

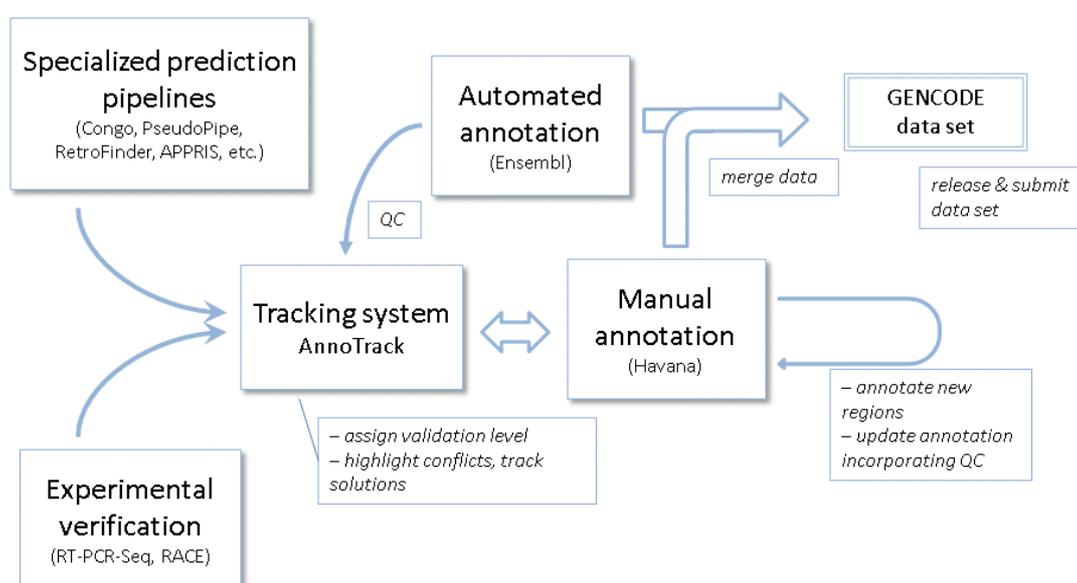
Overview of the GENCODE reference gene set

Aim

This module will give an overview of the GENCODE gene set that is available from the genome browsers and explain how ENCODE data is integrated to improve the set.

Introduction

Schematic showing interconnection between different GENCODE pipelines



HAVANA (Human and Vertebrate Analysis and Annotation) group at the Wellcome Trust Sanger Institute (WTSI) perform manual genome annotation. Finished genomic sequence is analysed on a clone by clone basis using a combination of similarity searches against DNA and protein databases (including cross-species) and a series of *ab initio* gene predictions. Annotation is based on supporting evidence, which is external sequence such as ESTs, cDNAs and protein. There are multiple biotypes that reflect confidence levels and there are additional data sources included as DAS tracks (e.g. CAGE tags, RNAseq).

The GENCODE reference gene set (<http://www.gencodegenes.org/>) is produced in collaboration between HAVANA and Ensembl and is available for human and mouse. The HAVANA group provides manual annotation of the

genes and transcripts onto the genome, whilst the Ensembl group performs automatic gene annotation. These data are merged to produce the GENCODE gene set, which is used by a variety of projects, such as Ensembl, ENCODE, 1000 genomes, UCSC and more.

The complete GENCODE set is available to view in genome browsers, such as Ensembl (<http://www.ensembl.org/>) and UCSC (<http://genome.ucsc.edu/>). The manually annotated genes only can be seen in the Vega (Vertebrate Genome Annotation) database (<http://vega.sanger.ac.uk/>).

The GENCODE gene set is an important contributor to the Consensus CDS (CCDS) project, which is a collaborative effort between the European Bioinformatics Institute (EBI), the National Centre for Biotechnology Information (NCBI), the Wellcome Trust Sanger Institute (WTSI) and the University of California at Santa Cruz (UCSC). The aim of the project is to identify a core set of human protein coding regions that are consistently annotated between the different institutes. The long-term goal is to support convergence towards a standard set of gene annotations. The CCDS gene set is generated by Ensembl and NCBI and there is extensive QC by WTSI, NCBI and UCSC. A set of guidelines have been developed for the annotation of coding sequence regions by the collaborating Institutes, and any changes to the CCDS set have to be agreed by all three sites.

Demo: Looking at GENCODE genes in genome browsers

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

The screenshot shows the Ensembl website homepage with several callout boxes pointing to specific features:

- Link back to homepage**: Points to the Ensembl logo.
- Ensembl tools**: Points to the navigation menu (BLAST/BLAT, BioMart, Tools, Downloads, Help & Documentation, Blog, Mirrors).
- Blue bar remains visible on every Ensembl page**: Points to the dark blue header bar.
- Search**: Points to the search bar at the top right.
- Search**: Points to the search input field in the main content area.
- News**: Points to the 'What's New in Release 73' section.
- Drop-down list of species**: Points to the 'Select a species' dropdown menu.
- How-tos for commonly used Ensembl features**: Points to the 'Did you know...?' section.

We're going to look at the human *ESPN* gene. This gene encodes a multifunctional actin-bundling protein with a major role in mediating sensory transduction in various mechanosensory and chemosensory cells. Mutations in this gene are associated with deafness (<http://tinyurl.com/espn-ncbi-gene>).

From ensembl.org, type *ESPN* into the search bar and click the **Go** button. You will get a list of hits with the human gene at the top.

Where you search for something without specifying the species, or where the ID is not restricted to a single species, the most popular species will appear first, in this case, human, mouse and zebrafish appear first. You can restrict your query to species or features of interest using the options on the left.

Restrict categories to:

Gene	53
Transcript	114
Variation	10
Marker	1
Somatic Mutation	22

Restrict species to:

Human	54
Mouse	41
Zebrafish	11

Search results for **espn**:

- ESPN (Human Gene)**
ENSG00000187017 1:6484848-6521430
espn [Source:HGNC Symbol;Acc:HGNC:13281]
Variation table • Location • Regulation • Orthologues • Gene tree
- Espn (Mouse Gene)**
ENSMUSG0000028943 4:152120331-152152371-1
espn [Source:MGI Symbol;Acc:MGI:1861630]
Variation table • Location • Regulation • Orthologues • Gene tree
- espn (Zebrafish Gene)**
ENSDARG00000076414 8:49095530-49250014-1
espn [Source:ZFIN;Acc:ZDB-GENE-081105-173]
Variation table • Location • Regulation • Orthologues • Gene tree

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

Gene tab

Option: Open table of transcripts

Information about annotation

ESPIN-001 transcript. Click for info

Gene views

Forward-stranded transcripts

Reverse-stranded transcripts

Blue bar is the genome

The screenshot shows the Ensembl Gene page for **Gene: ESPN ENSG00000187017**. The page includes a navigation menu on the left with categories like Gene-based displays, Comparative Genomics, and Genetic Variation. The main content area displays gene details such as Description, Location, INSDC coordinates, and Transcripts. A 'Show transcript table' button is highlighted. Below the text, a genomic track visualization shows the gene structure with exons and introns, and various transcripts (ESPIN-001 to ESPIN-008) are shown in different colors. A blue bar at the bottom represents the genome. Callouts provide additional context for these elements.

From this page we can see that the gene was annotated by both Ensembl automatic and Havana manual annotation. It is also a member of the CCDS

set. Click on [Show transcript table](#) to see which transcripts have a CCDS associated with them.

The first transcript has a CCDS associated with it

Show/hide columns		Filter					
Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CDS incomplete	CCDS
ESPN-001	ENST00000377828	3531	ENSP00000367059	854	Protein coding	-	CCDS70
ESPN-009	ENST00000461727	1869	ENSP00000465308	288	Protein coding	-	-
ESPN-201	ENST00000416731	1665	ENSP00000399239	288	Protein coding	-	-
ESPN-007	ENST00000434576	750	ENSP00000413621	188	Protein coding	5'	-
ESPN-002	ENST00000418286	641	ENSP00000401793	214	Protein coding	5' and 3'	-
ESPN-005	ENST00000478323	270	ENSP00000466437	28	Protein coding	3'	-
ESPN-004	ENST00000475228	813	No protein product	-	Processed transcript	-	-
ESPN-008	ENST00000468561	664	No protein product	-	Processed transcript	-	-
ESPN-006	ENST00000475479	360	No protein product	-	Processed transcript	-	-
ESPN-003	ENST00000477679	885	No protein product	-	Retained intron	-	-

Select the top transcript [ESPN-001](#) to go to the transcript tab.

This is the transcript associated with the CCDS transcript. It is shown in gold because it has identical annotation from the Ensembl automatic and Havana manual annotation.

Transcript: ESPN-001 ENST00000377828

Description: [espin](#) [Source:HGNC Symbol;Acc:13281]

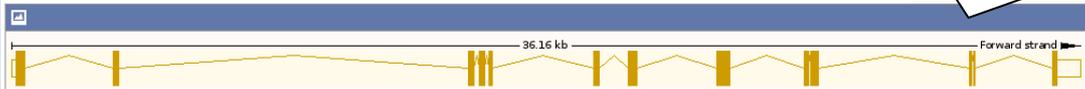
Location: [Chromosome 1: 6,484,848-6,521,004](#) forward strand.

Gene: This transcript is a product of gene [ENSG00000187017](#)

This gene has 10 transcripts (splice variants) [Hide transcript table](#)

Show/hide columns		Filter					
Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CDS incomplete	CCDS
ESPN-001	ENST00000377828	3531	ENSP00000367059	854	Protein coding	-	CCDS70
ESPN-009	ENST00000461727	1869	ENSP00000465308	288	Protein coding	-	-
ESPN-201	ENST00000416731	1665	ENSP00000399239	288	Protein coding	-	-
ESPN-007	ENST00000434576	750	ENSP00000413621	188	Protein coding	5'	-
ESPN-002	ENST00000418286	641	ENSP00000401793	214	Protein coding	5' and 3'	-
ESPN-005	ENST00000478323	270	ENSP00000466437	28	Protein coding	3'	-
ESPN-004	ENST00000475228	813	No protein product	-	Processed transcript	-	-
ESPN-008	ENST00000468561	664	No protein product	-	Processed transcript	-	-
ESPN-006	ENST00000475479	360	No protein product	-	Processed transcript	-	-
ESPN-003	ENST00000477679	885	No protein product	-	Retained intron	-	-

Transcript summary ⓘ



Statistics: Exons: 13 Coding exons: 13 Transcript length: 3,531 bps Translation length: 854 residues

CCDS: This transcript is a member of the Human CCDS set: [CCDS70](#)

Ensembl version: ENST00000377828.1

Type: Known protein coding

Prediction Method: Transcript where the Ensembl genebuild transcript and the [Vega](#) manual annotation have the same sequence, for [article](#).

Alternative transcripts: This transcript corresponds to the following database identifiers:
Transcript having exact match between ENSEMBL and HAVANA: [OTTHUMT0000001887](#) (version 3)

Graphic of the transcript model

Information about transcript annotation

Click on the [CCDS370](#) to go to the CCDS record. This will open in a new tab in your browser.

Report for CCDS70.1 (current version)

CCDS	Status	Species	Chrom.	Gene	CCDS Release	NCBI Annotation Release	Ensembl Annotation Release	Links
70.1	Public	<i>Homo sapiens</i>	1	ESPN	14	105	73	H G C G

Public since: CCDS release 1, NCBI annotation release 35.1, Ensembl annotation release 23

Sequence IDs included in CCDS 70.1

Original	Current	Source	Nucleotide ID	Protein ID	Status in CCDS	Seq. Status	Links
✓	✓	EBL,WTSI	ENST00000377828	ENSP00000367059	Accepted	alive	N P N P
✓	✓	EBL,WTSI	OTTHUMT0000001887	OTTHUMP0000000828	Accepted	alive	N P N P
✓	✓	NCBI	NM_031475.2	NP_113663.2	Accepted	alive	N P N P B

Chromosomal Locations for CCDS 70.1

Assembly GRCh37.p13 ([GCF_000001405.25](#))

On '+' strand of Chromosome 1 (NC_000001.10)

Genome Browser links: [N](#)[N](#)[U](#)[E](#)[E](#)[V](#)

Chromosome	Start	Stop	Links
1	6485016	6485309	N N U E E V
1	6488286	6488479	N N U E E V
1	6500314	6500500	N N U E E V
1	6500686	6500868	N N U E E V
1	6500994	6501125	N N U E E V
1	6504541	6504742	N N U E E V
1	6505724	6505995	N N U E E V
1	6508701	6509151	N N U E E V
1	6511663	6511808	N N U E E V
1	6511893	6512156	N N U E E V
1	6517244	6517323	N N U E E V
1	6517421	6517432	N N U E E V
1	6520059	6520206	N N U E E V

CCDS Sequence Data

Blue highlighting indicates alternate exons.
Red highlighting indicates amino acids encoded across a splice junction.

Mouse over the nucleotide or protein sequence below and click on the highlighted codon or residue to select the pair.

Nucleotide Sequence (2565 nt):

ATGGCCCTGGAGCAGGCGCTGCAGGCGGCGCGGCAGGGCGAGCTGGAGCTGCTGAGGTCGCTGCACGCCG
CAGGCCTCCTGGGGCCCTCGCTGCGCGACCCGCTGGACGCGCTGCCCGTGACACCACGCGGCCCGCGCTGC
GAAGCTGCACTGTCTGCGCTTCTGGTGGAGGAAGCGCCCTCCCGCCGCGGCCGCGCCGCAACGGC
GCCACACGGGCCACGACGCTCCGCCACCGGCCACCTCGCTGCGTGCAGTGGCTGCTGTCGAGGGCG
GCTGCAGAGTGCAGGACAAAGACAATTCTGGTGCACAGTCTTGCATCTGGTGCCTGCGCCCTGCGGCCACCC
CGAGGTGGTGAAGTGGCTCTTGCATCATGGCGGTGGGGACCCACCGCGGCCACAGACATGGGCGCCCTG
CCTATCCACTACGCTGCCCAAAGGAGACTTCCCTCCCTGAGGCTTCTCGTGCAGCACTACCCTGAGG
GAGTGAATGCCAAACCAAGAACGGTGCCACGCCCTGTACCTGGCGTGCCAGGAGGGCCACCTGGAGGT

Summary

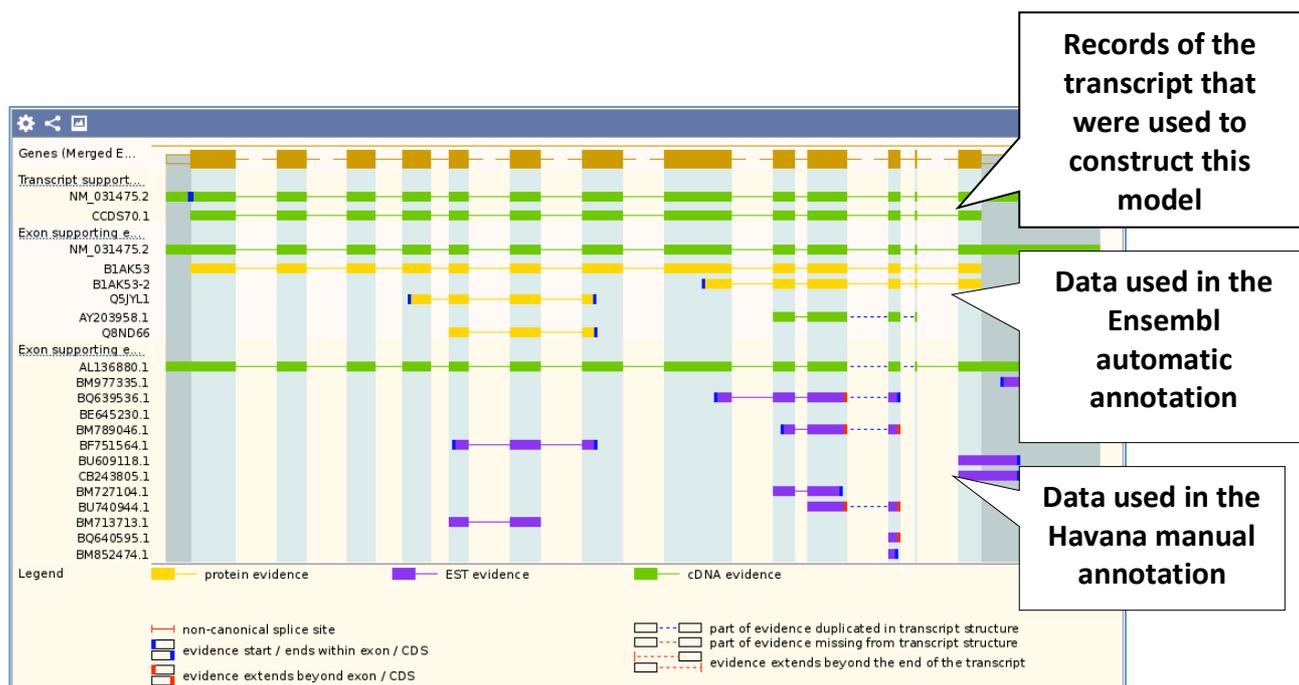
Links to the transcript in Ensembl, Vega and NCBI.

Summary of exons

Sequence

This page summarises the CCDS transcript.

Go back to the Ensembl page and click on [Supporting evidence](#) in the left-hand menu.



You can see that Ensembl and Havana used different pieces of evidence to construct their transcript model, yet still came up with the same model, demonstrating how reliable the model is.

Click on [General Identifiers](#) in the left-hand menu. This lists records of the transcript and its protein product in other databases.

General identifiers ⓘ

This transcript corresponds to the following database identifiers:

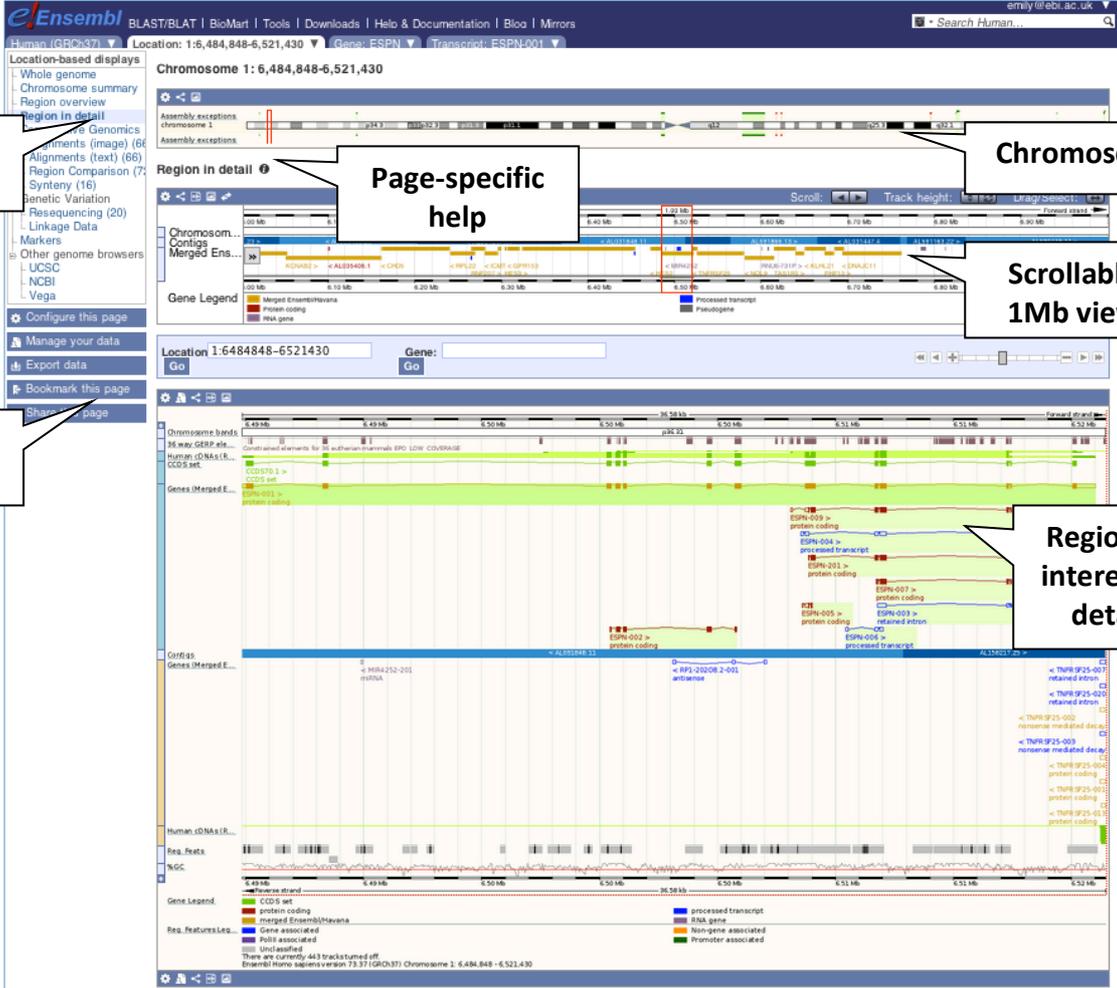
External database	Database identifier
HGNC Symbol	ESPN espinn [view all locations]
UniParc	UPI000013D2B6 [view all locations]
CCDS	CCDS70.1 [view all locations]
UniProtKB/Swiss-Prot	ESPN_HUMAN [align] Espinn [view all locations]
RefSeq peptide	NP_113663.2 (Target %id: 100; Query %id: 100) [align] espinn [view all locations]
RefSeq mRNA	NM_031475.2 [align] [view all locations]
UCSC Stable ID	uc001amy.3 [view all locations]
Human Protein Atlas	HPA028674 [view all locations] HPA028674 [view all locations]
European Nucleotide Archive	AF134401 [align] [view all locations] AL031848 [align] [view all locations] AL136880 [align] [view all locations] AL158217 [align] [view all locations] AY203958 [align] [view all locations] CH471130 [align] [view all locations]
HGNC transcript name	ESPN-001 espinn [view all locations]
INSDC protein ID	AAD24480.1 [align] [view all locations] AAP34481.1 [align] [view all locations] CAB66814.1 [align] [view all locations] CAI19773.1 [align] [view all locations] CAI22163.1 [align] [view all locations] EAW71537.1 [align] [view all locations]

We can also see genes and transcripts in a location. Click on the tab saying Location 1:6,484,848-6,521,530 at the top of the page.

Click on the button  to view page-specific help.

The help pages provide links to [Frequently Asked Questions](#), a [Glossary](#), [Video Tutorials](#), and a form to [Contact HelpDesk](#).

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

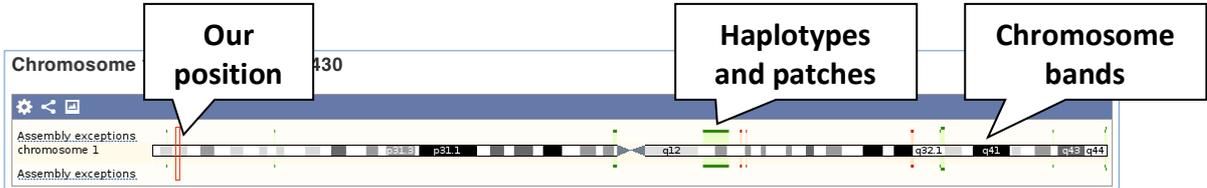


The screenshot shows the Ensembl genome browser interface. Callouts point to the following features:

- Location views**: Points to the left-hand navigation menu.
- Page-specific help**: Points to the 'Region in detail' tab.
- Chromosome**: Points to the 'Chromosome 1' header.
- Scrollable 1Mb view**: Points to the main genomic track with a scroll bar.
- Tool buttons**: Points to the 'Configure this page', 'Manage your data', 'Export data', 'Bookmark this page', and 'Share page' buttons.
- Region of interest in detail**: Points to the detailed view of the ESPN gene region.

The Region in detail page is made up of three images, let's look at each one on detail.

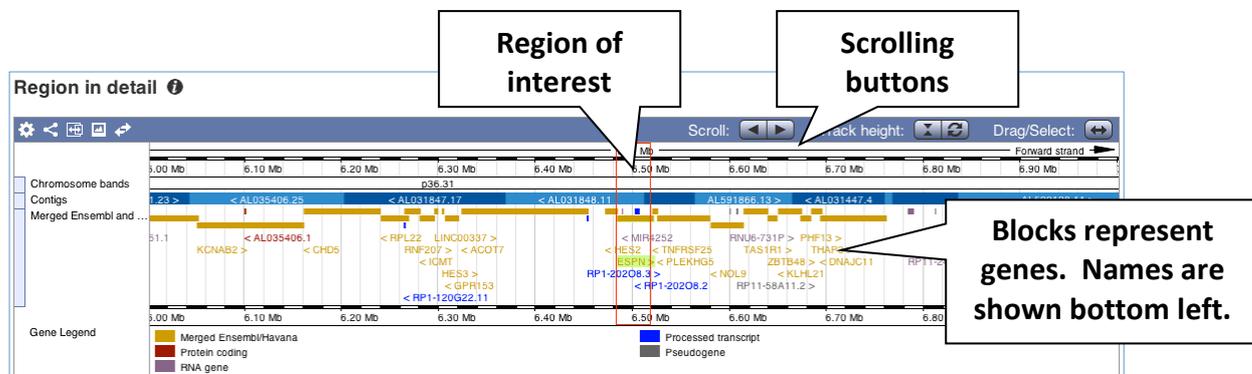
The first image shows the chromosome:



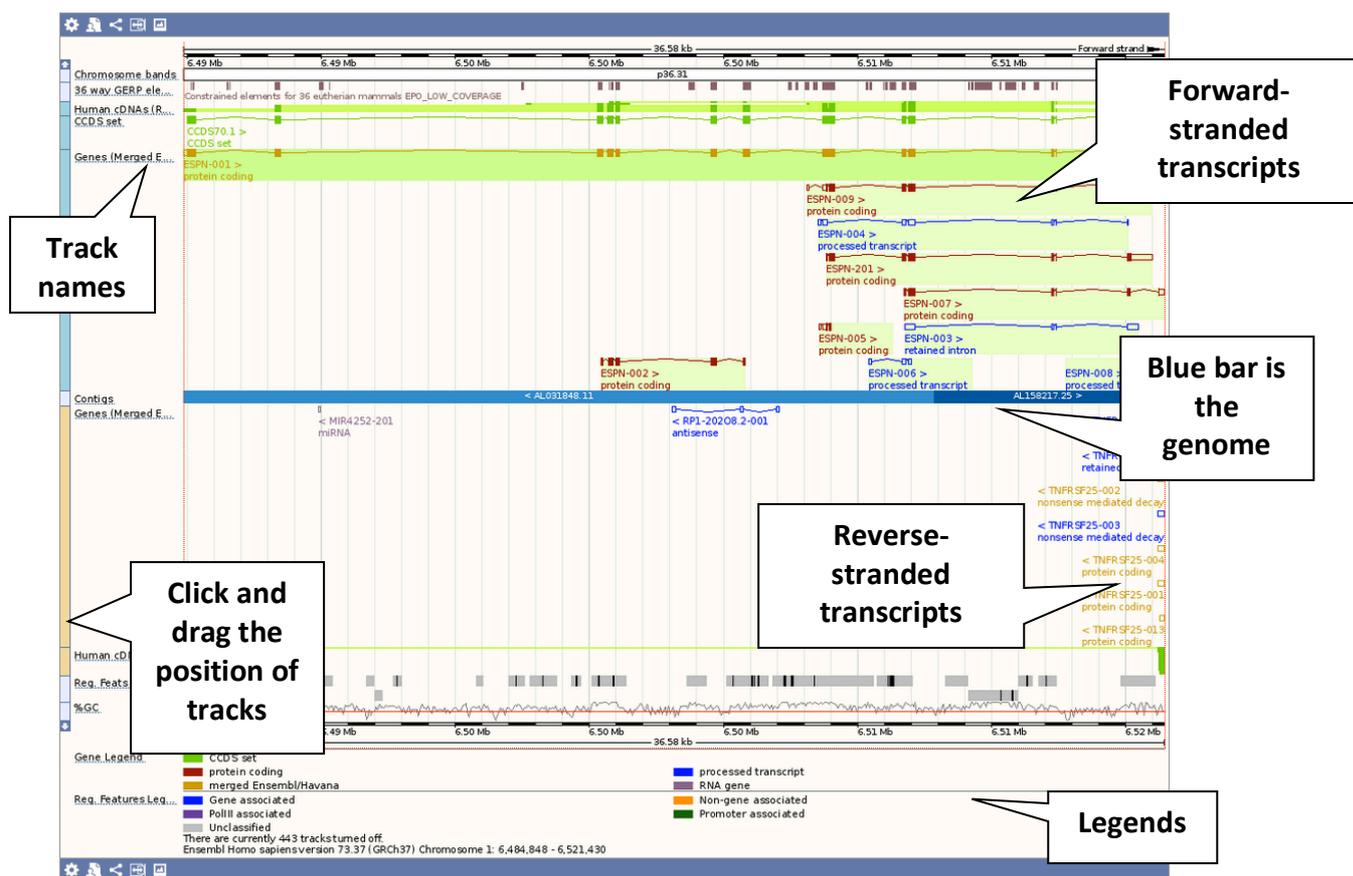
This close-up shows the chromosome overview track. Callouts point to:

- Our position**: Points to a red vertical line indicating the current location on the chromosome.
- Haplotypes and patches**: Points to the colored bars representing different haplotypes.
- Chromosome bands**: Points to the black and white bands representing cytogenetic bands (p31.1, q12, q21, q41, q44).

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

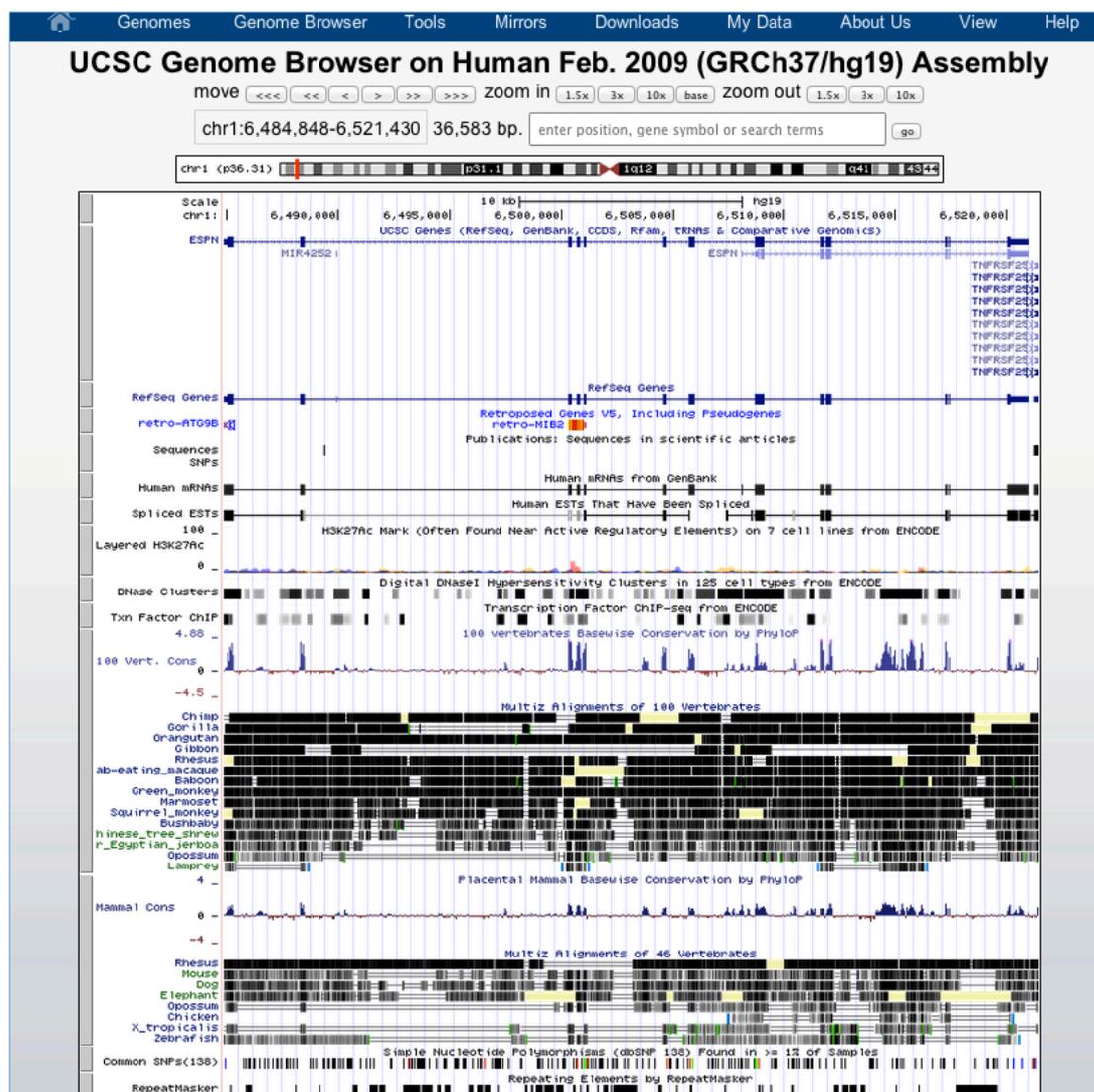


The third image is a detailed, configurable view of the region.

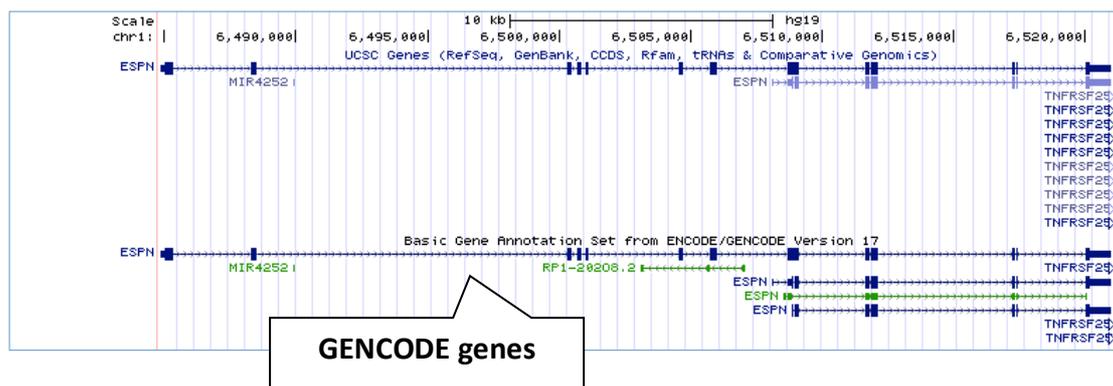


We can see the GENCODE genes in this view, with the CCDS plotted alongside (since CCDS transcripts lack UTRs).

Click on [UCSC](#) in the left-hand menu to see the same region in UCSC. This will open in a new tab.



The GENCODE gene set is not shown by default. Scroll down to [Genes and Prediction tracks](#) then select the drop down under [GENCODE](#) and chose [show](#), then click on [refresh](#).



Click on the transcripts to see information about them, including links to Ensembl and Vega.

GENCODE Transcript Annotation ENST00000377828.1 (ESPN)

	Transcript	Gene
Gencode id	ENST00000377828.1	ENSG00000187017.10
HAVANA manual id	OTTHUMT00000001887.3	OTTHUMG00000000753.5
Position	chr1:6484848-6521004	chr1:6484848-6521004
Strand	+	
Biotype	protein_coding	protein_coding
Status	KNOWN	KNOWN
Annotation Level	manual (2)	
Annotation Method	manual & automatic	manual & automatic
Transcription Support Level	ts1	
HGNC gene symbol	ESPN	
CCDS	CCDS70.1	
GeneCards	ESPN	
APPRIS	ENST00000377828.1	ENSG00000187017.10

Demo: Vega update genes

Start at the homepage for Vega (<http://vega.sanger.ac.uk>).

Click on human and search for the POLR2E gene.

POLR2E (Human Havana Gene)

OTTHUMG00000181873 19:1086594-1095598:-1

Polymerase (RNA) II (DNA directed) polypeptide E, 25kDa. *Havana annotation.*

[Location](#) • [Sequence](#)

Select the top result.

Human (VEGA54) Location: 19:1,086,594-1,095,598 Gene: POLR2E

Link to updated annotation

Gene: POLR2E OTTHUMG00000181873

Updated annotation available
There is updated annotation for this gene available here.

Description polymerase (RNA) II (DNA directed) polypeptide E, 25kDa
Location Chromosome 19: 1,086,594-1,095,598 reverse strand.
Transcripts This gene has 10 transcripts (splice variants) [Hide transcript table](#)

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
POLR2E-001	OTTHUMT00000458044	1504	OTTHUMP00000267716	210	Protein coding	CCDS12056
POLR2E-010	OTTHUMT00000458043	799	OTTHUMP00000267715	210	Protein coding	CCDS12056
POLR2E-002	OTTHUMT00000458046	1041	OTTHUMP00000267717	68	Nonsense mediated decay	-
POLR2E-006	OTTHUMT00000458049	588	OTTHUMP00000267718	51	Nonsense mediated decay	-
POLR2E-007	OTTHUMT00000458048	530	No protein product	-	Processed transcript	-
POLR2E-009	OTTHUMT00000458042	530	No protein product	-	Processed transcript	-
POLR2E-003	OTTHUMT00000458045	2441	No protein product	-	Retained intron	-
POLR2E-008	OTTHUMT00000458047	756	No protein product	-	Retained intron	-
POLR2E-005	OTTHUMT00000458050	685	No protein product	-	Retained intron	-
POLR2E-004	OTTHUMT00000458051	440	No protein product	-	Retained intron	-

The current version of Vega has 10 splice variants for POLR2E, but there is updated annotation available. Click on the update link.

Gene: POLR2E OTTHUMG00000181873

Vega update gene

This is a Havana update gene with newer annotation than the core Vega gene.

Description polymerase (RNA) II (DNA directed) polypeptide E, 25kDa

Location Chromosome 19: 1,086,578-1,095,379 reverse strand.

Transcripts This gene has 14 transcripts (splice variants) [Hide transcript table](#)

Show **All** entries Show/hide columns Filter

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CDS incomplete	CCDS
POLR2E-001	OTTHUMT00000458044	2831	OTTHUMP00000267716	210	Protein coding	-	-
POLR2E-013	OTTHUMT00000473950	1749	OTTHUMP00000274778	210	Protein coding	-	-
POLR2E-010	OTTHUMT00000458043	1096	OTTHUMP00000267715	210	Protein coding	-	-
POLR2E-009	OTTHUMT00000474120	932	OTTHUMP00000274865	184	Protein coding	-	-
POLR2E-014	OTTHUMT00000474115	487	OTTHUMP00000274861	134	Protein coding	5'	-
POLR2E-012	OTTHUMT00000473949	1286	OTTHUMP00000274777	204	Nonsense mediated decay	-	-
POLR2E-002	OTTHUMT00000458046	1214	OTTHUMP00000267717	68	Nonsense mediated decay	-	-
POLR2E-006	OTTHUMT00000458049	1105	OTTHUMP00000267718	51	Nonsense mediated decay	-	-
POLR2E-007	OTTHUMT00000458048	1052	OTTHUMP00000274864	83	Nonsense mediated decay	5'	-
POLR2E-011	OTTHUMT00000474116	589	No protein product	-	Processed transcript	-	-
POLR2E-003	OTTHUMT00000458045	2441	No protein product	-	Retained intron	-	-
POLR2E-008	OTTHUMT00000458047	756	No protein product	-	Retained intron	-	-
POLR2E-005	OTTHUMT00000458050	685	No protein product	-	Retained intron	-	-
POLR2E-004	OTTHUMT00000458051	440	No protein product	-	Retained intron	-	-

Variants 1-6 are unchanged, variant 7 is now protein coding, variant 8 is unchanged, variant 9 is now protein coding, variant 10 is unchanged. There has also been the addition of 4 new splice variants. Variant 11 is a non-coding transcript and variants 12 – 14 are protein coding.

Due to the complexity of the release process it can take up to 3 months for new annotation to be available in Vega. To address this, the Vega update track is run every two weeks for human and mouse, and so new annotation is publicly available much more quickly. These will be incorporated into the main Vega site when there is a new human gene release (approximately every 3 months), and then later into the Ensembl merge and the Gencode geneset.

Demo: Looking at GRC patches

We're now going to look at *ABO*, a protein-coding gene known to be involved in blood grouping. From the Ensembl homepage, search for **ABO** in **Human**.

The screenshot shows two search results for the gene *ABO* in Human. The top result is labeled "Gene on the primary assembly" and the bottom result is labeled "Gene on a patch".

ABO (Human Gene)
ENSG00000175164 9:136125788-136150617:-1
ABO blood group (transferase A, alpha 1-3-N-acetyl-3-galactosyltransferase) [Source:HGNC Symbol;Acc:79].
Variation table • Location • Regulation • Orthologues • Gene tree

ABO (Human Alternate Sequence Gene)
ENSG00000256062 HG79_PATCH:136125799-136150736:-1
ABO blood group (transferase A, alpha 1-3-N-acetylgalactosamin-3-galactosyltransferase) [Source:HGNC Symbol;Acc:79]. Not a A gene.
Variation table • Location • Regulation • Orthologues • Gene tree

This search yields two results, the gene on the primary assembly (above) and the gene on the patch (below). To find out more about this, click on the top gene [ABO \(Human Gene\)](#).

Gene: ABO ENSG00000175164

Description ABO blood group (transferase A, alpha 1-3-N-acetylgalactosaminyltransferase; transferase B, alpha 1-3-galactosyltransferase) [Source:HGNC Symbol;Acc:79]

Location [Chromosome 9: 136,125,788-136,150,617](#) reverse strand.

INSDC coordinates chromosome:GRCh37:CM000671.1:136125788:136150617:1

Transcripts This gene has 2 transcripts (splice variants) [Hide transcript table](#)

Name	Transcript ID	Length (bp)	Protein ID	Length (aa)	Biotype	CCDS
ABO-001	ENST00000453660	6341	No protein product	-	Processed transcript	-
ABO-201	ENST00000538324	937	No protein product	-	Processed transcript	-

Gene summary ⓘ

Name [ABO](#) (HGNC Symbol)

Synonyms A3GALNT, A3GALT1 [To view all Ensembl genes linked to the name [click here](#).]

Ensembl version ENSG00000175164.9

Gene type Known processed transcript

Prediction Method Annotation for this gene includes both automatic annotation from Ensembl and [Havana](#) manual curation, see [article](#).

Alternative genes This gene corresponds to the following database identifiers:
Havana gene: [OTTHUMG00000020872](#) (version 4)

[Go to Region in Detail for more tracks and navigation options \(e.g. zooming\)](#)

Gene Legend ■ processed transcript

This gene has two transcripts, both of which are non-coding. This does not fit with what we know about the gene, that it is protein coding.

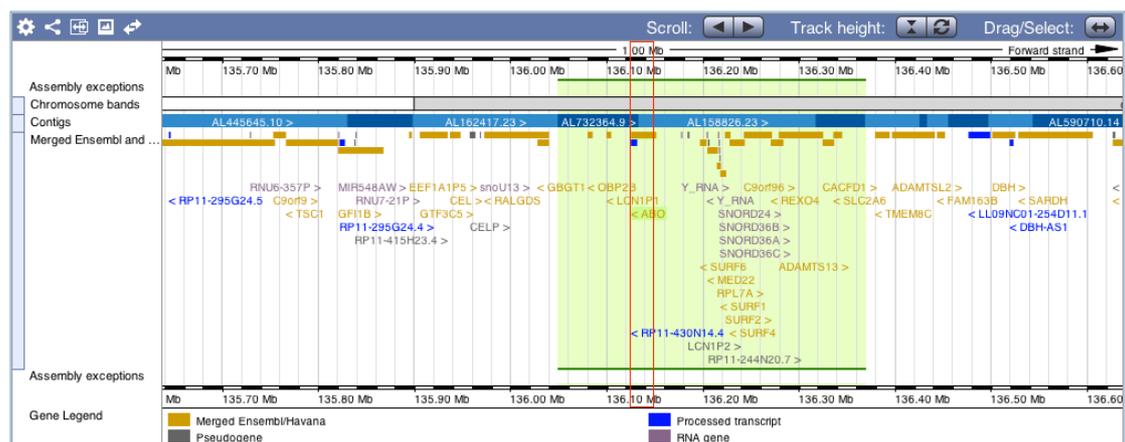
To understand what's going on, click on the Havana ID [OTTHUMG00000020872](#) to see the gene in Vega. Open the transcript table and click on the transcript [OTTHUMT00000054907](#).

Have a look at the Remarks:

Remarks	<p>ABO blood group (transferase A, alpha 1-3-N-acetylgalactosaminyltransferase; transferase B, alpha 1-3-galactosyltransferase), ABO-*O01 allele</p> <p>ABO blood group (transferase A, alpha 1-3-N-acetylgalactosaminyltransferase; transferase B, alpha 1-3-galactosyltransferase), ABO-*O02 allele</p> <p>The ABO gene in this individual produces a truncated protein without functional glycosyltransferase activity indicative of blood group O</p>
----------------	---

The gene lies between 2 BAC clones and each half of the gene represents a different allele. As a result there is no coding gene for this locus.

Go back into Ensembl and click onto the [location tab](#).



As we saw in the search results, this gene falls within a patch, shown in green. To find out why, we want to add a track. By default only a very limited number of tracks is shown (note that it says at the bottom the display that ‘There are currently 441 tracks turned off’). Additional tracks can be added using [Configure this page](#).



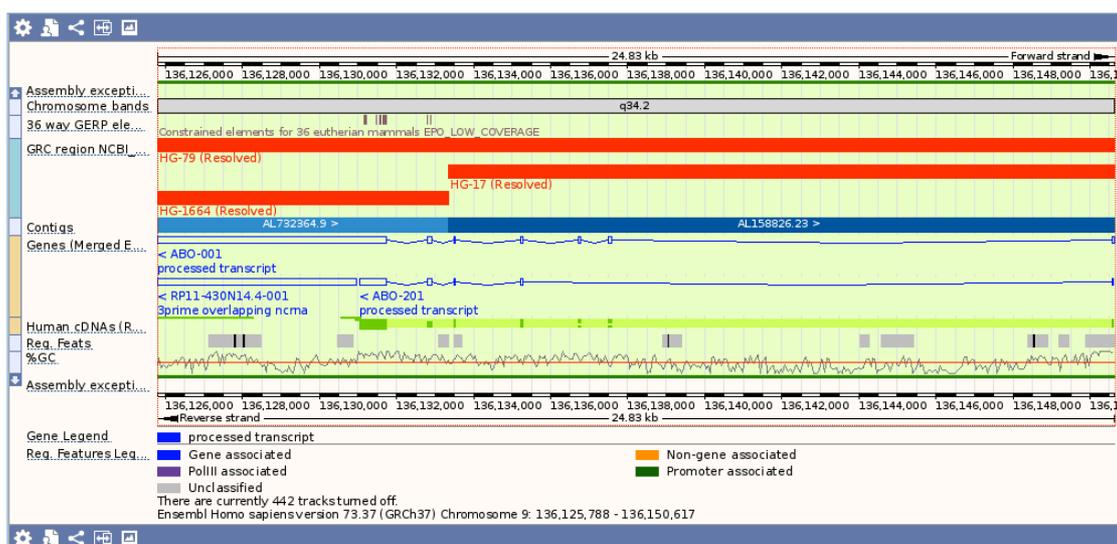
This will open a menu that allows you to change the image.

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

The screenshot shows the 'Configure Region Image' interface. At the top, there are tabs for 'Configure Region Image', 'Configure Overview Image', 'Manage Configurations', and 'Personal Data'. A search bar labeled 'Find a track' is located at the top right. On the left, a tree view lists track categories such as 'Sequence and assembly', 'Genes and transcripts', 'mRNA and protein alignments', 'Regulation', 'Comparative genomics', 'Information and decorations', 'Variation', and 'Regulation'. A callout box labeled 'Track categories' points to this list. The main area displays a list of tracks under the heading 'Active tracks', including 'Contigs', 'Sequence', 'Primary assembly mapping', 'Merged Ensembl and Havana genes (GENCODE)', 'CCDS set', 'Human cDNAs (RefSeq/ENA)', 'Reg. Feats', 'Constrained elements for 36 eutherian mammals EPO LOW COVERAGE', '%GC', 'Chromosome bands', 'Assembly exceptions', 'Scale bar', 'Ruler', 'Variation Legend', 'Structural Variation Legend', 'Alignment Difference Legend', 'Reg. Feats Legend', 'Reg. Features Legend', 'Met', and 'Dis'. Callouts include 'Configuration tabs' pointing to the top tabs, 'Search for tracks' pointing to the search bar, 'Track information' pointing to the star and info icons on the right of the track list, and 'Turn tracks on/off and change style' pointing to the track list area.

Open the [Configure this page](#) menu and select [Sequence and assembly](#) from the left. Turn on the track [GRC region NCBI_37](#) in [Labels](#), then close the menu by clicking anywhere outside the menu.

You can now see the track we added in red. These indicate problems in the primary assembly. All of the bars are labelled [\(Resolved\)](#) indicating that the problems have been fixed.



There are three red bars, one of them is labelled [HG-79 \(Resolved\)](#). Click on this red bar to open a pop-up.

HG-79 (Resolved)

Type: Variation; In this region, the ABO gene in the reference assembly reflects a haplotype of "Type O" not found in the human population.

Method: Variation; Status Resolved

Start: 136049442

End: 136317857

Strand: +

[GRC report for HG-79](#)

In this region, the ABO gene in the reference assembly reflects a haplotype of "Type O" not found in the human population.

The pop-up tells us that the reference assembly gives us a gene that doesn't exist in human populations. This is due to the fusion of two alleles. Since this is listed as Resolved, we know that it was fixed by a patch.

The thin dark green line at the top of the image indicates the position of the patch. Click on it to open a pop-up menu.

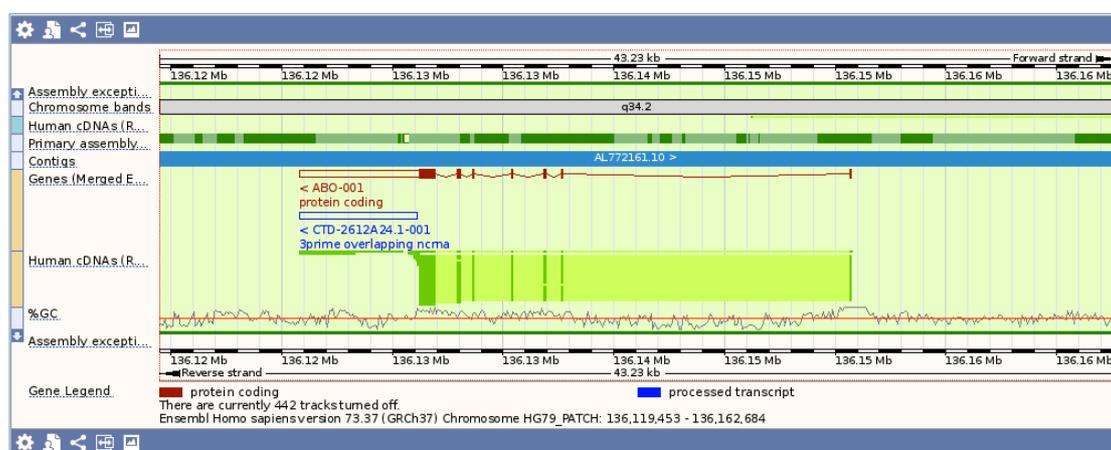
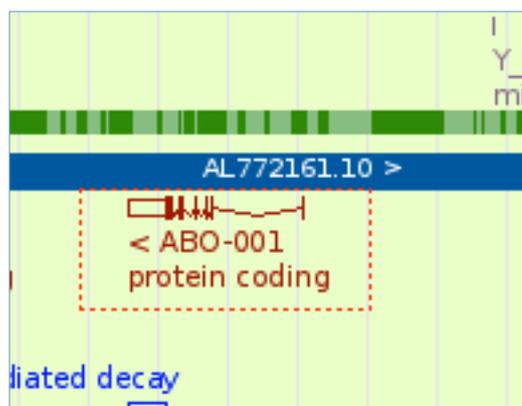
HG79_PATCH:136049442-136379605 

Synonyms: GL339450.1

[HG79_PATCH:136049442-136379605](#)

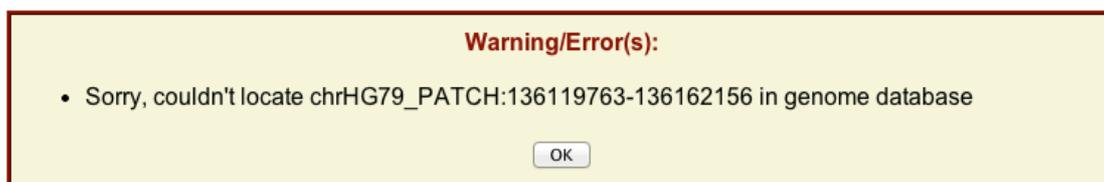
[Compare with reference](#)

Click on [HG79_PATCH:136049442-136379605](#) to go to the [Region in Detail](#) view for the patch. We have zoomed right out to view the whole patch. We can zoom back in on the ABO gene by dragging out a box around it.



We can now see the ABO gene labelled as `protein_coding`. We can also see that the whole of the gene is covered by a single clone. Using a single clone meant that all of the data came from a single genomic sequence, so the gene fusion problem is resolved.

Click on the link to [UCSC](#) at the left to open the patch in UCSC.



It is not possible to view patches in UCSC because they only have the primary assembly.

Exercises

Exercise 1 – Searching for splice variants of a gene

Search for the *BRAF* gene in human in Ensembl, which is an important gene in cancer. How many splice variants are there and what are their biotypes? How many merged transcripts are there? How many have a CCDS?

Exercise 2 – Searching for genes on haplotypes

(a) Search for the *HERC2* gene in human in Ensembl. How many genes are there called *HERC2*? Why is this? Take a look at both genes. Is there a difference between the number of transcripts?

(b) What is the name of the haplotype that the alternate sequence falls on? Go to the location view. Can you compare the haplotype with the primary assembly? What differences can you see?

Answers

Exercise 1 – Searching for splice variants of a gene

Go to ensembl.org.

Select [human](#) from the drop down list and type in [braf](#) then hit return.

Click [BRAF \(Human Gene\)](#).

If the transcript table is hidden, click on [Show transcript table](#).

There are five transcripts of *BRAF*. Two are protein coding, two are subject to nonsense-mediated decay and one has a retained intron. One has a CCDS.

Exercise 2 – Searching for genes on haplotypes

(a) Go to ensembl.org.

Select [human](#) from the drop down list and type in [herc2](#) then hit return.

Narrow down to genes only by selecting [Gene](#) from the left hand list.

The search returns 11 genes in total, of which two are called *HERC2*.

One of these is on a haplotype, where the genome has different sets of variants between individuals, and the other is on the primary sequence.

Open the two genes, [HERC2 \(Human Gene\)](#) and [HERC2 \(Human Alternate sequence Gene\)](#), in different tabs.

The primary assembly gene has twelve transcripts, whilst the haplotype gene has eight.

- (b) The alternate sequence gene is described as being on HRSCH15_1_CTG4.

Click on the [Location](#) tab in the top bar from either gene. You will see the haplotype represented as a red highlighted region. Click on the dark red line at the top or bottom of the region, then select [Compare with patch](#) or [Compare with reference](#), depending on which one you're looking at.

There is a short intergenic insertion in the haplotype compared to the primary assembly.