

# Access to genes and genomes with Ensembl



## Course Manual

Nov, 2008

## **CONTENTS**

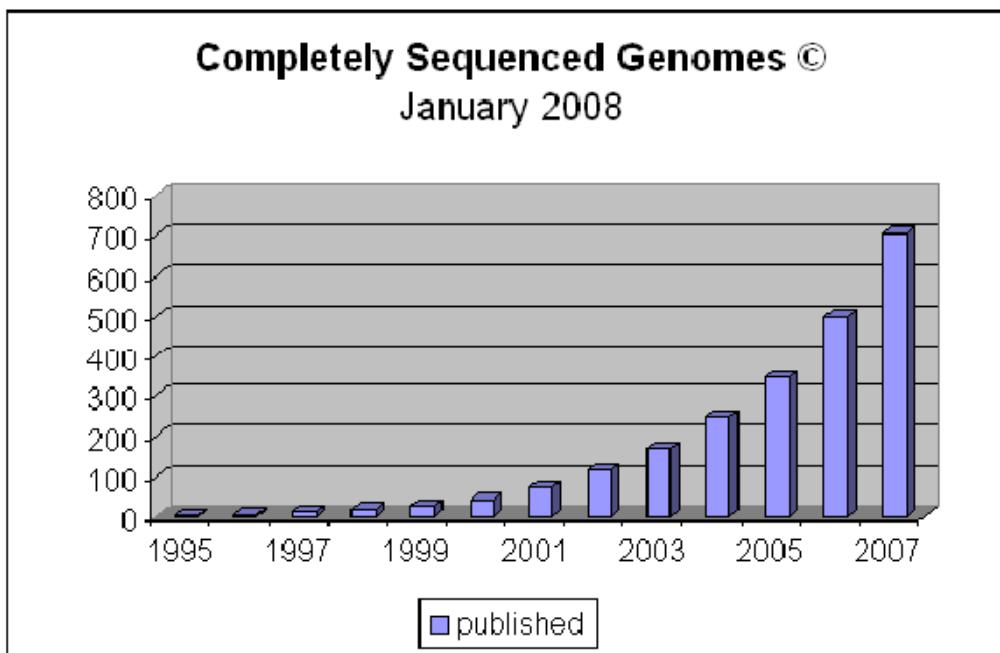
I) INTRODUCTION.....	2
II) BROWSING ENSEMBL – Worked example.....	7
III) BROWSING ENSEMBL – Exercises.....	27
Answers.....	28
IV) BIOMART – Worked example.....	31
V) BIOMART – Exercises.....	39
Answers.....	40
VI) EVALUATING GENES AND TRANSCRIPTS (GENEBUILD)	
Exercises.....	43
Answers.....	44
VII) COMPARATIVE GENOMICS	
Exercises.....	46
Answers.....	48
VIII) VARIATIONS	
Exercises.....	50
Answers.....	51
IX) TYING IT TOGETHER - CASE STUDIES.....	53
X) Quick fact sheet.....	55

## I) Introduction

Ensembl is one of the world's primary resources for genomic research, a resource through which scientists can access the human genome as well as the genomes of other model organisms. Because of the complexity of the genome and the many different ways in which scientists want to use it, Ensembl has to provide many levels of access with a high degree of flexibility. Through the Ensembl website a wet-lab researcher with a simple web browser can for example perform BLAST searches against chromosomal DNA, download a genomic sequence or search for all members of a given protein family. But Ensembl is also an all-round software and database system that can be installed locally to serve the needs of a genomic centre or a bioinformatics division in a pharmaceutical company enabling complex data mining of the genome or large-scale sequence annotation.

### The need for automatic annotation

Recent years have seen the release of huge amounts of sequence data from genome sequencing centres (figure 1). However, this raw sequence data is most valuable to the



*Figure 1. Completely sequenced genomes as of Jan, 2008 (figure taken from <http://www.genomesonline.org>).*

laboratory biologist when provided along with quality annotation of the genomic sequence. This information can be the starting point for planning experiments, interpreting Single Nucleotide Polymorphisms, inferring the function of gene products, predicting regulatory sites for gene expression and

so on. The currently agreed 'gold standard' for the annotation of eukaryotic genomes is annotation made by a human being. This so-called "manual annotation" is based on information derived from sequence homology searches, the results of various *ab initio* gene prediction methods and literature searches. Annotation of large genomes (such as mouse and human) that meet this standard is slow and labour intensive, taking large teams of annotators years to complete. As a result, the annotation can almost never be entirely up-to-date and free of inconsistencies (as the annotation process usually begins before the sequencing process is complete). Hence, an automated annotation system is desirable since it is a relatively rapid process that allows frequent updates to accommodate new data. To meet this need, we produced the Ensembl annotation system by observing how annotators build gene structures and condensing this process into a set of rules.

### The start of Ensembl

Ensembl's genesis was in response to the acceleration of the public effort to sequence the human genome in 1999. At that point it was clear that if annotation of the draft sequence was to be available in a timely fashion it would have to be automatically generated and that new software systems would be needed to handle genome data sets that were much larger, much more fragmented and much more rapidly changing than anything previous dealt with.

Ensembl was conceived in three parts: as a scalable way of storing and retrieving genomic data; as a web site for genome display; and as an automatic annotation method based around a set of heuristics. It was initially written for the draft human genome, which was sequenced clone-by-clone but has also been successfully used for whole genome shotgun assemblies. The storage and display parts of Ensembl are used for all the genomes currently present in Ensembl, while the automatic gene annotation has been run for most of the genomes with the exception of Takifugu, Tetraodon, Fruitfly, *C. elegans* and Yeast.

Over the past few years Ensembl has grown into a large scale enterprise, with substantial computing resources enabling it to process and provide live database access to currently more than 25 different genomes (figure 2) and a bimonthly update frequency to its website. It has a large community of users in both industry and academia, using it as a base for their individual organisation's experimental and computational genome based investigations, some of which maintain their own local installations.

Ensembl is a collaboration between the European Bioinformatics Institute (EBI) and the Wellcome Trust Sanger Institute, both located on the Wellcome Trust Genome Campus in Hinxton, Cambridge, UK. Ensembl is funded principally by the Wellcome Trust, with additional funding from the European Molecular Biology Laboratory (EMBL), the National Institutes of Health – National Institute of Allergy and Infectious Disease (NIH-NIAID) and the Biotechnology and Biological Sciences Research Council (BBSRC).

## The Ensembl software and database system

As a software/database system Ensembl can be best described as a hybrid of a scripting programming language (Perl) and a relational database (MySQL, pronounced "My Ess Que Ell").).

Ensembl Perl software inherits from a tradition of biological object-design developed through BioPerl (<http://www.bioperl.org/>). This means that developers at Ensembl aimed at creating reusable pieces of software that would faithfully describe biological entities such as gene, transcript, protein, genomic clone or chromosome. Rules of usage and design of Ensembl and BioPerl objects can be best learned while using them, browsing their code and through a bit of trial-and-error. There is a comprehensive BioPerl tutorial available at the BioPerl website.

The Ensembl database is based on a relational database called MySQL. SQL in MySQL stands for 'Structured Query Language', a universal database programming language shared by many relational databases. Because MySQL is available free of charge for non-commercial developers, every academic centre can install its own local copy of MySQL (MySQL server) and download Ensembl data from the Ensembl ftp site. Simple queries of the database can be handled using the SQL language (see appendix), but for complex queries demanded by most biological analyses the Ensembl MySQL server is best accessed using Ensembl Perl objects.

## The Ensembl annotation pipeline

The Ensembl analysis and annotation pipeline is based on a rule set of heuristics that a human annotator would use. All Ensembl gene predictions are based on experimental evidence, which is imported via manually curated UniProt/Swiss-Prot, partially manually curated NCBI RefSeq and automatically annotated UniProt/TrEMBL records. Untranslated regions (UTRs) are annotated to the extent supported by EMBL mRNA records. As there is no guarantee that UTR sequences in EMBL records are complete there is similarly no guarantee that the Ensembl genome analysis and annotation pipeline has enough biological evidence to predict complete UTR regions. For a limited number of species regulatory regions are annotated, but this annotation isn't very extensive yet as the set of well-characterised promoters is still small and there is currently no algorithm yielding reliable results on a genomic scale.

## The Ensembl website

Ensembl provides easy access to genomic information with a number of visualisation tools. The Ensembl website gives you for example the possibility to directly download data, whether it is a DNA sequence of a genomic contig you are trying to identify novel genes in, or positions of SNPs in a gene you are working on. The key Ensembl web pages are called Views (e.g. GeneView, ContigView and SNPView), and will all be introduced appropriately later on. An updated version of the website is released bimonthly. Old

versions are for at least two years accessible on the ‘Archive!’ website. Apart from that the ‘Pre!’ website provides displays of genomes that are still in the process of being annotated. There is also an ftp site to download large amounts of data from the Ensembl database, as well as the data-mining tool BioMart, that allows rapid retrieval of information from the databases, and BLAST/BLAT sequence searching and alignment.

## Further reading

Flicek, P. et al.

### **Ensembl 2008**

Nucleic Acids Res. Jan 2008; 36: D707 - D714

Spudich, G., Fernández-Suárez, X. M., and Birney, E.

### **Genome Browsing with Ensembl: a practical overview**

Brief Funct Genomic Proteomic, 2007 Sept; 6: 202-219

Fernández Suárez X. M. and Schuster M.

### **Using the Ensembl Genome Server to Browse Genomic Sequence Data.**

*Current Protocols in Bioinformatics*, UNIT 1.15, January 2007.

Hubbard, T.J.P. et al.

### **Ensembl 2007**

Nucleic Acids Res. 2007 (*Database Issue*)

Birney, E. et al.

### **Ensembl 2006.**

Nucleic Acids Res. 2006 Jan 34:D556-D561 (2006)

Hubbard, T. et al.

### **Ensembl 2005.**

Nucleic Acids Res. 2005 33 D447-D453 (2005)

Birney, E. et al.<sup>1</sup>

### **An Overview of Ensembl.**

Genome Research 14(5): 925-928 (2004)

Kasprzyk, A. et al.

### **EnsMart: a generic system for fast and flexible access to biological data.**

Genome Research (2004) 14:1, 160-9.

Ashurst, J. L. et al.

### **The Vertebrate Genome Annotation (Vega) database.**

Nucl. Acids Res. 33:D459-D465 (2005)

\* Additional references can be found here:

<http://www.ensembl.org/info/about/publications.html>

---

<sup>1</sup> This paper was part of the May 2004 issue of *Genome Research* which included an Ensembl special covering detailed aspects of the Ensembl web site, the underlying scalable database system for storing genome sequence and annotation information, as well as the automated genome analysis and annotation pipeline

SPECIES: Mammals		ASSEMBLY		GENEBUILD	
Armadillo	<i>Dasypus novemcinctus</i>	ARMA	May-05	Ensembl	Aug-06
Bushbaby	<i>Otolemur garnettii</i>	otoGar1	May-06	Ensembl	Feb-07
Cat	<i>Felis catus</i>	CAT	Feb-07	Ensembl	Jun-06
Chimpanzee	<i>Pan troglodytes</i>	PanTro 2.1	Mar-06	Ensembl	May-07
Cow	<i>Bos taurus</i>	Btau 3.1	Feb-07	Ensembl	Sep-06
Dog	<i>Canis familiaris</i>	CanFam 2.0	May-06	Ensembl	Dec-06
Elephant	<i>Loxodonta africana</i>	BROAD E1	May-05	Ensembl	Aug-06
Guinea pig	<i>Cavia porcellus</i>	cavPor2	Feb-07	Ensembl	Oct-06
Hedgehog	<i>Erinaceus europaeus</i>	eriEur1	Feb-07	Ensembl	Oct-06
Horse	<i>Equus caballus</i>	Equus1	Sep-07	Ensembl	Sep-07
Human	<i>Homo sapiens</i>	NCBI 36	Oct-05	Ensembl	Sep-07
Lesser hedgehog tenrec	<i>Echinops telfairi</i>	TENREC	May-05	Ensembl	Aug-06
Microbat	<i>Myotis lugigfugus</i>	myoLuc1	Mar-06	Ensembl	Jan-07
Mouse	<i>Mus musculus</i>	NCBI m37	Apr-07	Ensembl	Sep-07
Mouse Lemur	<i>Microcebus murinus</i>	micMur1	Jun-07	Ensembl	Jul-07
Opossum	<i>Monodelphis domestica</i>	MonDom 5.0	Oct-06	Ensembl	Feb-07
Orangutan	<i>Pongo pygmaeus</i>	PPYG2	Sep-07	Ensembl	Oct-07
Pig*	<i>Sus scrofa</i>	Sscrofa1			
Pika	<i>Ochotona princeps</i>	OchPri2.0	Jun-07	Ensembl	Jul-07
Platypus	<i>Ornithorhynchus anatinus</i>	OANA 5	Dec-05	Ensembl	Jan-07
Rabbit	<i>Oryctolagus cuniculus</i>	RABBIT	May-05	Ensembl	Aug-06
Rat	<i>Rattus norvegicus</i>	RGSC 3.4	Dec-04	Ensembl	Feb-06
Rhesus macaque	<i>Macaca mulatta</i>	MMUL 1	Feb-06	Ensembl	Aug-06
Shrew	<i>Sorex araneus</i>	sorAra1	Oct-05	Ensembl	Apr-07
Squirrel	<i>Spermophilus tridecemlineatus</i>	speTri1	Jun-06	Ensembl	Oct-06
Tree shrew	<i>Tupaia belangeri</i>	tupBel1	Feb-07	Ensembl	Oct-06
Species: Other					
Aedes	<i>Aedes aegypti</i>	AaegL 1	Oct-05	VectorBase	Jun-06
Anole Lizard*	<i>Anolis carolinensis</i>	AnoCar1.0	Feb-07		
Anopheles	<i>Anopheles gambiae</i>	AgamP 3	Feb-06	VectorBase	Jun-07
<i>C. elegans</i>	<i>Caenorhabditis elegans</i>	WS 180	Sep-07	WormBase	Sep-07
Chicken	<i>Gallus gallus</i>	WASHUC 2	May-06	Ensembl	Aug-06
<i>C. intestinalis</i>	<i>Ciona intestinalis</i>	JGI 2	Mar-05	Ensembl	Feb-06
<i>C. savignyi</i>	<i>Ciona savignyi</i>	CSAV 2.0	Oct-05	Ensembl	Apr-06
Fruitfly	<i>Drosophila melanogaster</i>	BDGP 5.4	Nov-07	FlyBase	Oct-07
Lamprey*	<i>Petromyzon marinus</i>	PMAL3			
Medaka	<i>Oryzias latipes</i>	HdrR 1	Oct-05	Ensembl	May-06
<i>S. cerevisiae</i>	<i>Saccharomyces cerevisiae</i>	SGD 1.01	Dec-06	SGD	Dec-06
Stickleback	<i>Gasterosteus aculeatus</i>	BROAD S1	Feb-06	Ensembl	Jun-06
Takifugu	<i>Takifugu rubripes</i>	FUGU 4.0	Jun-05	Ensembl	Nov-07
Tetraodon	<i>Tetraodon nigroviridis</i>	TETRAODON 7	Apr-03	Genoscope	Sep-04
<i>X. tropicalis</i>	<i>Xenopus tropicalis</i>	JGI 4.1	Aug-05	Ensembl	Nov-05
Zebrafish	<i>Danio rerio</i>	Zv 7	Apr-07	Ensembl	Jun-07

Table 1 – Species, assemblies and gene sets in Ensembl  
(\*currently available in Pre! website [pre.ensembl.org](http://pre.ensembl.org))

## II) WORKED EXAMPLE – A walk through the main pages of the Ensembl browser, using the IL2 (Interleukin-2 precursor) gene as an example.

**STEP 1:**  
Load Ensembl  
[www.ensembl.org](http://www.ensembl.org)

Navigation column

Search

**STEP 2:**  
Click on  
“Human”

Help pages  
and  
Documents

What's new

The screenshot shows the Ensembl homepage for release 47 (October 2007). The page includes a navigation column on the left with links for 'Your Ensembl' (Login or Register, About User Accounts), 'Help & Documentation' (About Ensembl, Genomic Data, Help & Information, Software), and 'Ensembl Archive' (View previous release of page in Archive!, Stable Archive! link for this page). The main content area features a 'Search Ensembl' bar at the top with a dropdown set to 'All species' and a 'Go' button. Below it is a 'Search Ensembl' section with a placeholder 'e.g. mouse chromosome 2 or rat X:10000..20000 or human gene BRCA2'. The 'Ensembl tools' section contains links for 'Start a sequence search', 'Mine Ensembl with BioMart', 'Customise Your Ensembl', and 'Fetch data with the Ensembl API'. The 'About Ensembl' section provides information about the project, its funding, and access to data. The 'Other Ensembl websites' section lists links to archipel, VEGA, Ensembl PreL, EBI Genome Reviews database, and Trace server. The 'What's new' section highlights recent news items: 'New Mouse 37 assembly and genebuild (Mus musculus)', 'New genebuild on human assembly NCBI 36 (Homo sapiens)', 'WormBase 180 (Caenorhabditis elegans)', 'Functional Genomics (H. sapiens, M. musculus)', and 'Variation updates (H. sapiens, M. musculus, R. norvegicus, O. anatinus)'. The right side of the page features a 'Help' section with a 'PreL species' dropdown set to 'Human' (NCBI m37 | Vega), 'Mouse' (NCBI m37 | Vega), and 'Zebrafish' (Zv6 | Vega). The footer includes a 'HOME - BLAST - SITEMAP HELP' menu and copyright information: '© 2007 WTSI / EBI. Ensembl is available to download for public use - please see the code licence for details.'

**STEP 3:**  
Type in 'IL2 Gene'.  
Click 'Go'.

**Karyotype**

Ensembl release 44 - April 2005

Your Ensembl

- Login or Register
- About User Accounts

Help & Documentation

- About Ensembl
- Genomic Data
- Help & Information
- Software

Select a species

- Mammals
- Other chordates
- Other eukaryotes

Ensembl Archive

- View previous release of page in Archive!
- Stable Archive! link for this page

Sanger EBI

**Microbat**  
*Myotis lucifugus*  
2X assembly and genebuild

Search Ensembl Homo sapiens

Search: e.g. chromosome X or 14:10000..200000 or BBCA2

Karyotype

Click on a chromosome for a closer view

About the Human genome

Assembly

This release is based on the NCBI 36 assembly of the [human genome](#) (November 2005). The data consists of a reference assembly of the complete genome plus the Celera vGSS and a number of alternative assemblies of individual haplotypes or regions.

[Full list of assemblies](#)

The International Human Genome Sequencing Consortium have published their scientific analysis of the finished human genome:

- [Nature 431, 931 - 945 \(21 October 2004\)](#)
- [WT Sanger Institute Press Release](#)

Annotation

The human genome sequence is now considered sufficiently stable that the three major genomic browsers have come together to produce a common set of gene IDs for their annotations. This Consensus CDS ID set has been incorporated into the Ensembl database alongside the existing identifiers.

- [More information about the CDS project](#)

The [ENCODE](#) (ENCYclopedia Of DNA Elements) project aims to find functional elements in the human genome.

- [More information about the ENCODE project](#)

Jump directly to sequence position

Chromosome:  or region

From (bp):

To (bp):  Go

What's New in Ensembl 44

**Homo sapiens News**

- Patch for Ensembl Human database
- Ensembl *Homo sapiens* has been patched with a few extra transcripts and some new CCDS IDs, and the corresponding xrefs and variation databases have been updated.
- cDNA Updates
- Ensembl human and mouse databases have received their usual cDNA updates.
- Vega updates
- Vega human and mouse have both been updated since the last release of Ensembl. See [VEGA](#) for more information.
- Variation updates
- All species with variation data now have a failed\_variation table. Also, in species that have duplicate variations (with the same mapping but different IDs), these have been put into the variation\_synonym table.

**General News**

- ncRNAs for Ensembl chordates
- All Ensembl genebuild databases (i.e. excluding the imported invertebrate databases for mosquitoes, fruit fly, worm and yeast) have been updated with new ncRNA data.

[More news...](#)

Statistics

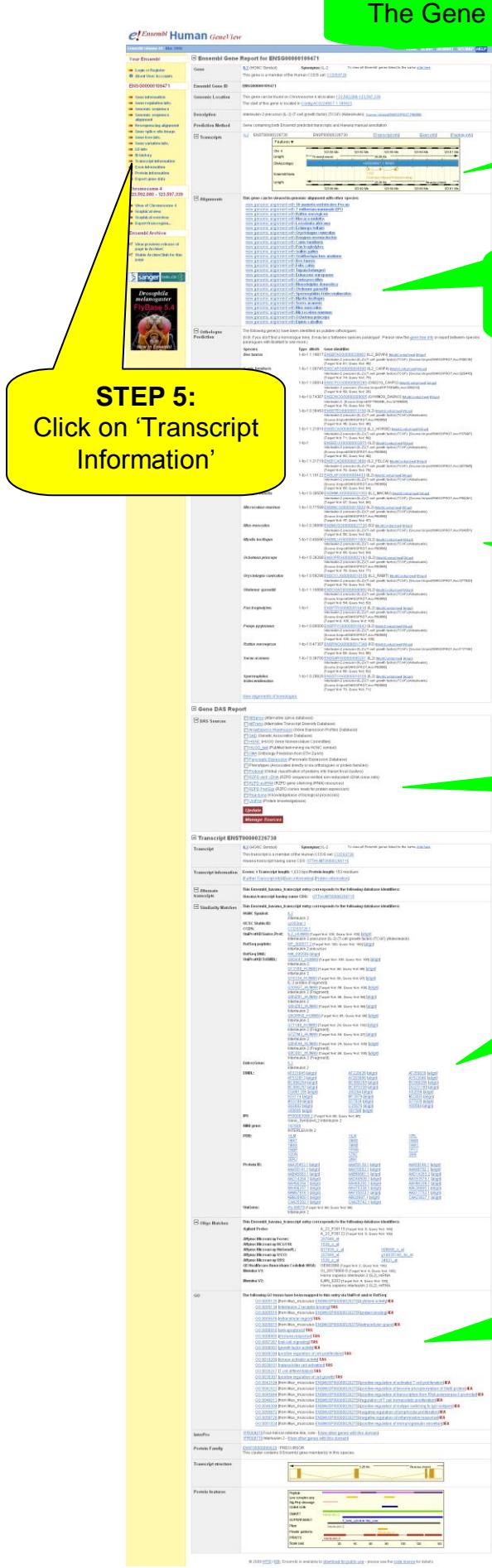
Assembly:	NCBI 36, Oct 2005
Genebuild:	Ensembl, Aug 2006
Database version:	44.36f
Known genes:	21,724
Novel genes:	1,017
Pseudogenes:	1,040
RNA genes:	4,113
Immunoglobulin/T-cell receptor gene segments:	388
Genscan gene predictions:	69,185
Gene exons:	270,214
Gene transcripts:	44,567
SNPs:	11,561,833
Base Pairs*:	3,253,037,807
Golden Path Length**:	3,093,120,360

\* Total number of base pairs = sum of lengths of DNA table  
\*\* Reference assembly (Golden path) length = sum of non-redundant top level seq regions

© 2007 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

**Source and version of assembly and genebuild**





## The Gene View Page

## Gene Model

# Whole Genome Alignments

## **STEP 5:** Click on 'Transcript Information'

## Orthologues in other species

## External Information (DAS)

IDs in other  
databases

GO  
(Gene Ontology)  
terms

## Protein domains

**STEP 7:** Click on 'Exon information'

A link back to the GeneView page.

Transcript report: each transcript for a gene has a 'TransView' page.

Spliced transcript sequence

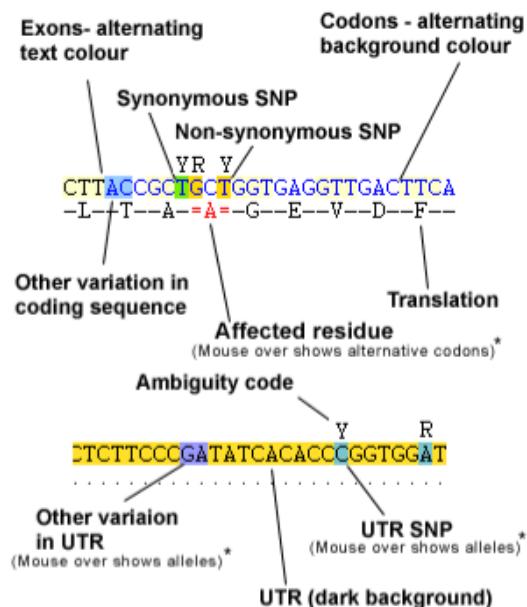
STEP 6:  
Select 'Exons, Codons, Translations and SNPs'.  
Select 'Number residues: Yes' and click on [Refresh]

## Result of STEP 6:

```

121 AAGGTAAATGTTTTTCAGACAGGTAAAGTCTTGAAAAATATGTGTAAATATGTAAAACATT
-----
181 TTGACACCCCCATAATATTTTCCAGAATT&ACAGTATAAATTGCATCTCTTGTTCAGA
-----
241 GTTCCCTATCACTCTTTAATCACTACTCACAGTAACCTCAACTCCTGCCACAATGTAC
-----M--Y-
301 AGGATGCAACTCCTGTCCTGCATTGCACTAAGTCTTGCACTTGTACAAACAGTGCACCT
3 -R--M--Q--L--L--S--C--I--A--L--S--L--A--L--U--T--N--S--A--P-
361 ACTTCAAGTTCTACAAAGAAAAACACAGCTACAACGGAGCATTTACTKKGCTGGATTTACAG
23 -T--S--S--S--T--K--T--Q--L--E--H--L--=L--L--D--L--Q-

```



## Result of STEP 7:

**STEP 9:**  
Click on  
'Graphical view'

**STEP 8:**  
Choose 'Flanking sequence at either end of transcript – 500', tick 'Show full intronic sequence' and click on [Go]

Flank (green)  
UTR (purple)  
Intron (blue)  
Coding sequence (black)

**Supporting Evidence**

The supporting evidence below consists of the sequence matches on which the exon predictions were based and are sorted by alignment score.  
There are a large number of supporting evidence hits for this transcript. Only the top ten 10 hits have been shown. [Click here to view all 22 supporting evidence hits.](#)

Score	100	>=99	>=97	>=90	>=75	>=50	<=50	NO EVIDENCE
BC070398.1	1	2	3	4				
BC066254.1								
AF228639.1								
AY229630.1								
BC066257.1								
BC066256.1								
BC066255.1								
AY523040.1								
AY236966.1								
CD559408.1								
CD52012.1								

© 2008 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

## **Result of STEP 8:**

## Result of STEP 9: ContigView

**Ensembl Human ContigView**

**Chromosome** Chromosome 4

**Markers**

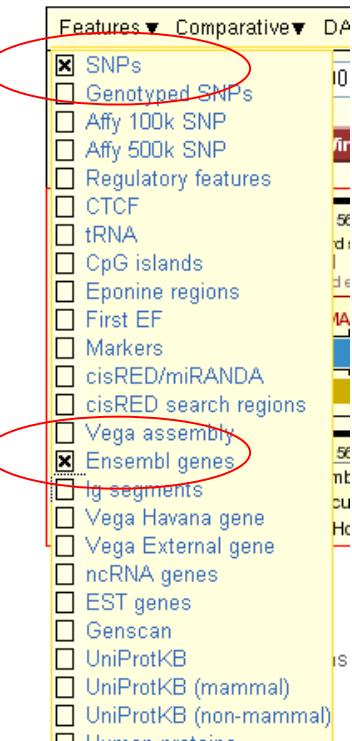
**1 Mb region**

**Assembly**

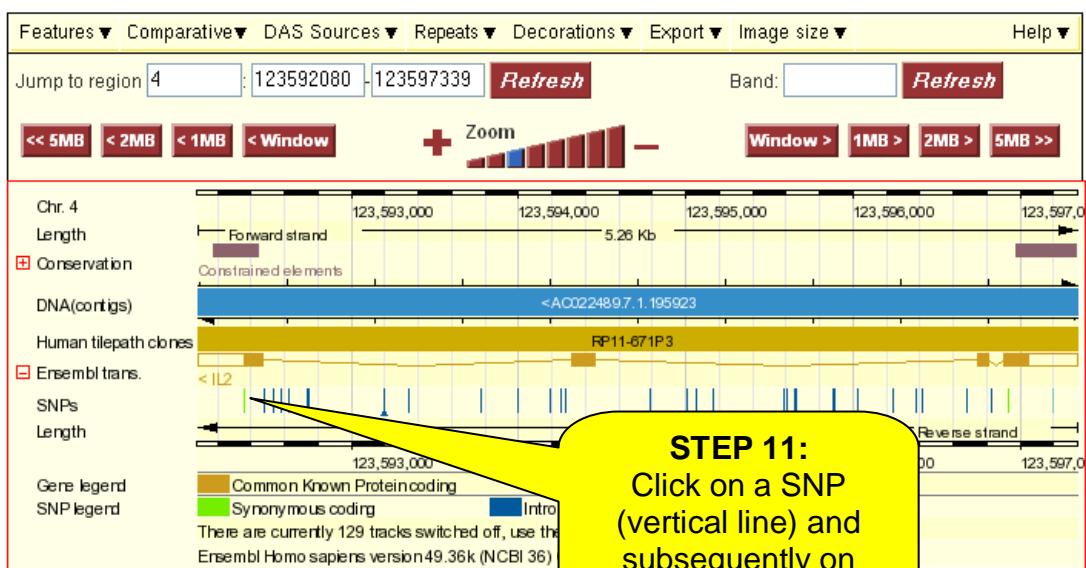
**IL2 transcript**

**STEP 10:**  
In the 'Features' drop-down menu select 'SNPs' and 'Ensembl genes' (deselect other options). Close the menu to view new tracks.

**(STEP 10):** Deselect all other options from the features menu and select 'SNPs' and 'Ensembl genes'.



## RESULT OF STEP 10



**Ensembl Human SNPView**

Ensembl release 49 - Mar 2008

Your Ensemble

- Login or Register
- About Us
- dbSNP: rs1051753
- GeneSNP Info
- rs1051753 - SNP Info
- rs1051753 - LD Info
- Chromosome 4 123,592,363
- View of Chromosome 4
- Graphical view
- Graphical overview
- Export from region...

Ensembl Archive

- View previous release of page in Archive!
- Stable Archive link for this page

Sanger EMBL-EBI

**dbSNP identifier**

**Alleles**

**Allele and genotype frequencies**

**SNP of interest (boxed) and other SNPs in the region**

rs1051753 (dbSNP128)

HGVSbase SNP000618875

C/T (ambiguity code: Y)  
Ancestral allele: C

Validation status: Unknown

Linkage disequilibrium data: No linkage data for this SNP

Flanking sequence: TAAATAGAAGGGCTGATATGTTTAACTGGGAAGCAGTTAATTACAAAGT**Y**AGCTGTGAG  
ATGATGCTTGGACAAAAGGTAACTCCACTGTTTCAGAAAATTTC (SNP highlighted)

SNP rs1051753 is located in the following transcripts:

Genomic location (strand)	Gene	Transcript	Position	Translation	AA Type	GeneSNPView
4: 123592363-123592363 (+)	ENSG00000109471	ENST00000226730	750-750	ENSP00000226730	152-152	L SYNONYMOUS_CODING SNPs in gene context

Population genotypes and allele frequencies:

Population	Alleles	Alleles	Description
SEQUENOM_CEPH	C	T	1,000 unknown comprised of UTAH (93%), French (4%), and Venezuelan (3%) samples were purchased from Genomic DNA samples were obtained for a panel of 92 unrelated individuals. Coriell Cell Repository and pooled in equimolar amounts.

Individual genotypes for SNP rs1051753

SNP Context - 4 123592363

Features ▾ Source ▾ SNP class ▾ SNP type ▾ Decorations ▾ Export ▾ Image size ▾

No EST transcripts in this region

No ncRNAs in this region

No Ensembl transcripts in this region

No Ensembl transcripts in this region

Forward strand 20.00 Kb Reverse strand

No Ensembl transcripts in this region

No ncRNAs in this region

18 of the 157 variations in this region have been filtered out by the Source, Class and Type menus.

SNPs

Gerotyped SNPs

SNP legend: Intergenic, Intronic, Upstream, Downstream, Non-synonymous coding, 5' UTR

123.58 Mb 123.59 Mb 123.59 Mb 123.60 Mb 123.60 Mb

© 2008 WTSI / EBI. Ensembl is available to download for public use - please see the code licence for details.

**STEP 12:**  
Go back to  
ContigView with the  
back button of the  
internet browser.

**STEP 13:**  
To see the same chromosomal region in the UCSC genome browser, click on 'View region at UCSC' on the left of the page. A new window will open.

UCSC Genome Browser on Human Mar. 2006 Assembly

position/search chr4:123,592,080-123,597,339 | jump | clear | size 5,260 bp | configure

IL2 gene

Chr4: 123,592,080 - 123,597,339

UCSC Gene Predictions Based on RefSeq, UniProt, GenBank, and Comparative Genomics

RefSeq Genes

Human mRNAs

Spliced ESTs

Mammal Cons

Rhesus

Mouse

Dog

Horse

African Green Monkey

Opossum

P. troglodytes

Chicken

X. tropicalis

Stickleback

SNPs (128)

RepeatMasker

Simple Nucleotide Polymorphisms (obs/Exp > 128)

Repeating Elements by RepeatMasker

Click on a feature for details. Click on base position to zoom in around cursor. Click gray/blue bars on left for track options and descriptions.

move start < 2.0 > move end < 2.0 >

default tracks hide all add custom tracks configure reverse refresh

Use drop-down controls below and press refresh to alter tracks displayed.

Tracks with lots of items will automatically be displayed in more compact modes.

**Mapping and Sequencing Tracks**

Base Position	Chromosome Band	STS Markers	FISH Clones	Recomb Rate
dense	hide	hide	hide	hide
Map Contigs	Assembly	Gap	Coverage	BAC End Pairs
hide	hide	hide	hide	hide
Fosmid End Pairs	GC Percent	Short Match	Restr Enzymes	
hide	hide	hide	hide	

**Phenotype and Disease Associations**

GAD	Case Control	NIMH Bipolar	RGD Human QTL	RGD Rat QTL
hide	hide	hide	hide	hide
MGI Mouse QTL				
hide				

**Genes and Gene Prediction Tracks**

UCSC Genes	Old Known Genes	Alt Events	CCDS	RefSeq Genes
pack	hide	hide	hide	dense
Other RefSeq	MGC Genes	ORFeome Clones	Ensembl Genes	AccView Genes
hide	pack	hide	hide	hide
SIB Genes	N-SCAN	CONTRAST	CDP Genes	Geneid Genes
hide	hide	hide	hide	hide
GenScan Genes	Exoniphy	Augustus	RNA Genes	Superfamily
hide	hide	hide	hide	hide
ACEScan	EvoFold	sno/mRNA		
hide	hide	hide		

**mRNA and EST Tracks**

Human mRNAs	Spliced ESTs	Human ESTs	Other mRNAs	Other ESTs
dense	dense	hide	hide	hide
H-Inv	UniGene	Gene Bounds	SIE Alt-Splicing	Poly(A)
hide	hide	hide	hide	hide

refresh

**STEP 14:**  
Turn on 'Ensembl genes' by changing 'hide' to 'full' and clicking 'refresh' at the bottom of the page.

**Close the window to return to ContigView.**

**e! Ensembl Human ContigView**

Ensembl release 49 - Mar 2008

**Your Ensembl**

- Login or Register
- About User Accounts

**Chromosome 4 123,592,080 - 123,597,339**

- View of Chromosome 4
- Graphical view
- Graphical overview
- Resequencing alignment
- View alignment with ...
- View alongside ...
- View Syntenic regions ...
- View region at UCSC
- View region at NCBI

**Export**

**STEP 15:**  
Click on 'View Syntenic regions ... with *Mus musculus*'

**Orangutan**  
*Pongo pygmaeus abelii*  
Now in Ensembl

**Chromosome 4**

**Overview**

**Detailed view**

**Basepair view**

Chr. 4 band

DNA(contigs)

Markers

Ensembl Genes

ncRNA Genes

Gene legend

Ensembl Known Protein in Coding

Merged Known Protein in coding

RNA Pseudogene (Novel)

Features ▾ Comparative ▾ DAS Sources ▾ Repeats ▾ Decorations ▾ Export ▾ Image size ▾ Help ▾

Jump to region 4 123592080 123597339 Refresh Band: Refresh

<< 5MB < 2MB < 1MB < Window + Zoom Window > 1MB > 2MB > 5MB >>

Chr. 4 Length

Conservation

DNA(contigs)

Human BACpath clones

Ensembl transcripts

SNPs Length

Gene legend

SNP legend

Common Known Protein coding

Synonymous coding

Intronic

5' UTR

There are currently 129 tracks switched off, use the menus above the image to turn them on.

Ensembl Homo sapiens version49 36k (NCBI 36) Chromosome 4 123,592,080 - 123,597,339

© 2008 WTSI / EBI. Ensembl is available to download for public use - please see the code licence for details.

Mouse chromosomes

e! Ensembl Human Synteny View

Ensembl release 49 - Mar 2008

Your Ensembl

Login or Register  
About User Accounts

Chromosome 4

View Chromosome 4  
View Chr 4 Synteny  
Map your data onto this chromosome

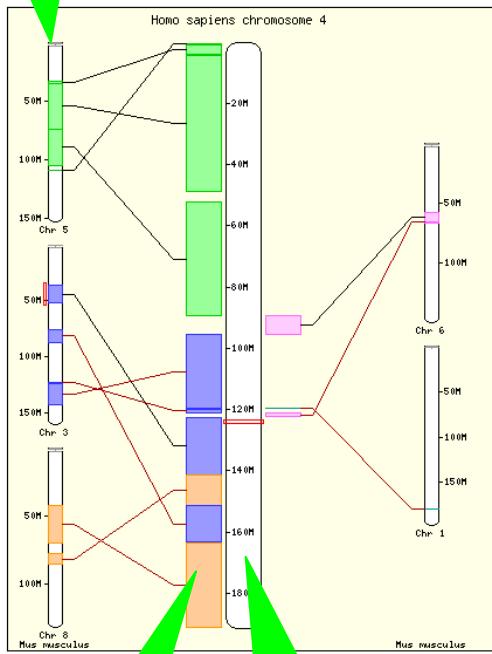
Ensembl Archive

View previous release of page in Archive!  
Stable Archive! link for this page

welcome trust sanger EMBL-EBI



Syntenic block



Human chromosome

Human genes

Mouse homologues

E.g. Chr 1's Synteny with Pan

Homology Matches

Homo sapiens Genes	Mus musculus Homologues
IL2 (0.12 Gb) [ContigView]	-> <a href="#">IL2</a> (3: 37.02 Mb) [ContigView] [MultiContigView]
IL21 (0.12 Gb) [ContigView]	-> <a href="#">IL21</a> (3: 37.12 Mb) [ContigView] [MultiContigView]
BBS12 (0.12 Gb) [ContigView]	-> <a href="#">Bbs12</a> (3: 37.21 Mb) [ContigView] [MultiContigView]
FGFR2 (0.12 Gb) [ContigView]	-> <a href="#">Fgf2</a> (3: 37.25 Mb) [ContigView] [MultiContigView]
NUDT6 (0.12 Gb) [ContigView]	-> <a href="#">Nudt6</a> (3: 37.30 Mb) [ContigView] [MultiContigView]
SPATA5 (0.12 Gb) [ContigView]	-> <a href="#">Spata5</a> (3: 37.31 Mb) [ContigView] [MultiContigView]
SPRY1 (0.12 Gb) [ContigView]	->
ANKRD50 (0.13 Gb) [ContigView]	->
BLYM_HUMAN (0.13 Gb) [ContigView]	No homologues
ENSG00000214786 (0.13 Gb) [ContigView]	No homologues
FAT4 (0.13 Gb) [ContigView]	-> <a href="#">Fat4</a> (3: 38.79 Mb) [ContigView] [MultiContigView]
INTU (0.13 Gb) [ContigView]	-> <a href="#">Intu</a> (3: 40.43 Mb) [ContigView] [MultiContigView]
SLC25A31 (0.13 Gb) [ContigView]	-> <a href="#">Sic25a31</a> (3: 40.51 Mb) [ContigView] [MultiContigView]
HSPA4L (0.13 Gb) [ContigView]	-> <a href="#">Hspa4l</a> (3: 40.55 Mb) [ContigView] [MultiContigView]
PLK4 (0.13 Gb) [ContigView]	-> <a href="#">Plk4</a> (3: 40.60 Mb) [ContigView] [MultiContigView]

STEP 16:  
Click on  
[MultiContigView]

Navigate Homology

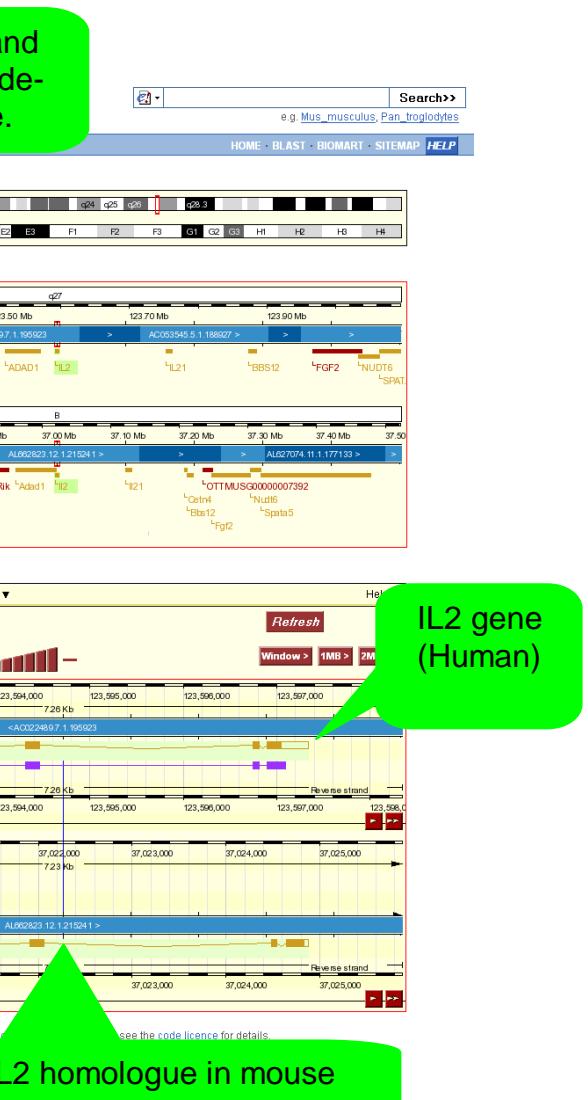
Upstream (<0.12 Gb) Downstream (>0.13 Gb)

Change Chromosome

Chromosome 4

Fields marked with \* are required

© 2008 WTSI / EBI. Ensembl is available to download for public use - please see the [code licence](#) for details.



## Select region/feature to Export

Choose at least one feature to export. Features must map to the current Ensembl Golden tile path. Please note we will not export more than 5Mb.

### Region

Chromosome name/fragment: 4

From (type): Base pair 123591080

To (type): Base pair 123598339

### Context

Bp upstream (to the left): [ ]

Bp downstream (to the right): [ ]

### Output

Output: FASTA sequence

**Continue >>**

**STEP 18:**  
Click on  
[Continue>>]

### Configure FASTA File output for FASTA sequence

You are exporting Chromosome 4 123,591,080 - 123,598,339.

This region is defined by: Chromosome 4, Bp 123591080, Bp 123598339

**Output format**  HTML

Text

Compressed text (.gz)

**Continue >>**

**STEP 19:**  
Click on  
[Continue>>]

ST · BIOMART · SITEMAP · HELP

**STEP 20:**  
Select and copy a part of the sequence

**STEP 21:**  
Click on 'BLAST'

**STEP 22:**  
Paste the copied sequence

**STEP 23:**  
Select 'Homo\_sapiens' and 'BLASTN' and click on [RUN>]

**e! Ensembl Human BlastView**

Ensembl release 44 - Apr 2007

Your Ensembl

- Login or Register
- About User Accounts

Sanger EBI

Bushbaby  
Otolemur garnettii

2X assembly and genebuild

new SETUP CONFIG RESULTS DISPLAY refresh Online Help

Retrieve result for ID: BLA\_APLFannCO Retrieve

Retrieving Results

Job pending results can be retrieved by clicking on the button above. Alternatively, this page can be bookmarked for later, or the ID noted and entered on the BLAST page.

Results are retained for 7 days. After this, they must be re-submitted.

1: unnamed (3300 letters) Vs. LATESTGP

Homo\_sapiens Job Queued

• -E: 10  
• -M: 1  
• -N: -3

**Summary**

- setup
  - Homo\_sapiens
  - Genomic sequence
  - BLASTN
  - Low sensitivity
- configure
  - E: 10
  - B: 100
  - filter: dust

**STEP 24:**  
Click on 'Retrieve'

new SETUP CONFIG RESULTS DISPLAY refresh Online Help

Retrieve result for ID: BLA\_APLFannCO Retrieve

Alignment Display Options:

Locations vs. Karyotype     Locations vs. Query  
 Summary Table

1: unnamed (3300 letters) Vs. LATESTGP

Homo\_sapiens 147 alignments, 52 hits [RawResult] view ►

**Summary**

- setup
  - Homo\_sapiens
  - Genomic sequence
  - BLASTN
  - Low sensitivity
- configure
  - E: 10
  - B: 100
  - filter: dust

**STEP 25:**  
Click on [VIEW]

**Ensembl Human BlastView**

Ensembl release 49 - Mar 2008

HOME · BLAST · BIOMART · SITEMAP · HELP

Your Ensembl

- Login or Register
- About User Accounts

Ensembl BLOG  
Keep up-to-date with news, views & announcements from the Ensembl team

Location of hits on the genome

Best hit (boxed)

Alignment of hits to query sequence

STEP 26: click on [C] in front of best (top) hit

Summary

- setup
  - Homo\_sapiens
  - Genomic sequence
  - BLASTN
  - Low sensitivity
- configure
  - E 10
  - B 100
  - filter\_dust
  - RepeatMasker
  - W 15
  - M 1
  - N -3
  - Q 3
  - R 3
- results
- display

Not yet initialised

new SETUP CONFIG RESULTS DISPLAY refresh Online Help

Displaying 4 sequence alignments vs Homo\_sapiens LATESTGP database

Showing top 1000 alignments of 5257, sorted by Raw Score

refresh

Chromosome (click arrow to hide)

Alignment Locations vs. Query (click arrow to hide)

+ hps

coverage

>4

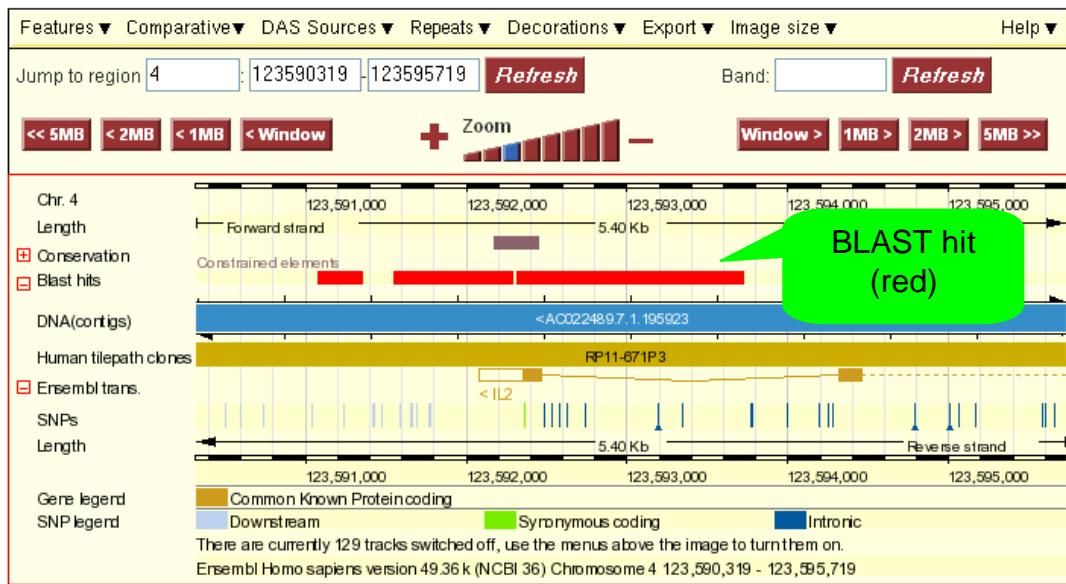
- hps

Alignment Summary (click arrow to hide)

Select rows to include in table, and type of sort (Use the 'ctrl' key to select multiples)

Query	Subject	Chromosome	Supercontig	Clone	Contig	Chromosome	Stats	Sort By				
_off_	_off_	_off_	_off_	_off_	_off_	_off_	_off_	_off_				
Name	Name	Name	Name	Name	Name	Name	Score	Eval				
Start	Start	Start	Start	Start	Start	Start	E-value	>Score				
[A]	[S]	[G]	1240	2640	+	Chr_4	123592319	123593719	1401	0.	100.00	1401
[A]	[S]	[G]	470	1206	+	Chr_4	123591543	123592289	737	0.	100.00	737
[A]	[S]	[G]	1	272	+	Chr_4	123591089	123591351	272	0.	100.00	272
[A]	[S]	[G]	193	233	+	Chr_8	124376539	124375679	30	1.4e-05	92.86	42
[A]	[S]	[G]	185	224	+	Chr_8	124376539	124375679	28	0.024	92.86	40
[A]	[S]	[G]	1121	1452	+	Chr_5	126452498	126452829	28	1.2	98.98	32
[A]	[S]	[G]	1395	1424	+	Chr_5	130242130	130242160	23	1.1e-05	93.55	31
[A]	[S]	[G]	2301	2323	+	Chr_2	240016107	240016129	23	0.0039	100.00	23
[A]	[S]	[G]	73	73	+	Chr_2	237422040	237422068	23	0.0017	96.30	27
[A]	[S]	[G]				Chr_2	237422040	237422068				
[A]	[S]	[G]	980	1002	+	Chr_X	11080600	11080622	23	0.73	100.00	23
[A]	[S]	[G]	1069	1083	+	Chr_X	11080600	11080621	23	1.1e-05	96.35	25
[A]	[S]	[G]	916	948	+	Chr_8	45498675	45498607	23	1.1e-05	18.34	34
[A]	[S]	[G]	2284	2324	+	Chr_8	88471051	88471096	22	0.10	98.96	46
[A]	[S]	[G]	532	553	+	Chr_5	71467068	71467087	22	3.1	100.00	22
[A]	[S]	[G]	986	1021	+	Chr_11	100577751	100577775	22	0.24	96.15	26
[A]	[S]	[G]	2237	2311	+	Chr_4	172121716	172121739	21	0.	100.00	21
[A]	[S]	[G]	1389	1419	+	Chr_4	172121716	172121739	21	0.	100.00	21
[A]	[S]	[G]	4107	4431	+	Chr_4	165382178	165382322	21	0.0012	98.00	25
[A]	[S]	[G]	1391	1415	+	Chr_4	31837404	31837428	21	5.4	96.00	25
[A]	[S]	[G]	1397	1417	+	Chr_8	94510145	94510165	21	1.2	100.00	21
[A]	[S]	[G]	176	196	+	Chr_8	19840422	19840442	21	0.00047	100.00	21
[A]	[S]	[G]	1996	2020	+	Chr_8	97952654	97952598	21	0.0012	98.00	25
[A]	[S]	[G]	1	25	+	Chr_2	30034541	30034541	21	0.0008	100.00	25
[A]	[S]	[G]	1388	1419	+	Chr_8	181716913	181716833	21	0.0017	100.00	21
[A]	[S]	[G]	2332	2356	+	Chr_5	142264541	142264564	21	0.030	96.00	25
[A]	[S]	[G]	1635	1659	+	Chr_5	67984528	67984552	21	0.0069	96.00	25
[A]	[S]	[G]	1514	1537	+	Chr_3	67251928	67251952	21	0.0081	96.00	25
[A]	[S]	[G]	1873	1937	+	Chr_3	87861931	87862005	21	0.0104	98.00	25
[A]	[S]	[G]	1260	1279	+	Chr_3	38043916	38043948	21	4.2	100.00	21
[A]	[S]	[G]	1333	1357	+	Chr_13	26132996	26134098	21	4.2	98.00	26

Back in the contigview page....



END of the  
Worked Example

## **EXERCISES and ANSWERS**

Note: the answers to these exercises correspond to current version (50) of Ensembl. If you use these exercises at a later date, please use the archive site for version 50.

### **III) BROWSING ENSEMBL**

These exercises address using the browser to determine a variety of gene-relevant information such as transcript number and size, protein domains, functional classes and sequence.

#### **1. Exploring features related to a gene**

*Exercise 1 begins with the TAC1 (tachykinin precursor 1) gene and moves into the browser from the main GeneView page.*

- (a) Open the home page of Ensembl ([www.ensembl.org](http://www.ensembl.org)). This is the current version. Search for the human TAC1 gene by typing ‘human TAC1 gene’ in the search window.
- (b) How many transcripts are predicted for this gene? What is the size of the longest predicted mRNA? How many exons does it have? How many amino acids does it code for?
- (c) Follow some of the links in the ‘Similarity Matches’ section of GeneView. What is a possible function of TAC1?
- (d) Which InterPro domains does the protein product contain?
- (e) Find the GO section of GeneView and follow some of the links to explore the ‘Gene ontology’ terms (describing gene and protein function) in Ensembl GOView.
- (f) In which chromosomal band and on which clone and contig in the genomic sequence assembly is the TAC1 gene located?
- (g) Go back to GeneView by clicking on ‘TAC1’ in the Overview panel and following the link for the gene. Is there a putative mouse orthologue? If so, where is it in the mouse genome?  
[http://www.ensembl.org/Homo\\_sapiens/glossaryview](http://www.ensembl.org/Homo_sapiens/glossaryview)
- (h) Go to the Ensembl main page. Look up ‘Ensembl genes’ in the glossary by following the ‘Using Ensembl’ link in the ‘Help and Documentation’ at the left. Follow the link: ‘Browse the provided data’ on the website.  
[http://www.ensembl.org/Homo\\_sapiens/glossaryview](http://www.ensembl.org/Homo_sapiens/glossaryview)

## 2. Exploring a region

*Exercise 2 begins with a search for a specific chromosomal region, rather than one gene.*

- (a) Click on the large ‘e!’ at the top left of the screen to start a new search. Go to the human homepage, and click on chromosome 12. From ‘MapView’, choose to display the region between markers D12S764 and D12S1871 (in the ‘Jump to ContigView’ section).
- (b) How many contigs are used to make this portion of the assembly? View the human tile path clones. Do they correspond to the assembly?
- (c) Click on a marker. What are other names for the marker? Is there an expected product size (what does this mean?)
- (d) In ContigView, zoom in three steps on the zoom triangle/ladder of ‘Detailed view’ (towards the ‘+’) and turn on the SNP track. Identify an intronic SNP and look at the corresponding SNPView page.

## Answers (Browsing Ensembl)

### 1. Exploring features related to a gene

- (a) Click on the identifier ENSG00000006128 from the search results. To ascertain it is indeed the TAC1 gene check that the HGNC symbol (the ‘official’ gene name given by the HUGO Gene Nomenclature Committee) is ‘TAC1’. You should now be in the GeneView page.
- (b) The TAC1 gene (ENSG00000006128) has 3 predicted transcripts, ENST00000319273, ENST00000346867 and ENST00000350485. Scroll down to the ‘Transcript’ sections for more information about these transcripts. The longest transcript is ENST00000319273. See this information in the heading of each ‘Transcript section’ of the GeneView page. The length of ENST00000319273 is 1060 bp. It has 7 exons and codes for 129 aa.
- (c) The TAC1 gene encodes for Protachykinin 1 precursor. Follow the links to MIM and EntrezGene or UniProt/Swiss-Prot in the ‘Similarity Matches’ section to learn more. Also the GO (Gene Ontology) and InterPro sections can give you clues about the biological and molecular function of the TAC1 protein. Tachikinins are neuropeptides. These hormones are thought to function as neurotransmitters which interact with nerve receptors and smooth muscle cells. They are known to induce behavioral responses and function as vasodilators and secretagogues.
- (d) Check the ‘InterPro’ section in GeneView. The domains include IPR013055 (Tachykinin/Neurokinin like), IPR002040 (Tachykinin/Neurokinin), IPR008215 (Tachykinin) and IPR008216 (Protachykinin).

(e) Clicking on a GO identifier gives you a GOView page showing the position of that term in the hierarchical ‘GO tree’ (note the number of Ensembl genes mapped to each term). Click [Help] to find out more about GOView.

(f) Go back to GeneView and click the ‘Graphical View’ link in the side menu to go to ContigView. In the ‘Overview’ panel you can see that TAC1 is located on band 7q21.3 (‘Chr.7 band’ track). In the ‘Detailed view’ panel you can see that it is located on contig AC004140.2.1.74918 (‘DNA(contigs)’ track) and clone RP5-841B21 (‘Human tilepath clones’ track).

(g) In GeneView, ENSMUSG00000061762 (Tac1) is named in the ‘Orthologue Prediction’ section. Click on it to go to its GeneView page to find that it is located on mouse chromosome 6 (band A1).

## 2. Exploring a region

(a) In the ‘Jump to ContigView’ section choose ‘From (type): Marker D12S764 To (type): Marker D12S1871’ and click [Go]. This leads you to ContigView.

(b) The region between the two markers will be displayed in ‘Detailed View’. This region includes sequence from 4 different contigs (one is quite small), displayed in light blue and dark blue in the ‘DNA(contigs)’ track. Clones are shown in gold and pink. Portions of the ‘Tile path clones’ sequences were used to form the assembly and correspond to ‘contigs’. The clones overlap each other whereas the contigs don’t.

(c) Click on a marker name (shown in pink at the top) and follow the link ‘Marker info’ to the MarkerView page. Other names in the UniSTS database will be at the top. Synonyms (or names in other databases) are listed further down. The expected product size is the calculated size of the fragment (in base pairs) if both primers are used against the genome.

(d) SNPs can be turned on using the ‘Features’ menu. Coding SNPs are shown in yellow (non-synonymous) and green (synonymous). Intronic SNPs are dark blue. Click on a SNP. Be careful to click exactly on the vertical bar representing the SNP, otherwise you will get the wrong pop-up menu. Follow the link ‘SNP properties’ to the SNPView page. Note the ‘SNP Context’ display in SNPView.

## IV) Data mining in Ensembl with BioMart Worked Example

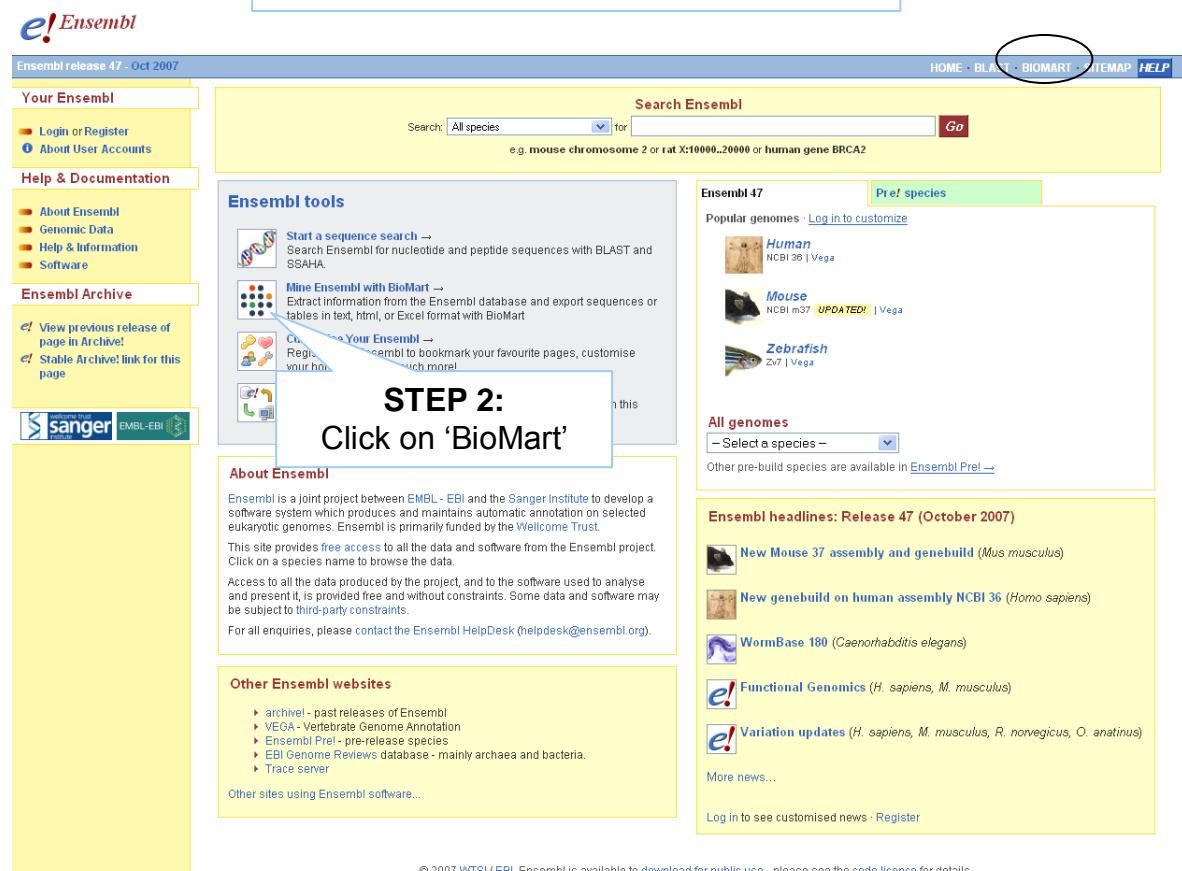
The human gene encoding Glucose-6-phosphate dehydrogenase (G6PD) is located on chromosome X in cytogenetic band q28.

Which other genes related to human diseases locate to the same band? What are their Ensembl Gene IDs and Entrez Gene IDs?

What are their cDNA sequences?

Follow the worked example below to answer these questions.

**STEP 1:**  
Go to the Ensembl main page  
[www.ensembl.org](http://www.ensembl.org)



The screenshot shows the Ensembl main page for release 47 from October 2007. At the top, there's a blue header bar with the Ensembl logo, the release number, and links for HOME, BLAST, BIOMART (which is highlighted with a red oval), TEMPLAR, and HELP. Below the header is a search bar labeled "Search Ensembl" with a dropdown menu set to "All species" and a "Go" button. To the right of the search bar is a "Pref. species" section for Human (NCBI 36 | Vega). The main content area has several sections: "Your Ensembl" (with links for Login/Register and About User Accounts), "Help & Documentation" (with links for About Ensembl, Genomic Data, Help & Information, and Software), "Ensembl Archive" (with links for previous releases and a Stable Archive link), and "Ensembl tools" (with links for Start a sequence search, Mine Ensembl with BioMart, and Create Your Ensembl). A large blue box labeled "STEP 2: Click on 'BioMart'" is overlaid on the "Mine Ensembl with BioMart" section. Other sections include "About Ensembl" (describing the project as a joint effort between EMBL-EBI and the Sanger Institute), "Other Ensembl websites" (links to ArchGen, VEGA, Ensembl PreL, EBI Genome Reviews, and Trace server), and "Other sites using Ensembl software..." (links to various bioinformatics databases). On the right side, there's a "All genomes" section with a dropdown for selecting a species, a "Ensembl headlines: Release 47 (October 2007)" section with links to news items about new mouse and human assemblies, and a "More news..." link at the bottom.

**New | Count | Results | XML | Perl | Help**

**Dataset**  
[None selected]

Ensembl 47

- CHOOSE DATASET -

**STEP 3:**  
Select the database:  
Ensembl genes (version 50)  
and the species of interest  
under 'Choose Dataset'.  
*(Homo sapiens)*

**New | Count | Results | XML | Perl | Help**

**Dataset**  
Homo sapiens genes (NCBI36)

**Filters**  
[None selected]

**Attributes**  
ensembl Gene ID  
ensembl Transcript ID

Please restrict your query using criteria below

REGION:  
 GENE:  
 GENE ONTOLOGY:  
 EXPRESSION:  
 MULTI SPECIES COMPARISONS:  
 PROTEIN:

**STEP 4:**  
Narrow the geneset by  
clicking '**Filters**' on the left.  
Click on the '+' in front of  
'REGION' to expand the  
choices.

**New** **Count** **Results**      **XML** **Perl** **Help**

<b>Dataset</b> Homo sapiens genes (NCBI36) <b>Filters</b> Chromosome: X Start : q28 End : q28 <b>Attributes</b> Ensembl Gene ID Ensembl Transcript ID  <b>Dataset</b> [None Selected]	<p>Please restrict your query using criteria below</p> <p><input type="checkbox"/> REGION:</p> <p><input checked="" type="checkbox"/> Chromosome      <input type="button" value="x"/> <input type="button" value="▼"/></p> <p><input type="checkbox"/> Base pair            Gene Start (bp)      <input type="text" value="1"/>            Gene End (bp)      <input type="text" value="10000000"/></p> <p><input checked="" type="checkbox"/> Band            Start      <input type="text" value="q28"/> <input type="button" value="▼"/>            End      <input type="text" value="q28"/> <input type="button" value="▼"/></p> <p><input type="checkbox"/> Marker            Start            End</p> <p><input type="checkbox"/> Encode type      <input type="button" value="manual_picks"/> <input type="button" value="▼"/></p> <p><input type="checkbox"/> Encode region</p>
--	---

**STEP 5:**  
Select 'Chromosome X'

**STEP 6:**  
Select 'Band Start q28'  
and 'End q28'

**STEP 7:**  
**Expand the 'GENE' panel and choose 'Limit to genes 'with MIM disease ID'.**  
 These associations have been determined using MIM (Online Mendelian Inheritance in Man).  
[www.ncbi.nlm.nih.gov/omim/](http://www.ncbi.nlm.nih.gov/omim/)

The filters have determined our gene set.  
 Click '**Count**' (at the top) to see how many genes have passed these filters.

**STEP 8:**  
 Click on 'Attributes' to select output options (i.e. what we would like to know about our geneset).

**STEP 9:**  
 Expand the 'GENE' panel.

**New | Count | Results**      **XML | Perl | Help**

**Dataset** 25 / 32584 Genes  
**Filters**  
 Chromosome: X  
 Start: q28  
 End: q28  
 with Disease association: Only

**Attributes**  
 Ensembl Gene ID  
 Ensembl Transcript ID  
 Gene name

**Dataset**  
 [None Selected]

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

**GENE:**

**Ensembl Attributes**

- Ensembl Gene ID
- Ensembl Transcript ID
- Ensembl Peptide ID
- Description
- Chromosome Name
- Gene Start (bp)
- Gene End (bp)
- Strand
- Band
- Transcript Start (bp)
- Transcript End (bp)
- Gene name

**EXTERNAL:**

- Gene DB
- External Transcript ID
- External Transcript DB
- Ensembl CDS length
- Ensembl cDNA length
- Ensembl Peptide length
- Transcript count
- % GC content
- Biotype
- Source
- Status (gene)
- Status (transcript)

**Note the summary of selected options.**

The order of attributes determines the order of columns in the result table.

**STEP 10:**  
 Select, along with the default options, '**Associated Gene name**' (this shows the gene symbol from HGNC).

**New | Count | Results**      **XML | Perl | Help**

**Dataset** 25 / 32584 Genes  
**Filters**  
 Chromosome: X  
 Start: q28  
 End: q28  
 with Disease association: Only

**Attributes**  
 Ensembl Gene ID  
 Ensembl Transcript ID  
 Gene name

**Dataset**  
 [None Selected]

**EXTERNAL:**

**GO Attributes**

- GO ID
- GO description
- GO evidence code

**External References (max 3)**

- CCDS ID
- EMBL (Genbank) ID
- EntrezGene ID
- Havana ID
- Havana transcripts (Shared)
- Havana transcripts (Identical)
- HGNC ID
- HGNC symbol
- PMID
- Imgt gene db
- Imgt ligm db
- Mim Gene Accession
- Mim Morbid accession
- Mirbase
- OTTP

**Microarray Attributes (max 2)**

**STEP 12:**  
 Click 'RESULTS' at the top to preview the output.

**STEP 11:**  
 Expand the 'EXTERNAL' panel to select External References (IDs outside of Ensembl). Select '**EntrezGene ID**' and '**MIM Gene Accession**' and '**MIM Morbid Accession**'. These are MIM phenotypes and diseases, respectively.  
[www.ncbi.nlm.nih.gov/omim/](http://www.ncbi.nlm.nih.gov/omim/)

**New | Count | Results**      **XML | Perl | Help**

**Dataset 25 / 32584 Genes**

**Filters**  
 Chromosome: X  
 Start: q28  
 End: q28  
 with Disease association: Only

**Attributes**  
 Ensembl Gene ID  
 Ensembl Transcript ID  
 Gene name  
 EntrezGene ID  
 Mim Gene Accession  
 Mim Morbid accession

**Dataset**  
 [None Selected]

Export all results to    Unique results only **Go**

Email notification to

View  rows as   Unique results only

Ensembl Gene ID	Ensembl Transcript ID	Gene name	EntrezGene ID	Mim Gene Accession
ENSG00000102081	ENST00000370475	FMR1	23	309550
ENSG00000102081	ENST00000370475	FMR1	23	309550
ENSG00000155966	ENST00000370460	AFF2	2334	309548
ENSG0000010404	ENST00000340855	IDS	3423	309900
ENSG00000013619	ENST00000370401	CXorf6	10046	300120
ENSG00000013619	ENST00000370401	CXorf6	7280	300120
ENSG00000013619	ENST00000370401	CXorf6	730	300120
ENSG00000013619	ENST00000262858	CXorf6	7280	300120
ENSG00000013619	ENST00000262858	CXorf6	730818	300120
ENSG00000013619	ENST00000262858	CXorf6	728030	300120
ENSG00000013619	ENST00000262858	CXorf6	730818	300120
ENSG000000171100	ENST00000306167	MTM1	4534	300415
ENSG00000147383	ENST00000370274	NSDHL	50814	300275
ENSG00000130821	ENST00000253122	SLC6A8	6535	300036
ENSG00000130821	ENST00000253122	SLC6A8	731026	300036
ENSG00000185825	ENST00000345046	BCAP31	10134	300398
ENSG00000185825	ENST00000370133	BCAP31	10134	300398
ENSG00000101986	ENST00000218104	ABCD1	215	300371
ENSG00000101986	ENST00000218104	ABCD1	642762	300371
ENSG00000101986	ENST00000218104	ABCD1	215	300371
ENSG00000101986	ENST00000218104	ABCD1	642762	300371
ENSG00000198910	ENST00000370058	L1CAM	3897	308840
ENSG00000198910	ENST00000370058	L1CAM	3897	308840
ENSG00000198910	ENST00000370058	L1CAM	3897	308840
ENSG00000198910	ENST00000370058	L1CAM	3897	308840
ENSG00000198910	ENST00000370060	L1CAM	3897	308840

To save a file of the complete table, click 'Go'. Or, email the results to any address.

**STEP 13:**  
 Go back and change Filters or Attributes if desired.  
 Or, View 'ALL' as HTML...

## Result Table 1

Ensembl Gene ID	Ensembl Transcript ID	Gene name	EntrezGene ID	Mim Gene Accession	Mim Morbid accession
ENSG00000102081	ENST00000370475	FMR1	2332	309550	300623
ENSG00000102081	ENST00000370475	FMR1	2332	309550	300624
ENSG00000155966	ENST00000370460	AFF2	2334	309548	309548
ENSG00000010404	ENST00000340855	IDS	3423	309900	309900
ENSG00000013619	ENST00000370401	CXorf6	10046	300120	300633
ENSG00000013619	ENST00000370401	CXorf6	728030	300120	300633
ENSG00000013619	ENST00000370401	CXorf6	730818	300120	300633
ENSG00000013619	ENST00000262858	CXorf6	728030	300120	300633
ENSG00000013619	ENST00000262858	CXorf6	730818	300120	300633
ENSG000000171100	ENST00000306167	MTM1	4534	300415	310400
ENSG00000147383	ENST00000370274	NSDHL	50814	300275	308050
ENSG00000130821	ENST00000253122	SLC6A8	6535	300036	300352
ENSG00000130821	ENST00000253122	SLC6A8	731026	300036	300352
ENSG00000185825	ENST00000345046	BCAP31	10134	300398	300475
ENSG00000185825	ENST00000370133	BCAP31	10134	300398	300475
ENSG00000101986	ENST00000218104	ABCD1	215	300371	300100
ENSG00000101986	ENST00000218104	ABCD1	642762	300371	300100
ENSG00000101986	ENST00000218104	ABCD1	215	300371	300475
ENSG00000101986	ENST00000218104	ABCD1	642762	300371	300475
ENSG00000198910	ENST00000370058	L1CAM	3897	308840	142623
ENSG00000198910	ENST00000370058	L1CAM	3897	308840	303350
ENSG00000198910	ENST00000370058	L1CAM	3897	308840	304100
ENSG00000198910	ENST00000370058	L1CAM	3897	308840	307000
ENSG00000198910	ENST00000370060	L1CAM	3897	308840	142623

**STEP 14:**  
To view sequences, go  
back to 'Attributes'

Chrom Start : q1  
End : q2  
with Disease association: Only

**Attributes**

- Ensembl Gene ID
- Ensembl Transcript ID
- External Gene ID
- EntrezGene ID
- Mim Gene Accession

**Dataset**  
[None Selected]

**XML | Perl | Help**

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

Features    Homologs  
 Structures    Sequences  
 SNPs

**GENE:**

Ensembl Gene ID  
 Ensembl Transcript ID  
 Ensembl Peptide ID  
 Description  
 Chromosome Name  
 Gene Start (bp)  
 Gene End (bp)  
 Strand  
 Band  
 Transcript Start (bp)  
 Transcript End (bp)  
 External Gene ID

**EXTERNAL:**

**GO Attributes**  
 GO ID  
 GO description    GO evidence code

**External References (max 3)**  
 CCDS ID  
 Codelink ID  
 EMBL ID  
 EntrezGene ID  
 Havana ID  
 HGNC Symbol  
 Illumina v1  
 Illumina v2  
 IPI ID  
 Imgt gene db  
 Imgt ligr db  
 Mim Gene Accession  
 Mim Morbid accession    Protein ID  
 RefSeq DNA ID  
 RefSeq Predicted DNA ID  
 RefSeq Peptide ID  
 Rfam ID  
 Unigene ID  
 Shares cds with enst  
 Shares cds with ott  
 UniProt/SPTREMBL ID  
 UniProt/Swiss-Prot ID  
 UniProt/Swiss-Prot Accession  
 Unified UniProt ID  
 Unified UniProt Accession

**STEP 15:**  
Select 'Sequences'

**New | Count | Results**

**Dataset** 25 / 32584 Genes

**Filters**  
Chromosome: X  
Start : q28  
End : q28  
with Disease association: Only

**Attributes**

- Ensembl Gene ID
- Chromosome
- Biotype
- cDNA sequences

**Dataset**  
[None Selected]

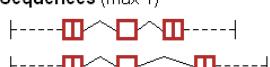
**XML | Perl | Help**

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

Features    Homologs  
 Structures    Sequences  
 SNPs

**SEQUENCES:**

**Sequences (max 1)**



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

**Upstream flank**  
 Upstream flank \_\_\_\_\_

**Downstream flank** \_\_\_\_\_

**STEP 16:**  
Expand the 'SEQUENCES' panel and  
select 'cDNA sequences'.

UTR  
'UTR  
Exon sequences  
 cDNA sequences  
 Coding sequence  
 Peptide

**STEP 17:**  
Expand the 'Header Information' to choose  
'Ensembl Gene ID',  
'Chromosome', and 'Biotype'

**Dataset 25 / 32584 Genes**

**Filters**  
Chromosome: X  
Start: q28  
End: q28  
with Disease association: Only

**Attributes**  
Ensembl Gene ID  
Chromosome  
Biotype  
cDNA sequences

**Header Information**

**Gene Attributes**

- Ensembl Gene ID
- Description
- External Gene ID
- External Gene DB

**Transcript Attributes**

- Ensembl Transcript ID
- Ensembl Peptide ID
- RefSeq ID
- Transcript Start (bp)
- Transcript End (bp)
- Ensembl CDNA length
- Ensembl CDS length
- 5' UTR Start (bp)
- 5' UTR End (bp)
- 3' UTR Start (bp)
- 3' UTR End (bp)

**STEP 18:  
Click 'Results'**

**Dataset**  
[None Selected]

Ensembl CDNA Start (bp)  
 Ensembl CDNA End (bp)  
 Ensembl CDS Start (bp)  
 Ensembl CDS End (bp)  
 Coding Start (bp)  
 Coding End (bp)

Exon Start (bp)  
 Exon End (bp)  
 Exon Strand  
 Exon Block in Transcript

**Dataset 25 / 32584 Genes**

**Filters**  
Chromosome: X  
Start: q28  
End: q28  
with Disease association: Only

**Attributes**  
Ensembl Gene ID  
Chromosome  
Biotype  
cDNA sequences

**Dataset**  
[None Selected]

Export all results to    Unique results only

Email notification to

View  rows as   Unique results only

**STEP 19:  
View all rows as  
FASTA'**

```
>ENSG00000196924|X|protein_coding
GCGCGTCGCGCAGCGGACGCCGACAGAATCCCTGGCGCTGGCGGGCGCGGGCG
CGAAGGCATCCGGCGCCACCCCGCGGTATCGG16CTGGCTCTCAGGAACAGCA
GCGCAACCTCTGCTCCCTGCCCCCTGGCTCCCGCGGGCTTAAGCTGGCTGGCTGG
AGGGCUCGCCCCCTCGCCGAGCTGGCTGGCGCCACCGCCCGCC
AGTAGCTCCCACTCTCGGGCGGGCCAGAGCCAGCA
GACACGGGGACGCCAGATGCCGGCCACCGAGAAC
AAGAAATCCAGAGAACACTTACCGCGCTGGTGC
AAGGCATCGCCAACCTTGAGACGGACCTGAGCAGAC
GAGGTGCTCAGCCAGAAGAGATGCCACCGCAAGCAC
ATGCCACCTTGAGAACCTGTCGGTGGCGCTCGAGTC
GTGTCATCGACAGCAAGGCCATCGTGGACGGGAACCTGAAGCTGATCTGGGCTCATC
TGGACCTGATCTGCACTACTCCATCTCATGCCATGTGGGACGAGGGAGGGATGAG
GAGGCCAAGAAGCAGACCCCCAAGCAGAGGGCTCTGGGCTGGATCCAGAACAGCTGCC
CAGCTGCCCATACCAAACCTCAGGCCAGAGCAGGGGGCCCTGGGCGCCCTG
GTGGACAGCTGTGCCCCGGCTGTGTCCTGACTGGGACTCTTGGGACGCCAGCAGGCC
GTTACCAATGCGCGAGAGGCCATGCAGCAGGCCGATGACTGGCTGGGATCCCCCAGGTG
```

## RESULTS

### Header: chromosome, Ensembl Gene ID, Biotype

```
>ENSG00000196924 |X|protein_coding
GCGCGTCGCGCGCAGCGGACGCCGACAGAATCCCGAGGCCTGGCGCGGGCGCGGGCG
CGAAGGCATCCGGCGCCACCCCGCGGTATCGGTACCGGCTCGACTCAGGAACAGCA
GCGCAACCTCTGCTCCCTGCCTCGCCTCCCAGCGCCTAGGGCTCGACTTAATTAA
AGGGCGTCCCCTCGCCGAGGCTGCAGCACGCCCGGCTTCGCGCCTCAAATG
AGTAGCTCCACTCTCGGGCGGGCCAGAGCGCAGCAGGCCGGCTCCGGCGGGCGTC
GACACGCGGGACGCCGAGATGCCGCCACCGAGAAGGACCTGGCGGAGGACGCCGTGG
AAGAAGATCCAGCAGAACACTTCACCGCCTGGTCAACGAGCACCTGAAGTGCCTGAGC
AAGCGCATGCCAACCTGCAGACGGACCTGAGCGACGGCTGCGGCTTATCGCGCTGTTG
GAGGTGCTCAGCCAGAAGAAGATGCACCGCAAGCACAACCAGCGGCCACTTCCGCAA
ATGCAGCTTGAGAACGTGTCGGTGGCGCTGAGTTCTGGACCGCGAGAAGGACCTGGCG
GTGTCCATCGACAGCAAGGCATCGTGGACGGAACCTGAAGCTGATCC
TGGACCCCTGATCTGCACTACTCCATCTCCATGCCATGTGGGACGAGG1 cDNA 1
GAGGCCAAGAAGCAGACCCCCAAGCAGAGGCTCTGGGCTGGATCCAGAACAGCTGCCG
CAGCTGCCATACCAACTTCAGCCGGACTGGCAGAGCGGCCGGCCCTGGCGCCCTG
GTGGACAGCTGTGCCCCGGGCTGTGCTGACTGGGACTCTGGGACGCCAGAACCC
GTTACCAATGCGCGAGAGGCCATGCAGCAGGCCGGATGACTGGCTGGGATCCCCCAGGTG
...
GACAAGGGGGAGTACACACTGGTGGTCAAATGGGGGACGAGCACATCCCAGGCAGCCCC
TACCGCGTTGTGGTGCCTGAGTCTGGGCCCCGTGCCAGCCGGCAGCCCCAACGCTGCC
CCGCTACCCAAGCAGCCCCGCCCTCTTCCCTCAACCCCGGCCAGGCCCTGGCGCC
CCGCTGTCACTGCAGCCGCCCTGCCCTGTGCCGTGCTGCCACCTGCCTCCCCAGC
CAGCCGCTGACCTCTCGGCTTCACTGGGAGAGGGAGCCA
>ENSG0000013619 |X|protein_coding
GGCGCGGAGCCGGCGGTGGGAATGGAGCGAGCAGATTGAGGCCACTGCAGCGCCGC
CAGCATGAACCTGGCCGCAGCTGAAGCGGCCGGCGGGCGGGCGGCC
CGCTAGCCAGGGGGTGAATCTGCAACAGGGCTGGGTTCTGGCGGCC
CTGCCCCGCCGGCGGCCGCCCTCGGACACTGCCCGGCCGCC
CAGCACGCCCTGTCTAGGCTTGGAAACGCCCTGGAGAGTCAGAACATAATTG
AGGTCAAACAATGGATGACTGGAAAAGCTGGCTTGTAAATCAAGAGCATGCTCCCCATT
CGCCATGGGGAAATCGTCAGGAGCCCAGAAAGCTCCAGGAATGGGAAAGAACGCCCTC
GTGGATGGAGGAAGAAGATTATCTTCTACAAGAGCAGCCCAGGAAGAACGATCA
GGGAACGTAAAGAGGAGACAAGAACGACACTTCCAGTTCCAGACATGGCTGATGG
GGGCTACCCATAAAATTAAAGAGGCCCTGCCTTGAAAGATGTCACCCCTGCAATGGGCC
AGGTGCTCATCCTAGTACTGCTTGTGCAGAACTGCAAGGCTCCATTGACAATAATCC
TAGCCCTGCGGCTATGGGAGTGGCTGCCAGTCATTACTGCTGGAGAATAACCCTATGAA
TGGCAACATCATGGGCTCACCATTTGTTAGTACCAACAGACTACAGAACAGACTGGGACTGAAAGG
GCCCACTGTTCTTACTATGAGAAAATCAACAGCGTGCCTGAGGAGCTTCA
AGAGCTGCTAGAGGAGCTACCAAAATTCAAGACCTTCTCAAATGAGCTAGATCTGA
GAAGATACTGGGACGAAGCCAGAAGAGCCACTGGTTTAGATCATCCCCAGGCAACCCT
AAGCACAACTCCCAAGCCTCGGTTAGATGTCACACTGGAGAGCCTGGCTCCAGCAA
GGAGTTGCTTAGTTGCAAGCTTACTGGCATGTCACCTCAGATCCCCTCC
CACAGGGATCAGCTATTGATTGCTTCCACCGTAAGCAGATAGTGTACCGAGTTCTC
AATGGCACAGTCAAGAGCCAGGTCCAGGCCATGCTCCCTGCGCTCTGCCCTTACC
AGTGCCTCAGTGGCATCACGCCACCAAGCTGAAGGCCAGGCAGCAAGCAGGGTC
TGCTACAAAGCAGCAAGGGCCCACCCCCAGTTGGCTGGCTGCCCTCCAGGACTCTC
TCCACCTTACCGCCCAGTGCACCATCACCAACCCACCACCGCTGCCACTGCCACCAC
```

**cDNA 1**

**cDNA 2**

## V) BIOMART

*These exercises have been designed to familiarise you with different questions you can answer with this tool, and the types of data you can retrieve with BioMart.*

1. Retrieve all SNPs for ‘novel’ human G-protein coupled receptor genes (GPCRs – use the InterPro domain ID: IPR000276) on chromosome 2.

*Note: As this is the first exercise we walk you this time through BioMart step-by-step (but of course you can also try to do this exercise without our help!)*

Start a new BioMart session by clicking ‘New’, or go back to the Ensembl homepage and click on ‘Mine Ensembl with Biomart’ under ‘Ensembl tools’.

Choose the **database** and the **dataset** for your query as follows:

- Select ‘Ensembl 50’
- Select ‘Homo sapiens genes (NCBI36)’.

Click on ‘**Filters**’ at the left. Filter this dataset to select your genes of interest as follows:

- Expand the ‘REGION’ section at the right by clicking on the ‘+’. Select ‘Chromosome 2’. Click [count] at the top of the panel and note the number of Ensembl genes on *Homo sapiens* chromosome 2.
- In the ‘GENE’ section, select ‘Status (gene)’ ‘NOVEL’.
- In the ‘PROTEIN’ section, select the second ‘Limit to genes with these family or domain IDs’ option. Select ‘Interpro ID(s)’ and enter ‘IPR000276’ in the box. Click [count] again and note that the number of genes is now 3.

Click on ‘**Attributes**’ (at the left). Select the output for your gene list as follows:

- Select the ‘SNPs’ Attribute Page.
- In the ‘GENE’ section ‘Ensembl Gene ID’ and ‘Ensembl Transcript ID’ are selected by default – also select ‘Ensembl Protein’.
- In the ‘GENE ASSOCIATED SNPs’ section ‘Reference ID’ is selected. Also select ‘Allele’, ‘Protein location (aa)’, ‘Location in Gene’ and ‘Protein Allele’.

*Note: Clicking on count now will not show an altered number. Attribute selections should not affect the count (i.e. the number of genes that have passed the filters).*

Click on ‘**Results**’ (at the top) to obtain the first 10 rows of your table. To obtain the entire table select ‘View all rows as HTML’ or export a file by clicking ‘Go’. Check the box ‘Unique results only’, otherwise you can end up with redundant rows!!

Note that the output for this query gives you one row for each SNP, and if there are alternative transcripts then SNP data is given for each. This means that a particular SNP may appear more than once.

Find the coding SNPs, and note that you have information about the effect of the SNP, and its location within the protein.

2. Click 'New' to start a new query. Retrieve the gene structure (i.e. start and end coordinates of exons) of the mouse gene ENSMUSG00000042351.
3. Retrieve all human genes that are located between p11.2 and q22 (these are bands on chromosome 1) and that have a MIM disease ID.
4. The file at [http://www.ebi.ac.uk/~xose/Affy\\_exercise.txt](http://www.ebi.ac.uk/~xose/Affy_exercise.txt) contains a list of probeset IDs from a microarray experiment using the Affymetrix array HG-U133 Plus 2.0 (human). Enter the probe list at this url into BioMart to retrieve the 500 bp upstream of the transcripts matching these probeset IDs.
5. Retrieve the sequences 5kb upstream of all human 'known' genes between D1S2806 and D1S464.
6. Retrieve sequence (including reference ID in the header) of all human SNPs on chromosome 6, band p25.3, that have an ID from Watson's genome. Export both the SNP sequence and 200 bases flanking (adjacent to) it.
7. Retrieve the mouse homologues of *Homo sapiens* genes CASP1, CASP2, CASP3, and CASP4 (these are the HGNC symbols for the genes).
8. Design your own query!

## Answers (BioMart)

1. You should find **three** novel genes on chromosome 2 with this InterPro domain. The result set has three transcripts (one for each gene) and a total of 375 rows of output (to see this, change the option from TSV to XLS under 'Export all results', select the option 'Unique results only' and click 'Go', then open in Excel so you don't have to count the rows manually). The three genes/transcripts have two, eight and one coding SNPs ('Location in Gene' is 'coding'), respectively. Most of these are non-synonymous and thus affect the amino acid sequence of the encoded protein. One allele is a stop codon (\*)- can you find it?

**2. Click New. Select: Database and dataset:** 'Ensembl 50' and '*Mus musculus* genes (NCBIM37)'.

**Filters: GENE:** 'ID list limit Ensembl Gene ID(s)': enter the mouse gene ID.

**Attributes: Structures:** select in the **EXON** panel: 'Ensembl Exon ID', 'Exon Start' and 'Exon End'.

Click '**Results**'.

You should find **8 exons**. Follow the link from the Ensembl Gene ID in your output back to the **GeneView** page to confirm the BioMart data with the gene structure displayed on this page.

**3. Database and dataset:** ‘Ensembl 50’ and ‘*Homo sapiens* genes (NCBI36)’.

**Filters:** **REGION:** ‘Chromosome 1’, ‘Band Start p11.2’, ‘Band End q22’, **GENE:** ‘Limit to genes ... with MIM disease ID(s) Only’.

**Attributes:** **Features:** **EXTERNAL:** select ‘MIM Morbid Description’ along with the default options (‘Ensembl Gene ID’ and ‘Transcript ID’).

**Results** should show **18 Ensembl genes** (some of which with multiple transcripts).

**4. Database and dataset:** ‘Ensembl 50’ and ‘*Homo sapiens* genes (NCBI36)’.

**Filters:** **GENE:** ‘ID list limit’: Affy hg u133 plus 2 ID(s) and enter the list of probeset IDs.

**Attributes:** **Sequences:** **SEQUENCES:** select ‘Flank (Transcript)’, ‘Upstream flank 500’. In ‘Header Information’, apart from the already default selected options, select ‘Ensembl Transcript ID’.

You should find upstream sequences for the transcripts of **26 genes** (Hint: click ‘count’ to see the number of genes!)

**5. Database and dataset:** ‘Ensembl 50’ and ‘*Homo sapiens* genes (NCBI36)’.

**Filters:** **REGION:** ‘Marker’ Start D1S2806 End D1S464, **GENE:** ‘Status: KNOWN’.

**Attributes:** **Sequences:** select, apart from the already default selected options, ‘Flank (Gene)’ and ‘Upstream flank 5000’.

You should find sequences for **41 genes**.

When you choose the option ‘Flank (Gene)’ you will see only one upstream sequence per gene in the output. In the case where a gene has multiple transcripts, the upstream sequence of the transcript that extends the furthest at the 5’ end is shown. If you want to export the upstream sequences for each transcript you should choose the option ‘Flank (Transcript)’.

‘Known’ genes are Ensembl gene predictions that could be matched to same-species external database entries (e.g. UniProt/SwissProt) with a high similarity score (i.e. with BLAST or a similar sequence identity-matching program)

**6. Database:** ‘SNP’ and **dataset:** ‘*Homo sapiens* SNPs

(dbSNP127;HGVbase 15; TSC 1; ENSEMBL, Affy 100K and 500K arrays)

**Filters: REGION:** ‘Chromosome 6’, ‘Band p25.3’, **GENERAL SNP FILTERS:** SNP source: ‘with Variation synonym ensemblwatson ID(s)’ Only.

**Attributes Sequences: SEQUENCES:** ‘SNP sequences’, ‘Upstream flank 200’, ‘Downstream flank 200’. **SNP: SNP attributes:** default is ‘Reference ID’.

You should find **2473 SNPs**.

## 7. Database: ‘Ensembl 50’ Dataset: *Homo sapiens* genes (NCBI36)

**Filters: GENE:** ‘ID list limit HGNC Symbol(s)’. Enter the human HGNC (HUGO) symbols in the box: CASP1, CASP2, CASP3, and CASP4.

**Attributes:** Under **Homologs**, select in the **MOUSE ORTHOLOGS** panel ‘Mouse Ensembl Gene ID’. Also, scroll back up to the **GENE** panel. Along with the default option ‘Ensembl gene ID’, select ‘Associated Gene Name’ and deselect ‘Ensembl Transcript ID’ (these attributes are for the starting dataset... i.e. Human).

**Results** displays the mouse orthologues of the 4 human CASP genes.

## VI) EVALUATING GENES AND TRANSCRIPTS

### (The ‘GeneBuild’)

#### Main Exercise: Examine the evidence for the FOXH1 gene.

*These exercises focus on the FOXH1 gene to demonstrate how the underlying protein and mRNA used to build an Ensembl gene can be seen. Is it a well-determined gene...?*

#### 1. Display the human FOXH1 gene in GeneView.

Enter FOXH1 into the text search box at the top of any human Ensembl page. Take the link to the **GeneView** page for the Ensembl Gene: ENSG00000160973.

**Q1:** What are the external database sources for the gene name and for the description?

#### 2. Examine the supporting evidence for the FOXH1 gene in ExonView.

Scroll down the **GeneView** page and have a look at the predicted exon structures. Note the 5' and 3'UTRs (untranslated regions). Compare the two transcripts using this view.

Click on ‘Exon information’ in the left-hand menu to go to **ExonView** for one of the transcripts. The bottom section of **ExonView** shows the supporting evidence that was used during the Ensembl transcript building process.

**Q2:** Which databases did the entries come from?

Click on a green box of a supporting evidence entry to see the alignments of the Ensembl predicted transcript against the supporting evidence.

*Optional: Click on the ‘Gene information’ link in the left-hand menu to return to the **GeneView** ‘Gene Report’ for FOXH1.*

#### 3. Examine other evidence for the FOXH1 gene in ContigView.

Click on the link “Graphical View” in the left-hand menu. This takes you to **ContigView** displaying only the region encompassing the gene. Zoom out by clicking on the ‘-’ button next to the Zoom triangle.

Look at other protein/mRNA tracks, e.g. Human proteins, Unigene (for all species), Human cDNAs, and EMBL mRNAs (again across species). These are proteins/mRNAs that align to the genome in this area. Note that the track labels are links to the help page. Clicking on a block drawn in the new tracks brings up a pop-up menu with a link to the database entry. Try some links.

Condense the Unigene and Protein track as well as the Overview section by clicking on the '-' boxes.

Under 'Decorations' select 'Show empty tracks'. This will help you remember which tracks you have selected. Examine the evidence for the FOXH1 gene in the new evidence tracks. Are there proteins and mRNAs that align to the Ensembl prediction?

#### **4. Compare transcript predictions made by other methods.**

In the 'Features' menu, turn off most of the evidence and make sure Ensembl genes, Vega Havana and Genscans are turned on.

Look at the Genscan track. Zoom out further.

**Q3:** How does the Genscan prediction differ from the Ensembl prediction? Note that this track shows *ab initio* Genscan predictions, not relying on supporting evidence.

**Q4:** What is a Vega Havana transcript?

Use the 'DAS sources' menu to turn on tracks showing transcript predictions from other groups (e.g. NCBI Gnomon). Compare and contrast!

#### **5. Look at the 'Similarity Matches' section.**

Go to **TransView** by clicking on the Ensembl transcript FOXH1 and taking the direct 'Transcr.' link from the pop-up menu.

'Known' Ensembl transcripts like this one (shown in red in **ContigView**) have been successfully mapped to external database entries such as UniProt or NCBI EntrezGene entries *for the same species*. (Note that this mapping is done *after* the genes have been built). Novel transcripts may match to mRNA and protein information in other databases, for an alternate species. Thus a novel transcript is novel for that species.

Matches to IDs in other databases are shown in the 'Similarity Matches' section of **TransView** (repeated in the 'Transcript' section of **GeneView**). Have a look at the types of databases linked out to.

**Q5:** What do the Target and Query % ids indicate? Check the online Help pages.

### **Answers (Evaluating Genes and Transcripts).**

**A1:** Name: HUGO Gene Nomenclature Committee (HGNC). Description: Uniprot/Swiss-Prot.

**A2:** Click on the ID number to go to the original database entry. You may have to scroll down to find the correct entry. The boxes represent the exons,

the darker green they are, the better the supporting evidence is. Click on a green box for the alignments with the transcript.

**A3:** The Genscan transcript prediction shows exons not present in the Ensembl transcripts for FOXH1 and doesn't predict UTRs. Genscan is an *ab initio* predictor, a program run on the sequence alone, without using protein and mRNA evidence. It has the tendency to overpredict exons.

**A4:** Havana is a subgroup of VEGA, the Vertebrate Genome Annotation consortium. Havana transcripts are manually curated (determined on a case-by-case basis) and are merged with Ensembl predictions if they match up exactly (leading to golden transcripts). For more, go to these links:

<http://www.sanger.ac.uk/HGP/havana/>

<http://vega.sanger.ac.uk/index.html>

**A5:** In **TransView**, under Similarity Matches, Target %ID indicates the percentage of the Ensembl prediction matching the external sequence database and Query %ID is the percentage of the external database sequence matching the Ensembl prediction! Can you find this in the help pages?

## VII) COMPARATIVE GENOMICS

### 1. Investigating a protein family.

This exercise focuses on protein families and orthologies, using AIM1, a cancer-related gene, as an example.

Find the **GeneView** page for human AIM1, a gene absent in a type of melanoma. Beware... there is a VEGA gene as well as an Ensembl gene! Let's work with the Ensembl gene for now.

**Q1:** What is the MIM ID for the AIM1 transcript? Hint... find this in 'Similarity Matches' for one of the transcripts. Click on the MIM ID...

Examine the protein family by taking the link to the associated Protein Family (the **FamilyView** page).

**Q2:** How many human Ensembl genes produce peptides in this family?

**Q3:** Are they all 'known' genes?

**Q4:** Are there peptides in the same family for mouse (*Mus musculus*), rat (*Rattus norvegicus*) and zebrafish (*Danio rerio*)?

**Q5:** Families are calculated using Ensembl peptides along with UniProt peptides. Can you find on this page where the UniProt peptides in this family are listed?

### 2. Genomic alignments and conserved regions

Find the **ContigView** page for human SNX5.

**Q6:** In which band is the human SNX5 gene located?

The 'constrained elements' track shows the most highly conserved regions calculated from the 10-way multi-species alignment.

**Q7:** Do any of these regions fall outside of exons?

Expand the 'Conservation' track by clicking on the '+' in front of it. This shows a score for each basepair based on how conserved it is across the alignments. Clusters of peaks correspond to the 'Constrained elements blocks' shown above.

Collapse the track if desired. Choose the pairwise alignment with mouse from the Comparative menu. Selecting the '+' in front of the track will expand the track so that all alignments on the same mouse chromosome and strand will be clustered together.

**Q8:** On which chromosome(s) and strand(s) are the mouse alignments for this region?

Turn on the frog (*X. tropicalis*) pairwise alignment with human, using the Comparative roll-down menu.

**Q9:** Are there fewer or more conserved regions between human and frog when compared with human and mouse? Would you expect an SNX5 frog gene in this region?

Click on a frog alignment. ‘Jump to Xenopus tropicalis’ to view ContigView for frog. Can you see the snx5 gene there?

### 3. Examine the syntenic regions

Syntenic regions in Ensembl are calculated from the pairwise alignments (whole genome alignments between two species) and correspond to long regions (100 kb or more) that have high sequence similarity. Gene order is expected to be conserved in these regions. Two pages are available to view these regions: **SyntenyView** and **CytoView**.

To go to SyntenyView, where two species may be compared, start by finding the **ContigView** page for human Myosin VI.

Click on ‘View syntenic regions with... *Mus musculus*’, a link at the left hand side of the page.

**Q10:** How many mouse chromosomes show syntenic regions with human chromosome 6?

**Q11:** Is there a mouse homologue to human Myosin VI? Can you see this without leaving the **SyntenyView** page?

Gene order is conserved in human and mouse for this syntenic region. Comparing homologous genes across both species can show where gene loss may have occurred.

Click on ‘**MultiContigView**’ next to ‘Myo6’ in mouse to reach a comparison of the human and mouse chromosomes in this region.

**Q12:** What does the vertical blue line connecting the lowest panels show?

**Q13:** The Myosin VI genes are golden in human and mouse. What does the colour signify?

Syntenic regions may also be viewed in the **CytoView** page. Take the ‘Graphical overview’ link at the left to go to **CytoView**. Turn on syntenic regions with mouse and opossum using the ‘Comparative’ roll-down menu.

Click on the red and green lines in the ‘Mouse’ and ‘Opossum’ tracks to see which base pairs and chromosome strand are syntenic to human in this region.

The ‘View alongside...’ link at the left will take you back to **MultiContigView**. Instead, take the ‘View alignment with...’ link and choose ‘*Mus musculus*’. This will bring you to the ‘**AlignSliceView**’ page, where again, human and mouse chromosomes can be compared side-by-side. However, as **MultiContigView** shows the chromosomes as they are, **AlignSliceView** shows the alignment, and the genomic assembly may be altered (gaps introduced) to fit the alignment in this view. Read the ‘Help’ page to understand the view.

Select the ‘12 amniota vertebrates’ in the ‘Comparative’ roll-down menu.

**Q14:** Is there a gene in this aligned region across the twelve species?

## Answers (Comparative Genomics)

### Main Exercise (1)

**A1:** Find this in the GeneView page. Make sure the ID at the top is ENSG00000112297. The MIM ID (corresponding to the entry on the ‘Online Mendelian Inheritance in Man database) is 601797, found in the ‘Similarity Matches’ section of the **GeneView** page. (See the report for the second, AIM1, transcript).

**A2:** Family IDs change with every release... in 50, it is fam50v00000001399. 3 human genes..

**A3:** They are all “Known” protein-encoding Ensembl genes (all genes have a name; none are ‘novel’).

**A4:** Yes. Scroll down the ‘FamilyView’ page to find the peptide lists for other species.

**A5:** In the ‘Other peptides...’ section, just under the gene names and karyotype, the UniProt identifiers clustered into this family are listed. 4 mouse, 4 rat, and 4 zfish proteins.

**A6:** Band p11.23, chromosome 20 (*hint: go to ContigView*)

**A7:** Yes. This track is turned on by default- if not, switch it on in the ‘Comparative’ roll-down menu under ‘Constrained elements’. Looking at the visual display in Detailed view, most of the conserved blocks line up with exons (if you are unsure, zoom in.) Clicking on a conserved element shows its precise location across the genomes aligned. There are some conserved

elements near the 3' end of the gene that fall outside of exons (especially in the long intron of the SNX5 transcript). Could these be important regulatory regions? Hint... turn on ncRNA (non-coding RNAs) in the Features menu.

**A8:** All are on chromosome 2, forward strand. In the expanded format, i.e. a (-) is found next to the track, all alignments appear grouped together on one line, indicating they are on the same chromosome and strand. Click on a block to find this is chromosome 2, forward strand.

**A9:** Fewer. You would expect a gene, because the conserved regions with frog line up to the human SNX5 exons, supporting a hypothesis in which those exons have been conserved across time in frog as well as in human. Thus, they are probably part of an expressed protein.

**A10:** 7 chromosomes in mouse are shown in **SyntenyView**. These are the chromosomes that share syntenic blocks with human chromosome 6.

**A11:** The homologues are listed at the right of the page. Mouse Myo6 is a homologue of human Myo6. It is on chromosome 9 at position 80.01 Mb in mouse.)

**A12:** The vertical blue line in the lowest, 'Detailed View' panel connects homologous genes across, in this case, human and mouse. Remember homologues are determined using the longest translation of a gene only.

**A13:** A golden gene has at least one transcript that is identical to (and merged with) a transcript determined by manual curation (the Havana team).

**A14:** Yes, a gene is present in these aligned regions for all species. Some exons align well, indicating they could be homologues.

## VIII) VARIATIONS

These exercises focus on SNPs (Single Nucleotide Polymorphisms) in Ensembl. Most of these are downloaded from dbSNP and placed along the sequence assembly using flanking sequence to the SNP to match it.

### 1. Display all SNPs for a region.

*This exercises focuses on SNPs in the **ContigView** page.*

Go to the **ContigView** page displaying Human chromosome 7: bp 116124438-116129149. Turn on SNPs.

**Q1:** Are there non-synonymous coding SNPs in the gene? (Remember, non-synonymous SNPs change the amino acid sequence).

You can see information about a SNP by clicking on it. Click on a non-synonymous SNP and follow the link 'SNP properties' to its **SNPView** page.

In the **SNPView** page, click on ‘SNPs in gene context’ at the right of the second panel. You will be taken to the **GeneSNPView** page, where you can see all the SNPs for the MET gene. Under ‘SNP type’ roll-down menu in the **GeneSNPView** page, select only ‘non-synonymous’.

## 2. Display SNPs linked to a disease within the transcript sequence.

*Starting with disease-causing SNPs, how can we view them within the sequence? Hint... focus on the TransView page.*

In the article “Screening of the delta-F508 mutation and analysis of two single nucleotide polymorphisms of the CFTR gene in a sample of the general population of Valparaiso, Chile” by L.A. Vera et al.” (Rev Med Chil 2005, 133:767-775.) the SNPs leading to peptide alleles M470V and T854T are studied.

**Q2:** Find these two SNPs in the **TransView** and **GeneSNPView** pages. Note, the positions refer to numbers in the amino acid sequence.

## 3. Display all SNPs for one gene.

Retrieve all the validated SNPs associated with the human CFTR gene as follows. Display only validated SNPs in the ‘**GeneSNPView**’ page used in question 2.

Now use BioMart to ask how many of these validated SNPs are coding.

**Q3:** How many validated SNPs are in the coding region of CFTR?

## Answers (Variations)

**A1:** Click on the ‘human’ icon on the Ensembl homepage. Enter the chromosome number and region in base pairs under the karyotype. This should take you to **ContigView**. You may have to turn on the ‘SNPs’ track from the ‘Features’ drop-down menu on top of ‘Detailed View’.

SNPs should now be shown as vertical lines of different colours along the chromosome. The SNP legend is shown at the bottom of the panel. This region contains one exon from the ‘MET’ gene. There are a few non-synonymous coding SNPs (shown in yellow) in this region.

**A2:** Perform a text search for ‘CFTR’ in human. This should lead you to ENSG00000001626. Go to the **GeneView** page and then the **TransView** page by clicking on ‘transcript information’ on the left hand navigation column. At the bottom of the transcript sequence are panels for customisation of the display. By selecting the options ‘Exons, Codons, Translations and SNPs’ and ‘Number residues: yes’ you can display the SNPs in the transcript sequence. Alleles and alternative codons are shown by pointing your mouse over the nucleotide and amino acid residues, respectively. M470V is shown as a red ‘V’ at position 470 in the peptide sequence (the second line of numbering in the sequence).

Continue on to the **GeneSNPView** page for this gene by clicking ‘Gene variation info.’ in the side menu.

The two SNPs are displayed in the ‘SNPs and variations’ figure and the ‘Variations and consequences’ table. In the figure, M470V is shown as ‘V/M’ in yellow (as it is a non-synonymous coding SNPs). T854T is shown as ‘T’ in green, as it is a synonymous coding SNP. The diagram only shows bp positions, so it could be easier to find the two allelic changes in the table. Note that you can use the ‘SNP class’ and ‘SNP type’ drop-down menus in the figure to configure (simplify) both figure and table by choosing only to see coding SNPs.

**A3:** Validated SNPs have been checked after submission to determine if they are SNPs or mutations. The validation status comes from NCBI, and explanation to the validation terms can be found in the ‘Help’ to SNP-based pages.

To view only validated SNPs, configure the ‘Validation’ roll-down menu in the **GeneSNPView** page by deselecting ‘no information’ and closing the menu. Now only validated SNPs are shown. How many are coding? Deselect, under the ‘SNP type’ roll-down menu in this page, upstream, downstream, UTR, intronic and intergenic SNPs. Close the menu. The display and table will be simplified. Count **14** coding SNPs that have been validated.

Go to BioMart. Select the SNP database and *Homo sapiens* SNPs. In **Filters**, expand the ‘GENERAL SNP FILTERS’ section. Select ‘SNPs that have been validated’ (Only). In the ‘GENE ASSOCIATED SNP FILTERS’ section, enter the Ensembl Gene ID ENSG00000001626. Select ‘Ensembl Gene Location’ as ‘Coding Only’.

‘Count’ gives you **14** validated, coding SNPs.

## **IX) TYING IT TOGETHER – CASE STUDIES (using Ensembl comprehensively)**

*Consider the following case studies. The answers to these questions tie together pages and tools on the Ensembl site.*

- 1) I work on non-coding RNA genes such as snoRNA (small nucleolar RNA) and miRNA. Are there any on mouse chromosome 1 (of any type), and where are they? How many are there (use BioMart). Could I view them graphically?
- 2) Many regulatory factors appear in the 5'UTR, and I am interested in searching for them using motif-scanning tools. How would I obtain the 5'UTR or upstream region in the IRX4 gene in human? What can I learn about the UTR and upstream regions using **ContigView**?
- 3) (*Using DAS*): We have determined two SNPs in the lab not found in Ensembl. We have not yet submitted to dbSNP, but can we just display these SNPs on our computer in Ensembl? Here is some information about the SNPs:

Names: 1 and 2

Positions: both on chromosome 3 (mouse)

base pairs: 75305500 and 75450001

Type and subtype are both 'a' (internal lab naming convention)

Phase is 0, Score is 100.

Hint: start from ContigView for mouse chromosome 3: 75300000 - 75500000

- 4) I would like to export the gene structure of zfish 'pp2ca2' in GFF format (I will use a program that requires data in this format).  
(Find out more about the GFF format here)  
[http://www.sanger.ac.uk/Software/formats/GFF/GFF\\_Spec.shtml](http://www.sanger.ac.uk/Software/formats/GFF/GFF_Spec.shtml)

### **Case Studies: Answers**

- 1) Using BioMart to find all the RNA genes will take 2 queries as the filters only allow a subset of ncRNA types to be selected. There are **192** total RNA genes on chromosome 1. You need only select the filters, then click 'count' to see how many genes there are.

From the result table, follow the link to one of these genes into the Ensembl browser (for example, ENSMUSG00000070103).

Alternatively, go to '**ContigView**' for a region on chromosome 1 and select, in the 'features' menu, 'ncRNA' to visualise the Ensembl ncRNAs in the chromosome. Also, select 'rFAM' to select RNAs predicted by this program. You can also select tRNA, miRNA, and other options.

- 2) **Step 1: Exporting UTR sequence.** Export the UTR information by clicking ‘Export gene data’ at the left of the **GeneView**, **TransView**, or other page and then selecting ‘5’ or 3’ UTR’ in the next window to export the sequence.

Alternatively, use **BioMart** to export the 5’ and 3’ UTRs.

**Step 2: Visual display and comparison of promoter elements.**

*What can you learn about the upstream regions of this gene in ContigView? Use the features menu and conserved tracks.*

In **ContigView** (the ‘graphical view’ link) select ‘CTCF’ for CTCF binding sites (click the blue ‘Help’ if you are unsure of what this is), CpG islands, and motifs in the CisRED/miRANDA database’ for sequences associated with regulatory regions. For genes that have no UTR annotated, it can be helpful to select “Eponine” (a program that predicts transcription start sites) under the ‘Features’ menu. Also, in the Comparative menu, select alignments such as ‘10 amniota vertebrates’. Constrained elements show the highest sequence similarity for the alignment.

- 3) **HINT: Upload as a DAS source from ContigView.**

Go to **ContigView** encompassing this region in mouse, then select ‘**DAS Sources**’ from the roll-down menu above the detailed-view panel and ‘**Manage Sources**’. Next, click on ‘**Upload your data**’ in the new window (at the left). Read the instructions by clicking on the link at the top of the page (especially the ‘formatted correctly’ link.) Enter in your email and a NON-SECURE password, then paste your data according to the format described in the link. (Note: columns must be separated by tabs, NOT SPACES! This will require making the data columns in another program such as Notepad).

Once the data is entered, click ‘Next’ and select **ContigView** (note you can select more than one page). Click ‘Next’ again, name the track if you’d like, and click ‘Finish’. Once you reach the DAS sources list (end page) ‘Close window’ to refresh **ContigView**... a new track should be displayed, along with SNPs 1 and 2!

Here is the correct format for the SNP information given in the question:

[http://www.ebi.ac.uk/~gspudich/workshop\\_presentations/snp\\_example.txt](http://www.ebi.ac.uk/~gspudich/workshop_presentations/snp_example.txt)

Note: you can deselect your track under the ‘DAS sources’ menu of ContigView.

- 4) To export in GFF format: Select ‘export gene data’ at the left of the **GeneView** page and select GFF under ‘Output format’. Or, export the chromosome and base pair start and stop of the exons using the ‘structures’ attribute page in **BioMart**. Select GFF format under ‘Display rows as GFF’ and click ‘Go’ to export the file.

# e! Ensembl Did you know?

## Focus

**1** Genome sequence assemblies determined by sequencing institutes are incorporated into Ensembl. View the assembly used for a species on the main page.



Homo sapiens  
NCBI 36

**2** The Ensembl gene set is based on mRNA and protein evidence and is aligned to the genomic sequence assembly in the 'genebuild'. Use 'ExonView' to view this information.

**3** Ensembl focuses on vertebrates however other model organisms are included such as yeast, fly, and soil worm. View newly available genomic sequences in the Pre! site.

## Browsing

**4** A new version of Ensembl is released every two months. View older versions in the Archive! site. New gene sets are determined when a new sequence assembly is available.



**5** ContigView and GeneView can be customised using the 'features' menu. Select options to visualise transcripts, variations, regulatory regions, etc. Use the DAS menu to show data in current databases external to Ensembl.



View a short video about ContigView here!  
[http://www.ensembl.org/commonWorkshops\\_Online](http://www.ensembl.org/commonWorkshops_Online)

## Tools

**6** Use BioMart to get tables of genes and annotation in Microsoft Excel, HTML, or txt format. You can also export sequences with this Web-based tool.  
[www.biomart.org](http://www.biomart.org)

**7** BLAST and SSAHA are alignment programs- use them to match a sequence on any genome.



**8** View variations for a gene or sequence in GeneSNPView or ContigView.

We calculate phylogenetic gene trees, **9** alignments, and homologies.



Send any questions or comments to

[helpdesk@ensembl.org](mailto:helpdesk@ensembl.org)