1      **Comparative population genomics of three related *Populus* species**

2

3      Jing Wang[1], Nathaniel R. Street[2], Douglas G. Scofield[1,3,4], Pär K. Ingvarsson[1*]

4

5      [1] Department of Ecology and Environmental Science, Umeå University, Umeå,

6      Sweden

7      [2] Umeå Plant Science Centre, Department of Plant Physiology, Umeå University,

8      Umeå, Sweden

9      [3] Department of Ecology and Genetics: Evolutionary Biology, Uppsala University,

10     Uppsala, Sweden

11     [4] Uppsala Multidisciplinary Center for Advanced Computational Science, Uppsala

12     University, Uppsala, Sweden

13

14     * Corresponding author

15     Email: par.ingvarsson@umu.se

16

17

18

19

20

21

22

23

24

## Abstract

A central aim of evolutionary genomics is to identify the relative roles that various evolutionary forces have played in generating and shaping genetic variation within and among species. Here we used whole-genome re-sequencing data from three related *Populus* species to characterize and compare genome-wide patterns of nucleotide polymorphism, site frequency spectrum, population-scaled recombination rate and linkage disequilibrium. Our analyses revealed that *P. tremuloides* has the highest level of genome-wide variation, skewed allele frequencies and population-scaled recombination rates, whereas *P. trichocarpa* harbors the lowest. Consistent with this, linkage disequilibrium decay was fastest in *P. tremuloides* and slowest in *P. trichocarpa*. Pervasive natural selection has been proven to be the primary force creating significant positive correlations between neutral polymorphism and recombination rate in all three species. Disparate effective population sizes and recombination rates among species, on the other hand, drive the distinct magnitudes and signatures of linked selection and consequent heterogeneous patterns of genomic variation among them. We find that purifying selection against slightly deleterious non-synonymous mutations is more effective in regions experiencing high recombination, which may provide one explanation for a partially positive association between recombination rate and gene density in these species. Moreover, distinct signatures of linked selection dependent on gene density are found between genic and intergenic regions within each species. To our knowledge, the present work is the first comparative population genomic study among forest tree species and represents an important step toward dissecting how the interactions of various evolutionary forces have shaped genomic variation within and among these ecologically and economically important tree species.

50   **Author Summary**

51   A fundamental goal of population genetics is to understand how various evolutionary

52   forces shape the heterogeneity of genomic variation within and among species. Here,

53   we characterize and compare genome-wide patterns of nucleotide diversity, site

54   frequency spectra, population-scaled recombination rates and linkage disequilibrium

55   among three related *Populus* species: *Populus tremula*, *P. tremuloides* and *P.*

56   *trichocarpa*. Pervasive natural selection, mediated by the local recombination

57   environments, is supposed to be the primary force shaping heterogeneous patterns of

58   neutral polymorphism throughout the genome. The disparate magnitudes and

59   signatures of linked selection among the three species, however, likely result from

60   either different effective population sizes and/or differences in recombination rates

61   among them. Moreover, we find distinct patterns of selection between genic and

62   intergenic regions in all three species, indicating these two types of sites may have

63   undergone independent evolutionary responses to selection in *Populus*. To our

64   knowledge, the present work provides the first phylogenetic comparative study of

65   genome-wide patterns of variation between closely related forest tree species. This

66   information will also improve our ability to understand how various evolutionary

67   forces have interacted to influence genome evolution among related species.

68

69

70

71

72

73

74

## Introduction

A major goal in evolutionary genetics is to understand how genomic variation is established, maintained and diverge within and between species [1, 2]. Various evolutionary forces are known to have substantial impacts in shaping genetic variation and linkage disequilibrium throughout the genome [3]. Under the neutral theory, genetic variation is the manifestation of the balance between mutation and genetic drift [4]. Demographic fluctuations, such as population expansion and/or bottlenecks, can cause patterns of genome-wide variation to deviate from standard neutral model in various ways [5]. Natural selection, via positive selection favoring beneficial mutations (genetic hitchhiking) and/or negative selection against deleterious mutations (background selection), plays an important role in sculpting the landscape of polymorphism across the genome [2, 6-8]. The signature and magnitude of apparent selection at linked sites depends heavily on the local environment of recombination [9, 10]. Linked selection is expected to remove more neutral polymorphism in low-recombination regions compared to high-recombination regions [10-12]. In addition to indirectly affecting genetic variation via linked selection, the rate of recombination can also shape the landscape of genomic features, such as base composition and gene density [13, 14]. However, there remains much to be learned about how these various evolutionary forces have shaped the heterogeneous patterns of genomic polymorphism within and between species [2, 6, 15]. With the advance of next-generation sequencing technology, sufficient genome-wide data among multiple related species are becoming available [16, 17]. Phylogenetic comparative approaches using these data will place us in a stronger position to understand the relative importance of mutation, genetic drift, natural selection and recombination in determining patterns of genome evolution [18, 19].

100        Thus far, genome-wide comparative studies have largely dealt with

101        experimental model species, mammals, and cultivated plants of either agricultural or

102        horticultural interest [19-21]. Forest trees, as a group, are characterized by extensive

103        geographical distributions and are of high ecological and economic value [22]. Most

104        forest trees have largely persisted in an undomesticated state and, until quite recently,

105        without anthropogenic influence [22]. Accordingly, in contrast to crop and livestock

106        lineages that have been through strong domestication bottlenecks, most extant

107        populations of forest trees harbor a wealth of genetic variation and they may be

108        excellent model systems for dissecting the dominant evolutionary forces that sculpt

109        patterns of variation throughout the genome [22, 23]. Among forest tree species, the

110        genus *Populus* represents a particularly attractive choice because of its wide

111        geographic distribution, important ecological role in a wide variety of habitats,

112        multiple economic uses in wood and energy products, and relatively small genome

113        size [24, 25]. Here, we studied three *Populus* species which differ in their morphology,

114        geographic distribution, population size and phylogenetic relationship (S1 Fig) [26,

115        27]. *P. tremula* and *P. tremuloides* (collectively 'aspens') have wide native ranges

116        across Eurasia and North America respectively, and are closely related, belonging to

117        the same section of the genus *Populus* (section *Populus*) [27]. In contrast, *P.*

118        *trichocarpa* belongs to a different section of the genus (section *Tacamahaca*) that is

119        reproductively isolated from members of the section *Populus* [27]. The distribution of

120        *P. trichocarpa* is restricted to western North America and it's range is considerably

121        smaller than the two aspen species [28]. Importantly, *P. trichocarpa* also represents

122        the first tree species to have its genome sequenced [29] and the genome sequence and

123        annotation have undergone continual improvement [http://phytozome.jgi.doe.gov].

124        This enables us to provide important context for our genome comparisons. The

125 phylogenetic relationship of the three species ((*P. tremula–P.tremuloides*) *P.*

126 *trichocarpa*) is well established by both chloroplast and nuclear DNA sequences [26,

127 30].

128     In this study, we used novel and existing Illumina short read ($2 \times 100$ bp)

129 datasets to characterize, compare and contrast genome-wide patterns of nucleotide

130 diversity, recombination rate, and linkage disequilibrium, and to infer contextual

131 patterns of selection throughout the genomes of all three species.

132

## Results

134 We generated whole-genome sequencing data of 24 genomes of *P. tremula* and 22

135 genomes of *P. tremuloides* (S1 Table) with all samples sequenced to relatively high

136 depth (24.2×-69.2×; S2 Table). We also downloaded 24 genomes from already

137 published data of *P. trichocarpa* [31]. After adapter removal and quality trimming,

138 949.2 Gb of high quality sequence data remained (S2 Table, S2 Fig). Reads from the

139 three species were mapped to the *P. trichocarpa* reference genome [29] using BWA-

140 MEM [32], with the mean mapping rate being 89.8% for individuals of *P. tremula*,

141 91.1% for individuals of *P. tremuloides*, and 95.2% for individuals of *P. trichcoarpa*

142 (S2 Table). On average, the genome-wide coverage of uniquely mapped reads was

143 more than 20× for each species (S2 Table). After excluding sites with extreme

144 coverage, low mapping quality, or those overlapping with annotated repetitive

145 elements separately in each species (see Materials and Methods), 42.8% of collinear

146 genomic sequences remained for downstream analyses. Among all retained genomic

147 regions, 54.9% were located within gene boundaries, which covers 70.1% of all genic

148 regions predicted from *P. trichocarpa* assembly, and the remainder (45.1%) was

149 located in intergenic regions.

6

150

151 **Negligible population substructure in samples of all three species**

152 Given the great dispersal capabilities of pollen and seed in *Populus* [24, 25],

153 population genetic structure appears to be generally weak in most *Populus* species

154 [31, 33]. In order to ascertain the population structure within and between species, we

155 used a model-based clustering algorithm, implemented in ADMIXTURE [34], to

156 cluster sampled individuals using only 4-fold synonymous single nucleotide

157 polymorphisms (SNPs) with minor allele frequency greater than 10%. When we

158 analyzed population structure between species, we found the model exhibits the

159 lowest cross-validation error when $K$=3 (S3b Fig), which clearly subdivides the three

160 species into three distinct clusters (S3a Fig). When we analyzed local population

161 structure within each species, $K$=1 minimized the cross-validation error in all species,

162 implying that extremely weak population structure in the samples of the three *Populus*

163 species (S3c-e Fig). Therefore, intra-species population structure likely play a

164 negligible role in our comparative population genetic analyses among the three

165 species.

166

167 **Patterns of divergence among the three *Populus* species**

168 We measured pairwise nucleotide divergence ($d_{xy}$) among pairs of the three species

169 across the genome in non-overlapping 100 kilobase pairs (Kbp) windows. Between

170 either of the two aspen species and *P. trichocarpa*, $d_{xy}$ was significantly higher than

171 between the two aspen species (Wilcoxon rank sum test, *P*-value<0.001) (S4 Fig). We

172 found extremely consistent patterns of divergence between the two aspen species and

173 *P. trichocarpa* (Spearman's $\rho$ = 0.994, *P*-value<0.001) (S5a Fig), reflecting the

174 historical divergence of the common ancestor of two aspen species from *P.*

175   *trichocarpa* and the relatively recent divergence of the two aspen species. In addition,

176   we found that the divergence was significantly correlated between the evolutionarily

177   independent lineages (*P. tremula-P. tremuloides)* vs. (aspens-*P. trichocarpa*) (S5b,c

178   Fig), suggesting that patterns of genome-wide variation of mutation rates and/or

179   selective constraints are relatively conserved across these three *Populus* species.

180

181   **Polymorphism varies, but is highly correlated, between species**

182   Fig 1 shows genome-wide estimates of nucleotide diversity among all three species

183   over non-overlapping 100 Kbp windows. We also performed the analyses using 1

184   megabase pair (Mbp) windows, with the results being nearly identical (S6a Fig). We

185   found that both aspen species harbor substantial levels of nucleotide diversity

186   ($\Theta_\Pi$=0.0133 in *P. tremula*; $\Theta_\Pi$=0.0144 in *P. tremuloides*), approximately two-fold

187   higher than the diversity in *P. trichocarpa* ($\Theta_\Pi$=0.0059) (Table 1; Fig 1; S6a Fig). The

188   overall nucleotide diversity we observe in *P. trichocarpa* was a slightly higher than

189   the value reported in [31]. This is likely due to differences in the methods used

190   between the two studies. In this study, we utilized the full information of the filtered

191   data and estimated the population genetic statistics directly from genotype

192   likelihoods, which takes statistical uncertainty of SNP and genotype calling into

193   account and should give more accurate estimates [35, 36]. In accordance with the

194   highly consistent genome-wide distribution of $\Theta_\Pi$ among the three *Populus* species

195   (Fig 1), we observed a significantly positive correlation of $\Theta_\Pi$ between each pair of

196   species across the whole genome (Fig 2a). Such strong correlations of polymorphism

197   suggest that mutation rates and/or selective constraints are highly conserved among

198   the species despite the clades represented by genetic sections *Populus* and

199   *Tacamahaca* having diverged for ~4.5 million years [37]. Not surprisingly, we found

8

200　higher correlation of $\Theta_\Pi$ between *P. tremula* and *P. tremuloides* (Spearman's

201　$\rho$=0.829, *P*-value<0.001, Fig 2a), which both belong to section *Populus*, compared to

202　the correlation between the two aspen species and *P. trichocarpa*, which is likely due

203　to the higher levels of shared ancestral polymorphism between the aspens [26].

204

205　**Table 1. Diversity statistics (median and central 95% range) for various genomic**

206　**contexts over 100 Kbp non-overlapping windows across genome**

| | Filtered bases (Mbp) | P. tremula | | P. tremuloides | | P. trichocarpa | |
|---|---|---|---|---|---|---|---|
| | | $\Theta_\Pi$ | Tajima's D | $\Theta_\Pi$ | Tajima's D | $\Theta_\Pi$ | Tajima's D |
| **Total** | 136.47 | 0.0133(0.0076-0.0236) | -0.2723(-0.7727-0.2941) | 0.0144(0.0091-0.0247) | -1.1688(-1.6003--0.4899) | 0.0059(0.0031-0.0125) | 0.0643(-0.8266-0.9496) |
| **0-fold[a]** | 16.52 | 0.0035[***](0.0011-0.0085) | -1.0913[***](-2.2240-0.0128) | 0.0044[***](0.0018-0.0091) | -2.1717[***](-2.7932--1.2275) | 0.0013[***](0.0003-0.0043) | -0.4090[***](-1.7625-1.4391) |
| **4-fold** | 3.40 | 0.0108(0.0035-0.0207) | -0.2220(-1.4676-0.9667) | 0.0120(0.0044-0.0214) | -1.3689(-2.2458--0.2602) | 0.0040(0.0008-0.0104) | 0.0084(-1.5663-1.7753) |
| **Introns[a]** | 31.89 | 0.0096[***](0.0038-0.0182) | -0.2669[**](-1.3490-0.7526) | 0.0106[***](0.0046-0.0184) | -1.4286[**](-2.2283--0.4173 ) | 0.0038[**](0.0013-0.0094) | -0.0245[*](-1.5429-1.6489) |
| **UTR 5'[a]** | 4.02 | 0.0091[***](0.0033-0.0192) | -0.5642[***](-1.7370-0.6994) | 0.0104[***](0.0041-0.0197) | -1.5829[***](-2.3688--0.4202) | 0.0038[**](0.0007-0.0102) | -0.1040[**](-1.6203-1.6915) |
| **UTR 3'[a]** | 7.19 | 0.0108(0.0039-0.0190) | -0.3081[**](-1.4577-0.7662) | 0.0121(0.0050-0.0204) | -1.3842[*](-2.2319--0.3766) | 0.0043(0.0012-0.0101) | -0.0033(-1.5265-1.6444) |
| **Intergenic[a]** | 73.46 | 0.0184(0.0110-0.0313) | -0.3062[**](-0.8360-0.4298) | 0.0198(0.0130-0.0326) | -1.1843(-1.6698--0.3969) | 0.0088(0.0044-0.0175) | 0.1042(-0.9211-1.1131) |

207　[a] One-sided paired Mann-Whitney U test in comparison to the 4-fold synonymous

208　* *P*<0.05

209　** *P*<0.001

210　*** *P*<$2.2\times10^{-16}$

211

212　**Fig 1. Genome-wide patterns of polymorphism among three *Populus* species.**

213　Nucleotide diversity ($\Theta_\pi$) was calculated over 100 Kbp non-overlapping windows in

9

214    *P. tremula* (orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line)

215    along the 19 chromosomes.

216

217    **Fig 2. Correlations of polymorphism, Tajima's D and population recombination**

218    **rate between species.** Distributions and correlations of (a) pairwise nucleotide

219    diversity ($\Theta_\pi$), (b) Tajima's D, (c) population-scaled recombination rate ($\rho$) between

220    pairwise comparisons of *P. tremula*, *P. tremuloides* and *P. trichocarpa*, over 100 Kbp

221    non-overlapping windows. The red to yellow to blue gradient indicates decreased

222    density of observed events at a given location in the graph. Spearman's rank

223    correlation coefficient ($\rho$) and the *P*-value are shown in each subplot. (*** *P*<2.2×10$^-$

224    $^{16}$, ***P*<0.001). The dotted grey line in each subplot indicates simple linear regression

225    line with intercept being zero and slope being one.

226

227            Along all chromosomes, the distribution of polymorphisms was more variable

228    (average coefficient of variation (CV) of $\Theta_\Pi$ among three species=0.3362) than was

229    divergence (average CV of $d_{xy}$ =0.1670) (Fig 1; S4 Fig). As the *Populus* karyotype

230    has not been established and thus the locations of centromeres and telomeres remains

231    unknown, we can only speculate that the genomic regions with long measurement

232    gaps that failed to pass our quality requirements may represent repetitive regions of

233    chromosomes located near centromeres, and with distal chromosomal regions being

234    the approximate locations of telomeres. Fig 1 shows that diversity generally declines

235    near the supposed locations of centromeres and telomeres in all three *Populus* species.

236    Divergence, however, did not show similar patterns of decline in such regions (S4

237    Fig), potentially indicating that the reduced polymorphism in these regions is most

238    likely due to a greater influence of selection at linked sites because of reduced

239   recombination in these regions (see below) rather than reduced neutral mutation rates

240   [2].

241

242   **Tajima's D varies, and is weakly correlated, between species**

243   Genome-wide allele frequency distributions can also help elucidate the relative

244   contributions of different evolutionary dynamics in charactering patterns of

245   polymorphism. We compared the site frequency spectrum among the three species

246   based on the Tajima's D statistic [38], which is the standardized difference between

247   the average pairwise sequence diversity ($\Theta_\Pi$) and the number of segregating sites

248   ($\Theta_W$). Under the standard neutral model, the expected value of Tajima's D is roughly

249   equal to 0 ($\Theta_\Pi = \Theta_W$). Negative Tajima's D ($\Theta_\Pi < \Theta_W$; an excess of rare alleles)

250   usually results from purifying selection, selective sweeps, or population expansion,

251   whereas positive Tajima's D ($\Theta_\Pi > \Theta_W$; an excess of common alleles) indicates either

252   balancing selection or a decrease in population size. We found dramatically different

253   patterns in the genome-wide distribution of Tajima's *D* among the three species (Fig

254   3). The genome-wide average of Tajima's D was slightly positive in *P. trichocarpa*,

255   whereas *P. tremula* had negative genome-wide averages of Tajima's D (Table1; Fig

256   3; S6c Fig). Compared to *P. trichocarpa* (average Tajima's D=0.064) and *P. tremula*

257   (average Tajima's D=-0.272), *P. tremuloides* (average Tajima's D=-1.169) showed

258   substantially more negative values of Tajima's D along all chromosomes (Fig 3; S6c

259   Fig), reflecting a large excess of low-frequency polymorphisms across the genome.

260   As natural selection is usually expected to act on a relatively small number of

261   genomic regions, the marked genome-wide negative Tajima's D is most likely to be

262   explained by a recent substantial expansion in population size in *P. tremuloides*. The

263   weakly negative Tajima's D in *P. tremula* could also reflect an increase in population

11

264    size although not as great as that experienced by *P. tremuloides*. The slightly positive

265    Tajima's D in *P. trichocarpa*, however, implies that it may have experienced a recent

266    population contraction as was also suggested by [39]. In contrast to the significantly

267    positive correlations of nucleotide diversity among the three species, the much weaker

268    correlations seen for Tajima's D (Fig 2b) could be ascribed to either the different

269    demographic histories of these species, or result from different targets of divergent

270    selection due to different environmental conditions experienced by the species since

271    their divergence [37].

272

273    **Fig 3. Genome-wide patterns of allele frequency distribution among three**

274    ***Populus* species.** Tajima's D was calculated over 100 Kbp non-overlapping windows

275    in *P. tremula* (orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line)

276    along the 19 chromosomes.

277

278    **Little confounding effects of population structure, biased sampling schemes and**

279    **hybridization**

280    The observed variation of intra-species population genetic patterns could also be

281    caused by other factors, such as population sub-structure, biased sampling schemes

282    and/or hybridization [40, 41]. We found no obvious population sub-structure in

283    samples among all three species (S3 Fig) and thus expect that any effect of population

284    structure is negligible in structuring polymorphisms in these three species. Biased

285    sampling schemes could lead to biased estimates of genome-wide diversity and allele

286    frequency spectrum as the samples of *P. tremula* and *P. trichocarpa* were all collected

287    from continuous local populations whereas those of *P. tremuloides* were collected

288    from two discrete populations (S1 Fig).  The lack of population substructure suggests

12

289    that bias is unlikely, but to be sure we tested this by calculating $\Theta_\Pi$ and Tajima's D

290    separately for the two local *P. tremuloides* populations (Alberta and Wisconsin). We

291    observed remarkably similar values of $\Theta_\Pi$ in both local population samples and the

292    pooled samples (S7a Fig), confirming that structured sampling in this species does not

293    affect our results. Values of Tajima's D were slightly skewed toward more negative

294    for the pooled samples compared with single sampling localities (S7b Fig), likely

295    reflecting low sharing of low-frequency polymorphisms between these localities

296    which is consistent with widespread population expansion. This seems unlikely to

297    influence the comparison of the overall patterns of genetic variation among the

298    species. An additional possibility is that the excess of rare alleles we observed in *P.*

299    *tremuloides* could be derived from one or few "outlier" individuals that are

300    misidentified or are recent inter-specific hybrids. To assess this possibility we

301    calculated the number of singletons contributed by each individual in the dataset. We

302    found an overall higher number of singletons in individuals of *P. tremuloides* relative

303    to the other two species, which was expected from patterns of Tajima's D, but there

304    were no outlier individuals in *P. tremuloides* that contribute disproportionally large

305    numbers of singletons (S8 Fig). Together these results indicate that the genome-wide

306    excess of rare variants we observed in *P. tremuloides* is a species-wide pattern rather

307    than being population or individual specific.

308

309    **Patterns of polymorphism and divergence vary by genomic contexts**

310    We compared patterns of nucleotide diversity and divergence across different

311    genomic contexts and in all comparisons levels of nucleotide diversity and divergence

312    were highest for intergenic sites, followed by 4-fold synonymous sites, 3'UTRs,

313    5'UTRs, introns and were lowest at 0-fold non-synonymous sites (Table 1; S3 Table;

13

314    S9 Fig; S11 Fig). The extremely high levels of diversity and divergence in intergenic

315    regions could arise due artifacts of mapping errors in repetitive sequences [42].

316    However, we applied the same strict filtering steps in both genic and intergenic

317    regions, making this error bias less likely. Therefore, the markedly higher levels of

318    diversity and divergence in intergenic regions probably result from a higher mutation

319    rate, a relaxed selective constraint or both [2]. If we assume that the mutation rate of

320    intergenic regions does not differ from that in genic regions, we could infer that there

321    is strong selective constraint on all genic features throughout the genomes.

322    Nevertheless, the relative contribution of alternative factors to the higher divergence

323    rate in intergenic regions requires further investigation.

324         Within genic regions, 3' UTRs showed only slightly lower levels of

325    divergence and similar levels of diversity and allele frequency distribution compared

326    to 4-fold synonymous sites (Table 1; S3 Table). This suggests that the large majority

327    of sites in 3' UTRs are effectively neutral or are subject to purifying selection to an

328    extent comparable to 4-fold synonymous sites. We found a slight, but significant,

329    reduction in diversity, Tajima's D and divergence in introns and 5' UTRs, consistent

330    with the notion that introns and 5'UTRs have undergone stronger selective constraint

331    than 4-fold synonymous sites (Table 1; S3 Table). Finally, both diversity and

332    divergence at 0-fold non-synonymous sites was nearly three times lower than 4-fold

333    synonymous sites. In accordance with this, we found significantly lower Tajima's D

334    at 0-fold non-synonymous sites compared to 4-fold synonymous sites ($P<0.001$,

335    Mann-Whitney U test) (Table 1), indicating that a large majority of amino acid

336    substitutions are under strong purifying selection [43].

337

338    **Linkage Disequilibrium (LD) and Recombination**

14

339    *Populus* species are predominantly outcrossing and thus the expectation is thus that

340    LD decays rapidly and that the rates of scaled recombination are high [44]. However,

341    a recent genome-wide analysis in *P. trichocarpa* has revealed more extensive LD

342    across the genome that was expected base on earlier studies [45]. We found that the

343    average LD ($r^2$) between pairs of SNPs fell to lower than 0.2 within approximately 6-

344    7 Kbp in *P. trichocarpa* (Fig 4), which is consistent with values previously reported in

345    this species [45]. In *P. tremula*, mean $r^2$ dropped below 0.2 within about 5 Kbp, which

346    is substantially greater than reported in earlier studies that were based on a small

347    number of candidate gene fragments [44]. Finally, LD decayed considerably more

348    rapidly in *P. tremuloides* compared to the other two species, with mean $r^2$ dropping

349    below 0.2 within ~2-3 Kbp (Fig 4).

350

351    **Fig 4. Decay of linkage disequilibrium (LD).** The decay of LD (estimated as $r^2$)

352    with physical distance in *P. tremula* (orange line), *P. tremuloides* (blue line) and *P.*

353    *trichocarpa* (green line).

354

355         We also estimated population-scaled recombination rates ($\rho$) in each species.

356    There was considerable large-scale variation in recombination rates throughout the

357    genomes of all three species, with $\rho$ in *P. tremuloides* consistently being higher than

358    in the other two species (Fig 5). In accordance with the genome-wide patterns of

359    diversity, we also found patterns of decreasing $\rho$ near the putative locations of

360    centromeres and telomeres in all three species (Fig 5). When we measured the average

361    $r^2$ over 100 Kbp non-overlapping windows across the genome, we found population

362    recombination rates were significantly correlated with the extent of LD (mean

363    pairwise $r^2$) in all species (S12 Fig). The mean $\rho$ computed from 100 Kbp windows in

15

364    *P. tremuloides* was 8.42 Kbp$^{-1}$ (standard deviation of 4.71 Kbp$^{-1}$), and the mean $\rho$ in

365    *P. tremula* was 3.23 Kbp$^{-1}$ (standard deviation: 1.66 Kbp$^{-1}$). The genome-wide

366    average $\rho$ in *P. trichocarpa* was 2.19 Kbp$^{-1}$ (standard deviation: 1.11 Kbp$^{-1}$), which is

367    consistent with the previously reported $\rho$ value estimated from exome re-sequencing

368    data [39]. Concordant $\rho$ values for all three species were also observed in 1Mbp

369    windows (S6d Fig). In comparison to the extremely high correlation of diversity and

370    low correlation of allele frequency spectrum among the three *Populus* species (Fig

371    2a,b), we found an intermediate correlation in recombination rates between species,

372    suggesting that the overall recombination environment is only partially conserved

373    among the three species (Fig 2c).

374

375    **Fig 5. Genome-wide patterns of population-scaled recombination rate among**

376    **three *Populus* species.** Population-scaled recombination rate ($\rho$) was averaged over

377    100 Kbp non-overlapping windows in *P. tremula* (orange line), *P. tremuloides* (blue

378    line) and *P. trichocarpa* (green line) along the 19 chromosomes.

379

380        For populations under drift-mutation-recombination equilibrium, $\rho = 4N_ec$

381    (where $N_e$ is the effective population size and $c$ is the recombination rate) and $\theta_W =$

382    $4N_e\mu$ (where $N_e$ is the effective population size and $\mu$ is the mutation rate). In order to

383    compare the relative contribution of recombination ($c$) and mutation ($\mu$) in shaping

384    genomic variation, we measured the ratio of population recombination rate to the

385    nucleotide diversity ($\rho/\theta_W$) across the genome (S13 Fig). The mean $c/\mu$ in *P.*

386    *tremuloides* and *P. trichocarpa* was 0.39 and 0.38 respectively, indicating that

387    mutations occur approximately two to three times more frequently than recombination

388    events. On the other hand, the average value of $c/\mu$ in *P. tremula* was 0.22, implying

16

389    that recombination is less important than mutation in generating diversity in *P.*

390    *tremula* compared to the other two *Populus* species.

391

392    **Neutral polymorphism, not divergence, is positively correlated with**

393    **recombination rate**

394    If natural selection is pervasive across the genome, positive correlations between

395    levels of neutral polymorphisms and recombination rates are expected since

396    demography alone is unlikely to generate these patterns [8]. If selection is the primary

397    force driving the association of neutral polymorphism and recombination rate, the

398    association should be stronger in genic regions of the genome than in intergenic

399    regions since genes are more likely to be targets of selection. In order to examine

400    these correlations, we first assumed that 4-fold synonymous sites in genic regions

401    represent selectively neutral sites, as every possible mutation in 4-fold degenerate

402    sites is synonymous. In the following we refer to the pairwise nucleotide diversity at

403    4-fold synonymous sites ($\theta_{4\text{-fold}}$) as "neutral polymorphism". We then measured the 4-

404    fold synonymous substitution rate ($d_{4\text{-fold}}$) between either of the two aspen species and

405    *P. trichocarpa* and used this to represent "neutral divergence", which was further

406    taken as a proxy for the neutral mutation rate [46]. As many other genomic features

407    may also influence the variation of neutral polymorphism, we also tabulated GC

408    content, gene density and the number of neutral bases covered by sequencing data for

409    all three species. All measurements were carried out in non-overlapping windows that

410    were either 100 Kbp or 1Mbp in size.

411         We found significantly positive correlations between the level of neutral

412    polymorphism ($\theta_{4\text{-fold}}$) and population recombination rate for the two aspen species

413    (Table 2), with correlations being stronger in *P. tremula* compared to *P. tremuloides*.

17

414 In *P. trichocarpa*, however, we found either no or weak correlation between diversity

415 and recombination (Table 2). Compared to 100 Kbp windows, the correlations were

416 stronger in 1Mbp windows among all species, which most likely results from the

417 higher signal-to-noise ratio provided by larger genomic regions (Table 2). In the

418 remainder of this paper we therefore focus our analyses primarily on data generated

419 using 1Mbp window size. We performed simple linear regression analysis between

420 recombination rate and diversity, and the recombination rate explained 45.8%, 21.3%,

421 and 3.9% of the amount of neutral genetic variation in *P. tremula*, *P. tremuloides* and

422 *P. trichocarpa*, respectively (Fig 6).

423

424 **Table 2.** Summary of the correlation coefficients (Spearman's $\rho$) between levels of

425 neutral polymorphism, divergence and recombination rate in all three *Populus*

426 species.

| Dataset | Species | $\rho$ vs. $\theta_{4\text{-fold}}$ | | $\rho$ vs. $d_{4\text{-fold}}$ | $\rho$ vs. $\theta_{\text{Intergenic}}$ | | $\rho$ vs. $d_{\text{Intergenic}}$ |
|---------|---------|---------|---------|---------|---------|---------|---------|
| | | Pairwise | Partial[a] | | Pairwise | Partial[b] | |
| 100Kbp | *P. tremula* | 0.339[***] | 0.309[***] | 0.043 | 0.062[**] | 0.142[***] | -0.077[**] |
| | *P. tremuloides* | 0.310[***] | 0.284[***] | 0.061[**] | -0.037 | 0.100[**] | -0.029 |
| | *P. trichocarpa* | 0.011 | -0.024 | 0.053[*] | -0.080[**] | -0.002 | -0.015 |
| 1Mbp | *P. tremula* | 0.647[***] | 0.573[***] | -0.070 | 0.201[**] | 0.348[**] | -0.209[**] |
| | *P. tremuloides* | 0.400[**] | 0.363[**] | -0.033 | 0.032 | 0.320[**] | -0.127[*] |
| | *P. trichocarpa* | 0.227[**] | 0.151[*] | -0.027 | -0.072 | 0.165[*] | -0.120[*] |

427 [a]Partial correlation controls for GC content, gene density, divergence of 4-fold synonymous sites

428 between aspen and *P. trichocarpa*, and coverage (the number of 4-fold synonymous bases covered by

429 sequencing data).

430 [b]Partial correlation controls for GC content, gene density, divergence of intergenic sites between aspen

431 and *P. trichocarpa*, and coverage (the number of intergenic bases covered by sequencing data).

432 *  *P*<0.05

18

433    ** $P<0.001$

434    *** $P<2.2\times10^{-16}$

435

436    **Fig 6. Correlations between estimates of neutral genetic diversity and divergence**

437    **with population recombination rates over 1Mbp non-overlapping windows.**

438    Correlations between estimates of 4-fold synonymous diversity ($\Theta_{4\text{-fold}}$) (left panel)

439    and divergence ($d_{4\text{-fold}}$) (right panel) with population-scaled recombination rate ($\rho$)

440    over 1Mbp non-overlapping windows. Linear regression lines are colored according

441    to species: (a) *P. tremula* (orange line), (b) *P. tremuloides* (blue line) and (c) *P.*

442    *trichocarpa* (green line).

443

444         If the relationship between diversity and recombination rate was merely

445    caused by the mutagenic effect of recombination, similar correlations should also be

446    observed between divergence and recombination rate. However, no such correlations

447    were observed in any of the three species (Table 2; Fig 6). The association between

448    recombination rate and nucleotide diversity, and not with divergence, is thus most

449    likely caused by the effects of linked natural selection, where the elimination of linked

450    polymorphisms caused by selection is disproportionally stronger in low-

451    recombination genomic regions relative to regions of high recombination [8-10, 47].

452    Moreover, among all three species, the correlations between neutral polymorphism

453    and recombination rate remained significant even after we performed partial

454    correlation analyses to control for several possible confounding factors such as GC

455    content, gene density, divergence at neutral sites, and the number of neutral bases

456    covered by sequencing data (Table 2).

457         In accordance with the view that genes represent the most likely targets of

458    natural selection, the correlations between intergenic diversity and recombination rate

19

459    were substantially weaker than those correlations in genic regions (Table 2). Only 7.3%

460    of intergenic genetic variation in *P. tremula* could be explained by recombination,

461    whereas the impact of recombination rate on intergenic diversity in *P. tremuloides*

462    and *P. trichocarpa* was <1% and could be considered negligible (Table 2; S14 Fig).

463    In addition, we found slightly negative correlation between the divergence and

464    recombination in intergenic regions (Table 2; S14 Fig). This pattern is likely to be

465    explained by Hill-Robertson interference where weakly deleterious intergenic

466    mutations would reach fixation due to ineffective purifying selection in regions of low

467    recombination [48]. Further investigation is required to support this assertion. Notably,

468    after controlling for GC content, gene density, divergence and the number of covered

469    intergenic bases using partial correlation analyses, the correlations between intergenic

470    diversity and recombination rate become significant in all species, but remained

471    relatively weak compared to the values for genic regions in the two aspen species

472    (Table 2).

473

474    **Inconsistent effect of gene density on patterns of polymorphism in genic vs.**

475    **intergenic regions**

476    Genome-wide signatures of linked selection are not only influenced by the local

477    environments of recombination rate, but also are sensitive to the density of

478    functionally important sites within specific genomic regions [14]. Genomic regions

479    with a high density of genes are therefore expected to have undergone stronger effects

480    of linked selection and should therefore exhibit lower levels of neutral polymorphism

481    [1, 14]. However, a positive or negative co-variation of gene density and

482    recombination rate would either act to obscure or strengthen the genome-wide

483    signatures of linked selection, respectively [7, 14, 49]. We measured gene density as

484    the number of protein-coding genes in each 1Mbp window, which was unsurprisingly

485    also found to be highly correlated with the proportion of coding bases in each window

486    (S15 Fig). In all three *Populus* species, we found significantly positive correlation

487    between population recombination rate and gene density (Fig 7a). However, rather

488    than being linear, the relationships between recombination rate and gene density were

489    found to be curvilinear in all three species, with a significant positive correlation

490    observed only in regions of lower gene density (gene number smaller than ~85 within

491    each 1Mbp window) (Table 3). In clear contrast, in high gene density regions (gene

492    number greater than ~85 within each 1Mbp window) we observed no correlations

493    between recombination rate and gene density in any of the two aspen species, and

494    only weak correlation in *P. trichocarpa* (Table 3; Fig 7a). These correlation patterns

495    persisted after controlling for the GC content and the number of bases covered by

496    sequencing data in each window (Table 3).

497

498    **Fig 7. Correlations between estimates of population recombination rates, genic**

499    **and intergenic genetic diversity with gene density over 1 Mbp non-overlapping**

500    **windows.** (a) Correlations between gene density and population-scaled recombination

501    rate ($\rho$) in *P. tremula* (left panel), *P. tremuloides* (middle panel) and *P. trichocarpa*

502    (right panel). (b) Correlations between gene density and neutral genetic diversity ($\Theta_{4\text{-}fold}$)

503    in *P. tremula*, *P. tremuloides* and *P. trichocarpa*. (c) Correlations between gene

504    density and intergenic genetic diversity ($\Theta_{\text{Intergenic}}$) in *P. tremula*, *P. tremuloides* and

505    *P. trichocarpa*. Grey points represent the statistics computed over 1Mbp non-

506    overlapping windows. Colored lines denote the lowess curves fit to the two variables

507    in each species.

508

509 **Table 3.** Summary of the correlation coefficients (Spearman's $\rho$) between gene

510 density and population recombination rate, neutral polymorphism in genic and

511 intergenic regions over 1 Mbp non-overlapping windows in three *Populus* species.

| Species | Correlation type | Gene density vs. $\rho$[a] | | Gene density vs. $\theta_{4\text{-fold}}$[b] | | Gene density vs. $\theta_{Intergenic}$[c] | |
|---|---|---|---|---|---|---|---|
| | | low | high | low | high | low | high |
| *P. tremula* | Pairwise | 0.674** | -0.112 | 0.601** | -0.180* | 0.431** | -0.605*** |
| | Partial | 0.516** | 0.263* | 0.191* | 0.110 | 0.263* | -0.438** |
| *P.tremuloides* | Pairwise | 0.527** | 0.006 | 0.576** | -0.077 | 0.419** | -0.600*** |
| | Partial | 0.315** | 0.048 | 0.407** | 0.280** | 0.363** | -0.444** |
| *P.trichocarpa* | Pairwise | 0.609** | 0.168* | 0.417** | -0.033 | 0.529** | -0.513*** |
| | Partial | 0.477** | 0.193* | 0.242* | 0.263** | 0.432** | -0.273** |

512 [a]Partial correlation controls for GC content and the number of bases covered by the data

513 [b]Partial correlation controls for GC content, population recombination rate, divergence of 4-fold

514 synonymous sites between aspen and *P. trichocarpa*, and coverage (the number of 4-fold synonymous

515 bases covered by sequencing data).

516 [c]Partial correlation controls for GC content, population recombination rate, divergence of intergenic

517 sites between aspen and *P. trichocarpa*, and coverage (the number of intergenic bases covered by

518 sequencing data).

519 * $P<0.05$

520 ** $P<0.001$

521 *** $P<2.2\times10^{-16}$

522

523     We then examined the correlation between neutral polymorphism and gene

524 density. Compared to the prediction of lower diversity in regions with higher

525 functional density [50], we found that the correlation pattern between genic diversity

526 and gene density was highly consistent with the pattern found in recombination rate,

527 where significantly positive correlations were found in regions of lower gene density

528 and either no correlation or weak negative correlations were found in regions of

529  higher gene density (Table 3; Fig 7b). After controlling for potential confounding

530  variables such as GC content, recombination rate, neutral divergence, and the number

531  of covered sites in each window, weaker but significant positive correlations between

532  neutral diversity and gene density remained in all three species in regions of low gene

533  density (Table 3). Positive associations between neutral diversity and gene density

534  were also found in high gene-density regions (Table 3).

535      Compared with genic regions, different correlation patterns between intergenic

536  diversity and gene density were found in all three species (Fig 7c). In accordance with

537  genic regions, we found significantly positive correlation between intergenic diversity

538  and gene density in regions of lower gene density. However, in regions of higher gene

539  density, strongly negative correlations between intergenic diversity and gene density

540  were observed in all three species (Table 3; Fig 7c). Due to the lack of correlation

541  between intergenic divergence and gene density (S16 Fig), our findings suggest that

542  the levels of intergenic polymorphism are also largely affected by natural selection,

543  with the intensity of selection increasing with an increase of gene density. These

544  correlations remained significant even after controlling for possible confounding

545  variables (Table 3).

546

547  **Lack of correlation between synonymous diversity and non-synonymous**

548  **divergence**

549  A distinctive signature of recurrent selective sweeps is the local reduction of linked

550  neutral polymorphism due to frequent adaptive substitutions [51]. Given amino acid

551  substitutions compose a substantial number of adaptive substitutions, negative

552  correlation between neutral polymorphism and non-synonymous divergence can be

553  particularly informative of the prevalence of selective sweeps [52]. However, in all

23

554    three species, we found either no or weak negative correlations between neutral

555    polymorphism ($\theta_{4\text{-fold}}$) and the rate of non-synonymous substitutions ($d_{0\text{-fold}}$) in both

556    100 Kbp and 1 Mbp windows (S4 Table). The correlational patterns did not change

557    after we controlled for GC content, recombination rate, gene density, neutral

558    divergence rate, and the number of 4-fold synonymous and 0-fold non-synonymous

559    sites covered by the data (S4 Table). This result contrasts with our previous study

560    reported from a small number of candidate genes, where we found a significant

561    negative correlation between polymorphism at synonymous sites and amino acid

562    divergence in *P. tremula* [53]. One possible explanation for the different patterns

563    between these two studies is that they are based on different scales of measurement,

564    from single genes to 100 Kbp and 1 Mbp windows [52]. Accordingly, additional

565    future analyses are still needed to examine the relationship between the synonymous

566    polymorphism and the rate of amino acid evolution on a genic scale.

567

568    **The effect of recombination on the efficacy of natural selection**

569    We next characterized the ratio of non-synonymous to synonymous polymorphism

570    ($\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$) and divergence ($d_{0\text{-fold}}/d_{4\text{-fold}}$) for each of the three *Populus* species in

571    order to assess whether there is a relationship between the efficacy of natural selection

572    and the rate of recombination (Table 4). Once GC content, gene density and number

573    of 4-fold synonymous and 0-fold non-synonymous sites were taken into account, we

574    found no correlation between recombination rate and $d_{0\text{-fold}}/d_{4\text{-fold}}$ in all three species

575    (Table 4). We did not observe any significant correlations between recombination rate

576    and $\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$ in the 1 Mbp windows after controlling for the various confounding

577    factors (Table 4). However, for 100 Kbp windows, we found significantly negative

578    correlations between recombination rate and $\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$ in *P. tremula* and *P.*

579 *tremuloides*, but not in *P. trichocarpa*. The relative overabundance of non-

580 synonymous polymorphism in regions of low recombination most likely suggests that

581 the effective elimination of weakly deleterious non-synonymous mutations was

582 reduced in low recombination regions in the two aspen species [12]. The lack of such

583 correlation in *P. trichocarpa* may reflect its lower effective population size and

584 accordingly weaker efficacy of selection across the genome [54]. In addition, since no

585 correlation between $\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$ and recombination rate was observed at a broad scale

586 (1 Mbp) in any of the three species, it is likely that interference between weakly

587 selected mutations is more easier to be detected at fine scales [54], although this

588 requires further investigation.

589

590 **Table 4.** Summary of the correlation coefficients (Spearman's $\rho$) between

591 recombination rate and the ratio of non-synonymous to synonymous polymorphism

592 ($\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$) and divergence ($d_{0\text{-fold}}/d_{4\text{-fold}}$).

| Dataset | Species | $\rho$ vs. $\theta_{0\text{-fold}}/\theta_{4\text{-fold}}$ | | $\rho$ vs. $d_{0\text{-fold}}/d_{4\text{-fold}}$ | |
| | | Pairwise | Partial[a] | Pairwise | Partial[a] |
|---|---|---|---|---|---|
| 100Kbp | *P. tremula* | -0.057[*] | -0.075[**] | -0.012 | -0.005 |
| | *P. tremuloides* | -0.118[**] | -0.122[**] | -0.003 | -0.002 |
| | *P. trichocarpa* | -0.004 | -0.002 | -0.026 | -0.020 |
| 1Mbp | *P. tremula* | -0.063 | -0.045 | -0.007 | 0.017 |
| | *P. tremuloides* | -0.142[*] | -0.092 | 0.014 | 0.020 |
| | *P. trichocarpa* | 0.035 | -0.002 | 0.030 | 0.036 |

593 [a]Partial correlation controls for GC content, gene density, and the number of 4-fold synonymous and 0-

594 fold non-synonymous bases covered by sequencing data.

595 *$P<0.05$

596 **$P<0.001$

597

## Discussion

599  We have characterized and compared genome-wide nucleotide polymorphism, site

600  frequency spectra, linkage disequilibrium (LD), and population-scaled recombination

601  rates among three related *Populus* species. Widespread variation in nucleotide

602  diversity is found throughout the genomes of all three species and we found

603  significant genome-wide correlations of diversity among the three *Populus* species.

604  This likely results from shared selective constraints and/or patterns of conserved

605  variation in mutation rate between these related species [55, 56]. Compared to *P.*

606  *trichocarpa*, levels of diversity in *P. tremula* and *P. tremuloides* are more than two-

607  fold higher throughout the genome. The higher diversity we find in both aspen species

608  is likely due to their larger effective population sizes ($N_e$) because consistent patterns

609  of interspecific sequence divergence between independent evolutionary lineages (S5

610  Fig) indicate that mutation rates are likely to be conserved among the three species

611  [4]. Larger effective population sizes in the two aspen species are also in agreement

612  with their larger current census population size and substantially more extensive

613  geographic ranges [24]. Assuming that mutation rates do not differ dramatically

614  among the three species, we could infer that the effective population size in the two

615  aspen species are more than twice as large as in *P. trichocarpa* [4]. However, the

616  relative importance of mutation rate variation in determining diversity levels across

617  related species obviously deservs to be studies further, particularly in light of very

618  recent results indicating that high levels of heterozygosity, as are observed in these

619  species, can increase local and genome-wide mutation rates [57].

620      Compared to the consistent patterns of diversity among species, the much

621  weaker correlations observed for allele frequency spectrum (Tajima's D) could either

622  be ascribed to divergent selective targets to different environments since their

26

623    divergence, or different demographic histories experienced by the three species during

624    the Quaternary ice ages [39, 44, 58]. In particular, the genome-wide excess of rare

625    frequency alleles in *P. tremuloides* is most likely explained by a recent substantial

626    population expansion that was specific to this species. Many other factors, such as

627    population structure, an unbalanced sampling scheme, and hybridization, can

628    influence estimates of genomic variation and may therefore contribute to the different

629    patterns observed among species [40, 41, 59], but we were able to exclude each of

630    these. First, in accordance with mating characteristics of the genus where the seed and

631    pollen are both wind-dispersed, we found little evidence of population structure in any

632    of the three *Populus* species in this study. Second, despite a potentially biased

633    sampling scheme in *P. tremuloides*, where the samples were collected from two

634    geographically distinct populations, when we analyzed the genome-wide patterns of

635    polymorphism separately we found the same patterns as those obtained when

636    analyzing the populations jointly. Third, with regard to the influence of hybridization,

637    it should be noted that there are no other species of *Populus* that occur naturally in the

638    regions from where the *P. tremula* samples were collected [60]. For *P. tremuloides*,

639    naturally occurring hybridization is only known to occur at very low levels with *P.*

640    *grandidentata* [61]. These two species occur sympatrically in central and eastern

641    North America, so in our study any possible hybridization in *P. tremuloides* would be

642    limited to samples from the Wisconsin population. Although hybridization with other

643    nearby *Populus* species is more frequent in *P. trichocarpa* [41], in this study we only

644    used individuals of *P. trichocarpa* that have previously been shown having no

645    evidence of admixture with other species [31].

646         The recombination rate and the extent of linkage disequilibrium (LD) are key

647    factors influencing the feasibility and power of genome-wide association studies [62,

27

648    63]. The three *Populus* species exhibit different patterns in the decay of LD ($r^2$) with

649    physical distance, with LD decaying fastest in *P. tremuloides* and slowest in *P.*

650    *trichocarpa*. This reflects the rank order of their population-scaled recombination rate

651    ($\rho=4N_e$c), for which *P. tremuloides* is the highest (8.42 Kbp$^{-1}$), followed by *P.*

652    *tremula* (3.23 Kbp$^{-1}$), and *P. trichocarpa* is the lowest (2.19 Kbp$^{-1}$). It is important to

653    note that differences in $\rho$ among the species cannot simply reflect differences in

654    species' $N_e$, because recombination rate correlations in 100 Kbp windows show that $\rho$

655    is only partially (not highly) conserved among these species (Fig 2c) [64]. This

656    suggests that even with conserved gene function and synteny, associations might be

657    more easily discovered in one *Populus* species than another.

658        The genome-wide ratio of recombination to mutation rate ($\rho/\theta_W$ or $c/\mu$) was

659    similar between *P. tremuloides* (0.39) and *P. trichocarpa* (0.38), but substantially

660    smaller in *P. tremula* (0.22). If mutation rate is indeed unchanged between species,

661    the lower estimate of $c/\mu$ in *P. tremula* indicates a considerably smaller recombination

662    rate relative to the other species. Nevertheless, these $c/\mu$ estimates are of the same

663    order of magnitude as recent genome-wide estimates of other plant species, such as

664    *Medicago truncatula* (0.29) [15], *Mimulus guttatus* (0.8) [65] and the tree *Eucalyptus*

665    *grandis* (0.65) [66]. However, the discrepant results obtained from patterns of

666    polymorphism and recombination between *P. tremula* and *P. tremuloides* are likely

667    due to differences in the effective population sizes influencing patterns of nucleotide

668    diversity and linkage disequilibrium [67]. These processes operate over different

669    time-scales and are therefore subject to temporal variation in the effective population

670    size [67, 68]. The recent population size expansion that we infer to have taken place

671    in *P. tremuloides* can thus also explain why its recombination rate is higher than *P.*

672    *tremula*, even if they share similar levels of genome-wide polymorphism.

28

673        In addition to the historical patterns of mutation, recombination and

674      demographic processes, patterns of genomic variation also contain much information

675      about natural selection [54]. In all three species, as expected we find 0-fold non-

676      synonymous sites exhibit significantly lower levels of polymorphism and divergence

677      compared to 4-fold synonymous sites. The 0-fold non-synonymous sites are likely

678      experiencing strong selective constraint, consistent with their excess of ultra-rare

679      variants as indicated by Tajima's D [69]. In addition, introns and 5' UTR sites are

680      also likely to be under some degree of selective constraint, although much weaker

681      than non-synonymous sites. The 3' UTR sites seem to be either neutral or under

682      comparable extent of selective constraint as 4-fold synonymous sites [70]. In contrast

683      to all genic categories, we find there are substantially higher levels of polymorphism

684      and divergence in intergenic regions throughout the genome, reflecting either higher

685      mutation rates, relaxed selective constraint or both in these regions [2].

686        Apart from strong selective constraints on protein-coding genes, multiple lines

687      of evidence indicate that genomic patterns of polymorphism have been primarily

688      shaped by widespread natural selection in all three *Populus* species. First, we find

689      significantly positive correlations between neutral polymorphism and population-

690      scaled recombination rate in both genic and intergenic regions, even after controlling

691      for the confounding variables such as GC content, gene density, mutation rate and

692      number of covered sites by the data. Such patterns could be explained by both

693      background selection and recurrent selective sweeps, where perturbations of linked

694      selection on neutral genetic variation are more drastic and extensive in regions of low

695      recombination compared to high recombination regions [8, 9, 71]. An alternative

696      explanation to natural selection would be that recombination itself has a mutagenic

697      effect [47]. In this case, the neutral theory predicts that we would also detect a

698    correlation between nucleotide divergence and recombination rate [10, 47] but this

699    relationship was not observed for any of the three species. Thus our findings support

700    the notion that ubiquitous linked selection, as selective sweeps of adaptive alleles

701    and/or background selection against deleterious alleles, is the dominant force shaping

702    the observed associations between recombination and neutral polymorphism in all

703    three species [72]. In addition, the extent of such associations can also reflect the

704    magnitude of the impact of linked selection on genomes [54, 73]. Here, we tried to

705    decipher the factors that may contribute to inconsistent signatures and magnitudes of

706    linked selection across the three species. First of all, the genome-wide effects of

707    linked selection ought to be influenced by effective population size ($N_e$) across

708    species, where the impact of selection at linked sites should be more severe in larger

709    populations [73, 74]. As a result, the substantially stronger signatures of linked

710    selection in *P. tremula* and *P. tremuloides* are most likely due to their larger $N_e$

711    compared to *P. trichocarpa*. Furthermore, just as the impact of natural selection at

712    linked sites depends on the local environment of recombination, we expect that the

713    disparate patterns of linked selection among species is also likely to be caused by the

714    various recombination rates across genomes [54, 71]. In particular, compared with *P.*

715    *tremuloides*, the stronger signature of linked selection in *P. tremula* is supposed to be

716    primarily driven by its lower average levels of recombination across the genome.

717    More broadly, the different magnitude of linked selection may provide one of the

718    major explanations for the disparate patterns of genomic variation across related

719    species [73].

720         In addition to the association between recombination and neutral

721    polymorphism, we find slightly negative correlations between recombination rate and

722    the ratio of non-synonymous- to synonymous- polymorphism, but not divergence, in

723 *P. tremula* and *P. tremuloides* after controlling for the confounding variables. This

724 pattern indicates a potential reduced efficacy of purifying selection at eliminating

725 weakly deleterious non-synonymous mutations in low recombination regions [7, 48].

726 As a consequence, such Hill-Robertson interference (HRI) may help to understand

727 patterns of partially positive correlations between gene density and recombination rate

728 among all three species [13]. Given the relaxed efficacy of purifying selection in

729 regions of low recombination where weakly deleterious mutations are more likely to

730 accumulate at a high rate, important functional elements should thus not cluster in

731 these regions, as has already been shown in several other plant species [14, 15, 75]

732 Consistent with this prediction [76], we find positive association between gene

733 density and recombination rate in regions that experience low rates of recombination.

734 In high-recombination regions where selection is more effective at eliminating

735 slightly deleterious mutations, the association between gene density and

736 recombination become much weaker in all three species. However, it remains unclear

737 whether it is the effects of recombination gradients that drive the functional

738 organization of genomes in response to selection, or it is the gradients of functional

739 genomic elements that modify the evolution of recombination rates in *Populus*.

740        By examining the relationship of neutral polymorphism, recombination rate

741 and gene density, we find that levels of neutral polymorphism in genic regions are

742 primarily dominated by local rates of recombination, regardless of the density of

743 functional genes nearby. This suggests that widespread selection might have

744 uniformly shaped the patterns of neutral polymorphism in genic regions across the

745 genome, with variation of genetic diversity primarily relying on the variation of local

746 recombination rates [7, 49, 71]. However, there is a more complex pattern in

747 intergenic regions where levels of intergenic polymorphism are mainly dominated by

31

748   recombination rates in regions of lower gene density, while in regions of higher gene

749   density, levels of intergenic diversity are primarily shaped by the density of genes

750   nearby. Patterns of polymorphism vary, on both quantitative and qualitative scales,

751   between genic and intergenic sequences, with the latter exhibiting substantially higher

752   diversity, divergence and more non-uniformly distributed selective effects compared

753   with the former [54]. In addition, 84.2% of the intergenic sites included in this study

754   are located within 5-Kbp upstream/downstream regions of functional genes. This

755   suggests that many of these intergenic regions may have important functions in gene

756   regulation, in accordance with the widespread signatures of linked selection as we

757   found in these regions [77]. In these cases, we could argue that the differences in

758   neutral mutation rate alone are not sufficient to explain the distinct patterns of genetic

759   variation between genic and intergenic sites. Various rates, distributions, and selective

760   coefficients for either adaptive or deleterious mutations, however, may at least in part

761   drive the distinct patterns of polymorphism and divergence between these different

762   genomic environments.

763        In conclusion, we have examined and compared the relative roles of mutation,

764   population history, recombination and natural selection in forging the landscape

765   heterogeneity of genomic variation within and among three related *Populus* species.

766   We find substantially different magnitudes and signatures of linked selection among

767   species, with selection effects being strongest in *P. tremula* and weakest in *P.*

768   *trichocarpa*. Various effective population sizes and genome-wide recombination rates

769   are likely to be the primary factors causing the disparate genome-wide signatures of

770   linked selection among species. By analyzing the ratio of non-synonymous- to

771   synonymous- polymorphism along recombination gradients, we find that purifying

772   selection at purging slightly deleterious non-synonymous mutations is more effective

32

773    in regions experiencing high recombination. Such selective interaction between

774    recombination and selection may provide one of the explanations for the co-varying

775    patterns of gene density and recombination in the *Populus* species, where functional

776    genes are more likely to cluster in high-recombination regions. Finally, we find

777    distinct genomic signatures of selection between genic and intergenic regions. The

778    recombination rate-dependent effect of selection dominates levels of polymorphism at

779    genic sites, while patterns of linked selection at intergenic sites are shaped by

780    interactions between recombination and local gene density. Thus, our study provides a

781    promising avenue to dissect how interactions of various evolutionary forces are

782    driving the evolution of genomes for even closely related species.

783

784    **Materials and Methods**

785

786    **Samples and sequencing**

787    Leaf samples were separately collected from 24 genotypes of *P. tremula* and 24

788    genotypes of *P. tremuloides* (S1 Table). Genomic DNA was extracted from leaf

789    samples, and paired-end sequencing libraries with insert sizes of 650bp were

790    constructed for all genotypes. Whole-genome sequencing with a minimum expected

791    depth of $20 \times$ was performed on the Illumina HiSeq 2000 platform, and $2 \times 100$-bp

792    paired-end reads were generated for all genotypes. As two samples of *P. tremuloides*

793    failed to obtain the expected coverage, all analyses are based on data from 24 *P.*

794    *tremula* genotypes and 22 *P. tremuloides* genotypes. All newly generated Illumina

795    reads from this study have been submitted to the Short Read Archive at NCBI under

796    accession IDs ranging from XXXXXX-XXXXXX. We obtained publicly available

797    short read Illumina data for 24 *P. trichocarpa* individuals from NCBI SRA (S1 Table).

798    Individuals were selected to have a similar read depth as the samples of the two aspen

799    species. The accession numbers of *P. trichcoarpa* samples can be found in [31]. These

800    data are paired-end 100bp reads generated on the Illumina HiSeq2000 platform.

801

802    **Raw read filtering, read alignment and post-processing alignment**

803    Prior to read alignment, we used Trimmomatic [78] to remove adapter sequences

804    from reads. Since the quality of reads always drops towards the end of reads, we used

805    Trimmomatic to cut bases off the start and end of each read when the quality values

806    dropped below 20. If the length of the processed reads was reduced to below 36 bases

807    after      trimming,      reads      were      completely      discarded.      FastQC

808    (http://www.bioinformatics.babraham.ac.uk/projects/fastqc/) was used to check and

809    compare the per base sequence quality of the raw sequence data and the filtered data.

810    After quality control, all paired-end and orphaned single-end reads of each sample

811    were mapped to the *P. trichocarpa* version 3 (v3.0) genome [29] using the BWA-

812    MEM algorithm with default parameters in bwa-0.7.10 [32].

813         Several post-processing steps of alignments were performed in order to

814    minimize the number of artifacts in downstream analysis: First, indel realignment was

815    performed as sequence reads are often mapped with mismatching bases in regions

816    with insertions and deletions (indels). The RealignerTargetCreator in GATK (The

817    Genome Analysis Toolkit) [79] was first used to find suspicious-looking intervals

818    which were likely in need of realignment. Then, the IndelRealigner was used to run

819    the realigner over those intervals. Second, as reads resulting from PCR duplicates can

820    arise during the sequencing library preparation, we used the MarkDuplicates methods

821    in the Picard package (http://picard.sourceforge.net) to remove those reads or read

822    pairs having identical external coordinates and the same insert length. In such cases

34

823     only the single read with the highest summed base qualities was kept for downstream

824     analysis. Third, in order to exclude genotyping errors caused by paralogous or

825     repetitive DNA sequences where reads were poorly mapped to the reference genome,

826     or by other genome feature differences between *P. trichocarpa* and *P. tremula* or *P.*

827     *tremuloides*, we removed sites with extremely low or high read depths. After

828     investigating the empirical distribution of read coverage, we filtered out sites with a

829     total coverage less than 100X or greater than 1200X across all samples per species.

830     When reads were mapped to multiple locations in the genome, they were randomly

831     assigned to one location with a mapping score of zero by BWA-MEM. In order to

832     account for such misalignment effects for each species, we removed those sites if

833     there were more than 20 mapped reads with mapping score equal to zero across all

834     individuals. Lastly, because the short read alignment is generally unreliable in highly

835     repetitive genomic regions, we filtered out sites that overlapped with known repeat

836     elements as identified by RepeatMasker [80]. In the end, the subset of sites that

837     passed all these filtering criteria in the three *Populus* species were used in all

838     following analyses.

839

840     **SNP and genotype calling**

841     We implemented two complementary bioinformatics approaches in downstream

842     analyses:

843     (i) Population genetic inferences that rely on the site frequency spectrum (SFS).

844     Recently, many studies pointed out the bias introduced in population genetic estimates

845     by inaccurate genotype calls from NGS data [35, 81]. Either the single-sample

846     genotype calling (calling genotypes for each individual separately and then merging

847     them later) or the multi-sample genotype calling (jointly calling genotypes for all

848     individuals) can result in a bias in the estimation of SFS, as the former method usually

849     leads to overestimation of rare variants, whereas the latter often leads to the opposite

850     [35]. Therefore, all the population genetic statistics that based on the SFS in this study

851     were estimated directly and jointly from filtered sites and individuals without calling

852     genotypes, as implemented in the software package Analysis of Next-Generation

853     Sequencing Data (ANGSD v0.602) [82].

854     (ii) Analyses based on accurate SNP and genotype calls. We performed SNP calling

855     with HaplotypeCaller of the GATK v3.2.2 [79], which called SNPs and indels

856     simultaneously via local re-assembly of haplotypes for each individual and created

857     single-sample gVCFs. GenotypeGVCFs in GATK was then used to merge multi-

858     sample records together, correct genotype likelihoods, and re-genotype the newly

859     merged record and perform re-annotation. Several filtering steps were then used to

860     reduce the number of false positive SNPs and retain high-quality SNPs: (1) we

861     removed all SNPs that were overlapped with sites excluded by the previous filtering

862     criteria. (2) only biallelic SNPs with a distance of more than 5bp away from indels

863     were retained for further analysis. (3) Genotypes were accepted for each SNP and

864     each individual only if the genotype quality score (GQ) was ≥10, otherwise that

865     specific genotype was treated as missing data. (4) SNPs with missing rate higher than

866     20% were removed from downstream analysis. (5) SNPs that showed significant

867     deviation from Hardy-Weinberg Equilibrium ($P<0.001$) were removed from further

868     downstream analysis.

869

870     **Population structure**

871     We used only 4-fold synonymous SNPs with minor allele frequency >0.1 to perform

872     population structure analyses with ADMIXTURE [34]. We ran ADMIXTURE

36

873    separately on all the sampled individuals among species and on the samples within

874    each species, varying the number of genetic clusters $K$ from 1 to 6. The most likely

875    number of genetic cluster was selected by minimizing the cross-validation error in

876    ADMIXTURE.

877

878    **Diversity and divergence - related summary statistics**

879    For nucleotide diversity and divergence estimates, only the reads with mapping

880    quality above 30 and the bases with quality score higher than 20 were used in all of

881    the following analyses with ANGSD [82] and ngsTools [83]. To infer the global SFS,

882    we firstly used the -doSaf implementation in ANGSD to calculate the site allele

883    frequency likelihood based on the SAMTools genotype likelihood model [84]. Then,

884    we used the –realSFS implementation in ANGSD to obtain an optimized folded

885    global SFS using Expectation Maximization (EM) algorithm for each species. Based

886    on the global SFS, we used the –doThetas function in ANGSD to estimate the per-site

887    nucleotide diversity from posterior probability of allele frequency based on a

888    maximum likelihood approach [36]. Two standard estimates of nucleotide diversity,

889    the average pairwise nucleotide diversity ($\Theta_\pi$) [38] and the proportion of segregating

890    sites ($\Theta_W$) [85], and one neutrality statistic test Tajima's D [38] were then

891    summarized along all 19 chromosomes using non-overlapping sliding windows of 100

892    Kbp and 1 Mbp. Windows with less than 10% covered sites left after previous quality

893    filtering steps were excluded. Accordingly, 3340 100-Kbp and 343 1-Mbp windows,

894    with an average of 50,538 and 455,910 covered bases per window, were respectively

895    included for downstream analyses. Based on posterior probabilities of sample allele

896    frequencies at each site, we further used the ngsTools [83] to calculate pairwise

37

897  nucleotide divergence, $d_{xy}$, between pairs of species over all non-overlapping 100-

898  Kbp and 1-Mbp windows.

899      All these statistics were also calculated for each type of functional element (0-

900  fold non-synonymous, 4-fold synonymous, intron, 3' UTRs, 5' UTR, and intergenic

901  sites) over the non-overlapping 100-Kbp and 1-Mbp windows in all three *Populus*

902  species. The category of gene models we used followed the gene annotation of *P.*

903  *trichocarpa* version 3.0 [29]. For protein-coding genes, we only included genes with

904  at least 90% covered sites left from previous filtering steps to ensure that the three

905  species have same gene structures. We also excluded genes overlapping with other

906  genes. For the remaining genes, we selected the transcript with the highest content of

907  protein-coding sites. For regions overlapped by different transcripts in each gene, we

908  classified each site according to the following hierarchy (from highest to lowest):

909  Coding regions (CDS), 3'UTR, 5'UTR, Intron. Thus, if a site resides in a 3'UTR in

910  one transcript and CDS for another, the site was classified as CDS. In the end, a

911  respective of 16.52, 3.4, 7.19, 4.02, 31.89, 73.46 megabases (Mbp) were partitioned

912  into 0-fold non-synonymous (where all DNA sequence changes lead to protein

913  sequence changes), 4-fold synonymous (where all DNA sequence changes lead to the

914  same protein sequences), 3'UTR, 5'UTR, intron, and intergenic categories. Windows

915  were not used if there were less than 100 sites left for any of the functional elements.

916

917  **Linkage disequilibrium (LD) and population-scaled recombination rate ($\rho$)**

918  A total of 1,409,377 SNPs, 1,263,661 SNPs and 710,332 SNPs with minor allele

919  frequency higher than 10% were used for the analysis of LD and $\rho$ in *P. tremula*, *P.*

920  *tremuloides* and *P. trichocarpa*, respectively. To estimate and compare the rate of LD

921  decay in the three *Populus* species, we firstly used PLINK 1.9 [86] to randomly thin

922    the number of SNPs to 100,000 in each species. Then we calculated the squared

923    correlation coefficients ($r^2$) between all pairs of SNPs within 50 Kbp windows using

924    PLINK 1.9. The decay of LD against physical distance was estimated using nonlinear

925    regression of pairwise $r^2$ vs. the physical distance between sites in base pairs [87].

926    Furthermore, we estimated the population-scaled recombination rate $\rho$ using the

927    Interval program of LDhat 2.2 [88] with 1,000,000 MCMC iterations sampling every

928    2,000 iterations and a block penalty parameter of five. The first 100,000 iterations of

929    the MCMC iterations were discarded as a burn-in. We then calculated the scaled value

930    of $\rho$ in each 100-Kbp and 1-Mbp window as the average across SNPs in that window.

931    In order to evaluate the extent of correlation between the estimated $\rho$ and the pattern

932    of LD, we also calculated the scaled $r^2$ by averaging $r^2$ over all pairwise SNPs in each

933    100 Kbp and 1 Mbp window. Only windows with more than 10,000 (in 100 Kbp

934    windows) and 100,000 bases (in 1 Mbp windows) and 100 SNPs left after previous

935    filtering steps were used for the estimation of $\rho$ and $r^2$.

936

937    **Genomic correlates of diversity**

938    Within each non-overlapping 100 Kbp or 1 Mbp window, levels of neutral

939    polymorphism in genic and intergenic regions were tabulated as the pairwise

940    nucleotide diversity ($\Theta_\pi$) at 4-fold synonymous and intergenic sites respectively. In

941    order to examine the factors influencing levels of neutral polymorphism in all three

942    *Populus* species, we further tabulated several genomic features within each window.

943    First, we summarized population-scaled recombination rate ($\rho$) as described above for

944    each species. Second, we tabulated GC content as the fraction of bases where the

945    reference sequence (*P. trichocarpa* v3.0) was a G or a C. Third, we measured the

946    gene density as the number of functional genes within each window according to the

947     gene annotation of *P. trichocarpa* version 3.0. Fourth, we accounted for the variation

948     of mutation rate by calculating the number of fixed differences between aspen and *P.*

949     *trichocarpa* per neutral site (either 4-fold synonymous site or intergenic site) within

950     each window. The reason why we used divergence between aspen and *P. trichocarpa*

951     to measure mutation rate is because they are distantly related [26], and thus the

952     estimate of divergence are unlikely to be influenced by shared ancestral

953     polymorphisms between species. Fifth, we tabulated the number of covered bases in

954     each window as those met the filtering criteria described above.

955         We used Spearman's rank-order correlation tests to examine pairwise

956     correlations between the variables as described above. In order to account for the

957     autocorrelation between many of these variables, we calculated partial correlations

958     between the interested variables [89], which simultaneously remove the confounding

959     effects of other variables. All statistical tests were performed using R version 3.2.0

960     unless stated otherwise.

961

962     **Reference**

963     1.     Nordborg M, Hu TT, Ishino Y, Jhaveri J, Toomajian C, Zheng H, et al. The

964     pattern of polymorphism in *Arabidopsis thaliana*. PLoS Biol. 2005;3:1289.

965     2.     Begun DJ, Holloway AK, Stevens K, Hillier LW, Poh Y-P, Hahn MW, et al.

966     Population genomics: whole-genome analysis of polymorphism and divergence in

967     *Drosophila simulans*. PLoS Biol. 2007;5:e310.

968     3.     Hellmann I, Prüfer K, Ji H, Zody MC, Pääbo S, Ptak SE. Why do human

969     diversity levels vary at a megabase scale? Genome Res. 2005;15:1222-1231.

970     4.     Kimura M. The neutral theory of molecular evolution: Cambridge University

971     Press; 1984.

972    5.        Li H, Durbin R. Inference of human population history from individual whole-

973    genome sequences. Nature. 2011;475:493-496.

974    6.        Mackay TF, Richards S, Stone EA, Barbadilla A, Ayroles JF, Zhu D, et al.

975    The *Drosophila melanogaster* genetic reference panel. Nature. 2012;482:173-178.

976    7.        Cutter AD, Choi JY. Natural selection shapes nucleotide polymorphism across

977    the genome of the nematode *Caenorhabditis briggsae*. Genome Res. 2010;20:1103-

978    1111.

979    8.        Begun DJ, Aquadro CF. Levels of naturally occurring DNA polymorphism

980    correlate with recombination rates in *D. melanogaster*. Nature. 1992;356:519 - 520.

981    9.        McGaugh SE, Heil CS, Manzano-Winkler B, Loewe L, Goldstein S, Himmel

982    TL, et al. Recombination modulates how selection affects linked sites in *Drosophila*.

983    PLoS Biol. 2012;10:e1001422.

984    10.      Kulathinal RJ, Bennett SM, Fitzpatrick CL, Noor MA. Fine-scale mapping of

985    recombination rate in *Drosophila* refines its correlation to diversity and divergence.

986    Proc Natl Acad Sci U S A. 2008;105:10051-10056.

987    11.      Campos JL, Halligan DL, Haddrill PR, Charlesworth B. The relation between

988    recombination rate and patterns of molecular evolution and variation in *Drosophila*

989    *melanogaster*. Mol Biol Evol. 2014;31:1010-1028.

990    12.      Charlesworth B, Campos JL. The relations between recombination rate and

991    patterns of molecular variation and evolution in *Drosophila*. Annu Rev Genet.

992    2014;48:383-403.

993    13.      Gaut BS, Wright SI, Rizzon C, Dvorak J, Anderson LK. Recombination: an

994    underappreciated factor in the evolution of plant genomes. Nat Rev Genet. 2007;8:77-

995    84.

996  14.   Flowers JM, Molina J, Rubinstein S, Huang P, Schaal BA, Purugganan MD.

997  Natural selection in gene-dense regions shapes the genomic pattern of polymorphism

998  in wild and domesticated rice. Mol Biol Evol. 2012;29:675-687.

999  15.   Branca A, Paape TD, Zhou P, Briskine R, Farmer AD, Mudge J, et al. Whole-

1000 genome nucleotide diversity, recombination, and linkage disequilibrium in the model

1001 legume *Medicago truncatula*. Proc Natl Acad Sci U S A. 2011;108:E864-E870.

1002 16.   Luikart G, England PR, Tallmon D, Jordan S, Taberlet P. The power and

1003 promise of population genomics: from genotyping to genome typing. Nat Rev Genet.

1004 2003;4:981-994.

1005 17.   Ellegren H. Genome sequencing and population genomics in non-model

1006 organisms. Trends Ecol Evol. 2014;29:51-63.

1007 18.   Lawrie DS, Petrov DA. Comparative population genomics: power and

1008 principles for the inference of functionality. Trends Genet. 2014;30:133-139.

1009 19.   Hufford MB, Xu X, Van Heerwaarden J, Pyhäjärvi T, Chia J-M, Cartwright

1010 RA, et al. Comparative population genomics of maize domestication and

1011 improvement. Nature Genet. 2012;44:808-811.

1012 20.   Locke DP, Hillier LW, Warren WC, Worley KC, Nazareth LV, Muzny DM, et

1013 al. Comparative and demographic analysis of orang-utan genomes. Nature.

1014 2011;469:529-533.

1015 21.   Liu S, Lorenzen ED, Fumagalli M, Li B, Harris K, Xiong Z, et al. Population

1016 genomics reveal recent speciation and rapid evolutionary adaptation in polar bears.

1017 Cell. 2014;157:785-794.

1018 22.   Neale DB, Kremer A. Forest tree genomics: growing resources and

1019 applications. Nat Rev Genet. 2011;12:111-122.

1020    23.    González‑Martínez SC, Krutovsky KV, Neale DB. Forest‑tree population

1021    genomics and adaptive evolution. New Phytol. 2006;170:227-238.

1022    24.    Eckenwalder JE. Systematics and evolution of *Populus*. In: Stettler RF,

1023    Bradshaw HD, Heilman PE, Hinckley TM, editors. Biology of *Populus* and its

1024    Implications for Management and Conservation. Ottawa: NRC Research Press;

1025    1996.pp.7-32.

1026    25.    Jansson S, Douglas CJ. *Populus*: a model system for plant biology. Annu Rev

1027    Plant Biol. 2007;58:435-458.

1028    26.    Wang Z, Du S, Dayanandan S, Wang D, Zeng Y, Zhang J. Phylogeny

1029    Reconstruction and Hybrid Analysis of *Populus* (Salicaceae) Based on Nucleotide

1030    Sequences of Multiple Single-Copy Nuclear Genes and Plastid Fragments. Plos One.

1031    2014;9:e103645.

1032    27.    Jansson S, Bhalerao RP, Groover AT. Genetics and genomics of *Populus*:

1033    Springer; 2010.

1034    28.    Dickmann DI, Kuzovkina J. Poplars and willows of the world, with emphasis

1035    on silviculturally important species. In: Isebrands JG, Richardson J, editors. Poplars

1036    and Willows: trees for society and the environment. Rome: The Food and Agriculture

1037    Organization of the United Nations (FAO) and CAB International (CABI); 2014:8-91.

1038    29.    Tuskan GA, Difazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, et al.

1039    The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). Science.

1040    2006;313:1596-604.

1041    30.    Hamzeh M, Dayanandan S. Phylogeny of *Populus* (Salicaceae) based on

1042    nucleotide sequences of chloroplast trnT-trnF region and nuclear rDNA. Am J Bot.

1043    2004;91:1398-408.

1044    31.    Evans LM, Slavov GT, Rodgers-Melnick E, Martin J, Ranjan P, Muchero W,

1045    et al. Population genomics of *Populus trichocarpa* identifies signatures of selection

1046    and adaptive trait associations. Nature Genet. 2014;46:1089–1096.

1047    32.    Li H. Aligning sequence reads, clone sequences and assembly contigs with

1048    BWA-MEM; 2013. Preprint. Available: arXiv:13033997.

1049    33.    Hall D, Luquez V, Garcia VM, St Onge KR, Jansson S, Ingvarsson PK.

1050    Adaptive population differentiation in phenology across a latitudinal gradient in

1051    European aspen (*Populus tremula*, L.): a comparison of neutral markers, candidate

1052    genes and phenotypic traits. Evolution. 2007;61:2849-2860.

1053    34.    Alexander DH, Novembre J, Lange K. Fast model-based estimation of

1054    ancestry in unrelated individuals. Genome Res. 2009;19:1655-1664.

1055    35.    Nielsen R, Korneliussen T, Albrechtsen A, Li Y, Wang J. SNP calling,

1056    genotype calling, and sample allele frequency estimation from New-Generation

1057    Sequencing data. PLoS One. 2011;7:e37558.

1058    36.    Kim SY, Lohmueller KE, Albrechtsen A, Li Y, Korneliussen T, Tian G, et al.

1059    Estimation of allele frequency and association mapping using next-generation

1060    sequencing data. BMC bioinformatics. 2011;12:231.

1061    37.    Ismail M, Soolanayakanahally RY, Ingvarsson PK, Guy RD, Jansson S, Silim

1062    SN, et al. Comparative nucleotide diversity across North American and European

1063    *Populus* species. J Mol Evol. 2012;74:257-272.

1064    38.    Tajima F. Statistical method for testing the neutral mutation hypothesis by

1065    DNA polymorphism. Genetics. 1989;123:585-595.

1066    39.    Zhou L, Bawa R, Holliday J. Exome resequencing reveals signatures of

1067    demographic and adaptive processes across the genome and range of black

1068    cottonwood (*Populus trichocarpa*). Mol Ecol. 2014;23:2486-2499.

1069    40.    Chikhi L, Sousa VC, Luisi P, Goossens B, Beaumont MA. The confounding

1070    effects of population structure, genetic diversity and the sampling scheme on the

1071    detection and quantification of population size changes. Genetics. 2010;186:983-995.

1072    41.    Huang DI, Hefer CA, Kolosova N, Douglas CJ, Cronk QC. Whole plastome

1073    sequencing reveals deep plastid divergence and cytonuclear discordance between

1074    closely related balsam poplars, *Populus balsamifera* and *P. trichocarpa* (Salicaceae).

1075    New Phytol. 2014;204:693-703.

1076    42.    Wang J, Scofield D, Street NR, Ingvarsson PK. Variant calling using NGS

1077    data in European aspen (*Populus tremula*). In: Sablok G, Kumar S, Ueno S, Kuo J,

1078    Varotto C, editors. Advances in the Understanding of Biological Sciences Using Next

1079    Generation Sequencing (NGS) Approaches. Springer; 2015. pp.43-61.

1080    43.    Begun DJ. Population genetics of silent and replacement variation in

1081    *Drosophila simulans* and *D. melanogaster*: X/autosome differences? Mol Biol Evol.

1082    1996;13:1405-1407.

1083    44.    Ingvarsson PK. Multilocus patterns of nucleotide polymorphism and the

1084    demographic history of *Populus tremula*. Genetics. 2008;180:329-340.

1085    45.    Slavov GT, DiFazio SP, Martin J, Schackwitz W, Muchero W, Rodgers‐

1086    Melnick E, et al. Genome resequencing reveals multiscale geographic structure and

1087    extensive linkage disequilibrium in the forest tree *Populus trichocarpa*. New Phytol.

1088    2012;196:713-725.

1089    46.    Comeron JM, Kreitman M. The correlation between synonymous and

1090    nonsynonymous substitutions in *Drosophila*: mutation, selection or relaxed

1091    constraints? Genetics. 1998;150:767-775.

1092    47.    Hellmann I, Ebersberger I, Ptak SE, Pääbo S, Przeworski M. A neutral

1093    explanation for the correlation of diversity with recombination rates in humans. Am J

1094    Hum Genet. 2003;72:1527-1735.

1095    48.    Hill WG, Robertson A. The effect of linkage on limits to artificial selection.

1096    Genetical Res. 1966;8:269-294.

1097    49.    Cutter AD, Payseur BA. Selection at linked sites in the partial selfer

1098    *Caenorhabditis elegans*. Mol Biol Evol. 2003;20:665-673.

1099    50.    Payseur BA, Nachman MW. Gene density and human nucleotide

1100    polymorphism. Mol Biol Evol. 2002;19:336-340.

1101    51.    Andolfatto P. Hitchhiking effects of recurrent beneficial amino acid

1102    substitutions in the *Drosophila melanogaster* genome. Genome Res. 2007;17:1755-

1103    1762.

1104    52.    Macpherson JM, Sella G, Davis JC, Petrov DA. Genomewide spatial

1105    correspondence between nonsynonymous divergence and neutral polymorphism

1106    reveals extensive adaptation in *Drosophila*. Genetics. 2007;177:2083-2099.

1107    53.    Ingvarsson PK. Natural selection on synonymous and nonsynonymous

1108    mutations shapes patterns of polymorphism in *Populus tremula*. Mol Biol Evol.

1109    2010;27:650-660.

1110    54.    Cutter AD, Payseur BA. Genomic signatures of selection at linked sites:

1111    unifying the disparity among species. Nat Rev Genet. 2013;14:262-274.

1112    55.    Hudson RR, Kreitman M, Aguadé M. A test of neutral molecular evolution

1113    based on nucleotide data. Genetics. 1987;116:153-159.

1114    56.    Charlesworth B, Morgan M, Charlesworth D. The effect of deleterious

1115    mutations on neutral molecular variation. Genetics. 1993;134:1289-1303.

1116    57.    Yang S, Wang L, Huang J, Zhang X, Yuan Y, Chen J-Q, et al. Parent-progeny

1117    sequencing indicates higher mutation rates in heterozygotes. Nature. 2015;523:463-

1118    467.

1119    58.    Callahan CM, Rowe CA, Ryel RJ, Shaw JD, Madritch MD, Mock KE.

1120    Continental‐scale assessment of genetic diversity and population structure in

1121    quaking aspen (*Populus tremuloides*). J Biogeogr. 2013;40:1780-1791.

1122    59.    Städler T, Haubold B, Merino C, Stephan W, Pfaffelhuber P. The impact of

1123    sampling schemes on the site frequency spectrum in nonequilibrium subdivided

1124    populations. Genetics. 2009;182:205-216.

1125    60.    Lexer C, Fay M, Joseph J, Nica MS, Heinze B. Barrier to gene flow between

1126    two ecologically divergent *Populus* species, *P. alba* (white poplar) and *P. tremula*

1127    (European aspen): the role of ecology and life history in gene introgression. Mol Ecol.

1128    2005;14:1045-1057.

1129    61.    Pregitzer KS, Barnes BV. Flowering phenology of *Populus tremuloides* and *P.*

1130    *grandidentata* and the potential for hybridization. Can J Forest Res. 1980;10:218-223.

1131    62.    Kim S, Plagnol V, Hu TT, Toomajian C, Clark RM, Ossowski S, et al.

1132    Recombination and linkage disequilibrium in *Arabidopsis thaliana*. Nature Genet.

1133    2007;39:1151-1155.

1134    63.    Neale DB, Ingvarsson PK. Population, quantitative and comparative genomics

1135    of adaptation in forest trees. Curr Opin Plant Biol. 2008;11:149-155.

1136    64.    Smukowski C, Noor M. Recombination rate variation in closely related

1137    species. Heredity. 2011;107:496-508.

1138    65.    Hellsten U, Wright KM, Jenkins J, Shu S, Yuan Y, Wessler SR, et al. Fine-

1139    scale variation in meiotic recombination in Mimulus inferred from population shotgun

1140    sequencing. Proc Natl Acad Sci U S A. 2013;110:19478-19482.

1141    66.    Silva‐Junior OB, Grattapaglia D. Genome‐wide patterns of recombination,

1142    linkage disequilibrium and nucleotide diversity from pooled resequencing and single

1143    nucleotide polymorphism genotyping unlock the evolutionary history of *Eucalyptus*

1144    *grandis*. New Phytol. 2015; doi: 10.1111/nph.13505.

1145    67.    Tenesa A, Navarro P, Hayes BJ, Duffy DL, Clarke GM, Goddard ME, et al.

1146    Recent human effective population size estimated from linkage disequilibrium.

1147    Genome Res. 2007;17:520-526.

1148    68.    Cutter AD, Jovelin R, Dey A. Molecular hyperdiversity and evolution in very

1149    large populations. Mol Ecol. 2013;22:2074-2095.

1150    69.    Larracuente AM, Sackton TB, Greenberg AJ, Wong A, Singh ND, Sturgill D,

1151    et al. Evolution of protein-coding genes in *Drosophila*. Trends Genet. 2008;24:114-

1152    123.

1153    70.    Andolfatto P. Adaptive evolution of non-coding DNA in *Drosophila*. Nature.

1154    2005;437:1149-1152.

1155    71.    Slotte T. The impact of linked selection on plant genomic variation. Brief

1156    Funct Genomics. 2014;13:268-275.

1157    72.    Hahn MW. Toward a selection theory of molecular evolution. Evolution.

1158    2008;62:255-265.

1159    73.    Corbett-Detig RB, Hartl DL, Sackton TB. Natural selection constrains neutral

1160    diversity across a wide range of species. PLoS Biol. 2015;13:e1002112.

1161    74.    Leffler EM, Bullaughey K, Matute DR, Meyer WK, Segurel L, Venkat A, et

1162    al. Revisiting an old riddle: what determines genetic diversity levels within species.

1163    PLoS Biol. 2012;10:e1001388.

1164    75.    Anderson LK, Lai A, Stack SM, Rizzon C, Gaut BS. Uneven distribution of

1165    expressed sequence tag loci on maize pachytene chromosomes. Genome Res.

1166    2006;16:115-122.

1167    76.    Haddrill PR, Halligan DL, Tomaras D, Charlesworth B. Reduced efficacy of

1168    selection in regions of the *Drosophila* genome that lack crossing over. Genome Biol.

1169    2007;8:R18.

1170    77.    Haygood R, Fedrigo O, Hanson B, Yokoyama K-D, Wray GA. Promoter

1171    regions of many neural-and nutrition-related genes have experienced positive

1172    selection during human evolution. Nature Genet. 2007;39:1140-1144.

1173    78.    Lohse M, Bolger A, Nagel A, Fernie AR, Lunn JE, Stitt M, et al. RobiNA: a

1174    user-friendly, integrated software solution for RNA-Seq-based transcriptomics.

1175    Nucleic Acids Res. 2012;40:W622-W627.

1176    79.    DePristo MA, Banks E, Poplin R, Garimella KV, Maguire JR, Hartl C, et al. A

1177    framework for variation discovery and genotyping using next-generation DNA

1178    sequencing data. Nature Genet. 2011;43:491-498.

1179    80.    Tarailo‑Graovac M, Chen N. Using RepeatMasker to identify repetitive

1180    elements in genomic sequences. Curr Protoc in Bioinformatics. 2009; 25:4.10.1-

1181    4.10.14.

1182    81.    Nevado B, Ramos‑Onsins S, Perez‑Enciso M. Resequencing studies of

1183    nonmodel organisms using closely related reference genomes: optimal experimental

1184    designs and bioinformatics approaches for population genomics. Mol Ecol.

1185    2014;23:1764-1779.

1186    82.    Korneliussen TS, Albrechtsen A, Nielsen R. ANGSD: analysis of next

1187    generation sequencing data. BMC bioinformatics. 2014;15:356.

1188    83.    Fumagalli M, Vieira FG, Linderoth T, Nielsen R. ngsTools: methods for

1189    population genetics analyses from next-generation sequencing data. Bioinformatics.

1190    2014;30:1486-1487.

1191    84.    Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The

1192    sequence alignment/map format and SAMtools. Bioinformatics. 2009;25:2078-2079.

1193    85.    Watterson G. On the number of segregating sites in genetical models without

1194    recombination. Theor Pop Biol. 1975;7:256-276.

1195    86.    Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D, et al.

1196    PLINK: a tool set for whole-genome association and population-based linkage

1197    analyses. Am J Hum Genet. 2007;81:559-575.

1198    87.    Remington DL, Thornsberry JM, Matsuoka Y, Wilson LM, Whitt SR,

1199    Doebley J, et al. Structure of linkage disequilibrium and phenotypic associations in

1200    the maize genome. Proc Natl Acad Sci U S A. 2001;98:11479-11484.

1201    88.    McVean GA, Myers SR, Hunt S, Deloukas P, Bentley DR, Donnelly P. The

1202    fine-scale structure of recombination rate variation in the human genome. Science.

1203    2004;304:581-584.

1204    89.    Kim S-H, Soojin VY. Understanding relationship between sequence and

1205    functional evolution in yeast proteins. Genetica. 2007;131:151-156.

1206

1207

1208

1209

1210

1211

1212

1213 **Supporting Information**

1214

1215 **S1 Fig. Sampling localities (details in S1 Table, black star symbols) and**

1216 **distributions of *P. tremula* (orange areas), *P. tremuloides* (blue areas) and *P.***

1217 ***trichocarpa* (green areas).**

1218

1219 **S2 Fig. Comparison of per-base sequence quality between raw and filtered**

1220 **sequence data.** Per-base sequence quality comparison between raw paired-end

1221 sequence data (forward reads: top left and reverse reads: top right), and filtered

1222 sequence data with both forward (bottom left) and reverse (bottom middle) reads left

1223 or only single-end (bottom right) reads left. The x-axis of the BoxWhisker plot shows

1224 the position in read, and the y-axis shows the quality scores. The higher the score the

1225 better the base call. The background of the plot divides the y axis into very good

1226 quality calls (green), calls of reasonable quality (orange), and calls of poor quality

1227 (red). The central red line is the median quality value, the yellow box represents the

1228 inter-quartile of quality, the upper and lower whiskers represent the 10% and 90%

1229 points, the blue line represents the mean quality. (a) Sample SwAsp009 of *Populus*

1230 *tremula*. (b) Sample Alb16-1 of *P. tremuloides*. (c) Sample GW-9772 (accession

1231 number in SRA: SRR1571518) of *P. trichocarpa*.

1232

1233 **S3 Fig. Population structure within and between species.** (a) Genetic structure of

1234 three *Populus* inferred using ADMIXTURE when it identifies three genetic clusters in

1235 the dataset. (b) The cross-validation error when $K$ varies from 1 to 6 across the three

1236 species. (c,d,e) The cross-validation error when $K$ varies from 1 to 6 separately in

1237 samples of *P. tremula*, *P. tremuloides*, *P. trichocarpa*.

1238

1239   **S4 Fig. Genome-wide patterns of divergence among three _Populus_ species.** Mean

1240   pairwise divergence ($d_{xy}$) between pairs of three _Populus_ species was calculated over

1241   100 Kbp non-overlapping windows along the 19 chromosomes. _P. tremula-_

1242   _P.tremuloies_: red, _P. tremula-P. trichocarpa_: light purple, _P. tremuloides-_

1243   _P.trichocarpa_: purple.

1244

1245   **S5 Fig. Correlations of divergence between independent pairs of the three**

1246   **_Populus_ species.** Spearman's correlations of pairwise nucleotide divergence ($d_{xy}$)

1247   between $d_{xy(P.\ tremula-P.\ trichocarpa)}$ and $d_{xy(P.\ tremulodies-P.\ trichocarpa)}$ (a); between $d_{xy(P.\ tremula-P.}$

1248   $_{tremuloides)}$ and $d_{xy(P.\ tremula-P.\ trichocarpa)}$ (b); and between $d_{xy(P.\ tremula-P.\ tremuloides)}$ and $d_{xy(P.}$

1249   $_{tremuloides-P.\ trichocarpa)}$ (c). All datasets are based on 100 Kbp non-overlapping windows

1250   across the genome.

1251

1252   **S6 Fig. The distributions of estimates of (a) pairwise sequence diversity ($\Theta_\Pi$), (b)**

1253   **the number of segregating sites ($\Theta_W$), (c) Tajima's D and (d) population-scaled**

1254   **recombination rate ($\rho$) over 1Mbp non-overlapping windows in _P. tremula_**

1255   **(orange), _P. tremuloides_ (blue) and _P. trichocarpa_ (green).**

1256

1257   **S7 Fig. Estimates of (a) pairwise sequence diversity ($\Theta_\Pi$) and (b) Tajima's D in**

1258   **samples of Alberta (light blue), Wisconsin (light green) and all samples of _P._**

1259   **_tremuloides_ (blue) over 100 Kbp non-overlapping windows.**

1260

1261   **S8 Fig. Number of singletons in samples of (a) _P. tremula_, (b) _P. tremuloides_, and**

1262   **(c) _P. trichocarpa_.**

1263

1264 **S9 Fig. The distributions of estimates of pairwise sequence diversity ($\Theta_\Pi$) in *P.***

1265 ***tremula* (orange), *P. tremuloides* (blue) and *P. trichocarpa* (green) over 1 Mbp**

1266 **non-overlapping windows in different site categories.**

1267

1268 **S10 Fig. The distributions of estimates of Tajima's D in *P. tremula* (orange), *P.***

1269 ***tremuloides* (blue) and *P. trichocarpa* (green) over 1 Mbp non-overlapping**

1270 **windows in different site categories.**

1271

1272 **S11 Fig. The distributions of estimates of nucleotide divergence ($d_{xy}$) between**

1273 **pairs of the three *Populus* species over 1 Mbp non-overlapping windows in**

1274 **different site categories.**

1275

1276 **S12 Fig. Relationship between population-scaled recombination rate and linkage**

1277 **disequilibrium.** Scatter plots display correlations between population-scaled

1278 population rates ($\rho$) and linkage disequilibrium ($r^2$) over 100 Kbp non-overlapping

1279 windows in (a) *P. tremula*, (b) *P. tremuloides*, and (c) *P. trichocarpa*. The red to

1280 yellow to blue gradient indicates decreased density of observed events at a give

1281 location in the graph.

1282

1283 **S13 Fig. Distributions of the ratio of population-scaled recombination rate to**

1284 **nucleotide diversity ($\rho/\theta_W$) over 100 Kbp non-overlapping windows in *P. tremula***

1285 **(orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line).**

1286

1287 **S14 Fig. Correlations between estimates of intergenic genetic diversity ($\Theta_{\text{Intergenic}}$)**

1288 **(left panel) and divergence ($d_{\text{Intergenic}}$) (right panel) with population-scaled**

1289 **recombination rate ($\rho$) over 1 Mbp non-overlapping windows.** Linear regression

1290 lines are colored according to species: (a) *P. tremula* (orange line), (b) *P. tremuloides*

1291 (blue line) and (c) *P. trichocarpa* (green line).

1292

1293 **S15 Fig. Relationship between gene number and the proportion of coding bases**

1294 **within 1 Mbp non-overlapping windows.**

1295

1296 **S16 Fig. Correlations between estimates of genic and intergenic genetic**

1297 **divergence with gene density over 1 Mbp non-overlapping windows.** Correlations

1298 between estimates of genetic divergence at 4-fold synonymous sites ($d_{4\text{-fold}}$) (left

1299 panel) and intergenic sites ($d_{\text{Intergenic}}$) (right panel) with gene density over 1Mbp non-

1300 overlapping windows. Linear regression lines are colored according to species: (a) *P.*

1301 *tremula* (orange line), (b) *P. tremuloides* (blue line) and (c) *P. trichocarpa* (green

1302 line).

1303

1304 **S1 Table. Samples used in this study.**

1305

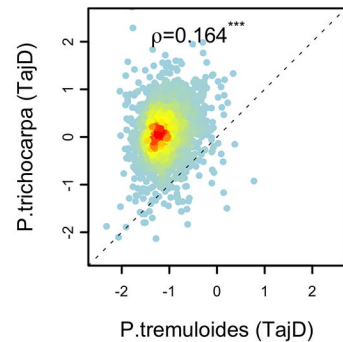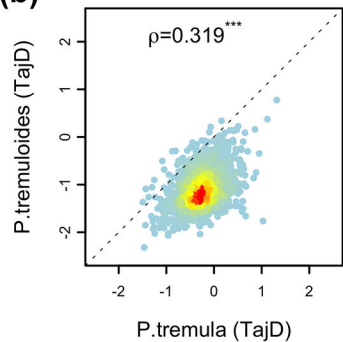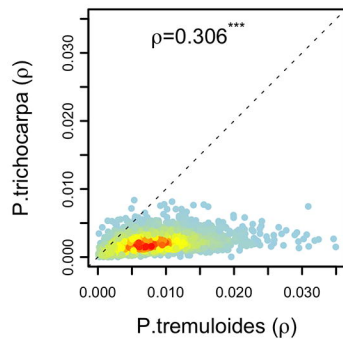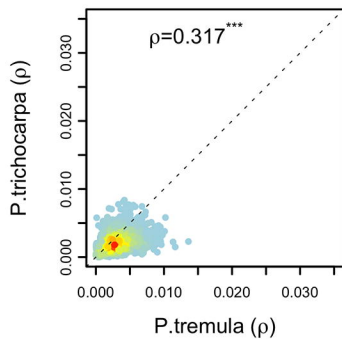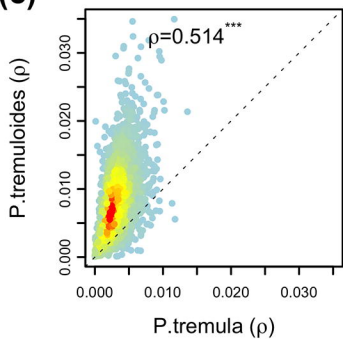1306 **S2 Table. Summary statistics of Illumina re-sequencing data per sample.**

1307

1308 **S3 Table. Pairwise divergence ($d_{xy}$) (median and central 95% range) between *P.***

1309 ***tremula*, *P. tremuloides* and P. *trichocarpa* for various genomic contexts over 100**

1310 **Kbp non-overlapping windows across genomes**.
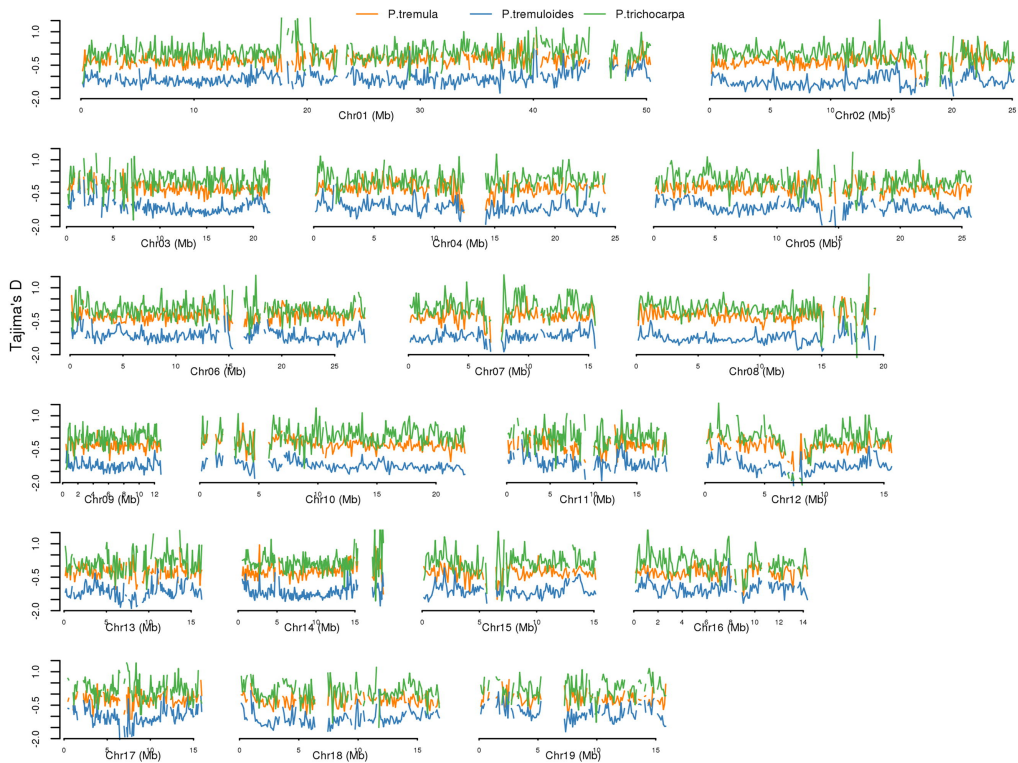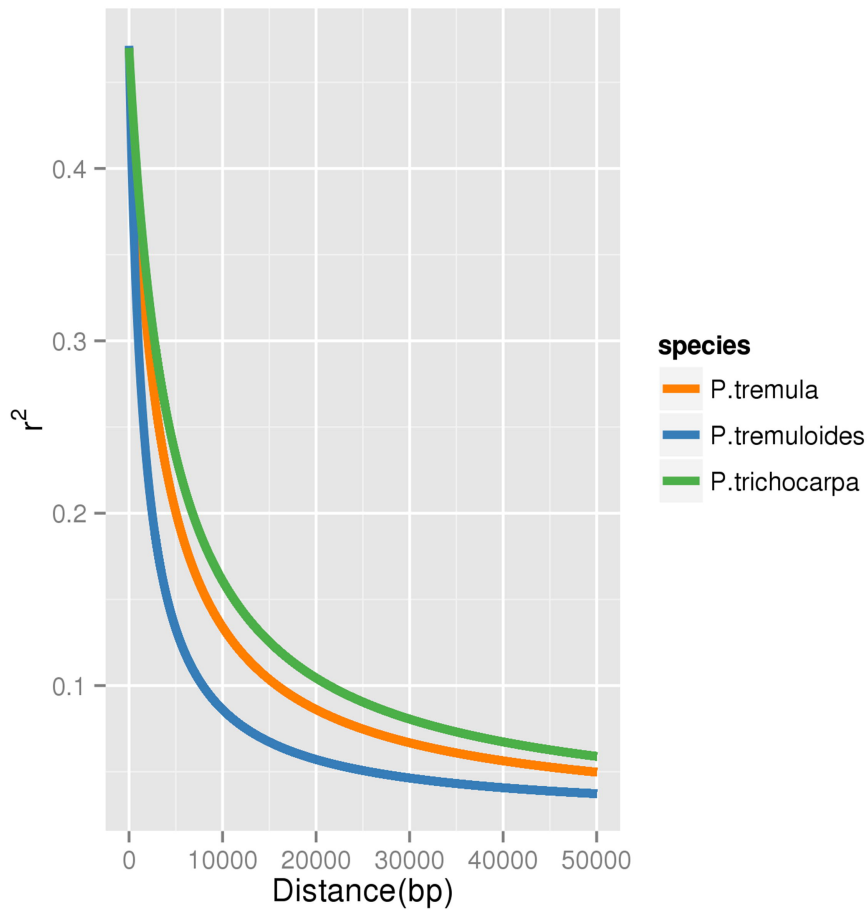
1311

54

1312 **S4 Table. Summary of the correlation coefficients (Spearman's $\rho$) between levels**

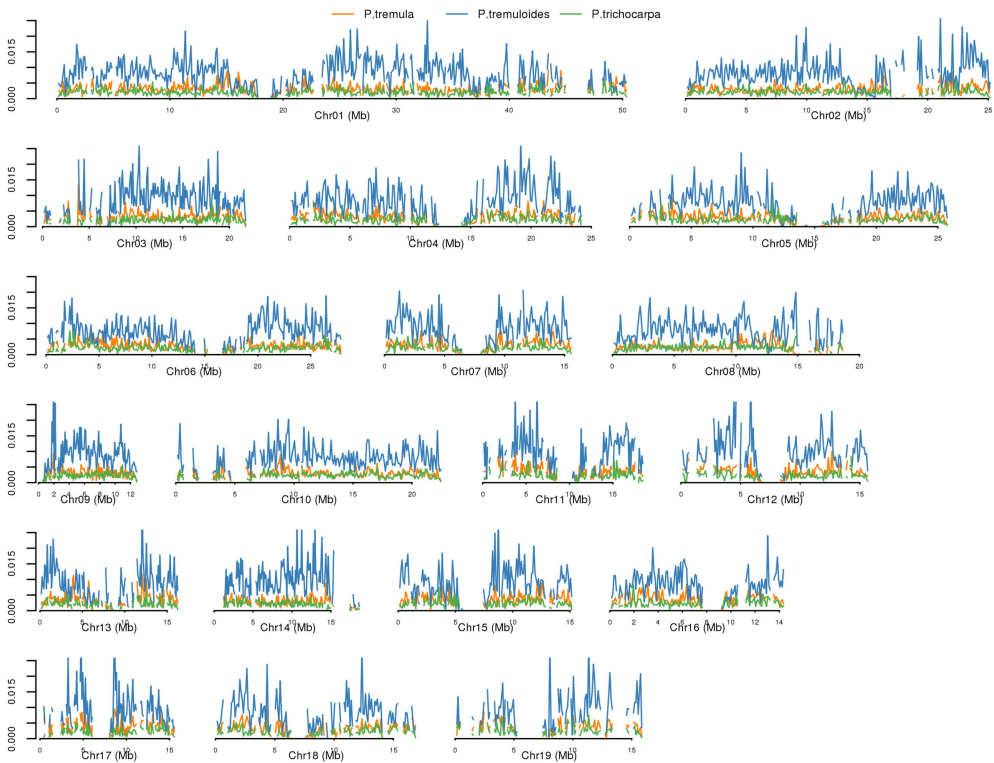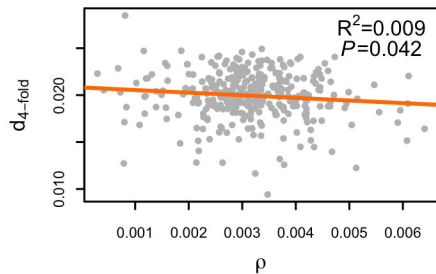1313 **of synonymous diversity and non-synonymous divergence.**

Legend: P.tremula (orange), P.tremuloides (blue), P.trichocarpa (green)

Y-axis: $\theta_\pi$

Chr01 (Mb), Chr02 (Mb), Chr03 (Mb), Chr04 (Mb), Chr05 (Mb), Chr06 (Mb), Chr07 (Mb), Chr08 (Mb), Chr09 (Mb), Chr10 (Mb), Chr11 (Mb), Chr12 (Mb), Chr13 (Mb), Chr14 (Mb), Chr15 (Mb), Chr16 (Mb), Chr17 (Mb), Chr18 (Mb), Chr19 (Mb)

(a)

R²=0.458
P<0.001

$\theta_{4-\text{fold}}$

*P. tremula*

$\rho$

R²=0.009
P=0.042

$d_{4-\text{fold}}$

$\rho$

(b)

R²=0.213
P<0.001

$\theta_{4-\text{fold}}$

*P. tremuloides*

$\rho$

R²=-0.002
P=0.695

$d_{4-\text{fold}}$

$\rho$

(c)

R²=0.039
P<0.001

$\theta_{4-\text{fold}}$

*P. trichocarpa*

$\rho$

R²=-0.001
P=0.431

$d_{4-\text{fold}}$

$\rho$

Figure showing three columns (*P. tremula*, *P. tremuloides*, *P. trichocarpa*) and three rows of scatter plots with trend lines. Row (a) shows $\rho$, row (b) shows $\theta_{4\text{-fold}}$, and row (c) shows $\theta_{\text{Intergenic}}$, each plotted against Gene number/Mbp.