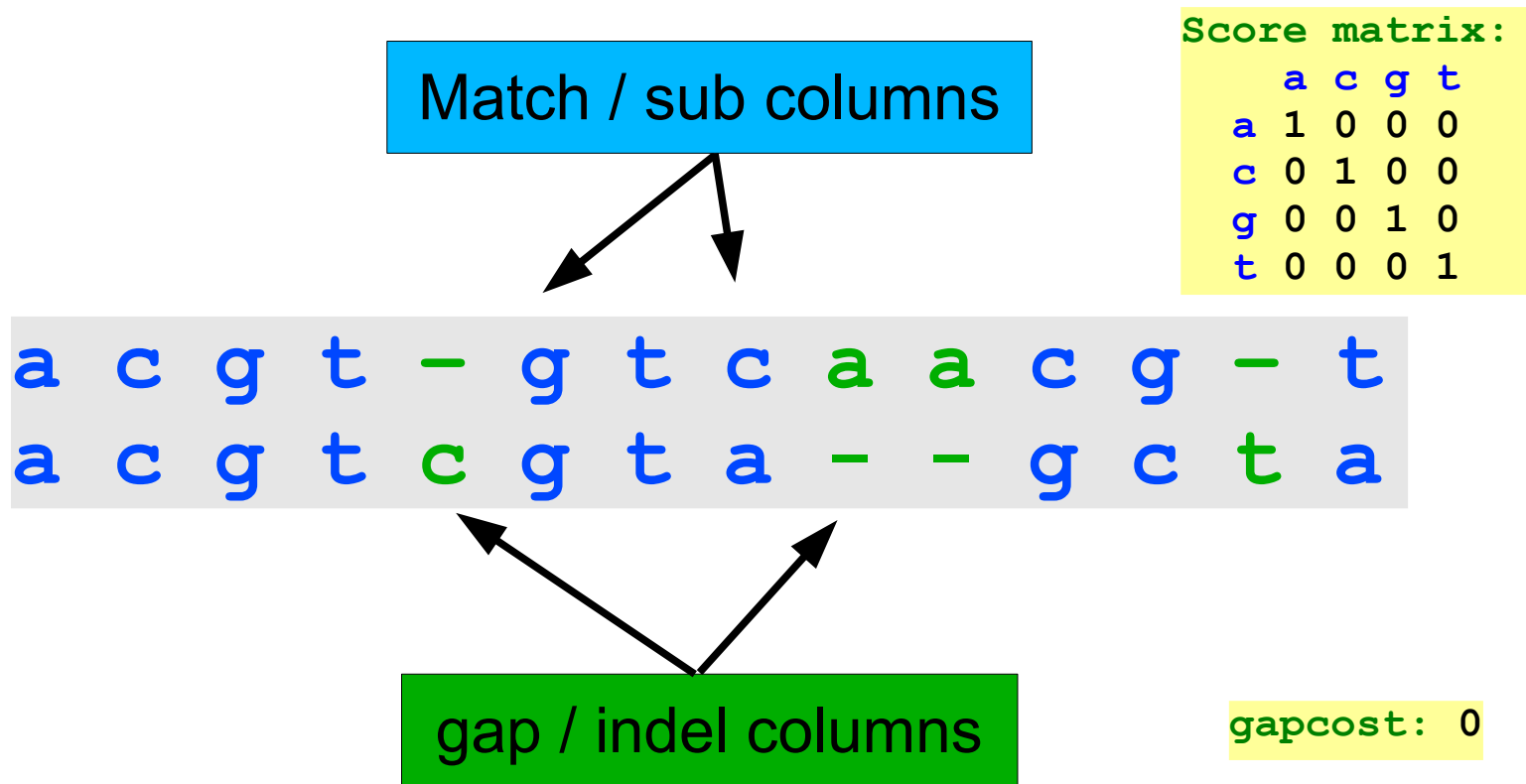


**Global alignment with  
general and affine gapcost**

# Global alignment



## Note about gap cost

In general: cost of "gap block" =  $g(k)$ , where  $k$  is the gap length

Our examples:  $g(k) = c$  "constant gap cost"

$g(k) = b \cdot k$  "linear gap cost"

Many programs:  $g(k) = a + b \cdot k$  "affine gap cost"

# Modeling gapcost

**Biological observation:** longer insertions and deletions (indels) are preferred to shorter indels, i.e. a “good” alignment tends to few long indels rather than many short indels ...

Can the simple algorithm for pairwise alignment be adapted to reflect this additional biological insight, i.e. a better model of biology?

Yes, we introduce the concept of a gapcost-function  $g(k)$  which gives the cost/penalty for a block of  $k$  consecutive insertions or deletions ...

## Example

A	T	A	C	A	-	-	-	C	G	C	A
A	T	-	C	T	C	C	A	C	-	C	T

$$\begin{aligned} & s(A,A) + s(T,T) + g(1) + s(C,C) + s(A,T) + g(3) + \\ & s(C,C) + g(1) + s(C,C) + s(A,T) \end{aligned}$$

# Modeling gapcost

**Biological observation:** longer insertions and deletions (indels) are preferred to shorter indels, i.e. a “good” alignment tends to few long indels rather than many short indels ...

## Computational challenge

Can the simple model be extended to reflect this additional biology? Can we compute an optimal global alignment with general (affine) gapcost efficiently?

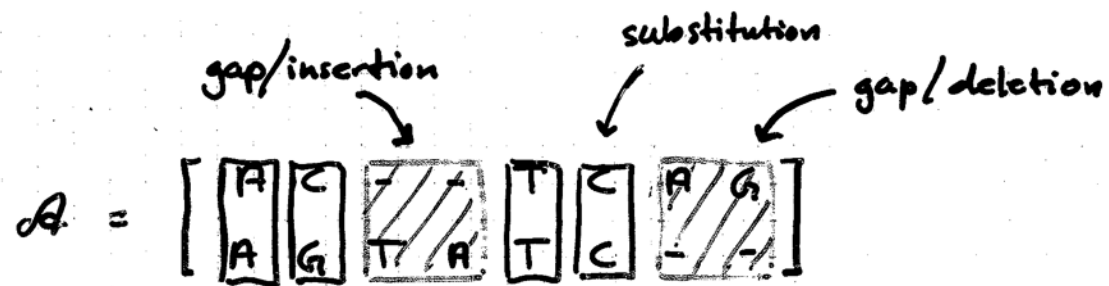
Yes, we introduce the concept of a gapcost-function  $g(k)$  which gives the cost/penalty for a block of  $k$  consecutive insertions or deletions ...

## Example

A	T	A	C	A	-	-	-	C	G	C	A
A	T	-	C	T	C	C	A	C	-	C	T

$$\begin{aligned} & s(A,A) + s(T,T) + g(1) + s(C,C) + s(A,T) + g(3) + \\ & s(C,C) + g(1) + s(C,C) + s(A,T) \end{aligned}$$

## Global alignment of two strings



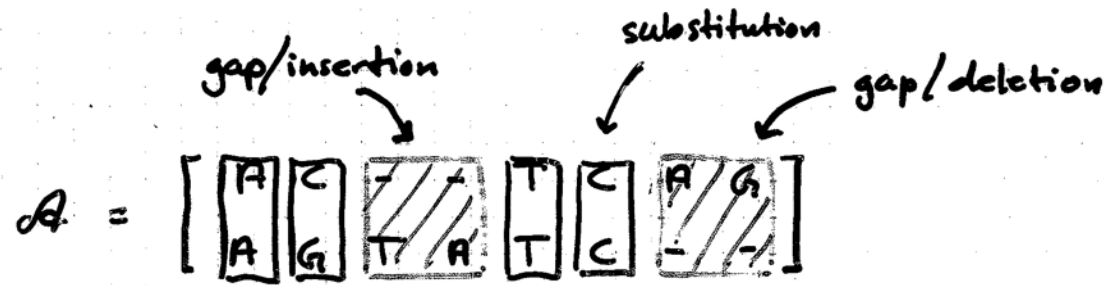
### General cost:

$$\text{cost}(A) = \sum_{\text{blocks}} \text{"cost of block"}$$

Substitution cost:  $s: \Sigma \times \Sigma \rightarrow \mathbb{R}$

gap cost:  $g: \mathbb{N}_+ \rightarrow \mathbb{R}$

## Global alignment of two strings



The book [BA] defines a gap penalty function in definition 8.1.25 as a function  $g: \mathbb{N} \rightarrow \mathbb{R}$ , where  $g(0)=0$  and  $g(k) \geq 0$ , that is **subadditive**, i.e. where  $g(k+l) \leq g(k) + g(l)$ . Why do you think that subadditivity is required?



$$\text{cost}(\mathcal{A}) = \sum_{\text{blocks}} \text{"cost of block"}$$

$$\text{Substitution cost: } s: \Sigma \times \Sigma \rightarrow \mathbb{R}$$

$$\text{gap cost: } g: \mathbb{N}_+ \rightarrow \mathbb{R}$$

# Computing optimal cost with general gap cost

## Intuition

$$S(i, j) = \max_{\text{last block}} \{ \text{"cost of last block"} + \text{"opt. cost of rest."} \}$$

## Formalization

$$S(i, j) = \max \begin{cases} 0 & \text{M}(i, j) \\ \underbrace{s(A[i], B[j]) + S(i-1, j-1)}_{D(i, j)} & \left[ \begin{array}{c} \text{---} A[i] \\ \text{---} B[j] \end{array} \right] \\ \underbrace{\max_{0 < k \leq i} [S(i-k, j) - g(k)]}_{I(i, j)} & \left[ \begin{array}{c} \text{---} A[i-k+1] \dots A[i] \\ \text{---} \end{array} \right] \\ \underbrace{\max_{0 < k \leq j} [S(i, j-k) - g(k)]}_{I(i, j)} & \left[ \begin{array}{c} \text{---} \\ \text{---} B[j-k+1] \dots B[j] \end{array} \right] \end{cases}$$

## Computing optimal cost with general gap cost

Can be implemented using dynamic programming

Running time:  $(\frac{n}{2})^2 \cdot \frac{n}{2} \leq T(n) \leq n^2 \cdot n \quad O(n^3)$  

Space consumption:  $O(n^2)$

$$S(i, j) = \max \left\{ \begin{array}{l} 0 \\ \underbrace{s(A[i], B[j]) + S(i-1, j-1)}_{D(i, j)} \quad \left[ \begin{array}{c} \text{--- } A[i] \\ \text{--- } B[j] \end{array} \right] \\ \underbrace{\max_{0 < k \leq i} [S(i-k, j) - g(k)]}_{I(i, j)} \quad \left[ \begin{array}{c} \text{--- } A[i-k+1] \dots A[i] \\ \text{--- } \phantom{A[i-k+1] \dots A[i]} \end{array} \right] \\ \underbrace{\max_{0 < k \leq j} [S(i, j-k) - g(k)]}_{I(i, j)} \quad \left[ \begin{array}{c} \text{--- } \phantom{B[j-k+1] \dots B[j]} \\ \text{--- } B[j-k+1] \dots B[j] \end{array} \right] \end{array} \right.$$



## Computing optimal cost with affine gap cost

$$g(k) = \alpha \cdot k + \beta, \quad \alpha, \beta \geq 0$$

### Intuition

... as general gap cost but reduce time to compute

$D(i, j)$  and  $I(i, j)$ . Let us consider  $D(i, j)$  ...

### Trick

Consider the best deletion-block, two possibilities:

① A new block

$$\begin{array}{cc} \overbrace{M(i-1, j) - (\alpha + \beta)} & \overbrace{I(i-1, j) - (\alpha + \beta)} \\ \left[ \begin{array}{cc} \sim AC[i-1] & \boxed{AC[i]} \\ \sim B[j] & \boxed{-} \end{array} \right] & \left[ \begin{array}{cc} \sim & \boxed{AC[i]} \\ \sim B[j] & \boxed{-} \end{array} \right] \end{array}$$

② Continuation of existing block

$$\begin{array}{c} \overbrace{D(i-1, j) - \alpha} \\ \left[ \begin{array}{cc} \sim AC[i-k] \dots AC[i-0] & \boxed{AC[i]} \\ \sim B[j] & \boxed{- \dots -} \end{array} \right] \end{array}$$

## Computing optimal cost with affine gap cost

$$g(k) = \alpha \cdot k + \beta$$

Because

$$S(i-1, j) = \max\{M(i-1, j), I(i-1, j), D(i-1, j)\}$$

$$\text{and } D(i-1, j) - (\alpha + \beta) \leq D(i-1, j) - \alpha$$

Hence,

$$D(i, j) = \max \{ M(i-1, j) - (\alpha + \beta), I(i-1, j) - (\alpha + \beta), D(i-1, j) - \alpha \}$$

$$= \max \{ S(i-1, j) - (\alpha + \beta), D(i-1, j) - \alpha \}$$

Trick

Consider the best deletion-block, two possibilities:

① A new block

$$\begin{array}{cc} \overbrace{M(i-1, j) - (\alpha + \beta)} & \overbrace{I(i-1, j) - (\alpha + \beta)} \\ \left[ \begin{array}{cc} \sim A[i-1] & A[i] \\ \sim B[j] & - \end{array} \right] & \left[ \begin{array}{cc} \sim & - \\ \sim & B[j] \end{array} \right] \end{array}$$

② Continuation of existing block

$$\overbrace{D(i-1, j) - \alpha}$$

$$\left[ \begin{array}{cc} \sim A[i-k] \dots A[i-1] & A[i] \\ \sim B[j] & - \end{array} \right]$$

## Computing optimal cost with affine gap cost

$$S(i, j) = \max \begin{cases} 0 & i = 0 \text{ and } j = 0 \\ S(i-1, j-1) + s(A[i], B[j]) & i > 0 \text{ and } j > 0 \\ D(i, j) & i > 0 \text{ and } j \geq 0 \\ I(i, j) & i \geq 0 \text{ and } j > 0 \end{cases}$$

$$D(i, j) = \max \begin{cases} S(i-1, j) - (\alpha + \beta) & i > 0 \text{ and } j \geq 0 \\ D(i-1, j) - \alpha & i > 1 \text{ and } j \geq 0 \end{cases}$$

$$I(i, j) = \max \begin{cases} S(i, j-1) - (\alpha + \beta) & i \geq 0 \text{ and } j > 0 \\ I(i, j-1) - \alpha & i \geq 0 \text{ and } j > 1 \end{cases}$$

Can be implemented using dyn. prog./memorization  
using 3 tables of size  $\Theta(n \cdot m)$

## Computing optimal cost with affine gap cost

$$S(i, j) = \max \begin{cases} 0 & i = 0 \text{ and } j = 0 \\ S(i-1, j-1) + s(A[i], B[j]) & i > 0 \text{ and } j > 0 \\ D(i, j) & i > 0 \text{ and } j \geq 0 \\ I(i, j) & i \geq 0 \text{ and } j > 0 \end{cases}$$

$$D(i, j) = \max \begin{cases} S(i-1, j) - (\alpha + \beta) & i > 0 \text{ and } j \geq 0 \\ D(i-1, j) - \alpha & i > 1 \text{ and } j \geq 0 \end{cases}$$

$$I(i, j) = \max \begin{cases} S(i, j-1) - (\alpha + \beta) & i \geq 0 \text{ and } j > 0 \\ I(i, j-1) - \alpha & i \geq 0 \text{ and } j > 1 \end{cases}$$

Can be implemented using dyn. prog./memorization  
using 3 tables of size  $\Theta(n \cdot m)$

## Computing optimal cost with affine gap cost

Time and space is  $O(nm)$ . An optimal alignment can be retrieved by backtracking in time  $O(n)$ .

$$S(i, j) = \max \begin{cases} 0 & i = 0 \text{ and } j = 0 \\ S(i-1, j-1) + s(A[i], B[j]) & i > 0 \text{ and } j > 0 \\ D(i, j) & i > 0 \text{ and } j \geq 0 \\ I(i, j) & i \geq 0 \text{ and } j > 0 \end{cases}$$

$$D(i, j) = \max \begin{cases} S(i-1, j) - (\alpha + \beta) & i > 0 \text{ and } j \geq 0 \\ D(i-1, j) - \alpha & i > 1 \text{ and } j \geq 0 \end{cases}$$

$$I(i, j) = \max \begin{cases} S(i, j-1) - (\alpha + \beta) & i \geq 0 \text{ and } j > 0 \\ I(i, j-1) - \alpha & i \geq 0 \text{ and } j > 1 \end{cases}$$

Can be implemented using dyn. prog./memorization  
using 3 tables of size  $\Theta(n \cdot m)$

# Global alignment with convex gap cost

$g(k)$  = "a convex function, fx  $\log(k)$  or  $\text{sqrt}(k)$ "

$$S(i, j) = \max \begin{cases} 0 \\ \underbrace{M(i, j)}_{D(i, j)} = s(A[i], B[j]) + S(i-1, j-1) \quad \left[ \begin{array}{c} \text{---} A[i] \\ \text{---} B[j] \end{array} \right] \\ \underbrace{I(i, j)}_{I(i, j)} = \max_{0 < k \leq i} [S(i-k, j) - g(k)] \quad \left[ \begin{array}{c} \text{---} A[i-k+1] \dots A[i] \\ \text{---} - \quad - \quad - \quad - \end{array} \right] \\ \max_{0 < k \leq j} [S(i, j-k) - g(k)] \quad \left[ \begin{array}{c} \text{---} - \quad - \quad - \quad - \\ \text{---} B[j-k+1] \dots B[j] \end{array} \right] \end{cases}$$

In cubic time using above recursion, but it can be done in  $O(nm \log(n))$