

# Population genomics of *Populus trichocarpa* identifies signatures of selection and adaptive trait associations

Luke M Evans<sup>1</sup>, Gancho T Slavov<sup>2</sup>, Eli Rodgers-Melnick<sup>1</sup>, Joel Martin<sup>3</sup>, Priya Ranjan<sup>4</sup>, Wellington Muchero<sup>4</sup>, Amy M Brunner<sup>5</sup>, Wendy Schackwitz<sup>3</sup>, Lee Gunter<sup>4</sup>, Jin-Gui Chen<sup>4</sup>, Gerald A Tuskan<sup>3,4</sup> & Stephen P DiFazio<sup>1</sup>

**Forest trees are dominant components of terrestrial ecosystems that have global ecological and economic importance. Despite distributions that span wide environmental gradients, many tree populations are locally adapted, and mechanisms underlying this adaptation are poorly understood. Here we use a combination of whole-genome selection scans and association analyses of 544 *Populus trichocarpa* trees to reveal genomic bases of adaptive variation across a wide latitudinal range. Three hundred ninety-seven genomic regions showed evidence of recent positive and/or divergent selection and enrichment for associations with adaptive traits that also displayed patterns consistent with natural selection. These regions also provide unexpected insights into the evolutionary dynamics of duplicated genes and their roles in adaptive trait variation.**

A suite of forces and factors, including mutation, recombination, selection, population history and gene duplication influence patterns of intraspecific genetic variation. Distinguishing which factors have shaped sequence variation across a genome requires extensive whole-genome sequencing of multiple individuals, which has only recently become tractable<sup>1</sup>. Most large-scale whole-genome resequencing studies have focused on model and domesticated species<sup>1–5</sup>. However, extensive sequencing of natural populations holds great promise for advancing understanding of evolutionary biology, including identifying functional variation and the molecular bases of adaptation. Recent work in a number of species has identified genomic regions that show signatures of positive selection, suggesting that such regions contain loci that control adaptive traits<sup>4,6–8</sup>. Relatively few studies, however, have combined genome-wide scans with phenotypic data to determine whether computationally identified selected regions influence adaptive phenotypic variation<sup>5,9–13</sup>. Genome-wide studies of large natural populations combined with phenotypic measurements are necessary to determine which factors shape patterns of genetic variation within species and, therefore, enhance understanding of adaptation.

With large geographic ranges spanning wide environmental gradients and a long history of research showing local adaptation<sup>14</sup>, forest trees are ideal for examining the processes shaping genetic variation in natural populations. Forest trees cover approximately 30% of terrestrial land area<sup>15</sup>, provide direct feedback to global climate<sup>15</sup> and are often foundation species that organize entire biotic communities and biogeochemical systems<sup>16,17</sup>. Clearly, biotic and abiotic interactions have influenced population sizes and distributions of forest trees, leaving diagnostic signatures in the genomes of present-day populations<sup>14,18,19</sup>. A deeper understanding of the evolutionary and

ecological forces that shaped these patterns will offer insights and options for ecosystem management, applied tree improvement and accelerated domestication efforts<sup>20</sup>.

Black cottonwood, *Populus trichocarpa* Torr. & Gray, is a dominant riparian tree that has become a model for the advancement of genome-level insights in forest trees<sup>21</sup>. The sequencing of 16 *P. trichocarpa* genomes revealed widespread patterns of linkage disequilibrium (LD) and population structure<sup>22</sup> and extensive genecological studies have revealed a high degree of adaptive phenotypic variation in growth, vegetative phenology and physiological traits such as water-use efficiency and photosynthesis<sup>23–25</sup>, suggesting that local adaptation is prevalent. To date, candidate gene–association analyses have revealed loci with significant effects on phenotypic traits<sup>26,27</sup>. However, thus far there have been no publications describing whole-genome associations for adaptive traits in *P. trichocarpa*, or their relationship to signatures of selection in any forest tree species.

One of the salient features of the *P. trichocarpa* genome is a remarkably well-conserved whole-genome duplication that is shared by all members of the Salicaceae and near relatives: the Salicoid duplication<sup>28,29</sup>. Despite the extensive occurrence of segments of collinear paralogous genes, more than two-thirds of the duplicated genes are lost after the duplication event, and there are substantial functional biases in the remaining gene pairs—in particular, there is an overabundance of gene categories with large numbers of protein–protein interactions<sup>30,31</sup>. A major unexplored question is whether the fundamental, diagnostic differences in diversity between retained duplicate pairs and genes lacking paralogs from the Salicoid duplication (singletons) are connected to patterns of natural selection and adaptive phenotypic variation.

<sup>1</sup>Department of Biology, West Virginia University, Morgantown, West Virginia, USA. <sup>2</sup>Institute of Biological, Environmental and Rural Sciences, Aberystwyth University, Aberystwyth, UK. <sup>3</sup>The Joint Genome Institute, Walnut Creek, California, USA. <sup>4</sup>Plant Systems Biology Group, BioSciences Division, Oak Ridge National Laboratory, Oak Ridge, Tennessee, USA. <sup>5</sup>Department of Forest Resources and Environmental Conservation, Virginia Tech, Blacksburg, Virginia, USA. Correspondence should be addressed to S.P.D. (spdifazio@mail.wvu.edu).

Received 22 May; accepted 30 July; published online 24 August 2014; doi:10.1038/ng.3075

Here we report the whole-genome resequencing of a collection of 544 *P. trichocarpa* individuals, spanning much of the species' natural latitudinal range, that have been clonally replicated in three contrasting environments. We use this resource to detect signatures of recent selection across the *Populus* genome and on adaptive traits themselves. We also show that the signals of association with adaptive traits are stronger in positively selected regions. Finally, we demonstrate that Salicoid duplicate genes have distinctive patterns of adaptive variation that reveal the evolutionary effects of dosage constraints.

## RESULTS

### Polymorphism and population structure

From high-coverage whole-genome sequencing of 544 unrelated *P. trichocarpa* individuals (Fig. 1a and Supplementary Table 1) we collected more than 3.2 terabases (Tb) of data that aligned to 394 Mb of the *P. trichocarpa* genome. Approximately 87.5% of the 3.2 Tb was accessible for analysis based on median sequencing depth across all samples (Supplementary Fig. 1). From these data, we detected 17,902,740 single nucleotide polymorphisms (SNPs).

Using this resource, nucleotide diversity was twofold higher in intergenic sequence than in genic sequence, largely consistent with purifying selection (Table 1). Diversity was particularly low in coding sequences, where nonsynonymous diversity was one-third that of synonymous diversity. Most SNPs were rare (minor allele frequency (MAF)  $\leq 0.01$ ), particularly those predicted to have major effects (for example, splice site mutations) (Table 1 and Supplementary Fig. 2). We also identified 5,660 large (>100 bp) and 254,464 small (<50 bp) insertion or deletion (indel) polymorphisms (unpublished data).

On the basis of principal component analysis (PCA) of these 17.9 million SNPs, we identified four major regional genetic groups corresponding to geographical origin (Fig. 1a). We also found spatial genetic structure within regional groupings that clustered as separate subgroups within source locations (Fig. 1b). These data indicate that there is genome-wide genetic structure at both broad latitudinal and local spatial scales.

### Phenotypic evidence of selection

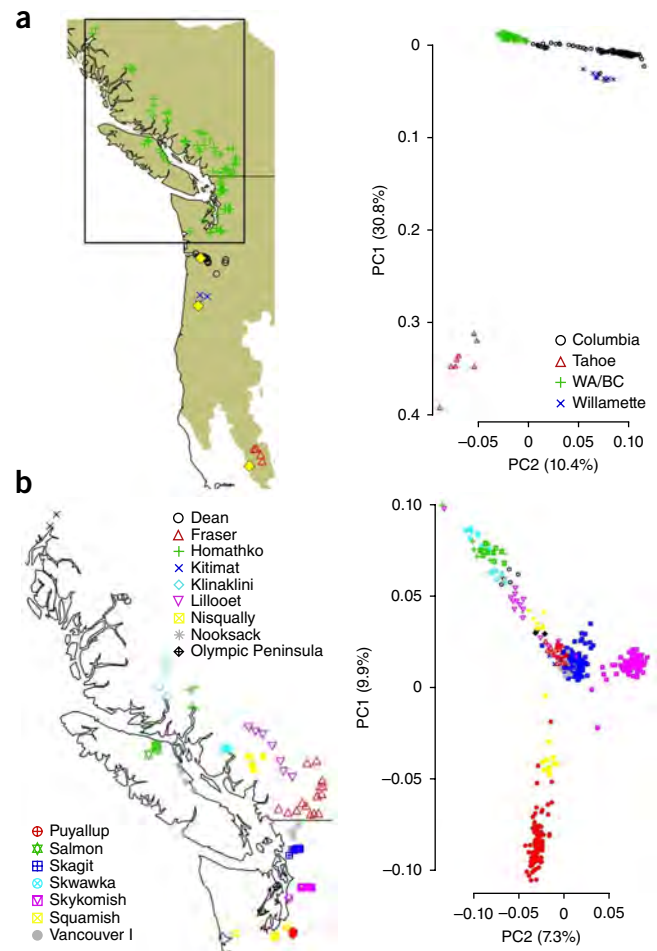
We examined two different indicators of selection using phenotypic data from three clonally replicated plantations representing the center and southern extent of the extant range of *P. trichocarpa*. We found that quantitative differentiation ( $Q_{ST}$ ) in height, spring bud flush and fall bud set among source rivers was greater than genome-wide marker differentiation ( $F_{ST}$ ) (Fig. 2a), consistent with spatially divergent selection<sup>32</sup>, as is commonly observed in forest trees<sup>14,24,25</sup>. Furthermore, at all three plantations, these same adaptive traits showed correlations with the first two principal components (PC) axes of multivariate climate variables (Fig. 2b–d and Supplementary Fig. 3). Warmer climates (negative PC1) are associated with earlier bud flush and later bud set, strongly supporting the hypothesis that climate is a major determinant of adaptive genetic variation throughout the sampled range of *P. trichocarpa*<sup>24,25</sup>.

**Figure 1** Geographic locations and genetic structure of the 544 *P. trichocarpa* individuals sequenced. (a) Map (left) of collection locations of the 544 *P. trichocarpa* genotypes sampled in this study from along the Northwest coast of North America (tan shading indicates species range) and PCA (right) of all 544 individuals color-coded by general geographic regions. Yellow diamonds represent plantation locations. (b) Map (left) and PCA (right) of the central Washington and British Columbia (WA/BC) group of individuals (outlined by box in a), color-coded according to collection river. The percentage of the variance explained by the PC1 and PC2 axes for the regional analysis and the WA/BC group is shown.

### Recent positive and divergent selection

We related the strong evidence of climate-driven, divergent selection on adaptive traits to genomic regions that also appeared to be affected by natural selection. We examined five distinct metrics of natural selection using 1-kb windows across the genome. These metrics included differentiation ( $F_{ST}$ ), allele frequency cline steepness across mean annual temperature and precipitation measurements (SPA)<sup>33</sup>, extended haplotype homozygosity around alleles from rapid allele frequency increase (iHS)<sup>8</sup> and allele frequency clines (bayenv)<sup>34</sup> with each of the first two climate PC axes (PC1 and PC2). From these data we classified the empirical top 1% of windows or regions as 'selection outliers', i.e., regions with unusually strong polymorphism patterns consistent with recent positive and/or divergent selection (Fig. 3, Supplementary Figs. 4 and 5 and Supplementary Tables 2–6). Most of the selection outlier regions occurred uniquely among selection scan metrics, suggesting that each metric provides a distinct view of selection and that different selective forces are shaping these genomic regions (Fig. 3a). However, we found 397 regions in the top 1% for at least two of the selection scan metrics; we termed these regions candidate selection regions (CSRs) (Supplementary Table 7).

We tested whether the genes spanning or nearest to these CSRs (452 genes) and the selection outliers (1,418; 1,718; 1,151; 257 and 312 genes for  $F_{ST}$ , SPA, iHS, bayenv PC1 and bayenv PC2, respectively) were overrepresented among annotation categories, gene families or genes with known involvement in several biological processes (Fig. 3 and Supplementary Tables 8–11). On the basis of Fisher exact tests,



**Table 1** Per-site nucleotide diversity estimated across the genome for all annotated features of the *P. trichocarpa* genome and the number of variants annotated in each class

Feature	$\pi$ (median and central 95% range)
Overall <sup>a</sup>	0.0041 (0.0004–0.01226)
Intergenic	0.0064 (0.0012–0.0125)
Genic <sup>b</sup>	0.003 (0.0006–0.0106)
5' UTR	0.0028 (0.0001–0.0114)
3' UTR	0.0033 (0.0001–0.0123)
Intron	0.0034 (0.0005–0.0114)
Coding sequence	0.002 (0.0002–0.0111)
Nonsynonymous	0.0018 (0–0.0122)
Synonymous	0.0054 (0–0.0348)
$\pi_{\text{nonsynonymous}}/\pi_{\text{synonymous}}$	0.3179 (0–14.5447)
Annotation	Number of variants <sup>c</sup>
Intergenic	14,520,224
Intron	1,962,848
Nonsynonymous coding	612,655
Nonsynonymous start	253
Start lost	1,631
Stop gained	18,702
Stop lost	2,175
Splice site acceptor	3,748
Splice site donor	4,449
Synonymous coding	386,103
Synonymous stop	959
3' UTR	389,771
5' UTR	169,083

Total is greater than total observed number of variants because some SNPs have multiple annotations for alternative transcripts.

<sup>a</sup>Based on *P. trichocarpa* version 3.0 reference genome. <sup>b</sup>Predicted transcript from 5' to 3' UTR. <sup>c</sup>Variants annotated using SnpEff<sup>60</sup>.

certain functional categories were overrepresented, including GO annotations related to response to stimuli; 1,3- $\beta$ -glucan (callose) synthesis; and metabolic processes, as well as PANTHER annotations for leucine-rich repeat receptor-like protein kinase and homeobox protein transcription factors (**Supplementary Tables 8–10**).

Despite some similarities, genes associated with the top 1% of each scan were generally overrepresented in unique categories (**Fig. 3**). For example, genes encoding transcription factors were, as a group, overrepresented among  $F_{ST}$  and SPA outliers; those encoding DELLA proteins (gibberellin-interacting transcriptional regulators; Pfam PF12041), among  $F_{ST}$  and bayenv PC2; and genes encoding phytochromes (Pfam PF00360) or involved in photoperiodic or circadian clock regulation, ATPase activity and transmembrane movement (for example, GO:0042626) were only overrepresented in  $F_{ST}$  (**Supplementary Tables 8 and 9**). Heat shock-related annotations were significantly overrepresented only in SPA (PANTHER PTHR10015 and PTHR11528), and genes encoding proteins induced by water stress or abscisic acid (Pfam PF02496) were overrepresented

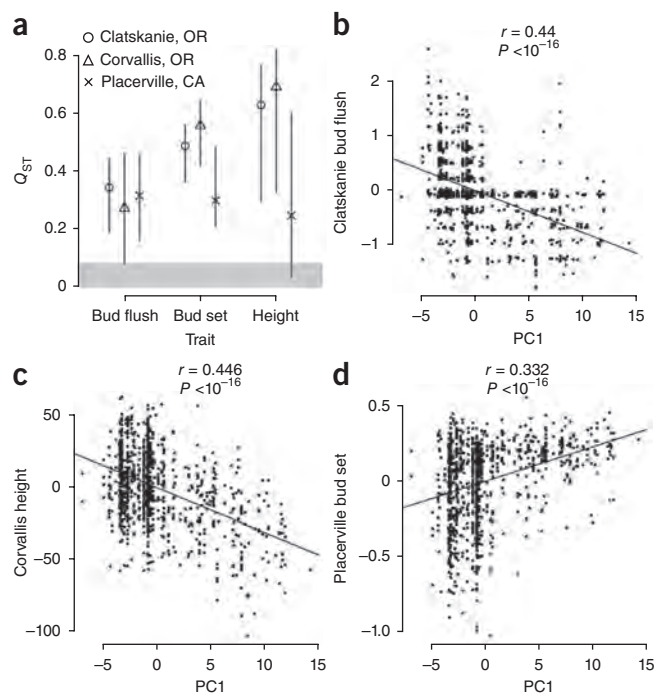
in bayenv PC2 and SPA outliers. Genes encoding the hydrolase 4-nitrophenylphosphatase (PTHR19288) were overrepresented among bayenv PC1 and weakly represented in  $F_{ST}$  (**Supplementary Table 9**). Genes encoding class III aminotransferases (PANTHER PTHR11986, involved in abiotic stress)<sup>35</sup> were overrepresented most strongly in bayenv PC2 (**Fig. 3**).

Intriguingly, although moderate-effect SNPs were underrepresented among genic regions of all selection scan outliers, owing presumably to purifying selection, SNPs predicted to have high impacts were overrepresented among strong sweep loci implicated by the iHS scans (**Supplementary Table 12**), potentially because SNPs with major, presumably beneficial effects are more likely to be swept to high frequency. Because different selection processes (for example, hard sweeps and subtle frequency shifts of standing variation) will influence diversity patterns differently, these five metrics reveal an assortment of potential selection pressures acting on *P. trichocarpa* through the largely nonoverlapping regions identified in each.

### Adaptive trait associations in candidate selected regions

If climate is a major force driving the signatures of positive selection, we predict polymorphisms in regions affected by selection to be associated with climate-related adaptive traits. Vegetative bud phenology in particular should be a major determinant of fitness in these perennial populations, as timing of the onset and release of dormancy is shaped largely by photoperiod and temperature regimes<sup>23,24</sup>. Indeed, genes related to photoperiod, drought and stress response were overrepresented among the selection outliers (**Supplementary Table 11**). To more directly test this hypothesis, we performed a genome-wide association study (GWAS) with spring bud flush, fall bud set and tree height measured at the three test sites, accounting for population stratification and background genetic effects in a mixed model framework for both univariate<sup>36</sup> and multivariate traits<sup>37</sup> (**Fig. 1b**, **Supplementary Tables 1 and 13** and **Supplementary Figs. 6–10**). More specifically, we found that those regions in the top 1% of selection scans had stronger adaptive trait association signals at all three test sites than would be expected by chance (i.e., the observed mean association signal was stronger than randomly resampled windows,

**Figure 2** Phenotypic evidence of climate-driven selection in *P. trichocarpa*. **(a)** Patterns of quantitative trait differentiation ( $Q_{ST}$ ) are stronger than genome-wide differentiation ( $F_{ST}$ ) among sampled geographic locations. Shaded area represents the 95% confidence interval (CI) of  $F_{ST}$ ; points and bars represent the point and 95% CI of  $Q_{ST}$ . **(b–d)** Genotypic estimates of best linear unbiased predictors (BLUPs) for adaptive traits growing in multiple plantation environments show strong correlations with the first principal component of 20 climate variables measured at the collection location. Negative PC1 values are associated with warmer conditions, and more positive bud flush and bud set BLUPs indicate earlier flush or set, respectively. Pearson correlation  $r$  and  $P$  values are shown.



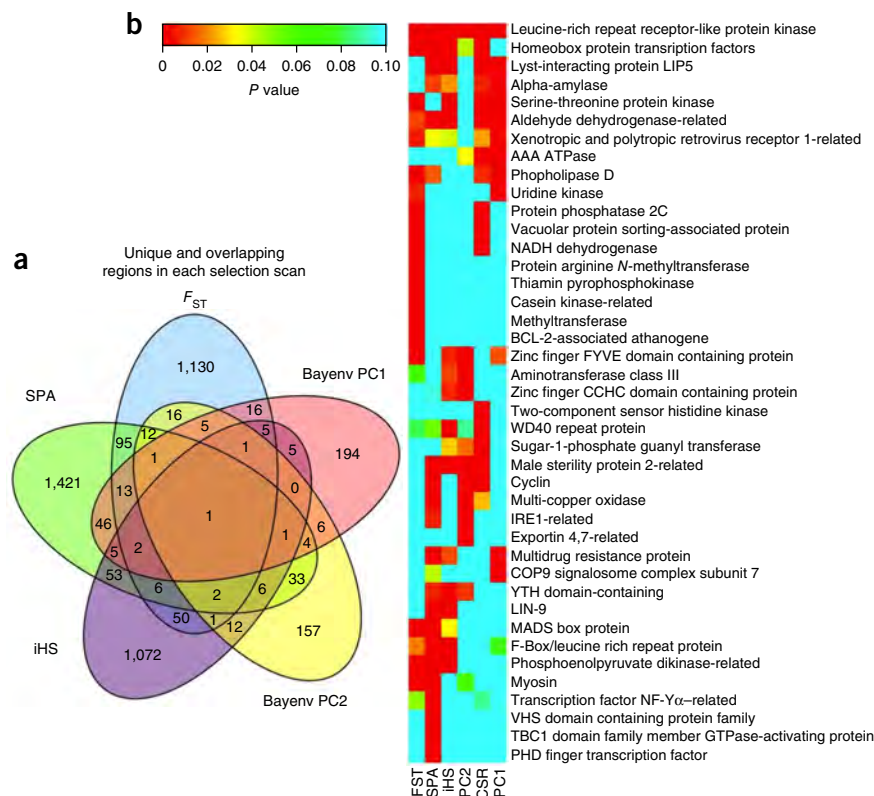
**Figure 3** Unique and shared genomic regions among five selection scans. (a) Venn diagram of the number of regions throughout the genome in the top 1% for each selection scan. (b) Overrepresentation  $P$  value (Fisher's exact test) for PANTHER annotation categories in selection outliers. Only the ten most strongly overrepresented categories for each selection scan are shown.

controlling for gene density;  $P < 0.00005$ ) (Fig. 4 and Supplementary Fig. 11). This was the case for all scans, including those based on spatial variation in allele frequency (for example,  $F_{ST}$  and bayenv) as well as those based on long haplotypes (iHS). This correspondence is therefore unlikely to be artifactual, supporting the hypothesis that these outlier regions are driven partly by selection on adaptive traits.

We found strong associations for both univariate analyses as well as the multi-trait GWAS for each trait among test sites (Supplementary Table 13). Though some of the strongest univariate associations were also identified in the multiple-plantation GWAS, many associations were nonoverlapping, perhaps owing to the strong environmental differences among the locations, which ranged from cool and wet (Clatskanie, Oregon) to hot and dry (Placerville, California). Strikingly few individual height-associated SNPs overlapped in comparisons between the Placerville plantation and the two other sites.

### Dormancy candidate genes in selection and GWAS regions

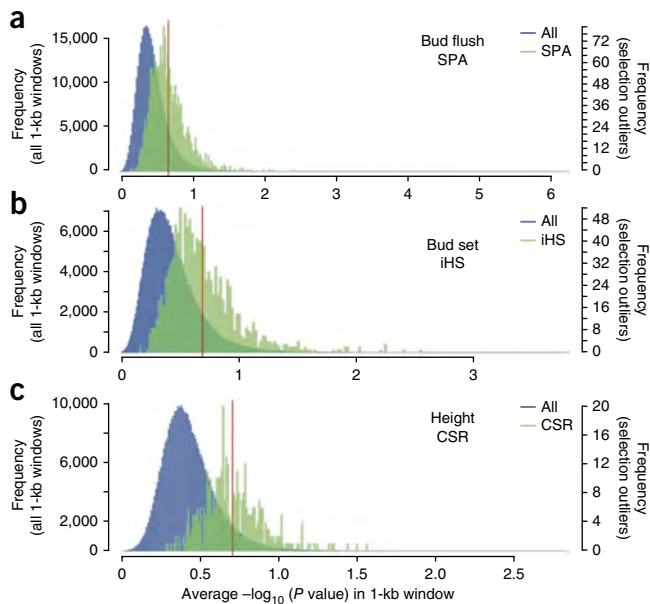
A number of dormancy-related genes were near the strongest GWAS and selection signals. A region on chromosome 10, characterized by high LD, was one of the CSRs and was associated with bud flush (mixed-model SNP association  $P = 5.19 \times 10^{-6}$ ) (Fig. 5). The strongest selection signal occurred near *Potri.010G079600*, encoding a DNA-damage repair protein, and a number of genes encoding lipid biosynthesis



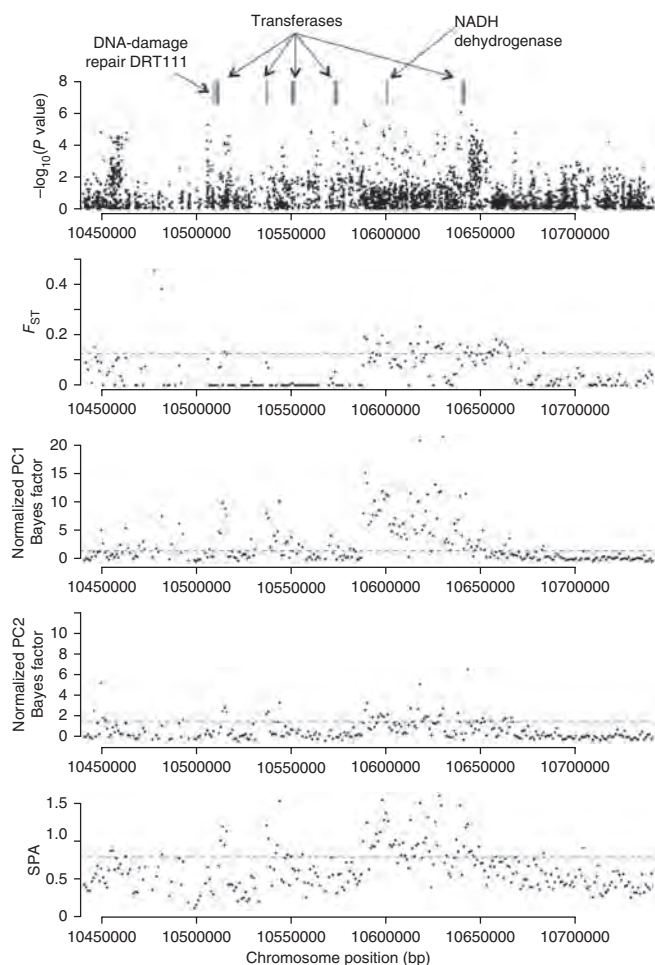
transferases. A strong bud set association also occurred near this region (Clatskanie and Corvallis, Oregon; Supplementary Fig. 12). The strongest association signal (mixed-model SNP association  $P = 5.69 \times 10^{-7}$ ), within 15 kb of a CSR, was just downstream of the coding region of *Potri.010G076100*, encoding a ureidoglycolate amidohydrolase (UAH), whose leaf and root expression is downregulated with short days<sup>38</sup>. Ureides are transportable intermediates of purine catabolism, and by catalyzing the final step in ureide catabolism, UAH has a role in the remobilization of nitrogen<sup>39</sup>. The ureide allantoin also influences abscisic acid metabolism and promotes abiotic stress tolerance in *Arabidopsis*<sup>39</sup>. However, to our knowledge, ureides and UAH have not previously been implicated as having important roles in seasonal nitrogen cycling or cold tolerance in *Populus*.

Among the photoperiodic and dormancy genes we found an  $F_{ST}$  outlier, *Potri.010G179700* (*FT2*), which influences growth cessation in *Populus*<sup>40</sup>. This gene had an intronic SNP strongly associated with bud set and height (mixed-model SNP association  $P < 0.00015$ , Supplementary Table 13) and was near strong SPA and bayenv outliers. A second gene, *Potri.008G117700* (similar to *PFT1*), occurred as an  $F_{ST}$  outlier region and was within 5 kb of several multitrait association signals ( $P = 7.17 \times 10^{-5}$ ). *Arabidopsis PFT1* is hypothesized to influence both defense and phytochrome B-mediated FT regulation<sup>41</sup>.

Among the strongest bud flush associations (mixed-model SNP association  $P = 2.72 \times 10^{-14}$ ) was a nonsynonymous mutation in *Potri.008G077400*, a 4-nitrophenylphosphatase-associated locus (Clatskanie and Corvallis, Fig. 6). This mutation is in high LD with many other significantly associated SNPs in the surrounding 40 kb,



**Figure 4** The selection outliers have a stronger association signal with adaptive traits than that expected by chance. (a–c) The genome-wide distribution of association signal in 1-kb windows through the genome (blue) and the association within the selection outliers (green; red line indicates mean) for three traits in different gardens.



**Figure 5** A region of chromosome 10 shows an abundance of bud flush association and strong evidence of selection from multiple selection scans. Top, mixed model SNP association  $P$  value for bud flush in Clatskanie, Oregon. Genes of interest identified by bars above. Dashed lines represent the 1% cutoff mark for selection scans.

including *Potri.008G076800*, (encoding FAR1 transcription factor) and *Potri.008G077300* (encoding UDP-galactose transporter), and is in an  $F_{ST}$  and bayenv PC1 outlier region. In this same region there is a bud flush association signal in all three test sites ( $P$  values ranging from  $2.01 \times 10^{-7}$  to  $1.08 \times 10^{-5}$ ) within *Potri.008G077700* (*FT1*), a gene previously implicated in *Populus* dormancy cycling<sup>42</sup>. However, it appears to be a weakly linked ( $r^2 = 0.14$ ), separate association signal from that in *Potri.008G077400*.

In summary, we detected genomic regions with patterns of diversity consistent with divergent and/or recent positive selection on a range of traits, and particularly on climate-related phenological and growth patterns. Although our selection scans and GWAS analyses identified genes previously known to influence adaptive traits, they also identified many loci of unknown function, which would not have been considered in any a priori candidate gene approach. Furthermore, the results and discussion presented above focus primarily on vegetative phenology, but many other traits are likely to be involved in determining fitness in these highly variable environments. In fact, the CSRs contained genes that have been implicated in controlling numerous other adaptive characteristics, including temperature stress tolerance, ion uptake and homeostasis, insect and pathogen defense and reproduction (**Supplementary Note**).

## Duplication and network connectedness

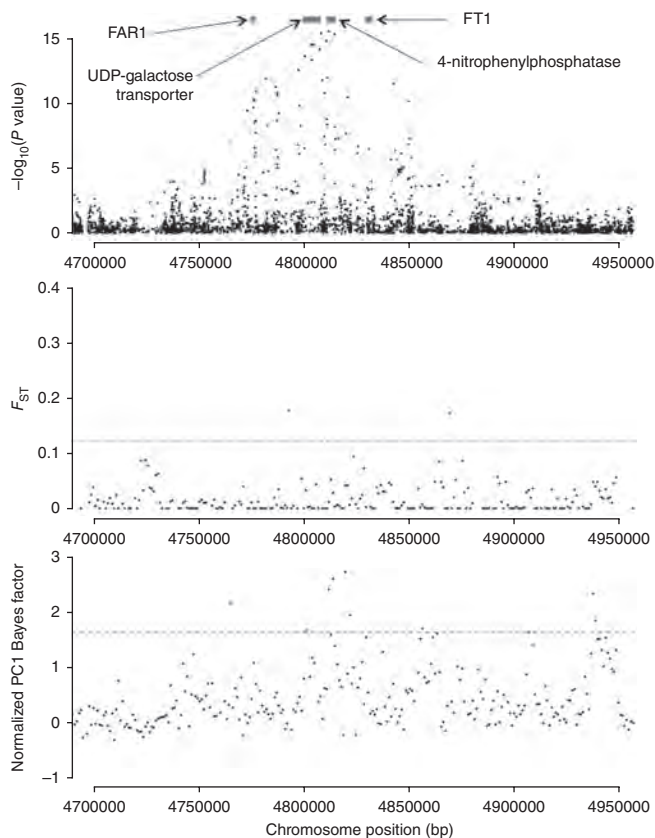
We tested whether genes associated with selection outliers were over- or underrepresented among the 7,906 identified gene pairs resulting from the Salicoid whole-genome duplication<sup>29,31</sup> (hereafter referred to as Salicoid duplicates), as compared to genes that occur as singletons (**Table 2**). These analyses suggest that recent positive selective sweeps (indicated by iHS) are less likely for retained Salicoid duplicates than for singleton genes, but when one occurs, the sweep tends to occur for both duplicates. We also found that genes nearest to the individual  $F_{ST}$ , SPA and iHS outliers had more predicted protein-protein interactions than genes in the rest of the genome (**Supplementary Fig. 13**;  $P \leq 0.05$ ). Furthermore, protein-protein interactions were negatively correlated with total and within-population nucleotide diversity and the ratio of nonsynonymous to synonymous diversity ( $r < -0.06$ ,  $P < 0.0001$ ). These results suggest that patterns of selection (both purifying and positive) are influenced by genomic context, including past whole-genome duplication events and gene or protein-protein interactions (**Supplementary Note**).

## DISCUSSION

A primary goal of evolutionary biology is to determine the influences of positive and purifying selection, as well as neutral forces in shaping genetic variation. Natural populations spanning wide climatic gradients offer an ideal opportunity to investigate these patterns. We sequenced more than 500 *P. trichocarpa* individuals from across much of the species range and identified more than 17 million SNPs (**Table 1** and **Fig. 1**). These polymorphisms revealed significant spatial and geographic structure, even at fine spatial scales. As previously suggested on the basis of small-scale sequencing and genotyping<sup>22</sup>, such patterns seem to have resulted from a complex demographic history.

Geographically structured, adaptive phenotypic variation is common among forest trees<sup>14,24,43</sup>. Climate is a fundamental driver of such variation<sup>14,24,25</sup>, and we identified quantitative trait differentiation and climate-related variation within our sample consistent with this pattern. However, the molecular and evolutionary processes underlying such adaptation often remain unknown. Although genome-wide polymorphism patterns suggest strong purifying selection throughout genic space, we also identified regions of the genome with unusually long haplotypes, among-population differentiation, and climatic gradients consistent with recent positive or divergent selection. Genes within these regions contain a variety of annotations plausibly related to local biotic and abiotic conditions, including photoperiod-responsive and dormancy-related loci, insect and pathogen defense, abiotic stress tolerance and phenylpropanoid metabolism. Such genes provide excellent targets for natural selection and for functional studies aimed at elucidating the drivers of local adaptation in black cottonwood and other species.

These largely nonoverlapping regions also provide insight into the variety of selection pressures and modes of selection acting within and among populations. For instance, classic, recent selective sweeps (iHS) are overrepresented among genes with annotations associated with heavy-metal homeostasis and symbiosis. But if climate-driven selection primarily acts on standing variation rather than new mutations, subtle allele frequency shifts among populations for many loci of small effects may be expected rather than hard selective sweeps. This is consistent with relatively little overlap among outlier regions identified with bayenv PC2 and iHS. Adaptation, therefore, probably occurs through different processes for different mutations—perhaps depending on mutation age, trait heritability and penetrance and number of loci involved, as has been suggested to occur in human populations<sup>44</sup>.



**Figure 6** A region of chromosome 8 shows multiple strong bud flush associations and evidence of positive selection. Top, the multi-environment mixed model SNP association  $P$  value for bud flush in all three plantations. Genes of interest identified by bars above. Dashed lines represent the 1% cutoff mark for selection scans.

Remarkably, the selection outlier loci were also enriched for polymorphisms associated with adaptive traits such as bud flush, bud set and height. Although factors such as stratification and linkage may produce erroneous associations<sup>45</sup>, mapping traits to computationally identified selection regions lends greater support to their functional significance. Similar patterns have been observed in the model annual plant *Arabidopsis*, where genomic regions showing signatures of selection are structured by climate variation<sup>9,12</sup> and collocated with adaptive trait associations<sup>9</sup>. Similar examples have been identified in domesticated crops<sup>5,11</sup>. However, to our knowledge, this is the first report of such concordance in a widespread, ecologically important undomesticated plant species.

We recognize that complex peaks of association may also be partially responsible for the overlap between selection scans and GWAS and differences in GWAS signal among gardens. LD combined with spurious patterns of random mutation or neutral stratification may produce synthetic associations<sup>45</sup> and/or composite phenotypes driven by multiple causal loci<sup>46</sup>. However, there is no reason to expect this correlative effect at high frequency on a genome-wide scale. Therefore, our findings suggest that the outliers contain variation relevant to adaptation on the basis of their statistically stronger-than-expected adaptive trait association signal.

The power of combining selection scans and association analyses is well illustrated by insights gained from our study into winter dormancy control in natural settings. Building on previous functional studies under highly controlled environments<sup>40–42,47</sup>, our results support a model of vegetative bud set and spring bud flush timing that

centers on regulation of expression and symplastic mobility of the FT1 and FT2 proteins. *FT1* is known to be transiently induced by chilling during winter and promotes the floral transition<sup>40</sup>. However, associations of *FT1* with vegetative bud flush suggest an additional function. Prolonged chilling releases endodormancy, the timing of which is correlated with bud flush through subsequent accumulation of warm-temperature units<sup>24</sup>. Moreover, the timing of the reopening of callose-plugged symplastic paths, endodormancy release, and *FT1* upregulation are correlated<sup>42</sup>. On the basis of our association results, we hypothesize that *FT1* is also involved in regulating endodormancy release and, hence, subsequent bud flush timing.

Reported studies of *Populus CEN1*, a flowering repressor and homolog of *TFL1*, encoding an FT antagonist, also provide support for this model<sup>48</sup>. Its winter expression is low when *FT1* expression is high, but *CEN1* is highly and transiently upregulated shortly before bud flush. However, constitutive overexpression of *CEN1* delays endodormancy release and bud flush<sup>48</sup>. In *Arabidopsis*, the balance between FT and *TFL1* seems to be central to the transition to flowering versus maintenance of indeterminate meristems<sup>49</sup>. Thus, *CEN1* might counterbalance *FT1* promotion of endodormancy release. In this model, the relative timing of *FT1* regulation could influence phenotypic variation observed in bud flush timing.

Patterns of adaptive variation are not independent of genomic history, and large-scale events such as whole-genome duplications can alter the evolutionary trajectories of certain loci. The deficiency of Salicoid duplicates among iHS outliers indicates that recent hard selective sweeps are less likely for genes retained from genome duplication, possibly because of fitness costs associated with altered function and/or stoichiometry of paralogs with large numbers of protein-protein interactions<sup>50,51</sup>. Furthermore, selective sweeps tend to affect both paralogs of a duplicated pair when they do occur, providing further support for the role of dosage constraints in duplicate gene evolution.

This is not to suggest that dosage constraints are the sole or even the primary drivers of the retention and evolution of duplicate genes. Abundant evidence supports subfunctionalization and neofunctionalization of Salicoid duplicates<sup>31</sup>. The case of the FT paralogs is again illustrative. *FT1* and *FT2* are Salicoid duplicates with divergent functions affecting distinct aspects of phenology and displaying diametrically opposed expression patterns in *Populus*<sup>40</sup>. Whereas *FT1* is expressed primarily during winter in dormant buds, *FT2* is expressed mainly during the growing season, maintaining vegetative growth<sup>40</sup>. Short days during fall lead to *FT2* suppression, in part through phytochrome influence on the transcription factor PFT1 (refs. 40,41). In support of this model, we found bud set associations with *FT2* and a *PFT1* paralog, and bud flush associations for *FT1*. This remarkable

**Table 2** Tests of over- and underrepresentation of retained Salicoid duplicate genes and pairs among the selection outliers

Selection scan	Duplicate genes in outlier regions	Duplicate pairs in outlier regions		
		$P$ value	$P$ value	
CSR	178	NS (0.623)	2	NS (0.263)
$F_{ST}$	674	Over ( $2.8 \times 10^{-9}$ )	27	Over (0.002)
SPA	741	Over (0.004)	24	NS (0.065)
iHS	348	Under ( $3.0 \times 10^{-12}$ )	8	Over (0.039)
bayenv PC1	100	NS (0.661)	1	NS (0.263)
bayenv PC2	134	NS (0.156)	0	NS (1)

Shown are the number of genes in each category and the associated  $P$  values (Fisher's exact test). 39,514 genes are found on the 19 chromosomes, with 7,609 pairs from 15,797 genes. NS, not significant; over, overrepresented compared to genome-wide expectation; under, underrepresented regions compared to genome-wide expectation.

divergence in function demonstrates the adaptive potential of Salicoid duplicate pairs, consistent with classic models of duplicate gene evolution<sup>52,53</sup>.

Intriguingly, a Salicoid duplicate pair that occurred in the CSRs are 1,3- $\beta$ -glucan (callose) synthase–encoding homologs (*Potri.002G058700* and *Potri.005G203500*). *Arabidopsis* callose synthases, when expressed in the phloem, deposit callose in the plasmodesmata, altering sugar and signaling molecule transport<sup>54,55</sup>. Returning to the phenological model outlined above, it has been hypothesized<sup>42</sup> that the formation and degradation of callose plugs are a control point for dormancy onset and release, possibly blocking translocation of *FT1* and *FT2*. These duplicates may also have divergent functions and expression patterns, similar to those observed for the *FT* paralogs.

Our findings have important implications for understanding mechanisms of adaptation of ecologically dominant plants with widespread distributions. Forestry trials have for more than 200 years indicated substantial local adaptation of dominant trees<sup>56</sup>, but ours is the first that we know of to explore the genomic legacy of this selection across the entire genome and highlight both the wide range of selection pressures as well as the climatic influence on phenological systems. These findings also have important implications for the management of natural populations in the face of environmental change. Seed-transfer zone guidelines have historically required large numbers of plantations to accurately estimate transfer parameters<sup>57</sup>. Computationally identifying adaptive variants through selection scans and genome-wide phenotypic prediction could provide information in the absence of extensive plantation trials, maximizing genetic diversity while matching germplasm to current and future environmental pressures. Management and modification of such genetic diversity will undoubtedly affect dependent biotic communities and ecosystem functioning, which are influenced by genetic variation in trees<sup>17</sup>.

The 17.9 million SNPs we identified represent naturally segregating variants in wild populations, which can be used for multiple objectives. Forest tree improvement has traditionally relied on natural variation in breeding programs through targeted crossing based on superior phenotypes<sup>20</sup>. The availability of whole-genome sequences can enable alternative breeding approaches, including genome-wide phenotypic prediction<sup>58</sup> and breeding with rare defective alleles, which relies on rare, recessive mutations of large effect that are commonly heterozygous and therefore masked from many approaches<sup>59</sup>. Most SNPs found here are intergenic and uncommon, but many have predicted major effects in genic regions. Several SNPs of the latter type are in the candidate selection regions, including altered start and stop codons and alternative splice variants, which could represent an immediate set of tractable targets for breeding programs constrained by long generation times. Several occur at high frequency in the isolated southern or northern populations, demonstrating that sampling populations throughout the range, including marginal populations, will yield many more variants of potential utility.

**URLs.** Phytozome, <http://www.phytozome.net/poplar.php>; Plant Transcription Factor Database (version 3.0), <http://plantfdb.cbi.pku.edu.cn/index.php?sp=Pth>; Picard package, <http://picard.sourceforge.net>.

## METHODS

Methods and any associated references are available in the [online version of the paper](#).

**Accession codes.** SNP data are also available through Phytozome (see URLs).

*Note: Any Supplementary Information and Source Data files are available in the online version of the paper.*

## ACKNOWLEDGMENTS

We thank the members of BioEnergy Science Center for their varied contributions to this work, and especially those involved in the collection, propagation and maintenance of the common gardens, including G. Howe, A. Groover, R. Stettler, J. Johnson and the staff at Mt. Jefferson Farms and Greenwood Resources. We thank the West Virginia University High Performance Computing facility, in particular N. Gregg and M. Carlise. *P. balsamifera* transcriptomes were provided by M. Olson (Texas Tech University). This work was supported by funding from the BioEnergy Science Center, a US Department of Energy (DOE) Bioenergy Research Center supported by the Office of Biological and Environmental Research in the DOE Office of Science. A.M.B. acknowledges support from the Virginia Agricultural Experiment Station and the McIntire Stennis Program of the National Institute of Food and Agriculture, US Department of Agriculture.

## AUTHOR CONTRIBUTIONS

G.A.T., S.P.D., G.T.S. and L.M.E. conceived and designed the study. All authors performed measurements. L.G., J.M. and W.S. performed sequencing. L.M.E., S.P.D., G.T.S., E. R.-M., J.M., P.R., W.M. and W.S. performed analyses. L.M.E., S.P.D. and A.M.B. drafted the manuscript. All authors read, revised, and approved the manuscript.

## COMPETING FINANCIAL INTERESTS

The authors declare no competing financial interests.

Reprints and permissions information is available online at <http://www.nature.com/reprints/index.html>.

1. The 1000 Genomes Project Consortium. An integrated map of genetic variation from 1,092 human genomes. *Nature* **491**, 56–65 (2012).
2. Cao, J. *et al.* Whole-genome sequencing of multiple *Arabidopsis thaliana* populations. *Nat. Genet.* **43**, 956–963 (2011).
3. Axelsson, E. *et al.* The genomic signature of dog domestication reveals adaptation to a starch-rich diet. *Nature* **495**, 360–364 (2013).
4. Hufford, M.B. *et al.* Comparative population genomics of maize domestication and improvement. *Nat. Genet.* **44**, 808–811 (2012).
5. Huang, X. *et al.* Genome-wide association study of flowering time and grain yield traits in a worldwide collection of rice germplasm. *Nat. Genet.* **44**, 32–39 (2012).
6. Zhao, S. *et al.* Whole-genome sequencing of giant pandas provides insights into demographic history and local adaptation. *Nat. Genet.* **45**, 67–71 (2013).
7. Miller, W. *et al.* Polar and brown bear genomes reveal ancient admixture and demographic footprints of past climate change. *Proc. Natl. Acad. Sci. USA* **109**, E2382–E2390 (2012).
8. Voight, B.F., Kudaravalli, S., Wen, X. & Pritchard, J.K. A map of recent positive selection in the human genome. *PLoS Biol.* **4**, e72 (2006).
9. Fournier-Level, A. *et al.* A map of local adaptation in *Arabidopsis thaliana*. *Science* **334**, 86–89 (2011).
10. Tishkoff, S.A. *et al.* Convergent adaptation of human lactase persistence in Africa and Europe. *Nat. Genet.* **39**, 31–40 (2007).
11. Jia, G. *et al.* A haplotype map of genomic variations and genome-wide association studies of agronomic traits in foxtail millet (*Setaria italica*). *Nat. Genet.* **45**, 957–961 (2013).
12. Hancock, A.M. *et al.* Adaptation to climate across the *Arabidopsis thaliana* genome. *Science* **334**, 83–86 (2011).
13. Grossman, S.R. *et al.* Identifying recent adaptations in large-scale genomic data. *Cell* **152**, 703–713 (2013).
14. Savolainen, O., Pyhäjärvi, T. & Knürr, T. Gene flow and local adaptation in trees. *Annu. Rev. Ecol. Evol. Syst.* **38**, 595–619 (2007).
15. Bonan, G.B. Forests and climate change: forcings, feedbacks, and the climate benefits of forests. *Science* **320**, 1444–1449 (2008).
16. Ellison, A.M. *et al.* Loss of foundation species: consequences for the structure and dynamics of forested ecosystems. *Front. Ecol. Environ.* **3**, 479–486 (2005).
17. Whitham, T.G. *et al.* Extending genomics to natural communities and ecosystems. *Science* **320**, 492–495 (2008).
18. Parmesan, C. Ecological and evolutionary responses to recent climate change. *Annu. Rev. Ecol. Evol. Syst.* **37**, 637–669 (2006).
19. Ingvarsson, P.K., García, M.V., Hall, D., Luquez, V. & Jansson, S. Clinal variation in *phyB2*, a candidate gene for day-length-induced growth cessation and bud set, across a latitudinal gradient in European aspen (*Populus tremula*). *Genetics* **172**, 1845–1853 (2006).
20. Neale, D.B. & Kremer, A. Forest tree genomics: growing resources and applications. *Nat. Rev. Genet.* **12**, 111–122 (2011).
21. Jansson, S. & Douglas, C.J. *Populus*: a model system for plant biology. *Annu. Rev. Plant Biol.* **58**, 435–458 (2007).
22. Slavov, G.T. *et al.* Genome resequencing reveals multiscale geographic structure and extensive linkage disequilibrium in the forest tree *Populus trichocarpa*. *New Phytol.* **196**, 713–725 (2012).

23. Pauley, S.S. & Perry, T.O. Ecotypic variation in the photoperiodic response in *Populus*. *J. Arnold Arbor.* **35**, 167–188 (1954).
24. Howe, G.T. *et al.* From genotype to phenotype: unraveling the complexities of cold adaptation in forest trees. *Can. J. Bot.* **81**, 1247–1266 (2003).
25. McKown, A.D. *et al.* Geographical and environmental gradients shape phenotypic trait variation and genetic structure in *Populus trichocarpa*. *New Phytol.* **201**, 1263–1276 (2014).
26. Wegrzyn, J.L. *et al.* Association genetics of traits controlling lignin and cellulose biosynthesis in black cottonwood (*Populus trichocarpa*, Salicaceae) secondary xylem. *New Phytol.* **188**, 515–532 (2010).
27. Porth, I. *et al.* Genome-wide association mapping for wood characteristics in *Populus* identifies an array of candidate single nucleotide polymorphisms. *New Phytol.* **200**, 710–726 (2013).
28. Tang, H. *et al.* Unraveling ancient hexaploidy through multiply-aligned angiosperm gene maps. *Genome Res.* **18**, 1944–1954 (2008).
29. Tuskan, G.A. *et al.* The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* **313**, 1596–1604 (2006).
30. Rodgers-Melnick, E., Culp, M. & DiFazio, S.P. Predicting whole genome protein interaction networks from primary sequence data in model and non-model organisms using ENTS. *BMC Genomics* **14**, 608 (2013).
31. Rodgers-Melnick, E. *et al.* Contrasting patterns of evolution following whole genome versus tandem duplication events in *Populus*. *Genome Res.* **22**, 95–105 (2012).
32. Spitze, K. Population structure in *Daphnia obtusa*: quantitative genetic and allozymic variation. *Genetics* **135**, 367–374 (1993).
33. Yang, W.-Y., Novembre, J., Eskin, E. & Halperin, E. A model-based approach for analysis of spatial structure in genetic data. *Nat. Genet.* **44**, 725–731 (2012).
34. Günther, T. & Coop, G. Robust identification of local adaptation from allele frequencies. *Genetics* **195**, 205–220 (2013).
35. Sun, J., Xie, D., Zhao, H. & Zou, D. Genome-wide identification of the class III aminotransferase gene family in rice and expression analysis under abiotic stress. *Genes Genomics* **35**, 597–608 (2013).
36. Kang, H.M. *et al.* Variance component model to account for sample structure in genome-wide association studies. *Nat. Genet.* **42**, 348–354 (2010).
37. Zhou, X. & Stephens, M. Efficient multivariate linear mixed model algorithms for genome-wide association studies. *Nat. Methods* **11**, 407–409 (2014).
38. Ruttink, T. *et al.* A molecular timetable for apical bud formation and dormancy induction in poplar. *Plant Cell* **19**, 2370–2390 (2007).
39. Werner, A.K. *et al.* The ureide-degrading reactions of purine ring catabolism employ three amidohydrolases and one aminohydrolase in *Arabidopsis*, soybean, and rice. *Plant Physiol.* **163**, 672–681 (2013).
40. Hsu, C.-Y. *et al.* FLOWERING LOCUS T duplication coordinates reproductive and vegetative growth in perennial poplar. *Proc. Natl. Acad. Sci. USA* **108**, 10756–10761 (2011).
41. Iñigo, S., Alvarez, M.J., Strasser, B., Califano, A. & Cerdán, P.D. PFT1, the MED25 subunit of the plant Mediator complex, promotes flowering through CONSTANS dependent and independent mechanisms in *Arabidopsis*. *Plant J.* **69**, 601–612 (2012).
42. Rinne, P.L.H. *et al.* Chilling of dormant buds hyperinduces FLOWERING LOCUS T and recruits GA-inducible 1,3- $\beta$ -glucanases to reopen signal conduits and release dormancy in *Populus*. *Plant Cell* **23**, 130–146 (2011).
43. Hall, D. *et al.* Adaptive population differentiation in phenology across a latitudinal gradient in European aspen (*Populus tremula*, L.): a comparison of neutral markers, candidate genes and phenotypic traits. *Evolution* **61**, 2849–2860 (2007).
44. Pritchard, J.K. & Di Rienzo, A. Adaptation—not by sweeps alone. *Nat. Rev. Genet.* **11**, 665–667 (2010).
45. Platt, A., Vilhjálmsson, B.J. & Nordborg, M. Conditions under which genome-wide association studies will be positively misleading. *Genetics* **186**, 1045–1052 (2010).
46. Atwell, S. *et al.* Genome-wide association study of 107 phenotypes in *Arabidopsis thaliana* inbred lines. *Nature* **465**, 627–631 (2010).
47. Böhlenius, H. *et al.* CO/FT regulatory module controls timing of flowering and seasonal growth cessation in trees. *Science* **312**, 1040–1043 (2006).
48. Mohamed, R. *et al.* *Populus CEN/TFL1* regulates first onset of flowering, axillary meristem identity and dormancy release in *Populus*. *Plant J.* **62**, 674–688 (2010).
49. Jaeger, K.E., Pullen, N., Lamzin, S., Morris, R.J. & Wigge, P.A. Interlocking feedback loops govern the dynamic behavior of the floral transition in *Arabidopsis*. *Plant Cell* **25**, 820–833 (2013).
50. Freeling, M. Bias in plant gene content following different sorts of duplication: tandem, whole-genome, segmental, or by transposition. *Annu. Rev. Plant Biol.* **60**, 433–453 (2009).
51. Birchler, J.A. & Veitia, R.A. The gene balance hypothesis: implications for gene regulation, quantitative traits and evolution. *New Phytol.* **186**, 54–62 (2010).
52. Lynch, M. & Force, A. The probability of duplicate gene preservation by subfunctionalization. *Genetics* **154**, 459–473 (2000).
53. Taylor, J.S. & Raes, J. Duplication and divergence: the evolution of new genes and old ideas. *Annu. Rev. Genet.* **38**, 615–643 (2004).
54. Vatén, A. *et al.* Callose biosynthesis regulates symplastic trafficking during root development. *Dev. Cell* **21**, 1144–1155 (2011).
55. Xie, B., Wang, X., Zhu, M., Zhang, Z. & Hong, Z. *CalS7* encodes a callose synthase responsible for callose deposition in the phloem. *Plant J.* **65**, 1–14 (2011).
56. Langlet, O. Two hundred years of genealogy. *Taxon* **20**, 653–722 (1971).
57. Wang, T., O'Neill, G.A. & Aitken, S. N. Integrating environmental and genetic effects to predict responses of tree populations to climate. *Ecol. Appl.* **20**, 153–163 (2010).
58. Grattapaglia, D. & Resende, M.D.V. Genomic selection in forest tree breeding. *Tree Genet. Genomes* **7**, 241–255 (2011).
59. Vanholme, B. *et al.* Breeding with rare defective alleles (BRDA): a natural *Populus nigra* HCT mutant with modified lignin as a case study. *New Phytol.* **198**, 765–776 (2013).
60. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80–92 (2012).



## ONLINE METHODS

**Sequencing, assembly and variant calling.** We obtained plant materials from 1,100 black cottonwood (*Populus trichocarpa* Torr & Gray) from wild populations in California, Oregon, Washington and British Columbia, as previously described<sup>22</sup>. We resequenced a set of 649 genotypes to a minimum expected depth of 15× using the Illumina Genome Analyzer, HiSeq 2000 and HiSeq 2500. Sequences were down-sampled for those individuals sequenced at greater depths to ensure even coverage throughout the population (**Supplementary Fig. 1a**). Short reads were then aligned to the *P. trichocarpa* version 3 (v3.0) genome using BWA 0.5.9-r16 with default parameters<sup>61</sup>. We corrected mate pair metadata and marked duplicate molecules using the FixMateInformation and MarkDuplicates methods in the Picard package (<http://picard.sourceforge.net>). Next, we called SNPs and small indels for the merged data set using SAMtools mpileup (-E -C 50 -DS -m 2 -F 0.000911 -d 50000) and bcftools (-bcgv -p 0.999089)<sup>62</sup>.

**Genotype validation.** We compared the samtools mpileup genotype calls for 649 individuals to 22,438 SNPs assayed on the *Populus* Illumina Infinium platform, which was designed on the basis of assembly version 2.0 (refs. 22,63). These were high-quality SNPs that we could confidently place on the v3 reference genome. The 649 individuals had, on average, a 97.9% match rate. SNPs with a minor allele frequency (MAF)  $\geq 0.05$  had a match rate of 98.1%, and those with MAF  $\leq 0.01$  ( $n = 159$  SNPs) had a match rate of 78.2%, similar to other published studies<sup>4,64,65</sup>. Stringent filtering had minimal impact on match rate, though it reduced substantially the number of known SNPs passing the filtering thresholds. For example, requiring an individual minimum depth of 3, minimum mapping quality of 30, minor allele count of 15 and minimum quality score of 30 increased the false negative rate by 3.9%, but only increased the match rate by 0.3%. Therefore, no additional filtering after samtools mpileup variant calling was performed.

Nisqually-1 was the original individual sequenced<sup>29</sup> using Sanger technology, and it was also resequenced during this study using the Illumina platform. 716,691 heterozygous polymorphisms found in the v3.0 reference genome assembly (<http://www.phytozome.net/poplar.php>) had at least three Sanger reads of each allele, and therefore had strong evidence of heterozygosity in the Sanger assembly. In the current study, we correctly identified 557,738 of these (77.82%), including 3,205 of 3,220 singleton variants in Nisqually-1 in the Illumina data, suggesting a 22.18% false negative rate. Conversely, of 1,115,963 heterozygous positions identified in Nisqually-1 in the current Illumina genotyping, 972,254 had at least one Sanger read supporting each allele, suggesting a 12.86% false positive rate. All of these comparisons were done with no filtering of the samtools mpileup genotype calls. It is important to note that errors occur in both the Sanger and Illumina methods, so these are likely to be overestimates of the true error rates in the resequencing SNP data.

**The accessible genome.** Next, we identified the *Populus trichocarpa* ‘accessible genome’ as those positions that had sufficient read depth across enough individuals to enable genotypes to be accurately determined (similar to the approach used in the 1000 Genomes Project)<sup>1</sup>. We estimated the median and interquartile range of depth for each position in the genome, for all sequenced individuals, using samtools mpileup. With our target of 15× coverage, ‘accessible’ positions were those with median depth between 5 and 45 (inclusive) and with an interquartile range less than or equal to 15 (**Supplementary Fig. 1a,b**). Of the 394,507,732 positions that were sequenced across all individuals, 345,217,484 met these criteria (~87.51%), 17,902,170 of which were single nucleotide polymorphisms (SNPs) (15,454,190 biallelic). We observed a slight deficiency of heterozygotes at lower depth positions; however, these positions cumulatively comprise only between 0.7% and 2.5% of positions at an uncorrected Hardy-Weinberg equilibrium  $P$  value threshold of 0.001 (**Supplementary Fig. 1c**). Furthermore, these cutoffs did not bias the outcomes of selection scans throughout the genome, as putative selection outliers (see below) had a very similar distribution of depth as the rest of the genome (**Supplementary Table 14**) and there was no relationship with association  $P$  value (see below; all Pearson  $|r| < 0.005$ , **Supplementary Fig. 1d**).

**Relatedness, hybridization and population structure.** We next identified individuals that showed evidence of admixture with other species of *Populus*

because hybridization is common within the genus<sup>66</sup>. We used seven additional individuals sequenced to at least 32× depth as above: three *Populus deltoides*, one *Populus fremontii*, one *Populus angustifolia*, one *Populus nigra* and one *Populus tremuloides*. These were aligned to the *P. trichocarpa* v3.0 reference genome using Bowtie2 in local alignment mode and default parameters<sup>67</sup>, and variants were called using the samtools mpileup function for each species separately. We then used smartpca<sup>68</sup> to identify sampled individuals in this study that were genetically similar to these alternative species. This method identified three individuals that appear intermediate between the *P. trichocarpa* cluster and an alternate species (**Supplementary Fig. 14**).

We performed similar analyses using overlapping genomic regions from 32 *Populus balsamifera* transcriptomes (provided by M. Olson) (**Supplementary Fig. 15**), and, separately, the Illumina Infinium array data, which contained additional individuals of alternative species<sup>63</sup>. These identified an additional three genetically intermediate individuals. These six potentially admixed individuals were removed from subsequent analyses.

We next identified and removed individuals more related than first cousins using the program GCTA<sup>69</sup>. Because this, like most other relatedness estimates, relies on allele frequency estimates within populations, it was necessary to first identify genetic clusters. We iteratively identified genetic clusters using PCA<sup>68</sup>, each representing a putative genetic group. We removed related individuals within each from further analyses, leaving a total of 544 individuals, which were used for all subsequent analyses.

To assess population structure, we used PCA analyses with these unrelated 544 individuals. This identified roughly four major groupings (**Fig. 1a**). We then performed PCA analysis using only those individuals from the Washington and British Columbia group to investigate finer-scale structure (**Fig. 1b**). PCA was performed using all 17.9 million SNPs.

**Phenotypic evidence of selection.** We investigated phenotypic evidence of selection using two methods. First, we compared neutral genetic differentiation among collection rivers or subpopulations ( $F_{ST}$ ; see below for details of estimation) to differentiation among rivers for second-year height and fall and spring phenology using data collected from three replicated plantations ( $Q_{ST}$ ). Briefly, more than 1,000 *P. trichocarpa* genotypes were planted in 2009 in three replicated common gardens (Clatskanie and Corvallis, OR, and Placerville, CA) in a randomized block design with three replicates of each genotype. In 2010, we measured spring bud flush, fall bud set, and total height in each garden. We removed within-garden microsite variation using thin-plate spline regression (*fields* R package), then estimated among river, among genotypes within rivers, and residual variance components ( $\sigma^2_R$ ,  $\sigma^2_G$ , and  $\sigma^2_e$ , respectively) using a linear mixed-model (*lmer* function of the *lme4* R package).  $Q_{ST}$  was estimated at the river level as  $\sigma^2_R/(\sigma^2_R + 2 \cdot \sigma^2_G)$  (ref. 32). A 95% confidence interval of  $Q_{ST}$  was estimated by resampling rivers, with replacement, 1,000 times and estimating  $Q_{ST}$  for each bootstrapped data set. We directly compared the 95% CIs for  $Q_{ST}$  and  $F_{ST}$ . We note that in using clonal replicates  $\sigma^2_G$  includes additive and non-additive genetic effects, rather than the additive genetic variance alone; however, simulations have shown that this approach lowers  $Q_{ST}$  estimates, and is therefore a conservative test of  $Q_{ST} > F_{ST}$ <sup>70</sup>.

Second, we tested for correlations between these adaptive traits and the climate of the source location. We tested correlations with mean annual temperature, mean annual precipitation, and the first two principal components (cumulatively >85% of variance explained) of 20 climate variables obtained using ClimateWNA<sup>71</sup>. We used the genotypic best linear unbiased predictors obtained from mixed model analysis (*lmer* function of the *lme4* R package) as the phenotypic traits. Climate variables were averaged within collection locations before correlation analysis.

**Genetic variation and signatures of recent positive selection throughout the genome.** We assessed species-wide nucleotide diversity ( $\pi$ )<sup>72</sup> using the MLE estimate of allele frequency from the samtools mpileup output<sup>62</sup> in all annotated regions (coding sequence, introns and 5′ and 3′ UTRs) of the v3.0 genome longer than 150 bp and with at least 95% accessibility.

We performed five genome-wide scans of recent positive selection, using four conceptually different approaches. First, we estimated genetic differentiation<sup>72</sup> among collection rivers as  $F_{ST}$  in 1-kb windows throughout the genome (again, requiring at least 95% accessibility and using the accessible positions

in a window as the window's full length). We restricted this analysis to rivers or subpopulations with at least eight individuals and randomly chose 20 individuals from those that contained >20 individuals (14 rivers total: Homathko, Skwawka, Lillooet, Squamish, Salmon, Fraser, Columbia, Nisqually, Nooksack, Puyallup, Skagit, Skykomish, Tahoe, Willamette). We estimated nucleotide diversity across all individuals ( $\pi_T$ ) and weighted within-river nucleotide diversity ( $\pi_S$ ), accounting for sequencing error<sup>73</sup>. We calculated  $F_{ST}$  as the difference between total and weighted within-river diversity, divided by the total diversity ( $\pi_{T-S} / \pi_T$ ) (ref. 72). We took the top 1% of the empirical distribution of  $F_{ST}$  as genomic regions representing unusually strong allele frequency differences among rivers and candidates of divergent selection.

The second selection scan quantified the steepness of allele frequency clines across two climate variables, using the program SPA<sup>33</sup>. SPA uses a logistic regression-based approach to model allele frequency clines, without a priori population assignment and represents a fundamentally different approach from that of the  $F_{ST}$  scan described above. We used mean annual temperature and mean annual precipitation of the source location for each sample, obtained using ClimateWNA<sup>71</sup>, because these variables are significantly correlated (Pearson  $|r| = 0.1-0.57$ ,  $P < 0.0001$ ) with growth and phenological traits. We averaged SPA in nonoverlapping 1-kb bins throughout the genome, requiring at least five SNPs in each window. We identified the top 1% of these windows as regions of the genome with unusually steep allele frequency clines across mean annual temperature and precipitation.

Third, we identified regions of the genome with recent, unusually rapid increases in allele frequency across the range. Strong, recent selective sweeps will result in long haplotypes associated with the selected allele<sup>8,74</sup>. First, we phased the 544 diploid individuals using SHAPEIT2 (ref. 75). Because we have no reference haplotype panels to test the accuracy of computationally determined haplotypes, we determined the optimal method by estimating the accuracy of imputed masked loci<sup>76</sup>. We used 10 Mb of chromosome 2 (5–15 Mb), using only variants with MAF >0.1 (307,123 sites). We randomly masked out 5% of the center 260,000 positions for each individual (avoiding the ends), treating them as missing for phasing. To determine the optimal number of hidden Markov states (K) and the window size (W) used in SHAPEIT2, we phased the data using combinations of parameters from K = 50–600 and W = 0.1–2 Mb (Supplementary Fig. 14), using the default effective population size ( $N_e$ ) = 15K, and run with four threads. The genetic position was determined through linear interpolation using a genetic map derived from a *P. trichocarpa* × *P. deltoides* pseudo-backcross pedigree and 3,559 Infinium SNP markers<sup>22</sup>. Genetic position and recombination rate were estimated using local linear regression with the *loess* function in R. For comparison, we also phased the same data using the default settings of BEAGLE<sup>77</sup>. We then determined the squared correlation coefficient ( $R^2$ ) between the known allele dosages (0, 1, or 2) and the imputed genotypes for masked positions in each individual. The average  $R^2$  is shown in Supplementary Figure 16, and peaks at approximately K = 350, W = 0.1 Mb. We varied  $N_e$  from 10,000–20,000, and found that  $N_e = 15,000$  gave the highest correlation between known and imputed allele dosage for masked missing data. Using the same 10 Mb region of chromosome 2, we tested whether the 0.1 MAF cutoff affected accuracy, and found that with no MAF cutoff accuracy was actually increased. We therefore phased all chromosomes using SHAPEIT2 with K = 350 states, W = 0.1 Mb window size, and  $N_e = 15,000$  effective population size, using all non-singleton and -private doubleton sites, parallelized using 24 threads.

We then estimated the integrated haplotype score (iHS)<sup>8</sup> for SNPs. Because the program is computationally intensive, we thinned the data set to SNPs separated by at least 100 bp and with a MAF of at least 0.05, resulting in 1,898,506 SNPs throughout the genome. In calculating iHS, we used the genetic distance as described above. iHS was standardized within allele frequency bins<sup>8</sup>, and |iHS| averaged within nonoverlapping 1-kb windows, again requiring at least five SNPs in a window. We took the top 1% of these bins as genomic regions that have experienced an unusually rapid allele frequency change, resulting in extended haplotype homozygosity, and potential targets of positive selection.

Finally, we used bayenv2.0 (ref. 34) to identify regions of the genome with unusually strong allele frequency clines along climatic gradients while controlling for background neutral population structure. We performed this analysis with 13 of the populations used in the  $F_{ST}$  analysis described

above. We excluded the Tahoe population because it was so divergent that the neutral model of bayenv2.0 had difficulty accounting for the covariance in allele frequencies among populations (data not shown). We used the first two PCs of the climate data from source locations, averaged within populations, which cumulatively explained >85% of the variance in the correlation matrix. Loadings showed that the first PC was strongly related to all climate WNA variables, while the second PC was more strongly related to precipitation, heat-moisture indices, and frost-free period metrics (Supplementary Fig. 17). To estimate the covariance matrix of allele frequency among populations, we used 19,420 genome-wide SNPs that were separated by at least 20 kb and with MAF > 0.01 across the 13 populations using bayenv2.0 with 100,000 steps through the chain, performed three times independently. The three runs were very similar (all Mantel  $R > 0.985$ ,  $P < 0.001$ ), and the difference in covariances among runs were always less than 3% of the smallest estimated covariance, indicating convergence<sup>78</sup>. We assessed the strength of the correlation of allele frequency and the climate variables, as estimated by the Bayes factor (BF) and Spearman correlation, for 9,519,343 SNPs (MAF > 0.01 across the 13 populations). We tested, for 20,000 randomly chosen SNPs, the effect of chain length on the BFs. Correlations of the individual SNPs among the different chain lengths and independent runs for each chain length indicated that ten chains of 50,000 steps were sufficient to ensure repeatability and accuracy (Supplementary Fig. 18) and remain tractable for millions of SNPs. For the final analysis of >9.5 million SNPs, we calculated the BF and Spearman correlation using 50,000 steps in each of ten independent runs. We averaged the  $\log_{10}(\text{BF})$  and the posterior Spearman correlation estimate for each SNP, normalized these values within MAF bins (0.05 bin size), and averaged these within 1-kb windows throughout the genome, requiring at least five SNPs per 1-kb window.

To identify regions of the genome with unusually strong allele frequency-climate correlations, we selected the windows in the top 1% of Spearman climate-allele frequency correlations and top 1% of BFs as those with unusually strong climate-related allele frequency clines. This process was done separately for the first and second PCs, resulting in two separate selection scans.

**Candidate selection regions and annotation analysis.** The selection scans represent five different approaches to identifying unusually strong patterns throughout the genome that are consistent with recent positive or divergent selection. Merging nearby windows ( $\leq 5$  kb), we found 397 regions that were in the top 1% of at least two of the five scans (the candidate selection regions (CSRs)), spanning or adjacent to 452 different genes. We identified the genes spanning or nearest to the CSRs and selection outlier regions. We used Fisher exact tests (FET) to determine whether GO, PANTHER and Pfam annotations were overrepresented in the genes associated with the CSRs and outlier regions.

We also tested whether these genes were overrepresented among lists from known gene families and pathways, and known to be responsive to drought and dormancy cycling. Families of transcription factors were identified using the Plant Transcription Factor Database v3.0 (<http://planttfdb.cbi.pku.edu.cn/index.php?sp=Pth>)<sup>79</sup>. Genes in additional pathways and families are listed in Supplementary Table 11. When necessary, we used the best reciprocal BLAST hit between the v1 and v3 genome assemblies to locate the gene models identified by previous studies for each set of published genes.

**Genome duplication and network connectedness.** First, we examined the genes spanning or nearest to the CSRs and the windows of the top 1% of each selection scan in the context of the Salicoid whole-genome duplication using the 7,936 duplicate pairs identified previously<sup>31</sup>. We used FETs to test whether these selection scan lists were under- or overrepresented among the duplicate pairs. To determine whether there were more duplicate pairs in which both genes of the pair were associated with the selection outliers than expected by chance, we used a random resampling procedure. For each selection scan, we resampled without replacement the same number of genes observed in that scan that were also retained duplicates from the total number of retained duplicates (15,812) 10,000 times and recorded how many complete pairs were resampled each time—i.e., how many times both genes of a pair were randomly sampled. We tested whether genes associated with selection outliers had more protein-protein interactions (PPIs) than expected. We used the number of connections in PPI networks with 65% confidence determined

by the ENTS random forest prediction program<sup>30</sup>. We tested whether PPIs of the genes in each scan were different from the genome-wide average using Wilcoxon two-sample tests. These analyses examined patterns of genes associated with the CSRs and the selection outlier regions.

We also examined patterns at the whole-genome level, by calculating  $\pi_S$ ,  $\pi_T$ , and the ratio of nonsynonymous to synonymous polymorphisms ( $\pi_{\text{Nonsynonymous}}/\pi_{\text{Synonymous}}$ ) for 39,514 genes on the 19 chromosomes using the same methods described above. We then calculated the correlation of each statistic between the 7,936 Salicoid duplicate pairs of genes. To determine whether the observed correlation was greater than expected by chance, we randomly chose 7,936 pairs of genes from all genes 10,000 times, as a null distribution of correlation between pairs of randomly chosen genes.

We also tested whether the mean observed selection statistic differed between Salicoid duplicates and nonduplicate genes using Wilcoxon two-sample tests. To test whether the connectedness of genes might influence patterns of selection, we examined correlations between PPI and the observed statistics. We assessed significance using 10,000 permutations of connectedness across the test statistic as above. We  $\log_{10}$ -transformed the data as necessary.

**Signal of association throughout the entire genome and within the CSRs.** To determine whether loci within the identified regions have functional significance, we tested for statistical associations with second-year height and fall and spring bud phenology using data collected from three replicated plantations. We estimated genotypic best linear unbiased predictors using linear mixed models (*lmer* function of the *lme4* R package, described above) as the phenotypes for GWAS. We used the same set of resequenced, unrelated individuals used described above, excluding the highly differentiated Tahoe, Willamette Valley, and far northern British Columbia samples because strong stratification can lead to spurious associations<sup>80</sup>, leaving 498 individuals. We tested phenotypic association only with SNPs having a MAF  $\geq 0.05$ , leaving 5,939,334 SNPs. The analysis was performed for single traits in each plantation using *emmax*<sup>36</sup> and the identity by state kinship matrix to account for background genetic effects. To account for population structure, for each trait we included as covariates the principal-component axes that were significant predictors of the trait, chosen using stepwise regression (model selection based on Akaike's Information Criterion in the *step* function in the R package). We used the gemma multi-trait association model<sup>37</sup> to test for SNP association with each trait across all three plantations simultaneously, and in a nine-trait (three traits and three plantations) model as well. We used the mixed-model framework incorporating kinship and principal component axes that were significant (nominal  $\alpha = 0.05$ ) in a multivariate multiple linear regression.

We estimated  $\alpha$  values for association  $P$  values by permutation<sup>81</sup>. We permuted individual alleles among individuals, randomly generating genotypes while mirroring exactly the true MAF distribution. We then tested for association of these random genotypes with the observed phenotype data using the actual kinship matrix and principal components as above, thereby testing only the effect of randomly assigned genotypes while holding the structure of population stratification, relatedness, and the phenotypes constant. For univariate analyses in *emmax* we performed  $10^8$  permutations. For gemma multitrait analyses, we used  $>10^8$  permutations for bud set and height and  $8\text{--}33 \times 10^6$  permutations for bud flush and the nine-trait model, which were computationally more intensive. For each trait, we then estimated the cutoffs at various  $\alpha$  levels (**Supplementary Table 15**).

To determine whether the observed associations within the selection outliers was greater than expected by chance, we used the  $-\log_{10}(P \text{ value})$  as the

association signal within each selection outlier, and used the average of these values for each trait. We then randomly sampled the same number of 1-kb bins from throughout the genome 20,000 times. The number of random samples with a mean equal to or greater than the observed for each trait represents the probability of finding a median association signal in the selection outliers by chance alone. We also calculated the empirical  $P$  value for each CSR using the distribution of average association  $P$  values within 1-kb windows throughout the genome. This was done while controlling for the distribution of gene density within the surrounding 100 kb of the selection scans (**Supplementary Fig. 11g**). We also repeated this with a 50-kb window and without controlling for gene density and found the same patterns (data not shown).

61. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
62. Li, H. A statistical framework for SNP calling, mutation discovery, association mapping and population genetical parameter estimation from sequencing data. *Bioinformatics* **27**, 2987–2993 (2011).
63. Geraldes, A. *et al.* A 34K SNP genotyping array for *Populus trichocarpa*: design, application to the study of natural populations and transferability to other *Populus* species. *Mol. Ecol. Resour.* **13**, 306–323 (2013).
64. Huang, X. *et al.* Genome-wide association studies of 14 agronomic traits in rice landraces. *Nat. Genet.* **42**, 961–967 (2010).
65. Jiao, Y. *et al.* Genome-wide genetic changes during modern breeding of maize. *Nat. Genet.* **44**, 812–815 (2012).
66. Eckenwalder, J.E. in *Biology of Populus and its Implications for Management and Conservation* (eds. Stettler, R.F. *et al.*) 7–32 (NRC Research Press, 1996).
67. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
68. Patterson, N., Price, A.L. & Reich, D. Population structure and eigenanalysis. *PLoS Genet.* **2**, e190 (2006).
69. Yang, J., Lee, S.H., Goddard, M.E. & Visscher, P.M. GCTA: a tool for genome-wide complex trait analysis. *Am. J. Hum. Genet.* **88**, 76–82 (2011).
70. Goudet, J. & Büchi, L. The effects of dominance, regular inbreeding and sampling design on  $Q(ST)$ , an estimator of population differentiation for quantitative traits. *Genetics* **172**, 1337–1347 (2006).
71. Wang, T., Hamann, A., Spittlehouse, D.L. & Murdock, T.Q. ClimateWNA—high-resolution spatial climate data for western North America. *J. Appl. Meteorol. Climatol.* **51**, 16–29 (2012).
72. Charlesworth, B. Measures of divergence between populations and the effect of forces that reduce variability. *Mol. Biol. Evol.* **15**, 538–543 (1998).
73. Johnson, P.L.F. & Slatkin, M. Accounting for bias from sequencing error in population genetic estimates. *Mol. Biol. Evol.* **25**, 199–206 (2008).
74. Sabeti, P.C., Reich, D. & Higgins, J. Detecting recent positive selection in the human genome from haplotype structure. *Nature* **419**, 832–837 (2002).
75. Delaneau, O., Zagury, J.-F. & Marchini, J. Improved whole-chromosome phasing for disease and population genetic studies. *Nat. Methods* **10**, 5–6 (2013).
76. Browning, S.R. & Browning, B.L. Haplotype phasing: existing methods and new developments. *Nat. Rev. Genet.* **12**, 703–714 (2011).
77. Browning, B.L. & Browning, S.R. A unified approach to genotype imputation and haplotype-phase inference for large data sets of trios and unrelated individuals. *Am. J. Hum. Genet.* **84**, 210–223 (2009).
78. Pyhäjärvi, T., Hufford, M.B., Mezouk, S. & Ross-Ibarra, J. Complex patterns of local adaptation in teosinte. *Genome Biol. Evol.* **5**, 1594–1609 (2013).
79. Zhang, H. *et al.* PlantTFDB 2.0: update and improvement of the comprehensive plant transcription factor database. *Nucleic Acids Res.* **39**, D1114–D1117 (2011).
80. Price, A.L., Zaitlen, N.A., Reich, D. & Patterson, N. New approaches to population stratification in genome-wide association studies. *Nat. Rev. Genet.* **11**, 459–463 (2010).
81. Dudbridge, F. & Gusnanto, A. Estimation of significance thresholds for genomewide association scans. *Genet. Epidemiol.* **32**, 227–234 (2008).