

Bioinformatics Chapter 8.

Sequence Alignment and Database Searching

오석준

장병탁

서울대 바이오정보기술연구소

서울대 컴퓨터공학부 바이오지능연구실 &

E-mail: sjaugh@cbit.snu.ac.kr

바이오정보기술연구소

E-mail: btzhang@cse.snu.ac.kr

Sirk-June Augh
Center for Bioinformation Technology (CBIT)
Seoul National University

Byoung-Tak Zhang
Center for Bioinformation Technology (CBIT) &
Biointelligence Laboratory
School of Computer Science and Engineering
Seoul National University

This material is available at <http://bi.snu.ac.kr/> &
<http://cbit.snu.ac.kr/>

Outline

- The Evolutionary Basis of Sequence Alignment
- The Modular Nature of Proteins
- Optimal Alignment Methods
- Substitution Scores and Gap Penalties
- Statistical Significance of Alignments
- Database Similarity Searching
- FASTA
- BLAST
- Database Searching Artifacts
- Position-Specific Scoring Matrices
- Spliced Alignments

(C) 2002 SNU CSE Biointelligence Lab (BI)

2

Similarity vs. Homology

- Similarity:
 - ◆ An observable quantity that might be expressed, as say, percent identity or some other suitable measure.
- Homology:
 - ◆ Refers to a conclusion drawn from these data that two genes share a common evolutionary history

(C) 2002 SNU CSE Biointelligence Lab (BI)

3

Conserved Positions are often of Functional Importance

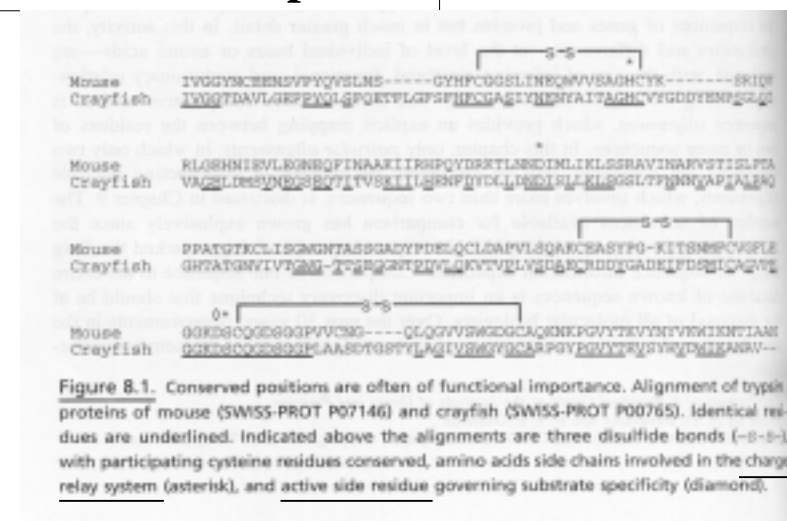


Figure 8.1. Conserved positions are often of functional importance. Alignment of tryptic proteins of mouse (SWISS-PROT P07146) and crayfish (SWISS-PROT P00765). Identical residues are underlined. Indicated above the alignments are three disulfide bonds (-S-S-) with participating cysteine residues conserved, amino acids side chains involved in the charge relay system (asterisk), and active site residue governing substrate specificity (diamond).

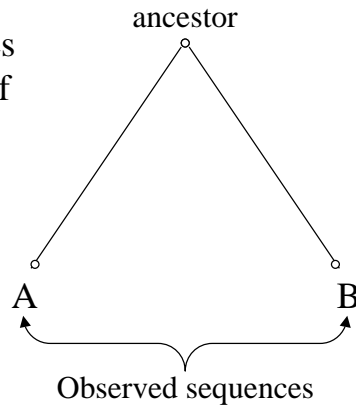
(C) 2002 SNU CSE Biointelligence Lab (BI)

4

Alignment Jargon

Evolutionarily related sequences differ from one other because of several processes:

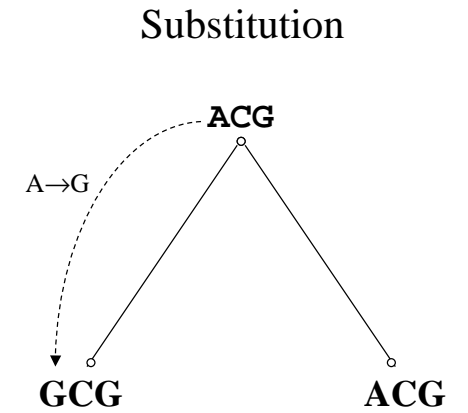
- *Substitutions*
- *Insertions*
- *Deletions*



Alignment Jargon

GCG
||
ACG

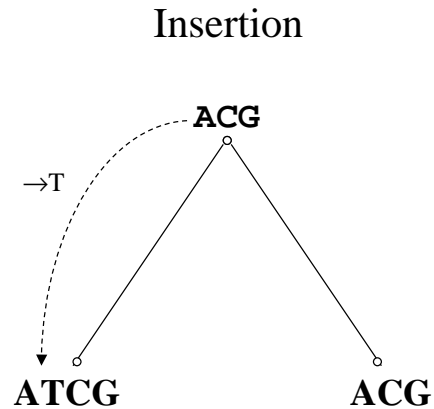
- 1 mismatch
- 2 matches



Alignment Jargon

ATCG
| ||
A-CG

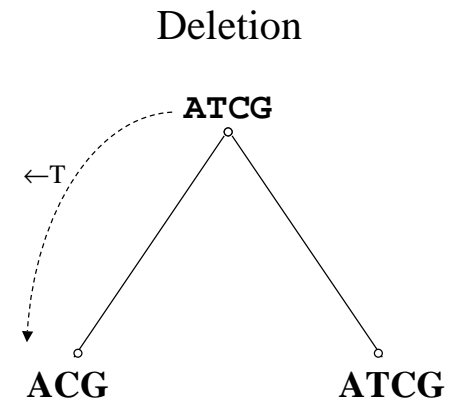
- 0 mismatches
- 3 matches
- 1 gap



Alignment Jargon

ATCG
| ||
A-CG

- 0 mismatches
- 3 matches
- 1 gap



Alignment vs. Prediction: When are Alignment Reliable?

- New sequences are aligned against all sequences in DB
 - ◆ Hints toward structural/functional relationships + functional insights
- Fundamental question
 - ◆ When is the sequence similarity high enough that one may infer a structural/functional similarity from the pairwise alignment of two sequences?
 - ◆ Given the detected overlap in a sequence segment, can a similarity threshold be defined that sifts out cases where the inference will be reliable?
- Answer
 - ◆ It depends on the structural/functional aspect one wants to investigate
 - ◆ Different for each task
- Prediction
 - ◆ When alignment alone is not enough to lead a reliable inference

Local Alignment vs. Global Alignment

- Local Alignment
 - ◆ Attempts to align regions of sequences with the highest density of matches. In doing so, one or more islands of subalignments are created in the aligned sequences.
- Global Alignment
 - ◆ Attempts to match as many characters as possible, from end to end, in a set of two or more sequences.

Local Alignment vs. Global Alignment

	A	C	C	A	C	A	C	A
	0	0	0	0	0	0	0	0
A	0	1	0	0	1	0	1	0
C	0	0	2	1	0	2	0	2
A	0	1	0	1	2	0	3	1
C	0	0	2	1	0	3	1	4
C	0	0	0	3	2	1	2	2
A	0	1	0	1	4	2	2	1
T	0	0	0	0	2	3	1	1
A	0	1	0	0	1	1	4	2

	A	C	C	A	C	A	C	A
	0	m	m	m	m	m	m	m
A	m	1	-1	-3	-5	-7	-9	-11
C	m	-1	2	0	-2	-4	-6	-8
A	m	-3	0	1	1	-1	-3	-5
C	m	-5	-2	1	0	2	0	-2
C	m	-7	-4	-1	0	1	1	-1
A	m	-9	-6	-3	0	-1	2	0
T	m	-11	-8	-5	-2	-1	0	1
A	m	-13	-10	-7	-4	-1	0	-1

Distribution of 2 Blast Hits on the Query Sequence

BLAST search result

Query: Soy bean (plant) leghemoglobin

Database: homo sapiens

Alignment result shows two merely matched sequences, but their functions and structures are surprisingly coincided.

Sequences producing significant alignments:

Sequence	Score	E Value
gi1120002 c1114002 1 Homo sapiens, theta 1 (Homo sapiens)	27.0	7.8
gi1120002 c1114002 1 Homo sapiens, theta 1 (Homo sapiens)	27.0	7.8

Alignments

Query: 1: APTKSLA... (length 12)

Subject: 1: ... (length 142)

Score = 27.0 bits (83), Expect = 7.0

Ident: 100% (208), Positives = 45/113 (40%), Gaps = 5/113 (4%)

Query: 2: ... (length 12)

Subject: 2: ... (length 142)

Score = 27.0 bits (83), Expect = 7.0

Ident: 100% (208), Positives = 45/113 (40%), Gaps = 5/113 (4%)

Optimal Global Sequence Alignment

```

Human-ZCr      MATGQKLMRAVRVFEFGGPEVLKLRSDIAVPIPKDHQVLIKVHACGVNPNVETVIRSGTYS
Ecoli-QOR      -----MATRIEFHKHGGPEVLQA-VEFTPADPAENEIQVENKAIGINFIDTYIRSGLYS
                *****
                * * * * *
Human-ZCr      RKPLLPTYPGSDVAGVIEAVGDNASAFKKGDRVFTSSTISGGYAEYALAADHTVYKLPK
Ecoli-QOR      -PPSLPSGLGTEAAGIVSKVSGVKHIKAGDRVVYAQSALGAYSSVHNIADKAAILPAA
                * * * * *
Human-ZCr      LDFKQGAAGIPIPYFTAYRALIHSACVKAGESVLVHGASGGVGLAACQIARAYGLKILGTA
Ecoli-QOR      ISFEQAAASFLKGLTVYLLRKYEIKPDEQLFHAAAGGVGLIACQWAKALGAKLIGTV
                * * * * *
Human-ZCr      GTEEGQKIVLQNGAHEVFNHREVNYIDKIKKYVGEKGIDIIEMLANVNLKDLSSLSSHG
Ecoli-QOR      GTAQKAQSALKAGAWQVINYREEDLVERLKEITGGKKVVRVYDSVGRDWTWERSLDCLQR
                * * * * *
Human-ZCr      GRVIVVG-SRGTIEINPRDTMAKES----SIIGVTLFSSTKEEFQYAAALQAGMEICWL
Ecoli-QOR      GLMVSFGNSSGAVTVNLGILNQKGLVYTRPQLQYITTRBELTEASNELFSLIASGVI
                * * * * *
Human-ZCr      KPVIGSQ--YPLEKVAEAHENI IHGSGATGKMILL
Ecoli-QOR      KVDVAEQQKYPKDAQRAHE-ILESRATQGSLLIP
                * * * * *
    
```

Figure 8.2. Optimal global sequence alignment. Alignment of the amino acid sequences of human zeta-crystallin (SWISS-PROT Q08257) and *E. coli* quinone oxidoreductase (SWISS-PROT P28304). It is an optimal global alignment produced by the CLUSTAL W program (Higgins et al., 1996). Identical residues are marked by asterisks below the alignment, and dots indicate conserved residues.

Prediction of Functional Features

- Two protein sequences share sequence similarity
 - ◆ These proteins share common function?
- New sequence identity threshold is required for prediction of function
 - ◆ Threshold used for structural problems can not be used.
- Solution
 - ◆ Split each sequence into a number of subsequences
 - ◆ The fraction of aligned site per alignment

Modular Nature of Proteins

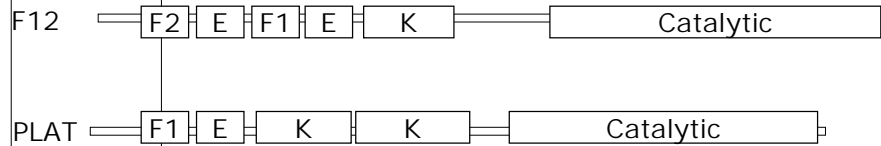


Figure 8.3. Modular structure of two proteins involved in blood clotting. Schematic representation of the modular structure of human tissue plasminogen activator and coagulation factor XII. A module labeled C is shared by several proteins involved in blood clotting. F1 and F2 are frequently repeated units that were first seen in fibronectin. E is a module resembling epidermal growth factor. A module known as a “kringle domain” is denoted K.

Measuring Alignment Quality

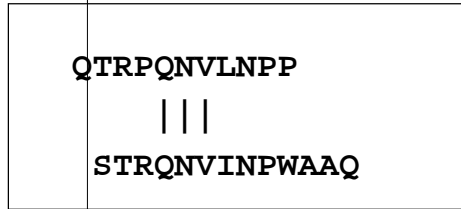
Good alignments should have ...

- “many” exact matches
- “few” mismatches
 - “many” of the mismatches should be similar residues
- “few” gaps

Measuring Alignment Quality

Begin with...

Longest Exact Match

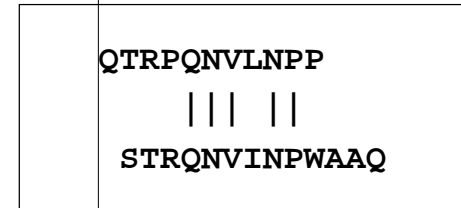


$$S = 3a$$

S=alignment score
a=match score

Measuring Alignment Quality

... allow some mismatches

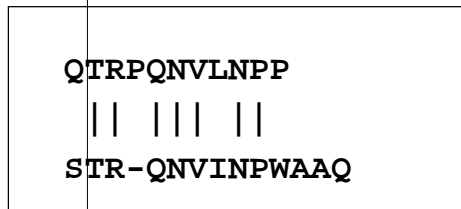


$$S = 5a - 1b$$

S=alignment score
a=match score
b=mismatch penalty

Measuring Alignment Quality

...and finally, introduce some gaps



$$S = 7a - 1b - 1c$$

S=alignment score
a=match score
b=mismatch penalty
c=gap penalty

Dot-matrix Representations

● Figure 8.4

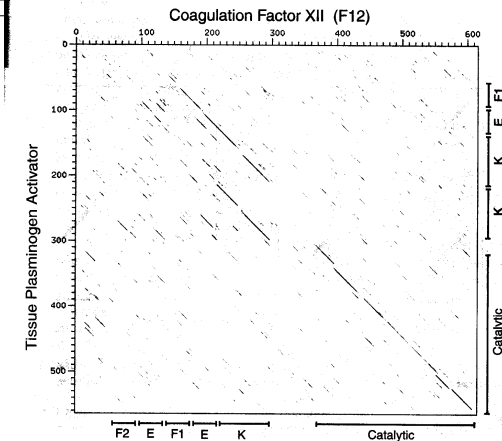
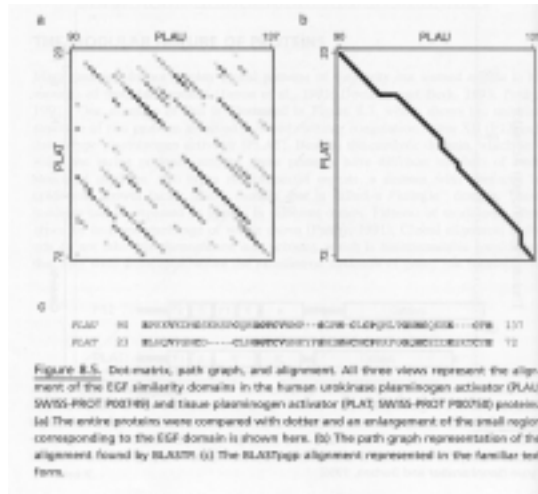


Figure 8.4. Dot matrix sequence comparison. Dot matrix comparison of the amino acid sequences of human coagulation factor XII (F12; SWISS-PROT P00748) and tissue plasminogen activator (PLAT; SWISS-PROT P00750). The figure was generated using the dotter program (Sonnhammer and Durban, 1996).

Dot-matrix Representations

- Figure 8.5



(C) 2002 SNU CSE Biointelligence Lab (BI)

21

Local Alignments

- Many problems in computer science can be reduced to the task of finding the optimal path through a graph
- Some positive incremental scores will be used for aligning identical residues, with negative scores used for substitutions and gaps.
- Finding optimal local alignment
 - ◆ Dynamic programming.
 - Needleman-Wunsch algorithm (Needleman and Wunsch, 1970)
 - Smith-Waterman algorithm (Smith and Waterman, 1981)
 - Viterbi decoding algorithm

(C) 2002 SNU CSE Biointelligence Lab (BI)

22

Dynamic Programming

- General optimization technique
- Application environment
 - ◆ Problem can be recursively subdivided into two similar subproblems of smaller size
 - ◆ Solution can be obtained by piecing together the solutions to the two subproblems
 - ◆ Example: finding shortest path
 - Problem: $[A \rightarrow C \rightarrow B] \Rightarrow [A \rightarrow C]$ and $[C \rightarrow B]$
 - Solution: $[A \rightarrow C] + [C \rightarrow B] \Rightarrow [A \rightarrow B]$

(C) 2002 SNU CSE Biointelligence Lab (BI)

23

Smith-Waterman Algorithm (1)

- Finding local alignments
- Using dynamic programming

		T	C	A	T	G
C	0	0	0	0	0	0
A	0	0	1	0	0	0
T	0	1	0	0	3	2
T	0	1	1	0	2	3
G	0	0	1	1	1	3

$$s(i, j) = \text{Max} \begin{cases} S(i-1, j-1) + A(i, j) \\ S(i-1, j) - GP \\ S(i, j-1) - GP \\ FG(i-2, j) - GEP \\ HG(i, j-2) - GEP \\ 0 \end{cases}$$

TCAT*G
*CATTG

(C) 2002 SNU CSE Biointelligence Lab (BI)

24

Smith-Waterman Algorithm (2)

- Best results and slow performance
- Can grasp results missed by BLAST or FASTA.
- Available on web:
spiral.genes.nig.ac.jp/homolgy/ssearch-e.shtml
- Ref : Smith, T. F. and Waterman, M. S. (1981)
Identification of common molecular subsequence.
Journal of Molecular Biology, **147**: 196-197.



Figure 8.6. Optimal and suboptimal local alignments. The three best alignments found when using Jalisco align the sequences of human coagulation factor III (PII; SWISS-PROT 598146) and coagulation factor III (F11; SWISS-PROT P00746).

Scoring Matrix

- A unitary matrix is used for base pairs
 - ◆ Each position can be given a score of +1 if it matches and a score of zero if it does not.
- Substitution matrices are used for amino acid alignments.
 - ◆ Certain amino acids can substitute easily for one another in related proteins because of their similar physicochemical properties.

Substitution Scores

- Point accepted mutation (PAM) model of evolution
 - ◆ A unit of evolutionary divergence in which 1% of the amino acids have been changed
 - ◆ Log-odds approach
 - The substitution scores in the matrix are proportional to the natural log of the ratio of target frequencies to background frequencies
- BLOSUM substitution matrices
 - ◆ Use of a different strategy for estimating the target frequencies
 - ◆ The underlying data are derived from the BLOCKS database, which contains local alignments (“blocks”)

The BLOSUM62 Scoring Matrix

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	1										
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	-2	-3	-3	-3	-1	-3	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	1	-1	-2	1	-1	-2	-2	0	-3	-1	4	

Statistical Significance: E-values

- Let P_i be the frequency of aa i . For unrelated sequences, an alignment of i with j has probability P_{ij} .
- Given P_{ij} and s_{ij} we can calculate normalized scores ("bit scores") from the raw score, S :

$$S' = \frac{\lambda S - \ln K}{\ln 2}$$

K is a function of the database size,

λ is a function of the scoring matrix

- When 2 random sequences of length m and n are aligned, the expected number of HSPs with normalized scores greater than S' is approximately

$$E = \frac{N}{2^{S'}}$$

where $N = nm$

BLAST & FASTA

- Using heuristic algorithms
- Word based match
- Faster than Smith & Waterman
- BLAST:** www.ncbi.nlm.nih.gov/BLAST
Ref: Altschul, S. F., et al. (1990) Basic local alignment search tool. *Journal of Molecular Biology*, **215**, 403-410.
- FASTA:** www.ebi.ac.uk/fasta3
Ref: Wilbur, W. J. and Lipman, D. J. (1983) Rapid similarity searches of nucleic acid and protein data banks. *Proceedings of the National Academy of Sciences USA* **80**, 726-30.

Figure 8.9

(a) The best scores are:

```

gi|1706794|sp|P49789|HIT_HUMAN FRAGILE HISTIDINE INITIN INIT1 OPT 2-SC R(59248)
gi|1703339|sp|P49776|APH1_SCHPO BIS(5'-NUCLEOSYL) 431 395 395 536.2 2.8e-23
gi|1723425|sp|P49775|VD15_YEAST HYPOTHETICAL 24.8 290 171 316 426.1 2.9e-17
gi|1724021|sp|Q11066|HIT_MYCTU HYPOTHETICAL 20.0 178 178 184 250.7 2.2e-07
gi|417124|sp|Q04344|HIT_YEAST HIT1 PROTEIN (ORF U 159 104 157 216.2 1.8e-05
gi|418447|sp|P52084|HIT_SYNP7 HYPOTHETICAL 12.4 139 139 140 195.0 0.00028
gi|1351828|sp|P47376|HIT_MYCOB HYPOTHETICAL 15.6 132 132 133 185.9 0.0012
gi|1169826|sp|P43424|GAL7_RAT GALACTOSE-1-PHOSPHA 97 97 128 169.7 0.0072
gi|418446|sp|P52083|HIT_MYCHR HYPOTHETICAL 13.1 102 102 119 166.8 0.01
gi|1708543|sp|P49773|HIT_HUMAN PROTEIN KINASE C 87 87 116 164.5 0.014
gi|1724020|sp|P49774|HIT_MYCLE HYPOTHETICAL 17.0 131 82 117 161.6 0.02
gi|1724019|sp|P53795|HIT_CAMEL HYPOTHETICAL HIT- 98 98 116 161.5 0.02
gi|1170581|sp|P16436|IPK1_BOVIN PROTEIN KINASE C 86 86 115 160.4 0.023
gi|1730188|sp|Q03249|GAL7_MOUSE GALACTOSE-1-PHOSP 87 87 120 159.3 0.027
gi|1177047|sp|P42856|ZBI4_MALIZ 14 KD ZINC-BINDIN 132 79 112 156.3 0.04
gi|120908|sp|P07902|GAL7_HUMAN GALACTOSE-1-PHOSPH 78 78 117 154.8 0.048
gi|1177046|sp|P42855|ZBI4_BRAJU 14 KD ZINC-BINDIN 115 76 110 154.5 0.05
gi|140775|sp|P26724|HIT_AZOBH HYPOTHETICAL 13.2 115 65 109 152.6 0.064
gi|1169825|sp|P51764|GAL7_HAEM GALACTOSE-1-PHOSP 82 82 104 137.9 0.42
gi|113999|sp|P16550|APA1_YEAST 5'.5'-P-1, P-4-TE 108 66 103 137.1 0.47
    
```

(b)

```

>>gi|1169826|sp|P43424|GAL7_RAT GALACTOSE-1-PHOSPHATE TR (379 aa)
initn: 97    init1: 128    z-score: 169.7    E(1): 0.0072
Smith-Waterman score: 128;    30.8% identity in 107 aa overlap

      HIT          MSFRFG-QHLIKPSVVELKTELSFALVNRKVP
      X.....X.....X.....X.....X.....X.....X.....
Galt VWASNFLPDIAQREERSSQQTTHNQHKPILLVGKGGKSRKRELRVTSEVIVVDFPFNAV
190   200   210   220   230   240

      HIT          40   50   60   70   80
      VFGHLVCLPRLPVERPHDLRPOEVADLPTQTTRGVTVKEKPHGTSLSFSM---QDG---
      . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Galt WPFQTLPLRRHQVQLPILTPAERDLDASTWKLKTYDNLFE-TSFPYEMSNHGAPMGL
250   260   270   280   290   300

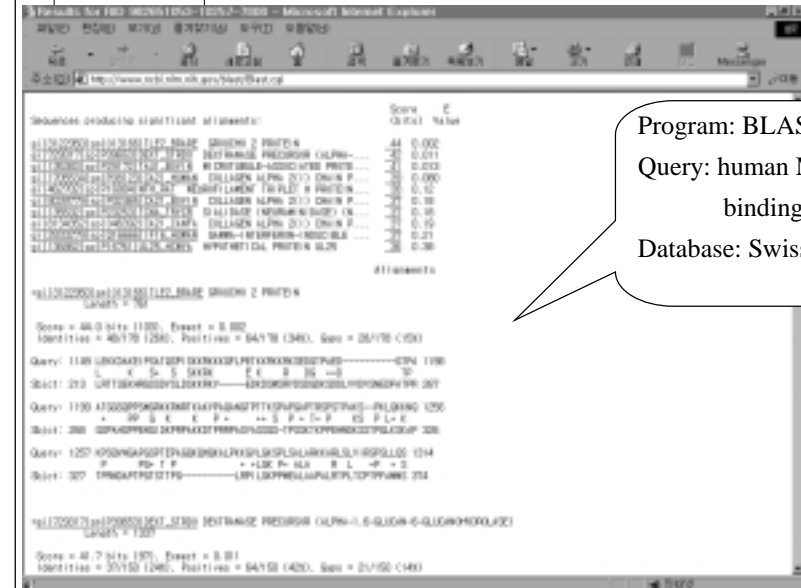
      HIT          90   100  110  120  130  140
      EAGQTVKX--VHVIVLPRKAGDPRNDSIYELQHKDEDPASWSEEMAAEAALRV
      . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . . .
Galt KTGATCDHWQLHANYIPLLRSAATVRFKFMVGYEMLAQRDLTPQQAERLRLPEVHYC
310   320   330   340   350   360
    
```

Figure 8.9. Output of a FASTA search. (a) Hit list from a FASTA search with human histidine triad (HIT) protein (SWISS-PROT P49789) as the query against the swissprot database. The search was performed using $k_{\text{cut}} = 1$. (b) Optimal local alignment of the query to one of the database entries (marked with arrow in hit list) containing the sequence of rat galactose-1-phosphate uridylyltransferase (GalT). Although the sequence similarity is weak, these proteins have been shown to share structural similarity.

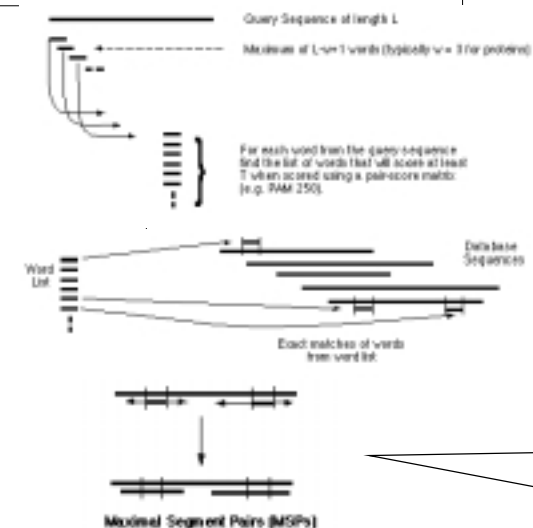
FASTA

- Score using a substitution matrix
 - ◆ Applies a substitution matrix to the 10 best regions
 - ◆ The matrix encapsulates the biological significance of word matches
 - ◆ Single best subalignment - init1
 - ◆ db strings with init1 < threshold are filtered

BLAST Result



The BLAST Search



For the query, find the list of high scoring words of length W.

Compares the word list to the database and identifies exact matches.

For each word match, extends the alignment in both directions to find alignments that score greater than a threshold of value S.

BLAST Programs

Program	Query	Database	Comments
BLASTP	Protein	Protein	Uses substitution matrix for finding distant relationship relationships; SEG filtering available
BLASTN	Nucleotide	Nucleotide	Tuned for very high-scoring matches, not distant relationships
BLASTX	Nucleotide (translated)	Protein	Useful for analysis of new DNA sequences and ESTs
TBLASTN	Protein	Nucleotide (translated)	Useful for finding unannotated coding regions in database sequences
TBLASTX	Nucleotide (translated)	Nucleotide (translated)	May be useful for EST analysis, but computationally intensive

Protein Sequence Databases for use with BLAST

Databas	Description
<i>nr</i>	Non-redundant merge of SWISS-PROT, PIR, PRF, and proteins derived form GenBank coding sequences and PDB atomic coordinates
<i>month</i>	Subset of <i>nr</i> witch is new or modified within the last 30 days
<i>swissprot</i>	The SWISS-PROT database
<i>pdb</i>	Amino acid sequences parsed from atomic coordinates of three-dimensional structures
<i>ecoli</i>	Complete set of proteins encoded by the <i>E. coli</i> genome
<i>yeast</i>	Complete set of proteins encoded by the <i>S. cerevisiae</i> genome
<i>drosoph</i>	Complete set of proteins encoded by the <i>E. melanogaster</i> genome

Nucleotide Sequence Databases for use with BLAST

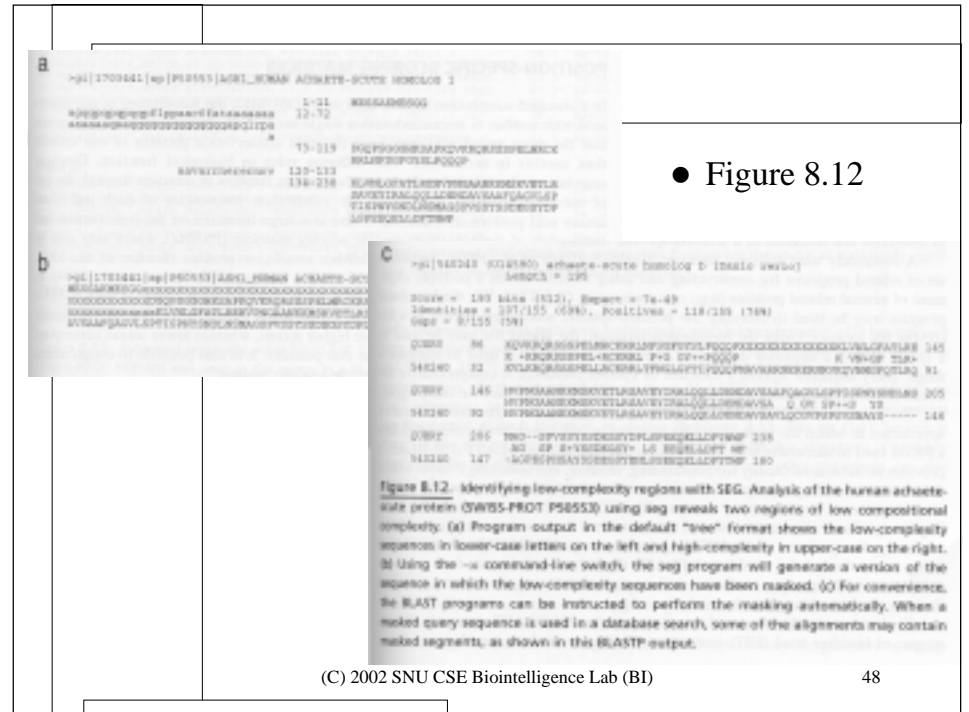
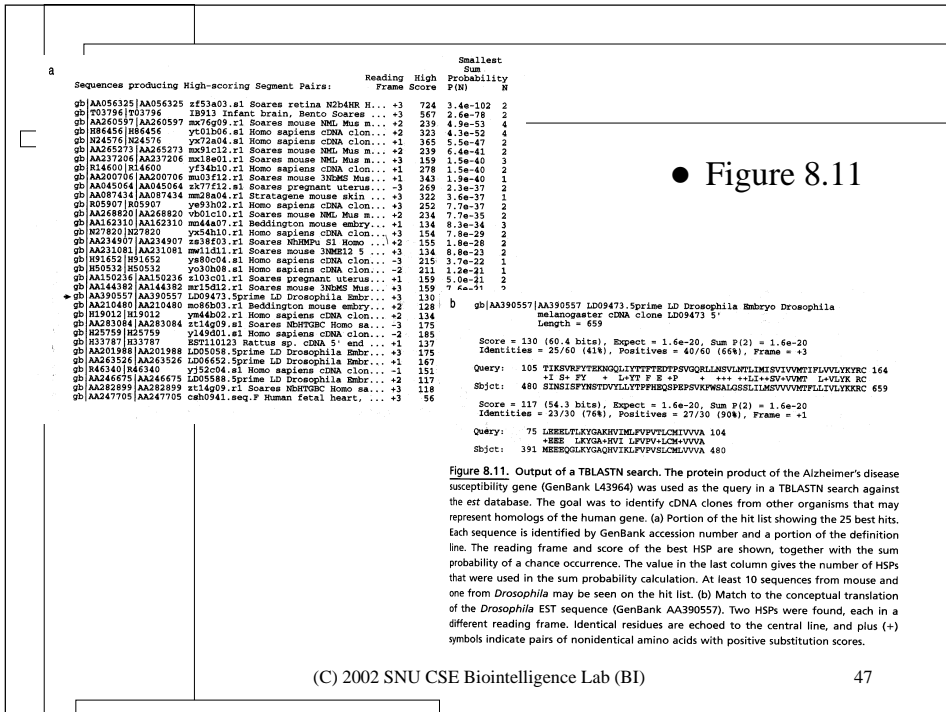
Databas	Description
<i>nr</i>	Non-redundant Genbank, excluding the EST, STS, and GSS divisions
<i>month</i>	Subset of <i>nr</i> , witch is new or modified within the last 30 days
<i>est</i>	Genbank EST division (expressed sequence tags)
<i>sts</i>	Genbank STS division (sequence tagged sites)
<i>htgs</i>	Genbank HTG division (high-throughput genomic sequences)
<i>gss</i>	Genbank GSS division (genome survey sequences)
<i>ecoli</i>	Complete genome sequence of <i>E. coli</i>
<i>yeast</i>	Complete genome sequence of <i>S. cerevisiae</i>
<i>drosoph</i>	Complete genome sequence of <i>E. melanogaster</i>
<i>mito</i>	Complete genome sequence of vertebrate mitochondria
<i>alu</i>	Collection of primate Alu repeat sequences
<i>vector</i>	Collection of popular cloning vectors

TBLASTN Search

- The protein product of the Alzheimer's disease susceptibility gene (GenBank L43964) was used as the query in a TBLASTN search against the *est* database. The goal was to identify cDNA clones from other organisms that may represent homologs of the human gene.

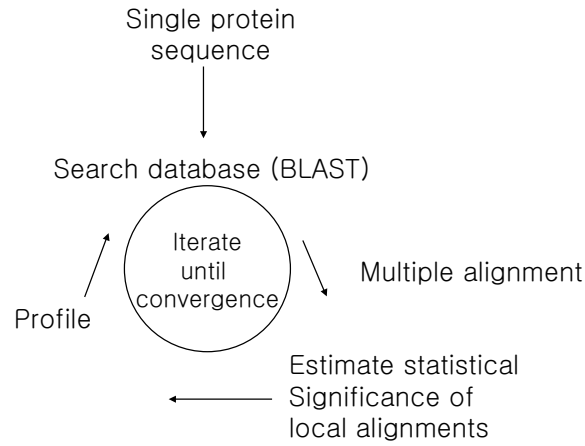


- Portion of the hit list showing the 25 best hits. Each sequence is identified by GenBank accession number and a portion of the definition line. The reading frame and score of the best HSP are shown, together with the sum probability of a chance occurrence. The value in the last column gives the number of HSPs that were used in the sum probability calculation. At least 10 sequence from mouse and one from *Drosophila* may be seen on the hit list.
- Match to the conceptual translation of the *Drosophila* EST sequence (GenBank AA390557). Two HSPs were found, each in a different reading frame. Identical residues are echoed to the central line, and plus (+) symbols indicate pairs of nonidentical amino acids with positive substitution scores.



Position Specific Scoring Matrices

- Position-Specific Iterated BLAST (PSI-BLAST)



(C) 2002 SNU CSE Biointelligence Lab (BI)

●Figure 8.13

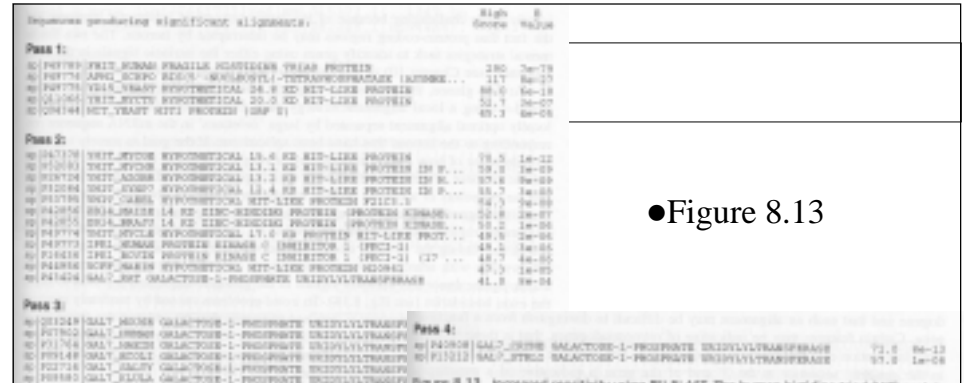
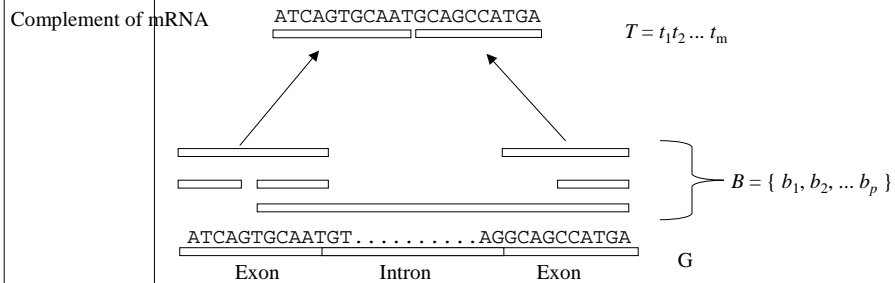


Figure 8.13. Increased sensitivity using PSI-BLAST. The human histidine triad (HT) protein (P60789) was used as the query in a BLASTP search with the PSI-BLAST functionality enabled. Definition lines, scores, and E values are shown for all statistically significant matches newly identified in each iteration.

(C) 2002 SNU CSE Biointelligence Lab (BI)

Spliced Alignments Problem

- Given strings G (genomic sequence) and T (target sequence), and a set B of substrings of G , find a set of non-overlapping strings from B whose concatenation C fits T the best; i.e., the edit distance between C and T is minimal among all sets of blocks from B .



(C) 2002 SNU CSE Biointelligence Lab (BI)

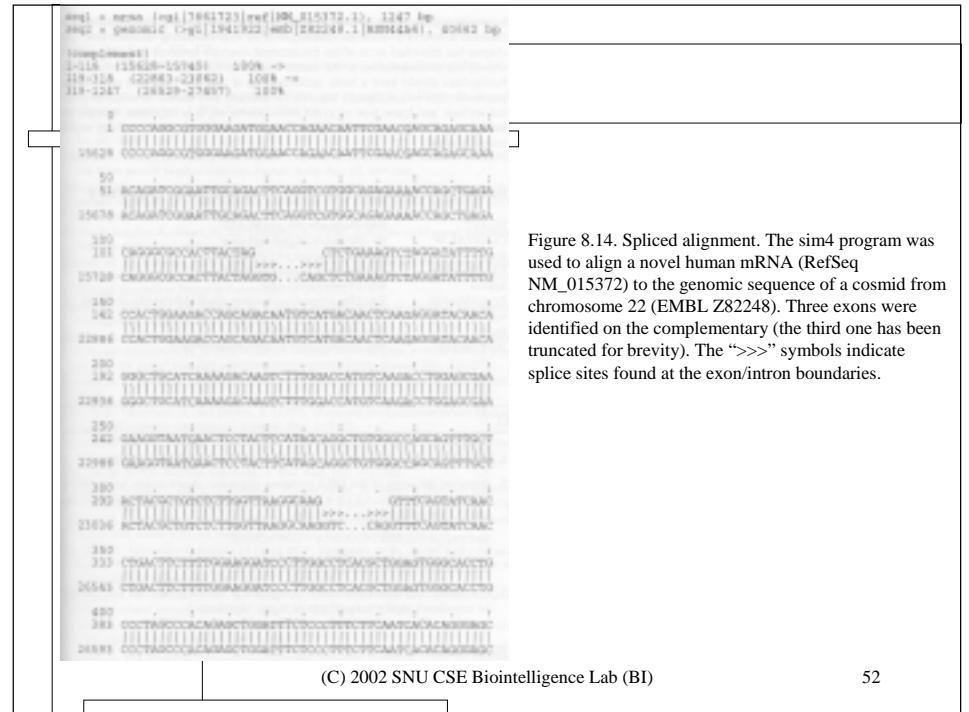


Figure 8.14. Spliced alignment. The sim4 program was used to align a novel human mRNA (RefSeq NM_015372) to the genomic sequence of a cosmid from chromosome 22 (EMBL Z82248). Three exons were identified on the complementary (the third one has been truncated for brevity). The ">>>" symbols indicate splice sites found at the exon/intron boundaries.

(C) 2002 SNU CSE Biointelligence Lab (BI)

BIOLOGY IN THE FUTURE

Systems biology
Neuroimmunology

Bioinformatics
Data mining
Combinatorial chemistry
Synthetic biology

Biochip & microfluidics (LOC)
Biocomputation