

Finding Function By Sequence Similarity

Concepts of Sequence Similarity Searching

- The major goal of sequence analysis is to predict the function and structure of genes and proteins from their sequence similarity.
- One sequence by itself is not informative
- Sequence must be analyzed by comparative methods against existing sequence databases to develop hypothesis concerning relatives and function.

BLAST algorithm

- **Basic Local Alignment Search Tools**
- widely used sequence similarity search tool
- set of sequence comparison algorithms used to search sequence databases
- Finds best local alignments to a query
- Heuristic approach based on Smith Waterman algorithm
- Provides statistical significance
- www, standalone and network clients

Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ (1990) "Basic local alignment search tool." J. Mol. Biol. 215:403-410.

Altschul SF, Madden TL, Schaeffer AA, Zhang J, Zhang Z, Miller W, Lipman DJ (1997) "Gapped BLAST and PSI-BLAST: a new generation of protein database search programs." NAR 25:3389-3402.

```
pgt115237380|ref|NP_197165.1| myb family transcription factor (MYB43) [Arabidopsis thaliana]
MGRQPCCKVGLKKGPMTEEDKKLINFILTNGKCRALPKL SGLLACGKSCRKWKINYLKPKGLL
SEYEEQVNLHAQLGNRWKIASLPGRTDNEKNMNTHTKXKLRKMGIDPLTKMPLSEGEASQQAQG
IKKSLVPHGKMPKQQTNDGQKHL EQALEKNTSVSGDGFCDVPLLNPHLIL101SSHHHSH
DQNVNNTSKFTSPSSSSSTSSCISVAVGQEFKFFDEMLLDLKKLSSDGLGDTISKGGKFNKSTV
DTMLWDINDLSSLQPMHEHGGFIDGNGKSRMVLQDQSWTFLL
```

Submit Query

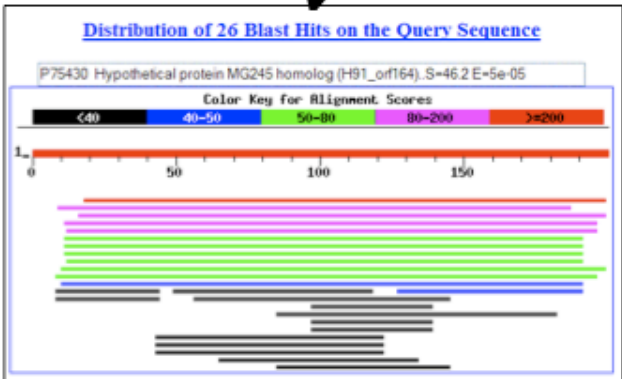


Request Results



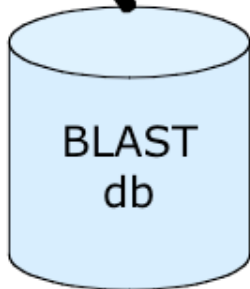
Return Formatted Results

Display Results



fetch ASN.1

fetch sequence



What BLAST tells you ...

- Assumptions
 - Random sequences
 - Constant composition
- reports surprising alignments
 - Different than chance
- Conclusions
 - Surprising similarities imply evolutionary homology

Evolutionary Homology: having a descent from a common ancestor does not always imply similar function

BLAST programs

nucleotide blast

Search a **nucleotide** database using a **nucleotide** query

Algorithms: blastn, megablast, discontinuous megablast

protein blast

Search **protein** database using a **protein** query

Algorithms: blastp, psi-blast, phi-blast

blastx

Search **protein** database using a **translated nucleotide** query

tblastn

Search **translated nucleotide** database using a **protein** query

tblastx

Search **translated nucleotide** database using a **translated nucleotide** query

BLAST Algorithm

- Scoring of matches done using scoring matrices
- Sequences are split into words (default $n=3$)
 - Speed, computational efficiency
- BLAST algorithm extends the initial “seed” hit into an HSP
 - HSP = high scoring segment pair = Local optimal alignment

the statistics

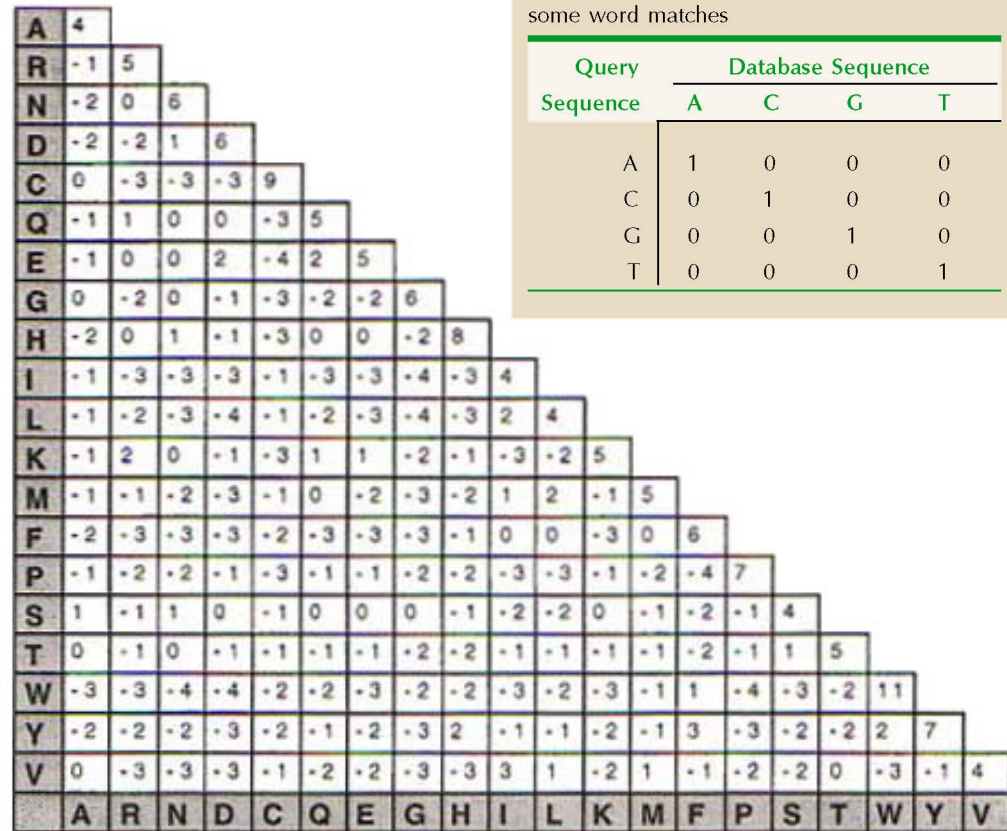
- Discriminating between real and artifactual matches is done using an estimate of probability that the match might occur by chance.
- meaning of the scores (S) and e-values (E) that are associated with BLAST hits

score (S) ?

- The quality of each pair-wise alignment is represented as a score and the scores are ranked.
- Scoring matrices are used to calculate the score of the alignment base by base (DNA) or amino acid by amino acid (protein).
- The alignment score will be the sum of the scores for each position.

What's a scoring matrix?

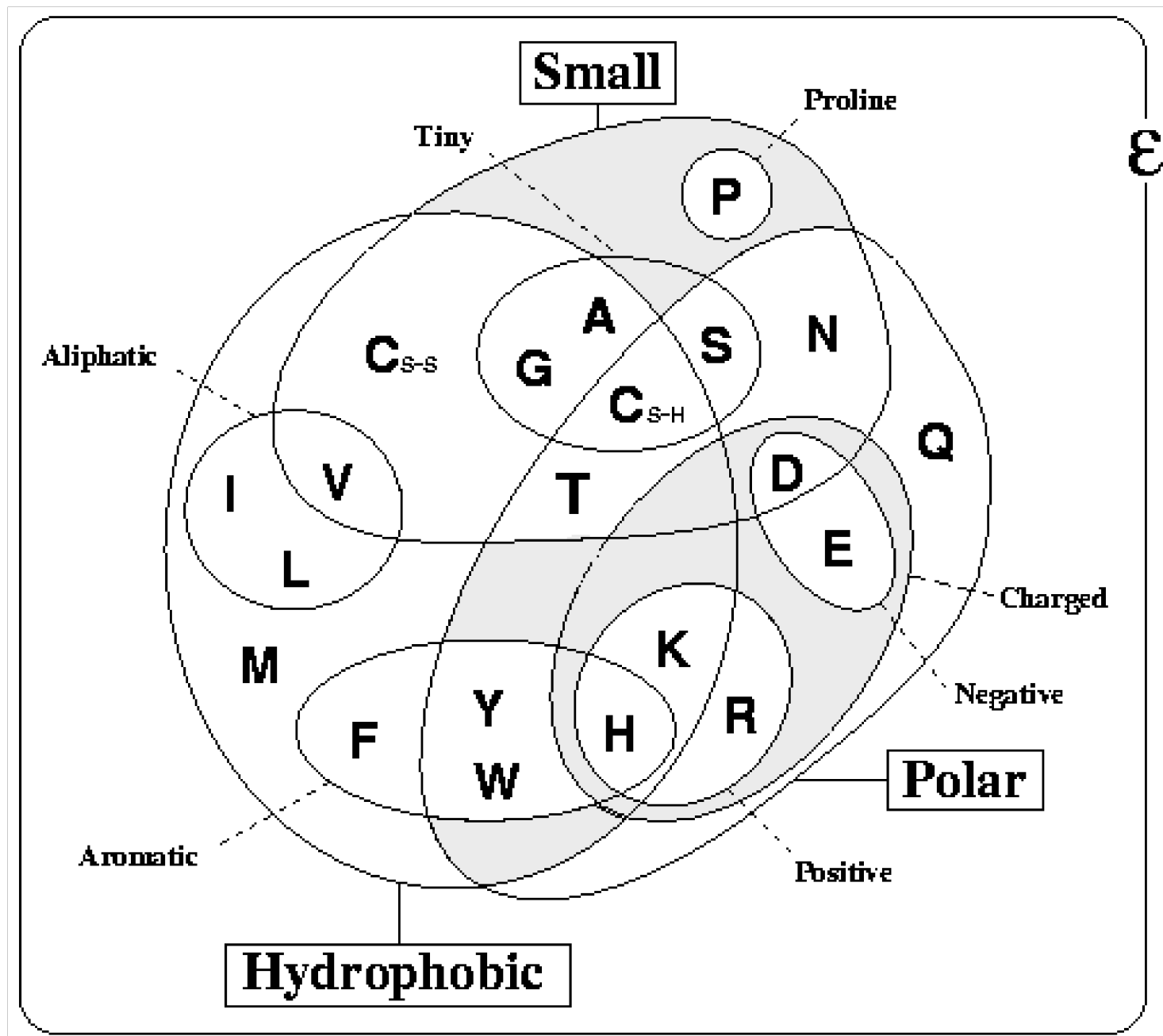
- Substitution matrices are used for amino acid alignments.
 - each possible residue substitution is given a score
- A simpler unitary matrix is used for DNA pairs (+1 for match, -2 mismatch)



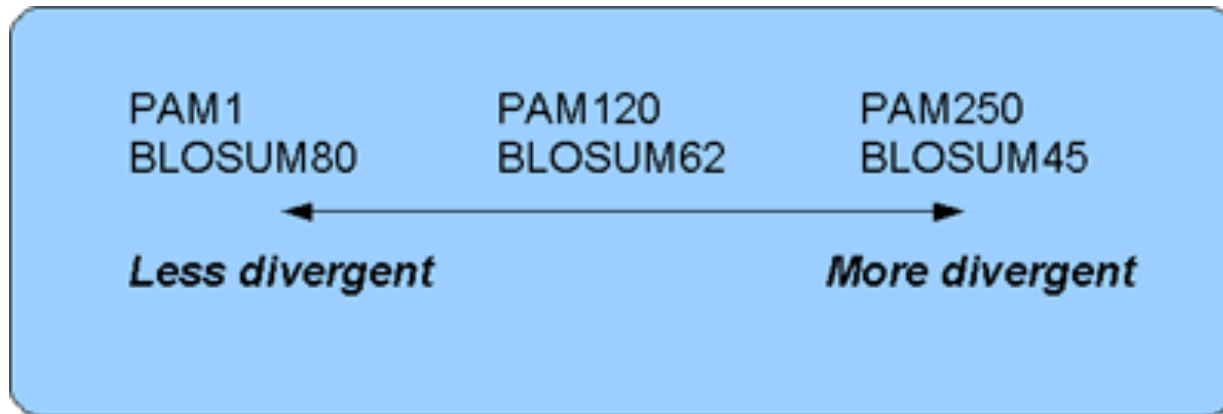
| | | | | | | | | | | | | | | | | | | | | |
|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|---|
| A | 4 | | | | | | | | | | | | | | | | | | | |
| R | -1 | 5 | | | | | | | | | | | | | | | | | | |
| N | -2 | 0 | 6 | | | | | | | | | | | | | | | | | |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

TABLE 27.3. DNA substitution matrix and some word matches

| Query Sequence | Database Sequence | | | |
|----------------|-------------------|---|---|---|
| | A | C | G | T |
| A | 1 | 0 | 0 | 0 |
| C | 0 | 1 | 0 | 0 |
| G | 0 | 0 | 1 | 0 |
| T | 0 | 0 | 0 | 1 |



BLOSUM vs PAM



- BLOSUM 62 is the default matrix in BLAST.
- This works well to identify moderately distant proteins, and performs well in detecting closer relationships.
- A search for distant relatives may be more sensitive with a different matrix.

Score and the e-value?

- The quality of the alignment is represented by the **Score (S)**.
- The score of an alignment is calculated as the sum of substitution and gap scores. Substitution scores are given by a look-up table (PAM, BLOSUM) whereas gap scores are assigned empirically .
- The significance of each alignment is computed as an **E value (E)**.
- Expectation value. The number of different alignments with scores equivalent to or better than S that are expected to occur in a database search by chance. The lower the E value, the more significant the score.

Notes on E-values

- Low E-values suggest that sequences are homologous
 - Can't show non-homology
- Statistical significance depends on both the size of the alignments and the size of the sequence database
 - Important consideration for comparing results across different searches
 - E-value increases as database gets bigger
 - E-value decreases as alignments get longer

Homology

- Similarity can be indicative of homology
 - if two sequences are significantly similar over entire length they are likely homologous
- Low complexity regions can be highly similar without being homologous
- Homologous sequences not always highly similar

BLAST Cutoffs

- nucleotide based searches
 - look for hits with E-values of 10^{-6} or less and sequence identity of 70% or more
- protein based searches
 - look for hits with E-values of 10^{-3} or less and sequence identity of 25% or more

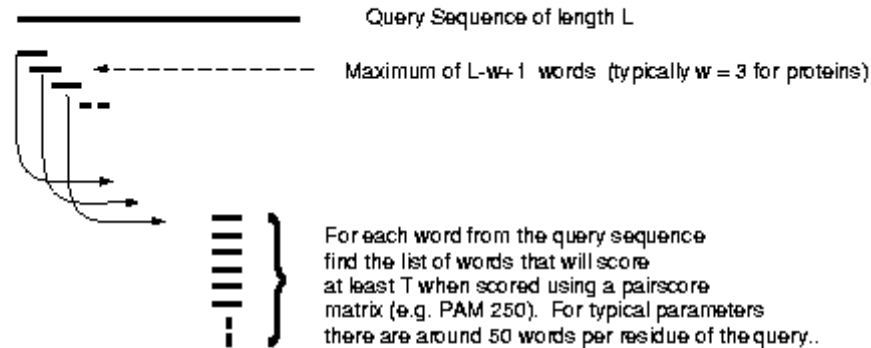
BLAST Algorithm

- Scoring of matches done using scoring matrices
- Sequences are split into words (default $n=3$)
 - Speed, computational efficiency
- BLAST algorithm extends the initial “seed” hit into an HSP
 - HSP = high scoring segment pair = Local optimal alignment

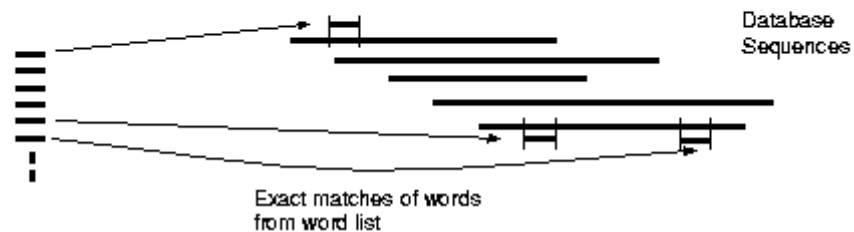
How does BLAST work?

BLAST Algorithm

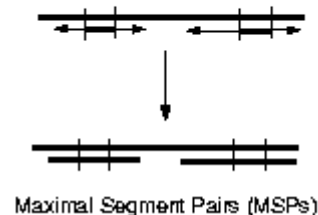
- (1) For the query find the list of high scoring words of length w .



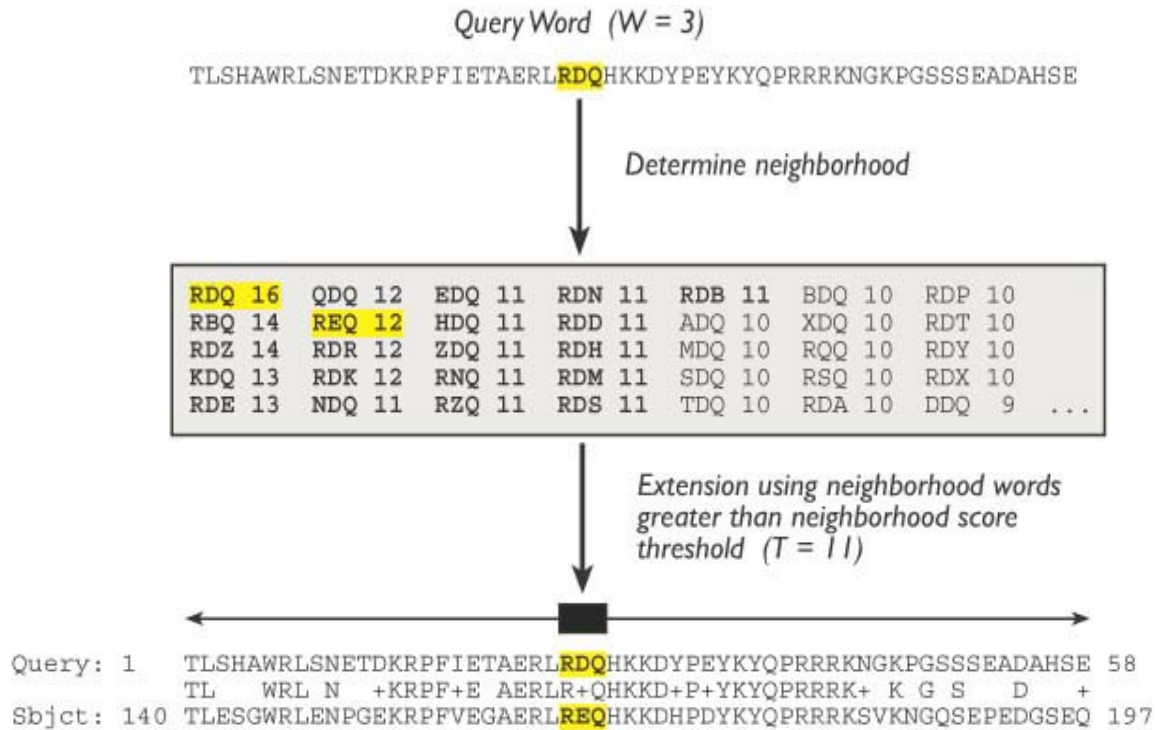
- (2) Compare the word list to the database and identify exact matches.



- (3) For each word match, extend alignment in both directions to find alignments that score greater than score threshold S .



BLAST Algorithm



>|gb|[AAL08419.1](#)| 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 IVSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI 61
+VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI
Sbjct 8 MVS RNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKI 67

Query 62 YNLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFKQN 101
YNLCAERHYD AKFNCRVAQYPPFEDHNPPQLELIKPF ++
Sbjct 68 YNLCAERHYDAAKFNCRVAQYPPFEDHNPPQLELIKPFCE D 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 KQNKMLKKDKMFHFWVNTFFIPGPEEV-----D 126
KQNK M+KKDKMFHFWVNTFFIPGPEE +
Sbjct 260 KQNKMMKKDKMFHFWVNTFFIPGPEESRDKLENGAVNNADSQQGVPAPGGQPQSAECRE 319

Query 127 NDKEYLVLTLTkndldkankdkanRYFSPNFKVKLYFTKTVEE 169
+D++YL+LTL+KND DKANKDKANRYFSPNFKVKL F+KTVEE
Sbjct 320 SDRDY LILTL SKNDRDKANKDKANRYFSPNFKVKLCFSKTVEE 362
```

>|gb|[AAH93110.1](#)| **UG** Ptenb 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 VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKKNHYKIY 62
VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHK+HYKIY
Sbjct 9 VSRNKRRYQEDGFDDLTYIYPNIIAMGFPAERLEGVYRNNIDDVVRFLDSKHKHDHYKIY 68

Query 63 NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFKQN 101
NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPF ++
Sbjct 69 NLCAERHYDTAKFNCRVAQYPPFEDHNPPQLELIKPFCE D 107
```

BLAST Algorithm

- Scoring of matches done using scoring matrices
- Sequences are split into words (default $n=3$)
 - Speed, computational efficiency
- BLAST algorithm extends the initial “seed” hit into an HSP
 - HSP = high scoring segment pair = Local optimal alignment

- Additional slides

Extending the High Scoring Segment Pair (HSP)

