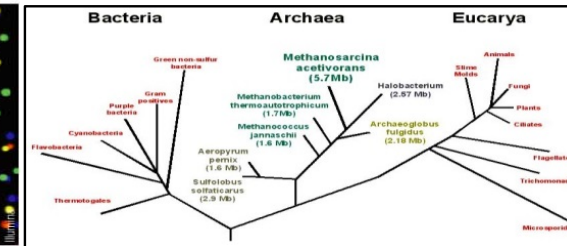
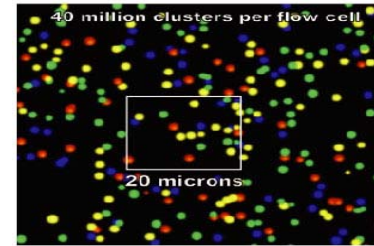


TAACCCTAACCCCTAACCCCTAACCCCTAACCCCTA
CCTAACCCCTAACCCCTAACCCCTAACCCCTAACCC
CCCTAACCCCTAACCCCTAACCCCTAACCCCTAAC
AACCCCTAACCCCTAACCCCTAACCCCTAACCCCTA
ACCCTAACCCCAACCCCAACCCCAACCCCAAC
CTACCCTAACCCCTAACCCCTAACCCCTAACCCCTA
ACCCTAACCCCTAACCCCTAACCCCTAACCCCTAA

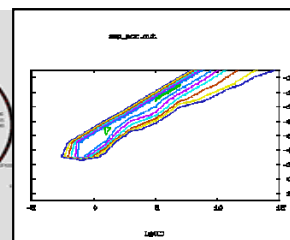
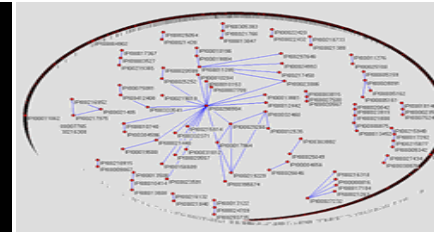
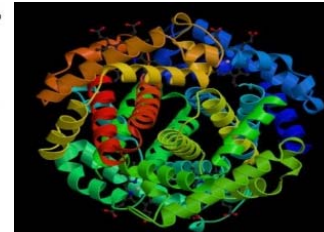
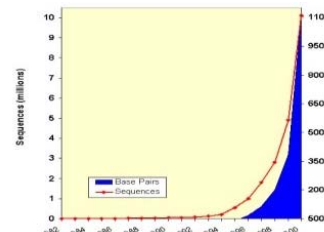


Unit 2: BLAST Algorithm: a Primer

北京大学生物信息学中心 高歌

Ge Gao, Ph.D.

Center for Bioinformatics, Peking University



Basic Local Alignment Search Tool

Stephen F. Altschul¹, Warren Gish¹, Webb Miller²
Eugene W. Myers³ and David J. Lipman¹

¹National Center for Biotechnology Information
National Library of Medicine, National Institutes of Health
Bethesda, MD 20894, U.S.A.

²Department of Computer Science
The Pennsylvania State University, University Park, PA 16802, U.S.A.

³Department of Computer Science
University of Arizona, Tucson, AZ 85721, U.S.A.

(Received 26 February 1990; accepted 15 May 1990)

A new approach to rapid sequence comparison, basic local alignment search tool (BLAST), directly approximates alignments that optimize a measure of local similarity, the maximal segment pair (MSP) score. Recent mathematical results on the stochastic properties of MSP scores allow an analysis of the performance of this method as well as the statistical significance of alignments it generates. The basic algorithm is simple and robust; it can be implemented in a number of ways and applied in a variety of contexts including straight-forward DNA and protein sequence database searches, motif searches, gene identification searches, and in the analysis of multiple regions of similarity in long DNA sequences. In addition to its flexibility and tractability to mathematical analysis, BLAST is an order of magnitude faster than existing sequence comparison tools of comparable sensitivity.

Basic local alignment search tool

SF Altschul, W Gish, W Miller, EW Myers... - Journal of molecular ..., 1990 - Elsevier

A new approach to rapid sequence comparison, **basic local alignment search tool** (BLAST), directly approximates alignments that optimize a measure of **local** similarity, the maximal **segment pair** (MSP) score. Recent mathematical results on the stochastic properties of ...

Cited by 47577 Related articles All 98 versions Cite



There are nm entries in the matrix.

Sequence X of length m

Sequence Y of length n

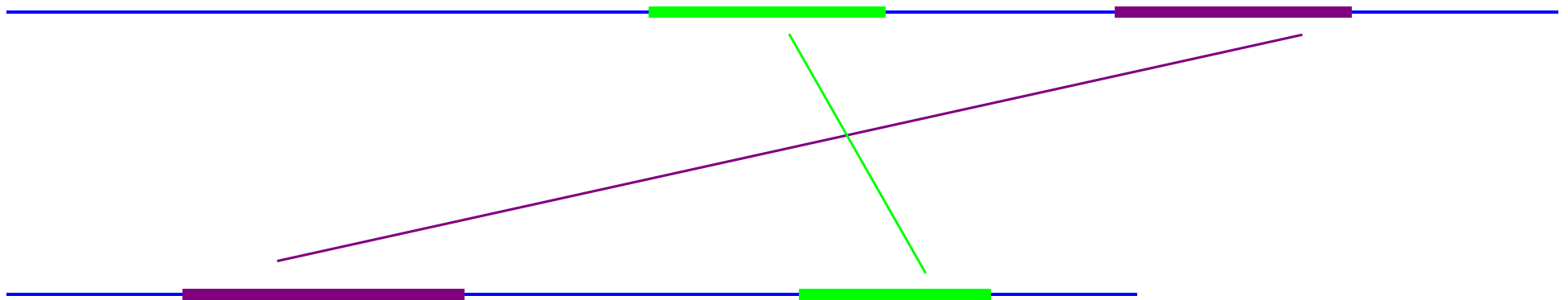
Each entry requires a constant number c of operation(s).

Dynamic programming matrix

$c * m * n$ operations needed in total, for one pair-wise alignment.

BLAST Ideas: Seeding-and-extending

1. Find matches (**seed**) between the query and subject
2. Extend seed into High Scoring Segment Pairs (**HSPs**)
 - Run Smith-Waterman algorithm on the specified region only.
3. Assess the reliability of the alignment.

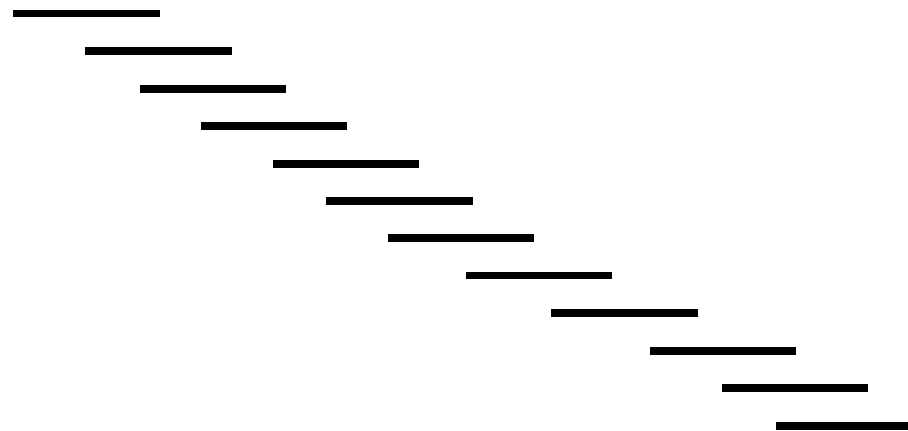


Seeding

For a given **word length w** (usually 3 for proteins and 11 for nucleotides), slicing the query sequence into multiple continuous “**seed words**”

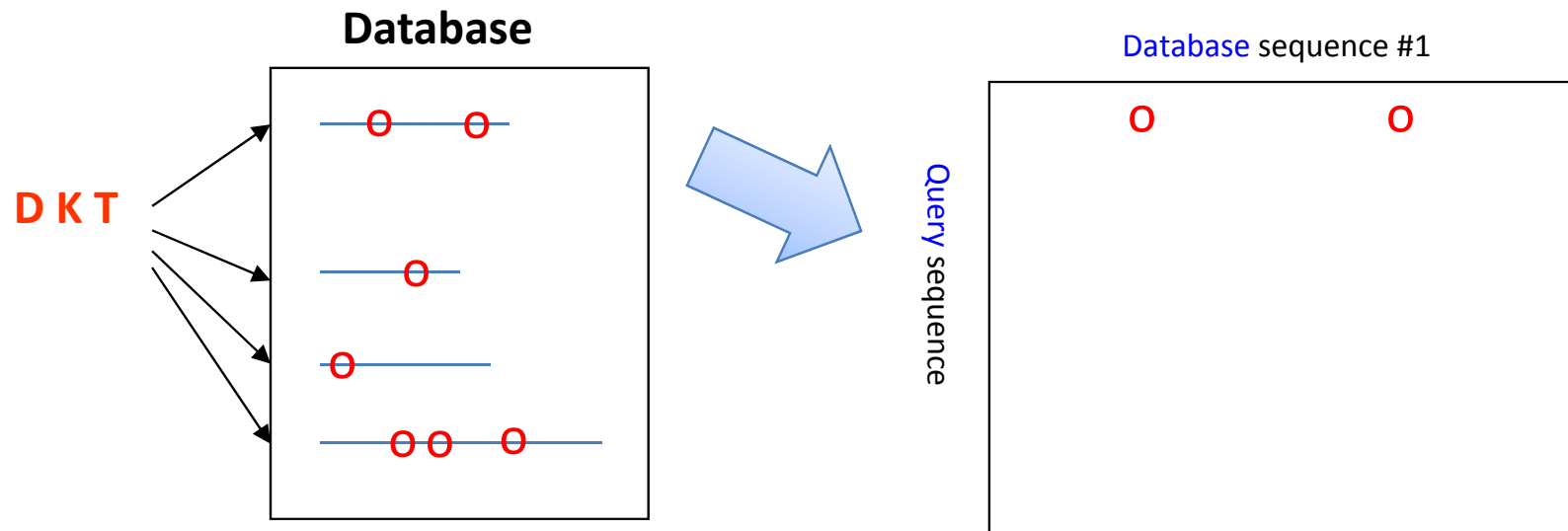
Query Sequence

M V L S P A D K T N V K A A W

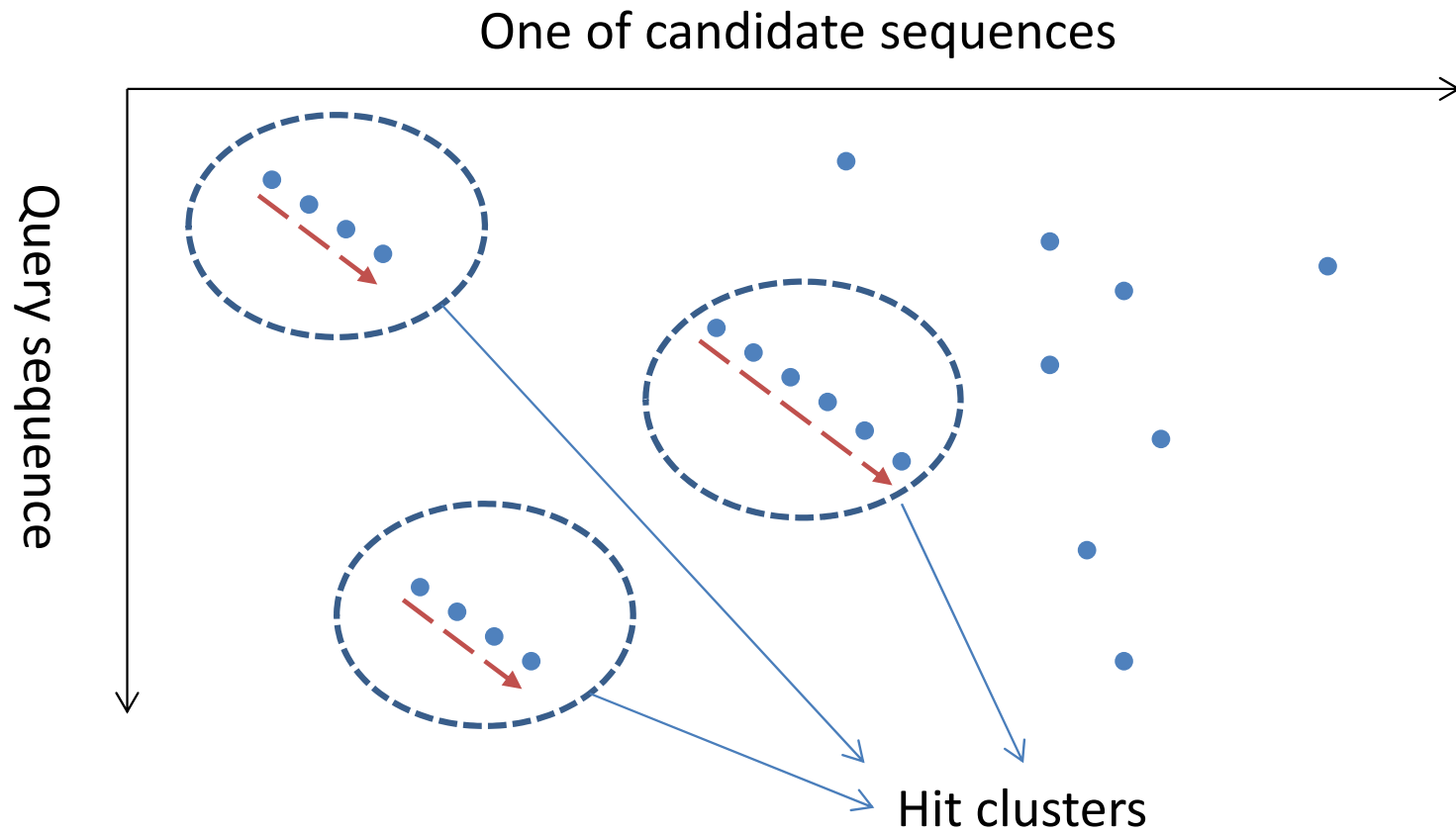


Speedup: Index database

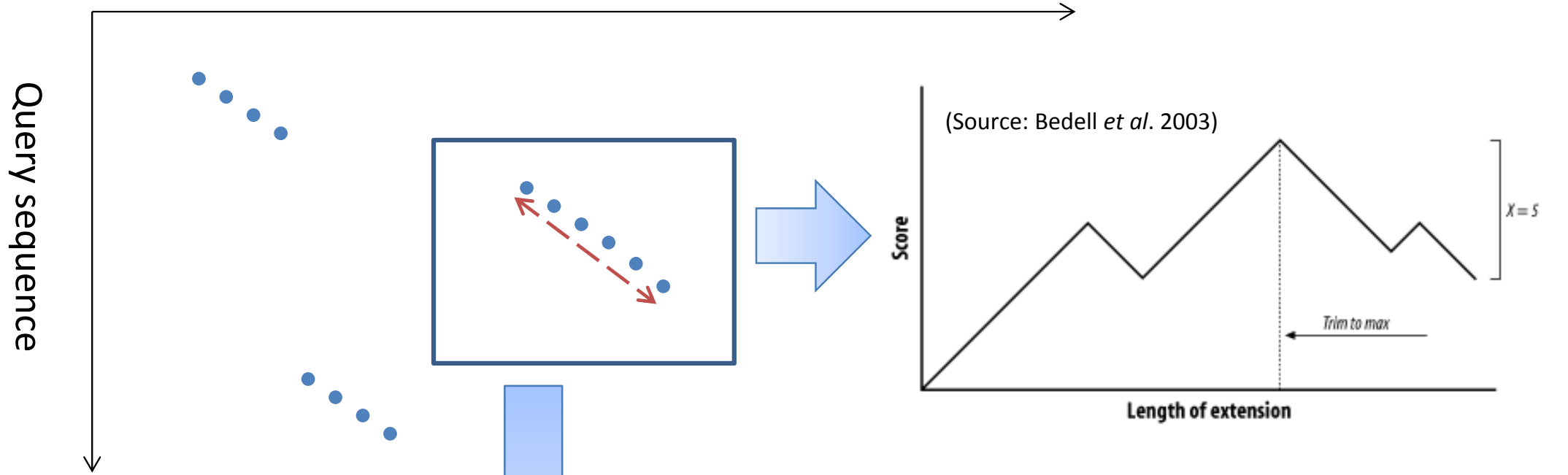
The database was pre-indexed to quickly locate all positions in the database for a given seed.



Diagonal *and* Two-hits



One of candidate sequences



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

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

Speedup: mask low-complexity

- Low complexity sequences yield false positives.

- CACACACACACACACA
- KKKLKKLKKLKKL

$$K = \frac{1}{L} \log_N \left(\frac{L!}{\prod_i n_i!} \right)$$

Diagram illustrating the formula for K (complexity) based on window length L and the frequency of letters n_i in the window. The alphabet size N is noted as 4 or 20.

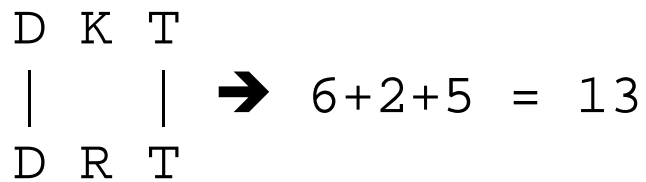
Labels in the diagram:

- Window length (points to L)
- Frequency of the i th letter (points to n_i)
- Alphabet size (4 or 20) (points to N)

For example, for typical microsatellite “CACACACACACACA”, with window length 6:

$$\begin{aligned} K &= \frac{1}{6} \log_4 \left(\frac{6!}{n_A! * n_C! * n_G! * n_T!} \right) \\ &= \frac{1}{6} \log_4 \left(\frac{6!}{3! * 3! * 0! * 0!} \right) \\ &= \frac{1}{6} \log_4 \left(\frac{6!}{3! * 3!} \right) \\ &= \frac{1}{6} \log_4 20 \\ &= 0.36 \end{aligned}$$

To improve sensitivity, in addition to the **seed word** itself, the BLAST also use these highly similar “**neighbourhood words**” (based on the substitution matrix) for seeding.

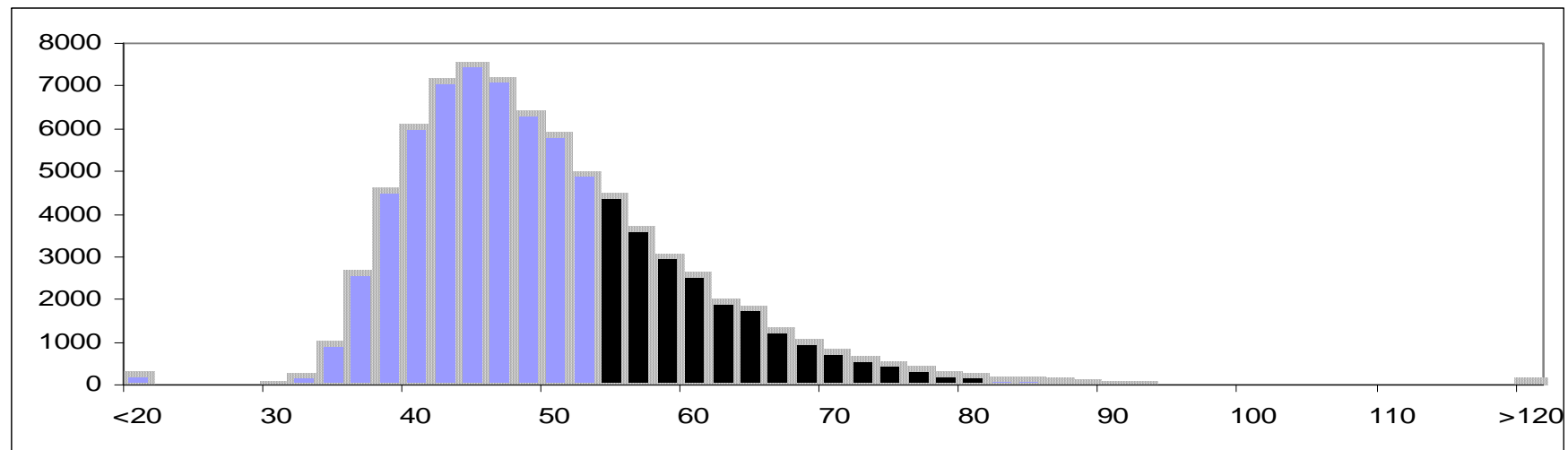


- DKT 16
- DRT 13
- DET 12
- DKS 12
- DQT 12
- EKT 12
- DKA 11
- DKN 11
- DKV 11
- DNT 11
- DST 11
- NKT 11
- DAT 10
- DDT 10
- DHT 10
- DKC 10
- DKD 10
- DKE 10
- DKI 10
- DKK 10
- DKL 10
- DKM 10
- DKP 10
- DKQ 10
- DKR 10
- DMT 10
- DPT 10
- DTT 10
- QKT 10
- SKT 10

	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W		
C	9																				C	
S	-1	4																				S
T	-1	1	5																			T
P	-3	-1	-1	7																		P
A	0	1	0	-1	4																	A
G	-3	0	-2	-2	0	6																G
N	-3	1	0	-2	-2	0	6															N
D	-3	0	-1	-1	-2	-1	1	6														D
E	-4	0	-1	-1	-1	-2	0	2	5													E
Q	-3	0	-1	-1	-1	-2	0	0	2	5												Q
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8											H
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5										R
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5									K
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5								M
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4							I
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4						L
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4					V
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6				F
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7			Y
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11		W
C		S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W		

Quality Assessment

Given the large data volume, it's critical to provide some measures for assessing the **statistical significance** of a given hit.



For a given amino acid, we have the chance of $1/20$ to have a random match (as there are just 20 different amino acids in total).

Thus, for an amino acid sequence with length L , the probability of having a random match across the full length is $(1/20)^L$

Say you have a 6 AA peptide (not so unusual, e.g. the Tryptophyllin-T2-6 in *Phyllomedusa azurea* or “Orange-legged monkey frog, 橙腿猴树蛙” is a 6 AA peptide), then the odd would be $(1/20)^6 = 1.56 * 10^{-8}$

Looks not so big, huh?

But you're searching Swiss-Prot databases which contains 192,206,270 amino acids in 540,958 sequences (Sept 18th, 2013). Then you would expect to have $(1/20)^6 * 192,206,270 = 3.00$ matches by chance, for any 6 AA peptide!

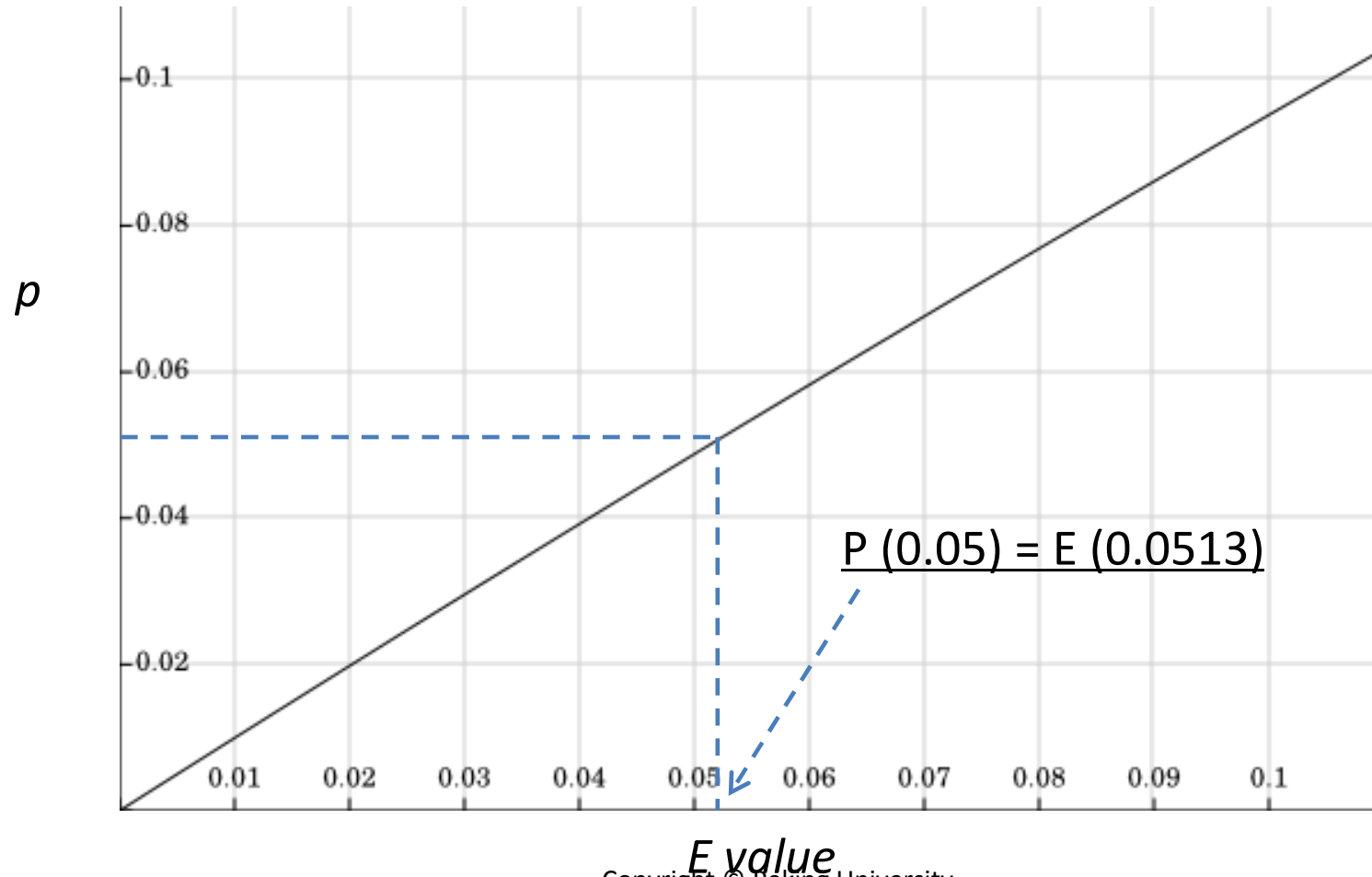
E-Value: How a match is likely to arise **by chance**

- **The expected number** of alignments with a given score that would be expected to occur **at random** in the database that has been searched
 - e.g. if $E=10$, 10 matches with scores this high are expected to be found **by chance**

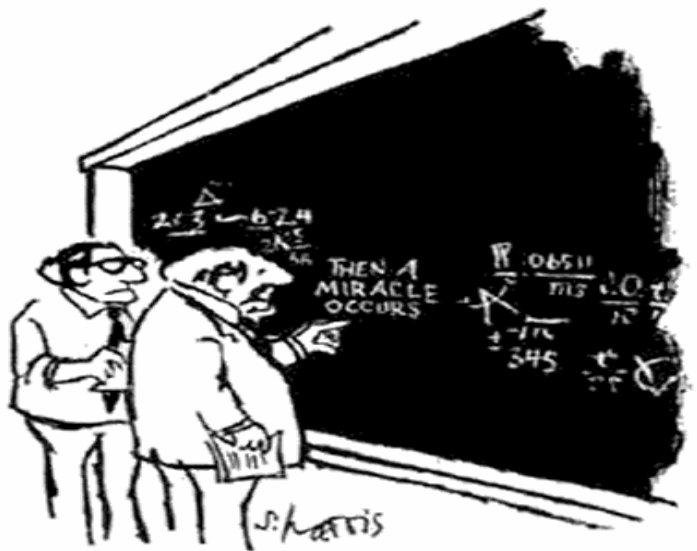
$$E = kmne^{-\lambda S}$$

$$E = kmne^{-\lambda S}$$

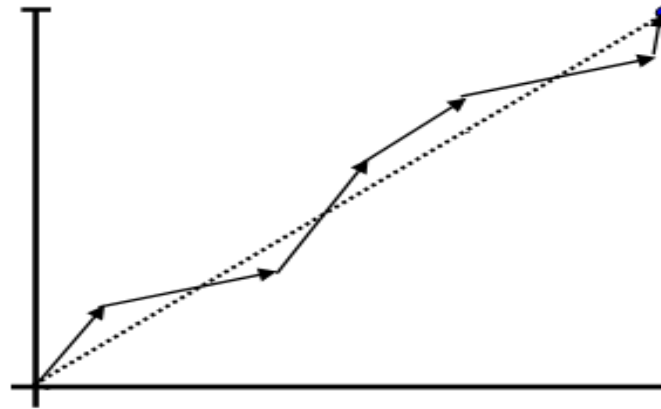
$$p = 1 - e^{-E}$$



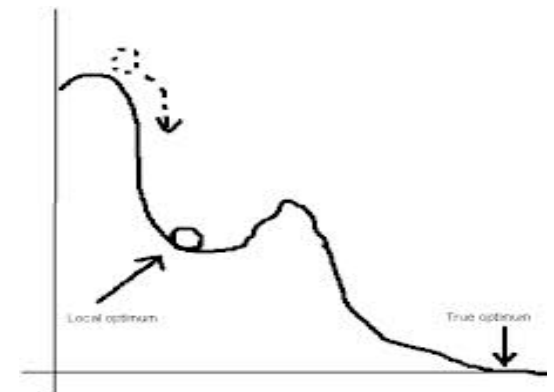
Heuristic (pronounced hyu-RIS-tik, Greek: "Εὕρισκω", "**find**" or "**discover**") refers to experience-based techniques for problem solving, learning, and discovery. (Source: Wikipedia)



"I THINK YOU SHOULD BE MORE EXPLICIT
HERE IN STEP TWO."



Not best,
but good *enough*



- Key heuristics in BLAST
 - Seeding-and-extending: looking for seeds of high scoring alignments ONLY
 - Use dynamic programming selectively
- Tradeoff: speed vs. sensitivity
 - Empirically, 1000 ~ 10000 times faster than plain Dynamic-Programming-based local alignment
 - But suffer from low sensitivity, especially for distant sequences (e.g. *E.coli* → human)

Summary Questions

- Why do you need E-value?
- Could you give an alignment that the Smith-Waterman algorithm, but not the BLAST algorithm, can identify? Explain your case.

生物信息学：导论与方法

Bioinformatics: Introduction and Methods



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