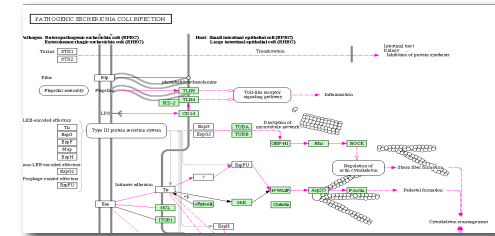# Grand Challenges in Computational Biology
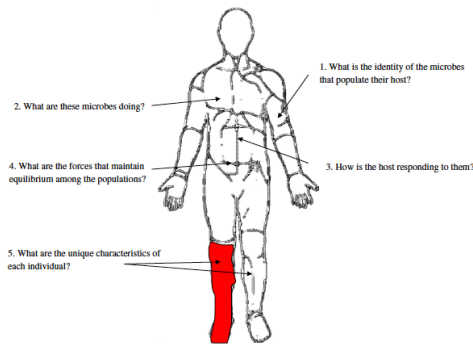


Reconstructing
the Tree of Life

### Kimmen Sjölander
### UC Berkeley

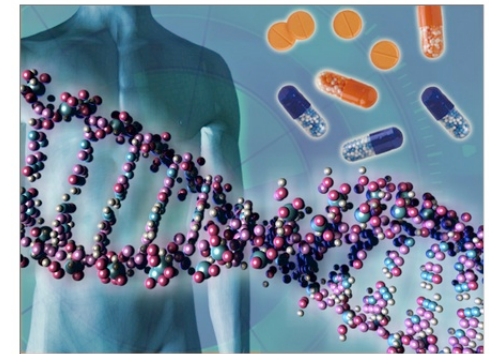### CITRIS-INRIA workshop
### 24 May, 2011

Prediction of biological
pathways and networks



Human microbiome and
metagenome dataset analysis

Infectious disease: new drugs
and diagnostics;
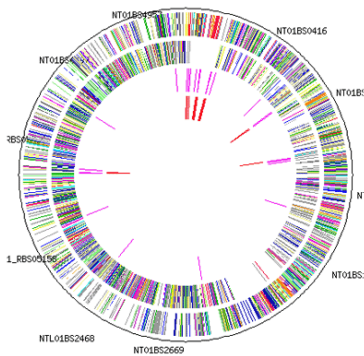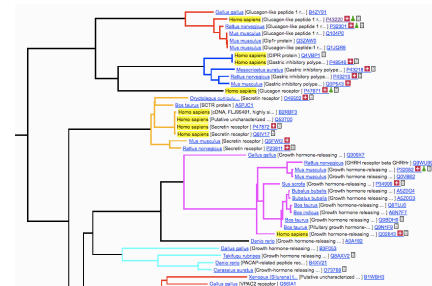pharmacogenomics

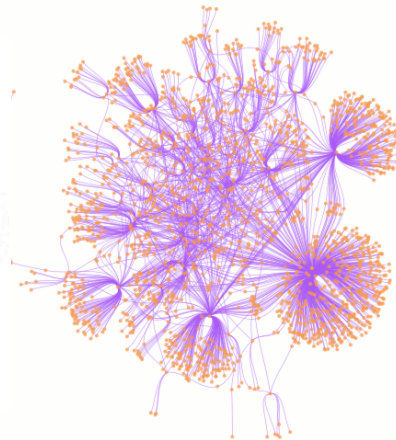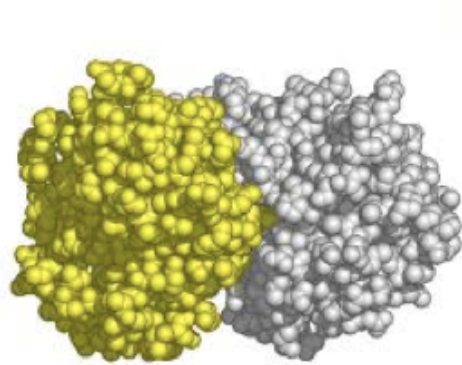Interpreting genetic variation

U.S. DEPARTMENT OF ENERGY

# The expanding genomics universe

- The situation now: huge quantities of noisy, error-ridden and poorly connected data
  - Experimental data are sparse: ~1% of sequences have experimental support for their assigned functions
  - Errors abound: Up to 25% of sequences are mis-annotated [1, 2]
  - The one-time static annotation protocol does not allow annotations to be modified in the light of new evidence [3]
  - Expert knowledge is critical to detecting and correcting annotation errors
    - But manual annotation is expensive and does not scale to the quantity of sequences being produced

1. "Annotation Error in Public Databases: Misannotation of Molecular Function in Enzyme Superfamilies," Schnoes et al, PLoS Computational Biology 2009
2. "Phylogenomic inference of protein molecular function: advances and challenges," Sjolander, Bioinformatics 2004
3. "Genome re-annotation: a wiki solution?" Salzberg, Genome Biology 2007

# Increasing the specificity of function prediction requires the integration of heterogeneous data & bioinformatics methods

Homology & orthology prediction
Genome neighbors
Expression data
Localization information
3D structure
Yeast-2-hybrid data
Phylogenetic profiles
Pull-down assays
Site-directed mutagenesis
Text-mining (co-occurrence in an abstract)
Etc.

Eisenberg et al, "Protein function in the post-genomic era" *Nature* 2000
Sjölander, K., "Phylogenomic inference of protein molecular function: advances and challenges," Bioinformatics 2004
Matthews et al, "Identification of Potential Interaction Networks Using Sequence-Based Searches for Conserved Protein-Protein Interactions or "Interologs"" *Genome Research* 2001
Troyanskaya et al, "A Bayesian framework for combining heterogeneous data sources for gene function prediction (in Saccharomyces cerevisiae)," PNAS, 2003
Myers et al, "Discovery of biological networks from diverse functional genomic data," Genome Biology 2005

# Data is not the same thing as information

# Biologists who need to use bioinformatics tools are divided by a huge gulf from the computer scientists who are creating these tools

# Automatic protein function prediction using a hyper-dimensional network

# Hyperdimensional information network

## for data integration, navigation & community annotation



**Nodes**: Genes/proteins
**Edges**: different types of connection between genes (e.g., orthology, similar structure, interaction, disease association, regulated by, adjacent in metabolic network, genome neighbor, etc.).
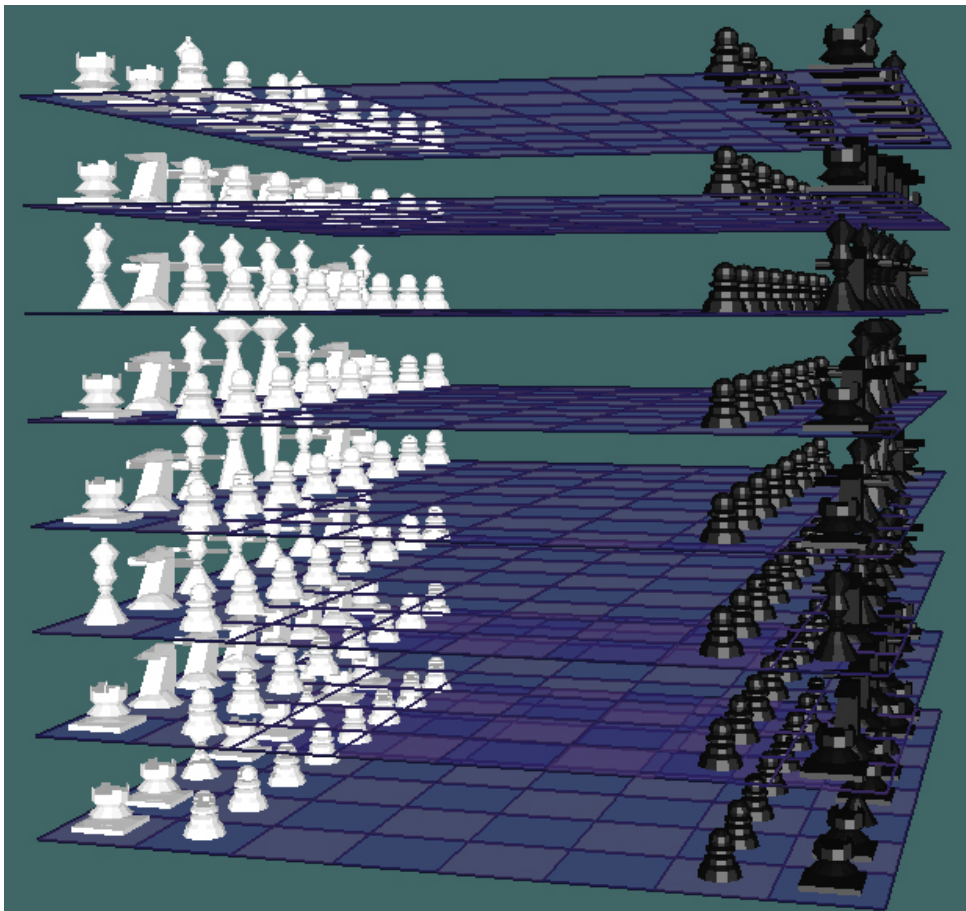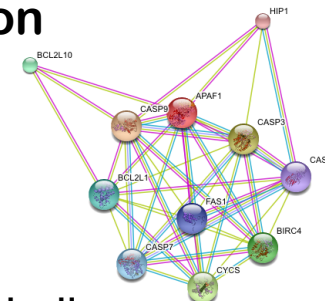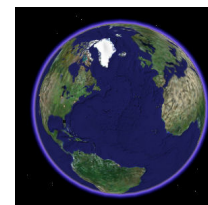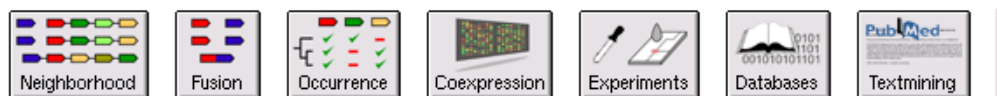**Edges have weights** proportional to confidence

**Experimental data can enter at any point in the graph, and be propagated to neighboring nodes based on learned rules:**

- Biological process for one gene can be made available to genome neighbors
- A protein-protein interaction between two genes in one species can be used to infer corresponding interaction between their orthologs in another
- Roles in a pathway (e.g., EC number) known for one gene can be assigned to an ortholog
- Participation in a biological process can be inferred based on genome neighbors
- 3D structure information can be propagated to all homologs
- Protein structure information can be propagated to all homologs

**Biologists can:** subscribe to news feeds arriving at their selected nodes, upload data, attach links to their papers, manually curate biological "functions"

**Manual annotations** from biologists will need to be weighted according to estimated confidence

# Phylogenomic tools for investigating and interpreting (meta)genome datasets
## (DOE Systems Biology Knowledgebase grant)



**THE METAGENOMICS PROCESS**

2. What are these microbes doing?

4. What are the forces that maintain equilibrium among the populations?

5. What are the unique characteristics of each individual?

**Extract all DNA from microbial community in sampled environment**

**DETERMINE WHAT THE GENES ARE**
**(Sequence-based metagenomics)**
- Identify genes and metabolic pathways
- Compare to other communities
- and more…

**DETERMINE WHAT THE GENES DO**
**(Function-based metagenomics)**
- Screen to identify functions of interest, such as vitamin or antibiotic production
- Find the genes that code for functions of interest
- and more…

**Challenges in metagenome data analysis:**
- Most tools designed for these data answer only "What species are present?" and do not answer the question, "What's going on?" (what processes & pathways are represented)
- Sequences are fragmentary and noisy, presenting additional challenges to bioinformatics methods
- **Huge** datasets (in the millions of reads)

"Harnessing the power of the human microbiome", Blaser, PNAS 2010
"The New Science of Metagenomics: Revealing the Secrets of Our Microbial Planet" Committee on Metagenomics: Challenges and Functional Applications, National Research Council. 2007.
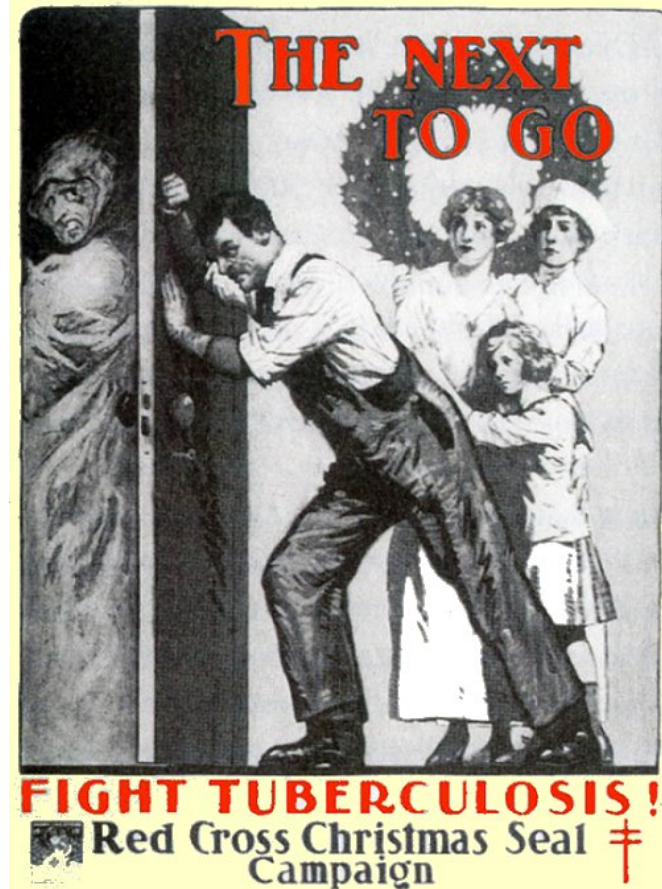
# SNP prioritization and interpreting human genetic variation



**Prediction of biological pathways and network alignment**

SNPs occurring in coding regions of the genome can be prioritized for investigation based on:

- Predicted biological process or function of gene containing SNP
- Predicted interactions (hubs of networks) of gene containing SNP
- Impact of mutation at that site (INTREPID and Discern methods)

9

# PhyloFacts Pathogen Commons





- **Drug target identification & prioritization**

- **Development of accurate diagnostics**

**TB collaborations**
- UC Berkeley Center for Emerging and Neglected Diseases (Tom Alber, Lee Riley, others)
- Royal Institute of Tropical Diseases, Amsterdam, Netherlands (Richard Anthony)
- Institute of Bioinformatics, Bangalore, India (Akhilesh Pande)
- IISc, Bangalore, India (Nagasuma Chandra)

**Was there really life before the web?**

How can we bring this to biology?