

Sequence alignment

Eva Stukenbrock & Julien Dutheil

`eva.stukenbrock@mpi-marburg.mpg.de`

`julien.dutheil@mpi-marburg.mpg.de`

Max Planck Institute for Terrestrial Microbiology – Marburg

19 February 2013

Sequence similarity

The concept of sequence similarity is used for:

- Predicting gene functions, by “homology”
- Find particular motives in a sequence
- Find particular sequences in a genome / collection
- Classify sequences into families
- Compute evolutionary distances and build phylogenies
- Assemble genomes
- *etc.*

An old problem

Definition: Edit distance

⇒ *a.k.a* Levenshtein distance • Comes for Information Theory • Is defined as:

“The minimal number of operations to turn a string into another”

An old problem

Definition: Edit distance

⇒ *a.k.a* Levenshtein distance • Comes for Information Theory • Is defined as:

“The minimal number of operations to turn a string into another”

where “operations” can be:

- A substitution of one letter with another
- An insertion of one or several letters
- A deletion of one or several letters

Application (to biology)

- Computing the edit distance between two sequences is equivalent to computing their alignment:

-I--NTERESTINGLY

<.<<.*.*.*.*>>

BIOINFORMATICS--

- ☞ 3 deletions (<), 5 substitutions (.) and 2 insertions (>).

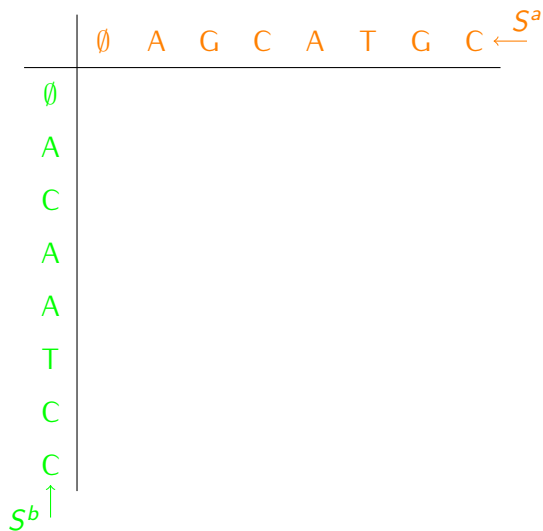
Application (to biology)

- Computing the edit distance between two sequences is equivalent to computing their alignment:
 - I--NTERESTINGLY
 - <.<<.*.*.*.*>>
 - BIOINFORMATICS--
- ☞ 3 deletions (<), 5 substitutions (.) and 2 insertions (>).
- We need to give weights to each operation to obtain a numerical value

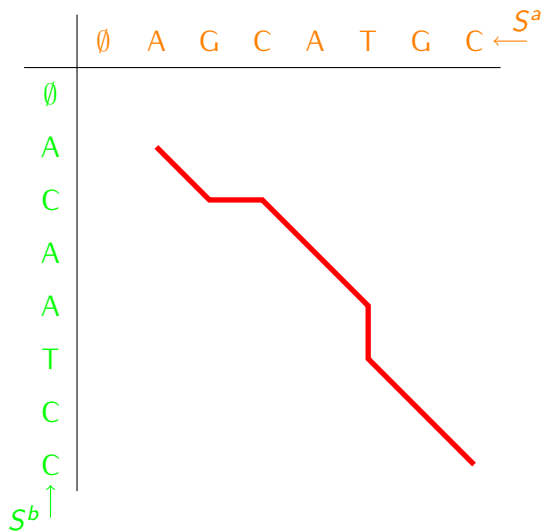
Application (to biology)

- Computing the edit distance between two sequences is equivalent to computing their alignment:
 - I--NTERESTINGLY
 - <.<<.*.*.*.*>>
 - BIOINFORMATICS--
- ☞ 3 deletions (<), 5 substitutions (.) and 2 insertions (>).
- We need to give weights to each operation to obtain a numerical value
- These operations have biological interpretation, they correspond to distinct mutational events

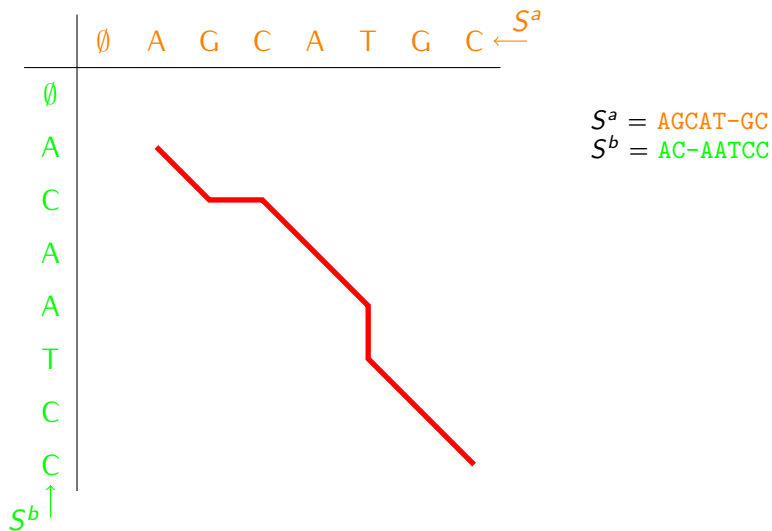
Another way to represent alignments



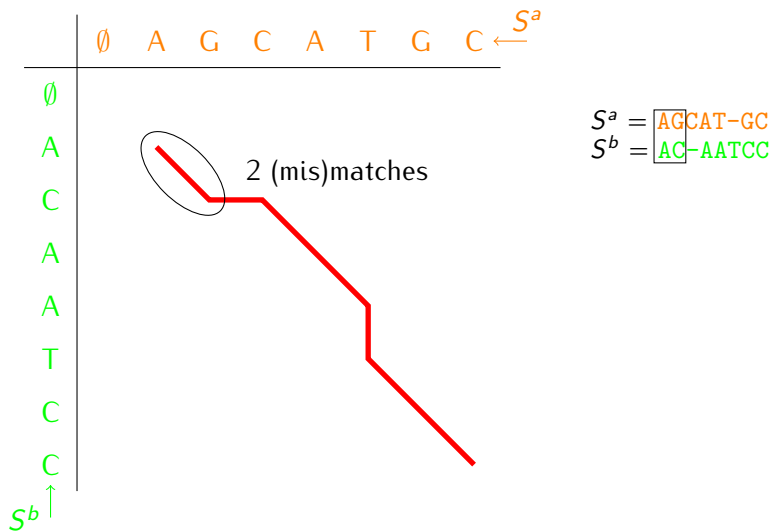
Another way to represent alignments



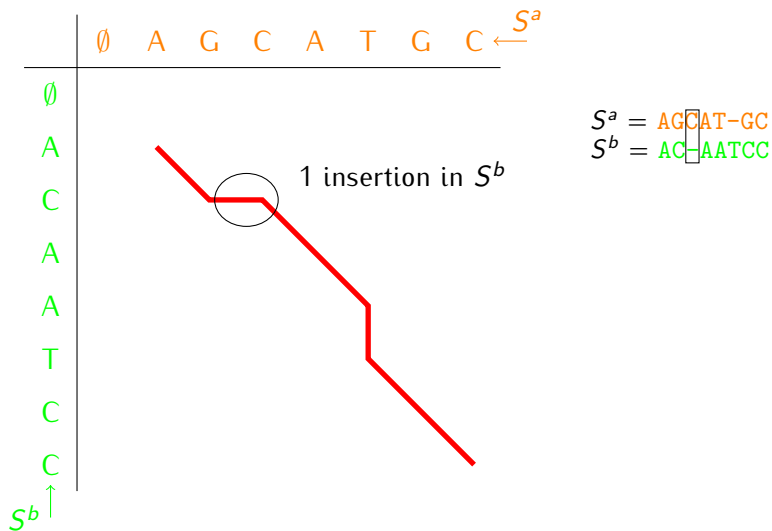
Another way to represent alignments



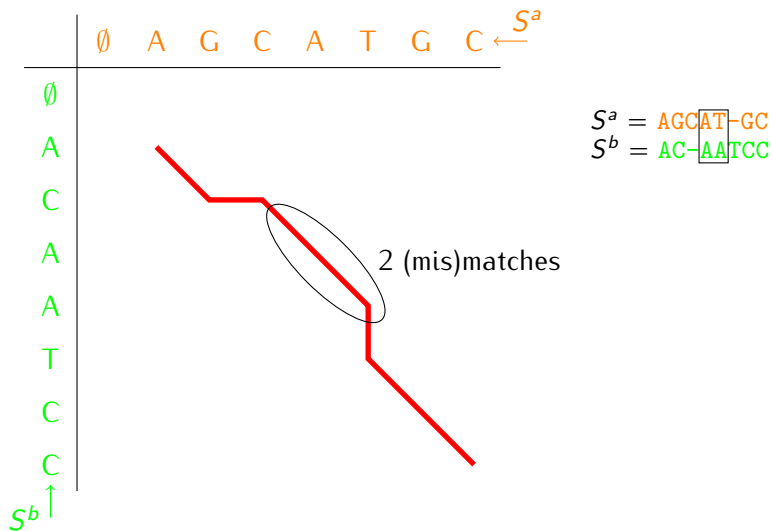
Another way to represent alignments



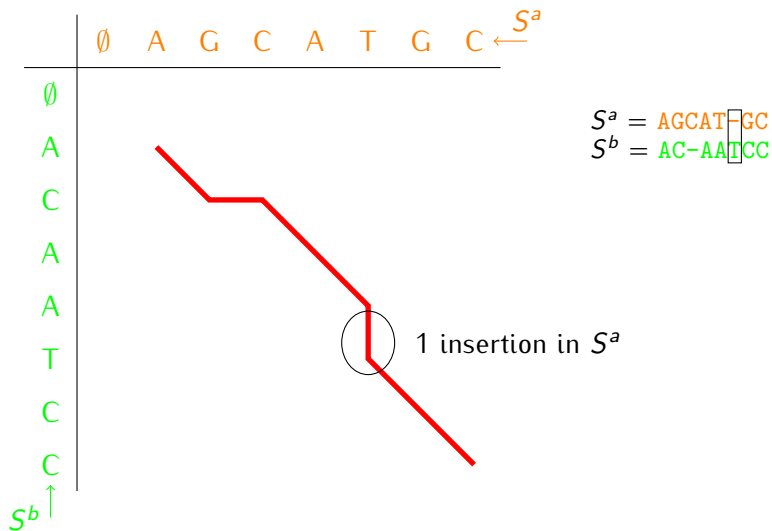
Another way to represent alignments



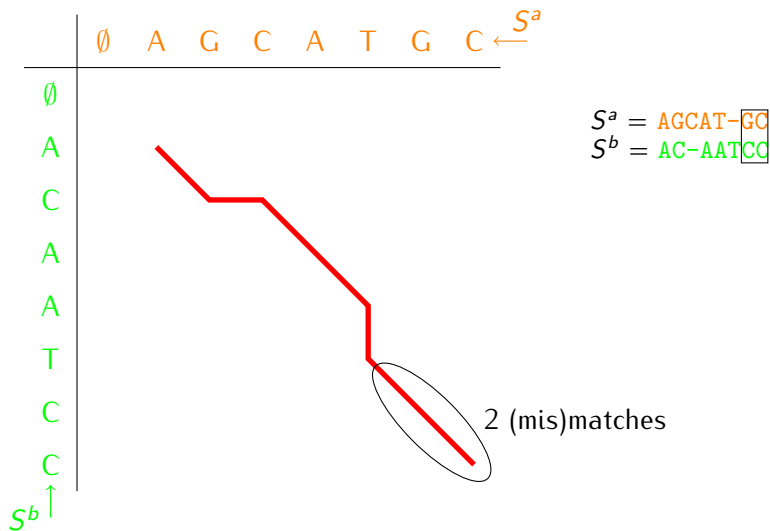
Another way to represent alignments



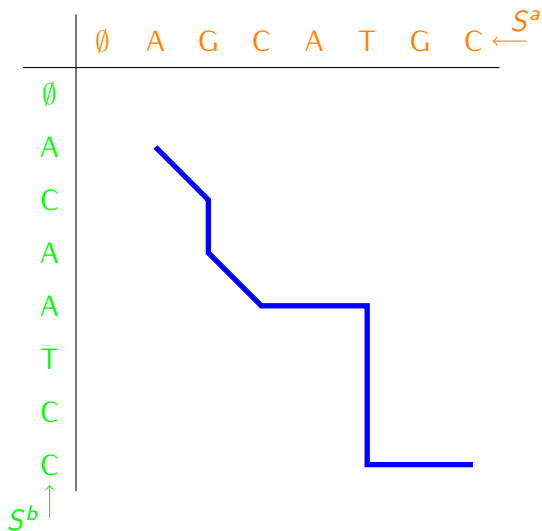
Another way to represent alignments



Another way to represent alignments

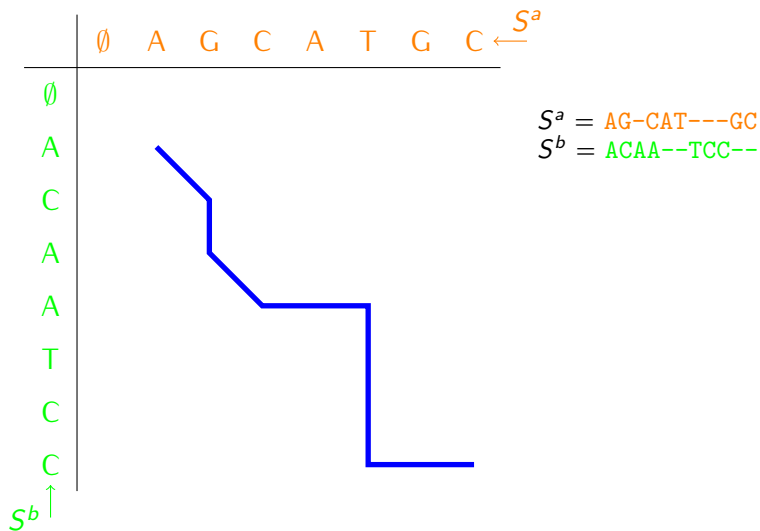


Another way to represent alignments

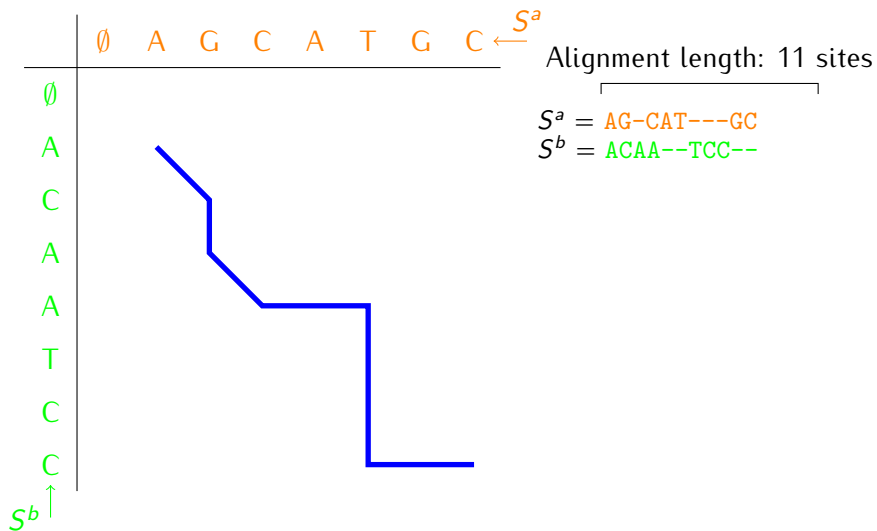


What is this alignment?

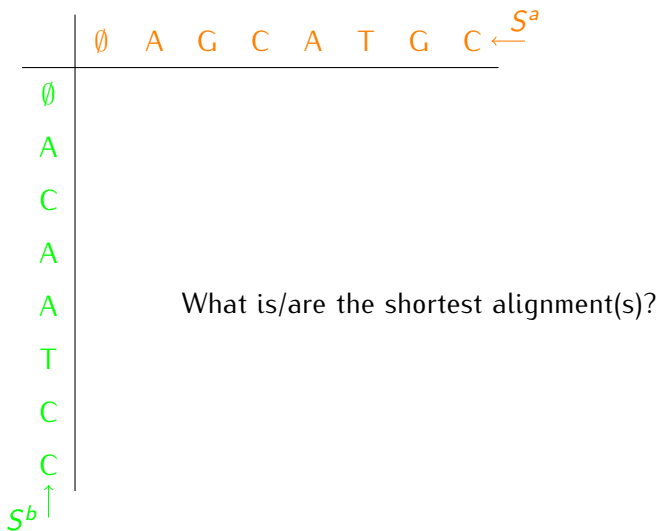
Another way to represent alignments



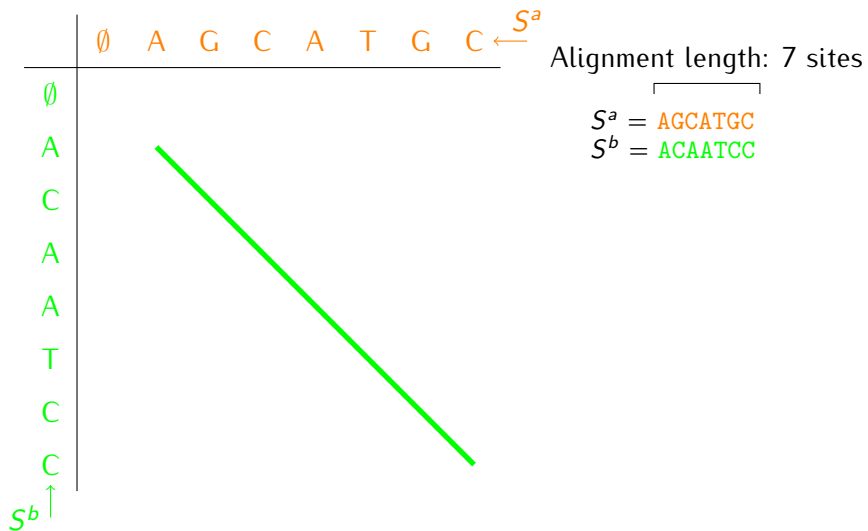
Another way to represent alignments



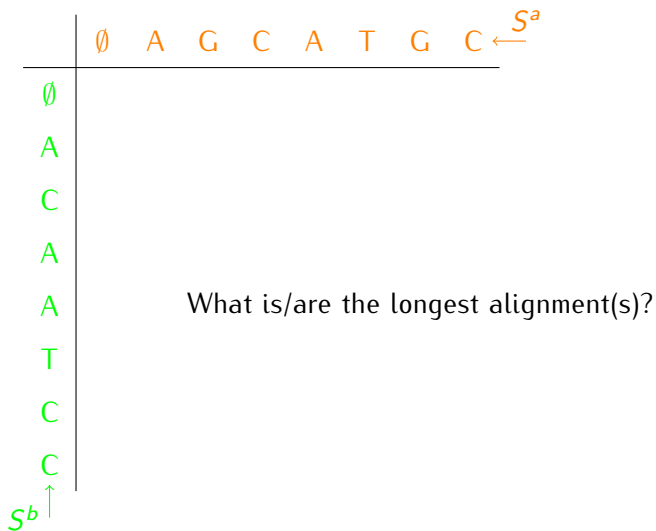
Another way to represent alignments



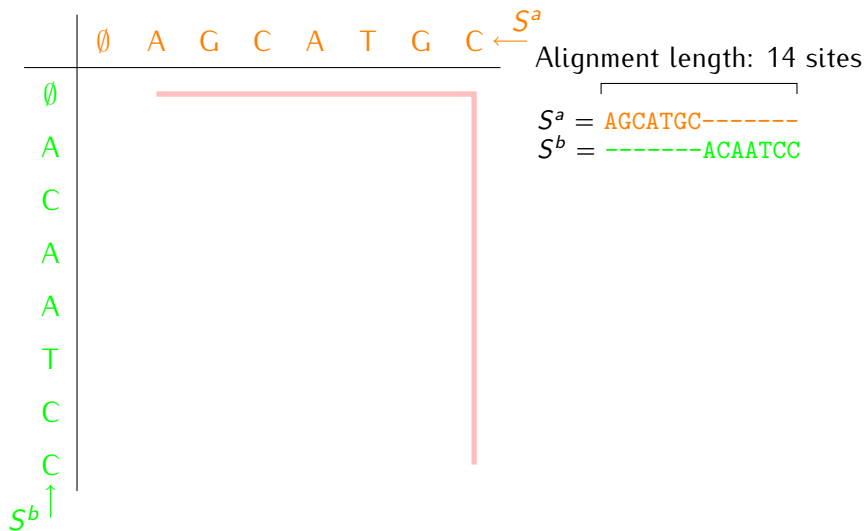
Another way to represent alignments



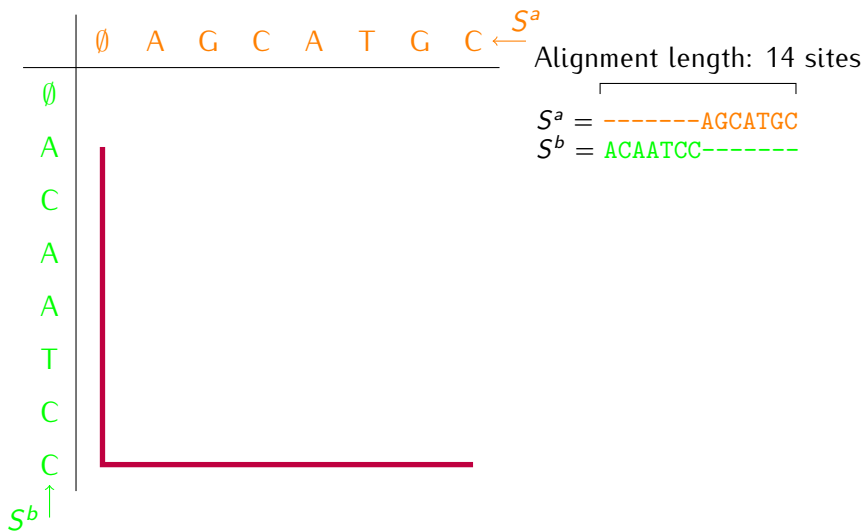
Another way to represent alignments



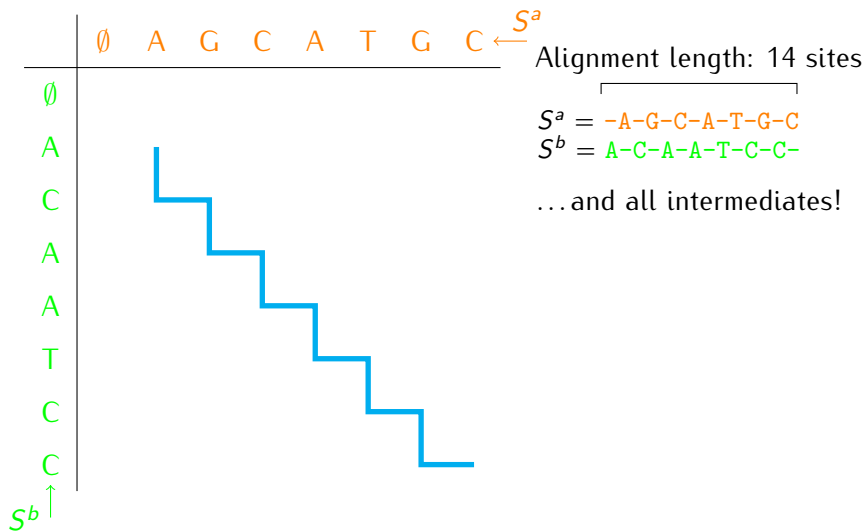
Another way to represent alignments



Another way to represent alignments

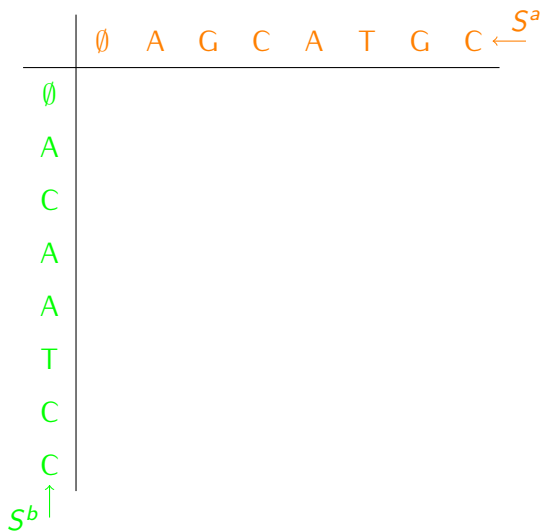


Another way to represent alignments



What is the best alignment?

Part1: alignment score

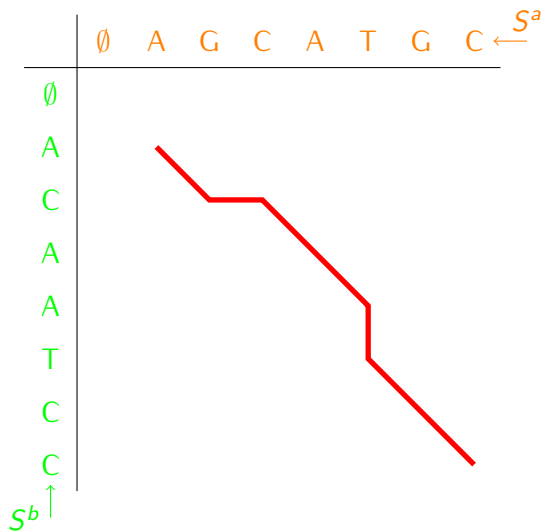


Scores:

	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

What is the best alignment?

Part1: alignment score



Scores:

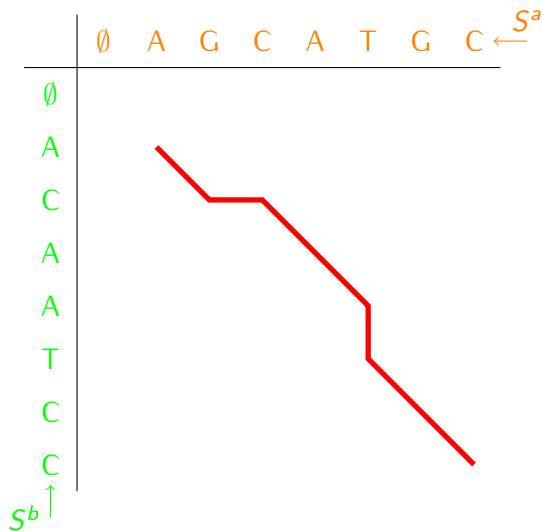
	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

$S^a =$ A G C A T - G C

$S^b =$ A C - A A T C C

What is the best alignment?

Part1: alignment score



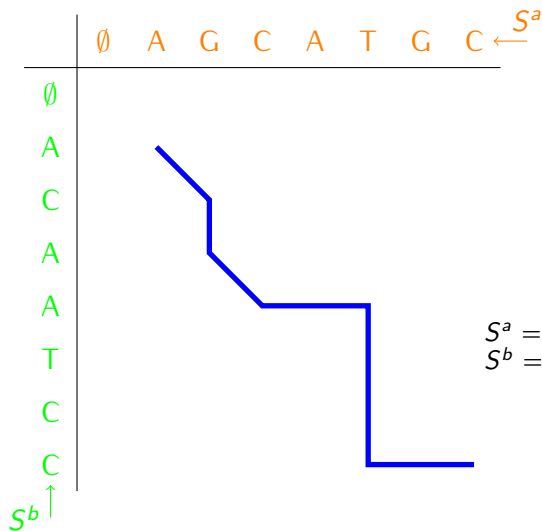
Scores:

	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

$$\begin{aligned} S^a &= \text{A G C A T - G C} \\ S^b &= \text{A C - A A T C C} \\ &+2-1-1+2-1-1-1+2 \\ \text{Score} &= 1 \end{aligned}$$

What is the best alignment?

Part1: alignment score



Scores:

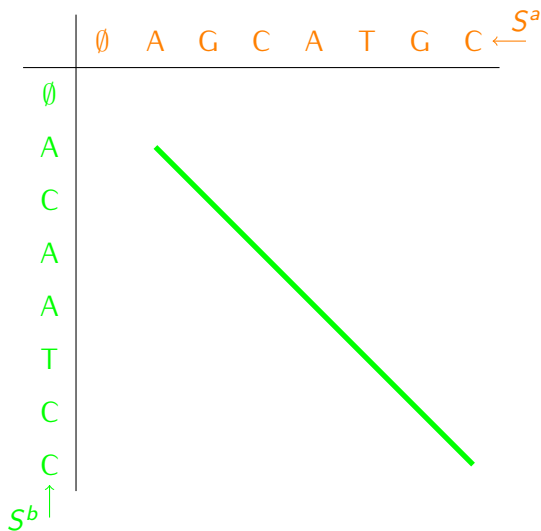
	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

$$\begin{aligned}
 S^a &= \text{A G - C A T - - - G C} \\
 S^b &= \text{A C A A - - T C C - -} \\
 &\quad +2-1-1-1-1-1-1-1-1-1-1-1
 \end{aligned}$$

$$\text{Score} = -9$$

What is the best alignment?

Part1: alignment score



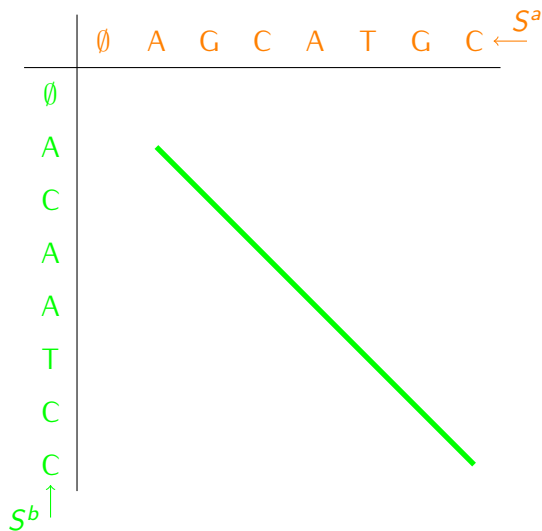
Scores:

	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

$S^a = \text{A G C A T G C}$
 $S^b = \text{A C A A T C C}$

What is the best alignment?

Part1: alignment score



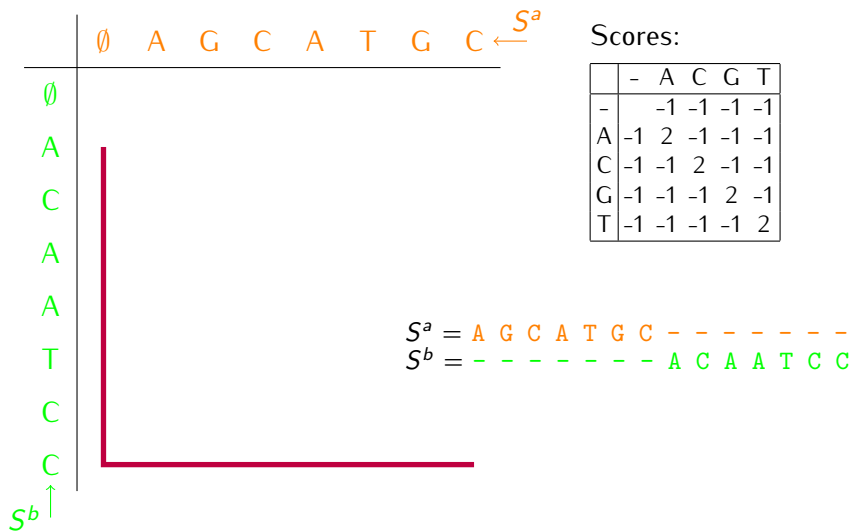
Scores:

	-	A	C	G	T
-		-1	-1	-1	-1
A	-1	2	-1	-1	-1
C	-1	-1	2	-1	-1
G	-1	-1	-1	2	-1
T	-1	-1	-1	-1	2

$$\begin{aligned} S^a &= \text{A G C A T G C} \\ S^b &= \text{A C A A T C C} \\ &+2-1-1+2+2-1+2 \\ \text{Score} &= 5 \end{aligned}$$

What is the best alignment?

Part1: alignment score



What is the best alignment?

Part2: How to find the best score

The brute force

- Define scores for insertions, deletions and substitutions

What is the best alignment?

Part2: How to find the best score

The brute force

- Define scores for insertions, deletions and substitutions
- Generate a possible alignment of the two sequences

What is the best alignment?

Part2: How to find the best score

The brute force

- Define scores for insertions, deletions and substitutions
- Generate a possible alignment of the two sequences
- Computes the edit distance for this alignment, and record it

What is the best alignment?

Part2: How to find the best score

The brute force

- Define scores for insertions, deletions and substitutions
- Generate a possible alignment of the two sequences
- Computes the edit distance for this alignment, and record it
- Generate another alignment, and compute its edit distance, *etc*

What is the best alignment?

Part2: How to find the best score

The brute force

- Define scores for insertions, deletions and substitutions
- Generate a possible alignment of the two sequences
- Computes the edit distance for this alignment, and record it
- Generate another alignment, and compute its edit distance, *etc*
- Do that for all possible alignments, and keep the one with the highest score (=minimal edit distance)!

What is the best alignment?

Part2: How to find the best score

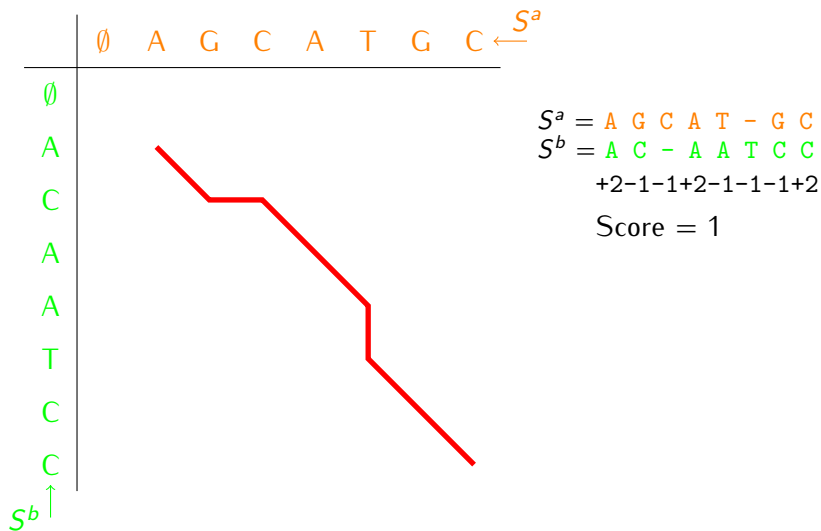
The brute force

- Define scores for insertions, deletions and substitutions
- Generate a possible alignment of the two sequences
- Computes the edit distance for this alignment, and record it
- Generate another alignment, and compute its edit distance, *etc*
- Do that for all possible alignments, and keep the one with the highest score (=minimal edit distance)!

☞ extremely inefficient, as the number for possible alignments is huge!

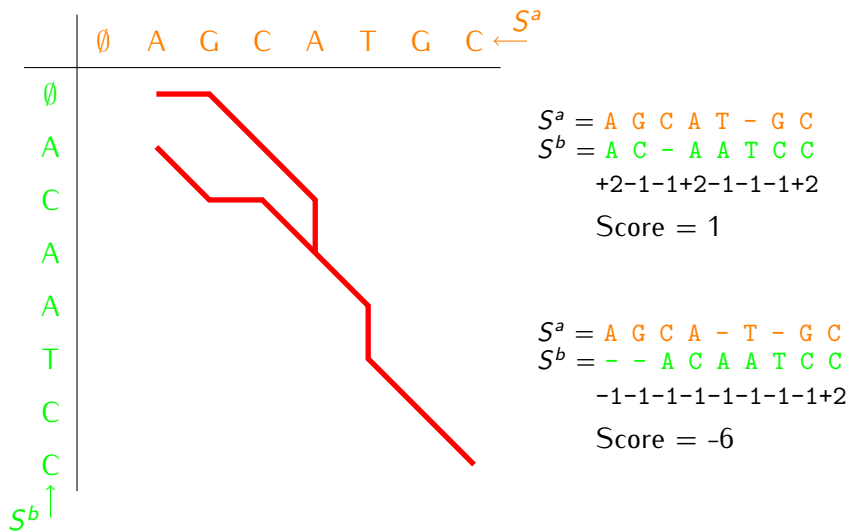
What is the best alignment?

Part2: Trivial non-optimal alignments



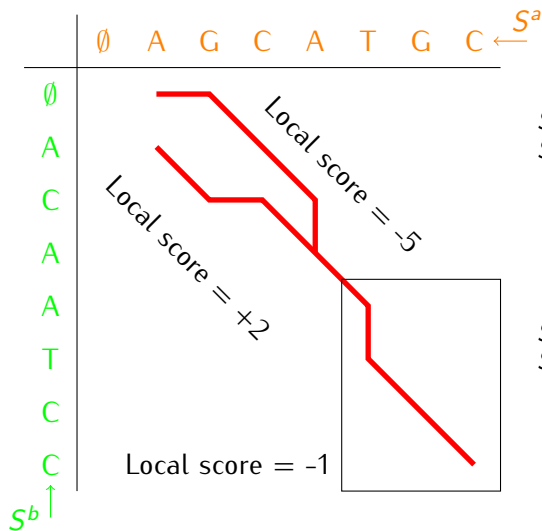
What is the best alignment?

Part2: Trivial non-optimal alignments



What is the best alignment?

Part2: Trivial non-optimal alignments



$$\begin{array}{r}
 S^a = \text{A G C A T - G C} \\
 S^b = \text{A C - A A T C C} \\
 +2-1-1+2-1-1-1+2 \\
 \text{Score} = 1
 \end{array}$$

$$\begin{array}{r}
 S^a = \text{A G C A - T - G C} \\
 S^b = \text{- - A C A A T C C} \\
 -1-1-1-1-1-1-1-1+2 \\
 \text{Score} = -6
 \end{array}$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0							
A								
C								
A								
A								
T								
C								
C								

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	\emptyset	A	G	C	A	T	G	C
\emptyset	0	$\leftarrow -1$						
A								
C								
A								
A								
T								
C								
C								

$$S_1^a = A$$

$$S_0^b = -$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2					
A								
C								
A								
A								
T								
C								
C								

$$S_2^a = \text{AG}$$

$$S_0^b = \text{--}$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3				
A								
C								
A								
A								
T								
C								
C								

$$S_3^a = \text{AGC}$$

$$S_0^b = \text{---}$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A								
C								
A								
A								
T								
C								
C								

$$S_7^a = \text{AGCATGC}$$

$$S_0^b = \text{-----}$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1							
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

$$S_0^a = \text{-----}$$
$$S_7^b = \text{ACAATCC}$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	?						
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_1^a = A$$

$$S_1^b = A$$

$$0 + 2 = 2$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	↑ -1						
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$$S_1^a = A$$

$$0 + 2 = 2$$

$$S_1^b = A$$

Insertion in S^a :

$$S_1^a = A-$$

$$-1 - 1 = -2$$

$$S_1^b = -A$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	← -1						
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$$S_1^a = A$$

$$0 + 2 = 2$$

$$S_1^b = A$$

Insertion in S^a :

$$S_1^a = A-$$

$$-1 - 1 = -2$$

$$S_1^b = -A$$

Insertion in S^b :

$$S_1^a = -A$$

$$-1 - 1 = -2$$

$$S_1^b = A-$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	2						
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$$S_1^a = A$$

$$S_1^b = A$$

$$0 + 2 = 2$$

Insertion in S^a :

$$S_1^a = A-$$

$$S_1^b = -A$$

$$-1 - 1 = -2$$

Insertion in S^b :

$$S_1^a = -A$$

$$S_1^b = A-$$

$$-1 - 1 = -2$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$$S_2^a = \text{AG}$$

$$S_1^b = \text{-A}$$

$$-1 - 1 = -2$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_2^a = AG$$

$$S_1^b = -A$$

$$-1 - 1 = -2$$

Insertion in S^a :

$$S_2^a = AG-$$

$$S_1^b = --A$$

$$-2 - 1 = -3$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_2^a = \text{AG}$$

$$S_1^b = \text{-A}$$

$$-1 - 1 = -2$$

Insertion in S^a :

$$S_2^a = \text{AG-}$$

$$S_1^b = \text{--A}$$

$$-2 - 1 = -3$$

Insertion in S^b :

$$S_2^a = \text{AG}$$

$$S_1^b = \text{A-}$$

$$2 - 1 = 1$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_2^a = \text{AG}$$

$$S_1^b = \text{-A}$$

$$-1 - 1 = -2$$

Insertion in S^a :

$$S_2^a = \text{AG-}$$

$$S_1^b = \text{--A}$$

$$-2 - 1 = -3$$

Insertion in S^b :

$$S_2^a = \text{AG}$$

$$S_1^b = \text{A-}$$

$$2 - 1 = 1$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$$S_3^a = \text{AGC}$$

$$S_1^b = \text{--A}$$

$$-2 - 1 = -3$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_3^a = \text{AGC}$$

$$S_1^b = \text{--A}$$

$$-2 - 1 = -3$$

Insertion in S^a :

$$S_3^a = \text{AGC-}$$

$$S_1^b = \text{---A}$$

$$-3 - 1 = -4$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_3^a = \text{AGC} \quad -2 - 1 = -3$$

$$S_1^b = \text{--A}$$

Insertion in S^a :

$$S_3^a = \text{AGC-}$$

$$S_1^b = \text{---A} \quad -3 - 1 = -4$$

Insertion in S^b :

$$S_3^a = \text{AGC}$$

$$S_1^b = \text{A--} \quad 1 - 1 = 0$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	2	← 1	← 0				
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$$S_3^a = \text{AGC}$$

$$S_1^b = \text{--A}$$

$$-2 - 1 = -3$$

Insertion in S^a :

$$S_3^a = \text{AGC-}$$

$$S_1^b = \text{---A}$$

$$-3 - 1 = -4$$

Insertion in S^b :

$$S_3^a = \text{AGC}$$

$$S_1^b = \text{A--}$$

$$1 - 1 = 0$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	2	← 1	← 0	← -1			
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$S_4^a = \text{AGCA}$

$S_1^b = \text{---A}$

$$-3 + 2 = -1$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	2	← 1	← 0	↑ -5			
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$$S_4^a = \text{AGCA}$$

$$S_1^b = \text{---A}$$

$$-3 + 2 = -1$$

Insertion in S^a :

$$S_4^a = \text{AGCA-}$$

$$S_1^b = \text{----A}$$

$$-4 - 1 = -5$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	-1	↑	← 2	← 1	← 0	← -1		
C	-2	↑						
A	-3	↑						
A	-4	↑						
T	-5	↑						
C	-6	↑						
C	-7	↑						

Match / Mismatch:

$$S_4^a = \text{AGCA}$$

$$S_1^b = \text{---A}$$

$$-3 + 2 = -1$$

Insertion in S^a :

$$S_4^a = \text{AGCA-}$$

$$S_1^b = \text{----A}$$

$$-4 - 1 = -5$$

Insertion in S^b :

$$S_4^a = \text{AGCA}$$

$$S_1^b = \text{A---}$$

$$0 - 1 = -1$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_4^a = \text{AGCA}$$

$$S_1^b = \text{---A}$$

$$-3 + 2 = -1$$

Insertion in S^a :

$$S_4^a = \text{AGCA-}$$

$$S_1^b = \text{----A}$$

$$-4 - 1 = -5$$

Insertion in S^b :

$$S_4^a = \text{AGCA}$$

$$S_1^b = \text{A---}$$

$$0 - 1 = -1$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$S_5^a = \text{AGCAT}$

$S_1^b = \text{-----A}$

$$-4 - 1 = -5$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$S_5^a = \text{AGCAT}$

$S_1^b = \text{-----A}$

$$-4 - 1 = -5$$

Insertion in S^a :

$S_5^a = \text{AGCAT-}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	↑ -1	2	← 1	← 0	← -1	← -2		
C	↑ -2							
A	↑ -3							
A	↑ -4							
T	↑ -5							
C	↑ -6							
C	↑ -7							

Match / Mismatch:

$S_5^a = \text{AGCAT}$

$S_1^b = \text{-----A}$

$$-4 - 1 = -5$$

Insertion in S^a :

$S_5^a = \text{AGCAT-}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^b :

$S_5^a = \text{AGCAT}$

$S_1^b = \text{A-----}$

$$-1 - 1 = -2$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	-1	↑	← 2	← 1	← 0	← -1	← -2	
C	-2	↑						
A	-3	↑						
A	-4	↑						
T	-5	↑						
C	-6	↑						
C	-7	↑						

Match / Mismatch:

$S_5^a = \text{AGCAT}$

$S_1^b = \text{-----A}$

$$-4 - 1 = -5$$

Insertion in S^a :

$S_5^a = \text{AGCAT-}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^b :

$S_5^a = \text{AGCAT}$

$S_1^b = \text{A-----}$

$$-1 - 1 = -2$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$S_6^a = \text{AGCATG}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$S_6^a = \text{AGCATG}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^a :

$S_6^a = \text{AGCATG-}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	-1	↑	← 2	← 1	← 0	← -1	← -2	← -3
C	-2	↑	↑					
A	-3	↑	↑					
A	-4	↑	↑					
T	-5	↑	↑					
C	-6	↑	↑					
C	-7	↑	↑					

Match / Mismatch:

$S_6^a = \text{AGCATG}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^a :

$S_6^a = \text{AGCATG-}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^b :

$S_6^a = \text{AGCATG}$

$S_1^b = \text{A-----}$

$$-2 - 1 = -3$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0	← -1	← -2	← -3	← -4	← -5	← -6	← -7
A	-1	↑	← 2	← 1	← 0	← -1	← -2	← -3
C	-2	↑						
A	-3	↑						
A	-4	↑						
T	-5	↑						
C	-6	↑						
C	-7	↑						

Match / Mismatch:

$S_6^a = \text{AGCATG}$

$S_1^b = \text{-----A}$

$$-5 - 1 = -6$$

Insertion in S^a :

$S_6^a = \text{AGCATG-}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^b :

$S_6^a = \text{AGCATG}$

$S_1^b = \text{A-----}$

$$-2 - 1 = -3$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^a :

$S_7^a = \text{AGCATGC-}$

$S_1^b = \text{-----A}$

$$-7 - 1 = -8$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^a :

$S_7^a = \text{AGCATGC-}$

$S_1^b = \text{-----A}$

$$-7 - 1 = -8$$

Insertion in S^b :

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{A-----}$

$$-3 - 1 = -4$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{-----A}$

$$-6 - 1 = -7$$

Insertion in S^a :

$S_7^a = \text{AGCATGC-}$

$S_1^b = \text{-----A}$

$$-7 - 1 = -8$$

Insertion in S^b :

$S_7^a = \text{AGCATGC}$

$S_1^b = \text{A-----}$

$$-3 - 1 = -4$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

Match / Mismatch:

$$S_1^a = -A$$

$$S_2^b = AC$$

$$-1 - 1 = -2$$

Insertion in S^a :

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_1^a = -A$$

$$S_2^b = AC$$

$$-1 - 1 = -2$$

Insertion in S^a :

$$S_1^a = A-$$

$$S_2^b = AC$$

$$2 - 1 = 1$$

Insertion in S^b :

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_1^a = -A \quad -1 - 1 = -2$$

$$S_2^b = AC$$

Insertion in S^a :

$$S_1^a = A-$$

$$S_2^b = AC \quad 2 - 1 = 1$$

Insertion in S^b :

$$S_1^a = --A \quad -2 - 1 = -3$$

$$S_2^b = AC-$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Match / Mismatch:

$$S_1^a = -A$$

$$-1 - 1 = -2$$

$$S_2^b = AC$$

Insertion in S^a :

$$S_1^a = A-$$

$$2 - 1 = 1$$

$$S_2^b = AC$$

Insertion in S^b :

$$S_1^a = --A$$

$$-2 - 1 = -3$$

$$S_2^b = AC-$$

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

Red arrows indicate the path of the Needleman-Wunsch algorithm, starting from the bottom-right cell (C, C) and moving towards the top-left cell (∅, ∅). The path is: (C, C) → (C, G) → (C, T) → (C, A) → (C, C) → (A, C) → (A, G) → (A, A) → (∅, A) → (∅, ∅).

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0 ← -1 ← -2 ← -3 ← -4 ← -5 ← -6 ← -7							
A	-1	2 ← 1 ← 0 ← -1 ← -2 ← -3 ← -4						
C	-2	1	1	3 ← 2 ← 1 ← 0 ← -1				
A	-3	0	0	2	5 ← 4 ← 3 ← 2			
A	-4	-1	-1	1	4	4 ← 3 ← 2		
T	-5							
C	-6							
C	-7							

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0 ← -1 ← -2 ← -3 ← -4 ← -5 ← -6 ← -7							
A	-1	2 ← 1 ← 0 ← -1 ← -2 ← -3 ← -4						
C	-2	1	1	3 ← 2 ← 1 ← 0 ← -1				
A	-3	0	0	2	5 ← 4 ← 3 ← 2			
A	-4	-1	-1	1	4	4 ← 3 ← 2		
T	-5	-2	-2	0	3	6 ← 5 ← 4		
C	-6							
C	-7							

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

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

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). *Journal of Molecular Biology*, 48:443-453.

	∅	A	G	C	A	T	G	C
∅	0 ← -1 ← -2 ← -3 ← -4 ← -5 ← -6 ← -7							
A	-1	2 ← 1 ← 0 ← -1 ← -2 ← -3 ← -4						
C	-2	1	1	3 ← 2 ← 1 ← 0 ← -1				
A	-3	0	0	2	5 ← 4 ← 3 ← 2			
A	-4	-1	-1	1	4	4 ← 3 ← 2		
T	-5	-2	-2	0	3	6 ← 5 ← 4		
C	-6	-3	-3	0	2	5	5 ← 7	
C	-7	-4	-4	-1	-1	4	4	7

The Needleman-Wunsch algorithm

Needleman SB and Wunsch CD (1970). Journal of Molecular Biology, 48:443-453.

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

$$S_7^a = \text{A G C A - T G C}$$

$$S_7^b = \text{A - C A A T C C}$$

$$+2-1+2+2-1+2-1+2$$

$$\text{Score} = 7$$

Variations on a theme

Scoring

Variations on a theme

Scoring

Match / Mismatch

Not all substitutions are as likely:

- Transitions / transversions for nucleic acids
- PAM, BLOSUM matrices for proteins

Variations on a theme

Scoring

Match / Mismatch

Not all substitutions are as likely:

- Transitions / transversions for nucleic acids
- PAM, BLOSUM matrices for proteins

Gap penalty

It is biologically more relevant to give a higher cost for gap opening than for gap extension. The algorithm can be generalized to account for a gap penalty, as a function of the gap length (the longer the gap, the easier it is to extend it).

Variations on a theme

Semi-global alignment

- Some issues arise when aligning a small sequence on a large one:

$$\begin{aligned} S^a &= \text{ATCCGAACATCCAATCGAAGC} \\ S^b &= \text{AGCATGCAAT} \end{aligned}$$

are aligned as

$$\begin{aligned} S^a &= \text{ATCCGAACATCCAATCGAAGC} \\ S^b &= \text{A---G--CATGCAAT-----} \end{aligned}$$

Variations on a theme

Semi-global alignment

Semi-global alignment

- Initialize first row and first column with 0 instead of negative scores
- Back-tracing: start from highest score in the last row and the last column

Variations on a theme

Semi-global alignment

Semi-global alignment

- Initialize first row and first column with 0 instead of negative scores
- Back-tracing: start from highest score in the last row and the last column

$$\begin{array}{l} S^a \\ S^b \end{array} = \begin{array}{l} \text{ATCCGAACATCCAATCGAAGC} \\ \text{-----AGCATGCAAT-----} \end{array}$$

Variations on a theme

Local alignment

- **(Semi)global alignment**: align two complete sequences
 - **Local alignment**: find the best alignment of two sub-sequences
- ☞ A local alignment starts and ends with a match

Variations on a theme

Local alignment

- **(Semi)global alignment**: align two complete sequences
 - **Local alignment**: find the best alignment of two sub-sequences
- ☞ A local alignment starts and ends with a match

Local alignment

- ☞ Smith TF and Waterman MS (1981). *Journal of Molecular Biology*, 147:195-197.
- Alignment with empty sequence costs 0 (as in semi-global)
 - No negative score! If a sub-alignment has negative score, then its score is set to 0
 - Back-tracing: we start from the highest value in the full table

Multiple alignment

☞ An alignment of at least 3 sequences!

```
KGVHSDLNET YFVGDINDVP KEGKELVETC YFSLMEAIAK CKPGMFYKNI GTLIDAYVSK KN-----F SVVRSYSGHH VGKLFHSNP
KGVHADLNET YFVGENIS-- NEAKQLVETC YFSLMEAIAK CKPGMLYKNL GNIIDAYVSK KH-----F SVIRTYSGHH VGKLFHSNP
QQMHADLNET YVVDGNIS-- KEALNTVETA RECLKIATKM CKPGVRFQDL GDAIEKHAKQ NK-----C SVVKTYCGHH VGKFFHCSP
QGYHADLNET YVVGENIS-- KEALNTTETS RECLKLAIKM CKPGTTFQEL GDHIEKHATE NK-----C SVVRTYCGHH VGEFFHCSP
LGFHADLNET YVVGDKAKCN PELVNLVETT RECLDLAIKH VKPGIAFREL GNIEKHASE NN-----C SVVRTYCGHH CGALFHCQP
NGFHDGDLNET YVVGDKAKAN PDLVCLVENT RIALDKAIAA VKPGVLFQEF GNIEKHTNS IT----EKQI SVVRTYCGHH INQLFHCSP
EGFHGDINET YVVG EKARSN PDAVRVVETA RECLDKSIEI VKPGLMFRDP GNVIEKHAKS RN-----C SVVKSYSCHGH INQLFHCAP
EGYHADLNET YYIGDKAKAD PDTVRVVETA RQCLDESIAK VKPGLTIREF GNIEKHAKQ HN-----C SVIRTYCGHH VGKLFHCPP
EGYHGDNET YYVGDKAKAD ADSVRVVETA RECLEEAIKL VKPGLTFRDF GNVIEAHAKS RG-----C SVIRTYVGHG INKTFHCPP
EGYHGDNET YYVGDKALAD PDVVRLVETT RECLDEAIKL VKPGLTFREF GNVIEKHAKA NN-----C SVIRTYVGHG INSVFHCPP
KGYHGDNET FVVGKKAEDD AESMKLRVA RECLDAAINI CGPGVPYGEI GRVIQPLAES QG-----C AVVKNYTGHG ISNCFHAAP
KGCHGDNET YFVGNDV--- EASRQLVKCT YECLEKATAI VKPGVRFREI GEIVNRHATM SG-----L SVVRSYCHGH IGDLFHCAP
RGFHGDNET FVVGNDV--- EKHKKLQVET HEALSKAIEF VRPGEKYRDI GNVIQYVAP HG-----F SVVRSYCHGH IHRVFHTAP
DSYHGDTST FVVGTPS--- PLAKRLVEVT EKCLMEAIKT VKPGSRIGDI GAAIQECAEP QG-----F SVVRDFVGHG VSKVFHTAP
DGYHGDTST FIVGNAA--- PKTKKLVEVT QECLNLGIAE VKPGAKIGDI GAAIQEYAEA QG-----F SVVRDFVGHG ISNIFHTAP
EGYHGDTST FVVGTPS--- PKTRKLVEVT EECLRLGIAE VKPGGRIGDI GAAIQEYAEQ QG-----F SVVRDFVGHG ISNIFHTAP
NGYSDASRM FIIGEAS--- ENAKRLVKVA KECLEKIEA VKPWGFLGDI GAAIQEHAEK NG-----Y SVVRDFVGHG VGLKFHEDP
KGYSDASRM FMIGDVS--- PEMRKLQVET KECMEIGIAA AQPWQQLGVDV GAAIQEHAEK NG-----F NVVRDLCHGH VGMQFHEAP
DGYHGDSST FIVGNTS--- SEVKTLVQDT EKAMFIGIEQ VRPGRNVHDI ANAIDDFLTP KG-----Y GIVRDLMGHH IGRGFHEDP
DGYFGDNSKM YIVGGETN-- IRSKKLVEAA QEALYVGIRT VKPDIRLNEI GKAVQKYTES QT-----F SVVREYCGHH VGTTEFHCEP
DGYFGDNSKM YIVGGETN-- IRSKKLVEAA QEALYVGIRT VKPDIRLNEI GKAVQKYTES QT-----F SVVREYCGHH VGTTEFHCEP
DGYHGDTSKM FLIGDVS--- IEDKRLCHVA QECLYLALKQ VKPGVQLGEI GTTIEKHIKT NNKNNPRFKF SIVRDYCGHH IGAEPHEEP
QGYHGDTSKM FLVG DVS--- PANKRLCMVA QEALYVGMRT VKPGSTVVDI GTAIEKYIKE NNKNNPRNKF SIVKDFCGHH IGDEFHEEP
QGFHGDTSKM FLVG DVS--- PANKRLCMVA QEALYIGMRQ VKPGATVVDI GTAIEKYIKD NNKNNPRNKF SIVKDFCGHH IGDEFHEEP
DGYHGDTSAM FIVGETT--- PLRRQLCKVA QESLYAAIKQ VRPGMCIETE GAVIQPIVEK AG-----F SVVRDYCGHH IGAEPHEEP
DGYHGDTSKM FVIGKTS--- ILSKRLCQVA RESLYLSLKL VKPGPILYKI GEIIQNYVES NN-----F SVVRDYCGHH IGRNFHEEP
DKYSDASKM FVVGKPT--- ELGKKLCYVA KKSLLYALYT IRPGINLQKL GKVIQNYVKK QN-----F SIVKEYCGHH IGRSFHEEP
```

Finding the best multiple alignment

- Score of a multiple alignment: sum of the scores of all pairwise alignments

Finding the best multiple alignment

- Score of a multiple alignment: sum of the scores of all pairwise alignments
- Several approaches to maximize the score:
 - Global maximization approaches
 - Agglomerative approaches
 - Iterative approaches

Finding the best multiple alignment

Global maximize approaches

Very slow as the number of sequences increases • several approximations exist • yet not usable for medium or large data sets

☞ Example methods: MSA, DCS, SAGA

Finding the best multiple alignment

Agglomerative approaches

Heuristic: use a series of pairwise alignments to build a multiple alignment • first proposed by Feng DF and Doolittle RF (1987). *Journal of Molecular Evolution*, 25:351–360.

Finding the best multiple alignment

Agglomerative approaches

Heuristic: use a series of pairwise alignments to build a multiple alignment • first proposed by Feng DF and Doolittle RF (1987). *Journal of Molecular Evolution*, 25:351–360.

- 1 Compute scores for all pairwise alignments
- 2 Build a clustering tree from the scores, to identify the closest sequences (guide tree)
- 3 Align sequences following the guide tree

Finding the best multiple alignment

Agglomerative approaches

Heuristic: use a series of pairwise alignments to build a multiple alignment • first proposed by Feng DF and Doolittle RF (1987). *Journal of Molecular Evolution*, 25:351–360.

- 1 Compute scores for all pairwise alignments
- 2 Build a clustering tree from the scores, to identify the closest sequences (guide tree)
- 3 Align sequences following the guide tree

Profile-profile alignment

Gather two alignments by inserting gap columns in each of them, without modifying the original alignments

☞ Can be done using a modification of the global alignment procedure for two sequences

Finding the best multiple alignment

Agglomerative approaches

Heuristic: use a series of pairwise alignments to build a multiple alignment • first proposed by Feng DF and Doolittle RF (1987). *Journal of Molecular Evolution*, 25:351–360.

- 1 Compute scores for all pairwise alignments
- 2 Build a clustering tree from the scores, to identify the closest sequences (guide tree)
- 3 Align sequences following the guide tree

Profile-profile alignment

Gather two alignments by inserting gap columns in each of them, without modifying the original alignments

☞ Can be done using a modification of the global alignment procedure for two sequences

☞ Example programs: ClustalW

Finding the best multiple alignment

Iterative approaches

Problems with the agglomerative approach:

- The existing sub alignments are never modified, only fused
- Highly depends on the guiding tree, any early error is reported to all further steps

Finding the best multiple alignment

Iterative approaches

Problems with the agglomerative approach:

- The existing sub alignments are never modified, only fused
- Highly depends on the guiding tree, any early error is reported to all further steps

- 1 Build a multiple alignment using an agglomerative approach
- 2 Use the resulting alignment to build a new guiding tree
- 3 Build a new multiple alignment
- 4 etc.

☞ Example software: MAFFT, MUSCLE

Other approaches to multiple alignment

Multiple Local Alignment:

- DIALIGN

Probabilistic methods:

- Joint estimation of phylogeny and alignment: StatAlign, BALi-Phy
- Use models to tune the alignment: PRANK, PROBCONS, ProAlign, MUMMALS, etc.

Meta-alignments:

- Take a consensus of several methods, example: M-Coffee

Evaluating the confidence of alignments

- ☞ All methods assume that sequences are homologous! There will always be an output, even from random input sequences!
- ☞ Apart from methodological artifacts, some regions are more difficult to align than others, and are therefore more prone to errors

Evaluating the confidence of alignments

- ☞ All methods assume that sequences are homologous! There will always be an output, even from random input sequences!
- ☞ Apart from methodological artifacts, some regions are more difficult to align than others, and are therefore more prone to errors
 - Methods based on block conservation: GBlocks
 - The Heads-or-Tails method
 - Methods assessing the uncertainty in the guide tree: GUIDANCE

Genome alignment

Genome alignment

- Complexity issue (memory, CPU)
- Caryotype
- Genome rearrangements
- Duplicates/Repeated regions
- Recombination
- etc.

Genome alignment

General strategy

- 1 **Find anchors** = highly conserved regions between two genomes. For example: shared k -mers,

Genome alignment

General strategy

- 1 **Find anchors** = highly conserved regions between two genomes. For example: shared k -mers,

Maximum Unique Match (MUM)

A MUM is a substring of a given minimum size (usually 20) that is common to two genomes, so that

- it cannot be extended on either side without introducing a mismatch,
- it is unique to both sequences.

Genome alignment

General strategy

- 1 **Find anchors** = highly conserved regions between two genomes. For example: shared k -mers,
- 2 **Define conserved regions**: sets of collinear, non-overlapping anchors, which form the basis of the alignment,

Maximum Unique Match (MUM)

A MUM is a substring of a given minimum size (usually 20) that is common to two genomes, so that

- it cannot be extended on either side without introducing a mismatch,
- it is unique to both sequences.

Genome alignment

General strategy

- 1 **Find anchors** = highly conserved regions between two genomes. For example: shared k -mers,
- 2 **Define conserved regions**: sets of collinear, non-overlapping anchors, which form the basis of the alignment,
- 3 **Close the gaps** between the anchors to build the alignment.

Maximum Unique Match (MUM)

A MUM is a substring of a given minimum size (usually 20) that is common to two genomes, so that

- it cannot be extended on either side without introducing a mismatch,
- it is unique to both sequences.