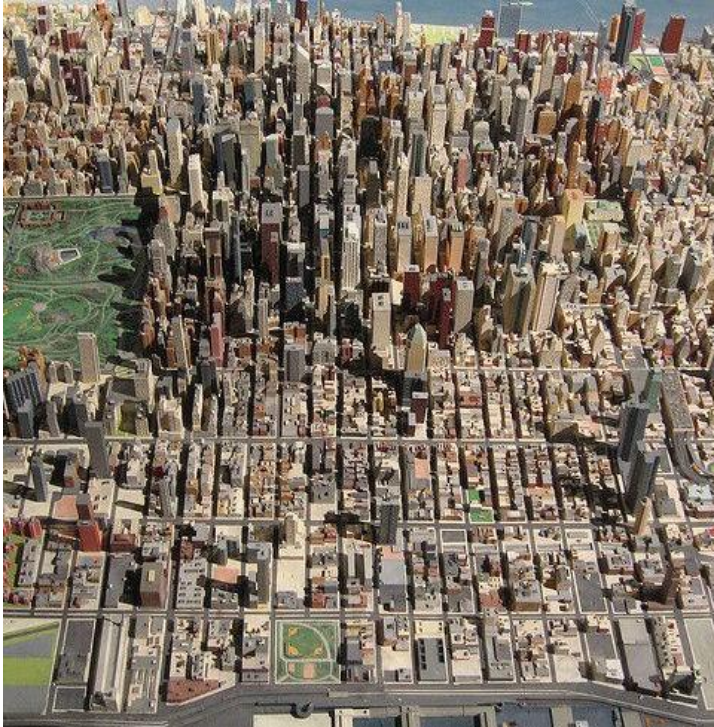
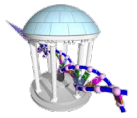


# Comp 555 - BioAlgorithms - Spring 2021



- *HOW OUR MANHATTAN TOUR RELATES TO SEQUENCES*
- *MIDTERM GRADES NEXT TUESDAY*
- *LOOK FOR PROBLEM SET #4 TOMORROW*

Sequence Alignment

# Comparing Sequences



- What makes two sequences similar?
- What is the best measure of similarity?
- Consider the two DNA sequences  $v$  and  $w$  :

$v$  : TAGACAAT  
 $w$  : AGAGACAT  
1111100 = 6

- The Hamming distance,  $d_H(v, w) = 6$ , is large but the sequences seem to have more similarity
- What if we allowed for insertions and deletions?

# Allowing Insertions and Deletions



- By shifting each sequence over one position:

	Shifts and gaps:	Another one:
v : _TAGACAAT	v : _TAGACAAT	v : T_AGACAAT
w : AGAGACAT_	w : AGAGAC_AT	w : AGAGACA_T
11000011 = 4	110000100 = 3	110000010 = 3

- The edit distance:  $d_H(v, w) = 3$ .
- Hamming distance neglects insertions and deletions

# Edit Distance

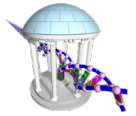


- Vladimir Levenshtein introduced the notion of an “edit distance” between two strings as the minimum number of elementary operations (insertions, deletions, and substitutions) to transform one string into the other in 1965.
- $d_L(v,w)$  = Minimum number of elementary operations to transform  $v \rightarrow w$
- Computing Hamming distance is a trivial task
- Computing edit distance is less trivial



**Vladimir Levenshtein**

1935 - 2017



# Edit Distance: Example

TGCATAT → ATCCGAT in 5 steps

TGCATAT → (DELETE last T)

TGCATA → (DELETE last A)

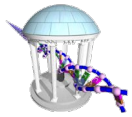
TGCAT → (INSERT A at front)

ATGCAT → (SUBSTITUTE C for G)

ATCCAT → (INSERT G before last A)

ATCCGAT (Done)

What is the edit distance? 5? (Recall it has to be the *minimum*)



# Edit Distance: Example (2<sup>nd</sup> Try)

TGCATAT → ATCCGAT in 4 steps

TGCATAT → (INSERT A at front)  
ATGCATAT → (DELETE 2nd T)  
ATGCAAT → (SUBSTITUTE G for 2nd A)  
ATGCGAT → (SUBSTITUTE C for 1st G)  
ATCCGAT (Done)

But is 4 the minimum edit distance? Is 3 possible?

- Edit sequences are invertible, i.e given  $v \rightarrow w$ , one can generate  $w \rightarrow v$ , without recomputing
- A little jargon: Since the effect of insertion in one string can be accomplished via a deletion in the other string these two operations are correlated. Often algorithms will consider them together as a single operation called INDEL

# An aside: Longest Common Subsequence



- A special case of alignment where only *matches, insertions, and deletions* are allowed
- A variant of Edit distance, sometimes called LCS distance, where only indels are allowed
- A subsequence need not be contiguous, but the symbol order must be preserved  
Ex. If  $v = \text{ATTGCTA}$  then  $\text{AGCA}$  and  $\text{TTTA}$  are subsequences of  $v$ , but  $\text{TGTT}$  and  $\text{ACGA}$  are not
- All substrings of  $v$  are subsequences, but not vice versa
- Edit distance,  $d_{LCS}$ , is related to the length of the LCS,  $s$ , by the following relationship:

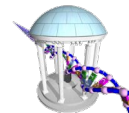
$$d_{LCS}(u,w) = \text{len}(v) + \text{len}(w) - 2s(u,w)$$

Example:

AN**UNCL**E**I**KE  
UN**CB**E**A**T**D**U**K**E

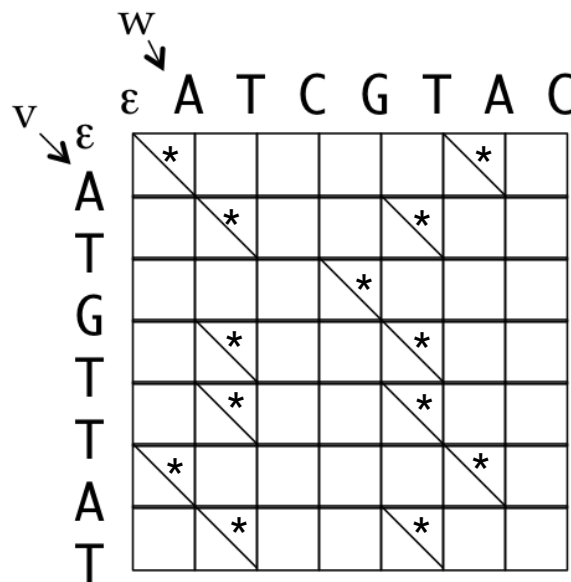
	len	LCS
anUNC_lE____iKE	10	- 6 = 4
__UNCb_Eatdu_KE	11	- 6 = 5

# LCS as a Dynamic Program



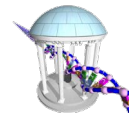
There are similarities between the LCS and MTP

- All possible possible alignments can be represented as a path from the string's beginning (source) to its end (destination)
- Horizontal edges add gaps in  $v$
- Vertical edges add gaps in  $w$
- Diagonal edges are a match
- Notice that we've only included valid diagonal edges for "matches" in our graph
- *An maximum LCS is a path from  $(\epsilon, \epsilon)$  to the end of both strings that matches the most bases (a.k.a. a Manhattan tour)*





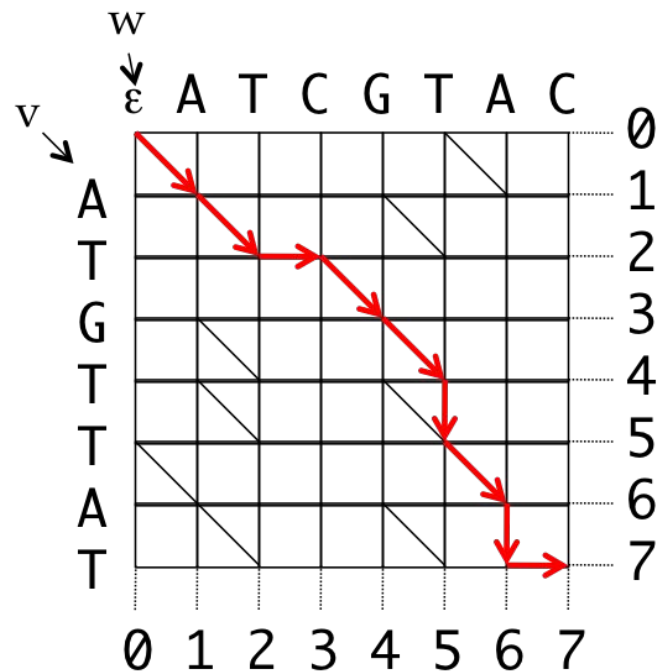
# The "Space" of All Alignments



- Introduce coordinates for the grid
- All valid paths from the source to the destination represent some alignment

$\emptyset$	1	2	2	3	4	5	6	7	7
<b>v</b>	A	T	_	G	T	T	A	T	_
<b>w</b>	A	T	C	G	T	_	A	_	C
$\emptyset$	1	2	3	4	5	5	6	6	7

- Path:  
(0,0), (1,1), (2,2), (2,3), (3,4), (4,5), (5,5), (6,6), (7,6), (7,7)



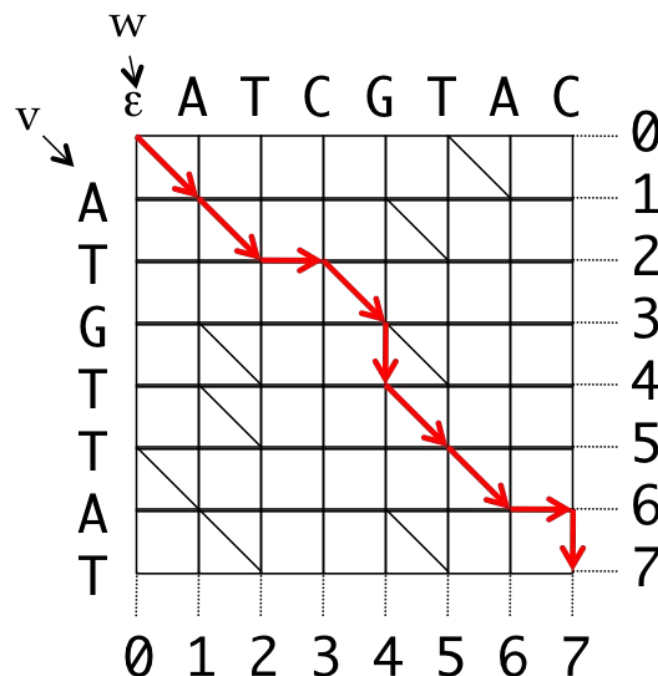
# Alternate Alignment

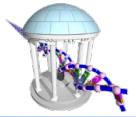


- Introduce coordinates for the grid
- All valid paths from the source to the destination represent some alignment

$\emptyset$	1	2	2	3	4	5	6	6	7
<b>v</b>	A	T	_	G	T	T	A	_	T
<b>w</b>	A	T	C	G	_	T	A	C	_
$\emptyset$	1	2	3	4	4	5	6	7	7

- Path:  
(0,0), (1,1), (2,2), (2,3), (3,4), (4,4), (5,5), (6,6), (6,7), (7,7)



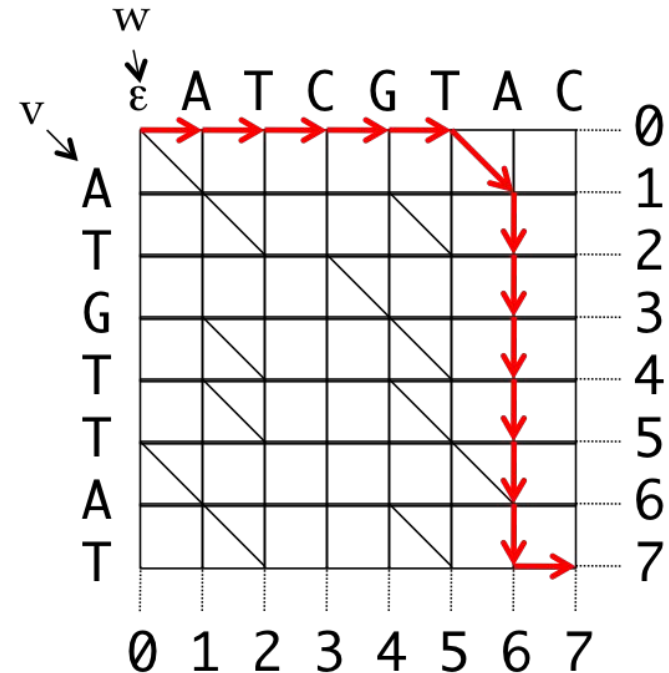


# Even Bad Alignments

- Introduce coordinates for the grid
- All valid paths from the source to the destination represent some alignment

	0	0	0	0	0	0	1	2	3	4	5	6	7	7
<b>v</b>	-	-	-	-	-	A	T	G	T	T	A	T	-	
<b>w</b>	A	T	C	G	T	A	-	-	-	-	-	-	C	
	0	1	2	3	4	5	6	6	6	6	6	6	6	7

- Path:  
(0,0), (0,1), (0,2), (0,3), (0,4), (0,5), (1,6),  
(2,6), (3,6), (4,6), (5,6), (6,6), (7,6), (7,7)

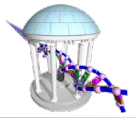


# What makes a good alignment?



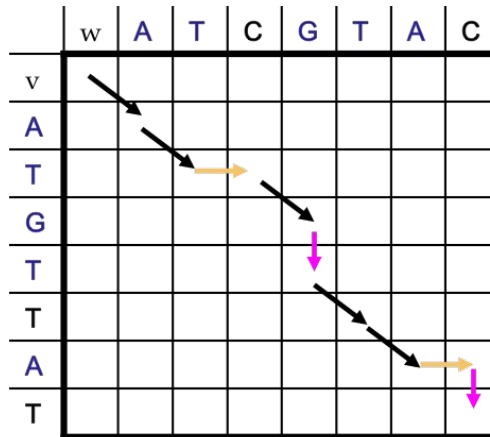
- Using as many diagonal segments, when they correspond to matches, as possible. Why?
- The end of a good alignment from  $(j\dots k)$  begins with a good alignment from  $(i..j)$
- Same as Manhattan Tourist problem, where the **sites** are only on the diagonal streets!
- Set diagonal street weights = 1, and horizontal and vertical weights = 0



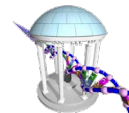


# LCS: Dynamic Program

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + 1 & \text{if } v_i = w_j \quad \searrow \\ s_{i-1,j} & \downarrow \\ s_{i,j-1} & \rightarrow \end{cases}$$



# Step 1



Initialize 1st row and 1st column to all zeroes.

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0							
T	0							
G	0							
T	0							
T	0							
A	0							
T	0							

- Note intersections/vertices are cells/entries of this matrix

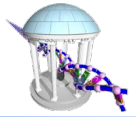
# Step 2



Evaluate recursion for next row and/or next column

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1						
G	0	1						
T	0	1						
T	0	1						
A	0	1						
T	0	1						

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



# Step 3

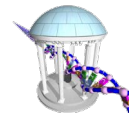
Continue recursion for next row and/or next column

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2					
T	0	1	2					
T	0	1	2					
A	0	1	2					
T	0	1	2					

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



# Step 4

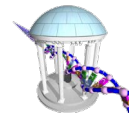


Then one more row and/or column

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2				
T	0	1	2	2				
A	0	1	2	2				
T	0	1	2	2				

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

# Step 5



And so on...

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2	3	4	4	4
T	0	1	2	2	3			
A	0	1	2	2	3			
T	0	1	2	2	3			

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + 1 & \text{if } v_i = w_j \quad \swarrow \\ s_{i-1,j} & \downarrow \\ s_{i,j-1} & \rightarrow \end{cases}$$

# Step 6



And so on...

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2	3	4	4	4
T	0	1	2	2	3	4	4	4
A	0	1	2	2	3	4		
T	0	1	2	2	3	4		

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

# Step 7

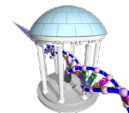


## Getting closer

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2	3	4	4	4
T	0	1	2	2	3	4	4	4
A	0	1	2	2	3	4	5	5
T	0	1	2	2	3	4	5	

$$s_{i,j} = \max \begin{cases} s_{i-1,j-1} + 1 & \text{if } v_i = w_j \quad \swarrow \\ s_{i-1,j} & \downarrow \\ s_{i,j-1} & \rightarrow \end{cases}$$

# Step 8

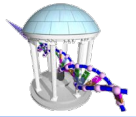


Until we reach the last row and column

	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2	3	4	4	4
T	0	1	2	2	3	4	4	4
A	0	1	2	2	3	4	5	5
T	0	1	2	2	3	4	5	5

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

# Finally



We reach the end, which corresponds to an LCS of length 5

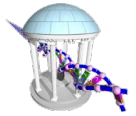
	w	A	T	C	G	T	A	C
v	0	0	0	0	0	0	0	0
A	0	1	1	1	1	1	1	1
T	0	1	2	2	2	2	2	2
G	0	1	2	2	3	3	3	3
T	0	1	2	2	3	4	4	4
T	0	1	2	2	3	4	4	4
A	0	1	2	2	3	4	5	5
T	0	1	2	2	3	4	5	5

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

**w = ATCGT\_A\_C**  
**v = AT\_GTTAT\_**  
**len(LCS) = 5**

Our answer includes both an optimal score, and a path back to both the LCS and an alignment

# LCS Code

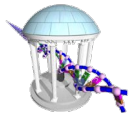


Let's see how well the code matches the approach we sketched out...

```
1 from numpy import *
2
3 def findLCS(v, w):
4     score = zeros((len(v)+1,len(w)+1), dtype="int32")
5     backt = zeros((len(v)+1,len(w)+1), dtype="int32")
6     for i in range(1,len(v)+1):
7         for j in range(1,len(w)+1):
8             # find best score at each vertex
9             if (v[i-1] == w[j-1]): # test for a match ("diagonal street")
10                score[i,j], backt[i,j] = max((score[i-1,j-1]+1,3), (score[i-1,j],1), (score[i,j-1],2))
11            else:
12                score[i,j], backt[i,j] = max((score[i-1,j],1), (score[i,j-1],2))
13        return score, backt
14
15 v = "ATGTTAT"
16 w = "ATCGTAC"
17 s, b = findLCS(v,w)
18 for i in range(len(s)):
19     print("%10s %-20s" % (" " if i else 'score =', str(s[i]), " " if i else 'backtrack =', str(b[i])))
```

score =	[0 0 0 0 0 0 0 0]	backtrack =	[0 0 0 0 0 0 0 0]
	[0 1 1 1 1 1 1 1]		[0 3 2 2 2 2 3 2]
	[0 1 2 2 2 2 2 2]		[0 1 3 2 2 3 2 2]
	[0 1 2 2 3 3 3 3]		[0 1 1 2 3 2 2 2]
	[0 1 2 2 3 4 4 4]		[0 1 3 2 1 3 2 2]
	[0 1 2 2 3 4 4 4]		[0 1 3 2 1 3 2 2]
	[0 1 2 2 3 4 5 5]		[0 3 1 2 1 1 3 2]
	[0 1 2 2 3 4 5 5]		[0 1 3 2 1 3 1 2]

- The same score matrix that we found by hand
- *"backtrack"* keeps track of the "arrow" used



# Backtracking

[0	0	0	0	0	0	0	0]
[0	3	2	2	2	2	3	2]
[0	1	3	2	2	3	2	2]
[0	1	1	2	3	2	2	2]
[0	1	3	2	1	3	2	2]
[0	1	3	2	1	3	2	2]
[0	3	1	2	1	1	3	2]
[0	1	3	2	1	3	1	2]



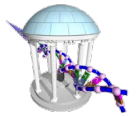
Our score table kept track of the longest common subsequence so far. How do we figure out what the subsequence is?

The **second** "arrow" table kept track of the decisions we made... and we'll use it to backtrack to our answer.

In our example we used arrows {↓, →, ↘}, which were represented in our matrix as {1,2,3} respectively. This numbering is *arbitrary*, except that it does break ties in our implementation (matches > w deletions > w insertions).

Now we need code that finds a path from the end of our strings to the beginning using our arrow matrix





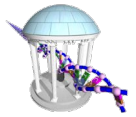
# Code to extract an answer

A simple recursive LCS() routine to return along the path of arrows that led to our best score.

```
In [7]: def LCS(b,v,i,j):  
        if ((i == 0) and (j == 0)):  
            return ''  
        elif (b[i,j] == 3):  
            return LCS(b,v,i-1,j-1) + v[i-1]  
        elif (b[i,j] == 2):  
            return LCS(b,v,i,j-1)  
        else:  
            return LCS(b,v,i-1,j)  
  
        print(LCS(b,v,b.shape[0]-1,b.shape[1]-1))
```

ATGTA

A	[0	0	0	0	0	0	0	0]
T	[0	3	2	2	2	2	3	2]
G	[0	1	3	2	2	3	2	2]
T	[0	1	1	2	3	2	2	2]
T	[0	1	3	2	1	3	2	2]
A	[0	1	3	2	1	3	2	2]
T	[0	3	1	2	1	1	3	2]



# But that's not an alignment

- Technically correct, ATGTA is the LCS

$$w = \text{ATcGT\_A\_c}$$
$$v = \text{AT\_GTtAt\_}$$

- Notice that we only need one of  $v$  or  $w$  since both contain the LCS
- But we would like to get more than just the LCS
- For example, the corresponding alignment.

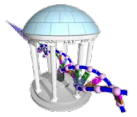
# An alignment of v and w



```
In [10]: def Alignment(b,v,w,i,j):
    if ((i == 0) and (j == 0)):
        return [' ','']
    if (b[i,j] == 3):
        result = Alignment(b,v,w,i-1,j-1)
        result[0] += v[i-1]
        result[1] += w[j-1]
        return result
    if (b[i,j] == 2):
        result = Alignment(b,v,w,i,j-1)
        result[0] += "_"
        result[1] += w[j-1]
        return result
    if (b[i,j] == 1):
        result = Alignment(b,v,w,i-1,j)
        result[0] += v[i-1]
        result[1] += "_"
        return result

align = Alignment(b,v,w,b.shape[0]-1,b.shape[1]-1)
print("v =", align[0])
print("w =", align[1])
```

```
v = AT_GTTAT_
w = ATCG_TA_C
```



# Next Time

- Convert LCS to a general purpose sequence aligner
- Scoring matrices
- Global vs. Local alignments
- Affine gap penalties

