

## A comprehensive study of the sugar pine (*Pinus lambertiana*) transcriptome implemented through diverse next-generation sequencing approaches

Pedro J. Martínez-García, Daniel González-Ibeas, Randi A. Famula, Annette Delfino-Mix, Kristian A. Stevens, Jeffrey D. Puryear, Carol A. Loopstra, Charles H. Langley, David B. Neale, Jill L. Wegrzyn



PAG XXIV Forest Tree Workshop January-10 San Diego

## Conifer Reference Assemblies

**Loblolly pine (*P. taeda*)**  
 Assembly LP\_v1.01  
 ~ 65X coverage  
 Total Sequence: 23.2 Gbp  
 N50 Scaffold: 66.9Kbp

**Sugar pine (*P. lambertiana*)**  
 (just released)  
 Assembly SP\_v1.0  
 ~ 53X (PE)  
 ~ 17.3X (MP)  
 0.5-1X (DiTags)  
 Total Sequence: 31Gbp  
 N50 Scaffold: 246.5Kbp

**Douglas fir (*Pseudotsuga menziesii*)**  
 Assembly Psme\_v1.0  
 (just released)  
 ~ 60X (PE)  
 11X (MP)  
 Total Sequence: 15.7Gbp  
 N50 Scaffold: 340.7Kbp

PAG XXIV Forest Tree Workshop January-10 San Diego

## Previous Genomic resources for white pines

Species	Technology	Reads	Tissue	Reads after QC
Sugar pine Jessica Wright	Illumina GA Ix	SE, 80bp (3 lanes)	needle	66,894,169
Sugar pine (Lorenz et al. 2012)	Roche 454	SE, 350 bp (avg)	stem needle	952,310
Limber pine Jeff Mitton	Illumina HiSeq	PE, 100bp (2 lanes)	needle	374,191,816
Whitebark pine Patricia Maloney	Illumina HiSeq	PE, 100bp (3 lanes)	needle	839,389,034
Western white pine (J-J. Liu et al 2013)	Illumina GA Ix	PE, 76bp	needle	208,059,003



PAG XXIV Forest Tree Workshop January-10 San Diego

## Generating new resources for sugar pine



PAG XXIV Forest Tree Workshop January-10 San Diego

## Tissues collection



PAG XXIV Forest Tree Workshop January-10 San Diego

## Stress Treatments

Needles

Stems

Roots

- ☼ Drought Stress
- ☼ Cold Stress
- ☼ Heat Shock
- ☼ Flooding
- ☼ Salt
- ☼ Wounding
- ☼ Jasmonic Acid
- ☼ Blister rust disease

PAG XXIV Forest Tree Workshop January-10 San Diego

**Difficulties Resolving Transcriptomes with Short Reads**

**Assessment of transcript reconstruction methods for RNA-seq**

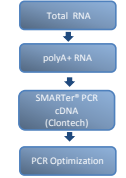
Tamara Stojicevic<sup>1</sup>, Joseph F. Ahehi<sup>1,2,3</sup>, Pir G. Engstrom<sup>1,2,3</sup>, Felix Kolkovskiy<sup>1,2,3</sup>, The BGASP Consortium<sup>1</sup>, Tim J. Hubbard<sup>4</sup>, Roderic Guigou<sup>5,6</sup>, Jennifer Harrow<sup>7</sup> & Paul Bertone<sup>1,2,3,8</sup>

We evaluated 23 protocol variants of 14 independent computational methods for transcript reconstruction and compared their performance to that of a single reference-based method. Expression analysis methods...  
 ...the complexity of higher eukaryotic genomes imposes severe limitations on transcript recall and splice product discrimination...  
 ...assembly of complete isoform structures poses a major challenge even when all constituent elements are identified...  
 ...Ultimately, the evolution of RNA-seq will move toward single-pass determination of intact transcripts....

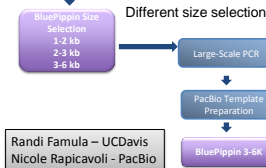
PAG XXIV Forest Tree Workshop January-10 San Diego

**Full-length reads with Iso-Seq**

**Library preparation**



**Bioinformatics workflow**



Clustering step ICE/Quiver:  
 Consensus isoforms (Pb1)  
 Low quality polished sequences (Pb2)  
 High quality polished sequences (Pb3)

Randi Famula - UC Davis  
 Nicole Raponi - PacBio

PAG XXIV Forest Tree Workshop January-10 San Diego

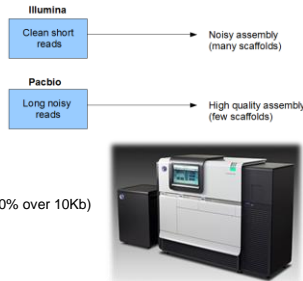
**Transcriptome Sequencing Strategy**

**Hybrid Approach to Sequencing:**

**Hi-Seq**  
 Average length=(100x2)  
 180 million reads/lane  
**Accuracy: 99%**

**Mi-Seq**  
 Average length=(300x2)  
 25 million reads/lane  
**Accuracy: 99.6%**

**PacBio SMRT II Iso-Seq**  
 Size selected lengths (5-6Kb, 10% over 10Kb)  
 40,000 reads/SMRT cell (run)  
**Accuracy: 86%**



PAG XXIV Forest Tree Workshop January-10 San Diego

**Illumina**

Description	Tissue	Technology	Number of reads after QC	Total trinity genes	Total trinity transcripts	iso/lane
Blister Resistant needles (LCO2-03)	needle	HiSeq	264207835	73681	112529	1
Seedling was not drought-stressed	Root	HiSeq	232870127	100792	150944	1
Seedling was not drought-stressed	Stem	HiSeq	272808440	64239	104148	1
Germinating sugar pine seed	Embryo	HiSeq	232340072	50340	79393	1
Methyl jasmonate treatment for 5 hrs	Stem	HiSeq	255543502	58239	92768	1
NaCl treatment	Root	HiSeq	237568138	68861	101178	1
Female cones (2cm)	Female cones	HiSeq	241875834	51363	79350	1
Female strobili near pollination	Female strobili	HiSeq	271361467	71438	110173	1
Wounding treatment	Stem	HiSeq	268217058	71226	110544	1
Blister Resistant stem (LCO2-03)	Stem	HiSeq	33374491	64531	94955	1
Seedling slowly drought-stressed	Needle	HiSeq	32947205	69877	95471	1
Pollen	Pollen	HiSeq	33131689	69359	102314	1
"Basket stage" seedling	Root, stem and needles	HiSeq	33725577	52998	80065	1
Germinating sugar pine seed	Embryo	HiSeq	28838541	48942	73271	1
Early female conelets before pollination	female conelets	HiSeq	37166765	115809	155024	1
Pollen cones	Pollen cones	HiSeq	27597667	71414	102389	1
		Total	2503693428	1103099	1644816	16

**2.5 billion reads**

PAG XXIV Forest Tree Workshop January-10 San Diego

**PacBio**

Description	Tissue	size selection	Number SMRT cells	Reads of insert after QC	Number of full-length non-chimeric reads	Number of consensus isoforms	Number of polished low-quality isoforms	Number of polished high-quality isoforms
Seedling was not drought-stressed	stem	1 Kb	4	201744	102092	77603	2564	6574
Seedling was not drought-stressed	stem	2 Kb	4	259668	103022	81386	2398	5646
Germinating sugar pine seed	Embryo	1 Kb	1	60846	14343	10890	497	835
Germinating sugar pine seed	Embryo	2 Kb	4	259631	125057	88710	2881	9854
Germinating sugar pine seed	Embryo	3-6 Kb	1	44815	19413	14100	1090	539
2 cm female cones	female cones	1 Kb	3	162898	61611	46531	2062	3790
2 cm female cones	female cones	2 Kb	4	224284	84540	64185	2563	5487
2 cm female cones	female cones	3-6 Kb	3	249945	85231	60218	41383	18834
Female strobili near pollination	female strobili	1 Kb	4	228401	111972	67126	3158	9830
Female strobili near pollination	female strobili	2 Kb	4	199885	83411	55594	3415	5117
Female strobili near pollination	female strobili	3-6 Kb	3	238685	82948	53428	34317	19111

**1.6 million reads**

PAG XXIV Forest Tree Workshop January-10 San Diego

**Eukaryote Non-Model Transcriptome Annotation Pipeline**

QC-sickle -> Trinity\_denovo\_assembly -> Transdecoder -> USEARCH/UCLUST + enTAP



(enTAP, <https://github.com/SamGinzburg/WegrzynLab>)

PAG XXIV Forest Tree Workshop January-10 San Diego

**Transcriptome statistics**

**Assembled transcripts (number of sequences)**

Total transcripts	278812
HiSeq	75175
MiSeq	45524
PacBio	158113

**Set of unique transcripts/full length transcripts**

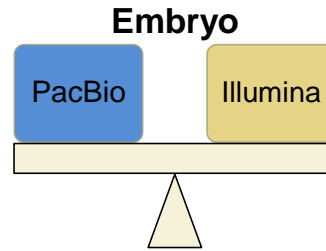
Number (scaffolding)	33113
Average length	1144
Shortest transcript	300
Largest transcript	13236
N50 Statistic	1386

**Functional annotation**

Annotated	30839
Informative	26568
Uninformative	3923
Unannotated	1243
Contaminants	1399

PAG XXIV Forest Tree Workshop January-10 San Diego

**Comparison of Sequencing Technologies**



**TRANSCRIPTS**  
Length - Completeness - Mapping rates

**TRANSCRIPTOME**  
Coverage - Diversity

PAG XXIV Forest Tree Workshop January-10 San Diego

**Comparison – Transcripts**

- PacBio highest number of complete coding regions (8940 PacBio > 8782 HiSeq > 7892 MiSeq)
- PacBio longer lengths before TS, BUT after TS lengths were similar between technologies
- Less 4% of transcripts were full-length (70%-70%) with either technology before TS, with a increment up to 21% of full-length after TS
- Less than 50% of PacBio transcripts mapped (Illumina >70%) (after TS: PacBio ~60%, Illumina ~90%)

PAG XXIV Forest Tree Workshop January-10 San Diego

**Comparison - Transcriptome – Coverage and Diversity**


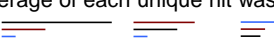
**Embryo transcriptome: 17505 unique embryo transcripts mapped against SP\_v1.0 (100%coverage/100%identity)**

- 9249 transcripts map uniquely
- 74% covered by PacBio
- 1615 unique hits covered by all three technologies

PAG XXIV Forest Tree Workshop January-10 San Diego

**Comparison - Transcriptome – Coverage and Diversity**

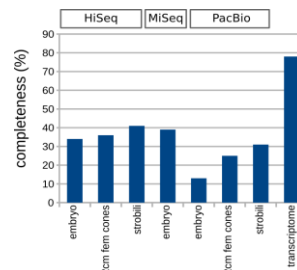
**1615 covered by all three technologies**

- Longest splicing variant (overlap) was provided by HiSeq 
- A single technology producing the longest splice variants (in number) was HiSeq
- Contribution to coverage of each unique hit was better by HiSeq 
- Highest number of non-redundant splice variants was provided by PacBio

PAG XXIV Forest Tree Workshop January-10 San Diego

**Comparison - Transcriptome - Completeness**

- Lower coverage and higher variation by PacBio



Benchmarking Universal Single-Copy Orthologs **BUSCO** (Simão et al 2015)

PAG XXIV Forest Tree Workshop January-10 San Diego

**Take-home message...**

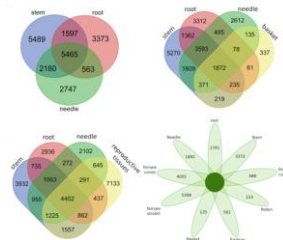
**PacBio: better in splice variant detection and highly productive for full-length (not necessarily the longest) final transcripts**

**Illumina: better in term of transcripts coverage/length, longest splice variant and high depth for expression studies**

**Differential Expression – without replicates**

**Gfold**  
(Feng et al. 2012)  
**-without replicates-**

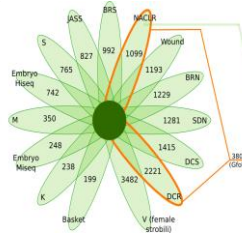
- Stem tissue more similar to needle tissue
- Low number in Basket
- Reproductive tissues more similar to stem
- High number provided by female reproductive tissue



PAG XXIV Forest Tree Workshop January-10 San Diego

**Differential Expression – without replicates**

- More than 10000 transcripts shared by all libraries
- Overall 5958 transcripts were differential expressed with fold change >2.0
- NACL-roots 1099 transcripts library specific, 233 differentially expressive



GO-ID	Term	Category	FDR
NACL-roots 1099 transcripts library specific, 233 differentially expressive			
GO:0046034	ATP metabolic process	P	6.53E-002
GO:0024689	protein-transporting two-sector ATPase complex	C	7.34E-002
NACL vs. untreated control differentially expressed transcripts (CK48)			
GO:0042555	MCM complex	C	5E-02
GO:0043168	anion binding	F	8E-08
GO:005524	ATP binding	F	2E-09
GO:0005524	ATPase activity	F	2E-03
GO:0005524	ATPase activity	F	5E-02
GO:0010817	regulation of hormone levels	P	4E-08
GO:0048765	root hair cell differentiation	P	2E-02
GO:0048767	root hair elongation	P	7E-03
GO:0009809	lignin biosynthetic process	P	7E-03

PAG XXIV Forest Tree Workshop January-10 San Diego

**Differential Expression – GO categories**

GO-ID	Term	Category	FDR
NACL treatment			
GO:0042555	MCM complex - mostly helicases	C	0.049359
GO:0043168	anion binding	F	7.84E-18
GO:0005524	ATP binding	F	2.14E-09
GO:0016887	ATPase activity	F	0.004976
GO:0005034	osmosensor activity - three different histidine kinases	F	0.046233
GO:0010817	regulation of hormone levels	P	6.04E-06
GO:0048765	root hair cell differentiation	P	0.018912
GO:0048767	root hair elongation	P	0.006801
GO:0009809	lignin biosynthetic process	P	0.023398
JASS treatment			
GO:0010533	response to cyclopentane - most topoisomerases	P	0.007481
GO:0009805	response to external stimulus	P	0.018085
GO:0043207	response to external biotic stimulus	P	0.021576
GO:0051707	response to other organism	P	0.021576
GO:0016458	gene silencing	P	0.023677
GO:0051567	histone H3-K9 methylation	P	7.68E-09
GO:0042742	defense response to bacterium	P	0.028652
GO:0010476	gibberellin mediated signaling pathway	P	0.011213
GO:0042221	response to chemical	P	1.90E-07

PAG XXIV Forest Tree Workshop January-10 San Diego

**Differential Expression – GO categories**

GO-ID	Term	Category	FDR
Wounding			
GO:0006950	response to stress	P	8.39E-09
GO:0006952	defense response	P	2.54E-06
GO:0006911	cell-cell junction	C	0.000132
GO:0030855	epithelial cell differentiation	P	0.010485
GO:0009913	epidermal cell differentiation	P	0.010485
GO:0008544	epidermis development	P	0.013693
GO:0060429	epithelium development	P	0.017064
GO:0042545	cell wall modification	P	0.035399
Reproductive tissue			
GO:0007169	transmembrane receptor protein tyrosine kinase signaling pathway oxidoreductase activity, acting on diphenols and related substances	P	2.86E-07
GO:0016882	as donors, oxygen as acceptor	F	0.002152
GO:0009751	response to salicylic acid	P	0.002236
GO:0010333	terpene synthase activity	F	6.87E-14
GO:0009740	gibberellic acid mediated signaling pathway	P	0.000720
GO:0048506	regulation of timing of meristematic phase transition	P	0.000373
GO:0007389	pattern specification process	P	0.002009
GO:0009955	adaxial/abaxial pattern specification	P	0.008520
GO:0007165	signal transduction	P	1.66E-17
GO:0050793	regulation of developmental process	P	2.67E-05
GO:0010476	gibberellin mediated signaling pathway	P	0.000174
GO:0009686	gibberellin biosynthetic process	P	0.000641

PAG XXIV Forest Tree Workshop January-10 San Diego

**Genome expansion in conifers**

- Expansion consequence of transposable element (TE) proliferation (80%) rather than genome duplications
- Unique small RNA profile (24-nt) in conifers associated with epigenetic processes and control of repetitive element proliferation
- Potential lineage-specific Dicer-like (DCL) (key in sRNA biogenesis) proteins were identified (Dolgosheina et al. 2008) in conifers.
- 12 transcripts in reproductive tissue with sequence similarity and domain topology matching DCL features expanding their characterization in sugar pine.
- 6 were supported by gene models (genome v1.0).

PAG XXIV Forest Tree Workshop January-10 San Diego

## Additional 1.7 billion Illumina reads

Samples	Technology	Reads of Insert after QC	Total trinity transcripts	Total trinity genes	line
Female strobili	MSSeq	32260898	218646	180581	89203
Female cones	MSSeq	31376822	128302	103083	59841
Female conelets (PM)	HiSeq	27460906	61700	44961	36784
"Basket stage" seedling	HiSeq	26048750	56835	40070	34728
pollen (PM)	HiSeq	29461236	56963	42176	32970
pollen cones (PM)	HiSeq	25211652	67374	49795	8081
Cold Root (Tree 1)	HiSeq	24537908	77775	57633	6942
Primary needle stage (is this tree #1)	HiSeq	24592668	61195	43902	37262
Susceptible WPBR (Stem) #2	HiSeq	99080770	112587	73700	28605
Susceptible WPBR (Root) #2	HiSeq	102280544	128879	88061	22143
Susceptible WPBR (Needle) #2	HiSeq	114994106	146059	97675	28718
Susceptible WPBR (Stem) #1	HiSeq	76938486	105459	69517	28291
Susceptible WPBR (Root) #1	HiSeq	217641798	199948	137558	26214
Susceptible WPBR (Needle) #1	HiSeq	142041716	177094	116622	31923
Resistance WPBR (stem) #1	HiSeq	92196090	111708	78176	16551
Resistance WPBR (root) #1	HiSeq	16724006	195158	138647	27903
Resistance WPBR (needle) #1	HiSeq	109706506	136890	96206	19621
Resistance WPBR (stem) #2	HiSeq	109114556	124300	85202	23869
Resistance WPBR (root) #2	HiSeq	100142516	134295	94199	21987
Resistance WPBR (needle) #2	HiSeq	157203340	147026	99169	23038
<b>Total</b>		<b>1726543178</b>	<b>2443193</b>	<b>1736993</b>	<b>607674</b>

PAG XXIV Forest Tree Workshop January-10 San Diego

## Acknowledgements

### University of Connecticut

- Jill Wegrzyn
- Daniel Gonzalez-Ibeas
- Ethan Baker
- Robin Paul

### University of California, Davis

- Randi Famula
- Hans Vasquez-Gross
- David Neale
- Kristian Stevens
- Charles Langley

### Texas A&M University

- Carol Loopstra
- Jeff Puryear

### USDA Forest Service

- Detlev Volger
- Annette Delfino-Mix

### Indiana University

- Keithanne Mockaitis

