

Browsing Genomic Information with Ensembl

Stefan Wyder

January 2018



Universität
Zürich ^{UZH}



URPP
**Evolution in
Action**

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
CAGAAGTTGTCCAACCTTTTTGGGAGTGCTTTTGTATTATG
TAGCTTACCATATTAGAAATTTATCCCAGCACTTTGGGA
GGCCGAGGTGGGCGGATCACTTCTATCTCTACTAAAAAT
ACAAAAAATGTGCTGCGTGTGGGAGAATCGCTTGAACCC
TGGAGGCAGAGGTTGCAGTGAGACTCTGTCTCAAAACAA
ACAAACAAACAAAAAACTAAGAAATTAAGTTAATTTACTTAAAAATAATGAAGCTAACCATTGCATATTATCACAACAT
TCTTAGGAAAAATAACTTTTTGAAAACAAGTGAGTGGAATAGTTTTTACATTTTTTGCAGTTCTCTTTAATGTCTGGCTAAAT
AGAGATAGCTGGATTCACTTATCTGTGTCTAATCTGTATTATTTGGTAGAAGTATGTGAAAAAAAATTAACCTCACGTTGAAA
AAAGGAATATTTTAATAGTTTTTCAGTTACTTTTTGGTATTTTTCCCTTGTA CTTTGCATAGATTTTTCAAAGATCTAATAGAT
ATACCATAGGTCTTTCCCATGTCGCAACATCATGCAGTGATTATTTGGAAGATAGTGGTGTTCTGAATTATACAAAGTTTCC
AAATATTGATAAATTGCATTAACTATTTTAAAAATCTCATTCAATTAATACCACCATGGATGTCAGAAAAGTCTTTTAAGAT
TGGGTAGAAATGAGCCACTGGAAATTCTAATTTTCATTTGAAAGTTCACATTTTGT CATTGACAACAACTGTTTTCCCTGC
AGCAACAAGATCACTTCATTGATTTGTGAGAAAAATGTCTACCAAATTATTTAAGTTGAAATAACTTTGTCAGCTGTTCTTTC
AAGTAAAAATGACTTTTTCATTGAAAAAATTGCTTGTTTCAGATCACAGCTCAACATGAGTGCTTTTCTAGGCAGTATTGTACT
TCAGTATGCAGAAGTGCTTTATGTATGCTTCCTATTTTGT CAGAGATTATTAAAAGAAGTGCTAAAGCATTGAGCTTCGAAA
TTAATTTTTTACTGCTTCATTAGGACATTCTTACATTAAACTGGCATTATTATTACTATTATTTTAAACAAGGACACTCAGTG
GTAAGGAATATAATGGCTACTAGTATTAGTTTGGTGCCACTGCCATAACTCATGCAAATGTGCCAGCAGTTTTACCCAGCAT

**Large amounts of raw
DNA sequence data**

We need to make the data mean something...



<http://ensembl.org>



<http://genome.ucsc.edu>

Ensembl

Ensembl is a data base + genome browser which integrates 4 types of resources:

- gene annotation
- variation (SNPs, structural, ...)
- regulation (for selected species)
- comparative genomics

Entry page

main menu






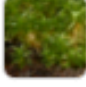
help

The screenshot shows the Ensembl Plants website interface. At the top is a green navigation bar with the Ensembl Plants logo on the left and a 'Login/Register' link on the right. Between them is a main menu with links: HMMER, BLAST, BioMart, Tools, Downloads, Documentation, and Website help. To the right of the menu is a search bar with a magnifying glass icon and the text 'Search Ensembl Plants...'. Below the navigation bar is a light blue search box with a dropdown menu set to 'All species', a text input field, and a 'Go' button. Below this is a section titled 'Favourite genomes' containing six entries, each with a small image, the species name, and a version number. Below the favourites is a section titled 'All genomes' with a dropdown menu set to '-- Select a species --'. To the right of the search box is a yellow callout bubble with the word 'search'. To the right of the 'All genomes' dropdown is a yellow callout bubble with the words 'select species'. On the right side of the page, there is a section titled 'New Ensembl Genomes Archive Sites' with a list of five URLs and a paragraph of text. Below this is a section titled 'New and updated genomes' with a paragraph of text. A yellow callout bubble with the word 'help' points to the 'Website help' link in the main menu.

Search: for

e.g. **Carboxy*** or **chx28**

Favourite genomes

 Arabidopsis thaliana TAIR10	 Oryza sativa Japonica IRGSP-1.0
 Triticum aestivum TGACv1	 Hordeum vulgare Hv_IBSC_PGSC_v2
 Zea mays AGPv4	 Physcomitrella patens ASM242v1

[Edit favourites](#)

All genomes

-- Select a species --

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

New Ensembl Genomes Archive Sites

Ensembl Genomes now has archive sites for all divisions:

- <http://oct2017-bacteria.ensembl.org>
- <http://oct2017-fungi.ensembl.org>
- <http://oct2017-metazoa.ensembl.org>
- <http://oct2017-plants.ensembl.org>
- <http://oct2017-protists.ensembl.org>

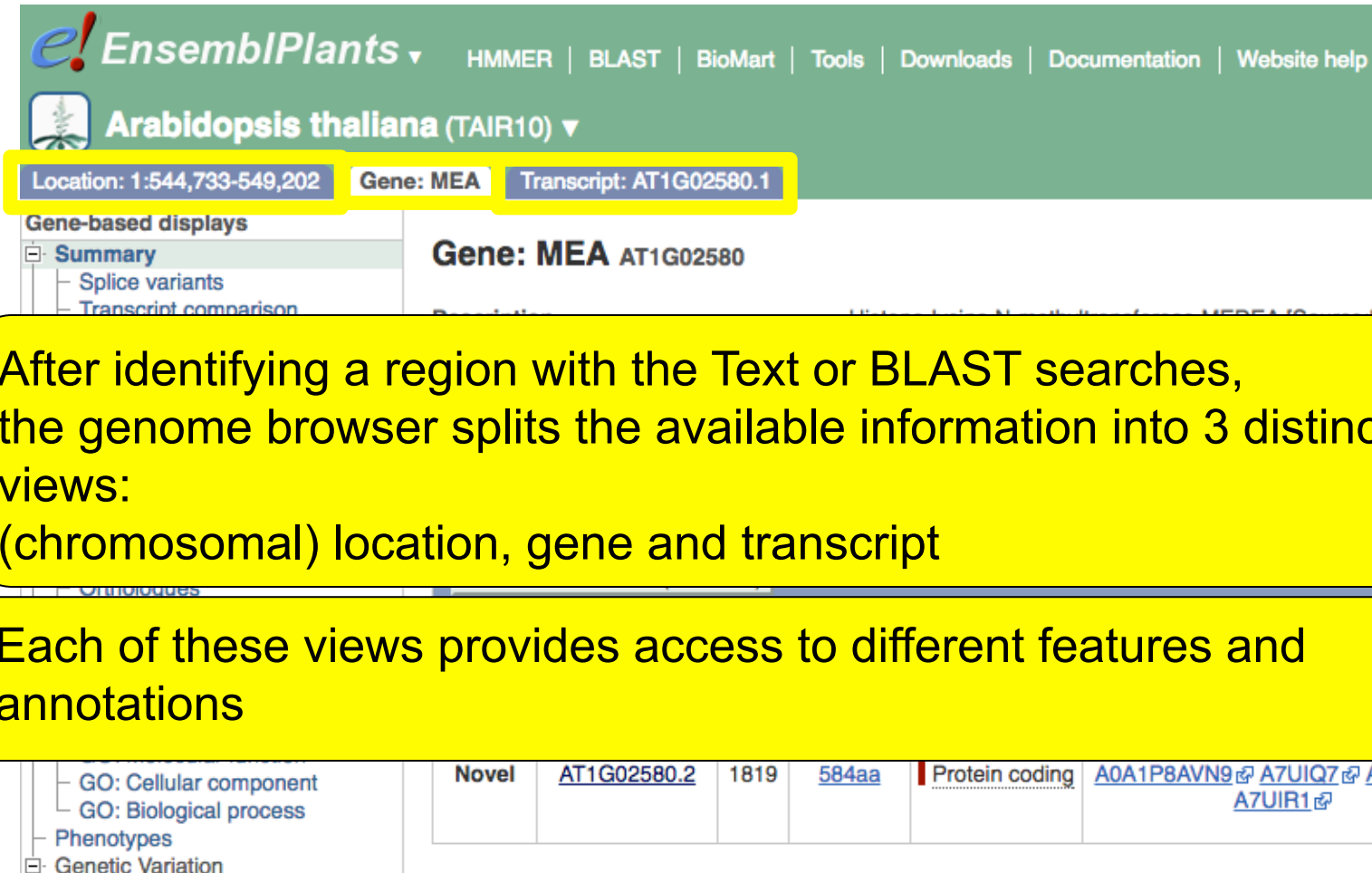
allowing researchers to access data from old releases and display track hubs for previous assemblies. Archive sites are searchable and have BioMarts available. Archival REST servers will not initially be available, but will be added in future.

The first release of the archive sites contains content from Release 37. New archive sites will be released at least once a year, under URLs indicating the date of the data they contain. The previously existing archive for Ensembl Plants, <http://archive.plants.ensembl.org>, will continue to be available at this URL, but also as <http://mar2016-plants.ensembl.org>, in accordance with the new naming scheme. As previously, data from all recent releases will continue to be available for download at <ftp://ftp.ensemblgenomes.org>.

New and updated genomes

This bumper release of [Ensembl Plants](#) brings eight new genomes, including [cassava](#), [cotton](#), [cucumber](#), [green bean](#), [sunflower](#) and [yam](#). In addition we have updated to latest and greatly improved (v3.0) [sorghum](#) genome assembly from JGI as well as updated

3 tabs with different views/information



EnsemblPlants | HMMER | BLAST | BioMart | Tools | Downloads | Documentation | Website help

Arabidopsis thaliana (TAIR10)

Location: 1:544,733-549,202 | Gene: MEA | Transcript: AT1G02580.1

Gene-based displays

- Summary
- Splice variants
- Transcript comparison

Gene: MEA AT1G02580

After identifying a region with the Text or BLAST searches, the genome browser splits the available information into 3 distinct views:
(chromosomal) location, gene and transcript

Each of these views provides access to different features and annotations

Novel	AT1G02580.2	1819	584aa	Protein coding	A0A1P8AVN9	A7UIQ7	A7UIR1

GO: Cellular component
GO: Biological process
Phenotypes
Genetic Variation

Chromosomal location view

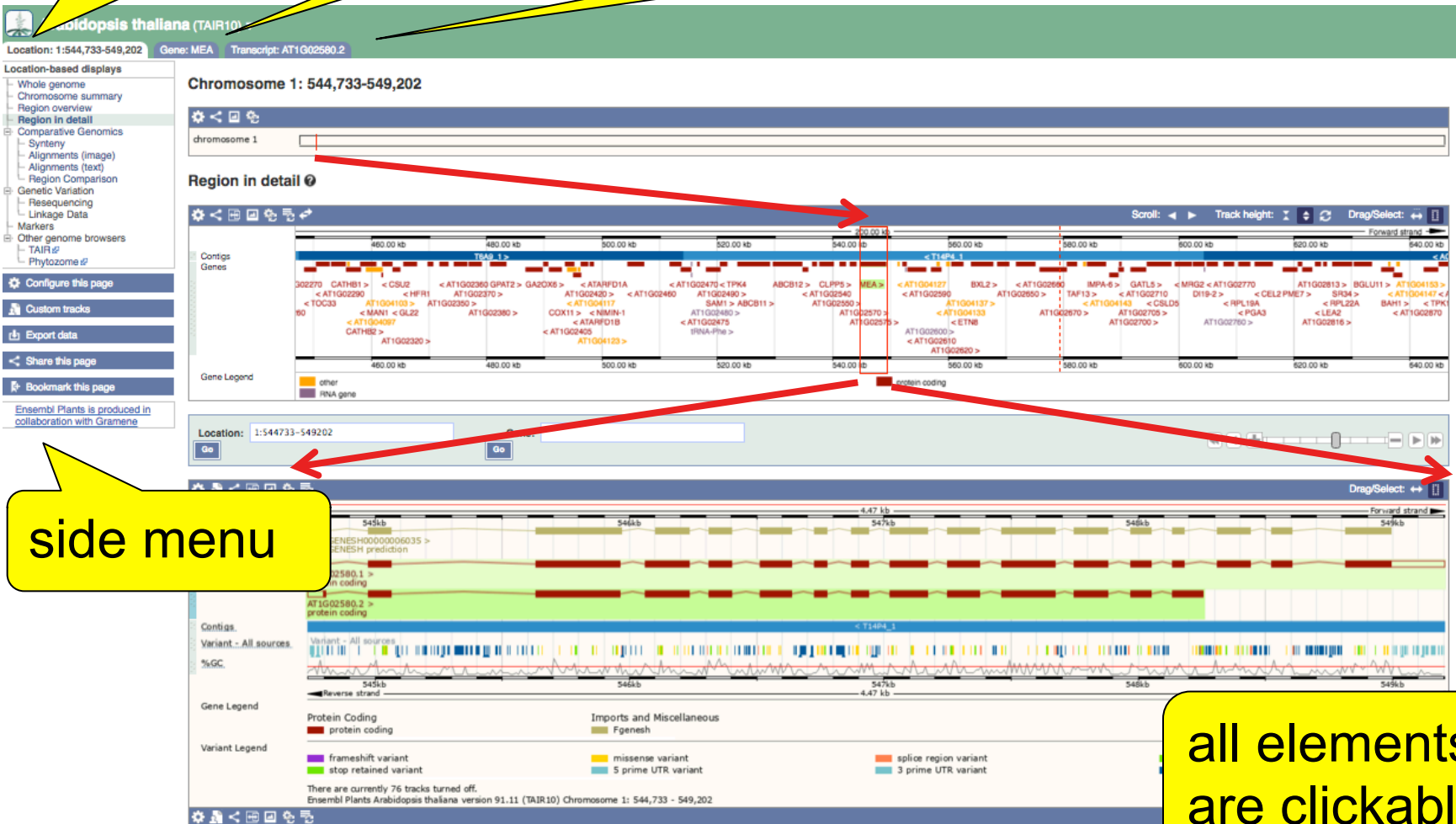
Location tab

Gene tab

Transcript tab

side menu

all elements
are clickable



Gene view

Location: 1:544,733-549,202

Gene: MEA

Transcript: AT1G02580.1

Gene-based displays

Summary

Splice variants

Transcript comparison

Gene alleles

Sequence

Secondary Structure

Gene families

Literature

Orthologues

Paralogues

Pan-taxonomic Comparison

Gene Tree

Orthologues

Ontologies

GO: Molecular function

GO: Cellular component

GO: Biological process

Phenotypes

Genetic Variation

Variant table

Variant image

Structural variants

Gene expression

Pathway

Regulation

External references

Supporting evidence

Gene: MEA AT1G02580

Description

Synonyms

Location

About this gene

Transcripts

Histone-lysine N-methyltransferase MEDEA [Source:UniProtKB/Swiss-Prot;Acc:[O65312](#)]

At1g02580, EMB173, FIS1, MEDEA, SDG5, SET5, T14P4.11

[Chromosome 1: 544,733-549,202](#) forward strand.

This gene has 2 transcripts ([splice variants](#)), [23 orthologues](#) and [8 paralogues](#).

Hide transcript table

lists isoforms

Name	Transcript ID	bp	Protein	Biotype	UniProt	RefSeq	Flags
Novel	AT1G02580.1	2341	689aa	Protein coding	A0A178WQQ4 A7UIQ9 A7UIR6 C0KG89 C0KG93 C0KG94 C0KG95 C0KG97 C0KG98 C0KG99 O65312	NM_001331345.1 NM_100139.4 NP_563658.1	
Novel	AT1G02580.2	1819	584aa	Protein coding	A0A1P8AVN9 A7UIQ7 A7UIQ8 A7UIR1	NM_001331345.1 NM_100139.4 NP_001322965.1	

Summary

Name

UniProtKB

Gene type

Annotation method

MEA (UniProtKB Gene Name)

This gene has proteins that correspond to the following UniProtKB identifiers: [O65312](#)

Protein coding

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

Transcript view

Location: 1:544,733-549,202

Gene: MEA

Transcript: AT1G02580.1

Transcript-based displays

Summary

Sequence

- Exons
- cDNA
- Protein

Protein Information

- Protein summary
- Domains & features
- Variants

Genetic Variation

- Variant table
- Variant image
- Population comparison
- Comparison image

External References

- General identifiers
- Oligo probes

Supporting evidence

ID History

- Transcript history
- Protein history

Configure this page

Custom tracks

Export data

Share this page

Transcript: AT1G02580.1

Description

Location

About this transcript

Gene


Histone-lysine N-methyltransferase MEDEA [Source:UniProtKB/Swiss-Prot;Acc:[O65312](#)]

[Chromosome 1: 544,733-549,202](#) forward strand.

This transcript has [17 exons](#), is annotated with [19 domains and features](#), is associated with [108 variations](#) and maps to [3 oligo probes](#).

This transcript is a product of gene [AT1G02580](#) [Show transcript table](#)

Summary ⓘ



Statistics

Uniprot

Version

Type

Annotation Method

Exons: 17, Coding exons: 17, Transcript length: 2,341 bps, Translation length: 689 residues

This transcript corresponds to the following Uniprot identifiers: [O65312](#)

AT1G02580.1.

Protein coding

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

Gene view | Gene tree

Location: 1:544,733-549,202 **Gene: MEA**

Gene-based displays

- Summary
 - Splice variants
 - Transcript comparison
 - Gene alleles
- Sequence
 - Secondary Structure
- Gene families
- Literature
- Plant Compara
 - Genomic alignments
 - Gene tree**
 - Gene gain/loss tree
 - Orthologs
 - Paralogous
- Pan-taxonomic
- Orthologs
- Paralogous
- Structural variants
- Gene expression
- Pathway
- Regulation
- External references
- Supporting evidence
- ID History
- Gene history

Gene: MEA AT1G02580

Description Histone-lysine N-methyltransferase MEDEA [Source:UniProtKB/Swiss-Prot;Acc:Q65312]

Synonyms AT1g02580, EMB173, FIS1, MEDEA, SDG5, SET5, T14P4.11

Location [Chromosome 1: 544,733-549,202](#) forward strand.

About this gene This gene has 2 transcripts ([splice variants](#)), [23 orthologues](#) and [8 paralogues](#).

Transcripts [Show transcript table](#)

Gene tree [GeneTree EPIGT00880000131274](#)

491
361
100
22
7

[Show annotations table](#)

Click on "Gene tree"

Gene tree

Opisthokonta: 4 homologs
Papilionoideae: 10 homologs
Brassica: 8 homologs
MEA, Arabidopsis thaliana
MEA, Arabidopsis lyrata
Land plants: 128 homologs
Eukaryotes: 339 homologs

LEGEND

Branch Length

- x1 branch length
- x10 branch length
- x100 branch length

Genes

- Gene ID gene of interest
- Gene ID within-sp. paralog

Nodes

- gene node
- speciation node
- duplication node
- ambiguous node
- gene split event

Collapsed Nodes

- collapsed sub-tree
- collapsed (paralog)
- collapsed (gene of interest)

Collapsed Alignments

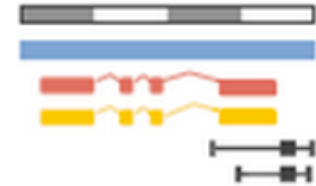
- 0 - 33% Aligned AA
- 33 - 66% Aligned AA
- 66 - 100% Aligned AA

Expanded Alignments

- Gap
- Aligned AA

Gene annotation

- Genomic sequence
- Gene / transcript / protein models
- External references
- Mapped sequences
 - cDNAs, proteins, repeats, markers, probes, etc.



Methods in gene annotation

Automatic annotation

Using known
proteins/ESTs/cDNAs

Homology annotation

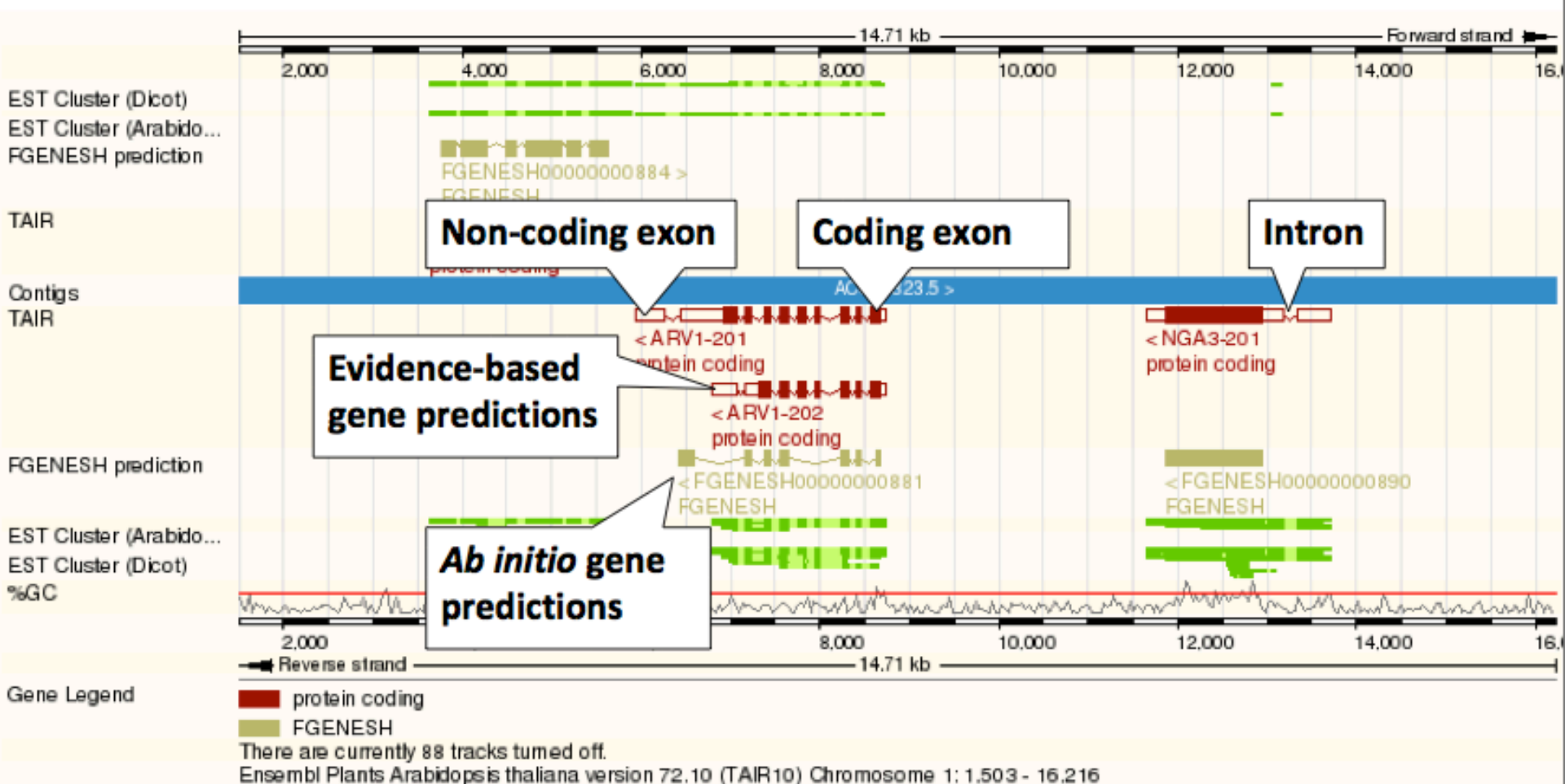
Using known
proteins/ESTs/cDNAs
from other species

Manual annotation

To correct errors in
automatic annotation

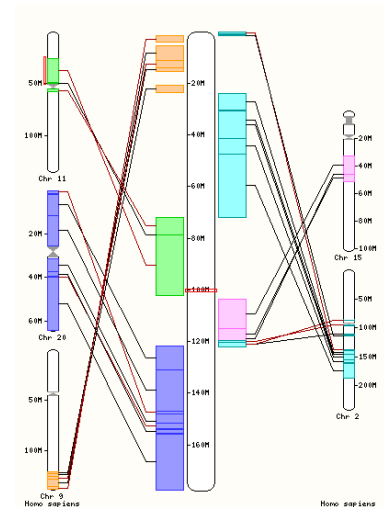
***Ab initio* predictions**
(finding apparent ORFs
in the sequence)

Gene view



Comparative data

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



Synteny between Rat chromosome 3
and Human

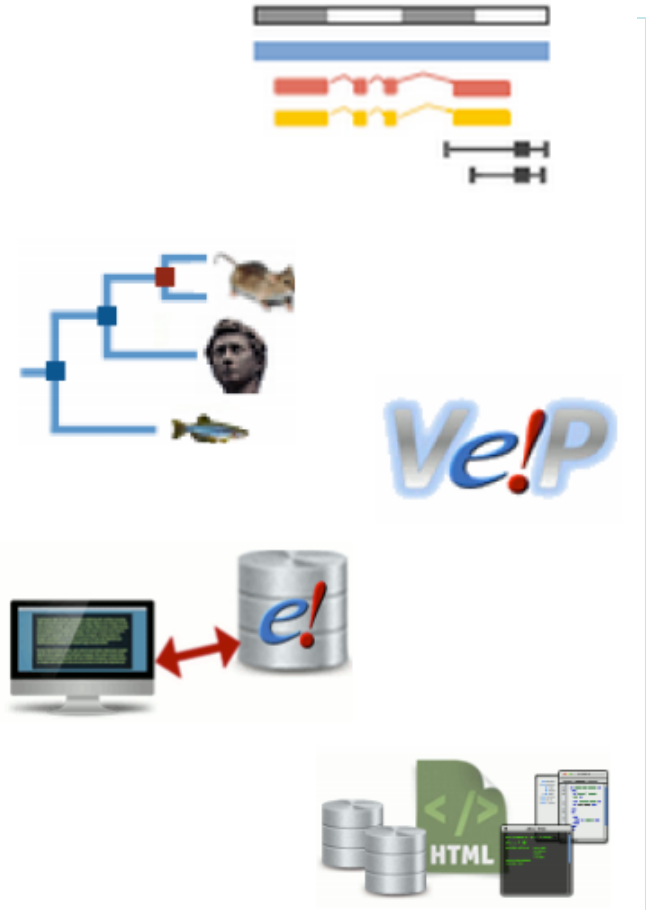
Names in Ensembl

- ENS**G**#### Ensembl **Gene** ID
- ENST**T**#### Ensembl **Transcript** ID
- ENS**P**#### Ensembl **Peptide** ID
- ENSE**E**#### Ensembl **Exon** ID
- For other species than human a suffix is added:

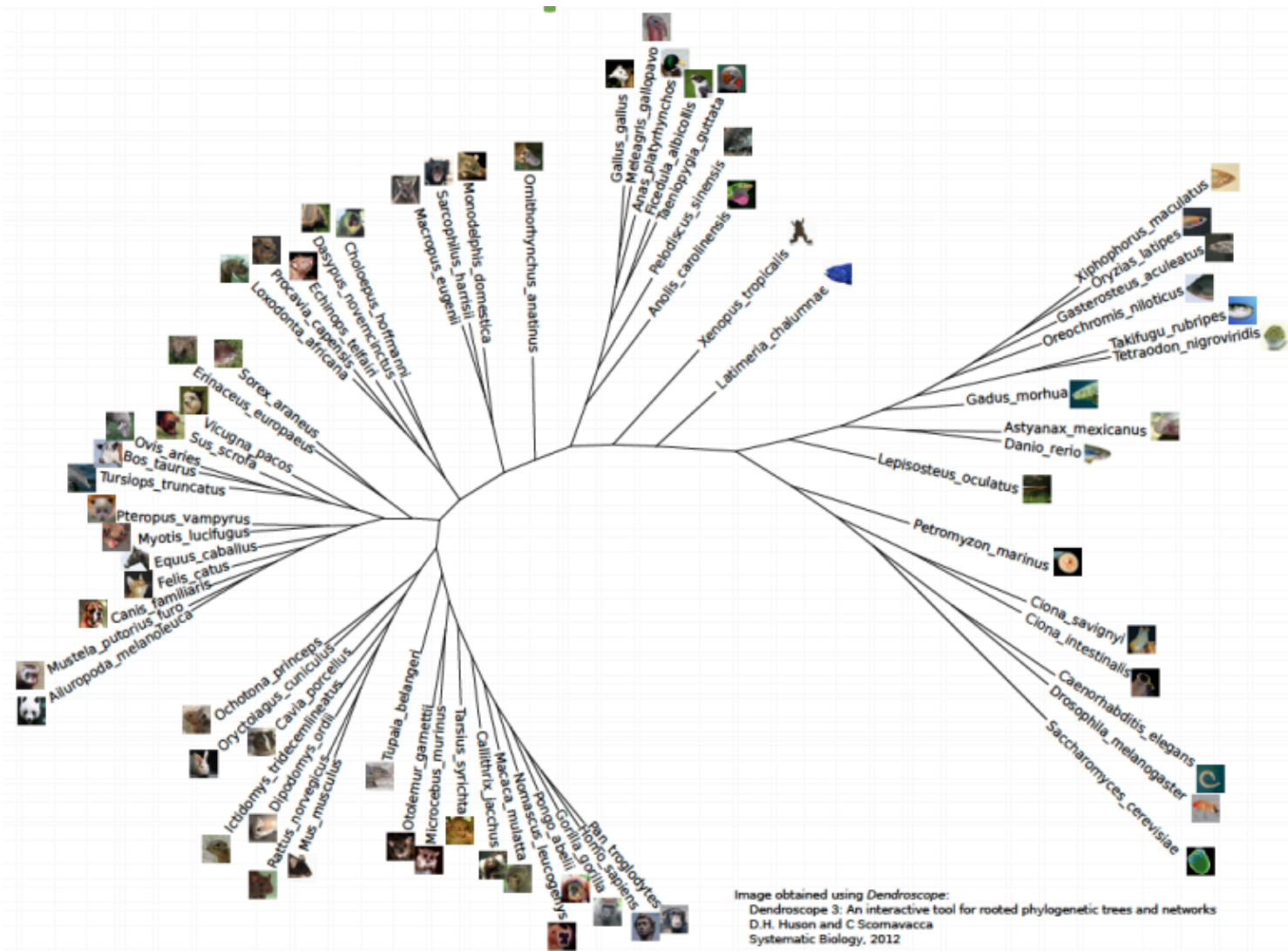
MUS (*Mus musculus*) for mouse: ENS**MUS**G####
DAR (*Danio rerio*) for zebrafish: ENS**DAR**G####, etc.
- Imported gene annotations keep the IDs (e.g. Arabidopsis, fly, C. elegans)

More Ensembl features

- Variation display and VEP
- Display of user data
- Completely open source
- 4-5 updates (versions/freezes) a year
- archived versions



124 vertebrate species in Ensembl v91



+ 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 (44,046 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 (189)

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

Fungi (811)

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 (69)

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 (53)

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, e.g.:

- Export list of genes in a region
- Convert IDs
- Retrieve fasta sequences of all introns
- Get orthologues and paralogues

Biomart: 4 steps



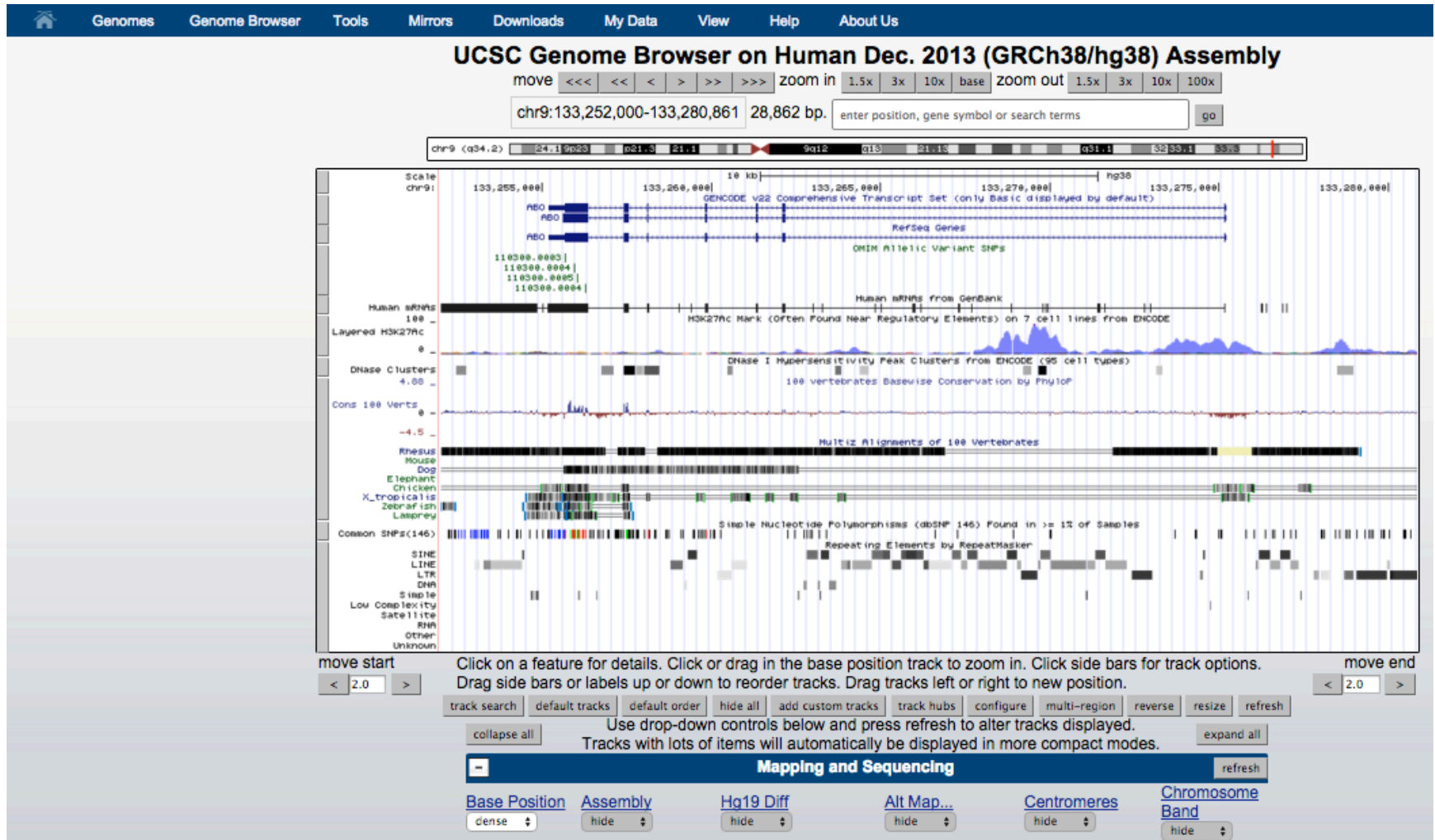
Dataset
choose
database
& species

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)

Sources

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

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