# Gene Sequence Analysis

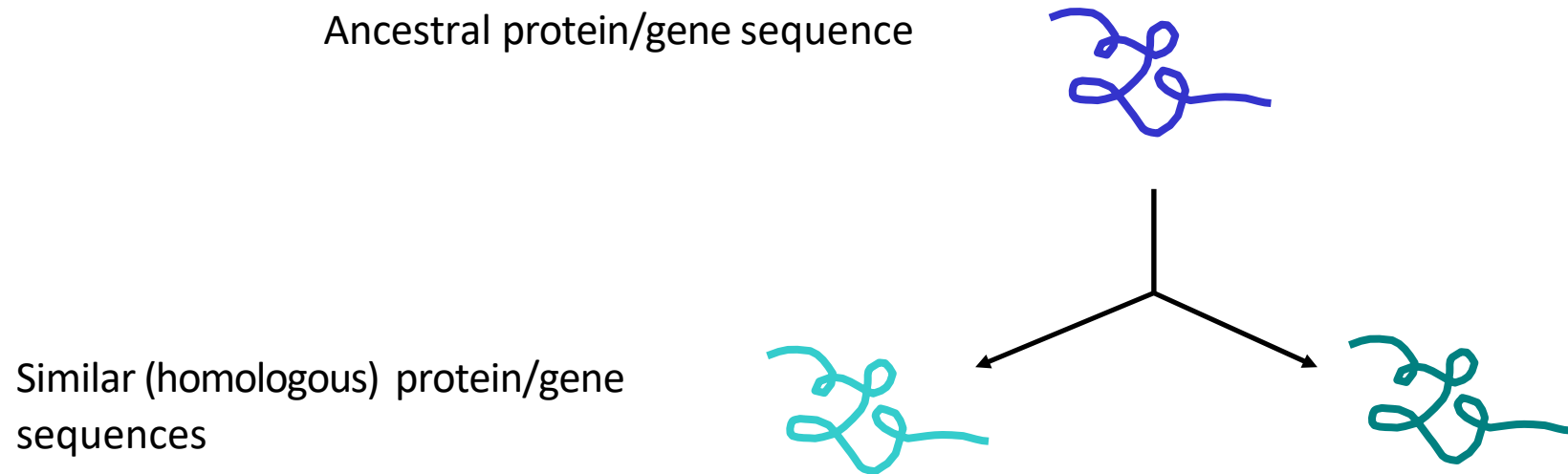# Lecture 3: **Multiple Alignment**

02/06/2024
Phuc Loi Luu, PhD
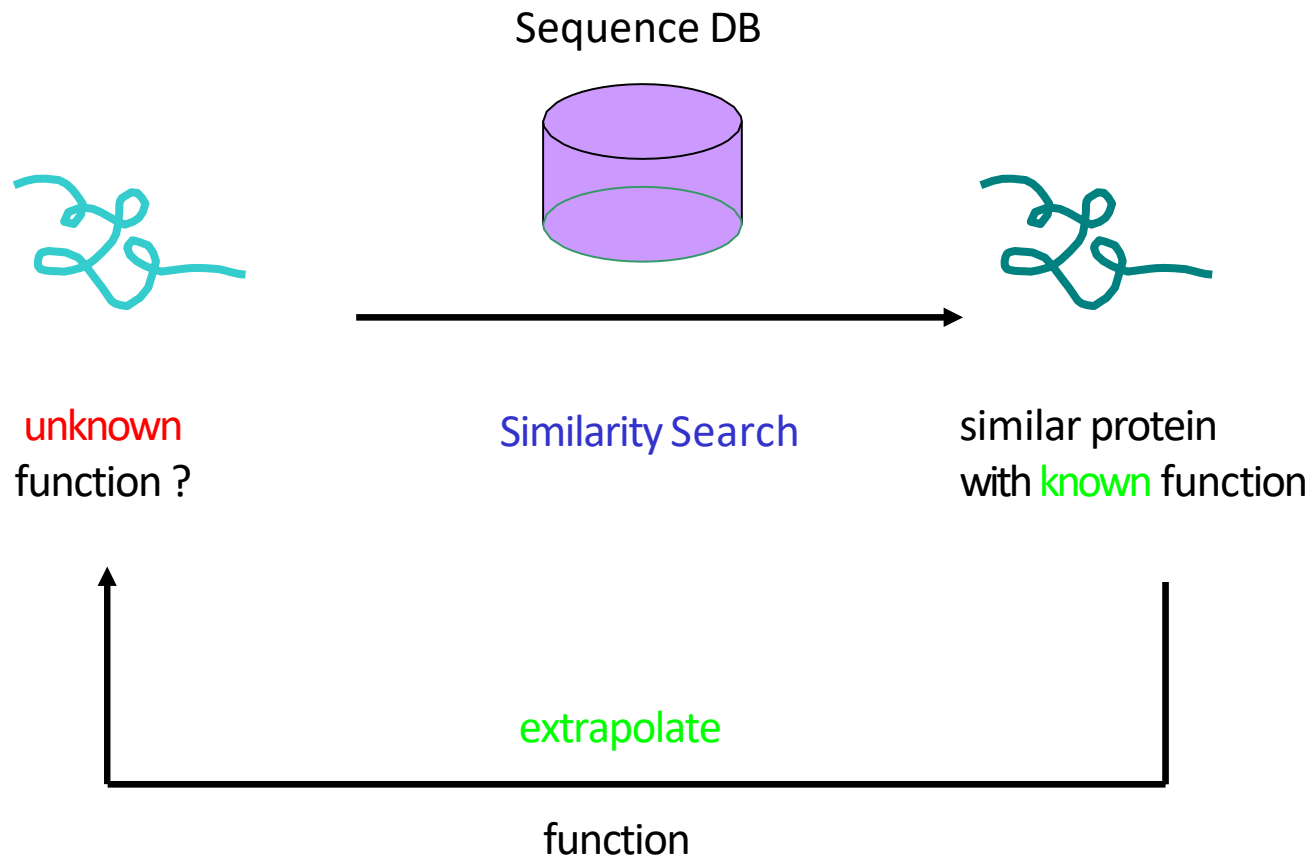Apdated from Dr. Morgan Langille

# Importance of Similarity

Ancestral protein/gene sequence

Similar (homologous) protein/gene sequences

Similar sequences: probably have the same ancestor, share the same structure, and have a similar biological function

# Importance of Similarity

Sequence DB



unknown
function ?

Similarity Search

similar protein
with known function

extrapolate

function

# Importance of Similarity

Rule-of-thumb:

If your sequences are more than **100 amino acids** long (or 100 nucleotides long) you can considered them as homologues if **25%** of the **aa** are identical (**70%** of **nucleotide** for DNA). Below this value you enter the twilight zone.
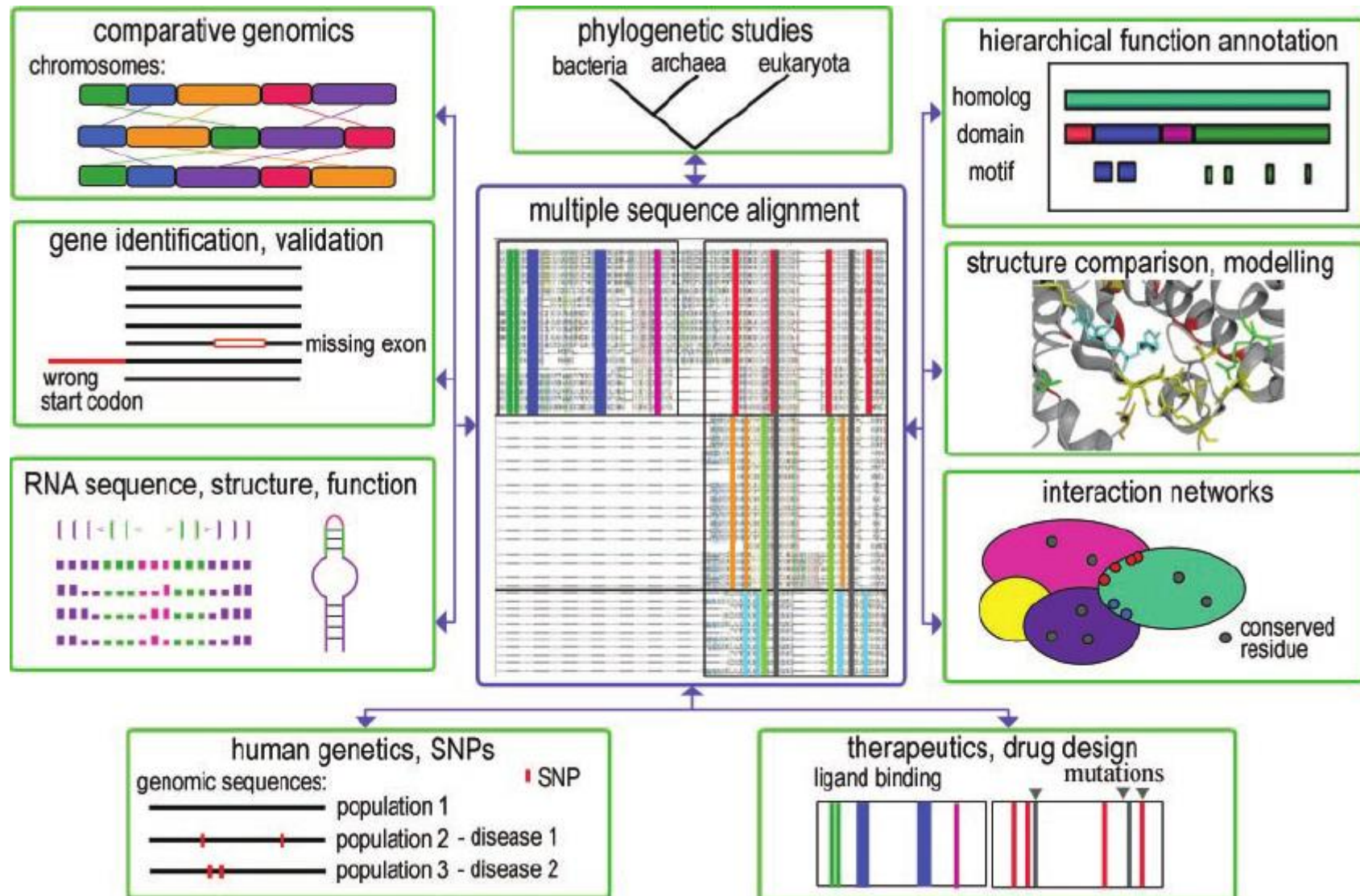
Twilight zone = protein sequence similarity between ~0-20% identity:
is not statistically significant, i.e. could have arisen by chance.

Beware:
- E-value (*Expectation value*)
- length of the segments similar between the two sequences
- The number of insertions/deletions

# Examples of molecular biology applications (shown in green boxes) that rely on multiple sequence alignments

# Outline

- What is a multiple alignment?

- Why do we need a multiple alignment?

- Characters of evolution

- How to create multiple alignments?

- Editing Alignments

- Viewing Alignments

# What is a multiple alignment?

- Simply an alignment of more than 2 sequences

- Sequences are aligned globally (end to end)

- Multiple Sequence Alignment (MSA) programs try to insert gaps in the sequences so that *homologous characters* are aligned
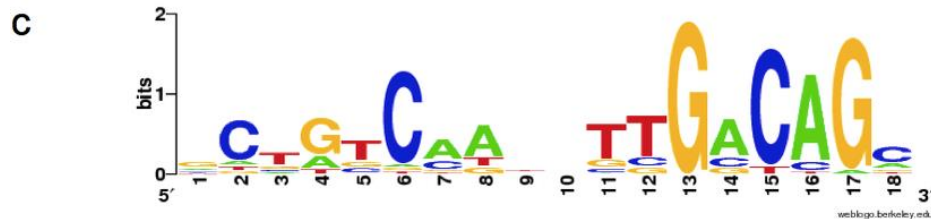
# What is a multiple alignment?

# What to do with a MSA?

- To look for sites of interest/conservation within a gene (motifs, binding sites, etc.)

- To build phylogenetic trees

- To identify positive selection

- To look for sites of interest/conservation within a gene (motifs, binding sites, etc.)



**(a)** Alignment including the complete version of the afe-box (30 nucleotides); **(b)** Alignment including the central palindromic module of the afe-box (18 nucleotides). PBS_24_OL and PBS_46 align with better scoring to the right half of the complete afe-box. Numbers correspond to PBSs listed in Supplementary Table 1;

**(c)** Sequence logo file created from the sequence alignment shown in b.

- To look for sites of interest/conservation within a gene (motifs, binding sites, etc.)

https://meme-suite.org/meme/index.html

- To look for sites of interest/conservation within a gene (motifs, binding sites, etc.)



Multiple Sequence Aligment

Sequence Fasta file
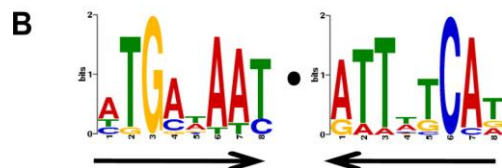
12

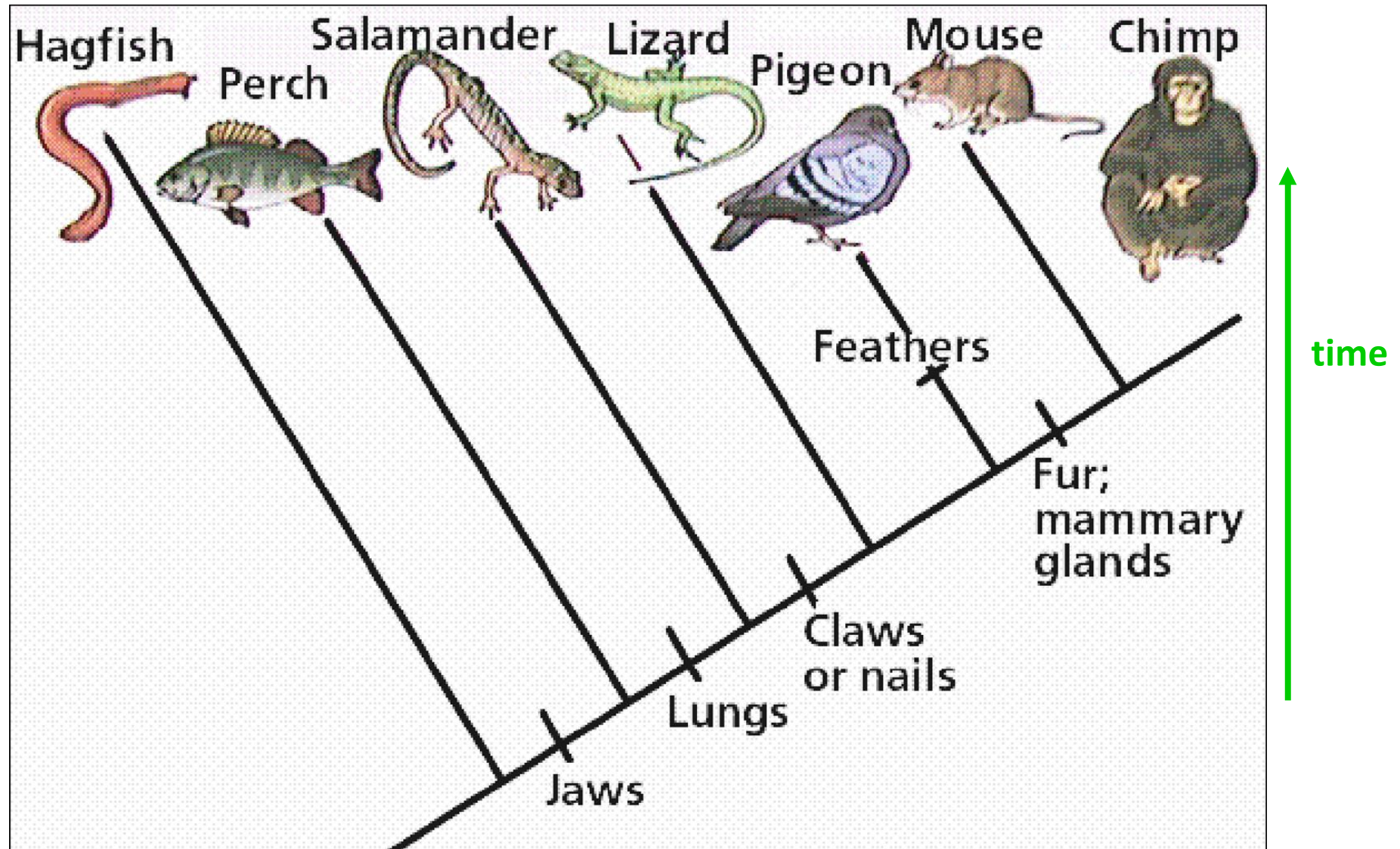https://meme-suite.org/meme/index.html

# Reconstruct phylogenetic trees

# Characters

- Heritable changes in features (morphology, DNA sequence etc…)

- The more similar characters you have, the more related you are

# Evolution and characters



Hagfish   Perch   Salamander   Lizard   Pigeon   Mouse   Chimp

Feathers

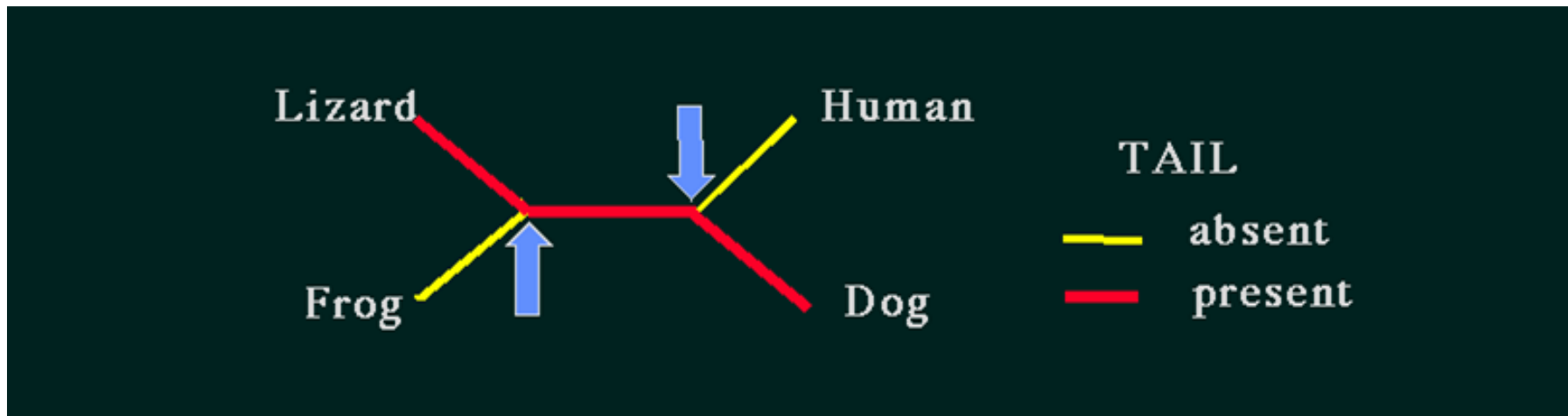Fur; mammary glands

Claws or nails

Lungs

Jaws

time

# A Unique Character:  Hair for Mammals

- Hair evolved only once and is "unreversed"
- Presence of hair is strong indication that  organism is a mammal

# Homoplasy: The formation of tails

- Tails evolved independently in the ancestors of frogs and humans
- Presence of a tail is no useful conclusions

# Classification according to characters – more characters can be good

|  | Colour | Skin | Cost |
|---|---|---|---|
| **Beef** | red | no | $$$ |
| **Duck** | red | yes | $$$ |
| **Pork** | white | no | $$ |
| **Chicken** | white | yes | $ |
| **Tofu** | white | sometimes | $ |

Is Chicken most similar to Tofu?

# Classification according to characters – increasing the number of characters

|  | Colour | Skin | Cost | Legs | Feathers | Hair |
|---|---|---|---|---|---|---|
| **Beef** | red | no | $$$ | four | no | yes |
| **Duck** | red | yes | $$$ | two | yes | no |
| **Pork** | white | no | $$ | four | no | yes |
| **Chicken** | white | yes | $ | two | yes | no |
| **Tofu** | white | sometimes | $ | none | no | no |

# Evolution and characters – the importance of comparing characters with common origins (homologous)
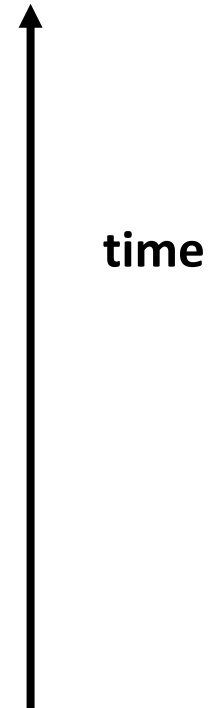
```
bioinformatics
bioinformatics
bioinformatios
oinformatios
informatios
information
information
```

time

# Evolution and characters

```
bioinformatics
bioinformatics
bioinformatios
--oinformatios
---informatios
---information
---information
```

time ↑

- Gaps represent non-homologous positions in the sequence.

- They reflect the occurrence of insertions/deletions or other rearrangements during the evolutionary process.

# Multiple Sequence Alignment

```
VTISCTGSSSNIGAG-NHVKWYQQLPG
VTISCTGTSSNIGS--ITVNWYQQLPG
LRLSCSSSGFIFSS--YAMYWVRQAPG
LSLTCTVSGTSFDD--YYSTWVRQPPG
PEVTCVVVDVSHEDPQVKFNWYVDG--
ATLVCLISDFYPGA--VTVAWKADS--
AALGCLVKDYFPEP--VTVSWNSG---
VSLTCLVKGFYPSD--IAVEWESNG--
```

The sole purpose of multiple sequence alignments is to place *homologous positions* of *homologous sequences* into the *same column*.

# Multiple sequence alignments and phylogenetic analysis

- *First step in any phylogenetic analysis*

- Phylogenetic analysis only as good as the alignment

in    -->    out!

# Steps in Multiple Alignment

## (A) Pairwise Alignment

Example – 4 sequences  $S_1$ $S_2$ $S_3$ $S_4$

$S_1$ ————————————

$S_2$ ————————————

$S_3$ ————————————

$S_4$ ————————————

6 pairwise comparisons
then cluster analysis

$S_2$

$S_4$

$S_1$

$S_3$

similarity →

## (B) Multiple alignment following the tree from A

$S_2$ ————————    ————

$S_4$ —  ————————————

align most similar pair

Gaps to optimize alignment

$S_1$ ——  ——————  ——

$S_3$ ——  ————————  ————

align next most similar pair

New gap to optimize
alignment of $(s_2 s_4)$ with $(s_1 s_3)$

$S_2$ ——————————  _▼__

$S_4$ ——  ——  ——————  ——

$S_1$ ——  ————————  ——

$S_3$ ——  ——————  ————

align alignments – preserve gaps

24

# Creating a MSA

- Clustal

  - Been around for ever and widely used

  - ClustalW (command line)

  - ClustalX (GUI)

  - Also available on many web servers

  - http://www.clustal.org/clustal2/

- Muscle

  - Faster and maybe more accurate than Clustal

  - Command line only

  - There are web servers  http://www.ebi.ac.uk/Tools/msa/muscle/

- T-Coffee

  - Most accurate, but also the slowest

  - Also has special variations for RNA, protein structure, etc.

  - http://tcoffee.crg.cat/

# Need something faster

- Clustal Omega:  HMM Based

- http://www.clustal.org/omega/

# Editing Alignments

- A MSA is rarely perfect
- Downstream tools will presume columns are homologous
- Remove unreliably aligned regions for phylogenetic analysis

# Editing Alignments

- A MSA is rarely perfect
- Downstream tools will presume columns are homologous
- Remove unreliably aligned regions for phylogenetic analysis

```
ILPITSPSKEGYESGKAPDEFSSGG
ILPEH--IKDDGELGAAPHSFSTAG
VLPLD-----S--AGRPADSFSAAG
VLPVDR-------DGQARDEYT-VG
VLPVDN-------KGEARDEYT-VG
LLPYDD-------QGRPQDDYSRAG
GIVSRSG---SNFDGEPKDSYGKVG
```

Delete?

# **BioEdit**: a biological sequence alignment editor

- An intuitive multiple document interface with convenient features
- Several sequence manipulation and analysis options and links to external analysis

https://bioedit.software.informer.com/7.2/

# JalView: Manual

# JalView: Manual

# GBlocks: Automatic

http://molevol.cmima.csic.es/castresana/Gblocks.html
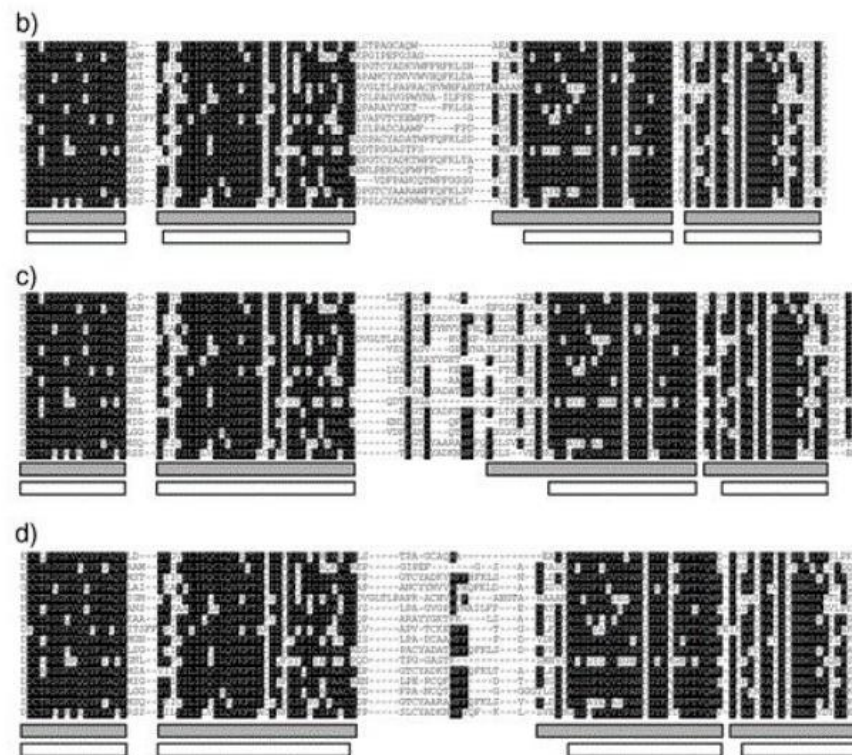
http://molevol.cmima.csic.es/castresana/Gblocks_server.html

The first one: GBLOCKS (http://molevol.cmima.csic.es/castresana/Gblocks.html)

Gblocks selects blocks in a similar way as it is usually done by hand but following a reproducible set of conditions. The selected blocks must fulfill certain requirements with respect to the lack of large segments of contiguous nonconserved positions, lack of gap positions and high conservation of flanking positions, making the final alignment more suitable for phylogenetic analysis.



The white and grey blocks under the alignments represent the parts of the alignment that Gblocks would keep using a more relaxed and a more stringent approach.

http://molevol.cmima.csic.es/castresana/Gblocks.html

http://molevol.cmima.csic.es/castresana/Gblocks_server.html

# GBlocks: Automatic

## How to run Gblocks: website



The is an on-line server that you can use if you only want to trim one alignment.

At the end of the alignment representation there's a link to obtain the trimmed alignment.

http://molevol.cmima.csic.es/castresana/Gblocks.html

http://molevol.cmima.csic.es/castresana/Gblocks_server.html

# How to run Gblocks: command line



```
mmarcet@saturn:~/Desktop/evomics$ Gblocks
*************************************************************
                      GBLOCKS 0.91b
SELECTION OF CONSERVED BLOCKS FROM MULTIPLE ALIGNMENTS
           FOR THEIR USE IN PHYLOGENETIC ANALYSIS
*************************************************************

o. Open File        Used to upload your alignment file (Fasta or NBRF/PIR format)

b. Block Parameters

s. Saving Options

g. (Get Blocks)

q. Quit


Your Choice: █
```

http://molevol.cmima.csic.es/castresana/Gblocks.html

http://molevol.cmima.csic.es/castresana/Gblocks_server.html

The probable truth: it depends on the dataset and the methodology used.

| Program | Number of citations (Google scholar) |
|---|---|
| BMGE | 381 |
| trimAl | 1737 |
| Gblocks | 5736 |

Trimming alignments tends to be part of a normal phylogenetic reconstruction pipeline.