# Browsing Genes and Genomes with Ensembl
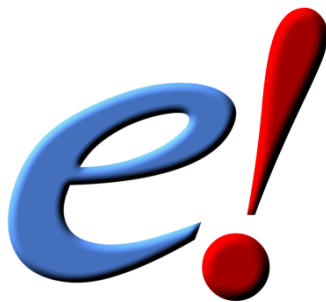


# Ensembl Browser Workshop

## Academia Sinica

## Taiwan
## 22nd February 2016

### Dr Denise Carvalho-Silva

## Notes:

This workshop is based on Ensembl release 83 (December 2015) and Ensembl Genomes release 30 (December 2015)

Some useful links:

1) Ensembl Browser website
www.ensembl.org

2) Ensembl Genomes Browser website
www.ensemblgenomes.org

3) Workshop materials (in pdf)
http://www.ebi.ac.uk/~denise/taiwan/sinica


Feel free to tackle questions relative to your own research instead of following the ones provided in our course booklet. The answers for the latter can be found here:

http://www.ebi.ac.uk/~denise/taiwan/answers

More exercises? http://tinyurl.com/e-exercises

Questions or comments?

helpdesk@ensembl.org
helpdesk@ensemblgenomes.org

# TABLE OF CONTENTS

## OVERVIEW

Ensembl provides annotation of genes and other genomic features such as sequence variants, conserved regions across species, and regulatory regions. The Ensembl gene set of the human, mouse and zebrafish genomes is based on protein and nucleotide evidence annotated by both automatic and manual means (the latter carried out by the **Havana** group).

All the data are freely available and can be accessed via the web browser or programmatically via our APIs (PERL or otherwise). Gene sequences can be downloaded from the Ensembl browser itself, or through the use of the **BioMart** web interface, which can extract information from the Ensembl databases without the need for programming knowledge.

Points covered:
- The need for genome browsers
- An introduction to the Ensembl browser
- Accessing genomic data in Ensembl
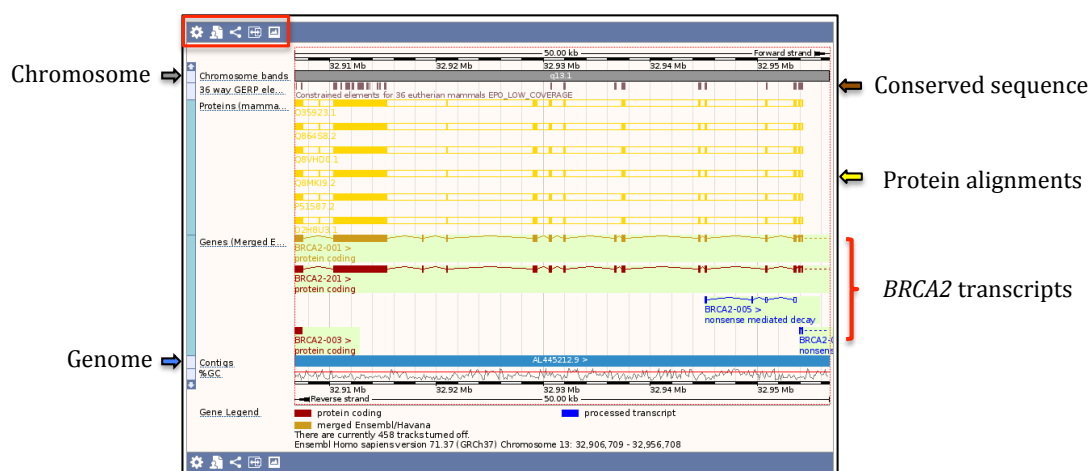- An overview of Ensembl tools

Check our video tutorial!
http://www.youtube.com/user/EnsemblHelpdesk

The Ensembl Genome browser
Introduction to BioMart

# INTRODUCTION TO ENSEMBL

Ensembl is a joint project between the EBI (European Bioinformatics Institute) and the Wellcome Trust Sanger Institute that annotates **chordate** genomes (i.e. vertebrates and closely related invertebrates with a notochord, such as sea squirt). Gene sets from model organisms (e.g. yeast, fruitfly and worm) are also imported for comparative analysis by the Ensembl Comparative Genomics team. Most annotations are updated every two months, leading to increasing Ensembl versions (such as version 83), however the gene sets are determined less frequently.
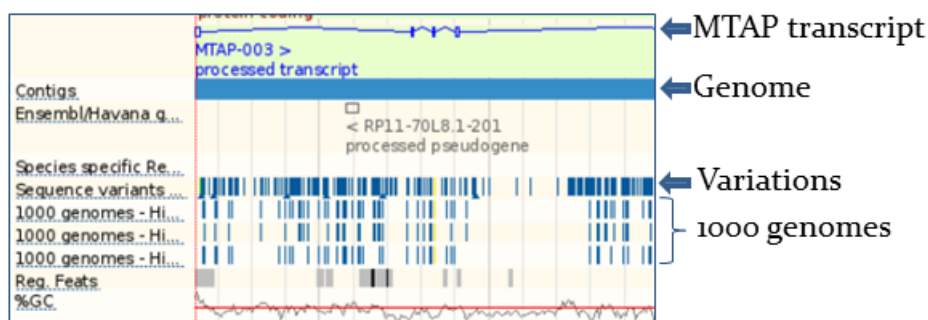


Click on the cog wheel icon ⚙ to add more data tracks to Ensembl views. Alternatively, you may want to click on the *Configure this page* button ⚙ Configure this page instead at the left hand side.

The vast amount of information associated with genomic sequences demands a way to organise and access that information. This is where genome browsers come in. Ensembl strives to display many layers of genome annotation into a simplified view for the ease of the user. Figure above shows the *Region in detail* page for the *BRCA2* gene in human. The example shows conserved sequence reflecting conservation on a base pair level across dozens of vertebrate species. Conserved regions are displayed as dark blocks that represent local regions of alignment.

In the above figure, proteins alignments from the UniProtKB have been added to the genomic location of the human *BRCA2* gene. Filled yellow blocks show where these UniProtKB proteins align to the genome, and gaps in the alignment are shown as empty yellow blocks. Note that the UniProtKB proteins support most of the exons shown in the Ensembl BRCA2-001 and BRCA2-201 transcripts.

Both **Ensembl** and **Havana** transcripts are displayed as exons (boxes) and introns (connecting lines). Filled boxes show coding sequence and empty boxes reflect **Un**Translated **R**egions (**UTRs**). This *Region in detail* view is useful for comparing Ensembl gene models with current proteins, mRNAs and ESTs from other databases, such as NCBI RefSeq, ENA, Unigene and UniProtKB. Everything in this view is aligned to the genome.



The *Region in detail* view can be configured (using the *Configure this page* button) to show regulatory features, sequence variation, and more! For example, click on dbSNP under the *Variation* menu and turn on the sequence variants (e.g. dbSNP) at the right side of this page. Save and close. Back to the Region in detail view, click on the sequence variation of interest. A pop-up box will show you a few variation properties, such as rs number, alleles, type, and others. Click on the *rs properties* link to take you to an information page for the genetic variation, including links to population frequencies, if available. You can do the same for regulatory features as well.

An index page is provided for each species with information about the source of the genomic sequence assembly, a karyotype (if available), and a link to our archives, which contain previous versions of the Ensembl Browser. The picture below shows the Ensembl homepage for human. Links to the human karyotype, to the previous

human assembly and a summary of gene and genome information are found in this index page.
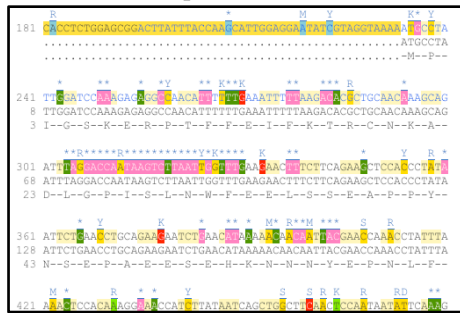


Ensembl uses a tabbed structure for separate pages and views in the browser to display different types of information. The three main entry points in Ensembl are the *Location*, *Gene* and *Transcript* tabs. We have also the *Species*, *Variation* and *Regulation* tabs.
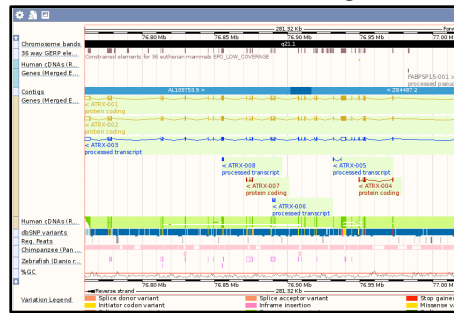


You can for example change the species of interest in the *Species* tab, view a chromosomal region in the *Location* tab (also known as *Region in detail*), visualise gene trees in the *Gene* tab, browse the cDNA sequence alongside the protein translation in the *Transcript* tab, look for genotype information in the *Variation* tab and get regulatory information in the *Regulation* tab. You can also perform a similarity search against any species in Ensembl by using our BLAST/BLAT tools.
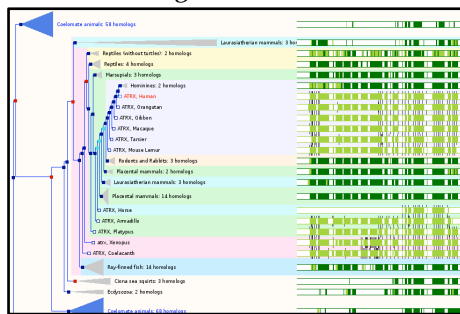
Transcript with variations



Genes and conserved regions



Homologues in Gene Trees



BLAST and BLAT searches



Let's now take a look at the Ensembl Genomes browser:
www.ensemblgenomes.org



Links to the taxa-specific sites

Link back to Ensembl

News

Click on the different taxa to see their homepages. Each of them is colour-coded as it follows:

| Protists | Fungi | Metazoa | Plants | Bacteria |
|----------|-------|---------|--------|----------|











You can navigate most of the taxa in the same way as you do with Ensembl, but since Ensembl Bacteria has a large number of genomes (>20,000), it needs slightly different methods for browser navigation. Let's look at this in more detail.

There's no full species list for bacteria as it would be hard to navigate with the number of species available in Bacteria. To find a species, start to type the species name (or any three letter code) into the species search box. A drop down list will appear with possible species

For example, to find substrains of *Clostridium difficile* type in Clostridium d to find the possible options.

The drop down contains various strains of *Clostridium difficile*. Let's choose Clostridium difficile 630. This will take us to another page, the species homepage, where we can explore various features of *C. difficile*.



## Retrieving Data from Ensembl

**BioMart** is a web-interface that can extract information from the Ensembl databases and present the user with a table of information without the need for programming. It can be used to output sequences or tables of genes along with gene positions (chromosome and base pair locations), single nucleotide polymorphisms (SNPs), homologues, and other annotation in HTML, text, or Microsoft Excel format. BioMart can also translate one type of ID to another, identify genes associated with **InterPro** domains or gene ontology (**GO**) terms, export gene expression data and much more.

Ensembl uses MySQL relational databases to store its information. A comprehensive set of Application Program Interfaces (APIs) serve as a middle-layer between underlying database schemes and more specific application programs. The API aims to encapsulate the database layout by providing efficient high-level access to data tables and isolate applications from data layout changes.

**Synopsis: what can I do with Ensembl?**

- View genes with different annotations along the chromosome;
- View alternative transcripts (i.e. splice variants) for a given gene;
- Explore homologues and phylogenetic trees across more than 70 chordate species for any gene;
- Compare whole genome alignments and conserved regions across species;
- View microarray sequences that match to Ensembl genes;
- View ESTs, clones, mRNA and proteins for any chromosomal region;
- Examine single nucleotide polymorphisms (SNPs) for a gene or chromosomal region;
- View SNPs across strains (rat, mouse) and human populations;
- View positions and sequence of mRNAs and proteins that align with an Ensembl genes;
- Display your own data on the Ensembl browser;
- Use BLAST or BLAT against any Ensembl genome;
- Export sequence or create a table of gene information with BioMart;
- Determine how your variants affect genes and transcripts using the Variant Effect Predictor;
- Share Ensembl views with your colleagues and collaborators;
- Retrieve our data using the Perl or REST APIs.

## Need more help?

- ❓ Check Ensembl [documentation](#)

- ❓ Watch [video tutorials](#) on YouTube

- ❓ View the [FAQs](#)

- ❓ Try some [exercises](#)

- ❓ Read some [publications](#)

- ❓ Go to our [online course](#)

## Stay in touch!

- ❖ Comments/questions [http://www.ensembl.org/Help/Contact](http://www.ensembl.org/Help/Contact)

- ❖ Read our Ensembl [blog](#)

- ❖ Follow us on Twitter [@ensembl](#) [@ensemblgenomes](#)

- ❖ Sign up to our [mailing lists](#)

## Further reading

Yates, A. *et al.*
**Ensembl 2016**
Nucleic Acids Res (Database Issue)
http://nar.oxfordjournals.org/content/early/2015/12/19/nar.gkv1157.full

Kersey, PJ. *et al*
Nucleic Acids Res (Database Issue)

**For a complete list of publications, see below**
http://www.ensembl.org/info/about/publications.html
http://ensemblgenomes.org/info/publications

# BROWSER WALKTHROUGH

We will guide you through the website using the human *ESPN* gene. This gene encodes a multifunctional actin-bundling protein with a major role in mediating sensory transduction in various mechanosensory and chemosensory cells. Mutations in this gene are associated with deafness (www.uniprot.org/uniprot/B1AK53).

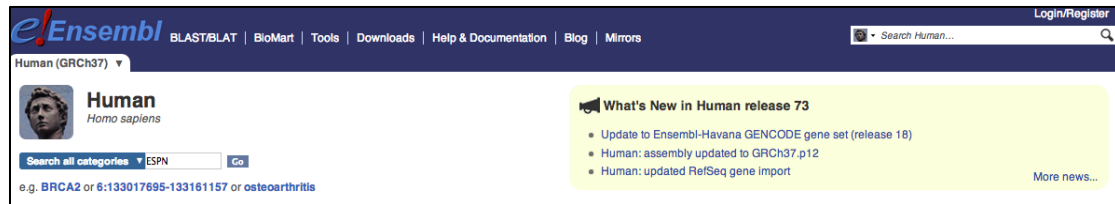The following points will be addressed during the walkthrough:

- **The Location tab and genomic location related links:**
    - How do I zoom out to change the gene focus?
    - How to add a track (e.g. protein alignments, variation data)?

- **The Gene tab and gene related links:**
    - Can I view the genomic sequence of my gene with its variations?
    - How to find orthologues and paralogues?

- **The Transcript tab and transcript related links:**
    - What is the protein sequence?
    - What proteins and mRNAs are found in other databases?

- **Exporting a sequence and running BLAT**

Go to www.ensembl.org and click on the human icon to open the human home page.
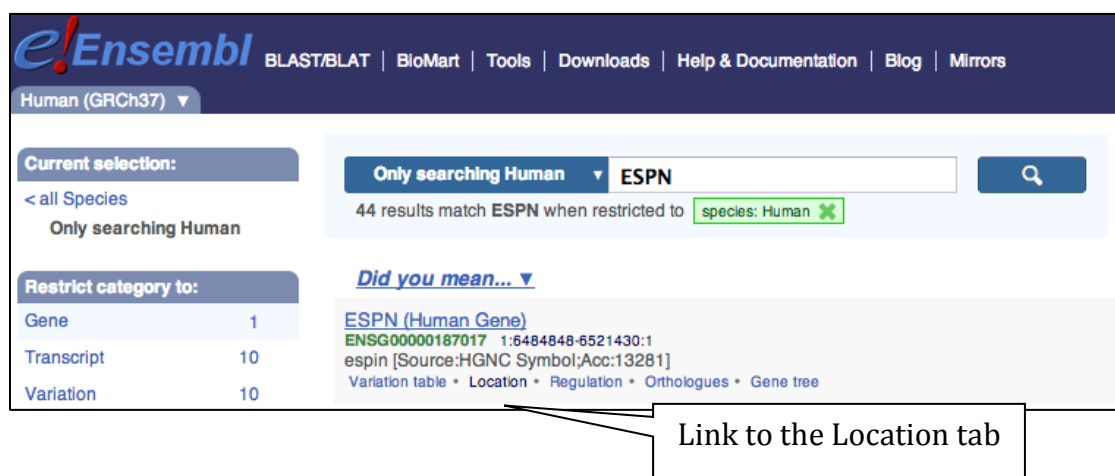
Type *ESPN* into the search bar and click the *Go* button.



One gene matches the query in human. Links to the *Gene* tab, *Variation* Table, *Location* tab, among others are provided.
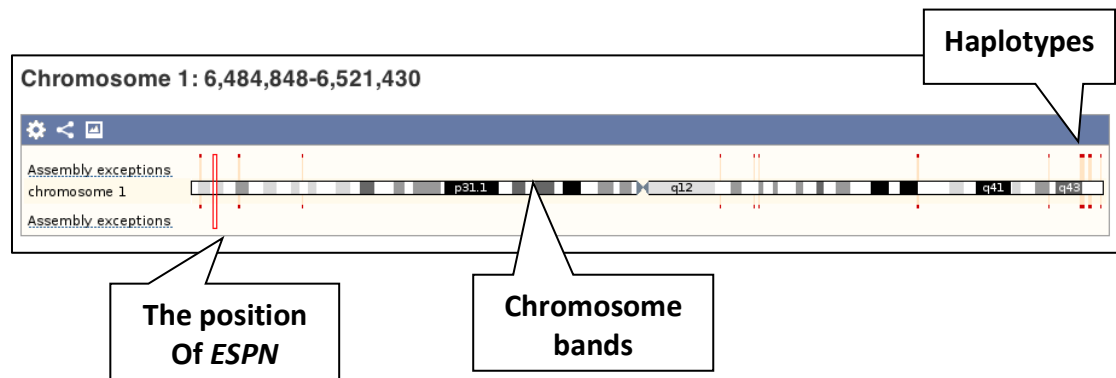


Let's view the genomic region in which this gene is located by clicking on the Location link. The Location tab should open.



The Location tab in Ensembl is also known as *Region in detail* view. There is a help video on this page at http://youtu.be/tTKEvgPUq94.
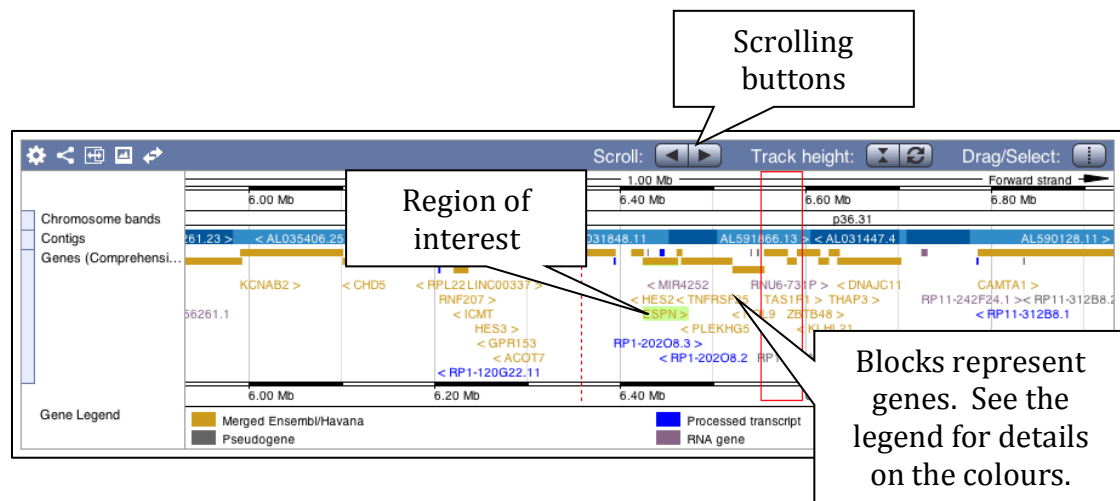
The *Location* tab contains three images.

The first image shows an overview of chromosome where the human *ESPN* gene has been mapped. In this image, you can see G banding pattern of the chromosome as well as the regions that correspond to the patches and haplotypes of the human chromosome 1.

**Haplotypes**

**The position Of *ESPN***

**Chromosome bands**

More details on haplotypes (e.g. MHC) can be found in the link below:

http://www.ensembl.org/Multi/Help/Movie?db=core;id=372

The second image shows a 1Mb region around the *ESPN* gene. This view allows scrolling back and forth along the chromosome. This view can also be configurable/customised.



**Scrolling buttons**

**Region of interest**

**Blocks represent genes. See the legend for details on the colours.**

At the moment the gene track is set to a fixed height. Click on the Automatic track height button to expand the image to include all possible data in the track.

Scroll along the chromosome by clicking and dragging within the image. As you do this, you'll see the image below grey out and two blue buttons appear. Clicking on *Update this image* will jump the lower image to the region central to the scrollable image. If you want to go back to where you started, click on Reset scrollable image.
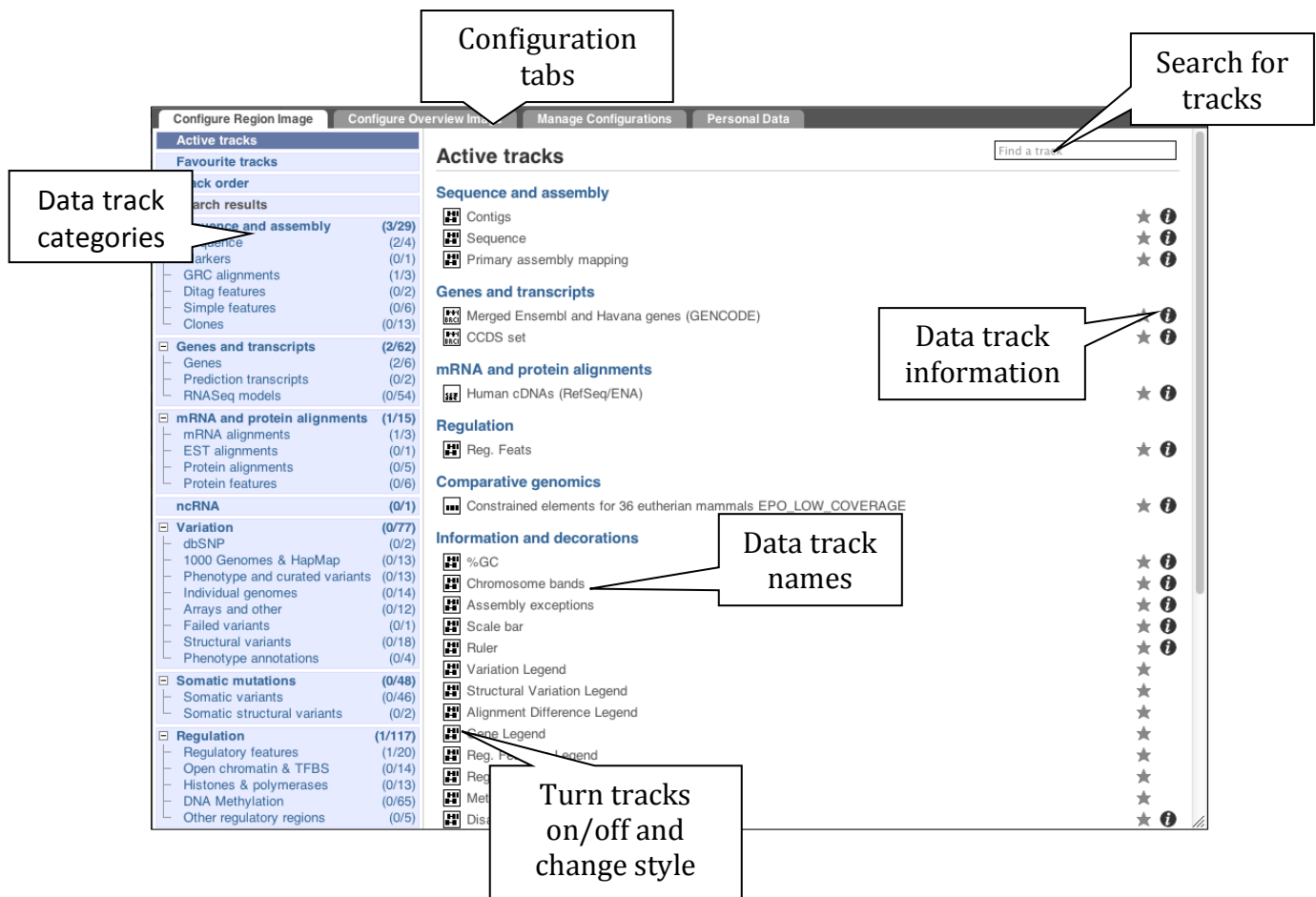
The third and final image is a detailed and highly configurable view of the region.

You can edit what you see on this page by clicking on the blue Configure this page menu at the left or click on the *cog wheel* in the bottom image.



This will open a menu that allows you to change the image and looks like the image below:
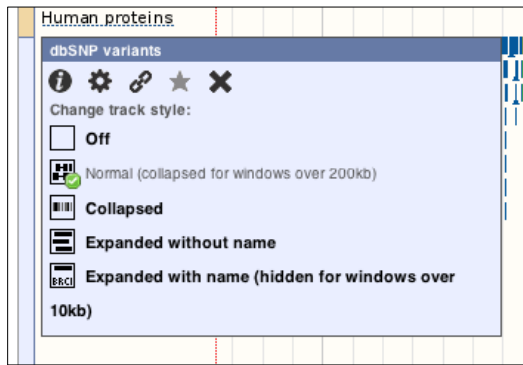
Let's add some tracks to this image, such as:
- Human ESTs – Labels

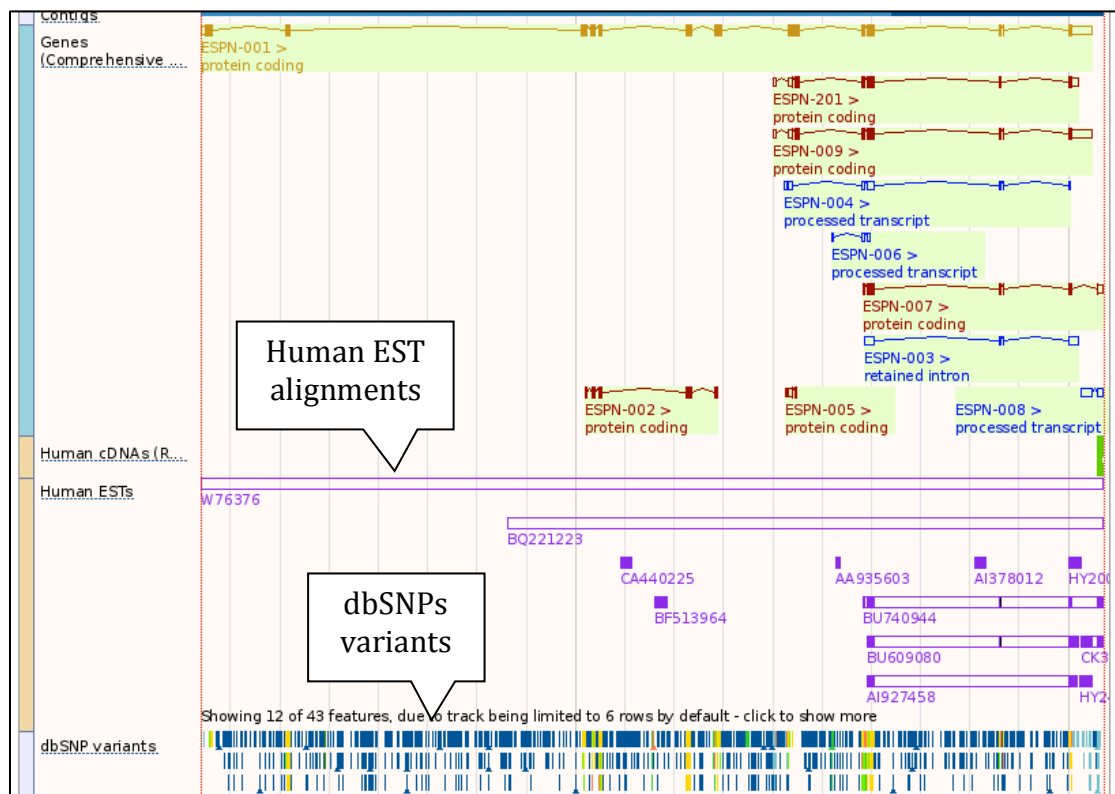- dbSNP variants – Normal

- 1000 Genomes – AMR – Collapsed

The data tracks can be displayed on different styles. For more details, have a look at our FAQ: http://www.ensembl.org/Help/Faq?id=335.

Now click on the tick in the top left hand to SAVE and close the menu with the newly added configuration. Alternatively, click anywhere outside of the menu.

We can also change the way the tracks appear by hovering over the track name, then the cog wheel to open a menu.

Have a look at the changes in the Location tab (sometimes referred to as *Region in detail* view). Click and drag tracks to reorder them, if it helps with comparing the data.



Now that you've got the view how you want it, you might like to show it to a colleague or collaborator. Click on the *Share this page* button to generate a link so that you can email it to someone.



They will see exactly the same view as you, including all the tracks you have added and moved up and/or down. These links contain the

Ensembl release number, so if a new release or even assembly comes out, your link will just take you to the archive site for the release it was made on.

To return to the default view, click on *Configure this page* and select Reset configuration at the bottom of the menu. You can also reset the track order.
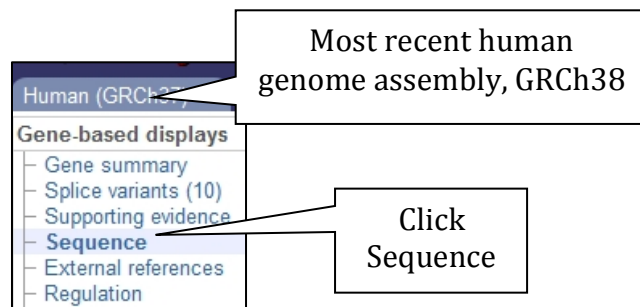
Let's now explore the Gene tab.



We will walk you through some of the links in the left hand navigation column. Note that the left hand side menu in the Gene tab differs from the one we saw previously in the Location tab.

How can we view the genomic sequence of the *ESPN* gene?

Click Sequence at the left of the page.



Most recent human genome assembly, GRCh38

Click Sequence

**Marked-up sequence**

Page-specific help

Download sequence

**Key**

Exons    All exons in this region    ESPN exons

>chromosome:GRCh38:1:6424188:6461970:1
AGCGCACTAGTGGTTCCCGTCCTGCCTTCCTGCCCTCCCCGC...GGCAG
TTCTCGGATCTGGAGGGACCCTGGAAGGCAGGGCTCTTTGCA...TTCGAC
CCAGAGCCCTTCAGGGACGTGGCAGGGCTGCTCCTGCCTCAGGG...GTTGTCCTCGTGCT
CCTCACCCCGCCTGGAATACCCTTCTCGCCGCTCAAACCCAGCCCCACGGCACCTCCTCA
GAGACCTTTCCCTGTCCGCCCACGCGGTCCCGACAATCACTCCCCATCAC CTCTGGAATT
GCGTCGCCGGCGCCTGGAACCGCAGTTAGCGGGCACTGGGCAGATGAATGA...
TGCCTGGACGGCTCTCCAATTCGAACCCAGTTTTGCTGCCCTCTGGGGTC...
CGTGAGGCAAATTAGGAGAGAAGCCCCTGGGCACCTTGCCCCAGTCGCACG
GCG TCGCGGCGGGGGCGGGCGGGGAACTCGGGCGGAGGCTGCGGGGCGGGG
GTGGGGGCGGGCCCGAGTCTTAAGCCGGCGTCCGCGGGCTCCGGCCCCAGAGCGCGCGGCGG
AGCGGAGCGCCAGGCAGCGCGGAGCGGAGGCCAGGCCCACAGCCGCTCCGCCTCCCGGCC
CGCAGATCCCCGACGGCCGCACCGCGGGCTCCTCTGGCCCGCAAGAACACGTGCATGGCC
TCCTGGGGAAGGCGCTGAGTGCGGAGTCGCGGCGCCGCACGCGGCACCATGGCCCTGGAG
CAGGCGCTGCAGGCGGCGCGGCAGGGCGAGCTGGACGTGCTGAGGTCGCTGCACGCCGCA
GGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGCGCTGCCCGTGCACCACGCGGCC
CGCGCTGGGAAGCTGCACTGTCTGCGCTTCCTGGTGGAGGAAGCCGCCCTCCCCGCCGCG
GCCCGCGCCCGCAACGGCGCCACACCGGCCCACGACGCCTCCGCCACCGGCCACCTCGCC
TGCCTGCAGTGGCTGCTGTCG AGGGCGGCTGCAGAGTGCAG GTGGGTCCGCGCGGTTCG
CCAGGGGCACTGAGGCTTCC...TCAGGACA GAGTCCTGGCCCAGAGTCCCCGGGGCTC
AAGGATGGGTGGGGTTT...CAGCTGAACCCTGCACGGAGCTCCTTCCA
GAGGCCCTCAAGTGAAT...GCCAGTACTGGGGCAGATGCCCTGGCGAG
CCTGGGTGCTCCCTGGAAGCGCACCTGGGTGATGGGAGCCAGAAGGGAGGGGCCTCCGTG
...

Upstream sequence
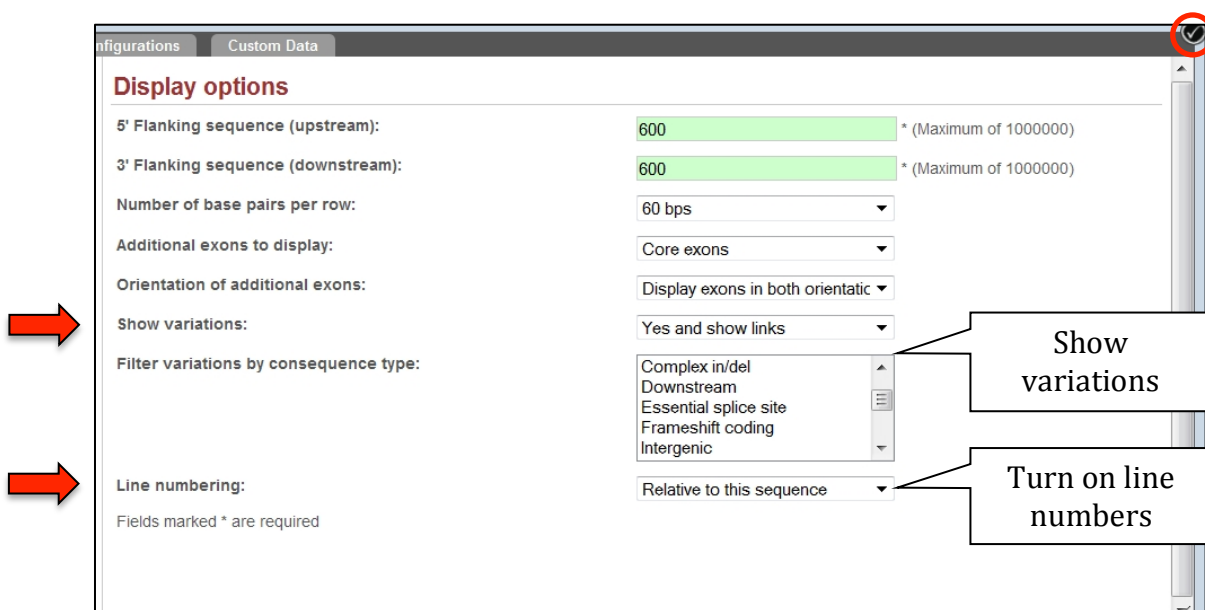
Exon of an overlapping gene

*ESPN* exon

Click on the button  *i*  to view page-specific help.

The help pages also provide links to Frequently Asked Questions, a Glossary, Video Tutorials, and a form to Contact HelpDesk.

The sequence is shown in FASTA format. Take a look at the FASTA header:



Exons are highlighted within the genomic sequence by default. Variations and other features can be added with the Configure this page link found at the left. Click on it now.



Once you have selected the new display options (in this example, Show variations and Line numbering) click at the top right of the image (circled in red above).

Links to the variation tab

To view all the sequence variations in this locus, click the Variation table link at the left of the gene tab.



The table is divided into consequence types.

Click on Show to expand a detailed table for any of the consequence types available such as Missense variants:



Let's have a look at some of the regulatory data that is available for the *ESPN* gene in human. You are still in the Gene tab.

In the left hand side menu of the Gene page, click on the *Regulation* link.

This page shows all the available data from the Regulatory build of Ensembl. The source of this data is mainly from the ENCODE project.

This view can also be configured to show cell specific data. Lets choose to display data for GM12878 and K562 cells only.

More information on the sources of this data (including what the cells are) can be found below:

http://www.ensembl.org/info/genome/funcgen/regulation_sources.html



The different blocks show the regions where the available biochemical data (e.g. ChIP-sequencing to assess Transcription Factor Binding Sites, or TFBS, and modified histones) maps to on the human genome. Check the different colours in the legend to find out more about this annotation e.g. Enhancer, Promoter, etc…

We provide a summary of these data across 18 cells in humans (MultiCell) and five cells in mouse as part of our Regulatory build:

http://www.ensembl.org/info/genome/funcgen/regulatory_build.html

Click on the regulatory features ( ) to learn more.

**Black lines indicate TF motifs**



**Regulatory Feature - MultiCell**

| | |
|---|---|
| Stable ID | ENSR00000529720 |
| Type | Promoter |
| Core bp | 1:6445384-6450983 |

**Attributes** -

**Motif Information**

| Name | PWM ID | Score |
|---|---|---|
| CTCF | MA0139.1 | 9.37 |
| Gabp | PB0020.1 | 6.641 |
| Gabp | PB0124.1 | 6.802 |
| Egr1 | MA0341.1 | 6.44 |
| Egr1 | MA0162.2 | 10.688 |
| Egr1 | PB0010.1 | 8.538 |
| Egr1 | MA0337.1 | 7.398 |
| Egr1 | MA0366.1 | 6.678 |
| CTCF | MA0139.1 | 9.156 |

**Regulatory feature ID, ENSR00000529720. Click on it to go to the Regulation tab**

**TF motifs within Regulatory feature ENSR00000529720**

Note the different colours to indicate specific annotation at the Gene Regulation level: regions can be associated with promoter, enhancer, CTCF, TFBS, open chromatin activity, etc. See the legend for more details:



We also provide cell–specific data (e.g. Reg. Feats for CD4, Reg. Feats GM12878, etc). Click on 'Configure this page' to add cells specific data.

Regulatory features annotated by Ensembl based on ChIP seq data and others

Inactive Regulatory feature

In addition to the Regulatory Features, Segmentation data from the ENCODE project can be displayed in this view. See the legend for more details on the different colours for this data track. For example, dark yellow represent regions annotated as predicted enhancers.



Segmentation data from the ENCODE project

Details on the Segmentation data in Ensembl can be found below:

http://www.ensembl.org/info/genome/funcgen/regulatory_segmentation.html

Let's now have a look at some of the views on Comparative Genomics in the Ensembl Browser. In the left hand menu in the Gene tab, click on the 'Gene tree' link. This will display the Ensembl gene in the context of a phylogenetic tree used to determine orthologues and paralogues.

Click the *Orthologues* link at the left of this page to view homologues detected by this tree. Note the links under Compare.



Let's have a look at the pairwise alignment between human and mouse. Click on *Genomic alignments*.

You can download the alignment in multiple formats:



Other comparative genomics views are available from the Location tab such as Region comparison.

Click on the link in the left hand menu and choose the species to compare against my human *ESPN* gene



for example mouse



This is what you will see:

You can zoom in and out in this view and configure it to show a blue line (under Comparative features) to connect the orthologous genes. Tick the 'Join genes' box:



Save and close the new configuration to view an image like the one below:

Click on the ID for ESPN-001 (ENST00000377828).

You are now in the Transcript tab for ESPN-001. The left hand navigation column provides several options for the transcript ESPN-001. Click on the *Exons* link.





You may want to change the display (to show more flanking sequence or to show full introns, for example). In order to do so, click on *Configure this page* and change the display options accordingly.

If you would like to export the sequence either as FASTA or RTF (Rich Text Format), click on





Now click on the *cDNA* link to see the spliced transcript sequence.

## cDNA sequence ⓘ

**Key**

| | | |
|---|---|---|
| Codons | Alternating codons | Alternating codons |
| Exons | Alternating exons | Alternating exons |
| Variations | 3 prime UTR | 5 prime UTR | Coding sequence |
| | Missense | Splice region | Stop retained |
| | Synonymous | | |
| Other features | UTR | | |

```
                                                    Y  R
      1  AGCGGAGCGCCAGGCAGCGCGGAGCGGAGGCCAGGCCCACAGCCGCTCCGCCTCCCGGCC
         ............................................................
         ............................................................

     61  CGCAGATCCCCGACGGCCGCACCGCGGGCTCCTCTGGCCCGCAAGAACACGTGCATGGCG
         ............................................................
         ............................................................

                                                 R
    121  TCCTGGGGAAGGCGCTGAGTGCGGAGTCGCGGCGCCGCACGCGGCACCATGGCCCTGGAG
         ....................................................ATGGCCCTGGAG
         .......................................................-M--A--L--E-

    181  CAGGCGCTGCAGGCGGCGCGGCAGGGCGAGCTGGACGTGCTGAGGTCGCTGCACGCCGCA
     13  CAGGCGCTGCAGGCGGCGCGGCAGGGCGAGCTGGACGTGCTGAGGTCGCTGCACGCCGCA
      5  -Q--A--L--Q--A--A--R--Q--G--E--L--D--V--L--R--S--L--H--A--A-

    241  GGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGCGCTGCCCGTGCACCACGCGGCC
     73  GGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGCGCTGCCCGTGCACCACGCGGCC
     25  -G--L--L--G--P--S--L--R--D--P--L--D--A--L--P--V--H--H--A--A-
```

UnTranslated Regions (UTRs) are highlighted in yellow, codons are highlighted in light yellow, and exon sequence is shown in black or blue letters to show exon divides. Sequence variants are represented by highlighted nucleotides and clickable IUPAC codes are above the sequence.

Next, follow the *General identifiers* link at the left. This page shows data from other resources, such as HGNC, RefSeq, CCDS, EntrezGene, OMIM, UniProtKB, and others, that match to the Ensembl transcript and protein.

**Transcript-based displays**
- Transcript summary
- Supporting evidence (20)
- Sequence
  - Exons (13)
  - cDNA
  - Protein
- External References
  - **General identifiers (21)**
  - Oligo probes (48)

Links to matches to this Ensembl *ESPN* transcript in other databases

**General identifiers** ℹ️

This transcript corresponds to the following database identifiers:

Show [All ▼] entries                    [Filter        ]  📊

| External database | Database identifier |
|---|---|
| HGNC Symbol | ESPN<br>espin [view all locations] |
| UniParc | UPI000013D2B6 [view all locations] |
| CCDS | CCDS70.1 [view all locations] |
| UniProtKB/Swiss-Prot | ESPN_HUMAN [align]<br>Espin [view all locations] |
| RefSeq peptide | NP_113663.2 [Target %id: 100; Query %id: 100] [align]<br>espin [view all locations] |
| RefSeq mRNA | NM_031475.2 [align] [view all locations] |
| UCSC Stable ID | uc001amy.3 [view all locations] |
| Human Protein Atlas | HPA028674 [view all locations]<br>HPA028674 [view all locations] |
| European Nucleotide Archive | AF134401 [align] [view all locations]<br>AL031848 [align] [view all locations]<br>AL136880 [align] [view all locations]<br>AL158217 [align] [view all locations]<br>AY203958 [align] [view all locations]<br>CH471130 [align] [view all locations] |
| HGNC transcript name | ESPN-001<br>espin [view all locations] |

Now, click on Ontology table to see GO terms from the Gene Ontology consortium (www.geneontology.org).

Click on the ![info icon] to see a guide to the three-letter Evidence codes.

Click on *Protein summary* to view domains from Pfam, PROSITE, Superfamily, InterPro, and more.





Clicking on *Domains & features* shows a table of this information.

| | | | | | |
|---|---|---|---|---|---|
| Prosite_profiles | 768 | 825 | - | PS50313 | - |
| Smart | 35 | 64 | Ankyrin_rpt | SM00248 | IPR002110 [Display all genes with this domain] |
| Prosite_profiles | 69 | 93 | Ankyrin_rpt | PS50088 | IPR002110 [Display all genes with this domain] |
| Smart | 69 | 99 | Ankyrin_rpt | SM00248 | IPR002110 [Display all genes with this domain] |
| Prosite_profiles | 103 | 127 | Ankyrin_rpt | PS50088 | IPR002110 [Display all genes with this domain] |

Our last task is to export genomic sequence and perform a similarity search using the new Ensembl BLAST/BLAT tool.

You can export a sequence from different views in Ensembl but lets go back to the Location tab.

Click on the Export data option, select the default parameters (e.g. Fasta sequence as output) and click Next then HTML. Please note that you can export the genomic sequence as either unmasked or masked (soft or hard masked).

This is a snapshot what you will see in your browser:

```
>1 dna:chromosome chromosome:GRCh37:1:6484848:6521430:1
AGCGGAGCGCCAGGCAGCGCGGAGCGGAGGCCAGGCCCACAGCCGCTCCGCCTCCCGGCCCGCAGATCCC
CGACGGCCGCACCGCGGGCTCCTCTGGCCCGCAAGAACACGTGCATGGCGTCCTGGGGAAGGCGCTGAG
TGCGGAGTCGCGGCGCCGCACGCGGCACCATGGCCCTGGAGCAGGCGCTGCAGGCGGCGCGGCAGGGCG
AGCTGGACGTGCTGAGGTCGCTGCACGCCGCAGGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGC
GCTGCCCGTGCACCACGCGGCCCGCGCTGGGAAGCTGCACTGTCTGCGCTTCCTGGTGGAGGAAGCCGCC
CTCCCCGCCGCGGCCCGCGCCCGCAACGGCGCCACACCGGCCCACGACGCCTCCGCCACCGGCCACCTCGC
CTGCCTGCAGTGGCTGCTGTCGCAGGGCGGCTGCAGAGTGCAGGTGGGTCCGCGCGGTTCGCCAGGGGC
ACTGAGGCTTCCTCCTCAGGACAGAGTCCTGGCCCAGAGTCCCCCGGGGCTCAAGGATGGGTGGGGTTT
GGCACCTCCTGGCCCAGCTGAACCCTGCACGGAGCTCCTTCCAGAGGCCCTCAAGTGAATGGGCTCCCTG
GCTTGCCAGTACTGGGGCAGATGCCCTGGCGAGCCTGGGTGCTCCCTGGAAGCGCACCTGGGTGATGGG
AGCCAGAAGGGAGGGGCCTCCGTG
```

To use this sequence for similarity searches, you can select the header and a few lines of the nucleotide sequence and then copy it on the clipboard. Click on the BLAST/BLAT link in the bar at the top of the main Ensembl homepage. Paste the sequence into the appropriate box and select BLAT as the search algorithm.

Click Run.

The table with results should look like these examples. The jobs are given a ticket number and highlighted as green when successfully completed.

Click on view results for a table or karyotype view.

In the Results table, click on the Genomic location link to view the BLAST/BLAT hit in the Region in detail (Locatin tab).



Click on the red bar for some summary statistics such as the score, %ID, and other BLAST/BLAT values.

Export Image for your lab notebook or publications, or Share it with your colleagues and collaborators!



<mark>END OF THE BROWSER WALKTHROUGH</mark>

# EXERCISES

## ENSEMBL BROWSER

### Exercise 1 – Exploring a genomic region in human

a) Go to the region from bp 31,873,863 to 32,623,863 on human chromosome 13. On which cytogenetic band is this region located? How many contigs make up this portion of the assembly (contigs are contiguous stretches of DNA sequence that have been assembled solely based on direct sequencing information)?

b) Zoom in on the *BRCA2* gene.

c) Are there any Tilepath clones (i.e. BAC clones upon which the current genomic assembly was based) that contain the complete *BRCA2* gene?

d) Add the track with RefSeq gene models. This track is known is Ensembl as 'RefSeq human import' Has RefSeq annotated the *BRCA2* gene? If so, how many transcripts have been annotated? Do they differ from the Ensembl transcripts?

e) Save a picture of the main panel so that you can use it in your publication (.png format).

f) Export the genomic sequence of the region you are looking at in FASTA format.

g) Delete all tracks you have added to the Location tab.

### Exercise 2 – Exploring a gene and its transcripts in human

a) Find the human *F9* (coagulation factor IX) gene. On which chromosome and which strand of the genome is this gene located? How many transcripts (splice variants) have been annotated for it?

b) What is the longest transcript? How long is the protein it encodes? Has this transcript been annotated automatically (by Ensembl) or manually (by Havana)? How many exons does it have? Are there exons completely or partially untranslated?

c) In which part (i.e. the N-terminal –start- or C-terminal –end) of the protein encoded by ENST00000218099 does its peptidase activity reside? Does the protein contain any transmembrane domains?

d) Have missense variants been discovered for the protein encoded by ENST00000218099?

e) What is the regulatory feature annotated in the region of this gene? Find its ID. Is this feature active in all cells Ensembl has got data for?

f) How many orthologues are predicted for this gene in rodents? How much sequence identity does the mouse protein have to the human one? Note the Target %id and Query %id. View the alignment (protein) between the two sequences.

g) Go to the orthologue in mouse and find the genomic alignment between mouse and human. Can you configure the view to show both START and STOP codons?

---

### ADDITIONAL EXERCISES: BROWSER

*If you have finished the exercises above, you may want to do these extra ones*

### Extra exercise 1 – Mouse assembly, protein domains

*a) What is the latest assembly of the mouse genome in Ensembl? When was it produced? What is the total number of non-coding genes in the current release? Can you find the previous assembly of the mouse genome (i.e. NCBIM37)?*

*b) Find the Brca1 gene in mouse. Does this gene have any transcript (s) with annotation that has been agreed between the EBI, Sanger, NCBI and UCSC? What is the protein sequence for this transcript? Are there any domains or features present in the amino acid sequence? Can you download this information as a table?*

## Extra exercise 2 – Exploring a gene in Ensembl Bacteria

*Start in http://bacteria.ensembl.org/index.html and select the Streptomyces lividans 1326 genome.*

*a) What GO: molecular function terms are associated with the 'era' gene?*

*b) What domains can be found in the protein product of this gene? How many different domain prediction methods agree with each of these domains?*

## Extra exercise 3 – miRNA genes in A. thaliana

*MicroRNAs (miRNAs) are small non-coding RNA molecules (ca. 22 nucleotides) found in plants and animals, which function in transcriptional and post-transcriptional regulation of gene expression. A well-studied miRNA family in plants is the MIR395 family (See also: http://en.wikipedia.org/wiki/MicroRNA and http://en.wikipedia.org/wiki/Mir-395_microRNA_precursor_family).*

*a) How many members does the MIR395 family in Arabidopsis thaliana have?*

*b) How are the MIR395 genes organised? Are they clustered? Are they all located on the same strand of the genome? How are they positioned relative to each other?*

# ENSEMBL TOOLS: BIOMART

*Follow these instructions to guide yourself through BioMart and answer the following query:*

You have three questions about this set of human genes: *ESPN, MYH9, USH1C, CHD7, CISD2, THRB, DFNB31*

Note: These are HGNC gene symbols.

1) What are the EntrezGene IDs for these genes?

2) Are there associated functions from the GO (gene ontology) project that might help describe their function?

3) What are their cDNA sequences?

**Step 1:** Click on BioMart in the top header of a www.ensembl.org page to go to: www.ensembl.org/biomart/martview



**STEP 2:**
Choose Ensembl Genes as the primary database.



**STEP 3:**
Choose *Homo sapiens* genes as the dataset.

**STEP 4:**
Click Filters at the left. Expand the GENE panel.

URL  XML  Perl  Help

Please restrict your query using criteria below

Filters
[None selected]
Attributes
Ensembl Gene ID
Ensembl Transcript ID
Dataset
[None Selected]

⊞ GENE:
⊞ TRANSCRIPT EVENT:
⊞ GENE ONTOLOGY:
⊞ EXPRESSION:
⊞ MULTI SPECIES COMPARISONS:
⊞ PROTEIN DOMAINS:
⊞ VARIATION:



New  Count  Results

URL  XML  Perl  Help

Please restrict your query using criteria below

**STEP 5:**
In Input external references ID list, paste in your gene symbols. Change the heading to read HGNC symbol(s).

genes ...
with ArrayExpress ID(s)
◉ Only
◯ Excluded

mit [Max 500 advised]
HGNC symbol(s) [e.g. ZFY]
ESPN, MYH9, USH1C, CHD7, CISD2, THRB, DFNB31

Dataset



New  Count  Results

Dataset 8 / 62352 Gene
Homo sapiens gen...
(GRCh37.p10)
Filters
HGNC symbol(s) [e.g. ...
[ID-list specified]
Attributes
Ensembl Gene ID
Ensembl Transcript ID

**STEP 6:**
Click Count to see BioMart is reading 7 genes out of 66,232 possible *H. sapiens* genes (this number includes ncRNA genes).

**STEP 7:**
Click on Attributes to select output options
(i.e. GO terms)

Features    Homologs
Structures    Variation
Transcript Event    Sequences

⊞ GENE:

⊞ EXTERNAL:

⊞ EXPRESSION:

⊞ PROTEIN DOMAINS:

HGNC symbol(s) [e.g. ZFY]:
[ID-list specified]
**Attributes**
Ensembl Gene ID
Ensembl Transcript ID

**STEP 8:**
Expand the EXTERNAL panel.

☐ LRG to Ensembl link transcript
☑ EntrezGene ID
☐ VEGA transcript ID(s) (OTTT)
☐ VEGA gene ID(s) (OTTG)
☐ Ensembl transcript (where OTTT shares CDS with ENST)
☐ HAVANA transcript (where ENST
☐ HAVANA transcript (where ENST
☐ HGNC ID(s)
☑ HGNC symbol

**STEP 9:**
Scroll down to select
EntrezGene ID
*(to answer question 1)*

**STEP 10:**
Also select HGNC symbol to see the input gene symbols we started with.

⊞ GENE:

⊟ EXTERNAL:
**GO**
☑ GO Term Accession
☑ GO Term Name
☑ GO Term Definition

**STEP 11:**
Scroll back up to select GO term fields
*(to answer question 2)*

**STEP 12:**
Click Results.

New    Count    Results    URL    XML    Perl    Help

Export all results to    File    TSV    ☐ Unique results only    Go

Email notification to

View    10 ▼ rows as HTML ▼    ☐ Unique results only

| Ensembl Gene ID | Chromosome Name | Associated Gene Name | EntrezGene ID | HGNC symbol | GO Term Accession | GO Term Name | GO Term Definition |
|---|---|---|---|---|---|---|---|
| ENSG00000187017 | 1 | ESPN | 83715 | ESPN | GO:0007605 | sensory perception of sound | "The series of events required for an organism to receive an auditory stimulus, convert it to a molecular signal, and recognize and characterize the signal. Sonic stimuli are detected in the form of vibrations and are processed to form a sound." [GOC:ai] |
| ENSG00000187017 | 1 | ESPN | 83715 | ESPN | GO:0007626 | locomotory behavior | "The specific movement from place to place of an organism in response to external or internal stimuli. Locomotion of a whole organism in a manner dependent upon some combination of that organism's internal state and external conditions." [GOC:dph] |
| ENSG00000187017 | 1 | ESPN | 83715 | ESPN | GO:0030046 | parallel actin filament bundle assembly | "Assembly of actin filament bundles in which the filaments are tightly packed (approximately 10-20 nm apart) and oriented with the same polarity." [GOC:mah, ISBN:0815316194] |

Ensembl Gene ID
Chromosome Name
Associated Gene Name
EntrezGene ID
HGNC symbol
GO Term Accession
GO Term Name
GO Term Definition

**Dataset**
[None Selected]

Why are there multiple rows for one gene ID? Look at the first few rows of your results table.

| Ensembl Gene ID | Ensembl Transcript ID | EntrezGene ID | GO Term Accession | GO Term Name | GO Term Definition | HGNC symbol |
|---|---|---|---|---|---|---|
| ENSG00000187017 | ENST00000377828 | 83715 | GO:0007605 | sensory perception of sound | "The series of events required for an organism to receive an auditory stimulus, convert it to a molecular signal, and recognize and characterize the signal. Sonic stimuli are detected in the form of vibrations and are processed to form a sound." [GOC:ai] | ESPN |
| ENSG00000187017 | ENST00000377828 | 83715 | GO:0007626 | locomotory behavior | "The specific movement from place to place of an organism in response to external or internal stimuli. Locomotion of a whole organism in a manner dependent upon some combination of that organism's internal state and external conditions." [GOC:dph] | ESPN |
| ENSG00000187017 | ENST00000377828 | 83715 | GO:0030046 | parallel actin filament bundle assembly | "Assembly of actin filament bundles in which the filaments are tightly packed (approximately 10-20 nm apart) and oriented with the same polarity." [GOC:mah, ISBN:0815316194] | ESPN |

**STEP 13:**
Click *Attributes* again

**STEP 14:**
Select Sequences at the top, then expand SEQUENCES and choose the option cDNA sequences (*to answer question 3*).

**STEP 15:**
Expand Header Information to select the Associated Gene Name

**STEP 16:**
Click Results to see the cDNA sequences in FASTA format.

>ENSG00000171316|ENST00000307121|CHD7
CGGCGGCGGCGGCAGCGGCGGCGGCGGCGGCGGCGGCGCGGGGGTTGAGTCGTGGTGGTGCGG
ACGCGCTCGTGCTCGGGAACTATCGGATTAAACTTGAATCGAGTGAAATTACACAAAGGA
GCGCCGCGGAGGAGCGGCCCGGGGACCCGGACACCCTGAAACTCACCAGAGACCCGTTC
GCCCCCGGCCAACTCCGTGCCCGTGGATTCAGCCCCCTGGCCGCAGCTGCCGAGCCAACT
CCGGAGCCCGCTCTGCGTTTTGTTTTCCCCTCGGCACTAGGCAGCGGAGGAGCCCGACCG
ACCCGGACCTATATCCAGACTTTGCCTGACACTGCAGGGTCCAAGAGAATTAAAGAAATA

**STEP 17:**
Change View **10** rows to View **All** rows so that you see the full table.

**Note**: Pop-up blocking must be switched off in your browser.

Note: you can use the Go button to export a file.

What did you learn about the human genes in this exercise?
Could you get these results using the Ensembl browser?
Would it take longer?

For more details on BioMart, have a look at these publications:

Smedley, D. *et al.*
BioMart – biological queries made easy
BMC Genomics 2009 Jan 14;10:22

Kinsella, R.J. *et al.*
Ensembl BioMarts: a hub for data retrieval across taxonomic space.
Database (Oxford) 2011:bar030

**Exercise 1 – Export sequences in FASTA format from mouse**

Retrieve the sequences of all mouse genes that are located on chromosome 17, that are protein coding and that encode for proteins containing transmembrane domains. Do a count after selection of each filter to check the number of genes remaining in your dataset. Export the results of the protein sequences (FASTA) as Compressed web file and get the results notified to you by email.

---

**Exercise 2 – Convert IDs of human genes**

BioMart is a very handy tool when you want to map IDs between different databases. The following is a list of 19 accession numbers from the UniProtKB/Swiss-Prot database (http://www.uniprot.org/) of human proteins that are supposedly involved in the sensory perception of pain.

Q99608, P34913, P28482, P28223, Q96LB1, P10997, P01210, P25929, P17481, P43681, P29460, Q9HC23, P20366, Q9Y2W7, Q99572, Q99571, Q00535, P01138, P17787

Can you convert this list of UniProtKB/Swisse-Prot IDs into HGNC symbols? You may want to include as attributes Ensembl Gene IDs and 'Description' (under the field GENE).

---

*ADDITIONAL EXERCISE: BIOMART*

*If you have finished the exercises above, you may want to do these extra ones*

*Extra exercise 4: Convert UniProt IDs into Ensembl IDs for Arabidopsis proteins*

*BioMart is a very handy tool when you want to map between different databases.*

*The following is a list of IDs from the UniProtKB/Swiss-Prot database*

*(http://www.uniprot.org/) of Arabidopsis thaliana proteins that are believed to be involved in flavonoid metabolism*

*http://amigo.geneontology.org/cgi-bin/amigo/term_details?term=GO:0009812:*

*P42813, Q9LS08, Q9ZST4, Q9SYM2, P51102, Q9LPV9, Q9FE25, Q96323, Q9FKW3, P13114, P41088, Q9S818, Q96330, O22203, Q39224, O22264, Q9SD85, Q9LYT3, Q9FJA2, Q43128, P43254, O04153, Q43125, Q9S9P6, Q94C57, Q9LNE6, Q9FK25, Q9SYM5, Q9ZQ95*

*Go to plants.ensembl.org and click on the link BioMart at the top of the page. Generate a list that shows, to which Ensembl Gene IDs these UniProtKB/Swiss-Prot IDs map to. Also include the Gene name, Gene description and Pfam ID.*

---

## Extra exercise 5: Retrieve a list of SNPs from the tomato genome (Solanum lycopersicum)

*The region between coordinates 21,394,819 and 21,397,868 on chromosome 6 in tomato contains a gene involved in oxidation-reduction process (GO:0055114).*

*Can I use BioMart to retrieve all the SNPs that cause a change at the amino acid level of this gene (those SNPs are known as missense variants) including their IDs and possible alleles?*

---

## Extra exercise 6: Find genes associated with array probes in Ensembl release 78

*In the paper 'Discovery of novel biomarkers by microarray analysis of peripheral blood mononuclear cell gene expression in benzene-exposed workers' (Forrest et al. Environ Health Perspect. 2005 June;*
*113(6): 801–807) the effect of benzene exposure on peripheral blood mononuclear cell gene expression in a population of shoe factory workers with well-characterized occupational exposures was examined using microarrays. The microarray used was the Affymetrix U133A/B GeneChip (also called 'U133 plus 2').*
*The top probe sets that are up-regulated by benzene exposure were:*

*207630_s_at, 221840_at, 219228_at, 204924_at 227613_at, 223454_at, 228962_at, 214696_at, 210732_s_at, 212371_at, 225390_s_at, 227645_at, 226652_at, 221641_s_at, 202055_at, 226743_at, 228393_s_at, 225120_at, 218515_at, 202224_at, 200614_at, 212014_x_at, 223461_at, 209835_x_at, 213315_x_at*

*Use the archive version of BioMart on the previous release of Ensembl, i.e. release 78. The list of our archives can be found below:*

*http://www.ensembl.org/info/website/archives/index.html*

*Choose Ensembl 78 and then select BioMart.*

*a) Generate a list of the genes to which these probe sets map. Include the Ensembl Gene ID, name and description as well as the probe set name.*

*b) As a first step towards analysing them for possible regulatory features they have in common, retrieve the 250 bp upstream of the transcripts of these genes. Include the Ensembl Gene and Transcript ID, name and description in the sequence header.*

*c) In order to be able to study these human genes in mouse, generate a list of the human genes and their mouse orthologues. Include the Ensembl Gene ID for both the human and mouse genes and the homology type in your list.*

*d) Generate the same list as in (c), but now also include the name and description of both the human and mouse genes. As the name and description of the mouse genes are not available as attributes in the Ensembl human genes dataset, you have to add the Ensembl mouse genes as a second dataset.*

# ENSEMBL TOOLS: THE VEP

You may have identified in your lab, previously unknown variants (e.g. SNPs, SVs) and would like to find out which genes they map to and what are the effects of them on genes and proteins.

In Ensembl, we have got a tool that allows you to do just that: annotate your own SNPs/SVs in any given genome, even if Ensembl has not got variation data for it (e.g. Bacteria).

This tool is called Variant Effect Predictor or VEP.

The VEP can be run on a web interface, via a Perl script or via our REST API. More information can be found below:

www.ensembl.org/vep

Let's try some exercises using the web interface of the VEP based on the talk and demo given by the instructor.

**Exercise 1 – Using VEP to predict the consequence of SNPs on the previous human assembly, GRCh37. Go to grch37.ensembl.org**

An analysis of 5,000 patients from a Taipei cohort has identified few variants associated with lung cancer:

chr 15, genomic coordinate 78889339, alleles G/A, forward strand
chr 22, genomic coordinate 30332586, alleles T/C, forward strand
chr 6, genomic coordinate 31721033, alleles G/A, forward strand
chr 5, genomic coordinate 1260624, alleles G/A, forward strand
chr 17, genomic coordinate 63554591, alleles A/G, forward strand
chr 5, genomic coordinate 1254510, alleles C/T, forward strand

Use the VEP to answer the following:

a) Which genes and transcripts do these variants map to?

b) What are the consequence terms for these variants?

c) Are there deleterious variants according to the SIFT/PolyPhen predictions? Are these two predictions always in agreement with each other?

---

## Exercise 2 – The VEP tool and variants on the bread wheat genome.

An analysis of 5,000 individuals from two different populations of bread wheat (*T. aestivum*) has identified thousands of polymorphic loci. See a list of a few of them below:

chr 2D, genomic coordinate 89551917, alleles G/A, forward strand
chr 2D, genomic coordinate148408765, alleles G/T, forward strand
chr 3D, genomic coordinate113574123, alleles C/A, forward strand
chr 3D, genomic coordinate 93827883, alleles G/A, forward strand
chr 3B, genomic coordinate 727928129, alleles C/T, forward strand
chr 3B, genomic coordinate 736734474, alleles C/T, forward strand
chr 6A, genomic coordinate 196872409, alleles T/G, forward strand
chr 6A, genomic coordinate 196153918, alleles A/G, forward strand
chr 6A, genomic coordinate 196774882, alleles G/C, forward strand

Can you use the VEP tool to answer the following?

a) Are any of these variants known in the public domain? Can you list a few of the IDs of these existing variants?

b) Which genes and transcripts do these variants map to?

c) Which consequence types can be found for these variants? Do any of them cause a change at the amino acid level?

---

## Exercise 3 – VEP in Bacterial genomes

Find the genome for *Bacteroides fragilis* 638R and launch the VEP tool. Use the VEP to predict the effects of a 7 bp deletion of TCTACAA on the supercontig FQ312004 at the position 258140-258146.

# QUICK GUIDE TO DATABASES AND PROJECTS

Here is a list of databases and projects you will come across in these exercises. Projects include many species, unless otherwise noted.

**Other help:**

**The Ensembl Glossary:**
http://www.ensembl.org/Help/Glossary

**Ensembl FAQs:**
http://www.ensembl.org/Help/Faq

## SEQUENCES

**EMBL-Bank, GenBank and DDBJ –** They contain nucleic acid sequences deposited by submitters, such as wet-lab biologists and gene sequencing projects. These three databases are synchronised with each other every day, so the same sequences should be found in all of them.

**CCDS** – Coding sequences that are agreed upon by Ensembl, Havana, UCSC and NCBI (human and mouse).

**NCBI Entrez Gene –** NCBI's gene collection.

**NCBI RefSeq –** NCBI's collection of 'reference sequences'. It includes genomic DNA, mRNAs and proteins. NM represents 'Known mRNA' (e.g. NM_005476) and NP (e.g. NP_005467) is 'Known proteins'.

**UniProtKB –** the "Protein knowledgebase" is a comprehensive set of protein sequences. It is divided into two parts: Swiss-Prot and TrEMBL

**UniProt Swiss-Prot –** the manually annotated, reviewed protein sequence set in the UniProtKB. High quality.

**UniProt TrEMBL –** the automatically annotated, unreviewed set of protein sequences (EMBL-Bank translated). Varying quality.**VEGA –** Vertebrate Genome Annotation database providing a selection of

manually-curated genes, transcripts, and proteins (human, mouse, zebrafish, gorilla, wallaby, pig, chimpanzee and dog).

**HAVANA –** Human and Vertebrate Analysis and annotation group at the Wellcome Trust Sanger Institute. It's the main contributor to the manual annotation presented in VEGA.

## GENE NAMES

**HGNC –** HUGO Gene Nomenclature Committee, a project assigning a unique name and symbol to every **human** gene.

**ZFIN –** The Zebrafish Model Organism Database
(Zebrafish International Resource Center).

**MGI –** Mouse Genome Informatics

**For other species you may want to try/use WikiGenes**

## ANNOTATION OF PROTEINS

**InterPro –** A collection of domains, motifs, and other protein signatures. Protein signature records are extensive, and combine information from individual projects such as UniProt, along with other databases such as SMART, PFAM and PROSITE.

**PFAM –** A collection of protein families

**PROSITE –** A collection of protein domains, families, and functional sites.

**SMART –** A collection of evolutionarily conserved protein domains.

## OTHER PROJECTS

**dbSNP –** A collection of sequence polymorphisms; mainly single nucleotide polymorphisms, along with small insertion and deletions.

**OMIM –** Online Mendelian Inheritance in Man – a resource for phenotypes and diseases related to genes *(human).*