

## Finding Signals in DNA

- The first signal in DNA was found in 1970 by Hamilton Smith after the discovery of the Hind II restriction enzyme.
- The palindromic site of the restriction enzyme is a signal that initiates DNA cutting.

(Fact) 1. Hamilton Smith was lucky: restriction sites are the simplest signals in DNA. (Reliably find in DNA.)

2. Most other signals (promoters, splicing sites, etc.) are so complicated that we don't yet have good models or reliable algorithms for their recognition.

(\*\*) Understanding gene regulation is a major challenge in computational biology. For example, regulation of gene expression may involve a protein binding to a region of DNA to affect transcription of an adjacent gene.

(Fact) Since protein-DNA binding mechanisms are still insufficiently understood to allow "in silico" prediction of binding sites, the common experimental approach is to locate the approximate position of the binding site.

(\*)\*) These experiments usually lead to identification of a DNA fragment of length  $n$  that contains a binding site (an unknown magic word) of length  $l \ll n$ .

•• Of course, one experiment is insufficient for finding the binding site, but a sample of experimentally found DNA fragments gives one hope of recovering the magic word.

GATTCTTAGGC  
 TATACGTTTGA  
 TGATTGACTTC  
 ⋮

Problem

Given a sample of  $K$  sequences where an (unknown) magic word appears at different (unknown) positions in these  $K$  sequences. Find the magic word!

Common Sense Approach (Staden, 1989; Wolfertstetter et al, 1996; Tompa, 1999)

Since  $l \ll n$ , we may test all words of length  $l$  and find those that appear in all (or almost all)  $K$  sequences.

Note  $l$  可能也不知道有多長; 所以可從較小的長度測試再增長; 一直到有"唯一"出現的 Magic word 為止。

- ① The described approach usually works fine for short continuous words such as GATTC, the restriction site of EcoRI.
- ② In the idea of randomness, this is also quite difficult, since we have  $4^6$  words to compete for the magic word.
- ③ The problem gets even more difficult (complicated) when the magic word has gaps.

(4)

For example,

$CCAN_9TGG$  (Xcm I restriction enzyme)

↑ a gap of length 9

$(P_u)^m C N_{40-2000} (P_u)^m C$  (McrBC Endonuclease)

↑  
A or G

$TTGACAN_{17}TATAAT$  (E. Coli promoters)

Enumeration and check of all patterns of the above types are hardly possible due to computational complexity.

### Remark

寻找 Magic Words 基本上就像密码学中寻找 Key 的概念,除了利用各种可能提供的讯息之外,还要发挥一些想象力来协助判断;把所有可能的答案都试一次绝非良策。

(5)

• DNA linguistics is at the heart of the pattern-driven approach to signal finding, which is based on enumerating all possible patterns and choosing the most frequent or the fittest among them.

(\*) The fitness measures vary from estimates of the statistical significance of discovered signals to the information content of the fragments that approximately match the signal.

Steps for pattern-driven approach

Step 1. Define the frequency or fitness measure.  
(f.f.m.)

Step 2. Calculate the f.f.m. of each word w.r.t. ~~the~~ sample DNA fragments.

Step 3. Report the fittest words as potential signals.

(Fact)

Note that if  $A$  denotes the set of alphabets, then the search space for patterns of length  $l$  is  $|A|^l$ .

① DNA texts are not easy to decipher, and there is little doubt that Nature can construct an "enigma" of the kind which human ingenuity may not resolve.

- A popular approach in DNA linguistics is based on the assumption that frequent or rare words may correspond to signals in DNA.

- A word occurs considerably more (or less) frequently than expected has the potential to become a "signal".

What is its biological meaning?

For example, Gelfand and Koonin (1997) showed that the most avoided 6-palindrome in the archaeon *M. jannaschii* is likely to be the recognition site of a restriction-modification system.

2

- ① To find frequent and rare words  $(W)$  in a text, one has to compute expected value  $E(W)$  and the variance  $\sigma^2(W)$  for the number of occurrences (frequency) of each word  $W$ .
- ② Afterwards, the frequent and rare words are identified as the words with significant deviations from expected frequencies.

### Remark

In many DNA linguistics papers, the variance  $\sigma^2(W)$  of the number of occurrences of a word in a text was erroneously assumed to be  $E(W)$ .

要探讨 word's occurrence 的机率就不是一件直观(或简单)的事。参考 "Correlation Polynomial" (Autocorrelation)

8

## The best bet for simpletons

The best bet for simpletons starts out with two players who select words of length  $l$  in 0-1 alphabet. Player I selects a sequence  $A$  of  $l$  alphabets from  $\{0, 1\}$ , and Player II, knowing what  $A$  is, selects another sequence  $B$  of length  $l$ . The players then flip a coin to obtain heads (1) or tails (0) in turns until either  $A$  or  $B$  appears as a block of  $l$  consecutive outcomes. If  $A$  comes first, then  $A$  wins the game. (B)

(Fact 1) The game will terminate with probability 1, assuming the coin is a fair one or at least both sides have positive probability.

(Fact 2)  $B$  has higher probability to win the game.

(Fact 3) The best bet for simpletons is a non-transitive game.



John Conway

(10)

The odds that B will win over A is

$$\frac{K_{AA} - K_{AB}}{K_{BB} - K_{BA}}$$

Martin Gardner, 1974 :

I have no idea why it works. It just cranks out the answer as if by magic, like so many of Conway's other algorithms.

$$A = 00$$

$$B = 10$$

$$K_{AB}(t) = (0, 0) \sim 0 + 0t$$

$$K_{AA}(t) = (1, 1) \sim 1 + t$$

$$K_{BB}(t) = (1, 0) \sim 1$$

$$K_{BA}(t) = (0, 1) \sim 0 + 0t$$

$$\frac{1 + \frac{1}{2}}{1 - \frac{1}{2}} = \frac{\frac{3}{2}}{\frac{1}{2}} = \frac{3}{1}$$

Proof. 1980, Li || 1981, Guibas and Odlyzko,

1993, Pevzner

shortest one !