**Jianlin Cheng, PhD**
Informatics Institute, Computer Science Department
University of Missouri, Columbia
Fall, 2011

# Objectives

- Walk students through the complete process of sequencing, assembling and annotating a genome. During the process, students lean 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)
- 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 (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 presentation

- **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. (chengji@missouri.edu)

# Evaluation and Grading

- literature reviews (individual, 10%), mid-term presentation (individual, 10%), class discussion (individual, 15%), topic presentations (group, 15%), topic plans and reports (i.e. progress and assessment) (group, 15%), topic implementation (group, 20%), a final presentation and report (group, 15%)
- Group components may be 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

# Introduction



- **Grow**
- **Sustain**
- **Adapt**
- **Reproduce**

- **Genome & Components**
- **Environment**

http://www.scq.ubc.ca/wp-content/uploads/2006/08/molecular-machine.gif

# Applications of Genome Knowledge

# Genome Sequencing – Cracking the Code

- Virus & Bacteria genomes (small)
- Human genome



Figure 1-41 Essential Cell Biology 3/e (© Garland Science 2010)

nucleotide pairs per haploid genome

# Human Genome Project

# Fun

# Genome Sequencing Routine

# MU Genome Sequencing

- Soybean genome: Gary Stacey, Dong Xu, Jay Thelen, Jianlin Cheng, Henry Nguyen, etc (Nature, 2010)
- Chris Pires – plant genomes

# Sequencing process

# Genome Sequencing Machine

# STRATEGIES FOR SEQUENCING THE HUMAN GENOME

## BY MAPPED CLONES

## BY WHOLE GENOME SHOTGUN

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).

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

# Protein Sequence, Structure, Function

AGCWY……



**Cell**

# Protein Structure Space

# Protein Data Bank
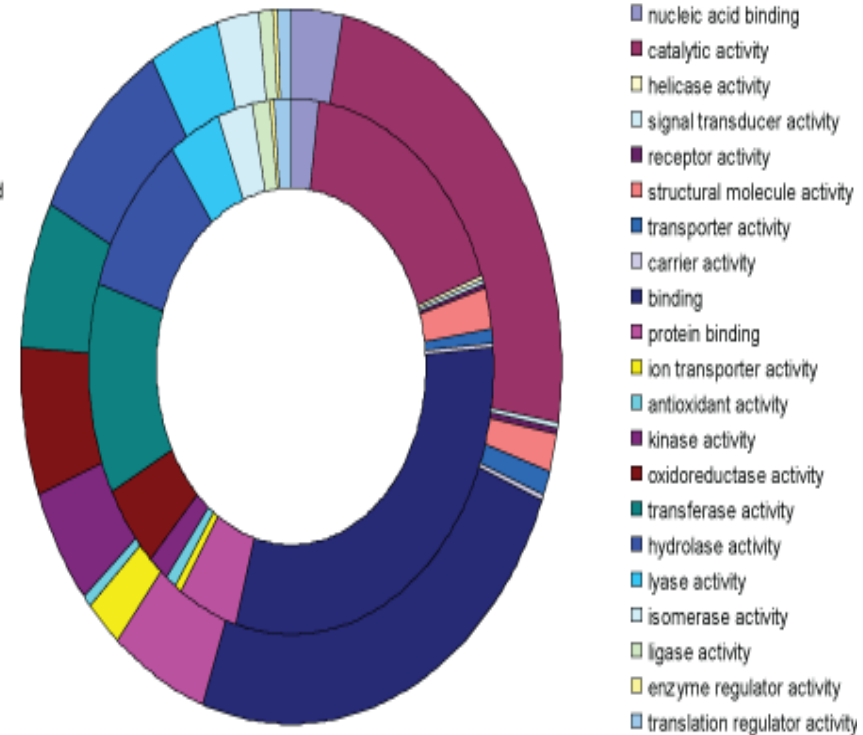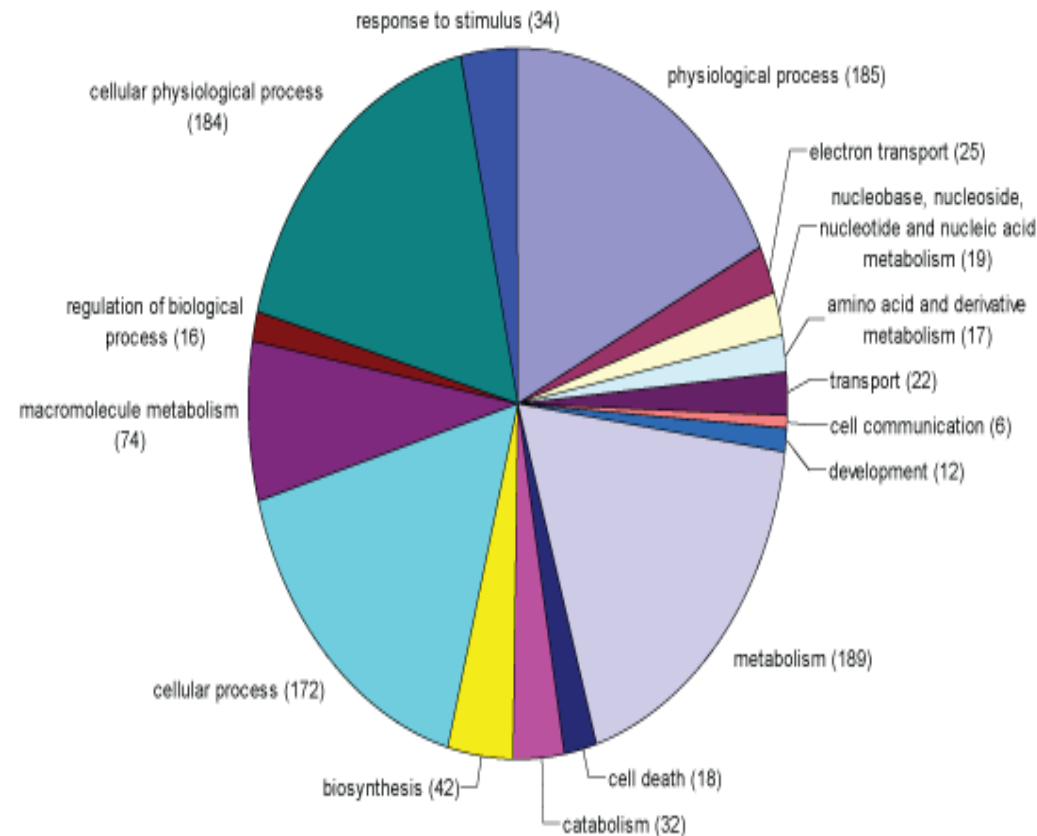
# Topic 3: Protein structure prediction

- Epstein & Anfinsen, 1961:
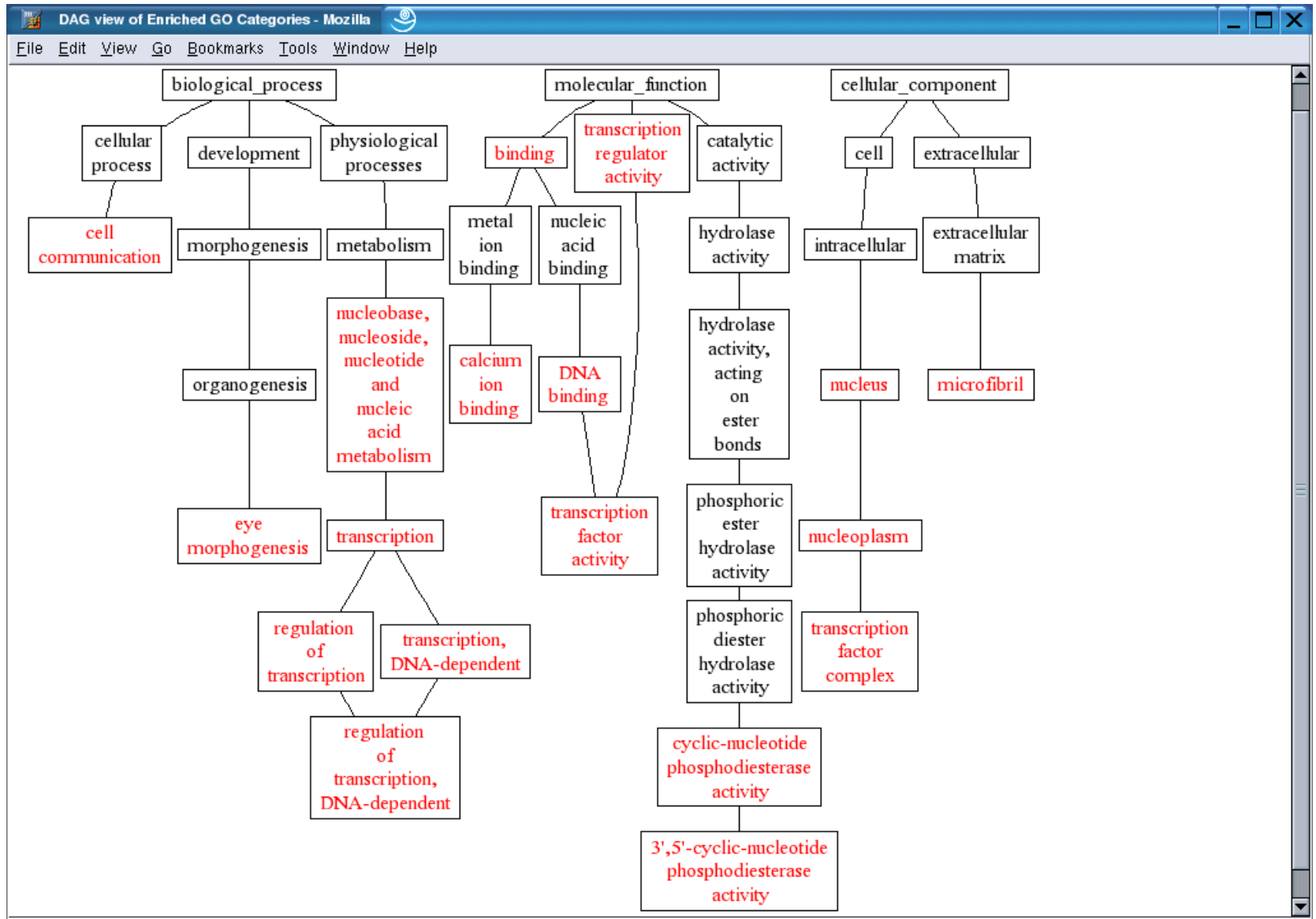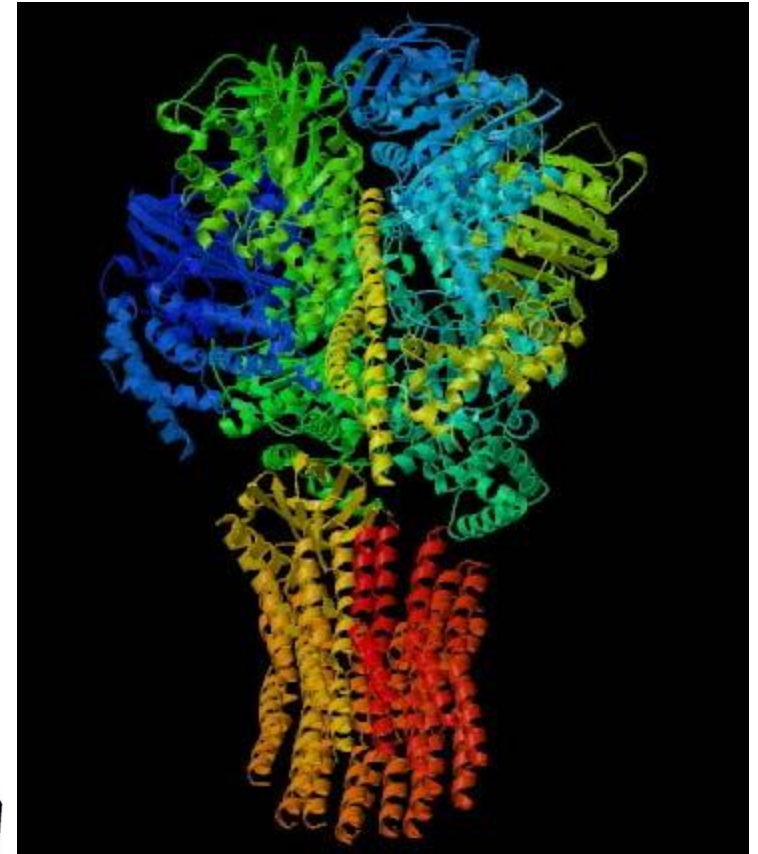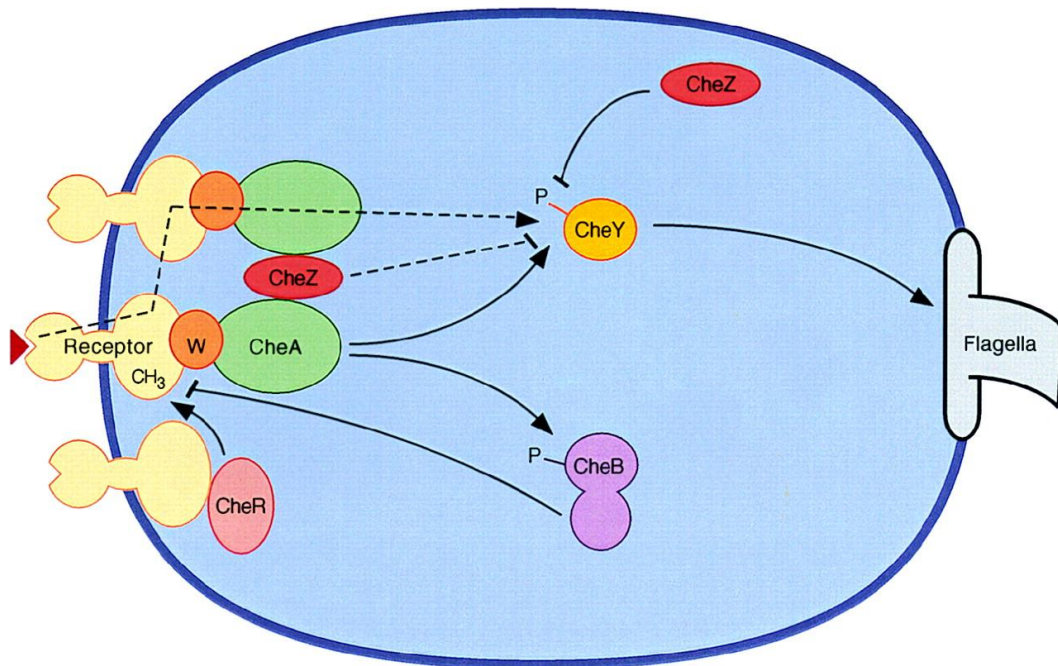  sequence uniquely determines structure
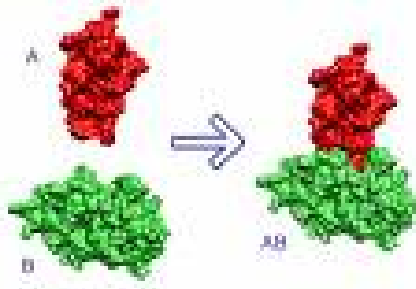
- INPUT:          sequence
- OUTPUT:

*3D structure and function*

B. Rost, 2005

# Topic 4: Protein Function Prediction
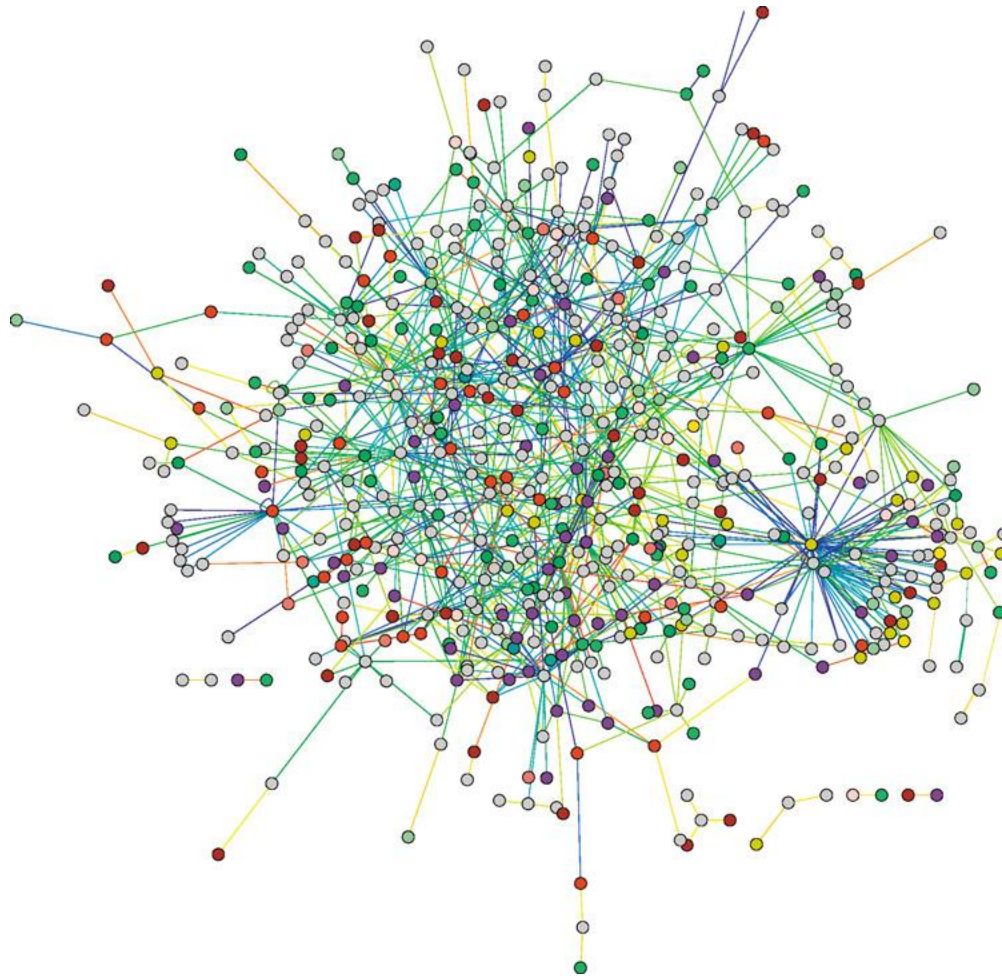
# Gene Ontology

# Topic 5: Protein-Protein Interaction Prediction
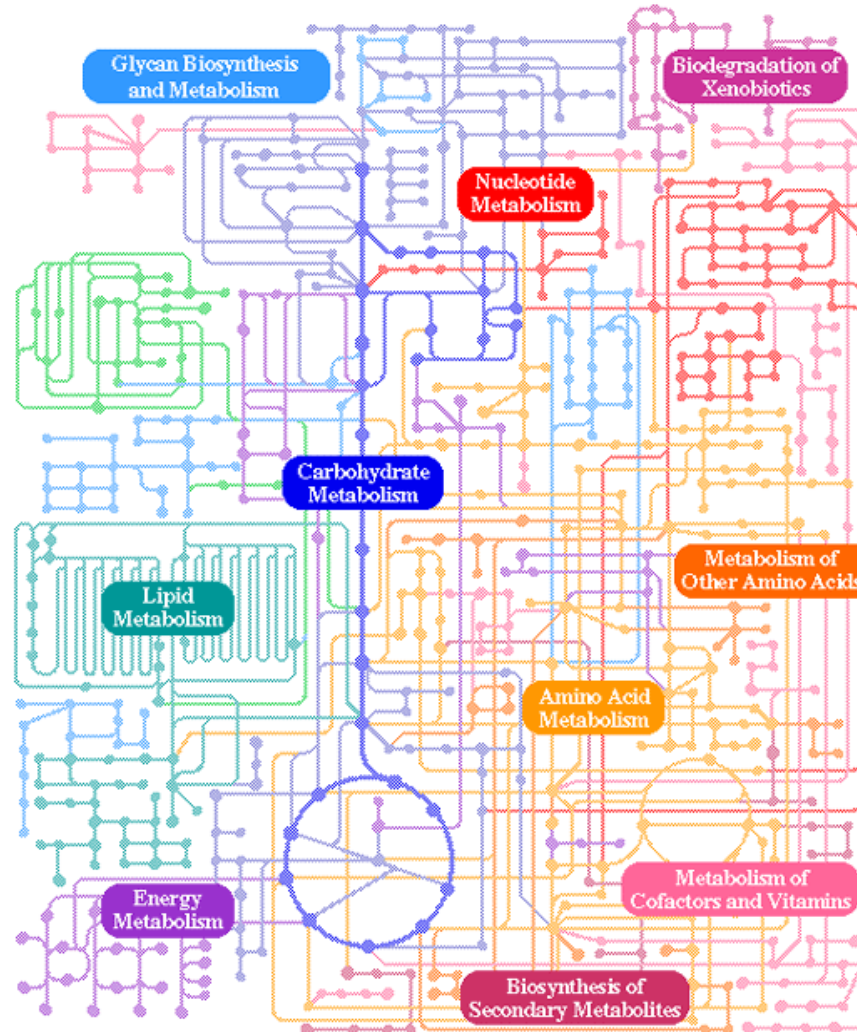






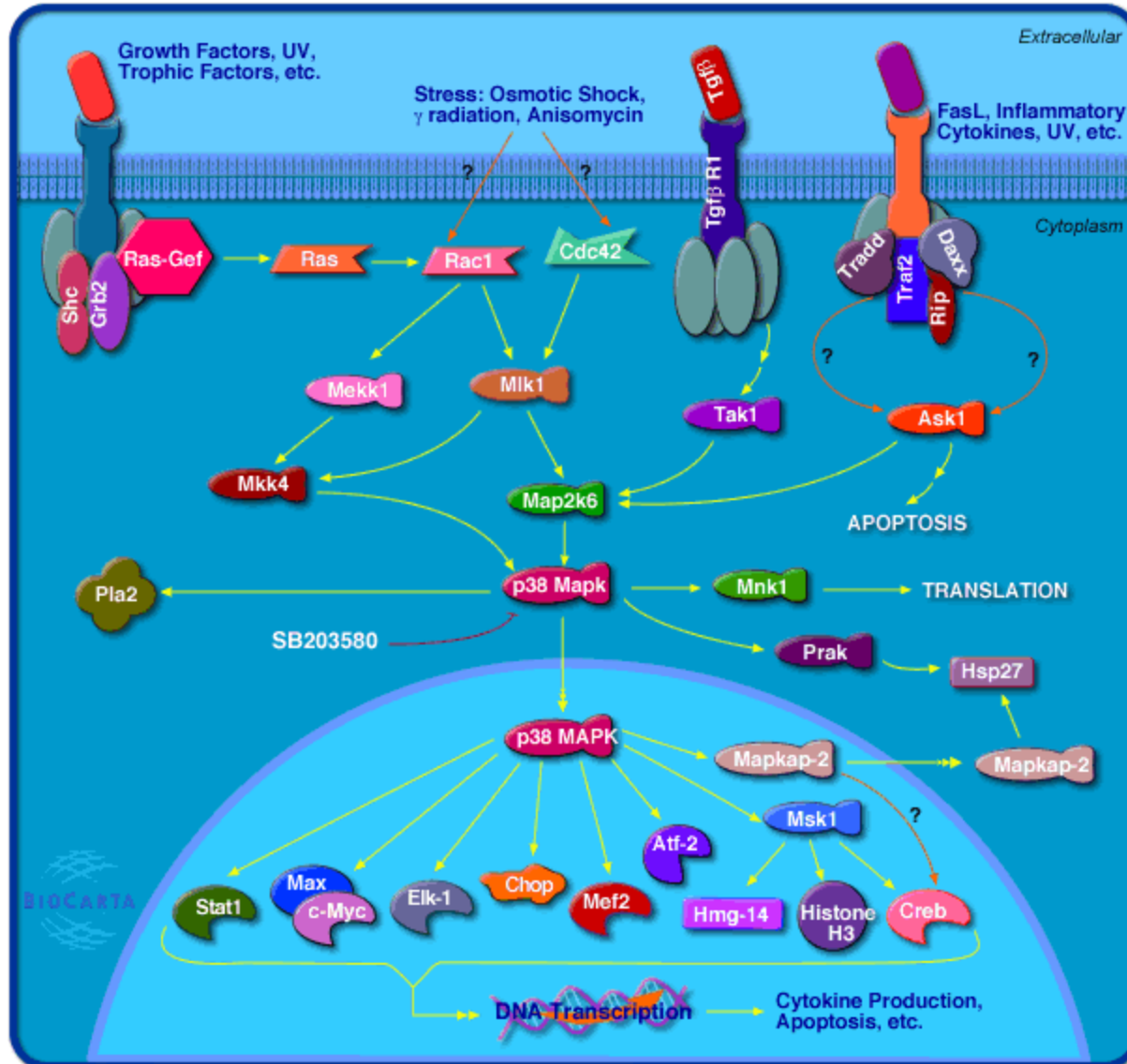**ATP Synthase**

# Protein Interaction Network

# Topic 6: Reconstruction of Biological Pathway and Networks

- Metabolic pathway
- Signal transduction pathway
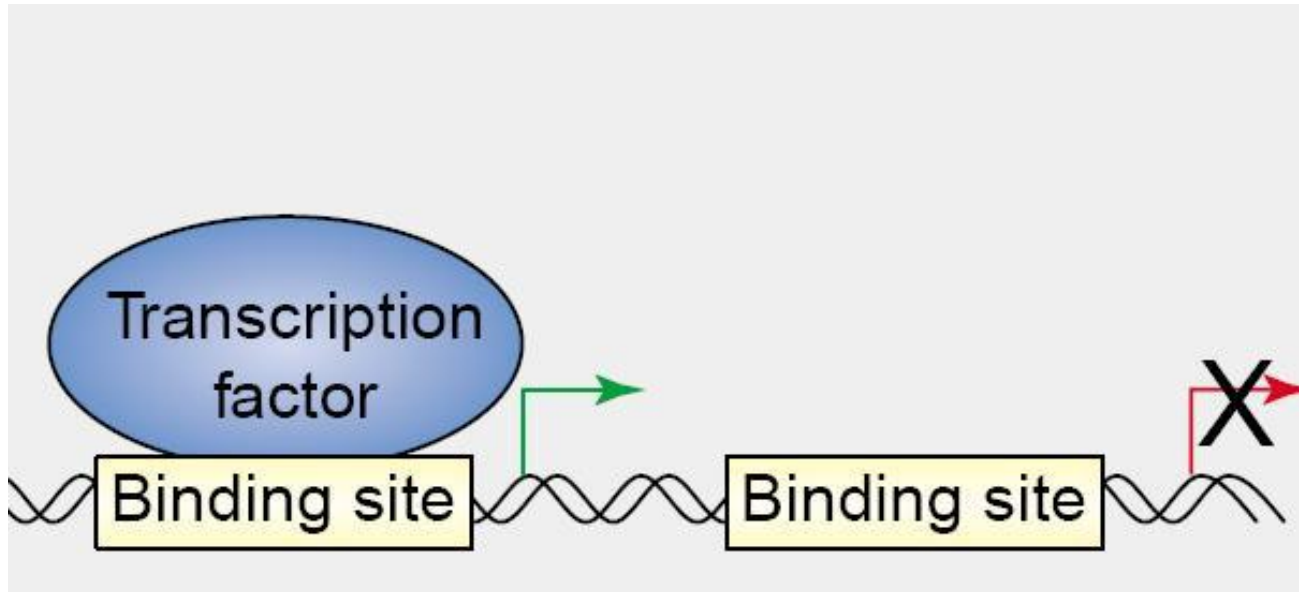- Gene regulatory pathway

# Metabolic Pathway (KEGG)



01100  5/31/04  Image source from KEGG
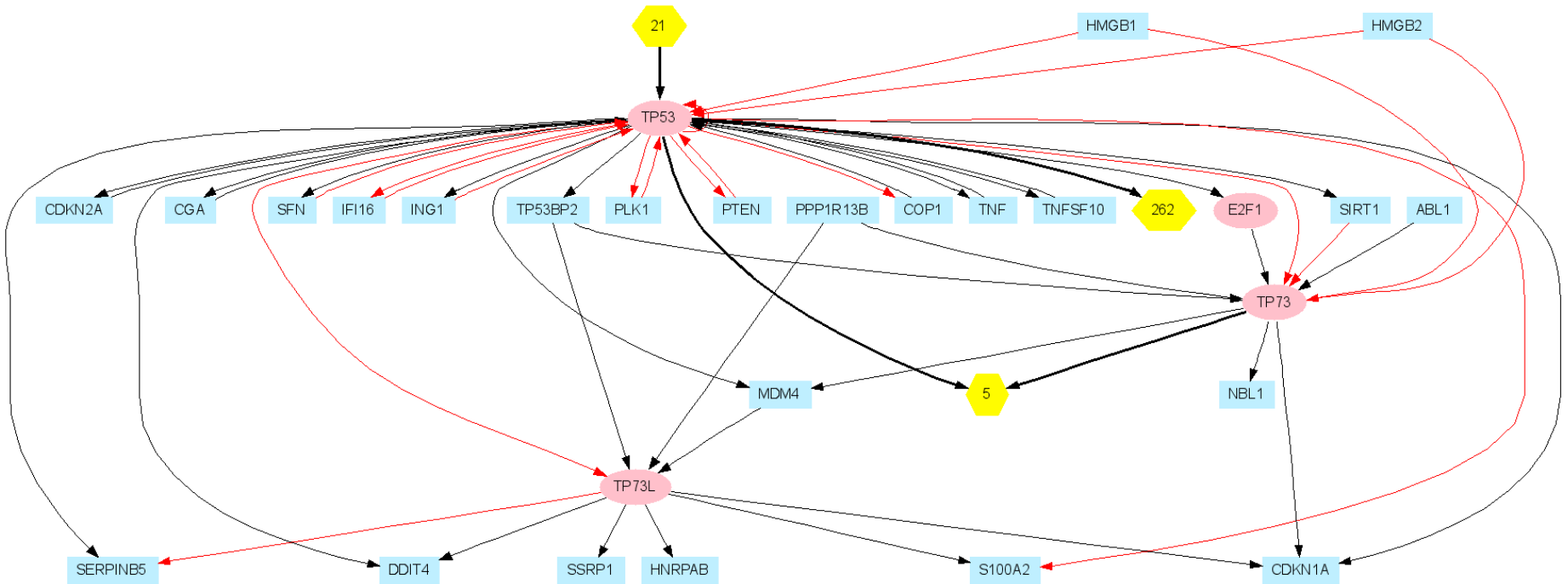
# Signal Transduction Pathway

# Gene Regulatory Pathway
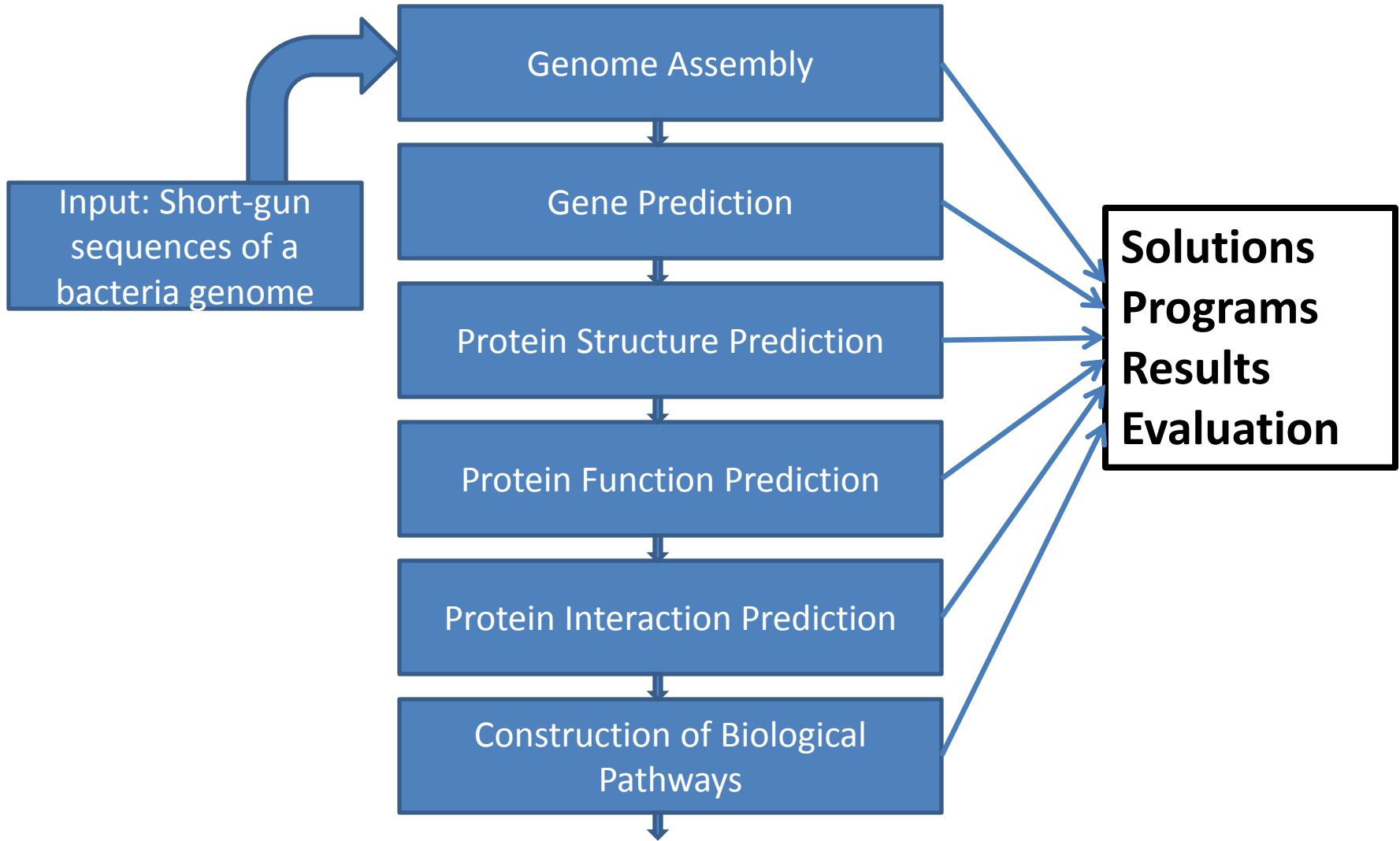
# Gene Regulatory Network

**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 page) to summarize the main problems, methods and results**