# Module 1: Manual Genome Annotation and *de novo* Analysis of Sequence
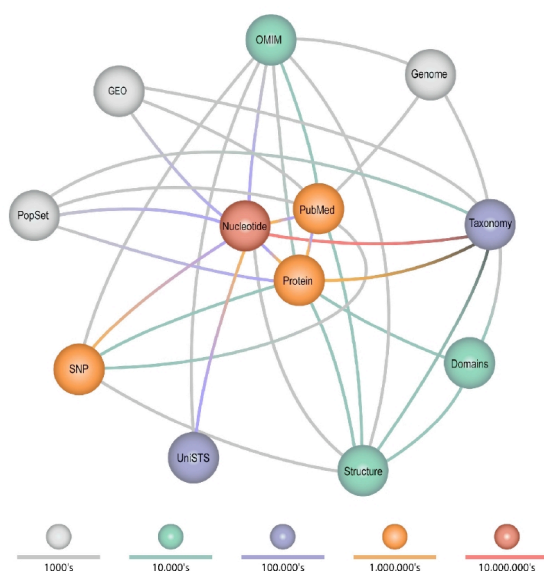
## Aims

- How to access sequence information using Entrez and UniProt.
- General *de novo* analysis of cDNAs, including ORF finder to highlight putative protein products of a cDNA, blastp link from within this to investigate the potential protein products, Spidey to align cDNA to genomic DNA. Clustalw to align similar sequences, view in Jalview and use GeneDoc to produce a graphics file.
- View manually annotated genes in the Vega browser.
- Perform a Blast search using sanger blast server on all finished and unfinished zebrafish clones

## General Introduction

Genomic and protein sequences can be accessed from various databases around the world. I will now introduce the major ways to access this sequence information, namely Entrez at NCBI and UniProt at EBI.

**Entrez** is a search and retrieval system that integrates information from databases at NCBI including the reference sequence (RefSeq) collection. The databases at NCBI include nucleotide sequences, protein sequences, macromolecular structures, whole genomes, and MEDLINE (through PubMed).



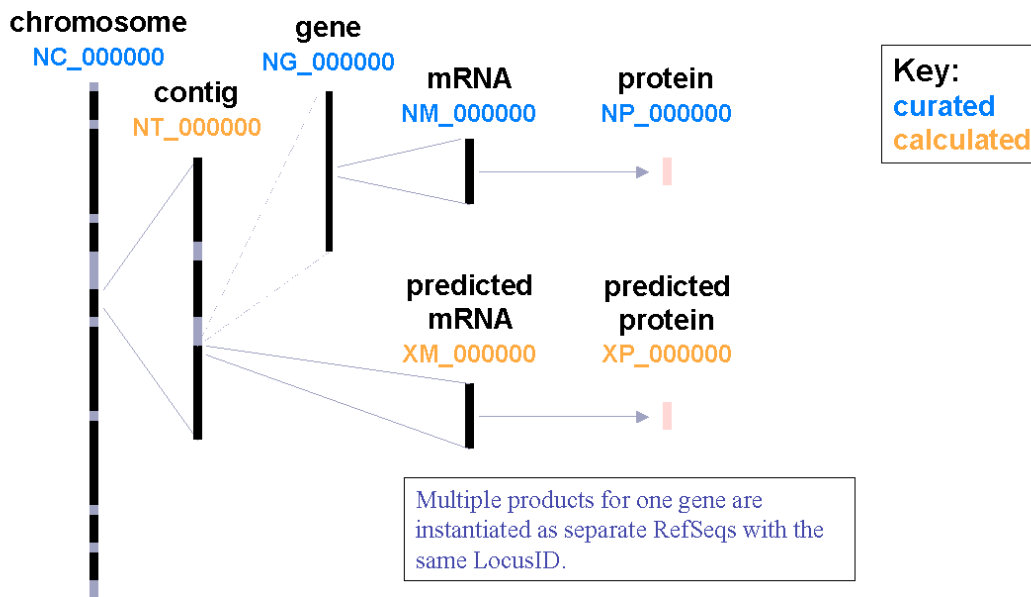Dataflow diagram of the Entrez retrieval system

_____

**Main Entrez entry point:**



**The Reference Sequence (RefSeq) database** provides a biologically non-redundant collection of DNA, RNA, and protein sequences. Each RefSeq represents a single, naturally occurring molecule from a particular organism. RefSeqs are frequently based on GenBank records but differ in that each RefSeq is a synthesis of information, not a piece of a primary research data in itself. Similar to a review article in literature, a RefSeq is an interpretation by a particular group at a particular time. RefSeqs can be retrieved in several different ways: by searching the Entrez Nucleotide or Protein database, by BLAST searching, by FTP, or through links from other NCBI resources.

# Reference Sequences

Goal: One sequence entry for each naturally occurring DNA, RNA and protein molecule



**For further information about RefSeq please visit: http://www.ncbi.nlm.nih.gov/RefSeq/key.html#accession**

## Curated Records

NC_123455          Complete genomic sequence (chromosome)

NG_123456          Incomplete genomic sequence

NM_123456          mRNA

NP_123456          Protein derived from NM

NR_123456          Non-coding RNA

## Model Records

NT_123456          Assembly of BAC data

NW_123456          Assembly of WGS data

NZ_ABCD12345678  Collection of WGS data

XM_123456          mRNA

XP_123456          Protein derived from XM

ZP_123456          Protein derived from NZ

XR_123456          Non-coding RNA

**The benefits of RefSeq:**

- non-redundancy
- explicitly linked nucleotide and protein sequences
- updates to reflect current knowledge of sequence data and biology
- data validation and format consistency
- distinct accession series
- ongoing curation by NCBI staff and collaborators, with reviewed records indicated

**UniProt** (Universal Protein Resource) is the world's most comprehensive catalogue of protein information. It is a central repository of protein sequence and function created by joining the information contained in Swiss-Prot (a database containing non-redundant, manual high-quality annotation that is cross-referenced to many other databases, with flatfiles containing features and information on individual protein sequences. Each entry has a unique accession number), TrEMBL (an automatic computer-annotated supplement to SwissProt, containing the translations of all coding sequences (CDS) present in the EMBL Nucleotide Sequence Database that are not yet integrated into SwissProt), and The Protein Information Resource (PIR), located at Georgetown University Medical Center (GUMC).

_____

**Module Introduction**

The Vertebrate Genome Annotation (Vega) database is a central repository for high quality, frequently updated, manual annotation of vertebrate finished genome sequence. The browser is based on the code from the Ensembl project and the data is produced by the Human and Vertebrate Analysis and Annotation (HAVANA) group at the Wellcome Trust Sanger Institute. Currently the available species are human, mouse, zebrafish, pig and dog. The annotation is undertaken in collaboration and synchronisation with the central zebrafish database ZFIN.

If you have generated a piece of sequence (either cDNA or genomic sequence) you are unlikely to want to wait for the sequence to be annotated automatically by genome viewers. It is possible to annotate the sequence manually, using freely available tools to access various databases. In addition, some of these tools, such as ORF finder, are not used by genome viewers so if you have found a region of interest in Ensembl or NCBI, you may wish to analyse an already annotated sequence using these programs to provide additional information. You may also want to perform *in silico* predictions of expression of a putative protein product, search for similar proteins or align similar sequences.

Such analyses are part of the following flow diagram:

The following is a summary of the DNA and protein analysis programs used in this module, together with a brief outline of their functions:

*ORF Finder (Open Reading Frame Finder)*
A graphical analysis tool that finds open reading frames in a sequence. The putative protein sequences may then be blasted by integral blastp.

*Spidey*
This aligns cDNA to genomic DNA and will output a list of exons together with an alignment and protein translation.

*CLUSTALW*
DNA and protein sequences can be aligned with clustalw and viewed in Jalview.

*GeneDoc*
A desktop package, used to produce graphical images of sequence alignments in many formats.
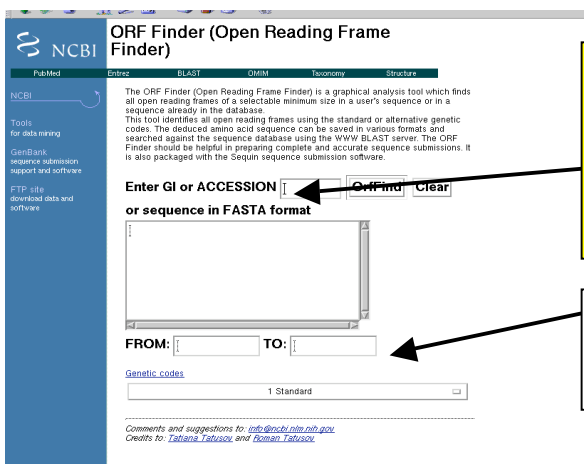
_____

**The Vega database**

Vega differs from Ensembl in that it shows annotation from the labour intensive process of manual curation produced by the Human and Vertebrate Analysis and Annotation (HAVANA) group at the Wellcome Trust Sanger Institute (WTSI). Finished genomic sequence is analysed on a clone by clone basis using a combination of similarity searches against DNA and protein databases and a series of *ab initio* gene predictions. Annotation is based on supporting evidence, which is external sequence such as ESTs, cDNAs and protein, and is performed to the standards agreed at the human annotation workshops (HAWK). Vega displays complete chromosomes and specific regions of interest. Grey shading indicates annotation status, with light grey showing partially annotated regions and dark grey showing regions with no annotation. Currently, human chromosomes 1, 3, 6, 7, 8, 9, 10, 13, 14, 16, 17, 18, 19, 20, 22, X and Y have full manual annotation, together with 44 genome-wide ENCODE regions and eight human haplotypes (6-COX, 6-QBL, 6-SSTO, 6-APD, 6-DBB, 6-MANN, 6-MCF, 6-PGF) of the chromosome 6 MHC region. There are also CORF genes on chromosomes 1, 2, 3, 4, 5, 8, 9, 11, 12, 15 and 17. The CORF project (at the WTSI) aims to produce a cDNA clone of each protein coding gene in the human genome, and manual annotation has been used to confirm ORFs and UTRs, to enable the design of PCR primers. Mouse currently has chromosomes 2, 4, 11 and X, together with candidate Insulin Dependent Diabetes (IDD) regions and the DeL36H regions of chromosome 13. Zebrafish has annotation for the majority of genes in ZFIN as well as full chromosome annotation on chrs 1, 2, 4, 5, 8, 9, 10, 13, 18, 19, 20, 22 and 23 currently displayed. Pig shows the MHC (SLA) region on chromosome 7 from Large White Boar and 8Mb of chromosome 17, and dog shows the MHC (DLA) class II region on chromosome 12 from Doberman.

Vega is an important contributor to the conserved 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.

_____

**1.** Does AB006190 have any open reading frames? Use NCBI ORF finder to find ORFs. What potential protein products does this cDNA code for? Use the link from ORF finder to investigate using blastp.

### ORF Finder: www.ncbi.nlm.nih.gov/gorf/gorf.html



Type in accession number AB006190 into box, or copy and paste sequence in fasta format. **Click on OrfFind**

Select bases and genetic code if required



Position and length of ORFs in the sequence

Shaded boxes indicate ORFs. AB006190 has 7 potential ORFs.

Click on 'SixFrames'

ORF Finder (Open Reading Frame Finder)

Blue lines = start codons
Pink lines=stop codons



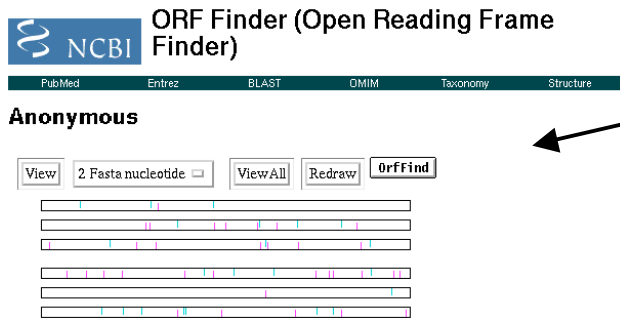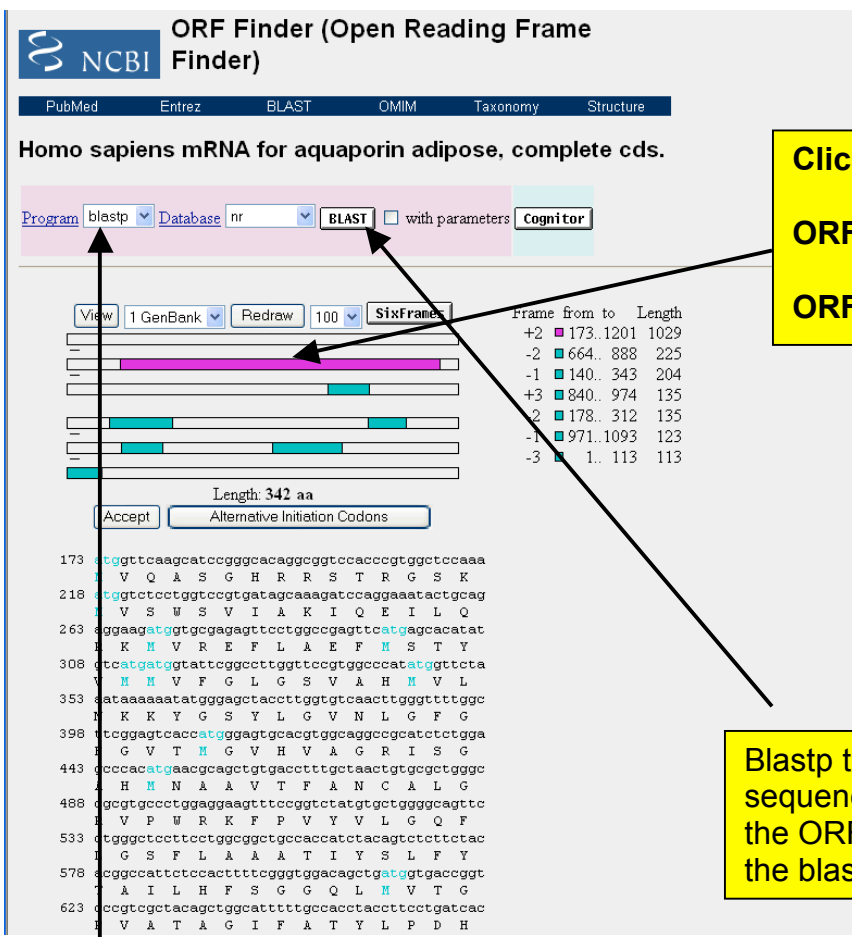ORF Finder (Open Reading Frame Finder)

Homo sapiens mRNA for aquaporin adipose, complete cds.

Click on a shaded
ORF to get this view.
ORF changes to pink

| Frame | from..to | Length |
|-------|----------|--------|
| +2 | 173..1201 | 1029 |
| -2 | 664.. 888 | 225 |
| -1 | 140.. 343 | 204 |
| +3 | 840.. 974 | 135 |
| 2 | 178.. 312 | 135 |
| -1 | 971..1093 | 123 |
| -3 | 1.. 113 | 113 |

Length: 342 aa

```
173  atggttcaagcatccgggcacaggcggtccacccgtggctccaaa
     M  V  Q  A  S  G  H  R  R  S  T  R  G  S  K
218  atggtctcctggtccgtgatagcaaagatccaggaaatactgcag
     M  V  S  W  S  V  I  A  K  I  Q  E  I  L  Q
263  aggaagatggtgcgagagttcctggccgagttcatgagcacatat
     R  K  M  V  R  E  F  L  A  E  F  M  S  T  Y
308  gtcatgatggtattcggccttggttccgtggcccatatggttcta
     V  M  M  V  F  G  L  G  S  V  A  H  M  V  L
353  aataaaaaatatgggagctaccttggtgtcaacttgggttttggc
     N  K  K  Y  G  S  Y  L  G  V  N  L  G  F  G
398  ttcggagtcaccatggagttgcacgtggcaggccgcatctctgga
     F  G  V  T  M  G  V  H  V  A  G  R  I  S  G
443  gcccacatgaacgcagctgtgacctttgctaactgtgcgctgggc
     A  H  M  N  A  A  V  T  F  A  N  C  A  L  G
488  gtcgtgccctggaggaagtttccggtctatgtgctggggcagttc
     V  V  P  W  R  K  F  P  V  Y  V  L  G  Q  F
533  gtgggctccttcctggccggctgccaccatctacagtctcttctac
     G  S  F  L  A  A  A  T  I  Y  S  L  F  Y
578  gcggccattctccacttttcggggtggacagctgatggtgaccggt
     A  I  L  H  F  S  G  G  Q  L  M  V  T  G
623  gccgtcgctacagctggcattttttgccacctaccttcctgatcac
     V  A  T  A  G  I  F  A  T  Y  L  P  D  H
```

Blastp the protein
sequence shown below
the ORF by clicking on
the blast button.

blastp or tblastn and
database can be selected

9

Length: 113 aa

Accept | ATG Initiation Codon

241 ctgtgctggcagctggcttctgtggcctggggttatgagtcttgtt
     L  C  W  Q  L  A  S  V  A  W  V  M  S  L  V
286 caatcgatagtccagacaccatccgcctccacttgccttctgt
     Q  S  I  V  Q  T  P  S  T  L  H  L  P  F  C
331 ccccaccagcagatagatgacttttttatgtgaggtcccatctctg
     P  H  Q  Q  I  D  D  F  L  C  E  V  P  S  L
376 attcgactctcctgtggagatacctcctacaatgaactccagttg
     I  R  L  S  C  G  D  T  S  Y  N  E  L  Q  L
421 gctgtgtccagtgtcatcttcgtggttgtgcctctcagcctcatc
     A  V  S  S  V  I  F  V  V  V  P  L  S  L  I
466 cttgcctcttatggagccactgcccaggcagtgctgaggattaac
     L  A  S  Y  G  A  T  A  Q  A  V  L  R  I  N
511 tctgccacagcatggagaaaggcctttgggacctgctcctcccat
     S  A  T  A  W  R  K  A  F  G  T  C  S  S  H
556 ctcactgtggtcaccctcttctacagc 582
     L  T  V  V  T  L  F  Y  S

Alternative initiation codons: highlights leucine codons as well as methionine codons.

Click on accept when you are happy with a CDS (the highlighted region will turn green).

NCBI    **ORF Finder (Open Reading Frame Finder)**
PubMed    Entrez    BLAST    OMIM    Taxonomy    Structure

**Homo sapiens mRNA for aquaporin adipose, complete cds.**

View | 3 Fasta protein | ▼ | Redraw | 100 | SixFrames |

| Frame | from | to | Length |
|---|---|---|---|
| +2 | 122.. | 1201 | 1080 |
| -2 | 664.. | 888 | 225 |
| -1 | 140.. | 343 | 204 |
| +3 | 840.. | 974 | 135 |
| -2 | 178.. | 312 | 135 |
| -1 | 971.. | 1093 | 123 |
| -3 | 1.. | 113 | 113 |

Select the view to be fasta protein and click on view to bring up the sequence in fasta format. This can now be copied and pasted into other applications.

>lcl|Sequence 1 ORF:173..1201 Frame +2
MVQASGHRRSTRGSKMVSWSVIAKIQEILQRKMVREFLAEFMSTYVMMVFGLGSVAHMVLNKKYGSYLGV
NLGFGFGVTMGVHVAGRISGAHMNAAVTFANCALGRVPWRKFPVYVLGQFLGSFLAAATIYSLFYTAILH
FSGGQLMVTGPVATAGIFATYLPDHMTLWRGFLNEAWLTGMLQLCLFAITDQENNPALPGTEALVIGILV
VIIGVSLGMNTGYAINPSRDLPPRIFTFIAGWGKQVFSNGENWWWVPVVAPLLGAYLGGIIYLVFIGSTI
PREPLKLEDSVAYEDHGITVLPKMGSHEPTISPLTPVSVSPANRSSVHPAPPLHESMALEHF*

10

Homo sapiens mRNA for aquaporin adipose, complete cds.

Program blastp ▾ Database nr ▾ BLAST ☐ with parameters Cognitor

**Click on BLAST button**

**Click on View report to view BLAST results**

Getting Started   Latest Headlines   Apple ▾   Amazon   eBay   Yahoo!   News ▾

BLAST                              *Basic Local Alignment Search Tool*

Home   Recent Results   Saved Strategies   Help

▸ NCBI/ BLAST/ Format Request

Query    lcl|2640 (342 letters)

Database   nr

Job title   lcl|2640 (342 letters)

Request ID   ZYCC0UMZ012        View report  ☐ Show results in a new window

Format

Show    Alignment ▾  as  HTML ▾  ☐ Advanced View    Reset form to defaults

Alignment View    Pairwise ▾

Display   ☑ Graphical Overview   ☑ Linkout   ☑ Sequence Retrieval   ☐ NCBI-gi

Masking Character: Lower Case ▾   Masking Color: Grey ▾

Limit results   Descriptions: 100 ▾   Graphical overview: 100 ▾   Alignments: 100 ▾

Organism    Type common name, binomial, taxid, or group name. Only 20 top taxa will be shown.
Enter organism name or id--completions will be suggested

Entrez query:

Expect Min:          Expect Max:

Format for   ☐ PSI-BLAST   with inclusion threshold:

Copyright | Disclaimer | Privacy | Accessibility | Contact | Send feedback   new interface

**Can also format for PSI-Blast if required by checking the PSI-BLAST box**

11

_____

BLASTP results

NCBI                    *results of* **BLAST**

**BLASTP 2.2.9 [May-01-2004]**

<u>Reference</u>:
Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schäffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1997),
"Gapped BLAST and PSI-BLAST: a new generation of protein database search
programs",  Nucleic Acids Res. 25:3389-3402.


RID: 1089369892-23111-86899017791.BLASTQ4


**Query=**
        (359 letters)


**Database:** All non-redundant GenBank CDS
translations+PDB+SwissProt+PIR+PRF excluding environmental samples
          1,921,851 sequences; 641,055,585 total letters


If you have any problems or questions with the results of this search
please refer to the **BLAST FAQs**


<u>Taxonomy reports</u>


```
                                                        Score    E
Sequences producing significant alignments:            (bits) Value

ref|NP_001161.1|  aquaporin 7; aquaporin adipose [Homo sapie...   695    0.0
dbj|BAC05693.1|  aquaporin adipose [Homo sapiens]                691    0.0
emb|CAD13298.1|  bA251017.3 (similar to aquaporin 7) [Homo s...   635    0.0
ref|XP_376852.1|  similar to bA251017.3 (similar to aquapori...   583    e-165
ref|XP_376866.1|  similar to bA251017.3 (similar to aquapori...   583    e-165
ref|XP_372110.2|  similar to bA251017.3 (similar to aquapori...   562    e-159
ref|NP_031499.1|  aquaporin 7 [Mus musculus] >gi|9910621|sp|...   442    e-123
dbj|BAC36431.1|  unnamed protein product [Mus musculus]          442    e-123
ref|NP_062030.2|  aquaporin 7 [Rattus norvegicus] >gi|321724...   431    e-119
gb|AAH62701.1|  AQP7 protein [Homo sapiens]                      396    e-109
ref|NP_956204.1|  Unknown (protein for MGC:63700); wu:fj98f0...   300    4e-80
emb|CAG01413.1|  unnamed protein product [Tetraodon nigrovir...   280    3e-74
ref|NP_004916.1|  aquaporin 3 [Homo sapiens] >gi|2497938|sp|...   275    9e-73
gb|AAP36954.1|  Homo sapiens aquaporin 3 [synthetic construct]   275    9e-73
emb|CAG46822.1|  AQP3 [Homo sapiens]                             274    3e-72
ref|NP_113891.1|  aquaporin 3 [Rattus norvegicus] >gi|135196...   273    3e-72
dbj|BAA04559.1|  aquaporin 3 [Rattus rattus]                     273    3e-72
```

**2.** Align the genomic sequence AL133548 against a piece of cDNA sequence, BC007459, using Spidey.

**SPIDEY:**
**http://www.ncbi.nlm.nih.gov/IEB/Research/Ostell/Spidey/**

Results:    Exon length    Percentage identity

Genomic: gi|7899142|emb|AL133548.6| Human DNA sequence from clone RP11-235O14 on chromosome 9q13-21.2 Contains a novel gene (FLJ20559), the OSTF1 gene for osteoclast stimulating factor 1 (OSF, SH3P2) and a CpG island

mRNA: gi|13938612|gb|BC007459.1| Homo sapiens osteoclast stimulating factor 1, mRNA (cDNA clone MGC:12220 IMAGE:4052054), complete cds

Alignment is on plus strand of genomic sequence and on plus strand of mRNA sequence
mRNA coverage: 99%
Overall percent identity: 99.5%
Non-aligning poly(A) tail: 30

Graphical overview

59766 ———————————————————————— 118358

| | Genomic coordinates | mRNA coordinates | length | identity | mismatches | gaps | Donor site | Acc. site |
|---|---|---|---|---|---|---|---|---|
| Exon 1 | 59766-59916 | 1-151 | 151 | 100.0% | 0 | 0 | d | |
| Exon 2 | 88685-88731 | 152-198 | 47 | 100.0% | 0 | 0 | d | a |
| Exon 3 | 98750-98800 | 199-249 | 51 | 100.0% | 0 | 0 | d | a |
| Exon 4 | 101759-101822 | 250-313 | 64 | 100.0% | 0 | 0 | d | a |
| Exon 5 | 102951-103004 | 314-367 | 54 | 100.0% | 0 | 0 | d | a |
| Exon 6 | 104475-104582 | 368-475 | 108 | 100.0% | 0 | 0 | d | a |
| Exon 7 | 105527-105576 | 476-525 | 50 | 100.0% | 0 | 0 | d | a |
| Exon 8 | 108719-108797 | 526-604 | 79 | 98.7% | 1 | 0 | d | a |
| Exon 9 | 112015-112113 | 605-703 | 99 | 100.0% | 0 | 0 | d | a |

Genomic co-ordinates            mRNA co-ordinates

Zebrafish Workshop

Module 1: *de novo* Analysis of Sequence
_____

Exon 3: 98750-98800 (genomic); 199-249 (mRNA)

```
98750     TTTCTTCTAGCCAGATGAATTATACTTTGAGGAAGGTGATATTATCTACA
                    |||||||||||||||||||||||||||||||||||||||||
199               CCAGATGAATTATACTTTGAGGAAGGTGATATTATCTACA
                   P  D  E  L  Y  F  E  E  G  D  I  I  Y


98790     TTACTGACATGGTAAGTCCAG
                    |||||||||||
239       TTACTGACATG
          I  T  D  M
```

Top

Exon 4: 101759-101822 (genomic); 250-313 (mRNA)

```
101759    CTTTTACCAGAGCGATACCAATTGGTGGAAAGGCACCTCCAAAGGCAGGA
                    |||||||||||||||||||||||||||||||||||||||||||
250               AGCGATACCAATTGGTGGAAAGGCACCTCCAAAGGCAGGA
                     S  D  T  N  W  W  K  G  T  S  K  G  R


101799    CTGGACTAATTCCAAGCAACTATGGTAAGTGTTG
                    ||||||||||||||||||||||||
290       CTGGACTAATTCCAAGCAACTATG
          T  G  L  I  P  S  N  Y
```

Top

Exon 5: 102951-103004 (genomic); 314-367 (mRNA)

```
102951    TTCAATCTAGTGGCTGAGCAGGCAGAATCCATTGACAATCCATTGCATGA
                    ||||||||||||||||||||||||||||||||||||||||||
314               TGGCTGAGCAGGCAGAATCCATTGACAATCCATTGCATGA
                  V  A  E  Q  A  E  S  I  D  N  P  L  H  E
```

Alignments for each exon:

There are other cDNA to genomic alignment programs available:

EST2GENOME
http://bioweb.pasteur.fr/seqanal/interfaces/est2genome.html


SIM4
http://pbil.univ-lyon.fr/sim4.html

**3.** Use clustalw to align similar sequences.

To fetch some sequences, used UniProtKB:

> Search UniProtKB for RPS6.
> Select a human, mouse, rat, cow and chicken sequence.



Click on retrieve

Select Fasta and open to display sequences in a new window

_____

# http://www.ebi.ac.uk/clustalw/index.htm



**Go into ClustalW. Copy and paste sequences and click on run…….**

**Top of results page**

**Alignment**

[ Show Colors ]   [ View Alignment File ]

**Scroll down to view alignment**

CLUSTAL W (1.82) multiple sequence alignment

```
swall|P18653|K6A1_MOUSE    MPLAQLKEPWPLMELVPLDPENGQTSGEEAGLQPS--------------- 35
swall|Q63531|K6A1_RAT      MPLAQLKEPWPLMELVPLDPENGQASGEEAGLQPS--------------- 35
swall|Q15418|K6A1_HUMAN    MPLAQLKEPWPLMELVPLDPENGQTSGEEAGLQPS--------------- 35
swall|P18652|K6AA_CHICK    MPLAQLAEPWPNMELVQLDTENGQAAPEEGGNPPCKAKSDITWVEKDLVD 50
                           ****** **** **** **.****:: **.*  *.

swall|P18653|K6A1_MOUSE    ---KDEAILKEISITHHVKAGSEKADPSQFELLKVLGQGSFGKVFLVRKV 82
swall|Q63531|K6A1_RAT      ---KDEGILKEISITHHVKAGSEKADPSHFELLKVLGQGSFGKVFLVRKV 82
swall|Q15418|K6A1_HUMAN    ---KDEGVLKEISITHHVKAGSEKADPSHFELLKVLGQGSFGKVFLVRKV 82
swall|P18652|K6AA_CHICK    STDKGEGVVKEINITHHVKEGSEKADPSQFELLKVLGQGSFGKVFLVRKI 100
                           *.*.::***.****** ********:****************:
```

**Click on show colours**

```
swall|P18653|K6A1_MOUSE    TRPDSGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|Q63531|K6A1_RAT      TRPDNGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|Q15418|K6A1_HUMAN    TRPDSGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|P18652|K6AA_CHICK    TPPDSNHLYAMKVLKKATLKVRDRVRTKIERDILADVNHPFVVKLHYAFQ 150
                           * **..********************:*******************

swall|P18653|K6A1_MOUSE    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|Q63531|K6A1_RAT      TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|Q15418|K6A1_HUMAN    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|P18652|K6AA_CHICK    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 200
                           *************************************************

swall|P18653|K6A1_MOUSE    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|Q63531|K6A1_RAT      IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|Q15418|K6A1_HUMAN    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|P18652|K6AA_CHICK    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 250
                           *************************************************

swall|P18653|K6A1_MOUSE    VVNRQGHTHSADWWSYGVLM-----------GKDRKETMTLILKAKLGMP 271
swall|Q63531|K6A1_RAT      VVNRQGHTHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 282
swall|Q15418|K6A1_HUMAN    VVNRQGHSHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 282
swall|P18652|K6AA_CHICK    VVNRQGHSHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 300
                           *******:************          ****************

swall|P18653|K6A1_MOUSE    QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHIFYSTIDWNKLYRR 321
swall|Q63531|K6A1_RAT      QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHIFYSTIDWNKLYRR 332
swall|Q15418|K6A1_HUMAN    QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHVFYSTIDWNKLYRR 332
swall|P18652|K6AA_CHICK    QFLSAEAQSLLRALFKRNPANRLGSGPDGAEEIKRHPFYSTIDWNKLYRR 350
                           ****:*************************** *************
```

**Alignment**

[ Hide Colors ]   [ View Alignment File ]

CLUSTAL W (1.82) multiple sequence alignment

```
swall|P18653|K6A1_MOUSE    MPLAQLKEPWPLMELVPLDPENGQTSGEEAGLQPS--------------- 35
swall|Q63531|K6A1_RAT      MPLAQLKEPWPLMELVPLDPENGQASGEEAGLQPS--------------- 35
swall|Q15418|K6A1_HUMAN    MPLAQLKEPWPLMELVPLDPENGQTSGEEAGLQPS--------------- 35
swall|P18652|K6AA_CHICK    MPLAQLAEPWPNMELVQLDTENGQAAPEEGGNPPCKAKSDITWVEKDLVD 50
                           ****** **** **** **.****:: **.*  *.

swall|P18653|K6A1_MOUSE    ---KDEAILKEISITHHVKAGSEKADPSQFELLKVLGQGSFGKVFLVRKV 82
swall|Q63531|K6A1_RAT      ---KDEGILKEISITHHVKAGSEKADPSHFELLKVLGQGSFGKVFLVRKV 82
swall|Q15418|K6A1_HUMAN    ---KDEGVLKEISITHHVKAGSEKADPSHFELLKVLGQGSFGKVFLVRKV 82
swall|P18652|K6AA_CHICK    STDKGEGVVKEINITHHVKEGSEKADPSQFELLKVLGQGSFGKVFLVRKI 100
                           *.*.::***.****** ********:****************:

swall|P18653|K6A1_MOUSE    TRPDSGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|Q63531|K6A1_RAT      TRPDNGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|Q15418|K6A1_HUMAN    TRPDSGHLYAMKVLKKATLKVRDRVRTKMERDILADVNHPFVVKLHYAFQ 132
swall|P18652|K6AA_CHICK    TPPDSNHLYAMKVLKKATLKVRDRVRTKIERDILADVNHPFVVKLHYAFQ 150
                           * **..********************:*******************

swall|P18653|K6A1_MOUSE    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|Q63531|K6A1_RAT      TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|Q15418|K6A1_HUMAN    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 182
swall|P18652|K6AA_CHICK    TEGKLYLILDFLRGGDLFTRLSKEVMFTEEDVKFYLAELALGLDHLHSLG 200
                           *************************************************

swall|P18653|K6A1_MOUSE    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|Q63531|K6A1_RAT      IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|Q15418|K6A1_HUMAN    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 232
swall|P18652|K6AA_CHICK    IIYRDLKPENILLDEEGHIKLTDFGLSKEAIDHEKKAYSFCGTVEYMAPE 250
                           *************************************************

swall|P18653|K6A1_MOUSE    VVNRQGHTHSADWWSYGVLM-----------GKDRKETMTLILKAKLGMP 271
swall|Q63531|K6A1_RAT      VVNRQGHTHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 282
swall|Q15418|K6A1_HUMAN    VVNRQGHSHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 282
swall|P18652|K6AA_CHICK    VVNRQGHSHSADWWSYGVLMFEMLTGSLPFQGKDRKETMTLILKAKLGMP 300
                           *******:************          ****************

swall|P18653|K6A1_MOUSE    QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHIFYSTIDWNKLYRR 321
swall|Q63531|K6A1_RAT      QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHIFYSTIDWNKLYRR 332
swall|Q15418|K6A1_HUMAN    QFLSTEAQSLLRALFKRNPANRLGSGPDGAEEIKRHVFYSTIDWNKLYRR 332
swall|P18652|K6AA_CHICK    QFLSAEAQSLLRALFKRNPANRLGSGPDGAEEIKRHPFYSTIDWNKLYRR 350
                           ****:*************************** *************
```
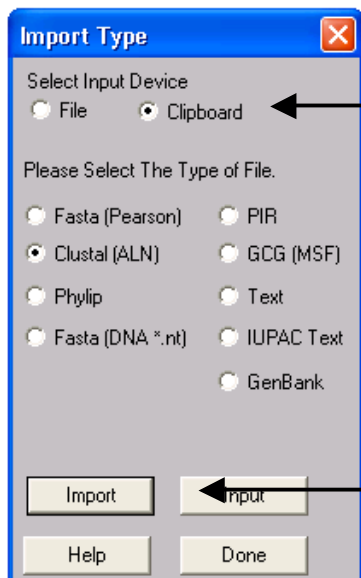
Click on Jalview to view coloured, shaded alignment.



To edit the alignment:
Click above the amino acid in the sequence to select the column. From the edit menu select remove sequence to left of selected column………see below
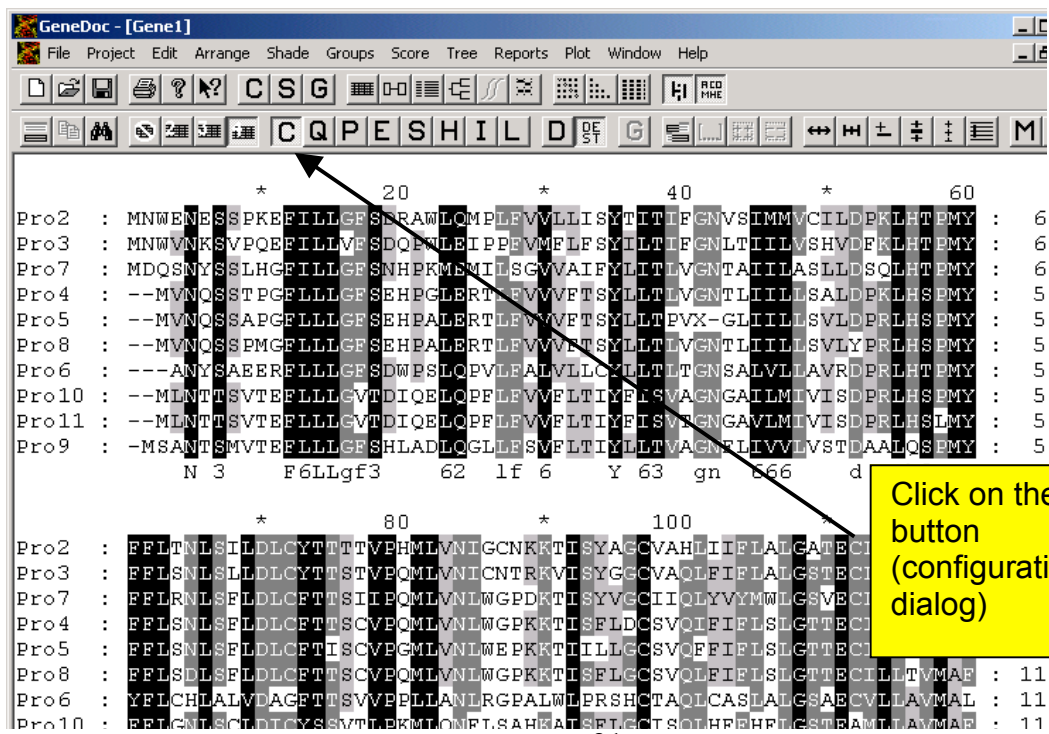
Create a coloured and shaded figure with GeneDoc (Not available for Mac):
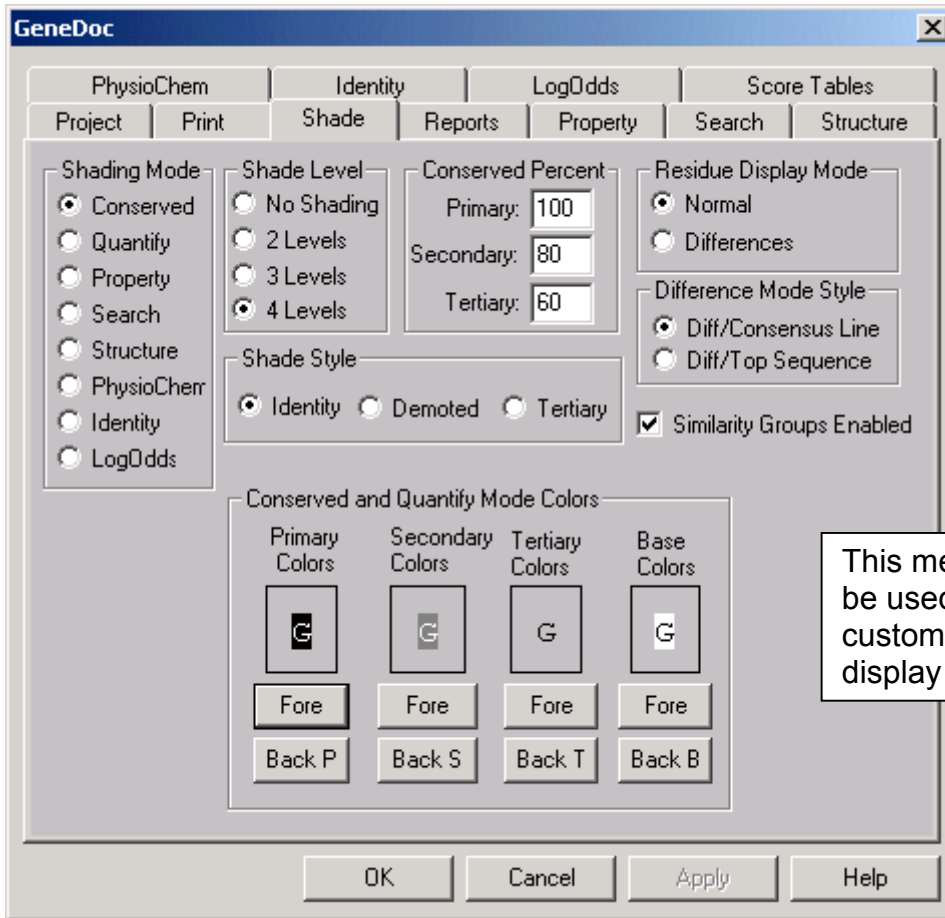GeneDoc can be downloaded from here: **http://www.psc.edu/biomed/genedoc/**



Copy the aln format alignment from clustalw (**make sure that you include the top line of text**), then open up GeneDoc and select file, import. Select clipboard as the input device and Clustal for the file type.

Click on import, then done

**The imported file will be displayed in black and white**



Click on the C button (configuration dialog)

21

This menu box can be used to customise the display



22

Change to colours of your choice

_____



From the edit menu click on Select Blocks for Copy

The background of each selected block will turn black. To copy the whole alignment, select all the blocks



From under the edit menu select Copy Selected Blocks to, RTF file

**Other file types are also available. Use RTF to import into Word, and metafile file to import into PowerPoint**

Import the alignment into Word. RTF format files may be edited in Word

```
Pro2  : MNWENESSPKEFILLGESDRAWLQMPLFVVLLISYTITLFGNVSIMM
Pro3  : MNWVNKSVPQEFILLVESDQPWLEIPPEVMFLFSYILTLFGNLTIIL
Pro7  : MDQSNYSSLHGFILLGESNHPKMEMILSGVVAIFYLITLVGNTAIILA
Pro4  : --MVNQSSTPGFLLLGESEHPGLERTLFVVVFTSYLLTLVGNTLIIL
Pro5  : --MVNQSSAPGFLLLGESEHPALERTLFVVVFTSYLLTPVX-GLIILLSVLDPRLLSEMY : 57
Pro8  : --MVNQSSPMGFLLLGESEHPALERTLFVVVFTSYLLTLVGNTLIILSVLYPRLSPMY : 58
Pro6  : ---ANYSAEERFLLLGESDWPSLQPVLFALVLLCYLLTLTGNSALVLAVREPRLLTPMY : 57
Pro10 : --MLNTTSVTEFLLLGVTDIQELQPFLFVVFLTIYFISVAGNGAIILMIVISDPRLLSEMY : 58
Pro11 : --MLNTTSVTEFLLLGVTDIQELQPFLFVVFLTIYFISVTGNGAVLMIVISDPRLISLMY : 58
Pro9  : -MSANTSMVTEFLLLGFSHLADLQGLLFSVFLTIYLLTVAGNFLIVVLVSTDAALQSFMY : 59
```

```
                *           80        *         100        *         120
Pro2  : FFLTNLSILDLCYTTTTVPHMLVNIGCNKKTLSYAGCVAHLIIFLALGATECLLALMSE : 120
Pro3  : FFLSNLSLLDLCYTISTVPQMLVNICNTRAVLSYGGCVAQLFIFLALGSTECLLALMCE : 120
Pro7  : FFLRNLSFLDLCFTTSIIPQMLVNLWGPDRTLSYVGCIIQLYVYMWLGSVECLLALMSY : 120
Pro4  : FFLSNLSFLDLCFTTSCVPQMLVNLWGPKKTLSFLDCSVQLFIFLSLGTTECLLLTVMAE : 118
Pro5  : FFLSNLSFLDLCFTISCVPGMIVNLWEPKKTLILLCCSVQFFIFLSIGTTECILLTVMAE : 117
Pro8  : FFLSDLSFLDLCFTTSCVPQMIVNLWGPKKTLSFLCCSVQTFIFLSIGTTECILLTVMAE : 118
Pro6  : YFLCHLALVTAGFTSVVPPLIANLRGPALWLPRSHCTAQLCASLALGSAECVLLAVMAL : 117
Pro10 : FFLGNLSCLDICYSSVTLPKMIQNFLSAHKALSFLCCISQLHFFHFLGSTEAMLAVMAE : 118
Pro11 : FFLGNLSYLDICYSTVTLPKMIQNFLSTHKALSFLCCISQLHFFHFLGSTESMLFAVMAE : 118
Pro9  : FFLRTLSAFEIGYTSVTVPLLIHHLLTGRRHTSRSGCALQMFFFIFFGATECCLLAAMAY : 119
```

You may also align sequences directly from UniProtKB by clicking on the Align button:



There are many other alignment programs available on the web:

**MULTIALIN** http://prodes.toulouse.inra.fr/multalin/multalin.html

**PIMA** http://bioweb.pasteur.fr/seqanal/interfaces/pima-simple.html

**DIALIGN** http://bioweb.pasteur.fr/seqanal/interfaces/dialign2-simple.html

**DCA** http://bioweb.pasteur.fr/seqanal/interfaces/dca-simple.html

_____

## Manual Genome Annotation

The underlying data for the Vega database is generated by the Havana group.

Below is a screen-shot of the crfb1 locus from the groups' annotation software. Protein coding genes are shown in red and green, whilst non-coding transcripts are shown in red.
Other columns show Blast hits to DNA and protein databases, repeats and gene models generated from the solexa sequenced transcriptome data.

# The Vertebrate Genome Annotation (VEGA) Database
## Worked example:
**1.** View the crfb1 locus. How many transcripts are there ? What is the supporting evidence ? Export peptide sequence.



STEP 1:
**Load** Vega:
**http://vega.sanger.ac.uk**

STEP 2:
**Select zebrafish genome annotation**

STEP 3:
**Search for gene symbol crfb1**

Ideograms of annotated chromosomes and additional information

Link to ZFIN gene entry

STEP 5:
**Click on transcript and link through to Supporting evidence**



The various proteins and EST's used to build variant

STEP 6:
Link through to
Export Peptide

STEP 7:
Choose Sequence
and format type

_____

# Sanger Blast: 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

**Worked Example:**

Find which clone the 5' end of the znf385b gene is on ?

5' end of the znf385b gene lies in a gap.

STEP 1:
**Select link to supporting evidence**

STEP 2:
**Select link to external record for BC059587.**



STEP 3:
**Copy mRNA sequence**

**STEP 4:**
**Paste sequence into Query Data box. Start Blast**

Finished and unfinished clones.

Choose method.

**Blast Server Results**

Retrieve result for id: `d01TrBOmfa9782hD3`  (retrieve)

Format: Graphical

Your BLAST query has been added to the queue of jobs.
The majority of BLASTs are completed within two minutes.

To retrieve your results, click the **retrieve** button above, or use the following
URL: http://www.sanger.ac.uk/cgi-bin/blast/getblast?
id=d01TrBOmfa9782hD3;format=graphic

Click here to start a new blast job

Options: cpus=1 -warnings B=100 -filter=dust V=100
Job id d01TrBOmfa9782hD3.1 status is DONE

BLASTN 2.0MP-WashU [04-May-2006] [linux26-i686-ILP32F64 2006-05-09T11:47:08]

Copyright (C) 1996-2006 Washington University, Saint Louis, Missouri USA.
All Rights Reserved.

Reference:  Gish, W. (1996-2006) http://blast.wustl.edu

Notice:  this program and its default parameter settings are optimized to find
nearly identical sequences rapidly.  To identify weak protein similarities
encoded in nucleic acid, use BLASTX, TBLASTN or TBLASTX.

Query=  UNKNOWN-QUERY
        (1767 letters)



Database:  zebrafish_all.fa
           33,354 sequences; 1,924,132,883 total letters.
Searching....10....20....30....40....50....60....70....80....90....100% done

```
                                                   Smallest
                                                     Sum
                                           High  Probability
Sequences producing High-scoring Segment Pairs:    Score  P(N)      N

AL732458.5                                         1609  4.7e-169  3
CU855879.5                                         1581  7.3e-156  4
AL713987.7                                         1581  5.3e-118  3
BX957295.5                                         1609  9.9e-103  2
zK5E8.00938    EMBL:CU693458   Unfinished sequence: zK5E8    Conti... 1601  1.4e-66   1
AL732411.14                                        1619  8.3e-66   1
AL732436.6                                          732  4.5e-23   1
BX571968.10                                         353  1.0e-20   3
BX663515.13                                         353  1.2e-20   3
CR848668.16                                         426  6.5e-09   1
zK5E8.00101    EMBL:CU693458   Unfinished sequence: zK5E8    Conti...  364  4.4e-06   1
BX470219.6                                          364  4.6e-06   1
```

Missing exons align with
CU855879 – finished clone

>AL732458.5
[Full Sequence] [EMBL:AL732458.5]



```
        Length = 84,824

  Plus Strand HSPs:

 Score = 1609 (247.5 bits), Expect = 4.7e-169, Sum P(3) = 4.7e-169
 Identities = 325/329 (98%), Positives = 325/329 (98%), Strand = Plus / Plus
[HSP Sequence]

Query:  1438 GCCAAACTCGCCCTGCAGAATGACTTGGTGAAGCCCATTTCACCAGCCTTCCTCCCGTCA 1497
             |||||||||||||||||||||||||||||||| ||||||||||||||||||||||||||
Sbjct: 12395 GCCAAACTCGCCCTGCAGAATGACTTGGTGAAGCCCATTTCACCAGCCTTCCTCCCGTCA 12454

Query:  1498 CCCTTCTCTACGACCACAGTCCCGTCCATCTCTCTCCACCCTCGCCCCAACACCTCCATC 1557
             |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 12455 CCCTTCTCTACGACCACAGTCCCGTCCATCTCTCTCCACCCTCGCCCCAACACCTCCATC 12514

Query:  1558 TTCCAGACGGCTTCACTTCCGCACTCGTTTCTCCGTGCCGCTCCCGGACCCATTCGACCC 1617
             |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 12515 TTCCAGACGGCTTCACTTCCGCACTCGTTTCTCCGTGCCGCTCCCGGACCCATTCGACCC 12574

Query:  1618 ACTACCGGCTCCATCCTCTTTGCGCCTTACTGAGCGGCTGGATATTGAAACAGTGGTGTA 1677
             |||||||||||||||||||||||||||||||||||||||||||||||||||||||||||
Sbjct: 12575 ACTACCGGCTCCATCCTCTTTGCGCCTTACTGAGCGGCTGGATATTGAAACAGTGGTGTA 12634

Query:  1678 TTTATTGAAACTTTATCCTCAACTGTCCTCACCATCACGAAAACCCAGGAAGGATATCGT 1737
             |||||||||||||||||||||||||||||||||||||||||  ||||||||||||||||
Sbjct: 12635 TTTATTGAAACTTTATCCTCAACTGTCCTCACCATCACGAAAACCCAATAAGGATATCGT 12694
```

_____

**Tasks**

1.
How many ORFs does U64564 have and are there any potential protein products and/or domains? (Use OrfFinder) Repeat this for AF177198 and blast the protein product from one of the longest ORFs. What protein hits do you get? Accept this protein and view in fasta format.
2.

Use ClustalW to produce an alignment of human olfactory receptor proteins. Search UniProtKB with OR6 in and select several sequences of your choice. For example: O76002, O76001, Q9GZK4, Q9Y3N9, O76000, P58173.  Do the alignment in ClustalW and then copy the aln format alignment into GeneDoc. Open this alignment in GeneDoc and produce an rtf file to import into Word.

3.
Find the rnd1l gene in Vega. Which end falls in the gap ? Can you find the clone that contains the missing exons ?

**Answers:** (These are likely to change due to database updates etc.)

1.

**ORF of U64564**:

There are 6 in the forward frames and 10 in the reverse frames. The longest is in frame one and when blasted there is a hit to myelin oligodendrocyte glycoprotein. The CD search has found an immunoglobulin-like domain.

ORF of AF177198:

There are 39 ORF in total. Blastp of the longest ORF on frame +1 will give hits to the TNL1 (Talin 1). Click on accept and view in fasta protein format. The sequence can be copied and pasted into Word for later use or directly into another program.

**Protein from ORF:**

Various Talin family proteins and B41 and ILWEQ domains are found. Further iterations bring up hits to more distantly related species e.g. Tetraodon, Apis mellifera, Anopheles

Q6NXR9 CDD domains:

DEADc and HELICc domains are detected.

*2.*
*Use ClustalW and GeneDoc to produce an olfactory receptor 6 protein alignment*

Search for OR6 in UniProtKB (or in Entrez).

In Entrez, click on proteins, then select your sequences of choice. Change the display to FASTA and click on 'send to' so the view refreshes.

Copy and paste the sequences into ClustalW. Make sure that the space at the header of each sequence is removed to ensure proper formatting if using Unix. Submit the job.

Copy the ALN format alignment from clustalw, and then open GeneDoc and select file, import. Select clipboard as the input device and aln for filetype. Click on import. Change the colours from within the configuration dialog box and anything else you wish. From the edit menu choose select blocks to copy and click on the blocks. Then select copy selected blocks to, RTF file. Save the file and open in Word.

3.

The 3' end of the rnd11 gene fall in a gap in Vega. Copy the supporting cDNA BC122138 sequence from the EMBL file and blast using Sanger Danio Rerio blast server. The missing exon aligns to finished clone CU862020.

Zebrafish Workshop

Module 1: *de novo* Analysis of Sequence
_____

37