# "Bringing Bioinformatics into the Biology Classroom"

Marie-Claude Blatter & Patricia Palagi

Marie-Claude.Blatter@isb-sib.ch

**SIB Swiss Institute of Bioinformatics**

Global Organisation for Bioinformatics Learning, Education & Training

# SIB Swiss Institute of Bioinformatics

- academic, non-profit foundation established in 1998

- coordinates **research** and **education** in bioinformatics throughout Switzerland

- provides high quality **bioinformatics services** to the national and international research community.

- helps shape the future of life sciences

- 52 groups, more than 600 scientists
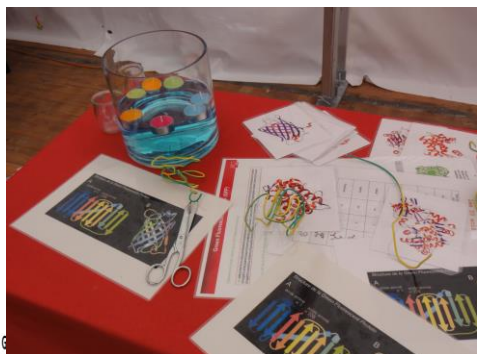
- GOBLET  member from the beginning

SIB
Swiss Institute of
Bioinformatics

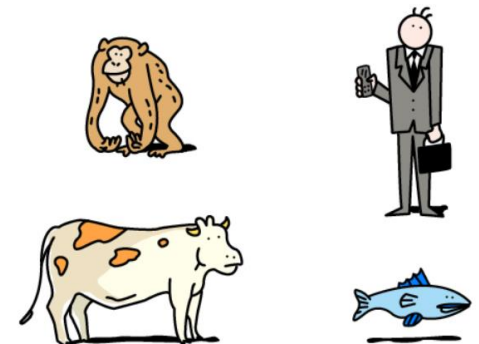# SIB outreach activities
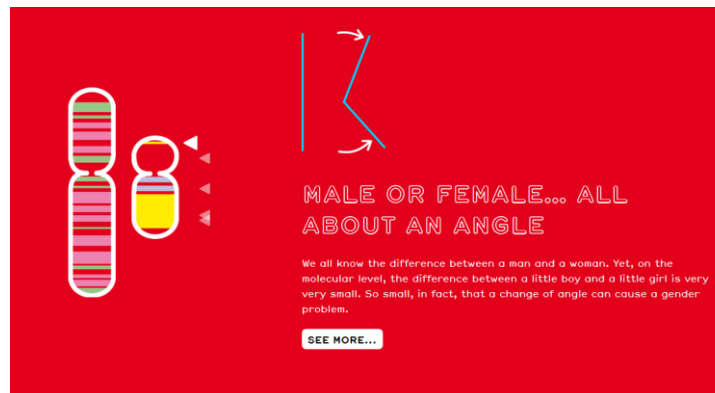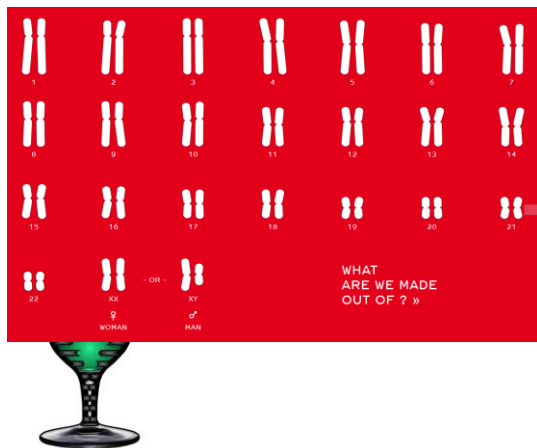## education & public at large

- Since 2000 (science fairs, electronic publication, exhibitions, hands on workshops, high school (HS) teacher continuing education training, etc.)
- Collaboration with public laboratories, didacticians and HS teachers
- **www.chromosomewalk.ch** (EN, FR, DE)
- **Protein Spotlight** (EN): http://web.expasy.org/spotlight/
- '**Ateliers de bioinformatique**' (FR): http://education.expasy.org/bioinformatique/
- *New project*: Drug Design and personalized medicine

# www.ChromosomeWalk.ch

- **a saunter along the human genome**

- …take a walk and discover the world of **genes**, **proteins** and **bioinformatics**.

- quizzes, videos, links to databases and bioinformatics tools



WHAT
ARE WE MADE
OUT OF ? »

22    XX    - OR -    XY
      ♀              ♂
      WOMAN          MAN

MALE OR FEMALE... ALL
ABOUT AN ANGLE

We all know the difference between a man and a woman. Yet, on the molecular level, the difference between a little boy and a little girl is very very small. So small, in fact, that a change of angle can cause a gender problem.

SEE MORE...

1    L W P P P P A R A F V N
2    L W G P D P A S A F V N
3    L W G P D P A A A F V N
4    F S G P G T S Y A A A N

# proteinspotlight

> ONE MONTH, ONE PROTEIN <

- http://web.expasy.org/spotlight/
- above 160 articles, informal tone (V. Gerritsen)

«*The German inventor Nikolaus Otto is credited with having invented the first automobile engine that ran on alcohol.*»

## Moving Forward

September 2014

Nature's imagination seems endless, and so is Man's. For as long as humans have existed, they have twisted Nature to meet their own needs. Wood has been used to keep them warm. Whale oil has been used to make light. Water has been harnessed to make electricity. And when the era of bio-engineering developed, it was not long before scientists found ways to tinker with an organism's genome for the benefits of mankind...

# 'Ateliers de Bioinformatique'

**http://education.expasy.org/bioinformatique/** (FR)

# Understanding a genetic disease thanks to Bioinformatics

**http://education.expasy.org/bioinformatique/Diabetes.html**
(Atelier 7: L'insuline de A à Z ; English version)

additional documents are available here:
**http://education.expasy.org/cours/Toronto**

SIB
Swiss Institute of
Bioinformatics

# Bioinformatics

This field of science designs software tools and databases for research in the life sciences.

Today, the quantity of biological data accumulated by laboratories is daunting.

As a result, the data can no longer be dealt with 'manually' and **bioinformatics** has become an essential ally.

http://www.chromosomewalk.ch/en/we-need-bioinformatics-to/

# Context

In a **special case** of type I diabetes
described in a Norwegian family,
a genetic variation has been found,
leading to the production of inactive insulin

http://www.ncbi.nlm.nih.gov/pubmed/18192540

## Mutations in the insulin gene can cause MODY and autoantibody-negative type 1 diabetes.

Molven A[1], Ringdal M, Nordbø AM, Raeder H, Støy J, Lipkind GM, Steiner DF, Philipson LH, Bergmann I, Aarskog D, Undlien DE, Joner G, Søvik O; Norwegian Childhood Diabetes Study Group, Bell GI, Njølstad PR.

⊕ Collaborators (27)

⊕ Author information

**Abstract**

**OBJECTIVE:** Mutations in the insulin (INS) gene can cause neonatal diabetes. We hypothesized that mutations in INS could also cause maturity-onset diabetes of the young (MODY) and autoantibody-negative type 1 diabetes.

**RESEARCH DESIGN AND METHODS:** We screened INS in 62 probands with MODY, 30 probands with suspected MODY, and 223 subjects from the Norwegian Childhood Diabetes Registry selected on the basis of autoantibody negativity or family history of diabetes.

**RESULTS:** Among the MODY patients, we identified the INS mutation c.137G>A (R46Q) in a proband, his diabetic father, and a paternal aunt. They were diagnosed with diabetes at 20, 18, and 17 years of age, respectively, and are treated with small doses of insulin or diet only. In type 1 diabetic patients, we found the INS mutation c.163C>T (R55C) in a girl who at 10 years of age presented with ketoacidosis and insulin-dependent, GAD, and insulinoma-associated antigen-2 (IA-2) antibody-negative diabetes. Her mother had a de novo R55C mutation and was diagnosed with ketoacidosis and insulin-dependent diabetes at 13 years of age. Both had residual beta-cell function. The R46Q substitution changes an invariant arginine residue in position B22, which forms a hydrogen bond with the glutamate at A17, stabilizing the insulin molecule. The R55C substitution involves the first of the two arginine residues localized at the site of proteolytic processing between the B-chain and the C-peptide.

**CONCLUSIONS:** Our findings extend the phenotype of INS mutation carriers and suggest that INS screening is warranted not only in neonatal diabetes, but also in MODY and in selected cases of type 1 diabetes.

**Comment in**
Insulin mutations in diabetes: the clinical spectrum.  [Diabetes. 2008]

PMID: 18192540 [PubMed - indexed for MEDLINE]     **Free full text**

*This publication is not available as free 'full text' in PubMed Central (PMC).*
*For full text:*
*http://education.expasy.org/cours/Toronto/*

Global Organisation for Bioinformatics Learning, Education & Training

**Activity 1: The insulin gene and the human genome**
(Genome browser (USCS), BLAT)

**Activity 2: Comparing DNA sequences - Diagnosing a rare genetic disease**
(alignment tool, database dbSNP)

**Activity 3: DNA translation -> protein**
(translate tool)

**Activity 4: 3D structure of insulin**
(database PDB, 3D visualization tool)

**Activity 5:  Is insulin specific to humans?**
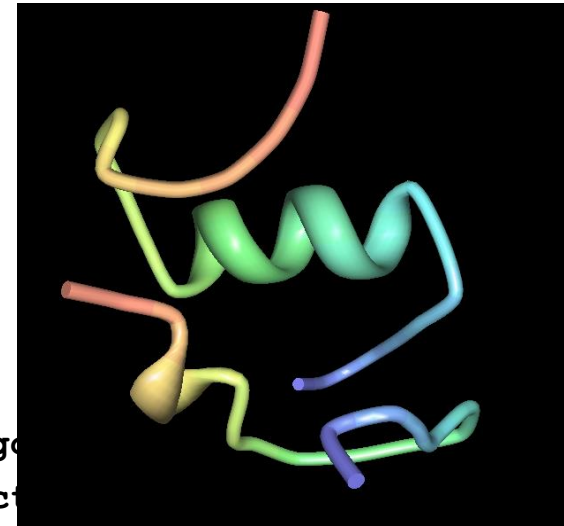(similarity search (BLAST), dabatase UniProtKB, alignment tool)

**&**

Global Organisation for Bioinformatics Learning, Education & Training

**biological function**



**protein**
**(amino acid)**

MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVE
ALYLVCGERGFFYTPKT**C**REAEDLQVGQVELGGGPGA
GSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

atggccctgtggatgcgcctcct...cgtgctggcg...
ccagccgcagcctttgtga...gtgcggct...

**gene**
**(DNA; nucleic acid)**

tagtgtgcggggaacgag...acacacccaagacccg**t**cgggaggcagaggac
tgcaggtggggcaggtgga...tgggcggggggccctggtgcaggcagcctgcagcccttg
gccctggaggggtccctgcagaagcgtggcattgtggaacaatgctgtaccagcatctgc
tccctctaccagctggagaactactgcaactag

*mutation*

**genome**

chr11:2,181,082-2,182,201 1,120 bp. | insulin | go

chr11 (p15.5) | p15.4 | p13 | p12 | q14.1 | q21 | q22.3 | 23.3 | 25

# Activity 1

## Activity 1: The insulin gene and the human genome

Bellow is a piece of the gene sequence that encodes for the insulin protein ('*wild sequence*')...

cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc

**Question:**

- **On which of our 23 chromosomes is this gene located?**

**Bioinformatics approach:**

**Use the tool 'BLAT'**
*Technical information*: 'BLAT' is a bioinformatics tool for comparing a DNA sequence against the whole genome sequence (the human genome has 3 billion nucleotides). *If the sequence exists, BLAT finds the sequence that is the most similar in just a few seconds. It's a bit like a small* **'google map'** *of the human genome.*

\* Copy the DNA sequence and paste it in the tool **'BLAT'**
\* Click on 'submit'
\* In the page 'BLAT Search Result': choose the best score and click 'browser'

- **On which chromosome is located the gene for insulin?**
- **What are the beginning and end positions of the sequence on the chromosome (nucleotide 'numbers')?**

- **For fun: write a random sequence (about 30 letters), always using the 4-letter alphabet (a, t, g, c) into 'BLAT': can you find it in the genome?**

**http://education.expasy.org/bioinformatique/Diabetes.html**

**Bioinformatics approach:**

**Use the tool 'BLAT' @ USCS**

_Technical information_: 'BLAT' is a bioinformatics tool for comparing a DNA sequence against a whole genome sequence.

_If the sequence exists, BLAT finds the sequence that is the most similar in just a few seconds. It's a bit like a small **'google map'** of the human genome._

**1. Google: look for 'BLAT UCSC'**

**2. Choose the latest release of the human genome (GRCh38)**

Genomes  Genome Browser  Tools  Mirrors  Downloads  My Data  Help  About Us

**Human BLAT Search**

## BLAT Search Genome

| Genome: | Assembly: | Query type: | Sort output: | Output type: |
|---|---|---|---|---|
| Human ▼ | Dec. 2013 (GRCh38/hg38) ▼ | BLAT's guess ▼ | query,score ▼ | hyperlink ▼ |

cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc

**3. Click on submit**

submit   I'm feeling lucky   clear

Paste in a query sequence to find its location in the the genome. Multiple sequences may be

**Bioinformatics – Genome Browser (Blat USCS)**
**http://genome.ucsc.edu/cgi-bin/hgBlat**

Global Organisation for Bioinformatics Learning, Education & Training

SIB
Swiss Institute of
Bioinformatics

**Human BLAT Results**

## BLAT Search Results

*human genome (GRCh37)*

| ACTIONS | | QUERY | SCORE | START | END | QSIZE | IDENTITY | CHRO | STRAND | START | END | SPAN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| browser | details | YourSeq | 60 | 1 | 60 | 60 | 100.0% | 11 | − | 2182081 | 2182140 | 60 |
| browser | details | YourSeq | 20 | 26 | 45 | 60 | 100.0% | 9 | + | 138953442 | 138953461 | 20 |

**At each genome release the positions may change**

**Human BLAT Results**

## BLAT Search Results

*human genome (GRCh38)*

| ACTIONS | | QUERY | SCORE | START | END | QSIZE | IDENTITY | CHRO | STRAND | START | END | SPAN |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| browser | details | YourSeq | 60 | 1 | 60 | 60 | 100.0% | 11 | − | 2160851 | 2160910 | 60 |
| browser | details | YourSeq | 20 | 26 | 45 | 60 | 100.0% | 9 | + | 136061596 | 136061615 | 20 |

*By default, choose the best score*
*Click on 'browser'*

Global Organisation for Bioinformatics Learning, Education & Training

SIB
Swiss Institute of Bioinformatics

The insulin DNA sequence is located on **chromosome 11 (11p15.5)**
(positions: 2,160,851-2'160,910 (GRCh38))



Zoom out 100 x

# UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

move `<<<` `<<` `<` `>` `>>` `>>>` zoom in `1.5x` `3x` `10x` `base` zoom out `1.5x` `3x` `10x` `100x`

chr11:2,157,881-2,163,880  6,000 bp.  [ enter position, gene symbol or search terms ]  `go`



**The insulin gene consists of 3 exons and 2 introns**

**A readthrough transcript INS-IGF2 involves INS and IGF2 genes (neighboring genes)**

TCCACACGCTCCTTGCGGCCTCATGGGTGTAGGGTCCAGCCCCACAGGGTCGGTGGGTCTCTCCCCGTG
GCAGAGACGAGAGAGTGTAGAAATAAAGACACAAGACAAAGAGATTAAAAAA
GGACCACTACCACCAATGCGCGGAGACCGGTAGTGGCCCCGAATGTCTGGCT
TACAAGGCAAAAGGGGCAGGGTAAAGAGTGTGAGTCATCTCCAATGATAGAT
TGTCCACTGGACAGGGGGCCCTTCCCTGCCTGGCAGCCGAGGCAGAGAGGGA
ATAGCTTACGCCATTATTTTTGTATATTAGAGACGTTTAGTACTTTCACTAA
AAGGCAGAGCCAGGTGCACAGGATGGAACATGAAGGAGGACTAGGAGCGTGA
ACAGGGAGACGGTTAGGCCTCCGGATAACTGCGGGCAGGTCTGACTGATGTC
GGAGGAGCAGAGTCTTCTCTAAACTCCCCCGGGGAAAGGGAGCCCCTCCTTT
GGGTGTTTTTCCTTGACACTTACGCTACCGCTAGACCACGGTCCGCTTGGCA
GCTGGCATCACCGCTAGACCAAGGAGCCCTCTAGTGGCCTTGTCCGGGCATA
TGTCTTCTGGTCACTCCTCACTATGTCCCCTCAGCTCCTATCTCTGTATGGCCTGGTTTTTCCTAGTTT
TGATTATAGAGCGAGGATTGTTATAATATTGGAATAAAGAGTAATTGCTACAAACTAATGATTAATGAT
TTCATATATAATCATATCTAAGATCTATATCTGGTGTAACTATTTTTATTTTATATTTTATTATACTGG
ACAGCTCGTGTCCTCAGTCTCTTGCCTCGGCACCTGGGTGGCTTGCCGCCCACAATGGGCAGCTTATTC
TCAGGGAAGGCCTTTGTCTCCACACCTGTGGGTGAAGACCATCGGGATGCTTTGCCTTCAACAGGCAAG
CCAACAATTCACCTTCACTTCCCTCCCTCCAGGAACACCAGCTCCCAGCTCAGAGTCATCGGCCTCGCT
ACAGGGACGTCACACTACCCGCTCTGTGGGGGGCATCGTGTGGTCTGGACTTGCTGAGCAGAAAGTAGC
GCTGCCCTCAACACCTCCCTAGAGCATCTGCGAGCCGAACACCTGGGGCCCACGCCTCCGGCACGTCTA
GGACCCAGTGGTCCATCCCTTCCCAAGCACAAGGCAAGTGGCTACCTCAGTCCCTTCCTCCACGAAGAA
GAGGCACGATGCCTAGTGCTGTAGGTCCCATGTTATTTGGGAAGCAACTTTTGCCCTATTTGGAAGTGC

http://www.ncbi.nlm.nih.gov/nuccore/71514639?report=fasta

Click on the link
This is part of the sequence of Craig Venter chromosome 11
(GenBank database;  3'852'046 bp over 135'006'516 bp)

**Select short sequences (about 40 bp)**
**and check with BLAT that they are located on chromosome 11 ….**

Global Organisation for Bioinformatics Learning, Education & Training

Write a random sequence (about 30 letters), always using the 4-letter alphabet (a, t, g, c)

| | Genomes | Genome Browser | Tools | Mirrors | Downloads | My Data | Help | About Us |

**Human BLAT Search**

# BLAT Search Genome

| Genome: | Assembly: | Query type: | Sort output: | Output type: |
|---|---|---|---|---|
| Human ▼ | Dec. 2013 (GRCh38/hg38) ▼ | BLAT's guess ▼ | query,score ▼ | hyperlink ▼ |

atgctagatcatgctagctagtcgatgctattgatcgatttagatcgat

*Click on 'submit'*

| submit | I'm feeling lucky | clear |

It is virtually impossible to match a randomly typed sequence (ATGC, n=30) on the human genome sequence, even on «junk» DNA regions (Application: PCR and primer selection)

Randomly selected letters (i.e. n=5) rarely create a correct word….

SIB
Swiss Institute of Bioinformatics

Write a random sequence (about 30 letters),
always using the 4-letter alphabet (a, t, g, c)

| 🏠 | Genomes | Genome Browser | Tools | Mirrors | Downloads | My Data | Help | About Us |

**White rhinoceros BLAT Search**

# BLAT Search Genome

| Genome: | Assembly: | Query type: | Sort output: | Output type: |
|---|---|---|---|---|
| White rhinoceros ▾ | May 2012 (CerSim Sim1.0/cerSim1) ▾ | BLAT's guess ▾ | query,score ▾ | hyperlink ▾ |

atgctagatcatgctagctagtcgatgctattgatcgatttagatcgat

*Select another
genome*

submit    I'm feeling lucky    clear

SIB
Swiss Institute of
Bioinformatics

# Activity 2

**Activity 2: Comparing DNA sequences - Diagnosing a rare genetic disease**

In 2008, scientists studied a Norwegian family in which several members had diabetes (type I or type II) (Molven et al., 2008).

All diabetic type I members of the family carry the same rare variation in the gene which encodes for insulin.

Here is the family's pedigree (phenotype and family relationship):



Family T1D-N781

■ affected male

□ male

● affected female

○ female

**Is this baby diabetic (type I) ?**

**Question:**

- Is this baby diabetic?

**http://education.expasy.org/bioinformatique/Diabetes.html**

**Biological context: Human genome & variations**

The human genome = a text of 3'000'000'000 pb

= a reference sequence

All the differences (also called *variations, variants, Single Nucleotide Polymorphisms (SNPs), ~mutations, ...*) between human subjects are described on the basis of this 'text'

**DNA sequence variants**
*1 in 1000 nts vary in two randomly selected genomes*

We are all different from each other

WE ARE ALL UNIQUE

~ 3.3 millions of SNPs between 2 individuals

~10 millions of SNPs in the human population
  neutral (the majority)
  associated with a particular phenotype…
  associated with a predisposition
  associated with a genetic disease

~10 – 30  new mutation at each new generation

**DNA sequence variants**

Neutral
Predisposing
Pathogenic

Pr S.E. Antonarakis (UNIGE)

# Biological context: Human genome & variations

1. ☐ Comprehensive characterization of **human genome** variation by high coverage whole-**genome** sequencing of **forty four Caucasians**.

Shen H, Li J, Zhang J, Xu C, Jiang Y, Wu Z, Zhao F, Liao L, Chen J, Lin Y, Tian Q, Papasian CJ, Deng HW.

PLoS One. 2013;8(4):e59494. doi: 10.1371/journal.pone.0059494. Epub 2013 Apr 5.

PMID: 23577066 [PubMed - indexed for MEDLINE]    **Free PMC Article**

Related citations

*Publication (free full text) in PubMed Central (PMC @NCBI) are freely available for everyone*

"On average, each individual genome carried ~3.3 million SNPs and ~492,000 indels/block substitutions, including approximately 179 variants that were predicted to cause loss of function of the gene products. "

PMID: 23577066

SIB
Swiss Institute of Bioinformatics

Global Organisation for Bioinformatics Learning, Education & Training

# SNPs are stored in the dbSNP database



**Bioinformatics – Biological Database: dbSNP @ NCBI**
**http://www.ncbi.nlm.nih.gov/SNP/**

Family T1D-N781

Question: is the baby diabetic ?

To answer this question, researchers extracted
DNA from 8 members of the Norwegian family
and sequenced part of the gene that encodes
for insulin.

*Compare these    sequences, and locate the common variation for diabetes.*

**'Paper and pencil' approach:**

... You can do it **manually which will help you better understand the principle of sequence comparison and alignment.**

Take into account all the given clues and play with our strips of DNA sequences...

- 8 family members - 4 DNA sequences - one allele
- 8 family members - 1 DNA sequence - two alleles
- 8 family members - 2 DNA sequences - two alleles (not easy)

**Bioinformatics approach:**

Build an alignment of these 8 sequences using a bioinformatics tool and look out for the common variation among those with diabetes

\* Copy these 8 sequences (including the lines starting with '>1') and paste them into the align tool

\* Click on the *Run Align* button.

\* On the results page, on the lefthand column 'Highlight': select 'Similarity'

Global Organisation for Bioinformatics Learning, Education & Training

SIB
Swiss Institute of
Bioinformatics

**'Paper and pencil' approach:**

... You can do it **manually which will help you better understand the principle of sequence comparison and alignment.** Take into account all the given clues and play with our strips of DNA sequences...

- 8 family members - 4 DNA sequences - one allele
- 8 family members - 1 DNA sequence - two alleles
- 8 family members - 2 DNA sequences - two alleles (not easy)

**2 different DNA sequences (INS gene)**

**8 subjects (same family)**

>1.1
tagtgtgcggggaacgaggcttcttctaca

>1.3
cacccaagacccgccgggaggcagagg

>1.2
Tagtgtgcggggaacgaggcttcttctaca

>1.4
cacccaagacctgccgggaggcagaggacc

>2.1
tagtgtgcggggaacgaggcttcttctaca

>2.3
cacccaagacccgccgggaggcagaggacc

**2 alleles (maternal / paternal)**

>2.2
tagtgtgcggggaacgaggcttcttctaca

cacccaagacccgccgggaggcagaggacc

>3.1
tagtgtgcggggaacgaggcttcttctaca

>3.3
cacccaagacccgccgggaggcagaggacc

>3.2
tagtgtgcggggaacgaggcttcttctaca

>3.4
cacccaagacctgccgggaggcagaggacc

>4.1
tagtgtgcggggaacgaggcttcttctaca

>4.3
cacccaagacccgccgggaggcagaggacc

>4.2
tagtgtgcggggaacgaggcttcttctaca

>4.4
cacccaagacccgccgggaggcagaggacc

>5.1
tagtgtgcggggaacgaggcttcttctaca

>5.3
cacccaagacccgccgggaggcagaggacc

>5.2
tagtgtgcggggaacgaggcttcttctaca

>5.4
cacccaagacccgccgggaggcagaggacc

# Sequence 1



```
1.1    tagtgtgcggggaacgaggcttcttctaca
1.2    tagtgtgcggggaacgaggcttcttctaca
2.1    tagtgtgcggggaacgaggcttcttctaca
2.2    tagtgtgcggggaacgaggcttcttctaca
3.1    tagtgtgcggggaacgaggcttcttctaca
3.2    tagtgtgcggggaacgaggcttcttctaca
4.1    tagtgtgcggggaacgaggcttcttctaca
4.2    tagtgtgcggggaacgaggcttcttctaca
5.1    tagtgtgcggggaacgaggcttcttctaca
5.2    tagtgtgcggggaacgaggcttcttctaca
6.1    tagtgtgcggggaacgaggcttcttctaca
6.2    tagtgtgcggggaacgaggcttcttctaca
7.1    tagtgtgcggagaacgaggcttcttctaca
7.2    tagtgtgcggagaacgaggcttcttctaca
8.1    tagtgtgcggggaacgaggcttcttctaca
8.2    tagtgtgcggggaacgaggcttcttctaca
```

**Where are the differences ?**

http://www2.grifil.com/album.html

**Sequence 1**

Family T1D-N781



| | |
|---|---|
| 1.1 | tagtgtgcggggaa |
| 1.2 | tagtgtgcggggaa |
| 2.1 | tagtgtgcggggaa |
| 2.2 | tagtgtgcggggaa |
| 3.1 | tagtgtgcggggaa |
| 3.2 | tagtgtgcggggaa |
| 4.1 | tagtgtgcggggaa |
| 4.2 | tagtgtgcggggaa |
| 5.1 | tagtgtgcggggaa |
| 5.2 | tagtgtgcggggaa |
| 6.1 | tagtgtgcggggaacgaggcttcttctaca |
| 6.2 | tagtgtgcggggaacgaggcttcttctaca |
| 7.1 | tagtgtgcggagaacgaggcttcttctaca |
| 7.2 | tagtgtgcggagaacgaggcttcttctaca |
| 8.1 | tagtgtgcggggaacgaggcttcttctaca |
| 8.2 | tagtgtgcggggaacgaggcttcttctaca |

**Sequence 1**

Family T1D-N781

| | |
|---|---|
| 1.1 | tagtgtgcggggaa |
| 1.2 | tagtgtgcggggaa |
| 2.1 | tagtgtgcggggaa |
| 2.2 | tagtgtgcggggaa |
| 3.1 | tagtgtgcggggaa |
| 3.2 | tagtgtgcggggaa |
| 4.1 | tagtgtgcggggaa |
| 4.2 | tagtgtgcggggaa |
| 5.1 | tagtgtgcggggaa |
| 5.2 | tagtgtgcggggaa |
| 6.1 | tagtgtgcggggaacgaggcttcttctaca |
| 6.2 | tagtgtgcggggaacgaggcttcttctaca |
| 7.1 | tagtgtgcggagaacgaggcttcttctaca |
| 7.2 | tagtgtgcggagaacgaggcttcttctaca |
| 8.1 | tagtgtgcggggaacgaggcttcttctaca |
| 8.2 | tagtgtgcggggaacgaggcttcttctaca |



**The SNP g -> a (homozygous; subject 7) is not associated with diabetes (neutral)**

**Sequence 2**


Family T1D-N781

```
11   cacccaagacccgccgggaggcagaggacc
12   cacccaagacctgccgggaggcagaggacc
21   cacccaagacccgccgggaggcagaggacc
22   cacccaagacccgccgggaggcagaggacc
31   cacccaagacccgccgggaggcagaggacc
32   cacccaagacctgccgggaggcagaggacc
41   cacccaagacccgccgggaggcagaggacc
42   cacccaagacccgccgggaggcagaggacc
51   cacccaagacccgccgggaggcagaggacc
52   cacccaagacccgccgggaggcagaggacc
61   cacccaagacccgccgggaggcagaggacc
62   cacccaagacccgccgggaggcagaggacc
71   cacccaagacccgccgggaggcagaggacc
72   cacccaagacccgccgggaggcagaggacc
81   cacccaagacccgccgggaggcagaggacc
82   cacccaagacccgccgggaggcagaggacc
```

**Where are the differences ?**

**Sequence 2**

Family T1D-N781



| | |
|---|---|
| 11 | cacccaagacccgccggg |
| 12 | cacccaagacctgccggg |
| 21 | cacccaagacccgccggg |
| 22 | cacccaagacccgccggg |
| 31 | cacccaagacccgccggg |
| 32 | cacccaagacctgccggg |
| 41 | cacccaagacccgccggg |
| 42 | cacccaagacccgccggg |
| 51 | cacccaagacccgccggg |
| 52 | cacccaagacccgccgggaggcagaggacc |
| 61 | cacccaagacccgccgggaggcagaggacc |
| 62 | cacccaagacccgccgggaggcagaggacc |
| 71 | cacccaagacccgccgggaggcagaggacc |
| 72 | cacccaagacccgccgggaggcagaggacc |
| 81 | cacccaagacccgccgggaggcagaggacc |
| 82 | cacccaagacccgccgggaggcagaggacc |

# Sequence 2

Family T1D-N781



| | |
|---|---|
| 11 | cacccaagacccgccggg |
| 12 | cacccaagacctgccggg |
| 21 | cacccaagacccgccggg |
| 22 | cacccaagacccgccggg |
| 31 | cacccaagacccgccggg |
| 32 | cacccaagacctgccggg |
| 41 | cacccaagacccgccggg |
| 42 | cacccaagacccgccggg |
| 51 | cacccaagacccgccggg |
| 52 | cacccaagacccgccggggaggcagaggacc |
| 61 | cacccaagacccgccggggaggcagaggacc |
| 62 | cacccaagacccgccggggaggcagaggacc |
| 71 | cacccaagacccgccggggaggcagaggacc |
| 72 | cacccaagacccgccggggaggcagaggacc |
| 81 | cacccaagacccgccggggaggcagaggacc |
| 82 | cacccaagacccgccggggaggcagaggacc |

**ANSWER: The SNP c -> t is present in subjects 3 and 1 (heterozygous) and is associated with Type I Diabetes ('all type I diabetic members carry the same variation in the INS gene')**

**Bioinformatics approach:**

Build an alignment of these 8 sequences using a bioinformatics tool and look out for the common variation among those with diabetes

* Copy these 8 sequences (including the lines starting with '>1') and paste them into the align tool
* Click on the *Run Align* button.
* On the results page, on the lefthand column 'Highlight': select 'Similarity'

**Bioinformatics – Alignment tool (UniProt)**
**http://www.uniprot.org/align/**

*This tool is used to align protein sequences, but it can also properly align short DNA sequences*

SIB
Swiss Institute of
Bioinformatics

Family T1D-N781

## Question: is the baby diabetic ?

To answer this question, researchers extracted DNA from 8 members of the Norwegian family and sequenced part of the gene that encodes for insulin.

```
>1
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc
>2
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
>3
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc
>4
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
>5
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
>6
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
>7
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggagaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
>8
cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc
tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc
```

## Where are the differences ?



HE' HE' HE

http://www2.grifil.com/album.html

# http://www.uniprot.org/align/

# Align

🛒 Basket **7** ▾

## Display  All None

☑ ALIGNMENT

☐ TREE

☐ RESULT INFO

## Highlight

**Annotation**

**Amino acid properties**

☐ Similarity
☐ Hydrophobic
☐ Negative
☐ Positive
☐ Aliphatic
☐ Tiny

⬇ **Download**    ❮ **Edit and resubmit**

## Alignment

🖨 How to print an alignment in color

```
1   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
2   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
3   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
4   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
5   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
6   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
7   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
8   1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
        ************************************************************

1   61  tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc  120
2   61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
3   61  tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc  120
4   61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
5   61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
6   61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
7   61  tagtgtgcggagaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
8   61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc  120
        **********.*************************************  ******************
```

Captu

Global Organisation for Bioinformatics Learning, Education & Training

SIB
Swiss Institute of Bioinformatics

**Bioinformatics – Alignment tool (UniProt)**
**http://www.uniprot.org/align/**

(only one DNA sequence (allele) per subject)

```
1  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
2  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
3  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
4  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
5  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
6  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
7  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
8  1   cagccgcagcctttgtgaaccaacacctgtgcggctcacacctggtggaagctctctacc   60
       ************************************************************

1  61  tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc   120
2  61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
3  61  tagtgtgcggggaacgaggcttcttctacacacccaagacctgccgggaggcagaggacc   120
4  61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
5  61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
6  61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
7  61  tagtgtgcggagaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
8  61  tagtgtgcggggaacgaggcttcttctacacacccaagacccgccgggaggcagaggacc   120
       *********.********************************** ****************
```

Family T1D-N781

- The subject **(1)** with the **c -> t** mutation (heterozygous mutation) is a girl who presented type I diabetes at the early age of 10.
- The girl's mother **(3)** has type I diabetes that was diagnosed when she was 13. Currently, she is being treated with insulin. She also carries the heterozygous mutation **c -> t** .
- The girl's maternal grandfather **(6)** has type 2 diabetes, which was diagnosed at the age of 40. He is currently being treated with insulin. Neither he nor the healthy maternal grandmother carry mutations.
- **-> Thus, the girl's mother is carrying a *de novo* c -> t  mutation, which must be a germline mutation since it has been inherited by her daughter.**

(Molven et al., 2008)

# Biology / Statistics / Bioinformatics

Molven et al., 2008

"We screened INS (*gene coding for insulin*) in 92 probands, and 223 subjects from the Norwegian Childhood **Diabetes** Registry selected on the basis of autoantibody negativity or family history of diabetes."

Identification of a rare genetic variation - rs121908261 – in the human INS gene which is the cause of type I diabetes in a Norwegian family

SIB
Swiss Institute of Bioinformatics

# Insulin gene
## (4992 pb ; chromosome 11)

In red, part of the DNA sequence translated into the protein sequence

rs121908261 *[Homo sapiens]*

1.

AGGCTTCTTCTACACACCCAAGACC [C/T] GCCGGGAGGCAGAGGACCTGCAGGG

Chromosome: 11:2160809
Gene: INS-IGF2 (GeneView) INS (GeneView)
Functional Consequence: missense,nc transcript variant
Clinical significance: Pathogenic
Validated: no info
HGVS: NC_000011.10:g.2160809G>A, NC_000011.9:g.2182039G>A, NG_007114.1:g.5386C>T, NM_000207.2:c.163C>T, NM_001042376.2:c.163C>T, NM_001185097.1:c.163C>T, NM_001185098.1:c.163C>T, NM_001291897.1:c.163C>T, NP_000198.1:p.Arg55Cys, NP_001035835.1:p.Arg55Cys, NP_001172026.1:p.Arg55Cys, NP_001172027.1:p.Arg55Cys, NP_001278826.1:p.Arg55Cys, NR_003512.3:n.222C>T

PubMed Varview Protein3D OMIM

rs121908261

c -> t

# Bioinformatics – Biological Database: dbSNP @ NCBI
## http://www.ncbi.nlm.nih.gov/SNP/

*Variant accession number in dbSNP*

---

**NCBI**    Resources ☑   How To ☑         Sign in to NCBI

**dbSNP**     [SNP ▾]   rs121908261        ⊗   **Search**

Save search   Advanced            Help

**Display Settings:** ☑ Summary         Send to: ☑

☐ rs121908261 *[Homo sapiens]*

1.

     AGGCTTCTTCTACACACCCAAGACC[C/T]GCCGGGAGGCAGAGGACCTGCAGGG

| | |
|---|---|
| Chromosome: | 11:2160809 |
| Gene: | INS-IGF2 (GeneView) INS (GeneView) |
| Functional Consequence: | missense,nc transcript variant |
| Clinical significance: | Pathogenic |
| Validated: | no info |
| HGVS: | NC_000011.10:g.2160809G>A, NC_000011.9:g.2182039G>A, NG_007114.1:g.5386C>T, NM_000207.2:c.163C>T, NM_001042376.2:c.163C>T, NM_001185097.1:c.163C>T, NM_001185098.1:c.163C>T, NM_001291897.1:c.163C>T, NP_000198.1:p.Arg55Cys, NP_001035835.1:p.Arg55Cys, NP_001172026.1:p.Arg55Cys, NP_001172027.1:p.Arg55Cys, NP_001278826.1:p.Arg55Cys, NR_003512.3:n.222C>T |

PubMed   Varview   Protein3D   OMIM

**Search details**

rs121908261[All Fields]

[Search]      See more...

**Recent activity**

Turn Off   Clear

🔍 rs121908261 (1)
               SNP

📄 Mutations in the insulin gene can cause MODY and autoantik PubMed

See more...

*Link to the publication of Molven et al (2008)*

**Bioinformatics – Biological Database: dbSNP @ NCBI**
**http://www.ncbi.nlm.nih.gov/SNP/**

Look for all the SNPs located within human INS gene

# Activity 3

## Activity 3: DNA translation -> protein

*Check the effect of the mutation 'R55C'...*

Like all proteins, insulin is composed of a sequence of amino acids. The order of the amino acids is determined by the nucleic acid sequence of the insulin gene. 3 letters of DNA (codon) correspond to one amino acid (symbolized by letters: K for lysine, M for methionine, etc.).

This is a piece of the DNA sequence of the normal insulin gene.

aag acc cgc cgg gag

This is a piece of the DNA sequence of the insulin gene with the c -> t variation, associated with type I diabetes.

aag acc tgc cgg gag

*Question:*

- Does the **c->t** mutation change the amino acid sequence of insulin?
- Does the **aag -> aaa mutation** change the amino acid sequence of insulin?

http://education.expasy.org/bioinformatique/Diabetes.html

Global Organisation for Bioinformatics Learning, Education & Training

couples de nucléotides possibles :

ADN

Transcription

codon-stop    codon-stop

couples de nucléotides possibles :

ARNm

Traduction

20 acides aminés possibles

Protéine

N1 Méthionine — Agrinine — Alanine C1

M    R    A

http://fr.wikipedia.org/wiki/Synth%C3%A8se_des_prot%C3%A9ines#mediaviewer/File:Proteines.png

## Activity 3: DNA translation -> protein

*Check the effect of the mutation 'R55C'...*

Like all proteins,insulin is composed of a sequence of amino acids. The order of the amino acids is determined by the nucleic acid sequence of the insulin gene.
3 letters of DNA (codon) correspond to one amino acid (symbolized by letters: K for lysine, M for methionine, etc.).

This is a piece of the DNA sequence of the normal insulin gene.

```
aag acc cgc cgg gag
```

This is a piece of the DNA sequence of the insulin gene with the c -> t variation, associated with type I diabetes.

```
aag acc tgc cgg gag
```

*Question:*

- **Does the c->t mutation change the amino acid sequence of insulin?**
- **Does the aag -> aaa mutation change the amino acid sequence of insulin?**

Cap

You could manually translate the nucleic acid sequences into amino acid sequences ('1 'letter code) using the genetic code below: :



Capture

SIB
Swiss Institute of
Bioinformatics

Global Organisation for Bioinformatics Learning, Education & Training

# Bioinformatics – Translate tool
## http://www.bioinformatics.org/sms2/translate.html



**Fasta format**

Translate results

>rf 1 normal
KTRRE

>rf 1 mutated
KTCRE

# Amino acid sequence of the 'normal' insulin

MALWMRLLPLLALLALWGPDPAAA**FVNQHLCGSHLVEALYLVCGERGFFYTPKT****R**REAEDLQVGQVELGGGPGAGSLQPLALEGSLQKR**GIVEQCCTSICSLYQLENYCN**

| Chain B | | Chain A |
|---|---|---|

**FVNQHLCGSHLVEALYLVCGERGFFYTPKT**          **GIVEQCCTSICSLYQLENYCN**

# Amino acid sequence of the 'mutated' insulin (variant c-> t; R 55 C)

MALWMRLLPLLALLALWGPDPAAA**FVNQHLCGSHLVEALYLVCGERGFFYTPKT****C**REAEDLQVGQVELGGGPGAGSLQPLALEGSLQKR**GIVEQCCTSICSLYQLENYCN**

**The mutation R -> C prevents insulin from being cut and thus from being biologically active**

## Mutations in the insulin gene can cause MODY and autoantibody-negative type 1 diabetes.

Molven A[1], Ringdal M, Nordbø AM, Raeder H, Støy J, Lipkind GM, Steiner DF, Philipson LH, Bergmann I, Aarskog D, Undlien DE, Joner G, Søvik O; Norwegian Childhood Diabetes Study Group, Bell GI, Njølstad PR.

⊕ Collaborators (27)

⊕ Author information

**Abstract**

**OBJECTIVE:** Mutations in the insulin (INS) gene can cause neonatal diabetes. We hypothesized that mutations in INS could also cause maturity-onset diabetes of the young (MODY) and autoantibody-negative type 1 diabetes.

**RESEARCH DESIGN AND METHODS:** We screened INS in 62 probands with MODY, 30 probands with suspected MODY, and 223 subjects from the Norwegian Childhood Diabetes Registry selected on the basis of autoantibody negativity or family history of diabetes.

**RESULTS:** Among the MODY patients, we identified the INS mutation c.137G>A (R46Q) in a proband, his diabetic father, and a paternal aunt. They were diagnosed with diabetes at 20, 18, and 17 years of age, respectively, and are treated with small doses of insulin or diet only. In type 1 diabetic patients, we found the INS mutation c.163C>T (R55C) in a girl who at 10 years of age presented with ketoacidosis and insulin-dependent, GAD, and insulinoma-associated antigen-2 (IA-2) antibody-negative diabetes. Her mother had a de novo R55C mutation and was diagnosed with ketoacidosis and insulin-dependent diabetes at 13 years of age. Both had residual beta-cell function. The R46Q substitution changes an invariant arginine residue in position B22, which forms a hydrogen bond with the glutamate at A17, stabilizing the insulin molecule. The R55C substitution involves the first of the two arginine residues localized at the site of proteolytic processing between the B-chain and the C-peptide.

**CONCLUSIONS:** Our findings extend the phenotype of INS mutation carriers and suggest that INS screening is warranted not only in neonatal diabetes, but also in MODY and in selected cases of type 1 diabetes.

**Comment in**
Insulin mutations in diabetes: the clinical spectrum. [Diabetes. 2008]

PMID: 18192540 [PubMed - indexed for MEDLINE]     **Free full text**

*This publication is not available as free 'full text' in PubMed Central (PMC).*
*For full text:*
*http://education.expasy.org/cours/Toronto/*

Global Organisation for Bioinformatics Learning, Education & Training

## Activity 4: 3D structure of insulin



Since 1958, researchers have been able to crystallize proteins and then 'take a picture' of them by using X-rays. The results of these experiments are then analyzed using bioinformatic programs which make it possible to **view** the 3D structure of proteins such as insulin.

### View the 3D structure of insulin

* Go to the PDB entry **2LWZ**
* Select the 3D viewer 'Protein Workshop'.
   *A Jmol application will be launched and you will be asked to accept it. Jmol is a viewer for chemical structures in 3D.*
   *The Jmol application requires Java to be installed in your computer. Both programs are free.*
* In Shortcuts: Recolor the backbone 'By compound' - and then look at the positions of the different amino acids (mouse over)
* In Tools: 'Surfaces' play with the Transparency slider
* In Tools: 'Visibility', 'atoms and bonds', click on 'Chain A: Insulin" and see the atoms (balls and sticks) that are displayed
* In Option: Reset - to go back to the original image

For fun, here are the raw experimental data, the spatial coordinates(X, Y, Z) of every atom in each amino of insulin (search ATOM in the page)

*Note: There is no 3D structure data for insulin with the R55C mutation.*

Global Organisation for Bioinformatics Learning, Education & Training

**http://education.expasy.org/bioinformatique/Diabetes.html**

http://www.pdb.org/pdb/explore/explore.do?structureId=2LWZ



….. requires Java to be installed in your computer.

# Bioinformatics – Insulin 3D structure in PDB database  (2LWZ) (Protein Workshop)
## http://www.pdb.org/pdb/explore/explore.do?structureId=2LWZ



**Visualization tool: Protein Workshop**

Protein Workshop: in Tools: 'Surfaces' play with the Transparency slider

Protein workshop:
In Shortcuts: Recolor the backbone 'By compound' - and then look at
the positions of the different amino acids (mouse over)

Protein workshop: In Tools: 'Visibility', 'atoms and bonds', click on 'Chain A: Insulin" and see the atoms (balls and sticks) that are displayed

# Activity 5

**Activity 5: Is insulin specific to humans?**

**BLAST**

This is the full sequence of human insulin amino acid (in UniProtKB):

MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

*Question:*

- **Is this protein specific to humans?**

**Bioinformatics approach:**
Do a 'BLAST' against a database of proteins called UniProtKB

<u>Technical information</u>: *BLAST is a bioinformatics tool that compares the sequence of a protein with millions of other sequences contained in a database. If they exist, it finds those that resemble a given sequence the most within a few seconds. We can thus find out quickly whether a protein exists in a given species, or not.*

* Copy the sequence and paste it into the tool 'BLAST'
* Select '**Target Database = UniProtKB/Swiss-Prot**'
* Click on 'Run BLAST'
* Check the conservation of amino acids ('View alignment') and the conservation of the disulfide bonds ('Highlight' 'disulfide bond', when available)
* Search on Google for images corresponding to the different Latin names of the species (example 'Octodon degus')

According to wikipedia, insulin is a very old protein that may have originated one billion years ago.
Apart from animals, insulin-like proteins are also known to exist in Fungi and Protista kingdoms.

* Select '**Target Database = ...Nematoda**' or '**Target Database = ...Arthropoda**'

**http://education.expasy.org/bioinformatique/Diabetes.html**

# Bioinformatics – BLAST Similarity search tool
## www.uniprot.org/blast/

**'reviewed' entries (UniProtKB/Swiss-Prot section) are manually reviewed**
**'unreviewed' entries (UniProtKB/TrEMBL section) are automatically annotated**

# BLAST

🛒 Basket 7 ▾

Identity %
100   80   60   40   20   0

## Filter by[i]

**Reviewed (33)** ✖
Swiss-Prot

With 3D structure (3)
Proteomes (27)

### Organisms
Fruit fly (3)
BOMMO (24)
AGRCO (2)
SAMCY (3)
LOCMI (1)

### Map To
UniProtKB
UniRef
UniParc

❮ Edit and resubmit   Order by: Score ▾   Limit to sequences from organism: All ▾

## Overview

Capture

Show all 33

Bombyxin A-3 (Bombyx mori)

| P33721 | Bombyxin B-1 homolog (Samia cynthia) | 43.0% |
| P33722 | Bombyxin B-2 homolog (Samia cynthia) | 45.0% |
| Q9VT52 | Probable insulin-like peptide 3 (Drosophila melanogaster) | 26.0% |
| P26733 | Bombyxin B-1 (Bombyx mori) | 30.0% |
| P26741 | Bombyxin B-7 (Bombyx mori) | 30.0% |
| P26729 | Bombyxin A-6 (Bombyx mori) | 30.0% |

## Alignments

✏ Columns   🔍 BLAST   ≡ Align   ⬇ Download   🛒 Add to basket

◀ 1 to 25 of 33 ▶   Show 25 ▾

**How similar are the human and drosophila sequences ?**

>sp|Q9VT52|INSL3_DROME Probable insulin-like peptide 3 OS=Drosophila melanogaster GN=Ilp3 PE=2 SV=2
MGIEMRCQDRRILLPSLLLLILMIGGVQATMKLCGRKLPETLSKLCVYGFNAMTKRTLDP
VNFNQIDGFEDRSLLERLLSDSSVQMLKTRRLRDGVFDECCLKSCTMDEVLRYCAAKPRT

>sp|P01308|INS_HUMAN Insulin OS=Homo sapiens GN=INS PE=1 SV=1
MALWMRLLPLLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAED
LQVGQVELGGGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYCN

# Bioinformatics – alignment tool
## www.uniprot.org/align/

# Many thanks to all of you

# and
# to Michelle Brazas