---

# Module 2 - Maps and Genome Sequence

<u>**Aims**</u>     ## iv – How do I find a zebrafish gene?

- Introduce the different search facilities in Vega/Ensembl
- Discuss strategies for locating genes in the zebrafish genome
- Present other resources like the trace repository

## Introduction

The genome sequence would not be of much use without annotation. The interfaces of the Vega/Ensembl browsers are designed to efficiently present users with relevant information, but the interpretation of much of the data is still in the user's domain. Searching for a region of interest can be a difficult task in its own right.

Every gene, transcript, exon and translation in Vega and Ensembl have an identifier, for example, ENSDARG00000021389 in Ensembl or OTTDARG00000005397 in Vega. These identifiers can be used as external references. In every new Ensembl release the set of identifiers from an old version are carried over wherever possible. In some cases, as the assembly is not finished yet, the identifiers cannot be mapped and they might vanish when moving to the latest release. Since October 2004 the old Ensembl releases are available from the Ensembl Archive site so old data can be checked and compared to the latest assembly. Links to the Ensembl Archive site can be found in the left-hand side menu bar in the Ensembl pages. In Vega the identifiers remain stable since genes are linked to finished clones.

## TextView

The simplest way of searching for a gene (or indeed any term or accession) is using the text-based searches. The Vega/Ensembl pages have text boxes where the user can enter a keyword to perform a search over a collection of pre-indexed items. If you know the name of a gene or a keyword that might be present in its description then you can use it in this kind of search. Try searching with the name jag2. The search result is displayed in a TextView page.

_____

**Ensembl text search**



Target:
Danio rerio

Query:
jag2

Stable identifier

Features

**1 matches in the *Danio rerio* Gene index [first 5 matches shown]:**

**1. Ensembl Gene:** ENSDARG00000021389
Ensembl gene ENSDARG00000021389 has 2 transcripts: ENSDART00000024922, ENSDART00000049586 and associated peptides: ENSDARP00000010799 ENSDARP00000049585
jagged 2 isoform 1 [Source:RefSeq_peptide;Acc:NP_571937]
The gene has the following external identifiers mapped to it:
Affymx Microarray Zebrafish: Dr.8287.1.S1_a_at
EMBL: AF090432, BX004766, AF229449, AF229450
EntrezGene: 140422
GO: GO:0016020, GO:0001889, GO:0005509, GO:0007154
IPI: IPI00500671, IPI00500671.1, IPI00501275, IPI00496898, IPI00496898.1, IPI00501275.2
Predicted UniProt/TrEMBL: Q90Y55_BRARE, Q90Y56_BRARE, Q5TZK8, Q90Y55, Q5TZK8_BRARE, Q90Y56, Q9YHU2_BRARE, Q9YHU2, Q5TZK7_BRARE, Q5TZK7
Protein ID: AAL08214.1, CAH69087, AAL08215, CAH69088, AAL08214, AAC98354, CAH69... C98354.1
RefSeq DNA: NM_131862.1, NM_131665.1, NM_131862, NM_131665
RefSeq peptide: NP_571937, NP_571937.1, NP_571740, NP_571740.1
UniGene: Dr.8287
ZFIN ID: jag2, ZDB-GENE-011128-3
Geneview: http://www.ensembl.org/Danio_rerio/geneview?gene=ENSDARG00000021389
ContigView: http://www.ensembl.org/Danio_rerio/contigview?gene=ENSDARG00000021389

A **TextView** page summarises the result of a text-based search. If the query appears under different indices then the TextView page organises the results in categories. For example the page below corresponds to the result of searching for jag2 in Vega:

_____



If a text-based search fails to return any meaningful output then we can instead use one of the available alignment algorithms. If a term does not return any output in a text-based searched does not mean that associated feature is missing. It might be that due to a difference in the annotation (like a missing exon) or because the term is not present in ZFIN it was not feasible to link it to an annotated feature.

## SSAHA and BLAST

The Vega/Ensembl browsers provide a page where you can search using different sequences  to query the databases. You can access the BLASTView page from any of the Vega/Ensembl views through the link labelled 'run a BLAST search' in Ensembl or 'BLAST' in Vega.

The **BLASTView** page is an interface to set up, run and visualise the output of a sequence-based search. It is designed to work in clear steps where the user can first enter the query and select the target for the search, configure the parameters for the algorithm and finally customise the format of the output.

Open a BLASTView page in the Vega browser. In order to specify the query for the search you have the option of either using the sequence(s) or using an EMBL/GenBank accession number.

_____



Enter the accession number AF229449 and click on the button 'Retrieve'. Verify that *Danio rerio* is selected as the target database. You can choose the method to be used for the search. If your query is DNA you can choose from:

- BLASTN - the well-known BLAST algorithm performing a DNA-DNA search.
- SSAHA - this tool runs a hash-based algorithm. It is very fast since most of the needed data structures are pre-loaded. It works very well when searching for near-exact matches. It only performs searches where the query is DNA.
- TBLASTX - the BLAST algorithm where query and target are translated to the six possible reading frames. This is recommended when the query sequence is from a different organism.

Select the SSAHA algorithm and run the search by clicking on the 'Run' button. After some seconds you will be presented with the result. Every search is identified with a ticket that can be used later to retrieve a result (the results are stored for a couple of days). A Blast search can take a few minutes and your job may have to queue until it gets executed.
The result of searching the *Danio rerio* Vega database with AF229449 is the following page:

The result page can be customised and the relevant hits sorted in different ways. The most relevant match is framed in the diagrammatic view of the chromosomes. There is also a diagram indicating the coverage of the query,

which can be relevant to identify a repetitive subsequence in the query. At the bottom of the page there is a list of all hits. You can change the order of this list using the toolbar provided. If you are looking for a gene it is helpful to order the hits by, for example, chromosome coordinates. The matches can also be displayed in a ContigView page by selecting the [C] link.



This ContigView page adds a new track with the relevant hits. In this example the alignment coincides with an exon of a manually annotated gene (perhaps jag2!).

**Important note**

The Ensembl database contains the latest zebrafish assembly with automatic annotation (currently version Zv6 – March 2006 ). The zebrafish assembly is obtained by integrating all the available sequenced clones with a whole genome shotgun assembly. When reading and interpreting the outcome of a search it is important to understand the quality of the underlying sequence. In particular remember that the current assembly still includes sequences that do not have a chromosome assigned. If the best hit of your search matches one of these 'floating' fragments it will not appear in a framed box (since this only covers chromosomes 1 to 25). Refer to the exercises for an example. In section 2 the structure of the zebrafish assemblies is explained in more detail.

_____

The Vega database contains all the finished clones featuring high-quality manual annotation. This database is updated more often than Ensembl in order to incorporate new annotation and reflect changes in the map. When searching Vega you should bear in mind that it currently covers part of the zebrafish genome. An unsuccessful search in Vega is not sufficient evidence to conclude that the query sequence is not present in zebrafish. Despite not being complete, Vega features the best sequence with the best annotation and should be your starting point when searching the zebrafish genome. As explained in section 3, the Vega database also contains sequenced clones from the AB strain (the AB chromosome). Chromosome U collects all the sequenced clones that have not been assigned to chromosomes.

**Searching for all Finished/Unfinished clones**

New sequenced clones come through the pipeline on a daily basis. These sequences are submitted to EMBL/GenBank. Although sequenced clones are made public as soon as possible it takes time until they appear in Vega or in a new assembly. The Sanger Institute offers a Blast search page whose target is all the available sequenced clones for zebrafish. This service can be accessed though the *Danio rerio* project page or directly at:
http://www.sanger.ac.uk/cgi-bin/blast/submitblast/d_rerio

_____

This search can be used to find out whether a clone containing your region of interest is covered by a sequenced clone. An unfinished clone might be submitted in several contigs with artificial gaps between them. Contigs from unfinished clones can be long enough to contain a gene but they will not present in Vega until they are properly finished.

This collection of all sequenced clones does not replace the assembly since it is incomplete and also lacks all the extra features like alignments and automatic gene predictions. Moreover the sequences in this collection are isolated without a tiling path or contextual information. If you want to learn more about a clone and its flanking regions you can query the FPC database at:

http://www.sanger.ac.uk/Projects/D_rerio/WebFPC/zebrafish/small.shtml



FPC contig 10, for example, contains 62 clones and 5 have been sequenced. Click on this FPC contig to see more information:

_____



## Trace repository

The whole genome shotgun assembly used to build the zebrafish assembly is based on a collection of reads from cosmids, fosmids and BAC ends. This collection currently gives around 7x coverage of the genome. All these reads are deposited into the trace repository at

http://trace.ensembl.org

This page contains a long list. Look for *Danio rerio* and then expand the list. The reads are sorted by type and the sequencing centre of origin. All these reads can be downloaded but you can also perform a SSAHA search using the server at:

http://trace.ensembl.org/perl/ssahaview

If you know the name of a read or the prefix from the plasmid, fosmid or BAC you can use the provided text boxes to search for them.

_____



## Exercises

1. A TextView search looks for data in different indices and includes stable ids and other text-based information. Try for example using the text "activator of transcription" (including the quotes) in Vega or Ensembl.

2. Perform a SSAHA search with AF229449 but using the *Danio rerio* Ensembl database as the target (in the example above the search was done with Vega as the target). Do you obtain the same results using SSAHA and BLAST?

3. Search the Ensembl/Vega database with a human/mouse cDNA. What is the best approach?

4. Search the Ensembl/Vega database with a human protein. Why is SSAHA not in the list of possible tools?

5. Use the sequence

   GCCCTTAACTGATCGGCTTCTTCAGCAAGAGAGTTGCAAAGTACAG

   to search in Ensembl using SSAHA. Where did you find the best match? Try using the same sequence with BLAST. Why do think you get more hits with BLAST? Try using the same sequence in Vega.