

Wheat Genomes in Ensembl Plants

Paul Kersey

EMBL-EBI

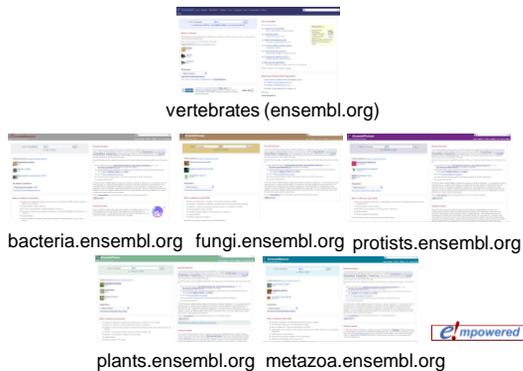
Ensembl 

- A platform for genome annotation, analysis and browsing, developed jointly by the EBI and Wellcome Trust Sanger Institute
- Modules for handling variation data, comparative genomics, gene prediction etc.
- Multiple points of access to data: web-based application, Perl and REST-ful APIs, public MySQL databases, BioMart data mining tool, DAS client and server, FTP
- Upload your own data and compare to reference annotation
- Originally developed for vertebrate genomes, subsequently extended to non-vertebrate species ("Ensembl Genomes")

2

26.02.2014

EMBL-EBI



vertebrates (ensembl.org)

bacteria.ensembl.org fungi.ensembl.org protists.ensembl.org

plants.ensembl.org metazoa.ensembl.org



03.01.12
pkersay@ebi.ac.uk

IWGSC - Standards and Protocols PAG XXX

EMBL-EBI

Data

- Genomic sequence
- Gene / transcript / protein models
- External references
- Mapped cDNAs, proteins, microarray probes, BAC clones, repeats, markers etc. etc.
- Variation data:
 - sequence variants
 - structural variants
- Comparative data:
 - gene trees, orthologues, paralogues
 - pairwise whole genomic alignments, syntenic regions



Wheat data in Ensembl Plants (release 21)

- Version 1.0 of the chromosome survey sequence assembly generated by the International Wheat Genome Sequencing Consortium.
- Version 2.0 of the IWGSC gene models called on the survey sequence
 - Version 2.1 due in next release (March 2014)
- Whole genome alignments against Brachypodium and rice
- Alignments to wheat Unigenes and RNA-seq data

5

19.02.2013 pkersay@ebi.ac.uk Plant and Animal Genomes 2014

EMBL-EBI

Chromosome survey sequence: plans for the next release (release 22, due March 2014)

- Whole genome alignments
 - Against barley, *Triticum uratu* and *Aegilops tauschcii*
 - Wheat AvB, AvD, BvD
- Alignment to wheat full length ESTs
- Additional RNA-seq alignments

6

19.02.2013 pkersay@ebi.ac.uk Plant and Animal Genomes 2014

EMBL-EBI

Genome assemblies from Brenchley *et al.*

- The wheat genome assemblies generated by Brenchley *et al.* (PMID: 23192148) have also been aligned to the survey sequence, Brachypodium and barley
 - Homeologous variants inferred between the 3 wheat genomes (from this work) are also displayed in the context of the gene models from Brachypodium and barley

Diploid progenitor genomes

- Aegilops tauschii* and *Triticum uratu* are also included in Ensembl Plants
- These genomes have been aligned to rice, and to relevant RNA-seq reads.
- These genomes and other more distant relatives (e.g. barley, Brachypodium, and rice) are all included in gene-centric comparative analyses

7 19.02.2013 jksrasy@ebi.ac.uk Plant and Animal Genomes 2014 EMBL-EBI

8 19.02.2013 jksrasy@ebi.ac.uk Plant and Animal Genomes 2014 EMBL-EBI

Sequence

- The chromosome survey sequence (and its annotations) are available to search via BLAST and other search alignment algorithms.
- EST-based search (with onward genomic mapping) also available

The screenshot shows the Ensembl Plants search interface. The search bar contains 'wheat'. Below the search bar, there are several sections: 'Popular genomes' with links to 'Aegilops tauschii', 'Triticum aestivum', and 'Brachypodium distachyon'; 'About this genome' with a link to 'Wheat genome'; and 'About wheat genome and gene annotations' with a link to 'Wheat genome and gene annotations'. The page also features a 'Read this about the assembly, annotation and usage of wheat genome provided by Ensembl Plants' section.

9 19.02.2013 jksrasy@ebi.ac.uk Plant and Animal Genomes 2014 EMBL-EBI

1 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI

The screenshot shows the 'Triticum aestivum Assembly and Gene Annotation' page. The page is divided into several sections: 'Chromosome survey sequences', 'About Triticum aestivum', and 'IWGSC Chromosome survey sequence'. The 'Chromosome survey sequences' section includes a table with the following data:

Assembly	Accession
Wheat	WGS11_L1_0110
Brachypodium	Brachypodium
Barley	Barley
Arabidopsis	Arabidopsis

The 'About Triticum aestivum' section provides information about the genome assembly and gene annotations. The 'IWGSC Chromosome survey sequence' section provides information about the survey sequence and its alignment to the wheat genome.

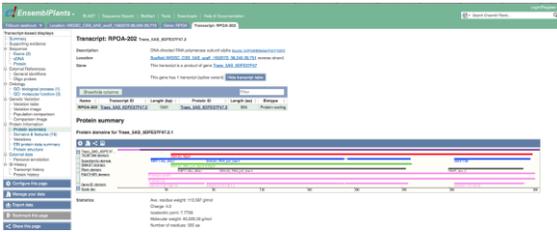
The screenshot shows the 'Triticum aestivum' page. The page is divided into several sections: 'Chromosome survey sequences', 'About Triticum aestivum', and 'IWGSC Chromosome survey sequence'. The 'Chromosome survey sequences' section includes a table with the following data:

Assembly	Accession
Wheat	WGS11_L1_0110
Brachypodium	Brachypodium
Barley	Barley
Arabidopsis	Arabidopsis

The 'About Triticum aestivum' section provides information about the genome assembly and gene annotations. The 'IWGSC Chromosome survey sequence' section provides information about the survey sequence and its alignment to the wheat genome.

1 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI

1 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI



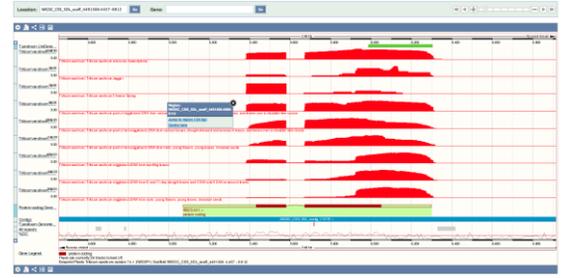
1 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI



2 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI



2 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI

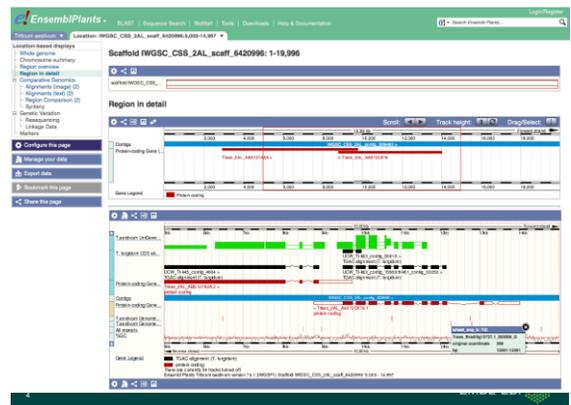


2 26.02.2014 IWGSC - Standards and Protocols PAG XXII EMBL-EBI

Inter-homeologous variants

- Data from Brenchley et al. mapped onto chromosome survey sequence
- Data from chromosome survey sequence itself due in next release

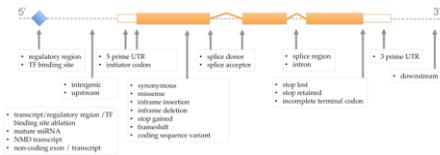
	Inter-homeologue variants (total)	Inter-homeologue variants (mapped)	% variants mapped
A/D	102,593	101,570	99
B	22,959	29,638	99
Total	132,552	131,208	99



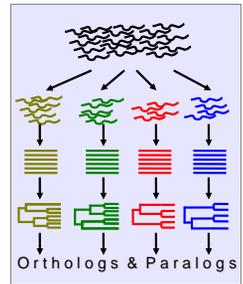
23 19.02.2013 jkarsky@ebi.ac.uk Plant and Animal Genomes 2014 EMBL-EBI

Variant Effect Predictor (VEP)

- Predicts functional consequences of known and unknown variants
- For substitutions, insertions, deletions and structural variants
- Web interface (for up to 750 variants), standalone Perl script, Perl API and REST API
- <http://www.ensembl.org/info/docs/variation/vep/index.html>



Gene tree pipeline



Take canonical protein for each gene belonging to one Ensembl Genomes clade

Cluster: WU-BLASTP + Smith-Waterman all-versus-all, hcluster_sg

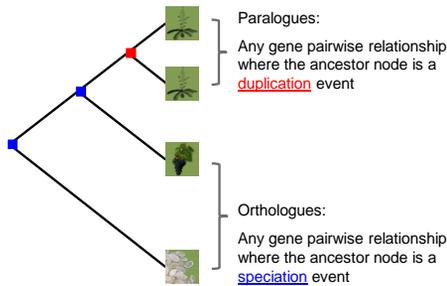
Align: multiple aligners consensified by M-Coffee

Build trees: PhyML-WAG + PhyML-HKY + NJ-p + NJ-dN + NJ-dS + species tree -> TreeBeST-merge

Infer orthologues and paralogues

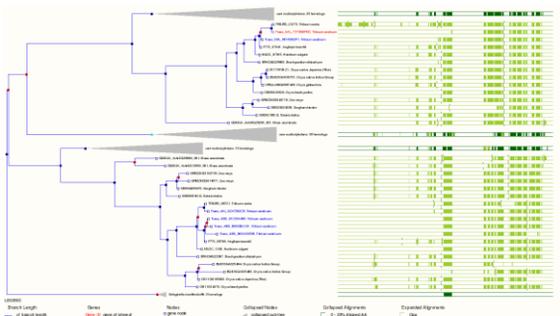
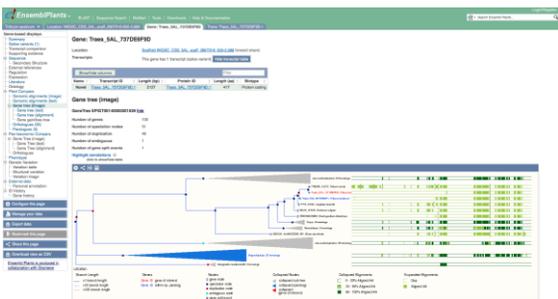


Orthologues and paralogues



Orthology / paralogy types

- ortholog_one2one
- ortholog_one2many
- ortholog_many2many
- apparent_ortholog_one2one
- possible_ortholog (weakly supported duplication node)
- within_species_paralog
- other_paralog (too distant to be in the same tree)
- contiguous_gene_split (artefact)
- putative_gene_split (artefact)
- 94,280 wheat genes in inferred orthology relationships with 270216 genes from other creals



Visualise your own data

Upload data

- Data saved by Ensembl
- 5 MB limit (and therefore not possible for large file formats)

Attach remote file

- URL-based (HTTP or FTP)
- No size limit
- Uploaded / attached tracks can be saved (to account)
- Uploaded / attached tracks can be shared with other users
- Only trivial security, do not use for sensitive data!
- <http://www.ensembl.org/info/website/upload/index.html>

EMBL-EBI

Possible formats

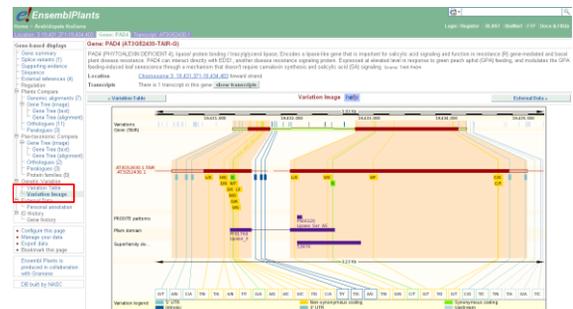
- BAM sequence alignments (no upload)
- BED genes / features
- BedGraph continuous-valued data
- BigBed genes / features (no upload)
- BigWig continuous-valued data (no upload)
- TrackHub collection of tracks
- Gbrowse genes / features
- GFF/GTF genes / features
- PSL sequence alignments
- VCF variants (no upload)
- WIG continuous-valued data

EMBL-EBI

Next release

- Improved representation of the polyploid genome
 - Preconfigured A v B v D comparative views
 - Separation of 3 wheat genomes in the gene tree analysis
- Representation of inter-homologue data as a full variation database
 - Data structure can also accommodate inter-varietal SNPs, links to phenotypes, haplotype structure etc.

EMBL-EBI

3 26.02.2014
4

EMBL-EBI

Funding

- Ensembl Genomes Funded by
 - EMBL
 - EU (INFRAVEC, Microme, transPLANT, AllBio)
 - BBSRC (PhytoPath, wheat/barley/midge sequencing, UK-US collaboration, RNAcentral)
 - Wellcome Trust (PomBase)
 - NIH/NIAID (VectorBase)
 - NSF (Gramene collaboration)
 - Bill and Melinda Gates Foundation (wheat rust)

35 26.02.2014

pkersley@ebi.ac.uk

EMBL-EBI

People

- James Allen, Irina Armean, **Dan Bolser**, Mikkel Christensen, Paul Davies, Lee Falin, Christoph Grabmueller, Kevin Howe, Malcolm Hinsley, Jay Humphrey, **Arnaud Kerhornou**, Julia Khobdova, Eugene Kulesha, Nick Langridge, Dan Lawson, Mark McDowall, Uma Maheswari, Gareth Maslen, Michael Nuhn, Chuang Kee Ong, Michael Paulini, Helder Pedro, Anton Petrov, Dan Staines, Mary Ann Tuli, **Brandon Walts**, Gary Williams
- The Ensembl teams @ EBI (Paul Flicek) and WTSI (Steve Searle)

36 26.02.2014

pkersley@ebi.ac.uk

EMBL-EBI