

Homology  $\longrightarrow$  Similarity

$x: \overbrace{A C A T G C G C A T}^n$   
 $y: A G T T C C A A T T$

mismatch

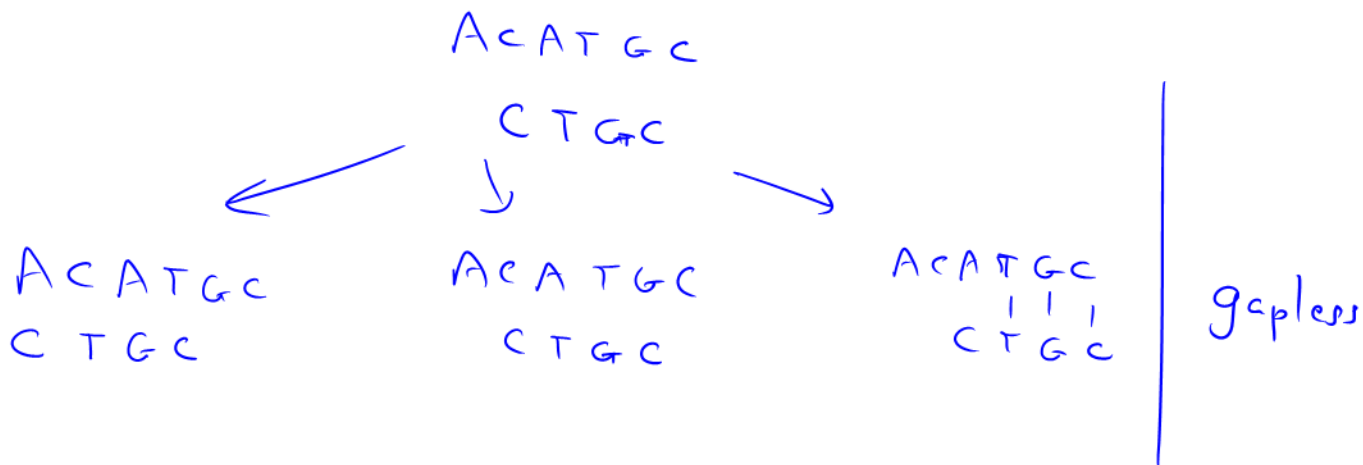
# of exact matches

score

	A	C	G	T
A	1			0.9
C		1		
G			1	
T				1

$$S = \sum_{i=1}^n s(x_i, y_i)$$

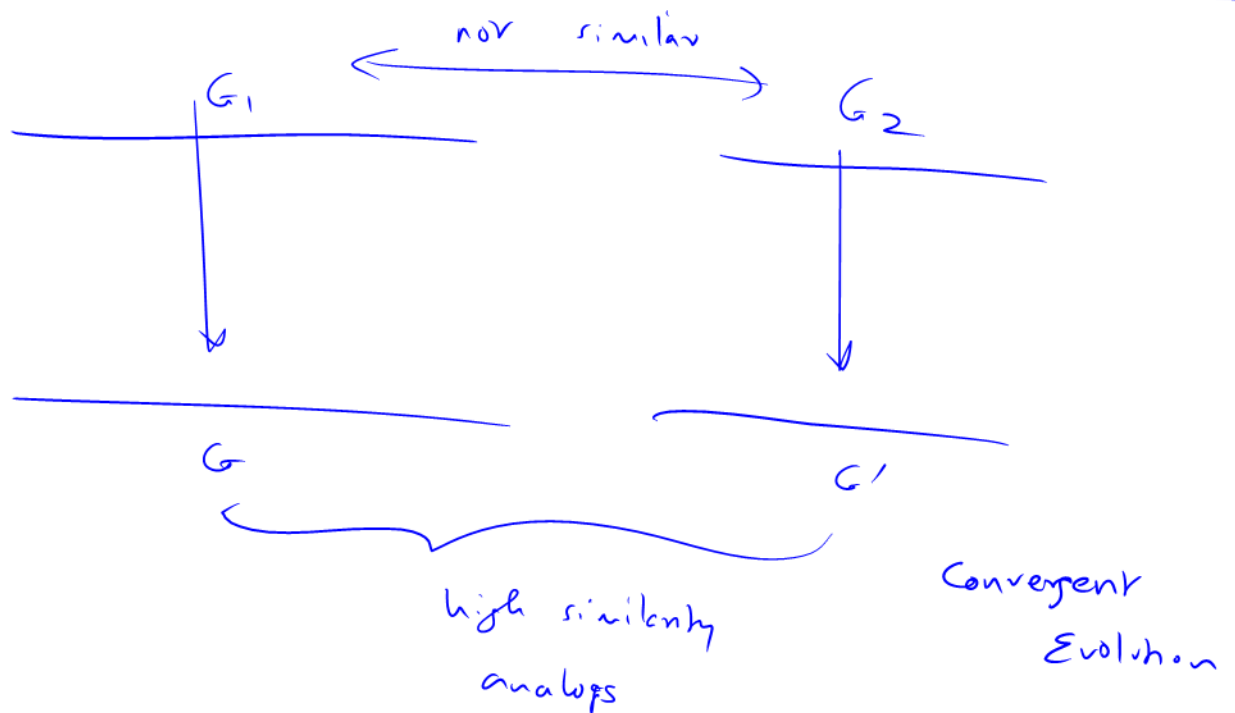
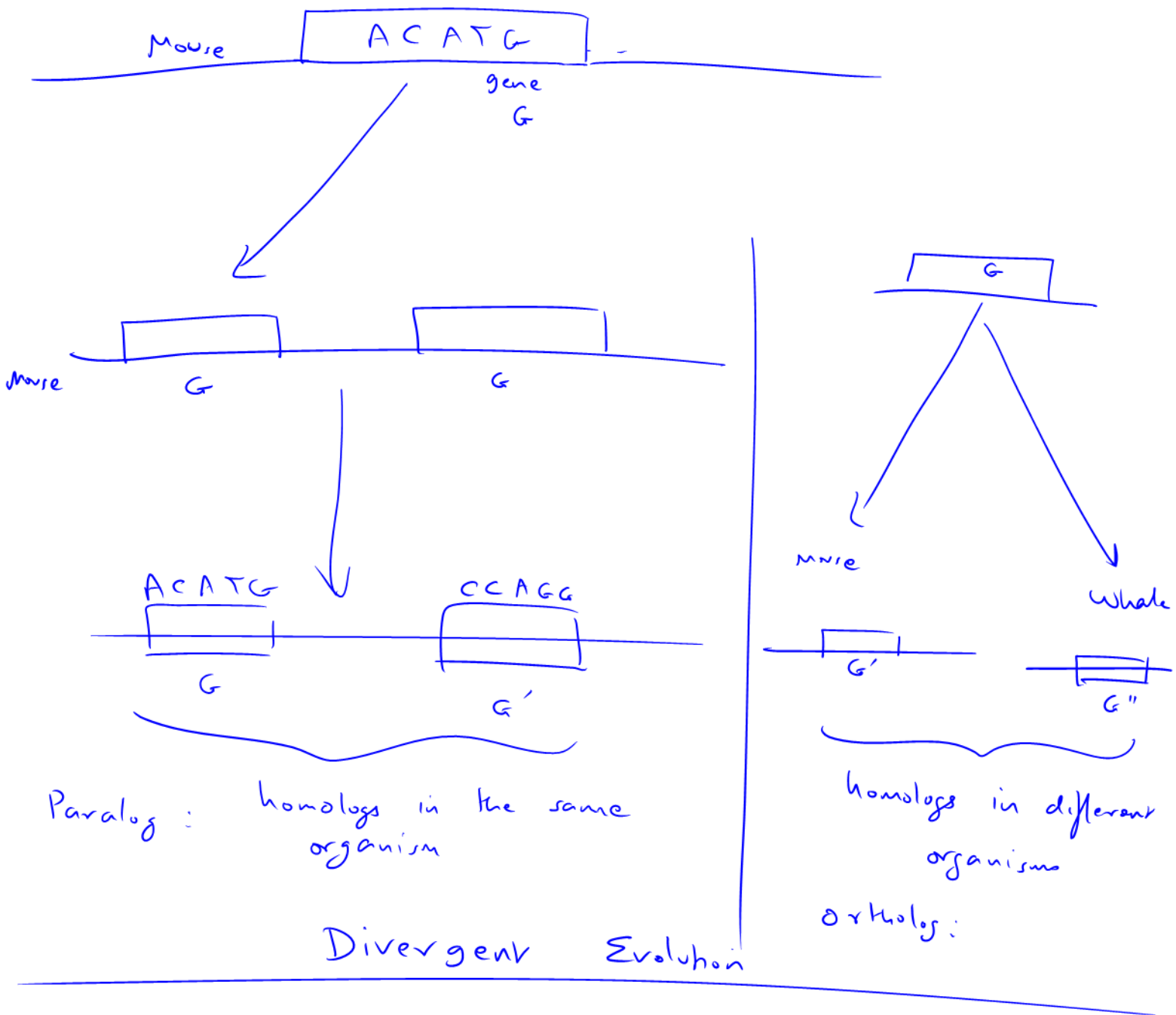
Sum of scores model



gapped alignment

A is an insertion in  $s_1$

A is a deletion in  $s_2$

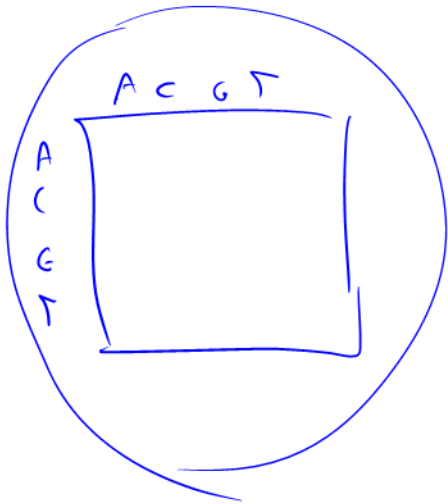


X: A C T A G

Y: A - T G G



$$S = s(A, A) + \underbrace{s(C, -)}_{\substack{\uparrow \\ \text{gap model}}} + s(T, T) + s(A, G) + s(G, G)$$



1) position independence

$$P(A \& B) = P(A) \times P(B)$$

if A & B are independent

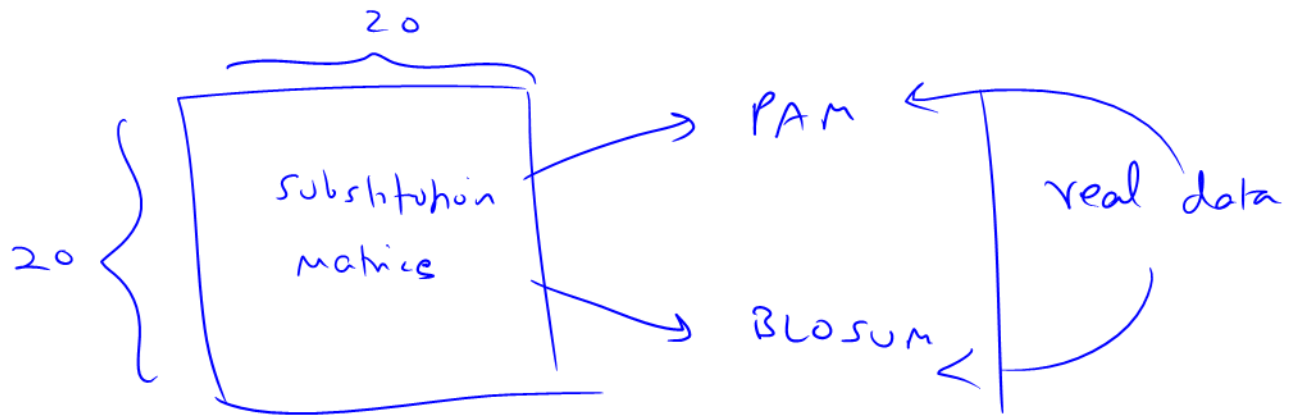
$$\log P(A \& B) = \log P(A) + \log P(B)$$

$$\log P(\text{alignment}) = \log \left( \prod_{i=1}^n \text{OR}(x_i, y_i) \right)$$

↑  
odds ratio

$$S = \sum_{i=1}^n s(x_i, y_i)$$

Compare / align protein sequence



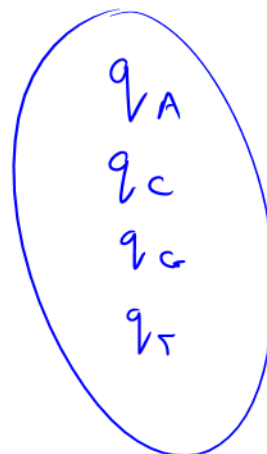
Score of the alignment  
Observe vs. random chance  $\rightarrow$  statistically significant alignment

Ungapped alignment of a given length  $n$ .

$\begin{matrix} & 1 & 2 & 3 & 4 & 5 \\ X: & A & C & A & T & G \\ Y: & C & G & A & G & G \end{matrix}$ 
 $n = 5$

$$P(x, y | R) = \prod_{i=1}^n q_{x_i} \times \prod_{i=1}^n q_{y_i}$$

null model  $\left\{ \begin{array}{l} \text{random model} \end{array} \right.$



background or random probabilities

$$P(x, y | M) = \prod_{i=1}^n p_{x_i y_i} \quad \leftarrow \text{how? PPM Blorum}$$

↑  
alternative  
hypothesis/  
model

Odds ratio:

$$\frac{P(x, y | M)}{P(x, y | R)} = \frac{\prod_{i=1}^n p_{x_i y_i}}{\prod_{i=1}^n q_{x_i} q_{y_i}} \quad \left| \begin{array}{c} \text{Observed} \\ \hline \text{expected} \end{array} \right.$$

$$= \prod_{i=1}^n \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}}$$

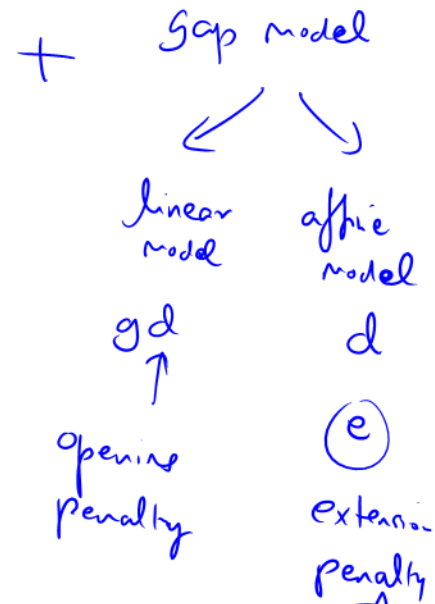
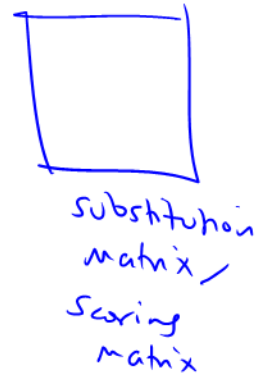
$$S = \log \left( \frac{P(x, y | M)}{P(x, y | R)} \right) = \sum_{i=1}^n \log \left( \frac{p_{x_i y_i}}{q_{x_i} q_{y_i}} \right)$$

$$S = \sum_{i=1}^n s(x_i, y_i)$$

# Alignment Algorithms

X : length n

Y : length m



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

$$g = 4$$

$$d = -10$$

$$d + e$$

$$-40$$

$$d + (g-1)e$$

# of possible alignments?

$$\binom{n+m}{n} \text{ or } \binom{n+m}{m}$$

e.g. if  $m = n$

then  $\binom{2n}{n}$

$$\approx \frac{2^{2n}}{\sqrt{\pi n}}$$

Exponential  
 # of  
 possible  
 alignments

# Dynamic Programming Algorithm

Optimal subproblem property

X: A C A G

Y: G G

	A	C	G	T
A	1		0	
C		1		
G	0		1	
T				1

gap = -1

$F(i, j)$  = optimal score of aligning  $X[1 \dots i]$  with  $Y[1 \dots j]$

then  $F(n, m)$  gives the final score

X: A C A G

Y: G G

$\begin{pmatrix} A C A G \\ G G \end{pmatrix} =$

Case 1

$\begin{pmatrix} A C A \\ G \end{pmatrix} \begin{matrix} i \\ - \\ j \end{matrix}$

$\begin{pmatrix} A C A \\ G G \end{pmatrix} \begin{matrix} G \\ - \\ j \end{matrix}$

$\begin{pmatrix} A C A G \\ G \end{pmatrix} \begin{matrix} i \\ - \\ j \end{matrix}$

$$F(i, j) = \max$$

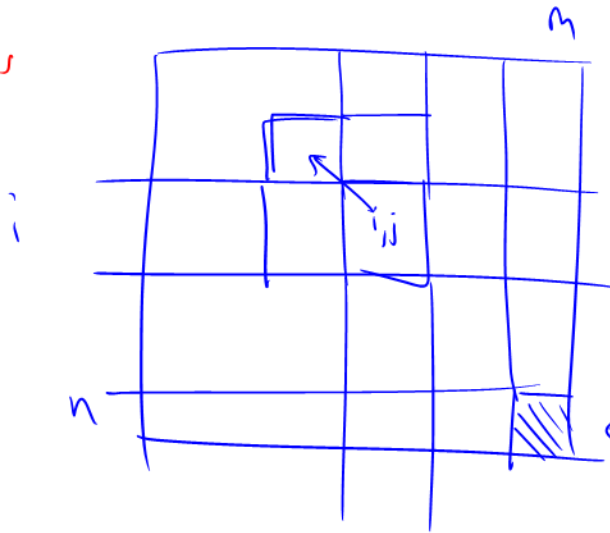
$$F(i-1, j-1) + \underbrace{s(x_i, y_j)}_{\text{case 1}}$$

$$F(i-1, j) + d \quad \leftarrow 2$$

$$F(i, j-1) + d \quad \leftarrow 3$$

$O(nm)$  time algorithm

$n \cdot m$  cells



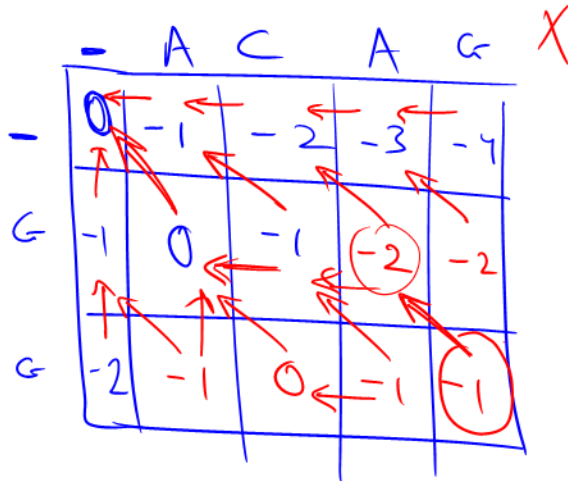
trace back marker

final score

match = 1

mismatch = 0

gap = -1



$$\frac{F}{(n+1) \times (m+1)}$$

size

$$s(A, G) = 0$$

best score = -1

A C A G \_ \_  
\_ \_ \_ \_ G G



One  
optimal  
alignment

← ← ↗ ↗  
A C A G

- - G G

-1 -1 0 1

⏟

-1

↗ ← ← ↗  
A C A G

G - - G

0 -1 -1 1

⏟

-1