

---

# Sequence Alignment

---

---

# Outline

- Global Alignment
  - Scoring Matrices
  - Local Alignment
  - Alignment with Affine Gap Penalties
-

## From LCS to Alignment: Change up the Scoring

- The Longest Common Subsequence (LCS) problem —the simplest form of sequence alignment – allows only insertions and deletions (no mismatches).
- In the LCS Problem, we scored 1 for matches and 0 for indels
- Consider penalizing indels and mismatches with negative scores
- Simplest *scoring schema*:
  - +1 : match premium
  - $\mu$  : mismatch penalty
  - $\sigma$  : indel penalty

# Simple Scoring

- When mismatches are penalized by  $-\mu$ , indels are penalized by  $-\sigma$ , and matches are rewarded with  $+1$ , the resulting score is:

$$\#matches - \mu(\#mismatches) - \sigma(\#indels)$$

# The Global Alignment Problem

Find the best alignment between two strings under a given scoring schema

Input : Strings  $\mathbf{v}$  and  $\mathbf{w}$  and a scoring schema

Output : Alignment of maximum score

$\uparrow \rightarrow = -\sigma$

$= 1$  if match

$\swarrow \left\{ \begin{array}{l} = -\mu \text{ if mismatch} \end{array} \right.$

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + 1 & \text{if } v_i = w_j \\ s_{i-1,j-1} - \mu & \text{if } v_i \neq w_j \\ s_{i-1,j} - \sigma \\ s_{i,j-1} - \sigma \end{cases}$$

$m$  : mismatch

penalty

$\sigma$  : indel penalty

# Scoring Matrices

To generalize scoring, consider a  $(4+1) \times (4+1)$  **scoring matrix**  $\delta$ .

In the case of an amino acid sequence alignment, the scoring matrix would be a  $(20+1) \times (20+1)$  size. The addition of 1 is to include the score for comparison of a gap character “-”.

This will simplify the algorithm as follows:

$$s_{i,j} = \max \begin{array}{l} s_{i-1,j-1} + \delta(v_i, w_j) \\ s_{i-1,j} + \delta(v_i, -) \\ s_{i,j-1} + \delta(-, w_j) \end{array}$$

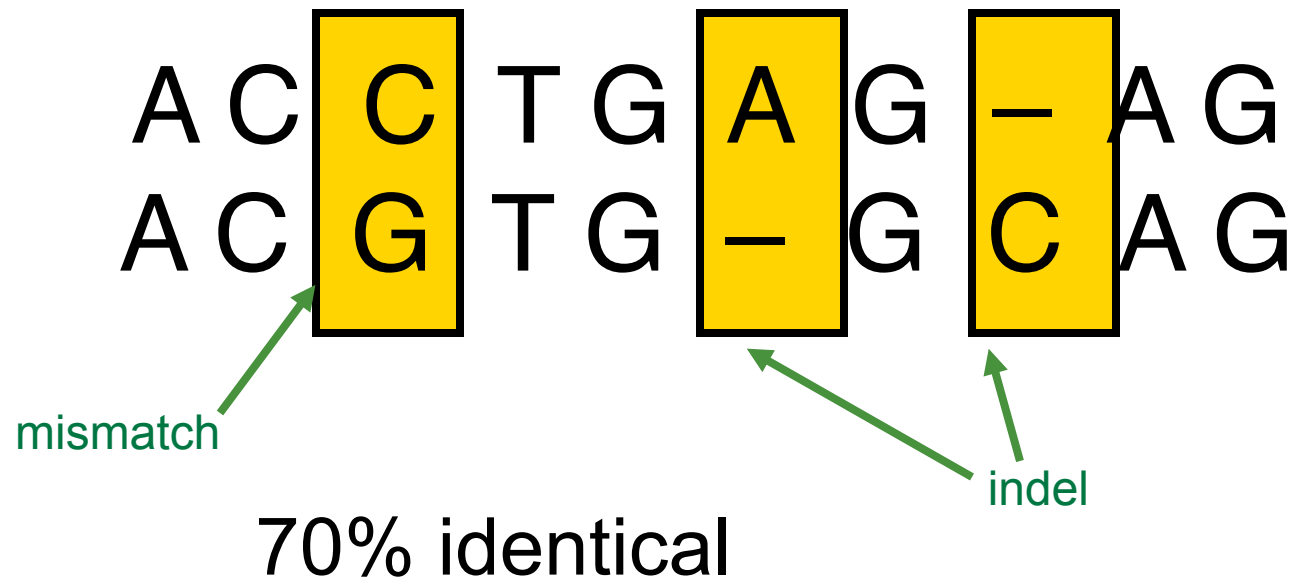
---

# Measuring Similarity

- Measuring the extent of similarity between two sequences
    - Based on percent sequence identity
    - Based on conservation
-

# Percent Sequence Identity

- The extent to which two nucleotide or amino acid sequences are invariant





# Making a Scoring Matrix

- Scoring matrices are created based on biological evidence.
- Alignments can be thought of as two sequences that differ due to mutations.
- Some of these mutations have little effect on the protein's function, therefore some penalties,  $\delta(v_i, w_j)$ , will be less harsh than others.

# Scoring Matrix: Example

	A	R	N	K
A	5	-2	-1	-1
R	-	7	-1	3
N	-	-	7	0
K	-	-	-	6

AKRANR  
KAAANK

- Notice that although R and K are different amino acids, they have a positive score.
- Why? They are both positively charged amino acidsà will not greatly change function of protein.

~~-1 + (-1) +~~

~~(-2) + 5 + 7 +~~

# Conservation

- Amino acid changes that tend to preserve the physico-chemical properties of the original residue
  - Polar to polar
    - aspartate à glutamate
  - Nonpolar to nonpolar
    - alanine à valine
  - Similarly behaving residues
    - leucine to isoleucine

# Scoring matrices

- Amino acid substitution matrices
  - PAM
  - BLOSUM
- DNA substitution matrices
  - DNA is less conserved than protein sequences
  - Less effective to compare coding regions at nucleotide level

# PAM

- **P**oint **A**ccepted **M**utation (Dayhoff et al.)
- 1 PAM = PAM<sub>1</sub> = 1% average change of all amino acid positions
  - After 100 PAMs of evolution, not every residue will have changed
    - some residues may have mutated several times
    - some residues may have returned to their original state
    - some residues may not changed at all

# PAM<sub>X</sub>

- $\text{PAM}_X = \text{PAM}_1^X$
- $\text{PAM}_{250} = \text{PAM}_1^{250}$
- $\text{PAM}_{250}$  is a widely used scoring matrix:

	Ala	Arg	Asn	Asp	Cys	Gln	Glu	Gly	His	Ile	Leu	Lys	...
	A	R	N	D	C	Q	E	G	H	I	L	K	...
Ala A	13	6	9	9	5	8	9	12	6	8	6	7	...
Arg R	3	17	4	3	2	5	3	2	6	3	2	9	
Asn N	4	4	6	7	2	5	6	4	6	3	2	5	
Asp D	5	4	8	11	1	7	10	5	6	3	2	5	
Cys C	2	1	1	1	52	1	1	2	2	2	1	1	
Gln Q	3	5	5	6	1	10	7	3	7	2	3	5	
...													
Trp W	0	2	0	0	0	0	0	0	1	0	1	0	
Tyr Y	1	1	2	1	3	1	1	1	3	2	2	1	
Val V	7	4	4	4	4	4	4	4	5	4	15	10	

---

# BLOSUM

- **B**locks **S**ubstitution **M**atrix
  - Scores derived from *observations* of the frequencies of substitutions in blocks of local alignments in related proteins
  - Matrix name indicates evolutionary distance
    - BLOSUM62 was created using sequences sharing no more than 62% identity
-





# Local vs. Global Alignment

- The Global Alignment Problem tries to find the longest path between vertices  $(0,0)$  and  $(n,m)$  in the edit graph.
- The Local Alignment Problem tries to find the longest path among paths between **arbitrary vertices**  $(i,j)$  and  $(i',j')$  in the edit graph.

# Local vs. Global Alignment

- The Global Alignment Problem tries to find the longest path between vertices  $(0,0)$  and  $(n,m)$  in the edit graph.
- The Local Alignment Problem tries to find the longest path among paths between **arbitrary vertices**  $(i,j)$  and  $(i',j')$  in the edit graph.
- In the edit graph with negatively-scored edges, Local Alignment may score higher than Global Alignment

# Local vs. Global Alignment (cont'd)

- **Global Alignment**

```

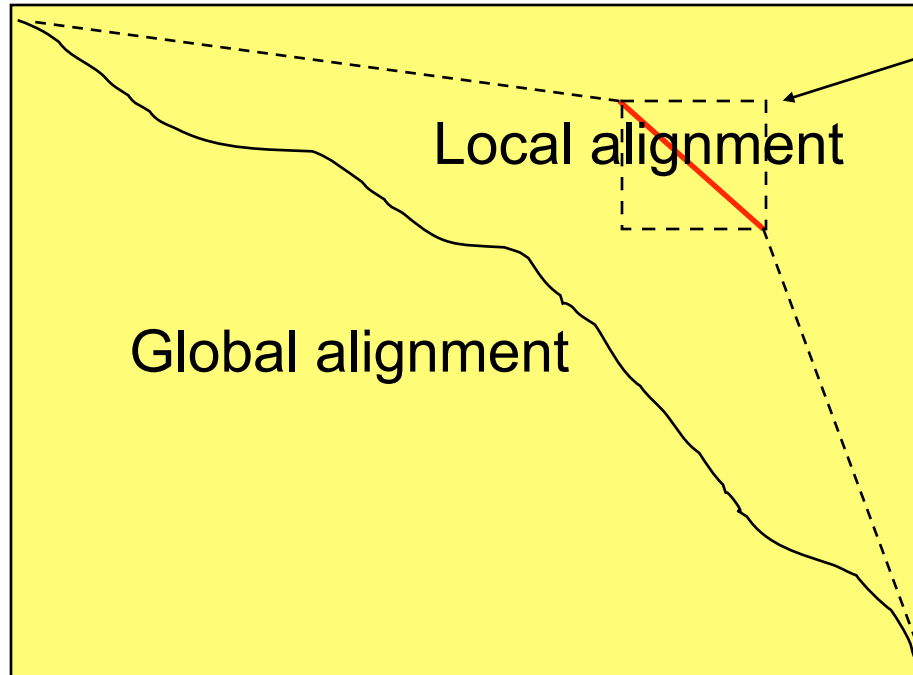
--T--CC-C-AGT--TATGT-CAGGGGACACG-A-GCATGCAGA-GAC
  |  ||  |  ||  |  |  |||  |  |  |  |  |  |||  |
AATTGCCGCC-GTCGT-T-TTCAG----CA-GTTATG-T-CAGAT--C
  
```

- **Local Alignment—better alignment to find conserved segment**

```

                                tccCAGTTATGTCAGgggacacgagcatgcagagac
                                |||
aattgccgccgctcgtttttcagCAGTTATGTCAGatc
  
```

# Local Alignment: Example



Compute a "mini"  
Global Alignment to  
get Local

# Local Alignments: Why?

- Two genes in different species may be similar over short conserved regions and dissimilar over remaining regions.
- Example:
  - Homeobox genes have a short region called the *homeodomain* that is highly conserved between species.
  - A global alignment would not find the homeodomain because it would try to align the ENTIRE sequence

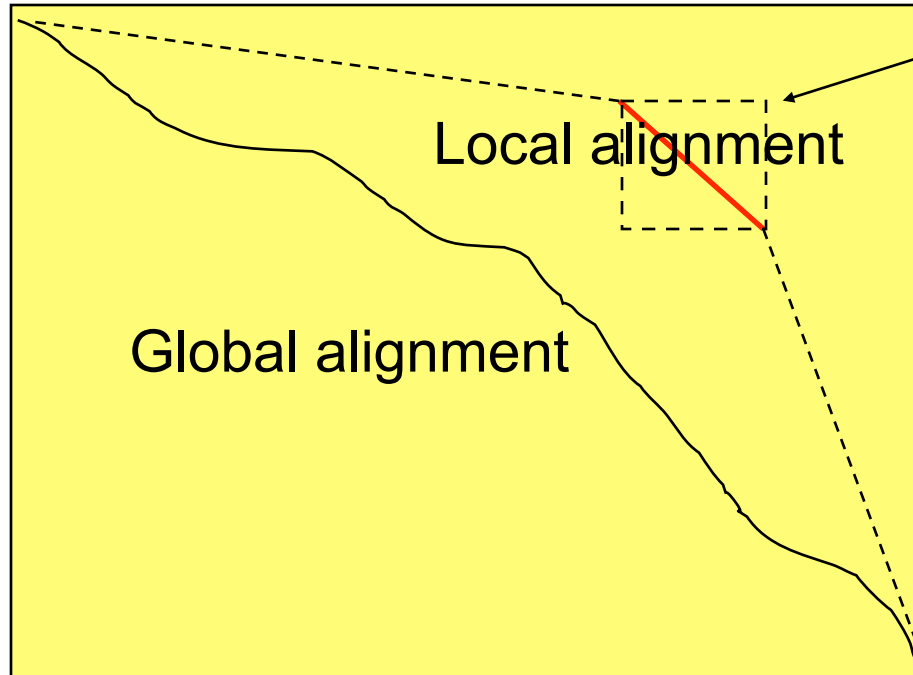
# The Local Alignment Problem

- Goal: Find the best local alignment between two strings
- Input : Strings  $\mathbf{v}$ ,  $\mathbf{w}$  and scoring matrix  $\delta$
- Output : Alignment of substrings of  $\mathbf{v}$  and  $\mathbf{w}$  whose alignment score is maximum among all possible alignment of all possible substrings

# The Problem with this Problem

- Long run time  $O(n^4)$ :
  - In the grid of size  $n \times n$  there are  $\sim n^2$  vertices  $(i,j)$  that may serve as a source.
  - For each such vertex computing alignments from  $(i,j)$  to  $(i',j')$  takes  $O(n^2)$  time.
- ~~This can be remedied by giving free rides~~

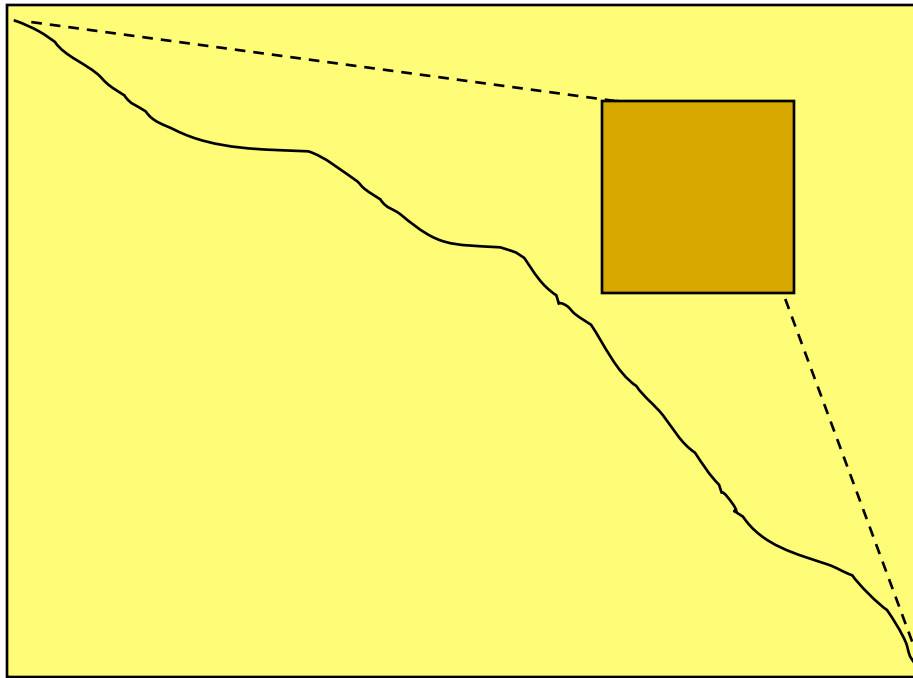
# Local Alignment: Example



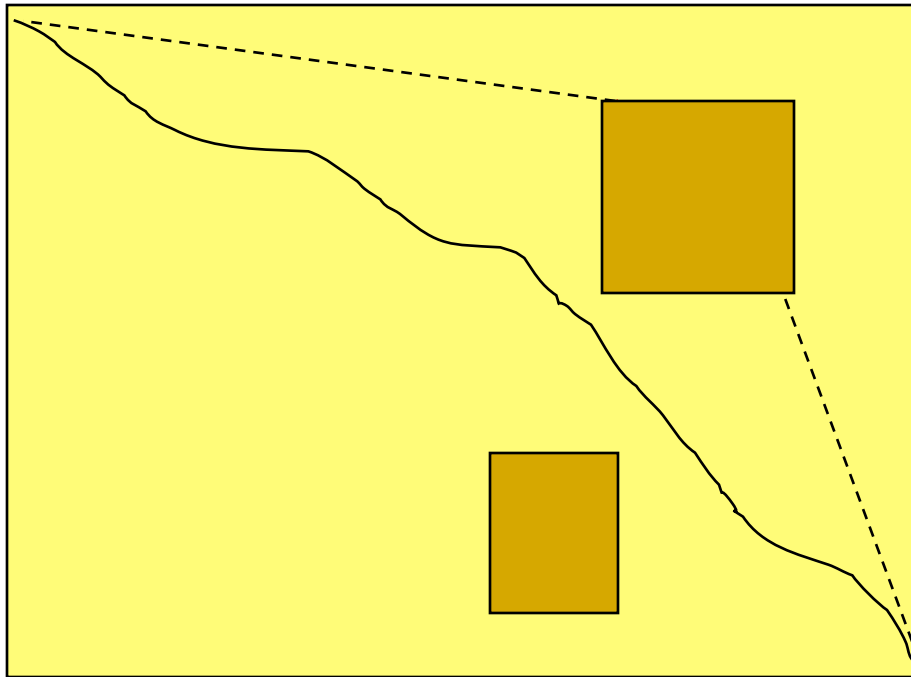
Compute a "mini"  
Global Alignment to  
get Local



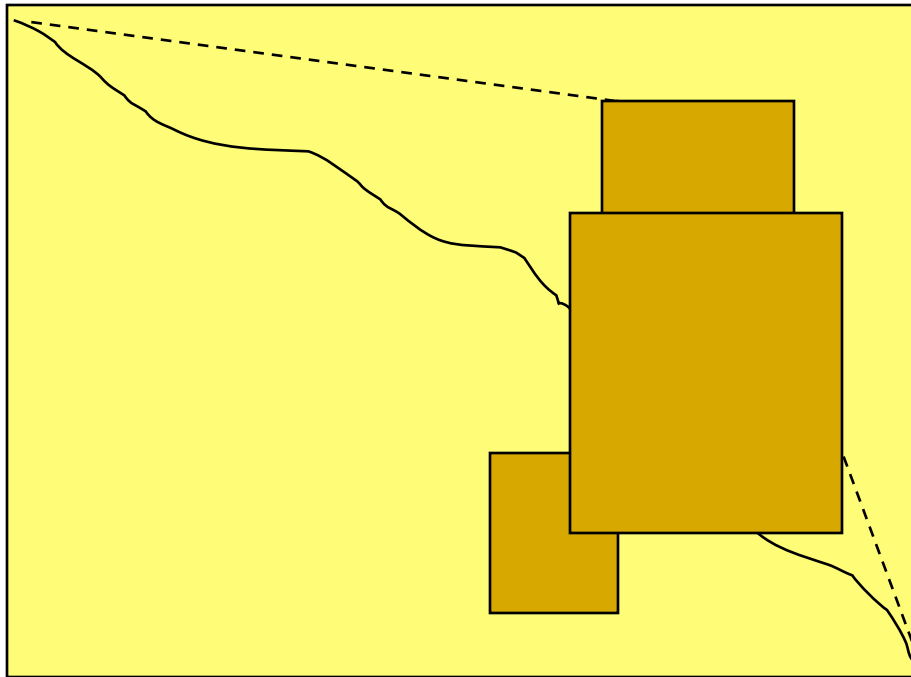
# Local Alignment: Example



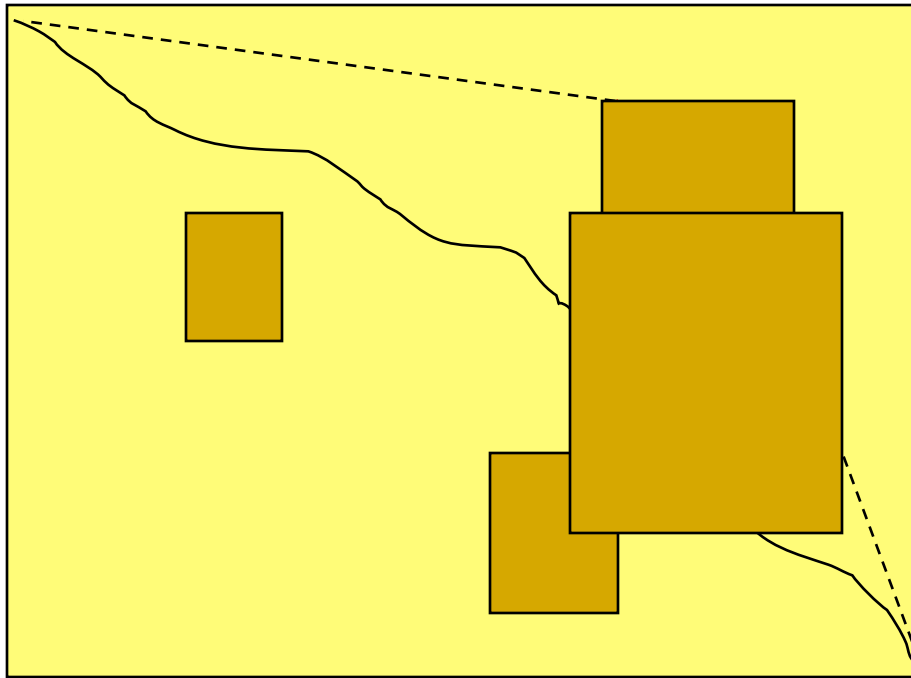
# Local Alignment: Example



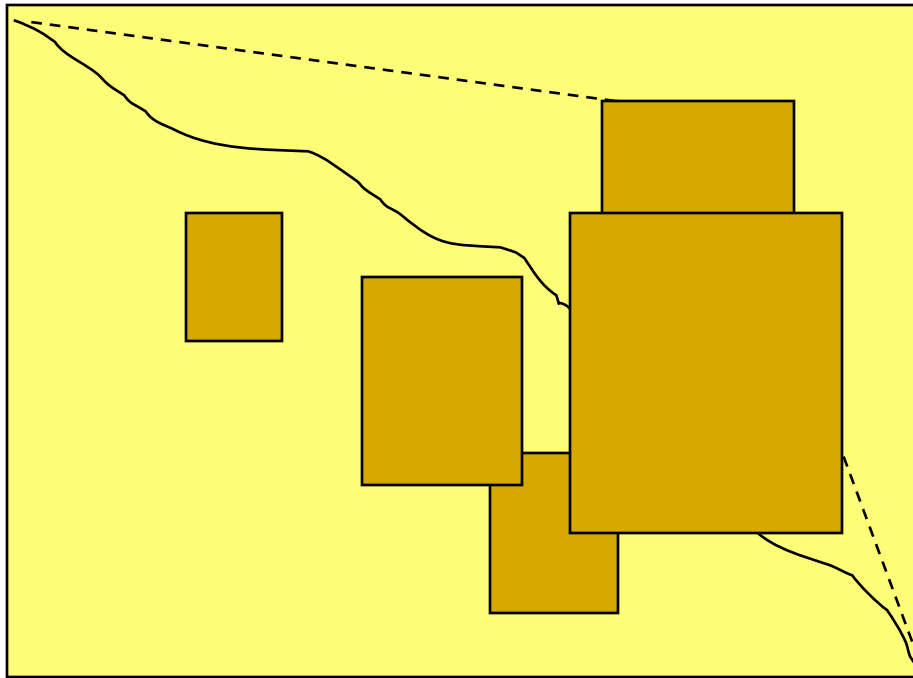
# Local Alignment: Example



# Local Alignment: Example

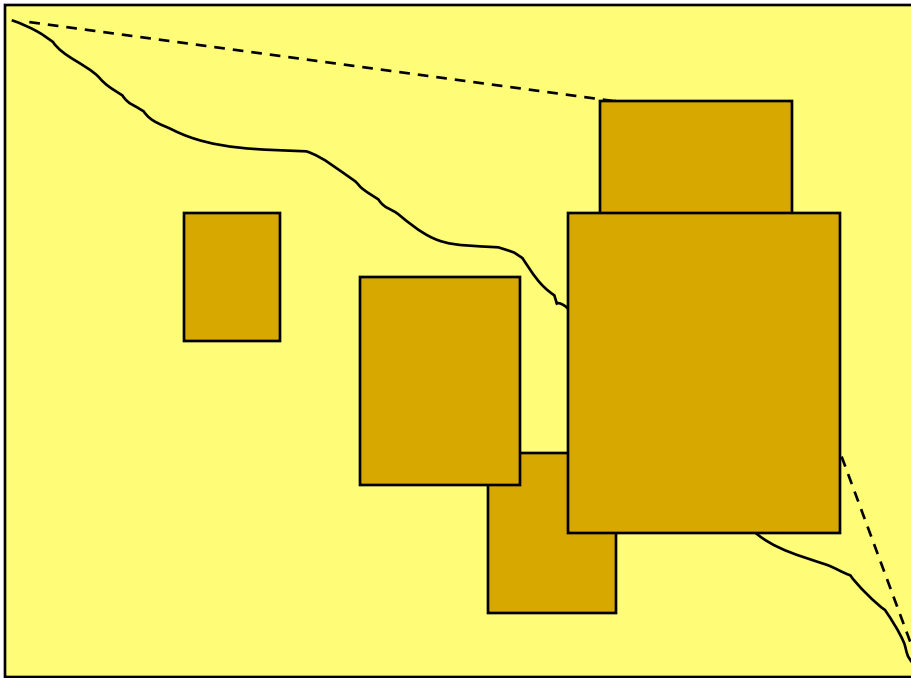


# Local Alignment: Example

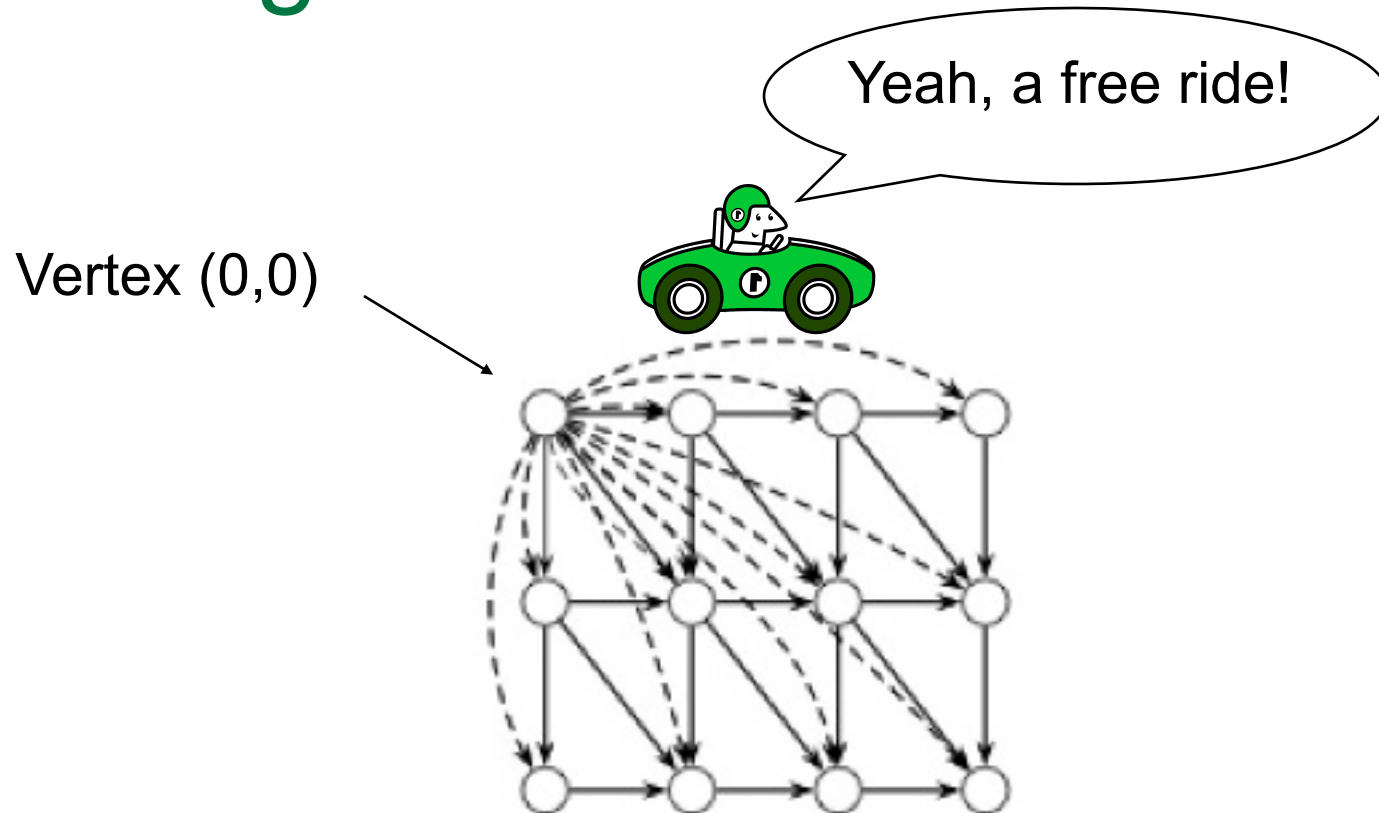


# Local Alignment: Running Time

- Long run time  $O(n^4)$ :
  - In the grid of size  $n \times n$  there are  $\sim n^2$  vertices  $(i,j)$  that may serve as a source.
  - For each such vertex computing alignments from  $(i,j)$  to  $(i',j')$  takes  $O(n^2)$  time.
- This can be remedied by giving free rides



# Local Alignment: Free Rides



The dashed edges represent the free rides from  $(0,0)$  to every other node.

# The Local Alignment Recurrence

- The largest value of  $s_{i,j}$  over the whole edit graph is the score of the best local alignment.
- The recurrence:

$$s_{i,j} = \max \begin{cases} 0 \\ s_{i-1,j-1} + (v_i, w_j) \\ s_{i-1,j} + (v_i, -) \\ s_{i,j-1} + (-, w_j) \end{cases}$$

Notice there is only this change from the original recurrence of a Global Alignment



# The Local Alignment Recurrence

- The largest value of  $s_{i,j}$  over the whole edit graph is the score of the best local alignment.

- The recurrence:

$$s_{i,j} = \max \begin{cases} 0 \\ s_{i-1,j-1} + (v_i, w_j) \\ s_{i-1,j} + (v_i, -) \\ s_{i,j-1} + (-, w_j) \end{cases}$$

**Power of ZERO:** there is only this change from the original recurrence of a Global Alignment - since there is only one “free ride” edge entering into every vertex

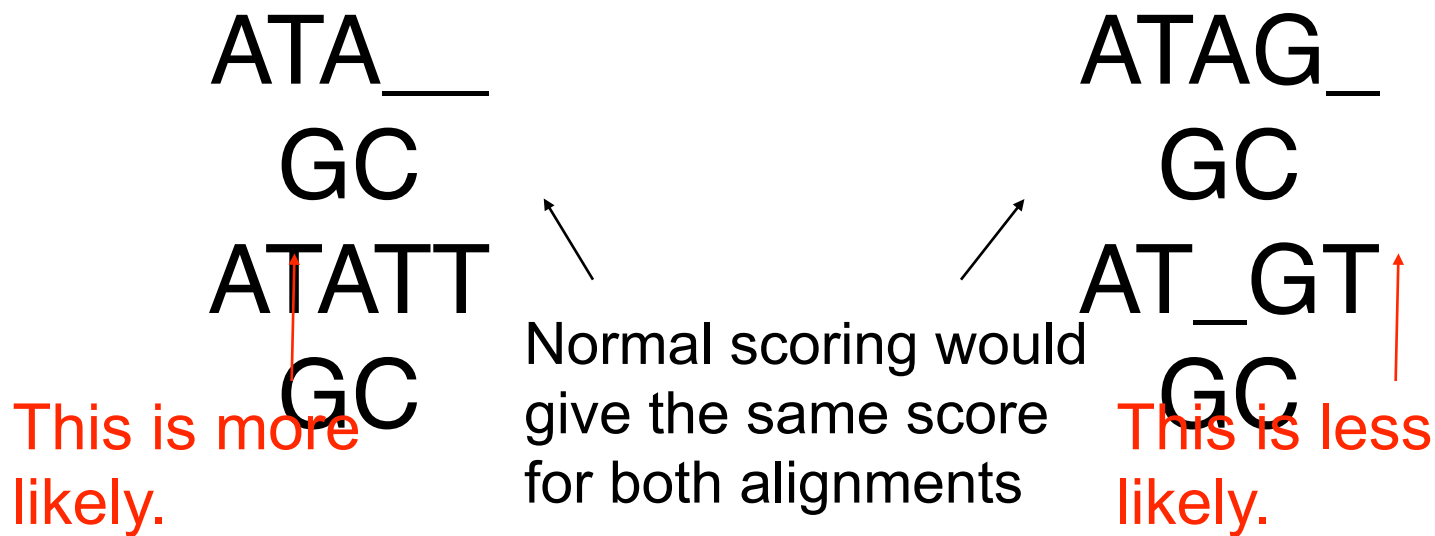
# Scoring Indels: Naive Approach

- A fixed penalty  $\sigma$  is given to every indel:
  - $-\sigma$  for 1 indel,
  - $-2\sigma$  for 2 consecutive indels
  - $-3\sigma$  for 3 consecutive indels, etc.

Can be too severe penalty for a series of 100 consecutive indels

# Affine Gap Penalties

- In nature, a series of  $k$  indels often come as a single event rather than a series of  $k$  single nucleotide events:



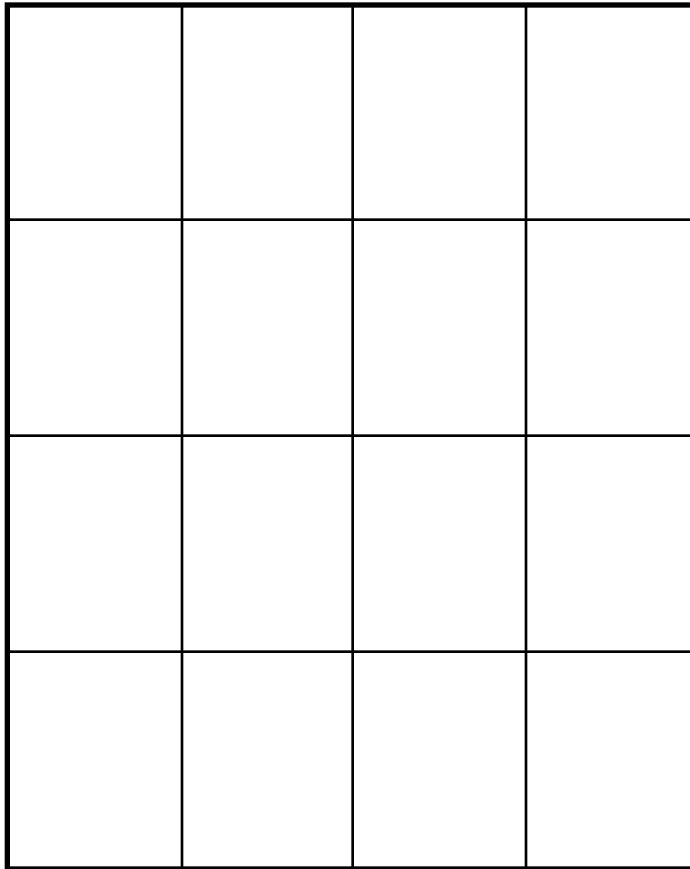
# Accounting for Gaps

- *Gaps*- contiguous sequence of spaces in one of the rows
- Score for a gap of length  $x$  is:  
$$-(\rho + \sigma x)$$
where  $\rho > 0$  is the penalty for introducing a gap:  
**gap opening penalty**  
 $\rho$  will be large relative to  $\sigma$ :  
**gap extension penalty**  
because you do not want to add too much of a penalty for extending the gap.

# Affine Gap Penalties

- Gap penalties:
  - $-\rho - \sigma$  when there is 1 indel
  - $-\rho - 2\sigma$  when there are 2 indels
  - $-\rho - 3\sigma$  when there are 3 indels, etc.
  - $-\rho - x \cdot \sigma$  (-gap opening -  $x$  gap extensions)
- Somehow reduced penalties (as compared to naïve scoring) are given to runs of horizontal and vertical edges

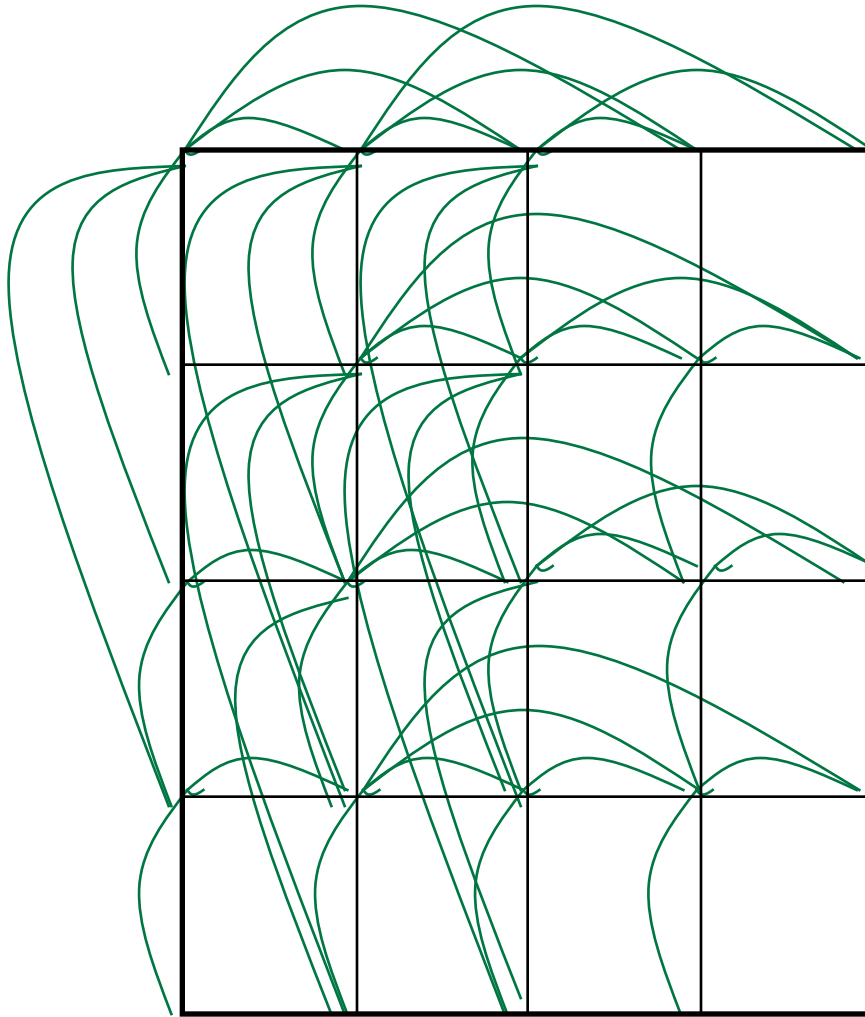
# Affine Gap Penalties and Edit Graph



To reflect affine gap penalties we have to add “long” horizontal and vertical edges to the edit graph. Each such edge of length  $x$  should have weight

$$-r - x * s$$

## Adding “Affine Penalty” Edges to the Edit Graph

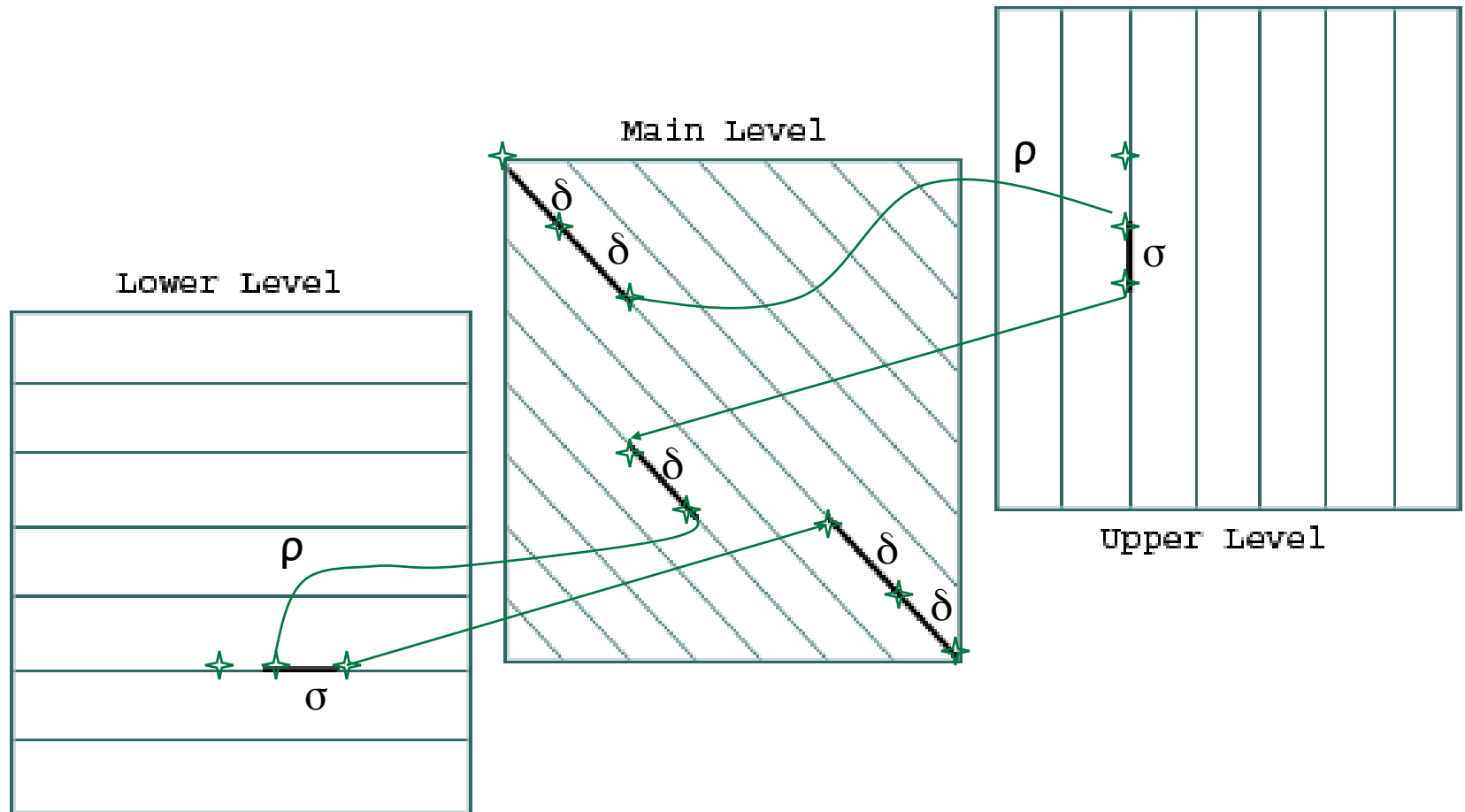


There are many such edges!

Adding them to the graph increases the running time of the alignment algorithm by a factor of  $n$  (where  $n$  is the number of vertices)

So the complexity increases from  $O(n^2)$  to  $O(n^3)$

# Manhattan in 3 Layers





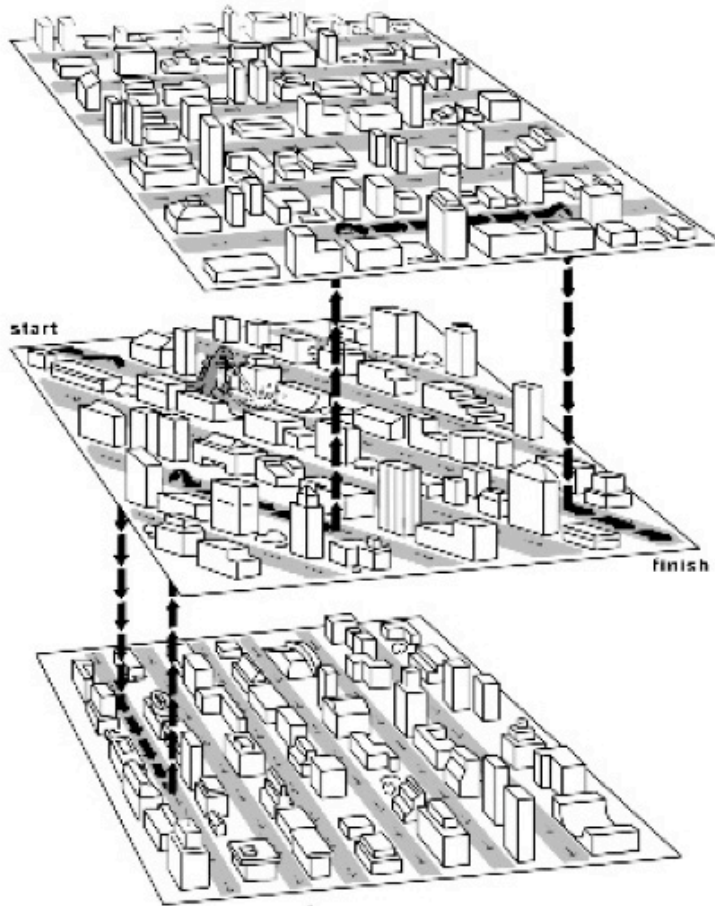
## Affine Gap Penalties and 3 Layer Manhattan Grid

- The three recurrences for the scoring algorithm creates a 3-layered graph.
  - The top level creates/extends gaps in the sequence  $w$ .
  - The bottom level creates/extends gaps in sequence  $v$ .
  - The middle level extends matches and mismatches.
-

# Switching between 3 Layers

- Levels:
    - The **main level** is for diagonal edges
    - The **lower level** is for horizontal edges
    - The **upper level** is for vertical edges
  - A jumping penalty is assigned to moving from the main level to either the upper level or the lower level ( $-r - s$ )
  - There is a gap extension penalty for each continuation on a level other than the main level ( $-s$ )
-

# The 3-leveled Manhattan Grid



Gaps in  $w$

Matches/  
Mismatch

es  
Gaps in  $v$

# Affine Gap Penalty Recurrences

$$\downarrow s_{i,j} = \max \begin{cases} \downarrow s_{i-1,j} - \sigma & \text{Continue Gap in } w \text{ (deletion)} \\ s_{i-1,j} - (\rho + \sigma) & \text{Start Gap in } w \text{ (deletion): from middle} \end{cases}$$

$$\rightarrow s_{i,j} = \max \begin{cases} \rightarrow s_{i,j-1} - \sigma & \text{Continue Gap in } v \text{ (insertion)} \\ s_{i,j-1} - (\rho + \sigma) & \text{Start Gap in } v \text{ (insertion): from middle} \end{cases}$$

$$s_{i,j} = \max \begin{cases} \downarrow s_{i-1,j-1} + \delta(v_i, w_j) & \text{Match or Mismatch} \\ \rightarrow \bar{s}_{i,j} & \text{End deletion: from top} \\ s_{i,j} & \text{End insertion: from bottom} \end{cases}$$