# GenSAS

### **Gen**ome **S**equence **A**nnotation Server

*Computational annotation and curation
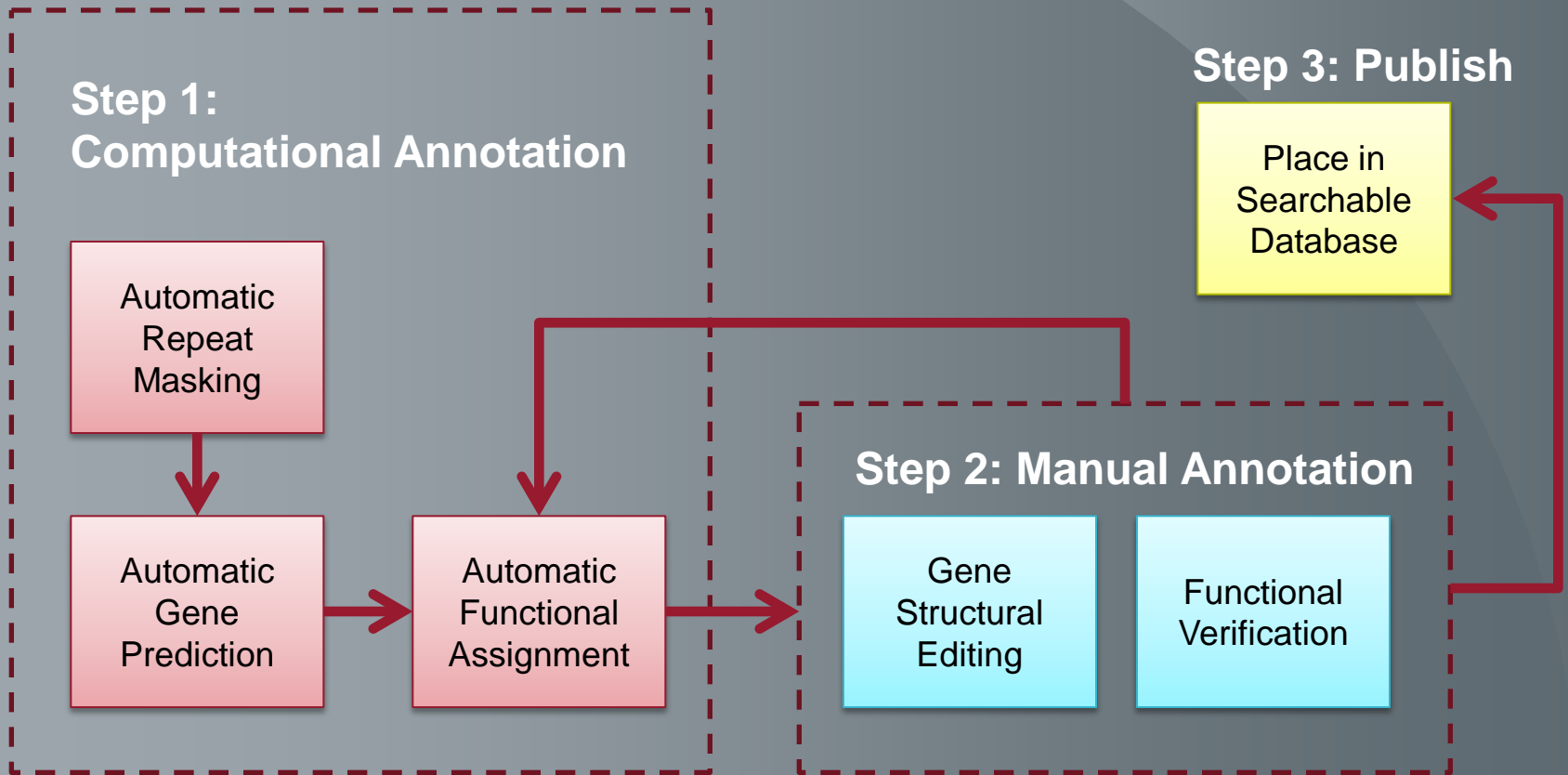of genome sequences*

Stephen Ficklin*, Taein Lee, Jodi Humann, Chun-Huai Cheng, Jill Wegrzyn, David Neale, Dorrie Main

*\* Presenting author*

*Pine Genome Workshop
PAG XXIII, San Diego CA.  Jan 10, 2015.*

**WASHINGTON STATE UNIVERSITY**

# Genome Annotation Workflow

**Step 1:**
**Computational Annotation**

**Step 3: Publish**

Place in Searchable Database

**Step 2: Manual Annotation**

Automatic Repeat Masking

Automatic Gene Prediction

Automatic Functional Assignment

Gene Structural Editing

Functional Verification

WASHINGTON STATE UNIVERSITY

# Existing Annotation Tools

- Computational Annotation (Automatic)
  - Maker
  - Ergatis
  - JAMg (Just Annotate My Genome)
  - iPlant DNA Subway

- Manual Annotation (Curation)
  - Apollo
  - WebApollo
  - Manatee

- These tools provide excellent annotation support but an integrated online solution to bridge all three workflow steps is lacking.

# GenSAS Goals

- Streamline genomic annotation
  - Integrate all three stages of the annotation workflow
  - An easy-to-use online application

- Simplify pipeline setup for Bioinformaticist

- Support integration of final results into an online genome database for publication.

- Support analysis of very large genomes
  - Such as 22Gb *Pinus taeda* (Loblolly pine)
  - Using high-performance computing clusters

- Educate students and researchers
  - Inline instructions and best practices
  - Online videos & tutorials

# GenSAS Implementation

- ## A Drupal module
  - Drupal is an open-source, free Content Management System (CMS)
  - Used by millions of websites world-wide
  - High profile sites include:



**US White House**
http://www.whitehouse.gov

**British Medical Journal**
http://www.bmj.com/

**Popular Science**
http://www.popsci.com/

  - A Drupal site can support social  and content needs for a community in addition to functionality provided by GenSAS.
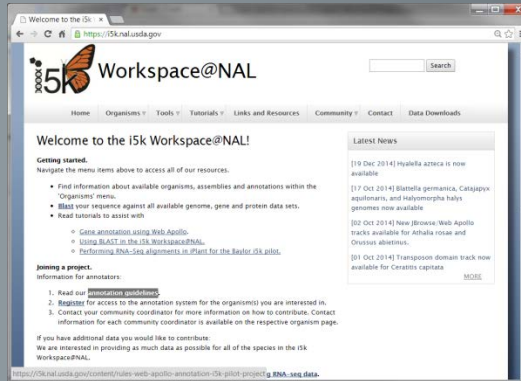  - Integrates easily with Tripal

# Tripal

- Construction of online genomic websites
  - A suite of Drupal modules
  - Uses the GMOD Chado schema for data storage
  - Decreases cost and man-power required to create online genomics, genetics and breeding databases
  - Provides an Application Programming Interface (API)
  - Tripal development is currently funded by:
    - NSF DIBBs Award #1443040 (2015-2018)
    - USDA NRSP10 Award (2014-2019)
    - USDA SCRI Award # 2014-51181-22376 (2014-2019)
    - Pending support from several NSF PGRP proposals (include funds to convert Dendrome/TreeGenes to Tripal)

http://tripal.info

# Example Sites Using Tripal



**Banana Genome Hub**
http://banana-genome.cirad.fr/

**i5K Workspace@NAL**
http://i5k.nal.usda.gov/

**Genome Database for Rosaceae**
http://www.rosaceae.org

**Knowpulse: pulse crop genomics & breeding**
http://knowpulse2.usask.ca/portal

**Legume Information System**
http://legumeinfo.org/

**CottonGen**
http://www.cottongen.org

# GenSAS + Tripal

- Completes the "publish" stage
    - GenSAS fully integrates into the same Drupal site that uses Tripal
    - Results generated by GenSAS are 100% compatible
    - Versioning of results in GenSAS ensures no conflict with future annotation updates.
    - Tripal provides searching and pages (e.g. gene pages)

- Three options:
    - GenSAS expanded later with Tripal for full publication of the genomic annotations
    - GenSAS added to an existing Tripal site
    - GenSAS can be installed as a separate tool for non Tripal sites.

Tripal

GenSAS

# GenSAS Availability

- ## Public server:
  - http://gensas.bioinfo.wsu.edu
  - Resources are currently limited, thus accounts are given on a case-by-case basis
  - Anyone can request an account
  - Currently houses *Pinus taeda* Maker annotations



- ## Source Code
  - http://drupal.org/project/gensas
  - Available soon…
  - download and install on local computational resources.

# Functionality Overview

- ## GenSAS v3.0 provides:
  - ### Repeat Finding & Masking:  repeat libraries and *de novo*
    - *RepeatMasker, RepeatModeler*
  - ### Gene Prediction
    - #### Intrinsic
      - using tools the use heuristics and the genomic sequence
      - *Augustus, FgeneSH (parsing only), Genscan, Glimmer3, GlimmerM, SNAP, tRNAScan, getorf*
    - #### Extrinsic
      - using tools that use transcript or protein libraries to assist with gene prediction.
      - *BLAT, BLAST*
  - ### Gene Consensus Prediction
    - *EvidenceModeler*
  - ### Gene Visualization and Curation
    - *WebApollo, Jbrowse*
  - ### Publish
    - Generates GFF3 & FASTA files with properly versioned naming.

# Using GenSAS



View the GenSAS Computer Demo C03 at 2:10pm today in California Room for further detail, plus Poster # P1153

# Step 1: Start a project

# Step 2: Upload Sequences

# Step 3: Upload Supporting Files

**Evidence Files**

One of the best ways to identify genes is to provide transcript (e.g. ESTs, full length cDNA, RNA-seq) and protein files from the species or from closely related species. These files should be provided in FASTA or FASTQ format, and when aligned to the genomic sequence serve as direct evidence of gene expression. Transcripts and proteins of closely related species can be provided when species-specific files are not available or are insufficient.

Supplying a species-specific (or closely-related species) file of known repeats, in FASTA format, can help with identification of repetitve elements more accurate for the species.

**Uploaded Evidence Files**

| File | Type | Size | Action |
|------|------|------|--------|
| There are currently no files uploaded. | | | |

▶ Repeat Library

▶ EST Evidence

▶ Protein Evidence

▶ RNA-seq Evidence (currently unavailable)

# Step 3: Upload Supporting Files



Current *Pinus taeda* Maker predictions loaded in this way

WASHINGTON STATE UNIVERSITY

# Step 4: Repeat Finding



▾  RepeatMasker

*RepeatMasker is a program that screens DNA sequences for interspersed repeats and low complexity DNA sequences. The output of the program is a detailed annotation of the repeats that are present in the query sequence as well as a modified version of the query sequence in which all the annotated repeats have been masked (default: replaced by Ns). On average, almost 50% of a human genomic DNA sequence currently will be masked by the program. Sequence comparisons in RepeatMasker are performed by one of several popular search engines including, cross_match, ABBlast/WUBlast, RMBlast and Decypher.*

**Job Name**

RepeatMasker

Please provide a name for this RepeatMasker job. It is used to distinguish between two or more RepeatMasker jobs.

**Search Engine**

ncbi ▾

The search engine to use (e.g. wublast, cross_match)

**Speed / Sensitivity**

quick ▾

The speed vs. sensitivity parameter

**DNA Source**

Select an option ▾

Specify the species or clade of the input sequence. The species name must be a valid NCBI Taxonomy Database species name and be contained in the RepeatMasker repeat database.

**Repeat Options**

Mask interspersed and simple repeat ▾

Select the types of repeats you would like to mask.

# Step 5: Masking

**Repeat Masking**

▸ Instructions

**Final Masking Selection**

Select the repeat masking jobs that you want to include in masking of the genomic sequence prior to gene discovery.

| Masking Job | Status |
|---|---|
| ☐ RepeatMasker-slow | Completed |
| ☐ RepeatMasker | Completed |

[ Mask Sequences ]        [ Skip Repeat Masking Step ]

# Step 6: Gene Prediction

# View Results

# Step 7: Consensus

Project  Sequences  Files  Repeats  Masking  Genes  **Consensus**  Refine  Functional  Annotate  Publish

Welcome!  x  Gene Consensus

## Identification of Genes and other Features

▶ Instructions

Select the gene finding jobs that you want to include in building the consensus gene predictions and provide a numerical weight indicating the expected accuracy of the predictions. Higher weights indiciate higher accuracy.

| Gene Finding Job | Status | Weight |
|---|---|---|
| **Gene Prediction** | | |
| Augustus | Completed | |
| Augustus-complete genes only (Augustus) | Completed | |
| Genscan | Completed | |
| GlimmerM | Completed | |
| SNAP | Completed | |
| **Protein Alignments** | | |
| BLAST (proteins) | Completed | |
| **Transcript Alignments** | | |
| BLAST (nucleotide) | Completed | |
| BLAT | Completed | |

# Step 8: Manual Curation



Curation Track

# Step 9: Publish

## Available Results for Publishing

*Please choose the jobs to be included in the published release for this project.*

**Consensus Gene Predictions**

☑ Consensus Gene Set (EvidenceModeler)

The consensus gene set will be the primary set of genes in your published annotation set, therefore, it is selected by default to be published.

**Consensus Masking**

☑ Masked Consensus

The repeat masked consensus job created the FASTA sequence on which all other predictions were made. This job should be included in any published release.

**Gene Predictions**

☐ Augustus (Augustus)

☐ Genscan (Genscan)

☐ GlimmerM (GlimmerM)

☐ SNAP (SNAP)

Because you have a consensus gene prediction set you do not need to publish additional singular gene prediction results. The consensus will be the set used by others and will be considered candidates for future annotations. However, these gene predictions help provide evidence for how the consensus genes were constructed. If you would like to provide this information please include any gene prediction jobs as desired.

**Repeats & Masking**

☑ RepeatMasker (RepeatMasker)

☐ RepeatModeler (RepeatModeler)

Jobs that were used in construction of the consensus masked sequence are selected by default.

**Protein Alignments**

☐ BLAST (proteins) (BLAST (proteins))

Protein alignments help provide evidence for the predicted gene models and provide clues for functional assignment. Consider including protein alignments that help provide support. Avoid including any with an overwhelming set of alignments such as alignments against an all-inclusive protein database. It is best to include alignments to species-specific or species-related alignments.

# Step 9: Publish

**Annotation version number**

`1`

In order to prevent naming conflicts with past or future annotations for the same genomic sequence assembly version, the annotions should have a separate version. If this is the first time that the genomic sequence will be annotated then the annotation version should be 1. When a new set of predictions are created for the same genomic sequence then the annotation version number should be incremented by 1. GenSAS will include both the sequence version number and the annotation verson number in the names of predicted features. For example, if the genomic sequence version is v1.0 and this is the first annotation set, GenSAS will include "v1.0.a1" in the feature names to indicate the genomic sequence and annotation set to which the predicted features belong.

Publish

# Ongoing Work

- Functional Annotation
  - Does not yet fully support functional annotation.
  - Requires addition of functional tools and result viewers.
  - Currently under development
  - Will be available in a future release of GenSAS

- Support of High-Performance Computing
  - Currently executes on stand-alone server
  - Will integrate with two types of HPC schedulers: PBS, GE
  - Currently under development
  - Will be available in a future release of GenSAS

- Support of RNA-Seq Datasets

- Integration with PASA for gene refinements

- Full Integration with Tripal
  - Currently only exported files are compatible with Tripal loaders
  - Integrate so that publish button can automatically import into Tripal.

# Funding

- ## Current Funding:
  - The work presented herein was funded by the SDA/NIFA (2011-67009-30030) subaward to Dorrie Main at Washington State University (PI: David Neale at University of California, Davis)

- ## Continued Funding
  - Further funding for development and implementation of GenSAS is provided by an USDA National Research Support Project (NRSP10) to WSU (Dorrie Main PI) through 2019.

# GenSAS for Pine Annotation

- ## Current capabilities

  - *Pinus Taeda v1.01* Maker annotations loaded
  - Pine community can annotate the predicted genes using integrated WebApollo
  - WebApollo manages user access (via GenSAS integration)
  - Annotation training may be available via WebApollo outreach for community.

- ## Future capabilities

  - Functional Annotation:

    - Execute functional tools inside GenSAS

    - Result viewers for tools such as InterProScan, SignalP, TargetP, Pfam, blastp, etc.

  - Tools will execute on high-performance computing to decrease time requirements
  - Future pine genomes can be annotated directly in GenSAS

# Loblolly Pine in GenSAS

# Thank You!



## More GenSAS at PAG

- Computer Demo: C03 at 2:10pm today, Jan 10th in California Room for further details

- Poster:  P1153.  Meet w/ Jodi Humann from 3.00-4.30 PM Monday, Jan 12th

- Contact Info: dorrie@wsu.edu, stephen.ficklin@wsu.edu, jhumann@wsu.edu

WASHINGTON STATE UNIVERSITY