

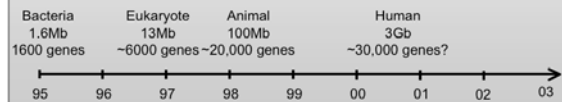
Graphical Models in Computational Molecular Biology

Nir Friedman
Hebrew University

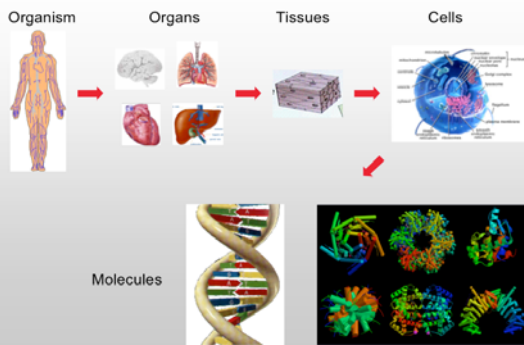
Includes slides by:

Yoseph Barash, Nebojsa Jojic, Tommy Kaplan,
Daphne Koller, Iftach Nachman, Dana Pe'er, Tal
Pupko, Aviv Regev, Eran Segal

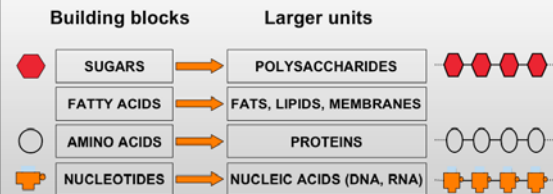
The Age of Genomes



Biology → Molecular Biology



Macro Molecules as Sequences

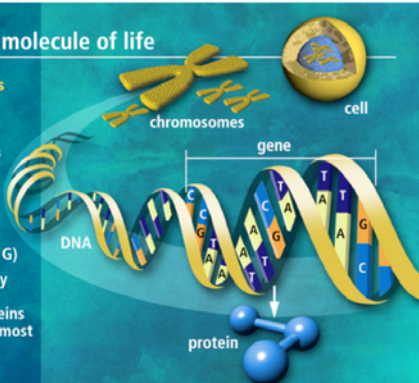


DNA the molecule of life

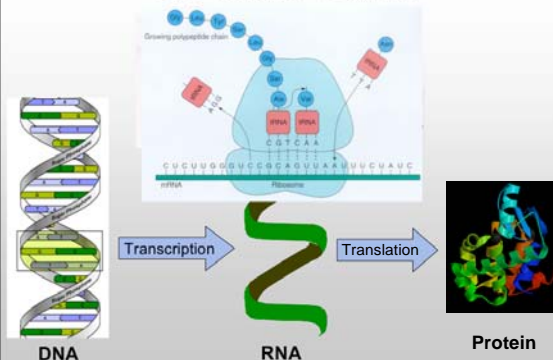
Trillions of cells

Each cell:

- 46 human chromosomes
- 2 meters of DNA
- 3 billion DNA subunits (the bases: A, T, C, G)
- Approximately 30,000 genes code for proteins that perform most life functions



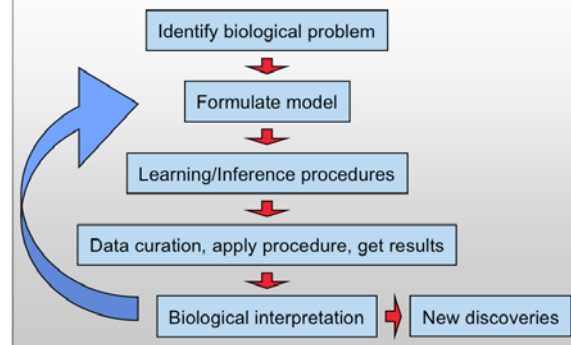
The Central Dogma



What Biologists Want?

- ◆ Identify Components
 - Find genes in DNA sequences
 - ◆ Associate them with function
 - Understand what a gene does
 - ◆ Interactions between components
 - Which genes work together?
 - ◆ Dynamics of systems
 - When genes are activated, and by what?
 - ◆ How did we get here
 - How did genes evolve, and why this way?
- (And many other things)

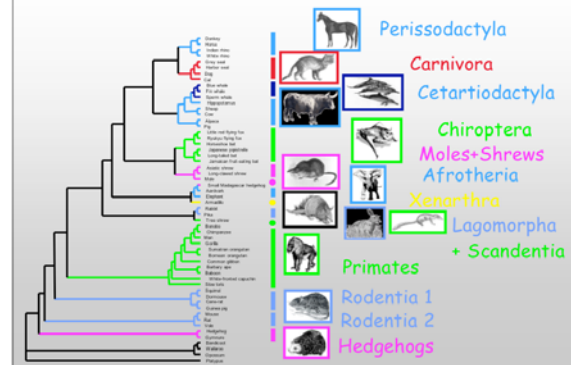
Philosophy



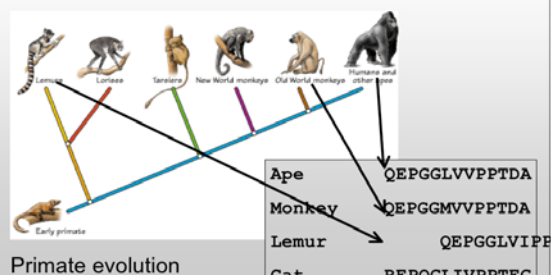
Outline

- ◆ Sequence evolution
- ◆ Protein families
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Discussion

Evolutionary History



Phylogeny



Related Proteins

Conserved Position	Insertion	Deletion
QEPGGLVVPPTDA		
QEPGGMVVPPTDA		
QEPGGLVIPE		
REPQGLIVPPTG		

Globins <http://www.sanger.ac.uk/cgi-bin/Pfam/getacc?PF00042>

Evolution

Evolution = Neutral Variation + Selection

Neutral variation

- Random changes
Mutation, duplication, deletions, rearrangements, ...

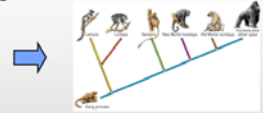
Selection

- Preference for variants with better fit

Challenges

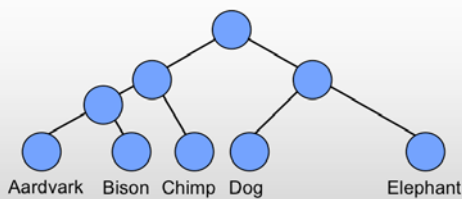
- ♦ From sequences to phylogenies

Ape	QEPGGLVVPPTDA
Monkey	QEPGGMVVPPTDA
Lemur	QEPGGLVIPPTDA
Cat	REPQGLIVPPTeg



- ♦ Understand selection
Selection → Function & Structure
- ♦ Reconstruct history of evolutionary events

Probabilistic Model of Evolution



Random variables – sequence at current day taxa or at ancestors

Potentials/Conditional distribution – represent the probability of evolutionary changes along each branch

Parameterization of Phylogenies

Assumptions:

- ♦ Positions (columns) are independent of each other
- ♦ Each branch is a reversible continuous time discrete state Markov process

$$P(a \rightarrow c | t + t') = \sum_b P(a \rightarrow b | t) P(b \rightarrow c | t')$$

$$P(a)P(a \rightarrow b | t) = P(b)P(b \rightarrow a | t)$$

governed by a **rate matrix** Q

Parameterization:

- ♦ Rate matrix + tree topology + branch lengths

Computational Tasks

Likelihood computation, inference of ancestral states

- ♦ Inference (dynamic programming, belief propagation)

Branch length estimation:

- ♦ Parameter estimation (EM)

Reconstruction:

- ♦ Structure learning

Felsenstien, JME 1981

Maximum Likelihood Reconstruction

Observed data: (D)

- ♦ N sequences of length M
- ♦ Each position: an independent sample from the marginal distribution over N current day taxa

Likelihood:

- ♦ Given a tree (T, t) :

$$l(T, t : D) = \log P(D | T, t)$$

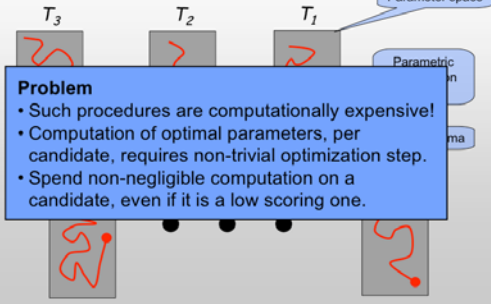
$$= \sum_m \log P(x_{[1, \dots, N]}^m | T, t)$$

Goal:

- ♦ Find a tree (T, t) that maximizes $l(T, t : D)$.

Current Approaches

- ◆ Perform search over possible topologies



SEMPHY: Structural EM Phylogenetic Reconstruction

Borrow the idea of Structural EM:
Use parameters found for current topology to help evaluate new topologies.

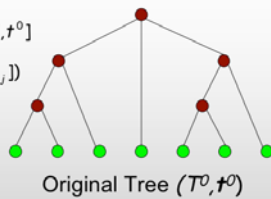
Outline:

- ◆ Perform search in (T, t) space, using EM-like iterations:
 - **E-step:** use current solution to compute expected sufficient statistics for **all topologies**
 - **M-step:** select new topology based on these expected sufficient statistics

Friedman et al, JCB 2002

Algorithm Outline

- Compute: $E[S_{(i,j)}(a,b) | D, T^0, t^0]$
- Weights: $w_{i,j} = \max_{\tau} F(\tau, E[S_{(i,j)}])$



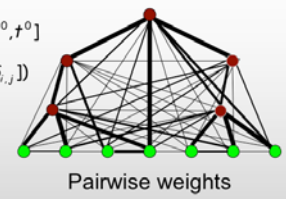
Unlike standard EM for trees, we compute all possible pairwise statistics

Time: $O(N^2M)$

Friedman et al, JCB 2002

Algorithm Outline

- Compute: $E[S_{(i,j)}(a,b) | D, T^0, t^0]$
- Weights: $w_{i,j} = \max_{\tau} F(\tau, E[S_{(i,j)}])$
- Find: $T' = \arg \max_{\tau} \sum_{(i,j) \in \mathcal{T}} w_{i,j}$



This stage also computes the branch length for each pair (i,j)

Friedman et al, JCB 2002

Algorithm Outline

- Compute: $E[S_{(i,j)}(a,b) | D, T^0, t^0]$
- Weights: $w_{i,j} = \max_{\tau} F(\tau, E[S_{(i,j)}])$
- Find: $T' = \arg \max_{\tau} \sum_{(i,j) \in \mathcal{T}} w_{i,j}$
- Construct bifurcation T_1



Fast greedy procedure to find tree
By construction:

$$Q(T', t') \geq Q(T_0, t_0)$$

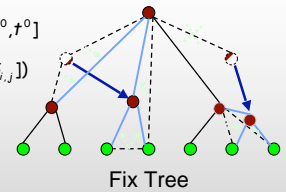
Thus,

$$l(T', t') \geq l(T_0, t_0)$$

Friedman et al, JCB 2002

Algorithm Outline

- Compute: $E[S_{(i,j)}(a,b) | D, T^0, t^0]$
- Weights: $w_{i,j} = \max_{\tau} F(\tau, E[S_{(i,j)}])$
- Find: $T' = \arg \max_{\tau} \sum_{(i,j) \in \mathcal{T}} w_{i,j}$
- Construct bifurcation T_1



Remove redundant nodes
Add nodes to break large degree

This operation preserves likelihood

$$l(T_1, t') = l(T', t') \geq l(T_0, t_0)$$

Friedman et al, JCB 2002

Algorithm Outline

→ Compute: $E[S_{(i,j)}(a,b) | D, T^0, t^0]$

→ Weights: $w_{i,j} = \max F(t, E[S_{(i,j)}])$

→ Find: $T' = \operatorname{argmax}_T \sum_{(i,j) \in T} w_{i,j}$

→ Construct bifurcation T_1

→ Thm: $I(T_1, t_1) \geq I(T_0, t_0)$

These steps are then repeated until convergence



Friedman et al, JCB 2002

Beyond Vanilla Flavor

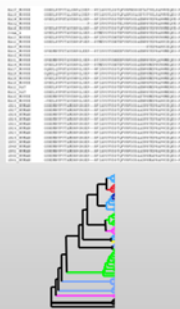
Rate variation:

- ◆ Different positions evolve at different rate

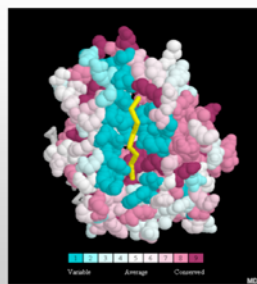


Yang, JME 1994

Application: Map Conservation



Calculate rates

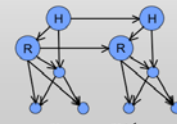


<http://consurf.tau.ac.il/>

Beyond Vanilla Flavor

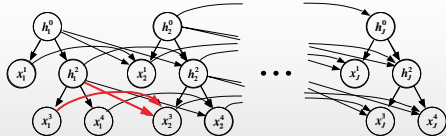
Dependencies between Positions:

- ◆ Dependent rate
- ◆ Dependent state



Felsenstein and Churchill, MBE 1996

Phylogenetic HMMs



- ◆ Each substitution depends on the substitution at the previous position
- ◆ This structure captures **context specific** effects during evolution

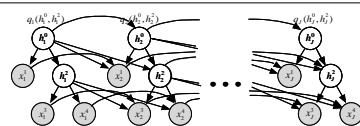
Problem:

- ◆ Resulting model is intractable

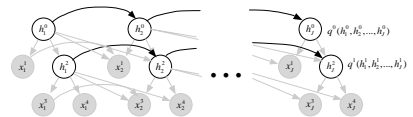
Siepel & Haussler, MBE 2004

Variational approximations

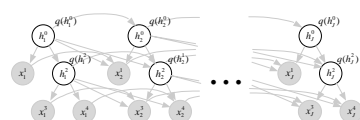
Product of trees



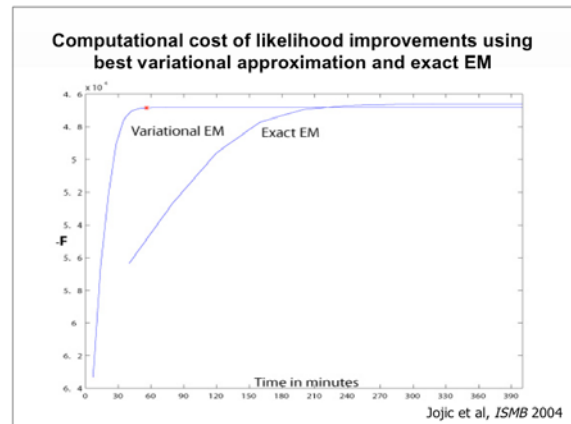
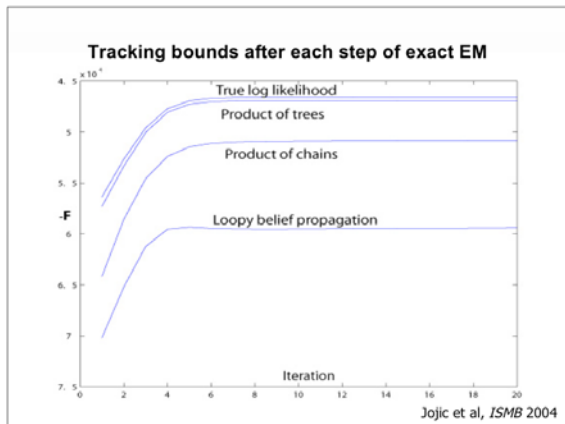
Product of chains



Mean field (product of "nodes")



Jojic et al, ISMB 2004



Remaining Challenges

- ◆ To learn phylogenies, need to align sequences
 - ◆ For good alignment, need to know which positions are conserved
- Can we do both at the same time?

Requires

- ◆ Dealing with insertions + deletions
 - ◆ Different ways of shifting each sequence
- Hard problem

Outline

- ◆ Sequence evolution
- ◆ **Protein families**
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Discussion

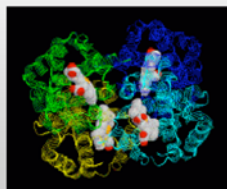
Proteins

Sequence

- Structure
- Function

Major question:

- ◆ How to annotate structure and function of new protein sequences?



Hemoglobin β

Protein Families

Idea:

- ◆ Use knowledge about function of known proteins
 - ◆ Sequence conservation
- similar structure & function

Protein Family:

A collection of proteins sequences that have a common

- ◆ Evolutionary ancestor
- ◆ Structure
- ◆ Function

•

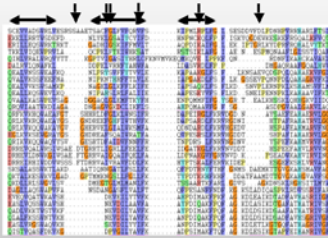
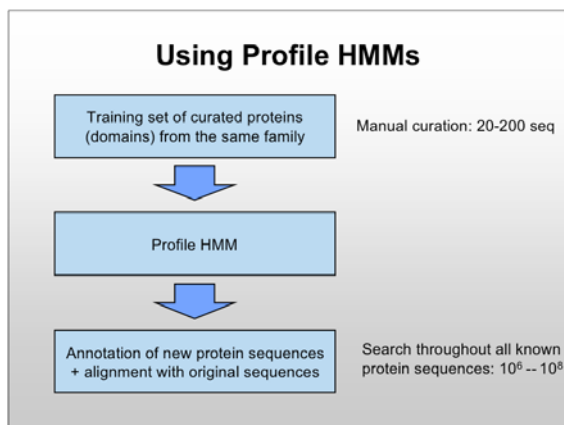
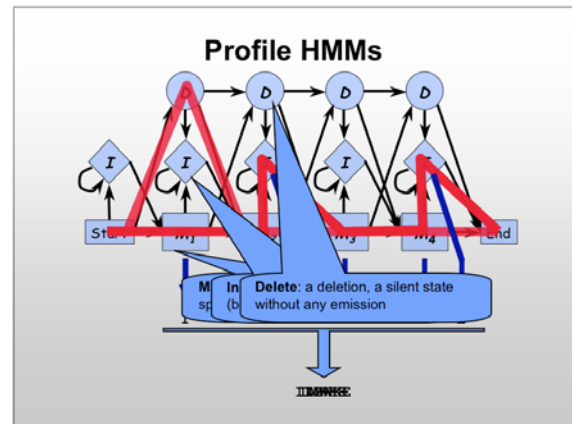
•

Protein Family


Multiple Sequence alignment → detect new proteins

Need to represent

- Commonalities
- Variations
- Conserved stretches
- Potential insertions
- Potential deletions
- Preferred letters at each position

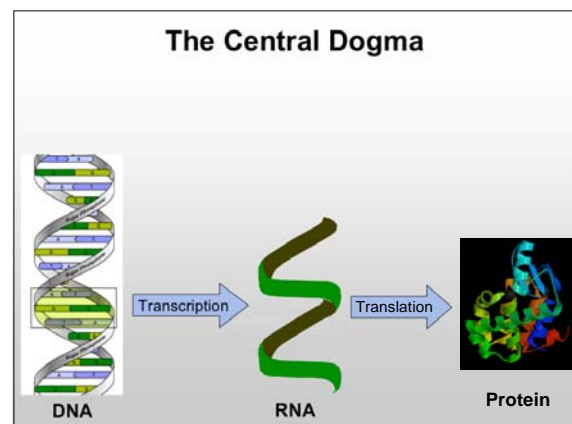



PFAM Database

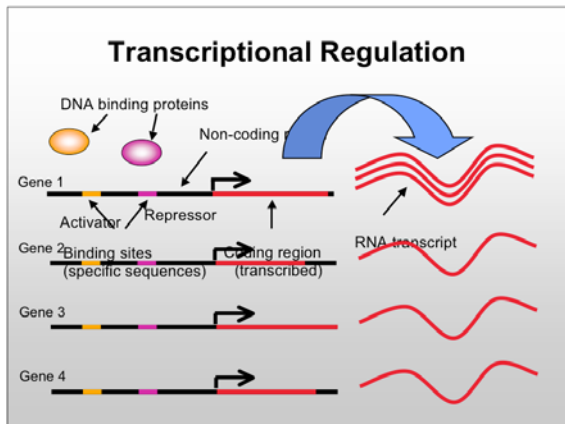


Outline

- Sequence evolution
- Protein families
- Transcriptional regulation
 - Transcription factor binding sites
 - Combinatorial regulation
- Gene expression
- Discussion



•



Transcription Factor Binding Sites

- Gene regulatory proteins contain structural elements that can "read" DNA sequence "motifs"
- The amino acid – DNA recognition is not straightforward
- Experiments can pinpoint binding sites on DNA

Zinc finger, Helix-Turn-Helix, Leucine zipper

Modeling Binding Sites

Given a set of (aligned) binding sites ...

- Consensus sequence
- Probabilistic model (profile of a binding site)

	A	C	G	T		A	C	G	T		A	C	G	T		A	C	G	T
A	4	3	1	1	0	0	1	0	0	1	0	0	1	0	1	0	0	1	0
C	1	3	0	0	0	0	0	13	6	0	0	0	1	9	0	0	0	0	1
G	5	5	13	13	14	14	0	8	14	12	13	1	0	0	0	0	0	0	0
T	4	3	0	0	0	0	0	0	0	0	0	1	0	3	0	0	0	0	0

Is this sufficient?

How to model binding sites ?

$P(X_1 X_2 X_3 X_4 X_5) = ?$ represents a distribution of binding sites

- Profile:** Independency model
- Tree:** Direct dependencies
- Mixture of Profiles:** Global dependencies
- Mixture of Trees:** Both types of dependencies

Barash et al, RECOMB 2002

Arabidopsis ABA binding factor 1

Profile

Test LL per instance -19.93

Tree

Test LL per instance -18.47 (+1.46)
(improvement in likelihood > 2.5-fold)

Mixture of Profiles

76%

24%

Test LL per instance -18.70 (+1.23)
(improvement in likelihood > 2-fold)

Barash et al, RECOMB 2002

Likelihood improvement over profiles

TRANSFAC: 95 aligned data sets

Fold-change in likelihood

Zinc finger, bZIP, bHLH, Helix Turn Helix, beta Sheet, others, ???

Barash et al, RECOMB 2002

Motif finding problem

Input: A set of potentially co-regulated genes
Output: A common motif in their promoters

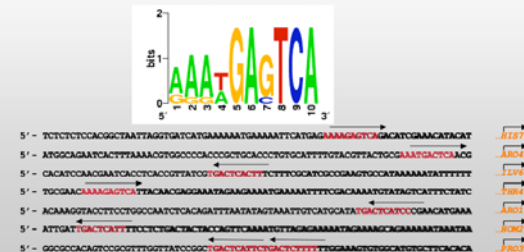
Sources of data:

- ◆ Gene annotation (e.g. Hughes et al, 2000)
- ◆ Gene expression (e.g. Spellman et al, 1998; Tavazoie et al, 2000)
- ◆ ChIP (e.g. Simon et al, 2001; Lee et al, 2002)

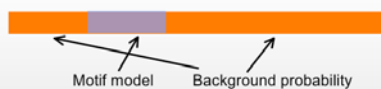


Example

- ◆ Upstream regions from yeast *Sacharomyces cerevisiae* genes (300-600bp)



Probabilistic Model



- ◆ Background probability: given
- ◆ Motif model – parameters being learned

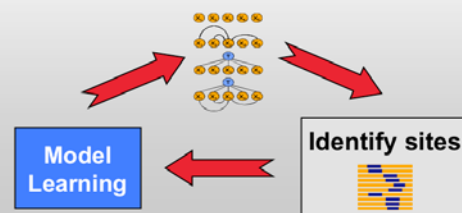
Hidden variable:

- ◆ Location of motif within each sequence

Learning models: unaligned data

EM (MEME)

- ◆ Identify binding site positions
- ◆ Learn a dependency model



Learning models: unaligned data

EM (MEME)

Gibbs Sampling (AlignACE)

Discriminative Sampling

- Find motif that best separates positive examples from rest of promoter sequences

Challenges

Small training sets

- ◆ 10-500 sequences (out of 1000s genes)

Short motifs within long sequences

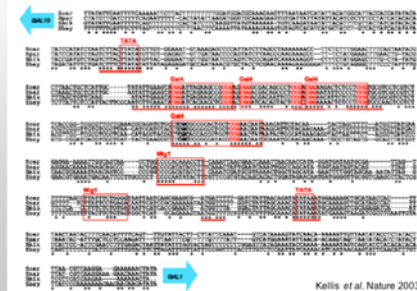
- ◆ Motifs are 6-20bp, promoters are 500-5000bps

Motifs are not perfect words

- ◆ Mismatches are allowed

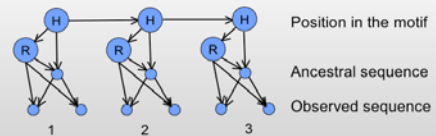
Comparative Genomics

Functional areas should be conserved



Modeling Conserved Motifs

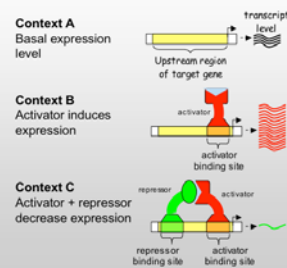
- ♦ Parallel search in orthologous promoters
- ♦ Model evolution at each position



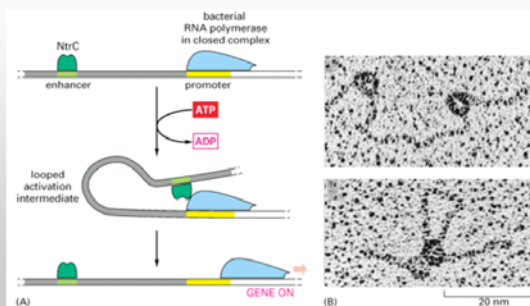
Outline

- ♦ Sequence evolution
- ♦ Protein families
- ♦ Transcriptional regulation
 - Transcription factor binding sites
 - **Combinatorial regulation**
- ♦ Gene expression
- ♦ Discussion

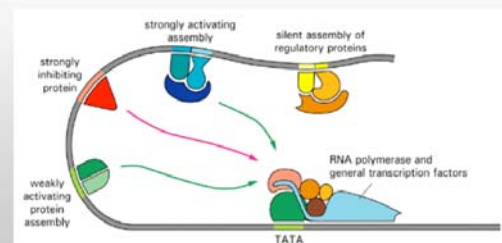
Combinatorial Regulation

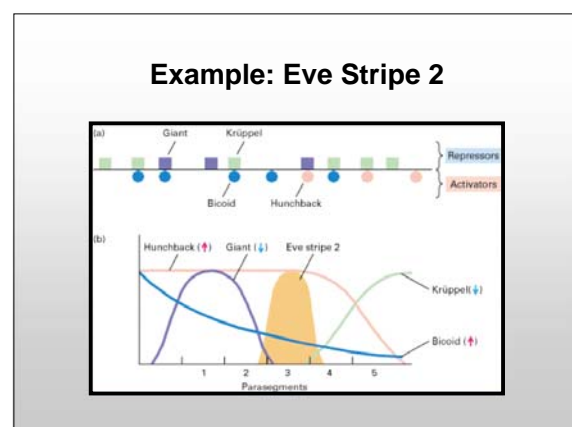
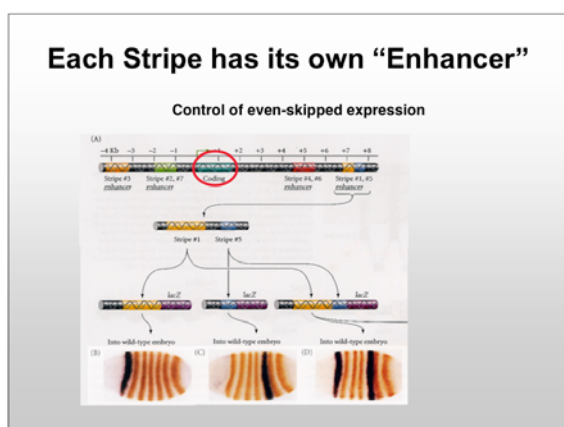
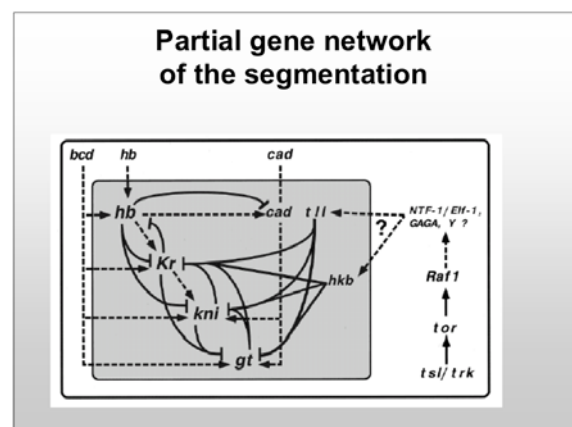
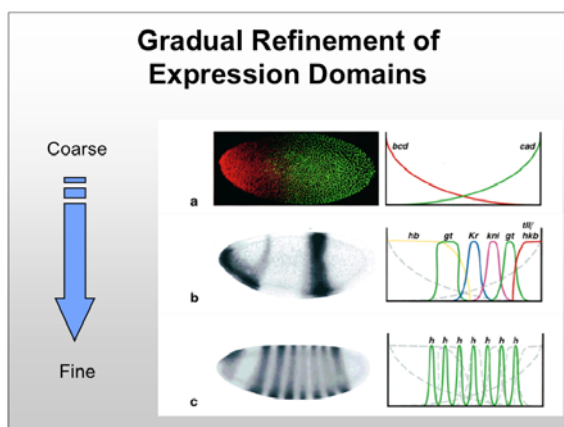
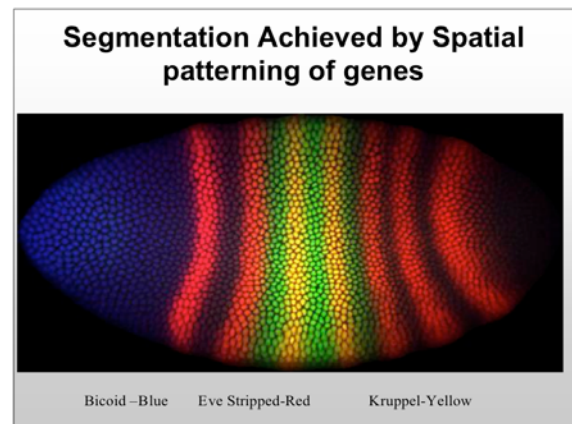
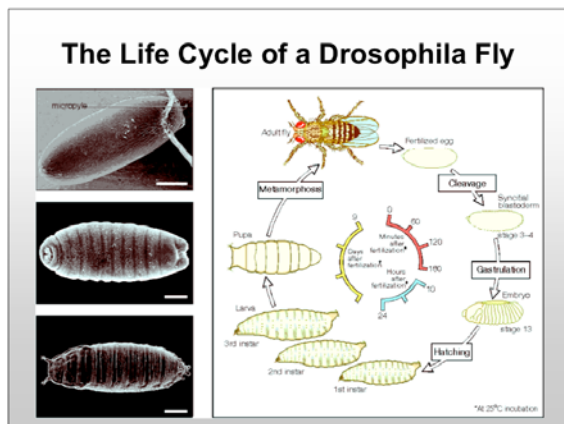


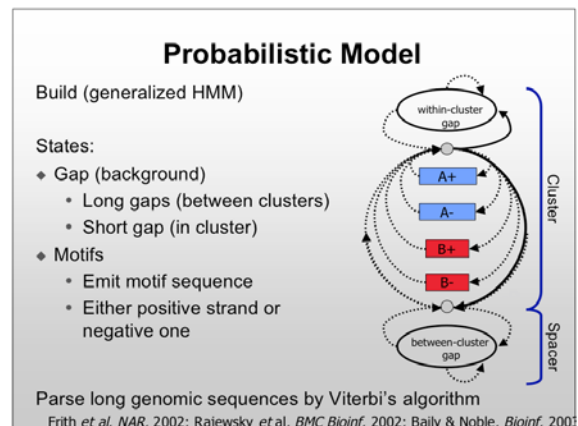
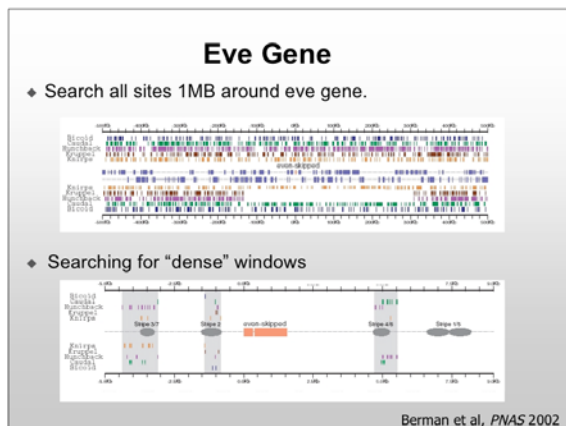
Effect at a Distance



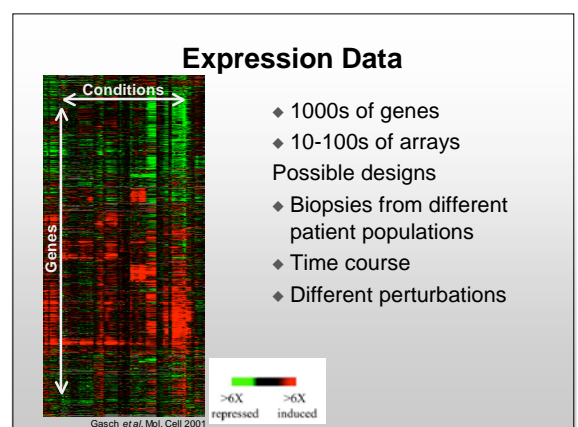
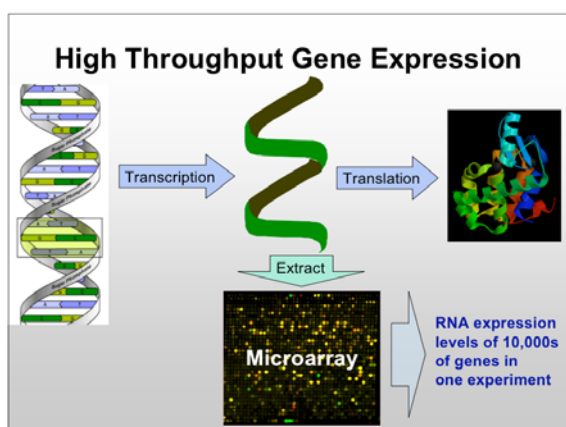
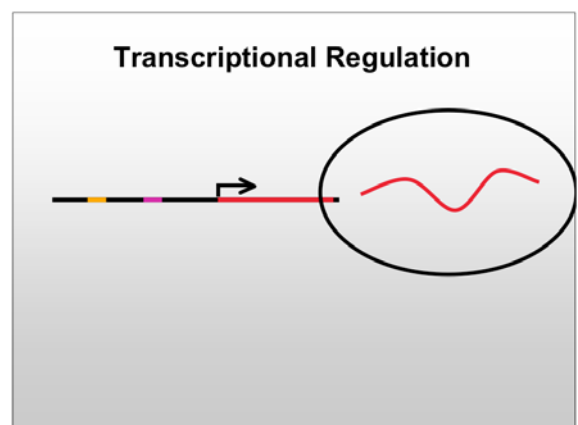
Eukaryotic Regulation

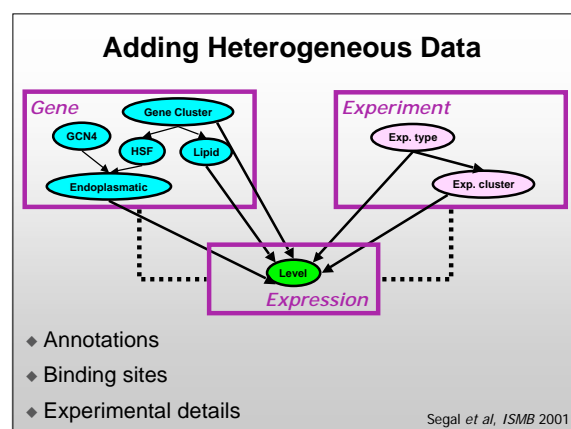
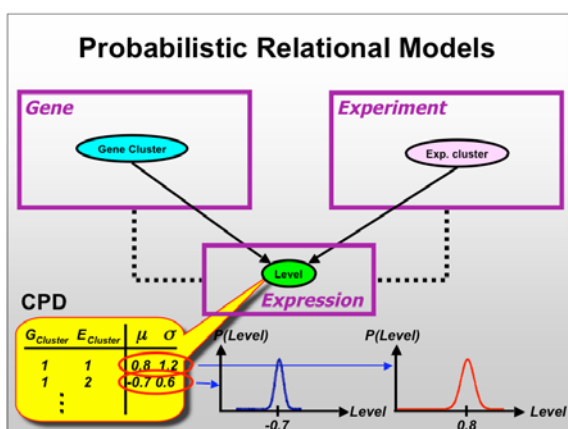
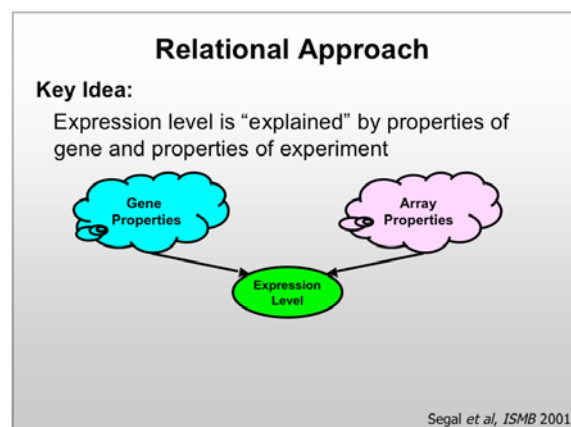
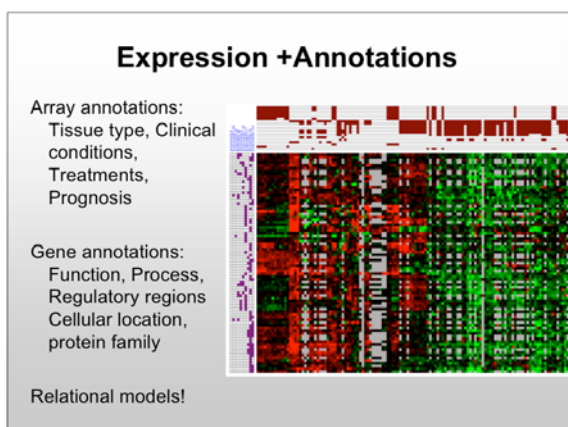
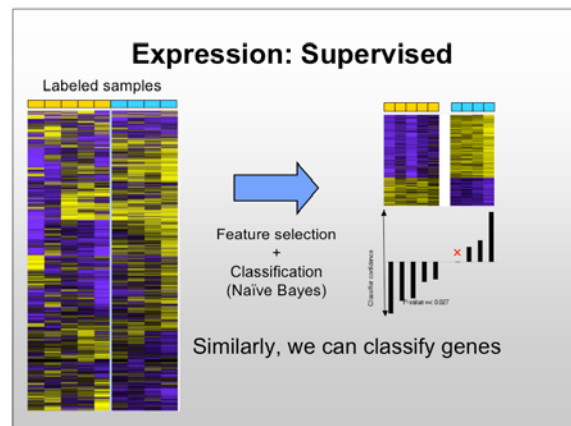
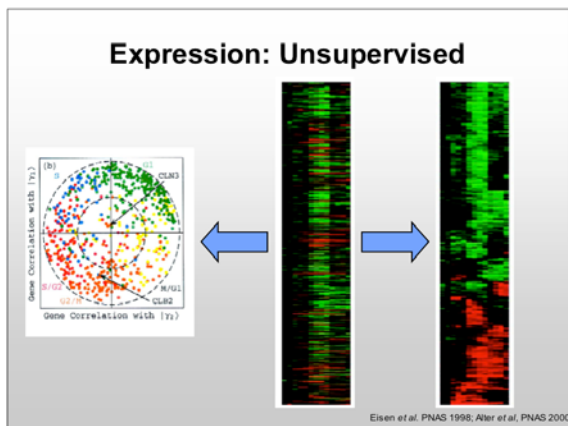


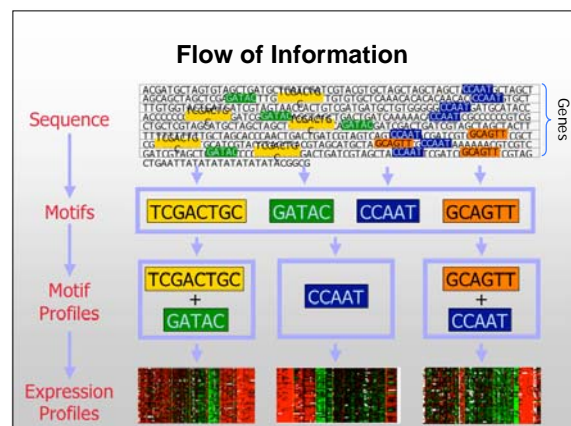
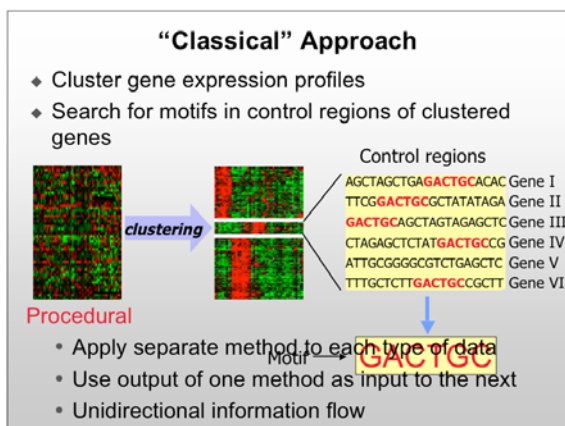
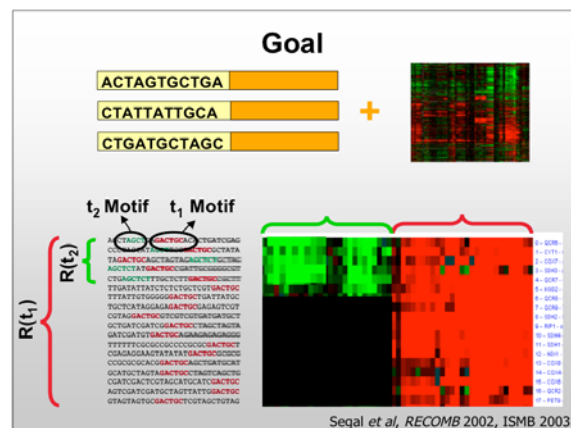
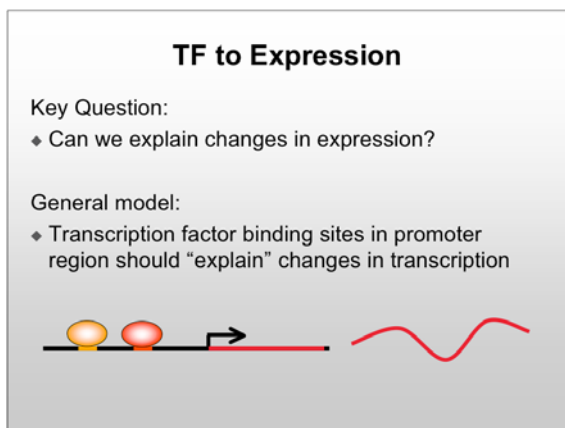
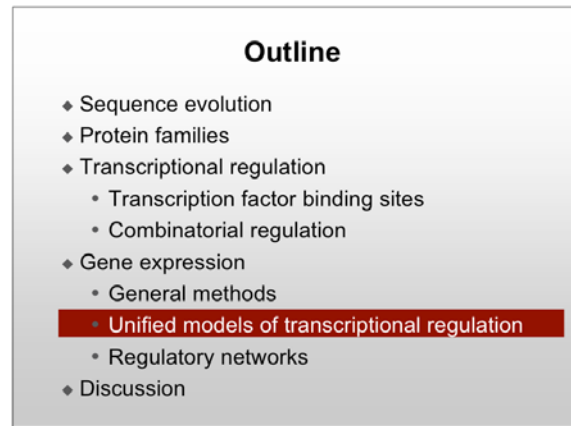
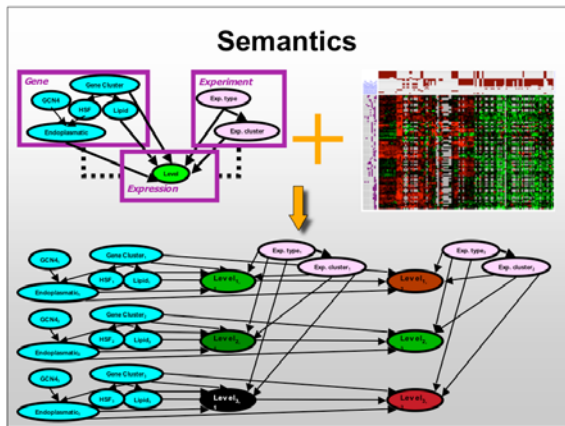


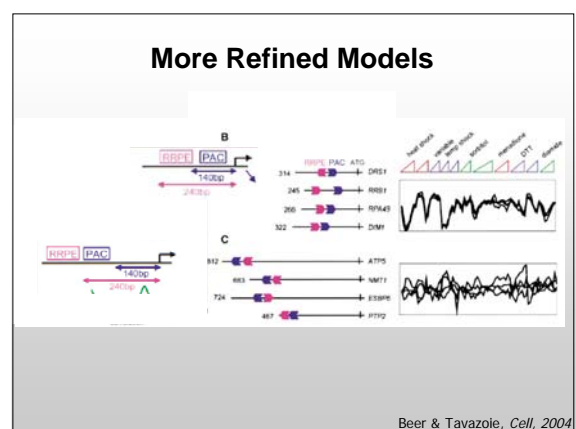
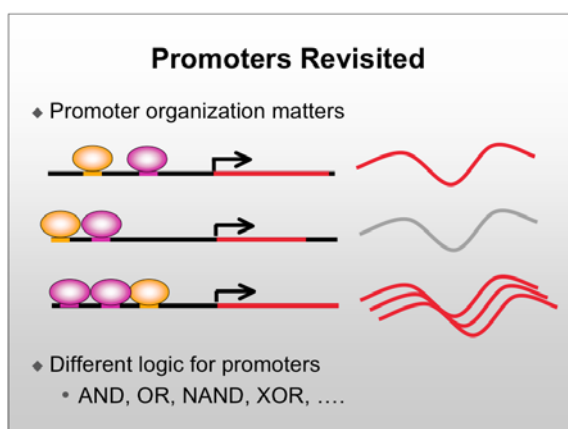
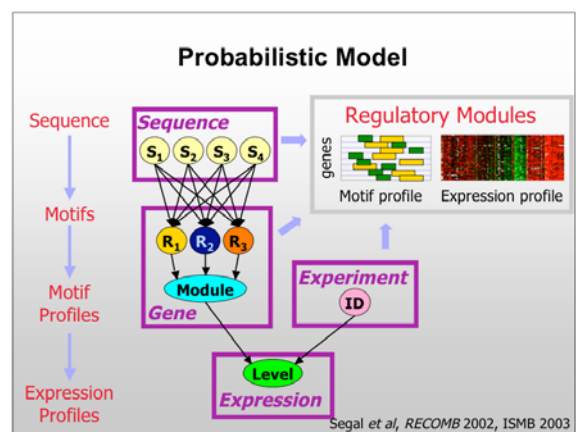
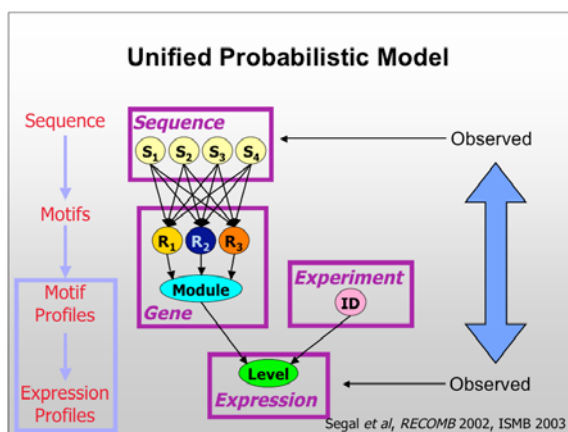
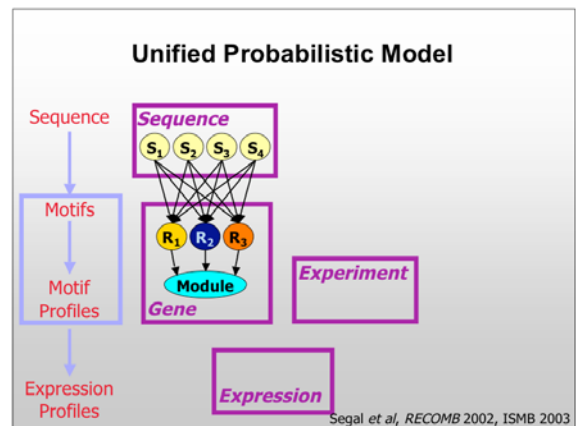
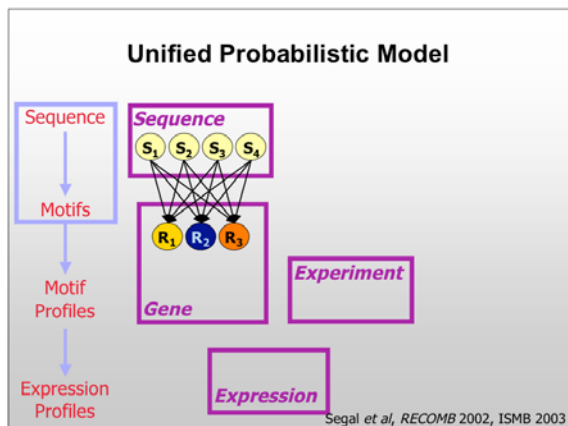


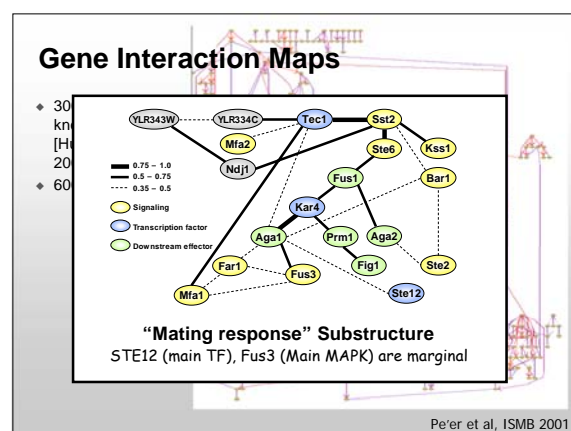
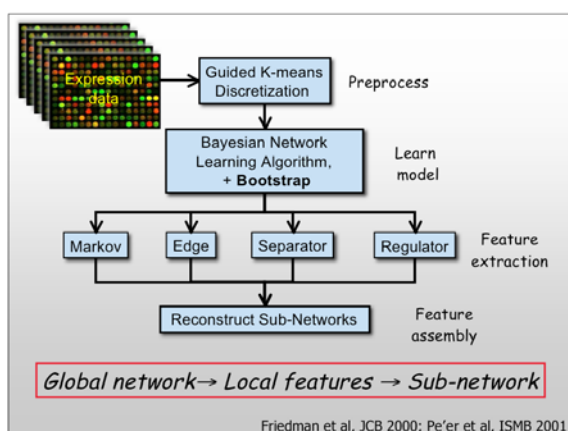
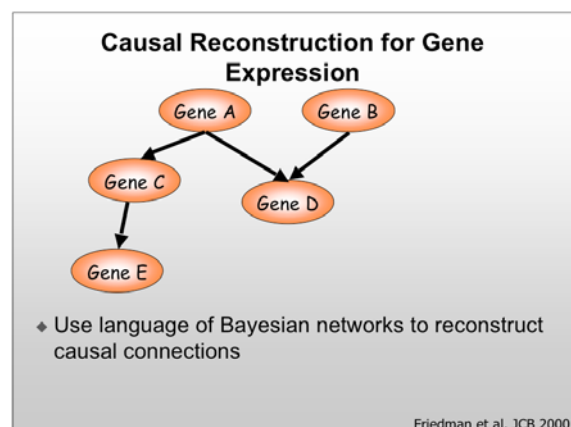
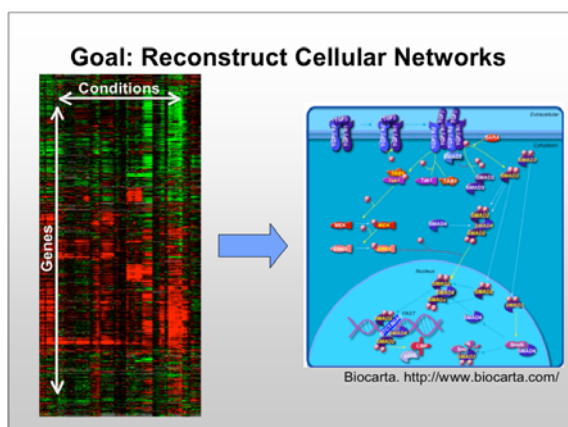
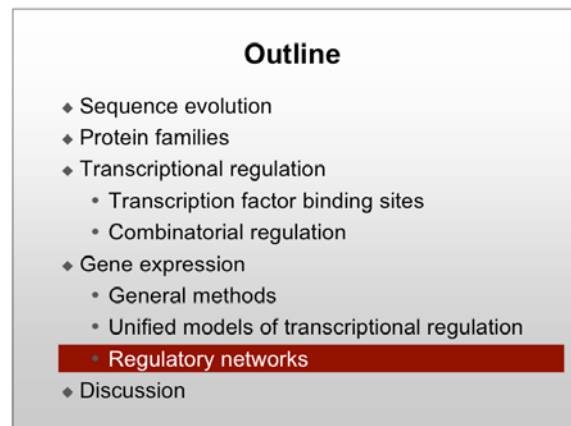
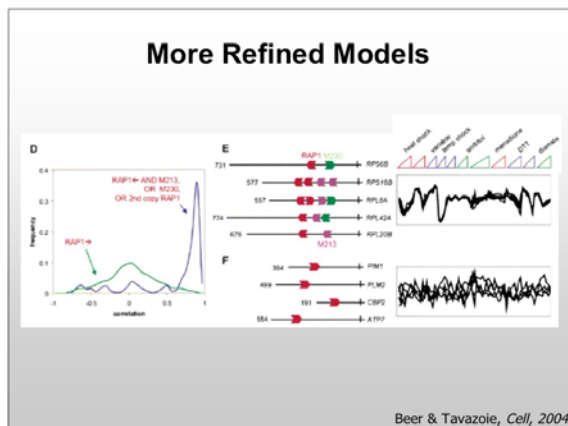
- ### Outline
- Sequence evolution
 - Protein families
 - Transcriptional regulation
 - Transcription factor binding sites
 - Combinatorial regulation
 - Gene expression
 - General methods**
 - Unified models of transcriptional regulation
 - Regulatory networks
 - Discussion



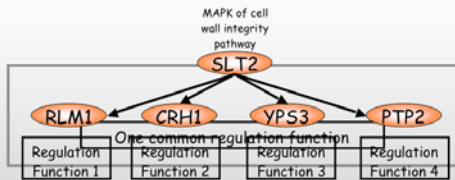






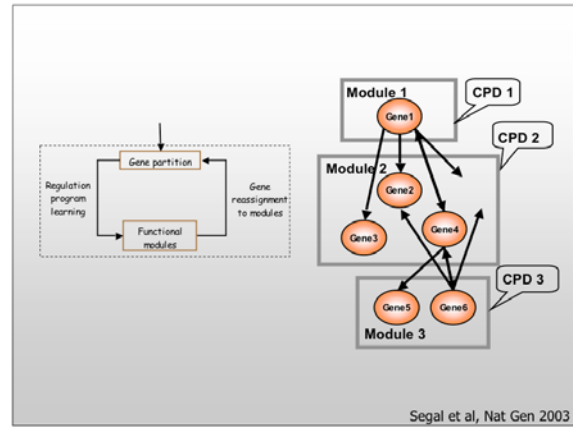


From Networks to Modules



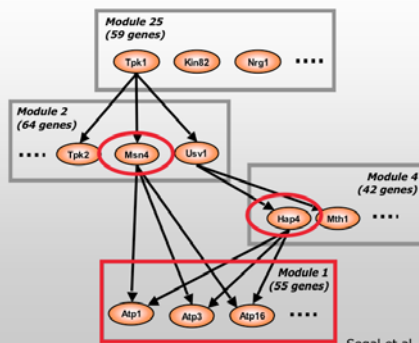
Idea: enforce common regulatory program

- Statistical robustness: Regulation programs are estimated from $m \times k$ samples
- Organization of genes into regulatory modules: Concise biological description

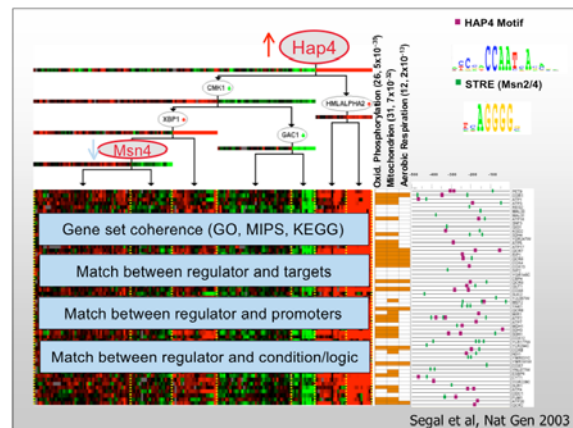


Segal et al, Nat Gen 2003

Learned Network (fragment)

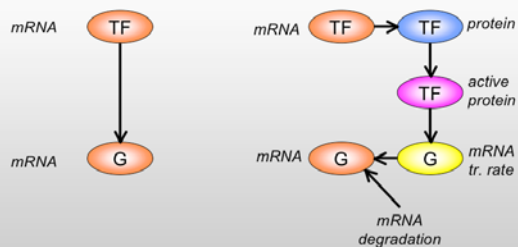


Segal et al, Nat Gen 2003



Segal et al, Nat Gen 2003

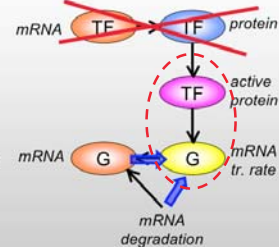
A Major Assumption

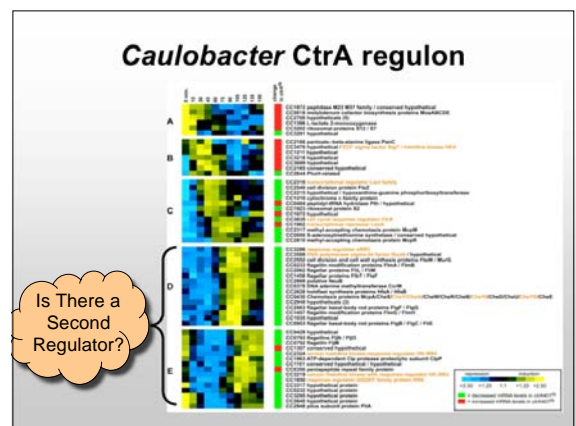
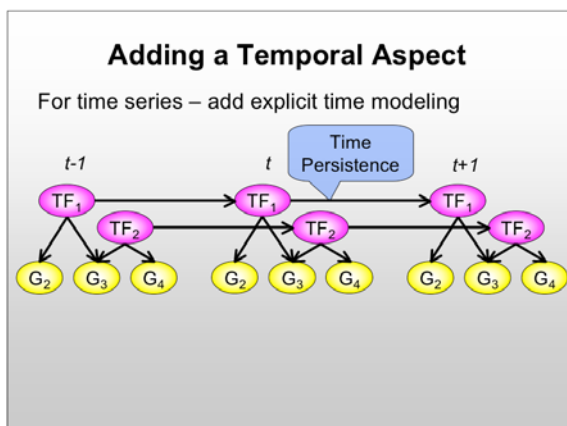
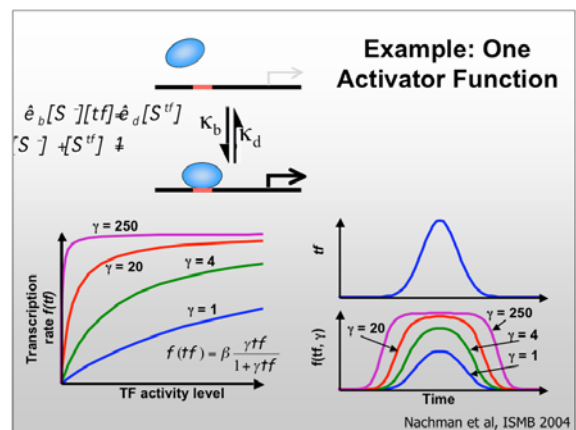
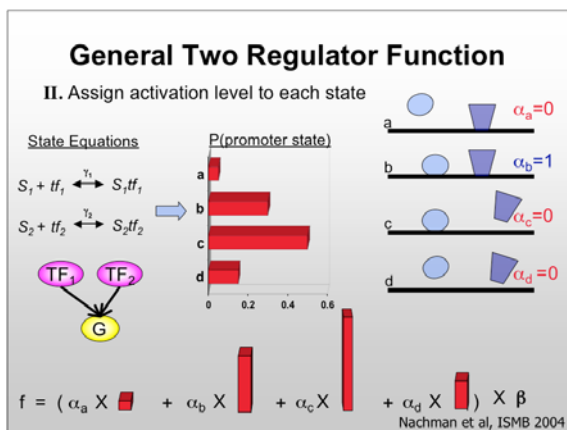
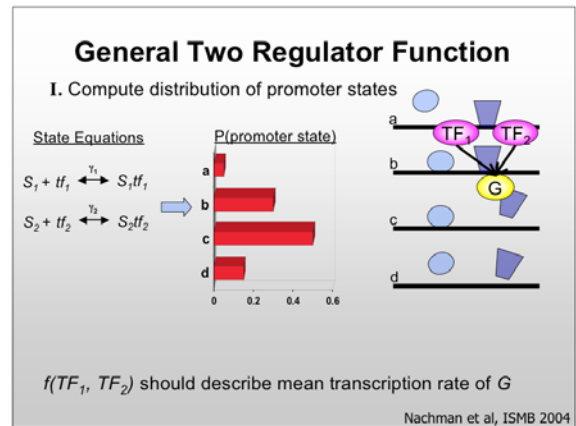
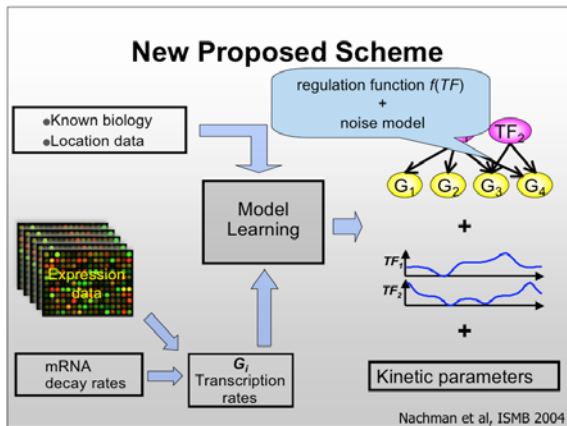


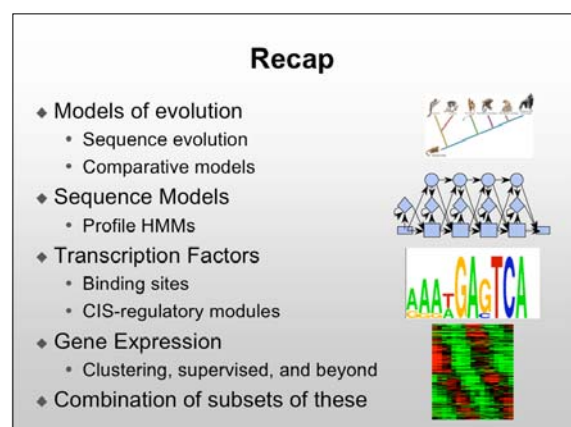
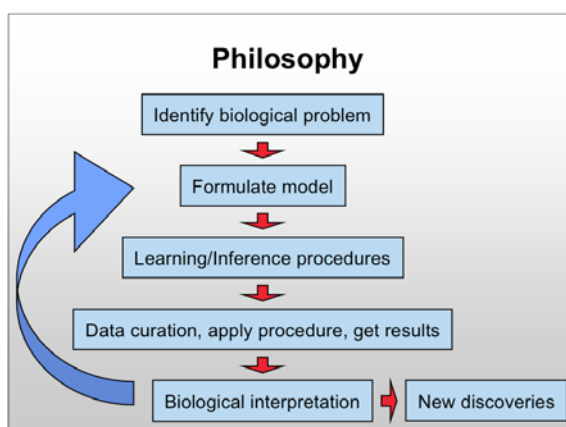
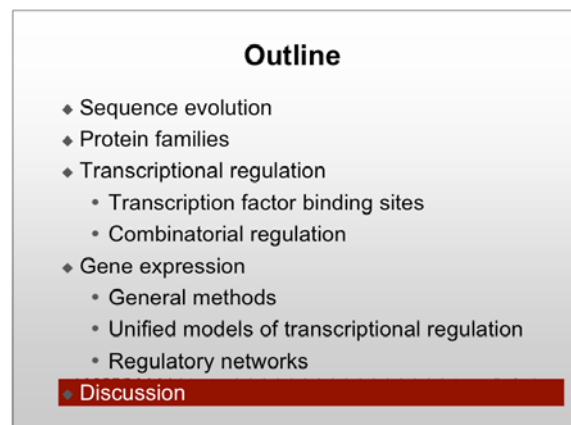
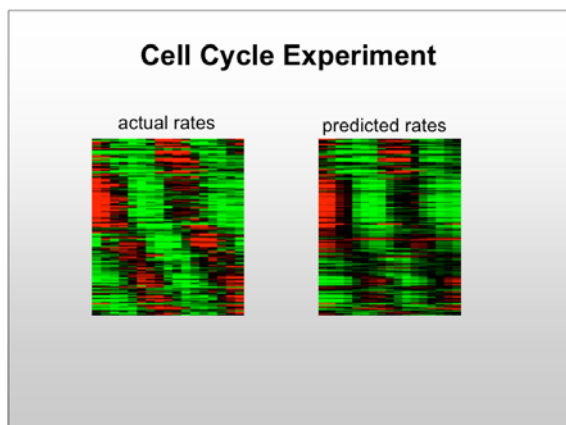
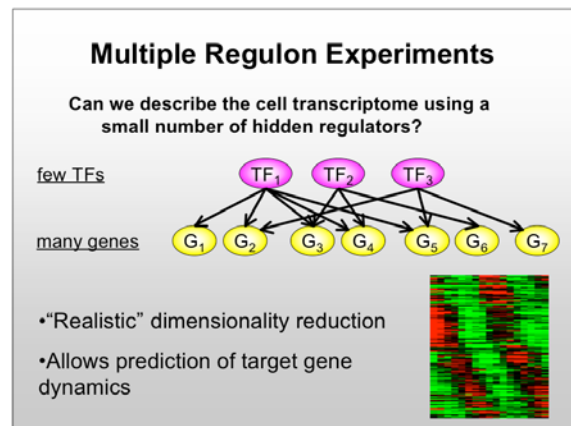
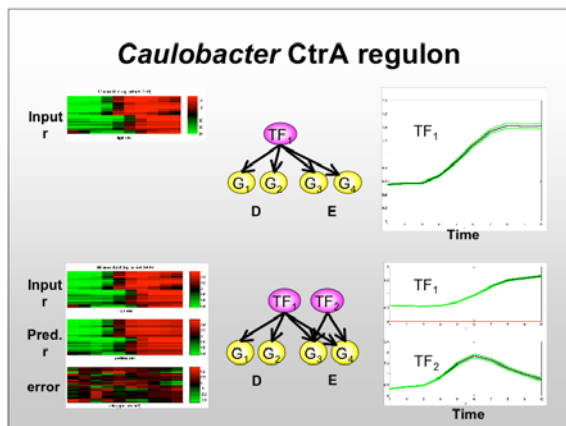
Realistic Regulation Modeling

- Model the closest connection

- Active protein levels are not measured
- Transcript rates are computed from expression data and mRNA decay rates







Additional Areas

- ◆ Gene finding
 - Extended HMMs + evolutionary models
- ◆ Analysis of genetic traits and diseases
 - Linkage analysis
 - SNPs, haplotypes, and recombination
- ◆ Interaction networks (protein-DNA, protein-protein)
 - Relational models
- ◆ Protein structure
 - 2nd-ary and 3rd-ary structure, molecular recognition

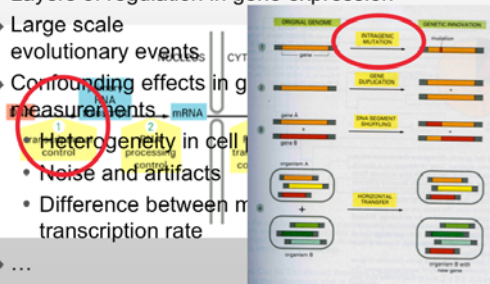
Take Home Message

- ◆ **Graphical models as a methodology**
 - Modeling language
 - Foundations & algorithms for learning
 - Allows to incorporate prior knowledge about biological mechanisms
 - Learning can reveal "structure" in data
- ◆ **Exploring unified system models**
 - Learning from heterogeneous data
 - ◆ Not simply combining conclusions
 - Combine weak evidence from multiple sources
 - ⇒ detect subtle signals
 - Get closer to **mechanistic** understanding of the signal

What About Biology?

We skimmed over many details

- ◆ Layers of regulation in gene expression
- ◆ Large scale evolutionary events
- ◆ Confounding effects in gene measurements
- ◆ Heterogeneity in cell
 - Noise and artifacts
 - Difference between mRNA transcription rate



The END