IPMB
Institut für Pharmazie und
Molekulare Biotechnologie

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

# Help!! What is ....

## ... an ORF ??

# ORF = open reading frame

open reading frame



**Standard genetic code**
start codon = ATG (AUG)
stop codon = TAG / TAA / TGA

ORF = sequence stretch between
start codon and stop codon

IPMB
Institut für Pharmazie und
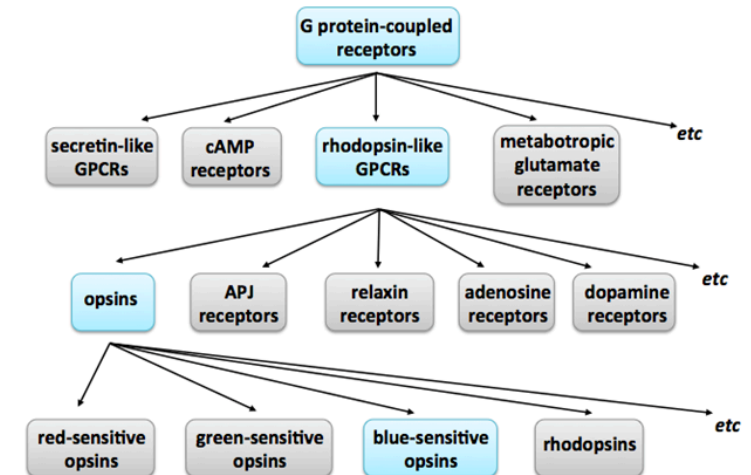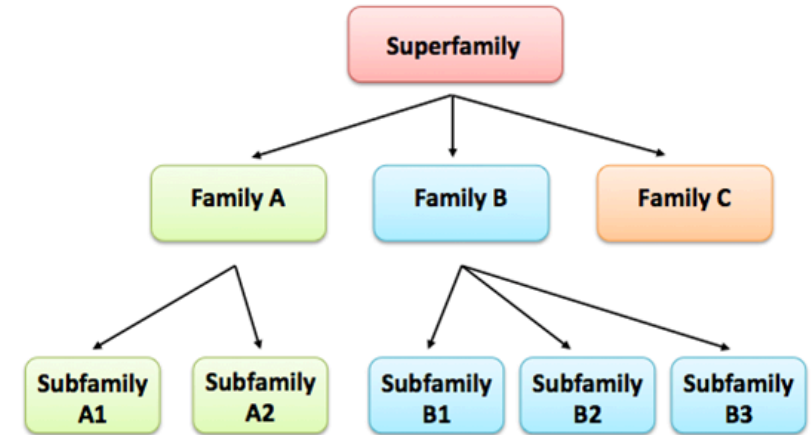Molekulare Biotechnologie

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

# Help!! What is ....

# ... a protein (super)family ??

# Families, superfamilies,...

- **Protein family** = group of proteins that share a **common evolutionary origin**, reflected by their related functions and similarities in sequence or structure.
  - ◉ superfamily = large group of distantly related proteins
  - ◉ subfamily = small group of closely related proteins

- **Protein families** are often arranged into hierarchies, with proteins that share a common ancestor subdivided into smaller, more closely related groups.
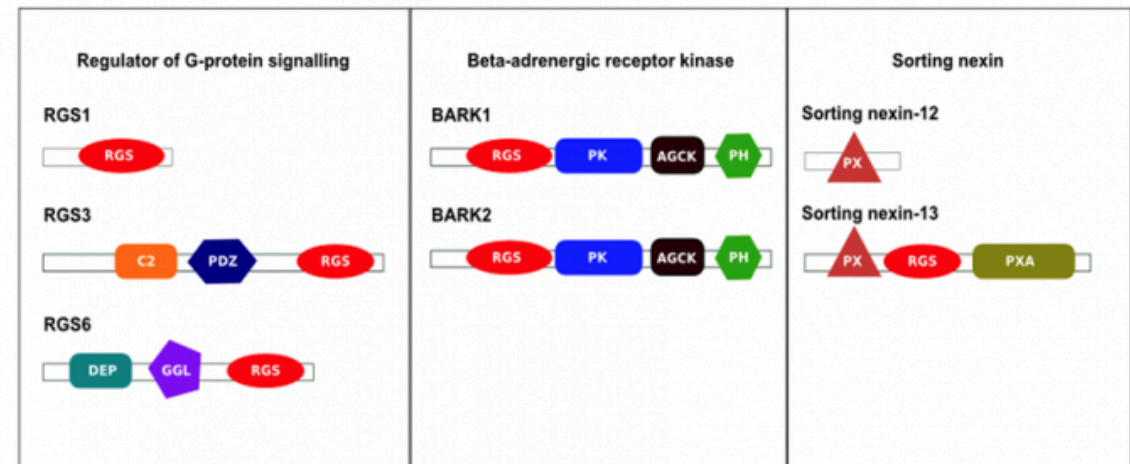
# Protein domains

- Protein domains = **functional** and/or **structural** units in a protein

- Protein usually contain **several protein domains**



SH3 domain



RGS family     beta-adrenergic receptor kinase family     Sorting nexin family

RGS = Regulator of G-protein signalling

IPMB

Institut für Pharmazie und
Molekulare Biotechnologie

UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

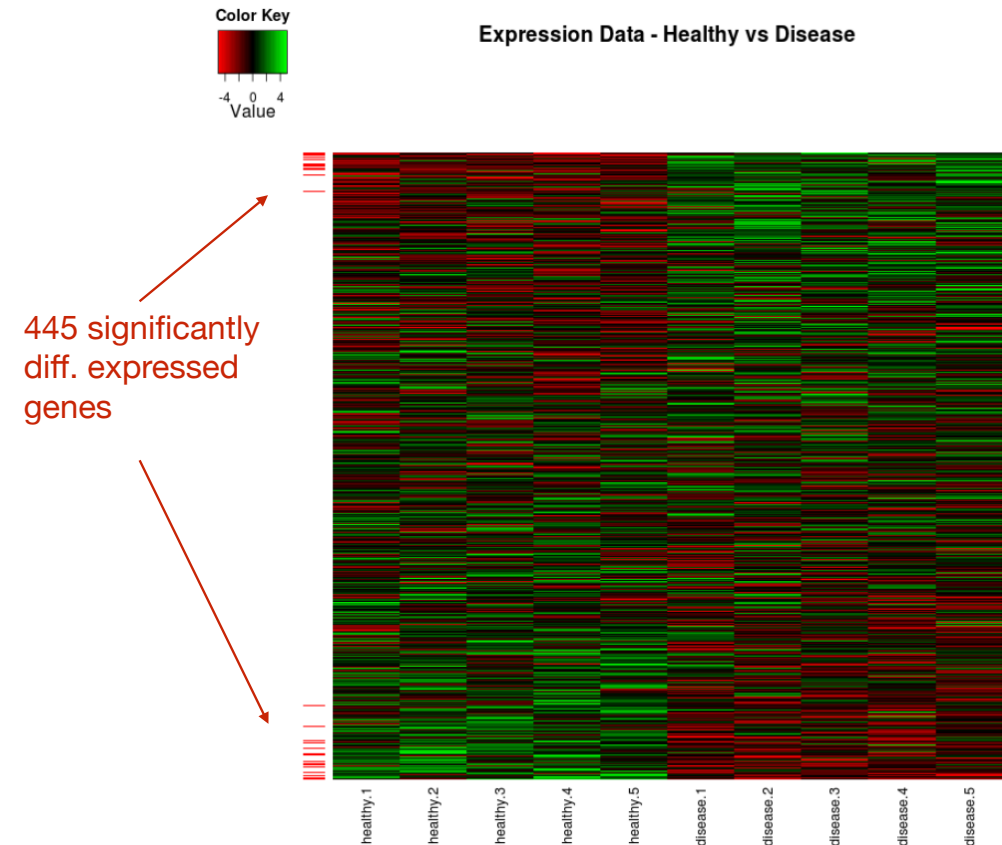# Help!! What is ....

# ... an E-value ??

# Fake gene expression data

- Finding differentially expressed genes between healthy and disease patients

- t-test with $\alpha = 5\%$

- $H_0$: non-significant expression difference between the two groups

**This dataset contains only random numbers**
**→ $H_0$ holds for all 10.000 "genes"**
**→ all the 445 genes are false-positives**



445 significantly diff. expressed genes

# Multiple testing

- **Significance level** $\alpha$: level at which to **reject** $(p < \alpha)$ or **accept** $(p > \alpha)$ the Null hypothesis

- **P-value**: probability to observe a more extreme effect if H0 is true ("risk of a false-positive by random chance")

- **E-value**: expected number of false-positive events when N tests are performed

$$E = p \cdot N$$

# Help!! What is ....

## ... BLAST ??

# Why sequence alignments ?

```
>Protein sequence
MLCPISGWAIYSKDNSIRIGSKGDVFVIREPFISCSHLECRTFFLTQGALLNDKHSN
GTVKDRSPYRTLMSCPVGEAPSPYNSRFESVAWSASACHDGISWLTIGISGPDNGAV
AVLKYNGIITDTIKSWRNNTLRTQESECACVNGSCFTVMTDGPSNEQASYKIFKIEK
```

- Open questions
  - ◉ **Homologues**: *are there related sequences in other organisms?*
  - ◉ **Function**: *possible biological/molecular/enzymatic function?*
  - ◉ **Origin**: *from which organisms / clade? (Example: Metagenomics)*
- **Global** (= whole sequence) / **local** (= parts of the sequence) comparisons

house    house

hanse    haus

# Lost in translation...

house     house

| ||     | ||

hanse     haus

3 Matches     3 Matches

2 Mismatches     1 Mismatch

1 Ins./Del. (Indel)

# Lost in translation...

**house**          **house**

|   | |          | | |

**hanse**          **haus**

3 Matches          3 Matches

2 Mismatches          1 Mismatch

1 Ins./Del. (Indel)

- which comparison is better?
  → *Scoring-Method*
- **Score** should take into account...

  ◉ Matches (+)

  ◉ Mismatches (-)

    ▸ Vokal/Vokal oder Kons./Kons. (-)

    ▸ Vokal/Kons. (--)

  ◉ Insertion/Deletion (-)

# Lost in translation...

**house**
|   | |
**hanse**

3 x (+**1**) - 1 x (-**1**)
-1 x (-**2**)
= **0**

**house**
|   | |
**haus**

3 x (+**1**) - 1 x (-**1**)
-1 x (-**1**)
= **+1**

- which comparison is better?
  → *Scoring-Method*
- **Score** should take into account...

  ◉ Matches (+)

  ◉ Mismatches (-)

    ▸ Vokal/Vokal oder Kons./Kons. (-)

    ▸ Vokal/Kons. (--)

  ◉ Insertion/Deletion (-)

# Scoring Verfahren

- **Matches/mismatches**
  subtitution matrix: values represent frequencies of observed substitutions in homologous sequences (ex. BLOSUM62)



**Substitution matrix BLOSUM62**

- Gaps
  mostly affine cost

  ‣ gap opening (O)
  ‣ gap extensions (E)

```
AITP--------VPQ
AVTPQSLPCSSLQQ
```

indels I = 7

  ‣ opening: O = -11

  ‣ extension: E = -2

here: $-11 + 7 \times (-2) = -25$

# Which is the best alignment?

- Given 2 sequences, which is the best alignment?

$$\overset{m}{\longleftrightarrow} \qquad \overset{n}{\longleftrightarrow}$$

`AVTPQSLPCSSLQQ`     `AITPVPQ`

```
AITP-------VPQ          AITP--VPQ-----
AVTPQSLPCSSLQQ          AVTPQSLPCSSLQQ
```

BLOSUM62 Matrix
O = -11; E = -2

$$S = -1 \qquad\qquad S = -12$$

$$\binom{m+n}{n} = \binom{21}{7} = \text{116.280 possible alignments...}$$

# Dynamical programming (DP)

- Dynamical programming allows to determine exactly the best alignment

- Alignment = path in the score matrix

- Best alignment is obtained by determinung at each step the best alignment

- Needleman & Wunsch = **global** alignment
  Smith & Watermann = **local** Alignment

**Complexity**

$$\mathcal{O}(m \cdot n)$$

|   |   | V | L | Q | S |
|---|---|---|---|---|---|
| I |   |   |   |   |   |
| V |   |   |   |   |   |
| Q |   |   |   |   |   |
| P |   |   |   |   |   |
| T |   |   |   |   |   |

V-L-QS        VLQ-S
IVQP-T        IVQPT

→ gap in 2nd sequence

↓ gap in 1st sequence

↘ (mis)match

# Needlman & Wunsch: globales Alignment

**1. Phase**
Füllen der Matrix

|  |  | V | L | Q | S |
|---|---|---|---|---|---|
| | | | | | |
| I | | | | | |
| V | | | | | |
| Q | | | | | |
| P | | | | | |
| T | | | | | |

Gap = -8

| 2 | -6 |
|---|---|
| -4 | |

↓ -6-8 = -14
→ -4-8 = -12
↘ **2-0 = 2**

**2. Phase**
Backtracking von unten rechts nach oben links

|  |  | V | L | Q | S |
|---|---|---|---|---|---|
| | 0 | –8 | –16 | –24 | –32 |
| I | –8 | 2 | –6 | –14 | –22 |
| V | –16 | –4 | 2 | –6 | –14 |
| Q | –24 | –12 | –6 | 7 | –1 |
| P | –32 | –20 | –14 | –1 | 7 |
| T | –40 | –28 | –22 | –9 | –1 |

```
V L Q S –
I V Q P T
```

# Smith Watermann: lokales Alignment

**1. Phase**
Füllen der Matrix (negative Werte werden durch Null ersetzt)

**2. Phase**
Backtracking vom **höchsten Wert** bis zur ersten Null

Gap = -8

|   |   | V | L | Q | S |
|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 2 | 0 | 0 | 0 |
| V | 0 | 4 | 2 | 0 | 0 |
| Q | 0 | 0 | 0 | 7 | 0 |
| P | 0 | 0 | 0 | 0 | 7 |
| T | 0 | 0 | 0 | 0 | 2 |

|   |   | V | L | Q | S |
|---|---|---|---|---|---|
|   | 0 | 0 | 0 | 0 | 0 |
| I | 0 | 2 | 0 | 0 | 0 |
| V | 0 | 4 | 2 | 0 | 0 |
| Q | 0 | 0 | 0 | 7 | 0 |
| P | 0 | 0 | 0 | 0 | 7 |
| T | 0 | 0 | 0 | 0 | 2 |

```
V L Q S
I V Q P
```

# Problem solved?

**Enumeration of all possible alignments**



$$\begin{pmatrix} m+n \\ n \end{pmatrix} = 10^{29}$$

**complexity**
**n= m = 50**

**Optimal alignment using DP**



$$m \cdot n = 2500$$

**Optimal alignment with database**



*N ~ 230 million sequences*

$$N \cdot m \cdot n = 5.8 \cdot 10^{11}$$

- **Problem**: DP alignment cannot be computed for all target sequences (too long!)
- **Solution**: select most promising sequences first... then do DP

# Heuristic = short-cut

house ⟷ DUDEN    *N* ~150.000 words

**all words starting with h**

house ⟷     *N* ~1.000 words

- **Advantage**: much faster!
- **Disadvantage**: maybe the right translation starts with a different letter ...

# BLAST: basic local alignment search tool

## *Heuristic*

1. homologous sequences share **very similar short words**

   (Protein: $k$=3; DNA: $k$=11)

2. these words reside in **longer homologous sequences without gaps**

   (HSP = high scoring pairs)

3. starting from HSP **longer alignments with gaps** can be obtained using DP.

   $\rightarrow$ final raw score $S$ depends on the subsitution matrix, number of matches / mismatches / gaps

BLAST 1.0 (Altschul et al., 1990)
Alignments **ohne** Gaps

BLAST 2.0 (Altschul et al., 1997)
Alignments **mit** Gaps

*Raw score obtained from the dynamical programming*

*Parameter (= black magic)*

$$S' = \frac{\lambda S - \ln K}{\ln 2} \text{ (bits)}$$

*Size of the database*

*E = expected number of false-positive in a database of the same size*

$$E = \frac{Q}{2^{S'}}$$

# NCBI BLAST

E-value = number of false-positives with equal scores

$S = 238$

$S' = 96.3$ bits

**RecName: Full=Non-symbiotic hemoglobin 0; AltName: Full=Non-vascular plant hemoglobin**

Sequence ID: Q9M630.1   Length: 180   Number of Matches: 1

Range 1: 24 to 169 GenPept   Graphics       ▼ Next Match   ▲ Previous Match

| Score | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|
| 96.3 bits(238) | 1e-24 | Compositional matrix adjust. | 58/146(40%) | 84/146(57%) | 5/146(3%) |

```
Query   3    FTEKQEALVNSSSQLFKQNPSNYSVLFYTIILQKAPTAKAMFSFLKDSA-GVVDSPKLGA   61
             ++++ E LV  S ++ K++      + F+  + + AP AKAM+SFL+DS       ++PK+
Sbjct   24   YSKENEQLVKQSWEILKKDAQRNGINFFRKVFEIAPGAKAMYSFLRDSTIPFEENPKVKN   83

Query   62   HAEKVFGMVRDSAVQLRATGEV-VLDGKD---GSIHIQKGVLDPHFVVVKEALLKTIKEA   117
             HA  VF M  D+AVQL    G   VL+ K     + H+  GV D  F +VKEA+L  I+
Sbjct   84   HARYVFMMTGDAAVQLGEKGAYQVLESKLQKLAATHVNAGVTDDQFEIVKEAILYAIEMG   143

Query   118  SGDKWSEELSAAWEVAYDGLATAIKA   143
              D WS EL +AW  AYD LA  +KA
Sbjct   144  VPDLWSPELKSAWGDAYDMLAEQVKA   169
```

# Help!! What is ....

## ... phylogenetic tree ??

# Phylogenetic tree

- Represent the evolutionary history of a set of sequences or organisms

- Beware the a tree built from a single gene can differ from the evolutionnary tree of the species!

- Trees are constructed based on multiple alignments, from which a distance matrix can be built

Bacteria

Zixibacteria
Cloacimonetes
Fibrobacteres
*Gemmatimonadetes*
WOR-3
TA06
Poribacteria
Latescibacteria
BRC1

Atribacteria
*Aquificae*
Calescamantes
Caldiserica
Dictyoglomi
Thermotogae
*Deinococcus-Therm.*
*Synergistetes*
*Fusobacteria*

*Actinobacteria*

*Armatimonadetes*

(*Tenericutes*)

*Chloroflexi*

*Firmicutes*

*Cyanobacteria*

Melainabacteria
RBX1
WOR1

Nomurabacteria
Kaiserbacteria
Adlerbacteria
Campbellbacteria

Giovannonibacteria
Wolfebacteria
Jorgensenbacteria

Azambacteria
Parcubacteria

Yanofskybacteria
Moranbacteria
Magasanikbacteria
Uhrbacteria
Falkowbacteria

Candidate
Phyla Radiation

SM2F11
Peregrinibacteria
Gracilibacteria BD1-5, GN02
Absconditabacteria SR1
Saccharibacteria
Berkelbacteria

Marinimicrobia
*Ignavibacteria*
Caldithrix

*Bacteroidetes*
Chlorobi

PVC
superphylum

*Planctomycetes*
Chlamydiae,
Lentisphaerae,
Verrucomicrobia

*Elusimicrobia*

Omnitrophica

Aminicentantes  Rokubacteria     NC10
*Acidobacteria*
Tectomicrobia, Modulibacteria
*Nitrospinae*
*Nitrospirae*
Dadabacteria
*Deltaproteobacteria*
(*Thermodesulfobacteria*)
Chrysiogenetes
*Deferribacteres*
Hydrogenedentes NKB19
*Spirochaetes*
TM6
*Epsilonproteobacteria*

Wirthbacteria

Woesebacteria
Shapirobacteria
Amesbacteria
Collierbacteria
Pacebacteria
Beckwithbacteria
Roizmanbacteria
Gottesmanbacteria
Levybacteria
Daviesbacteria
Curtissbacteria

Dojkabacteria WS6
CPR1
CPR3
Katanobacteria
WWE3

Microgenomates

*Alphaproteobacteria*

Zetaproteo.
Acidithiobacillia

*Betaproteobacteria*

0.4

*Gammaproteobacteria*

*Major lineages with isolated representative: italics*
Major lineage lacking isolated representative: ●

Eukaryotes

*Micrarchaeota*
Diapherotrites

Nanohaloarchaeota
Aenigmarchaeota
Parvarchaeota

Loki.
Thor.

DPANN

Pacearchaeota
Nanoarchaeota
Woesearchaeota

Altiarchaeales
Z7ME43
*Methanopyri*
*Methanococci*
*Hadesarchaea*
*Thermococci*
*Methanobacteria*
*Thermoplasmata*
Archaeoglobi
*Methanomicrobia*

Korarch.
Crenarch.
Bathyarc.
YNPFFA
Aigarch.

*Halobacteria*

TACK

*Thaumarchaeota*

Opisthokonta

Excavata

Archaeplastida

Chromalveolata

Amoebozoa

Archaea

Carl Herrmann

# Outgroup

- IN an unrooted tree, one cannot tell which is the evolutionary origin

- If you know that a group of sequences is more distant than the rest ('outgroup'), then the root of the tree can be set on the branch separating the outgroup from the rest!