1 **Coalescent-based analysis distinguishes between allo- and autopolyploid origin in**

2 **shepherd's purse** *(Capsella bursa-pastoris).*

3

4 Kate St.Onge*[1], John Paul Foxe*[4], Li Junrui*[5,7], Li Haipeng[5], Karl Holm[1], Pádraic

5 Corcoran[2], Tanja Slotte[2], Martin Lascoux[1,5,†], Stephen Wright[3,6†]

6

7 *These authors contributed equally to this work

8

9 [1]Department of Evolutionary Functional Genomics, Uppsala University, 752 36 Uppsala,

10 Sweden

11 [2]Department of Evolutionary Biology, Uppsala University 752 36 Uppsala, Sweden

12 [3]Department of Ecology and Evolutionary Biology, University of Toronto, 25 Willcocks

13 Street, Toronto, On M5S 3B2, Canada

14 [4]Department of Biology, York University, 4700 Keele St. Toronto, Ontario M3J 1P3,

15 Canada

16 [5]Laboratory of Evolutionary Genomics, CAS Key Laboratory of Computational Biology,

17 CAS-MPG Partner Institute for Computational Biology, Chinese Academy of Sciences,

18 Shanghai, China

19 [6]Centre for Analysis of Genome Evolution and Function, University of Toronto

20 [7]Graduate School of the Chinese Academy of Sciences, Beijing 100039, China

21

22 [†]Corresponding authors: Stephen Wright and Martin Lascoux

23 E-mail: Stephen.Wright@utoronto.ca, Martin.Lascoux@ebc.uu.se

24 Phone: +1 (416) 946-8508, +46 (18) 4716416

25 Fax: +1 (416) 978-5878, +46 (18) 4716457

26

27 Keywords: autopoplyploidy, coalescent, Brassicaceae, Approximate bayesian

28 computation, IM model

29 Running title: Autopolyploid speciation in *Capsella bursa-pastoris*

30

31

32

33   ABSTRACT

34

35   Polyploidization plays an important role in plant speciation. The most recent estimates

36   report that up to 15% of angiosperm speciation events and 31% in ferns are accompanied

37   by changes in ploidy level. Polyploids can arise either through autopolyploidy, when the

38   sets of chromosomes originate from a single species, or through allopolyploidy, when

39   they originate from different species. In this study we used two different coalescent-based

40   methods to determine the date and mode of the polyploidization event that led to the

41   tetraploid cosmopolitan weed, *Capsella bursa-pastoris*. We sampled 78 *C. bursa-pastoris*

42   accessions, and 53 and 43 accessions from the only two other members of this genus, *C.*

43   *grandiflora* and *C. rubella,* respectively, and sequenced these accessions at 14 unlinked

44   nuclear loci with locus-specific primers in order to be able to distinguish the two

45   homeologues in the tetraploid. A large fraction of fixed differences between

46   homeologous genes in *C. bursa-pastoris* are segregating as polymorphisms in *C.*

47   *grandiflora*, consistent with an autopolyploid origin followed by disomic inheritance. To

48   test this, we first estimated the demographic parameters of an isolation-with-migration

49   model in a pairwise fashion between *C. grandiflora* and both genomes of *C. bursa-*

50   *pastoris* and used these parameters in coalescent simulations to test the mode of origin of

51   *C. bursa-pastoris*. Secondly we used Approximate Bayesian Computation to compare an

52   allopolyploid and an autopolyploid model. Both analyses led to the conclusion that *C.*

53   *bursa-pastoris* originated less than one million years ago by doubling of the *C.*

54   *grandiflora* genome.

55

56

3

57    INTRODUCTION

58

59    Polyploidy, or whole genome duplication, is widespread in plants. Polyploidy occurs in

60    virtually all groups of vascular plants including ferns, mosses and algae (Otto and

61    Whitton 2000). The most recent estimate of the prevalence of polyploids using

62    phylogenetic data reports that 15% of speciation events in angiosperms and 31% in ferns

63    are accompanied by changes in ploidy level (Wood et al. 2009), over 4 times higher than

64    previous estimates of 2-4% in angiosperms and 7% in ferns (Otto and Whitton 2000).

65    Analysis of fossil and genomic data estimate that 47%-100% of angiosperms have a

66    polyploidy event at some point in their histories (Masterson 1994; Cui et al. 2006) and

67    genomic studies have revealed that chromosomally diploid plant species, such as

68    *Arabidopsis*, *Populus*, *Vitis*, and *Oryza* went through one or many rounds of

69    polyploidisation during their evolution (e.g. Fawcett et al. 2009).

70

71    Typically, polyploids are divided into two categories based on their mode of origin,

72    allopolyploids and autopolyploids. Allopolyploids have two full genome complements

73    originating from two different species. These polyploids are expected to display disomic

74    inheritance and form bivalents at meiosis, although disomic inheritance is not a strict

75    indicator of allopolyploidy. All polyploid Brassica species studied so far are

76    allopolyploids: *Brassica carinata*, *B. juncea*, and *B. napus* are tetraploids created from

77    hybridization of the species *B. nigra*, *B. rapa* and *B. oleracea* in different combinations

78    (U 1935). On the other hand, autopolyploids result from genome doubling within a

79    species. Genome doubling can occur spontaneously or following the fusion of unreduced

80    diploid gametes.  Examples of autopolyploid plants include alfalfa and potato, and it was

81    recently shown that the domesticated apple had an ancient autopolyploid origin (Velasco

82    et al. 2010). Autopolyploids are typically expected to display polysomic inheritance and

83    form multivalents at meiosis, although the generality of this rule has started to be

84    questioned. Some autopolyploids are known to display disomic inheritance and this is

85    probably more frequent than previously assumed (Soltis et al. 2010). However, how this

86    occurs or how quickly disomic inheritance can evolve from polysomic inheritance is still

87    poorly known (Cifuentes et al. 2010), although indirect evidence suggests that it can take

88   place rapidly (Parisod et al. 2010). If disomic inheritance follows a period of polysomic

89   inheritance, divergence times estimated from duplicated genes will reflect the time of

90   onset of disomic inheritance, rather than the time of polyploidization (Gaut and Doebley

91   1997).

92

93   Determining the origin of polyploid species is an important aspect of speciation genetics

94   and is central to our understanding of the mechanisms of formation of polyploids. While

95   issues such as multiple origins of polyploid species, extinction of parental lineages and

96   sampling of standing variation from progenitor species complicate this task (Doyle and

97   Egan 2009; Soltis et al. 2010), recent advances in coalescent modeling have meanwhile

98   facilitated it (Noor and Feder 2006; Becquet and Przeworski 2007; Hey and Nielsen

99   2007; Hey 2010).   In particular, models of isolation-with-migration (IM) allow the

100  differentiation of ancestral polymorphism from introgression and provide statistically

101  sound estimates of divergence events (Wakeley and Hey 1997; Nielsen and Wakeley

102  2001). Using these models, diploid speciation processes have been studied in many

103  organisms including *Drosophila* (Wang et al. 1997; Hey and Nielsen 2007), *Arabidopsis*

104  (Ramos-Onsins et al. 2004), *Oryza* (Zhang and Ge 2007) and *Capsella* (Foxe et al. 2009).

105  However, the use of coalescent-based models to study polyploidy and speciation has so

106  far been limited with the notable exception of the studies of Jakobsson et al. (2006) in *A.*

107  *suecica*, where an allopolyploid origin from *A. thaliana* and *A. arenosa* was known, and

108  of *Capsella bursa-pastoris* in Slotte et al. (2008) and of *Arabidopsis lyrata* ssp.

109  *kamchatica* of Taiwan in Wang et al. (2010).

110

111  The genus *Capsella* belongs to the mustard family (*Brassicaceae*) and is an attractive

112  model genus because it is a young genus that contains few species with different mating

113  systems and ploidy levels. The genus includes three species: *C. bursa-pastoris* (L.)

114  Medik., a selfing tetraploid that displays a disomic inheritance and two diploid species,

115  the outcrosser *C. grandiflora* (Fauché & Chaub.) Boiss., and the selfer *C. rubella* Reuter

116  (Shull 1929; Hurka and Neuffer 1997). Previous studies suggested that *C. grandiflora* is

117  ancestral to *C. bursa-pastoris* and *C. rubella* (Hurka and Neuffer 1997) and more recent

118  findings confirmed that *C. rubella* diverged from *C. grandiflora* as recently or more

119    recently than the Last Glacial Maximum  (LGM, 18,000 years ago, St.Onge et al. 2011;

120    13,500 years ago, Foxe et al. 2009). *C. bursa-pastoris* has a worldwide distribution that

121    can partly be explained anthropogenically. In contrast to *C. grandiflora* and *C. rubella*, *C.*

122    *bursa-pastoris* can be found on each continent and thrives in a wide range of climates

123    (Hurka and Neuffer 1997).

124

125    It is still unknown if *C. bursa-pastoris* is of autopolyploid or allopolyploid origin, and

126    both possibilities have been suggested in previous work. Early isozyme electrophoresis

127    indicated that *C. bursa-pastoris* shared alleles with both *C. grandiflora* and *C. rubella*

128    and was hence thought to be an allopolyploid between these two species (Hurka et al.

129    1989). Later, evidence from restriction site variation in the chloroplast genome indicating

130    that *C. rubella* was a more recently derived species led to the suggestion that *C. bursa-*

131    *pastoris* was an ancient autopolyploid of *C. grandiflora* (Hurka and Neuffer 1997),

132    despite the fact that *C. bursa-pastoris* displays disomic inheritance. Most recently,

133    phylogenetic analysis suggested again that *C. bursa-pastoris* may be an allopolyploid,

134    although not between *C. grandiflora* and *C. rubella* (Slotte et al. 2006).

135

136    A major limitation of these past studies is that they lack comprehensive data from all

137    three *Capsella* species. In particular, the lack of large population data from *C.*

138    *grandiflora*, the species of the genus known to harbor the most genetic variation, makes it

139    difficult to conclusively determine the polyploid origin of *C. bursa-pastoris*. Here, we use

140    DNA sequence data from 14 unlinked nuclear loci from large samples of all three

141    *Capsella* species. In the absence of linkage data, assigning homeologues to particular

142    genome copies in *C. bursa-pastoris* is not possible. To address this, we took the extreme

143    possibility that more divergent copies from *C. grandiflora* all come from the same

144    lineage. Since this would be most likely under an allopolyploid model, this allows us to

145    explicitly test the plausibility of this model compared to autopolyploidy. To compare the

146    fit of the data to allopolyploid vs. autopolyploid models of speciation we used a novel

147    coalescent-based approach. First, we estimate the parameters of an Isolation-with-

148    Migration model for pairs of species and then use these parameters in coalescent

149    simulations to test the fit of the data to different models. Second, we use Approximate

150     Bayesian Computation (ABC, Beaumont 2010) to implement a two-split model and test

151     our two competing hypotheses, the allopolyploid and the autopolyploid models. As

152     Figure 1 shows, if *C. bursa-pastoris* has an autopolyploid origin we would expect the

153     divergence time between the two homeologues to be as recent as, or (if there was an

154     initial period of polysomic inheritance) more recent than the time at which *C. bursa-*

155     *pastoris* derived from *C. grandiflora*, suggesting a simple way to test whether *C. bursa-*

156     *pastoris* is of auto- or allopolyploid origin.

157

158     MATERIALS AND METHODS

159

160     ***Sample collection***

161     Genetic data was collected from 78 accessions of *C. bursa-pastoris* from China, Taiwan,

162     Israel and Europe, 43 accessions of *C. rubella* from Africa, South America, Europe and

163     Israel and 53 accessions of *C. grandiflora* from Greece, covering a large portion of the

164     narrow distribution of this species. Because this study focuses on the origin of *C. bursa-*

165     *pastoris*, we have excluded samples from the Americas as Capsella species are a recent

166     introduction there (Hurka and Neuffer, 1997). All our accessions come from natural

167     populations from which we have collected seeds. In this study, we used a single accession

168     per sampled population in most cases (see Table S1). Genetic data was also collected

169     from one accession of *Neslia paniculata*, which was used as an outgroup in some

170     analyses. *Neslia* is more recently diverged from *Capsella* than *Arabidopsis* (Bailey et al.

171     2006), providing a closer outgroup for inferences about *Capsella* divergence. Plants were

172     grown in standard long-day conditions and DNA was extracted from fresh tissue of each

173     individual using the QIAgen DNeasy Plant Mini Kit (QIAGEN, Valencia, California,

174     USA). Accessions and their geographic origins are given in Table S1.

175

176     ***PCR and Sequencing***

177     Fourteen gene fragments were selected for sequencing in this panel of individuals. These

178     genes were found to be single copy in both diploids and duplicated in *C. bursa-pastoris,*

179     as expected in a tetraploid. For four of the loci (At1g77120 (ADH), At5g10140 (FLC),

180     At4g00650 (FRI) and At4g02560 (LD)), PCR primers for the diploid species and

181    homeologue-specific primers for *C. bursa-pastoris* were designed as described by Slotte

182    et al. (2006) and Slotte et al. (2008). For eight genes (At1g01040; At1g03560,

183    At1g15240, At1G65450, At2g26730, At4g14190, At5g51670, At5g53020) primers for

184    the diploid species were designed as described in Ross-Ibarra et al. (2008) and Foxe et al.

185    (2009). For two additional loci (At2g18790 (PHYB) and At5g42800 (DFR)), primers

186    were designed following a similar strategy. For all loci, initial primers were designed

187    using Primer3 version 0.4.0 (Rozen and Skaletsky 2000) or PrimerQuest (Integrated

188    DNA Technologies, Inc.) to amplify between 400-1000 bps using the *A. thaliana* genome

189    sequence. The *A. thaliana* sequences were aligned to other Brassicaceae sequences when

190    available to identify conserved regions. Both forward and reverse strands of the

191    amplicons were sequenced directly at Lark Technologies (Houston, Texas), the Genome

192    Quebec Innovation Centre (McGill University, Canada) or the Macrogen sequencing

193    facility in Korea (Macrogen, Korea). Sequences were aligned and checked manually for

194    heterozygous sites using either Sequencher version 4.7 (Gene Codes, Ann Arbor, MI) and

195    Genedoc (Nicholas *et al.* 1997) or Codoncode Aligner version 2.0.6 (CodonCode,

196    Dedham, MA). To differentiate the two homeologues of *C. bursa-pastoris*, the resulting

197    sequences were used to design new homeologue-specific primers as in Slotte et al.

198    (2006). In particular, we designed primers specific to SNPs showing fixed

199    'heterozygosity' amongst all of our samples, representing fixed SNP differences between

200    homeologues. Each homeologue-specific amplicon was then sequenced directly and

201    aligned as above. Based on direct sequencing of these samples only a single haplotype

202    per homeologue was found for all of our primer pairs, implying homozygosity of our

203    inbred samples. Details of the new primers for this study are shown in supplementary file

204    S1. Sites with indels were removed before proceeding with analysis. The program

205    PHASE 2.1 (Stephens et al. 2001), implemented in DnaSP 5.0 (Librado and Rozas 2009)

206    was used to infer haplotypes in *C. grandiflora*. Additionally, each gene fragment was

207    aligned with the homologous *A. thaliana* gene to infer the ancestral state of polymorphic

208    sites. Loci and accessions where only one homeologue amplified were removed. New

209    nucleotide sequences generated in this study that are greater than 200bp in length have

210    been deposited in GenBank (accession numbers JQ418636-JQ419488). Complete

211  sequence alignments, and sequence data from regions less than 200bp in length, are

212  available upon request to the corresponding authors.

213

214  ***Summary statistics and estimation of species trees***

215  A central challenge for our study is the difficulty in assigning homeologous genes to

216  separate genomes of origin, designated as the *C. bursa-pastoris* A and B genomes.

217  Homeologues were assigned to A and B genomes based upon the minimum number of

218  synonymous substitutions between *C. grandiflora* and each homeologue as estimated

219  using DnaSP version 5.0 (Librado and Rozas 2009). The most distant homeologue was

220  assigned to the B genome while the other was assigned to A (Table S2; similar to Slotte

221  *et al.* 2006 and Slotte *et al.* 2008; however, in these papers classification was based on all

222  sites and *C. rubella* was used instead of *C. grandiflora*). These putative genomes were

223  analysed separately for all subsequent analyses. Importantly, this classification effectively

224  biases our analysis toward rejecting the hypothesis of the autopolyploid origin of *C.*

225  *bursa-pastoris*. In particular, if the allopolyploid model is correct, the A and B

226  homeologues likely represent distinct genomes with different parental origins, while

227  under the autopolyploid model their difference is simply due to stochastic noise in the

228  coalescent process, and the sorting does not reflect genome structure.

229

230  Classic genetic diversity summary statistics π (Tajima 1983) and Tajima's D (Tajima

231  1989) were calculated for synonymous sites in each species using a modified version of

232  the Polymorphorama perl code

233  (http://ib.berkeley.edu/labs/bachtrog/data/polyMORPHOrama/polyMORPHOrama.html)

234  written by D. Bachtrog (UC Berkeley) and P. Andolfatto (Princeton University). The

235  joint frequency spectra of derived polymorphic variants and the number of shared derived

236  polymorphisms, unique polymorphisms, and fixed differences between each of the four

237  genomes (Wakeley and Hey 1997) were calculated separately in a pairwise fashion using

238  a Perl script written by S. Wright and a C program written by J. Li.

239

240  The molecular phylogenetic program BEST v. 1 (Bayesian estimation of species trees)

241  (Liu 2008), which implements a Bayesian hierarchical model while accounting for the

242    presence of deep coalescent events, was used to estimate the *Capsella* genus species tree

243    using our multi-locus dataset (Liu 2008). Models within the BEST program assume (i)

244    No population substructure within each population, (ii) No gene flow after species

245    divergence and, (iii) No recombination within loci. Some of these assumptions, in

246    particular the last one, will likely be violated. For example, recombination will be present

247    in *C. grandiflora* and will make the length of terminal branches and the total branch

248    length larger, and the time to the most recent common ancestor smaller (Schierup and

249    Hein 2000). The program reportedly works best using concatenated alignments with little

250    missing data. Consequently, we ran BEST using the 7 loci in this dataset that had the

251    most consistent sampling of individuals across loci (At1g03560, At1g15240, At1g65450,

252    At2g26730, At4g14190, At5g51670 and At5g53020). Alignments were concatenated

253    using MacClade version 4.08 (available from http://macclade.org/). BEST was run in two

254    ways, once using *A. thaliana* as an outgroup and again including both *A. thaliana* and *N.*

255    *paniculata* (where available). In each case BEST was run twice, with 4 chains for a

256    maximum of 2 million generations, with a burnin of 200,000 generations, sampling every

257    100 generations.

258

259    ***MIMAR and coalescent simulations***

260    A first test of the null hypothesis that *C. bursa-pastoris* is an autopolyploid of *C.*

261    *grandiflora* was done by first estimating the parameters of an isolation-with-migration

262    model using the program MIMAR (Becquet and Przeworski 2007), and then performing

263    coalescent simulations based on these parameters to test the null hypothesis (Hudson

264    2002). Because previous studies showed that *C. rubella* diverged very recently from *C.*

265    *grandiflora* (Foxe et al. 2009; St.Onge et al. 2011), *C. rubella* was initially not included

266    in this analysis. Furthermore, sites with >2 segregating bases were also excluded.

267    MIMAR simulations were run in a pairwise fashion using *C. bursa-pastoris* A, *C. bursa-*

268    *pastoris* B and *C. grandiflora* and allowing for three different models of migration

269    between genomes: 1) absence of migration 2) symmetrical migration and 3) asymmetrical

270    migration. Additionally, all analyses were run both with the ancestral effective

271    population size unconstrained or assumed to be identical to the effective size of *C.*

272    *grandiflora*. Prior limits for all parameters can be found in Table S3; these priors were set

273  based on short initial runs with very wide priors. The program was run as described in

274  Foxe et al. (2009), with the exception that each simulation was run for a total of 10,080

275  min. (1 week). We note that the model implemented by MIMAR does not allow a

276  temporary reduction in Ne at the polyploid origin and so effectively allows multiple

277  polyploid origins. Thus, inferences of effective population sizes should be considered a

278  weighted average since divergence, rather than a direct estimate of the number of

279  founders.

280

281  Because MIMAR simulations only model two taxa at a time, it does not on its own

282  provide an explicit test of the mode of polyploid speciation. We therefore conducted

283  coalescent simulations using MIMAR parameter estimates under models of both

284  autopolyploidy and allopolyploidy. These models are depicted in Figure 1. Importantly,

285  the differences between these models are the split times between *C. grandiflora* and the

286  two genomes of *C. bursa-pastoris*. Under the autopolyploid model all three divergence

287  times are the same, or the divergence time of the A and B homeologues is shorter than the

288  time of either to *C. grandiflora*, if there was a period of polysomic inheritance. Under

289  allopolyploidy the divergence time between *C. grandiflora* and *C. bursa-pastoris* B and

290  between *C. bursa-pastoris* A and B are much longer than that between *C. grandiflora* and

291  *C. bursa-pastoris* A.  To model autopolyploid speciation, we used the inferred divergence

292  time from MIMAR runs considering the two homeologues of *C. bursa-pastoris*, since this

293  should provide an estimate of the lower bound for the time of autopolyploid origin. Under

294  the allopolyploid model, the A and B copies in *C. bursa-pastoris* truly represent distinct

295  genomes with different parental origins, and we used the two inferred divergence times

296  from the MIMAR runs of *C. grandiflora* with the two distinct homeologue sets.

297

298  To compare our simulated data to our empirical data we used summary statistics

299  introduced by Wakeley and Hey (1997) for each locus: the number of polymorphisms

300  specific to the samples from populations 1 and 2 (called s1 and s2, respectively), where

301  the population pairs correspond to *C. grandiflora/C. bursa-pastoris* A, *C. grandiflora/ C.*

302  *bursa-pastoris* B and *C. bursa-pastoris* A */C. bursa-pastoris* B, the number of shared

303  polymorphisms between two samples (sp), and the number of sites fixed in either sample

304    (f1 and f2, depending on which of the two species carries the ancestral state). We

305    conducted simulations under both auto- and allopolyploidy models using the program *ms*

306    (Hudson 2002) and the demographic parameters inferred with MIMAR. Namely, in the

307    autopolyploid model we used $T_1$ as the divergence time between the two genomes and *C.*

308    *grandiflora* (Figure 1) and in the allopolyploid model we used $T_2$' for the divergence

309    between *C. bursa-pastoris* A and *C. grandiflora* and $T_2$" for the divergence between *C.*

310    *bursa-pastoris* B and *C. grandiflora* (Figure 1). For each of the 14 genes we assumed that

311    10 chromosomes were sampled in each species and ran 10,000 simulations. We then

312    calculated shared and fixed sites for each run and the mean over runs for each locus

313    (additional information is available in Supplementary file S2 where the same analysis

314    was carried out but considering both *C. grandiflora* and *C. rubella*).

315

316    Using these simulations we determined which of the summary statistics described above

317    were informative in differentiating the two models. We found that unique polymorphisms

318    (s1 and s2) did not differ between the two models. This may seem intuitive given that

319    unique polymorphisms mostly reflect genealogies within that species, and therefore give

320    limited information about speciation and divergence between the species in a genus. We

321    therefore did not use these sites further. In contrast fixed sites (f1 and f2, depending on

322    which of the two species carries the ancestral state) and shared polymorphisms (sp) did

323    differ between ploidy models. Again this is intuitive as fixed differences correspond to

324    mutations that happened in the early stages of speciation and are closely associated to

325    divergence time, whereas shared polymorphisms, assuming they represent shared

326    ancestral polymorphism and not recent introgression, represent polymorphism that were

327    segregating in the ancestor and therefore give information about ancestral effective

328    population sizes and divergence times. To make use of these two informative statistics we

329    calculated their difference in the following way. If f1(C.bp B, C.g) and f1(C.bp A, C.g)

330    are the number of fixed sites between *C. grandiflora*  and *C. bursa-pastoris* B and *C.*

331    *bursa-pastoris* A, respectively then their difference, fix_diff = f1(C.bp A, C.g) - f1(C.bp

332    B, C.g). For convenience we used f1 to define fix_diff but the conclusions were the same

333    when we used f2 (data not shown). If  sp(C.bp A, C.g) and sp(C.bp B, C.g) represent the

334    number of shared sites between *C. grandiflora* and *C. bursa-pastoris* A and *C. bursa-*

335     *pastoris* B, respectively, then the difference between them, shared_diff = sp(C.bp A, C.g)

336     - sp(C.bp B, C.g).

337

338     We calculated the two differences defined above in our observed data and used them to

339     test the null hypothesis that *C. bursa-pastoris* is an autopolyploid. Essentially we used a

340     goodness-of-fit test to test the fit of our null hypothesis to the empirical data. We

341     compared the observed values of the mean of fix_diff and shared_diff over the fourteen

342     loci with the distribution of the same mean for the 10,000 simulation runs obtained under

343     the autopolyploid model. When calculating the *p*-values for the autopolyploid model we

344     used a two-tailed test for both test statistics. The *p*-value is therefore the fraction of

345     simulations in which the absolute value of the mean is higher than the observed mean,

346     with p-values of <0.05 indicating that our empirical values lie in the tails of the

347     simulated distributions. We also assessed the fit of our alternative hypothesis,

348     allopolyploidy, in a similar way except that in this case the tests were one-tailed because

349     of the bias we created in our dataset. Therefore the *p*-value of the fixed_diff statistics is

350     the fraction of simulations in which the mean is higher than the observed mean while in

351     the case of shared_diff it is the fraction of simulations in which the mean is lower than

352     the observed mean.

353

354     ***Approximate Bayesian Computation***

355     An Approximate Bayesian Computation (ABC) analysis was used to evaluate two-split

356     models resulting in *C. grandiflora* and the A and B genomes of *C. bursa-pastoris*. This

357     analysis was performed using the program Seqlib-1.6 (De Mita et al., 2007)

358     (http://sourceforge.net/projects/seqlib/), on the silent sites of the dataset. Because *C.*

359     *rubella* has recently evolved from *C. grandiflora* and the variation in the species is more

360     or less a subset of the variation found in *C. grandiflora* (Foxe et al. 2009; Guo et al. 2009,

361     St Onge et al. 2011), we chose not to include it in the present analysis (see results and

362     discussion). We evaluated two possible arrangements of coalescent events involving the

363     three lineages 1) the A and B genomes of *C. bursa-pastoris* coalesce first, followed by

364     this lineage coalescing with *C. grandiflora* 2) *C. grandiflora* coalesces first with the A

365     genome of *C. bursa-pastoris*, followed by coalescence with the B genome. Model 1

13

represents an autopolyploidy event, where the divergence times of both *C. bursa-pastoris* genomes from to *C. grandiflora* are the same (Figure 1). Under model 2, the *C. bursa-pastoris* B genome is more diverged from *C. grandiflora* than *C. bursa-pastoris* A, representing an allopolyploidy event (Figure 1). These models have 9 parameters: the population mutation rates of each lineage, θ1, θ2 and θ3 where θ2 and θ3 are relative to θ1, the population recombination rate, ρ, a combined migration rate between all lineages, the dates of each divergence event, where the second event is additive to the first and the population sizes after each coalescent event, relative to θ1. It should be noted again that, for each gene, the *C. bursa-pastoris* allele most divergent from *C. grandiflora* was assigned to the B genome, effectively biasing our analysis towards model 2.

The ABC analysis of our two two-split models was performed using a set of 13 summary statistics; the number of shared, fixed and unique polymorphisms in all possible configurations with the three populations. We first performed initial runs with 1,000,000 samples using wide priors (Table S3). Using the local linear regression method described by Beaumont *et al.* (2002), 0.1% of the samples best fitting our empirical data were selected and used to create a prior for the ABC run. This allows us to explore the region of high probability identified in the initial run. 500,000 samples are taken in the ABC run, and 0.2% of the best fitting samples were used to estimate model parameters. A goodness-of-fit (GoF) test was used to validate the results of the ABC analysis. This test consisted of two sets of simulations, one using the point estimates for each parameter estimated in the ABC and one using the posterior distributions of each parameter. Further details on the goodness-of-fit test are in the Supplementary File S3..

A second analysis was performed using the same method but with only *C. bursa-pastoris* accessions from China. This was done to assess the influence of putatively introgressed alleles from *C. rubella* that only occurred in Europe (Slotte et al. 2008).

### Testing for interlocus gene conversion in C. bursa-pastoris

We followed the approach of Slotte and colleagues (2008) to test for gene conversion between homeologues. In particular, we calculated the minimum number of

397 recombination events, Rm, between homeologues (Hudson and Kaplan 1985) using

398 DNAsp 5.0 (Librado and Rozas 2009), and tested for gene conversion using the geneconv

399 software (Sawyer, 1989).

400

401 RESULTS

402

403 *Patterns of polymorphism and phylogeny*

404 Synonymous site diversity, measured as $\pi$, was higher in *C. grandiflora* than in *C.*

405 *rubella* and *C. bursa-pastoris* A and B; median values were 0.028 for *C. grandiflora,*

406 while they were zero for the latter two species (Figure S1). This is in agreement with

407 expectations based on the respective mating systems of these species and previous studies

408 (Slotte et al. 2008; Foxe et al. 2009; St.Onge et al. 2011). In particular the low level of

409 nucleotide diversity observed in *C. rubella* is consistent with the presence of a severe

410 population bottleneck associated to the shift to selfing (Foxe et al. 2009; Guo et al. 2009,

411 St.Onge et al. 2011). The reduction in diversity seen in both *C. bursa-pastoris* A and B

412 may also be the result of a recent bottleneck at speciation and transition to selfing.

413 However, all species showed a high variance in diversity, with *C. rubella* showing the

414 most extreme variance, with synonymous $\pi$ values varying from 0 to the extremely high

415 value of 0.15 for the DFR (At5g42800) locus. Resequencing of the full *C. rubella*

416 genome and mRNA resequencing indicate that the high variance in diversity is a

417 genomewide characteristic of the species (Wright et al, unpublished; D. Weigel, pers

418 comm). This locus also showed high, but less elevated, polymorphism in *C. grandiflora*

419 (0.08). Excluding DFR, the average synonymous diversity was 0.027 in *C. grandiflora*,

420 0.004 in *C. rubella*, 0.003 in *C. bursa-pastoris A*, and 0.003 in *C. bursa-pastoris B*. The

421 average Tajima's D values at synonymous sites were negative for *C. bursa-pastoris* A (-

422 0.19) and B genomes (-0.9), possibly reflective of recent population expansion. In *C.*

423 *grandiflora*, synonymous Tajima's D was close to zero (-0.08), consistent with previous

424 conclusions suggesting that this species is close to demographic equilibrium (Foxe et al.

425 2009; St.Onge et al. 2011). In *C. rubella*, synonymous Tajma's D was slightly negative (-

426 0.2.).

427

15

428    The minimum number of synonymous substitutions was calculated in a pairwise fashion

429    between *C. grandiflora* and *C. bursa-pastoris* A and B (Table S2). Under an

430    allopolyploidy model we would expect a higher number of synonymous substitutions

431    between *C. grandiflora* and *C. bursa-pastoris* B than between *C. grandiflora* and *C.*

432    *bursa-pastoris* A. We do of course observe this since we have used the minimum number

433    of synonymous substitutions to *C. grandiflora* to assign alleles to the A and B genomes,

434    assigning the more distant allele to the B genome. However, for most loci, there is only a

435    slight difference in this quantity between homeologues, suggesting that the two

436    homeologues are nearly equal in their distance from standing *C. grandiflora* haplotype

437    variation. Furthermore, the minimum number of synonymous substitutions between the

438    two *C. bursa-pastoris* genomes is higher than either comparison with *C. grandiflora* as

439    previously observed (Slotte et al. 2006).  Likewise, we observe 29 fixed synonymous

440    differences between *C. bursa-pastoris* A and B compared with 2 between *C. grandiflora*

441    and *C. bursa-pastoris* A and 19 between *C. grandiflora* and *C. bursa-pastoris* B (Figure

442    2). The cause of this large difference in fixed sites observed between the *C. bursa-*

443    *pastoris* genomes is likely their small effective population size causing alleles to drift to

444    fixation quickly. On the other hand, the large effective population size of *C. grandiflora*

445    would allow the maintenance of many shared alleles with both *C. bursa-pastoris*

446    genomes.

447

448    Looking at the pattern of fixed differences between homeologues in *C. bursa-pastoris*

449    reveals a striking pattern; 43% of fixed differences between homeologues are segregating

450    with our *C. grandiflora* sample. Furthermore, if we restrict this to the 7 genes with large

451    *C. grandiflora* samples (>20 chromosomes), this fraction increases to 52%. This retention

452    of *C. grandiflora* polymorphism as fixed differences between homeologues in *C. bursa-*

453    *pastoris* is consistent with an autopolyploid model, where distinct haplotypes sampled

454    from the ancestral *C. grandiflora* population were 'frozen' as gene duplicates during

455    polyploidization. Under this scenario, the remaining fixed differences would reflect rare

456    SNPs not sampled in *C. grandiflora* and/or new mutations and fixation events following

457    speciation. Considerably fewer fixed differences between homeologues are still

458    segregating in *C. rubella* (20%).

459

460   In terms of identical haplotypes, we identified identical haplotypes between *C. bursa-*

461   *pastoris* A and the other two species for all but three of our loci. Of the genes showing

462   haplotype sharing 4 loci showed sharing with both species, 4 showed sharing only with

463   *Capsella rubella*, and 2 showed sharing with *C. grandiflora* alone. Although the excess

464   haplotype sharing in *C. rubella* is consistent with the inference of introgression (Slotte et

465   al. 2008), it is important to note that the requirement of inferring phase in *C. grandiflora*

466   and extensive recombination may erode some of the signal of haplotype sharing. Indeed,

467   for the seven loci where we have relatively large *C. grandiflora* sample sizes for better

468   inferences of phased haplotypes, only one locus shows *C. rubella* only haplotype sharing,

469   and for this one it is only a single *C. rubella* individual that shows the shared haplotype.

470

471   We estimated the species tree of the *Capsella* genus using the program BEST, which

472   implements a Bayesian hierarchical model while accounting for the presence of deep

473   coalescences (Liu 2008). The analysis was performed twice, first by including *A. thaliana*

474   as an outgroup (Figure S2-A) and second by including both *A. thaliana* and *Neslia*

475   *paniculata*, where available, as outgroups (Figure S2-B). *C. grandiflora* was not shown to

476   be more closely related to either *C. bursa-pastoris* A or B in either of the resulting trees.

477   In fact, the tree resulting from the first analysis is the expected tree under an

478   autopolyploidy model where the branch lengths between the two *C. bursa-pastoris*

479   genomes and *C. grandiflora* are equal. Despite biasing our analysis toward the

480   allopolyploidy model our results thus lend support to the autopolyploidy hypothesis.

481

482

483   ***Demographic model fitting: MIMAR and ms simulations***

484   We used the program MIMAR (Becquet and Przeworski 2007) to fit models of isolation

485   with migration in a pairwise fashion to *C. grandiflora* and *C. bursa-pastoris* A and B.

486   The model assumes that a single ancestral population of size $N_a$ splits into two

487   descendant populations at time t, and the two descendant populations have distinct

488   population sizes. Models including symmetric migration, asymmetric migration and no

489   migration between the two derived populations were analysed for all three species pairs.

490    The results, however, show no evidence for migration between *C. grandiflora* and *C.*

491    *bursa-pastoris*, so we only report the results from analyses assuming no migration

492    between descendant populations.

493

494    Mimar runs that included gene flow, both between the *C. bursa-pastoris* homeologues

495    and from *C. bursa-pastoris* to *C. grandiflora*, showed modes that approached zero (Table

496    S4), providing little evidence for extensive gene conversion between homeologues and/or

497    introgression from *C. grandiflora* following divergence. We therefore focus the

498    presentation of the results on the no-migration model, although all results are reported in

499    Table S4. C. *bursa-pastoris* A and B show a 5- and 7- fold decrease in effective

500    population size, respectively, compared with *C. grandiflora* (Figure 3A; Table S4), with

501    effective population sizes around 50,000-80,000 for *C. bursa-pastoris* A and B and

502    values around 410,000 for *C grandiflora*, if we assume a mutation rate of $1.5 \times 10^{-8}$

503    /site/year (Koch et al. 2000). The estimated time of divergence between each pair of

504    genomes were 278,000 years between *C. grandiflora* and *C. bursa-pastoris* A, 1.1 million

505    years between *C. grandiflora* and *C. bursa-pastoris* B and 563,000 years between the two

506    *C. bursa-pastoris* genomes (Figure 3). It is not unexpected that the divergence time is

507    much older between *C. grandiflora* and *C. bursa-pastoris* B compared with the *C.*

508    *grandiflora* and *C. bursa-pastoris* A divergence time, since we have biased our analysis

509    toward finding this result. What is striking is that the divergence time estimate between

510    the two *C. bursa-pastoris* genomes is intermediate between the other two estimates, and

511    the 90% highest posterior density (HPD) overlaps the HPD intervals between *C.*

512    *grandiflora* and both *C. bursa-pastoris* homeologues. Under an allopolyploidy model the

513    divergence time between the two *C. bursa-pastoris* genomes should be the same as the

514    divergence between *C. grandiflora* and *C. bursa-pastoris* B, and significantly different

515    from divergence between *C. grandiflora* and *C. bursa-pastoris* A. This suggests that the

516    true divergence between C. *grandiflora* and both *C. bursa-pastoris* copies reflects an

517    autopolyploid event about 563,000 years ago.

518

519    To further test whether the data fit an autopolyploid model we used test statistics based

520    on shared and fixed sites. We calculated these summary statistics for both the observed

521    data and the data simulated under both models. Most of the differences between the two

522    models are confined to the fixed and shared sites (Table S5). We calculated two further

523    statistics, the differences in both the number of fixed and the number of shared

524    polymorphic sites between *C. grandiflora* and *C. bursa-pastoris* B, on the one hand and

525    *C. grandiflora* and *C. bursa-pastoris* A, on the other hand. We used our two statistics to

526    test for significant departures from the autopolyploid and allopolyploid models. Neither

527    statistics in our observed data depart significantly from the simulated values under the

528    autopolyploid model ($P = 0.4339$ for fixed differences and $P = 0.3673$ for shared

529    differences) while both depart significantly under the allopolyploid model ($P = 0.0032$

530    for fixed differences and $P = 0.0008$ for shared differences) (Figure 4). We therefore

531    cannot reject the autopolyploid hypothesis, while we can reject the allopolyploid model.

532

533    ***Demographic model fitting: Approximate Bayesian computation***

534    Model 2 (allopolyploidy) of our two-split analysis failed to converge in the initial run,

535    making it impossible to continue on to the ABC run. Model 1 (autopolyploidy), however,

536    did produce usable samples indicating that this model fits better our data than model 2.

537    Furthermore, the posterior distributions of most parameters have clear modes, showing

538    that the data is informative for this model (Figure 5). The point estimates of the current

539    population size of the A and B genomes of *C. bursa-pastoris* are 15,000 and 22,000

540    respectively (90% CR: 12,000-23,000 for *C. bursa-pastoris* A and 1,500-43,400 for *C.*

541    *bursa-pastoris* B), while the estimate for *C. grandiflora* is 91,000 (90% CR: 32,600-

542    162,000). The date of the first divergence event, between *C. bursa-pastoris* A and *C.*

543    *bursa-pastoris* B, is 649,000 years (90% CR: 314,000-1,187,000 years), when assuming a

544    generation time of 1 year and a mutation rate of $1.5 \times 10^{-8}$. The date of the second

545    divergence event (739,000 years, 90%CR 361,000-1,443,000) is close to the time of the

546    first divergence event suggesting that the A and B genomes of *C. bursa-pastoris* diverged

547    from each other at a relatively similar time to when they diverged from *C. grandiflora*,

548    thereby strongly supporting an autopolyploid origin of *C. bursa-pastoris*. It was not

549    possible to estimate the population sizes after each coalescent event as the posteriors of

550    these parameters were not informative. Goodness-of-fit tests indicate that the resulting

551    model fit our data reasonably well. We calculated Tajima's D, θw and θπ for each

552 genome from our goodness-of-fit simulations and $S_{nn}$, $G_{ST}$ and $K_{ST}$ among the genomes

553 using Seqlib's build-in goodness-of-fit test and found that all summary statistics fit our

554 data (two-tailed P-value > 0.05) except for Tajima's D (supplementary file S3). The

555 reduced fit to Tajima's D may be reflective of population expansion following

556 divergence.

557 To explore the possible impact of introgression events between *C. bursa-pastoris* and *C.*

558 *rubella* on our inferences, the same analysis, using model 1, was performed using only

559 Chinese *C.bursa-pastoris* samples, which were previously inferred to not be subject to

560 introgression (Slotte et al. 2008). Introgressed alleles would be expected to decrease the

561 divergence time between *C. grandiflora* and *C. bursa-pastoris*. Although the point

562 estimates of the two divergence times were older for this analysis than for the total

563 dataset, the 90% CR was extremely wide and overlapping with time estimates from the

564 full dataset. However, this analysis was not very informative because the divergence time

565 parameters and several other parameter estimates from this analysis had very wide 90%

566 CRs, or/and had no clear mode. Importantly the posterior of the date of the first

567 coalescent event encompasses the prior for this parameter (Figure S3). This may be due

568 to lack of data in the Chinese samples, which have much less diversity than the European

569 ones. This is probably due to the recent origin of the Chinese *C. bursa-pastoris*

570 populations (Slotte et al. 2008). In fact, this reduction in diversity is supported by our

571 Chinese-only ABC analysis, as θ is one of the few well-inferred parameters of the model

572 (effective population size of Chinese *C. bursa-pastoris* 4550, 90% CR: 3,383-12,033)

573

### *Gene conversion and interlocus recombination*

575 The results indicating a lack of gene flow between homeologues suggest that there has

576 not been extensive gene conversion and/or historical recombination events, but we also

577 conducted explicit tests for this. None of our loci showed evidence for gene conversion

578 between *C. bursa-pastoris* homeologues using the geneconv software. However, two

579 highly polymorphic loci, DFR ($R_m$=4) and At4g14190 ($R_m$=2), showed non-zero

580 minimum number of recombination events between the two homeologues, suggesting the

581 possibility of some level of interlocus gene conversion. Given that these loci, particularly

582 DFR, are highly polymorphic in the diploid species, it is possible that the recombination

583    events may have originated in the ancestral population rather than be due to homeologous

584    gene conversion. Indeed, one of the recombination events in At4G14190 is also present in

585    *C. grandiflora* (data not shown).

586

587

588 DISCUSSION

589

590 So far, it has proven difficult to establish whether *C. bursa-pastoris* is an allopolyploid or

591 an autopolyploid. Various studies have resulted in often-conflicting theories as to the

592 evolutionary origins of *C. bursa-pastoris,* some lending support to an allopolyploid origin

593 (Hurka et al. 1989; Slotte et al. 2006) and others to an autopolyploid one (Hurka and

594 Neuffer 1997). Because divergence is recent and extensive shared polymorphisms persist,

595 a coalescent-based approach incorporating population samples and multilocus nuclear

596 data becomes crucial to accurately distinguish models of polyploid speciation. In the

597 present study we used sequence polymorphism and divergence at 14 nuclear loci and two

598 different coalescent-based approaches to test whether *C. bursa-pastoris* had an

599 autopolyploid or allopolyploid origin.

600

601 We conducted three types of analysis to investigate the two possible origins of *C. bursa-*

602 *pastoris*. First we examined the diversity among the three *Capsella* species and inferred

603 their phylogeny using the program BEST. Second we estimated the parameters of an

604 isolation-with-migration model with the program MIMAR and used these estimates to

605 conduct coalescent simulations under both models. Finally, we used Approximate

606 Bayesian Computation to estimate parameters of two-split models representing our null

607 and alternative hypotheses. We could not reject an autopolyploid origin of *C. grandiflora*

608 in any of these analyses, whereas our results were inconsistent with an allopolyploid

609 model. Based on our analyses, the lower bound of the time of origin of *C. bursa-pastoris*

610 is between 270,000 and 700,000 years ago. *C. bursa-pastoris* would thus still be much

611 older than *C. rubella* which most likely diverged from *C. grandiflora* less than 50,000

612 years ago (Foxe et al. 2009; StOnge et al. 2011) allowing us to rule out the suggestion

613 that *C. bursa-pastoris* could be an allopolyploid of *C. rubella* and *C. grandiflora* (Hurka

614 et al. 1989) in agreement with the conclusion of Slotte et al. (2006). Even if these time

615 estimates should be taken with a grain of salt given the uncertainty around mutation rates

616 (Beilstein et al. 2010; Ossowski et al. 2010) a rather recent autopolyploid origin would be

617 consistent with the low level of diversity in *C. bursa-pastoris*, and it would also mean

618 that disomic inheritance has evolved quite rapidly in this species. The ABC analysis

619   indicates that the divergence time of the two homeologueous chromosomes of *C. bursa-*
620   *pastoris* is very close to the divergence between *C. grandiflora* and *C. bursa-pastoris*,
621   suggesting that if there was a period of tetrasomic inheritance it was short relative to the
622   age of the tetraploid species. It has been shown in other species that polyploids with
623   tetrasomic segregation (pairing of four homologous chromosomes during meiosis) tend to
624   rediploidize over time as mutations accumulate and chromosomes diverge (Ramsey and
625   Schemske 1998; Soltis et al. 2010). This process can indeed occur rather quickly and
626   diploidization can proceed through structural rearrangements within 30 generations in *A.*
627   *thaliana* (Parisod et al. 2010). Furthermore, autopolyploids with small chromosomes or
628   low chiasma frequencies may exhibit disomic inheritance immediately after their
629   formation (Stebbins 1971). It is also possible that autopolyploid formation from a highly
630   diverse ancestral population such as *C. grandiflora*, may enhance the speed at which
631   disomic inheritance can occur.

632

633   Many polyploid species have multiple origins (Soltis et al. 2003). In a previous study
634   Slotte et al. (2006) argued that the low nucleotide diversity observed for cpDNA
635   sequences and at seven chloroplast microsatellite loci supports a single origin of *C.*
636   *bursa-pastoris*. The chloroplast sequences resulted in a strongly supported phylogeny in
637   which *C. bursa-pastoris* is sister to both diploid species. This topology is consistent with
638   an ancient origin of *C. bursa-pastoris* from *C. grandiflora* given the fact that *C. rubella*
639   derived from *C. grandiflora* much more recently. The level of variation in *C. bursa-*
640   *pastoris* across the 14 loci is similarly low, and is a consequence of a 5-7 fold decrease of
641   the effective population size compared to *C. grandiflora*. This reduction is not as severe
642   as the reduction in population size observed in *C. rubella* (100-1,500 fold reduction, Foxe
643   et al. 2009; 18 fold reduction, St.Onge et al. 2011). This may be the result of a
644   combination of factors. Recurrent polyploid formation would increase genetic variation
645   but would not leave such a strong bottleneck signature; while this might seem to
646   contradict the lack of variation observed in cpDNA, this could reflect subsequent
647   coalescent events in the chloroplast following species formation (Ceplitis et al. 2005;
648   Slotte et al. 2006). Alternatively, the severity of the bottleneck could have been lessened
649   by early gene flow from *C. grandiflora* via pollen, which would not affect diversity in

650   cpDNA. A third alternative is that the evidence for a severe population bottleneck might

651   simply have eroded with time as the divergence of *C. bursa-pastoris* from *C. grandiflora*

652   is much older than the divergence of *C. rubella* from *C. grandiflora*; a more detailed

653   model of small founding population size followed by a recovery in population size is

654   likely equally consistent with the data, and might explain our observed negative values of

655   Tajima's D.

656

657   Gene conversion can have a strong impact on the histories of duplicated genes and

658   genomes (e.g. Osada and Innan 2008) and, in principle, extensive gene conversion in *C.*

659   *bursa-pastoris* could also have affected our results. Extensive gene conversion could

660   theoretically cause an allopolyploid genome to appear as an autopolyploid under our

661   analysis. However, for this to have happened in *C. bursa-pastoris* the amount of gene

662   conversion would have had to be very extensive, which seems highly unlikely. We

663   identified only two of our loci with evidence of interlocus recombination using the

664   minimum number of recombination events, and no evidence for gene conversion using

665   geneconv. Furthermore, the loci showing gene conversion are highly polymorphic in the

666   diploid species, raising the possibility that the identified recombination events could be

667   due to their retention from ancestral polymorphism and/or due to introgression events.

668   Even though gene conversion is unlikely to have been potent enough to alter our

669   conclusion it might still have contributed to the pattern of divergence among the different

670   genomes. Assuming autopolyploidization and speciation occurred simultaneously we

671   would expect the A and B genomes of *C. bursa-pastoris* to split from *C. grandiflora* at

672   the same time. However, we observe a slight gap in the mean values of these dates. This

673   could be caused by early gene conversion between the A and B genomes, making them

674   appear to be slightly more recently diverged from each other than either is to *C.*

675   *grandiflora* although a period of initial tetrasomic inheritance, as discussed previously,

676   might be a more parsimonious explanation. Overall, the similar divergence times between

677   homeologues and *C. grandiflora* make long periods of disomy and/or gene conversion

678   unlikely.

679

680  Another factor that might have influenced our results is introgression. Previous work has

681  identified evidence of introgression from *C. rubella* to *C. bursa-pastoris* (Slotte et al.

682  2008). Evidence for introgression was detected in European populations of *C. bursa-*

683  *pastoris* but was absent in China where *C. rubella* is absent. Since these introgressed

684  alleles would generally be grouped with the A genome, they are expected to increase the

685  divergence between the A and B genomes of *C. bursa-pastoris* and thereby favor our

686  allopolyploid hypothesis. Introgression is therefore not expected to alter our conclusion

687  that *C. bursa-pastoris* has an autopolyploid origin. It would, however, be expected to

688  cause the inferred divergence date between *C. bursa-pastoris* and *C. grandiflora* to be

689  younger. To examine the possible role of introgression from *C. rubella*, and to confirm

690  our general conclusions using *C. rubella* instead of *C. grandiflora*, we conducted *mimar*

691  analysis with asymmetrical gene flow for *C. rubella* and both *C. bursa-pastoris* A and B.

692  Parameter estimates for these runs had particularly wide confidence intervals, likely due,

693  at least in part, to the loss of information on ancestral polymorphism caused by the severe

694  bottleneck in *C. rubella*. Nevertheless, the results are consistent with our previous

695  conclusions: divergence estimates between *C. bursa-pastoris* A and B fall in between the

696  divergence times estimated between *C. bursa-pastoris* A and *C. rubella* (mode: 66,066,

697  95% HPD: 22022-3.9 million years) and *C. bursa-pastoris* B and *C. rubella* (mode: 3.1

698  million years, 95% HPD: 2.2 million-4.0 million years). Furthermore, simulations of

699  autopolyploid models of the observed data conform well to our observed comparisons of

700  *C. rubella* to *C. bursa-pastoris*, while we get higher rejection rates for the allopolyploid

701  model (Supplementary file S2). To further test if the inferred divergence times were

702  being affected by putatively introgressed alleles we conducted an ABC analysis using

703  only the Chinese samples. Although this analysis was not very informative, the 90% CR

704  of the first inferred divergence time using only China's *C. bursa-pastoris* samples was

705  overlapping with the estimate from the total dataset, suggesting that introgression from *C.*

706  *rubella* into *C. bursa-pastoris* did not have a strong impact on our conclusion. Finally, the

707  patterns of haplotype sharing do not indicate that extensive introgression from *C. rubella*

708  is likely to greatly influence our analysis; haplotype sharing was generally comparable

709  for both diploid species. With genome-wide data from large samples of all three species,

710 it will be interesting to re-examine the extent to which haplotype sharing reflects

711 ancestral polymorphism vs. gene flow following speciation.

712

713      It is important to note that all of our modelling approaches focus on a simplified

714 model of speciation and divergence, and it is possible that additional model mis-

715 specifications, particularly in the allopolyploid model, could be leading to a higher

716 rejection rate. For example, Mimar assumes a single population size change following

717 divergence and a constant migration rate, and subsequent population size changes and/or

718 changes in gene conversion rates between homeologues over time could be complicating

719 our inferences. However, our simulations lead us to conclude that the autopolyploid

720 model can explain our data quite well, and it is not obvious why model mis-specification

721 would be a problem specific only to the allopolyploidy model. Nevertheless, it will be

722 important to confirm out conclusions with large-scale genomic data, where the patterns of

723 haplotype structure and divergence across chromosomes can also be incorporated into

724 these analyses.

725

726 **CONCLUSIONS**

727

728 Our study confirms the usefulness of coalescent-based approaches when studying the

729 mode of origin of polyploids, although as pointed by Doyle and Egan (2009) precise time

730 estimates remain elusive and are highly dependent on demographic details and on

731 assumptions on mutation rates. While these results shed much light on the evolutionary

732 origin of *C. bursa-pastoris*, little is still known about the extensive phenotypic changes

733 that have occurred in both *C. bursa-pastoris* and *C. rubella*. Understanding the genomic

734 context and underlying evolutionary forces that have promoted these changes will be of

735 considerable interest in future studies.

736

737

746

747

748

749

750

751 **References**

752

753 Bailey CD, Koch MA, Mayer M, Mummenhoff K, O'Kane Jr SL, Warwick SI,

754 Windham MD, Al-Shehbaz IA. 2006. Toward a global phylogeny of the

755 Brassicaceae. Mol. Biol. Evol. 23:2142–2160.

756 Beaumont MA. 2010. Approximate Bayesian Computation in evolution and ecology.

757 Annu. Rev. Ecol. Evol. Syst. 41: 379–406.

758 Beaumont M, Zhang W, Balding D. 2002. Approximate Bayesian computation in

759 population genetics. Genetics 162:2025-2035.

760 Beilstein MA, Nagalingum NS, Clements MD, Manchester SR, Mathews S. 2010. Dated

761 molecular phylogenies indicate a Miocene origin for Arabidopsis thaliana. Proc.

762 Natl. Acad. Sci., USA 107: 18724-18728.

763 Becquet C, Przeworski M. 2007. A new approach to estimate parameters of speciation

764 models with application to apes. Genome Res 17:1505-1519.

765 Ceplitis A, Su Y, Lascoux M. 2005. Bayesian inference of evolutionary history from

766 chloroplast microsatellites in the cosmopolitan weed *Capsella bursa-pastoris*

767 (Brassicaceae). Molecular Ecology 14:4221-4233

768 Cifuentes M., Grandont L, Moore G, Chèvre AM, Jenczewski E. 2010. Genetic

769 regulation of meiosis in polyploid species: new insights into an old question. New

770 Phytol. 186: 29-36.

771 Cui LY, Wall PK, Leebens-Mack JH, Lindsay BG, Soltis DE, Doyle JJ, Soltis PS,

772 Carlson JE, Arumuganathan K, Barakat A, Albert VA, Ma H, dePamphilis CW.

773 2006. Widespread genome duplications throughout the history of flowering plants.

774 Genome Res 16:738 –749.

775 De Mita S, Ronfort J, McKhann HI, Poncet C, El Malki $_R$, Bataillon T. 2007.

776 Investigation of the demographic and selective forces shaping the nucleotide

777 diversity of genes involved in nod factor signaling in *Medicago truncatula*.

778 Genetics 177: 2123–2133.

779 Doyle JJ, Egan AN. 2009. Dating the origins of polyploidy events. New Phytol. 186:73-

780 85

781    Fawcett JA, Maere S, Van de Peer Y. 2009. Plants with double genomes might have had
782         a better chance to survive the Cretaceous–Tertiary extinction event. Proc Natl
783         Acad Sci, USA. 106: 5737–5742.

784    Foxe JP, Slotte T, Stahl EA, Neuffer B, Hurka H, Wright SI. 2009. Rapid morphological
785         evolution and speciation associated with the evolution of selfing in *Capsella*. Proc
786         Natl Acad Sci, USA.  106:5241-5245.

787    Gaut  BS, Doebley JF. 1997. DNA sequence evidence for the segmental allotetraploid
788         origin of maize. Proc Natl Acad Sci, USA. 94: 6809-6814.

789    Guo YL, Bechsgaard J, Slotte T, Neuffer B, Lascoux M, Weigel D, Schierup MH. 2009.
790         Recent speciation of *Capsella rubella* from *Capsella grandiflora*, associated with
791         loss of self-incompatibility and an extreme bottleneck. Proc. Natl. Acad. Sci.,
792         USA 106:5246-5251.

793    Hey J. 2010. Isolation with migration models for more than two populations. Mol Biol
794         Evol. 27: 905-920.

795    Hey J, Nielsen R. 2007. Integration within the Felsenstein equation for improved Markov
796         chain Monte Carlo methods in population genetics. Proc Natl Acad Sci, USA.
797         104: 2785-2790.

798    Hudson RR 2002. Generating samples under a Wright-Fisher neutral model of genetic
799         variation. Bioinformatics 18: 337-338.

800    Hudson, RR, Kaplan N. 1985. Statistical properties of the number of recombination
801         events in the history of a sample of DNA-sequences. Genetics, 111: 147–164.

802    Hurka H, Freundner S, Brown AH, Plantholt U. 1989. Aspartate aminotransferase
803         isozymes in the genus *Capsella* (Brassicaceae): subcellular location, gene
804         duplication, and polymorphism. Biochemical genetics 27:77-90.

805    Hurka H, Neuffer B. 1997. Evolutionary processes in the genus *Capsella* (Brassicaceae).
806         Plant Syst Evol. 206:295-316.

807    Jakobsson M, Hagenblad J, Tavaré S, Säll T, Halldén C, Lind-Halldén S, Nordborg M.
808         2006. A unique recent origin of the allotetraploid species *Arabidopsis suecica*:
809         Evidence from nuclear DNA markers. Mol Biol Evol. 23:1217-1231.

810    Koch M, Haubold B, Mitchell-Olds T. 2000. Comparative evolutionary analysis of

811         chalcone synthase and alcohol dehydrogenase loci in Arabidopsis, Arabis, and

812         related genera (Brassicaceae). Mol Biol Evol 17:1483-1498.

813    Librado, P, Rozas, J 2009. DnaSP v5: a software for comprehensive analysis of DNA

814         polymorphism data. Bioinformatics 25: 1451-1452.

815    Liu L. 2008. BEST: Bayesian estimation of species trees under the coalescent model.

816         Bioinformatics 24: 2542-2543.

817    Masterson J. 1994. Stomatal size in fossil plants: evidence for polyploidy in the majority

818         of angiosperms. Science 264: 421–423.

819    Nicholas KB, Nicholas HB Jr., Deerfield DW II. 1997 GeneDoc: Analysis and

820         Visualization of Genetic Variation, EMBNEW.NEWS 4:14

821    Nielsen R, Wakeley J. 2001. Distinguishing migration from isolation: a Markov chain

822         Monte Carlo approach. Genetics 158:885-896.

823    Noor MA, Feder JL. 2006. Speciation genetics: evolving approaches. Nature Reviews

824         Genetics 7:851-861.

825    Osada, N., Innan I. 2008. Duplication and gene conversion in the *Drosophila*

826         *melanogaster* genome. PLoS Genet 4(12): e1000305

827    Ossowski S, Schneeberger K, Lucas-Lledó JI*,* Warthmann N, Clark RM, Shaw RG,

828         Weigel D, Lynch M. 2010.  The rate and molecular spectrum of spontaneous

829         mutations in Arabidopsis thaliana. Science 327: 92-94.

830    Otto SP, Whitton J. 2000. Polyploid incidence and evolution. Ann. Rev. Genet. 34:401-

831         437.

832    Parisod C., Holderegger R, Brochmann C. 2010. Evolutionary consequences of

833         autopolyploidy. New Phytol. 186:5-17.

834    Ramos-Onsins SE, Stranger BE, Mitchell-Olds T, Aguadé M. 2004. Multilocus analysis

835         of variation and speciation in the closely related species *Arabidopsis halleri* and

836         *A. lyrata*. Genetics 166:373-388.

837    Ramsey J, Schemske D. 1998. Pathways, mechanisms and rates of polyploid formation in

838         flowering plants. Ann. Rev Ecol Syst. 29:467-501.

839    Ross-Ibarra J, Wright SI, Foxe JP, Kawabe A, DeRose-Wilson L, Gos G, Charlesworth

840        D, Gaut BS. 2008. Patterns of polymorphism and demographic history in natural
841        populations of *Arabidopsis lyrata*. PLoS One. 6:e2411.

842    Rozen S, Skaletsky H. 2000. Primer3 on the WWW for general users and for biologist
843        programmers. Methods Mol Biol. 132:365–386.

844    Sawyer S. 1989. Statistical test for detecting gene conversion. Mol. Biol. Evol. 6: 526-
845        538.

846    Schierup MH, Hein J. 2000. Consequences of recombination on traditional phylogenetic
847        analysis. Genetics 156: 879-891.

848    Shull GH. 1929. Species hybridization among old and new species of Shepherd's Purse.
849        Int Congr Plant Sci. 1:837–888.

850    Slotte T, Ceplitis A, Neuffer B, Hurka H, Lascoux M. 2006. Intrageneric phylogeny of
851        *Capsella* (*Brassicaceae*) and the origin of the tetraploid *C. bursa-pastoris* based
852        on chloroplast and nuclear DNA sequences. Am J Bot 93:1714-1724.

853    Slotte T, Huang H, Lascoux M, Ceplitis A. 2008. Polyploid speciation did not confer
854        instant reproductive isolation in *Capsella* (Brassicaceae). Mol Biol Evol. 25:1472-
855        1481.

856    Soltis DE, Soltis PS, Tate JA. 2003. Advances in the study of polyploidy since Plant
857        Speciation. New Phytol 161:173-191.

858    Soltis DE, Buggs RJA, Doyle JJ, Soltis PS. 2010.  What we still don't know about
859        polyploidy. Taxon. 59:1387-1403.

860    Stebbins GL. 1971. Chromosomal Evolution in Higher Plants. Edward Arnold: London.

861    Stephens M, Smith NJ, Donnelly P. 2001. A new statistical method for haplotype
862        reconstruction from population data. Am J Hum Genet 68: 978-989.

863    StOnge K, Källman T, Slotte T, <u>Lascoux M</u>, Palmé AE. 2011. Divergent population
864        history and structure in two closely related species (*Capsella rubella* and *C.*
865        *grandiflora*) with different mating systems. Molecular Ecology 20: 3306–3320.

866    Tajima F. 1983. Evolutionary relationships of DNA-sequences in finite populations.
867        Genetics  105: 437-460.

868    Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by DNA
869        polymorphism. Genetics 123: 585-595.

870    U, N. 1935. Genome analysis in Brassica with special reference to the experimental

871            formation of *B. napus* and peculiar mode of fertilization. Jpn J Bot 7: 389–452.

872     Velasco R, and 86 co-authors 2010. The genome of the domesticated apple (*Malus* •

873            *domestica* Borkh.). Nature Genetics 42:833-841.

874     Wakeley J., Hey J. 1997. Estimating ancestral population parameters. Genetics 145:847-

875            855.

876     Wang RL, Wakeley J, Hey J. 1997. Gene flow and natural selection in the origin of

877            *Drosophila pseudoobscura* and close relatives. Genetics 147:1091-1106.

878

879     Wang W-K, Ho C-W, Hung K-H, Wang, K-H, Huang, C-C, Araki, H,  Hwang, C-C,

880            Hsu, T-W, Osada N, Chiang, T-Y (2010) Multilocus analysis of genetic

881            divergence between outcrossing Arabidopsis species: evidence of genome-wide

882            admixture. New Phytol 188: 488-500

883

884     Watterson GA. 1975. On the number of segregating sites in genetical models without

885            recombination. Theor Pop Biol 7: 256-276.

886     Wood TE, Takebayashi N, Barker MS, Mayrose I, Greenspoon PB, Rieseberg LH. 2009.

887            The frequency of polyploid speciation in vascular plants. Proc Natl Acad Sci USA

888            106:13875-13879

889     Zhang L-B, Ge S. 2007. Multilocus analysis of nucleotide variation and speciation in

890            *Oryza officinalis* and its close relatives. Mol Biol Evol 24:769-783.

891

892

893    **Figure text**

894

895    Figure 1. Model diagrams of the null hypothesis, autopolyploidy, and alternative

896    hypothesis, allopolyploidy.

897

898    Figure 2. Number of synonymous fixed differences between all pairs of *C. bursa-pastoris*

899    A, *C. bursa-pastoris* B, *C. grandiflora* and *C. rubella.*

900

901    Figure 3. Marginal posterior distributions of speciation parameters estimated by MIMAR,

902    with posterior modes showing good fit to data summaries. $\theta = 4Ne\mu$ where *Ne* is the

903    effective population size and $\mu$ is the mutation rate ($1.5 \times 10^{-8}$/site/year)

904    A) Constrained model: the model assumes equal effective population sizes in the ancestor

905    as in present-day *C.grandiflora*: Model 1; Species 1 = *C. grandiflora*, Species 2 = *C.*

906    *bursa-pastoris A*. The model is represented by continuous lines. Model 2; Species 1 = *C.*

907    *grandiflora*, Species 2 = *C. bursa-pastoris B.* The model is shown by a dotted line. Tgen

908    Divergence time (years) between *C. grandiflora* and *C. bursa-pastoris* A and between *C.*

909    *grandiflora* and *C. bursa-pastoris* B

910    B) Unconstrained model: θA ancestral *C. grandiflora*, θ1 *C. bursa-pastoris* A

911    (continuous line), θ2 *C. bursa-pastoris* B (dotted line). Tgen Divergence time (years)

912    between *C. bursa-pastoris* A and *C. bursa-pastoris* B.

913

914    Figure 4: Density distribution of the simulated values of the summary statistics under (A)

915    autopolyploidy and (B) allopolyploidy. The left column gives the distribution of the mean

916    of *fix_diff* over the fourteen genes, where *fix_diff* is the difference between the number of

917    fixed sites of each of the homoelogues when it is compared to *C. grandiflora*. The right

918    column gives the same for *shared_diff,* the difference between the number of shared

919    polymorphic sites of each of the homoelogues to when it is compared to *C. grandiflora*.

920    The blue vertical line is the observed value. P values are given in the upper right corner

921    of each plot. See text for details.

922

923    Figure 5. Posterior distributions of informative parameters in the two-split model where

924     the two *C. bursa-pastoris* genomes coalesce first, followed by coalescence of their

925     common ancestor with *C. grandiflora*. $\theta=4N_e\mu$ where $N_e$ is the effective population size

926     of *C. bursa-pastoris* A and $\mu$ is the mutation rate ($1.5 \times 10^{-8}$/site/year). Other effective

927     population sizes and the divergence times are relative to this first estimate. Divergence

928     times are on a scale of $4N_e\times$generations and the second date parameter is additive to the

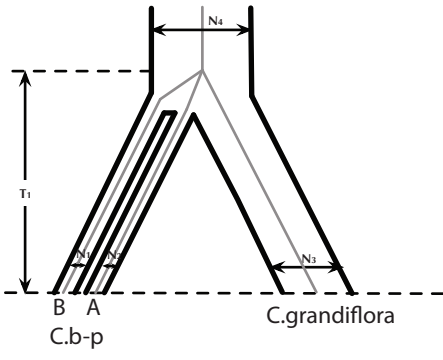929     first.

930

931

932

933

**autopolyploid**

**allopolyploid**

$N_4$

$T_1$

$N_1$ $N_2$

$N_3$

B A

C.b-p

C.grandiflora

$T_2''$

$N_4$

$T_2'$

$N_1$ $N_2$

$N_3$

B A

C.b-p

C.grandiflora

**A**

θ1 θ2 TGen

**B**

θA θ1/θ2 TGen

**A**

**fixed**

*P* =0.4339

differences in number of fixed sites

**Shared**

*P* =0.3673

differences in number of Shared sites

**B**

**fixed**

*P* =0.0032

differences in number of fixed sites

**Shared**

*P* =8e−04

differences in number of Shared sites