_____

# Module 2 – Maps and Genome Sequence
## i. The Ensembl Genome Browser

**Caveat**: At the time of writing this tutorial, Zv6 had not been released with a full gene build yet. All following examples are therefore taken from the Zv5 Ensembl. If you're trying to work through the examples yourself, please be aware of the difference in the scaffold naming ('Zv5_...' versus 'Zv6_...').

### Aims

- Explain the source for the data in Ensembl
- Introduce the Ensembl browser
- Show the different Ensembl views with examples

### Introduction

Ensembl is a joint project of the European Bioinformatics Institute (EBI) and the Wellcome Trust Sanger Institute, funded mainly by the Wellcome Trust, with additional funding from EMBL and NIH-NIAID. Ensembl provides easy access to genomic information with a number of visualisation tools.

The Ensembl site provides automatic baseline annotation of the latest assembly sequence, including gene, transcript and protein predictions.The annotation is integrated with external data sources, such as ZFIN for the zebrafish site. The latest zebrafish assembly is Zv6, which was released on March 31st, 2006.

The key Ensembl web pages are called Views (e.g. GeneView, TextView, MapView, and ContigView). The Ensembl web site gives you the opportunity to directly download data, whether it is a DNA sequence of a genomic contig you are trying to identify novel genes in, or positions of SNPs in a gene you are working on. There is also an FTP site which you can use to download large amounts of data from the Ensembl database, as well as a data mining tool (BioMart, see section 6) which allows flexible and rapid retrieval of information from the databases. There are many ways you can access the data in Ensembl depending on your needs and these are explained here and in other sections.

The Ensembl site is at:

<div align="center">

http://www.ensembl.org

</div>

On this page you will find links to all Ensembl species, documentation, search facilities, downloads and other related links. All Ensembl pages have a tool bar on the left-hand side with quick-access links to several resources and facilities.
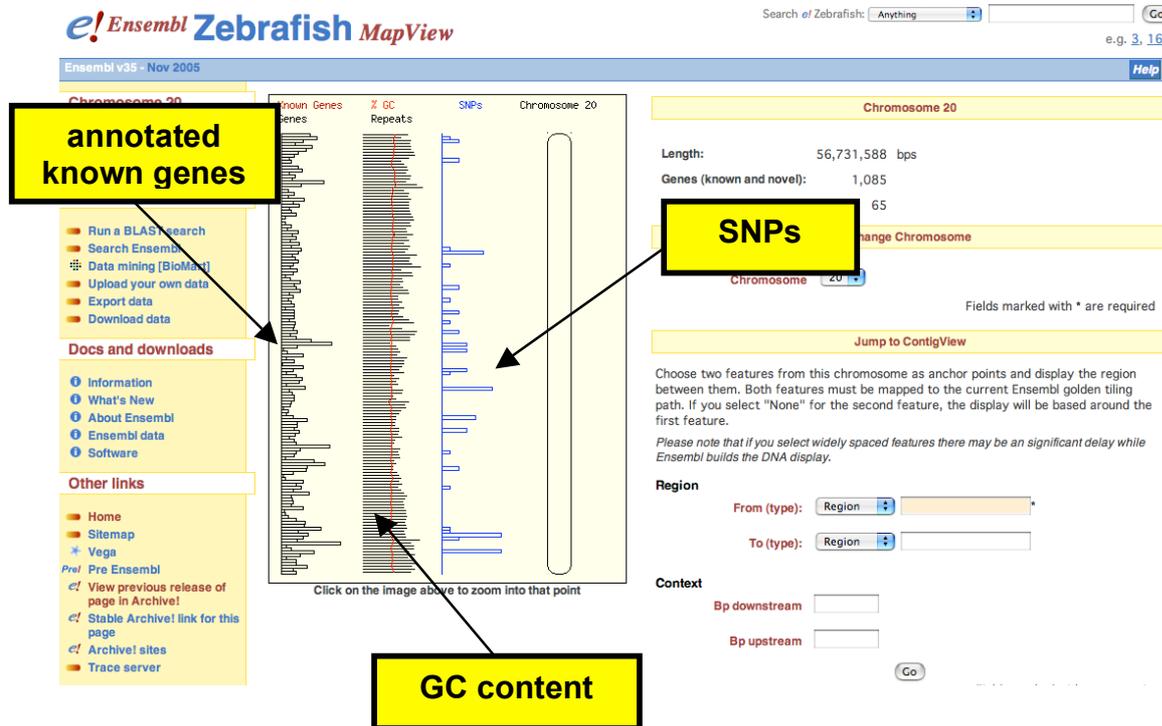
From the main Ensembl site you can access the zebrafish site by clicking on the appropriate species button. As soon a new assembly is released the sequence is made available as a pre-Ensembl site. This includes valuable information such as EST and UniProt alignments and *ab initio* predictions. The main missing data are the Ensembl genes and Ensembl ESTgenes. A full Ensembl dataset for a new assembly is typically made public a couple of months after the assembly release date.



## MapView and ContigView

This zebrafish Ensembl page provides various access points to the assembly sequence. For example you can browse a particular chromosome. The

chromosomes are linked to the **MapView** pages. The figure below shows the MapView for chromosome 20.



A MapView page plots the gene and SNP density and GC content. From this page you can zoom in to a more detailed display called ContigView by clicking on the schematic figure representing the chromosome.

**ContigView** can be considered the central view of the Ensembl web site. It shows the fragments (contigs, clones, etc) that make up a genome assembly. It allows you to scroll along entire chromosomes, whilst viewing the annotated features within a selected region in detail.

A ContigView page is divided into four panels: a chromosome overview, a zoomed-in **overview** of the region in the chromosome you are browsing, a **detailed view** showing features and a **basepair view** that goes down to individual bases. In order to continue with this module, jump to the region under the accession BX004766 (in chromosome 20) with start coordinate 1 and end coordinate 200000. (Use the text box provided to enter these coordinates.)

_____

The Features menu in the detailed view controls the tracks you can visualise in the panel. Tracks can be turned on and off and the features can be collapsed to simplify the view. Spend some time on this page trying the different menus and studying the displayed features. Observe that there are two tracks for predicted genes: Ensembl transcripts and EST transcripts. (If these features are not visible verify that the corresponding tracks are selected in the menu.)

### GeneView, TransView, ExonView and ProtView

Another important view in Ensembl I are the **GeneView** pages with information about the  Ensembl predicted genes. In the ContigView page above there is a predicted transcript on the forward strand called **jag2**. Clicking on this transcript displays a pop-up window with several options. Follow the ink labelled  Ensembl Gene: ENSDARG00000021389. Below we only show the top of the GeneView page for jag2; scroll down to view all the information available.

GeneView provides annotation and supporting evidence for the selected gene. The annotation consists of transcripts, homologies to other species, known and predicted proteins and domains, and links to external documentation. In this example, jag2 is a gene known to ZFIN and so a link to the corresponding external page is provided. The annotation for jag2 is based on 2 transcripts. In the Transcripts sections there are links to the corresponding TransView pages. Click on the link labelled "Transcript info"  for the first one with identifier ENSDART00000024922.
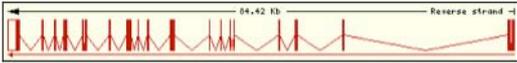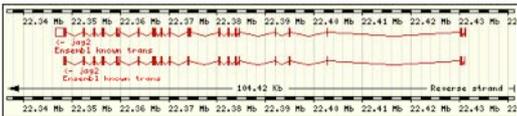
**TransView** provides annotation and supporting evidence for the selected transcript (structure, transcribed proteins, Gene Ontology and InterPro associated entries). The Transcript report panel provides a top-level summary of the transcript, with links to its genomic location, alignments to sequences in external databases, and export options. Underneath the report, the cDNA sequence of the transcript can be shown with codons, peptide sequence and/or SNPs highlighted.

From the GeneView page there are also links to the ExonViews labelled as "Exon info".



**ExonView** provides annotation and supporting evidence for the exons of a selected transcript. Ensembl gene predictions are based on aligned evidence from external databases like UniProt and RefSeq. At the bottom of an ExonView page you can find the evidence linked to this prediction.

Finally from the links labelled "Peptide info" in the GeneView page we can visit the ProtView page for the associated translation.



**ProtView** shows information about the structure and function of the encoded protein in the transcript's report with external links to various databases like Pfam, Prosite, etc…

_____

### ExportView

**ExportView** lets you download/dump data. All the features for a genomic region may be downloaded or exported to several formats (for example, FASTA, GenBank or EMBL-style flat file, as a feature list or an image). The ExportView pages are accessible from the link 'Export data'  in the left-hand side menu from any of the pages above.



### Zebrafish assembly in Ensembl

The sequence in the *Danio rerio* Ensembl database is the latest assembly release with automatic annotation. The genomic sequence released is based on all the sequenced clones with remaining gaps covered by contigs from a whole genome shotgun (WGS) assembly. The WGS fragments are placed in those gaps using a mixed strategy that looks at sequence similarity and other anchors as BAC-ends and  markers. This placement is hard to perform without errors - mainly due to the presence of mis-joins in the WGS assembly and duplicates. It is even more difficult to place sequence where there is no sequenced clone or marker to use as an anchor.

In this context the user has to evaluate the data with a critical eye. In particular when the sequence of interest is known to the community but it is wrong in the assembly. There are three kinds of scaffolds and these are, in order of quality from best to worst:

1.  scaffolds that have been attached to chromosomes (they may contain sequenced clones),
2.  scaffolds that can be aligned to clones but the physical map cannot assign a chromosome yet (they may contain sequenced clones), and
3.  NA (non-attached) scaffolds that corresponds to WGS contigs that could not be placed in the map (they don not contain sequenced clones).

Zv5_scaffold1699 is an example of category 1 above. This scaffold is placed in chromosome 20.



In the detailed view for this page there is a genomic region labelled BX470265.3. This is the accession number of a sequenced clone. The region labelled Zv5_scaffold1699 is a WGS supercontig. This is of lower quality than the sequenced clone (and may contain gaps represented by a sequence of Ns).

Zv5_scaffold935 is an example of a region that is part of the map but, when the assembly was built, did not have a placement in a chromosome (category 2). This example shows that the region contains some sequenced clones as shown by the presence of their accession numbers.



Finally a scaffold from category 3 is Zv5_NA10. This region does not contain any finished clones.

_____

### **Exercises**

This section introduces the Ensembl browser and some of its basic views. In other section we will study more advanced features like the compara database and Blast/SSAHA search facilities. The user is encouraged to navigate the site and experiment with the different views discussed above.

1. Find the GeneView page for jag2 (Ensembl gene), and scroll down to the first 'Transcript/Translation Summary'. As jag2 has been identified in Zv6 you can use this gene name in a text search box.

2. Examine the genomic context. From GeneView, follow the link 'View gene in genomic location' to ContigView.

3. Customise the display of ContigView selecting different tracks and comparing the data from different tracks.

4. In ContigView zoom in to examine the data in more detail.

5. Export a file containing the cDNA of one of the predicted transcripts for jag2.

6. One of the Ensembl tracks displays probes for which ZFIN  has a expression pattern page. Search for the mapping of the EST with accession CK685476 and open the corresponding ContigView page. Make sure that the 'expression pattern' track is selected in the 'features' menu. The ContigView page displays a link to the expression pattern page in ZFIN, try it.