

How Do We Compare Biological Sequences?

*Dynamic Programming and
Divide-and Conquer Algorithms*

Phillip Compeau and Pavel Pevzner.

Bioinformatics Algorithms: An Active Learning Approach

©2018 by Compeau and Pevzner. All rights reserved.

How Do We Compare Biological Sequences

- **From Sequence Comparison to Biological Insights**
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

RNA Tie Club

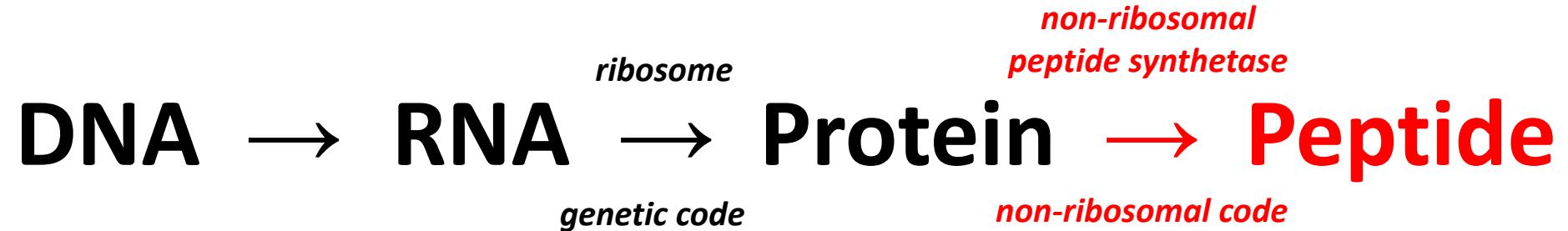
DNA → RNA → Protein



RNA tie club

Bioinformatics Algorithms: An Active Learning Approach.
Copyright 2018 Compeau and Pevzner.

From Genetic Code to Non-Ribosomal Code



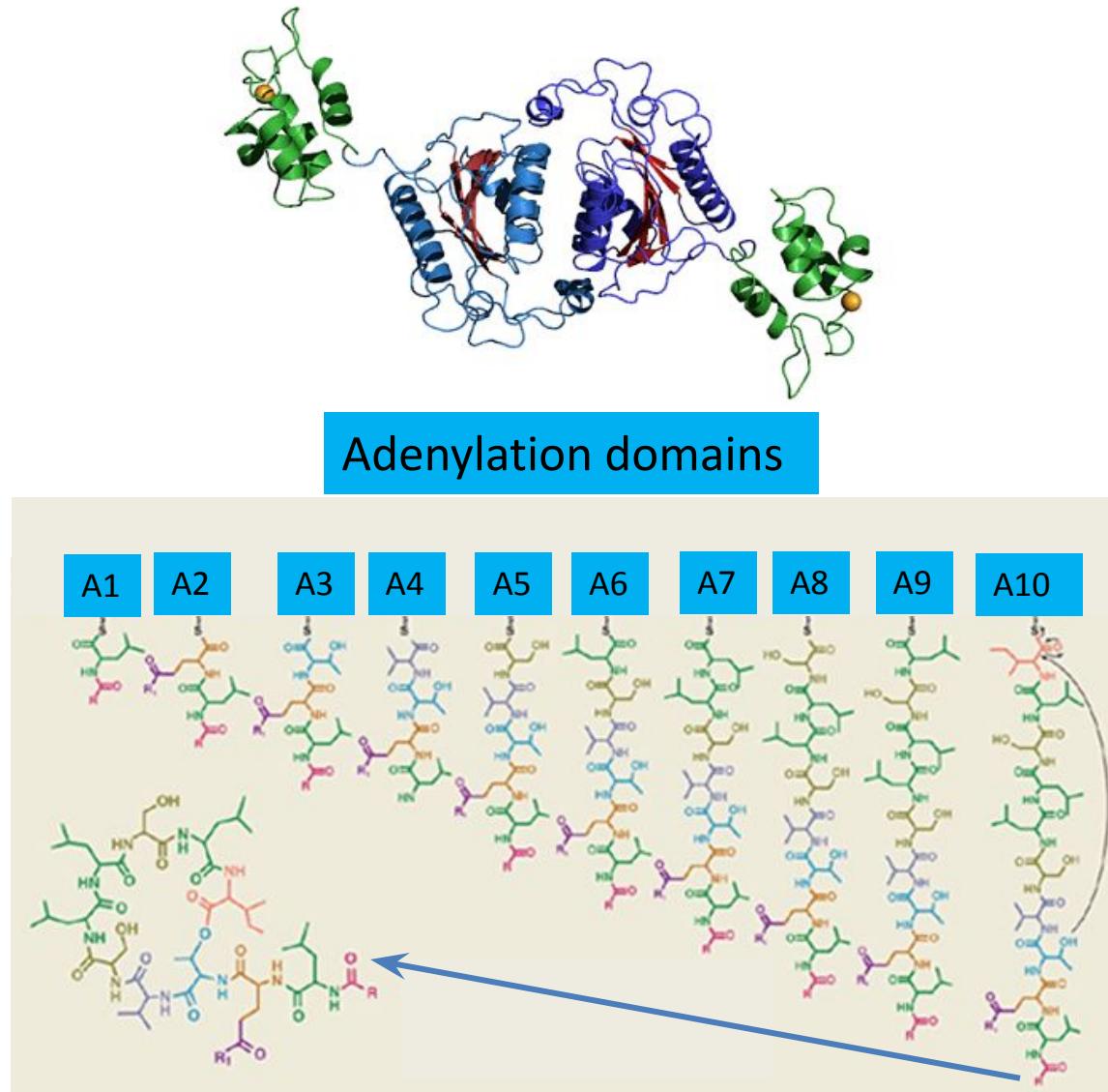
RNA tie club

Marahiel

Bioinformatics Algorithms: An Active Learning Approach.

Copyright 2018 Compeau and Pevzner.

NRP Synthetase: A Giant Molecular Assembly Line



Bioinformatics Algorithms: An Active Learning Approach.
NRP synthetase adds one amino acid at a time
Copyright 2018 Compeau and Pevzner.

These Three A-domains Do Not Look Similar

YAFDLGYTCMFVLLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA
AFDVSAGDFARALLTGGQLIVCPNEVKMDPASLYAIKKYDITIFEA TPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRGSTIRIVNSYGVTEACIDS
IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPPALLQCLVSAPTMISSLEILFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

These Three A-domains Do Not Look Similar

YAFDLGYTCMFPVLGGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPSLFHTIVNTASFAFDANFESLRLIVLGGEKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA
AFDVSAGDFARALLTGSQLIVCPNEVKMDPASLYAIIKKYDITIFEATPALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRGSTIRIVNSYGVTEACIDS
IAFDASSWEIYAPLLNGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLPPALLQCLVSAPTMISSLIELFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

just 3 conservative columns

Do they look similar now?

Y**A****F**D LGYTCMF**P**V**L**GG**E**L HIVQKETYTAPDEIAHYI**K**EHG**I**TYIKLT**PS**LFHTIVNTASFAFDANFESLRLIVLG**G**EKIIPIDVIAFRKMYGHTEFINHYGPTEATIGA
-**A****F**D VSAGDFAR**A****L****L**T**G**QLIVCPNEVKMDPASLYAI**I****K**YD**I**TIFEA**T****P**ALVIPLMEYIYEQKLDISQLQILIVGSDSCSMEDFKTLVSRGSTIRIVNSYGVTEACIDS
I**A****F**D ASSWEIYAP**L****L****N****G**TVVCIDYYTTIDIKALEAVF**K**QHH**I**RGAMLP**P**ALLKQCLVSAPTMISSLIELFAAGDRLSSQDAILARRAVGSGVYNAYGPTENTVLS

11 conservative columns

And now?

Y**A****F**D~~L~~G~~Y~~T~~C~~M~~F~~P~~V~~**L****L****G**~~G~~E~~H~~I~~V~~Q~~K~~E~~T~~Y~~T~~A~~P~~E~~I~~A~~H~~Y~~I~~**K**~~E~~H~~G~~**I**~~T~~Y~~I~~K~~L~~T~~P~~S~~L~~F~~H~~T~~I~~V~~N~~T~~A~~S~~F~~A~~F~~D~~A~~N~~F~~E~~S~~**L**R~~L~~I~~V~~L~~G~~G~~E~~K~~I~~I~~P~~I~~D~~V~~I~~A~~F~~R~~K~~M~~Y~~**G**~~H~~T~~E~~~~-~~**F**~~I~~N~~H~~**Y**~~G~~**P**~~T~~E~~A~~T~~I~~G~~A~~
~~-~~**A****F**D~~V~~S~~A~~G~~D~~F~~A~~R~~A~~**L**~~L~~**T****G**~~G~~Q~~L~~I~~V~~C~~P~~N~~E~~V~~K~~M~~D~~P~~A~~S~~L~~Y~~A~~I~~I~~**K**~~K~~Y~~D~~**I**~~T~~I~~F~~E~~A~~T~~P~~A~~L~~V~~I~~P~~L~~M~~E~~Y~~I~~~~-~~Y~~E~~Q~~K~~L~~D~~I~~S~~Q~~L~~Q~~I~~I~~V~~G~~S~~D~~C~~S~~M~~E~~D~~F~~K~~T~~L~~V~~S~~R~~F~~**G**~~S~~T~~I~~R~~I~~V~~N~~**S**~~Y~~**G**~~V~~~~T~~E~~A~~C~~I~~D~~S~~
~~I~~**A****F**D~~A~~S~~S~~W~~E~~I~~Y~~A~~P~~**L**~~L~~**N****G**T~~V~~V~~C~~I~~D~~Y~~Y~~T~~I~~D~~I~~K~~A~~L~~E~~A~~V~~F~~K~~Q~~H~~H~~I~~R~~G~~A~~M~~L~~P~~**P**~~A~~~~L~~K~~Q~~C~~L~~V~~S~~A~~---~~~~---~~~~PT~~MI~~S~~**L**E~~I~~F~~A~~A~~G~~D~~R~~L~~S~~**S**~~Q~~D~~A~~I~~L~~A~~R~~R~~A~~V~~G~~SG~~V~~~~-~~**Y**~~-~~**N**~~A~~**Y**~~G~~**P**~~T~~E~~N~~V~~L~~S

19 conservative columns!

Red Positions Encode Conservative Core of A-domains

Y**A**FDLGYTCMFV**L**GG**E**LHVQKETYTAPDEIAHYI**K**EHG**I**TYIKLT**PS**L**F**H**T**IVNTASFAFDANFES**L**R**L**IVLG**G**E**K**I**I**P**I**D**V**IAFRK**M**Y**G**H**T**E**-****F****I**N**H****Y****G****P****T****E**ATIGA
-AFDVSAGDFAR**A**LL**T****GG**QLIVCPNEVKMDPASLYAI**I****K**KYD**I**T**I**FEAT**P****A**LV**I**PLMEY**I****-**YEQ**K**L**D**IS**Q****L**Q**I**LVGS**D****S****C****S****M****E****D****F****K****T****L****V****S****R****F****G****S****T****I****R****I****V****N****S****Y****G****V****T****E**ACIDS
I**A**FDASSWEIYAP**L****L****N****GG**TVVCIDYYTTIDIKALEAVF**K**QHH**I**RGAMLPP**A**LL**K****Q****C****L****V****S****A****--****--****P****T****M****I****S****S****L****E****I****L****F****A****G****D****R****L****S****S****Q****D****A****I****L****A****R****R****A****V****G****S****G****V****-****Y****-****N****A****Y****G****P****T****E**NTVLS

Which positions are responsible for encoding **different** amino acids Asp, Orn, Val?

Blue Positions in A-domains Define Non-Ribosomal Code

YAFD**L**GY**T**CMFPV**L**GG**E**LHIVQ**K**ETYTAPDEIAHY**I**KEHG**I**TY**I**KLT**PSL**FHTIVNTASFAFDANFES**L**R**L****I****V****L****G**GEK**I****I****P****I**DVIAFRKMY**G**HTE-FIN**H****Y****G**P**T**EAT**I**GA
-AFD**V**SAGDFAR**A****L****L**TGG**Q**LVCPNEVKMDPASLY**A****I****I**KKYD**I**T**I****F**EAT**P****A****L****V****I****P****L****M****E****Y****I**-YEQ**K**L**D**ISQL**Q****I****L****I****V****G**SDSCSMEDFKTLVSRF**G****S****T****I****R****I****V****N****S****Y****G****V****T**EAC**I****D****S**
IAFD**A****S****S****W****E****I****Y****A****P****L****L****N****G****G****T****V****V****C****I****D****Y****Y****T****T****I****D****I****K****A****L****E****A****V****F****K****Q****H****H****I****R****G****A****M****L****P****P****A****L****L****K****Q****C****L****V****S****A**---PTMISS**L****E****I****L****F****A****A****G****D****R****L****S****S****Q****D****A****I****L****A****R****R****A****V****G****S****G****V**-Y-NAYGP**T****E****N****T****V****L****S**

LTKVGHIG	Asp
VGEIGSID	Orn
AWMFAAVL	Val

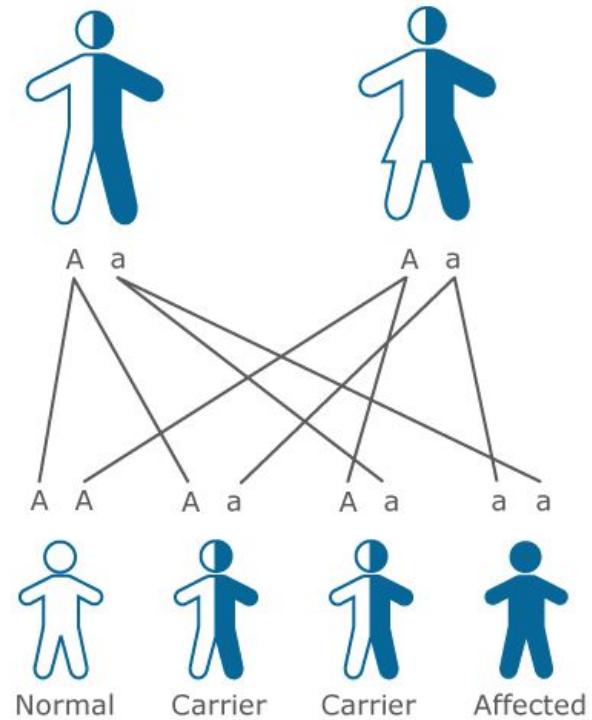
Cystic Fibrosis

- **Cystic fibrosis (CF):** An often fatal disease which affects the respiratory system and produces an abnormally large amount of mucus.
 - Mucus is a slimy material that coats epithelial surfaces and is secreted into fluids such as saliva.



Approximately 1 in 25 Humans Carry a Faulty CF Gene

- In the early 1980s biologists hypothesized that CF is caused by mutations in an unidentified gene.
- When BOTH parents carry a faulty gene, there is a 25% chance that their child will have cystic fibrosis.



Where Is the Cystic Fibrosis Gene?

- In the late 1980s, biologists narrowed the search for the CF gene to a small region on chromosome 7.
- One of these genes was **similar** to **ATP binding proteins** that act as transport channels responsible for secretion.
- **Hint:** cystic fibrosis involves sweet secretion with abnormally high sodium levels.



Chromosome 7

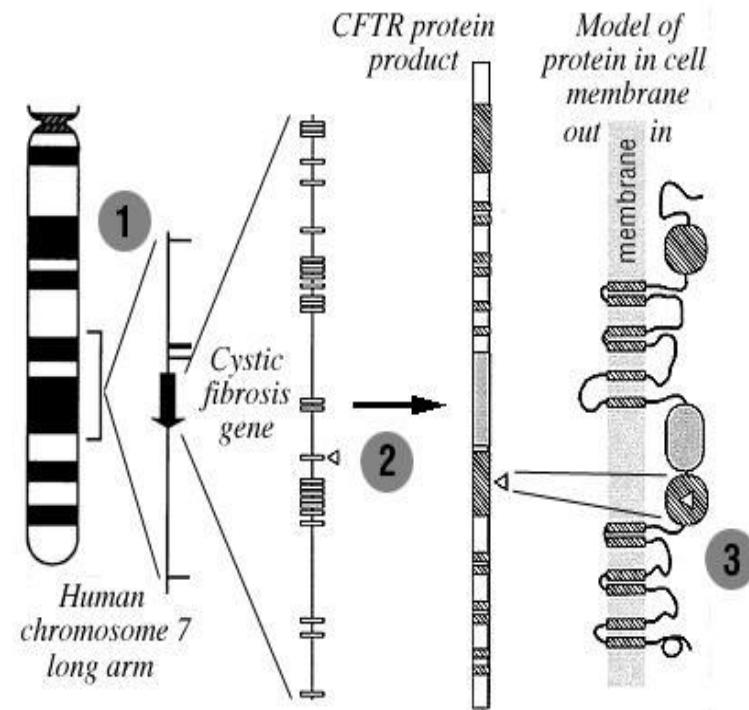
Adenosine Triphosphate (ATP) transports chemical energy within cells.

Bioinformatics Algorithms: An Active Learning Approach.

Copyright 2018 Compeau and Pevzner.

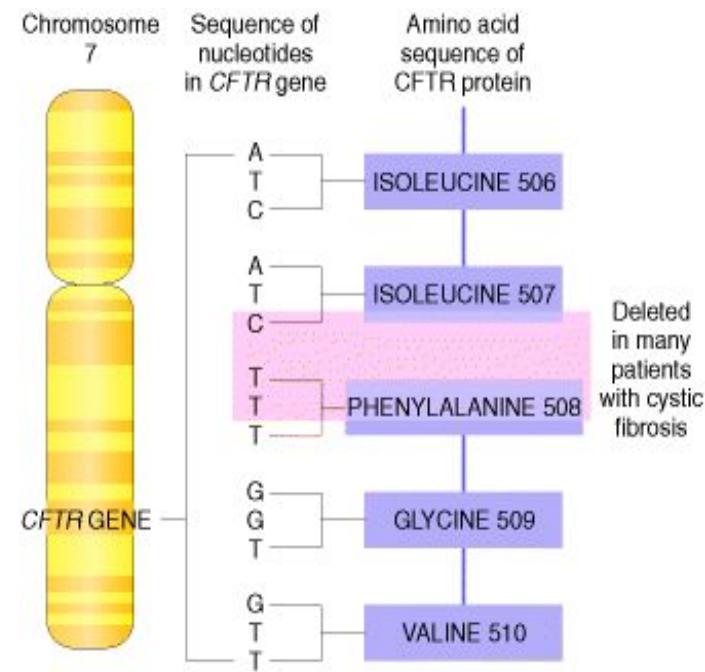
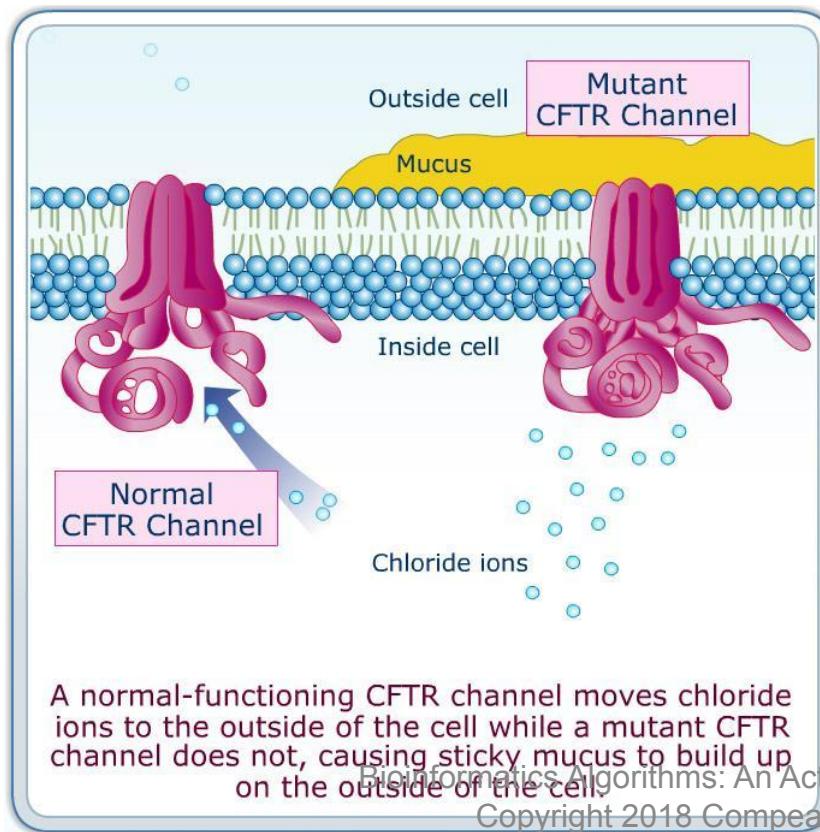
Where Is the Cystic Fibrosis Gene?

- In the late 1980s, biologists narrowed the search for the CF gene to a small region on chromosome 7.
- One of these genes was **similar** to **ATP binding proteins** that act as transport channels responsible for secretion.
- **Hint:** cystic fibrosis involves sweet secretion with abnormally high sodium levels.



CFTR: Cystic Fibrosis Transmembrane Conductance Regulator

- The CFTR protein controls the flow of ions in and out of cells inside the lungs.



How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- **The Alignment Game and the Longest Common Subsequence**
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

The Alignment Game

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

Alignment Game (maximizing the number of points):

- Remove the 1st symbol from each sequence
 - 1 point if the symbols match, 0 points if they don't match
- Remove the 1st symbol from one of the sequences
 - 0 points

The Alignment Game

A T G T T A T A
A T C G T C C
+1

The Alignment Game

A	T	G	T	T	A	T	A
A	T	C	G	T	C	C	
+1+1							

The Alignment Game

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

The Alignment Game

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

The Alignment Game

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

The Alignment Game

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

The Alignment Game

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

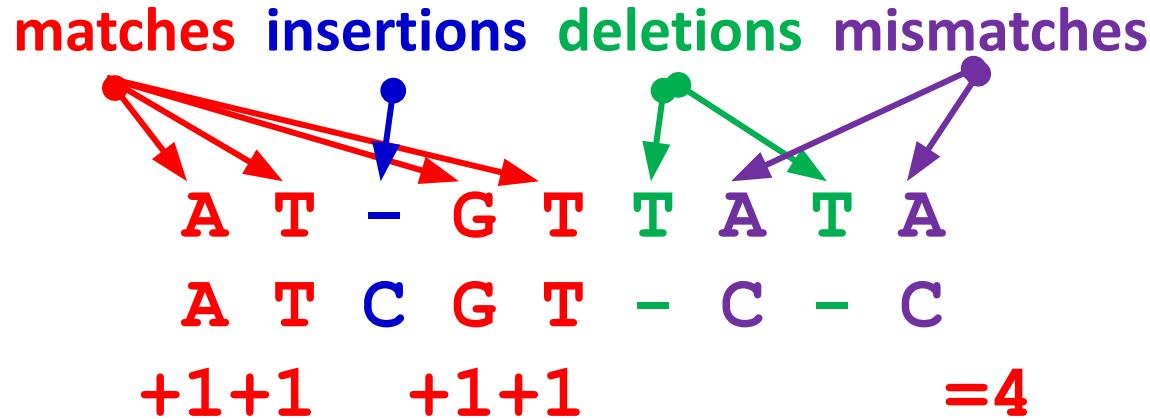
The Alignment Game

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

The Alignment Game

A	T	-	G	T	T	A	T	A
A	T	C	G	T	-	C	-	C
+1+1	+1+1						=4	

What Is the Sequence Alignment?



Alignment of two sequences is a two-row matrix:

1st row: symbols of the 1st sequence (in order) interspersed by “-”
2nd row: symbols of the 2nd sequence (in order) interspersed by “-”

Longest Common Subsequence

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

Matches in alignment of two sequences (**ATGT**) form their
Common Subsequence

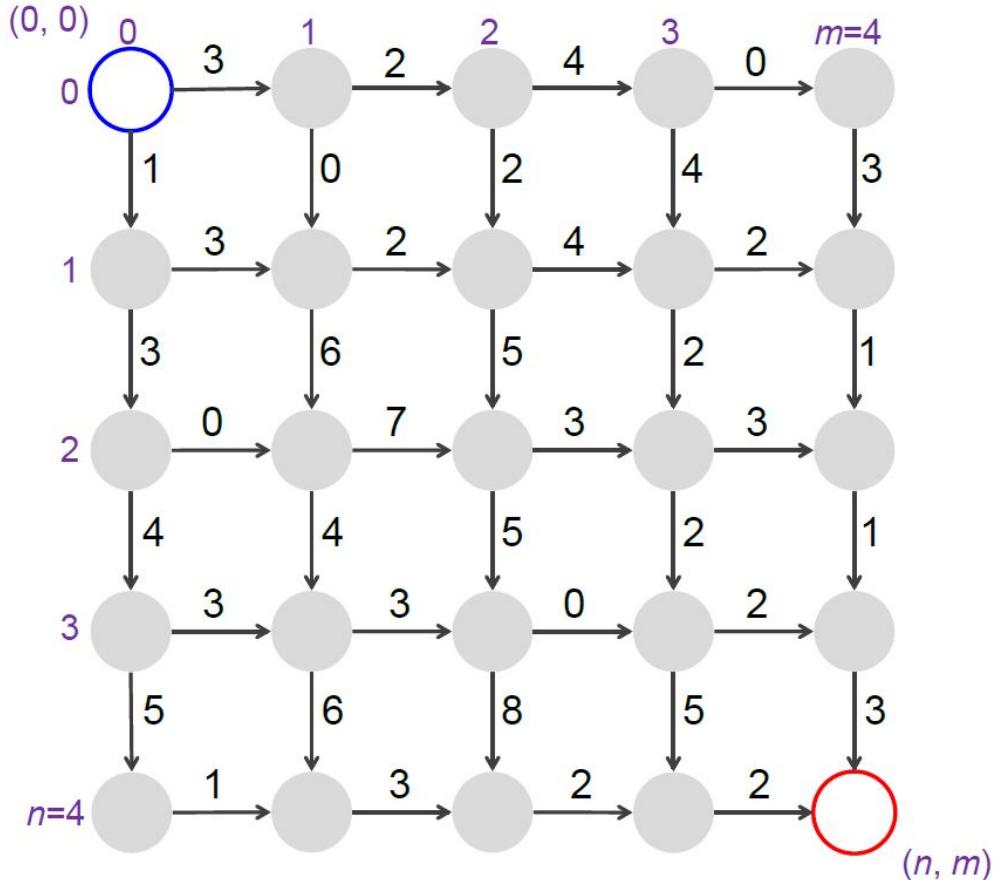
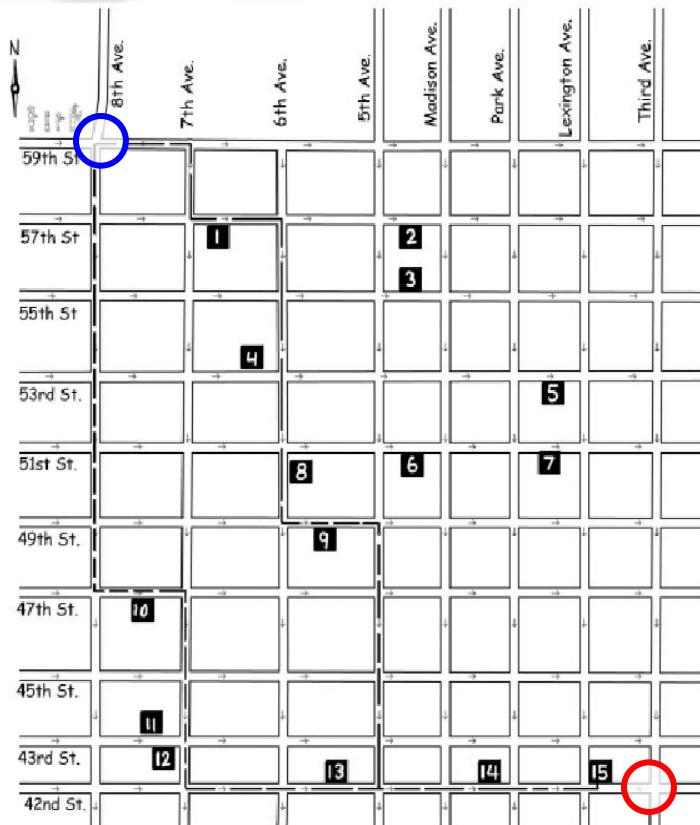
Longest Common Subsequence Problem: Find a longest common subsequence of two strings.

- **Input:** Two strings.
- **Output:** A longest common subsequence of these strings.

How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- **The Manhattan Tourist Problem**
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

From Manhattan to a Grid Graph

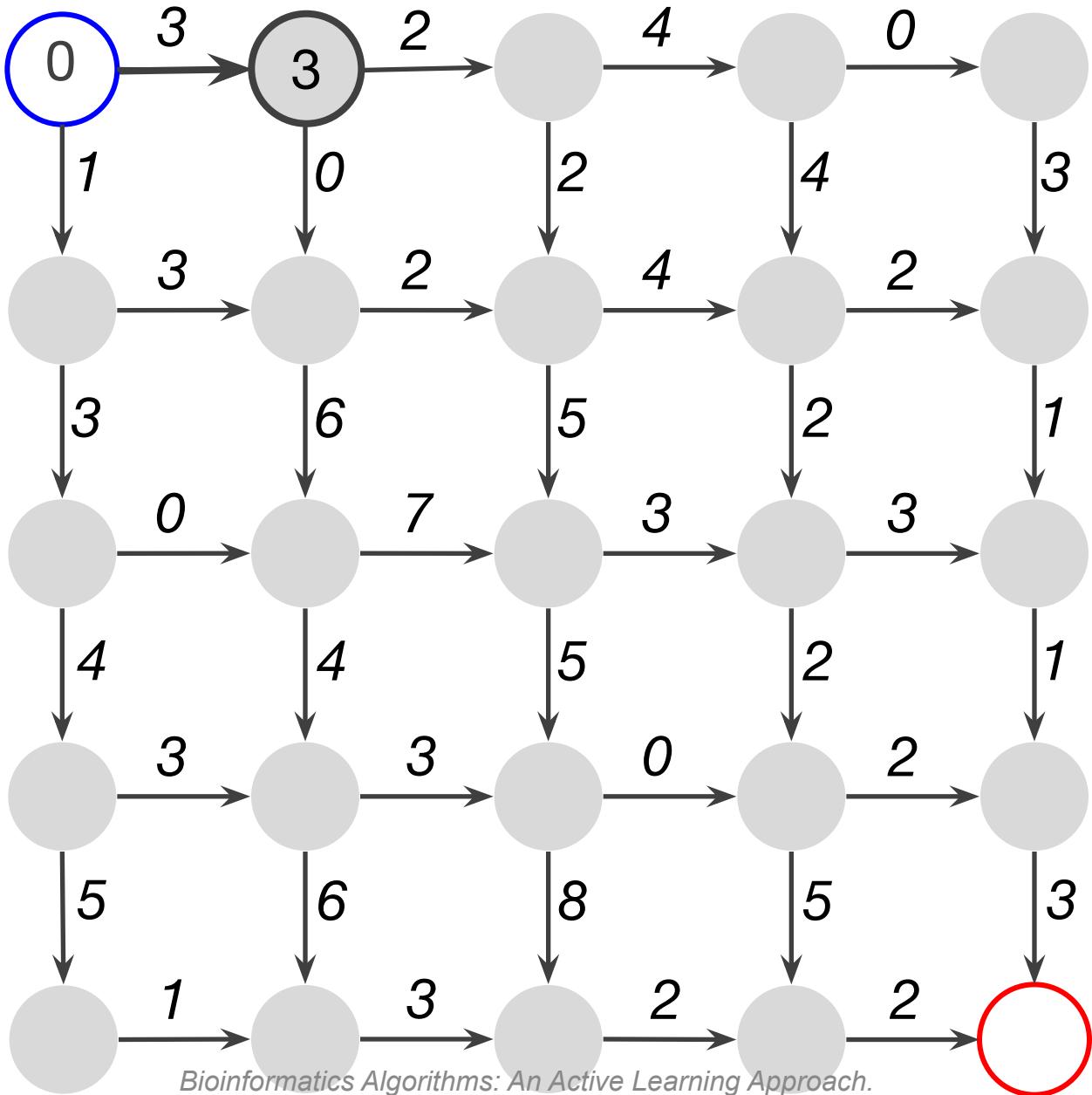


Manhattan Tourist Problem

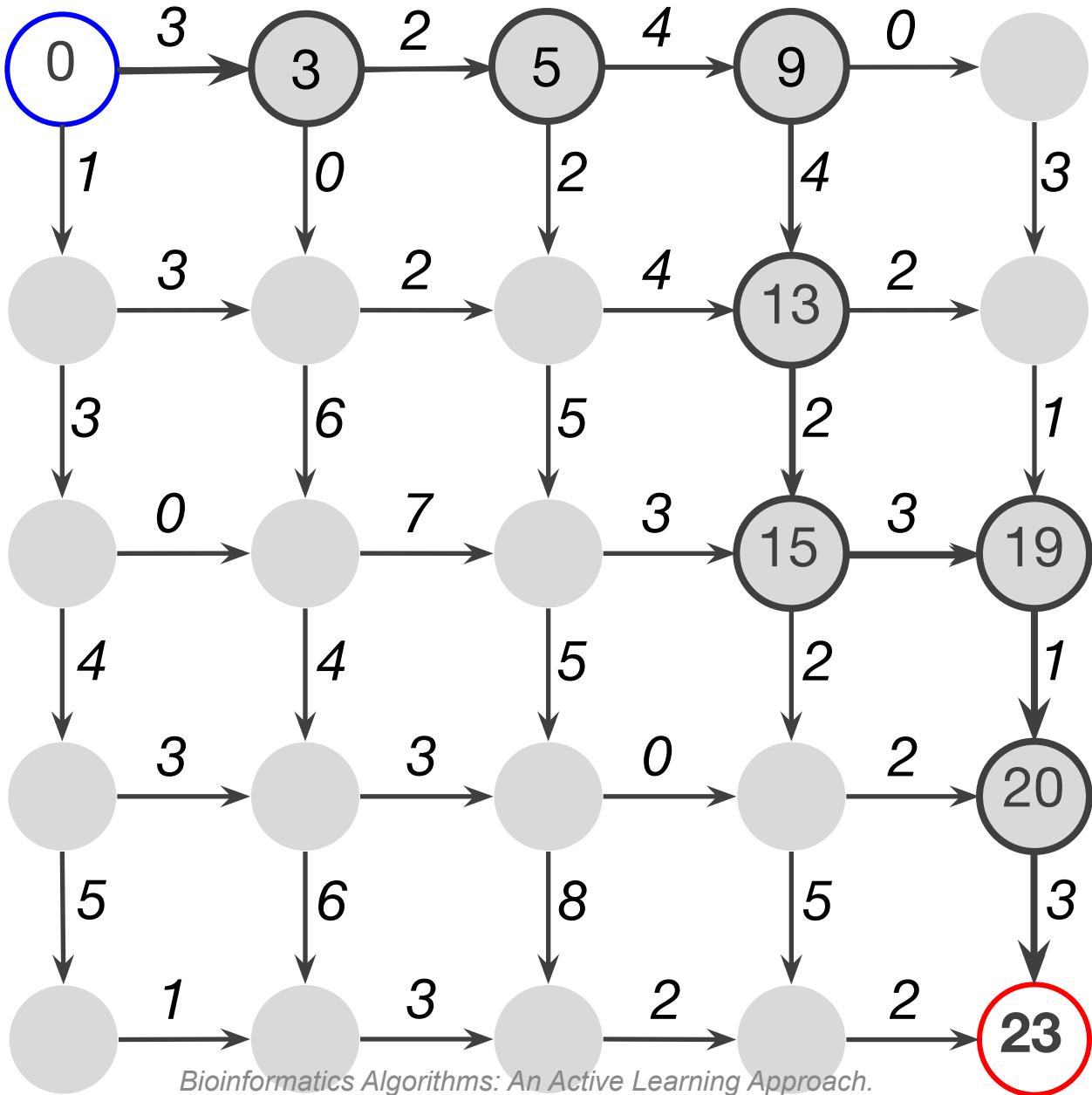
Manhattan Tourist Problem: Find a longest path in a rectangular city grid.

- **Input:** A weighted rectangular grid.
- **Output:** A longest path from the source to the sink in the grid.

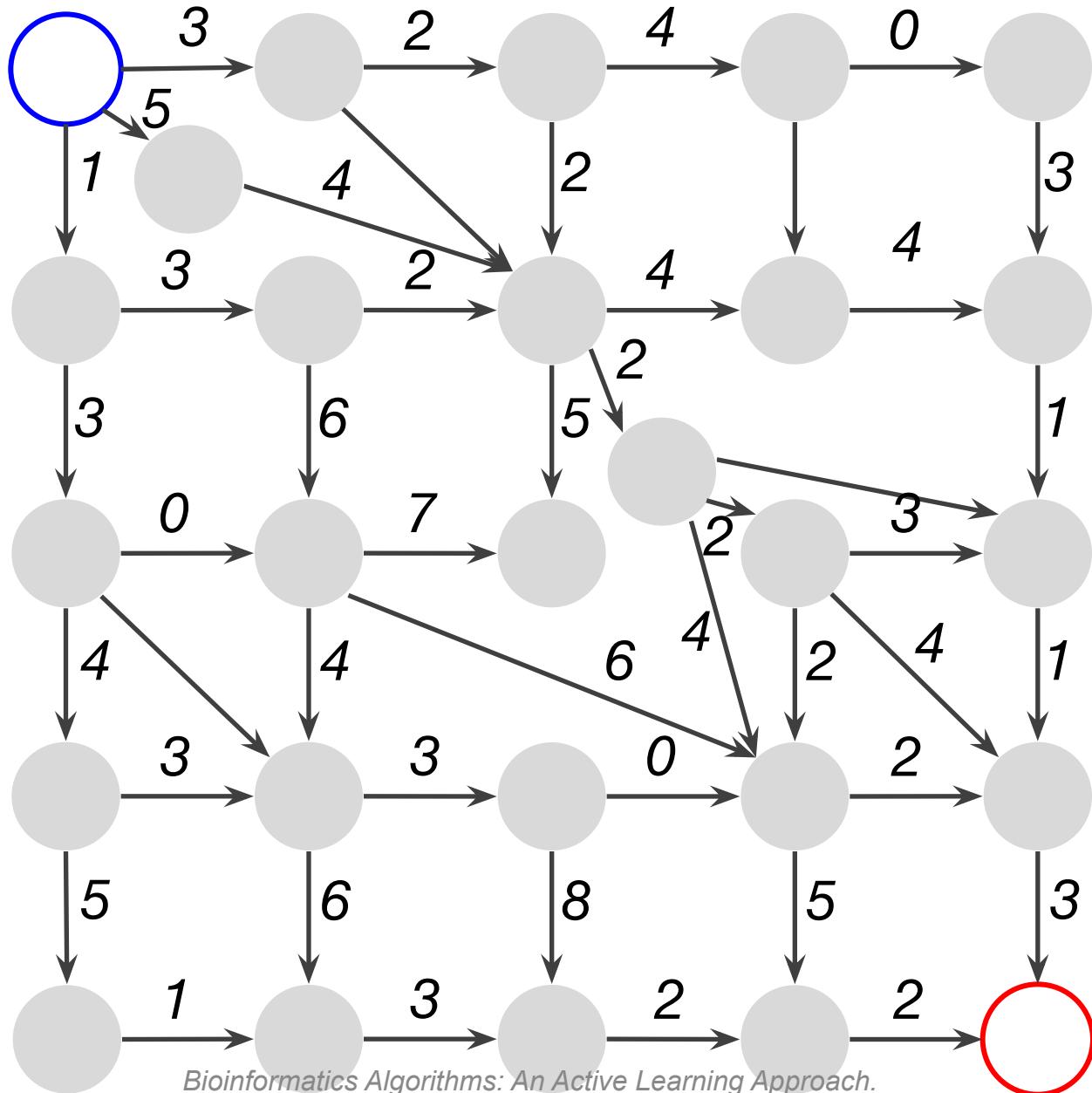
**Greedy
algorithm?**



**Greedy
algorithm?**



**From a
regular to an
irregular grid**



Search for Longest Paths in a Directed Graph

Longest Path in a Directed Graph Problem: Find a longest path between two nodes in an edge-weighted directed graph.

- **Input:** An edge-weighted directed graph with source and sink nodes.
- **Output:** A longest path from source to sink in the directed graph.

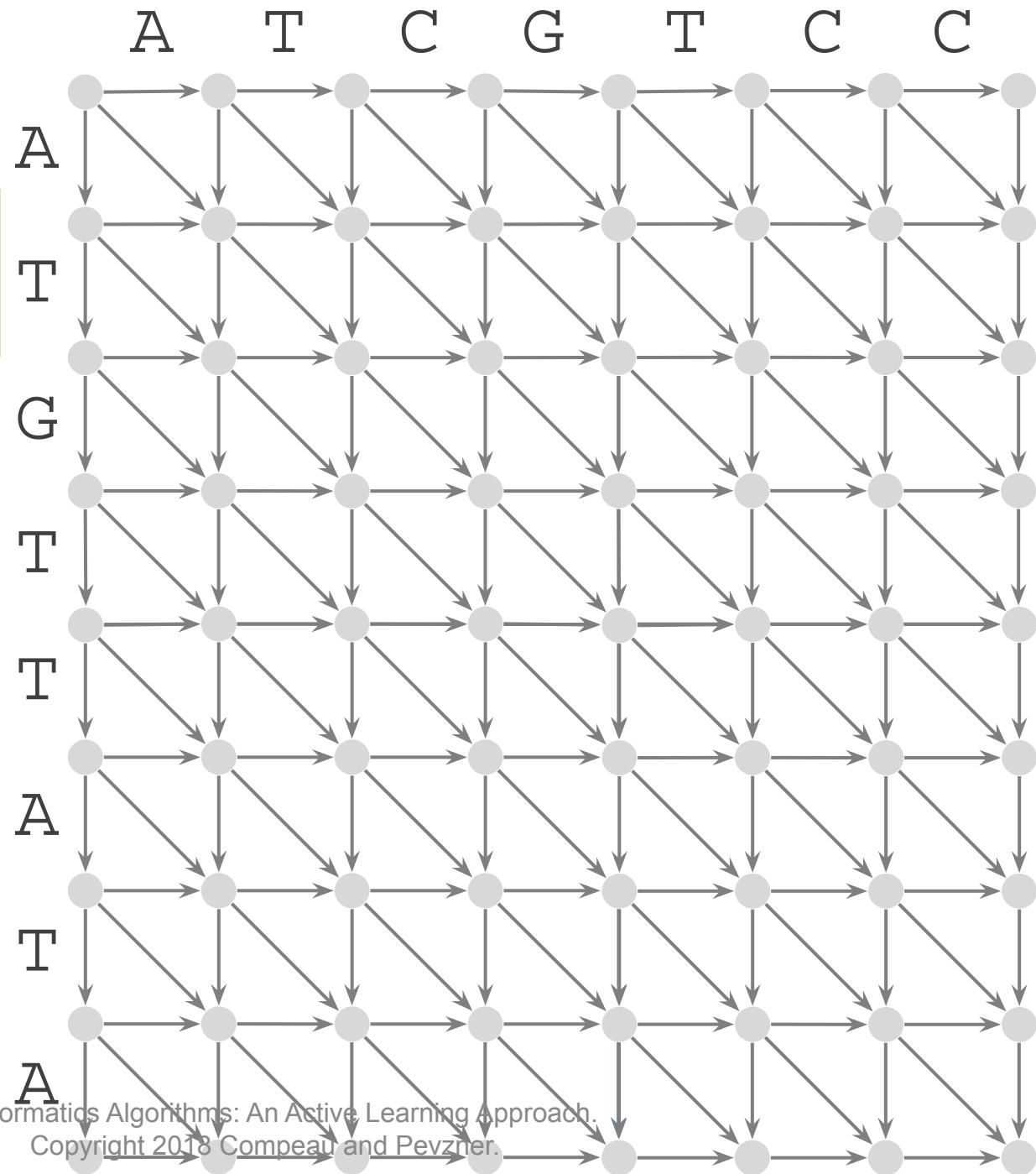
Do You See a Connection between the Manhattan Tourist and the Alignment Game?

A	T	-	G	T	T	A	T	A
A	T	C	G	T	-	C	-	C
↓	↓	→	↓	↓	↓	↓	↓	↓

?

alignment → path

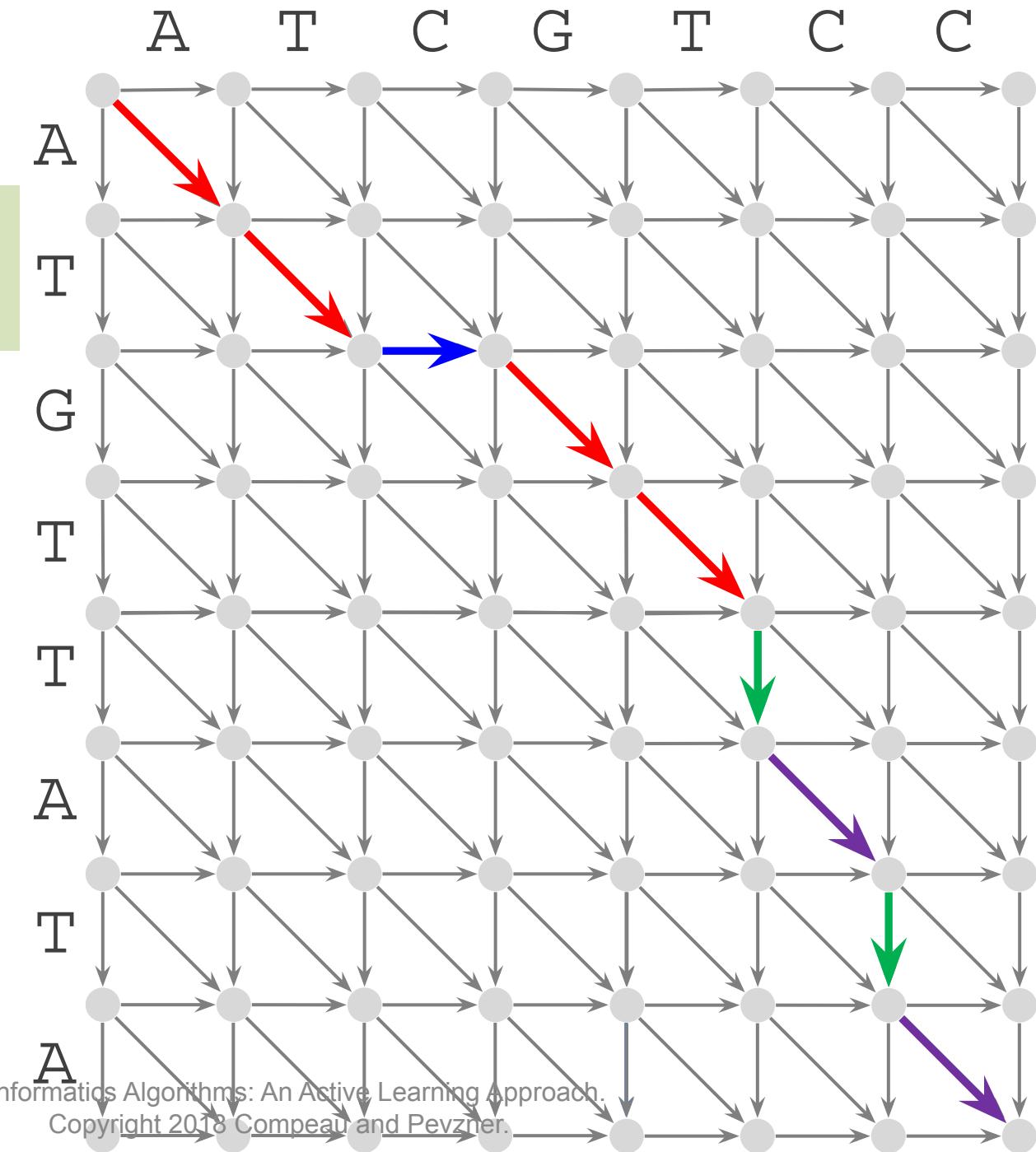
A T - G T T A T A
A T C G T - C - C
↓ ↓ → ↓ ↓ ↓ ↓ ↓



?

alignment → path

A	T	-	G	T	T	A	T	A
A	T	C	G	T	-	C	-	C
↓	↓	→	↓	↓	↓	↓	↓	↓



?

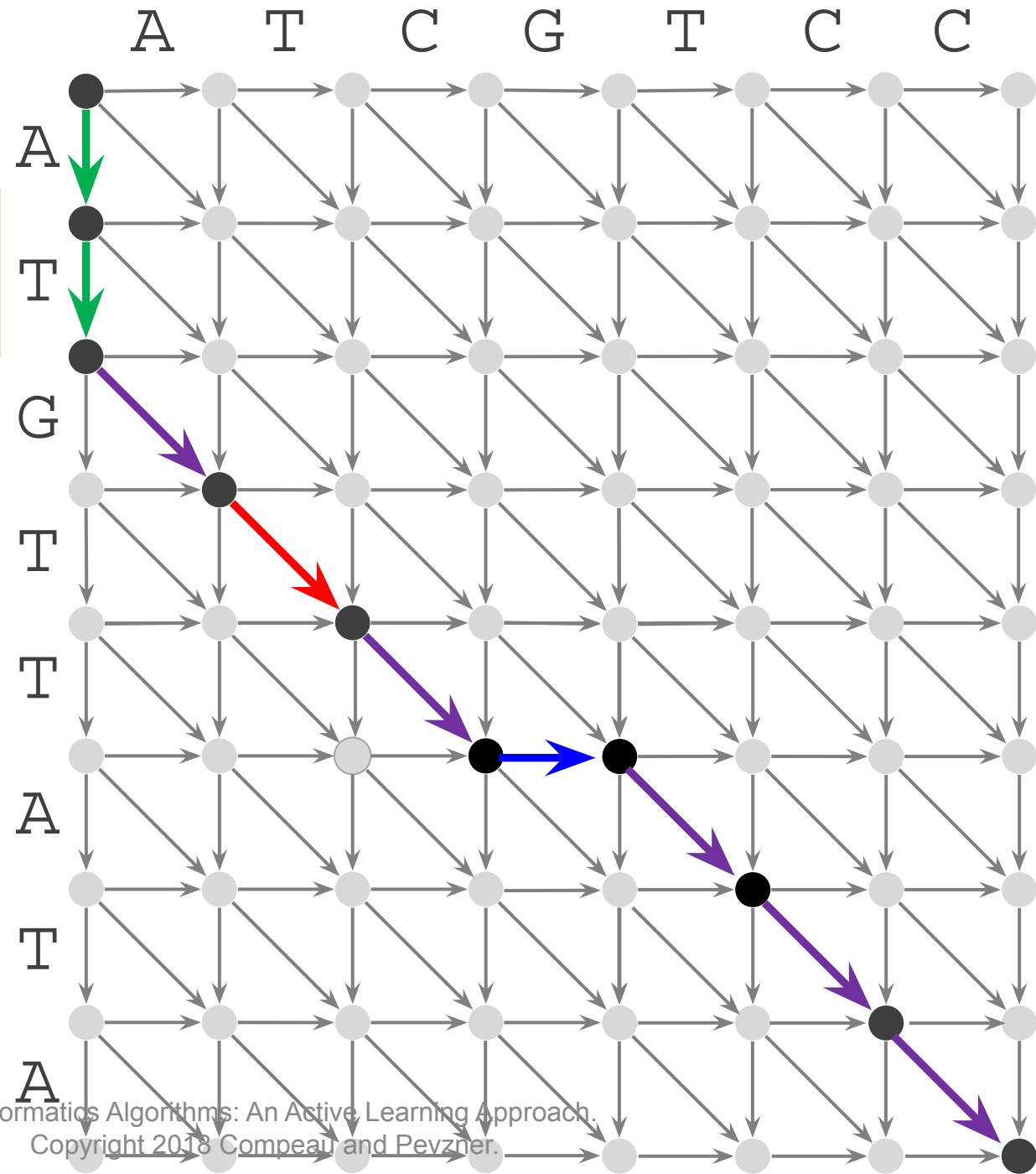
path → alignment

A T G T T - A T A
- - A T C G T C C
↓ ↓ ↓ ↘ ↓ → ↓ ↓ ↓

highest-scoring alignment

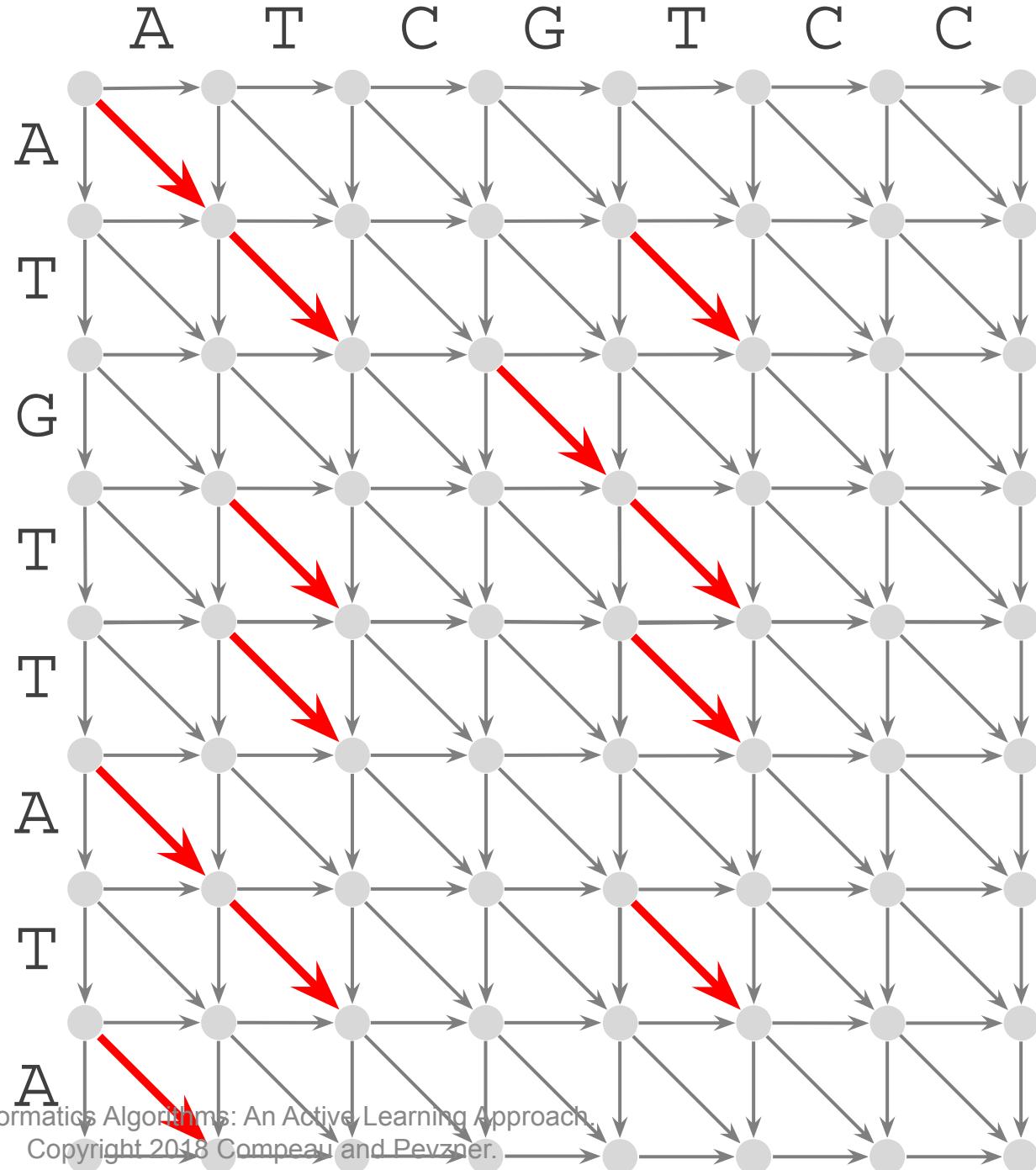
=

longest path in a properly built Manhattan



How to built a Manhattan for the Alignment Game and the Longest Common Subsequence Problem?

Diagonal red edges correspond to matching symbols and have scores 1



How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- **The Change Problem**
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

The Change Problem

Change Problem: Find the minimum number of coins needed to make change.

- **Input:** An integer *money* and an array of positive integers ($coin_1, \dots, coin_d$).
- **Output:** The minimum number of coins with denominations ($coin_1, \dots, coin_d$) that changes *money*.



Changing Money in a Greedy Way

GreedyChange(*money*)

change \leftarrow empty collection of coins

while *money* $>$ 0

coin \leftarrow largest denomination that does not exceed *money*

 add *coin* to **change**

money \leftarrow *money* – *coin*

return **change**

Changing Money in Tanzania



40 cents =
25+10+5
Greedy



Changing Money in Tanzania: GreedyChange Fails



40 cents = $25+10+5$ = 20+20
Greedy is not **Optimal**



Recursive Change

Given the denominations 6, 5, and 1, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>									?			

$\text{MinNumCoins}(9) =$

?

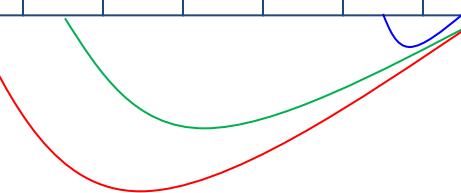
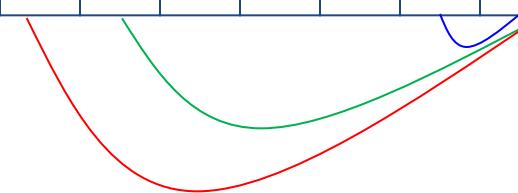
Recursive Change

Given the denominations **6**, **5**, and **1**, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>			?	?				?	?			

$\text{MinNumCoins}(9) =$

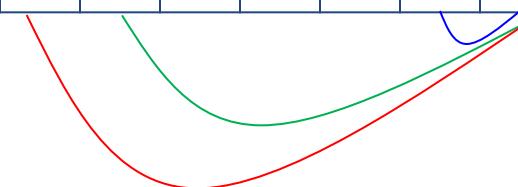
?



Recursive Change

Given the denominations **6**, **5**, and **1**, what is the minimum number of coins needed to change 9 cents?

money	1	2	3	4	5	6	7	8	9	10	11	12
MinNumCoins			?	?				?	?			



$$\begin{aligned} \text{MinNumCoins}(9) = \min \{ & \text{MinNumCoins}(9-6)+1 = \text{MinNumCoins}(3)+1 \\ & \text{MinNumCoins}(9-5)+1 = \\ & \text{MinNumCoins}(4)+1 \\ & \text{MinNumCoins}(9-1)+1 = \text{MinNumCoins}(8)+1 \end{aligned}$$

Recursive Change

Given the denominations 6, 5, and 1, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>			?	?				?	?			

$$\text{MinNumCoins}(3) =$$

$$\text{MinNumCoins}(4) =$$

$$\text{MinNumCoins}(8) =$$

?

Recursive Change

Given the denominations 6, 5, and 1, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>		?	?	?			?	?	?			

$$\text{MinNumCoins}(3) =$$

$$\text{MinNumCoins}(4) =$$

$$\text{MinNumCoins}(8) =$$

?

Recursive Change

Given the denominations 6, 5, and 1, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>		?	?	?			?	?	?			

$$\text{MinNumCoins}(\text{money}-6) \min \{ \text{MinNumCoins}(\text{money}) = \text{MinNumCoins}(\text{money}-5) + 1 \\ \text{MinNumCoins}(\text{money}-1) + 1 \}$$

Recursive Change

Given the denominations 6, 5, and 1, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>		?	?	?			?	?	?			

$$\begin{aligned} \text{MinNumCoins}(\text{money} - \text{coin}_1) &+ 1 \\ \text{MinNumCoins}(\text{money}) = &\dots \\ &\dots + \text{MinNumCoins}(\text{money} - \text{coin}_d) + 1 \end{aligned}$$

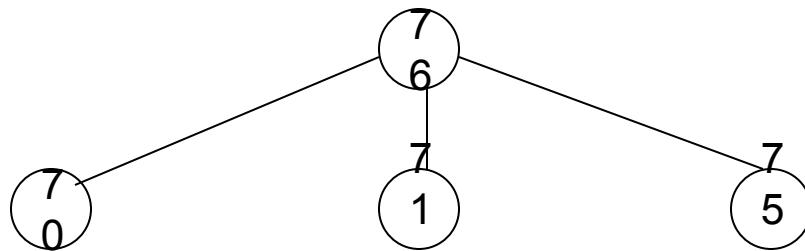
RecursiveChange

```
RecursiveChange(money, coins)
    if money = 0
        return 0
    MinNumCoins ← infinity
    for i ← 1 to |coins|
        if money ≥ coini
            NumCoins ←
RecursiveChange(money-coini, coins)
    if numCoins + 1 < MinNumCoins
        MinNumCoins ← numCoins + 1
return MinNumCoins
```

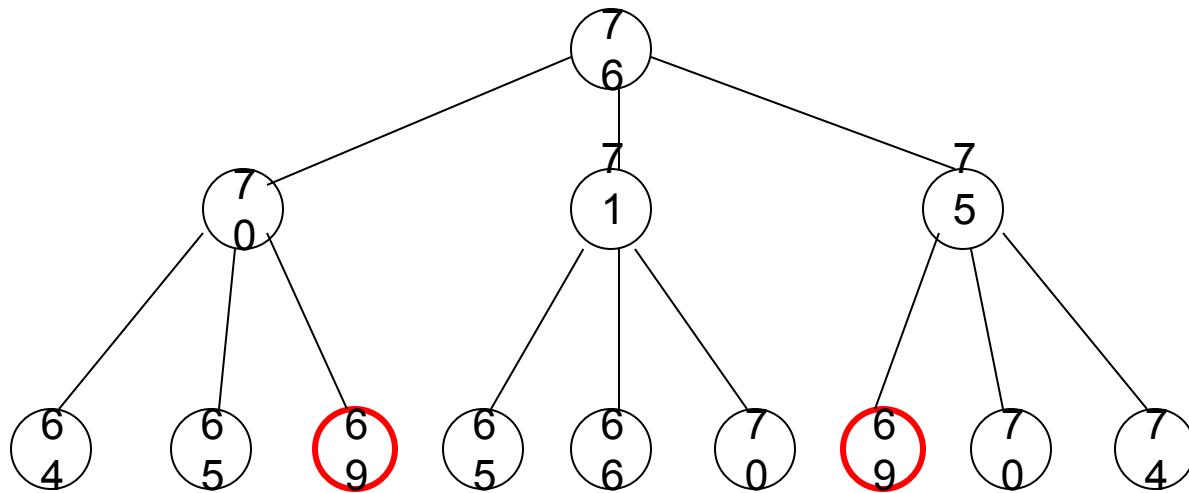
How Fast is the RecursiveChange?

7
6

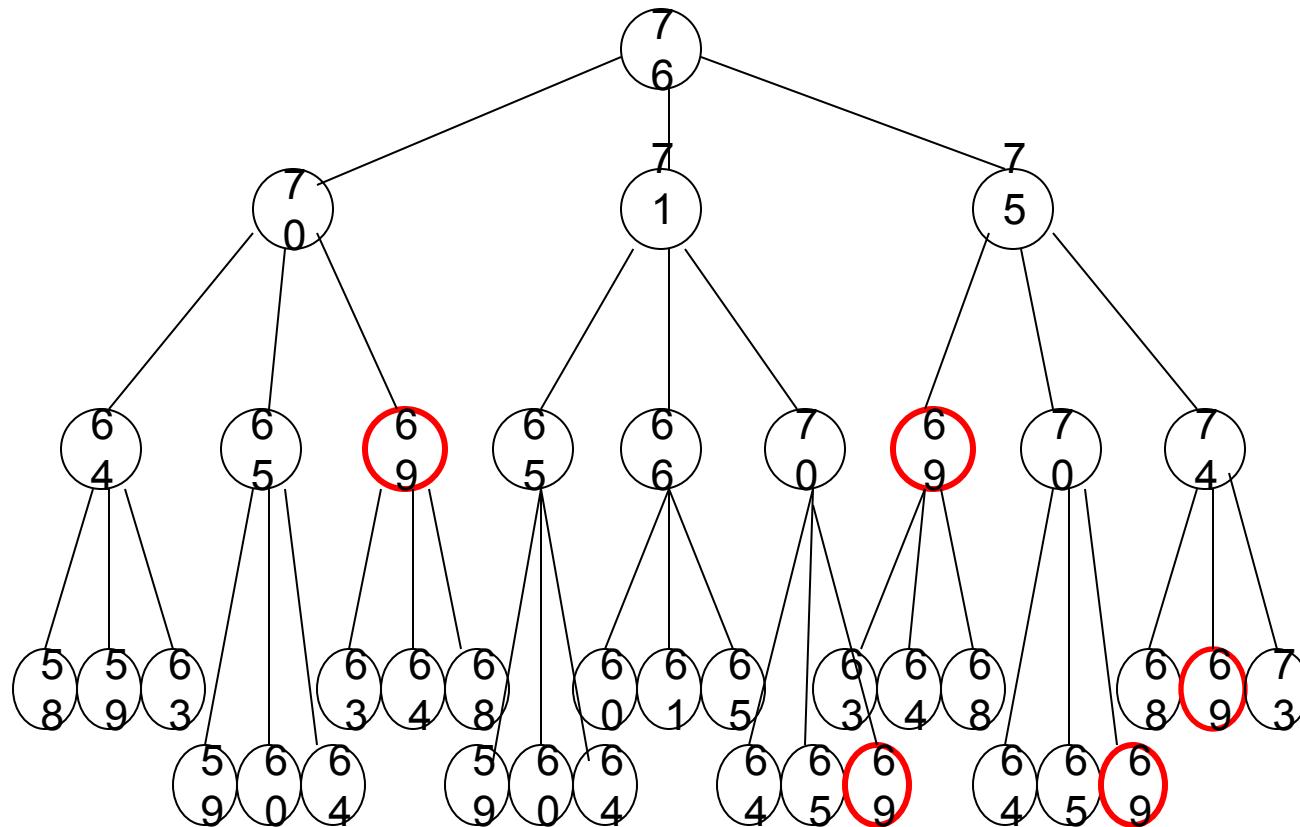
The Recursive Tree



The Recursive Tree



The Recursive Tree



the optimal coin combination for 69 cents is computed **6** times!

the optimal coin combination for 30 cents is computed **trillions** of times!

Changing Money by Dynamic Programming

Hint. Wouldn't it be nice to know all the values of

$$\text{MinNumCoins}(\text{money} - \text{coin}_i)$$

by the time we need to compute

$$\text{MinNumCoins}(\text{money})?$$



Richard
Bellman

Instead of the time-consuming calls:

$$\text{RecursiveChange}(\text{money}-\text{coin}_i, \text{Coins}, d)$$

we would simply look up the values of

$$\text{MinNumCoins}(\text{money} - \text{coin}_i)$$

Changing Money by Dynamic Programming

Given the denominations **6**, **5**, and **1**, what is the minimum number of coins needed to change 9 cents?

<i>money</i>	0	1	2	3	4	5	6	7	8	9	10	11	12
<i>MinNumCoins</i>	0	1	2	3	4	1	1	2	3	?			

The diagram illustrates the dynamic programming solution. It shows three overlapping bell-shaped curves representing the coin denominations: a red curve for 1 cent, a green curve for 5 cents, and a blue curve for 6 cents. These curves are centered at their respective values (1, 5, and 6) and overlap significantly. Below the curves, a series of colored lines connect the points on the curves to the corresponding cells in the 'MinNumCoins' row, indicating the paths of coin selection for each value from 0 to 9.

DPChange

DPChange(*money, coins*)

MinNumCoins(0) ← 0

for *m ← 1* to *money*

MinNumCoins(m) ← infinity

for *i ← 1* to $|coins|$

if $m \geq coin_i$

if $MinNumCoins(m - coin_i) + 1 < MinNumCoins(m)$

$MinNumCoins(m) \leftarrow MinNumCoins(m - coin_i) + 1$

return *MinNumCoins(money)*

“Programming” in “Dynamic Programming” Has Nothing to Do with Programming!

Richard Bellman developed this idea in 1950s working on an Air Force project.

At that time, his approach seemed completely impractical.

He wanted to hide that he is really doing math from the Secretary of Defense.



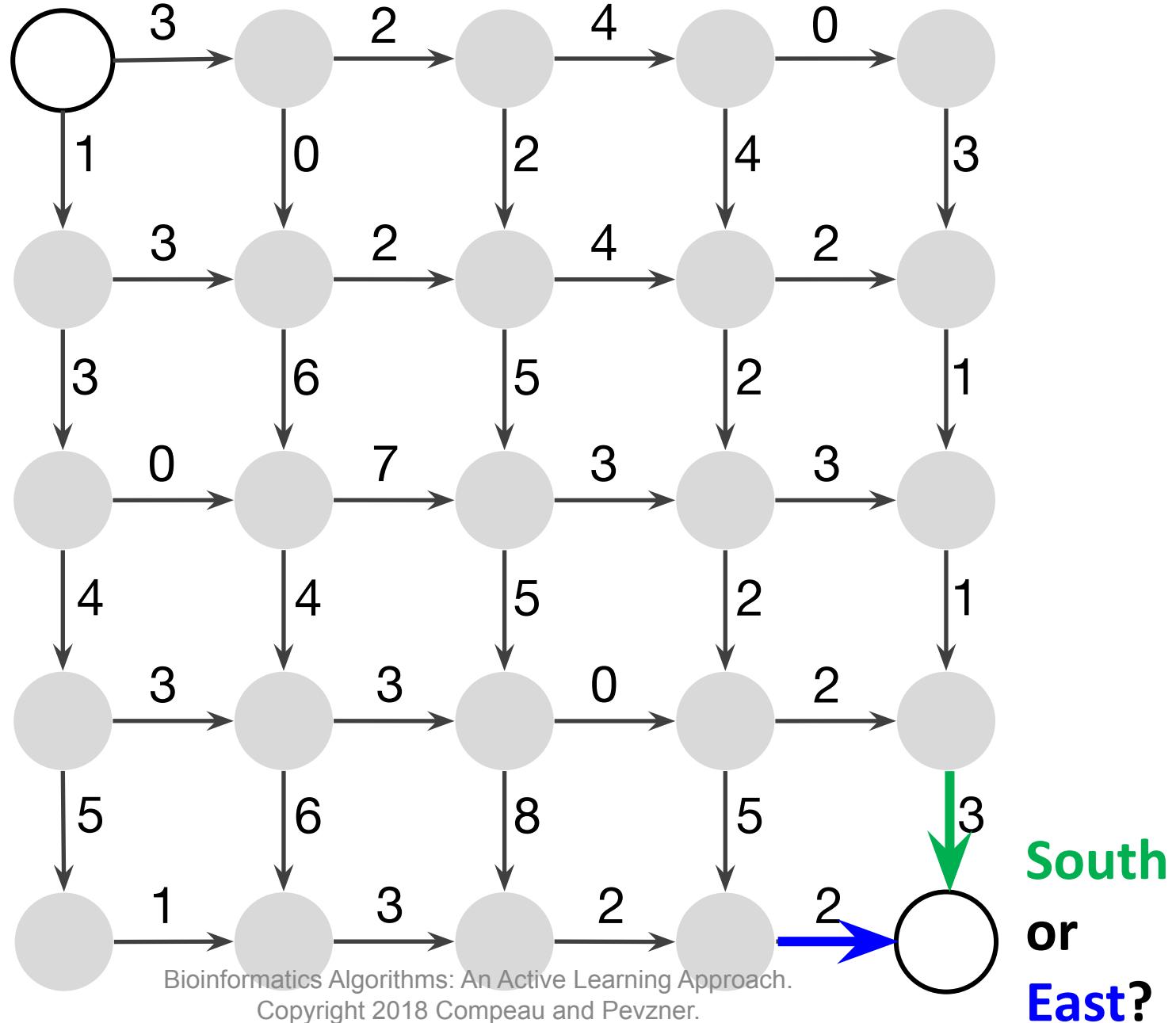
Richard
Bellman

“...What name could I choose? I was interested in planning but planning, is not a good word for various reasons. I decided therefore to use the word, “programming” and I wanted to get across the idea that this was dynamic. **It was something not even a Congressman could object to. So I used it as an umbrella for my activities.**”

How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- **Dynamic Programming and Backtracking Pointers**
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

There are
only 2 ways
to arrive to
the sink:
by moving
South ↓
or by moving
East →



South or East?

SouthOrEast(n,m)

if $n=0$ and $m=0$

return 0

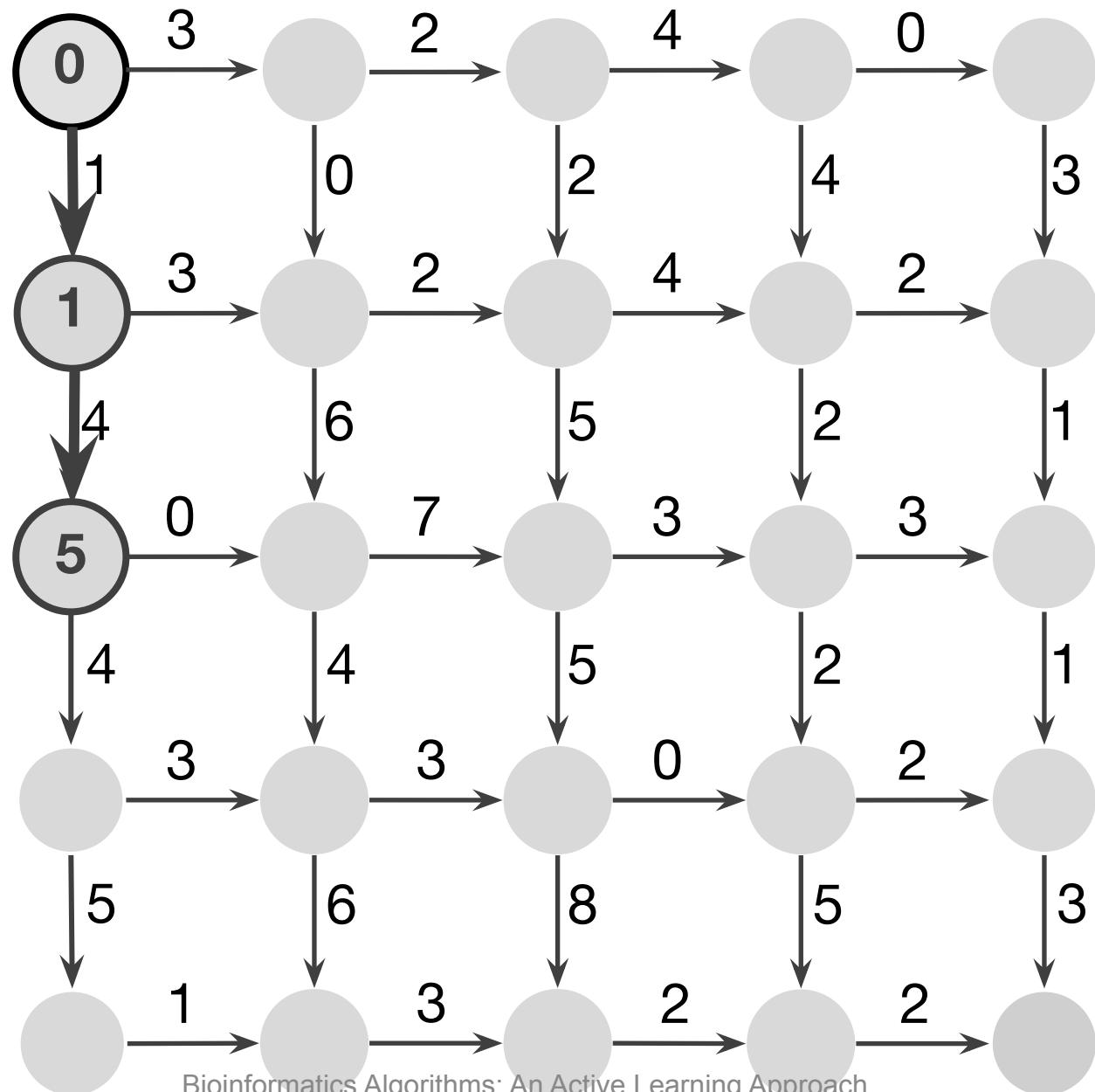
if $n>0$ and $m>0$

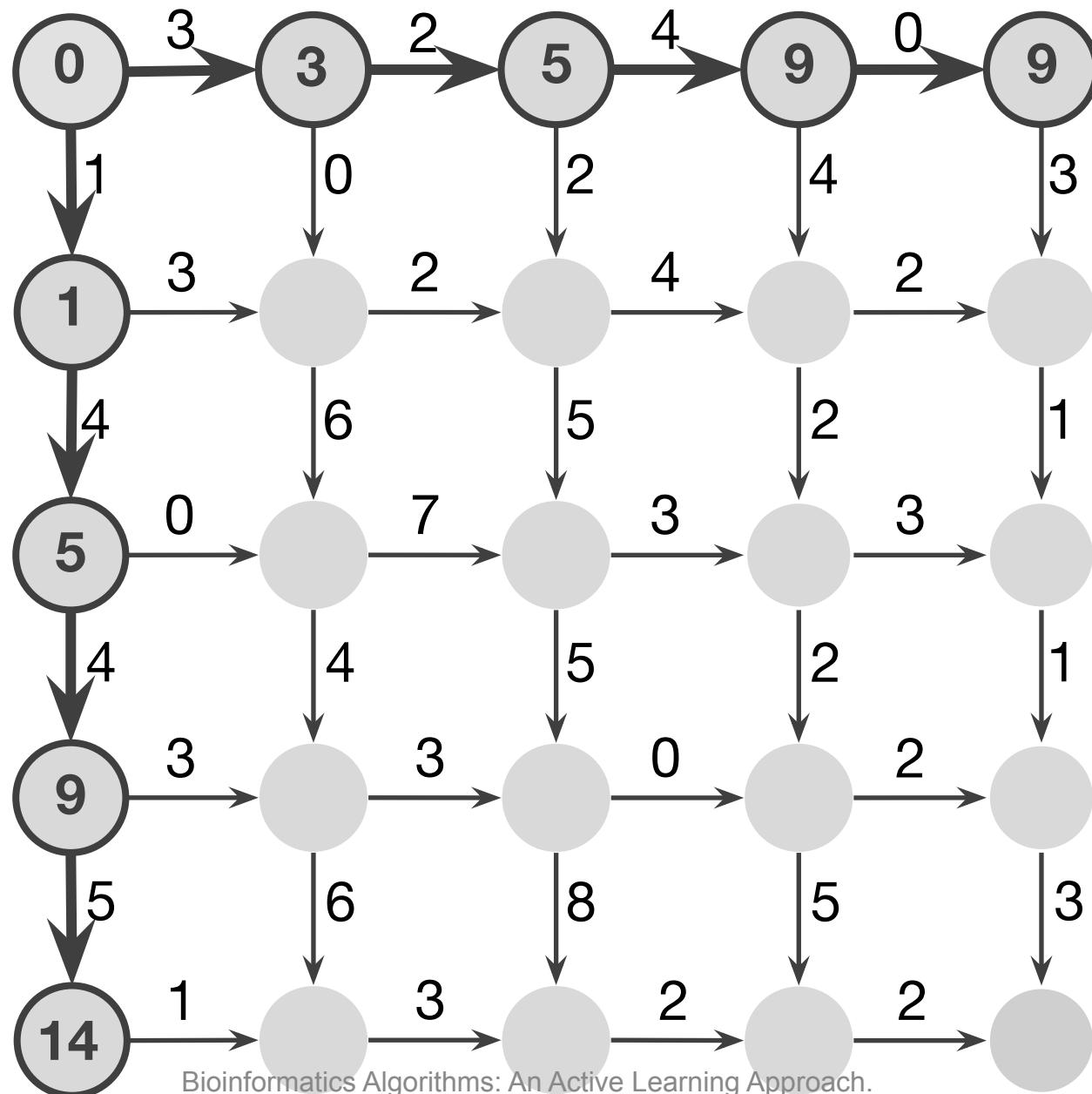
$x \square$ **SouthOrEast($n-1,m$)**+weight of edge “

$y \square$ **SouthOrEast($n,m-1$)**+ weight of edge “

return $\max\{x,y\}$

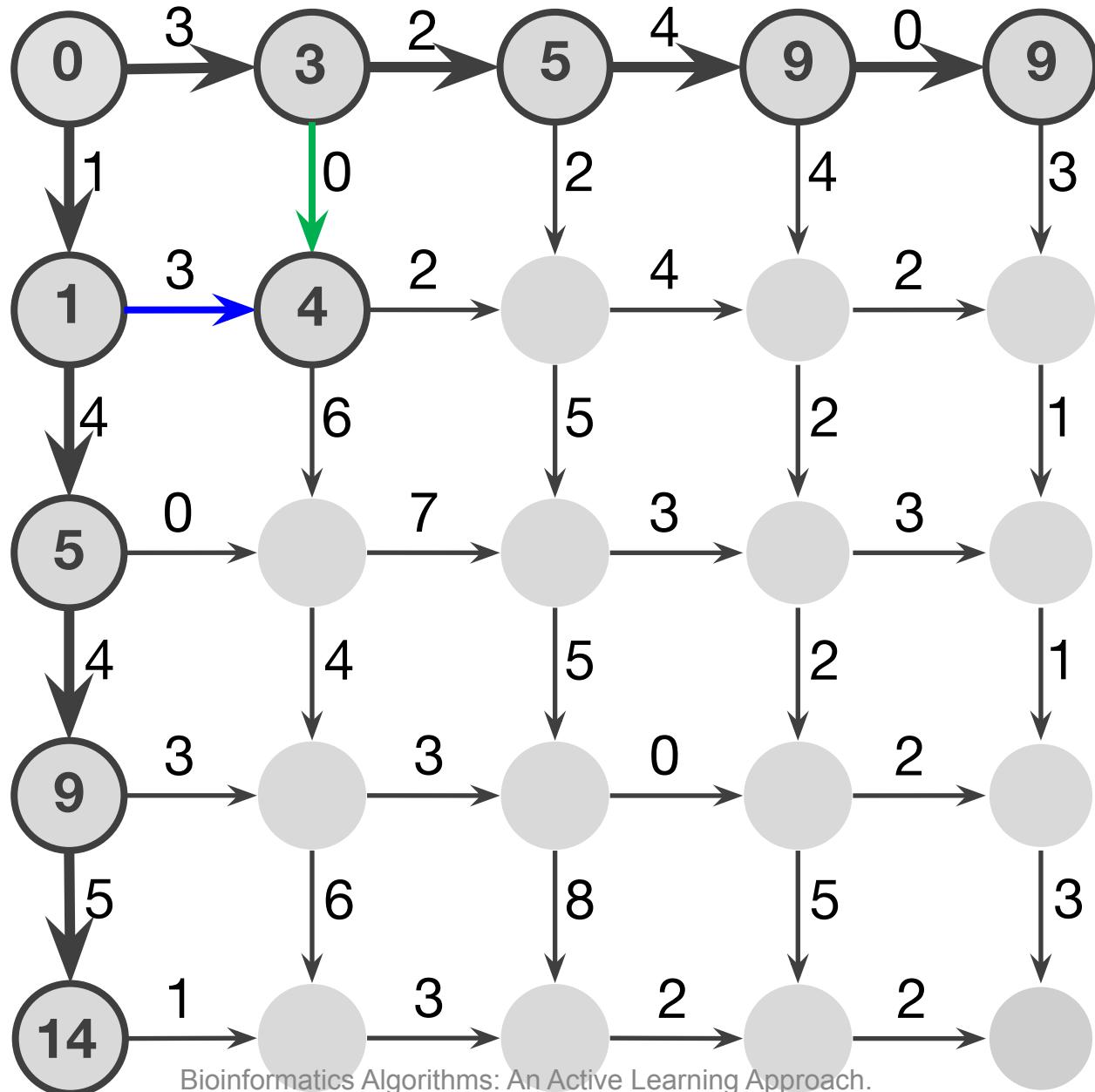
return -infinity



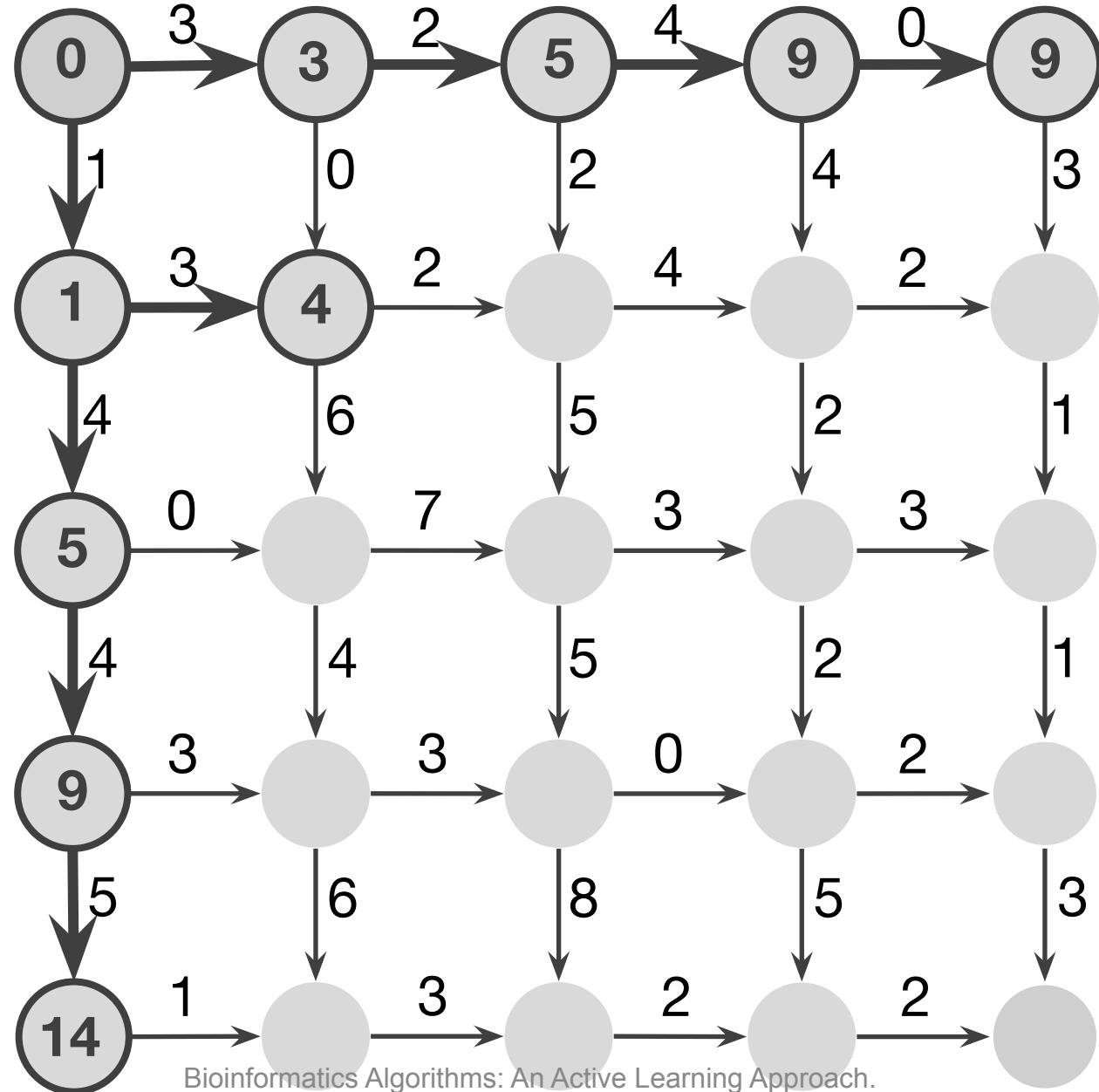


**South
or
East?**

$$1+3 > 3+0$$

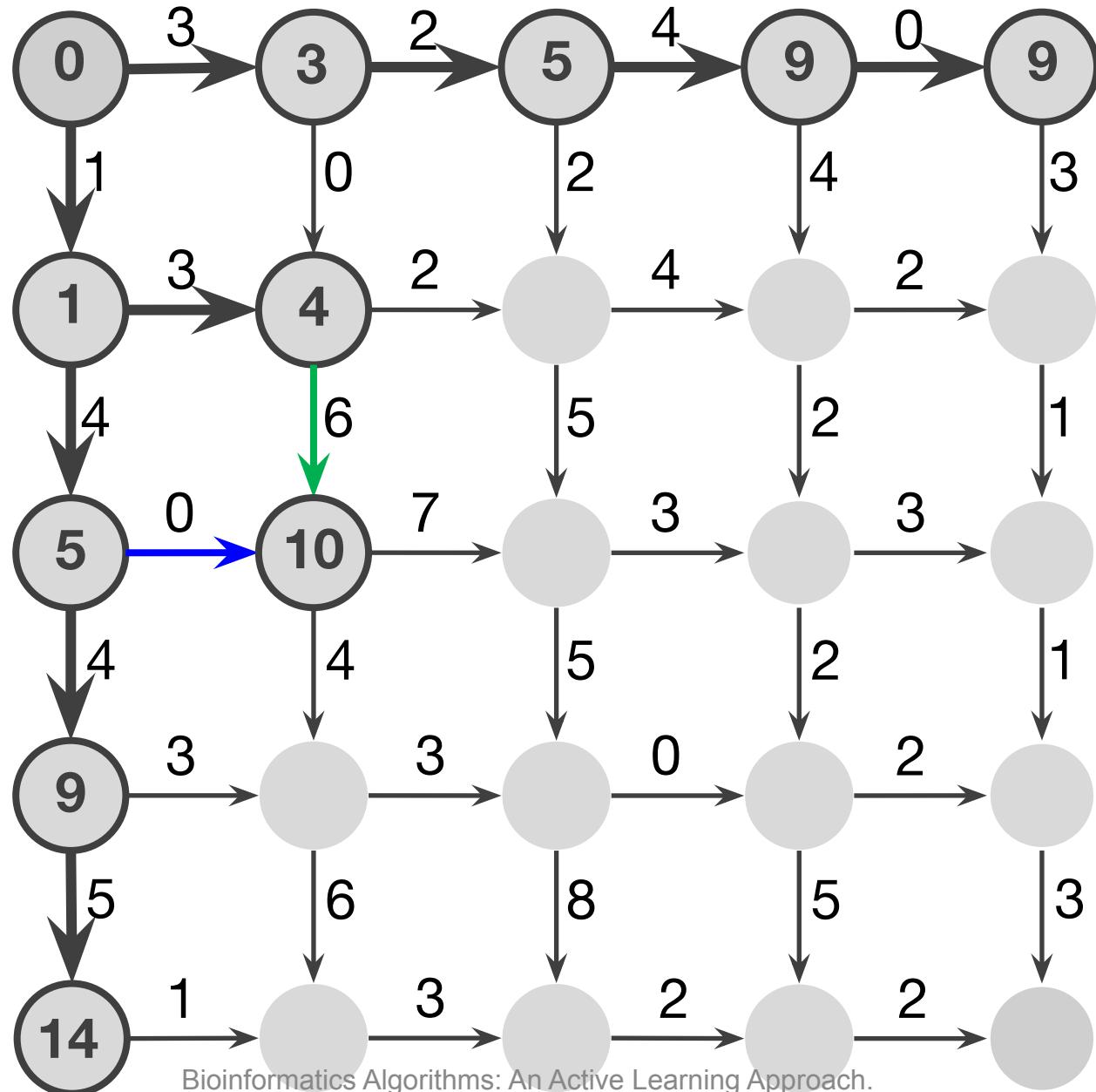


We arrived
to (1,1)
by the **bold**
edge:

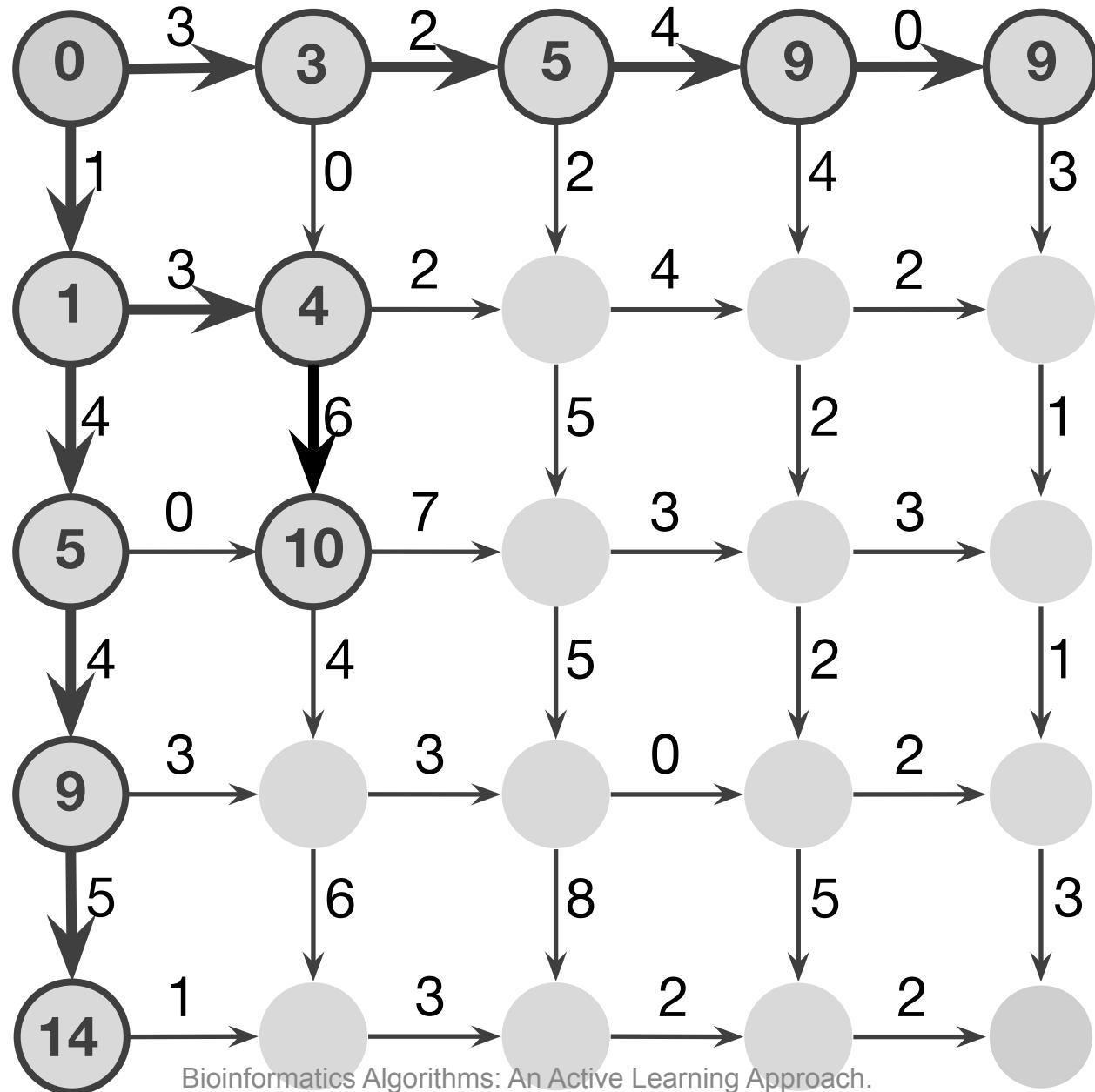


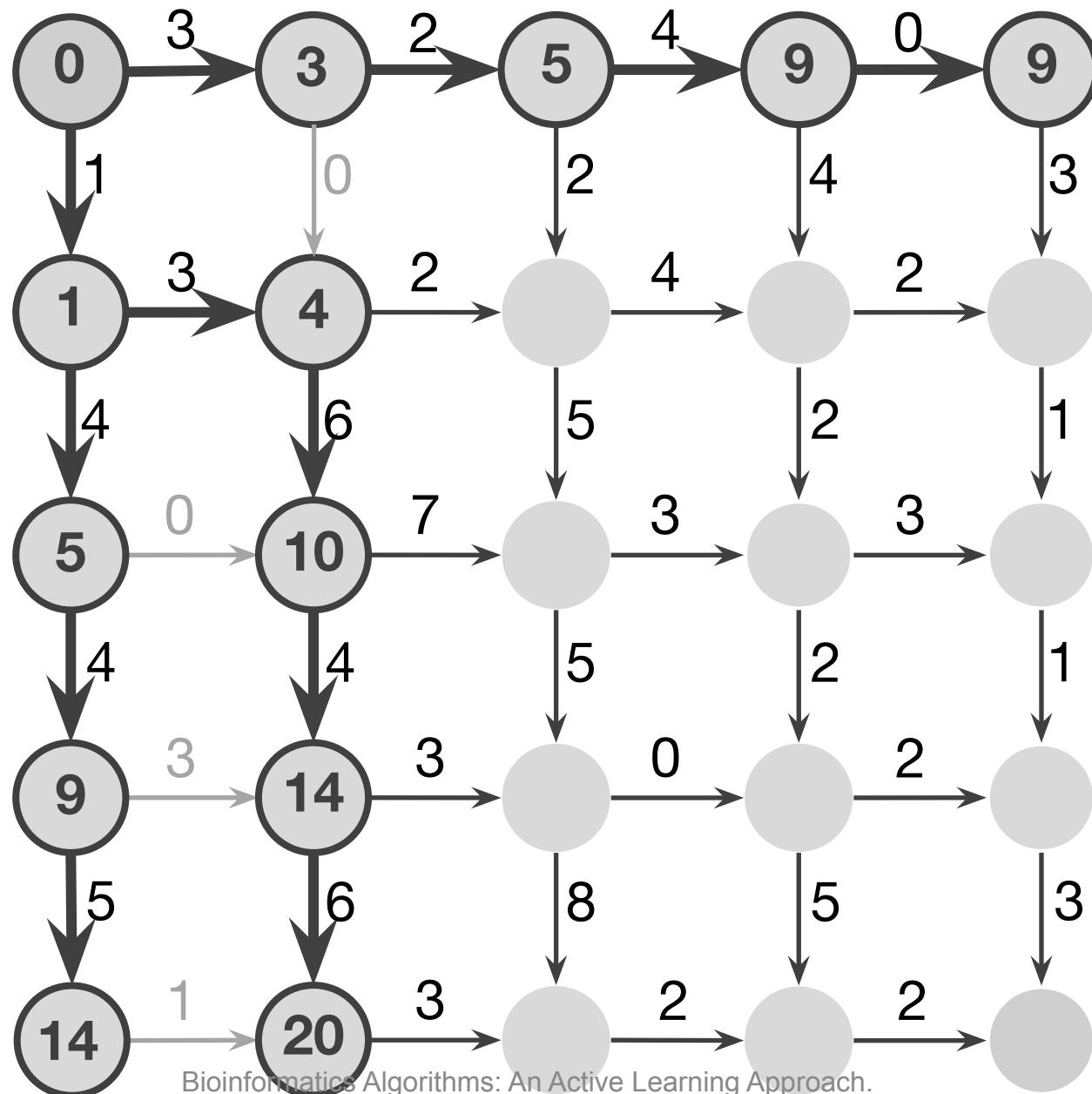
**South
or
East?**

$5+0 < 4+6$



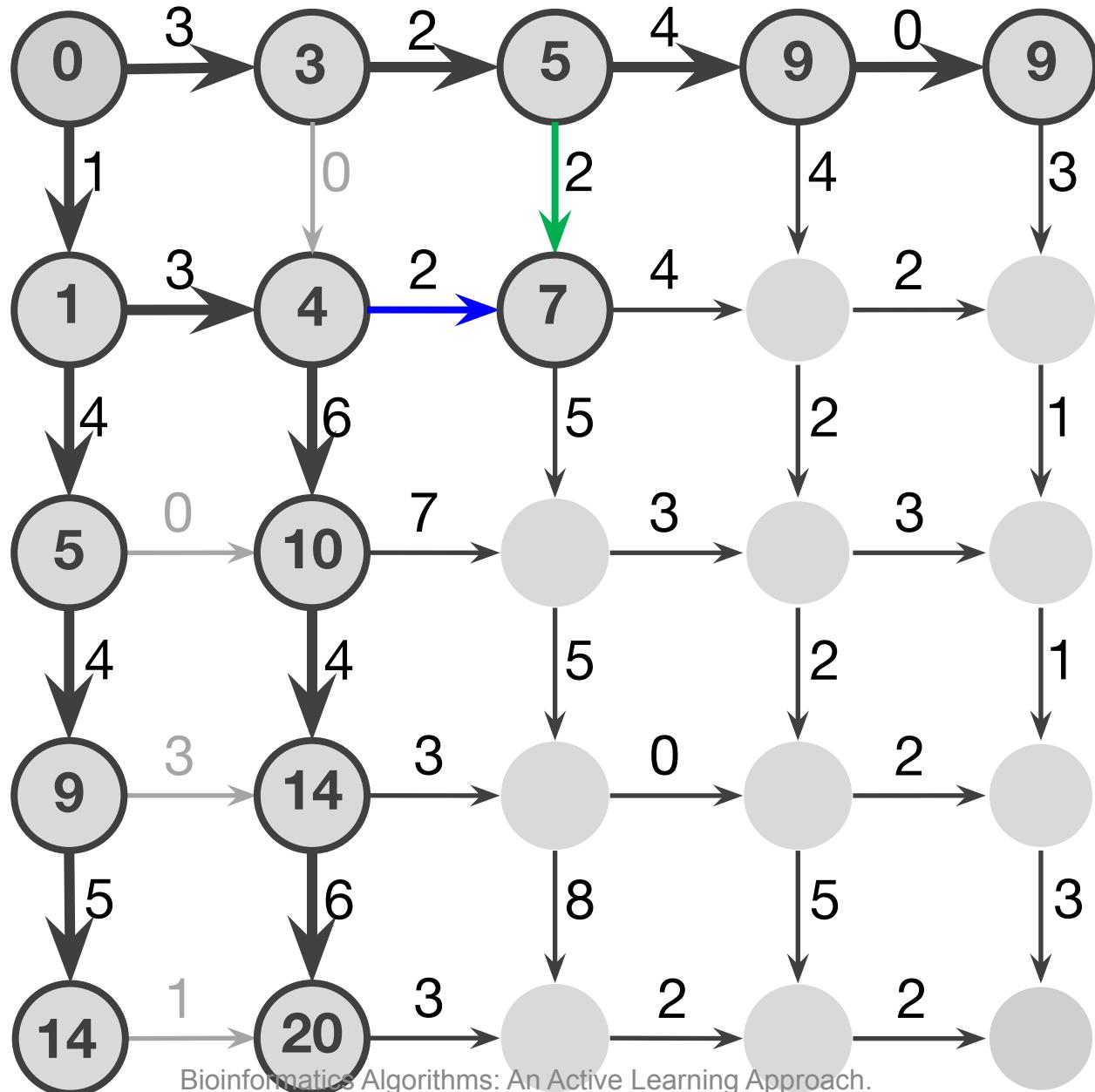
We arrived
to (2,1)
by the **bold**
edge:





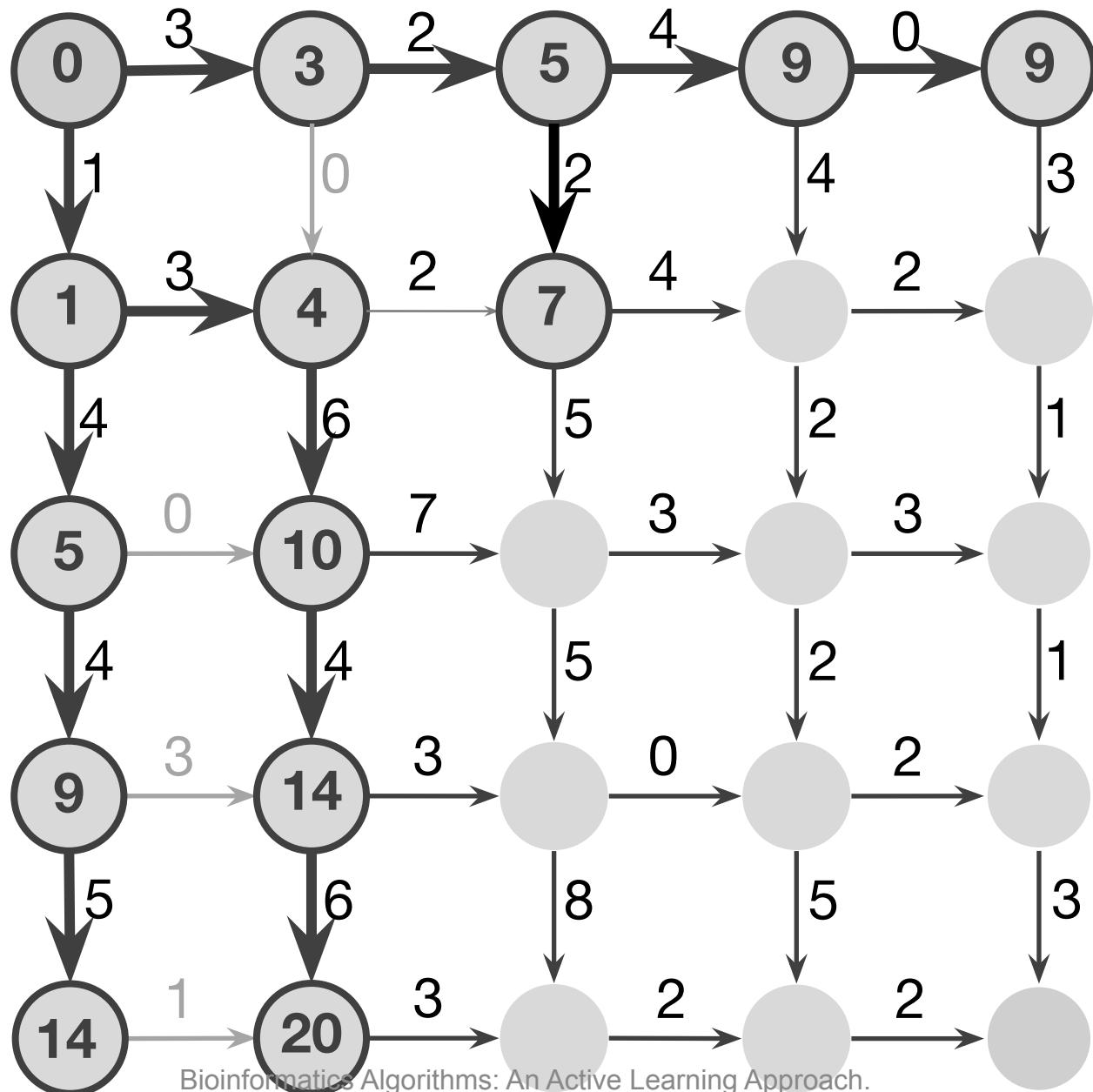
**South
or
East?**

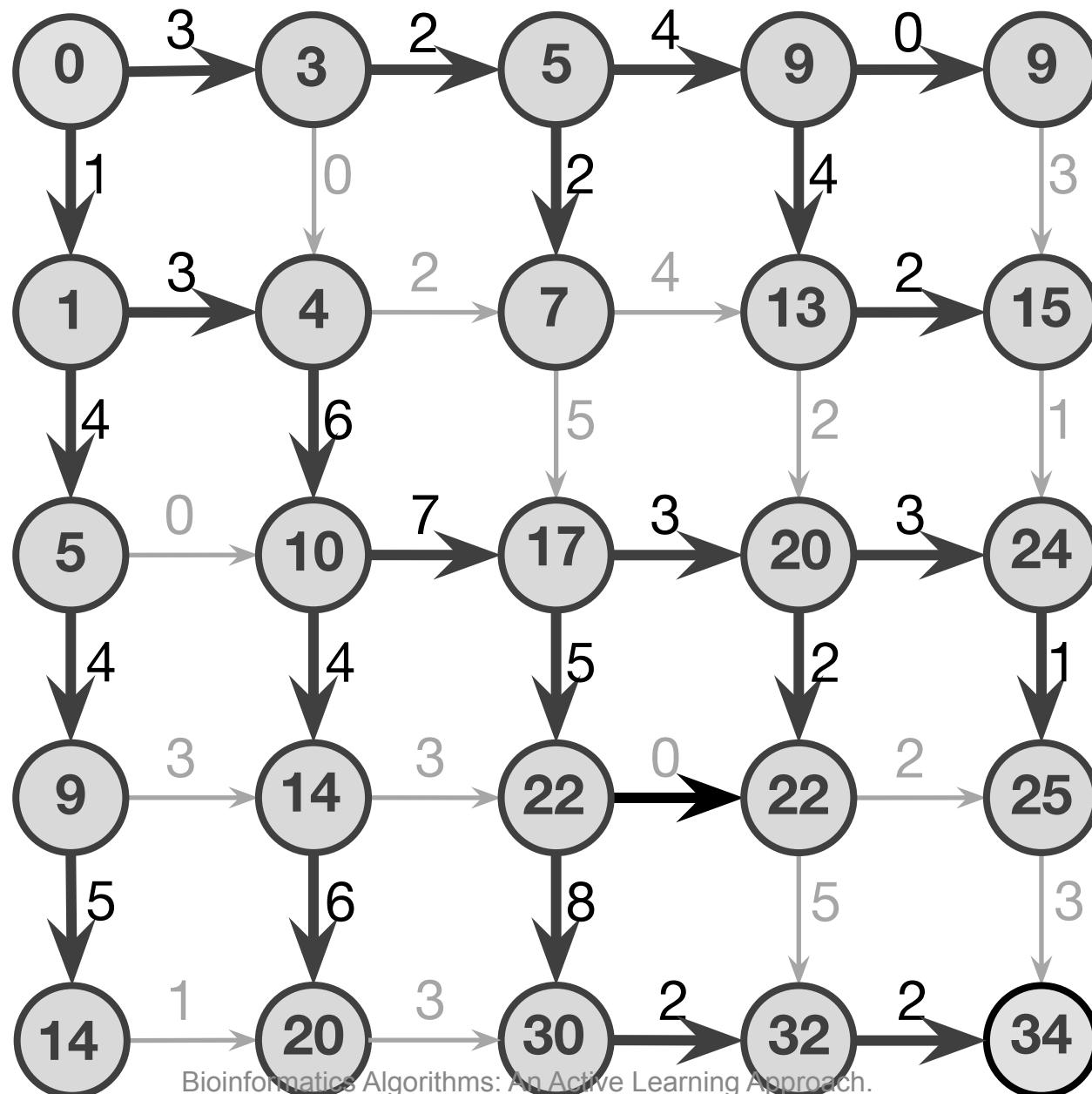
$$5+2 > 4+2$$



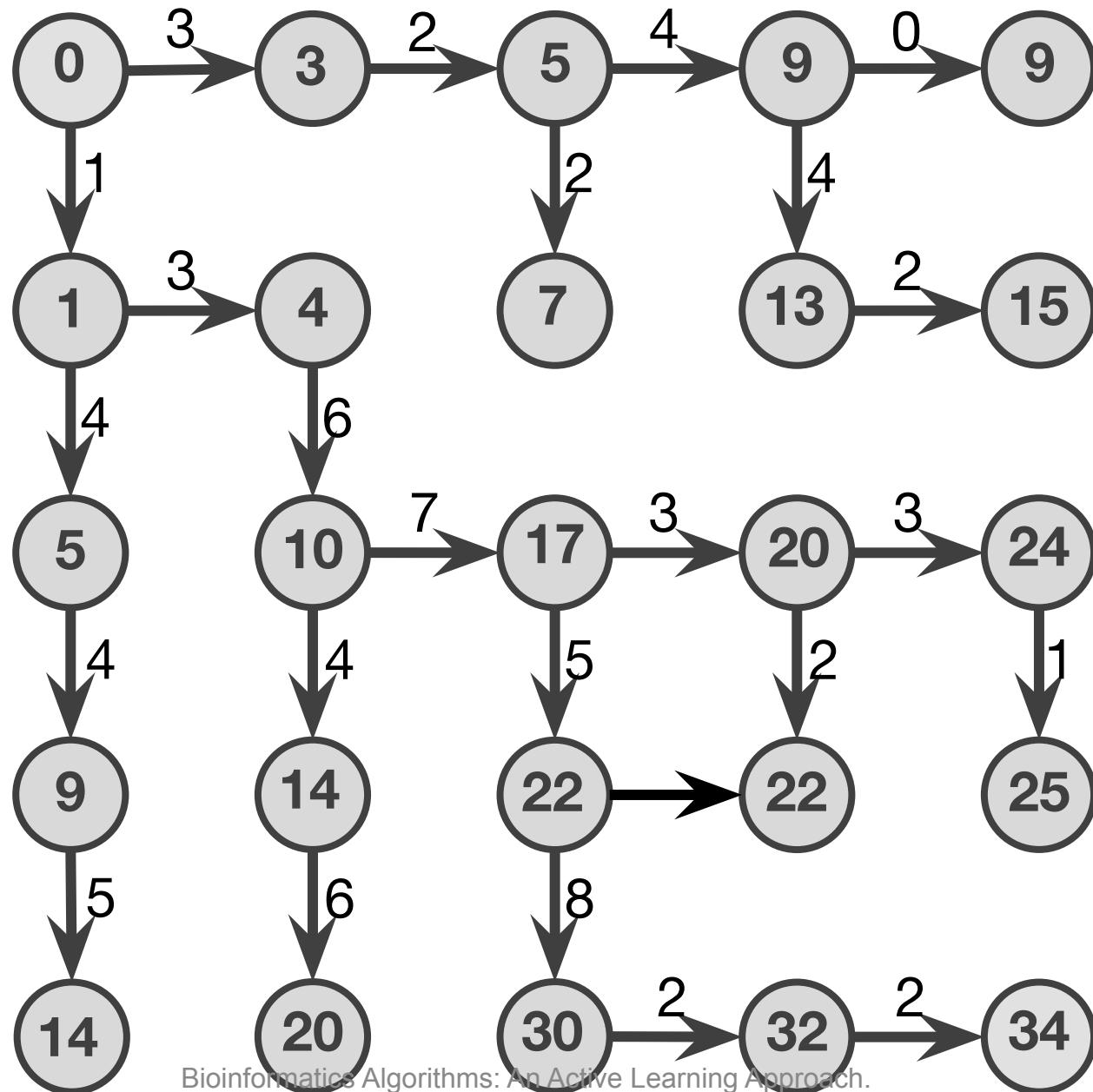
**South
or
East?**

5+2 > 4+2





**Backtracking
pointers:**
the best way
to get to
each node



Dynamic Programming Recurrence

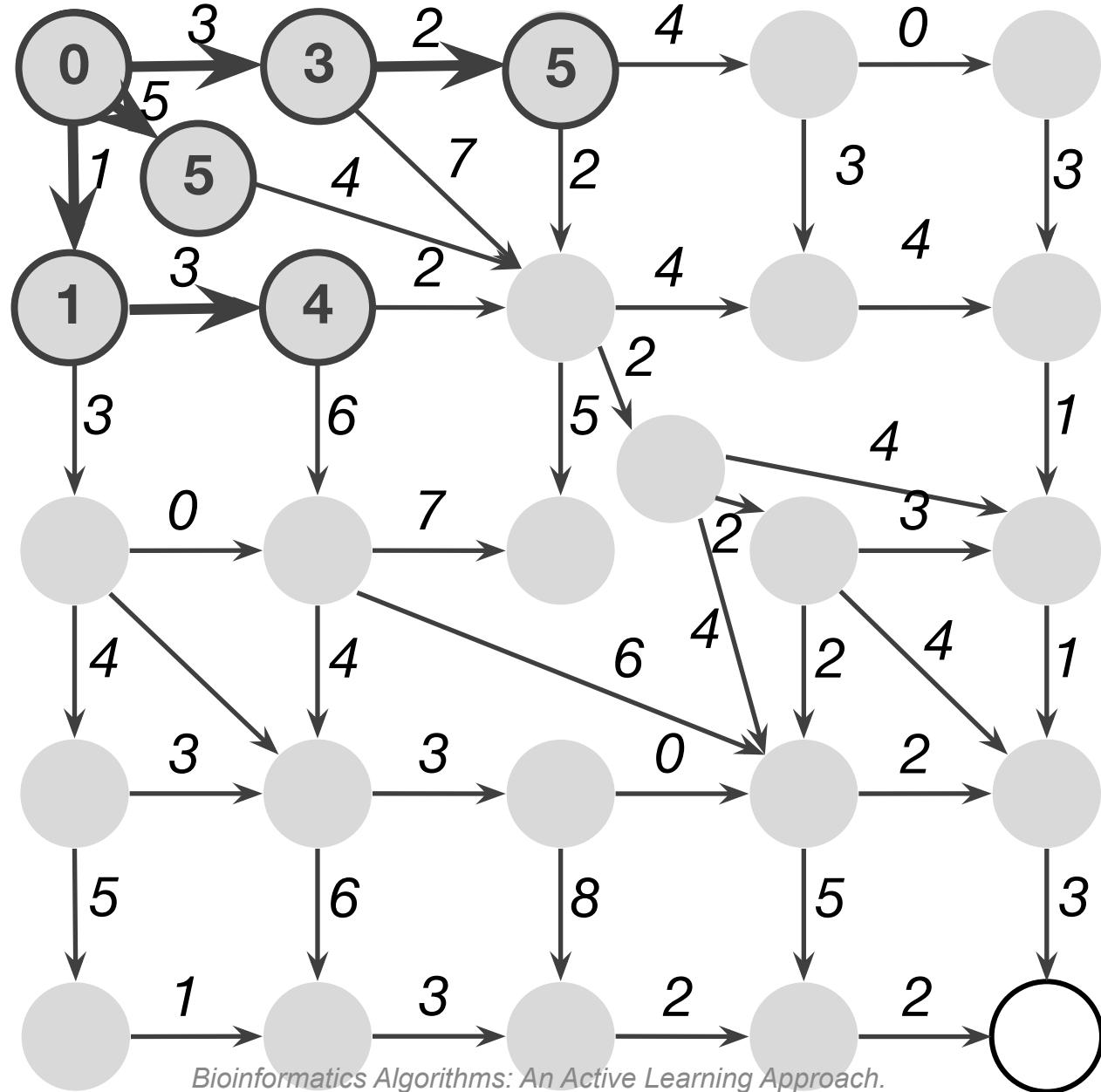
$s_{i,j}$: the length of a longest path from (0,0) to (i,j)

$$s_{i,j} = \max \{ \begin{array}{l} s_{i-1,j} + \text{weight of edge "↓" into } (i,j) \\ s_{i,j-1} + \text{weight of edge "→" into } (i,j) \end{array}$$

How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- **From Manhattan to the Alignment Graph**
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

How does
the
recurrence
change for
this graph?



$$s_a = \max_{\text{all predecessors } b \text{ of node } a} \{s_b + \text{weight of edge from } b \text{ to } a\}$$

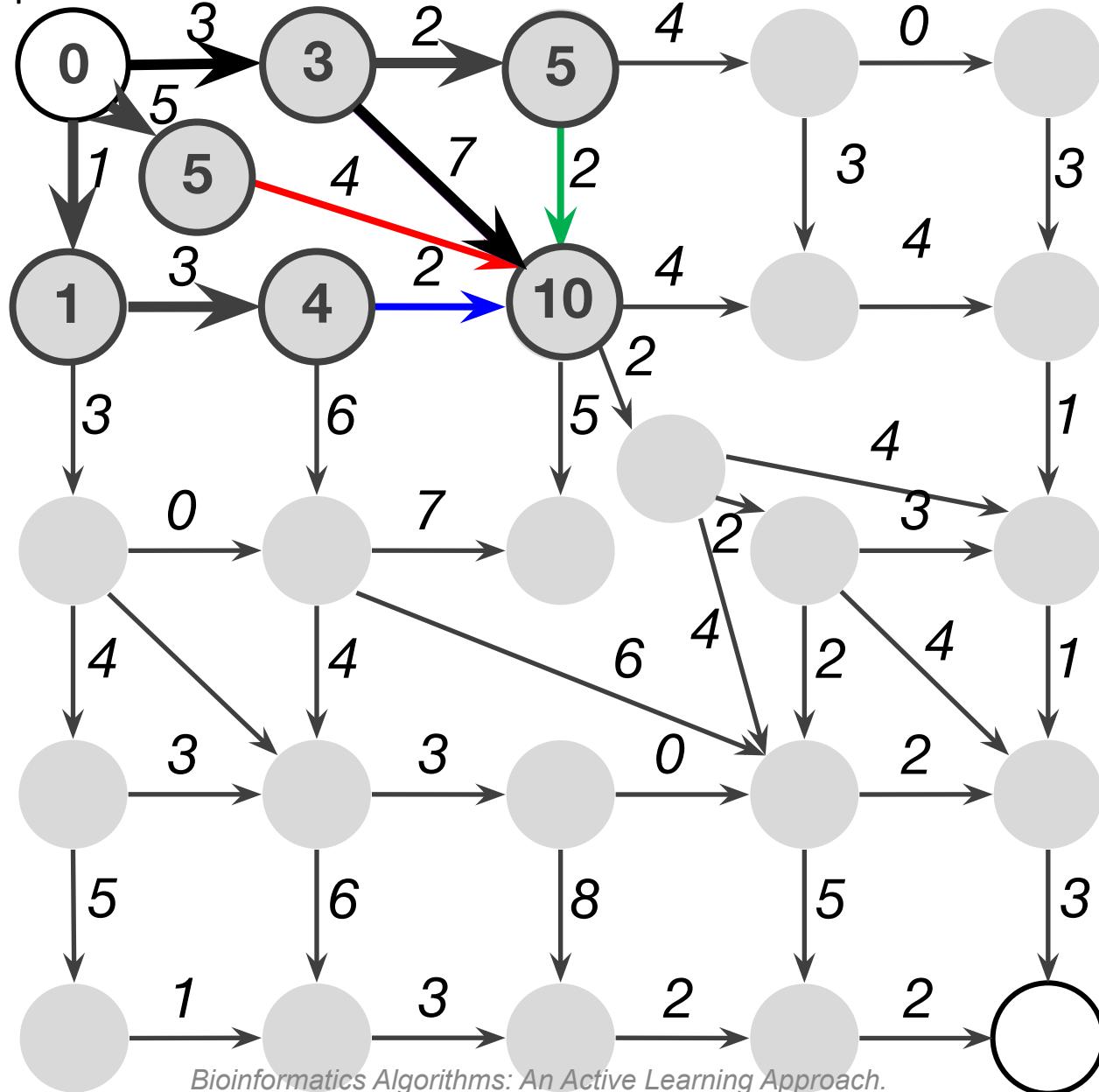
4 choices:

$5 + 2$

$3 + 7$

$5 + 4$

$4 + 2$



$$s_a = \max_{\text{all predecessors } b \text{ of node } a} \{s_b + \text{weight of edge from } b \text{ to } a\}$$

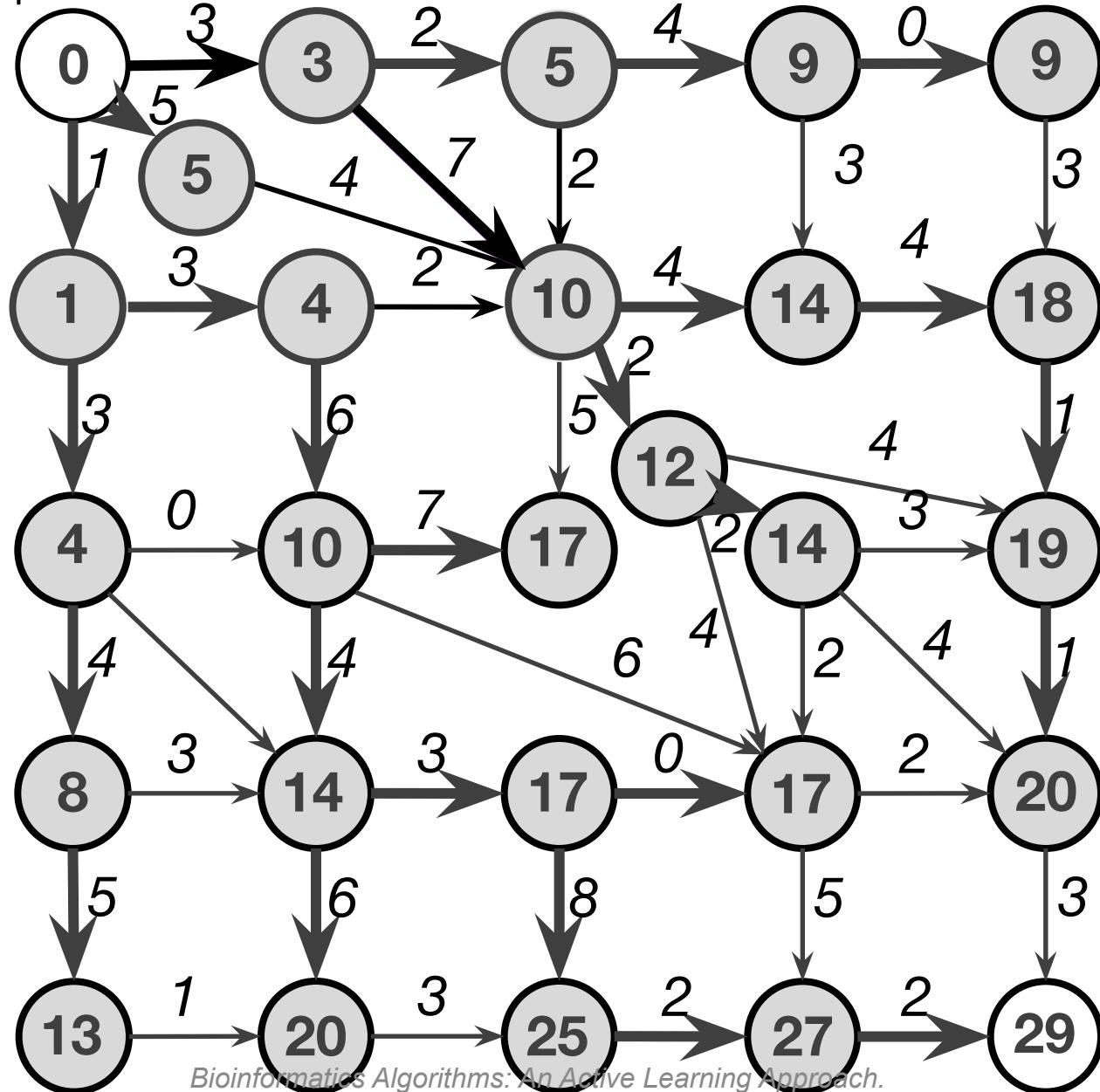
4 choices:

$5 + 2$

$3 + 7$

$5 + 4$

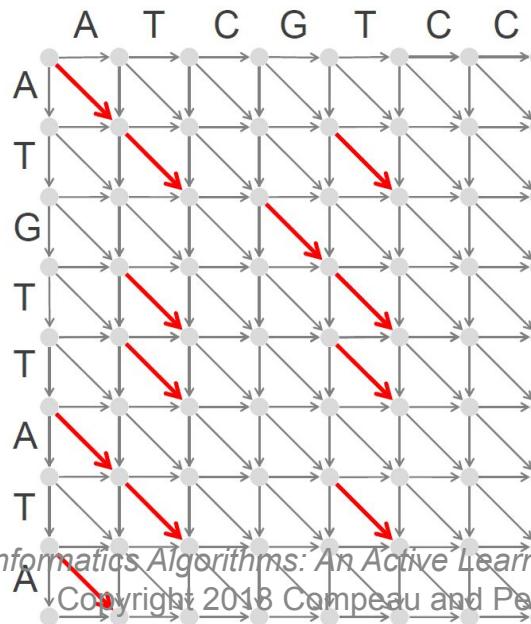
$4 + 2$



Dynamic Programming Recurrence for the Alignment Graph

$s_{i,j}$: the length of a longest path from $(0,0)$ to (i,j)

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge "↓" into } (i,j) \\ s_{i,j-1} + \text{weight of edge "→" into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge "↖" into } (i,j) \end{cases}$$

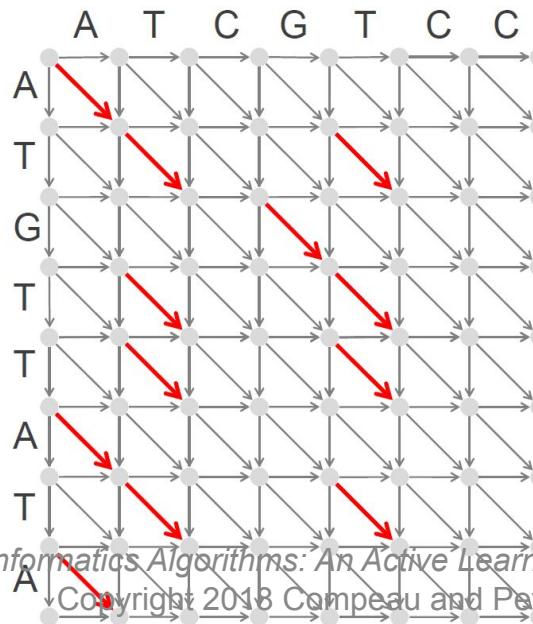


red edges \searrow – weight 1
other edges – weight 0

Dynamic Programming Recurrence for the Longest Common Subsequence Problem

$s_{i,j}$: the length of a longest path from (0,0) to (i,j)

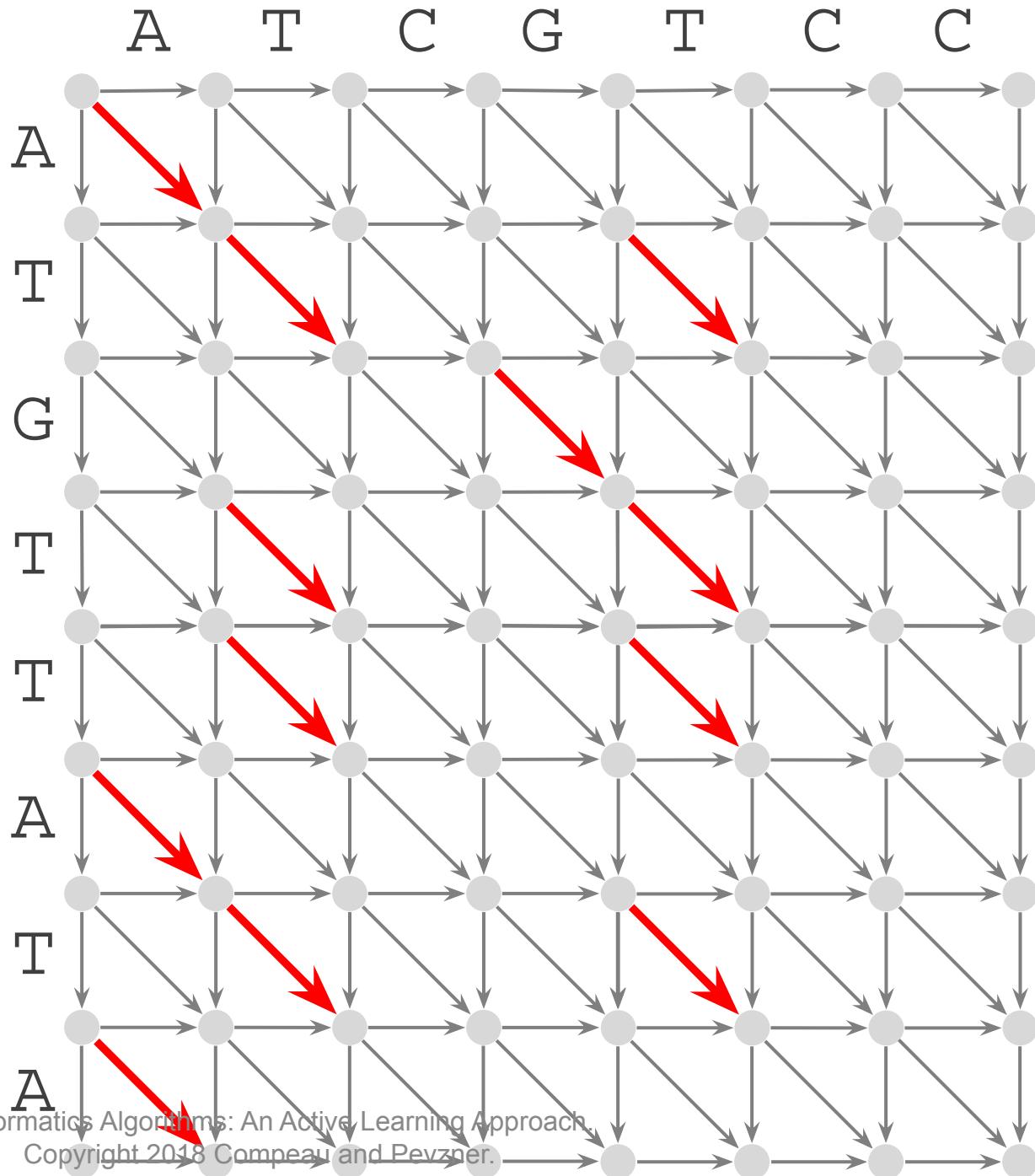
$$s_{i,j} = \max \begin{cases} s_{i-1,j} + 0 \\ s_{i,j-1} + 0 \\ s_{i-1,j-1} + 1, \text{ if } v_i = w_j \end{cases}$$



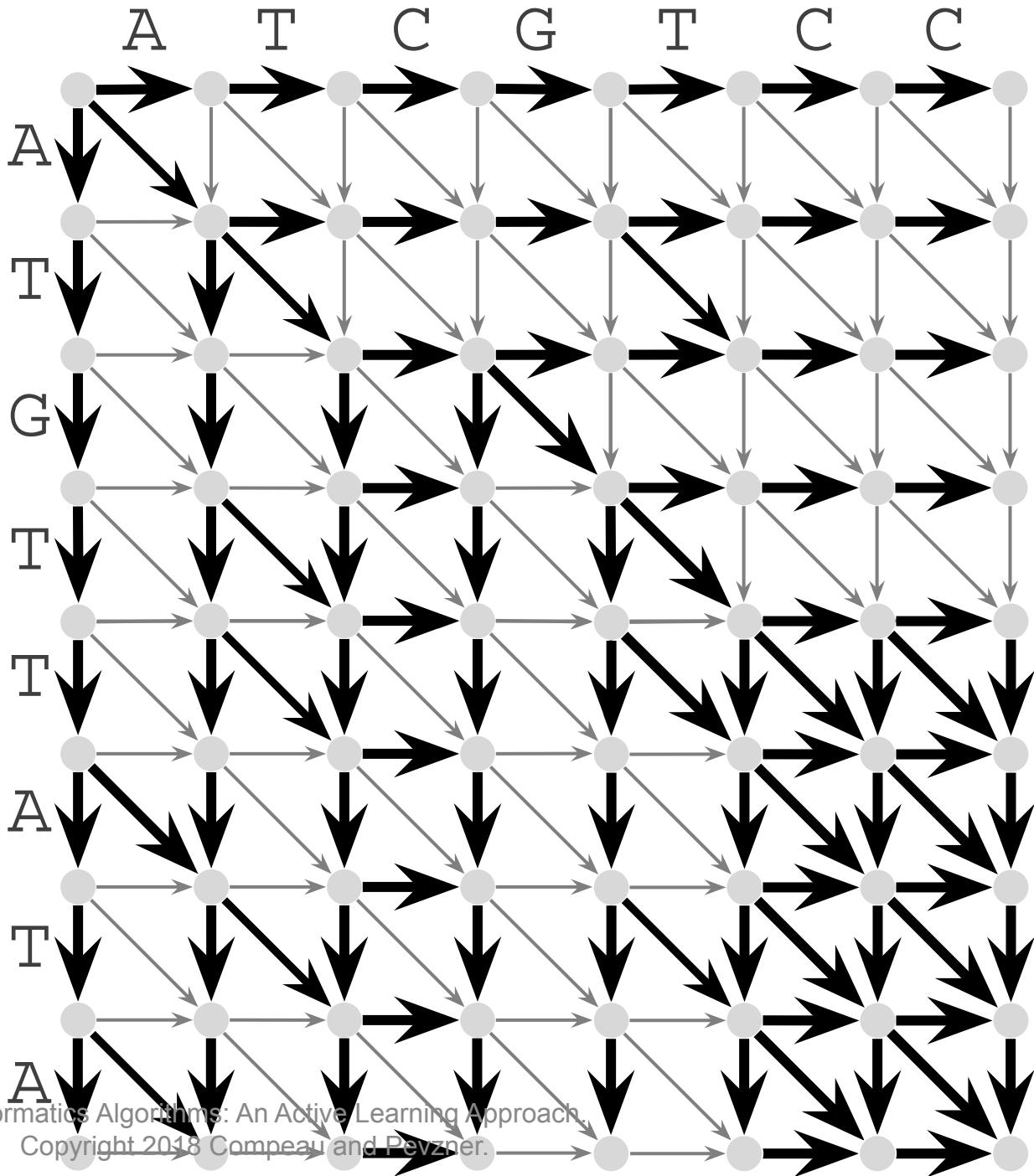
red edges \rightarrow – weight 1
other edges – weight 0

backtracking pointers
for the Longest
Common Subsequence

red edges → – weight 1
other edges – weight 0



backtracking pointers for the Longest Common Subsequence

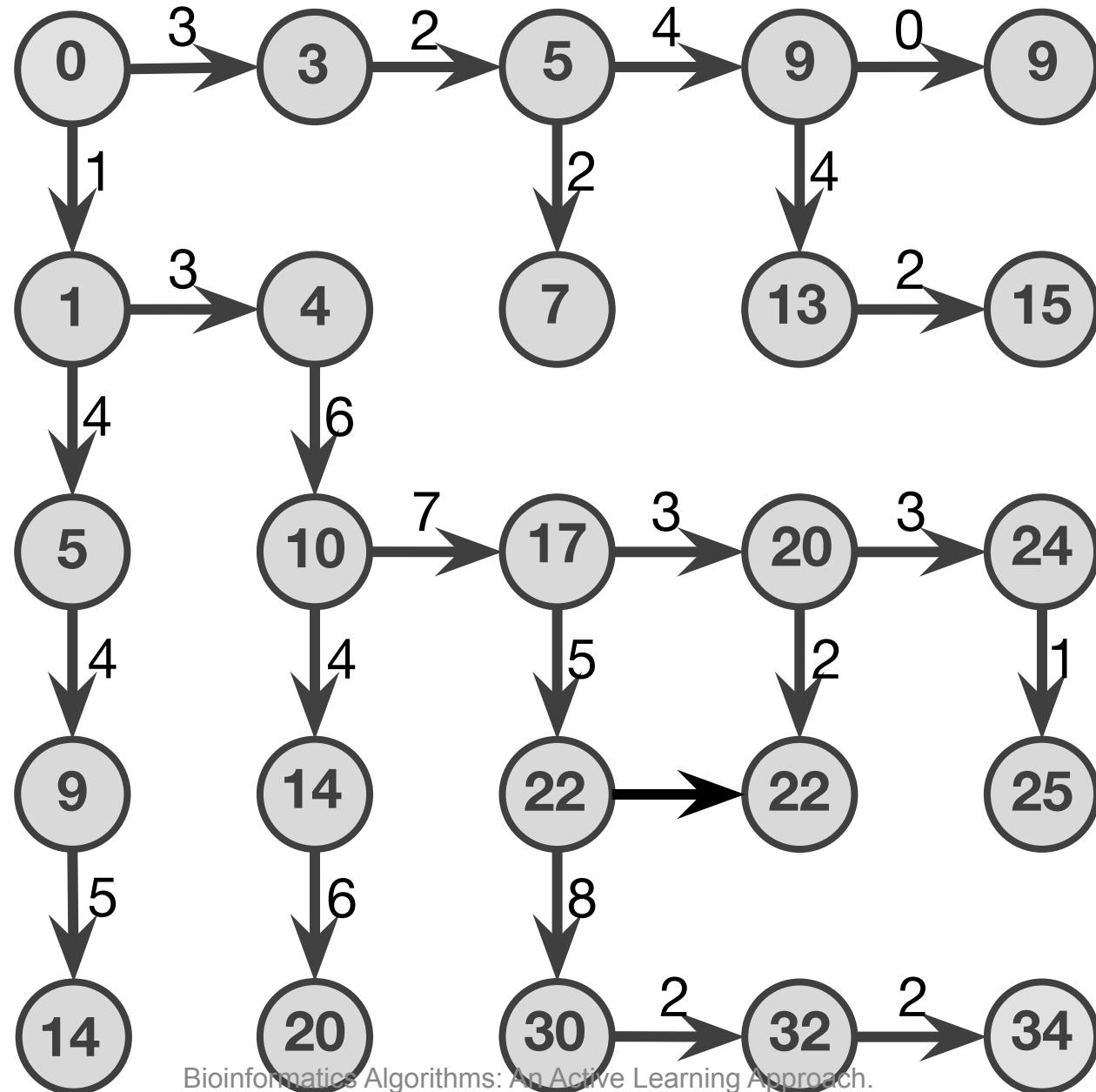


Computing Backtracking Pointers

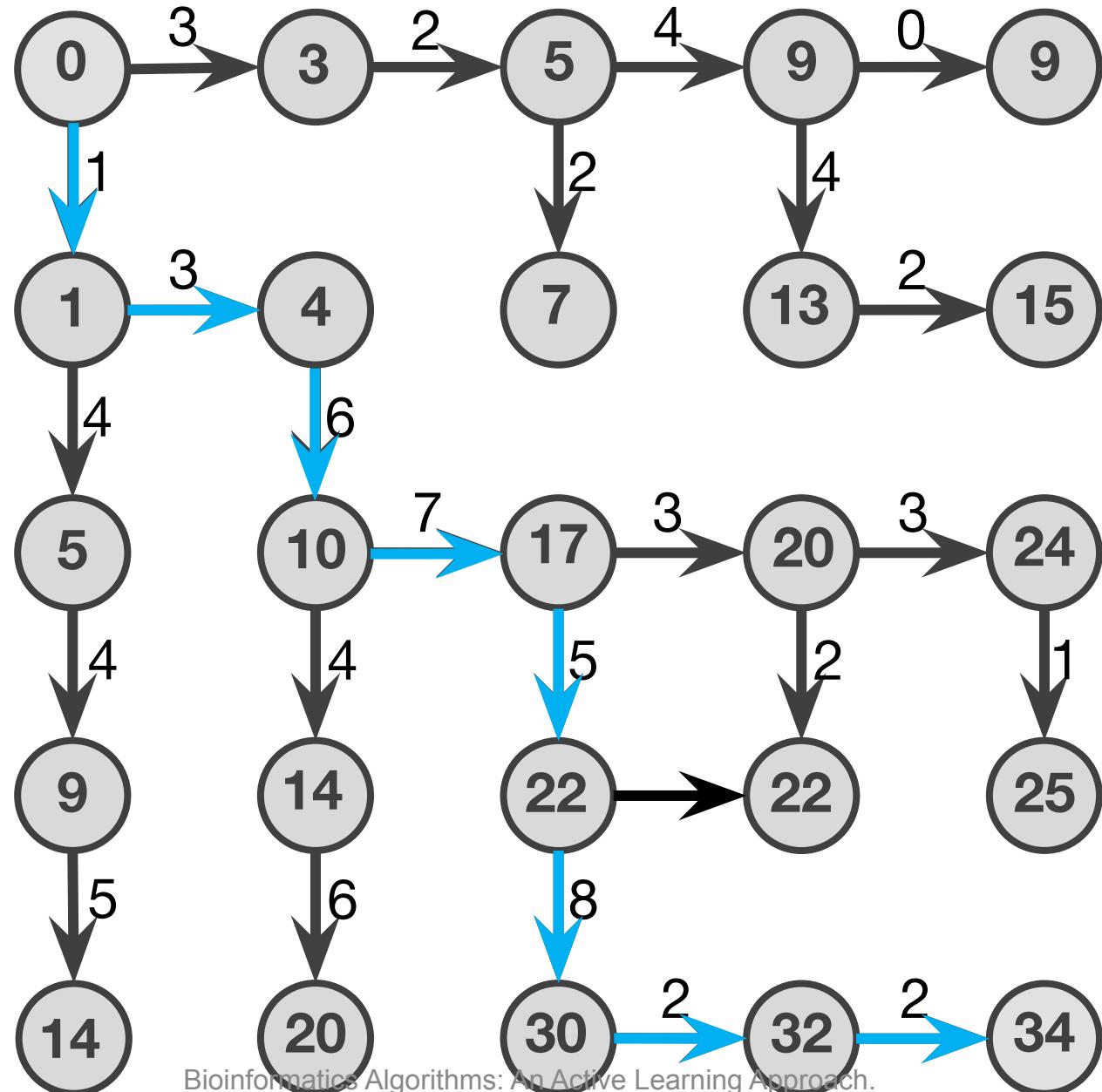
$$s_{i,j} \leftarrow \max\{ \begin{array}{l} s_{i,j-1} + 0 \\ s_{i-1,j} + 0 \\ s_{i-1,j-1} + 1, \text{ if } v_i = w_j \end{array} \}$$

$$\text{backtrack}_{i,j} \leftarrow \{ \begin{array}{l} " \rightarrow ", \text{ if } s_{i,j} = s_{i,j-1} \\ " \downarrow ", \text{ if } s_{i,j} = s_{i-1,j} \\ " \searrow ", \text{ if } s_{i,j} = s_{i-1,j-1} + 1 \end{array} \}$$

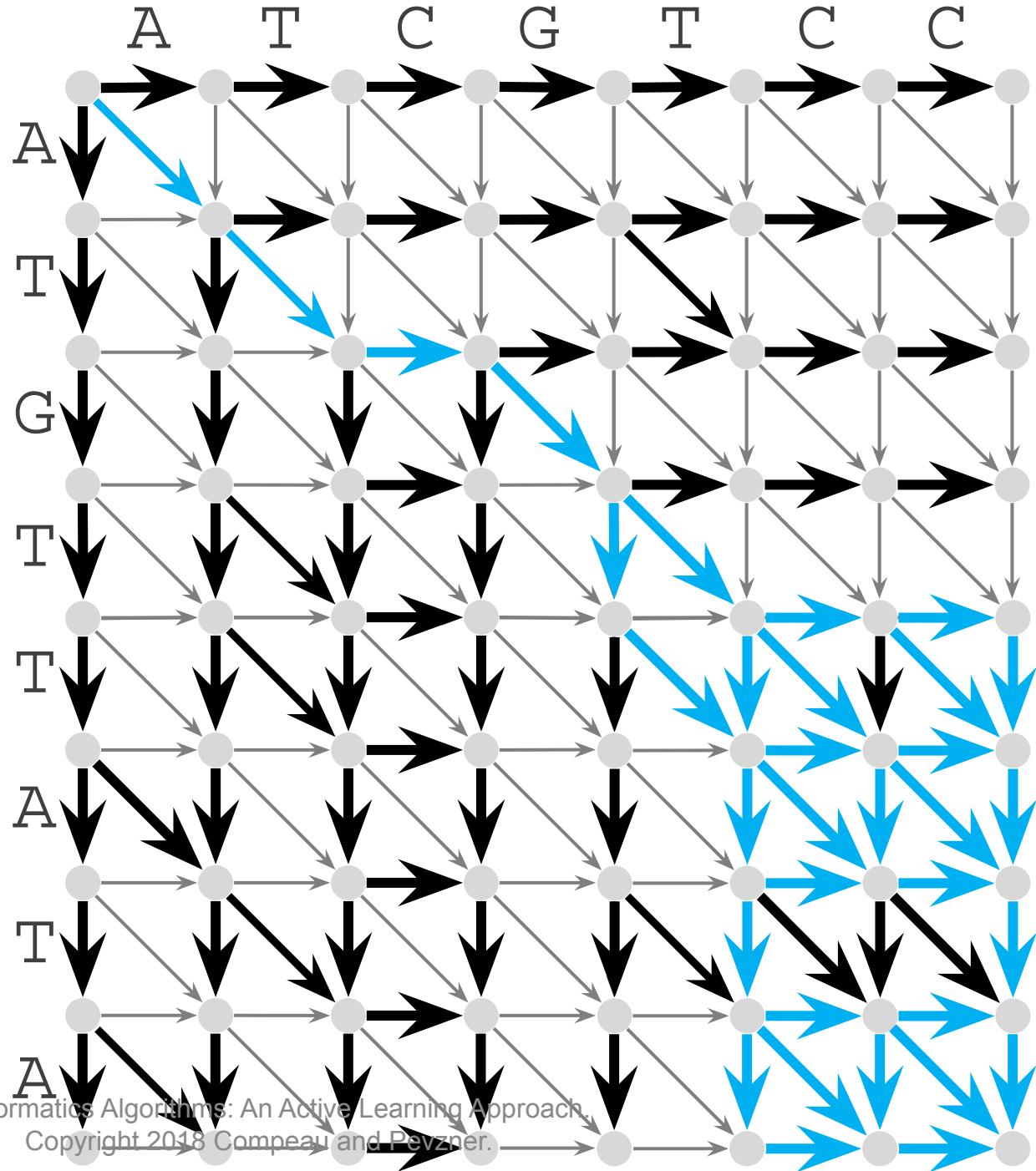
Why did we
store the
backtracking
pointers?



What is the optimal alignment path?



backtracking pointers for the Longest Common Subsequence



Using Backtracking Pointers to Compute LCS

OutputLCS (*backtrack*, v , i , j)

if $i = 0$ **or** $j = 0$

return

if $\text{backtrack}_{i,j} = “\rightarrow”$

OutputLCS (*backtrack*, v , i , $j-1$)

else if $\text{backtrack}_{i,j} = “\downarrow”$

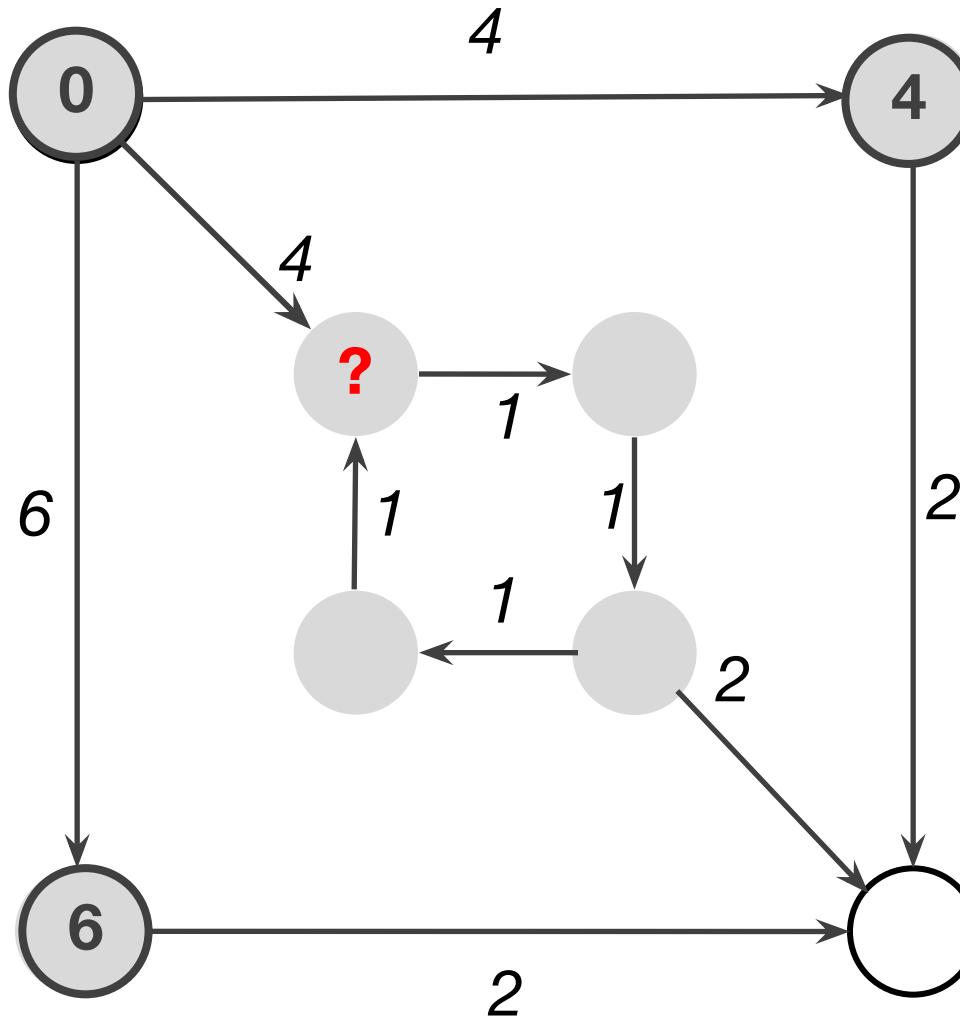
OutputLCS (*backtrack*, v , $i-1$, j)

else

OutputLCS (*backtrack*, v , $i-1$, $j-1$)

output v_i

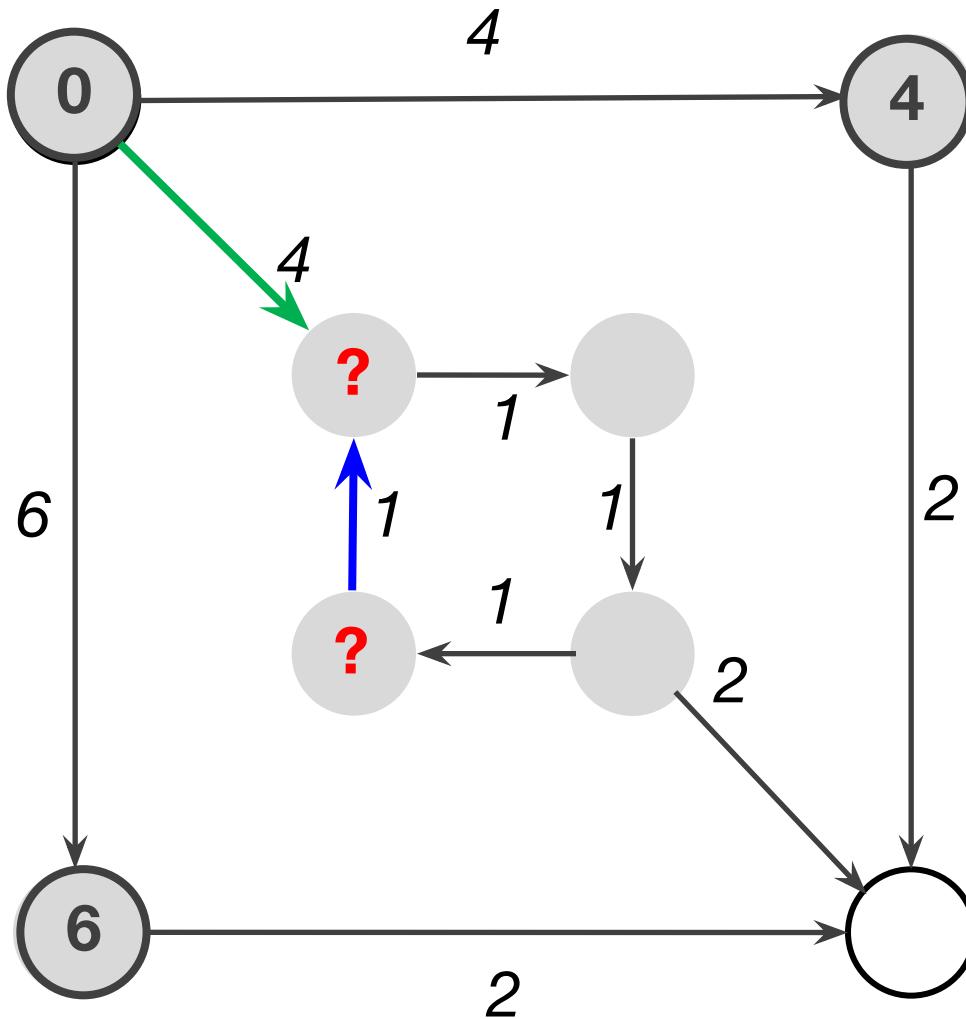
Computing Scores of **ALL** Predecessors

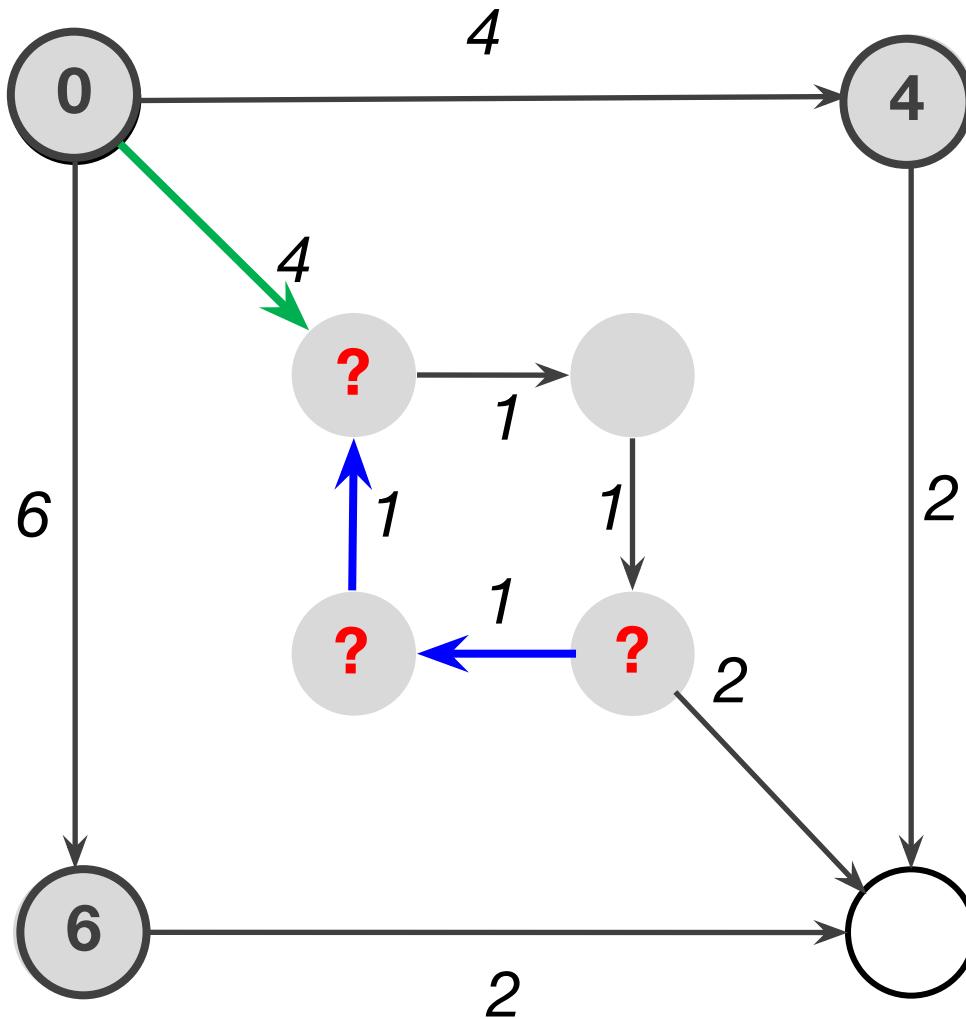


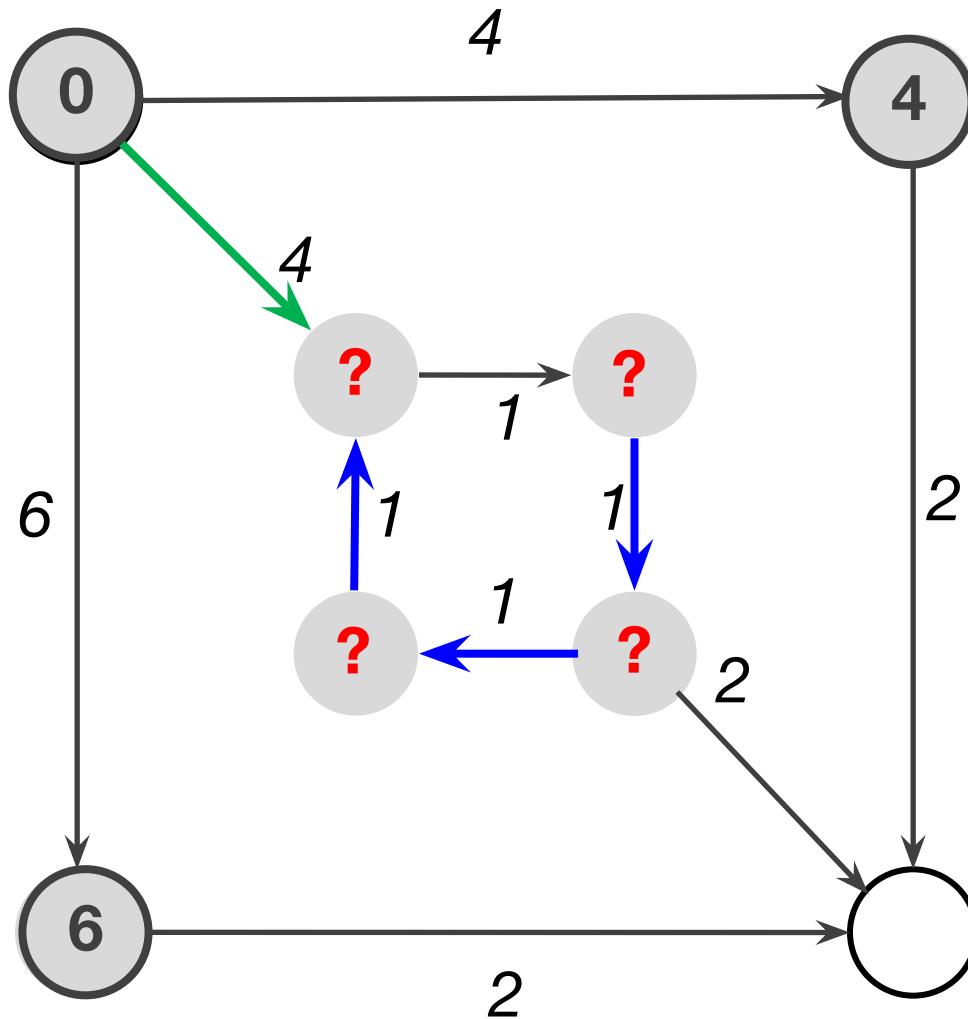
$$s_a = \max_{\text{ALL predecessors } b \text{ of node } a} \{s_b + \text{weight of edge from } b \text{ to } a\}$$

Bioinformatics Algorithms: An Active Learning Approach.

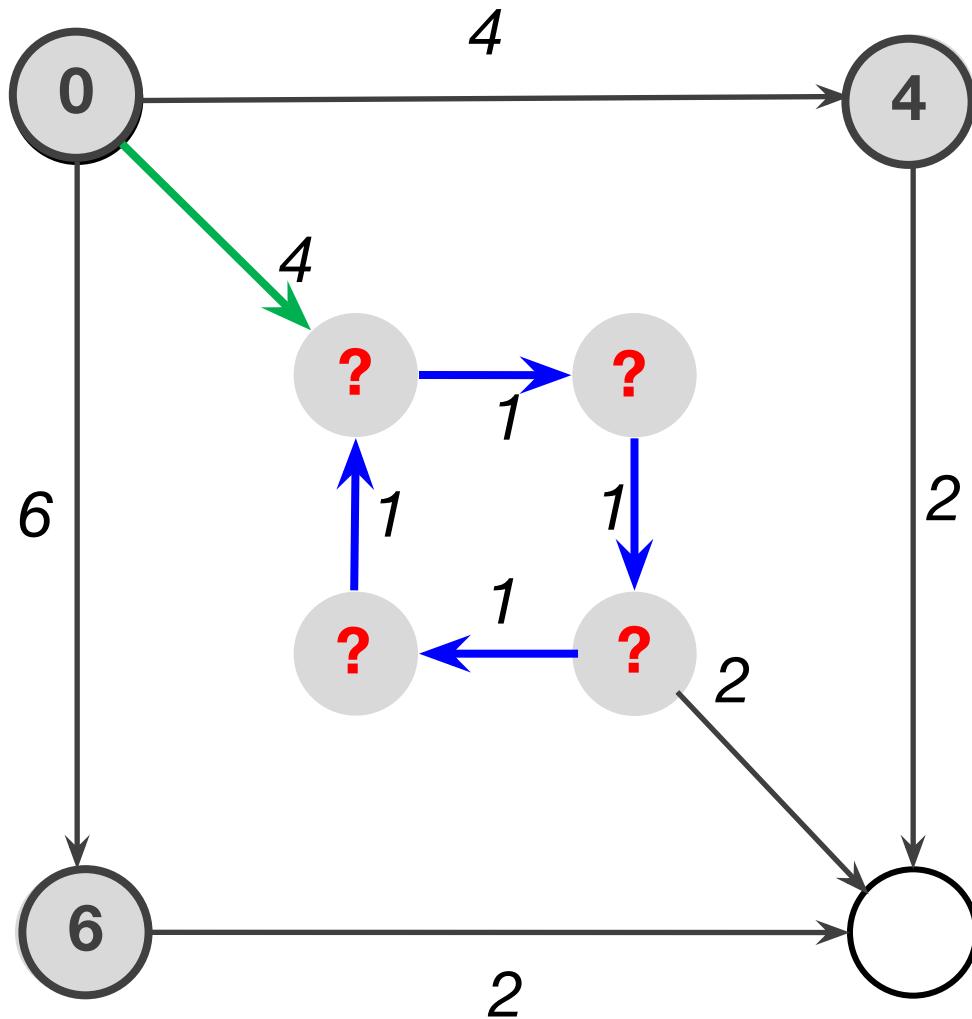
Copyright 2018 Compeau and Pevzner.







A Vicious Cycle



In What Order Should We Explore Nodes of the Graph?

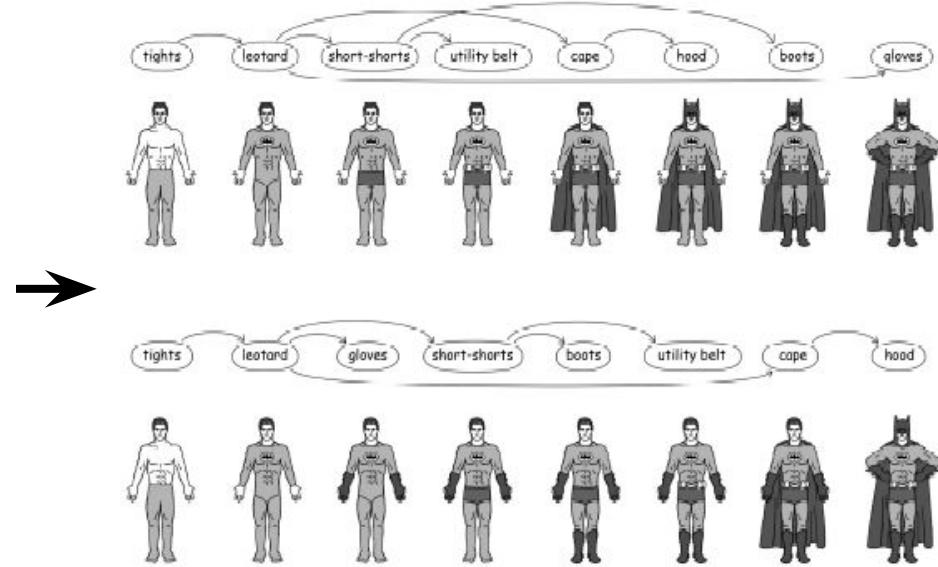
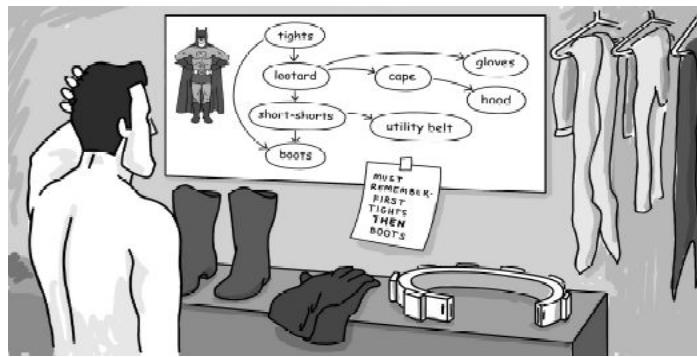
$$s_a = \max_{\text{ALL predecessors } b \text{ of node } a} \{s_b + \text{weight of edge from } b \text{ to } a\}$$



- By the time a node is analyzed, the scores of all its predecessors should already be computed.
- If the graph has a **directed cycle**, this condition is impossible to satisfy.
- **Directed Acyclic Graph (DAG)**: a graph without directed cycles.

Topological Ordering

- **Topological Ordering:** Ordering of nodes of a DAG on a line such that all edges go from left to right.



- **Theorem:** Every DAG has a topological ordering.

LongestPath

LongestPath(*Graph*, *source*, *sink*)

for each node *a* in *Graph*

$s_a \leftarrow -\infty$

$s_{\text{source}} \leftarrow 0$

topologically order *Graph*

for each node *a* (from *source* to *sink* in topological order)

$s_a \leftarrow \max_{\text{all predecessors } b \text{ of node } a} \{s_b + \text{weight of edge from } b \text{ to } a\}$

return s_{sink}

How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- **From Global to Local Alignment**
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

Current (Primitive) Scoring

#matches

Mismatches and Indel Penalties

$$\text{#matches} - \mu \cdot \text{#mismatches} - \sigma \cdot$$

#indels

$$\begin{array}{ccccccccc} A & T & - & G & T & T & A & T & A \\ A & T & C & G & T & - & C & - & C \\ +1 & +1 & -2 & +1 & +1 & -2 & -3 & -2 & -3 = -7 \end{array}$$

	A	C	G	T	-
A	+1	$-\mu$	$-\mu$	$-\mu$	$-\sigma$
C	$-\mu$	+1	$-\mu$	$-\mu$	$-\sigma$
G	$-\mu$	$-\mu$	+1	$-\mu$	$-\sigma$
T	$-\mu$	$-\mu$	$-\mu$	+1	$-\sigma$
-	$-\sigma$	$-\sigma$	$-\sigma$	$-\sigma$	

	A	C	G	T	-
A	+1	-3	-5	-1	-3
C	-4	+1	-3	-2	-3
G	-9	-7	+1	-1	-3
T	-3	-5	-8	+1	-4
-	-4	-2	-2	-1	

Scoring matrix

Even more general scoring matrix

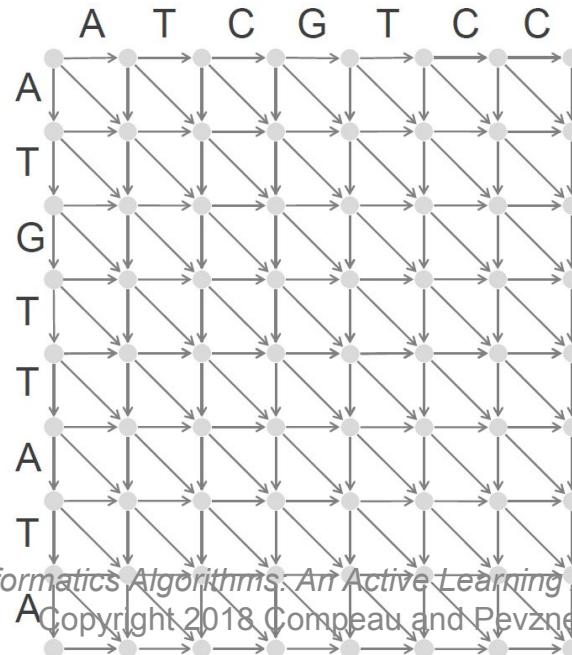
Scoring Matrices for Amino Acid Sequences

C	Cys	12																	
S	Ser	0	2																
T	Thr	-2	1	3															
P	Pro	-3	1	0	6														
A	Ala	-2	1	1	1	2													
G	Gly	-3	1	0	-1	1	5												
N	Asn	-4	1	0	-1	0	0	2											
D	Asp	-5	0	0	-1	0	1	2	4										
E	Glu	-5	0	0	-1	0	0	1	3	4									
Q	Gln	-5	-1	-1	0	0	-1	1	2	2	4								
H	His	-3	-1	-1	0	-1	-2	2	1	1	3	6							
R	Arg	-4	0	-1	0	-2	-3	0	-1	-1	1	2	6						
K	Lys	-5	0	0	-1	-1	-2	1	0	0	1	0	3	5					
M	Met	-5	-2	-1	-2	-1	-3	-2	-3	-2	-1	-2	0	0	6				
I	Ile	-2	-1	0	-2	-1	-3	-2	-2	-2	-2	-2	-2	2	5				
L	Leu	-6	-3	-2	-3	-2	-4	-3	-4	-3	-2	-2	-3	-3	4	2	6		
V	Val	-2	-1	0	-1	0	-1	-2	-2	-2	-2	-2	-2	2	4	2	4		
F	Phe	-4	-3	-3	-5	-5	-5	-4	-6	-5	-5	-2	-4	-5	0	1	2	-1	9
Y	Tyr	0	-3	-3	-5	-3	-5	-2	-4	-4	-4	0	-4	-4	-2	-1	-1	-2	7
W	Trp	-8	-2	-5	-6	-6	-7	-4	-7	-7	-5	-3	2	-3	-4	-5	-2	-6	0
		C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F
																			Y
																			W

Y (Tyr) often mutates into F (score +7)
but rarely mutates into P (score -5)

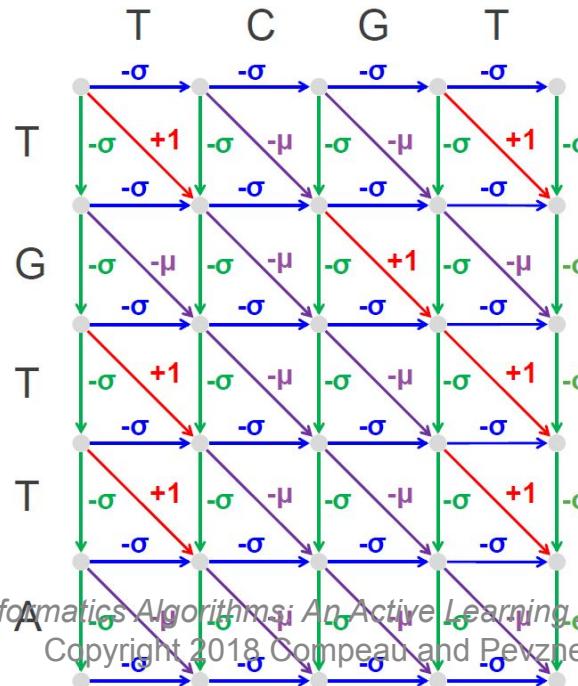
Dynamic Programming Recurrence for the Alignment Graph

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge "↓" into } (i,j) \\ s_{i,j-1} + \text{weight of edge "→" into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge "↘" into } (i,j) \end{cases}$$



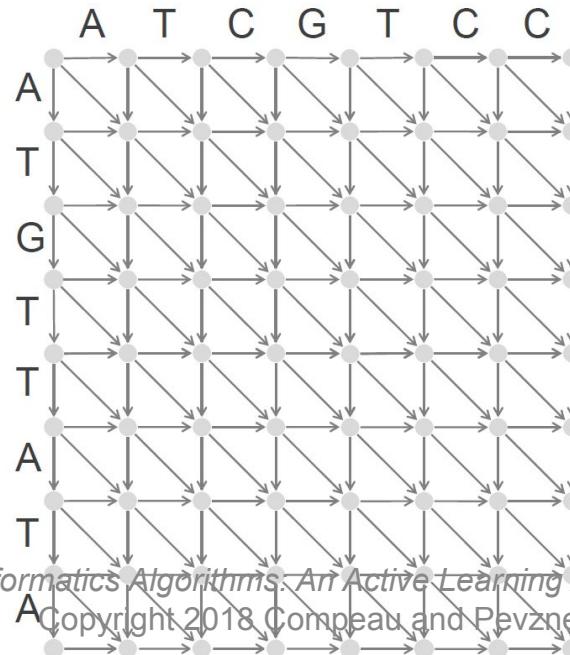
Dynamic Programming Recurrence for the Alignment Graph

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



Dynamic Programming Recurrence for the Alignment Graph

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{score}(v_r, -) \\ s_{i,j-1} + \text{score}(-, w_j) \\ s_{i-1,j-1} + \text{score}(v_r, w_j) \end{cases}$$



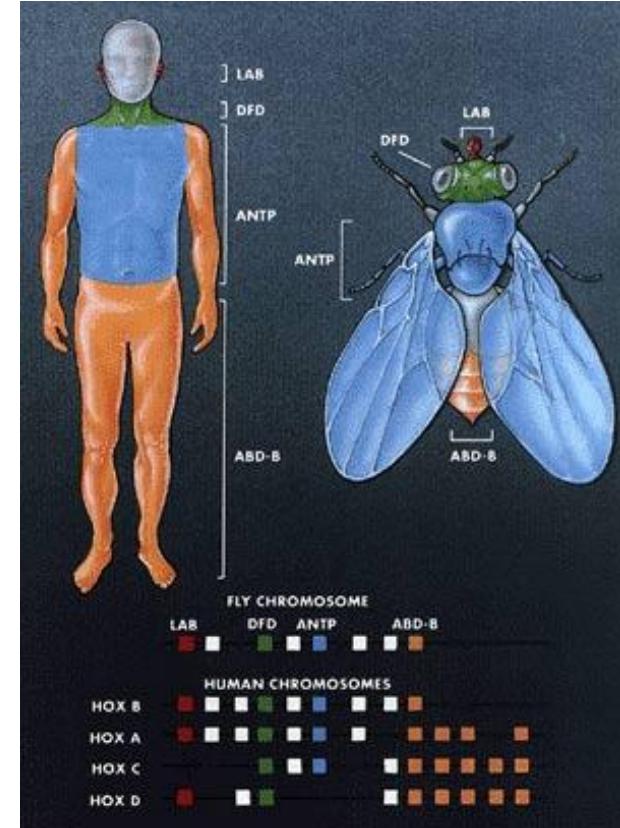
Global Alignment

Global Alignment Problem: Find the highest-scoring alignment between two strings by using a scoring matrix.

- **Input:** Strings v and w as well as a matrix $score$.
- **Output:** An alignment of v and w whose alignment score (as defined by the scoring matrix $score$) is maximal among all possible alignments of v and w .

Homeobox Genes

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



Which Alignment is Better?

- Alignment 1: score = 22 (matches) - 20 (indels)=2.

GCC-C-A**GT**--**TATGT**-CAGGGGG**CACG**--A-G**CATGCAGA**-
GCCGCC-**GTCGT**-T-**TTCAG**----CA-GTT**ATG**--T-CAGAT

- Alignment 2: score = 17 (matches) - 30 (indels)=-13.

---**G**---C-----C--**CAGTTATGTCAGGGGGCACGAGCA****TGCAGA**
GCCGCCGT**CGTTTCAGCAGTTATGTCAG**-----A-----T-----

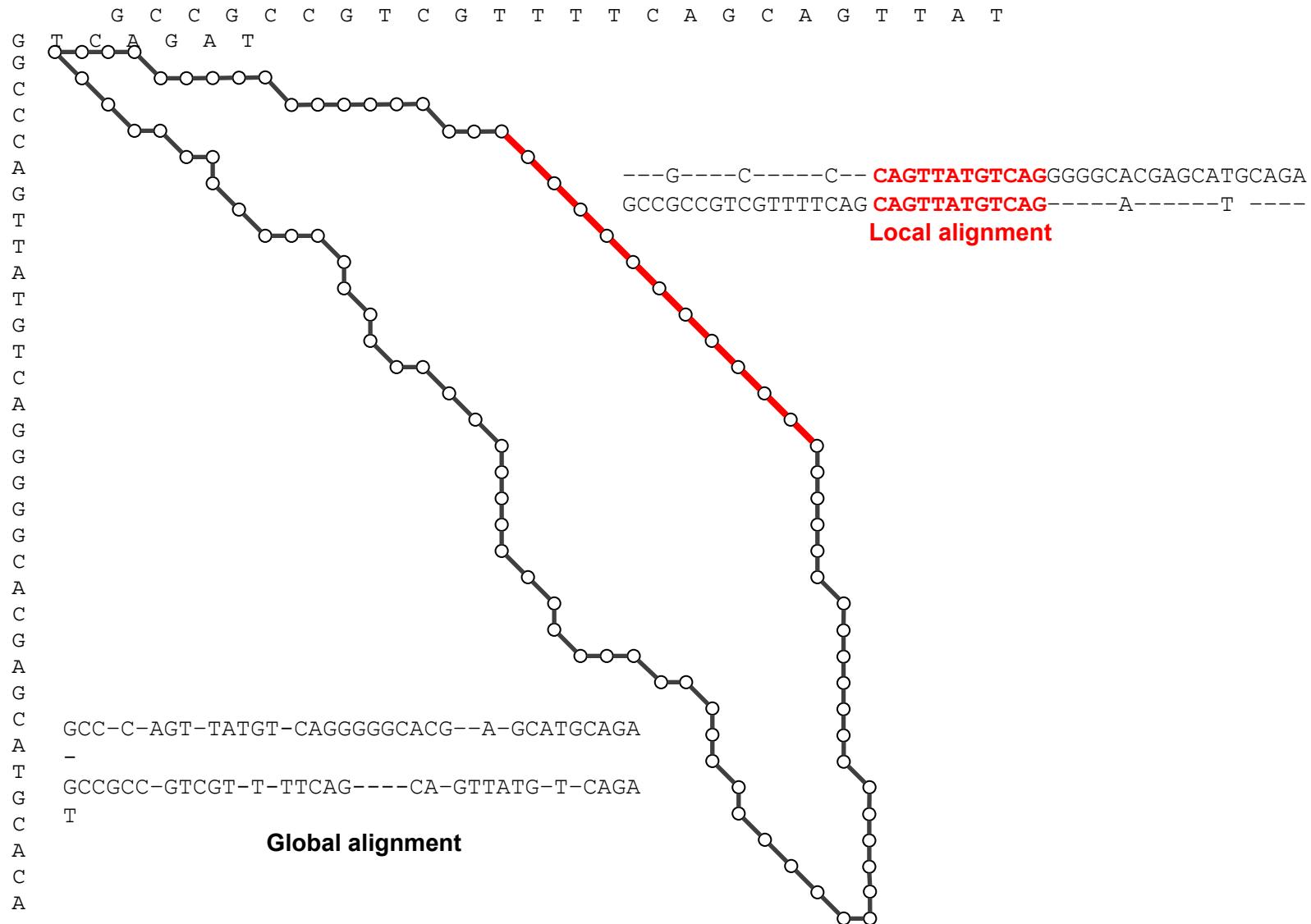
Which Alignment is Better?

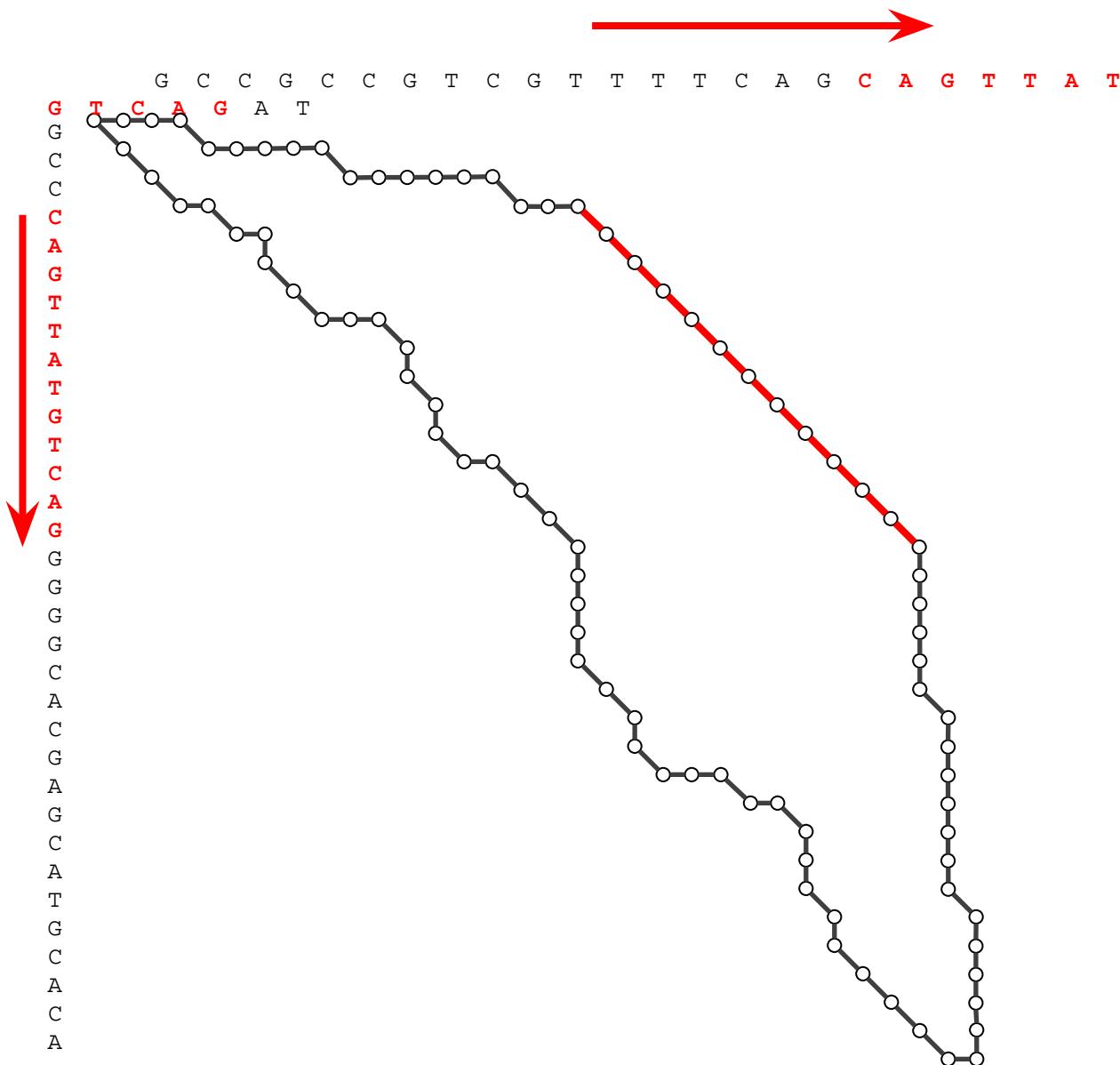
- Alignment 1: score = 22 (matches) - 20 (indels)=2.

GCC-C-A**GT**--**TATGT**-CAGGGG**CACG**--A-G**CATGCAGA**-
GCCGCC-**GTCGT**-T-**TTCAG**----CA-GTT**ATG**--T-CAGAT

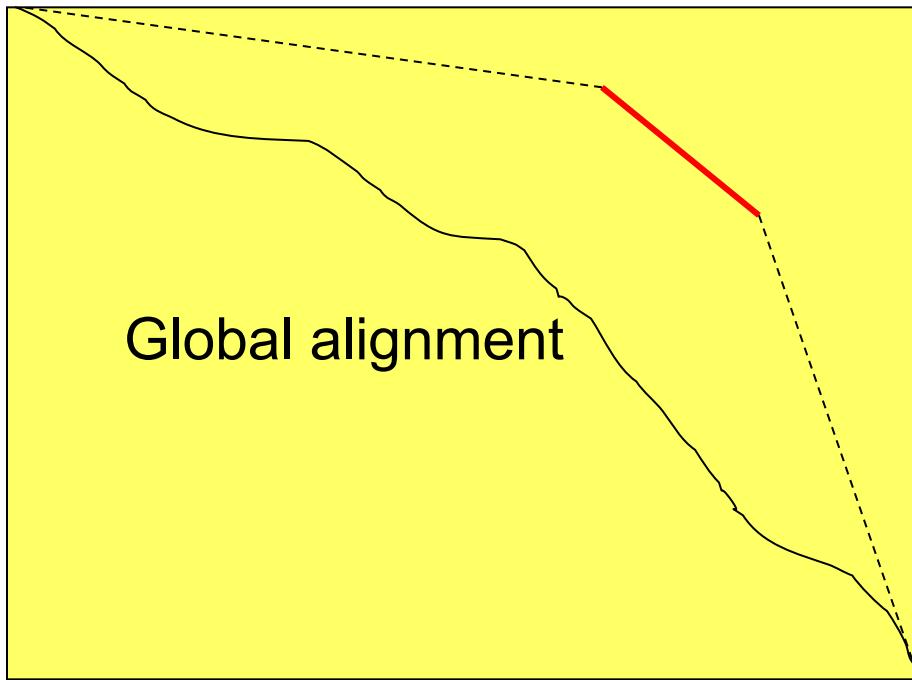
- Alignment 2: score = 17 (matches) - 30 (indels)=-13.

---G-----C-----C--**CAGTTATGTCAG**GGGGCACGAGCATGCAGA
GCCGCCGTCGTTTCAG**CAGTTATGTCAG**-----A-----T-----
local alignment

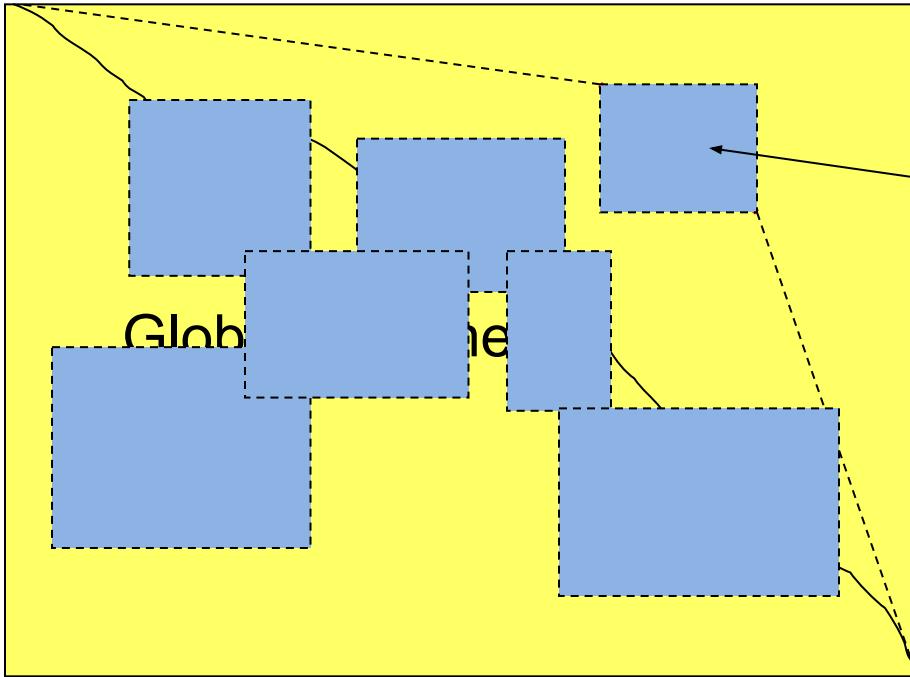




Local Alignment



Local Alignment= Global Alignment in a Subrectangle



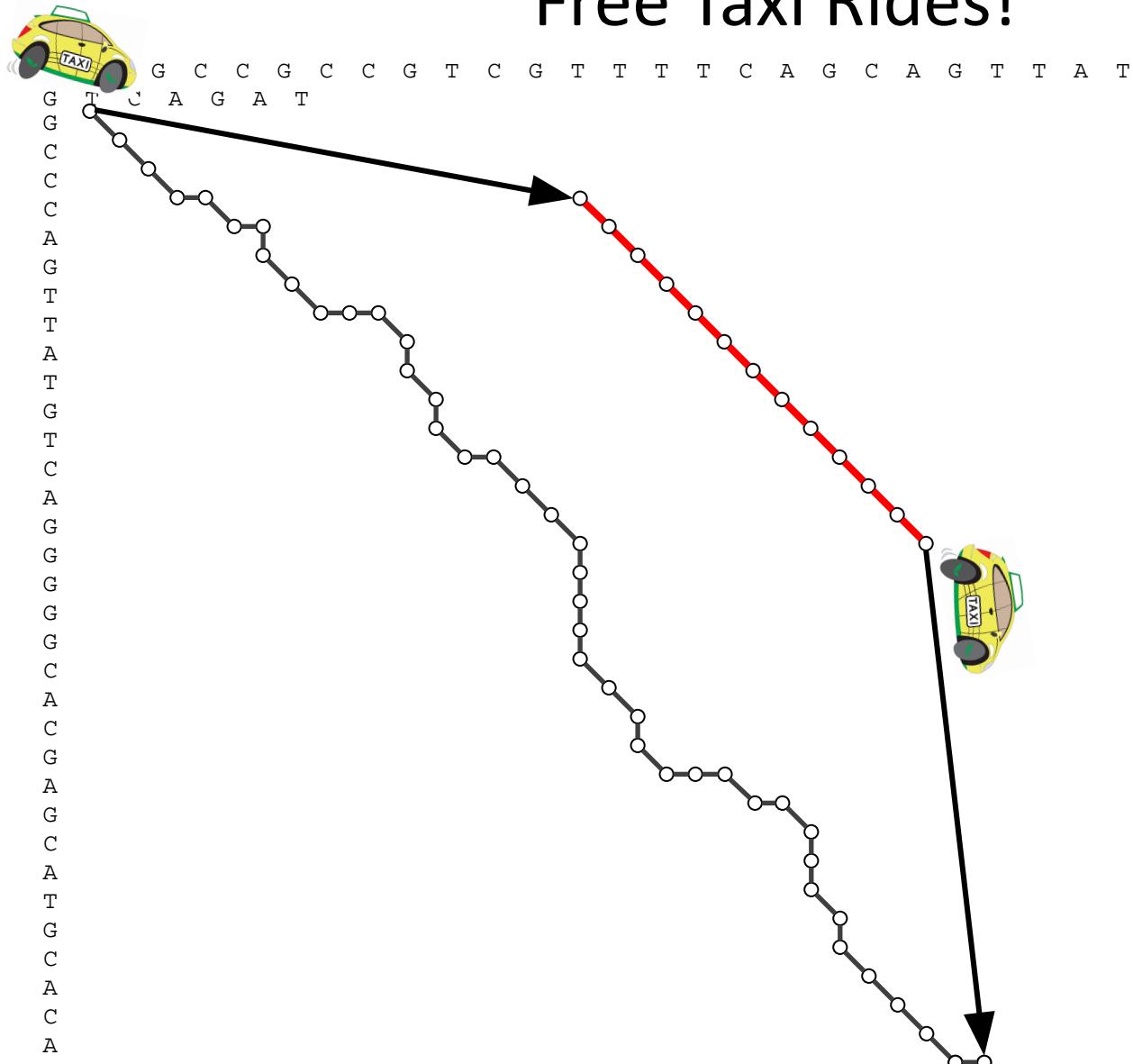
Compute a Global Alignment within each rectangle to get a Local Alignment

Local Alignment Problem

Local Alignment Problem: Find the highest-scoring local alignment between two strings.

- **Input:** Strings v and w as well as a matrix $score$.
- **Output:** Substrings of v and w whose global alignment (as defined by the matrix $score$), is maximal among all global alignments of all substrings of v and w .

Free Taxi Rides!



GCC-C-AGT-TATGT-CAGGGGGCACG--A-GCATGCACA

-
GCCGCC-GTCGT-T-TTCAG---CA-GTTATG-T-CAGA

T

---G-----C-----C-- CAGTTATGTCAG GGGGCACGAGCATGCACA

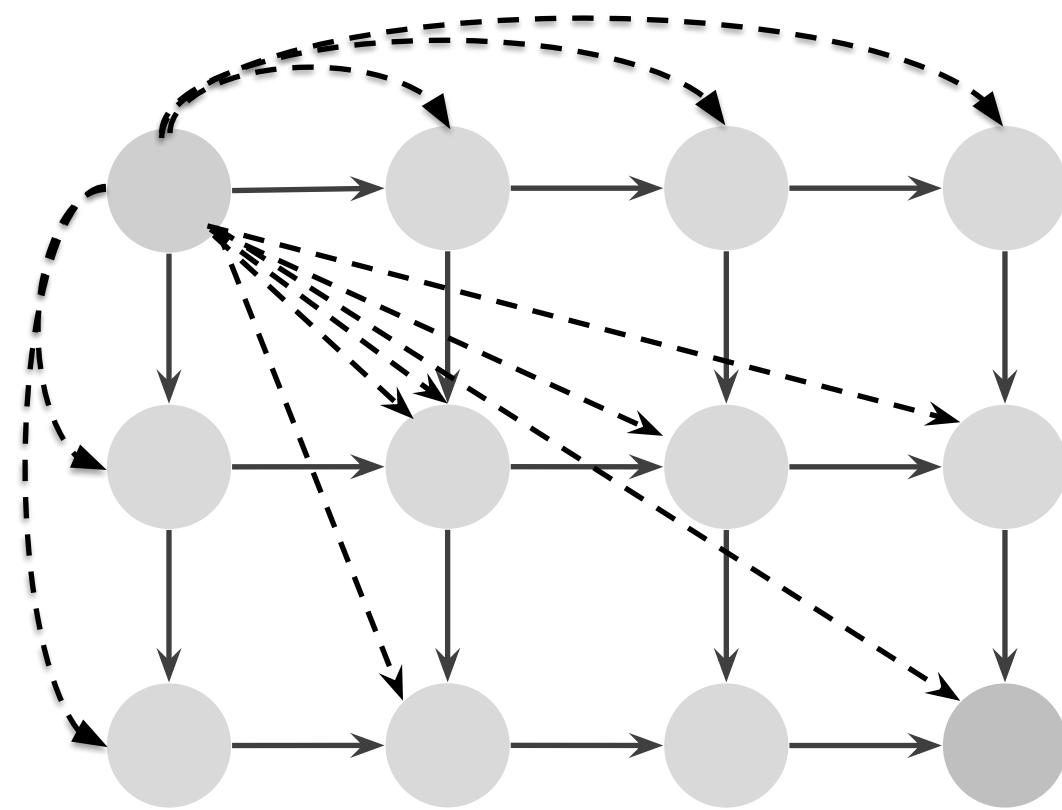
CCCAGGGCTGGTTTTGAC CAGTTATGTCAG -----A-----T -----

Local alignment

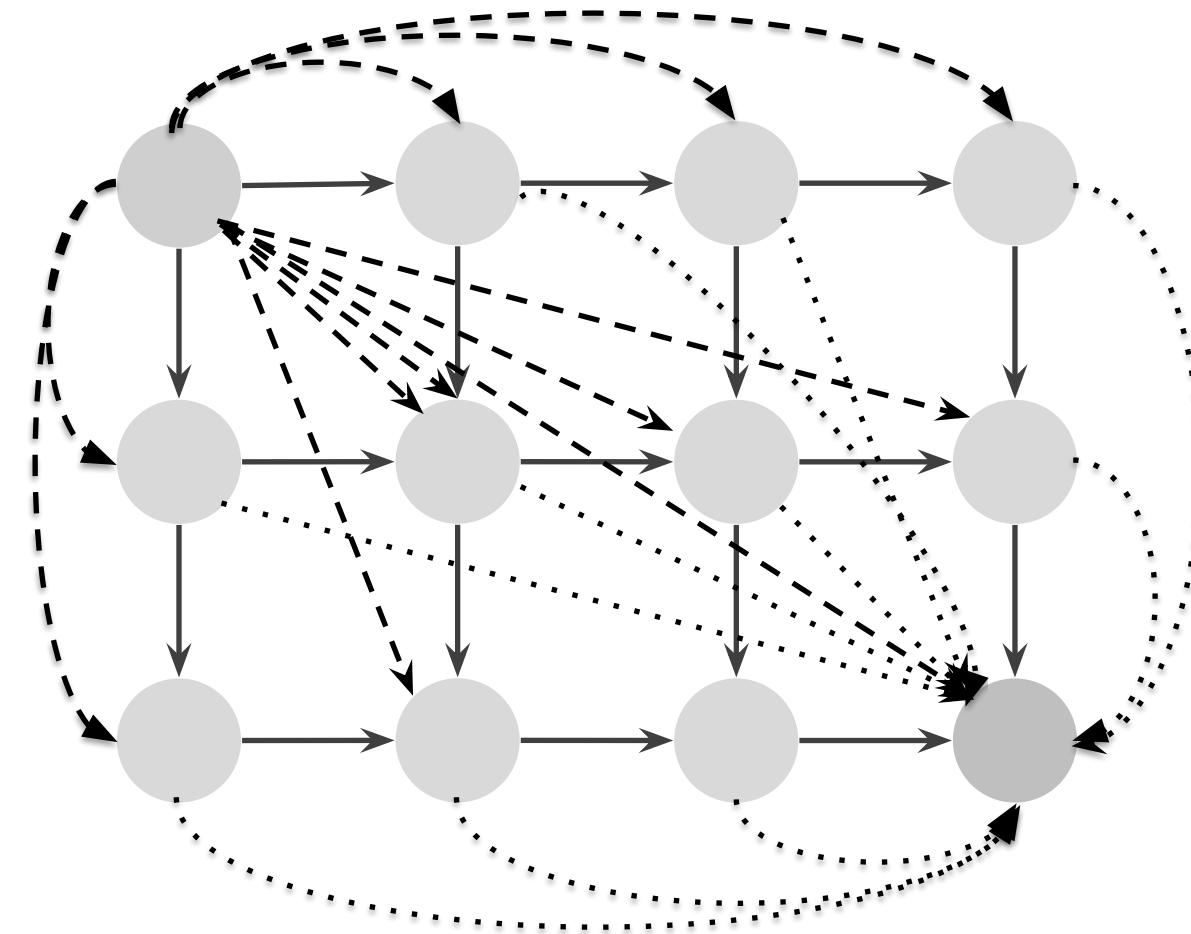
Bioinformatics Algorithms: An Active Learning Approach.

Copyright 2018 Compeau and Pevzner.

What Do Free Taxi Rides Mean in the Terms of the Alignment Graph?



Building Manhattan for the Local Alignment Problem

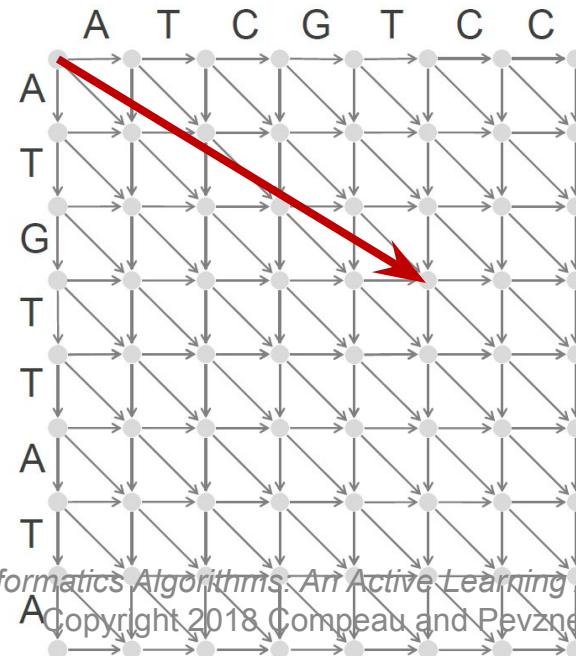


How many edges have we added?

Dynamic Programming for the Local Alignment

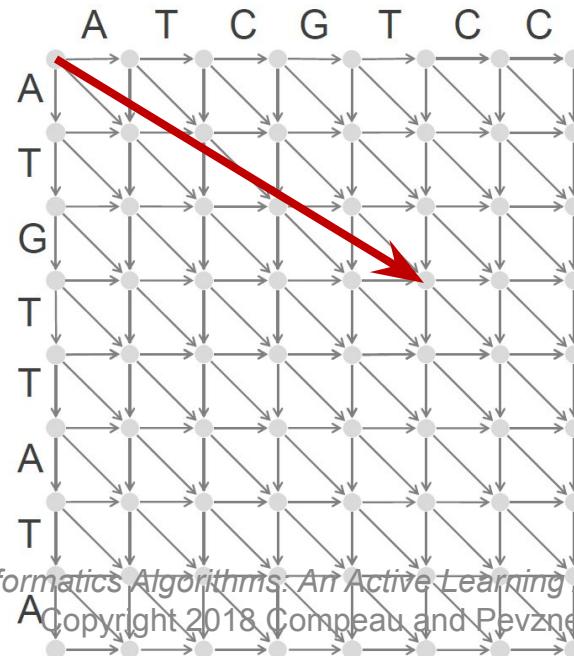
weight of edge from (0,0) to (i,j)

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge "↓" into } (i,j) \\ s_{i,j-1} + \text{weight of edge "→" into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge "↘" into } (i,j) \end{cases}$$



Dynamic Programming for the Local Alignment

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge "↓" into } (i,j) \\ s_{i,j-1} + \text{weight of edge "→" into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge "↖" into } (i,j) \end{cases}$$



How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to the Alignment Graph
- From Global to Local Alignment
- **Penalizing Insertions and Deletions in Sequence Alignment**
- Space-Efficient Sequence Alignment
- Multiple Sequence Alignment

Scoring Gaps

- We previously assigned a fixed penalty σ to each indel.
- However, this fixed penalty may be too severe for a series of 100 consecutive indels.
- A series of k indels often represents a single evolutionary event (**gap**) rather than k events:

two gaps (lower score)	GATCCAG GA-C-AG	GATCCAG GA--CAG	a single gap (higher score)
---------------------------	----------------------------------	----------------------------------	--------------------------------

More Adequate Gap Penalties

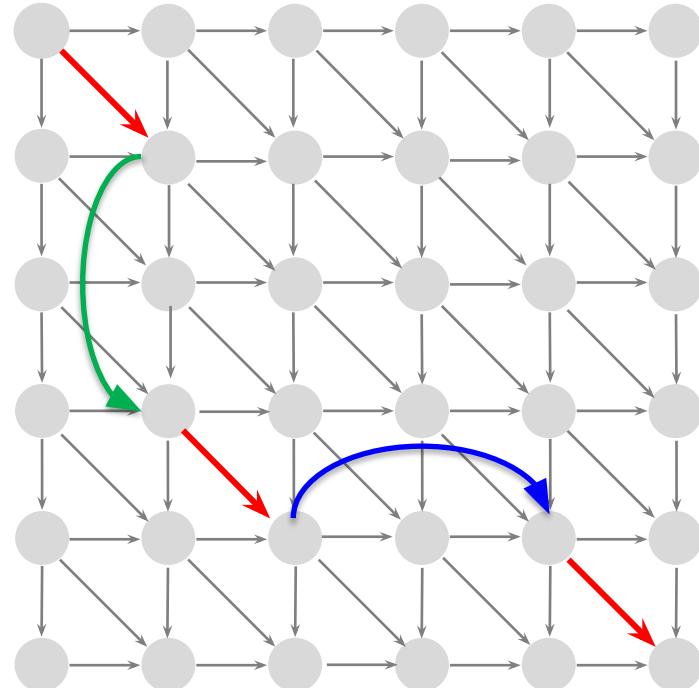
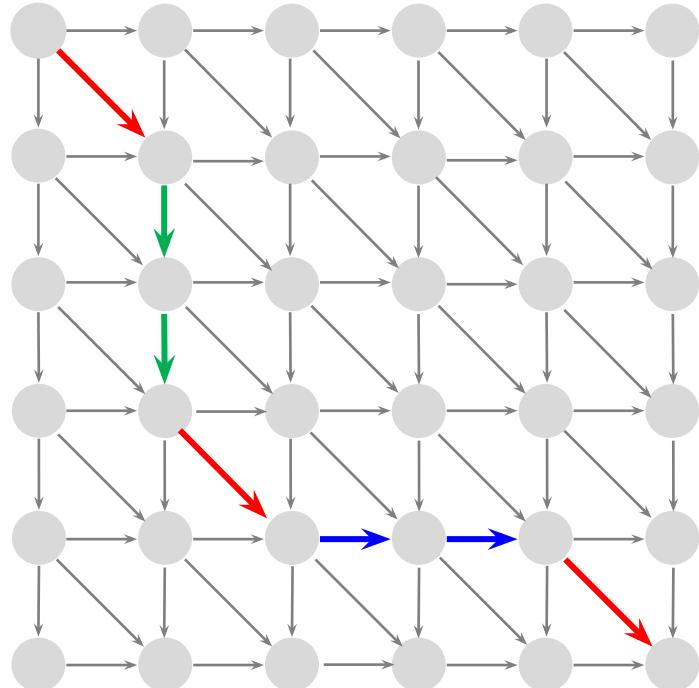
Affine gap penalty for a gap of length k : $\sigma + \varepsilon \cdot (k-1)$

σ - the **gap opening penalty**

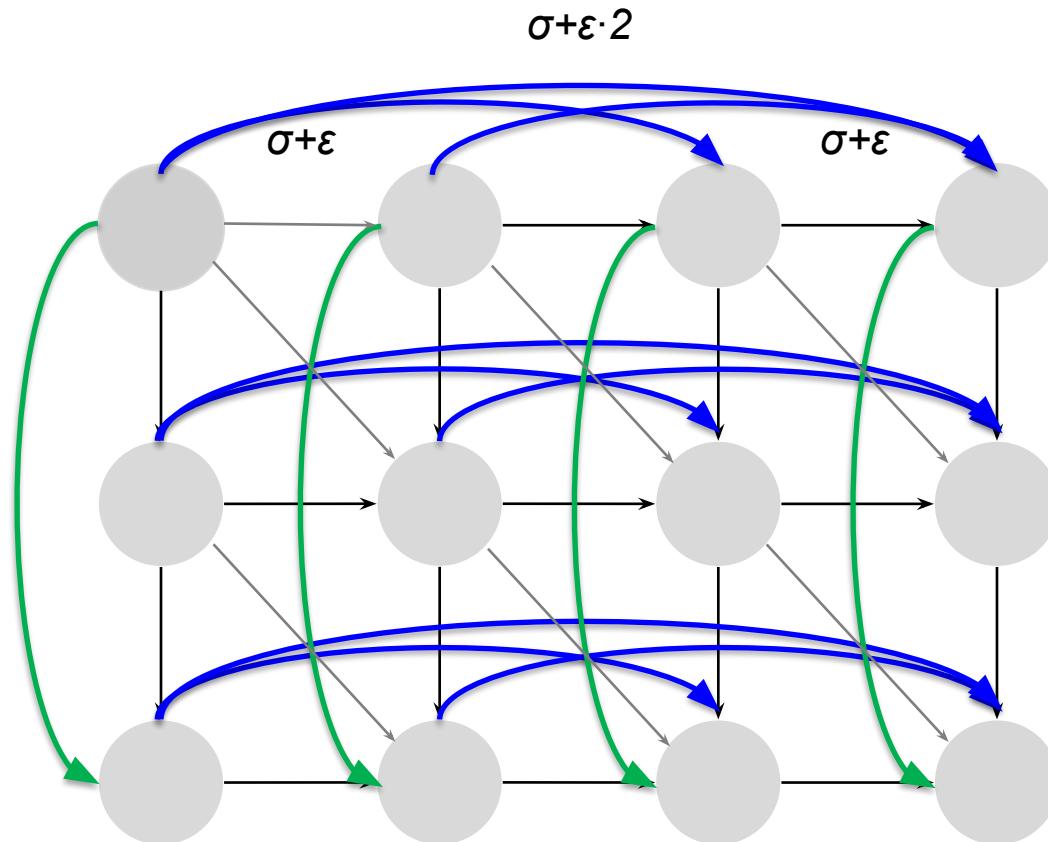
ε - the **gap extension penalty**

$\sigma > \varepsilon$, since starting a gap should be penalized more than extending it.

Modelling Affine Gap Penalties by Long Edges



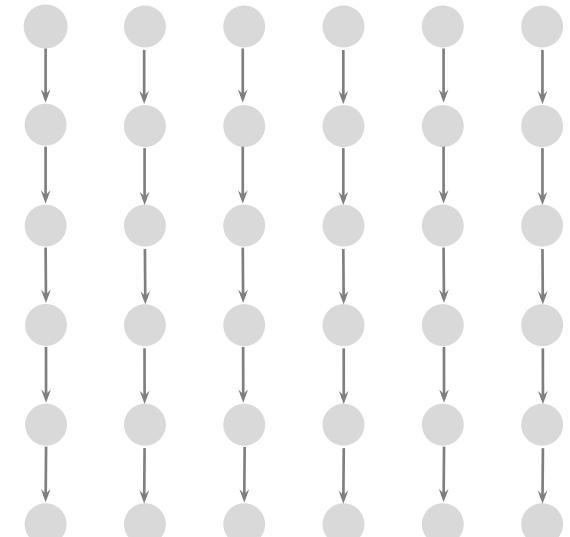
Building Manhattan with Affine Gap Penalties



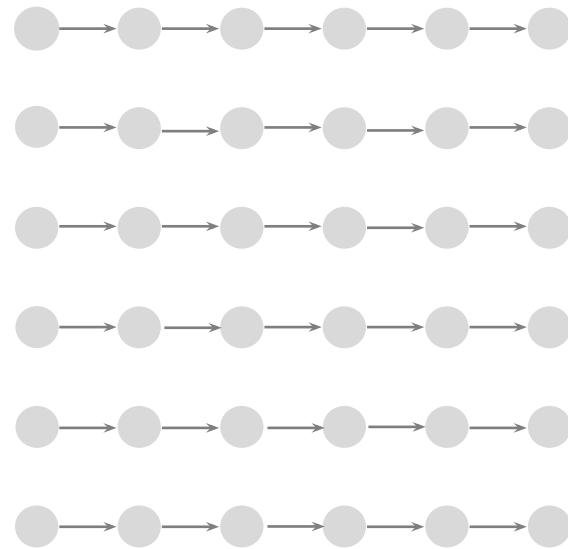
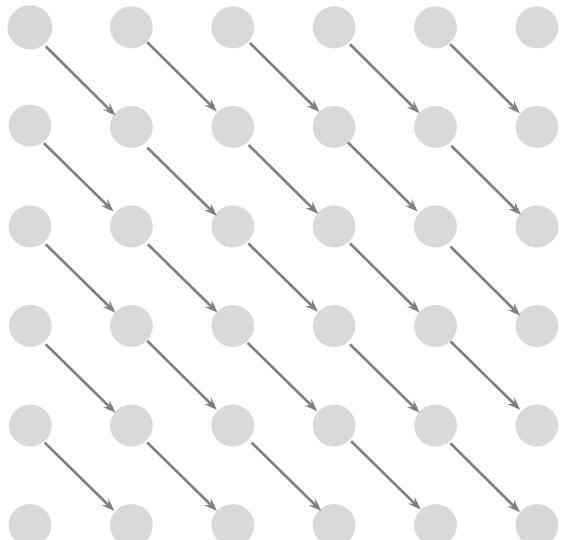
We have just added $O(n^3)$ edges to the graph...

Building Manhattan on 3 levels

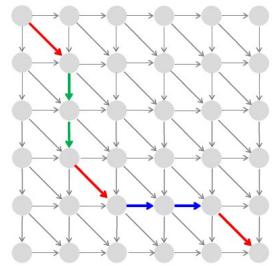
bottom level
(insertions)



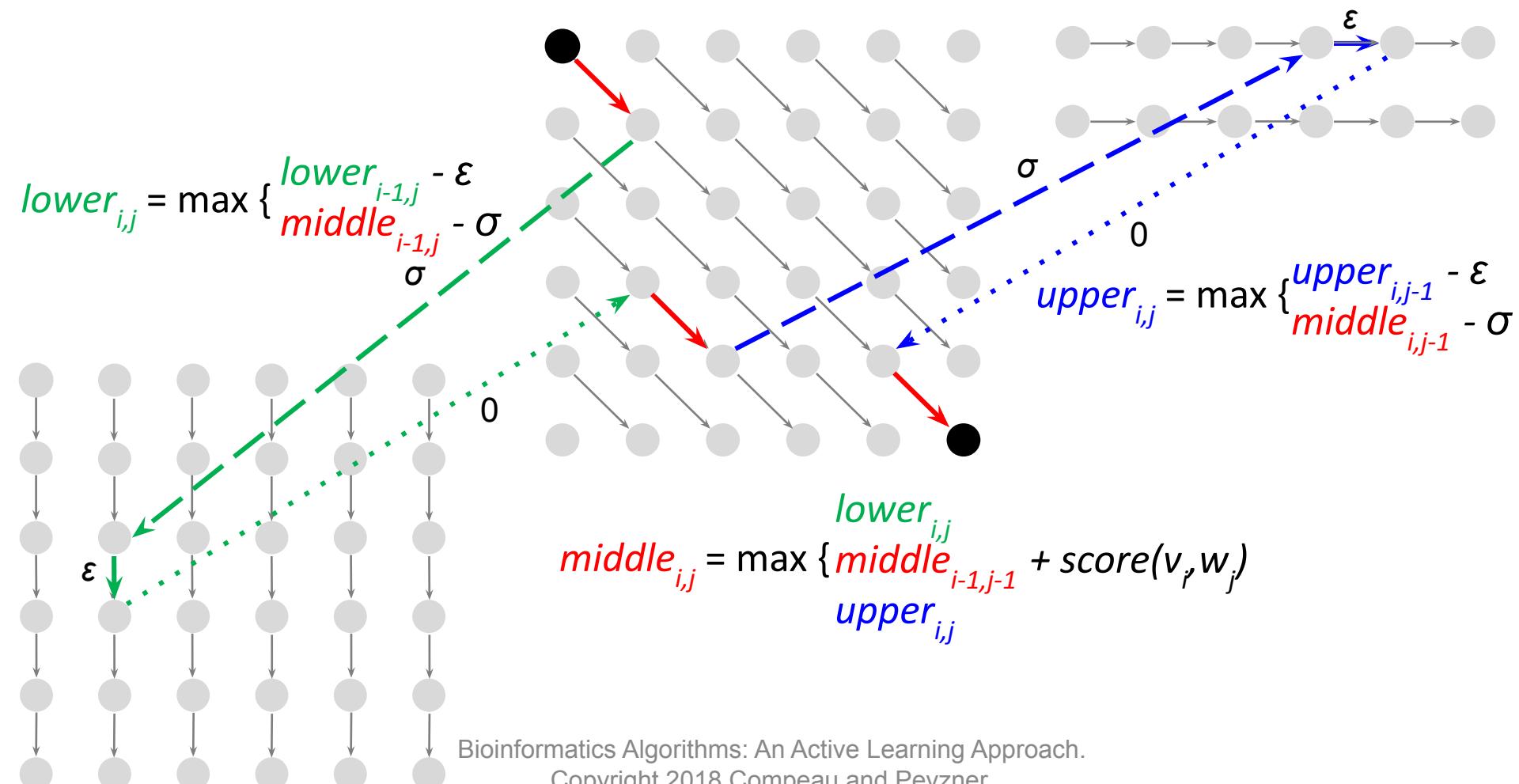
middle level
(matches/mismatches)



upper level
(deletions)



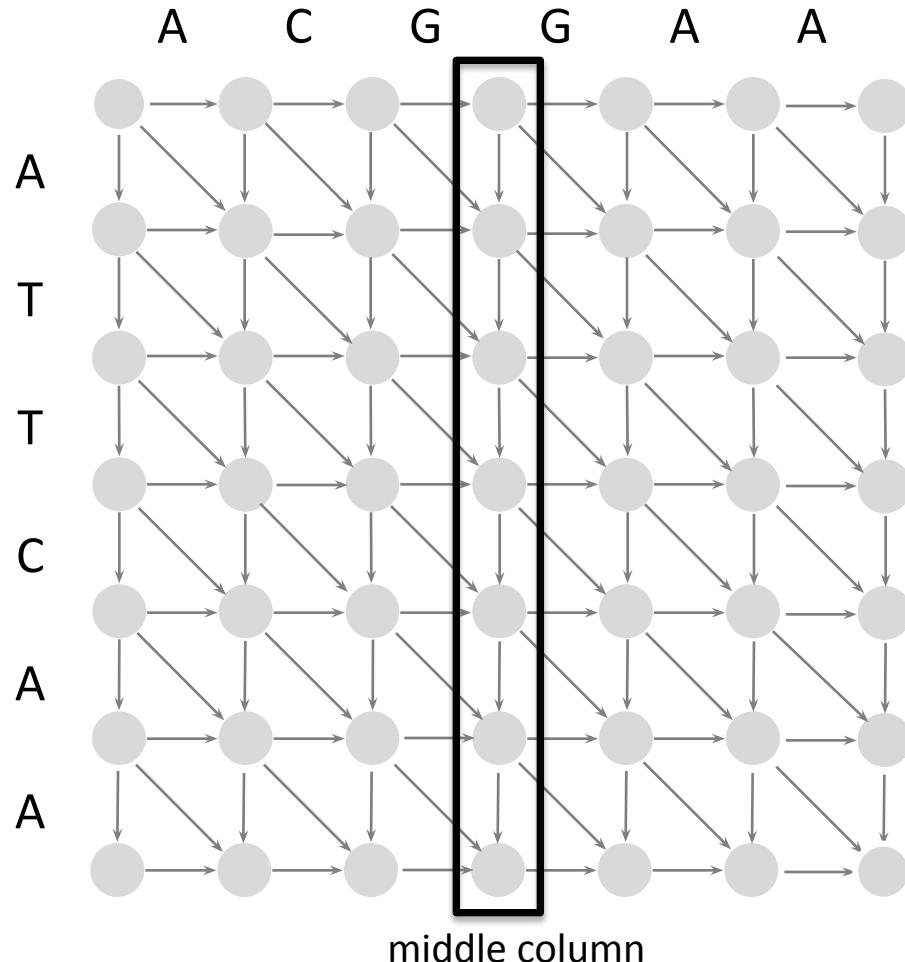
How can we emulate
this path in the 3-level
Manhattan?



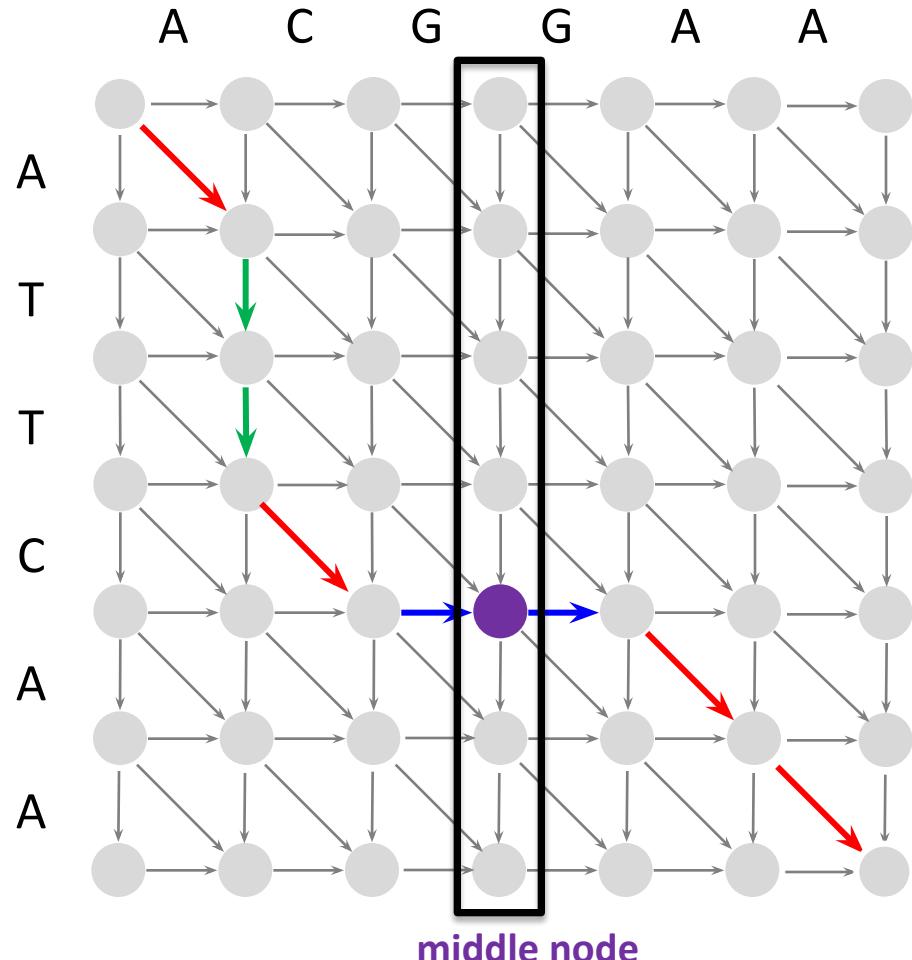
How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to an Arbitrary Directed Acyclic Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- **Space-Efficient Sequence Alignment**
- Multiple Sequence Alignment

Middle Column of the Alignment



Middle Node of the Alignment



(a node where an optimal alignment path crosses the middle column)

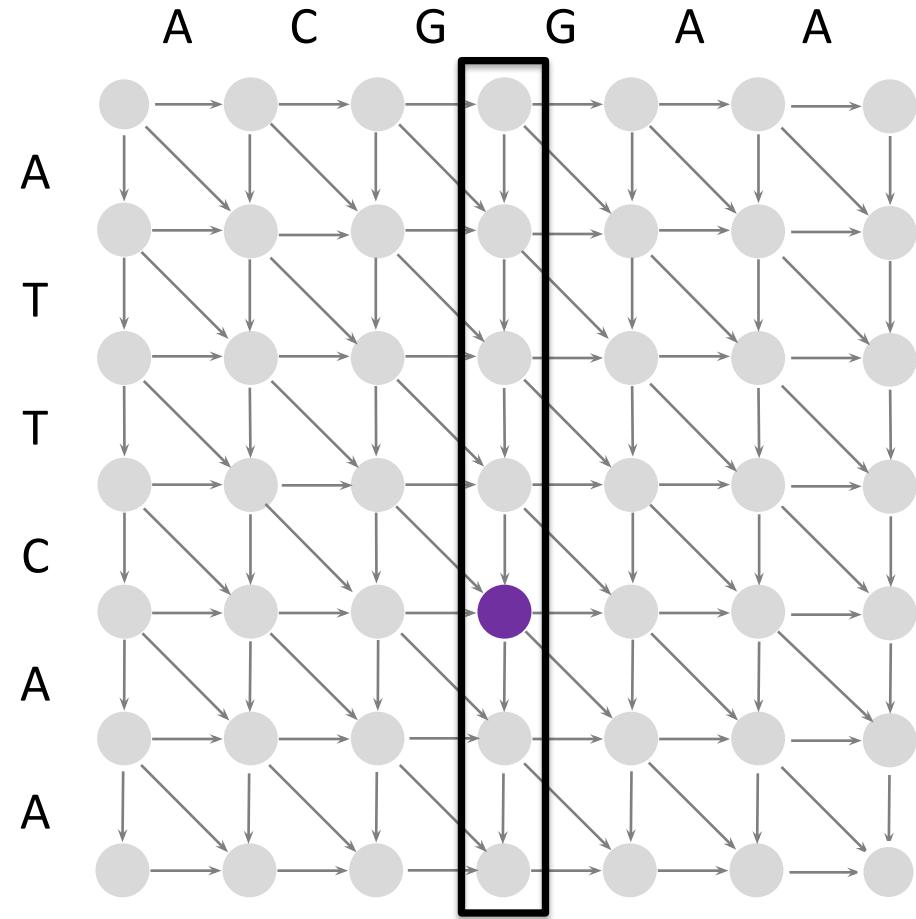
Bioinformatics Algorithms: An Active Learning Approach.

Copyright 2018 Compeau and Pevzner.

Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source, sink*)

find *MiddleNode*

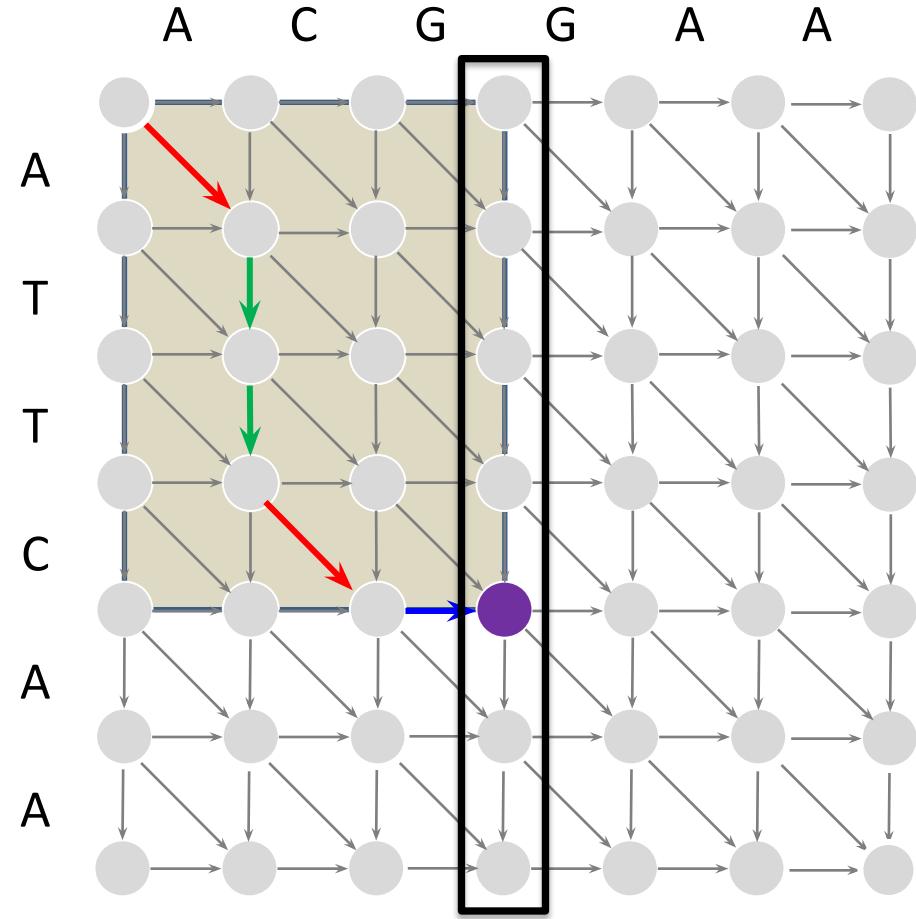


Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source, sink*)

find *MiddleNode*

AlignmentPath(*source, MiddleNode*)



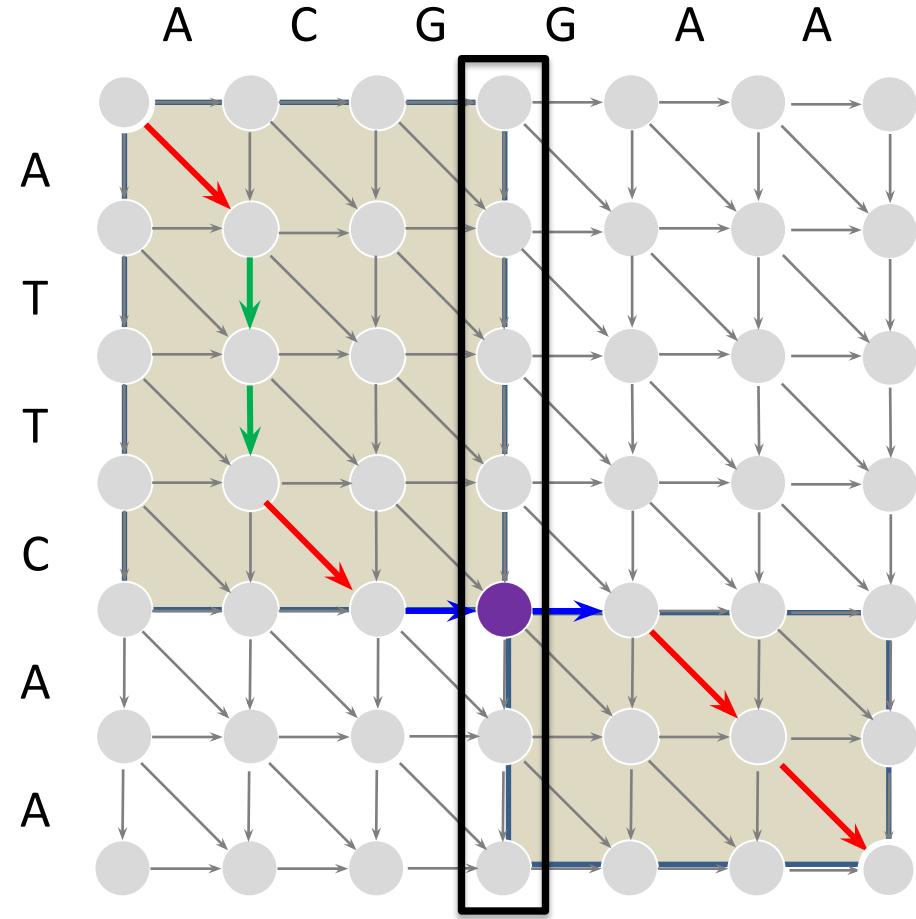
Divide and Conquer Approach to Sequence Alignment

AlignmentPath(*source, sink*)

 find *MiddleNode*

AlignmentPath(*source, MiddleNode*)

AlignmentPath(*MiddleNode, sink*)



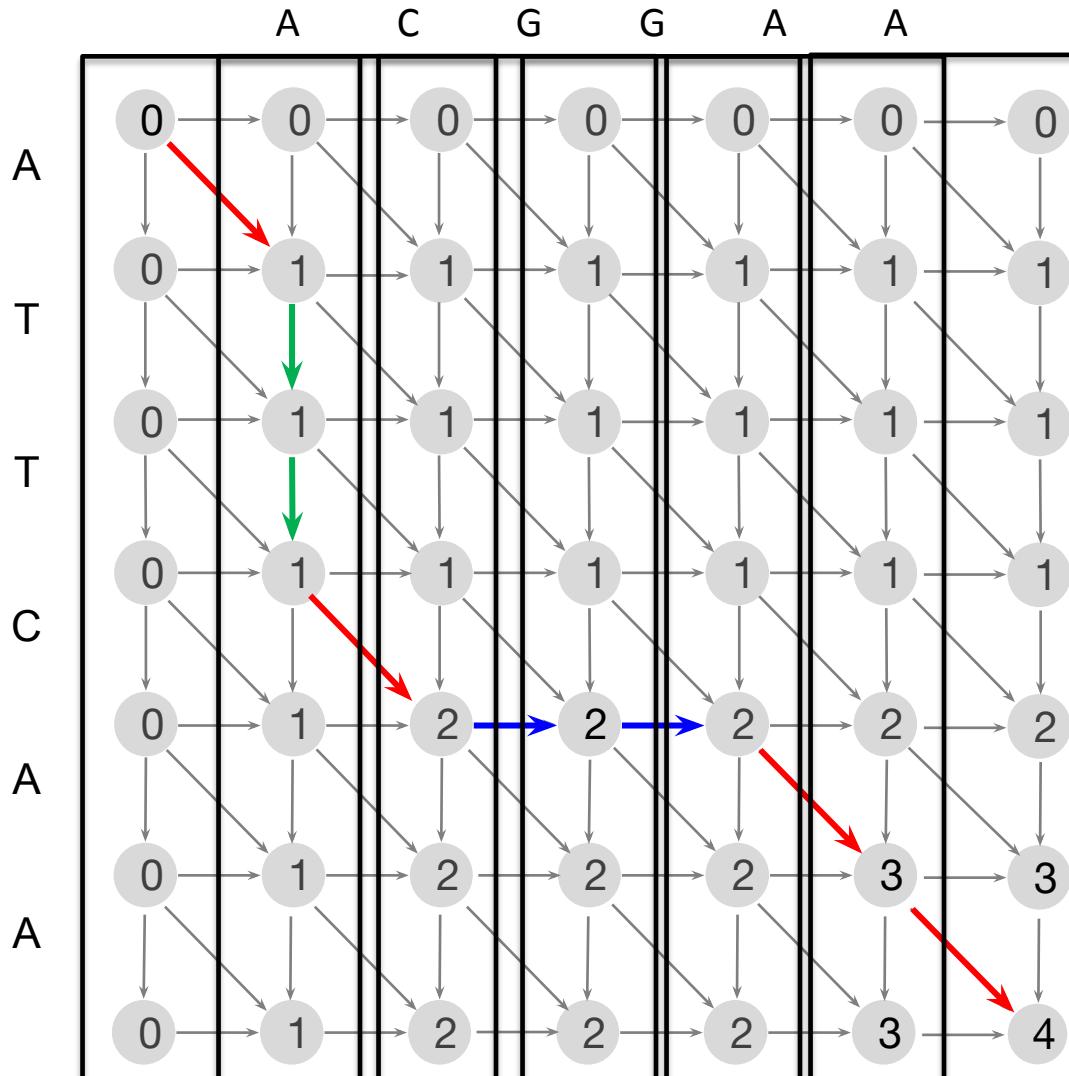
The only problem left is how to find this middle node in **linear space**!

Computing Alignment Score in Linear Space

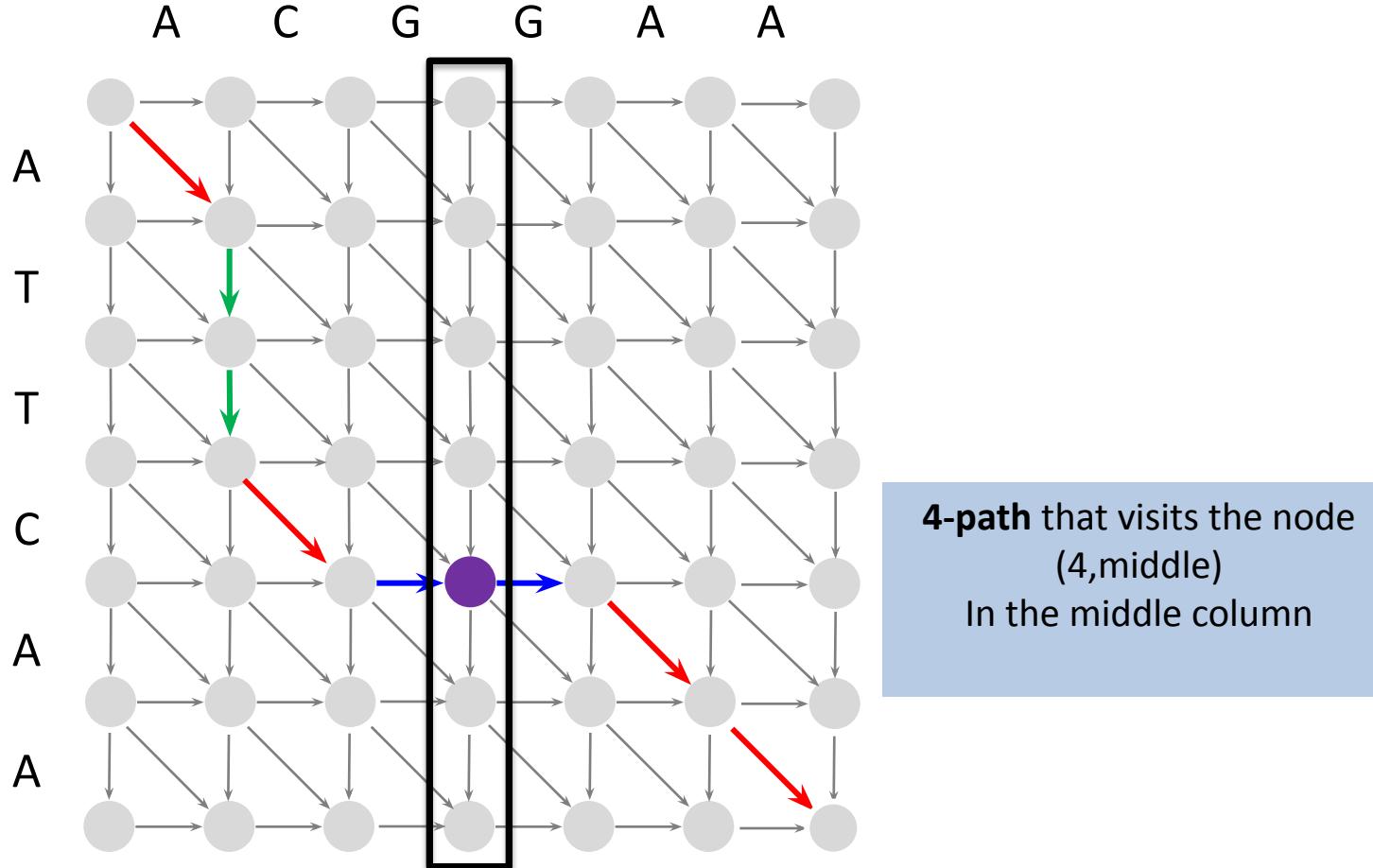
Finding the **longest path** in the alignment graph **requires** storing all backtracking pointers – $O(nm)$ memory.

Finding the **length of the longest path** in the alignment graph **does not require** storing any backtracking pointers – $O(n)$ memory.

Recycling the Columns in the Alignment Graph

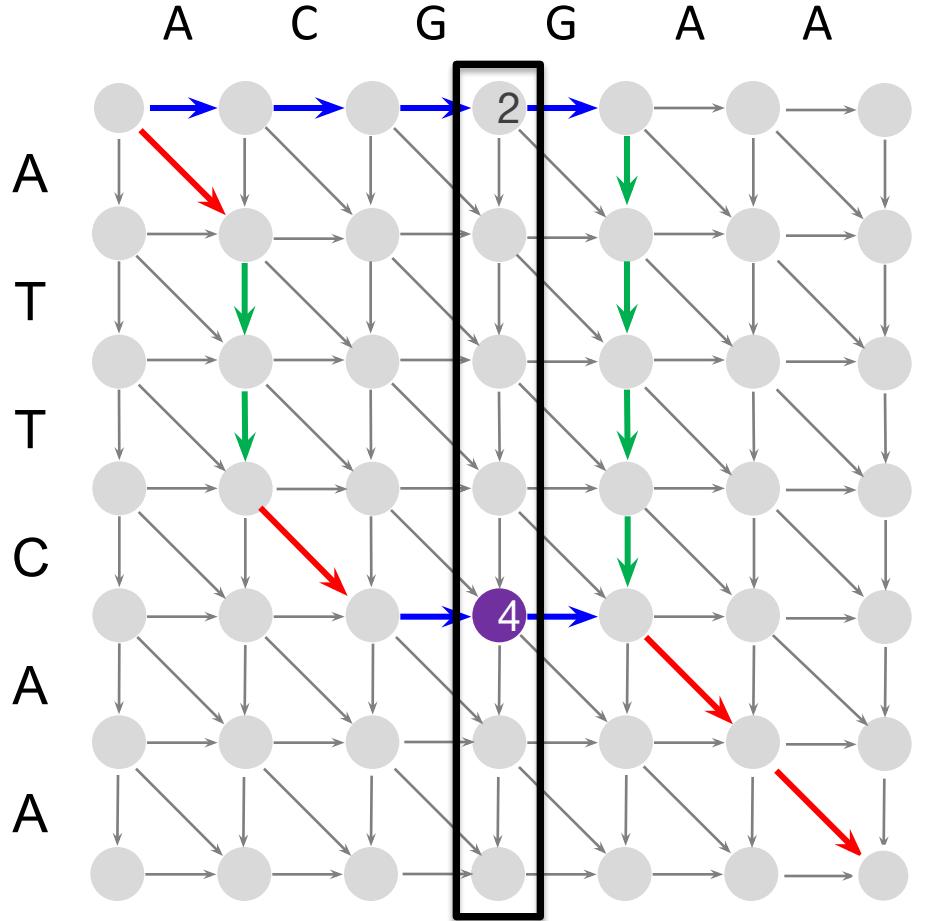


Can We Find the Middle Node without Constructing the Longest Path?



i-path – a longest path among paths that visit the *i*-th node in the middle column

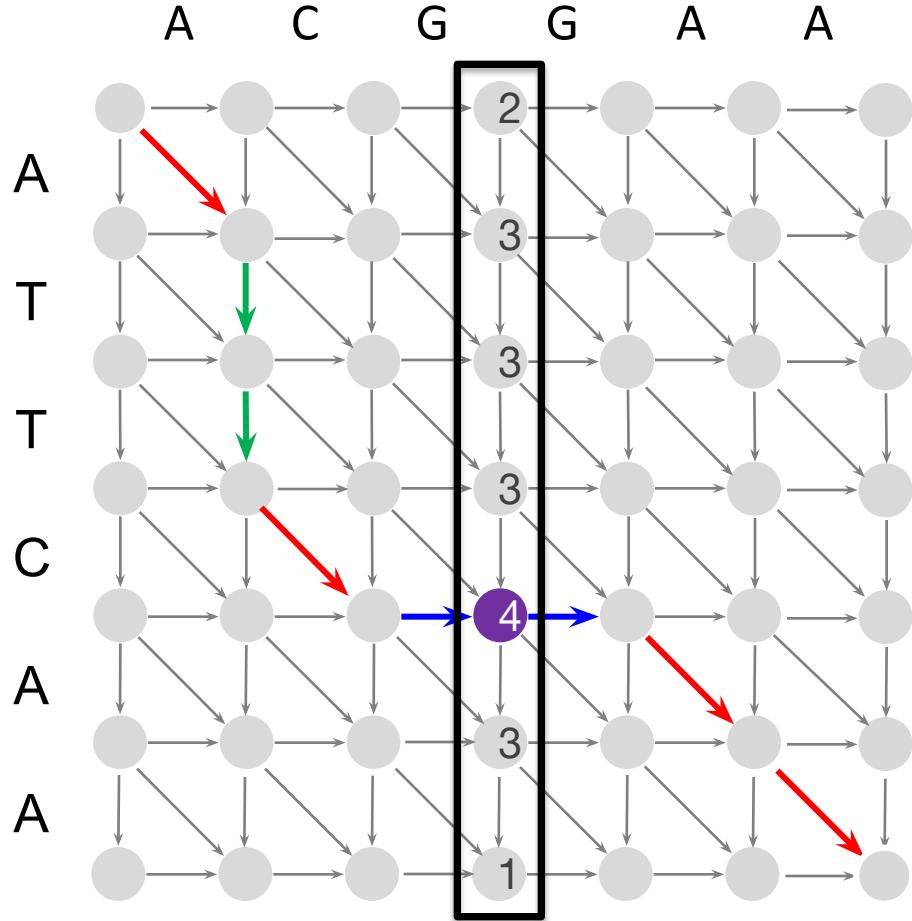
Can We Find The Lengths of All i -paths?



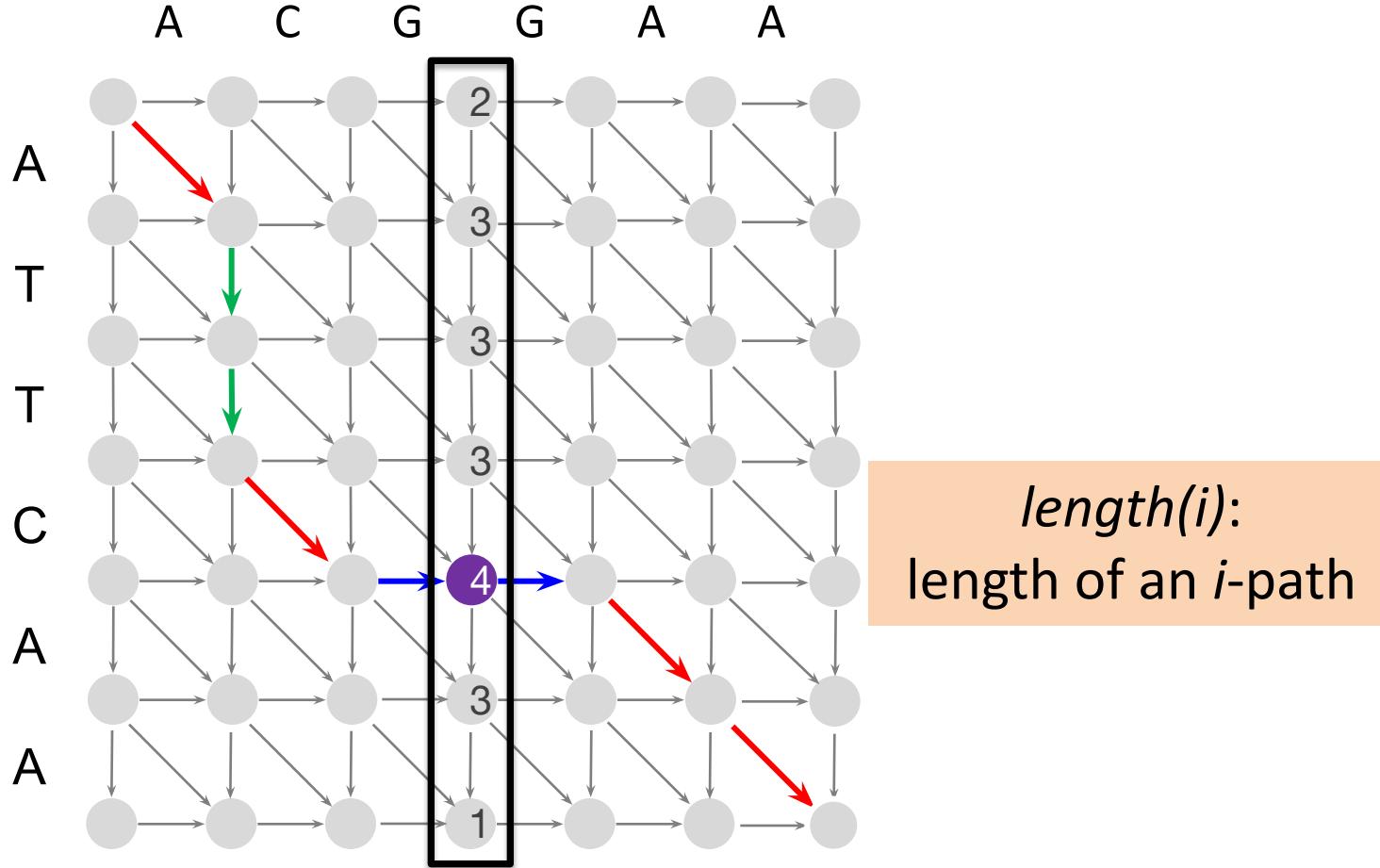
$\text{length}(i)$:
length of an i -path:

$$\begin{aligned}\text{length}(0) &= 2 \\ \text{length}(4) &= 4\end{aligned}$$

Can We Find The Lengths of All i -paths?



Can We Find The Lengths of i -paths?

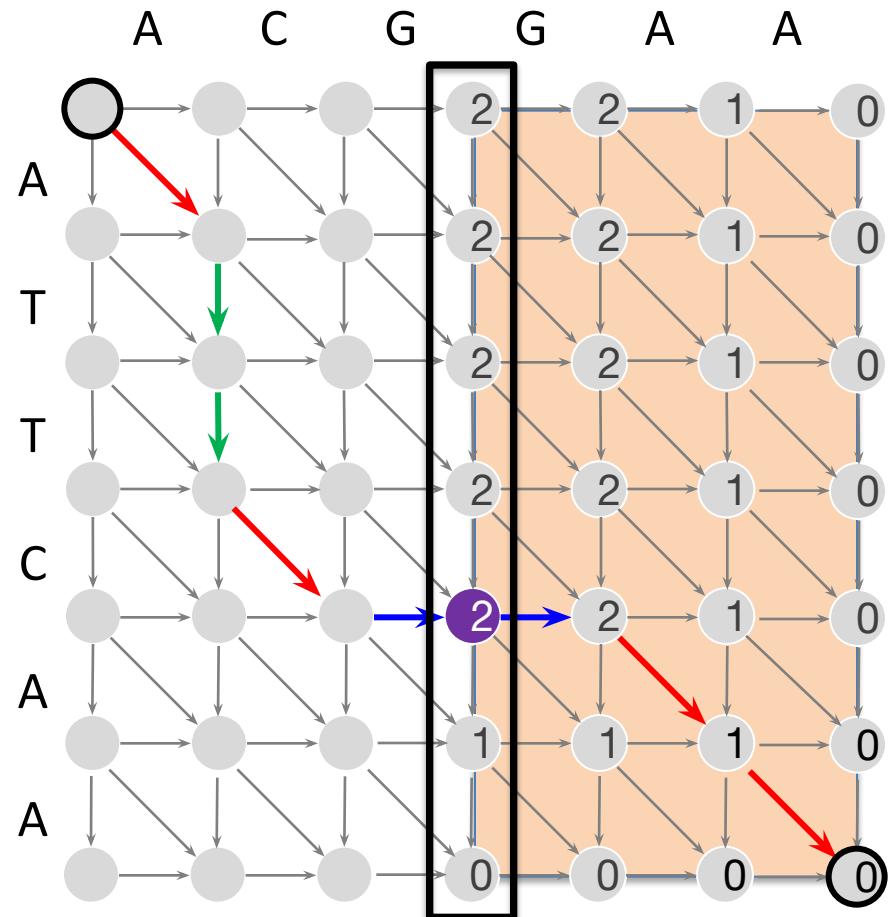
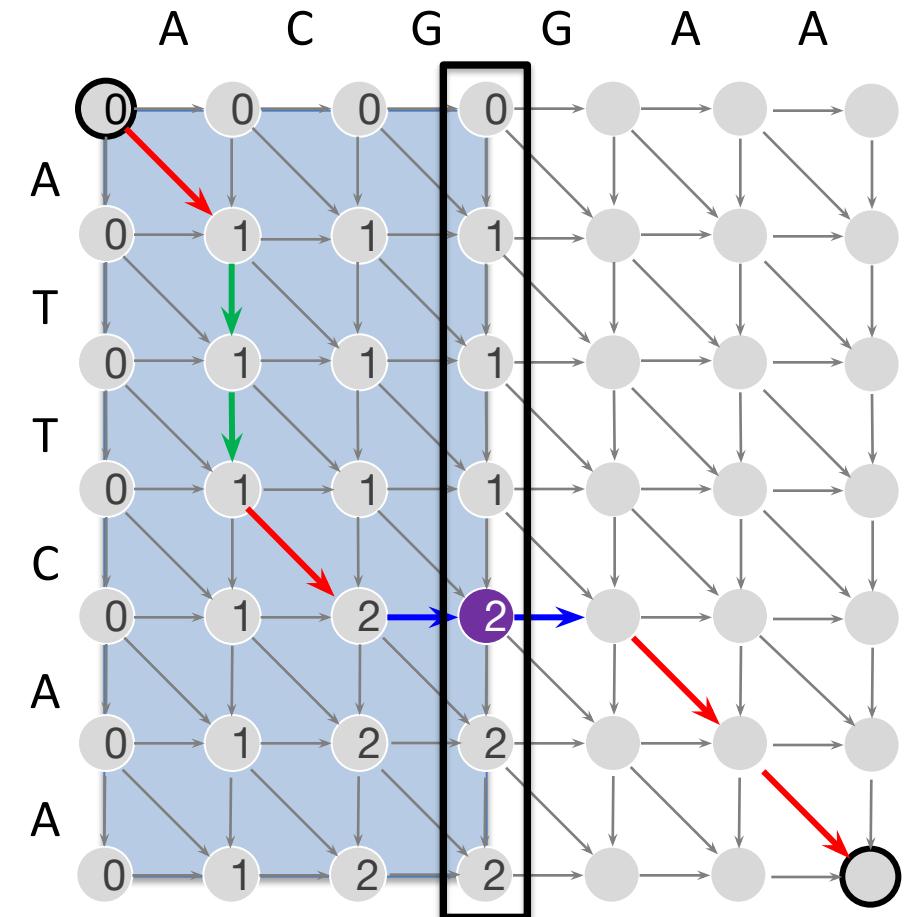


$$length(i) = fromSource(i) + toSink(i)$$

Bioinformatics Algorithms: An Active Learning Approach.

Copyright 2018 Compeau and Pevzner.

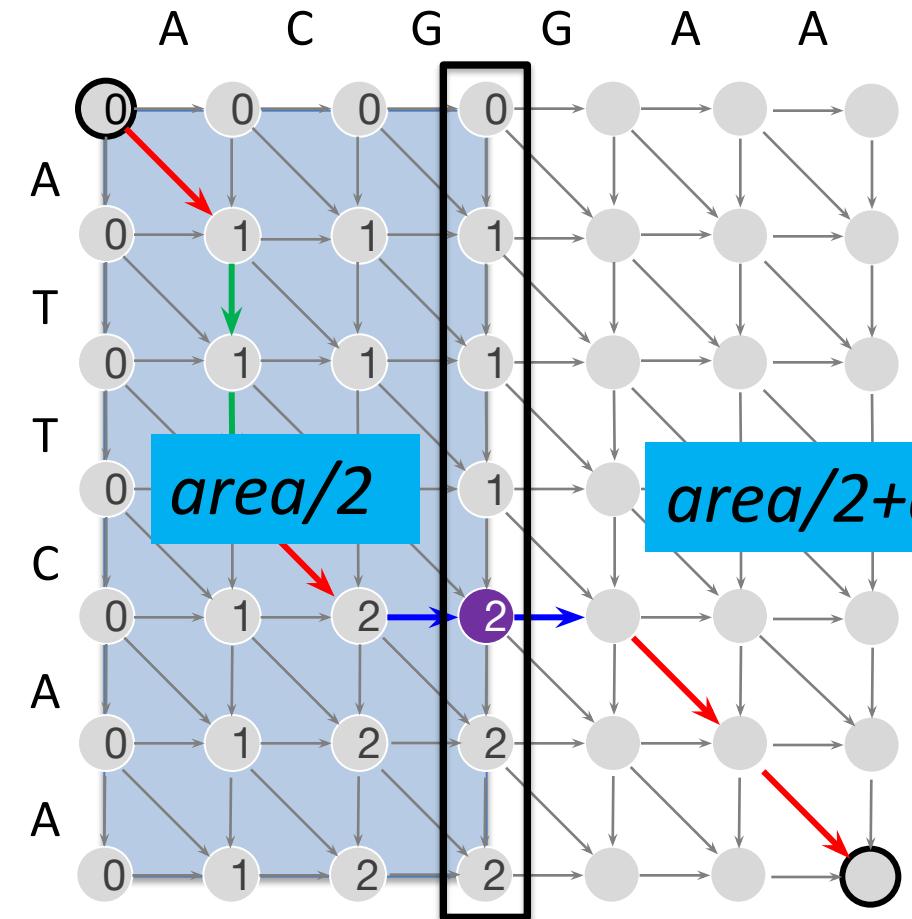
Computing *FromSource* and *toSink*



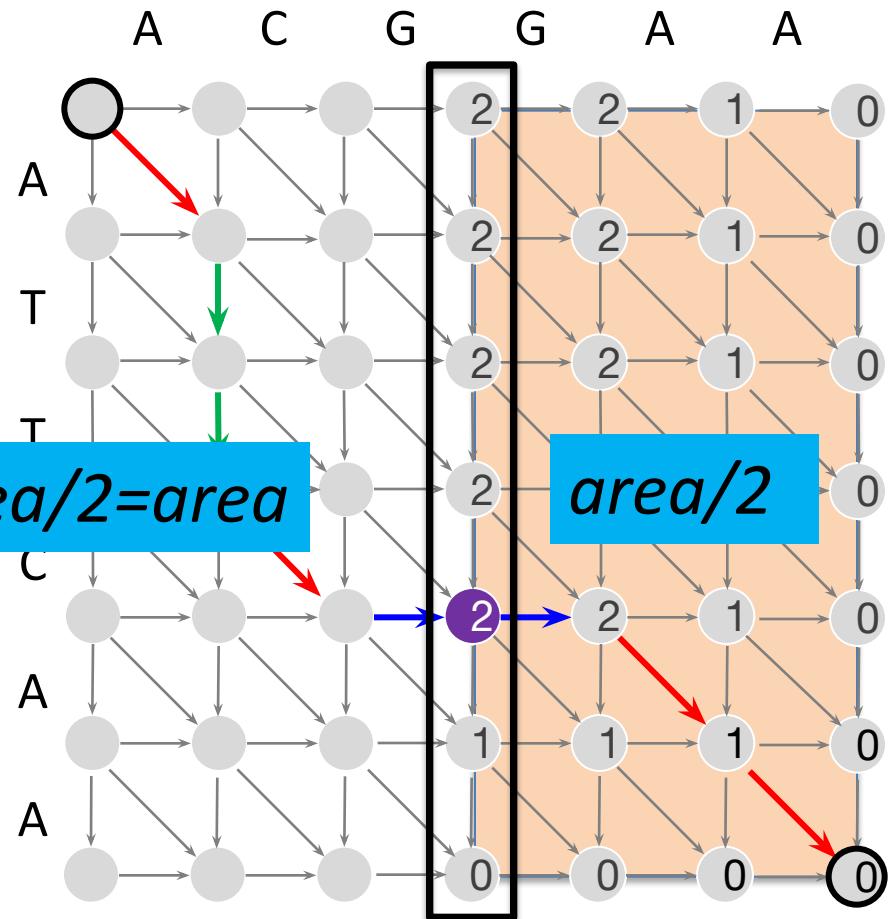
fromSource(i)

toSink(i)

How Much Time Did It Take to Find the Middle Node ?

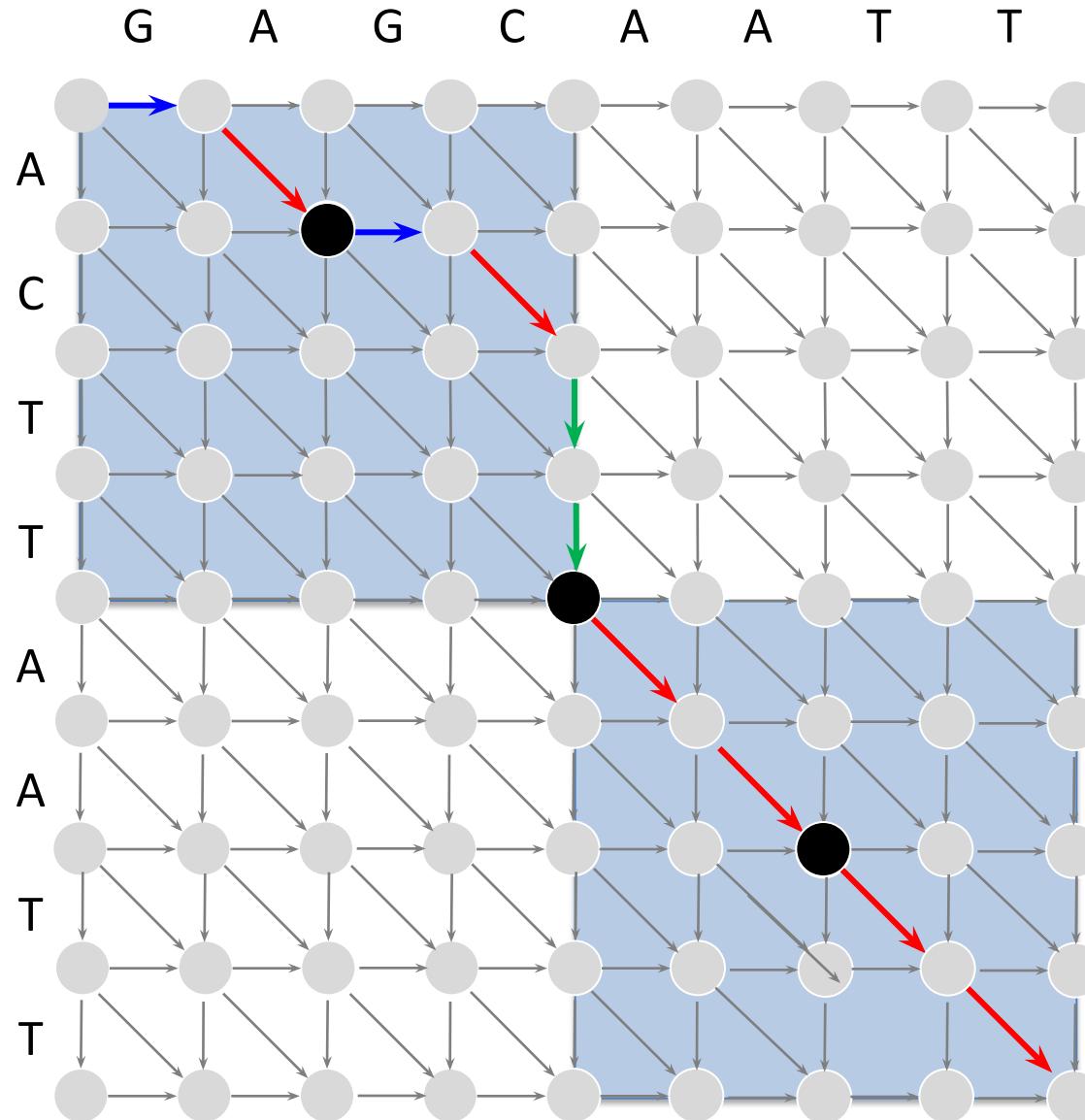


fromSource(i)



toSink(i)

Laughable Progress: $O(nm)$ Time to Find **ONE** Node!

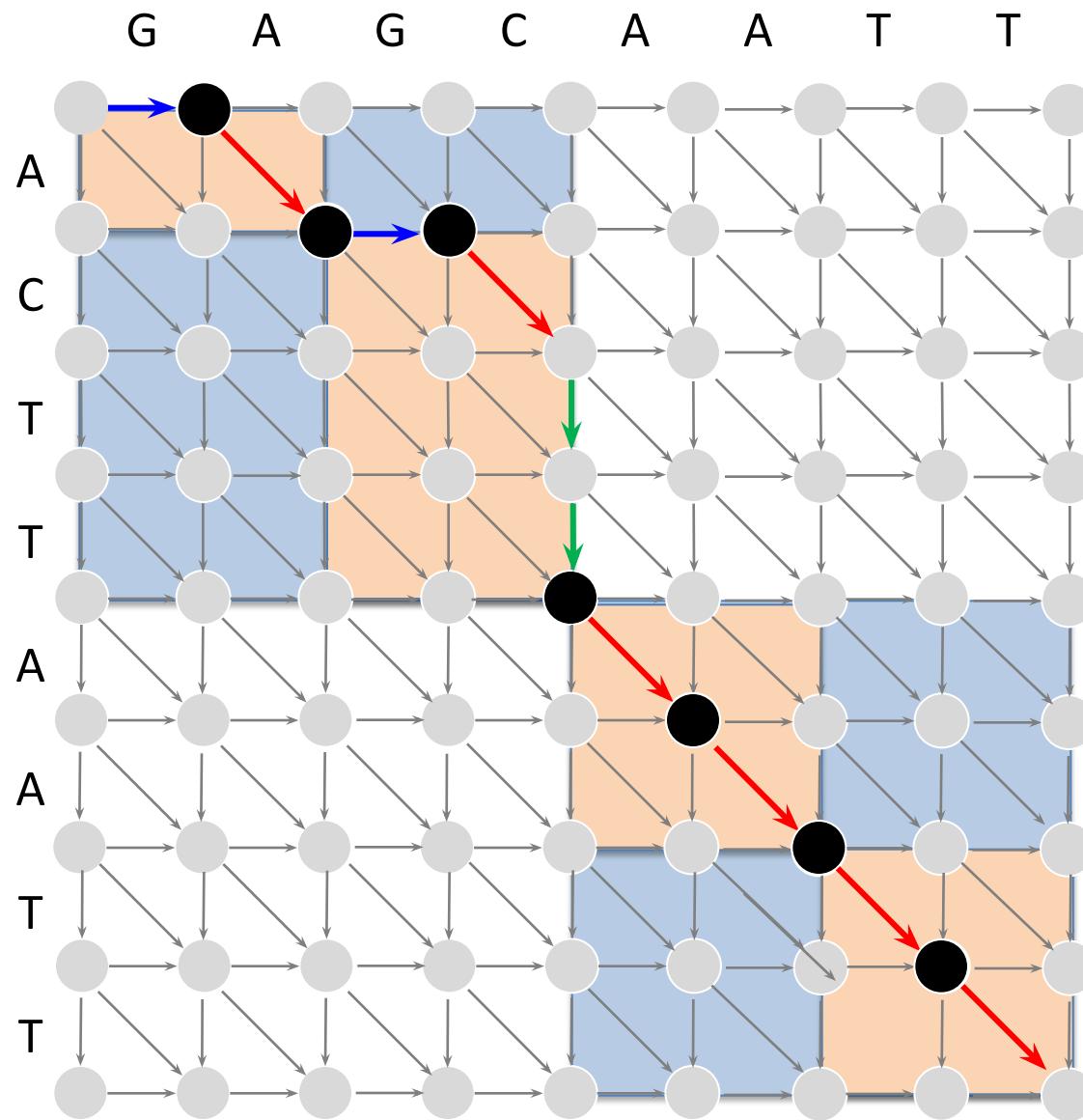


Each subproblem
can be conquered
in time
proportional to its
area:

$$\text{area}/4 + \text{area}/4 = \text{area}/2$$

How much time would it take to conquer 2 subproblems?

Laughable Progress: $O(nm+nm/2)$ Time to Find THREE Nodes!

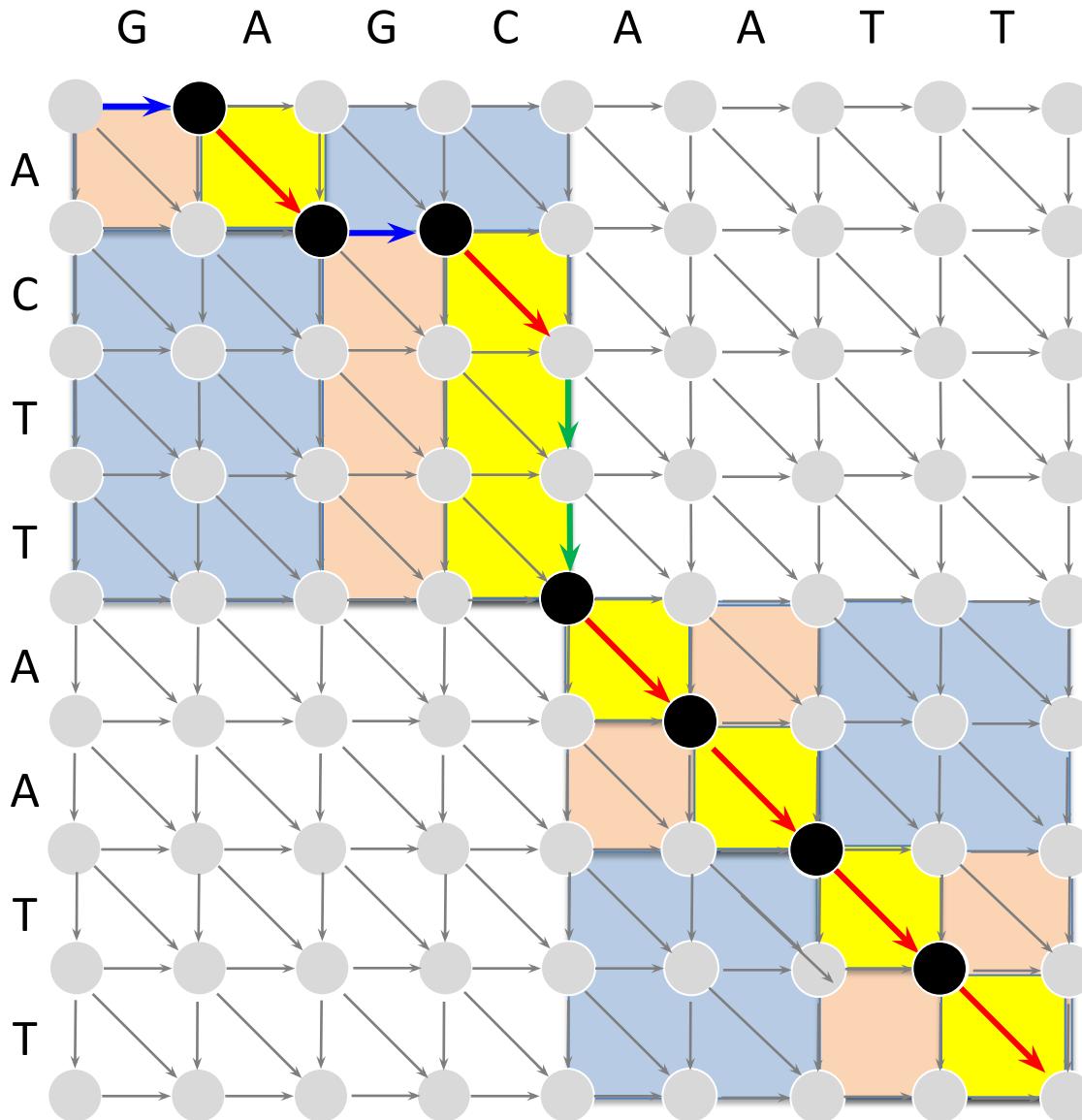


Each subproblem
can be conquered
in time
proportional to its
area:

$$\text{area}/8 + \text{area}/8 + \text{area}/8 + \text{area}/8 = \text{area}/4$$

How much time would it take to conquer 4 subproblems?

$O(nm+nm/2+nm/4)$ Time to Find NEARLY ALL Nodes!



$area +$
 $area/2 +$
 $+ area/4 +$
 $+ area/8 +$
 $+ area/16 +$
 $\dots +$
 $<$
 $2 \cdot area$

How much time would it take to conquer ALL subproblems?

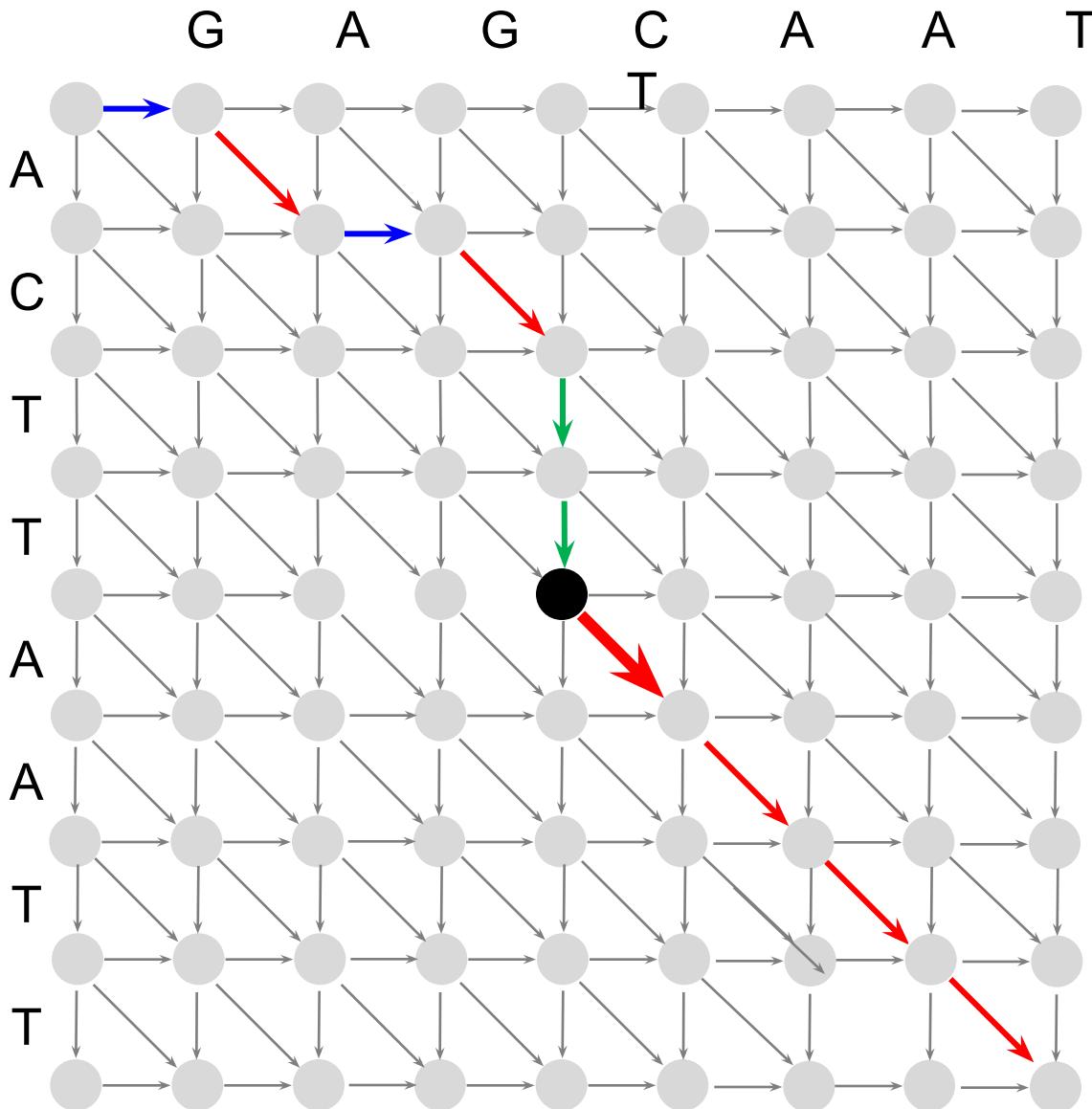
Total Time:

$$area + area/2 + area/4 + area/8 + area/16 + \dots$$



$$1 + \frac{1}{2} + \frac{1}{4} + \dots < 2$$

The Middle Edge

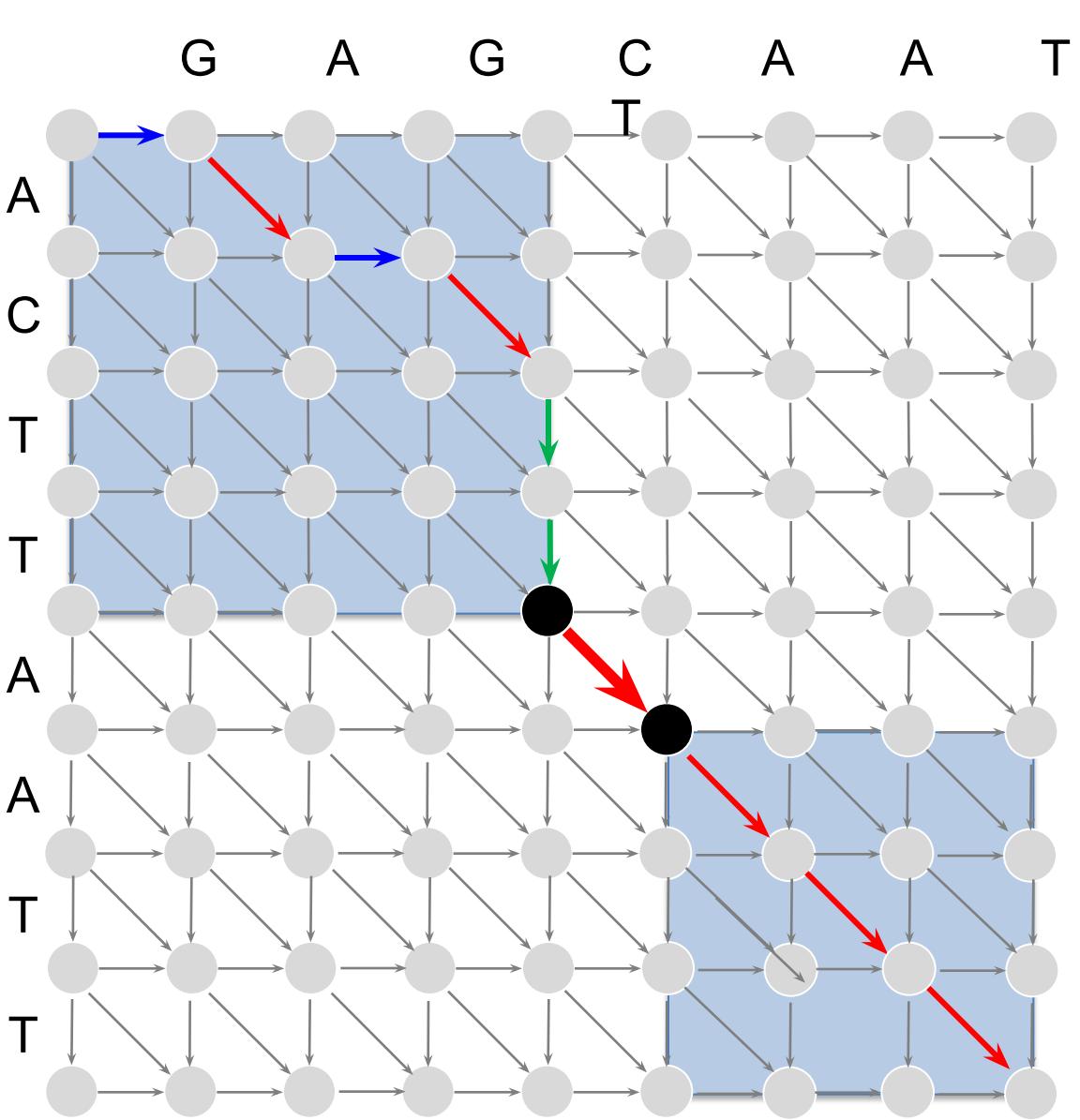


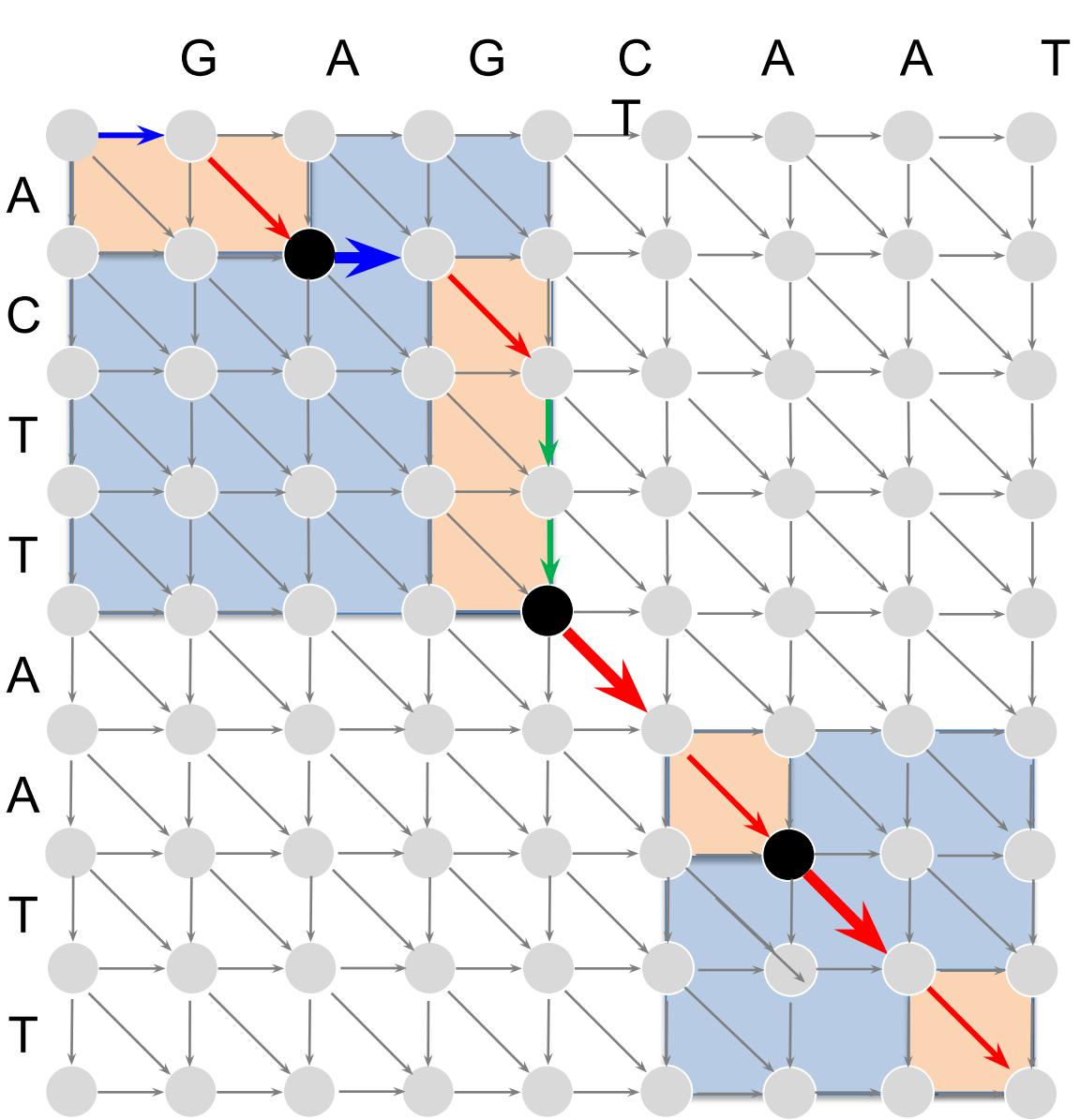
Middle Edge:
an edge in an
optimal
alignment path
starting at the
middle node

The Middle Edge Problem

Middle Edge in Linear Space Problem. Find a middle edge in the alignment graph in linear space.

- **Input:** Two strings and matrix *score*.
- **Output:** A middle edge in the alignment graph of these strings (as defined by the matrix *score*).





Recursive LinearSpaceAlignment

LinearSpaceAlignment(*top, bottom, left, right*)

if *left* = *right*

return alignment formed by *bottom-top* edges “↓”

middle $\leftarrow \lfloor (left+right)/2 \rfloor$

midNode \leftarrow **MiddleNode**(*top, bottom, left, right*)

midEdge \leftarrow **MiddleEdge**(*top, bottom, left, right*)

LinearSpaceAlignment(*top, midNode, left, middle*)

output *midEdge*

if *midEdge* = “→” **or** *midEdge* = “↖”

middle \leftarrow *middle*+1

if *midEdge* = “↓” **or** *midEdge* = “↙”

midNode \leftarrow *midNode*+1

LinearSpaceAlignment(*midNode, bottom, middle, right*)

How Do We Compare Biological Sequences

- From Sequence Comparison to Biological Insights
- The Alignment Game and the Longest Common Subsequence
- The Manhattan Tourist Problem
- The Change Problem
- Dynamic Programming and Backtracking Pointers
- From Manhattan to an Arbitrary Directed Acyclic Graph
- From Global to Local Alignment
- Penalizing Insertions and Deletions in Sequence Alignment
- Space-Efficient Sequence Alignment
- **Multiple Sequence Alignment**

From Pairwise to Multiple Alignment

- Up until now we have only tried to align two sequences.
- A faint (and statistically insignificant) similarity between two sequences becomes significant if it is present in many other sequences.
- Multiple alignments can reveal subtle similarities that pairwise alignments do not reveal.



Alignment of Three A-domains

Y**A**FDLGYTCMPV**L**GGELHIVQKETYTAPDEIAHYI**K**EHG**I**TYIKLT**PS**LFTIVNTASFAFDANFES**L**R LIVLGGEKIIPIDVIAFRKMY**G**HTE-F**N**HYG**P**TEATIGA
-**A**FDVSAGDFARAL**L**LT**GG**QLIVCPNEVKMDPASLYAI**I****K**KYD**I**TIFEA**T****P**ALVIPLMEYI-YEQKLDISQLQILIVGSDSCS**M**EDFKTLVSR**F****G**STIRIV**N****S****Y****G****V****T**EA**C**IDS
I**A**FDASSWEIYAP**L**LN**GG**TVVCIDYYTTIDIKALEAVF**K**QHH**I**RGAMLP**P**ALLKQCLVSA---PTMISS**L**EILFAAGDRLSSQ**D**AILARRAV**G**SGV-Y-**N****A**Y**G****P**TENTVLS

Alignment of Three A-domains

Y**A**FDLGYTCMF~~PV~~**LGG**ELHIVQKETYTAPDEIAHYI**K**EHG**I**TYIKLT**PSL**FHTIVNTASFAFDANFES**L**R~~L~~IVLGGEKIIPIDVIAFRKMY**G**HTE-F**NHYG**P**T**EATIGA
-**A**FD**V**SAGDFARALL**TGG**QLIVCPNEVKMDPASLY**A**II**K**KYD**I**TIFEA**T****P**ALVIPLMEYI-YEQKLD**I****S**QL**Q****I**LIVGSDSCSMEDFKTLVSRF**G****S**TIRIV**N****S****Y****G****V****T**EACIDS
I**A**FD**A**S**S**WEIYAP**LLNGG**TVVCIDYYTTIDIKALE**A****V****F****K**QHH**I**RGAMLP**P****A****L****L****K**QCLVSA---PTM**I****S****S****L****E****I****L****F**AAAGDRLSQ**D****A****I****L****A****R****R****A****V****G****S****G****V**-Y-**NAYG**P**T**E**N****T****V****L****S**

Alignment of Three A-domains

YAFDLYTCMFPVLLGGELHIVQKETYTAPDEIAHYIKEHGITYIKLTPLSFHTIVNTASFAFDANFESLRLIVLGGEKIIPIIDVIAFRKMYGHTE-FINHYGPTEATIGA
-AFDVSAAGDFARALLTGGQLIVCPNEVKMDPASLYAIKKYDITIFEATPALVIPLMEYI-YEQKLDISQLQILIVGSDSCSMEDFKTLVSRFGSTIRIVNSYGVTEACIDS
IAFDASSWEIYAPLLNGGTVVCIDYYTTIDIKALEAVFKQHHIRGAMLP~~P~~ALLKQCLVSA---PTMISSLEILFAAGDRLSQDAILARRAVGSGV-Y-NAYGPENTVLS

Generalizing Pairwise to Multiple Alignment

- Alignment of 2 sequences is a 2-row matrix.
- Alignment of 3 sequences is a 3-row matrix

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

- Our scoring function should score alignments with conserved columns higher.

Alignments = Paths in 3-D

- Alignment of ATGC, AATC, and ATGC

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

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

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

Alignments = Paths in 3-D

- Alignment of ATGC, AATC, and ATGC

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

#symbols up to a given position

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

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

Alignments = Paths in 3-D

- Alignment of ATGC, AATC, and ATGC

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

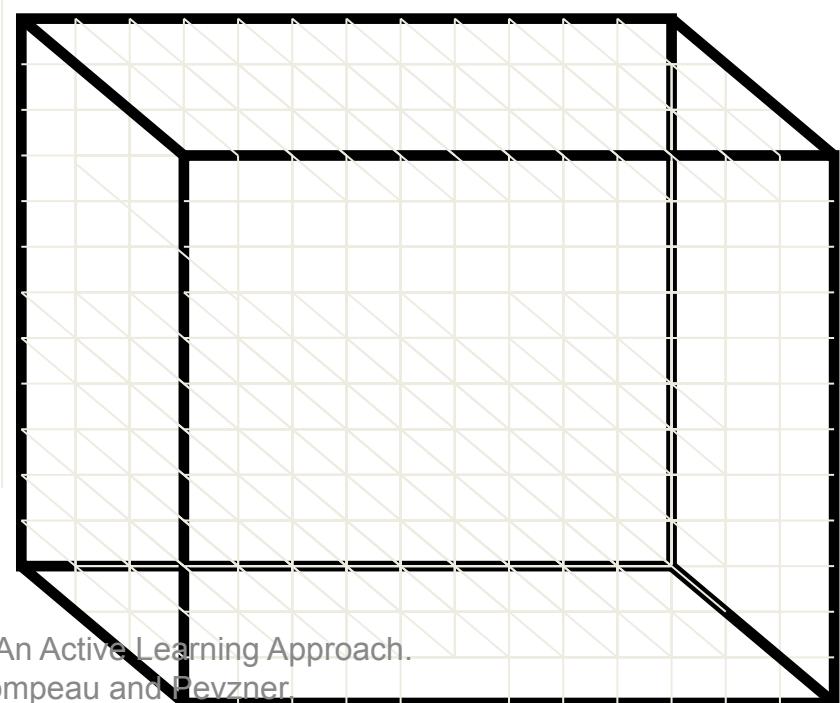
#symbols up to a given position

Alignments = Paths in 3-D

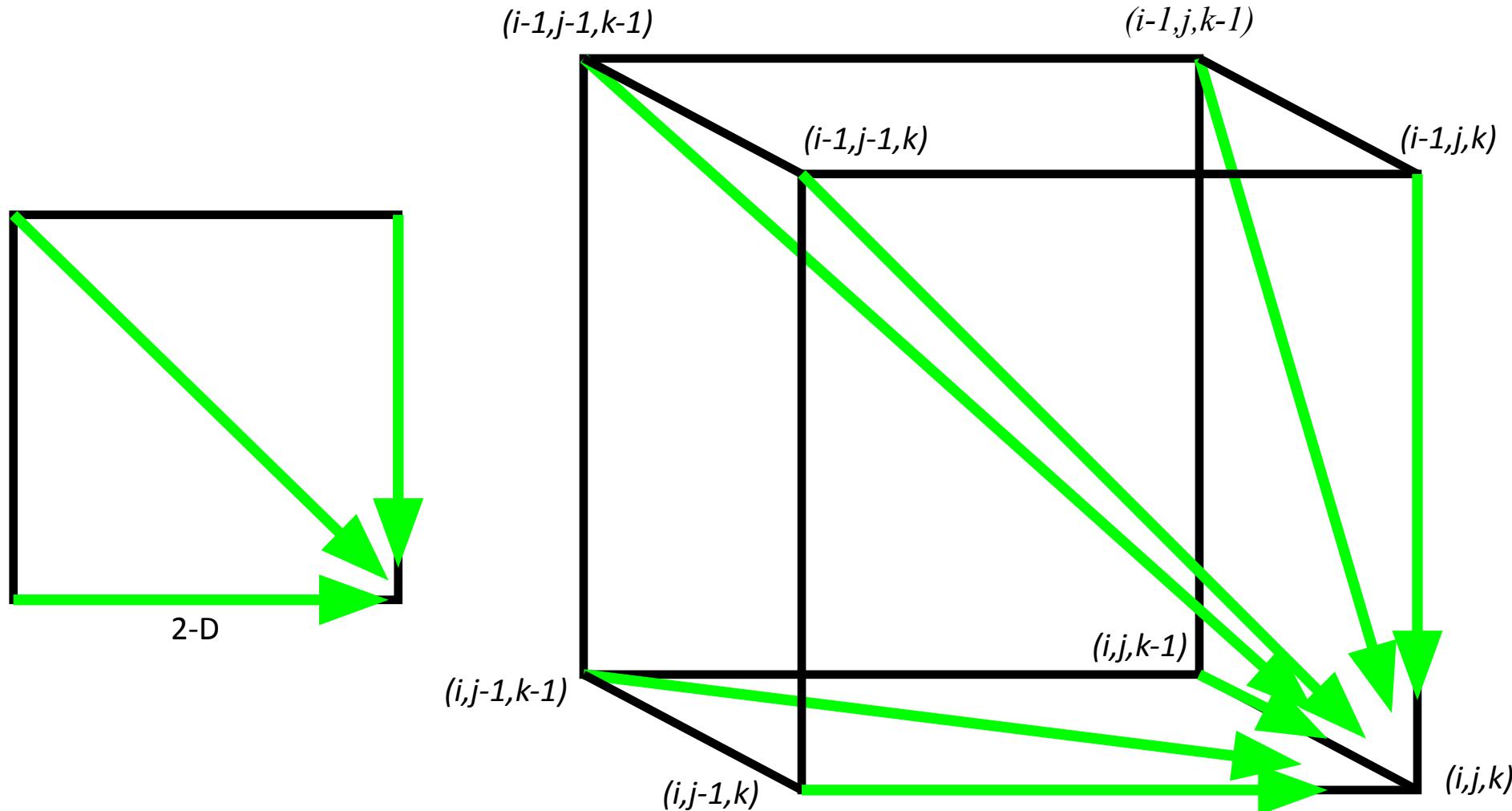
- Alignment of ATGC, AATC, and ATGC

$(0,0,0) \rightarrow (1,1,0) \rightarrow (1,2,1) \rightarrow (2,3,2) \rightarrow (3,3,3) \rightarrow (4,4,4)$

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



2-D Alignment Cell versus 3-D Alignment Cell



Multiple Alignment: Dynamic Programming

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

- $\delta(x, y, z)$ is an entry in the 3-D scoring matrix.

Multiple Alignment: Running Time

- For 3 sequences of length n , the run time is proportional to $7n^3$
- For a k -way alignment, build a k -dimensional Manhattan graph with
 - n^k nodes
 - most nodes have $2^k - 1$ incoming edges.
 - Runtime: $O(2^k n^k)$

Multiple Alignment Induces Pairwise Alignments

Every multiple alignment induces pairwise alignments:

AC-GCGG-C

AC-GC-GAG

GCCGC-GAG



ACGCGG-C

AC-GCGG-C

AC-GCGAG

ACGC-GAC

GCCGC-GAG

GCCGCGAG

Idea: Construct Multiple from Pairwise Alignments

Given a set of **arbitrary** pairwise alignments, can we construct a multiple alignment that induces them?

AAAATTTTT----

----**TTTG**GGGG

----**AAAAT**TTT

GGGAAAAA----

TTT**GGG**----

----**GGGAA**AA

Profile Representation of Multiple Alignment

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G
A	0	1	0	0	0	0	1	0	0	.8	0	0	0
C	.6	0	0	0	1	0	0	.4	1	0	.6	.2	0
G	0	0	1	.2	0	0	0	0	.2	0	0	.4	1
T	.2	0	0	0	0	1	0	.6	0	0	0	.2	0
-	.2	0	0	.8	0	0	0	0	0	.4	.8	.4	0

Aligning Sequence Against Sequence

- In the past we were aligning a **sequence against a sequence**.

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G
A	0	1	0	0	0	0	1	0	0	.8	0	0	0
C	.6	0	0	0	1	0	0	.4	1	0	.6	.2	0
G	0	0	1	.2	0	0	0	0	.2	0	0	.4	1
T	.2	0	0	0	0	1	0	.6	0	0	0	.2	0
-	.2	0	0	.8	0	0	0	0	0	.4	.8	.4	0

Aligning Sequence Against Profile

- In the past we were aligning a **sequence against a sequence**.
 - Can we align a **sequence against a profile**?

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G
A	0	1	0	0	0	0	1	0	0	.8	0	0	0
C	.6	0	0	0	1	0	0	.4	1	0	.6	.2	0
G	0	0	1	.2	0	0	0	0	.2	0	0	.4	1
T	.2	0	0	0	0	1	0	.6	0	0	0	.2	0
-	.2	0	0	.8	0	0	0	0	0	.4	.8	.4	0

Aligning Profile Against Profile

- In the past we were aligning a **sequence against a sequence**.
 - Can we align a **sequence against a profile**?
 - Can we align a **profile against a profile**?

-	A	G	G	C	T	A	T	C	A	C	C	T	G
T	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	C	C	A	-	-	-	G
C	A	G	-	C	T	A	T	C	A	C	-	G	G
C	A	G	-	C	T	A	T	C	G	C	-	G	G
A	0	1	0	0	0	1	0	0	.8	0	0	0	0
C	.6	0	0	0	1	0	0	.4	1	0	.6	.2	0
G	0	0	1	.2	0	0	0	0	.2	0	0	.4	1
T	.2	0	0	0	0	1	0	.6	0	0	0	.2	0
-	.2	0	0	.8	0	0	0	0	.4	.8	.4	.4	0

Multiple Alignment: Greedy Approach

- Choose the most similar sequences and combine them into a profile, thereby reducing alignment of k sequences to an alignment of $k - 2$ sequences and 1 profile.
- Iterate

Greedy Approach: Example

- Sequences: GATTCA, GTCTGA, GATATT, GTCAGC.
- 6 pairwise alignments (premium for **match** +1, penalties for **indels** and mismatches -1)

s_2 GTCTGA

s_4 GTCAGC (score = 2)

s_1 GATTCA--

s_4 G-T-CAGC (score = 0)

s_1 GAT-TCA

s_2 G-TCTGA (score = 1)

s_2 G-TCTGA

s_3 GATAT-T (score = -1)

s_1 GAT-TCA

s_3 GATAT-T (score = 1)

s_3 GAT-ATT

s_4 G-TCAGC (score = -1)

Greedy Approach: Example

- Since s_2 and s_4 are closest, we consolidate them into a profile:

$$\left. \begin{array}{l} s2\text{GTCTGA} \\ s4\text{GTCAGC} \end{array} \right\} s_{2,4} = \text{GTCt/aGa/cA}$$

- New set of 3 sequences to align:

s_1	GATTCA
s_3	GATATT
$s_{2,4}$	GTCt/aGa/c