

BI694

# Bioinformatics & Phylogenetics

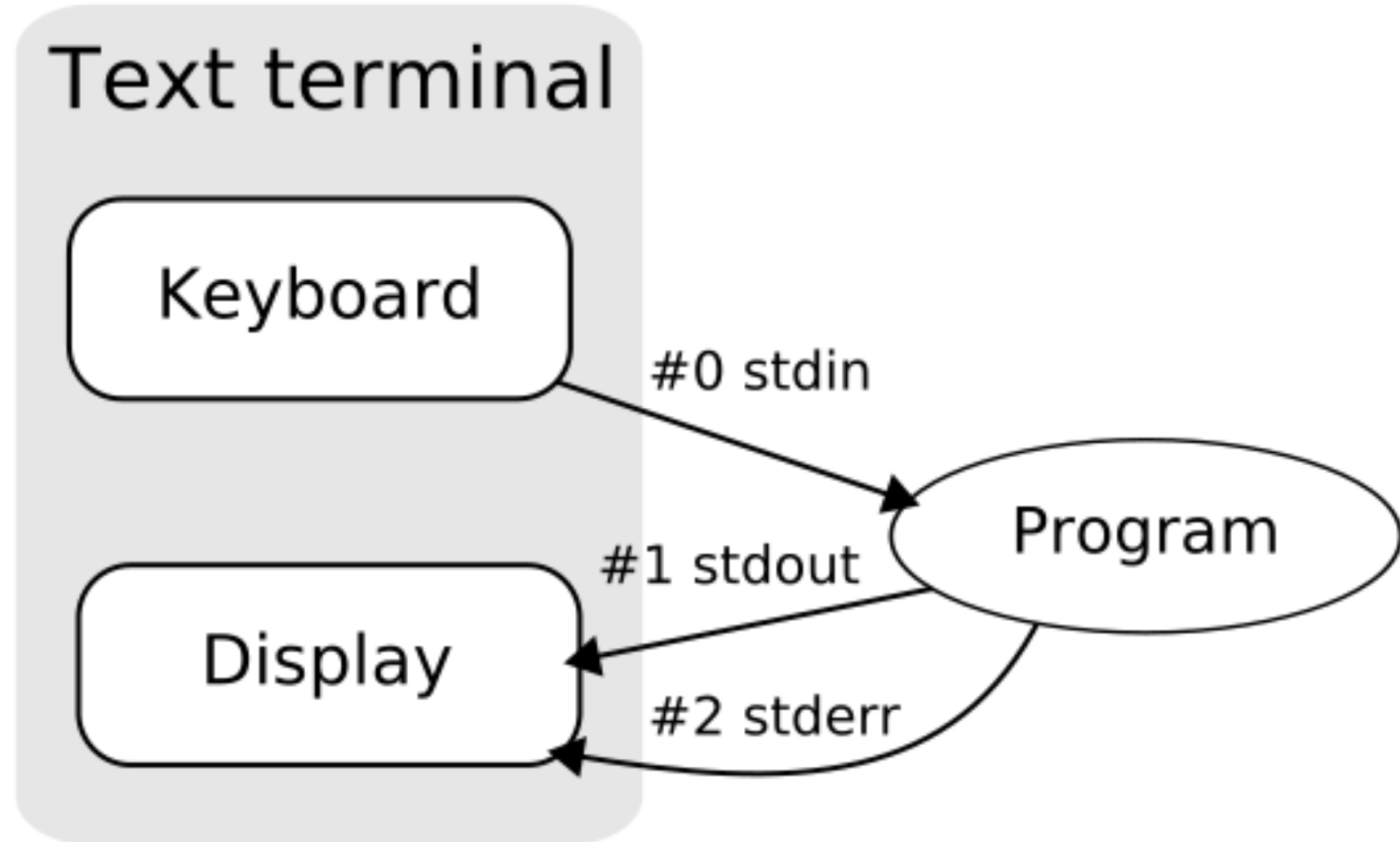
Winter Semester 2017

WEEK 3

(several slides courtesy of Maya Schushan and John R. Rose)

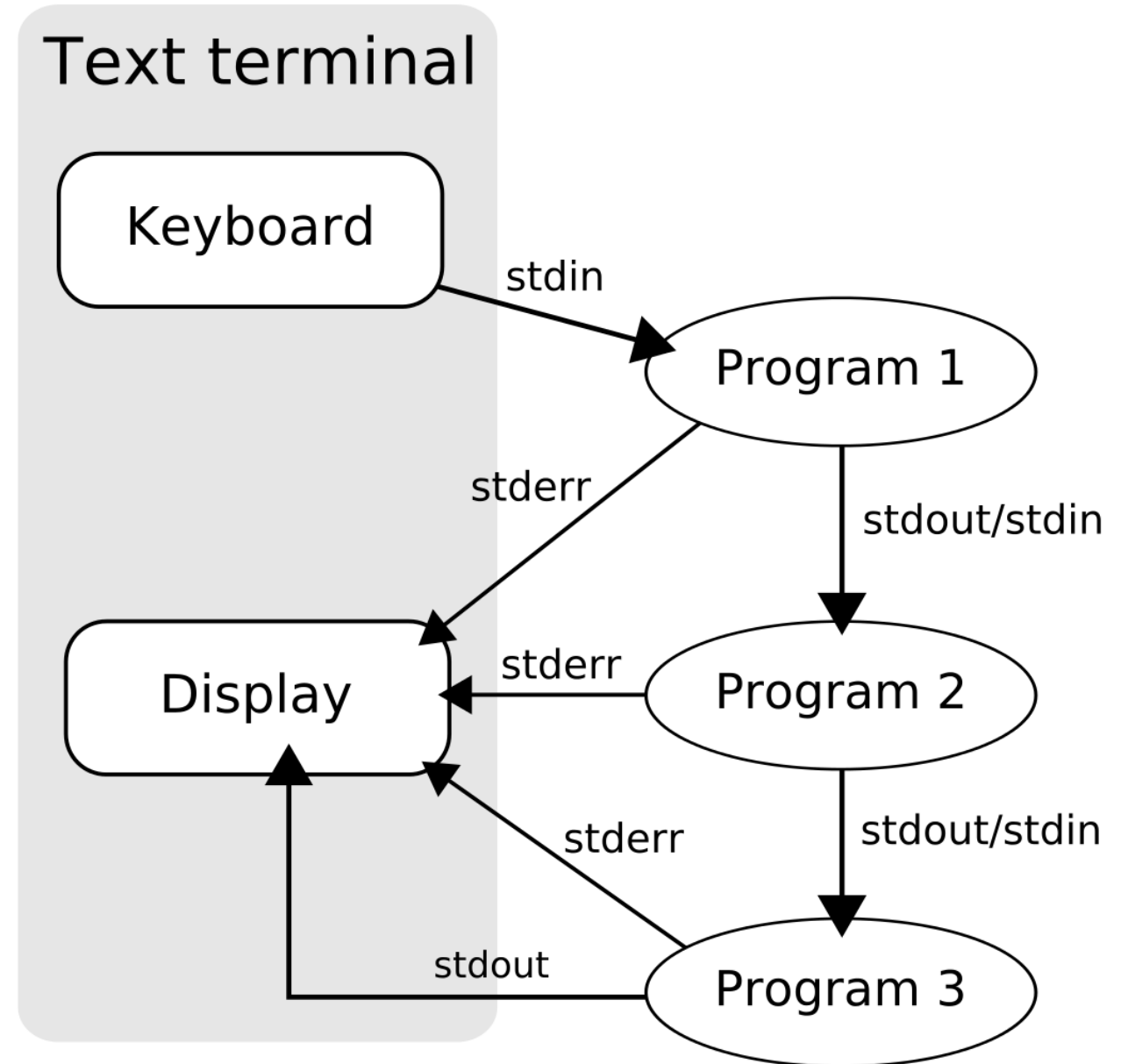
# Streams and pipes

- pipes – |
- redirects - > < << >>
- Separating streams – 1> 2>



# Streams and pipes

- pipes – |
- redirects - > < << >>
- Separating streams – 1> 2>



# Text editors

- **less** – text viewer
- nano – intuitive and easy but limited
- **vi**(m) – available everywhere and powerful
- emacs – powerful but too bloated for my taste
- geany – light weight IDE

# Let's get some data...

- Navigate to week 3 of the github course repository
- Download the **coral catalog csv** (through a browser or the terminal with wget)
- Let's look at it with **less, head, tail**
- How can we extract information quickly? (grep and awk)

# grep, sed, awk

- **grep** – powerful pattern searches, including regex
- **sed** – stream editor, search and replace (perl...)
- **awk** – scanning and processing programming language

# Let's look at some sequence data

- Navigate to week 3 of the github course repository
- Download OTUs.txt
- Let's examine it? What file format is this?
- `grep -c '>'`
- Seaview
- Phylip and Nexus formats

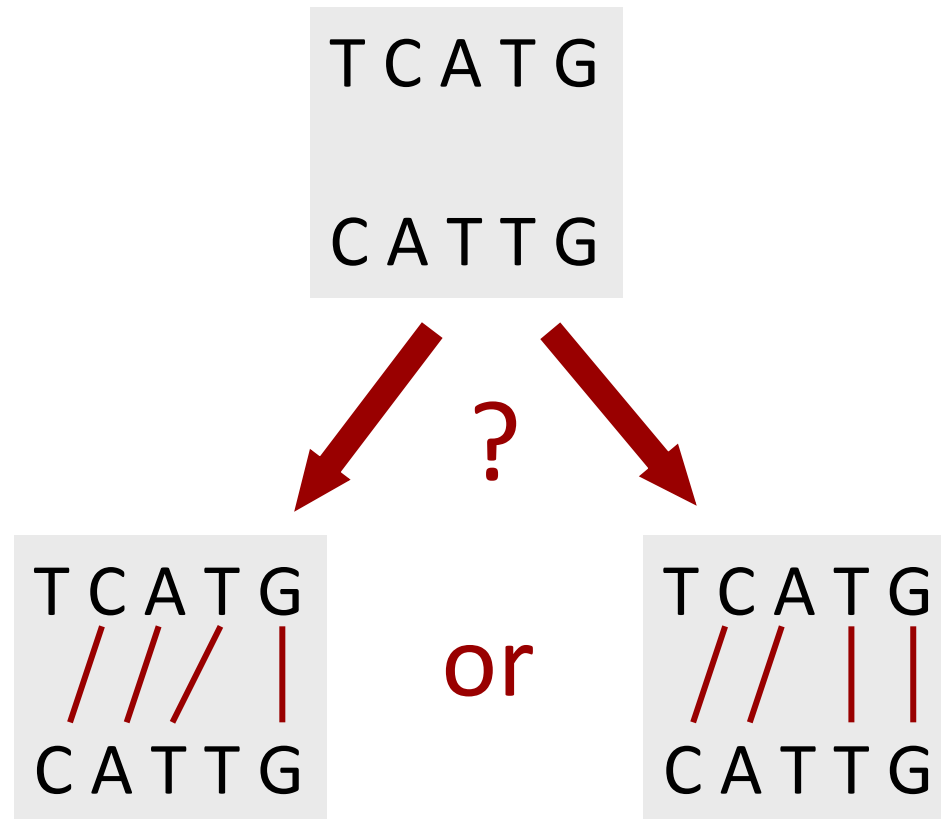
# Aligning sequences

- Navigate to week 3 of the github course repository
- Download Jellys.fst
- Are these sequences aligned?



# What is a sequence alignment?

Process of lining-up 2 or more sequences to achieve maximum level of identity, in order to find homologies.



# What is a multiple sequence alignment?

- Comparing 2 (pairwise) or more (multiple) sequences.
- Searching for a series of identical or similar characters in the sequences.

<b>VLSPADKTNVKA</b>	<b>AWAKVGAHAAGHG</b>
<b>VLSEAEWQLVLH</b>	<b>WAKVEADVAGHG</b>

# Defining Terms

- **Homology:**

Relation of sequences which is a result of shared from common ancestry

- **Identity:**

Sequences or Sub-sequences that are invariant.

- **Similarity:**

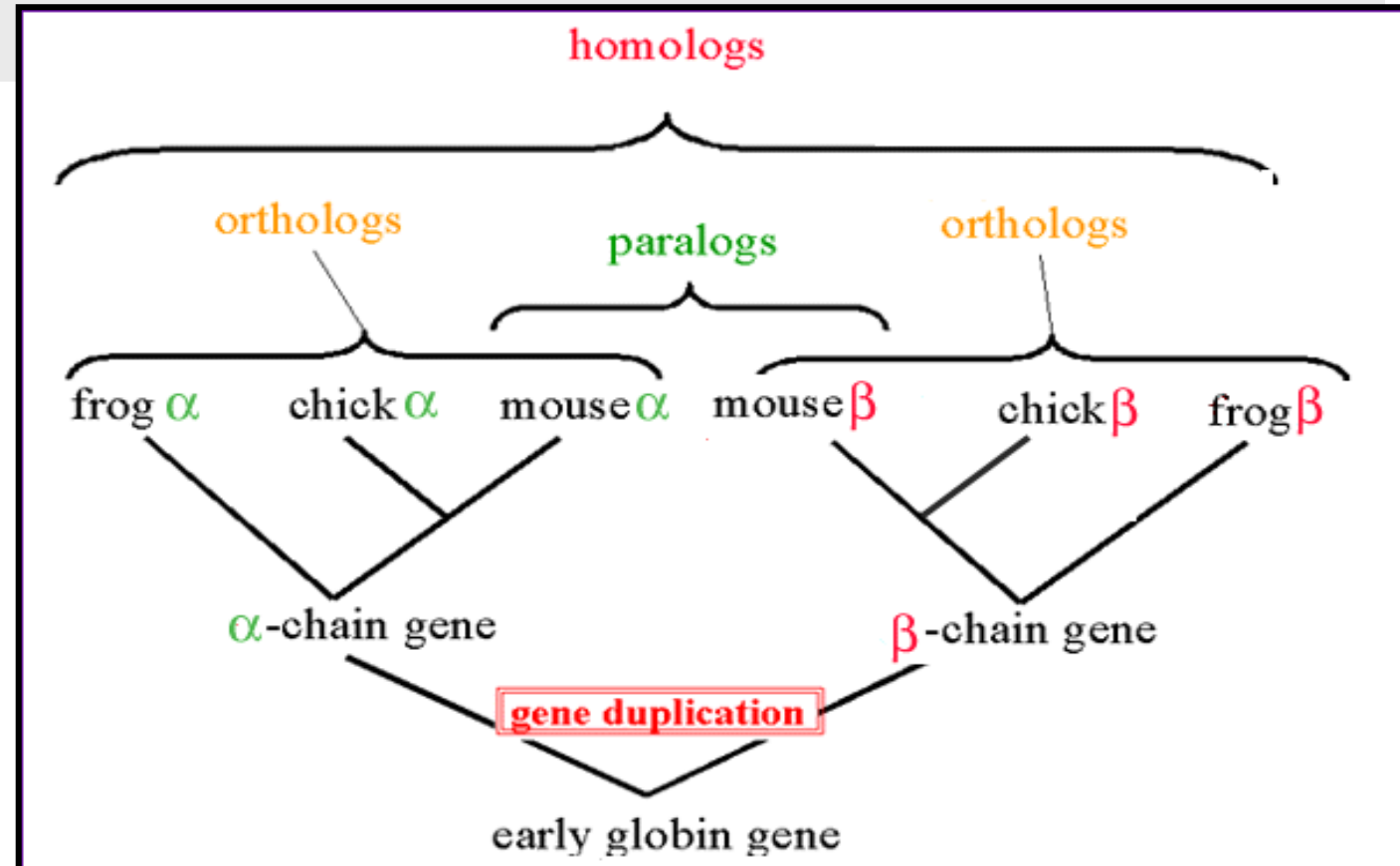
Sequences or Sub-sequences that are related.

C	A	G
		⋮
C	A	T

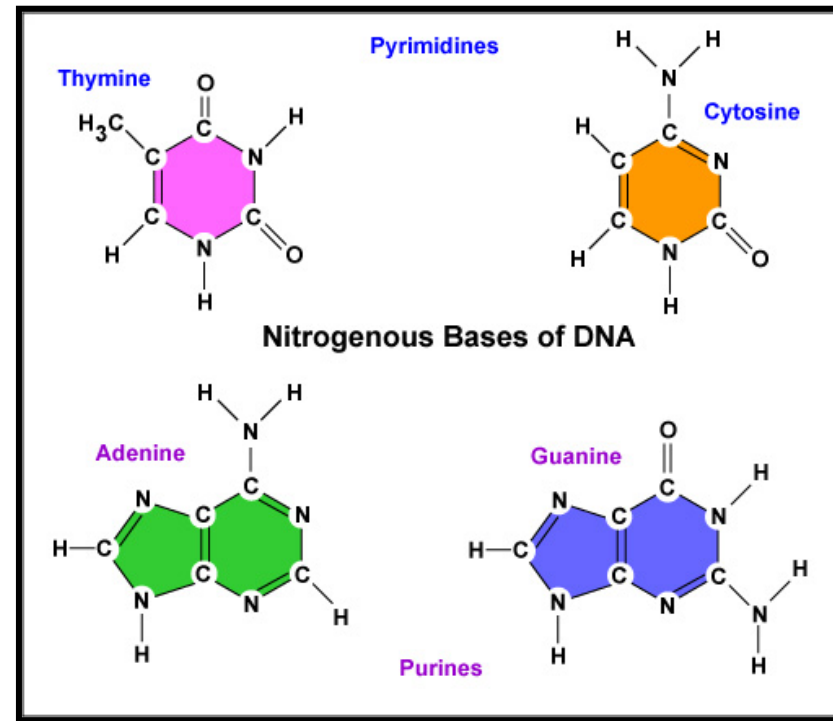
# Defining Terms

- **Homology:**

Relation of sequences which is a result of shared from common ancestry



# Defining Terms

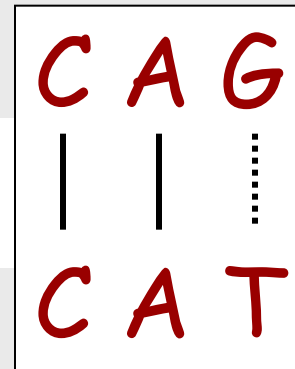


- **Identity:**

Sequences or Sub-sequences that are invariant.

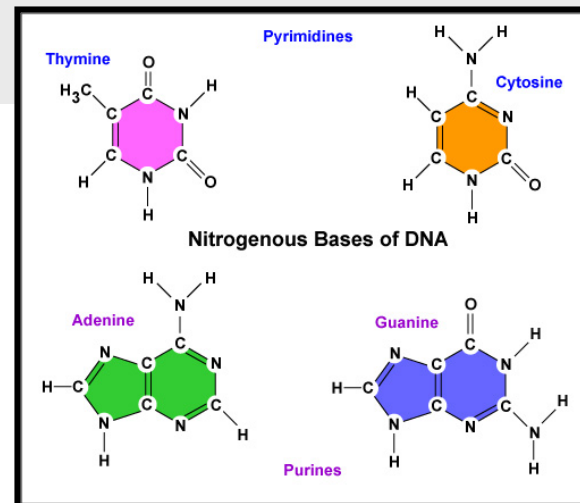
- **Similarity:**

Sequences or Sub-sequences that are related.



# DNA (RNA) similarity scoring

- Transitions – **purine to purine** or **pyrimidine to pyrimidine**  
(4 possibilities)
- Transversions – **purine to pyrimidine** or **pyrimidine to purine**  
(8 possibilities)
- By chance alone transversions should occur twice as often as transitions.
- De-facto **transitions** are more frequent than **transversions**.



# DNA (RNA) similarity scoring

From To	A	G	C	T
A	2			
G	-4	2		
C	-6	-6	2	
T	-6	-6	-4	2

Transversion

Transition

Match

# Protein similarity scoring

- **Observation:** some substitutions are more frequent than others, e.g., chemically similar amino acids
- As for DNA, protein matrices define the **probabilities of change** between the different amino acids
- Popular matrices are based on **empirical data: PAM & BLOSUM**



# Protein similarity scoring

Ala	4																			
Arg	-1	5																		
Asn	-2	0	6																	
Asp	-2	-2	1	6																
Cys	0	-3	-3	-3	9															
Gln	-1	1	0	0	-3	5														
Glu	-1	0	0	2	-4	2	5													
Gly	0	-2	0	-1	-3	-2	-2	6												
His	-2	0	1	-1	-3	0	0	-2	8											
Ile	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
Leu	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
Lys	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
Met	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
Phe	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
Pro	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
Ser	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
Thr	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
Trp	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
Tyr	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
Val	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	Met	Phe	Pro	Ser	Thr	Trp	Tyr	Val

# Why align sequences?

Predict characteristics of a protein

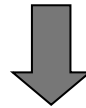
```
VTISCTGSSSNIGAG-NHVKWYQQLPG  
VTISCTGSSSNIGS--ITVNWYQQLPG  
LRLSCTGSGFIFSS--YAMYWYQQAPG  
LSLTCTGSGTSFDD-QYYSTWYQQPPG
```

# Why align sequences?

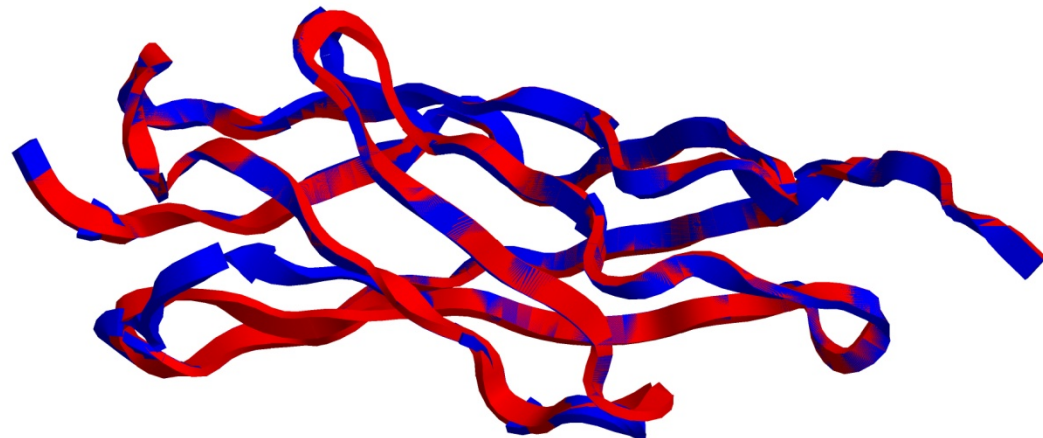
```
cad23_dom3  NLAIIDVQMDPIFINLPYSTNIYEHSPPGTTVRIITAIDQDKGR---PRGIGYTIVSGN
1EDH        EIVITDQNDNRPEFTQEVFEGSVAEGAVPGTSMKVSATDADDDVNTYNAAIAYTIVSQD

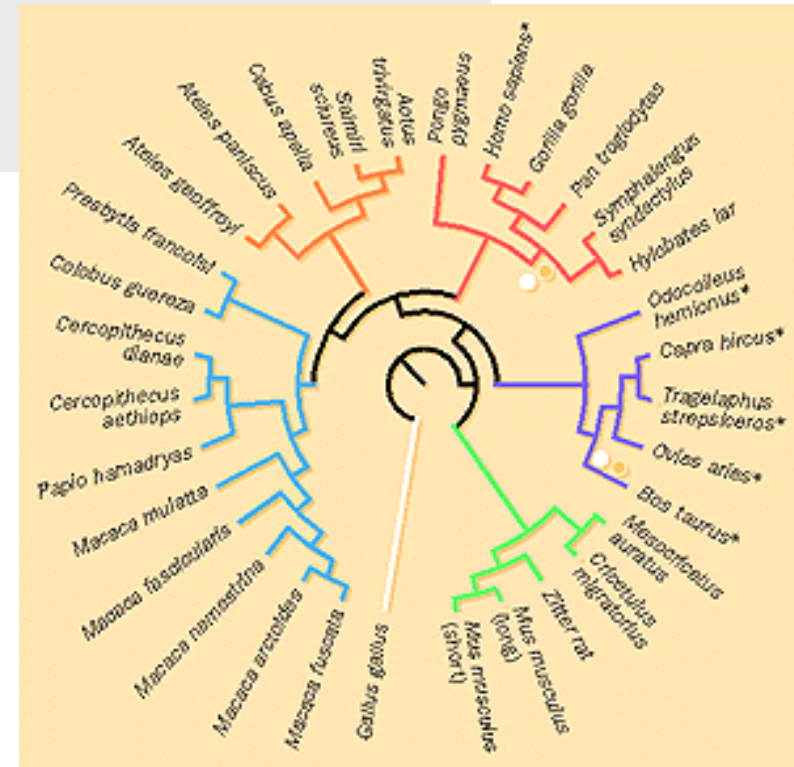
cad23_dom3  ----TNSIFALDYISGVL-TLNGLLDRENPLYSHGFILTVKGTELNDDRTPSDATVTTTF
1EDH        PELPHKMMFTVNRDTGVISVLTSGLDRES--YPT-YTLVVQAADLQGE----GLSTTAKA

cad23_dom3  NILVIDINDNAPEFNSSEYSVAITELAQVGFALPLF
1EDH        VITVKDINDNAPVFNP-----
```



A **model** is generated  
according to a  
**template** structure of a  
homologous protein





# Pairwise versus Multiple Sequence Alignment

Pairwise:

For 2 sequences

<b>F</b>	<b>G</b>	<b>K</b>	<b>-</b>	<b>G</b>	<b>K</b>	<b>G</b>
<b>F</b>	<b>G</b>	<b>K</b>	<b>F</b>	<b>G</b>	<b>K</b>	<b>G</b>



MSA:

For more than 2 sequences

<b>F</b>	<b>G</b>	<b>K</b>	<b>-</b>	<b>G</b>	<b>K</b>	<b>G</b>
<b>F</b>	<b>G</b>	<b>K</b>	<b>F</b>	<b>G</b>	<b>K</b>	<b>G</b>
<b>-</b>	<b>G</b>	<b>K</b>	<b>Q</b>	<b>G</b>	<b>K</b>	<b>G</b>
<b>-</b>	<b>-</b>	<b>K</b>	<b>F</b>	<b>G</b>	<b>K</b>	<b>G</b>

# Multiple sequence alignment

- By definition a multiple sequence alignment (MSA) is an alignment of 3 or more sequences (amino acid or nucleotide)
- Pairwise alignments align 2 sequences – useful for searching a database to find sequence matching query best (eg, BLAST)
- Generally 3 approaches to MSA:
  - direct
  - progressive (hierarchical)
  - iterative

# Direct sequence alignment

- Amino acids (AAs) are aligned using gap penalty and substitution model based on chemical properties of AAs
- Nucleotides are aligned using gap penalty and mismatch score
- Problem:
  - computationally expensive
  - for  $n > 2$  sequences  $n$ -dimensional matrix of pairwise alignments
  - scales exponentially

# Progressive alignment

Basic method:

1. Distance matrix is calculated
  - Distances are pairwise alignment scores
  - Gives divergence of each pair of sequences
2. Guide tree built from distance matrix
3. Progressive alignment according to guide tree
  - Branching order of tree specifies alignment order
  - Alignment progresses from leaves to root.



# Progressive alignment

## Distance matrix/pairwise alignments phase

- Fast approximation:
  - $k$ -tuple (ordered list of length  $k$ ) matches for identical residues (AAs or nucleotides), typically
    - 1 to 2 for proteins
    - 2 to 4 for nucleotide sequences
  - Scores are calculated as: ( $k$ -tuple matches) – fixed penalty per gap
  - Score is initially calculated as a percent identity score.
  - Distance =  $1.0 - (\text{score}/100)$

# Clustal alignments

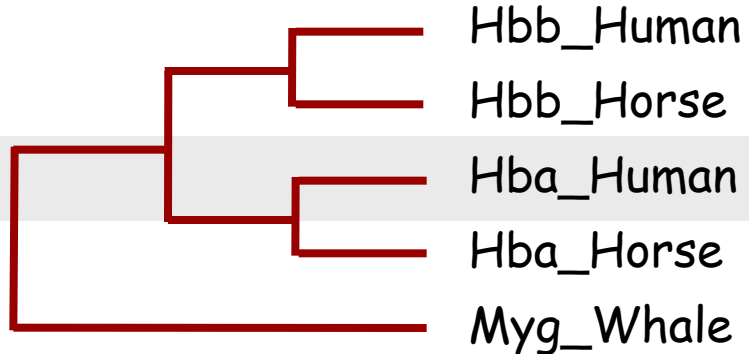
- Based on the idea that the sequences we want to align are **phylogenetically related**: a pairwise alignment algorithm is used iteratively, first to align the most closely related pair of sequences, then the next most closely related one to that pair.
- Rule “**once a gap, always a gap**”: The gaps between more similar pairs of sequences should not be affected by more distantly related ones.

# Clustal alignments

Hbb\_Human 1  
Hbb\_Horse 2  
Hba\_Human 3  
Hba\_Horse 4  
Myg\_Whale 5

	1	2	3	4	5
1	-				
2	17	-			
3	59	60	-		
4	59	59	13	-	
5	77	77	75	75	-

1. Quick pairwise alignment  
calculate distance matrix

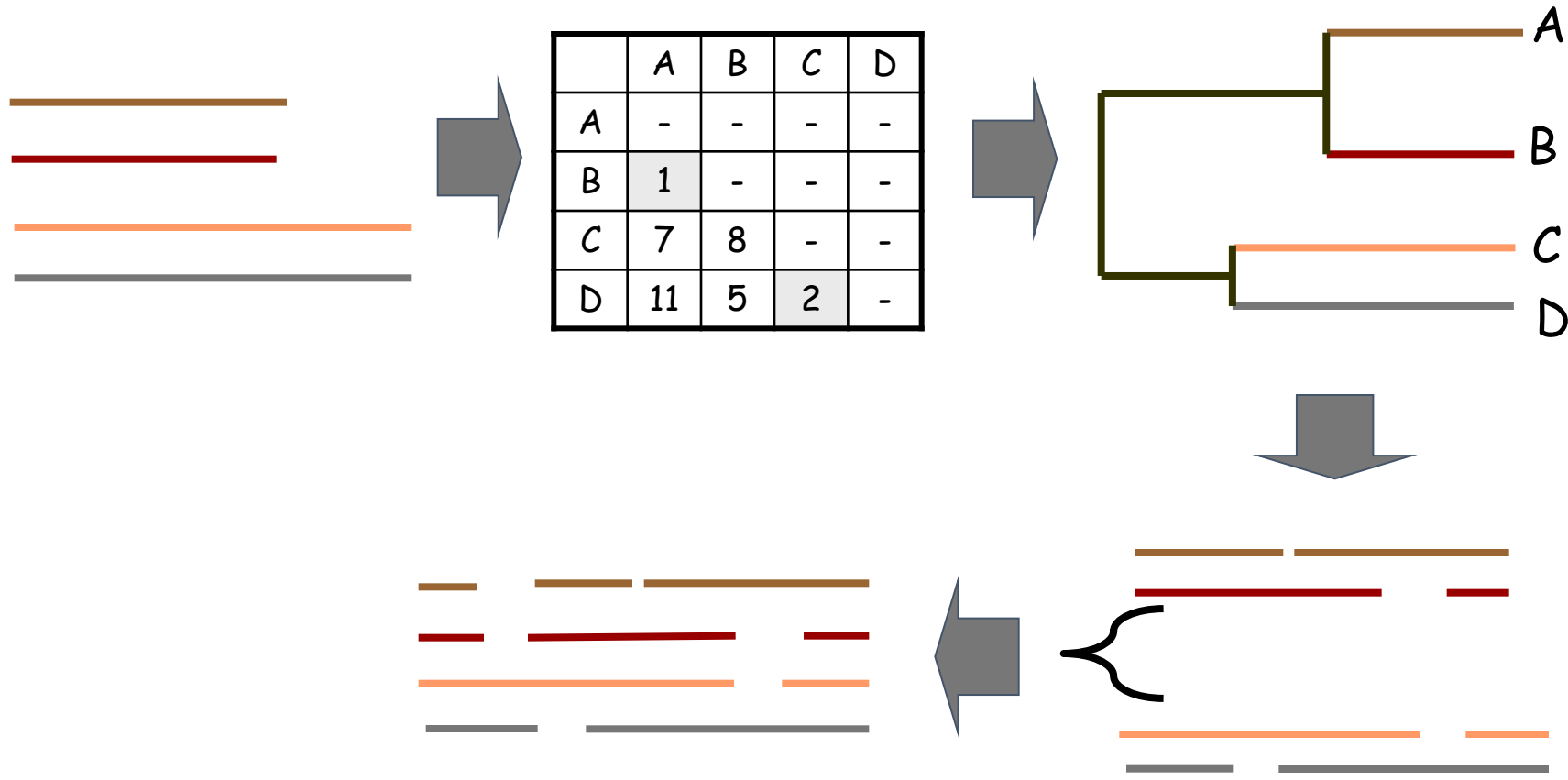


2. Build a guide tree using clustering  
or NJ

1	PEEKSAVTALWGKVN--VDEVGG	2	<div style="display: inline-block; vertical-align: middle; border-left: 2px solid orange; padding-left: 5px;"> <div style="border-top: 2px solid orange; height: 100px; width: 10px;"></div> </div>	3	<div style="display: inline-block; vertical-align: middle; border-left: 2px solid orange; padding-left: 5px;"> <div style="border-top: 2px solid orange; height: 100px; width: 10px;"></div> </div>	4
2	GEEKA AVLALWDKVN--EEEVGC					
3	PADKTNVKA AWGKVG AHAGEYGA					
4	AADKTNVKA AWSKVGGHAGEYGA	1				
5	EHEWQLVLHVWAKVEADVAGHGQ					

3. Progressive alignment  
following guide tree

# Clustal alignments

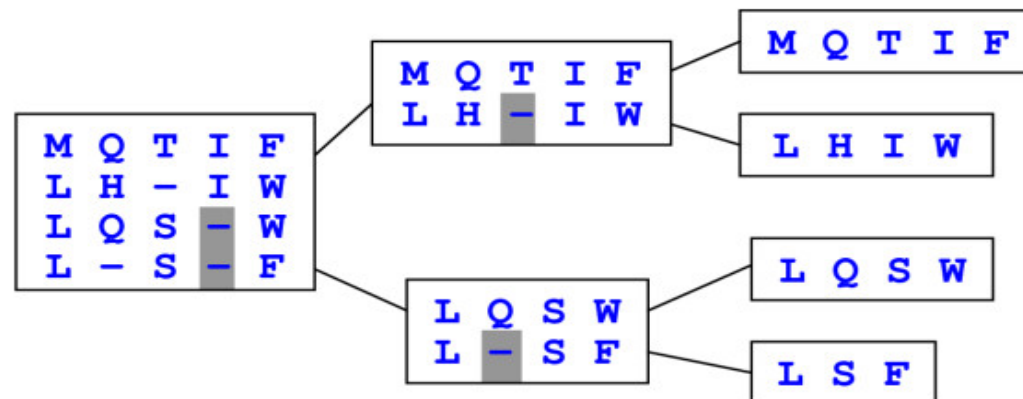


# Clustal alignments

- Progressive alignment not guaranteed to find global optimum
- Errors made at any stage in growing the MSA are propagated through to the final result

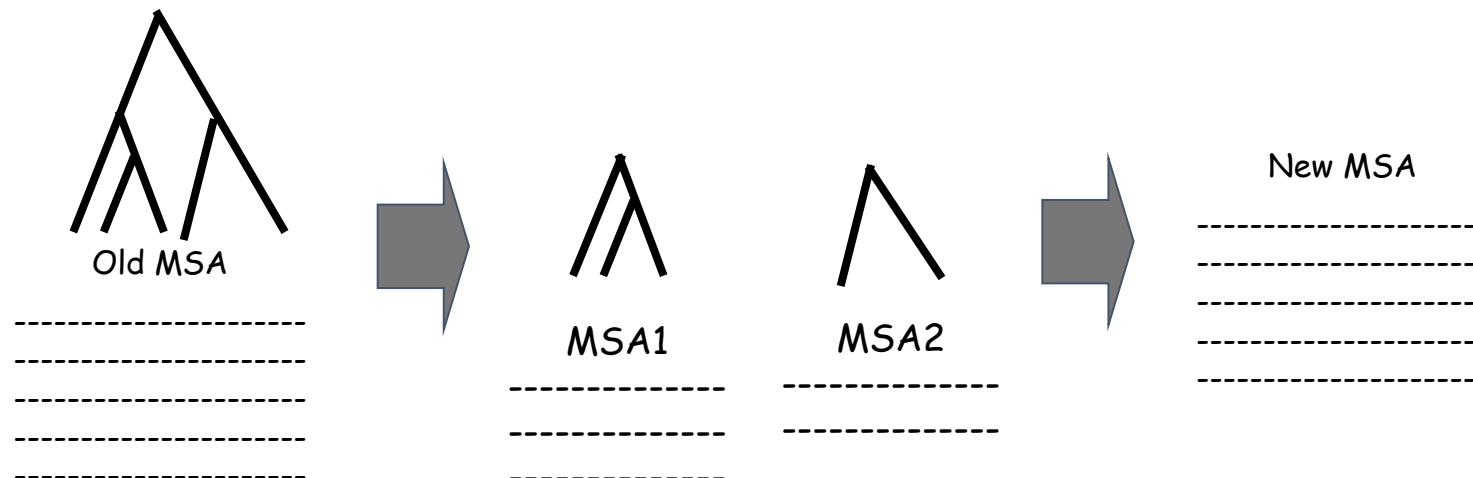
# Iterative alignment

- MUSCLE: **popular** MSA software
- Considered highly **accurate** MSA software
- The basic idea: **iterative progressive alignment** → refinement of initial results



# MUSCLE alignment

- An edge is chosen from the progressive alignment tree.
- The tree is divided into two subtrees by deleting this edge.
- The MSA from each subtree is computed by progressive alignment.
- The two MSAs are aligned, generating an entire new MSA
- If the new MSA achieves higher score than the previous → keep it



# MUSCLE alignment

It's a bit more complicated...

