# GMOD Project Update
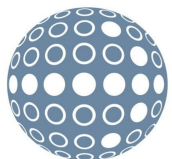
Scott Cain
GMOD Project Coordinator
Ontario Institute for Cancer Research
scott@scottcain.net

GMOD Meeting
San Diego, CA
January 14-15, 2010

Ontario Institute
for Cancer Research

# Introduction: GMOD is …

- A set of interoperable open-source **software** components for visualizing, annotating, and managing biological data.

- An active **community** of developers and users asking diverse questions, and facing common challenges, with their biological data.

# Who uses GMOD?



Plus hundreds of others

# Software

GMOD components can be categorized as

**V** Visualization

**D** Data Management

**A** Annotation

# Visualization: GBrowse

## GBrowse

JBrowse

GBrowse_syn

CMap



**Releases**
  1.70 released
  2.0, 1.71 in the pipe

**AJAX/Interface:**
  Rubberband region selection, drag and drop track ordering, collapsible tracks, popup balloons, asynchronous rendering (2.0)
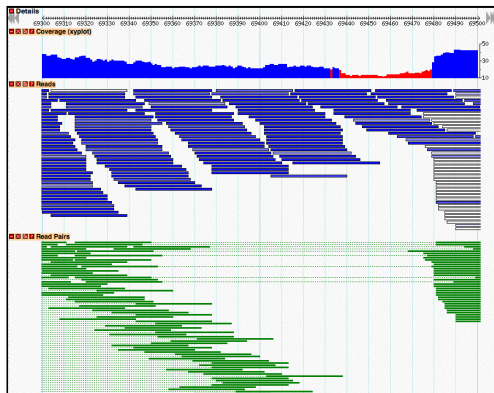
**Biology:**
  Allele/gentotype frequency, LD glyphs, geolocation popups, circular genome support (1.71)

**Infrastructure**
  User logins, server multiplexing (2.0), SQLite and SAMtools (NGS) adaptors

**modENCODE Fly:**
  http://modencode.oicr.on.ca/cgi-bin/gb2/gbrowse/fly/

The generic genome browser: a building block for a model organism system database. Stein LD et al. (2002) *Genome Res* 12: 1599-610
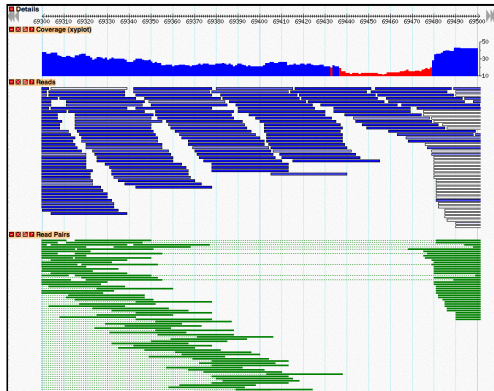
# Visualization

## GBrowse

JBrowse

GBrowse_syn

CMap



### Resources

Tutorials (http://gmod.org/wiki/GBrowse_Tutorial):
  GBrowse User Tutorial at OpenHelix.com
  GBrowse Admin Tutorial
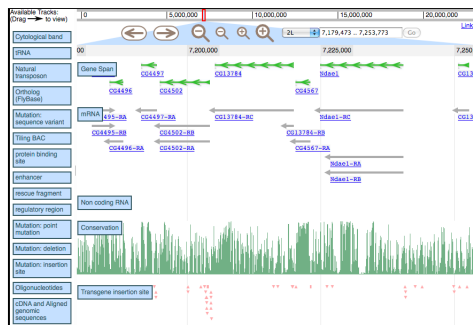  NGS in GBrowse and SAMtools Tutorial

Web Sites:
  GMOD          http://gmod.org/wiki/GBrowse
  WebGBrowse    http://webgbrowse.cgb.indiana.edu/
  GBrowse.org   http://gbrowse.org

Mailing List:
  https://lists.sourceforge.net/lists/listinfo/gmod-gbrowse

# Visualization

## GBrowse

## JBrowse

## GBrowse_syn

## CMap

### GMOD's 2nd Generation Genome Browser
### It's *fast*

Completely new genome browser implementation:
Client side rendering
Heavy use of AJAX
Uses JSON and Nested Containment Lists

JBrowse Fly:
http://jbrowse.org/genomes/dmel/

Web Sites:
GMOD          http://gmod.org/wiki/JBrowse
JBrowse       http://jbrowse.org

Mailing List:
https://lists.sourceforge.net/lists/listinfo/gmod-ajax

JBrowse: A next-generation genome browser, Mitchell E. Skinner, Andrew V. Uzilov, Lincoln D. Stein, Christopher J. Mungall and Ian H. Holmes, *Genome Res.* 2009. 19: 1630-1638

# Visualization

GBrowse

JBrowse

**GBrowse_syn**

CMap

GBrowse based comparative genomics viewer
Shows a reference sequence compared to 2+ others
Can also show any GBrowse-based annotations



wormbase.org

Syntenic blocks do not have to be colinear
Can also show duplications

http://gmod.org/wiki/GBrowse_syn

Sheldon McKay, Cold Spring Harbor Laboratory

# Visualization

GBrowse

JBrowse

GBrowse_syn

CMap

Web based comparative map viewer
CMap is data type agnostic:
  Can link sequence, genetic, physical, QTL, deletion,
    optical, …
CMap 2.0 coming
  Faster, internals cleanup
  Circos export

# Data Management

## Chado

Tripal

TableEdit

BioMart

InterMine

A extensible, modular database schema for storing biological data

1.0 release:
- Stable schema
- Tools for data in/out

1.1 release (soon):
- Stable schema (minor, nondestructive changes)
- Improvements to data loading scripts
- Additional modules: cell line, natural diversity
- Tool for managing materialized views
- Tool for creating ontology-based views

# Data Management

Chado

Tripal

TableEdit

BioMart

InterMine

New web front end for Chado databases

Set of Drupal modules

Modules approximately correspond to Chado modules

Easy to create new modules

Includes user authentication, job management,
        curation support



Stephen Ficklin, Meg Staton, Chun-Huai Cheng, …
Clemson University Genomics Institute

# Data Management

Chado

Tripal

**TableEdit**

BioMart

InterMine

MediaWiki extension
 MediaWiki software
 used at Wikipedia,
 GMOD.org, …

GUI to wiki tables
Also a GUI to
 database tables
Work in progress to
 use this with Chado

**Potential to give
wiki access to
Chado databases**



Example: GONUTS (http://gowiki.tamu.edu/)

Jim Hu, Daniel Renfro, *et al*., Texas A&M

# Data Management

Chado

Tripal

TableEdit

**BioMart**

InterMine



New GUIs - more configurable and easier to use

Virtual Marts - marts running off source schema without materializing

Improved scalability

Security and access control

Improved federation

New configuration tool



??

# Data Management

Chado

Tripal

TableEdit

BioMart

**InterMine**

Data integration and web-based query package

Now supports ~20 common data formats:
GFF3, Chado, GO annotation, biopax, BioGrid, TreeFam, PubMed, Ensembl, …

Interfaces: RESTful web service, Java & Perl APIs

Upload & analyse gene lists with graphical and statistical widgets



FlyMine: an integrated database for *Drosophila* and *Anopheles* genomics, Rachel Lyne, *et al.*, *Genome Biol.* 2007; 8(7): R129.

# Annotation

**MAKER**

DIYA

Galaxy

Ergatis

Apollo

## MAKER
### Annotate this!

Genome annotation pipeline for creating gene models

Output can be loaded into GBrowse, Apollo, Chado, …

Incorporates

      SNAP, RepeatMasker, exonerate, BLAST,

      Augustus, FGENESH, GeneMark, MPI

Other capabilities

      Map existing annotation onto new assemblies

      Merge multiple legacy annotation sets into a consensus set

      Update existing annotations with new evidence

      Integrate raw InterProScan results

MAKER Web Annotation Service - MAKER online

MAKER

**DIYA**

Galaxy

Ergatis

Apollo



Lightweight, modular, and configurable Perl-based pipeline framework

Initial application is gene prediction for prokaryotes

Working on integration of Amos assembly tools

MAKER

DIYA

Galaxy

Ergatis

Apollo

NGS tools support
QC and Manipulation: FASTQ, 454, SOLiD support
Mapping: Bowtie or BWA, Megablast
SAMtools: Web interface to SAMtools scripts

LIMS system in beta.
Import data from your sequencer into Galaxy



DIYA: a bacterial annotation pipeline for any genomics lab, Andrew C. Stewart, Brian Osborne and Timothy D. Read, *Bioinformatics* 2009 25(7):962-963

# Annotation

MAKER

DIYA

Galaxy

Ergatis

Apollo



- Currently up to 162 analysis tool / components for use in pipelines

- Updated prokaryotic annotation pipeline template

- Updated comparative annotation pipeline template

- Lots of work for use on Amazon EC2

- Now the engine behind the new CloVR cloud computing project (http://clovr.igs.umaryland.edu/)

http://ergatis.sourceforge.net/

# Annotation

MAKER

DIYA

Galaxy

Ergatis

Apollo

Better Chado support

  including DBMS independent support)

GFF3 support

GUI based configurations

Multiple alignment transcript viewer and editor
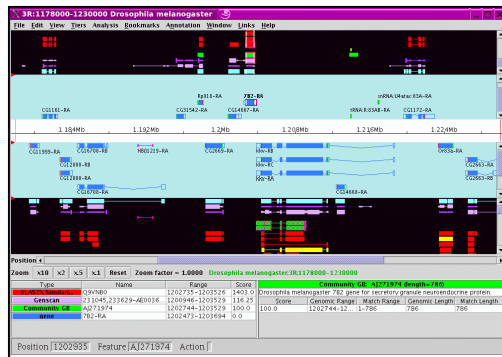
Continuous data display

  sgr, wiggle

Remote analysis to NCBI services

  BLAST, Primer-BLAST

Undo support

More robust Java Web Start support

DIYA: a bacterial annotation pipeline for any genomics lab, Andrew C. Stewart, Brian Osborne and Timothy D. Read, *Bioinformatics* 2009 25(7):962-963