# Comparative Genomics with GBrowse_syn

**Sheldon McKay,**
**Cold Spring Harbor Laboratory**
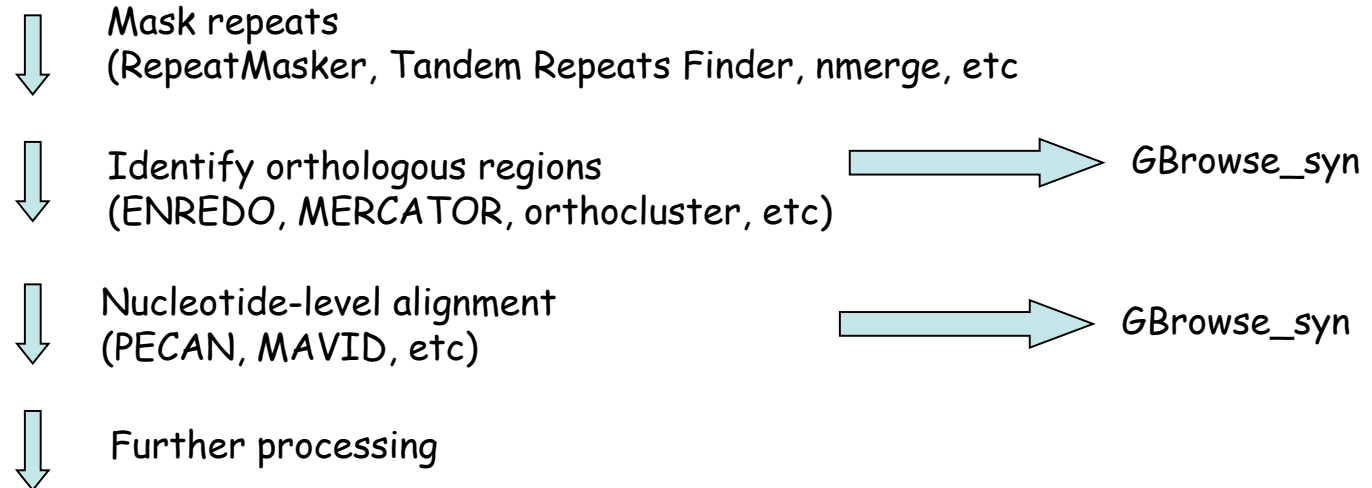
## Outline

A few words on data

A brief survey of synteny browsers

A few challenges of rendering comparative data

Comparative genome browsing with GBrowse_syn

# Hierarchical Genome Alignment Strategy

Raw genomic sequences

⬇ Mask repeats
(RepeatMasker, Tandem Repeats Finder, nmerge, etc

⬇ Identify orthologous regions
(ENREDO, MERCATOR, orthocluster, etc)  ➡ GBrowse_syn

⬇ Nucleotide-level alignment
(PECAN, MAVID, etc)  ➡ GBrowse_syn

⬇ Further processing

GBrowse

# A Few Use Cases

- Multiple sequence alignment data from whole genome

- Synteny or Co-linearity data without alignments

- Gene orthology assignments based on proteins

- Self vs. Self comparison of duplications, homeologous regions, etc

- Others

**What is a Synteny Browser?**

- Has display elements in common with genome browsers

- Uses sequence alignments, orthology or co-linearity data, to highlight different genomes, strains, etc.

- Usually displays co-linearity relative to a reference genome.

# An Embarrassment of Riches*

## A Brief Survey of GMOD-friendly Synteny Browsers

*From John Ozell's 1738 translation of a French play, L'Embarras des richesses (1726)

SynView: A Simple Approach to Visualizing Comparative Genome Data

Wang H, Su Y, Mackey AJ, Kraemer ET and JC Kissinger . SynView: a GBrowse-compatible approach to visualizing comparative genome data  Bioinformatics 2006 22:2308-2309
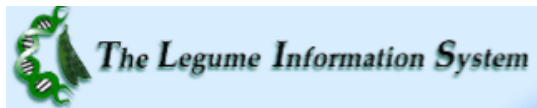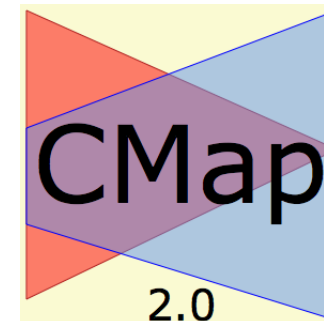
Pan, X., Stein, L. and Brendel, V. 2005. SynBrowse: a Synteny Browser for Comparative Sequence Analysis. Bioinformatics 21: 3461-3468

# Sybil: Web-based software for comparative genomics

Crabtree, J., Angiuoli, S. V., Wortman, J. R., White, O. R. Sybil: methods and software for multiple genome comparison and visualization Methods Mol Biol. 2007 Jan 01; 408: 93-108.

Youens-Clark K, Faga B, Yap IV, Stein LD, Ware, D. 2009.
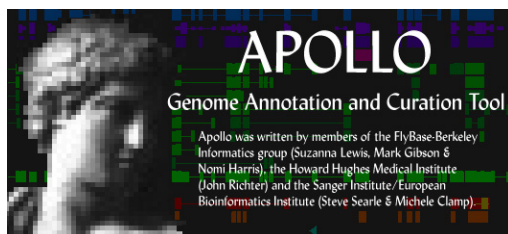CMap 1.01: A comparative mapping application for the Internet.  doi:10.1093

+ others...

# GBrowse_syn

Branding ideas...

GBrowse_syn

GBROWSE_SYN
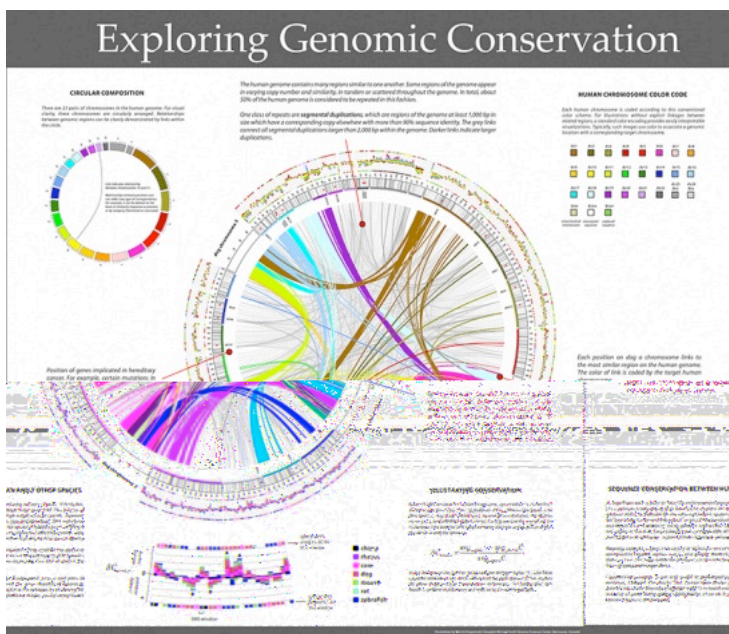
When one genome just isn't enough...

Desktop Synteny Viewers: Apollo and Artemis

# Debating the relative merits of Apollo* and Artemis‡

# Other non-GMOD Browsers



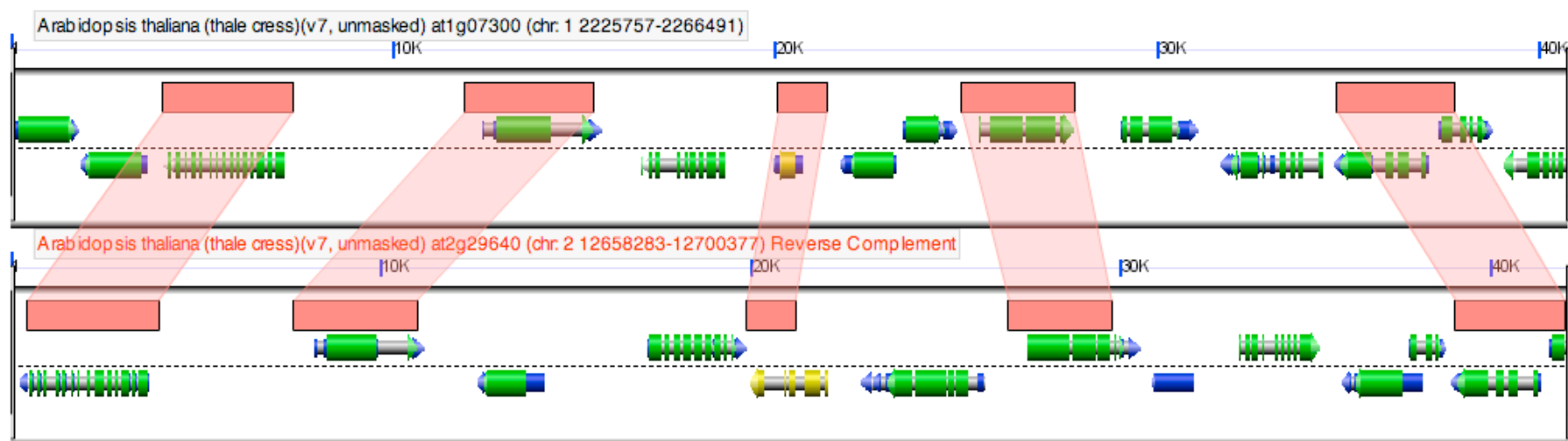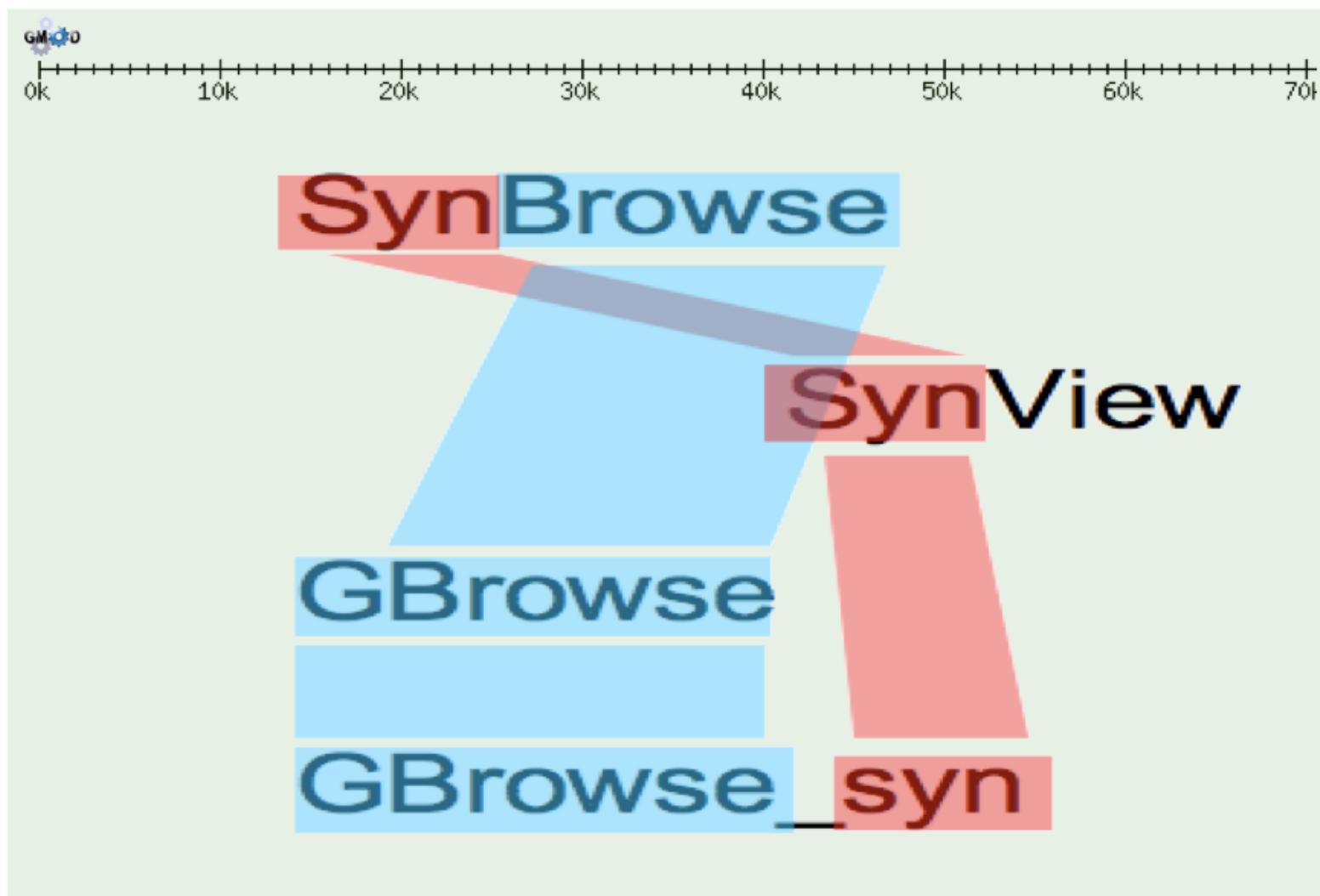http://mkweb.bcgsc.ca/circos/



http://www.mizbee.org

# Other non-GMOD Browsers



http://synteny.cnr.berkeley.edu/CoGe/

## GMOD Browser branding/nomenclature issues...

**SynView:**
- Add-on to native GBrowse package
- Uses GFF3 or DAS1 compliant data adapters
- GFF requires special tags (allowed in spec.)
- Reference panel on top

**SynBrowse:**
- Uses same core libraries as Gbrowse
- Uses GFF database adapter
- GFF2 uses standard 'Target' syntax
- Currently only supports two species
- Central reference panel?

**Sybil:**

- Not GBrowse-based
- Uses chado database
- Whole genome and detailed views

**GBrowse_syn:**

- Part of GBrowse distribution
- Uses native GFF2/3 or chado adapters for species' data
- Synteny data are stored in a separate joining database

# How is GBrowse_syn different?

- Does not rely on perfect co-linearity across the entire displayed region (no orphan alignments)
- Offers on the fly alignment chaining
- No upward limit on the number of species
- Used grid lines to trace fine-scale sequence gain/loss
- Seamless integration with GBrowse data sources
- Ongoing support and development
- Some people think it looks nice

# GBrowse-like interface



**PECAN alignments for _Caenorhabditis_ (WS197)**

■ **Instructions**

Select a Region to Browse and a Reference species:

**Examples**: c_elegans X:1050001..1150000, c_briggsae chrX:620000..670000, c_elegans R193.2.

■ **Search**

**Landmark**:
[X:1050001..1150000]  [Search]  [Reset]

**Reference Species**:
[C. elegans ▾]

[«« ‹ ━ Show 100 kbp ▾ ✛ › ››]

**Aligned Species**:
☑C. briggsae ☑C. remanei ☑C. brenneri ☑C. japonica

**Data Source :**
[PECAN alignments for Caenorhabditis ▾]

**Display Mode**:
Three species/panel Click to show all species in one panel

■ **Overview**

Reference genome: _C. elegans_

X

0M  1M  2M  3M  4M  5M  6M  7M  8M  9M  10M  11M  12M  13M  14M  15M  16M  17M

GBrowse_syn

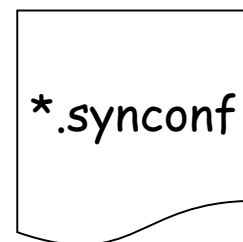GBrowse
Databases*

GBrowse_syn
alignment
database

*.syn
or
*.conf

*.synconf

Species config.

Master config.

# GBrowse_syn Architecture

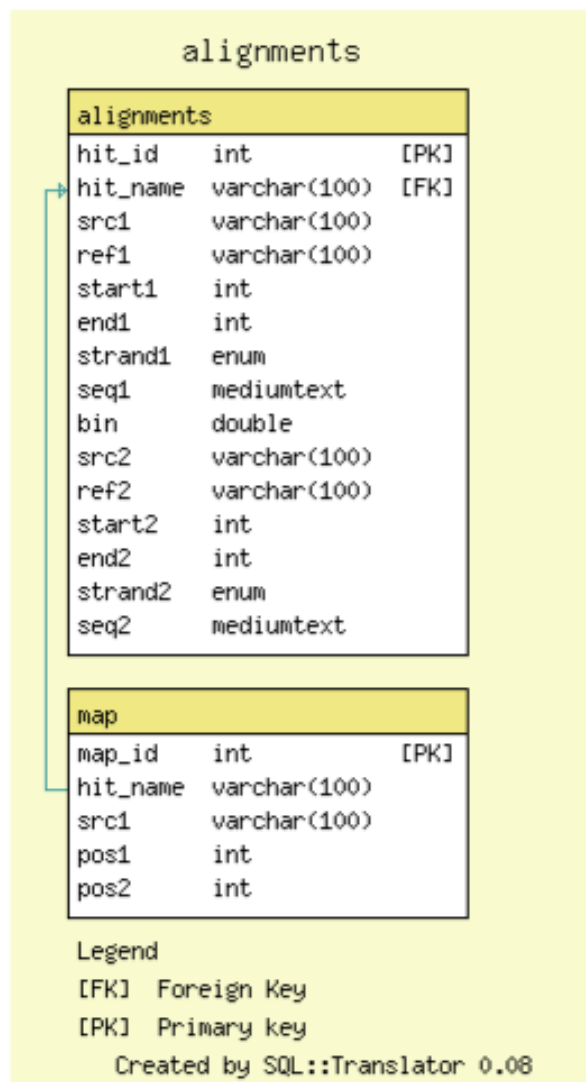[GBrowse]                                                          [GBrowse]

**Bio::DB::GFF**
**species1**  ⟷

**Bio::DB::GFF**
**species3**  ⟷

alignments

| alignments | | |
|---|---|---|
| hit_id | int | [PK] |
| hit_name | varchar(100) | [FK] |
| src1 | varchar(100) | |
| ref1 | varchar(100) | |
| start1 | int | |
| end1 | int | |
| strand1 | enum | |
| seq1 | mediumtext | |
| bin | double | |
| src2 | varchar(100) | |
| ref2 | varchar(100) | |
| start2 | int | |
| end2 | int | |
| strand2 | enum | |
| seq2 | mediumtext | |

| map | | |
|---|---|---|
| map_id | int | [PK] |
| hit_name | varchar(100) | |
| src1 | varchar(100) | |
| pos1 | int | |
| pos2 | int | |

Legend
[FK]  Foreign Key
[PK]  Primary key
     Created by SQL::Translator 0.08

**Bio::DB::GFF**
**species2**  ⟷

**Bio::DB::GFF**
**species4**  ⟷

[GBrowse]                                                          [GBrowse]

# Getting Data into GBrowse_syn

CLUSTALW    PECAN
MSF     *ad hoc* tab-delimited
FASTA    STOCKHOLM
  GFF3   etc…

Loading scripts

# Gbrowse_syn: quick tour

# Gbrowse_syn: quick tour (shaded alignments)

Gbrowse_syn: quick tour (strand correction)

# Optional "All in one" view

# Adding markup to the annotations

# Problem : How to use Insertions/Deletion data

# Tracking Indels with grid lines
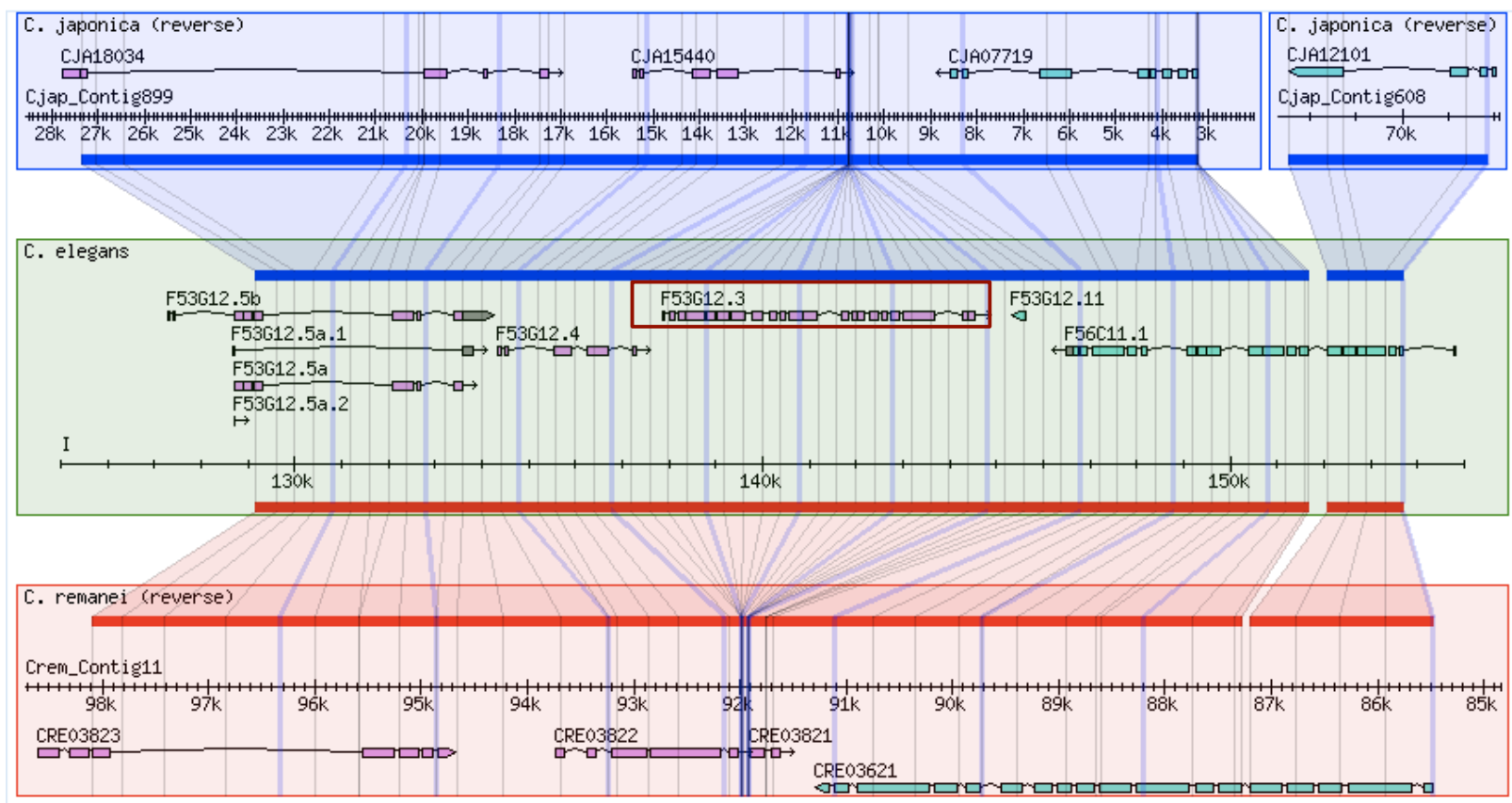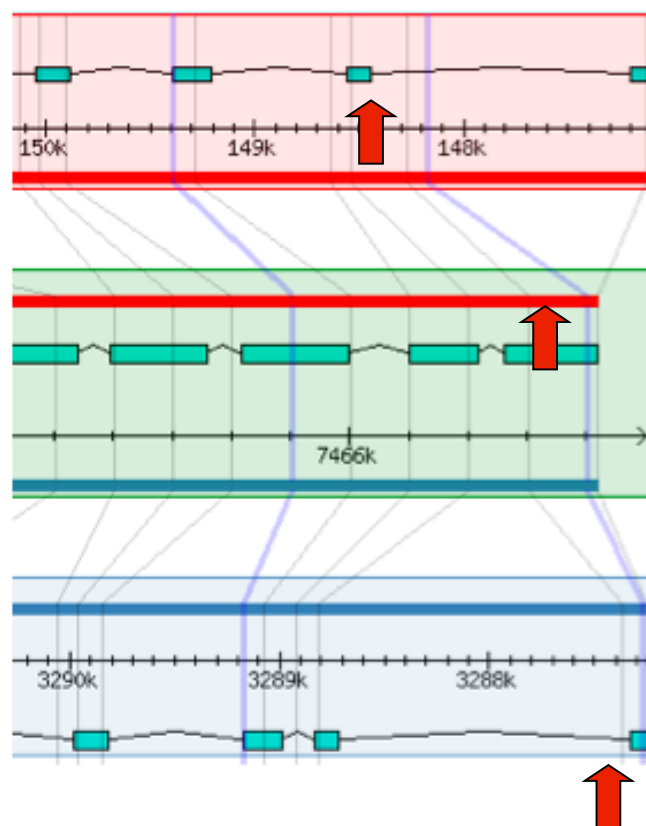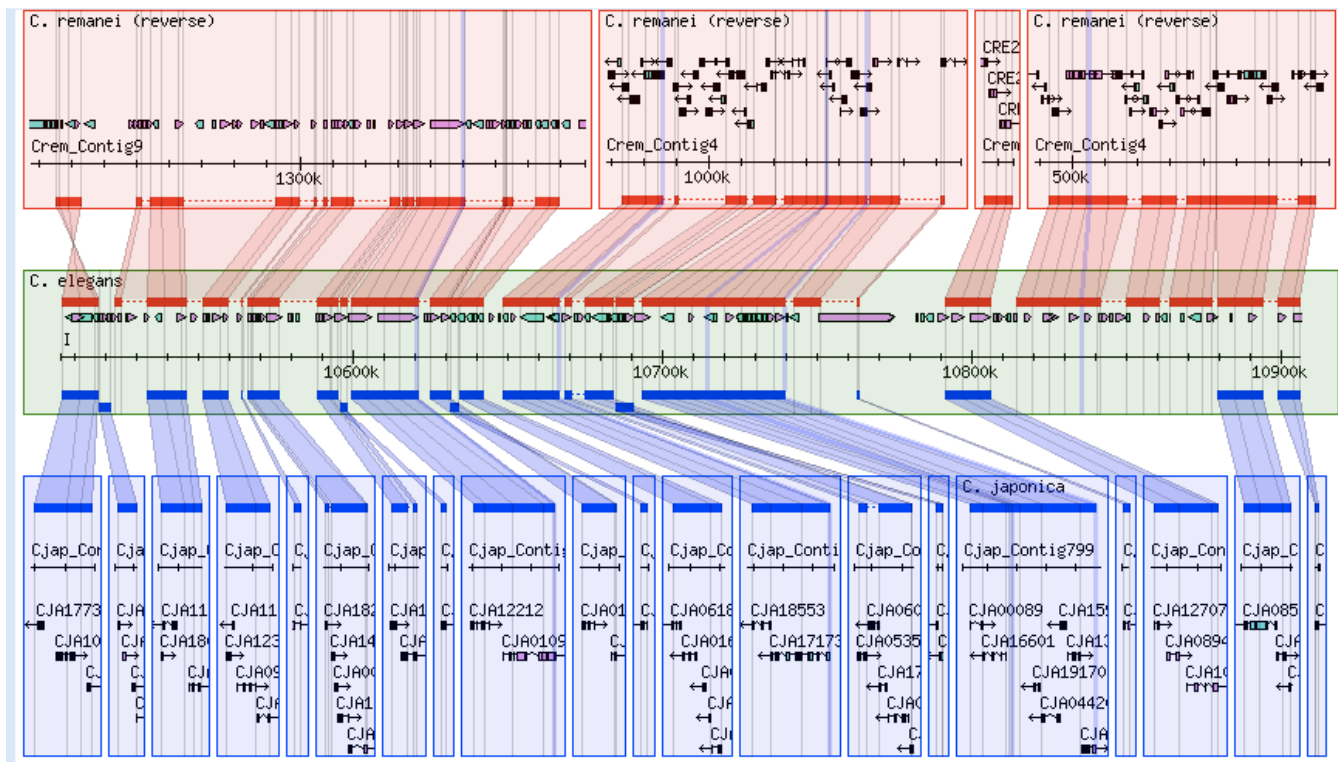
# Evolution of Gene Structure

# Putative gene or loss

# Comparing gene models
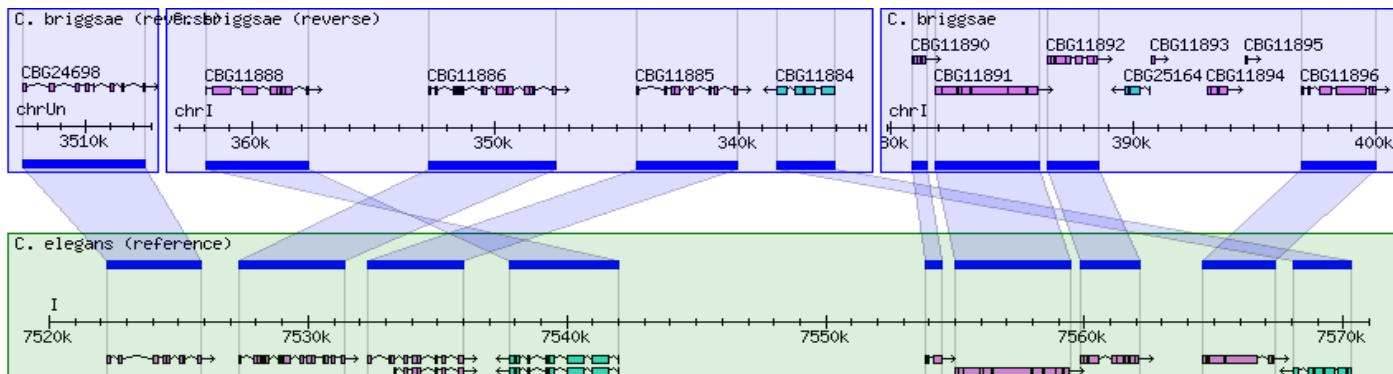
# Comparing assemblies
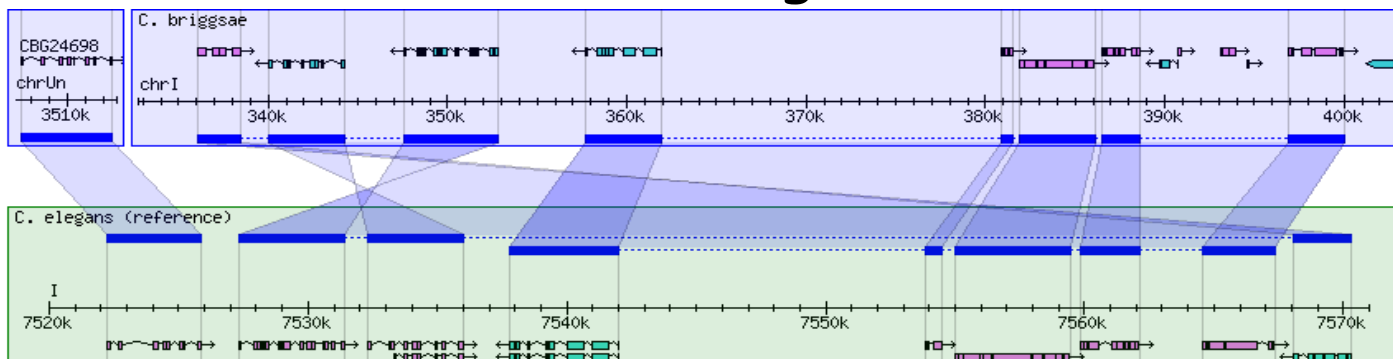


Not bad

Needs work

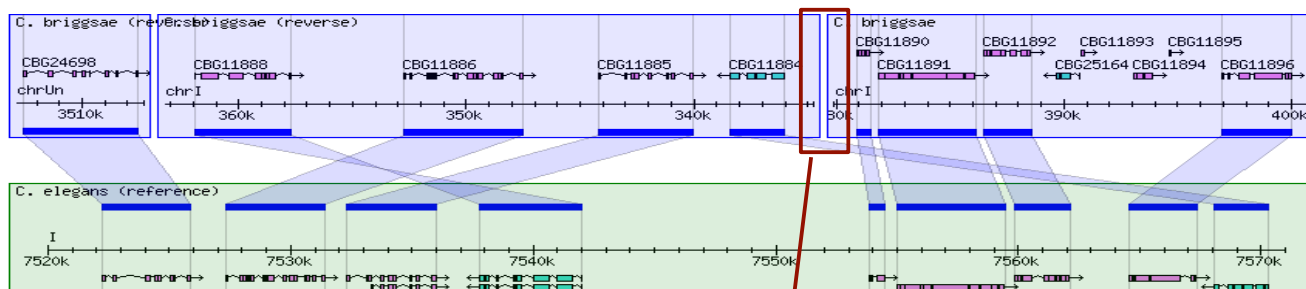# Getting the most out of small aligned regions or orthology-only data
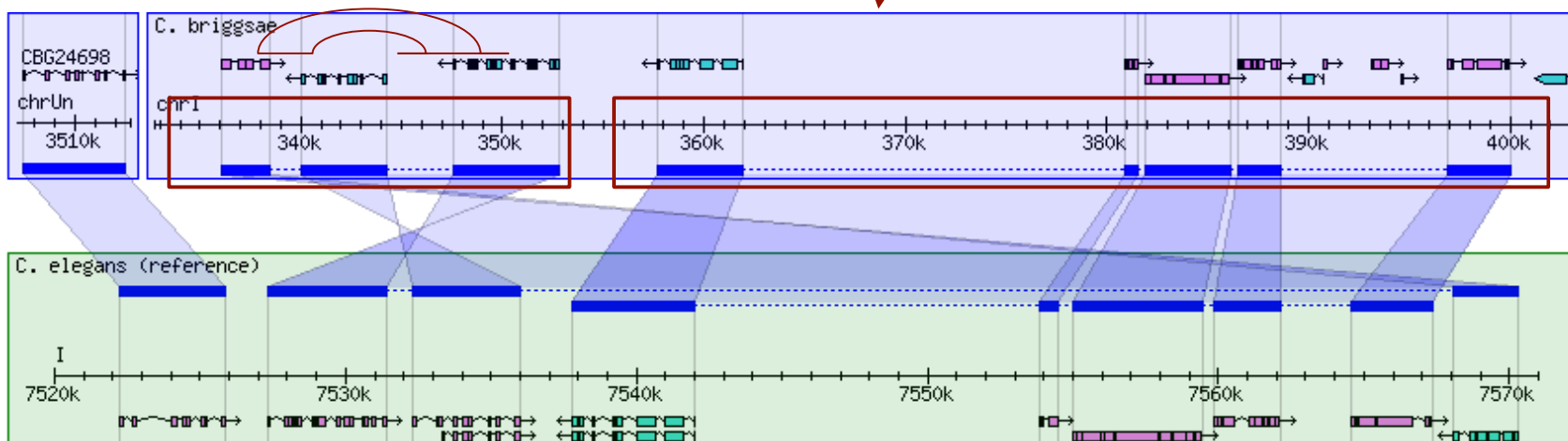
# Gene Orthology



# Chained Orthologs

Inversion + translocation?

2 panels merged

# What about synteny blocks that fall off the ends of the displayed reference sequence?
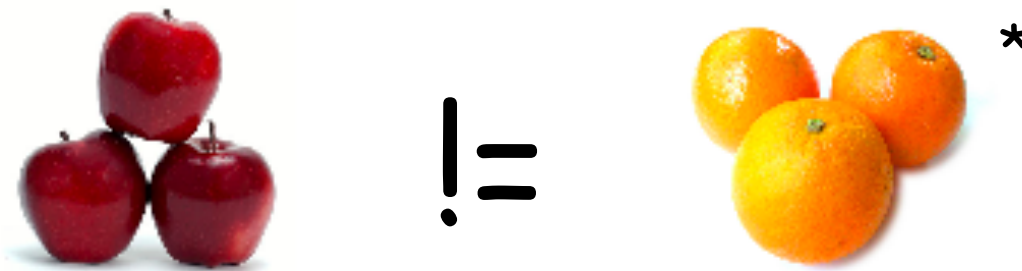
Solution 1 : With multiple sequence alignment data,
          calculate many anchor points (done anyway
          for grid lines)


Solution 2 : For orthology-based synteny blocks, use
          individual start and end coordinates of orthologs
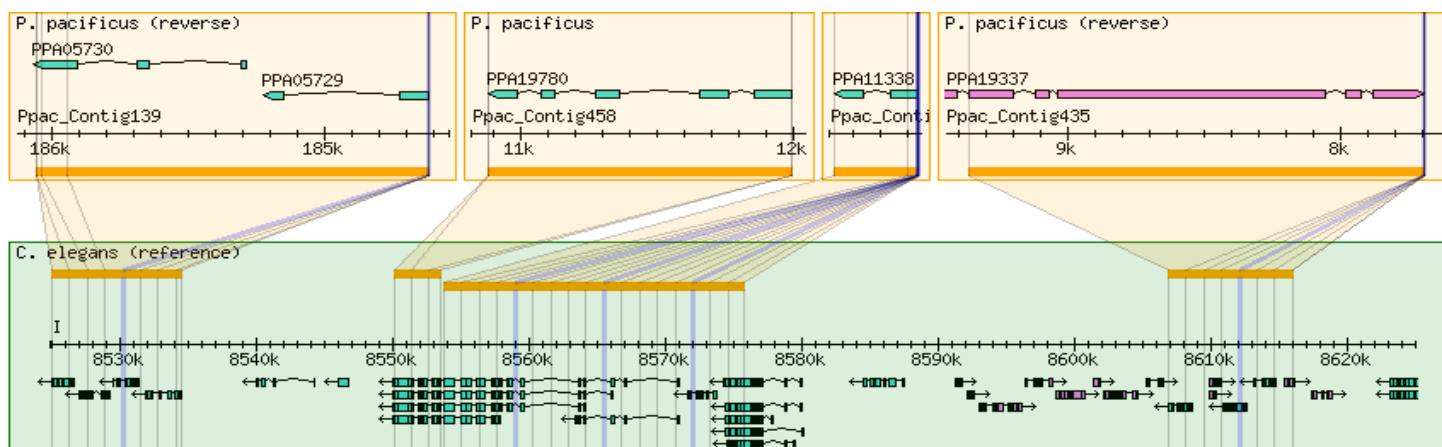          as anchor points.


Solution 3: If all else fails, guess the end of the target block
          based on the overall length ratio.

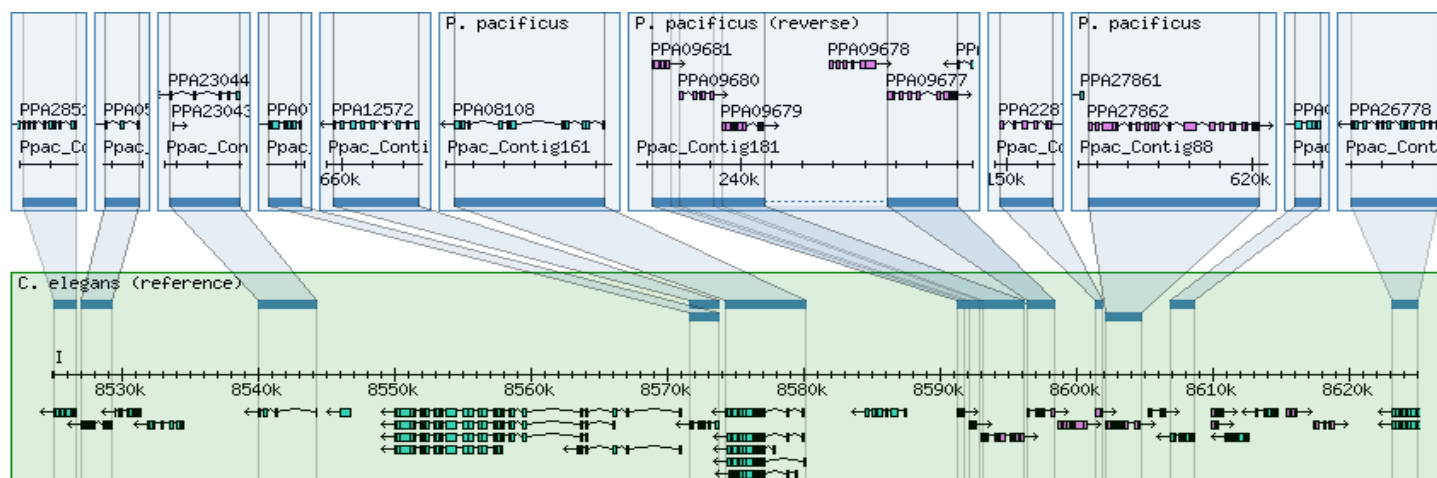length displayed target = (length target/length reference)* length displayed reference
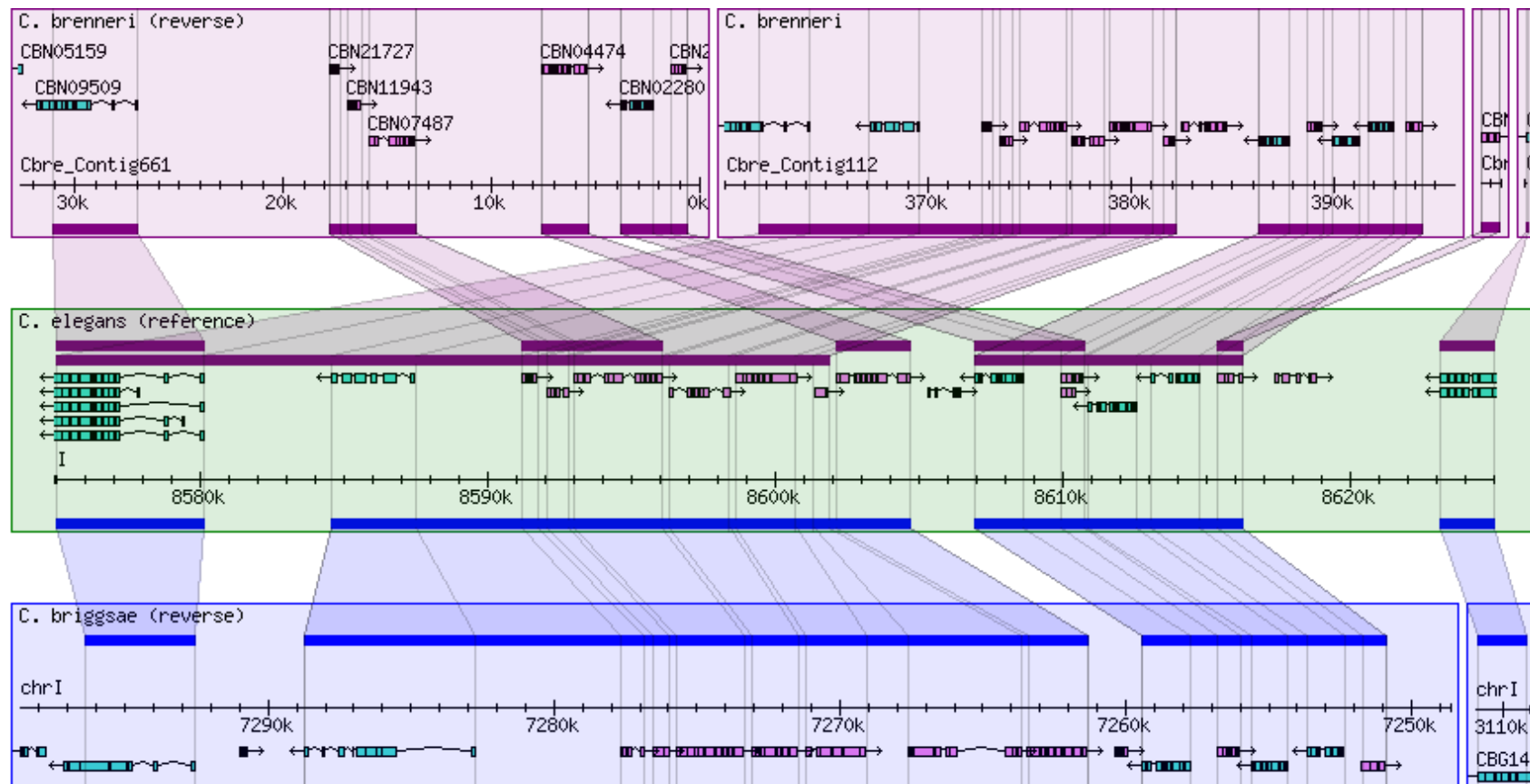
# Pecan alignments



# Protein orthology based Synteny blocks

# What about segmental duplications?

# The Future of GBrowse_syn*

- Integration with GBrowse 2.0

- "On the fly" sequence alignment view

- AJAX-based user interface and navigation (Jbrowse_syn)

- Suggestions?

*GBrowse_syn has a future

# Acknowledgments

Lincoln Stein
Dave Clements
Scott Cain
Jason Stajich
Bonnie Hurwitz
Eva Huala
Cynthia Lee
Jack Chen
Ismael Verga
Michael Han
WormBase Curators

Projects

Funding