

File S1 Materials and Methods

The following text represents supplemental information with respect to the **Materials and Methods**. Citations are found at the end of the Supplemental Text in this document. Supplemental figures and tables are found after the Supplemental Text in this document.

Development and application of genetic markers

Construction of EST clusters: An internally developed primer design package (*wt_primer*) was used to design polymerase chain reaction (PCR) primers from the 20,500 unique cDNA cluster consensus sequences. In total, 14,000 primer pairs were successfully designed with the following default conditions: primer length of 16-28 bp, a maximum melting temperature of 66°C, and a maximum difference of 5°C between melting temperatures of forward (F) and reverse (R) primers to maximize the likelihood that amplicons would amplify using standard PCR pipelines. Amplicon size was set to 450 bp, and for longer contigs, multiple overlapping amplicons were designed and the primer pair with the best score was selected for validation. If a predicted contig was less than 450 bp, the largest possible amplicon was chosen for further investigation. The best-scoring oligo pairs were tagged with M13F (GTAAAACGACGGCCAGT) and M13R (CAGGAAACAGCTATGACC) primers for high-throughput sequencing. Primers were validated using loblolly pine DNA from a single tree at a concentration of 2.5 ng/μl. The resulting PCR product was sequenced with M13F and M13R. Sequence quality was assessed using overall PHRED quality score of 20 and signal strength. Passing sequences were compared to the loblolly pine cluster consensus sequences using BLAST to ensure specificity.

Re-sequencing and SNP discovery: Genomic DNA was amplified in 384-well format PCR setup. Each PCR reaction contained 10 ng DNA, 1x HotStar buffer, 0.8 mM dNTPs, 1 mM MgCl₂, 0.2U HotStar enzyme (Qiagen) and 0.2 uM F and R primers in a 10 ul total reaction volume. PCR cycling parameters were: one cycle of 95°C for 15 min, 35 cycles of 95°C for 20 s, 60°C for 30 s and 72°C for one min, followed by one cycle of 72°C for three min. The resultant PCR products were purified using solid phase reversible immobilization chemistry followed by dye-terminator fluorescent sequencing with universal M13 primers. Sequencing reactions proceeded as follows: 95°C for 15 min to start followed by 40 cycles of 95°C for 10 s, 50°C for five seconds, 60°C for 2.5 min. Reactions were cleaned using solid phase reversible immobilization (Beckman Coulter Genomics) and the resulting sequencing fragments were detected via capillary electrophoresis using ABI Prism 3730xl DNA analyzers (Applied Biosystems, Foster City CA).

A customized pipeline, PineSAP (Wegrzyn *et al.* 2009), which employs PHRED/PHRAP (Ewing *et al.* 1998, Ewing and Green 1998), CONSED (Gordon *et al.* 1998), POLYBAYES (Marth *et al.* 1999), POLYPHRED (Nickerson *et al.* 1997), and machine

learning tools was used to generate sequence alignments and identify polymorphisms for these data. Custom scripts were added to the PHRED/PHRAP pipeline to identify the downstream amplification primer regions and to quantify the presence of secondary signal in the chromatograms. This allowed for improved screening of sequencing primer sequence and sequence due to mispriming in the PCR product and allowed us to reject chromatograms with high secondary signal that could indicate the presence of signal due to unintentional amplification of a paralogous locus.

Specifically, an integrated automatic and human input pipeline was designed to identify amplicons as putatively paralogous using the ratio of primary to secondary signal at each peak in a chromatogram. One sign that a pair of amplification primers is amplifying more than one site, given that we are working with haploid material, is the presence of secondary peaks in the chromatograms. We generated our assemblies using a single forward and reverse read, whereas programs that are designed to look for the presence of secondary peaks and call heterozygous SNPs often require more than two reads for input. Samples with many low secondary peaks can also indicate the presence of a paralog that may not be amplifying as strongly as the primary sequence. To address both of these issues, we ran PHRED and instructed it to generate a .poly file for each chromatogram (output intended for the POLYPHRED program) that contains information about the signal strength of each of the four fluorophores at the time each base was called. We then use the information in the .poly file to create custom tags that appear in CONSED that indicate the relative intensity of the most prominent secondary signal present. These tags are then visible when reviewing the trace files, and can also be detected automatically. If we look at the ratio of secondary to primary signal, if it was greater than 0.6, it was marked as 'high', if not but greater than 0.47, it was marked as 'medium', and if not but greater than 0.35, it was marked as 'low'. We then rejected any sample that had one or more 'high' sites, more than 2 'medium' sites or more than 6 'low' sites.

SNP genotyping: Total genomic DNA for each sample for genotyping was obtained from either pooled megagametophytes or needle tissue using Qiagen 96-well DNeasy Plant Mini Kits. Arrays were imaged on a Bead Array reader (Illumina) and genotype calling was performed using BeadStudio v. 3.1.3.0 (Illumina). Sample independent controls were assessed on each array to ensure assay integrity. We used a threshold of 55% for the call rate (CR) and 0.15 for the GenCall50 (GC50) scores for inclusion of SNP amplicons in the final dataset (see Eckert *et al.* 2009). More information about these data is available in Eckert *et al.* (2010a).

Generation of linkage maps

We utilized a maximum likelihood method that allows for genotyping errors (Cartwright *et al.* 2007) to create linkage maps comprised of SNPs and restriction fragment length polymorphism (RFLP) framework markers (see Eckert *et al.* 2009b and references therein) following a double pseudo-testcross strategy (Grattapaglia and Sederoff 1994). Each marker was tested for

segregation distortion using a goodness-of-fit test (significance threshold: $P = 0.01$). Once distorted markers were removed from each pedigree ($n = 35$ QTL, $n = 32$ base), phasing was performed for each pedigree separately, markers were grouped using a minimum LOD of 5.0 and a maximum distance of 20 cM, and marker ordering was improved within groups following the methods described by Cartwright *et al.* (2007). The Kosambi mapping function was used subsequently to transform recombination frequencies into map distances (Kosambi 1944). Resulting linkage maps were deposited in the TreeGenes CMAP database (<http://dendrome.ucdavis.edu/cmap/>) and first appeared in Eckert *et al.* (2010a,b).

Phenotypic trait analysis

Expression levels of lignin and cellulose related genes (expression): Transcript levels of 112 genes putatively involved with lignin and cellulose production were determined for each of the replicated genotypes using quantitative real-time PCR (qRT-PCR). Transcript levels were quantified with duplicate reactions (i.e. technical replicates) carried out on a GeneAmp 7900HT Sequence Detection System (Applied Biosystems, Carlsbad, CA, USA) using SYBER-Green PCR Master Mix (Applied Biosystems). Two ramets per genotype were used as biological replicates, and the final estimate of transcript level per gene for each genotype was the arithmetic average of the biological and technical replicates. All expression levels were standardized to 18S rRNA and β -actin controls. This process quantified transcript levels for 112 genes putatively involved with lignin and cellulose production for 400 of the 498 targeted genotypes (Palle *et al.* 2011).

Primary metabolite concentrations (metabolite): Gas chromatography coupled with time-of-flight mass spectrometry (GC-TOF-MS) was used to determine primary metabolite concentrations from pulverized xylem tissue collected from each replicated genotype established in the NCSU common garden. All analysis was conducted at the University of California Davis Genome Center Metabolomics Core Facility (<http://www.metacore.ucdavis.edu/techno1>). Two technical and biological replicates were used in these analyses. Resulting GC-TOF-MS data were processed following the methods outlined by Fiehn *et al.* (2008). Mixed linear models were used subsequently to adjust clonal least-square means for evaluation dates and experimental design. This process resulted in concentration estimates for 292 primary metabolites, of which 82 were known, assayed in 297 of the 498 targeted genotypes (Eckert *et al.* 2012).

Drought-tolerance and growth (drought): Estimates of carbon isotope ratios, height and foliar nitrogen content were assayed in each of the replicated genotypes established at the NCSU common garden. Isotope (^{13}C and ^{12}C) and nitrogen content (%N) analyses were based on 3 mg of needle tissue and were carried out at the COIL (<http://www.cobsil.com/>) stable isotope facility located at Cornell University. Total tree height (cm) was determined at the end of the second growing season. Phenotypic values for each genotype were estimated using mixed linear models that accounted for experimental design and

spatial heterogeneity in the common garden. This process resulted in estimates of $\delta^{13}\text{C}$, foliar nitrogen content and height after the second year for 425 of the 498 targeted genotypes (Cumbie *et al.* 2011).

Disease resistance (disease): Lengths of lesions (mm) produced in response to inoculation with pitch canker (*Fusarium circinatum* Nirenberg & O'Donnell) were estimated at four, eight and 12 weeks post inoculation for each of the replicated genotypes established in the NCSU common garden. These estimates were taken as a measure of disease resistance, and phenotypic values for each genotype were estimated as best linear unbiased predictors (BLUPs) using mixed linear models incorporating effects due to experimental design. This process resulted in lesion length estimates at multiple time points for 404 of the targeted 498 genotypes of which the estimates at 12 weeks post inoculation were used in association analyses (Quesada *et al.* 2010).

Identification of phenotypic associations: Phenotypic associations were identified using a two-stage approach where clonal values were predicted and then associated with SNPs using linear models. Most often these were general or mixed linear models (Yu *et al.* 2006) with fixed effects for SNPs as implemented in the program TASSEL (Bradbury *et al.* 2007), although Eckert *et al.* (2012) used the method of Price *et al.* (2006). All analyses included effects due to population structure, as described in Eckert *et al.* (2010a), and were largely based on single locus tests of fixed effects for each SNP. When kinship was included, the kinship matrix was estimated using EMMA (Kang *et al.* 2008) and included as a matrix of random effects. Exceptions to this were the multilocus models used by Quesada *et al.* (2010) and Eckert *et al.* (2012). Multiple testing was accounted for during single locus testing using the positive false discovery rate (FDR) with a threshold of $Q = 0.10$ (Storey and Tibishirani 2003). Additional information about the statistical methodologies used for each phenotypic trait is given in the original publications (disease resistance: Quesada *et al.* 2010; drought-tolerance and growth: Cumbie *et al.* 2011; primary metabolites: Eckert *et al.* 2012) or in the supplemental online materials (gene expression: Palle *et al.* 2013).

Results

The following text represents supplemental information with respect to the **Results**. Citations are found at the end of the Supplemental Text in this document. References to figures and tables are for those in the main text unless noted otherwise. Supplemental figures and tables are found after the Supplemental Text in this document.

Re-sequencing data summary

The number of amplicons passing design thresholds decreased from approximately 7,900 to 7,413 after requiring both F and R reads to be present for each sample followed by a further decrease to 6,669 amplicons after screening for amplification primers in both reads. A total of 5,773 amplicons passed our final quality thresholds, which also included screens for organellar contamination. The average (± 1 standard deviation [sd]) sample size per amplicon was 12 (± 6), with the frequency distribution of sample sizes being skewed towards larger sample sizes (Figure S1). The average sample size also

changed little across different categories of amplicons. Of these 5,773 amplicons, 1,306 could be positioned on the linkage map (22.6%), 2,626 could be annotated to level of coding and noncoding regions (45.5%), 3,484 had a putative ortholog for radiata pine (60.3%) and 950 had a putative ortholog for sugar pine (16.4%). Only a moderate fraction of the total number of amplicons (45.4%) and the number of amplicons that were annotated (44.6%) were represented by at least one SNP on the genotyping chip. Of the total number of amplicons represented on the genotyping chip ($n = 2,619$), 689 (26.3%) unique amplicons had at least one SNP associated to at least one phenotype. The fraction of annotated and unannotated amplicons on the genotyping chip ($n = 1,173$ and $1,446$, respectively) with at least one SNP associated to at least one phenotype were similar, with 309 (26.3%) and 380 (26.3%) being associated to at least one phenotype. For the 689 amplicons containing at least one SNP associated to at least one phenotype, 76 were associated to expression, 576 to metabolite, 12 to drought, and nine to disease related phenotypes. The remaining 16 amplicons were associated to environmental variables (see Eckert *et al.* 2010a, 2010b). One hundred and ninety-five (28.3%) out of the 689 amplicons were also associated with more than one phenotype and/or environmental variable (range: 1 to 6).

At the level of individual sites, a total of 2,135,607 aligned sites were analyzed across the 5,773 amplicons. The average (± 1 sd) length of amplicons was 370 (± 126) bp, with longer amplicons being more likely to have at least one SNP genotyped (Table 1) and more likely to be annotated. Of these sites, 1,161,888 could be annotated (54.4%), with 583,159 (50.2%), 160,814 (13.8%), and 417,915 (36.0%) sites being nonsynonymous, synonymous, and noncoding, respectively. The higher percentage of annotated sites relative to annotated amplicons is accounted for by the observation that annotated amplicons were longer than the genome-wide average (432 bp versus 370 bp). This difference, however, was not statistically significant ($P_{\text{perm}} = 0.348$). The same patterns were observed for classes of amplicons (Table 1).

A total of 22,621 SNPs were detected across the 5,773 amplicons. There was little to no effect of sequence quality on the number of SNPs per amplicon, as the correlation between the number of SNPs called for each amplicon using a PHRED threshold of 30 versus a PHRED threshold of 40 was large (Pearson's $r > 0.85$ for all, nonsynonymous, synonymous and noncoding SNPs) and the slope was almost equal to one (Figure S2). Coverage, however, affected several alignment metrics related to nucleotide diversity and divergence (Figures S3-S5; Tables S1-S3). This level of polymorphism was similar to that reported previously for loblolly pine (Brown *et al.* 2004, González-Martínez *et al.* 2006), with one SNP per 94 bp on average. Of these 22,621 SNPs, 10,591 could be annotated as nonsynonymous ($n = 2,915$), synonymous ($n = 3,233$) and noncoding ($n = 4,443$). On a per site basis, SNPs were more common at synonymous sites (one SNP per 50 bp on average) and noncoding (one SNP per 94 bp on average) relative to nonsynonymous (one SNP per 200 bp on average) sites. Patterns were similar across different categories of amplicons (Table 1). At the level of SNPs selected genotyping ($n = 7,216$), there was no enrichment of

certain types of SNPs in the set of associated amplicons ($n = 873$ SNPs associated to at least one phenotype), so that the numbers of nonsynonymous ($n = 127$), synonymous ($n = 160$) and noncoding SNPs ($n = 201$) associated to at least one phenotype were no different than those expected by randomly subsampling the annotated SNPs on the Illumina genotyping array ($P_{\text{perm}} > 0.15$).

Linkage Disequilibrium

Genome-wide patterns: Intragenic linkage disequilibrium, as assessed with Kelly's Z_{ns} (Kelly 1997), was positively correlated with nucleotide diversity (Spearman's $\rho > 0.30$, $P_{\text{perm}} < 0.005$), while it was negatively correlated with the number of haplotypes (Spearman's $\rho = -0.427$, $P_{\text{perm}} = 0.008$). Breaking the range of observed values for Z_{ns} into high and low categories, the correlation becomes significantly negative between nucleotide diversity and linkage disequilibrium when $Z_{\text{ns}} > 0.50$ (Spearman's $\rho < -0.20$, $P_{\text{perm}} < 0.05$). Correlations with nucleotide diversity at different categories of sites were approximately 2.5-fold smaller, yet still positive, and non-significant (Spearman's $\rho < 0.05$, $P_{\text{perm}} > 0.20$). These correlations, however, changed when considering only amplicons with $Z_{\text{ns}} > 0.50$, so that nucleotide diversity at nonsynonymous and noncoding sites was significantly, negatively correlated with linkage disequilibrium when $Z_{\text{ns}} > 0.50$ (Spearman's $\rho < -0.25$, $P_{\text{perm}} < 0.05$). Correlations of linkage disequilibrium with nucleotide divergence were close to zero and non-significant ($-0.05 < \text{Spearman's } \rho < 0.05$, $P_{\text{perm}} > 0.40$), even when breaking Z_{ns} into low and high categories ($-0.10 < \text{Spearman's } \rho < 0.10$, $P_{\text{perm}} > 0.35$).

Comparisons across categories of amplicons: Linkage disequilibrium varied across sets of amplicons defined by whether or not they were located on a linkage map (Mann-Whitney U-test: $P = 0.0371$, $P_{\text{perm}} = 0.011$), whether or not they were annotated (Mann-Whitney U-test: $P = 1.486e-05$, $P_{\text{perm}} < 0.001$), and whether or not they were associated with at least one phenotype (Mann-Whitney U-test: $P = 0.0493$, $P_{\text{perm}} = 0.026$). On average (± 1 sd), linkage disequilibrium was higher for amplicons that were mapped ($Z_{\text{ns}}: 0.327 \pm 0.276$ vs. 0.304 ± 0.292), while it was lower for those that were annotated ($Z_{\text{ns}}: 0.302 \pm 0.287$ vs. 0.336 ± 0.288) and for those associated with at least one phenotype ($Z_{\text{ns}}: 0.282 \pm 0.263$ vs. 0.313 ± 0.277). Significant differences in the level of linkage disequilibrium were also noted among amplicons grouped into categories based on the types of phenotypes to which they were associated (Kruskal-Wallis rank sum test: $P = 0.032$, $P_{\text{perm}} = 0.012$), with amplicons associated with disease phenotypes having the lowest ($Z_{\text{ns}} = 0.144$) and amplicons associated with expression phenotypes having the highest ($Z_{\text{ns}} = 0.308$) average levels of linkage disequilibrium. In general, correlations between levels of linkage disequilibrium and diversity and divergence estimates within different categories of amplicons were similar to genome-wide patterns.

Literature Cited

- Bradbury, P. J., Z. Zhang, D. E. Kroon, T. M. Casstevens, Y. Ramdoss, *et al.*, 2007 TASSEL: software for association mapping of complex traits in diverse samples. *Bioinformatics* 23: 2633–2635.
- Brown, G. R., G. P. Gill, R. J. Kuntz, C. H. Langley, and D. B. Neale, 2004 Nucleotide diversity and linkage disequilibrium in loblolly pine. *Proc. Natl. Acad. Sci. USA* 101: 15255–15260.
- Cartwright, D. A., M. Troggio, R. Velasco, and A. Gutin, 2007 Genetic mapping in the presence of genotyping errors. *Genetics* 176: 2521–2527.
- Cumbie, W. P., A. J. Eckert, J. L. Wegrzyn, R. Whetten, D. B. Neale, *et al.*, 2011 Association genetics of carbon isotope discrimination, height, and foliar nitrogen in a natural population of *Pinus taeda* L. *Heredity* 107: 105–114.
- Eckert, A. J., B. Pande, E. S. Ersöz, M. H. Wright, V. K. Rashbrook, *et al.*, 2009b High-throughput genotyping and mapping of single nucleotide polymorphisms in loblolly pine (*Pinus taeda* L.). *Tree Genet. Genomes* 5: 225–234.
- Eckert, A. J., J. van Heerwaarden, J. L. Wegrzyn, C. D. Nelson, J. Ross-Ibarra, *et al.*, 2010a Patterns of population structure and environmental associations to aridity across the range of loblolly pine (*Pinus taeda* L., Pinaceae). *Genetics* 185: 969–982.
- Eckert, A. J., A. D. Bower, S. C. González-Martínez, J. L. Wegrzyn, G. Coop, *et al.*, 2010b Back to nature: Ecological genomics of loblolly pine (*Pinus taeda*, Pinaceae). *Mol. Ecol.* 19: 3789–3805.
- Eckert, A. J., J. L. Wegrzyn, W. P. Cumbie, B. Goldfarb, D. A. Huber, *et al.*, 2012 Association genetics of the loblolly pine (*Pinus taeda*, Pinaceae) metabolome. *New Phytol.* 193: 890–902.
- Ersöz, E. S., M. H. Wright, S. C. González-Martínez, C. H. Langley, and D. B. Neale, 2010 Evolution of disease response genes in loblolly pine: Insights from candidate genes. *PLoS ONE* 5: e14234.
- Ewing, B., and P. Green, 1998 Base-calling of automated sequencer traces using PHRED. II. Error probabilities. *Genome Res.* 8: 186–194.
- Ewing, B., L. Hillier, M. C. Wendl, and P. Green, 1998 Base-calling of automated sequencer traces using PHRED. I. Accuracy assessment. *Genome Res.* 8: 175–185.
- Fiehn, O., G. Wohlgemuth, M. Scholz, T. Kind, D. Y. Lee, *et al.*, 2008 Quality control for plant metabolomics: reporting MSI-compliant studies. *Plant J.* 53: 691–704.
- González-Martínez, S. C., E. Ersöz, G. R. Brown, N. C. Wheeler, and D. B. Neale, 2006a DNA sequence variation and selection of tag SNPs at candidate genes for drought-stress response in *Pinus taeda*. *Genetics* 172: 1915–1926.
- Gordon, D., C. Abajian, and P. Green, 1998 Consed: A graphical tool for sequence finishing. *Genome Res.* 8: 195–202.

- Grattapaglia, D., and R. Sederoff, 1994 Genetic linkage maps of *Eucalyptus grandis* and *E. urophylla* using a pseudo-testcross mapping strategy and RAPD markers. *Genetics* 137: 1121–1137.
- Hudson, R. R., 1990 Gene genealogies and the coalescent process. *Oxford Surv. Evol. Biol.* 7: 1–44.
- Kang, H. M., N. A. Zaitlen, C. M. Wade, A. Kirby, D. Heckerman, *et al.*, 2008 Efficient control of population structure in model organism association mapping. *Genetics* 178: 1709–1723.
- Kelly, J., 1997 A test of neutrality based on interlocus associations. *Genetics* 146: 1197–1206.
- Kosambi, D. D., 1944 The estimation of map values from recombination values. *Ann. Eugen.* 12: 172–175.
- Marth, G. T., I. Korf, M. D. Yandell, R. T. Yeh, Z. Gu, *et al.*, 1999 A general approach to single nucleotide polymorphism discovery. *Nat. Genet.* 23: 452–456.
- Nickerson, D. A., V. O. Tobe, and S. L. Taylor, 1997 PolyPHRED: Automating the detection and genotyping of single nucleotide substitutions using fluorescence-based re-sequencing. *Nucleic Acids Res.* 25: 2745–2751.
- Palle, S. R., C. M. Seeve, A. J. Eckert, W. P. Cumbie, B. Goldfarb, *et al.*, 2011 Natural variation in expression of genes involved in xylem development in loblolly pine (*Pinus taeda* L.). *Tree Genet. Genomes* 7: 193–206.
- Palle, S. R., C. M. Seeve, A. J. Eckert, J. L. Wegrzyn, D. B. Neale, and C. A. Loopstra, 2013 Association of loblolly pine xylem development gene expression with single nucleotide polymorphisms. *Tree Physiol.* 33: 763–774.
- Price, A. L., N. J. Patterson, R. M. Plenge, M. E. Weinblatt, N. A. Shadick, *et al.*, 2006 Principal components analysis corrects for stratification in genome-wide association studies. *Nat. Genet.* 38: 904–909.
- Quesada, T., V. Gopal, W. P. Cumbie, A. J. Eckert, J. L. Wegrzyn, *et al.*, 2010 Association mapping of quantitative disease resistance in a natural population of loblolly pine (*Pinus taeda* L.). *Genetics* 186: 677–686.
- 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.
- Stoletzki, N., and A. Eyre-Walker, 2011 Estimation of the neutrality index. *Mol. Biol. Evol.* 28: 63–70.
- Storey, J. D., and R. Tibshirani, 2003 Statistical significance for genome-wide studies. *Proc. Natl. Acad. Sci. USA* 100: 9440–9445.
- Wegrzyn, J. L., J. M. Lee, J. D. Liechty, and D. B. Neale, 2009 PineSAP - Pine alignment and SNP Identification Pipeline. *Bioinformatics* 25: 2609–2610.
- Yu, J., G. Pressoir, W. H. Briggs, I. V. Bi, M. Yamasaki, *et al.*, 2006 A unified mixed-model method for association mapping that accounts for multiple levels of relatedness. *Nat. Genet.* 38: 203–208.