

# GRAMENE COMPARA: A PHYLOGENOMICS RESOURCE FOR PLANTS

Sharon Wei<sup>1</sup>, Zhenyuan Lu<sup>1</sup>, William Spooner<sup>1</sup>, Joshua C. Stein<sup>1</sup>, Andy Yates<sup>3</sup>, Paul Kersey<sup>3</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 Robert W. Holley Center for Agriculture and Health, Cornell University, Ithaca, NY, USA 148533

<sup>3</sup> EMBL, The European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton Cambridge CB10 1SD UK

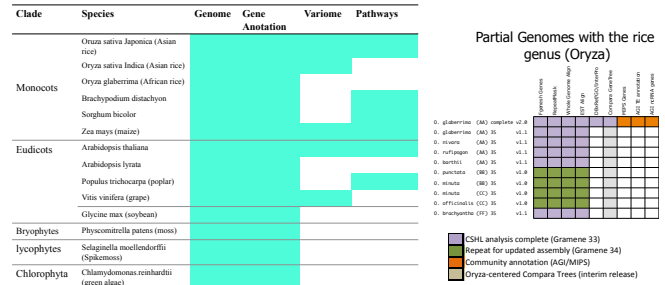
<http://www.gramene.org>  
<http://plants.ensembl.org>

European Bioinformatics Institute  
Cold Spring Harbor Laboratory  
Oregon State University  
Cornell University

## Abstract

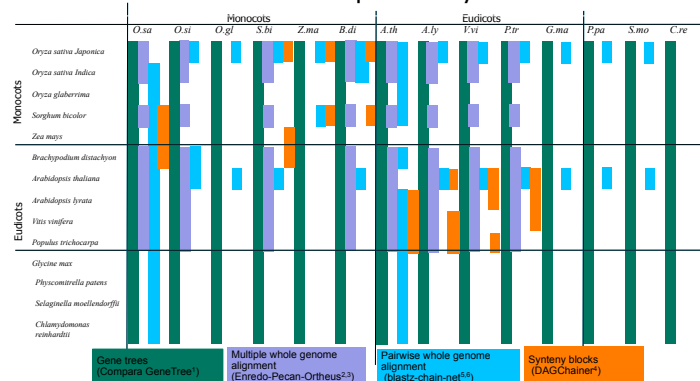
The integration of genome annotation with evolutionary analysis, often referred to as phylogenomics, is a powerful strategy in the study of gene structure and function, and is a compelling motivation for acquiring complete genome sequences. The Gramene Project ([www.gramene.org](http://www.gramene.org)) in collaboration with Ensembl Plants (<http://plants.ensembl.org>) provides a comprehensive platform for comparative genomics in plants, utilizing the Ensembl Compara pipelines and database structure. The site offers data and visualizations of whole genome alignments, synteny analysis, phylogenetic trees, and ortholog/paralog designations. Release 34 includes the whole genomes of six monocots (rice japonica, rice indica, African rice, sorghum, Brachypodium, and maize), five dicots (Arabidopsis, A. lyrata, grape, soybean and poplar), two basal land plants (Physcomitrella, and Selaginella) and the green alga Chlamydomonas. Through collaboration with AGI (<http://www2.genome.arizona.edu/>), Gramene also hosts partial genomes of seven non-cultivated members of the Oryza genus. Comparative analyses include gene tree construction with ortholog/paralog inference, pairwise synteny maps based on phylogenetically-determined orthologs, 8-way multi-species whole genome alignments with ancestor reconstruction using the Enredo/Pecan/Orthus pipeline and pairwise whole genome alignments using blastz-chain-net. These data are fully integrated with other Gramene resources, including gene and protein-level annotations, GO ontology, genome browsers, diversity data, and pathways. In addition, Rice japonica, Arabidopsis, Physcomitrella and Chlamydomonas are included as plant representatives in the Clade Compara to offer a broad view of homologous relationships from across a wide range of taxonomic groups from eukaryotes to prokaryotes. These data are accessible through genome browser, via the public MySQL database, and ftp flat file downloads. We describe details of this resource and demonstrate its use in multiple applications, including the definition of duplication events, large and small-scale rearrangements, annotation inconsistencies, and comparison of gene-family diversity across species. The availability of this platform provides unique opportunities to elucidate the evolutionary history of flowering plants.

## Genome-scale Datasets in Gramene

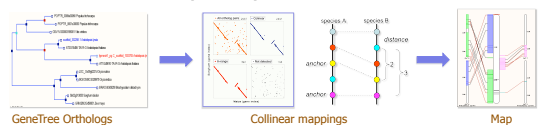


Gramene 34 hosts 14 whole genomes and 8 partial genomes of chromosome 3 short arm from Oryza Genomes Evolution project (previously known as OGPAP <http://www.ogpap.org/>). All of the 14 whole genomes have community gene annotations, most of them have finished predictions, protein domain annotation, cross references to other databases, GO annotations, some have variation data and pathways. The OGPAP partial genomes have finished gene predictions. They have compared gene tree built within the oryza clade, and whole genome alignments against the same reference genome Oryza sativa Japonica.

## Gramene Compara Analyses

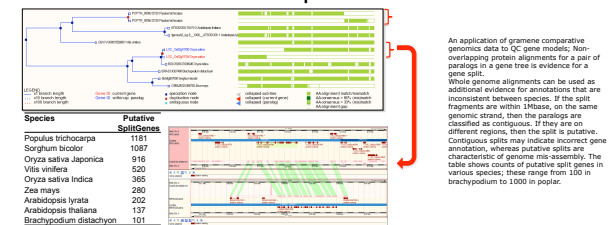


## Synteny Detection

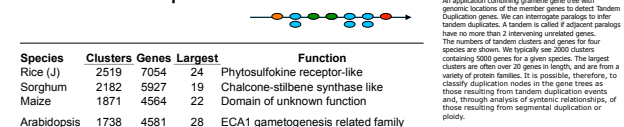


Synteny: genes and regions were calculated using gene orthologs predicted from Gramene Protein Compara analysis. DAGChainer<sup>4</sup> was used to identify collinear chains. To increase sensitivity, chaining was created on the basis of gene order, excluding positions of non-orthologous genes rather than using gene coordinates. Chains were required to have at least five collinear genes with no more than ten intervening genes between neighbors. Resulting gene-pairs were classified as 'syntentic-collinear'. Next, we searched for additional synteny among non-collinear genes to account for small-scale rearrangements and assembly errors. This also allowed members of tandemly duplicated clusters to be classified as syntentic, as they were often missed when strict collinearity was enforced. For a non-collinear gene in species, nearest flanking orthologs were identified and mapped to the respective positions in species. If the non-collinear ortholog in species was located within five genes of either anchor, then the non-collinear pair of genes was classified as 'syntentic-in-range'. Next of these occurred within defined collinear chains. The coordinates of collinear chains were used to define syntentic regions.

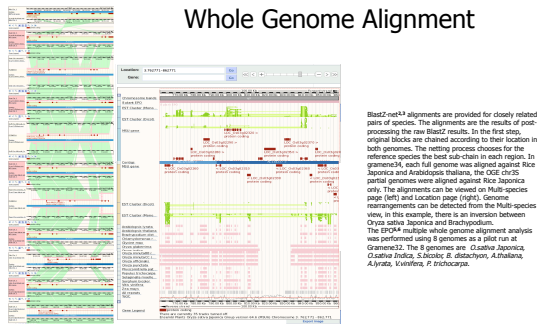
## Automated Detection of Split Genes



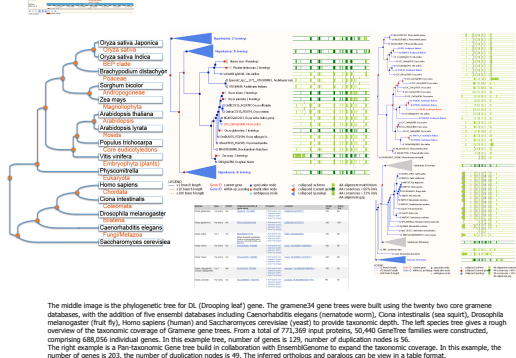
## Tandem Duplicate Detection



## Whole Genome Alignment



## Gene Tree



## References

- Vielze A.J., Severin J., Ureta-Vidal A., Durbin R., Heng L., Birney E. EnsemblCompara GeneTrees: Analysis of complete, duplication aware phylogenetic trees in vertebrates. Genome Research 2009; 19:327-335
- Paton B., Hurren J., Beal K., Fitzgerald S., Birney E. Enredo and Pecan: genome-wide mammalian consistency-based multiple alignment with paralogs. Genome Res. 2008;18:1614-1628
- Paton B. et al. Genome-wide nucleotide-level mammalian ancestor reconstruction. Genome Res. 2008;18:1829-1843
- Huan B. et al. 2004. DAGChainer: a tool for mining segmental genome duplications and synteny. Bioinformatics 20(16):3634-6
- Schwarz S. et al. Human-Mouse Alignments with BLASTZ. Genome Res. 13(1):103-7
- Kent WJ et al. Evolution's cauldron: Duplication, deletion, and rearrangement in the mouse and human genomes. Proc Natl Acad Sci U S A. 2003;100(20):11484-9

## The Gramene-Ensembl Plants collaboration

Gramene and Ensembl Plants are collaborating to maintain a common set of reference databases for plant genomic sequence and annotation, integrating data for important plant species generated from across the globe.



[www.gramene.org](http://www.gramene.org)  
Sharon Wei  
516-367-8828

Cold Spring Harbor Laboratory