# Module 3: Genome Browsing

**Aims**

- Briefly present the main web-based genome browsers.

- Using Ensembl, demonstrate some of the features and applications of genome browsers.

- Introduce the BioMart data retrieval system.

- Create files with your own data to upload to a genome browser.

**Introduction**

Web-based genome browsers have been developed to make it easier to access comprehensive information about regions of the human genome and about the whole human gene set.  They help you to:

- Explore what is in a chromosomal region

- See features in and around a specific gene

- Search & retrieve data across the whole genome

- Investigate genome organisation

- Compare to other genomes

Browsers display the location and structure of known genes and predicted novel genes along with information about the mRNA transcripts and may also include information about protein products.  Information about genes is integrated with information about other genomic features (e.g. variation data, markers, repeated sequences, regions homologous to other species) and displayed alongside the genomic sequence assembly.  Protein, mRNA and EST entries from various sequence databases may also be shown aligned to the chromosomes.

In addition to providing annotation across the whole genome, browsers provide other resources. The browsers differ in what is provided and how it is presented. Resources that can be found include:

- **Links** to other databases and resources
- **Text Searching**
- **BLAST** and other sequence similarity searching
- **Download** of genomic sequence, gene information and other data
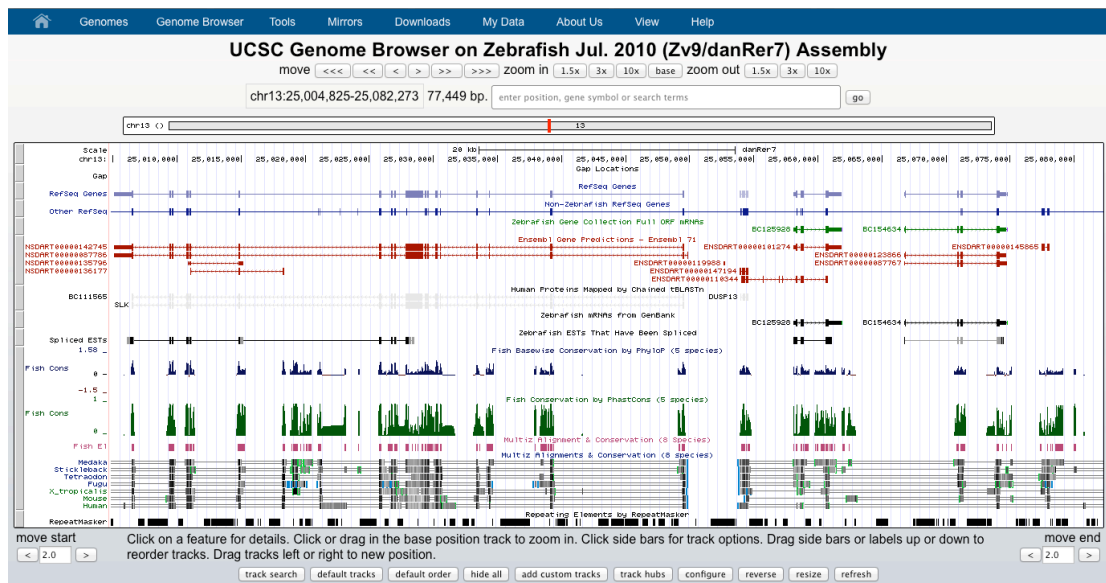- **Data mining** facilities

*Browsers (and some of their strengths)*

- **NCBI Map Viewer** – maintained by NCBI
  http://www.ncbi.nlm.nih.gov/mapview/
- **UCSC Genome Browser** – maintained by UCSC
  http://genome.ucsc.edu/cgi-bin/hgGateway
- **Ensembl** – maintained by EBI / Sanger Institute
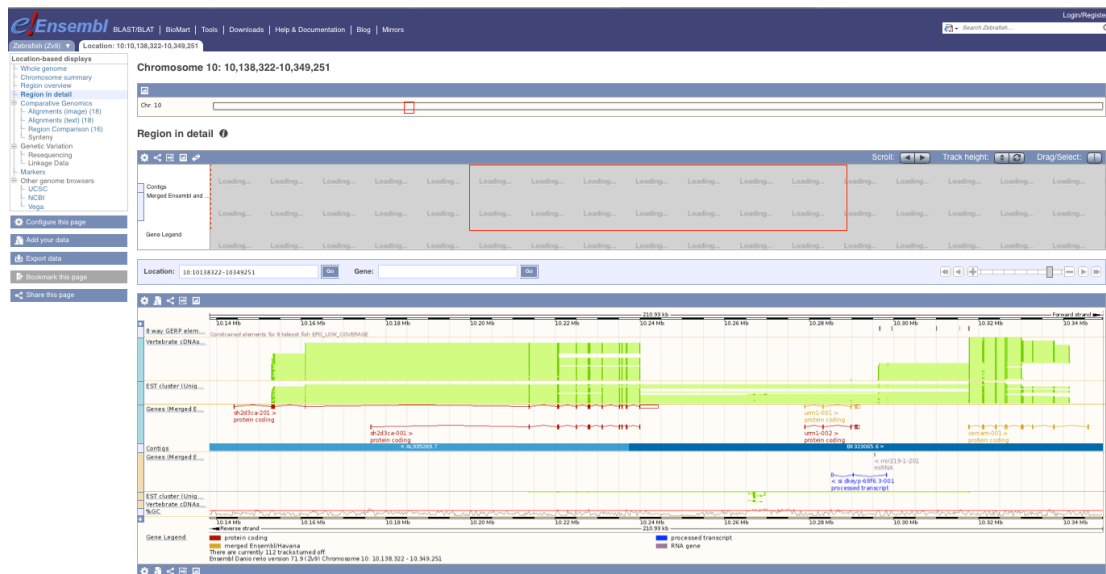  http://www.ensembl.org



**NCBI Map Viewer**

- Good integration with other NCBI resources

## UCSC Genome Browser

- Straightforward feature display
- Old assemblies available
- Wide range of tracks supplied by other groups
- Trackhub support



## Ensembl

- Well-supported gene set with evidence
- Range of different views
- Archive available
- Trackhub support

While browsers can be very useful tools, they do not provide the definitive answer to every question!  Remember, new data and updates make genome browsing a fluid, changing, and improving, process.

## Data retrieval and data mining

Genomic annotation data, due to its complexity and volume, does not lend itself to easy access.  Presenting it on a web site is important, but so is providing simple but flexible ways to select and retrieve specific sets of data.  NCBI has the Entrez query system and UCSC has its Table Browser.

In Ensembl, BioMart facilitates rapid retrieval of richly annotated gene lists, sequences, and variations, among other annotation, integrated with third party data and applications.  Genes can be selected by chromosome region, protein domains, associated external identifiers or SNP properties, and these filters can be combined to group and refine biological data, including cross-species analyses, disease links, sequence variations and expression patterns.

BioMart is built upon a query-optimised relational database schema allowing quick and efficient access to voluminous data through a user-friendly, interactive web interface.  After selecting the biological object and the species, the results can be refined using a set of pre-defined filters.  After each navigation event, the user is provided with immediate feedback on the number of matches found.  Output can consist of annotated gene lists, gene structures, SNP details or various kinds of sequence sets.  Output can be in HTML, text, Microsoft Excel and compressed formats.

*Further reading*

Ensembl Help and Documentation http://www.ensembl.org/info/index.html

Cunningham F. *et al*.
**Ensembl 2015**
*Nucleic Acids Res*. 2013 Jan;43(Database issue):D662-9

Spudich GM, Fernandez-Suarez XM
**Touring Ensembl: a practical guide to genome browsing.**
*BMC Genomics*, 2010 May 11;11:295

Meyer LR, et al.
**The UCSC Genome Browser database: extensions and updates 2013.**
*Nucleic Acids Res*. 2012 Nov 15.

Karolchik, D *et al.*
**The UCSC Genome Browser Database.**
*Nucl. Acids Res*. 2003 **31**, 51-54

Dombrowski, S M and Maglott, D.
**Using the Map Viewer to Explore Genomes**
in The NCBI Handbook
http://www.ncbi.nlm.nih.gov/books/bookres.fcgi/handbook/ch20d1.pdf

## WALKING THROUGH THE WEBSITE

The instructor will guide you through the release 80 Ensembl website using the **thyrotrophic embryonic factor a (tefa)** gene. The following points will be addressed:

- **The Gene Summary tab and gene-related links:**
    - Are there splice variants?
    - Can I view the genomic sequence with variations?
    - Find orthologues and paralogues, show alignments with other genomes

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

- **The Location tab and related links:**
    - How do I zoom in and change the gene focus.
    - Un-stacking a track (e.g. human cDNAs)
    - Adding a track (i.e. variations)

- **Exporting a sequence and running BLAT/BLAST**

Start by going to **www.ensembl.org**

Click on 'Zebrafish', or the picture circled above, which brings us to the species index page.



Type 'tefa' into the search bar circled above and click the 'Go' button.

The search will return everything that matches tefa, the first result is the *tefa* gene. You can navigate to the 'Gene' page, to a 'Location view' page (coordinates hyperlink) or to any other View' type as listed below the search result.



Click on 'tefa (Zebrafish Gene)' and the following 'Gene' tab will open:

Let's walk through some of the links in the left hand navigation column.



How can we view the genomic sequence? Click **'Sequence'** at the left

'**Configure this page**' in the left hand menu allows you to make changes to the display, e.g. add coordinates and make variation visible.



Once you have selected changes (in this example, we display variations and show chromosome coordinates) click the **tick** at the top right.



Now variations in the sequence are highlighted. Coordinates have been added.

Now let's click on **'Gene tree'**, which will display the current gene in the context of a phylogenetic tree of orthologous and paralogous genes.



Use the mouse over and 'expand sub-tree' to get to the view displayed above. Note that there are two *tef* genes in fish species. In zebrafish these are annotated as *tefa* and *tefb*.

Click on **'Variation image'** to display genetic variation mapped onto all transcripts of a gene.

Click any variation, then **'Variation properties'** to learn more about it.   A fourth tab will open:



To find **orthologs** of the *tefa* gene in other species, return to the 'Gene' tab and select 'Orthologues'. You will find that the zebrafish 'tefa' gene has one human ortholog, TEF.



The relationship to the ortholog is '1-to-many', meaning that several zebrafish paralogs share the close relationship with one human TEF gene (compare answer to **Exercise 1**). You can check this by navigating to the human gene and looking at its zebrafish orthologs:



In order to find out how the TEF and *tefa* loci compare, go back to the *tefa* location tab and click '**Comparative Genomics**' in the left hand menu.

**'Region Comparisons'** provides pre-computed alignments between species. Select 'Human' to see the below.

Now, we would like to work with a **transcript** of this gene. Return to the *tefa* 'Gene' page, select 'Show transcript table' and and click tefa-001. This will lead to the transcript summary display.



Note that this is a 'merged' or 'golden' transcript, i.e. the automated (Ensembl) and manual (Vega, vega.sanger.ac.uk) annotation are identical.

Again, the left hand navigation column provides several options for this particular transcript.

Choose the **'Exons'** option first, which displays exon sequences in full and introns in a configurable context. Use the **'Configure this page'** link to change the display (for example, show more flanking sequence, show full introns).



Green: Flanking sequence

Purple: UTR

Grey: Introns

Blue: Coding sequence

Next, follow the **'Supporting Evidence'** link, which shows which biological evidence has been used for the annotation of this transcript.



'Golden' boxes: exons in the transcript

Alignments of cDNA and protein to the Ensembl exons.

Other transcript-specific displays include the cDNA sequence, general identifiers and gene ontology terms from the GO consortium (www.geneontology.org).

Ensembl **'Location'** displays are also highly configurable. To enter the configuration dialogue, use the **'Configure this page'** link. As an exercise, add **all variations** to the **'Region in detail'** display and view the **zebrafish cDNA** track in **'normal'** expanded form.

Rather than reconfiguring your preferred view with every new visit to the site, you can preserve your configuration by registering and logging into **your account** (upper right corner).

After investigating the **'Location'** display, we would like to export genomic sequence. Click the **'Export data'** option and select the **'FASTA'** sequence format.

```
>12 dna:chromosome chromosome:GRCz10:12:19022959:19031303:1
TAAAACAAATCAGAGTTTTTTAATAATATGTACAACATGAAACAGCCAAGGCATTCAACA
TTTACAAGTCAAAAACATCTAAAAGACAATCCAAAAAGAAAAACTTTTCAAAACGGACAA
TTAGAAAAAGACTTTTTGTACAGGAACAAATTACATATTCAACCAAGATTCATGGAGTTA
GTCGCTGTTCCTCTCAAATCAGCCTGCATTGGTTTGCTTTCTCGTCTTACACTTAGGAAG
CTTGAAGAGAAACGTCAAACTGAAGGTGCGCTTGTAAAGAGATTCTAGTACAGACCGCAA
ATATGCACATCGTAAAACCTGATTTAAGTGAAGTCATTGTGTACTAGCCACAAAAACTAA
GCTTCAGAGTTTCAATCAATTAGCTTTGGTCAAGCTCACATTACAGCAACTGCCATTCTG
AAAAAAAAACACATGAAAAAAATACAAACAGGCCTCACAATGAGTACTGCAACATTAGAC
TGCTAGCCTACTTCACAACAATACACAGAGCAAAATACACGACATATAGTGAATCTTAAG
AACTGACCCTTTATCTTCTTTCCCTGATGATAAAAAGGACAGAGACTAAATATGTGCCT
ACAATAAAAAAAAAAGCAGGAAGGTGGCAAGCAGAGGTTTGAATTCAGTTTCTGATATGT
AACACACTAGTGCTTTTATAAACAGTCACGCTGCTCAAGACACTTATGGATCACTGCAAA
ACCACTTTTGTTAAACATTTGATGGTTATCTGAAAAGATAAATGATGCATACGAAAACGA
```

Select the header and a few lines of sequence and then follow the 'BLAST/BLAT' link in the blue header bar. Paste the sequence into the appropriate box and select 'BLAT' as the search algorithm and 'Danio_rerio' as species. Finally, click 'Run'.



Wait a bit until:



Note that the Results Table can be configured/reordered to display the desired data. Follow links to the Location View ('Genomic bp')

Export or share the image using the links at the bottom. The 'share' option will preserve your configuration and might be more helpful than sharing a 'picture'.

## EXERCISES and ANSWERS

Note: The answers to these exercises correspond to version 80 of Ensembl featuring the GRCz10 assembly.  If you use a different version and your answer doesn't correspond with the given answer, please consult the instructors. Note that certain versions are preserved on the Ensembl Archive site.

### Exercise 1 – Exploring a gene

(a) Search for the zebrafish *tead1a* gene. On which chromosome is this gene located? How many transcripts (splice variants) has Ensembl annotated for it? Are these transcribed from the forward or from the reverse strand of the genome assembly?

(b) What is the longest transcript? How long is the protein it encodes? How many exons does it have? Are any of the exons completely or partially untranslated? How do the transcripts differ?

(c) Have a look at the General identifiers for one of the *tead1a* transcripts. Click on some of the links. What is the function of *tead1a*?

(d) Which PFAM domains do the proteins encoded by *tead1a* contain?

(e) Is there a human ortholog predicted for the zebrafish *tead1a* gene? What 'type' does it have? Why?

(f) If you have yourself a gene of interest, explore what information Ensembl displays about it!

Advanced questions drawing in other modules:

 (1) What does ZFIN say about *tead1a*?

(2) What are the paralogs of *tead1a*?

*Answers*

(a)

✌ Go to http://www.ensembl.org.
✌ Under 'Search' select 'Zebrafish' and type 'tead1a'.
✌ Click [Go].
✌ On the page with search results follow 'gene' -> 'Zebrafish' and click the gene ID of *tead1a*.

The zebrafish *tead1a* gene is located on linkage group 25. Ensembl has 3 transcripts annotated for this gene. The transcripts are transcribed from the forward strand of the genome assembly.

(b)

✌ Have a look at the transcript table at the top of the page.

The longest transcript is ENSDART00000125925. The length of this transcript is 1744 base pairs and the length of the encoded protein 422 amino acids.

✌ Click on 'ENSDART00000125925'.
✌ Click on 'Exons' in the side menu.

ENSDART00000125925 has 13 exons, of which the first two are untranslated and the third and the last one are partially translated.

✌ Click on the 'Location' tab and zoom in on different areas of the transcripts
✌ Click on the transcripts and in in the pop-up menu check the 'Analysis' entry.

*tead1a*-001 is longer than the other transcripts. *tead1a*-001 was annotated both by the automated Ensembl pipeline as well by manual annotation (Havana), resulting in the same structure, and was therefore merged. *tead1a*-002 was manually annotated, and found to be protein-coding. *tead1a*-003 was manually annotated and found to contain a retained intron; it is non-coding.

(c)

✌ Click on 'General identifiers' in the side menu of a Transcript tab or 'External References' from a Gene tab.
✌ Explore some of the links (a good place to start is 'ZFIN').

*tead1a* encodes an DNA-binding transcription factor involved in the hippo signalling cascade.

(d)

🖱 Select a protein.

E.g. the *tead1a*-001 protein contains a TEA/ATTS domain.

(e)

🖱 Click on the 'Gene: *tead1a*' tab.
🖱 Click on 'Orthologues' in the side menu.

There is one human ortholog predicted for zebrafish *tead1a*, TEAD1 (ENSG000000187079). It has the type '1-to-many'.

🖱 Explore the 'Help & Documentation' pages, a definition of homology types can be found at
http://www.ensembl.org/info/genome/compara/homology_method.html

Human TEAD1 is the ortholog of the zebrafish genes *tead1a* and *tead1b*.

_____

## Exercise 2 – Exploring a region

(a) Go to the region from bp 33100000 to 33350000 on zebrafish chromosome 13. How many contigs make up this portion of the assembly (contigs are contiguous stretches of DNA sequence that have been assembled solely based on direct sequencing information, in the zebrafish assembly there are finished clones and whole genome shotgun contigs)?

(b) Make the tilepath clones (i.e. the BAC clones that were sequenced to generate the sequence for the human genome assembly) visible, what are the clone names in this region? Note that these clones are not shown by default! Which clone library does the clone containing the *btbd6a* gene come from?

(c) Zoom in on the *btbd6a* transcript, including a bit of flanking sequence on both sides. Which markers are located close by? Do the markers appear anywhere else in the genome?

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

(e) Is this region being worked on by the Genome Reference Consortium?

(f) If you have yourself a genomic region of interest, explore what information Ensembl displays about it!

_____

_____

*Answer*

(a)

ᐧᐤ Go to the Ensembl homepage.
ᐧᐤ Under 'Search Ensembl' type 'zebrafish 13: 33100000-33350000'.
ᐧᐤ Click [Go].

This genomic region is made up of 8 contigs, indicated by the alternatingly light and dark blue coloured bars in the 'Contigs' track.

(b)

ᐧᐤ Click on 'Configure this page' in the side menu.
ᐧᐤ Click on the individual contigs to see more details.
ᐧᐤ Follow the EMBL link to the submission record to find out about the name.

The tilepath clones in this region are DKEY-71P21, CH1073-380H14, CH1073-224M6 and CH1073-127N13. There are also 4 whole genome shotgun contigs, their accessions start with CABZ.

ᐧᐤ Click on 'CU855940.5', then follow the EMBL link
ᐧᐤ Read the last lines of the comments.

CU855940.5  (CH1073-380H14) is from the CHORI-1073 Zebrafish double haploid fosmid library.

(c)

ᐧᐤ Switch on the 'Markers' track under 'Sequence and Assembly'.
ᐧᐤ Draw a box around the  transcript.
ᐧᐤ Click on 'Jump to region' in the pop-up menu.
ᐧᐤ Click on the markers and 'Marker info'

Gene *btbd6a* is e.g. close to the fc21e08.y1 marker. This marker is only placed in this location.

 (d)

ᐧᐤ Click on 'Export data' in the side menu.
ᐧᐤ Click on [Next>].
ᐧᐤ Click on 'HTML'.

21

Note that the sequence has a header that provides information about the genome assembly (GRCz10), the nature of the sequence (dna), the coordinate system (chromosome), the coordinate system descriptor (13), the start and end coordinates (e.g. 33249320:33297062) and the strand (1):

>13 dna:chromosome chromosome:GRCz10:13:33249320:33297062:1

(e)

ᗰ Go to 'Configure this page' and switch on 'Sequence and assembly' -> 'GRC alignments' and switch on 'Genome curation'

This reveals that there is indeed an issue with the assembly here, registered as ZG-6933. The report states that a gap exists between clones components BX284673.9 and CU855940.5. The GRC is currently awaiting sequence data to close the gap that is filled with WGS contigs.

## Data mining

On top of visualising genome data, all browsers offer data to be extracted and stored. This ranges from a simple download of sequence or features for a certain region to genome-wide pre-prepared data collections. However, if you are only interested in data that passes certain filters it becomes trickier. Sometimes, this can be extracted from the offered downloads, but this might require advanced bioinformatics skills and additional data.

UCSC and Ensembl offer a service to filter and extract data according to your needs. UCSC provides the **Table Browser** (https://genome.ucsc.edu/goldenPath/help/hgTablesHelp.html) and Ensembl provides **Biomart** which we explain in detail below.

## Mining data using BioMart - worked example

- Find all protein-coding zebrafish genes on linkage group 1 that have a human orthologue.

- Display the Ensembl IDs of the zebrafish and human genes plus the chromosomal location of the human gene.

- Download the sequence of all available 5' UTRs of these genes.

Note that the below example was created on Ensembl version 80. Since the gene set gets adapted to the ongoing manual gene annotation with every other release, the results might differ with a different release, even on the same genome assembly.

**STEP 1:** Click on 'BioMart' in the top header bar of the Ensembl home page.



**STEP 2:** Choose 'Ensembl Genes 80' as the primary database.

**STEP 3:** Choose '*Danio rerio GRCz10*' as the species of interest.

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

**STEP 5:**
Select 'Chromosome 1'

**STEP 6:**
Expand the 'GENE' panel.

**STEP 7:**
Expand the 'MULTI SPECIES COMPARISON' panel.

**STEP 8:**
Limit to genes of type **'protein coding'**

**STEP 9:**
Limit to **'Orthologous Human Genes Only'**



**STEP 10:**
The filters have determined our gene set.
Click **'Count'** to see how many genes have passed these filters.

**The 'Count' results show 797 zebrafish genes out of 31,953 total genes passed the filters.**

**STEP 11:**
Click on **'Attributes'** to select output options (i.e. what we would like to know about our gene set).

**STEP 12:**
Expand the 'GENE' panel. Deselect the **Transcript ID**

**STEP 13:**
Expand the 'Homologs' panel and select 'Orthologs'

**STEP 14:**
Select 'Human Ensembl Gene ID' and 'Human Chromosome' from 'Human Orthologs'

**STEP 15:**
Click 'RESULTS' at the top to preview the output.

Note the summary of selected options.
The order of attributes determines the order of columns in the result table.

And here you have the first 10 results; you can change the number of displayed results in the drop down menu. Expanded to 'all' this gives you a nice overview of possible syntenic regions in the two genomes.

In order to obtain all 5'UTRs of these genes, go back to the '**Attributes**'.



Click '**Results**' and you will get the required list. Note that not all genes have 5' UTRs annotated.

_____

## *EXERCISES and ANSWERS*

Note: The answers to these exercises correspond to version 80 of Ensembl.
_____

**Exercise 1**

Generate a list of all zebrafish protein coding genes on chr1 with a ZFIN ID
that have more than one splice variants and that are causing the caudal fin to
be absent when mutated. Download the peptide sequences and make sure
the header states the Ensembl ID, a description, the associated gene name
and the associated gene DB.
_____


*Answer*

⍾ Go to the Ensembl homepage.
⍾ Click the BioMart link on the toolbar.

Start with all the zebrafish Ensembl genes:

⍾ Choose the 'Ensembl 80' database.
⍾ Choose the 'Danio rerio genes GRCz10' dataset.

Now filter for the genes on chromosome 1:

⍾ Click on 'Filters' in the left panel.
⍾ Expand the 'REGION' section by clicking on the + box.
⍾ Select 'Chromosome – 1'. Make sure the check box in front of the filter
is ticked, otherwise the filter won't work.
⍾ Click the [Count] button on the toolbar.

This should give you 1,386 / 31,953 Genes. Now for genes with a ZFIN ID

⍾ Click 'Limit to genes...' and choose 'with ZFIN IDs'
⍾ Click the [Count] button on the toolbar.


1,186 genes have ZFIN IDs.

Now filter further for genes that are protein coding:

⍾ Expand the 'GENE' section by clicking on the + box.
⍾ Select 'Gene type – protein_coding'.
⍾ Click the [Count] button on the toolbar.

This should give you 1,081 / 31,953 Genes.

Now only select those genes with at least 2 alternative splice variants.

🖰 Expand the 'GENE' section again by clicking on the + box.
🖰 Select "Transcript count >=' and enter '2'
🖰 Click the [Count] button on the toolbar.

661 genes left. Let's see what disrupts the caudal fin development:

🖰 Expand the 'Phenotype' filter.
🖰 Select 'caudal fin absent' under 'Phenotype'.
🖰 Select 'ZFIN' under 'Phenotype source'.

One gene left!

Now download the cDNA sequences with the Ensembl gene and transcript IDs, the associated gene name, gene DB and a description.

🖰 Click on 'Attributes' in the left panel.
🖰 Select the 'Sequences' attributes page.
🖰 Select 'cDNA'.
🖰 Expand the 'Header Information' section and select 'Ensembl Gene ID', 'Description', 'Associated Gene Name' and 'Associated Gene Source' plus 'Ensembl Transcript ID'
🖰 Click the [Results] button on the toolbar.

```
>ENSDARG00000031894|ENSDART00000047876|lymphoid enhancer-binding factor 1
[Source:ZFIN;Acc:ZDB-GENE-990714-26]|lef1|ZFIN
GGAGCACGACACAGACCTGATGCACATGAAACCTCAGCACGAGCAGAGAAAGGAGCAGGA
GCCCAAAAGACCTCACATCAAGAAACCTCTAAACGCTTTCATGCTGTATATGAAAGAGAT
GCGCGCCAATGTGGTGGCCGAATGCACGCTGAAGGAGAGCGCCGCTATCAATCAGATCCT
CGGCCGGAGGTGGCATGCTTTATCTCGGGAAGAGCAAGCTAAGTATTACGAATTAGCCCG
CAAGGAACGGCAGCTCCATATGCAGCTTTACCCAGGATGGTCTGCCAGAGACAATTATGG
AAAGAAAAAAAGCGGAAGAGGGAAAAGATCCAGGAACCTGCTTCAGATGGAAATGGCTT
TTTCTTTTATGGAACACAAAAGGTACAGGCCAGAGAATGAAAACGGCGTACATCTGAACA ATGGTAAGAG
>ENSDARG00000031894|ENSDART00000132405|lymphoid enhancer-binding factor 1
[Source:ZFIN;Acc:ZDB-GENE-990714-26]|lef1|ZFIN
GTAGTCAGTCAGAGATCAGGGGGAGGAGTACAGCACTACACTCTCTCCAGCCCAACATTA
CTCTCAGTCTCTGCTGAGCTCATTTCTGAAGAGGGACACCTTTTTTTACCCAACAAACCAA
ACGGGAATGACACACACCATCTGAACTCCAACATTTCTTTTTTTGTTGTTGTTGCTTTTA
TTTTGAAACAAGTGAAACTGTCCTTTTCTGAACTTTAAGTTCCAACTTTTTCCTTCCACC
AAAGGATTCGTATTTTAACTTTTTCCCCAAACCCGCTATTTTTCTTCCTCTGGATTCCCG
AGAGTTTTTCCACCGGACGCGCGCGCTCTTGTTACCGTAAACCAAACACACTCACGCGCG
CGTGTTCGGAGTGCGCGAGCTGACCAGAAACAAAACAACTATACGGGGGGTTTAATTTCA
ATTGCACGCGTTTGCGTCCCTGGCGTTTGTAGGGTGAGGAGGACTTTCATTCACCCGAGA
```

…

**Exercise 2**

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

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

Generate a list that shows to which Ensembl Gene IDs and to which HGNC symbols these RefSeq IDs correspond. Which of these genes have a zebrafish ortholog?
_____

*Answer*

☝ Click [New].
☝ Choose the 'Ensembl 80' database.
☝ Choose the 'Homo sapiens genes (GRCh38.p2)' dataset.

☝ Click on 'Filters' in the left panel.
☝ Expand the 'GENE' section by clicking on the + box.
☝ Select 'Input external references ID list - Refseq protein ID(s)'
☝ Enter the list of IDs in the text box below (either comma separated or as a list).

☝ Click on 'Attributes' in the left panel.
☝ Select the 'Features' attributes page.
☝ Expand the 'GENE' section by clicking on the + box.
☝ Deselect 'Ensembl Transcript ID'.
☝ Expand the 'External' section by clicking on the + box.
☝ Select 'HGNC symbol' and 'RefSeq Protein ID'.

☝ Click the [Results] button on the toolbar.
☝ Select 'View All rows as HTML' or export all results to a file. Tick the box 'Unique results only'.

Note: BioMart is 'transcript-centric', which means that it will give a separate row of output for each transcript of a gene, even if you don't include the Ensembl Transcript ID in your output. When you don't want this, use the 'Unique results only' option.

_____

Your results should show 11 genes, most of them Caspase (CASP) genes. Several RefSeq IDs map to the same Ensembl Gene ID and HGNC symbol.

Now narrow down to genes with zebrafish orthologs.

🖰 Click on 'Filters' in the left panel.
🖰 Expand the 'MULTI SPECIES COMPARISONS' section by clicking on the + box.
🖰 Select 'Homolog filters' and select 'Orthologous Zebrafish Genes - Only'
🖰 Click the [Count] button on the toolbar.

You will be left with 8 genes.

_____


**Exercise 3**

Generate a list of all zebrafish genes on chr 1 that have an human ortholog on human chr 13. Display the gene names, are they the same? Note: This requires you to select an additional data set.

🖰 Choose database 'Ensembl 80' and dataset 'Danio rerio genes (GRCz10).
🖰 Narrow down by filtering for 'REGION' 'Chromosome - 1' and 'MULTI SPECIES COMPARISONS' selecting 'Homolog Filters' 'Orthologuos Human genes -Only'
🖰 Click on "Attributes', then 'Features' , deselect 'Ensembl Transcript ID', select 'Associated Gene Name'

🖰 Click on 'Dataset' (bottom left) and select '[Ensembl 80] Homo sapiens genes (GRCh38.p2)'
🖰 Narrow down by filtering for 'REGION' 'Chromosome - 13' and 'MULTI SPECIES COMPARISONS' selecting 'Homolog Filters' 'Orthologuos Zebrafish genes -Only'
🖰 Click on "Attributes', then 'Features' , deselect 'Ensembl Transcript ID', select 'Associated Gene Name'

You will end up with a list where quite a lot of names are identical in zebrafish and human. Note that the uppercase zebrafish gene names are projected from the most likely human ortholog whereas lowercase names are given by ZFIN.

_____

**Exercise 4**

Design your own query!
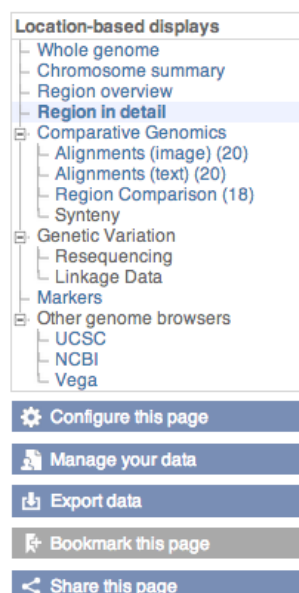
# Make your own data visible with BED files

BED files have a very simple format (tab delimited, only chr, start and stop required) and can be quickly created to provide a list of features you want to make visible in a genome browser like e.g. Ensembl or the UCSC browser. The BED format also allows for optional extension, which can make these data collections very powerful tools, even providing whole trackhubs. The format is described in detail at

http://genome.ucsc.edu/FAQ/FAQformat.html#format1

Here is an example for a bed file content (the header is optional):

```
track name=random description="random collection of features for
workshop"
chr10 0       50000 gene1
chr10 69999 100000        gene2
chr10 109999        120000        gene2
```

Go to the Ensembl browser, bring up a zebrafish location view and choose 'Add your data' and select 'Data format – BED'

Paste the bed file data and upload.



Navigate to the start of chromosome 10 and you should see the two features added.



Here is a more advanced example:

```
track name=random description="random collection of features for
workshop" useScore=1 itemRgb="On"
chr10 0      50000 feature1    1     +      0      50000 255,255,0    3
       12000,5000,3000    0,22999,46999
chr10 59999 65000 feature2     0.5   +      59999 65000 0,255,0
chr10 69999 100000          feature3    1      -      69999 100000
       0,0,255
```

resulting in



**Exercises:**

1. Play around with the data. What do you need to do to display separate items and what to display exon/intron like structures?

2. Create your own data collection in a bed file and make it visible in Ensembl.

3. Try to query biomart for a certain range of features and adapt the output to bed format. A text editor should help with e.g. the 'chr' prefix. Make the features visible in Ensembl.

4. Where exactly does gene1 in the first example start and end? Why?

**Answers:**

1. Column 4 accepts a name for the feature. If you give the same name for more than one feature, those with the same name will be drawn in exon/intron style.

4. The bed format is requiring a 0-based start and a 1-based end, so using the range of 0-50,000 in our example bed file specifies to use the bp starting after 0 (i.e. bp 1), ending at 50,000. A good explanation of coordinate systems can be found here: https://www.biostars.org/p/84686/