

BIOINFORMATICS

How do we compare biological sequences?

Marco Beccuti

Università degli Studi di Torino

Dipartimento di Informatica

March 2019



Outline

- 1 Introduction to Sequence Alignment
- 2 Hamming distance for similarity between sequences
- 3 Alignment Game and the Longest Common Subsequence;
- 4 The Manhattan Tourist Problem;
- 5 The Change Problem;
- 6 Dynamic programming and backtracking pointers;
- 7 From Manhattan to Alignment Graph;
- 8 From Global to Local Alignment;
- 9 Penalizing Insertions and Deletions in Sequence Alignment;
- 10 Space-Efficient Sequence Alignment;
- 11 Multiple Sequence Alignment.

Chapter 5 in *Bioinformatics Algorithms: An active Learning Approach (Vol.1)*.



Part 2

Dynamic Programming and Backtracking Pointers

Dynamic Programming and Backtracking Pointers

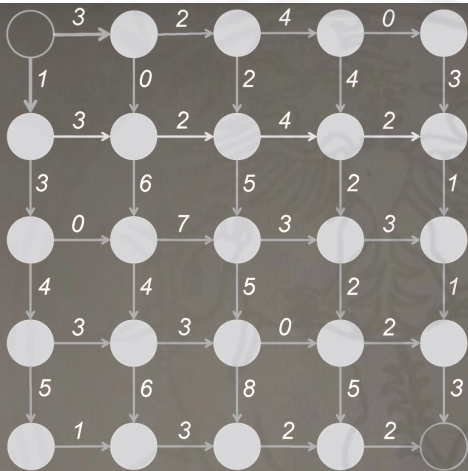
- We come back to Longest Path Problem in a grid;

There are only 2
ways to arrive to
the sink:

by moving
South ↓

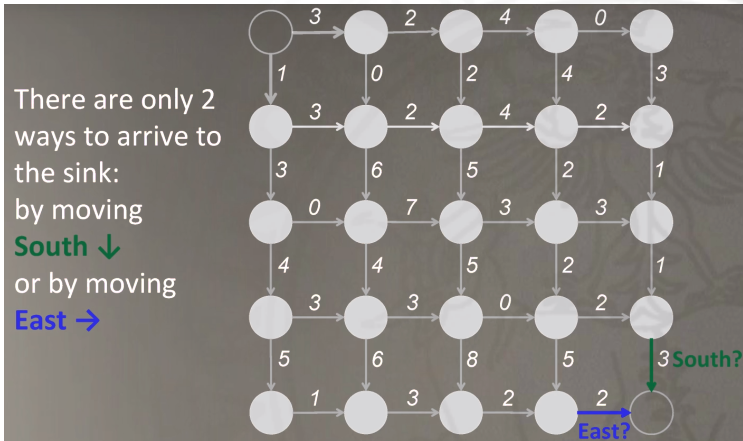
or by moving

East →



Dynamic Programming and Backtracking Pointers

- We come back to Longest Path Problem in a grid;

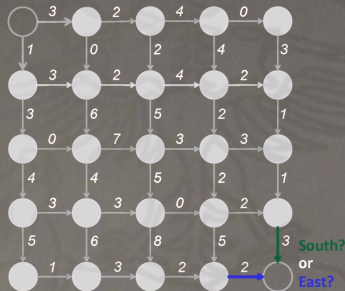


Dynamic Programming and Backtracking Pointers

- We come back to Longest Path Problem in a grid;

South or East?

SouthOrEast(i,j):
the length of the longest path from
(0,0) to (i,j)



SouthOrEast(n,m)=

$$\text{MAX} \begin{cases} \text{SouthOrEast}(n-1,m) + \text{weight of edge } \downarrow \text{ into } (n,m) \\ \text{SouthOrEast}(n,m-1) + \text{weight of edge } \rightarrow \text{ into } (n,m) \end{cases}$$

Dynamic Programming and Backtracking Pointers

- A recursive algorithm can easily define to compute the Longest Path Problem in a grid;

```
SouthOrEast(i,j)
```

```
  If i=0 and j=0
```

```
    return 0
```

```
  x = -infinity, y = -infinity
```

```
  If i>0
```

```
    x ← SouthOrEast(i-1,j) + weight of the vertical edge into (i,j)
```

```
  If j>0
```

```
    y ← SouthOrEast(i,j-1) + weight of the horizontal edge into (i,j)
```

```
  return max{x,y}
```

Dynamic Programming and Backtracking Pointers

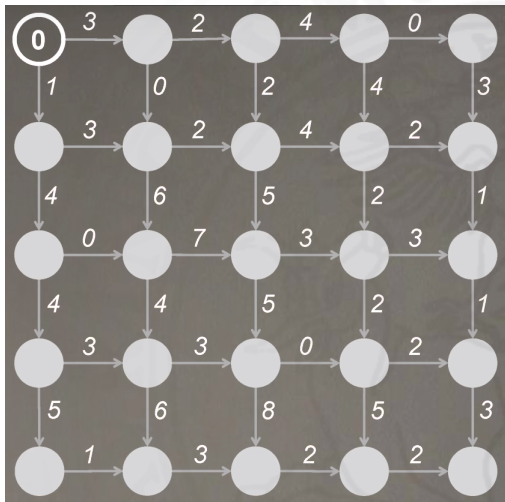
- A recursive algorithm can easily define to compute the Longest Path Problem in a grid;

```
SouthOrEast(i,j)
If i=0 and j=0
  return 0
x= -infinity, y= -infinity
If i>0
  x ← SouthOrEast(i-1,j) + weight of the vertical edge into (i,j)
If j>0
  y ← SouthOrEast(i,j-1) + weight of the horizontal edge into (i,j)
return max{x,y}
```

- It visits all the possible paths: ***it is correct. but it is too expensive!!***

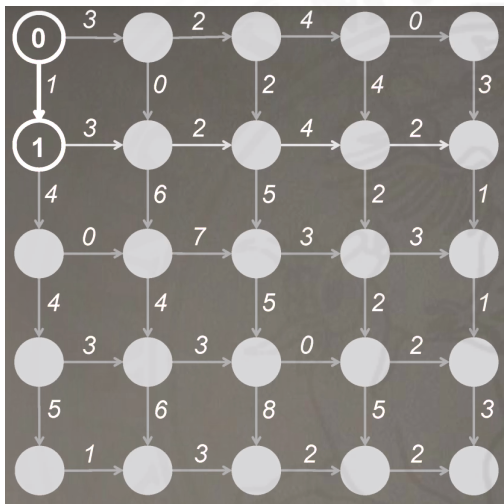
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



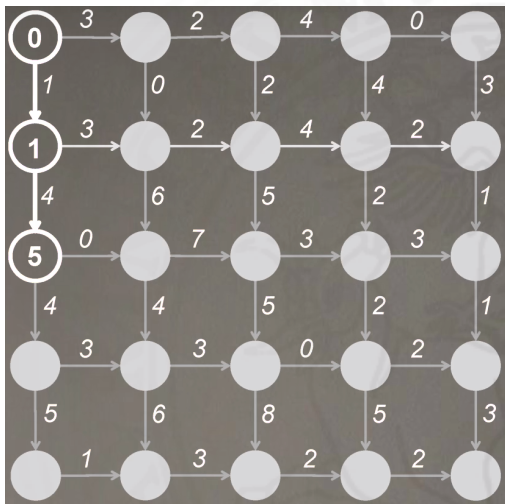
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



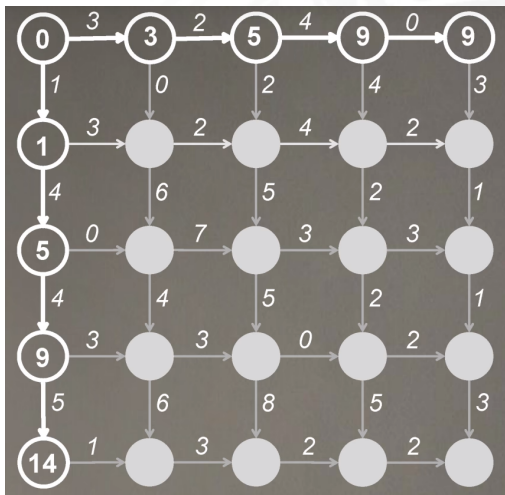
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



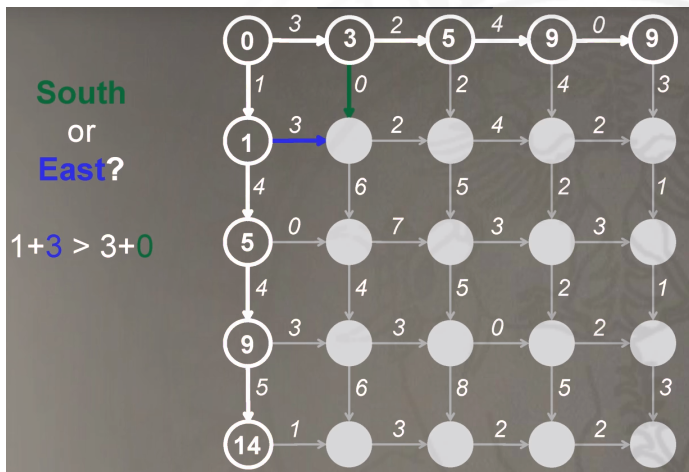
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



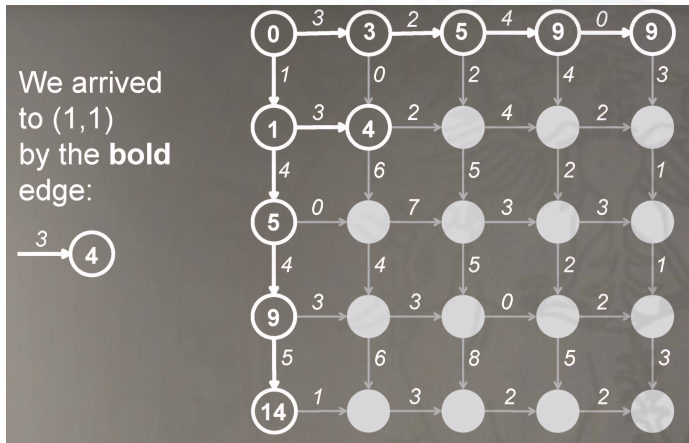
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



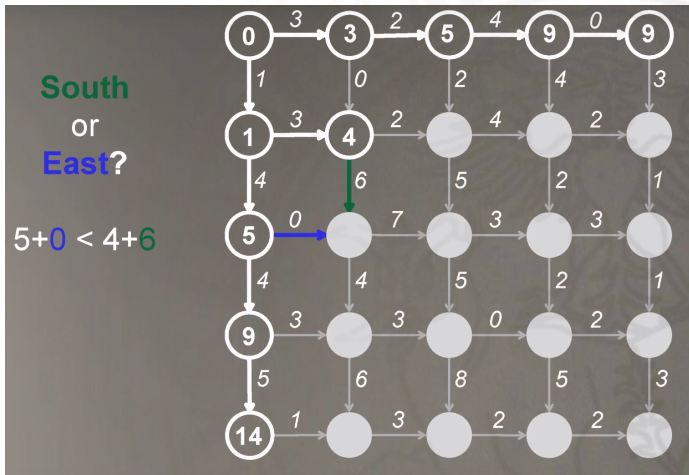
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



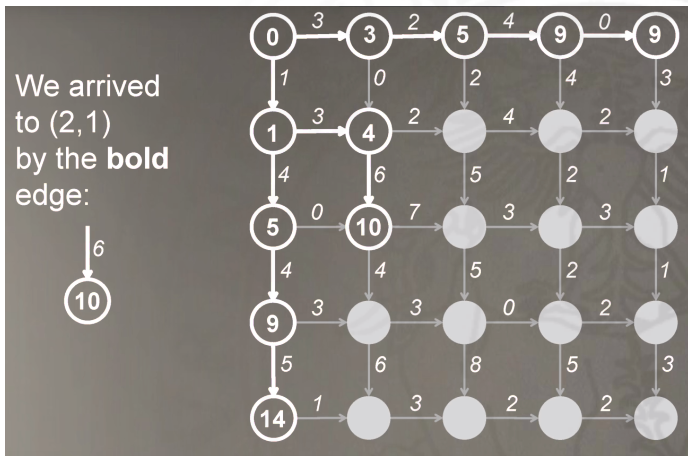
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



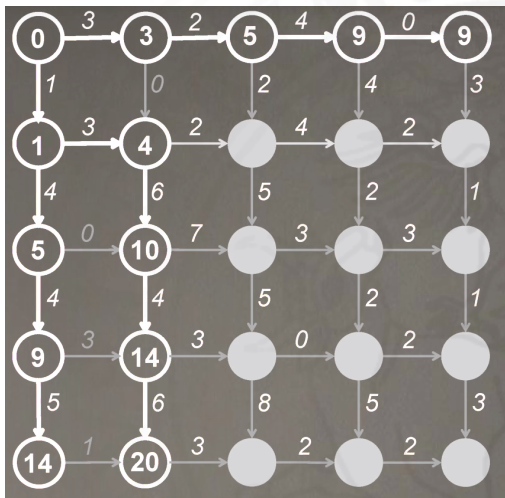
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



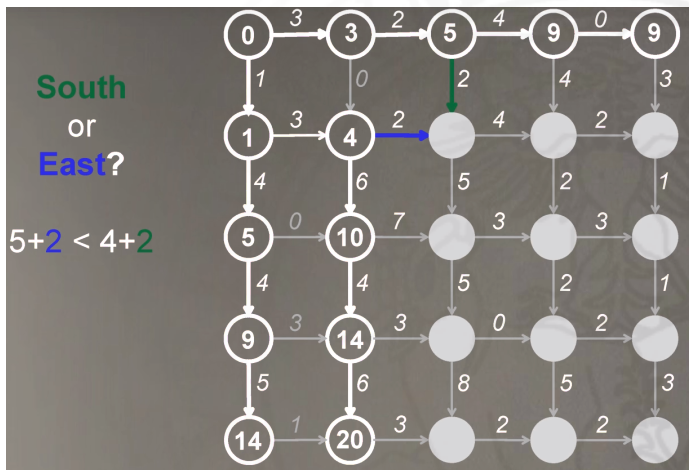
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



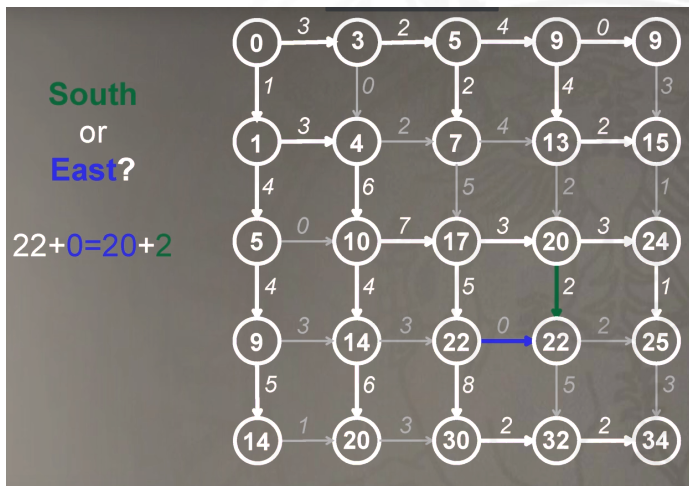
Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



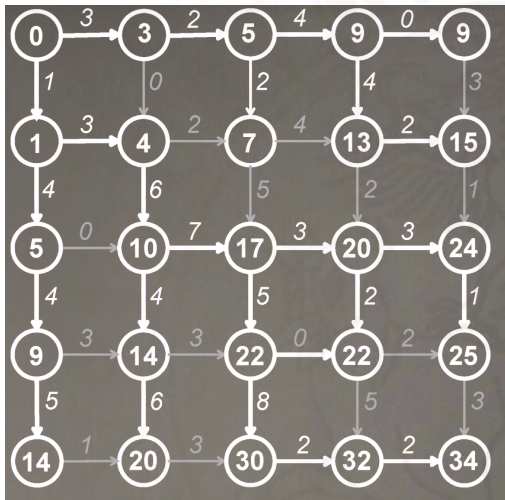
Dynamic Programming and Backtracking Pointers

- We can have cases in which both choices provide a similar results;



Dynamic Programming and Backtracking Pointers

- Now, we exploit Dynamic Programming to compute Longest Path in grid;



Dynamic Programming and Backtracking Pointers

- Now, we can define the Dynamic Programming Recurrence as follows:

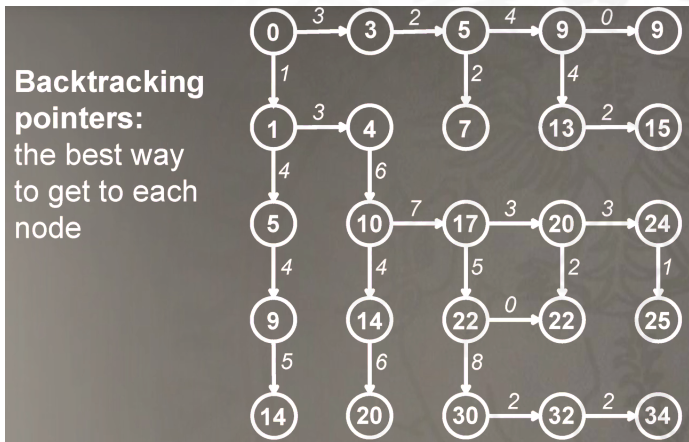
$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge } \downarrow \text{ into } (i,j) \\ s_{i,j-1} + \text{weight of edge } \rightarrow \text{ into } (i,j) \end{cases}$$

where $s_{i,j}$ is the length of the longest path from $(0,0)$ to (i,j) .

Dynamic Programming and Backtracking Pointers

- Backtracking pointer can be defined as:

the best way to get to each node



- the value in the sink node is the *value of the longest path*.

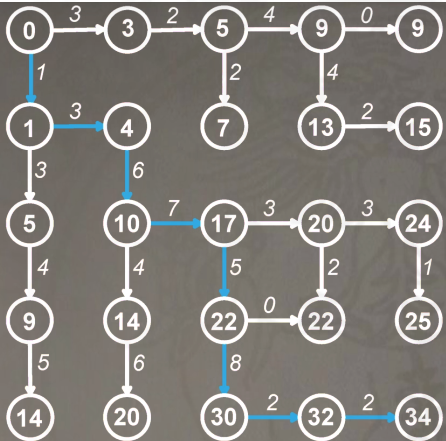
Dynamic Programming and Backtracking Pointers

How can we derive the optimal path from source to sink?

Dynamic Programming and Backtracking Pointers

How can we derive the optimal path from source to sink?

What is the optimal path from *source* to *sink*?





Part 2

From Manhattan to Alignment Graph

From Manhattan to Alignment Graph

- Extending the approach for DAGs;

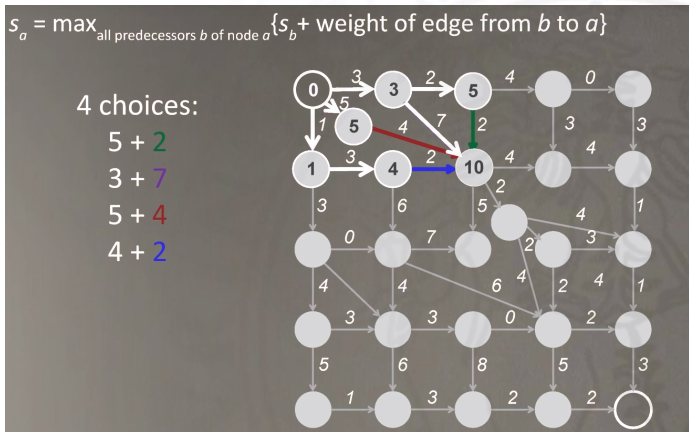
How does the recurrence change for this graph?

The diagram shows a transition from a Manhattan map to a directed acyclic graph (DAG). The Manhattan map on the left shows a grid of streets with a highlighted path. The DAG on the right consists of nodes arranged in a grid, with edges and weights. The highlighted nodes in the DAG are 0, 1, 3, 4, and 5. The graph starts at node 0 and ends at a white circle at the bottom right.

- all the highlighted nodes can be reached in only one-way.

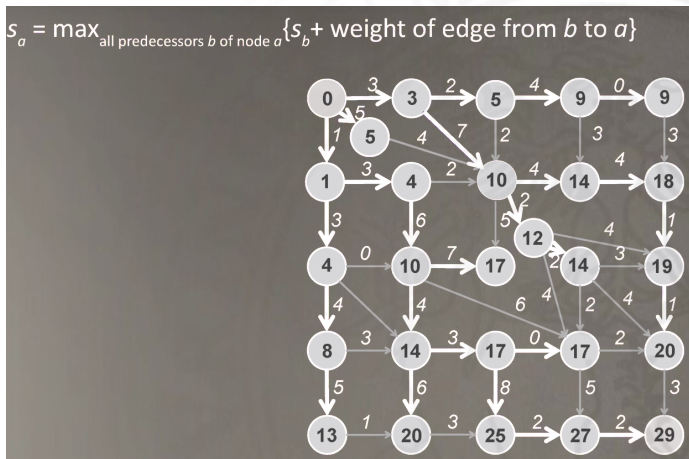
From Manhattan to Alignment Graph

- Extending the approach for DAGs;



From Manhattan to Alignment Graph

- Extending the approach for DAGs;

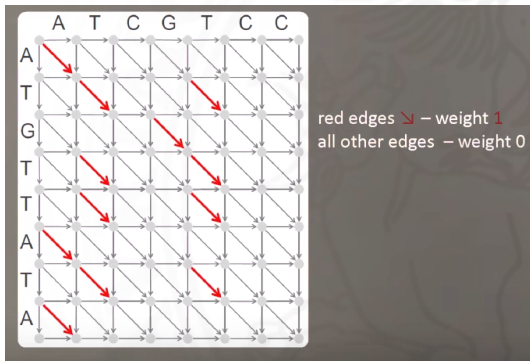


From Manhattan to Alignment Graph

- Dynamic Programming Recurrence for Alignment Graph

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge } \swarrow \text{ into } (i,j) \\ s_{i,j-1} + \text{weight of edge } \rightarrow \text{ into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge } \searrow \text{ into } (i,j) \end{cases}$$

where $s_{i,j}$ is the length of the longest path from $(0,0)$ to (i,j) .

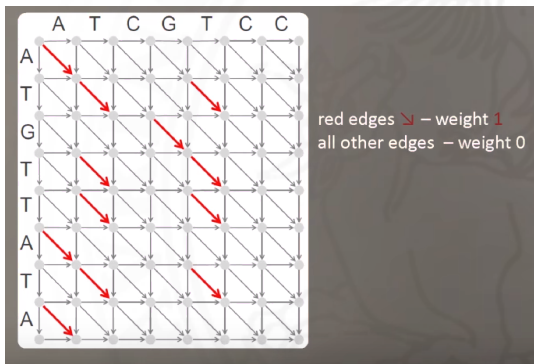


From Manhattan to Alignment Graph

- Dynamic Programming Recurrence for Longest Subsequence Problem

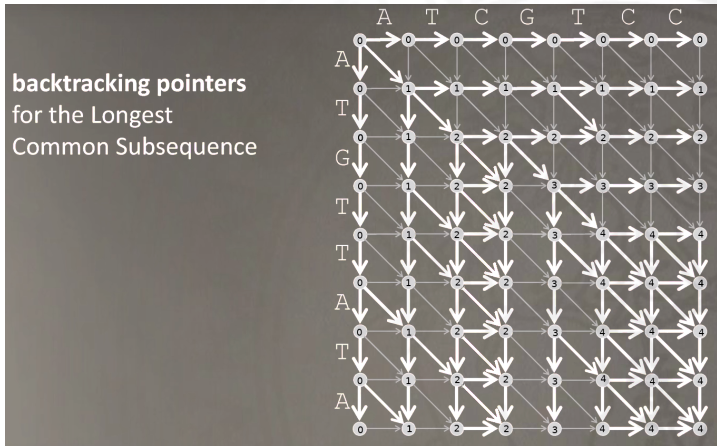
$$s_{i,j} = \max \begin{cases} s_{i-1,j} + 0 \\ s_{i,j-1} + 0 \\ s_{i-1,j-1} + 0 & \text{iff } V_i \neq W_j \\ s_{i-1,j-1} + 1 & \text{iff } V_i = W_j \end{cases}$$

where $s_{i,j}$ is the length of the longest path from $(0,0)$ to (i,j) .



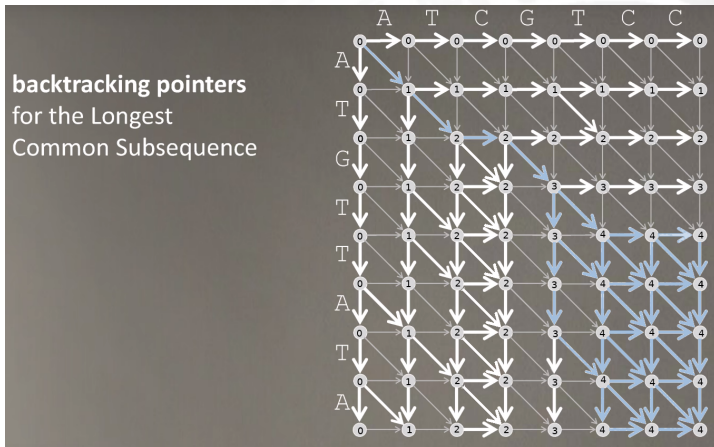
From Manhattan to Alignment Graph

- Dynamic Programming Recurrence for Longest Subsequence Problem

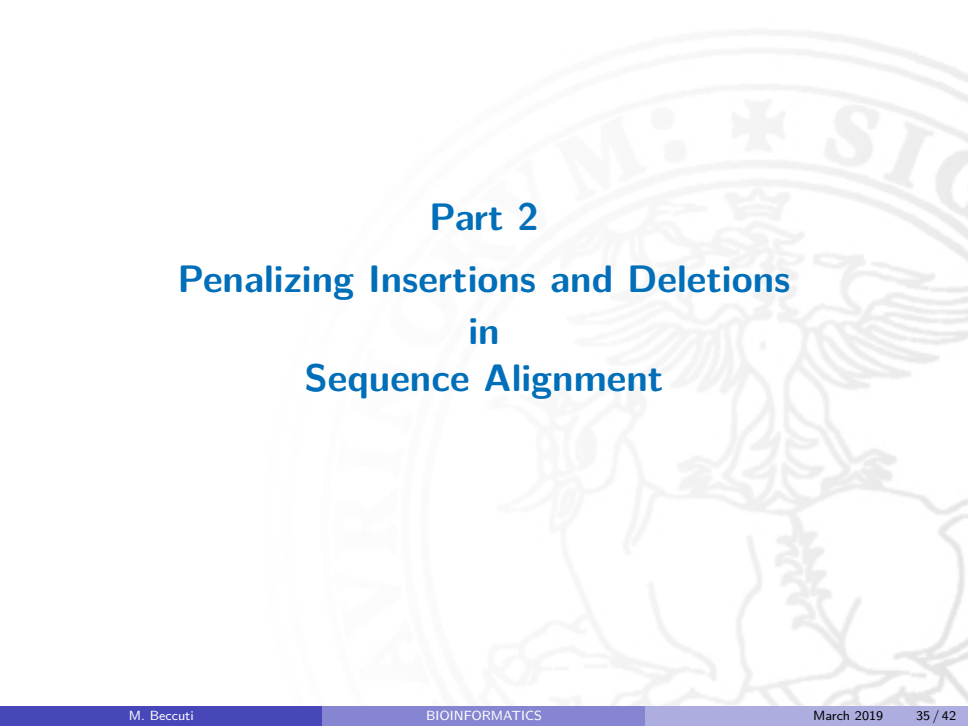


From Manhattan to Alignment Graph

- Dynamic Programming Recurrence for Longest Subsequence Problem



- Blue arrows show the set of optimal alignments



Part 2
Penalizing Insertions and Deletions
in
Sequence Alignment

Penalizing Insertions and Deletions in Sequence Alignment

- it is not difficult to construct an alignment having a lot of matches at the expense of introducing more indels, ...
- but more indels we add, then less biologically relevant the alignment becomes;

How can we cope with this point?

Penalizing Insertions and Deletions in Sequence Alignment

How can we cope with this point?

- we can extend our score function

$\|matches\|$

into

$\|matches\| - \mu\|mismatches\| - \sigma\|indels\|$

so that we take into account penalties associated with mismatches and indels

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

Penalizing Insertions and Deletions in Sequence Alignment

- we can define a scoring matrix as follows:

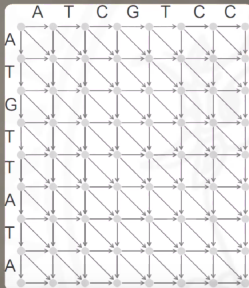
	A	C	G	T	-		A	C	G	T	-
A	+1	- μ	- μ	- μ	- σ	A	+1	-3	-5	-1	-3
C	- μ	+1	- μ	- μ	- σ	C	-4	+1	-3	-2	-3
G	- μ	- μ	+1	- μ	- σ	G	-9	-7	+1	-1	-3
T	- μ	- μ	- μ	+1	- σ	T	-3	-5	-8	+1	-4
-	- σ	- σ	- σ	- σ		-	-4	-2	-2	-1	

scoring matrix scoring matrix with arbitrary values

Penalizing Insertions and Deletions in Sequence Alignment

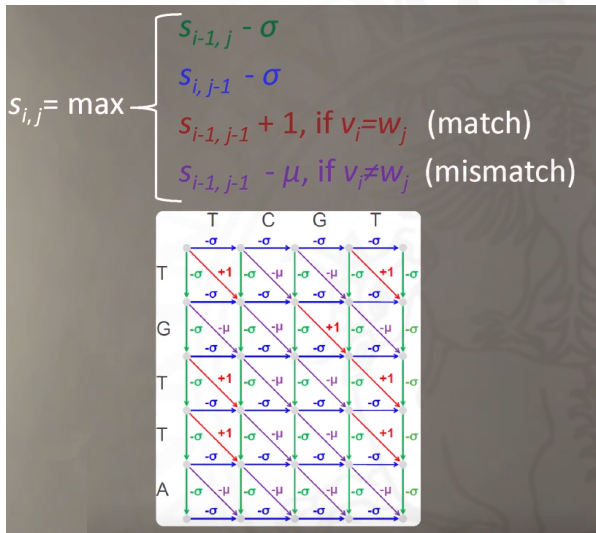
- How to change the dynamic programming recurrence:

$$s_{i,j} = \max \begin{cases} s_{i-1,j} + \text{weight of edge } \downarrow \text{ into } (i,j) \\ s_{i,j-1} + \text{weight of edge } \rightarrow \text{ into } (i,j) \\ s_{i-1,j-1} + \text{weight of edge } \searrow \text{ into } (i,j) \end{cases}$$



Penalizing Insertions and Deletions in Sequence Alignment

- How to change the dynamic programming recurrence:



Penalizing Insertions and Deletions in Sequence Alignment

- How to change the dynamic programming recurrence in general way:

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