

# CMSC 423:

# Sequence Alignment

Part5

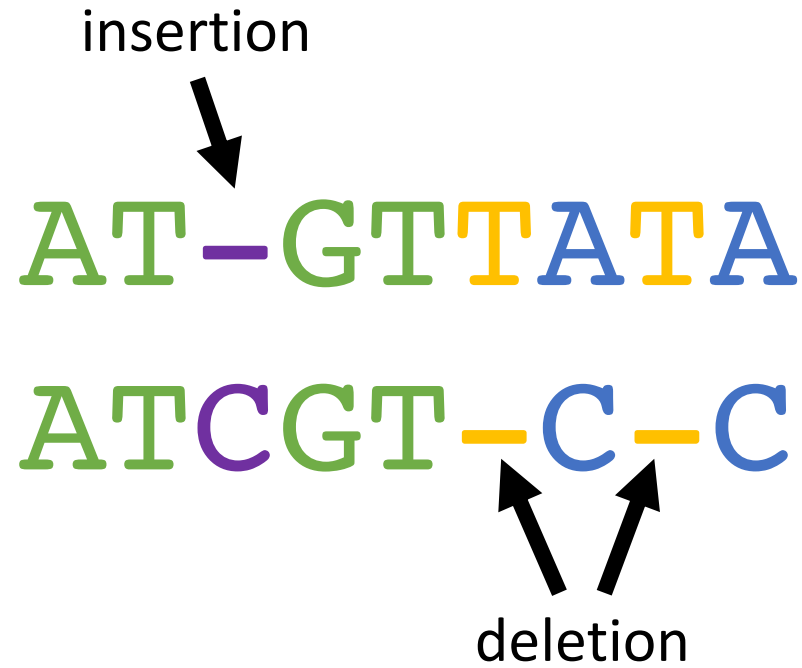
# Inexact matching: why?

- Redundancy in the genetic code: nucleotide sequence may differ, but proteins are the same
- Different amino-acid sequences can still fold the same way: function is unchanged
- Aligning RNA sequences to DNA- need to account for gaps corresponding to exons
- Sequencing errors

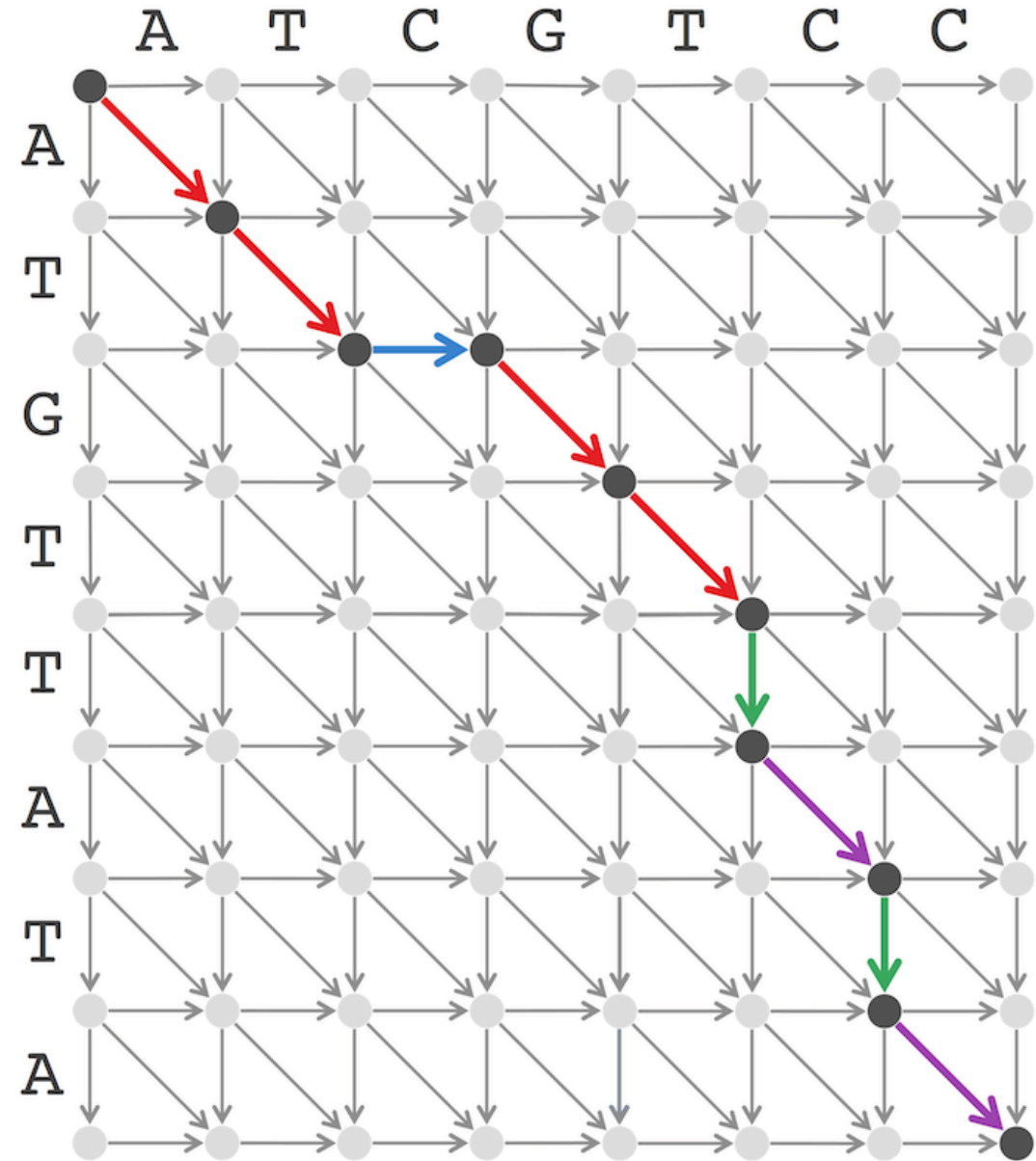
# Alignment of sequences $v$ and $w$

$v = \text{ATGTTATA}$

$w = \text{ATCGTCC}$



Sequence  
Alignment is  
the  
Manhattan  
Tourist  
Problem in  
Disguise



		i							
		-	A	G	C	G	T	A	G
j	-	0							
	G								
	T								
	C								
	A								
	G								
	A								

$v = \text{AGCGTAG}$

$w = \text{GTCAGA}$

$$s_{i,j} = \max \left\{ \begin{array}{l} s_{i-1,j-1} + 1, \\ \quad \text{if } v_i = w_j \\ s_{i-1,j} + 0 \\ s_{i,j-1} + 0 \end{array} \right.$$

	-	A	G	C	G	T	A	G
-	0	0	0	0	0	0	0	0
G	0	0	1	0	1	0	0	1
T	0	0	1	1	1	2	2	2
C	0	0	1	2	2	2	2	2
A	0	1	1	2	2	2	3	3
G	0	0	2	2	3	3	3	4
A	0	1	2	2	3	3	4	4

$v = \text{AGCGTAG}$

$w = \text{GTCAGA}$

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

# Local alignment

- What if we just want a region of similarity?
- Change the first row and column in the dynamic programming table to 0s
- Allow the alignment to start anywhere

$$\text{Score}[i,j] = \max\{0, \text{case 1}, \text{case 2}, \text{case 3}\}$$

- Answer is the location in the matrix with the highest score

# Extending to sequence alignment

- When solving for the LCS, mis-alignments are free
- What happens if we pay for our "mistakes"? (this also allows us to account for "similar" amino acids)

Value[Match] = 10

Value[Mismatch] = -5

Value[Gap] = -2

Match = [A,A], [C,C], [G,G], [T,T]

Mismatch = [A,G], [A,C], [A,T], ...

Gap = [A,-], [-,A], ...

The same dynamic programming algorithm works!

# Penalizing insertions and Deletions

- Linear scoring model
  - $\sigma$  = penalty for insertion or deletion of a single symbol
  - $\sigma \cdot k$  = penalty for insertion or deletion of  $k$  symbols

GATCAG

GA-C-AG

GATCAG

GA--CAG

# Penalizing insertions and Deletions

Mutations are often caused by errors in DNA replication that insert or delete an entire interval of  $k$  nucleotides as a single event (instead of  $k$  independent insertions or deletions)

GATCAG

GA—C—AG

GATCAG

GA——CAG

# Penalizing insertions and Deletions

Mutations are often caused by errors in DNA replication that insert or delete an entire interval of  $k$  nucleotides as a single event (instead of  $k$  independent insertions or deletions)

## Affine gap penalties

$$\begin{aligned}\text{Cost}(k \text{ gaps in a row}) &= \text{Cost}(\text{gap open}) + (k-1) * \text{Cost}(\text{gap}) \\ &= \sigma + (k-1) \cdot \varepsilon\end{aligned}$$

Gap opening penalty is high and gap extension penalty is low (once we start a gap we might as well pile more gaps on top)



and Think

How are the following alignments penalized using the new affine gap penalties?

GATCAG  
GA-C-AG

GATCAG  
GA--CAG



and Think

How are the following alignments penalized using the new affine gap penalties?

$$\begin{aligned}\text{Cost}(k \text{ gaps in a row}) &= \text{Cost}(\text{gap open}) + (k-1) \cdot \text{Cost}(\text{gap}) \\ &= \sigma + (k-1) \cdot \varepsilon\end{aligned}$$

GATCAG

GA-C-AG

$$= \sigma + (1-1) \cdot \varepsilon + \sigma + (1-1) \cdot \varepsilon$$

$$= 2 \sigma$$

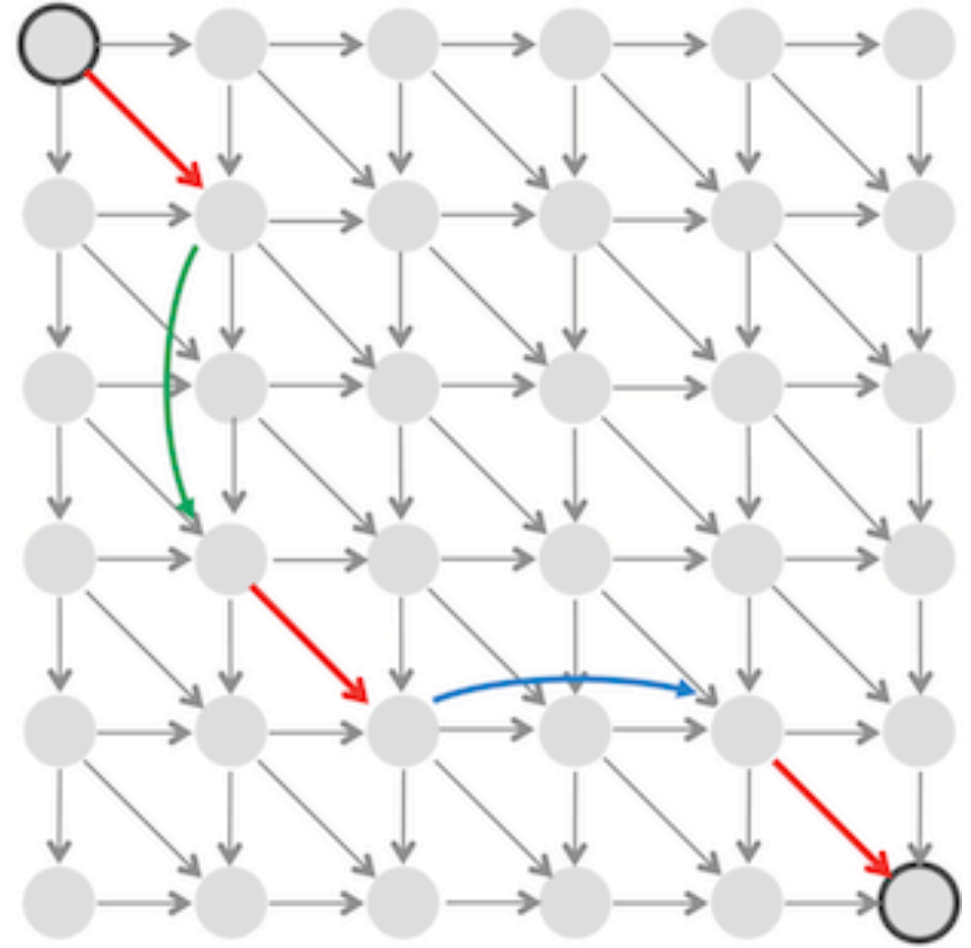
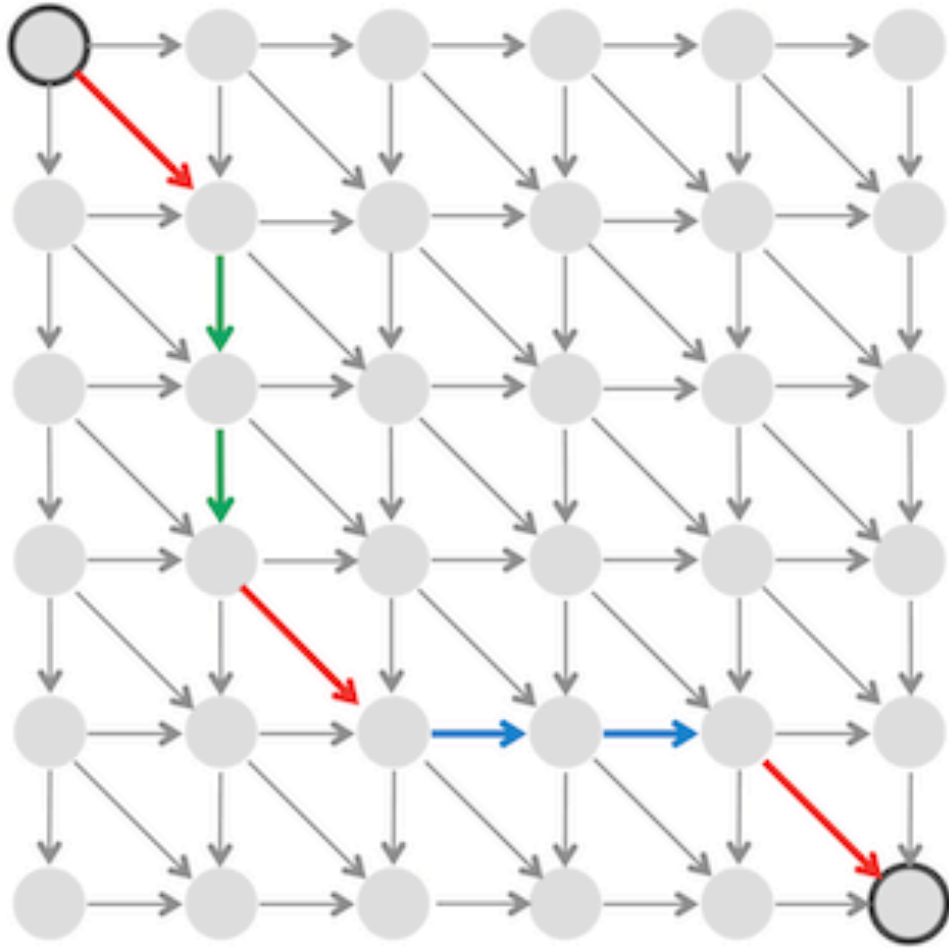
GATCCAG

GA--CAG

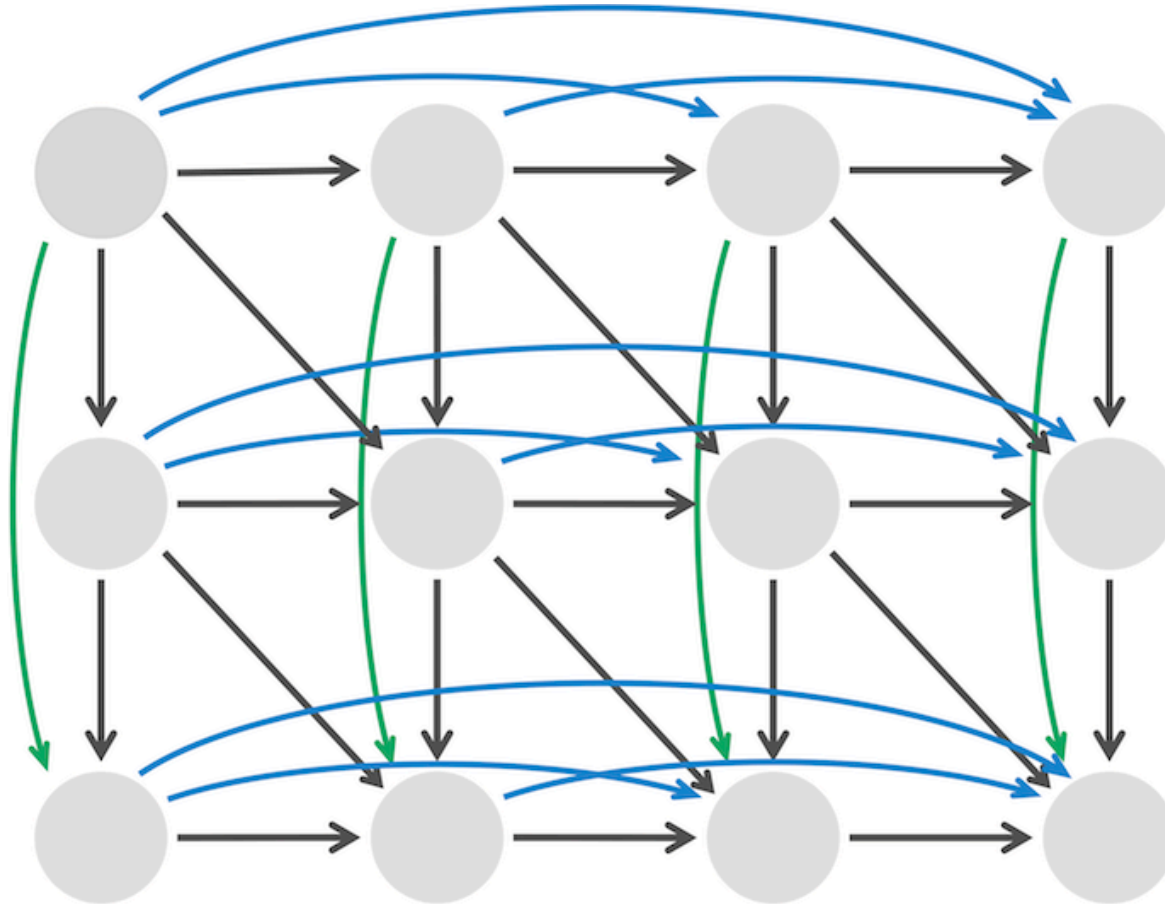
$$= \sigma + (2-1) \cdot \varepsilon$$

$$= \sigma + \varepsilon$$

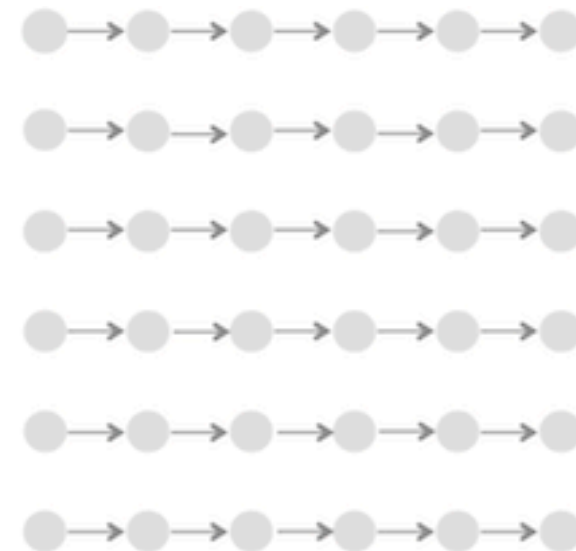
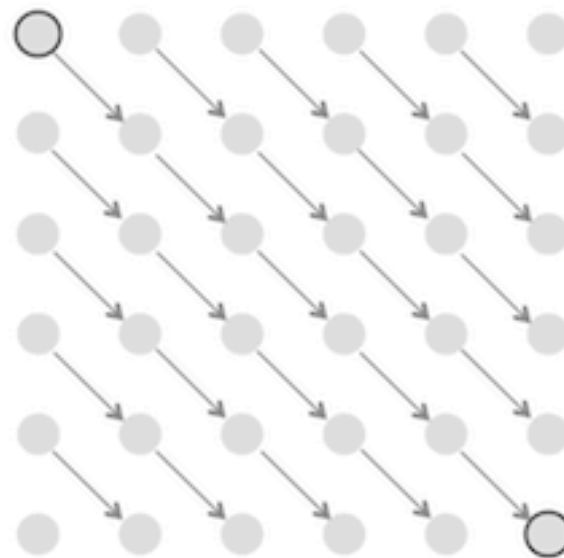
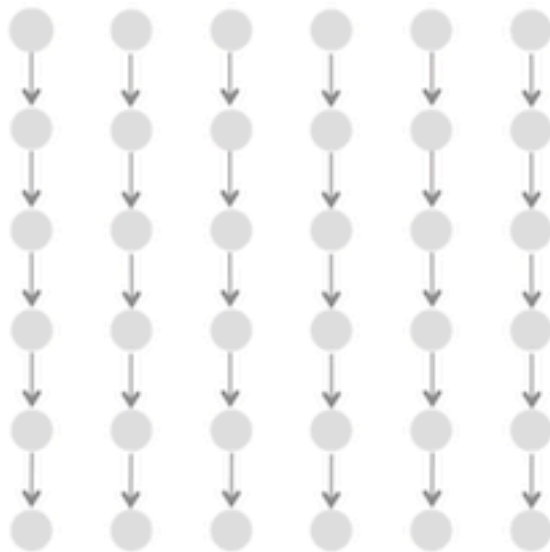
# Adding Affine Gap Penalties to the Alignment Graph



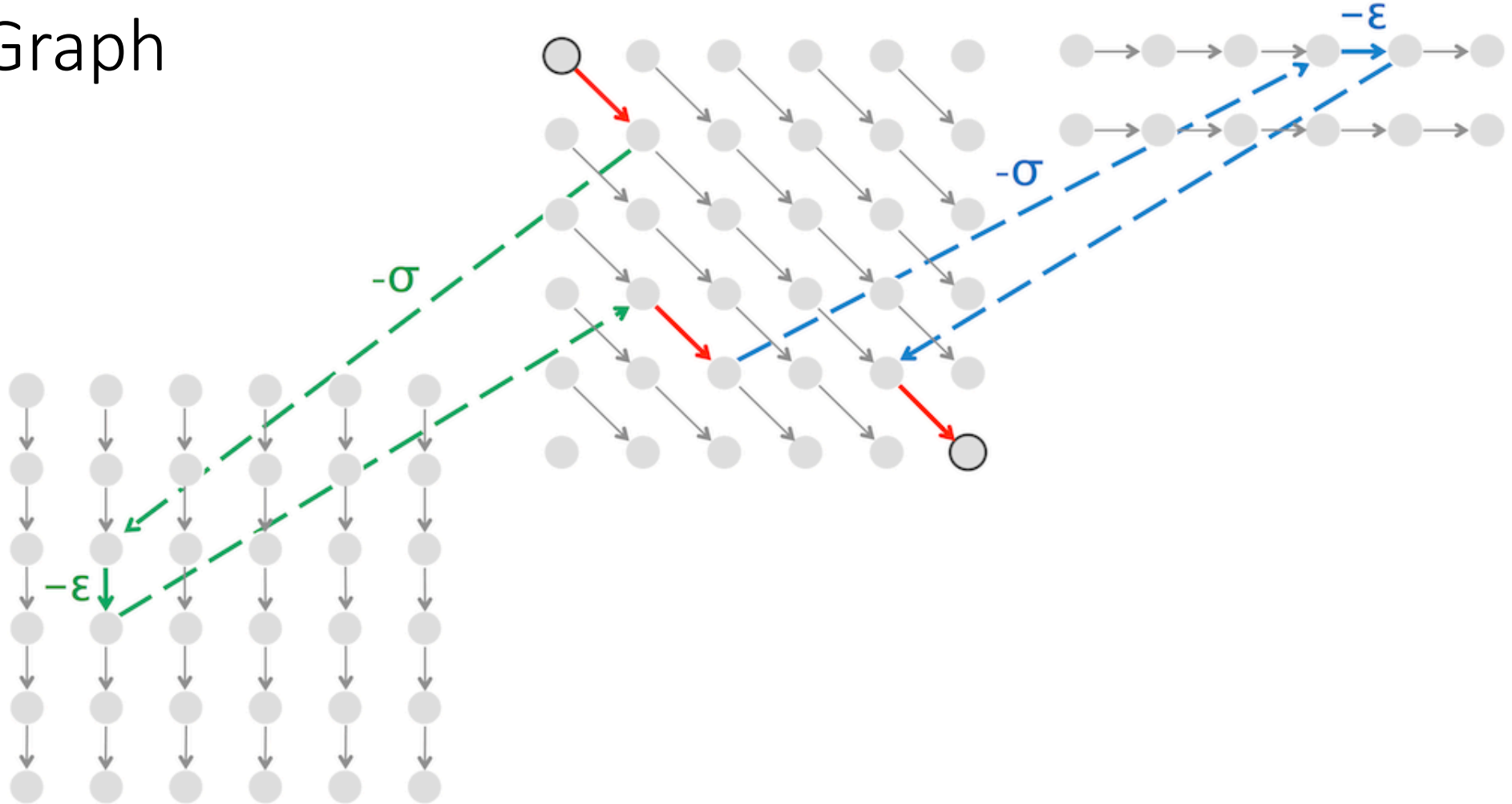
# Adding Affine Gap Penalties to the Alignment Graph



# Adding Affine Gap Penalties to the Alignment Graph



# Adding Affine Gap Penalties to the Alignment Graph



# Where do the alignment scores come from?

- PAM matrices –
  - PAM1 – based on frequency of mutations between closely related proteins (within 1 "evolutionary step")
  - PAM 2 - ... within 2 evolutionary steps – ...
  - PAM 250 – commonly used
- BLOSUM matrices – Frequency of mutations between proteins that are x% similar
  - BLOSUM100 – based on proteins that are exactly the same (e.g.  $\text{score}(A,A)$  is defined but not  $\text{score}(A,G)$  )
  - BLOSUM62 – commonly used
- gap scores usually determined empirically

# BLOSUM62

<b>Ala</b>	4																			
<b>Arg</b>	-1	5																		
<b>Asn</b>	-2	0	6																	
<b>Asp</b>	-2	-2	1	6																
<b>Cys</b>	0	-3	-3	-3	9															
<b>Gln</b>	-1	1	0	0	-3	5														
<b>Glu</b>	-1	0	0	2	-4	2	5													
<b>Gly</b>	0	-2	0	-1	-3	-2	-2	6												
<b>His</b>	-2	0	1	-1	-3	0	0	-2	8											
<b>Ile</b>	-1	-3	-3	-3	-1	-3	-3	-4	-3	4										
<b>Leu</b>	-1	-2	-3	-4	-1	-2	-3	-4	-3	2	4									
<b>Lys</b>	-1	2	0	-1	-3	1	1	-2	-1	-3	-2	5								
<b>Met</b>	-1	-1	-2	-3	-1	0	-2	-3	-2	1	2	-1	5							
<b>Phe</b>	-2	-3	-3	-3	-2	-3	-3	-3	-1	0	0	-3	0	6						
<b>Pro</b>	-1	-2	-2	-1	-3	-1	-1	-2	-2	-3	-3	-1	-2	-4	7					
<b>Ser</b>	1	-1	1	0	-1	0	0	0	-1	-2	-2	0	-1	-2	-1	4				
<b>Thr</b>	0	-1	0	-1	-1	-1	-1	-2	-2	-1	-1	-1	-1	-2	-1	1	5			
<b>Trp</b>	-3	-3	-4	-4	-2	-2	-3	-2	-2	-3	-2	-3	-1	1	-4	-3	-2	11		
<b>Tyr</b>	-2	-2	-2	-3	-2	-1	-2	-3	2	-1	-1	-2	-1	3	-3	-2	-2	2	7	
<b>Val</b>	0	-3	-3	-3	-1	-2	-2	-3	-3	3	1	-2	1	-1	-2	-2	0	-3	-1	4
	<b>Ala</b>	<b>Arg</b>	<b>Asn</b>	<b>Asp</b>	<b>Cys</b>	<b>Gln</b>	<b>Glu</b>	<b>Gly</b>	<b>His</b>	<b>Ile</b>	<b>Leu</b>	<b>Lys</b>	<b>Met</b>	<b>Phe</b>	<b>Pro</b>	<b>Ser</b>	<b>Thr</b>	<b>Trp</b>	<b>Tyr</b>	<b>Val</b>