# Proteomics & Bioinformatics Part II
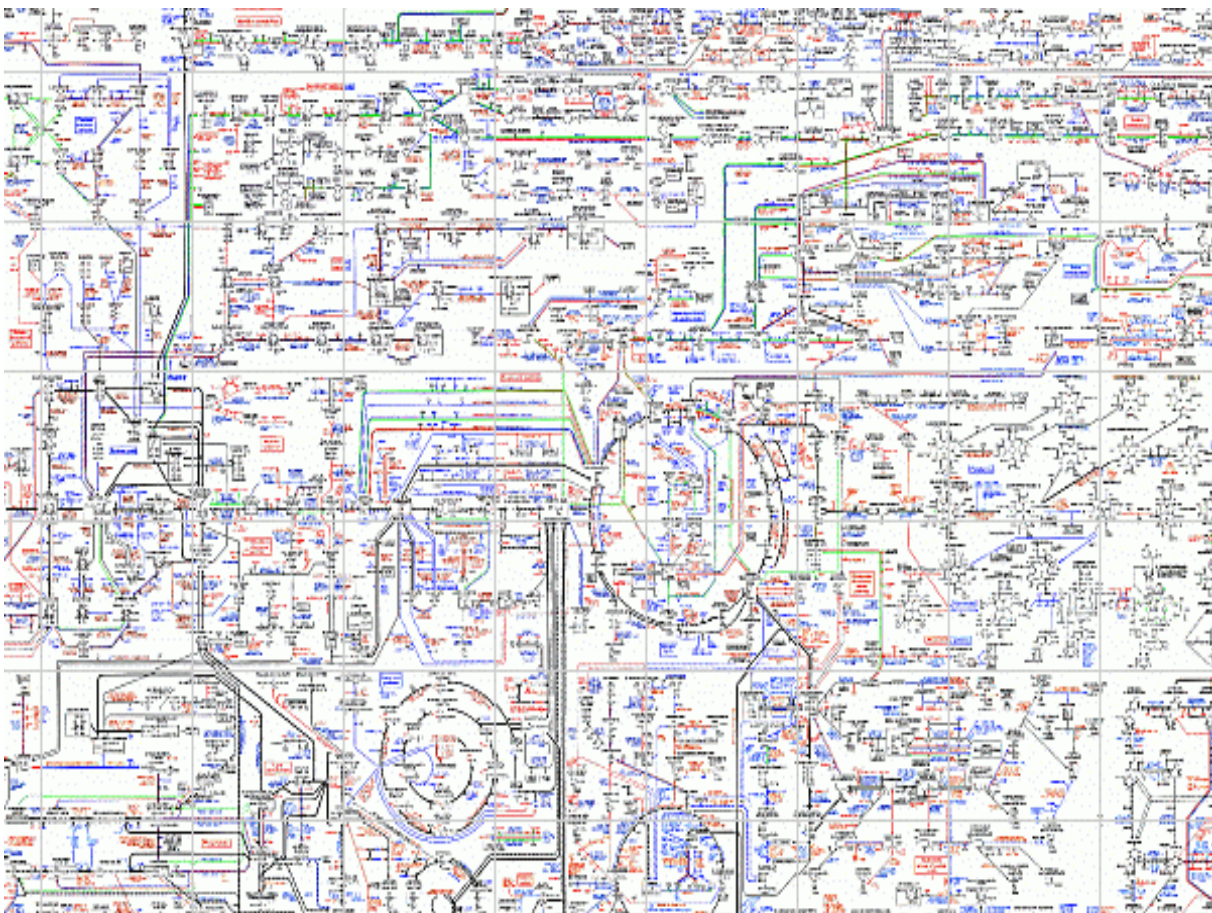
**David Wishart**

**University of Alberta**

---

# 3 Kinds of Proteomics

- **Structural Proteomics**
  - **High throughput X-ray Crystallography/Modelling**
  - **High throughput NMR Spectroscopy/Modelling**
- **Expressional or Analytical Proteomics**
  - **Electrophoresis, Protein Chips, DNA Chips, 2D-HPLC**
  - **Mass Spectrometry, Microsequencing**
- **Functional or Interaction Proteomics**
  - **HT Functional Assays, Ligand Chips**
  - **Yeast 2-hybrid, Deletion Analysis, Motif Analysis**
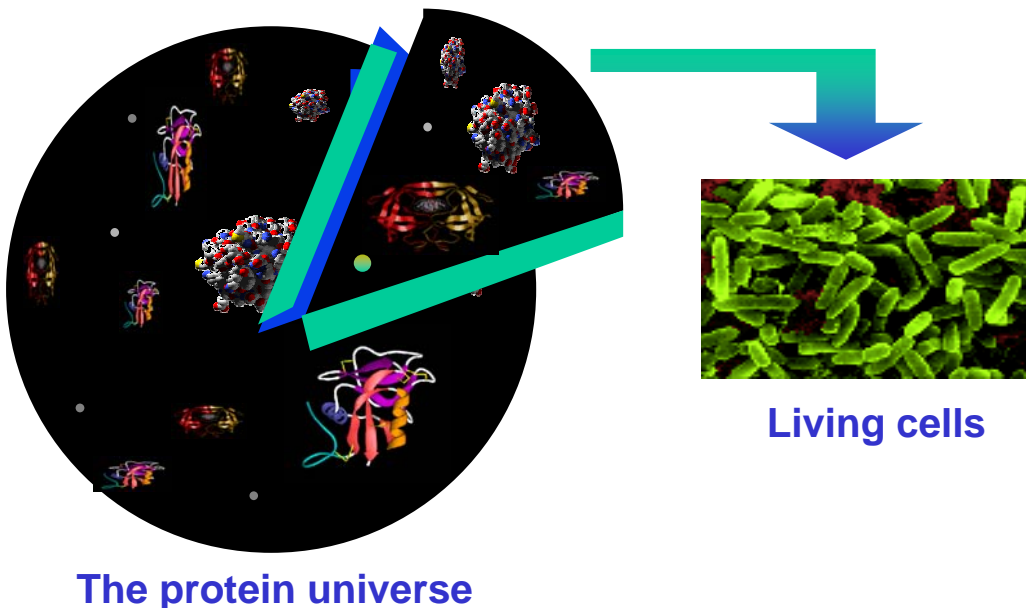
# Historically...

- **Most of the past 100 years of biochemistry has focused on the analysis of small molecules (i.e. metabolism and metabolic pathways)**

- **These studies have revealed much about the processes and pathways for about 400 metabolites which can be summarized with this...**
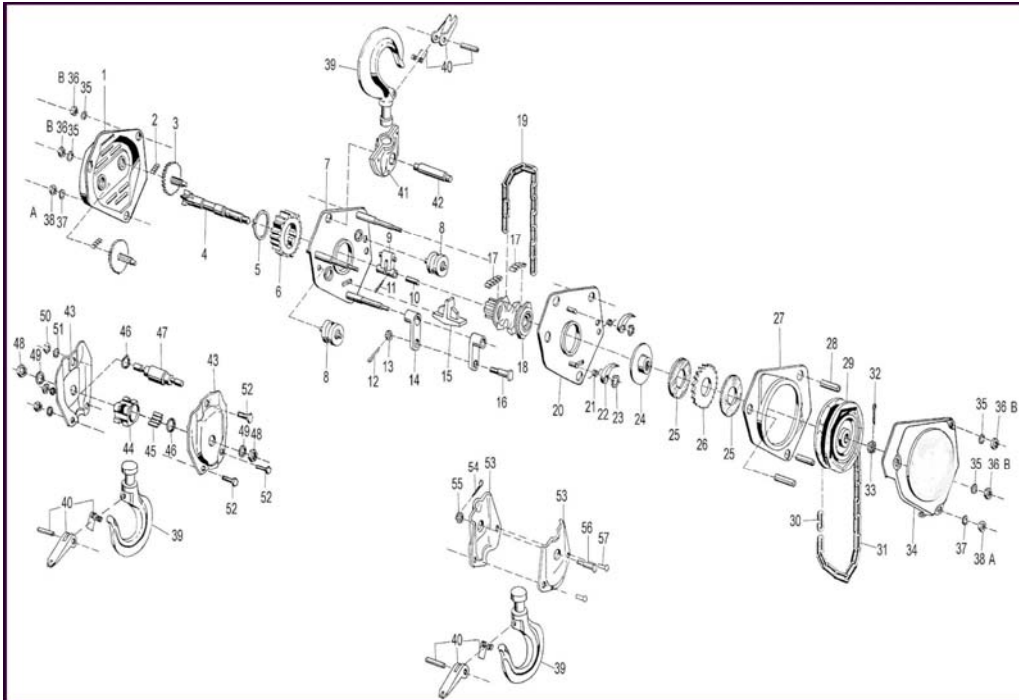
# More Recently...

- **Molecular biologists and biochemists have focused on the analysis of larger molecules (proteins and genes) which are much more complex and much more numerous**
- **These studies have primarily focused on identifying and cataloging these molecules (Human Genome Project)**
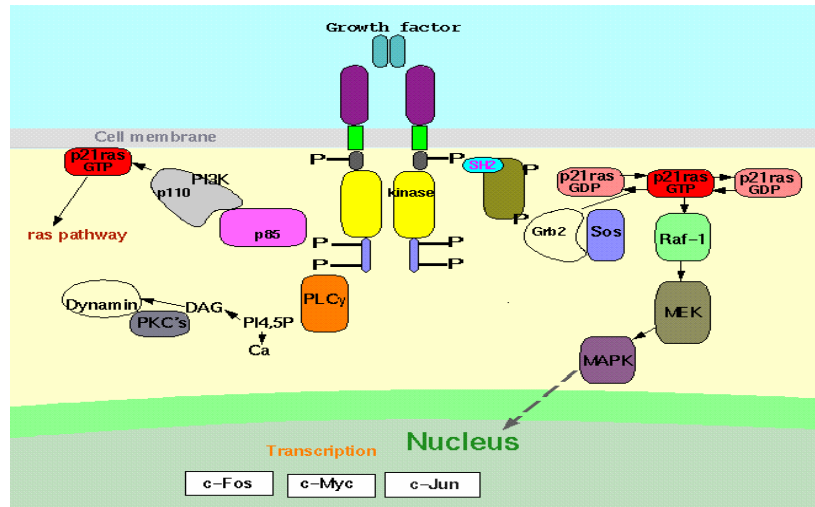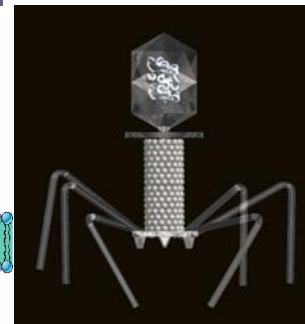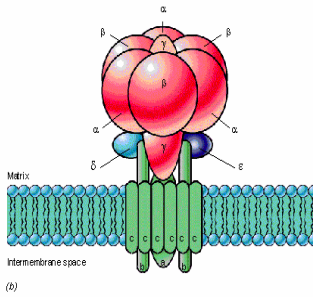
# Nature's Parts Warehouse



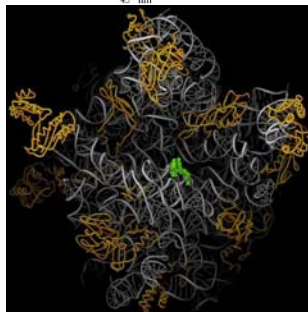**The protein universe**
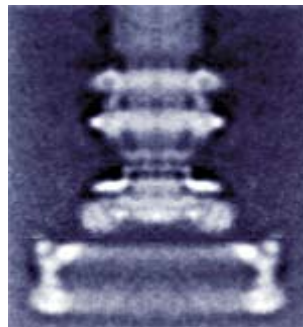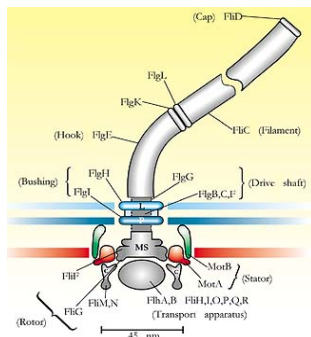
**Living cells**

# The Protein Parts List



# However...

- **This cataloging (which consumes most of bioinformatics) has been derogatively referred to as "stamp collecting"**
- **Having a collection of parts and names doesn't tell you how to put something together or how things connect --** *this is biology*
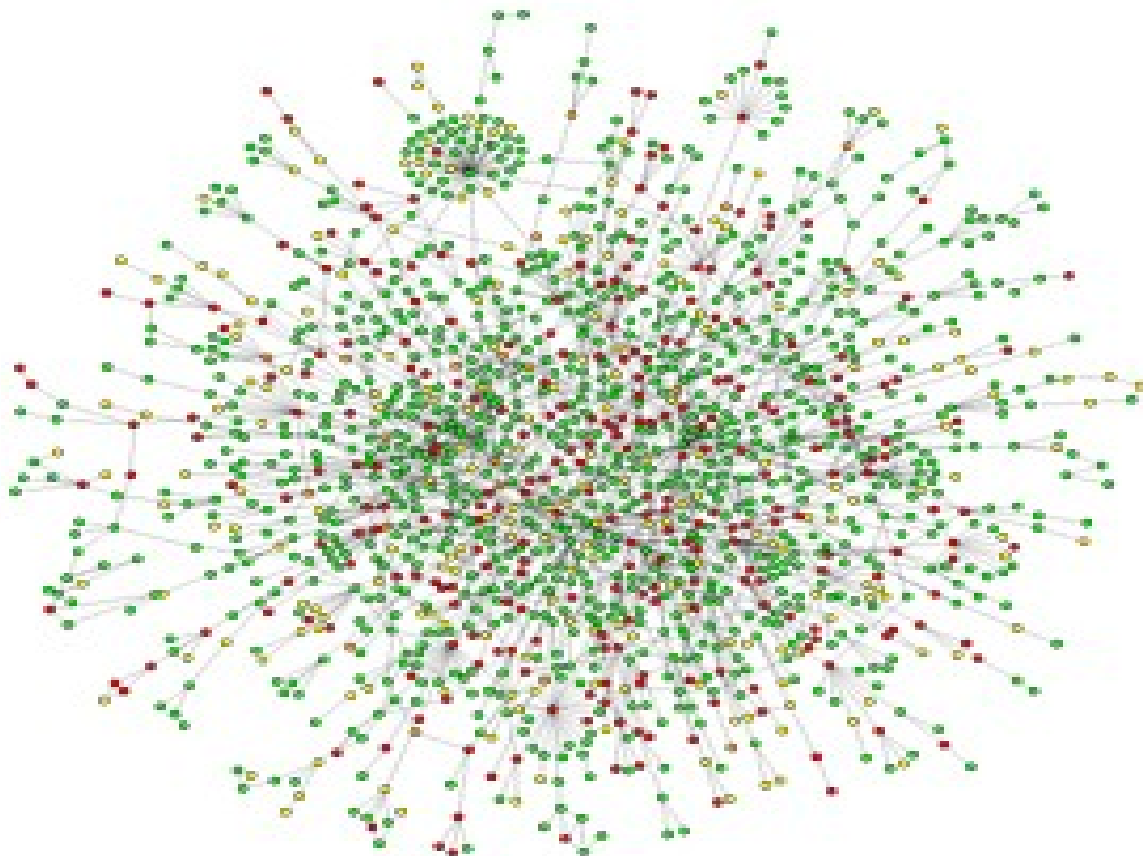
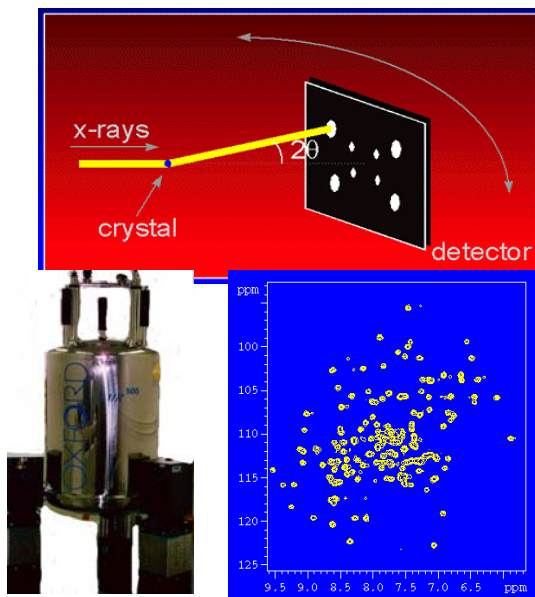# Remember: *Proteins Interact*



# Proteins *Assemble*

# For the Past 10 Years...

- **Scientists have increasingly focused on "signal transduction" and transient protein interactions**
- **New techniques have been developed which reveal which proteins and which parts of proteins are important for interaction**
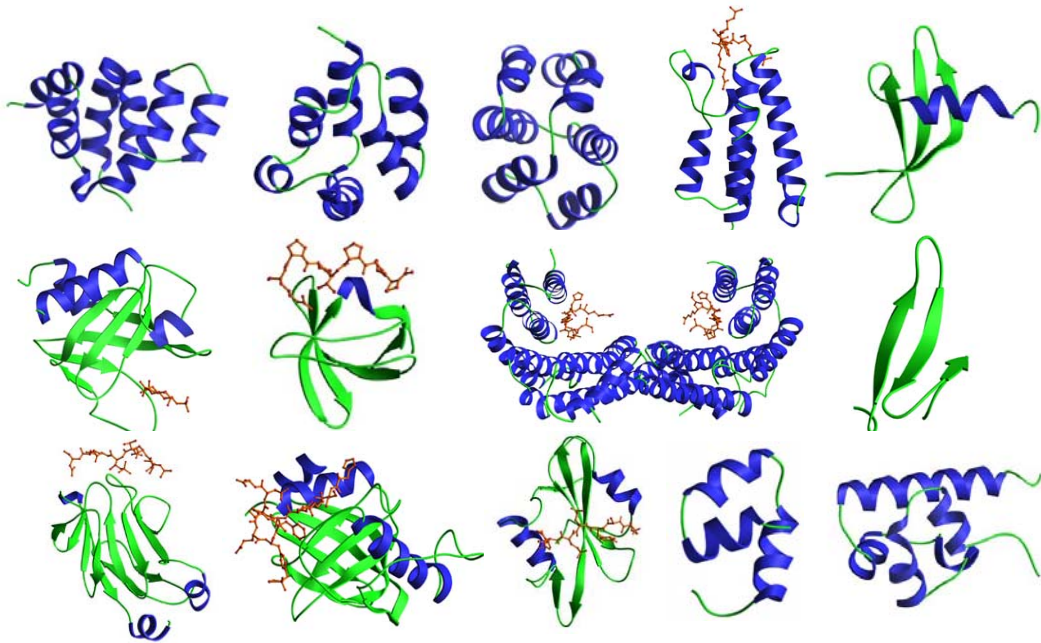- **The hope is to get something like this..**

# Protein Interaction Tools and Techniques - Experimental Methods
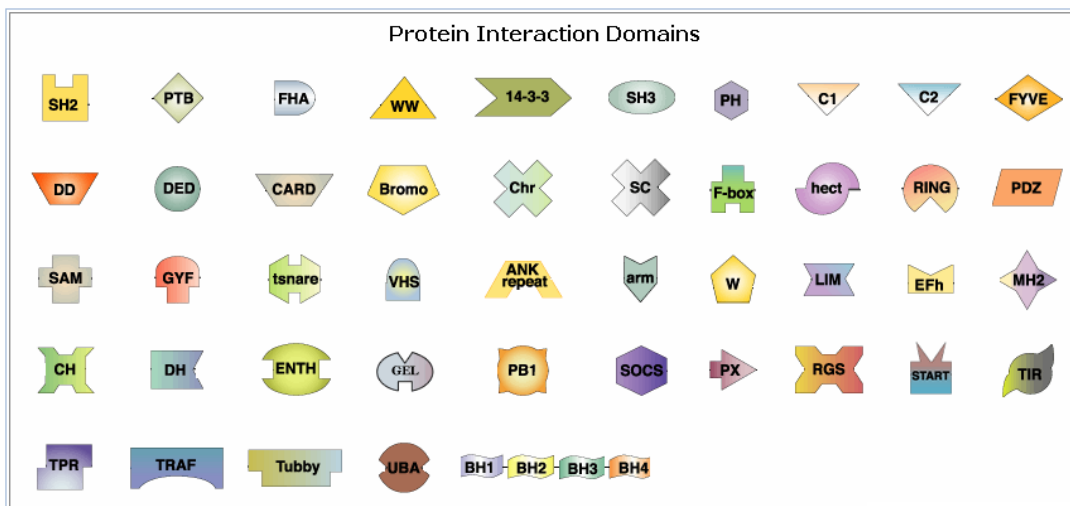
---

# 3D Structure Determination



- **X-ray crystallography**
  - **grow crystal**
  - **collect diffract. data**
  - **calculate e- density**
  - **trace chain**
- **NMR spectroscopy**
  - **label protein**
  - **collect NMR spectra**
  - **assign spectra & NOEs**
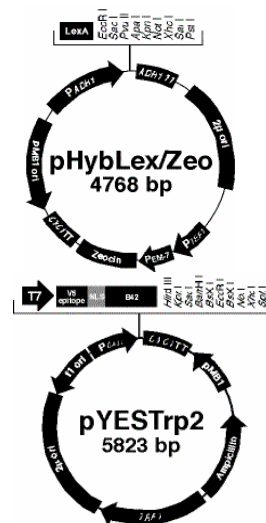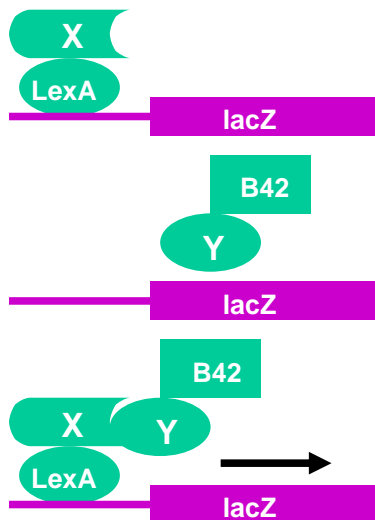  - **calculate structure using distance geom.**

# Protein Interaction Domains



**http://www.mshri.on.ca/pawson/domains.html**

# Protein Interaction Domains



http://www.mshri.on.ca/pawson/domains.html

# Yeast Two-Hybrid Analysis



- Yeast two-hybrid experiments yield information on protein protein interactions
- GAL4 Binding Domain
- GAL4 Activation Domain
- X and Y are two proteins of interest
- If X & Y interact then reporter gene is expressed

# Invitrogen Yeast 2-Hybrid

# Example of 2-Hybrid Analysis

- **Uetz P. et al., "*A Comprehensive Analysis of Protein-Protein Interactions in Saccharomyces cerevisiae*" Nature 403:623-627 (2000)**
- **High Throughput Yeast 2 Hybrid Analysis**
- **957 putative interactions**
- **1004 of 6000 predicted proteins involved**

# Example of 2-Hybrid Analysis

- **Rain JC. et al., "*The protein-protein interaction map of Helicobacter pylori*" Nature 409:211-215 (2001)**
- **High Throughput Yeast 2 Hybrid Analysis**
- **261 H. pylori proteins scanned against genome**
- **>1200 putative interactions identified**
- **Connects >45% of the H. pylori proteome**

# Another Way?

- **Ho Y, Gruhler A, et al. *Systematic identification of protein complexes in Saccharomyces cerevisiae by mass spectrometry*. Nature 415:180-183 (2002)**

- **High Throughput Mass Spectral Protein Complex Identification (HMS-PCI)**

- **10% of yeast proteins used as "bait"**

- **3617 associated proteins identified**

- **3 fold higher sensitivity than yeast 2-hybrid**

# Affinity Pull-down

# Transposon Tagging



# Protein Arrays



H Zhu, J Klemic, S Chang, P Bertone, A Casamayor, K Klemic, D Smith, M Gerstein, M Reed, & M Snyder (2000).**Analysis of yeast protein kinases using protein chips**. Nature Genetics 26: 283-289

# Protein Arrays



# Protein Interaction Tools and Techniques - Computational Methods

# Sequence Searching Against Known Domains



Protein Interaction Domains

**http://www.mshri.on.ca/pawson/domains.html**

# Motif Searching Using Known Motifs



Protein Interaction Domains

# Text Mining

- **Searching Medline or Pubmed for words or word combinations**
- **"X binds to Y"; "X interacts with Y"; "X associates with Y" etc. etc.**
- **Requires a list of known gene names or protein names for a given organism**
- **Sometimes called "Textomy"**



**http://textomy.iit.nrc.ca/**

# Pre-BIND

- *Donaldson et al. BMC Bioinformatics 2003 4:11*
- **Used Support Vector Machine (SVM) to scan literature for protein interactions**
- **Precision, accuracy and recall of 92% for correctly classifying PI abstracts**
- **Estimated to capture 60% of all abstracted protein interactions for a given organism**

# Rosetta Stone Method

Monomeric proteins that are fused in other organisms tend to be functionally related and physically interacting.

For example, using the Rosetta Stone™ method, it was found that human Nit and Fhit proteins are:

→ fused in invertebrates
→ form a heterocomplex in mammals

Human Nit       Human Fhit

Invertebrate NitFhit

Fhit

Nit

# Interologs, Homologs, Paralogs...

- **Homolog**
  - **Common Ancestors**
  - **Common 3D Structure**
  - **Common Active Sites**
- **Ortholog**
  - **Derived from Speciation**
- **Paralog**
  - **Derived from Duplication**

YM2

- **Interolog**
  - **Protein-Protein Interaction**



# Finding Interologs

- **If A and B interact in organism X, then if organism Y has a homolog of A (A') and a homolog of B (B') then A' and B' should interact too!**

- **Makes use of BLAST searches against entire proteome of well-studied organisms (yeast, E. coli)**

- **Requires list of known interacting partners**

# A Flood of Data

- **High throughput techniques are leading to more and more data on protein interactions**
- **This is where bioinformatics can play a key role**
- **Some suggest that this is the "future" for bioinformatics**

# Interaction Databases

- **BIND**
  - **http://www.blueprint.org/bind/bind.php**
- **DIP**
  - **http://dip.doe-mbi.ucla.edu/**
- **MINT**
  - **http://mint.bio.uniroma2.it/mint/**
- **PathCalling**
  - **http://portal.curagen.com/extpc/com.curagen.portal.servlet.Yeast**

# The BIND Database

- **BIND - Biomolecular Interaction Network Database**
- **Conceived and Developed by Chris Hogue, Tony Pawson, Francis Ouellette**
- **Designed to capture almost all interactions between biomolecules (large and small)**
- **Largest database of its kind**

# BIND Data Model

$$S \xrightarrow{\text{E}} P$$

$$E+S \rightleftharpoons E\text{-}S$$

*Interaction Record*

$$\longrightarrow P$$

*Chemical State Data*

$$S \rightleftharpoons P$$

*Chemical Action Data*

# BIND Can Encode...

- **Simple binary interactions**
- **Enzymes, substrates and conformational changes**
- **Restriction enzymes**
- **Limited proteolysis**
- **Phosphorylation (reversible)**
- **Glycosylation**
- **Intron splicing**
- **Transcriptional factors**

# BIND

# BIND Query Result



# BIND Details

# BIND Details



# BIND Details

# DIP Database of Interacting Proteins



**http://dip.doe-mbi.ucla.edu/**

# DIP Query Page

# DIP Results Page



click

# DIP Results Page

# MINT Molecular Interaction Database



http://mint.bio.uniroma2.it/mint/

# MINT Results

click

# KEGG Kyoto Encyclopedia of Genes and Genomes



**http://www.genome.ad.jp/kegg/kegg2.html**

# KEGG



# KEGG

# TRANSPATH



http://www.biobase.de/pages/products/transpath.html

# BIOCARTA

- **www.biocarta.com**
- **Go to "Pathways"**
- **Web interactive links to many signalling pathways and other eukaryotic protein-protein interactions**

# Other Databases



http://www.hgmp.mrc.ac.uk/GenomeWeb/prot-interaction.html

# Functional Proteomics: A Three-Pronged Process



**Data Mining Backfilling**



**Exp. Data Collection**



**Computer Simulation**

---

# Simulation: Three Types of Data (Models)



**Atomic Scale**
**0.1 - 1.0 nm**
**Coordinate data**
**Dynamic data**
**0.1 - 10 ns**
**Molecular dynamics**



**Meso Scale**
**1.0 - 10 nm**
**Interaction data**
**Kon, Koff, Kd**
**10 ns - 10 ms**
**Mesodynamics**



**Continuum Model**
**10 - 100 nm**
**Concentrations**
**Diffusion rates**
**10 ms - 1000 s**
**Fluid dynamics**

# Cell Simulation with DEs

$$\frac{dx_1}{dt} = k_{11}x_1 + k_{21}x_2 + k_{31}x_3 + \ldots$$

$$\frac{dx_2}{dt} = k_{12}x_1 + k_{22}x_2 + k_{32}x_3 + \ldots$$

$$\frac{dx_3}{dt} = k_{13}x_1 + k_{23}x_2 + k_{33}x_3 + \ldots$$
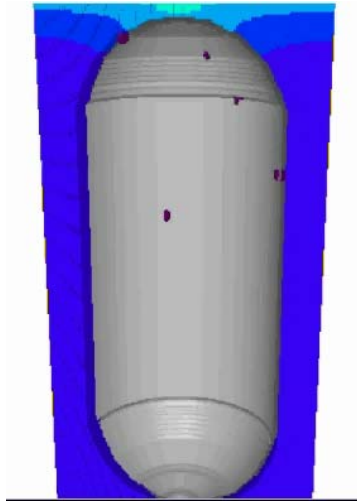
$$\frac{dx_4}{dt} = k_{13}x_1 + k_{24}x_2 + k_{34}x_3 + \ldots$$



---

# Continuum Modelling

- **Desire to simulate spatially and temporally (to make movies)**
- **Use techniques developed for oil and gas resevoir simulation (pumping, diffusion, reaction, pressure -- CMG Inc.)**
- **Uses theory of non-turbulent fluid dynamics, discretized over small volumes**
- **Based on measured parameters of real cells, real metabolites, proteins**
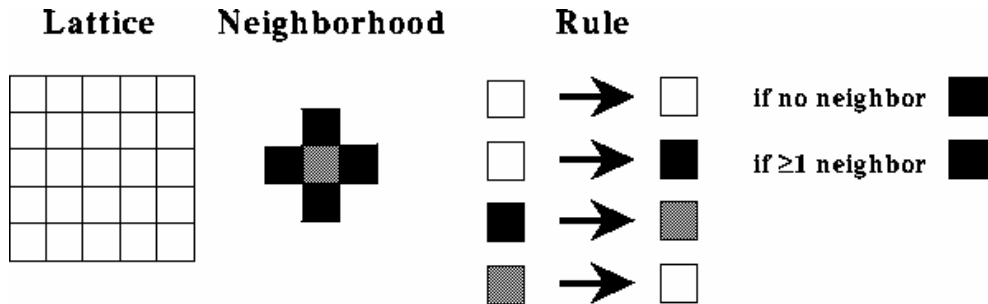
# Continuum Simulation
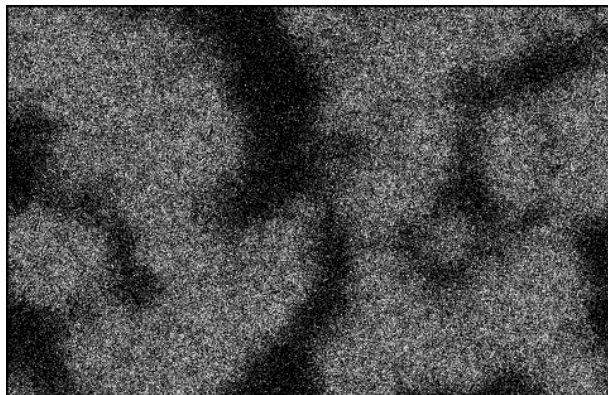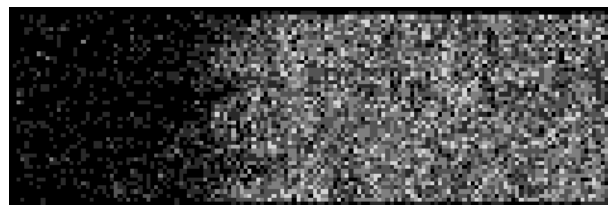


movie

# Cellular Automata (CA)

- **Computer modelling method that uses lattices and discrete state "rules" to model time dependent processes**
- **No differential equations to solve, easy to calculate, more phenomenological**
- **Simple unit behavior -> complex group behavior**
- **Can be used to create Mandelbrot figures**
- **Used to model fluid flow, percolation, reaction + diffusion, traffic flow, ecology**

# Cellular Automata



## Can be extended to 3D lattice

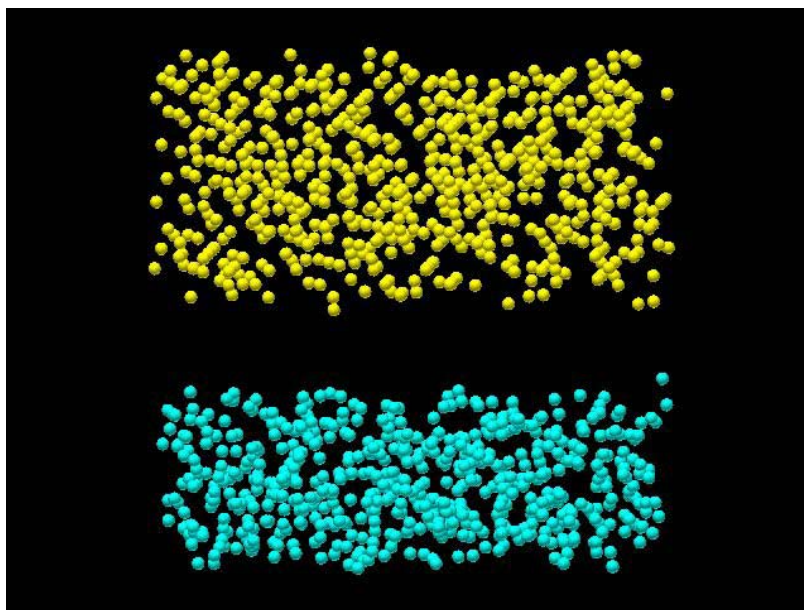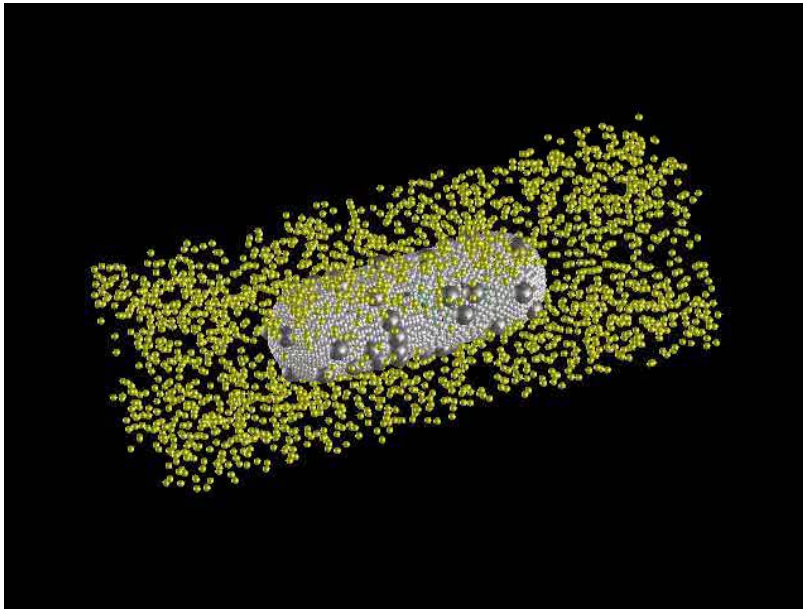# Reaction/Diffusion with Cellular Automata

# Another Example of CA
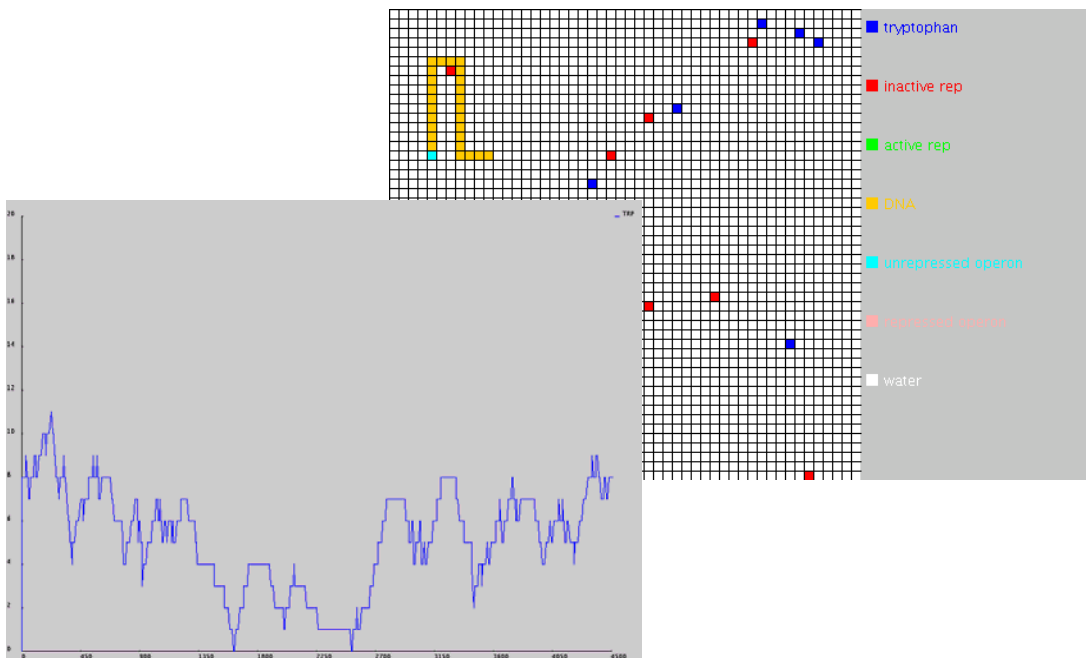


**SimCity 2000**

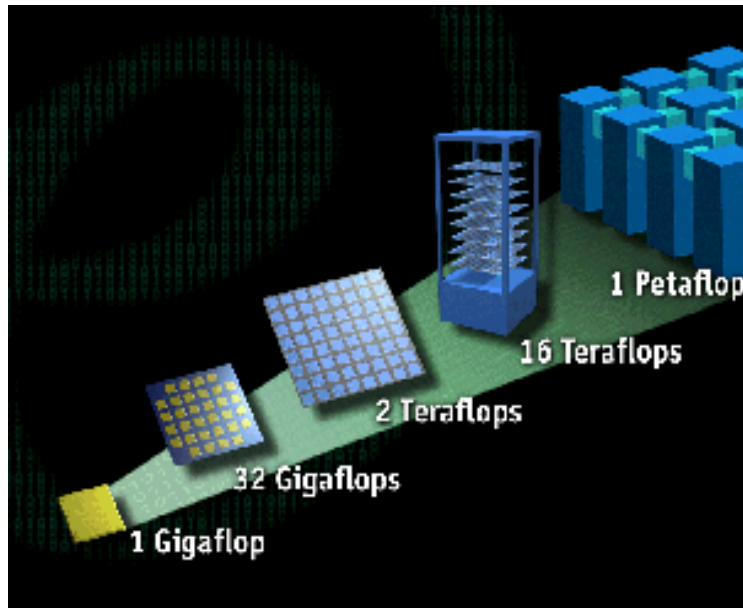# CA Simulations of Diffusion + Reaction

# CA Simulations of Transport



# CA for Trp Repressor

# How Big A Computer?



# Functional Proteomics

- Mixture of experimental and computational techniques
- Trying to reach a point where functions and interactions can be predicted and modelled
- The future of proteomics (and bioinformatics)