# Browsing Genomes with Ensembl

**www.ensembl.org**
**www.ensemblgenomes.org**

**Coursebook v72**

http://tinyurl.com/Camb0613

Cambridge – 21st June 2013

# TABLE OF CONTENTS

# Introduction to Ensembl

## Getting started with Ensembl
### www.ensembl.org

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 such as yeast and worm are also imported for comparative analysis by the Ensembl 'compara' team. Most annotation is updated every two months, leading to increasing Ensembl versions (such as version 70), however the gene sets are determined less frequently. A sister browser at www.ensemblgenomes.org is set up to access non-chordates, namely bacteria, plants, fungi, metazoa, and protists.

Ensembl provides genes and other **annotation** such as regulatory regions, conserved base pairs across species, and sequence variations. The Ensembl gene set is based on protein and mRNA evidence in **UniProtKB** and **NCBI RefSeq** databases, along with manual annotation from the **VEGA/Havana** group. All the data are freely available and can be accessed via the web browser at www.ensembl.org. Perl programmers can directly access Ensembl databases through an Application Programming Interfaces (**Perl APIs**). 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 by the user.

**Synopsis – What can I do with Ensembl?**

- View genes with other annotation along the chromosome.
- View alternative transcripts (i.e. splice variants) for a given gene.
- Explore homologues and phylogenetic trees across more than 60 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), populations (human), or breeds (dog).
- View positions and sequence of mRNAs and proteins that align with Ensembl genes.
- Upload your own data.
- 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.

## 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!

- ❖ [Email](#) the team with comments or questions at [helpdesk@ensembl.org](mailto:helpdesk@ensembl.org)

- ❖ Follow the Ensembl [blog](#)

- ❖ Sign up to a [mailing list](#)

## Further reading

Flicek, P. *et al*
**Ensembl 2013**
Nucleic Acids Res. Advanced Access (Database Issue)
[http://www.ncbi.nlm.nih.gov/pubmed/23203987](http://www.ncbi.nlm.nih.gov/pubmed/23203987)

**Ensembl Methods Series**
[http://www.biomedcentral.com/series/ENSEMBL2010](http://www.biomedcentral.com/series/ENSEMBL2010)

Xosé M. Fernández-Suárez and Michael K. Schuster
**Using the Ensembl Genome Server to Browse Genomic Sequence Data.**
UNIT 1.15 in Current Protocols in Bioinformatics, Jun 2010.

Giulietta M Spudich and Xosé M Fernández-Suárez
**Touring Ensembl: A practical guide to genome browsing**
BMC Genomics 2010, 11:295 (11 May 2010)

# Exploring the Ensembl genome browser

## Demo: Ensembl species

The front page of Ensembl is found at ensembl.org. It contains lots of information and links to help you navigate Ensembl:



Click on View full list of all Ensembl species.

Click on the common name of your species of interest to go to the species homepage. We'll click on Human.

**Search**

**Information and statistics**

**News**

**Links to example features in Ensembl**

To find out more about the genome assembly and genebuild, click on More information and statistics.



**Information**

**Tables of statistics**

Let's take a look at the Ensembl Genomes homepage at ensemblgenomes.org.



Click on the different taxa to see their homepages. Each one is colour-coded.



Protists

Fungi

Metazoa



Plants



Bacteria

You can navigate most of the taxa in the same way as you would with Ensembl, but Ensembl Bacteria has a large number of genomes, so needs slightly different methods. Let's look at it in more detail.



**Search for a gene**

**Search for a species**

**Information on Ensembl Bacteria**

There's no full species list for bacteria as it would be hard to navigate with the number of species. To find a species, start to type the species name into the species search box. A drop down list will appear with possible species.

For example, to find a substrain of *Clostridium difficile* type in Clostridium d.



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

**Exercises: Ensembl species**

**Exercise 1 – Panda**

(a) Go to the species homepage for Panda. What is the name of the genome assembly for Panda?

(b) Click on More information and statistics. How long is the Panda genome (in bp)? How many genes have been annotated?

**Exercise 2 – Zebrafish**

(a) What's new in release 71 for zebrafish?

(b) What previous release is available for zebrafish?

**Exercise 3 – Mosquitos**

(a) Go to Ensembl Metazoa. How many species of the genus *Anopheles* are there?

(b) Who published the genome sequence for *Anopheles gambiae*?

**Exercise 4 – Bacteria**

Go to Ensembl Bacteria and find the species *Belliella baltica*. How many coding and non-coding genes does it have?

**Demo: The Region in detail view**

Start at the Ensembl front page, ensembl.org. You can search for a region by typing it into a search box, but you have to specify the species.

Type (or copy and paste) human 4:123792818-123867893 into either search box.

Press Enter or click Go to jump directly to the **Region in detail** Page.

Click on the button  to view page-specific help.

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

There is a help video on this page at http://youtu.be/tTKEvgPUq94.

The Region in detail page is made up of three images, let's look at each one on detail.
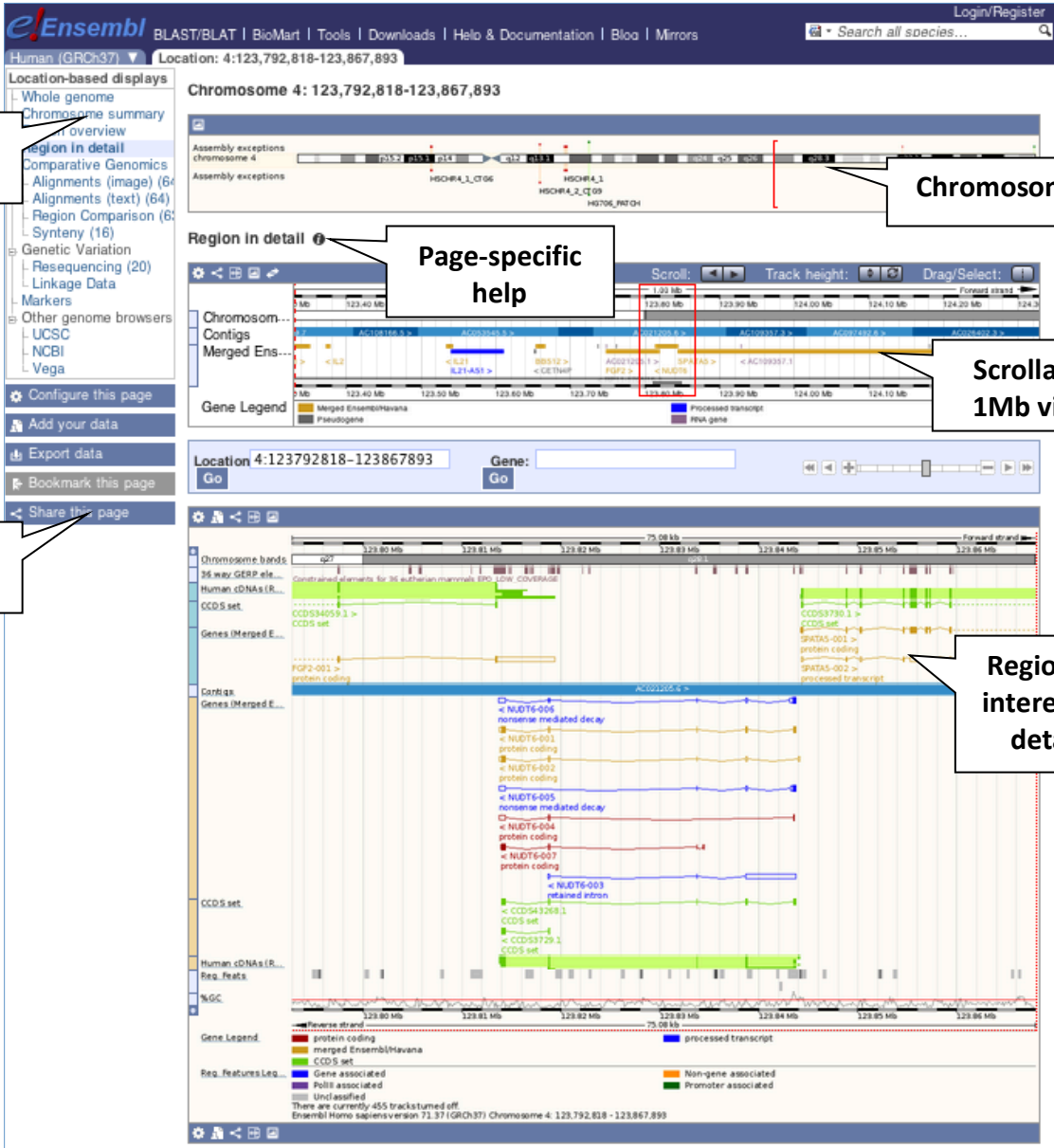
The first image shows the chromosome:



You can jump to a different region by dragging out a box in this image. Drag out a box on the chromosome, a pop-up menu will appear.



If you wanted to move to the region, you could click on Jump to region (###bp). For now, we'll close the pop-up by clicking on the X on the corner.

The second image shows a 1Mb region around our selected region. This view allows you to scroll back and forth along the chromosome.

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 would jump the lower image to the region central to the scrollable image. We want to go back to where we started, so we'll click on Reset scrollable image.



You can also drag out and jump to a region. Either hold down shift and drag in the image, or click on the Drag/Select button  to change the action of your mouse click, and drag out a box.



Click on the X to close the pop-up menu.

The third image is a detailed, configurable view of the region.

We can edit what we see on this page by clicking on the blue Configure this page menu at the left.



This will open a menu that allows you to change the image.

You can put some tracks on in different styles; more details are in this FAQ: http://www.ensembl.org/Help/Faq?id=335.

Let's add some tracks to this image. Add:

- Human proteins – Labels

- dbSNP variants – Normal

- 1000 Genomes – AMR – Collapsed

Now click on the tick in the top left hand to save and close the menu. Alternatively, click anywhere outside of the menu. We can now see the tracks in the image.

We can also change the way the tracks appear by hovering over the track name then the cog wheel to open a menu. We can move tracks around by clicking and dragging on the bar to the left of the track name.

Now that you've got the view how you want it, you might like to show something you've found to a colleague or collaborator. Click on the Share this page button to generate a link. Email the link to someone else, so that they can see the same view as you, including all the tracks you've added. 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 this to the default view, go to Configure this page and select Reset configuration at the bottom of the menu.

**Exercises: The Region in Detail view**

**Exercise 5 – Exploring a genomic region in human**

(a) Go to the region from 32,448,000 to 33,198,000 bp 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 that contain the complete *BRCA2* gene?

(d) Create a Share link for this display. Email it to yourself and open the link.

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

(f) Turn off all tracks you added to the Region in detail page.

## Exercise 6 – Exploring patches and haplotypes in human

(a) Go to the region 6:112294691-112624977 in human. What is the green highlighted region? (Tip: you can search for help terms in the Ensembl search boxes.)

(b) Can you see the patches in the chromosome view? Drag out a box to jump to a region containing the patch labelled HG27_patch. What are the coordinates of the patch?
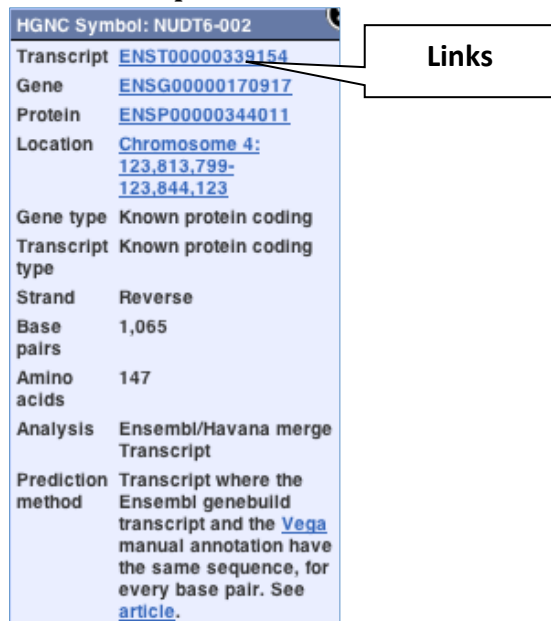
(c) Can you compare this patch with the reference? What has changed between this patch and the sequence it replaced?

(d) Go back to the previous view and scroll to the right in the 1Mb view until you reach a red highlighted region. What is this?

# Genes and transcripts

## Demo: The gene tab

If you click on any one of the transcripts in the Region in detail image, a pop-up menu will appear, allowing you to jump directly to that gene or transcript.
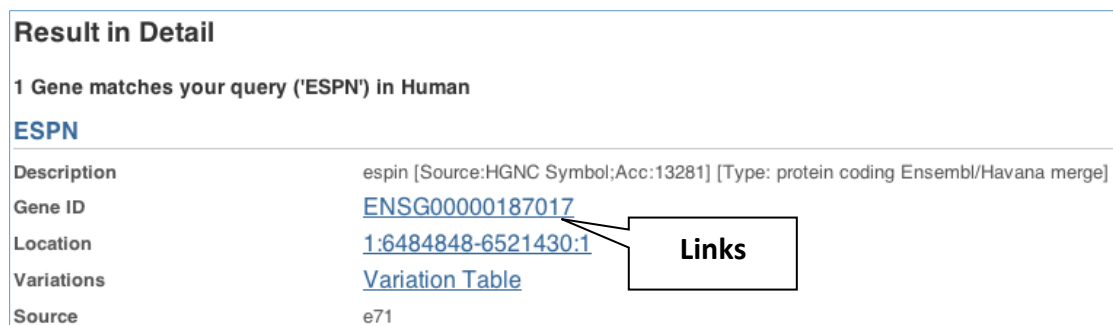
| | | |
|---|---|---|
| HGNC Symbol: NUDT6-002 | | |
| Transcript | ENST00000339154 | Links |
| Gene | ENSG00000170917 | |
| Protein | ENSP00000344011 | |
| Location | Chromosome 4: 123,813,799-123,844,123 | |
| Gene type | Known protein coding | |
| Transcript type | Known protein coding | |
| Strand | Reverse | |
| Base pairs | 1,065 | |
| Amino acids | 147 | |
| Analysis | Ensembl/Havana merge Transcript | |
| Prediction method | Transcript where the Ensembl genebuild transcript and the Vega manual annotation have the same sequence, for every base pair. See article. | |

Another way to go to a gene of interest is to search directly for it.

We're going to look at 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 (http://tinyurl.com/espn-ncbi-gene).

From ensembl.org, type *ESPN* into the search bar and click the Go button. Click on Gene and Human to find the hits.

**Result in Detail**

1 Gene matches your query ('ESPN') in Human

**ESPN**

| | |
|---|---|
| Description | espin [Source:HGNC Symbol;Acc:13281] [Type: protein coding Ensembl/Havana merge] |
| Gene ID | ENSG00000187017 |
| Location | 1:6484848-6521430:1    Links |
| Variations | Variation Table |
| Source | e71 |

Click on the gene name or Ensembl ID. The **Gene tab** should open:



Let's walk through some of the links in the left hand navigation column. How can we view the genomic sequence? Click Sequence at the left of the page.

**Most recent human genome assembly GRCh37 = hg19**

Human (GRCh37) ▼

Gene-based displays
- **Gene summary**
- Splice variants (10)
- Transcript comparison
- Supporting evidence
- Sequence
- External references
- Regulation
- Expression

**Click Sequence**

## Marked-up sequence ❶

Key

Exons    All exons in this region    ESPN exons

```
>chromosome:GRCh37:1:6484248:6522030:1
AGCGCACTAGTGGTTCCCGTCCTGCCTTCCTGCCCTCCCCGCTGGAACCTCTG
TTCTCGGATCTGGAGGGACCCTGGAAGGCAGGGCTCTTTGCAATCTCCGGGGA
CCAGAGCCCTTCAGGGACGTGGCAGGGCTGCTCCTGCCTCAGGGCCGTTGTCC
CCTCACCCCGCCTGGAATACCCTTCTCGCCGCTCAAACCCAGCCCCACGGCACCTCCTCA
GAGACCTTTCCCTGTCCGCCCACGCGGTCCCGACAATCACTCCCCATCACCTCTGGAATT
GCGTCGCCGGCGCCTGGAACCGCAGTTAGCGGGCACTGGGCAGATGAATGAAT
TGCCTGGACGGCTCTCCAATTCGAACCCAGTTTTGCTGCCCTCTGGGGTCTCAA
CGTGAGGCAAATTAGGAGAGAAGCCCCTGGGCACCTTGCCCCAGTCGCACGAG
GCGTCGCGGCGGGGGCGGGCGGGGAACTCGGGCGGAGGCTGCGGGGCGGGGCGGGGCGGG
GTGGGGGCGGGCCCGAGTCTTAAGCCGGCGTCCGCGGGCTCCGGCCCCAGAGCGCGGCGG
AGCGGAGCGCCAGGCAGCGCGCGGAGCGGAGGCCAGGCCCACAGCCGCTCCGCCTCCCGGCC
CGCAGATCCCCGACGGCCGCACCGCGGGCTCCTCTGGCCCGCAAGAACACGTGCATGGCG
TCCTGGGGAAGGCGCTGAGTGCGGAGTCGCGGCGCCGCACGCGGCACCATGCC
CAGGCGCTGCAGGCGGCGCGGCAGGGCGAGCTGGACGTGCTGAGGTCGCTGCAC
GGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGCGCTGCCCGTGCACCACGCGGCC
CGCGCTGGGAAGCTGCACTGTCTGCGCTTCCTGGTGGAGGAAGCCGCCCTCCCCGCCGCG
GCCCGCGCCCGCAACGGCGCCCACACCGGCCCACGACGCCTCCGCCACCGGCCACCTCGCC
TGCCTGCAGTGGCTGCTGTCGCAGGGCGGCTGCAGAGTGCAGGTGGGTCCGCGCGGTTCG
CCAGGGGCACTGAGGCTTCCTCCTCAGGACAGAGTCCTGGCCCAGAGTCCCCGGGGCTC
AAGGATGGGTGGGGTTTGGCACCTCCTGGCCCAGCTGAACCCTGCACGGAGCTCCTTCCA
```

**Upstream sequence**

**Exon of an overlapping gene**

***ESPN* Exon**

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

**name of genome assembly** | **chromosome** | **base pair start** | **base pair end** | **forward strand (-1 is reverse)**

```
>chromosome:GRCh37:1:6484248:6522030:1
AGCGCACTAGTGGTTCCCGTCCTGCCTTCCTGCCCTCCCCGCTGGAACCTCTGGGGGCAG
TTCTCGGATCTGGAGGGACCCTGGAAGGCAGGGCTCTTTGCAATCTCCGGGGATTTCGAC
CCAGAGCCCTTCAGGGACGTGGCAGGGCTGCTCCTGCCTCAGGGCCGTTGTCCTCGTGCT
```

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



Once you have selected changes (in this example, Show variations and Line numbering) click at the top right.



Let's look at where our gene is expressed. Click on Expression in the left-hand menu.

## Expression ℹ

Expression data is available for the following tissues:

| Tissue | All data | RNASeq gene models | Intron-spanning reads | RNASeq alignments |
|--------|----------|--------------------|-----------------------|-------------------|
| Adipose | View in location | Models built using Human adipose total RNA, lot 05060581, caucasian female, throat cancer, Illumina Human Bodymap 2.0 Data | Y | Y |
| Adrenal | View in location | Models built using Human adrenal total RNA, lot 0812003, caucasian male, cerebral vascular accident, Illumina Human Bodymap 2.0 Data | Y | Y |
| Blood | View in location | Models built using Human white blood cell, caucasian male, healthy, Illumina Human Bodymap 2.0 Data | Y | Y |
| Brain | View in location | Models built using Human brain total RNA, lot 03070051, caucasian female, copd, Illumina Human Bodymap 2.0 Data | Y | Y |

Hover over the column titles for a pop-up definition.

Can our gene be found in other databases? Go up the left-hand menu to External references:



## External references ℹ

This gene corresponds to the following database identifiers:

| | Filter |
|---|---|

| External database | Database identifier |
|-------------------|---------------------|
| HGNC Symbol | ESPN<br>espin [view all locations] |
| EntrezGene | ESPN<br>espin [view all locations] |
| MIM disease | DEAFNESS, AUTOSOMAL RECESSIVE 36, [#609006]<br>DEAFNESS, AUTOSOMAL RECESSIVE 36, WITH OR WITHOUT VESTIBULAR INVOLVEMENT; [view all locations] |
| UniProtKB Gene Name | ESPN [view all locations] |
| Orphanet | Autosomal recessive nonsyndromic sensorineural deafness type DFNB<br>Autosomal recessive nonsyndromic sensorineural deafness type DFNB [view all locations] |
| WikiGene | ESPN<br>espin [view all locations] |
| MIM gene | ESPIN, MOUSE, HOMOLOG OF [*606351]<br>ESPIN, MOUSE, HOMOLOG OF; ESPN [view all locations] |
| UniGene | Hs.652319 [Target %id: 99; Query %id: 94]<br>Transcribed locus [view all locations]<br>Hs.744222 [Target %id: 99; Query %id: 99]<br>Espin [view all locations] |
| ArrayExpress | ENSG00000187017 [view all locations] |

This contains links to the gene in other projects, such as Uniprot.

To find out more about the individual transcripts of this gene, click on Transcript comparison in the left-hand menu.

You must now choose the transcripts you'd like to see, click on the blue Select transcripts button.



Let's select all the protein-coding transcripts, then close the menu.



## Demo: The transcript tab

Let's now explore one splice isoform. Click on Show transcript table at the top.



Click on the ID for the largest one, 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, circled in the image below.



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

**Display options**

| | |
|---|---|
| Flanking sequence at either end of transcript: | 50 |
| Number of base pairs per row: | 60 bps |
| Intron base pairs to show at splice sites: | 25 |
| Show full intronic sequence: | ☑ |
| Show exons only: | ☐ |
| Line numbering: | None |
| Show variations: | In exons only |
| Filter variations by consequence type: | No filter<br>3 prime UTR variant<br>5 prime UTR variant<br>Coding sequence variant<br>Downstream gene variant |

If you would like to export the sequence, including the colours, click Download view as RTF. A Rich Text Format document will be generated that can be opened in word processor such as MS Word.



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

```
                  R         YR  Y  Y      M
481 GCCACAGTCTTGCATCTGGCTGCCCGCTTCGGCCACCCCGAGGTGGTGAACTGGCTCTTG
313 GCCACAGTCTTGCATCTGGCTGCCCGCTTCGGCCACCCCGAGGTGGTGAACTGGCTCTTG
105 -A--T--V--L--H--L--A--A--R--F--G--H--P--E--V--V--N--W--L--L-

            R      R  B    M YRYR        R        K                    Y
541 CATCATGGCGGTGGGGACCCCACCGCGGCCACAGACATGGGCGCCCTGCCTATCCACTAC
373 CATCATGGCGGTGGGGACCCCACCGCGGCCACAGACATGGGCGCCCTGCCTATCCACTAC
125 -H--H--G--G--G--D--P--T--A--A--T--D--M--G--A--L--P--I--H--Y-
```
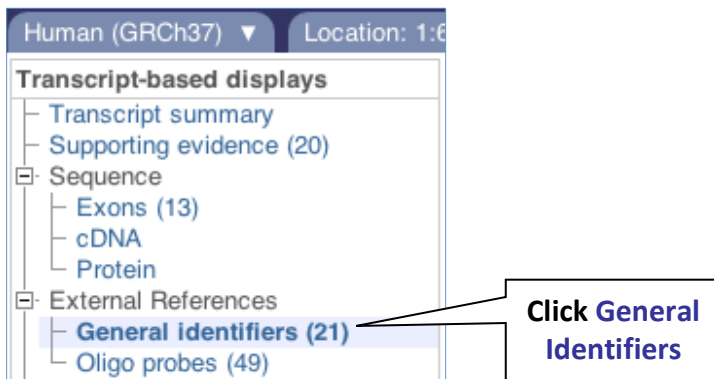
UnTranslated Regions (UTRs) are highlighted in dark yellow, codons are highlighted in light yellow, and exon sequence is shown in black or blue letters to show exon divides. Sequence variants are represented by highlighted nucleotides and clickable IUPAC codes are above the sequence.

Next, follow the General identifiers link at the left.

```
Human (GRCh37) ▼   Location: 1:6
Transcript-based displays
 ├ Transcript summary
 ├ Supporting evidence (20)
 ⊟ Sequence
 │  ├ Exons (13)
 │  ├ cDNA
 │  └ Protein
 ⊟ External References              Click General
    ├ General identifiers (21)         Identifiers
    └ Oligo probes (49)
```

This page shows information from other databases such as RefSeq, EntrezGene, OMIM, UniProtKB, and others, that matches to the Ensembl transcript and protein

## General identifiers ℹ

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] |
| INSDC protein ID | AAD24480.1 [align] [view all locations]<br>AAP34481.1 [align] [view all locations]<br>CAB66814.1 [align] [view all locations]<br>CAI19773.1 [align] [view all locations]<br>CAI22163.1 [align] [view all locations]<br>EAW71537.1 [align] [view all locations] |

Click on Ontology table to see GO terms from the Gene Ontology consortium. www.geneontology.org

Click on the 🛈 to see a guide to the three-letter Evidence codes.

Now click on Protein summary to view domains from Pfam, PROSITE, Superfamily, InterPro, and more.



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

## Exercises: Genes and transcripts

### Exercise 7 – Exploring the human *MYH9* gene

(a) Find the human *MYH9* (myosin, heavy chain 9, non-muscle) gene, and go to the Gene tab.

- On which chromosome and which strand of the genome is this gene located?

- How many transcripts (splice variants) are there?

- How many of these transcripts are protein coding?

- What is the longest transcript, and how long is the protein it encodes?

- Which transcript has a CCDS record associated with it?

❓ Why is the CCDS important – what does it tell us?

(b) Click on Phenotype at the left side of the page. Are there any diseases associated with this gene, according to MIM (Mendelian Inheritance in Man)?

(c) In the transcript table, click on the transcript ID for MYH9-001, and go to the Transcript tab.

- How many exons does it have?

- Are any of the exons completely or partially untranslated?

- Is there an associated sequence in UniProtKB/Swiss-Prot? Have a look at the General identifiers for this transcript.

- What are some functions of MYH9-001 according to the Gene Ontology consortium? Have a look at the Ontology table for this transcript.

(d) Are there microarray (oligo) probes that can be used to monitor ENST00000216181 expression?


## Exercise 8 – Finding a gene associated with a phenotype

Phenylketonuria is a genetic disorder caused by an inability to metabolise phenylalanine. This results in an accumulation of phenylalanine causing seizures and mental retardation.

(a) Search for phenylketonuria from the Ensembl homepage. What gene is associated with this disorder?

(b) What tissues is this gene expressed in? Is this surprising, given the gene's role in disease? What is meant by "Intron-spanning reads" and "RNASeq alignments"?

(c) How many protein coding transcripts does this gene have? View all of these in the transcript comparison view.

(d) What is the MIM disease identifier for this gene? Does Orphanet list any other disorders associated with this gene?


## Exercise 9 – Exploring a plant gene (*Vitis vinifera*, grape)

Start in http://plants.ensembl.org/index.html and select the *Vitis vinifera* genome.

(a) What GO: biological process terms are associated with the *MADS4* gene?

(b) Go to the transcript tab for the only transcript, Vv01s0010g03900.t01. How many exons does it have? Which one is the longest? How much of that is coding?

(c) What domains can be found in the protein product of this transcript? How many different groups agree with each of these domains?

# BioMart

## Demo: BioMart

Follow these instructions to guide you through BioMart to answer the following query:

> You have three questions about a set of human genes: *ESPN, MYH9, USH1C, CHD7, CISD2, THRB, DFNB31*
> *(*these are HGNC gene symbols. More details on the HUGO Gene Nomenclature Committee can be found on http://www.genenames.org)
>
> 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

NOTE: These answers were determined using BioMart Ensembl 71.



STEP 2:
Choose Ensembl Genes 72 as the primary database.



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

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



**STEP 5:**
In ID List Limit, paste in your gene symbols. Change the heading to read HGNC symbol(s) [e.g. ZFY].



**STEP 6:**
Click Count to see BioMart is reading 8 genes out of 63,253 possible *H. sapiens* genes (this number includes ncRNA genes) (note that *CHD7* has an Ensembl copy and an LRG copy (http://www.ensembl.org/Help/Glossary?id=406), hence 8 counts for 7 HGNC genes).

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

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

**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.

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

**STEP 12:**
Click Results.

*Why are there multiple rows for one gene ID? For example, look at the first few rows.*

| 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

⭐ URL  ➡ XML  📜 Perl  ⊙ Help

...ults

Please select columns to be included in the output and hit 'Results'

○ Features    ○ Homologs
○ Structures    ○ Variation
○ Transcript Event   ◉ Sequences

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

⊟ SEQUENCES:

Sequences (max 1)

○ Unspliced (Transcript)
○ Unspliced (Gene)
○ Flank (Transcript)
○ Flank (Gene)
○ Flank-coding region (Transcript)
○ Flank-coding region (Gene)

○ 3' UTR
○ Exon sequences
◉ cDNA sequences
○ Coding sequence
○ Protein

...37.p3)
...ters
...GNC symbol(s) [e.g. ZFY]: [ID-list specified]
**Attributes**
Ensembl Gene ID
Ensembl Transcript ID
cDNA sequences

**Dataset**
[None Selected]

⊟ Header Information
**Gene Information**
☑ Ensembl Gene ID
☐ Description
☑ Associated Gene Name
☐ Associated Gene DB
☐ Chromosome Name

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

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

**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 learn these things from 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

**Exercises: BioMart**

**Exercise 10 – Finding genes by protein domain**

Find mouse proteins with transmembrane domains located on chromosome 9.

## Exercise 11 – Convert IDs

BioMart is a very handy tool when you want to convert IDs from different databases. The following is a list of 29 IDs of **human proteins** from the NCBI **RefSeq** database (http://www.ncbi.nlm.nih.gov/projects/RefSeq/):

NP_001218, NP_203125, NP_203124, NP_203126, NP_001007233, NP_150636, NP_150635, NP_001214, NP_150637, NP_150634, NP_150649, NP_001216, NP_116787, NP_001217, NP_127463, NP_001220, NP_004338, NP_004337, NP_116786, NP_036246, NP_116756, NP_116759, NP_001221, NP_203519, NP_001073594, NP_001219, NP_001073593, NP_203520, NP_203522

Generate a list that shows to which Ensembl Gene IDs and to which HGNC symbols these RefSeq IDs correspond. Do these 29 proteins correspond to 29 genes?

Hint: For this exercise, it's easier to copy and paste the IDs from the online exercise booklet. One copy is here:
URL

## Exercise 12 – Export homologues

For a list of *Ciona savignyi* Ensembl genes, export the human orthologues.

ENSCSAVG00000000002, ENSCSAVG00000000003, ENSCSAVG00000000006, ENSCSAVG00000000007, ENSCSAVG00000000009, ENSCSAVG00000000011

## Exercise 13 – Export structural variants

You can use BioMart to query variants, not just genes.

(a) Export the study accession, source name, chromosome, sequence region start and end (in bp) of human structural variations (SV) on chromosome 1, starting at 130,408 and ending at 210,597.

(b) In a new BioMart query, find the alleles, phenotype descriptions, and associated genes for rs1801500 and rs1801368. Can you view this same information in the Ensembl browser?


## Exercise 14 – Find genes associated with array probes

Forrest *et al* performed a microarray analysis of peripheral blood mononuclear cell gene expression in benzene-exposed workers (Environ Health Perspect. 2005 June; 113(6): 801–807). The microarray used was the human Affymetrix U133A/B (also called U133 plus 2) GeneChip. The top 25 up-regulated probe-sets were:

207630_s_at, 221840_at, 219228_at, 204924_at, 227613_at, 223454_at, 228962_at, 214696_at, 210732_s_at, 212370_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

(a) Retrieve for the genes corresponding to these probe-sets the Ensembl Gene and Transcript IDs as well as their HGNC symbols and descriptions.

(b) In order to analyse these genes for possible promoter/enhancer elements, retrieve the 2000 bp upstream of the transcripts of these genes.

(c) In order to be able to study these human genes in mouse, identify their mouse orthologues. Also retrieve the genomic coordinates of these orthologues.


## Exercise 15 – BioMart in Ensembl Fungi

You can use BioMart for the non-vertebrate species hosted at www.ensemblgenomes.org. Export a list of Gene IDs (from the PomBase project only) for *S. pombe* that are protein coding and located on Chromosome III.

(Start at http://fungi.ensembl.org)

## Variation

### Demo: Finding variants in Ensembl

In any of the sequence views shown in the Gene and Transcript tabs, you can view variants on the sequence. You can do this by clicking on Configure this page from any of these views.

Let's take a look at the Gene sequence view for *MCM6* in human. Search for *MCM6* and go to the Sequence view.

If you can't see variants marked on this view, click on Configure this page and select Show variations: Yes and show links.



Find out more about a variant by clicking on it.



You can add variants to all other sequence views in the same way.

You can go to the Variation tab by clicking on the variant ID. For now, we'll explore more ways of finding variants.

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



The table is divided into consequence types. You can also view the consequence types as an ontology tree. Click on Switch to tree view at the right hand side to change the view.



Click on Show to expand a detailed table for any of the consequence types available.

Let's expand Missense variants.



The table contains lots of information about the variants. You can click on the IDs here to go to the Variation tab too.

Let's look at Structural Variation in the Gene Tab. You'll find it in the left-hand menu.



You can click on the structural variants (SVs) in the image, or on their IDs in the table to go to the SV tab.

Let's have a look at variants in the Location tab. Click on the Location tab in the top bar.



Configure this page and open Variation from the left-hand menu.



There are various options for turning on variants. You can turn on variants by source, by frequency, presence of a phenotype or by individual genome they were isolated from. Turn on the following sequence variants in Normal.

- 1000 genomes – All
- 1000 genomes – All – common
- All phenotype-associated variants
- ENSEMBL:Venter

Also turn on Larger and Smaller Structural variants (all sources) in Expanded.

Click on a variant to find out more information. It may be easier to see the individual variants if you zoom in.

Let's zoom in on the region 2:136607850-136609811 by typing it into the Location box.

Change the track style on one of the SNP variant tracks to Expanded with name. Click on the variant rs4988235 to open a pop-up, then click on the ID to open the Variation tab.

The icons show you what information is available for this variant. Click on Genes and regulation, or follow the link at the left.



This variant is found in three transcripts of the *MCM6* gene. It has not been associated with any regulatory features or motifs.

Let's look at population genetics. Either click on Explore this variant in the left hand menu then click on the Population genetics icon, or click on Population genetics in the left-hand menu.

Population genetics ⓘ

**1000 Genomes allele frequencies**

| ALL | AFR | AMR | ASN | EUR |
|-----|-----|-----|-----|-----|
| A: 23%<br>G: 77% | A: 4%<br>G: 96% | A: 26%<br>G: 74% | G: 100% | A: 52%<br>G: 48% |
| Sub-populations ⊞ | Sub-populations ⊞ | Sub-populations ⊞ | Sub-populations ⊞ | Sub-populations ⊞ |

**1000 Genomes (19)**

Show [All ⬍] entries        Show/hide columns        Filter

| Population | Alleles A | Alleles G | Genotypes A\|A | Genotypes A\|G | Genotypes G\|G | Allele count | Genotype count | |
|---|---|---|---|---|---|---|---|---|
| 1000GENOMES:phase_1_ALL | 0.234 | 0.766 | 0.135 | 0.198 | 0.668 | 510 (A) / 1674 (G) | 147 (A\|A) / 216 (A\|G) / 729 (G\|G) | Show |
| 1000GENOMES:phase_1_AFR | 0.037 | 0.963 | 0.012 | 0.049 | 0.939 | 18 (A) / 474 (G) | 3 (A\|A) / 12 (A\|G) / 231 (G\|G) | Show |
| 1000GENOMES:phase_1_AMR | 0.262 | 0.738 | 0.094 | 0.337 | 0.569 | 95 (A) / 267 (G) | 17 (A\|A) / 61 (A\|G) / 103 (G\|G) | Show |

These data are mostly from the **1000 genomes** and **HapMap** projects in human.

There are big differences in allele frequencies between populations. Let's have a look at the phenotypes associated with this variant to see if they are known to be specific to certain human populations. Either click on Explore this variant in the left hand menu then click on the Phenotype data icon, or click on Phenotype Data in the left-hand menu.



Phenotype Data ⓘ

Show/hide columns        Filter

| Disease/Trait | Source(s) | Study | Reported gene(s) | Associated variant(s) | Most associated allele | P value |
|---|---|---|---|---|---|---|
| LACTASE PERSISTENCE [View on Karyotype] | OMIM | MIM:601806 | MCM6 | rs4988235 | 0001 | |

This variant is associated with lactase persistence, which is known to be common in European populations, and rare in Asian populations, exactly as we saw in the allele frequencies in these populations.

Are there other variants in the genome that also cause lactase persistence? Click on View on Karyotype to find out.

Locations of features associated with LACTASE PERSISTENCE

**Hits on the karyotype**

Click on the image above to jump to a chromosome, or click and drag to select a region

Key

| Feature type | Colour |
|---|---|
| Variation | 1.0  3.0  4.0  5.0  6.0  7.0  8.0 |

Less significant -log(p-values) ← → More significant -log(p-values)

**Legend showing hit significance**

Features associated with phenotype LACTASE PERSISTENCE

Show/hide columns

**Table of variants**

| Genomic location (strand) | Ensembl ID | Feature type | Reported gene(s) | Associated phenotype(s) | Annotation source(s) | |
|---|---|---|---|---|---|---|
| 2:136603646-136613646(1) | rs4988235 | Variation | MCM6 | LACTASE PERSISTENCE | OMIM | |
| 2:136611754-136621754(1) | rs182549 | Variation | MCM6 | LACTASE PERSISTENCE | OMIM | - |

Two variants are known to be associated with this phenotype. Both are found with the *MCM6* gene.

Click back to the Variation Tab. Click on Phylogenetic context to see the variant in other species.



**Choose your alignment**

Phylogenetic Context ⓘ

Alignment: 6 primates EPO    Go

Key

Variations    Intronic

● **Focus variant**

**Aligned regions**

| Homo sapiens › | chromosome:GRCh37:2:136608636:136608656:1 |
| Pan troglodytes › | chromosome:CHIMP2.1.4:2B:139811109:139811129:1 |
| Gorilla gorilla gorilla › | chromosome:gorGor3.1:2b:22952947:22952967:1 |
| Pongo abelii › | chromosome:PPYG2:2b:24805669:24805689:1 |
| Macaca mulatta › | chromosome:MMUL_1:13:116302799:116302819:-1 |
| Callithrix jacchus › | chromosome:C_jacchus3.2.1:...676:-1 |

**SNP of interest**

```
                          SR  R
Homo sapiens              GAGGCCAGGGGTACA
Pan troglodytes           GAGGCCAGGGGCTACATTATC
Gorilla gorilla gorilla   GAGGCCAGGGGCTACATTATC
Pongo abelii              GAGGCCAGGGGCTACATTATC
Macaca mulatta            GAGGCCAGGGGCTACATTATC
Callithrix jacchus        GAGGCCAGGGGCTACATTATC
```

**Alignment between species**

48

The variant is not marked in the other species. This means that the variant arose in humans.

Another way to look at variation is using the **Resequencing** view. Click on the Location tab, then choose Resequencing from the left-hand menu.



This view is used to look at sequence between individuals. Craig Venter and James Watson are shown by default. You can change the individuals shown by clicking on Configure this page.

It can also be used to look at different mouse strains or dog breeds. Can you find our variant rs4988235 in this view? Is the alternate allele present in James Watson or Craig Venter?

**Exercises: Finding variants in Ensembl**

**Exercise 16 – Human population genetics and phenotype data**

The SNP rs1738074 in the 5' UTR of the human *TAGAP* gene has been identified as a genetic risk factor for a few diseases.

(a) In which transcripts is this SNP found?

(b) What is the least frequent genotype for this SNP in the Yoruba (YRI) population from the HapMap set?

(c) What is the ancestral allele? Is it conserved in the 36 eutherian mammals?

(d) With which diseases is this SNP associated? Are there any known risk (or associated) alleles?

## Exercise 17 – Exploring a SNP in human

The missense variation rs1801133 in the human *MTHFR* gene has been linked to elevated levels of homocysteine, an amino acid whose plasma concentration seems to be associated with the risk of cardiovascular diseases, neural tube defects, and loss of cognitive function. This SNP is also referred to as 'A222V', 'Ala222Val' as well as other HGVS names.

(a) Find the page with information for rs1801133.

(Note: a bug in the current release means that these alleles are erroneously reported as G/A/CT/CT. Please ignore these extra CT alleles.)

(b) Is rs1801133 a Missense variation in all transcripts of the *MTHFR* gene?

(c) Why are the alleles for this variation in Ensembl given as G/A and not as C/T, as in dbSNP and literature?
(http://www.ncbi.nlm.nih.gov/projects/SNP/snp_ref.cgi?rs=180113
3)

(d) What is the major allele in rs1801133?

(e) In which paper is the association between rs1801133 and homocysteine levels described?

(f) According to the data imported from dbSNP, the ancestral allele for rs1801133 is G. Ancestral alleles in dbSNP are based on a comparison between human and chimp. Does the sequence at this same position in four other primates, i.e. gorilla, orangutan, macaque and marmoset, confirm that the ancestral allele is G?

(g) Were both alleles of rs1801133 already present in Neanderthal? To answer this question, have a look at the individual reads at its genomic position in the Neanderthal Genome Browser (http://neandertal.ensemblgenomes.org/).

## Exercise 18 – Structural variation in human

In the paper 'The influence of *CCL3L1* gene-containing segmental duplications on HIV-1/AIDS susceptibility' (Gonzalez *et al* Science. 2005 Mar 4; 307(5704):1434-40) it is shown that a higher copy number of the *CCL3L1* (Chemokine (C-C motif) ligand 3-like 1) gene is associated with lower susceptibility to HIV infection.

(a) Find the human *CCL3L1* gene.

(b) Have any CNVs been annotated for this gene? Note: In Ensembl, CNVs are classified as structural variants**.**

## Exercise 19 – Exploring a SNP in mouse

Madsen *et al* in the paper 'Altered metabolic signature in pre-diabetic NOD mice' (PloS One. 2012; 7(4): e35445) have described several regulatory and coding SNPs, some of them in genes residing within the previously defined *insulin dependent diabetes (IDD)* regions. The authors describe that one of the identified SNPs in the murine *Xdh* gene (rs29522348) would lead to an amino acid substitution and could be damaging as predicted as by SIFT (http://sift.jcvi.org/).

(a) Which chromosome and coordinates in the SNP located?

(b) What is the HGVS recommendation nomenclature for this SNP?

(c) Why does Ensembl put the C allele first (C/T)?

(d) Are there differences between the genotypes reported in NOD/LTJ and BALB/cByJ?

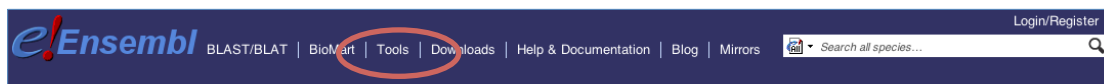## Demo: The Variant Effect Predictor (VEP)

We have analysed a samples from a patient with a genetic disorder. The patient presents with facial and limb deformities, mental

retardation and gastrointestinal reflux. Our genotyping has identified a mutation that may be responsible for the phenotype:
*An A->G mutation on chromosome 5 at 37,017,205 on the + strand.*

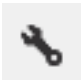We will use the **Ensembl VEP** to determine:

- Has my variant already been annotated in Ensembl?

- What genes are affected by my variant?

- Does my variant result in a protein change?

From any page in Ensembl click on Tools in the top bar.



This provides a table showing the tools that are available in Ensembl. There's a short description of what each of the tools do. All are available as both an online tool, and code that you can download. We're going to using the Variant Effect Predictor online tool.



Click on the spanner ![spanner] beside the Variant Effect Predictor (VEP).

This will open up a dialogue box. This allows us to input data on our variant.



The data is in the format:
Chromosome      Start  End   alleles (reference/mutation)   strand

Delete the writing already in the Paste data box and type in:
5     37017205 37017205 A/G  +

Scroll down to see some of the options we can also choose.

**Choose which database to map your variant to.**

**Find out if variants already exist in our database.**

**Choose to see scores for protein changes.**

**Choose to only see common or rare variants**

Select Prediction and Score for SIFT predictions and PolyPhen predictions. These are algorithms that predict how deleterious a mutation will be on a protein.

When you've selected everything you need, scroll right to the bottom and click Next.

Click HTML to view your results with clickable links.



| Uploaded Variation | Location | Allele | Gene | Feature | Feature type | Consequence | Position in cDNA | Position in CDS | Position in protein | Amino acid change | Codon change | Co-located Variation | Extra |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 5_37017205_A/G | 5:37017205 | G | ENSG00000164190 | ENST00000448238 | Transcript | missense_variant | 5329 | 4861 | 1621 | R/G | Aga/Gga | rs62654861 | **PolyPhen**=probably_damaging(1); **SIFT**=deleterious(0); **GMAF**=- |
| 5_37017205_A/G | 5:37017205 | G | ENSG00000164190 | ENST00000282516 | Transcript | missense_variant | 5360 | 4861 | 1621 | R/G | Aga/Gga | rs62654861 | **PolyPhen**=probably_damaging(0.994); **SIFT**=deleterious(0); **GMAF**=- |

**Our mutation affects two transcripts of one gene**

**Our mutation causes an amino acid change**

**Our mutation is already in the Ensembl database**

**Exercise: The Variant Effect Predictor (VEP)**

**Exercise 20 – VEP**

Resequencing of the genomic region of the human *CFTR* (cystic fibrosis transmembrane conductance regulator (ATP-binding cassette sub-family C, member 7) gene (ENSG00000001626) has revealed the following variants (alleles defined in the forward strand):
- G/A at 7:117,171,039
- T/C at 7:117,171,092
- T/C at 7:117,171,122

(a) Use the VEP tool in Ensembl and choose the options to see SIFT and PolyPhen predictions. Do these variants result in a change in the proteins encoded by any of the Ensembl genes? Which gene? Have the variants already been found?

(b) Go to Region in detail for *CFTR*. Do you see the VEP track?

# Comparative genomics

## Demo: Gene trees and homologues

Let's look at the homologues of human *BRCA2*. Search for the gene and go to the **Gene tab**.

Click on Gene tree (image), which will display the current gene in the context of a phylogenetic tree used to determine orthologues and paralogues.



Funnels indicate collapsed nodes. We can expand them by clicking on the node and selecting Expand this sub-tree from the pop-up menu.

We can look at homologues in the Orthologues and Paralogues pages, which can be accessed from the left-hand menu. The numbers of orthologues or paralogues available are indicated in brackets alongside the name. If there are none, then the name will be greyed out. Paralogues is greyed out for *BRCA2* indicating that there are no paralogues available.

Click on Orthologues to see the 56 orthologues available.



Choose to see only **Rodent** orthologues by selecting the box. The table below will now only show details of rodent orthologues. Let's look at mouse.



Links from the orthologue allow you to go to alignments of the orthologous proteins and cDNAs. Click on Alignment (protein) for the mouse orthologue.

**Orthologue alignment** ⓘ

Orthologue type: 1 to 1 orthologue

| Species | Gene ID | Peptide ID | Peptide length | % identity | Genomic location |
|---|---|---|---|---|---|
| Homo sapiens | ENSG00000139618 | ENSP00000439902 | 3418 aa | 50 % | 13:32889611-32973805 |
| Mus musculus | ENSMUSG00000041147 | ENSMUSP00000038576 | 3329 aa | 52 % | 5:150522630-150569746 |

*Information on orthologue pair*

*Alignment in Clustal W format*

```
CLUSTAL W(1.81) multiple sequence alignment


ENSP00000439902/1-3418      MPIGSKERPTFFEIFKTRCNKADLGPISLNWFEELSSEAPPYNSEPAEE  NNNYEPN
ENSMUSP00000038576/1-3329   MPVEYKRRPTFWEIFKARCSTADLGPISLNWFEELSSEAPPYNSEP  SEYKPHGYEPQ
                            **:   *.****:****:**..********************.****:*  :.***:

       9902/1-3418  LFKTPQRKP-SYNQLASTPIIFKEQGLTLPLYQSPVKELDKFKLDLGRNV-PNSR--HKS
       0038576/1-3329 LFKTPQRNPP-YHQFASTPIMFKERSQTLPLDQSPFREL-------GKV-VAS--SKHKT
                       *******:*   *:*:*****:***:. **** ***.:**       *:   ..    **:
```

*Protein IDs*

## Exercises: Gene trees and homologues

### Exercise 21 - Orthologues, paralogues and genetrees for the human *BRAF* gene.

(a) How many orthologues are predicted for this gene in primates? Note the Target %id and Query %id.
How much sequence identity does the *Tarsius syrichta* protein have to the human one? Click on the Alignment link next to the Ensembl identifier column to view a protein alignment in Clustal format.

(b) Go to the orthologue in marmoset. Is there a genomic alignment between marmoset and human? Is there a gene for both species in this region?

### Exercise 22 – Zebrafish orthologues

Go to www.ensembl.org to find the *dbh* gene on the zebrafish genome.
(a) Go to the Location page for this gene. View the Alignments (image) and Alignments (text) for the 5 teleost fish. Which fish genomes are represented in the alignment? Do all the fish show a gene in these alignments?

(b) Export the alignments (as Clustal).

(c) Click on the Region in detail link at the left and turn on the tracks for multiple alignments and conservation score for the 5 teleost fish EPO by configuring the page.

What is the difference between the 5 teleost fish EPO multiple alignment track and the Constrained elements already turned on by default? Which regions of the gene, do most of the constrained element blocks match up to?

Can you find more information on how the constrained elements track was generated?


**Demo: Whole genome alignments**

Let's look at some of the comparative genomics views in the Location tab. Go to the region 2:176914144-177094980 in human, which contains the *HoxD* cluster which is involved in limb development and is highly conserved between species.

In the **Region in detail** view, we can already see the Constrained elements for 36 eutherian mammals EPO_LOW_COVERAGE track by default. This track indicates regions of high conservation between species, considered to be "constrained" by evolution.



This track has a matching conservation score track. Click on Configure this page, then Comparative genomics and turn on the track for Conservation score for 36 eutherian mammals EPO_LOW_COVERAGE. Save and close the menu.



You can now see the conservation scores that were used to determine the peaks indicated in the constrained elements track.

We can also look at individual species comparative genomics tracks in this view by clicking on Configure this page.

Select BLASTz/LASTz alignments from the left-hand menu to choose alignments between closely related species. Turn on the alignments for Mouse and Chimpanzee in Normal. Go to Translated blat alignments and turn on alignments with Zebrafish and Xenopus in Normal. Save and close the menu.



Nucleotide alignments in baby pink

Protein alignments in magenta

Filled boxes are aligned sequences. Empty boxes are no alignments

The alignment is greatest between closely related species.

We can also look at the alignment between species or groups of species as text. Click on Alignments (text) in the left hand menu.



Choose an alignment from the drop-down

Multiple alignments

Pairwise alignments

Select Mouse from the alignments list then click Go.

You will see a list of the regions aligned, followed by the sequence alignment. Exons are shown in red.

This can also be viewed graphically. Click on Alignments (image) in the left-hand menu.



In both alignment views the contig is the compared species is rearranged to align to the species of interest. To compare with both contigs in their natural order, go to Region comparison.

To add species to this view, click on the blue Select species or regions button. Choose Mouse from the list then close the menu.

We can view large scale syntenic regions from our chromosome of interest. Click on Synteny in the left hand menu.

Sy<table>
</table>

Exercises: Whole genome alignments

Exercise 23 – Synteny

Go to www.ensembl.org
Find the Rhodopsin (*RHO*) gene for Human. Go to the Location tab.

(a) Click Synteny at the left. Are there any syntenic regions in dog? If so, which chromosomes are shown in this view?

(b) Stay in the Synteny view. Is there a homologue in dog for human *RHO*? Are there more genes in this syntenic block with homologues?

**Exercise 24 – Whole genome alignments**

(a) Find the Ensembl *BRCA2* (Breast cancer type 2 susceptibility protein) gene for human and go to the Region in detail page.

(b) Turn on the BLASTZ alignment tracks for chicken, chimp, mouse and platypus and the Translated BLAT alignment tracks for anole lizard and zebrafish. Does the degree of conservation between human and the various other species reflect their evolutionary relationship? Which parts of the *BRCA2* gene seem to be the most conserved? Did you expect this?

(c) Have a look at the Conservation score and Constrained elements tracks for the set of 36 mammals and the set of 19 vertebrates. Do these tracks confirm what you already saw in the tracks with pairwise alignment data?

(d) Retrieve the genomic alignment for a constrained element. Highlight the bases that match in >50% of the species in the alignment.

(e) Retrieve the genomic alignment for the *BRCA2* gene for primates. Highlight the bases that match in >50% of the species in the alignment.

# Regulation

## Demo: Raw ChIPSeq data

We're going to add some regulation data to the **Region in detail** view. We'll start at the human region 11:2012486-2030153, which contains the imprinted *H19* gene.

Add regulation tracks using Configure this page. First, we're going to add ChIP-seq data for histone modifications and polymerase binding. Click on Histones & polymerases under Regulation in the left-hand menu.

You can turn on a single track by clicking on the box in the matrix. Note that certain tracks are selected for all cell lines by default (PolII, PolIII, H3K27me3, H3K36me3, H3K4me3, H3K9me3). These will appear in the Region in detail view only if you specify a track style for the cell lines.

Turn on all the tracks for GM12878. Hover over the cell line name then select All.



Now choose the track style for the tracks you've switched on. Click on the track style box for GM12878 and select Both.



There is a similar matrix for Open chromatin &TFBS. Use this to turn on all tracks for GM12878 in Both.

Close the menu to see the tracks in the browser.

## Demo: Regulatory features and segmentation

These data are used to construct the **Reg-feats** and **Segmentation features**. The merged Reg-feats are switched on in the Region in detail view by default.

Click on Configure this page. Then select Regulatory features. Turn on the Reg. Feats: GM12878 and Reg. Segs: GM12878 tracks.

Save and close the menu.



Can you see correlations between the different kinds of regulatory data representation?

You can also add methylation data using Configure this page. Find it under DNA methylation and turn on GM12878 RRBS ENCODE and GM12878 WGBS ENCODE.

Our regulatory data incorporates the ENCODE data. To see the raw ENCODE data and the ENCODE segmentation, you need to add the ENCODE hub.

From ensembl.org, click on the ENCODE icon.



This page contains information about the ENCODE data and how it is incorporated into Ensembl.

Add the ENCODE hub by clicking on the Link to add the ENCODE track hub.

This will take you directly to the matrices for adding ENCODE data to the Region in detail view. The ENCODE matrices work in the same way as the Open chromatin &TFBS and Histones & polymerases matrices, except that some have multiple options (indicated by numbers within the boxes).

**Exercises: Regulation**

**Exercise 25 – Gene regulation: Human *STX7***

(a) Find the Location tab (Region in detail page) for the *STX7* gene. Are there regulatory features in this gene region? If so, where in the gene do they appear?

(b) Click Configure this page and on the Regulatory features menu in the left hand side. Turn on Segmentation features for HUVEC, HeLa-S3, and HepG2 cell types. Do any of these cells show predicted enhancer regions in the *STX7* region?

(c) Use Configure this page to add supporting data indicating open chromatin for HeLa-S3 cells. Are there sites enriched for marks of

open chromatin (DNase1 and FAIRE) in HeLa cells at the 5' end of *STX7*?

(d) Configure this page once again to add histone modification supporting data for the same cell type as above (e.g.HeLa-S3). Which ones are present at the 5' end of *STX7*?

(e) Is there any data to support methylated CpG sites in this region (5' end) of *STX7* in B-cells?

(f) Create a Share link for this display. Email it to yourself then open the link.


## Exercise 26 – Regulatory features in human

The *HLA-DRB1* and *HLA-DQA1* genes are part of the human major histocompatibility complex class II (MHC-II) region and are located about 44 kb from each other on chromosome 6. In the paper 'The human major histocompatibility complex class II *HLA-DRB1* and *HLA-DQA1* genes are separated by a CTCF-binding enhancer-blocking element' (Majumder *et al* J Biol Chem. 2006 Jul 7;281(27):18435-43) a region of high acetylation located in the intergenic sequences between *HLA-DRB1* and *HLA-DQA1* is described. This region, termed XL9, coincided with sequences that bound the insulator protein CCCTC-binding factor (CTCF). Majumder *et al* hypothesise that the XL9 region may have evolved to separate the transcriptional units of the *HLA-DR* and *HLA-DQ* genes.

(a) Go to the region from 32,540,000 to 32,620,000 bp on human chromosome 6

(b) Is there a regulatory feature annotated in the intergenic region between the *HLA-DRB1* and *HLA-DQA1* genes that has CTCF binding supporting data as (part of) its core evidence?

(c) Has the CTCF binding detected at this position been observed in all cell/tissue types analysed?

(d) Have a look at the Regulatory supporting evidence - Histones & Polymerases configuration matrix. For which cell/tissue type are the most histone acetylation data sets available? In this cell/tissue type,

is the region that shows CTCF binding also a region of high acetylation, as found by Majumder *et al*?

# Quick Guide to Databases and Projects

Here is a list of databases and projects you will come across in these exercises. Google any of these to learn more. 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, NCBI GenBank, DDBJ –** 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 each.

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

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

**NCBI RefSeq –** NCBI's collection of 'reference sequences', includes genomic DNA, transcripts and proteins. NM stands for 'Known mRNA' (eg NM_005476) and NP (eg NP_005467) are 'Known proteins'.

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

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

**UniProt TrEMBL –** the automatically annotated, unreviewed set of proteins (EMBL-Bank translated). Varying quality.

**VEGA –** Vertebrate Genome Annotation, a selection of manually-curated genes, transcripts, and proteins. *(human, mouse, zebrafish, gorilla, wallaby, pig, and dog).*

**VEGA-HAVANA –** The main contributor to the VEGA project, located at the Wellcome Trust Sanger Institute, Hinxton, UK.

## GENE NAMES

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

**ZFIN –** The Zebrafish Model Organism Database. Gene names are only one part of this project. *(Z-fish).*

## PROTEIN SIGNATURES
**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 (explained below).

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

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