

Access to genes and genomes with **Ensembl**



Introduction and Worked Example

January 2007

CONTENTS

INTRODUCTION.....	2
WORKED EXAMPLE	7
Exercises.....	26
Answers.....	27

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

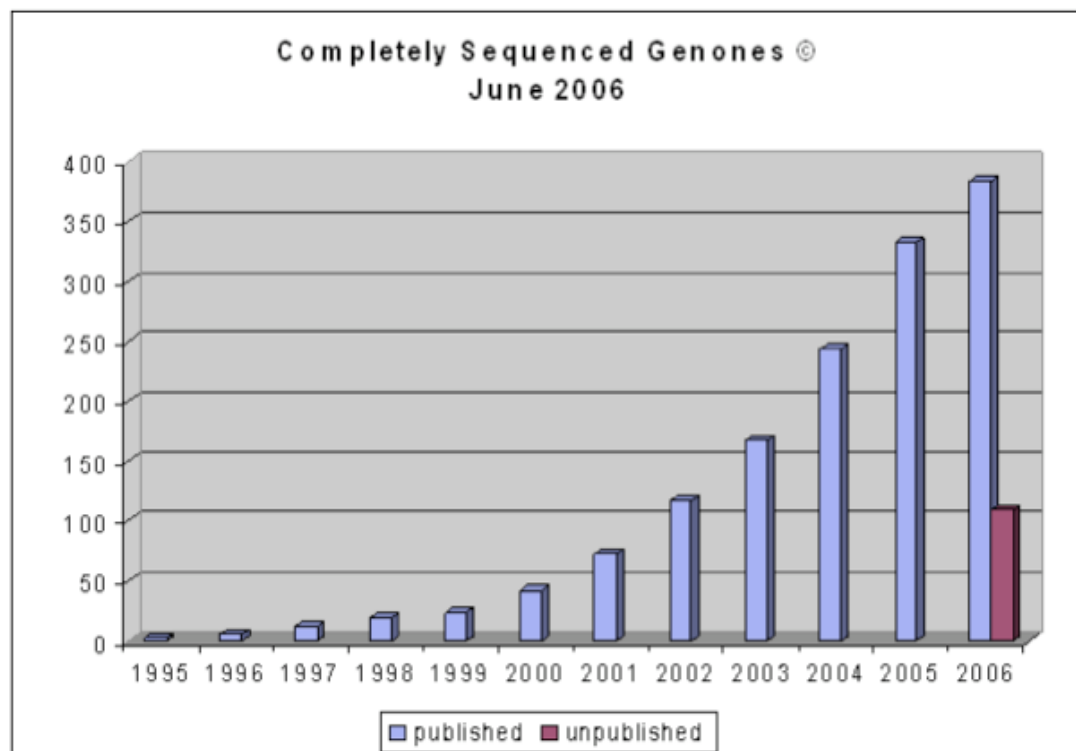


Figure 1. Completely sequenced genomes as of June 2006 (figure taken from <http://www.genomesonline.org>).

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. Finally, Ensembl BLAST offers the possibility to perform sequence searches against genomes and Ensembl gene and peptide sets.

Further reading

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.* *

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)

* 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		ASSEMBLY		GENEBUILD	
Mammals					
Human	<i>Homo sapiens</i>	NCBI 36	oct 2005	Ensembl	jul 2006
Chimpanzee	<i>Pan troglodytes</i>	PanTro 2.1	mar 2006	Ensembl	mar 2005
Rhesus macaque	<i>Macaca mulatta</i>	MMUL 1	feb 2006	Ensembl	aug 2006
Bushbaby*	<i>Otolemur garnettii</i>	BUSHBABY1			
Mouse	<i>Mus musculus</i>	NCBI m36	dec 2005	Ensembl	apr 2006
Rat	<i>Rattus norvegicus</i>	RGSC 3.4	dec 2004	Ensembl	feb 2006
Rabbit	<i>Oryctolagus cuniculus</i>	RABBIT	may 2005	Ensembl	aug 2006
Dog	<i>Canis familiaris</i>	CanFam 1.0	jul 2004	Ensembl	nov 2004
Cat*	<i>Felis catus</i>	CAT			
Cow	<i>Bos taurus</i>	Btau 2.0	mar 2005	Ensembl	dec 2005
Pig**	<i>Sus scrofa</i>				
Shrew*	<i>Sorex araneus</i>	sorAra1			
Hedgehog*	<i>Erinaceus europaeus</i>	eriEur1			
Microbat*	<i>Myotis lucifugus</i>	MICROBAT1			
Armadillo	<i>Dasypus novemcinctus</i>	ARMA	may 2005	Ensembl	aug 2006
Elephant	<i>Loxodonta africana</i>	BROAD E1	may 2005	Ensembl	aug 2006
Lesser hedgehog tenrec	<i>Echinops telfairi</i>	TENREC	may 2005	Ensembl	aug 2006
Opossum	<i>Monodelphis domestica</i>	MonDom 4.0	jan 2006	Ensembl	feb 2006
Platypus*	<i>Ornithorhynchus anatinus</i>	OANA 5			
Other species					
Chicken	<i>Gallus gallus</i>	WASHUC 1	mar 2004	Ensembl	dec 2005
<i>X. tropicalis</i>	<i>Xenopus tropicalis</i>	JGI 4.1	aug 2005	Ensembl	nov 2005
Zebrafish	<i>Danio rerio</i>	Zv 6	mar 2006	Ensembl	aug 2006
Fugu	<i>Takifugu rubripes</i>	FUGU 4.0	jun 2005	IMCB/JGI	may 2005
Tetraodon	<i>Tetraodon nigroviridis</i>	TETRAODON 7	apr 2003	Genoscope	sep 2004
Stickleback	<i>Gasterosteus aculeatus</i>	BROAD S1	feb 2006	Ensembl	aug 2006
Medaka	<i>Oryzias latipes</i>	HdrR 1	oct 2005	Ensembl	may 2006
<i>C. intestinalis</i>	<i>Ciona intestinalis</i>	JG 12	mar 2005	Ensembl	feb 2006
<i>C. savignyi</i>	<i>Ciona savignyi</i>	CSAV 2.0	oct 2005	Ensembl	apr 2006
Fruitfly	<i>Drosophila melanogaster</i>	BDGP 4	jul 2005	FlyBase	mar 2006
Anopheles	<i>Anopheles gambiae</i>	AgamP 3	feb 2006	VectorBase	oct 2005
Aedes	<i>Aedes aegypti</i>	AaegL 1	aug 2005	VectorBase	jun 2006
<i>C. elegans</i>	<i>Caenorhabditis elegans</i>	WS 150	nov 2005	WormBase	nov 2005
<i>S. cerevisiae</i>	<i>Saccharomyces cerevisiae</i>	SGD 1	nov 2005	SGD	nov 2005

Figure 2 – Species in Ensembl, including name and date of their genome assembly and source and date of the genebuild. * = currently only available on the Pre! website, ** = only clone information available.

WORKED EXAMPLE – A walk through the main pages of the Ensembl browser, using the EPO (Erythropoietin precursor) gene as an example.

STEP 1:
Load Ensembl
www.ensembl.org

The screenshot shows the Ensembl genome browser homepage. Several green callout boxes highlight key features:

- Navigation column:** Points to the left sidebar containing links like 'Your Ensembl', 'Healthchecks', 'Help & Documentation', 'Select a species', and 'Ensembl Archive'.
- Search:** Points to the 'Search Ensembl' box at the top right, which includes a search input field and a 'Go' button.
- Help:** Points to the 'Help' link in the top navigation bar.
- Help pages and Documentation:** Points to the 'Helpdesk' link in the left sidebar.
- What's new:** Points to the 'Ensembl headlines: Release 42 (December 2006)' section in the main content area.

STEP 2: Click on "Homo sapiens"

The 'Popular genomes' section on the right lists various species, with 'Homo sapiens' highlighted. A yellow arrow points from the 'Homo sapiens' link to the 'STEP 2' callout box.

At the bottom of the page, there is a footer with copyright information: © 2006 WTSI / EBI. Ensembl is available to download for public use - please see the code licence for details.



STEP 3:
Type in 'EPO Gene'.
Click 'Go'.

Karyotype

e!Ensembl Human

Ensembl release 42

Search Ensembl *Homo sapiens*

Search: Go

e.g. chromosome X or 14:10000..200000 or BRCA2

Karyotype

Click on a chromosome for a closer view

Jump directly to sequence position

Chromosome: or region

From (bp):

To (bp): Go

Ensembl headlines: Release 42 (December 2006)

- New - User accounts** (all species)
- New species - Duck-billed Platypus** (*Ornithorhynchus anatinus*)
- New Dog assembly and genebuild** (*Canis familiaris*)
- New Chicken assembly and genebuild** (*Gallus gallus*)
- New Human Ensembl-Vega** (*Homo sapiens*)

[More news...](#)

Go to your account to customise this news panel

Statistics

Assembly:	NCBI 36, Oct 2005
Genebuild:	Ensembl, Aug 2006
Database version:	42.36d
Known genes:	21,774
Novel genes:	1,036
Pseudogenes:	1,069
RNA genes:	3,976
GENSCAN gene predictions:	69,195
Gene exons:	270,661
Gene transcripts:	44,676
Base Pairs[†]:	3,253,037,807
Golden Path Length^{††}:	3,093,120,360
Most common InterPro domains:	Top 40 Top 500

[†] Total number of base pairs = sum of lengths of DNA table

^{††} Reference assembly (Golden path) length = sum of non-redundant top level seq regions

You are using the web team's integration server. [More...](#)

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

**Source and version of
assembly and genebuild**

STEP 4:
Click on 'ENSG00000130427'

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

The Gene View Page

Gene Model

STEP 5:
Click on 'Transcript Information'

Orthologues in other species

Matches in other databases

GO
(Gene Ontology)
terms

Your Ensembl

- Show account - Log out
- Save bookmark
- Save configuration as...

ENST00000252723

- Gene information
- Gene splice site image
- Gene regulation info.
- Genomic sequence
- Gene variation info.
- ID history
- Transcript information**
- Exon information
- Protein information
- Export transmembrane

Chromosome 7
100,156,359 - 100,156,359

- View of Chir
- Graphical v
- Graphical o
- Export info
- Export info in region
- Export sequ
- Export EMB
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

Healthchecks

- Health checks
- Old Health checks

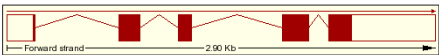
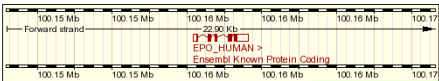
Ensembl Archive

- [View previous release of page in Archive!](#)
[Stable Archive! link for this page](#)



- Logins**
Bookmarks
Settings
Groups
- User accounts**
New in Ensemble!

Ensembl Transcript Report

Transcript	<p>EPO_HUMAN (UniProt/KB/Swiss-Prot) To view all Exon links linked to the name click here.</p> <p>This transcript is a member of the Human CCDS set: CCDS5705</p>		
Ensembl Transcript ID	ENST00000252723		
Transcript information	<p>Exons: 5 Transcript length: 1,328 bps Translation length: 193 residues</p> <p>This transcript is a product of gene: ENSG00000130427</p>		
Genomic Location	<p>This transcript can be found on Chromosome 7 at location 100,156,359-100,159,257.</p> <p>The start of this transcript is located in Contig AC009488.5.1.98876.</p>		
Description	Erythropoietin precursor (Epoetin). Source: UniProt/KB/SwissProt P01488		
Prediction Method	<p>Genes were annotated by the Ensembl automatic analysis pipeline using either a GeneWise/Exonerate model from a database protein or a set of aligned cDNAs followed by an ORF prediction. GeneWise/Exonerate models are further combined with available aligned cDNAs to annotate UTRs (For more information see V.Curwen et al., Genome Res. 2004 14:942-50)</p>		
Similarity Matches	<p>This Ensembl entry corresponds to the following database identifiers:</p> <p>CCDS: CCDS5705.1</p> <p>UniProt: EPO_HUMAN [Target Xid: 100; Query Xid: 100] [align] NP_000790.2 [Target Xid: 100; Query Xid: 100] [align] NM_000799.2 [align] Q2MZL6_HUMAN [Target Xid: 100; Query Xid: 100] [align] 2056 A_14_P113914 [Target Xid: 3; Query Xid: 100] A_23_P145664 [Target Xid: 4; Query Xid: 100] A_23_P145669 [Target Xid: 4; Query Xid: 100] AC009488 [align] AF053356 [align] AF202306 [align] AF202307 [align] AF202308 [align] AF202309 [align] AF202310 [align] AF202311 [align] AF202312 [align] AF202314 [align] AF202314 [align] BC093628 [align] BC111937 [align] M11318 [align] S65458 [align] X02157 [align] X02158 [align] IP00307.226.3 [Target Xid: 100; Query Xid: 100] MIM gene: 133170 PDB: 1BUV 1CN4 1EER Protein ID: AAA52400.1 [align] AAC78791.1 [align] AAD13964.1 [align] AAF17572.1 [align] AAF23132.1 [align] AAF23133.1 [align] AAF23134.1 [align] AAH93628.1 [align] AAI1938.1 [align] AAP2357.1 [align] CAA26094.1 [align] CAA26095.1 [align] UniGene: Hs.2303 [Target Xid: 99; Query Xid: 98] Affymx Microarray Focus: 207257_at Affymx Microarray HG C110: 1023_at Affymx Microarray HUGeneFL: X02158_ma1_at Affymx Microarray U133: 207257_at 207257_at 207257_at 217254_s_at 217254_s_at 217254_s_at Affymx Microarray U95: 1023_at GE Healthcare/Amersham Codewalk WGA: GE79554 [Target Xid: 2; Query Xid: 100] Illumina V1: GI_4503589-S [Target Xid: 3; Query Xid: 98]</p>		
GO	<p>The following GO terms have been mapped to this entry via UniProt and/or RefSeq:</p> <p>GO:001666 [response to hypoxia] IEA GO:0005128 [erythropoietin receptor binding] IEA GO:0005173 [hormone activity] IEA GO:0005576 [extracellular region] IEA GO:0005615 [extracellular space] TAS GO:0006950 [response to stress] TAS GO:0007165 [signal transduction] NAS GO:0007267 [cell-cell signaling] NR GO:0007275 [development] NR GO:0008015 [circulation] NAS GO:0030218 [from] [erythrocyte differentiation] IEA GO:0043249 [erythrocyte maturation] IEA</p>		
InterPro	<p>IPR003013 Erythropoietin - View other genes with this domain IPR001323 Erythropoietin/thrombopoietin - View other genes with this domain</p>		
Protein Family	<p>ENSF00000006225 : ERYTHROPOIETIN PRECURSOR This cluster contains 1 Ensembl gene member(s) in this species.</p>		
Transcript structure			
Transcript neighbourhood			
Transcript sequence	<pre> CCCGAGCCGACGCGGGCCACGCGCCGCTCTGCTCCGACACGCGCCCCCTGGACAG CGCCCTCTCTCCAGGCCCTGGGCTGGCTCCAGCGGACGCTCCCGGATAGAG GCCGCCGCTGTCTCCAGCGCGCGCCGCGAGCTGCTCAGGACGACGCGCCGACGCGGGA GAGGGGGTGCAGCAATGTCTCTGGGTGGGCTCTCTGTGCTGGTGGTGGTGGTGGTGGT TCTGGGCTCTCAATCTGGGCGCCGACGAGCTCTATCTGGGAGGAGGAGCTCTGGA GAGTACTCTTGGAGGCGCGGAGAGATATCAGCAGGCGCTGTCTCAACACTG CACTTGAATGAGATATCACTCTCCGACACGAGGAGTAACTTCTATCTCTGGAGAG GATGAGATCTGGGAGAGGAGCTGAGAGTCTGGGAGGCGCTCTGTCTGGGAGCTGG TGCTCTGGGAGCGAGCGCTGTGGTCAACTCTTCCAGCGTGGGAGCGCTCTGGAGCT GCATGATAGAGCTCAATGAGCTCTGGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAG AGCCGACAGGAGGAGGAGCTCTCCCTCTCAGATAGGCGCTCAGCTGCTCAGCTCAGAC CACTGCTGACACTTCCGAGCAATCTTCCGAGTCTACTCAAACTCTCTCCGGGAGAGCT GAGGCTCTACAGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG GGATATCCACAGCTCTCTCAGCAACTTGTCTGGCAGGAGCTCTCCCGCGCACTCT GAACCGCTCCGAGGAGGCTCTCAGCTCAGCGCGAGCTGTCCATGGAGCACTCTAGTCC GCAGAGGAGATCTCGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG TCAGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAG TCAGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAG AGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAG TGGATAGTAACTGGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAG GGTGGGAGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGGAGG GCTCTGAGTCTCAGGGGGCCAGAGTCTTGTCTTCTCAAGCTCAGGAGGAGGAGGAGG AAGCCAC</pre>		

Spliced transcript
sequence

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

The diagram illustrates the analysis of a DNA sequence for various genetic variations. The sequence is shown in a multi-line format with different background colors for exons and introns. Annotations include:

- Exons - alternating text colour**: Exons are highlighted with alternating text colors (yellow and blue).
- Introns - alternating background colour**: Introns are highlighted with alternating background colors (yellow and blue).
- Synonymous SNP**: A single nucleotide change that does not alter the amino acid sequence.
- Non-synonymous SNP**: A single nucleotide change that alters the amino acid sequence.
- Other variation in coding sequence**: A variation in the coding sequence that is not a SNP.
- Affected residue**: A residue that is affected by a variation. (Mouse over shows alternative codons)*
- Ambiguity code**: A code used to represent a variation in the sequence.
- Other variation in UTR**: A variation in the untranslated region (UTR) of the sequence. (Mouse over shows alleles)*
- UTR SNP**: A single nucleotide change in the UTR. (Mouse over shows alleles)*
- UTR (dark background)**: The untranslated region (UTR) of the sequence, highlighted with a dark background.

Result of STEP 7:

e!Ensembl Human ExonView

Ensembl release 42 - Dec 2006 (ecs33307 - homo_sapiens_core_42_36d)

Search e!Human:

e.g. ENSE00001428812, ENSE00000837374

HOME - BLAST - BIOMART - SITMAP - HELP

Your Ensembl

- Show account - Log out
- Save bookmark
- Save configuration as...

ENSE00000252723

- Gene information
- Gene splice site image
- Gene regulation info
- Genomic tracks
- Gene variants
- ID history
- Transcript
- Exon information
- Protein info
- Export transcript

Chromosome 7
100,156,359 - 100,156,552

- View of Chromosome 7
- Graphical view
- Graphical overview
- Export information about region
- Export sequence as FASTA
- Export EMBL file
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

Healthchecks

- Health checks
- Old Health checks

Ensembl Archive

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

Ensembl Exon Report

Transcript: **EPO_HUMAN** (UniProt/Swiss-Prot) To view all Ensembl genes linked to this transcript is a member of the Human CCDS set: [CCDS5705](#)

Ensembl Transcript ID: **ENSE00000252723**

Transcript information: Exons: 5 Transcript length: 1,328 bps Translation length: 193 residues
This transcript is a product of gene: [ENSG00000130427](#)
This transcript can be found on Chromosome 7 at location [100,156,359 - 100,156,552](#)
This transcript is located in [Contig AC009488.5:1.98875](#)
Protein precursor (Epoetin). Source: UniProt/SWISSPROT P01988

Flanking sequence at either end of transcript: 50
Intron base pairs to show at splice sites: 25

Show full intronic sequence ☐
Show exons only ☐

Exon Information

No.	Exon / Intron	Chr	Strand	Start	End	Start Phase	End Phase	Length	Sequence
1	5' upstream sequence	7	1	100,156,359	100,156,552	-	-	-ccttgggcccacccggcgctgcgtgcgtgcgcgcacgcgcgtgtcct
2	Intron 1-2	7	1	100,156,553	100,157,116	-	1	-	gagcgggctggggcgtc.....ctctcagcctgggtatctgttctag
3	Intron 2-3	7	1	100,157,521	100,157,607	0	0	87	ACGGGCTGTGCTGAACAACCTGCAGCTTGAATGAGAATATCACTGTCCAGACACCAAGTT
4	Intron 3-4	7	1	100,157,608	100,157,721	0	0	114	gttccctttttttttttttt.....gactccagcagtcacactccctgtag
5	Intron 4-5	7	1	100,158,537	100,159,257	0	0	721	gagtgaggagcgacacttctgct.....cgacactcgttttctccttggcag
6	3' downstream sequence	7	1	100,159,258	100,159,552	-	-	-	aatatgactcttggttttctgttttctgggaacctccaaatcccctggc.....

Supporting Evidence

The supporting evidence below consists of the sequence matches on which the exon predictions were based and are sorted by alignment score.

Score: 100 99 97 95 90 85 80 75 70 65 60 50 40 30 20 10 0 NO EVIDENCE

X02157.1
X02157.1 Human mRNA for fetal erythropoietin
P01988
P01988.1 EPO_HUMAN Erythropoietin precursor (Epoetin)

You are using the web team's integration server. [More](#)
© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

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

STEP 9:
Click on 'Graphical view'

Flank (green)

Intron (blue)

UTR (purple)

Coding sequence (black)

Supporting evidence

Result of STEP 8:

No. Exon / Intron	Chr	Strand	Start	End	Start Phase	End Phase	Length	Sequence
5' upstream sequence								
1	ENSE00001130431	7	1	100,156,359	100,156,552	-	1	194
Intron 1-2								
2	ENSE00001144077	7	1	100,157,117	100,157,262	1	0	146
Intron 2-3								
3	ENSE00001130423	7	1	100,158,223	100,158,402	0	0	180
Intron 3-4								
4	ENSE00001130416	7	1	100,158,403	100,158,536	0	0	134
5	ENSE00000894545	7	1	100,158,537	100,159,257	0	-	721
3' downstream sequence								

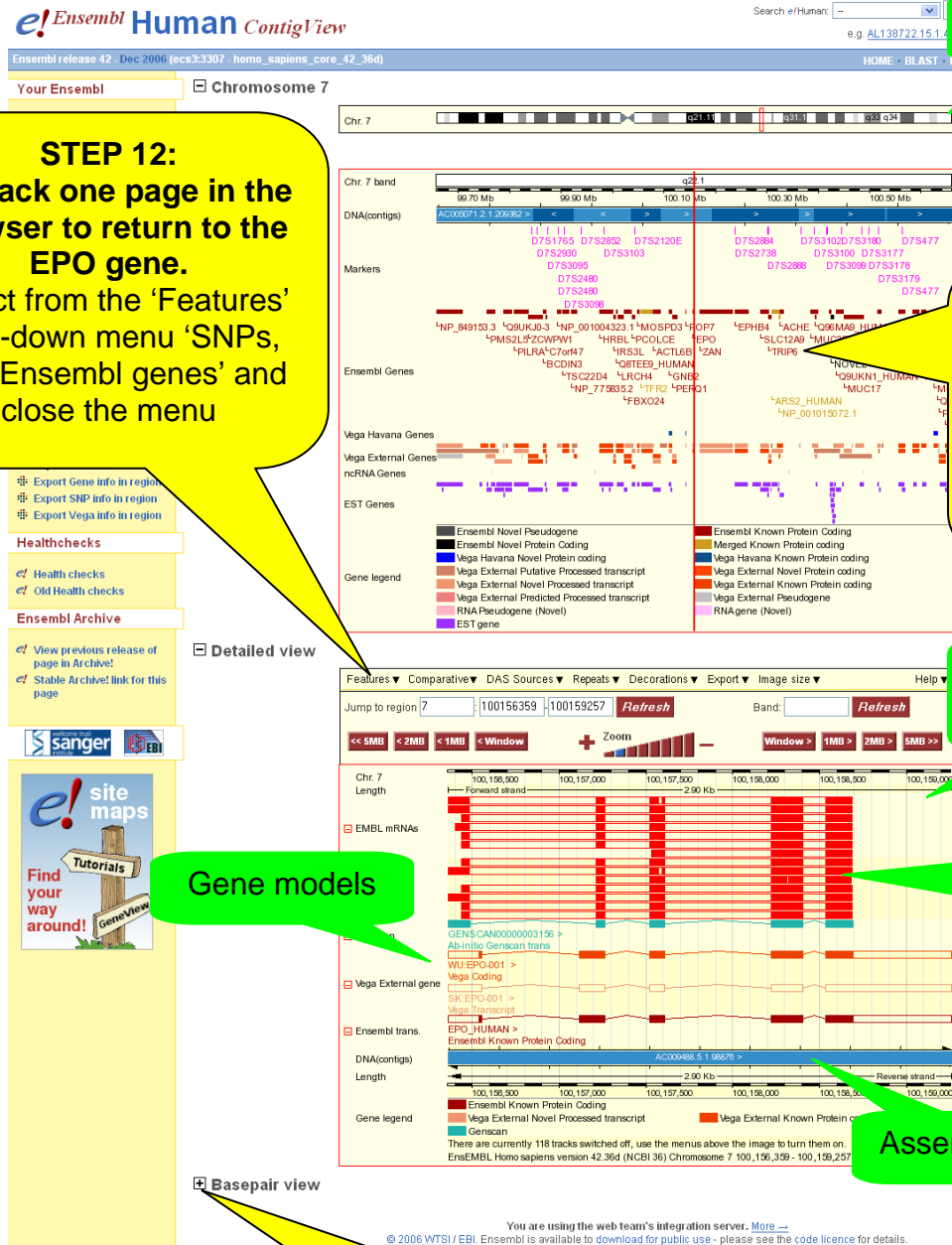
Flank
(green)

Intron
(blue)

Coding sequence
(black)

UTR
(purple)

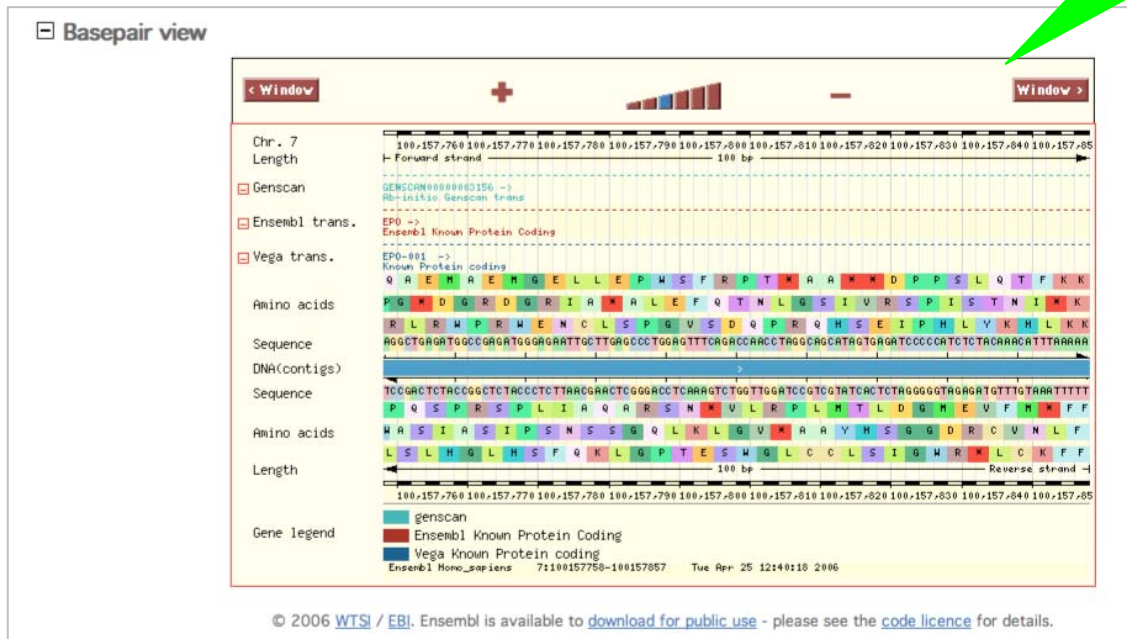
Result of STEP 9:



STEP 10:
Click on the '+' in front of 'Basepair view'

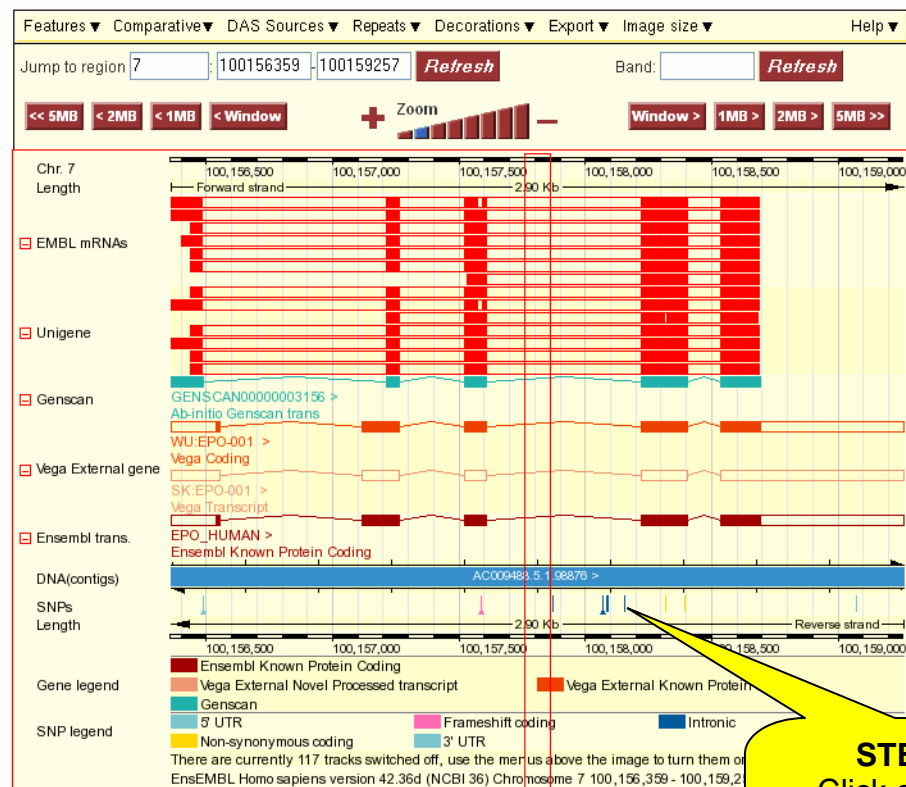
25 – 500 bp
region

Result of STEP 10:



Result of STEP 12:

Detailed view



STEP 13:
Click on a SNP
(vertical line) and
subsequently on
'SNP properties'

Your Ensemble

- Show account
- Save bookmarks
- Save configurations

dbSNP identifier

dbSNP: rs7789679

- rs7789679 - SNP info
- rs7789679 - LD info

Chromosome 7
100,158,157

- View of Chromosome 7
- Graphical view
- Graphical overview
- Export information about region
- Export sequence as FASTA
- Export EMBL file
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

Healthchecks

- Health checks
- Old Health checks

Ensembl Archive

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

Logins

Bookmarks

Settings

Groups

User accounts

New in Ensembl!

rs7789679 (dbSNP126)

None currently in the database

G/A (ambiguity code **R**)

Alleles

Validation status: Unknown

Linkage disequilibrium data: **No linkage data for this SNP**

Sequence region:

```

GAAGAGAG AGAGGT GAATT CATT TTTT TTTT TTTT CTTT CTTT TTTT GGAAGT CATT
TGCAGCCT GATTTT GGATG AAAAGGGA AATGAT CAGGAGAAAGT AAATG AGACAGCA
GAGAT GAGCTT GCTT GGGCG CAGAGGCT ACAGT CT ATT ATC CCGAGCT GAGAT GGG CAG
ATGGGAAATTT GCTT GAGGCTCT GAGTTT TGAACACACCT AGAGCAGT AGT GAT AAT CCG
CCATCT CT AC AAGCATT AAAA AATTT AGT GAGT GAGGT GGT GC ATG GGT CCA
GAT ATTT GAGAGGCT GAGGC GGGAGAGT CTTT GAGCC CAGAG AATTT GAGGCT GCACTGA
GCT GT GAT CAGACCAT CT GAGCT CAGCCTT ATG TAT GAGT GAGGCT CTT CT CAGAAA
GAAAAA AAAAAA AAAAAAT AAT GAGGCT CT AT GGAAT AC ATT ATT ATT CACTCA
CT CACT CACT CACT CATT CATT CATT CATT CATT CAGCAAGT CTT ATT GCAT CACTT CTG
TTT GCT CAGCTT GGT GCTT GAGAGT CTTT GAGGCG GAGAGGAGAGT GAGT ATG GGT CCA
GCT GACT CCGAGAT CCACCT CCT GT AGT GGGAG CAGAGGCT GT AGAAT CT GGCAGGG
CCT GGG CCT GCT GT CCGAGCT GT CTT CCGGGG CAGAGGCT GT GGT CAGT CT CTT CCA
GCCCT GGGAGG CTT CAGCT CAGT GT GAT AAAAGGCT CAGT (2BP highlighted)
          
```

SNP rs7789679 is located in the following transcripts

Genomic location (strand)	Transcript: relative SNP position	Translation: relative SNP position
7: 100158157-100158157 (1)	ENST00000252723: n/a	ENSP00000252723: n/a

Population genotypes and allele frequencies

This SNP has no allele or genotype frequencies per population.

Individual genotypes for SNP rs7789679

SNP Context - 7 100158157

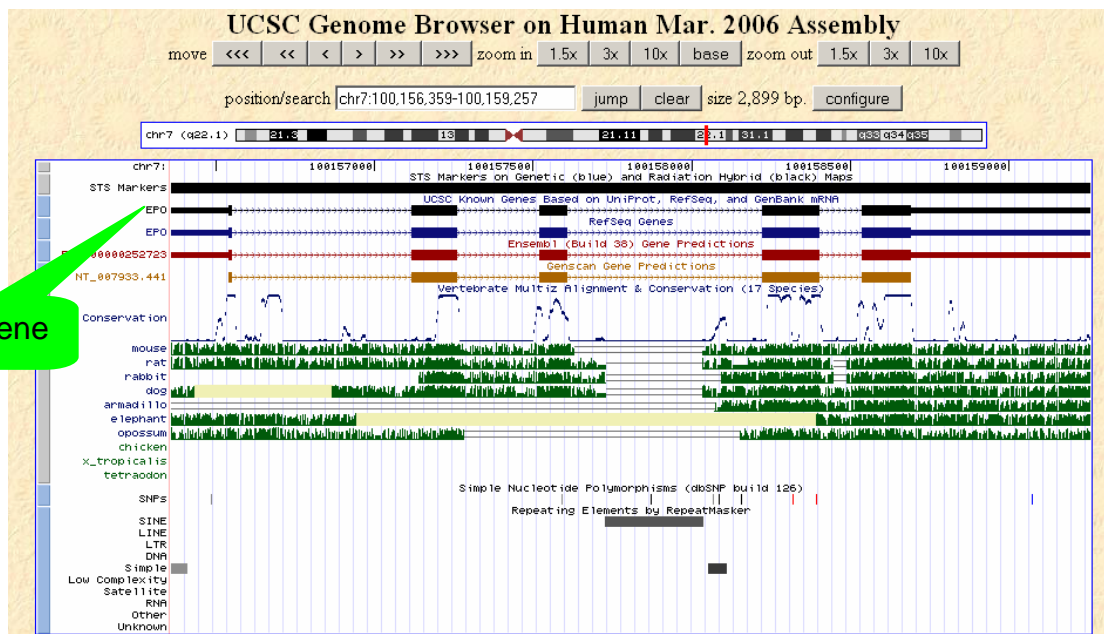
Allele and genotype frequencies

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

You are using the web team's integration server. [More →](#)
 © 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

STEP 14:
Go back to
ContigView with the
back button of the
internet browser.

STEP 15:
To see the same
chromosomal region
in the UCSC
genome browser,
click on 'Show in
UCSC browser' on
the left of the page.
A new window will
open.



Ensembl Human CytosView

Ensembl release 48

Your Ensembl

- Show account
- Save bookmarks
- Save configurations

Chromosome 7
99,657,808

- View of Chromosome
- Graphical overview
- View alignment with
- View alongside
- View Syntenic regions ...
- View region at UCSC
- View region in NCBI browser

Export data

- Export information about region
- Export sequence as FASTA
- Export Ensembl file
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

Healthchecks

- Health checks
- Old Health checks

Ensembl Archive

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

Sanger EBI

Logins

Bookmarks

Settings

Groups

User accounts

New in Ensembl!

STEP 18:
Click on 'View Syntenic regions ... with *Mus musculus*'

STEP 17:
Make sure '1Mb clones', '30k TPA clones', '32k clones' and 'Human tilepath clones' are selected under 'Decorations.' Zoom out 2 steps.

200 kb – 50 Mb region

BAC clones

Tiling path clones

Export cloneset information

Export data

Select Set of features to render: *select*


Output format: HTML

Select type to export: *select*


Export

Fields marked with * are required

You are using the web team's integration server. [More ...](#)
© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.



Human genes



Mouse homologues

e!Ensembl

Ensembl release 42

Your Ensembl

- Show account · Log out
- Save bookmark
- Save configuration as...

Chromosome 7


- View Chromosome 7
- View Chr 7 Synteny
- Map your data onto this chromosome

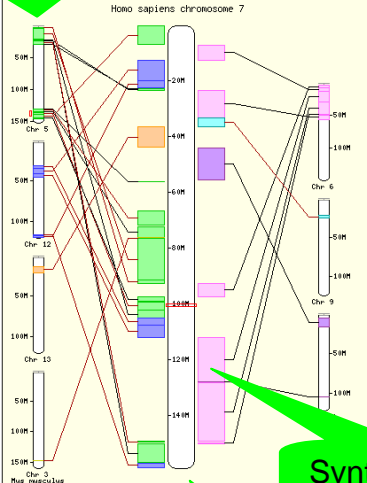
Healthchecks

- Health checks
- Old Health checks

Ensembl Archive

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





Mouse chromosomes

Human chromosome

Syntenic block

Homology Matches

Homo sapiens Genes	Mus musculus Homologues
EPO (0.10 Gb) [ContigView]	-> Epo (5: 137.71 Mb) [ContigView] [MultiContigView]
ZAN (0.10 Gb) [ContigView]	No homologues
EPHB4 (0.10 Gb) [ContigView]	-> Ephb2 (5: 137.58 Mb) [ContigView] [MultiContigView]
SLC12A9 (0.10 Gb) [ContigView]	-> Slc12a2 (18: 58.00 Mb) [ContigView] [MultiContigView]
	-> Slc12a9 (5: 137.54 Mb) [ContigView] [MultiContigView]
	-> Slc12a3 (8: 97.22 Mb) [ContigView] [MultiContigView]
	-> Slc12a4 (8: 108.83 Mb) [ContigView] [MultiContigView]
	-> Slc12a7 (13: 74.20 Mb) [ContigView] [MultiContigView]
	-> Slc12a8 (16: 33.44 Mb) [ContigView] [MultiContigView]
	-> Slc12a1 (1: 124.84 Mb) [ContigView] [MultiContigView]
	-> Slc12a5 (2: 164.66 Mb) [ContigView] [MultiContigView]
	-> Slc12a6 (2: 112.07 Mb) [ContigView] [MultiContigView]
	-> Zfx (6: 42.28 Mb) [ContigView] [MultiContigView]
	-> Wtp (7: 33.82 Mb) [ContigView] [MultiContigView]
	-> Trip6 (5: 137.54 Mb) [ContigView] [MultiContigView]
	-> Jub (14: 53.52 Mb) [ContigView] [MultiContigView]
	-> Lpp (16: 24.31 Mb) [ContigView] [MultiContigView]
	-> Lmd1 (9: 123.33 Mb) [ContigView] [MultiContigView]
	-> Fhlh1 (4: 140.85 Mb) [ContigView] [MultiContigView]
ARS2_HUMAN (0.10 Gb) [ContigView]	-> Ars2 (5: 137.53 Mb) [ContigView] [MultiContigView]
NP_001015072.1 (0.10 Gb) [ContigView]	-> 2700038N03Rik (5: 137.52 Mb) [ContigView] [MultiContigView]
	-> 1810047C23Rik (8: 47.47 Mb) [ContigView] [MultiContigView]
ACHE (0.10 Gb) [ContigView]	-> Ache (5: 137.52 Mb) [ContigView] [MultiContigView]
	-> Bche (3: 73.72 Mb) [ContigView] [MultiContigView]
ENSG00000208819 (0.10 Gb) [ContigView]	No homologues
MUC3B (0.10 Gb) [ContigView]	No homologues
Q96MA9_HUMAN (0.10 Gb) [ContigView]	No homologues
MUC12 (0.10 Gb) [ContigView]	No homologues
ENSG00000205277 (0.10 Gb) [ContigView]	No homologues
Q9UKH1_HUMAN (0.10 Gb) [ContigView]	No homologues
MUC17 (0.10 Gb) [ContigView]	No homologues

Navigate Homology

[Upstream](#) (<0.10 Gb) [Downstream](#) (>0.10 Gb)

Change Chromosome

Chromosome 7

Fields marked with * are required

You are using the web team's integration server. [More](#) →

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

Your Ensembl

- Show account : Log out
- Save bookmark
- Save configuration as...

Chromosome 7
100,155,359 - 100,160,257

- View of Chromosome 7
- Graphical view of...
- Graphical overview
- View alignment with ...
- View alongside ...
- View Syntenic regions ...
- View region at UCSC
- View region in NCBI browser

Export data

- Export information about region
- Export sequence as FASTA
- Export EMBL file
- Export Gene info in region
- Export NP info in region
- Export info in region

STEP 20:
Click on 'Export
sequence as FASTA'

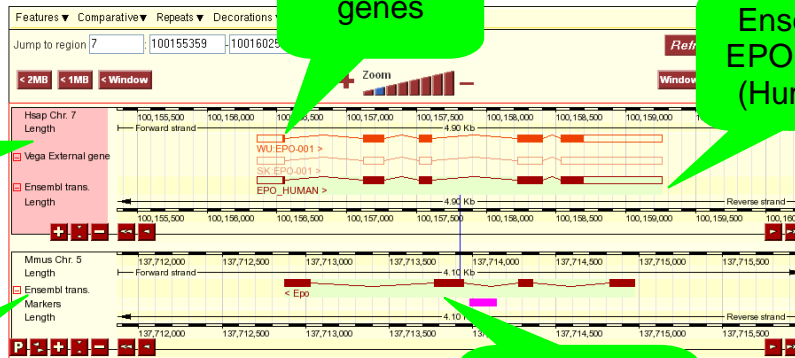
Top level



Navigational overview



Detailed View



Human

Mouse

Vega EPO genes

Ensembl EPO gene (Human)

Mouse EPO homologue



e!Ensembl Human ExportView

Search e!Human:

e.g. AL138722.15.1.44776, ENSG00000139618

STEP 21:
Click on
[Continue>>]

File output for FASTA format text file

Chromosome 7 100,155,359 - 100,160,257.

Continue >>

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

e!Ensembl Human ExportView

Search e!Human:

e.g. AL138722.15.1.44776, ENSG00000139618

STEP 22:
Select and copy a part of the sequence

STEP 23:
Click on 'BLAST'

Chromosome 7
100,155,359 - 100,160,257

View of Chromosome 7
Graphical View
Graphical overview
View alignment with ...
View alongside ...
View Syntenic regions ...
View region at UCSC
View region in NCBI
browser

Export data

- Export information about region
- Export sequence as FASTA
- Export EMBL file
- Export Gene info in region
- Export SNP info in region
- Export Vega info in region

Healthchecks

- Health checks
- Old Health checks

Ensembl Archive

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

Saccharomyces cerevisiae
Yeast
Revised assembly

You are using the web team's integration server. [More ...](#)
© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

STEP 24:
Paste the copied sequence

STEP 25:
Select 'Homo_sapiens' and 'BLASTN' and click on [RUN>]

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.

STEP 26:
Click on [Retrieve] to check for results

STEP 27:
Click on [VIEW]

Summary of BLAST search

© 2006 WTSI / EBI. Ensembl is available to [download for public use](#) - please see the [code licence](#) for details.



e!Ensembl Human BlastView

Search e!Human: Anything

e.g. AL138722.15.1.44776, ENSG00000139618

Use Ensembl

- Run a search
- Search
- Data mining (BioMart)
- Upload your own data
- Export data
- Download data

Docs and downloads

- Information
- What's New
- About Ensembl
- Ensembl data
- Software

Other links

- Home
- Sitemap
- Vega
- Pre Ensembl
- View previous release of page in Archive!
- Stable Archive! link for this page
- Archive! sites
- Trace server

Location of hits on the genome

alignments vs Homo_sapiens LATESTGP database

alignments of 107, sorted by Raw Score

Alignment Locations vs. Karyotype (click arrow to hide)

Best hit

Alignment Locations vs. Query (click arrow to hide)

Alignment Summary (click arrow to hide)

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

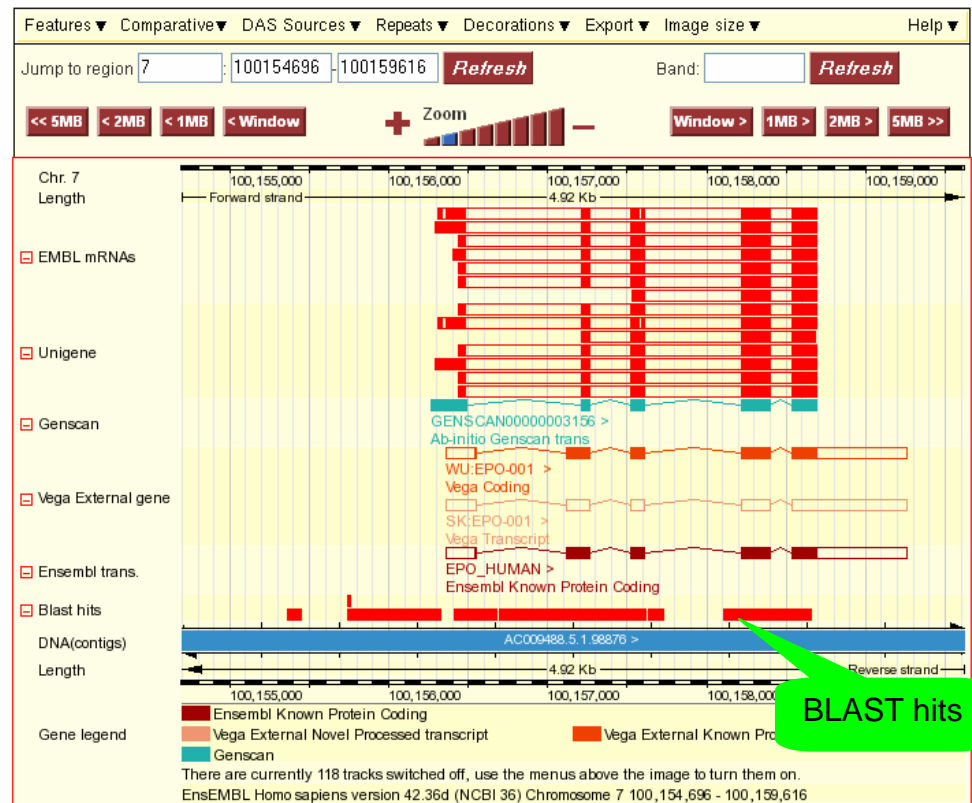
Query	Subject	Chromosome	Supercontig	Clone	Contig	Chromosome	Stats	Sort By
off Name Start	_off_ Name Start	_off_ Name Start	_off_ Name Start	_off_ Name Start	_off_ Name Start	_off_ Name Start	_off_ Name Start	
Links	Query	Start End Ori	Chromosome	Name Start End Ori	Chromosome	Name Start End Ori	Score E-val W/D Length	
[A] [S] [G] [C]	1538 2258 +	Chr-7	100156696 100157616 +	Chr-7	100156696 100157616 +	921 0. 100.00 921		
[A] [S] [G] [C]	384 144 +	Chr-7	100155745 100156329 +	Chr-7	100155745 100156329 +	585 0. 100.00 585		
[A] [S] [G] [C]	1051 1788 +	Chr-7	100156409 100156678 +	Chr-7	100156409 100156678 +	270 0. 100.00 270		
[A] [S] [G] [C]	2753 2880 +	Chr-7	100155111 100158238 +	Chr-7	100155111 100158238 +	128 0. 100.00 128		
[A] [S] [G] [C]	2274 2360 +	Chr-7	100157632 100157727 +	Chr-7	100157632 100157727 +	96 0. 100.00 96		
[A] [S] [G] [C]	1 90 +	Chr-7	100155449 100155449 +	Chr-7	100155449 100155449 +	91 0. 100.00 91		
[A] [S] [G] [C]	1571 1611 +	Chr-7	5074992 5074992 +	Chr-7	5074992 5074992 +	22 0.84 100.00 22		
[A] [S] [G] [C]	2852 2912 +	Chr-7	24002631 24002631 +	Chr-7	24002631 24002631 +	22 0.84 93.33 30		
[A] [S] [G] [C]	422 440 +	Chr-7	29854974 29854974 +	Chr-7	29854974 29854974 +	21 0.25 100.00 21		
[A] [S] [G] [C]	2335 2412 +	Chr-7	121454290 121454290 +	Chr-7	121454290 121454290 +	21 0.88 100.00 21		
[A] [S] [G] [C]	740 770 +	Chr-7	156849279 156849279 +	Chr-7	156849279 156849279 +	21 2.7 96.00 25		
[A] [S] [G] [C]	1212 1212 +	Chr-7	62162973 62162973 +	Chr-7	62162973 62162973 +	21 3.2 96.00 25		
[A] [S] [G] [C]	877 880 +	Chr-7	154980882 154980882 +	Chr-7	154980882 154980882 +	21 3.3 100.00 21		
[A] [S] [G] [C]	807 810 +	Chr-7	245694404 245694404 +	Chr-7	245694404 245694404 +	21 3.3 96.00 25		
[A] [S] [G] [C]	2814 2814 +	Chr-7	43166599 43166599 +	Chr-7	43166599 43166599 +	21 3.3 100.00 21		
[A] [S] [G] [C]	1434 1434 +	Chr-7	86415907 86415907 +	Chr-7	86415907 86415907 +	21 3.3 96.00 25		
[A] [S] [G] [C]	423 445 +	Chr-7	69322913 69322913 +	Chr-7	69322913 69322913 +	21 3.3 100.00 21		
[A] [S] [G] [C]	1806 1826 +	Chr-3	47263809 47263829 +	Chr-3	47263809 47263829 +	21 3.3 100.00 21		
[A] [S] [G] [C]	877 897 +	Chr-11	45147216 45147236 +	Chr-11	45147216 45147236 +	21 4.1 100.00 21		
[A] [S] [G] [C]	899 918 +	Chr-4	1010012 1010031 +	Chr-4	1010012 1010031 +	20 0.074 100.00 20		
[A] [S] [G] [C]	563 582 +	Chr-1	23045906 23045925 +	Chr-1	23045906 23045925 +	20 2.6 100.00 20		
[A] [S] [G] [C]	875 897 +	Chr-17	2245254 2245277 +	Chr-17	2245254 2245277 +	20 2.8 95.83 24		
[A] [S] [G] [C]	1072 1094 +	Chr-11	2404836 2404859 +	Chr-11	2404836 2404859 +	20 3.1 95.83 24		
[A] [S] [G] [C]	2235 2254 +	Chr-2	55009141 55009160 +	Chr-2	55009141 55009160 +	20 3.2 100.00 20		
[A] [S] [G] [C]	873 896 +	Chr-19	18119374 18119397 +	Chr-19	18119374 18119397 +	20 3.4 95.83 24		
[A] [S] [G] [C]	809 828 +	Chr-12	21163734 21163753 +	Chr-12	21163734 21163753 +	20 4.0 100.00 20		
[A] [S] [G] [C]	878 897 +	Chr-2	43261275 43261294 +	Chr-2	43261275 43261294 +	20 4.1 100.00 20		
[A] [S] [G] [C]	661 679 +	Chr-1	18606381 18606399 +	Chr-1	18606381 18606399 +	19 0.19 100.00 19		
[A] [S] [G] [C]	1768 1797 +	Chr-1	18656691 18656720 +	Chr-1	18656691 18656720 +	19 0.19 90.32 31		
[A] [S] [G] [C]	1741 1770 +	Chr-10	45140401 45140450 +	Chr-10	45140401 45140450 +	19 1.2 90.32 31		
[A] [S] [G] [C]	418 436 +	Chr-20	44598942 44598960 +	Chr-20	44598942 44598960 +	19 2.4 100.00 19		
[A] [S] [G] [C]	2209 2227 +	Chr-14	46044909 46044927 +	Chr-14	46044909 46044927 +	19 3.2 100.00 19		
[A] [S] [G] [C]	2804 2822 +	Chr-3	12983509 12983509 +	Chr-3	12983509 12983509 +	19 4.3 100.00 19		
[A] [S] [G] [C]	881 899 +	Chr-14	104884074 104884092 +	Chr-14	104884074 104884092 +	19 5.3 100.00 19		
[A] [S] [G] [C]	881 899 +	Chr-14	104884028 104884046 +	Chr-14	104884028 104884046 +	19 5.3 100.00 19		
[A] [S] [G] [C]	881 899 +	Chr-14	104883982 104884000 +	Chr-14	104883982 104884000 +	19 5.3 100.00 19		
[A] [S] [G] [C]	69 86 +	Chr-7	100155741 100155758 +	Chr-7	100155741 100155758 +	18 0. 100.00 18		
[A] [S] [G] [C]	1905 1922 +	Chr-X	39633546 39633563 +	Chr-X	39633546 39633563 +	18 0.70 100.00 18		
[A] [S] [G] [C]	1962 1979 +	Chr-X	70306710 70306727 +	Chr-X	70306710 70306727 +	18 3.9 100.00 18		
[A] [S] [G] [C]	1700 1720 +	Chr-X	70247948 70247969 +	Chr-X	70247948 70247969 +	18 3.9 95.45 22		
[A] [S] [G] [C]	2306 2323 +	Chr-1	150284632 150284649 +	Chr-1	150284632 150284649 +	18 6.0 100.00 18		
[A] [S] [G] [C]	1801 1818 +	Chr-2	73022513 73022530 +	Chr-2	73022513 73022530 +	18 9.3 100.00 18		
[A] [S] [G] [C]	74 90 +	Chr-2	29824827 29824843 +	Chr-2	29824827 29824843 +	17 0.25 100.00 17		
[A] [S] [G] [C]	881 897 +	Chr-X	39611091 39611107 +	Chr-X	39611091 39611107 +	17 0.70 100.00 17		
[A] [S] [G] [C]	881 897 +	Chr-1	20541390 20541406 +	Chr-1	20541390 20541406 +	17 1.1 100.00 17		
[A] [S] [G] [C]	1798 1814 +	Chr-1	20511096 20511112 +	Chr-1	20511096 20511112 +	17 1.1 100.00 17		
[A] [S] [G] [C]	1487 1503 +	Chr-10	43031647 43031665 +	Chr-10	43031647 43031665 +	17 2.2 100.00 17		
[A] [S] [G] [C]	2100 2119 +	Chr-14	104633755 104633775 +	Chr-14	104633755 104633775 +	17 1.8 95.24 21		
[A] [S] [G] [C]	1564 1580 +	Chr-9	139178040 139178056 +	Chr-9	139178040 139178056 +	17 2.2 100.00 17		
[A] [S] [G] [C]	883 899 +	Chr-9	139163518 139163534 +	Chr-9	139163518 139163534 +	17 2.2 100.00 17		
[A] [S] [G] [C]	2805 2825 +	Chr-14	46020025 46020044 +	Chr-14	46020025 46020044 +	17 3.2 95.24 21		
[A] [S] [G] [C]	1577 1593 +	Chr-16	88233355 88233371 +	Chr-16	88233355 88233371 +	17 3.5 100.00 17		
[A] [S] [G] [C]	1615 1629 +	Chr-19	47367644 47367660 +	Chr-19	47367644 47367660 +	17 5.4 100.00 17		
[A] [S] [G] [C]	1805 1821 +	Chr-1	150273546 150273562 +	Chr-1	150273546 150273562 +	17 6.0 100.00 17		
[A] [S] [G] [C]	737 756 +	Chr-6	34166916 34166936 +	Chr-6	34166916 34166936 +	17 7.0 95.74 21		

STEP 28:
Click on [C] in front of best hit

Alignment of hits to query sequence

Back in the contigview page...

Detailed view



END of the
Worked Example

EXERCISES and ANSWERS

(NOTE: please use release 42 (from the archive site) as the answers may have changed in newer releases. Thank you.)

1. Exploring features related to a gene

(a) 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) Look up 'Similarity Match' in the glossary:

Follow the 'Help and Documentation' link, 'HelpDesk section' to this url:

http://www.ensembl.org/Homo_sapiens/glossaryview

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?

2. Exploring a region

(a) Display the region between markers D12S764 and D12S1871 in ContigView. Start on the human homepage, and click on chromosome 12.

(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) What is the closest marker to the TENC1 gene? How many synonyms does this marker have?

(e) Zoom in (towards the '+') three steps on the zoom triangle and turn on the SNP track. Identify an intronic SNP and look at the corresponding SNPView page.

3. Exploring the zebrafish (*Danio rerio*) genome with Ensembl

- (a) Bring up a ContigView display of zebrafish (*Danio rerio*) chromosome 1 between 64.0 Mb and 64.5 Mb.
- (b) How many 'known' and 'novel' genes are predicted in this region? For one of the known genes, find some information about its function, and look at an entry for it in EntrezGene, UniProt/Swiss-Prot or the ZFIN site.
- (c) Can you find out anything about the possible functions of one of the novel genes? For this, try looking at homologues in other species, at other members of protein families and InterPro domains.

Answers (Browsing Ensembl)

1. Exploring features related to a gene

- (a) A 'Vega' gene and 'Ensembl' gene will be shown. VEGA (Vertebrate Genome Annotation) is a consortium of manual curators for certain chromosomes in human, mouse, zebrafish, pig and dog. However, we would like to explore the 'Ensembl Gene: ENSG00000006128'. 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'. Click on the 'Ensembl Gene: ENSG00000006128' link to go to the GeneView page for this gene.
- (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. The length of this transcript is 1060 bp. It has 7 exons and codes for 129 aa.
- (c) The TAC1 gene is Protachykinin 1 precursor. Follow the links to MIM and EntrezGene or UniProt/Swiss-Prot in the 'Similarity Matches' section to learn more. Choose 'UniProt' under 'DAS Sources' to see references in the literature (click 'Update' after making the selection). 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 (loading of the page can take a while) showing the position of that term in the GO structure (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). If you click on the contig and follow the link to the EMBL source (or if you turn on the 'Human tile path clones' track from the 'Decorations' menu of ContigView) you can see that this sequence is derived from clone RP5-841B21.

(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.

2. Exploring a region

(a) Start on the homepage for human and click on chromosome 12 to go MapView. 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 displayed region in the Overview panel is larger than the area between the two markers. The red line or small box is drawn over the first marker (D12S764). Zoom into the region between the two markers by drawing a box with the mouse around it.

The region will be displayed below, 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. To see also the 4 clones that make up this region, select the 'Human tilepath clones' track from the 'Decorations' menu. Clones are shown in gold and pink. Portions of the 'Tile path clones' were used to form the assembly and correspond to 'contigs'. The clones overlap each other whereas the contigs don't.

(c) Marker D12S2110 is closest to the TENC1 gene. Click on the marker (i.e. on the vertical bar representing it, not on its name) and follow the link 'Marker info' to the MarkerView page. There are 2 synonyms listed.

(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.

3. Exploring the zebrafish (*Danio rerio*) genome with Ensembl

(a) Start on the homepage for zebrafish (*Danio rerio*). In the 'Karyotype' section choose 'Chromosome: 1', 'From (type): Base pair: 64000000 To (type) Base pair: 64500000' and click [Go]. This leads you to ContigView for a larger region. Type in the base pairs in Detailed View (change the second number to 64500000).

(b) On the 'Overview' panel of ContigView Ensembl known and novel genes are displayed in the 'Ensembl Genes' track in reddish brown and black, respectively. Known genes are Ensembl gene predictions that match species specific entries in the UniProt and/or RefSeq database, while novel genes map back to entries from other species. There are 12 known and 9 novel genes. Click on one of the known genes to go to its GeneView page and explore the links in the 'Similarity Matches' section.

(c) To find out more about the possible function of a novel gene there are many options. Click on the gene in ContigView to go to its GeneView page. If the gene has orthologues in other species (shown in the 'Orthologue Prediction' section) and these orthologues are better characterized than the novel zebrafish gene this can give a clue about the possible function of this gene. If the gene belongs to a family (shown in the 'Protein Family' section) other family members may provide information. InterPro domains (shown in the 'InterPro' section) may also provide clues.