

Generation and analysis of an *Eucalyptus globulus* cDNA library constructed from seedlings subjected to low temperature conditions

Susana Rasmussen-Poblete

Departamento de Ciencias Biológicas
Universidad Andrés Bello
República 252, Santiago, Chile
Tel: 56 2 2383178
Fax: 56 2 2372259
E-mail: susana.rasmussen@bionova.cl

Jorge Valdés

Departamento de Ciencias Biológicas
Universidad Andrés Bello
República 252, Santiago, Chile
Tel: 56 2 2383178
Fax: 56 2 2372259
E-mail: jorge.valdes@gmail.com

Maria Cecilia Gamboa

Departamento de Ciencias Biológicas
Universidad Andrés Bello
República 252, Santiago, Chile
Tel: 56 2 2383178
Fax: 56 2 2372259
E-mail: cecy_gamboa@yahoo.es

Pablo D.T. Valenzuela

Departamento de Ciencias Biológicas
Universidad Andrés Bello
República 252, Santiago, Chile
Tel: 56 2 2398969
Fax: 56 2 2372259
E-mail: pvalenzu@bionova.cl

Erwin Krauskopf*

Departamento de Ciencias Biológicas
Universidad Andrés Bello
República 252, Santiago, Chile
Tel: 56 2 2383178
Fax: 56 2 2383178
E-mail: ekrauskopf@unab.cl

Financial support: This work was partially funded by Universidad Andrés Bello (DI Proyect: 04-05/1) and MIFAB (Proyect: P04-071-F) and by the Microsoft Joint Research Program.

Keywords: cellulose, cold-temperature, EST database, forest biotechnology, lignin.

Note: The sequences have been deposited in GenBank. Accession numbers: ES588357-ES597093

Abbreviations: 4CL: CoA ligase
AUX/IAA: auxin/indole-3-acetic acid
bZIP: basic leucine zip
C3H: p-coumarate 3-hydroxylase
C4H: 4-hydroxylase
CAD: cinnamyl alcohol dehydrogenase
CCR: cinnamoyl CoA reductase
CCoAOMT: caffeoyl-CoA 3-O-methyltransferase
ESTs: expressed sequence tags
F5H: ferulate 5-hydroxylase
GO: gene ontology
HCT: hydroxycinnamoyltransferase
PAL: phenylalanine ammonia lyase

***Eucalyptus globulus* is the most important commercial temperate hardwood in the world because of its wood properties and due to its characteristics for biofuel production. However, only a very low number of expressed sequence tags (ESTs) are publicly available for this tree species. We constructed a cDNA from *E. globulus* seedlings subjected to low temperature and sequenced 9,913 randomly selected clones, generating 8,737 curated ESTs. The assembly produced 1,062 contigs and 3,879 singletons forming a *Eucalyptus* unigene set. Based on BLASTX analysis, 89.3% of the contigs and 88.5% of the singletons had significant similarity to known genes in the non-redundant database of GenBank. The *Eucalyptus* unigene set generated is a valuable public resource that provides an initial model for genes and regulatory pathways involved in cell wall biosynthesis at low temperature.**

Forests cover nearly 30% of the earth surface, nearly 4 billion hectares, serving multiple functions including conservation of biological diversity, renewing the oxygen supply of the atmosphere, preventing soil erosion and supplying pulp and wood (FAO, 2005). Forest tree breeding aims to improve the quality of trees by the selection of individuals with desirable traits that will later be used to produce trees with improved genotype. Genetic improvement programs such as controlled cross-pollination breeding have been used since the 1950s. Nevertheless, phenotype assessment is a complex process due to the long generation times of woody species (Grattapaglia, 2004). It is within the context of reducing this time-frame that functional genomics has become a powerful tool in forestry.

In the last few years functional genomics has been used extensively for gene discovery in species whose genomes have not been completely sequenced. A cost-effective and rapid way to obtain new data from an organism is through partial sequencing of randomly selected cDNA clones (Braütigam et al. 2005). The resulting collection of expressed sequence tags (ESTs) reveals a portion of genes in an organism expressed under a particular condition. Using this approach, several traits have been analyzed in trees, such as wood formation (Allona et al. 1998; Sterky et al. 1998; Israelsson et al. 2003) or cold tolerance (Nanjo et al. 2004; Sterky et al. 2004). Unfortunately, these studies have focused on gene expression profiles having a direct effect on the particular trait studied, without expanding the range of effects that the set condition might have on other metabolic pathways. In fact, cold stress in poplar cuttings (*Populus tremula* x *Populus tremuloides* cv. Mush1) has been shown to produce variations in parameters such as sucrose concentration and lignin content, illustrating the direct effect of cold conditions on wood quality (Hausman et al. 2000).

The amount and type of lignin and cellulose are important

in the timber and pulp industry as they have a direct effect on the chemical properties of the wood produced by the tree (Jung and Ni, 1998; Fukushima, 2001; Plomion et al. 2001). For the production of biofuels, cellulose needs to be separated from lignin so it can be made available for enzyme hydrolysis. Therefore, several research groups have studied different ways by which to modify lignin and cellulose content on the plant cell wall. As a result, various studies have shown a co-regulation of these two compounds (Hu et al. 1999; Li et al. 2003; Rastogi and Dwivedi, 2006). For instance, the down-regulation of a single lignin biosynthetic gene resulted in a decrease of lignin production by the plant, while exhibiting an increase in cellulose production (Hu et al. 1999). Hence, the modification of plant cell wall composition in trees may provide a way to engineer wood for biofuel production.

E. globulus is considered the most important temperate hardwood plantation species in the world due to its combination of wood properties suitable for the pulp and paper industry (Jones et al. 2002; Grattapaglia, 2004). This tree species has fast growth rates and an ability to adapt to a broad range of geographic locations (ranging from latitude 35°S to 42°S), even though its growth rate diminishes due to frost conditions (Jones et al. 2002; Miranda and Pereira, 2002). Most importantly, *Eucalyptus* has been listed as one of the candidate biomass energy crops by the U.S. Department of Energy (U.S. Department of Energy, 2007). Nevertheless, public genomic information from *E. globulus* is limited. In fact, an analysis of publicly available *E. globulus* ESTs at the GenBank EST repository (by July 06, 2007) registered only 3,953 ESTs for *E. globulus* compared to the mostly represented tree, *Pinus taeda* (329,469 ESTs). Thus, in this study we provide and describe the first publicly available cDNA library from cold-treated *E. globulus* seedlings, paying particular attention to genes predicted to be involved in cell wall biosynthesis and the transcription factors suggested to be involved in their regulation).

MATERIALS AND METHODS

Plant material

E. globulus seeds were germinated in a soil mixture and grown in a culture cabinet with a 16 hrs day/8 hrs night photoperiod at a temperature of 21°C. The library was constructed from 3-month old *Eucalyptus globulus* plants maintained at 4°C degrees for 30 min. After cold treatment, *E. globulus* leaves were collected and frozen in liquid nitrogen until use.

RNA extraction and cDNA library construction

Total RNA was extracted according to the method described by Chang and colleagues (Chang et al. 1993). RNA integrity was confirmed by gel electrophoresis and 1

*Corresponding author

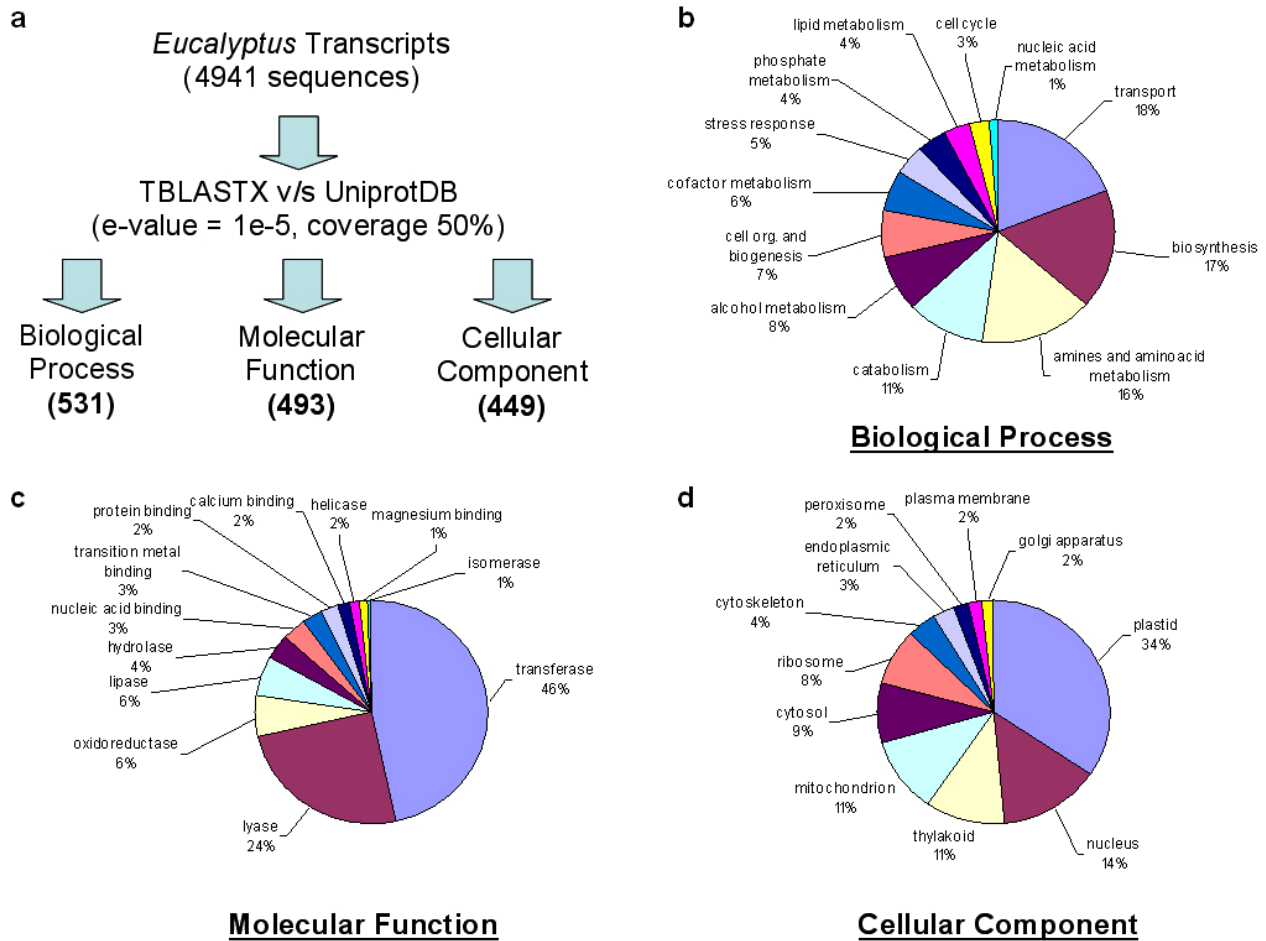


Figure 1. Functional categorization of *E. globulus* unigenes. (a) Schematic representation of the functional categorization process. (b-d) Distribution of *E. globulus* unigenes across GO categories. Parent categories and their percentages are shown in bold, sub-categories and the number of deduced proteins is shown in parenthesis.

mg was quantified using a RNA standard (Invitrogen, Cat 15620-016). Poly (A) mRNA was isolated from total RNA with the Stratagene Poli (A) Quick mRNA Isolation Kit (Stratagene, La Jolla, CA, USA). cDNA was prepared and cloned using the vector pExpress I exploiting the Not I and Eco RV restriction sites. The cDNA library was not normalized, *i.e.* no attempt was made to reduce the redundancy of highly expressed transcripts.

EST sequencing, filtering and assembly

In total, 9,913 bacterial colonies were randomly picked and single-pass sequence reactions performed. These sequences were analyzed using Phred base calling software (with $Q > 20$) (Ewing et al. 1998). All traces were subjected to a trimming process for the removal of ribosomal RNA, poly (A) tails, low-quality sequences, vector and adapter regions. Sequences with 94% of identity over 40 or more nucleotides were assembled using the CAP3 software (Huang and Madan, 1999).

Unigene function assignment and categorization

The unigene set was classified and analyzed according to gene ontology (GO) terms (Ashburner et al. 2000) across functional categories. The complete unigene set was compared against the protein non-redundant database using BLASTX (Altschul et al. 1997) and analyzed with the InterProScan program (Zdobnov and Apweiler, 2001) to assign a putative function. GO terms were extracted from the best hits obtained from the BLASTX comparison against SwissProt-Trembl database (Fleischmann et al. 1999) ($E\text{-value} < E-15$ and $>70\%$ of alignment coverage) and compared to the InterProScan GO suggestions. All the GO assignments were curated manually (Ashburner et al. 2000). The unigene dataset was compared to other *Eucalyptus* cDNA libraries available in Genbank through BlastN program using an e-value cutoff of $E-5$.

RESULTS AND DISCUSSION

Analysis of *E. globulus* cDNA library

Table 1. Contigs with ESTs highly represented. Assigned function is indicated in contigs with more than 50 ESTs.

Number ESTs	Contigs length (nt)	Assigned function	Relative organism	Similarity (%)	Accession number (GI)
118	1604	Plastidic aldolase	<i>Nicotiana paniculata</i>	94.8	4827253
86	1509	Chloroplast latex aldolase-like protein	<i>Manihot esculenta</i>	90.9	56122688
81	979	Ribulose-1,5-bisphosphate carboxylase/oxygenase small subunit	<i>Panax ginseng</i>	93.2	77157637
80	1693	Ribulose-1,5-bisphosphate coarboxylase/oxygenase activase precursor	<i>Malus x domestica</i>	93.1	415852
73	1590	Glyceraldehydes-3-phosphate dehydrogenase A subunit	<i>Glycine max</i>	87.8	77540210
51	1721	AAA ATPase, central region; homeodomain-like	<i>Medicago truncatula</i>	89.1	92870561

of lignin monomers and cellulose. All of the genes known to participate in the lignin biosynthetic pathway are represented in our cDNA library. Two of the predicted gene products, *p-coumarate 3-hydroxylase* (C3H) and *CoA:shikimate/ quinate hydroxycinnamoyltransferase* (HCT) had not been described previously in any *Eucalyptus* species. However, genes encoding *trans-cinnamate 4-hydroxylase* (C4H), *ferulate 5-hydroxylase* (F5H) and *4-coumarate: CoA ligase* (4CL) had been described in other *Eucalyptus* species but not in *E. globulus* (Harakava, 2004). The remainder of the genes found had been previously described for *E. globulus* and published in GenBank, including *phenylalanine ammonia lyase* (PAL), *cinnamoyl CoA reductase* (CCR), *cinnamyl alcohol dehydrogenase* (CAD), *caffeic acid O-methyltransferase* (COMT) and *caffeoyl-CoA 3-O-methyltransferase* (CCoAOMT) (Figure 2) (Supplementary data 2).

The assembly of the C3H and HCT ESTs showed that two isoforms of their gene-products are represented in our cDNA library. C3H and HCT participate in the process of converting p-coumaroyl CoA into caffeoyl-CoA, resulting in the production of coniferyl (G) and sinapyl (S) lignin units. Down-regulation of C3H in transgenic alfalfa plants and *Arabidopsis* mutants resulted in a significant difference in lignin composition due to an alteration in the number and nature of the monolignol monomers (Franke et al. 2002; Ralph et al. 2006). The characterization of the *Arabidopsis reduced epidermal fluorescence* (*ref8*) mutant defective in C3H suggested that the genetic modification of this gene may not be appropriate for the reduction of lignin content in forest species because the mutant plants generated exhibited vascular collapse, developmental abnormalities and increased susceptibility to pathogen attack (Boerjan et al. 2003; Cooke et al. 2004).

Three unigenes exhibited similarity to known cellulose synthase genes. Analysis of their predicted domains by InterProScan revealed that all of them contained the cellulose synthase domain that is composed of three

aspartic residues and a QXXRW motif, playing a significant role in the catalytic activity of this enzyme (Krauskopf et al. 2005). However, the zinc finger domains (IPR001841 and IPR011011) present in cellulose synthase proteins were not found in our sequences since the sequences were not full-length. The deduced *E. globulus* proteins were compared with the ones previously described for *E. grandis* (Ranik and Myburg, 2006) as no sequences were available for *E. globulus* (Supplementary data 3).

Transcription factors involved in wood formation

Of the 56 transcription factor families described in *Arabidopsis* and 63 in rice (Guo et al. 2005; Gao et al. 2006), 11 of them were represented in our library: auxin/indole-3-acetic acid (AUX/IAA) family, B3 family, basic/helix-loop-helix (bHLH) family, basic leucine zip (bZIP) family, GRAS family, homeodomain-leucine zipper (HD-Zip) (HB) family, heat shock family (HSF), MYB family, WRKY family, zinc finger homeobox (ZF-HD) family and ZIM family. Transcription factors families such as AUX/IAA, MYB and HD containing domains (zinc finger proteins and homeodomain-leucine zipper) regulate the expression of genes that participate in xylem development and secondary wall formation (lignin and cellulose biosynthesis) (Oh et al. 2003; Cánovas et al. 2004).

Many of the genes encoding the enzymes of general phenylpropanoid metabolism, such as PAL, C4H, 4CL, COMT and CAD contain conserved motifs within their promoters that are recognized by plant MYB transcription factors (Tamagnone et al. 1998). Twelve members of the MYB family were found in our library. Some of them had a best BLASTX hit with *GOLDEN2-like 1* gene, *LHY-CCA1-like 5* gene and *DIVARICATA* gene. The coverage of the sequences with their best BLASTX hit ranged from 25% and 100%. Two *E. gunnii* MYB transcription factors sequences were found in GenBank [GenBank: AJ576023-AJ576024] (Goicoechea et al. 2005). Based on BLASTN

analysis, these sequences were different from the ones obtained in our library. Others families less represented in our library belonged to the ZF family and bZIP (with seven members each), WRKY family (five members with coverage of their best BLASTX hit between 12% and 50%) and one member of the AUX/IAA family, (Supplementary data 4).

In addition, the data gathered through these analyses was compared with the few existing *Eucalyptus* cDNA libraries currently found in GenBank. The comparison was made against *Eucalyptus gunnii* (8,538 ESTs), *Eucalyptus globulus subsp. bicostata* (2,685 ESTs), *Eucalyptus grandis* (1,574 ESTs) and *Eucalyptus globulus 'blue gum'* (1,266 ESTs). BlastN comparisons against our *E. globulus* database revealed a low level of similarity between our sequenced library and the available datasets. The number of sequences that have at least one match with E-values better than 1E-5 for each library were 1,335 ESTs for *E. gunnii* (15%), 464ESTs for *E. globulus subsp. bicostata* (17%), 267 for *E. grandis* (17%) and 261 ESTs for *E. globulus 'blue gum'* (17%).

In conclusion, a unigene set of approximately 4900 unigenes was obtained from our *E. globulus* cDNA library. Analysis of its content has provided valuable data for the future metabolic engineering of plant cell walls by identifying new potential targets that will allow future modification for biofuel production and industrial use. In addition, our results will be useful for comparative genomic studies among hardwoods and softwoods.

ACKNOWLEDGMENTS

We thank Dr. Danilo González who provided the computer cluster to generate the unigene set and Dr. David Holmes for critical reading of the manuscript.

REFERENCES

ALLONA, I.; QUINN, M.; SHOOP, E.; SWOPE, K.; ST. CYR, S.; CARLIS, J.; RIEDL, J.; RETZEL, E.; CAMPBELL, M.; SEDEROFF, R. and WHETTEN, R. Analysis of xylem formation in pine by cDNA sequencing. *Proceedings of the National Academy of Sciences of the United States of America*, August 1998, vol. 95, no. 16, p. 9693-9698.

ALTSCHUL, S.; MADDEN, T.; SCHAFFER, A.; ZHANG, J.; ZHANG, Z.; MILLER, W. and LIPMAN, D. Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research*, September 1997, vol. 25, no. 17, p. 3389-3402.

ASHBURNER, M.; BALL, C.; BLAKE, J.; BOTSTEIN, D.; BUTLER, H.; CHERRY, J.; DAVIS, A.; DOLINSKI, K.; DWIGHT, S.; EPPIG, J.; HARRIS, M.; HILL, D.; ISSEL-TARVER, L.; KASARSKIS, A.; LEWIS, S.; MATESE, J.; RICHARDSON, J.; RINGWALD, M.;

RUBIN, G. and SHERLOCK, G. Gene Ontology: tool for the unification of biology. *Nature Genetics*, May 2000, vol. 25, p. 25-29.

BOERJAN, W.; RALPH, J. and BAUCHER, M. Lignin biosynthesis. *Annual Review of Plant Biology*, June 2003, vol. 54, p. 519-546.

BRAÛTIGAM, M.; LINDLÖF, A.; ZAKHRABETKOVA, S.; GHARTI-CHHETRI, G.; OLSSON, B. and OLSSON, O. Generation and analysis of 9792 EST sequences from cold acclimated oat, *Avena sativa*. *BMC Plant Biology*, September 2005, vol. 5, p. 18.

CÁNOVAS, F.; DUMAS-GAUDOT, E.; RECORBET, E.; JORRIN, J.; MOCK, H.-P. and ROSSIGNOL, M. Plant proteome analysis. *Proteomics* 2004, vol. 4, p.285-298.

CHANG, S.; PURYEAR, J. and CAIRNEY, J. A simple and efficient method for isolating RNA from pines trees. *Plant Molecular Biology Reporter*, June 1993, vol. 11, p. 113-116.

COOKE, J.; MORSE, A. and DAVIS, J. Forestry. In: KLEE and P. CRISTOU eds. *Handbook of Plant Biotechnology*. United Kingdom. John Wiley and Sons, 2004, vol. 2, p. 881-904.

EWING, B.; HILLIER, L.; WENDL, M. and GREEN, P. Basecalling of automated sequencer traces using phred I. Accuracy Assessment. *Genome Research*, March 1998, vol. 8, no. 3, p. 175-185.

FAO. Progress towards sustainable forest management. In: FOOD AND AGRICULTURE ORGANIZATION OF THE UNITED NATIONS eds. *Global forest resources assessment*. Rome. FAO forestry paper, 2005, p. 147.

FLEISCHMANN, W.; MOELLER, S.; GATEAU, A. and APWEILER, R. A novel method for automatic functional annotation of proteins. *Bioinformatics*, March 1999, vol. 15, no. 3, p. 228-233.

FRANKE, R.; HUMPHREYS, J.M.; HEMM, M.R.; DENAULT, J.W.; RUEGGER, M.O.; CUSUMANO, J.C. and CHAPPLE, C. Changes in secondary metabolism and deposition of an unusual lignin in the *ref8* mutant of *Arabidopsis*. *The Plant Journal*, April 2002, vol. 30, no. 1, p. 33-45.

FUKUSHIMA, K. Regulation of syringyl to guaiacyl ratio in lignin biosynthesis. *Journal of Plant Research*, February 2001, vol. 114, no. 4, p. 499-508.

GAO, G.; ZHONG, Y.; GUO, A.; ZHU, Q.; TANG, W.; ZHENG, W.; GU, X.; WEI, L. and LUO, J. DRTF: a database of rice transcription factors. *Bioinformatics*, March 2006, vol. 22, no. 10, p. 1286-1287.

- GOICOECHEA, M.; LACOMBE, E.; LEGAY, S.; MIHALJEVIC, S.; RECH, P.; JAUNEAU, A.; LAPIERRE, C.; POLLET, B.; VERHAEGEN, D.; CHAUBET-GIGOT, N. and GRIMA-PETTENATI, J. EgMYB2, a new transcriptional activator from *Eucalyptus* xylem, regulates secondary cell wall formation and lignin biosynthesis. *The Plant Journal*, August 2005, vol. 43, no. 4, p. 553-567.
- GRATTAPAGLIA, D. Integrating genomics into *Eucalyptus* breeding. *Genetics and Molecular Research*, September 2004, vol. 3, no. 3, p. 369-379.
- GUO, A.; HE, K.; LIU, D.; BAI, S.; GU, X.; WEI, L. and LUO, J. DATF: a database of *Arabidopsis* transcription factors. *Bioinformatics*, February 2005, vol. 21, no. 10, p. 2568-2569.
- HARAKAVA, R. Genes encoding enzymes of the lignin biosynthesis pathways in *Eucalyptus*. *Genetics and Molecular Biology*, 2004, vol. 28, no. 3, p. 601-607.
- HAUSMAN, J.F.; EVERS, D.; THIELLEMENT, H. and JOUVE, L. Compared responses of poplar cuttings and in vitro raised shoots to short-term chilling treatments. *Plant Cell Reports*, October 2000, vol. 19, no. 10, p. 954-960.
- HU, W.J.; HARDING, S.A.; LUNG, J.; POPKO, J.L.; RALPH, J.; STOKKE, D.D.; TSAI, C.J. and CHIANG, V.L. Repression of lignin biosynthesis promotes cellulose accumulation and growth in transgenic trees. *Nature Biotechnology*, August 1999, vol. 17, p. 808-812.
- HUANG, X. and MADAN, A. CAP3 A DNA sequence assembly program. *Genome Research*, September 1999, vol. 9, no. 9, p. 868-877.
- ISRAELSSON, M.; ERIKSSON, M.; HERTZBERG, M.; ASPEBORG, H.; NILSSON, H. and MORITZ, T. Changes in gene expression in the wood-forming tissue of transgenic hybrid aspen with increased secondary growth. *Plant Molecular Biology*, July 2003, vol. 52, no. 4, p. 893-903.
- JONES, R.C.; STEANE, D.A.; POTTS, B.M. and VAILLANCOURT, R.E. Microsatellite and morphological analysis of *Eucalyptus globulus* populations. *Canadian Journal of Forest Research*, January 2002, vol. 32, no. 1, p. 59-66.
- JUNG, H.-J. and NI, W. Lignification of plant cell walls. Impact of genetic manipulation. *Proceedings of the National Academy of Sciences of the United States of America*, October 1998, vol. 95, no. 22, p. 12742-12743.
- KRAUSKOPF, E.; HARRIS, P. and PUTTERILL, J. The cellulose synthase gene PrCESA10 is involved in cellulose biosynthesis in developing tracheids of the gymnosperm *Pinus radiata*. *Gene*, May 2005, vol. 350, no. 2, p. 107-116.
- LI, L.; ZHOU, Y.H.; CHENG, X.F.; SUN, J.Y.; MARITA, J.M.; RALPH, J. and CHIANG V.L. Combinatorial modification of multiple lignin traits in trees through multigene cotransformation. *Proceedings of the National Academy of Sciences of the United States of America*, April 2003, vol. 100, no. 8, p. 4939-4944.
- MIRANDA, I. and PEREIRA, H. Variation of pulpwood quality with provenances and site in *Eucalyptus globulus*. *Annals of Forest Science*, 2002, vol. 59, p. 283-291.
- NANJO, T.; FUTAMURA, N.; NISHIGUCHI, M.; IGASAKI, T.; SHINOZAKI, K. and SHINOHARA, K. Characterization of full-length enriched expressed sequence tags of stress-treated poplar leaves. *Plant and Cell Physiology*, 2004, vol. 45, no. 12, p. 1738-1748.
- OH, S.; PARK, S. and HAN, K.-H. Transcriptional regulation of secondary growth in *Arabidopsis thaliana*. *Journal of Experimental Botany*, December 2003, vol. 54, no. 393, p. 2709-2722.
- PLOMION, C.; LEPROVOST, G. and STOKES, A. Wood formation in trees. *Plant Physiology*, December 2001, vol. 127, p. 1513-1523.
- QUACKENBUSH, J.; LIANG, F.; HOLT, I.; PERTEA, G. and UPTON, J. The TIGR Gene Indices: reconstruction and representation of expressed gene sequences. *Nucleic Acids Research*, 2000, vol. 28, no. 1, p. 141-145.
- RALPH, J.; AKIYAMA, T.; KIM, H.; LU, F.; SCHATZ, P.F.; MARITA, J.M.; RALPH, S.A.; SRINIVASA-REDDY, M.S.; CHEN, F. and DIXON, R.A. Effects of coumarate 3-hydroxylase down-regulation lignin structure. *Journal of Biological Chemistry*, March 2006, vol. 281, no. 13, p. 8843-8853.
- RANIK, M. and MYBURG, A. Six new cellulose synthase genes from *Eucalyptus* are associated with primary and secondary cell wall biosynthesis. *Tree Physiology*, May 2006, vol. 26, no. 5, p. 545-556.
- RASTOGI, S. and DWIVEDI, UN. Down-regulation of lignin biosynthesis in transgenic *Leucaena leucocephala* harboring O-methyltransferase gene. *Biotechnology Progress*, 2006, vol. 22, p. 609-616.
- STERKY, F.; REGAN, S.; KARLSSON, J.; HERTZBERG, M.; ROHDE, A.; HOLMBERG, A.; AMINI, B.; BHALERAO, R.; LARSSON, M.; VILLAROEEL, R.; VAN MONTAGU, M.; SANDBERG, G.; OLSSON, O.; TEERI, T.; BOERJAN, W.; GUSTAFSSON, P.; UHLÉN, M.; SUNDBERG, B. and LUNDEBERG, J. Gene discovery in the wood forming tissues of poplar: analysis of 5,692 expressed sequence tags. *Proceedings of the National Academy of Sciences of the United States of America*, October 1998, vol. 95, no. 22, p. 13330-13335.
- STERKY, F.; BHALERAO, R.; UNNEBERG, P.; SEGERMAN, B.; NILSSON, P.; BRUNNER, A.;

CHARBONNEL-CAMPAA, L.; LINDVALL, J.; TANDRE, K.; STRAUSS, S.; SUNDBERG, B.; GUSTAFSSON, P.; UHLÉN, M.; BHALERAO, R.; NILSSON, O.; SANDBERG, G.; KARLSSON, J.; LUNDEBERG, J. and JANSSON, S. A populus EST resource for plant functional genomics. *Proceedings of the National Academy of Sciences of the United States of America*, September 2004, vol. 101, no. 38, p. 13951-13956.

TAMAGNONE, L.; MERIDA, A.; PARR, A.; MACKAY, S.; CULIANEZ-MACIA, F.; ROBERTS, K. and MARTIN, C. The AmMYB308 and AmMYB330 transcription factors from *Antirrhinum* regulate phenylpropanoid and lignin biosynthesis in transgenic tobacco. *Plant Cell*, February 1998, vol. 10, no. 2, p. 135-154.

U.S. DEPARTMENT OF ENERGY. OIT Agriculture Plants/crop-based renewable resources 2020. [cited date July 2007], 2007. Available from Internet: <http://www.energy.gov>.

ZDOBNOV, E.M. and APWEILER, R. InterProScan-an integration platform for the signature-recognition methods in InterPro. *Bioinformatics*, 2001, vol. 17, no. 9, p. 847-848.

APPENDIX SUPPLEMENTARY DATA

Supplementary data 1. Contigs with ESTs highly represented. Assigned function is indicated in contigs with more than 20 ESTs.

N° ESTs	Contigs Length (nt)	Assigned Function	Relative Organism	Similarity (%)	Accession Number (gi)
118	1604	plastidic aldolase	<i>Nicotiana paniculata</i>	94.8	4827253
86	1509	chloroplast latex aldolase-like protein	<i>Manihot esculenta</i>	90.9	56122688
81	979	ribulose-1,5-bisphosphate carboxylase/oxygenase small subunit	<i>Panax ginseng</i>	93.2	77157637
80	1693	ribulose-1,5-bisphosphate carboxylase/oxygenase activase precursor	<i>Malus x domestica</i>	93.1	415852
73	1590	glyceraldehyde-3-phosphate dehydrogenase A subunit	<i>Glycine max</i>	87.8	77540210
51	1721	AAA ATPase, central region; Homeodomain-like	<i>Medicago truncatula</i>	89.1	92875540
49	1320	glyceraldehyde-3-phosphate dehydrogenase B subunit	<i>Glycine max</i>	79.5	77540212
47	1051	Ferritin-related	<i>Medicago truncatula</i>	86.7	92870561
45	951	glycolate oxidase	<i>Mesembryanthemum crystallinum</i>	97.4	1773330
43	1085	unnamed protein product type I (CAB-21) (LHCP)	<i>Nicotiana tabacum</i>	96.7	19823
41	940	serine-glyoxylate aminotransferase	<i>Spirodela polyrrhiza</i>	84.6	74229863
39	1068	thi	<i>Citrus sinensis</i>	85.1	2582665
36	1468	putative galactinol synthase	<i>Pisum sativum</i>	83.2	541885
35	916	myo-inositol 1-phosphate synthase (MI-1-P) synthase (IPS)	<i>Sesamum indicum</i>	94.0	9858816
33	1346	catalase	<i>Soldanella alpina</i>	97.9	2462661
32	1368	oxygen evolving enhancer protein 1 precursor	<i>Bruguiera gymnorrhiza</i>	88.8	9229957
32	1261	glycine hydroxymethyltransferase methylase	<i>Solanum tuberosum</i>	95.6	438247
31	1286	galactinol synthase, isoform GolS-1	<i>Ajuga reptans</i>	83.6	5608497
31	961	BURP domain-containing protein	<i>Bruguiera gymnorrhiza</i>	59.0	14422444
28	1118	ferredoxin-NADP(+) reductase	<i>Nicotiana tabacum</i>	90.3	2225993
27	1761	S-adenosylmethionine decarboxylase	<i>Ipomoea nil</i>	82.9	1498080
27	2319	PAPS-reductase-like-protein	<i>Catharanthus roseus</i>	85.7	12831474
26	1256	Phosphoribulokinase, chloroplast precursor	<i>Mesembryanthemum crystallinum</i>	92.7	125578
25	1192	AlaT1	<i>Vitis vinifera</i>	95.4	71842524
24	975	polyubiquitin	<i>Pinus sylvestris</i>	96.1	1332579
23	843	S-adenosylmethionine decarboxylase	<i>Prunus persica</i>	58.0	47232488
23	1119	putative RNA binding protein	<i>Arabidopsis thaliana</i>	84.7	3850621
23	1320	translation elongation factor 1A-9	<i>Gossypium hirsutum</i>	95.4	74486744
23	1214	AT4g04640	<i>Arabidopsis thaliana</i>	85.4	7267222
22	1074	OJ000223_09.15	<i>Oryza sativa</i> (japonica cultivar-group)	91.1	38346061
21	1038	aminomethyltransferase system T protein	<i>Arabidopsis thaliana</i>	95.3	15221119
21	1066	NADPH-protochlorophyllide oxidoreductase	<i>Cucumis sativus</i>	90.5	2244614
21	1288	carbonic anhydrase	<i>Vigna radiata</i>	79.5	8954289
21	956	NAD-dependent epimerase/dehydratase family protein-like protein	<i>Solanum tuberosum</i>	93.7	82400136
20	1979	P-Protein-like protein precursor	<i>Arabidopsis thaliana</i>	84.5	7270248
20	934	S-adenosylmethionine decarboxylase	<i>Prunus persica</i>	86.7	47232488
20	1155	sedoheptulose-1,7-bisphosphatase precursor	<i>Oryza sativa</i> (indica cultivar-group)] cultivar-group)	91.8	27804772
20	1238	Malate dehydrogenase, glyoxysomal precursor	<i>Citrullus lanatus</i>	83.3	126894
20	1857	Transketolase, C-terminal-like	<i>Medicago truncatula</i>	88.4	92892666

^aTemperature incubation at 30°C, ^bInitial pH at 7.
μ: mean value and standard deviation of three determinations.

Supplementary data 2. Analysis of *E. globulus* unigenes corresponding to enzymes involved in wood formation.

Gene Name	Best Blastx Hit	Score	E-value	Coverage	Domains
EgPAL1	PAL <i>Daucus carota</i>	387	8e-113	29.09%	IPR001106 IPR008948 PF00221
EgPAL2	PAL <i>Camellia sinensis</i>	306	8e-82	29.13%	IPR001106 IPR008948 PF00221
EgPAL3	PAL <i>Citrus limon</i>	278	2e-73	25.76%	IPR001106 IPR008948 PF00221
EgC4H	C4H <i>Populus kitakamiensis</i>	482	3e-142	53.07%	IPR001128 IPR002401 PF00067
Eg4CL1	4CL <i>Populus balsamifera subsp. trichocarpa x Populus deltoides</i>	246	6e-64	33.16%	IPR000873 PF00501
Eg4CL2	4CL <i>Eucalyptus camaldulensis</i>	260	4e-68	36.95%	IPR000873 PF00501
EgHCT1	HCT <i>Nicotiana tabacum</i>	331	4e-89	49.66%	IPR003480 IPR02458
EgHCT2	HCT <i>Oryza sativa (japonica cultivar-group)</i>	186	1e-45	45.13%	IPR003480 IPR02458
EgC3H1	C3H-1 <i>Ocimum basilicum</i>	249	3e-72	29.49%	IPR001128 IPR002401 PF00067
EgC3H2	C3H <i>Ocimum basilicum</i>	375	1e-102	44.41%	IPR001128 IPR002401 PF00067
EgCOMT	COMT <i>Eucalyptus gunnii</i>	451	1e-125	63.66%	IPR001077 PF00891 IPR011991
EgCCoAOMT1	CCoAOMT <i>Plantago major</i>	96.7	6e-19	47.27%	IPR002935 PF01596
EgCCoAOMT2	CCoAOMT <i>Ammi majus</i>	121	2e-26	42.32%	IPR002935 PF01596
EgCCR1	CCR <i>Eucalyptus gunnii</i>	443	1e-123	63.00%	IPR001509 PF01370
EgCCR2	CCR <i>Arabidopsis thaliana</i>	256	5e-67	61.68%	IPR001509 PF01370
EgF5H1	F5H <i>Camptotheca acuminata</i>	539	1e-151	69.26%	IPR001128 IPR002401 PF00067
EgF5H2	F5H <i>Camptotheca acuminata</i>	126	8e-28	14.00%	IPR001128 IPR002401 PF00067
EgCAD1	CAD <i>Eucalyptus globulus</i>	530	1e-149	88.76%	IPR002085 IPR002328 IPR011032
EgCAD2	CAD <i>Arabidopsis thaliana</i>	220	4e-56	57.91%	IPR001509

Supplementary data 3. Analysis of *E. globulus* cellulose synthase unigenes.

Gene Name	Best Hit with <i>E. grandisi</i>	Best Blastx Hit	Score	E-value	Coverage	Domains
EgCesA1	EgrCesA5	CesA <i>Populus tremula x Populus tremuloides</i>	400	1e-110	24.11%	IPR005150 PF03552
EgCesA2	No EgrCesA related	CesA <i>Medicago truncatula</i>	238	2e-62	28.34%	IPR005150 PF03552
EgCesA3	EgrCesA4	CesA4 <i>Eucalyptus grandis</i>	309	9e-83	17.4%	IPR005150 PF03552

Supplementary data 4. Analysis of *E. globulus* unigenes corresponding to transcription factors.

Best Blast Hit	Related Organism	Score	E-value	Coverage	Domains
putative MYB transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)] cultivar-group]]	130	1e-28	30.8%	IPR001005 IPR006447 IPR009057 IPR012287
MYB transcription factor LHY-CCA1-like5	<i>Arabidopsis thaliana</i>	241	5e-62	90.1%	IPR001005 IPR006447 IPR009057 IPR012287
GPR11 (GOLDEN2-LIKE 1); transcription factor	<i>Arabidopsis thaliana</i>	239	2e-61	77.9%	IPR000183 IPR001005 IPR006447 IPR009057
MYB-like transcription factor DIVARICATA	<i>Antirrhinum majus</i>	328	2e-88	73.9%	IPR001005 IPR006447 IPR009057 IPR012287
MYB-like transcription factor 6	<i>Gossypium raimondii</i>	288	2e-76	103.1%	IPR001005 IPR009057 IPR012287
MYB11	<i>Malus x domestica</i>	230	4e-59	50.3%	IPR001005 IPR009057 IPR012287
MYBR2	<i>Malus x domestica</i>	75.9	2e-12	41.6%	No related InterPro
MYB-like DNA-binding protein	<i>Catharanthus roseus</i>	95.5	2e-18	24.4%	No related InterPro
transcription factor MYB1	<i>Malus xiaojinensis</i>	88.6	2e-16	47.0%	No related InterPro
MYB, DNA-binding	<i>Medicago truncatula</i>	78.6	2e-13	28.3%	No related InterPro
MYB transcription factor LHY-CCA1-like5	<i>Arabidopsis thaliana</i>	147	5e-34	35.5%	IPR001005 IPR006447 IPR009057 IPR012287
MYB-related protein	<i>Arabidopsis thaliana</i>	218	2e-55	48.8%	IPR001005 IPR009057 IPR012287
IAA18; transcription factor	<i>Arabidopsis thaliana</i>	55.5	2e-06	33.3%	No related InterPro
Transcriptional factor B3	<i>Medicago truncatula</i>	102	2e-20	32.7%	IPR003340
Transcriptional factor B3	<i>Medicago truncatula</i>	188	2e-46	22.7%	No related InterPro
bZIP transcription factor protein	<i>Capsicum annuum</i>	75.1	3e-12	45.5%	IPR004827 IPR008917 IPR011616
Putative ripening-related bZIP protein	<i>Vitis vinifera</i>	171	3e-41	38.7%	No related InterPro
bZIP transcription factor ZIP-2	<i>Nicotiana tabacum</i>	80.5	7e-14	32.9%	IPR004827 IPR008917
Putative ripening-related bZIP protein	<i>Vitis vinifera</i>	123	7e-27	37.4%	No related InterPro
ATBZIP60	<i>Arabidopsis thaliana</i>	82	2e-14	61.7%	IPR004827 IPR008917 IPR011616
Putative ripening-related bZIP protein	<i>Vitis vinifera</i>	176	9e-43	45.2%	No related InterPro
Putative ripening-related bZIP protein	<i>Vitis vinifera</i>	84.3	6e-15	16.3%	No related InterPro
GRAS transcription factor	<i>Medicago truncatula</i>	67.8	5e-10	10.2%	No related InterPro
GRAS transcription factor	<i>Medicago truncatula</i>	252	9e-66	33.4%	IPR005202
ATHB-7	<i>Arabidopsis thaliana</i>	142	2e-32	33.3%	IPR000047 IPR001356 IPR003106 IPR009057 IPR012287
Heat shock transcription factor	<i>Phaseolus acutifolius</i>	272	1e-71	60.7%	IPR000232 IPR002341 IPR011991
WRKY9; transcription factor	<i>Arabidopsis thaliana</i>	127	1e-27	75.4%	IPR003657
Putative WRKY-type DNA binding protein	<i>Glycine max</i>	380	1e-104	50.5%	IPR003657 IPR000583
Putative WRKY4 transcription factor	<i>Vitis aestivalis</i>	74.3	4e-12	34.4%	No related InterPro
DNA-binding WRKY	<i>Medicago truncatula</i>	168	2e-40	38.3%	IPR003657

Generation and analysis of an *Eucalyptus globulus* cDNA

Putative WRKY4 transcription factor	<i>Vitis aestivalis</i>	240	5e-62	66.6%	IPR003657
Putative zinc finger transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	242	1e-62	35.1%	IPR000571
Putative zinc finger transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	243	4e-63	35.1%	IPR000571
Zinc finger protein	<i>Malus x domestica</i>	134	3e-30	26.1%	IPR007087
Zinc finger protein, putative	<i>Arabidopsis thaliana</i>	137	1e-30	52.4%	IPR000315
putative zinc finger protein	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	133	9e-30	73.5%	IPR007087
Putative zinc finger transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	225	2e-57	31.3%	IPR000571
Zinc finger protein-like	<i>Arabidopsis thaliana</i>	123	6e-27	5.4%	IPR003349
Zinc finger protein, putative	<i>Plasmodium falciparum</i> 3D7	60.1	8e-08	28.8%	IPR002653
Putative zinc finger transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	239	9e-62	34.8%	IPR000571
Zinc finger protein, putative	<i>Arabidopsis thaliana</i>	109	2e-22	21.6%	IPR000315 IPR002906
Zinc finger protein, putative	<i>Arabidopsis thaliana</i>	152	1e-35	36.3%	IPR000315
Putative zinc finger transcription factor	<i>Oryza sativa</i> (japonica cultivar-group)cultivar-group	218	2e-55	28.4%	IPR000571