

## Getting To Know Your Protein

### Comparative Protein Analysis: Part II. Protein Domain Identification & Classification

Robert Latek, PhD  
Sr. Bioinformatics Scientist  
Whitehead Institute for Biomedical Research

## Comparative Protein Analysis

- **Part I. :**
  - **Phylogenetic Trees and Multiple Sequence Alignments** are important tools to understand global relationships between sequences.
  - **Tree Building Tools with Different Algorithms**
    - <http://bioweb.pasteur.fr/seqanal/phylogeny/intro-uk.html>
    - <http://evolution.genetics.washington.edu/phylip/software.xref.html>
  - **Tree Reliability**
    - Bootstrapping 1. Randomly re-sample MSA columns to produce a random alignment (equal length as original MSA), 2. Build tree based on random alignment, 3. Predicted branches are significant if they occur in ~ >70% of the trees from multiple, randomized alignments.
    - Use a several tree building algorithms to determine whether they produce similar trees as the original.

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

2

## Comparative Protein Analysis

- **Part II. :**
  - How do you identify sequence relationships that are restricted to localized regions?
  - Can you apply phylogenetic trees and MSAs to only sub-regions of sequences?
  - How do you apply what you know about a group of sequences to finding additional, related sequences?
  - What can the relationship between your sequences and previously discovered ones tell you about their function?
- Assigning sequences to **Protein Families**

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

3

## Syllabus

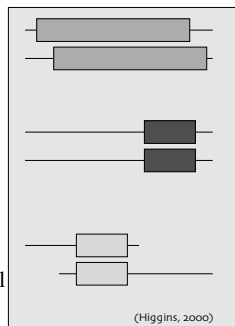
- **Protein Families**
  - Identifying Protein Domains
  - Family Databases & Searches
- **Searching for Family Members**
  - Pattern Searches
    - Patscan
  - Profile Searches
    - PSI-BLAST/HMMER2

WIBR Bioinformatics Courses, © Whitehead Institute, 2005

4

## Proteins As Modules

- Proteins are derived from a limited number of basic building blocks (**Domains**)
- Evolution has shuffled these modules giving rise to a diverse repertoire of protein sequences
- As a result, proteins can share a global or local relationship



WIBR Bioinformatics Courses, © Whitehead Institute, 2005

5

## Protein Domains

SH2 Motif

```

BLK_MOUSE 117-138  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
LCK_MOUSE 126-208  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
LYN_MOUSE 128-210  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
FGR_HUMAN 144-226  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
SRC_RBDP 144-230  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
NFKB_HUMAN 282-376  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
VAV_MOUSE 671-745  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
BCL2_HUMAN 573-648  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
P53A_HUMAN 624-698  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
SRC_HUMAN 688-858  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
ITK_HUMAN 235-323  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
SYK_HUMAN 381-562  WFFPTTIRBAPGCLLAPWKADEPFLIRREINRQAPLVEVETLTVQV...VVEVVEIIEIDW...DEVEEPEET...PPELQALVQRY
  
```

Janus Kinase 2 Modular Sequence Architecture

JAK2 Ubiquitin-Like Domain Acyl-CoEnzyme A Binding Domain PH Domain SH2 Kinase

Motifs describe the domain

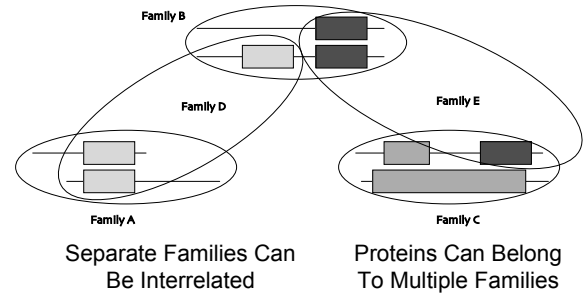
WIBR Bioinformatics Courses, © Whitehead Institute, 2005

6

# Protein Families

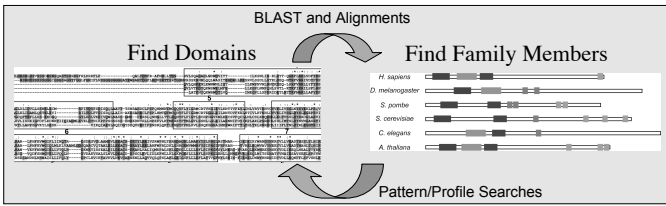
- **Protein Family** - a group of proteins that share a common function and/or structure, that are potentially derived from a common ancestor (set of homologous proteins)
- **Characterizing a Family** - Compare the sequence and structure patterns of the family members to reveal shared characteristics that potentially describe common biological properties
- **Motif/Domain** - sequence and/or structure patterns common to protein family members (trait/feature/characteristic)

# Protein Families



# Creating Protein Families

- Use domains to identify family members
  - Use a sequence to search a database and characterize a pattern/profile
  - Use a specific pattern/profile to identify homologous sequences (family members)



# Family Database Resources

- **Curated Databases\***
  - Proteins are placed into families with which they share a specific sequence pattern
- **Clustering Databases\***
  - Sequence similarity-based without the prior knowledge of specific patterns
- **Derived Databases\***
  - Pool other databases into one central resource
- **Search and Browse**
  - **InterPro** <http://www.ebi.ac.uk/interpro/> \*(Higgins, 2000)

# Curated Family Databases

- **Pfam** (<http://pfam.wustl.edu>) \*\*
  - Uses manually constructed seed alignments and PSSM to automatically extract domains
  - db of protein families and corresponding profile-HMMs of prototypic domains
  - Searches report e-value and bits score
- **Prosite** (<http://www.expasy.ch/tools/scanprosite/>)
  - Hit or Miss -> no stats
- **PRINTS** (<http://www.bioinf.man.ac.uk/fingerPRINTScan/>)

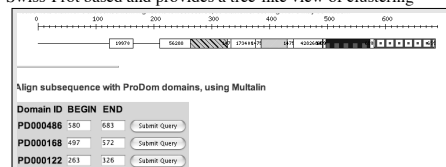
Pfam HMM search results, glocal+local alignments merged (Pfam\_ls+Pfam\_fs)

[Go here for an explanation of the format of the results]

Model	Seq-from	Seq-to	HMM-from	HMM-to	Score	E-value	Alignment	Description
GTP_EFTU	258	483	1	298	315.7	5.5e-92	glocal	Elongation factor Tu GTP binding domain
GTP_EFTU_D2	502	570	1	75	46.1	8e-11	glocal	Elongation factor Tu domain 2
GTP_EFTU_D3	576	684	1	112	142.9	6.1e-40	glocal	Elongation factor Tu C-terminal domain

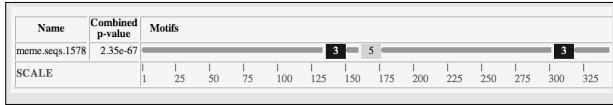
# Clustering Family Databases

- Search a database against itself and cluster similar sequences into families
- **ProDom** (<http://prodes.toulouse.inra.fr/prodom/current/html/home.php>)
  - Searchable against MSAs and consensus sequences
- **Protomap** (<http://protomap.cornell.edu/>)
  - Swiss-Prot based and provides a tree-like view of clustering



# Derived Family Databases

- Databases that utilize protein family groupings provided by other resources
- Blocks** - Search and Make (<http://blocks.fhere.org/blocks/>)
  - Uses Protomap system for finding blocks that are indicative of a protein family (GIBBS/MOTIF)
- Proclass** (<http://pir.georgetown.edu/gfserver/proclass.html>)
  - Combines families from several resources using a neural network-based system (relationships)
- MEME** (<http://meme.sdsc.edu/meme/website/intro.html>)



# Syllabus

- Protein Families**
  - Identifying Protein Domains
  - Family Databases & Searches
- Searching for Family Members**
  - Pattern Searches
    - Patscan
  - Profile Searches
    - PSI-BLAST/HMMER2

# Searching Databases By Family

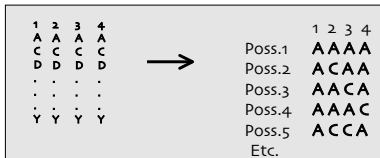
- BLAST searches provide a great deal of information, but it is difficult to select out the important sequences (listed by score, not family)
- Family searches can give an immediate indication of a protein's classification/function
- Use Family Database search tools to identify domains and family members

# Patterns & Profiles

- Techniques for searching sequence databases to uncover common domains/motifs of biological significance that categorize a protein into a family
- Pattern** - a deterministic syntax that describes multiple combinations of possible residues within a protein string
- Profile** - probabilistic generalizations that assign to every segment position, a probability that each of the 20 aa will occur

# Pattern Discovery Algorithms

- Pattern Driven Methods**
  - Enumerate all possible patterns in solution space and try matching them to a set of sequences



# Pattern Discovery Algorithms

- Sequence Driven Methods**
  - Build up a pattern by pair-wise comparisons of input sequences, storing positions in common, removing positions that are different



# Pattern Building

- Find patterns like “pos1 xx pos2 xxxx pos3”
  - Definition of a non-contiguous motif

```

1. C Y D - - C A F T L R Q S A V M H K H A R E H
2. C A T Y - C R T A I D T V K N S L K H H S A H
3. C W D G G C G I S V E R M D T V H K H D T V H
4. C Y C - - C S D H M K K D A V E R M H K K D H
5. C N M F - C M P I F R Q N S L A R E H E R M H
6. C L N N T C T A F W R Q K K D D T V H N S L H

C xxxx C xxxx [LIVMFYW] xxxxxxxx H xxxxx H
    
```

Define/Search A Motif <http://us.expasy.org/tools/scanprosite/>

# Pattern Properties

- Specification**
  - a single residue K, set of residues (KPR), exclusion {KPR}, wildcards X, varying lengths x(3,6) -> variable gap lengths
- General Syntax**
  - C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H
- Patscan Syntax**
  - <http://jura.wi.mit.edu/bio/education/bioinfo/homework/hw8/patscan.txt>
  - C 2...4 C 3...3 any(LIVMFYWC) 8...8 H 3...5 H
- Pattern Database Searching**
  - %scan\_for\_matches -p pattern\_file < nr > output\_file

# Sequence Pattern Concerns

- Pattern descriptors must allow for approximate matching by defining an acceptable distance between a pattern and a potential hit
  - Weigh the sensitivity and specificity of a pattern
- What is the likelihood that a pattern would randomly occur?

# Sequence Profiles

- Consensus** - mathematical probability that a particular aa will be located at a given position
- Probabilistic** pattern constructed from a MSA
- Opportunity to assign penalties for insertions and deletions, but not well suited for variable gap lengths
- PSSM** - (Position Specific Scoring Matrix)
  - Represents the sequence profile in tabular form
  - Columns of weights for every aa corresponding to each column of a MSA

# Profile Discovery/Analysis

- Perform global MSA on group of sequences
- Move highly conserved regions to smaller MSAs
- Generate scoring table with log odds scores
  - Each column is independent
  - Average Method: profile matrix values are weighted by the proportion of each amino acid in each column of MSA
  - Evolutionary Method: calculate the evolutionary distance (Dayhoff model) required to generate the observed amino acid distribution



# PSSM Example

( i.e. Distribution of aa in an MSA column)

← Target sequences      Resulting Consensus: I T L S

PSSM

P	O	S	A	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y
1	8	-2	5	4	5	5	-4	24	0	15	13	1	1	1	-7	2	22	21	-18	-6		
2	13	-5	24	18	-18	19	7	1	7	-7	-4	14	11	10	-1	9	22	3	-28	-14		
3	5	-5	3	4	13	4	2	8	-4	14	12	8	-5	0	-10	0	10	10	-1	5		
4	17	17	13	10	-12	29	-5	-5	6	-14	-9	12	10	0	-2	34	19	1	-8	-15		

## PSSM Properties

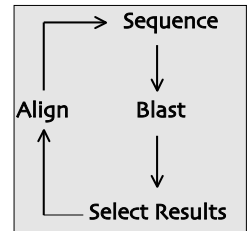
- Score-based sequence representations for searching databases
- Goal
  - Limit the diversity in each column to improve reliability
- Problems
  - Differing length gaps between conserved positions (unlike patterns)

## PSI-BLAST Implementation

### PSI-BLAST

<http://www.ncbi.nlm.nih.gov/BLAST/>

- Start with a sequence, BLAST it, align select results to query sequence, estimate a profile with the MSA, search DB with the profile - constructs PSSM
- Iterate until process stabilizes
- Focus on domains, not entire sequences
- Greatly improves sensitivity (but may affect specificity)



## PSI-BLAST Sample Output

```

Sequences with E-value WORSE than threshold

gi196290991.ref|NP_044074.1| (NC_001771) KC12R [Molluscum contag... 32 0.16
gi181755541.db|AA035409.2| (S79774) bile salt-dependent lipase: B... 34 0.25
gi145027711.ref|NP_001798.1| (NM_001807) carboxyl ester lipase (b... 35 0.86
gi12316291.sp|P19035|EAL_HUMAN Bile-salt-activated lipase precurs... 35 0.89
gi1152429291.ref|NP_200612.1| (NM_125189) putative protein (Arabi... 34 1.1
gi197695291.db|I8AB10595.1| (AB024029) gene_id:K21L19.3-unknown p... 34 1.3
gi11804821.sp|AA052014.1| (M85201) cholesterol esterase [Homo sap... 33 1.8
gi1187061.sp|P211731|DNAA_M1CLU Chromosomal replication initiator... 32 4.6
gi1126761.sp|P161101|EG3_MOUSE GALECTIN-3 (GALACTOSE-SPECIFIC LE... 32 4.9
gi1528511emb|CA034206.1| (X16074) L-34 protein (AA 1-264) [Mus sp.] 32 5.0
gi15399071.p1|L1A5903| lactose-binding lectin Mac-2 - mouse 32 5.0
gi11871111.db|AA037311.1| (J03723) carbohydrate binding protein 3... 32 5.4
gi195064271.ref|NP_065019.1| (NM_019146) bassoon [Rattus norvegic... 32 5.5
  
```

## HMM Building

- **Hidden Markov Models** are Statistical methods that consider all the possible combinations of matches, mismatches, and gaps to generate a consensus (Higgins, 2000)
- Sequence ordering and alignments are not necessary at the onset (but in many cases alignments are recommended)
- Ideally use at least 20 sequences in the training set to build a model
- Calibration prevents over-fitting training set (i.e. Ala scan)
- Generate a model (profile/PSSM), then search a database with it

## HMM Implementation

- **HMMER2** (<http://hmmer.wustl.edu/>)
  - Determine which sequences to include/exclude
  - Perform alignment, select domain, excise ends, manually refine MSA (pre-aligned sequences better)
  - Build profile
    - `%hmmbuild [-options] <hmmfile output> <alignment file>`
  - Calibrate profile (re-calc. Parameters by making a random db)
    - `%hmmcalibrate [-options] <hmmfile>`
  - Search database
    - `%hmmsearch [-options] <hmmfile> <database file> > out`

## HMMER2 Output

- Hmmssearch returns e-values and bits scores
- Repeat process with selected results
  - Unfortunately need to extract sequences from the results and manually perform MSA before beginning next round of iteration

```

HMMER 2.2g (August 2001)
Copyright (C) 1992-2001 HHMI/Washington University School of Medicine
Freely distributed under the GNU General Public License (GPL).

HMM file: pfam_h3d.hmm [Hydrolase]
Sequence database: /cluster/d0/Data/nr
per-sequence score cutoff: [none]
per-domain score cutoff: [none]
per-sequence Eval cutoff: <= 10
per-domain Eval cutoff: [none]

-----
Query HMM: Hydrolase
Accession: PF00702
Description: halosialid algalactonase-like hydrolase
[HMM has been calibrated: E-values are empirical estimates]

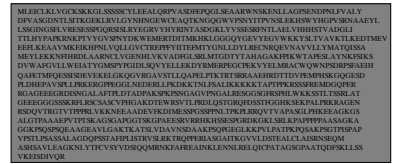
Scores for complete sequences (score includes all domains):
Sequence Description Score E-value N
-----
gi16131263.ud|NP_417844.1| phosphoglycolat 168.4 2.9e-45 1
gi24114648.ud|NP_709158.1| phosphoglycolat 167.8 4.2e-45 1
gi15803888.ud|NP_289924.1| phosphoglycolat 167.8 4.2e-45 1
gi26249979.ud|NP_756019.1| Phosphoglycolat 166.4 1.1e-44 1
  
```

# Patterns vs. Profiles

- **Patterns**
  - Easy to understand (human-readable)
  - Account for different length gaps
- **Profiles**
  - Sensitivity, better signal to noise ratio
  - Teachable

# Domain ID & Searching

- Family/Domain Search
  - <http://pfam.wustl.edu>
- Pattern Search
  - scan\_for\_matches (Patscan)
    - scan\_for\_matches -p pattern\_file </cluster/db0/Data/yeast.aa > output\_file
- Profile Search
  - HMMER2
    - hmmbuild [-options] <hmmfile output> <alignment file>



```
any(CF) any(HE) any(GK) 1...1 any(LL) 4...4 any(AS) 3...3 any(LJ) 3...3 any(GA) 3...3 G 1...1 any(YF) 1...1 any(LJ) R
```

# Exercises

- Use PFAM to identify domains within your sequence
- Scan your sequences with ProSite to find a pattern to represent the domain
- Use the ProSite pattern to search the non-redundant db
- Use PSI-BLAST to build a sequence profile and search the non-redundant db

# References

- Bioinformatics: A Practical Guide to the Analysis of Genes and Proteins. Andreas D. Baxevanis and B.F. Francis Ouellete. Wiley Interscience, 2001.
- Bioinformatics: Sequence, structure, and databanks. Des Higgins and Willie Taylor. Oxford University Press, 2000.