**November 29, 2010**

**BLAST:**
**Basic local alignment**
**search tool**

**BLAST!**

Jonathan Pevsner, Ph.D.
Bioinformatics
pevsner@kennedykrieger.org
Johns Hopkins School of Medicine

---

**Copyright notice**

Many of the images in this powerpoint presentation are from *Bioinformatics and Functional Genomics* by Jonathan PevsnerCopyright © 2009 by John Wiley & Sons, Inc.

These images and materials may not be used without permission from the publisher. We welcome instructors to use these powerpoints for educational purposes, but please acknowledge the source.

The book has a homepage at http://www.bioinfbook.org including hyperlinks to the book chapters.

---

**Outline of today's lecture**

BLAST
    Practical use
    Algorithm
    Strategies

Finding distantly related proteins:
    PSI-BLAST
    Hidden Markov models

BLAST-like tools for genomic DNA
    PatternHunter
    Megablast
    BLAT, BLASTZ

---

**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.

page 101

---

**Why use BLAST?**

BLAST searching is fundamental to understanding the relatedness of any favorite query sequence to other known proteins or DNA sequences.

Applications include
• identifying orthologs and paralogs
• discovering new genes or proteins
• discovering variants of genes or proteins
• investigating expressed sequence tags (ESTs)
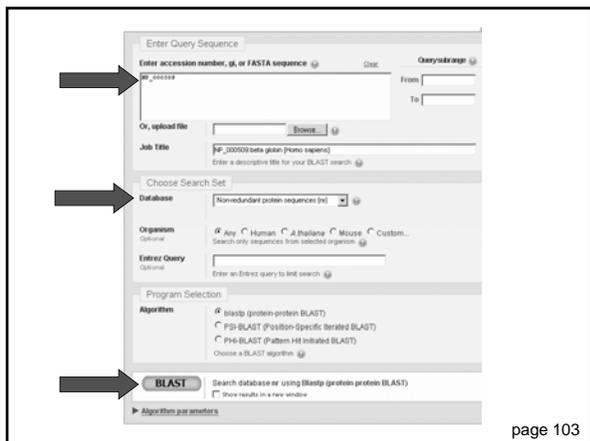• exploring protein structure and function

page 102

---

**Four components to a BLAST search**

(1) Choose the sequence (query)

(2) Select the BLAST program

(3) Choose the database to search

(4) Choose optional parameters

Then click "BLAST"

page 102

page 103

**Step 1: Choose your sequence**

Sequence can be input in FASTA format or as accession number

page 103

**Example of the FASTA format for a BLAST query**



Fig. 2.9
page 32

**Step 2: Choose the BLAST program**



page 104

**Step 2: Choose the BLAST program**

blastn (nucleotide BLAST)

blastp (protein BLAST)

blastx (translated BLAST)

tblastn (translated BLAST)

tblastx (translated BLAST)

page 104

**Choose the BLAST program**

| Program | Input | | Database |
|---------|-------|-----|----------|
| blastn | DNA | 1 | DNA |
| blastp | protein | 1 | protein |
| blastx | DNA | 6 | protein |
| tblastn | protein | 6 | DNA |
| tblastx | DNA | 36 | DNA |

page 104

## DNA potentially encodes six proteins

```
                5' CAT CAA
               5' ATC AAC
              5' TCA ACT
5' CATCAACTACAACTCCAAAGACACCCTTACACATCAACAAACCTACCCAC 3'
3' GTAGTTGATGTTGAGGTTTCTGTGGGAATGTGTAGTTGTTTGGATGGGTG 5'
                                    5' GTG GGT
                                   5' TGG GTA
                                  5' GGG TAG
```

page 105

## Step 3: choose the database

nr = non-redundant (most general database)
dbest = database of expressed sequence tags
dbsts = database of sequence tag sites
gss = genomic survey sequences



protein databases

nucleotide databases

page 106

## Step 4a: Select optional search parameters



organism

Entrez!

algorithm

page 107

## Step 4a: optional blastp search parameters



Expect

Word size

Scoring matrix

Filter, mask

page 108

## Step 4a: optional blastn search parameters



Expect

Word size

Match/mismatch scores

Filter, mask

page 108

## Step 4: optional parameters

You can...
• choose the organism to search
• turn filtering on/off
• change the substitution matrix
• change the expect (e) value
• change the word size
• change the output format

page 106

(a) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Default settings:
Unfiltered ("composition-based statistics")

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 32.7 bits (73),  Expect = 0.034, Method: Composition-based stats.
Identities = 30/101 (29%), Positives = 41/101 (40%), Gaps = 14/101 (13%)

Query  10  LLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELG  69
               LA+L L  P P+ A +    LCGS L    L  VC +          +R A+
Sbjct  16  FLAILLLSSFTPSDASI--RLCGSRLTTTLLAVCRNQLCTGLTAFKRSADQSY-------  66

Query  70  GGPGAGSLQPLALEGSLQKRG-IVEQCCTSICSLYQLENYC  109
                A + L       QKRG I  +CC   CS   L+ +C
Sbjct  67  ----APTTRDLFHIHHQQKRGGIATECCEKRCSFAYLKTFC  103

Our starting point: search human insulin against worm
RefSeq proteins by blastp using default parameters

page 109

---

(b) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Option: No compositional adjustment

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 34.7 bits (78),  Expect = 0.009, Method: Composition-based stats.
Identities = 30/100 (30%), Positives = 41/100 (41%), Gaps = 14/100 (14%)

Query  11  LALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELGG  70
               LA+L L  P P+ A +    LCGS L    L  VC +          +R A+
Sbjct  17  LAILLLSSFTPSDASIR--LCGSRLTTTLLAVCRNQLCTGLTAFKRSADQSY-------  66

Query  71  GPGAGSLQPLALEGSLQKRG-IVEQCCTSICSLYQLENYC  109
                A + L       QKRG I  +CC   CS   L+ +C
Sbjct  67  ---APTTRDLFHIHHQQKRGGIATECCEKRCSFAYLKTFC  103

Note that the bit score, Expect value, and percent identity
all change with the "no compositional adjustment" option

page 109

---

(c) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Option: conditional compositional score matrix adjustment

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 33.5 bits (75),  Expect = 0.020, Method: Compositional matrix adjust.
Identities = 27/100 (27%), Positives = 39/100 (39%), Gaps = 12/100 (12%)

Query  10  LLALLALWGPDPAAAFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELG  69
               LA+L L  P P+ A +    LCGS L    L  VC +          +R A+
Sbjct  16  FLAILLLSSFTPSDASIR--LCGSRLTTTLLAVCRNQLCTGLTAFKRSADQ-------S  65

Query  70  GGPGAGSLQPLALEGSLQKRGIVEQCCTSICSLYQLENYC  109
                P     L    +  ++ GI  +CC   CS   L+ +C
Sbjct  66  YAPTTRDL--FHIHHQQKRGGIATECCEKRCSFAYLKTFC  103

Note that the bit score, Expect value, and percent identity
all change with the compositional score matrix adjustment

page 109

---

(d) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Option: Filter low complexity regions

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 25.4 bits (54),  Expect = 6.3, Method: Composition-based stats.
Identities = 11/24 (45%), Positives = 14/24 (58%), Gaps = 1/24 (4%)

Query  87  QKRG-IVEQCCTSICSLYQLENYC  109
           QKRG I  +CC   CS   L+ +C
Sbjct  80  QKRGGIATECCEKRCSFAYLKTFC  103

Note that the bit score, Expect value, and percent identity
all change with the filter option

page 109

---

(e) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Option: Mask for lookup table only

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 32.7 bits (73),  Expect = 0.034, Method: Composition-based stats.
Identities = 30/101 (29%), Positives = 41/101 (40%), Gaps = 14/101 (13%)

Query  10  llallalwgpdpaaaFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELG  69
               LA+L L  P P+ A +    LCGS L    L  VC +          +R A+
Sbjct  16  FLAILLLSSFTPSDASI--RLCGSRLTTTLLAVCRNQLCTGLTAFKRSADQSY-------  66

Query  70  GGPGAGSLQPLALEGSLQKRG-IVEQCCTSICSLYQLENYC  109
                A + L       QKRG I  +CC   CS   L+ +C
Sbjct  67  ----APTTRDLFHIHHQQKRGGIATECCEKRCSFAYLKTFC  103

**Filtering**
**(the filtered sequence is the query**
**in lowercase and grayed out)**

page 109

---

(e) Query: human insulin NP_000198
Program: blastp
Database: *C. elegans* RefSeq
Option: Mask for lookup table only

>ref|NP_501926.1| UG INSulin related family member (ins-1) [Caenorhabditis elegans]
Length=109

Score = 32.7 bits (73),  Expect = 0.034, Method: Composition-based stats.
Identities = 30/101 (29%), Positives = 41/101 (40%), Gaps = 14/101 (13%)

Query  10  llallalwgpdpaaaFVNQHLCGSHLVEALYLVCGERGFFYTPKTRREAEDLQVGQVELG  69
               LA+L L  P P+ A +    LCGS L    L  VC +          +R A+
Sbjct  16  FLAILLLSSFTPSDASI--RLCGSRLTTTLLAVCRNQLCTGLTAFKRSADQSY-------  66

Query  70  GGPGAGSLQPLALEGSLQKRG-IVEQCCTSICSLYQLENYC  109
                A + L       QKRG I  +CC   CS   L+ +C
Sbjct  67  ----APTTRDLFHIHHQQKRGGIATECCEKRCSFAYLKTFC  103

Note that the bit score, Expect value, and percent identity
could change with the "mask for lookup table only" option

page 109

## BLAST search output: top portion

BLAST — Basic Local Alignment Search Tool — My NCBI

Home | Recent Results | Saved Strategies | Help

NCBI/BLAST/blastp suite/Formatting Results - GS1F74BK011

Edit and Resubmit   Save Search Strategies   ►Formatting options   ►Download

**NP_000509:beta globin [Homo sapiens]**

Query ID  gi|4504349|ref|NP_000509.1|  ← **query**
Description  beta globin [Homo sapiens]
>gi|55635219|ref|XP_508242.1| PREDICTED:
hypothetical protein [Pan troglodytes]
>gi|56749856|sp|P68871.2|HBB_HUMAN RecName:
Full=Hemoglobin subunit beta; AltName:
Full=Hemoglobin beta chain; AltName: Full=Beta-

hemoglobin, beta [synthetic construct]
>gi|189053145|dbj|BAG34767.1| unnamed protein
product [Homo sapiens]  ← **taxonomy**

Molecule type  amino acid
Query Length  147

**database** →
Database Name  nr
Description  All non-redundant GenBank CDS
translations+PDB+SwissProt+PIR+PRF excluding
environmental samples from WGS projects
Program  BLASTP 2.2.22+ ►Citation  ← **program**

Other reports: ►Search Summary [Taxonomy reports] [Distance tree of results] [Related Structures] [Multiple alignment] NEW

▼ Graphic Summary

▼ Show Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.

Query seq.
Specific hits            globin
Superfamilies      globin_like superfamily

page 112

## BLAST search output: taxonomy report summarizes species with matches



## BLAST search output: graphical output

Distribution of 17 Blast Hits on the Query Sequence

NP_058652 hemoglobin, beta adult minor chain [Mus musculus] S=244 E=1.7e-65

**Color key for alignment scores**

| <40 | 40-50 | 50-80 | 80-200 | >=200 |

Query  0    20    40    60    80    100    120    140

page 112

## BLAST search output: tabular output

Distance tree of results NEW

| Sequences producing significant alignments: | Score (Bits) | E Value |  |
|---|---|---|---|
| ref|NP_058652.1| hemoglobin, beta adult minor chain [Mus musculu | 244 | 2e-65 | UG |
| ref|NP_032246.2| hemoglobin, beta adult major chain [Mus musculu | 228 | 2e-60 | UG |
| ref|XP_970992.1| PREDICTED: similar to Hemoglobin epsilon-Y2 ... | 226 | 3e-60 | G |
| ref|NP_032247.1| hemoglobin Y, beta-like embryonic chain [Mus mu | 223 | 4e-59 | G |
| ref|NP_032245.1| hemoglobin Z, beta-like embryonic chain [Mus mu | 223 | 6e-59 | G |
| ref|XP_908314.1| PREDICTED: similar to Hemoglobin beta-H1 sub... | 203 | 4e-53 | G |
| ref|XP_978924.1| PREDICTED: similar to Hemoglobin epsilon-Y2 ... | 187 | 2e-48 | G |
| ref|XP_912634.1| PREDICTED: similar to Hemoglobin beta-2 subu... | 161 | 2e-40 | G |
| ref|XP_488069.1| PREDICTED: similar to Hemoglobin beta-2 subu... | 154 | 3e-38 | UG |
| ref|NP_032244.1| hemoglobin alpha 1 chain [Mus musculus] | 105 | 1e-23 | UG |
| ref|XP_994669.1| PREDICTED: similar to Hemoglobin alpha subun... | 101 | 3e-22 | G |
| ref|XP_356935.3| PREDICTED: similar to Hemoglobin alpha subun... | 100 | 4e-22 | UG |
| ref|NP_034535.1| hemoglobin X, alpha-like embryonic chain in ... | 94.0 | 4e-20 | UG |
| ref|NP_001029153.1| similar to hemoglobin, theta 1 [Mus musculus] | 88.2 | 2e-18 | G |
| ref|NP_778165.1| hemoglobin, theta 1 [Mus musculus] | 73.9 | 5e-14 | UG |
| ref|XP_978150.1| PREDICTED: similar to hemoglobin, beta adult... | 41.6 | 2e-04 | UG |
| ref|NP_795942.2| 5'-nucleotidase, cytosolic II-like 1 protein [M | 28.9 | 1.5 | UG |

**High scores low E values**

**Cut-off: .05? $10^{-10}$?**

page 113

## BLAST search output: alignment output

>ref|NP_000510.1| UG delta globin [Homo sapiens]
sp|P02042.2|HBD_HUMAN G RecName: Full=Hemoglobin subunit delta; AltName: Full=Hemoglobin
delta chain; AltName: Full=Delta-globin
emb|CAA23763.1| G delta globin [Homo sapiens]
►24 more sequence titles
Length=147

GENE ID: 3045 HBD | hemoglobin, delta [Homo sapiens] (Over 100 PubMed links)

Score = 284 bits (727), Expect = 2e-75, Method: Compositional matrix adjust.
Identities = 137/147 (93%), Positives = 142/147 (96%), Gaps = 0/147 (0%)

Query  1    MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPK  60
            MVHLTPEEK+AV ALWGKVNVD VGGEALGRLLVVYPWTQRFFESFGDLS+PDAVMGNPK
Sbjct  1    MVHLTPEEKTAVNALWGKVNVDAVGGEALGRLLVVYPWTQRFFESFGDLSSPDAVMGNPK  60

Query  61   VKAHGKKVLGAFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFG  120
            VKAHGKKVLGAFSDGLAHLDNLKGTF+ LSELHCDKLHVDPENFRLLGNVLVCVLA +FG
Sbjct  61   VKAHGKKVLGAFSDGLAHLDNLKGTFSQLSELHCDKLHVDPENFRLLGNVLVCVLARNFG  120

Query  121  KEFTPPVQAAYQKVVAGVANALAHCYH  147
            KEFTP +QAAYQKVVAGVANALAHCYH
Sbjct  121  KEFTPQMQAAYQKVVAGVANALAHCYH  147

## Outline of today's lecture

BLAST
    Practical use
    Algorithm
    Strategies

Finding distantly related proteins:
    PSI-BLAST
    Hidden Markov models

BLAST-like tools for genomic DNA
    PatternHunter
    Megablast
    BLAT, BLASTZ

## 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").

page 115

## 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.

BLAST is a heuristic approximation to local alignment. It examines only part of the search space.

page 115; 84

## 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)

(page 115)

## 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 …FSGTWYA… (query word is in yellow)

A list of words (w=3) is:
```
FSG SGT GTW TWY WYA
YSG TGT ATW SWY WFA
FTG SVT GSW TWF WYS
```

Fig. 4.11
page 116

## Phase 1: compile a list of words (w=3)

neighborhood word hits > threshold

(T=11)

| | |
|---|---|
| GTW 6,5,11 | 22 |
| GSW 6,1,11 | 18 |
| ATW 0,5,11 | 16 |
| NTW 0,5,11 | 16 |
| GTY 6,5,2 | 13 |
| GNW | 10 |
| GAW | 9 |

neighborhood word hits < below threshold

Fig. 4.11
page 116

**Pairwise alignment scores are determined using a scoring matrix such as Blosum62**

| | A | R | N | D | C | Q | E | G | H | I | L | K | M | F | P | S | T | W | Y | V |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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 | -1 | 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 | 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 | 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 |

Page 73

## How a BLAST search works: 3 phases

Phase 2:

Scan the database for entries that match the compiled list.

This is fast and relatively easy.

Fig. 4.11
page 116

## How a BLAST search works: 3 phases

Phase 3: when you manage to find a hit
(i.e. a match between a "word" and a database
entry), extend the hit in either direction.

Keep track of the score (use a scoring matrix)

Stop when the score drops below some cutoff.

```
KENFDKARFSGTWYAMAKKDPEG 50  RBP (query)
MKGLDIQKVAGTWYSLAMAASD. 44  lactoglobulin (hit)
```

← **extend**  **Hit!**  **extend** →

page 116

## How a BLAST search works: 3 phases

Phase 3:

In the original (1990) implementation of BLAST, hits were extended in either direction.

In a 1997 refinement of BLAST, two independent hits are required. The hits must occur in close proximity to each other. With this modification, only one seventh as many extensions occur, greatly speeding the time required for a search.

page 116

## How a BLAST search works: threshold

You can modify the threshold parameter.

The default value for blastp is 11.

To change it, enter "-f 16" or "-f 5" in the advanced options of BLAST+.

(To find BLAST+ go to BLAST → help → download.)

page 117



Fig. 4.12
page 118

## Phase 1: compile a list of words (w=3)

| neighborhood word hits > threshold | GTW 6,5,11 | 22 |
| | GSW 6,1,11 | 18 |
| | ATW 0,5,11 | 16 |
| | NTW 0,5,11 | 16 |
| | GTY 6,5,2 | 13 |
| (T=11) | | |
| | GNW | 10 |
| neighborhood word hits < below threshold | GAW | 9 |

Fig. 4.11
page 116

For blastn, the word size is typically 7, 11, or 15 (EXACT match). Changing word size is like changing threshold of proteins.
w=15 gives fewer matches and is faster than w=11 or w=7.

For megablast (see below), the word size is 28 and can be adjusted to 64. What will this do? Megablast is VERY fast for finding closely related DNA sequences!

---

### How to interpret a BLAST search: expect value

It is important to assess the statistical significance of search results.

For global alignments, the statistics are poorly understood.

For local alignments (including BLAST search results), the statistics are well understood. The scores follow an extreme value distribution (EVD) rather than a normal distribution.

page 118

---



normal distribution

Fig. 4.13
page 119

---

**The probability density function of the extreme value distribution (characteristic value u=0 and decay constant λ=1)**



normal distribution

extreme value distribution

Fig. 4.13
page 119

---

### How to interpret a BLAST search: expect value

The expect value $E$ is the number of alignments with scores greater than or equal to score $S$ that are expected to occur by chance in a database search.

An $E$ value is related to a probability value $p$.

The key equation describing an $E$ value is:

$$E = Kmn\, e^{-\lambda S}$$

page 120

---

$$E = Kmn\, e^{-\lambda S}$$

This equation is derived from a description of the extreme value distribution

$S$ = the score

$E$ = the expect value = the number of high-scoring segment pairs (HSPs) expected to occur with a score of at least $S$

$m$, $n$ = the length of two sequences

$\lambda$, $K$ = Karlin Altschul statistics

page 120

8

## Some properties of the equation $E = Kmn\ e^{-\lambda S}$

• The value of E decreases exponentially with increasing $S$ (higher S values correspond to better alignments). Very high scores correspond to very low $E$ values.

• The $E$ value for aligning a pair of random sequences must be negative! Otherwise, long random alignments would acquire great scores

• Parameter $K$ describes the search space (database).

• For E=1, one match with a similar score is expected to occur by chance. For a very much larger or smaller database, you would expect E to vary accordingly

page 120

## From raw scores to bit scores

• There are two kinds of scores:
  raw scores (calculated from a substitution matrix) and
  bit scores (normalized scores)

• Bit scores are comparable between different searches because they are normalized to account for the use of different scoring matrices and different database sizes

S' = bit score = $(\lambda S - \ln K) / \ln 2$

The $E$ value corresponding to a given bit score is:
$E = mn\ 2^{-S'}$

Bit scores allow you to compare results between different database searches, even using different scoring matrices.

page 121

## How to interpret BLAST: $E$ values and $p$ values

The expect value $E$ is the number of alignments with scores greater than or equal to score $S$ that are expected to occur by chance in a database search. A $p$ value is a different way of representing the significance of an alignment.

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

page 121

## How to interpret BLAST: $E$ values and $p$ values

Very small $E$ values are very similar to $p$ values.
$E$ values of about 1 to 10 are far easier to interpret than corresponding $p$ values.

| E | p |
|---|---|
| 10 | 0.99995460 |
| 5 | 0.99326205 |
| 2 | 0.86466472 |
| 1 | 0.63212056 |
| 0.1 | 0.09516258 (about 0.1) |
| 0.05 | 0.04877058 (about 0.05) |
| 0.001 | 0.00099950 (about 0.001) |
| 0.0001 | 0.0001000 |

Table 4.3
page 122

## How to interpret BLAST: overview



## Search Parameters

| | |
|---|---|
| Program | blastp |
| Word size — **word size w = 3** | 3 |
| Expect value — **10 is the E value** | 10 |
| Hitlist size | 100 |
| Gapcosts — **gap penalties** | 11,1 |
| Matrix — **BLOSUM matrix** | BLOSUM62 |
| Threshold — **threshold score = 11** | 11 |
| Composition-based stats | 2 |
| Filter string | F |
| Genetic Code | 1 |
| Window Size | 40 |

## Database

| | |
|---|---|
| Posted date | Jun 12, 2009 5:40 PM |
| Number of letters | 2,279,144,659 — **length of database** |
| Number of sequences | 6,500,228 |
| Entrez query | none |

Fig. 4.14
page 122

## Slide 1 (Fig. 4.14)

| Karlin-Altschul statistics | **EVD parameters** | |
|---|---|---|
| Params | Ungapped | Gapped |
| Lambda | 0.320339 | 0.267 |
| K | 0.136843 | 0.041 |
| H | 0.422367 | 0.14 |
| **Results Statistics** | | |
| Length adjustment | 111 | |
| Effective length of query **147 − 111 = 36** | 36 | **m** |
| Effective length of database | 1557619351 | **n** |
| Effective search space | 56074296636 | **mn** |
| Effective search space used | 56074296636 | |

**Effective search space**
**= mn**
**= length of query x db length**

Fig. 4.14
page 122

## Slide 2

### Why set the E value to 20,000?

Suppose you perform a search with a short query (e.g. 9 amino acids). There are not enough residues to accumulate a big score (or a small E value).

Indeed, a match of 9 out of 9 residues could yield a small score with an E value of 100 or 200. And yet, this result could be "real" and of interest to you.

By setting the E value cutoff to 20,000 you do not change the way the search was done, but you do change which results are reported to you.

## Slide 3

### Outline of today's lecture

BLAST
    Practical use
    Algorithm
    Strategies

Finding distantly related proteins:
    PSI-BLAST
    Hidden Markov models

BLAST-like tools for genomic DNA
    PatternHunter
    Megablast
    BLAT, BLASTZ

## Slide 4

### BLAST search strategies

General concepts

How to evaluate the significance of your results

How to handle too many results

How to handle too few results

BLAST searching with HIV-1 pol, a multidomain protein

page 123

## Slide 5 (Fig. 4.16)

### Sometimes a real match has an E value > 1



real match?

**…try a reciprocal BLAST to confirm**

Fig. 4.16
page 125

## Slide 6 (Fig. 4.17)

### Sometimes a similar *E* value occurs for a short exact match and long less exact match



short, nearly exact

long, only 31% identity, similar E value

Fig. 4.17
page 125

## Assessing whether proteins are homologous

```
>gi|4505583|ref|NP_002562.1| progestagen-associated endometrial protein (placental protein 14,
         pregnancy-associated endometrial alpha-2-globulin, alpha
         uterine protein); Progestagen-associated endometrial
         protein (placental protein 14) [Homo sapiens]
 gi|190215|gb|AAA60147.1| (J04129) placental protein 14 [Homo sapiens]
         Length = 162

 Score = 32.0 bits (71), Expect = 0.49
 Identities = 26/107 (24%), Positives = 48/107 (44%), Gaps = 11/107 (10%)

Query: 26  RVKENFDKARFSGTWYAMAKKDPEGLFLQDNIVAEFSVDETGQMSATAKGRVLLNNWD- 84
           + K++ + + +GTW++HA      + L    + A   V  T +      +L+ W+
Sbjct: 5   QTWQDLELPKLAGTWHSMAMAT-NNISLMATLKAPLRVHITSLLPTPEDNLEIVLHRWEN 63

Query: 85  -VCADMVGTFTDTEDPAKFKMKYWGVASFLQKGNDDHWIVDTDYDTY 130
            C +       T +P KFK+ Y  VA      ++  ++DTDYD +
Sbjct: 64  NSCVEKKVLGEKTGNPKKFKINY-TVA-------NEATLLDTDYDNF 102
```

RBP4 and PAEP:
Low bit score, E value 0.49, 24% identity ("twilight zone").
But they are indeed homologous. Try a BLAST search
with PAEP as a query, and find many other lipocalins.

~Fig. 4.18
page 126

---

**The universe of lipocalins (each dot is a protein)**



Fig. 5.7
Page 151

---

## BLAST search with PAEP as a query finds many other lipocalins



Fig. 4.19
page 127

---

Using human beta globin as a query, here are the blastp results
searching against human RefSeq proteins (PAM30 matrix).
Where is myoglobin? It's absent! We need to use **PSI-BLAST**.



---

## Outline of today's lecture

BLAST
      Practical use
      Algorithm
      Strategies

Finding distantly related proteins:
      PSI-BLAST
      Hidden Markov models

BLAST-like tools for genomic DNA
      PatternHunter
      Megablast
      BLAT, BLASTZ

---

## Two problems standard BLAST cannot solve

[1] Use human beta globin as a query against human
RefSeq proteins, and blastp does not "find" human
myoglobin. This is because the two proteins are too
distantly related. PSI-BLAST at NCBI as well as hidden
Markov models easily solve this problem.

[2] How can we search using 10,000 base pairs as a
query, or even millions of base pairs? Many BLAST-
like tools for genomic DNA are available such as
PatternHunter, Megablast, BLAT, and BLASTZ.

Page 141

## Position specific iterated BLAST: PSI-BLAST

The purpose of PSI-BLAST is to look deeper into the database for matches to your query protein sequence by employing a scoring matrix that is customized to your query.

Page 146

## PSI-BLAST is performed in five steps

[1] Select a query and search it against a protein database

Page 146

## PSI-BLAST is performed in five steps

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)

Page 146

### Inspect the blastp output to identify empirical "rules" regarding amino acids tolerated at each position



| | | |
|---|---|---|
| 730496 | 66 | FTVDENGQMSATAKGRVRLFNNWDVCADMIGSFTDTEDPAKFKMKYWGVASFLQKGNDDH 125 |
| 200679 | 63 | FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDPAKFKMKYWGVASFLQRGNDDH 122 |
| 206589 | 34 | FSVDEKGHMSATAKGRVRLLSNWEVCADMVGTFTDTEDPAKFKMKYWGVASFLQRGNDDH 93 |
| 2136812 | 2 | MSATAKGRVRLLNNWDVCADMVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 53 |
| 132408 | 65 | FKIEDNGKTTATAKGRVRILDKLELCANMVGTFIETNDPAKYRMKYHGALAILERGLDDH 124 |
| 267584 | 44 | FSVDESGKVTATAMGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAASYLQTGNDDH 103 |
| 267585 | 44 | FSVDGSGKVTATAQGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAAAYLQSGNDDH 103 |
| 8777608 | 63 | FTIHEDGANTATAKGRVIILNNWEMCADMMATFETTPDPAKFKMRYWGAASYLQTGNDDH 122 |
| 6607453 | 60 | FKVEEDGTMTATAIGRVIILNNWEMCANMFGTFEDTEDPAKFKMKYWGAAAYLQTGYDDH 119 |
| 10697027 | 81 | FKVQEDGTMTATATGRVIILNNWEMCANMFGTFEDTEEPARFKMKYWGAAAYLQTGYDDH 140 |
| 13645517 | 1 | MVGTFTDTEDPAKFKMKYWGVASFLQKGNDDH 32 |
| 13925316 | 38 | FSVDGSGKMTATAQGRVIILNNWEMCANMFGTFEDTPDPAKFKMRYWGAAAYLQSGNDDH 97 |
| 131649 | 65 | YTVEEDGTMTASSKGRVKLFGFWVICADMAAQYTDPTTPAKMYMYTYQGLASYLSSGGDNY 126 |

R,I,K    C    D,E,T    K,R,T    N,L,Y,G

Fig. 5.4
Page 147

|       | A  | R  | N  | D  | C  | Q  | E  | G  | H  | I  | L  | K  | M  | F  | P  | S  | T  | W  | Y  | V  |
|-------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 M   | -1 | -2 |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 2 K   | -1 | 1  | 0  | 1  | -4 | 2  | 4  | -2 | 0  | -3 | -3 | 3  | -2 | -4 | -1 | 0  | -1 | -3 | -2 | -3 |
| 3 W   | -3 | -3 | -4 | -5 | -3 | -2 | -3 | -3 |    |    |    |    |    | -4 | -3 | -3 | 12 | 2  | -3 |    |
| 4 V   | 0  | -3 | -3 | -4 | -1 | -3 | -3 | -4 |    |    |    |    | -3 | -2 | 0  | -3 | -1 | 4  |    |    |
| 5 W   | -3 | -3 | -4 | -5 | -3 | -2 | -3 | -3 | -3 | -3 | -2 | -3 | -2 | 1  | -4 | -3 | -3 | 12 | 2  | -3 |
| 6 A   | 5  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |
| 7 L   | 2  | -2 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 1  |
| 8 L   | 1  | -3 | -3 | -4 | -1 | -3 | -3 | -4 | -3 | 2  | 2  | -3 | 1  | 3  | -3 | -2 | -1 | -2 | 0  | 3  |
| 9 L   | 1  | -3 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 2  |
| 10 L  | 2  | -2 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 1  |
| 11 A  | 5  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |
| 12 A  | 5  |    |    |    |    |    |    |    |    |    | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |
| 13 W  | 2  |    |    |    |    |    |    |    | 1  | 4  | -3 | 2  | 1  | -3 | -3 | -2 | 7  | 0  | 0  |    |
| 14 A  | 3  |    |    |    |    |    |    |    | 2  | -2 | -1 | -2 | -3 | -1 | 1  | -1 | -3 | -3 | -1 |    |
| 15 A  | 2  |    |    |    |    |    |    |    | 0  | -3 | 0  | -2 | -3 | -1 | 3  | 0  | -3 | -2 | 2  |    |
| 16 A  | 2  |    |    |    |    |    |    |    | 2  | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | -1 |    |
| ...   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 37 S  | 2  | -1 | 0  | -1 | -1 | 0  | 0  | 0  | -1 | -2 | -3 | 0  | -2 | -3 | -1 | 4  | 1  | -3 | -2 | -2 |
| 38 G  | 0  | -3 | -1 | -2 | -3 | -2 | -2 | 6  | -2 | -4 | -4 | -2 | -3 | -4 | -2 | 0  | -2 | -3 | -3 | -4 |
| 39 T  | 0  | -1 | 0  | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -2 | -1 | 1  | 5  | -3 | -2 | 0  |
| 40 W  | 3  | -3 | -4 | -5 | -3 | -2 | -3 | -3 | -3 | -3 | -2 | -3 | -2 | 1  | -4 | -3 | -3 | 12 | 2  | -3 |
| 41 Y  | 2  | -2 | -2 | -3 | -3 | -2 | -2 | -3 | 2  | -2 | -1 | -2 | -1 | 3  | -3 | -2 | -2 | 2  | 7  | -1 |
| 42 A  | 4  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |

20 amino acids

all the amino acids from position 1 to the end of your PSI-BLAST query protein

Fig. 5.5
Page 149

|       | A  | R  | N  | D  | C  | Q  | E  | G  | H  | I  | L  | K  | M  | F  | P  | S  | T  | W  | Y  | V  |
|-------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| 1 M   | -1 | -2 | -2 | -3 | -2 | -1 | -2 | -3 | -2 | 1  | 2  | -2 | 6  | 0  | -3 | -2 | -1 | -2 | -1 | 1  |
| 2 K   | -1 | 1  | 0  | 1  | -4 | 2  | 4  | -2 | 0  | -3 | -3 | 3  | -2 | -4 | -1 | 0  | -1 | -3 | -2 | -3 |
| 3 W   | -3 | -3 | -4 | -5 | -3 | -2 | -3 | -3 | -3 | -3 | -2 | -3 | -2 | 1  | -4 | -3 | -3 | 12 | 2  | -3 |
| 4 V   | 0  | -3 | -3 | -4 | -1 | -3 | -3 | -4 | -4 | 3  | 1  | -3 | 1  | -1 | -3 | -2 | 0  | -3 | -1 | 4  |
| 5 W   | -3 | -3 | -4 | -5 | -3 | -2 | -3 | -3 | -3 | -3 | -2 | -3 | -2 | 1  | -4 | -3 | -3 | 12 | 2  | -3 |
| 6 A   | 5  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |
| 7 L   | -2 | -2 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 1  |
| 8 L   | -1 | -3 | -3 | -4 | -1 | -3 | -3 | -4 | -3 | 2  | 2  | -3 | 1  | 3  | -3 | -2 | -1 | -2 | 0  | 3  |
| 9 L   | -1 | -3 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 2  |
| 10 L  | -2 | -2 | -4 | -4 | -1 | -2 | -3 | -4 | -3 | 2  | 4  | -3 | 2  | 0  | -3 | -3 | -1 | -2 | -1 | 1  |
| 11 A  | 5  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |
| 12 A  | 5  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |
| 13 W  | -2 | -3 | -4 | -4 | -2 | -2 | -3 | -4 | -3 | 1  | 4  | -3 | 2  | 1  | -3 | -3 | -2 | 7  | 0  | 0  |
| 14 A  | 3  | -2 | -1 | -2 | -1 | -1 | -2 | 4  | -2 | -2 | -1 | -2 | -3 | -1 | 1  | -1 | -3 | -3 | -1 |    |
| 15 A  | 2  | -1 | 0  | -1 | -2 | 2  | 0  | 2  | -1 | -3 | 0  | -2 | -3 | -1 | 3  | 0  | -3 | -2 | 2  |    |
| 16 A  | 4  | -2 | -1 | -2 | -1 | -1 | -1 | 3  | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | -1 |    |
| ...   |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |    |
| 37 S  | 2  | -1 | 0  | -1 | -1 | 0  | 0  | 0  | -1 | -2 | -3 | 0  | -2 | -3 | -1 | 4  | 1  | -3 | -2 | -2 |
| 38 G  | 0  | -3 | -1 | -2 | -3 | -2 | -2 | 6  | -2 | -4 | -4 | -2 | -3 | -4 | -2 | 0  | -2 | -3 | -3 | -4 |
| 39 T  | 0  | -1 | 0  | -1 | -1 | -1 | -1 | -2 | -2 | -1 | -1 | -1 | -1 | -2 | -1 | 1  | 5  | -3 | -2 | 0  |
| 40 W  | -3 | -3 | -4 | -5 | -3 | -2 | -3 | -3 | -3 | -3 | -2 | -3 | -2 | 1  | -4 | -3 | -3 | 12 | 2  | -3 |
| 41 Y  | -2 | -2 | -2 | -3 | -2 | -2 | -2 | -3 | 2  | -2 | -1 | -2 | -1 | 3  | -3 | -2 | -2 | 2  | 7  | -1 |
| 42 A  | 4  | -2 | -2 | -2 | -1 | -1 | 0  | -2 | -2 | -2 | -1 | -1 | -3 | -1 | 1  | 0  | -3 | -2 | 0  |    |

Fig. 5.5
Page 149

```
         A  R  N  D  C  Q  E  G  H  I  L  K  M  F  P  S  T  W  Y  V
 1 M    -1 -2 -2 -3 -2 -1 -2 -3 -2  1  2 -2  6  0 -3 -2 -1 -2 -1  1
 2 K    -1  1  0  1 -4  2  4 -2  0 -3 -3  3 -2 -4 -1  0 -1 -3 -2 -3
 3 W    -3 -3 -4 -5 -3 -2 -3 -3 -3 -3 -2 -3 -2  1 -4 -3 -3 12  2 -3
 4 V     0 -3 -3 -4 -1 -3 -3 -4 -4  3  1 -3  1 -1 -3 -2  0 -3 -1  4
 5 W    -3 -3 -4 -5 -3 -2 -3 -3 -3 -3 -2 -3 -2  1 -4 -3 -3 12  2 -3
 6 A     5 -2 -2 -2 -1 -1 -1  0 -2 -2 -2 -1 -1 -3 -1  1  0 -3 -2  0
 7 L    -2 -2 -4 -4 -1 -2 -3 -4 -3  2  4 -3  2  0 -3 -3 -1 -2 -1  1
 8 L    -1 -3 -3 -4                              3 -3 -2 -1 -2  0  3
 9 L    -1 -3 -4 -4                              0 -3 -3 -1 -2 -1  2
10 L    -2 -2 -4 -4                              0 -3 -3 -1 -2 -1  1
11 A     5 -2 -2 -2                             -3 -1  1  0 -3 -2  0
12 A     5 -2 -2 -2                             -3 -1  1  0 -3 -2  0
13 W    -2 -3 -4 -4                              1 -3 -3 -2  7  0  0
14 A     3 -2 -1 -2                             -3 -1  1 -1 -3 -3 -1
15 A     2 -1                                   -3 -1  3  0 -3 -2 -2
16 A     4 -2 -1 -2                             -3 -1  1  0 -3 -2 -1
...
37 S     2 -1  0 -1                              2 -3 -1  4  1 -3 -2 -2
38 G     0 -3 -1 -2                             -4 -2  0 -2 -3 -3 -4
39 T     0 -1  0 -1                              1 -3 -3  1  5 -3 -2  0
40 W    -3 -3 -4 -5 -3 -2 -3 -3 -3 -3 -2 -3     1 -4 -3 -3 12  2 -3
41 Y    -2 -2 -2 -3 -3 -2 -2 -3  2 -2 -1 -2 -1  3 -3 -2 -2  2  7 -1
42 A     4 -2 -2 -2 -1 -1 -1  0 -2 -2 -2 -1 -1 -3 -1  1  0 -3 -2  0
```

note that a given amino acid (such as alanine) in your query protein can receive different scores for matching alanine— depending on the position in the protein

Fig. 5.5
Page 149

## PSI-BLAST is performed in five steps

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)

[3] The PSSM is used as a query against the database

[4] PSI-BLAST estimates statistical significance (E values)

Page 146



## PSI-BLAST is performed in five steps

[1] Select a query and search it against a protein database

[2] PSI-BLAST constructs a multiple sequence alignment then creates a "profile" or specialized position-specific scoring matrix (PSSM)

[3] The PSSM is used as a query against the database

[4] PSI-BLAST estimates statistical significance (E values)

[5] Repeat steps [3] and [4] iteratively, typically 5 times. At each new search, a new profile is used as the query.

Page 146

## Results of a PSI-BLAST search

| Iteration | # hits | # hits > threshold |
|---|---|---|
| 1 | 104 | 49 |
| 2 | 173 | 96 |
| 3 | 236 | 178 |
| 4 | 301 | 240 |
| 5 | 344 | 283 |
| 6 | 342 | 298 |
| 7 | 378 | 310 |
| 8 | 382 | 320 |

Table 5-2
Page 146

**PSI-BLAST search: human RBP versus RefSeq, iteration 1**



See Fig. 5.6
Page 150

**PSI-BLAST search: human RBP versus RefSeq, iteration 2**



See Fig. 5.6
Page 150

**PSI-BLAST search: human RBP versus RefSeq, iteration 3**



See Fig. 5.6
Page 150

**RBP4 match to ApoD, PSI-BLAST iteration 1**
**E value 3e-07**



Fig. 5.6
Page 150

**RBP4 match to ApoD, PSI-BLAST iteration 2**
**E value 1e-42**
**Note that PSI-BLAST E values can improve dramatically!**



Fig. 5.6
Page 150

**RBP4 match to ApoD, PSI-BLAST iteration 3**
**E value 6e-34**



Fig. 5.6
Page 150

**The universe of lipocalins (each dot is a protein)**

retinol-binding protein

apolipoprotein D

odorant-binding protein

Fig. 5.7
Page 151

**Scoring matrices let you focus on the big (or small) picture**

retinol-binding protein

your RBP query

Fig. 5.7
Page 151

**Scoring matrices let you focus on the big (or small) picture**

PAM250

PAM30

retinol-binding protein

Blosum80

Blosum45

Fig. 5.7
Page 151

**PSI-BLAST generates scoring matrices more powerful than PAM or BLOSUM**

retinol-binding protein

Fig. 5.7
Page 151

## PSI-BLAST: performance assessment

Evaluate PSI-BLAST results using a database in which protein structures have been solved and all proteins in a group share ≤ 40% amino acid identity.

Page 150

## PSI-BLAST: the problem of corruption

PSI-BLAST is useful to detect weak but biologically meaningful relationships between proteins.

The main source of false positives is the spurious amplification of sequences not related to the query. For instance, a query with a coiled-coil motif may detect thousands of other proteins with this motif that are not homologous.

Once even a single spurious protein is included in a PSI-BLAST search above threshold, it will not go away.
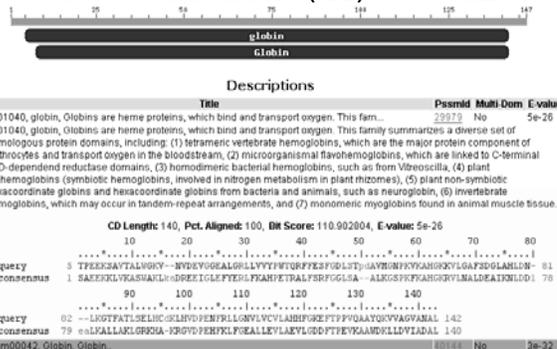
Page 152

## PSI-BLAST: the problem of corruption

Corruption is defined as the presence of at least one false positive alignment with an E value $< 10^{-4}$ after five iterations.

Three approaches to stopping corruption:

[1] Apply filtering of biased composition regions

[2] Adjust E value from 0.001 (default) to a lower value such as E = 0.0001.

[3] Visually inspect the output from each iteration. Remove suspicious hits by unchecking the box.

Page 152

---

## Conserved domain database (CDD) uses RPS-BLAST



### Descriptions

| Title | Pssmid | Multi-Dom | E-value |
|---|---|---|---|
| cd01040, globin, Globins are heme proteins, which bind and transport oxygen. This fam... | 29979 | No | 5e-26 |

cd01040, globin, Globins are heme proteins, which bind and transport oxygen. This family summarizes a diverse set of homologous protein domains, including: (1) tetrameric vertebrate hemoglobins, which are the major protein component of erythrocytes and transport oxygen in the bloodstream, (2) microorganismal flavohemoglobins, which are linked to C-terminal FAD-dependend reductase domains, (3) homodimeric bacterial hemoglobins, such as from Vitreoscilla, (4) plant leghemoglobins (symbiotic hemoglobins, involved in nitrogen metabolism in plant rhizomes), (5) plant non-symbiotic hexacoordinate globins and hexacoordinate globins from bacteria and animals, such as neuroglobin, (6) invertebrate hemoglobins, which may occur in tandem-repeat arrangements, and (7) monomeric myoglobins found in animal muscle tissue.

CD Length: 140, Pct. Aligned: 100, Bit Score: 110.902804, E-value: 5e-26

```
              10        20        30        40        50        60        70        80
      ....*....|....*....|....*....|....*....|....*....|....*....|....*....|....*....|
query    5 TPEEKSAYTALVGKY--NVDEVGGEALGRLLVVYPWTQRFFESFGDLSTpdAVMGNPKVKAHGKKVLGAFSDGLAHLDN- 81
consensus 1 SAEEKKLVKASWAKLadDREEIGLEFYERLFKAMPETRALFSRFGGLSA--ALKGSPKFKAHGKRVLNALDEAIKNLDD1 78

              90       100       110       120       130       140
      ....*....|....*....|....*....|....*....|....*....|....*....|
query   82 --LKGTFATLSELHCdKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVANAL 142
consensus 79 eaLKALLAKLGKKHA-KRGVDPKHFKLFVEALLEVLAEVLGDDFTPEVKAAWDKLLDVIADAL 140
```

pfam00042, Globin, Globin. | 40144 | No | 3e-32 |

| Main idea: you can search a query protein against a database of position-specific scoring matrices | Fig. 5.8 Page 153 |

---

## Outline of today's lecture

BLAST
 Practical use
 Algorithm
 Strategies

Finding distantly related proteins:
 PSI-BLAST
 Hidden Markov models

BLAST-like tools for genomic DNA
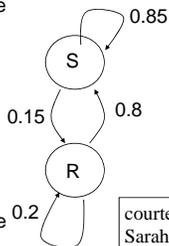 PatternHunter
 Megablast
 BLAT, BLASTZ

---

## Multiple sequence alignment to profile HMMs

• in the 1990's people began to see that aligning sequences to profiles gave much more information than pairwise alignment alone.

• Hidden Markov models (HMMs) are "states" that describe the probability of having a particular amino acid residue at arranged in a column of a multiple sequence alignment

• HMMs are probabilistic models

• Like a hammer is more refined than a blast, an HMM gives more sensitive alignments than traditional techniques such as progressive alignments

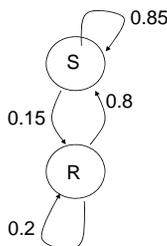Page 155

---

## Simple Markov Model



Rain = dog may not want to go outside

Sun = dog will probably go outside

Markov condition = no dependency on anything but nearest previous state ("memoryless")

courtesy of Sarah Wheelan

---

## Simple Hidden Markov Model



P(dog goes out in rain) = 0.1
P(dog goes out in sun) = 0.85

Observation: YNNNYYNNNYN

(Y=goes out, N=doesn't go out)

What is underlying reality (the hidden state chain)?

| | | | position | | | |
|---|---|---|---|---|---|---|
| | Probability | 1 | 2 | 3 | 4 | 5 |
| | p(H) | 1.0 | | | | |
| | p(A) | | 0.4 | | | |
| | p(I) | | 0.2 | | | |
| 1D8U HAMSV | p(G) | | 0.4 | | | |
| 1OJ6A HIRKV | p(M) | | | 0.2 | | |
| 2hhbB HGKKV | p(R) | | | 0.2 | | |
| 1FSL HAEKL | p(K) | | | 0.2 | | |
| 2MM1 HGATV | p(E) | | | 0.2 | | |
| | p(A) | | | 0.2 | | |
| | p(S) | | | | 0.2 | |
| | p(K) | | | | 0.6 | |
| | p(T) | | | | 0.2 | |
| | p(V) | | | | | 0.8 |
| | p(L) | | | | | 0.2 |

Fig. 5.11
Page 157

$p(HARTV) = (1.0)(0.4)(0.2)(0.2)(0.8) = 0.0128$
Log odds score $= \ln(1.0) + \ln(0.4) + \ln(0.2) + \ln(0.2) + \ln(0.8) =$



Fig. 5.11
Page 157



Fig. 5.12
Page 158



Fig. 5.12
Page 158



Fig. 5.15
Page 160

## HMMER: build a hidden Markov model

Determining effective sequence number   ... done. [4]
Weighting sequences heuristically        ... done.
Constructing model architecture          ... done.
Converting counts to probabilities       ... done.
Setting model name, etc.                 ... done. [x]

Constructed a profile HMM (length 230)
Average score:     411.45 bits
Minimum score:     353.73 bits
Maximum score:      460.63 bits
Std. deviation:     52.58 bits

Fig. 5.13
Page 159

## HMMER: calibrate a hidden Markov model

HMM file:            lipocalins.hmm
Length distribution mean: 325
Length distribution s.d.: 200
Number of samples:    5000
random seed:          1034351005
histogram(s) saved to:   [not saved]
POSIX threads:        2
- - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

HMM  : x
mu   : -123.894508
lambda :   0.179608
max  : -79.334000

Fig. 5.13
Page 159

## HMMER: search an HMM against GenBank

```
Scores for complete sequences (score includes all domains):
Sequence                        Description      Score    E-value  N
--------                        -----------      -----    -------  ---
gi|20888903|ref|XP_129259.1|    (XM_129259) ret  461.1    1.9e-133  1
gi|132407|sp|P04916|RETB_RAT    Plasma retinol-  458.0    1.7e-132  1
gi|20548126|ref|XP_005907.5|    (XM_005907) sim  454.9    1.4e-131  1
gi|5803139|ref|NP_006735.1|     (NM_006744) ret  454.6    1.7e-131  1
gi|20141667|sp|P02753|RETB_HUMAN Plasma retinol-  451.1    1.9e-130  1
.
.
gi|16767588|ref|NP_463203.1|    (NC_003197) out  318.2    1.9e-90   1


gi|5803139|ref|NP_006735.1|: domain 1 of 1, from 1 to 195: score 454.6, E = 1.7e-131
                  *->mkwVMkLLLLaALagvfgaAErdAfsvgkCrvpsPPRGfrVkeNFDv
                     mkwV++LLLLaA +  +aAErd     Crv+s    frVkeNFD+
   gi|5803139    1   MKWVWALLLLAA--W--AAAERD------CRVSS----FRVKENFDK 33

                  erylGtWYeIaKkDprFErGLllqdkItAeySleEhGsMsataeGrirVL
                  +r++GtWY++aKkDp  E GL+lqd+I+Ae+S++E+G+Msata+Gr+r+L
   gi|5803139   34   ARFSGTWYAMAKKDP--E-GLFLQDNIVAEFSVDETGQMSATAKGRVRLL 80

                  eNkelcADkvGTvtqiEGeasevfLtadPaklklKyaGvaSflqpGfddy
                  +N+++cAD+vGT+t++E        dPak+k+Ky+GvaSf1q+G+dd+
   gi|5803139   81   NNWDVCADMVGTFTDTE----------DPAKFKMKYWGVASFLQKGNDDH 120
```

Fig. 5.13
Page 159

## PFAM is a database of HMMs and an essential resource for protein families
### http://pfam.sanger.ac.uk/

welcome trust
sanger
institute

HOME | SEARCH | BROWSE | FTP | HELP
ABOUT

Pfam
keyword search  Go

Pfam 24.0 (October 2009, 11912 families)

The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). More...

QUICK LINKS        QUERY PFAM BY KEYWORD
SEQUENCE SEARCH    Search for keywords in text data in the Pfam database.
VIEW A PFAM FAMILY                        Go  Example
VIEW A CLAN        You can also use the keyword search box at the top of every page.
VIEW A SEQUENCE
VIEW A STRUCTURE
KEYWORD SEARCH
JUMP TO

## Outline of today's lecture

BLAST
        Practical use
        Algorithm
        Strategies

Finding distantly related proteins:
        PSI-BLAST
        Hidden Markov models

BLAST-like tools for genomic DNA
        PatternHunter
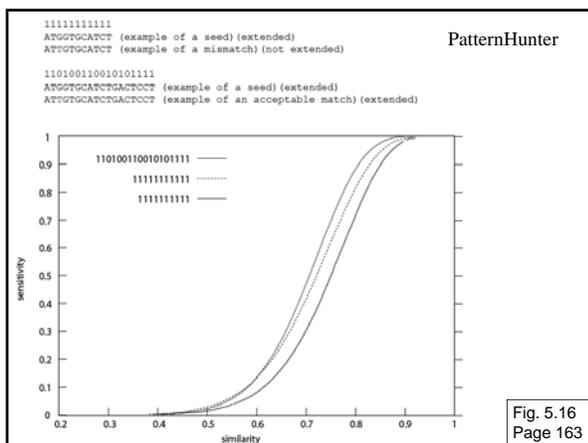        Megablast
        BLAT, BLASTZ

## BLAST-related tools for genomic DNA

The analysis of genomic DNA presents special challenges:
• There are exons (protein-coding sequence) and
  introns (intervening sequences).
• There may be sequencing errors or polymorphisms
• The comparison may between be related species
  (e.g. human and mouse)

Page 161

## BLAST-related tools for genomic DNA

Recently developed tools include:

• MegaBLAST at NCBI.

• BLAT (BLAST-like alignment tool). BLAT parses an entire
  genomic DNA database into words (11mers), then
  searches them against a query. Thus it is a mirror image
  of the BLAST strategy. See http://genome.ucsc.edu

• SSAHA at Ensembl uses a similar strategy as BLAT.
  See http://www.ensembl.org

Page 162

## Slide 1 (PatternHunter)

```
1111111111
ATGGTGCATCT (example of a seed)(extended)
ATTGTGCATCT (example of a mismatch)(not extended)

110100110010101111
ATGGTGCATCTGACTCCT (example of a seed)(extended)
ATTGTGCATCTGACTCCT (example of an acceptable match)(extended)
```

PatternHunter



Fig. 5.16
Page 163

## Slide 2 (MegaBLAST at NCBI)

MegaBLAST at NCBI

--very fast
--uses very large word sizes (e.g. W=28)
--use it to align long, closely related sequences



Fig. 5.19
Page 167

## Slide 3 (MegaBLAST output)

MegaBLAST output



Sequences producing significant alignments:
(Click headers to sort columns)

| Accession | Description | Max score | Total score | Query coverage | E value | Max ident |
|---|---|---|---|---|---|---|
| M92296.1 | Pongo pygmaeus gamma-1 and gamma-2 globin genes, co | 1.805e+04 | 2.046e+04 | 26% | 0.0 | 95% |
| M18038.1 | Orangutan (P.pygmaeus) beta- and eta-globin pseudogene | 1.095e+04 | 1.156e+04 | 15% | 0.0 | 94% |
| X05035.1 | Orangutan epsilon-globin gene with Alu repeats in flanking | 6547 | 8190 | 10% | 0.0 | 96% |
| M18796.1 | Orangutan beta- and delta-globin gene intergenic region wi | 5171 | 5889 | 7% | 0.0 | 96% |
| M21825.1 | Orangutan delta globin gene, complete cds | 3616 | 4516 | 5% | 0.0 | 97% |
| M16209.1 | Orangutan gamma-2-fetal globin gene, complete cds | 2950 | 6424 | 9% | 0.0 | 94% |
| M16200.1 | Orangutan gamma-1-fetal globin gene, complete cds | 2935 | 6667 | 9% | 0.0 | 94% |

Fig. 5.19
Page 167

## Slide 4 (To access BLAT)

**To access BLAT, visit http://genome.ucsc.edu**



"BLAT on DNA is designed to quickly find sequences of 95% and greater similarity of length 40 bases or more. It may miss more divergent or shorter sequence alignments. It will find perfect sequence matches of 33 bases, and sometimes find them down to 20 bases. BLAT on proteins finds sequences of 80% and greater similarity of length 20 amino acids or more. In practice DNA BLAT works well on primates, and protein blat on land vertebrates."          --BLAT website

## Slide 5 (BLAT Search Genome)



Paste DNA or protein sequence here in the FASTA format

Fig. 5.20
Page 167

## Slide 6 (BLAT output)

**BLAT output includes browser and other formats**

Blastz



Fig. 5.17
Page 165

Blastz (laj software): human versus rhesus duplication



Fig. 5.18
Page 166

Blastz (laj software): human versus rhesus gap



Fig. 5.18
Page 166

BLAT



Fig. 5.20
Page 167

BLAT



Fig. 5.20
Page 167

LAGAN



Fig. 5.21
Page 168

SSAHA



## Outline of today's lecture

BLAST
    Practical use
    Algorithm
    Strategies

Finding distantly related proteins:
    PSI-BLAST
    Hidden Markov models

BLAST-like tools for genomic DNA
    PatternHunter
    Megablast
    BLAT, BLASTZ

## Where we are in the course

--We started with "access to information" (Chapter 2)

--We next covered pairwise alignment (Chapter 3), then BLAST in which one sequence is compared to a database (Chapters 4, 5)

--Next we'll describe multiple sequence alignment (Chapter 6)

--We'll then visualize multiple sequence alignments as phylogenetic trees (Chapter 7). That topic spans molecular evolution.