

## Protein sequence alignment and evolution

Tuesday, April 5, 2005

Protein Bioinformatics  
260.841  
Jonathan Pevsner  
pevsner@jhmi.edu

### Outline: entire course

T Mar. 29 Th Mar. 31	Introduction to physical properties of amino acids Protein Structure (level of Branden and Tooze)	Prigge Prigge
T Apr. 5 Th Apr. 7	Protein sequence alignment and evolution Principles of mass spectrometry	Pevsner Cotter
T Apr. 12 Th Apr. 14	Applications of mass spectrometry to proteomics Applications of mass spectrometry to proteomics	Pandey Pandey
T Apr. 19 Th Apr. 21	Protein structure determination Protein databases, structural classification of proteins, visualization	Prigge Ruczinski
T Apr. 26 Th Apr. 28	Protein secondary structure prediction Protein structure prediction	Ruczinski Ruczinski
T May 3 Th May 5	Protein structure prediction (CASP) Protein networks	Ruczinski Bader
T May 10 Th May 12	To be announced Protein-protein docking	Gray
T May 17 Th May 19	To be announced Final exam	

### Outline: entire course

T Mar. 29 Th Mar. 31	Introduction to physical properties of amino acids Protein Structure (level of Branden and Tooze)	Prigge Prigge
T Apr. 5 Th Apr. 7	Protein sequence alignment and evolution Principles of mass spectrometry	Pevsner Cotter
T Apr. 12 Th Apr. 14	Applications of mass spectrometry to proteomics Applications of mass spectrometry to proteomics	Pandey Pandey
T Apr. 19 Th Apr. 21	Protein structure determination Protein databases, structural classification of proteins, visualization	Prigge Ruczinski
T Apr. 26 Th Apr. 28	Protein secondary structure prediction Protein structure prediction	Ruczinski Ruczinski
T May 3 Th May 5	Protein structure prediction (CASP) Protein networks	Ruczinski Bader
T May 10 Th May 12	High throughput approaches to proteomics Protein-protein docking	Boeke Gray
T May 17 Th May 19	Lab Final exam	

### Outline: today's topic

1. How to access the sequence and structure of a protein at NCBI and the Protein Data Bank (PDB)
2. Overview of databases of all proteins: NCBI and SwissProt
3. How to align the sequences of two proteins: Dayhoff's evolutionary perspective
4. How to align the sequences of two proteins: pairwise alignment

Many of the powerpoints for today's lecture are from *Bioinformatics and Functional Genomics* (J. Pevsner, 2003). The powerpoints are available on-line at [www.bioinfbook.org](http://www.bioinfbook.org)

Chapter 2: Access to sequence data  
Chapter 3: Pairwise sequence alignment  
Chapter 4: Basic Local Alignment Search Tool (BLAST)  
Chapter 8: Protein analysis and proteomics  
Chapter 9: Protein structure

### Outline: today's topic

1. How to access the sequence and structure of a protein at NCBI and the Protein Data Bank (PDB)
2. Overview of databases of all proteins: NCBI and SwissProt
3. How to align the sequences of two proteins: Dayhoff's evolutionary perspective
4. How to align the sequences of two proteins: pairwise alignment

**NCBI National Center for Biotechnology Information**  
National Library of Medicine National Institutes of Health

PubMed All Databases BLAST OMM Books TaxBrowser Structure

Search All Databases for amyloid Go

**SITE MAP**  
Alphabetical List  
Resource Guide

**About NCBI**  
An introduction to NCBI

**Get/Bank**  
Sequence submission support and software

**Literature databases**  
PubMed, OMM, Books, and PubMed Central

**Molecular databases**  
Sequences, structures, and taxonomy

**Genomic biology**  
The human genome, whole genomes, and related resources

**Tools**  
Data mining

**Research at NCBI**  
People projects

**What does NCBI do?**  
Established in 1988 as a national resource for molecular biology information, NCBI creates public databases, conducts research in computational biology, develops software tools for analyzing genome data, and disseminates biomedical information - all for the better understanding of molecular processes affecting human health and disease. [More...](#)

**Hot Spots**  
► Assembly Archive  
► Clusters of orthologous groups  
► Coffee Break, Genes & Disease, NCBI Handbook  
► Electronic PCR  
► Entrez Home  
► Entrez Tools  
► Gene expression omnibus (GEO)  
► Human genome resources  
► Malaia genetics & genomics  
► Map Viewer  
► dbMHC  
► Mouse genome resources  
► My NCBI  
► ORF finder  
► Rat genome

**Influenza Virus Resource**  
The Influenza Virus Resource enables comparison of influenza virus strains and provides a reference for viral structures. The resource contains data from the NIAID Influenza Genome Sequencing Project and GenBank, as well as sequence-based alignments of flu sequences.

**Entrez Gene**  
You can now use Entrez to search for information centered on the concept of a gene, and connect to many sources of related information both within and outside NCBI.

**PubMed Central**  
An archive of life sciences journals  
► Free fulltext  
► Over 300,000 articles from over 150 journals  
► Linked to PubMed and fully searchable  
Use of PubMed Central requires no registration or fee. Access it from any computer with an internet connection.

[www.ncbi.nlm.nih.gov](http://www.ncbi.nlm.nih.gov)

**NCBI Entrez, The Life Sciences Search Engine**

Search across databases amyloid

25512	PubMed: biomedical literature citations and abstracts	165	Books: online books
1484	PubMed Central: free, full text journal articles	192	ONHM: online Mendelian Inheritance in Man
		10	Site Search: NCBI web and FTP sites
6450	Nucleotide: sequence database (GenBank)	219	UniGene: gene-oriented clusters of transcript sequences
3419	Protein: sequence database	14	KDD: conserved protein domain database
2	Genome: whole genome sequences	447	3D Domains: domains from Entrez Structure
125	Structure: three-dimensional macromolecular structures	353	UniSTS: markers and mapping data
8086	Taxonomy: organisms in GenBank	4	PopSet: population study data sets
6199	SNP: single nucleotide polymorphism	36203	GEO Profiles: expression and molecular abundance profiles
534	Gene: gene-centered information	14	GEO Datasets: experimental sets of GEO data
307	Homolog: eukaryotic homolog groups	8086	Cancer Chromosomes: cytogenetic databases
1	PubChem Compound: small molecule chemical structures	8086	PubChem Bioassay: bioactivity screens of chemical substances
8086	PubChem Substance: chemical substances screened for toxicity	70	GENAT: gene expression atlas of mouse central nervous system
8086	Genome Project: genome project information		

<http://www.expasy.ch> allows queries of Swiss-Prot

Site Map Search Expasy Contact us

Search Swiss-Prot/TrEMBL for amyloid Go Clear

**ExpASY Proteomics Server**

The ExpASY (Expert Protein Analysis System) proteomics server of the Swiss Institute of Bioinformatics (SIB) is dedicated to the analysis of protein sequences and structures as well as 2-D PAGE.

[Announcements](#) [Link opening](#) [Mirror Sites](#)

Databases	Tools and software packages
<ul style="list-style-type: none"> <li>Swiss-Prot and TrEMBL - Protein knowledgebase</li> <li>PROSITE - Protein families and domains</li> <li>SWISS-2D PAGE - Two-dimensional polyacrylamide gel electrophoresis</li> <li>ENZYME - Enzyme nomenclature</li> <li>SWISS-3DIMAGE - 3D images of proteins and other biological macromolecules</li> <li>SWISS-MODEL Repository - Automatically generated protein models</li> <li>GermOnLine - Knowledgebase on germ cell differentiation</li> <li>Ashbya Genome Database</li> <li>Links to many other molecular biology databases</li> </ul>	<ul style="list-style-type: none"> <li>Proteomics and sequence analysis tools               <ul style="list-style-type: none"> <li>Proteomics (Aldente (PMP), PeptideMaps, ...)</li> <li>DNA -&gt; Protein (Translat)</li> <li>Similarity searches (BLAST)</li> <li>Pattern and profile searches (ScanProsite)</li> <li>Post-translational modification and topology prediction</li> <li>Primary structure analysis (ProtParam, pI/MW, ProtScale)</li> <li>Secondary and tertiary structure prediction (SWISS-MODEL, Swiss-PaTriViewer)</li> <li>Alignment (L-COFFEE, SIM)</li> <li>Biological text analysis</li> </ul> </li> <li>ImageMaster / Melanie - Software for 2-D PAGE analysis</li> <li>MSight - Mass Spectrometry Imager</li> <li>Roche Applied Science's Biochemical Pathways</li> </ul>

### Search in Swiss-Prot and TrEMBL for: amyloid

Swiss-Prot Release 46.4 of 29-Mar-2005  
TrEMBL Release 29.4 of 29-Mar-2005

- Number of sequences found in [Swiss-Prot](#) (309) and [TrEMBL](#) (16) 319
- Note that the selected sequences can be saved to a file to be later retrieved, to do so, go to the [bottom](#) of this page.
- For more directed searches, you can use the [Sequence Retrieval System](#) [SRS](#)

#### Search in Swiss-Prot: There are matches to 103 out of 178022 entries

- A4\_BOVIN (Q28053)**  
Alzheimer's disease amyloid A4 protein homolog [Contains: Beta-amyloid protein (Beta-AFF) (A-beta)] (Fragment). (GENE: Name=APP) - Bos taurus (Bovine)
- A4\_CAEEL (Q10651)**  
Beta-amyloid-like protein precursor. (GENE: Name=apl-1, ORFNames=C42D8.8) - *Caenorhabditis elegans*
- A4\_CAFEA (Q28280)**  
Alzheimer's disease amyloid A4 protein homolog [Contains: Beta-amyloid protein (Beta-AFF) (A-beta)] (Fragment). (GENE: Name=APP) - *Canis familiaris* (Dog)
- A4\_CAFPS (Q04945)**  
Amyloid beta A4 protein precursor (APP) (ARFF) (Alzheimer's disease amyloid protein homolog) [Contains: Soluble APP-alpha (S-APP-alpha), Soluble APP-beta (S-APP-beta), CTF-alpha, CTF-beta, Beta-amyloid protein 42 (Beta-APP42), Beta-amyloid protein 40 (Beta-APP40), P3(42), P3(40), Gamma-CTF(59) (Gamma-secretase C-terminal fragment 59), Gamma-CTF(57) (Gamma-secretase C-terminal fragment 57), C31] (GENE: Name=APP) - *Canis porcellus* (Canine pup)
- A4\_DEOME (P14599)**  
Beta-amyloid-like protein precursor. (GENE: Name=Appl, Synonyms=VND, ORFNames=CQ7727) - *Drosophila melanogaster* (Fruit fly)
- A4\_FUGSI (Q93259)**  
Alzheimer's disease amyloid A4 protein homolog precursor [Contains: Beta-amyloid protein (Beta-AFF) (A-beta)] (GENE: Name=APP) - *Fugu rubripes* (Japanese pufferfish) (Takifugu rubripes)
- A4\_HUMAN (P05067)**  
Amyloid beta A4 protein precursor (APP) (ARFF) (Alzheimer's disease amyloid protein) (Cerebral vascular amyloid peptide) (CVAP) (Protease neuron-ID) (PH-ID) (APP) (PreA4) [Contains: Soluble APP-alpha (S-APP-alpha), Soluble APP-beta (S-APP-beta), C99, Beta-amyloid protein 42 (Beta-APP42), Beta-amyloid protein 40 (Beta-APP40), C83, P3(42), P3(40), Gamma-CTF(59) (Gamma-secretase C-terminal fragment 59) (Amyloid intracellular domain 59) (AID(59)), Gamma-CTF(57) (Gamma-secretase C-terminal fragment 57) (Amyloid intracellular domain 57) (AID(57)), Gamma-CTF(50) (Gamma-secretase C-terminal fragment 50) (Amyloid intracellular domain 50) (AID(50)), C31] (GENE: Name=APP, Synonyms=A4, AD1) -

### Protein Data Bank (PDB) (<http://www.pdb.org>)

**RCSB PDB PROTEIN DATA BANK**

Welcome to the PDB, the single worldwide repository for the processing and distribution of 3-D biological macromolecular structure data.

[ABOUT PDB](#) | [NEW FEATURES](#) | [USER GUIDES](#) | [FILE FORMATS](#) | [DATA INTEGRITY](#) | [STRUCTURAL GENOMICS](#) | [SOFTWARE](#) | [PUBLICATIONS](#) | [REGISTRATION](#)

**Search the Archive**  
Enter a PDB ID or keyword  
amyloid Search

**PDB Mirrors**  
"Please bookmark a mirror site"  
San Diego Supercomputer Center, UCSD  
Rutgers University  
Center for Advanced Research in Biotechnology, NIST  
Cambridge Crystallographic Data Centre, UK  
National University of Singapore  
Osaka University, Japan  
Max Delbrück Center for Molecular Medicine, Germany

**News**  
29-Mar-2005  
RCSB PDB Education Activities, ASSEMB and NSTA Members of the RCSB PDB will be participating in a variety of upcoming education-based meetings. [MORE...](#)

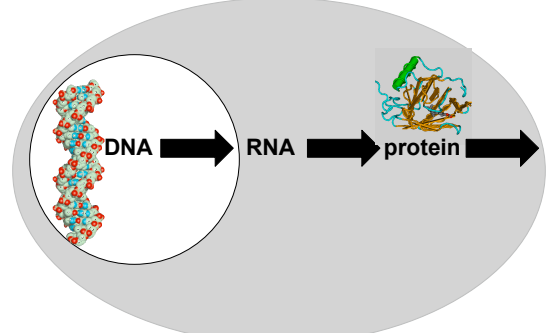
**Current Holdings**  
30263 Structures  
Last Update: 29-Mar-2005  
PDB Statistics

**We are building a new home for your molecules.**

**Molecule of the Month**  
T-cell Receptor

The Protein Data Bank (PDB) is operated by Rutgers, The State University of New Jersey, the San Diego Supercomputer Center at the University of California, San Diego, and the Center for Advanced Research in Biotechnology/UMINIST - three members of the Structural Collaborations for Structural Bioinformatics (SCSB).

### Central dogma of molecular biology



genome → transcriptome → proteome

Central dogma of bioinformatics and genomics

## Accession numbers are labels for sequences

NCBI includes databases (such as GenBank) that contain information on DNA, RNA, or protein sequences. You may want to acquire information beginning with a query such as the name of a protein of interest, or the raw nucleotides comprising a DNA sequence of interest.

DNA sequences and other molecular data are tagged with accession numbers that are used to identify a sequence or other record relevant to molecular data.

## What is an accession number?

An accession number is a label that used to identify a sequence. It is a string of letters and/or numbers that corresponds to a molecular sequence.

Examples (all for retinol-binding protein, RBP4):

X02775	GenBank genomic DNA sequence	<b>DNA</b>
NT_030059	Genomic contig	
Rs7079946	dbSNP (single nucleotide polymorphism)	
N91759.1	An expressed sequence tag (1 of 170)	<b>RNA</b>
NM_006744	RefSeq DNA sequence (from a transcript)	
NP_007635	RefSeq protein	<b>protein</b>
AAC02945	GenBank protein	
Q28369	SwissProt protein	
1KT7	Protein Data Bank structure record	

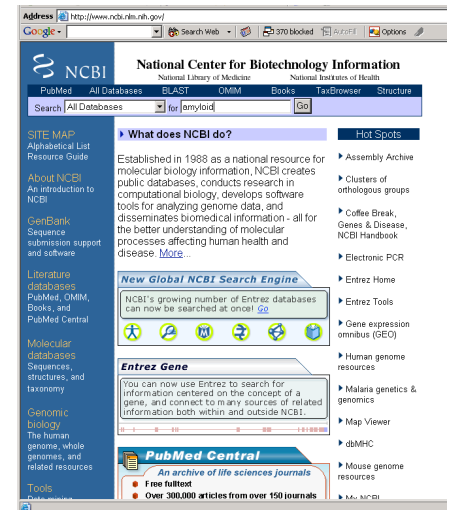
## NCBI's important RefSeq project: best representative sequences

RefSeq (accessible via the main page of NCBI) provides an expertly curated accession number that corresponds to the most stable, agreed-upon "reference" version of a sequence.

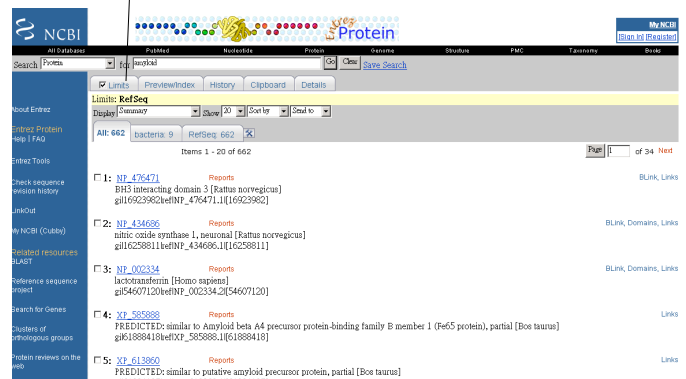
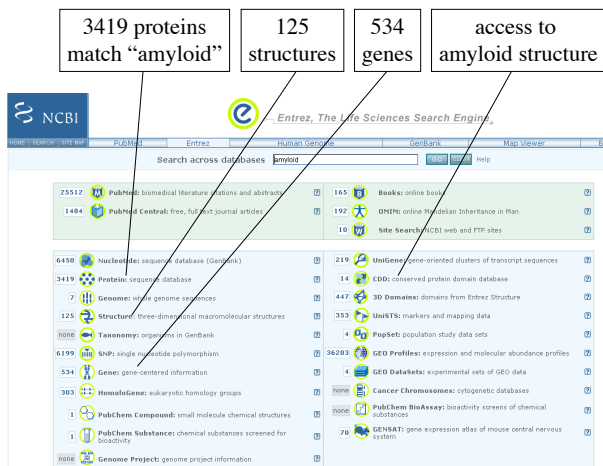
RefSeq identifiers include the following formats:

Complete genome	NC_#####
Complete chromosome	NC_#####
Genomic contig	NT_#####
mRNA (DNA format)	NM_##### e.g. NM_006744
Protein	NP_##### e.g. NP_006735

Example: type "amyloid" at NCBI



Click "protein" to find 3419 records for amyloid. Further limit the search to RefSeq only, then to human.



Your query found 354 structures in the current PDB release and you have selected 0 structures so far. (There are currently 1 structures being processed. You can select specific structures by clicking on the checkbox next to their id. If you do not select any structures, certain options will default to all structures. Use the Explore link!

Pull down to select option:

1-20

KEY:  = Download compressed (GNU zipped) PDB file  = View PDB file  = Structure viewing options

<b>133L</b>	Deposited: 01-Jun-1993 Exp. Method: X-ray Diffraction Resolution: 1.77 Å
<b>Title</b>	Role of Arg115 in the catalytic action of human lysozyme. X-ray structure of His115 and Glu115 mutants.
<b>Classification</b>	Hydrolase(O-Glycosyl)
<b>Compound</b>	Lysozyme (E.C. 3.2.1.17) Mutant With Arg 115 Replaced By His (R115H)
<b>134L</b>	Deposited: 01-Jun-1993 Exp. Method: X-ray Diffraction Resolution: 1.77 Å
<b>Title</b>	Role of Arg115 in the catalytic action of human lysozyme. X-ray structure of His115 and Glu115 mutants.
<b>Classification</b>	Hydrolase(O-Glycosyl)
<b>Compound</b>	Lysozyme (E.C. 3.2.1.17) Mutant With Arg 115 Replaced By Glu (R115E)
<b>1AAP</b>	Deposited: 14-Sep-1990 Exp. Method: X-ray Diffraction Resolution: 1.50 Å
<b>Title</b>	X-ray crystal structure of the protease inhibitor domain of Alzheimer's amyloid β-protein precursor.
<b>Classification</b>	Protease Inhibitor (Trypsin)
<b>Compound</b>	Protease Inhibitor Domain Of Alzheimer'S Amyloid β-Protein Precursor (APPI)
<b>1AMB</b>	Deposited: 21-Oct-1994 Exp. Method: NMR
<b>Title</b>	Solution structure of residues 1-28 of the amyloid β-peptide.
<b>Classification</b>	Protease Inhibitor(Trypsin)
<b>Compound</b>	Alzheimer'S Disease Amyloid β-Peptide (Residues 1 - 28) (E.C. Number Not Assigned) (NMR, Minimized Average Structure)
<b>1AMC</b>	Deposited: 14-Nov-1994 Exp. Method: NMR
<b>Title</b>	Solution structure of residues 1-28 of the amyloid β-peptide.
<b>Classification</b>	Protease Inhibitor(Trypsin)

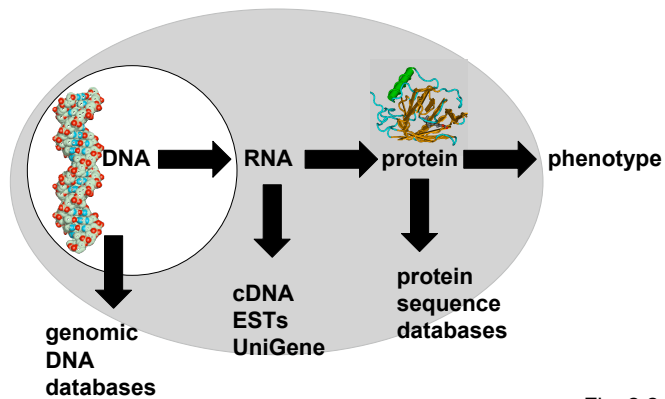
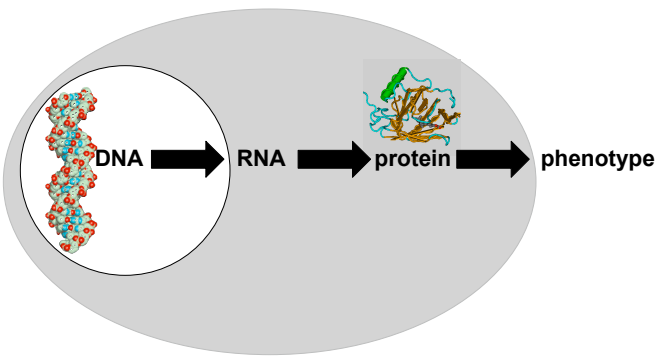


Fig. 2.2  
Page 20

**DNA**

GenBank		

**protein**

UniProt (www.uniprot.org)			
	European Bioinformatics Institute	Protein Information Resource	
			Protein Data Bank

**Growth of GenBank**

Release 146 (Feb 2005) has 46,849,831,226 base pairs

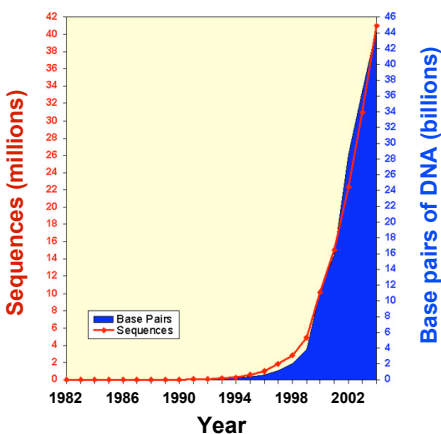
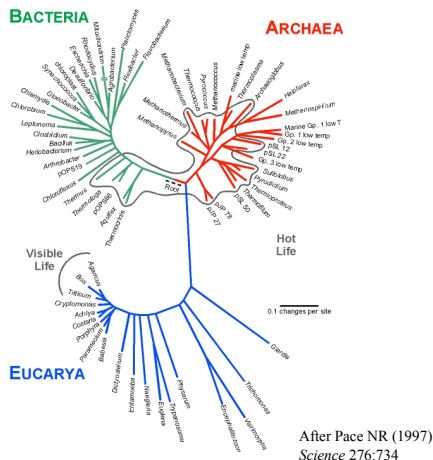


Fig. 2.1  
Page 17



Page 6

## The most sequenced organisms in GenBank

<i>Homo sapiens</i>	10.7 billion bases
<i>Mus musculus</i>	6.5b
<i>Rattus norvegicus</i>	5.6b
<i>Danio rerio</i>	1.7b
<i>Zea mays</i>	1.4b
<i>Oryza sativa</i>	0.8b
<i>Drosophila melanogaster</i>	0.7b
<i>Gallus gallus</i>	0.5b
<i>Arabidopsis thaliana</i>	0.5b

Updated 8-12-04  
GenBank release 142.0

Table 2-2  
Page 18

**www.uniprot.org**

SwissProt: 178,022 entries  
TrEMBL: 1,647,645 entries  
3-29-05 update

## PDB content growth (www.pdb.org)

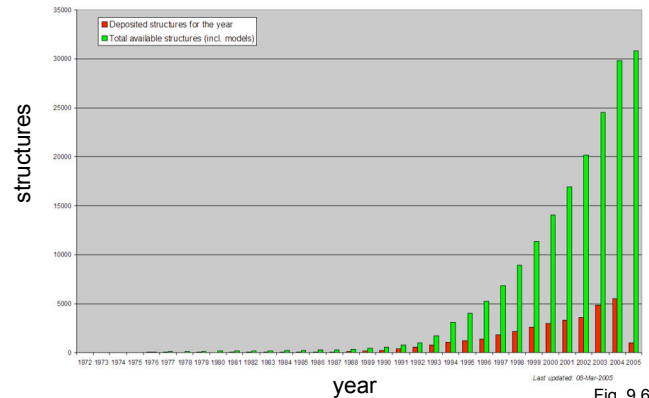


Fig. 9.6  
Page 281

## Outline: today's topic

1. How to access the sequence and structure of a protein at NCBI and the Protein Data Bank (PDB)
2. Overview of databases of all proteins: NCBI and SwissProt
3. How to align the sequences of two proteins: Dayhoff's evolutionary perspective
4. How to align the sequences of two proteins: pairwise alignment

## Definitions

Signature:

- a protein category such as a domain or motif

## Definitions

### Signature:

- a protein category such as a domain or motif

### Domain:

- a region of a protein that can adopt a 3D structure
- a fold
- a family is a group of proteins that share a domain
- examples: zinc finger domain  
immunoglobulin domain

### Motif (or fingerprint):

- a short, conserved region of a protein
- typically 10 to 20 contiguous amino acid residues

Page 225

## 15 most common domains (human)

Zn finger, C2H2 type	1093 proteins
Immunoglobulin	1032
EGF-like	471
Zn-finger, RING	458
Homeobox	417
Pleckstrin-like	405
RNA-binding region RNP-1	400
SH3	394
Calcium-binding EF-hand	392
Fibronectin, type III	300
PDZ/DHR/GLGF	280
Small GTP-binding protein	261
BTB/POZ	236
bHLH	226
Cadherin	226

Table 8-3  
Page 227

Source: Integr8 program at [www.ebi.ac.uk/proteome/](http://www.ebi.ac.uk/proteome/)

## Pairwise alignments in the 1950s

**$\beta$ -corticotropin (sheep)**      ala gly glu asp asp glu  
**Corticotropin A (pig)**      asp gly ala glu asp glu

**Oxytocin**            CYIQNCPLG  
**Vasopressin**        CYFQNCPRG

### Early alignments revealed

- differences in amino acid sequences between species
- differences in amino acids responsible for distinct functions

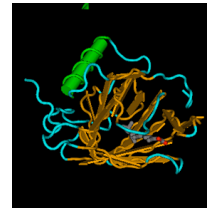
Page 40

## Pairwise sequence alignment is the most fundamental operation of bioinformatics

- It is used to decide if two proteins (or genes) are related structurally or functionally
- It is used to identify domains or motifs that are shared between proteins
- It is the basis of BLAST searching
- It is used in the analysis of genomes

Page 41

## RBP and $\beta$ -lactoglobulin are homologous proteins that share related three-dimensional structures



retinol-binding protein  
(NP\_006735)



$\beta$ -lactoglobulin  
(P02754)

Figure 3.1  
Page 42

Page 73















## The PAM250 mutation probability matrix

The PAM250 matrix is of particular interest because it corresponds to an evolutionary distance of about 20% amino acid identity (the approximate limit of detection for the comparison of most proteins).

Note the loss of information content along the main diagonal, relative to the PAM1 matrix.

## PAM250 mutation probability matrix

	A	R	N	D	C	Q	E	G	H	I	L	K	M	F	P	S	T	W	Y	V
A	13	6	9	9	5	8	9	12	6	8	6	7	7	4	11	11	11	2	4	9
R	3	17	4	3	2	5	3	2	6	3	2	9	4	1	4	4	3	7	2	2
N	4	4	6	7	2	5	6	4	6	3	2	5	3	2	4	5	4	2	3	3
D	5	4	8	11	1	7	10	5	6	3	2	5	3	1	4	5	5	1	2	3
C	2	1	1	1	52	1	1	2	2	2	1	1	1	1	2	3	2	1	4	2
Q	3	5	5	6	1	10	7	3	7	2	3	5	3	1	4	3	3	1	2	3
E	5	4	7	11	1	9	12	5	6	3	2	5	3	1	4	5	5	1	2	3
G	12	5	10	10	4	7	9	27	5	5	4	6	5	3	8	11	9	2	3	7
H	2	5	5	4	2	7	4	2	15	2	2	3	2	2	3	3	2	2	3	2
I	3	2	2	2	2	2	2	2	10	6	2	6	5	2	3	4	1	3	9	9
L	6	4	4	3	2	6	4	3	5	15	34	4	20	13	5	4	6	6	7	13
K	6	18	10	8	2	10	8	5	8	5	4	24	9	2	6	8	8	4	3	5
M	1	1	1	1	0	1	1	1	1	1	2	3	2	6	2	1	1	1	1	2
F	2	1	2	1	1	1	1	1	3	5	6	1	4	32	1	2	2	4	20	3
P	7	5	5	4	3	5	4	5	5	3	3	4	3	2	20	6	5	1	2	4
S	9	6	8	7	7	6	7	9	6	5	4	7	5	3	9	10	9	4	4	6
T	8	5	6	6	4	5	5	6	4	6	4	6	5	3	6	8	11	2	3	6
W	0	2	0	0	0	0	0	0	1	0	1	0	0	1	0	1	0	55	1	0
Y	1	1	2	1	3	1	1	1	3	2	2	1	2	15	1	2	2	3	31	2
V	7	4	4	4	4	4	4	5	4	15	10	4	10	5	5	5	7	2	4	17

Top: original amino acid

Side: replacement amino acid

Fig. 3.13

Page 57

A	2																			
R	-2	6																		
N	0	0	2																	
D	0	-1	2	4																
C	-2	-4	-4	-5	12															
Q	0	1	1	2	-5	4														
E	0	-1	1	3	-5	2	4													
G	1	-3	0	1	-3	-1	0	5												
H	-1	2	2	1	-3	3	1	-2	6											
I	-1	-2	-2	-2	-2	-2	-2	-3	-2	5										
L	-2	-3	-3	-4	-6	-2	-3	-4	-2	-2	6									
K	-1	3	1	0	-5	1	0	-2	0	-2	-3	5								
M	-1	0	-2	-3	-5	-1	-2	-3	-2	2	4	0	6							
F	-3	-4	-3	-6	-4	-5	-5	-5	-2	1	2	-5	0	9						
P	1	0	0	-1	-3	0	-1	0	0	-2	-3	-1	-2	-5	6					
S	1	0	1	0	0	-1	0	1	-1	-1	-3	0	-2	-3	1	2				
T	1	-1	0	0	-2	-1	0	0	-1	0	-2	0	-1	-3	0	1	3			
W	-6	2	-4	-7	-8	-5	-7	-7	-3	-5	-2	-3	-4	0	-6	-2	-5	17		
Y	-3	-4	-2	-4	0	-4	-4	-5	0	-1	-1	-4	-2	7	-5	-3	-3	0	10	
V	0	-2	-2	-2	-2	-2	-2	-1	-2	4	2	-2	2	-1	-1	-1	0	-6	-2	4
A																				
R																				
N																				
D																				
C																				
Q																				
E																				
G																				
H																				
I																				
L																				
K																				
M																				
F																				
P																				
S																				
T																				
W																				
Y																				
V																				

PAM250 log odds scoring matrix

Fig. 3.14  
Page 58

## Why do we go from a mutation probability matrix to a log odds matrix?

- We want a scoring matrix so that when we do a pairwise alignment (or a BLAST search) we know what score to assign to two aligned amino acid residues.
- Logarithms are easier to use for a scoring system. They allow us to sum the scores of aligned residues (rather than having to multiply them).

## How do we go from a mutation probability matrix to a log odds matrix?

- The cells in a log odds matrix consist of an "odds ratio":

$\frac{\text{the probability that an alignment is authentic}}{\text{the probability that the alignment was random}}$

The score S for an alignment of residues a,b is given by:

$$S(a,b) = 10 \log_{10} (M_{ab}/p_b)$$

As an example, for tryptophan,

$$S(a,\text{tryptophan}) = 10 \log_{10} (0.55/0.010) = 17.4$$

## What do the numbers mean in a log odds matrix?

$$S(a,\text{tryptophan}) = 10 \log_{10} (0.55/0.010) = 17.4$$

A score of +17 for tryptophan means that this alignment is 50 times more likely than a chance alignment of two Trp residues.

$$S(a,b) = 17$$

$$\text{Probability of replacement } (M_{ab}/p_b) = x$$

Then

$$17 = 10 \log_{10} x$$

$$1.7 = \log_{10} x$$

$$10^{1.7} = x = 50$$



Two randomly diverging protein sequences change in a negatively exponential fashion

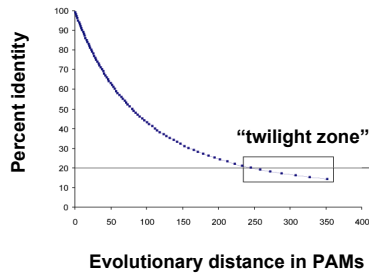


Fig. 3.19  
Page 62

At PAM1, two proteins are 99% identical  
At PAM10.7, there are 10 differences per 100 residues  
At PAM80, there are 50 differences per 100 residues  
At PAM250, there are 80 differences per 100 residues

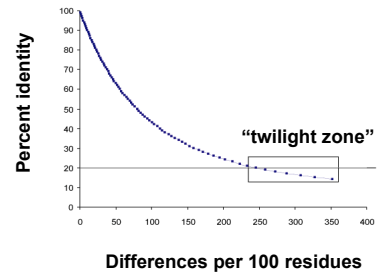
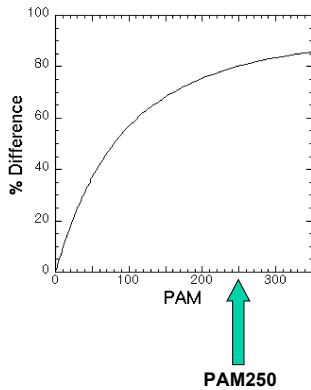


Fig. 3.19  
Page 62

PAM matrices reflect different degrees of divergence



### PAM: "Accepted point mutation"

- Two proteins with 50% identity may have 80 changes per 100 residues. (Why? Because any residue can be subject to back mutations.)
- Proteins with 20% to 25% identity are in the "twilight zone" and may be statistically significantly related.
- PAM or "accepted point mutation" refers to the "hits" or matches between two sequences (Dayhoff & Eck, 1968)

### Ancestral sequence

ACCCTAC

A	no change	A
C	single substitution	C --> A
C	multiple substitutions	C --> A --> T
C --> G	coincidental substitutions	C --> A
T --> A	parallel substitutions	T --> A
A --> C --> T	convergent substitutions	A --> T
C	back substitution	C --> T --> C

Sequence 1  
ACCGATC

Sequence 2  
AATAATC

Fig. 11.11  
Page 374

### Percent identity between two proteins: What percent is significant?

100%  
80%  
65%  
30%  
23%  
19%

## Outline: today's topic

1. How to access the sequence and structure of a protein at NCBI and the Protein Data Bank (PDB)
2. Overview of databases of all proteins: NCBI and SwissProt
3. How to align the sequences of two proteins: Dayhoff's evolutionary perspective
4. How to align the sequences of two proteins: pairwise alignment

## General approach to pairwise alignment

- Choose two sequences
- Select an algorithm that generates a score
- Allow gaps (insertions, deletions)
- Score reflects degree of similarity
- Alignments can be global or local
- Estimate probability that the alignment occurred by chance

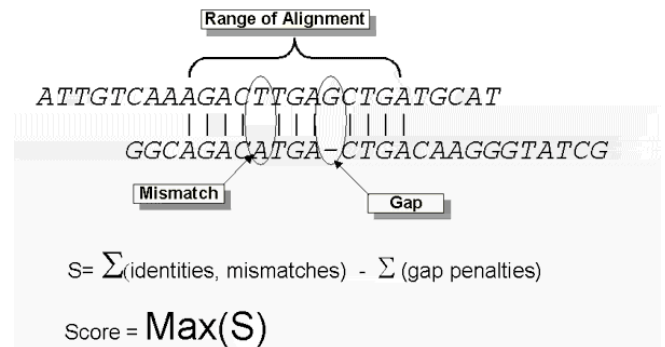
Page 50

## An alignment scoring system is required to evaluate how good an alignment is

- positive and negative values assigned
- gap creation and extension penalties
- positive score for identities
- some partial positive score for conservative substitutions
- global versus local alignment
- use of a substitution matrix

Page 62

## Calculation of an alignment score



[http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/Alignment\\_Scores2.html](http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/Alignment_Scores2.html)

## Two kinds of sequence alignment: global and local

We will first consider the global alignment algorithm of Needleman and Wunsch (1970).

We will then explore the local alignment algorithm of Smith and Waterman (1981).

Finally, we will consider BLAST, a heuristic version of Smith-Waterman.

Page 63

## Global alignment with the algorithm of Needleman and Wunsch (1970)

- Two sequences can be compared in a matrix along x- and y-axes.
- If they are identical, a path along a diagonal can be drawn
- Find the optimal subpaths, and add them up to achieve the best score. This involves
  - adding gaps when needed
  - allowing for conservative substitutions
  - choosing a scoring system (simple or complicated)
- N-W is guaranteed to find optimal alignment(s)

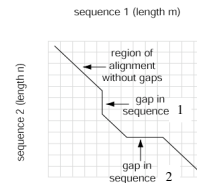
Page 63



## Three steps to global alignment with the Needleman-Wunsch algorithm

- [1] set up a matrix
- [2] score the matrix
- [3] identify the optimal alignment(s)

## Four possible outcomes in aligning two sequences



- [1] identity (stay along a diagonal)
- [2] mismatch (stay along a diagonal)
- [3] gap in one sequence (move vertically!)
- [4] gap in the other sequence (move horizontally!)

Fig. 3.20  
Page 64

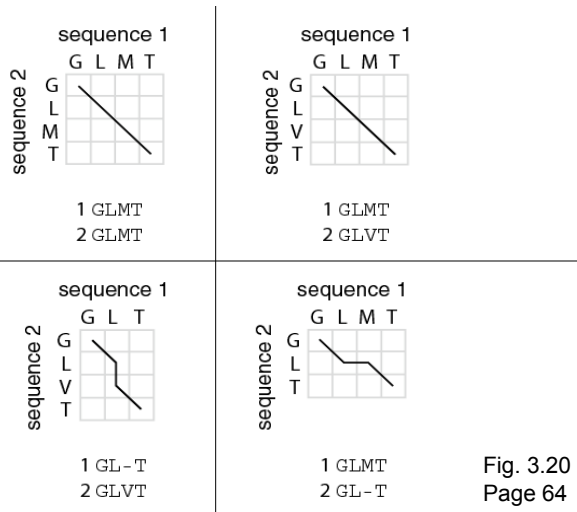


Fig. 3.20  
Page 64

## Start Needleman-Wunsch with an identity matrix

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1												
J		1											
C			1										
N				1									
R					1								
Q						1							
C							1						
L								1					
C									1				
R										1			
P											1		

Fig. 3.21  
Page 65

## Start Needleman-Wunsch with an identity matrix

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1												
J		1											
C			1										
N				1									
R					1								
Q						1							
C							1						
L								1					
C									1				
R										1			
P											1		

sequence 1 **A**BCNJ-RQCLCR-PM  
sequence 2 **A**JC-JNR-CKCRBP-

sequence 1 **A**BC-NJRQCLCR-PM  
sequence 2 **A**JCJN-R-CKCRBP-

Fig. 3.21  
Page 65

## Fill in the matrix starting from the bottom right

	A	B	C	N	J	R	Q	C	L	C	R	P	M
A	1												
J		1											
C			1										
N				1									
R					1								
Q						1							
C							1						
L								1					
C									1				
R										1			
P											1		

Fig. 3.21  
Page 65

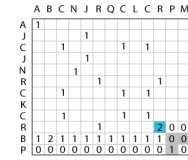
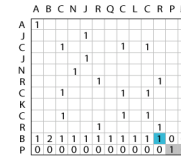
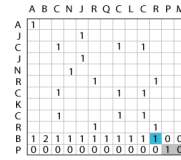
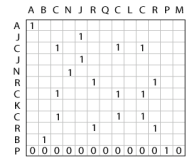
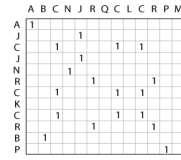
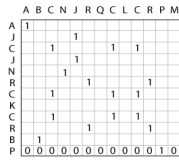
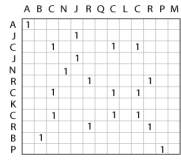


Fig. 3.21  
Page 65

Fig. 3.21  
Page 65

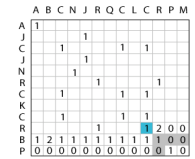
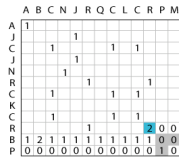
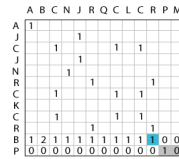
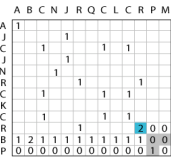
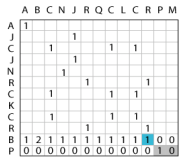
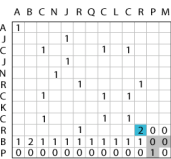
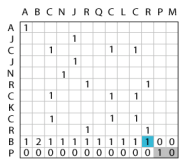
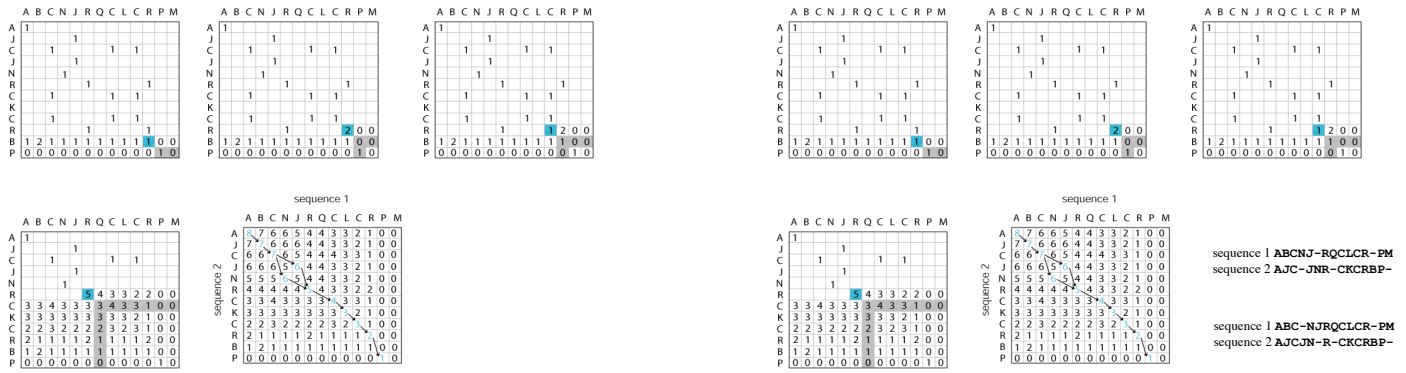


Fig. 3.21  
Page 65

Fig. 3.21  
Page 65





After you've filled in the matrix, find the optimal path(s) by a "traceback" procedure

Fig. 3.22  
Page 66

## Needleman-Wunsch: dynamic programming

N-W is guaranteed to find optimal alignments, although the algorithm does not search all possible alignments.

It is an example of a dynamic programming algorithm: an optimal path (alignment) is identified by incrementally extending optimal subpaths. Thus, a series of decisions is made at each step of the alignment to find the pair of residues with the best score.

```
> gap
Gap uses the algorithm of Needleman and Wunsch to find the alignment of
two complete sequences that maximizes the number of matches and minimizes
the number of gaps.

GAP of what sequence 1 ? hsrbp.pep
  Begin (* 1 *) ?
  End (* 199 *) ?

to what sequence 2 (* hsrbp.pep *) ? btllacto.pep
  Begin (* 1 *) ?
  End (* 178 *) ?

What is the gap creation penalty (* 8 *) ?
What is the gap extension penalty (* 2 *) ?
What should I call the paired output display file (* hsrbp.pair *) ?

Aligning .....-
Aligning .....-

          Gaps:      8
          Quality:    37
Quality Ratio: 0.208
% Similarity: 31.902
          Length:    214
```

Fig. 3.23  
Page 68

```
Gap Weight:      8      Average Match: 2.912
Length Weight:   2      Average Mismatch: -2.003

Quality:         37      Length:      214
Ratio:          0.208    Gaps:       8
Percent Similarity: 31.902  Percent Identity: 26.300

Match display thresholds for the alignment(s):
: = IDENTITY
: = 2
: = 1

hsrbp.pep x btllacto.pep  July 16, 2001 14:45 ..

1  MKUWUALLLLAAWAAARDCRUSSFUKEHFKRAFSGTUYAMAKKDEG 50
   |||  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
1  ...MKCLLLALALTCGAQALIUT..QTKMGLDIKUAGTIYSLAMASD..44
51  LFLQDNIVAREFSUDETGQMSATAGKAVR..LLNND..UCADMUGTFTDTE 97
   |||  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
45  ISLLDQASPLRV..YUEELKPTPEGOLEILLQKWENGECAQKIIAEKTK 93
98  DPAKFKMKYVUGVASFQKGNDDHVIUDTDYDYAV.....QYSC 136
   |||  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
94  IPAUFKIDALHENKUL.....ULDTVKKYKYLFCMENSAREPEQSLAC 135
137  RLLNLGDTCADSYSFUSADPNGLPPEAQKIURAQ..EELCLARQYALIU 185
   |||  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |  |
136  QCLURTPEVDDALEKFKALKALPMHIALSFNPTQLEEQCHI..... 178
```

Fig. 3.24  
Page 69

```
> bestfit
BestFit makes an optimal alignment of the best segment of similarity
between two sequences. Optimal alignments are found by inserting gaps to
maximize the number of matches using the local homology algorithm of
Smith and Waterman.

BESTFIT of what sequence 1 ? hsrbp.pep
  Begin (* 1 *) ?
  End (* 199 *) ?

to what sequence 2 (* hsrbp.pep *) ? btllacto.pep
  Begin (* 1 *) ?
  End (* 178 *) ?

What is the gap creation penalty (* 8 *) ?
What is the gap extension penalty (* 2 *) ?
What should I call the paired output display file (* hsrbp.pair *) ?

Aligning .....-
Aligning .....-

          Gaps:      5
          Quality:    59
Quality Ratio: 0.621
% Similarity: 39.130
          Length:    105
```

Fig. 3.26  
Page 71



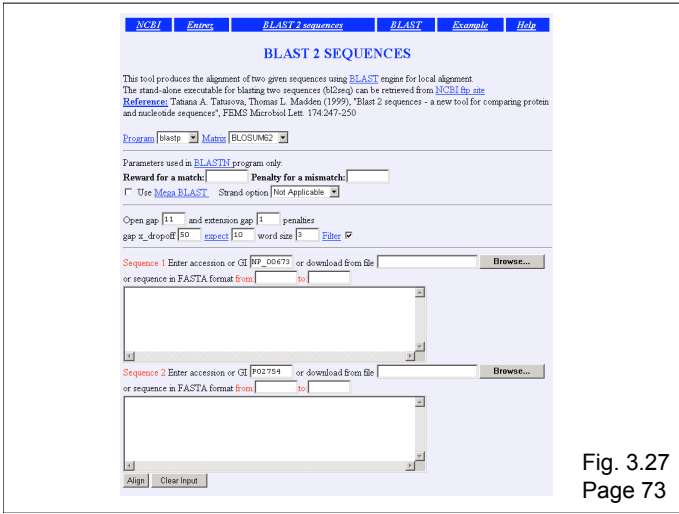


Fig. 3.27  
Page 73

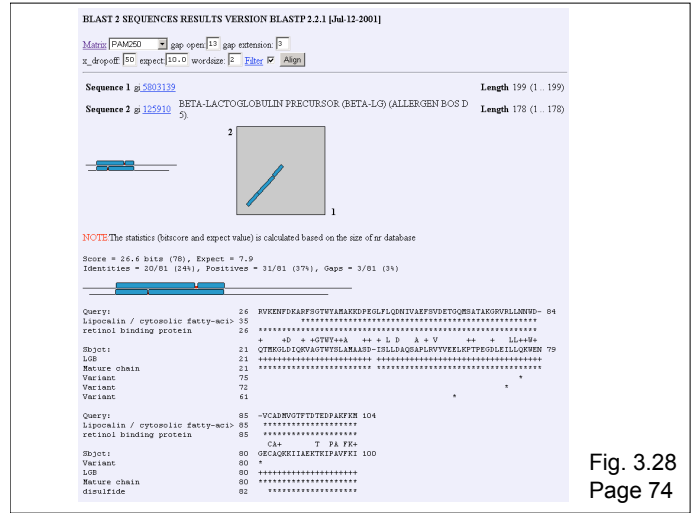


Fig. 3.28  
Page 74

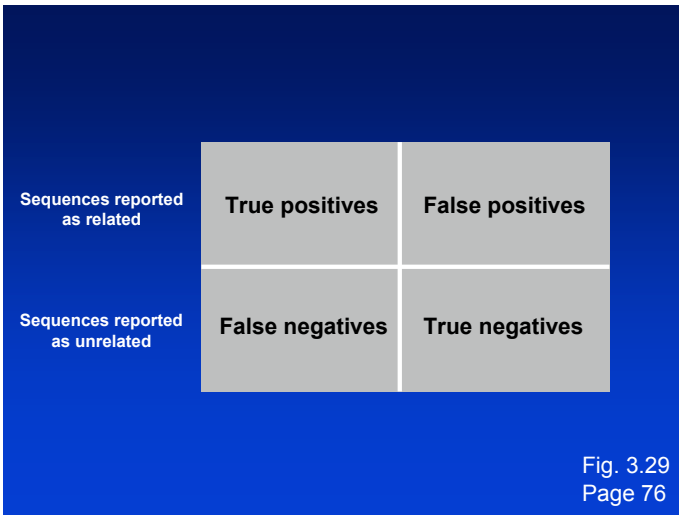


Fig. 3.29  
Page 76

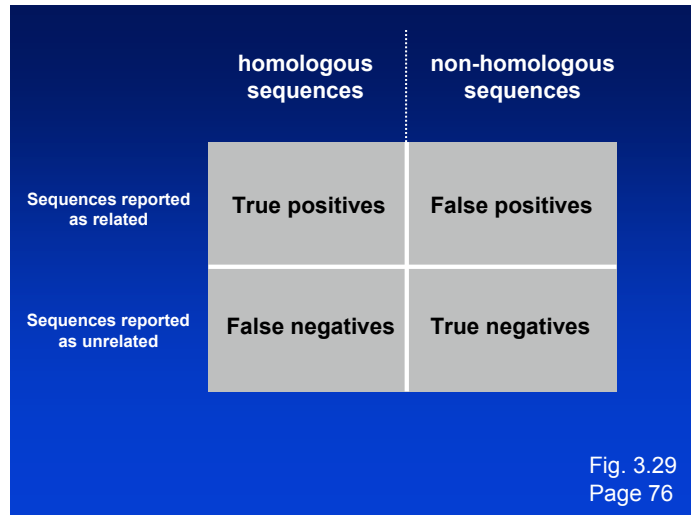


Fig. 3.29  
Page 76

