

# Lecture 2 Sequencing

Date

No. L

A sequence of DNA is denoted by  $x_1 x_2 \dots x_k$  where  $x_i \in \{A, T, G, C\}$ . If  $S = x_1 x_2 \dots x_k$ , we denote  $x_i = S(i)$ .

The first problem we study about DNA sequences, is the comparison of two sequences and determine how "similar" they are.

In general, the sequences we compare may not of the same length, i.e., two sequences may have different number of terms, but not too "different" from each other. So, we may use the so-called alignment to adjust their length (the number of terms).

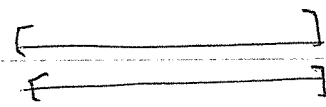
e.g.

A G A C C T A G

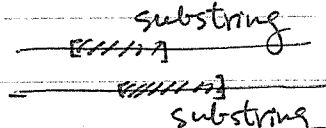
G C A C C T G C A G

Types

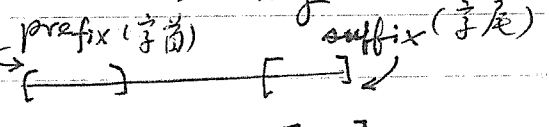
① Global alignment



② Local alignment



③ Semi-global alignment



(\*) Definition (Alignment)

The alignment is "the comparison of two or more nucleotide or protein sequences (strings) to determine the degree of similarity.

(\*) Commonly used to deduct functional or evolutionary relationships between genes and proteins.

(\*) Global alignments attempt to align every base (or amino acid) in each aligned sequence (string).

(\*) Local alignments will align only similar regions between sequences (strings), and leave regions with too many differences unaligned. (相差太多的部份没有调整的必要!)

(\*) Multiple Sequence alignment refers to the process of aligning three or more nucleotide or protein sequences to identify similarities between the sequences.  
(among)

Many algorithms are known so far!

## More details

3

### Definition (Alignment)

An alignment of two sequences  $S_1$  and  $S_2$  is a pair of sequences  $(S'_1, S'_2)$  obtained by insertion of spaces in  $S_1$  and  $S_2$  (respectively) such that

(1)  $|S'_1| = |S'_2|$ , and

(2)  $\forall i, S'_1(i)$  is aligned with  $S'_2(i)$  and either

$S'_1(i)$  or  $S'_2(i)$  is not a space. (不可以同時為 space!)

(Note)  $S_1$  can be viewed as a subsequence of  $S'_1$  if  
( $S_2$ ) ( $S'_2$ )  
we consider the insertions as elements.

### Example

$$S_1 = \text{AGAC} \leftarrow \langle A, G, A, C \rangle$$

$$S_2 = \text{TACCC} \leftarrow \langle T, A, G, C, C \rangle$$

$$S'_1 = *AGAC$$

$$S'_2 = TACCC$$

or

$$S'_1 = *AGAC*$$

$$\downarrow \quad \downarrow \downarrow$$

$$S'_2 = TACCC * C$$

### Score of Alignment

$$p \geq 0 \quad \text{if } x=y$$

two sequences

4

Similarity of  $S_1$  and  $S_2$

$|S_1| = |S_2|$

$$\text{Sim}(S_1, S_2) = \max_{\substack{(S'_1, S'_2) \\ \text{all alignments}}} \left\{ \alpha(S'_1, S'_2) = \sum_{i=1}^k \alpha(S'_1(i), S'_2(i)) \right\}$$

Definition (Optimal Alignments)

The pair  $(S'_1, S'_2)$  with largest score is called an optimal alignment.

For example: (For AGC and AAAC)

\* AGC            A \* GC            AG \* C  
 AAAC            AAAC            AAAC

are optimal alignments with  $\text{Sim}(AGC, AAAC) = -1$ .

$p=1, q=-1, r=-2$

Dynamic Algorithm

↓  
 $\frac{1}{\sqrt{6}} \frac{1}{\sqrt{10}}$

$a[i, j] = \text{Sim}(S_1[0:i], S_2[0:j])$

$$= \max \begin{cases} a[i, j-1] + q \\ a[i-1, j-1] + p \\ a[i-1, j] + r \end{cases}$$

$a[i, 0] = i \cdot q$

$a[0, j] = j \cdot r$

	*	A	G	C
*	0 → -2 → -4 → -6			
A	-2	1 → -1 → -3		
A	-4	-1	0 → -2	
A	-6	-3	-2	-1
r	-8	-5	-4	-1

Theorem The time to accomplish "finding optimal alignment" is  $O(|S_1| |S_2|)$  and the space we need is also  $O(|S_1| |S_2|)$ .

Proof. Since constant work is required per entry in the matrix, the proof follows. ■

Note No algorithm is known that uses asymptotically less time and has the same generality. Although, there are algorithms for more specific problems.

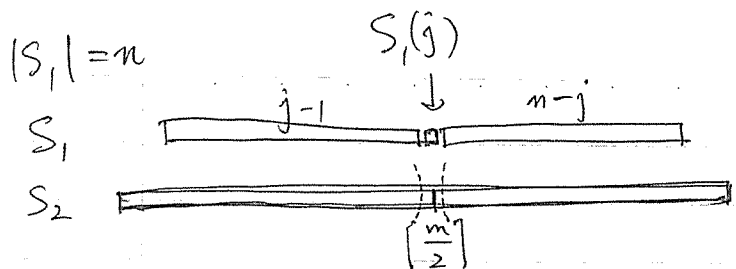
Note 儲存兩序列的 sequences 並不難; 但是, 如果每個 sequence 都有 10,000 字, 則寫出矩陣就要填入  $(10,000)^2 = 100,000,000$  壹佰萬個數, 是龐大的空間。

(\*) We are aiming to find the "optimal alignments" instead of "presenting the matrix".

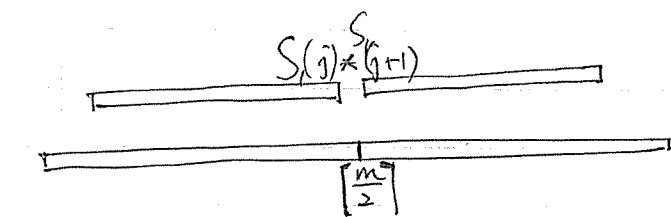
Another idea (By Bonnie Berger)

(\*)  $A(i, j)$ : Score of optimal alignment of  $S_1[1, i]$  and  $S_2[1, j]$ .

6



$$S_{im}(S_1, S_2) = A(|S_1|, |S_2|).$$



$\rightarrow A$   $\leftarrow B$  (把两 sequences 都引到这边)

Step 1 选择中间的字母 (从  $S_2$  中).

(这个字母对应的位置不外乎是  $S_1(j)$  或  $S_1(j)$  或  $S_1(j+1)$  中间的空格)

Step 2 计算 ( $\forall j$ )

$$\max \left\{ \begin{aligned} & (A(\lfloor \frac{m}{2} \rfloor - 1, j - 1) \pm 1) + B(\lfloor \frac{m}{2} \rfloor, n - j); \text{ and} \\ & (A(\lfloor \frac{m}{2} \rfloor - 1, j) - 2) + B(\lfloor \frac{m}{2} \rfloor, n - j). \end{aligned} \right.$$

Step 3

The maximum tell you the best score for aligning the  $\frac{m}{2}$  (th) character of  $S_2$  with some character of  $S_1$  or gap. (2n values)

Time we need

$$T(m, n) = 2cmn.$$

$$T(m, 1) = m$$

$$T(1, n) = n$$

$$T(m, n) \leq cmn + \max_{j=1}^n \left( 2c \left( \left\lceil \frac{m}{2} \right\rceil - 1 \right) \cdot j + 2c \left\lfloor \frac{m}{2} \right\rfloor (n-j) \right).$$

By induction

$$\leq cmn + \max_{j=1}^n (cmj + cm(n-j)).$$

$$\leq cmn + cmn.$$

$$\leq 2cmn.$$

## Global Alignments

Semiglobal Comparison :

To find the optimal alignment between suffix of  $S_1$ ,  
(prefix)

with prefix of  $S_2$ .

(suffix)

prefix

suffix

Example

$S_1$ : CAGCA\*CTTGG ATTCTCGG

$S_2$ : CAGCGTGG

sub-string

## Multiple Sequences Alignment (MSA)

### Definition (Sum-of-pairs, SP-score)

Let  $\alpha$  be a collection of alignments of  $S_1, S_2, \dots, S_n$ :  
 $S'_1, S'_2, \dots, S'_n$ . Then, the SP-score of  $\alpha$ , denoted by  
 $SP(\alpha) = \sum_{1 \leq i < j \leq n} \sigma(S'_i, S'_j)$ ,  $\sigma(S'_i, S'_j)$  is called the score  
of  $\alpha_{ij}$ , denoted by  $s(\alpha_{ij})$

(\*) Since MSAs are applied in finding the similarity of protein sequences, we use characters of proteins for examples.

$S_1$ : RCTLEE	$S'_1$ : RCTLEE
$S_2$ : RCLEE	$\alpha \Rightarrow S'_2$ : RC*LEE
$S_3$ : CTLEE	$S'_3$ : *CTLEE
$S_4$ : CTEE	$S'_4$ : *CT*EE

(With  $p=1$  and  $q=-2$ )

\*  $\leftrightarrow$  \*  
SCORE 0

$$SP(\alpha) = s(\alpha_{1,2}) + s(\alpha_{1,3}) + s(\alpha_{1,4}) + s(\alpha_{2,3}) + s(\alpha_{2,4}) + s(\alpha_{3,4})$$



MSA Problem

Given  $k$  sequences, find the optimal alignment of these  $k$  sequences.

Dynamic Programming

$$\text{Time-complexity} = O(k^2 n^k)$$

(Exponential in  $k$ !)

观察: 每一次调整都要和其之前的 sequences 配对一次!

研究进展

1. MSA problem is NP-complete

Wang and Jiang (1994) + Bonizzoni and Vedova (2001)

2. Gusfield (1993):  $(2 - \frac{2}{k})$ -approximation algorithm  
(If  $k=2$ , then the algorithm is working.)

3. Pezner (1992):  $(2 - \frac{3}{k})$ -approx. algorithm.

(\*) 4. Bafna, Lawler and Pezner (1997): Approximation algorithms for MSA, Theoretical Computer Science, 182, 233-244.  
 $(2 - \frac{l}{k})$ -approx. algorithm for  $l < k$ .

For references[Home](#)[Sequence Alignment Software](#)[Sequence Alignment Web Resources](#)[Glossary](#)[About](#)[Contact](#)

## DNA Sequence Alignment

Sequence alignment describes the way of aligning DNA, RNA, or protein sequences to highlight or identify similarities between DNA sequences. Typically, gaps have to be inserted into sequences so that identical or similar nucleotides or amino acids are aligned in columns. Here is an example:

Scarites	C	T	T	G	A	T	C	G	T	A	C	C	A	A	-	-	-	A	A	T	A	T	T	A	C	
Carenum	C	T	T	G	A	T	C	G	T	A	C	C	A	C	A	-	T	A	C	-	T	T	T	A	C	
Pasimachus	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	T	A	T	A	A	G	T	T	A	C	
Pheropsophus	C	T	T	G	A	T	C	G	T	A	C	C	A	C	-	-	-	C	A	T	A	T	T	A	C	
Brachinus armiger	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	T	T	A	T	T	T	A	C
Brachinus hirsutus	A	T	T	A	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	T	A	T	T	T	A	C
Aptinus	C	T	T	G	A	T	C	G	T	A	C	C	A	C	-	-	-	C	A	T	A	T	T	A	C	
Pseudomorpha	C	T	T	G	A	T	C	G	T	A	C	C	A	C	-	-	-	A	C	A	A	T	T	A	C	

Sequence alignment is often used in several biomedical research fields, including phylogenetics, biogeography, and evolution research.

Here is an example of a phylogenetics experiment that includes DNA sequence alignments. Our researcher works for a museum of natural history, and wants to study the evolutionary relation between several animal species he studies. A typical approach would be:

- Isolate DNA from each of the species studied
- Generate DNA sequences for a gene region, an entire gene, or several genes
- Assemble the DNA sequences for each species into separate "contigs"
- Edit the sequence assembly results to remove any errors
- Align the contig sequences to each other using a multiple sequence alignment program
- Verify any observed differences by going back to the original DNA sequences
- Correct the placement of gaps in the aligned sequences, if necessary
- Export the sequence alignment for further analysis with phylogenetics software, for example to generate phylogenetic trees

This web site provides links to commonly used programs and web resources for DNA sequence alignments. Since hundreds of different programs and relevant web sites exist, the goal is *not* to provide lists, but rather to concentrate on the most commonly used and the most useful sequence alignment software.

[Sequence Alignment Terms Explained](#)

[Sequence Alignment Software](#)

[Sequence Alignment Web Resources](#)

[Glossary](#)

[About](#)

[Contact](#)

## Sequence Alignment Terms

### Alignment

The comparison of two or more nucleotide or protein sequences to determine the degree of similarity. Commonly used to deduct functional or evolutionary relationships between genes and proteins.

### Assembly

The process of combining short DNA sequence fragments into larger units by looking for overlaps between different fragments. Often required because the length of the genes studied exceeds the length of the sequence fragments produced by DNA sequencing machines. Also used to combine several fragments that cover the same region, for example in forward and reverse direction, with the goal to reduce errors in the consensus sequence.

### Consensus Sequence

A single sequence generated from an alignment or assembly of sequence fragments that is the "best fit" for the given sequences. Historically, majority ("vote based") and inclusive methods were most commonly used to determine consensus sequence. For sequence assemblies, these methods have often been replaced by quality-based consensus methods. Quality-based consensus sequences are typically more accurate than majority-based sequences, and can reduce the need for manual editing of sequence assemblies drastically.

### Contig

The result of a sequence assembly or alignment that shows the arrangement of the fragments to form a contiguous large sequence.

### Dynamic Programming

A computer-science based [method](#) to find the optimal alignment between sequences. For two sequences, this algorithm creates a two-dimensional matrix based on identity or similarity of bases (or amino acids) in both sequences, and then finds the highest-scoring path to obtain the alignment. A commonly used dynamic programming method is the [Needleman-Wunsch algorithm](#). A nice graphical display of the dynamic programming methods for sequence alignments can be found [here](#).

### Global Alignments

Global-alignments attempt to align every base (or amino acid) in each aligned sequence.

### Local Alignments

Local alignments will align only similar regions between sequences, and leave regions with too many differences unaligned. Local alignments can be better suited for the alignment of very dissimilar sequences. In sequence assembly, the program [Phrap](#) demonstrate that local alignments can be used to reduce or eliminate the need to remove low-quality sequence (end clipping) before assembly.

### Multiple Sequence Alignment

[Multiple sequence alignment](#) refers to the process of aligning three or more nucleotide or protein sequences to identify similarities between the sequences. Alignments that include many sequences can be computational intensive, and require more sophisticated algorithms than pairwise alignments.

### Pairwise Alignment

In pairwise sequence alignment, exactly two nucleotide or protein sequences are aligned to each other to determine the similarity between the two sequences.

### Word-based Alignment Methods

Word-based alignment methods are an optimization often used in sequence alignment and assemblies. Instead of examining every single nucleotide or amino acid, "words" of a fixed length are analyzed. This can lead to substantial reductions in memory use and alignment times. One common application is to use the number of shared words between two sequences to estimate the similarity in early phases of sequence alignments, or to identify sequences that share overlaps in sequence assembly.