

Browsing Genomic Information with Ensembl

Stefan Wyder

April 2016



Objectives

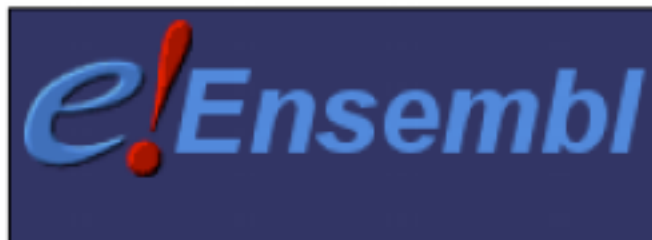
- What is **Ensembl**?
- How to navigate the **Ensembl** browser
- What type of data can you get in **Ensembl**?
- How to 'data mine' **Ensembl**
- Where to go for **help** and **documentation**

Why do we need genome browsers?

CGGCCTTTGGGCTCCGCCTTCAGCTCAAGACTTAACTTCCCTCCCAGCTGTCCCAGATGACGCCATCTGAAATTTCTTGGA
ACACGATCACTTTAACGGAATATTGCTGTTTTGGGGAAGTGTTTTACAGCTGCTGGGCACGCTGTATTTGCCTTACTTAAGC
CCCTGGTAATTGCTGTATTCCGAAGACATGCTGATGGGAATTACCAGGCGGCGTTGGTCTCTAACTGGAGCCCTCTGTCCCC
ACTAGCCACGCGTCACTGGTTAGCGTGATTGAAACTAAATCGTATGAAAATCCTCTTCTCTAGTCGCACTAGCCACGTTTCG
AGTGCTTAATGTGGCTAGTGGCACC GGTTTGGACAGCACAGCTGTAAAATGTTCCCATCCTCACAGTAAGCTGTTACCGTTC
CAGGAGATGGGACTGAATTAGAATTCAAACAAATTTTCCAGCGCTTCTGAGTTTTACCTCAGTCACATAATAAGGAATGCAT
CCCTGTGTAAGTGCATTTTGGTCTTCTGTTTTTGCAGACTTATTTACCAAGCATTGGAGGAATATCGTAGGTA AAAATGCCTA
TTGGATCCAAAGAGAGGCCAACATTTTTTTGAAATTTTTTAAGACACGCTGCAACAAAGCAGGTATTGACAAATTTTATATAAC
TTTATAAATTACACCGAGAAAGTGTTTTTCTAAAAAATGCTTGCTAAAAACCCAGTACGTCACAGTGTTGCTTAGAACCATAA
ACTGTTCCCTTATGTGTGTATAAATCCAGTTAACAACATAATCATCGTTTGCAGGTTAACCACATGATAAATATAGAACGTCT
AGTGGATAAAGAGGAAACTGGCCCCCTTGACTAGCAGTAGGAACAATTACTAACAAATCAGAAGCATTAAATGTTACTTTATGG
CAGAAGTTGTCCAACCTTTTTTGGGAGTGCTTTTGTATTATG
TAGCTTACCATATTAGAAATTTATCCCAGCACTTTGGGA
GGCCGAGGTGGGCGGATCACTTCTATCTCTACTAAAAAT
ACAAAAAATGTGCTGCGTGTGGGAGAATCGCTTGAACCC
TGGAGGCAGAGGTTGCAGTGAGACTCTGTCTCAAAACAA
ACAAACAAACAAAAAATAAGAAATTAAGTTAATTTACTTAAAAATAATGAAAGCTAAGCCATTGCATATTATCACAACAT
TCTTAGGAAAAATAACTTTTTTGAAAACAAGTGAGTGGAATAGTTTTTTACATTTTTTGCAGTTCTCTTTAATGTCTGGCTAAAT
AGAGATAGCTGGATTCACTTATCTGTGTCTAATCTGTATTATTTGGTAGAAGTATGTGAAAAAAAATTAACCTCACGTTGAAA
AAAGGAATATTTTAATAGTTTTTCAGTTACTTTTTGGTATTTTTCCCTTGTA CTTTGCATAGATTTTTCAAAGATCTAATAGAT
ATACCATAGGTCTTTCCCATGTCGCAACATCATGCAGTGATTATTTGGAAGATAGTGGTGTTCTGAATTATACAAAGTTTCC
AAATATTGATAAATTGCATTAACTATTTTAAAAATCTCATTCAATTAATACCACCATGGATGTCAGAAAAGTCTTTTAAGAT
TGGGTAGAAATGAGCCACTGGAAATTCTAATTTTCATTTGAAAGTTCACATTTTGT CATTGACAACAACTGTTTTCCCTGC
AGCAACAAGATCACTTCATTGATTTGTGAGAAAAATGTCTACCAAATTATTTAAGTTGAAATAACTTTGTCAGCTGTTCTTTC
AAGTAAAAATGACTTTTTCATTGAAAAAATTGCTTGTTTCAGATCACAGCTCAACATGAGTGCTTTTCTAGGCAGTATTGTACT
TCAGTATGCAGAAGTGCTTTATGTATGCTTCCTATTTTGT CAGAGATTATTAAAAAGAAGTGCTAAAGCATTGAGCTTCGAAA
TTAATTTTTTACTGCTTCATTAGGACATTCTTACATTAAACTGGCATTATTATTACTATTATTTTAAACAAGGACACTCAGTG
GTAAGGAATATAATGGCTACTAGTATTAGTTTGGTGCCACTGCCATAACTCATGCAAATGTGCCAGCAGTTTTTACCCAGCAT

**Large amounts of raw
DNA sequence data**

We need to make the data mean something...



<http://ensembl.org>




<http://genome.ucsc.edu>

Ensembl features

- Gene builds
- Gene trees
- Variation display and VEP
- Display of user data
- BioMart (data export)
- Programmatic access via API
- Completely open source



Exploring the Ensembl Genome Browser

 [BLAST/BLAT](#) | [BioMart](#) | [Tools](#) | [Downloads](#) | [Help & Documentation](#) | [Blog](#) | [Mirrors](#)

Search: for


Go


e.g. [BRCA2](#) or [rat 5:62797383-63627669](#) or [rs699](#) or [coronary heart disease](#)


Browse a Genome


The Ensembl project produces genome databases for vertebrates and other eukaryotic species, and makes this information freely available online.

Popular genomes

 **Human**
GRCh38.p5

 **Human**
GRCh37

 **Mouse**
GRCm38.p4

 **Zebrafish**
GRCz10

★ [Log in to customize this list](#)


All genomes

-- Select a species --


[View full list of all Ensembl species](#)

Other species are available in [Ensembl PreRelease](#) and [Ensembl Genomes](#)

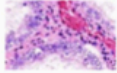
Still using Human GRCh37?

Go to 


Variant Effect Predictor



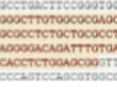
Gene expression in different tissues




Find SNPs and other variants for my gene



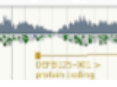
Retrieve gene sequence




Compare genes across species




Use my own data in Ensembl




ENCODE data in Ensembl

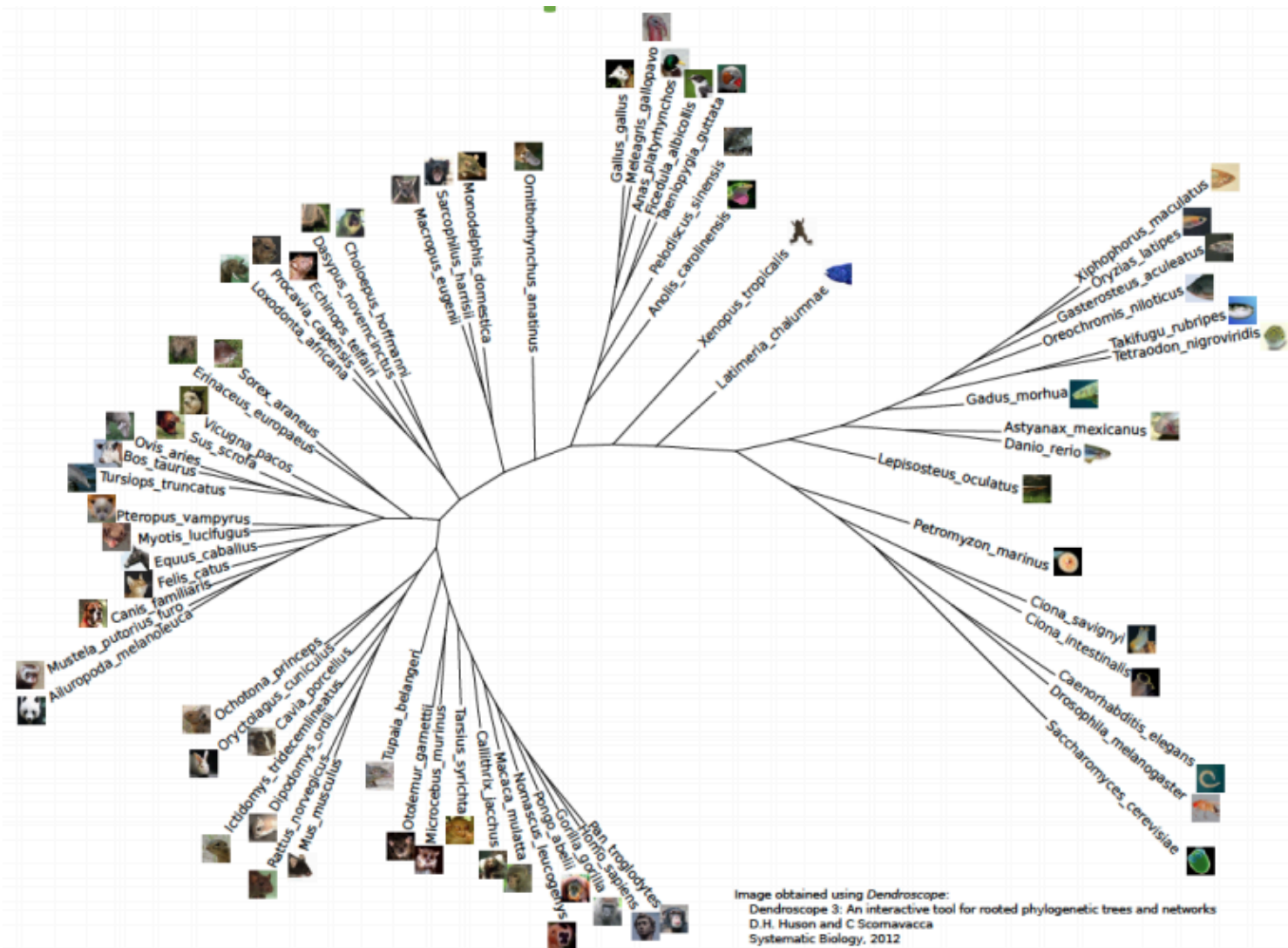


Did you know...?

 Access archive sites to see previous releases of Ensembl. Add e## to the



87 Vertebrate Species in Ensembl v84



+ early access Pre! websites

Non-vertebrates on Ensembl Genomes

The screenshot shows the EnsemblBacteria interface. At the top, there's a search bar and navigation links. Below, the 'Find a Species' section lists various bacterial species, including *Bacillus anthracis* and *Bacillus cereus*, with links to their respective genome pages.

Bacteria (39,584 genomes)

The screenshot shows the EnsemblProtists interface. It features a search bar and a list of protist species, including *Plasmodium falciparum* and *Toxoplasma gondii*, with links to their genome pages.

Protists (158)

The screenshot shows the EnsemblFungi interface. It displays a search bar and a list of fungal species, including *Aspergillus fumigatus* and *Neurospora crassa*, with links to their genome pages.

Fungi (589)

The screenshot shows the EnsemblMetazoa interface. It includes a search bar and a list of metazoan species, primarily focusing on *Drosophila* species like *Drosophila melanogaster* and *Drosophila obscura*, with links to their genome pages.

Metazoa (65)

The screenshot shows the EnsemblPlants interface. It features a search bar and a list of plant species, including *Oryza sativa* (rice) and *Arabidopsis thaliana*, with links to their genome pages.

Plants (39)

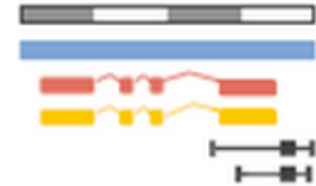
Ensembl vs Ensembl Genomes

	Ensembl	EnsemblGenomes
Released	2000	2009
Species	Vertebrates (fly, worm and yeast as outgroups)	Non-vertebrates (protists, plants, fungi, metazoa, bacteria)
Annotation	by Ensembl	in collaboration with the scientific communities
URL	www.ensembl.org	www.ensemblgenomes.org

- Joint project between EMBL-EBI and Sanger
- Funded primarily by the Wellcome Trust

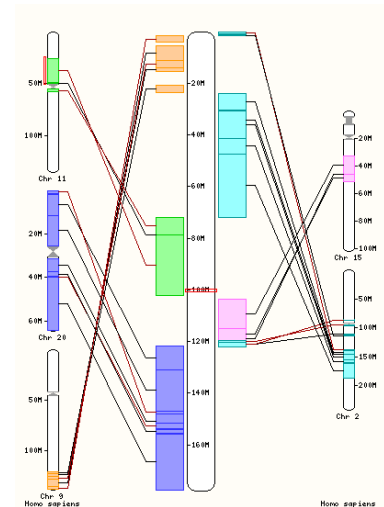
Data

- Genomic sequence
- Gene / transcript / protein models
- External references
- Mapped sequences
 - cDNAs, proteins, repeats, markers, probes, etc.
- Variation data:
 - sequence variants
 - structural variants (selected species)



Data

- Comparative data:
 - Orthologues and paralogues (between plants and pan-taxonomic)
 - Protein families
 - Whole genome pairwise alignments (selected species)
 - Synteny (selected species)
 - whole genome multiple alignment



Synteny between Rat chromosome 3
and Human

Methods in gene annotation

Automatic annotation

Using known
proteins/ESTs/cDNAs

Homology annotation

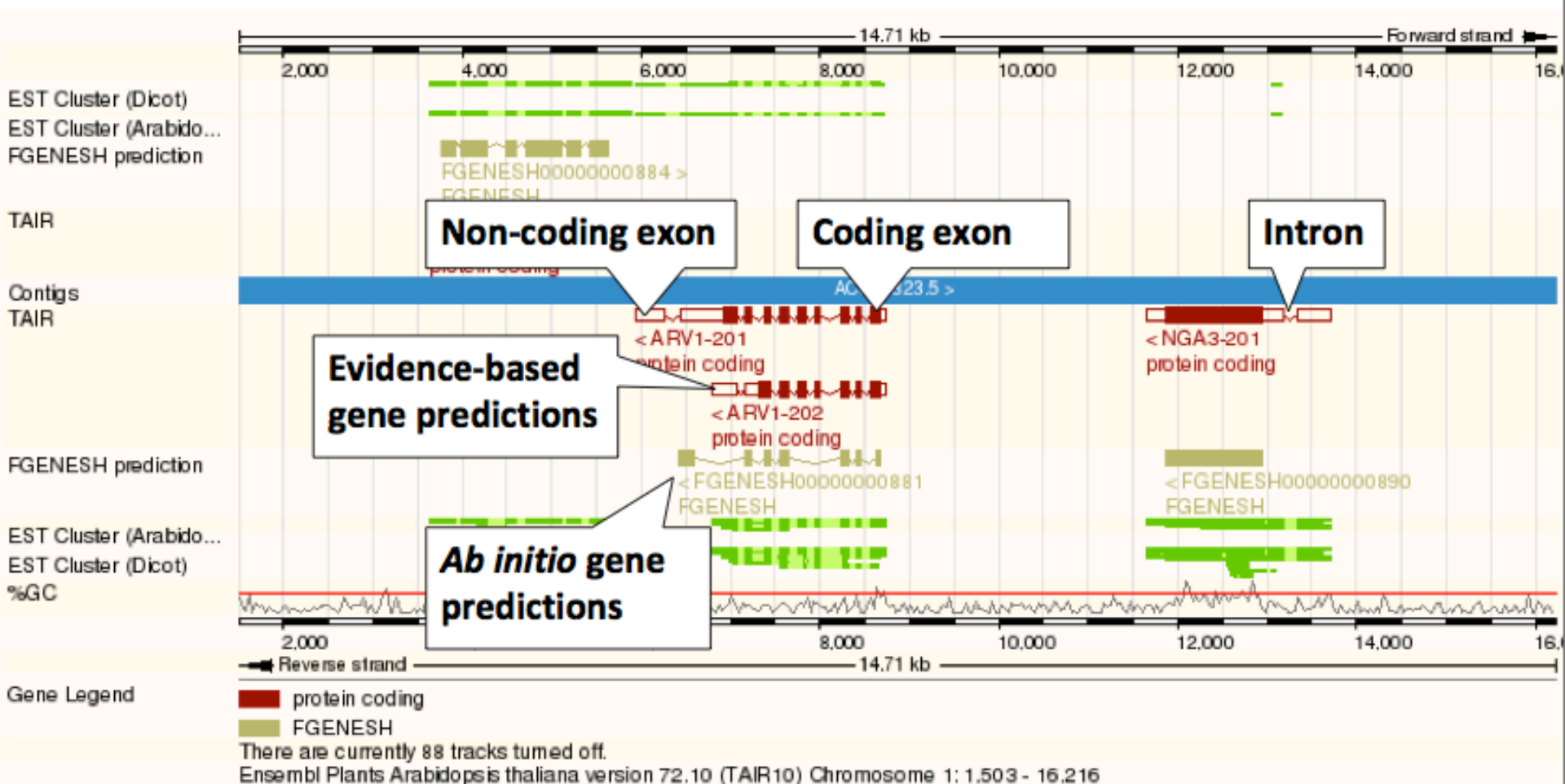
Using known
proteins/ESTs/cDNAs
from other species

Manual annotation

To correct errors in
automatic annotation

***Ab initio* predictions**
(finding apparent ORFs
in the sequence)

Gene view



Transcript view

Arabidopsis thaliana Location: 5:17,164,141-17,165,918 Gene: DFR Transcript: AT5G42800.1

Transcript-based displays

- Transcript summary
- Supporting evidence
- Sequence
 - Exons (6)
 - cDNA
 - Protein
- EBI Protein Summary
- Protein Structure
- External References
 - General identifiers (142)
 - Oligo probes (2)
- Ontology
 - Ontology graphs
 - plant anatomical entity (15)
 - biological process (4)
 - plant structure development
 - cellular component (1)
 - molecular function (6)
 - Ontology table (37)
- Genetic Variation
 - Variation Table
 - Population comparison
 - Comparison image
- Protein Information
 - Protein summary
 - Domains & features (5)
 - Variations (71)
- External Data
 - Personal annotation
- ID History
 - Transcript history
 - Protein history

Description dihydroflavonol-4-reductase [Source: EMBL](#)


Location [Chromosome 5: 17,164,141-17,165,918](#)

Gene This transcript is a product of gene [AT5G42800](#) - This gene has 1 transcript

Show/hide columns Filter

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype
AT5G42800.1	AT5G42800.1	1358	AT5G42800.1	382	Protein coding

Transcript summary [help](#)



Reverse strand 1.78 Kb Export image

Statistics Exons: 6 Transcript length: 1,358 bps Translation length: 382 residues

Ensembl version AT5G42800.1.1

Type Known protein coding

Prediction Method Gene annotation by [TAIR](#) through a process of automatic and manual curation.

Ensembl Plants release 14 - May 2012 © [EBI](#) [About Ensembl Genomes](#) | [Contact Us](#) | [EMBL-EBI Terms of use](#) | [Privacy](#) | [Cookies](#) | [Help](#)

Configure this page Manage your data Export data Bookmark this page

Transcript tab

Changed side menu

Names in Ensembl

- ENS**G**#### Ensembl **G**ene ID
- ENST**T**#### Ensembl **T**ranscript ID
- ENS**P**#### Ensembl **P**eptide ID
- ENSE**E**#### Ensembl **E**xon ID

- For other species than human a suffix is added:

MUS (*Mus musculus*) for mouse: ENS**MUS**G####
DAR (*Danio rerio*) for zebrafish: ENS**DAR**G####, etc.

- Imported gene annotations keep the IDs (e.g. Arabidopsis, fly, C. elegans)

Access to data

- Web browser
 - <http://ensemblgenomes.org/>
 - <http://ensembl.org>
- BioMart 'Data mining tool'
 - <http://ensembl.org/biomart/martview/>
- FTP download site
 - <ftp://ftp.ensemblgenomes.org/pub/>
 - <http://ensembl.org/info/data/ftp/>
- Public MySQL server
 - `mysql.ebi.ac.uk:4157:anonymous`
- Ensembl APIs/ via your favourite programming language
 - <http://ensembl.org/info/docs/api/>
 - <http://rest.ensembl.org/>

BioMart

- Data export tool
- Quick table generator (e.g. Excel)
- Web interface to mine ensembl data
- Integrated with BioConductor

What can I do with BioMart?

Extract data for a large number of genes, eg:

- Export list of genes in a region
- Convert IDs
- Extract sequences
- Get orthologues and paralogues

4 steps



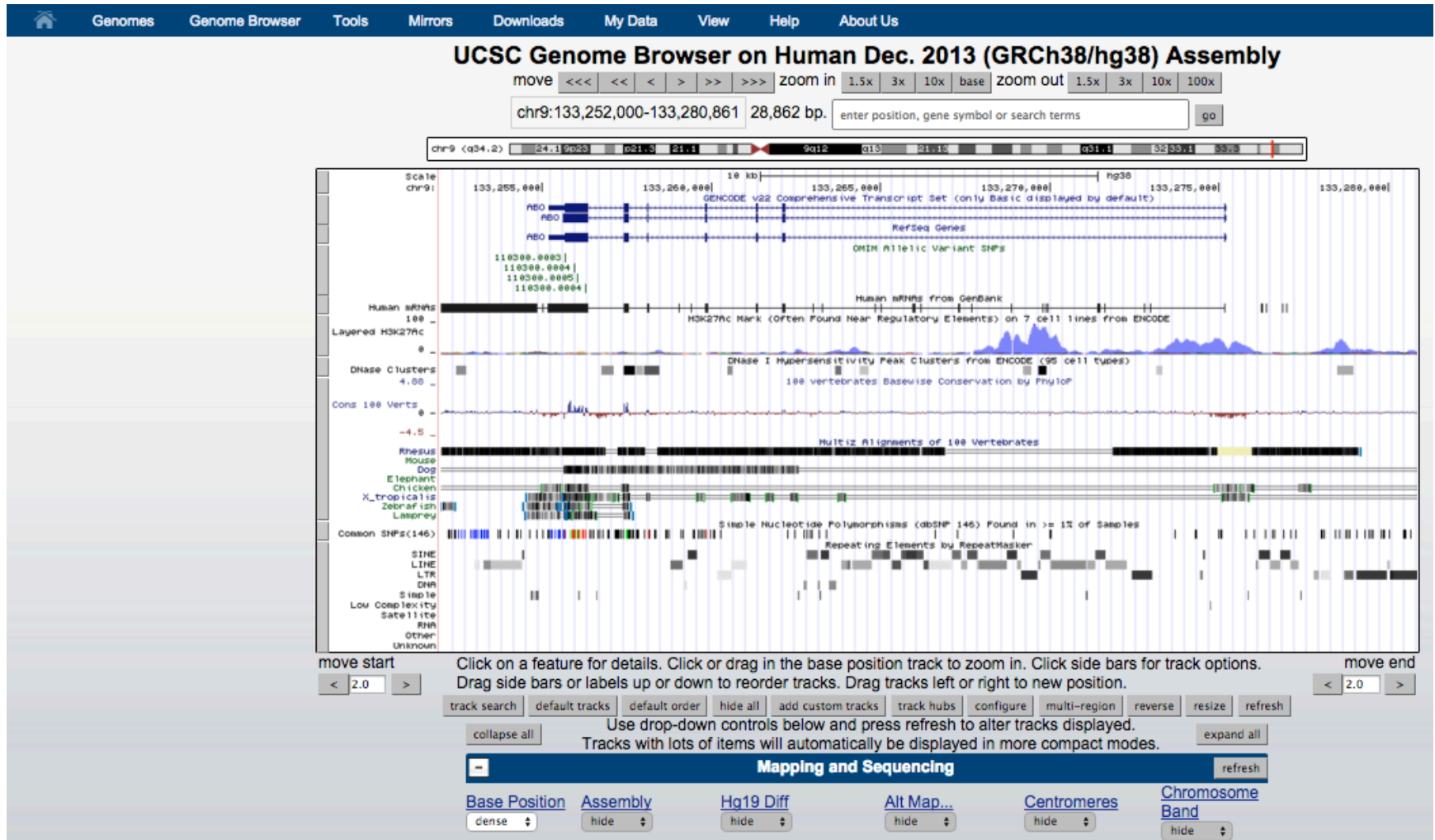
Dataset
choose
database
& species

Filters
what you
know

Attributes
what you want
to know

Results
table/
sequences

UCSC Genome Browser



<http://genome.ucsc.edu>

What Distinguishes Ensembl from the UCSC Browser?

- The gene set. Automatic annotation based on mRNA and protein information.
- Species coverage (UCSC only animals)
- Comparative analysis (gene trees)
- BioMart (vs Table Browser)
- Programmatic access via the Perl API (open source)
- Integration with other databases (DAS)

Summary

Ensembl is a genome browser which integrates:

- gene annotation
- variation
- regulation (for selected species)
- comparative genomics

Sources

- slides by Dan Bolser / Bert Overduin
- <http://www.ensembl.org/info/website/tutorials/index.html>