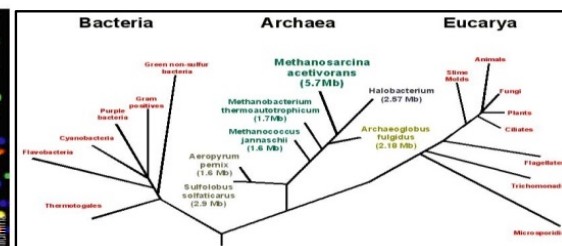
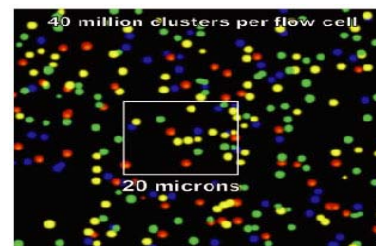




TAACCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
CCTAACCCCTAACCCCTAACCCCTAACCCCTAACCC  
CCCTAACCCCTAACCCCTAACCCCTAACCCCTAAC  
AACCCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
ACCCTAACCCCAACCCCAACCCCAACCCCAAC  
CTACCCTAACCCCTAACCCCTAACCCCTAACCCCTA  
ACCCTAACCCCTAACCCCTAACCCCTAACCCCTAA

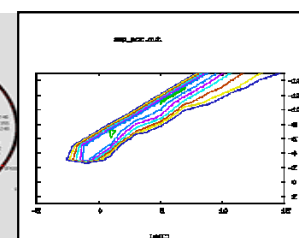
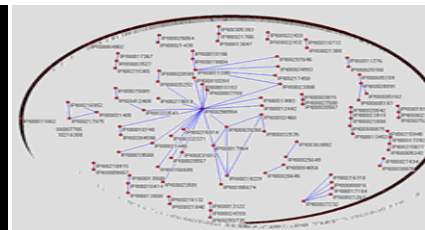
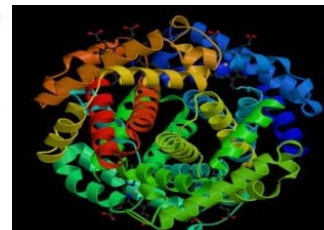
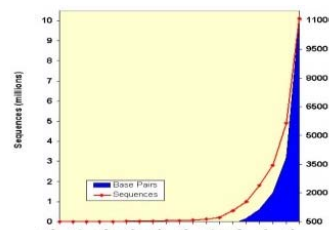
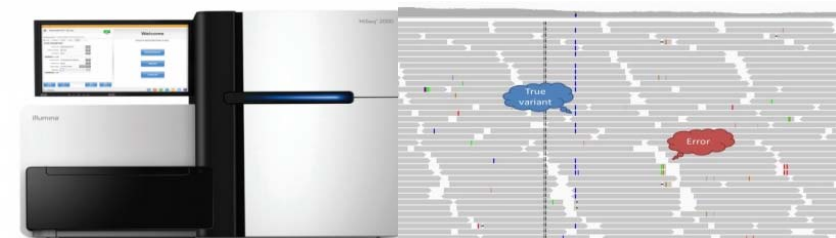


# Sequence Alignment

北京大学生物信息学中心 高歌

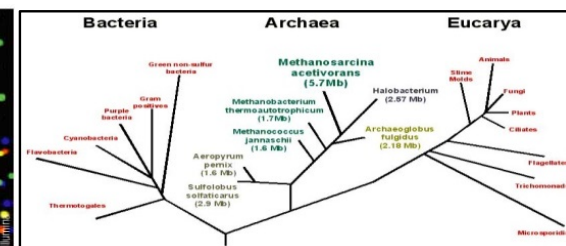
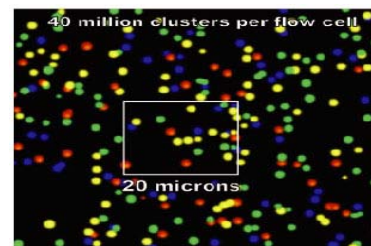
Ge Gao, Ph.D.

Center for Bioinformatics, Peking University





TAACCCTAACCCTAACCCTAACCCTAACCCTA  
CCTAACCCTAACCCTAACCCTAACCCTAACC  
CCCTAACCCTAACCCTAACCCTAACCCTAAC  
AACCCTAACCCTAACCCTAACCCTAACCCTA  
ACCCTAACCCTAACCCTAACCCTAACCCTAAC  
CTACCCTAACCCTAACCCTAACCCTAACCCTA  
ACCCTAACCCTAACCCTAACCCTAACCCTA

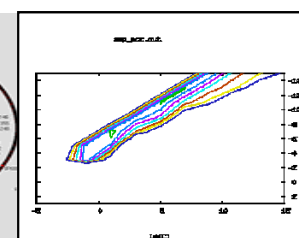
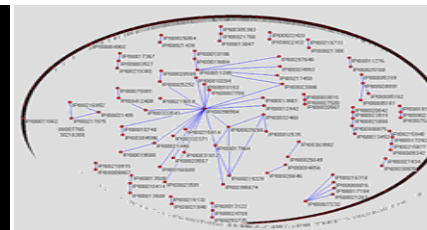
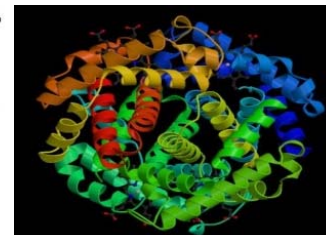
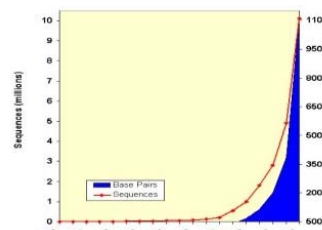
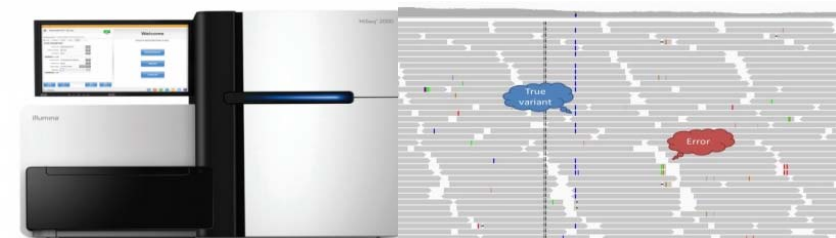


# Unit 3: From Global to Local

北京大学生物信息学中心 高歌

Ge Gao, Ph.D.

Center for Bioinformatics, Peking University



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

		A	A	G
	0	-5		
A		2	-3	
G				-1
C				-6

End-to-end: Global Alignment

# Global Alignment: End-to-end

*J. Mol. Biol.* (1970) 48, 443–453

## Needleman–Wunsch algorithm

### A General Method Applicable to the Search for Similarities in the Amino Acid Sequence of Two Proteins

SAUL B. NEEDLEMAN AND CHRISTIAN D. WUNSCH

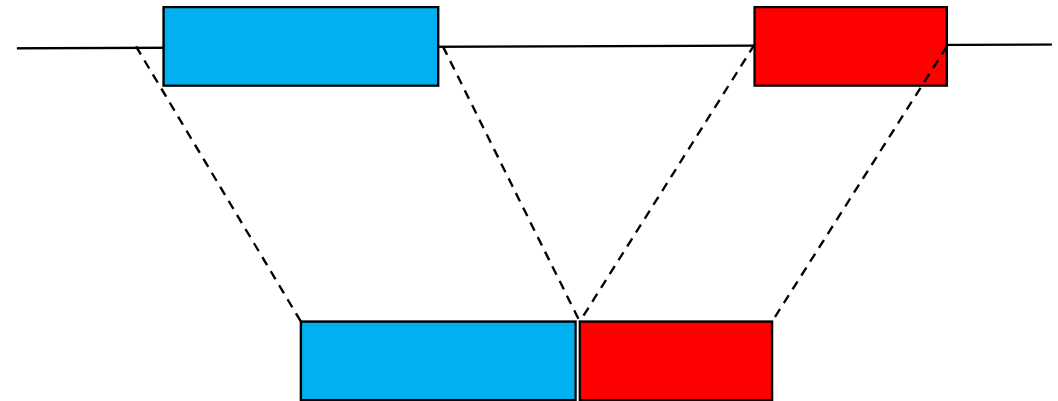
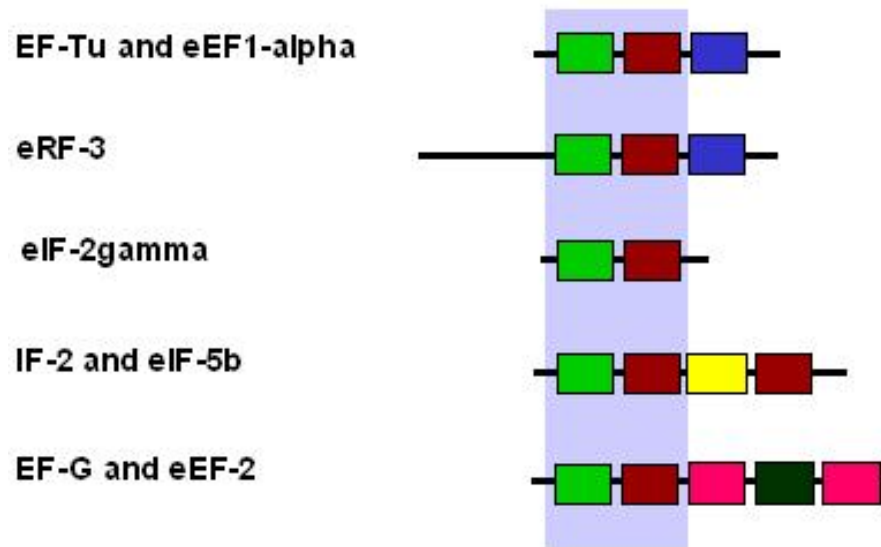
*Department of Biochemistry, Northwestern University, and  
Nuclear Medicine Service, V. A. Research Hospital  
Chicago, Ill. 60611, U.S.A.*

*(Received 21 July 1969)*

A computer adaptable method for finding similarities in the amino acid sequences of two proteins has been developed. From these findings it is possible to determine whether significant homology exists between the proteins. This information is used to trace their possible evolutionary development.

The maximum match is a number dependent upon the similarity of the sequences. One of its definitions is the largest number of amino acids of one protein that can be matched with those of a second protein allowing for all possible interruptions in either of the sequences. While the interruptions give rise to a very large number of comparisons, the method efficiently excludes from consideration those comparisons that cannot contribute to the maximum match.

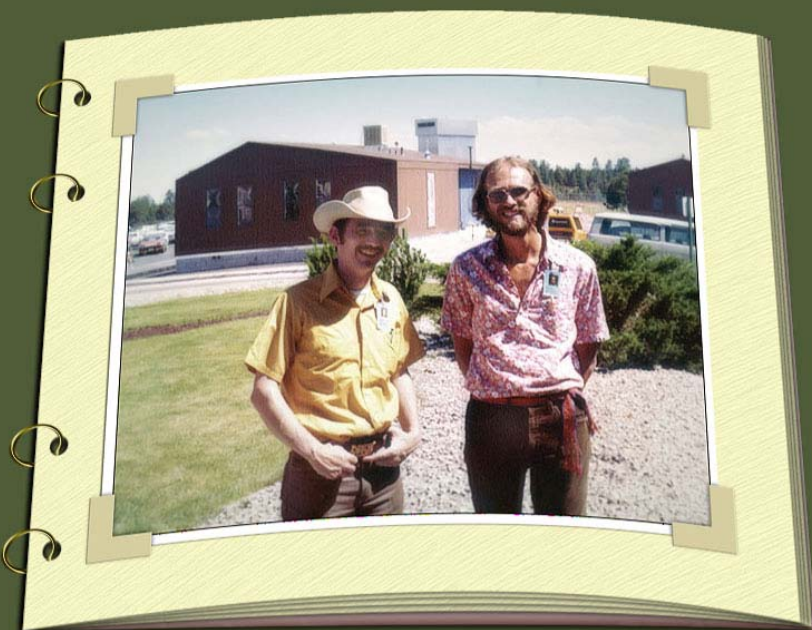
Comparisons are made from the smallest unit of significance, a pair of amino acids, one from each protein. All possible pairs are represented by a two-dimensional array, and all possible comparisons are represented by pathways through the array. For this maximum match only certain of the possible pathways must be evaluated. A numerical value, one in this case, is assigned to every cell in the array representing like amino acids. The maximum match is the largest number that would result from summing the cell values of every pathway.



Identify similar **sub-sequence**



Smith and Waterman at Los Alamos, New Mexico  
Photo by David Lipman, taken summer of 1980



(<http://www.cmb.usc.edu/people/msw/SmithWaterman.html>)

*J. Mol. Biol.* (1981), **147**, 195–197

### Identification of Common Molecular Subsequences

The identification of maximally homologous subsequences among sets of long sequences is an important problem in molecular sequence analysis. The problem is straightforward only if one restricts consideration to contiguous subsequences (segments) containing no internal deletions or insertions. The more general problem has its solution in an extension of sequence metrics (Sellers 1974; Waterman *et al.*, 1976) developed to measure the minimum number of “events” required to convert one sequence into another.

These developments in the modern sequence analysis began with the heuristic homology algorithm of Needleman & Wunsch (1970) which first introduced an iterative matrix method of calculation. Numerous other heuristic algorithms have been suggested including those of Fitch (1966) and Dayhoff (1969). More mathematically rigorous algorithms were suggested by Sankoff (1972), Reichert *et al.* (1973) and Beyer *et al.* (1979), but these were generally not biologically satisfying or interpretable. Success came with Sellers (1974) development of a true metric measure of the distance between sequences. This metric was later generalized by Waterman *et al.* (1976) to include deletions/insertions of arbitrary length. This metric represents the minimum number of “mutational events” required to convert one sequence into another. It is of interest to note that Smith *et al.* (1980) have recently shown that under some conditions the generalized Sellers metric is equivalent to the

$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \end{cases}$$

Global alignment

$$F(0,0) = 0$$

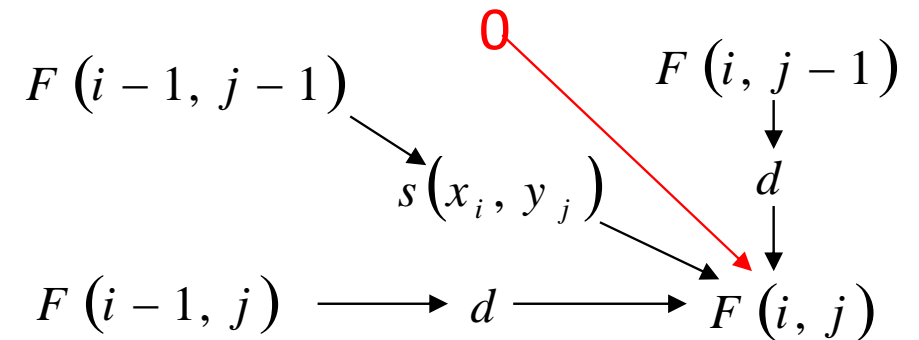
$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

Local alignment

# DP for Local alignment: Formula

$$F(0, 0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$



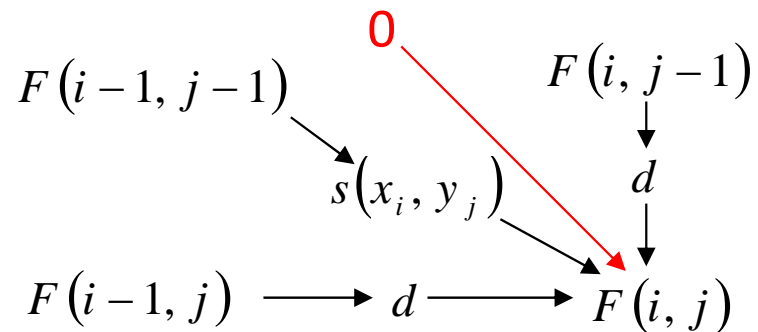


# DP for Local alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal **local alignment** of AAG and AGC.  
Use a linear gap penalty of  $d = -5$ .

		A	A	G
A				
G				
C				

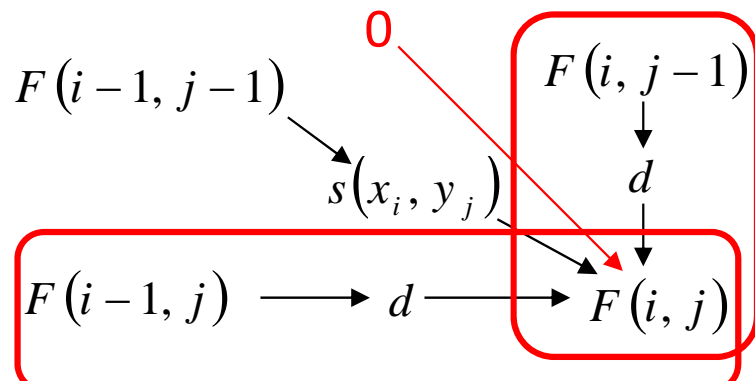


# DP for Local alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal **local alignment** of AAG and AGC.  
Use a linear gap penalty of  $d = -5$ .

		A	A	G
	0	0	0	0
A	0			
G	0			
C	0			

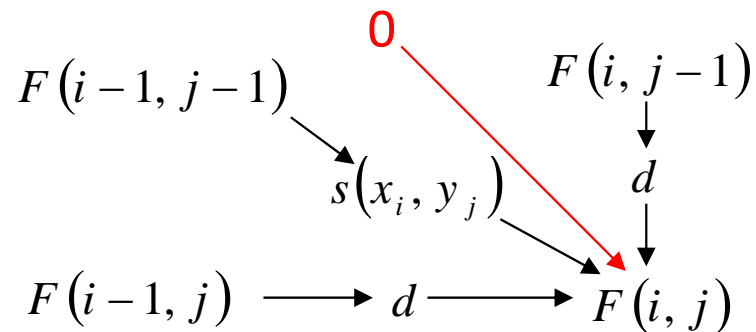


# DP for Local alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal **local alignment** of AAG and AGC.  
Use a linear gap penalty of  $d = -5$ .

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

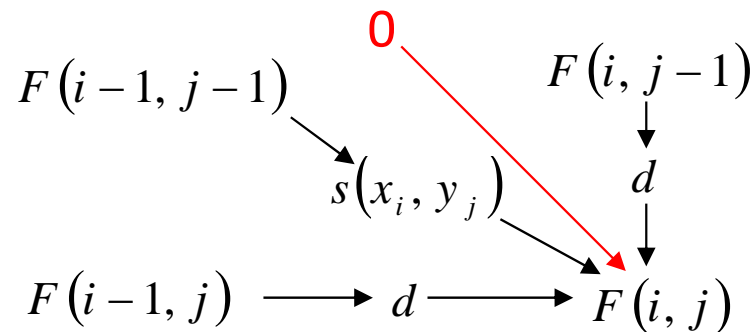


# DP for Local alignment: Example

	A	C	G	T
A	2	-7	-5	-7
C	-7	2	-7	-5
G	-5	-7	2	-7
T	-7	-5	-7	2

Find the optimal **local alignment** of AAG and AGC.  
Use a linear gap penalty of  $d = -5$ .

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0



# Traceback: Decode the Local Alignment

- Trace back begins at **the highest score** in the matrix and continues **until you reach 0**.

A G  
A G

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

# Traceback: Decode the Local Alignment

- And also the **secondary best** alignment

A  
A

		A	A	G
	0	0	0	0
A	0	2	2	0
G	0	0	0	4
C	0	0	0	0

# Global vs. Local

$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \end{cases}$$

A	A	G	-
-	A	G	C

A	A	G	-
A	-	G	C

$$F(0,0) = 0$$

$$F(i, j) = \max \begin{cases} F(i-1, j-1) + s(x_i, y_j) \\ F(i-1, j) + d \\ F(i, j-1) + d \\ 0 \end{cases}$$

A	G
A	G

A
A



# 生物信息学：导论与方法

## Bioinformatics: Introduction and Methods



<https://www.coursera.org/course/pkubioinfo>