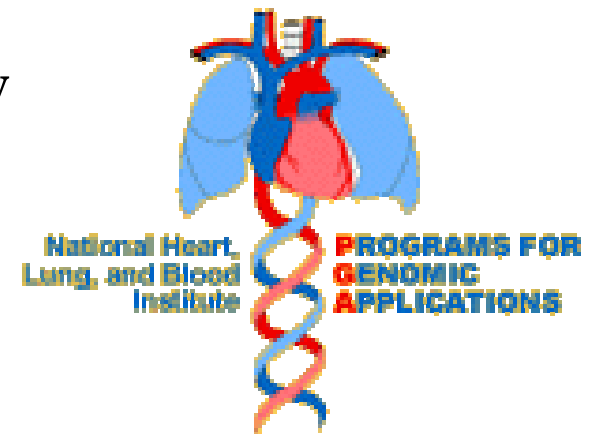


# BLAST:

## Basic Local Alignment Search Tool

Jonathan M. Urbach  
Bioinformatics Group  
Department of Molecular Biology



# Topics to be covered:

- " BLAST as a Sequence Alignment Tool
- " Uses of BLAST
- " Types of BLAST
- " How BLAST works
  - Scanning for 'hits'
  - Scoring with Substitution Matrices
- " Common Databases for Use with BLAST available at NCBI
- " Interpretation of Blast Results
- " Blast options: on the net or on **your** computer
- " Learning More About BLAST,
- " A BLAST demo

[gi|13325078|gb|AAG33875.2](#) (AF232004) HrpL [*Pseudomonas syringae* pv. tomato]

Length = 184

Score = 347 bits (889), Expect = 4e-95

Identities = 182/184 (98%), Positives = 183/184 (98%)

Query: 1 MFQKIVILDSTQPRQPSSSAGIRQMTADQIQMLRAFIQKRVMNPDDVDDILQCVFLEALR 60

MFQKIVILDSTQPRQPSSSAGIRQMTADQIQMLRAFIQKRVMNPDDVDDILQCVFLEALR

Sbjct: 1 MFQKIVILDSTQPRQPSSSAGIRQMTADQIQMLRAFIQKRVMNPDDVDDILQCVFLEALR 60

Query: 61 NEHKFQHASKPQTWLCGIALNLIRNHFRKMYRQPYQESWEDEVHSELEGHGVDVSHQVDGH

NEHKFQHASKPQTWLCGIALNLIRNHFRKMYRQPYQESWEDEVHSELEGHGVDVSHQV+GH

Sbjct: 61 NEHKFQHASKPQTWLCGIALNLIRNHFRKMYRQPYQESWEDEVHSELEGHGVDVSHQVEGH 121

Query: 121 RQLARVIQAIDCLPSNMQKVLEVSLEMDGNYQETANSLGVPIGTVRSRLSRARVQLKQQI 180

RQLARVIQAIDCLPSNMQKVLEVSLEMDGNYQETANSLGVPIGTVRSRLS ARVQLKQQI

Sbjct: 121 RQLARVIQAIDCLPSNMQKVLEVSLEMDGNYQETANSLGVPIGTVRSRLSGARVQLKQQI 180

Query: 181 DPFA 184

DPFA

Sbjct: 181 DPFA 184

ClustalX (1.81)

File Edit Alignment Trees Colors Quality Help

Multiple Alignment Mode Font Size: 14

|    |        |         |             |       |       |                |     |         |      |        |       |          |           |          |     |
|----|--------|---------|-------------|-------|-------|----------------|-----|---------|------|--------|-------|----------|-----------|----------|-----|
| 1  | 100285 | RARCQS  | ----        | SG    | RG    | SATGASHSGH     | LP  | AVEHAKS | AGS  | VAGDGR | LSGN  | SEQPRRAD | -----     |          |     |
| 2  | 100287 | RARCQS  | ----        | SG    | RR    | SATGASHSGH     | LP  | AVEHAKS | AGS  | VAGDGR | LSGN  | SEQPRRAD | -----     |          |     |
| 3  | 067651 | EWRVHS  | ----        | SG    | RR    | CAIGTTRHRSH    | LL  | AVEHAKK | GPGS | VAGDGW | LSGY  | GEHARRSD | WHCALATVP |          |     |
| 4  | 100279 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 5  | 100281 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 6  | 100277 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 7  | 100275 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 8  | 864041 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 9  | 864039 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 10 | 100283 | ERRYHP  | ----        | SG    | RR    | PPTAGTTRHRGHRL | PA  | DEHAEG  | SGS  | FSGNGW | LLSGN | RQHARCSY | -----     |          |     |
| 11 | 100273 | ERRYHP  | ----        | SG    | RR    | PTAGTTRHRGHRL  | PA  | DKHAEG  | SGS  | FSGNGW | LSGN  | RQHARCSY | -----     |          |     |
| 12 | 864037 | ERRYHP  | ----        | SG    | RR    | SPTAGTTRHRGHRL | PA  | DKHAEG  | SGS  | FSGNGR | LSGN  | RQHARCS  | D-----    |          |     |
| 13 | 360504 | EDEVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOK | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 14 | 360496 | EDEVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOK | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 15 | 360480 | EDEVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOK | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 16 | 360392 | EDEVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOK | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 17 | 360280 | EDEVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOK | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 18 | 360360 | EDDVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOR | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 19 | 360352 | EDDVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOR | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 20 | 360336 | EDDVHSE | ELEGHGDVSHQ | VDGHR | LARV  | IQ             | AID | CLPS    | NMOR | VLEVS  | LEMDG | NYQETAN  | SLGVPIGT  | VRS      |     |
| 21 | 360416 | EDDVHSE | ELEWNGD     | ITHQ  | VDGHR | LARV           | IA  | AID     | CLPS | NMOK   | VLEVS | LEMDG    | NYQDTANT  | ILGVPIGT | VRS |

180.....190.....200.....210.....220.....230.....240.....

File /tmp/blastres\_parsed\_PROT.aln loaded.

# Sequence Alignment Tools

## Database Searching:

### **BLAST:**

NCBI, Web Interface: <http://www.ncbi.nlm.nih.gov/BLAST/>

WuBLAST <http://blast.wustl.edu>

**FASTA:** <http://www.ebi.ac.uk/fasta3/>

### **Smith-Waterman**

Par-Align: <http://dna.uio.no/search/>

## Multiple Sequence Alignment:

**CLUSTALW:** <http://www-igbmc.u-strasbg.fr/BioInfo/ClustalX/Top.html>

**DiAlign**, Web Interface: <http://genomatix.gsf.de/cgi-bin/dialign/dialign.pl>

**MSA:** <http://www.ncbi.nlm.nih.gov/CBBresearch/Schaffer/msa.html>

Web Interface: <http://bioweb.pasteur.fr/seqanal/interfaces/msa-simple.html>

# **Uses of BLAST:**

**Query a database for sequences similar to an input sequence.**

# Uses of BLAST:

**Query a database for sequences similar to an input sequence.**



**Identify previously characterized sequences.**

# Uses of BLAST:

**Query a database for sequences similar to an input sequence.**



- „ Identify previously characterized sequences.**
- „ Find phylogenetically related sequences.**



# Uses of BLAST:

**Query a database for sequences similar to an input sequence.**

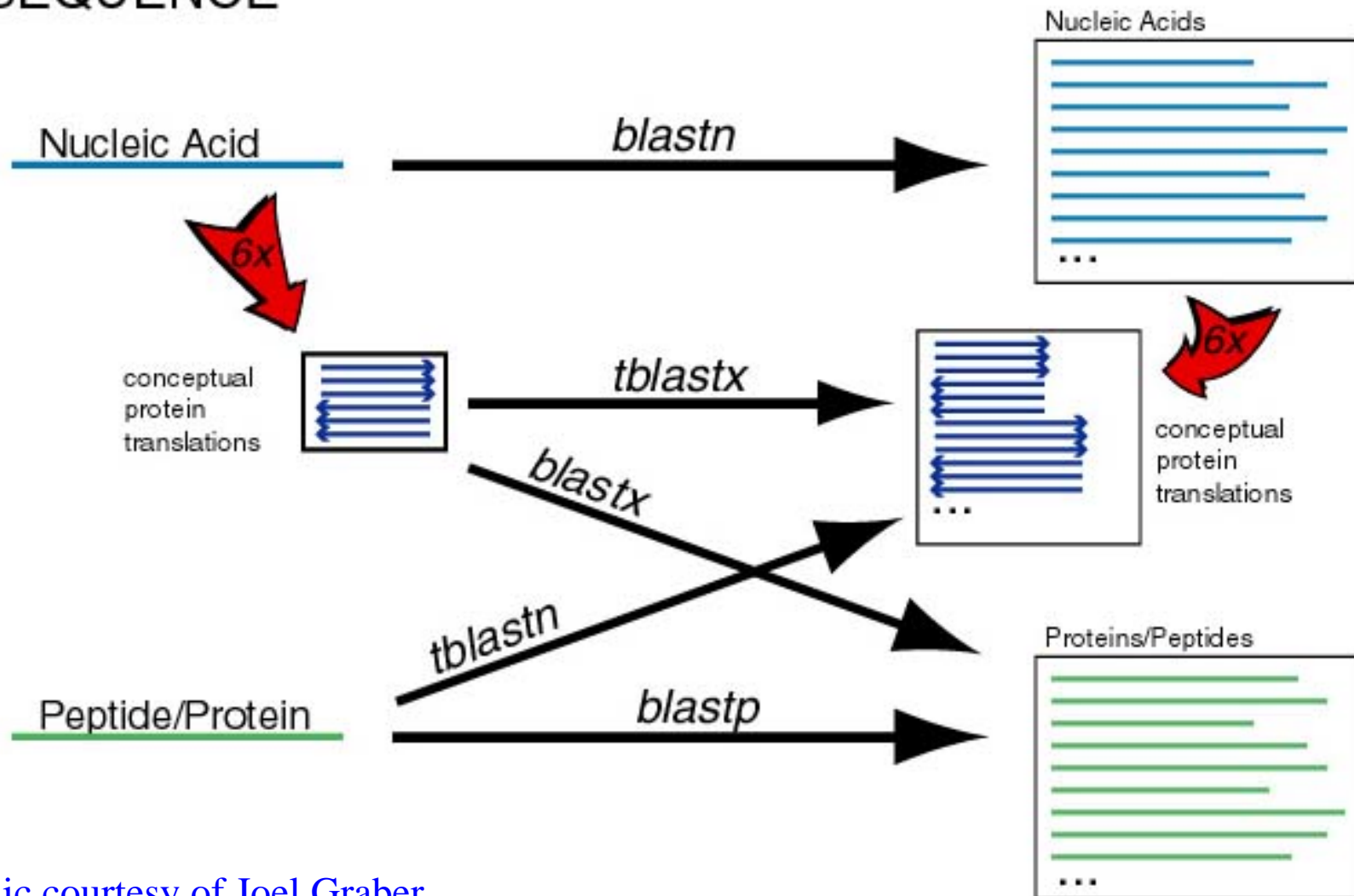


- „ Identify previously characterized sequences.**
- „ Find phylogenetically related sequences.**
- „ Identify possible functions based on similarities to known sequences.**

# Types of BLAST:

QUERY  
SEQUENCE

DATABASE



# How BLAST Works

- (1) BLAST scans database for 'words' of a predetermined length (a 'hit') with some minimum threshold parameter,  $T$ .**
- (2) BLAST then extends the hit until the score falls below the maximum score yet attained minus some value  $X$ .**

**Query:**

**MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT**

Query:

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT

▼ Use 2 or 3-letter words...

~~MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT~~

~~MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT~~

~~MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT~~

**Query:**

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



**Scan against subject sequence:**

>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPSVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQSLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA

Query:

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



A hit!

>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQSLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA

**Query:**

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPYFPGALGDEKTTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTT



>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQSLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA



**Extension:**

**Query:**

**YFP**

**Y P**

**Sbjct:**

**YLP**



**Query:**

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPSVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQSLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA



**Extension:**

VNNSYFPGAL

V Y P L

VQEQYLPSVL

**Query:**

**Sbjct:**

**Query:**

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHVVEWWGGEEARPTLADVQEQYLPSVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQSLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA



**Extension:**

**Query:**

VSMWSAESCRTPLCSVNNSYFPGALGDEKTT

V W E R L V Y P L E T

**Sbjct:**

VEWWGGEEARPTLADVQEQYLPSVLAQESVT

**Query:**

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT

MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTT



>gi|507311|gb|AAA25685.1| aminoglycoside 6'-N-acetyltransferase

MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPSVLAQESVTPYIAMLNGEPIG

SGDGWWEETDPGVRGIDQLANASQLGKGLGTKLVRALVELLFNDPEVTKIQTDPSPSNLR

GFERQGTVTTDPGPAVYMVQTRQAFERTRSDA



**Extension:**

**Query: MVSHSAAQAYSMLTNSEFVSMWSAESCRTPLCSVNNSYFPGALGDEKTTKVI**

**M H A Y L S V W E R L V Y P L E T I L**

**Sbjct: MTEHDLAMLYEWLNRSHIVEWWGGEEARPTLADVQEQYLPSVLAQESVTPY**

**HSP: A High-Scoring Segment Pair**

# **Towards BLAST Scoring**

- „ **Expected negative score for alignment of two random residues.**
- „ **Maximal score for a perfect match.**
- „ **Combinations of residues that can commonly substitute for one another in proteins may have positive score.**

```

# Matrix made by matblas from blosum62.ij
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units
# Blocks Database = /data/blocks_5.0/blocks.dat
# Cluster Percentage: >= 62
# Entropy = 0.6979, Expected = -0.5209
  A R N D C Q E G H I L K M F P S T W Y V B Z X *
A 4 -1 -2 -2 0 -1 -1 0 -2 -1 -1 -1 -1 -2 -1 1 0 -3 -2 0 -2 -1 0 -4
R -1 5 0 -2 -3 1 0 -2 0 -3 -2 2 -1 -3 -2 -1 -1 -3 -2 -3 -1 0 -1 -4
N -2 0 6 1 -3 0 0 0 1 -3 -3 0 -2 -3 -2 1 0 -4 -2 -3 3 0 -1 -4
D -2 -2 1 6 -3 0 2 -1 -1 -3 -4 -1 -3 -3 -1 0 -1 -4 -3 -3 4 1 -1 -4
C 0 -3 -3 -3 9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
Q -1 1 0 0 -3 5 2 -2 0 -3 -2 1 0 -3 -1 0 -1 -2 -1 -2 0 3 -1 -4
E -1 0 0 2 -4 2 5 -2 0 -3 -3 1 -2 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
G 0 -2 0 -1 -3 -2 -2 6 -2 -4 -4 -2 -3 -3 -2 0 -2 -2 -3 -3 -1 -2 -1 -4
H -2 0 1 -1 -3 0 0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2 -3 0 0 -1 -4
I -1 -3 -3 -3 -1 -3 -3 -4 -3 4 2 -3 1 0 -3 -2 -1 -3 -1 3 -3 -3 -1 -4
L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -2 -1 1 -4 -3 -1 -4
K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 5 -1 -3 -1 0 -1 -3 -2 -2 0 1 -1 -4
M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5 0 -2 -1 -1 -1 -1 1 -3 -1 -1 -4
F -2 -3 -3 -3 -2 -3 -3 -3 -1 0 0 -3 0 6 -4 -2 -2 1 3 -1 -3 -3 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 7 -1 -1 -4 -3 -2 -2 -1 -2 -4
S 1 -1 1 0 -1 0 0 0 -1 -2 -2 0 -1 -2 -1 4 1 -3 -2 -2 0 0 0 -4
T 0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 5 -2 -2 0 -1 -1 0 -4
W -3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1 1 -4 -3 -2 11 2 -3 -4 -3 -2 -4
Y -2 -2 -2 -3 -2 -1 -2 -3 2 -1 -1 -2 -1 3 -3 -2 -2 2 7 -1 -3 -2 -1 -4
V 0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4 -3 -2 -1 -4
B -2 -1 3 4 -3 0 1 -1 0 -3 -4 0 -3 -3 -2 0 -1 -4 -3 -3 4 1 -1 -4
Z -1 0 0 1 -3 3 4 -2 0 -3 -3 1 -1 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4

```

# BLAST Scoring

• **Nominal HSP scores ( $S$ ) are sums of scores from substitution matrices.**

• **Nominal scores are normalized to give 'bit scores' ( $S'$ ):**

$$(I) \quad S' = \frac{\lambda S B \ln K}{\ln 2}$$

*K and  $\lambda$  are statistical parameters that relate the calculated score to the probability finding a hit with at least that score.*

**1 Allows comparison of alignments scored by different methods**

# Relating Scores to Probability: E-Values and P-Values

*The expected number of HSPs with scores of at least  $S$  is given by the following equations:*

$$(II) \quad E = Kmn e^{B \leftarrow S}$$

$$(III) \quad E = mn 2^{BS'}$$

$K$  and  $\lambda$  are statistical  
'normalization' parameters.  
 $m$  and  $n$  are the lengths of the  
query  
sequence and database.  
 $S$  is a nominal score.

$S'$  is a bit score.

P-Values are the likelihood of finding a match with a score of at least  $S$ :

$$(IV) \quad P = 1B e^{BE}$$

# Substitution Matrices

**Scores in the substitution matrix are expressed in 'log-odds' format:**

$$(V) \quad s_{ij} = \ln \left( \frac{q_{ij}}{p_i p_j} \right) / \lambda$$

$q_{ij}$  = target frequency  
 $p_i, p_j$  = frequency those residues appear by chance  
 $\lambda$  = normalization parameter

- 1 The more frequently the substitution occurs, the higher the score.**
- 1 The less frequently the residue occurs in the sequence as a whole, the higher the score.**



# Substitution Matrices

- „ Derived from empirically observed substitution frequencies
- „ Higher scores for substitution with similar residues.
- „ Random substitutions give negative scores

# Types of Substitution Matrices

Each tailored to a specific degree of evolutionary divergence.

PAM Matrices:

'Percent Accepted Mutation'

start with closely related sequences, and extrapolate substitution probabilities for more distantly related sequences.

1 PAM unit=1 mutation event per 100 bases.

e.g.: PAM 100 tailored for 100 mutation events per 100 bases.

*Barker, W.C. & Dayhoff, M.O. Atlas of Protein Sequence and Structure, pp 101-110, National Biomedical Research Foundation (1972).*

# Types of Substitution Matrices

## BLOSUM Matrices

'BLOck SUbstitution Matrix'

Values inferred from sequences sharing a maximum of the given value.

e.g.: BLOSUM62 derived from sequences no more than 62% identical.

Henikoff, S. & Henikoff, J.G., *Proc. Natl. Acad. Sci., USA*,  
89, 10915-10919 (1992).

# Comparing Substitution Matrices

Similar Evolutionary Distances

PAM 120 <-----> BLOSUM80

PAM160 <-----> BLOSUM62

PAM250 <-----> BLOSUM45

BLOSUM more tolerant to hydrophobic than PAM  
but less tolerant to hydrophilic substitutions.

```

# Matrix made by matblas from blosum62.ij
# * column uses minimum score
# BLOSUM Clustered Scoring Matrix in 1/2 Bit Units
# Blocks Database = /data/blocks_5.0/blocks.dat
# Cluster Percentage: >= 62
# Entropy = 0.6979, Expected = -0.5209
  A R N D C Q E G H I L K M F P S T W Y V B Z X *
A 4 -1 -2 -2 0 -1 -1 0 -2 -1 -1 -1 -1 -2 -1 1 0 -3 -2 0 -2 -1 0 -4
R -1 5 0 -2 -3 1 0 -2 0 -3 -2 2 -1 -3 -2 -1 -1 -3 -2 -3 -1 0 -1 -4
N -2 0 6 1 -3 0 0 0 1 -3 -3 0 -2 -3 -2 1 0 -4 -2 -3 3 0 -1 -4
D -2 -2 1 6 -3 0 2 -1 -1 -3 -4 -1 -3 -3 -1 0 -1 -4 -3 -3 4 1 -1 -4
C 0 -3 -3 -3 9 -3 -4 -3 -3 -1 -1 -3 -1 -2 -3 -1 -1 -2 -2 -1 -3 -3 -2 -4
Q -1 1 0 0 -3 5 2 -2 0 -3 -2 1 0 -3 -1 0 -1 -2 -1 -2 0 3 -1 -4
E -1 0 0 2 -4 2 5 -2 0 -3 -3 1 -2 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
G 0 -2 0 -1 -3 -2 -2 6 -2 -4 -4 -2 -3 -3 -2 0 -2 -2 -3 -3 -1 -2 -1 -4
H -2 0 1 -1 -3 0 0 -2 8 -3 -3 -1 -2 -1 -2 -1 -2 -2 2 -3 0 0 -1 -4
I -1 -3 -3 -3 -1 -3 -3 -4 -3 4 2 -3 1 0 -3 -2 -1 -3 -1 3 -3 -3 -1 -4
L -1 -2 -3 -4 -1 -2 -3 -4 -3 2 4 -2 2 0 -3 -2 -1 -2 -1 1 -4 -3 -1 -4
K -1 2 0 -1 -3 1 1 -2 -1 -3 -2 5 -1 -3 -1 0 -1 -3 -2 -2 0 1 -1 -4
M -1 -1 -2 -3 -1 0 -2 -3 -2 1 2 -1 5 0 -2 -1 -1 -1 -1 1 -3 -1 -1 -4
F -2 -3 -3 -3 -2 -3 -3 -3 -1 0 0 -3 0 6 -4 -2 -2 1 3 -1 -3 -3 -1 -4
P -1 -2 -2 -1 -3 -1 -1 -2 -2 -3 -3 -1 -2 -4 7 -1 -1 -4 -3 -2 -2 -1 -2 -4
S 1 -1 1 0 -1 0 0 0 -1 -2 -2 0 -1 -2 -1 4 1 -3 -2 -2 0 0 0 -4
T 0 -1 0 -1 -1 -1 -1 -2 -2 -1 -1 -1 -1 -2 -1 1 5 -2 -2 0 -1 -1 0 -4
W -3 -3 -4 -4 -2 -2 -3 -2 -2 -3 -2 -3 -1 1 -4 -3 -2 11 2 -3 -4 -3 -2 -4
Y -2 -2 -2 -3 -2 -1 -2 -3 2 -1 -1 -2 -1 3 -3 -2 -2 2 7 -1 -3 -2 -1 -4
V 0 -3 -3 -3 -1 -2 -2 -3 -3 3 1 -2 1 -1 -2 -2 0 -3 -1 4 -3 -2 -1 -4
B -2 -1 3 4 -3 0 1 -1 0 -3 -4 0 -3 -3 -2 0 -1 -4 -3 -3 4 1 -1 -4
Z -1 0 0 1 -3 3 4 -2 0 -3 -3 1 -1 -3 -1 0 -1 -3 -2 -2 1 4 -1 -4
X 0 -1 -1 -1 -2 -1 -1 -1 -1 -1 -1 -1 -1 -2 0 0 -2 -1 -1 -1 -1 -1 -4

```

# BLAST Databases at NCBI:

| Database            | Description   | DNA | Protein |
|---------------------|---|-----|---------|
| nr                  | All non-redundant GenBank+EMBL+DDBJ+PDB sequences (but no EST, STS, GSS, or HTGS sequences).  | 4   | 4       |
| month               | All new or revised GenBank+EMBL+DDBJ+PDB sequences released in the last 30 days.  | 4   | 4       |
| dbest               | Non-redundant database of GenBank+EMBL+DDBJ EST Divisions.  | 4   |         |
| dbsts               | Non-redundant database of GenBank+EMBL+DDBJ STS Divisions.  | 4   |         |
| mouse ests          | The non-redundant Database of GenBank+EMBL+DDBJ EST Divisions limited to the organism mouse.  | 4   |         |
| human ests          | The Non-redundant Database of GenBank+EMBL+DDBJ EST Divisions limited to the organism human.  | 4   |         |
| other ests          | The non-redundant database of GenBank+EMBL+DDBJ EST Divisions all organisms except mouse and human.   | 4   |         |
| yeast               | Yeast ( <i>Saccharomyces cerevisiae</i> ) genomic nucleotide sequences. Not a collection of all Yeast nucleotide sequences, but the sequence fragments from the Yeast complete genome.  | 4   | 4       |
| E. coli             | E. coli ( <i>Escherichia coli</i> ) genomic nucleotide sequences.   | 4   | 4       |
| pdb                 | Sequences derived from the 3-dimensional structure of proteins.   | 4   | 4       |
| kabat<br>[kabatnuc] | Kabat's database of sequences of immunological interest. For more information <a href="http://immuno.bme.nwu.edu/">http://immuno.bme.nwu.edu/</a>   | 4   | 4       |
| patents             | Nucleotide sequences derived from the Patent division of GenBank.   | 4   | 4       |
| vector              | Vector subset of GenBank(R), NCBI, ( <a href="ftp://ncbi.nlm.nih.gov/pub/blast/db/">ftp://ncbi.nlm.nih.gov/pub/blast/db/</a> directory).  | 4   |         |
| swiss prot          | The last major release of the SWISS-PROT protein sequence database (no updates). These are uploaded to our system when they are received from EMBL.   |     | 4       |
| alu                 | Translations of select Alu repeats from REPBASE, suitable for masking Alu repeats from query sequences. It is available at <a href="ftp://ncbi.nlm.nih.gov/pub/jmc/alu">ftp://ncbi.nlm.nih.gov/pub/jmc/alu</a> . See "Alu alert" by Claverie and Makalowski, Nature vol. 371, page 752 (1994) . |     | 4       |

Adapted From:[http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/query\\_tutorial.html](http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/query_tutorial.html)

# Interpreting Blast Results

>gi|6580755|gb|AAF18265.1|U22895\_1 (U22895) alternative sigma factor AlgU [Azotobacter vinelandii]  
Length = 193

Score = 334 bits (857), Expect = 2e-91  
Identities = 180/192 (93%), Positives = 189/192 (97%)

Query: 1  
MLTQEQDQQLVERVQRGDKRAFDLLVLKYQHKILGLIVRFVHDAQEAQDVAQEAFIKAYR 60  
ML QEQDQQLVERVQRGD+RAFDLLVLKYQHKILGLIVRFVHDA EAQDVAQEAFIKAYR  
Sbjct: 1 MLNQEQDQQLVERVQRGDRRAFDLLVLKYQHKILGLIVRFVHDAHEAQDVAQEAFIKAYR  
60

Query: 61 ALGNFRGDSAFYTWLYRIANTAKNHLVARGRRPPDSVTAEDAEEFFEGDHALKDIESPE  
120  
ALGNFRGDSAFYTWLYRIANTAKNHLVARGRRPPDSV+A DAEF+EGDHALKDIESPE  
Sbjct: 61 ALGNFRGDSAFYTWLYRIANTAKNHLVARGRRPPDSVSAAGDAEFYEGDHALKDIESPE  
120

Query: 121 RAMLRDEIEATVHQTIQQLPEDLRTALTLREFEGLSYEDIA TVMQCPVGTVRSRIFRARE  
180  
R++LRDEIEATVH+TIQQLPEDLRTALTLREF+GLSYEDIA+VMQCPVGTVRSRIFRARE  
Sbjct: 121 RSLLRDEIEATVHRTIQQLPEDLRTALTLREFDGLSYEDIASVMQCPVGTVRSRIFRARE 180

Query: 181 AIDKALQPLLRE 192  
AIDKALQPLL+E

# Interpreting Blast Results

Hit name

>gi|6580755|gb|AAF18265.1|U22895\_1 (U22895) alternative sigma factor AlgU [Azotobacter vinelandii]  
Length = 193

Score = 334 bits (857), Expect = 2e-91  
Identities = 180/192 (93%), Positives = 189/192 (97%)

Query: 1  
MLTQEQQQLVERVQRGDKRAFDLLVLKYQHKILGLIVRFVHDAQEAQDVAQEAFIKAYR 60  
ML QEQQDQQLVERVQRGD+RAFDLLVLKYQHKILGLIVRFVHDA EAQDVAQEAFIKAYR  
Sbjct: 1 MLNQEQQDQQLVERVQRGDRRAFDLLVLKYQHKILGLIVRFVHDAHEAQDVAQEAFIKAYR  
60

Query: 61 ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSDVTAEDAEEFEGDHALKDIESPE  
120  
ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSDV+A DAEF+EGDHALKDIESPE  
Sbjct: 61 ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSDVSAGDAEFYEGDHALKDIESPE  
120

Query: 121 RAMLRDEIEATVHQTIIQQLPED VMQCPVGTVRSRIFRARE  
180  
R++LRDEIEATVH+TIQQLPEDLRTALTLKEF+GLSYEDIA+VMQCPVGTVRSRIFRARE  
Sbjct: 121 RSLLRDEIEATVHRTIIQQLPEDLRTALTLREFDGLSYEDIASVMQCPVGTVRSRIFRARE 180

Alignment with query  
sequence

Query: 181 AIDKALQPLLRE 192  
AIDKALQPLL+E



# Interpreting Blast Results

Normalized bit scores
Nominal HSP scores
Expectation value
for AlgU [Azotobacter vinelandii]

Score = 334 bits (857), Expect = 2e-91  
Identities = 180/192 (93%), Positives = 189/192 (97%)

**Query: 1**  
 MLTQEQDQOLVERVORCDKRAFDLLVLDAQEAQDVAQEAFIKAYR 60  
 RGD+RAFDL FVHDA EAQDVAQEAFIKAYR  
**Sbjct:** RVRQGRDRA VRFVHDAHEAQDVAQEAFIKAYR  
 60

Number of Identities
Number of Identities

**Query: 61** ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSVTAEDAEFFEGDHALKDIESPE 120

ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSV+A DAEF+EGDHALKDIESPE

**Sbjct: 61** ALGNFRGDSAFYTWLYRIAINAKNHLVARGRRPPDSVSAAGDAEFYEGDHALKDIESPE 120

**Query: 121** RAMLRDEIEATVHQTIQQLPEDLRTALTLREFEGLSYEDIA+VMQCPVGTVRSRIFRARE 180

R++LRDEIEATVH+TIQQLPEDLRTALTLREF+GLSYEDIA+VMQCPVGTVRSRIFRARE

**Sbjct: 121** RSLLRDEIEATVHRTIQQLPEDLRTALTLREFDGLSYEDIASVMQCPVGTVRSRIFRARE 180

**Query: 181** AIDKALQPLLRE 192  
AIDKALQPLL+E

# **BLAST: On the Net, and On Your Computer**

Advantages/Disadvantages of Net Based Blast:

- (1) Use databases hosted remotely at NCBI.
- (2) Little/No setup required.
- (3) But, Cannot use a customized database.

Advantages/Disadvantages of Local Microcomputer-Based Blast:

- (1) Can Use a Customized Database.
- (2) Better suited to scripting / automation or when a large number of queries will be performed (UNIX).
- (3) But, Requires some setup and computer expertise.

# **BLAST: On the Net, and On Your Computer**

On the Net:

<http://www.ncbi.nlm.nih.gov/BLAST/>

On Your Computer:

UNIX/MacOS/Windows

<ftp://ncbi.nlm.nih.gov/blast/executables/>

NCBI Tools for UNIX

<ftp://ncbi.nlm.nih.gov/toolbox/>

WUBLAST

<http://blast.wustl.edu>

# Learning More about BLAST

How Blast Works:

Altschul, S.F. et al., *Nucleic Acids Research*, **25**, 3389-3402 (1997).

Scoring Schemes:

Karlin, S., and Altschul, S.F., *Proc. Natl. Acad. Sci.*, **87**, 2264-2268 (1990).

*Henikoff, S., and Henikoff, J.G., Proc. Natl. Acad. Sci.*, **89**, 10915-10919 (1992).

<http://www.ncbi.nlm.nih.gov/BLAST/tutorial/Altschul-1.html>

**Online Tutorial**

<http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/information3.html>