

# Natural Selection and Recombination Rate Variation Shape Nucleotide Polymorphism Across the Genomes of Three Related *Populus* Species

Jing Wang,\* Nathaniel R. Street,<sup>†</sup> Douglas G. Scofield,<sup>\*,\*,§</sup> and Pär K. Ingvarsson<sup>\*,1</sup>

\*Department of Ecology and Environmental Science, and <sup>†</sup>Department of Plant Physiology, Umeå Plant Science Centre, Umeå University, Umeå SE 90187, Sweden, <sup>‡</sup>Department of Ecology and Genetics: Evolutionary Biology, and <sup>§</sup>Uppsala Multidisciplinary Center for Advanced Computational Science, Uppsala University, Uppsala SE 75105, Sweden

ORCID IDs: 000-0002-3793-3264 (J.W.); 0000-0001-6031-005X (N.R.S.); 0000-0001-5235-6461 (D.G.S.); 0000-0001-9225-7521 (P.K.I.)

**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 use whole-genome resequencing data to characterize and compare genome-wide patterns of nucleotide polymorphism, site frequency spectrum, and population-scaled recombination rates in three species of *Populus*: *Populus tremula*, *P. tremuloides*, and *P. trichocarpa*. We find 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. Our findings highlight multiple lines of evidence suggesting that natural selection, due to both purifying and positive selection, has widely shaped patterns of nucleotide polymorphism at linked neutral sites in all three species. Differences in effective population sizes and rates of recombination largely explain the disparate magnitudes and signatures of linked selection that we observe among species. The present work provides the first phylogenetic comparative study on a genome-wide scale in forest trees. This information will also improve our ability to understand how various evolutionary forces have interacted to influence genome evolution among related species.

**KEYWORDS** *Populus*; whole-genome resequencing; nucleotide polymorphism; recombination; natural selection

A major goal in evolutionary genetics is to understand how genomic variation is established and maintained within and between species (Nordborg *et al.* 2005; Begun *et al.* 2007) and how different evolutionary forces have substantial impacts in shaping genetic variation throughout the genome (Hellmann *et al.* 2005). Under the neutral theory, genetic variation is the manifestation of the balance between mutation and genetic drift (Kimura 1983). Demographic fluctuations, such as population expansion and/or bottlenecks, can cause patterns of genome-wide variation deviating from the standard neutral model in various ways (Li and Durbin 2011). It is now clear, however, that natural selection—via positive selection favoring beneficial mutations (genetic hitchhiking) and/or purifying selection against deleterious

mutations (background selection)—plays an important role in molding the landscape of nucleotide polymorphism in many species (Begun and Aquadro 1992; Begun *et al.* 2007; Cutter and Choi 2010; Mackay *et al.* 2012).

If natural selection is pervasive across the genome, patterns of genetic variation at linked neutral sites can be influenced by selection in a number of ways. First, positive correlations between levels of neutral polymorphism and recombination rates are expected since linked selection is expected to remove more neutral polymorphism in low-recombination regions compared to high-recombination regions and such a pattern is unlikely to be generated by demographic processes alone (Begun and Aquadro 1992; Kulathinal *et al.* 2008; McGaugh *et al.* 2012; Campos *et al.* 2014; Charlesworth and Campos 2014). Second, in addition to influencing the level of neutral variability, recombination rate can affect the efficacy of selection through the process known as Hill–Robertson interference (HBI) (Hill and Robertson 1966). If HRI is operating, genetic linkage effects in regions of low recombination will reduce the local effective population size ( $N_e$ ) and accordingly reduce the efficacy of selection ( $N_e s$ ), since the effects of

Copyright © 2016 by the Genetics Society of America  
doi: 10.1534/genetics.115.183152

Manuscript received September 25, 2015; accepted for publication December 24, 2015; published Early Online December 30, 2015.

Available freely online through the author-supported open access option.

Supporting information is available online at [www.genetics.org/lookup/suppl/doi:10.1534/genetics.115.183152/-/DC1](http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.115.183152/-/DC1).

<sup>1</sup>Corresponding author: Department of Ecology and Environmental Science, Umeå University, Umeå SE 90187, Sweden. E-mail: [par.ingvarsson@umu.se](mailto:par.ingvarsson@umu.se)

selection are determined by the product of  $N_e$  and the selection coefficient on a mutation ( $s$ ) (Kimura 1983). We would therefore expect both a reduced fixation of favorable mutations and an increased frequency of deleterious mutations in these regions (Hill and Robertson 1966; Haddrill *et al.* 2007; Campos *et al.* 2014). Third, signatures and magnitudes of linked selection are sensitive to the density of important functional sites (*e.g.*, gene density) within specific genomic regions (Flowers *et al.* 2012). In accordance with the view that genes represent the most likely targets of natural selection, regions with a high density of genes are expected to have undergone stronger effects of linked selection and exhibit lower levels of neutral polymorphism (Nordborg *et al.* 2005; Flowers *et al.* 2012). Therefore, a positive or negative covariation of recombination rate and gene density would act to either obscure or strengthen the signatures of linked selection across the genome (Cutter and Payseur 2003; Cutter and Choi 2010; Flowers *et al.* 2012). Finally, a distinctive signature of recurrent selective sweeps is the local reduction of linked neutral polymorphism in regions experiencing frequent adaptive substitutions (Andolfatto 2007). A substantial number of adaptive substitutions are likely composed of amino acid substitutions, and a negative correlation between neutral polymorphism and nonsynonymous divergence can thus be particularly informative of the prevalence of selective sweeps (Macpherson *et al.* 2007). With the advance of next-generation sequencing technology, sufficient genome-wide data among multiple related species are becoming available (Luikart *et al.* 2003; Ellegren 2014). Phylogenetic comparative approaches will thus place us in a stronger position to understand how various evolutionary forces have interacted to shape the heterogeneous patterns of nucleotide polymorphism across the genome (Hufford *et al.* 2012; Cutter and Payseur 2013; Lawrie and Petrov 2014).

Thus far, genome-wide comparative studies have largely dealt with experimental model species, mammals, and cultivated plants of either agricultural or horticultural interest (Locke *et al.* 2011; Hufford *et al.* 2012; Liu *et al.* 2014). Forest trees, as a group, are characterized by extensive geographical distributions and are of high ecological and economic value (Neale and Kremer 2011). Most forest trees have largely persisted in an undomesticated state and, until quite recently, without anthropogenic influence (Neale and Kremer 2011). Accordingly, in contrast to crop and livestock lineages that have been through strong domestication bottlenecks, most extant populations of forest trees harbor a wealth of genetic variation, and they are thus excellent model systems for dissecting the dominant evolutionary forces that sculpt patterns of variation throughout the genome (González-Martínez *et al.* 2006; Neale and Kremer 2011). Among forest tree species, the genus *Populus* represents a particularly attractive choice because of its wide geographic distribution, important ecological role in a wide variety of habitats, multiple economic uses in wood and energy products, and relatively small genome size (Eckenwalder 1996; Jansson and Douglas 2007). Here, we studied three *Populus* species that differ in

morphology, geographic distribution, population size, and phylogenetic relationship (Supporting Information, Figure S1) (Jansson *et al.* 2010; Wang *et al.* 2014). *P. tremula* and *P. tremuloides* (collectively, “aspens”) have wide native ranges across Eurasia and North America, respectively; are closely related; and belong to the same section of the genus (section *Populus*) (Jansson *et al.* 2010). In contrast, *P. trichocarpa* belongs to a different section of the genus (section *Tacamahaca*) that is reproductively isolated from members of the *Populus* section (Jansson *et al.* 2010). The distribution of *P. trichocarpa* is restricted to western regions of North America, and its distribution range is considerably smaller than the two aspen species (Dickmann and Kuzovkina 2014). Importantly, *P. trichocarpa* also represents the first tree species to have its genome published (Tuskan *et al.* 2006), and the genome sequence and annotation have undergone continual improvement (<http://phytozome.jgi.doe.gov>). This enables us to provide important context for our genome comparisons. The phylogenetic relationship of the three species [*(P. tremula–P. tremuloides) P. trichocarpa*] is well established by both chloroplast and nuclear DNA sequences (Hamzeh and Dayanandan 2004; Wang *et al.* 2014).

In this study, we used data sets generated by Next-Generation Sequencing (NGS) to characterize, compare, and contrast genome-wide patterns of nucleotide diversity, site frequency spectrum, and recombination rate and to infer contextual patterns of selection throughout the genomes for all three species.

## Materials and Methods

### Samples and sequencing

Leaf samples were collected from 24 genotypes of *P. tremula* and 24 genotypes of *P. tremuloides* (Table S1). Genomic DNA was extracted from leaf samples, and paired-end sequencing libraries with insert sizes of 650 bp were constructed for all genotypes. Whole-genome sequencing with a minimum expected depth of 20× was performed on the Illumina HiSeq-2000 platform at the Science for Life Laboratory, Stockholm, and 2× 100-bp paired-end reads were generated for all genotypes. Two samples of *P. tremuloides* failed to yield the expected coverage and were therefore removed from subsequent analyses. We obtained publicly available short-read Illumina data of 24 *P. trichocarpa* individuals from the National Center for Biotechnology Information (NCBI) Short Read Archive (SRA) (Table S1). Individuals were selected to have a similar read depth as the samples of the two aspen species. The accession numbers of *P. trichocarpa* samples can be found in Evans *et al.* (2014). All analyses are thus based on data from 24 *P. tremula*, 22 *P. tremuloides*, and 24 *P. trichocarpa* genotypes.

### Raw read filtering, read alignment, and postprocessing alignment

Prior to read alignment, we used Trimmomatic (Lohse *et al.* 2012) to remove adapter sequences from reads. Since the quality of reads always drops toward the end of reads, we

used Trimmomatic to cut off bases from the start and/or end of reads when the quality values were  $<20$ . If the length of the processed reads was reduced to  $<36$  bases after trimming, reads were completely discarded. FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>) was used to check and compare the per-base sequence quality between the raw sequence data and the filtered data. After quality control, all paired-end and orphaned single-end reads from each sample were mapped to the *P. trichocarpa* version 3 (v3.0) genome (Tuskan *et al.* 2006) using BWA-MEM with default parameters in bwa-0.7.10 (Li 2013).

Several postprocessing steps of alignments were performed to minimize the number of artifacts in downstream analysis: First, we performed insertions and deletion (indel) realignment since mismatching bases are usually found in regions with indels (Wang *et al.* 2015). The RealignerTarget-Creator in The Genome Analysis Toolkit (GATK) (DePristo *et al.* 2011) was first used to find suspicious-looking intervals that were likely in need of realignment. Then the IndelRealigner was used to run the realigner over those intervals. Second, as reads resulting from PCR duplicates can arise during the sequencing library preparation, we used the MarkDuplicates methods in the Picard package (<http://broadinstitute.github.io/picard/>) to remove those reads or read pairs having identical external coordinates and the same insert length. In such cases, only the single read with the highest summed base qualities was kept for downstream analysis. Third, to exclude genotyping errors caused by paralogous or repetitive DNA sequences where reads were poorly mapped to the reference genome or by other genome feature differences between *P. trichocarpa* and *P. tremula* or *P. tremuloides*, we removed sites with extremely low- and extremely high-read depths after investigating the empirical distribution of read coverage. We filtered out sites with a total coverage  $<100\times$  or  $>1200\times$  across all samples per species. When reads were mapped to multiple locations in the genome, they were randomly assigned to one location with a mapping score of zero by BWA-MEM. To account for such misalignment effects, we removed those sites if there were  $>20$  mapped reads with a mapping score equaling to zero across all individuals in each species. Finally, because the short-read alignment is generally unreliable in highly repetitive genomic regions, we filtered out sites that overlapped with known repeat elements as identified by RepeatMasker (Tarailo-Graovac and Chen 2009). In the end, the subset of sites that passed all these filtering criteria in the three *Populus* species was used in downstream analyses.

### Single nucleotide polymorphism and genotype calling

We implemented two complementary bioinformatics approaches: First, many studies have pointed out the bias inherent in population genetic estimates using genotype calling approach from NGS data (Nielsen *et al.* 2011; Nevado *et al.* 2014). Single- or multiple-sample genotype calling can result in a bias in the estimation of site frequency spectrum (SFS), as the former usually leads to overestimation of rare variants,

whereas the latter often leads to the opposite (Nielsen *et al.* 2011). Therefore, in this study we employed a method implemented in the software package—Analysis of Next-Generation Sequencing Data (ANGSD v0.602) (Korneliussen *et al.* 2014)—to estimate the SFS and all population genetic statistics derived from the SFS without calling genotypes. Second, for those analyses that require accurate single nucleotide polymorphism (SNP) and genotype calls, we performed SNP calling with HaplotypeCaller of the GATK v3.2.2 (DePristo *et al.* 2011), which called SNPs and indels simultaneously via local reassembly of haplotypes for each individual and created single-sample genomic VCFs (gVCFs). gVCFs in GATK were then used to merge multi-sample records together, correct genotype likelihoods, re-genotype the newly merged record, and perform re-annotation. The following filtering steps were then used to reduce the number of false-positive SNPs and retain high-quality SNPs: (1) We removed all SNPs that overlapped with sites excluded by all previous filtering criteria. (2) We retained only biallelic SNPs with a distance of  $>5$  bp away from any indels. (3) We treated genotypes with a genotype quality score (GQ)  $<10$  as missing and then removed those SNPs with a genotype missing rate  $>20\%$ . (4) We removed SNPs that showed significant deviation from Hardy–Weinberg Equilibrium ( $P < 0.001$ ). After all filtering, 8,502,169 SNPs were detected among the three *Populus* species and were used in downstream analyses.

### Population structure

We used fourfold synonymous SNPs with minor allele frequency  $>0.1$  to perform population structure analyses with ADMIXTURE (Alexander *et al.* 2009). We ran ADMIXTURE on all the sampled individuals among species and on the samples within each species separately. The number of genetic clusters ( $K$ ) varied from 1 to 6. The most likely number of genetic clusters was selected by minimizing the cross-validation error in ADMIXTURE.

### Diversity and divergence: related summary statistics

For nucleotide diversity and divergence estimates, only the reads with mapping quality  $>30$  and the bases with a quality score  $>20$  were used in all downstream analyses with ANGSD (Korneliussen *et al.* 2014). First, we used the -doSaf implementation in ANGSD to calculate the site-allele-frequency likelihood based on the SAMTools genotype likelihood model (Li *et al.* 2009). Then, we used the -realSFS implementation in ANGSD to obtain an optimized folded global SFS using the expectation maximization algorithm for each species. Based on the global SFS, we used the -doThetas function in ANGSD to estimate the per-site nucleotide diversity from posterior probability of allele frequency based on a maximum-likelihood approach (Kim *et al.* 2011). Two standard estimates of nucleotide diversity, the average pairwise nucleotide diversity ( $\Theta_{\pi}$ ) (Tajima 1989) and the proportion of segregating sites ( $\Theta_W$ ) (Watterson 1975), and one neutrality statistic test, Tajima's  $D$  (Tajima 1989),

were summarized along all 19 chromosomes using nonoverlapping sliding windows of 100 kilobases (kbp) and 1 megabases (Mb). Windows with <10% of covered sites left from previous quality filtering steps were excluded. In the end, 3340 100-kbp and 343 1-Mb windows, with an average of 50,538 and 455,910 covered bases per window, respectively, were included.

All these statistics were also calculated for each type of functional element (0-fold nonsynonymous, fourfold synonymous, intron, 3' UTR, 5' UTR, and intergenic sites) over nonoverlapping 100-kbp and 1-Mb windows in all three *Populus* species. The category of gene models followed the gene annotation of *P. trichocarpa* version 3.0 (Tuskan *et al.* 2006). For protein-coding genes, we included only genes with at least 90% of covered sites left from previous filtering steps to ensure that the three species have the same gene structures. For regions overlapped by different transcripts in each gene, we classified each site according to the following hierarchy (from highest to lowest): coding regions (CDS), 3' UTR, 5' UTR, and intron. Thus, if a site resides in a 3' UTR in one transcript and in a CDS for another, the site was classified as CDS. We used the transcript with the highest content of protein-coding sites to categorize synonymous and nonsynonymous sites within each gene. A total of 16.52, 3.4, 7.19, 4.02, 31.89, and 73.46 Mb was partitioned into 0-fold nonsynonymous (where all DNA sequence changes lead to protein sequence changes), fourfold synonymous (where all DNA sequence changes lead to the same protein sequences), 3' UTR, 5' UTR, intron, and intergenic categories.

#### **Linkage disequilibrium and population-scaled recombination rate**

A total of 1,409,377 SNPs, 1,263,661 SNPs, and 710,332 SNPs with minor allele frequency >10% were used for the analysis of linkage disequilibrium (LD) and population-scaled recombination rate ( $\rho$ ) in *P. tremula*, *P. tremuloides*, and *P. trichocarpa*, respectively. To estimate and compare the rate of LD decay among the three *Populus* species, we first used PLINK 1.9 (Purcell *et al.* 2007) to randomly thin the number of SNPs to 100,000 in each species. We then calculated the squared correlation coefficients ( $r^2$ ) between all pairs of SNPs that were within a distance of 50 kbp using PLINK 1.9. The decay of LD against physical distance was estimated using nonlinear regression of pairwise  $r^2$  vs. the physical distance between sites in base pairs (Remington *et al.* 2001).

We estimated the population-scaled recombination rate  $\rho$  using the Interval program of LDhat 2.2 (McVean *et al.* 2004) with 1,000,000 MCMC iterations sampling every 2000 iterations and a block penalty parameter of 5. The first 100,000 iterations of the MCMC iterations were discarded as a burn-in. We then calculated the scaled value of  $\rho$  in each 100-kbp and 1-Mb window by averaging over all SNPs in that window. Only windows with >10,000 (in 100-kbp windows) and 100,000 sites (in 1-Mb windows) and 100 SNPs left from previous filtering steps were used for the estimation of  $\rho$ .

#### **Estimating the distribution of fitness effects of new amino acid mutations and the proportion of adaptive amino acid substitutions**

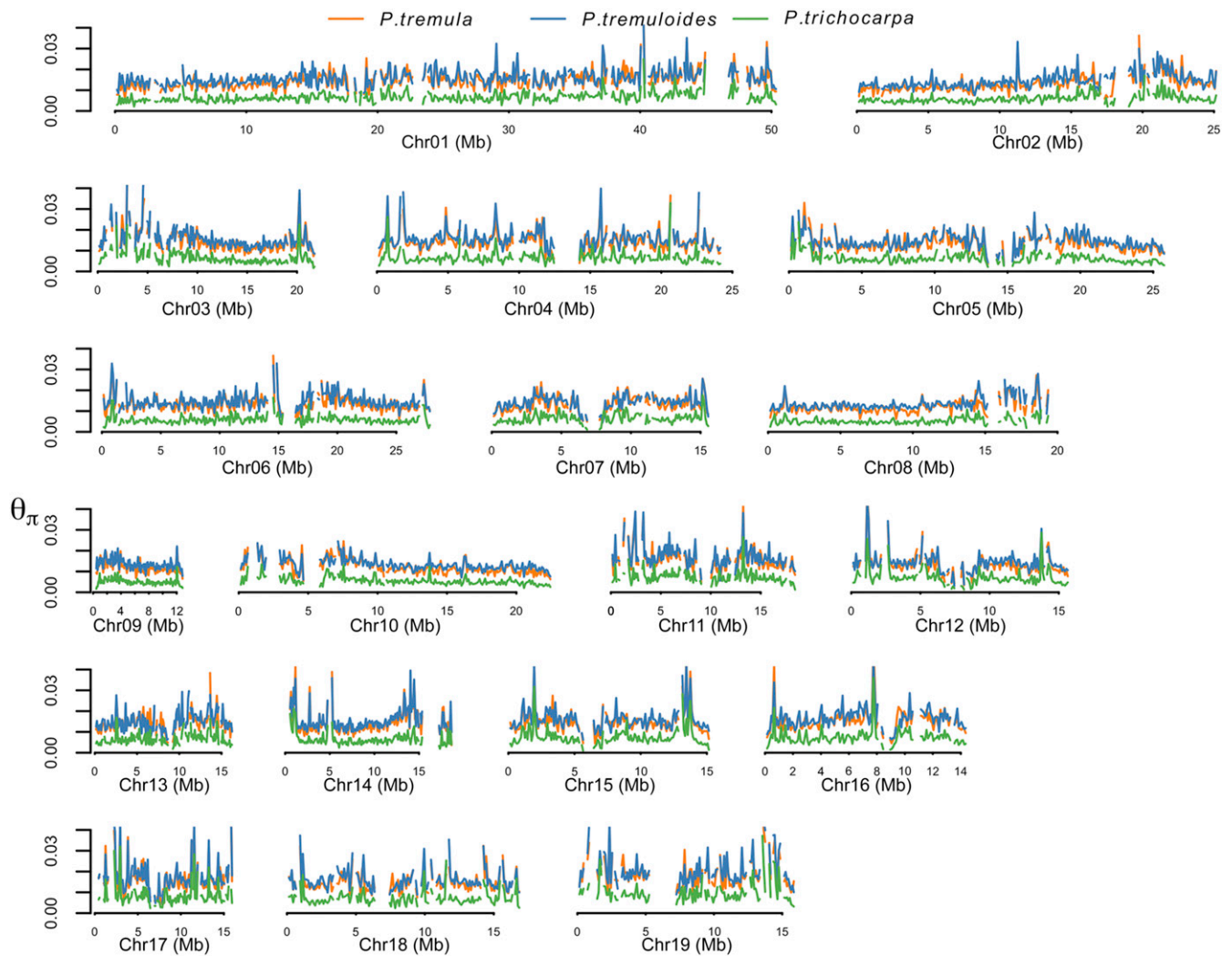
We generated the folded SFS in each species for a class of selected sites (0-fold nonsynonymous sites) and a class of putatively neutral reference sites (fourfold synonymous sites) from SNP data using a custom Perl script. We employed a maximum-likelihood (ML) approach as implemented in the program distribution of fitness effects (DFE)- $\alpha$  (Keightley and Eyre-Walker 2007; Eyre-Walker and Keightley 2009) to fit a demographic model with a step of population size change to the neutral SFS. Fitness effects of new deleterious mutations at the selected site class were sampled from a gamma distribution after incorporating the estimated parameters for the demographic model. This method assumes that fitness effects of new mutations at neutral sites are zero and unconditionally deleterious at selected sites since it assumes that advantageous mutations are too rare to contribute to polymorphism (Keightley and Eyre-Walker 2007). We report the proportion of amino acid mutations falling into different effective strengths of selection ( $N_e s$ ) range: 0–1, 1–10, and >10, respectively.

From the estimated DFE, the proportion of adaptive amino acid substitutions ( $\alpha$ ) and the relative rate ( $\omega$ ) of adaptive substitution at 0-fold nonsynonymous sites were estimated using the method of Eyre-Walker and Keightley (2009). This method explicitly accounts for past changes in population size and the presence of slightly deleterious mutations. Among the total of 8,502,169 SNPs detected by GATK, on average <1% were shared between either of the two aspen species and *P. trichocarpa* (Figure S2). We therefore used the aspen species and *P. trichocarpa* as each other's outgroup species to calculate between-species nucleotide divergence at fourfold synonymous and 0-fold nonsynonymous sites since it is unlikely to be influenced by shared ancestral polymorphisms. Jukes–Cantor multiple hits correction was applied to the divergence estimates (Jukes and Cantor 1969). For the parameters of  $N_e s$ ,  $\alpha$ , and  $\omega$ , we generated 200 bootstrap replicates by resampling randomly across all SNPs in each site class using R (R Development Core Team 2014). We excluded the top and bottom 2.5% of bootstrap replicates and used the remainder to represent the 95% confidence intervals for each parameter.

#### **Genomic correlates of diversity**

To examine the factors influencing levels of neutral polymorphism in all three *Populus* species, we first assumed that the fourfold synonymous sites in genic regions were selectively neutral as every possible mutation at fourfold degenerate sites is synonymous. In the following we refer to the pairwise nucleotide diversity at fourfold synonymous sites ( $\theta_{\text{fourfold}}$ ) as “neutral polymorphism.” As a comparison to genic region, we also estimated levels of nucleotide diversity at intergenic sites ( $\theta_{\text{Intergenic}}$ ). Then we tabulated several other genomic features within each 100-kbp and 1-Mb window that may correlate with patterns of polymorphism. First, we summarized population-scaled recombination rate ( $\rho$ ) as





**Figure 1** Genome-wide patterns of polymorphism among three *Populus* species. Nucleotide diversity ( $\theta_{\pi}$ ) was calculated over 100-kbp nonoverlapping windows in *P. tremula* (orange line), *P. tremuloides* (blue line), and *P. trichocarpa* (green line) along the 19 chromosomes.

described above for each species. Second, we tabulated GC content as the fraction of bases where the reference sequence (*P. trichocarpa* v3.0) was a G or a C. Third, we measured the gene density as the number of functional genes within each window according to the gene annotation of *P. trichocarpa* version 3.0. Any portion of a gene that fell within a window was counted as a full gene. Fourth, we accounted for the variation of mutation rate by calculating the number of fixed differences per neutral site (either four-fold synonymous sites or intergenic sites) between aspen and *P. trichocarpa* within each window, which was performed in the ngsTools (Fumagalli *et al.* 2014). The reason why we used divergence between aspen and *P. trichocarpa* to measure mutation rate is because they are distantly related (Wang *et al.* 2014) and the estimate of divergence is unlikely to be influenced by shared ancestral polymorphisms between species as shown above. Fifth, we tabulated the number of covered bases in each window as those left from all previous filtering criteria.

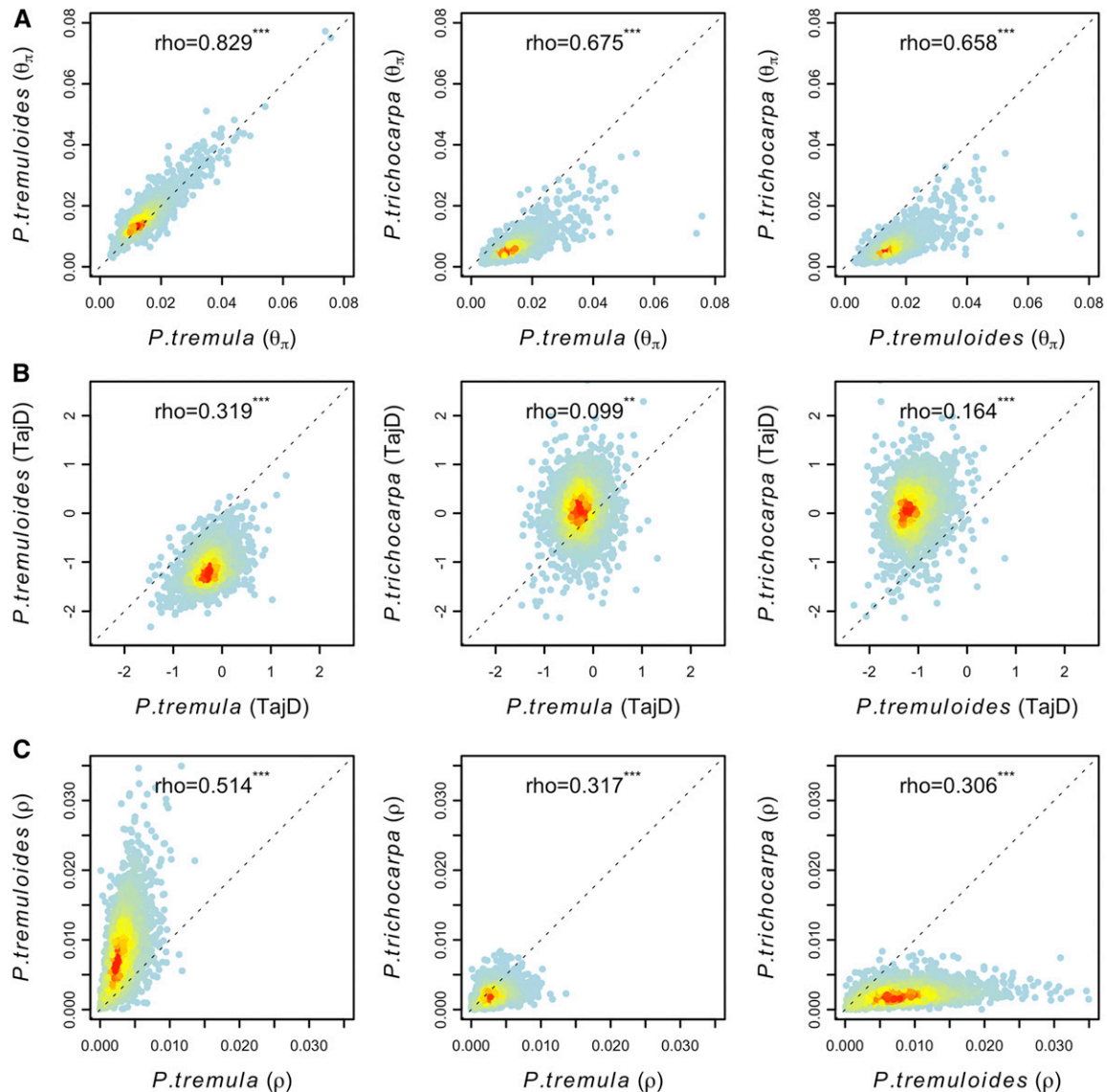
We used Spearman's rank-order correlation test to examine pairwise correlations between the variables of interest. To account for the autocorrelation between variables, we further calculated partial correlations between the variables of interest by removing the confounding effects of other variables (Kim and Soojin 2007). All statistical tests were performed using R version 3.2.0 unless stated otherwise.

#### Data availability

All newly generated Illumina reads of 24 *P. tremula* and 22 *P. tremuloides* from this study have been submitted to the SRA at NCBI. All accession numbers can be found in Table S1.

#### Results

We generated whole-genome sequencing data for 24 *P. tremula* and 22 *P. tremuloides* (Table S1) with all samples sequenced to relatively high depth (24.2×–69.2×; Table S2). We also downloaded whole-genome resequencing data for



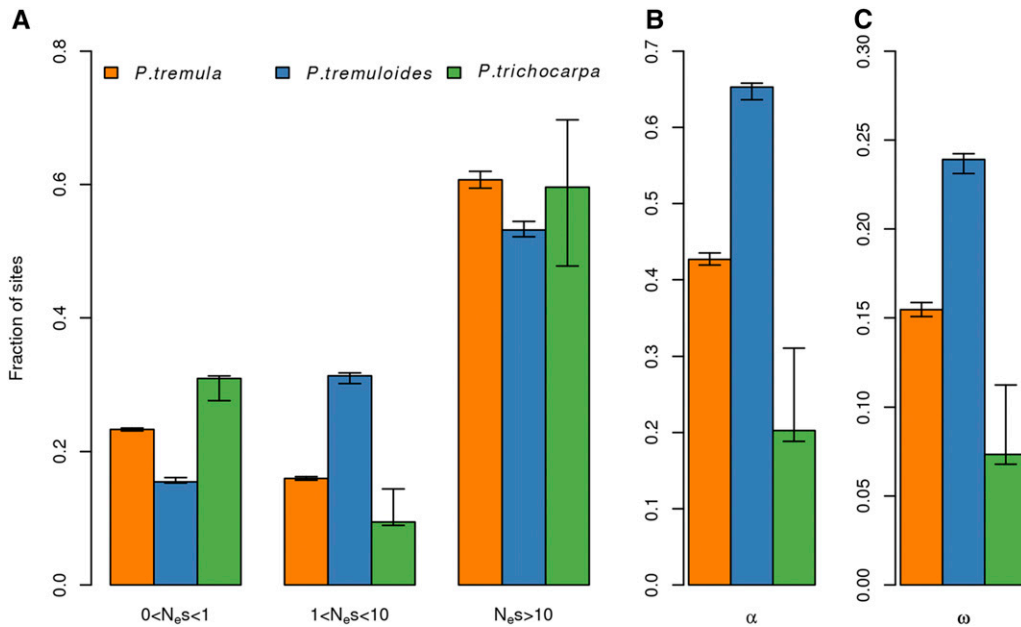
**Figure 2** Distribution and correlations of (A) polymorphism ( $\theta_{\pi}$ ), (B) Tajima's  $D$ , and (C) population-scaled recombination rate ( $\rho$ ) between pairwise comparisons of *P. tremula*, *P. tremuloides*, and *P. trichocarpa* over 100-kbp nonoverlapping windows. The red-to-yellow-to-blue gradient indicates decreased density of observed events at a given location in the graph. Spearman's rank correlation coefficient ( $\rho$ ) and the  $P$ -value are shown in each subplot. (\*\*\* $P < 2.2 \times 10^{-16}$ , \*\* $P < 0.001$ ). The dashed gray line in each subplot indicates a simple linear regression line with intercept being zero and slope being one.

24 samples of *P. trichocarpa* from the NCBI SRA (Evans *et al.* 2014). After adapter removal and quality trimming, 949.2 Gb of high quality sequence data remained (Figure S3; Table S2). The mean mapping rate of reads to the *P. trichocarpa* reference genome was 89.8% for *P. tremula*, 91.1% for *P. tremuloides*, and 95.2% for *P. trichocarpa* (Table S2). On average, the genome-wide coverage of uniquely mapped reads was  $>20\times$  for each species (Table S2). After excluding sites with extreme coverage, low mapping quality, or those overlapping with annotated repetitive elements (see *Materials and Methods*), 42.8% of collinear genomic sequences remained for downstream analyses. Of these sites, 54.9% were found within gene boundaries, covering 70.1% of all genic regions predicted from the *P. trichocarpa* assembly.

The remaining sites (45.1%) were located in intergenic regions.

#### Genome-wide patterns of polymorphism, site frequency spectrum, and recombination among the three *Populus* species

When analyzing population structure between species, we found that the model exhibited the lowest cross-validation error when the number of ancestral populations ( $K = 3$ ) (Figure S4B), which clearly subdivided the three species into three distinct clusters (Figure S4A). When we analyzed population structure within each species separately, we found that the cross-validation error increased linearly with increasing  $K$ , with  $K = 1$  minimizing the cross-validation error for all



**Figure 3** Estimates of purifying and positive selection at 0-fold nonsynonymous sites in three *Populus* species. (A) The distribution of fitness effects of new amino acid mutations, (B) the proportion of adaptive substitution ( $\alpha$ ), and (C) the rate of adaptive nonsynonymous-to-synonymous substitutions ( $\omega$ ) in *P. tremula* (orange bar), *P. tremuloides* (blue bar), and *P. trichocarpa* (green bar). Error bars represent 95% bootstrap confidence intervals.

three species (Figure S4, C–E). Our results are slightly different from several earlier studies that have documented population subdivision in these species (De Carvalho *et al.* 2010; Evans *et al.* 2014), but it is likely due to the small sample sizes used in this study (22–24 individuals). In addition, it could also be caused by the low power of model-based approaches in detecting population structure when it is very weak (Alexander *et al.* 2009). More generally, population structure is expected to be weak in *Populus*, given the great dispersal capabilities of both pollen and seeds (Eckenwalder 1996; Jansson and Douglas 2007).

The two aspen species harbor substantial levels of nucleotide diversity across the genome ( $\Theta_{\Pi} = 0.0133$  in *P. tremula*;  $\Theta_{\Pi} = 0.0144$  in *P. tremuloides*), approximately two- to threefold higher than diversity in *P. trichocarpa* ( $\Theta_{\Pi} = 0.0059$ ) (Figure 1; Table S3). Among various genomic contexts, we found that the levels of nucleotide diversity were highest at intergenic sites, followed by fourfold synonymous sites, 3' UTRs, 5' UTRs, and introns and were lowest at 0-fold nonsynonymous sites (Figure S5; Table S3). In accordance with the view that the large majority of amino acid mutations are selected against (Larracuent *et al.* 2008), we found significantly lower Tajima's  $D$  at 0-fold nonsynonymous sites compared to fourfold synonymous sites ( $P < 0.001$ , Mann–Whitney  $U$ -test) (Figure S6; Table S3). In addition, we observed significantly positive correlations of  $\Theta_{\Pi}$  between each pair of the three species across the whole genome (Figure 2A). The overall nucleotide diversity estimated in *P. trichocarpa* was slightly higher than the value reported in Evans *et al.* (2014) ( $\Theta_{\Pi} = 0.0041$ ), but this likely only reflects differences between the methods used in the two studies. In this study, we utilized the full information of the filtered data and estimated the population genetics statistics directly from genotype likelihoods, which take

statistical uncertainty of SNP and genotype calling into account and should give more accurate estimates (Kim *et al.* 2011; Nielsen *et al.* 2011).

Compared to patterns of polymorphism, we observed much weaker correlations of the site frequency spectrum, summarized using the Tajima's  $D$  statistic (Tajima 1989), between species (Figure 2B). *P. tremuloides* (average Tajima's  $D = -1.169$ ) showed substantially greater negative values of Tajima's  $D$  along all chromosomes compared to both *P. trichocarpa* (average Tajima's  $D = 0.064$ ) and *P. tremula* (average Tajima's  $D = -0.272$ ) (Figure S7; Table S3), reflecting a large excess of low-frequency polymorphisms segregating in this species. Furthermore, the three *Populus* species showed different extents of genome-wide LD decay (Figure S8), with LD decaying fastest in *P. tremuloides* and slowest in *P. trichocarpa* (Figure S8). This reflects the rank order of their population-scaled recombination rates ( $\rho = 4N_e c$ ) (Figure S9), for which the mean  $\rho$  over 100-kbp nonoverlapping windows was highest in *P. tremuloides* ( $8.42 \text{ kbp}^{-1}$ ), followed by *P. tremula* ( $3.23 \text{ kbp}^{-1}$ ), and lowest in *P. trichocarpa* ( $2.19 \text{ kbp}^{-1}$ ). Intermediate correlations of recombination rates were observed between species (Figure 2C). In addition, concordant values of  $\Theta_{\Pi}$ , Tajima's  $D$  and  $\rho$  for all three species were also observed in 1-Mb windows (Figure S10). For populations under drift-mutation-recombination equilibrium,  $\rho = 4N_e c$  (where  $N_e$  is the effective population size and  $c$  is the recombination rate) and  $\theta_W = 4N_e \mu$  (where  $N_e$  is the effective population size and  $\mu$  is the mutation rate). To compare the relative contribution of recombination ( $c$ ) and mutation ( $\mu$ ) in shaping genomic variation, we measured the ratio of population recombination rate to the nucleotide diversity ( $\rho/\theta_W$ ) across the genome (Figure S11). The mean  $c/\mu$  in *P. tremula*, *P. tremuloides*, and *P. trichocarpa* was 0.22, 0.39, and 0.38, respectively.

**Table 1 Summary of the correlation coefficients (Spearman's rank correlation coefficient) between levels of neutral polymorphism ( $\theta$ ), divergence ( $d$ ), and recombination rate ( $\rho$ ) in genic and intergenic regions among all three *Populus* species**

Data set	Species	$\rho$ vs. $\theta_{\text{fourfold}}$		$\rho$ vs. $d_{\text{fourfold}}$	$\rho$ vs. $\theta_{\text{intergenic}}$		$\rho$ vs. $d_{\text{intergenic}}$
		Pairwise	Partial <sup>a</sup>		Pairwise	Partial <sup>b</sup>	
100 kbp	<i>P. tremula</i>	0.339***	0.309***	0.043	0.062**	0.142***	-0.077**
	<i>P. tremuloides</i>	0.310***	0.284***	0.061**	-0.037	0.100**	-0.029
	<i>P. trichocarpa</i>	0.011	-0.024	0.053*	-0.080**	-0.002	-0.015
1 Mb	<i>P. tremula</i>	0.647***	0.573***	-0.070	0.201**	0.348**	-0.209**
	<i>P. tremuloides</i>	0.400**	0.363**	-0.033	0.032	0.320**	-0.127*
	<i>P. trichocarpa</i>	0.227**	0.151*	-0.027	-0.072	0.165*	-0.120*

<sup>a</sup> Partial correlation controls for GC content, gene density, divergence of fourfold synonymous sites between aspen and *P. trichocarpa*, and coverage (the number of fourfold synonymous bases covered by sequencing data).

<sup>b</sup> Partial correlation controls for GC content, gene density, divergence of intergenic sites between aspen and *P. trichocarpa*, and coverage (the number of intergenic bases covered by sequencing data).

\*  $P < 0.05$ .

\*\*  $P < 0.001$ .

\*\*\*  $P < 2.2 \times 10^{-16}$ .

### Distribution of fitness effects and proportion of adaptive amino acid substitutions

We quantified the efficacy of both purifying and positive selection using the information of polymorphism and divergence among the three species. The estimated distribution of fitness effects of new 0-fold nonsynonymous mutations indicates that the majority of new amino acid mutations were strongly deleterious ( $N_{\text{es}} > 10$ ) and likely to be under strong levels of purifying selection in all three species (Figure 3; Table S4). There was a greater proportion of amino acid mutations under moderate levels of purifying selection ( $1 < N_{\text{es}} < 10$ ) in *P. tremuloides* (~31%), compared to *P. tremula* (~16%) and *P. trichocarpa* (~10%). In comparison, we found a higher proportion of weakly deleterious mutations that behave as effectively neutral ( $N_{\text{es}} < 1$ ) in *P. trichocarpa* (~31%) relative to *P. tremula* (~23%) and *P. tremuloides* (~16%) (Figure 3; Table S4).

Using fourfold synonymous sites as a neutral reference, we employed an extension of the McDonald–Kreitman test (Eyre-Walker and Keightley 2009) to estimate the fraction of adaptive amino acid substitutions ( $\alpha$ ) and the rate of adaptive substitution relative to the rate of neutral substitution ( $\omega$ ) in all three species. Both  $\alpha$  and  $\omega$  were highest in *P. tremuloides* [ $\alpha$ : ~65% (95% C.I.: 63.6–65.8%);  $\omega$ : ~0.24 (95% C.I.: 0.231–0.242)], intermediate in *P. tremula* [ $\alpha$ : ~43% (95% C.I.: 41.9–43.5%);  $\omega$ : ~0.16 (95% C.I.: 0.151–0.159)], and lowest in *P. trichocarpa* [ $\alpha$ : ~20% (95% C.I.: 18.8–31.1%);  $\omega$ : ~0.07 (95% C.I.: 0.068–0.112)] (Figure 3; Table S4).

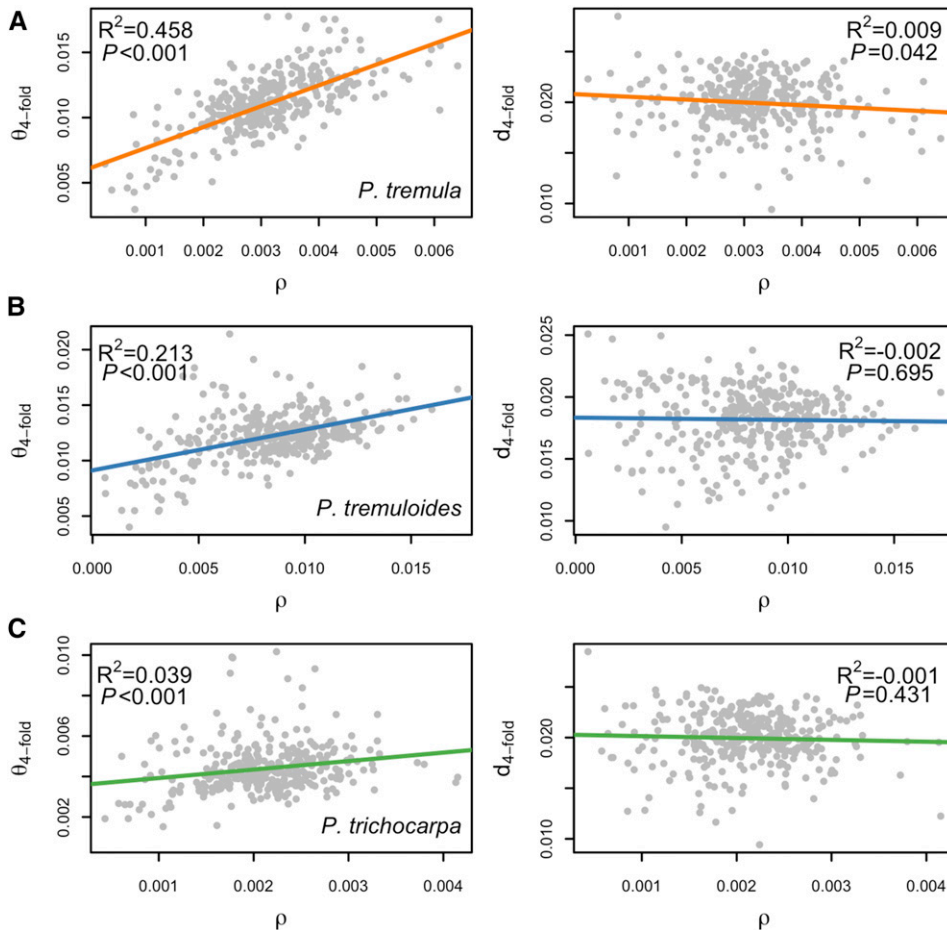
### Neutral polymorphism, but not divergence, is positively correlated with recombination rate

If natural selection (either purifying or positive selection) occurs throughout the genome at similar rates, they should leave a stronger imprint on patterns of neutral polymorphism in regions experiencing low recombination (Begun and Aquadro 1992). In accordance with this expectation, we found significantly positive correlations between levels of neutral polymorphism ( $\theta_{\text{fourfold}}$ ) and population recombination rates in both aspen species (Table 1), with correlations being

stronger in *P. tremula* than in *P. tremuloides*. In *P. trichocarpa*, however, we found either no or a weak correlation between diversity and recombination rate (Table 1). Compared to 100-kbp windows, correlations were stronger for 1-Mb windows in all species, which most likely results from the higher signal-to-noise ratio provided by larger genomic windows (Table 1). In the remainder of this article we thus focus our analyses primarily on data generated with a 1-Mb window size. When performing simple linear regression analysis between diversity and recombination rate over 1-Mb windows, the recombination rate explained 45.8, 21.3, and 3.9% of the amount of neutral genetic variation in *P. tremula*, *P. tremuloides*, and *P. trichocarpa*, respectively (Figure 4). If the positive relationship between diversity and recombination rate was merely caused by the mutagenic effect of recombination, similar patterns should also be observed between divergence and recombination rate (Kulathinal *et al.* 2008). However, no such correlations were observed in any of the three species (Table 1; Figure 4). The correlations between neutral diversity and recombination rate were slightly lower, but still significant, after using partial correlations to control for possible confounding factors such as GC content, gene density, divergence at neutral sites, and the number of neutral bases covered by sequencing data (Table 1).

In accordance with the view that genes represent the most likely targets of natural selection (Lohmueller *et al.* 2011), the correlations between intergenic diversity and recombination rate were substantially weaker than those correlations in genic regions (Table 1). Only 7.3% of the variation in intergenic nucleotide diversity in *P. tremula* could be explained by variation in the recombination rate, whereas the impact of recombination rate variation on intergenic diversity in *P. tremuloides* and *P. trichocarpa* was negligible (<1%; Figure S12; Table 1). However, after using partial correlation analyses to control for possible confounding factors, the correlations between intergenic diversity and recombination rate became significant in all species. Compared to genic regions, these correlations were slightly higher in *P. trichocarpa*, of similar magnitude in *P. tremuloides*, and weak in *P. tremula* (Table 1).





**Figure 4** Correlations of estimates between neutral genetic diversity ( $\theta_{\text{fourfold}}$ ) (left), neutral genetic divergence ( $d_{\text{fourfold}}$ ) (right), and population-scaled recombination rates ( $\rho$ ) over 1-Mb nonoverlapping windows. Linear regression lines are colored according to species: (A) *P. tremula* (orange line), (B) *P. tremuloides* (blue line), and (C) *P. trichocarpa* (green line).

### Effect of recombination on the efficacy of natural selection

We characterized the ratio of nonsynonymous-to-synonymous polymorphism ( $\theta_{0\text{-fold}}/\theta_{\text{fourfold}}$ ) and divergence ( $d_{0\text{-fold}}/d_{\text{fourfold}}$ ) to assess whether there was a relationship between the efficacy of natural selection and the rate of recombination (Table 2). Once GC content, gene density, and the number of fourfold synonymous and 0-fold nonsynonymous sites were taken into account, we found no correlation between recombination rate and  $d_{0\text{-fold}}/d_{\text{fourfold}}$  in any of the three species (Table 2). We also did not observe any significant correlations between recombination rate and  $\theta_{0\text{-fold}}/\theta_{\text{fourfold}}$  over 1-Mb windows after controlling for confounding factors (Table 2). However, when using 100-kbp windows, we found significantly negative correlations between recombination rate and  $\theta_{0\text{-fold}}/\theta_{\text{fourfold}}$  in *P. tremula* and *P. tremuloides*, but not in *P. trichocarpa* (Table 2).

### Inconsistent effect of gene density on patterns of polymorphism in genic vs. intergenic regions

We measured gene density as the number of protein-coding genes in each 1-Mb window, which in turn was highly correlated with the proportion of coding bases in each window (Figure S13). For all three species, we found significantly positive correlations between population recombination rate

and gene density (Figure 5A; Table 3). However, rather than being linear, the relationships between recombination rate and gene density were curvilinear with a significant positive correlation observed only in regions of low gene density (gene number  $< \sim 85$  within each 1-Mb window) (Table 3). For regions of high gene density (gene number  $> \sim 85$  within each 1-Mb window), we found no correlations between recombination rate and gene density in both aspen species and only a weak, positive correlation in *P. trichocarpa* (Figure 5A; Table 3). After controlling for GC content and the number of bases covered by sequencing data, the correlation became significant in regions of high gene density for *P. tremula*, but remained nonsignificant for *P. tremuloides* (Table 3).

We then examined the relationship between neutral polymorphism and gene density. Compared to the prediction of lower diversity in regions with higher functional density (Payseur and Nachman 2002), we found that the correlation pattern between gene density and levels of neutral polymorphism in genic regions ( $\theta_{\text{fourfold}}$ ) was highly consistent with the pattern found in the recombination rate, where significantly positive correlations were found in regions of low gene density and either no or weak negative correlation was found in regions of high gene density (Figure 5B; Table 3). After again controlling for potential confounding variables, the positive correlations remained significant in regions of low

**Table 2 Summary of the correlation coefficients (Spearman's rank correlation coefficient) between recombination rate ( $\rho$ ) and the ratio of nonsynonymous to synonymous polymorphism ( $\theta_{0\text{-fold}}/\theta_{\text{fourfold}}$ ) and divergence ( $d_{0\text{-fold}}/d_{\text{fourfold}}$ )**

Data set	Species	$\rho$ vs. $\theta_{0\text{-fold}}/\theta_{\text{fourfold}}$		$\rho$ vs. $d_{0\text{-fold}}/d_{\text{fourfold}}$	
		Pairwise	Partial <sup>a</sup>	Pairwise	Partial <sup>a</sup>
100 kbp	<i>P. tremula</i>	-0.057*	-0.075**	-0.012	-0.005
	<i>P. tremuloides</i>	-0.118**	-0.122**	-0.003	-0.002
	<i>P. trichocarpa</i>	-0.004	-0.002	-0.026	-0.020
1 Mb	<i>P. tremula</i>	-0.063	-0.045	-0.007	0.017
	<i>P. tremuloides</i>	-0.142*	-0.092	0.014	0.020
	<i>P. trichocarpa</i>	0.035	-0.002	0.030	0.036

<sup>a</sup>Partial correlation controls for GC content, gene density, and the number of fourfold synonymous and 0-fold nonsynonymous bases covered by sequencing data.

\*  $P < 0.05$ .

\*\*  $P < 0.001$ .

gene density among all three species (Table 3), as well as in high-gene-density regions in *P. tremuloides* and *P. trichocarpa* (Table 3). We did not find any significant relationships between neutral divergence and gene density in any of the three species (Figure S14).

Compared with genic regions, correlations between intergenic diversity and gene density followed a different pattern in the three species (Figure 5C). In intergenic regions, nucleotide diversity and gene density were positively correlated in regions of low gene density but negatively correlated in regions of high gene density (Figure 5C; Table 3). These correlations remained significant even after controlling for possible confounding variables (Table 3). No relationship between intergenic divergence and gene density was found in any species (Figure S14).

#### Negative correlations between synonymous diversity and nonsynonymous divergence at small physical scales

A negative relationship between synonymous diversity and nonsynonymous divergence has been suggested to be a strong evidence of the occurrence of recurrent selective sweeps (Andolfatto 2007), and such a pattern has previously been observed in *P. tremula* using data from a small number of candidate genes (Ingvarsson 2010). Here, however, we found either no or very weak negative correlations between neutral polymorphism ( $\theta_{\text{fourfold}}$ ) and the rate of nonsynonymous substitutions ( $d_{0\text{-fold}}$ ) in all three species for both 100-kbp and 1-Mb windows, and these correlations did not change after controlling for possible confounding factors (Table 4). However, the effects of recurrent selective sweeps on synonymous nucleotide diversity are thought to be highly localized within genes (Andolfatto 2007), and we therefore examined the association between  $\theta_{\text{fourfold}}$  and  $d_{0\text{-fold}}$  at smaller physical scales, using data from 20,759 genes that retained >90% of bases after all filtering steps. In contrast to the lack of correlations observed across larger scales (100 kbp or 1 Mb), we found a significantly negative correlation between  $\theta_{\text{fourfold}}$  and  $d_{0\text{-fold}}$  in all three species when assessed within genes (Table 4). After accounting for the

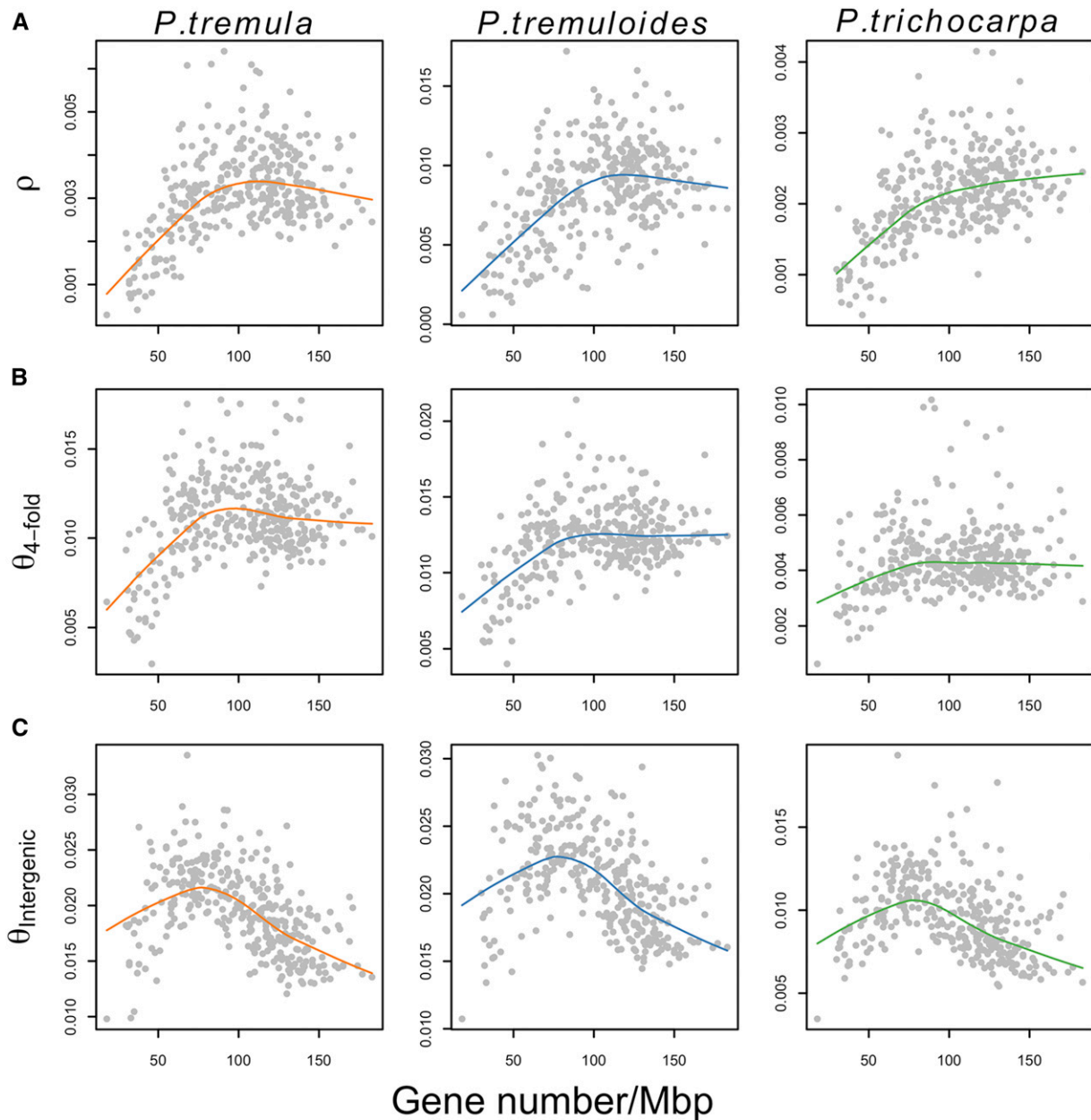
possible influence of mutation rate variation among genes by normalizing  $\theta_{\text{fourfold}}$  by neutral divergence rate ( $d_{\text{fourfold}}$ ), the negative correlations became stronger in all species (Figure S15; Table 4).

## Discussion

### Genome-wide patterns of nucleotide polymorphism, site frequency spectrum, and recombination

We have characterized and compared genome-wide patterns of nucleotide polymorphism, site frequency spectra, and recombination rates in three species of *Populus*: *P. tremula*, *P. tremuloides*, and *P. trichocarpa*. Although levels of nucleotide diversity varied greatly throughout the genome in all three species, we find strong genome-wide correlations of nucleotide diversity among species. This likely reflects conserved variation in mutation rates and/or shared selective constraints across the genomes in these closely related species during the time since their last common ancestor (Hudson *et al.* 1987; Charlesworth *et al.* 1993). Levels of nucleotide diversity are slightly higher in *P. tremuloides* than in *P. tremula*, and the two aspen species collectively harbor greater than twofold levels of genome-wide diversity compared to *P. trichocarpa*. In accordance with the larger current census population size and substantially more extensive geographic range (Eckenwalder 1996), the higher genetic diversity in both aspen species most likely reflects their larger effective population sizes ( $N_e$ ) compared to *P. trichocarpa*. Nevertheless, interspecific variation in the mutation rate also deserves further study, particularly in light of recent results showing a feed-forward effect of genome-wide levels of heterozygosity and mutation rates (Lynch 2015; Yang *et al.* 2015). Compared to the consistent pattern of nucleotide diversity between species, the weak correlations in the allele frequency spectrum (Tajima's  $D$ ) likely reflect different demographic histories for the three species during the Quaternary Ice Ages (Ingvarsson 2008; Callahan *et al.* 2013; Zhou *et al.* 2014). For example, the genome-wide excess of rare frequency alleles we observe in *P. tremuloides* is likely explained by a recent, substantial population expansion that was specific to this species.

In contrast to the mutation rate, recombination rates are only partially conserved among the three species (Figure 2C). The genome-wide average of the ratio of recombination to mutation ( $\rho/\theta_W$  or  $c/\mu$ ) was similar in *P. tremuloides* (0.39) and *P. trichocarpa* (0.38), but substantially smaller in *P. tremula* (0.22). If mutation rates are indeed unchanged between species, as suggested above, the lower estimate of  $c/\mu$  in *P. tremula* indicates considerably lower recombination rates in *P. tremula* relative to the other two species. These discrepant results obtained from patterns of polymorphism and recombination between *P. tremula* and *P. tremuloides* likely stem from different effects of effective population size on nucleotide diversity and linkage disequilibrium (Tenesa *et al.* 2007). These effects are known to operate over different timescales and therefore are likely differentially affected by



**Figure 5** Correlations of estimates between (A) population-scaled recombination rates ( $\rho$ ), (B) genic genetic diversity ( $\Theta_{\text{fourfold}}$ ), (C) intergenic genetic diversity ( $\Theta_{\text{Intergenic}}$ ), and gene density over 1-Mb nonoverlapping windows in *P. tremula* (left), *P. tremuloides* (middle), and *P. trichocarpa* (right). Gray points represent the statistics computed over 1-Mb nonoverlapping windows. Colored lines denote the lowest curves fit to the two analyzed variables in each species.

temporal variation in the effective population size (Tenesa *et al.* 2007; Cutter *et al.* 2013). The recent population size expansion that we infer to have taken place in *P. tremuloides* can thus also explain why its recombination rate is seemingly higher than in *P. tremula*, even if both species share similar levels of genome-wide polymorphism. Finally, the  $c/\mu$  estimates that we have estimated for *Populus* are in line with recent genome-wide estimates from several other plant species, such as *Medicago truncatula* (0.29) (Branca *et al.* 2011), *Mimulus guttatus* (0.8) (Hellsten *et al.* 2013), and *Eucalyptus grandis* (0.65) (Silva-Junior and Grattapaglia 2015).

### **Pervasive signatures of purifying and positive selection across the *Populus* genome**

In line with results from most other plant species (Gossmann *et al.* 2010), a majority (>50–60%) of new amino-acid-altering mutations are subject to strong purifying selection (defined as  $N_e s > 10$ ) in *Populus*. We find that the efficacy of purifying selection on weakly deleterious mutations is positively correlated with the inferred  $N_e$ , with purifying selection acting more efficiently in *P. tremuloides* that has the largest  $N_e$  compared to the other two species. The same pattern is also found for rates of adaptive evolution, where estimates of the proportion of

**Table 3 Summary of the correlation coefficients (Spearman's rank correlation coefficient) between gene density and population recombination rate ( $\rho$ ), neutral polymorphism in genic ( $\Theta_{\text{fourfold}}$ ), and intergenic regions ( $\Theta_{\text{intergenic}}$ ) over 1-Mb nonoverlapping windows in three *Populus* species**

Species	Correlation type	Gene density vs. $\rho^a$		Gene density vs. $\Theta_{\text{fourfold}}^b$		Gene density vs. $\Theta_{\text{intergenic}}^c$	
		Low	High	Low	High	Low	High
<i>P. tremula</i>	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**
<i>P. tremuloides</i>	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**
<i>P. trichocarpa</i>	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**

<sup>a</sup> Partial correlation controls for GC content and the number of bases covered by the data.

<sup>b</sup> Partial correlation controls for GC content, population recombination rate, divergence of fourfold synonymous sites between aspen and *P. trichocarpa*, and coverage (the number of fourfold synonymous bases covered by sequencing data).

<sup>c</sup> Partial correlation controls for GC content, population recombination rate, divergence of intergenic sites between aspen and *P. trichocarpa*, and coverage (the number of intergenic bases covered by sequencing data).

\*  $P < 0.05$ .

\*\*  $P < 0.001$ .

\*\*\*  $P < 2.2 \times 10^{-16}$

amino acid substitutions driven to fixation by positive selection are highest in *P. tremuloides* (65%), lowest in *P. trichocarpa* (20%), and intermediate in *P. tremula* (43%). The prevalence of adaptive evolution in *Populus* contrasts markedly with the estimates in most plant species, where little evidence of widespread adaptive evolution is found (Gossmann *et al.* 2010). However, *Populus* is not unique among plants showing high rates of adaptive evolution, and similar estimates have recently been reported in both *Capsella grandiflora* (Slotte *et al.* 2010; Williamson *et al.* 2014) and a number of *Helianthus* species (Strasburg *et al.* 2011). Most earlier studies on such estimation have been based toward subsets of genes rather than genome-wide data, and more estimates from other plant species would be valuable to assess whether the high rate of adaptive evolution that we find in *Populus* is widespread or exceptional.

Patterns of genomic variation contain abundant information on the relative importance of natural selection vs. neutral processes in the evolutionary process (Cutter and Payseur 2013). We find that 0-fold nonsynonymous sites exhibit significantly lower levels of polymorphism compared to fourfold synonymous sites, and combined with an excess of rare variants found at 0-fold nonsynonymous sites, our results suggest that the vast majority of amino acid mutations in *Populus* are under purifying selection (Larracuent *et al.* 2008). In addition, introns and 5' UTR sites are also under some degree of selective constraint, although this constraint is much weaker than what we observe at nonsynonymous sites. The 3' UTR sites seem to be either neutral or at least under comparable extents of selective constraint as fourfold synonymous sites are (Andolfatto 2005). In contrast to genic categories, we find substantially higher levels of polymorphism in intergenic regions in all three species. Although an artifact of mapping errors due to a greater fraction of repetitive sequences in intergenic regions could not be entirely excluded, the marked increase in diversity may also reflect either higher mutation rates or relaxed selective constraint in these regions. Future investigations are required to assess the relative contribution of these alternative factors (Kimura 1983; Begun *et al.* 2007).

Apart from strong selective constraints on protein-coding genes, multiple lines of evidence suggest that genome-wide patterns of polymorphism have been shaped by widespread natural selection in all three *Populus* species. First, we find significantly positive correlations between neutral polymorphism and population-scaled recombination rates in both genic and intergenic regions, even after controlling for confounding variables such as GC content, gene density, mutation rate, and the number of sites covered by the data. While such a pattern is indicative of the action of natural selection, it could be explained by either background selection or selective sweeps. Both of these selective forces affect neutral sites through linkage, and the impact of selection on linked neutral diversity is more drastic and extensive in regions of low recombination (Begun and Aquadro 1992; McGaugh *et al.* 2012; Slotte 2014). The differences in the strength of the association between recombination and levels of neutral polymorphism likely reflect differences in the effective population size between species (Cutter and Payseur 2013; Corbett-Detig *et al.* 2015), as we observe substantially stronger signatures of linked selection in *P. tremula* and *P. tremuloides* compared to *P. trichocarpa*, matching the larger  $N_e$  inferred for these species. However, the impact of natural selection at linked sites also depends greatly on the local environment of recombination (Cutter and Payseur 2013; Slotte 2014), and in line with this we observe the strongest signatures of linked selection in *P. tremula* instead of *P. tremuloides*, consistent with the lower levels of genome-wide recombination rates that we find in *P. tremula*. Different magnitudes of linked selection provide one of the major explanations for the disparate patterns of genomic variation among even closely related species (Corbett-Detig *et al.* 2015), and we find that this also holds true for the three species of *Populus* that we have investigated.

Second, we find slightly negative correlations between recombination rate and the ratio of nonsynonymous-to-synonymous polymorphism, but not divergence, in *P. tremula* and *P. tremuloides*, a pattern that suggests a reduced efficacy of purifying selection at eliminating weakly deleterious



**Table 4 Summary of the correlation coefficients (Spearman's rank correlation coefficient) between levels of synonymous diversity ( $\theta_{\text{fourfold}}$ ) and nonsynonymous divergence ( $d_{0\text{-fold}}$ ) at different physical scales in three *Populus* species**

Data set	Species	$d_{0\text{-fold}}$ vs. $\theta_{\text{fourfold}}$	
		Pairwise	Partial
100 kbp <sup>a</sup>	<i>P. tremula</i>	-0.029	-0.032
	<i>P. tremuloides</i>	-0.021	-0.025
	<i>P. trichocarpa</i>	-0.053*	-0.051*
1 Mb <sup>a</sup>	<i>P. tremula</i>	-0.049	0.043
	<i>P. tremuloides</i>	-0.069	-0.008
	<i>P. trichocarpa</i>	-0.086	-0.006
Single genes <sup>b</sup>	<i>P. tremula</i>	-0.087**	-0.185**
	<i>P. tremuloides</i>	-0.087**	-0.192**
	<i>P. trichocarpa</i>	-0.148**	-0.218**

<sup>a</sup> Partial means partial correlation controls for GC content, gene density, population recombination rate, divergence of fourfold synonymous sites between aspen and *P. trichocarpa*, the number of fourfold synonymous bases and 0-fold nonsynonymous bases covered by sequencing data.

<sup>b</sup> Partial means correlation between  $d_{0\text{-fold}}$  and  $\theta_{\text{fourfold}}/d_{\text{fourfold}}$ .

\*  $P < 0.05$ .

\*\*  $P < 2.2 \times 10^{-16}$

mutations in low-recombination regions (Hill and Robertson 1966; Cutter and Choi 2010). The reduction of the efficacy of natural selection in regions of low recombination, known as Hill–Robertson interference, may help to understand patterns of partially positive correlations between gene density and recombination rate in these species (Gaut *et al.* 2007). Given the relaxed efficacy of purifying selection in regions of low recombination where weakly deleterious mutations are more likely to accumulate at a high rate, important functional elements are unlikely to cluster in these regions, as has already been shown in several other plant species (Anderson *et al.* 2006; Branca *et al.* 2011; Flowers *et al.* 2012). Consistent with this prediction (Haddrill *et al.* 2007), we find positive association between gene density and recombination rate in regions that experience low rates of recombination. In high-recombination regions where selection is more effective at eliminating slightly deleterious mutations, the association becomes much weaker in all three species. However, it remains unclear whether it is the recombination gradients that drive the functional organization of genomes in response to selection, or whether it is the gradients of functional genomic elements that in turn modify the evolution of recombination rates in *Populus*.

Third, by examining the relationship of neutral polymorphism, recombination rate, and gene density, we find that levels of neutral polymorphism in genic regions are primarily driven by local rates of recombination, regardless of the density of functional genes. In contrast, we observe a more complex pattern in intergenic regions where levels of intergenic polymorphism are driven mainly by recombination rates in regions of low gene density, while in regions of high-gene-density levels of intergenic diversity are primarily shaped by the density of nearby genes. As we find that gene density and recombination rates covary in all three species, the signatures of linked selection associated with gene density could thus become obscured by rates of recombination, especially in

regions of low gene density (Flowers *et al.* 2012). As shown in most plants studied so far (Nordborg *et al.* 2005; Slotte 2014), a negative relationship between gene density and levels of neutral polymorphism is more likely attributed to more intense purifying selection against deleterious mutations in regions of greater gene density, and the magnitude of such effects depends on the strength of purifying selection (Sella *et al.* 2009). In accordance with this expectation, most new mutations in genic regions are strongly deleterious and would be eliminated too quickly to remove large amounts of genetic variation at linked neutral loci. Thus even in regions of high gene density, we do not find negative correlations between gene density and genetic diversity in genic regions. However, background selection due to deleterious mutations of moderate effect in intergenic regions could account for the negative association that we observe between levels of intergenic polymorphism and gene density in regions of high gene density. It is apparent that the extent to which natural selection is acting on noncoding regions of the genome in *Populus* (e.g., intergenic regions) will be an interesting avenue for future studies.

Finally, in all three *Populus* species we find significantly negative correlations between levels of synonymous polymorphism and the rate of amino acid substitution at the scale of single genes. This pattern could be driven by either recurrent selective sweeps or background selection (Charlesworth *et al.* 1993; Andolfatto 2007). However, background selection reduces local  $N_e$  due to the removal of weakly deleterious mutations and is therefore expected to result in both reduced levels of nucleotide polymorphism and an increase of the fixation rate of slightly deleterious mutations (Charlesworth *et al.* 1993). Background selection is thus expected to affect the rates of both synonymous and nonsynonymous substitutions equally, but when variation in the rates of synonymous substitution is taken into account, we find a substantially stronger (rather than weaker) negative correlation between levels of synonymous polymorphism and the rate of protein evolution. This suggests that the negative relationship that we observe between nonsynonymous substitution rate and levels of variation at synonymous sites is most likely driven by effects of recurrent selective sweeps in all three species (Andolfatto 2007; Sella *et al.* 2009). Furthermore, the physical scale at which these signatures of natural selection are detected carries valuable information about the strength of positive selection at the genomic level (Macpherson *et al.* 2007). Since the signatures of recurrent selective sweeps are detectable only on a genic scale, it mostly reflects relatively weak selection on the majority of adaptive amino acid substitutions and may thus explain why we do not observe the effects at either 100-kbp or 1-Mb scales (Macpherson *et al.* 2007; Sella *et al.* 2009).

### Conclusion and perspectives

In summary, our findings highlight multiple lines of evidence suggesting that natural selection, due to both purifying and positive selection, has shaped patterns of nucleotide polymorphism at linked neutral sites in all three *Populus* species. Compared to the predictions of the Neutral Theory that

suggest that adaptations contribute negligibly to divergence between species (Kimura 1983), we find that ~20–65% of all amino acid substitutions are driven to fixation by adaptive evolution in *Populus*. These estimates are in accordance with the results from a number of other organisms with large effective population sizes, such as *Drosophila* (Sella *et al.* 2009), mammals (Halligan *et al.* 2010; Carneiro *et al.* 2012), and a few plant species (Slotte *et al.* 2010; Strasburg *et al.* 2011), but substantially higher than in species with relatively small effective population sizes, such as humans and most other plant species, where little evidence of adaptive evolution has been detected (Eyre-Walker and Keightley 2009; Gossmann *et al.* 2010). Given that all three *Populus* species share similar life-cycle characteristics, such as an outcrossing mating system, relatively large  $N_e$ , and limited population subdivision, future studies from other long-lived forest trees are needed to investigate whether these are characteristics more generally influencing genome-wide patterns of selection in plants (Hough *et al.* 2013). Furthermore, differences in  $N_e$  and rates of recombination among the three *Populus* species largely explain differences in the magnitude of linked selection that we observe between them.

Our analyses suggest pervasive adaptive evolution in all three species of *Populus*, and although alternative hypotheses such as demographic effects could lead to spurious evidence of natural selection (Fay *et al.* 2001), the presence of linked selection could also bias inferences of demographic history (Slotte 2014). Due to the pervasive effects of linked selection that we have documented in these species, our findings suggest that more attention should be paid to the process of choosing neutral sites for demographic inferences. Alternatively, new methods that allow for the joint estimation of demography and selection from genome-wide data are urgently needed.

## Acknowledgments

We thank Rick Lindroth for providing access to the samples of *P. tremuloides* used in this study; Carin Olofsson for extracting DNA for all samples used in this study; and Robert J. Williamson for sharing a data analysis script. We also thank both the editor and two anonymous referees for useful comments on the manuscript. The research has been funded through grants from Vetenskapsrådet and a Young Researcher Award from Umeå University (to P.K.I.). J.W. was supported by a scholarship from the Chinese Scholarship Council.

## Literature Cited

- Alexander, D. H., J. Novembre, and K. Lange, 2009 Fast model-based estimation of ancestry in unrelated individuals. *Genome Res.* 19: 1655–1664.
- Anderson, L. K., A. Lai, S. M. Stack, C. Rizzon, and B. S. Gaut, 2006 Uneven distribution of expressed sequence tag loci on maize pachytene chromosomes. *Genome Res.* 16: 115–122.
- Andolfatto, P., 2005 Adaptive evolution of non-coding DNA in *Drosophila*. *Nature* 437: 1149–1152.
- Andolfatto, P., 2007 Hitchhiking effects of recurrent beneficial amino acid substitutions in the *Drosophila melanogaster* genome. *Genome Res.* 17: 1755–1762.
- Begun, D. J., and C. F. Aquadro, 1992 Levels of naturally occurring DNA polymorphism correlate with recombination rates in *D. melanogaster*. *Nature* 356: 519–520.
- Begun, D. J., A. K. Holloway, K. Stevens, L. W. Hillier, Y.-P. Poh *et al.*, 2007 Population genomics: whole-genome analysis of polymorphism and divergence in *Drosophila simulans*. *PLoS Biol.* 5: e310.
- Branca, A., T. D. Paape, P. Zhou, R. Briskine, A. D. Farmer *et al.*, 2011 Whole-genome nucleotide diversity, recombination, and linkage disequilibrium in the model legume *Medicago truncatula*. *Proc. Natl. Acad. Sci. USA* 108: E864–E870.
- Callahan, C. M., C. A. Rowe, R. J. Rye, J. D. Shaw, M. D. Madritch *et al.*, 2013 Continental-scale assessment of genetic diversity and population structure in quaking aspen (*Populus tremuloides*). *J. Biogeogr.* 40: 1780–1791.
- Campos, J. L., D. L. Halligan, P. R. Haddrill, and B. Charlesworth, 2014 The relation between recombination rate and patterns of molecular evolution and variation in *Drosophila melanogaster*. *Mol. Biol. Evol.* 31: 1010–1028.
- Carneiro, M., F. W. Albert, J. Melo-Ferreira, N. Galtier, P. Gayral *et al.*, 2012 Evidence for widespread positive and purifying selection across the European rabbit (*Oryctolagus cuniculus*) genome. *Mol. Biol. Evol.* 29: 1837–1849.
- Charlesworth, B., and J. L. Campos, 2014 The relations between recombination rate and patterns of molecular variation and evolution in *Drosophila*. *Annu. Rev. Genet.* 48: 383–403.
- Charlesworth, B., M. Morgan, and D. Charlesworth, 1993 The effect of deleterious mutations on neutral molecular variation. *Genetics* 134: 1289–1303.
- Corbett-Detig, R. B., D. L. Hartl, and T. B. Sackton, 2015 Natural selection constrains neutral diversity across a wide range of species. *PLoS Biol.* 13: e1002112.
- Cutter, A. D., and J. Y. Choi, 2010 Natural selection shapes nucleotide polymorphism across the genome of the nematode *Caenorhabditis briggsae*. *Genome Res.* 20: 1103–1111.
- Cutter, A. D., and B. A. Payseur, 2003 Selection at linked sites in the partial selfer *Caenorhabditis elegans*. *Mol. Biol. Evol.* 20: 665–673.
- Cutter, A. D., and B. A. Payseur, 2013 Genomic signatures of selection at linked sites: unifying the disparity among species. *Nat. Rev. Genet.* 14: 262–274.
- Cutter, A. D., R. Jovelin, and A. Dey, 2013 Molecular hyperdiversity and evolution in very large populations. *Mol. Ecol.* 22: 2074–2095.
- De Carvalho, D., P. K. Ingvarsson, J. Joseph, L. Suter, C. Sedivy *et al.*, 2010 Admixture facilitates adaptation from standing variation in the European aspen (*Populus tremula* L.), a widespread forest tree. *Mol. Ecol.* 19: 1638–1650.
- DePristo, M. A., E. Banks, R. Poplin, K. V. Garimella, J. R. Maguire *et al.*, 2011 A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* 43: 491–498.
- Dickmann, D. I., and J. Kuzovkina, 2014 Poplars and willows of the world, with emphasis on silviculturally important species, pp. 8–91 in *Poplars and Willows: Trees for Society and the Environment*, edited by J. D. Isebrands and J. Richardson. The Food and Agriculture Organization of the United Nations and CAB International, Rome.
- Eckenwalder, J. E., 1996 Systematics and evolution of *Populus*, pp. 7–32 in *Biology of Populus and Its Implications for Management and Conservation (Part I)*, edited by R. F. Stettler, H. D. Bradshaw, P. E. Heilman, and T. M. Hinckley. NRC Research Press, Ottawa.
- Ellegren, H., 2014 Genome sequencing and population genomics in non-model organisms. *Trends Ecol. Evol.* 29: 51–63.

- Evans, L. M., G. T. Slavov, E. Rodgers-Melnick, J. Martin, P. Ranjan *et al.*, 2014 Population genomics of *Populus trichocarpa* identifies signatures of selection and adaptive trait associations. *Nat. Genet.* 46: 1089–1096.
- Eyre-Walker, A., and P. D. Keightley, 2009 Estimating the rate of adaptive molecular evolution in the presence of slightly deleterious mutations and population size change. *Mol. Biol. Evol.* 26: 2097–2108.
- Fay, J. C., G. J. Wyckoff, and C.-I. Wu, 2001 Positive and negative selection on the human genome. *Genetics* 158: 1227–1234.
- Flowers, J. M., J. Molina, S. Rubinstein, P. Huang, B. A. Schaal *et al.*, 2012 Natural selection in gene-dense regions shapes the genomic pattern of polymorphism in wild and domesticated rice. *Mol. Biol. Evol.* 29: 675–687.
- Fumagalli, M., F. G. Vieira, T. Linderöth, and R. Nielsen, 2014 ngsTools: methods for population genetics analyses from next-generation sequencing data. *Bioinformatics* 30: 1486–1487.
- Gaut, B. S., S. I. Wright, C. Rizzon, J. Dvorak, and L. K. Anderson, 2007 Recombination: an underappreciated factor in the evolution of plant genomes. *Nat. Rev. Genet.* 8: 77–84.
- González-Martínez, S. C., K. V. Krutovsky, and D. B. Neale, 2006 Forest-tree population genomics and adaptive evolution. *New Phytol.* 170: 227–238.
- Gossmann, T. I., B.-H. Song, A. J. Windsor, T. Mitchell-Olds, C. J. Dixon *et al.*, 2010 Genome wide analyses reveal little evidence for adaptive evolution in many plant species. *Mol. Biol. Evol.* 27: 1822–1832.
- Haddrill, P. R., D. L. Halligan, D. Tomaras, and B. Charlesworth, 2007 Reduced efficacy of selection in regions of the *Drosophila* genome that lack crossing over. *Genome Biol.* 8: R18.
- Halligan, D. L., F. Oliver, A. Eyre-Walker, B. Harr, and P. D. Keightley, 2010 Evidence for pervasive adaptive protein evolution in wild mice. *PLoS Genet.* 6: e1000825.
- Hamzeh, M., and S. Dayanandan, 2004 Phylogeny of *Populus* (Salicaceae) based on nucleotide sequences of chloroplast trnT-trnF region and nuclear rDNA. *Am. J. Bot.* 91: 1398–1408.
- Hellmann, I., K. Prüfer, H. Ji, M. C. Zody, S. Pääbo *et al.*, 2005 Why do human diversity levels vary at a megabase scale? *Genome Res.* 15: 1222–1231.
- Hellsten, U., K. M. Wright, J. Jenkins, S. Shu, Y. Yuan *et al.*, 2013 Fine-scale variation in meiotic recombination in *Mimulus* inferred from population shotgun sequencing. *Proc. Natl. Acad. Sci. USA* 110: 19478–19482.
- Hill, W. G., and A. Robertson, 1966 The effect of linkage on limits to artificial selection. *Genet. Res.* 8: 269–294.
- Hough, J., R. J. Williamson, and S. I. Wright, 2013 Patterns of selection in plant genomes. *Annu. Rev. Ecol. Evol. Syst.* 44: 31–49.
- Hudson, R. R., M. Kreitman, and M. Aguadé, 1987 A test of neutral molecular evolution based on nucleotide data. *Genetics* 116: 153–159.
- Hufford, M. B., X. Xu, J. Van Heerwaarden, T. Pyhäjärvi, J.-M. Chia *et al.*, 2012 Comparative population genomics of maize domestication and improvement. *Nat. Genet.* 44: 808–811.
- Ingvarsson, P. K., 2008 Multilocus patterns of nucleotide polymorphism and the demographic history of *Populus tremula*. *Genetics* 180: 329–340.
- Ingvarsson, P. K., 2010 Natural selection on synonymous and nonsynonymous mutations shapes patterns of polymorphism in *Populus tremula*. *Mol. Biol. Evol.* 27: 650–660.
- Jansson, S., and C. J. Douglas, 2007 *Populus*: a model system for plant biology. *Annu. Rev. Plant Biol.* 58: 435–458.
- Jansson, S., R. P. Bhalerao, and A. T. Groover, 2010 *Genetics and genomics of Populus*, Springer, New York.
- Jukes, T. H., and C. R. Cantor, 1969 Evolution of protein molecules, pp. 21–132 in *Mammalian Protein Metabolism*, edited by H. N. Munro. Academic Press, New York.
- Keightley, P. D., and A. Eyre-Walker, 2007 Joint inference of the distribution of fitness effects of deleterious mutations and population demography based on nucleotide polymorphism frequencies. *Genetics* 177: 2251–2261.
- Kim, S. H., and V. Y. Soojin, 2007 Understanding relationship between sequence and functional evolution in yeast proteins. *Genetica* 131: 151–156.
- Kim, S. Y., K. E. Lohmueller, A. Albrechtsen, Y. Li, T. Korneliussen *et al.*, 2011 Estimation of allele frequency and association mapping using next-generation sequencing data. *BMC Bioinformatics* 12: 231.
- Kimura, M., 1983 *The Neutral Theory of Molecular Evolution*. Cambridge University Press, Cambridge, UK.
- Korneliussen, T. S., A. Albrechtsen, and R. Nielsen, 2014 ANGSD: analysis of next generation sequencing data. *BMC Bioinformatics* 15: 356.
- Kulathinal, R. J., S. M. Bennett, C. L. Fitzpatrick, and M. A. Noor, 2008 Fine-scale mapping of recombination rate in *Drosophila* refines its correlation to diversity and divergence. *Proc. Natl. Acad. Sci. USA* 105: 10051–10056.
- Larracuente, A. M., T. B. Sackton, A. J. Greenberg, A. Wong, N. D. Singh *et al.*, 2008 Evolution of protein-coding genes in *Drosophila*. *Trends Genet.* 24: 114–123.
- Lawrie, D. S., and D. A. Petrov, 2014 Comparative population genomics: power and principles for the inference of functionality. *Trends Genet.* 30(4): 133–139.
- Li, H., 2013 Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. arXiv: 1303.3997. Available at: <http://arxiv.org/abs/1303.3997>.
- Li, H., and R. Durbin, 2011 Inference of human population history from individual whole-genome sequences. *Nature* 475: 493–496.
- Li, H., B. Handsaker, A. Wysoker, T. Fennell, J. Ruan *et al.*, 2009 The sequence alignment/map format and SAMtools. *Bioinformatics* 25: 2078–2079.
- Liu, S., E. D. Lorenzen, M. Fumagalli, B. Li, K. Harris *et al.*, 2014 Population genomics reveal recent speciation and rapid evolutionary adaptation in polar bears. *Cell* 157: 785–794.
- Locke, D. P., L. W. Hillier, W. C. Warren, K. C. Worley, L. V. Nazareth *et al.*, 2011 Comparative and demographic analysis of orangutan genomes. *Nature* 469: 529–533.
- Lohmueller, K. E., A. Albrechtsen, Y. Li, S. Y. Kim, T. Korneliussen *et al.*, 2011 Natural selection affects multiple aspects of genetic variation at putatively neutral sites across the human genome. *PLoS Genet.* 7: e1002326.
- Lohse, M., A. Bolger, A. Nagel, A. R. Fernie, J. E. Lunn *et al.*, 2012 RobiNA: a user-friendly, integrated software solution for RNA-Seq-based transcriptomics. *Nucleic Acids Res.* 40: W622–W627.
- Luikart, G., P. R. England, D. Tallmon, S. Jordan, and P. Taberlet, 2003 The power and promise of population genomics: from genotyping to genome typing. *Nat. Rev. Genet.* 4: 981–994.
- Lynch, M., 2015 Genetics: feedforward loop for diversity. *Nature* 523: 414–416.
- Mackay, T. F., S. Richards, E. A. Stone, A. Barbadilla, J. F. Ayroles *et al.*, 2012 The *Drosophila melanogaster* genetic reference panel. *Nature* 482: 173–178.
- Macpherson, J. M., G. Sella, J. C. Davis, and D. A. Petrov, 2007 Genomewide spatial correspondence between nonsynonymous divergence and neutral polymorphism reveals extensive adaptation in *Drosophila*. *Genetics* 177: 2083–2099.
- McGaugh, S. E., C. S. Heil, B. Manzano-Winkler, L. Loewe, S. Goldstein *et al.*, 2012 Recombination modulates how selection affects linked sites in *Drosophila*. *PLoS Biol.* 10: e1001422.
- McVean, G. A., S. R. Myers, S. Hunt, P. Deloukas, D. R. Bentley *et al.*, 2004 The fine-scale structure of recombination rate variation in the human genome. *Science* 304: 581–584.

- Neale, D. B., and A. Kremer, 2011 Forest tree genomics: growing resources and applications. *Nat. Rev. Genet.* 12: 111–122.
- Nevado, B., S. Ramos-Onsins, and M. Perez-Enciso, 2014 Resequencing studies of nonmodel organisms using closely related reference genomes: optimal experimental designs and bioinformatics approaches for population genomics. *Mol. Ecol.* 23: 1764–1779.
- Nielsen, R., T. Korneliussen, A. Albrechtsen, Y. Li, and J. Wang, 2011 SNP calling, genotype calling, and sample allele frequency estimation from New-Generation Sequencing data. *PLoS One* 7: e37558.
- Nordborg, M., T. T. Hu, Y. Ishino, J. Jhaveri, C. Toomajian *et al.*, 2005 The pattern of polymorphism in *Arabidopsis thaliana*. *PLoS Biol.* 3: 1289.
- Payseur, B. A., and M. W. Nachman, 2002 Gene density and human nucleotide polymorphism. *Mol. Biol. Evol.* 19: 336–340.
- Purcell, S., B. Neale, K. Todd-Brown, L. Thomas, M. A. Ferreira *et al.*, 2007 PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am. J. Hum. Genet.* 81: 559–575.
- R Development Core Team, 2014 *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna.
- Remington, D. L., J. M. Thornsberry, Y. Matsuoka, L. M. Wilson, S. R. Whitt *et al.*, 2001 Structure of linkage disequilibrium and phenotypic associations in the maize genome. *Proc. Natl. Acad. Sci. USA* 98: 11479–11484.
- Sella, G., D. A. Petrov, M. Przeworski, and P. Andolfatto, 2009 Pervasive natural selection in the *Drosophila* genome. *PLoS Genet.* 5: e1000495.
- Silva-Junior, O. B., and D. Grattapaglia, 2015 Genome-wide patterns of recombination, linkage disequilibrium and nucleotide diversity from pooled resequencing and single nucleotide polymorphism genotyping unlock the evolutionary history of *Eucalyptus grandis*. *New Phytol.* 208: 830–845.
- Slotte, T., 2014 The impact of linked selection on plant genomic variation. *Brief. Funct. Genomics* 13: 268–275.
- Slotte, T., J. P. Foxe, K. M. Hazzouri, and S. I. Wright, 2010 Genome-wide evidence for efficient positive and purifying selection in *Capsella grandiflora*, a plant species with a large effective population size. *Mol. Biol. Evol.* 27: 1813–1821.
- Strasburg, J. L., N. C. Kane, A. R. Raduski, A. Bonin, R. Michelmore *et al.*, 2011 Effective population size is positively correlated with levels of adaptive divergence among annual sunflowers. *Mol. Biol. Evol.* 28: 1569–1580.
- Tajima, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* 123: 585–595.
- Tarailo-Graovac, M., and N. Chen, 2009 Using RepeatMasker to identify repetitive elements in genomic sequences. *Curr. Protoc. in Bioinformatics Chapter 4: Unit 4.10*.
- Tenesa, A., P. Navarro, B. J. Hayes, D. L. Duffy, G. M. Clarke *et al.*, 2007 Recent human effective population size estimated from linkage disequilibrium. *Genome Res.* 17: 520–526.
- Tuskan, G. A., S. Difazio, S. Jansson, J. Bohlmann, I. Grigoriev *et al.*, 2006 The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* 313: 1596–1604.
- Wang, J., D. Scofield, N. R. Street, and P. K. Ingvarsson, 2015 Variant calling using NGS data in European aspen (*Populus tremula*), pp. 43–61 in *Advances in the Understanding of Biological Sciences Using Next Generation Sequencing (NGS) Approaches*, edited by G. Sablok, S. Kumar, S. Ueno, J. Kuo, and C. Varotto. Springer, New York.
- Wang, Z., S. Du, S. Dayanandan, D. Wang, Y. Zeng *et al.*, 2014 Phylogeny reconstruction and hybrid analysis of *Populus* (Salicaceae) based on nucleotide sequences of multiple single-copy nuclear genes and plastid fragments. *PLoS One* 9: e103645.
- Watterson, G., 1975 On the number of segregating sites in genetical models without recombination. *Theor. Popul. Biol.* 7: 256–276.
- Williamson, R. J., E. B. Josephs, A. E. Platts, K. M. Hazzouri, A. Haudry *et al.*, 2014 Evidence for widespread positive and negative selection in coding and conserved noncoding regions of *Capsella grandiflora*. *PLoS Genet.* 10: e1004622.
- Yang, S., L. Wang, J. Huang, X. Zhang, Y. Yuan *et al.*, 2015 Parent-progeny sequencing indicates higher mutation rates in heterozygotes. *Nature* 523: 463–467.
- Zhou, L., R. Bawa, and J. Holliday, 2014 Exome resequencing reveals signatures of demographic and adaptive processes across the genome and range of black cottonwood (*Populus trichocarpa*). *Mol. Ecol.* 23: 2486–2499.

Communicating editor: S. C. Gonzalez-Martinez



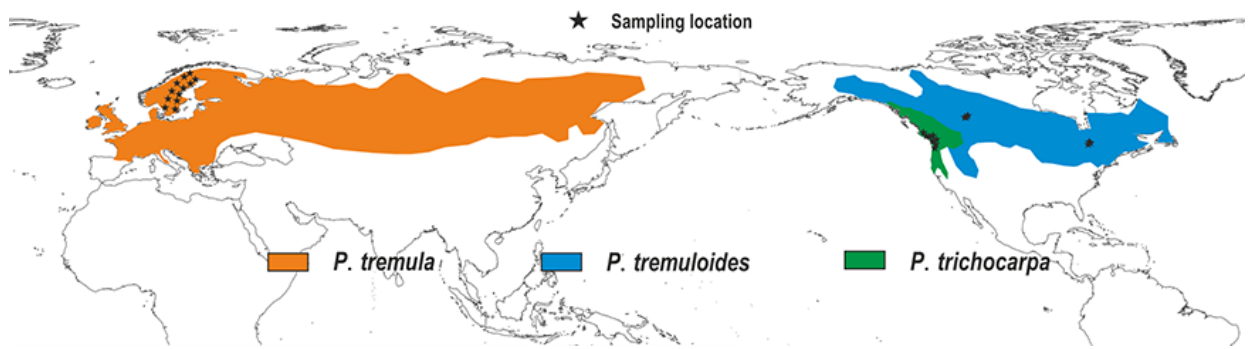
# GENETICS

Supporting Information

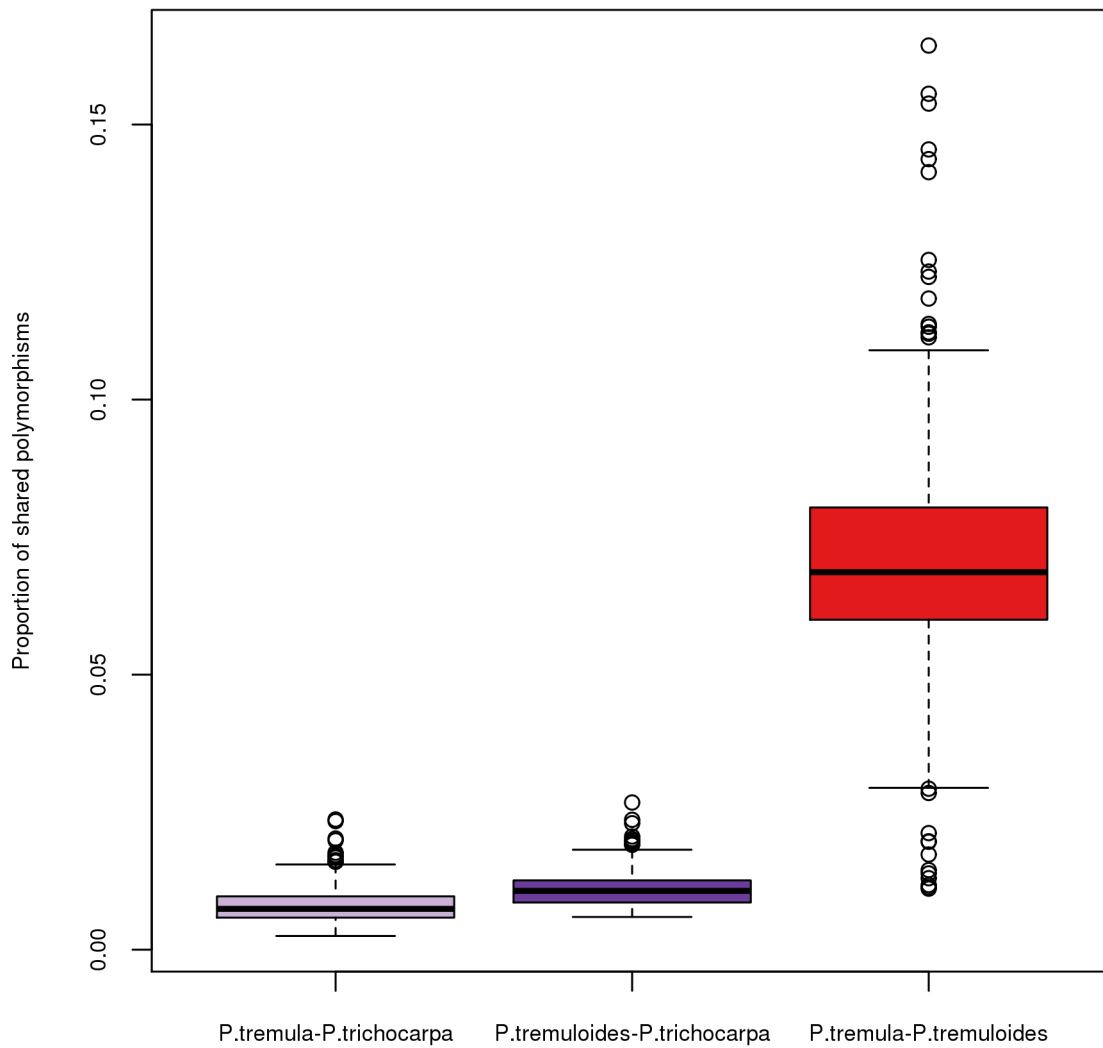
[www.genetics.org/lookup/suppl/doi:10.1534/genetics.115.183152/-/DC1](http://www.genetics.org/lookup/suppl/doi:10.1534/genetics.115.183152/-/DC1)

## Natural Selection and Recombination Rate Variation Shape Nucleotide Polymorphism Across the Genomes of Three Related *Populus* Species

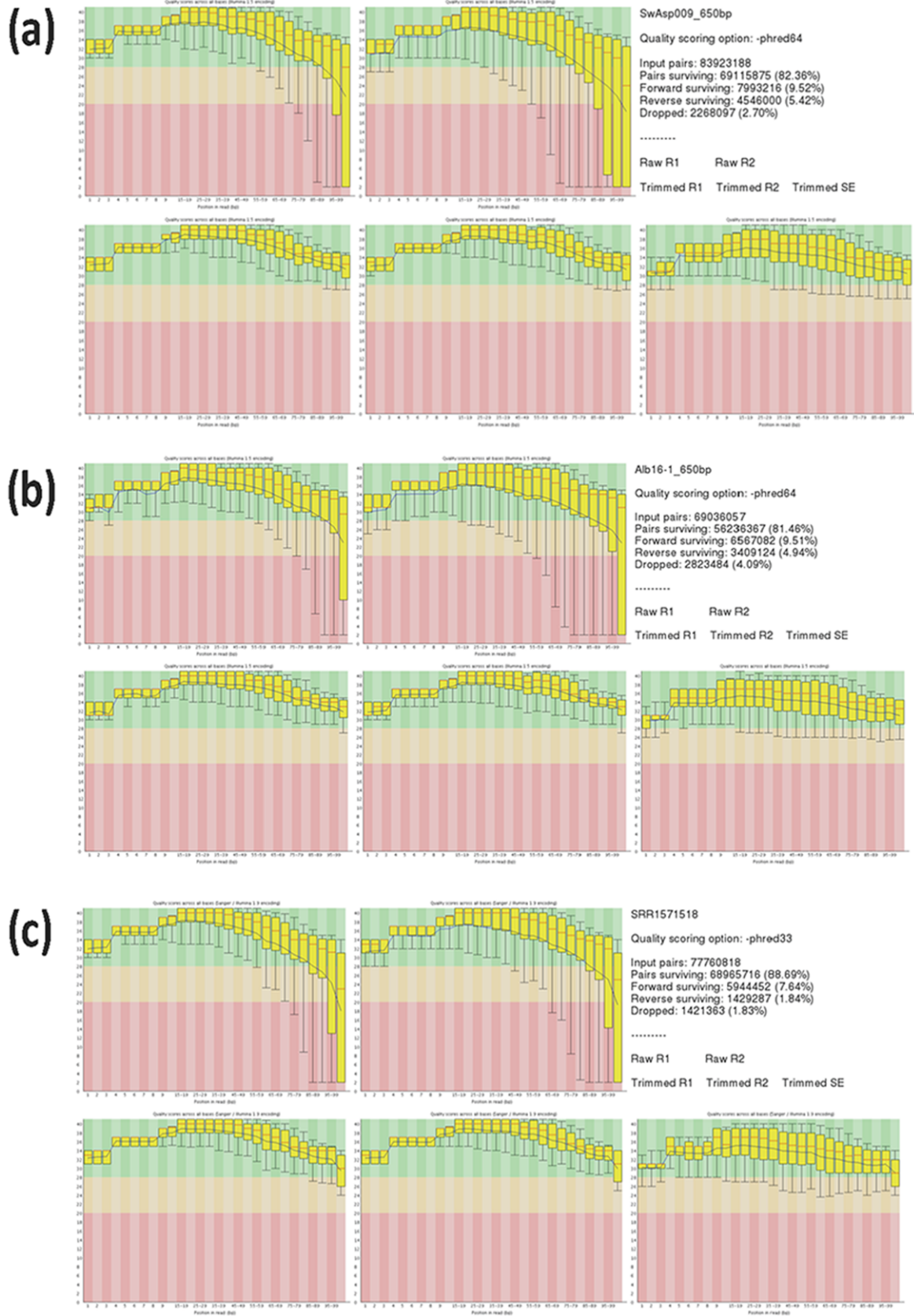
Jing Wang, Nathaniel R. Street, Douglas G. Scofield, and Pär K. Ingvarsson



**Figure S1.** Sampling localities (details in Table S1, black star symbols) and distributions of *P. tremula* (orange areas), *P. tremuloides* (blue areas) and *P. trichocarpa* (green areas).



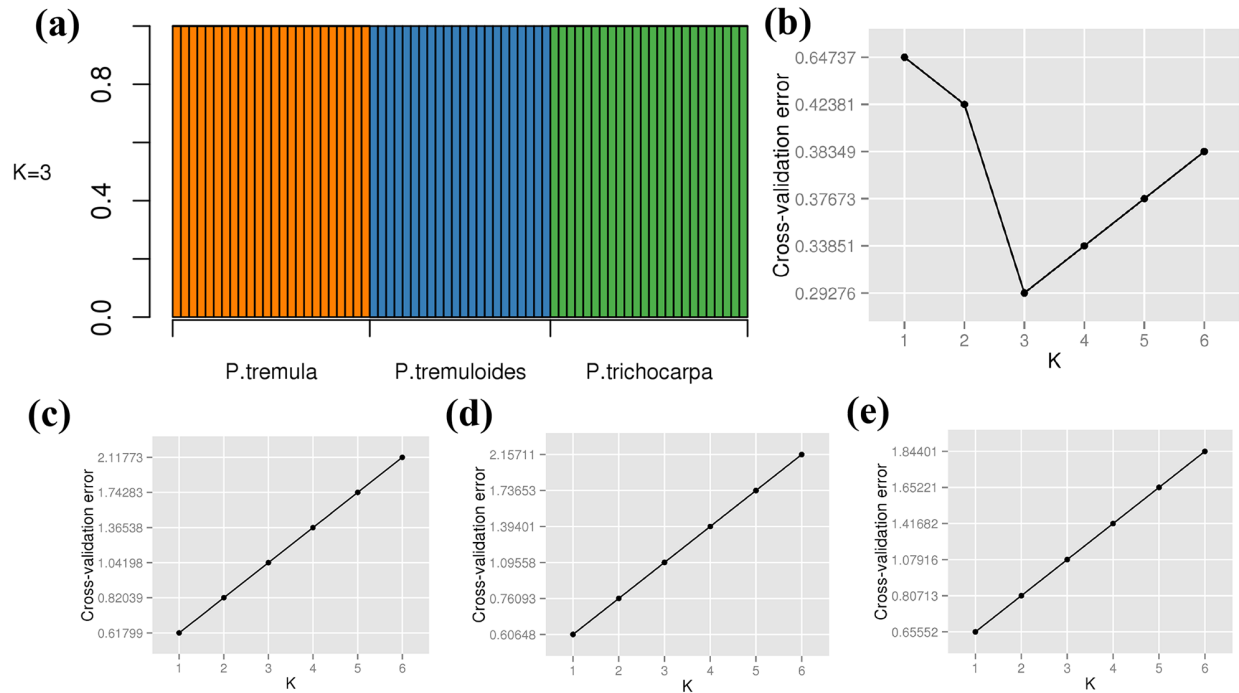
**Figure S2.** The distributions of estimates of the proportion of shared polymorphisms between pairs of three *Populus* species among the total number of polymorphisms in all three species. It was calculated over 1 Mbp non-overlapping windows across the genome. *P. tremula*-*P.tremuloies*: red, *P. tremula*-*P. trichocarpa*: light purple, *P. tremuloides*-*P.trichocarpa*: purple.



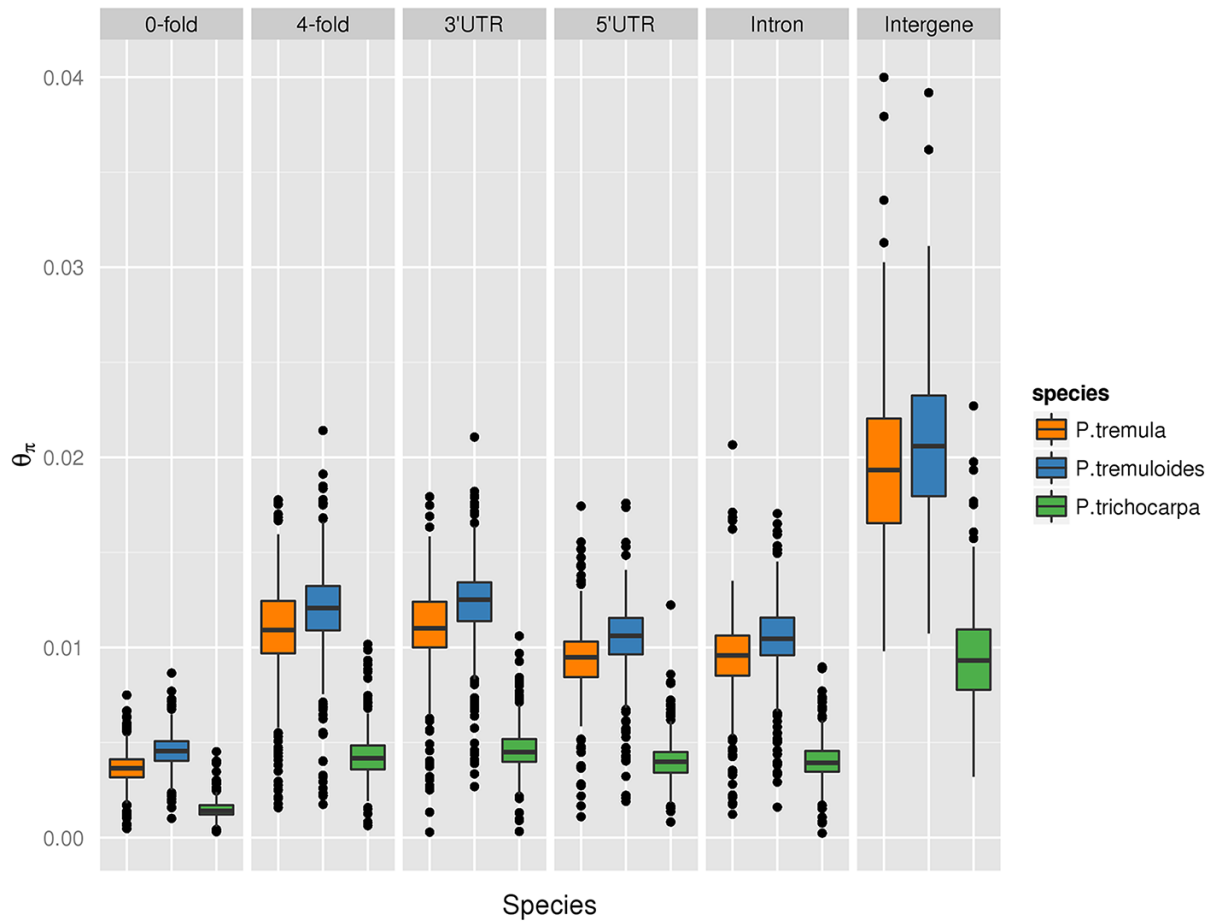
**Figure S3.** Comparison of per-base sequence quality between raw and filtered sequence data. Per-base sequence quality comparison between raw paired-end sequence data (forward reads: top



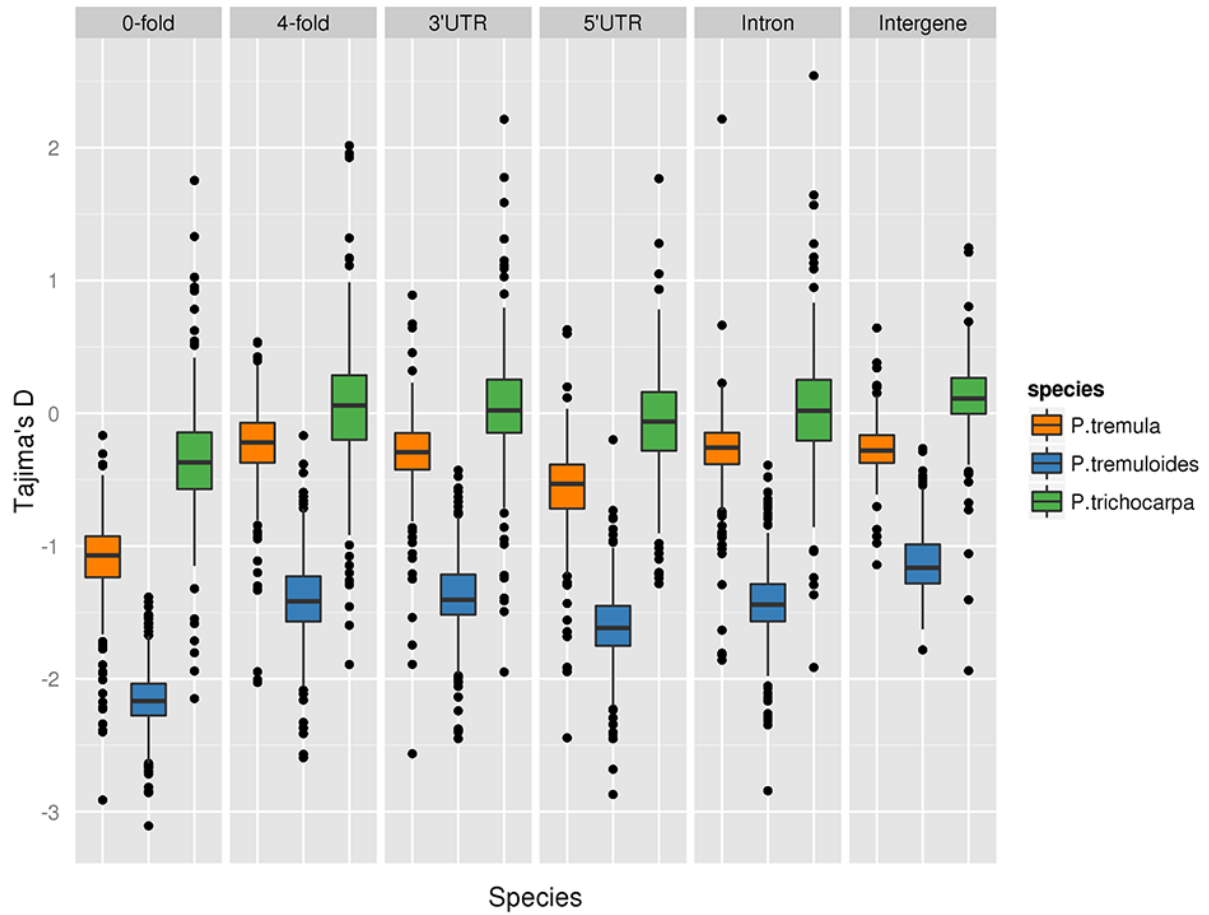
left and reverse reads: top right), and filtered sequence data with both forward (bottom left) and reverse (bottom middle) reads left or only single-end (bottom right) reads left. The x-axis of the BoxWhisker plot shows the position in read, and the y-axis shows the quality scores. The higher the score the better the base call. The background of the plot divides the y axis into very good quality calls (green), calls of reasonable quality (orange), and calls of poor quality (red). The central red line is the median quality value, the yellow box represents the inter-quartile of quality, the upper and lower whiskers represent the 10% and 90% points, the blue line represents the mean quality. (a) Sample SwAsp009 of *Populus tremula*. (b) Sample Alb16-1 of *P. tremuloides*. (c) Sample GW-9772 (accession number in SRA: SRR1571518) of *P. trichocarpa*.



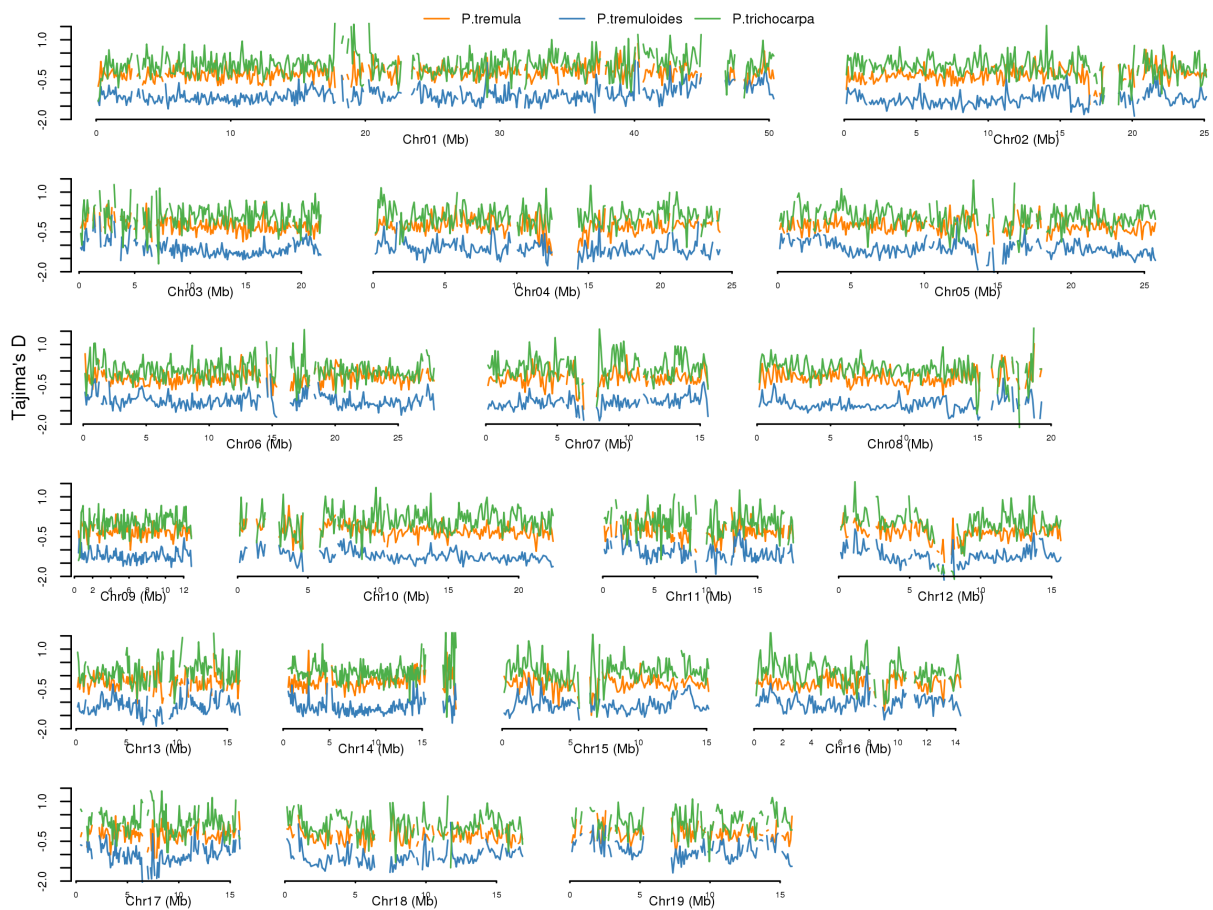
**Figure S4.** Population structure within and between species. (a) Genetic structure of three *Populus* inferred using ADMIXTURE when it identifies three genetic clusters in the dataset. (b) The cross-validation error when  $K$  varies from 1 to 6 across the three species. (c,d,e) The cross-validation error when  $K$  varies from 1 to 6 separately in samples of *P. tremula*, *P. tremuloides*, *P. trichocarpa*.



**Figure S5.** The distributions of estimates of pairwise sequence diversity ( $\Theta_{\pi}$ ) in *P. tremula* (orange), *P. tremuloides* (blue) and *P. trichocarpa* (green) over 1 Mbp non-overlapping windows in different site categories.

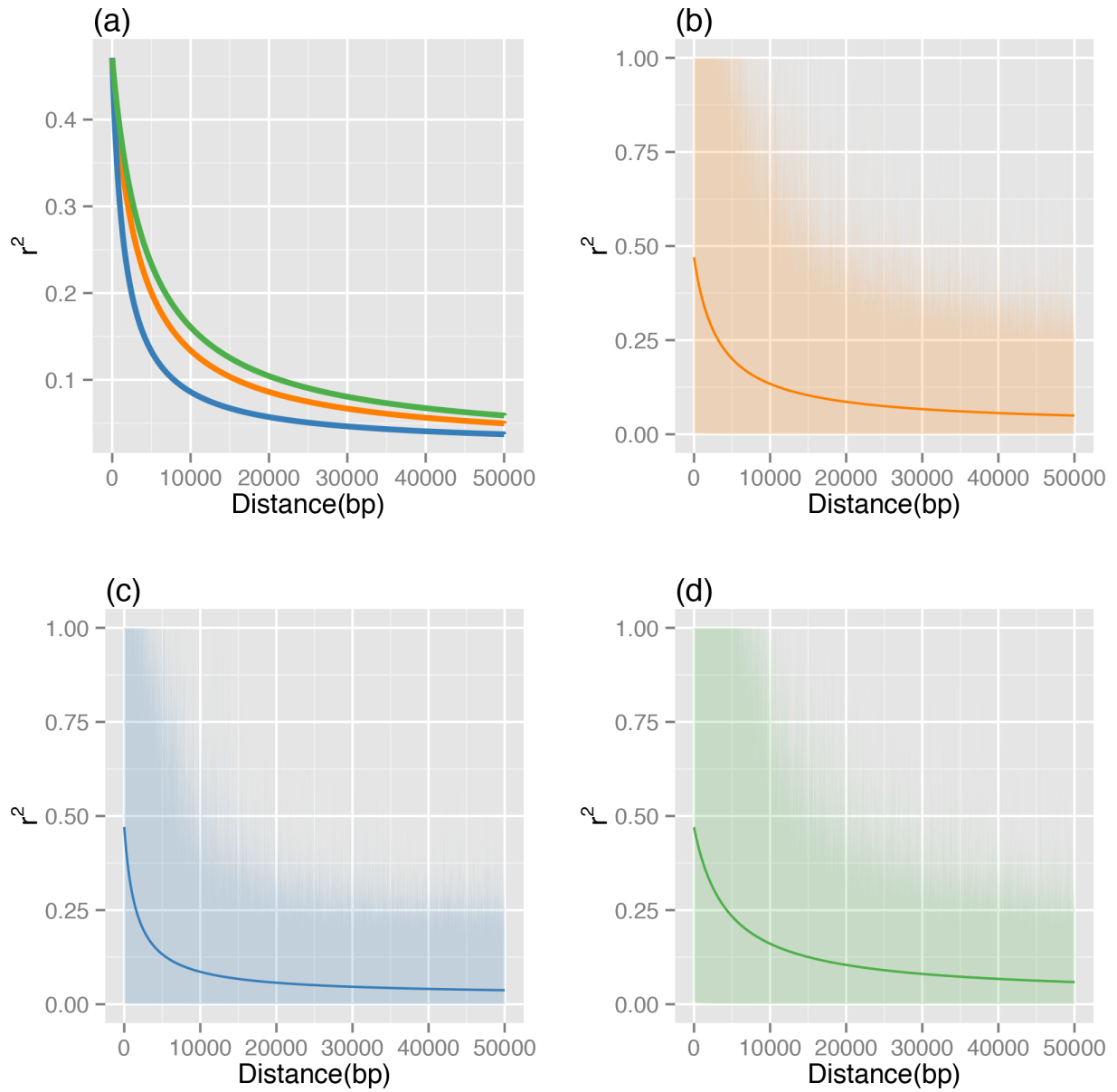


**Figure S6.** The distributions of estimates of Tajima's  $D$  in *P. tremula* (orange), *P. tremuloides* (blue) and *P. trichocarpa* (green) over 1Mbp non-overlapping windows in different site categories.



**Figure S7.** Genome-wide patterns of allele frequency distribution among three *Populus* species. Tajima's D was calculated over 100 Kbp non-overlapping windows in *P. tremula* (orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line) along the 19 chromosomes.

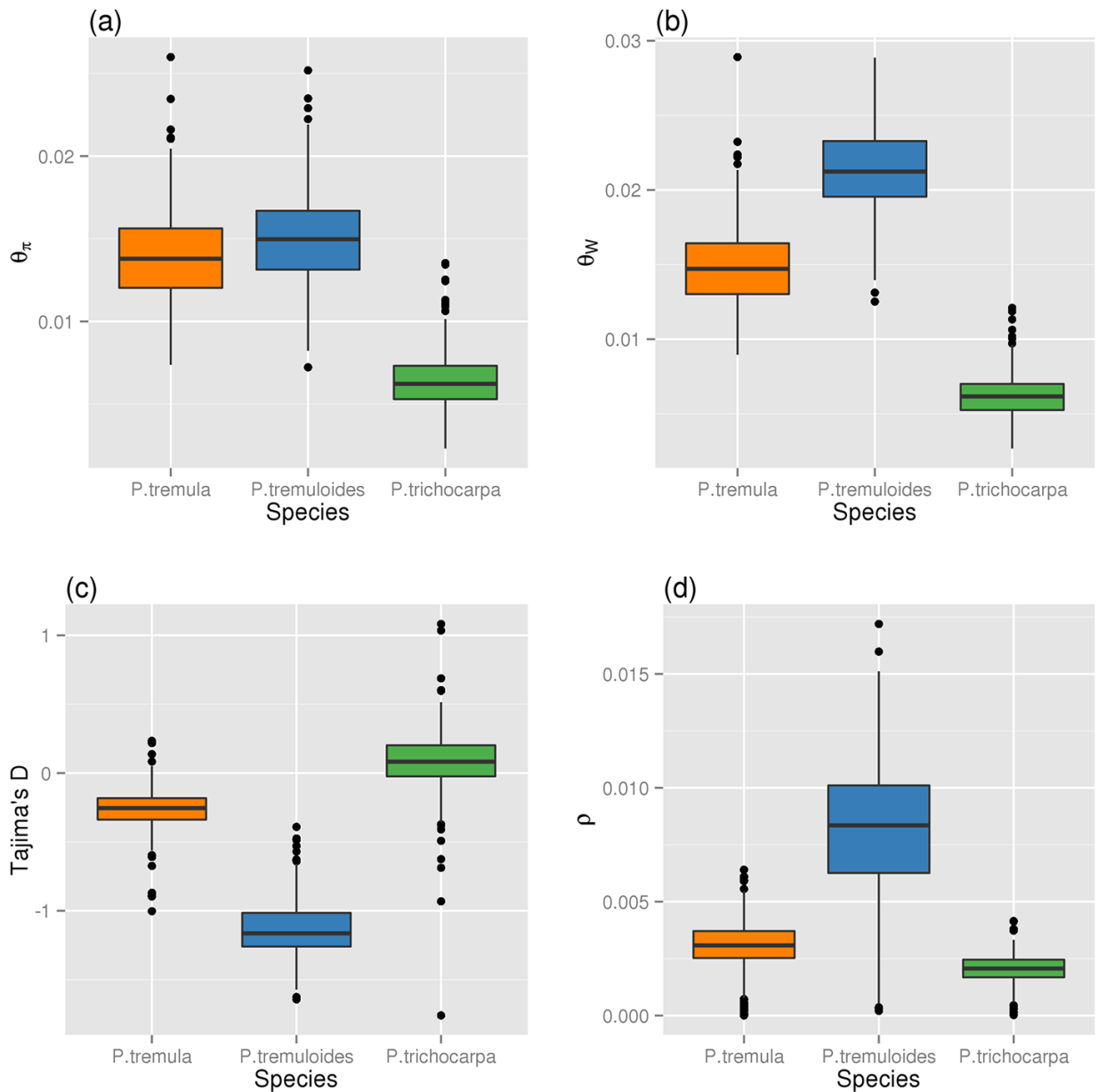




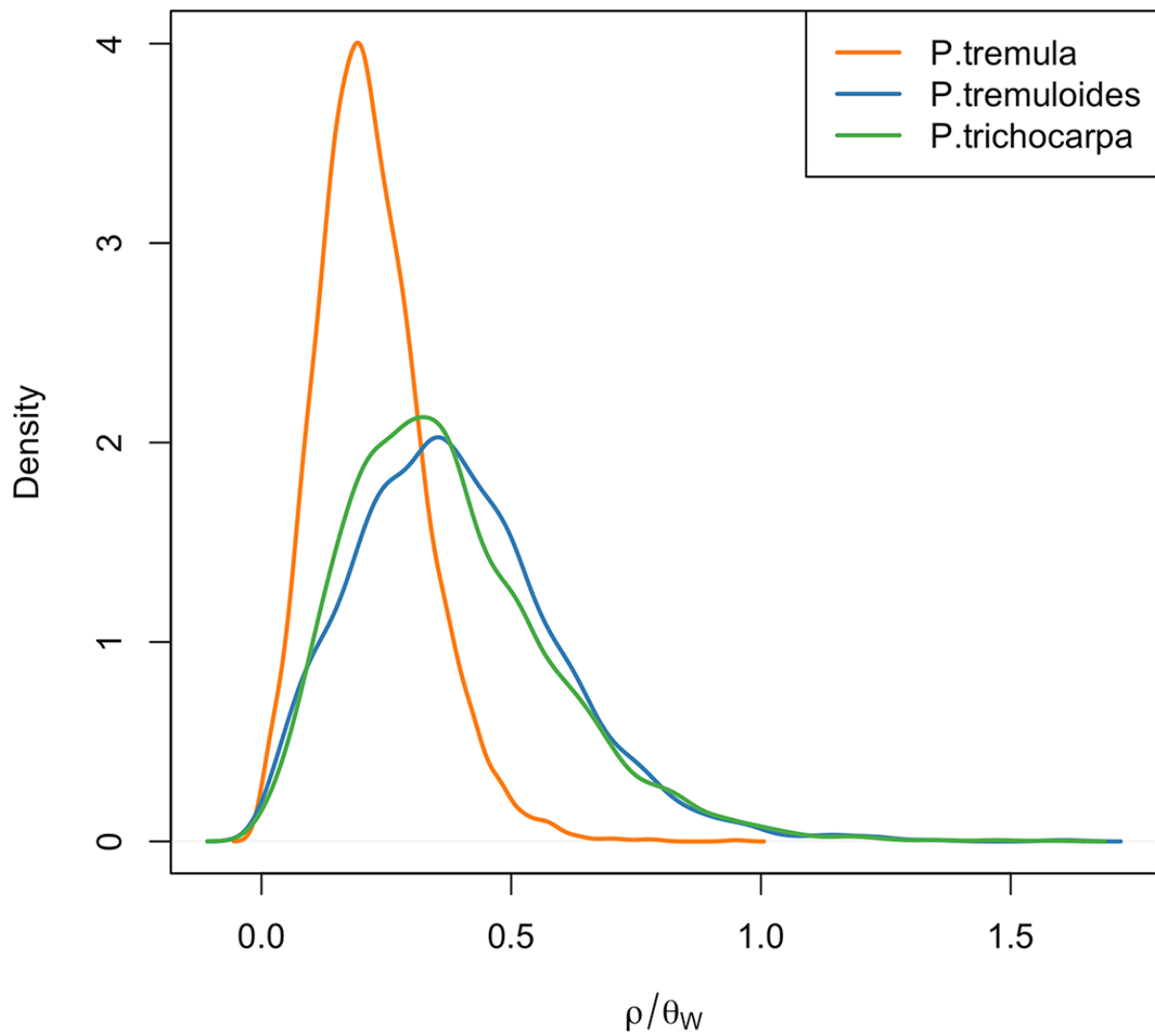
**Figure S8.** Decay of linkage disequilibrium (LD). The comparison of mean LD decay (estimated as  $r^2$ ) with physical distance among the three *Populus* species (a), with the 90% ranges of  $r^2$  values shown in *P. tremula* (orange line) (b), *P. tremuloides* (blue line) (c) and *P. trichocarpa* (green line) (d), respectively.



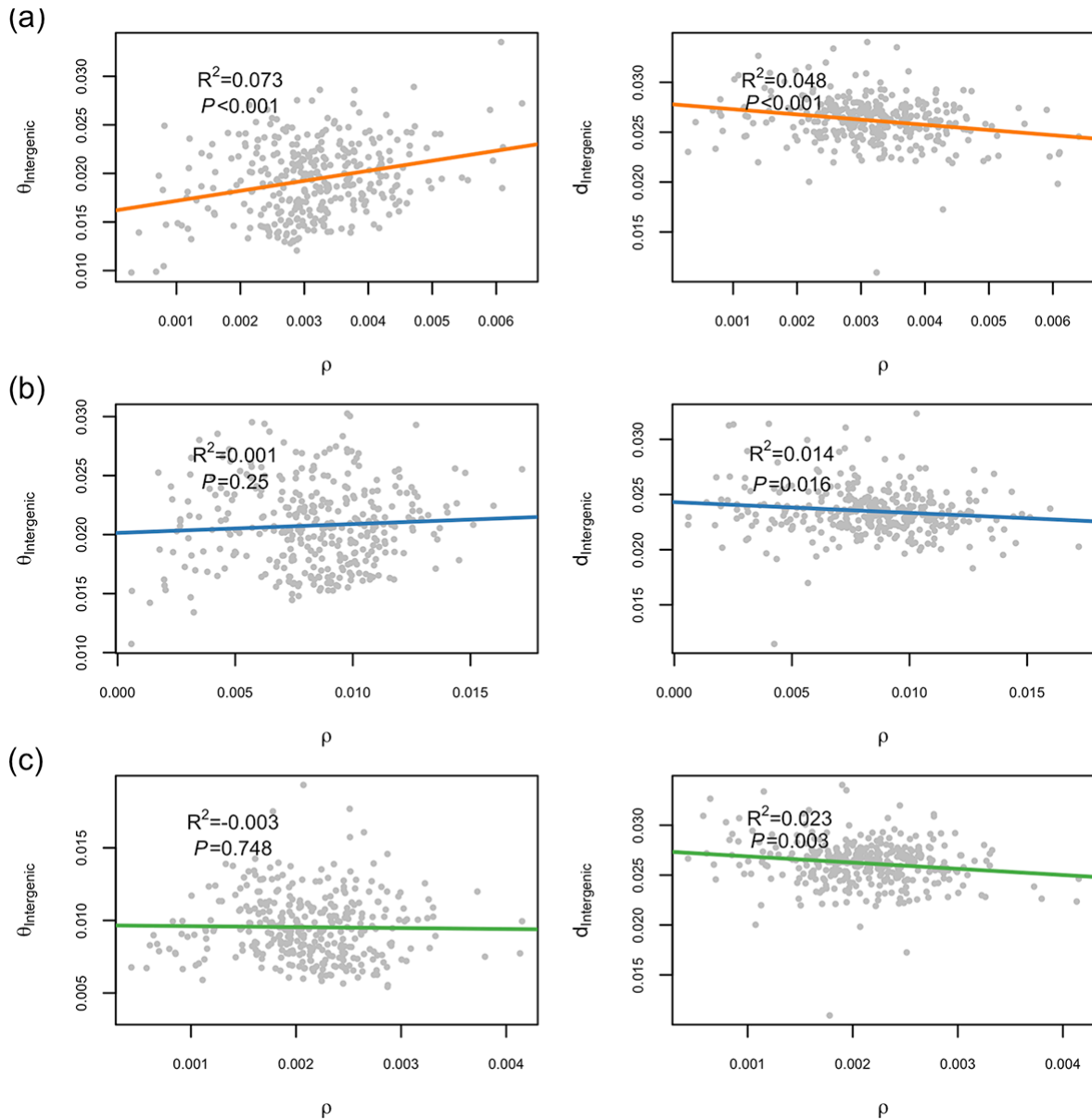
**Figure S9.** Genome-wide patterns of population-scaled recombination rate among three *Populus* species. Population-scaled recombination rate ( $\rho$ ) was averaged over 100 Kbp non-overlapping windows in *P. tremula* (orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line) along the 19 chromosomes.



**Figure S10.** The distributions of estimates of (a) pairwise sequence diversity ( $\Theta_{\pi}$ ), (b) the number of segregating sites ( $\Theta_W$ ), (c) Tajima's D and (d) population-scaled recombination rate ( $\rho$ ) over 1Mbp non-overlapping windows in *P. tremula* (orange), *P. tremuloides* (blue) and *P. trichocarpa* (green).

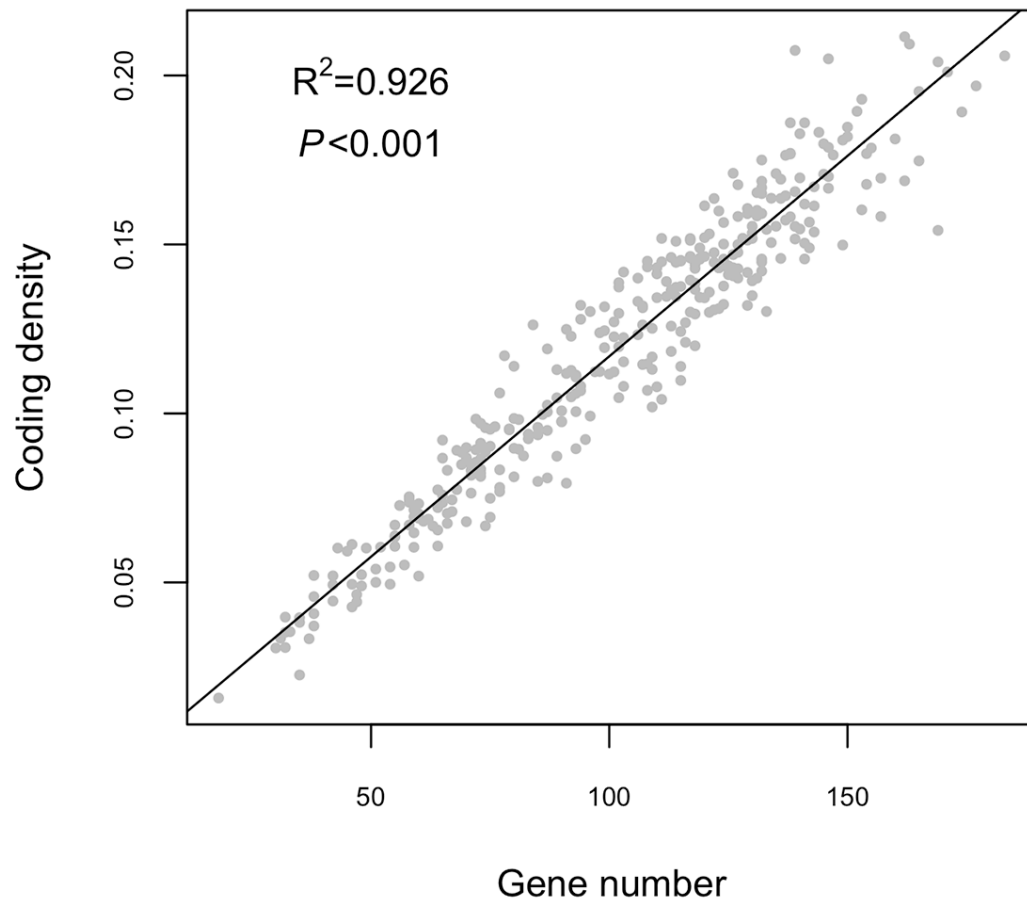


**Figure S11.** Distributions of the ratio of population-scaled recombination rate to nucleotide diversity ( $\rho/\theta_w$ ) over 100 Kbp non-overlapping windows in *P. tremula* (orange line), *P. tremuloides* (blue line) and *P. trichocarpa* (green line).

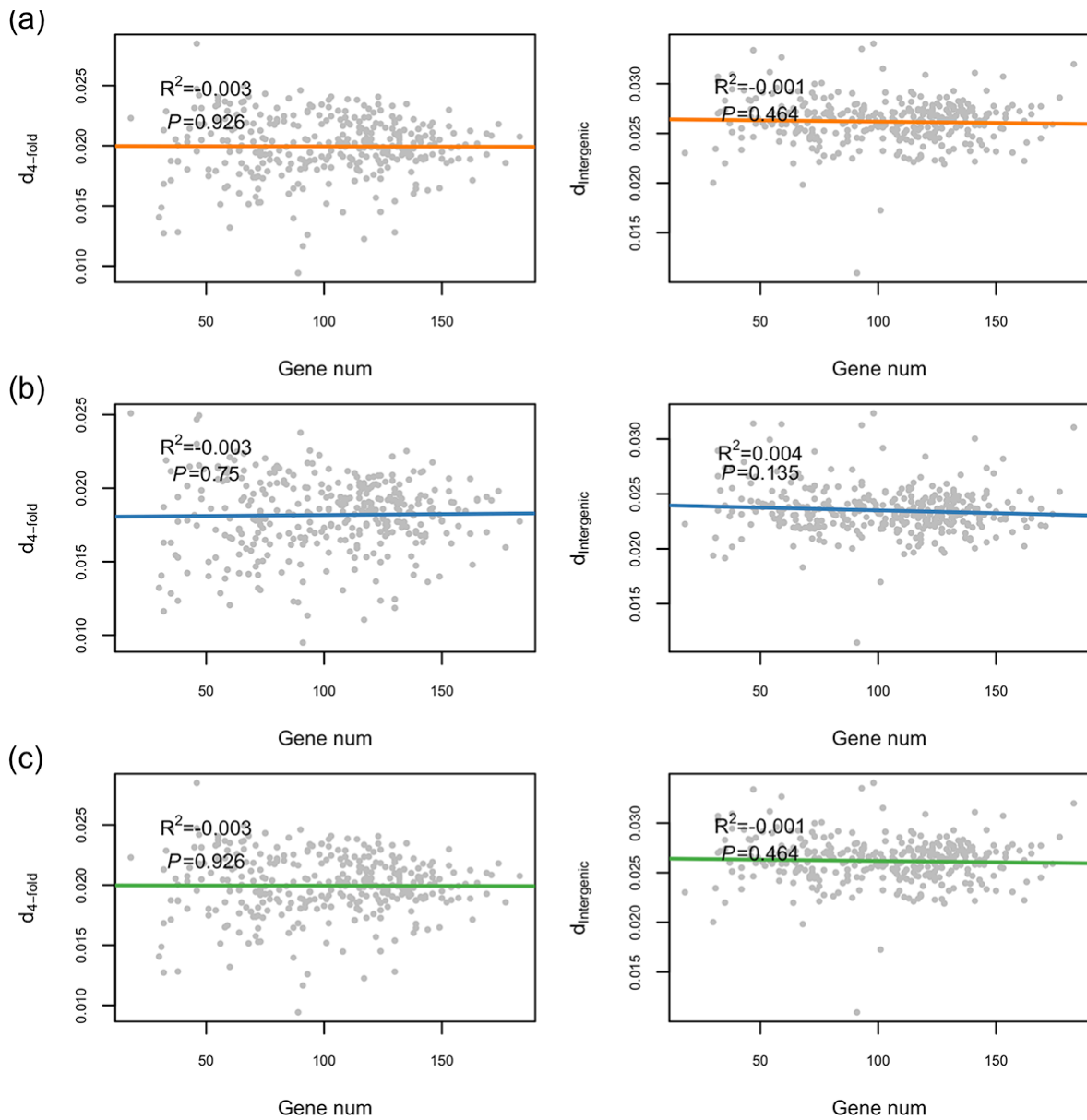


**Figure S12.** Correlations between estimates of intergenic genetic diversity ( $\Theta_{\text{Intergenic}}$ ) (left panel) and divergence ( $d_{\text{Intergenic}}$ ) (right panel) with population-scaled recombination rate ( $\rho$ ) over 1 Mbp non-overlapping windows. Linear regression lines are colored according to species: (a) *P. tremula* (orange line), (b) *P. tremuloides* (blue line) and (c) *P. trichocarpa* (green line).

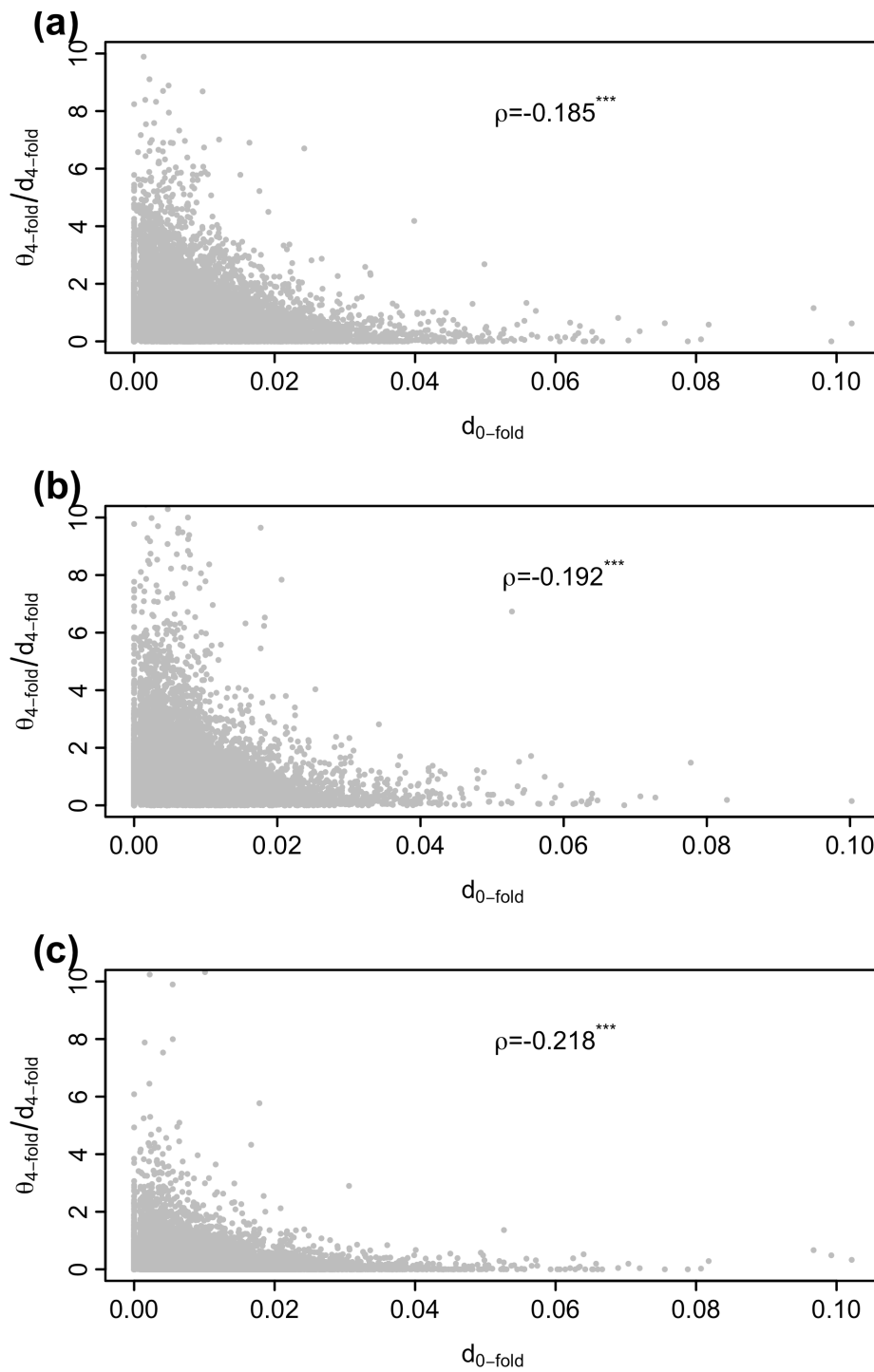




**Figure S13.** Relationship between gene number and the proportion of coding bases over 1Mbp non-overlapping windows.



**Figure S14.** Correlations between estimates of genic and intergenic genetic divergence with gene density over 1 Mbp non-overlapping windows. Correlations between estimates of genetic divergence at 4-fold synonymous sites ( $d_{4\text{-fold}}$ ) (left panel) and intergenic sites ( $d_{\text{Intergenic}}$ ) (right panel) with gene density over 1Mbp non-overlapping windows. Linear regression lines are colored according to species: (a) *P. tremula* (orange line), (b) *P. tremuloides* (blue line) and (c) *P. trichocarpa* (green line).



**Figure S15.** 4-fold synonymous diversity corrected for synonymous divergence ( $\theta_{4\text{-fold}}/d_{4\text{-fold}}$ ) plotted against 0-fold non-synonymous divergence ( $d_{0\text{-fold}}$ ) in *P. tremula* (a), *P. tremuloides* (b) and *P. trichocarpa* (c).

**Table S1.** Samples in this study

<b>SampleID</b>	<b>Site</b>	<b>Latitude</b>	<b>Longitude</b>	<b>Short Read Archive ID</b>
<b><i>P. tremula</i></b>				
SwAsp001	Simlang	56.6925	13.2147	SRR2744682
SwAsp009	Simlang	56.7336	13.2517	SRR2745308
SwAsp011	Ronneby	56.3478	15.025	SRR2745904
SwAsp014	Ronneby	56.3081	15.1269	SRR2745905
SwAsp021	Vargarda	57.9917	12.9119	SRR2745906
SwAsp025	Vargarda	57.9869	12.9358	SRR2745907
SwAsp032	Ydre	57.8492	15.3217	SRR2745908
SwAsp033	Ydre	57.8281	15.3103	SRR2745909
SwAsp045	Brunsborg	59.6425	12.9408	SRR2745910
SwAsp047	Brunsborg	59.6308	12.9608	SRR2745911
SwAsp055	Uppsala	59.8131	17.9817	SRR2745912
SwAsp057	Uppsala	59.7761	17.9889	SRR2745961
SwAsp067	Alvdalen	61.1978	13.8092	SRR2745962
SwAsp068	Alvdalen	61.3017	13.7222	SRR2745964
SwAsp076	Delsbo	61.7106	16.7311	SRR2745965
SwAsp078	Delsbo	61.6925	16.67	SRR2746465
SwAsp087	Dorotea	64.3406	16.3992	SRR2746774
SwAsp088	Dorotea	64.3358	16.3736	SRR2746831
SwAsp096	Umea	63.9781	20.7056	SRR2746685
SwAsp098	Umea	63.8656	20.4986	SRR2746760
SwAsp103	Arjeplog	66.0247	18.5742	SRR2746804
SwAsp110	Arjeplog	66.2592	18	SRR2746848
SwAsp111	Lulea	65.6703	21.8986	SRR2746921
SwAsp114	Lulea	65.5544	22.3939	SRR2746922
<b><i>P. tremuloides</i></b>				
Alb10-3	Alberta	51.0718	-115.0044	SRR2748652
Alb13-1	Alberta	51.0479	-115.0232	SRR2748654
Alb16-1	Alberta	51.0838	-115.3892	SRR2748656
Alb17-4	Alberta	51.0809	-115.3946	SRR2748657
Alb18-3	Alberta	51.0686	-115.3516	SRR2748658
Alb25-4	Alberta	51.0524	-114.9131	SRR2748659
Alb27-1	Alberta	51.0405	-114.8939	SRR2748660
Alb31-1	Alberta	51.0435	-114.8352	SRR2748661
Alb33-2	Alberta	51.0431	-114.7568	SRR2748662
Alb35-2	Alberta	51.0234	-115.0640	SRR2748693
Alb6-3	Alberta	51.1324	-115.0664	SRR2749823
Albb15-3	Alberta	51.0811	-115.3767	SRR2749863
Dan1-	UW Arboretum	43.0520	-89.4242	SRR2749866

1C13

Dan2-1B7	UW Arboretum	43.0526	-89.4253	SRR2749867
PG1-1B4	Parfrey's Glenn	43.4249	-89.6445	SRR2749956
PG2-1B9	Parfrey's Glenn	43.4184	-89.6417	SRR2751048
PG3-1B6	Parfrey's Glenn	43.4198	-89.6532	SRR2751050
PI3-1B3	Pine Island area	43.5402	-89.5666	SRR2751053
Sau1-1B10	Boxter's Hollow	43.4036	-89.8176	SRR2751064
Sau2-1B2	Boxter's Hollow	43.4048	-89.8243	SRR2751068
Sau3-1B13	Boxter's Hollow	43.4053	-89.8141	SRR2751087
Wau1-1B5	Waushara country	44.1314	-89.2082	SRR2751102
<b><i>P. trichocarpa</i></b>				
BESC-56	Talley_Way	46.099	-122.878	SRR1571263
BESC-108	Rainier	46.114	-122.99	SRR1571274
BESC-264	Monroe	47.842	-121.979	SRR1571343
BESC-281	Monroe	47.851	-121.962	SRR1571362
BESC-374	Skiou_Island	48.489	-122.16	SRR1571397
BESC-840	Orting	47.042	-122.209	SRR1571416
BESC-873	Sultan	47.856	-121.811	SRR1571152
BESC-884	Gold_Bar	47.84	-121.691	SRR1571432
DEND-17- 2	DEND	52.817	-126.95	SRR1569629
FNYI-28-4	Fanny_Bay_	52.817	-126.95	SRR1569762;SRR1569763
GW-11031	Corvallis	52.817	-126.95	SRR1569519;SRR1569520
GW-7983	Longview	46.117	-123	SRR1571215
GW-9595	Turner	44.75	-122.867	SRR1571497
GW-9598	Nisqually_River	47.067	-123.733	SRR1571500
GW-9772	Acme	48.717	-122.2	SRR1571518
GW-9792	Sedro_Woolley	48.717	-122.2	SRR1571532
GW-9920	Orting	47.1	-122.2	SRR1571552
GW-9959	Skamania	45.95	-121.95	SRR1571572
HARC-26- 1	Harrison	49.767	-122.217	SRR1571016
HOMC- 21-3	Homathko	51.233	-124.95	SRR1569781
LILC-26-3	Harrison	51.233	-124.95	SRR1569814
SLMB-28- 4	Salmon	50.217	-125.817	SRR1570762
SQMB-25- 4	Squamish	50.217	-125.817	SRR1571038
WHTE- 28-1	Salmon	50.133	-126.05	SRR1570189

---



**Table S2.** Summary statistics of Illumina re-sequencing data per sample

<b>SampleID</b>	<b>Raw bases (Gb)</b>	<b>Filtered bases (Gb)</b>	<b>Uniquely mapped bases (Gb)</b>	<b>Mapping rate (%)</b>	<b>Mean Coverage</b>	<b>Proportion of covered genome</b>
<i>Populus tremula</i>						
SwAsp001	21.39	12.02	8.90	91.34%	20.50	74.44%
SwAsp009	16.95	13.73	10.29	90.52%	23.70	75.32%
SwAsp011	34.59	29.08	21.09	90.06%	48.58	78.50%
SwAsp014	18.48	14.46	10.65	88.78%	24.53	75.82%
SwAsp021	13.80	10.73	8.03	89.54%	18.50	74.51%
SwAsp025	13.65	10.08	7.49	88.62%	17.25	74.27%
SwAsp032	15.23	12.00	8.82	88.56%	20.32	75.25%
SwAsp033	15.78	11.89	8.73	90.27%	20.11	74.53%
SwAsp045	15.31	11.85	8.79	88.90%	20.25	74.93%
SwAsp047	20.77	11.83	8.45	89.72%	19.46	75.30%
SwAsp055	17.01	13.14	9.66	88.66%	22.25	75.43%
SwAsp057	18.44	14.12	10.37	88.43%	23.88	75.95%
SwAsp067	16.52	13.05	9.66	88.56%	22.25	75.13%
SwAsp068	17.18	13.47	9.67	88.78%	22.27	75.39%
SwAsp076	23.91	21.32	15.95	90.34%	36.74	76.63%
SwAsp078	15.61	12.54	9.12	88.88%	21.01	75.34%
SwAsp087	17.27	14.21	10.61	90.00%	24.44	75.36%
SwAsp088	16.47	13.51	10.08	89.94%	23.22	74.94%
SwAsp096	16.99	13.66	10.17	90.48%	23.43	75.01%
SwAsp098	15.94	14.14	10.68	90.27%	24.60	74.60%
SwAsp103	14.59	12.76	9.65	91.27%	22.23	74.52%
SwAsp110	24.49	21.60	15.99	90.62%	36.83	77.82%
SwAsp111	20.95	18.23	13.63	90.68%	31.39	76.20%
SwAsp114	22.96	20.22	15.03	91.25%	34.62	75.94%
Mean	18.51	14.74	10.90	89.77%	25.10	75.46%
<i>P. tremuloides</i>						
Alb10-3	14.99	12.66	9.15	90.47%	21.08	75.75%
Alb13-1	16.04	13.46	10.08	91.66%	23.22	75.02%
Alb16-1	13.94	11.45	8.52	90.75%	19.63	74.70%
Alb17-4	12.54	10.41	7.87	90.82%	18.13	74.51%
Alb18-3	14.60	12.38	9.33	91.20%	21.49	75.47%
Alb25-4	12.99	10.62	7.94	90.32%	18.29	74.40%
Alb27-1	20.07	17.27	12.60	91.50%	29.02	76.46%
Alb31-1	16.29	13.38	10.01	91.09%	23.06	75.36%
Alb33-2	19.44	16.09	12.05	90.91%	27.76	76.25%
Alb35-2	12.57	10.30	7.76	90.58%	17.87	74.41%
Alb6-3	19.47	16.67	12.49	91.92%	28.77	77.19%
Albb15-3	15.14	12.77	9.61	91.79%	22.14	75.05%

Dan1-1C13	16.42	14.23	10.47	91.08%	24.12	76.04%
Dan2-1B7	22.00	18.48	13.29	90.95%	30.61	76.94%
PG1-1B4	12.11	10.68	8.10	91.59%	18.66	74.70%
PG2-1B9	16.53	14.13	10.48	91.65%	24.14	76.11%
PG3-1B6	14.93	12.78	9.44	91.27%	21.74	75.76%
PI3-1B3	16.75	13.37	10.04	90.99%	23.13	75.87%
Sau1-1B10	13.30	10.61	7.81	90.05%	17.99	74.89%
Sau2-1B2	13.74	6.82	5.12	90.79%	11.79	73.55%
Sau3-1B13	12.13	9.86	7.16	91.25%	16.49	74.43%
Wau1-1B5	26.40	21.73	15.91	91.01%	36.65	77.67%
Mean	16.02	13.19	9.78	91.07%	22.53	75.48%
Total	796.67	643.79	476.74			
<b><i>P. trichocarpa</i></b>						
BESC-56	14.51	12.44	11.43	97.06%	25.43	94.08%
BESC-108	15.97	13.58	12.39	94.99%	27.75	93.75%
BESC-264	19.78	17.03	15.42	95.22%	34.35	94.01%
BESC-281	12.21	10.83	10.13	97.31%	22.89	93.48%
BESC-374	19.63	16.48	14.94	95.33%	33.64	94.18%
BESC-840	16.74	13.38	12.08	94.53%	27.56	94.03%
BESC-873	13.17	11.76	10.9	97.20%	24.36	93.68%
BESC-884	14.32	12.46	11.52	96.66%	25.89	93.83%
DEND-17-2	8.33	7.32	6.73	97.14%	15.26	93.08%
FNYI-28-4	15.71	14.14	12.79	96.58%	28.29	93.82%
GW-11031	19.01	16.86	15.24	96.35%	34.23	94.28%
GW-7983	10.68	9.58	8.91	97.18%	20.14	93.22%
GW-9595	16.17	13.77	12.12	93.13%	27.18	93.76%
GW-9598	12.36	11.09	9.67	91.74%	21.84	93.69%
GW-9772	15.55	13.35	12.12	96.31%	26.9	93.89%
GW-9792	14.53	12.84	11.78	96.88%	26.51	93.66%
GW-9920	12.67	11.1	9.36	88.71%	21.23	93.59%
GW-9959	14.41	9.75	8.3	89.35%	18.64	93.71%
HARC-26-1	12.07	10.2	9.56	97.41%	21.76	93.63%
HOMC-21-3	18.36	17.15	15.32	97.21%	34.05	93.63%
LILC-26-3	16.83	14.29	12.74	95.11%	28.53	93.89%
SLMB-28-4	13.43	11.06	9.84	92.43%	22.33	93.81%
SQMB-25-4	16	13.73	12.52	95.03%	28.03	94.03%
WHTE-28-1	12.75	11.24	10.28	96.69%	23.13	93.38%
Mean	14.80	12.73	11.50	95.23%	25.83	93.75%

**Table S3.** Diversity statistics (median and central 95% range) for various genomic contexts over 100 Kbp non-overlapping windows across genome

	Filtered bases (Mbp)	<i>P. tremula</i>		<i>P. tremuloides</i>		<i>P. trichocarpa</i>	
		$\Theta_{\pi}$	Tajima's D	$\Theta_{\pi}$	Tajima's D	$\Theta_{\pi}$	Tajima's D
<b>Total</b>	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)
<b>0-fold<sup>a</sup></b>	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)
<b>4-fold</b>	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)
<b>Introns<sup>a</sup></b>	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)
<b>UTR 5'<sup>a</sup></b>	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)
<b>UTR 3'<sup>a</sup></b>	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)
<b>Intergenic<sup>a</sup></b>	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)

<sup>a</sup> One-sided paired Mann-Whitney U test in comparison to the 4-fold synonymous sites

\*  $P < 0.05$

\*\*  $P < 0.001$

\*\*\*  $P < 2.2 \times 10^{-16}$

**Table S4.** Estimates of the distribution of fitness effects of new amino acid mutations falling in different  $N_e$ s ranges and associated demographic parameters, and proportion of amino acid substitution driven to fixation by positive selection ( $\alpha$ ) and the rate of adaptive substitution relative to neutral divergence ( $\omega$ ). 95% CIs are shown in parentheses.

Species	N2/N1	t/N2	$\beta$	Percentage of mutations in $N_e$ s range			$\alpha$	$\omega$
				0-1	1-10	>10		
<i>P. tremula</i>	0.12	34.776	0.227	0.233	0.160	0.607	0.427	0.155
	(0.12-0.12)	(34.730-47.278)	(0.223-0.231)	(0.231-0.235)	(0.157-0.162)	(0.594-0.620)	(0.419-0.435)	(0.151-0.159)
<i>P. tremuloides</i>	0.31	33.945	0.507	0.155	0.314	0.532	0.653	0.239
	(0.31-0.34)	(33.930-34.128)	(0.481-0.515)	(0.153-0.161)	(0.302-0.318)	(0.522-0.545)	(0.636-0.658)	(0.231-0.242)
<i>P. trichocarpa</i>	0.37	9.369	0.116	0.310	0.094	0.596	0.203	0.073
	(0.31-0.67)	(0.179-37.531)	(0.109-0.181)	(0.276-0.313)	(0.089-0.144)	(0.478-0.697)	(0.188-0.311)	(0.068-0.112)