA sequencing data. *Genome Res* **20**: 1297–1303.
- 891 Ossowski S, Schneeberger K, Lucas-Lledó JI, Warthmann N, Clark RM, Shaw RG, Weigel D, Lynch
892 M. 2010. The rate and molecular spectrum of spontaneous mutations in *Arabidopsis thaliana*.
893 *Science* **327**: 92–94.
- 894 Otto SP, Whitton J. 2000. Polyploid incidence and evolution. *Annu Review Gen* **34**: 401–437.
- 895 Pandit M, Tan H, Bisht M. 2006. Polyploidy in invasive plant species of Singapore. *Bot J Linn Soc*
896 **151**: 395–403.
- 897 Pandit MK, Pocock MJ, Kunin WE. 2011. Ploidy influences rarity and invasiveness in plants. *J*

- 898 *Ecology* **99**: 1108–1115.
- 899 Paradis E, Claude J, Strimmer K. 2004. APE: analyses of phylogenetics and evolution in R language.
900 *Bioinformatics* **20**: 289–290.
- 901 Peischl S, Dupanloup I, Bosshard L, Excoffier L. 2016. Genetic surfing in human populations: from
902 genes to genomes. *Cur Opin Gen Dev* **41**: 53–61.
- 903 Petit C, Thompson JD. 1999. Species diversity and ecological range in relation to ploidy level in
904 the flora of the Pyrenees. *Evol Ecol* **13**: 45–65.
- 905 Prentis PJ, Wilson JR, Dormontt EE, Richardson DM, Lowe AJ. 2008. Adaptive evolution in
906 invasive species. *Trends Plant Sci* **13**: 288–294.
- 907 Price AL, Tandon A, Patterson N, Barnes KC, Rafaels N, Ruczinski I, Beaty TH, Mathias R, Reich
908 D, Myers S. 2009. Sensitive detection of chromosomal segments of distinct ancestry in admixed
909 populations. *PLoS Gen* **5**: 1–18.
- 910 Ramsey J. 2011. Polyploidy and ecological adaptation in wild yarrow. *Proc Natl Acad Sci* **108**:
911 7096–7101.
- 912 Robertson K, Goldberg EE, Igic B. 2011. Comparative evidence for the correlated evolution of
913 polyploidy and self-compatibility in Solanaceae. *Evolution* **65**: 139–155.
- 914 Schnable JC, Springer NM, Freeling M. 2011. Differentiation of the maize subgenomes by genome
915 dominance and both ancient and ongoing gene loss. *Proc Natl Acad Sci* **108**: 4069–4074.
- 916 Skelly DA, Johansson M, Madeoy J, Wakefield J, Akey JM. 2011. A powerful and flexible statistical
917 framework for testing hypotheses of allele-specific gene expression from RNA-seq data. *Genome*
918 *Res* **21**: 1728–1737.
- 919 Slotte T, Hazzouri KM, Ågren JA, Koenig D, Maumus F, Guo YL, Steige K, Platts AE, Escobar
920 JS, Newman LK, et al.. 2013. The *Capsella rubella* genome and the genomic consequences of
921 rapid mating system evolution. *Nat Genet* **45**: 831–835.
- 922 Slotte T, Huang H, Lascoux M, Ceplitis A. 2008a. Polyploid speciation did not confer instant
923 reproductive isolation in *Capsella* (brassicaceae). *MBE* **25**: 1472–1481.
- 924 Slotte T, Huang H, Lascoux M, Ceplitis A. 2008b. Polyploid speciation did not confer instant

- 925 reproductive isolation in *Capsella* (Brassicaceae). *MBE* **25**: 1472–1481.
- 926 Soltis D, Soltis P, Rieseberg LH. 1993. Molecular data and the dynamic nature of polyploidy. *Crit*
927 *Rev Plant Sci* **12**: 243–273.
- 928 Soltis DE, Buggs RJ, Doyle JJ, Soltis PS. 2010. What we still don't know about polyploidy. *Taxon*
929 1387–1403.
- 930 Soltis DE, Soltis PS. 1999. Polyploidy: recurrent formation and genome evolution. *Trends Ecol*
931 *Evol* **14**: 348–352.
- 932 Soltis DE, Visger CJ, Soltis PS. 2014. The polyploidy revolution then... and now: Stebbins
933 revisited. *Amer J Botany* **101**: 1057–1078.
- 934 Soltis PS, Soltis DE. 2000. The role of genetic and genomic attributes in the success of polyploids.
935 *Proc Natl Acad Sci* **97**: 7051–7057.
- 936 Soltis PS, Soltis DE. 2012. *Polyploidy and genome evolution*. Springer.
- 937 Stamatakis A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
938 phylogenies. *Bioinformatics* **30**: 1312–1313.
- 939 Steige K, Reimegård J, Rebernic CA, Köhler C, Scofield DG, Slotte T. 2016. The role of trans-
940 posable elements for gene expression in capsella hybrids and allopolyploids. *bioRxiv* doi:
941 10.1101/044016.
- 942 te Beest M, Le Roux JJ, Richardson DM, Brysting AK, Suda J, Kubešová M, Pyšek P. 2011. The
943 more the better? the role of polyploidy in facilitating plant invasions. *Ann Botany* **109**: 19–45.
- 944 Terhorst J, Kamm JA, Song YS. 2017. Robust and scalable inference of population history from
945 hundreds of unphased whole genomes. *Nature Gen* **49**: 303–309.
- 946 Thompson A, Zakon HH, Kirkpatrick M. 2016. Compensatory drift and the evolutionary dynamics
947 of dosage-sensitive duplicate genes. *Genetics* **202**: 765–774.
- 948 Treier UA, Broennimann O, Normand S, Guisan A, Schaffner U, Steinger T, Müller-Schärer H.
949 2009. Shift in cytotype frequency and niche space in the invasive plant *Centaurea maculosa*.
950 *Ecology* **90**: 1366–1377.
- 951 Vaser R, Adusumalli S, Leng SN, Sikic M, Ng PC. 2016. SIFT missense predictions for genomes.

- 952 *Nature Prot* **11**: 1–9.
- 953 Weiss-Schneeweiss H, Emadzade K, Jang TS, Schneeweiss G. 2013. Evolutionary consequences,
954 constraints and potential of polyploidy in plants. *Cytog Genome Res* **140**: 137–150.
- 955 Woodhouse MR, Cheng F, Pires JC, Lisch D, Freeling M, Wang X. 2014. Origin, inheritance,
956 and gene regulatory consequences of genome dominance in polyploids. *Proc Natl Acad Sci* **111**:
957 5283–5288.
- 958 Yang J, Liu D, Wang X, Ji C, Cheng F, Liu B, Hu Z, Chen S, Pental D, Ju Y, et al.. 2016. The
959 genome sequence of allopolyploid *Brassica juncea* and analysis of differential homoeolog gene
960 expression influencing selection. *Nature Gen* **48**: 1225–1232.
- 961 Yang Z. 1997. PAML: a program package for phylogenetic analysis by maximum likelihood.
962 *CABIOS* **13**: 555–556.