

# Bioinformatics

Jiří Damborský  
National Center for Biomolecular Research

jiri@chemi.muni.cz, ph. 41129 377, Kottlarska 2, bld. 7, 2nd floor

---

---

---

---

---

---

---

---

## Bioinformatics - what is it?

- The term bioinformatics is used to encompass almost all computer applications in biological sciences.
- Information technology applied to the management and analysis of biological data
- originally - analysis of sequence data (80s)
- presently - also analysis of 3D-structures

---

---

---

---

---

---

---

---

## Bioinformatics - study material

- Introduction to bioinformatics, T.K. Attwood and D.J. Parry-Smith, Longman, Essex, 1999.
- copy of the slides
- <http://www.chemi.muni.cz/~jiri>
- <http://www.bioinf.man.ac.uk/dbbrowser/bioactivity/prefacefrm.html>

---

---

---

---

---

---

---

---

## Bioinformatics - composition

12 lectures per semester

3 hours per week

- 1<sup>st</sup> and 2<sup>nd</sup> hour = lectures - theory
- 3<sup>rd</sup> hour = practical course on computers

---

---

---

---

---

---

---

---

## Bioinformatics - lectures

- Introduction
- Information networks
- Protein information resources
- Genome information resources
- DNA sequence analysis
- Pairwise sequence alignment
- Multiple sequence alignment
- Secondary database searching
- Analysis packages
- Protein structure modelling

---

---

---

---

---

---

---

---

## Bioinformatics - practical training

- Biological databases
- Searching and modelling servers
- Building a sequence search protocol
- Case examples
- Protein structure prediction
- Protein modelling
  
- Follow-up of lectures

---

---

---

---

---

---

---

---

## Introduction

- history of sequencing
- what is it Bioinformatics?
- sequence to structure deficit
- genome projects
- why is Bioinformatics important?
- patten recognition and prediction
- folding problem
- sequence analysis
- homo/analogy and ortho/paralogy

---

---

---

---

---

---

---

---

## History of sequencing

- Protein sequencing
  - > separation of peptides, identification and quantification of amino acids
  - > Edman degradation
  - > mass-spectrometry - advantage in identification of post-translational modifications
  - > 1955 sequencing of peptide insuline
  - > 1960 sequencing of enzyme ribonuclease
  - > 1980s automated sequencers

---

---

---

---

---

---

---

---

## History of sequencing

- Nucleic acid sequencing
  - > tRNA - short, could be purified
  - > DNA - large (human chromosome 55-250 x 10<sup>6</sup> bp); the longest fragment for sequencing is 500 bp; purification is problematic
  - > advent of gene cloning and PCR
  - > 1972 DNA cloning
  - > 1975 DNA sequencing
  - > 1980s and 1990s sequence revolution

---

---

---

---

---

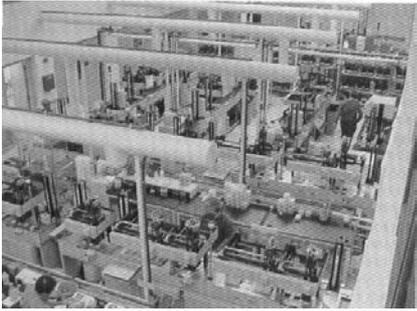
---

---

---



### Automated production line in sequencing "factory"



Whitehead Institute, Center for Genome Research, USA

---

---

---

---

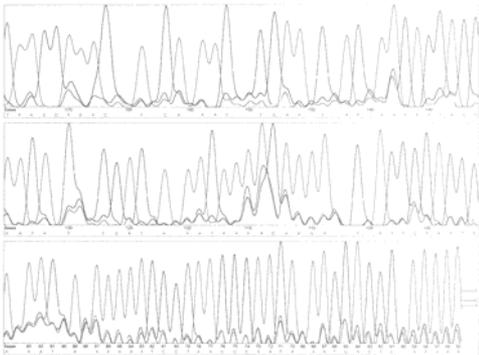
---

---

---

---

### Sequencing chromatogram



---

---

---

---

---

---

---

---

### What is Bioinformatics?

- improvements in DNA sequencing technologies and computer-based technologies
- originally - analysis of sequence data (1980s)
- presently - also analysis of 3D-structures
- The term bioinformatics is used to encompass almost all computer applications in biological sciences.
- Information technology applied to the management and analysis of biological data.

---

---

---

---

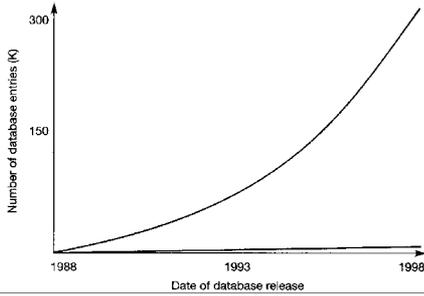
---

---

---

---

### The sequence to structure deficit




---

---

---

---

---

---

---

---

### Genome projects

- 1977 first complete genome - virus  $\phi$ X174, 5000 nucleotides; 11 genes
- 1995 first complete genome of living organism *Haemophilus influenzae*, 1.8 million nucleotides and 1700 genes
- sequencing of model systems: *Escherichia coli*, *Saccharomyces cerevisiae*, *Cernorhabditis elegans*, *Drosophila melanogaster*, *Arabidopsis thaliana*, *Canis familiaris*, *Mus Musculus*

---

---

---

---

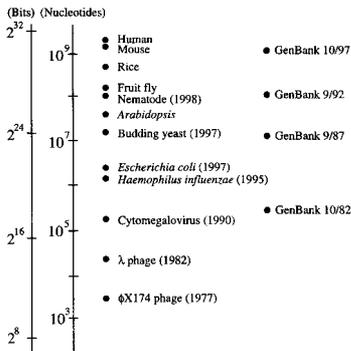
---

---

---

---

### The genome size of various species




---

---

---

---

---

---

---

---

## Comparative genomic analysis of model organisms

	Genome size (Mb)	Gene number	Haploid chromosome number
Bacterium ( <i>Escherichia coli</i> )	~4	4,403	1
Yeast ( <i>Saccharomyces cerevisiae</i> )	~12	6,190	16
Worm ( <i>Caenorhabditis elegans</i> )	97	19,730	6
Fruit Fly ( <i>Drosophila melanogaster</i> )	120	13,601	4
Mouse ( <i>Mus Musculus</i> )	3,454	~50,000 (estimated)	20
Human ( <i>Homo sapiens</i> )	2,910	33,609	23

---

---

---

---

---

---

---

---

## Human Genome Project

- in mid-1980s initiated Human Genome Project
- estimated 100.000 genes and completion in 2005
- need for automated sequencing and improved computational techniques
- shotgun method
  
- sequencing of rough draft first
- first draft completed in 2000 by publicly funded the International Consortium Human Genome Project and the company Celera Genomics

---

---

---

---

---

---

---

---

## Human Genome Project

- ~33.000 genes
- genes are complex due to alternative splicing
- >1.000.000 proteins (estimated)
- hundreds of genes resulted from horizontal transfer from bacteria (in vertebrate lineage)
- dozen of genes derived from transposable elements (their activity however has declined)
- the mutation rate in male is two-times higher than in female
- >1.400.000 single point polymorphisms (SNPs)

---

---

---

---

---

---

---

---

## Why is bioinformatics important?

- last 20-30 years - structural biology
- new era - bioinformatics - due to genome projects and sequence/structure deficit
- biological function is not known for about 50% of all genes in every sequenced genome
- role of bioinformatics
  - > data management and storage
  - > data analysis = conversion of primary sequence to biological knowledge

---

---

---

---

---

---

---

---

## Pattern recognition *versus* prediction

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



---

---

---

---

---

---

---

---

## Levels of protein structure

- Primary structure:** the linear sequence of amino acids in a protein molecule
- Secondary structure:** regions of local regularity within a protein fold (e.g.,  $\alpha$ -helices,  $\beta$ -turns,  $\beta$ -strands)
- Super-secondary structure:** the arrangement of  $\alpha$ -helices and/or  $\beta$ -strands into discrete folding units (e.g.,  $\beta$ -barrels,  $\beta\alpha\beta$ -units, Greek keys, etc.)
- Tertiary structure:** the overall fold of a protein sequence, formed by the packing of its secondary and/or super-secondary structure elements
- Quaternary structure:** the arrangement of separate protein chains in a protein molecule with more than one subunit
- Quinternary structure:** the arrangement of separate molecules, such as in protein-protein or protein-nucleic acid interactions

---

---

---

---

---

---

---

---

## Homology and analogy

- Sequences are said to be homologous if they are related by divergence from a common ancestor.
- Proteins can share similar folds (e.g.,  $\beta$ -barrel) or similar catalytic residues (e.g., serine proteases) without any sequential similarity. Convergence to similar biological solutions from different evolutionary starting points results in analogy.
- Sequence analysis assumes homologous proteins.
- Homology is not a measure of similarity.

---

---

---

---

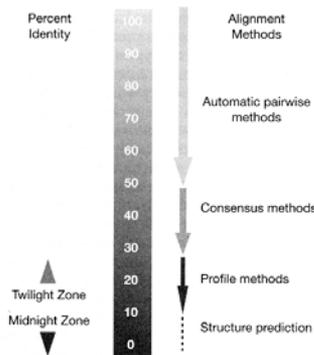
---

---

---

---

## Application areas of different analysis methods



---

---

---

---

---

---

---

---

## Orthology and paralogy

- Proteins performing the same function in different species - orthologues.
- Proteins performing different, but related functions within same organism - paralogues.
- Sequence comparison of orthologous proteins - phylogenetic analysis.

---

---

---

---

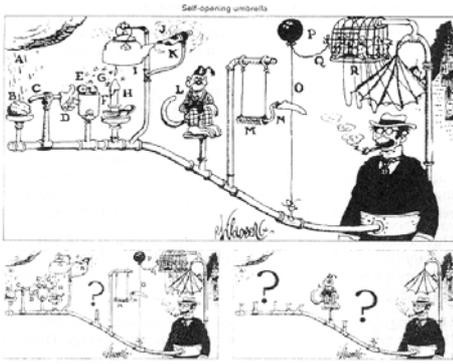
---

---

---

---

Modularity of proteins = difficulties of homology searches



---

---

---

---

---

---

---

---

## Information networks

- what is the Internet?
- how do computers find each other?
- FTP and Telnet
- what is the Worl Wide Web?
- HTTP, HTML and URL
- EMBnet, EBI, NCBI
- SRS and ENTREZ

---

---

---

---

---

---

---

---

## What is the Internet?

- Global network of computer networks that link government, academic and business institutions.
- communication by TCP/IP  
(Transmission Control Protocol/Internet Protocol)
- computers - nodes, data - packets
- packets may not be transferred directly from one computer to another

---

---

---

---

---

---

---

---

## How do computers find each other?

- Each computer is assigned IP address  
147.251.28.2  
machine.site.domain  
bilbo.chemi.muni.cz
- FTP - File Transfer Protocol
- Telnet - remote connection

---

---

---

---

---

---

---

---

## Example of Internet domains and subdomains

<i>Country-based domains</i>	<i>Other domains</i>	<i>Subdomains</i>
Australia	.au	Educational .edu Academic .ac
Denmark	.dk	Commercial .com Company .co
Finland	.fi	Governmental .gov Other organisation .org
France	.fr	Military .mil General .gen
Germany	.de	
Greece	.gr	
Hungary	.hu	
Ireland	.ie	
Israel	.il	
Italy	.it	
Netherlands	.nl	
New Zealand	.nz	
Poland	.pl	
Portugal	.pt	
South Africa	.za	
Spain	.es	
Sweden	.se	
Switzerland	.ch	
United Kingdom	.uk	
USA	.us	

---

---

---

---

---

---

---

---

---

---

## What is the World Wide Web?

- Developed at CERN - the European Laboratory of Particle Physics.
- The purpose was sharing of information.
- Hypermedia based information system.
- The most advanced information system found on the web.
- Very popular - almost synonymous with the Internet.

---

---

---

---

---

---

---

---

---

---

## Web browsers

- Browser is the client communicating with servers using standard protocols.
- Home page is the first point of contact between browser and the server.

Lynx - academic, VT100 terminal  
Mosaic - academic, X-windows  
Netscape Navigator - commercial  
Internet Explorer - commercial

---

---

---

---

---

---

---

---

---

---

## HTTP, HTML and URL

- HTTP - HyperText Transport Protocol  
documents exploited by browsers are written in hypertext and transferred by HTTP
- HTML - HyperText markup Language  
standard language for writing a hypertext
- URL - Uniform Resource Locator  
unique address for a document  
example: <http://www.chemi.muni.cz/~jiri>

---

---

---

---

---

---

---

---

## EMBnet, EBI, NCBI

- 1988 established the network of European biocomputing and bioinformatics laboratories.
- Eliminates the need for multicopies of biology databases and retrieval software.
- Hinxton Hall = Sanger Centre + MRC Human Genome Mapping Project Resource Centre + European Bioinformatics Institute (EBI)
- National Center for Biotechnology Information (NCBI)

---

---

---

---

---

---

---

---

## SRS, ENTREZ and LinkDB

### SRS - The Sequence Retrieval System

- > maintained by EBI
- > network browser for databases in molecular biology
- > allows indexation of flat-file databases
- > allows customised search of selected databases
- > link databanks: sequence, structure, bibliography, etc.

### ENTREZ

- > integrates databases of NCBI
- > less flexible than SRS
- > valuable concept of neighbouring
- > link databanks: DNA and protein sequences, genome data, structural data, PubMed bibliography

---

---

---

---

---

---

---

---

# SRS, ENTREZ and LinkDB

## LinkDB

- > maintained by Institute for Chemical Research, Japan
- > network browser for databases in DBGET and KEGG (Kyoto encyclopedia of genes and genomes)
- > link databanks: sequence, motifs, structure, amino acid properties, ligands, metabolic pathways

---

---

---

---

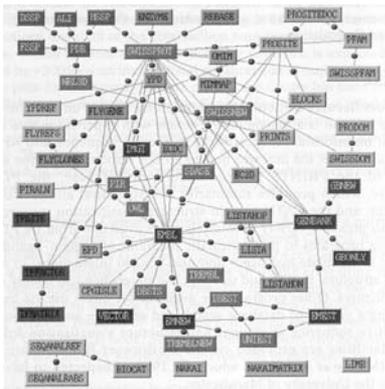
---

---

---

---

Network of databases linked via SRS



---

---

---

---

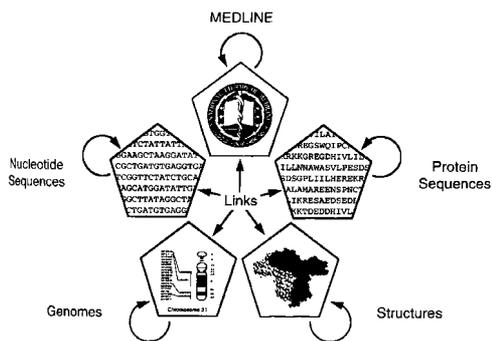
---

---

---

---

Network of databases linked via ENTREZ



---

---

---

---

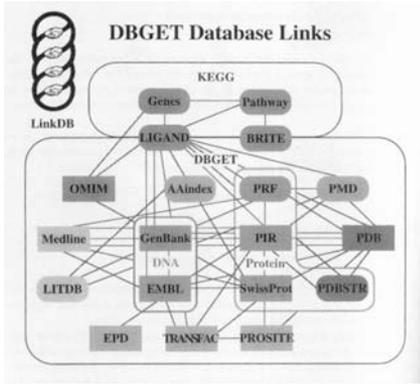
---

---

---

---

# Network of databases linked via LinkDB



---

---

---

---

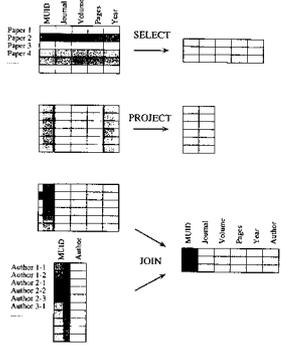
---

---

---



## Relational database




---

---

---

---

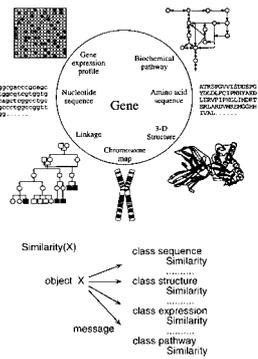
---

---

---

---

## Object-oriented database




---

---

---

---

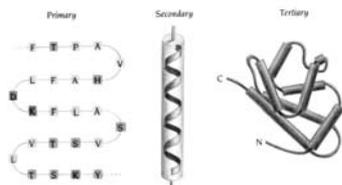
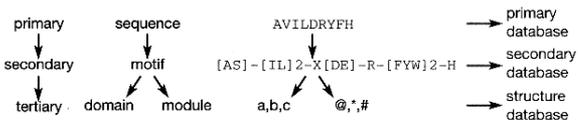
---

---

---

---

## Levels of protein structure and corresponding databases




---

---

---

---

---

---

---

---

## Primary protein sequence databases

- PIR
- MIPS
- SWISS-PROT
- TrEMBL
- NRL-3D

Store biomolecular sequences and annotations.

---

---

---

---

---

---

---

---

## Primary protein sequence databases

- PIR - Protein Sequence Database
  - > 1960s by Margaret Dayhoff
  - > maintained by international consortium
  - > four sections PIR1-PIR4
    - PIR1 - fully classified and annotated entries
    - PIR2 - preliminary entries
    - PIR3 - unverified entries
    - PIR4 - conceptual translations of artefactual sequences, non-transcribed, non-translated
- MIPS - Martinsried Institute for Protein Sequences
  - > collects and processes sequence data for PIR

---

---

---

---

---

---

---

---

## Primary protein sequence databases

- SWISS-PROT
  - > University Geneva → EBI → Swiss Inst. of Bioinformatics
  - > high-level annotations including description of the function, structure and domains, post-translational modifications, variants, etc.
  - > annotated manually (high quality)
  - > automatically annotated = TrEMBL
  - > minimally redundant
  - > interlinked with many other sources
  - > efficient searching of selected fields only
  - > most widely used protein sequences database

---

---

---

---

---

---

---

---



## Composite protein sequence databases

### ■ NRDB - Non-Redundant DataBase

- > developed and maintained by NCBI
- > composite: GenPept (CDS translations of GenBank), GenPeptupdate, PDB sequences, SWISS-PROT, SWISS-PROTupdate, RIR
- > advantages: comprehensive and up-to date
- > disadvantages: not fully redundant (only identical copies removed), occurrence of multiple entries due to polymorphism, incorrect sequences amended in SWISS-PROT re-introduced by translation of GenBank
- > default database of the NCBI BLAST (ENTREZ/NCBI)

---

---

---

---

---

---

---

---

## Composite protein sequence databases

### ■ OWL

- > developed and maintained by University of Leeds
- > composite: SWISS-PROT, PIR1-4, GenBank, NRL-3D
- > SWISS-PROT the highest priority for annotation
- > advantages: less redundant, fully indexed (fast)
- > disadvantages: not up-to-date (released every 6-8 weeks), incorrect sequences
- > available from SEQNET of UK EMBnet

---

---

---

---

---

---

---

---

## Composite protein sequence databases

### ■ MIPSX

- > developed by Max-Planck Institute in Martinsried
- > composite: PIR1-4, MIPS, NRL-3D, SWISS-PROT, TrEMBL, GenPept, Kabat, PSeqIP
- > identical entries and subsequences removed

### ■ SWISS-PROT+TrEMBL

- > developed and maintained by EBI
- > composite: SWISS-PROT, TrEMBL
- > advantages: comprehensive, minimally redundant, fewer errors
- > disadvantages: not as up-to-date as NRDB
- > available from SRS of EBI

---

---

---

---

---

---

---

---

## Overview of primary sources of composite databases

<i>NRDB</i>	<i>OWL</i>	<i>MIP SX</i>	<i>SP + TrEMBL</i>
PDB	SWISS-PROT	PIR1-4	SWISS-PROT
SWISS-PROT	PIR	MIPSOwn	TrEMBL
PIR	GenBank	MIPSTrn	
GenPept	NRL-3D	MIPSH	
SWISS-PROTupdate		PIRMOD	
GenPeptupdate		NRL-3D	
		SWISS-PROT	
		EMTrans	
		GBTrans	
		Kabat	
		PseqIP	

---

---

---

---

---

---

---

---

## Secondary databases

- Contains information derived from primary sequence data, typically in the form of abstractions: regular expressions, fingerprints, blocks, profiles or Hidden Markov Models.
- These abstractions represent distillations of the most conserved features of multiple alignments.
- The abstractions are useful for discrimination of family membership for newly determined sequences.

---

---

---

---

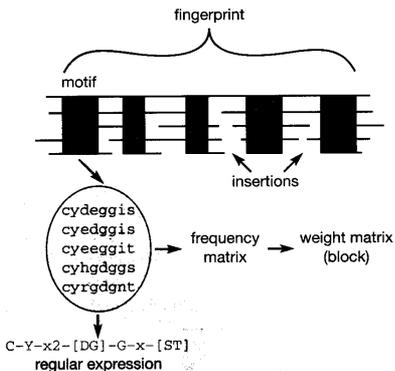
---

---

---

---

## Terms used in sequence analysis methods




---

---

---

---

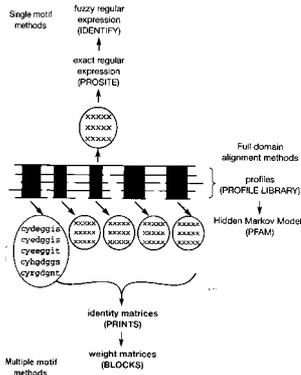
---

---

---

---

## Three principal methods for building secondary databases




---

---

---

---

---

---

---

---

---

---

---

---

## Implication of function from a sequence

### DNA-binding proteins

Name	Helix-loop-helix (Myc type)	Cys-His zinc finger	Leucine zipper
Sequence	[DENSTAP]K-[LYVMWAGN]- [FYWCHKR]-[LIVT]-[LV]-a(2)- [STAV]-[LV]MSTAC]-x-[VMPYH]- [LV]MTA)-[F]-[F]-[LV]NSR	C-x(2,4)-C-x(3)-[LV]MPFYWC]- x(5)-H-x(3,5)-H	L-x(6)-L-x(6)-L-x(6)-L
Structure			
Function	DNA Binding	DNA Binding	DNA Binding
Example	 3CRO	 2DRP	 1YSA

---

---

---

---

---

---

---

---

---

---

---

---

## Secondary databases

- PROSITE
- PRINTS
- BLOCKS
- Profiles
- Pfam
- IDENTIFY

---

---

---

---

---

---

---

---

---

---

---

---

## Secondary databases

### ■ PROSITE

- > historically the first secondary database
- > maintained by Swiss Institute of Bioinformatics
- > motivation: identification of protein families
- > abstraction: regular expressions (patterns)
- > construction: automatic multiple alignment and manual extraction of conserved regions
- > ideally patterns should identify only true-positives (not false-positives)
- > entries deposited as two distinct files: pattern file and documentation files
- > primary source: SWISS-PROT

### Pattern file of an entry from the PROSITE database

```
ID OPSIN_PATTERN.
AC P040238;
DF APR-1990 (CREATED); NOV-1997 (DATA UPDATE); NOV-1997 (INFO UPDATE).
DE Visual pigments (opsins) retinal binding site.
FA [LIVNM]-[PGC]-x(3)-[SAC]-K-[STALIM]-[GSACHYV]-x(2)-[IDNYT]-[AP]-
FA x(2)-[LV]-.
NR /RELEASE=32,49340;
NR /TOTAL=53(53) /POSITIVE=53(53) /UNKNOWN=0(0) /FALSE_POS=0(0);
NR /FALSE_NEG=0 /PARTIAL=1;
CC /TAXO-RANGE=?*?/? /MAX-REPEAT=1;
CC /SITE=5,retinal1;
DE P06602, OPS1_DROME, T; P28678, OPS1_DROPS, T; P22269, OPS1_CALV1, T;
DE P08699, OPS2_DROME, T; P28679, OPS2_DROPS, T; P04950, OPS3_DROME, T;
DE P28689, OPS3_DROPS, T; P08255, OPS4_DROME, T; P29404, OPS4_DROPS, T;
DE P17644, OPS4_DROV1, T; P35362, OPS5_SPHSP, T; P41591, OPS6_ANOCA, T;
DE P41590, OPS6_ASTFA, T; P03499, OPS2_BOVIN, T; P33308, OPS2_CANFA, T;
DE P32309, OPS2_CARAU, T; P22328, OPS2_CHICK, T; P26681, OPS2_CHIGH, T;
DE P08100, OPS2_HUMAN, T; P15409, OPS2_MOUSE, T; P35403, OPS2_POMMI, T;
DE P02700, OPS2_SHEEP, T; P29403, OPS2_XENLA, T; P22671, OPS2_LAMCA, T;
DE P31355, OPS2_BAFT1, T; P24661, OPS2_LOLPO, T; P05241, OPS2_OCTEO, T;
DE P35356, OPS2_PROCL, T; P31356, OPS2_PODFA, T; P35360, OPS1_LIMPO, T;
DE P35361, OPS2_LIMPO, T; P23210, OPSB_CARAU, T; P26682, OPSB_CHICK, T;
DE P35357, OPSB_GECOE, T; P03999, OPSB_HUMAN, T; P26684, OPSV_CHICK, T;
DE P22330, OPS6_ASTFA, T; P22331, OPS6_ASTFA, T; P23211, OPS6_CARAU, T;
DE P32312, OPS6_CARAU, T; P28683, OPS2_CHICK, T; P35359, OPS2_GECOE, T;
DE P04001, OPSG_HUMAN, T; P41592, OPSR_ANOCA, T; P22332, OPSR_ASTFA, T;
DE P22323, OPSR_CARAU, T; P22329, OPSR_CHICK, T; P04000, OPSR_HUMAN, T;
DE P34889, OPSL_CALAIA, T; P35359, OPSL_HARE, T; P23820, REIS_TODPA, T;
DE P47805, RGR_BOVIN, T; P47804, RGR_HUMAN, T;
DE P17445, OPS3_DROV2, T;
DO P0000211;
```

### Documentation file of an entry from the PROSITE database

```
!P040238
(P040238; OPSIN)
!OPSIN
*****
*Visual pigments (opsins) retinal binding site*
*****
Visual pigments [2] are the light-absorbing molecules that mediate vision. They
consist of an opsin protein, covalently linked to the chromophore cis-retinal.
Vision is effected through the absorption of a photon by cis-retinal which is
isomerized to trans-retinal. This isomerization leads to a change of conformation of
the protein. Opsins are integral membrane proteins with seven transmembrane
regions that belong to family 1 of G protein-coupled receptors (see <P0000210>).
In vertebrates, four different pigments are generally found. Rod cells, which mediate
vision in dim light, contain the pigment rhodopsin. Cone cells, which function in
bright light, are responsible for color vision and contain three or more color pig-
ments (for example, in mammals: red, blue and green).
In invertebrates, the eye is composed of 800 facets or ommatidia. Each ommatidium
contains eight photoreceptor cells (R1-R8) (in R1 and R6 cells are cone cells, R7
and R8 rod cells). Each of the three types of cells (R1, R6, R7 and R8) expresses a
specific opsin.
Previous evolutionary related to opsins include squid retinochrome, also known as
retinal photolipomerase, which converts various sources of retinal into 11-cis-retinal
and bacteriorhodopsin, pigment epithelium (RPE) RGR [3], a protein that may also
act in retinal isomerization.
The attachment site for retinal in the above protein is a conserved lysine residue in
the middle of the seventh transmembrane helix. The pattern we developed includes
this residue.
Consensus pattern: [LIVNM]-[PGC]-x(3)-[SAC]-K-[STALIM]-[GSACHYV]-x(2)-[IDNYT]-[AP]-
x(2)-[LV]-
[3] in the retinal binding site
Sequences known to belong to this class detected by the pattern: ALL
Other sequences detected in SWISS-PROT/NCBI:
-Last update: November 1997 / Pattern and text revised
[1] Applebury M.L., Hargrave P.A.
Vision Res. 26:1801-1805(1986).
[2] Pavesi G.L., Meryemova E.M.
Mol. Biol. 33:367-370(1991).
[3] Shen D., Jiang H., Hao W., Cho J., Salazar M., Fong H.K.W.
Biochemistry 33:13147-13122(1994).
!END!
```

## Secondary databases

### ■ PRINTS

- > developed at University College London
- > motivation: identification of protein families by more than one pattern
- > abstraction: fingerprints (aligned motifs)  
fingerprints store original sequence information
- > construction: sequence information in a seed motifs are augmented through iterative database scanning
- > construction of fingerprints done manually
- > primary source (original): OWL
- > primary source (new): SWISS-PROT and SP-TrEMBL

---

---

---

---

---

---

---

---

### Pattern file of a entry from the PRINTS database (I)

```
(a)
OPSIN          OPSIN SIGNATURE
Type of fingerprint: COMPOUND with 3 elements
Links:
PRINTS: PR00237 GPCRSHODAPSN; PR00247 GPCRAMP; PR00248 GPCRMR
PRINTS: PR00249 GPCRSECRETIN; PR00250 GPCRTEE; PR00251 BACTRLOPSTN
PROSITE: PS00238 OPSIN; PS00237 G_PROTEIN_RECEPTOR
BLOCKS: BL00238
SRASE: OPSD_HUMAN
CCDB: OCL_0085
Creation Date 20-DEC-1993; UPDATE 2-JUL-1996
1. APPELBERY, M.L. and HARGRAVE, P.A.
Molecular biology of the visual pigments.
VISION RES. 26 (12) 1881-1895 (1986).

(b)
SUMMARY INFORMATION
73 codes involving 3 elements
1 codes involving 2 elements
COMPOSITE FINGERPRINT INDEX
31 73 73 73
21 0 1 1
-----
1 1 2 3
```

---

---

---

---

---

---

---

---

### Pattern file of a entry from the PRINTS database (II)

```
(c)
INITIAL MOTIF SETS
OPSIN1 Length of motif = 13 Motif number = 1
Opsin motif I - 1
PCODE ST INT
YVTVQHKKLRTPPL OPSD_BOVIN 60 60
YVTVQHKKLRTPPL OPSD_HUMAN 60 60
YVTVQHKKLRTPPL OPSD_SHEEP 60 60
AATMFKKLRHPL OPSD_HUMAN 76 76
AATMFKKLRHPL OPSD_HUMAN 76 76
YIPATFKSLRTPA OPS1_DROME 73 73
VATLYKXKLRQPL OPSD_HUMAN 57 57
YIPGYSKLRTPA OPS2_DROME 80 80
WVFSAAKSLRTPS OPS3_DROME 81 81
WIFSFKSLRTPS OPS4_DROME 77 77
YLFSTKSLQTPA OPSD_OCTDO 58 58
YLFSTKSLQTPA OPSD_LOLFO 57 57

OPSIN2 Length of motif = 13 Motif number = 2
Opsin motif II - 1
PCODE ST INT
GHSRYIFEGMQCS OPSD_BOVIN 174 101
GHSRYIFEGMQCS OPSD_HUMAN 174 101
GHSRYIFQKQWCS OPSD_SHEEP 174 101
GHSRYIFHGLKTS OPSD_HUMAN 190 101
GHSRYIFHGLKTS OPSD_HUMAN 190 101
GHSRYVPEGNITS OPS1_DROME 187 101
GHSRYIFEGMQCS OPSD_HUMAN 171 101
GHSRYIFEGNITPA OPS2_DROME 194 101
TMGRFVPEGYLTES OPS3_DROME 194 100
FMGRFVPEGYLTES OPS4_DROME 190 100
MGAYVPEGLITS OPSD_OCTDO 174 103
MGAYTLBGLVLCN OPSD_LOLFO 173 103
```

---

---

---

---

---

---

---

---

## Secondary databases

- BLOCKS (abstraction: blocks)
- Profiles (abstraction: profiles)
- Pfam (abstraction: Hidden Markov Models)
  
- IDENTIFY
  - > developed at Stanford University
  - > abstraction: motifs encoded by fuzzy approach (alternative residues are tolerated in motifs)
  - > construction: automatically derived using the program eMOTIF
  - > primary sources: PRINTS and BLOCKS

---

---

---

---

---

---

---

---

## Properties of amino acids used in eMOTIF

<i>Residue property</i>	<i>Residue groups</i>
Small	Ala, Gly
Small hydroxyl	Ser, Thr
Basic	Lys, Arg
Aromatic	Phe, Tyr, Trp
Basic	His, Lys, Arg
Small hydrophobic	Val, Leu, Ile
Medium hydrophobic	Val, Leu, Ile, Met
Acidic/amide	Asp, Glu, Asn, Gln
Small/polar	Ala, Gly, Ser, Thr, Pro

---

---

---

---

---

---

---

---

## Overview of primary sources and stored information in secondary databases

<i>Secondary database</i>	<i>Primary source</i>	<i>Stored information</i>
PROSITE	SWISS-PROT	Regular expressions (patterns)
Profiles	SWISS-PROT	Weighted matrices (profiles)
PRINTS	OWL*	Aligned motifs (fingerprints)
Pfam	SWISS-PROT	Hidden Markov Models (HMMs)
BLOCKS	PROSITE/PRINTS	Aligned motifs (blocks)
IDENTIFY	BLOCKS/PRINTS	Fuzzy regular expressions (patterns)

---

---

---

---

---

---

---

---

## Composite secondary databases

- INTERPRO - Integrated resource of Protein Families, Domains and Sites
  - > developed by EBI, SIB, University of Manchester, Sanger Centre, GENE-IT, CNRS/INRA, LION Bioscience AG and University of Bergen (European Research Project)
  - > provides an integrated view of the commonly used secondary databases: PROSITE, PRINTS, SMART, Pfam and ProDom
  - > accessible by ftp, www and via member databases

---

---

---

---

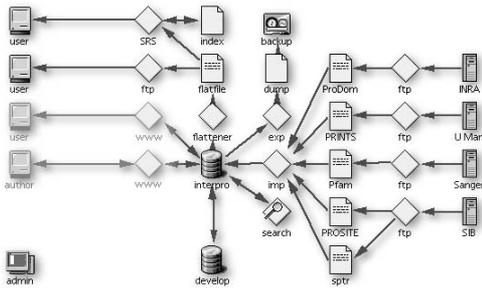
---

---

---

---

InterPro dataflow scheme



---

---

---

---

---

---

---

---

## Protein structure databases

- PDB
- PDBsum

## Protein structure classification databases

- SCOP
- CATCH

---

---

---

---

---

---

---

---

## Genome information resources

- primary DNA sequence databases
- specialised DNA sequence databases

---

---

---

---

---

---

---

---

## Primary DNA sequence databases

- EMBL
- DDBJ
- GenBank
- dbEST
- GSDB

Store DNA sequences and annotations.

---

---

---

---

---

---

---

---

## Primary DNA sequence databases

- EMBL - European Molecular Biology Laboratory
  - > European Bioinformatics Institute (EBI)
  - > collaboration with DDBJ and GenBank - exchange of new entries on daily basis
  - > source of sequences: direct author submissions, genome projects, scientific literature, patents
  - > rate of growth is exponential with doubling time ~9-12 months
  - > most entries from model organisms
  - > retrieval through SRS

---

---

---

---

---

---

---

---



## Primary DNA sequence databases

- dbEST
  - > National Center for Biotechnology Information (NCBI)
  - > maintains only Expressed Sequence Tag (EST) data
  
- GSDB - Genome Sequence DataBase
  - > National Center for Genome Resources
  - > complete collections of DNA sequence for genome-sequencing laboratories
  - > on-line submission of large-scale data
  - > quality checks
  - > format consistent with GenBank + GSDBID

---

---

---

---

---

---

---

---

## Specialised DNA sequence databases

- SGD
  
- UniGene
  
- TDB
  
- ACeDB

Store species-specific and technique-specific DNA sequences.

---

---

---

---

---

---

---

---

## Specialised DNA sequence databases

- SGD - *Saccharomyces* Genome Database
  - > molecular biology and genetics of *S. cerevisiae*
  - > complete genome, genes, proteins, phenotypes
  - > first eukaryotic genome sequenced (1998)
  - > sequence analysis, register of genes, 3D structural data, primer sequences for cloning
  
- UniGene
  - > collection of genes encoding proteins (transcript map)
  - > non-redundant; derived from GenBank
  - > data organised in clusters (1 cluster = 1 unique gene)
  - > gene-mapping projects and gene expression analysis

---

---

---

---

---

---

---

---

## Specialised DNA sequence databases

### ■ TDB - TIGR Database

- > suite of databases: DNA and protein sequences, gene expression, protein families, taxonomic data
- > links: TIGR microbial genome sequencing projects, parasite databases, gene index projects, *A. thaliana* database, human genomic dataset

### ■ ACeDB - A *Cernorhabditis elegans* DataBase

- > *C. elegans* genome project
- > restriction maps, gene structural information, cosmid maps, sequence data, bibliographic information
- > software to organise data ACeDB: CGI script and perl

---

---

---

---

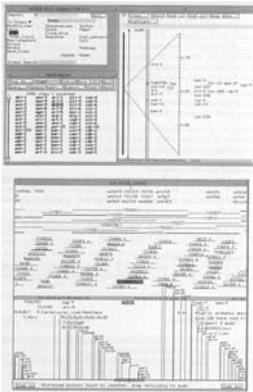
---

---

---

---

## ACeDB software for organisation of genomic data



---

---

---

---

---

---

---

---

## DNA sequence analysis

- why to analyse DNA?
- gene structure
- gene sequence analysis
- expression profile, cDNA, EST
- EST sequences analysis

---

---

---

---

---

---

---

---

## Why to analyse DNA?

- The most sensitive comparisons between sequences are on protein level because of redundancy of the genetic code.
- The loss of degeneracy is accompanied by a loss of information directly linked to the evolution - proteins are only functional abstractions of genetic events at DNA level.
- Silent mutations, important for phylogenetic analysis, can not be detected at protein level.
- Exon/intron analysis, open reading frame [ORF] analysis can not be performed at protein level.

---

---

---

---

---

---

---

---

The genetic code

	T	C	A	G					
T	TTC	Phe	TCT	Ser	TAT	Try	TGT	Cys	T
	TTC		TCC		TAC		TGC		C
	TTA	Leu	TCA		TAA	Stop	TGA	Stop	A
	TTG		TCG		TAG		TGG	Trp	G
C	CTT	Leu	CCT	Pro	CAT	His	CGT	Arg	T
	CTC		CCC		CAC		CGC		C
	CTA		CCA		CAA	Gln	CGA		A
	CTG		CCG		CAG		CGG		G
A	ATT	Ile	ACT	Thr	AAT	Asn	AGT	Ser	T
	ATC		ACC		AAC		AGC		C
	ATA		ACA		AAA	Lys	AGA	Arg	A
	ATG	Met	ACG		AAG		AGG		G
G	GTT	Val	GCT	Ala	GAT	Asp	GGT	Gly	T
	GTC		GCC		GAC		GGC		C
	GTA		GCA		GAA	Glu	GGA		A
	GTG		GCG		GAG		GGG		G

---

---

---

---

---

---

---

---

## Gene structure

- Eukaryotic genes are more complex than prokaryotic due to presence of introns.
- DNA databases typically contain genomic data: untranslated sequences, introns+exons, mRNA, cDNA.
- Gene products (proteins) can be of different length, because not all exons can be present in final mRNA.
- The proteins of different length originating from single sequence are called splice variants.

---

---

---

---

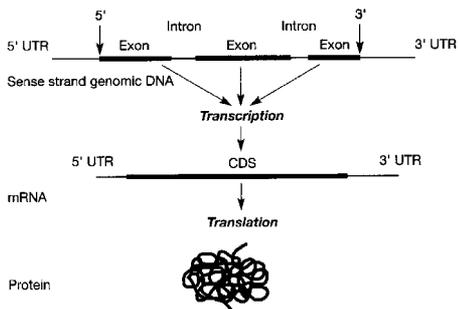
---

---

---

---

## Central dogma of molecular biology



---

---

---

---

---

---

---

---

## Gene structure

- Untranslated regions (UTRs)
  - > portions of the sequence flanking the coding sequence (CDS) not translated into protein
  - > UTRs (especially 3' end) is highly gene/species specific
- Exons
  - > protein-coding DNA sequences of a gene
- Introns
  - > DNA sequences interrupting protein-coding DNA sequence of a gene
  - > transcribed into RNA but are edited out during post-transcriptional modifications

---

---

---

---

---

---

---

---

## Gene sequence analysis

- Conceptual translation - theoretical translation of the DNA sequence to the protein sequence using DNA code without biochemical support.
- Six-frame translation results in six potential protein sequences (ORF analysis).
- ORF analysis
  - > codon for methionine - initial codon in the CDS
  - > sufficient CDS length - long CDS are rare
  - > pattern of codon usage - species specific
  - > bias towards G/C in the third base of a codon - species specific

---

---

---

---

---

---

---

---

## Expression profile, cDNA, EST

- Hierarchy of genomic information
  - > human genome consists of ~3 billion bp
  - > ~3% of the DNA is coding sequence → mRNA → protein
  - > rest of the genome needed for compact structure of chromosomes, replication, control of transcription, etc.
  - > 1. chromosomal genome (genome) - genetic information common to every cell in the organism
  - > 2. expressed genome (transcriptome) - part of genome expressed in a cell at specific stage in its development
  - > 3. proteome - protein molecules that interact to give the cell its individual character

---

---

---

---

---

---

---

---

## Expression profile, cDNA, EST

- Expression profile
  - > characteristic range of genes expressed at particular stage of development and functioning
  - > goal of genome projects is to sequence entire (chromosomal) genome
  - > having complete sequences and knowing what they mean - two distinct stages of understanding genome
  - > alternative approach is analysis of parts of genome expressed in a cell at specific stage in its development
  - > comparison of expression profiles: identification of abnormal expressions, expression levels
  - > interesting for industry - gene discovery, drug design

---

---

---

---

---

---

---

---

## Expression profile, cDNA, EST

### ■ Complementary DNA (cDNA)

- > DNA that is synthesised from a messenger RNA template using the enzyme reverse transcriptase
- > cDNA captures expression profile
- > preparation: cultivation/isolation of cells, mRNA extraction, reverse transcription of mRNA to cDNA, transformation of cDNA into library, sequencing of randomly chosen clones (100.000 out of 2 mil.)
- > ideally 100.000 sequences 200-400 bp length - expressed sequence tags (ESTs)
- > in reality many failures, number of sequences lower
- > number of clones constructed and sequenced must be large enough to represent expression profile

---

---

---

---

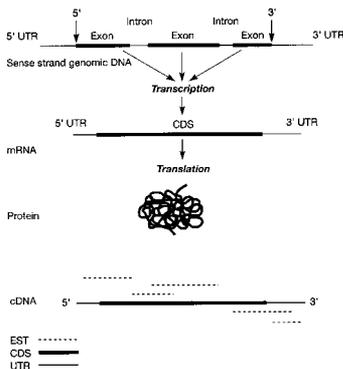
---

---

---

---

## Origin of complementary DNA and expression sequence tags



---

---

---

---

---

---

---

---

## Expression profile, cDNA, EST

### ■ Libraries of ESTs

- > Merck/IMAGE - 300 000 ESTs from a variety of normalised libraries - higher chance to capture different genes; expression levels not known; sequences deposited to dbEST
- > Incyte - quantitative information on expression levels - standardised libraries; expression profiles in healthy and diseased tissues; sequences form the commercial database LifeSeq
- > TIGR - TIGR Human Gene Index - integrates results from human gene projects [dbEST+GenBank] - purpose is to identify all possible human genes by sequence assembly - creates Tentative Human Consensus (THC) sequences and contigs

---

---

---

---

---

---

---

---



## List of base-ambiguity symbols defined by IUB-IUPAC

<i>IUB symbol</i>	<i>Represented bases</i>
A	A
C	C
G	G
T/U	T
M	A or C
R	A or G
W	A or T
S	C or G
Y	C or T
K	G or T
V	A or C or G
H	A or C or T
D	A or G or T
B	C or G or T
X/N	G or A or T or C

---

---

---

---

---

---

---

---

## EST sequences analysis

### ■ Splice variants

- > splice variants are represented by deletions arising from non-inclusion of exons
- > in EST maybe missing bases due to sequencing errors
- > partially good match = splice form or sequence error?

### ■ Non-coding regions

- > question: does this EST represent a new gene?
- > search of DNA database for similar non-coding regions
- > no hit found = the EST represents a new gene (CDS) or the EST represents non-coding sequence not present in the database

---

---

---

---

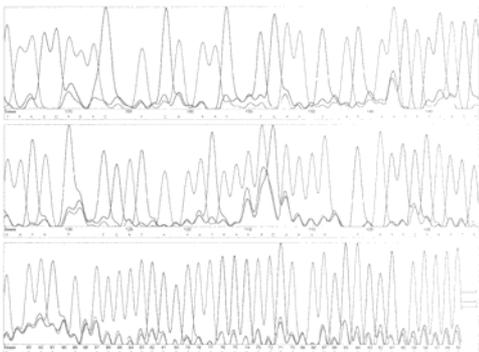
---

---

---

---

## Sequencing chromatogram



---

---

---

---

---

---

---

---

## EST sequences analysis

Three categories of EST analysis tools:

- Sequence similarity search tools
- Sequence assembly tools
- Sequence clustering tools

---

---

---

---

---

---

---

---

## EST sequences analysis

- Sequence similarity search tools
  - > current database search programs are designed to cope with EST: TBLASTN (translate DNA databases), BLASTX (translate input sequence), TBLASTX (translate both)
- Sequence assembly tools
  - > search of the databases reveals several ESTs matching the query sequence
  - > alignment of hits and construction of consensus
  - > search with consensus, alignment, ....
  - > iterative sequence alignment = sequence assembly

---

---

---

---

---

---

---

---

## EST sequences analysis

- Sequence clustering tools
  - > clustering of EST sequences reduces redundancy and saves the search time
  - > enables estimation of genes in the EST database
  - > approach 1: clustering based on sequences from comprehensive DNA database
  - > approach 2: clustering of all ESTs, construction of consensus sequences representing each cluster, DNA database search using consensus sequences only
  - > result = ESTs that do not match any of the database sequences

---

---

---

---

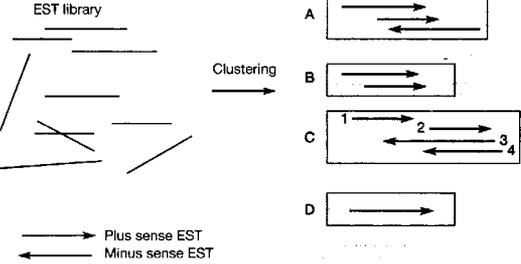
---

---

---

---

Clustering of EST library



---

---

---

---

---

---

---

---

## Pairwise sequence alignment

- database searching
- alphabets and complexity
- algorithms and programs
- sequences and sub-sequences
- identity and similarity
- dotplot
- local and global similarity
- pairwise database searching

---

---

---

---

---

---

---

---

## Database searching

- Database search can take a form of text queries or sequence similarity searches.
- Text queries are problematic due to missing annotations in many sequences.
- query sequence = probe  
searched sequence = subject
- The purpose of searches is to identify evolutionary relationships (homology) from sequence similarity. Important for search of analogous family members in different species.

---

---

---

---

---

---

---

---

## Alphabets and complexity

- A sequence consists of letters from an alphabet.
- The complexity of the alphabet is defined by the number of letters it contains:
  - > DNA = 4
  - > EST = 5
  - > proteins = 20
- Special letters can be used for ambiguous bases (N) or residues (X). Sequence searching programs must be able to deal with them.

---

---

---

---

---

---

---

---

## Algorithms and programs

- Algorithm is a set of steps that define a certain computational process.
- Program is a the implementation of the algorithm.
- Same algorithm may be implemented in many programs.

---

---

---

---

---

---

---

---

## Sequences and sub-sequences

- Alignment of two short sequences:

```
Unaligned                               score = 6
Sequence 1 (query)      AGGVLLIQVG
                        | | | | |
Sequence 2 (subject)    AGGVLLIQVG

Aligned                               score = 9
Sequence 1 (query)      AGGVLLIQVG
                        | | | | | | |
Sequence 2 (subject)    AGGVLI-QVG
```

- Score increases by the insertion of a gap. The gap increases the number of aligned identical residues.

---

---

---

---

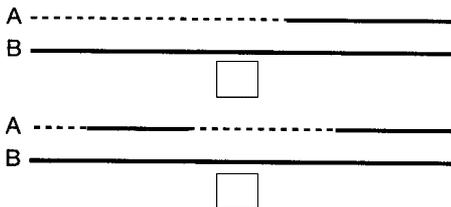
---

---

---

---

### Alignment of a sub-sequence with full sequence



---

---

---

---

---

---

---

---

## Identity and similarity

- Introduction of gaps solely to maximise identities is not biologically meaningful.
- Scoring penalties are introduced to minimise opening and extension of gaps.
- Unitary matrix (counting identities) is replaced by similarity matrix (counting similarities) = high-scoring matches are replaced by biologically meaningful low-scoring matches.
- Diagnostic power of similarity matrices is higher.

---

---

---

---

---

---

---

---

### Unitary scoring matrices: (a) DNA and (b) protein

(a)

	A	C	G	T
A	1	0	0	0
C	0	1	0	0
G	0	0	1	0
T	0	0	0	1

(b)

	C	S	T	P	A	G	N	D	E	H	R	K	M	I	L	V	F	W	B	Z	X
C	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
S	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
T	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
P	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
A	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
G	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0
D	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0	0
E	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0
Q	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0
H	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0
R	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0	0
K	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0	0
M	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0	0
I	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0	0
L	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0	0
V	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0	0
F	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0	0
Y	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0	0
W	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1	0
B	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1
Z	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
X	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	1

---

---

---

---

---

---

---

---

## Identity and similarity

- Dayhoff Mutation Data Matrix
  - > score is based on the concept of Point Accepted Mutation (PAM)
  - > evolutionary distance 1 PAM = probability of a residue mutating during a distance in which 1 point mutation is accepted per 100 residues
  - > 250 PAM matrix - similarity score equivalent to 20% matches remaining between two sequences = suitable for identification of similarities in twilight zone
  - > limitation: derived from alignment of sequences >85% identical

---

---

---

---

---

---

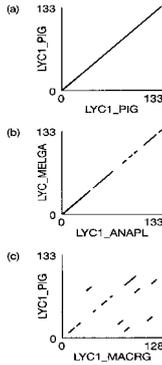
---

---





Dotplot of (a) identical, (b) similar and (c) related sequences




---

---

---

---

---

---

---

---

### Local and global similarity

- Alignments are mathematical models whose behaviour can be modified through the use of adjustable parameters. The models constructed by dynamic programming algorithms - finding solution of a problem by solving smaller, but similar sub-problems.
- Global alignment - considers similarity across the entire sequence.
- Local alignment - considers similarity in parts of sequences only.

---

---

---

---

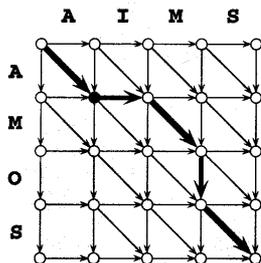
---

---

---

---

Path matrix with optimal path by dynamic programming



Alignment **AIM-S**  
**A-MOS**

---

---

---

---

---

---

---

---

## Local and global similarity

### ■ Global alignment

- > Needleman and Wunsch algorithm
- > suitable for sequences similar across most of their length (usually closely related)
  
- > 1. construction of 2D similarity matrix ("dotplot")
- > 2. successive summation of the cells in the matrix starting from N-terminal end → progressing through the sequence
- > 3. construction of maximum-match path through the entire sequence

---

---

---

---

---

---

---

---

## Local and global similarity

### ■ Local alignment

- > Smith-Waterman algorithm
- > suitable for distantly related sequences displaying local regions of similarity (functionally-relevant or structurally-relevant)
- > each point of the matrix defines the end point of a potential alignment = edge cells of the matrix are initialised to 0
- > possibility for ending the alignment are calculated for every cell
- > algorithm is much faster compared to global similarity algorithms

---

---

---

---

---

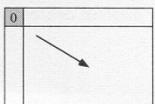
---

---

---

## Concepts of global and local optimality in the pairwise sequence alignment

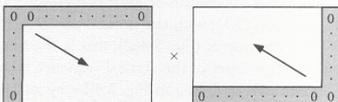
(a) Global vs. Global



(b) Local vs. Global



(c) Local vs. Local



---

---

---

---

---

---

---

---

## Pairwise database searching

- Extension of the pairwise sequence alignments.
- Large database searches can not be performed using the original Needleman and Wunsch or Smith-Waterman algorithms due to time limitations.
- Very fast local-similarity search methods employing heuristics = FastA and BLAST. These methods concentrates on finding short identical matches.

---

---

---

---

---

---

---

---

---

---

## Pairwise database searching

- FastA
  - > algorithm by Lipman and Pearson (1985)
  - > identifies short words (k-tuples) common to both sequences
  - > k-tuples for proteins: 1-2 residues
  - > k-tuples for DNA: up to 6 bases
  - > k-tuples lying close to each other on the same diagonal joined by heuristics → gapped alignments computed by dynamic programming

---

---

---

---

---

---

---

---

---

---

### Output from FastA search

```
FASTA searches a protein or DNA sequence data bank
version 3.3009 May 18, 2001
Please cite:
W.B. Pearson & D.J. Lipman PNAS (1988) 85:2444-2448

E1:1: 296 aa
EMBOSS_001
vs SWISS-PROT Protein Sequence Database library
searching /usr/local/ncbi/data/FASTA/EMBOSS/FASTA/prot library

37135523 residues in 101247 sequences
statistics extrapolated from 60000 to 101082 sequences
Expectation = 1e-11; rho(ln(x)) = 3.2158e+/-0.000184; mme 4.0373e+/- 0.010
mean_val=74.4386e+/-4.720, 0's: 132 2's:11 3's: 0 4's: 0 5's: 0 in 0/65
Lambd= 0.1487

FASTA (3.39 May 2001) function [optimized, BL50 matrix (151-51)] kdup: 2
join: 36, opt: 24, gap-pen: -12/-2, width: 16
Scan time: 1.930
The best scores are:
SW:LINE_PEEPA P51698 1,3,4,6-TETRACHLORO-1,4-CYCL ( 296) 2041 447 2.4e-125
SM:YFP9_MCTU Q50642 HYDROTHERICAL 33.7 KDa PROTEI ( 300) 1494 330 5.1e-90
SM:LUCC_PEESE P21620 BENTILIA-LUCIFERIN 2-NONOXYG ( 311) 744 169 1.4e-41
SM:MPD_PEEEP P19076 2-HYDROXYMICONIC SEMIALDEHYD ( 283) 169 46 0.00037
SM:PRC_PEEPL Q31136 NON-HEME CHLOROPEROXIDASE (E ( 273) 168 45 0.00019
SM:PRC_PEEPL P49323 NON-HEME CHLOROPEROXIDASE (E ( 273) 140 39 0.032
SM:PFP_BACCO P46041 PROLINE IMINOPEPTIDASE (EC 3. ( 288) 140 39 0.013
SM:PRC_PEEVY Q50921 PUTATIVE NON-HEME CHLOROPERO ( 276) 123 36 0.11
SM:PFP_HEIGO P42786 PROLINE IMINOPEPTIDASE (EC 3. ( 310) 122 35 0.2

>>SM:LINE_PEEPA P51698 1,3,4,6-TETRACHLORO-1,4-CYCLOHEXA (296 aa)
Inlin: 2041 init1: 2041 opt: 2041 E-score: 2372.6 bits: 467.0 E(): 2.4e-125
Smith-Waterman score: 2041; 100.00% identity (100.00% unaligned) in 296 aa overlap
(i1=296;i2=296)

      10      20      30      40      50
EMBOSS MSILGAKFFGKFKFIRKGRMAYIDETCTDPIFQGNPTSYLWNNIMFHCAGLRLIA
      10      20      30      40      50      60
SM:LIN MSILGAKFFGKFKFIRKGRMAYIDETCTDPIFQGNPTSYLWNNIMFHCAGLRLIA
      10      20      30      40      50      60
```

---

---

---

---

---

---

---

---

---

---

# Pairwise database searching

## ■ BLAST

- > Basic Local Alignment Search Tool
- > algorithm by Altschul *et al.* (1990)
- > identifies short ungapped sub-sequences (segment pairs) of the same length
- > sub-sequences are extended using dynamic programming to obtain local alignments - high scoring pairs (HSPs)
- > improved algorithm by Altschul *et al.* (1997) - produces gapped alignments
- > algorithm very fast - most commonly used for databases searching

---

---

---

---

---

---

---

---

---

---

## Output from BLAST search

```
BLASTF 2.0.14 (Jun-29-2000)

Reference: Altschul, Stephen F., Thomas L. Madden, Alejandro A. Schaffer,
Jinghui Zhang, Zheng Zhang, Webb Miller, and David J. Lipman (1990),
"Simple BLAST and PSI-BLAST: a new generation of protein database search
programs", Nucleic Acids Res. 25:3389-3402.

Query= /net/nts0/vol1/production/w3nobody/tmp/918495.5350-
80758.blastall.a [Unknown form], 297 bases, 818F03BD checksum.
(296 letters)

Database: swissprot
101,247 sequences; 37,135,523 total letters

Searching.....done

Sequences producing significant alignments:

Score E
(bits) Value
SW:LNIB_FSEPA P51698 1,3,4,6-TETRACHLORO-1,4-CYCLOHEXADIENE ... 616 e-176
SW:YF79_MYCTU Q06642 HYPOTHETICAL 33.7 KDA PROTEIN RV2579. 450 e-126
SW:LUCT_BENRE P27652 BENZILLA-LUCIFERIN 2-MONOXYGENASE (EC 1... 218 2e-56
SW:IMFP_FSEFP P19076 2-HYDROXYMETHYL BEMALISERVOSE HYDROLASE... 50 2e-06
SW:FMXC_FSEFL Q31158 NON-HEME CHLOROPEPOXIDASE (EC 1.11.1.10... 45 2e-04
SW:HWAL_STRAU P23719 NON-HAEM BROMOPEROXIDASE BPO-HA (EC 1.1... 39 0.011
SW:FIF_BACCO P46541 PROLINE IMINOPEPTIDASE (EC 3.4.11.5) (PI... 39 0.014
SW:FIF_MEIMS Q93286 PROLINE IMINOPEPTIDASE (EC 3.4.11.5) (PI... 36 0.016

>SW:LNIB_FSEPA P51698 1,3,4,6-TETRACHLORO-1,4-CYCLOHEXADIENE HYDROLASE (EC
3.4.11.-) [1,4-TCOH CHLOROHYDROLASE].
Length = 296

Score = 616 bits (1572), Expect = e-176
Identities = 296/296 (100%), Positives = 296/296 (100%)

Query: 1 MELGANFGEKFFIEIGRRMAYIDGTDGPIFGQGNFTSYLWNIWMPHCALGLRIA 60
MELGANFGEKFFIEIGRRMAYIDGTDGPIFGQGNFTSYLWNIWMPHCALGLRIA
Sbjct: 1 MELGANFGEKFFIEIGRRMAYIDGTDGPIFGQGNFTSYLWNIWMPHCALGLRIA 60
```

---

---

---

---

---

---

---

---

---

---

## Multiple sequence alignment

- multiple sequence alignment
- consensus sequence
- manual methods
- simultaneous and progressive methods
- databases of multiple sequence alignments
- hybrid approach for database searching

---

---

---

---

---

---

---

---

## Multiple sequence alignment

- Multiple sequence alignment is a 2D table in which the rows represent individual sequences and the columns the residue positions.
- Multiple sequence alignments are essential for analysis of sets of gene families.
- Sequence-based multiple sequence alignments - constructed according to similar strings of amino acid residues.
- Structure-based multiple sequence alignments - constructed according to structural evidence.

---

---

---

---

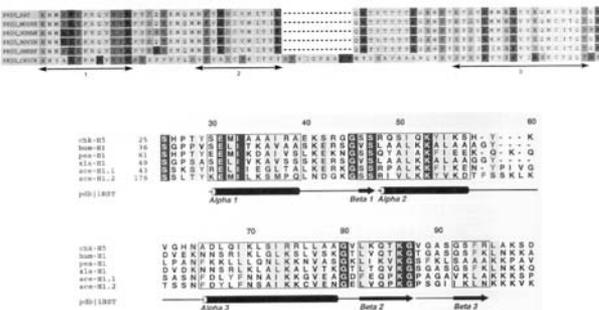
---

---

---

---

## Colour-coded multiple sequence alignments



---

---

---

---

---

---

---

---

## Multiple sequence alignment

- Construction of a multiple sequence alignment:
  - > positioning of residues within any sequence is preserved (absolute positions)
  - > similar residues in all sequences are brought into vertical register (relative positions)
- All residues in any single column of an alignment will have the same relative position but different absolute position (unless the sequences are identical).

---

---

---

---

---

---

---

---

## Consensus sequence

- The alignment table can be summarised by:
  - > a single line: pseudo-sequence
  - > unweighted matrix: fingerprint
  - > ungapped block of residues (weighted): block
  - > weighted matrix: profile

---

---

---

---

---

---

---

---

### Multiple alignment and the consensus sequence

	1	2	3	4	5	6	7	8	9	10
I	Y	D	G	G	A	V	-	E	A	L
II	Y	D	G	G	-	-	-	E	A	L
III	F	E	G	G	I	L	V	E	A	L
IV	F	D	-	G	I	L	V	Q	A	V
V	Y	E	G	G	A	V	V	Q	A	L
	y	d	G	G	A/I	V/L	V	e	A	l

---

---

---

---

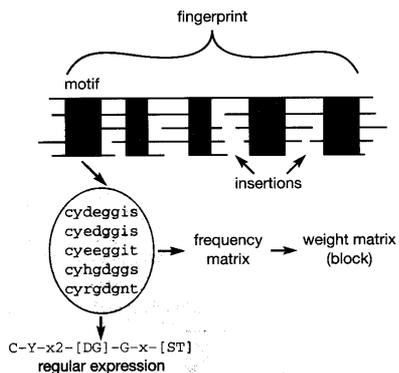
---

---

---

---

## Multiple alignment and the profile, block and fingerprint




---

---

---

---

---

---

---

---

## Manual methods

- Manual methods are subjective however they enable to incorporate experimental evidences (e.g., mutagenesis data, structural knowledge) into the multiple alignment.
- Manual modification of the multiple alignments from automatic methods is the best approach.
- Intuitive colouring schemes assist the eye in spotting similarities.
- Quantitative evaluation of relatedness through calculation of residue identities/similarities.

---

---

---

---

---

---

---

---

## Amino acid property groupings and colouring

<i>Residue</i>	<i>Property</i>	<i>Colour</i>
Asp, Glu	Acidic	red
His, Arg, Lys	Basic	blue
Ser, Thr, Asn, Gln	Polar neutral	green
Ala, Val, Leu, Ile, Met	Hydrophobic aliphatic	white
Phe, Try, Trp	Hydrophobic aromatic	purple
Pro, Gly	Special structural properties	brown
Cys	Disulphide bond former	yellow

---

---

---

---

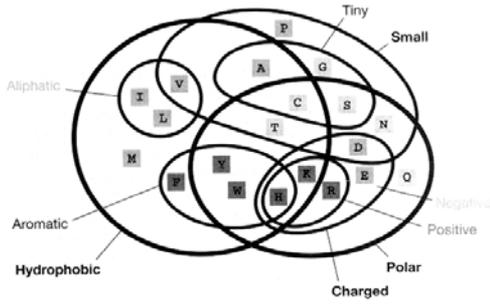
---

---

---

---

### Venn diagram grouping properties of the amino acids



---

---

---

---

---

---

---

---

### Simultaneous methods

- Simultaneous methods align all sequences in a given set at once, rather than aligning pairs of sequences or building sequence clusters.
- Extension of 2D dynamic programming matrix to more dimensions.
- Number of dimensions = number of sequences.
- Suitable only for small sets of short sequences.

---

---

---

---

---

---

---

---

### Progressive methods

- Multi-dimensional programming matrix is not applicable to realistic problems - larger sets of longer sequences.
- CLUSTAL
  - > 1. construction of evolutionary tree
  - > 2. pairwise alignment of two the most closely related sequences, addition of less related sequences
  - > 3. final alignment, final evolutionary tree
- CLUSTALW
  - > positioning of gaps in closely related sequences according to their variability

---

---

---

---

---

---

---

---

## Databases of multiple alignments

- Multiple alignments bring together sequences from different species. This important evolutionary information can enhance sensitivity of database searches.
- Various abstractions (regular expressions, profiles, blocks, fingerprints or HMMs) can be searched against sequence databases. More information used in a query - higher sensitivity.
- Results of the searches using the multiple alignments are more difficult to interpret.

---

---

---

---

---

---

---

---

## Databases of multiple alignments

- Multiple alignments databases available via Web are produced automatically (e.g., PFAM) or manually (e.g., PRINTS).
- Iterative automatic methods may include false-positive sequences in the alignment which will corrupt it by insertion of many unrealistic gaps.

---

---

---

---

---

---

---

---

### Example entry from PFAM database

 **Pfam**  
Protein families database of alignments and HMMs  
Home | Search | About | Help | Contact | Feedback | Privacy Policy

**zf-C2H2**

Accession number: PF00996  
Zinc finger, C2H2 type

The C2H2 zinc finger is the classical zinc finger domain. The two conserved cysteines and histidines co-ordinate a zinc ion. The following pattern describes the zinc finger:  $[X-C-X(1-6)-G-C-X(1-4)-X(1-2)-H-X(1-4)-G(1-4)-H(1-4)]$  Where X can be any amino acid, and numbers in brackets indicate the number of residues. The positions marked # are those that are important for the stable fold of the zinc finger. The first position can be either His or Cys. The C2H2 zinc finger is composed of two short beta strands followed by an alpha helix. The amino terminal part of the helix binds the major groove in DNA binding zinc fingers.

**INTERPRO description (entry IPR000121)**

Zinc finger domain [MEDLINE: 86151016, P183005929] are nucleic acid-binding protein structures first identified in the *Xenopus* transcription factor TFIIIA. These domains have since been found in numerous nucleic acid-binding proteins. A zinc finger domain is composed of 25 to 30 amino-acid residues including 2 conserved Cys and 2 conserved His residues in a C-2-C-12-H-3-H type motif. The 12 residues separating the second Cys and the first His are mainly polar and basic, implicating this region in particular in nucleic acid binding. The zinc finger motif is an unusually small, self-folding domain in which Zn is a crucial component of its tertiary structure. All bind 1 atom of Zn in a tetrahedral array to yield a finger-like projection, which interacts with nucleotides in the major groove of the nucleic acid. The Zn binds to the conserved Cys and His residues. Zingers have been found to bind to about 5 base pairs of nucleic acid containing short runs of guanine residues. They have the ability to bind to both DNA and RNA, a versatility not demonstrated by the helix-turn-helix motif. The zinc finger may thus represent the original nucleic acid binding protein. It has also been suggested that a C2H2-control domain could be used in a protein interaction, e.g. in protein kinase C. Many classes of zinc fingers are characterized according to the number and positions of the histidine and cysteine residues involved in the zinc atom coordination. In the first class to be characterized, called C2H2, the first pair of zinc-coordinating residues are cysteines, while the second pair are histidines.

**Figure 1: 1ath**  
Complex zinc finger dimer  
Cys9 (C83at variant) zinc finger-dna complex  
(cysic site)

For additional annotation, see the PROSITE document P0000018 [Sequence: EEE-SFJ EEE-SDA]

---

---

---

---

---

---

---

---



## Secondary database searching

- why to search secondary databases?
- secondary databases
- regular expressions
- fingerprints
- blocks
- profiles
- Hidden Markov Models

---

---

---

---

---

---

---

---

## Why to search secondary databases?

- Interpretation of the results from primary database searches is sometimes difficult:
  - > X.000.000 sequences from XX.000 organisms
  - > complex and redundant search outputs
  - > irrelevant matches of low-complexity sequences, repetitive sequences, modular sequences
  - > local regions of similarity in multi-domain proteins
  - > truncated description lines
- Secondary database searches enable to identify both homology and more exacting orthology.

---

---

---

---

---

---

---

---

## Secondary databases

- Contains information derived from primary sequence data, typically in the form of abstractions: regular expressions, fingerprints, blocks, profiles or Hidden Markov Models.
- These abstractions represent distillations of the most conserved features of multiple alignments.
- The abstractions are useful for discrimination of family membership for newly determined sequences.

---

---

---

---

---

---

---

---

## Secondary databases

- PROSITE - regular expressions
- PRINTS - fingerprints
- BLOCKS - blocks
- Profiles - profiles
- Pfam - Hidden Markov Models
- IDENTIFY - fuzzy regular expressions

---

---

---

---

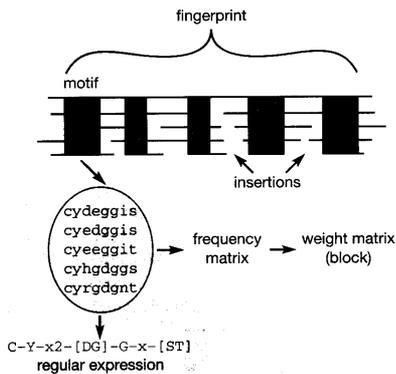
---

---

---

---

## Terms used in sequence analysis methods



---

---

---

---

---

---

---

---

## Regular expressions

- Regular expression reduces the sequence data to the most conserved residue information.

Multiple alignment	Regular expression
ADLGAVFALCDRYFQ	[AS]-D-[IVL]-G-X5-C-[DE]-R-[FY]2-Q
SDVGFPSFCERFYQ	
ADLGRTQNRCDRYFQ	
ADIGQP HSLCERYFQ	

- Limitations:
  - > stringent pattern - retrieves only identical matches and can miss remote relatives
  - > fuzzier pattern - better chance to detect remote relatives, but results in more noisy output
  - > single motif may not be sufficient to infer the function

---

---

---

---

---

---

---

---

## Regular expressions

- Regular expressions works most effectively when a particular protein family can be characterised by a highly conserved motif (10-20 residues).
- Limitation: short patterns (3-4 residues) are not sufficiently discriminative.

Asp-Ala-Val-Ile-Asp (DAVID)      71 exact matches in OWL29.6  
Asp-Ala-Val-Glu (DAVE)        1088 exact matches in OWL29.6

---

---

---

---

---

---

---

---

## Regular expressions

- Rules - short patterns that can be used to provide a guide to possible existence of functional sites:

Functional site	Regular expression
N-glycosylation site	N-(P)-[ST]-{P}
Protein kinase C phosphorylation site	[ST]-X-[RK]
Casein kinase II phosphorylation site	[ST]-X(2)-[DE]
Asp adn Asn hydroxylation site	C-X-[DN]-X(4)-[FY]-X-C

---

---

---

---

---

---

---

---

## Regular expressions

- Fuzzy regular expressions - regular expressions with introduced fuzziness into patterns using groups of amino acids with similar biochemical properties (FYW - aromatic, HKR - basic, etc.).

Multiple alignment	Fuzzy regular expression
ADLGAVFALCDRYFQ	[ASGPT]-D-[IVLM]-G-X5-C-[DENQ]-R-[FYW]2-Q
SDVGFPSFCERFYQ	
ADLGRQNRCDRYFQ	
ADIGQPHSLCERYFQ	

---

---

---

---

---

---

---

---

## Amino acid property groupings and colouring

<i>Residue</i>	<i>Property</i>	<i>Colour</i>
Asp, Glu	Acidic	red
His, Arg, Lys	Basic	blue
Ser, Thr, Asn, Gln	Polar neutral	green
Ala, Val, Leu, Ile, Met	Hydrophobic aliphatic	white
Phe, Try, Trp	Hydrophobic aromatic	purple
Pro, Gly	Special structural properties	brown
Cys	Disulphide bond former	yellow

---

---

---

---

---

---

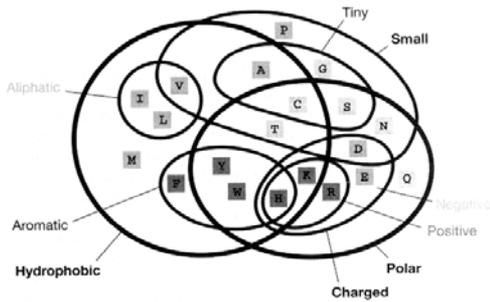
---

---

---

---

## Venn diagram grouping properties of the amino acids




---

---

---

---

---

---

---

---

---

---

## Regular expressions

- Introduction fuzziness into regular expressions increases the number of matches retrieved from the sequence database:

Regular expression	No. of exact matches (OWL29.6)
D-A-V-I-D	71
D-A-V-I-[DENQ]	252
[DENQ]-A-V-I-[DENQ]	925
[DENQ]-A-[VLI]-I-[DENQ]	2739
[DENQ]-[AQ]-[VLI]2-[DENQ]	51506

---

---

---

---

---

---

---

---

---

---



# Blocks

- Conserved motifs are located by a first motif-finding algorithm: search for the spaced residue triplets (e.g., Ala-X-X-Val-X-Trp); a block score is weighted using BLOSUM 62 substitution matrix.
- Validation of blocks by a second motif-finding algorithm: search for the highest-scoring set of blocks in the correct order without overlapping.
- Sequences are clustered to avoid a bias due to identical sequences.

---

---

---

---

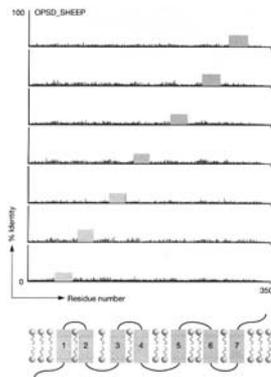
---

---

---

---

## Visualisation of protein fingerprints and blocks




---

---

---

---

---

---

---

---

## Block with clustered sequences and weighted scores

```

OCCR_HUMAN ( 302) SSCVRFZICVNRKPKF 3
OCCR_RAT ( 378) SSCVRFZICVNRKPKF 3
PRLA_HUMAN ( 294) HICLARNLVYVVGQDF 4
PRLA_HUMAN ( 293) HICLARNLVYVVGQDF 4
PRLA_MOUSE ( 304) HICLARNLVYVVGQDF 4
PRLA_RABIT ( 293) HICLARNLVYVVGQDF 4
GARR_CANYA ( 388) ZACVRFVYVCMRKP 5
GARR_HUMAN ( 382) ZACVRFVYVCMRKP 5
GARR_PANNA ( 385) ZACVRFVYVCMRKP 5
GARR_RABIT ( 387) ZACVRFVYVCMRKP 5
GARR_RAT ( 389) ZACVRFVYVCMRKP 5
ETFR_BOVIN ( 363) HICIRFALVYVGRK 9
ETFR_HUMAN ( 363) HICIRFALVYVGRK 9
ETFR_MOUSE ( 379) HICIRFALVYVGRK 9
ETFR_HUMAN ( 378) HICIRFALVYVGRK 9
ETFR_PIG ( 379) HICIRFALVYVGRK 9
ETFR_RAT ( 378) HICIRFALVYVGRK 9
OFSD_GALPO ( 397) SAIDNRKIVYVGRK 12
OFSD_GOTCO ( 398) SAIDNRKIVYVGRK 12
OFSD_TOGAA ( 396) SAIDNRKIVYVGRK 12
FOIR_HUMAN ( 286) HICLQVFLYVLAQGLV 13
FOIR_MOUSE ( 288) HICLQVFLYVLAQGLV 13
FOIR_RAT ( 287) HICLQVFLYVLAQGLV 13
SHL_RAT ( 312) HIRVRFVYVGRK 16
ENGL_HUMAN ( 302) HIRVRFVYVGRK 16
ERIL_HUMAN ( 305) HIRVRFVYVGRK 16
ORVY_HUMAN ( 321) HICIRNIVYVGRK 14
ORVY_PIG ( 323) HICIRNIVYVGRK 14
VLAR_HUMAN ( 345) HICIRNIVYVGRK 18
VLAR_RAT ( 344) HICIRNIVYVGRK 18
PERI_BOVIN ( 337) HICLQVFLYVLAQGLV 13
PERI_HUMAN ( 338) HICLQVFLYVLAQGLV 13
YHRA_CARRE ( 331) SCVRFVYVGRK 100
    
```

---

---

---

---

---

---

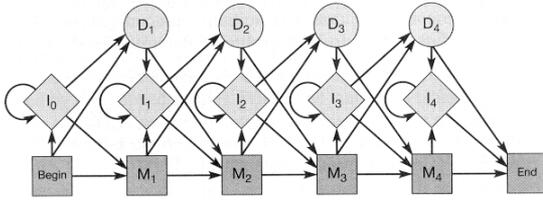
---

---



### Scheme of the linear Hidden Markov Model

(M) match, (I) insertion, (D) deletion



---

---

---

---

---

---

---

---

## Analysis packages

- commercial databases
- commercial software
- comprehensive packages
- packages for DNA analysis
- intranet packages
- Internet packages

---

---

---

---

---

---

---

---

## Comprehensive packages

- GCG
- EGCG/EMBOSS
- Staden
- Lasergene

## Packages for DNA analysis

- Sequencher
- VectorNTI
- MacVector

---

---

---

---

---

---

---

---

## Intranet packages

- SYNERGY
- GeneMill, GeneWorld, GeneThesaurus

## Internet packages

- CINEMA
- EGCG/EMBOSS
- Alfresco

---

---

---

---

---

---

---

---

## Protein structure modelling

- protein structure
- protein structure databases
- prediction of secondary structure
- prediction of protein fold
- prediction of tertiary structure
- modelling of protein-ligand complexes

---

---

---

---

---

---

---

---

## Protein structure

- Proteins are built up by amino acids that are linked by peptide bonds. The 20 different amino acids occur naturally in proteins.
- Protein structure can be experimentally determined by X-ray crystallography, nuclear magnetic resonance (NMR) or by electron crystallography.
- Levels of protein structure:
  - > primary structure
  - > secondary structure
  - > supersecondary structure
  - > tertiary structure
  - > quaternary structure

---

---

---

---

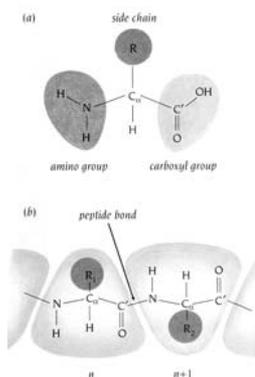
---

---

---

---

## Scheme of an amino acids (a) and polypeptide chain (b)



---

---

---

---

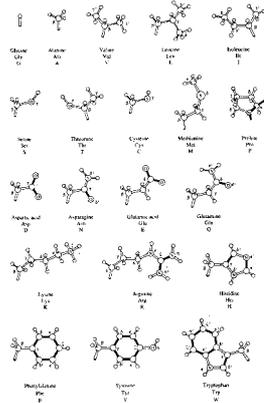
---

---

---

---

## Side chains of 20 different amino acids that occur in proteins




---

---

---

---

---

---

---

---

## Levels of protein structure

- Primary structure:** the linear sequence of amino acids in a protein molecule
- Secondary structure:** regions of local regularity within a protein fold (e.g.,  $\alpha$ -helices,  $\beta$ -turns,  $\beta$ -strands)
- Super-secondary structure:** the arrangement of  $\alpha$ -helices and/or  $\beta$ -strands into discrete folding units (e.g.,  $\beta$ -barrels,  $\beta\alpha\beta$ -units, Greek keys, etc.)
- Tertiary structure:** the overall fold of a protein sequence, formed by the packing of its secondary and/or super-secondary structure elements
- Quaternary structure:** the arrangement of separate protein chains in a protein molecule with more than one subunit
- Quinternary structure:** the arrangement of separate molecules, such as in protein-protein or protein-nucleic acid interactions

---

---

---

---

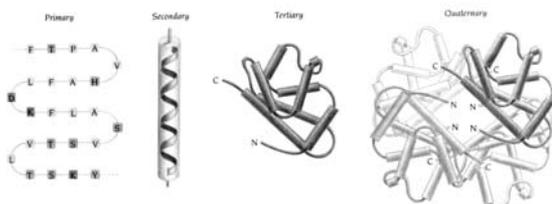
---

---

---

---

## Levels of protein structure




---

---

---

---

---

---

---

---

## Synchrotron radiation facility



European Synchrotron Radiation Facility at Grenoble, France

---

---

---

---

---

---

---

---

## Protein structure databases

- PDB
- PDBsum

## Protein structure classification databases

- SCOP
- CATCH

---

---

---

---

---

---

---

---

## Protein structure databases

- PDB - Protein Data Bank
  - > developed at Brookhaven National Laboratory
  - > currently maintained by Research Collaboratory for Structural Bioinformatics (RCSB)
  - > world repository of three-dimensional protein structures
  - > entries from crystallographic analysis (80%), nuclear magnetic resonance (16%) and modelling (2%)
  - > entries stored as flat files composed of section for information records and section for co-ordinates
  - > entries identified by unique PDB-ID code (e.g., 1EDE)
  - > searchable by keywords
  - > interactive visualization of structures

---

---

---

---

---

---

---

---

# Information on entry from the PDB database

**PDB**  
PROTEIN DATA BANK

**Structure Explorer - 1CV2**

**Summary Information**

**Summary Information**  
 Title: Hydrolytic Haloalkane Dehalogenase Lmb From *Sphingomonas Paucimobis* UT26 At 1.6 Å Resolution  
 Compound Mol ID: 1; Molecule: Haloalkane Dehalogenase; Chain: A; Synonym: Lmb; 1,3,4,6-Tetrachloro-1,4-Cyclohexadiene Hydrolase; EC: 3.8.1.5; Engineered: Yes  
 Authors: J. Marek, J. Vavrdova, J. Damborsky, I. Smatanova, L. A. Svensson, J. Newman, Y. Nagata, M. Takagi  
 Exp. Method: X-ray Diffraction  
 Classification: Hydrolase  
 EC Number: 3.8.1.5  
 Source: *Sphingomonas Paucimobis*  
 Primary Citation: Marek, J., Vavrdova, J., Smatanova, I., Nagata, Y., Svensson, L. A., Newman, J., Takagi, M., Damborsky, J.: Crystal Structure of the Haloalkane Dehalogenase from *Sphingomonas Paucimobis* Ut26 *Biochemistry* 39 pp. 14072 (2000) [Medline]

Deposition Date: 22-Aug-1999 Release Date: 11-Sep-2000

Resolution (Å): 1.58 R-Value: 0.149  
 Space Group: P 21 21 2  
 Unit Cell: a 90.26 b 71.67 c 72.70  
 angles (°): alpha 90.00 beta 90.00 gamma 90.00

Polymer Chains: A Residues: 296  
 Atoms: 2750  
 HET groups: HOH



# Entry from the PDB database (header)

```

HEADER HYDROLASE 22-AUG-99 1CV2
TITLE HYDROLYTIC HALOALKANE DEHALOGENASE LMB FROM SPHINGOMONAS
TITLE 2 PAUCIMOBIS UT26 AT 1.6 Å RESOLUTION
COMPND MOL ID: 1
COMPND 2 MOLECULE: HALOALKANE DEHALOGENASE;
COMPND 3 CHAIN: A;
COMPND 4 SYNONYM: LMB, 1,3,4,6-TETRACHLORO-1,4-CYCLOHEXADIENE
COMPND 5 HYDROLASE;
COMPND 6 EC: 3.8.1.5;
COMPND 7 ENGINEERED: YES;
COMPND 8 BIOLOGICAL UNIT: MONOMER
SOURCE MOL ID: 1
SOURCE 2 ORGANISM SCIENTIFIC: SPHINGOMONAS PAUCIMOBIS;
SOURCE 3 STRAIN: UT26;
SOURCE 4 EXPRESSION SYSTEM: ESCHERICHIA COLI;
SOURCE 5 EXPRESSION SYSTEM STRAIN: HB101;
SOURCE 6 EXPRESSION SYSTEM VECTOR TYPE: PLASMID;
SOURCE 7 EXPRESSION SYSTEM PLASMID: pFUS1
KEYWDS DEHALOGENASE, LINDANE, BIODEGRADATION, ALPHA/BETA-HYDROLASE
EXPTA X-RAY DIFFRACTION
AUTHOR J.MAREK, J.VAVRDOVA, J.DAMBORSKY, I.SMATANOVA, L.A.SVENSSON,
AUTHOR 2 J.NEWMAN, Y.NAGATA, M.TAKAGI
REMARK 1 REFERENCE 1
REMARK 1 AUTH 1, SMATANOVA, Y., NAGATA, L. A., SVENSSON, M., TAKAGI, J., MAREK
REMARK 1 TITL CRYSTALLIZATION AND PRELIMINARY X-RAY DIFFRACTION
REMARK 1 TITL 2 ANALYSIS OF HALOALKANE DEHALOGENASE LMB FROM
REMARK 1 TITL 3 SPHINGOMONAS PAUCIMOBIS UT26
REMARK 1 REF ACTA CRYST. D V. D53 1231 1999
REMARK 1 REFIN DR ISSN 0907-4449
REMARK 2
REMARK 2 RESOLUTION. 1.58 ANGSTROMS.
REMARK 3
REMARK 3 REFINEMENT.
REMARK 3 PROGRAM : SHELXL-97
REMARK 3 AUTHORS : G.M.SHELDRICK
REMARK 3
  
```



# Entry from the PDB database (crystallographic info)

```

REMARK 3 DATA USED IN REFINEMENT.
REMARK 3 RESOLUTION RANGE HIGH (ANGSTROMS) : 1.58
REMARK 3 RESOLUTION RANGE LOW (ANGSTROMS) : 20.0
REMARK 3 DATA CUTOFF (SIGMA(F)) : 0.000
REMARK 3 COMPLETENESS FOR RANGE (%) : 94.2
REMARK 3 CROSS-VALIDATION METHOD : THROUGHOUT
REMARK 3 FREE R VALUE TEST SET SELECTION : RANDOM
REMARK
REMARK 290 CRYSTALLOGRAPHIC SYMMETRY
REMARK 290 SYMMETRY OPERATORS FOR SPACE GROUP: P 21 21 2
REMARK 290
REMARK 290 SYMOP SYMMETRY
REMARK 290 NNNMMM OPERATOR
REMARK 290 1555 X, Y, Z
REMARK 290 2555 -X, -Y, Z
REMARK 290 3555 1/2-X, 1/2+Y, -Z
REMARK 290 4555 1/2+X, 1/2-Y, -Z
REMARK 290
REMARK 290 WHERE NNN -> OPERATOR NUMBER
REMARK 290 MMM -> TRANSLATION VECTOR
REMARK 290
REMARK 290 CRYSTALLOGRAPHIC SYMMETRY TRANSFORMATIONS
REMARK 290 THE FOLLOWING TRANSFORMATIONS OPERATE ON THE ATOM/RESIDUE
REMARK 290 RECORDS IN THIS ENTRY TO PRODUCE CRYSTALLOGRAPHICALLY
REMARK 290 RELATED MOLECULES.
REMARK 290 SMTRY1 1 1.000000 0.000000 0.000000 0.000000
REMARK 290 SMTRY2 1 0.000000 1.000000 0.000000 0.000000
REMARK 290 SMTRY3 1 0.000000 0.000000 1.000000 0.000000
REMARK 290 SMTRY1 2 -1.000000 0.000000 0.000000 0.000000
REMARK 290 SMTRY2 2 0.000000 -1.000000 0.000000 0.000000
REMARK 290 SMTRY3 2 0.000000 0.000000 1.000000 0.000000
  
```



## Entry from the PDB database (sequence, sec. elements)

```

DREF 1CV2 A 1 296 DBJ BAA03443 BAA03443 1 296
SEQRES 1 A 296 MET SER LEU GLY ALA LYS PRO PHE GLY GLU LYS LYS PHE
SEQRES 2 A 296 ILE GLU ILE LYS GLY ARG ARG MET ALA TYR ILE ASP GLU
SEQRES 3 A 296 GLY THR GLY ASP PRO ILE LEU PHE GLN HIS GLY ASN PRO
SEQRES 4 A 296 THR SER SER TYR LEU TRP ARG ASN ILE MET PRO HIS CYS
SEQRES 5 A 296 ALA GLY LEU GLY ARG LEU ILE ALA CYS ASP LEU ILE GLY
SEQRES 6 A 296 MET GLY ASP SER ASP LYS LEU ASP PRO SER GLY PRO GLU
SEQRES
HELIX 1 1 SER A 42 ALA A 53
HELIX 2 2 TYR A 82 LEU A 96
HELIX 3 3 TRP A 109 ARG A 120
HELIX 4 4 GLU A 145 ARG A 155
HELIX 5 5 GLY A 159 LEU A 164
HELIX 6 6 VAL A 168 LEU A 177
HELIX 7 7 GLU A 184 GLU A 192
HELIX 8 8 ARG A 202 ILE A 211
HELIX 9 9 ALA A 218 SER A 234
HELIX 10 10 THR A 250 ARG A 258
HELIX 11 11 ILE A 274 ASP A 277
HELIX 12 12 SER A 278 LEU A 293
SHEET 1 81 8 LYS A 12 ILE A 14 O
SHEET 2 81 8 MET A 21 GLU A 26 -1 N MET A 21 O ILE A 14
SHEET 3 81 8 ARG A 57 ASP A 62 -1 N ALA A 60 O ILE A 24
SHEET 4 81 8 ASP A 30 HIS A 36 1 N ILE A 32 O ARG A 57
SHEET 5 81 8 VAL A 102 HIS A 107 1 N VAL A 103 O PRO A 31
SHEET 6 81 8 VAL A 125 MET A 131 1 N ALA A 129 O LEU A 104
SHEET 7 81 8 LYS A 238 PRO A 245 1 N ILE A 241 O TYR A 130
SHEET 8 81 8 GLN A 263 GLY A 270 1 N THR A 264 O LYS A 238
CISPEP 1 ASN A 38 PRO A 39 0 -2.50
CISPEP 2 ASP A 73 PRO A 74 0 -2.40
CISPEP 3 THR A 216 PRO A 217 0 -3.84
CISPEP 4 GLU A 244 PRO A 245 0 3.01
CISPEP 5 PRO A 295 ALA A 296 0 20.14

```

## Entry from the PDB database (co-ordinates)

```

ATOM 1 N GLY A 4 7.096 3.531 6.684 1.00 20.14 N
ATOM 2 CA GLY A 4 7.885 3.530 5.461 1.00 21.48 C
ATOM 3 C GLY A 4 9.300 3.633 1.00 20.62 C
ATOM 4 O GLY A 4 9.599 4.865 6.501 1.00 12.02 O
ATOM 5 N ALA A 5 10.240 3.571 4.814 1.00 15.21 N
ATOM 6 CA ALA A 5 11.609 4.057 4.935 1.00 10.23 C
ATOM 7 C ALA A 5 11.883 5.182 3.955 1.00 11.96 C
ATOM 8 O ALA A 5 12.950 5.809 3.978 1.00 13.69 O
ATOM 9 CB ALA A 5 12.621 2.943 4.674 1.00 10.47 C
ATOM 10 N LYS A 6 10.929 5.437 3.056 1.00 12.31 N
ATOM 11 CA LYS A 6 11.251 6.452 2.053 1.00 17.87 C
ATOM 12 C LYS A 6 11.223 7.850 2.660 1.00 9.09 C
ATOM 13 O LYS A 6 10.310 8.161 3.422 1.00 10.53 O
ATOM 14 CB LYS A 6 10.274 6.419 0.870 1.00 20.08 C
ATOM 15 CG LYS A 6 10.901 6.898 -0.436 1.00 47.79 C
ATOM 16 CD LYS A 6 10.695 8.377 -0.703 1.00 63.02 C
ATOM 17 CE LYS A 6 11.654 8.950 -1.734 1.00 62.33 C
ATOM 18 NE LYS A 6 11.574 10.435 -1.832 1.00 50.12 N
ATOM 19 N PRO A 7 12.171 8.696 2.307 1.00 12.48 N
ATOM 20 CA PRO A 7 12.108 10.087 2.748 1.00 15.50 C
ATOM 21 C PRO A 7 10.895 10.808 2.144 1.00 16.27 C
ATOM 22 O PRO A 7 10.244 10.396 1.170 1.00 15.29 O
ATOM 23 CB PRO A 7 13.394 10.717 2.217 1.00 12.40 C
ATOM 24 CG PRO A 7 13.877 9.803 1.151 1.00 23.02 C
ATOM 25 CD PRO A 7 13.347 8.427 1.456 1.00 21.52 C
ATOM 26 N PHE A 8 10.408 11.842 2.751 1.00 10.86 N
ATOM 27 CA PHE A 8 9.557 12.848 2.302 1.00 6.69 C
ATOM 28 C PHE A 8 10.134 13.914 1.384 1.00 17.38 C
ATOM 29 O PHE A 8 11.121 14.590 1.716 1.00 17.05 O
ATOM 30 CB PHE A 8 8.912 13.490 3.531 1.00 6.78 C
ATOM 31 CG PHE A 8 7.776 14.444 3.183 1.00 12.53 C
ATOM 32 CD1 PHE A 8 6.526 13.921 2.874 1.00 18.59 C
ATOM 33 CD2 PHE A 8 7.984 15.811 3.166 1.00 15.33 C
ATOM 34 CE1 PHE A 8 5.494 14.797 2.547 1.00 19.62 C
ATOM 35 CE2 PHE A 8 6.961 16.701 2.851 1.00 15.73 C
ATOM 36 CE PHE A 8 5.718 16.160 2.537 1.00 22.66 C
ATOM 37 N GLY A 9 9.544 14.110 0.215 1.00 18.38 N

```

## Protein structure databases

### ■ PDBsum

- > developed at University College London
- > summaries and analyses of protein structures (secondary database derived from PDB)
- > summary of PDB entries: resolution, R-factor, # protein chains, topology, ligands, metal ions, etc.
- > analysis of PDB entries: protein-metal and protein-ligand interactions, protein validation
- > provides links to many related databases

## Information on entry from the PDBsum database

**PDBsum** 1cv2



**PDB id: 1cv2**  
**Hydrolase**  
**Title:** Hydrolytic halalkane dehalogenase *hhd* from sphingomonas paucimobilit ut29 at 1.6 Å resolution  
**Structure:** Halalkane dehalogenase. Chain: a. Synonym: *hhd*, 2,3,4,5-tetrachloro-1,4-cyclohexadiene hydrolase. Engineered yes  
**Source:** Sphingomonas paucimobilit. Strain: ut29. Expressed in: escherichia coli.  
**Resolution:** 1.50Å. **R-factor:** 0.152. **R-free:** 0.211.  
**Authors:** J Marek, J Vervandvo, J D Ambarsky, L Smatanova, L A Svensson, J Niemann, Y Nagata, M Takagi  
**Date:** 22-Aug-00

[PDB header records](#)

Enzyme class from PDB file: [EC:3.1.1.5](#) [E.C.-->PDB](#)

**Links:**

- [PDB](#)
- [DCA](#)
- [MMDB](#)
- [TM6 Jena](#)
- [STRINE](#)
- [GRASS](#)
- [PQS](#)
- [CATH](#)
- [SCOP](#)
- [FSSP](#)
- [PROCHECK](#)
- [PROCHECK2](#)
- [PROSAWEB](#)

---

---

---

---

---

---

---

---

---

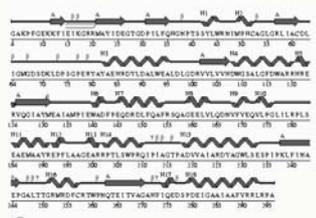
---

## Entry from the PDBsum database (secondary elements)

**Chain A (299 residues)**

- CATH** structural classification (1 domain):

Links	CATH no.	Class	Architecture
<a href="#">CATH file</a>	<a href="#">3.40.50.950</a>	Alpha Beta	3-Layer(alpha) Sandwich



**PROMOTIF summary:**  
1 sheet, 3 strands, 10 helices, 19 beta turns, 3 gamma turns, 1 beta bulge, 1 beta hairpin, 4 beta alpha beta units, 1 psi-loop

---

---

---

---

---

---

---

---

---

---

## Protein structure classification databases

- Classification attempts to capture the structural similarities among proteins.
- The structural similarities relate to the evolution.
- The structural similarities may imply function.
- The classification scheme is dependent on the underlying philosophy.

---

---

---

---

---

---

---

---

---

---

## Protein structure classification databases

### ■ SCOP - Structural Classification of Proteins

- developed at MRC Laboratory of Molecular Biology
- construction: combination of manual and automatic methods (complicated by multidomain proteins)
  
- fold = same secondary elements in same arrangement, independently of common evolutionary origin
- superfamily = low identity but common evolutionary origin implied from common structure and function
- family = sequence identity >30%

---

---

---

---

---

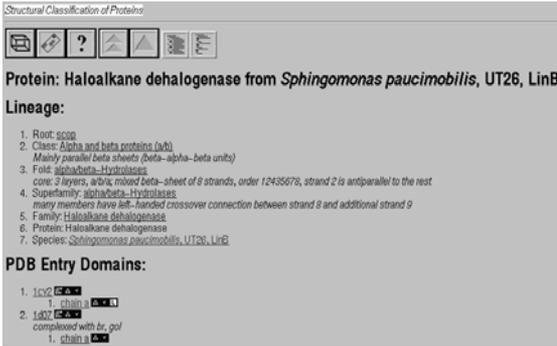
---

---

---

### Information on entry from the SCOP database

Structural Classification of Proteins



**Protein:** Haloalkane dehalogenase from *Sphingomonas paucimobilis*, UT26, LinB

**Lineage:**

1. Root: scop
2. Class: Alpha and beta proteins (αβ)  
Mainly parallel beta sheets (beta-alpha-beta units)
3. Fold: alpha/beta-Hydrolases  
core: 2 layers, alpha, mixed beta-sheet of 8 strands, order 12435678, strand 2 is antiparallel to the rest
4. Superfamily: alpha/beta-Hydrolases  
many members have left-handed crossover connection between strand 8 and additional strand 9
5. Family: Haloalkane dehalogenase
6. Protein: Haloalkane dehalogenase
7. Species: *Sphingomonas paucimobilis*, UT26, LinB

**PDB Entry Domains:**

1. 1cxf (E.C. 3.1.1.1)  
1, chain a (E.C. 3.1.1.1)
2. 1d02 (E.C. 3.1.1.1)  
complexed with ac. gal  
1, chain a (E.C. 3.1.1.1)

---

---

---

---

---

---

---

---

## Protein structure classification databases

### ■ CATCH - Class, Architecture, Topology, Homology

- developed at University College London
- construction: mostly automatic
- unique numbering scheme analogous to Enzyme Classification (E.C.) scheme
  
- class = gross secondary structure content
- architecture = gross secondary structure arrangement
- topology = shape and connectivity of secondary structures (60% of larger protein matches smaller one)
- homology = sequence identity >35%, common ancestry
- sequence = clustering based on sequence identity

---

---

---

---

---

---

---

---

## Information on entry from the CATCH database

Home > Top > C[3] > A[4] > T[5] > H[5] > S[1] > N[2] > [1]

**Domain 1cv2A0**

- Alpha Beta
- 3-Layer[alpha] Sandwich
- Rossmann fold
- Alpha/beta hydrolase
- HYDROLASE
- HYDROLASE
- HYDROLASE

**Fold relatives**

These are either no other non-identical relatives within this fold group or the structural comparisons for this domain have not yet been calculated.

View as XML Search

1cv2A0 View Rasmol

Goto...  
SSAP  
DHS  
Gene3D  
PDBsum

Navigation  
Top of hierarchy  
Up one level

Help  
Select a topic

---

---

---

---

---

---

---

---

## Prediction of secondary structure

- Algorithms assign probability for occurrence of  $\alpha$ -helix,  $\beta$ -strand, turn and random coil at particular position in the sequence.
- Methods: statistical, stereochemical and homology/neural networks based. All methods rely on information derived from known 3D structures. Most recent methods use the information from multiple alignments.
- Reliability of the best current methods is >70%.

---

---

---

---

---

---

---

---

## Prediction of secondary structure

- Chou-Fasman and GOR
  - statistical - amino acids show preference for particular secondary structure elements
- PHD and NNPredict
  - neural networks - the rules for prediction are not defined in advance, they are created by training
- NNSSP and PREDATOR
  - nearest neighbour approach
- JPRED
  - consensus approach - utilises multiple alignments and state-of-art method - makes consensus

---

---

---

---

---

---

---

---



## Prediction of protein fold

### ■ Bioinbgu

- > consensus method utilising predictions from five different algorithms

### ■ 3D-PSSM

- > scoring functions: 1D-PSSMs (sequence profiles built from relatively close homologues), 3D-PSSMs (more general profiles containing more remote homologues), matching of secondary structure elements, and propensities of the residues for solvent accessibility

### ■ GenThreader

- > hybrid method: profile-based alignment, evaluation of alignments by threading, evaluation of threaded models by neural network

---

---

---

---

---

---

---

---

## Prediction of tertiary structure

### ■ *Ab initio*

- > 3D structure of a protein is predicted from first principles (search for global minimum structure)
- > current algorithms are not very reliable

### ■ Homology modelling

1. alignment of modelled sequence against sequences of structurally similar proteins (templates)
2. "extraction" of the backbone from template structure and positioning of side-chain
3. modelling of loops
4. structure refinement and validation

---

---

---

---

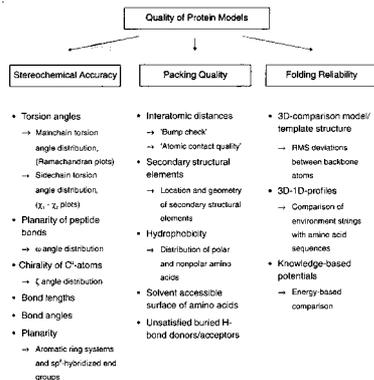
---

---

---

---

## Validation of protein models



---

---

---

---

---

---

---

---

## Prediction of tertiary structure

### ■ SWISS-MODEL

- > fully automated modelling server
- > input = protein sequence; output = PDB file
- > 1. search of ExNRL-3D using BLASTP for potential templates; 2. select all templates with sequence identities above 25%; 3. Generate structures of 3D models; 4. energy minimise models using GROMOS 96
- > first approach and optimise mode (Swiss-PDBViewer)

### ■ MODELLER

- > most widely used academic program for homology modelling (satisfaction of spatial restrains)

---

---

---

---

---

---

---

---

## Modelling of protein-ligand complexes

### ■ Docking

- > positioning of small organic molecules (ligands) inside the protein active site
- > different orientations and conformations of the ligand are evaluated using geometric or energetic scoring functions
- > Protein-ligand interaction energy = van der Waals term + electrostatic term + H-bond term + entropic term
- > flexible docking - considers different conformation of ligand; different rotamers of protein side chains

### ■ Software: DOCK, AUTODOCK, FLEX

---

---

---

---

---

---

---

---