**JKU**
JOHANNES KEPLER
UNIVERSITÄT LINZ

# Genome data analysis
## Computer lab session 3

Theresa Schwarz MSc

theresa.schwarz@jku.at

Institute of Biophysics, JKU

Gruberstraße 40, 4020 Linz

❑ **Genome browsers**

➢ **UCSC**
➢ **ENSEMBL**

❑ **BLAST**

➢ **Pairwise alignments**
➢ **Database alignments**
➢ **Primer-BLAST**
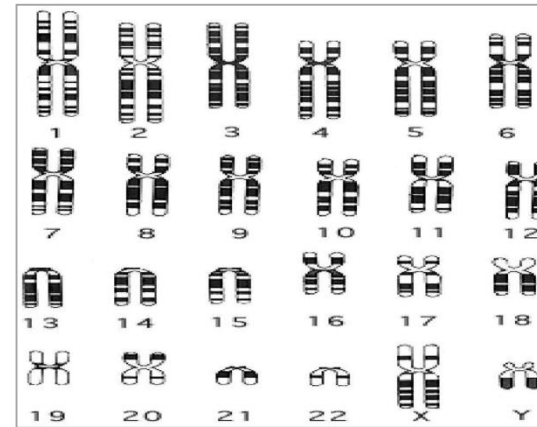
❑ **Genome browsers**

  ➢ **UCSC**
  ➢ **ENSEMBL**

❑ **BLAST**

  ➢ **Pairwise alignments**
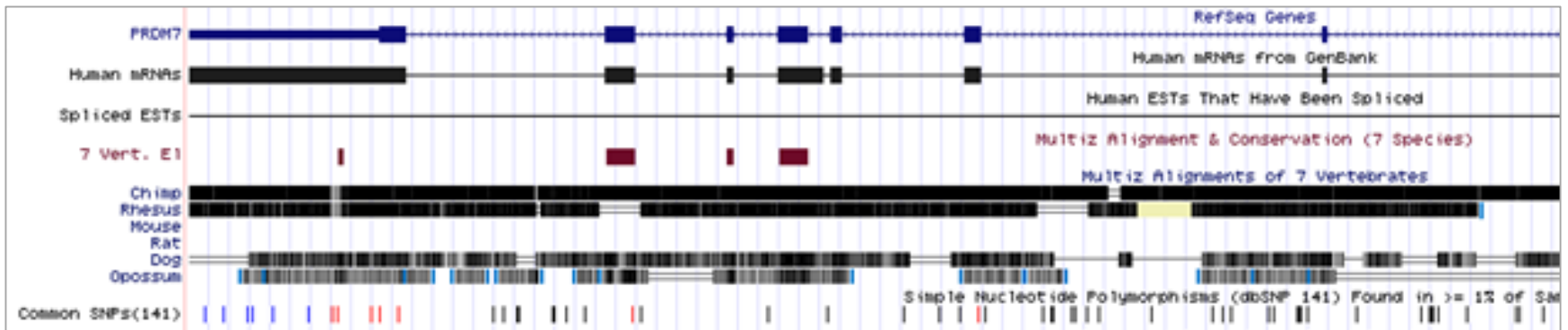  ➢ **Database alignments**
  ➢ **Primer-BLAST**

Theresa Schwarz

# Genome browsers

- **UCSC** Genome Browser (University of California Santa Cruz)
  *https://genome.ucsc.edu/*

- **ENSEMBL** (EMBO-Heidelberg/EBI-Cambridge)
  *http://www.ensembl.org/*

- **NCBI** (NIH, US) Genome Map Viewer
  *https://www.ncbi.nlm.nih.gov/mapview/*

Theresa Schwarz

# Genome browsers

- Genomic DNA is organized in chromosomes.



- Genome browsers display ideograms (pictures) of chromosomes. Users can select '**annotation tracks**' that display many kinds of information.
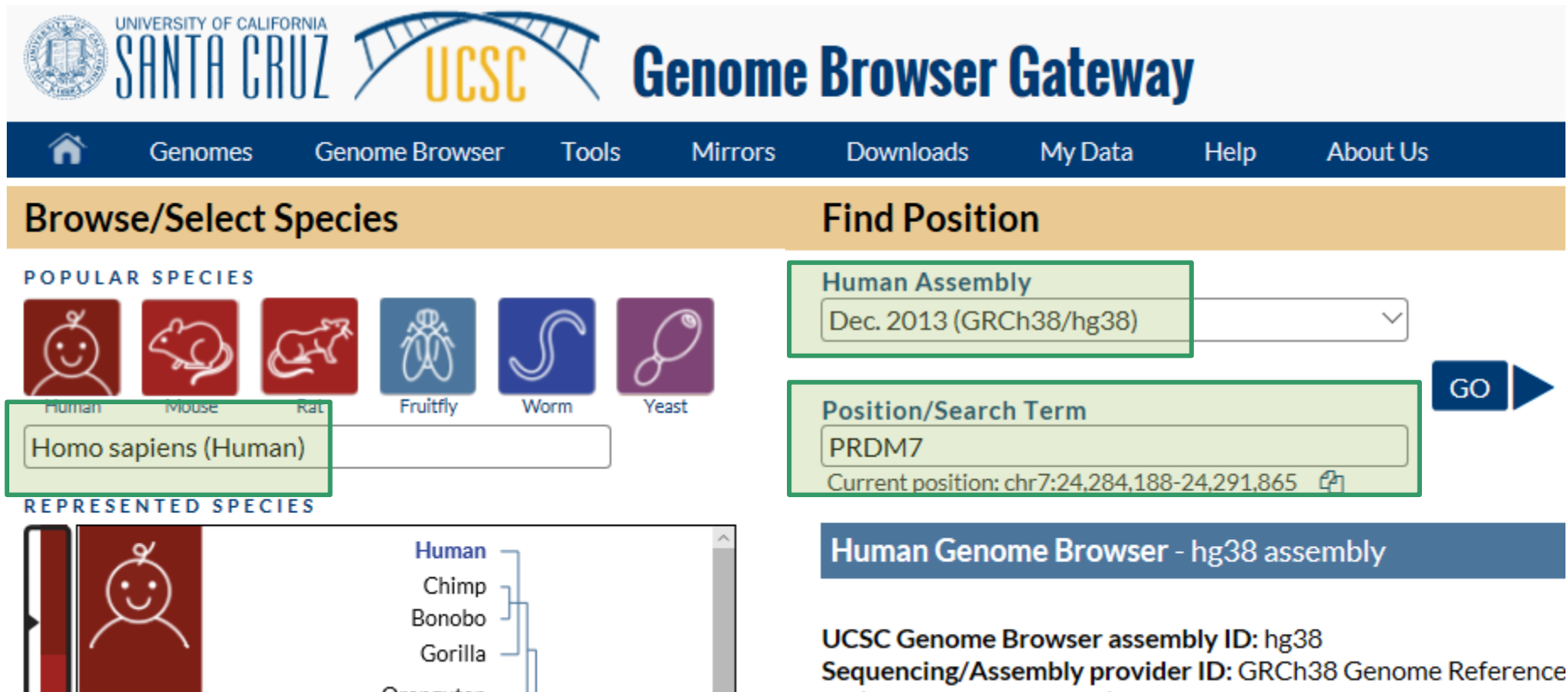


Theresa Schwarz

# UCSC Genome Browser

# UCSC Genome Browser

Search for the **human** protein **PRDM7** using the **newest genome assembly**.



Theresa Schwarz

# UCSC Genome Browser

**Known Genes**

PRDM7 **(uc010cje.4)** at chr16:90057780-90075930 - Homo sapiens PR domain 7 (PRDM7), mRNA. (from RefSeq NM_001098173
PRDM7 (uc059ywn.1) at chr16:90075005-90092072 - PR domain containing 7 (from HGNC PRDM7)
PRDM7 (uc059ywm.1) at chr16:90061452-90075930 - The sequence shown here is derived from an Ensembl automatic ana.
PRDM7 (uc059ywl.1) at chr16:90061452-90075925 - PR domain containing 7 (from HGNC PRDM7)
PRDM7 (uc002fqo.4) at chr16:90056566-90062325 - PR domain containing 7 (from HGNC PRDM7)
TRAF1 **(uc010mvl.2)** at chr9:120902393-120929173 - Homo sapiens TNF receptor associated factor 1 (TRAF1), transcript
TRAF1 (uc011lyg.2) at chr9:120902393-120914572 - Homo sapiens TNF receptor associated factor 1 (TRAF1), transcript
TRAF1 (uc004bku.3) at chr9:120902393-120928769 - Homo sapiens TNF receptor associated factor 1 (TRAF1), transcript

**RefSeq Genes**

PRDM7 at chr16:90056566-90075930 - (NM_001098173) probable histone-lysine N-methyltransferase PRDM7

**Basic Gene Annotation Set from GENCODE Version 24 (Ensembl 83)**

PRDM7 at chr16:90057780-90075930

**Comprehensive Gene Annotation Set from GENCODE Version 24 (Ensembl 83)**

PRDM7 at chr16:90056566-90062325
PRDM7 at chr16:90057780-90075930
PRDM7 at chr16:90061452-90075925
PRDM7 at chr16:90061452-90075930
PRDM7 at chr16:90075005-90092072

# UCSC Genome Browser
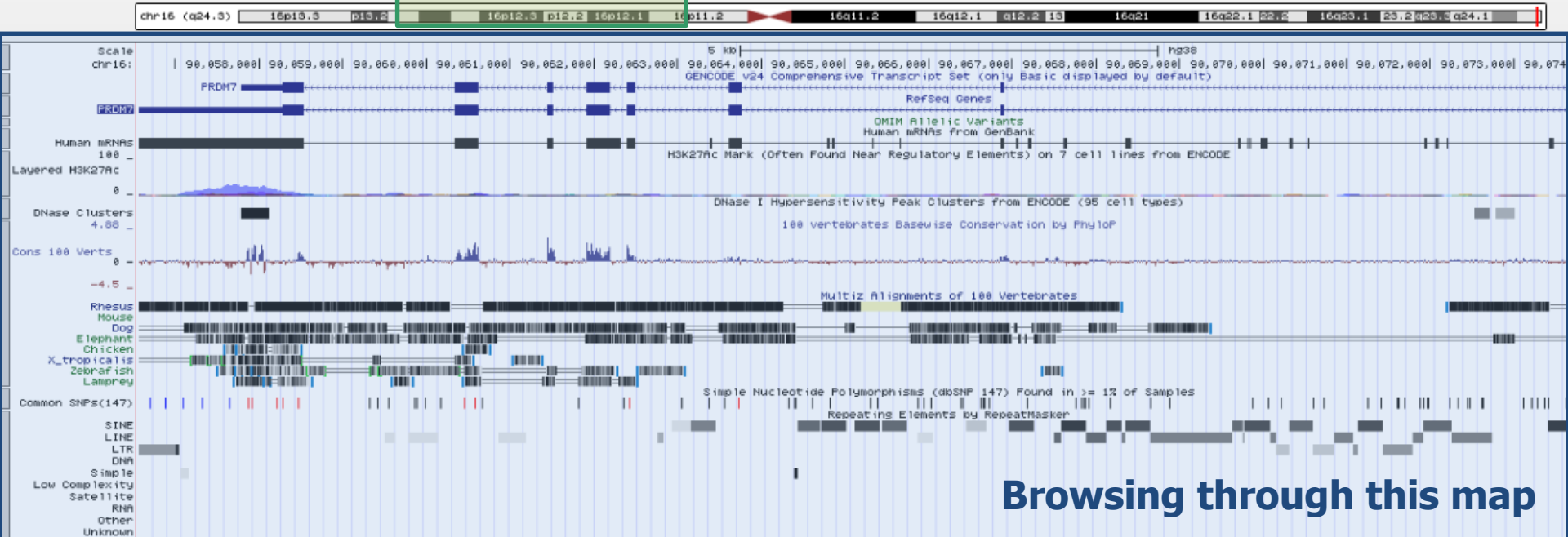


**Chromosome** and **coordinates**

**Browsing through this map**

You can **hide all** tracks

→ **Scroll down to select your settings**

Theresa Schwarz

# UCSC Genome Browser

You can choose all 'tracks' (settings) you want.

- Choose:

  - **NCBI RefSeq (full)**

  - **CCDS (full)**

- To update your settings select '**refresh**'



Theresa Schwarz

# UCSC Genome Browser



You can see your selected tracks

Boxes = Exons; Lines = Introns

Navigate through the map by **moving** and **zooming**.

# UCSC Genome Browser



If you would like to get the DNA sequence click on **View → DNA**

# UCSC Genome Browser

- You can add bases up- and/or downstream of the DNA: e.g. 500bp

- You can highlight the '**tracks**' you selected by '**extended case/color options**'

CCDS is highlighted in red.

Theresa Schwarz

❑ **Genome browsers**

➢ **UCSC**

➢ **ENSEMBL**

❑ **BLAST**

➢ **Pairwise alignments**

➢ **Database alignments**

➢ **Primer-BLAST**

# ENSEMBL Genome Browser



- **ENSEMBL** is a project between EMBL-EBI (European Bioinformatics Institute) and the Wellcome Trust Sanger Institute.

- This Genome Browser provides information about **eukaryotic genomes.**

- You can find an accurate description of **protein coding genes, promoters, exons, introns, transcripts, …**

Theresa Schwarz

# ENSEMBL Genome Browser

*http://www.ensembl.org/*



**Search again for human PRDM7**

Theresa Schwarz

# ENSEMBL

# ENSEMBL



→ Click on **Show transcript table**

# ENSEMBL

**Transcripts**

Hide transcript table

Show/hide columns (1 hidden)

Filter

| Name | Transcript ID | bp | Protein | Biotype | CCDS | UniProt | RefSeq | Flags | | |
|------|--------------|-----|---------|---------|------|---------|--------|-------|---|---|
| PRDM7-002 | ENST00000449207.6 | 2008 | 492aa | Protein coding | CCDS45557 | Q9NQW5 | NM_001098173 NP_001091643 | TSL:1 | GENCODE basic | APPRIS P1 |
| PRDM7-003 | ENST00000564210.2 | 714 | 73aa | Nonsense mediated decay | - | H3BUJ3 | - | | TSL:5 | |
| PRDM7-005 | ENST00000568473.5 | 706 | 138aa | Nonsense mediated decay | - | A4Q9G9 | - | | TSL:5 | |
| PRDM7-009 | ENST00000569206.1 | 693 | No protein | Processed transcript | - | - | - | | TSL:5 | |
| PRDM7-001 | ENST00000325921.10 | 2442 | No protein | Retained intron | - | - | - | | TSL:1 | |

Here you can see all transcript versions and some links to other databases like CCDS, UniProt.

For the first transcript you can also get the NCBI's **RefSeq** sequences for nucleotide and protein.

→ **Select the first entry and scroll down!**

Theresa Schwarz

**Statistics**    Exons: 10, Coding exons: 10, Transcript length: 2,008 bps, Translation length: 492 residues

CCDS    This transcript is a member of the Human CCDS set: CCDS45557

Uniprot    This transcript corresponds to the following Uniprot identifiers: Q9NQW5

Transcript Support Level (TSL)    TSL:1

Ensembl version    ENST00000449207.6

Type    Known protein coding

Annotation Method    Transcript where the Ensembl genebuild transcript and the Vega manual annotation have the same sequence, for every base pair. See article.

Alternative transcripts    This transcript corresponds to the following database identifiers: Havana transcript: OTTHUMT00000420560

GENCODE basic gene    This transcript is a member of the Gencode basic gene set.

**HGNC Symbol: PRDM7-002**

| | |
|---|---|
| Gene | PR/SET domain 7 |
| | ENSG00000126856 |
| Location | Chromosome 16: 90,057,780-90,075,930 |
| Exon | 9 of 10 |
| Transcript | ENST00000449207.6 |
| | Exons |
| | cDNA Sequence |
| Protein | ENSP00000396732 |
| | Protein Variations |
| Gene type | Known protein coding |
| Transcript type | Known protein coding |
| Strand | Reverse |
| Base pairs | 2,008 |
| Amino acids | 492 |
| Source | Ensembl/Havana merge |

→ **To get more information about exons and introns, click on a box and select 'Exons'**

# ENSEMBL

Here you can obtain sequence information on **exons** and **introns.**



Theresa Schwarz

# ENSEMBL

- Go back and select **Location**.



You can view:

the chromosome

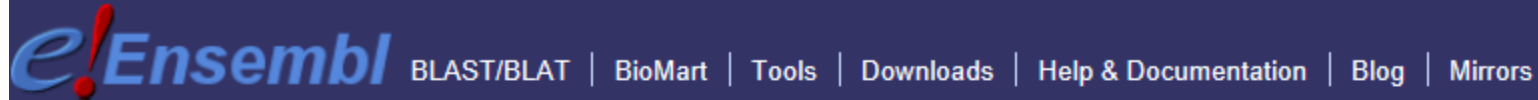a specific region of interest

a detailed view with certain tracks

❑ **Genome browsers**

   ➢ **UCSC**
   ➢ **ENSEMBL**

❑ **BLAST**

   ➢ **Pairwise alignments**
   ➢ **Database alignments**
   ➢ **Primer-BLAST**

Theresa Schwarz

# BLAST

- **BLAST** = Basic Local Alignment Search Tool

- **BLAST** is a nice tool to compare biological sequences and to find regions of identity or similarity.



Theresa Schwarz

- **Pairwise alignment:** process of lining up two sequences to achieve maximal levels of identity

- **Database alignments:** input sequence is aligned to similar sequences of an entire database

# Overview

☐ **Genome browsers**

➢ **UCSC**
➢ **ENSEMBL**

☐ **BLAST**

➢ **Pairwise alignments**
➢ **Database alignments**
➢ **Primer-BLAST**

Theresa Schwarz

- Where do you get the highest level of identity when comparing two sequences?

glu glu ala gly glu asp asp glu
asp gly ala glu asp glu asn asn ➢ **1**

glu glu ala gly glu asp asp glu ➢ **2**
   asp gly ala glu asp glu asn asn

glu glu ala gly glu asp asp glu ➢ **3**
   asp gly ala glu asp glu asn asn

Theresa Schwarz

- Aims (a few examples):

  - assess the **degree of similarity** of 2 sequences

  - search for **conservation** (e.g. protein domains or sequence motifs)

  - find **functionally or structurally** related proteins

  - assess the **possibility of homology**

Theresa Schwarz

When are two genes/proteins **<u>homologous</u>**, **<u>paralogous</u>** or **<u>orthologous</u>**?

- **Homologs** are related genes that descended from a common ancestral gene. Two genes can be separated by the event of **speciation** (see ortholog) or **gene duplication** (see paralog).

- **Paralogs** are related genes <u>in the same species</u> that have been separated by a duplication event within a genome. Paralogs mostly evolve <u>new functions</u>.

- **Orthologs** are related genes <u>in different species</u> that evolved from a common ancestral gene by speciation. Normally, orthologs retain the <u>same function</u> in the course of evolution.

Theresa Schwarz

Figure 12-4 Human Molecular Genetics, 3/e. (© Garland Science 2004)

- Hemoglobin tetramer oxygen transport in blood
- Myoglobin monomer oxygen transport in muscle
- Neuroglobin monomer oxygen transport in CNS
- Cytoglobin monomer oxygen transfer blood-brain

Theresa Schwarz

- Neuroglobin [Homo sapiens]

- Neuroglobin [Mus musculus]


- Alignment of those two neuroglobins

  - on level of amino acid sequence

    $\rightarrow$ **92% identities** $\rightarrow$ **difference in 13 amino acids**

  - on level of nucleotide sequence

    $\rightarrow$ **79% identities** $\rightarrow$ **difference in 193 bases**

    (divided by 3 = ~64 aa)

- When you want to find homologs or conserved domains the **amino acid sequence is much more informative** than the nucleotide sequence !

- Because: the genetic code is **redundant**

  (codons are degenerate: changes in the

  3rd position often do not change the aa)

  **3rd position = "wobble base"**



- Because: more characters in proteins: 20 amino acids vs. 4 bases

Theresa Schwarz

- **Protein alignments** are used

  - to find common ancestors million or billion of years ago (Amino acid sequences offer a longer **"look-back" time**)

  - DNA sequences can be translated into protein and then used in pairwise alignments to find homologs

- **DNA alignments** are used

  - to study DNA **polymorphisms** (SNPs, insertions, deletions, microsatellites,...)

  - to study **non-coding** regions of DNA

  - to confirm the **identity of a cDNA**

```
Query: 181 catcaactacaactccaaagacacccttacacccactaggatatcaacaaacctacccac 240
            ||||||||| |||| |||||| ||||| | ||||||||||||||||||||||||||||||||
Sbjct: 189 catcaactgcaaccccaaagccacccct-cacccactaggatatcaacaaacctacccac 247
```

Theresa Schwarz

# Pairwise alignment - BLAST

- You can do a pairwise alignment with nucleotide and protein sequences using BLAST – a tool of NCBI.

- *http://www.ncbi.nlm.nih.gov/*

→ Let's compare **Hemoglobin** and **Myoglobin** by using **Protein BLAST**.

# Pairwise alignment - BLAST

- Select **Align two or more sequences**



Theresa Schwarz

# Pairwise alignment - BLAST



Enter **Accession number** or sequence in **FASTA format** of the two proteins you want to compare:

Hemoglobin: **NP_000509.1**

Myoglobin: **NP_005359**

# Do they look similar ?

Human
Hemoglobin subunit beta
(NP_000509.1)

Human
Myoglobin
(NP_005359)

# Pairwise alignment - BLAST

- **No significant similarity found!**



- **We can go back and change some parameters**

Theresa Schwarz

# Pairwise alignment - BLAST



- Select **Algorithm parameters** (and scroll down)

- Change the Matrix (Scoring Parameters) from BLOSUM62 to **BLOSUM45**

**Repeat the BLAST search!**

Theresa Schwarz

**There are two kinds of sequence alignments using different matrices:**

**GLOBAL** alignment algorithm

- Needleman and Wunsch (1970)

**LOCAL** alignment algorithm

- Smith and Waterman (1981)

**GLOBAL** alignment extends from one end of each sequence to the other.

→ **PAM** matrices

Global

**LOCAL** alignment finds optimally matching regions within two sequences ("subsequences").

→ **BLOSUM** matrices

Local

Theresa Schwarz

# BLOSUM Matrices

- BLOSUM matrices are based on **LOCAL** alignments.

- BLOSUM stands for **blocks substitution matrix**.

- BLOSUM**80** is a matrix to compare sequences with similarities of **>80%**.

- BLOSUM**62** is a matrix to compare sequences with similarities of **>62%**.

- BLOSUM**45** is a matrix to compare sequences with similarities of **>45%**.

**The <u>higher</u> – the better !**

Theresa Schwarz

# PAM Matrices

- PAM matrices are based on **GLOBAL** alignments.

- PAM stands for **point accepted mutation**.

- PAM matrices are used to assess the **relatedness** of two proteins

- PAM matrices work like this:
  **How many differences are allowed per 100 amino acids?**

  **- PAM1            1 difference per 100 amino acids**
  **- PAM10.7         10 differences per 100 amino acids**
  **- PAM80           50 differences per 100 amino acids**
  **- PAM250          80 differences per 100 amino acids**

## The <u>lower</u> – the better !

Theresa Schwarz

| BLOSUM 80 | BLOSUM 62 | BLOSUM 45 |
| PAM 1 | PAM 120 | PAM 250 |

*Less divergent* ←————————————→ *More divergent*

**Closely related**                    **Distantly related**

**rat vs. mouse globin**               **rat vs. bacterial globin**

By choosing the matrix you can chose **which part** of the sequence should be used (global, local) and **how stringent** the alignment should be done.

# Pairwise alignment - BLAST

- When using a different **Scoring Matrix** an alignment is possible



- Scroll down to get more information about your results

Theresa Schwarz

# How does your BLAST result look like?

General information:

- Score

- Query coverage

- Expected (E) value

- Ident

- Accession

**How does your BLAST result look like?**

- **Score**:

  - ➤ a measure for the **quality** of the alignment
  - ➤ it is calculated by the scoring matrix and reflects the **degree of similarity**
  - ➤ <u>Max score</u>: Score of single best aligned sequence
  - ➤ <u>Total score</u>: Sum of scores of all aligned sequences
  - ➤ The higher the better!

⊟ **Descriptions**

Sequences producing significant alignments:

Select: <u>All</u> <u>None</u>   Selected:0

⫲ Alignments 🖫Download ⌄  GenPept  Graphics  Multiple alignment                                    ⚙

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| ☐ myoglobin [Homo sapiens] | 46.8 | 46.8 | 97% | 2e-12 | 26% | NP_005359.1 |

## How does your BLAST result look like?

- **Query coverage**:
  - ➢ Information on **how much of a sequence** is used for the alignment
  - ➢ Always check the query coverage to see whether the alignment is meaningful
  - ➢ The higher the better!

---

⊟ **Descriptions**

**Sequences producing significant alignments:**

Select: All None    Selected:0

⇅ Alignments  ▤ Download ⌄  GenPept  Graphics  Multiple alignment                    ⚙

| Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|
| ☐ myoglobin [Homo sapiens] | 46.8 | 46.8 | 97% | 2e-12 | 26% | NP_005359.1 |

## How does your BLAST result look like?

- **Expected ( E ) value**:

  - ➢ Represents the **significance** of a result
  - ➢ Probability of a random alignment
  - ➢ The lower the E-value the more significant
  - ➢ The lower the better!

⊟**Descriptions**

Sequences producing significant alignments:

Select: All None   Selected:0

⇅ Alignments  ⊟Download ⌄  GenPept  Graphics  Multiple alignment                    ⚙

| Description | Max score | Total score | Query cover | E value | dent | Accession |
|---|---|---|---|---|---|---|
| ☐ myoglobin [Homo sapiens] | 46.8 | 46.8 | 97% | 2e-12 | 26% | NP_005359.1 |

## How does your BLAST result look like?

- **Ident**:

  - Shows how many amino acids of the two sequences match perfectly

### Descriptions

**Sequences producing significant alignments:**

Select: All None   Selected:0

Alignments  Download  GenPept  Graphics  Multiple alignment

| | Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|---|
| ☐ | myoglobin [Homo sapiens] | 46.8 | 46.8 | 97% | 2e-12 | 26% | NP_005359.1 |

**How does your BLAST result look like?**

- **Accession**:
  - ➢ Protein accession number is directly linked to myoglobin entry in protein database

# Pairwise alignment - BLAST

Scroll down to **Alignments** for **details**.

- **Query** = **Hemoglobin**

- **Sbjct** = **Myoglobin**

- Sequence of alignment **hemoglobin vs. myglobin**

# Pairwise alignment - BLAST

- You can see again the '**Score**', '**E-value**' and '**Identities**'

- Total number of aligned amino acids: 145

- Positives

- Gaps

## What is the difference between <u>Identities</u> & <u>Positives</u>?



| | | | |
|---|---|---|---|
| Range 1: 3 to 147 GenPept Graphics | | | Next Match ▲ Previous Match |

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 46.8 bits(144) | 2e-12 | Compositional matrix adjust. | 37/145(26%) | 61/145(42%) | 2/145(1%) |

```
Query    4    LTPEEKSAVTALWGKVNVDEVG--GEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKV   61
              L+  E    V  +WGKV  D  G    E L RL+  +P T    F+ F  L + D +  +  +
Sbjct    3    LSDGEWQLVLNVWGKVEADIPGHGQEVLIRLFKGHPETLEKFDKFKHLKSEDEMKASEDL   62

Query   62    KAHGKKVLGAFSDGLAHLDNLKGTFA                                  ...
              K HG  VL A+    L     + +
Sbjct   63    KKHGATVLTALGGILKKKGHHEAEIK                                  ...

Query  122    EFTPPVQAAYQKVVAGVANALAHKY                                   ...
              +F    Q A  K +        +A  Y
Sbjct  123    DFGADAQGAMNKALELFRKDMASNY   147
```

Middle row displays identical and conserved amino acids; (+ sign for conserved amino acids)

**Identities:** amino acids that are **identical** at a specific position in the two sequences

**Similarity/Conservation:** amino acids at a specific position in the two sequences are **not identical**, however they share the same chemical properties (see amino acid classification on the next slide) = they are **similar**

**Positives:** The **sum** of identical and similar amino acids.

Theresa Schwarz

# Classification of amino acids



NONPOLAR

Glycine (Gly) — G
Alanine (Ala) — A
Valine (Val) — V
Leucine (Leu) — L
Isoleucine (Ile) — I
Methionine (Met) — M
Tryptophan (Trp) — W
Phenylalanine (Phe) — F
Proline (Pro) — P

POLAR

Serine (Ser) — S
Threonine (Thr) — T
Cysteine (Cys) — C
Tyrosine (Tyr) — Y
Asparagine (Asn) — N
Glutamine (Gln) — Q

Electrically Charged

Acidic
Aspartic Acid (Asp) — D
Glutamic Acid (Glu) — E

Basic
Lysine (Lys) — K
Arginine (Arg) — R
Histidine (His) — H

Dept. Biol. Penn State ©2002

Theresa Schwarz

# Pairwise alignment - BLAST

- You can see again the '**Score**', '**E-value**' and '**Identities**'

- The alignment ranges from aa 3-147 = 145 (relative to the Sbjct)

- **Positives**: aligned amino acids which are either **ident** or **similar**

- **Gaps**: NO alignment with any amino acid. Signed with a ' **–** '

How is the score calculated?

The score is a sum of match, mismatch and gap.

# Pairwise alignment – Summary

- Choose two sequences

- Select an **algorithm** that generates a score

- This algorithm can be used for <u>global</u> or <u>local</u> alignments

- **Score** reflects degree of similarity (quality control)

- **E-value** tells you the significance of the alignment

- Check whether your results are meaningful → **Query coverage**

Theresa Schwarz

❑ **Genome browsers**

➢ **UCSC**
➢ **ENSEMBL**


❑ **BLAST**

➢ **Pairwise alignments**
➢ **Database alignments**
➢ **Primer-BLAST**

- Despite pairwise alignment you can also align a sequence against an entire database

# BLAST

1.  Choose the BLAST program


2.  Choose sequence (query)


3.  Choose the database to search


4.  Choose optional parameters


Then click "BLAST"

Theresa Schwarz

# 1. Choose the BLAST program

- BLAST hemoglobin **NP_000509.1** against a very general protein database

- Therefore, first select **Protein BLAST** again.

Sequence can be entered in FASTA format or as accession number

Theresa Schwarz

# 3. Choose the database



protein databases

nucleotide databases

**nr = non-redundant (most general database)**

You can...

- choose the organism to search

- turn filtering on/off

- change the substitution matrix

- change the expect (e) value

- change the word size

- change the output format

# 4. Choose optional parameters



**Organism**

**Algorithm**

# 4. Choose optional parameters



**Algorithm parameters**    Note: Parameter values that differ from the default are highlighted in yellow and marked with ♦ sign

### General Parameters

| | |
|---|---|
| **Max target sequences** ♦ `20000 ▾` | **Number of outputs** |
| Select the maximum number of aligned sequences to display | |
| **Short queries** ☑ Automatically adjust parameters for short input sequences | |
| **Expect threshold** `10` | **Expect threshold = set a max. E-value** |
| **Word size** ♦ `3 ▾` | **Word size** |
| **Max matches in a query range** `0` | |

### Scoring Parameters

| | |
|---|---|
| **Matrix** `BLOSUM62 ▾` | **Scoring matrix** |
| **Gap Costs** `Existence: 11 Extension: 1 ▾` | |
| **Compositional adjustments** `Conditional compositional score matrix adjustment ▾` | |

### Filters and Masking

| | |
|---|---|
| **Filter** ☐ Low complexity regions | **Filter, mask** |
| **Mask** ☐ Mask for lookup table only | |
| ☐ Mask lower case letters | |

**BLAST**    Search **database Non-redundant protein sequences (nr)** using **Blastp (protein-protein BLAST)**
☑ Show results in a new window

1. The query sequence is cut into pieces ("**words**")
   e.g. one piece consists of **3** amino acids when the word size is **3**

2. The BLAST algorithm uses those "words" to find similar regions in sequences which are present in the chosen database

3. The default word size = 3 amino acids or 11 nucleic acids

**Word size = 3**

```
SWVSQA = Query
SWV
 WVS
  VSQ
   SQA
```

Theresa Schwarz

# BLAST Results



You can find:

- **general information** about your search

- **Graphic Summary** with the **Conserved Domains** of your Query sequence

- click on **search summary** to see all the chosen parameters

# BLAST Results

| Search Parameters | |
|---|---|
| Program | blastp |
| Word size | 3 |
| Expect value | 10 |
| Hitlist size | 20000 |
| Gapcosts | 11,1 |
| Matrix | BLOSUM62 |
| Filter string | F |
| Genetic Code | 1 |
| Window Size | 40 |
| Threshold | 11 |
| Composition-based stats | 2 |

| Database | |
|---|---|
| Posted date | Nov 18, 2016 8:24 AM |
| Number of letters | 38,985,428,197 |
| Number of sequences | 106,376,657 |
| Entrez query | txid9606 [ORGN] |

| Karlin-Altschul statistics | | |
|---|---|---|
| Lambda | 0.320339 | 0.267 |
| K | 0.136843 | 0.041 |
| H | 0.422367 | 0.14 |
| Alpha | 0.7916 | 1.9 |
| Alpha_v | 4.96466 | 42.6028 |
| Sigma | | 43.6362 |

Schwarz

# BLAST Results

- Scroll down to **Descriptions** to get an overview of all the results

Direct links to protein database

Sequences producing significant alignments:

Select 'All' to see how many results you got = 270

Select: All None   Selected:270

Alignments   Download   GenPept  Graphics  Distance tree of results  Multiple alignment

Score
Query coverage
E value
% Ident

| | Description | Max score | Total score | Query cover | E value | Ident | Accession |
|---|---|---|---|---|---|---|---|
| ☑ | hemoglobin subunit beta [Homo sapiens] | 301 | 301 | 100% | 4e-106 | 100% | NP_000509.1 |
| ☑ | beta globin chain variant [Homo sapiens] | 299 | 299 | 100% | 2e-105 | 99% | AAN84548.1 |
| ☑ | beta globin [Homo sapiens] | 299 | 299 | 100% | 2e-105 | 99% | AAZ39780.1 |
| ☑ | beta-globin [Homo sapiens] | 299 | 299 | 100% | 2e-105 | 99% | ACU56984.1 |
| ☑ | hemoglobin beta chain [Homo sapiens] | 299 | 299 | 100% | 2e-105 | 99% | AAD19696.1 |
| ☑ | Chain B, Structure Of Haemoglobin In The Deoxy Quaternary State With Ligand Bound At The Alpha Haems | 298 | 298 | 99% | 3e-105 | 100% | 1COH_B |
| ☑ | hemoglobin beta subunit variant [Homo sapiens] | 298 | 298 | 100% | 4e-105 | 99% | AAF00489.1 |
| ☑ | Chain B, Human Hemoglobin D Los Angeles: Crystal Structure | 298 | 298 | 99% | 6e-105 | 99% | 2YRS_B |
| ☑ | Chain B, Structure Of Aquomet Hemoglobin Bristol-alesha Alphawtbetav67m | 297 | 297 | 99% | 8e-105 | 99% | 4MQI_B |
| ☑ | Chain B, High-Resolution X-Ray Study Of Deoxy Recombinant Human Hemoglobins Synthesized From Beta-Globins Having Mutated Amino Termini | 297 | 297 | 99% | 8e-105 | 99% | 1DXU_B |
| ☑ | Chain B, Analysis Of The Crystal Structure, Molecular Modeling And Infrared Spectroscopy Of The Distal Beta-Heme Pocket Valine67(E11)-Threonine Mutation Of Hemoglobin | 297 | 297 | 99% | 9e-105 | 99% | 1HDB_B |
| ☑ | Chain B, High-resolution X-ray Study Of Deoxy Recombinant Human Hemoglobins Synthesized From Beta-globins Having Mutated Amino Termini | 297 | 297 | 98% | 9e-105 | 100% | 1DXV_B |
| ☑ | Chain B, Crystal Structure Of Deoxygenated Hemoglobin In Complex With An Allosteric Effector And Nitric Oxide | 297 | 297 | 98% | 1e-104 | 100% | 5E29_B |
| ☑ | Chain C, Room Temperature Time-Of-Flight Neutron Diffraction Study Of Deoxy Human Normal Adult Hemoglobin | 297 | 297 | 98% | 1e-104 | 100% | 3KMF_C |
| ☑ | mutant beta-globin [Homo sapiens] | 297 | 297 | 100% | 1e-104 | 99% | AAL68978.1 |
| ☑ | Chain B, Crystal Structure Of Human Hemoglobin E At 1.73 A Resolution | 297 | 297 | 99% | 1e-104 | 99% | 1NQP_B |

# BLAST Results

- Scroll down to **Alignments** to get **details** for single alignments



Theresa Schwarz

# BLAST programs

- What other BLAST programs can we use?

**Basic Local Alignment Search Tool**

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. Learn more

**N E W S** — **October 26th NCBI Minute**

NCBI staff will introduce two new BLAST databases: the RefSeq Representative Genomes database and the Model Organisms or Landmark protein database. Fri, 07 Oct 2016 18:00:00 EST

More BLAST news...

**Web BLAST**

**Nucleotide BLAST**
nucleotide ▶ nucleotide

**blastx**
translated nucleotide ▶ protein

**tblastn**
protein ▶ translated nucleotide

**Protein BLAST**
protein ▶ protein

- Nucleotide BLAST = blastn

- Protein BLAST = blastp

- **blastx**

- **tblastn**

Theresa Schwarz

# BLAST programs

| Program | Input | | Database |
|---------|-------|---|----------|
| blastn | nt | ➡ | nt |
| blastp | protein | ➡ | protein |
| blastx | nt | ➡ | protein |
| tblastn | protein | ➡ | nt |

- When using blastx the input is a **nucleotide sequence**

- Then the program translates this sequence into a **protein sequence**

- Since the program does not know where the translation starts there are 6 possibilities

```
       5' CAT CAA
         5' ATC AAC
           5' TCA ACT


5' CATCAACTACAACTCCAAAGACACCCTTACACATCAACAAACCTACCCAC 3'
3' GTAGTTGATGTTGAGGTTTCTGTGGGAATGTGTAGTTGTTTGGATGGGTG 5'

                                      5' GTG GGT
                                        5' TGG GTA
                                          5' GGG TAG
```

Theresa Schwarz

# BLAST programs



| Program | Input | | Database | |
|---------|-------|---|----------|---|
| blastn | nt | → | nt | |
| blastp | protein | → | protein | |
| blastx | nt ← | → | protein | Search protein database using a translated nt query |
| tblastn | protein | → → | nt | Search translated nt database using a protein query |

Theresa Schwarz

❑ **Genome browsers**

➢ **UCSC**
➢ **ENSEMBL**

❑ **BLAST**

➢ **Pairwise alignments**
➢ **Database alignments**
➢ **Primer-BLAST**

# Primer-BLAST at NCBI

With **Primer-BLAST** you can check your primers that you designed to use them in a PCR

- do they amplify the desired product

- do they also bind to other regions in the given template

Go to BLAST (NCBI)

*https://blast.ncbi.nlm.nih.gov/Blast.cgi*

Under '**Specialized searches**' you can find **Primer-BLAST**



Theresa Schwarz

# Primer-BLAST at NCBI



- Primers designed to amplify **hemoglobin subunit beta** of **401bp**

- Enter the **primer sequences** (You can download the sequences from MOODLE)

- You can select a **minimum** and **maximum product size** …

Theresa Schwarz

# Primer-BLAST at NCBI



- As template I planned to use **human gDNA**

- Select the **database** and **organism** you want to check

- **Get Primers**

Theresa Schwarz

# Primer-BLAST at NCBI

## Output

**General information about your primers**

### Primer pair 1

| | Sequence (5'->3') | Length | Tm | GC% | Self complementarity | Self 3' complementarity |
|---|---|---|---|---|---|---|
| **Forward primer** | TCTGTCCACTCCTGATGCTG | 20 | 59.10 | 55.00 | 2.00 | 1.00 |
| **Reverse primer** | AAAAATTGCGGAGAAGAAAA | 21 | 53.08 | 28.57 | 4.00 | 0.00 |

### Products on target templates

>NC_000011.10 Homo sapiens chromosome 11, GRCh38.p7 Primary Assembly

```
product length = 401
Features associated with this product:
   hemoglobin subunit beta

Forward primer   1        TCTGTCCACTCCTGATGCTG   20
Template         5226748  ....................   5226729

Reverse primer   1        AAAAATTGCGGAGAAGAAAA   21
Template         5226348  ....................   5226368
```

Product length of the amplicon = **401bp**
The product is **hemoglobin subunit beta**
Is this the region we expected? ➔ **YES**

```
product length = 868
Features associated with this product:
   protocadherin Fat 3 isoform X1

   protocadherin Fat 3 isoform X5

Forward primer   1         TCTGTCCACTCCTGATGCTG   20
Template         92718795  .G.A....G.........T    92718814

Reverse primer   1         AAAAATTGCGGAGAAGAAAA   21
Template         92719662  ....TA.TT...........   92719642
```

Do we want to amplify this region? ➔ **NO**
The primers can also anneal to a different region in our genome. However, there are some mismatches in the annealing sequence.
**Is this now a problem for our PCR?**
**Do we have to design new primers?**

# Mismatches in primer sequence

```
5'-ATCGGGGCCCAC-3'
    |||||||||||
3'-TAGCCCCGGGTATATTTAAACGGGCCCAAATTTTCTAGGTACCTAGGGTAAATTCG-5'
```

- In most cases, a mismatch at the 3'-end of the primer (where the polymerase attaches the next nucleotide) impairs the elongation. (Because the primer-template complex is destabilized at a crucial position)

```
5'-CTCGAAGCCCAT-3'
    |||  |||||||
3'-TAGCCCCGGGTATATTTAAACGGGCCCAAATTTTCTAGGTACCTAGGGTAAATTCG-5'
```

- In most cases, one or more mismatches in the middle or the 5'-end of the primer do not affect the binding of the polymerase and the DNA can be amplified.

Theresa Schwarz

## What kind of **mismatches** do we have in our second BLAST-result?



```
product length = 868
Features associated with this product:
    protocadherin Fat 3 precursor

    protocadherin Fat 3 isoform X1

Forward primer    1          TCTGTCCACTCCTGATGCTG    20
Template          92718795   .G.A....G..........T    92718814

Reverse primer    1          AAAAATTGCGGAGAAGAAAAA   21
Template          92719662   ....TA.TT............   92719642
```

- In most cases, a mismatch at the 3'-end of the primer impairs the elongation.

- In most cases, one or more mismatches in the middle or the 5'-end of the primer do NOT affect PCR.

**→ It is unlikely that the primers will amplify the WRONG product, therefore we don't have to design new ones.**

Theresa Schwarz

- **QUESTIONS?**


- Please, download **Exercises #3** from MOODLE and upload until next Monday 8:00 a.m.


# GOOD LUCK!

Theresa Schwarz