

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34

Discoveries Article

Variation in linked selection and recombination drive genomic divergence during allopatric speciation of European and American aspens

Jing Wang¹, Nathaniel R. Street², Douglas G. Scofield^{1,3,4}, Pär K. Ingvarsson¹

¹ Department of Ecology and Environmental Science, Umeå University, SE-90187, Umeå, Sweden

² Umeå Plant Science Centre, Department of Plant Physiology, Umeå University, SE-90187, Umeå, Sweden

³ Department of Ecology and Genetics: Evolutionary Biology, Uppsala University, Uppsala, Sweden

⁴ Uppsala Multidisciplinary Center for Advanced Computational Science, Uppsala University, Uppsala, Sweden

Corresponding author:

Dr Pär K. Ingvarsson, Department of Ecology and Environmental Science, Umeå University, Umeå, SE 90187, Sweden. Phone: +46907867414; Fax: +46-(0)-90-786-6705; E-mail: par.ingvarsson@emg.umu.se

Keywords: *Populus tremula*, *P. tremuloides*, Whole-genome re-sequencing, demographic histories, heterogeneous genomic differentiation, linked selection, recombination

35 **Abstract**

36

37 Despite the global economic and ecological importance of forest trees, the genomic
38 basis of differential adaptation and speciation in tree species is still poorly understood.
39 *Populus tremula* and *P. tremuloides* are two of the most widespread tree species in the
40 Northern Hemisphere. Using whole-genome re-sequencing data of 24 *P. tremula* and
41 22 *P. tremuloides* individuals, we find that the two species diverged ~2.2-3.1 million
42 years ago, coinciding with the severing of the Bering land bridge and the onset of
43 dramatic climatic oscillations during the Pleistocene. Both species have experienced
44 substantial population expansions following long-term declines after species
45 divergence. We detect widespread and heterogeneous genomic differentiation
46 between species, and in accordance with the expectation of allopatric speciation,
47 coalescent simulations suggest that neutral evolutionary processes can account for
48 most of the observed patterns of genomic differentiation. However, there is an excess
49 of regions exhibiting extreme differentiation relative to those expected under
50 demographic simulations, which is indicative of the action of natural selection.
51 Overall genetic differentiation is negatively associated with recombination rate in
52 both species, providing strong support for a role of linked selection in generating the
53 heterogeneous genomic landscape of differentiation between species. Finally, we
54 identify a number of candidate regions and genes that may have been subject to
55 positive and/or balancing selection during the speciation process.

56

57

58

59

60

61

62

63

64

65

66

67

68 **Introduction**

69

70 Understanding how genomes diverge during the process of speciation has received a
71 great deal of attention in the evolutionary genetics literature in recent years (Nosil et
72 al. 2009; Strasburg et al. 2012; Seehausen et al. 2014). Under strict neutrality,
73 differentiation is expected to accumulate as a result of the stochastic fixation of
74 polymorphisms by genetic drift (Coyne and Orr 2004). Demographic processes, such
75 as population bottlenecks or expansions, can accelerate or decelerate the rate of
76 differentiation through changes in the effective population sizes of nascent daughter
77 species (Avice 2000). Random genetic drift and demographic processes are both
78 expected to affect the entire genome (Luikart et al. 2003). Natural selection, however,
79 only influence loci involved in ecological specialization and/or reproductive isolation,
80 resulting in patterns of polymorphisms and divergence that deviate from neutral
81 predictions (Luikart et al. 2003; Via 2009). The functional architectures of genomes,
82 e.g. mutation and recombination rates, are also important factors in determining
83 genomic landscape of differentiation (Noor and Bennett 2009; Nachman and Payseur
84 2012; Renaut et al. 2013). For example, suppressed recombination could increase
85 genetic differentiation either by limiting inter-species gene flow to prevent the break-
86 up of co-adapted alleles, or through the diversity-reducing effects of linked selection
87 (Noor and Bennett 2009). However, disentangling the relative importance of these
88 evolutionary forces when interpreting patterns of genomic divergence remains a
89 challenge in speciation genetics.

90 With the advance of next generation sequencing (NGS) technologies, a
91 growing number of studies have found highly heterogeneous patterns of genomic
92 differentiation between recently diverged species (Turner et al. 2005; Ellegren et al.
93 2012; Renaut et al. 2013; Carneiro et al. 2014; Feulner et al. 2015). A common
94 explanation for these patterns is that levels of gene flow between species differ across
95 the genome. Increased genetic divergence is usually observed in a small number of
96 regions containing loci involved in reproductive isolation ('speciation islands'), where
97 as the remainder of the genome is still permeable to ongoing gene flow and therefore
98 shows lower levels of differentiation (Nosil et al. 2009; Sousa and Hey 2013).
99 However, some recent studies have argued that highly differentiated regions represent
100 'incidental islands' that are not tied to the speciation processes per se. Rather they are

101 seen simply as a result of the diversity-reducing effects of linked selection that
102 accelerate lineage sorting of ancestral variation and increase interspecific
103 differentiation, especially in regions of reduced recombination (Turner and Hahn
104 2010; Cruickshank and Hahn 2014). In addition, long-term balancing selection is
105 supposed to maintain stable trans-species polymorphisms and leaves signatures of
106 unusually low genetic differentiation between species (Charlesworth 2006). Under
107 these scenarios, natural selection alone is sufficient to generate patterns of
108 heterogeneous genomic differentiation even under complete allopatry (Noor and
109 Bennett 2009; Turner and Hahn 2010). Finally, strictly neutral forces, such as
110 stochastic genetic drift and complex demographic processes, can also create
111 heterogeneous genomic divergence and generate patterns of divergence and
112 polymorphism that mimic the effects of selection (Nosil et al. 2009; Campagna et al.
113 2015). In general, the three hypotheses listed above are not mutually exclusive and
114 exhaustive examination of these hypotheses requires detailed information on the
115 speciation process, such as the timing of speciation, the geographic and demographic
116 context in which it occurred (Nosil and Feder 2012).

117 Although largely understudied compared to other model species, forest trees
118 represent a promising system to understand the genomic basis of species divergence
119 and adaptive evolution; as a group they have developed diverse strategies to adapt and
120 thrive across a wide range of climates and environments (Neale and Kremer 2011).
121 *Populus tremula* (European aspen) and *P. tremuloides* (American aspen) are two of
122 the most ecologically important and geographically widespread tree species of the
123 Northern Hemisphere (Figure 1a). Both are keystone species, display rapid growth,
124 with high tolerance to environmental stresses and long-distance pollen and seed
125 dispersal via wind (Eckenwalder 1996; Müller et al. 2012). In addition, they both
126 harbor among the highest level of intraspecific genetic diversity reported in plant
127 species so far (Wang et al. forthcoming). Based on their morphological similarity and
128 close phylogenetic relationships, they are considered to be sister species, or less
129 commonly, conspecific subspecies (Eckenwalder 1996; Wang et al. 2013). They can
130 readily cross and artificial hybrids usually show high heterosis (Hamzeh and
131 Dayanandan 2004; Tullus et al. 2012). A recent study based on a handful of nuclear
132 and chloroplast loci suggests that the first opening of the Bering land bridge may have
133 driven the allopatric speciation of the two species (Du et al. 2015).

134 Due to their continent-wide distributions, extraordinary levels of genetic and
135 phenotypic diversity, along with the availability of a high-quality reference genome in
136 the congener, *P. trichocarpa* (Tuskan et al. 2006), *P. tremula* and *P. tremuloides*
137 represent a promising system for evaluating how various evolutionary processes have
138 shaped the patterns of genomic divergence during speciation. In this study, we use
139 whole-genome re-sequencing data from both species to estimate and infer their
140 divergence time and historical demographic processes of the two species. Explicit
141 characterizations of the demographic history not only allow us to estimate historical
142 population size fluctuations in both species, but also increase the accuracy of
143 identifying regions or genes that have been under natural selection. By incorporating
144 the inferred demographic scenarios into the null model, we investigate the extent to
145 which demographic and selective events have contributed to the overall patterns of
146 genomic differentiation between the two species. We also identify a number of outlier
147 regions and genes that likely have evolved in response to positive and/or balancing
148 selection during the speciation process.

149

150 **Results**

151

152 We generated whole-genome resequencing data for 24 *P. tremula* and 22 *P.*
153 *tremuloides*. The high extent of conserved synteny between the genomes of aspen and
154 *P. trichocarpa* (Pakull et al. 2009; Robinson et al. 2014) allowed us to map all reads
155 to the *P. trichocarpa* reference genome (v3.0) (Tuskan et al. 2006) after adapter
156 removal and quality trimming (see Materials and Methods). More than 88% of
157 sequenced reads were aligned and the mean coverage of uniquely mapped reads per
158 site was 25.1 and 22.5 in samples of *P. tremula* and *P. tremuloides*, respectively
159 (Table S1). Two complementary bioinformatics approaches were used in this study
160 (Figure S1): (1) For those population genetic statistics that relied on inferred site-
161 frequency-spectrum (SFS), estimation was performed directly from genotype
162 likelihoods without calling genotypes (Nielsen et al. 2011) as implemented in
163 ANGSD (Korneliussen et al. 2014). (2) For those estimations that required accurate
164 genotype calls, single nucleotide polymorphisms (SNPs) and genotypes were called
165 with HaplotypeCaller in GATK (Danecek et al. 2011). In total, we identified

166 5,894,205 and 6,281,924 SNPs passing filtering criteria (see Materials and Methods)
167 across the 24 *P. tremula* samples and 22 *P. tremuloides* samples, respectively.

168

169 **Population structure**

170 We used NGSadmix (Skotte et al. 2013) to infer individual ancestry based on
171 genotype likelihoods, which takes the uncertainty of genotype calling into account. It
172 clearly sub-divided all sampled individuals into two species-specific groups when the
173 number of clusters (K) was 2 (Figure 1b). When $K = 3$, there was evidence for further
174 population sub-structuring in *P. tremuloides*, where individuals from populations of
175 Alberta and Wisconsin clustered into two subgroups. With $K = 4$, most individuals of
176 *P. tremula* were inferred to be a mixture of two genetic components, showing slight
177 clinal variation with latitude. No further structure was found when $K = 5$ (Figure 1b).
178 A principal component analysis (PCA) further supported these results (Figure 1c).
179 Only the first two components were significant based on the Tracy-Widom test (Table
180 S2), which explained 21.4% and 2.1% of total genetic variance, respectively (Figure
181 1c). Among the total number of polymorphisms in the two species, fixed differences
182 between *P. tremula* and *P. tremuloides* accounted for 1.1%, whereas 16.7% of
183 polymorphisms were shared between species, with the remaining polymorphic sites
184 being private in either of the two species (Figure 1d).

185 To further examine the extent of population subdivision in *P. tremuloides*, we
186 measured F_{ST} and d_{xy} between the two subpopulations (Alberta and Wisconsin) along
187 individual chromosomes (Table S3). We found low levels of genetic differentiation
188 (average F_{ST} : 0.0443 ± 0.0325) between the two subpopulations (Table S3). Total
189 sequence differentiation in the inter-population comparison (mean d_{xy} =
190 0.0165 ± 0.0083) was similar to mean sequence differences in intra-population
191 comparisons (π_{Alberta} : 0.0161 ± 0.0081 ; $\pi_{\text{Wisconsin}}$: 0.0157 ± 0.0080 , Table S3), indicating
192 that individuals of the two populations were genetically not more different from each
193 other than individuals within each population. Based on the summaries of site
194 frequency spectrum (Tajima's D and Fay & Wu's H), both populations exhibited
195 strong skews toward low-frequency variants (negative D) and intermediate skews
196 toward high frequency-derived variants (negative H) (Table S3), suggesting that they
197 likely experienced similar species-wide demographic events.

198

199 **Demographic histories**

200 We used *fastsimcoal2* (Excoffier et al. 2013), a coalescent simulation-based method,
201 to infer the past demographic histories of *P. tremula* and *P. tremuloides* from the joint
202 site frequency spectrum. Eighteen divergence models were evaluated (Figure S2;
203 Table S4), and all models began with the split of the ancestral population into two
204 derived populations and differed in terms of (i) whether post-divergence gene flow
205 was present or not, (ii) levels and patterns of gene flow between the two species, and
206 (iii) how population size changes occurred, either at the time of species divergence or
207 afterwards (Figure S2). The best-fitting model was a simple isolation-with-migration
208 model where, after the two species diverged, *P. tremuloides* experienced exponential
209 growth and whereas a stepwise population size change occurred in *P. tremula* (Figure
210 2a). The exact parameter estimates of divergence time, gene flow, population sizes
211 and their associated 95% confidence intervals (CIs) are given in Table 1. The
212 estimated divergence time between *P. tremula* and *P. tremuloides* (T_{DIV}) was ~2.3
213 million years ago (Mya) (bootstrap range [BR]: 2.2-3.1 Mya). The contemporary
214 effective population sizes (N_e) of *P. tremula* ($N_{P.tremula}$) and *P. tremuloides*
215 ($N_{P.tremuloides}$) were 102,814 (BR: 93,688-105,671) and 309,500 (BR: 247,321-310,105)
216 respectively, with both being larger than the effective population size of their
217 common ancestor ($N_{ANC} = 56,235$ [48,012-69,492]). Gene flow ($2N_e m$, where N_e is
218 the effective population size and m is the migration rate) from *P. tremuloides* to *P.*
219 *tremula* was higher (0.202 migrants per generation [0.156-0.375]) than in the opposite
220 direction (0.053 [0.052-0.117]), most likely reflecting the higher N_e in *P. tremuloides*
221 than in *P. tremula* (Slatkin 1985). Overall, the migration rates in both directions were
222 fairly low given the large N_e of both species (Morjan and Rieseberg 2004), which is
223 not unexpected given the large geographical distance and disjunct distributions
224 between the two species.

225 We employed the multiple sequential Markovian coalescent (MSMC) method
226 to investigate changes of N_e over time based on inferring the time to the first
227 coalescence between pairs of haplotypes (Schiffels and Durbin 2014). Higher
228 resolution of recent population size changes is expected when more haplotypes are
229 used (Schiffels and Durbin 2014). We therefore applied MSMC to phased whole-
230 genome sequences from one (two haplotypes), two (four haplotypes) and four (eight
231 haplotypes) individuals in each species, respectively. We did not include more
232 haplotypes because of the high computational cost of larger samples. The MSMC-

233 based estimates of N_e for both *P. tremula* (60,796) and *P. tremuloides* (49,701) at the
234 beginning of species divergence (around 2.3 Mya) were very similar to the
235 *fastsimcoal2*-based estimates of N_e for their ancestral population (Figure 2). The two
236 species experienced similar magnitudes of population decline following their initial
237 divergence (Figure 2b). Population expansion in *P. tremuloides* occurred around
238 50,000-70,000 years ago and continued up to the present (Figure 2b), whereas *P.*
239 *tremula* experienced a population expansion following a substantially longer periods
240 of bottleneck (Figure 2b).

241 To assess the possible confounding effects of population subdivision and
242 biased sampling scheme on demographic inferences in both species, we first applied
243 MSMC analysis using *P. tremuloides* individuals originating from populations in
244 Alberta and Wisconsin separately and compared them to the result obtained from the
245 pooled samples (Figure S3). Although the Wisconsin population was found to have
246 undergone a decline in population size during the last 2000-3000 years ago (Figure
247 S3), both local populations of *P. tremuloides* show evidence for a longer period of
248 species-wide expansion when compared with *P. tremula*, which was in accordance
249 with the results observed from the pooled samples. Second, demographic inferences
250 of both species were also supported by the summary statistics based on the nucleotide
251 diversity (θ_π) and site frequency spectrum (Tajima's D and Fay & Wu's H) (Figure
252 S4). The θ_π in the two subpopulations of *P. tremuloides* were all marginally higher
253 than in *P. tremula*, suggesting that the large effective population size found in *P.*
254 *tremuloides* is not influenced by the presence of intraspecific population subdivision
255 (Figure S4a). In addition, the signal of more negative values of Tajima's D in both
256 local and pooled samples of *P. tremuloides* (Figure S4b) suggest that it may have
257 gone through a more pronounced and/or longer period of population expansion
258 compared to *P. tremula*. The lower values of the genome-wide Fay & Wu's H in *P.*
259 *tremula* (Figure S4c), on the other hand, might reflect the relatively longer period of
260 low population size during the bottleneck. Taken together, these results suggest that
261 population subdivision of *P. tremuloides* and the unbalanced sampling schemes
262 between the two species have negligible effects on our demographic inferences.

263

264 **Genome-wide patterns of differentiation and identification of outlier regions**
265 **against the best-fitting demographic model**

266 To investigate patterns of interspecific genetic differentiation across the genome, we
267 calculated the standard measure of genetic divergence, F_{ST} , between *P. tremula* and *P.*
268 *tremuloides* over non-overlapping 10 Kbp windows (Figure 3). Levels of genetic
269 differentiation varied greatly throughout the genome, with the majority of windows
270 showing high genetic differentiation (mean F_{ST} = 0.386) between species (Figure 3).

271 In order to test the extent to which historical demographic events can explain
272 the observed patterns of genetic divergence between the two species, we used
273 coalescent simulations performed in *msms* (Ewing and Hermisson 2010) to compare
274 the observed patterns of differentiation to that expected under three demographic
275 models (Figure S5). The demographic scenario in model 1 was same as the best-
276 fitting model inferred by *fastsimcoal2* (Figure S5a; Table S5). In another two models,
277 we incorporated the population subdivision of *P. tremuloides* into the best-fitting
278 demographic model. In model 2 (Figure S5b; Table S5), we assume there was no gene
279 flow between the two subpopulations of *P. tremuloides* and explored different values
280 of their divergence time until the simulated F_{ST} values between the two
281 subpopulations matched those observed (Figure S6). The same procedure was applied
282 to model 3 (Figure S5c; Table S5), except that we there assume the per-generation
283 gene flow between the two subpopulations of *P. tremuloides* ($4N_e m$) was equal to 10
284 and increased their divergence time in tandem with gene flow. To assess the fit of
285 these models, we compared two summary statistics, θ_π and Tajima's D, between the
286 simulated and observed data for both species. As can be seen from Figure S7 and
287 Figure S8, there was generally a good agreement between observed and simulated
288 data sets for all three models. In addition, the above three models showed consistent
289 distributions of simulated F_{ST} values between *P. tremula* and *P. tremuloides*,
290 indicating that the presence of population subdivision in *P. tremuloides* has little
291 effect on the overall patterns of genomic divergence that we observe between the two
292 species (Figure S9).

293 Comparing the observed distribution of inter-specific F_{ST} with that obtained
294 from simulations based on the best-fitting demographic model, we found that the
295 observed distribution was flatter and contained greater proportions of regions falling
296 in the extremes of distribution (Figure 4a). We used the distribution of inter-specific
297 F_{ST} based on 500,000 coalescent simulations to identify outlier windows that were
298 likely targets of natural selection by performing a false discovery rate (FDR) of 1% as

299 our cut-off. With this approach, 674 and 262 windows ($FDR < 0.01$) were, respectively,
300 identified as showing exceptionally high and low F_{ST} between the two species (Figure
301 4a). After examining the genomic distribution, physical sizes and overlaps of these
302 outlier windows, we found that both highly and lowly differentiated regions were
303 randomly distributed across the genome (Figure 3) and that the sizes of these regions
304 appeared to be rather small, with the majority occurring on a physical scale smaller
305 than 10 Kbp (Figure S10).

306

307 **Signatures of selection in outlier regions**

308 As F_{ST} is a relative measure of differentiation and is thus sensitive to any processes
309 that alter intra-species genetic variation (Charlesworth 1998; Cruickshank and Hahn
310 2014), we quantified and compared inter-specific genetic differentiation between two
311 unions of outlier windows and the rest of the genome using three additional
312 approaches: (i) pairwise nucleotide divergence between species (d_{xy}), which is a
313 measure that is independent of within-species diversity (Nei 1987); (ii) relative node
314 depth (RND) (Feder et al. 2005), which takes into account possible variation in the
315 mutation rate among genomic regions by dividing d_{xy} of the two aspen species with
316 d_{xy} between aspens and a third more distantly related species (*P. trichocarpa*); and
317 (iii) the proportion of inter-specific shared polymorphisms. Compared to the genomic
318 background average, both d_{xy} and RND revealed much greater divergence between
319 the two species in regions of high differentiation (Figure 4b; Table S6) and, in
320 accordance with these patterns, the proportion of inter-specific shared polymorphisms
321 was significantly lower in these regions (Figure 4b; Table S6). In addition, these
322 regions are characterized by multiple signatures of positive selection within one or
323 both species, including significantly reduced levels of polymorphism (θ_π), skewed
324 allele frequency spectrum towards rare alleles (more negative Tajima's D), increased
325 high-frequency derived alleles (more negative Fay & Wu's H), and stronger signals of
326 linkage disequilibrium (LD) ($P < 0.001$, Mann-Whitney U test) (Figure 4c; Table S6).
327 Relative to genome-wide averages, these regions also contained significantly higher
328 proportions of fixed differences that were caused by derived alleles fixed in either
329 species (Figure 4c; Table S6).

330 In contrast to patterns found in regions of high differentiation, regions of low
331 differentiation showed significantly higher levels of polymorphism, excesses of

332 intermediate-frequency alleles (higher Tajima's D and Fay & Wu's H values), higher
333 proportions of inter-specific shared polymorphisms and negligible proportions of
334 fixed differences compared to the genomic background (Figure 4b,c; Table S6). It is
335 therefore likely that some of these regions have been targets of long-term balancing
336 selection in both species (Charlesworth 2006). Consistent with this prediction, we
337 found slightly lower or comparable levels of LD in these regions (Figure 4c; Table
338 S6), which is likely due to the long-term effects of recombination on old balanced
339 polymorphisms (Leffler et al. 2013). The higher d_{xy} and RND values we observe in
340 these regions may, however, be a consequence of the higher levels of ancestral
341 polymorphisms that were maintained predating the split of the two species (Figure 4b;
342 Table S6) (Cruickshank and Hahn 2014).

343

344 **Impact of recombination rate on patterns of genetic differentiation**

345 We examined relationships between the scaled recombination rates ($\rho=4N_e c$) and
346 levels of inter-species divergence over non-overlapping 10 Kbp windows across the
347 genome (Figure S11). We found a significant negative correlation between relative
348 divergence, measured as F_{ST} that depends on genetic diversity within species, and
349 recombination rates in both *P. tremula* (Spearman's $\rho=-0.121$, P -value <0.001) and *P.*
350 *tremuloides* (Spearman's $\rho=-0.157$, P -value <0.001) (Figure S11a). In contrast to F_{ST} ,
351 we observed significantly positive correlations between absolute divergence d_{xy} and
352 recombination rates in both species (*P. tremula*: Spearman's $\rho=0.199$, P -value <0.001 ;
353 *P. tremuloides*: Spearman's $\rho=0.140$, P -value <0.001) (Figure S11b).

354 Because $\rho=4N_e c$, where c is the per-generation recombination rate and N_e is
355 the effective population size, a reduction of N_e in regions linked to selection will
356 lower local estimates of ρ even if local c is identical to other regions in the genome. In
357 order to account for such effects and to obtain a measure of recombination that is
358 independent of local N_e , we compared ρ/θ_π between regions with extreme genetic
359 differentiation and the remainder of the genome. Relative to the genomic background,
360 our results showed significantly suppressed recombination in outlier regions
361 displaying either exceptionally high or low inter-specific differentiation (Figure 4c).

362

363 **Genes under selection**

364 The availability of the annotated *P. trichocarpa* genome enabled functional analyses
365 of candidate target genes within regions that were likely under selection. In total, 722
366 and 391 genes were located in outlier windows displaying exceptionally high and low
367 differentiation (Table S7 and S8), respectively. Compared to the genome overall, we
368 did not find significantly higher gene density in these outlier windows ($P>0.05$,
369 Mann-Whitney U test; Figure S12). We used the Gene Ontology (GO) assignments
370 of those candidate genes to assess whether specific GO terms were significantly over-
371 represented. After accounting for multiple comparisons, we did not detect over-
372 representation of any functional category among the candidate genes within regions of
373 high differentiation. However, we identified 60 significantly overrepresented GO
374 terms for genes located within regions showing significantly low genetic
375 differentiation and that were likely candidates for being under the influence of
376 balancing selection. Most of these GO categories were associated with immune and
377 defense responses, signal transduction or apoptosis (Table S9). Nevertheless, some
378 caution should be applied when interpreting these results since we observed a skewed
379 pattern of low coverage breadth in outlier windows displaying significantly low
380 differentiation compared to either the genomic background or to those highly
381 differentiated windows (Figure S13). Such unequal coverage breadth likely results
382 from the inherent technical hurdle of short-read sequencing technologies and likely
383 represents difficulties of mapping short reads to a reference genome in highly
384 polymorphic regions (Brandt et al. 2015). After stringent quality filtering, more reads
385 were discarded in these regions, which further decreased the amount of usable
386 information. Therefore, future studies, incorporating a combination of careful
387 experimental design and long-read sequencing technologies, are needed to evaluate
388 the accuracy of the evolutionary significance of balancing selection found in the
389 candidate genes shown here.

390

391 **Discussion**

392

393 We use a population genomic approach to resolve the evolutionary histories of two
394 widespread and closely related forest tree species, and to highlight how genome-wide
395 patterns of differentiation have been influenced by a variety of evolutionary
396 processes. Our simulation-based analyses indicate that *P. tremula* and *P. tremuloides*

397 diverged around 2.2-3.1 Mya during the Late Pliocene and/or Early Pleistocene. This
398 timing corresponds closely with the first opening of the Bering Strait, which occurred
399 3.1-5.5 Mya and broke up the overland intercontinental migration route of terrestrial
400 floras between Eurasia and North America (Marincovich and Gladenkov 1999;
401 Gladenkov et al. 2002). This may have been less of an immediate barrier to wind-
402 dispersed *Populus* than some other tree species, but the severing of the Bering land
403 bridge associated with the onset of dramatic climatic oscillations throughout the
404 Pleistocene were likely the principal drivers for initial divergence between *P. tremula*
405 and *P. tremuloides* (Comes and Kadereit 1998; Milne and Abbott 2002; Du et al.
406 2015). Given the modern-day geographic isolation, disjunct distributions and
407 extremely low rates of gene flow, our results support an allopatric model of speciation
408 for these two aspen species (Morjan and Rieseberg 2004). The coalescent-based,
409 intra-specific demographic analyses using MSMC demonstrate that both species have
410 experienced substantial population expansions following long-term declines after
411 species divergence. The population expansion of *P. tremuloides* has occurred over the
412 last 50,000-70,000 years, following the retreat of the penultimate glaciation and has
413 continued up to the present (Kaufman and Manley 2004). *P. tremula*, in contrast,
414 experienced a more extended population contraction and, consistent with many other
415 forest trees in Europe, the initiation of the population expansion in *P. tremula*
416 coincided with the end of the Last Glacial Maximum (Hewitt 2000; Hewitt 2004).

417 Consistent with the expectation for allopatric speciation, where the absence of
418 gene flow allowed for the accumulation of inter-specific differentiation due to
419 stochastic genetic drift (Coyne and Orr 2004), we detect widespread genomic
420 differentiation between the two species. Although neutral processes are likely
421 responsible for the majority of the genetic differentiation we observe between the two
422 species on a genome-wide scale (Coyne and Orr 2004; Strasburg et al. 2012), a
423 number of regions exhibit more extreme genetic differentiation compared to
424 expectations based on demographic simulations and these regions also show multiple
425 evidences of the action of natural selection (Nielsen et al. 2009). If natural selection
426 has truly been one of the dominant evolutionary forces shaping genetic differentiation
427 between the two species, regions of low recombination would be expected to show
428 increased F_{ST} values, but not increased d_{xy} values (Noor and Bennett 2009;
429 Cruickshank and Hahn 2014). This occurs because natural selection (through either
430 selective sweeps and/or background selection) removes neutral variation over longer

431 distances in regions of low recombination (Begun and Aquadro 1992). As a
432 consequence, relative measures of divergence (e.g. F_{ST}) that rely on within-species
433 diversity are expected to be higher in regions with restricted recombination (Noor and
434 Bennett 2009; Nachman and Payseur 2012). In contrast, increased absolute
435 divergence (e.g. d_{xy}) is only expected if reduced gene flow occurred in regions of low
436 recombination (Nachman and Payseur 2012). In accordance with this view, we
437 observe a significant negative relationship between population-scaled recombination
438 rates (ρ) and F_{ST} , but not with d_{xy} , in both species (Noor and Bennett 2009; Keinan et
439 al. 2010). Contrary to expectations of heterogeneous gene flow, our findings highlight
440 a significant effect of linked selection in generating the heterogeneous differentiation
441 landscape we observe between the two *Populus* species (Turner and Hahn 2010;
442 Cruickshank and Hahn 2014; Burri et al. 2015).

443 Rather than being physically clustered into just a few large, discrete genomic
444 ‘islands’ as expected when species diverge in the presence of gene flow (Turner et al.
445 2005), differentiation islands in our study system are distributed throughout the
446 genome, being narrowly defined and mostly located in regions with substantially
447 suppressed recombination. Combined with the multiple signatures of positive
448 selection we observe in these regions, these islands of divergence likely represent
449 candidate regions harboring loci closely tied to species-specific adaptation rather than
450 loci that are “resistant” to gene flow (Turner and Hahn 2010). In addition, we find no
451 functional over-representation for candidate genes located in these regions, suggesting
452 that a wide range of genes and functional categories have likely been involved in the
453 adaptive evolution of these two species (Wolf et al. 2010).

454 In addition to the highly differentiated regions that show signs of species-
455 specific positive selection, we also identify a number of lowly differentiated regions
456 that are candidates for being affected by balancing selection in both species
457 (Charlesworth 2006). Apart from low inter-specific divergence and high intra-specific
458 diversity, these regions also contain an excess of sites at intermediate frequencies, a
459 greater proportion of shared polymorphisms between the two species and lack of fixed
460 inter-specific differences. Genes within these regions are enriched for immune and
461 defense response, signal transduction and apoptosis, highlighting the influence of co-
462 evolutionary arms races between hosts and natural enemies on the persistence of
463 functional genetic diversity in immunity and defense-related genes (Tiffin and

464 Moeller 2006; Salvaudon et al. 2008). That said, caution is required when interpreting
465 the functional properties of the regions identified here, and future studies of these
466 candidate genes are clearly needed to better assess the adaptive genetic potential of
467 these widespread forest tree species to current and future climate change.

468 A number of factors may have influenced our demographic inferences and the
469 detection of natural selection in the two species. First, the presence of within-species
470 population subdivision could have magnified the inference of demographic expansion
471 in *P. tremuloides*, because pooling samples from populations of Alberta and
472 Wisconsin skews the SFS toward low-frequency polymorphism (more negative
473 Tajima's D) (Figure S4) and results in larger estimates of effective population sizes
474 than estimates obtained from local samples (Figure S3). However, all our analyses
475 suggest that this confounding effect is very weak (see Results) and the divergence
476 between the two subpopulations of *P. tremuloides* likely is too recent to have any
477 major effects on the demographic reconstruction and tests of selection in these species
478 (Chikhi et al. 2010). Another caveat concerns the sampling scheme used in the two
479 species. Local samples in *P. tremula* may not adequately reflect species-wide
480 demography compared to the pooled samples in *P. tremuloides*. However, the extent
481 to which this might influence the estimates of inter-specific F_{ST} deserves further
482 study. More generally, sampling should likely be more extensive in both species to
483 capture a greater proportion of the species-wide diversity, although local sampling is
484 expected to only have small effects in species with high gene flow like *Populus*
485 (Wakeley 2000). Finally, inter-specific hybridization in either species could
486 potentially bias our results. However, there are no other species of *Populus* occurring
487 in regions where the *P. tremula* individuals were sampled. For *P. tremuloides*,
488 naturally occurring hybridization is only known to occur with *P. grandidentata* in
489 central and eastern North America where the two species co-occur (Pregitzer and
490 Barnes 1980). Therefore, any possible hybridization in our study would be limited to
491 samples from the Wisconsin population of *P. tremuloides* but, as noted above, we did
492 not detect any major differences in patterns of genetic variation between the two
493 subpopulations, suggesting little or no effect of hybridization.

494

495 **Conclusion**

496

497 Here we provide insights into the speciation and recent evolutionary histories of two
498 closely related forest tree species, *P. tremula* and *P. tremuloides*. Our study supports
499 an allopatric model of speciation for the two species, which are estimated to have
500 diverged around 2.2-3.1 Mya as a result of the first opening of Bering Strait.
501 Coalescent simulations suggest that stochastic genetic drift and historical
502 demographic processes can largely explain the genome-wide patterns of
503 differentiation between species. However, there is an excess of regions displaying
504 extreme inter-specific genetic differentiation in the observed data compared with
505 demographic simulations, which is likely indicative of the action of natural selection.
506 In addition, we find that heterogeneous genomic divergence is strongly driven by
507 linked selection and variation in recombination rate in the two species. Instead of
508 being clustered into a few large genomic “islands” as is expected under a model of
509 speciation with gene flow, regions of pronounced differentiation are characterized by
510 multiple signatures of positive selection in one or both species, and are distributed
511 throughout the genome at many small, independent locations. Regions displaying
512 exceptionally low differentiation are likely candidate targets of long-term balancing
513 selection, which are strongly enriched for genes involved in immune and defense
514 response, signal transduction and apoptosis, suggesting a possible link to long-term
515 co-evolutionary arms races with pest and pathogens. Our study highlights that future
516 work should integrate more information on the natural histories of speciation, such as
517 divergence time, geographical context, magnitudes of gene flow, demographic
518 histories and sources of adaptation, when interpreting the meaning of observed
519 patterns of genomic divergence between closely related species.

520

521 **Materials and Methods**

522

523 **Population samples, sequencing, quality control and mapping**

524 A total of 24 individuals of *P. tremula* and 22 individuals of *P. tremuloides* were
525 collected and sequenced (Figure 1a and Table S1). For each individual, genomic DNA
526 was extracted from leaf samples and paired-end sequencing libraries with an insert
527 size of 600bp were constructed according to the Illumina library preparation protocol.
528 Sequencing was carried out on the Illumina HiSeq 2000 platform at the Science for
529 Life Laboratory in Stockholm, Sweden. All samples were sequenced to a target

530 coverage of 25_x. The sequencing data has been deposited in the Short Read Archive
531 (SRA) at NCBI under accession numbers SRP065057 and SRP065065 for samples of
532 *P. tremula* and *P. tremuloides*, respectively.

533 For raw sequencing reads (Wang et al. 2015), we used Trimmomatic (Lohse et
534 al. 2012) to remove adapter sequences and cut off bases from either the start or the
535 end of reads when the base quality was lower than 20. Reads were completely
536 discarded if there were fewer than 36 bases remaining after trimming. We then
537 mapped all reads to the *P. trichocarpa* reference genome (v3.0) (Tuskan et al. 2006)
538 with default parameters implemented in bwa-0.7.10 using the BWA-MEM algorithm
539 (Li unpublished data, <http://arxiv.org/abs/1303.3997>, last accessed May 26, 2013).
540 Local realignment was performed to correct for the mis-alignment of bases in regions
541 around insertions and/or deletions (indels) using RealignerTargetCreator and
542 IndelRealigner in GATK v3.2.2 (DePristo et al. 2011). In order to account for the
543 occurrence of PCR duplicates introduced during library construction, we used
544 MarkDuplicates in Picard (<http://picard.sourceforge.net>) to remove reads with
545 identical external coordinates and insert lengths. Only the read with the highest
546 summed base quality was kept for downstream analyses.

547

548 **Filtering sites**

549 Prior to variant and genotype calling, we employed several filtering steps to exclude
550 potential errors caused by paralogous or repetitive DNA sequences. First, after
551 investigating the empirical distribution, we removed sites showing extremely low
552 (<100 reads across all samples per species) or high (>1200 reads across all samples
553 per species) read coverage. Second, as a mapping quality score of zero is assigned for
554 reads that could be equally mapped to multiple genomic locations, we removed sites
555 containing more than 20 such reads among all samples in each species. Third, we
556 removed sites that overlapped with known repeat elements as identified by
557 RepeatMasker (Tarailo-Graovac and Chen 2009). After all filtering steps, there were
558 42.8% of sites across the genome left for downstream analyses. Among them, 54.9%
559 were found within gene boundaries, and the remainder (45.1%) was located in
560 intergenic regions.

561

562 **SNP and genotype calling**

563 We employed two complementary approaches for SNP and genotype calling (Figure
564 S1): (i) Direct estimation without calling genotypes was implemented in the software
565 ANGSD v0.602 (Korneliussen et al. 2014). Only reads with a minimal mapping
566 quality of 30 and bases with a minimal quality score of 20 were considered. For all
567 filtered sites in both species, we defined the alleles that were the same as those found
568 in the *P. trichocarpa* reference genome as the ancestral allelic state. We used the -
569 doSaf implementation to calculate the site allele frequency likelihood based on the
570 SAMTools genotype likelihood model at all sites (Li et al. 2009), and then used the -
571 realSFS implementation to obtain a maximum likelihood estimate of the unfolded
572 SFS using the Expectation Maximization (EM) algorithm (Kim et al. 2011). Several
573 population genetic statistics were then calculated based on the global SFS (Figure S1).
574 (ii) Multi-sample SNP and genotype calling was implemented in GATK v3.2.2 with
575 HaplotypeCaller and GenotypeGVCFs (Figure S1) (DePristo et al. 2011). A number
576 of filtering steps were performed to reduce false positives from SNP and genotype
577 calling: (1) Remove SNPs that were located in regions not passing all previous
578 filtering criteria; (2) Removed SNPs with more than 2 alleles in both species; (3)
579 Removed SNPs at or within 5bp from any indels; (4) Assigned genotypes as missing
580 if their quality scores (GQ) were lower than 10, and then removed SNPs with more
581 than two missing genotypes in each species; (5) Removed SNPs showing significant
582 deviation from Hardy-Weinberg Equilibrium ($P < 0.001$) in each species.

583

584 **Population structure**

585 Population genetic structure was inferred using the program NGSadmix (Skotte et al.
586 2013), with only sites containing lower than 10% of missing data being used. We used
587 the SAMTools model (Li et al. 2009) in ANGSD to estimate genotype likelihoods and
588 then generated a beagle file for the subset of the genome that was determined as being
589 variable using a likelihood ratio test (P -value $< 10^{-6}$) (Kim et al. 2011). We predefined
590 the number of genetic clusters K from 2 to 5, and the maximum iteration of the EM
591 algorithm was set to 10,000.

592 As another method to visualize the genetic relationships among individuals,
593 we performed principal component analysis (PCA) using ngsTools accounting for
594 sequencing errors and uncertainty in genotype calls (Fumagalli et al. 2014). The

595 expected covariance matrix across pairs of individuals in both species was computed
596 based on the genotype posterior probabilities across all filtered sites. Eigenvectors and
597 eigenvalues from the covariance matrix were generated with the R function `eigen`, and
598 significance levels were determined using the Tracy-Widom test as implemented in
599 EIGENSOFT version 4.2 (Patterson et al. 2006).

600

601 **Demographic history**

602 We inferred demographic histories associated with speciation for *P. tremula* and *P.*
603 *tremuloides* using a coalescent simulation-based method implemented in *fastsimcoal*
604 2.5.1 (Excoffier et al. 2013). Two-dimensional joint site frequency spectrum (2D-SFS)
605 was constructed from posterior probabilities of sample allele frequencies by *ngsTools*
606 (Fumagalli et al. 2014). 100,000 coalescent simulations were used for the estimation
607 of the expected 2D-SFS and log-likelihood for a set of demographic parameters in
608 each model. Global maximum likelihood estimates for each model were obtained
609 from 50 independent runs, with 10-40 conditional maximization algorithm cycles.
610 Model comparison was based on the maximum value of likelihood over the 50
611 independent runs using the Akaike information criterion (AIC) and Akaike's weight
612 of evidence (Excoffier et al. 2013). The model with the maximum Akaike's weight
613 value was chosen as the optimal one. We assumed a mutation rate of 2.5×10^{-9} per site
614 per year and a generation time of 15 years in *Populus* (Koch et al. 2000) when
615 converting estimates to units of years and individuals. Parameter confidence intervals
616 of the best model were obtained by 100 parametric bootstraps, with 50 independent
617 runs in each bootstrap.

618 We then employed multiple sequential Markovian coalescent (MSMC)
619 method to estimate variation of scaled population sizes (N_e) over historical time in
620 both species (Schiffels and Durbin 2014), which is an extension of a pairwise
621 sequential Markovian coalescent (PSMC) method (Li and Durbin 2011). Prior to the
622 analysis, all segregating sites within each species were phased and imputed using
623 *fastPHASE* v1.4.0 (Scheet and Stephens 2006). A generation time of 15 years and a
624 rate of 2.5×10^{-9} mutations per nucleotide per year (Koch et al. 2000) were used to
625 covert the scaled times and population sizes into real times and sizes.

626

627 **Genome-wide patterns of differentiation**

628 We have previously shown that linkage disequilibrium (LD) decays within 10
629 kilobases (Kbp) in both *P. tremula* and *P. tremuloides* (Wang et al. forthcoming), and
630 we thus divided the genome into 39,406 non-overlapping windows of 10 Kbp in size
631 to investigate patterns of genomic differentiation between species. For a window to be
632 included in the downstream analyses, we required there to be at least 1000 bases left
633 after all above filtering steps. Levels of genetic differentiation between species at each
634 site were estimated using method-of-moments F_{ST} estimators implemented in ngsFST
635 from the ngsTools package (Fumagalli et al. 2014), which calculates indices of the
636 expected genetic variance between and within species from posterior probabilities of
637 sample allele frequencies, without relying on SNP or genotype calling (Fumagalli et
638 al. 2013). We then averaged F_{ST} values across all sites within each 10 Kbp non-
639 overlapping window.

640

641 **Coalescent simulations for detecting outlier windows**

642 In order to examine thresholds for detection of outlier windows that may have been
643 targets of natural selection, we conducted coalescent simulations to compare observed
644 patterns of genetic differentiation (F_{ST}) to those expected under different demographic
645 models (see Results). All simulations were performed using the program *msms* v3.2rc
646 (Ewing and Hermisson 2010) based on demographic parameters derived from the
647 best-fitting model inferred by *fastsimcoal2.5.1* (Excoffier et al. 2013). Population-
648 scaled recombination rates (ρ) were assumed to be between 1 Kbp⁻¹ and 5 Kbp⁻¹ given
649 the large variation we found in both species (Wang et al. forthcoming). We simulated
650 genotypes corresponding to a 10 Kbp region with the same sample size as the real
651 data for 100,000 replications, from where we simulate genotype likelihoods using the
652 program msToGlf in ANGSD (Korneliussen et al. 2014) by assuming a mean
653 sequencing depth of 20_x and an error rate of 0.5%. We estimated two summary
654 statistics, nucleotide diversity (θ_{π}) and Tajima's D, from sample allele frequency
655 likelihoods in ANGSD for all simulation replicates to test whether the simulated data
656 matches the observed data. To assess whether any of the observed windows display
657 F_{ST} values deviating significantly from neutral expectations, we determined the
658 conditional probability (*P*-value) of observing more extreme inter-specific F_{ST} values
659 among simulated data sets than among the observed data. Our significance was based
660 on running 500,000 coalescent simulations of the most acceptable demographic null

661 model (see Results). We then corrected for multiple testing by using the False
662 Discovery Rate (FDR) adjustment, and only windows with FDR lower than 1% were
663 considered as candidate regions affected by selection (Storey 2002).

664

665 **Molecular signatures of selection in outlier regions**

666 To assess the occurrence of selection in outlier windows displaying either
667 exceptionally high or low differentiation, we compared these two unions of outlier
668 windows to the remaining portion of the genome by a variety of additional population
669 genetic statistics in both species. First, θ_π , Tajima's D (Tajima 1989) and Fay & Wu's
670 H (Fay and Wu 2000) were calculated from sample allele frequency likelihoods in
671 ANGSD over non-overlapping 10 Kbp windows. Second, levels of LD and
672 population-scaled recombination rates (ρ) were estimated and compared. To evaluate
673 levels of LD within each 10 Kbp window, the correlation coefficients (r^2) between
674 SNPs with pairwise distances larger than 1 Kbp were calculated using VCFtools
675 v0.1.12b (Danecek et al. 2011). Population-scaled recombination rates (ρ) were
676 estimated using the Interval program of LDhat 2.2 (McVean et al. 2004) with
677 1,000,000 MCMC iterations sampling every 2,000 iterations and a block penalty
678 parameter of five. The first 100,000 iterations of the MCMC iterations were discarded
679 as burn-in. Resulting estimates of r^2 and ρ were averaged over each 10 Kbp window.
680 In both species, windows were discarded in the estimation of r^2 and ρ if there were
681 less than 3 Kbp and/or 10 SNPs left from previous filtering steps. Finally, we used the
682 program ngsStat (Fumagalli et al. 2014) to calculate several additional measures of
683 genetic differentiation: (1) with *P. trichocarpa* as an outgroup, the proportion of fixed
684 differences that is caused by derived alleles fixed in either *P. tremula* or *P.*
685 *tremuloides* among all segregating sites; (2) the proportion of inter-species shared
686 polymorphisms among all segregating sites; (3) d_{xy} , which was calculated from
687 sample allele frequency posterior probabilities at each site and was then averaged over
688 each 10 Kbp window; and (4) the relative node depth (RND), which was calculated
689 by dividing the d_{xy} of the two aspen species with d_{xy} between aspen (represented by
690 24 samples of *P. tremula* in this study) and *P. trichocarpa* (24 samples; see Wang et
691 al. forthcoming). Significance of the differences between outlier windows and the
692 genome-wide averages for all above mentioned population genetic statistics were
693 examined using one-sided Wilcoxon ranked-sum tests.

694

695 **Gene ontology (GO) enrichment**

696 To determine whether any functional classes of genes were overrepresented among
697 regions that were candidates for being under selection, we performed functional
698 enrichment analysis of GO using Fisher's exact test by agriGO's Term Enrichment
699 tool (<http://bioinfo.cau.edu.cn/agriGO/index.php>) (Du et al. 2010). GO groups with
700 fewer than two outlier genes were excluded from this analysis. *P*-values of Fisher's
701 exact test were further corrected for multiple testing with Benjamini-Hochberg false
702 discovery rate (Benjamini and Hochberg 1995). GO terms with a corrected *P*-value
703 <0.05 were considered to be significantly enriched.

704

705 **Acknowledgements**

706

707 We are grateful to Rick Lindroth for providing access to the samples of *P.*
708 *tremuloides* used in this study. We thank Carin Olofsson for extracting DNA for all
709 samples used in this study. We thank both the editor and two anonymous referees for
710 their useful comments on the manuscript. The research has been funded through
711 grants from Vetenskapsrådet and a Young Researcher Award from Umeå University
712 to PKI. JW was supported by a scholarship from the Chinese Scholarship Council.
713 The authors also would like to acknowledge support from Science for Life
714 Laboratory, the National Genomics Infrastructure (NGI), and Uppmax for providing
715 assistance in massive parallel sequencing and computational infrastructure.

716

717 **References**

718

- 719 Avise JC. 2000. Phylogeography: the history and formation of species. Cambridge:
720 Harvard University Press.
- 721 Benjamini Y, Hochberg Y. 1995. Controlling the false discovery rate: a practical and
722 powerful approach to multiple testing. *J R Stat Soc Series B Methodol.* 57:289-
723 300.
- 724 Brandt DY, Aguiar VR, Bitarello BD, Nunes K, Goudet J, Meyer D. 2015. Mapping
725 bias overestimates reference allele frequencies at the HLA genes in the 1000
726 genomes project phase I data. *G3.* 5:931-941.
- 727 Burri R, Nater A, Kawakami T, Mugal CF, Olason PI, Smeds L, Suh A, Dutoit L,
728 Bureš S, Garamszegi LZ, et al. 2015. Linked selection and recombination rate
729 variation drive the evolution of the genomic landscape of differentiation across
730 the speciation continuum of *Ficedula* flycatchers. *Genome Res.* 25:1656-1665.

- 731 Campagna L, Gronau I, Silveira LF, Siepel A, Lovette IJ. 2015. Distinguishing noise
732 from signal in patterns of genomic divergence in a highly polymorphic avian
733 radiation. *Mol Ecol.* 24:4238-4251.
- 734 Carneiro M, Albert F, Afonso S, Pereira R, Burbano H, Campos R, Melo-Ferreira J,
735 Blanco-Aguiar J, Villafuerte R, Nachman M, et al. 2014. The genomic
736 architecture of population divergence between subspecies of the European
737 Rabbit. *PLoS Genet.* 10:e1003519.
- 738 Charlesworth B. 1998. Measures of divergence between populations and the effect of
739 forces that reduce variability. *Mol Biol Evol.* 15:538-543.
- 740 Charlesworth D. 2006. Balancing selection and its effects on sequences in nearby
741 genome regions. *PLoS Genet.* 2:e64.
- 742 Chikhi L, Sousa VC, Luisi P, Goossens B, Beaumont MA. 2010. The confounding
743 effects of population structure, genetic diversity and the sampling scheme on
744 the detection and quantification of population size changes. *Genetics* 186:983-
745 995.
- 746 Comes HP, Kadereit JW. 1998. The effect of Quaternary climatic changes on plant
747 distribution and evolution. *Trends Plant Sci.* 3:432-438.
- 748 Coyne JA, Orr HA. 2004. Speciation. Sunderland: Sinauer Associates.
- 749 Cruickshank TE, Hahn MW. 2014. Reanalysis suggests that genomic islands of
750 speciation are due to reduced diversity, not reduced gene flow. *Mol Ecol.*
751 23:3133-3157.
- 752 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,
753 Lunter G, Marth GT, Sherry ST, et al. 2011. The variant call format and
754 VCFtools. *Bioinformatics* 27:2156-2158.
- 755 DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, Philippakis
756 AA, Del Angel G, Rivas MA, Hanna M, et al. 2011. A framework for
757 variation discovery and genotyping using next-generation DNA sequencing
758 data. *Nat Genet.* 43:491-498.
- 759 Du S, Wang Z, Ingvarsson PK, Wang D, Wang J, Wu Z, Tembrock LR, Zhang J.
760 2015. Multilocus analysis of nucleotide variation and speciation in three
761 closely related *Populus* (Salicaceae) Species. *Mol Ecol.* 24:4994-5005.
- 762 Du Z, Zhou X, Ling Y, Zhang Z, Su Z. 2010. agriGO: a GO analysis toolkit for the
763 agricultural community. *Nucleic Acids Res.* 38:W64-W70.
- 764 Eckenwalder JE. 1996. Systematics and evolution of *Populus*. In Stettler RF,
765 Bradshaw HD, Heilman PE, Hinckley TM, editors. *Biology of Populus and its*
766 *Implications for Management and Conservation.* Ottawa: NRC Research
767 Press.p.7-32.
- 768 Ellegren H, Smeds L, Burri R, Olason PI, Backström N, Kawakami T, Künstner A,
769 Mäkinen H, Nadachowska-Brzyska K, Qvarnström A, et al. 2012. The
770 genomic landscape of species divergence in *Ficedula* flycatchers. *Nature*
771 491:756-760.
- 772 Ewing G, Hermisson J. 2010. MSMS: a coalescent simulation program including
773 recombination, demographic structure and selection at a single locus.
774 *Bioinformatics* 26:2064-2065.
- 775 Excoffier L, Dupanloup I, Huerta-Sánchez E, Sousa V, Foll M. 2013. Robust
776 demographic inference from genomic and SNP data. *PLoS Genet.* 9:e1003905.
- 777 Fay JC, Wu C-I. 2000. Hitchhiking under positive Darwinian selection. *Genetics*
778 155:1405-1413.

- 779 Feder JL, Xie X, Rull J, Velez S, Forbes A, Leung B, Dambroski H, Filchak KE,
780 Aluja M. 2005. Mayr, Dobzhansky, and Bush and the complexities of
781 sympatric speciation in *Rhagoletis*. *Proc Natl Acad Sci USA*. 102:6573-6580.
- 782 Feulner P, Chain F, Panchal M, Huang Y, Eizaguirre C, Kalbe M, Lenz T, Samonte I,
783 Stoll M, Bornberg-Bauer E, et al. 2015. Genomics of divergence along a
784 continuum of parapatric population differentiation. *PLoS Genet*. 11:e1004966.
- 785 Fumagalli M, Vieira FG, Korneliussen TS, Linderoth T, Huerta-Sánchez E,
786 Albrechtsen A, Nielsen R. 2013. Quantifying population genetic
787 differentiation from next-generation sequencing data. *Genetics* 195:979-992.
- 788 Fumagalli M, Vieira FG, Linderoth T, Nielsen R. 2014. ngsTools: methods for
789 population genetics analyses from next-generation sequencing data.
790 *Bioinformatics* 30:1486-1487.
- 791 Gladenkov AY, Oleinik AE, Marincovich L, Barinov KB. 2002. A refined age for the
792 earliest opening of Bering Strait. *Palaeogeogr Palaeoclimatol Palaeoecol*.
793 183:321-328.
- 794 Hamzeh M, Dayanandan S. 2004. Phylogeny of *Populus* (Salicaceae) based on
795 nucleotide sequences of chloroplast trnT-trnF region and nuclear rDNA. *Am J*
796 *Bot*. 91:1398-1408.
- 797 Hewitt G. 2004. Genetic consequences of climatic oscillations in the Quaternary.
798 *Philos Trans R Soc Lond B Biol Sci*. 359:183-195.
- 799 Hewitt G. 2000. The genetic legacy of the Quaternary ice ages. *Nature* 405:907-913.
- 800 Kaufman D, Manley W. 2004. Quaternary glaciations—Extent and chronology Part
801 II: North America. In Ehlers J, Gibbard PL. editors. *Developments in*
802 *Quaternary Science*.p.1-440.
- 803 Keinan A, Reich D, Begun DJ. 2010. Human population differentiation is strongly
804 correlated with local recombination rate. *PLoS Genet*. 6:e1000886.
- 805 Kim SY, Lohmueller KE, Albrechtsen A, Li Y, Korneliussen T, Tian G, Grarup N,
806 Jiang T, Andersen G, Witte D, et al. 2011. Estimation of allele frequency and
807 association mapping using next-generation sequencing data. *BMC*
808 *bioinformatics* 12:231.
- 809 Koch MA, Haubold B, Mitchell-Olds T. 2000. Comparative evolutionary analysis of
810 chalcone synthase and alcohol dehydrogenase loci in *Arabidopsis*, *Arabis*,
811 and related genera (Brassicaceae). *Mol Biol Evol*. 17:1483-1498.
- 812 Korneliussen TS, Albrechtsen A, Nielsen R. 2014. ANGSD: analysis of next
813 generation sequencing data. *BMC bioinformatics* 15:356.
- 814 Leffler EM, Gao Z, Pfeifer S, Ségurel L, Auton A, Venn O, Bowden R, Bontrop R,
815 Wall JD, Sella G, et al. 2013. Multiple instances of ancient balancing
816 selection shared between humans and chimpanzees. *Science* 339:1578-1582.
- 817 Li H, Durbin R. 2011. Inference of human population history from individual whole-
818 genome sequences. *Nature* 475:493-496.
- 819 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G,
820 Durbin R. 2009. The sequence alignment/map format and SAMtools.
821 *Bioinformatics* 25:2078-2079.
- 822 Lohse M, Bolger A, Nagel A, Fernie AR, Lunn JE, Stitt M, Usadel B. 2012. RobiNA:
823 a user-friendly, integrated software solution for RNA-Seq-based
824 transcriptomics. *Nucleic Acids Res*. 40: W622-W627.
- 825 Luikart G, England PR, Tallmon D, Jordan S, Taberlet P. 2003. The power and
826 promise of population genomics: from genotyping to genome typing. *Nat*
827 *Rev Genet*. 4:981-994.

- 828 Marincovich L, Gladenkov AY. 1999. Evidence for an early opening of the Bering
829 Strait. *Nature* 397:149-151.
- 830 McVean GA, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. 2004. The
831 fine-scale structure of recombination rate variation in the human genome.
832 *Science* 304:581-584.
- 833 Milne RI, Abbott RJ. 2002. The origin and evolution of Tertiary relict floras. *Adv Bot*
834 *Res.* 38:281-314.
- 835 Morjan CL, Rieseberg LH. 2004. How species evolve collectively: implications of
836 gene flow and selection for the spread of advantageous alleles. *Mol Ecol.*
837 13:1341-1356.
- 838 Müller A, Leuschner C, Horna V, Zhang C. 2012. Photosynthetic characteristics and
839 growth performance of closely related aspen taxa: on the systematic
840 relatedness of the Eurasian *Populus tremula* and the North American *P.*
841 *tremuloides*. *Flora* 207:87-95.
- 842 Nachman MW, Payseur BA. 2012. Recombination rate variation and speciation:
843 theoretical predictions and empirical results from rabbits and mice. *Philos*
844 *Trans R Soc Lond B Biol Sci* 367:409-421.
- 845 Neale DB, Kremer A. 2011. Forest tree genomics: growing resources and
846 applications. *Nat Rev Genet.* 12:111-122.
- 847 Nei M. 1987. Molecular evolutionary genetics. New York: Columbia University
848 Press.
- 849 Nielsen R, Hubisz MJ, Hellmann I, Torgerson D, Andrés AM, Albrechtsen A,
850 Gutenkunst R, Adams MD, Cargill M, Boyko A, et al. 2009. Darwinian and
851 demographic forces affecting human protein coding genes. *Genome Res.*
852 19:838-849.
- 853 Nielsen R, Korneliussen T, Albrechtsen A, Li Y, Wang J. 2011. SNP calling,
854 genotype calling, and sample allele frequency estimation from New-
855 Generation Sequencing data. *PLoS One* 7:e37558.
- 856 Noor MA, Bennett SM. 2009. Islands of speciation or mirages in the desert?
857 Examining the role of restricted recombination in maintaining species.
858 *Heredity* 103:439-444.
- 859 Nosil P, Feder JL. 2012. Genomic divergence during speciation: causes and
860 consequences. *Philos Trans R Soc Lond B Biol Sci* 367:332-342.
- 861 Nosil P, Funk DJ, Ortiz-Barrientos D. 2009. Divergent selection and heterogeneous
862 genomic divergence. *Mol Ecol.* 18:375-402.
- 863 Pakull B, Groppe K, Meyer M, Markussen T, Fladung M. 2009. Genetic linkage
864 mapping in aspen (*Populus tremula* L. and *Populus tremuloides* Michx.). *Tree*
865 *Genet Genomes* 5:505-515.
- 866 Patterson N, Price A, Reich D. 2006. Population structure and eigenanalysis. *PLoS*
867 *Genet.* 2:e190.
- 868 Pregitzer KS, Barnes BV. 1980. Flowering phenology of *Populus tremuloides* and *P.*
869 *grandidentata* and the potential for hybridization. *Can J Forest Res.* 10:218-
870 223.
- 871 Renaut S, Grassa C, Yeaman S, Moyers B, Lai Z, Kane N, Bowers J, Burke J,
872 Rieseberg L. 2013. Genomic islands of divergence are not affected by
873 geography of speciation in sunflowers. *Nature Commun.* 4:1827.
- 874 Robinson KM, Delhomme N, Mähler N, Schiffthaler B, Önskog J, Albrechtsen BR,
875 Ingvarsson PK, Hvidsten TR, Jansson S, Street NR. 2014. *Populus tremula*
876 (European aspen) shows no evidence of sexual dimorphism. *BMC Plant Biol.*
877 14:276.

- 878 Salvaudon L, Giraud T, Shykoff JA. 2008. Genetic diversity in natural populations: a
879 fundamental component of plant–microbe interactions. *Curr Opin Plant Biol.*
880 11:135-143.
- 881 Scheet P, Stephens M. 2006. A fast and flexible statistical model for large-scale
882 population genotype data: applications to inferring missing genotypes and
883 haplotypic phase. *Am J Hum Genet.* 78:629-644.
- 884 Schiffels S, Durbin R. 2014. Inferring human population size and separation history
885 from multiple genome sequences. *Nat Genet.* 46:919-925.
- 886 Seehausen O, Butlin RK, Keller I, Wagner CE, Boughman JW, Hohenlohe PA,
887 Peichel CL, Saetre G-P, Bank C, Brännström Å, *et al.* 2014. Genomics and the
888 origin of species. *Nat Rev Genet.* 15:176-192.
- 889 Skotte L, Korneliussen TS, Albrechtsen A. 2013. Estimating individual admixture
890 proportions from next generation sequencing data. *Genetics* 195:693-702.
- 891 Slatkin M. 1985. Gene flow in natural populations. *Annu Rev Ecol Syst.* 16:393-430.
- 892 Sousa V, Hey J. 2013. Understanding the origin of species with genome-scale data:
893 modelling gene flow. *Nat Rev Genet.* 14:404-414.
- 894 Storey JD. 2002. A direct approach to false discovery rates. *J R Stat Soc Series B Stat*
895 *Methodol.* 64:479-498.
- 896 Strasburg JL, Sherman NA, Wright KM, Moyle LC, Willis JH, Rieseberg LH. 2012.
897 What can patterns of differentiation across plant genomes tell us about
898 adaptation and speciation? *Philos Trans R Soc Lond B Biol Sci.* 367:364-373.
- 899 Tajima F. 1989. Statistical method for testing the neutral mutation hypothesis by
900 DNA polymorphism. *Genetics* 123:585-595.
- 901 Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to Identify Repetitive
902 Elements in Genomic Sequences. *Curr Protoc Bioinformatics* 4:1-4.10.
- 903 Tiffin P, Moeller DA. 2006. Molecular evolution of plant immune system genes.
904 *Trends Genet.* 22:662-670.
- 905 Tullus A, Rytter L, Tullus T, Weih M, Tullus H. 2012. Short-rotation forestry with
906 hybrid aspen (*Populus tremula* L. × *P. tremuloides* Michx.) in Northern
907 Europe. *Scand J Forest Res.* 27:10-29.
- 908 Turner T, Hahn M, Nuzhdin S. 2005. Genomic islands of speciation in *Anopheles*
909 *gambiae*. *PLoS Biol.* 3:e285.
- 910 Turner TL, Hahn MW. 2010. Genomic islands of speciation or genomic islands and
911 speciation? *Mol Ecol.* 19:848-850.
- 912 Tuskan G, Difazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N,
913 Ralph S, Rombauts S, Salamov A, *et al.* 2006. The genome of black
914 cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* 313:1596-1604.
- 915 Via S. 2009. Natural selection in action during speciation. *Proc Natl Acad Sci USA.*
916 106:9939-9946.
- 917 Wakeley J. 2000. The effects of subdivision on the genetic divergence of populations
918 and species. *Evolution* 54:1092-1101.
- 919 Wang J, Scofield D, Street N, Ingvarsson P. 2015. Variant calling using NGS data in
920 European aspen (*Populus tremula*). In: Sablok G, Kumar S, Ueno S, Kuo J,
921 Varotto C, editors. *Advances in the Understanding of Biological Sciences*
922 *Using Next Generation Sequencing (NGS) Approaches.* Springer. p.43-61.
- 923 Wang J, Street NR, Scofield DG, Ingvarsson PK. forthcoming. Natural selection and
924 recombination rate variation shape nucleotide polymorphism across the
925 genomes of three related *Populus* species. *Genetics* doi: genetics.115.183152.
- 926 Wang Z, Du S, Dayanandan S, Wang D, Zeng Y, Zhang J. 2013. Phylogeny
927 reconstruction and hybrid analysis of *Populus* (salicaceae) based on nucleotide

928 sequences of multiple single-copy nuclear genes and plastid fragments. *PloS*
929 *One* 9:e103645.

930 Wolf JB, Lindell J, Backström N. 2010. Speciation genetics: current status and
931 evolving approaches. *Philos Trans R Soc Lond B Biol Sci.* 365:1717-1733.

932

933

934

935

936

937

938

939

940

941

942

943

944

945

946

947

948

949

950

951

952

953

954

955

956

957

958 **Figures and Tables**

959

960 **Figure 1. Geographic distribution and genetic structure of 24 *Populus tremula***

961 **and 22 *P. tremuloides*.** (a) Map showing the current geographic distribution of *P.*
962 *tremula* (red) and *P. tremuloides* (blue). Yellow circles and triangles indicate the
963 locations where the 24 individuals of *P. tremula* and 22 individuals of *P. tremuloides*
964 were sampled. (b) Genetic structure of the two species inferred using NGSadmix. The
965 y-axis quantifies subgroup membership, and the x-axis shows the sample ID for each
966 individual. (c) Principal component analysis (PCA) plot based on genetic covariance
967 among all individuals of *P. tremula* (red circle) and *P. tremuloides* (green square and
968 blue triangle). The first two principle components (PCs) are shown, with PC1
969 explaining 21.04% ($P=2.51 \times 10^{-19}$, Tacey-Widom test) of the overall genetic variation
970 and separating the two species and PC2 explaining 2.09% ($P=9.65 \times 10^{-4}$, Tracy-
971 Widom test) of the overall variation and separating samples from Wisconsin (blue
972 triangle) and Alberta (green square) in *P. tremuloides*. (d) Pie chart summarizing the
973 proportion of fixed, shared and exclusive polymorphisms of the two species.

974

975 **Figure 2. Demographic history of *Populus tremula* and *P. tremuloides*.** (a)

976 Simplified graphical summary of the best-fitting demographic model inferred by
977 *fastsimcoal2*. (b) Multiple sequential Markovian coalescent (MSMC) estimates of the
978 effective population size (N_e) changes for *P. tremula* (red line) and *P. tremuloides*
979 (blue line) based on the inference from two (dashed), four (dotted) and eight (solid)
980 phased haplotypes in each species. Time scale on the x-axis is calculated assuming a
981 neutral mutation rate per generation (μ) = 3.75×10^{-8} and generation time (g) = 15
982 years. The grey bar indicates the speciation time inferred by *fastsimcoal2*.

983

984 **Figure 3. Genome-wide divergence.** Chromosomal distribution of genetic

985 differentiation (F_{ST}) between *Populus tremula* and *P. tremuloides*. The small, light
986 blue dots indicate F_{ST} values estimated over 10 Kbp non-overlapping windows. Grey
987 lines indicate F_{ST} values estimated over 100 Kbp non-overlapping windows.
988 Locations for windows displaying extreme differentiation relative to demographic
989 simulations are highlighted with colored bars above the plot. Among them, candidate
990 windows displaying significantly high differentiation (orange bars) are located on the

991 topside; candidate windows displaying significantly low differentiation (green bars)
992 are located at the bottom.

993

994 **Figure 4.** Identification of outlier windows that are candidates for being affected by
995 natural selection. (a) Distribution of genetic differentiation (F_{ST}) between *P. tremula*
996 and *P. tremuloides* from the observed (blue bar) and simulated data sets (black line).
997 The dashed lines indicate the thresholds for determining significantly (False
998 Discovery Rate <1%) high (orange bars) and low (green bars) inter-specific
999 differentiation based on coalescent simulations. (b) Comparisons of d_{xy} , relative node
1000 depth (RND) and the proportion of inter-specific shared polymorphisms among
1001 regions displaying significantly high (orange boxes) and low (green boxes)
1002 differentiation versus the genomic background (blue boxes). (c) Comparisons of
1003 multiple population genetic statistics, nucleotide diversity (θ_π), Tajima's D, Fay
1004 & Wu's H, linkage disequilibrium (r^2), recombination rate (ρ/θ_π), the proportion of
1005 fixed differences caused by derived alleles fixed in either *P. tremula* or *P. tremuloides*,
1006 among regions displaying significantly high (orange boxes) and low differentiation
1007 (green boxes) versus the genomic background (blue boxes). Asterisks designate
1008 significant differences between the outlier and the rest of genomic regions by Mann-
1009 Whitney U test (* P -value < 0.05; ** P -value < 1e-4; *** P -value < 2.2e-16).

1010

1011

1012

1013

1014

1015

1016

1017

1018

1019

1020

1021

1022

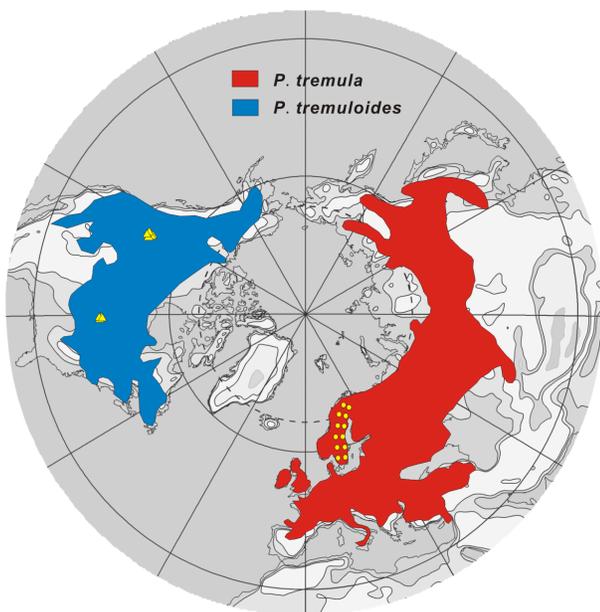
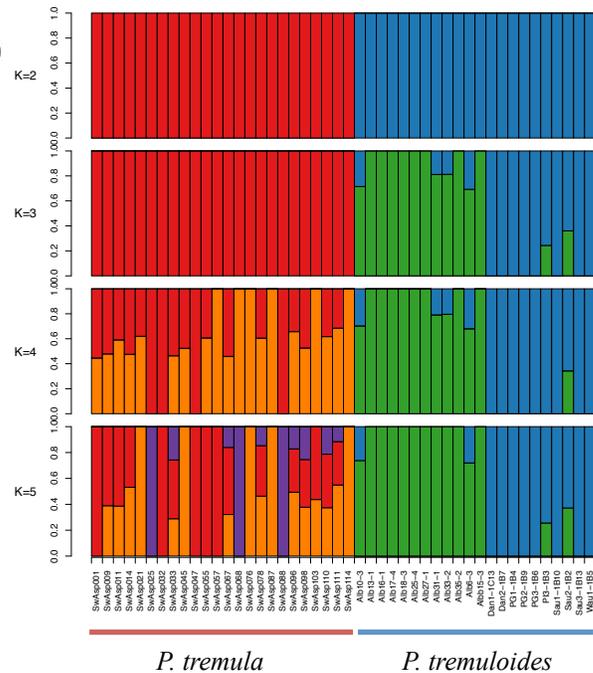
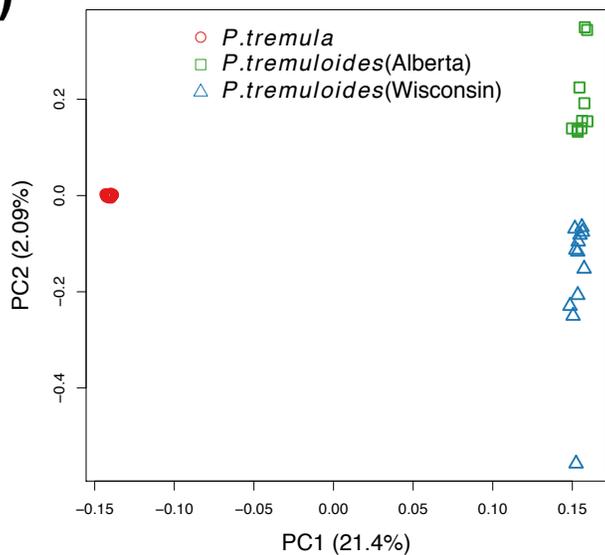
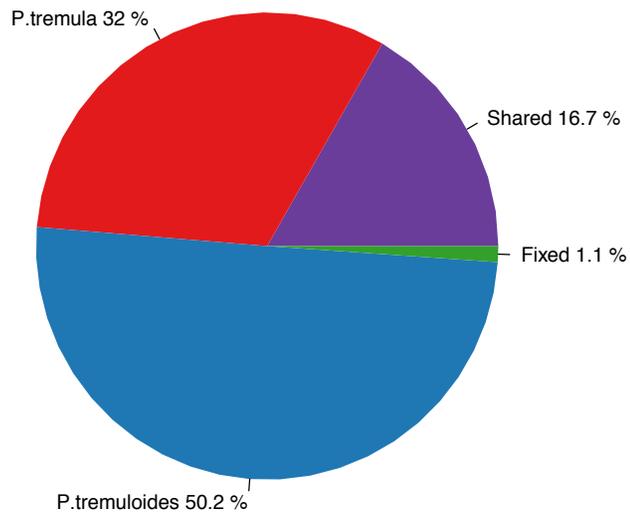
1023

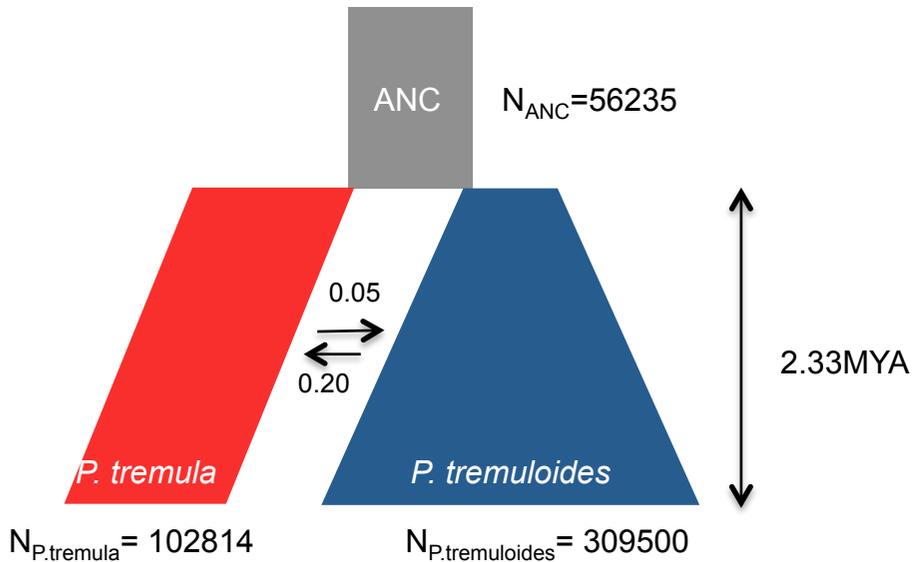
1024 **Table 1.** Inferred demographic parameters of the best-fitting demographic model
1025 shown in Figure 2a.
1026

Parameters	Point estimation	95% CI^a	
		Lower bound	Upper bound
N_{ANC}	56235	48012	69492
$N_{P.tremula}$	102814	93688	105671
$N_{P.tremuloides}$	309500	247321	310105
$2Nm_{P.tremuloides \rightarrow P.tremula}$	0.202	0.156	0.375
$2Nm_{P.tremula \rightarrow P.tremuloides}$	0.053	0.052	0.117
T_{DIV}	2332410	2186760	3113520

1027
1028 Parameters are defined in Figure 2a. N indicates the effective population size of *P.*
1029 *tremula*, *P. tremuloides* or their ancestral population, m indicates the migration rates
1030 between species on either direction, T_{DIV} indicates the estimated divergence time
1031 between the two species obtained from *fastsimcoal2*.
1032 ^aParametric bootstrap estimates obtained by parameter estimation from 100 data sets
1033 simulated according to the overall maximum composite likelihood estimates shown in
1034 point estimation columns. Estimations were obtained from 100,000 simulations per
1035 likelihood.

1036
1037

(a)**(b)****(c)****(d)**

(a)**(b)**