

# Gramene GeneTrees:

## A comprehensive phylogenomics database in plants and other model Eukaryotes

William Spooner<sup>1</sup>, Joshua C. Stein<sup>1</sup>, Sharon Wei<sup>1</sup>, Doreen Ware<sup>1,2</sup>

<sup>1</sup>Cold Spring Harbor Laboratory, Cold Spring Harbor, NY 11724

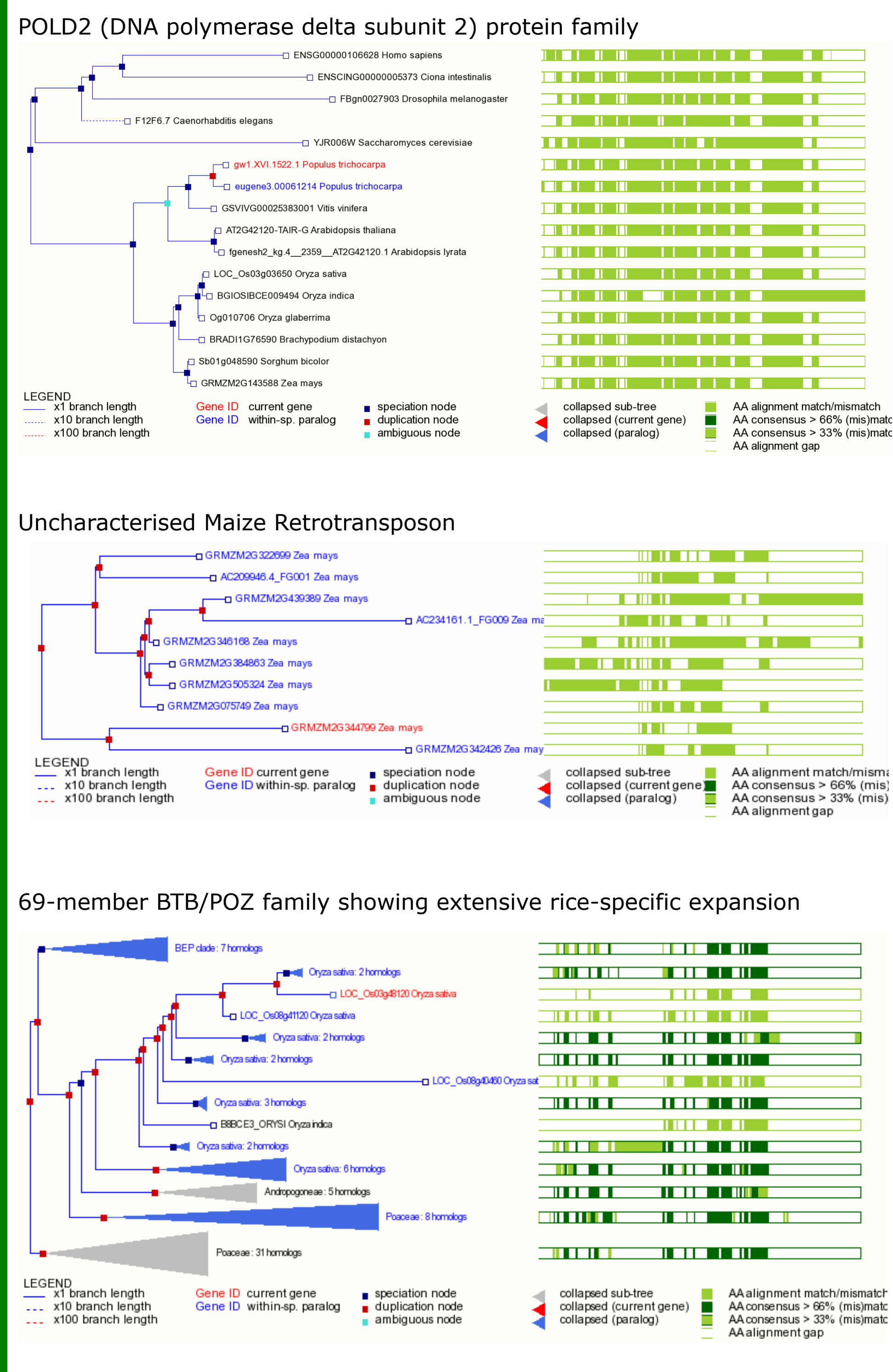
<sup>2</sup>USDA-ARS NAA Plant, Soil & Nutrition Laboratory Research Unit, Ithaca, NY 14853

## Abstract

Since the completion of the *Arabidopsis thaliana* genome in 2000, more than 20 plant genomes have been sequenced, with the number set to increase rapidly in the coming years. It is now becoming possible to apply phylogenomics-based techniques from metazoan biology to the plant arena. We have applied an automated methodology, EnsemblCompara GeneTrees, to proteins predicted from a wide range of genomes to develop the first comprehensive plant phylogenomics resource. This consists of protein-level phylogenetic trees between twelve whole genomes; four dicotyledon plants (*Arabidopsis lyrata*, *Arabidopsis thaliana*, *Populus trichocarpa*, *Vitis vinifera*), three monocotyledon plants (*Oryza sativa* Japonica Group, *Oryza sativa Indica* Group, *Sorghum bicolor*), and five model metazoa/fungi (*Caenorhabditis elegans*, *Ciona intestinalis*, *Drosophila melanogaster*, *Homo sapiens*, *Saccharomyces cerevisiae*). Validation of our data through comparison data with similar resources and consistent results both for summaries of the database as a whole and for individual example trees. The GeneTrees form a component of the Gramene database (<http://www.gramene.org>), an established resource for comparative plant genomics, and form a useful platform for the study of plant molecular evolution as well as functional annotation of newly sequenced genomes.

## Visualisation

### Example GeneTrees



## Data

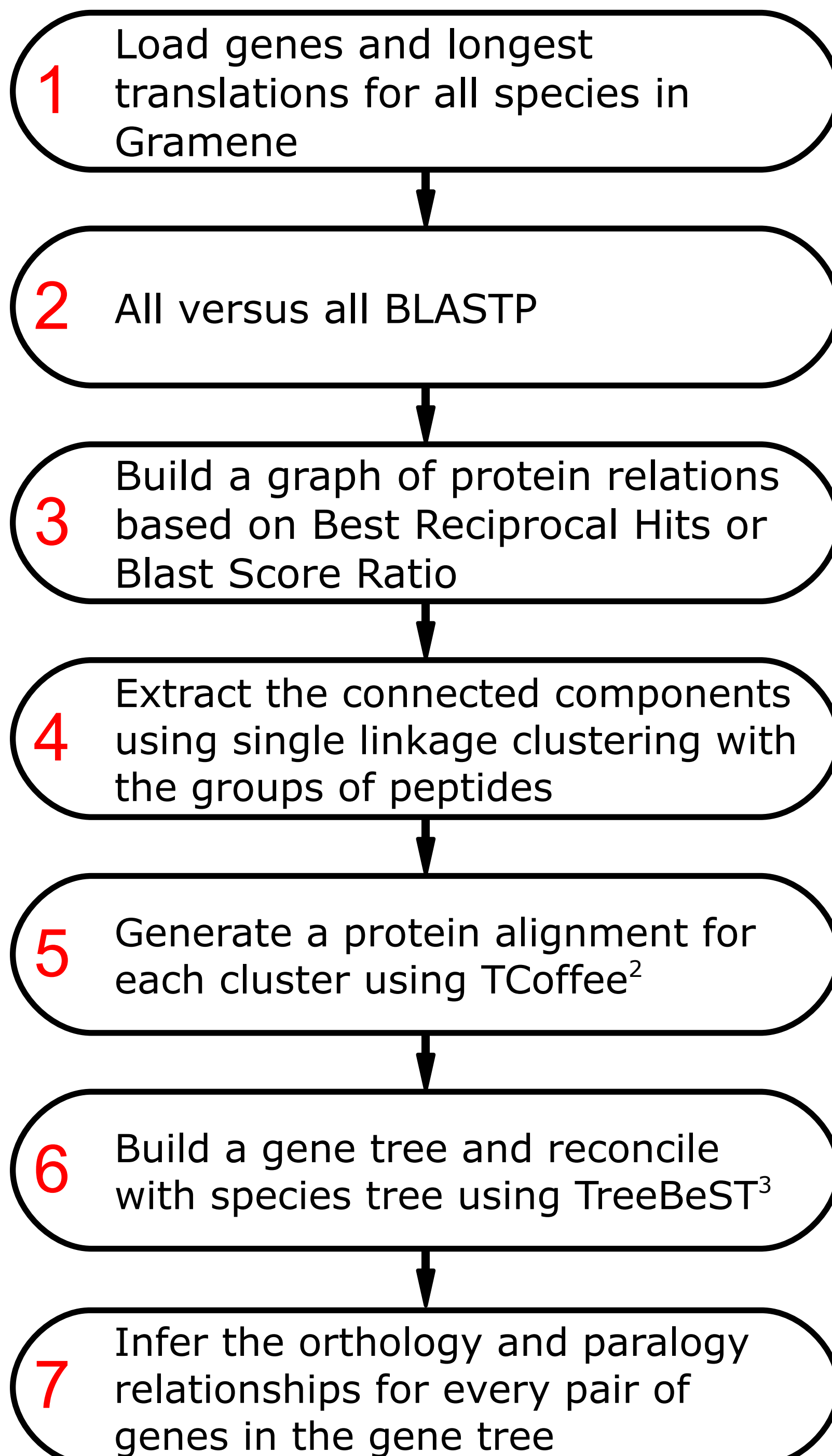
Protein-coding gene predictions from eight fully-sequenced plant species were supplemented with similar data from five fungi/metazoan species. Predictions from a second *Oryza sativa* cultivar (Indica Group) and an incomplete assembly from *Oryza glaberrima* added to a total species count of fifteen, and a total gene count of 402,192. The EnsemblCompara GeneTree method produced a total of 27,031 individual GeneTrees comprising of a total of 687,302 nodes; 357,182 leaves representing the same number of genes (88.8% of input genes), 180,383 nodes representing gene speciation events, and 149,741 nodes representing gene duplication events.

## Species

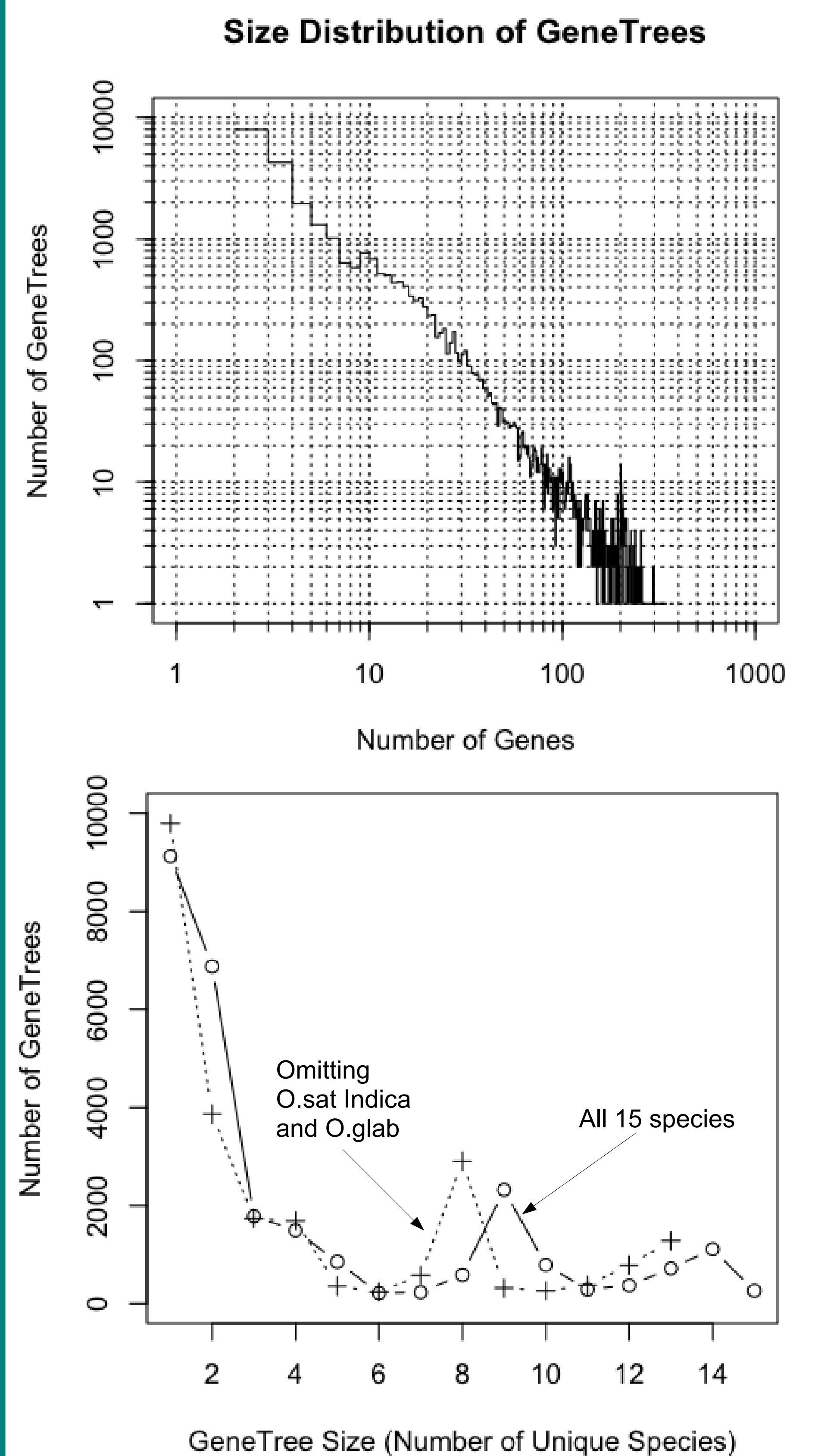
Clade	Species	Total Genes	Genes in trees	Genes in single-species tree	Singleton Genes
Monocotyledon Plants	<i>Oryza sativa</i> (Japonica rice)	57,995	51,359 (90%)	11,254 (19%)	6,636 (11%)
	Non TE-genes	41,775	35,793 (86%)	1,969 (4%)	5,982 (14%)
	Annotated TE-genes	16,220	15,566 (96%)	9,285 (57%)	654 (4%)
	<i>Oryza sativa</i> (Indica rice)	38,861	35,071 (89%)	511 (1%)	3,790 (10%)
	<i>Oryza glaberrima</i> (African rice)	2,467	2,120 (86%)	81 (3%)	347 (14%)
	<i>Brachypodium distachyon</i> (false brome)	25,532	24,564 (96%)	384 (2%)	968 (4%)
	<i>Sorghum bicolor</i> (sorghum)	34,496	32,730 (95%)	1,584 (5%)	1,766 (5%)
	<i>Zea mays</i> (maize)	32,540	30,258 (93%)	1,025 (3%)	2,282 (7%)
	<i>Arabidopsis thaliana</i> (thale cress)	31,280	29,550 (94%)	2,468 (8%)	1,730 (6%)
	Non TE-genes	27,379	26,030 (95%)	244 (1%)	1,349 (5%)
Eudicotyledon Plants	Annotated TE-genes	3,901	3,520 (90%)	2,244 (63%)	381 (10%)
	<i>Arabidopsis lyrata</i> (lyrate rockcress)	32,667	30,136 (92%)	1,839 (6%)	2,531 (8%)
	<i>Populus trichocarpa</i> (poplar)	38,449	33,903 (88%)	3,483 (9%)	4,546 (12%)
	<i>Vitis vinifera</i> (grape)	30,434	26,794 (88%)	1,979 (7%)	3,640 (12%)
	<i>Homo sapiens</i> (human)	22,294	18,774 (84%)	3,610 (16%)	3,520 (16%)
	<i>Ciona intestinalis</i> (sea squirt)	14,180	11,182 (79%)	1,930 (14%)	2,998 (21%)
Fungi/Metazoa	<i>Drosophila melanogaster</i> (fruit fly)	14,141	11,223 (79%)	2,084 (15%)	2,918 (21%)
	<i>Caenorhabditis elegans</i> (nematode)	20,158	15,415 (76%)	6,775 (34%)	4,743 (24%)
	<i>Saccharomyces cerevisiae</i> (yeast)	6,698	4,103 (61%)	989 (15%)	2,595 (39%)

## Data Generation

### Ensembl Compara Gene Tree Pipeline<sup>1</sup>



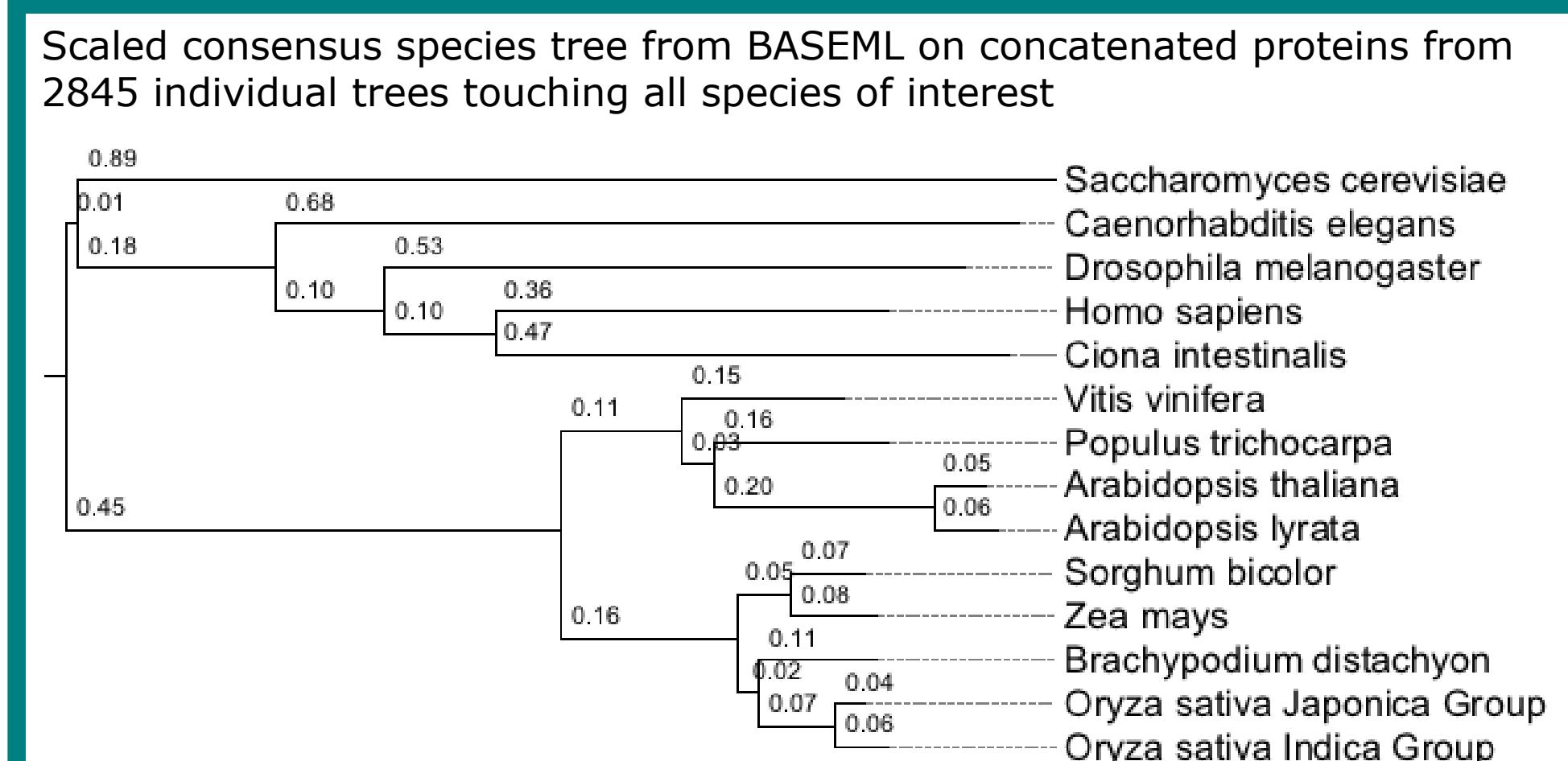
## Tree Sizes



## Phylogenetic Context

Phylogeny	Taxon	Trees with taxon	Nodes at taxon	Average nodes per Tree	% Duplication Nodes	Species Intersection Score
Oryza sativa Japonica	<i>Oryza sativa</i> Japonica	14,310	66,666	4.66	22.96	1.00
	<i>Oryza sativa</i>	11,100	36,273	3.27	18.56	0.84
	<i>Oryza sativa</i> Indica	11,735	37,572	3.20	6.66	1.00
Oryza sativa Indica	<i>Oryza sativa</i> Indica	11,735	37,572	3.20	6.66	1.00
	<i>Oryza sativa</i>	11,100	36,273	3.27	18.56	0.84
	<i>Oryza sativa</i> Japonica	14,310	66,666	4.66	22.96	1.00
Oryza glaberrima	<i>Oryza glaberrima</i>	1,302	1,997	1.53	17.68	0.55
	<i>Oryza sativa</i>	11,100	36,273	3.27	18.56	0.84
	<i>Oryza sativa</i> Japonica	14,310	66,666	4.66	22.96	1.00
BEP clade	BEP clade	7,548	19,217	2.55	3.44	0.62
	<i>Brachypodium distachyon</i>	7,953	28,195	3.55	12.88	1.00
	<i>Populus trichocarpa</i>	7,505	23,719	3.16	39.68	0.63
Brachypodium distachyon	<i>Brachypodium distachyon</i>	7,953	28,195	3.55	12.88	1.00
	<i>Populus trichocarpa</i>	7,505	23,719	3.16	39.68	0.63
	BEP clade	7,548	19,217	2.55	3.44	0.62
Sorghum bicolor	<i>Sorghum bicolor</i>	9,026	39,526	4.38	17.19	1.00
	<i>Oryza sativa</i>	11,100	36,273	3.27	18.56	0.84
	<i>Oryza sativa</i> Japonica	14,310	66,666	4.66	22.96	1.00
Zea mays	<i>Zea mays</i>	8,281	38,127	4.60	20.64	1.00
	<i>Oryza sativa</i>	11,100	36,273	3.27	18.56	0.84
	<i>Oryza sativa</i> Japonica	14,310	66,666	4.66	22.96	1.00
Magnoliophyta	Magnoliophyta	7,505	23,719	3.16	39.68	0.63
	<i>Arabidopsis thaliana</i>	8,793	33,971	3.86	13.01	1.00
	<i>Arabidopsis lyrata</i>	8,849	35,074	3.96	14.08	1.00
Arabidopsis thaliana	<i>Arabidopsis thaliana</i>	8,793	33,971	3.86	13.01	1.00
	<i>Arabidopsis lyrata</i>	8,849	35,074	3.96	14.08	1.00
	Magnoliophyta	7,505	23,719	3.16	39.68	0.63
Arabidopsis lyrata	<i>Arabidopsis lyrata</i>	8,849	35,074	3.96	14.08	1.00
	<i>Arabidopsis thaliana</i>	8,793	33,971	3.86	13.01	1.00
	Magnoliophyta	7,505	23,719	3.16	39.68	0.63
Rosids	Rosids	6,177	11,765	1.90	4.56	0.59
	<i>Populus trichocarpa</i>	8,552	49,174	5.75	31.06	1.00
	Core eudicotyledons	7,056	22,066	3.13	28.23	0.50
Populus trichocarpa	<i>Populus trichocarpa</i>	8,552	49,174	5.75	31.06	1.00
	Core eudicotyledons	7,056	22,066	3.13	28.23	0.50
	Rosids	6,177	11,765	1.90	4.56	0.59
Core eudicotyledons	Core eudicotyledons	7,056	22,066	3.13	28.23	0.50
	Rosids	6,177	11,765	1.90	4.56	0.59
	<i>Populus trichocarpa</i>	8,552	49,174	5.75	31.06	1.00
Vitis vinifera	<i>Vitis vinifera</i>	7,749	35,496	4.58	25.52	1.00
	Eukaryota	3,358	5,765	1.72	25.92	0.49
	<i>Homo sapiens</i>	6,393	27,981	4.38	32.90	1.00
Eukaryota	Eukaryota	3,358	5,765	1.72	25.92	0.49
	<i>Homo sapiens</i>	6,393	27,981	4.38	32.90	1.00
	Chordata	4,142	6,349	1.53	10.11	0.71
Homo sapiens	<i>Homo sapiens</i>	6,393	27,981	4.38	32.90	1.00
	Chordata	4,142	6,349	1.53	10.11	0.71
	Ciona intestinalis	5,068	15,040	2.97	25.65	1.00
Ciona intestinalis	<i>Ciona intestinalis</i>	5,068	15,040	2.97	25.65	1.00
	Coleomata	4,371	6,815	1.56	13.21	0.57
Coleomata	Coleomata	4,371	6,815	1.56	13.21	0.57
	<i>Drosophila melanogaster</i>	5,189	15,063	2.90	25.49	1.00
Drosophila melanogaster	<i>Drosophila melanogaster</i>	5,189	15,063	2.90	25.49	1.00
	Bilateria	3,882	7,725	1.99	29.45	0.56
Bilateria	Bilateria	3,882	7,725	1.99	29.45	0.56
	<i>Caenorhabditis elegans</i>	5,236	23,668	4.52	34.87	1.00
Caenorhabditis elegans	<i>Caenorhabditis elegans</i>	5,236	23,668	4.52	34.87	1.00
	Fungi/Metazoa	2,158	2,398	1.11	9.55	0.41
Fungi/Metazoa	Fungi/Metazoa	2,158	2,398	1.11	9.55	0.41
	<i>Saccharomyces cerevisiae</i>	2,838	5,344	1.88	23.22	1.00

## Consensus Tree



## Funding

This work was initially supported (2001-2004) by the USDA Initiative for Future Agriculture and Food Systems (IFAFS) (grant no. 00-52100-9622) and a Cooperative State Research and Education Service (CSREES) agreement through the [USDA Agricultural Research Service](#) (grant no. 58-1907-0-041). For the years 2004-2007 this work was supported by the [National Science Foundation](#) (NSF) PGI grant award #0321685. Current work is being supported by the NSF Plant Genome Research Resource grant award #0703908.

## References

1. Vilella A.J., et al. (2008). *Genome Res.* Pre-print: doi:10.1101/gr.073585.107
2. Edgar, R.C. (2004). *Nucleic Acids Res* **32**: 1792-1797.
3. Li, H. (2008). <http://treesoft.sourceforge.net/treebest.shtml>
4. Kent, W.J., et al. (2003). *Proc Natl Acad Sci U S A* **100**: 11484-11489.
5. Kent, W.J., (2002). *Genome Res.* 12: 656-664.