What is the Difference between a Human and a Chimp?

T. M. Murali

murali@cs.vt.edu http://bioinformatics.cs.vt.edu/~murali

Introduction to Computational Biology and Bioinformatics (CS 3824) March 25 and April 1, 2010

The Story Began in 1953



T. M. Murali

CS 3824: March 25 and April 1, 2010

Human vs. Chimp

The Human Genome

- DNA is a (very long) string containing letters A, T, C, and G.
- Length of human genome is 3 billion base pairs



The Human Genome

- DNA is a (very long) string containing letters A, T, C, and G.
- Length of human genome is 3 billion base pairs (1.8 metres).
- The Human Genome Project determined the spelling of the genome.



The Human Genome

- DNA is a (very long) string containing letters A, T, C, and G.
- Length of human genome is 3 billion base pairs (1.8 metres).
- The Human Genome Project determined the spelling of the genome.
- Eric Lander (Nano-Lecture, 2003 lg Nobel Prize Ceremony):

Genome. Bought the book, hard to read.



The Human Genome Project





Other Projects

The Human Genome Project

Before: human genome has about 100,000 genes.





The Human Genome Project

Before: human genome has about 100,000 genes.





After: human genome has about 30,000 genes.

CS 3824: March 25 and April 1, 2010

Shock and Dismay

- The New York Times: Genome Analysis Shows Humans Survive on Low Number of Genes The two teams report that there are far fewer human genes than thought—probably a mere 30,000 or so—only a third more than those found in the roundworm. ... The impact on human pride is another matter.
- Washington Post: It also raises new and difficult questions, such as how human beings—with all their passions and fears, their capacity for art, music, culture and war—can be all that they are with just 30,000 or so genes, only five times as many as in baker's yeast.

Genome size comparison

	Species C	hromosome	s Genes	Base pairs
X	Human (Homo sapiens)	46 (23 pairs)	28-35,000	3.1 billion
	Mouse (Mus musculus)	40	22.5-30,000	2.7 billion
6	Puffer fish (Fugu rubripes)	44	31,000	365 million
~	Malaria mosquito (Anopheles gambiae)	6	14,000	289 million
PR	Fruit fly (Drosophila melanogaster)	8	14,000	137 million
と	Roundworm (C. elegans)	12	19,000	97 million
•	Bacterium * (E. coli)	1	5,000	4.1 million

*Bacterial chromosomes are chromonemes, not true chromosomes

JOHN BLANCHARD / The Chronicle











Chimp and chump genomes are only about 1.2% different!

Molecular Biology

 Genomes provide the parts lists (e.g., genes and proteins) but do not directly tell us how these parts fit.

- Genomes provide the parts lists (e.g., genes and proteins) but do not directly tell us how these parts fit.
- We need to understand how genes, proteins, and other molecules interact with other in different cell states, different tissues, and under different external conditions.
- Study only of individual elements is unlikely to reveal higher-order organisation of cellular interaction networks.



Sea Urchin (Strongylocentrotus purpuratus)



Sea Urchin (Strongylocentrotus purpuratus)



- Very important in developmental biology.
- Many principles of embryo development were discovered in the sea urchin.



A Cell is a Modular



A Cell is a Modular



A Cell is a Modular Network



T. M. Murali

CS 3824: March 25 and April 1, 2010

Human vs. Chimp

A Cell is a Modular Network



C Module A functions:

Vegetal plate expression in early development:

Synergism with modules B and G enhancing endoderm expression in later development:

Repression in ectoderm (modules E and F) and skeletogenic mesenchyme (module DC):

Modules E, F and DC with LiCI treatment:

A Cell is a Modular Network that Computes



-			
if (F = 1 or E = 1 or CD = 1) and (Z = 1)		Repression functions of modules F, E, and	
	α = 1	DC mediated by Z site	
else	$\alpha = 0$		
if (P = 1 and CG, = 1)		Both P and CG, needed for synergistic lin	
	$\beta = 2$	with module B	
else	$\beta = 0$		
if (CG ₂ = 1 and CG ₃ = 1 and CG ₄ = 1)		Final step up of system output	
	γ = 2		
else	γ = 1		
$\delta(t) = B(t) + G(t)$		Positive input from modules B and G	
$\varepsilon(t) = \beta^* \delta(t)$		Synergistic amplification of module B output by CG,-P subsystem	
if $(\varepsilon(t) = 0)$		Switch determining whether Otx site in	
	ξ(t) = Otx(t)	module A, or upstream modules (i.e., mainly module B), will control level of activity	
else	$\xi(t) = \varepsilon(t)$		
if (α = 1)		Repression function inoperative in	
	η(t) = 0	endoderm but blocks activity elsewhere	
else	$\eta(t) = \xi(t)$		
$\Theta(t) = \gamma^* \eta(t)$		Final output communicated to BTA	

Network is Complex



Network is Complex



Network is Complex but Very Poorly Understood



Challenges with Molecular Interaction Networks

- Biological data sets and networks are large.
- They are intricate and of very diverse types.
- They are noisy: experiments are error-prone.
- They are highly incomplete. We barely know which genes interact, let alone the detailed kinetics of each interaction.

My Research: Understanding Interaction Networks



- Automatically find modules of coherently acting molecules from data.
- Build predictive module-level models of the cell.
- Emphasise a phenomenological approach to systems biology.
- Develop techniques in graph and discrete algorithms, data mining, and machine learning and apply them to solve specific biological questions.

Research Applications

- Predict the functions that genes and proteins perform in the cell.
- Develop drugs that may be effective against multiple pathogens.
- Build models of how cells communicate with each other.
- Zero in on the molecules and interactions that are active in cancer.
- Develop biologically-relevant representations of molecular interaction networks.

Functions of Many Genes are Unknown

- ▶ We have the sequences of 100s of genomes.
- ▶ We know the locations of genes in these genomes.
- ▶ Functions of over 50% of the genes are unknown!

Functions of Many Genes are Unknown

- ▶ We have the sequences of 100s of genomes.
- We know the locations of genes in these genomes.
- ► Functions of over 50% of the genes are unknown!
- Genes with similar sequences in different organisms are likely to have the same function.
- We need techniques for function prediction that go beyond sequence similarity.

Exploit Social Network of Genes



 Exploit network structure to determine whether gray genes have the same function as the red genes or the blue genes.

(Karaoz, Murali, Letovsky, Zheng, Ding, Cantor and Kasif, *PNAS*, 2004; Murali, Wu, and Kasif, *Nature Biotech.*, 2006.)

Node States



- Interaction network is a graph G = (V, E).
- Each node i has an associated state s_i:
 - $s_i = 1$: gene *i* is annotated with current function.
 - $s_i = -1$: gene *i* is annotated with another function.
 - $s_i = 0$: otherwise.
- An edge between nodes *i* and *j* has a weight w_{ij} .

Maximally-Consistent Assignments



- An edge is *consistent* if it is incident on nodes with the same state.
- Maximally-consistent assignment: number of consistent edges is maximised.

Maximally-Consistent Assignments



- An edge is *consistent* if it is incident on nodes with the same state.
- Maximally-consistent assignment: number of consistent edges is maximised.

Computational goal: Assign state of -1 or +1 to nodes with initial state 0 to achieve maximal consistency by maximising

$$\sum_{(u,v) \text{ is an edge}} w_{uv} s_u s_v$$

Activation rule is

$$s_i = \operatorname{sgn}\left(\sum_{j\in N_i} w_{ij}s_j\right),$$

Activation rule is

$$s_i = \operatorname{sgn}\left(\sum_{j\in N_i} w_{ij}s_j\right),$$

where N_i = neighbours of node *i*.

Applying this rule:

Activation rule is

$$s_i = \operatorname{sgn}\left(\sum_{j\in N_i} w_{ij}s_j\right),$$

- Applying this rule:
 - Parallel update: each node updates itself in parallel with the other nodes.
 - Serial update: go through each node in sequence.

Activation rule is

$$s_i = \operatorname{sgn}\left(\sum_{j\in N_i} w_{ij}s_j\right),$$

- Applying this rule:
 - Parallel update: each node updates itself in parallel with the other nodes.
 - Serial update: go through each node in sequence.
- Stopping criterion: converge when no node's state changes.

Activation rule is

$$s_i = \operatorname{sgn}\left(\sum_{j\in N_i} w_{ij}s_j\right),$$

- Applying this rule:
 - Parallel update: each node updates itself in parallel with the other nodes.
 - Serial update: go through each node in sequence.
- Stopping criterion: converge when no node's state changes.
- ▶ Proof of convergence: show (∑_(u,v) is an edge w_{uv}s_us_v) always increases when you change a node's state.















T. M. Murali

Human Proteins Interacting with Pathogens



Dyer, Murali, and Sobral, PLoS Pathog, 2008.

T. M. Murali

Pathogens are Becoming Drug-Resistant

Pssst! Hey kid! Wanna be a **Superbug**...? Stick some of **this** into your genome... Even **penicillin** won't be able to harm you...! 00000

It was on a short-cut through the hospital kitchens that Albert was first approached by a member of the Antibiotic Resistance.

$\textbf{One-Bug-One-Drug} \Rightarrow \textbf{Many-Bugs-One-Drug}$

- Develop drugs that target human proteins.
- Prioritize human proteins interacting with *multiple* pathogens.



T. M. Murali

CS 3824: March 25 and April 1, 2010

Viral Dependency Factors



- RNA viruses like HIV have very few genes.
- Viral dependency factor (VDF): human protein that virus needs to replicate and propagate.

Viral Dependency Factors



- RNA viruses like HIV have very few genes.
- Viral dependency factor (VDF): human protein that virus needs to replicate and propagate.
- Recent genome-wide experiments have discovered dozens of VDFs for HIV, flu virus, West Nile virus, and Hepatitis C virus.
- ► Different experiments for HIV show very little overlap.

Viral Dependency Factors



- RNA viruses like HIV have very few genes.
- Viral dependency factor (VDF): human protein that virus needs to replicate and propagate.
- Recent genome-wide experiments have discovered dozens of VDFs for HIV, flu virus, West Nile virus, and Hepatitis C virus.
- Different experiments for HIV show very little overlap.
- Proteins in one experiment interact with proteins detected in other experiments.

Collaboration with Michael Katze (Dept of Microbiology, Univ. of Washington) and Brett Tyler (VBI)

- Treat the problem as one of predicting gene function: which human genes have the function of "being used by viruses to propagate?"
- Modify previous approach: for every node v compute a function 0 ≤ f(v) ≤ 1

Collaboration with Michael Katze (Dept of Microbiology, Univ. of Washington) and Brett Tyler (VBI)

- Treat the problem as one of predicting gene function: which human genes have the function of "being used by viruses to propagate?"
- Modify previous approach: for every node v compute a function 0 ≤ f(v) ≤ 1 to minimise

$$\sum_{(u,v)} w_{uv}(f(u) - f(v))^2$$

Collaboration with Michael Katze (Dept of Microbiology, Univ. of Washington) and Brett Tyler (VBI)

- Treat the problem as one of predicting gene function: which human genes have the function of "being used by viruses to propagate?"
- Modify previous approach: for every node v compute a function 0 ≤ f(v) ≤ 1 to minimise

$$\sum_{(u,v)} w_{uv}(f(u) - f(v))^2 + \lambda \sum_{v} f^2(v)$$

Collaboration with Michael Katze (Dept of Microbiology, Univ. of Washington) and Brett Tyler (VBI)

- Treat the problem as one of predicting gene function: which human genes have the function of "being used by viruses to propagate?"
- Modify previous approach: for every node v compute a function 0 ≤ f(v) ≤ 1 to minimise

$$\sum_{(u,v)} w_{uv}(f(u) - f(v))^2 + \lambda \sum_{v} f^2(v)$$

Solve linear system of equations:

$$f(v) = \frac{\sum_{u} w_{uv} f(u)}{\lambda + \sum_{u} w_{uv}}$$

Cross-Validation Results



CS 3824: March 25 and April 1, 2010

Independent Evaluation of Predictions

- ► VDFs for HIV were discovered in cell lines.
- Protein-protein interactions discovered in a wide variety of experiments.

Independent Evaluation of Predictions

- ► VDFs for HIV were discovered in cell lines.
- Protein-protein interactions discovered in a wide variety of experiments.
- ► Used gene expression data from SIV-infected non-human primates.
 - African Green Monkeys (AGMs) are natural hosts for SIV.
 - Pig-tailed Macaques (PMs) are susceptible to SIV.
- Computed which genes were differentially expressed between AGMs and PMs.
- Computed the statistical significance of the overlap between differentially-expressed genes and known/predicted HDFs.





Other Projects

- Building blocks of molecular interaction networks.
- Models to enable explicit comparisons between cell's response to different conditions.

What Makes Systems Biology Exciting?

- Use principles in computer science to solve problems that impact human life.
- Research is inter-disciplinary: collaborate closely with biochemists, bioengineers, geneticists, doctors, and plant biologists.
- Seek synergistic collaborations between computer science and biology.
- Train students in both computer science and biology.

Michael Katze	RNA viruses (HIV, flu, etc.)
Padma Rajagopalan	Liver tissue engineering
Bruno Sobral	Host-pathogen interactions
Brett Tyler	Predicting gene function, plant pathogens

Other Projects

How Can You Contribute?



Other Projects

How Can You Contribute?



- Curiosity about how the cell works, how it becomes ill, how does it survive attacks.
- Curious to learn how computer science helps to discover new things about the cell.

How Can You Contribute?



- Curiosity about how the cell works, how it becomes ill, how does it survive attacks.
- Curious to learn how computer science helps to discover new things about the cell.
- ▶ We need students who will learn both computer science and biology.
- They will form the next generation of scientists in systems biology.

Nothing in Biology Makes Sense Except in the Light of Evolution

