

Current Topics in Genome Analysis
Fall 2003

Week 5
Biological Sequence Analysis II

Andy Baxevanis, Ph.D.



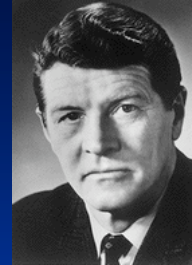
Overview

- Week 4: Comparative methods and concepts
 - Similarity vs. Homology
 - Global vs. Local Alignments
 - Dotplots
 - Scoring Matrices
 - BLAST
- Week 5: Predictive methods and concepts
 - Profiles, patterns, motifs, and domains
 - Secondary structure prediction
 - Structures: VAST, Cn3D, and *de novo* prediction

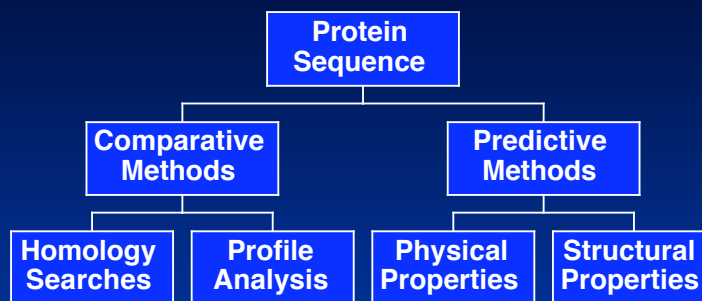


Protein Conformation

- Christian Anfinsen
Studies on reversible denaturation →
“Sequence specifies conformation”
- Chaperones and disulfide interchange enzymes:
involved but not controlling final state
- “Starting with a newly-determined sequence,
what can be determined computationally about
its possible function and structure?”



Protein Sequence Analysis



- *Common structure?*
- *Common function?*
- *Evolutionary relationship?*
- *Global or local similarity?*



Sequence Comparisons

- Homology searches
 - Usually “one-against-one” BLAST
FASTA
 - Allows for comparison of individual sequences against databases comprised of individual sequences
- Profile searches
 - Uses collective characteristics of a family of proteins
 - Search can be “one-against-many” ProfileScan
CDD
 - or “many-against-one” PSI-BLAST



Profiles

- Numerical representations of multiple sequence alignments
- Depend upon *patterns* or *motifs* containing conserved residues
- Represent the common characteristics of a protein family
- Can find similarities between sequences with little or no sequence identity
- Allow for the analysis of distantly-related proteins



Profile Construction

APHIIVATPG
 GCEIVIAATPG
 GVEICIAATPG
 GVDILIGATPG
 RPHIIVATPG
 KPHEIIAATPG
 KVQLIIATPG
 RPDIVIAATPG
 APHIIIVATPG
 APHIIIVATPG
 GCHVVIAATPG
 NQDIVVATPG

- Which residues are seen at each position?
- What is the frequency of observed residues?
- Which positions are conserved?
- Where can gaps be introduced?

Position-Specific Scoring Table

Cons	A	B	C	D	E	F	G	H	I	K	L	M	N	P	Q	R	S	T	V	W	Y	Z
G	17	18	0	19	14	-22	31	0	-9	12	-15	-5	15	10	9	6	18	14	1	-15	-22	11
P	10	0	10	0	0	12	10	0	0	0	0	0	0	23	2	-2	12	11	17	-31	-8	1
H	5	24	-12	29	25	-20	8	32	-9	9	-10	-9	22	7	30	10	0	4	-8	-20	-7	27
I	-1	-12	6	-13	-11	33	-12	-13	63	-11	40	29	-15	-9	-14	-15	-6	7	50	-17	8	-11
V	3	-11	1	-11	-9	22	-3	-11	46	-9	37	30	-13	-3	-9	-13	-6	6	50	-19	2	-8
V	5	-9	9	-9	-9	19	-1	-13	57	-9	35	26	-13	-2	-11	-13	-4	9	58	-29	0	-9
A	54	15	12	20	17	-24	44	-6	-4	-1	-11	-5	12	19	9	-13	21	19	9	-39	-20	10
T	40	20	20	20	20	-30	40	-10	20	20	-10	0	20	30	-10	-10	30	150	20	-60	-30	10
P	01	0	7	0	0	0	10	0	0	0	0	0	0	89	17	17	24	22	9	-50	-48	12
G	70	60	20	70	50	0	150	-20	-30	-10	-50	-30	40	30	20	-30	60	40	20	-100	-70	30



Patterns

[FY]-x-C-x(2)-[VA]-x-H(3)

reads as:

Phe or Tyr
 any amino acid
 Cys
 any two amino acids
 Val or Ala
 any amino acid
 three His

followed by
 followed by
 followed by
 followed by
 followed by



ProfileScan

- Search sequence against a collection of profiles and patterns
- Databases available
 - PROSITE profiles
 - PROSITE patterns
 - PfamA
 - PfamB
- <http://hits.isb-sib.ch/cgi-bin/PFSCAN>



Motif Scan in a Protein Sequence – Netscape

http://hits.isb-sib.ch/cgi-bin/PFSCAN_parser

hits Motif Scan in a Protein Sequence

- Databases: Prosite patterns (Hits-synchronized), Prosite profiles (Hits-synchronized), Pfam collection of hidden Markov models (Hits-synchronized)

...scanning for hits_pattern
 ...scanning for hits_profile
 ...scanning for hits_pfam
 ...confirming pfam

Result

- Summary:
 - pat:CYTOCHROME_P450 pos. 449 - 458
 - ! pfam:P450 pos. 41 - 506 E-value=4.2e-138
- Match Location:

query	MAFSQYISLAPPELLATAIFCLVFWLGRTRTQVPGKLSPPQ
pfam:P450	STPVVVL SGLNTIKQALVKQGD DFKGRPDLYSFTLITNGKSMTFNPDSPVWAARRRLAQDALKSPSIA SDPTSVSSCYL
query	EEHVSKEANHLISKFKLMAEVGHFEPVNVQVSVANVIGAMCPGKNFRKSEEMLNLVKSSKDFVENVTSGNAVDFFPV
pfam:P450	LRYLNPALKRFKFNDFVLSLQKTQVQEHYQDFNKNSIQDITGALFKHSENYKDNGLIPQEKIVNIVNDIFGAGFTV
query	TTAIFWSLLLVTPEKVQRKIHEELDTVIGRDRQPLSRDLPYLEAFLEIYRVTSFVFPFTIPHSSTRDTSLNQFHP
pfam:P450	KECCIFINQWVNHDEKQWKDFVFRPERFLTNDNTAIDKTLSEKVMFLGKRRRCIGEIPAKWEVFLFLALLHQLLEFT
query	pat:CYTOCHROME_P450
pfam:P450	VPFGVKVDLTPSYGLTMKPRTCHEVQAWPRFSK

status: !
 pos.: 41-506
 raw-score = 472.2
 N-score = 144.703
 E-value = 4.2e-138

pfam: P450
 Cytochrome P450
 PF00067
 Pfam-site
 InterPro.

Pfam: p450 – Netscape

http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00067

Pfam: p450

Wellcome Trust
 Sanger Institute

Home Keyword Search Protein Search Browse Pfam DNA Search Taxonomy ftp Help p450 domain

p450

Accession number: PF00067

Cytochrome P450 [Add Annotation](#)

Cytochrome P450s are involved in the oxidative degradation of various compounds. Particularly well known for their role in the degradation of environmental toxins and mutagens. Structure is mostly alpha, and binds a heme cofactor.

INTERPRO description (entry IPR001128)

The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [MEDLINE:93135827].

QuickGO

PROCESS : electron transport (GO:0006118)

→ Alignments
 → Domain Structure

Figure 1: 2bzh Oxidoreductase(oxygenase)

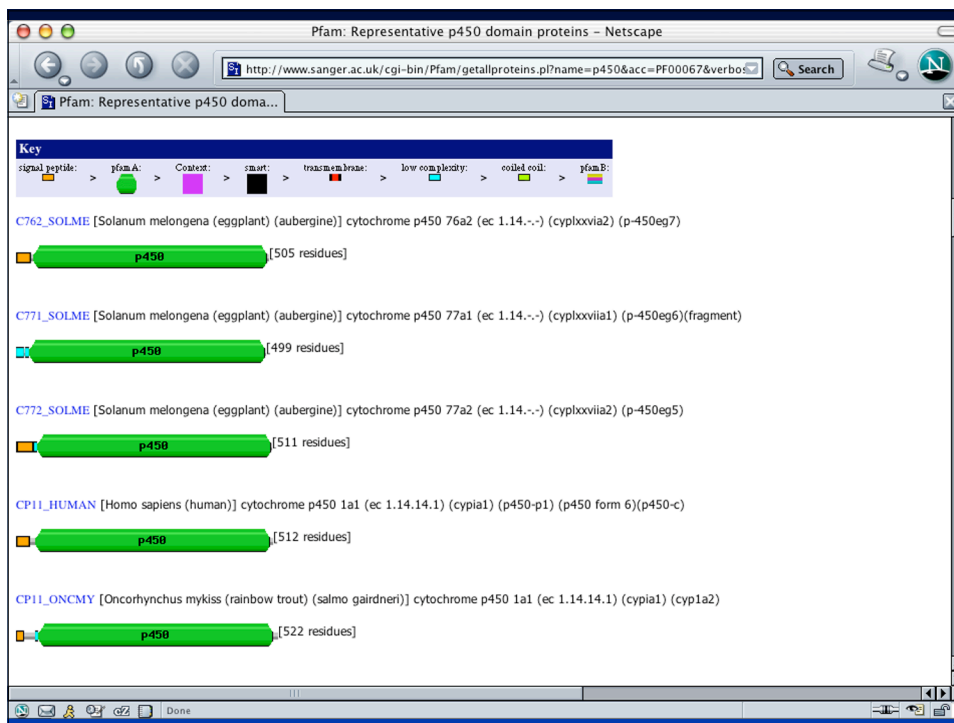
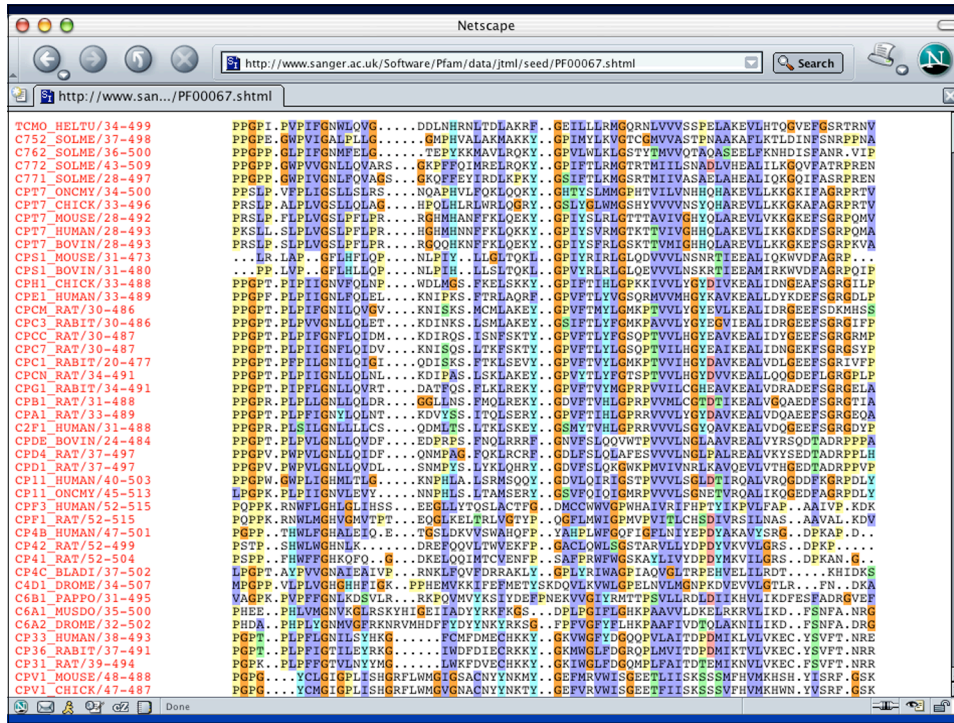
Key:

Domain	Chain	Start Residue	End Residue
p450	A	5	448
p450	B	5	448

The Swissprot/PDB mapping was provided by MSD

For additional annotation, see the PROSITE document P00C00081 [Expasy] SRS-UK | SRS-USA

NHGRI Current Topics in Genome Analysis 2003
 Biological Sequence Analysis II



NHGRI Current Topics in Genome Analysis 2003
 Biological Sequence Analysis II

InterPro: Cytochrome P450 – Netscape

http://www.ebi.ac.uk/interpro/DisplayProEntry?ac=IPR001128

InterPro: Cytochrome P450

InterPro Cytochrome P450 [? = help](#)

IPR001128 Cytochrome_P450 Matches: 3032 proteins
 View matches: [\[Overview\]](#), [\[sorted by Name\]](#), [\[of known structure\]](#), [\[Detailed view\]](#), [\[Table view\]](#)

Name [?](#) Cytochrome P450

Signatures [?](#) PF00067:p450 (2764 proteins)
 PR00385:P450 (2234 proteins)
 PS00086:CYTOCHROME_P450 (2315 proteins)
 SSF48264: Cytochrome_P450 (2931 proteins)

Type [?](#) Family

Dates [?](#) 1999-10-08 17:07:25.0 (created)
 2000-02-17 17:11:42.0 (modified)

Children [?](#) [\[tree\]](#) IPR002397: B-class P450
 IPR002399: Mitochondrial P450
 IPR002401: E-class P450, group I
 IPR002402: E-class P450, group II
 IPR002403: E-class P450, group IV

Process [?](#) electron transport ([GO:0006118](#))

Abstract [?](#) The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [1].

Structural links [?](#) PDB [1f4p](#)
 PDB [1jfb](#)
 SCOP [a_104.1](#)
 SCOP [c_23.5](#)

Database links [?](#) Blocks [IPB001128](#)
 PROSITE doc [PDOC00081](#)

Taxonomy [?](#) 3 Saccharomyces cerevisiae / Unclassified

InterPro: Cytochrome P450 – Netscape

http://www.ebi.ac.uk/interpro/DisplayProEntry?ac=IPR001128#

InterPro: Cytochrome P450

Dates [?](#) 1999-10-08 17:07:25.0 (created)
 2000-02-17 17:11:42.0 (modified)

Children [?](#) [\[tree\]](#) IPR002397: B-class P450
 IPR002399: Mitochondrial P450
 IPR002401: E-class P450, group I
 IPR002402: E-class P450, group II
 IPR002403: E-class P450, group IV

Process [?](#) electron transport ([GO:0006118](#))

Abstract [?](#) The cytochrome P450 enzymes constitute a superfamily of haem-thiolate proteins. P450 enzymes usually act as terminal oxidases in multicomponent electron transfer chains, called P450-containing monooxygenase systems and are involved in metabolism of a plethora of both exogenous and endogenous compounds. P450-containing monooxygenase systems primarily fall into two major classes: bacterial/mitochondrial (type I), and microsomal (type II). All P450 enzymes can be categorised into two main groups, the so-called B- and E-classes: P450 proteins of prokaryotic 3-component systems and fungal P450nor (CYP55) belong to the B-class; all other known P450 proteins from distinct systems are of the E-class [1].

Structural links [?](#) PDB [1f4p](#)
 PDB [1jfb](#)
 SCOP [a_104.1](#)
 SCOP [c_23.5](#)

Database links [?](#) Blocks [IPB001128](#)
 PROSITE doc [PDOC00081](#)

Taxonomy [?](#) 3 Saccharomyces cerevisiae / Unclassified

Species	Count
Saccharomyces cerevisiae	3
Caenorhabditis elegans	84
Fruit Fly	1620
Arthropoda	428
Chordata	830
Mouse	152
Human	169
Eukaryota	2612
Virus	3
Archaea	6
Bacteria	411
Cyanobacteria	10
Synechocystis PCC 6803	1
Rice spp.	166
Arabidopsis thaliana	351
Green Plants	977
Plastid Group	979
Other Eukaryotes	13

Conserved Domain Database (CDD)

- Identify conserved domains in a protein sequence
- “Secondary database”
 - Pfam A and B
 - Simple Modular Architecture Research Tool (SMART)
- Search performed using RPS-BLAST
 - Query sequence is used to search a database of precalculated position-specific scoring tables
 - *Not* the same method used by ProfileScan
- <http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml>



The screenshot shows the NCBI Conserved Domain Database (CDD) website. The browser window title is "NCBI Conserved Domain Database - Netscape". The address bar shows the URL "http://www.ncbi.nlm.nih.gov/Structure/cdd/cdd.shtml". The page features a navigation menu with links to PubMed, Entrez, BLAST, OMIM, Books, TaxBrowser, and Entrez Structure. A search bar is present with the text "Search Entrez Structure for" and a "Go" button. The main content area is titled "A Conserved Domain Database and Search Service, v1.62". It includes a "Run CD-Search:" section with a dropdown menu for "Search Database" set to "CDD v1.62 - 11088 PSSMs" and a "Submit Query" button. Below this is a text input field for a protein query in FASTA format, containing the sequence: ">deleted in colorectal cancer\nMENSLRCVWVPKLAFLVLFASLLSAHLQVTGFQIKAFATAI\nVIKWKKGDIHLALGMDERKQQLSNGSLLIQNILHSRHHK". A "Find CDs" button is located below the input field. To the right of the input field is a "by keyword:" label and an empty text box. The left sidebar contains links for "CDD Help", "NCBI Handbook", "CD-Search", "CDART", "Smart", "Pfam: US / UK", and "NEW COG".

NCBI CD-Search - Netscape
 http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi

NCBI Conserved Domain Search

RPS-BLAST 2.2.6 [Apr-09-2003]
 Query= local sequence: DELETED IN COLORECTAL CANCER
 (750 letters)
 Database: #cdd.v1.62
 11,088 PSSMs; 2,717,223 total columns

Masked-out region, low complexity

.. This CD alignment includes 3D structure. To display structure, download [Cn3D!](#)

PSSMs producing significant alignments:

Accession	Score	E (bits) value
gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of...	73.4	9e-14
gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of...	72.7	2e-13
gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of...	68.0	3e-12

NCBI CD-Search - Netscape
 http://www.ncbi.nlm.nih.gov/Structure/cdd/wrpsb.cgi#all2078917053

gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of internal repeats found in the plasma protein fibronectin. Its tenth fibronectin type III repeat contains an RGD cell recognition sequence in a flexible loop between 2 strands. Approximately 2% of all animal proteins contain the FN3 repeat; including extracellular and intracellular proteins, membrane spanning cytokine receptors, growth hormone receptors, tyrosine phosphatase receptors, and adhesion molecules. FN3-like domains are also found in bacterial glycosyl hydrolases.

CD-Length = 93 residues, 100.0% aligned
 Score = 73.4 bits (179), Expect = 9e-14

Query: 429 PSAPRDVVPVLVSSRFVRLSWRPPAEAKGNIQTFTVFFSREGDNREALNTTQPGSLQLT 488
 Sbjct: 1 P S P P T N L R V T D V T S T S V T L S W T P P E D D G G P I T G Y V V E Y R E K G S G D W K E V E T P G S E T S Y T 60

Query: 489 VGNLKP E A M Y T F R V V A Y N E W G P G E S S Q P I K V A T 521
 Sbjct: 61 L T G L K P G T E Y E F R V R A V N G G E S P P S E S V T V T T 93

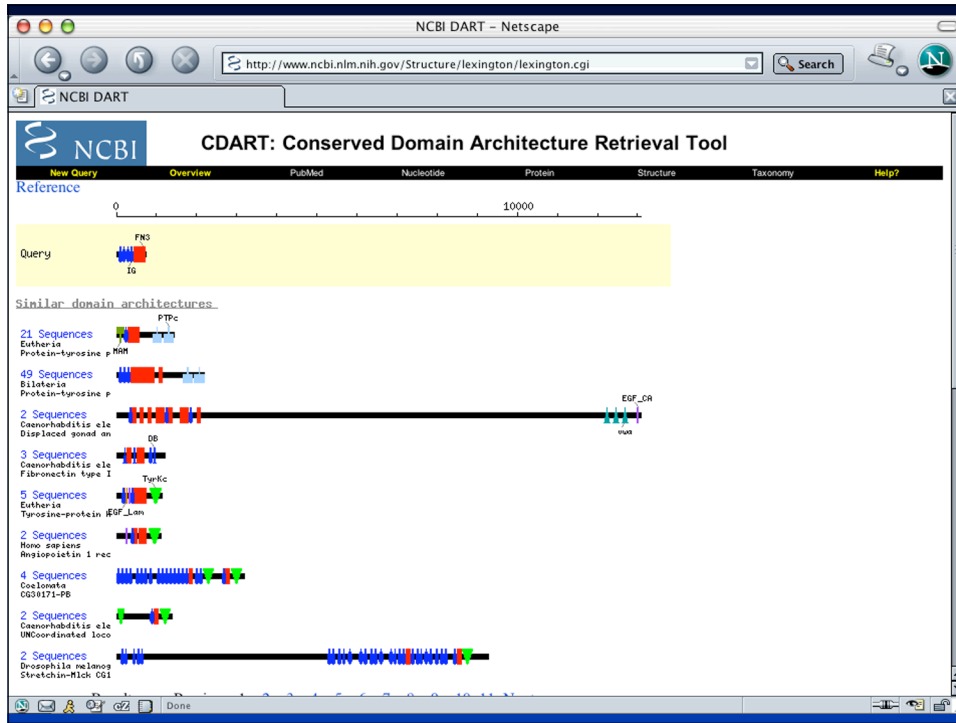
gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of internal repeats found in the plasma protein fibronectin. Its tenth fibronectin type III repeat contains an RGD cell recognition sequence in a flexible loop between 2 strands. Approximately 2% of all animal proteins contain the FN3 repeat; including extracellular and intracellular proteins, membrane spanning cytokine receptors, growth hormone receptors, tyrosine phosphatase receptors, and adhesion molecules. FN3-like domains are also found in bacterial glycosyl hydrolases.

CD-Length = 93 residues, 100.0% aligned
 Score = 72.7 bits (177), Expect = 2e-13

Query: 528 PGPVENLQAVSTSPSTSLITWEPAYANGPVQGYRFLCTEVS TGKEQNI EV---DGLSYK 584
 Sbjct: 1 P S P P T N L R V T D V T S T S V T L S W T P P E D D G G P I T G Y V V E Y R E K G S G D W K E V E T P G S E T S Y T 60

Query: 585 L E G L K K F T E Y S L R F L A Y N R Y G P G V S T D D I T V V T 617
 Sbjct: 61 L T G L K P G T E Y E F R V R A V N G G E S P P S E S V T V T T 93

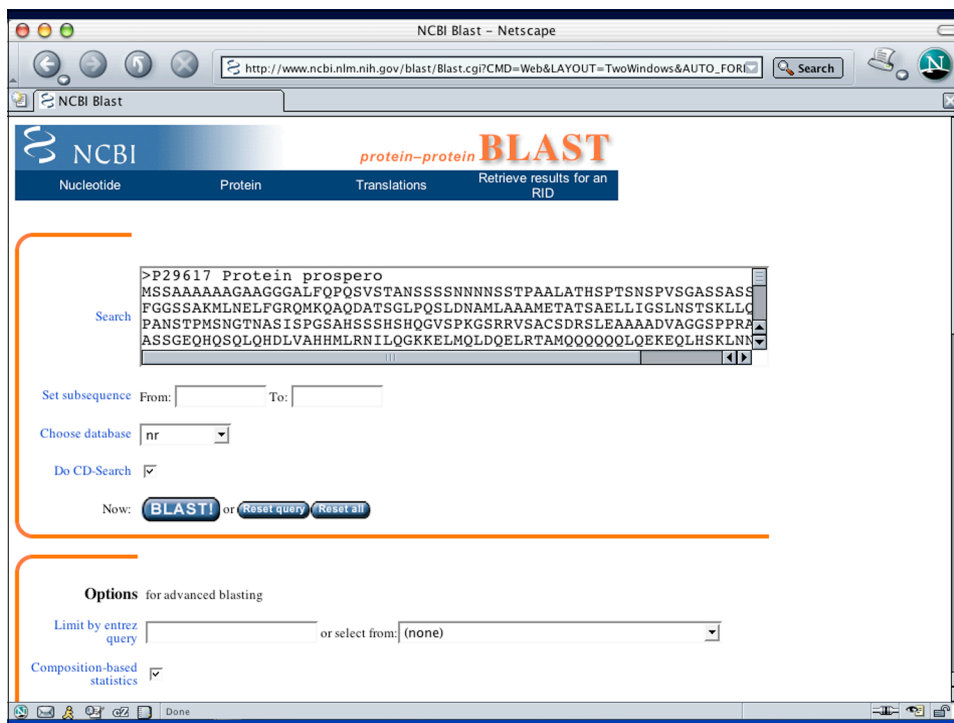
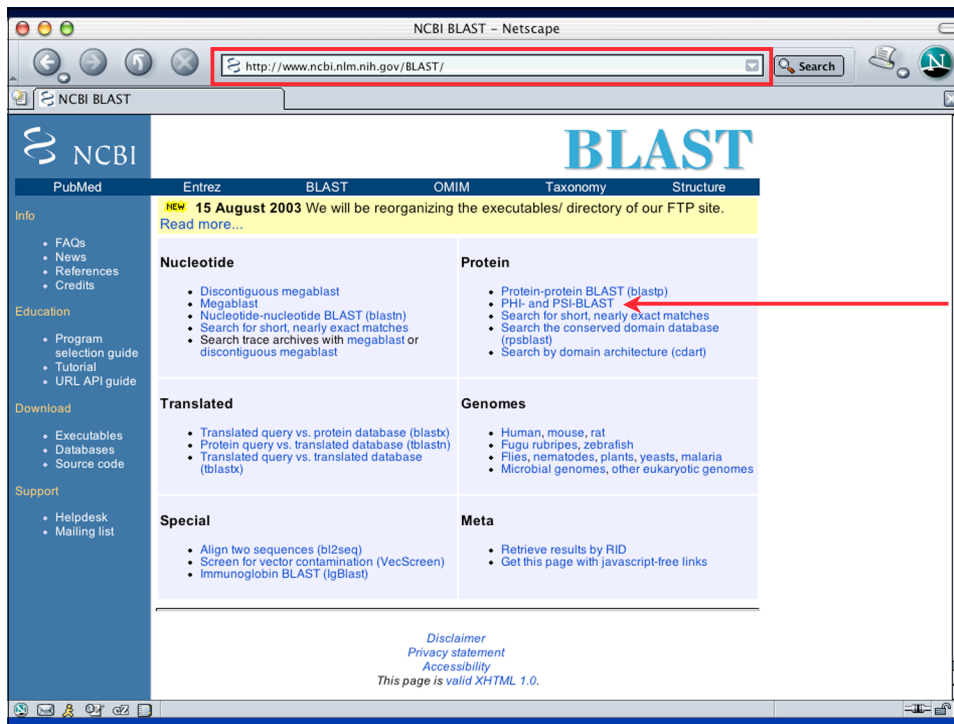
gnlCDDI14799.cd00063.FN3.Fibronectin type 3 domain: One of three types of internal repeats found in the plasma protein fibronectin. Its

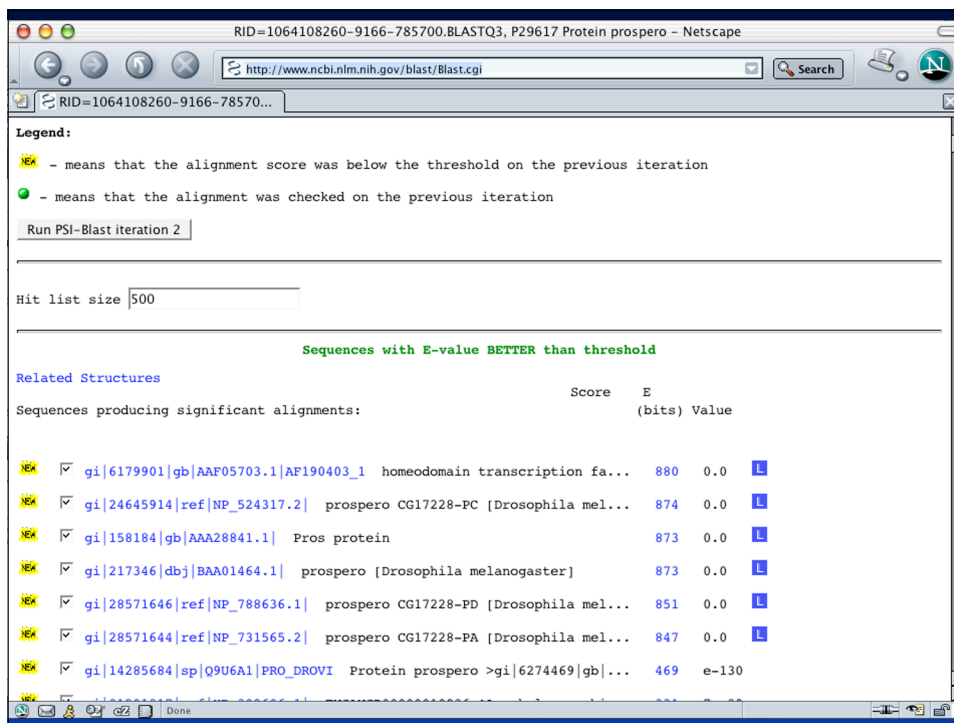
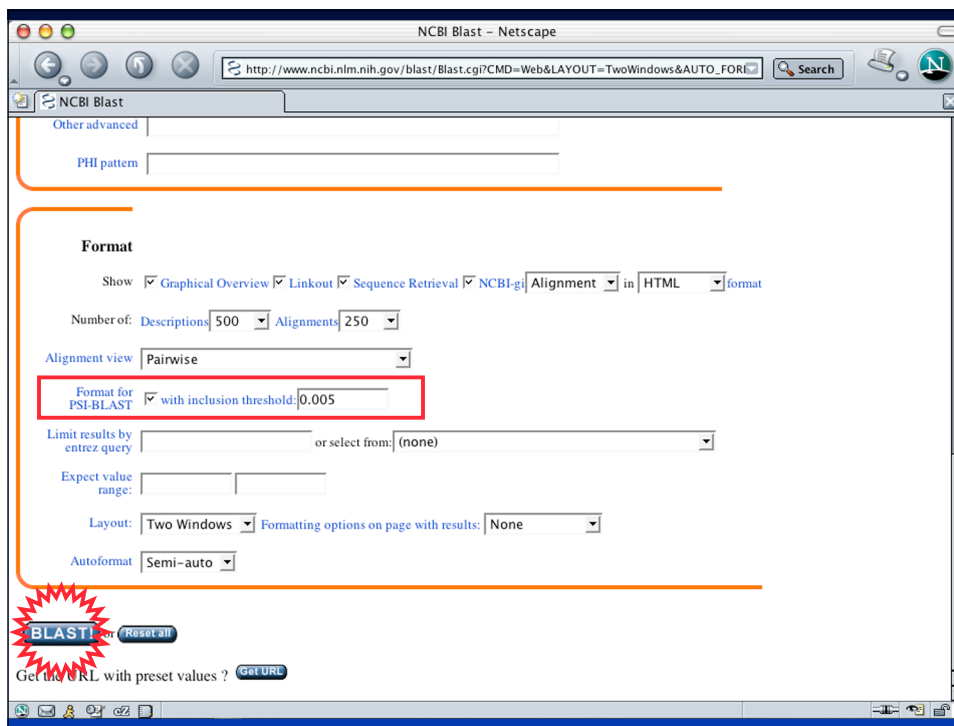


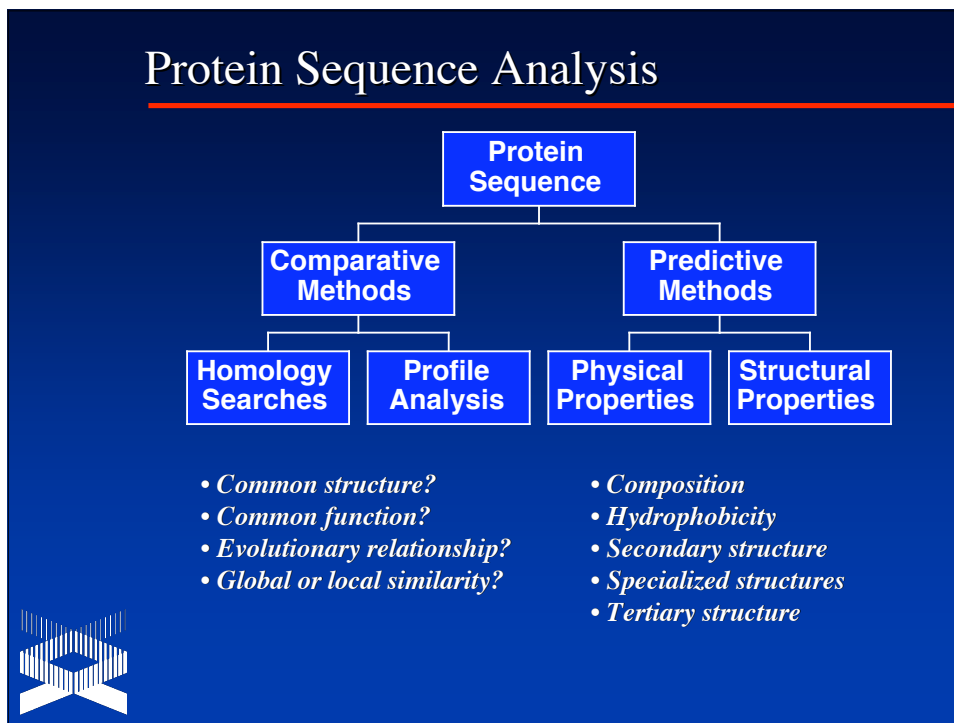
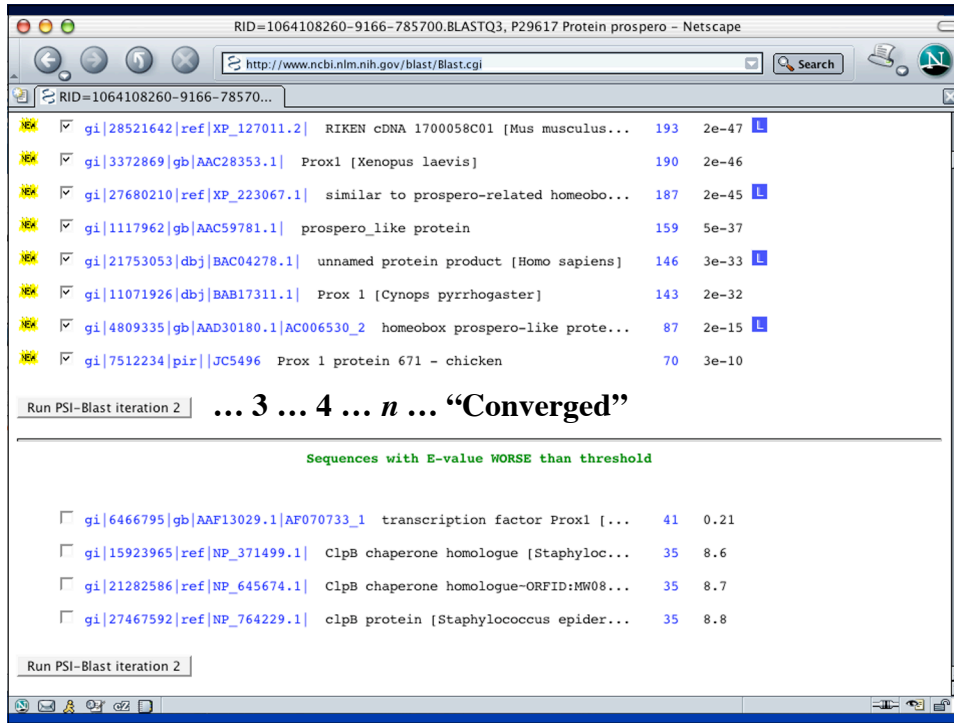
PSI-BLAST

- Position-Specific Iterated BLAST search
- Easy-to-use version of a profile-based search
 - Perform BLAST search against protein database
 - Use results to calculate a position-specific scoring matrix
 - PSSM replaces query for next round of searches
 - May be iterated until no new significant alignments are found
 - Convergence – all related sequences deemed found
 - Divergence – query is too broad, make cutoffs more stringent

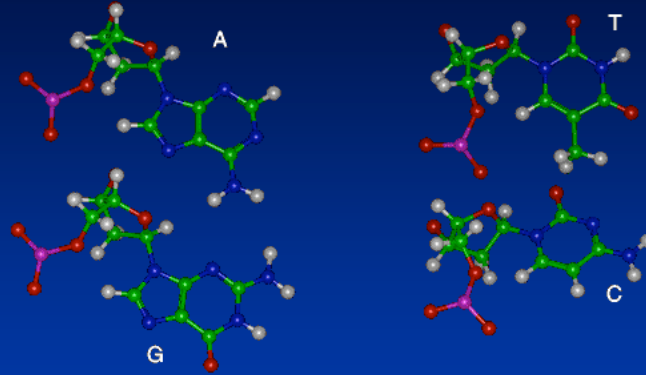






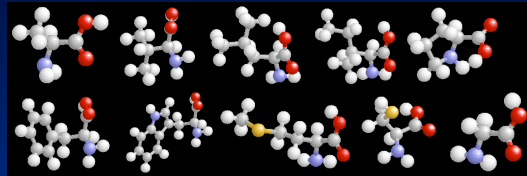


Information Landscape

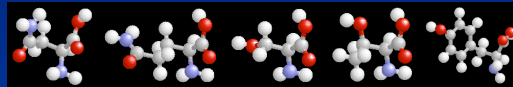


Information Landscape

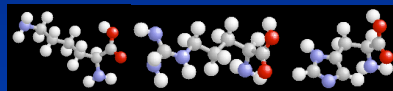
Nonpolar



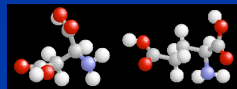
Polar Neutral



Polar Basic



Polar Acidic



ProtParam

- Computes physicochemical parameters
 - Molecular weight
 - Theoretical pI
 - Amino acid composition
 - Extinction coefficient
- Simple query
 - SWISS-PROT accession number
 - User-entered sequence, in single-letter format
- <http://www.expasy.ch/tools/protparam.html>



ProtParam Query

```
MNGEADCPTDLEMAAPKGDWRWSQEDMLTLLCEMKNLPSNDSKFKTTESHMDWEKVAFKDFSGDMCKL
KWVEISNEVRKFRITLTELILDAQEHVKNPYKGGKLLKHPDFPKKPLTPYFRFFMEKRAKYAKLHPM...
```

↓ Compute parameters

```
Number of amino acids: 727
Molecular weight: 84936.8
Theoretical pI: 5.44

Amino acid composition:

Ala (A) 35      4.8%      Leu (L) 57      7.8%
Arg (R) 39      5.4%      Lys (K) 97     13.3%
Asn (N) 28      3.9%      Met (M) 25      3.4%
Asp (D) 58      8.0%      Phe (F) 18      2.5%
Cys (C)  6      0.8%      Pro (P) 39      5.4%
Gln (Q) 36      5.0%      Ser (S) 67      9.2%
Glu (E) 98     13.5%     Thr (T) 22      3.0%
Gly (G) 26      3.6%      Trp (W) 11      1.5%
His (H) 11      1.5%      Tyr (Y) 20      2.8%
Ile (I) 18      2.5%      Val (V) 16      2.2%

Asx (B)  0      0.0%
Glx (Z)  0      0.0%
Xaa (X)  0      0.0%

Total number of negatively charged residues (Asp + Glu): 156
Total number of positively charged residues (Arg + Lys): 136
```



PROPSEARCH

- Uses amino acid composition to detect weak relationships
- Can be used to discern members of the same protein family
- 144 physical properties used in analysis (“vector”)
- Molecular weight
- Bulky residue content
- Average hydrophobicity and charge
- Search against “database of vectors” (PIR and SWISS-PROT)
- <http://www.embl-heidelberg.de/prs.html>



PROPSEARCH Query

```
>S18193 autoantigen NOR-90 - human
MNGEADCPDLEMAAPKGDWRWSQEDMLTLECMKNNLPSNDSSKFKTTESHMDWEKVAFKDFSGDMCKL
KWVEISNEVRKFRLLTELILDAQEHVKNPYKGGKLLKHPDFPKKPLTPYFRFFMEKRAKYAKLHPM...
```

↓ Vector search

Rank	ID	DIST	LEN2	POS1	POS2	pI	DE
1	>pl;s18193	0.00	727	1	727	5.33	autoantigen NOR-90 - human
2	ubf1_human	1.36	764	1	764	5.62	NUCLEOLAR TRANSCRIPTION FACTOR 1
3	ubf1_mouse	1.40	765	1	765	5.55	NUCLEOLAR TRANSCRIPTION FACTOR 1
4	ubf1_rat	1.57	764	1	764	5.61	NUCLEOLAR TRANSCRIPTION FACTOR 1
5	ubf1_xenla	3.95	677	1	677	5.79	NUCLEOLAR TRANSCRIPTION FACTOR 1
6	ubf2_xenla	4.18	701	1	701	6.05	NUCLEOLAR TRANSCRIPTION FACTOR 2
7	>pl;s57552	7.72	606	1	606	6.63	hypothetical protein YPR018w - yeast
8	>pl;150463	8.49	772	1	772	5.71	protein kinase - chicken
9	>pl;h54024	8.83	768	1	768	5.27	protein kinase (EC 2.7.1.37) cdc2-related
10	>pl;b54024	8.87	777	1	777	5.27	protein kinase (EC 2.7.1.37) cdc2-related
11	>pl;g54024	8.90	766	1	766	5.21	protein kinase (EC 2.7.1.37) cdc2-related
12	>pl;a55817	9.00	783	1	783	5.19	cyclin-dependent kinase p130-PITSLRE - mouse
13	>pl;f54024	9.11	777	1	777	5.30	protein kinase (EC 2.7.1.37) cdc2-related
14	>pl;e54024	9.11	779	1	779	5.42	protein kinase (EC 2.7.1.37) cdc2-related
15	yaa5_schpo	9.45	598	1	598	4.78	HYPOTHETICAL 69.5 KD PROTEIN C22G7.05
16	>pl;s62449	9.45	598	1	598	4.78	hypothetical protein SPAC22G7.05 - fission
17	>fl;158390	9.45	920	1	920	5.00	retinoblastoma binding protein 1 isoform I
18	>pl;s63193	9.58	590	1	590	6.15	hypothetical protein YNL227c - yeast
19	ynw7_yeast	9.58	590	1	590	6.15	HYPOTHETICAL 68.8 KD PROTEIN IN URE2-SSU72
20	>pl;s49634	9.74	899	1	899	4.79	hypothetical protein YML093w - yeast
21	ymj3_yeast	9.74	899	1	899	4.79	HYPOTHETICAL 103.0 KD PROTEIN IN RAD10-PRS4
22	radi_human	9.76	583	1	583	6.33	RADIXIN.
23	radi_pig	9.81	583	1	583	6.21	RADIXIN (MOESIN B).
24	>fl;178883	9.83	866	1	866	4.77	retinoblastoma binding protein 1 isoform II
25	>pl;b42997	9.87	754	1	754	5.17	retinoblastoma-associated protein 2 - human
26	>pl;a57467	9.91	647	1	647	5.74	RalBP1 - rat



PROPSEARCH Query

```
>S18193 autoantigen NOR-90 - human
MNGEADCPDLEMAAPKGGDRWSQEDMLTLECEMKNLPSNDSSKFKTTESHMDWEKVAFKDFSGDMCKL
KWWEISNEVRKFRFTLTELILDAQEHVKNPYKGGKLLKHPDFPKKPLTPYFRFFMEKRKRAYAKLHPM...
```

↓ Vector search

DIST Odds
 < 10 87.0%
 < 8.7 94.0%
 < 7.5 99.6%

	DIST	LEN2	POS1	POS2	pI	DE
	0.00	727	1	727	5.33	autoantigen NOR-90 - human
	1.36	764	1	764	5.62	NUCLEOLAR TRANSCRIPTION FACTOR 1
	1.40	765	1	765	5.55	NUCLEOLAR TRANSCRIPTION FACTOR 1
	1.57	764	1	764	5.61	NUCLEOLAR TRANSCRIPTION FACTOR 1
	3.95	677	1	677	5.79	NUCLEOLAR TRANSCRIPTION FACTOR 1
	4.18	701	1	701	6.05	NUCLEOLAR TRANSCRIPTION FACTOR 2
6	ubf2_wla	4.18	701	701	6.05	NUCLEOLAR TRANSCRIPTION FACTOR 2
7	>pl;s5752	7.72	606	606	6.63	hypothetical protein YPR018w - yeast
8	>pl;i50463	8.49	772	772	5.71	protein kinase - chicken
9	>pl;h54024	8.83	768	768	5.27	protein kinase (EC 2.7.1.37) cdc2-related
10	>pl;h54024	8.87	777	777	5.27	protein kinase (EC 2.7.1.37) cdc2-related
11	>pl;g54024	8.90	766	766	5.21	protein kinase (EC 2.7.1.37) cdc2-related
12	>pl;a55817	9.00	783	783	5.19	cyclin-dependent kinase p130-PITSLRE - mouse
13	>pl;f54024	9.11	777	777	5.30	protein kinase (EC 2.7.1.37) cdc2-related
14	>pl;e54024	9.11	779	779	5.42	protein kinase (EC 2.7.1.37) cdc2-related
15	yaa5_schpo	9.45	598	598	4.78	HYPOTHETICAL 69.5 KD PROTEIN C22G7.05
16	>pl;s62449	9.45	598	598	4.78	hypothetical protein SPAC22G7.05 - fission
17	>fl;i58390	9.45	920	920	5.00	retinoblastoma binding protein 1 isoform I
18	>pl;s63193	9.58	590	590	6.15	hypothetical protein YNL227c - yeast
19	yw7_yeast	9.58	590	590	6.15	HYPOTHETICAL 68.8 KD PROTEIN IN URE2-SSU72
20	>pl;s49634	9.74	899	899	4.79	hypothetical protein YML093w - yeast
21	ymj3_yeast	9.74	899	899	4.79	HYPOTHETICAL 103.0 KD PROTEIN IN RAD10-PRS4
22	radi_human	9.76	583	583	6.33	RADIXIN.
23	radi_pig	9.81	583	583	6.21	RADIXIN (MOESIN B).
24	>fl;i78883	9.83	866	866	4.77	retinoblastoma binding protein 1 isoform II
25	>pl;b42997	9.87	754	754	5.17	retinoblastoma-associated protein 2 - human
26	>pl;a57467	9.91	647	647	5.74	Ra1BP1 - rat

Secondary Structure Prediction

- Deduce the most likely position of alpha-helices and beta-strands
- Confirm structural or functional relationships when sequence similarity is weak
- Determine guidelines for rational selection of specific mutants for further laboratory study
- Basis for further structure-based studies

Folding Classes



Cyt c

Globins
 Orthogonal
 EF-hand
 Up-Down
 Cytochrome



CD4

Orthogonal
 Super-barrel
 Greek key
 Sandwich
 Jelly roll



Staph nuclease

Split sandwich
 Meander
 Metal-rich
 Open roll
 OB/UB roll



Triose phosphate isomerase

TIM barrel
 Doubly-wound



Neural Network

Output layer



Hidden layer



Input layer



nnpredict

- Neural network approach to making predictions
(Kneller et al., 1990)
- Best-case accuracy > 65%
- Search engines
 - E-mail nnpredict@celeste.ucsf.edu
 - Web <http://www.cmpharm.ucsf.edu/~nomi/nnpredict.html>



nnpredict Query

```
option: a/b
>flavodoxin - Anacystis nidulans
AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGYI
DDLDSVNFQGGKVVAYFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQPDLTKNRIKTWVSQKSEFGL
```

↓ / folding class

Tertiary structure class: alpha/beta

Sequence:
AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVG
ELQSDWEGYIDDLDSVNFQGGKVVAYFGAGDQVGYSDNFQDAMGILEEKISSLGSQTVGYW
PIEGYDFNESKAVRNNQFVGLAIDEDNQPDLTKNRIKTWVSQKSEFGL

Secondary structure prediction (H = helix, E = strand, - = no prediction):
---EEE-----EEHHHHHHH-----EEH-----EEEE-----
-----HHHH--EEEE-----H--HHHHHHH-----E--E-
-E-----HH--E-----EHHHH-----



PredictProtein

- Multi-step predictive algorithm (*Rost et al., 1994*)
 - Protein sequence queried against SWISS-PROT
 - MaxHom used to generate iterative, profile-based multiple sequence alignment (*Sander and Schneider, 1991*)
 - Multiple alignment fed into neural network (PHDsec)
- Accuracy
 - Average > 70%
 - Best-case > 90%
- Search engines
 - E-mail predictprotein@embl-heidelberg.de
 - Web <http://www.embl-heidelberg.de/predictprotein/>



PredictProtein Query

```

Joe Buzzcut
National Human Genome Research Institute, NIH
buzzcut@nhgri.nih.gov
# flavodoxin - Anacyctis nidulans
AKIGLFYGTQTGTQIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVGELQSDWEGIY
DDLDSVNFQGKVAIFGAGDQVGYSDNFQDAMGILEEKISSLSGQTVGYWPIEGYDFNESKAVRNNQFVG
LAIDEDNQPDLTKNRIKTWVSLKSEFGL
    
```

↓
Secondary structure

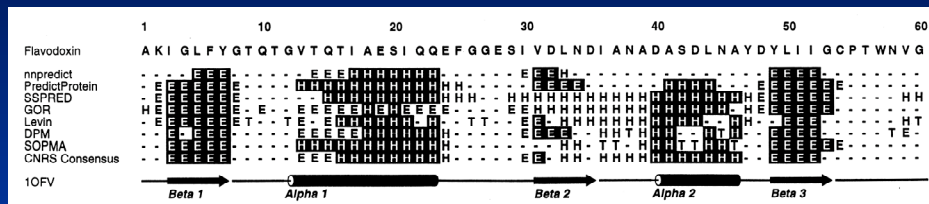
```

.....1.....2.....3.....4.....5.....6
AA      |AKIGLFYGTQTGTQIAESIQQEFGGESIVDLNDIANADASDLNAYDYLIIGCPTWNVG|
PHD sec |EEEEEE  HHHHHHHHHHHH  EEEEE HHH HHHH  EEEEE|
Rel sec |938999736982489999999999767982443213241278631241999861547765|
Detail:
prH sec |0000000001468999999999982100001111256538876432100000111111|
prE sec |0589988520000000000000000000000003665542100000000014899874120002|
prL sec |931000137985310000000000178985222344324511234554000114667776|
    
```

- SWISS-PROT hits
- Multiple alignment
- PDB homologues



Accuracy of Predictions



SignalP

- Neural network trained based on phylogeny
 - Gram-negative prokaryotic
 - Gram-positive prokaryotic
 - Eukaryotic
- Predicts secretory signal peptides
 (*not* those involved in intracellular signal transduction)
- <http://www.cbs.dtu.dk/services/SignalP/>



SignalP Query

```
>sp|P05019|IGFB_HUMAN INSULIN-LIKE GROWTH FACTOR IB PRECURSOR
MGKISSLPTQLFKCCFCDFLKVKMHMTSSSHLFYLALCLLFTSSATAGPETLCGAEIVDALQFVCGDRG
```

↓ N-terminal end only
 Eukaryotic set

```
***** SignalP predictions *****
Using networks trained on euk data

>IGF-IB      length = 195

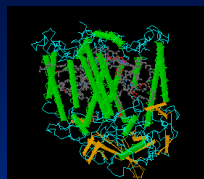
# pos aa      C      S      Y
.
.
.
46  A  0.365  0.823  0.495
47  T  0.450  0.654  0.577
48  A  0.176  0.564  0.369
49  G  0.925  0.205  0.855
50  P  0.185  0.163  0.376
.
.
.

< Is the sequence a signal peptide?
# Measure  Position  Value  Cutoff  Conclusion
max. C     49        0.925  0.37    YES
max. Y     49        0.855  0.34    YES
max. S     37        0.973  0.88    YES
mean S     1-48      0.550  0.48    YES
# Most likely cleavage site between pos. 48 and 49: ATA-GP
```

C = cleavage site score
 S = signal peptide score
 Y = combined score



Transmembrane Classes

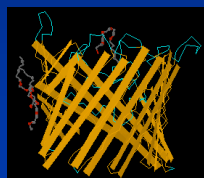


- Helix bundles
 Long stretches of apolar amino acids
 Fold into transmembrane alpha-helices
 “Positive-inside rule”

Cell surface receptors

Ion channels

Active and passive transporters



- Beta-barrel
 Anti-parallel sheets rolled into cylinder
Outer membrane of Gram-negative bacteria
Porins (passive, selective diffusion)



PHDtopology

- Approach similar to PredictProtein (PHDsec)
- Overall two-state accuracy **94.7%**
 - Accuracy of predicting helix **92.0%**
 - Accuracy of predicting loop **96.0%**
- Includes topology prediction
- Search engines
 - E-mail predictprotein@embl-heidelberg.de
 - Web <http://www.embl-heidelberg.de/predictprotein/>



PHDtopology Query

```

Joe Buzzcut
National Human Genome Research Institute, NIH
buzzcut@nhgri.nih.gov
predict htm topology
# pendrin
MAAPGGRSEPPQLPEYSCSYMVRPVSSELAFOQQHERRLQERKTLRESLAKCCSCSRKRAFGLVLTLPVILEWLPKYRV
KEWLLSDVISGVSTGLVATLQGMAYALLAAVPVGYGLYSAFFPILTYFIPGTSRHSVSGPPVVSMLVGSVVLSMAP...
    
```



```

          .....37.....38.....39.....40.....41.....42
AA      |YSLKYDYPLDGNQELIALGLGNIVCGVFRGFAGSTALRSRAVQESTGGKTQIAGLIGAI|
PHD htm |          HHHHHHHHHHHHH          HHHHHHHHHH|
Rel htm |36889999999999999986411046677776554312577888777621467788888|
detail: |
prH htm |31000000000000000012445788888877765321110000111135788899999|
prL htm |689999999999999999875542111111122234678889999888864211100000|
          :
          :
PHDThm |iiiiiiiiiiiiiiiiTTTTTTTTTTTTTTTTTTTTTccccccccccccTTTTTTT|
    
```



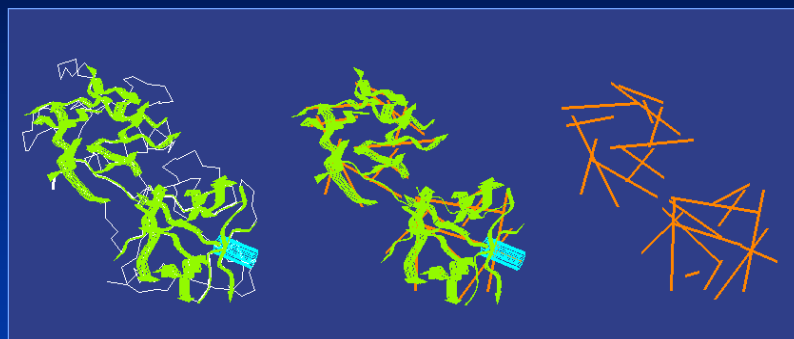
Predicting Tertiary Structure

- Sequence specifies conformation, *but* conformation does *not* specify sequence
- Structure is conserved to a much greater extent than sequence
 - Limited number of protein folds
- Similarities between proteins may not necessarily be detected through “traditional” methods



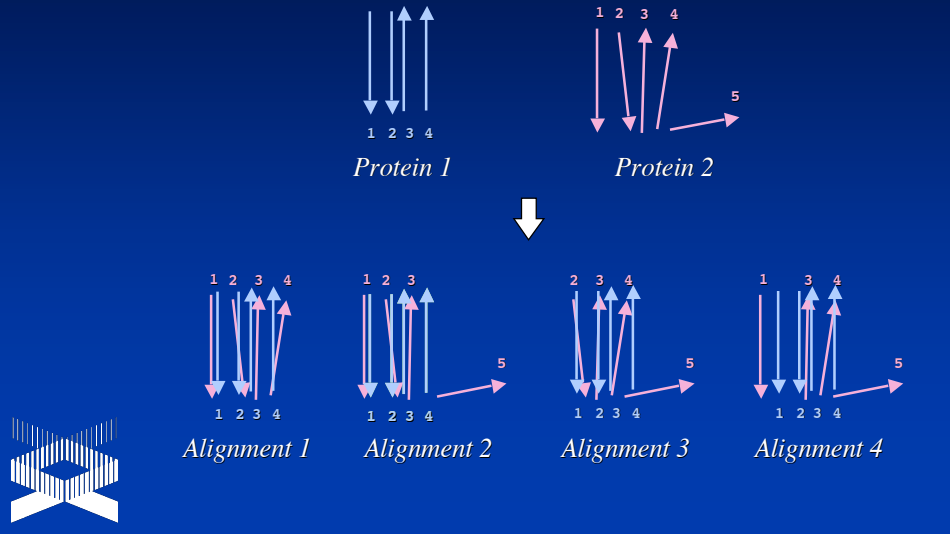
VAST Structure Comparison

Step 1: Construct vectors for secondary structure elements



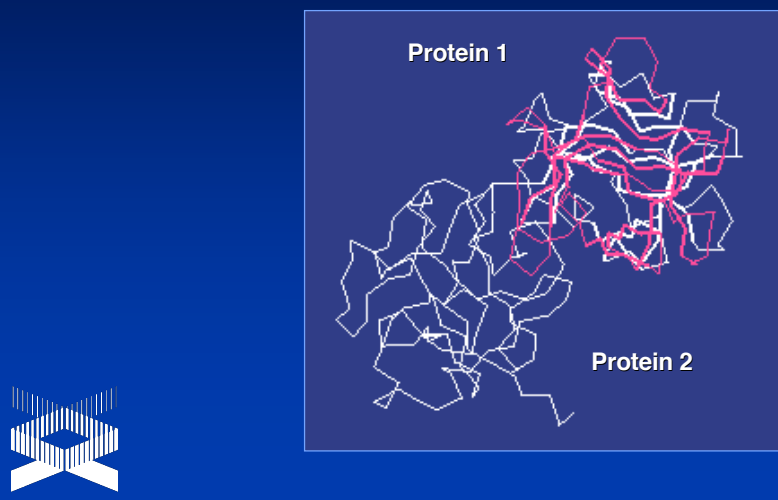
VAST Structure Comparison

Step 2: Optimally align structure element vectors



VAST Structure Comparison

Step 3: Refine residue-by-residue alignment



VAST Shortcomings

- Not the best method for determining structural similarities
- Reducing a structure to a series of vectors necessarily results in a loss of information (less confidence in prediction)
- Regardless of the “simplicity” of the method, provides a simple and fast first answer to the question of structural similarity



NCBI Structure Group - Netscape
http://www.ncbi.nlm.nih.gov/Structure/

NCBI Structure Group

PubMed Entrez BLAST OMIM Books TaxBrowser Entrez Structure

Search Entrez Structure for 2LIV Go

What's New?

- MMDB**
NCBI's structure database
- Cn3D v4.1**
3D-structure viewer
- CDD**
Conserved Domain Database
- VAST**
Structure comparisons
- VAST Search**
Submit structure database searches

The NCBI Structure Group

.. maintains MMDB, a database of macromolecular 3D structures, as well as tools for their visualization and comparative analysis. MMDB, the Molecular Modeling Database, contains experimentally determined biopolymer structures obtained from the Protein Data Bank (PDB).

Structure highlights

3D structure comparisons provide the basis for structure-neighboring in Entrez. The VAST database of structure neighbors contains millions of 3D superpositions and the corresponding alignments for proteins in MMDB. Structures not yet included in MMDB may be compared to database structures using the VAST Search service. More..

Text searches in MMDB (use the search toolbar at the top of this page) will yield *Structure Query* pages, providing access to entries which matched the keywords. *Structure Summary* pages for

try also:

- Structure summary via PDB/MMDB-Id: Go
- Read more about: MMDB, WWW-Entrez, VAST
- Resources: MMDB's FTP-site (including the MMDB database), NCBI Toolkit containing the MMDB-API and Cn3D source code.
- Publications in Entrez
- Research software: PKB and Threading (requires Splus).

Structure Summary - Netscape
 http://www.ncbi.nlm.nih.gov/Structure/mmdb/mmdbsrv.cgi?form=6&db=t&Dopt=s&uid=2

NCBI **MMDB** Structure Summary

PubMed BLAST Structure Taxonomy OMIM Help? Cn3d

Description: Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP).
Deposition: J.S.Sack, M.A.Saper & F.A.Quiocho, 10-Apr-89
Taxonomy: *Escherichia coli*
Reference: PubMed MMDB: 2778 PDB: 2LIV

View 3D Structure of Best Model with Cn3D Display [NEW Get Cn3D 4.1!](#)

Protein Chain 1-344
 3d Domains 1 2 1 2
 CDs LivK

Disclaimer | Write to the Help Desk
 NCBI | NLM | NIH

VAST Summary - Netscape
 http://www.ncbi.nlm.nih.gov/Structure/vast/vastsvr.cgi?did=6728

NCBI **VAST** Structure Neighbors

PubMed BLAST Structure Taxonomy OMIM Help? Cn3D

Query: MMDB 2778, 2LIV
Description: Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP)

View 3D Structure of All Atoms with Cn3D Display [NEW Get Cn3D 4.1!](#)

View Alignment using Hypertext for Selected VAST neighbors

List subset NR, Blast_p<10e-40 sorted by Aligned Length page 1 in Graphics

Find MMDB or PDB ids: or 3D-Domain ids:

60 out of 236 structure neighbors displayed.

ID	Length	Chain	CDs
2LIV	344	Chain	LivK
2LBP	344		
1EHT	332		
10P4	312		
1000	285		
1GUD	253		

VAST Summary - Netscape

http://www.ncbi.nlm.nih.gov/Structure/vast/vastsv.cgi

NCBI VAST Structure Neighbors

Query: MIMM 2778, 2LIV
 Description: Leucine(Slash)Isoleucine(Slash)Valine-Binding Protein (LIVBP)

View 3D Structure of All Atoms with Cn3D Display **NEW** Get Cn3D

View Alignment using Hypertext for Selected VAST neighbors

List subset NR, Blast_p<10e-7 sorted by Vast P_value page 1 in Table

Find MIMM or PDB ids: of 3D-Domain ids:

60 out of 146 structure neighbors displayed.

PDB	C	D	Ali.	Len.	SCORE	P-VAL	RMSD	%Id	Description
2LBP			344	43.5	10e-38.0	0.9	79.4	Leucine-Binding Protein (LBP)	
1DP4	C		312	29.0	10e-20.5	4.1	16.3	Dimerized Hormone Binding Domain Of The Atrial Natriuretic Peptide Receptor	
1Q00	B		285	32.6	10e-12.8	4.1	13.3	Amide Receptor Of The Amidase Operon Of Pseudomonas Aeruginosa (Amic) Complexed With The Negative Regulator Amir.	
1Q00	B 1		222	20.2	10e-12.6	2.8	15.8	Amide Receptor Of The Amidase Operon Of Pseudomonas Aeruginosa (Amic) Complexed With The Negative Regulator Amir.	

Cn3D Viewer

Rendering: Tubes

Coloring: Identity

Red – matches

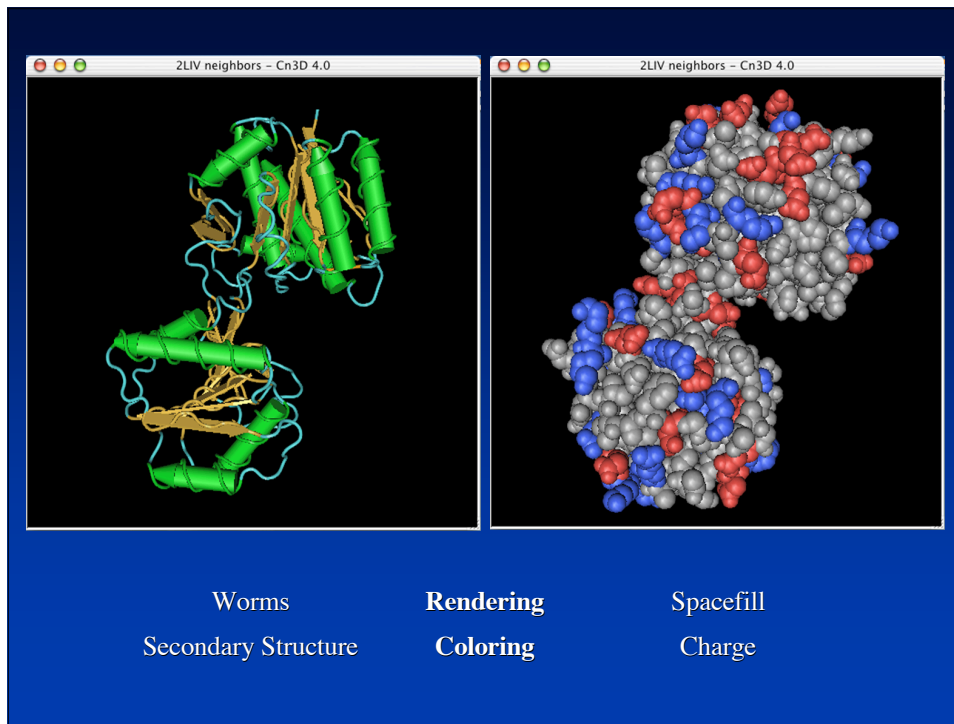
Blue – mismatches

2LIV neighbors - Cn3D 4.0

2LIV neighbors - Sequence/Alignment Viewer

```

2LIV  ...D I K V A V V G M S G P V A Q Y G D Q E F T G A E Q A V A D I N A K G G I K G N K L O I A K Y D D A C D P K Q A V A V A N K I V N D G I K Y V I G H L C S S S T Q P A S D I Y E D E G I L M I T P A
2LBP  ...D I K V A V V G M S G P I A Q W G I M E F N G A E Q A I K D I N A R G G I K G D K L V G V E Y D D A C D P K Q A V A V A N K I V N D G I K Y V I G H L C S S S T Q P A S D I Y E D E G I L M I S P G
    
```



SWISS-MODEL

- Automated comparative protein modelling server
- Web front-end at
<http://www.expasy.org/swissmod>
Results returned by E-mail

BLAST search to find similarities in PDB *by sequence*

Select templates with sequence identity > 25% and
projected model size > 20 amino acids

Generate models

Do energy minimization

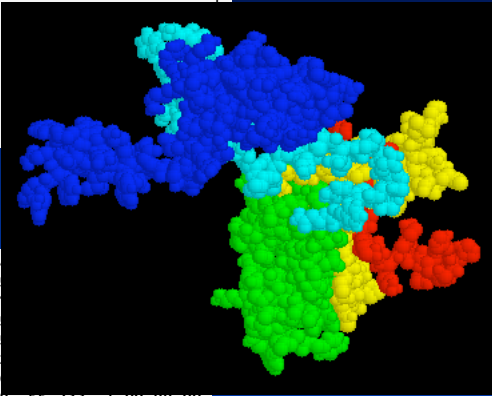
Generate PDB file for new protein model




```


21DJH.pdb: 42.77 % identity
21DJG.pdb: 42.77 % identity
11DJG.pdb: 42.22 % identity
11QAS.pdb: 44.17 % identity
11QAT.pdb: 43.52 % identity
21QAT.pdb: 43.52 % identity
21QAS.pdb: 43.52 % identity

Target:
21DJH.pdb
21DJG.pdb
11DJG.pdb
11QAS.pdb
11QAT.pdb
21QAT.pdb
21QAS.pdb
                
```





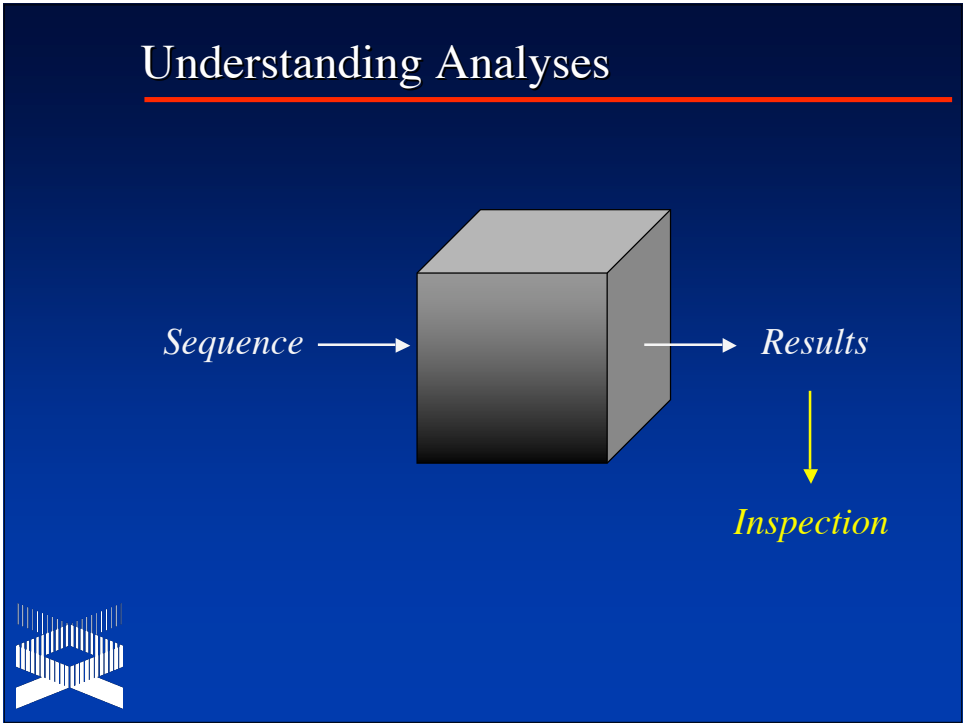
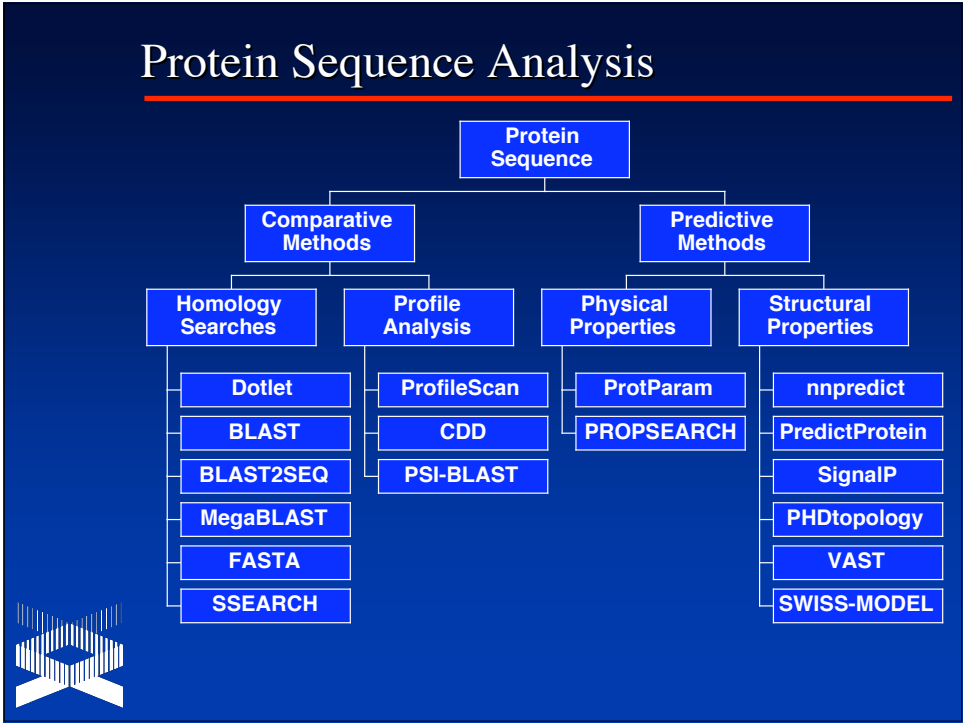
ATOM	1	H1	SER	1	24.219	22.95			
ATOM	2	H2	SER	1	24.770	21.43			
ATOM	3	N	SER	1	24.355	22.18			
ATOM	4	H3	SER	1	23.466	21.92			
ATOM	5	CA	SER	1	25.266	22.67			
ATOM	6	CB	SER	1	24.826	24.07			
ATOM	7	OG	SER	1	24.857	25.00			
ATOM	8	HG	SER	1	24.717	25.925	-55.253	1.00	99.00
ATOM	9	C	SER	1	25.471	21.750	-53.751	1.00	25.00
ATOM	10	O	SER	1	25.923	22.169	-52.684	1.00	25.00
ATOM	11	N	LYS	2	25.227	20.460	-53.972	1.00	25.00
ATOM	12	H	LYS	2	24.961	20.142	-54.878	1.00	99.00
ATOM	13	CA	LYS	2	25.366	19.408	-52.943	1.00	25.00
ATOM	14	CB	LYS	2	24.003	18.772	-52.622	1.00	25.00



Structural Modeling Software

- 3D-JIGSAW
<http://www.bmm.icnet.uk/servers/3djigsaw>
- ESyPred3D
<http://www.fundp.ac.be/urbm/bioinfo/esypred>
- MODELLER
<http://www.salilab.org/modeller/modeller.html>
- Protinfo
<http://protinfo.compbio.washington.edu>





Some lessons learned by bioinformaticians –
sometimes, the hard way

Simple Searches using Different Methods

- Pfam and SMART can be searched using
 - Pfam's search engine
 - SMART's search engine
 - NCBI's CDD search engine
- The methods used to search differ
 - Pfam and SMART use HMMs
 - NCBI-CDD uses RPS-BLAST
- The results obtained differ in a significant number of cases
 - CDD will miss entries representing short domains, repeats, and motifs



“Short Motif Pitfall”

- The level of sequence identity required for significant homology is much higher for smaller regions
- Two proteins may share a common domain while still being dissimilar elsewhere
- For very short motifs, homology *cannot* be inferred by sequence identity

→ *short motifs may not be helpful in describing what a protein does*



Immunoglobulin Signature

- Signature defined: [FY] – x – C – x – [VA] – x – H
- Precision
 - Total: 480 hits in 436 sequences
 - True positives: 390 hits
 - **False positives: 90 hits**
 - False negatives: 23 known

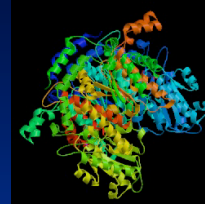
Acyl-CoA dehydrogenase
Acyl-amino acid-releasing enzyme
Alpha-adaptin A
GDP-mannose 6-dehydrogenase
Membrane alanyl aminopeptidase
Phosphatidyl cytidyl transferase
D-lactate dehydrogenase
DNA polymerase B
Hemerythrin
Anterior-restricted homeobox protein
Mast-stem cell growth factor
Limulus clotting factor C
Arachidonate 12-lipoxygenase

Amino adipate-semialdehyde dehydrogenase
DNA replication licensing factor
Neprin A
Cytochrome C-522
Phosphatidylinositol 3-kinase
Origin recognition complex subunit 2
Para-aminobenzoate synthase
Alpha-platelet-derived growth factor
Serine-threonine protein kinase
Photosystem II 44 kDa reaction center protein
DNA-directed RNA polymerase II (subunits)
Chloroplast 30S ribosomal protein S4
Titin



100% identity, but...

- **Phosphoglucose isomerase**
catalyzes interconversion of D-glucose-6-phosphate and D-fructose-6-phosphate
- **Neuroleukin**
secreted by T-cells, promotes survival of some embryonic spinal neurons and sensory nerves; B-cell maturation
- **Autocrine motility factor**
tumor cell product that stimulates cancer cell migration (metastasis?)
- **Differentiation and maturation mediator**
In vitro differentiation of human myeloid leukemia HL-60 cells to terminal monocytes



Jeffery et al., *Biochemistry* 39, 955-964, 2000

Proteins with Multiple Functions

Thymidine phosphorylase	Endothelial cell growth factor
Thymidylate synthase	Translation inhibitor
birA biotin synthase	bir operon repressor
Cystic fibrosis transmembrane conductance regulator (CFTR)	Regulates other ion channels
Crystallin	Enolase
	Lactate dehydrogenase
	Heat shock protein

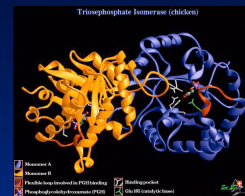


Does sequence similarity imply common function?

Maybe.

Structural Superfamilies: TIM Barrel

- Minimum 200 residues required for structure, with 160 residues structurally equivalent
- Structures mediate a wide variety of chemical reactions critical to biological survival
- May account for up to
 - 10% of all soluble enzymes
 - 10% of all proteins



Triose phosphate isomerase
Ribulose-phosphates
Thiamin phosphate synthase
FMN-linked oxidoreductases
NAD(P)-linked oxidoreductase
Glycosyltransferases
Metallo-dependent hydrolases
Aldolase
Enolase
Phosphoenol pyruvate
Malate synthase G
RuBisCo
Xylose isomerase-like proteins
Bacterial luciferase-like proteins
Quinolinic acid phosphoribosyltransferases
Cobalamin (B12)-dependent enzymes
tRNA-guanine transglycosylase
Dihydropterolate synthetase
Uroporphyrinogen decarboxylase
Methylenetetrahydrofolate reductase
Phosphoenolpyruvate mutase



Does structural similarity imply common function?

It depends.

Annotations

- Sequence-based annotations
 - Computational predictions on raw sequence data
 - Predictions are sometimes inconsistent
- Functional annotations
 - Review of functional annotations in *Mycoplasma* (Brenner, 1999)
 - 8% of functional annotations incorrect, with many inconsistent with known *Mycoplasma* biology and metabolism



A generalized strategy for avoiding the pitfalls

Predicting Function

- Understand the limitations of the programs being used
- Identify any special features in sequence
- Identify homologous proteins
- Identify protein family members based on sequence
- Look for structural homology
- Attempt to predict the function of the protein, with appropriate cautions in mind



Identify Special Features in Sequence

- Mask sequence to reduce biologically insignificant hits
 - Non-globular regions
 - Compositionally-biased regions
 - Coiled-coil regions
- Assay for putative transmembrane regions
 - May only be significant when similarity is global
- Perform secondary structure prediction
 - Best corroboration for BLAST or other homology-based match



Identify Homologous Proteins

- Search using specialized domain databases
 - Databases include PROSITE and Pfam
 - Short motifs are of limited utility in assessing function
 - Use multiple methods (ProfileScan, SMART, CDD)
 - Look for linkage of individual domains between proteins
- Search using BLAST
 - Use appropriate weight matrix
 - Using smaller subsequences of longer proteins reduces spurious matches (*e.g.*, against kinases)
 - Use known motifs or low-complexity regions as breakpoints



Identify Homologous Proteins

*Do NOT use sequence-based search methods as a
“black box.”*

*You MUST understand the methods and
optimize them on a case-by-case basis*



Identify Protein Family Members

- Perform iterative database searches to identify closely- and distantly-related family members
 - PSI-BLAST
 - MoST
- Construct a multiple sequence alignment
 - Look for conservation pattern between the unknown and the balance of the family to confirm presence of unique sequence features
 - Allows for assignment to family when there are few yet important sequence determinants
 - Keep in mind that sequence similarity is intransitive
 - $AB \sim BC$, and $BC \sim CD$, but $AB \not\sim CD$



Identify Protein Family Members

- “Prediction by Analogy”
 - Catalytic site residues are almost invariably polar
 - Large aromatic residues are often found to be involved in protein-ligand interactions
 - Zinc ions are coordinated by several residue types and, often, water molecules
 - Calcium ions are often bound by acidic residues and amides, although additional interactions occur with backbone atoms
 - Mn/Mg ions are often bound by two acidic residues separated by a hydrophobic residue in nucleases and glycotransferases



— Ponting, *Briefings in Bioinformatics* 2, 19-29, 2001

Identify Protein Family Members

- “Prediction by Analogy”
 - Phosphate and sulfate groups are found bound to the amino-terminus of alpha-helices in approximately half of all cases
 - Distinction can be made between disulfide-rich secreted proteins and zinc fingers, since the former occur in proteins with signal peptides and never possess substitutions of cysteine for histidine, or vice-versa



— Ponting, *Briefings in Bioinformatics* 2, 19-29, 2001

Look for Structural Homology

- ***Structure is more conserved than sequence***
- Comparison of two known structures
 - Vector-based (NCBI VAST)
 - Energy minimization methods
- Predictive modeling methods (sequence *vs.* structure)
 - SWISS-MODEL
 - Homology model building (“threading”)
 - *De novo* structure prediction



Predicting Function: Considerations

- The protein may actually have more than one function within the cell
- ***Never*** use database annotation as evidence of function...
 - when there are few homologues
 - when the homologues are not consistent
- Annotations are intransitive!
- Confirm database annotations in the literature
- Predictions are subjective



Predicting Function: Considerations

- Assure that the database hits and predictive methods based on sequence yield information that *make biological sense*
 - Predicted motifs or features biologically correct
 - Consistency with findings at the bench
- *Even if one is able to predict function, the prediction can indeed turn out to be incorrect – experimental proof is absolutely essential!*



Gene-Function Analysis

Genome	Complete set of genes of an organism	Systematic DNA sequencing
Transcriptome	Complete set of mRNA molecules present in a cell, tissue, or organ	Hybridization arrays SAGE High-throughput Northern
Proteome	Complete set of protein molecules present in a cell, tissue, or organ	2D gel electrophoresis Peptide mass fingerprinting Two-hybrid analysis
Metabolome	Complete set of metabolites (low-MW intermediates) in a cell, tissue, or organ	IR spectroscopy Mass spectroscopy NMR spectroscopy

