# Lawrence Berkeley National Laboratory
## LBL Publications

**Title**

JGI Plant Genomics Gene Annotation Pipeline

**Permalink**

https://escholarship.org/uc/item/8j463183

**Authors**

Shu, Shengqiang
Rokhsar, Dan
Goodstein, David
et al.

**Publication Date**

2014-07-15

# JGI Plant Genomics Gene Annotation Pipeline

Shengqiang Shu[1], David M. Goodstein[1], David Hayes[1], Therese Mitros[2] & Dan Rokhsar[1]

[1]Lawrence Berkeley National Laboratory/DOE Joint Genome Institute, Walnut Creek, California 94598, USA.
[2]University of California – Berkeley, Berkeley, California

*To whom correspondence should be addressed*:  S. Shu (sqshu@lbl.gov @lbl.gov)

July 2014

## ACKNOWLEDGMENTS:

## DISCLAIMER:

# JGI Plant Genomics Gene Annotation Pipeline

## Shengqiang Shu

### Richard Hayes, Therese Mitros, David Goodstein and Daniel Rohksar

**JGI** JOINT GENOME INSTITUTE DEPARTMENT OF ENERGY
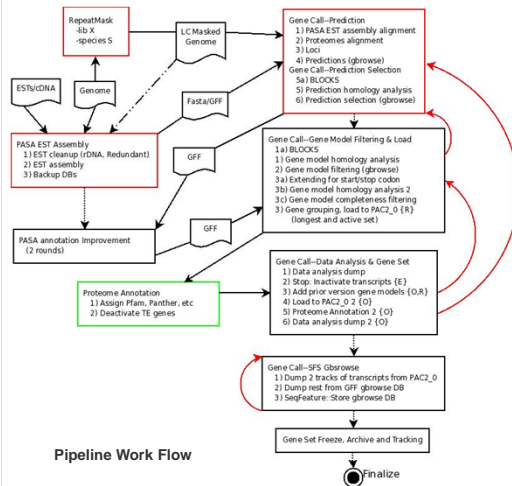
## Abstract

JGI plant genomics gene annotation pipeline is a streamlined pipeline integrating multiple components. It provides one or more optional redo operations in each component to allow parameter tweaking to suit each genome characteristics so a naïve or experienced user can make high quality gene annotation in a timely manner. It has been used for gene annotation/re-annotation of many JGI flagship plant genomes and other genomes, producing stable and high quality gene models.

## Methods

**Prerequisites**: genome fasta, transcriptome fasta, related species proteomes, repeat library

**Steps** (detailed work flow shown below):
1) Run RepeatModeler to get repeat library if not available
2) RepeatMask genome with repeat library
3) Run PASA to assembly transcripts
4) Run BLASTX and EXONERATE on homolog proteomes
5) Define locus based on EXONERATE hits and transcript assembly alignments
6) Make initial gene predictions for each locus by GenomeScan, FGENESH+ and FGENESH_EST
7) Select one best prediction for each locus
8) Run PASA on the selected predictions for annotation improvement like adding UTRs and alternative transcripts
9) Filter gene models
10) Run protein domain analysis on the filtered gene models and deactivate TE genes
11) Build browser database/datastore for visualization

**Pipeline Work Flow**

Contact: S. Shu (sqshu@lbl.gov)

## Introduction

Plant genomes vary in size and are highly complex with a high amount of repeats, genome duplication and tandem duplication. Gene encodes a wealth of information useful in studying organism and it is critical to have high quality and stable gene annotation. Thanks to advancement of sequencing technology, many plant species genomes have been sequenced and transcriptomes are also sequenced. To use these vastly large amounts of sequence data to make gene annotation or re-annotation in a timely fashion, an automatic pipeline is needed. JGI plant genomics gene annotation pipeline, called integrated gene call (IGC), is our effort toward this aim with aid of a RNA-seq transcriptome assembly pipeline. It utilizes several gene predictors based on homolog peptides and transcript ORFs. See Methods for detail.

Here we present genome annotation of JGI flagship green plants produced by this pipeline plus Arabidopsis and rice except for chlamy which is done by a third party. The genome annotations of these species and others are used in our gene family build pipeline and accessible via JGI Phytozome portal whose URL and front page snapshot are shown below.

### http://phytozome.jgi.doe.gov/pz/portal.html

## Results

Land plant protein coding gene structure is shown to be remarkably stable in term of exon number and exon length, and to some extent intron length as well (shown in far right tables).

Using percentage of gene models with PFAM domain as a yardstick, our automatic pipeline produced results approaching human curated Arabidopsis and many genomes are better than rice, also a well-known annotated genome. Except for chlamy (a non land plant), genome annotation with lower assigned PFAM domain reflects largely on quality of inputs.

| | Gene | Trans | %PFAM | G. Size | %Exonic |
|---|---|---|---|---|---|
| G. max | 56,044 | 88,647 | 75.1 | 979M | 9.66 |
| P. trichocarpa | 41.335 | 73,013 | 74.9 | 434M | 15.33 |
| B. distachyon | 31,694 | 42,868 | 70.7 | 272M | 19.38 |
| S. italica | 35,471 | 40,599 | 66.4 | 406M | 12.01 |
| S. bicolor | 33,032 | 39,441 | 67.3 | 727M | 8.48 |
| P. virgatum | 98,007 | 125,439 | 55.6 | 1,698M | 7.69 |
| P. patens | 26,610 | 42,392 | 68.4 | 473M | 10.44 |
| C. reinhardtii | 17,741 | 19,526 | 51.8 | 111M | 52.03 |
| A. thaliana | 27,416 | 35,386 | 77.4 | 120M | 34.19 |
| O. sativa | 39,049 | 49,061 | 62.0 | 375M | 15.21 |

## Results

| | Number of Exons | | | | |
|---|---|---|---|---|---|
| | 25% | Median | Mean | 75% | Max |
| G. max | 3 | 5 | 6.5 | 9 | 78 |
| P. trichocarpa | 2 | 5 | 6.3 | 9 | 76 |
| B. distachyon | 2 | 4 | 5.8 | 8 | 72 |
| S. italica | 2 | 3 | 5.0 | 7 | 60 |
| S. bicolor | 2 | 4 | 5.3 | 7 | 65 |
| P. virgatum | 2 | 3 | 4.5 | 6 | 60 |
| P. patens | 2 | 5 | 6.3 | 8 | 70 |
| C. reinhardtii | 4 | 7 | 8.9 | 11 | 173 |
| A. thaliana | 2 | 4 | 5.9 | 8 | 79 |
| O. sativa | 2 | 3 | 4.9 | 7 | 78 |

| | Exon Length | | | | |
|---|---|---|---|---|---|
| | 25% | Median | Mean | 75% | Max |
| G. max | 90 | 151 | 284 | 330 | 9,509 |
| P. trichocarpa | 89 | 149 | 277 | 322 | 7,911 |
| B. distachyon | 92 | 158 | 306 | 366 | 9,439 |
| S. italica | 91 | 157 | 288 | 345 | 7,851 |
| S. bicolor | 94 | 173 | 361 | 451 | 14,531 |
| P. virgatum | 94 | 171 | 316 | 385 | 18,512 |
| P. patens | 93 | 162 | 313 | 378 | 10,268 |
| C. reinhardtii | 95 | 156 | 369 | 351 | 12,274 |
| A. thaliana | 89 | 147 | 262 | 300 | 7,761 |
| O. sativa | 90 | 159 | 318 | 369 | 15,363 |

| | Intron Length | | | | |
|---|---|---|---|---|---|
| | 25% | Median | Mean | 75% | Max |
| G. max | 105 | 225 | 519 | 599 | 18,215 |
| P. trichocarpa | 101 | 180 | 380 | 483 | 10,053 |
| B. distachyon | 94 | 149 | 419 | 500 | 16,886 |
| S. italica | 92 | 136 | 334 | 418 | 6,930 |
| S. bicolor | 96 | 150 | 480 | 491 | 18,859 |
| P. virgatum | 92 | 137 | 391 | 422 | 18,637 |
| P. patens | 153 | 215 | 274 | 312 | 7,589 |
| C. reinhardtii | 164 | 227 | 173 | 314 | 82,837 |
| A. thaliana | 86 | 100 | 165 | 168 | 11,602 |
| O. sativa | 96 | 168 | 416 | 501 | 18,327 |

## References

Haas, B.J., Delcher, A.L., Mount, S.M., Wortman, J.R., Smith Jr, R.K., Jr., Hannick, L.I., Maiti, R., Ronning, C.M., Rusch, D.B., Town, C.D. et al. (2003) Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. Nucleic Acids Res, 31, 5654-5666.

Smit, AFA, Hubley, R & Green, P. RepeatMasker Open-3.0. 1996-2011 <http://www.repeatmasker.org>.

Yeh, R.-F., Lim, L. P., and Burge, C. B. (2001) Computational inference of homologous gene structures in the human genome. Genome Res. 11: 803-816.

Salamov, A. A. and Solovyev, V. V. (2000). Ab initio gene finding in Drosophila genomic DNA. Genome Res 10, 516-22.