



FASTA [Pearsor	n & Lipman 1988]
1. k-tup (k=6 fo	r DNA, 2 for proteins)
rule of thumb sensitive the s	the smaller the word, the slower and more earch is
2. build a lookup	table (also called as hash table or dictionary)
word	positions_in_the_query
АААААА	XXX
AAAAAT	XXX
AAAAAC	XXX
	CISC889, S04, Lec5-6, Liao

## FASTA cont'd

3. scan through the database: for each every word of size k, look up in the lookup table in step 2. The offsets between the positions of the word in the query and the database entry are calculated and saved.

(implication: same offsets suggest segment of similarity.)

4. join nearby *contiguous* stretches of similarity (diagonal) =>scores init1.

5. join adjacent diagonals into a single long region (by introducing gaps) => scores initN.

6. do a dynamic programming algorithm for regions with high initN score to determine the OPT score.

Q: what is the FASTA format?

Basic Local Alignment Search Toolkit [Altschul et al, 1990]

- 1. A list of neighborhood words of fixed length (3 for protein and 11 for DNA) that match the query with score > a threshold.
- 2. Scan the database sequences and look for words in the list; once find a spot, try a "hit extension" process to extend the possible match as an ungapped alignment in both directions, stopping at the maximum scoring extension.



Sig	mificance of scores
Pairv	wise alignment goals:
(1)	) whether and
(2)	) how two sequences are related.
It is	rare that you have just two particular sequences to compare. More often, you have one query sequence and a large database of sequences.
Data	base searching: find all sequences in the database that are related to the query sequence.
Solu	tion:
(1)	For each sequence in the database, use Smith-Waterman/FASTA/BLAST to align with the query sequence and return the score of the optimal alignment.
(2) F	Rank the sequences by the score.
Q: h	ow good is a score?
	CISC889, S04, Lec5-6, Liao



Score Statistics (cont'd) The probability that the maximum S is smaller than x is

 $P(S < x) = \prod_i [1 - Pr(\sigma_i > x)) ] \rightarrow exp[-\kappa e^{-\lambda x}]$  when  $\kappa \rightarrow \infty$ .

This is a form of **Extreme Value Distribution**.

Definition: Take N samples from the density g(u), the probability that the largest amongst them is less than x is G(x)N, where  $G(x) = \int_{-\infty}^{x} g(u) \, du$ . The density is Ng(x)G(x)(N-1). The limit for large N of Ng(x)G(x)(N-1) is called the EVD for g(x). (Draw a EVD curve).

Once the score distribution is known, for a given score, the percentage of the curve is to its right gives the p-value of that score.

**p-value** = probability of at least one sequence scoring above S in the given database.

 $P(S < x) = 1 - exp[-\kappa e^{-\lambda x}].$ 

**E-value** = expected number of scores better than S in a database search.

 $E(S) = kmn e^{-\lambda S}$ .

	• all of the above discussions only applicable to local alignments.
	•For gapped local alignments, same statistics, although not proved, are believed to apply.
	• The trick is to learn $\lambda$ and K. These values depend upon the substitution matrix and can be estimated from randomly generated data. The parameter $\lambda$ is the positive root of the equation
	$\sum_{a,b \in alphabet} s(a,b)p(a)b(b) = 1$
	• score statistics for global alignments are not well known.
Q:	what is a bit score in the blast search result?
A:	bit score is defined as $S' = (\lambda S - \ln K)/\ln 2$
	it is then convenient to calculate the e-value
	$E(S) = mn \ 2^{-S'}$
	CISC889, S04, Lec5-6, Liao



CISC 889, 504, Lecs-0, Liao

Bayesian perspective 
$$\begin{split} P(M|x,y) &= P(x,y|M)P(M)/P(x,y) \\ &= P(x,y|M)P(M) / [P(x,y|R)P(R) + P(x,y|M)P(M)] \\ &= \frac{[P(x,y|M)P(M)/P(x,y|R)P(R)]}{[1+P(x,y|M)P(M)/P(x,y|R)P(R)]} \end{split}$$
if 
$$\begin{split} S &= \log[P(x,y|M)/P(x,y|R)], \text{ this is the score returned from a SW alignment.} \\ S' &= S + \log[P(M)/P(R)] \\ \text{then} \\ P(M|x,y) &= \exp(S') / [1 + \exp(S')] \\ (\text{draw a curve and explain}) \end{split}$$
CISC889, S04, Lec5-6, Liao

## Gap penalties

Linear

 $\gamma(g) = \text{-} g \; d \label{eq:gamma}$  where g is the gap length and d is the penalty for a gap of one base

Affine

 $\gamma(g) = -d - (g-1)e$ 

where d is gap-open penalty and e is gap-extension penalty

Note: gap penalty is a sort of gray area due to less knowledge about gap distribution.