## Protein Bioinformatics
### Part I: Access to information

260.655
March 30, 2010
Jonathan Pevsner, Ph.D.
pevsner@kennedykrieger.org

---

## Outline for today

Introduction

Accessing information
- Entrez Gene
- Accession numbers and RefSeq
- Protein Databases: UniProt, ExPASy
- Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
- Perspective 1: Protein families (domains and motifs)
- Perspective 2: Physical properties (3D structure)
- Perspective 3: Localization
- Perspective 4: Function

---

## Course objectives

To provide students with the ability to analyze and understand data from high-throughput proteomics experiments. At the conclusion of the course the students will be able to:

(a) Define protein physical properties and analyze protein structure.

(b) Explain how proteins are studied experimentally and how data are generated in high-throughput experiments.

(c) Describe the computational methods used to study protein structure and interactions.

(d) Explain the algorithms, statistical techniques and software tools used to analyze high-throughput proteomics data.

## Syllabus (through April)

| | |
|---|---|
| Tues 3/30 | Protein bioinformatics I (Pevsner) |
| Thurs 4/1 | Protein bioinformatics II: Evolution (Pevsner) |
| | |
| Tues 4/6 | Physical properties of amino acids (Prigge) |
| Thurs 4/8 | Protein structure essentials (Prigge) |
| | |
| Tues 4/13 | How to visualize proteins (Prigge) |
| Thurs 4/15 | Why proteins fold (Prigge) |
| | |
| Tues 4/20 | Structure determination and databases (Prigge) |
| Thurs 4/22 | Crystallography practicum (Prigge/Bosch) |
| | |
| Tues 4/27 | Quantitative proteomics (Cole) |
| Thurs 4/29 | Proteomics and systems biology (Bosch) |

## Syllabus (through May)

| | |
|---|---|
| Tues 5/4 | Protein Structure: Databases & classification (Ruczinski) |
| Thurs 5/6 | Protein secondary struct. prediction (Ruczinski) |
| | |
| Tues 5/11 | Protein tertiary structure prediction (Ruczinski) |
| Thurs 5/13 | Protein structure prediction (CASP) (Ruczinski) |
| | |
| Tues 5/18 | Review (Prigge/Ruczinski/Pevsner) |
| Thurs 5/20 | Final Exam + Practicum |

## Website

**The course website is:**
  http://www.biostat.jhsph.edu/~iruczins/teaching/260.655/

(or Google "ingo teaching")

## Literature references

You are encouraged to read original source articles. They will enhance your understanding of the material. Readings are optional but recommended.

## Computer labs

There are several computer labs (details to follow).

## Grading

Grading is based on assignments and on a final exam.

## What is bioinformatics?

• Interface of biology and computers

• Analysis of proteins, genes and genomes using computer algorithms and computer databases

• Genomics is the analysis of genomes. The tools of bioinformatics are used to make sense of the billions of base pairs of DNA that are sequenced by genomics projects.

• Protein bioinformatics refers to the use of computational biology tools to understand protein structure and function, including high throughput approaches

---

## Protein bioinformatics spans the central dogma…
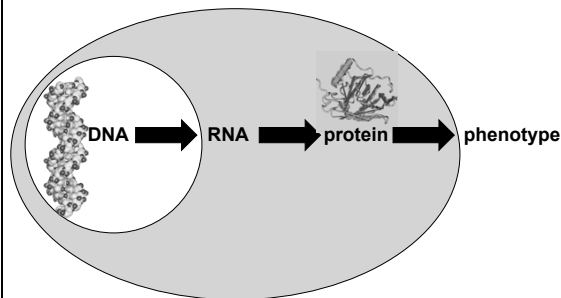


DNA ➡ RNA ➡ protein ➡ phenotype

Page 5

---

## …Protein bioinformatics spans the tree of life



BACTERIA

ARCHAEA

Hot Life

Visible Life

EUCARYA

0.1 changes per site

After Pace NR (1997)
*Science* 276:734

Page 6

4

**Growth of GenBank + Whole Genome Shotgun**
**(1982-November 2008)**

Number of sequences in GenBank (millions) □

Base pairs of DNA in GenBank (billions) ◆
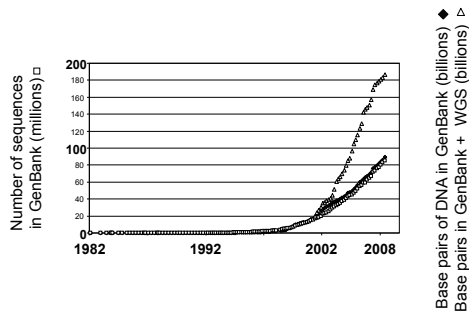Base pairs in GenBank + WGS (billions) △

200
180
160
140
120
100
80
60
40
20
0

1982    1992    2002  2008

Fig. 2.1
Page 15

---

Arrival of next-generation sequencing:
approaching 100 terabases (100,000 gigabases) in 2009

Sequencing center (2009)

GenBank + WGS

3 papers (Nov. 2008)

1000 Genomes

Terabases of DNA sequence (Tb)

12
10
8
6
4
2
0

Fig. 2.1
Page 15

---

DNA → RNA → protein → phenotype

cDNA
ESTs
UniGene

protein
sequence
databases

genomic
DNA
databases
(infer protein
sequences)

Fig. 2.2
Page 18

**Outline for today**

Introduction

Accessing information
   Entrez Gene
   Accession numbers and RefSeq
   Protein Databases: UniProt, ExPASy
   Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
   Perspective 1: Protein families (domains and motifs)
   Perspective 2: Physical properties (3D structure)
   Perspective 3: Localization
   Perspective 4: Function

---

**New NCBI homepage (November 2009):
To study a protein, try starting with Entrez Gene**



www.ncbi.nlm.nih.gov

---

**From the NCBI home page, type "beta globin" and hit "Search"**

**Follow the link to "Gene"**

Fig. 2.5
Page 28



Entrez Gene is in the header
Note the "Official Symbol" HBB for beta globin
Note the "limits" option



Entrez Gene (top of page)

Note that links to many other HBB database entries are available

Page 30

Entrez Gene (bottom of page): RefSeq accession numbers



## Outline for today

Introduction

Accessing information
      Entrez Gene
      Accession numbers and RefSeq
      Protein Databases: UniProt, ExPASy
      Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
      Perspective 1: Protein families (domains and motifs)
      Perspective 2: Physical properties (3D structure)
      Perspective 3: Localization
      Perspective 4: Function

## Access to sequences: Entrez Gene at NCBI

Entrez Gene is a great starting point: it collects
key information on each gene/protein from
major databases. It covers all major organisms.

RefSeq provides a curated, optimal accession number
for each DNA (NM_000518 for beta globin DNA
corresponding to mRNA) or protein (NP_000509)

Page 29

## 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 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 | **DNA** |
| NT_030059 | Genomic contig | |
| Rs7079946 | dbSNP (single nucleotide polymorphism) | |
| N91759.1 | An expressed sequence tag (1 of 170) | **RNA** |
| NM_006744 | RefSeq DNA sequence (from a transcript) | |
| NP_007635 | RefSeq protein | **protein** |
| 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 |

## Slide 1

Entrez Gene (bottom of page): non-RefSeq accessions
(it's unclear what these are, highlighting usefulness of RefSeq)



| | | |
|---|---|---|
| Genomic | M36640.1 | AAA52829.1 |
| Genomic | S82767.1 | AAD14420.1 |
| Genomic | U01317.1 | AAA16334.1 |
| | | AAA16335.1 |
| Genomic | U01317.1 | AAA16334.1 |
| | | AAA16335.1 |
| Genomic | U00223.1 | AAB60348.1 |
| Genomic | V00498.1 | CAA23757.1 |
| Genomic | V00499.1 | CAA23758.1 |
| mRNA | AF117710.1 | AAD19696.1 |
| mRNA | AF181832.1 | AAF00408.1 |
| mRNA | AF181989.1 | AAF00409.1 |
| mRNA | AF349114.1 | AAK29638.1 |
| mRNA | AY311605.1 | BAG34767.1 |
| mRNA | AY136510.1 | AHN11300.1 |
| mRNA | AY509193.1 | AAR96398.1 |
| mRNA | BC007075.1 | AAH07075.1 |
| mRNA | CR536530.1 | CAG38767.1 |
| mRNA | CR541913.1 | CAG46711.1 |
| mRNA | CR590940.1 | None |
| mRNA | CR594264.1 | None |
| mRNA | CR603426.1 | None |
| mRNA | CR609101.1 | None |
| mRNA | CR621601.1 | None |
| mRNA | EU694432.1 | ACD29349.1 |
| mRNA | M11428.1 | AAA52633.1 |
| mRNA | M25079.1 | AAA35597.1 |
| mRNA | M25113.1 | AAA35966.1 |
| mRNA | V00497.1 | CAA23756.1 |
| mRNA | V00500.1 | CAA23759.1 |
| Synthetic | AM292537.1 | CAL37415.1 |
| Synthetic | AM393351.1 | CAL38229.1 |
| Synthetic | DQ893159.2 | ABM64085.1 |
| Synthetic | DQ896453.2 | ABM67452.1 |
| Synthetic | EU176774.1 | ABH03575.1 |

Protein Accession — Links
O95408 — GenPept UniProtKB/TrEMBL

## Slide 2

Entrez Protein:
accession,
organism,
literature…



```
LOCUS       NP_000509            147 aa            linear   PRI 12-OCT-2008
DEFINITION  beta globin [Homo sapiens].
ACCESSION   NP_000509
VERSION     NP_000509.1  GI:4504349
DBSOURCE    REFSEQ: accession NM_000518.4
KEYWORDS    .
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
REFERENCE   1  (residues 1 to 147)
  AUTHORS   Bernaudin,F., Verlhac,S., Chevret,S., Torres,M., Coic,L.,
            Arnaud,C., Kamdem,A., Hau,I., Neonato,M.G. and Delacourt,C.
  TITLE     G6PD deficiency, absence of alpha-thalassemia and hemolytic rate at
            baseline are significant independent risk factors for abnormally
            high cerebral velocities in patients with sickle cell anemia
  JOURNAL   Blood (2008) In press
  PUBMED    18772456
  REMARK    GeneRIF: Observational study of gene-disease association. (HuGE
            Navigator)
            Publication Status: Available-Online prior to print
REFERENCE   2  (residues 1 to 147)
  AUTHORS   Crompton,P.D., Traore,B., Kayentao,K., Doumbo,S., Ongoiba,A.,
            Diakite,S.A., Krause,M.A., Doumtabe,D., Kone,Y., Weiss,G.,
            Huang,C.Y., Doumbia,S., Guindo,A., Fairhurst,R.M., Miller,L.H.,
            Pierce,S.K. and Doumbo,O.K.
  TITLE     Sickle Cell Trait is Associated with a Delayed Onset of Malaria:
            Implications for Time-to-Event Analysis in Clinical Studies of
            Malaria
```

Fig. 2.8
Page 31

## Slide 3

Entrez Protein:
…features of a protein, and its sequence
in the one-letter amino acid code

```
     Site            94
                     /site_type="modified"
                     /experiment="experimental evidence, no additional details
                     recorded"
                     /note="S-nitrosylation site"
                     /citation=[7]
     Site            121
                     /site_type="glycosylation"
                     /experiment="experimental evidence, no additional details
                     recorded"
                     /note="glycation site"
                     /citation=[9]
     CDS             1..147
                     /gene="HBB"
                     /gene_synonym="CD113t-C"
                     /coded_by="NM_000518.4:51..494"
                     /db_xref="CCDS:CCDS7753.1"
                     /db_xref="GeneID:3043"
                     /db_xref="HGNC:4827"
                     /db_xref="HPRD:00786"
                     /db_xref="MIM:141900"
ORIGIN
        1 mvhltpeeks avtalwgkvn vdevggealg rllvvypwtq rffesfgdls tpdavmgnpk
       61 vkahgkkvlg afsdglahld nlkgtfatls elhcdklhvd penfrllgnv lvcvlahhfg
      121 keftppvqaa yqkvvagvan alahkyh
//
```

Fig. 2.8
Page 31

Entrez Protein:
You can change the display (as shown)…



Page 31
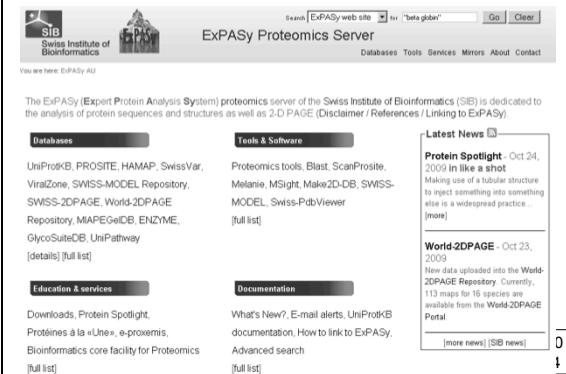
**FASTA format:**
**versatile, compact with one header line**
**followed by a string of nucleotides or amino acids**
**in the single letter code**



>gi|4504349|ref|NP_000509.1| beta globin [Homo sapiens]
MVHLTPEEKSAVTALWGKVNVDEVGGEALGRLLVVYPWTQRFFESFGDLSTPDAVMGNPKVKAHGKKVLG
AFSDGLAHLDNLKGTFATLSELHCDKLHVDPENFRLLGNVLVCVLAHHFGKEFTPPVQAAYQKVVAGVAN
ALAHKYH

Fig. 2.9
Page 32

## Outline for today

Introduction

Accessing information
  Entrez Gene
  Accession numbers and RefSeq
  Protein Databases: UniProt, ExPASy
  Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
  Perspective 1: Protein families (domains and motifs)
  Perspective 2: Physical properties (3D structure)
  Perspective 3: Localization
  Perspective 4: Function

UniProt:
a centralized
protein
database
(uniprot.org)



Page 33

---

ExPASy: vast proteomics resources (www.expasy.ch)



---

**Outline for today**

Introduction

Accessing information
        Entrez Gene
        Accession numbers and RefSeq
        Protein Databases: UniProt, ExPASy
        Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
        Perspective 1: Protein families (domains and motifs)
        Perspective 2: Physical properties (3D structure)
        Perspective 3: Localization
        Perspective 4: Function

Ensembl genome browser (www.ensembl.org)



Ensembl output for beta globin includes views of chromosome 11 (top), the region (middle), and a detailed view (bottom).

There are various horizontal annotation tracks.



[1] Visit http://genome.ucsc.edu/, click Genome Browser



[2] Choose organisms, enter query (beta globin), hit submit

Page 36

[3] Choose the RefSeq beta globin gene

**UCSC Genes**

HBB (uc009yem.1) at chr11:5204383-5212336 - Hemoglobin Lepore-Baltimore (Fragment).
HBB (uc001mae.1) at chr11:5203272-5204877 - beta globin
HBD (uc001maf.1) at chr11:5210635-5212434 - delta globin
RBM17 (uc001isb.1) at chr10:6171013-6198847 - RNA binding motif protein 17
HBA1 (uc002cfx.1) at chr16:166679-167520 - alpha 1 globin
HBA2 (uc002cfv.2) at chr16:162846-163709 - alpha 2 globin
HBA1 (uc002cfw.2) at chr16:162846-167520 - alpha 1 globin
HBBP1 (uc001maq.1) at chr11:5219761-5221398 - Homo sapiens hemoglobin, beta pseudoge

**RefSeq Genes**

HBB at chr11:5203272-5204877 - (NM_000518) beta globin
HBBP1 at chr11:5219761-5221398 - (NR_001589)

---

[4] The UCSC Genome Browser is an essential resource
       --choose which tracks to display
       --add custom tracks
       --the Table Browser is complementary



---

**Example of how to access sequence data:
HIV-1 *pol***

There are many possible approaches. Begin at the main page of NCBI, and type an Entrez query: hiv-1 pol

Page 36

14

**Searching for HIV-1 *pol*: >130,000 nucleotide, protein hits**



**Searching for HIV-1 *pol*:**
**using the command hiv-1[organism] limits the**
**output to just one entry**



Try Taxonomy Browser to easily limit your query to your favorite organism(s). *Example*:
NCBI home → Taxonomy → Taxonomy browser → human → protein to find a human protein



over 300,000 nucleotide entries for HIV-1

only 1 RefSeq

## Example of how to access sequence data: histone

| query for "histone" | # results |
|---|---|
| protein records | 85,000 |
| RefSeq entries | 32,000 |
| RefSeq (limit to human) | 1129 |
| NOT deacetylase | 863 |

At this point, select a reasonable candidate (e.g. histone 2, H4) and follow its link to Entrez Gene. There, you can confirm you have the right protein.

11-09



---

## Entrez Gene result for a histone



---

## Outline for today

Introduction

Accessing information
- Entrez Gene
- Accession numbers and RefSeq
- Protein Databases: UniProt, ExPASy
- Three genome browsers: NCBI, UCSC, Ensembl

Four perspectives on individual proteins
- Perspective 1: Protein families (domains and motifs)
- Perspective 2: Physical properties (3D structure)
- Perspective 3: Localization
- Perspective 4: Function

**Perspective 1:**
**Protein domains and motifs**

---

**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

## 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 |

Source: Integr8 at EBI website

Page 391

---

## 15 most common domains (various species)

The European Bioinformatics Institute (EBI) offers many key proteomics resources at the Integr8 site:

http://www.ebi.ac.uk/proteome/

Page 391

---

1. Go to the Integr8 site: http://www.ebi.ac.uk/proteome/

2. Browse species; choose *Homo sapiens*.

3. Click "Proteome analysis"

4. Obtain a variety of statistics, such as common repeats, domains, average protein length

**Amino acid composition**
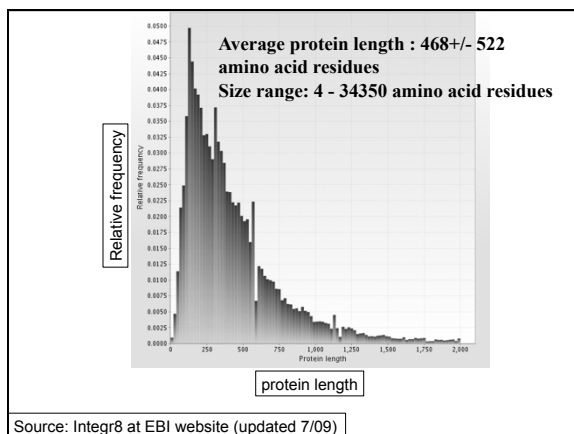


frequency / amino acid

| Amino acid | Total content | Frequency [%] |
|---|---|---|
| A | 1949808 | 6.95 |
| B | 67 | 0.00 |
| C | 627455 | 2.24 |
| D | 1330337 | 4.75 |
| E | 1996677 | 7.12 |
| F | 999412 | 3.56 |
| G | 1847364 | 6.59 |
| H | 731493 | 2.61 |
| I | 1207363 | 4.31 |
| K | 1596772 | 5.70 |
| L | 2761668 | 9.85 |
| M | 602128 | 2.15 |
| N | 998945 | 3.56 |
| P | 1797847 | 6.41 |
| Q | 1346588 | 4.80 |
| R | 1598017 | 5.70 |
| S | 2362325 | 8.43 |
| T | 1520103 | 5.42 |
| U | 57 | 0.00 |
| V | 1680528 | 5.99 |
| W | 348369 | 1.24 |
| X | 1920 | 0.01 |
| Y | 730054 | 2.60 |
| Z | 67 | 0.00 |

Source: Integr8 at EBI website (updated 7/09)

---



**Average protein length : 468+/- 522 amino acid residues**
**Size range: 4 - 34350 amino acid residues**

Relative frequency / protein length

Source: Integr8 at EBI website (updated 7/09)

---

### Definition of a domain

According to InterPro at EBI (http://www.ebi.ac.uk/interpro/):

A domain is an independent structural unit, found alone or in conjunction with other domains or repeats. Domains are evolutionarily related.

According to SMART (http://smart.embl-heidelberg.de):

A domain is a conserved structural entity with distinctive secondary structure content and a hydrophobic core. Homologous domains with common functions usually show sequence similarities.

Page 390

**Varieties of protein domains**

Extending along the length of a protein

Occupying a subset of a protein sequence

Occurring one or more times

Page 393

---

**Example of a protein with domains:**
**Methyl CpG binding protein 2 (MeCP2)**

| | MBD | | TRD | |

The protein includes a methylated DNA binding domain (MBD) and a transcriptional repression domain (TRD). MeCP2 is a transcriptional repressor.

Mutations in the gene encoding MeCP2 cause Rett Syndrome, a neurological disorder affecting girls primarily.

Page 393

---

**Result of an MeCP2 blastp search:**
**A methyl-binding domain shared by several proteins**

links to Smart00391, LOAD, and Pfam01429 database entries for methyl-CpG binding domain

Color Key for Alignment Scores

MeCP2

MBD4
MBD2
MBD3

MBD1

**domain**

Page 393

20

## Are proteins that share only a domain homologous?

MeCP2 — MBD — 486 aa
MBD1 — MBD — 605 aa
MBD2 — MBD — 411 aa
MBD2 (testis) — MBD — 302 aa
MBD3 — MBD — 291 aa
MBD4 — MBD — 580 aa

---

## Proteins can have both domains and motifs (patterns)

**Motif (several residues)**

**Motif (several residues)**

rvp    rvt    rnaseH    rve

**Domain (aspartyl protease)**

**Domain (reverse transcriptase)**

---

### Eukaryotic and viral aspartyl proteases signature and profile

**PROSITE cross-reference(s)**

| | |
|---|---|
| PS00141; ASP_PROTEASE | Retrieve an alignment of Swiss-Prot true positive hits: [Clustal format, color, condensed view] [Clustal format, color] [Clustal format, p] |
| PS50175; ASP_PROT_RETROV | Retrieve an alignment of Swiss-Prot true positive hits: [Clustal format, color, condensed view] [Clustal format, color] [Clustal format, p] |

**Documentation**

Aspartyl proteases, also known as acid proteases, (EC 3.4.23.-) are a widely distributed family of proteolytic enzymes [1,2,3] known to exist in vertebrates, fungi, plants, retroviruses and some plant viruses. Aspartate proteases of eukaryotes are monomeric enzymes which consist of two domains. Each domain contains an active site centered on a catalytic aspartyl residue. The two domains most probably evolved from the duplication of an ancestral gene encoding a primordial domain. Currently known eukaryotic aspartyl proteases are:

- Vertebrate gastric pepsins A and C (also known as gastricsin).
- Vertebrate chymosin (rennin), involved in digestion and used for making cheese.
- Vertebrate lysosomal cathepsins D (EC 3.4.23.5) and E (EC 3.4.23.34).
- Mammalian renin (EC 3.4.23.15) whose function is to generate angiotensin I from angiotensinogen in the plasma.
- Fungal proteases such as aspergillopepsin A (EC 3.4.23.18), candidapepsin (EC 3.4.23.24), mucoropepsin (EC 3.4.23.23) (mucor rennin), endothiapepsin (EC 3.4.23.22), polyporopepsin (EC 3.4.23.29), and rhizopuspepsin (EC 3.4.23.21).
- Yeast saccharopepsin (EC 3.4.23.25) (proteinase A) (gene PEP4). PEP4 is implicated in posttranslational regulation of vacuolar hydrolases.
- Yeast barrierpepsin (EC 3.4.23.35) (gene BAR1); a protease that cleaves alpha-factor and thus acts as an antagonist of the mating pheromone.
- Fission yeast sxa1 which is involved in degrading or processing the mating

| Consensus pattern | [LIVMFGAC]-[LIVMTADN]-[LIVFSA]-D-[ST]-G-[STAV]-[STAPDENQ]-x-[LIVMFSTNC]-x-[LIVMFGTA] [D is the active site residue] |
|---|---|
| Sequences known to belong to this class detected by the pattern | ALL. |
| Other sequence(s) detected in Swiss-Prot | 37. |
| Sequences known to belong to this class detected by the profile | ALL viral-type proteases. |

## Definition of a motif

A motif (or fingerprint) is a short, conserved region of a protein. Its size is often 10 to 20 amino acids.

Simple motifs include transmembrane domains and phosphorylation sites. These do not imply homology when found in a group of proteins.

PROSITE (www.expasy.org/prosite) is a dictionary of motifs (there are currently 1600 entries). In PROSITE, a <u>pattern</u> is a qualitative motif description (a protein either matches a pattern, or not). In contrast, a <u>profile</u> is a quantitative motif description. We will encounter profiles in Pfam, ProDom, SMART, and other databases.

Page 394

---

## Summary of Perspective 1: Protein domains and motifs

A signature is a protein category such as a domain or motif.

You can learn about domains at Integr8, and at databases such as InterPro and Pfam.

A motif (or fingerprint) is a short, conserved sequence. You can study motifs at Prosite at ExPASy.

---

## Perspective 2:
## Physical properties of proteins

Page 397

palmitate     farnesyl     myristate     inositol glycolipid

transferrin receptor    ankyrin     PKA     N-CAM
SNAP-25                             thy-1
                                       5'-nucleotidase

Page 398

---

## Physical properties of proteins

Many websites are available for the analysis of individual proteins. ExPASy and ISREC are two excellent resources.

The accuracy of these programs is variable. Predictions based on primary amino acid sequence (such as molecular weight prediction) are likely to be more trustworthy. For many other properties (such as posttranslational modification of proteins by specific sugars), experimental evidence may be required rather than prediction algorithms.

Page 399

---

## Access a variety of protein analysis programs from the top right of the ExPASy home page

**Compute pI/Mw**

**RETB_HUMAN (P02753)**

DE    Plasma retinol-binding protein precursor (PRBP) (RBP).
OS    Homo sapiens (Human).

The computation has been carried out on the complete sequence.

**Molecular weight:** 22867.85

**Theoretical pI:** 5.48

Page 399

23

**Coils output for syntaxin**

[ISREC-Server] Date: Sat Oct 27 20:52:40 MET 2001

```
# COILS version 2.1
# using MTIDK matrix
# no weights
# Input file is ../wwwtmp/.COILS.27003.1040.seq
#>syntaxin, 288 bases, 5DFEEF76 checksum.
```

Coils output for syntaxin

---

## Protein secondary structure

Protein secondary structure is determined by the amino acid side chains.

Myoglobin is an example of a protein having many $\alpha$-helices. These are formed by amino acid stretches 4-40 residues in length.

Thioredoxin from *E. coli* is an example of a protein with many $\beta$ sheets, formed from $\beta$ strands composed of 5-10 residues. They are arranged in parallel or antiparallel orientations.

---



**Myoglobin
(John Kendrew, 1958)**

Fig. 11.3
Page 427

**Thioredoxin**

## Secondary structure prediction

Chou and Fasman (1974) developed an algorithm based on the frequencies of amino acids found in $\alpha$ helices, $\beta$-sheets, and turns.

Proline: occurs at turns, but not in $\alpha$ helices.

GOR (Garnier, Osguthorpe, Robson): related algorithm

Modern algorithms: use multiple sequence alignments and achieve higher success rate (about 70-75%)

## Secondary structure prediction

Web servers:

GOR4
Jpred
NNPREDICT
PHD
Predator
PredictProtein
PSIPRED
SAM-T99sec

```
                    10        20        30        40        50        60        70
                     |         |         |         |         |         |         |
3DSEQ|pdb1pboA|pdb1pboA AQEEEAEQNLSELSGPWRTVYIGSTNPEKIQENGPFRTYFRELVFDDEKGTVDFYFSVKRDGKWKNVHVK
DPM                    cchhhhhhchcctctcceeeeetttctccchtcctccccehhhehhtchcccceeeeeeetttcctcccceeh
DSC                    cccccccccceeeeeccccceeeeecccceeeecccccccceeeeeeccccceeeeeecccccceeeee
GOR4                   cccccchhhhhhccccceeeeccccccceeeecccccccccceeeeeecccccceeeeeccccccceeeee
HNNC                   cchhhhhhhhhhcccceeeeccccccccccccccehehehecccccceeeeeeccccceeeeecccceeeee
PHD                    ccccccccccceeeeecccccccccccccchhheeeeeeeecccccceeeeeeccccceeeeecccceeeee
Predator               ccchhhhhhhcccceeeeeccccccccccccchhhhhhhhccccceeeeeeccccceeeeeccccceeeeee
SIMPA96                cchhhhhhhhhhcccceeeeeccccccchhhhhhheeccccceeeeeeeccccceeeeeccccceeeee
SOPM                   hhhhhhhhhhhhhcccceeeeeeccccttcccttcccchhhhheeecttccceeeeeecctttcccceeee
Sec.Cons.              cc?hhhhhhhhhcccceeeeecccccccccccch?e?heeccccccccceeeeeeccccccc?eeeee

                          80        90        100       110       120       130       140
                           |         |         |         |         |         |         |
3DSEQ|pdb1pboA|pdb1pboA ATKQDDGTYVADYEGQNVFKIVSLSRTHLVAHNIVDKHGQKTELTGLFVKLNVEDEDLEKFWKLTEDKG
DPM                    hhttttccccchctctcceeeeeeeeecehechhcstetcccccchhecceehhhhhhhhhhhhhcccc
DSC                    eeccccccceeeeecccceeeeeccceeeeeeeeeccccceeeeeecccccccchhhhheeeccc
GOR4                   eecccccceeeecccceeeeecccccecccchhhhhcccccccchhhhheeeccccchhhhhhhhccc
HNNC                   eecccccceeeecccccceeeeccceeeheeccccccccccccccchhhhhhhhhhhhhhhcc
PHD                    eeecccceeeeeccceeeeeeeccccceeeeeecccceeeeeeeeeccccchhhhhhhhhhhhhhc
Predator               eeecccceeeecccceeeeeccccccccccccccccchhhhhhcccchhhhhhhhhhcc
SIMPA96                cccccccceeeecccceeeeccccceeeccccccchhhhhhheeecccchhhhhhhhhhhcctt
SOPM                   eccttteeeeeettcceeeeeeccceeeeeeeettccchhhheeeeeecchhhhhhhhhhcctt
Sec.Cons.              eeccccceeeeecccceeeee?ccceeeeec?c?cccccc?hh?eeeeeccccchhhhhhhhhhhccc

                          150
                           |
3DSEQ|pdb1pboA|pdb1pboA IDKKNVVHFLENEDHPH
DPM                    ctcccceehhhht
DSC                    cccceeeeeccc
GOR4                   cccceeeeeecc
HNNC                   ccccchhhcccc
PHD                    cccceeeeccc
Predator               ccccceeeccc
SIMPA96                cchhhhhhhhhc
SOPM                   cchhhhhhhhhtc
Sec.Cons.              cccc?eeee?ccccccc
```

**Go to http://pbil.univ-lyon1.fr/, click "Secondary structure prediction" to access this prediction tool**

---

## Tertiary protein structure: protein folding

Main approaches:

[1] Experimental determination
(X-ray crystallography, NMR)

[2] Prediction

► Comparative modeling (based on homology)

► Threading

► *Ab initio* (de novo) prediction

Page 430

---

## Experimental approaches to protein structure

[1] X-ray crystallography
-- Used to determine 80% of structures
-- Requires high protein concentration
-- Requires crystals
-- Able to trace amino acid side chains
-- Earliest structure solved was myoglobin

[2] NMR
-- Magnetic field applied to proteins in solution
-- Largest structures: 350 amino acids (40 kD)
-- Does not require crystallization

Page 430

**Steps in obtaining a protein structure**

Target selection

Obtain, characterize protein

Determine, refine, model the structure

Deposit in repository

Fig 11.5
page 431

---

**The Protein Data Bank (PDB)**

- PDB is the principal repository for protein structures
- Established in 1971
- Accessed at http://www.rcsb.org/pdb or simply http://www.pdb.org
- Currently contains 64,000 structure entities

Updated 3/26/10

Page 434

---



Fig. 11.7
Page 435

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



Updated 3/26/10

Fig. 9.6
Page 281

---

**PDB holdings (12/08)**

| | |
|---|---|
| 50,621 | proteins, peptides |
| 2,225 | protein/nucl. complexes |
| 1,946 | nucleic acids |
| 33 | other; carbohydrates |
| 54,825 | total |

Table 11-4
Page 435

---



Structure Explorer - 1PBO

Fig. ~11.10
Page 436

## Viewing hemoglobin (accession 2H35) at PDB



## Viewing structures at PDB: WebMol



Fig. 11.11
Page 437

**gateways to access PDB files**

**Swiss-Prot, NCBI, EMBL**

**Protein Data Bank**

**CATH, Dali, SCOP, FSSP**

**databases that interpret PDB files**

## The CATH Hierarchy



C
α    α&β    β

A
TIM barrel    Sandwich    Roll

T
flavodoxin
(4fxn)
β–lactamase
(1mblA1)

Fig. 11.18
Page 444

---

## Access to PDB through NCBI

You can access PDB data at the NCBI several ways.

• Go to the Structure site, from the NCBI homepage
• Use Entrez
• Perform a BLAST search, restricting the output
  to the PDB database

Page 437

---



Fig. ~11.12
Page 438

**Do a blastp search;**
**set the database to pdb (Protein Data Bank)**



Structure
links

Structure accession
(e.g. 2JTZ)



Fig. 9.14
Page 289

Structure Summary
MMDB



**Access to PDB structures through NCBI**

Molecular Modeling DataBase (MMDB)

Cn3D ("see in 3D" or three dimensions):
structure visualization software

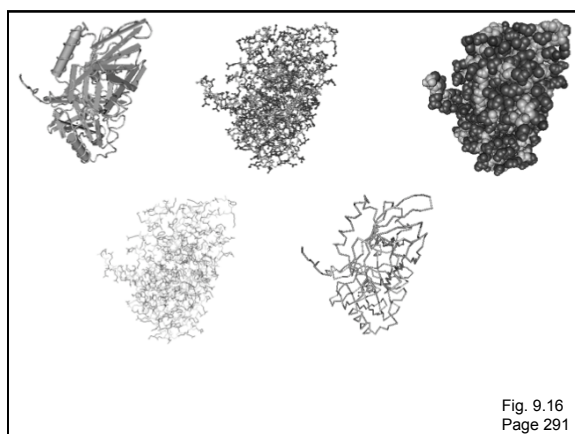Vector Alignment Search Tool (VAST):
view multiple structures

Page 291



Fig. 9.16
Page 291

**Introduction to Perspectives 3 and 4:**
**Gene Ontology (GO) Consortium**

Page 237

---

**The Gene Ontology Consortium**

An ontology is a description of concepts. The GO
Consortium compiles a dynamic, controlled vocabulary
of terms related to gene products.

There are three organizing principles:
    Molecular function
    Biological process
    Cellular compartment

You can visit GO at http://www.geneontology.org.
There is no centralized GO database. Instead, curators
of organism-specific databases assign GO terms
to gene products for each organism.

Page 237

---

**GO terms are assigned to Entrez Gene entries**

GeneOntology                                          Provided by GOA

| Function | | Evidence | |
|---|---|---|---|
| heme binding | | IEA | |
| hemoglobin binding | | IDA | PubMed |
| iron ion binding | | IEA | |
| metal ion binding | | IEA | |
| molecular_function | | ND | |
| oxygen binding | | IDA | PubMed |
| oxygen binding | | IEA | |
| oxygen transporter activity | | IEA | |
| oxygen transporter activity | | NAS | PubMed |

| Process | | Evidence | |
|---|---|---|---|
| biological_process | | ND | |
| nitric oxide transport | | NAS | PubMed |
| oxygen transport | | IEA | |
| oxygen transport | | NAS | PubMed |
| oxygen transport | | TAS | PubMed |
| positive regulation of nitric oxide biosynthetic process | | NAS | PubMed |
| regulation of blood pressure | | IEA | |
| regulation of blood vessel size | | IEA | |
| transport | | IEA | |

| Component | | Evidence | |
|---|---|---|---|
| hemoglobin complex | | IEA | |
| hemoglobin complex | | NAS | PubMed |
| hemoglobin complex | | TAS | PubMed |

Page 241

## HBB

protein from *Homo sapiens* (human)

Term associations ↓ Gene product information → Peptide Sequence → Sequence information →

### Term Associations

gene association format  RDF-XML

▼ Filter associations displayed

Filter Associations
Ontology: All / biological process / cellular component / molecular function
Evidence Code: All / IC / IDA / EXP

[Set filters] [Remove all filters]

[Select all] [Clear all] Perform an action with the selected terms... [Go!]

| | Accession, Term | Ontology | Qualifier | Evidence | Reference | Assigned by |
|---|---|---|---|---|---|---|
| ☐ | 165985 gene products / view in tree / GO:0008150 : biological_process | biological process | | ND | UniProtKB:Q9UPB1 | UniProtKB |
| ☐ | GO:0005833 : 27 gene products / hemoglobin / view in tree / complex | cellular component | | NAS | UniProtKB:Q9UPB1 | UniProtKB |
| ☐ | 164057 gene products / view in tree / GO:0003674 : molecular_function | molecular function | | ND | UniProtKB:Q9UPB1 | UniProtKB |

[Select all] [Clear all] Perform an action with the selected terms... [Go!]

Page 241

---

## The Gene Ontology Consortium: Evidence Codes

IC    Inferred by curator
IDA    Inferred from direct assay
IEA    Inferred from electronic annotation
IEP    Inferred from expression pattern
IGI    Inferred from genetic interaction
IMP    Inferred from mutant phenotype
IPI    Inferred from physical interaction
ISS    Inferred from sequence or structural similarity
NAS    Non-traceable author statement
ND    No biological data
TAS    Traceable author statement

Page 240

---

## Perspective 3:
## Protein localization

Page 242

34

## Protein localization



protein

## Protein localization

Proteins may be localized to intracellular compartments, cytosol, the plasma membrane, or they may be secreted. Many proteins shuttle between multiple compartments.

A variety of algorithms predict localization, but this is essentially a cell biological question.

# Results of Subprograms
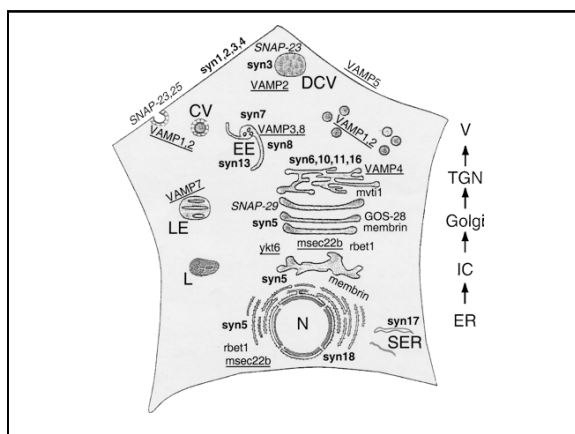
PSG:   a new signal peptide prediction method
       N-region:   length 2;   pos.chg 1;   neg.chg 0
       H-region:   length 14;   peak value   10.03
       PSG score:    5.63

GvH:   von Heijne's method for signal seq. recognition
       GvH score (threshold: -2.1):    3.93
       possible cleavage site: between 16 and 17

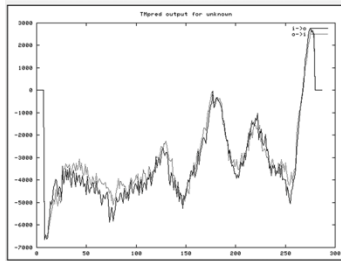>>> Seems to have a cleavable signal peptide (1 to 16)

Page 242

---

2 possible models considered, only significant TM-segments used

-----> slightly prefered model: N-terminus inside
1 strong transmembrane helices, total score : 2757
# from   to length score orientation
1  266  284 (19)    2757 i-o

-----> alternative model
1 strong transmembrane helices, total score : 2690
# from   to length score orientation
1  266  288 (23)    2690 o-i



Page 244