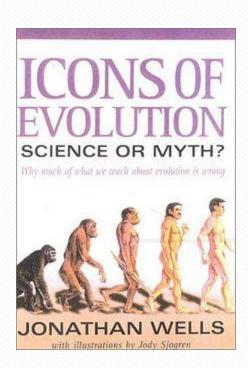


#### Why alignments?

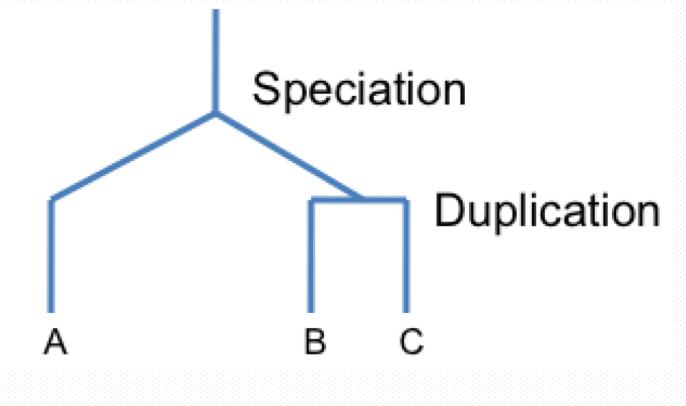
- Detect homology (similarity)
- Study evolution
- Predict functions
- Model 3D-structure



#### Sequence similarity

- High sequence similarity indicates homology
- Homologs have similar 3D-structure
- Homologs have a common ancestor
- Complexity arises due to
  - Speciation
    - Breed vs specie
    - Anole Lizards example video

Gene Duplication

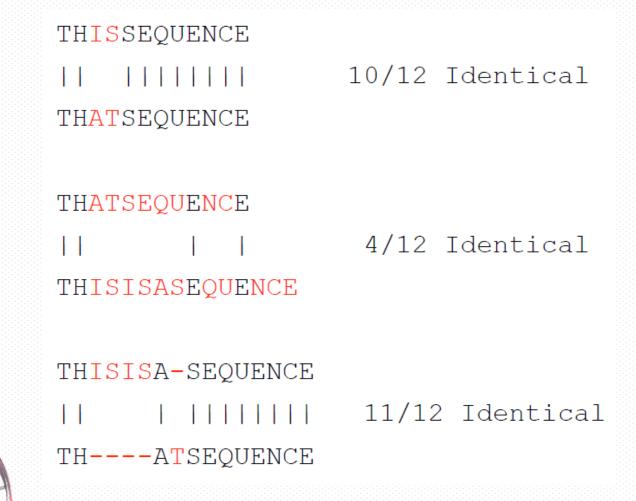


## Convergent evolution



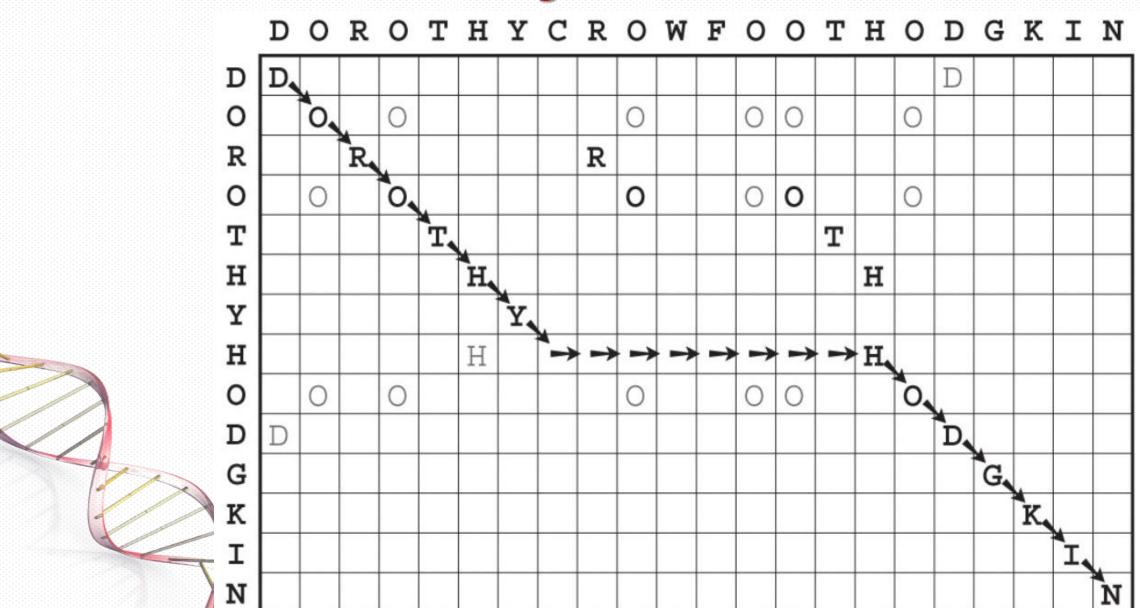


#### What is an alignment

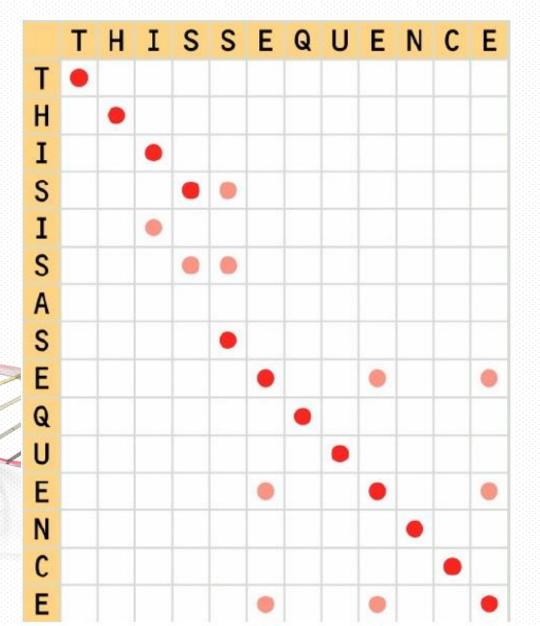


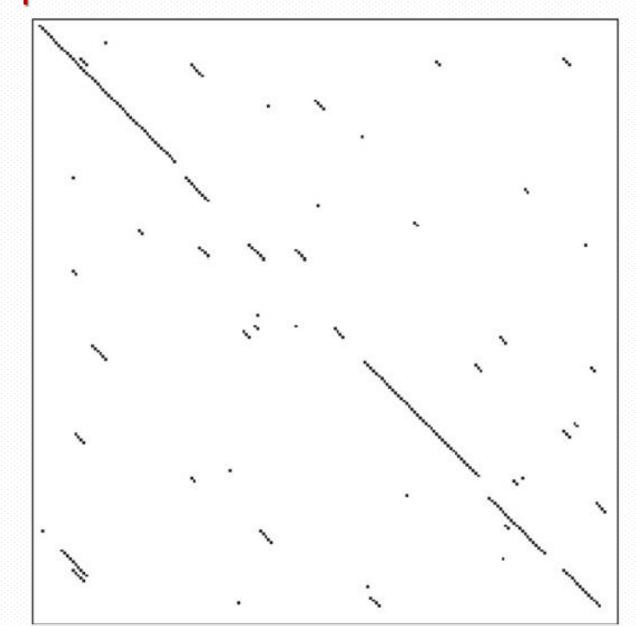
- Match
- Mismatch
- Insertions/deletions (Indels)

### An alignment matrix



### **Dotplots**





### Types of alignment

Global alignment

FTFTALILLAVAV F--TAL-LLA-AV

Local alignment



#### Inserting gaps

	Triserting gaps
(A)	
Bovine PI-3Kinase p110a	LNWENPDIMSELLFQNNEIIFKNGDDLRQDMLTLQIIRIMENIWQNQGLDLRMLPYGCLSIGDCVGLIEVVRNSHTIMQIQCKGGLKGAL
cAMP-dependent protein kinase	WENPAQNTAHLDQFERIKTLGTGSFGRV <mark>ML</mark> VKHMETGNHYAMKILDKQKVVKLKQIEHTLNEKRILQA <mark>V</mark> NFPFLVKLEFSFKDNSNLY
Bovine PI-3Kinase p110a	QFNSHTLHQWLKDKNKGEIYDAAIDLFTRSCAGYCVATFILGIGDRHNSNIMVKDDGQLFHIDFGHFLDHK <mark>K</mark> KKF <mark>G</mark> YKRERVPFVLTQDF
cAMP-dependent protein kinase	MVMEYVPGGEMFSHLRRIGRFSEPHARFYAAQIVLTFEYLHSLDLIYRDLKPENLLIDQQGYIQVTDFGFAKRVKGRTWXLCGTPEYLAP
Bovine PI-3Kinase p110a	LIVI <mark>SKG</mark> AQECTKTREFERFQEMCYKAYLAIRQHANLF <mark>I</mark> NLFSMMLGSGMPELQ <mark>SFD</mark> DIAYI <mark>R</mark> KTLALDKTEQEALEYFMKQMNDA <mark>H</mark> HGG
cAMP-dependent protein kinase	EIIL <mark>skg</mark> ynkavdwwalgvliyemaagyppffadqpiq <mark>i</mark> yekivsgkvrf <mark>p</mark> shf <mark>ssd</mark> lkdll <mark>rnll</mark> qv <mark>dlt</mark> krfgnlkngvndikn <mark>h</mark> kwf
Bovine PI-3Kinase p110a	WTTKMDWIFHTIKQHALN
cAMP-dependent protein kinase	ATT DWIAI YQRKVEAPFIPKFKGPGDTSNFDDYEEEEIRVXINEKCGKEFSEF
40.4	
(B)	

camp-dependent protein kinase GVNDIKNHKWFATTDWIAIYQRKVEAPFIPKFKGPGDTSNFDDYEEEEIRVXINEKCGKEFSEF

### What is an optimal alignment?

```
THISSEQUENCE
                 10/12 Identical
THATSEQUENCE
THATSEQUENCE
             4/12 Identical
THISISASEQUENCE
THISISA-SEQUENCE
        ||||||| 11/12 Identical
TH----ATSEQUENCE
```

#### Different scoring

```
THISSEQUENCE
                                              T-T=5
5 8-1 1 4 5 6 0 5 6 9 5
                            Score = 52
                                              H-H= 8
                                              S - S = 4
THATSEQUENCE
                                              E-E=5
                                              Q-Q=6
                                              U - U = 0
THATSEQUENCE
                                              N - N = 6
5 8-1-1-2 0-1 0 5 0 0 5
                            Score = 18
                                              C - C = 9
 HISISASEQU
                                              Gap=0
                    ENCE
```

THISISA-SEQUENCE

5 8 0 0 0 4 0 4 5 6 0 5 6 9 5 Score = 56

TH----ATSEQUENCE

#### With Gap cost

Score = 52

Score = 18

THISSEQUENCE

5 8-1 1 4 5 6 0 5 6 9 5

THATSEQUENCE

THATSEQUENCE

5 8-1-1-2 0-1 0 5 0 0 5

HISISASEQUENCE

HISISA-SEQUENCE

5 8-1-1-1-1 4-1 4 5 6 0 5 6 9 5

- - A T S E UΕ  $N \subset E$  T-T=5H-H= 8

S - S = 4E-E=5

Q-Q=6U - U = 0

N - N = 6

C - C = 9

Gap=-1

Score = 51

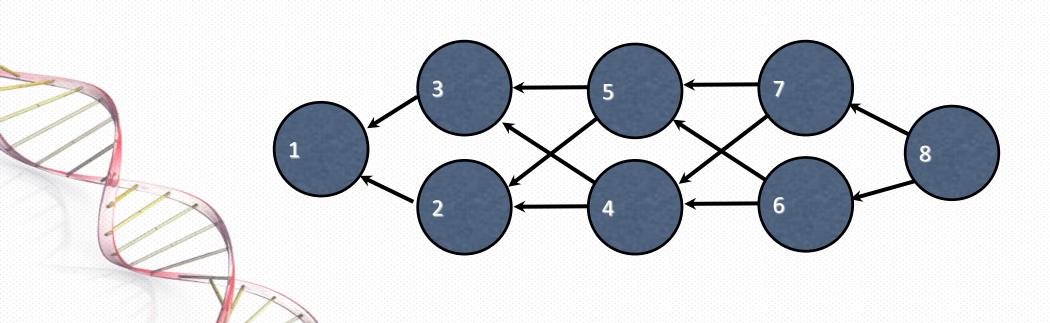
#### **Dynamic Programming**

- Algorithmic technique for optimization problems that have two properties:
  - Optimal substructure: Optimal solution can be computed from optimal solutions to sub-problems

 Overlapping sub-problems: Sub-problems overlap such that the total number of distinct sub-problems to be solved is relatively small

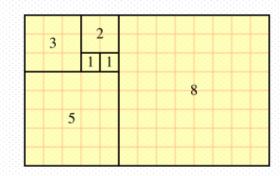
#### **Dynamic Programming**

- Break problem into overlapping subproblems
- use *memoization*: remember solutions to subproblems that we have already seen



#### Fibonacci example

- 1,1,2,3,5,8,13,21,...
- fib(n) = fib(n-2) + fib(n-1)



- Could implement as a simple recursive function
- However, complexity of simple recursive function is exponential in *n*

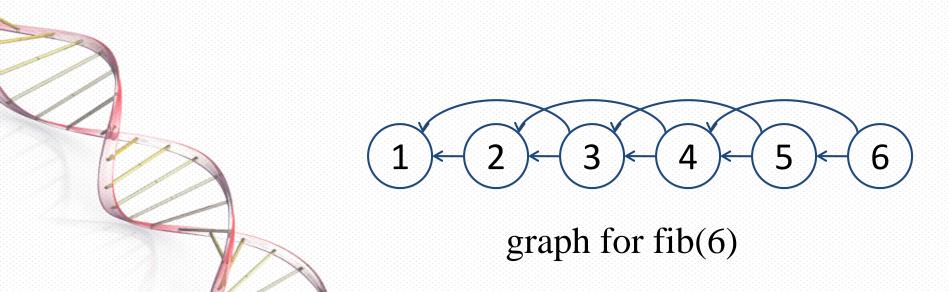
## Fibonacci dynamic programming

- Two approaches
  - 1. Memoization: Store results from previous calls of function in a table (top down approach)
  - 2. Solve subproblems from smallest to largest, storing results in table (bottom up approach)
- Both require evaluating all (n-1) subproblems only once: O(n)



#### **Dynamic Programming Graphs**

- Dynamic programming algorithms can be represented by a directed acyclic graph
  - Each subproblem is a vertex
  - Direct dependencies between subproblems are edges

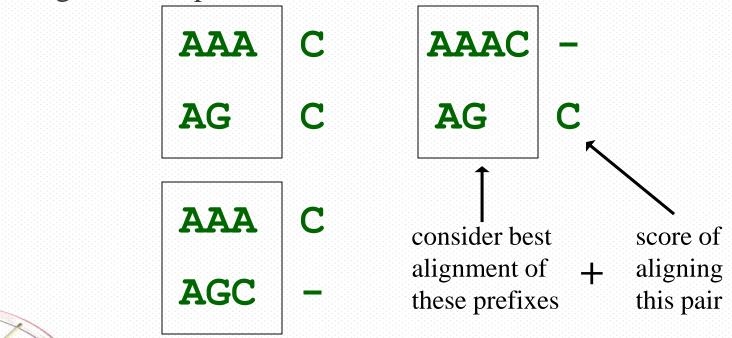


#### Pairwise Alignment Via Dynamic Programming

- first algorithm by Needleman & Wunsch, Journal of Molecular Biology, 1970
- dynamic programming algorithm:
  determine best alignment of two sequences
  by determining best alignment of all
  prefixes of the sequences

#### Dynamic Programming Idea

- consider last step in computing alignment of AAAC with AGC
- three possible options; in each we'll choose a different pairing for end of alignment, and add this to the best alignment of previous characters



## DP Algorithm for Global Alignment with Linear Gap Penalty

• Subproblem: F(i,j) = score of best alignment of the length i prefix of x and the length j prefix of y.

$$F(i,j) = \max_{i} F(i-1,j-1) + S(x_i, y_j)$$

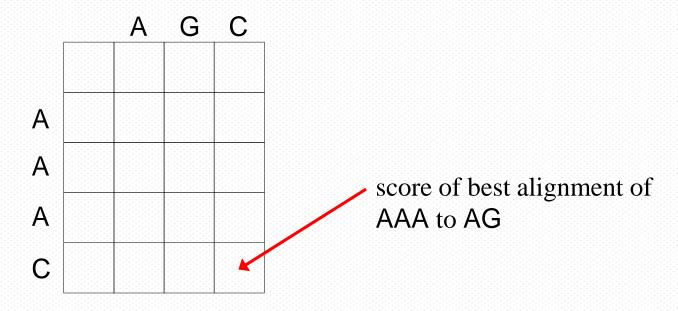
$$F(i,j) = \max_{i} F(i-1,j) + S$$

$$F(i,j-1) + S$$

#### **Dynamic Programming Implementation**

- given an *n*-character sequence *x*, and an *m*-character sequence *y*
- construct an  $(n+1) \times (m+1)$  matrix F
- F(i, j) = score of the best alignment of

x[1...i] with y[1...j]

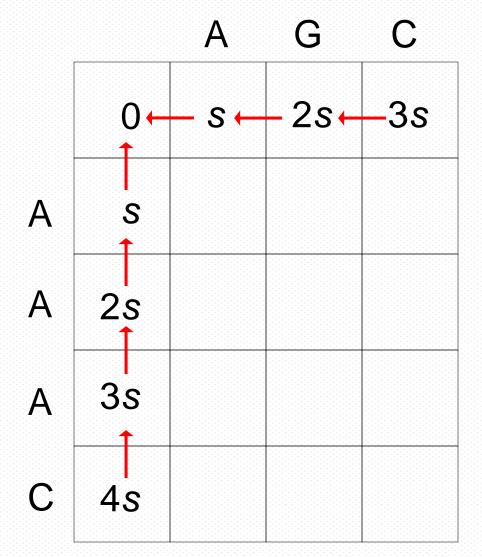


#### Global Alignment Example

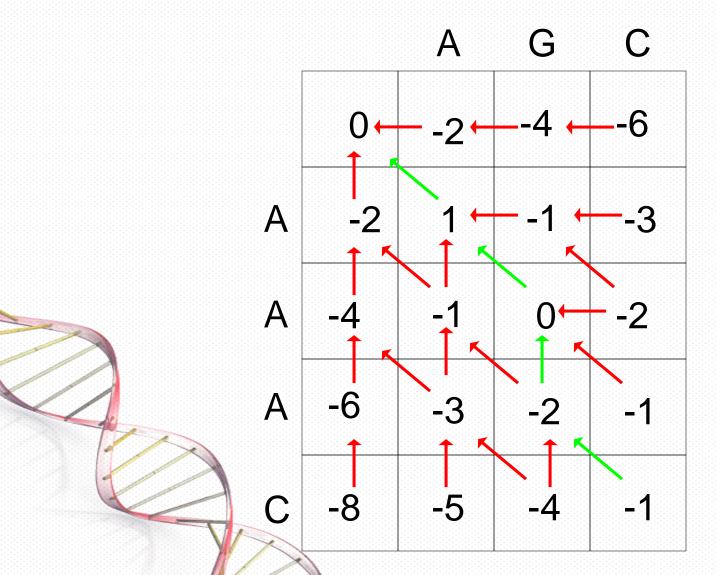
• suppose we choose the following scoring scheme:

$$S(x_i, y_i) =$$
+1 when  $x_i = y_i$ 
-1 when  $x_i^{-1} y_i$ 
 $s$  (penalty for aligning with a space) = -2

# Initializing Matrix: Global Alignment with Linear Gap Penalty



#### Global Alignment Example



#### one optimal alignment

x: A A A C v: A G - C

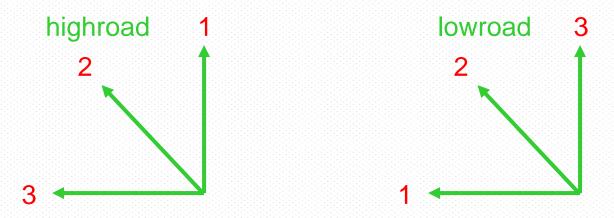
but there are three optimal alignments here (can you find them?)

## DP Algorithm Sketch: Global Alignment

- initialize first row and column of matrix
- fill in rest of matrix from top to bottom, left to right
- for each F(i, j), save pointer(s) to cell(s) that resulted in best score
- F(m, n) holds the optimal alignment score; trace pointers back from F(m, n) to F(0, 0) to recover alignment

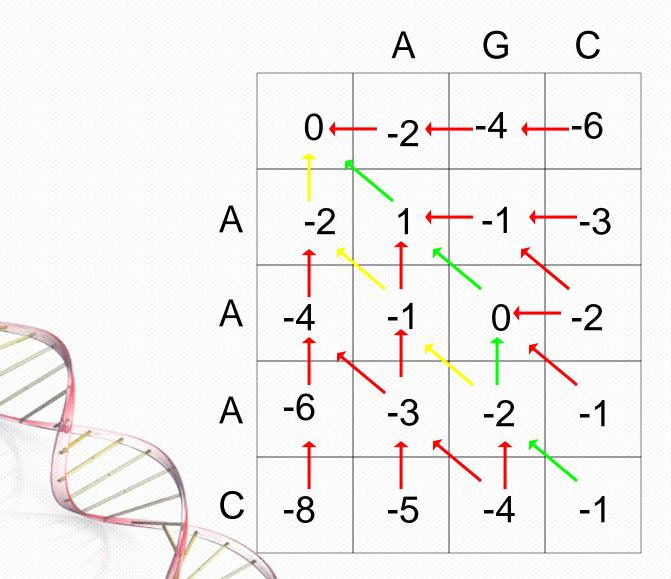
#### **Equally Optimal Alignments**

- many optimal alignments may exist for a given pair of sequences
- can use preference ordering over paths when doing traceback



 High road and low road alignments show the two most different optimal alignments

#### High road & Low road Alignments



#### High road alignment

x: A A A C v: A G - C

#### Low road alignment

x: A A A C y: - A G C

- Initialization step: Create Matrix with M + 1 columns
- and N + 1 rows. First row and column filled with 0.

			G	A	Α	T	T	С	A	G	T	T	A
		0	0	0	0	0	0	0	0	0	0	0	0
(	3[	0											
(	3[	0											
	A	0											
	г[	0											
	디	0											
	3[	0											
	A[	0											

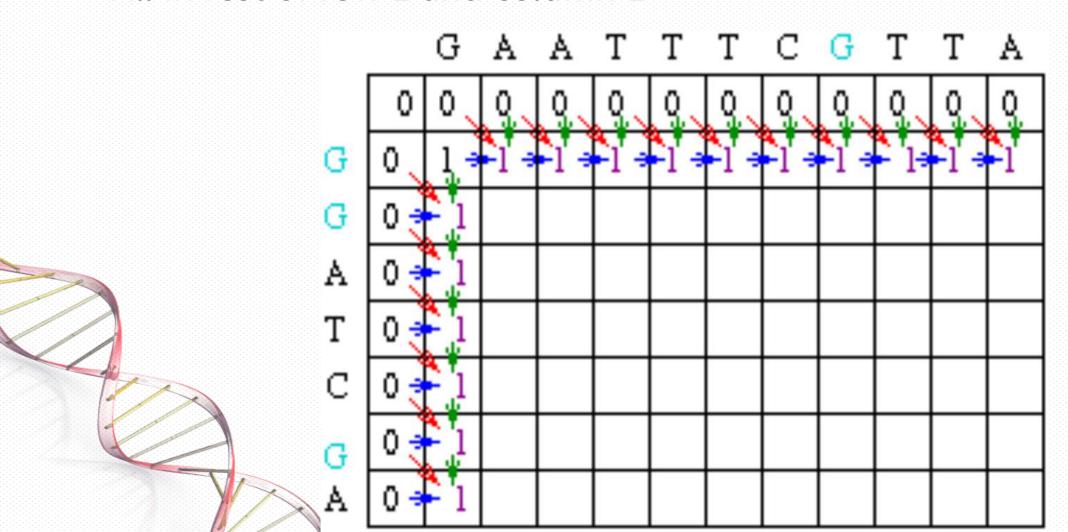
Matrix fill step: Each position Mi,j is defined to be the

Mismatch = 0

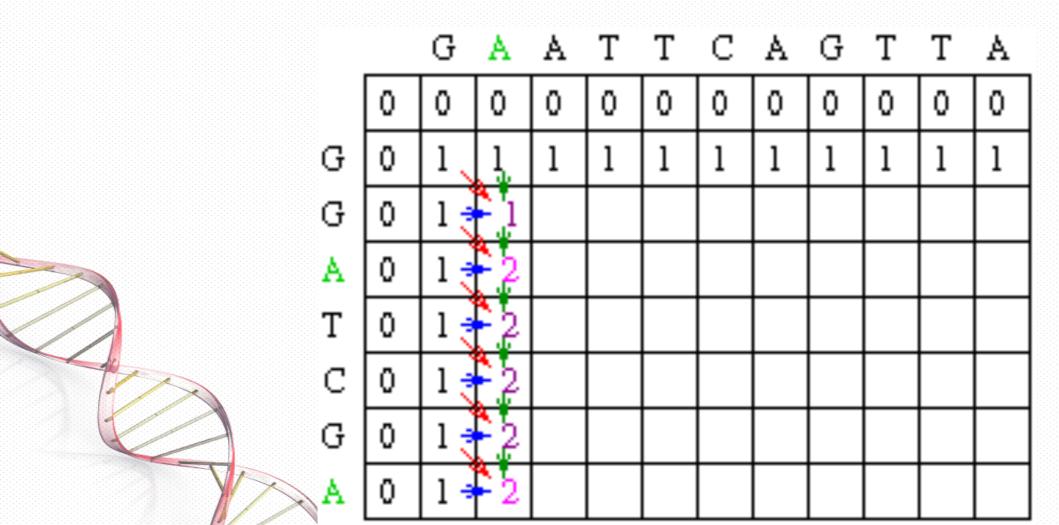
Space = 0

MAXIMUM score at position i,j  $M_{i,i} = MAXIMUM$  [  $M_{i-1, j-1} + s_{i,j}$  (match/mismatch in the diagonal  $M_{i, i-1}$  + w (gap in sequence #1)  $M_{i-1,i}$  + w (gap in sequence #2) Match = +1

Fill in rest of row 1 and column 1



• Fill in column 2

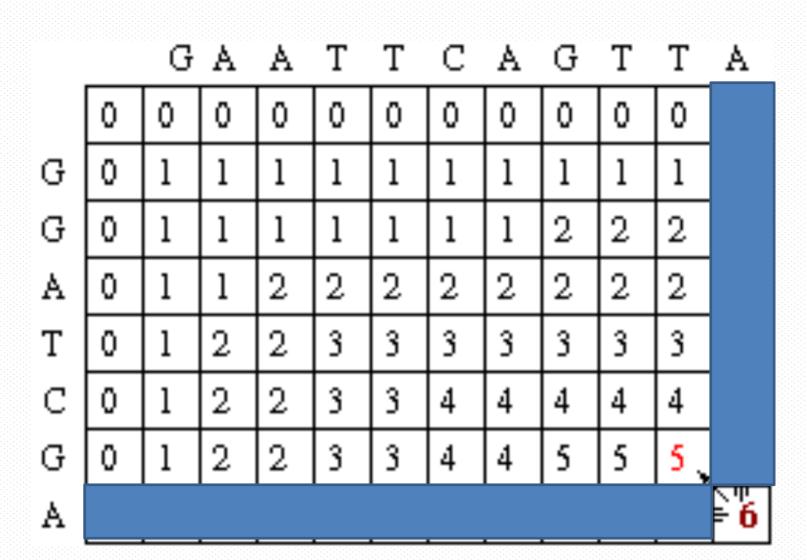


• Fill in column 3

		G	A	A	T	T	Ċ,	A	Ğ	T	T	À	
	0	0	0	0	0	0	0	0	0	0	0	0	
	0	1	1	1	1	1	1	1	1	1	1	1	
	0	1	1	1	1	1	1	1	2	2	2	2	
	0	1	2	2	2	2	2	2	2	2	2	3	
	0	1	2	2	3	3	3	3	3	3	3	3	
,	0	1	2	2	3	3	3	4	4	4	4	4	
	0	1	2	2	3	3	3	4	4	5	5	5	
	0	1	2	3	3	3	3	4	5	5	5	6	

- Traceback step:
  - Position at current cell and look at direct predecessors





- Traceback step:
  - Position at current cell and look at direct predecessors

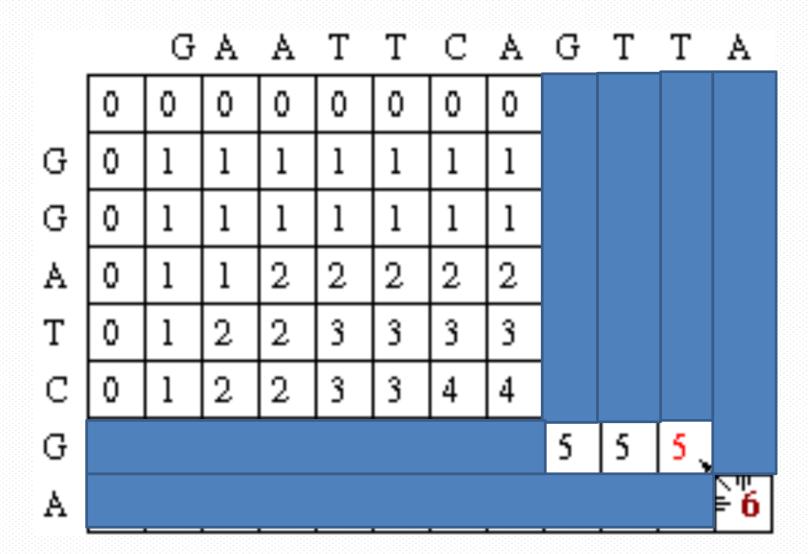
		G	Α	A	T	T	С	A	G	T	Т	A
	0	0	0	0	0	0	0	0	0	0		
G	0	1	1	1	1	1	1	1	1	1		
G	0	1	1	1	1	1	1	1	2	2		
A	0	1	1	2	2	2	2	2	2	2		
T	0	1	2	2	3	3	3	3	3	3		
С	0	1	2	2	3	3	4	4	4	4		
G	0	1	2	2	3	3	4	4	5	5	5	
Α												- 6

- Traceback step:
  - Position at current cell and look at direct predecessors

Seq # 1 T T A |
Seq # 2 - - A
Seq # 1 T A |
Seq # 2 - A

		G	Α	Α	T	T	С	Α	G	T	T	A
	0	0	0	0	0	0	0	0	0			
G	0	1	1	1	1	1	1	1	1			
G	0	1	1	1	1	1	1	1	2			
A	0	1	1	2	2	2	2	2	2			
Т	0	1	2	2	3	3	3	3	3			
С	0	1	2	2	3	3	4	4	4			
G	0	1	2	2	3	3	4	4	5	5	5	
Α												÷ 6

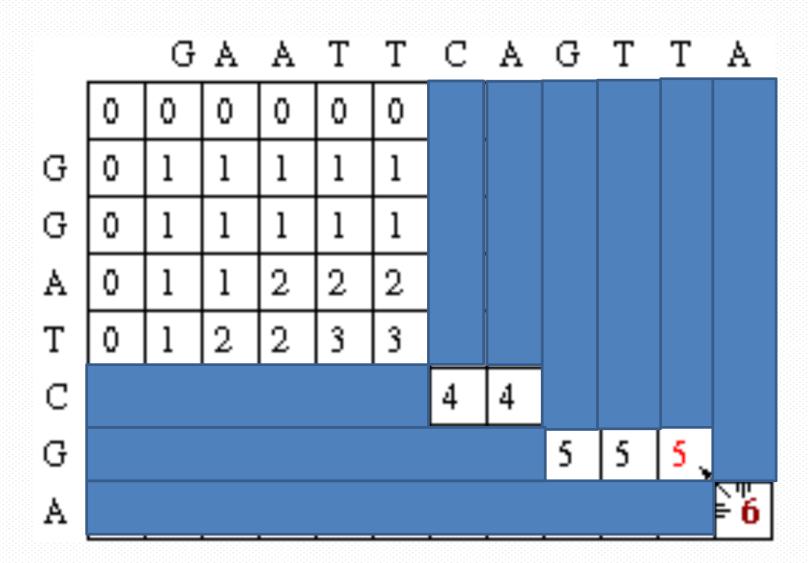
- Traceback step:
  - Position at current cell and look at direct predecessors



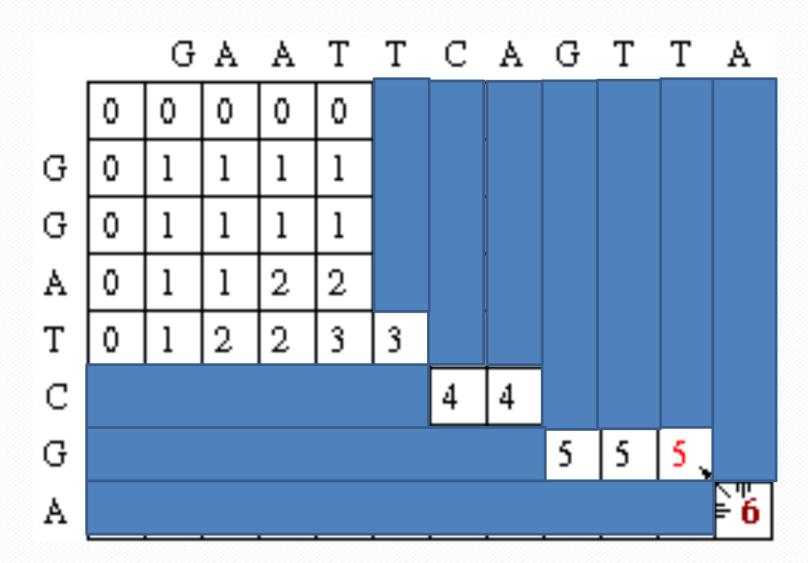
- Traceback step:
  - Position at current cell and look at direct predecessors

		G	A	Α	T	T	С	A	G	Т	T	A
	0	0	0	0	0	0	0					
G	0	1	1	1	1	1	1					
G	0	1	1	1	1	1	1					
Α	0	1	1	2	2	2	2					
Т	0	1	2	2	3	3	3					
С	0	1	2	2	3	3	4	4				
G									5	5	5	
Α												÷ 6

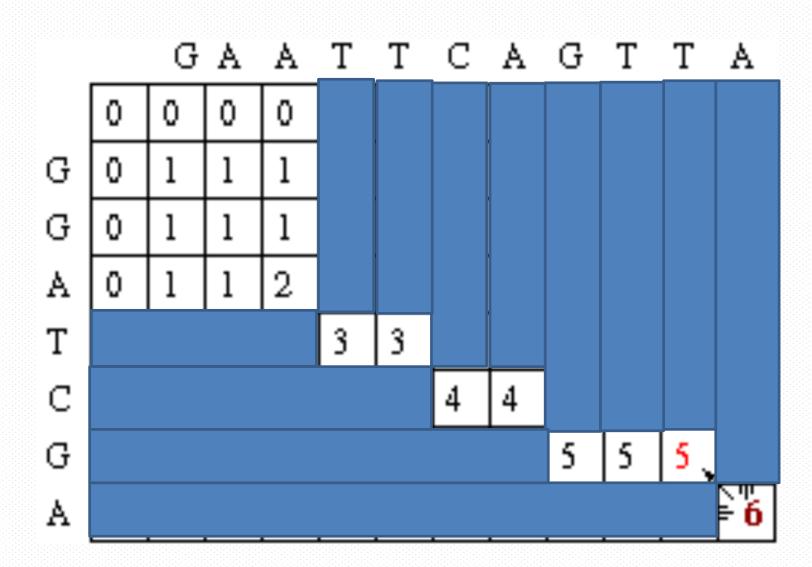
- Traceback step:
  - Position at current cell and look at direct predecessors



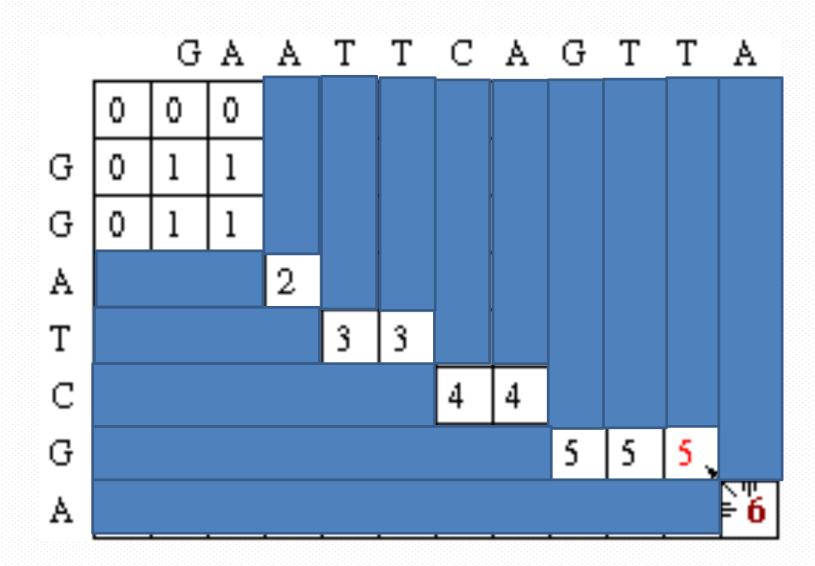
- Traceback step:
  - Position at current cell and look at direct predecessors



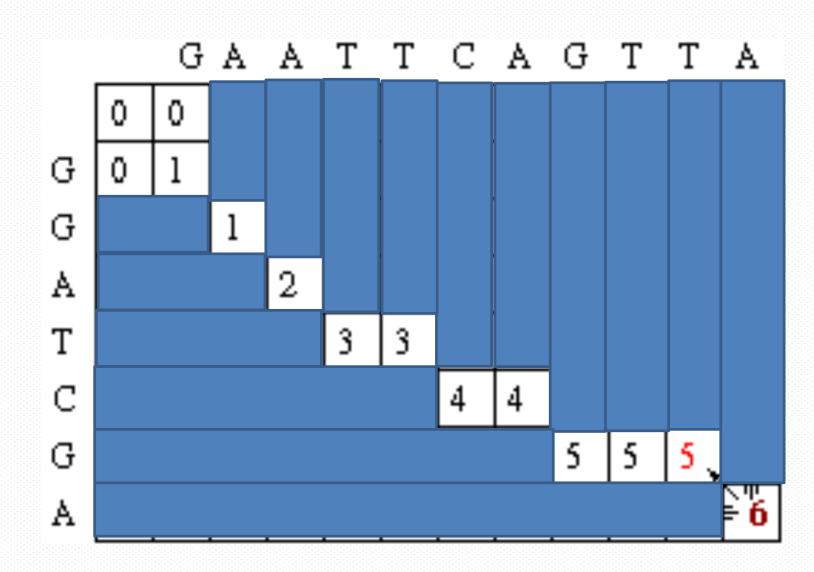
- Traceback step:
  - Position at current cell and look at direct predecessors



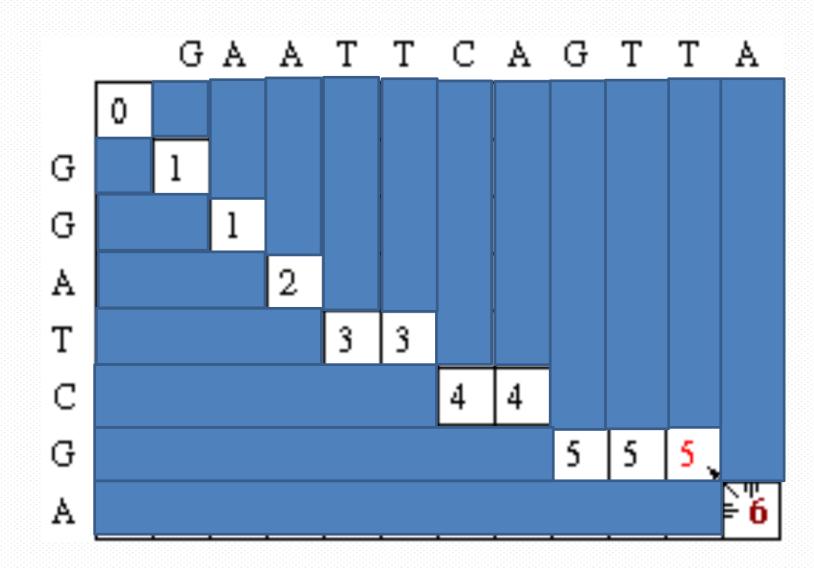
- Traceback step:
  - Position at current cell and look at direct predecessors



- Traceback step:
  - Position at current cell and look at direct predecessors



- Traceback step:
  - Position at current cell and look at direct predecessors



#### Pseudocode

```
for i=0 to length(A)
  F(i,0) \leftarrow d*i
for j=0 to length(B)
  F(0,j) \leftarrow d*j
for i=1 to length(A)
  for j=1 to length(B)
     Match \leftarrow F(i-1,j-1) + S(A<sub>i</sub>, B<sub>j</sub>)
     Delete \leftarrow F(i-1, j) + d
     Insert \leftarrow F(i, j-1) + d
     F(i,j) \leftarrow max(Match, Insert,
Delete)
```

#### Computational Complexity

- initialization: O(m), O(n) where sequence lengths are m, n
- filling in rest of matrix: O(mn)
- traceback: O(m + n)
- hence, if sequences have nearly same length, the computational complexity is



 $O(n^2)$