



Dept. of Computer Science and Engineering
 University of Rajshahi
 www.ru.ac.bd

Dr. Shamim Ahmad

BLAST:
 Basic local alignment
 search tool

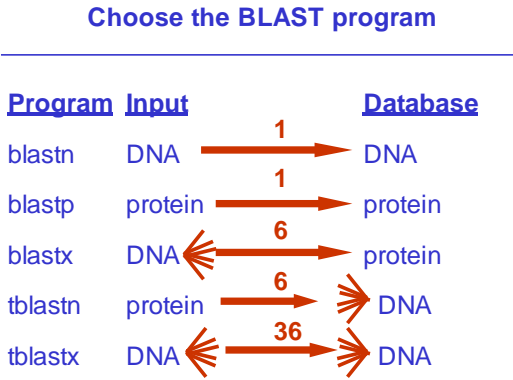
BLAST!

Courtesy of Jonathan Pevsner
 Johns Hopkins U.

BLAST

BLAST (Basic Local Alignment Search Tool) allows rapid sequence comparison of a query sequence against a database.

The BLAST algorithm is fast, accurate, and web-accessible.



BLAST: background on sequence alignment

There are two main approaches to sequence alignment:

[1] Global alignment (Needleman & Wunsch 1970) using dynamic programming to find optimal alignments between two sequences. (Although the alignments are optimal, the search is not exhaustive.) Gaps are permitted in the alignments, and the total lengths of both sequences are aligned (hence "global").

BLAST: background on sequence alignment

[2] The second approach is local sequence alignment (Smith & Waterman, 1980). The alignment may contain just a portion of either sequence, and is appropriate for finding matched domains between sequences. S-W is guaranteed to find optimal alignments, but it is computationally expensive (requires $(O)n^2$ time).

BLAST and FASTA are heuristic approximations to local alignment. Each requires only $(O)n^2/k$ time; they examine only part of the search space.

How a BLAST search works

"The central idea of the BLAST algorithm is to confine attention to segment pairs that contain a word pair of length w with a score of at least T ."

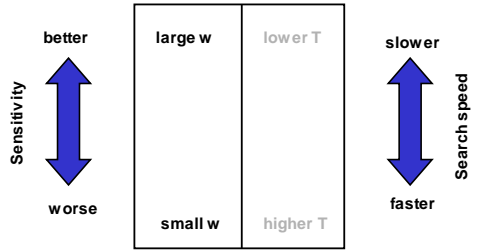
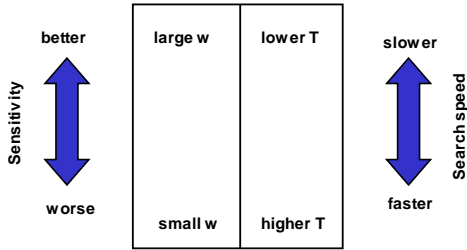
Altschul et al. (1990)

How the original BLAST algorithm works: three phases

Phase 1: compile a list of word pairs ($w=3$) above threshold T

Example: for a human RBP query
...FSG**GTW**YA... (query word is in yellow)

A list of words ($w=3$) is:
FSG SGT **GTW** TWY WYA



For proteins, default word size is 3.
(This yields a more accurate result than 2.)

BLAST Algorithm

RDQ	16	QDQ	12	EDQ	11	RDN	11	RDB	11	BDQ	10	RDQ	10
RBQ	14	REQ	12	HDQ	11	RDD	11	ADQ	10	XDQ	10	RDT	10
RDZ	14	RDR	12	ZDQ	11	RDH	11	MDQ	10	RQQ	10	RDY	10
KDQ	13	RDK	12	RNQ	11	RDM	11	SDQ	10	RSQ	10	RDX	10
RDE	13	NDQ	11	RZQ	11	RDS	11	TDQ	10	RDA	10	DDQ	9

Extension using neighborhood words greater than neighborhood score threshold ($T = 11$)

Query: 1 TL⁺SHAWRLSN⁺ETDKRPF⁺IETAERL⁺RDQ⁺HK⁺KD⁺Y⁺PEY⁺KY⁺Q⁺PRR⁺RK⁺NG⁺K⁺PG⁺SS⁺EAD⁺AH⁺SE 58
 TL⁺ WRL⁺ N⁺ +KRPF+E⁺ AERLR+OHK⁺KD+P+YKYQ⁺PRR⁺K+ K⁺ G⁺ S⁺ D⁺ +
 Sbjct: 140 TLES⁺GWLEN⁺PGE⁺KRPFV⁺GEAERL⁺RBQ⁺HK⁺KD⁺HPD⁺YKYQ⁺PRR⁺RK⁺V⁺KNG⁺QSE⁺PED⁺GSE⁺BQ 197

```
>|_gb[AA08419.1] PTEN [Takifugu rubripes]
Length=412

Score = 197 bits (501), Expect = 2e-49, Method: Composition-based stats.
Identities = 95/100 (95%), Positives = 98/100 (98%), Gaps = 0/100 (0%)

Query 2 IVSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI 61
+VSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI
Sbjct 8 IVSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI 67

Query 62 YNLCAERHYDTAK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FPK+ON 101
YNLCAERHYD+AK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FPK+ON ++
Sbjct 68 YNLCAERHYDAK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FP+CD 107

Score = 83.6 bits (205), Expect = 4e-15, Method: Composition-based stats.
Identities = 60/103 (58%), Positives = 68/103 (66%), Gaps = 32/103 (31%)

Query 99 EQNMLKEDKMF+FW+NTFF+IPG+EEV+-----D 126
+
Sbjct 260 EQNMLKEDKMF+FW+NTFF+IPG+EEV+-----D 126
EQNMLKEDKMF+FW+NTFF+IPG+EEV+-----D 126

Query 127 NDRKYLVL+LT+PK+sd+l+sk+sk+sk+sk+RY+FS+FN+KV+LY+TR+TV+EE 169
+D+++EL-LTL+RMD+ DRANK+DRARY+FS+FN+V+EL+ F+RYVE
Sbjct 320 SDRD+YL+L+L+SK+NR+DR+ANK+DR+ARY+FS+FN+V+EL+ F+RYVE 362

>|_gb[AA03110.1] UG pten protein [Danio rerio]
Length=289

Score = 197 bits (500), Expect = 2e-49, Method: Composition-based stats.
Identities = 95/99 (95%), Positives = 98/99 (98%), Gaps = 0/99 (0%)

Query 3 VSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI 62
VSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI
Sbjct 9 VSRNKRRYQEDGF+DLDTIY+PN+IANG+FAER+LEGV+RNNID+VV+FLDS+KHN+HYKI 68

Query 63 YNLCAERHYDTAK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FPK+ON 101
YNLCAERHYD+AK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FPK+ON ++
Sbjct 69 YNLCAERHYDTAK+FR+CRVAQ+YFFED+HN+FPQ+LEL+K+FP+CD 107
```

STEP 1**Remove low-complexity region or sequence repeats in the query sequence**

- "Low-complexity region" means a region of a sequence composed of few kinds of elements
- The regions will be marked with an X (protein sequences) or N (nucleic acid sequences) and then be ignored by the BLAST program

Low-complexity sequence can often be recognized by visual inspection.

For example,

Protein sequence

PPCDPPPPPKDKKKKDDGPP

Nucleotide sequence

AAATAAAAAAAAAATAAAAAAT.

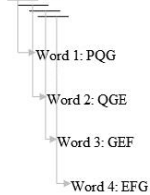
To filter out the low-complexity regions,
SEG program is used for protein sequences
DUST program is used for DNA sequences

STEP 2**Make a k-letter word list of the query sequence.**

Take $k=3$ for example, we list the words of length 3 in the query protein sequence

(k is usually 11 for a DNA sequence)

Query sequence: PQGEFG



STEP 3

List the possible matching words.

- BLAST only cares about the high-scoring words.
- Example: the score obtained by comparing PQG with PEG and PQA is respectively 15 and 12
- For DNA words, a match is scored as +5 and a mismatch as -4, or as +2 and -3

STEP 4

Organize the remaining high-scoring words into an efficient search tree.

- This allows the program to rapidly compare the high-scoring words to the database sequences.

STEP 5

Repeat step 3 to 4 for each k-letter word in the query sequence

STEP 6

Scan the database sequences for exact matches with the remaining high-scoring words.

The BLAST program scans the database sequences for the remaining high-scoring word,

STEP 7

Extend the exact matches to high-scoring segment pair (HSP).

- The original version of BLAST stretches a longer alignment between the query and the database sequence in the left and right directions

```

Query sequence: R P P Q G L F
Database sequence: D P P E G V V
                    |
                    |→Exact match is scanned
Score: -2 7 7 2 6 1 -1
                    |
                    |→HSP
Optimal accumulated score = 7+7+2+6+1 = 23

```

STEP 8

List all of the HSPs in the database whose score is high enough to be considered.

- List the HSPs whose scores are greater than the empirically determined **cutoff score S**.
- By examining the distribution of the alignment scores modeled by comparing random sequences, a cutoff score S can be determined such that its value is large enough to guarantee the significance of the remaining HSPs

STEP 9

Evaluate the significance of the HSP score.