

# Introduction to Bioinformatics for Computer Scientists

Lecture 3

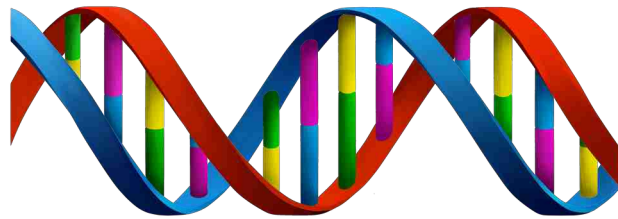
## Pair-wise Sequence Alignment

Lukas Hübner, PhD student  
lukas.huebner@h-its.org

slides by  
Benoit Morel, postdoc

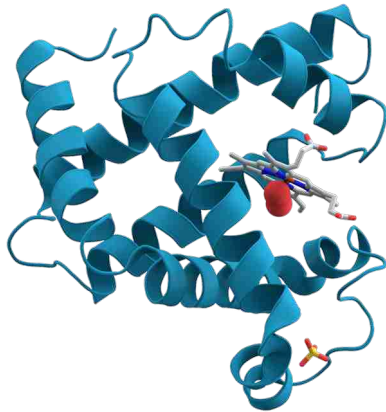
# DNA and protein sequences are strings

- DNA:



AACCTGTTGTCAAATG

- Protein:



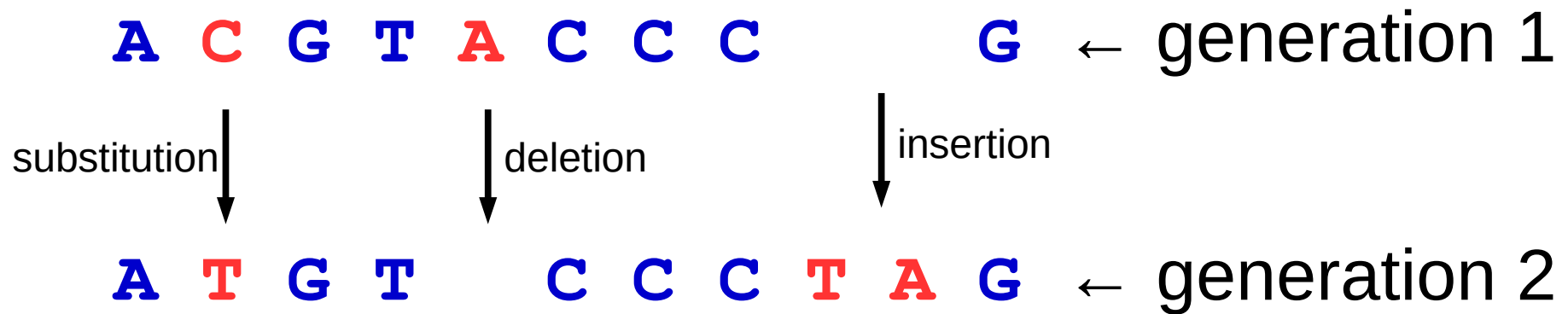
TTETTSFLIFETAVKNT

# Sequences evolve

**A C G T A C C C G**

← generation 1

# Sequences evolve



(and their lengths change)

# Pair-wise sequence alignment

Compare two sequences to infer their similarity

Example: alignment between 'GCGACGTCC'  
and 'GCGATAC'

|       | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 |
|-------|---|---|---|---|---|---|---|---|---|
| $x =$ | G | C | G | A | C | G | T | C | C |
|       |   |   |   |   |   |   |   | . |   |
| $y =$ | G | C | G | A | — | — | T | A | C |

# Example 1: Measure DNA similarity

How similar are human and chimpanzee at the DNA level?



Human DNA

Chimpanzee DNA

# Example 1: Measure DNA similarity

How similar are human and chimpanzee at the DNA level?



We first need to align their DNA sequences!



Human DNA



Chimpanzee DNA

(now nucleotids at the same position are comparable)

# Example 2: genome assembly



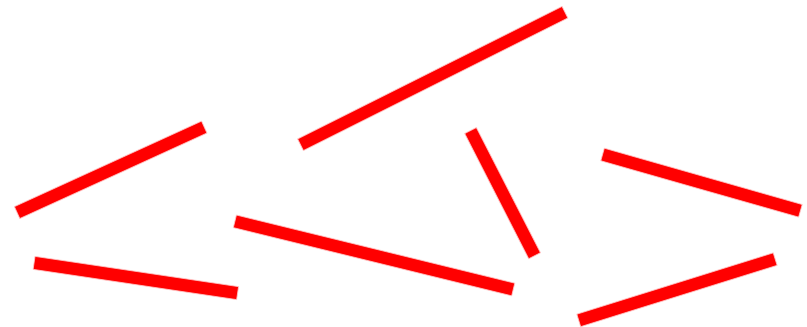
Individual to sequence



Sequencing machine



Assembled genome



Unordered reads  
(short DNA sequences)



# Example 2: genome assembly

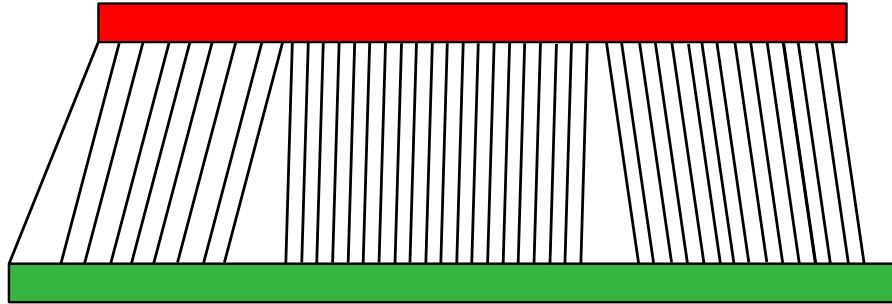
Reference genome



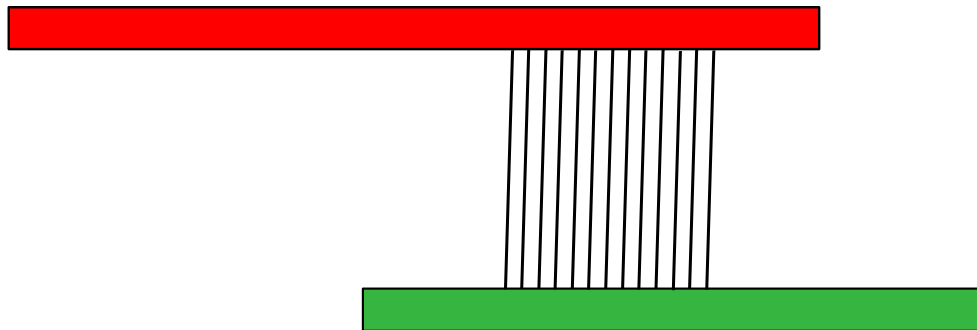
Reads to map

# Global and local alignments

Global alignment: align the full strings



Local alignment: align similar substrings



# Different approaches to alignment

Two approaches:

- **Online algorithms** ← **(today)**
- Offline algorithms

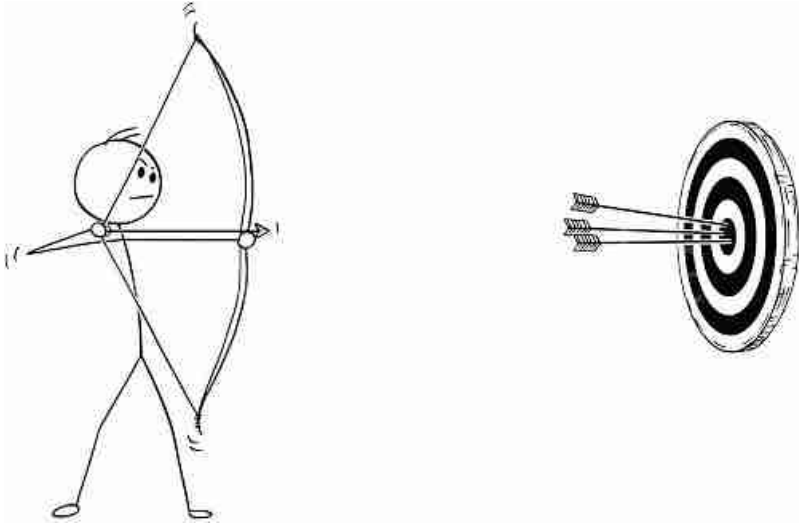
Four main approaches to online sequence alignment:

- **Dynamic programming** ← **(today)**
- Automata
- Word-level parallelism
- Filtering

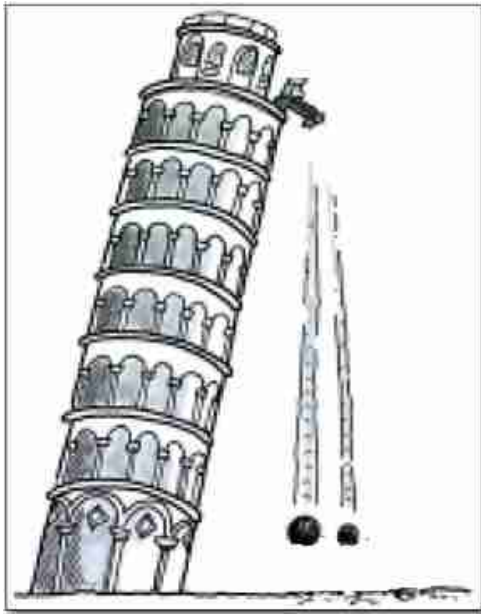
# Online VS offline

- **Online algorithm:** can process the input data piece by piece in a serial fashion
- **Offline algorithm:** needs the entire input data to start processing it

# Online VS offline example



**Online:** we don't need all arrows to start throwing them (assume they cannot collide)



**Offline:** we need both balls to start the experiment (check that they hit the ground at the same time)

# Dynamic programming

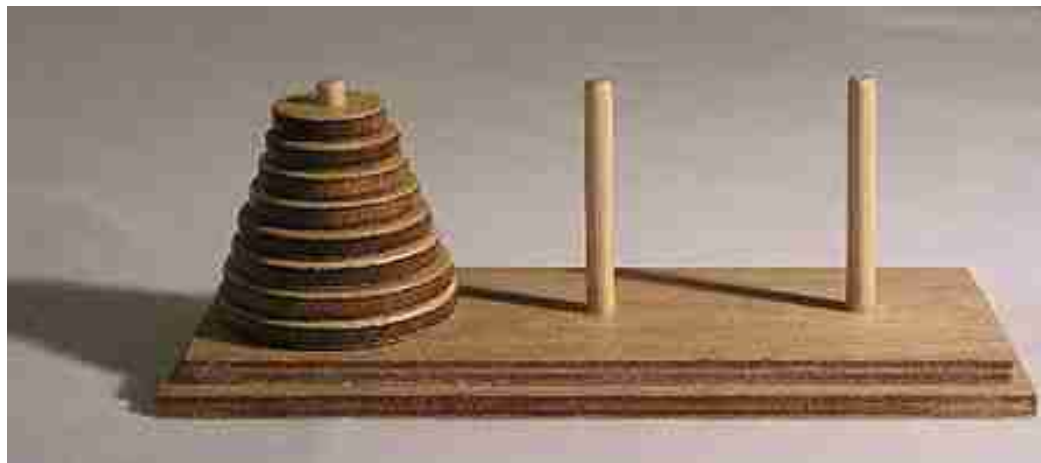
- Break down a complicated problem into simpler subproblems (typically with recursion)
- Cache the results of the recursive calls

# Dynamic programming

Break down a complicated problem into simpler subproblems (typically with recursion)

Example:

Tower of Hanoi algorithm



# Definitions



# Alphabet

An **alphabet**  $\Sigma$  is a finite non empty set whose elements are called **letters**

Examples :

- DNA:  $\Sigma = \{ 'A', 'C', 'G', 'T' \}$
- Numerics:  $\Sigma = \{ '0', '1', '2', '3', '4', '5', '6', '7', '8', '9' \}$
- German alphabet, Chinese alphabet...

# String

- A **string** is a finite sequence of elements from an alphabet
- The empty string is denoted  $\varepsilon$
- Notation:  $x[i]$  is the  $i$ th element of the string  $x$

Examples:

- “ACCAGTGGTGTGGCA” ( $x[0] = 'A'$ ,  $x[1] = 'C'$  etc.)
- “31349756984”
- $\varepsilon$

# String length

- The length of a string  $x$  is defined as the length of its sequence of letters.
- Notation:  $|x|$

Examples:

- $|\text{"pizza"}| = 5$
- $|\epsilon| = 0$

# Identity between strings

$$x = y$$

if and only if

$$|x| = |y|$$

$$x[i] = y[i] \text{ for all } i, 0 \leq i < |x|:$$

Examples:

- “42” = “42”
- “Potato”  $\neq$  “Tomato”

# String concatenation

- The concatenation of two strings  $x$  and  $y$  is the string of the letters of  $x$  followed by the letters of  $y$ .
- Notation:  $xy$

Example:

$x = \text{"James\_"}$

$y = \text{"Bond"}$

$xy = \text{"James\_Bond"}$

# Substrings

- A string  $x$  is a substring (or factor) of  $y$  if there exist  $u$  and  $v$ , such that  $y = uxv$

Examples:

- $\epsilon$  is a substring of any string
- “**logic**” is a substring of  
“**biological**” ( $u$ =“**bio**”  $v$ =“**al**”)

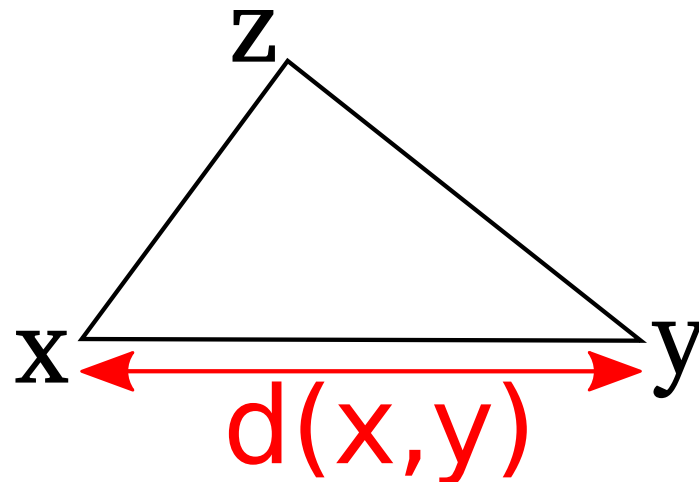
# Occurrence of a string

X is an **occurrence** of y (X **occurs** in Y) if X is a substring of Y

“GA” occurs 3 times in “AGATCTGACGA”

# Distance

- A function  $d$  is a distance if the following conditions are satisfied for all element  $x$  and  $y$ :
  - Positivity:  $d(x,y) \geq 0$
  - Separation:  $d(x,y) = 0 \iff x = y$
  - Symmetry:  $d(x,y) = d(y,x)$
  - Triangle inequality:  $d(x,y) \leq d(x,z) + d(z,y)$  for all elements  $z$





# Hamming distance

- Defined for strings of same length
- Counts the number of characters that differ
- Linear complexity

**P**ING

**H**AMMING

AC**G**T**T**GG**G**TT

**P**ONG

**L**EMMING

AC**G**AT**G**C**A**TT

$d = 1$

$d = 2$

$d = 3$

# Is hamming biologically relevant?

**x** = ACTATATATACG

**y** = CTATATATACGT

$d = 12$ , the distance between **x** and **y** is maximal

# Is hamming biologically relevant?

**x** = ACTATATATACG

**y** = CTATATATACGT

$d = 12$ , the distance between **x** and **y** is maximal

... but **x** and **y** are very similar!

**x** = ACTATATATACG-

**y** = -CTATATATACGT

(Hamming distance is actually relevant, but not to compare “raw” sequences)

# Edit operations: substitution, insertion and deletion

**Substitution** of one letter  $x$  by another letter  $y$

A C G T G C  
↓  
A C G A G C

# Edit operations: substitution, insertion and deletion

Insertion of one letter

A C G - G C  
↓  
A C G A G C

# Edit operations: substitution, insertion and deletion

**Deletion** of one letter

A C G T G C  
↓  
A C G - G C

# Edit distance

$\delta e(x,y)$  = minimum number of (edits) operations to transform  $x$  into  $y$

# Edit distance

Example: compute the edit distance between  
“SALADS” and “BALLAD”



# Edit distance

Example: compute the edit distance between  
“SALADS” and “BALLAD”

|                  |            |
|------------------|------------|
| SALADS → BALADS  | Subs S → B |
| BALADS → BALLADS | Insert L   |
| BALLADS → BALLAD | Del S      |

Edit distance = 3

# Alignment definition

Result of inserting gaps in both strings such that they have the same length

Alignments between

x='ACGA' and y='ATGCTA':

**ACGA--**

**A-----CGA**

**A--CGA**

**ATGCTA**

**-ATGCTA--**

**ATGCTA**

# What is a good alignment?

We can (for instance) score an alignment with the hamming distance.

ACGA--

A-----CGA

A--CGA

ATGCTA

-ATGCTA---

ATGCTA

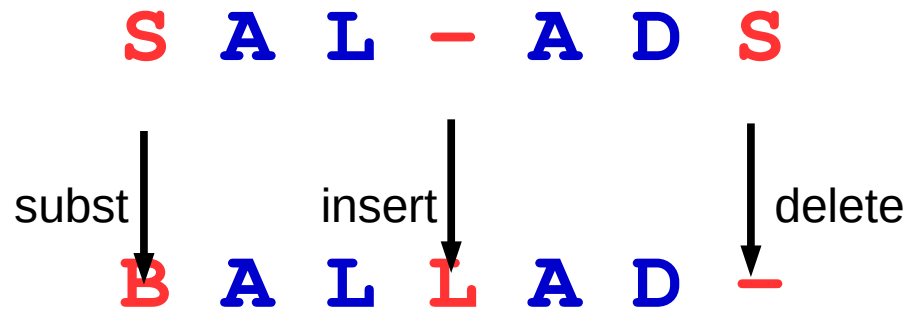
d=4

d=10

d=3

# Aligning sequences using edit distance

Computing the alignment with minimum hamming distance = Computing the edit distance and storing the sequence of edit operations



hamming(SAL-ADS,BALLAD-) = edit(SALADS, BALLAD)

How to compute the edit distance between  
two strings **X** and **Y**?

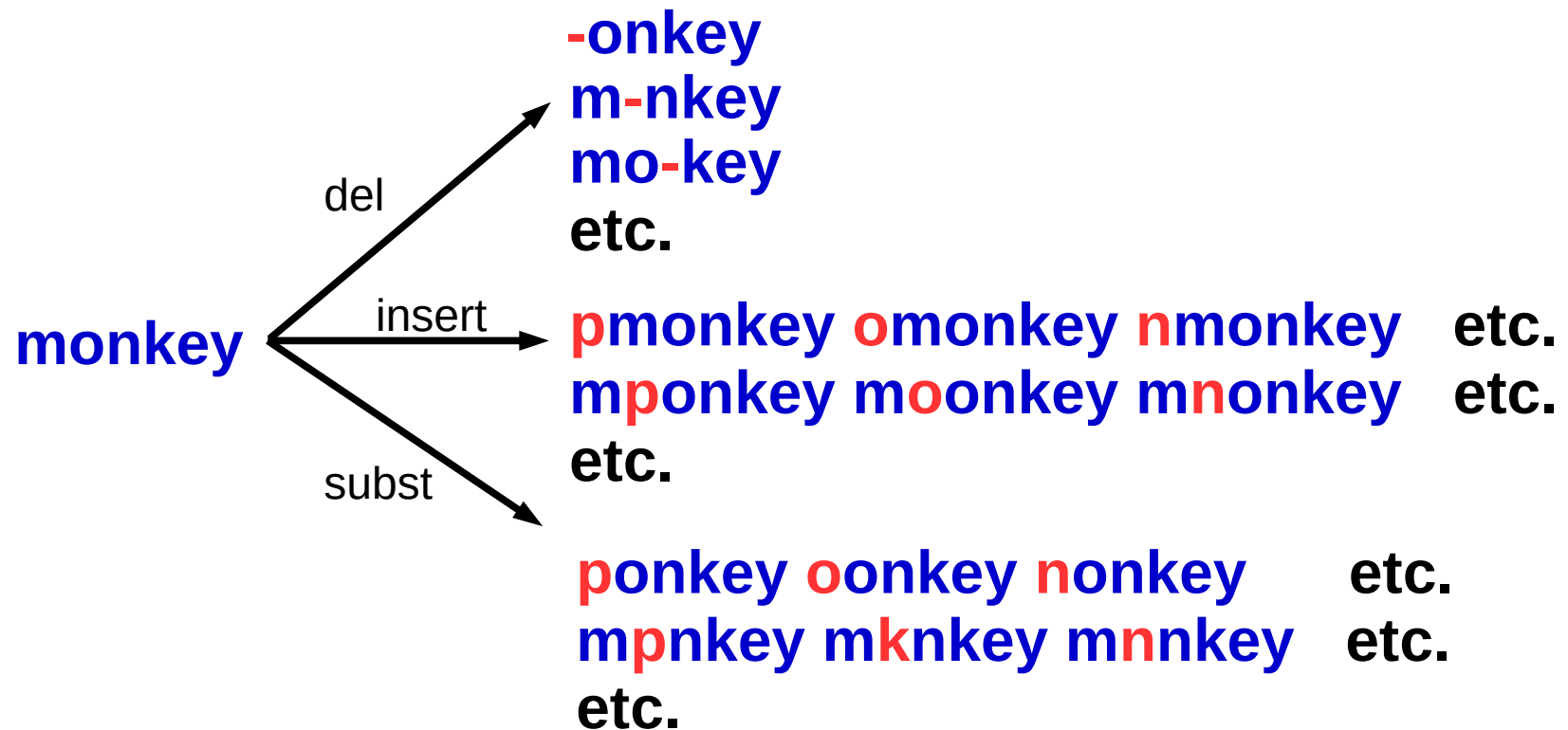
For instance, how to compute the distance between **poney** and **monkey**?

# Naive bruteforce

For each distance from 1 to  $|X|$ , try every possible combinations of edit operations

# Naive bruteforce

## First iteration:





# Naive bruteforce

**Second iteration:** repeat the same procedure from each output of the previous iteration

→ Non-polynomial time complexity (with respect to  $|X|$  and  $|Y|$ )

# Dynamic programming



Recursion on the prefixes of  $X$  and  $Y$

We need:

- A recursion formula
- Trivial edge case to end the recursion

# Edge case

$\delta e(\mathbf{X}, \varepsilon) = |\mathbf{X}|$  (delete all letters of  $\mathbf{X}$ )

$\delta e(\varepsilon, \mathbf{Y}) = |\mathbf{Y}|$  (insert all letters of  $\mathbf{Y}$ )

# Dynamic programming



Recursion on the prefixes of  $x$  and  $y$

We need:

- A recursion formula
- Trivial edge case to end the recursion

# Levenstein formula

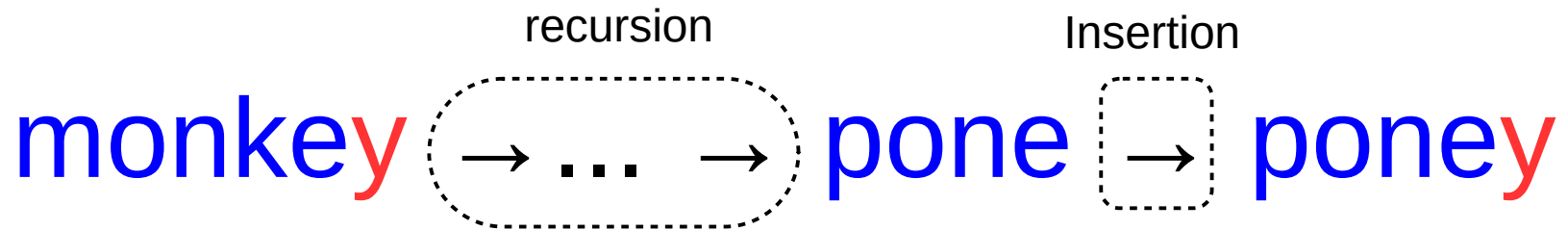
Let  $F(i,j)$  be the edit distance between the prefixes of  $X$  of size  $i$  and the prefix of  $Y$  of size  $j$ .

$$F(i,j) = \min ( \begin{array}{l} F(i, j - 1) + 1, \\ F(i - 1, j) + 1, \\ F(i - 1, j - 1) + 1 * (X[i] != Y[j]) \end{array} )$$

# Levenstein formula

$$\delta e(\text{monkey}, \text{poney}) = \min ($$
$$\delta e(\text{monkey}, \text{pone}) + 1,$$
$$\delta e(\text{monke}, \text{poney}) + 1,$$
$$\delta e(\text{monke}, \text{pone}) + 1 * (y \neq y)$$
$$)$$

# First term



$$\delta e(\text{monkey}, \text{poney}) + 1$$

$$F(i, j - 1) + 1$$

# Second term

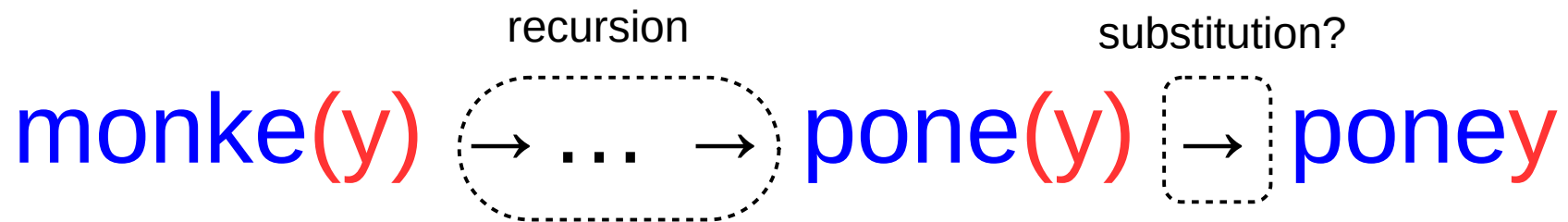


$$\delta e(\text{monke}, \text{poney}) + 1$$

$$F(i - 1, j) + 1$$



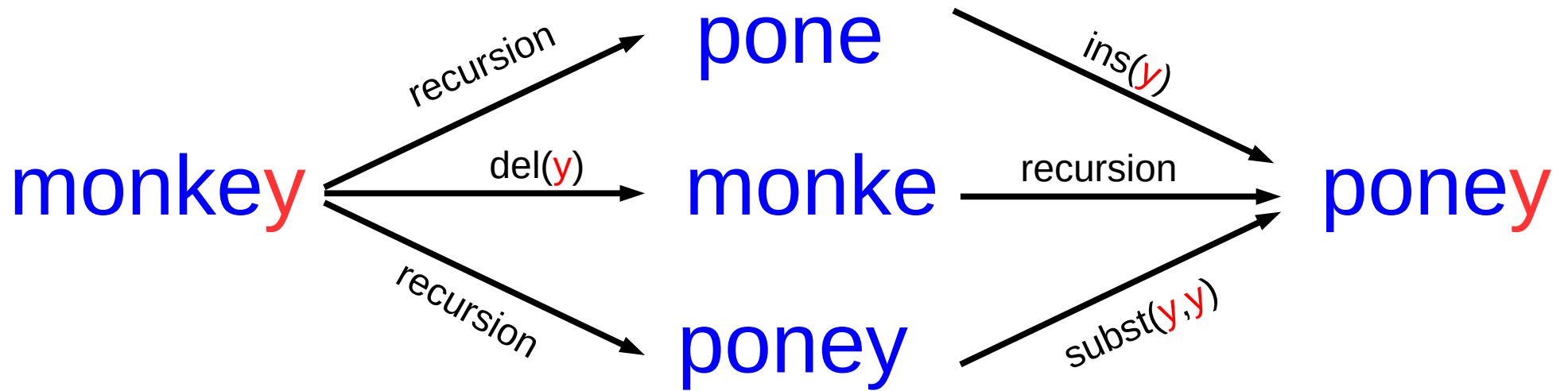
# Third term



$$\delta e(\text{monke}, \text{pone}) + 0$$

$$F(i - 1, j - 1) + 1 * (X[i] \neq Y[j])$$

# Take the best of the three paths



# When to stop?

Stop the recursion when one prefix is empty:

$$F(0, j) = j$$

$$F(i, 0) = i$$

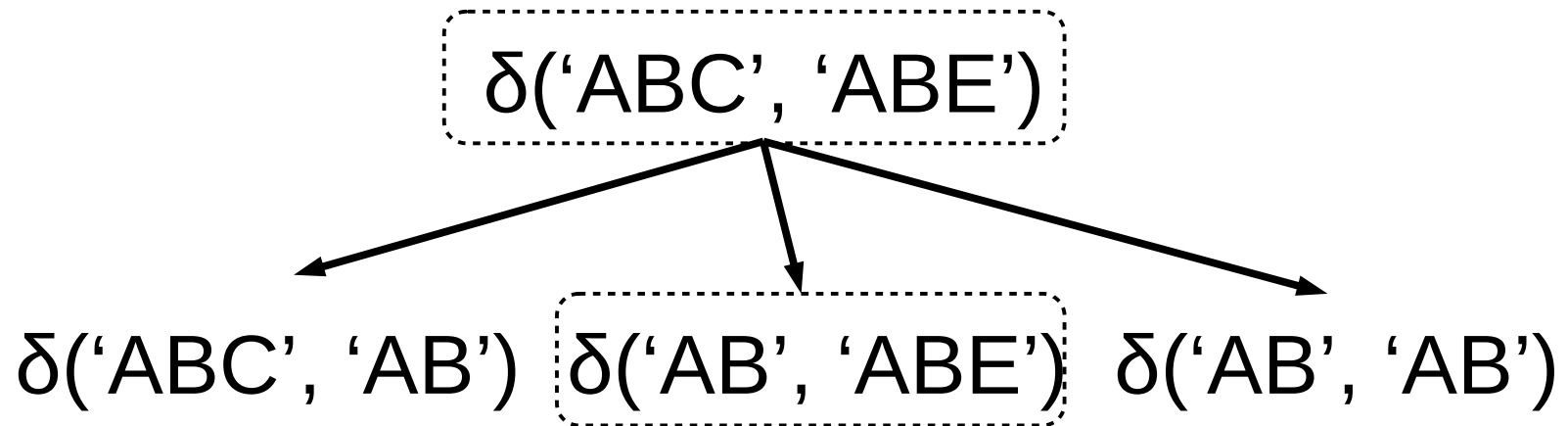
remember:  $\delta e(\mathbf{X}, \boldsymbol{\varepsilon}) = |\mathbf{X}| = \delta e(\boldsymbol{\varepsilon}, \mathbf{X})$

# Levenstein formula recursion convergence

$$F(i,j) = \min ( F(i, j - 1) + 1, \\ F(i - 1, j) + 1, \\ F(i - 1, j - 1) + 1 * (X[i] != Y[j]) )$$

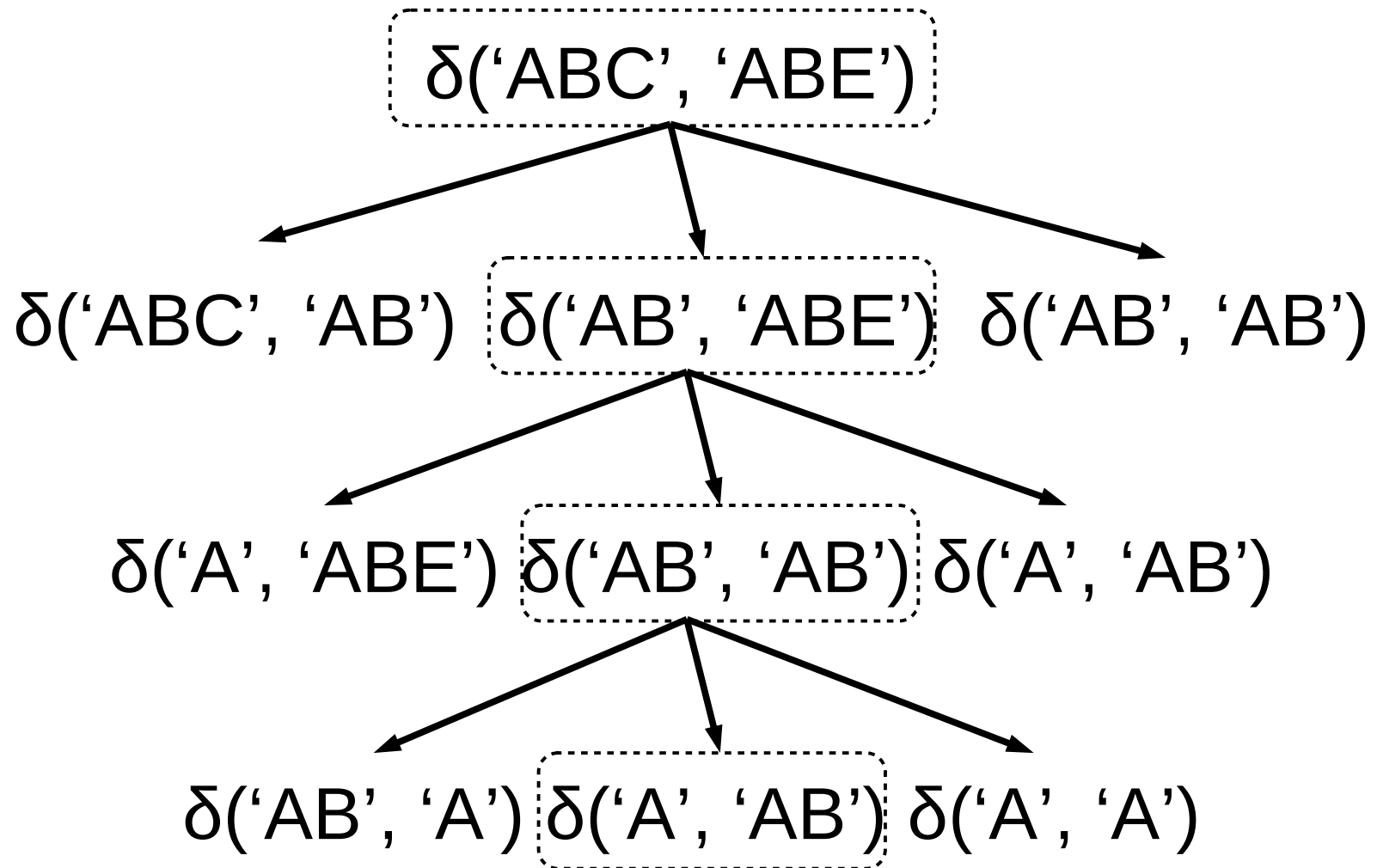
The quantity (i+j) decreases at each step, until i=0 or j=0, which ends the recursion

Let's have a look at  
the functions calls

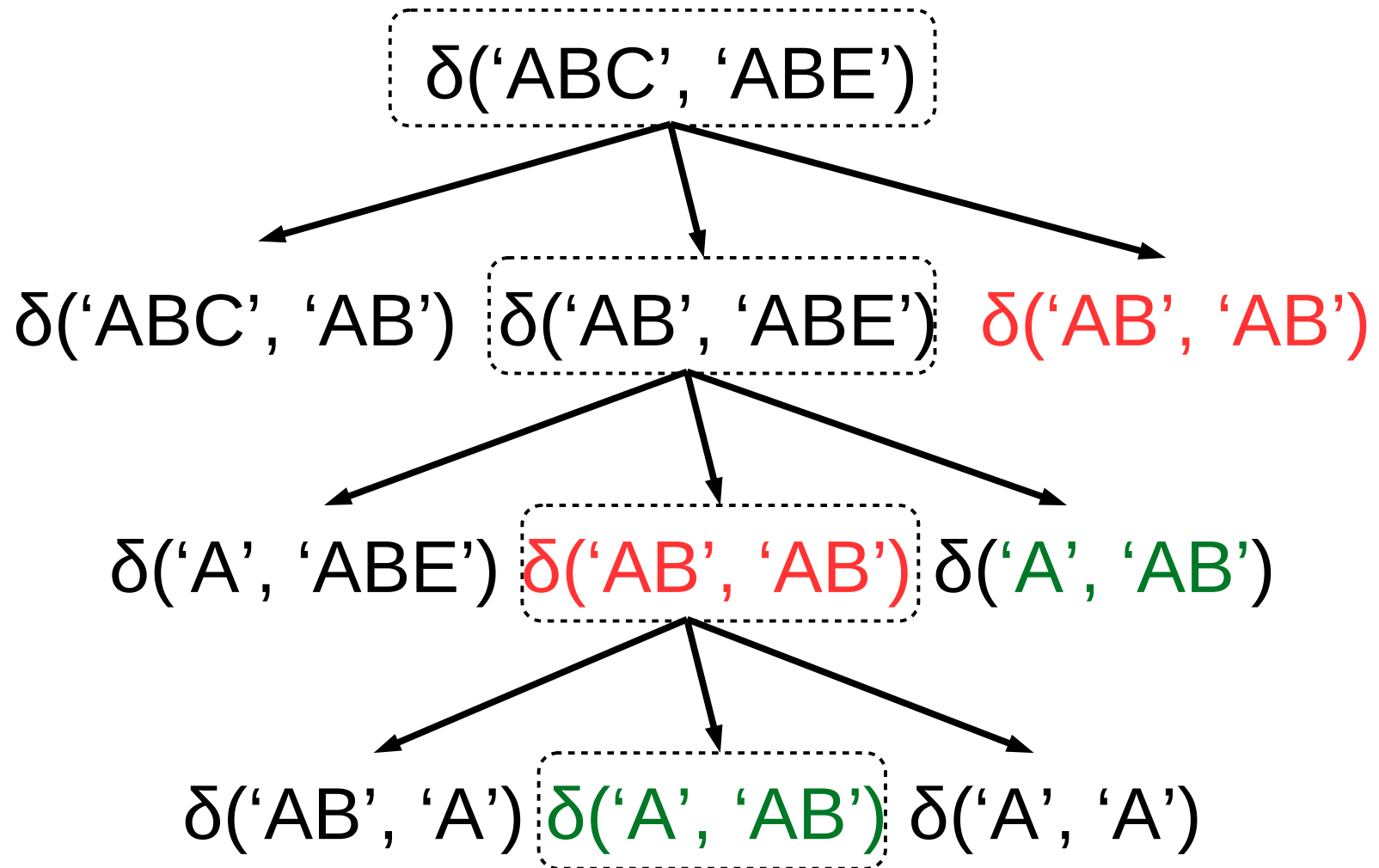


# Naive implementation

→ exponential complexity



# Redundant computations !



# Dynamic programming

Be careful not to blindly implement the recursion!

Store intermediate results in a table to avoid redundant computations



# Needleman-Wunsch algorithm


- Computes the edit distance
- Finds the best alignment

Stores intermediate results in a table to save computations

# Needleman-Wunsch algorithm

Store in a table T the edit distance between all the possible prefixes.

|   |  | M | O | N | K | E | Y |
|---|--|---|---|---|---|---|---|
|   |  |   |   |   |   |   |   |
| M |  |   |   |   |   |   |   |
| O |  |   |   |   |   |   |   |
| N |  |   |   |   |   |   |   |
| E |  |   |   |   |   |   |   |
| Y |  |   |   |   |   |   |   |



$T[4][4]$  = Edit distance between **MONE** and **MONK**

# Needleman-Wunsch algorithm

Edit distance between  $x=\text{money}$  and  $y=\text{monkey}$  with dynamic programming

Let's compute T:

|   |  | M | O | N | K | E | Y |
|---|--|---|---|---|---|---|---|
| M |  |   |   |   |   |   |   |
| O |  |   |   |   |   |   |   |
| N |  |   |   |   |   |   |   |
| E |  |   |   |   |   |   |   |
| Y |  |   |   |   |   |   |   |

# Needleman-Wunsch algorithm

Initialization of first row and first column: trivial.

Rationale:  $\delta e(X, \epsilon) = |X| = \delta e(\epsilon, X)$

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 |   |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

$T[0][3] = \text{Edit distance between } \epsilon \text{ and } \text{MON} = 3$

# Needleman-Wunsch algorithm

Fill the rest of the table :

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \\ T[i][j-1] + 1, \\ T[i-1][j-1] + \text{subst}(i,j) \end{array} \begin{array}{l} 1 + 1 = 2 \\ 1 + 1 = 2 \\ 0 + 0 = 0 \end{array} )$$

|   |   | M  | O | N | K | E | Y |
|---|---|----|---|---|---|---|---|
|   | 0 | 1  | 2 | 3 | 4 | 5 | 6 |
| M | 1 | ?? |   |   |   |   |   |
| O | 2 |    |   |   |   |   |   |
| N | 3 |    |   |   |   |   |   |
| E | 4 |    |   |   |   |   |   |
| Y | 5 |    |   |   |   |   |   |

# Needleman-Wunsch algorithm

Fill the rest of the table :

$$T[i][j] = \min(\begin{array}{l} T[i-1][j] + 1, \\ T[i][j-1] + 1, \\ T[i-1][j-1] + \text{subst}(i,j) \end{array} \quad \begin{array}{l} 1 + 1 = 2 \\ 1 + 1 = 2 \\ 0 + 0 = 0 \end{array})$$

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

# Needleman-Wunsch algorithm

After filling a cell, keep track of the best path

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

Diagonal path: substitution  $M \rightarrow M$

# Needleman-Wunsch algorithm

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \\ T[i][j-1] + 1, \\ T[i-1][j-1] + \text{subst}(i,j) \end{array} \quad \begin{array}{l} 0 + 1 = 1 \\ 2 + 1 = 3 \\ 1 + 1 = 2 \end{array} )$$

|   |   | M | O  | N | K | E | Y |
|---|---|---|----|---|---|---|---|
|   | 0 | 1 | 2  | 3 | 4 | 5 | 6 |
| M | 1 | 0 | ?? |   |   |   |   |
| O | 2 |   |    |   |   |   |   |
| N | 3 |   |    |   |   |   |   |
| E | 4 |   |    |   |   |   |   |
| Y | 5 |   |    |   |   |   |   |



# Needleman-Wunsch algorithm

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \\ T[i][j-1] + 1, \\ T[i-1][j-1] + \text{subst}(i,j) \end{array} \quad \begin{array}{l} 0 + 1 = 1 \\ 2 + 1 = 3 \\ 1 + 1 = 2 \end{array} )$$

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

# Needleman-Wunsch algorithm

After filling a cell, keep track of the best path

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

Horizontal path: insertion of 'O'

# Needleman-Wunsch algorithm

Iterate...

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

# Needleman-Wunsch algorithm

Iterate...

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

# Needleman-Wunsch algorithm

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \quad 2 + 1 = 3 \\ T[i][j-1] + 1, \quad 0 + 1 = 1 \\ T[i-1][j-1] + \text{subst}(i,j) \quad 1 + 1 = 2 \end{array} )$$

|   |   | M  | O | N | K | E | Y |
|---|---|----|---|---|---|---|---|
|   | 0 | 1  | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0  | 1 | 2 | 3 | 4 | 5 |
| O | 2 | ?? |   |   |   |   |   |
| N | 3 |    |   |   |   |   |   |
| E | 4 |    |   |   |   |   |   |
| Y | 5 |    |   |   |   |   |   |

Diagram illustrating the Needleman-Wunsch algorithm's dynamic programming table. The table shows the edit distance between prefixes of the strings "MONKE" (rows) and "MONKY" (columns). The cell at row 'O', column 'M' (index [2][1]) is highlighted in cyan and contains '??', indicating the current state of the algorithm. Arrows indicate the backpointers for the cells: a black arrow from (1,1) to (0,0), a blue arrow from (1,1) to (1,0), a red arrow from (1,1) to (0,1), and a black arrow from (1,0) to (1,1).

# Needleman-Wunsch algorithm

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \quad 2 + 1 = 3 \\ T[i][j-1] + 1, \quad 0 + 1 = 1 \\ T[i-1][j-1] + \text{subst}(i,j) \quad 1 + 1 = 2 \end{array} )$$

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

Diagram illustrating the Needleman-Wunsch algorithm's dynamic programming table. The table shows the edit distance between prefixes of the strings "MONKE" (rows) and "MONKY" (columns). The cell (2,2) containing '1' is highlighted in cyan. Arrows indicate the backpointers: a black arrow from (1,2) to (0,1), a red arrow from (2,2) to (1,1), a blue arrow from (2,2) to (1,2), and a yellow arrow from (2,2) to (2,1).

# Needleman-Wunsch algorithm

After filling a cell, keep track of the best path

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

Vertical path: deletion of 'O'

# Needleman-Wunsch algorithm

Iterate...

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |



# Needleman-Wunsch algorithm

Iterate...

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 |   |   |
| Y | 5 |   |   |   |   |   |   |

The diagram illustrates the Needleman-Wunsch algorithm for sequence alignment. It shows a dynamic programming table with rows labeled M, O, N, E, Y and columns labeled 0, M, O, N, K, E, Y. The table contains numerical values representing the alignment score at each position. Arrows indicate the path of maximum alignment, starting from the bottom-right cell (E, K) and moving back to the top-left cell (0, 0). The cell (E, K) is highlighted in cyan.

# Needleman-Wunsch algorithm

Iterate...

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 |   |   |   |

The diagram illustrates the Needleman-Wunsch algorithm's iteration step. It shows a dynamic programming table with rows labeled M, O, N, E, Y and columns labeled 0, M, O, N, K, E, Y. The table contains numerical values representing the maximum score for each subproblem. Arrows indicate the path of maximum score, starting from the top-left cell (0,0) and moving towards the bottom-right cell (5,5). The path is: (0,0) → (1,1) → (2,2) → (3,3) → (4,4) → (5,5). The cells (4,5) and (5,6) are shaded grey, indicating they are not part of the current iteration.

# Needleman-Wunsch algorithm

$$T[i][j] = \min( \begin{array}{l} T[i-1][j] + 1, \\ T[i][j-1] + 1, \\ T[i-1][j-1] + \text{subst}(i,j) \end{array} \begin{array}{l} 2 + 1 = 3 \\ 1 + 1 = 2 \\ 1 + 1 = 2 \end{array} )$$

|   |   | M | O | N | K  | E | Y |
|---|---|---|---|---|----|---|---|
|   | 0 | 1 | 2 | 3 | 4  | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3  | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2  | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1  | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1  | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | ?? |   |   |

# Needleman-Wunsch algorithm

Sometimes, there are several best paths

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 |   |   |

Ambiguity!

# Needleman-Wunsch algorithm

Edit distance between monkey and money = 1

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

# Backtrace the best alignment

Follow the lines!

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

The table shows the backtrace path for the best alignment. The path starts at the bottom-right cell (Y, Y) with a value of 1, which is circled with a dashed line. The path follows the arrows: (5, 7) to (4, 7), (4, 7) to (3, 6), (3, 6) to (2, 5), (2, 5) to (1, 4), (1, 4) to (0, 3), (0, 3) to (0, 2), (0, 2) to (0, 1), and finally (0, 1) to (0, 0).

# Backtrace the best alignment

- Vertical line: deletion
- Horizontal line: insertion
- Diagonal line: substitution

M O N ~~E~~ Y  
M O N K E Y

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

# Recap: Needleman-Wunsch algorithm

- Initialize first row and first column

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 |   |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |



# Recap: Needleman-Wunsch algorithm

- Initialize first row and first column
- Fill each element with the minimum of the three previous values. Store the best path.

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 |   |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | ? |   |   |
| Y | 5 |   |   |   |   |   |   |

# Recap: Needleman-Wunsch algorithm

- Initialize first row and first column
- Fill each element with the minimum of the three previous values. Store the best path.
- Backtrace to get the chain of operations

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 |   |   |   |   |   |   |
| O | 2 |   |   |   |   |   |   |
| N | 3 |   |   |   |   |   |   |
| E | 4 |   |   |   |   |   |   |
| Y | 5 |   |   |   |   |   |   |

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | ? |   |   |
| Y | 5 |   |   |   |   |   |   |

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |



Let  $n=|X|$  and  $m=|Y|$ . What is the complexity of the Needleman-Wunsch algorithm?

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

# Complexity

- Most expensive step: filling the table
- $O(n * m)$  where  $n = |X|$  and  $m = |Y|$

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

# Complexity

- Most expensive step: filling the table
- $O(n * m)$  where  $n = |X|$  and  $m = |Y|$
- Much quicker than naive recursion or bruteforce!

|   |   | M | O | N | K | E | Y |
|---|---|---|---|---|---|---|---|
|   | 0 | 1 | 2 | 3 | 4 | 5 | 6 |
| M | 1 | 0 | 1 | 2 | 3 | 4 | 5 |
| O | 2 | 1 | 0 | 1 | 2 | 3 | 4 |
| N | 3 | 2 | 1 | 0 | 1 | 2 | 3 |
| E | 4 | 3 | 2 | 1 | 1 | 1 | 2 |
| Y | 5 | 4 | 3 | 2 | 2 | 2 | 1 |

The table shows the dynamic programming values for the word 'MONKEY'. The rows are labeled M, O, N, E, Y and the columns are labeled 0, 1, 2, 3, 4, 5, 6. The values in the cells represent the number of ways to form the prefix of 'MONKEY' up to that point. Arrows indicate the path from the bottom-right cell (5,5) to the top-left cell (0,0). The value 1 in the bottom-right cell is circled.

# Ambiguities

|            | $\epsilon$ | A | T | G | C | T | A |
|------------|------------|---|---|---|---|---|---|
| $\epsilon$ | 0          | 1 | 2 | 3 | 4 | 5 | 6 |
| A          | 1          | 0 | 1 | 2 | 3 | 4 | 5 |
| C          | 2          | 1 | 1 | 2 | 2 | 3 | 4 |
| G          | 3          | 2 | 2 | 1 | 2 | 3 | 4 |
| A          | 4          | 3 | 3 | 2 | 2 | 3 | 3 |

$\begin{pmatrix} A--CGA \\ ATGCTA \end{pmatrix}$ 
 $\begin{pmatrix} ACG--A \\ ATGCTA \end{pmatrix}$

Should we really assign the same cost to substitutions, deletions and insertions?

Should we really assign the same cost to substitutions, deletions and insertions?

→ No, there is not reason to think that these events are equally likely.



# Adding weights

- We can penalize some less likely operations with weights.

Example: assume that substitutions happen 5 times more often than additions and deletions:

$$F(i,j) = \min ($$
$$F(i, j - 1) + 5,$$
$$F(i - 1, j) + 5,$$
$$F(i - 1, j - 1) + 1 * (X[i] \neq Y[j])$$
$$)$$

# Remark

If the insertion and deletion costs are different, the edit distance is not a distance anymore:

The symmetry  $d(X, Y) = d(Y, X)$

is not respected anymore



In coding genes, what could make an insertion less likely than a substitution?

(coding genes code for a protein)

# Insertion VS substitution in coding genes

Insertions are less likely to produce a fit organism

- Insertion of a non-multiple of three number of nucleotids → shift the whole sequence!

TTT CCC AAA GGG → TAT TCC CAA AGG

- Some nucleotide substitutions are silent. For instance, TTT and TTC both code for the same amino acid Lys

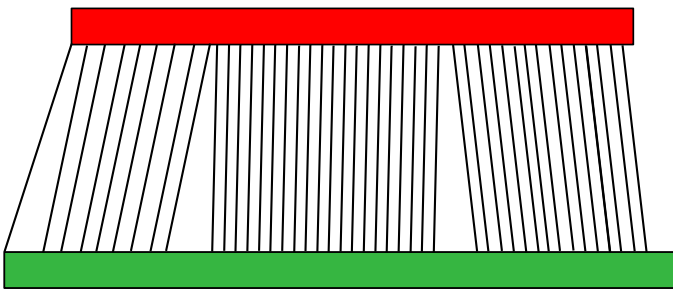
# What is the ratio between substitutions and insertions/deletions?

It depends!

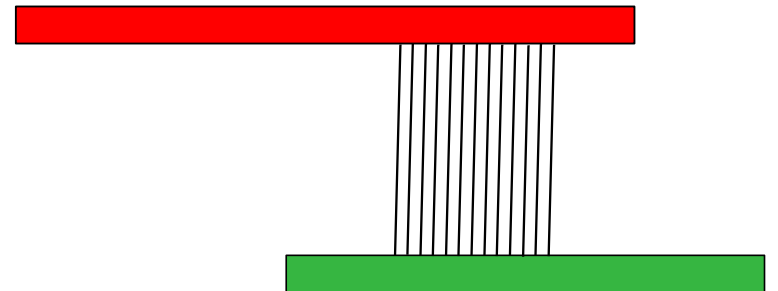
- Type of organism
- Coding/non coding genes
- Etc.



What kind of alignment problem did we solve with the Needleman-Wunsch algorithm?

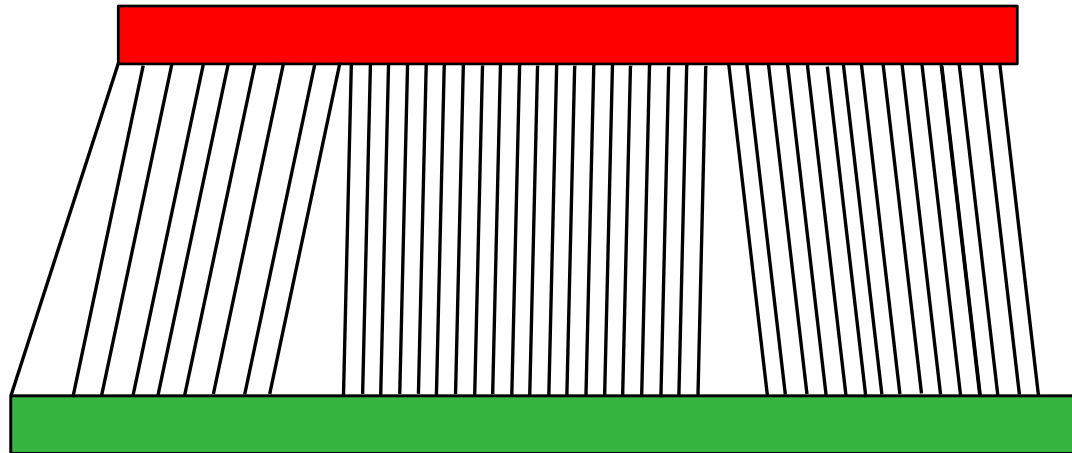


vs



# Answer: global alignment

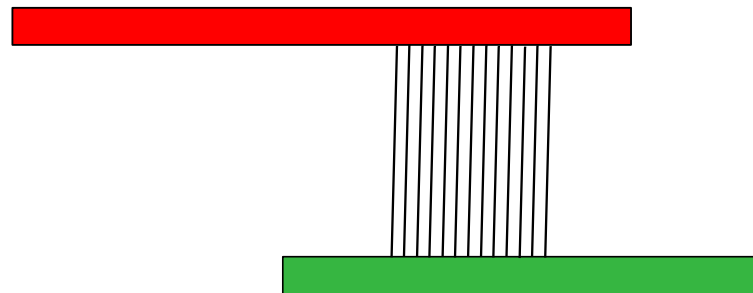
- Needleman-Wunsch gives us the best sequence of operations to get Y from X.
- This corresponds to the best global alignment.



# Local alignment

Input: two strings X and Y

Output: two aligned substrings





# Local alignment problem

Example: TTTACCACT and  
GACCAACGGG

|   |   |   |          |          |          |          |          |          |          |          |          |   |   |   |   |  |  |
|---|---|---|----------|----------|----------|----------|----------|----------|----------|----------|----------|---|---|---|---|--|--|
| T | T | T | <b>A</b> | <b>C</b> | <b>C</b> | <b>A</b> | -        | <b>C</b> | <b>A</b> | <b>A</b> | <b>C</b> | T |   |   |   |  |  |
|   |   | G | <b>A</b> | <b>C</b> | <b>C</b> | <b>A</b> | <b>T</b> | <b>C</b> | <b>A</b> | <b>A</b> | <b>C</b> | G | G | G | G |  |  |

# Formulate local alignment problem with distances

*“Find the substrings with minimum distance”*

Issue: distances increase with the lengths of the strings. Short strings will be selected

$$\delta_e(\text{ACCACAAC}, \text{ACCA}\color{red}{T}\text{CAAC}) = 1$$

$$\delta_e(\text{ACCA}, \text{ACCA}) = 0$$

# Similarity functions

- Assign positive and negative weights, to favor similarities
  - Insertion, deletion, non-identity substitution have a negative weight
  - Identity substitution has a positive weight

Similarity functions are NOT distances. They do not verify:

$$\text{Positivity: } s(x,y) \geq 0$$

# Formulate local alignment problem with similarities

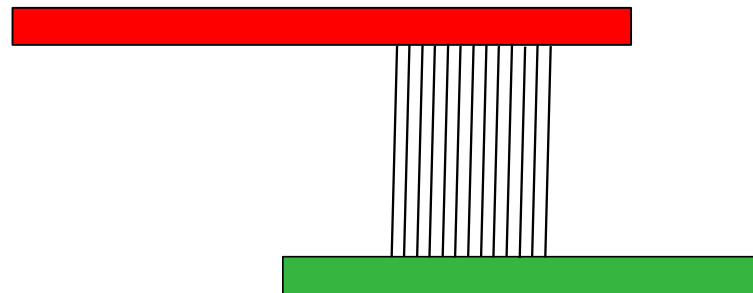
*“Find the substrings with maximum similarity”*

$$S(\text{ACCACAAC}, \text{ACCA}\color{red}{\text{T}}\text{CAAC}) = 8 - 1 = 7$$

$$S(\text{ACCA}, \text{ACCA}) = 4$$

# Local alignment problem

Find the two substrings of  $X$  and  $Y$  with the maximal similarity.



# Smith-Waterman algorithm to solve local alignment

Same as Needleman-Wunsch algorithm, but:

- Use a similarity function for the Levenstein formula
- Negative values are set to 0.
- Initialize first row and column with 0.
- Find the largest value in the table and traceback until a 0 is reached.

# Align EAWACQGKL and ERDAWCQPGKWY

| S | - | E | R | D | A | W | C | Q | P | G | K | W | Y |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| E | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| W | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| A | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 |
| Q | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 3 | 2 | 1 | 0 | 0 | 0 |
| G | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 | 1 | 3 | 2 | 1 | 0 |
| K | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 2 | 4 | 3 | 2 |
| L | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 3 | 2 | 1 |

(AWACQ-GK)  
(AW-CQPGK)

# Gap penalties

One large deletion is more likely than many small deletions:

GAAAAAT

GAA---T

GAAAAAT

G-A-A-T

→ different scores for gap-start and gap-extension



# Substitution matrices

- Are all substitution equally likely?
- Is a nucleotid **A** more likely to be replaced with a **G** or a **T**?

# Substitution matrices

- Are all substitution equally likely?

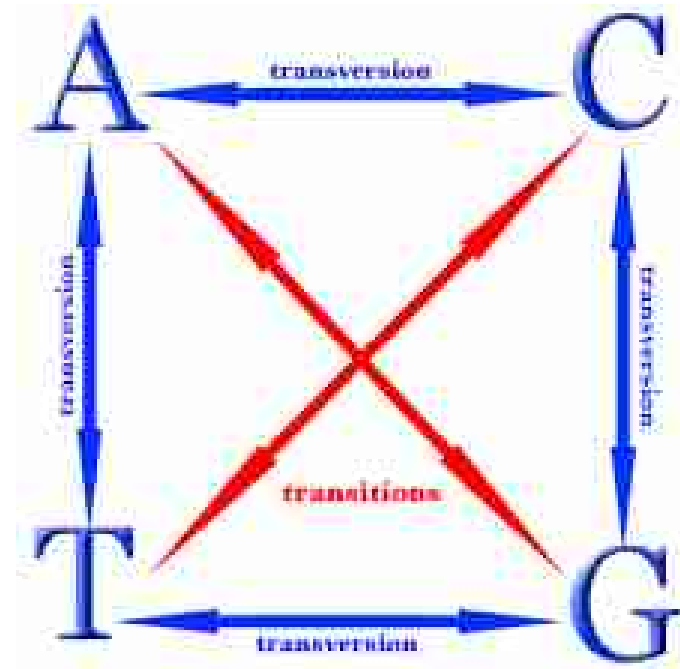
No!

- Is a nucleotid **A** more likely to be replaced with a **G** or a **T**?

Yes!

# Substitution matrices

For instance, transitions (A ↔ G and C ↔ T) happen more often than transversions



# Substitution matrices

- The substitutions costs in our alignment algorithms should take these rates into account:

$$\text{subst}(A,G) < \text{subst}(A,C)$$

- $F(i,j) = \min ($   
     $F(i, j - 1) + \text{ins}(\mathbf{b}),$   
     $F(i - 1, j) + \text{del}(\mathbf{a}),$   
     $F(i - 1, j - 1) + \text{subst}(\mathbf{a},\mathbf{b}))$

# Substitution matrices

A substitution matrix describes the rate at which one character in a sequence changes to other character states over the time.

|   | A  | G  | C  | T  |
|---|----|----|----|----|
| A | 10 | -1 | -3 | -4 |
| G | -1 | 7  | -5 | -3 |
| C | -3 | -5 | 9  | 0  |
| T | -4 | -3 | 0  | 8  |

← Example of a DNA substitution matrix

# Substitution matrices

Substitution matrices are crucial to build reliable alignments!

|   | A  | G  | C  | T  |
|---|----|----|----|----|
| A | 10 | -1 | -3 | -4 |
| G | -1 | 7  | -5 | -3 |
| C | -3 | -5 | 9  | 0  |
| T | -4 | -3 | 0  | 8  |

# BLOSUM matrices

- BLOSUM: BLOcks of Amino Acid SUBstitution Matrix
- The BLOSUM matrices are reference substitution matrices used for protein sequence alignment
- Computed from empirical local alignments

# BLOSUM matrices

- There are several BLOSUM matrices named with numbers.
- High numbers (**BLOSUM80**) → close sequences
- Low numbers (**BLOSUM45**) → distant sequences



# Log-odds scores

Observed frequency of substitution

$$S_{i,j} = \left[ \frac{1}{\lambda} \log \left( \frac{p_{i,j}}{\pi_i \cdot \pi_j} \right) \right]$$

Expected frequency of substitution

scaling factor to obtain easily computable integer values.

( $\pi_i$ : frequency of character  $i$  in the sequences)

# Log-odds scores

Observed frequency of substitution

$$S_{i,j} = \left[ \frac{1}{\lambda} \log \left( \frac{p_{i,j}}{\pi_i \cdot \pi_j} \right) \right]$$

Expected frequency of substitution

scaling factor to obtain easily computable integer values.

Very likely substitution → high log-odds score  
Unlikely substitution → low log-odds score

# BLOSUM62

(values are rounded to obtain integers)

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

# Hamming distance to handle mismatches

- Let  $X$  be a long string and  $y$  be a short string.
- Find all the substrings  $x$  of  $X$  such that
  - $|x| = |y|$
  - $\delta h(x, y) \leq k$  (hamming distance between  $x$  and  $y$  does not exceed  $k$ )

(Example: quickly align a read to a reference sequence)

# Example

- $X = \text{"ADCABCAABADBBCA"}\text{"}$ ,  $Y = \text{"ADBBCA"}\text{"}$ ,  $k = 3$
- 2 solutions:

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| D | C | A | B | C | A |
| . | . | . |   |   |   |
| A | D | B | B | C | A |

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| A | D | B | B | C | A |
|   |   |   |   |   |   |
| A | D | B | B | C | A |

# Example

$X = \text{"ADCABCAABADBBCA"} , Y = \text{"ADBBCA"} , k = 3$

Two solutions:

ADC**ABCA**ABADBBCA

**ADBBCA**

ADCABCAAB**ADBBCA**

**ADBBCA**

# Recursion/DP formula

- Initialize first row with 0
- Initialize first column with  $k + 1$
- $T[i][j] = T[i-1][j-1] + \text{subs}(i,j)$
- Look at the last line and keep all alignment with a score that does not exceed  $k$

# Hamming distance to handle mismatches

- $X = \text{"ADCABCAABADBBCA"}\text{"}$ ,  $Y = \text{"ADBBCA"}\text{"}$ ,  $k = 3$

| D | - | A | D | C | A | B | C | A | A | B | A | D | B | B | C | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 4 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 |
| D | 4 | 5 | 0 | 2 | 2 | 1 | 2 | 2 | 1 | 1 | 2 | 0 | 2 | 2 | 2 | 2 |
| B | 4 | 5 | 6 | 1 | 3 | 2 | 2 | 3 | 3 | 1 | 2 | 3 | 0 | 2 | 3 | 3 |
| B | 4 | 5 | 6 | 7 | 2 | 3 | 3 | 3 | 4 | 3 | 2 | 3 | 3 | 0 | 3 | 4 |
| C | 4 | 5 | 6 | 6 | 8 | 3 | 3 | 4 | 4 | 5 | 4 | 3 | 4 | 4 | 0 | 4 |
| A | 4 | 4 | 6 | 7 | 6 | 9 | 4 | 3 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 0 |

D C A B C A  
 . . . | | |  
 A D B B C A

A D B B C A  
 | | | | | |  
 A D B B C A



# Remarks

For this problem:

- we don't need to keep the whole table in memory (we only use diagonal paths)
- We can implement early stopping

| D | - | A | D | C | A | B | C | A | A | B | A | D | B | B | C | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| - | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| A | 4 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 1 | 0 |
| D | 4 | 5 | 0 | 2 | 2 | 1 | 2 | 2 | 1 | 1 | 2 | 0 | 2 | 2 | 2 | 2 |
| B | 4 | 5 | 6 | 1 | 3 | 2 | 2 | 3 | 3 | 1 | 2 | 3 | 0 | 2 | 3 | 3 |
| B | 4 | 5 | 6 | 7 | 2 | 3 | 3 | 3 | 4 | 3 | 2 | 3 | 3 | 0 | 3 | 4 |
| C | 4 | 5 | 6 | 6 | 8 | 3 | 3 | 4 | 4 | 5 | 4 | 3 | 4 | 4 | 0 | 4 |
| A | 4 | 4 | 6 | 7 | 6 | 9 | 4 | 3 | 4 | 5 | 5 | 5 | 4 | 5 | 5 | 0 |

# Key points of the lecture (exam-relevant!)

- Local VS global alignment
- Hamming/Edit distances, similarity functions
- Needleman-Wunsch algorithm (and variations):
  - Write the recursion formula
  - Fill a dynamic programming table
  - Backtrace to build the alignment
- Substitution matrices, Blosum, gap penalties

Questions, feedback?

[lukas.huebner@h-its.org](mailto:lukas.huebner@h-its.org)