

A black and white photograph of a cemetery at night. In the center of the sky, a glowing, saucer-shaped object with a central light source is visible, resembling a UFO. The cemetery is filled with various tombstones and a fence in the foreground. A tall palm tree stands on the right side of the frame. The overall atmosphere is mysterious and eerie.

# *Optimal Sequence Alignment*

*(low-budget production version)*

# Overview

- The alignment problem
- The dynamic programming solution
- Pairwise alignment: exact global and local solutions
- Multiple alignment and the cost of perfection

# Recap: protein scoring

$$\log\left(\frac{C_{a,b}}{f(a)f(b)}\right)$$

*C* matrix – scaled frequencies of change from amino acid *a* to amino acid *b*  
(based on observed changes in some set)

Expectation based solely on frequencies of amino acids (changes not favoured / disfavoured)

Better than random: ratio > 1

Random: ratio = 1

Worse than random: ratio < 1

# Recap: protein scoring

$$D_{a,b} = S \cdot \log \left( \frac{C_{a,b}}{f(a)f(b)} \right)$$

Better than random: ratio > 1  
Random: ratio = 1  
Worse than random: ratio < 1

Magic (why?)

log (why?)

Better than random:  $D_{a,b} > 0$

Random:  $D_{a,b} = 0$

Worse than random:  $D_{a,b} < 0$

# PAM scoring matrix

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

PAM250 matrix (S = 2, log base 2)

# DNA matrix

- Something like this usually works:

	A	G	C	T
A	1	-1	-1	-1
G	-1	1	-1	-1
C	-1	-1	1	-1
T	-1	-1	-1	1

- Or this:

	A	G	C	T
A	1	0.5	-1	-1
G	0.5	1	-1	-1
C	-1	-1	1	0.5
T	-1	-1	0.5	1

# Back to the alignment problem

Given a scoring scheme  $S$

and a set of homologous sequences, uh,  $S$

introduce gaps if necessary to generate an alignment (let's call it  $S$ ) that optimizes the score

# So let's make some alignments!

Sequence  $S_1$ : length  $m$

Sequence  $S_2$ : length  $n$

In total, there are  $\binom{n+m}{m}$  possible alignments of these sequences

$n = m = 2$ :  
 $4!/2!2! = 6$  possibilities

AB--	AB-	AB-	AB	A-B	-AB
--CD	-CD	C-D	CD	-CD	CD-

$n = m = 10$ : 184,756 possible alignments





Alignment of 2 sequences, each 100 amino acids in length:

=  $9.05485147 \times 10^{58}$  possibilities

Brute force is *\*not\** going to work here...

# The Key to Alignment

- If we were given the **midpoint X** within an optimal alignment of  $S_1$  and  $S_2$ , we could **split on X** and solve each problem independently

MEH..K <b>N</b> P..TYL
MDH..K <b>Q</b> P..SYI

MEH..K	+	P..TYL
MDH..K		P..SYI

- But we **don't know** any X, so divide and conquer isn't going to work

# However...

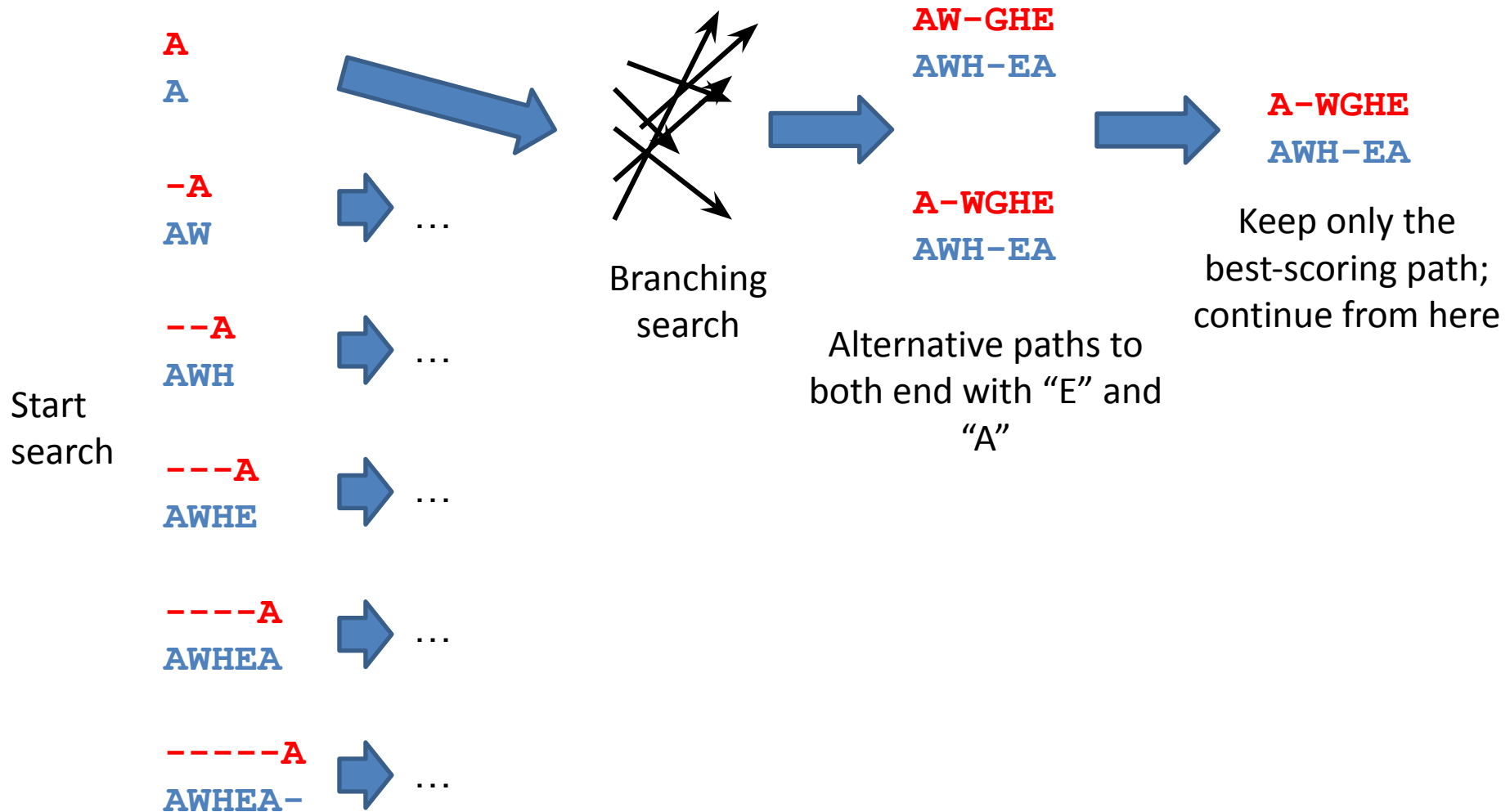
In searching for the best alignment:

- Start at the beginning of the sequences and consider **every possible X**

- BUT -

- Store only the best path (series of matches and gaps) that leads us to X

# Consider an alignment of **AWGHE** vs **AWHEA**:



(these continue as well)

**= Dynamic Programming**

Consider an alignment of AWGHE vs AWHEA:

Sequence 1

Sequence 2

	A	W	G	H	E
A					
W					
H					
E					
A					

# Every possible X

AWGHE vs. AWHEA	A	W	G	H	E
A	Best → (A,A)				
W					
H					
E		Best → (E,W)			
A					Best → (A,E)

# Filling the matrix

We need our substitution matrix  $S$  and gap penalty scheme  $G$

(we'll start with a linear gap penalty  $G = -gd$ )

For each possible  $X$ , consider the three immediate precursors



S = PAM250  
g = 5

AWGHE vs. AWHEA		A	W	G	H	E
	0					
A						
W						
H						
E						
A						

S = PAM250  
g = 5

-AWHEA  
AWGHE

---AWHEA  
AWGHE

-----AWHEA  
AWGHE

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5					
W	-10					
H	-15					
E	-20					
A	-25					

AWHEA  
-AWGHE

AWHEA  
--AWGHE

AWHEA  
---AWGHE

AWHEA  
----AWGHE

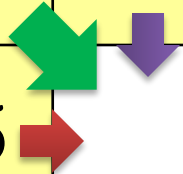
AWHEA  
-----AWGHE

insert gap in AWGHE

insert gap in AWHEA

S = PAM250  
g = 5

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5					
W	-10					
H	-15					
E	-20					
A	-25					



insert gap in AWGHE

insert gap in AWHEA

match

$$S(A,A) = 2$$

Therefore:

Insert -10

Insert -10

Match 2

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5	2				
W	-10					
H	-15					
E	-20					
A	-25					

	A	W	G	H	E
A					
W		F(2,2)	F(2,3)		
H		F(3,2)	F(3,3) = ?		
E					
A					

$$F(3,3) = \max \begin{cases} F(2,2) + S(G,H) \\ F(2,3) - d \\ F(3,2) - d \end{cases}$$

match

insert gap in AWGHE

insert gap in AWHEA

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5	2	-3	-8	-13	-18
W	-10	-3	19	14	9	4
H	-15	-8	14	17	20	15
E	-20	-13	9	14	18	24
A	-25	-18	4	10	13	19

Remember paths INTO  
(not out of)  
each cell

# Global Exact Alignment: Needleman-Wunsch

Since we have retained the best path to each  $F(x,y)$  in the matrix, we can trace back from  $F(m,n)$  to the origin and retrieve the optimal alignment path

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5	2	-3	-8	-13	-18
W	-10	-3	19	14	9	4
H	-15	-8	14	17	20	15
E	-20	-13	9	14	18	24
A	-25	-18	4	10	13	19

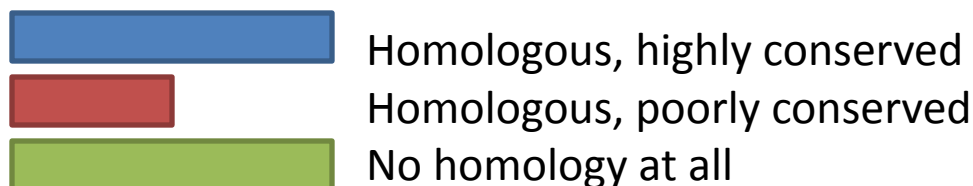


AWGHE-  
AW-HEA

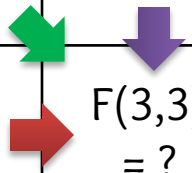
AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5	2	-3	-8	-13	-18
W	-10	-3	19	14	9	4
H	-15	-8	14	17	20	15
E	-20	-13	9	14	18	24
A	-25	-18	4	10	13	19

# Local Exact Alignment: Smith-Waterman

- Only return 'good' sub-alignments of the whole problem
- Useful, for instance, when



	A	W	G	H	E
A					
W		F(2,2)	F(3,2)		
H		F(2,3)	F(3,3) = ?		
E					
A					



$$F(3,3) = \max \left\{ \begin{array}{l} F(2,2) + S(G,H) \\ F(3,2) - d \\ F(2,3) - d \\ 0 \end{array} \right.$$

match

insert gap in AWGHE

insert gap in AWHEA

Nothing is particularly good

This is  
Needleman-Wunsch  
again

AWGHE-  
AW-HEA

AWGHE vs. AWHEA		A	W	G	H	E
	0	-5	-10	-15	-20	-25
A	-5	2	-3	-8	-13	-18
W	-10	-3	19	14	9	4
H	-15	-8	14	17	20	15
E	-20	-13	9	14	18	24
A	-25	-18	4	10	13	19

Slightly modified  
(non-trivial) S-W  
example

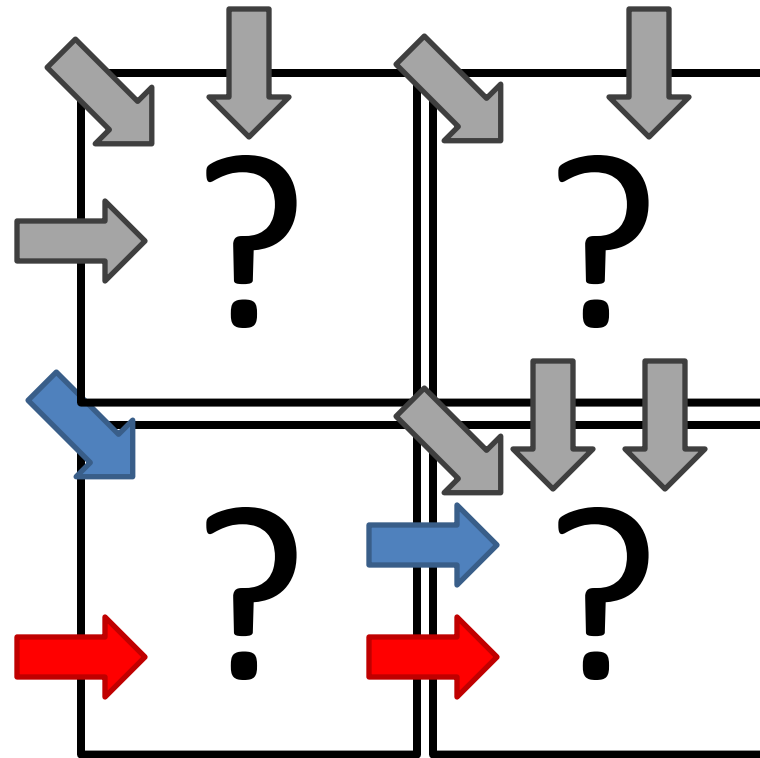
Find the **largest**  
value in the matrix,  
and trace back from  
there to 0

HE

HE

AWGHE vs. AYHEA		A	W	G	H	E
	0	0	0	0	0	0
A	0	2	0	1	0	0
Y	0	0	2	0	0	0
H	0	0	0	0	6	1
E	0	0	0	0	1	10
A	0	2	0	1	0	5

# Affine Gap Penalties



Opening a new gap  
(cost =  $d$ )

Extending a gap  
(cost =  $e$ )

A horizontal move now has two possible costs; we need to consider both alternatives

(and therefore store the best scores for each box given horizontal, vertical, or diagonal entry)

# Significance of S-W Alignments

RANDOMIZE n times

Compute Z-score for each replicate

$$Z(A,B) = \frac{S(A,B) - \tilde{m}}{\tilde{\sigma}}$$

Curve = null model of Z-score fit to Gumbel extreme value distribution

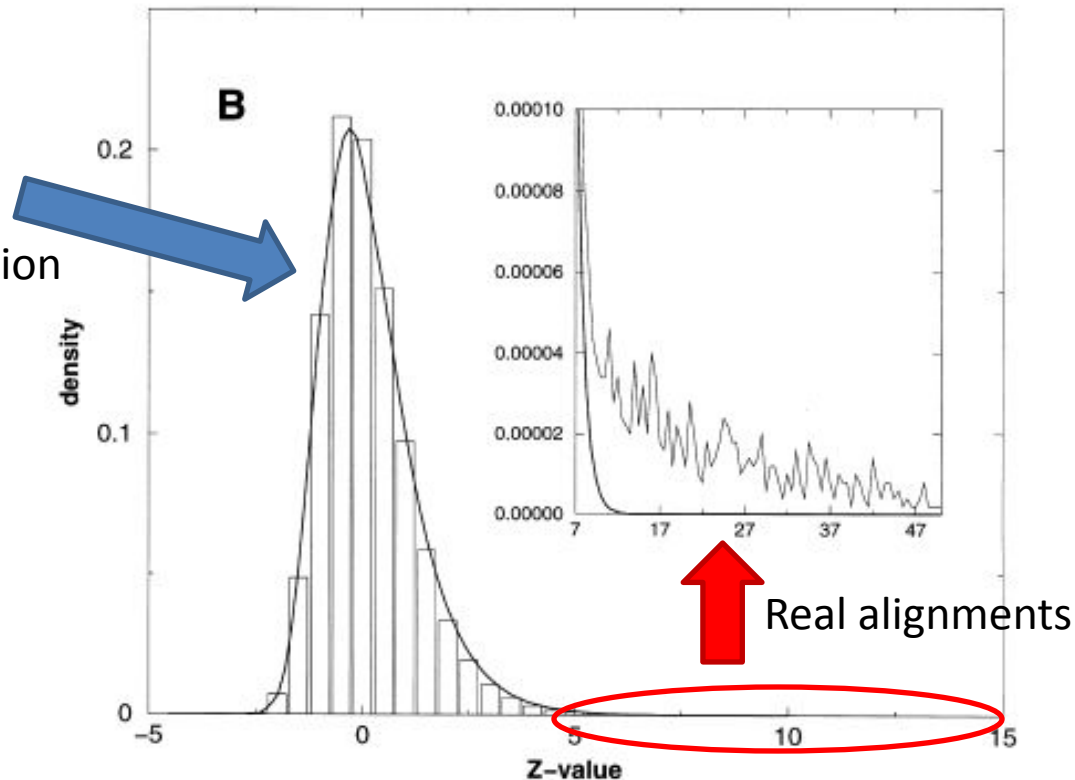


Fig. 6. Distribution of Z-values: (A) empirical distribution (rectangles) and Gumbel model (solid line) for quasi-real sequences. (Insert) the Gumbel model fits the experimental distribution for high Z-values. (B) empirical and Gumbel model for real sequences. (Insert) the Gumbel model (thick line) does not fit the experimental distribution (thin line) for high Z-values.

# Alignment Complexity

- For each possible matching of a residue from sequence  $S_1$  with a residue from  $S_2$ , we need to carry out a constant number of computations and comparisons
- Total =  $3 \times m \times n$
- =  $O(mn)$
- $\sim O(n^2)$  if we assume  $m \cong n$

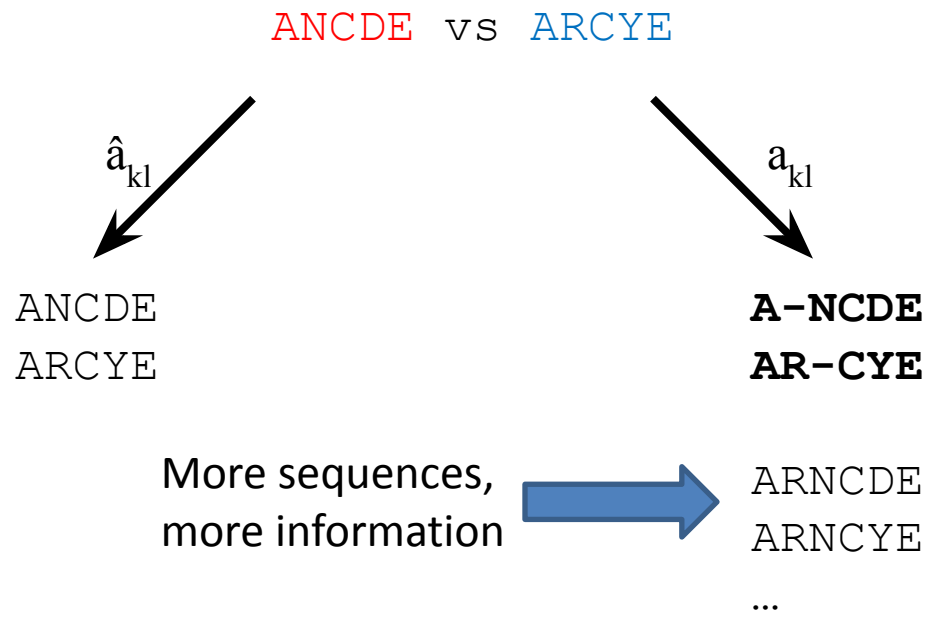


# Multiple Sequence Alignment

- In pairwise alignment, we are optimizing the score between two sequences
- When aligning 3 or more sequences, instead optimize the **sum of pairs** score:

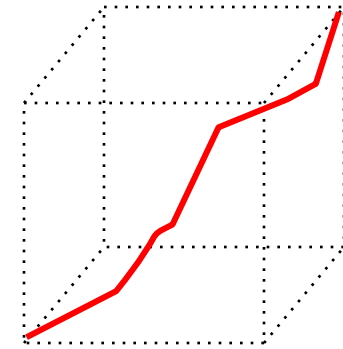
1	N		$2 \times S(N,Q)$
2	Q		$+ 2 \times S(D,Q)$
3	Q	$SP(N,Q,Q,D) =$	$+ S(Q,Q)$
4	D		$+ S(N,D)$

The best alignment between a **pair of sequences** may not appear in the optimal **multiple alignment**



# Multiple Sequence Alignment

- Dynamic programming on  $k$  sequences, each of length  $n$  requires construction of a  $k$ -dimensional matrix with  $n^k$  entries



- =  $O(n^k)$

- Therefore **exponential** in the number of sequences!

# MSA (Carrillo and Lipman, 1988)

- The score of the optimal multiple alignment  $S(a)$  can be no greater than the sum of optimal pairwise alignments  $S(\hat{a}^{kl})$

$$\sum_{k < l} S(a^{kl}) \leq \sum_{k < l} S(\hat{a}^{kl})$$

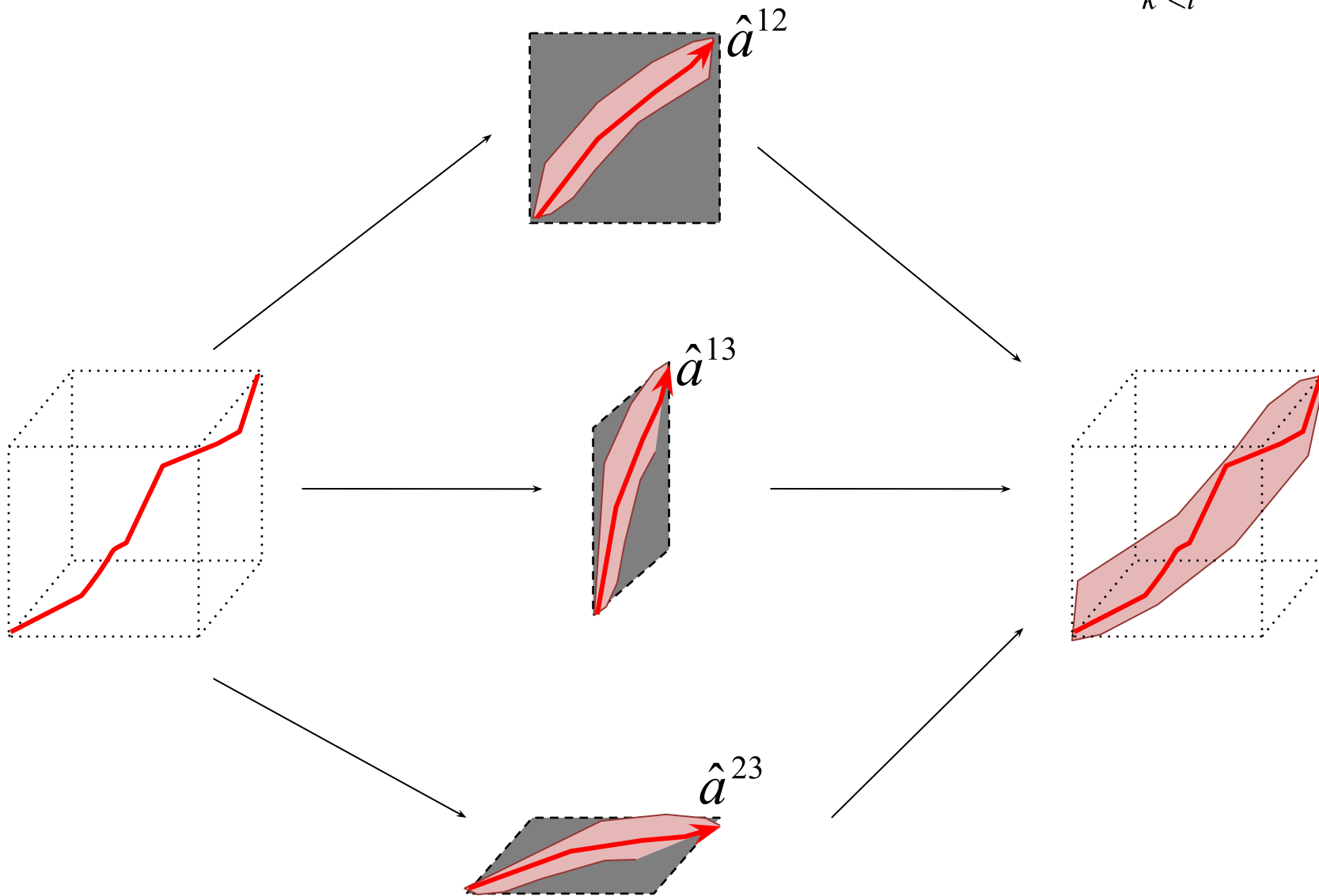
- If we can establish a lower bound  $\sigma$  on the multiple alignment score, then we constrain each  $S(a^{kl})$ :

$$S(\hat{a}^{kl}) - S(a^{kl}) \leq \sum_{k' < l'} S(\hat{a}^{k'l'}) - \sigma$$

Remember: sum of all  
optimal pairwise alignments!

$\sigma$  high:  $S(a^{kl})$  must be close to  $S(\hat{a}^{kl})$

Constrain each pairwise alignment to score no less than  $\sigma + S(\hat{a}^{kl}) - \sum_{k' < l'} S(\hat{a}^{k'l'})$



So we need all optimal pairwise alignments

We also need  $\sigma$ . Where can we find it?

# Types of multiple alignment

## A. Block alignment

```
VRALPDF KGDILRI WNA GMIPVPYV
FVALYDF KGEKLRV WCEA GWVPSNYI
VQALFDF RGDFIHV WWKG GMFPRNYV
VVALYDY KGDEYFI WRA GYIPSNYV
FRAMYDY DGDALIN WMYG GMLPANYV
VKALFDY KSALIQN WWRG LWFP SNYV
YRALYDY LGDILTV WLNQ GDFPGTYV
```

## B. Segment alignment

```
EYVRALPDFNGNDEEDLPFKKGDILRIRDKPEEQ.....WNAEDSEGKR.GMIPVPYVEK
NLFVALYDFVASGDNTLSITKGEKLRVLGYNHNGE.....WCEAQTNGQ..GWVPSNYITP
TYVQALFDFDPQEDGELGFRRGDFIHVMDNSDPN.....WWKGACHGQT..GMFPRNYVTP
KKVVALYDYMPMNANDLQLRKGDEYFILEESNLP.....WWRARDKNGQE.GYIPSNYVTE
KIFRAMYDYMAADADEVSPFKDGDALINVQAIDEG.....WMYGTVQRTGRTGMLPANYVEA
CAVKALFDYKAQREDELTFIKSALIQNVEKQEGG.....WWRGDYGGKKQ.LWFP SNYVEE
YQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSQGQEARPEEIGWLNQYNETTGERGDFPGTYVEY
```

## C. Local alignment

```
.....aeyVRALPDFngndeedlpfkKGDILRIrdkpeeQ.....WNAedsegkr.GMIPVPYVek.....
.....nlFVALYDFvasgdntlsitKGEKLRVLgynhngE.....WCEAqtknqQ..GWVPSNYItpvns.....
lvdyhrstsvsrnqqiflrldieqvpqqptyVQALFDFdpqedgelgfrRGDFIHVmdnsdpn.....WWKGachgqt..GMFPRNYVtpvnrnv.....
.....gsmstselkkVVALYDYmpmnandlqlrKGDEYFIleesnlP.....WWRARDkngqe.GYIPSNYVteaeds.....
.....tagkiFRAMYDYmaadaevsfkDGDALINVQAIDEG.....WMYgtvqrtgrtGMLPANYVeai.....
.....gsptfkcaVKALFDYkaqredeltfiKSALIQNVEKQEGG.....WWRGDyggkkq.LWFP SNYVeemvnpegihrd
.....gyqYRALYDYkkereedidlhlGDILTVnkgslvalgfsdqgearpeeigWLNQynettgerGDFPGTYVeyigrkkisp..
```

## D. Global alignment

```
.....AEYVRALPDFNGNDEEDLPFKKGDILRIRDKP.....EEQWNAEDS.EGKRGMIPVPYVEK.....
.....NLFVALYDFVASGDNTLSITKGEKLRVLGYN.....HNGEWCEAQTK..NGQGWVPSNYITPVNS.....
LVDYHRSTSVSRNQQIFLRDIEQVPQQPTYVQALFDFDPQEDGELGFRRGDFIHVMDNS.....DPNWWKGACH..GQTGMFPRNYVTPVNRNV.....
.....GSMSTSELKKVVALYDYMPMNANDLQLRKGDEYFILEES.....NLPWWRARDK.NGQEGYIPSNYVTEAEDS.....
.....TAGKIFRAMYDYMAADADEVSPFKDGDALINVQAI.....DEGWMYGTVQRTGRTGMLPANYVEAI.....
.....GSPTFKCAVKALFDYKAQREDELTFIKSALIQNVEKQ.....EGGWWRGDY.GKKQLWFP SNYVEEMVNPEGIHRD
.....GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSQGQEARPEEIGWLNQYNETTGERGDFPGTYVVEYIGRKKISP..
```

From Lecompte et al. (2001) *Gene*



# Summary

- Dynamic programming allows the calculation of optimal pairwise alignments (for a given scoring scheme!)
- As soon as we go from 2 to  $>2$  sequences, the exponential time complexity of the algorithm makes it impractical
- Need heuristics!