# Smith-Waterman Algorithm

**AMPP 0708-Q1**

Eduard Ayguade

Juan J. Navarro

Dani Jimenez-Gonzalez

October 4, 2007

# Why compare sequences of aminoacids?

- Proteins are made by aminoacid sequences

  t:c g g g t a t c c a a
- Similar sequences of aminoacids $\rightarrow$ similar protein structures

  t:c g g g t a t c c a a

  s:c c c t a g g t c c c a
- Evolutionary perspective: Mutations?, insertions?, etc.
  - $t_1 = g$ mutated to $s_1 = c$ ?
  - $s_1 = c$ has been an insertion?
- Some evolution are more important/likely than others

# How to compare sequences? Alignment

An alignment of two sequences $t$ and $s$ must satisfy:

- All symbols (residues) in the two sequences have to be in the alignment, and in the same order they appear in the sequences
- We can align one symbol from one sequence with one from the another
- A symbol can be aligned with a blank ('-')
- Two blanks cannot be aligned

```
t: c g g g t a t c c a a
s: c c c t a g g t c c c a

t: c g g g t a - - t - c c a a
s: c c c - t a g g t c c c - a
```

# What is the BEST alignment?

### Example

```
t:  c  g  g  t  a  t  c  c  a  a
s:  c  c  c  t  a  g  g  t  c  c  c  a
```

# What is the BEST alignment?

## Example

```
t:  c g g g t a t c c a a
s:  c c c t a g g t c c c a

t:  c g g g t a - - t - c c a a
s:  c c c - t a g g t c c c - a
```

# What is the BEST alignment?

## Example

```
t: c g g g t a t c c a a
s: c c c t a g g t c c c a

t: c g g g t a - - t - c c a a
s: c c c - t a g g t c c c - a

t: c g g g t a - - - t c c a a
s: c c c - - c t a g g t c c c a
```

# What is the BEST alignment?

### Example

```
t: c g g g t a t c c a a
s: c c c t a g g t c c c a

t: c g g g t a - - t - c c a a
s: c c c - t a g g t c c c - a

t: c g g g t a - - - t c c a a
s: c c - - c t a g g t c c c a

t: c - g g g t a - - t c c a a
s: c c - - c t a g g t c c c a
```

Which is the best?

# Scoring the relationships

- Needed a scoring matrix
- We will be able to find a optimal solution for the scoring matrix at hand



Figure: BLOSUM scoring matrix, S.

## What is the BEST alignment (for that Score Matrix)?

### Example

```
t: c g g g t a t c c a a
s: c c c t a g g t c c c a
```

| t : | c | g | g | g | t | a | − | − | t | − | c | c | a | a | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | +12 | −3 | −3 | −1 | +5 | +5 | −1 | −1 | +5 | −1 | +12 | +12 | −1 | +5 | 45 |
| s : | c | c | c | − | t | a | g | g | t | c | c | c | − | a | |

| t : | c | g | g | g | t | a | − | − | − | t | c | c | a | a | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | +12 | −3 | −1 | −1 | −1 | +0 | −1 | −1 | −1 | +5 | +12 | +12 | −1 | +5 | 36 |
| s : | c | c | − | − | c | t | a | g | g | t | c | c | c | a | |

| t : | c | − | g | g | g | t | a | − | − | t | c | c | a | a | |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| | +12 | −1 | −1 | −1 | −3 | +5 | +5 | −1 | −1 | +5 | +12 | +12 | −1 | +5 | 47 |
| s : | c | c | − | − | c | t | a | g | g | t | c | c | c | a | |

# How to find the best alignment?

- Homology search methods begin with DP algorithms
  - Needleman-Wusch: global search
  - Smith-Waterman (SW): local search
- Faster but less sensitive for larger datasets
  - FASTA
  - BLAST
- Optimal spaced seeds of pattern-writer increase
  - Speed and sensitivity
  - Similar to SW
  - Examples: Pattern Hunter and BLAT
    - SW sensitivity
    - BLAST speed

# Smith-Waterman Algorithm

Figure: Computation Matrix alginment, M

- $N \times N$ integer matrix
- $N$ is sequence length (both $s$ and $t$)
- Compute $M[i][j]$ based on Score Matrix and optimum score compute so far (DP)

|   | 0 | C | G | G | G | T | A | T | C | C | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |   |
| A |   |   |   |   |   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |   |   |   |   |   |
| G |   |   |   |   |   |   |   |   |   |   |   |   |
| T |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |   |   |

# Smith-Waterman Algorithm: Understanding Matrix

Alignment

$t$ :    —    —    —    —    —    —    —    —

$s$ :    c    c    c    t    a    g    g    t



Figure: Aligning $s$ to gaps

# Smith-Waterman Algorithm: Understanding Matrix

Alignment

$t:$    $c$    $g$    $g$    $g$    $t$    $a$    $t$    ...

$s:$    –    –    –    –    –    –    –    ...

|   | 0 | C | G | G | G | T | A | T | C | C | A | A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| T | 0 |   |   |   |   |   |   |   |   |   |   |   |
| A | 0 |   |   |   |   |   |   |   |   |   |   |   |
| G | 0 |   |   |   |   |   |   |   |   |   |   |   |
| G | 0 |   |   |   |   |   |   |   |   |   |   |   |
| T | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |
| C | 0 |   |   |   |   |   |   |   |   |   |   |   |

Figure: Aligning $t$ to gaps

# Smith-Waterman Algorithm: How to compute cell score?

How to find $M[i][j]$?

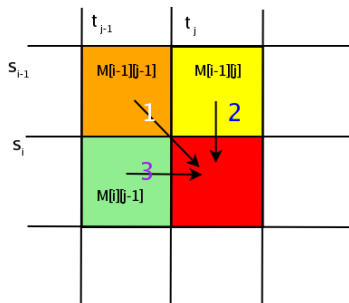Three ways to finish the alignment of $s_{0..i}$ and $t_{0..j}$

1. Score     $s_i$
              $t_j$

2. Gap       $s_i$
   in $t$    $t_j$    —

3. Gap     $s_i$    —
   in $s$         $t_j$
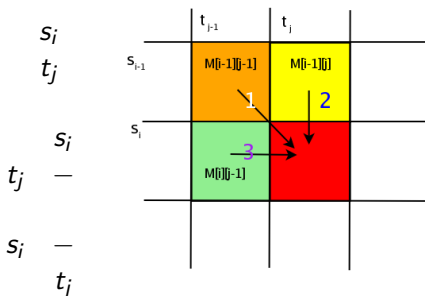
# Smith-Waterman Algorithm: How to compute cell score?

How to find $M[i][j]$?

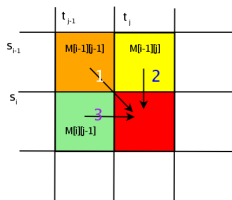Three ways to finish the alignment of $s_{0..i}$ and $t_{0..j}$

① $M[i-1][j-1] + S[s_i][t_j]$

$s_i$
$t_j$

② $M[i-1][j] - g$

$s_i$
$t_j \quad —$

③ $M[i][j-1] - g$

$s_i \quad —$
$t_j$

# Smith-Waterman Algorithm: Scoring Process

Element Computation $M[i][j]$:



$$M[i][0] = 0 \quad M[0][j] = 0$$

$$M[i][j] = \max \begin{cases} 0 \\ M[i-1][j-1] + S[s_i][t_j] & \text{if } s_i \ t_j \\ M[i-1][j] - d & \text{if } s_i \ - \\ M[i][j-1] - d & \text{if } - \ t_j \end{cases}$$

# Smith-Waterman Algorithm: Backtracking Process

If we want to find **BEST** local alignment...

- Find $Score_{opt}$ and then traceback

$$Score_{opt} = \max_{i,j=1}^{N} M[i][j]$$