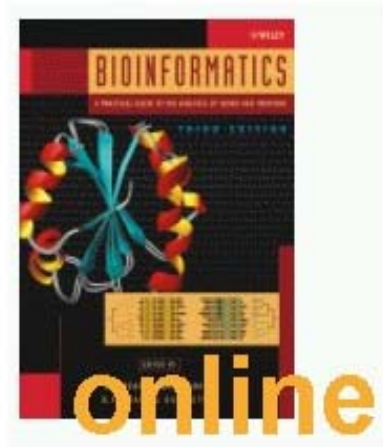


5. Predictive Methods Using Protein Sequences (2)

薛佑玲 Yow-Ling Shiue
國立中山大學生物醫學研究所
✉ ylshiue@mail.nsysu.edu.tw



Select a Chapter:

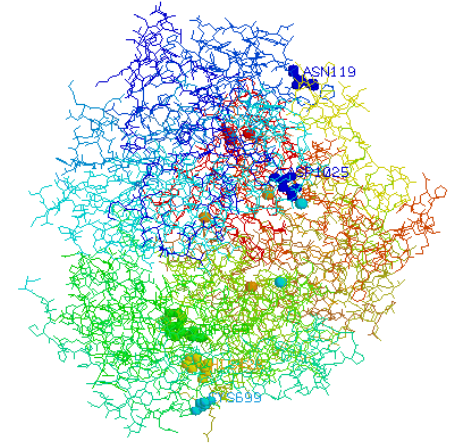
Chapter 8: Predictive Methods Using Protein Sequences

- [Sample Data for Problem Sets](#)
- [Internet Resources](#)

Predicting Function

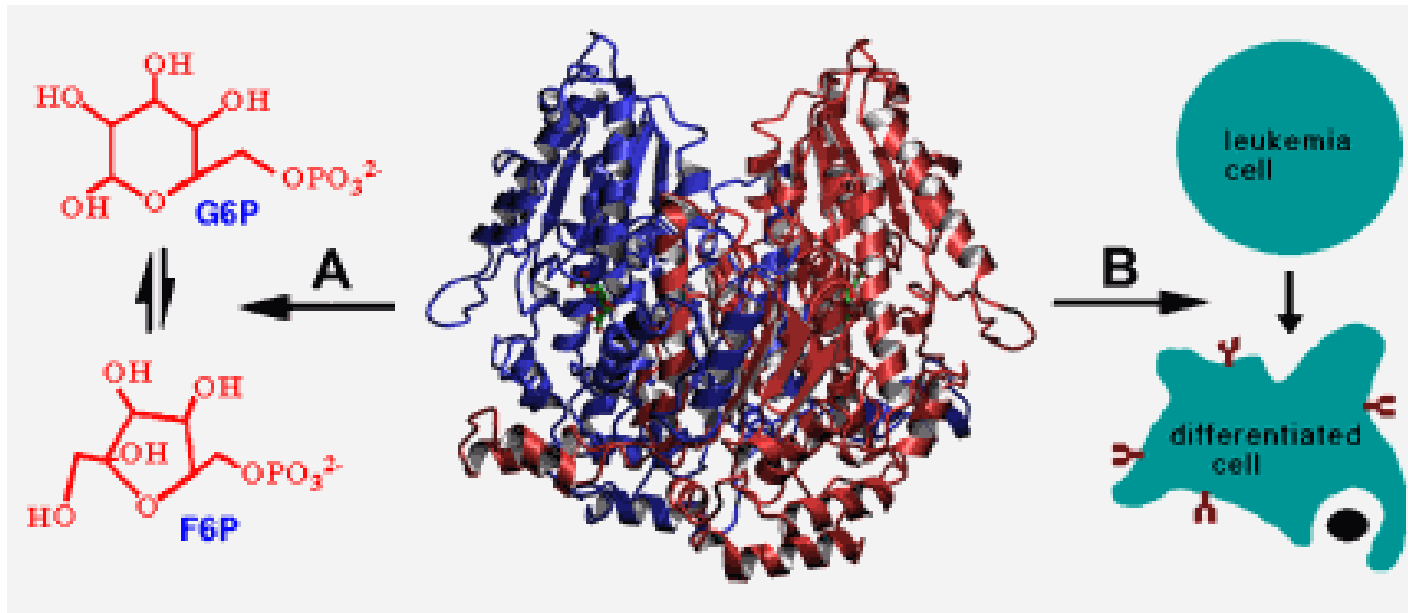
- × To extract **biologically important information** from a protein sequence
 - × **Structure prediction**
 - × To predict **secondary structure** + **solvent accessibility** + **transmembrane helices**...(above)
 - × **Function prediction**
 - × **Case-specific prediction**
 - × Did **not** result in automated tools for **function prediction**, but by
 - × **Annotation transfer**
 - × **Motifs & patterns**

Annotation Transfer (2)

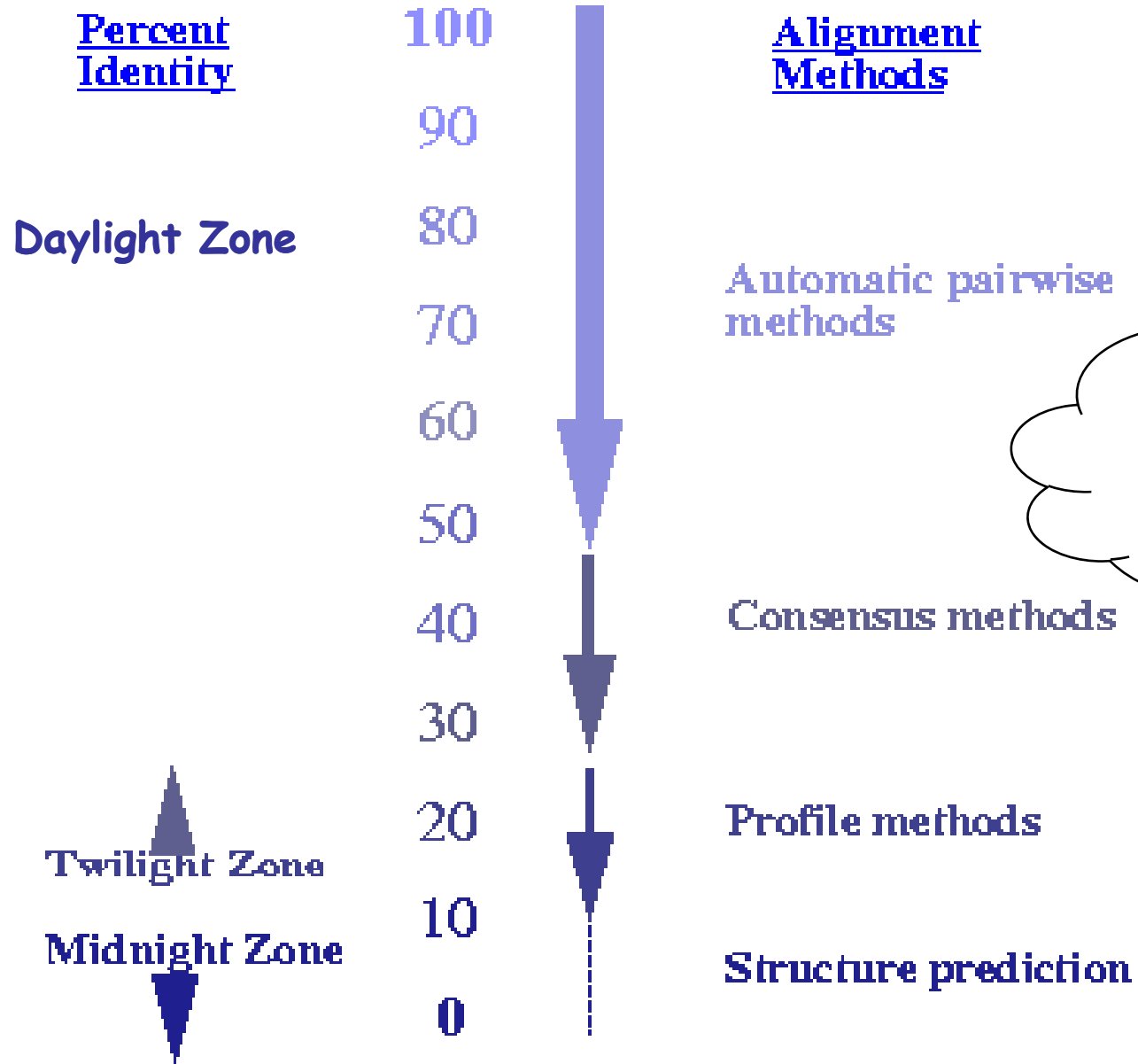


- x Thousands of proteins have been characterized **experimentally** in the lab
 - x Their functions have been recorded in **various protein databases**
 - x **Similarity at the sequence level** implies **similarity of function** (Next slide)
 - x If a **newly discovered protein** bears high **similarity** to another, well-characterized one, it is reasonable to assume that they have a similar function

- x Some **notable exceptions**
 - x Proteins have **different functions** depending on **their cellular location**
 - x **“moonlighting proteins”** (Jeffery 2003)
 - x **Not one-gene-one-protein** (follow the changes of the organism in its **physiological & pathological conditions**)



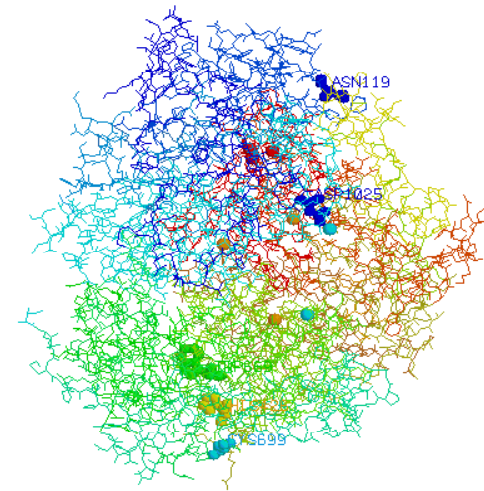
[Moonlighting proteins](#). Moonlighting proteins have multiple, seemingly **unrelated functions** not due to gene fusions or alternative splicing. Like **PGI**, which is a **cytosolic enzyme** and an **extracellular cytokine**, dozens of other proteins have been found to **moonlight**. Connie coined the term **moonlighting proteins** and has written several review articles that develop the **idea of moonlighting proteins** and describe **additional moonlighting proteins** from the **literature**, how **they switch between functions**, how they might have evolved, and how they might **benefit the cell**. She is currently writing two additional invited articles and planning **computational studies** of the **sequences** and **structures** of known moonlighting proteins.



Structures are more conserved than Sequences

Annotation Transfer (3)

- × **Automatic predictions, several important points**
 - × To define **a statistical threshold** of sequence similarity that permits the annotation transfer
 - × This **threshold** differs from **one biological function to another**, needs to be determined *ad hoc* (specifically) for
 - × **Each protein family &**
 - × **Each biological function** (Rost et al. 2003)



Motifs & Patterns (1)

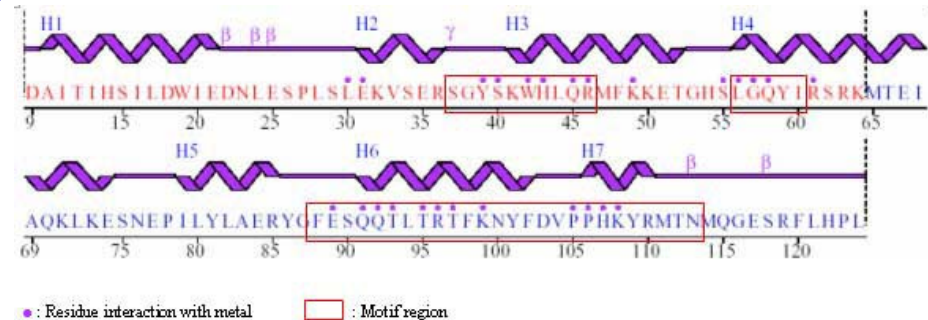
- x To develop various methods **to identify functionally important residues** based on their **conservation** throughout **evolution**
 - x Casari et al. 1995
- x Conserved **residues** & conserved **sequence elements** , likely to be **important functionally**
 - x But **no putative function** yet
- x Sometimes, there is **divergence** between the **sequence** of a **newly discovered** protein & any **other previously annotated** protein ⇒ **too wide** to establish relatedness

Levels of Protein Sequence and Structure Organization

Level/ Database	Content	Example
Primary	Sequence	"AVILDRYFH"
→ Secondary	Motif	[AS]-[IL]2-X[DE]-R- [FYW]2-H
Tertiary	Domain, module	a,b,c or @, *, #

Motifs & Patterns (2)

- × The existence of a **relatively short sequence motif** that is highly **conserved evolutionarily** & highly specific **functionally**
- × Helpful to reveal the **putative function**
 - × E.g., a sequence element that appears in **many known DNA-binding proteins**



(a)



(b)

× **Homeodomain-like proteins: 3 helices containing a helix-turn-helix (HTH) DNA binding motif (homeodomain in eukaryotes)**

Motifs & Patterns (3)

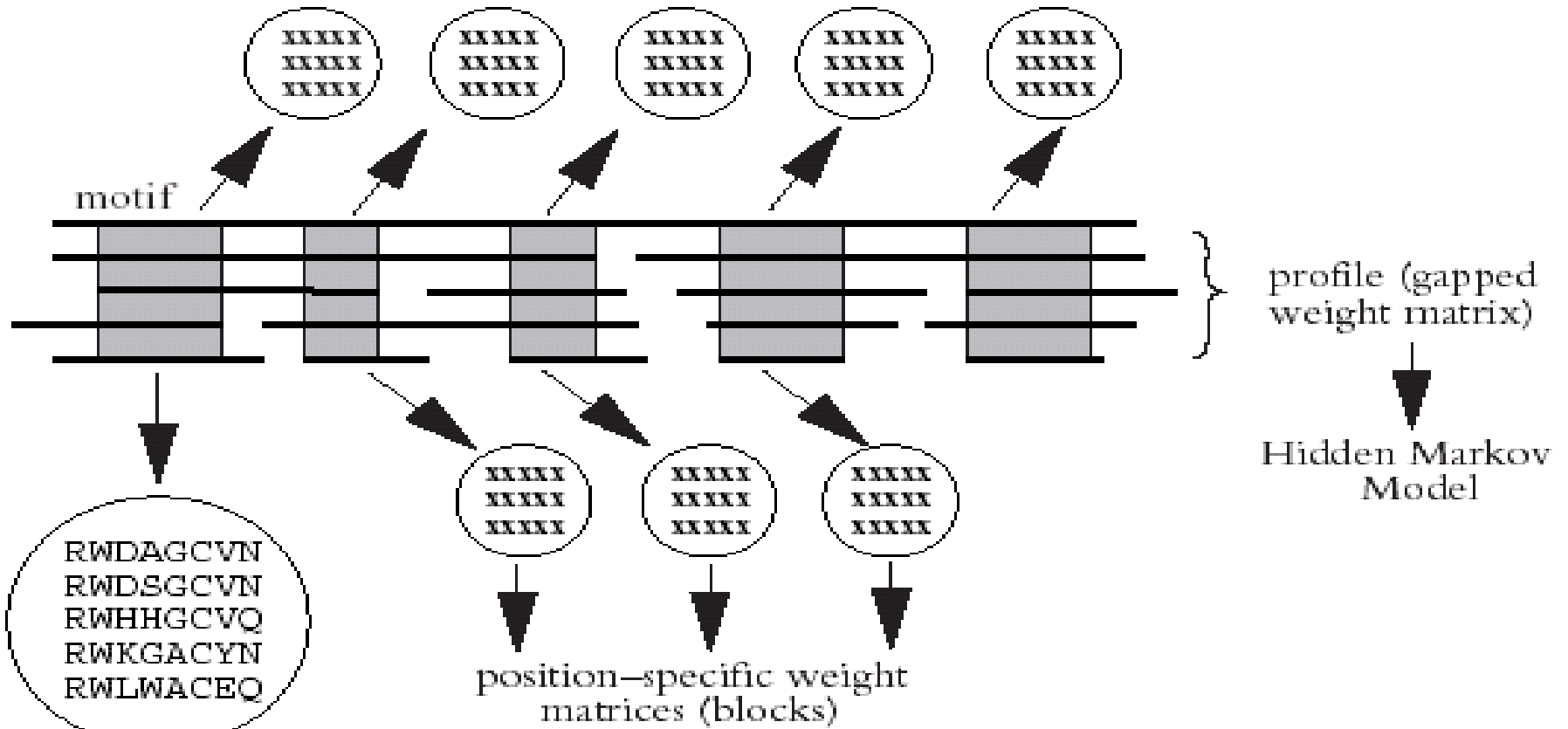
- × Several **databases** offer **large libraries** of these characterized sequence motifs that have been collected either
 - × **Manually** by experts or
 - × Automatically by **pattern-searching algorithms**
- × Many of these **libraries** include **a searching tool**
 - × To submit a sequence of interest to determine whether **an interesting motif** is contained within their protein
 - × Some important insight into **its structure & function**

Terminology - Profile

- × A **numerical** representation of a multiple sequence alignment (MSA)
 - × **Intrinsic sequence information** that represents the common characteristics of that particular collection of sequences, frequently a **protein family**
- × By using a profile, one is able to find **similarities** between sequences **with little or no absolute sequence identity** ⇒ to identify **distantly related proteins**

fingerprint

frequency (identity) matrices



profile (gapped weight matrix)

Hidden Markov Model

position-specific weight matrices (blocks)

RWDAGCVN
RWDSGCVN
RWHHGCVQ
RWKGACYN
RVLWACEQ

R-W-x(2)-[AG]-C-x-[NQ]
regular expression (pattern)

TK Attwood 2000 Briefings in Bioinformatics 1, 45-59

Terminology - Pattern/Signature

- x The common characteristics of a **protein family** or a multiple sequence alignment (MSA)
 - x Simply providing **a shorthand notation** for what residues can be present at any given position
- x E.g., the pattern **[IV]-G-x-G-T-[LIVMF]-x(2)-[GS]**
 - x The first position could contain either an isoleucine or a valine, the second position could contain only a glycine...
 - x {I}: residues **forbidden**
 - x An **x**: any residue can appear at that position
 - x **x(2)**: two positions can be occupied by any amino acid, the number reflecting simply the length of the nonspecific run

Methods - PROSITE (1)

- x Catalog of biologically **significant sites** through the use of **motif, sequence profiles, patterns** (Falquet et al. 2002; Hulo et al. 2004)

- x **Each entry in PROSITE**
 - x Linking to an "annotation **document**"
 - x Detailed information on the protein family that is characterized by **the profile or pattern**
 - x **Annotations**
 - x **Taxonomic range**
 - x **Domain architecture views**
 - x Known **biological** functions
 - x Any **experimentally determined three-dimensional** structures
 - x Relevance **references** from the literature

Methods - PROSITE (2)

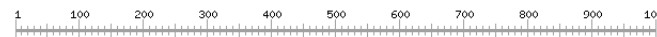
- x Each entry in PROSITE
 - x Information about how well the **pattern or profile** characterizes the **protein family** being described
 - x The number of **true & false positives**, as well as false negatives
 - x Which allow the users to determine the **specificity** of a given **profile or pattern**
- x A complete list **Swiss-Prot** sequences matching the profile or pattern

Cdc13p_yeast

USERSEQ1 (924 aa)

```
MDTLEEPECPPHKNRIFVSSSKDFEGYPSKAIVPVQFVALLTSIHLTETKCLLGF'SNFERRGDQSQ  
EDQYLKIKLKFDRGSERLARITISLLCQYFDIELPDLSDSGASPTVILRD IHLERLCF'SSCKALY  
VSKHGNYTLFLEDIKPLDLVSVISTISTRKSTNSKHSSELISECDLNNSLVDIFNNLIEMNRDEK  
NRFKFKLIHYDIELKKFVQDQKVLSSQSKAAAINPFFVFNRLGIPYIESQNEFNSQLMTLNVDE  
PTTDISNMGEEMHDSADPIEDSDSSTTSSTGKYFSSKSYIQSQTPERKTSVPPNNWHDDSDGSKRRR  
KLSFHSFNASSIRKAISYEQLSLASVGSVERLEGGIVGMNPPQFASINEFKYCTLKLKLYFTQLLPNV  
PDKVLVPGVNCIEIVIPTRERICELFGVLCQSDKISDILLLEKPDRISEVEVERILWDNDKTASPG  
MAVWSLKNISTDTQAQAQVQVPAQSSASIDPSRTRMSKMARKDPTIEFCQLGLDTEFKYITMFGM  
LVSCSFDKPAFISFVFSDFTKNDIVQNYLYDRYLIDYENKLELNEGFKAIMYKNQFETPDSKLRKI  
FNNGLRDLQNGRDENLSQYGVCKMNIKVKMYNGKLNIVRECEPVPHSQISSIASPSQCEHLRLE  
YQRAFKRIGESAISRYPFEYRRFFPIHRNGSHLAKLRFDEVKHEPKKSPPTPALAEHIPDLNADVS  
SFDVKFTDISLLDSARLPRPQQTHKSNLYSCEGRIIAIEYHASDLCFHTNELPLLQTRGLAP  
ERVLQLHIITSKNFAYFFNRSAYLQRQPLEEKYTLAQLGHSFKFNITSSLTLFPDITVALQIIN  
CPIECTFRELQQQLAHPKVAAPDSSGLDCAINATVNPRLRLAAQNGVTVKKEEDNDDDAGAVPTS
```

ruler:



hits by patterns: [1 hit (by 1 pattern) on 1 sequence]

Hits by **PS00141 ASP_PROTEASE** *Eukaryotic and viral aspartyl proteases active site* :

USERSEQ1



103 - 114:

LDSDSGASPTVI

Methods - Pfam (1)

- x A collection of **protein domain families** (Sonnhammer et al. 1997; Bateman et al. 2004)

- x **Each Pfam entry**
 - x **A multiple sequence alignment (MSA)** of a protein domain or conserved region

- x **Two Pfam databases**
 - x **Pfam-A**
 - x Manually

 - x **Pfam-B**
 - x Automatic clustering of **the ProDom database**

Methods - Pfam (2)

- x More than **80%** of all **Swiss-Prot & TrEMBL entries** are associated with a **Pfam entry**
- x **Profile HMMs** derived from the **multiple sequence alignments** associated with each entry can be used to deduce whether a **new sequence of interest** can be assigned to an already-characterized family
- x Particularly useful for when **similarity** can not be deduced through the use of simple **BLAST** or **FASTA** search
- x **Keyword & sequence searches**



Trusted matches - domains scoring higher than the gathering threshold (A)

Domain	Start	End	Bits	Evalue	Alignment	Mode
CDC13-DNA-bind	504	686	442.60	4.5e-130	Align	ls

Matches to Pfam-B

Domain	Start	End	Alignment
Pfam-B_137977	687	872	Align

Potential matches - Domains with Evalues above the cutoff

Domain	Start	End	Bits	Evalue	Alignment	Mode
DUF560	12	24	11.00	0.0078	Align	fs
UPF0227	41	69	4.80	0.36	Align	fs
gp32	325	335	5.00	0.87	Align	fs
ArsC	334	345	4.70	0.65	Align	fs
ArsC	678	685	0.30	12	Align	fs

[CDC13_YEAST](#)

Methods - InterPro (1)

- x An integrated resource for information about **protein families, domains and functional sites**, bring together information from a number of protein domain-based resources (APweiler et al. 2000)
 - x PROSITE
 - x PRINTS
 - x Pfam
 - x ProDom
- x To provide a **unified launching point** from which users can explore **protein families** further, moving from the *InterPro* entry to any of the sources from which *InterPro* entries are derived

Methods - InterPro (2)

- x Searches
 - x Text
 - x Sequence

- x Example (InterPro Scan)
 - x An [InterPro summary page](#) for CDC13_YEAST
 - x Click on IPR001969
 - x All InterPro pages give the name of **the family** links to
 - x **Signatures** that characterize members of the family
 - x Information on "**children**" of this protein family
 - x Used to indicate **protein family** & **subfamily** relationships

 - x Links to other databases
 - x **PDB, CATH, SCOP, BLOCKS and PROSITE**

Methods - InterPro (3)

- × **Radial taxonomy display**
 - × To provide a quick overview of **the taxonomic range** in which members of this protein family are observed
- × **Graphical view**
 - × The presence of **discrete motifs** in proteins belonging to this family
- × **References**
 - × More **in-depth** information about the domain

Methods - BLOCKS (1)

- x The concept of **blocks** to identify **a family proteins**, rather than relying on the individual sequences themselves (Petrovsky et al. 1996)
- x **Conserved motifs** from proteins **in the same family** are aligned **without introducing gaps** ⇒ **block**
 - x **Block = alignment** (not the individual sequence themselves)
- x An individual protein can **contain one or more blocks**, corresponding to each of its **functional & structural motifs**
- x The BLOCKS database is derived from the entries in ***InterPro***

Methods - BLOCKS (2)

- x **BLOCKS search**
 - x The query sequence is **aligned against all blocks in the database at all possible positions**
 - x For each alignment, **a score** is calculated using a **position specific scoring matrix (PSSM)**
 - x Searches can be performed **optionally** against **the PRINTS database**, including information on more than **300 families** that do not have corresponding entries in BLOCKS
 - x **Better way is to search both**
- x **BLOCK Maker** (the same used to construct the database)
 - x To submit **a set (≥ 3) of unaligned**, related protein sequences to **detect conserved blocks** within the sequence set

Evaluation of Performance

- x To **analyze**, not predict sequences
 - x **Impossible to assess**
 - x Accuracy & predictive power are **irrelevant**
 - x **All complementary to each other**, the output could be **integrated** more straightforwardly

- x **Pfam vs. PROSITE**
 - x **Pfam**: find **a few long motifs** that can help associate the query sequence with a **well-characterized family** & predict its function

 - x **PROSITE**: more short ones

Subcellular Localization

- x Annotation transfer yield no results...
- x **History**
 - x 2nd structure or topology
 - x Structural families or folds
 - x **Subcellular localization**
 - x To narrow down its putative function

Subcellular Localization - Methods (1)

- x PSORT
 - x Nakai & Horton 1999
- x Searches for **features** characterizing that **taxonomic group** that may help deduce the subcellular localization of the protein
 - x A library of **known signal peptides** & searches for them in the query sequence
 - x Also checks **predicted structural features** (e.g., **topology**) that may indicate whether the protein is **soluble** or **embedded in membrane**

Subcellular Localization - Methods (2)

- x *SUBLOC*
 - x Hua & Sun 2001
- x Using **the amino acid composition alone**, applies support vector machine (**SVM**) to predict the subcellular locale of a protein
- x **Prokaryotes**
 - x Extracellular, periplasmic, or cytoplasmic
- x **Eukaryotes**
 - x Extracellular, mitochondrial, cytoplasmic or nuclear

Subcellular Localization - Methods (3)

- x TargetP
 - x Emanuelsson et al. 2000
- x Signal peptides at **the N-terminal end** of a protein
- x A series of machine-learning algorithms, including **neural networks & SVMs**, to identify **signal peptides** of three types
 - x **Chloroplast** transit peptides
 - x **Mitochondrial** targeting peptide
 - x **Secretory** pathway signal peptide

Subcellular Localization - Methods (4)

- x LOC3D
 - x Nair & Rost 2003
- x A database of predicted subcellular localizations for eukaryotic proteins of **known three-dimensional structure (PDB)**, three underlying methods
 - x **PredictNLS**
 - x Searching for a known **nuclear localization signal**
 - x **LOChom**
 - x Using homology to determine localization
 - x **LOC3D**
 - x A neural network-based prediction method

Subcellular Localization - LOCkey

- × A related service of *LOC3D*
- × Using **keywords** in **Swiss-Prot annotation** to predict the subcellular localization
- × A **comprehensive** coverage of the methods & approaches available for the prediction of subcellular localization

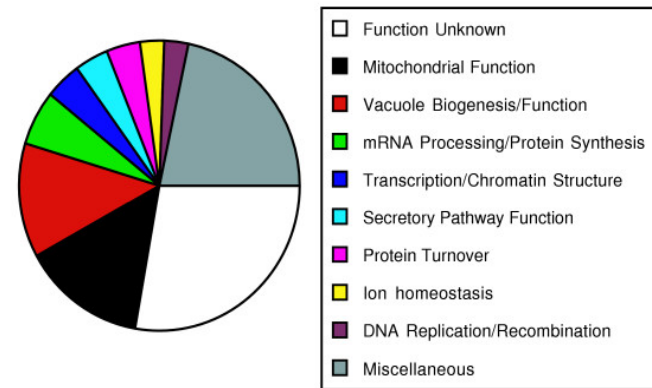
Evaluation of Performance

- x **No** continuous, large-scale system that assesses and compares the performance of the different available methods
- x Using a **large test set** & applying the **same statistical analysis** to all the methods
 - x Nair & Rost (2003) found **a large variability in the accuracy** of the methods
 - x Best: extracellular >80%; cytoplasmic proteins > 50%; nuclear > 70%; mitochondrial > 60%

Functional Class (1)

- x The prediction of **subcellular localization** is part of the process of reducing the ambiguous term *function* to one that is better define
- x The attempt to break down the notion of function into a series of **well-defined categories**
 - x It is facilitated by **the definition** of **a set of functional classes** to which a protein can be assigned
 - x **Software** were developed with the goal of assigning proteins to each one of the functional classes

Functional Class (2)



x Example

- x A scheme of functional classes to annotate E. coli (Riley 1993; next slide)
- x **Some genome projects** adopted this general scheme, eventually making this the most widely used **terminology for the assignment of functional class**

ENERGY	INFORMATION	COMMUNICATION
METABOLISM	DNA & RNA	SIGNAL
Amino Acid Biosynthesis.	Replication.	Regulatory Functions.
Biosynthesis of Cofactors, Prosthetic grps & Carriers.	Transcription	Cell Division.
Central Intermediary Metabolism.	TRANSLATION	Cell Killing.
Energy Metabolism.	Translation.	TRANSPORT
Fatty Acid and Phospholipid Biosynthesis.	PROTEINS	Transport and Binding Proteins.
Purines, Pyrimidines, Nucleosides, Nucleotides.	Chaperones.	ENVIRONMENT
	Detoxification.	Adaptation to Atypical Conditions and Others.
	Protein & Peptide Secretion and Transformation.	

Equivalence between the classification of *E. coli* proteins (*Riley, 1993*) and the three- and eight- categories classification (STRUCTURAL proteins are not included in this schema)

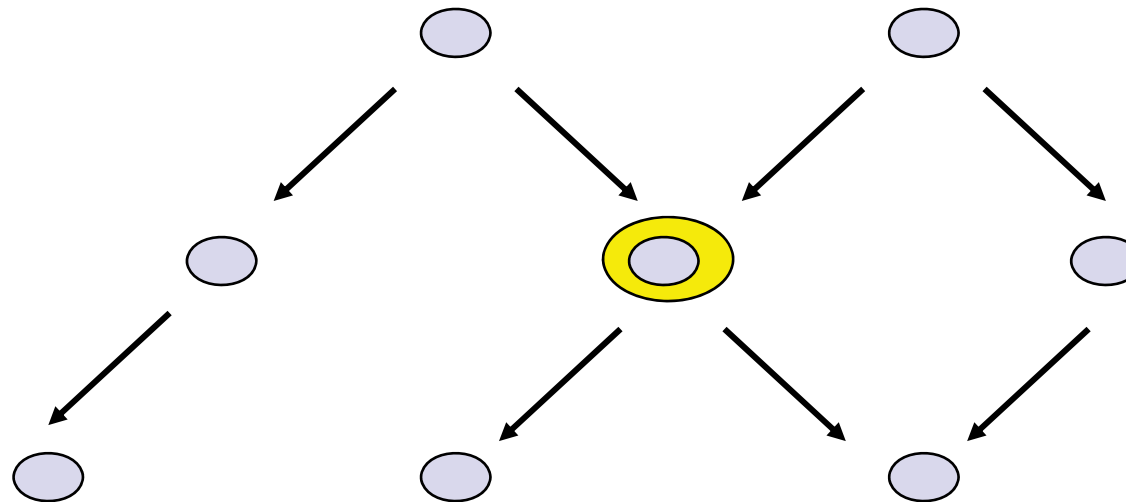
Prediction Methods - EUCLID

- x Tamames et al. 1998
 - x A basic **machine-learning algorithm** that learns, based on a **manually curated training set**, which combination of keywords is mostly likely to indicate that the protein belong to a certain functional type
 - x **Keywords** in **Swiss-Prot** to assign a protein to one of Riley's functional class
 - x >90% of the cases, the functional type that was determined by the automated method was identical to the one that was assigned to it **by a human expert**
 - x *EUCLID* requires that some annotation (**Swiss-Prot keywords**) already be assigned to the sequence
 - x Not really a method for prediction strictly from sequences

Prediction Methods - ProtFun (1)

- x Jensen et al. 2003
- x A recent & promising step toward **the prediction of function from sequence**
- x **Gene Ontology** (GO; Ashburner et al. 2000) in attempt to build a controlled vocabulary to describe genes & gene products
 - x Each protein can be assigned to a certain molecular function
 - x **Molecular function**
 - x **Biological process**
 - x **Cellular component**

DAG Structure

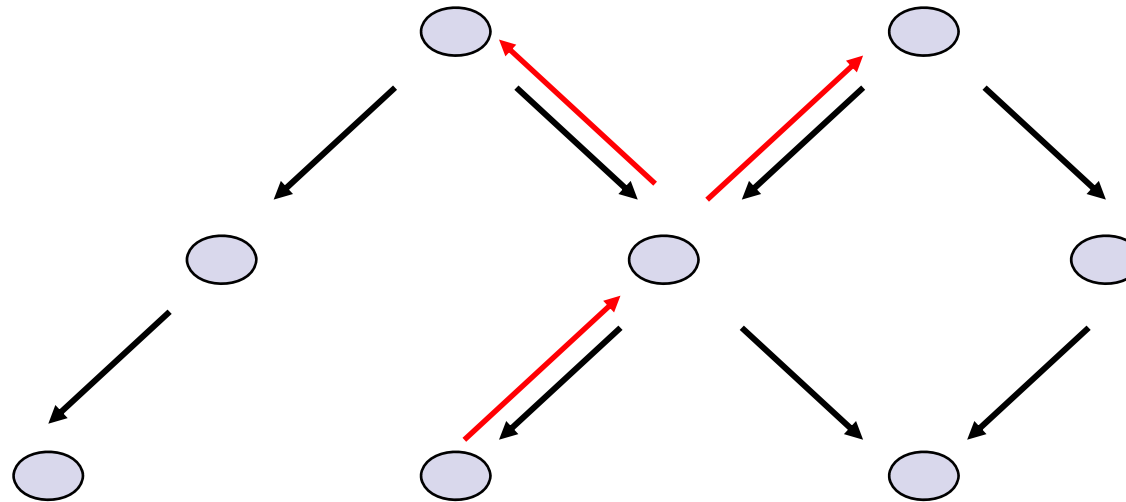


- * Directed acyclic (環式) graph
 - * Each **child** may have **one or more parents**

Relationship Types

- × **Is-a**
 - × Subclass; a **is** one type of b
- × **Part-of**
 - × Physical part of (component)
 - × Sub-process of (process)

The True Path Rule





- × Every **path** from a **node back** to the **root** must be **biologically accurate**

the Gene Ontology website - Microsoft Internet Explorer 05-17-2005 05:01:56

檔案(F) 編輯(E) 檢視(V) 我的最愛(A) 工具(T) 說明(H)

← 上一頁 → 搜尋 我的最愛 媒體

網址(D) http://www.geneontology.org/ 連結 >>

Open menus
[Site Map](#)
[Home](#) [New](#) [FAQ](#)

[Downloads](#)
[Ontologies](#)
[Annotations](#)
[Database](#)
[Mappings to GO](#)
[GO Tools](#)

[Documentation](#)

[About GO](#)
[Terms Of Use](#)
[Contact GO](#)
[Report Errors](#)

[Archives](#)

Search GO

 GO terms
 gene or protein name

Search Site

<http://www.geneontology.org/> **GENE ONTOLOGY**

[What is the Gene Ontology?](#)[Download the Ontologies](#)

The goal of the Gene Ontology project is to produce a controlled vocabulary that can be applied to all organisms even as knowledge of gene and protein roles in cells is accumulating and changing. GO provides three structured networks of defined terms to describe gene product attributes. GO is one of the controlled vocabularies of the [Open Biomedical Ontologies](#).

- Submit new GO term suggestions via the [Curator Requests Tracker](#) at [SourceForge](#). [Help with new term submission](#) is available.
- Send comments and questions to go@geneontology.org.

Search GO

GO term gene or protein name

This search uses the browser [AmiGO](#). See the [GO tools section](#) for other GO browsers.

What's New?

- We are pleased to announce the [2005 GO Users Meeting](#) will be held as part of the [MGED 8 meeting](#) in Bergen, Norway, in mid September. See the [meeting information page](#) for more on the meeting and to submit your abstract. *(posted May 2 2005)*

網際網路

GO Evidence Code

From reviews or introductions

- × IDA
 - × Inferred from Direct Assay
- × IMP
 - × Inferred from Mutant Phenotype
- × IGI
 - × Inferred from Genetic Interaction
- × IPI
 - × Inferred from Physical Interaction
- × IEP
 - × Inferred from Expression Pattern

- × TAS
 - × Traceable Author Statement
- × NAS
 - × Non-traceable Author Statement

- × IC
 - × Inferred by Curator
- × ISS
 - × Inferred from **Sequence** or **Structural Similarity**

- × IEA
 - × Inferred from Electronic Annotation
- × ND
 - × Not Determined

From Primary Literature

Automated

Prediction Methods - ProtFun (2)

- × Currently, there are hundreds of GO categories
- × *ProtFun* focuses on 347 of them & uses **complex systems of neural networks** to predict the **GO functional classification** of a protein from its sequence
 - × Impressive accuracy >90%
- × **Current coverage is only partial** & many of the query proteins are returned without any prediction at all

```
>Sequence
# Functional category      Prob      Odds
Amino_acid_biosynthesis   0.158     7.190
Biosynthesis_of_cofactors 0.087     1.203
Cell_envelope             0.065     1.061
Cellular_processes        0.027     0.373
Central_intermediary_metabolism => 0.241     3.830
Energy_metabolism         0.023     0.259
Fatty_acid_metabolism     0.019     1.436
Purines_and_pyrimidines   0.431     1.775
Regulatory_functions      0.334     2.075
Replication_and_transcription 0.352     1.313
Translation               0.063     1.425
Transport_and_binding     0.165     0.402

# Enzyme/nonenzyme        Prob      Odds
Enzyme                    => 0.439     1.534
Nonenzyme                 0.561     0.786

# Enzyme class            Prob      Odds
Oxidoreductase (EC 1.--.-) 0.024     0.114
Transferase (EC 2.--.-)    0.195     0.564
Hydrolase (EC 3.--.-)     0.160     0.505
Lyase (EC 4.--.-)         0.020     0.430
Isomerase (EC 5.--.-)     0.012     0.366
Ligase (EC 6.--.-)       => 0.118     2.318

# Gene Ontology category  Prob      Odds
Signal_transducer         0.063     0.293
Receptor                  0.006     0.033
```

CDC13_YEAST