# VU MEDICAL ZONE

# Admin: Amaan Khan

# **Introduction to Bioinformatics**

# **BIF101** Final Merge ppt

# Lecture 96 to 150

#### Introduction

Protein databases store

- Protein sequences
- Motif
- Structure
- Structure alignments

#### Origin

First sequences to be collected were Proteins using **Sanger and Tupy's** methods (1951) where Common protein families like cytochromes were sequenced

#### Origin

Atlas of protein sequences (mainly cytochromes) was assembled by Margret Dayhoff and her collaborators at National Biomedical Research Foundation (NBRF) in 1960s

**PIR ((Protein Information Resource**) The collection of Dayhoff and co became **PIR** which is now a collaboration of NBRF, Munich Center for Protein Sequences (MIPS) and **J**apan International **Protein Information** Database (JIPID)

#### **Protein Sequences**

- Swiss-Prot is Collaboration between the SIB and EBI
- Weekly releases from about 50 servers across the world, the main source being ExPASy in Geneva

Sib B	oinformatics Resource Portal			Home	About Contact
	Query all databases 🗾	×	search help		
Visual Guidance	ExPASy is the SIB Bioinformatics Resource P	ortal which provides a	ccess to	Popular resources	
Categories	scientific databases and software tools (i.e., reso sciences including proteomics, genomics, phyloge	ources) in different area	opulation	UniProtKB	
proteomics	genetics, transcriptomics etc. (see Categories in th	e left menu). On this p	ortal you	🥑 SWISS-MODEL	
genomics	find resources from many different SIB groups as we	ell as external institutions	s.	STRING	
structural bioinformatics				ROSITE	
systems biology	Featuring today				
phylogeny/evolution	AACompSim	K. ALC - SPIE PRAFT		Latest News	5
population genetics	Compare amine acid composition of a	All States and States and States and States		Latest News	<b>2</b>
transcriptomics	UniProtKB entry with UniProtKB entries	1 10 10 10 10 10 10 10 10 10 10 10 10 10		Protein Spotlight: The h	idden things
biophysics	[details]			Nature has its secret way	s. During the
imaging				course of the 19th century Augustinian friar Gregor N	y, the ⁄lendel
IT infrastructure		0		worked out the basics of g inheritance as he crossbr	genetic ed pea plants.
drug design				More	

1.3

http://www.expasy.org/

#### **Protein Sequences**

- International partnership between PIR, EBI and SIB
- Created UniProt, by unifying the PIR-PSD, Swiss-Prot, and TrEMBL databases



About PIR     Databases     Search/Analysi       PRO     Image: PRO Hierarchy     (Note that the imp       6 shown of 223047 records     Image: Provide the imp       6 shown of 223047 records     Image: Provide the imp       Image: Provide the imp     Image: Provide the imp       6 shown of 223047 records     Image: Provide the imp       Image: Provide the imp     Image: Provide the imp       Image:	s     Download     Support       blicit relationship is is_a, whereas d indicat <t< th=""><th>tes <i>derives_from</i> relationship</th></t<>	tes <i>derives_from</i> relationship
PRO Hierarchy     (Note that the imp     6 shown of 223047 records     Sort (re)     Sort (sre)     G0:0032991 macromolecular complex     PR:000025493 LPS:GPI-anchored CD14     PR:000025494 LPS:secreted CD14	PMID       Taxon       PANTHER       EcoCyc         Synonym       Gene       MGI       HGNC       Pfam       PIRSF	tes <i>derives_from</i> relationship
6 shown of 223047 records	/PMID     Taxon     PANTHER     EcoCyc       /Synonym     Gene     MGI     HGNC     Pfam     PIRSF	V Definition Reactome VIniProtKB Category
• expand         • sort (in)         • sort (sin)           -         GO:0032991 macromolecular complex           +         PR:000025493 LPS:GPI-anchored CD14           +         PR:000025494 LPS:secreted CD14	Shortyme Center ( Hort C Hort C France	Category
GO:0032991 macromolecular complex     PR:000025493 LPS:GPI-anchored CD14     PR:000025494 LPS:secreted CD14		
+ PR:000025493 LPS:GPI-anchored CD14 + PR:000025494 LPS:secreted CD14		
PR:000025494 LPS:secreted CD14		complex
		complex
GO:0043234 protein complex		
PR:000018263 amino acid chain		
+ PR:00000001 protein		
Home   About PIR   Databases   Search/Analysis   Download   Sup ©2014 Protein Information Resource University of Delaware 15 Innovation Way, Suite 205 Newark, DE 19711, USA	e Georgetown University Medical Center 3300 Whitehaven Street, NW, Suite 1200 Washington, DC 20007, USA	SE



http://www.uniprot.org/

#### ProLINK

*i*ProLINK (*integrated* **P**rotein Literature, **IN**formation and **K**nowledge) has been developed as a resource to facilitate text mining in the area of literature-based database curation, named entity recognition, and protein ontology development. The collection of data sources can be utilized by computational and biological researchers to explore literature information on proteins and their features or properties (Hu *et al.*, 2004).





#### Bibliography Mapping/ Annotation Extraction

- Bibliography Display/Submission
- Annotation-Tagged Corpora
- RLIMS-P Text Mining Tool
- eFIP Text Mining Tool
- eGIFT Text Mining Tool
- iSimp Sentence Simplification System

#### Entity Recognition/ Ontology Development

- Name Tagging Guidelines/Corpora
- Protein Ontology Development

#### Collaborators

#### iProLINK Paper

http://www.uniprot.org/

#### **Tutorial**





January Molecule of the Month



http://www.rcsb.org/pdb/home/home.do

#### SCOPe: Structural Classification of Proteins — extended. Release 2.04 (July 2014, new entries added 2014-12-18)

Browse Stats & History ASTRAL Subsets Downloads Related Resources References Help About

#### Welcome to SCOPe!

SCOPe is a database developed at the Berkeley Lab and UC Berkeley to extend the development and maintenance of SCOP.

SCOP was conceived at the MRC Laboratory of Molecular Biology, and developed in collaboration with researchers in Berkeley.

Work on SCOP (version 1) concluded in June 2009 with the release of SCOP 1.75.

SCOPe classifies many newer structures through a combination of automation and manual curation, and corrects some errors in SCOP,

aiming to have the same accuracy as the hand-curated SCOP releases. SCOPe also incorporates and updates the ASTRAL database.

For prior releases, click on the <u>Stats & History</u> tab above. For more info, click on the <u>About</u> tab above. New PDB entries were last added on **2014-12-18**; for more info on periodic updates click on the <u>Help</u> tab above.

Scarch SCOPe (example): Search

#### **Classes in SCOPe 2.04:**

- 1. a: All alpha proteins [46456] (285 folds)
- 2. 3 b: All beta proteins [48724] (176 folds)
- 3. 🔅 c: Alpha and beta proteins (a/b) [51349] (148 folds)
- 4. d: Alpha and beta proteins (a+b) [53931] (380 folds)
- 5. de: Multi-domain proteins (alpha and beta) [56572] (68 folds)
- 6. f: Membrane and cell surface proteins and peptides [56835] (57 folds)
- 7. 23 g: Small proteins [56992] (91 folds)
- 8. h: Coiled coil proteins [57942] (7 folds)
- 9. We will be structures [58117] (25 folds)





#### http://scop.mrc-Imb.cam.ac.uk/scop/



1

#### http://scop.mrc-Imb.cam.ac.uk/scop/

#### Conclusions

- First sequences to be collected were Protein sequences
- Protein databases are classified on the basis of sequences, motifs and structures
- Growth of Sequence in Databases is exponential

#### Origin

- First attempt to sequence free living organism was launched in late 1990's
  - (Blattner et al. 1997)
- Viruses had already been sequenced

#### Origin

- Haemophilus influenzae was the first published genome

   (Fleischmann et al. 1995)
- The project initiated at The Institute of
   Genome Research
   (TIGR) under
   leadership of Craig
   Ventor

#### Origin

- Use of shotgun
   sequencing method
   speeded up the process
- 1.8 million bp
- Took 9-months and cost was 1 million
- Paved the way for sequencing of many other organisms

#### **Examples**

- AceDB (A C. elegans
   DataBase) was the first genome database
   developed in 1989
  - by Richard durbin and Thierry-Miegi

#### AceDB



#### http://www.acedb.org/

#### **Examples**

 TAIR (The Arabidopsis Information Resource)
 http://www.arabidopsis.org/ and
 SGB (Saccharomyces
 Genome Database)
 http://www.yeastgenome.org/
 utilized AceDB system

#### **Human Genome Project**

 Pilot project, the Human Genome Initiative begun by Department Of Energy (DOE) in 1986



#### **Human Genome Project**

- National Human Genome Research Initiative (NHGRI), a federally funded organization (NIH) started in 1988 (Francis Collin)
- Celera, a commercial vendor (Craig Venter) joined the venture in 1998



#### Human Genome Project

- Both simultaneously announced the completion of the project in 2000
- Total 3.4 billion bases sequenced at a cost of \$1/base

#### **UCSC** Genome Bioinformatics Genomes -Blat - Tables -Gene Sorter - PCR - VisiGene -Session -FAQ - Help About the UCSC Genome Bioinformatics Site Genome Browser Welcome to the UCSC Genome Browser website. This site contains the reference sequence and working draft assemblies for a large collection of genomes. It also provides portals to ENCODE data at UCSC (2003 to 2012) and to the Neandertal project. Download or purchase the Genome Ebola Browser source code, or the Genome Browser in a Box (GBiB) at our online store. Blat We encourage you to explore these sequences with our tools. The Genome Browser zooms and scrolls over chromosomes, showing the work of annotators worldwide. The Gene Sorter shows expression, homology and other information on groups of genes that can be related in many ways. Table Browser Blat guickly maps your sequence to the genome. The Table Browser provides convenient access to the underlying database. VisiGene lets you browse through a large collection of in situ mouse and frog images to examine expression patterns. Genome Graphs allows you to upload and Gene Sorter display genome-wide data sets. In Silico PCR The UCSC Genome Browser is developed and maintained by the Genome Bioinformatics Group, a cross-departmental team within the UC Santa Cruz Genomics Institute and the Center for Biomolecular Science and Engineering (CBSE) at the University of California Santa Cruz (UCSC). If you Genome have feedback or questions concerning the tools or data on this website, feel free to contact us on our public mailing list. Graphs The Genome Browser project team relies on public funding to support our work. Donations are welcome -- we have many more ideas Galaxy DONATE NOW than our funding supports! If you have ideas, drop a comment in our suggestion box. VisiGene Utilities f News News Archives **Downloads** To receive announcements of new genome assembly releases, new software features, updates and training seminars by email, subscribe to the genome-announce mailing list. Please see our blog for posts about Genome Browser tools, features, projects and more. Release Log

#### http://genome.ucsc.edu/



http://genome.ucsc.edu/

#### Conclusion

- Success of *Haemophilus influenzae* paved the way for other genome sequencing projects
- Human Genome Project was accomplished by NHGRI and Celera
- Genome browsers help viewing Genome features

#### Gene Expression Omnibus (GEO)

 A public repository for the archiving and distribution of gene expression data submitted by the scientific community

Gene Expression Omnibus (GEO)

- MIAME compliant data
  - Minimum Information About a Microarray Experiment <u>http://www.mged.org/</u> <u>Workgroups/MIAME/</u> <u>miame.html</u>

Gene Expression Omnibus (GEO)

- Convenient for deposition of gene expression data, as required by funding agencies and journals
- Curated resource for gene expression data
- Browsing, querying, analysis and retrieval

Home - GEO - NCBI	* S GEO Accessio	n viewer	× S GEO Accession viewer × +			H <sub>2</sub>
Image: Second				∀ C' (8.+	♀ ☆ 自 ♣	<b>☆</b> Z ≡
🔯 Most Visited 👻 🍓 Getting Started						
SNCBI Resources 🕑 How To 🖸					S	ign in to NCBI
GEO Home Documentation -	Query & Browse 🔻	Email GE	0			
	Search GEO DataSet	s				
	Search GEO Profiles				Image: Content         Repository Browser         DataSets:       3848         Series:       51804	
Gene Expression	Analyze with GEO2R					
-	GEO BLAST					
GEO is a public functional genomics da	Programmatic Access	<i>0</i>	ompliant data submissions. Array- and		Gene Expre	ssion Omnibus
sequence-based data are accepted. It	Repository Browser		y and download experiments and curated			
gene expression premes.	FTP Site				Keyword or GEO Accession	Search
Getting Started		Tools		Browse	Content	
Overview Searc		Search	for Studies at GEO DataSets	Repository	Browser	
FAQ		Search	for Gene Expression at GEO Profiles	DataSets:	3848	
About GEO DataSets		Search GEO Documentation		Series: 🖾	51804	Search
About GEO Profiles Analyz		Analyze	a Study with GEO2R	Platforms:	13522	
About GEO2R Analysis		GEO BL	LAST	Samples:	1259935	
How to Construct a Query		Program	nmatic Access			
How to Download Data		ETP Site	8			

#### http://www.ncbi.nlm.nih.gov/geo/

#### **Gene Architecture**

GEO has four kinds of records

- Sample (GSM) preparation and description of the samples
- Platform (GPL) technology used and the features detected
   microarray or RNASeq

#### **Gene Architecture**

GEO has four kinds of records

### • Series (GSE)

defines a set of samples and how they are related

## Datasets (GDS)

sample data collections assembled by GEO

	Home – GEO – NCBI	× S GEO Accessio	on viewer 🗙 😒 GEO Accession viewer 🗙 🕂			
♦ ♦ ⊗ www	/.ncbi.nlm. <b>nih.gov</b> /geo/			∀ C' (8 -	오 ☆ 습 ♣	ΛZ
🔄 Most Visited 🔻	📵 Getting Started					
S NCBI RE	esources 🕑 How To 🕑				<u>S</u>	ign in to NCB
GEO Home	Documentation -	Query & Browse 🔻	Email GEO			
GEO is a publi sequence-base gene expressio	c functional genomics da ed data are accepted. To on profiles.	Omnibus Ita repository supporting Iols are provided to help	MIAME-compliant data submissions. Array- and users query and download experiments and curated	1	Gene Expres	Ssion Omnibus
Getting Sta	rted		Tools	Browse Cont	ent	
Overview			Search for Studies at GEO DataSets	Repository Brows	Ser	
FAO			Search for Gene Expression at GEO Profiles	DataSets:	3848	
FAQ About GEO Da	taSets		Search for Gene Expression at GEO Profiles Search GEO Documentation	DataSets:	3848 51804	
FAQ About GEO Da About GEO Pro	ntaSets ofiles		Search for Gene Expression at GEO Profiles Search GEO Documentation Analyze a Study with GEO2R	DataSets: Series: 🔊 Platforms:	3848 51804 13522	
FAQ About GEO Da About GEO Pri About GEO2R	ntaSets ofiles Analysis		Search for Gene Expression at GEO Profiles Search GEO Documentation Analyze a Study with GEO2R GEO BLAST	DataSets: Series: S Platforms: Samples:	3848 51804 13522 1259935	
FAQ About GEO Da About GEO Pri About GEO2R How to Constri	ntaSets ofiles Analysis Jot a Query		Search for Gene Expression at GEO Profiles Search GEO Documentation Analyze a Study with GEO2R GEO BLAST Programmatic Access	DataSets: Series: 🔊 Platforms: Samples:	3848 51804 13522 1259935	
FAQ About GEO Da About GEO Pro About GEO2R How to Constru How to Downlo	ataSets ofiles Analysis uct a Query ad Data		Search for Gene Expression at GEO Profiles Search GEO Documentation Analyze a Study with GEO2R GEO BLAST Programmatic Access FTP Site	DataSets: Series: 🔊 Platforms: Samples:	3848 51804 13522 1259935	
FAQ About GEO Da About GEO Pri About GEO2R How to Constri How to Downlo Information Login to Subm	ataSets ofiles Analysis uct a Query vad Data for Submitters it		Search for Gene Expression at GEO Profiles Search GEO Documentation Analyze a Study with GEO2R GEO BLAST Programmatic Access FTP Site Submission Guidelines	DataSets: Series: Series: Series: Series: Series: Samples: Samples: MIAME Standard	3848 51804 13522 1259935 s	

) 🕙 www.ncbi.nlm.nih.gov/geo/summary/		∀ C⁴ (8 ▼					
Most Visited 👻 🧕 Getting Started							
blic holdings			Total holdings				
Series Platforms Samples Organisms History		Series	Public 51,806	Unreleased 7,896	Total 59,702		
Series type	Count	Platforms Samples	13,522 1,259,985	438 192,394	13,960 1,452,37		
Expression profiling by array 3				14-144.5	1.200		
Exoression profiling by genome tiling array	612	-					
Expression profiling by high throughput sequencing	3,446						
Expression profiling by SAGE	241						
Expression profiling by MPSS	20						
Expression profiling by RT-PCR	269						
Expression profiling by SNP array	12						
Genome variation profiling by array	557						
Genome variation profiling by genome tiling array	978						
Genome variation profiling by high throughput sequencing	56						
Genome variation profiling by SNP array	726						
Genome binding/occupancy profiling by array	156						
Genome binding/occupancy profiling by genome tiling array	2,042						
Genome binding/occupancy profiling by high throughput sequencing	3,317						
Genome binding/occupancy profiling by SNP array	11						
Methylation profiling by array	476						
Methylation profiling by genome tiling array	579						
Methylation profiling by high throughput sequencing	563						
Methylation profiling by SNP array	9						
· · · · · · · · · · · · · · · · · · ·			and the second				
--	---	---	--	----------------------	--------------		
( www.ncbi.nlm.nih.go	v/gds/?term=colon+cancer+RNASeq 📴 🧟 🗸 sada hai daasta	aane haram	Q ☆ 自	. ♦	z =		
🔯 Most Visited 🔻 📵 Getting	Started						
SNCBI Resources 🖸	How To 🖸		w_haider	My NCBI	Sign Out		
GEO DataSets	GEO DataSets  Colon cancer RNASeq Save search Advanced	0	Search		Help		
Show additional filters Entry type Series (4)	Display Settings: ♥ Summary, Sorted by Default order Send to: ♥ ★ Did you mean: colon cancer ma seq (60 items) Results: 4	Filters: Mana Top Org Homo sap Mus muso	age Filters ganisms [Tree piens (4) culus (2)	<u>e]</u>			
Organism Select Study type Expression profiling by array More Author	<ul> <li>Tumor cell-specific inhibition of MYC function using small molecule inhibitors of the HUWE1</li> <li>ubiquitin ligase         (Submitter supplied) Deregulated expression of MYC is a driver of colorectal carcinogenesis, necessitating novel strategies to inhibit MYC function. The ubiquitin ligase HUWE1 (HECTH9, ARFBP1, MULE) associates with both MYC and the MYC-associated protein MIZ1. We show here that HUWE1 is required for growth of colorectal cancer cells in culture and in orthotopic xenograft models. Using high throughput screening, we identify small molecule inhibitors of HUWE1, which inhibit MYC-dependent transactivation in colorectal     </li> </ul>	Find related Database:	d data Select	\$			
Select Attribute name tissue strain More	cancer cells, but not in stem and normal colon epithelial cells. more Organism: Homo sapiens Type: Expression profiling by high throughput sequencing; Genome binding/occupancy profiling by high throughput sequencing Platform: GPL10999 22 Samples Download data: GEO (TXT, WIG), SRA SRP044171 Series Accession: GSE59223 ID: 200059223	Search deta ("colonic OR colon o RNASeg[Al	ails neoplasms" cancer[All l Fields]	[MeSH Te Fields])	rms ] AND		
Publication dates 30 days 1 year Custom range	<u>PubMed</u> Similar studies     Sequential ChIP-bisulfite sequencing enables direct genome-scale investigation of chromatin and <u>DNA methylation cross-talk</u> (Submitter supplied) Cross-talk between DNA methylation and histone modifications drives the establishment     of composite epigenetic signatures and is traditionally studied using correlative rather than direct     approaches. Here we present sequential ChIP-bisulfite-sequencing (ChIP- BS-seq) as an approach to	Recent acti	<b>ivity</b> Incer RNA Seq	<u>Turr</u> (60)	See more		

BI > GEO > Accer	GEO Publications F	AQ MIAME Email GEO Not logged in Login	2
ope: Self	Format: HTML + Amount: Quick + GEO accession: GSE57	043 GO	
Series GSE5704	Query DataSets for	GSE57043	
1840 AND	Public on Apr 25, 2014		
îtle	Dicer knockout NSCLC RNAseq and miRseq		
Organism	Mus musculus		
xperiment type	Expression profiling by high throughput sequencing Non-coding RNA profiling by high throughput sequencing		
jummary	Dicer knockout NSCLC mRNAseq profiles the transcriptome, Dice NSCLC miRseq profiles the miRnome	er knockout	
)verall design	DicerHet and DicerKO NSCLC, 2 biological reps each genotype for r biological rep each for miRseq	mRNAseq, 1	
Contributor(s)	Sharp PA, Chen S		
litation(s)	Chen S, Xue Y, Wu X, Le C et al. Global microRNA depletion suppres angiogenesis. Genes Dev 2014 May 15;28(10):1054-67. PMID: 247	sses tumor 188094	
Submission date	Apr 24, 2014		
ast update date	Oct 14, 2014		
Contact name	Sidi Chen		
-mail	chensidi@mit.edu		
hone	7734144158		
Organization name	MIT		
Pepartment	Biology		
ab	Sharp		
street address	77 Mass Ave, 76-461		
City	Cambridge		
state/province	MA		
IP/Postal code	02139		

Data is submitted to GEO as a Series, which represents the experiment design

Platforms (1)	GPL13112 Illumina HiSeq 2000 (Mus musculus)						
Samples (6)         GSM1373652         DcrHet_1_mRN/           ■ Less         GSM1373653         DcrHet_2_mRN/           GSM1373654         DcrKO_1_mRN/           GSM1373655         DcrKO_2_mRN/		<sup>IA</sup> In <sup>A</sup> sa	Individual samples		submissions need to be associated with		
	GSM1373656 DcrHet_1_miR GSM1373657 DcrKO_1_miR	IN	a serie	es	a platio	Jimme.	
Relations							
BioProject	PRJNA245291						
SRA	SRP041414				I.		
Download fami	ly			Format		Click?	
SOFT formatted	family file(s)			SOFT 🛛		One by	
Series Matrix File	d family file(s) e(s)					one	
Su	pplementary file	Size	Download	File typ	e/resource		
GSE57043_dicer	ko_fpkm.txt.gz	875.5 Kb	(ftp)(http)	TXT		Normalized	
GSE57043_dicer	ko_hairpin_rpm.txt.gz	4.1 Kb	(ftp)(http)	TXT		counts	
GSE57043 dicer	ko_mature_rpm_txt.gz	1.6 Kb	(ftp)(http)	ТХТ			
SRP/SRP041/SR	2041414		(ftp)	SRA Stud	у	Daw Doado	
Raw data provide	ed as supplementary file					naw neaus	

Processed data is available on Series record

#### An expression signature for p53 status in human breast cancer predicts mutation status, transcriptional effects, and patient survival

Lance D. Miller\*<sup>+</sup>, Johanna Smeds<sup>‡</sup>, Joshy George\*, Vinsensius B. Vega\*, Liza Vergara\*, Alexander Ploner<sup>§</sup>, Yudi Pawitan<sup>§</sup>, Per Hall<sup>§</sup>, Sigrid Klaar<sup>‡</sup>, Edison T. Liu\*<sup>+</sup>, and Jonas Bergh<sup>‡</sup>

SANC

\*Genome Institute of Singapore, 60 Biopolis Street, #02-01, Singapore 138672; <sup>4</sup>Department of Oncology and Pathology, Radiumhemmet, Karolinska Institute and Hospital, S-17176 Stockholm, Sweden; and <sup>5</sup>Department of Medical Epidemiology and Biostatistics, Karolinska Institutet, 17177 Stockholm, Sweden

**GENETICS.** For the article "An expression signature for p53 status in human breast cancer predicts mutation status, transcriptional effects, and patient survival," by Lance D. Miller, Johanna Smeds, Joshy George, Vinsensius B. Vega, Liza Vergara, Alexander Ploner, Yudi Pawitan, Per Hall, Sigrid Klaar, Edison T. Liu, and Jonas Bergh, which appeared in issue 38, September 20, 2005, of *Proc. Natl. Acad. Sci. USA* (**102**, 13550-13555; first published September 2, 2005; 10.1073/pnas.0506230102), the breast cancer microarray data discussed in this publication have been deposited in the National Center for Biotechnology Information's Gene Expression Omnibus database (GEO, www.ncbi.nlm.nih.gov/geo/) and are accessible through GEO Series accession no GSE3494 [NCBI GEO].

### Conclusion

- GEO is a public repository for the archiving and distribution of gene expression data
- Best resource to get microarray and Next Generation Sequencing (RNASeq) data

#### Introduction

Informatics in health care may be called as health informatics

 It deals with the resources, devices, and methods required to optimize the acquisition, storage, retrieval, and use of information in health and biomedicine.

(wiki)

#### Introduction

Medical databases store and provide medical information

 The premier database for biomedical literature is the National Library of Medicine (NLM)'s MEDLINE, accessible through **PubMed**

## PUBMED

- Comprises of more than 24 million citations for biomedical literature from MEDLINE, life science journals, and online books
- Citations may include links to full-text content from PubMed Central and publisher web sites

S NCBI Resources 🖂 H	How To 🕑		Sign in to NCB
Publiced.gov US National Library of Medicine National Institutes of Health	PubMed   Advan	ced	Search
	PubMed PubMed comprise literature from MEL Citations may inclu publisher web sites	s more than 24 million citations for biomedical DLINE, life science journals, and online books. ude links to full-text content from PubMed Central and s.	PubMed Commons Featured comment - Dec 26, 2014 Regulating ribosome recruitment? I Shatsky critiques proposed RNA regulon mechanism. <u>1.usa.gov/1zfZekD</u>
Using PubMed		PubMed Tools	More Resources
PubMed Quick Start Guide		PubMed Mobile	MeSH Database
Full Text Articles		Single Citation Matcher	Journals in NCBI Databases
PubMed FAQs		Batch Citation Matcher	Clinical Trials
PubMed Tutorials		Clinical Queries	E-Utilities (API)
New and Noteworthy		Topic-Specific Queries	LinkOut

1

## MEDLINE

- MEDLINE is the primary resource for biomedical journal articles
- Millions of citations to articles in biomedical journals
- MEDLINE uses the MeSH vocabulary

#### **Other Databases**

MEDLINE is the primary resource, but other databases may also be helpful

- Academic OneFile
- CINAHL (Cumulated Index of Nursing and Allied Health Literature)
- PsycINFO
- Web of Knowledge

## Academic OneFile

- Academic OneFile lists articles from journals covering a broad range of subjects
- While it does not primarily focus on medical topics, useful articles can still be found here



## **PsycINFO**

- PsycINFO searches the psychological literature
- While it does not primarily focus on medical topics, useful articles can still be found here http://www.apa.org/pubs/databases

/psycinfo/coverage.aspx

More APA W	ebsites 👻   Home   Help   Log In   🏣 Cart (0)
Merican Psychological Association	SEARCH Q Entire Site
About APA Topics Publications & Databases Psychology Help Center News & Events Research	h Education Careers Membership
Home // Publications & Databases // APA Databases // PsycINFO® Homepage // PsycINFO® Journal Coverage List	
Publications: Books Children's Books Databases Journals Magazines & Newsletters Reports & Broch	ures Software Videos Merchandise
PsycINFO <sup>®</sup> Journal Coverage List	More About the Database
October 2014 Update	PsycINFO <sup>®</sup> Home
Currently, there are 2,562 journals covered in the PsycINFO <sup>®</sup> database. The list changes continuously as journals are added ar discontinued throughout the year, so it is updated online monthly.	FAQs     Coverage List
<ul> <li>Download the journal coverage list in Excel format (2.358KB)</li> </ul>	<ul> <li>Sample Records</li> </ul>
<ul> <li>View a list of the currently covered neuroscience titles (PDF, 625KB)</li> </ul>	Publisher Relations
Updated February 2010	Cited Reference Facts
View journal coverage facts	<ul> <li>Subsets &amp; Special Collections</li> </ul>
<ul> <li>View journal coverage policy</li> </ul>	Archive of Sample Searches Podcasts
View journals added in 2013	
	Librarians
Other Information	<ul> <li>Institutional Access</li> </ul>
Visit the journals added page to see a list of the journals we've added to the Journal Coverage List since the last update.	Institutional Pricing
With the exception of journals indexed cover-to-cover, not all articles from each journal are included in the database. PsycINFO staff examine each article and select only those that have psychological relevance.	Free Trial for Institutions

http://www.apa.org/pubs/databases/psycinfo/coverage.aspx

## Web of Science

- Major source for articles in a wide range of fields, including the sciences, social sciences, and humanities.
- Excellent place to find articles from scientific journals that may not be included in MEDLINE

### Conclusions

Informatics in health care may be called as health informatics

 Medical databases deal with the acquisition, storage, retrieval, and use of information in health and biomedicine.

#### Introduction

- Sequences are submitted to the databases in order to share them with the scientific community
- Generally sequences are submitted at the time of publication and are reviewed by peers

#### Caution

- It is important to ensure that sequence files do not contain any special characters
- Control characters in addition to standard ASCII characters must be removed else they might mess up the analysis or data transfer

Table 2.1.	Base–nucleic acid code	es
Symbol	Meaning	Explanation
G	G	Guanine
Α	Α	Adenine
Т	Т	Thymine
С	С	Cytosine
R	A or G	puRine
Y	C or T	pYrimidine
М	A or C	aMino
K	G or T	Keto
S	C or G	Strong interactions
		3 h bonds
W	A or T	Weak interactions
		2 h bonds
Н	A, C or T	H follows G in
	not G	alphabet
В	C, G or T	B follows A in
	not A	alphabet
v	A, C or G	V follows U in
	not T (not U)	alphabet
D	A, G or T	D follows C in
	not C	alphabet
N	A,C,G or T	Anybase
A damtad f	NC III (1094)	

committee of International Union of Biochemistry (IUB) has established standard codes to represent ambiguous bases or aminoacids

Nomenclature

Adapted from NC-IUB (1984).

Mount, pg 28

1-letter code	3-letter code	Amino acid	
A <sup>a</sup>	Ala	alanine	
C	Cys	cysteine	
D	Asp	aspartic acid	
E	Glu	glutamic acid	
F	Phe	phenylalanine	
G	Gly	glycine	
н	His	histidine	
I	Ile	isoleucine	
K	Lys	lysine	
L	Leu	leucine	00
M	Met	methionine	
N	Asn	asparagine	8
P	Pro	proline	
Q	Gln	glutamine	Ľ.
R	Arg	arginine	2
S	Ser	serine	2
Т	Thr	threonine	Ĕ
v	Val	valine	2
w	Trp	tryptophan	
x	Xxx	undetermined amino acid	
Y	Tyr	tyrosine	
Zb	Glx	either glutamic acid or glutamine	

Table 2.2. Table of standard amino acid code letters

Adapted from IUPAC-IUB (1969, 1972, 1983).

<sup>a</sup> Letters not shown are not commonly used.

<sup>b</sup> Note that sometimes when computer programs translate DNA sequences, they will put a "Z" at the end to indicate the termination codon. This character should be deleted from the sequence.

#### NCBI

NCBI has two options

## BANKIt

- For simple
  - sequences and annotations
- For submission through web
- Do not require advanced tools

### NCBI

- Sequin
  - For Complex
     sequences and
     annotations
  - For off-line submissions
  - Do require
     advanced tools
     and graphical
     reports

# SNCBI Banklt

#### http://www.ncbi.nlm.nih.gov/WebSub/?tool =genbank

S NCBI	Se an	SequinA DNA Sequence Submission and Update Tool						
Sequin	Entrez	BLAST	OMIM	Taxonomy	Structure			
Sequin home	Sequin 1	3.05 is now ava	ailable.					

#### UniProt

**SPIN** web-based tool for submitting

- directly sequenced protein sequences
- biological annotations to the knowledgebase





https://www.ebi.ac.uk/swissprot/Submissions/spin

### Conclusion

 Sequences are stored in databases in specific format



#### Introduction

- Databases not merely collect and organize data but allow intelligent data retrieval
- And/or analysis



S	ave search Advanced			
Display Settings	💼 🖂 Tabular, 20 per page,	Sorted by Relevance	<u>s</u>	end to:
Did you mean Search Gene	n p53 as a gene symbol? for <u>p53</u> as a symbol.			
Results: 1 to	20 of 9031	<< First	< Prev Page 1 of 452 Next	> Last
Tillers activat	ted: Current only. <u>Clear all</u>	to show 9271 items.		
Name/Gene ID	Description	Location	Aliases	MIM
Name/Gene ID p53 ID: 2768677	Description CG33336 gene product from transcript CG33336-RB [ <i>Drosophila</i> <i>melanogaster</i> (fruit fly)]	Location Chromosome 3R, NT_033777.3 (2304965723054082, complement)	Aliases Dmel_CG33336, CG10873, CG31325, CG33336, D- DMP53, Dm-P53, DmP53, Dmel\CG33336, Dmp53, Dp53, dmp53, dp53, prac	MIM

http://www.ncbi.nlm.nih.gov/

#### p53 [Drosophila melanogaster (fruit fly)]

Gene ID: 2768677, updated on 4-Jan-2015

#### Summary

Official Symbol	p53 provided by <u>FlyBase</u>
Primary source	FLYBASE:FBgn0039044
Locus tag	Dmel_CG33336
Gene type	protein coding
RefSeq status	REVIEWED
Organism	Drosophila melanogaster (old-lineage: Eukaryota; Metazoa; Arthropoda; Hexapoda; Insecta; Pterygota; Neoptera;
	Endopterygota; Diptera; Brachycera; Muscomorpha; Ephydroidea; Drosophilidae; Drosophila; Sophophora)
Lineage	Eukaryota; Metazoa; Ecdysozoa; Arthropoda; Hexapoda; Insecta; Pterygota; Neoptera; Endopterygota; Diptera;
	Brachycera; Muscomorpha; Ephydroidea; Drosophilidae; Drosophila; Sophophora
Also known as	CG10873; CG31325; CG33336; D-p53; Dm-P53; Dmel\CG33336; dmp53; Dmp53; DmP53; DMP53; dp53; Dp53; prac

#### http://www.ncbi.nlm.nih.gov/

☆?

Genomic con	text			
Location: 94D1 Exon count: 10	0-94D10			See p53 in Epigenomics, MapViewer
Annotation release	Status	Assembly	Chr	Location
Release 6.01	current	Release 6 plus ISO1 MT (GCF_000001215.4)	3R	NT_033777.3 (2304965723054082, complement)
Release 5.57	previous assembly	Release 5 (GCF_000001215.2)	3R	NT_033777.2 (1887537918879804, complement)



http://www.ncbi.nlm.nih.gov/



#### http://www.ncbi.nlm.nih.gov/

#### Drosophila melanogaster chromosome 3R

NCBI Reference Sequence: NT\_033777.3

FASTA Graphics

LOCUS	NT_033777	4426 bp	DNA	linear	INV	05-AUG-2014	
DEFINITION	Drosophila melanogaster	chromosome	3R.				
ACCESSION	NT_033777 REGION: comple	ement(23049	657230	54082)			
VERSION	NT_033777.3 GI:67116212	22					
DBLINK	BioProject: PRJNA164						
	BioSample: SAMN02803731						
KEYWORDS	RefSeq.						
SOURCE	Drosophila melanogaster	(fruit fly	)				
ORGANISM	Drosophila melanogaster						
	Eukaryota; Metazoa; Ecd	ysozoa; Artl	hropoda;	Hexapoda	; Ins	secta;	
	Pterygota; Neoptera; End	dopterygota	; Diptera	a; Brachy	cera	;	
	Muscomorpha; Ephydroidea	a; Drosophi	lidae; D	rosophila	; Sop	phophora.	_
REFERENCE	1 (bases 1 to 4426)						
AUTHORS	Hoskins, R.A., Carlson, J	.W., Kenned	y,C., Ac	evedo,D.,	Evar	ns-Holm,M.,	
	Frise,E., Wan,K.H., Parl	k,S., Mende	z-Lago,M	., Rossi,	F.,		
	Villasante, A., Dimitri,	P., Karpen,	G.H. and	Celniker	,S.E.		
TITLE	Sequence finishing and m	mapping of 1	Drosophi	la melano	gaste	er	
	heterochromatin						
JOURNAL	Science 316 (5831), 162	5-1628 (200	7)				
PUBMED	17569867						
			htt	p://ww	w.no	bi.nlm.nih.g	;ov/

Location/Oualifiers FEATURES 1..4426 source /organism="Drosophila melanogaster" /mol type="genomic DNA" /db xref="taxon:7227" /chromosome="3R" /genotype="y[1]; Gr22b[1] Gr22d[1] cn[1] CG33964[R4.2] bw[1] sp[1]; LysC[1] MstProx[1] GstD5[1] Rh6[1]" 1..4426 gene /gene="p53" /locus tag="Dmel CG33336" /gene\_synonym="CG10873; CG31325; CG33336; D-p53; Dm-P53; Dmel\CG33336; dmp53; Dmp53; DmP53; DMP53; dp53; Dp53; prac" /map="94D10-94D10" /db xref="FLYBASE:FBgn0039044" /db xref="GeneID:2768677" mRNA join(1..118,178..501,884..964,1035..1071,1135..1161, 2959...3268,3333...3579,3642...4036,4096...4426) /gene="p53" /locus\_tag="Dmel CG33336" /gene synonym="CG10873; CG31325; CG33336; D-p53; Dm-P53; Dmel\CG33336; dmp53; Dmp53; DmP53; DMP53; dp53; Dp53; prac" /product="p53, transcript variant B" /note="p53-RB; Dmel\p53-RB; CG33336-RB; Dmel\CG33336-RB" /transcript id="NM 206544.2" /db xref="GI:281362333" /db xref="FLYBASE:FBtr0084360" /db xref="FLYBASE:FBgn0039044" /db xref="GeneID:2768677" http://www.ncbi.nlm.nih.gov/

CDS

join(75..118,178..501,884..964,1035..1071,1135..1161, 2959...3268,3333...3579,3642...4036,4096...4118) /gene="p53" /locus tag="Dmel CG33336" /gene synonym="CG10873; CG31325; CG33336; D-p53; Dm-P53; Dmel\CG33336; dmp53; Dmp53; DmP53; DMP53; dp53; Dp53; prac" /note="CG33336 gene product from transcript CG33336-RB; CG33336-PB; p53-PB; p53-like regulator of apoptosis and cell cycle; Dmp53; protein 53; drosophila p53" /codon start=1 /product="p53, isoform B" /protein id="NP 996267.1" /db xref="GI:45553461" /db xref="FLYBASE:FBpp0083753" /db xref="FLYBASE:FBqn0039044" /db xref="GeneID:2768677" /translation="MSLHKSASFSLTFNQNTSIVSRSNSRTIFEAFKEFLDFWDIGNE VSAESAVRVSSNGAFNLPQSFGNESNEYAHLATPVDPAYGGNNTNNMMQFTNNLEILA NNNSDGNNKINACNKFVCHKGTDSEDDSTEVDIKEDIPKTVEVSGSELTTEPMAFLOG LNSGNLMOFSOOSVLREMMLODIQIOANTLPKLENHNIGGYCFSMVLDEPPKSLWMYS IPLNKLYIRMNKAFNVDVOFKSKMPIOPLNLRVFLCFSNDVSAPVVRCONHLSVEPLT ANNAKMRESLLRSENPNSVYCGNAOGKGISERFSVVVPLNMSRSVTRSGLTROTLAFK FVCQNSCIGRKETSLVFCLEKACGDIVGQHVIHVKICTCPKRDRIQDERQLNSKKRKS VPEAAEEDEPSKVRRCIAIKTEDTESNDSRDCDDSAAEWNVSRTPDGDYRLAITCPNK EWLLOSIEGMIKEAAAEVLRNPNOENLRRHANKLLSLKKRAYELP"

#### http://www.ncbi.nlm.nih.gov/

#### ORIGIN

1	cctggagcac	ggaagattct	tgcggacaca	aatcgcaact	gctaaataaa	atttatttat
61	ttgagtgcac	agccatgagt	cttcacaagt	ccgcgtcgtt	tagcttgact	tttaaccagt
121	gagcggagat	attttattcg	gtcttaccca	acaaaataat	gttgcgcctt	tttgcagaaa
181	cacttcgatt	gtttcgcgta	gcaatagtcg	cacaatttt	gaagctttca	aggagttcct
241	ggatttttgg	gatatcggca	acgaagtttc	tgcagagtca	gcagttcggg	tctccagcaa
301	cggagctttc	aacttgccgc	agagttttgg	caacgaatcc	aacgaatatg	cccacctggc
361	tacgcctgtg	gatccagcct	acggaggcaa	caacacgaac	aacatgatgc	agttcacgaa
421	caatctggaa	attttggcca	acaataattc	cgatggcaat	aacaaaatta	atgcatgcaa
481	caaattcgtc	tgccacaagg	ggtgagcaaa	ttcaaaacac	gcgctccaat	cgataaacat
541	tggctacggc	gattgttcgc	gctgcgtggc	gaatggcaaa	atccaaatag	tcggtggcca
601	ctacgattct	gtagttttt	gttagcgaat	ttttaatatt	tagcctcctt	ccccaacaag
661	atcgcttgat	cagatatagc	cgactaagat	gtatatatca	cagccaatgt	cgtggcacaa
721	agaaaggtac	agtgcggcaa	caaattgatg	atcgaacagt	agaaaccttg	catgtagcaa

4261 ggcatgttcg atggccgaaa agaaaacatt tttatatttt tgatagtata ctgttgttaa 4321 ctgcāgttct atgtgactac gtaacttttg tctaccacaa caaacatact ctgtacaaaa 4381 aagccaaaag tgaatttatt aaagagttgt catattttgc aaacat

11

http://www.ncbi.nlm.nih.gov/
## **DNA Sequence Retrieval**

#### Conclusions

- DNA Sequences are stored in DNA sequence databases in specified formats
- Genebank format is a standard format



#### FASTA Sequence Format

- DNA sequences are stored in specified formats
- Different softwares need sequences in different formats

#### FASTA Sequence Format

- FASTA is the most frequently used format to present DNA and Protein sequences
- It is recommended that all lines of text be shorter than 80 characters in length

>gi[120407068]ref[NP\_000537.3] cellular tumor antigen p53 isoform a [Homo sapiens] MEEPQSDPSVEPPLSQETFSDLWKLLPENNVLS PLAPPVLGFLHSGTAKSVTCTYSPALNKMFCQLA KT--\*

Sometimes a \* might be present in the end

#### **Genebank Format**

- A sequence file in GenBank format can contain several sequences
- One sequence starts with a line containing the word LOCUS and a number of annotation lines

#### **Genebank Format**

 The start of the sequence is marked by a line containing "ORIGIN" and the end of the sequence is marked by two slashes ("//")

#### **Genebank Format**

LOCUS AAU03518 237 bp DNA PRI 04-FEB-1995 DEFINITION Aspergillus awamori internal transcribed spacer 1 (ITS1) and 18S rRNA and 5.8S rRNA genes, partial sequence. ACCESSION U03518 BASE COUNT 41 a 77 c 67 g 52 t ORIGIN 1 aacctgcgga aggatcatta ccgagtgcgg gtcctttggg cccaacctcc catccgtgtc 61 tattgtaccc tgttgcttcg gcgggcccgc cgcttgtcgg ccgccggggg ggcgcctctg 121 ccccccgggc ccgtgcccgc cggagacccc aacacgaaca ctgtctgaaa gcgtgcagtc 181 tgagttgatt gaatgcaatc agttaaaact ttcaacaatg gatctcttgg ttccggc

#### **EMBLFormat**

• An example sequence in EMBL format is:

ID AA03518 standard; DNA; FUN; 237 BP.

XX

AC U03518;

XX

DE Aspergillus awamori internal transcribed spacer 1 (ITS1) and 18S

DE rRNA and 5.8S rRNA genes, partial sequence.

XX

SQ Sequence 237 BP; 41 A; 77 C; 67 G; 52 T; 0 other;

aacctgcgga aggatcatta ccgagtgcgg gtcctttggg cccaacctcc catccgtgtc60tattgtaccc tgttgcttcg gcgggcccgc cgcttgtcgg ccgccgggg ggcgcctctg120ccccccgggc ccgtgcccgc cggagacccc aacacgaaca ctgtctgaaa gcgtgcagtc180tgagttgatt gaatgcaatc agttaaaact ttcaacaatg gatctcttgg ttccggc237

 $\parallel$ 

#### **SwissProt Format**

- SwissProt protein sequence format is similar to EMBL format
- There is considerably more information about physical and biochemical properties of a protein is provided

#### **SwissProt Format**

- ID Identification.
- AC Accession number(s).
- DT Date.
- DE Description.
- GN Gene name(s).
- OS Organism species.
- OG Organelle.
- OC Organism classification.
- RN Reference number.

- RP Reference position.
- RC Reference comments.
- RX Reference cross-references.
- RA Reference authors.
- RL Reference location.
- CC Comments or notes.
- DR Database cross-references.
- KW Keywords.
- FT Feature table data.
- SQ Sequence header.
  - // Termination line.

#### **XML Format**

- Extensible Markup Language
- Readable by both man and machine
- Becoming standard data format for transferring genome data

#### XML Format

<xsd:annotation> <xsd:documentation> XML Schema for SBOL core data model compatible with RDF/XML serialization. <dc:date>2012-01-19</dc:date> <dc:creator>Evren Sirin</dc:creator> <dc:contributor>Michal Galdzicki</dc:contributor> </xsd:documentation> </xsd:annotation>

#### **Sequence converters**

- READSEQ is a useful sequence converter developed by D.G.Gilbert at Indiana University, USA
- Recognizes DNA or Protein sequence file and interconvert them

#### Conclusions

- Databases store sequences in specified formats
- Genebank, DDBJ and EMBL has similar formats
- Different softwares need sequences in different formats

#### **Data Retrieval**

- Nearly all biological databases are available for download as simple text (flat) files
- A local version of the database allows one greater freedom in processing the data

#### Entrez

- is an integrated search engine which allows users to search and retrieve different data from the NCBI
- It can be accessed from the site
   www.ncbi.nlm.nih.gov/E
   ntrez/

#### Entrez

- integrates PubMed and 39 other scientific
   literatures, nucleotide
   and protein databases
- protein domain data, population studies, expression data, pathways, genome details and taxonomic information



#### SNCBI Resources 🖸 How To 🖸

S NCBI Resources How To Sign i			
Search NCBI database	9S		Help
Literature Books MeSH NLM Catalog	books and reports ontology used for PubMed indexing books, journals and more in the NLM Collections	Genes EST Gene GEO DataSets	expressed sequence tag sequences collected information about gene loci functional genomics studies
PubMed PubMed Central	scientific & medical abstracts/citations full-text journal articles	GEO Profiles HomoloGene	gene expression and molecular abundance profiles homologous gene sets for selected organisms sequence sets from phylogenetic and population
Health ClinVar dbGaP	human variations of clinical significance genotype/phenotype interaction studies genetic testing registry medical genetics literature and links online mendelian inheritance in man clinical effectiveness, disease and drug reports	UniGene Proteins	studies clusters of expressed transcripts
GTR MedGen OMIM PubMed Health Genomes		Conserved Domains Protein Protein Clusters Structure	conserved protein domains protein sequences sequence similarity-based protein clusters experimentally-determined biomolecular structures
Assembly BioProject BioSample	genomic assembly information biological projects providing data to NCBI descriptions of biological source materials	Chemicals  BioSystems  PubChem BioAssay	molecular pathways with links to genes, proteins and chemicals bioactivity screening studies

0

#### **Bulk Data Retrieval**

- The best option is to use ftp (File transfer protocol)
- The File Transfer Protocol (FTP) is a standard network protocol used to transfer files
- Via command line or application programs like FTP clients

#### **Bulk Data Retrieval**

- Data needs to be transformed or processed using programming languages
- PERL and Python are good for processing Biological data

#### Conclusions

- Data is transferred over the internet
- Data needs to be transformed or processed





## why Proteomics

Lecture-2

## What do mean by proteomics?

- Proteomics is the large-scale study of proteins, usually by biochemical methods.
   The word proteomics has been associated
- traditionally with displaying a large number of proteins from a given
- cell line or organism on two -dimensional polyacrylamide gels









## **Scope of Proteomics**

The identification, characterization and quantification of all proteins involved in a particular pathway, organelle, cell, tissue, organ or organism that can be studied in concert to provide accurate and comprehensive data about that system.

**Or - A complete description of proteins expressed in any given cell at any given time** 



## Why should we study Proteomics?

- directly contributes to drug development
- Verification of a gene product by proteomic methods
- Modifications of the proteins
- Protein expression level d



### WHY PROTEOMICS?

Many types of information cannot be obtained from the study of genes alone. For example, proteins, not genes, are responsible for the phenotypes of cells. It is impossible to elucidate mechanisms of disease, aging, and effects of the environment solely by studying the genome.



## Advantages of study of proteomics

- Shows that genetic alterations are not the reason for all types of diseases
- Helps in determining the proper treatment of diseases
- With the help of three dimensional analysis of proteins we have found that HIV protease is the enzyme which is responsible for AIDS.
- One of the most important use of proteomics in diagnosis is the identification of
- biomarkers. The study of drugs in proteomics

is called pharmacoproteomics.

#### **Proteomics aims**

- Genomics integrated strategies
- Study of post-translational modifications
- Identification of novel protein targets for drugs
- Analysis of tumor tissues
- Comparison between normal and diseased tissues
- Comparison between diseased and pharmacologically treated

tissues









## Limitations of Genomics Challenge of Proteomics

- co-translational modifications
  - differential mRNA splicing
- post-translational modifications (PTMs)
  - C-terminal GPI anchor
  - phosphorylation
  - sulfation
  - glycosylation
  - N-myristoylation
  - hydroxylation
  - N-methylation
  - carboxymethylation
  - signal peptidase site......





# Proteomics Introducción BIO 601

Lecture-1

## The Central Dogma



The central dogma states that information in nucleic acid can be perpetuated or transferred, but the transfer of information form into protein is irreversible.

## **Central Dogma of Biology**

- DNA -> RNA -> Protein Synthesis
- Transcription:
  - Process of DNA serving as a template for RNA synthesis
- Translation:
  - Process of RNA serving as a template for protein synthesis



## How do DNA and genes relate to proteins?

- DNA provides the genes, or genetic code, for protein synthesis
- Genes are
  expressed because
  DNA codes for RNA
  which then codes
  for ALL of our
  proteins




## **Background: 3 Types of RNA**

- mRNA: Messenger RNA
  - 1st RNA's made DIRECTLY from DNA template
  - Travel from nucleus to ribosome
- rRNA: Ribosomal RNA
  - Helps form ribosomes in cytoplasm
- tRNA: Transfer RNA
  - Brings amino acids from cytoplasm to ribosome so proteins can be made



## **Step 1: Transcription**

- INSIDE of the nucleus DNA is used to make mRNA
- DNA is unzipped then RNA polymerase makes an mRNA strand from the DNA template
- New mRNA strand then leaves the nucleus and travels into the cytoplasm
- DNA is ALWAYS left protected in the nucleus



## **Step 1: Transcription**

- DNA: 5' AAA TTT GGG CCC ATC GCA 3'
- mRNA: 3' UUU AAA CCC GGG UAG CGU 5'
- DNA: CTA GTT CCC TAA AAG GAG
- mRNA: GAU CAA GGG AUU UUC CUC
- DNA: TAC CGA GGT TTA ACT
- mRNA: AUG UGA CCA AAU UGA



## **Step 2: Translation**

- Each nucleotide
  sequence serves as
  a code for what
  amino acid will be
  added to the
  protein being made
- Nucleotides read in triplets, or codons





### **Step 2: Translation**

	U	С	A	G	
U	Phe Phe Leu Leu	Ser Ser Ser Ser	Tyr Tyr STOP STOP	Cys Cys <mark>STOP</mark> Trp	UCAG
с	Leu Leu Leu Leu	Pro Pro Pro Pro	His His GIn GIn	Arg Arg Arg Arg	UCAG
A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	UCAG
G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	UCAG

Third base in codon

#### Second base in codon

First base in codon



## **Step 2: Translation**

- mRNA is now connected to the ribosome
- tRNA has a corresponding anti-codon and brings over the

corresponding amino acid

toot/com



## The end result...

- An amino acid sequence that makes a protein
- GENES code for proteins/enzymes
- We NEED proteins to function
- The shape of the protein determines its function





#### **TYPES OF PROTEOMICS**

Lecture-3

### **Scope of Proteomics**

- Expression proteomics
- Structural proteomics
- Functional proteomics



## **EXPRESSION PROTEOMICS**

- Expression proteomics is used to study the qualitative and quantitative expression of total proteins under two different conditions.
  - □ Normal and diseased state.
  - $\Box$  E.g. :tumor or normal cell.
  - □ It studied that protein is over expressed or under expressed.
  - □ 2-D electrophorasis.



#### STRUCTURAL PROTEOMICS

- Structural proteomics helps to understand three dimensional shape and structural complexities of functional proteins
  - functional proteins.
  - □ It determine either by amino acid sequence in protei
  - from a gene this process is known as **homology modeling**.
  - $\hfill\square$  It identify all the protein present in complex system protein-protein interaction.
  - □ Mass spectroscopy is used for structure determination



# **FUNCTIONAL PROTEOMICS**

□ Functional proteomics explains understanding the protein functions well as unrevealing molecular mechanisms within the cell that depend on identification of the interacting protein partners.

 $\Box$  So that detailed description of the cellular signaling pathways migreatly benefit from the elucidation of protein- protein interaction



## Limitations of Genomics Challenge of Proteomics

- co-translational modifications
  - differential mRNA splicing
- post-translational modifications (PTMs)
  - C-terminal GPI anchor
  - phosphorylation
  - sulfation
  - glycosylation
  - N-myristoylation
  - hydroxylation
  - N-methylation
  - carboxymethylation
  - signal peptidase site......



### **Structural Proteomics**

Lecture-1

#### What is structural proteomics/genomics?

- High-throughput determination of the 3D structure of proteins
- Goal: to be able to determine or predict the structure of every protein.
  - Direct determination X-ray crystallography and nuclear magnetic resonance (NMR).
  - Prediction
    - Comparative modeling -
    - Threading/Fold recognition
    - Ab initio



#### What is structural proteomics/genomics?

- High-throughput determination of the 3D structure of proteins
- Goal: to be able to determine or predict the structure of every protein.
  - Direct determination X-ray crystallography and nuclear magentic resonance (NMR).
  - Prediction
    - Comparative modeling -
    - Threading/Fold recognition
    - Ab initio



## Why structural proteomics?

- To study proteins in their active conformation.
  - Study protein: drug interactions
  - Protein engineering
- Proteins that show little or no similarity at the primary sequence level can have strikingly similar structures.



#### An example

- FtsZ protein required for cell division in prokaryotes, mitochondria, and chloroplasts.
- Tubulin structural component of microtubules important for intracellular trafficking and cell division.
- FtsZ and Tubulin have limited sequence similarity and would not be identified as homologous proteins by sequence analysis.







FtsZ and tubulin have little similarity at the amino acid sequence level

Burns, R., Nature **391**:121-123 Picture from E. Nogales



## Are FtsZ and tubulin homologous?

 Yes! Proteins that have conserved secondary structure can be derived from a common ancestor even if the primary sequence has diverged to the point that no similarity is detected.



#### Current state of structural proteomics

- As of Feb. 2002 16,500 structures
   Only 1600 non-redundant structures
- To identify all possible folds predicted another 16,000 novel sequences needed for 90% coverage.
  - Of the 2300 structures deposited in 2000, only 11% contained previously unidentified folds.
- Overall goal directly solve enough structures directly to be able to computationally model all future proteins.



#### Protein domains - structure

- "clearly recognizable portion of a protein that folds into a defined structure"
  - Doesn't have to be the same as the domains we have been investigating with CDD.
  - RbsB proteins as an example.



### Main secondary structure elements

 $\alpha$ -helix - right handed helical structure

β-sheet - composed of two or more βstrands, conformation is more "zig-zag" than helical.



## Folds/motifs

- How these secondary structure elements come together to form structure.
   – Helix-turn-helix
- Determining the structure of nearly all folds is the goal of structural biology



## X-Ray Crystallography

- Make crystals of your protein
  - 0.3-1.0mm in size
  - Proteins must be in an ordered, repeating pattern.
- X-ray beam is aimed at crystal and data is collected.
- Structure is determined from the diffraction data.



### Origin of proteomics

Lecture-1

- In 1975, the introduction of the 2D gel by O'Farrell who began mapping proteins from *E. coli.*
- Although many proteins could be separated and visualized, they could not be identified.



## Human protein index

- Despite these limitations, shortly thereafter a large-scale analysis of all human proteins was proposed.
- The goal of this project, termed the human protein index, was to use two-dimensional protein electrophoresis (2-DE) and other methods to catalog all human proteins.
- However, lack of funding and technical limitations prevented this project from continuing.

- The first major technology to emerge for the identification of proteins was the sequencing of proteins by Edman degradation.
- A major breakthrough was the development of microsequencing techniques for electroblotted proteins.



- Microsequencing technique was used for the identification of proteins from 2-D gels to create the first 2-D databases.
- Improvements in microsequencing technology resulted in increased sensitivity of Edman sequencing in the 1990s to high picomole amounts.



- One of the most important developments in protein identification has been the development of MS technology.
- The sensitivity of analysis and accuracy of results for protein identification by MS have increased by several orders of magnitude.
- It is now estimated that proteins in the femtomolar range can be identified in gels.
- Because MS is more sensitive, can tolerate protein mixtures, and is amenable to highthroughput operations

#### Genomics vs. proteomics

Lecture-1

## **Genomics vs. proteomics**

- Genomics and proteomics are closelyrelated fields.
- The main difference between genomics and proteomics is that genomics is the study of the entire set of genes in the genome of a cell whereas proteomics is the study of the entire set of proteins produced by the cell.



## Nature of study material

**Genomics:** The genome is constant. Every cell of an organism has the same set of genes.

**Proteomics:** The proteome is dynamic and varies. The set of proteins produced in different tissues varies according to the gene expression.



### Use of High throughput techniques

**Genomics:** High throughput techniques are used in the genomics to map, sequence, and analyze genomes.

**Proteomics:** In proteomics, characterization of the 3D structure and the function of proteins is carried out by the use of high throughput methods.



## **Techniques** involved

**Genomics:** The techniques involved in genomics include gene sequencing strategies such as directed gene sequencing, whole-genome shotgun sequencing, construction of expressed sequence tags (ESTs), identification of single nucleotide polymorphisms (SNPs)


# **Techniques** involved

## **Proteomics:**

- Extraction and electrophoretic separation of proteins.
- Digestion of proteins with the use of trypsin into small fragments.
- Determination of the amino acid sequence by mass spectrometry.
- Identification of proteins using the information in the protein databases.



# Importance

**Genomics:** Genomic studies are important to understand the structure, function, location, regulation of the genes of an organism.

**Proteomics:** The study of the entire set of proteins produced by a cell type is done in order to understand its structure and function.





#### Life and Death of a Protein

Proteins are synthesized by the translation of mRNAs into polypeptides on ribosomes.

In most cases, the initial polypeptide-translation product undergoes some type of modification before it assumes its functional role in a living system.

These changes are broadly termed "posttranslational modifications" and encompass a wide variety of reversible and irreversible chemical reactions.

Approximately 200 different types of posttranslational modifications have been reported. Some of these are summarized in Fig. 1, which depicts the life cycle of a prototypical protein.





# Life Cycle of the Cell





## **Modifications during Protein Cycle**

Modifications those occur early in the life of the protein

- Carboxylation of glutamate residues
- Removal of the N-terminal methionine
- Glycosylation
- Addition of Prosthetic groups
- Formation of multisubunit complexes
- Prenylation of cysteine residues assists anchoring of proteins in or on membranes.

These more or less "permanent" modifications and transport ultimately result in the delivery of functional proteins to specific locations in cells.





- The activities of many proteins are then controlled by posttranslational modifications.
- The most prominent and best-understood of these is phosphorylation of serine, threonine, or tyrosine residues.
- Phosphorylation may activate or inactivate enzymes, alter proteinprotein interactions and associations, change protein structures, and target proteins for degradation.
- Protein phosphorylation regulates protein function in diverse contexts and appears to be a key switch for rapid on-off control of signaling cascades, cell-cycle control, and other key cellular functions.



## **Degradation of Proteins**

- Protein modifications appear to be critical to initiating processes that ultimately degrade proteins.
- Phosphorylation of some proteins is rapidly followed by conjugation with ubiquitin, which leads to degradation by the 26S proteasomal complex.
- There evidently are other stimuli for protein ubiquitination and turnover, including oxidative damage and other protein modifications.
- Proteins also undergo degradation by lysosomal enzymes.
- Any protein may be present in many forms at any one time in a cell.
- Collectively, the proteome of a cell comprises all of these many forms of all expressed proteins. This certainly makes the proteome bewilderingly

# Proteins as Modular Stru BIO 601

Lecture-1

## **Proteins as Modular Structures**

- Segments of amino acid sequences can be considered as functional building blocks or modules.
- The modular units in proteins that confer specific properties and functions are referred to as "motifs" or "domains".
- Motifs and domains are recognizable sequences that confer similar properties or functions when they occur in a variety of proteins.
- In some cases, amino acid sequences within motifs and domains are highly conserved and do not vary from protein to protein.
- In other cases, some key amino acids occur in a reproducible relationship to each other in a sequence, even though various substitutions in other amino acids occur









#### **Proteins as Modular Structures**

Longer amino acid sequences often form domains, which confer specific properties or functions on a protein.

Some domain structures refer simply to sequences that confer a bulk physical property to a segment of the polypeptide, such as transmembrane domains, which simply form helices that span a lipid bilayer membrane.

Other domain structures provide hydrogen bonding or other contacts for key enzyme substrates or prosthetic groups.

In many cases, domains are made up of combinations of units of secondary structure, such as helix-loop-helix domains.





Lecture-1

# **Protein Localization**

- Any process in which a protein is transported to, or maintained in, a specific location.
- Cells are organized into many different compartments such as the cytosol, nucleus, endoplasmic reticulum (ER), and mitochondria. Almost all proteins are made in the cytosol, yet each cellular compartment requires a specific set of proteins.
- How does the cell regulate protein localization to be sure that proteins end up where they should?



# **Compartments of an Animal Cell**



Figure 12–1. Molecular Biology of the Cell, 4th Edition.



## **Functions of major intracellular compartments:**

- Nucleus contains main genome, DNA and RNA synthesis.
- Cytosol most protein synthesis, glycolysis and metabolic pathways synthesizing amino acids, nucleotides.
- Endoplasmic reticulum synthesis of membrane proteins, lipid synthesis.
- Golgi apparatus covalent modification of proteins from ER, sorting of proteins for transport to other parts of the cell.



# Cont...

- Mitochondria and chloroplasts (plants) ATP synthesis.
- Lysosomes degradation of defunct intracellular organelles and material taken in from the outside of the cell by endocytosis.
- Endosomes sorts proteins received from both the endocytic pathway and from the Golgi apparatus.
- Peroxisomes oxidize a variety of small molecules.



## **Three basic modes of protein Localization**

- 1. Gated transport
- 2. Transmembrane transport
- 3. Vesicular transpo

fppt/com



©1998 GARLAND PUBLISHING

# **Roadmap of protein traffic**



The genesis and function of internal compartments depends on the appropriate targeting of proteins

## Some of the green route illustrated.



Figure 12-6. Molecular Biology of the Cell, 4th Edition. Figure 12-5. Molecular Biology of the Cell, 4th Edition.



# Chemical Composition of **BIO 601**

Lecture-

- Proteins are polymers of amino acids.
- They range in size from small to very large.
- All the proteins are made up of Twenty different types of amino acids. So these amino acids are called standard amino acids.







## **Chemical composition of**

## proteins









• In a protein molecule, each amino acid residue is joined to its neighbour by a specific type of covalent bond which is called Peptide Bond.







 Amino acids can successively join to form dipeptides, tripeptides, tetrapeptides, oligo peptides and polypeptides.







## **Homology Modelling**

## Introduction to Homology Modelling

## Background

 Proteins are 3D molecules with their own unique structures

 Protein structure is reflective of the protein function

 Protein structure includes 1', 2', 3' and 4' structures

#### Background

1' structure of proteins is the sequence of proteins and can be obtained by mass spectrometry

 2' structures formed by proteins are the helices, beta sheets, loops and coils

## Background

- 3' structure of proteins is the combination of 2' structures such that the overall protein structure is formed
- 4' protein structure is formed when two or more proteins complex together

## Background

- X-Ray Crystallography and NMR Spectroscopy are used to find the structures of proteins
- However, these methods are difficult and expensive
- Solution: Prediction of structures

## Introduction

- Protein sequence gives rise to its structure
- If another protein which has a <u>similar</u> <u>sequence</u> also has its <u>structure</u> also has its <u>structure known</u>, the structure of an unknown protein can be predicted based on that similar protein

#### Introduction

So, it is then possible to identify unknown protein structures by just examining the homologous protein sequences

#### Conclusions

- Sequence Identity
- Alignment Length

Which combination of identity and alignment length is suitable for best for structure prediction?
# **Homology Modelling**

Homology, Paralogy and Orthology

# Background

- In homology modelling, proteins with similar 1' sequences are considered
- Given that one of them has its 3' structure known, then the 3' structure of other protein can be predicted



http://bioweb.uwlax.edu/GenWeb/Molecular/Bioinformatics/Unit\_4/Lab\_4-2/lab\_4-2.htm

#### How much homology is required or better?



http://www.cmbi.ru.nl/gvteach/astra/lectures/homology\_modelling.ppt

# Conclusions

- Good sequence alignment and identity ensures that homology modelling will give accurate results
- Next, what is the workflow for homology modelling?

# **Homology Modelling**

# Workflow of Structural Modelling

# Background

- Homology modelling
  is used to predict
  structures of
  proteins having high
  sequence similarity
  with other proteins
  with known
  structures!
- Let's consider the workflow of homology modelling



# Introduction

Overall, there are three different strategies for structure prediction

- 1. Homology Modelling
- 2. Threading/Fold Recognition
- 3. Ab Initio Modelling

# Conclusions

 Next, we will proceed to perform homology modelling

For that there is a seven step procedure which we will see in the next module.

# **Homology Modelling**

Seven Steps to Homology Modelling – I

# **Background**

Protein structure can be predicted by 3 methods:

1. Homology Modelling

2. Fold Recognition / Threading

3. Ab Initio Modelling

# Introduction

- Let's start by looking at Homology Modelling
- There are seven salient steps in any Homology Modelling pipeline

 Definition of Template (known) & Target (unknown) Homology modeling of the target structure can be done as follows:

- 1. Template recognition and initial alignment
- 2. Alignment correction
- 3. Backbone generation
- 4. Loop modeling
- 5. Side-chain modeling
- 6. Model optimization
- 7. Model validation



# Conclusions

- Homology modelling works in seven steps
- It is a repetitive process
- Next, we will look at each step in detail!

# **Homology Modelling**

Seven Steps to Homology Modelling – II

# Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation



Compare the sequence of the unknown protein with all the sequences of known structures stored in the Protein Data Bank (PDB).

BLAST this sequence against PDB sequences – Obtain a list of known protein structures that match the sequence.

BLAST uses a residue exchange scoring matrix. Residues that are easily exchanged (e.g. Ile to Leu) get a better score than residues that have different properties (for example Glu to Trp). Function specific conserved residues get the best score (e.g. Cys to Cys).

BLAST will provide a list of possible templates for the unknown structure To make the best initial alignment, BLAST uses an alignment-matrix based on the residue exchange matrix and adds extra penalties for opening and extension of a gap between residues.

The target-sequence is sent to a BLAST server, which searches the PDB to obtain a list of possible templates and their alignments. The best hit has to be chosen, which is not necessarily the first one.

### Conclusions

 Now the template and target are selected

Next, we perform fine-tuning of alignment and introduce corrections to ready the mismatches and gaps

# **Homology Modelling**

Seven Steps to Homology Modelling – III

# Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation

Fine tune and adjust the BLAST alignments

Example: Ala -> Glu is possible but unlikely in a hydrophobic core, so these residues should not be aligned.

Use MSA tools (e.g. ClustalW), to find the residues and properties that need to be conserved.

Examine the template structure to check which residues are in the core hence less likely to change than the residues at the outside.

Insertions and deletions can be made in those parts of the sequence which are highly variable .

Note: MSA can be helpful to find these places.

Gaps have to be shifted around until they are as small as possible.



Note: After a deletion of 3 residues a big gap occurs in the RED structure, which was the best alignment.



Remember that deletion of 3 residues left a big gap in the RED structure, which was the best alignment.

After shifting several residues, the gap is much smaller (blue structure) and more likely to be correct.

### Conclusions

 The alignment now stands fine tuned and corrected

Gaps and mismatches have been evaluated and adjusted

 Next step, using this alignment, assemble the backbone

# **Homology Modelling**

Seven Steps to Homology Modelling – V

# Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation

Note that the conserved residues were already copied! Now, we just need to place the side chains

Copy the torsion angles C-alpha/beta to the target! Rotamers tend to be conserved in homologous proteins and can be predicted as backbone configurations strongly prefer a specific rotamer.

Moreover, libraries of flanking residues can also help estimate the side chain positioning.

The backbone of tyrosine strongly prefers two rotamers and the real side-chain may fit one of them!



# Next....

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation

<u>What is the need for further optimization?</u> Because the updated side-chains can effect the backbone, and this can effect the structure prediction.

Optimization can be done by performing refinements using Molecular Dynamics simulations of the model.

The model is placed in a force-field and the movements of the molecules are followed in time, this mimics the folding of the protein.

<u>The large anomalies like bumps will be removed but new</u> <u>smaller errors can be introduced.</u>

The calculated energy should be as low as possible.
### MD Example Sim: Crambin

(Ethiopian Cabbage Protein)



### Conclusions

- Now we have minimized large errors
- However, smaller errors may still exist
- Next step, validate the model that we have constructed!

### **Homology Modelling**

Seven Steps to Homology Modelling – VI

### Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation

The model with the lowest energy might still be folded completely wrong.

So, the model should be checked again for normal ranges of: Bumps, Bond angles, Torsion angles, Bond lengths Other properties, like the distribution of polar/apolar residues, can be compared with real structures.

#### **Important Points:**

An error occurs far away from the active site, is tolerable. But when an error occurs in the active site, one should reconsider the template and/or alignment.

Limitations of Homology Modelling

Large Bias towards structure of template

Cannot study conformational changes

Cannot elicit new catalytic/binding sites

### Conclusions

- So how can we overcome such limitations?
- Other strategies include: Threading, and Ab Initio Modelling
- We will also examine online tools for each

### **Homology Modelling**

# MODELLER for Homology Modelling

### Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation

## Background

- Modeller is a software for homology modelling
- salilab.org/modeller
- Inputs: Python script file, Sequence alignment & Template (PDB)

from modeller import *		>P:	1; <b>1q</b>	50							
from modeller.automodel import *			structureX: 1q5o : 443 : A : 644 : A ::::								
log.verbose()			DSSRRQYQEKYKQVEQYMSFHKLPADFRQKIHDYYEHRYQ-GKMFDEDSILGELNGPLRE								
env = environ()			$\verb"EIVNFNCRKLVASMPLFANADPNFVTAMLTKLKFEVFQPGDYIIREGTIGKKMYFIQHGV"$								
<pre>env.io.atom_files_directory = './'</pre>			VSVLTKGNKEMKLSDGSYFGEICLLTRGRRTASVRADTYCRLYSLSVDNFNEVLEEYP								
			MMRRAFETVAIDRLDRIGKKNSIL. *								
a = automodel(											
env,			>Pl;herg								
alnfile = <b>'herg.</b> ali <b>'</b> ,											
knowns = $'1q5o'$ ,			ISGTAKINTQMLKVKEFIKFNQIPNPLKQKLEEIFQHAWSITNGIDMNAVLKGFPECLQA								
sequence = 'herg'											
)			EFSDHFWSSLEITFNLRDTN-MIP. *								
, ,											
a starting model= 1											
a ending model = $1$			Converse Alignment (* eli)								
a make()			Sequence Alignment (^.ali)								
	ATOM	1	Ν	ASP A	443	-15.943	41.425	44.702	1.00 44.68	3	
	ATOM	2	CA	ASP A	443	-15.424	42.618	45.447	1.00 43.15	5	
	ATOM	3	С	ASP A	443	-14.310	43.306	44.686	1.00 41.81		
	ATOM	4	0	ASP A	443	-14.298	44.528	44.539	1.00 42.61	-	
		etc									
Innut Duthon Covint											
input Python Script L											

Template Structure (\*.pdb)

(\*.py)



- .log : log output from the run.
- **.B\*** : model generated in the PDB format.
- **.D**\* : progress of optimisation.
- **.V**\* : violation profile.
- .ini : initial model that is generated.
- .rsr : restraints in user format.
- .sch : schedule file for the optimisation process.

**Automated Modelling Servers** 

Swiss Model <u>http://swissmodel.expasy.org//SWISS-MODEL.html</u>

Robetta

http://robetta.bakerlab.org/

3D Jigsaw http://www.bmm.icnet.uk/servers/3djigsaw/

Phyre

http://www.sbg.bio.ic.ac.uk/phyre/

#### Conclusions

- Homology modelling helps predict protein structures by using prior structural information
- Several tools are available to perform homology modelling in a programmatic or automated way!

## **Homology Modelling**

# Fold Recognition – Threading I

## Background

**Template recognition** and initial alignment **Alignment correction Backbone generation** Loop modeling Side-chain modeling **Model optimization** Model validation



#### When should we use Fold Recognition?



http://www.cmbi.ru.nl/gvteach/astra/lectures/homology\_modelling.ppt

### Introduction

- A protein fold is defined by the way the secondary structure elements of the structure are arranged relative to each other in space.
- Common folds include 4-helix bundle and the TIM barrel.

## Introduction

 5,000 stable folds in nature

• Fold recognition: Finding the best fit of a sequence to a set of candidate folds

### Conclusions

- Fold recognition or Threading is a technique for predicting protein structures
- It is useful in cases where homology modelling fails to predict quality structures

### **Homology Modelling**

# Fold Recognition – Threading II

## Background

- Fold recognition is also called Threading
- Technique for predicting protein structures
- Employed when homology modelling cannot predict quality structures

Find the best way to

#### "mount" the residue sequence of one protein

#### onto a known protein structure!



#### The process of threading

- In the process of "Threading", we mount an amino acid sequence on to the backbone of template structures in a folds library
- Each step is "drag" along the sequence (MQVKLFTY...) through each location of each template fold
- Then, for each fold, we must compute the fitness of sequence matching that fold!

### **Inputs and outputs of threading**



### Conclusions

- Threading involves "passing" the amino acid sequence through each fold in the database
- The best match is computed using a scoring function

### **Homology Modelling**

# Fold Recognition – Threading III

### Background

- Threading involves "passing" the amino acid sequence through each fold in the database
- The best match is computed using a scoring function

Flowchart

### **Inputs and outputs of threading**







www.upch.edu.pe

### Conclusions

- Combinations of secondary structures come together to form the best prediction
- Scoring typically involves using a Z-Score function based on energy of a molecule

### **Homology Modelling**

# Online Tools for Threading iTasser

#### **Online Tools for Threading - iTasser**

### Background

- Threading involves "passing" the amino acid sequence through each fold in the database
- The best match is computed using a scoring function

iTASSER
## **Inputs and outputs of threading**



#### Iterative threading assembly refinement (I-TASSER) server

- Software for automated protein structure & function prediction based on the sequence-to-structure-to-function.
- Steps:
  - Starts from amino acid sequence
  - i-TASSER first generates 3D atomic models from multiple threading alignments and iterative structural assembly simulations.
  - The function of the protein is then inferred by structurally matching the 3D models with other known proteins.
  - Outputs full-length secondary & tertiary structures and functional annotations on ligand-binding sites
  - An estimate of accuracy of the predictions is provided based on the confidence score of the modeling

#### Link: http://zhanglab.ccmb.med.umich.edu/I-TASSER/

9	Zin <b>¢</b> i	ng L	20			
Home	Research	Services	Publications	People	Teaching	Job Opening Ne
Online Services			The man	_	~~	_
I-TASSER				-TA	55	EK
QUARK	Protein Structure & Function Predictions					
LOMETS	(The server completed predictions for 275701 proteins submitted by 68550 users from 12					
COACH	(The template library was updated on 2016/05/12)					
COFACTOR	I-TASSER (Iterative Threading ASSEmbly Refinement) is a hierarchical approach to protein structur Structural templates are first identified from the PDB by multiple threading approach <u>LOMETS</u> ; full-leng constructed by iterative template fragment assembly simulations. Finally, function inslights of the targe the 3D models through protein function database <u>BioLiP</u> . I-TASSER (as 'Zhang-Server') was ranked as structure prediction in recent community-wide <u>CASP7</u> , <u>CASP8</u> , <u>CASP9</u> , <u>CASP10</u> , and <u>CASP11</u> experi as the best for function prediction in <u>CASP9</u> . The server is in active development with the goal to p structural and function predictions using state-of-the-art algorithms. The server is only for non-comm					
MUSTER						
SEGMER						
FG-MD						
ModRefiner						
REMO						
SPRING	about the server)					
COTH		, ,		1		
BSpred	through clou	ncial users or tr ud computing.	lose who want to si	uomit a large hu	imper of jobs, pl	ease go to <u>DINAStar</u> which
SVMSEQ	10					
ANGLOR	[Queue] [Forum] [Download] [Search] [Registration] [Statistics] [Remove] [Potential] [Decoys] [Ne [FΔΩ]					

Copy and paste your sequence below ([10, 1500] residues in FASTA format). Click here for a sample input:

Or upload the sequence from your local computer: Choose File No file chosen

Email: (mandatory, where results will be sent to)

Password: (mandatory, please click here if you do not have a password)

ID: (optional, your given name of the protein)

Option I: Assign additional restraints & templates to guide I-TASSER modeling.

Option II: Exclude some templates from I-TASSER template library.

Option III: Specify secondary structure for specific residues.

Keep my results public (uncheck this box if you want to keep your job private. A key will be assigned for you to access the results)

Run I-TASSER Clear form

(Please submit a new job only after your old job is completed)

(If you want to submit multiple urgent jobs, please go to DNAStar which implements I-TASSER through cloud computing)

## Conclusions

- iTASSER helps thread amino acid sequences on fold and secondary structure databases
- It also helps predict function of structures output.

## **Homology Modelling**

# Advantages and Disadvantages of Threading

#### **Advantages and Disadvantages of Threading**

## Background

- Fold recognition or Threading is a technique for predicting protein structures
- It is useful in cases where homology modelling fails to predict quality structures

#### **Advantages and Disadvantages of Threading**

Inputs and outputs of threading



www.upch.edu.pe

#### **Advantages and Disadvantages of Threading**

### **Advantages**

- Threading helps predict secondary structures of proteins towards tertiary structure prediction
- For the "Twilight Zone" with low alignment quality and identity, threading is useful

## **Disadvantages**

- Novel proteins cannot be predicted using threading
- Fewer than 30% of the predicted first hits are true remote homologues
- Validation of each
   result is necessary

## **Homology Modelling**

# **3D-1D Bowie Algorithm**

#### **3D-1D Bowie Algorithm**

## Background

 Homology employed high alignment scores

 Threading worked by creating combinations of primary sequences and corresponding secondary structures

#### **3D-1D Bowie Algorithm**

## Introduction

- Proposed by Bowie et al in 1991
- Converts 3D structure into a 1-D string profile for each structure in the fold library
- Align the target sequence to these profiles

### **3D-1D Bowie Algorithm**

#### Inputs and outputs of 3D-1D

- Identify amino acids based on: protein core, side chain positioning, solubility etc. (6 in all)
- Part of secondary structure including  $\alpha$ -helix,  $\beta$ -sheet etc (3 in all)
- Total of 3 x 6 = 18 distinct states

### **Inputs and outputs of 3D-1D**

P<sub>a: i</sub>= prob. of finding amino acid (a) in environment (j)

**P**<sub>a</sub>=probability of finding (a) anywhere

Maximize sum of scores for the fold:

$$s_{aj} = \log\left(\frac{P_{a:j}}{P_a}\right)$$

## Conclusion

- 3D-1D methods convert structure and environment information into "profiles"
- Score for each amino acid is computed for each profile

## **Homology Modelling**

# Introduction to Ab Initio Modelling

## Background

- Ab initio methods have Anfinsen's thermodynamic hypothesis at the center
- These methods
   attempt to identify
   the structure with
   minimum free
   energy

## Need for Ab Initio Modelling

- Applicable to any sequence
- Not very accurate biologically
- Accuracy and applicability are limited by our understanding of the protein folding problem

## Limitation

Computationally expensive

 Suitable for proteins with less than 100 residues



## Conclusion

 Ab initio methods rely on computing the energies of folded proteins

The protein structures with the lowest energy are deemed as plausible predictions

## **Homology Modelling**

# Rationale of Ab Initio Modelling

## Background

 Ab initio methods rely on computing the energies of folded proteins

 The protein structures with the lowest energy are declared as plausible predictions

## Rationale

- Sometimes it so happens that even slightly homologous proteins may not be available.
- This renders
   homology modelling
   and threading/fold
   recognition as futile

## Rationale

- Also, newer protein structures continue to be discovered every day
- These could not
  have been identified
  by methods which
  only rely on
  matching with
  available structures

## Rationale

Lastly, homology / fold recognition predict protein structures without computing fundamental physical/chemical properties of the mechanisms and driving forces in structure formation

## Conclusion

- Ab initio methods, in contrast, base their predictions on physical models for these mechanisms
- Energy released during the folding process is computed for predicting structure

## **Homology Modelling**

# Strategies for Ab Initio Modelling

#### **Strategies for Ab Initio Modelling**

## Background

- Ab initio methods base their predictions on physical models of folding mechanisms
- Stabilization is measured by energy released during the folding process

**Energy Optimization in Ab Initio Modelling** 

1. Start with a rough initial model.

- 2. Define an energy function mapping structures to energy values. We have to minimize this later!!
- 3. Solve the computational problem of finding the global minimum.

## **Strategies for Ab Initio Modelling**

## **Simulation of the Folding Process**

- 1. Build an accurate initial model (including energy and forces).
- 2. Accurately simulate the dynamics of the protein folding process.
- 3. The native structure will steadily emerge.

### **Strategies for Ab Initio Modelling**

## Conclusion

- Start with an energy function
- Fold structures in order to obtain the most stable structure
- This structure will have the minimum energy

## **Homology Modelling**

## **Energy States of Folded Proteins**

#### **Energy States of Folded Proteins**

## Background

Ab initio methods predict protein structures by folding proteins based on each constituent atom's volume, charge, mass etc.

#### **Energy States of Folded Proteins**

#### **Energies of Bonded Atoms vs. Nonbonded Atoms**


### **Energy States of Folded Proteins**

Force Fields for Energy Calculations ->Start with an initial structure



### **Energy States of Folded Proteins**

#### **Force Fields for Energy Calculations**



### **Energy States of Folded Proteins**

# Conclusion

- The protein structure reporting lowest energy is selected to be the optimal structure
- How easy is it to compute the "really" lowest energy of a folded protein?

# **Homology Modelling**

# Local versus Global Minima

# Background

- The protein structure reporting lowest energy is selected to be the optimal structure
- How easy is it to compute the "really" lowest energy of a folded protein?

#### **Energies of Bonded Atoms vs. Nonbonded Atoms**



#### **Force Fields for Energy Calculations**



Global Energy Optimization helps find global minimum



Global Energy Optimization helps find global minimum



# **Best Case Energy Function**

- Clear energy minimum in the native structure
- Viable path towards this minimum
- Global optimization finds the most stable structure

# **Optimal Energy Function**

- Easier to design and compute
- Native structure not always at the global minimum
- No clear way of choosing among alternative structures that are generated

# **Homology Modelling**

# Pros and Cons of Ab Initio Modelling

## Background

- Native structure not always at the global minimum
- No clear way of choosing among alternative structures that are generated

### **Advantages**

- Ab Initio methods can fold any target sequence using only physical atomic properties
- Predictions are mostly accurate and correctly describe the natural folding process

### **Disadvantages**

- Ab initio methods are the very difficult to design (energy function)
- These methods are slow due to the huge possibilities

### **Disadvantages**

• An order of 10<sup>12</sup> steps are needed to simulate protein folding for medium sized protein structures

# Challenges in Ab Initio Modelling

 Very hard to accurately describe energy functions that can reliably discriminate native and non-native structures.

• Enormous amount of computations.

# **Homology Modelling**

Summary of Structural Modelling – I

Strategies for Structural Modelling

- Homology Modelling
- Fold Recognition
- Ab Initio Modelling



Homology modeling of the target structure can be done as follows:

- 1. Template recognition and initial alignment
- 2. Alignment correction
- 3. Backbone generation
- 4. Loop modeling
- 5. Side-chain modeling
- 6. Model optimization
- 7. Model validation

# Conclusion

- Homology modelling is performed in cases of high identity and alignment score
- For the "Twilight zone", other strategies are employed

# **Homology Modelling**

# Summary of Structural Modelling – II

Strategies for Structural Modelling

- Homology Modelling
- Fold Recognition
- Ab Initio Modelling



# **Inputs and outputs of threading**





# Conclusion

- For low identity and alignment scores, a "Twilight zone" for structure prediction exists
- Fold recognition / threading is useful in such cases

# **Homology Modelling**

# Summary of Structural Modelling – III

Strategies for Structural Modelling

- Homology Modelling
- Fold Recognition
- Ab Initio Modelling



**Energy Optimization in Ab Initio Modelling** 

1. Start with a rough initial model.

- 2. Define an energy function mapping structures to energy values. We have to minimize this later!!
- 3. Solve the computational problem of finding the global minimum.

# **Simulation of the Folding Process**

- 1. Build an accurate initial model (including energy and forces).
- 2. Accurately simulate the dynamics of the protein folding process.
- 3. The native structure will steadily emerge.

# Conclusion

- For cases where even the fold
  libraries do not give any high scoring matches, Ab Initio
  strategies can help
  model the structure
- However, this is a complex and computationally expensive process

## **Conclusion of the Course**

Review of Sequence Analysis

### **Review of Sequence Analysis**

# Important Concepts

How do we sequence:

- Genomes
- Proteomes
# Important Concepts

How do we compare sequences:

- Pair-wise Sequence
  Alignment
- Multiple Sequence
  Alignment

## Important Concepts

**Types of Alignments:** 

**Global Alignment** (Needle Wunsch)

Local Alignment (Smith Waterman)

## Important Concepts

**Advanced Tools:** 

Fast Alignment (FASTA)

Basic Local Alignment Search Tool (BLAST)

# Important Concepts

**Databases:** 

GenBank

UniProt

Important Concepts

**Online Portals:** 

Ensemble

**Expasy** 

**UniProtKB** 

## **Conclusion of the Course**

# **Review of Phylogenetics**

## **Review of Phylogenetics**

# Important Concepts

Molecular Evolution

- Insertions
- Deletions
- Substitutions

## **Review of Phylogenetics**

# Important Concepts

Phylogenetic Trees

- Scaled Trees
- Unscaled Trees

#### Edge length with and without a clock!



# Since the rate of evolution in species is different, it needs to be considered!

## **Review of Phylogenetics**

# Important Concepts

Phylogenetic Trees

- Rooted Trees
- Unrooted Trees

# **UPGMA: Unweighted Pair – Group Method using arithmetic Averages**

- Two sequences with with the shortest evolutionary distance between them are considered
- These sequences will be the last to diverge, and represented by the most recent internal node.



## **Review of Phylogenetics**

## Important Concepts

**Clustering Vs. Nonclustering Methods:** 

UPGMA is a clustering method

Maximum Parsimony etc are non-clustering methods (not included in this course).

## **Conclusion of the Course**

# Review of Protein Sequencing

# Important Concepts

Techniques of
 protein sequencing

- Edman Degradation
- Mass Spectrometry

# Important Concepts

- Protein Ionization
- Mass Analysis

Protein
 Fragmentation

# Important Concepts

• MS1

• MS2

# Important Concepts

 Estimating and scoring whole protein mass

# Important Concepts

 Extracting & Scoring Peptide Sequence Tags

# Important Concepts

 Searching Posttranslational Modifications

# Important Concepts

**Composite Scoring Schemes** 

**Online tools:** 

- Mascot
- Sequest
- Prosight PC

## **Conclusion of the Course**

# **Review of Protein Structures**

## Important Concepts

Protein Structures are generally of four types:

- Primary
- Secondary
- Tertiary
- Quaternary

## Important Concepts

Techniques for determining protein structures

 X-Ray Crystallography

NMR Spectroscopy

# Important Concepts

Why number of known protein sequences is much larger as compared to known proteins structures?

## Important Concepts

Types of Protein Secondary Structures

- Helices
- Beta Sheets
- Coils
- Loops

## Important Concepts

 Foundation of structure prediction algorithms

 Propensities of certain amino acids to form specific secondary structures

# Important Concepts

- Algorithm for predicting protein structures
- Chou Fasman Algorithm

# Important Concepts

Protein Structure
 Database - PDB

 Online tools for predicting structures by using proteins sequences

## **Conclusion of the Course**

# **Review of Protein Structures**

## Important Concepts

Protein Structures are generally of four types:

- Primary
- Secondary
- Tertiary
- Quaternary

## Important Concepts

Techniques for determining protein structures

 X-Ray Crystallography

NMR Spectroscopy

# Important Concepts

Why number of known protein sequences is much larger as compared to known proteins structures?

## Important Concepts

Types of Protein Secondary Structures

- Helices
- Beta Sheets
- Coils
- Loops

## Important Concepts

 Foundation of structure prediction algorithms

 Propensities of certain amino acids to form specific secondary structures

# Important Concepts

- Algorithm for predicting protein structures
- Chou Fasman Algorithm

# Important Concepts

Protein Structure
 Database - PDB

 Online tools for predicting structures by using proteins sequences
# **Conclusion of the Course**

Review of Homology Modelling

# Important Concepts

Four Strata of Protein Structures

- Primary
- Secondary
- Tertiary
- Quaternary

- Justification for homology modelling
- Number of known protein sequences is much larger as compared to known proteins structures

# Important Concepts

**Three Strategies for Structure Prediction** 

- Homology Modelling
- Fold Recognition
- Ab Initio Modelling



http://www.cmbi.ru.nl/gvteach/astra/lectures/homology\_modelling.ppt

#### **Review of Protein Structures**

# Important Concepts

Protein Structure
 Database - PDB

 Online tools for predicting structures such as MODELLER and iTASSER

# **Conclusion of the Course**

# Conclusions from this Course

- Definition of Bioinformatics
- Need for Bioinformatics
- Areas within Bioinformatics

- Bioinformatics as an interdisciplinary area
- Need to store, process and analyze biological data
- Requirement of newer faster algorithms

# Important Concepts

Specific areas focused were:

- Comparing sequences
- Comparing structures
- Predicting
   structures

Important Concepts

We looked at:

- Algorithms
- Databases
- Online Tools

for each topic.

- We studied the basic algorithms for each topic
- With evolution and growth of
  Bioinformatics,
  newer and better
  algorithms are now
  also available!

# **Conclusion of the Course**

Advanced Follow-up Courses

- We looked into the foundations of Bioinformatics
- However, each topic that was studied has a undergone a lot of development

- For advanced study in Genomics, you may take "Computational Genomics" course
- Topics:
   Genome Assembly,
   Gene Finding,
   Annotation, GWAS
   etc

# Important Concepts

 For advanced study in Proteomics, you may take "Computational Proteomics" course.

Topics:
 Protein Sequencing,
 PTM search,
 Structure Modelling
 and PPI studies

# Important Concepts

 For advanced study in Integrative Biology, you may take "Systems Biology" course.

 Topics: Metabolomics, Transcriptomics, Network Biology etc

# Important Concepts

Also, now there are cutting edge courses on:

- Nano-Bio-IT
- Computational Drug Design
- Personalized
   Medicine

# **Conclusion of the Course**

# **Careers in Bioinformatics**

# Background

- Pakistan as an infrastructurelimited country
- The onset of digital revolution
- Emergence of data as the most precious commodity, globally

# Background

 Specifically, health data as a key commodity of the future

 Health and disease as the primordial challenge of mankind

# Background

- Unique opportunity for us in Pakistan
- Bioinformatics requires two things
- 1. Smart mind
- 2. Internet connected computer

# **One man company**

 You can take public databases and design drugs!

One man vs. Roche?

# BIGDATA

- You can make a startup company which manages and process health BIGDATA!
- All it needs is basic software development skills coupled with Bioinformatics

# The next disruption

- The next Google, Facebook and Uber is going to emerge from Health and Bioinformatics
- Pharmaceutical companies are investing into bioinformatics human resource development

# **Jobs Market**

- Pharmaceutical Giants
- Research Centers & Universities
- Hospital & Diagnostic IT departments
- Your own startup company

# Final Term Syllabus

# **Effort By**

# Amaan Khan