

# Problem Solving in Bioinformatics



**Jianlin Cheng, PhD**

Informatics Institute, Computer Science Department

University of Missouri, Columbia

Fall, 2009

# Objectives

- Walk students through the complete process of sequencing, assembling and annotating a genome. During the process, students learn key bioinformatics techniques for analyzing a genome and its components (i.e. gene, RNA, protein, and pathway).
- By working on a comprehensive genome annotation project, students develop practical skills to apply bioinformatics methods to solve major problems in genome assembly and annotation.

# Instructors

- Jianlin Cheng, PhD (coordinator)
- Toni Kazic, PhD (?)
- Dmitry Korkin, PhD
- Chi-Ren Shyu, PhD
- Dong Xu, PhD

# Topics

- Introduction to the course and a project (Jianlin Cheng)
- Genome sequencing and assembly (Dong Xu)
- Gene prediction (Dong Xu)
- Protein structure prediction (Jianlin Cheng)
- Protein function prediction (Toni Kazic or Jianlin Cheng?)
- Protein interaction prediction (Dmitry Korkin)
- Biological pathways and networks (Chi-Ren Shyu)

# Course Format

- **Theory phase** (one month): lecturing, literature review, mid-term exam
- **Practice phase** (two and a half months): discussion, planning (group), presentation, programming (group), results (group), assessment (group), and report (group)
- Teamwork & leadership (two groups)
- See syllabus for details

# Assignments

- Literature review, topic plan (in presentation style), implementations of genome assembly and annotations (programs and results), topic report, and final report and presentation.
- All the assignments should be posted to the project web site or emailed to me by deadlines.

# Evaluation and Grading

- literature reviews (individual, 15%), mid-term exam (individual, 15%), class discussion (individual, 10%), presentations (individual, 10%), topic plans and reports (i.e. progress and assessment) (group, 15%), topic implementation (group, 25%), a final presentation and report (group, 10%)
- Group components graded by both instructors and group peers

# Course Web & Class Schedule

- Course web (**demo**):

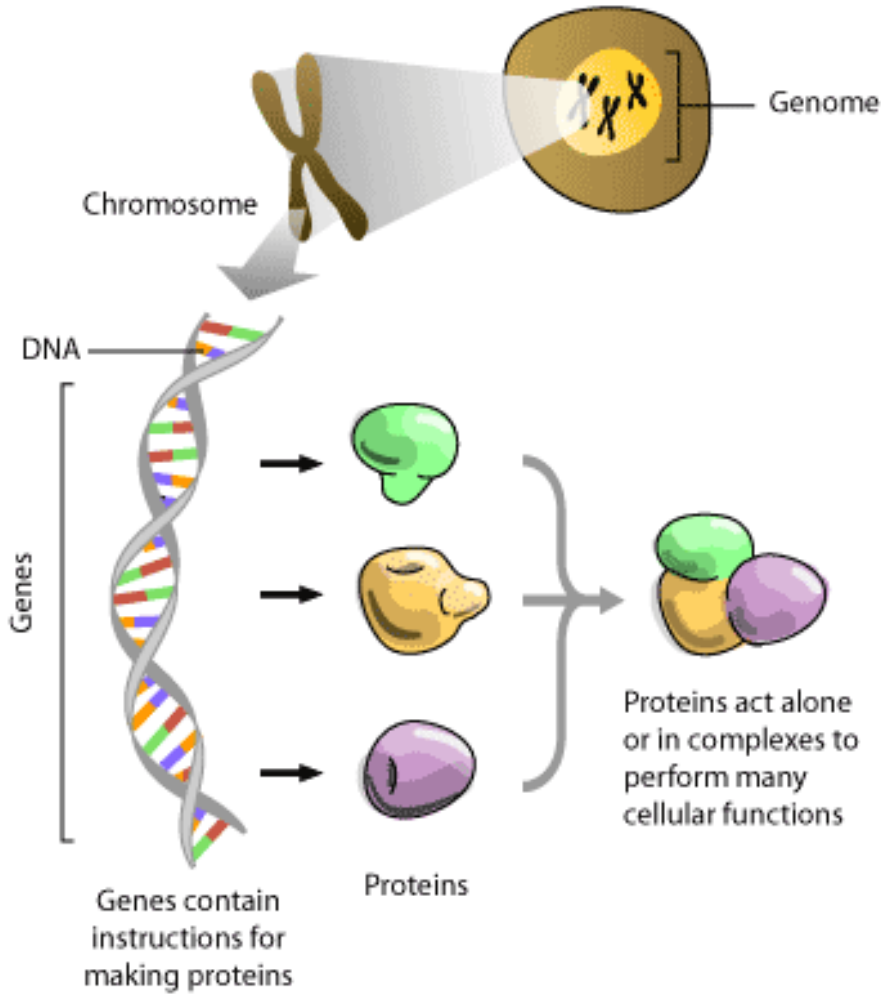
<http://www.cs.missouri.edu/~chengji/infoinst8010/>

- Class schedule and assignments:

[http://www.cs.missouri.edu/~chengji/infoinst8010/8010\\_schedule.htm](http://www.cs.missouri.edu/~chengji/infoinst8010/8010_schedule.htm)

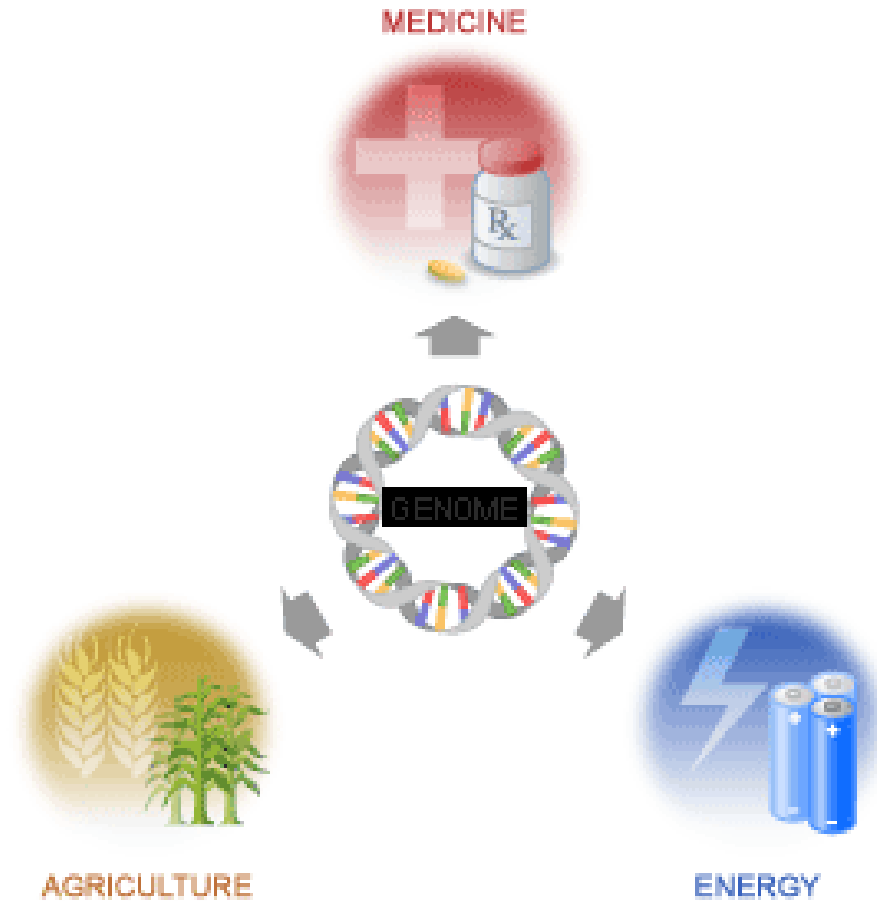


# Introduction



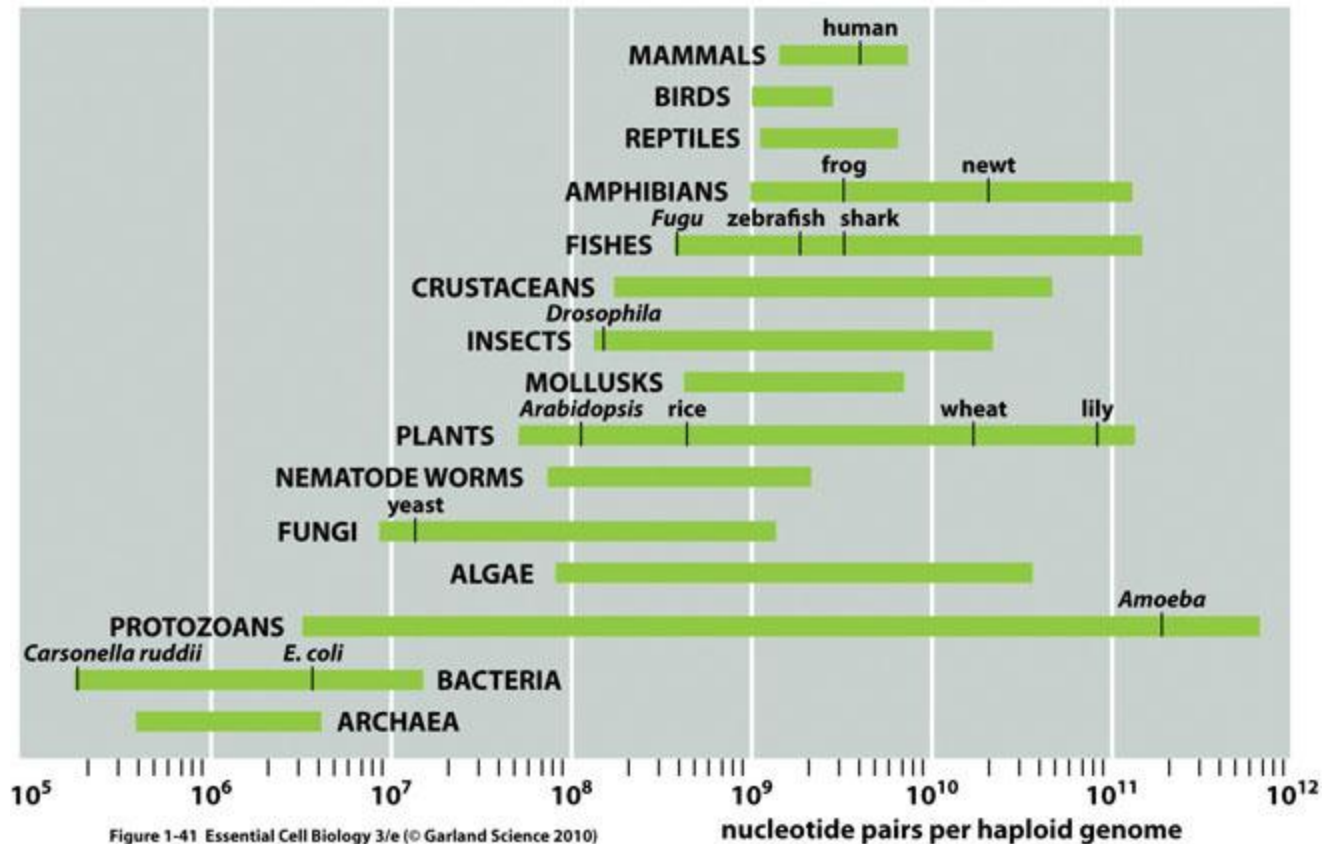
- **Grow**
  - **Sustain**
  - **Adapt**
  - **Reproduce**
- 
- **Genome & Components**
  - **Environment**

# Applications of Genome Knowledge



# Genome Sequencing – Cracking the Code

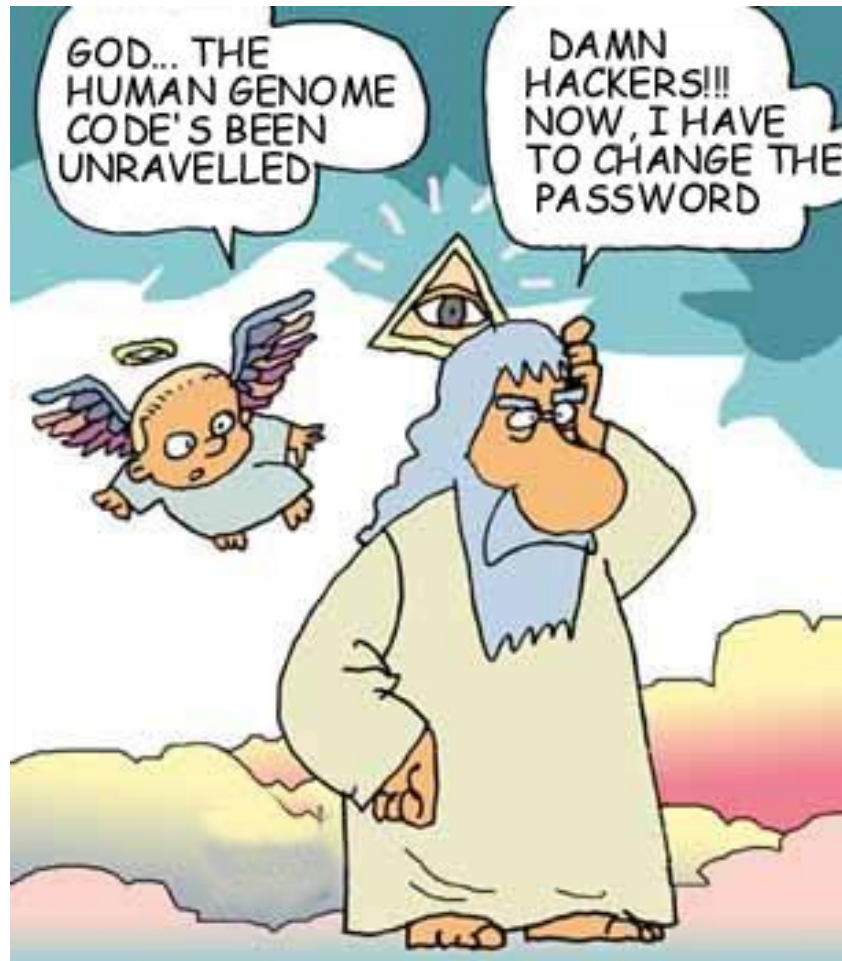
- Virus & Bacteria genomes (small)
- Human genome



# Human Genome Project



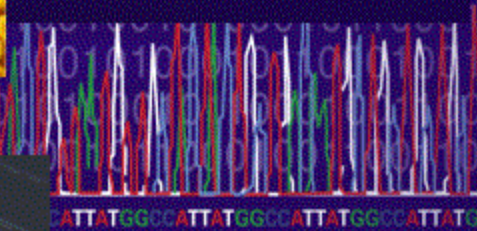
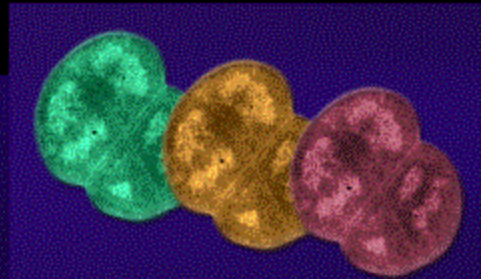
# Fun



# Genome Sequencing Routine



## BEYOND THE HUMAN GENOME PROJECT New Discovery Paths and Diverse Applications





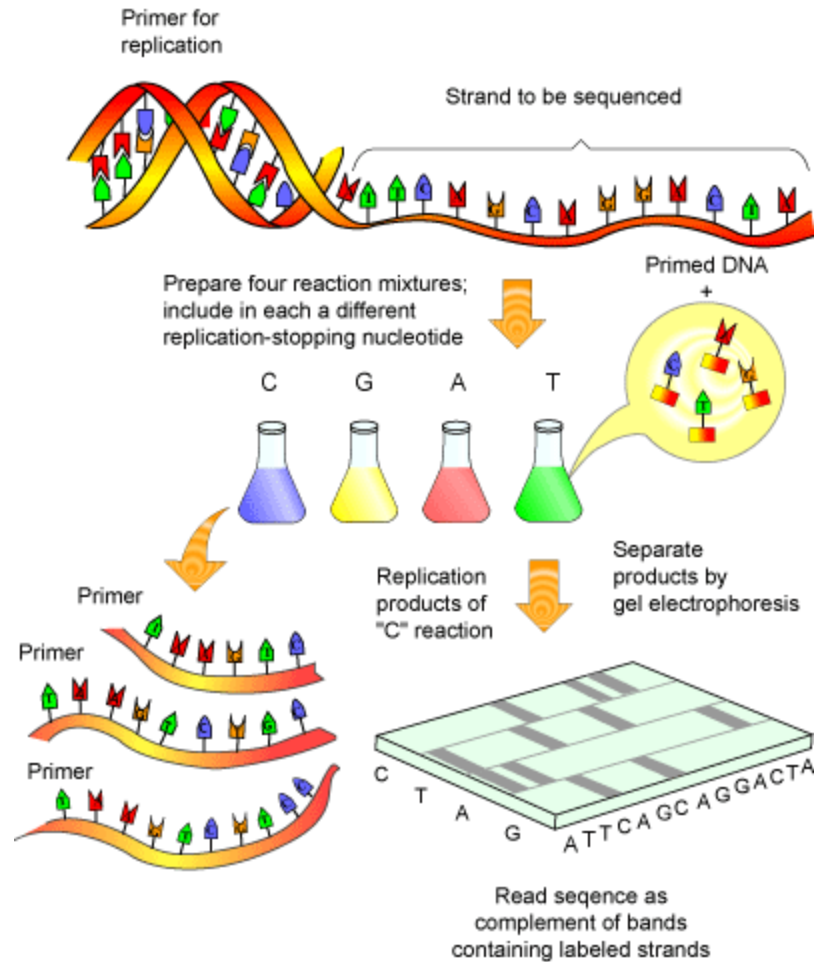
# MU Genome Sequencing

- Soybean genome: Gary Stacey, Dong Xu, Jay Thelen, Jianlin Cheng, etc (to be published in Nature)
- Chris Pires – plant genomes

**DOE**   
**Bioenergy**  
**Research**  
**Centers**



# Sequencing process





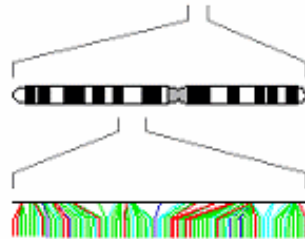
# Genome Sequencing Machine



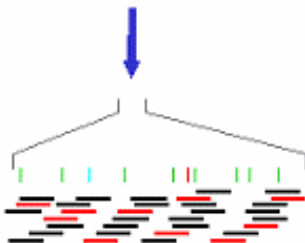
# STRATEGIES FOR SEQUENCING THE HUMAN GENOME

## BY MAPPED CLONES

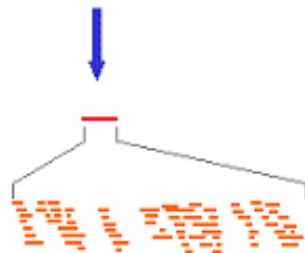
1. Construction of maps of ordered landmarks (genetic markers, genes); provides long-range map and organisation into individual chromosomes.



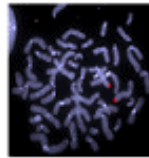
2. Physical maps of overlapping clones anchored to the landmark maps.



3. Selection of tile path (clones in red)



4. Shotgun sequencing and assembly (for working draft); subsequent directed finishing (for reference sequence).



## BY WHOLE GENOME SHOTGUN

1. Shotgun sequencing of short-insert clones



2. Paired end sequencing of large-insert clones



3. Assembly of seed contigs (unitigs)



4. Incorporation of other sequences, and integration of long-range data.



# Topic 1: Genome Assembly

a) Multiple copies of genome



b) Sheared random fragments



c) Size fractionated fragments



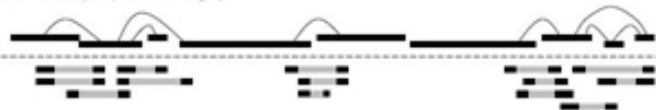
d) Reads



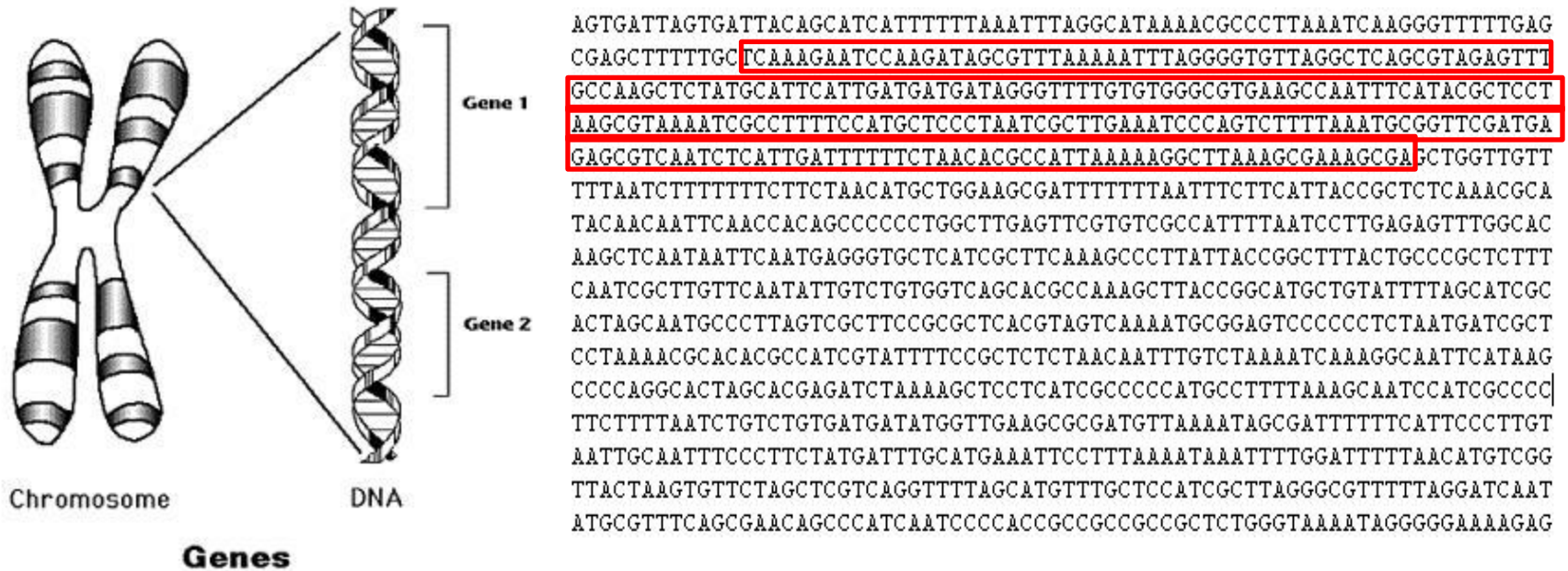
e) Contigs



f) Scaffolds(Super contigs)



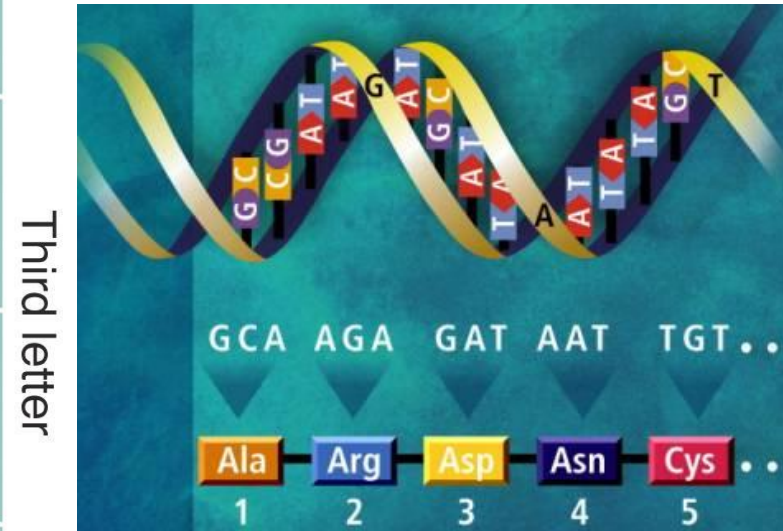
# Topic 2. Gene Prediction



## Pattern Recognition Problem

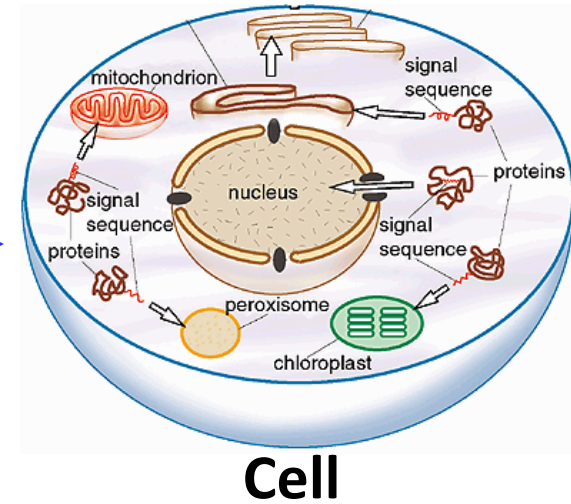
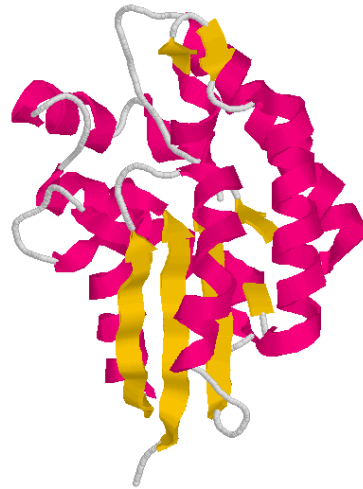
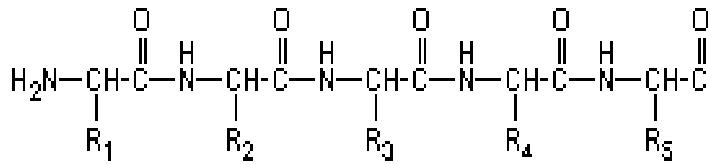
# Gene Product - Protein

		Second letter					
		U	C	A	G		
U	UUU } Phe	UCU } Ser	UAU } Tyr	UGU } Cys	U		
	UUC } Phe		UAC } Tyr			UGC } Cys	
	UUA } Leu		UAA Stop				UGA Stop
	UUG } Leu		UAG Stop				UGG Trp
C	CUU } Leu	CCU } Pro	CAU } His	CGU } Arg	U		
	CUC } Leu		CAC } His				
	CUA } Leu		CAA } Gln				
	CUG } Leu		CAG } Gln				
A	AUU } Ile	ACU } Thr	AAU } Asn	AGU } Ser	U		
	AUC } Ile		AAC } Asn				
	AUA } Met		AAA } Lys				
	AUG } Met		AAG } Lys				
G	GUU } Val	GCU } Ala	GAU } Asp	GGU } Gly	U		
	GUC } Val		GAC } Asp				
	GUA } Val		GAA } Glu				
	GUG } Val		GAG } Glu				



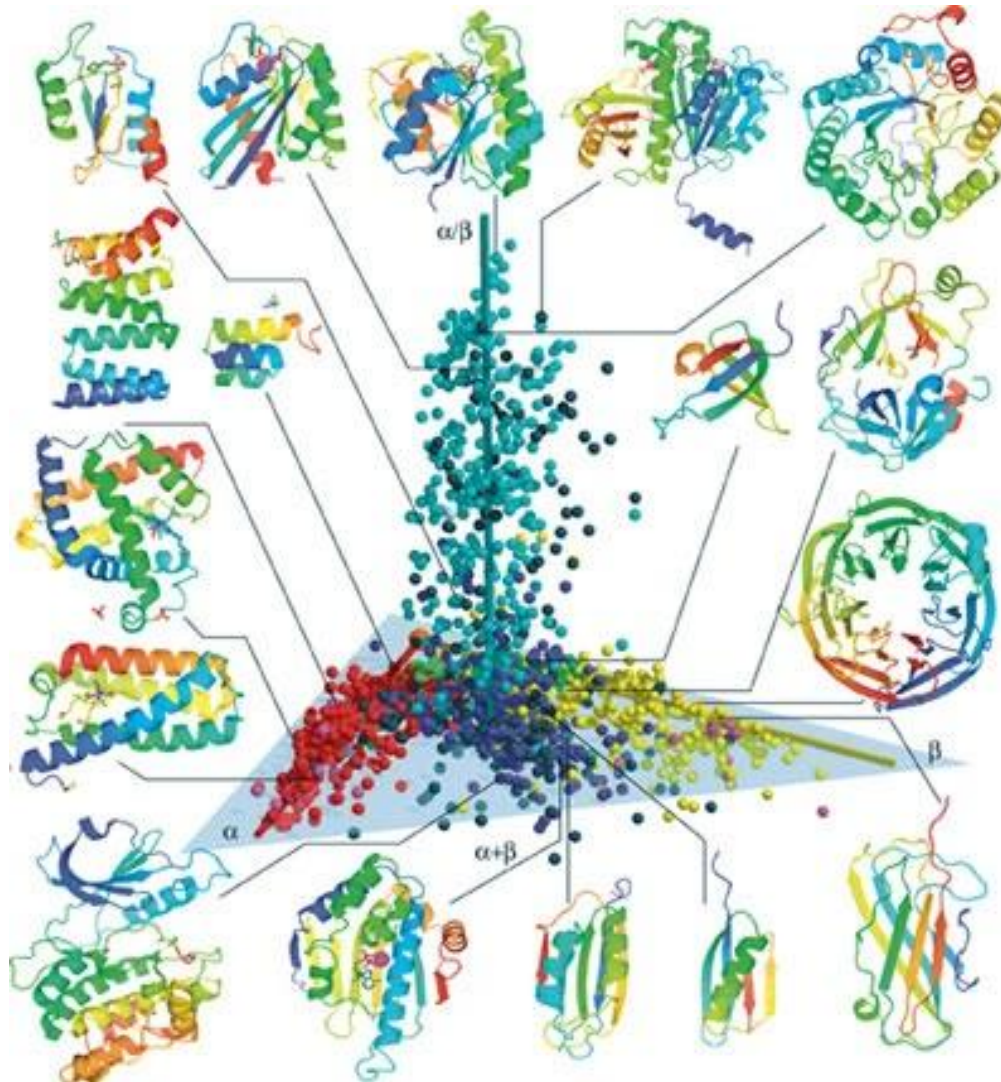
# Protein Sequence, Structure, Function

AGCWY.....





# Protein Structure Space



# Protein Data Bank

### Quick Tips :

Want to search by sequence? Click here.

Are you missing data updates? The PDB archive has moved to <ftp://ftp.wwpdb.org>. For more information click [here](#).

## Welcome to the RCSB PDB

The RCSB PDB provides a variety of tools and resources for studying the structures of biological macromolecules and their relationships to sequence, function, and disease.

The RCSB is a member of the [wwPDB](#) whose mission is to ensure that the PDB archive remains an international resource with uniform data.

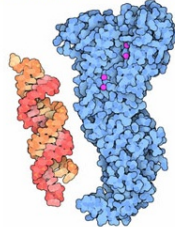
This site offers tools for browsing, searching, and reporting that utilize the data resulting from ongoing efforts to create a more consistent and comprehensive archive.

Information about compatible browsers can be found [here](#).

A [narrated tutorial](#) illustrates how to search, navigate, browse, generate reports and visualize structures using this new site. [This requires the Macromedia Flash player download.]

Comments? [info@rcsb.org](mailto:info@rcsb.org)

### Molecule of the Month: Small Interfering RNA (siRNA)



Double-stranded RNA is often a sign of trouble. Our transfer RNA and ribosomes do contain little hairpins that are double-stranded, but most of the free forms of RNA, messenger RNA molecules in particular, are single strands. Many viruses, however, form long stretches of double-stranded RNA as they replicate their genomes. When our cells find double-stranded RNA, it is often a sign of an infection, and they mount a vigorous response that often leads to death of the entire cell. However, plant and animal cells also have a more targeted defense that attacks the viral RNA directly, termed RNA interference.

- More ...
- Previous Features

The RCSB PDB is managed by two members of the RCSB: Rutgers, The State University of New Jersey and the San Diego Supercomputer Center and Skaggs School of Pharmacy and Pharmaceutical Sciences at the University of California, San Diego. It is supported by funds from the National Science Foundation (NSF), the National Institute of General Medical Sciences (NIGMS), the Office of Science, Department of Energy (DOE), the National Library of Medicine (NLM), the National Cancer Institute (NCI), the National Center for Research Resources (NCRR), the National Institute of Biomedical Imaging and Bioengineering (NIBIB), National Institute of Neurological Disorders and Stroke (NINDS), and the National Institute of Diabetes and Digestive and Kidney Diseases (NIDDK).

### News

- Complete News
- Newsletter
- Discussion Forum
- Job Listings

05-February-2008

#### Historical Look at the PDB Published in a Special Issue of *Acta Crystallographica*

The PDB archive has grown from its early beginnings in 1971 as a handwritten petition signed by crystallographers to its current status as an online biological database and resource used by a diverse community of teachers, students, and researchers in academia and industry worldwide.

- Full article ...

04-December-2007

#### Announcement: Experimental Data Will Be Required for Depositions Starting February 1, 2008

Effective February 1, 2008, structure factor amplitudes/intensities (for crystal structures) and restraints (for NMR structures) will be a mandatory requirement for PDB deposition.

- Full article ...

In citing the PDB please refer to: H.M. Berman, J. Westbrook, Z. Feng, G. Gilliland, T.N. Ehtai, H. Weissig, I.N. Shindyalov, P.E. Bourne. The Protein Data Bank. *Nucleic Acids Research*, 28 pp. 235-242 (2000).

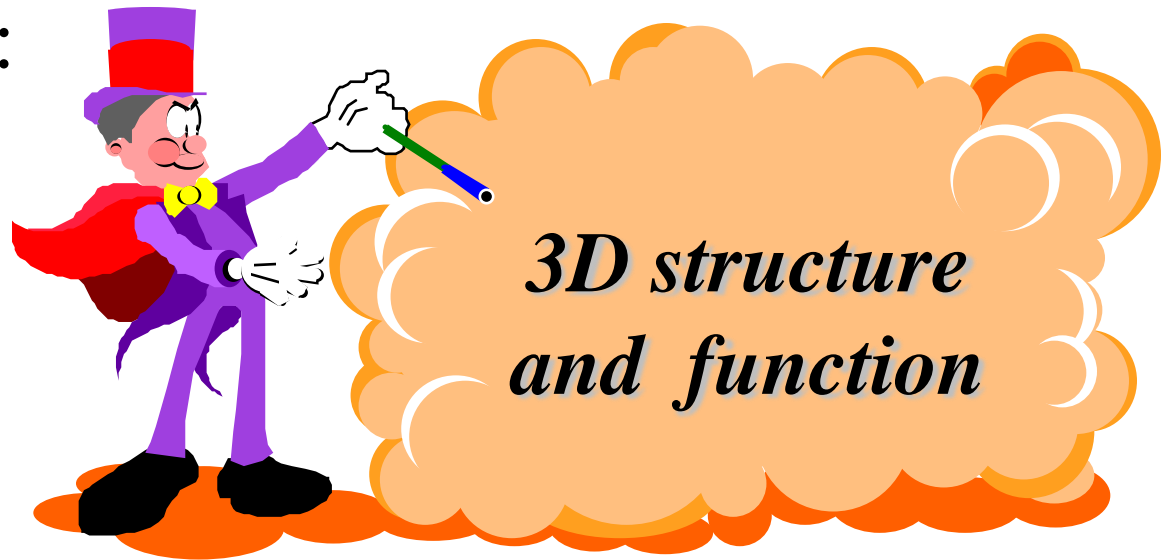


# Topic 3: Protein structure prediction

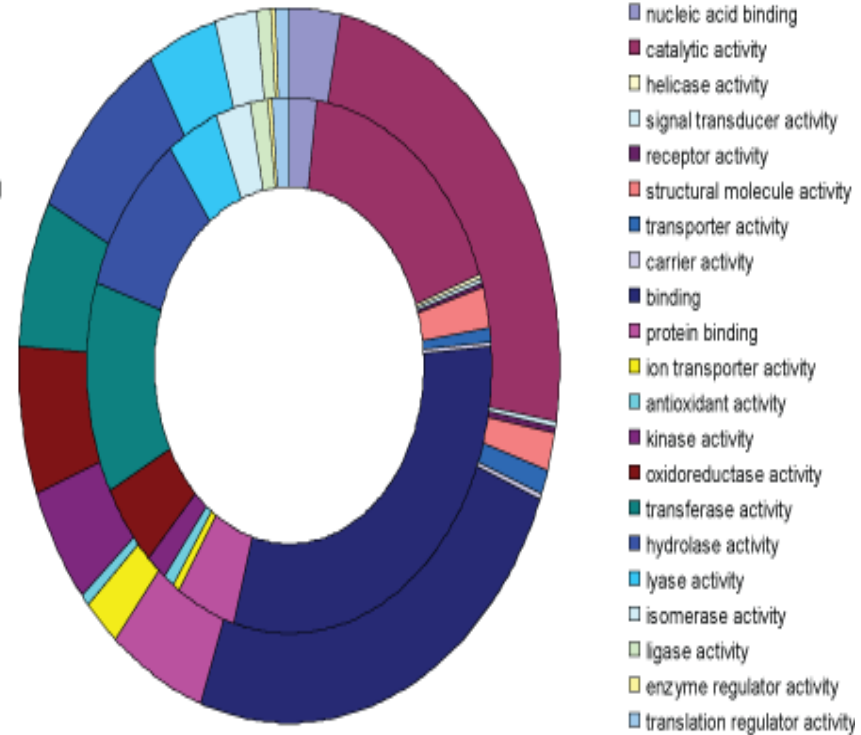
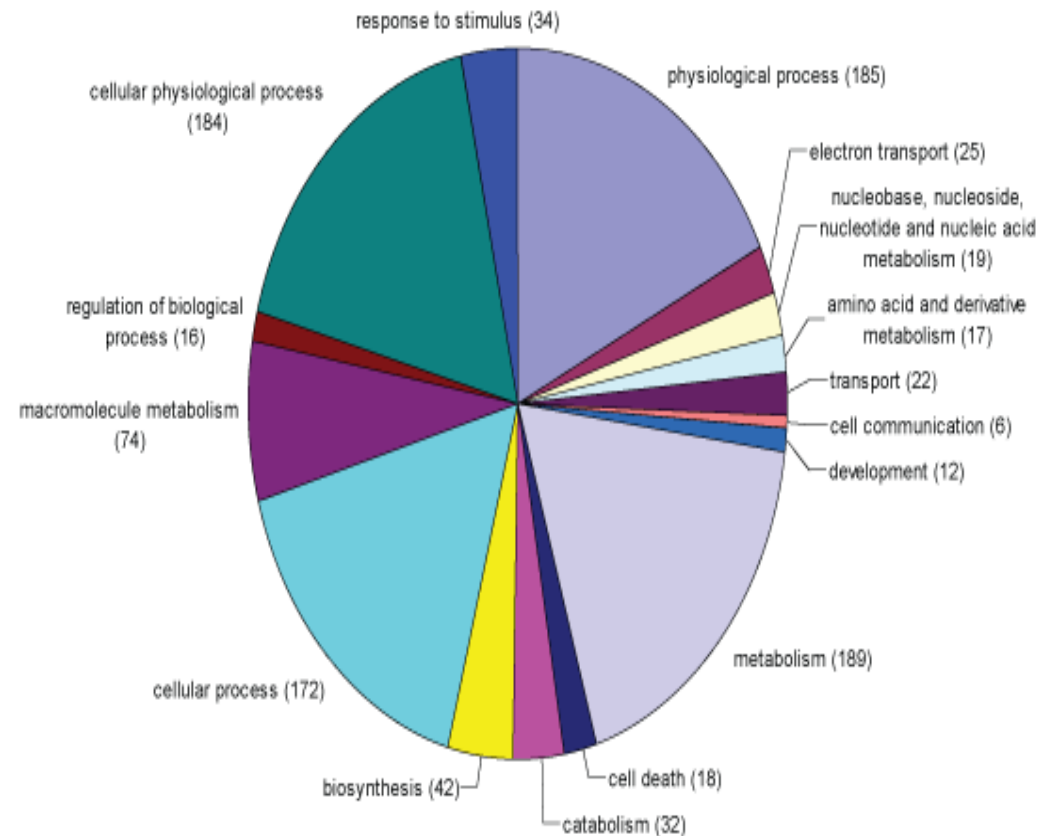
- Epstein & Anfinsen, 1961:  
sequence uniquely determines structure

- INPUT: sequence

- OUTPUT:

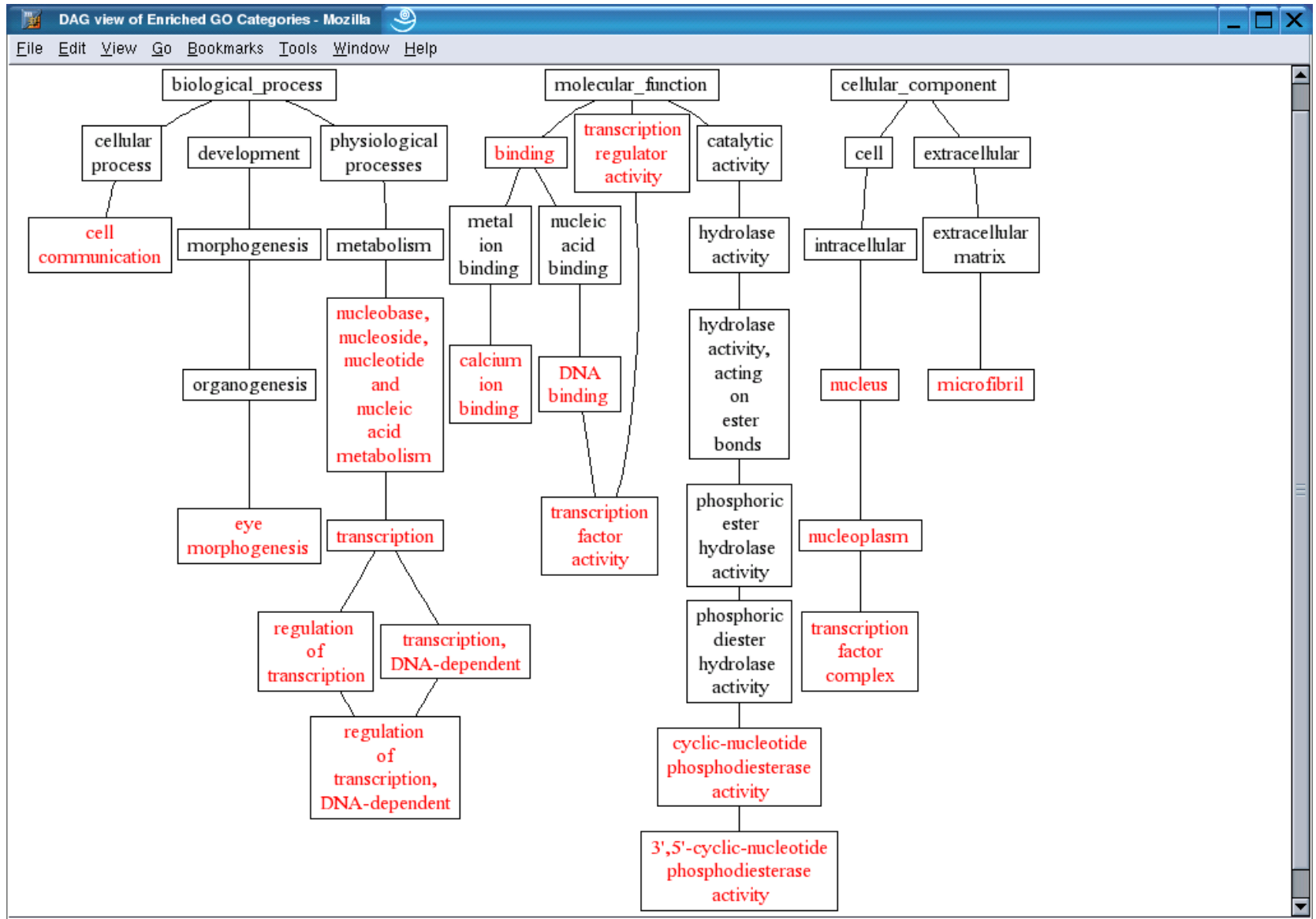


# Topic 4: Protein Function Prediction

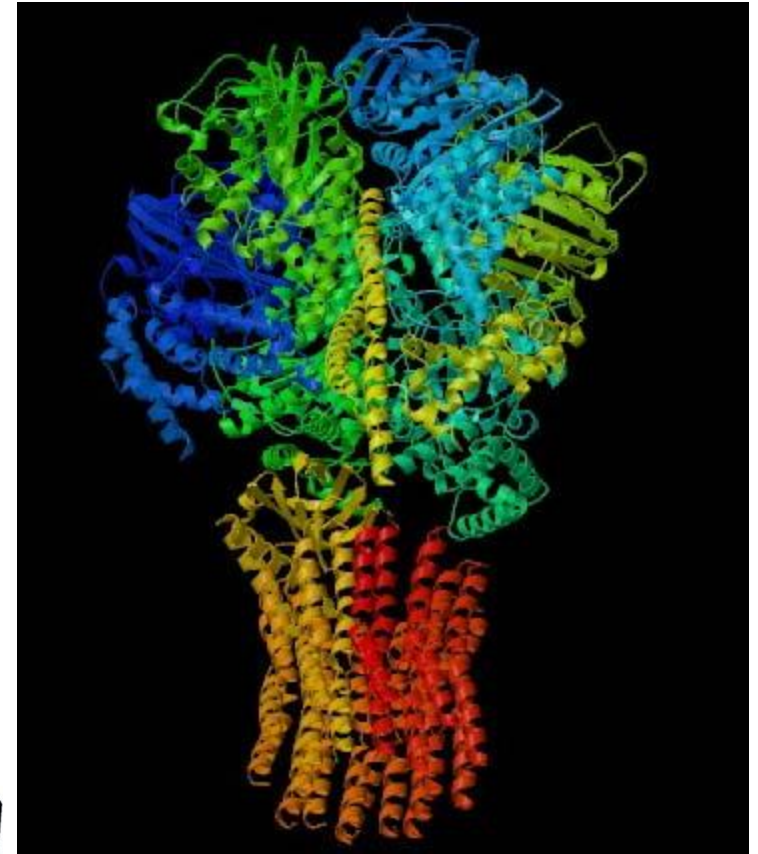
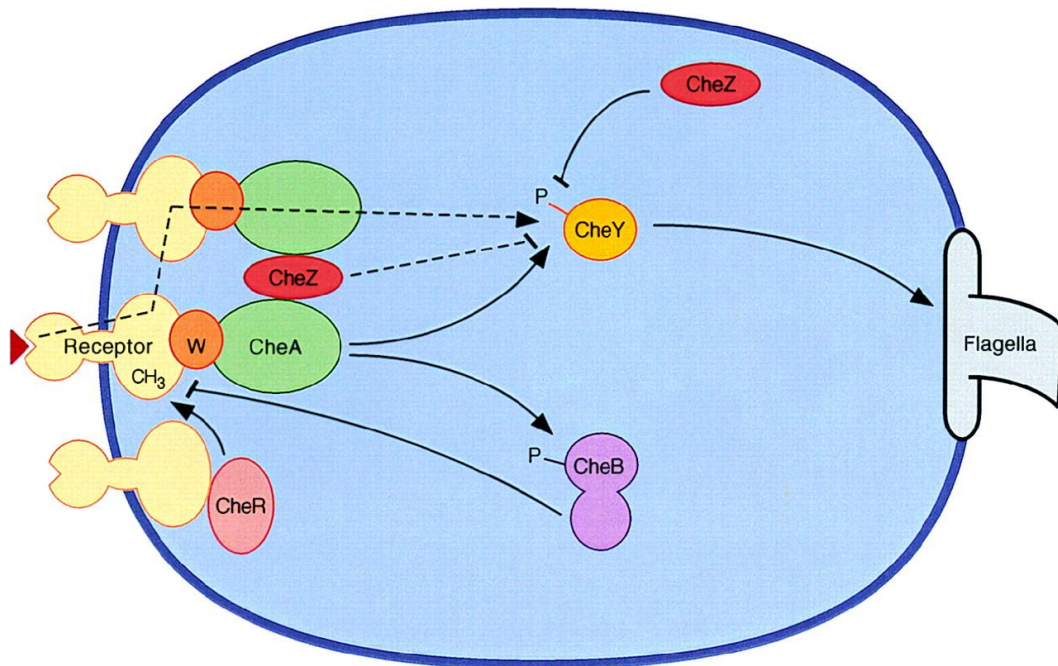
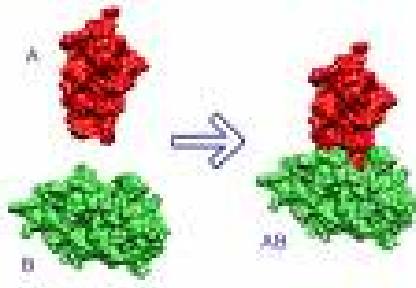


- nucleic acid binding
- catalytic activity
- helicase activity
- signal transducer activity
- receptor activity
- structural molecule activity
- transporter activity
- carrier activity
- binding
- protein binding
- ion transporter activity
- antioxidant activity
- kinase activity
- oxidoreductase activity
- transferase activity
- hydrolase activity
- lyase activity
- isomerase activity
- ligase activity
- enzyme regulator activity
- translation regulator activity

# Gene Ontology

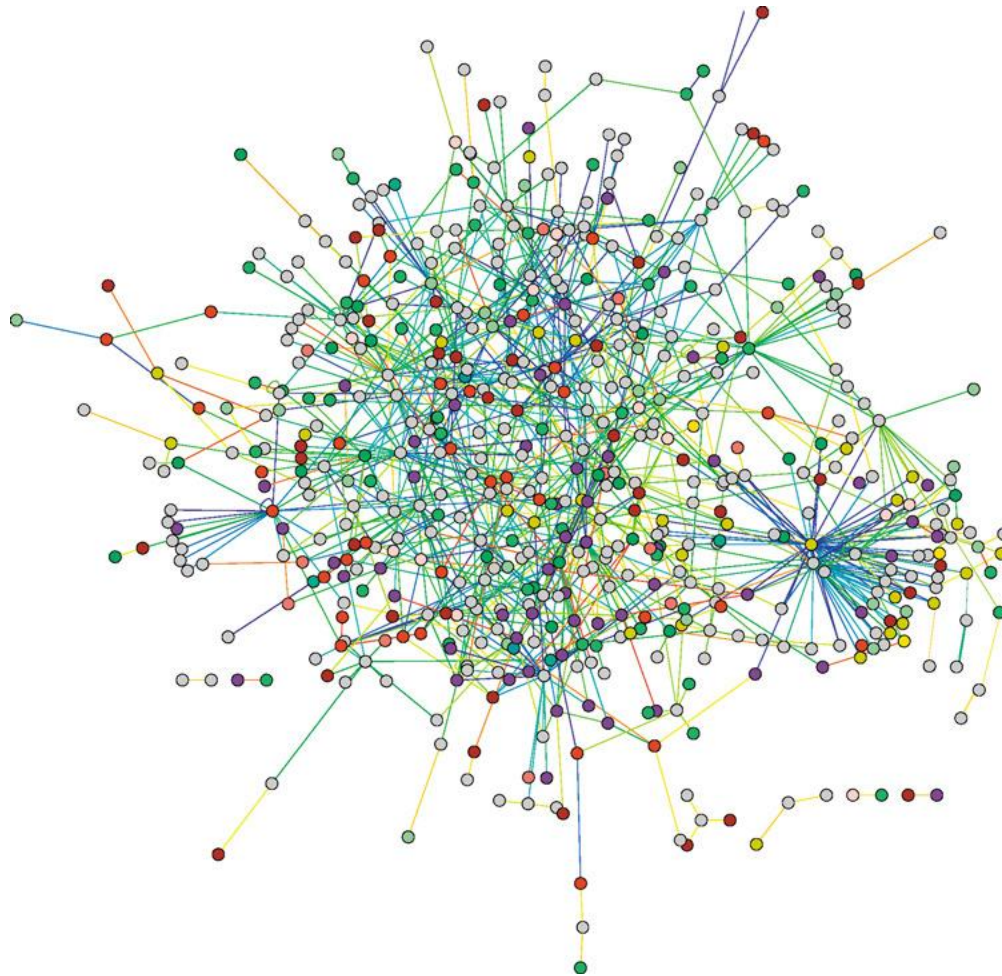


# Topic 5: Protein-Protein Interaction Prediction



ATP Synthase

# Protein Interaction Network

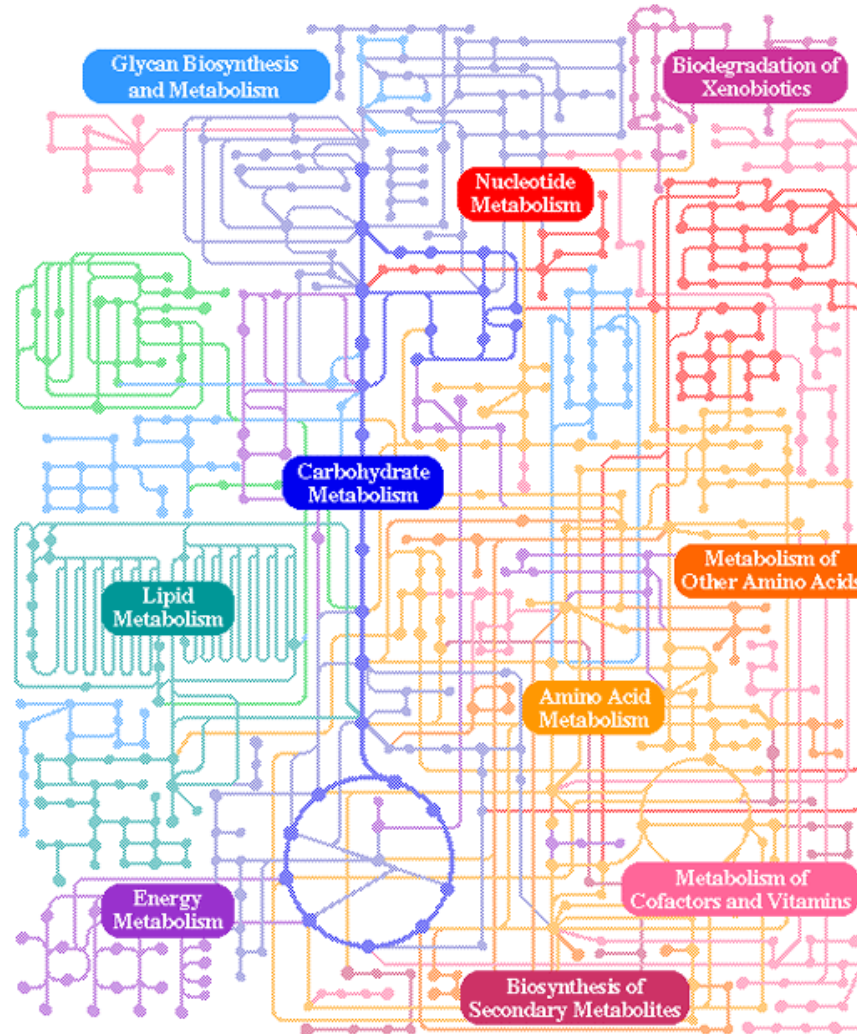


# Topic 6: Reconstruction of Biological Pathway and Networks

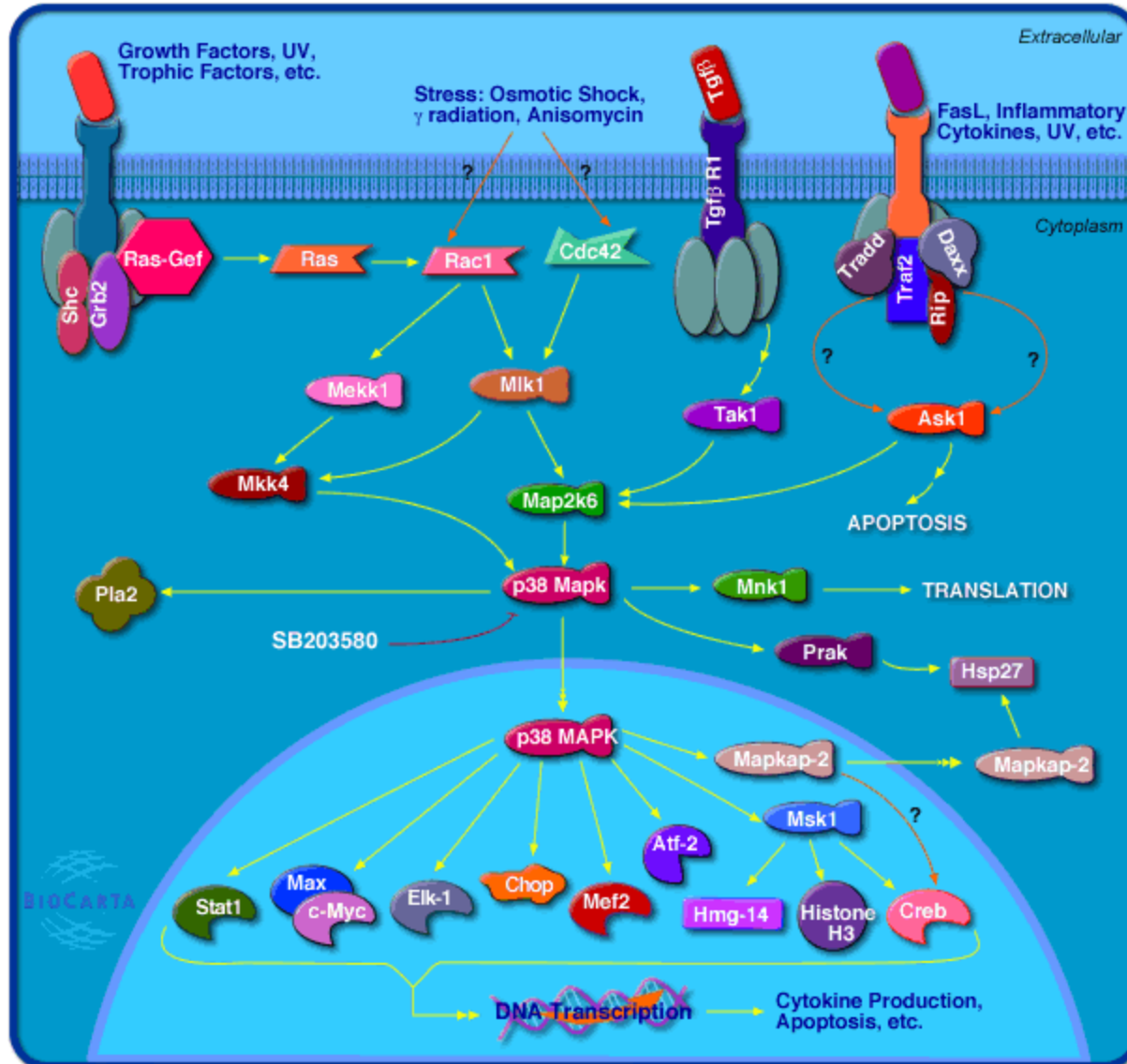
- Metabolic pathway
- Signal transduction pathway
- Gene regulatory pathway



# Metabolic Pathway (KEGG)

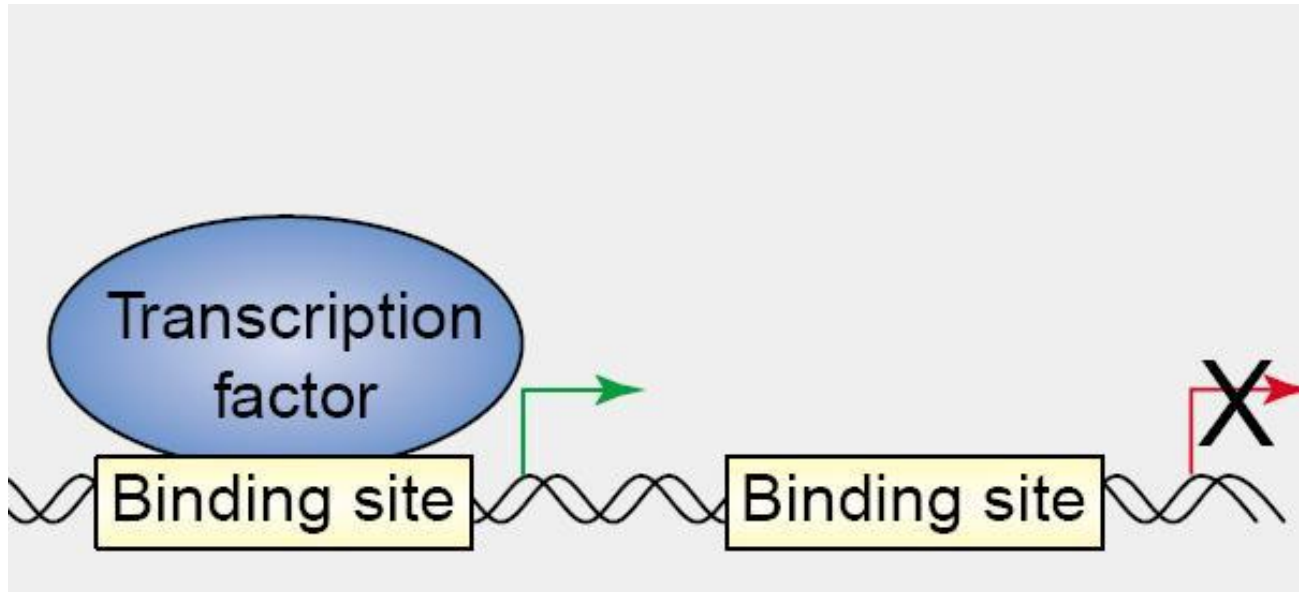


# Signal Transduction Pathway



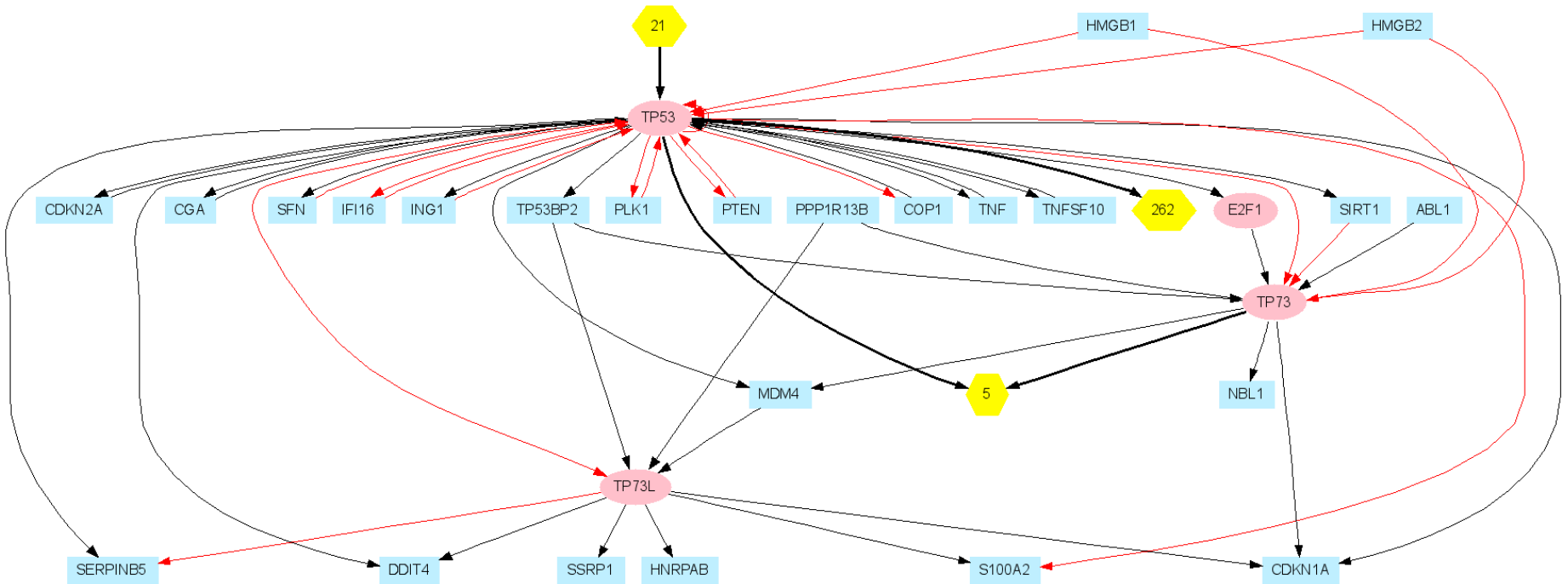


# Gene Regulatory Pathway



# Gene Regulatory Network

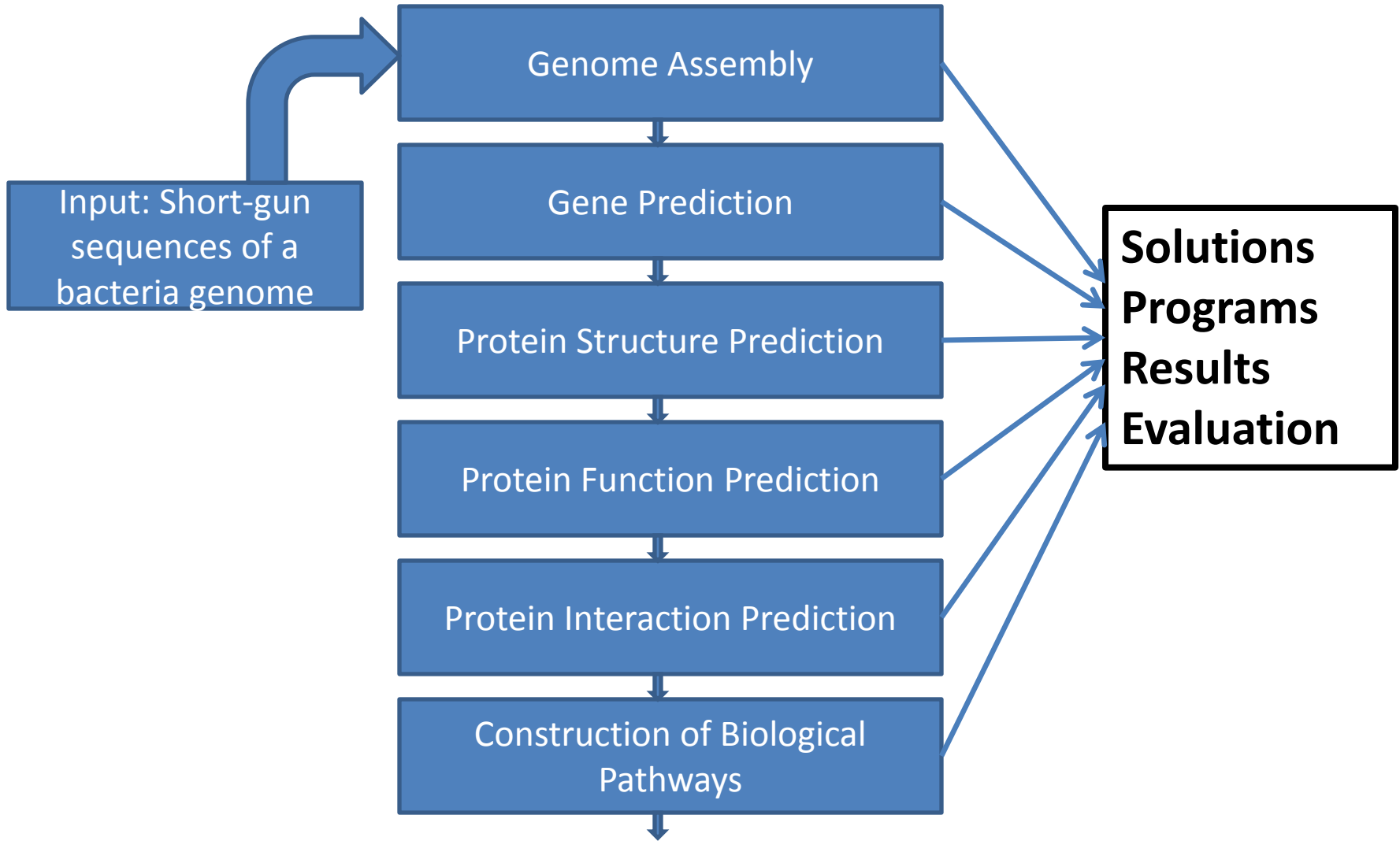
## Gene Regulatory Network of TF family p53 in human



**A lot of techniques and challenges,  
how can we get it done in one  
semester?**

Novel learning technique: doing one  
genome assembly and annotation  
project in six steps

# Group Project



# Reading Assignment

**J.C. Venter et al.. The sequence of the Huan Genome. Science. 291:1304, 2001**

**Read: Introduction and first three sections:**

**<http://www.sciencemag.org/cgi/reprint/291/5507/1304.pdf>**

**Write a review (one to two pages) to summarize the main problems, methods and results**