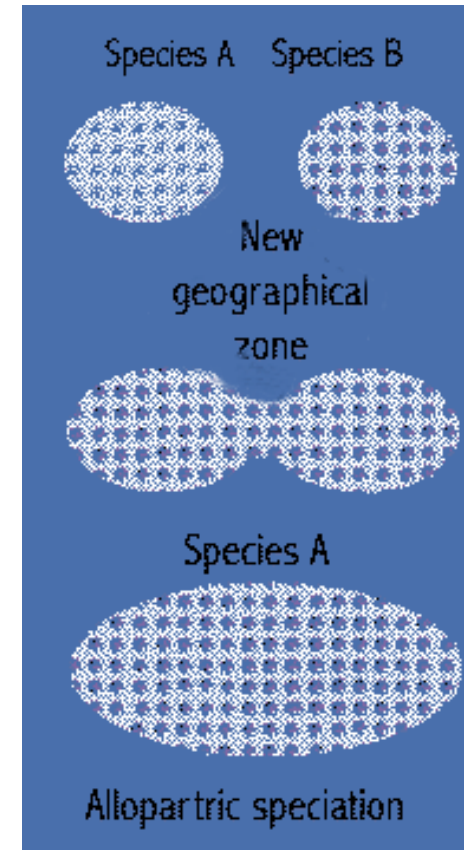
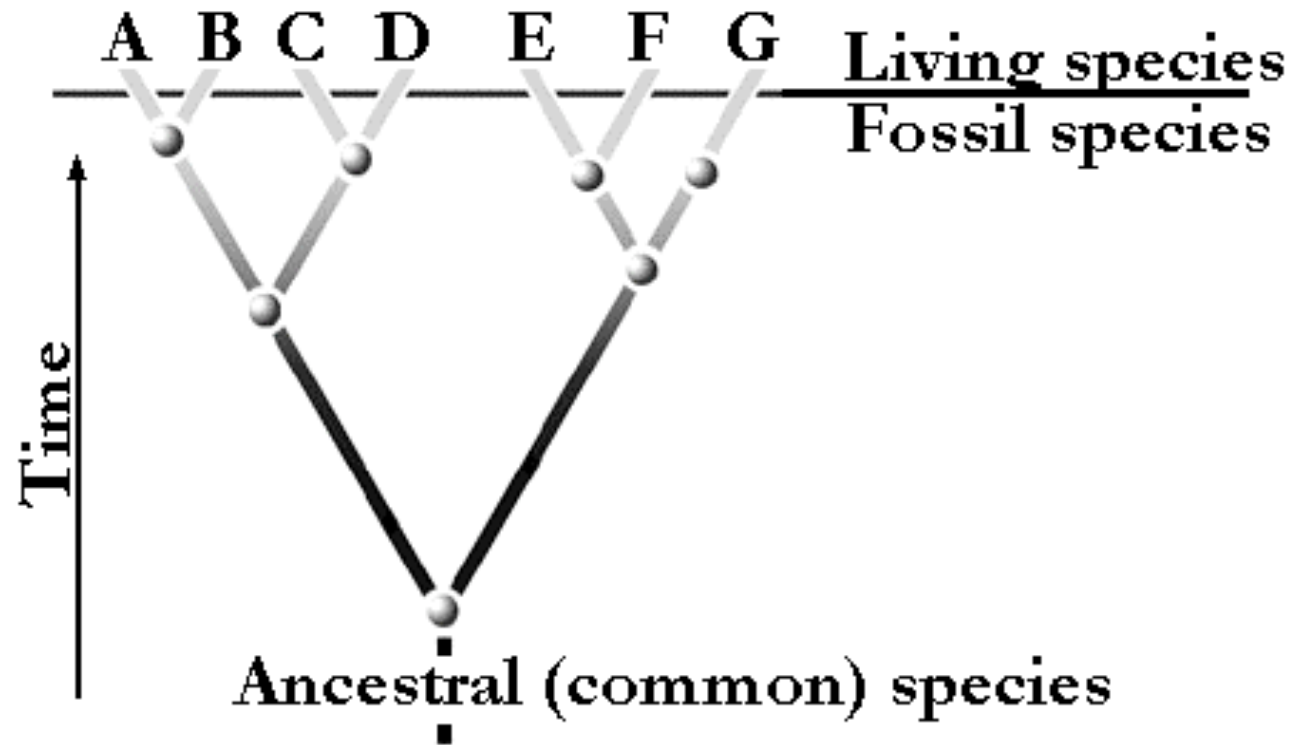
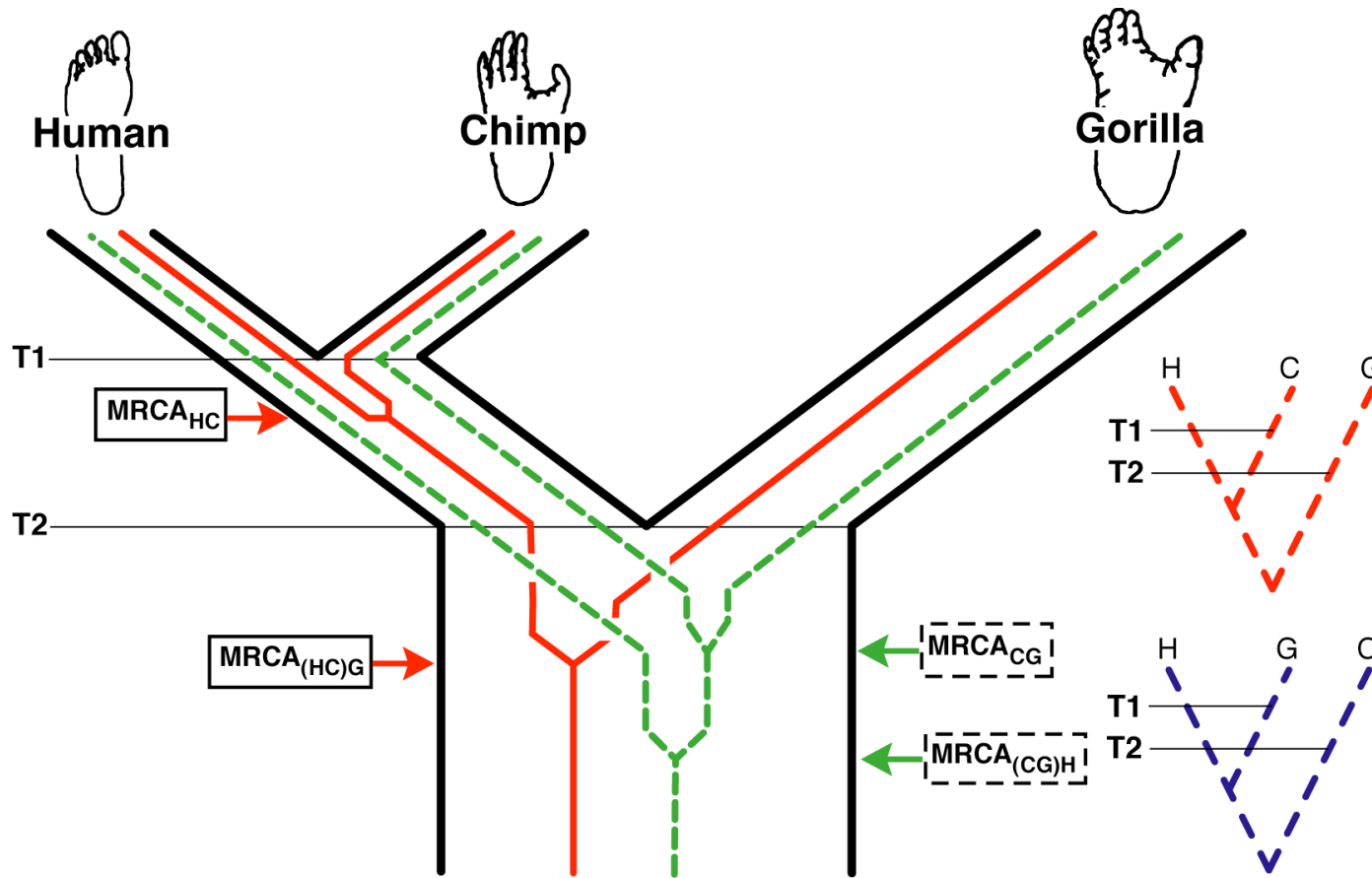


PHYLOGENY RECONSTRUCTION: THE BASICS

A Simple Concept of Speciation

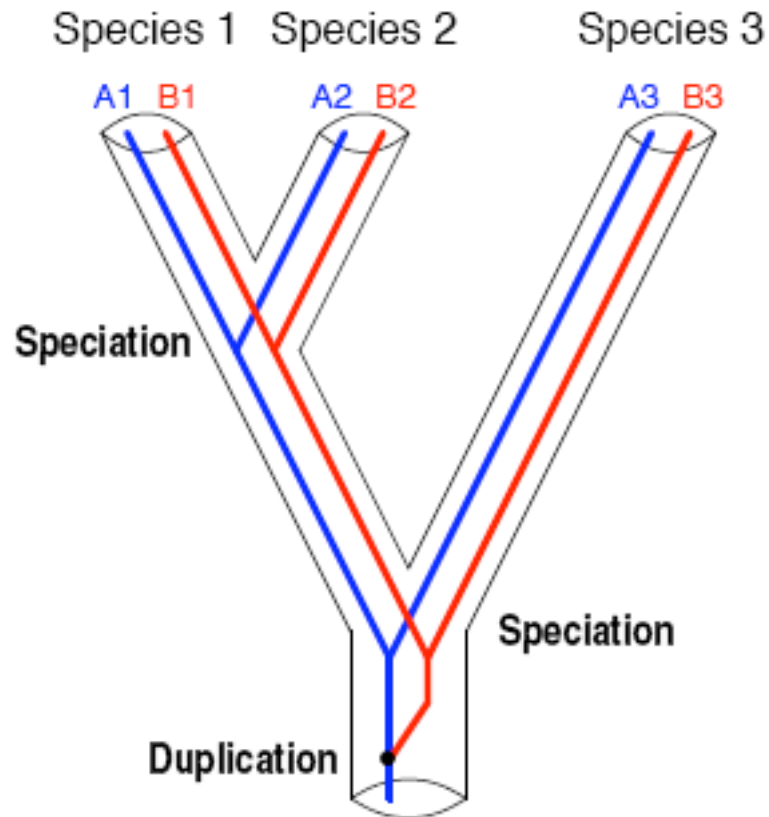


The Distinct History of Species and their DNA Sequences



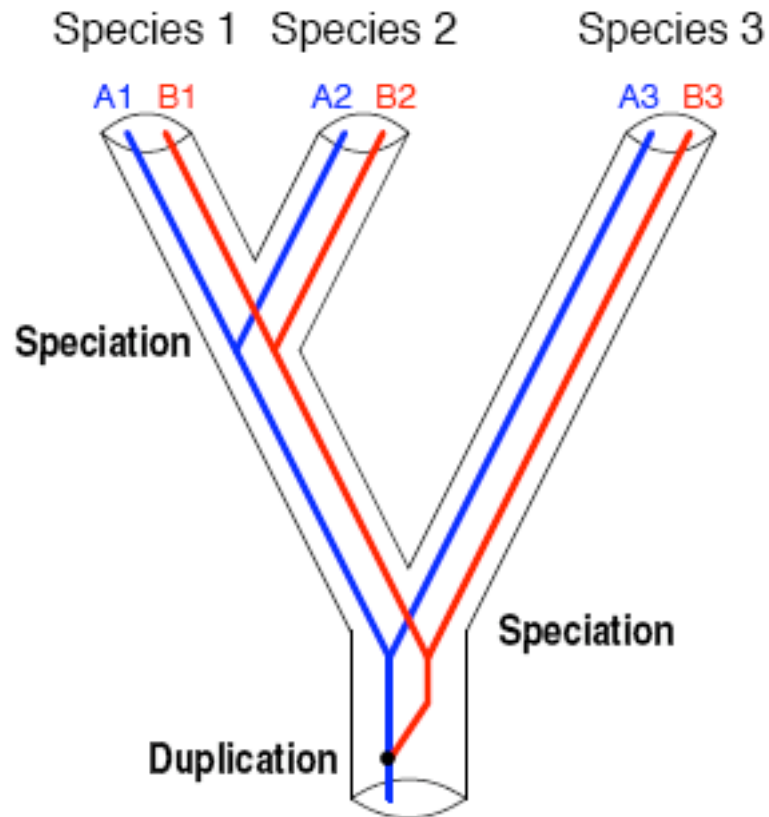
$$P_{old} = e^{-(T2-T1)/(2N_e \times g)}$$

Orthologous Sequences, Please!!



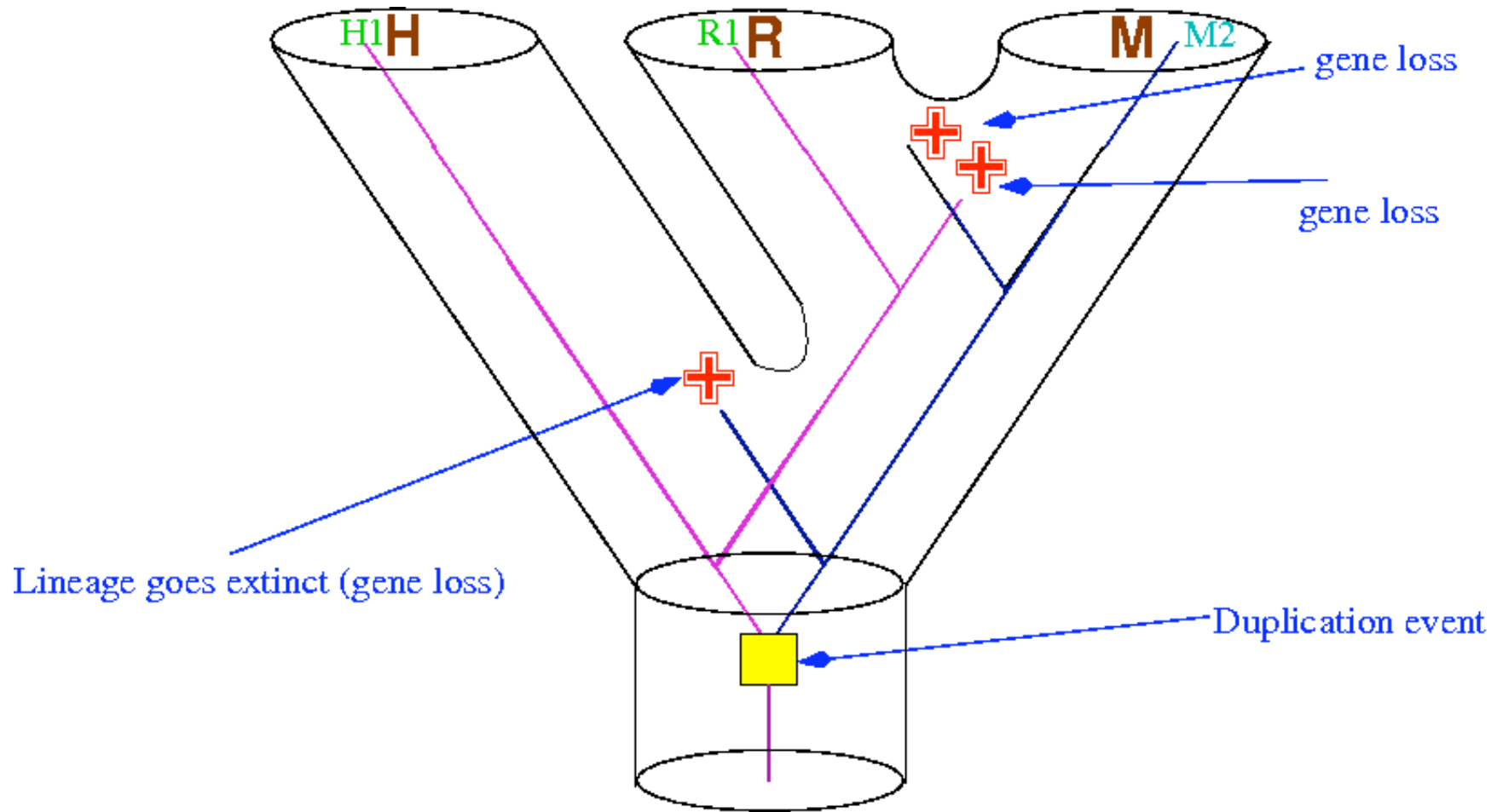
- Arguments for orthology assumption:
- a sequence tree that is congruent to the species tree
- conservation of genomic position
- sequence similarity (typically, reciprocal best blast hit)
- similarity of function

Orthologous Sequences, Please!!

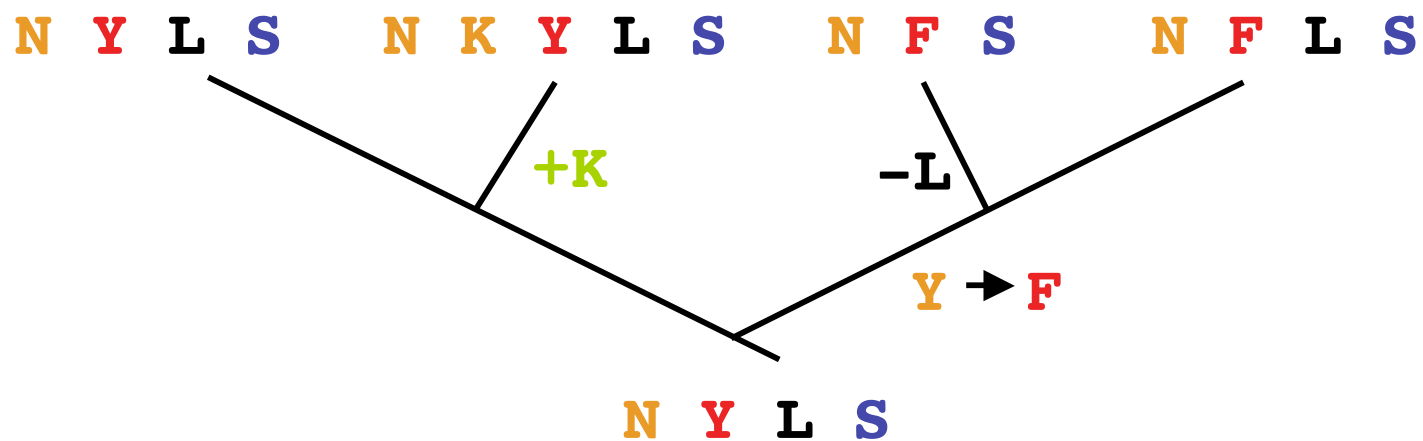


- Arguments for orthology assumption:
- a sequence tree that is congruent to the species tree
- conservation of genomic position
- sequence similarity (typically, reciprocal best blast hit)
- similarity of function

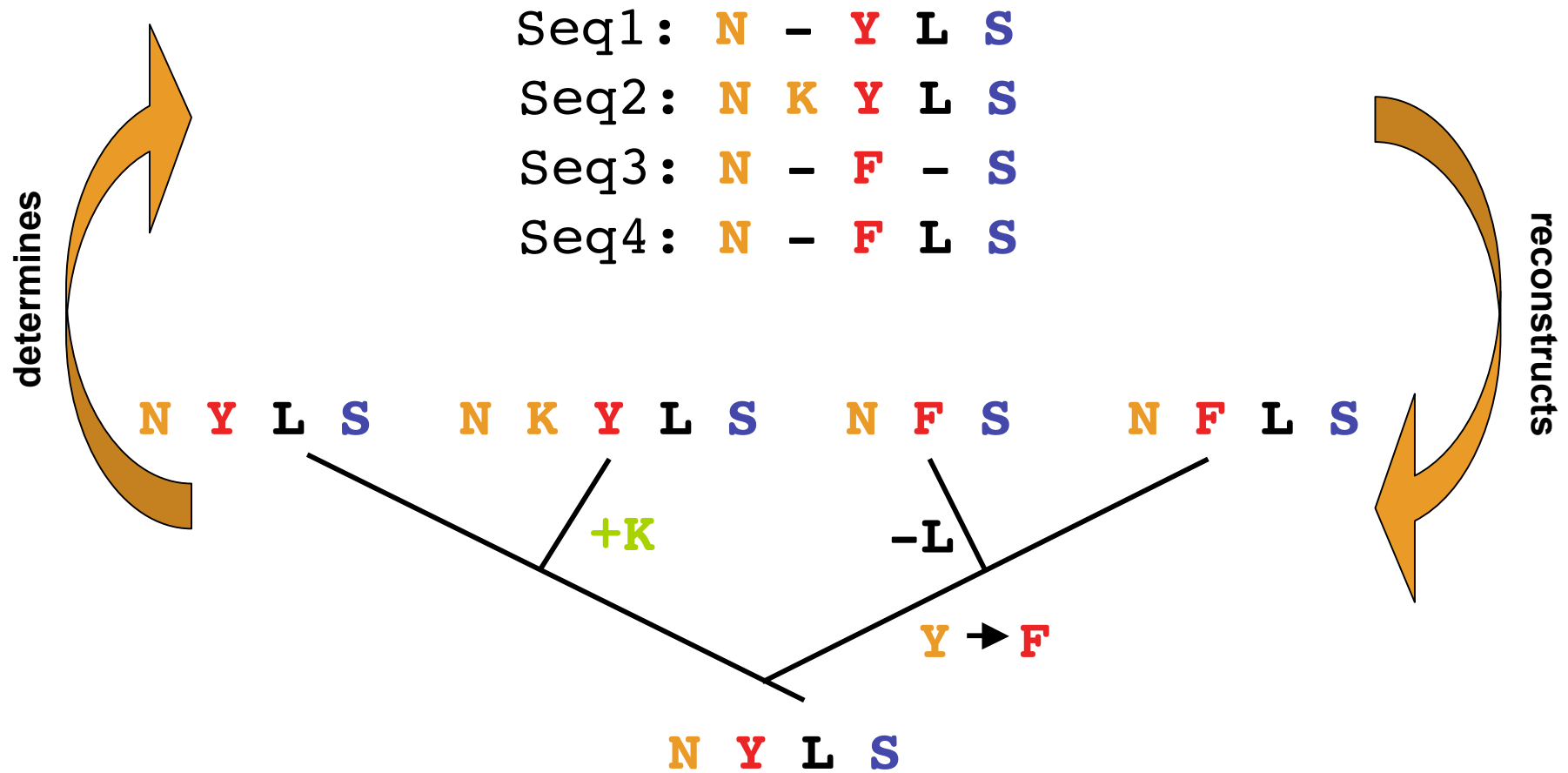
Hidden paralogy mimics orthology



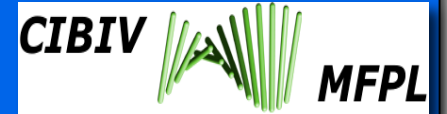
Sequence evolution in a nutshell



Sequence evolution in a nutshell



The Problem: Finding the homologous positions



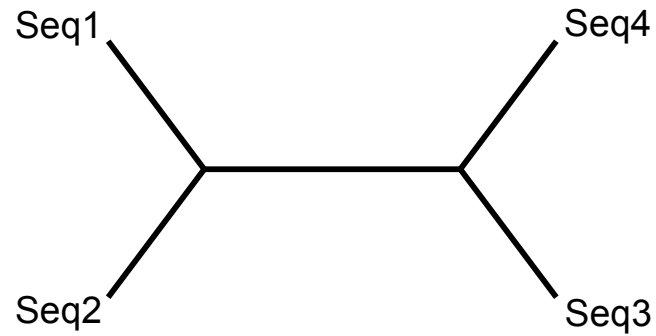
N **Y** **L** **S**

N **K** **Y** **L** **S**

N **F** **S**

N **F** **L** **S**

The Problem: Finding the homologous positions



Seq1: - N Y L S
Seq2: N K Y L S
Seq3: - N F - S
Seq4: - N F L S



N Y L S

N K Y L S

N F S

N F L S

The objective function



An mathematical function able to measure the biological quality of an alignment...

An mathematical function able to measure the biological quality of an alignment...

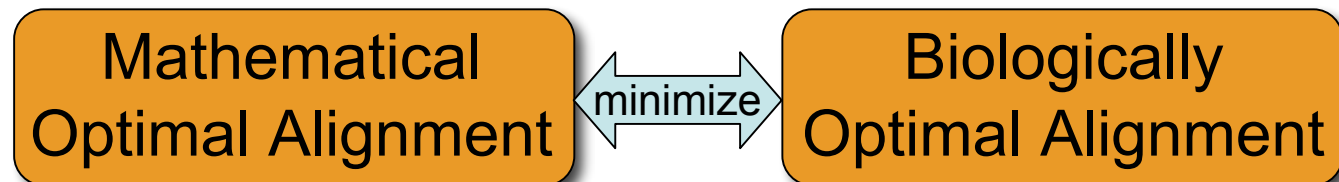
Related questions:

- What should a biologically correct alignment look like?
- To what extent can we define and formalize its properties?

An mathematical function able to measure the biological quality of an alignment...

Related questions:

- What should a biologically correct alignment look like?
- To what extent can we define and formalize its properties?



A mathematical function ment to measure the biological quality of an alignment...

$$\sigma(\alpha) = \sum_{i=1}^n S(a_i, b_i)$$

$\sigma(\alpha)$: the score of the pairwise alignment α

n : length of α

a_i : letter of sequence A at position i in α

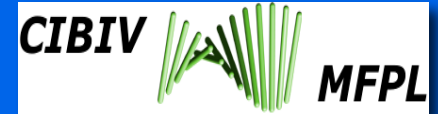
b_i : letter of sequence B at position i in α

A mathematical function ment to measure the biological quality of an alignment...

$$\sigma(\alpha) = \sum_{i=1}^n S(a_i, b_i)$$

Objective: find α that maximizes $\sigma(\alpha)$!

The scoring function S , *an example*

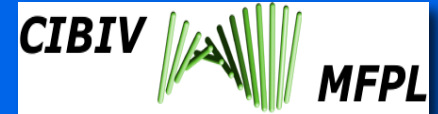


Given two sequences $A = \{a_1, a_2, \dots, a_n\}$ and $B = \{b_1, b_2, \dots, b_m\}$ and a scoring function S such that

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

then we look for that alignment, that gives us the highest score by summing up the column scores $S(a_i, b_j)$ for all columns of the alignment.

The scoring function S , an example



Given two sequences $A = \{a_1, a_2, \dots, a_n\}$ and $B = \{b_1, b_2, \dots, b_m\}$ and a scoring function S such that

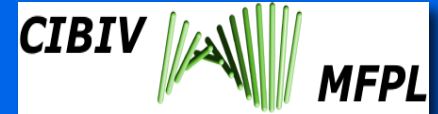
$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

then we look for that alignment, that gives us the highest score by summing up the column scores $S(a_i, b_j)$ for all columns of the alignment.

For example:

T	G	C	T	C	G	T	A	
T	-	-	T	C	A	T	A	
+5	-6	-6	+5	+5	-2	+5	+5	= 11

Why not just scoring all alignments?

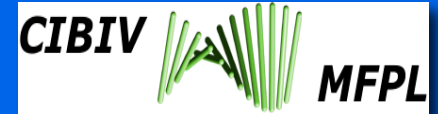


$$\begin{array}{cccccccc} \mathbf{A1:} & T & G & C & T & C & G & T & A \\ & T & - & - & T & C & A & T & A \\ & +5 & -6 & -6 & +5 & +5 & -2 & +5 & +5 & = & 11 \end{array}$$

$$\begin{array}{cccccccc} \mathbf{A2:} & T & G & C & T & C & G & T & A \\ & T & - & T & - & C & A & T & A \\ & +5 & -6 & -2 & -6 & +5 & -2 & +5 & +5 & = & 4 \end{array}$$

etc...

Why not just scoring all alignments?



- There are far too many
 - number of possible pairwise alignments: $\binom{2n}{n}$
 - for two sequences of length N there are 10^{179} possibilities

- There are far too many
 - number of possible pairwise alignments: $\binom{2n}{n}$
 - for two sequences of length N there are 10^{179} possibilities

Hence, we need a smart way to cut the computation short, like the **dynamic programming** approach for pairwise alignments by *Needleman and Wunsch* (1970).

Re-use of previous results

A1:

T	G	C	T	C	G	T	A	= 11
T	-	-	T	C	A	T	A	
+5	-6	-6	+5	+5	-2	+5	+5	

A2:

T	G	C	T	C	G	T	A	= 4
T	-	T	-	C	A	T	A	
+5	-6	-2	-6	+5	-2	+5	+5	

etc...

A **dynamic programming** approach usually includes:

- A mathematical description of the (biological) quality of an solution, i.e. an recursive objective function
- The computation of all intermediate values needed to obtain the globally optimal solution, thereby avoiding double-computations
- The reconstruction of the globally optimal solution from the values obtained in the previous step (backtracking)

The Needleman-Wunsch pair-wise alignment

	0	1	2	3	4	5	6	7	8
		T	G	C	T	C	G	T	A
0									
1	T								
2	T								
3	C								
4	A								
5	T								
6	A								

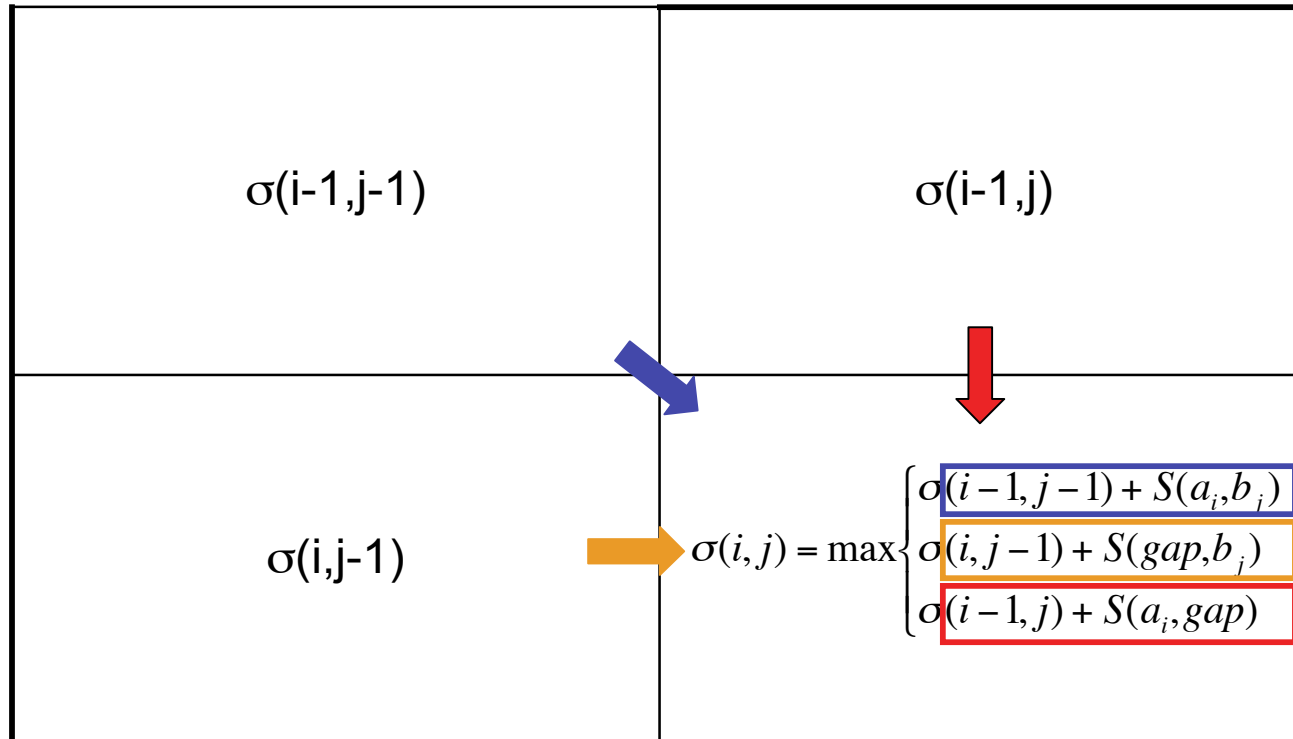
Scoring function

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

Objective function

$$\sigma(i, j) = \max \begin{cases} \sigma(i-1, j-1) + S(a_i, b_j) \\ \sigma(i, j-1) + S(\text{gap}, b_j) \\ \sigma(i-1, j) + S(a_i, \text{gap}) \end{cases}$$

The Needleman-Wunsch algorithm



➤ $\sigma(i,j)$ is the optimal alignment score up to and including a_i and b_j

$$S(a_i,b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

Needleman-Wunsch algorithm: Initialization

	0	1	2	3	4	5	6	7	8
		T	G	C	T	C	G	T	A
0	0	-6	-12	-18	-24	-30	-36	-42	-48
1	T	-6							
2	T	-12							
3	C	-18							
4	A	-24							
5	T	-30							
6	A	-36							

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

The Needleman-Wunsch algorithm: Recursion

	0	1	2	3	4	5	6	7	8
		T	G	C	T	C	G	T	A
0	0	-6	-12	-18	-24	-30	-36	-42	-48
1	T -6	5							
2	T -12								
3	C -18								
4	A -24								
5	T -30								
6	A -36								

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

The Needleman-Wunsch algorithm: Recursion

	0	1	2	3	4	5	6	7	8
		T	G	C	T	C	G	T	A
0	0	-6	-12	-18	-24	-30	-36	-42	-48
1	T -6	5	-1						
2	T -12								
3	C -18								
4	A -24								
5	T -30								
6	A -36								

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

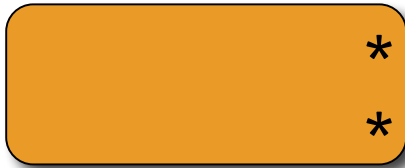
The Needleman-Wunsch algorithm: Recursion

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

Needleman-Wunsch algorithm: Backtrack

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11



Needleman-Wunsch algorithm: Backtrack

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11

A*

A*

Needleman-Wunsch algorithm: Backtrack

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11

TA*
TA*

Needleman-Wunsch algorithm: Backtrack

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11

TGCTCGTA
T--TCATA

Alignment Score: 11

Smith-Waterman pairwise local alignment

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	0	0	0	0	0	0	0	0	
1	T	0	5	0	0	5	0	0	5	0
2	T	0	5	3	0	5	3	0	5	3
3	C	0	0	3	8	2	10	4	0	3
4	A	0	0	0	2	6	4	8	2	5
5	T	0	5	0	0	7	4	2	13	7
6	A	0	0	3	0	1	5	2	7	18

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

$$\sigma(i, j) = \max \begin{cases} \sigma(i-1, j-1) + S(a_i, b_j) \\ \sigma(i, j-1) + S(\text{gap}) \\ \sigma(i-1, j) + S(\text{gap}) \\ 0 \end{cases}$$

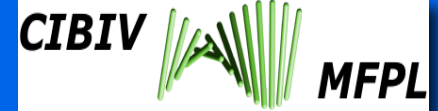
Smith-Waterman pairwise local alignment

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	0	0	0	0	0	0	0	0	
1	T	0	5	0	0	5	0	0	5	0
2	T	0	5	3	0	5	3	0	5	3
3	C	0	0	3	8	2	10	4	0	3
4	A	0	0	0	2	6	4	8	2	5
5	T	0	5	0	0	7	4	2	13	7
6	A	0	0	3	0	1	5	2	7	18

TCGTA
TCATA

Alignment Score: 18

Alternative Scoring Functions



Blosum62:

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	9																				C
S	-1	4																			S
T	-1	1	5																		T
P	-3	-1	-1	7																	P
A	0	1	0	-1	4																A
G	-3	0	-2	-2	0	6															G
N	-3	1	0	-2	-2	0	6														N
D	-3	0	-1	-1	-2	-1	1	6													D
E	-4	0	-1	-1	-1	-2	0	2	5												E
Q	-3	0	-1	-1	-1	-2	0	0	2	5											Q
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8										H
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5									R
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5								K
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							M
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4						I
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4					L
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4				V
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			F
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		Y
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	W

Many others...

PAM250:

C Cys	12																					
S Ser	0	2																				
T Thr	-2	1	3																			
P Pro	-3	1	0	6																		
A Ala	-2	1	1	1	2																	
G Gly	-3	1	0	-1	1	5																
N Asn	-4	1	0	-1	0	0	2															
D Asp	-5	0	0	-1	0	1	2	4														
E Glu	-5	0	0	-1	0	0	1	3	4													
Q Gln	-5	-1	-1	0	0	-1	1	2	2	4												
H His	-3	-1	-1	0	-1	-2	2	1	1	3	6											
R Arg	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6										
K Lys	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5									
M Met	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6								
I Ile	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	-2	2	5							
L Leu	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6						
V Val	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	-2	2	4	2	4					
F Phe	-4	-3	-3	-5	-5	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9				
Y Tyr	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7	10			
W Trp	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0	0	17		
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W		

Both, Needleman-Wunsch and Smith-Waterman alignment methods are **exact** methods since they guarantee a globally optimal solution for the optimization problem!

Drawback: Computational expensive, i.e. $O(nm)$ in time and memory

Exact vs. Heuristic searches

Solutions:

- omit regions from the grid, that cannot contribute to the optimal alignment (reduction of the search space, by remaining exact)

	0	1	2	3	4	5	6	7	8	
		T	G	C	T	C	G	T	A	
0	0	-6	-12	-18	-24	-30	-36	-42	-48	
1	T	-6	5	-1	-7	-13	-19	-25	-31	-37
2	T	-12	-1	3	-3	-2	-8	-14	-20	-26
3	C	-18	-7	-3	8	2	3	-3	-9	-15
4	A	-24	-13	-9	2	6	0	1	-5	-4
5	T	-30	-19	-15	-4	7	4	-2	6	0
6	A	-36	-25	-21	-10	1	5	2	0	11

Solutions:

- use of heuristics (more rigorous reduction of the search space, sacrificing the guaranteed optimal solution for search speed)

- Lookup method for finding an alignment

Pos: 1 2 3 4 5 6 7 8 9 10 11
 Seq 1: k c s p t a
 Seq 2: a c s p r k

Amino acid	Pos in Seq 1	Pos in Seq 2	Offset
k	1	11	10
c	2	7	-5
s	3	8	-5
p	4	9	-5
t	5	-	-
a	6	6	0
r	-	10	-

What we are really looking for:

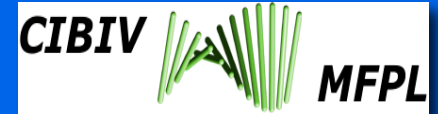


αA αB αC
 1 10 20 30 40 50 60
 Ther_tengcongensis MKG T I V G T W I K T L R D L Y G N D V V D E S L K S V G W E P D R V I T P L E D I D D D E V R R I F A K V S E K T G K N V N E
 Clos_acetobutylicum MKG T V V G T W V K T C K R L Y G E T V V E N A L E K V G F E R K K I F S P F E D V E D S K V N N F I E D I S K K V N E E K S I
 Clos_tetani MKG T I V A T W M R T C R K L Y N D D V V N K A M S S V G W D S N K I F K P T E N V E D S D L K K V I E Y I A K S E K L E L G H
 Desu_desulfuricans MRG I L P K I F M N F I K E I Y G D D V F A H V S K T M G . . . E P V F M P G N S Y P D Q V L R Q M A E I V C Q R T G E Q P K L
 Vibr_vulnificus MKG I I F T E F L E L V E E K F G L T V L D D I L D R A G D . . . E G V Y T A V G S Y D H R K L V S L I V H L S Q V T G L S V E Q
 Caul_crescentus MKG V I F N L L Q E V V S A A H G A D A W D D I L D E A G V . . . S G A Y T S L G S Y D D E E W E T L V E T A S A R L S L S R G E
 Micr_degradans MKG A V L I A L N D M V E E V F S M A V D D Q V L A K V K P D S E G I Y I S A E S Y D D A E V V G L V A L S E L T G V P V N E
 Vibr_cholerae M Q G I I Y T V L S D M V I E K F G V L F W D Q M L E D L K P S S E G V Y T S G Q Q Y N D D E L L A M V G Y L S E K A Q I P A P D
 Shew_oneidensis MKG I I F N V L E D M V V A Q C G M S V W N E L L E K H A P . K D R V Y V S A K S Y A E S E L F S I V Q D V A Q R L N M P I Q D
 Rat_beta1_sGC MYG F V N H A L E L L V I R N Y G P E V W E D I K K E A Q L D E E G Q F L V R I I Y D D S K T Y D L V A A A S K V L N L N A G E
 Rat_beta2_sGC MYG F I N T C L Q S L V T E K F G E E T W E K L K A P A E V Q D V . . . F M T Y T V Y D D I I T I K L I Q E A C K V L D V S M E A
 Nost_punctiforme MYG L V N K A I Q D M V C S R F G E E T W K Q I K H K A E V . D V D V F L S M E G Y P D D I T H K L V K A A S V I L S L S P K Q
 Nost_sp. MYG L V N K A I Q D M I S K H G E D T W E A I K Q K A G L E D I D F V G M E A Y S D D V T Y H L V G A A S E V L G K P A E E
 consensus>50 MkG.i....qdmv...ygedvwd dil...g.e.e.vf...e.ydd....lv...se.....e

αD αE αF $\beta 1$
 70 80 90 100 110 120
 Ther_tengcongensis IWRE V G R Q N I K T F S E W F P S Y F A G R . . . R L V N F L M M D E . V H L Q L T K M I K G A T P P R L I A K P V A K D .
 Clos_acetobutylicum IW E K I G E D N V I A F H K D F P A F F E H E . . . N L Y S F F K S M F D . V H V V M T K K F P G A K P P L I L I K P I S K R .
 Clos_tetani L W R Q I G K D N L V S F Y N D F P A F F Q H E . . . N L Y S F F N S L F D . I H V V M T K K F P G A K P P L V T I E P I S S K .
 Desu_desulfuricans F F E K A G R A S L Q A F N R M Y R Q Y F K G E . . . T L K E F L L A M N D . I H R H L T K D N P G V R P P K F E Y D D . Q G D .
 Vibr_vulnificus L Q E V F G E A V F D N L L A S I S N R S S L H Q C H S T F Q F I R H V E E Y I H V E V K K L Y P D A K P P E F I F I E Q D R M .
 Caul_crescentus L L R W F G Q E A M P H L A R A Y P V F F E G H V . . . S S R S F L A G V N D I I H A E V H K L Y A G A A C P H L K L R A I D A G .
 Micr_degradans L V R S F G T Y L F H Q L N S K F P I F C D L H T . . . N I F D L L S S I H G V I H K E V D K L Y S N A S L P T I N C T K L S D S .
 Vibr_cholerae L V R A Y G E Y L F T H L F N S L P E N Y P H K S . . . D L K T F L L S V D K V I H K E V Q R L Y P D A Y L P Q F E . N R V E E K .
 Shew_oneidensis V V K A F G Q F L F N G L A S R H T D V V D K F D . . . D F T S L V M G I H D V I H L E V N K L Y H E P S L P H I N G Q L L P N N .
 Rat_beta1_sGC I L Q M F G K M F F V F C Q E S G Y D T I L R V L G S N V R E F L Q N L D A . L H D H L A T I Y P G M R A P S F R C T D A E K G K
 Rat_beta2_sGC I L K L F G E Y F F K F C K M S G Y D R M L R T L G G N L T F I E N L D A . L H S Y L A L S Y Q E M N A P S F R V E E G A D G .
 Nost_punctiforme I M Q A F G E F W V Q Y T A Q E G Y G E M L D M S G D T L P E F L E N L D N . L H A R V G V S F P K L Q P P S F E C T D M E E N .
 Nost_sp. L L I A F G E Y W V T Y T S E E G Y G E L L A S A G D S L P E F M E N L D N . L H A R V G L S F P Q L R P P A F E C Q H T S S K .
 consensus>50fGe.....ll.....nl.efl...ldd.iH..v.k.y.p.a.p.p.f.....

$\beta 2$ αG $\beta 3$ $\beta 4$
 130 140 150 160 170 180
 Ther_tengcongensis A T E M E Y V S K R K . M Y D Y F L G L I E G S S K F F . K E E T S V E E V E R G E K D G F S R L K V R I K F K N P V F E Y K K N
 Clos_acetobutylicum E A I F T Y R S K R G . M F D Y L K G L I K G S A N H F . N E K I E I E V E K T K E S . . . V V L K F T F D K D I Y Y K K S F
 Clos_tetani E A I F Y Y E S K R G . M F D Y L L G L I E G S I K Y F . K E D I E I E E L E R T N E S . . . L K L K L K F Q K N I Y L K K E F
 Desu_desulfuricans T L V M T Y K S Q R D . Y G E Y F V G I I K A A A E F K . K E K V R I S S E H A G K G . . . R T T A R V T F I K
 Vibr_vulnificus K M V F D Y K S A R C . M G H V C L G L M R G C A K H F . G E E L A I Q M E T L N P T G . . . S H V R F N V A L V K G Q D G . . .
 Caul_crescentus G V A M A Y T S Q R R . M C A L A G G F T E G A A R Q F . H E V I T F E H A A C V E K G D . S A C V F H I G W P S L E A A A N D .
 Micr_degradans H L Q M R Y Y S P R K . L C V L A E G L I I G A A E H Y . K A D V S V S Q C C V H Q G A . D E C L I D V K I I
 Vibr_cholerae T L T M S Y Y S K R Q . L C A A A E G L I L G A A K Q F . N Q P V K I T Q P V C M H C G A . D H C E I V E F L P S
 Shew_oneidensis Q I A L R Y S P R R . L C F C A E G L L F G A A Q H F . Q Q K I Q I S H D T C M H T G A . D H C M L I I E L Q N D
 Rat_beta1_sGC G L I L H Y Y S E R E G L Q D I V G I I K T V A Q Q I H G T E I D M K V I Q Q R S E E C D H T Q F L I E E K E S K E E
 Rat_beta2_sGC A M L L H Y Y S D R H G L C H I V P G I E A V A K D F F D T V A M S I L D M N E E V E R T G K K E H V V F L V V Q K A H R Q I
 Nost_punctiforme S L S L H Y R S D R E G L T P M V I G L I K G L G T R F . D T E V H I T Q T Q N R D E G A E H D E F L V I Y K P N . . .
 Nost_sp. S M E L H Y Q S T R C G L A P M V L G L L H G L G K R F . Q T K V E V T Q T A F R E T G E D H D I F S I K Y E D S N L Y
 consensus>50 .l.m.Y.S.R..l.....Gli.g.a..f..eei.i.q.e.....v.f.....

How to construct Multiple Sequence Alignments?



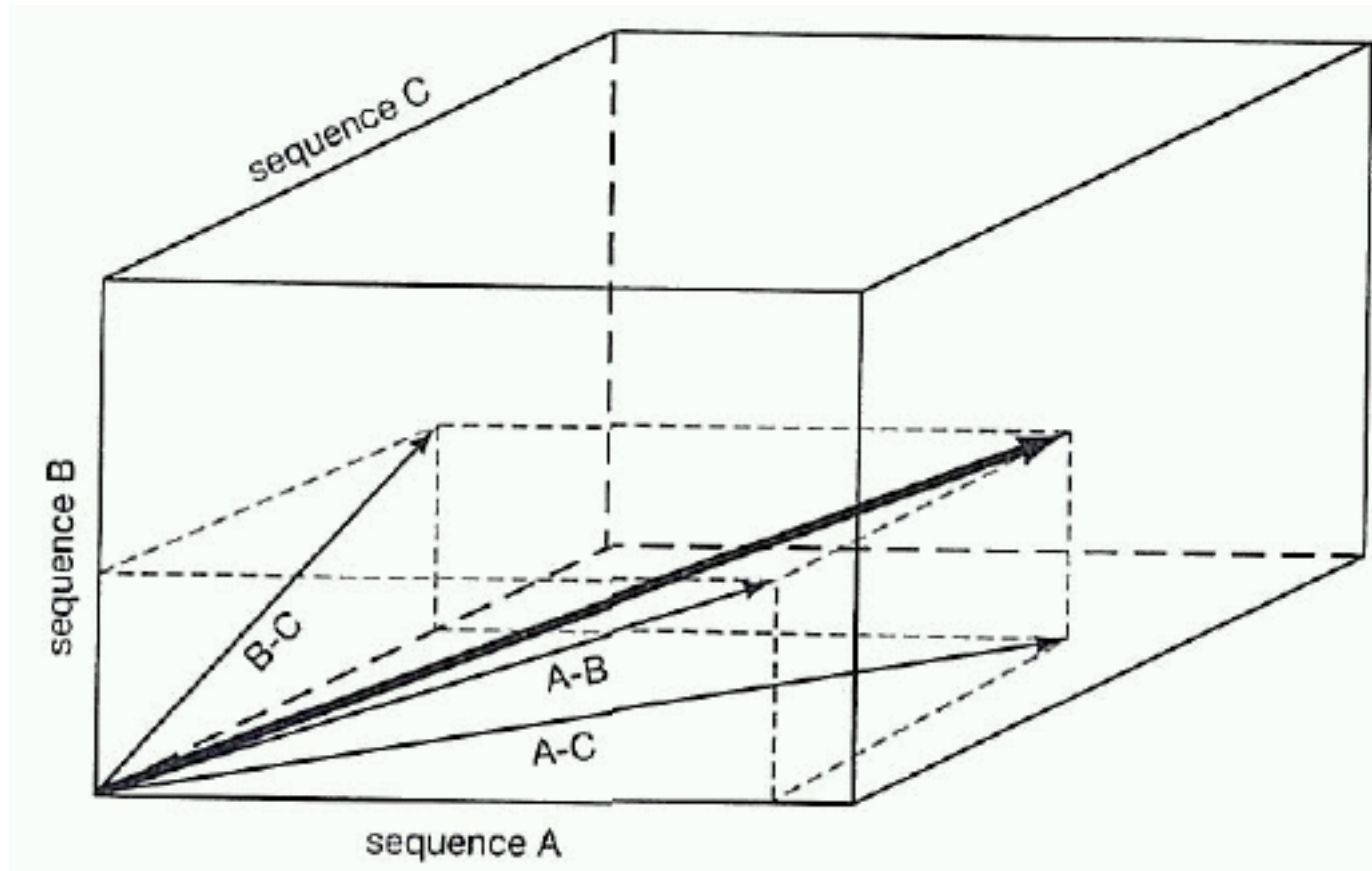
Optimal Solution:

Extend Needleman-Wunsch or Smith-Waterman to multiple sequences

How to construct Multiple Sequence Alignments?

Optimal Solution:

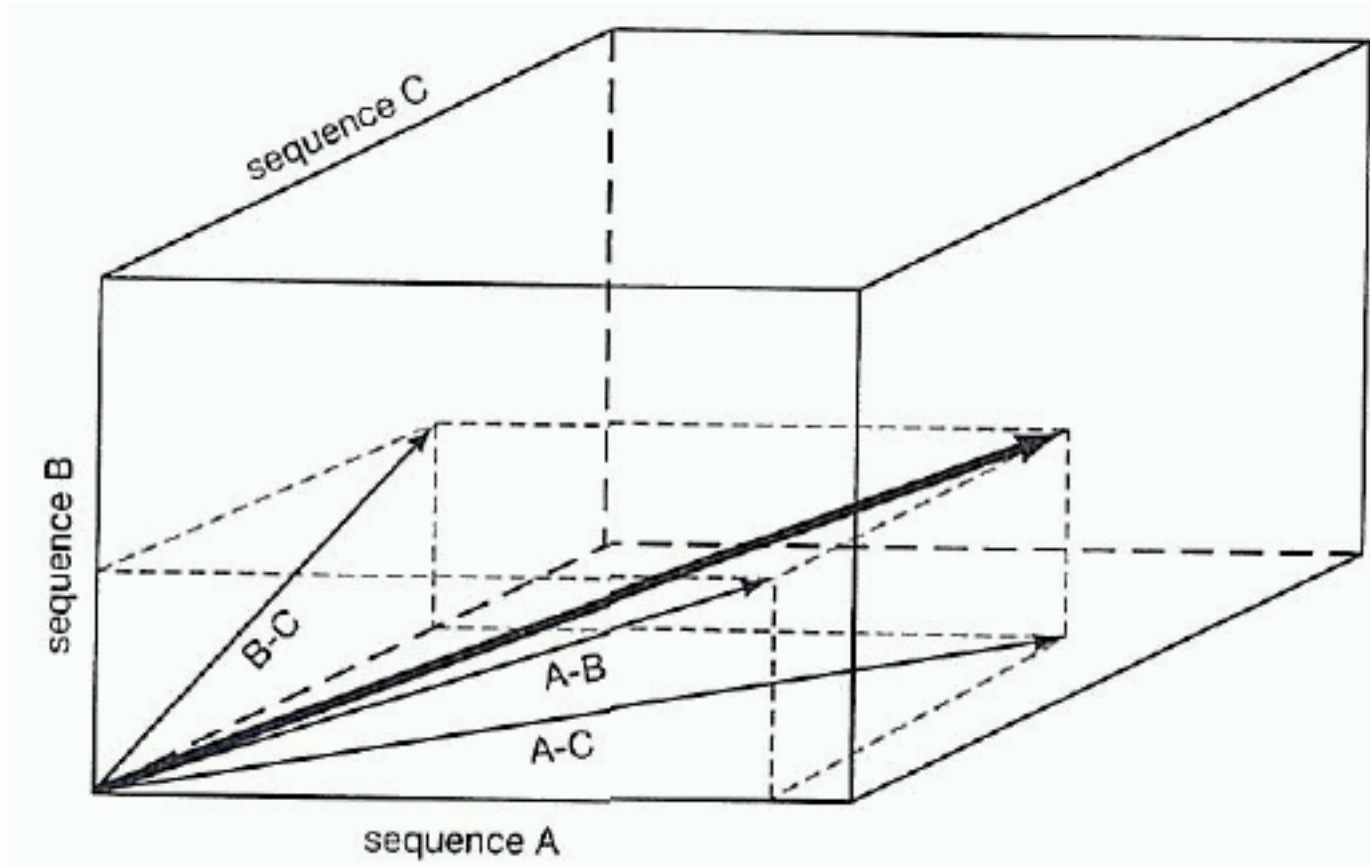
Extend Needleman-Wunsch or Smith-Waterman to multiple sequences



How to construct Multiple Sequence Alignments?

Optimal Solution:

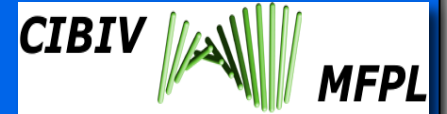
Extend Needleman-Wunsch or Smith-Waterman to multiple sequences



But $O(n^m)$ in time and memory:

Computationally not feasible... 4 sequences of length 1000 \rightarrow 1TB RAM

A new objective function: Sum of Pairs



Seq1 : AGA--CTA

Seq2 : G-A--CTT

Seq3 : AGAACTT

A new objective function: Sum of Pairs

Seq1: AGA--CTA

Seq2: G-A--CTT

Seq3: AGAACTT

Seq1: AGA--CTA
Seq2: G-A--CTT

Seq1: AGA--CTA
Seq3: AGAACTT

Seq2: G-A--CTT
Seq3: AGAACTT

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

Seq1: AGA--CTA
Seq2: G-A--CTT
Score: +5

Seq1: AGA--CTA
Seq3: AGAACTT
Score: +11

Seq2: G-A--CTT
Seq3: AGAACTT
Score: 0

A new objective function: Sum of Pairs

Seq1: AGA--CTA

Seq2: G-A--CTT

Seq3: AGAACTT

Seq1: AGA--CTA
Seq2: G-A--CTT

Seq1: AGA--CTA
Seq3: AGAACTT

Seq2: G-A--CTT
Seq3: AGAACTT

$$S(a_i, b_j) = \begin{cases} +5, & \text{if } a_i = b_j \\ -2, & \text{if } a_i \neq b_j \\ -6, & \text{for introduction of a gap} \end{cases}$$

Seq1: AGA--CTA
Seq2: G-A--CTT
Score: +5

Seq1: AGA--CTA
Seq3: AGAACTT
Score: +11

Seq2: G-A--CTT
Seq3: AGAACTT
Score: 0

SUM OF PAIRS SCORE: 16

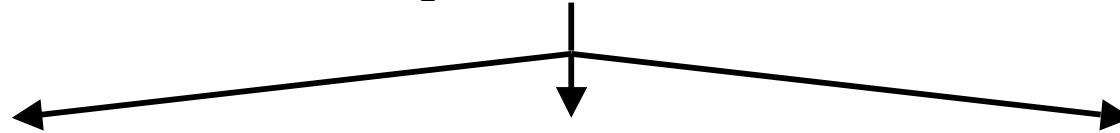
A typical variant: Weighted Sum of Pairs

Seq1 : AGA--CTA

Seq2 : AGA--CTA

Seq3 : G-A--CTT

Seq4 : AGAAACTT



Seq1 : AGA--CTA
Seq2 : AGA--CTA

Score: +30

Seq1 : AGA--CTA
Seq3 : G-A--CTT

Seq2 : AGA--CTA
Seq3 : G-A--CTT

Score: 2*(+5)

Seq1 : AGA--CTA
Seq4 : AGAAACTT

Seq2 : AGA--CTA
Seq4 : AGAAACTT

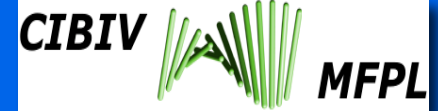
Score: 2*(+11)

Seq3 : G-A--CTT
Seq4 : AGAAACTT

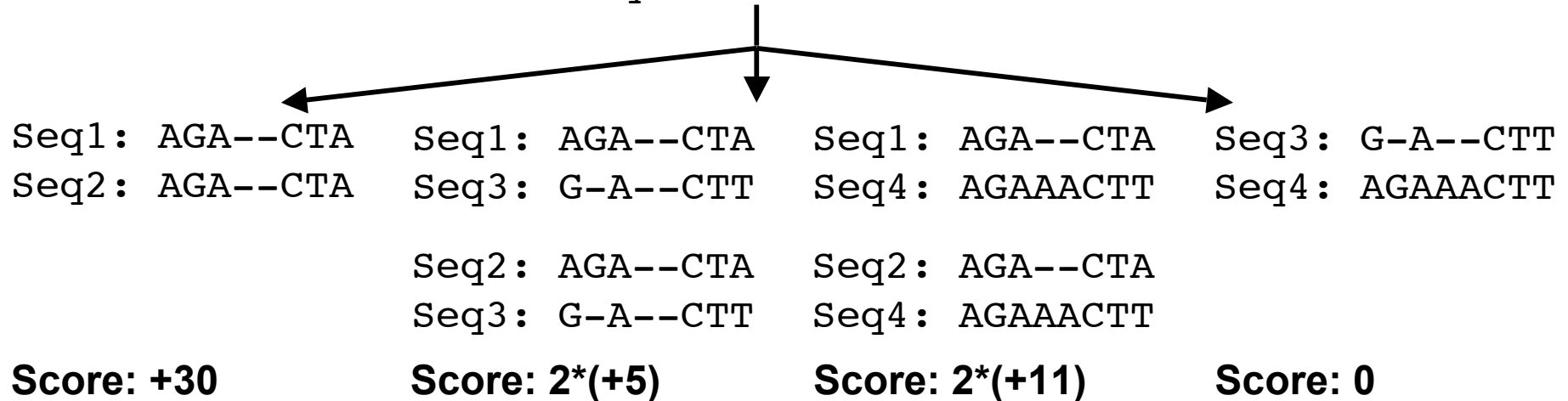
Score: 0

SUM OF PAIRS SCORE: 62

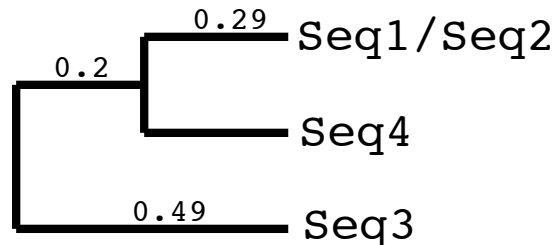
A typical variant: Weighted Sum of Pairs



Seq1 : AGA--CTA
 Seq2 : AGA--CTA
 Seq3 : G-A--CTT
 Seq4 : AGAACTT



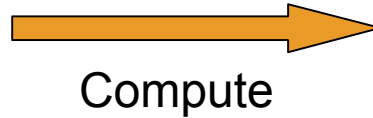
SUM OF PAIRS SCORE: 62



Weighting of sequences: one variant

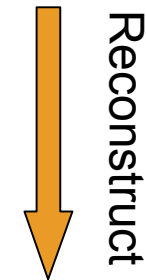
Dataset:

Seq1: AGACTA
 Seq2: AGACTA
 Seq3: GACTT
 Seq4: AGAACTT

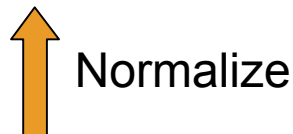


Pairwise Distance Matrix

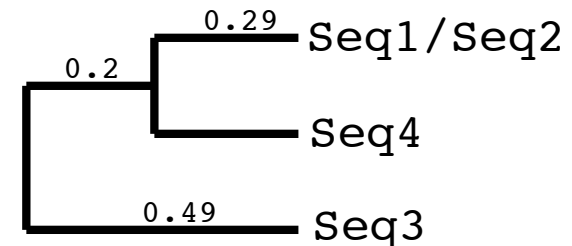
	1	2	3	4
1	-			
2		-		
3			-	
4				-



Seq1: 0.43
 Seq2: 0.43
 Seq3: 1
 Seq4: 0.73



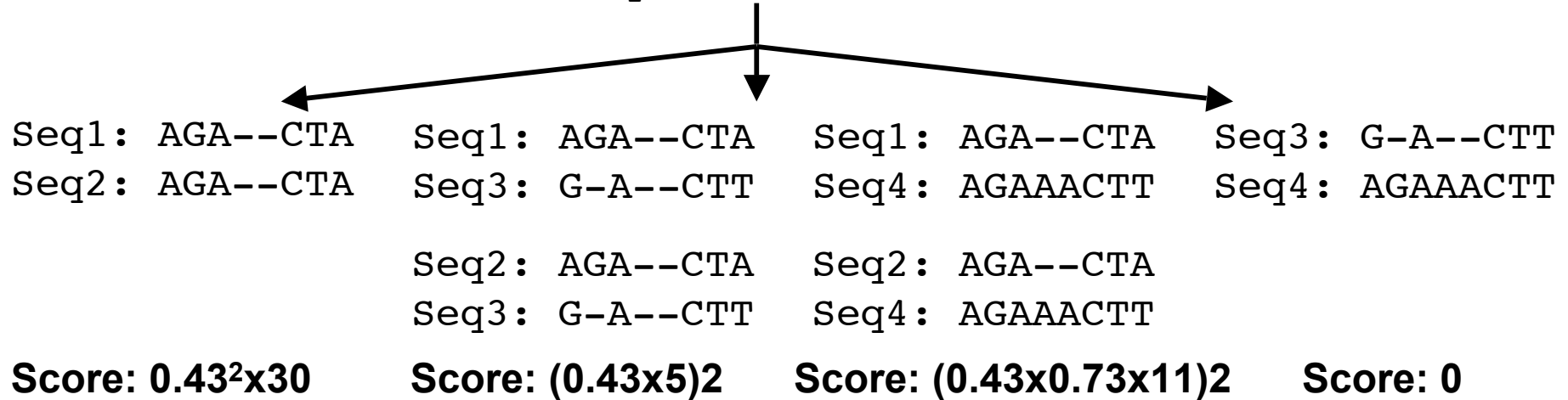
Seq1: $(0.29/2 + 0.2/3) = 0.21$
 Seq2: $(0.29/2 + 0.2/3) = 0.21$
 Seq3: 0.49
 Seq4: $(0.29 + 0.2/3) = 0.36$



A typical variant: Weighted Sum of Pairs

$$\sigma_{wsop}(\alpha) = \sum_{i < j} \omega_i \omega_j S(\alpha_i, \alpha_j)$$

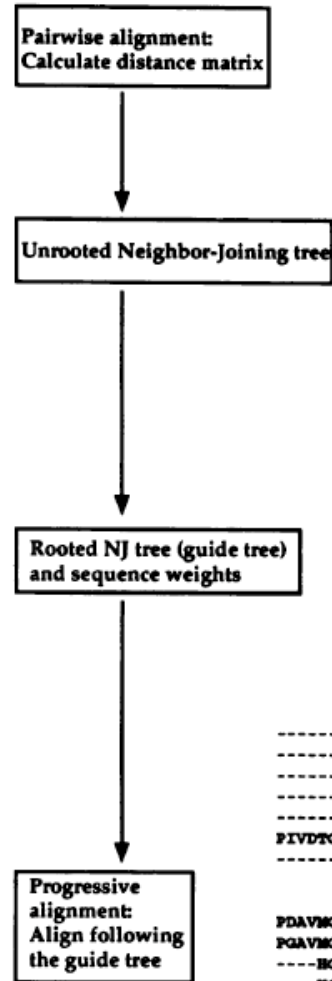
Seq1: AGA--CTA
 Seq2: AGA--CTA
 Seq3: G-A--CTT
 Seq4: AGAACTT



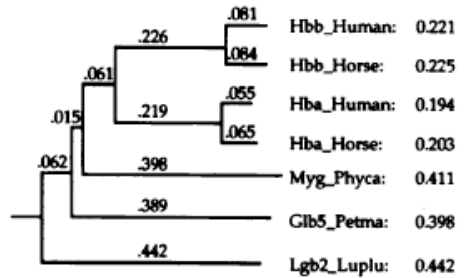
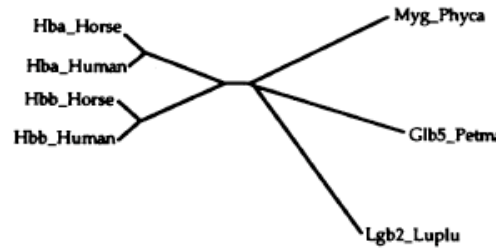
SUM OF PAIRS SCORE: 16.7

- The sequences are added stepwise. Thus, never more than two sequences (or multiple sequence alignments) are simultaneously aligned
- Sequences or MSAs are aligned using **Dynamic Programming**

Progressive Alignment Strategies (ClustalW)



Hbb_Human	1	-					
Hbb_Horse	2	.17	-				
Hba_Human	3	.59	.60	-			
Hba_Horse	4	.59	.59	.13	-		
Myg_Phyca	5	.77	.77	.75	.75	-	
Glb5_Petma	6	.81	.82	.73	.74	.80	-
Lgb2_Luplu	7	.87	.86	.86	.88	.93	.90
		1	2	3	4	5	6



```
-----VHLFPEKKAAYVIALNDKVN-----VDEVGGEALGQLLVVFFDQVFFSPQDLST
-----VQLSDEEKAAVIALNDKVN-----KEEVGGALGQLLVVFFDQVFFDQDLSE
-----VLSPADKTYFVKAANGKVDARAGHTGAALEKMFLLFFTKKFFPFDLS--
-----VLSAADKTYFVKAANGKVDARAGHTGAALEKMFLLFFTKKFFPFDLS--
-----VLSDSNQVLVAVVAKVADVAGSQDILIRLFLKSHFTLAKKDFDKELKT
P I V D T G S V A P L S A A E K T I R S A N A P V Y S E T S G V D I L V K F T T P L A Q V F F P F K O L T T
-----GALTEPQAAVLVKSSEKEMAMVPEKTRFFLIVLEIDAPAKDLSFLKGTSE
..*

PDAVNGQFVKKAKGKKVLEGGVSDQALHLD-----NLGTFPAALSRLKCDKLVLEFENFRL
PGAVNGQFVKKAKGKKVLESPGGVSDQALHLD-----NLGTFPAALSRLKCDKLVLEFENFRL
----HGSQVKKAKGKKVADALTAVAVVD-----DLPALSAALSRLKCDKLVLEFENFRL
----HGSQVKKAKGKKVGDALTAVAVVD-----DLPALSAALSRLKCDKLVLEFENFRL
EAEMKASNDLKKSGVTVLTAIGAILKQKQ-----EEMKAKPLAQSHATSKKIKPKYKAF
ADQLKKAADVGHAEKRIIMAVNDVAVSDDT--EKSEMKLRDLSGKHAESPQVLPQYKVV
VF--QNEFELQAEKGVKFLVYKAMQLQVTGVVVVTDATLKHGKQVYVYKGVVAQAEFFV
..*
```

Scoring for the alignment of two alignments

$$\sigma(a^i, b^j) = \frac{1}{n + m} \sum_{x=1}^n \sum_{y=1}^m S(a_x^i, b_y^j) \times \omega_x \times \omega_y$$

$\sigma(a^i, b^j)$: score for aligning column i from alignment (or sequence) \mathbf{a} to column j from alignment or sequence \mathbf{b}

n, m number of sequences in alignments \mathbf{a} and \mathbf{b} , respectively

$S(a_x^i, b_y^j)$ score for aligning position i in sequence \mathbf{x} from alignment \mathbf{a} to position j in sequence \mathbf{y} from alignment \mathbf{b}

ω_x, ω_y respective weights of the sequences \mathbf{x} and \mathbf{y}

Scoring for the alignment of two alignments

$$\sigma(a^i, b^j) = \frac{1}{n+m} \sum_{x=1}^n \sum_{y=1}^m S(a_x^i, b_y^j) \times \omega_x \times \omega_y$$

1 peeksavtal
2 geekaavllal
3 padktnvkaa
4 aadktnvkaa

4 egewglvlhv
5 aaektkirsa



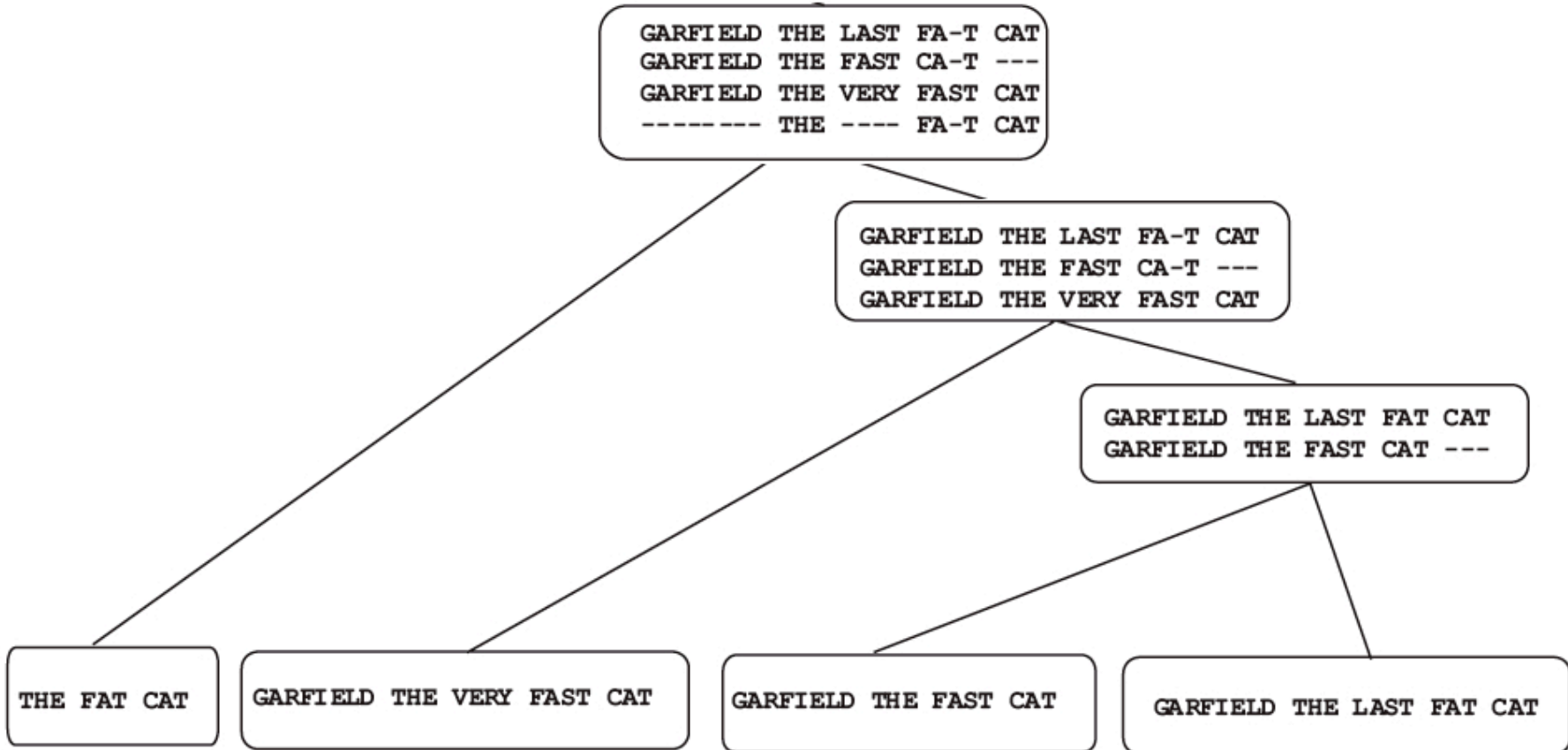
With sequence weights:

$$\begin{aligned} \text{Score} = & (S(t,v) * \omega_1 \omega_5 \\ & + S(t,i) * \omega_1 \omega_6 \\ & + S(l,v) * \omega_2 \omega_5 \\ & + S(l,i) * \omega_2 \omega_6 \\ & + S(k,v) * \omega_3 \omega_5 \\ & + S(k,i) * \omega_3 \omega_6 \\ & + S(k,v) * \omega_4 \omega_5 \\ & + S(k,i) * \omega_4 \omega_6) / 8 \end{aligned}$$

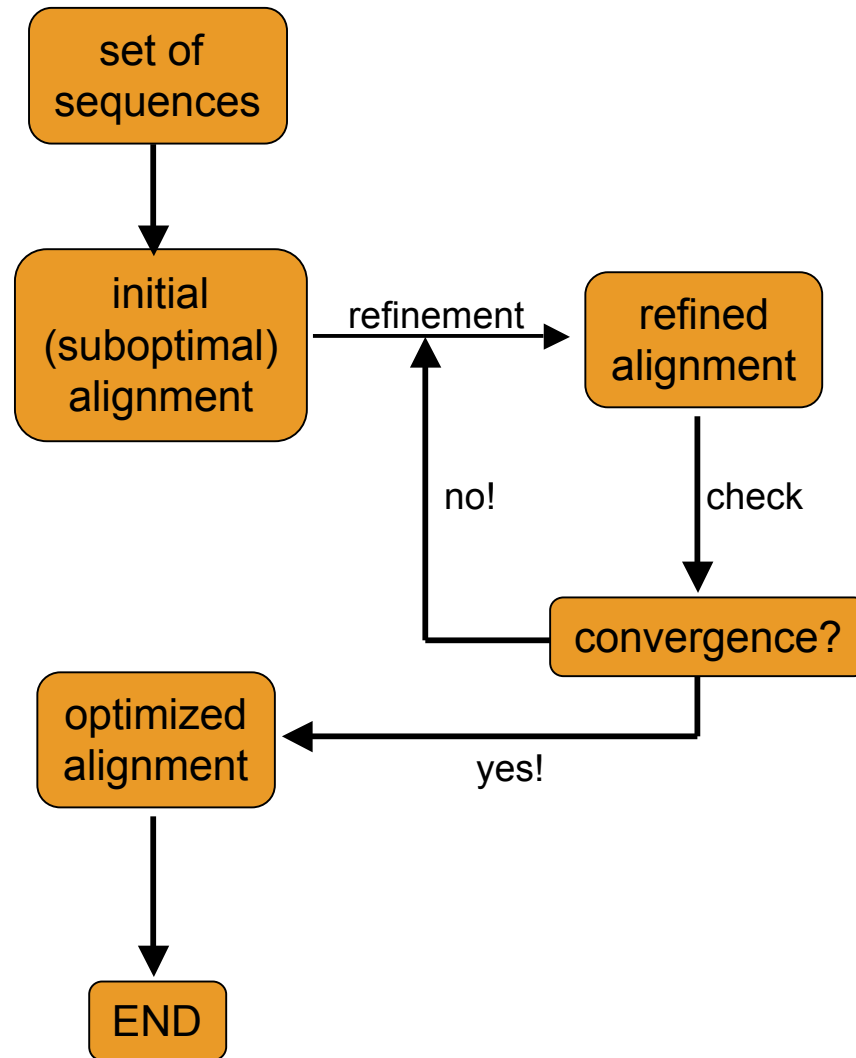
- progressive strategy
- Distance based generation of a guide tree (approximative or exact)
- tree-guided (NJ) alignment
- change of the scoring matrix as the alignment proceeds (adaptation to increasing divergence of the sequences)
- dynamic variation of gap penalties in position- and residue-specific manner
 - gap opening penalties are locally reduced in stretches of 5 or more hydrophilic residues (indicative of loop or random coil regions).
 - gap penalties are locally increased within eight residues of existing gaps.
- sequence weighting

(Known) Problem of ClustalW: Local Optima

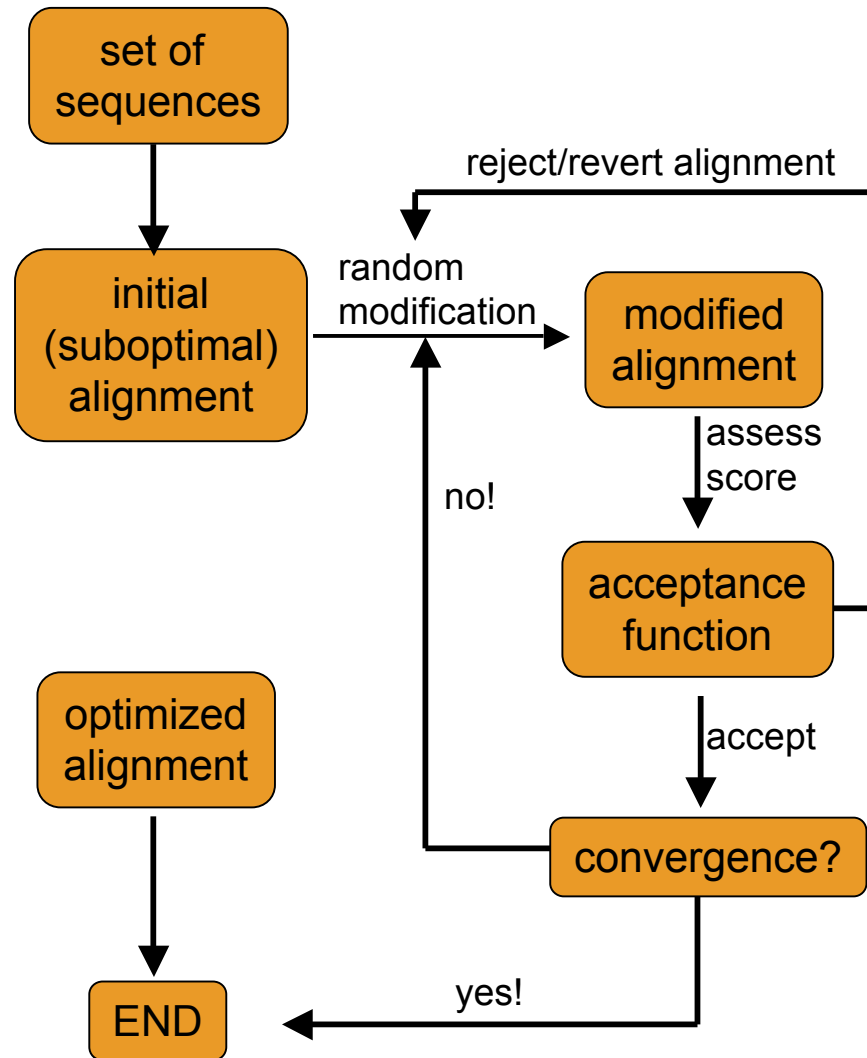
a.k.a: Once a gap always a gap



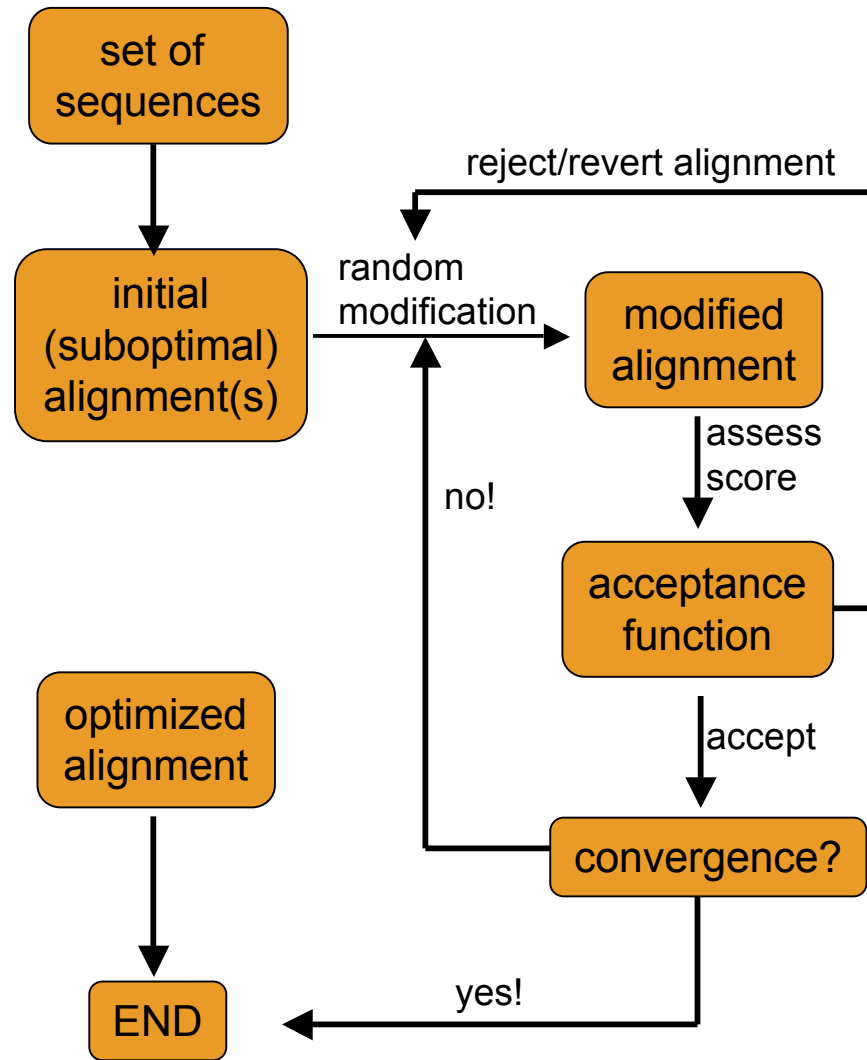
Iterative Alignment Strategy



Stochastic Iterative Alignment

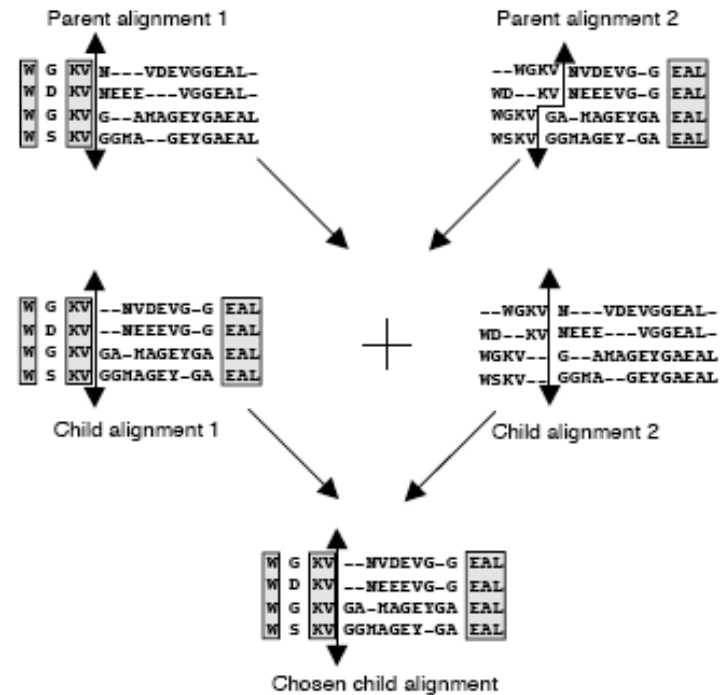


Stochastic Iterative Alignment (SAGA)



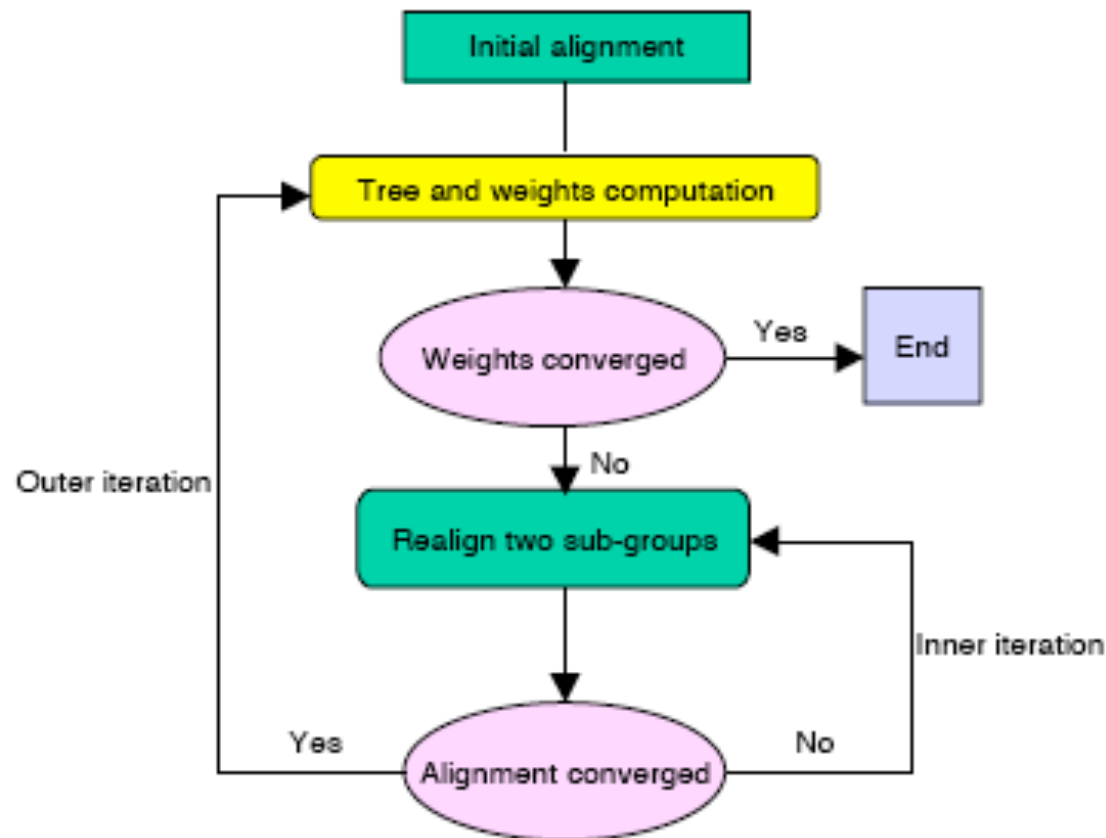
Genetic Algorithm:

- Alignments evolve by 'mutation' and crossing over
- alignments score determines fitness
- over the generations, alignments survive and reproduce or die



Non-Stochastic Iterative Alignment

Point: The initial alignment is modified by splitting it into two groups and re-aligning them with dynamic programming.



Example: Prp, both, alignment (inner loop) and tree/weight (outer loop) are optimized.

Point: The optimal MSA is defined as the one that agrees the most with all optimal pair-wise alignments

Features:

- does not depend on a specific substitution rate
- can apply any method capable to align two sequences
- position dependant, i.e. the score associated with the alignment of two residues depends on their position within the sequence rather than their individual nature
- rationale: given a set of independent observations, the constellation most often observed is often closer to the truth

Consistency based Objective Function For alignment Evaluation (COFFEE)

The Principle of T-Coffee

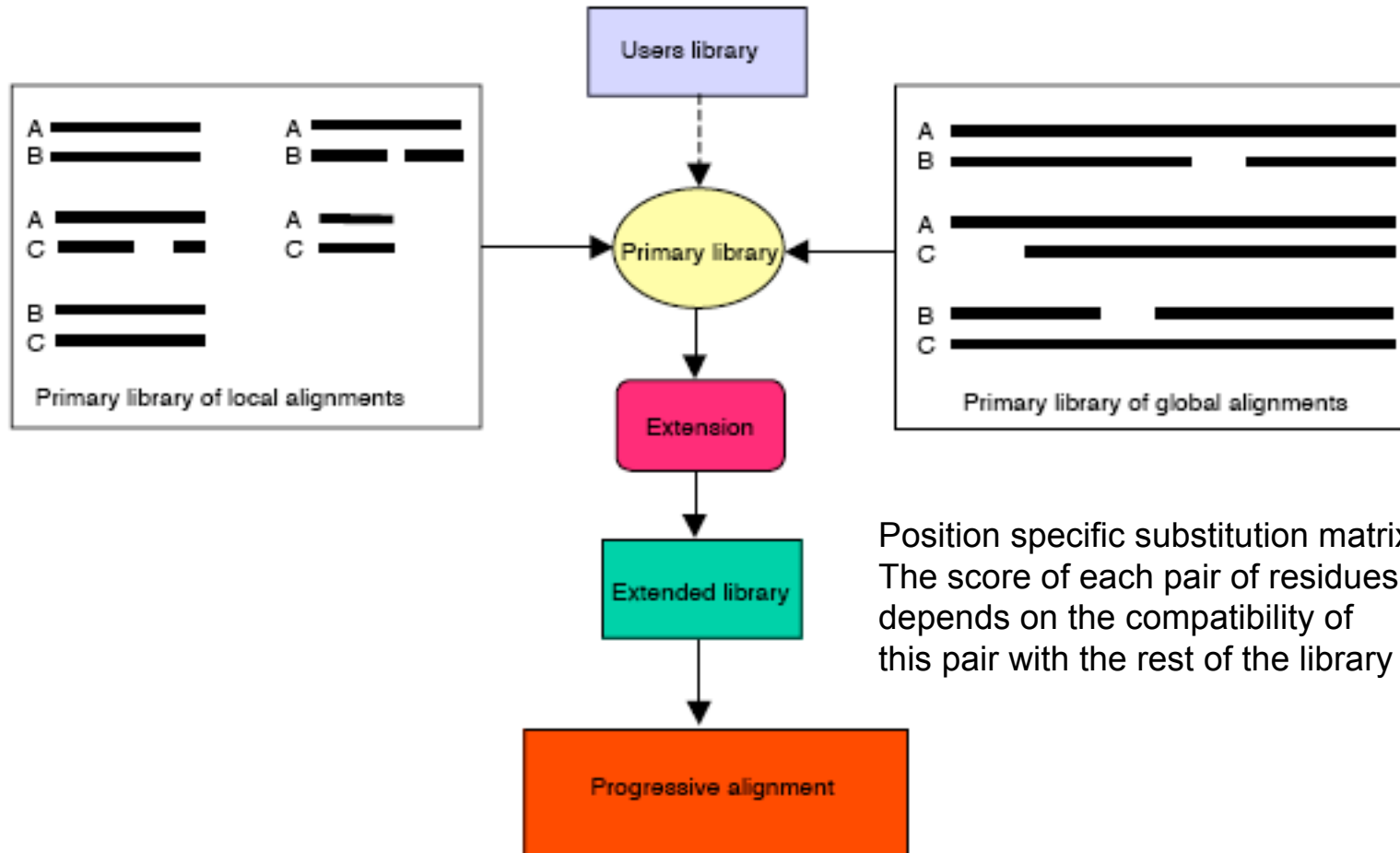
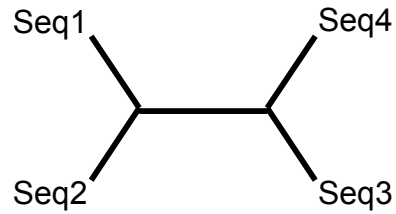


Table 2. Some elements of validation on BALiBASE.

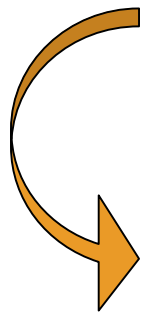
Method	Ref1	Ref2	Ref3	Ref4	Ref5	Total
DiAlign	71.0	25.2	35.1	74.7	80.4	57.3
ClustalW	78.5	32.2	42.5	65.7	74.3	58.7
Prrp	78.6	32.5	50.2	51.1	82.7	59.0
T-Coffee	80.7	37.3	52.9	83.2	88.7	68.7

Each method in the Method column was used to align the 141 test-sets contained in BALiBASE. The alignments were then compared with the reference BALiBASE alignment using `aln_compare` [34]. Ref1–5 indicates the five BALiBASE categories. Results obtained in each category were averaged. All the observed differences are statistically significant, as assessed by the Wilcoxon rank-based test [34,47]. Ref1 contains a homogenous set of sequences, ref2 contains a homogenous group of sequences and an outlayer, ref3 contains two distantly related groups of sequences. Ref4 contains sequences that require long internal gaps to be properly aligned and ref5 contains sequences that require long-terminal gaps to be properly aligned. Total is the average of ref1–5.

The Problem: Different alignments, different trees



Seq1: - N Y L S
Seq2: N K Y L S
Seq3: - N F - S
Seq4: - N F L S



N Y L S

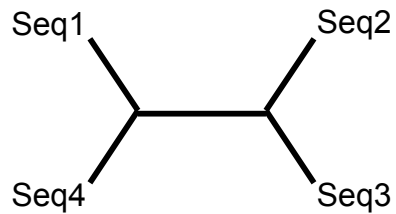
N K Y L S

N F S

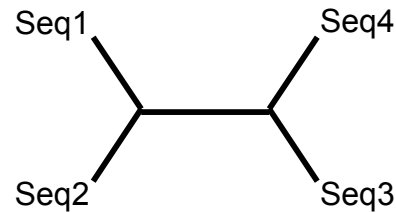
N F L S



Seq1: N - Y L S
Seq2: N K Y L S
Seq3: N - F - S
Seq4: N - F L S



The Problem: Different alignments, different trees



Seq1: - N Y L S
Seq2: N K Y L S
Seq3: - N F - S
Seq4: - N F L S

N Y L S

N K Y L S

N F S

N F L S



The alignment strategy may have more impact on the reconstructed tree than does the type of tree building method.

Morrison and Ellis (1997) Mol. Biol. Evol.
14:428-441

Gblocks (Castresana (2000) Mol. Biol. Evol. 17:540-552

Objective:

Define a set of conserved blocks from an alignment to be used in phylogeny reconstruction

Approach:

1) Classification of Columns

- non-conserved : $< n/2 + 1$ identical residues, or a gap
- conserved : $\geq n/2 + 1$ and $< 85\%$ identical residues
- highly conserved : $> 85\%$ identical residues

2) discard contiguous stretches of non-conserved positions (default $I = 8$)

3) from remaining blocks: remove flanking positions until blocks begin and end with highly conserved positions, i.e. selected blocks are anchored by positions that can be aligned with high confidence

4) discard blocks with $I < 15$

5) remove all positions with gaps together with adjacent positions until a conserved position is reached

6) discard blocks with $I < 10$

Note: all given values are the program defaults as given in the original publication

Focussing on stable parts of the alignment

