#### Next

#### Local and global sequence alignment

# Local vs. Global Alignment

• The <u>Global Alignment Problem</u> 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 <u>Global Alignment Problem</u> 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-CAGGGGGACACG-A-GCATGCAGA-GAC

 Local Alignment—better alignment to find conserved segment tccCAGTTATGTCAGgggacacgagcatgcagagac

aattgccgccgtcgttttcagCAGTTATGTCAGatc



EECS 4425, Fall 2018

# 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**

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

#### The Problem with this Problem

• Long run time O(n<sup>4</sup>):

- 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



EECS 4425, Fall 2018

















# Local Alignment: Running Time



• Long run time O(n<sup>4</sup>):

- In the grid of size *n* x *n* there are ~*n*<sup>2</sup> 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



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

EECS 4425, Fall 2018

#### **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:

$$s_{i,j} = max \begin{cases} 0 \\ s_{i-1,j-1} + \delta(V_i, W_j) \\ s_{i-1,j} + \delta(V_i, -) \\ s_{i,j-1} + \delta(-, W_j) \end{cases}$$

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:

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

**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:



20

# **Accounting for Gaps**

- Gaps- contiguous sequence of spaces in one of the rows

#### **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

-ho - X \* $\sigma$ 

#### 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  $O(n^3)$ 

EECS 4425, Fall 2018

#### 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 (- $\rho$   $\sigma$ )
- There is a gap extension penalty for each continuation on a level other than the main level (- $\sigma$ )

#### **The 3-leveled Manhattan Grid**



Gaps in w

#### Matches/Mismatches

Gaps in v

# **Affine Gap Penalty** Recurrences $\begin{array}{c} \downarrow \\ S_{i,j} \end{array} = \left\{ \begin{array}{c} \downarrow \\ S_{i-1,j} \end{array} - \sigma \\ S_{i-1,j} \end{array} \right. \left\{ \begin{array}{c} Continue \ Gap \ in \ w \ (deletion) \end{array} \right\}$ Continue Gap in $w \ (deletion)$ Start Gap in $w \ (deletion)$ : from middle $\vec{s}_{i,j} = \sigma$ $max \begin{bmatrix} \vec{s}_{i,j-1} & -\sigma \\ s_{i,j-1} & -(\rho+\sigma) \end{bmatrix}$ Continue Gap in *v* (insertion) $max \begin{bmatrix} \vec{s}_{i,j-1} & -(\rho+\sigma) \\ s_{i,j-1} & -(\rho+\sigma) \end{bmatrix}$ Start Gap in *v* (insertion):from middle $S_{i,j} = \begin{cases} s_{i-1,j-1} + \delta(v_i, w_j) & \text{Match or Mismatch} \\ s_{i,j} & \text{End deletion: from top} \\ s_{i,i} & \text{End insertion: from bot} \end{cases}$ End insertion: from bottom