# Database Tools
# for Biologists

Dave Clements
GMOD Help Desk
US National Evolutionary Synthesis Center (NESCent)
clements@nescent.org

Texas A&M University

Biochemistry and Biophysics

8 December 2009

**Sponsored by**

**Alliance for
Bioinformatics,
Computational Biology,
and Systems Biology**

**Laboratory for
Genome Bioinformatics**

NESCent

A&M

# Agenda

| | |
|---|---|
| 3:00 | Introduction |
| 3:10 | Software<br><br>    Visualization<br>        GBrowse<br>        JBrowse<br>        GBrowse_syn<br>        Sybil & SynView<br>        CMap |

| | |
|---|---|
| 4:00 | Software, cont.<br><br>    Data Management<br>        Chado, Tripal, GMODWeb<br>        BioMart and InterMine<br>        GFF3<br>    Annotation<br>        MAKER & DIYA<br>        Apollo<br>        Textpresso<br>        Community Annotation<br>        Pipelines & Workflows |
| 4:40 | Community |
| 5:00 | Finish |

# This Talk

## Algorithms?



```
INEXACTSEARCH(W,z)
    CALCULATED(W)
    return INEXRECUR(W,|W|-1,z,1,|X|-1)

CALCULATED(W)
    k ← 1
    l ← |X|-1
    z ← 0
    for i=0 to |W|-1 do
        k ← C(W[i])+O'(W[i],k-1)+1
        l ← C(W[i])+O'(W[i],l)
        if k > l then
            k ← 1          NO!
            l ← |X|-1
            z ← z+1
        D(i) ← z

INEXRECUR(W,i,z,k,l)
    if z < D(i) then
        return ∅
    if i < 0 then
        return {[k,l]}
    I ← ∅
    I ← I ∪ INEXRECUR(W,i-1,z-1,k,l)
    for each b ∈ {A,C,G,T} do
        k ← C(b)+O(b,k-1)+1
        l ← C(b)+O(b,l)
        if k ≤ l then
            I ← I ∪ INEXRECUR(W,i,z-1,k,l)
            if b=W[i] then
                I ← I ∪ INEXRECUR(W,i-1,z,k,l)
            else
                I ← I ∪ INEXRECUR(W,i-1,z-1,k,l)
    return I
```

## Plumbing!

# 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.

# Agenda

| 3:00 | Introduction |
|------|-------------|
| 3:10 | **Software**<br><br>Visualization<br>  GBrowse<br>  JBrowse<br>  GBrowse_syn<br>  Sybil & SynView<br>  CMap |

| 4:00 | Software, cont.<br><br>Data Management<br>  Chado, Tripal, GMODWeb<br>  BioMart and InterMine<br>  GFF3<br>Annotation<br>  MAKER & DIYA<br>  Apollo<br>  Textpresso<br>  Community Annotation<br>  Pipelines & Workflows |
|------|-------------|
| 4:40 | Community |
| 5:00 | Finish |

# Software

GMOD components can be categorized as

**V** Visualization

**D** Data Management

**A** Annotation

# Software

## You have

Sequence
Gene models

Mapping data

Alternative
  transcripts
Expression

SNP / variation

Methylation
GO terms
Stocks / lines
Publications /
  Attribution
Orthology

## GMOD Has

[A] MAKER
[A] DIYA
[A] Galaxy
[A] Ergatis

[A] Textpresso
[A] Apollo
[V][A] Table Edit

[V] GBrowse
[V] JBrowse
[V] CMap
[V] GBrowse_syn
[V] Sybil
[V] SynView

[D] Chado
[A][V] Tripal
[V] GMODWeb

[D] BioMart
[D] InterMine

[A] Annotation    [D] Data Management    [V] Visualization

# GMOD Requirements

- ## Server
  - Most use Linux

- ## GMOD Systems Administrator
  - Understands Linux package management, a scripting language, command line interfaces, relational databases, …
  - Grad/Undergrad, half time when starting up.

http://gmod.org/wiki/Computing_Requirements

# GBrowse

GMOD's leading genome browser

Landing page for *E. coli* example

Overview: chromosome / contig wide

Region: intermediate zoom

Details: current area

Tracks: current configuration



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

# GBrowse Example: modENCODE

- Uses GBrowse 2

http://www.modencode.org/gb2/gbrowse/fly/

# GBrowse Tutorials

- GBrowse User Tutorial at OpenHelix
  - Flash based, has handouts, very snazzy and thorough
  - Great resource for your users

- GBrowse Admin Tutorial
  - HTML based, written by Lincoln Stein, mostly
  - Excellent way to learn how to configure GBrowse

- GBrowse Admin Tutorial w/ VMware Image
  - From the 2009 GMOD Summer Schools
  - Gives you a system to start with

- NGS in GBrowse and SAMtools Tutorial
  - From Bioinformatics Australia 2009, October
  - Gives you a system to start with

http://gmod.org/wiki/GBrowse_Tutorial

# GBrowse Future Plans

- Circular genome support
  - Work done by Nathan Liles at Texas A&M
- 2.0, Release in 2010
  - Database and rendering multiplexing
  - Asynchronous track loading
  - GBrowse in the cloud
  - User authentication
- 1.x has a few more maintenance releases left.

# GBrowse Resources

| | |
|---|---|
| Home Page | http://gmod.org/wiki/GBrowse |
| User Tutorial | http://www.openhelix.com/gbrowse |
| Admin Tutorial | http://gmod.org/wiki/GBrowse_Tutorial |
| Configuration | http://gmod.org/wiki/GBrowse_Configuration_HOWTO |
| **WebGBrowse** | http://webgbrowse.cgb.indiana.edu/ |
| **GBrowse.org** | http://gbrowse.org |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-gbrowse |

# JBrowse

- GMOD's 2nd generation genome browser
- It's fast
- Completely new
  - Client side rendering
  - Heavily AJAX
  - JSON, Nested Containment Lists



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

# JBrowse Demo

http://jbrowse.org

# JBrowse Future Plans

- Tools for migrating from GBrowse
- An ecosystem comparable to GBrowse
  - Glyph library, user defined glyphs, callbacks, track sharing, …
- Comparative genomics (more on that later)
- Community Annotation
  - User authentication
  - User uploadable and sharable tracks and annotation

# JBrowse Resources

| | |
|---|---|
| Home Page | http://jbrowse.org |
| Getting Started | http://jbrowse.org/code/jbrowse-master/docs/tutorial/ |
| Admin Tutorial | http://gmod.org/wiki/JBrowse_Tutorial |
| Configuration | http://jbrowse.org/code/jbrowse-master/docs/config.html |
| Demo | http://jbrowse.org/genomes/dmel/ |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-ajax |

# GBrowse or JBrowse

**GBrowse**

Robust ecosystem

Feature rich

Large and growing user base

Track sharing

**JBrowse**

Very fast

Rapidly growing user base

Lots of future development

Easy to configure

# GBrowse_syn

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



Example comparing *C. elegans* to 4 other species at WormBase

Sheldon McKay, Cold Spring Harbor Laboratory

# GBrowse_syn



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

Sheldon McKay, Cold Spring Harbor Laboratory

# GBrowse_syn Future Work

- Integration with GBrowse 2

- High-level graphical overview

- AJAX based user interface and navigation.
  - Submitted grant last month proposing implementing a JBrowse based synteny browser based on GBrowse_syn

# GBrowse_syn Resources

| | |
|---|---|
| Home Page | http://gmod.org/wiki/GBrowse_syn |
| Tutorial | http://gmod.org/wiki/GBrowse_syn_Tutorial |
| User Help | http://gmod.org/wiki/GBrowse_syn_Help |
| Configuration | http://gmod.org/wiki/GBrowse_syn_Configuration |
| Example | http://www.wormbase.org/cgi-bin/gbrowse_syn/ |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-gbrowse |

# SynView and Sybil

## Sybil



Whole Genome Gradient Display



Cluster Report

## SynView

Sybil: Methods and Software for Multiple Genome Comparison and Visualization. Crabtree, *et al.*; in Gene Function Analysis, ed. by Michael F. Ochs (2007)

SynView: a GBrowse-compatible approach to visualizing comparative genome data. Haiming Wang, *et al.*; in Bioinformatics 22 (18)

# GBrowse_syn or Sybil or SynView?

**GBrowse_syn**
Most actively developed
Scalable
Familiar interface
Extensive documentation
Growing user community

**SynView**
Scalable
Runs inside GBrowse

**Sybil**
Scalable
Whole genome and
other unique visualizations
Built on Chado

# CMap



Web based comparative map viewer

CMap is data type agnostic: Can link sequence, genetic, physical, QTL, deletion, optical, …

Particularly popular in plant community

# CMap Future Work

- Streamline the database
- Faster access
- Display in SVG
- Save in Circos / MizBee format
- CMap3D?

# CMap Resources

| | |
|---|---|
| Home Page | http://gmod.org/wiki/CMap |
| User Tutorial | http://www.gramene.org/tutorials/cmap.html |
| Admin Guide | http://gmod.svn.sourceforge.net/viewvc/gmod/cmap/trunk/docs/ADMINISTRATION.pod |
| Example | http://www.gramene.org/cmap/ |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-cmap |

# Agenda

| | |
|---|---|
| 3:00 | Introduction |
| 3:10 | Software<br><br>    Visualization<br>        GBrowse<br>        JBrowse<br>        GBrowse_syn<br>        Sybil & SynView<br>        CMap |

| | |
|---|---|
| 4:00 | Software, cont.<br><br>    Data Management<br>        Chado, Tripal, GMODWeb<br>        BioMart and InterMine<br>        GFF3<br>    Annotation<br>        MAKER & DIYA<br>        Apollo<br>        Textpresso<br>        Community Annotation<br>        Pipelines & Workflows |
| 4:40 | Community |
| 5:00 | Finish |

# Chado: A database schema for biological data

- ## A *schema* is a database design
  - Blueprint for a database, a way of organizing data

- ## Independent of specific data
  - Chado provides structure
  - You provide the hard work and data

# Why use Chado?

- Very good at genomic data
- Widely used
  - AphidBase, BeetleBase, dictyBase, FlyBase, SGN, SpBase, VectorBase, wFleaBase, …
- Integrates with other GMOD tools
- Community of support
- Modular, flexible and extensible

# Chado Modules

# CVs and Ontologies in Chado

- **Controlled vocabularies and ontologies are key in Chado**
- **Maximally used for**
  - Integrity
  - Interoperability
- **Can create your own, *but* …**
  - Please use standard ontologies when they exist
  - See OBO: http://www.obofoundry.org/

# Chado Future Developments

Flexibility means core schema changes *slowly*

    That's a feature.

- Natural Diversity module
  - Better support for phenotypes, crosses, individuals, geolocation, …
  - Based on GDPDM from Cornell University, Terry Casstevens, *et al.* (http://www.maizegenetics.net/gdpdm/)

- Expression / Anatomy / Cell Fate Atlas support
  - Aniseed (http://aniseed-ibdm.univ-mrs.fr/) converting to Chado and extending it to better support atlases
  - Will have a web front end for atlases

# Chado Resources

| | |
|---|---|
| Home Page | http://gmod.org/wiki/Chado |
| Tutorial | http://gmod.org/wiki/Chado_Tutorial |
| Introduction | http://gmod.org/wiki/Introduction_to_Chado |
| Manual | http://gmod.org/wiki/Chado_Manual |
| Modules | http://gmod.org/wiki/GBrowse_Modules |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-schema |

# Chado Web Front Ends

- Chado is a schema, a server side technology
- It is not a web front end or a desktop client

- Options for Chado web front ends:
  - Do it yourself
  - GMODWeb
  - Tripal

# Do it yourself

- GMOD provides some support in form of libraries
- Perl
  - Chado::AutoDBI
  - Modware → Bio:Chado:Schema
- Java
  - Two projects under development

# GMODWeb

- A Chado specific set of templates for the generic Turnkey web site generation system

- Written in Perl

- Lots of Perl module dependencies



ParameciumDB, a website built with GMODWeb
http://paramecium.cgm.cnrs-gif.fr/

GMODWeb: a web framework for the generic model organism database, O'Connor *et al*., Genome Biology 2008, 9:R102.

# Tripal

- Added to GMOD this year
- Set of Drupal modules
  - Feature, Organism, Library, Analysis
  - Modules roughly correspond to Chado modules
  - Easy to create new modules
- Includes user authentication, job management, and data entry support
- Developed by Clemson University Genomics Institute



MarineGenomics.org

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

# Tripal Resources

| | |
|---|---|
| Home Page | http://gmod.org/wiki/Tripal |
| Tutorial | http://gmod.org/wiki/Tripal_Tutorial |
| User Guide | http://gmod.org/wiki/Media:TripalUsersGuideJune2009.pdf |
| Example | http://marinegenomics.org |
| Mailing List | https://lists.sourceforge.net/lists/listinfo/gmod-tripal |

# Chado Web: DIY or GMODWeb or Tripal?

**GMODWeb**
Complete
Requires some tuning
Perl

**Tripal**
User authentication
Data entry
Actively developed
Well documented
Easy to extend
Drupal

**Do It Yourself**
More work
Get exactly what you want

What really made us decide to switch over to Drupal was that we needed authentication mechanisms, customized data entry mechanisms, and the ability to add social networking features and other non-biological components to our sites. Drupal supported all of this and was widely used, well documented, and well supported.

Stephen Ficklin, CUGI

# BioMart and InterMine

- Chado well-suited for setting up organism databases that have
  - Easy to use query interface to support common types of questions
  - Unified, coherent presentation of information
- BioMart and InterMine
  - Allow users to ask complex queries on all data
  - At the expense of having to do more work

# GFF3

- The common file format of GMOD for genomic annotation

- Supported by Chado, GBrowse, JBrowse, CMap, Apollo, ....

# Agenda

| | |
|---|---|
| 3:00 | Introduction |
| 3:10 | Software<br><br>    Visualization<br>        GBrowse<br>        JBrowse<br>        GBrowse_syn<br>        Sybil & SynView<br>        CMap |

| | |
|---|---|
| 4:00 | Software, cont.<br><br>    Data Management<br>        Chado, Tripal, GMODWeb<br>        BioMart and InterMine<br>        GFF3<br>    Annotation<br>        MAKER & DIYA<br>        Apollo<br>        Textpresso<br>        Community Annotation<br>        Pipelines & Workflows |
| 4:40 | Community |
| 5:00 | Finish |

# MAKER

- Genome annotation pipeline for creating gene predictions
- 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 Online in beta

MAKER: An easy-to-use annotation pipeline designed for emerging model organism genomes, Brandi L. Cantarel, *et al.*, *Genome Res*. 2008. 18: 188-196

# MAKER Resources

| | |
|---|---|
| Home Page | http://www.yandell-lab.org/software/maker.html |
| Tutorial | http://gmod.org/wiki/MAKER_Tutorial |
| PAG Workshop | http://gmod.org/wiki/MAKER_PAG_2010_Workshop |
| Mailing List | http://yandell-lab.org/mailman/listinfo/maker-devel_yandell-lab.org |

# DIYA

- Lightweight, modular, and configurable Perl-based pipeline framework.

- Initial application is gene prediction pipeline for prokaryotes

- Working on integration of Amos assembly tools.

# Ergatis

- Web interface to the TIGR-Workflow engine
- Create, run and monitor reusable computational analysis pipelines
- Manage compute clusters or single machines
- Comes with several pre-configured pipelines

# Galaxy

- ## Web portal
  - Search remote resources, combine data from independent queries and visualize results

- ## Queries / pipelines can be saved and referenced in papers or rerun later.

- ## Supports set-theory operations on results

- ## Links to outside tools, including GBrowse

- ## Can use central server or install locally

**Galaxy**

**Tools**

Get Data
Send Data
ENCODE Tools
Lift-Over
Text Manipulation
Convert Formats
FASTA manipulation
Filter and Sort
Join, Subtract and Group
Extract Features
Fetch Sequences
Fetch Alignments
Get Genomic Scores
Operate on Genomic Intervals
Statistics
Graph/Display Data
Regional Variation
Multiple regression
Evolution
Metagenomic analyses
EMBOSS

NGS TOOLBOX BETA

NGS: QC and manipulation
NGS: Mapping
NGS: SAM Tools

GM D

# Apollo

- GMOD's genome annotation editor
- Add and refine annotations.
- Java desktop client
- Widely used
- Read/write in multiple formats
- Keep track of evidence, curator
- Used in several community annotation efforts

# Apollo Future Work

- Berkeley Bioinformatics Open-source Projects (BBOP)
  - Current developers of Apollo
  - Submitted a grant proposal for
    - Apollo on the web
    - Using same underlying tools as JBrowse
- Meanwhile, CCG/ABF
  - Is using Apollo (and Chado) for genome annotation
  - ABF is exploring the possibility of developing a web-based application to complement Apollo
  - NCRIS 5.1 funding for a 6 month project
- These two groups are talking to each other

# Apollo Resources

| | |
|---|---|
| Home Page | http://apollo.berkeleybop.org/ |
| Tutorial | http://gmod.org/wiki/Apollo_Tutorial |
| User Guide | http://apollo.berkeleybop.org/current/userguide.html |
| Mailing List | http://mail.fruitfly.org/mailman/listinfo/apollo |

# Textpresso

- **Text mining system for scientific papers**
- **Analyzes full article text**
- **Indexes articles by keywords and by category tags**
- **Stand alone search engine w/ web interface**
- **Curation tool**



Textpresso for *E. coli*, extensions by Nathan Liles, Hu Lab

Textpresso: an ontology-based information retrieval and extraction system for biological literature, Muller HM, Kenny EE, Sternberg PW, *PLoS Biol.* 2004 Nov;2(11):e309

# Community Annotation

- How do you get others to contribute?

- Social:
  - Sticks
    - Work well if your database is already the authority on your topic/organism and you have curators and a huge community
  - Carrots
    - Give people credit
    - Give people ownership
    - Seek mutually beneficial relationships
  - Comfort Level
    - Recent popularity of social computing

# Community Annotation

- Technological
  - Make it easier to fix something then it is to be irritated by its error or absence.
    - The Wikipedia model.
  - Make it relatively easy for people who really care to contribute significant content
    - Also the Wikipedia model.

# Community Annotation: GMOD Technology

- ## Apollo
  - Several projects use Apollo to distribute genome annotation efforts
  - Apollo infrastructure supports:
    - Read from Chado → Save to XML → Review → Upload to Chado
  - But
    - Java application; Infrequent Apollo users forget a lot.
    - Web Apollo will help some, maybe a lot

- ## Tripal
  - Supports update interfaces for data in Chado databases.
  - Has access to all of Drupal's social networking.

# Community Annotation: GMOD Technology

- Table Edit
  - A MediaWiki extension that provides a GUI interface to updating MediaWiki tables.
  - MediaWiki software used at Wikipedia
  - Has been extended to update and render database tables through a MediaWiki interface.
  - Work is in progress to apply it to Chado.
  - See http://ecoliwiki.net
  - Has potential to turn Chado into a wiki.

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

# Agenda

| 3:00 | Introduction |
|------|--------------|
| 3:10 | Software |
|      | Visualization |
|      | GBrowse |
|      | JBrowse |
|      | GBrowse_syn |
|      | Sybil & SynView |
|      | CMap |

| 4:00 | Software, cont. |
|------|-----------------|
|      | Data Management |
|      | Chado, Tripal, GMODWeb |
|      | BioMart and InterMine |
|      | GFF3 |
|      | Annotation |
|      | MAKER & DIYA |
|      | Apollo |
|      | Textpresso |
|      | Community Annotation |
|      | Pipelines & Workflows |
| 4:40 | Community |
| 5:00 | Finish |

# Who uses GMOD?

Plus hundreds of others

# GMOD Project

- Open Source
- Two full time project staff:
  - Project Coordinator: Scott Cain
  - Help Desk: Dave Clements
- Components
  - Some have dedicated funding
  - Others are contributed
  - New components must have:
    - An open source license
    - Interoperability with other GMOD components
    - A good faith commitment of at lest 2 years of support

# GMOD.org

A wiki, of course. GMOD.org is the hub for all things related to the project:

- – Documentation
- – News
- – Links
- – Calendar
- – Tutorials
- – HOWTOs
- – Glossary
- – Overview
- – …

# Mailing Lists

- Several project lists

- Many component-specific lists

- 3100 messages in last 12 months on the 7 lists managed by GMOD staff

- Up 69% from previous year

- Mailing lists are very active



http://gmod.org/wiki/GMOD_Mailing_Lists

# Meetings, Training and Outreach

- **Semi-annual community meetings**
  - Next Meeting:
    - January 2010, San Diego, after PAG
- **GMOD Summer Schools**
  - 2009
    - July, NESCent, North Carolina, US
    - August, Oxford, UK
  - 2010
    - ??, NESCent, North Carolina, US
    - ??, Asia / Pacific, maybe
- **Outreach**
  - BA, SMBE, PAG, Arthropod Genomics, …

http://gmod.org/wiki/Training_and_Outreach

# Tutorials

- **Summer school sessions become online tutorials with**
  - Starting VMware images
  - Step by step instructions
  - Example datasets
  - Ending VMware images
- **Topics:**
  - Apollo, Artemis-Chado Integration, BioMart, Chado, CMap, GBrowse, GBrowse_syn, JBrowse, MAKER, Tripal, GBrowse NGS



http://gmod.org/wiki/Training_and_Outreach#Online_Tutorials

# Agenda

| | |
|---|---|
| 3:00 | Introduction |
| 3:10 | Software<br><br>Visualization<br>　GBrowse<br>　JBrowse<br>　GBrowse_syn<br>　Sybil & SynView<br>　CMap |

| | |
|---|---|
| 4:00 | Software, cont.<br><br>Data Management<br>　Chado, Tripal, GMODWeb<br>　BioMart and InterMine<br>　GFF3<br>Annotation<br>　MAKER & DIYA<br>　Apollo<br>　Textpresso<br>　Community Annotation<br>　Pipelines & Workflows |
| 4:40 | Community |
| 5:00 | Finish |

# Acknowledgements

**NESCent**
Todd Vision
Hilmar Lapp

**BBOP**
Ed Lee

**OICR**
Scott Cain
Lincoln Stein

**Oregon**
Patrick Phillips
Phillips Lab

**Texas A&M**
Jim Hu
Daniel Renfro
Nathan Liles
Brenley McIntosh

**CBRG, Oxford**
Simon McGowan

**Broad**
Heng Li

**CUGI**
Stephen Ficklin

**CSHL**
Sheldon McKay
Ken Youens-Clark

# Thank You!



Dave Clements
GMOD Help Desk

US National Evolutionary
     Synthesis Center
     http://nescent.org

clements@nescent.org
help@gmod.org

http://gmod.org/wiki/GMOD_Help_Desk

# SAMtools

## Introduction

SAM (Sequence Alignment/Map) format is a generic format for storing large nucleotide sequence alignments. SAM aims to be a format that:

- Is flexible enough to store all the alignment information generated by various alignment programs;
- Is simple enough to be easily generated by alignment programs or converted from existing alignment formats;
- Is compact in file size;
- Allows most of operations on the alignment to work on a stream without loading the whole alignment into memory;
- Allows the file to be indexed by genomic position to efficiently retrieve all reads aligning to a locus.

SAM Tools provide various utilities for manipulating alignments in the SAM format, including sorting, merging, indexing and generating alignments in a per-position format.

SAMtools is hosted by SourceForge.net. The project page is here. The source codes are available from the download page. You can check out the latest source codes with:

```
svn co https://samtools.svn.sourceforge.net/svnroot/samtools/trunk
/samtools
```

### General Information

SAM Format Specification
SF Project Page
SF Download Page
Mailing Lists
SVN Browse
Related Software

### SAMtools in C

General Introduction
Manual Page
Pileup Format
Consensus/Indel Calling
Text Alignment Viewer
API Documentation
Example C Program
Working on a Stream
Open Tasks

### Variant Call Format

### Other Lang-bindings

Picard (Java)
Bio-SamTools (Perl)

## Platform neutral set of programs and file formats specifically for short reads.
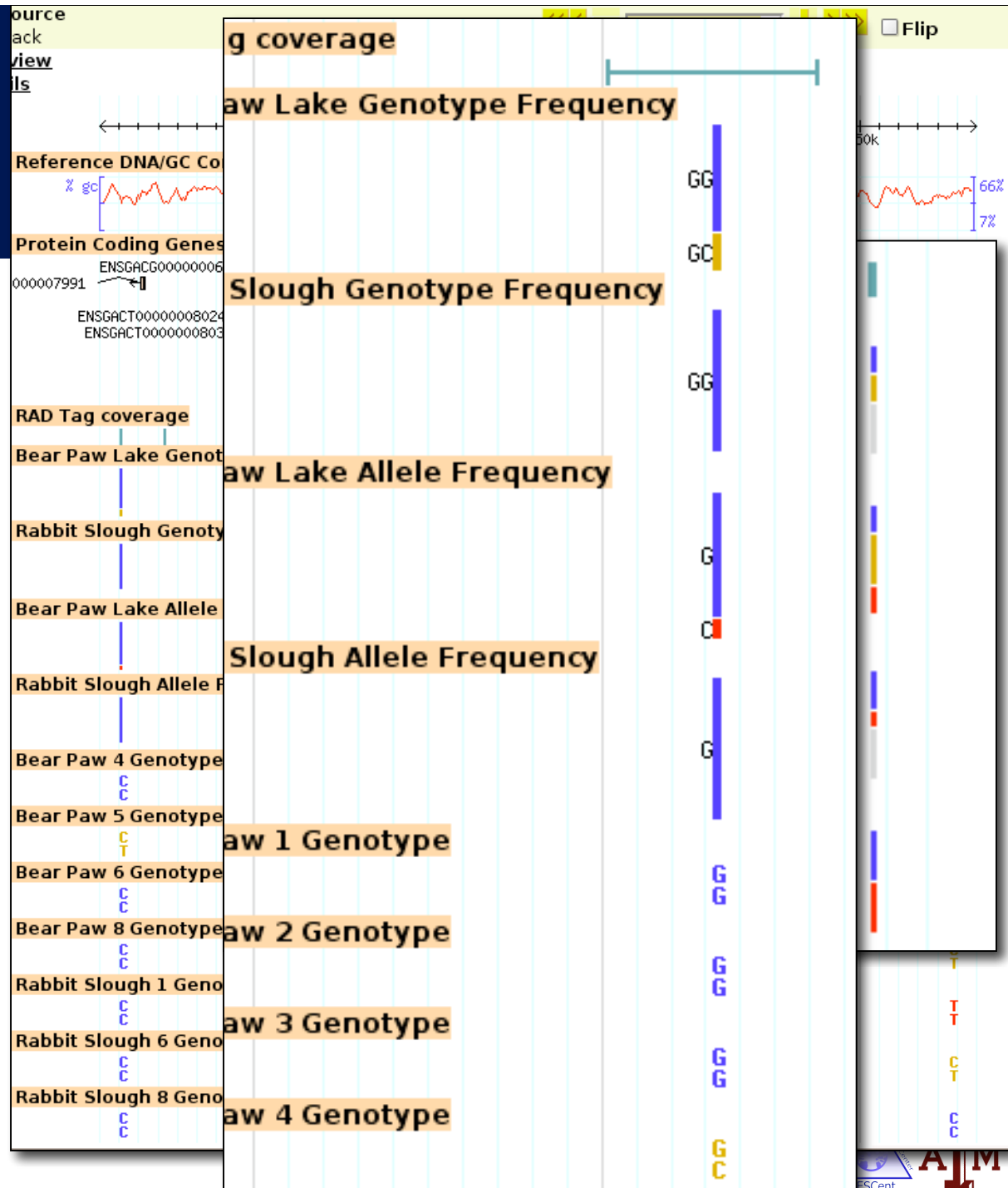
# GBrowse for Population Genetics

Shows
- Where we looked
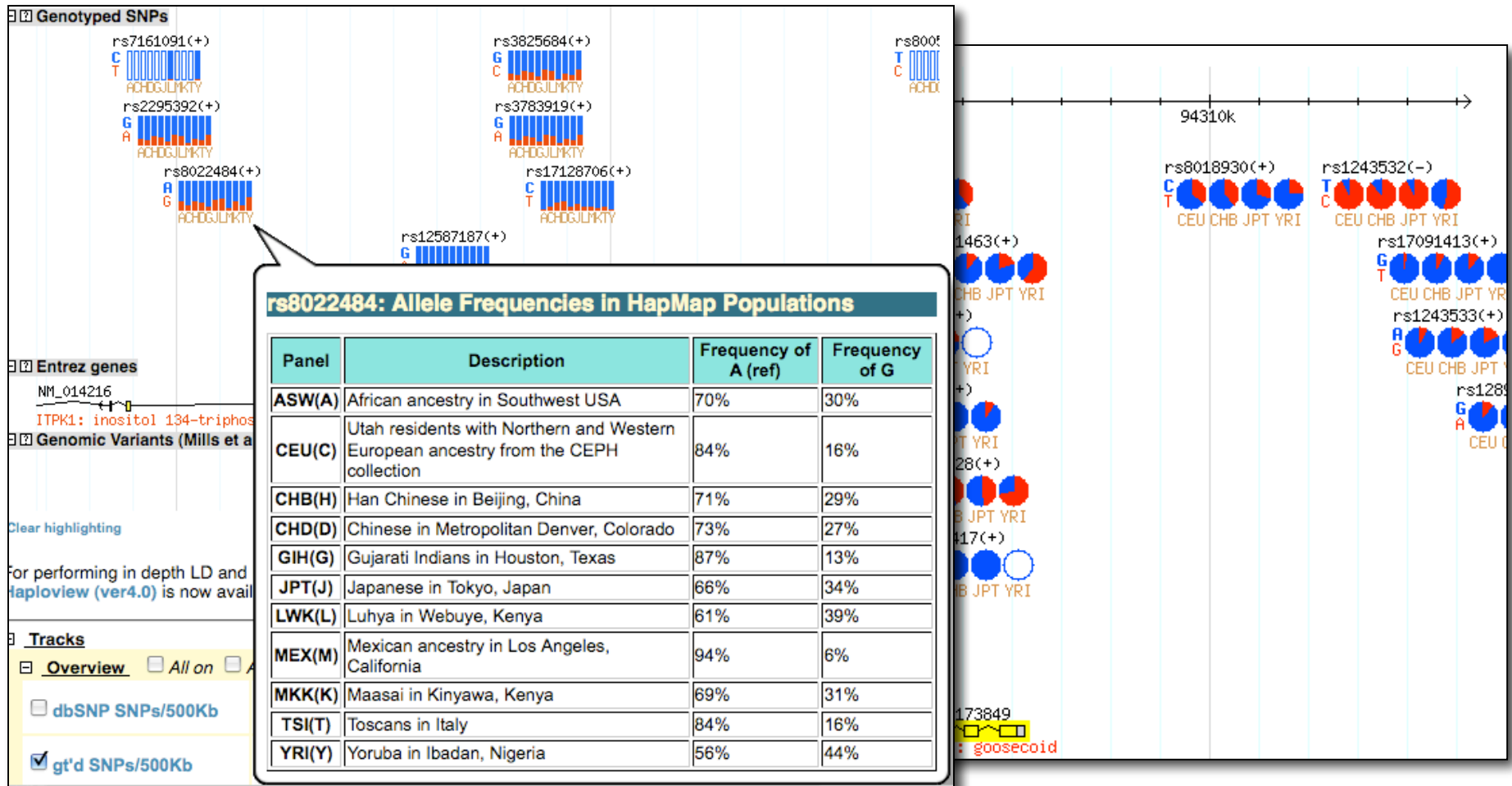- Allele & genotype frequencies
  - By population
  - Individual genotypes

Could also show:
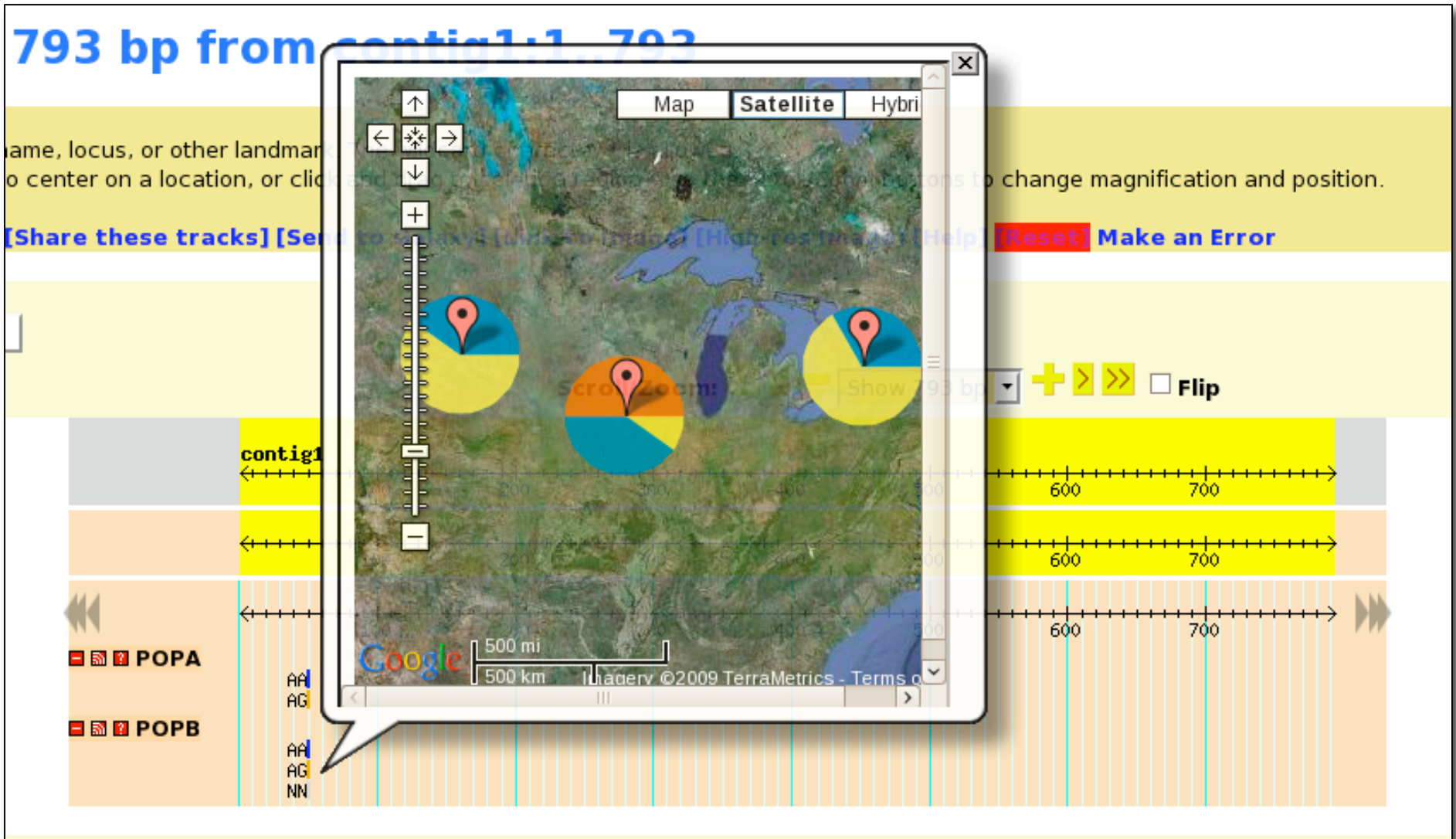- Frequency by phenotype or any other characteristic
- Sliding window stats

# HapMap Allele Frequencies



http://hapmap.org

# Geolocation data



Yi-Hsin Erica Tsai Ben Faga