

# Non-linear gap penalties

Waterman-Smith-Beyer (1976)

Gotoh (1982)

Rolf Backofen

Lehrstuhl für Bioinformatik

Institut für Informatik

Albert-Ludwigs-Universität Freiburg

Course Bioinformatics I

built on November 24, 2020

- **Observation:** gaps are different in nature - given a fixed number of gaps, a "small number of long gaps" is biologically likelier than a "big number of small gaps"

- **Observation:** gaps are different in nature - given a fixed number of gaps, a "small number of long gaps" is biologically likelier than a "big number of small gaps"
- **Problem:** with linear gap score not distinguishable:

$$\text{score} \left( \begin{array}{c} \text{AAA---} \\ \text{---TTT} \end{array} \right) = \text{score} \left( \begin{array}{c} \text{A-A-A-} \\ \text{-T-T-T} \end{array} \right)$$

- **Observation:** gaps are different in nature - given a fixed number of gaps, a "small number of long gaps" is biologically likelier than a "big number of small gaps"

- **Problem:** with linear gap score not distinguishable:

$$\text{score} \left( \begin{array}{c} \text{AAA---} \\ \text{---TTT} \end{array} \right) = \text{score} \left( \begin{array}{c} \text{A-A-A-} \\ \text{-T-T-T} \end{array} \right)$$

- **Solution:** use explicit gap function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$

- **Observation:** gaps are different in nature - given a fixed number of gaps, a "small number of long gaps" is biologically likelier than a "big number of small gaps"

- **Problem:** with linear gap score not distinguishable:

$$\text{score} \left( \begin{array}{c} \text{AAA---} \\ \text{---TTT} \end{array} \right) = \text{score} \left( \begin{array}{c} \text{A-A-A-} \\ \text{-T-T-T} \end{array} \right)$$

- **Solution:** use explicit gap function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$

- **Example:**

length 2   length 3  
 A   A   A - - - G T A  
 -   -   A T T T G T -  
 length 1

⇒ Wanted: gap cost  $g(2) + g(3) + g(1) \leq (2+3+1) * g(1)$

# Gaps

- **Observation:** gaps are different in nature - given a fixed number of gaps, a "small number of long gaps" is biologically likelier than a "big number of small gaps"

- **Problem:** with linear gap score not distinguishable:

$$\text{score} \begin{pmatrix} \text{AAA---} \\ \text{---TTT} \end{pmatrix} = \text{score} \begin{pmatrix} \text{A-A-A-} \\ \text{-T-T-T} \end{pmatrix}$$

- **Solution:** use explicit gap function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$

- **Example:**

length 2   length 3

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

length 1

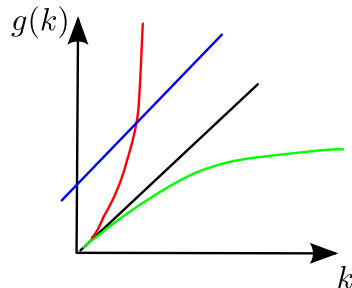
gap score difference  
translates to  
alignment scores  
 $\Rightarrow$  distinguishable

$$\Rightarrow \text{Wanted: gap cost } g(2) + g(3) + g(1) \leq (2+3+1) * g(1)$$

## Definition (Gap penalty)

A *gap penalty* is a function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$  that is subadditive, i.e.,

$$\forall k, l : g(k + l) \leq g(k) + g(l).$$

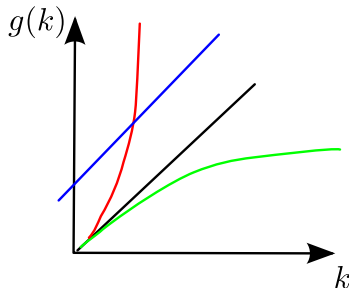


## Definition (Gap penalty)

A *gap penalty* is a function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$  that is subadditive, i.e.,

$$\forall k, l : g(k + l) \leq g(k) + g(l).$$

$g(k) = \alpha + \beta k^2$  ⚡ not subadditive!  
 $\Rightarrow$  unintended behaviour





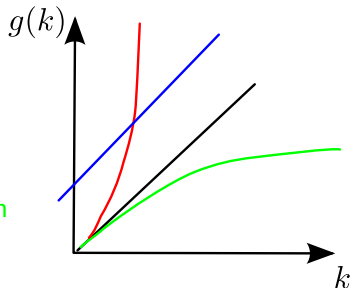
## Definition (Gap penalty)

A *gap penalty* is a function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$  that is subadditive, i.e.,

$$\forall k, l : g(k + l) \leq g(k) + g(l).$$

$g(k) = \alpha + \beta k^2$  ⚡ not subadditive!  
⇒ unintended behaviour

$g(k) = \alpha + \beta \ln(k)$  ✓  
⇒ biologically the best approximation



## Definition (Gap penalty)

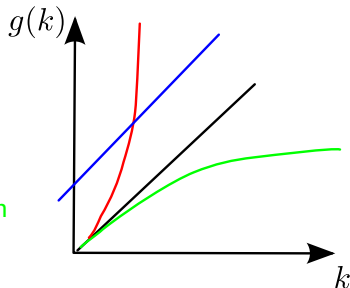
A *gap penalty* is a function  $g : \mathbb{N}^+ \rightarrow \mathbb{R}$  that is subadditive, i.e.,

$$\forall k, l : g(k + l) \leq g(k) + g(l).$$

$g(k) = \alpha + \beta k^2$  ⚡ not subadditive!  
⇒ unintended behaviour

$g(k) = \alpha + \beta \ln(k)$  ✓  
⇒ biologically the best approximation

$g(k) = \alpha + \beta k$   
⇒ affine, very common



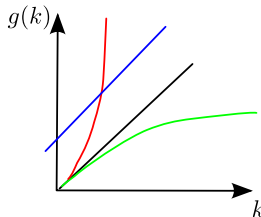
# Affine gap penalties

**Idea:** make initiation of gaps more expensive than extension of an existing gap

## Definition (Affine gap penalty)

A gap penalty is called *affine* if there are  $\alpha, \beta \in \mathbb{R}$  such that

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



- **Now:** How to calculate optimal alignments using a gap function?
- **Problem:** split like Needleman/Wunsch not working

$$w \begin{pmatrix} A & A & - & - \\ - & A & G & T \end{pmatrix}$$

$$g(1) + w(A, A) + g(2)$$

- **Now:** How to calculate optimal alignments using a gap function?
- **Problem:** split like Needleman/Wunsch not working

$$w \left( \begin{array}{cc|c} A & A & - \\ - & A & G \end{array} \right)$$

$$w \left( \begin{array}{cc|c} A & A & - \\ - & A & G \end{array} \right) + w \left( \begin{array}{c} - \\ T \end{array} \right)$$

$$g(1) + w(A, A) + g(2) \quad \neq \quad (g(1) + w(A, A) + g(1)) + g(1)$$

④ Substitution: ✓

$$\begin{array}{lcl} a^\diamond & = & \dots \quad | \quad a_i \\ b^\diamond & = & \dots \quad | \quad b_j \end{array} \quad \Rightarrow \quad D_{i,j} = D_{i-1,j-1} + w(a_i, b_j)$$

# Solution: more distinctions

- ① Substitution: ✓

$$\begin{array}{lcl} a^\diamond & = & \dots \quad | \quad a_i \\ b^\diamond & = & \dots \quad | \quad b_j \end{array} \Rightarrow D_{i,j} = D_{i-1,j-1} + w(a_i, b_j)$$

- ② Delete 1 position:

$$\begin{array}{lcl} a^\diamond & = & \dots \quad ? \quad | \quad a_i \\ b^\diamond & = & \dots \quad b_j \quad | \quad - \end{array} \Rightarrow D_{i,j} = D_{\underline{i-1},j} + g(1)$$

# Solution: more distinctions

- ① Substitution: ✓

$$\begin{array}{lcl} a^\diamond & = & \dots \quad | \quad a_i \\ b^\diamond & = & \dots \quad | \quad b_j \end{array} \quad \Rightarrow \quad D_{i,j} = D_{i-1,j-1} + w(a_i, b_j)$$

- ② Delete 1 position:

$$\begin{array}{lcl} a^\diamond & = & \dots \quad ? \quad | \quad a_i \\ b^\diamond & = & \dots \quad b_j \quad | \quad - \end{array} \quad \Rightarrow \quad D_{i,j} = D_{\underline{i-1},j} + g(1)$$

- ③ Delete 2 positions:

$$\begin{array}{lcl} a^\diamond & = & \dots \quad ? \quad | \quad a_{i-1} \quad a_i \\ b^\diamond & = & \dots \quad b_j \quad | \quad - \quad - \end{array} \quad \Rightarrow \quad D_{i,j} = D_{\underline{i-2},j} + g(2)$$

- ④ ...

⇒ Algorithm of Waterman-Smith-Beyer (1976)



**Theorem (Waterman, Smith, and Beyer (1976))**

Let  $g : \mathbb{N} \rightarrow \mathbb{R}$  be a gap penalty and  $w$  be a distance function on  $\Sigma \times \Sigma$ . Let  $a = a_1 \dots a_n$  and  $b = b_1 \dots b_m$  be two words in  $\Sigma^*$ . We define  $(D_{i,j})$  with  $1 \leq i \leq n$  and  $1 \leq j \leq m$  by

$$D_{0,0} = 0,$$

$$D_{0,j} = g(j),$$

$$D_{i,0} = g(i),$$

$$D_{i,j} = \min \left\{ \begin{array}{l} \min_{1 \leq k \leq j} \{D_{i,j-k} + g(k)\}, \\ D_{i-1,j-1} + w(a_i, b_j), \\ \min_{1 \leq k \leq i} \{D_{i-k,j} + g(k)\} \end{array} \right\}.$$

Then  $D_{i,j} = \text{opt. prefix alignment score for } a_1 \dots a_i \text{ and } b_1 \dots b_j$ .

⇒ on average a cell cost  $O(n)$  for filling

⇒ total:  $O(n^2)$  space and  $O(n^3)$  time

- Example:
  - 2 RNA sequences with  $n = 30\,000 = 3 \cdot 10^4$

- **assume:** computer with 1 Ghz

+ 1 operation per unit

$$\begin{aligned}\Rightarrow \frac{27 \cdot 10^{12} \text{ops}}{10^9 \text{ops/s}} &= 27 \cdot 10^3 \text{ s} \\ &/ 3,600 \frac{\text{s}}{\text{h}} \\ &\approx 10.5 \text{ h}\end{aligned}$$

- How much time a *quadratic* algorithm would have taken?

- W-S-B problem: arbitrarily long gaps tested in each step
- Solution: restrict to affine gap penalties  $g(k) = \alpha + \beta k$

- W-S-B problem: arbitrarily long gaps tested in each step
- Solution: restrict to affine gap penalties  $g(k) = \alpha + \beta k$

Analyzing gap case:

$$\begin{array}{rcl} a^\diamond & = & \dots \\ b^\diamond & = & \dots \end{array} \left| \begin{array}{l} a_i \\ - \end{array} \right. \Rightarrow \text{look at subcases}$$

# Gotoh's idea

- W-S-B problem: arbitrarily long gaps tested in each step
- Solution: restrict to affine gap penalties  $g(k) = \alpha + \beta k$

## Analyzing gap case:

$$\begin{array}{rcl} a^\diamond & = & \dots \mid a_i \\ b^\diamond & = & \dots \mid - \end{array} \Rightarrow \text{look at subcases}$$

$$\textcircled{1} \quad \begin{array}{rcl} \dots & ? & a_i \\ \dots & b_j & - \end{array} \Rightarrow D_{i,j} = D_{i-1,j} + g(1)$$

- W-S-B problem: arbitrarily long gaps tested in each step
- Solution: restrict to affine gap penalties  $g(k) = \alpha + \beta k$

## Analyzing gap case:

$$\begin{array}{lcl} a^\diamond & = & \dots \mid a_i \\ b^\diamond & = & \dots \mid - \end{array} \quad \Rightarrow \quad \text{look at subcases}$$

$$\textcircled{1} \quad \begin{array}{lcl} \dots & ? & a_i \\ \dots & b_j & - \end{array} \quad \Rightarrow \quad D_{i,j} = D_{i-1,j} + g(1)$$

$$\begin{array}{lcl} \textcircled{2} & \begin{array}{lcl} \dots & a_{i-1} & a_i \\ \dots & - & - \end{array} & \Rightarrow \quad \begin{aligned} D_{i,j} &= \star - g(k-1) + g(k) \\ &= \star - \alpha - (k-1)\beta + \alpha + k\beta \\ &= \star + \beta \end{aligned} \end{array}$$

$\underbrace{\hspace{10em}}_{k \text{ gaps}}$

$(\star : \text{best score ending in gap } \begin{smallmatrix} a_{i-1} \\ - \end{smallmatrix})$

... analogously for end gaps in  $a^\diamond$

# Gotoh's idea

- ⇒ recursion cases:
- a. no gap  $\Rightarrow +w(..)$
  - b. starting a **new gap**  $\Rightarrow +g(1)$
  - c. **elongate** an existing gap  $\Rightarrow +\beta$
- ⇒ length of the gap doesn't matter
- ⇒ saving time because:
- W-S-B: test with all possible gap lengths
  - Gotoh: just add  $\beta$  if a gap is elongated

- ⇒ recursion cases:
- a. no gap ⇒  $+w(..)$
  - b. starting a **new gap** ⇒  $+g(1)$
  - c. **elongate** an existing gap ⇒  $+\beta$

⇒ length of the gap doesn't matter

⇒ saving time because:

- W-S-B: test with all possible gap lengths
  - Gotoh: just add  $\beta$  if a gap is elongated
- 
- **Comment:** if gap penalty is not affine (e.g.  $g(k) = \alpha + \beta \cdot \ln(k)$ ) then
$$\begin{aligned} D_{i,j} &= \star - g(k-1) + g(k) \\ &= \dots \end{aligned}$$



# Gotoh's idea

- ⇒ recursion cases:
- a. no gap ⇒  $+w(..)$
  - b. starting a **new gap** ⇒  $+g(1)$
  - c. **elongate** an existing gap ⇒  $+\beta$
- ⇒ length of the gap doesn't matter
- ⇒ saving time because:
- W-S-B: test with all possible gap lengths
  - Gotoh: just add  $\beta$  if a gap is elongated
- **Comment:** if gap penalty is not affine (e.g.  $g(k) = \alpha + \beta \cdot \ln(k)$ ) then
- $$\begin{aligned} D_{i,j} &= \star - g(k-1) + g(k) \\ &= \star - \alpha - \ln(k-1) + \alpha + \ln(k) \\ &= \star + \ln(k) - \ln(k-1) \\ &= \star + \ln\left(\frac{k}{k-1}\right) \end{aligned}$$
- ⇒ depends on  $k$  ⇒ **Gotoh's idea doesn't work**

⇒ further matrices needed

- $(D_{i,j})$  cost for alignment of prefixes  $(a_1 \dots a_i, b_1 \dots b_j)$
- $(P_{i,j})$  cost for alignment of prefixes  $(a_1 \dots a_i, b_1 \dots b_j)$  that ends with a gap in  $b^\diamond$  (i.e., last column is  $\begin{pmatrix} a_i \\ - \end{pmatrix}$ )
- $(Q_{i,j})$  cost for alignment of prefixes  $(a_1 \dots a_i, b_1 \dots b_j)$  that ends with a gap in  $a^\diamond$  (i.e., last column is  $\begin{pmatrix} - \\ b_j \end{pmatrix}$ )

## Gotoh (1982)

- affine gap penalty  $g(k) = \alpha + k\beta$ ; distance function  $w : \Sigma \times \Sigma \rightarrow \mathbb{R}$
- recursive definition of matrices  $(D_{i,j})$ ,  $(P_{i,j})$ , and  $(Q_{i,j})$ :

# Gotoh (1982)

- affine gap penalty  $g(k) = \alpha + k\beta$ ; distance function  $w : \Sigma \times \Sigma \rightarrow \mathbb{R}$
- recursive definition of matrices  $(D_{i,j})$ ,  $(P_{i,j})$ , and  $(Q_{i,j})$ :

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j-1} + w(a_i, b_j) \\ P_{i,j} \\ Q_{i,j} \end{array} \right\},$$

with  $i, j \geq 1$ , where for  $1 \leq i \leq |a|$  and  $1 \leq j \leq |b|$ ,

- affine gap penalty  $g(k) = \alpha + k\beta$ ; distance function  $w : \Sigma \times \Sigma \rightarrow \mathbb{R}$
- recursive definition of matrices  $(D_{i,j})$ ,  $(P_{i,j})$ , and  $(Q_{i,j})$ :

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j-1} + w(a_i, b_j) \\ P_{i,j} \\ Q_{i,j} \end{array} \right\},$$

with  $i, j \geq 1$ , where for  $1 \leq i \leq |a|$  and  $1 \leq j \leq |b|$ ,

$$P_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j} + g(1) \\ P_{i-1,j} + \beta \end{array} \right\}$$

$$Q_{i,j} = \min \left\{ \begin{array}{l} D_{i,j-1} + g(1) \\ Q_{i,j-1} + \beta \end{array} \right\}$$

- affine gap penalty  $g(k) = \alpha + k\beta$ ; distance function  $w : \Sigma \times \Sigma \rightarrow \mathbb{R}$
- recursive definition of matrices  $(D_{i,j})$ ,  $(P_{i,j})$ , and  $(Q_{i,j})$ :

$$D_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j-1} + w(a_i, b_j) \\ P_{i,j} \\ Q_{i,j} \end{array} \right\},$$

with  $i, j \geq 1$ , where for  $1 \leq i \leq |a|$  and  $1 \leq j \leq |b|$ ,

$$P_{i,j} = \min \left\{ \begin{array}{l} D_{i-1,j} + g(1) \\ P_{i-1,j} + \beta \end{array} \right\}$$

$$Q_{i,j} = \min \left\{ \begin{array}{l} D_{i,j-1} + g(1) \\ Q_{i,j-1} + \beta \end{array} \right\}$$

order of calculation:  $P_{i,j}$ ,  $Q_{i,j}$  before  $D_{i,j}$

- $D$  as usual:  
 $D_{0,0} = 0,$   
 $D_{0,j} = g(j),$   
 $D_{i,0} = g(i)$

# Gotoh (1982) - Initialization

- $D$  as usual:
$$\begin{aligned}D_{0,0} &= 0, \\D_{0,j} &= g(j), \\D_{i,0} &= g(i)\end{aligned}$$
- $P$ :
  - recursion only on the **first** index ( $P_{i,j} \rightarrow P_{i-1,j} \rightarrow \dots \rightarrow P_{0,j}$ )
  - **hence:** only initialization for  $P_{0,j} \Rightarrow P_{i,0}$  not used



# Gotoh (1982) - Initialization

- $D$  as usual:  $D_{0,0} = 0$ ,  
 $D_{0,j} = g(j)$ ,  
 $D_{i,0} = g(i)$
- $P$ :
  - recursion only on the **first** index ( $P_{i,j} \rightarrow P_{i-1,j} \rightarrow \dots \rightarrow P_{0,j}$ )
  - **hence:** only initialization for  $P_{0,j} \Rightarrow P_{i,0}$  not used
  - **but:**  $P_{0,j}$  is best alignment of  $\epsilon$  and  $b_1 \dots b_j$  that ends with gap in  $b^\diamond \Rightarrow$  the only possible alignment would be:

$$\begin{array}{ccccccc} - & - & \dots & - & - & - \\ b_1 & b_2 & \dots & b_{j-1} & b_j & - \end{array}$$

$\updownarrow$

**disallowed in alignments!**

$\Rightarrow P_{0,j} = \infty$

# Gotoh (1982) - Initialization

- $D$  as usual:  $D_{0,0} = 0$ ,  
 $D_{0,j} = g(j)$ ,  
 $D_{i,0} = g(i)$
- $P$ :
  - recursion only on the **first** index ( $P_{i,j} \rightarrow P_{i-1,j} \rightarrow \dots \rightarrow P_{0,j}$ )
  - **hence:** only initialization for  $P_{0,j} \Rightarrow P_{i,0}$  not used
  - **but:**  $P_{0,j}$  is best alignment of  $\epsilon$  and  $b_1 \dots b_j$  that ends with gap in  $b^\diamond \Rightarrow$  the only possible alignment would be:

$$\begin{array}{ccccccc} - & - & & \dots & & - & - & - \\ b_1 & b_2 & & \dots & & b_{j-1} & b_j & - \end{array}$$

$\updownarrow$

**disallowed in alignments!**

$\Rightarrow P_{0,j} = \infty$

- $Q$  analogously ...

- given:  $a = \text{CC}$  and  $b = \text{ACCT}$ .
- cost functions:
  - substitutions:  $w(x, y) = \begin{cases} 0 & \text{if } x = y \\ 1 & \text{else} \end{cases}$
  - gap penalty:  $g(k) = 4 + k$  ( $\beta = 1$ ).
- **wanted:** optimal alignment using Gotoh

# Example - filling of matrices

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—				
C	—				

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$				
C	$\infty$				

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5				
C	6				

# Example - filling of matrices

$(P_{i,j}) =$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10			
C	—				

$(Q_{i,j}) =$

		A	C	C	T
	0	—	—	—	—
C	$\infty$				
C	$\infty$				

$(D_{i,j}) =$

		A	C	C	T
	0	5	6	7	8
C	5				
C	6				

$$P_{1,1} = \min \left\{ \begin{array}{ll} D_{1,0} + g(1) & (= 5 + 5 = 10) \\ P_{1,0} + \beta & (= \infty + 1 = \infty) \end{array} \right\} = 10$$

# Example - filling of matrices

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—				
C	—				

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	$\overset{a}{\leftarrow} 10$			
C	$\infty$	$\overset{b}{\rightarrow}$			

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5				
C	6				

$$Q_{1,1} = \min \left\{ \begin{array}{ll} D_{0,1} + g(1) & (= 5 + 5 = 10) \\ Q_{0,1} + \beta & (= \infty + 1 = \infty) \end{array} \right\} = 10$$

# Example - filling of matrices

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10			
C	—				

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10			
C	$\infty$				

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5	1	6		
C	6				

$$D_{1,1} = \min \left\{ \begin{array}{ll} D_{0,0} + w(C, A) & (= 1) \\ P_{1,1} & (= 10) \\ Q_{1,1} & (= 10) \end{array} \right\}$$

$= 1$

# Example - filling of matrices

$(P_{i,j}) =$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10			
C	—				

$(Q_{i,j}) =$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10			
C	$\infty$				

$(D_{i,j}) =$

		A	C	C	T
	0	5	6	7	8
C	5	1			
C	6				

$$P_{1,1} = \min \left\{ \begin{array}{ll} D_{1,0} + g(1) & (= 5 + 5 = 10) \\ P_{1,0} + \beta & (= \infty + 1 = \infty) \end{array} \right\} = 10$$

$$Q_{1,1} = \min \left\{ \begin{array}{ll} D_{0,1} + g(1) & (= 5 + 5 = 10) \\ Q_{0,1} + \beta & (= \infty + 1 = \infty) \end{array} \right\} = 10$$

$$D_{1,1} = \min \left\{ \begin{array}{ll} D_{0,0} + w(C, A) & (= 1) \\ P_{1,1} & (= 10) \\ Q_{1,1} & (= 10) \end{array} \right\} = 1$$



# Example - traceback

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10	11	12	13
C	—	6	10	11	13

Final matrices

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10	6	7	8
C	$\infty$	11	11	6	7

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5	1	5	6	8
C	6	6	1	5	7

# Example - traceback

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10	11	12	13
C	—	6	10	11	13

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10	6	7	8
C	$\infty$	11	11	6	7

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5	1	5	6	8
C	6	6	1	5	7

Final matrices and tracebacks

1.  $D \nwarrow D \nwarrow D \leftarrow Q \leftarrow Q \bullet$

C	C	—	—
A	C	C	T

# Example - traceback

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10	11	12	13
C	—	6	10	11	13

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10	6	7	8
C	$\infty$	11	11	6	7

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5	1	5	6	8
C	6	6	1	5	7

Final matrices and tracebacks

1.  $D \nwarrow D \nwarrow D \leftarrow Q \leftarrow Q \bullet$

C	C	—	—
A	C	C	T

2.  $D \leftarrow D \leftarrow D \nwarrow D \nwarrow$

—	—	C	C
A	C	C	T

# Example - traceback

$$(P_{i,j}) =$$

		A	C	C	T
	0	$\infty$	$\infty$	$\infty$	$\infty$
C	—	10	11	12	13
C	—	6	10	11	13

$$(Q_{i,j}) =$$

		A	C	C	T
	0	—	—	—	—
C	$\infty$	10	6	7	8
C	$\infty$	11	11	6	7

$$(D_{i,j}) =$$

		A	C	C	T
	0	5	6	7	8
C	5	1	5	6	8
C	6	6	1	5	7

Final matrices and tracebacks

1.  $D \nwarrow D \nwarrow D \leftarrow Q \leftarrow Q \bullet$

C	C	—	—
A	C	C	T

2.  $D \leftarrow D \leftarrow D \nwarrow D \nwarrow$

—	—	C	C
A	C	C	T

- Problem: long gaps are often single evol. events
  - ⇒ linear gap scoring leads to gap distribution
  - ⇒ too high penalty for long gaps (need subadditive gap scoring)

- Problem: long gaps are often single evol. events
  - ⇒ linear gap scoring leads to gap distribution
  - ⇒ too high penalty for long gaps (need subadditive gap scoring)
- Waterman-Smith-Beyer: optimization with subadditive gap score
  - ⇒ explicit consideration of every gap length
  - ⇒ increases runtime complexity to  $O(n^3)$  compared to linear gap cost

- Problem: long gaps are often single evol. events
  - ⇒ linear gap scoring leads to gap distribution
  - ⇒ too high penalty for long gaps (need subadditive gap scoring)
- Waterman-Smith-Beyer: optimization with subadditive gap score
  - ⇒ explicit consideration of every gap length
  - ⇒ increases runtime complexity to  $O(n^3)$  compared to linear gap cost
- Gotoh: optimization with affine gap score
  - ⇒ distinction of gap start and extension
  - ⇒ auxiliary matrices needed for gap handling
  - ⇒ runtime complexity again  $O(n^2)$