

Probabilistic Graphical Models in Systems Biology

Nir Friedman
Hebrew University

Includes slides by:

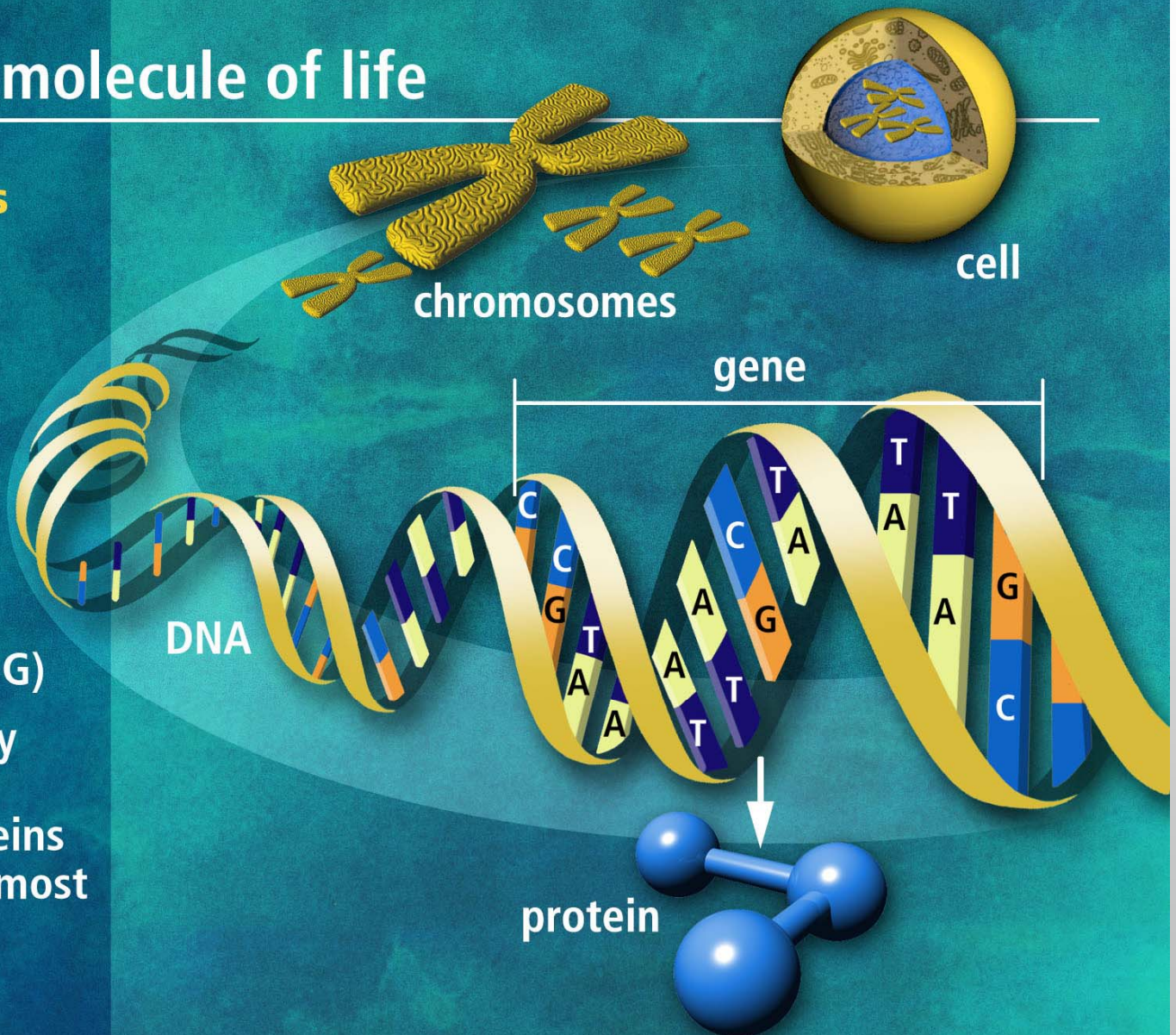
Yoseph Barash, Nebojsa Jojic, Ariel Jaimovich,
Tommy Kaplan, Daphne Koller, Iftach Nachman,
Dana Pe'er, Tal Pupko, Aviv Regev, Eran Segal

DNA the molecule of life

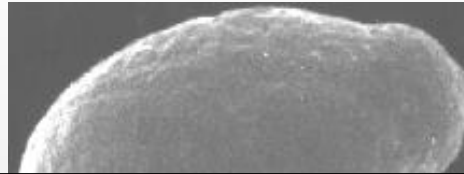
Trillions of cells

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

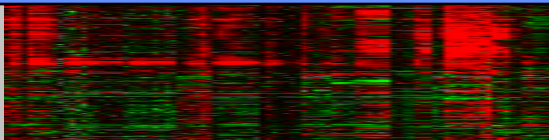


Challenges of The Post-Genome Era



High-throughput assays:

- Observations about **one** aspect of the system
- Often **noisy** and less reliable than traditional assays
- Provide partial account of the system



Challenges of The Post-Genome Era

Issues:

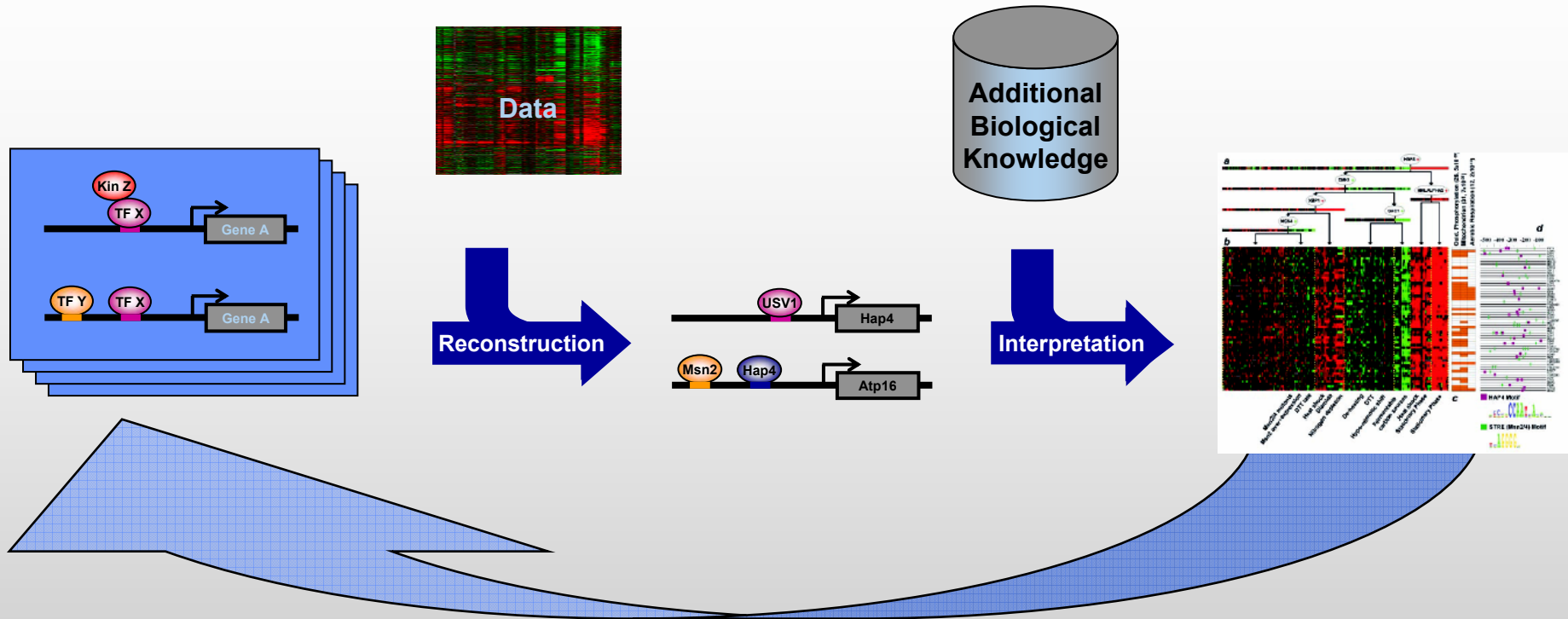
- ◆ Measurement noise
 - ⇒ Conclusions supported by more than one assay
- ◆ Each assay provides a view of a single aspect
 - ⇒ Combine multiple types of assays for more coherent reconstruction
- ◆ Combinatorial explosion of assay combinations
 - ⇒ Principles for integrating results from new assays

Solution Strategies

Procedural

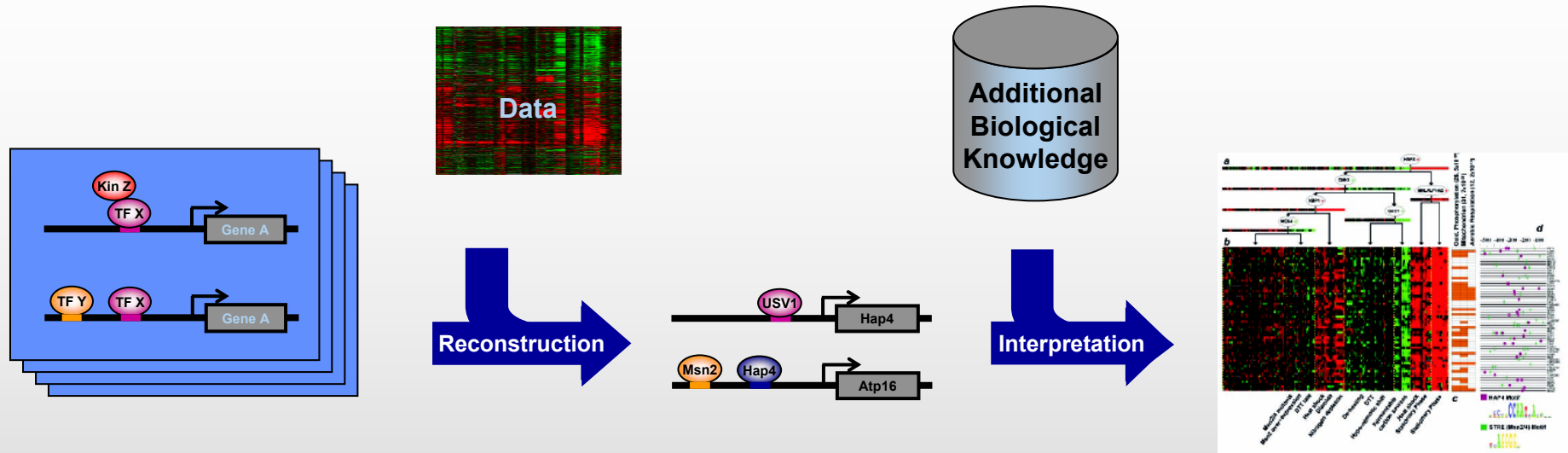
- ◆ Specify a set of steps for reaching biological conclusions from experimental data
 - Example
 - ◆ Cluster gene expression profiles
 - ◆ Search for enriched motif in each cluster
 - ◆ ...
- ◆ emphasis on the computational procedure and the order of data manipulation steps

Model Based Approach



- ◆ **Step 1:** define class of potential models
- ◆ **Step 2:** reconstruct a specific model
- ◆ **Step 3:** visualization & testable hypotheses
- ◆ Emphasis on the choice of model and how to use it
 - The data manipulation steps are derived from the model

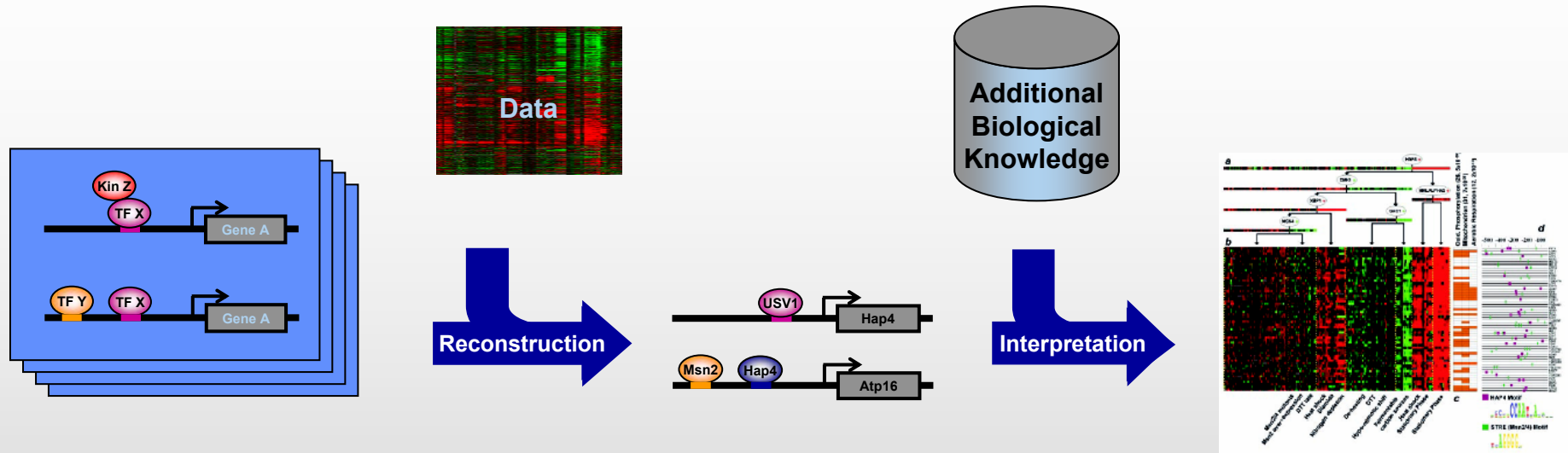
Model Based Approach



Representation – defining the class of models

- What entities to involve
- Model granularity
- Identifiably

Model Based Approach



Interpretation – what do they tell us about system

- Relation between components in the model to biological entities/mechanisms
- What predictions can be made with the model

Why Model-Based?

Declarative

- Explicit statement of the assumptions made
- Closer connection between biological principles and the solution
- Decouple the “what” (model) from the “how”

Flexibility

- Can use different computational approaches to perform the task specified by the model

Reusability

- Modifications & extensions are specified the level of the model

Stochastic Models

Use the probability theory to describe the system

- ◆ **State of the system:** assignment of value to all the attributes of all the relevant entities
- ◆ A **distribution** over these states describe which states are achievable and which ones are abnormal

Extensions:

- ◆ Modeling inputs: interventions, conditions
- ◆ Modeling outputs: phenotype, behavior, assays

Why Stochastic Models?

- ◆ Inherent noise in the system
- ◆ Uncertainty due to granularity of the model
- ◆ Noise in sensors
- ◆ Imperfect modeling --- noise as slack variable

What Can We Do with a Model?

◆ Inference

- Set some evidence, compute posterior over unobserved variables

◆ Estimation/Learning

- “Fill in the gaps” in the model based on empirical data

The Representation Hurdle

- ◆ Joint distributions grow large
 - Exponential in the number of attributes
 - Problem for inference & learning

We need to find **compact** representation

Strategy:

- Impose constraints
- Exploit these constraints for compact representation

Probabilistic Graphical Models

- ◆ Language(s) for representing complex joint distributions
- ◆ Generic methods for performing tasks with these representations

In this tutorial we will examine these in the context of modeling biological systems.

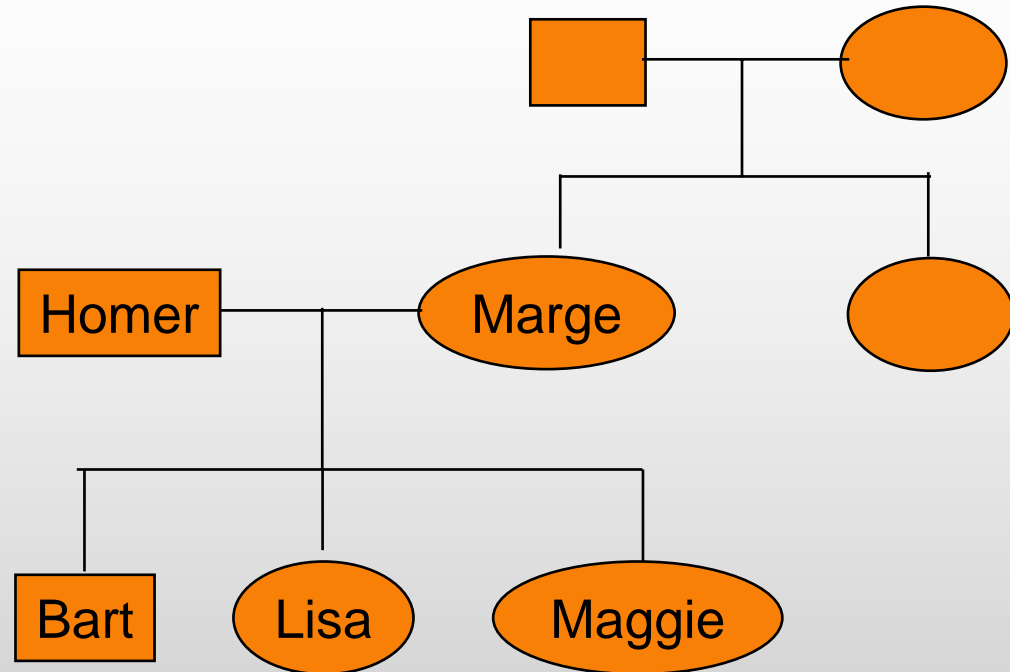
Outline

- ◆ Introduction
- ◆ **Bayesian Networks**
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

Bayesian Networks by Example

Example: Pedigree

- ◆ A node represents an individual's genotype



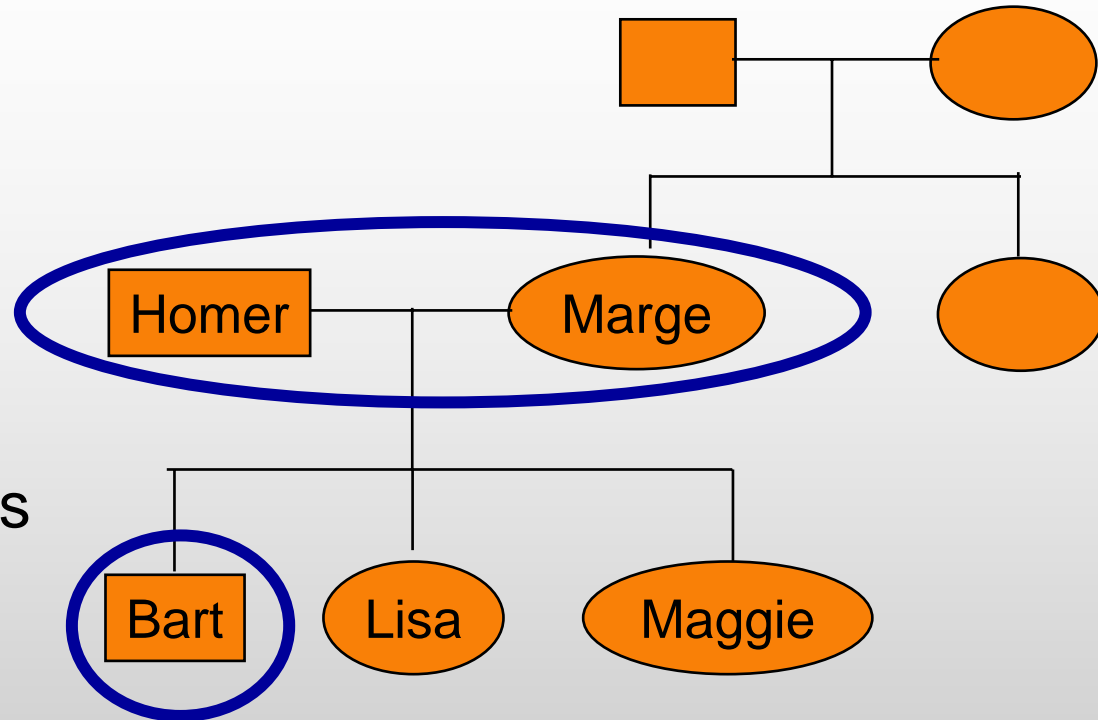
Joint distribution

$$\begin{aligned} &P(G_{\text{Bart}}, G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &= P(G_{\text{Bart}} \mid G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &P(G_{\text{Lisa}} \mid G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &\dots \end{aligned}$$

Bayesian Networks by Example

Modeling assumption:

- ◆ Ancestors can effect descendants' genotype only by passing genetic materials through intermediate generations

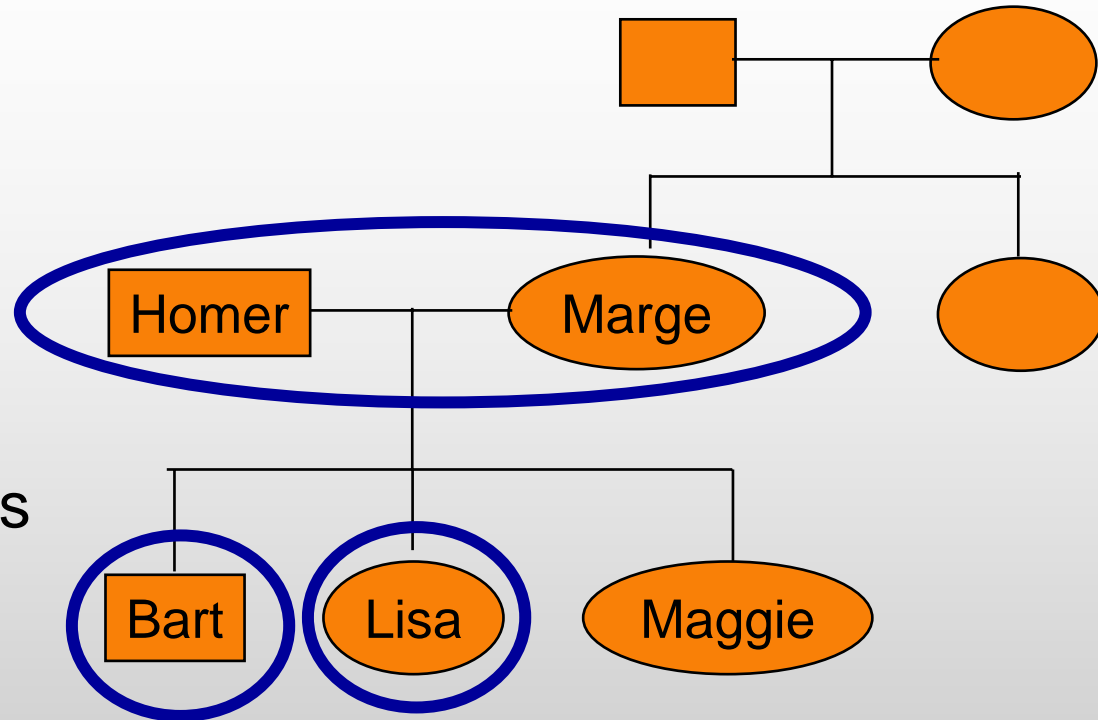


$$\begin{aligned} &P(G_{\text{Bart}}, G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &= P(G_{\text{Bart}} \mid G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &P(G_{\text{Lisa}} \mid G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &\dots \end{aligned}$$

Bayesian Networks by Example

Modeling assumption:

- ◆ Ancestors can effect descendants' genotype only by passing genetic materials through intermediate generations



$$P(G_{\text{Bart}}, G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots)$$

$$= P(G_{\text{Bart}} \mid G_{\text{Homer}}, G_{\text{Marge}})$$

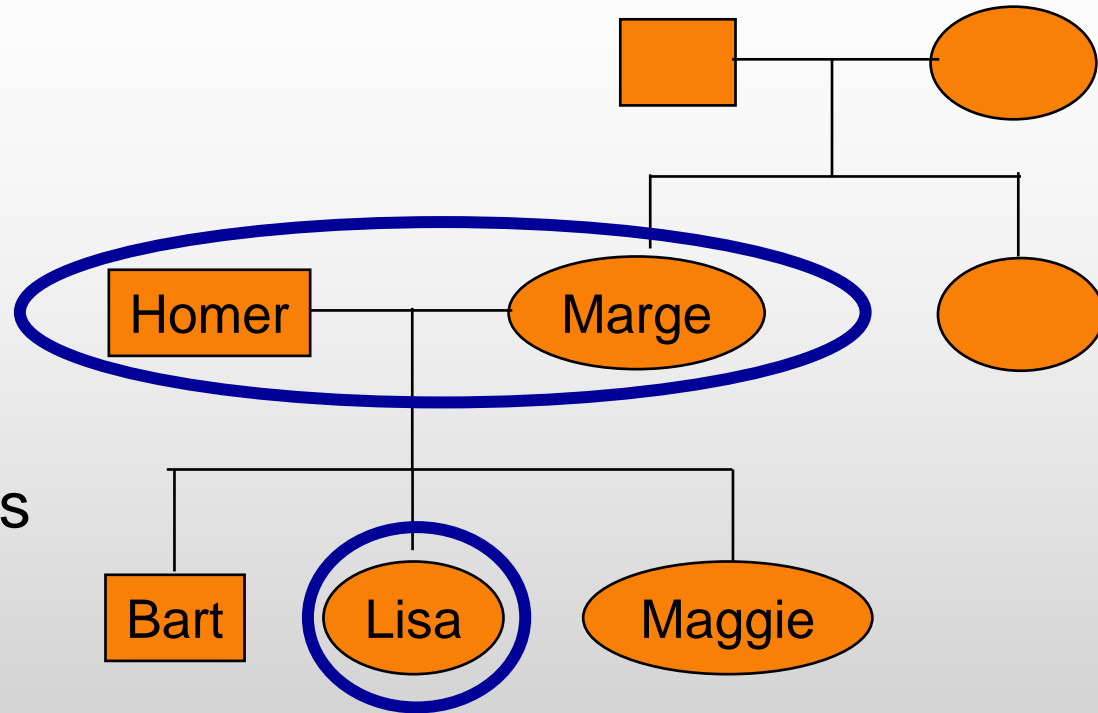
$$P(G_{\text{Lisa}} \mid G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots)$$

...

Bayesian Networks by Example

Modeling assumption:

- ◆ Ancestors can effect descendants' genotype only by passing genetic materials through intermediate generations

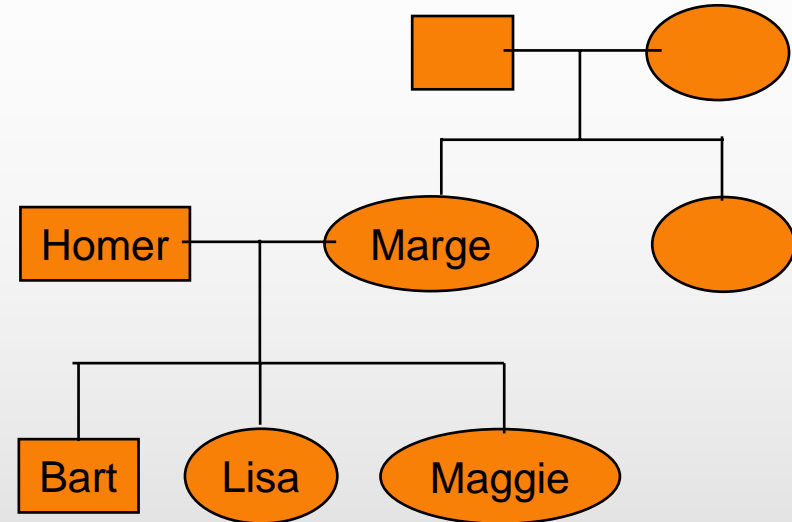


$$\begin{aligned} &P(G_{\text{Bart}}, G_{\text{Lisa}}, G_{\text{Maggie}}, G_{\text{Homer}}, G_{\text{Marge}}, \dots) \\ &= P(G_{\text{Bart}} \mid G_{\text{Homer}}, G_{\text{Marge}}) \\ &P(G_{\text{Lisa}} \mid G_{\text{Homer}}, G_{\text{Marge}}) \\ &\dots \end{aligned}$$

Bayesian Networks by Example

Extending this argument, we can derive a functional form for general pedigrees

Descendants



$$P(G_1, G_2, \dots) = \left(\prod_{j \in \text{Ancestors}} P(G_j) \right) \left(\prod_{i \in \text{Descendants}} P(G_i \mid G_{\text{father}(i)}, G_{\text{mother}(i)}) \right)$$

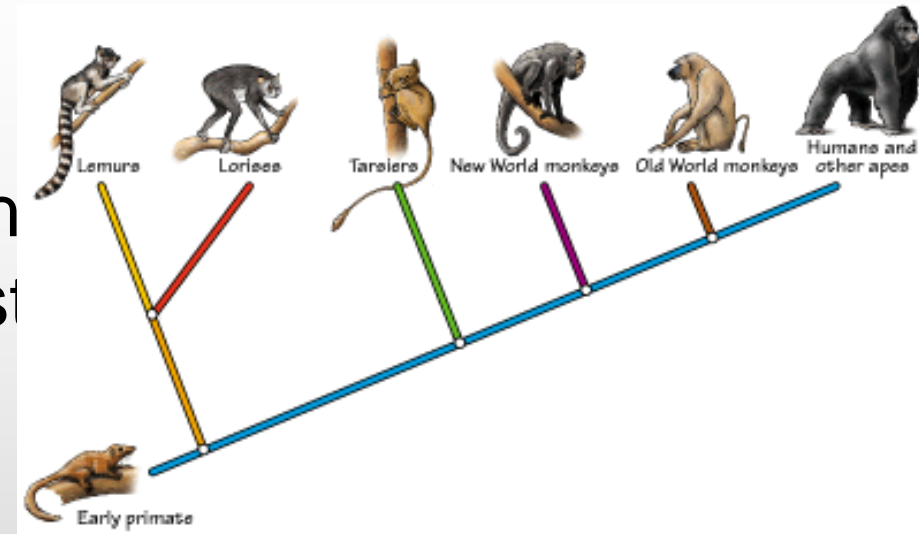
Probability of random genotype in population

Probability of genetic transmission within family

Bayesian Networks by Example II

Sequence evolution

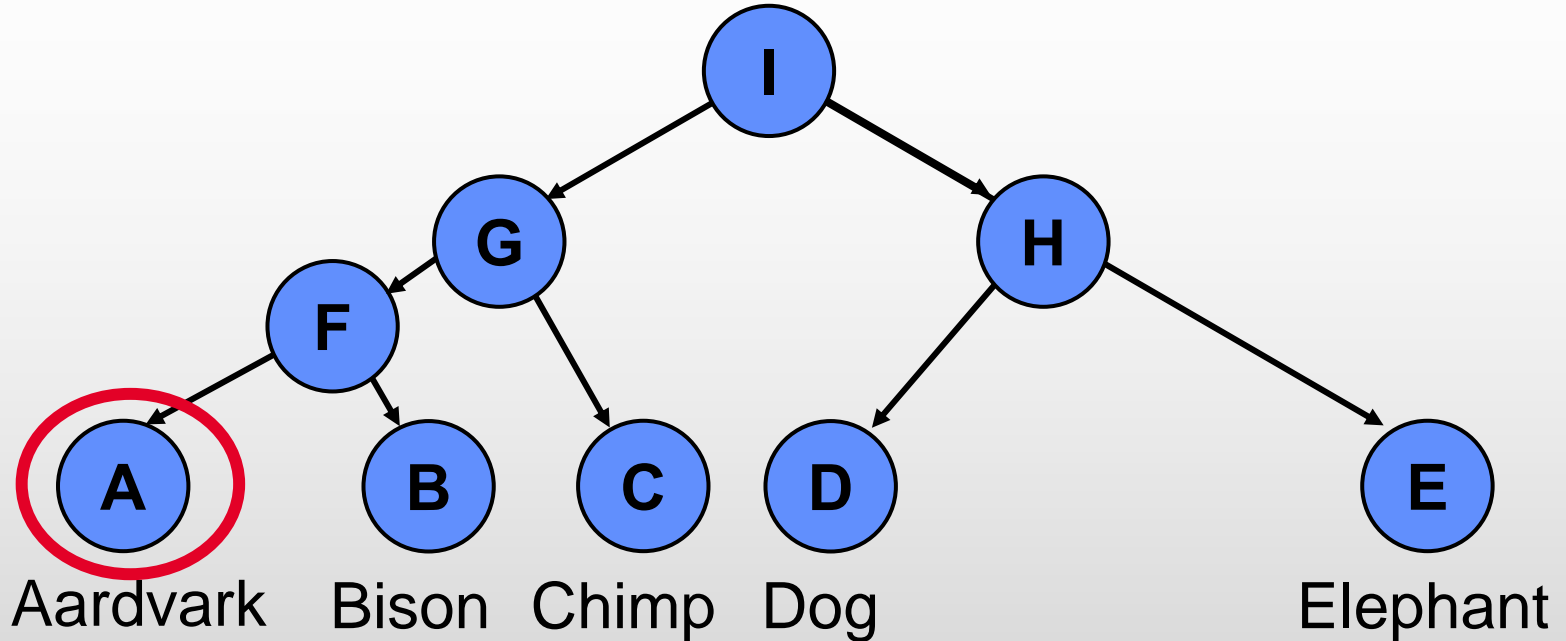
- ◆ Each random variable is the sequence of a taxa (ancestor or current day)



Assumption (neutral changes):

- ◆ Past history does not affect how the sequence will change in the future

Bayesian Networks by Example II



$$P(S_A, S_B, S_C, \dots, S_I) = P(S_A | S_B, S_C, \dots, S_I)$$

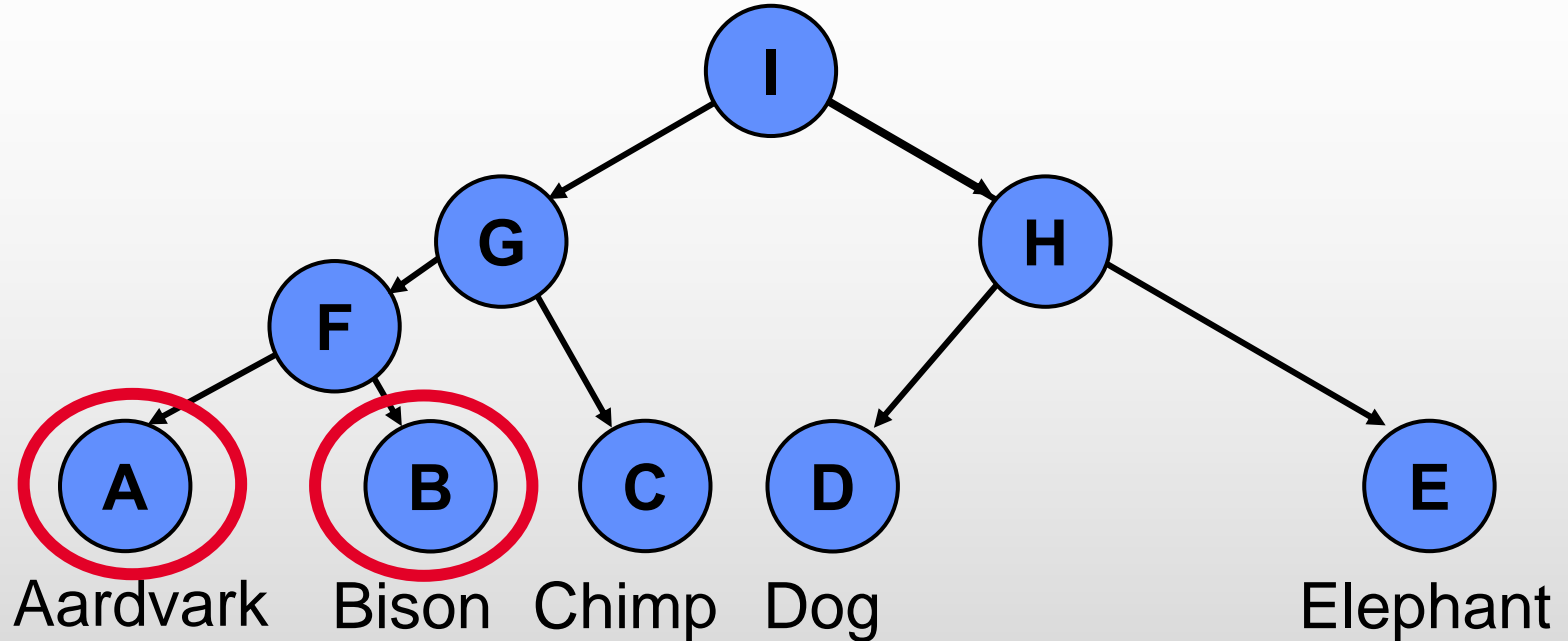
$$P(S_B | S_C, \dots, S_I)$$

...

$$P(S_F | S_G, S_H, S_I)$$

...

Bayesian Networks by Example II



$$P(S_A, S_B, S_C, \dots, S_I) = P(S_A | S_F)$$

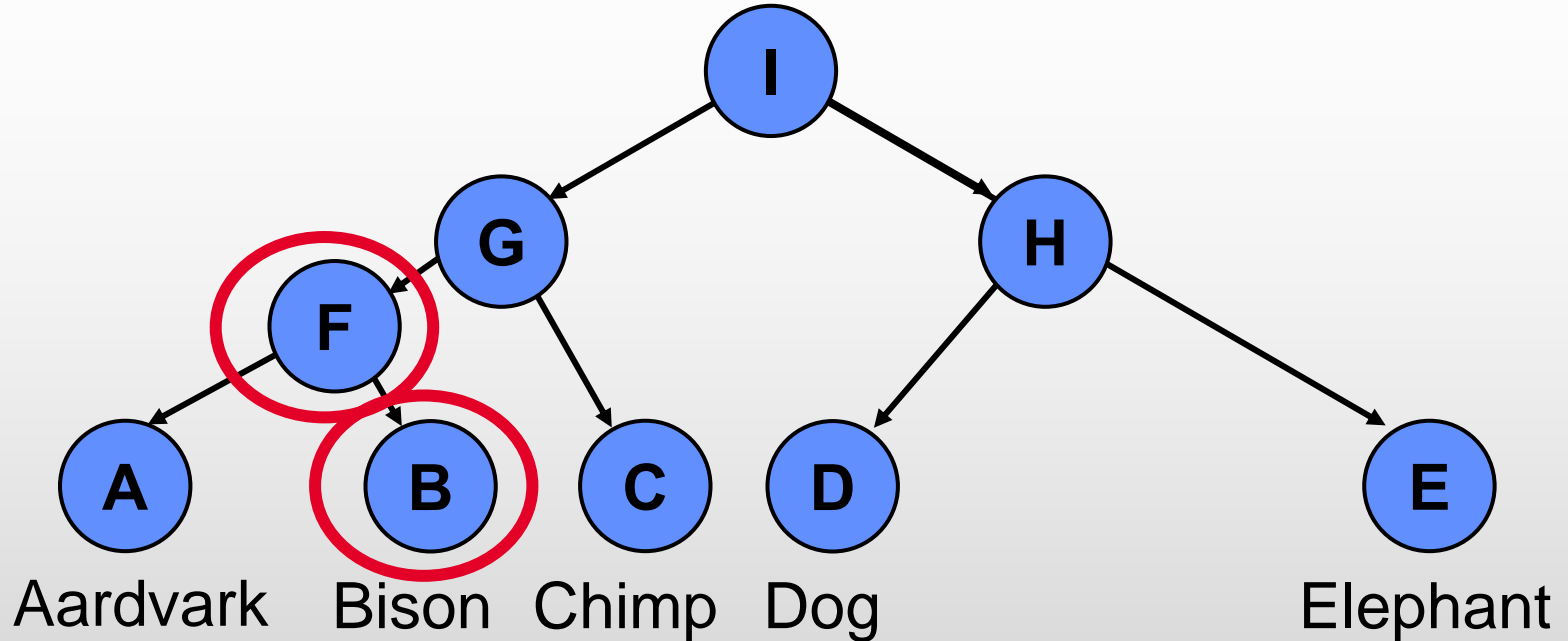
$$P(S_B | S_C, \dots, S_I)$$

...

$$P(S_F | S_G, S_H, S_I)$$

...

Bayesian Networks by Example II



$$P(S_A, S_B, S_C, \dots, S_I) = P(S_A | S_F)$$

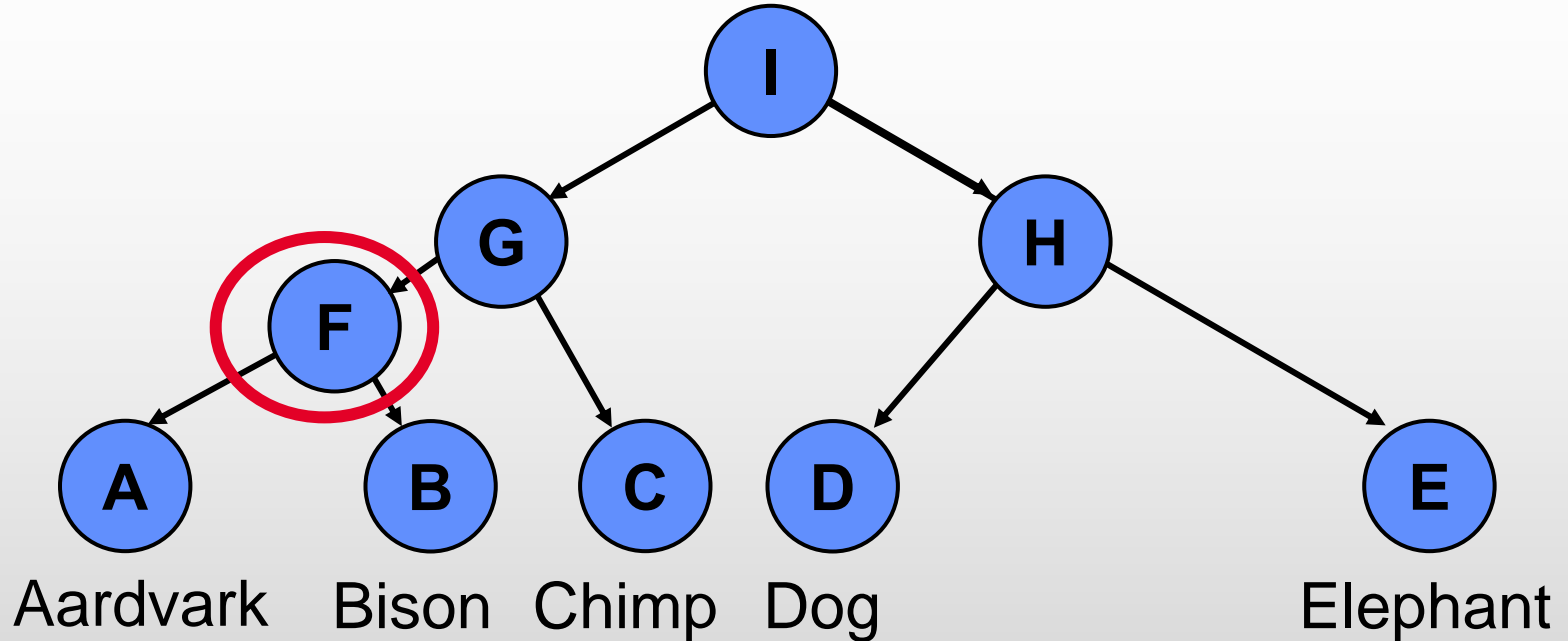
$$P(S_B | S_F)$$

...

$$P(S_F | S_G, S_H, S_I)$$

...

Bayesian Networks by Example II



$$P(S_A, S_B, S_C, \dots, S_I) = P(S_A | S_F)$$

$$P(S_B | S_F)$$

...

$$P(S_F | S_G)$$

...

Probability of mutations
over the given time
period

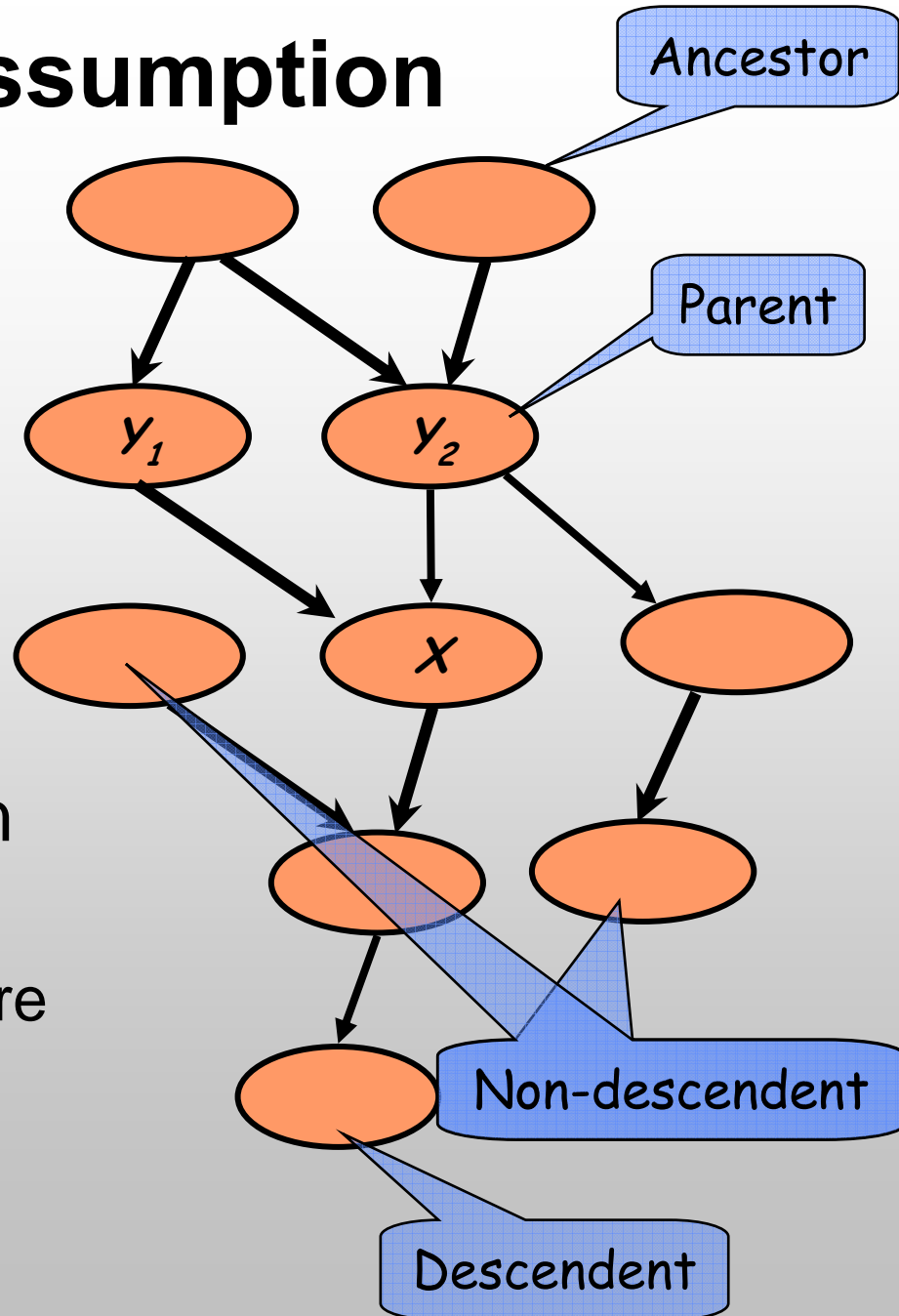
Markov Assumption

Generalizing to DAGs:

- ◆ A child is **conditionally independent** from its non-descendants, given the value of its parents

Often a natural assumption for **causal** processes

- if we believe that we capture the relevant state of each intermediate stage



Bayesian Networks

$$P(A, B, C, D, E) = P(A) \\ P(B | A) \\ P(C | A, B) \\ P(D | A, B, C) \\ P(E | A, B, C, D)$$

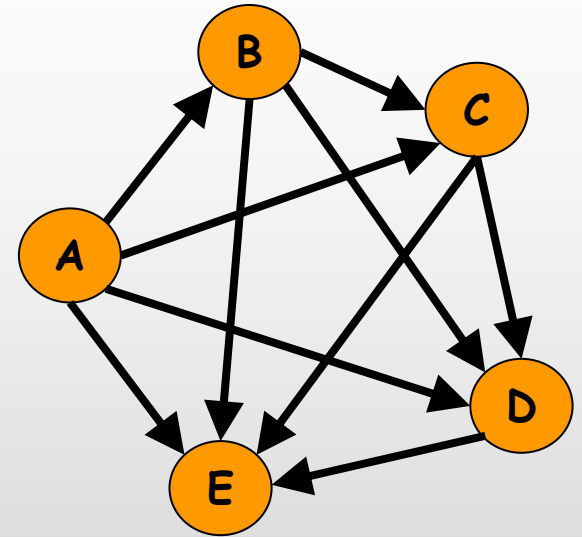
Bayesian Networks

$Ind(C ; B | A)$

$Ind(D ; A, B | C)$

$Ind(E ; B, C | A, D)$

...



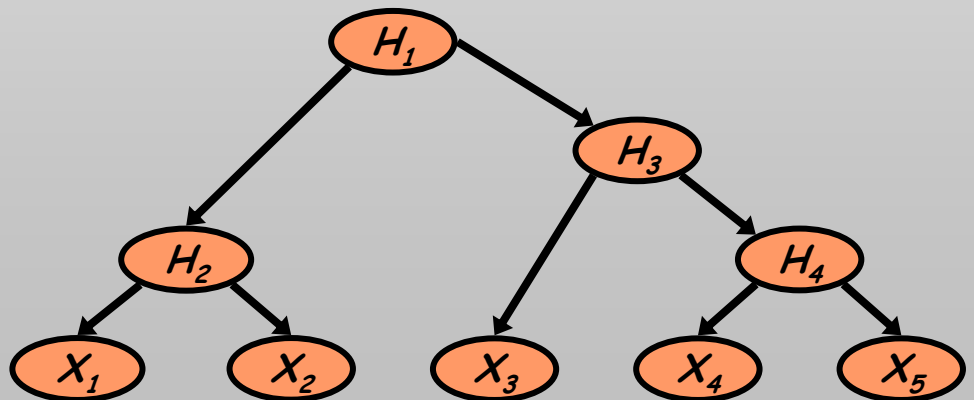
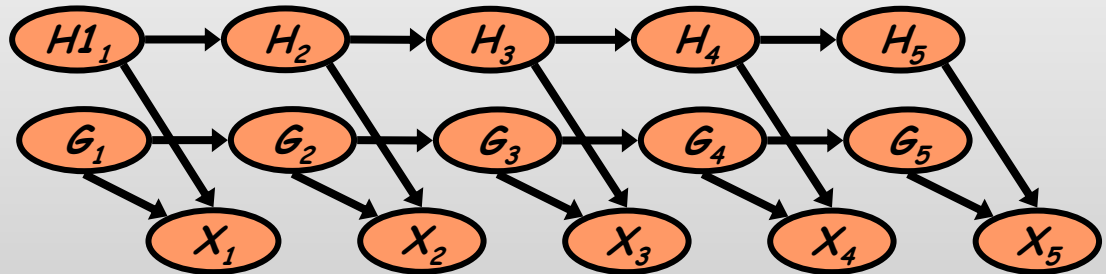
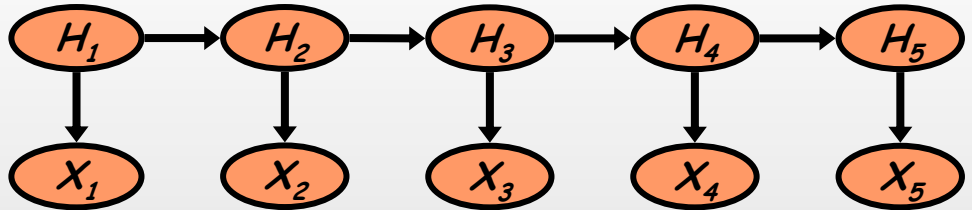
$$P(A, B, C, D, E) = P(A) \\ P(B | \cancel{A}) \\ P(C | \cancel{A}, \cancel{D}) \\ P(D | \cancel{A}, \cancel{B}, \cancel{C}) \\ P(E | \cancel{A}, \cancel{B}, \cancel{C}, \cancel{D})$$

Bayesian Networks

- ◆ Flexible language to capture a range
 - Maximal independence \rightarrow Full dependence
- ◆ Formal correspondence between
 - Acyclic directed graph structure
 - Factorization of joint distribution as a product of conditional probabilities
 - A set of (conditional) independence statements

Example Structures

- ◆ Markov chain
- ◆ Hidden Markov Model (HMM)
- ◆ Factorial HMM
- ◆ Tree



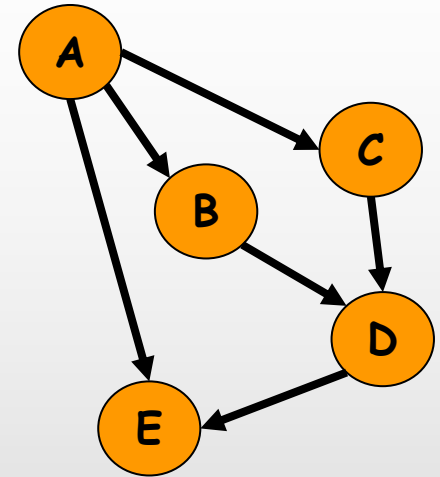
Local Probability Models

Bayesian Network Structure

⇒ Simpler product form

$$P(A, B, C, D, E) =$$

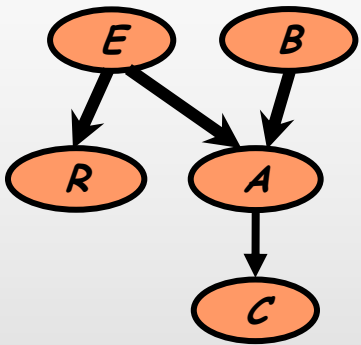
$$P(A) P(B | A) P(C | A) P(D | B, C) P(E | A, D)$$



To specify a distribution we need to supply these conditional probabilities

◆ Describe to “local” stochastic effects

Bayesian Network Semantics



Qualitative part

conditional
independence
statements
in BN structure

Quantitative part

+ local
probability
models =

Unique joint
distribution
over domain

Compact & efficient representation:

◆ nodes have $\leq k$ parents $\Rightarrow O(2^k n)$ vs. $O(2^n)$
params

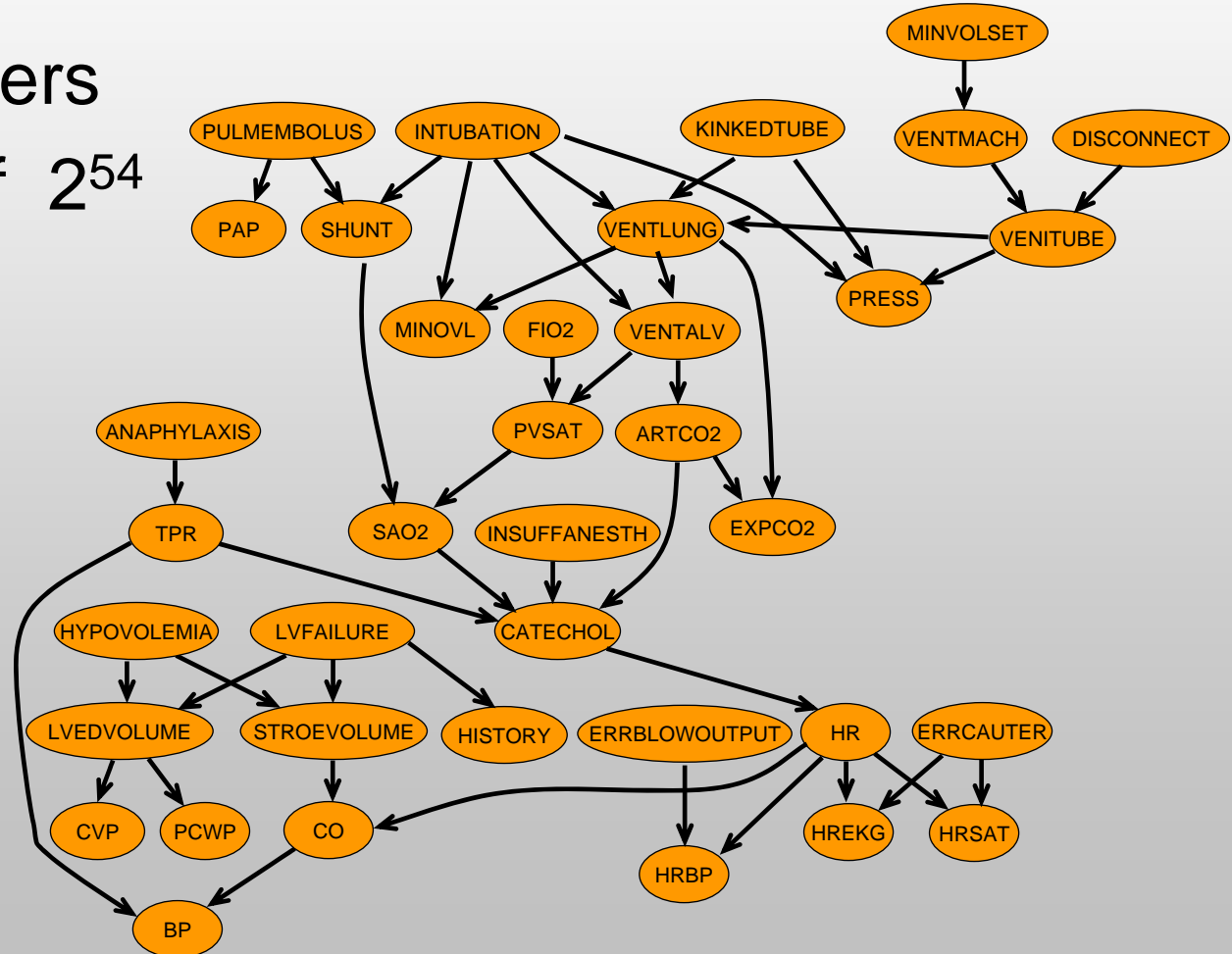
Example: “ICU Alarm” network

Domain: Monitoring Intensive-Care Patients

◆ 37 variables

◆ 509 parameters

...instead of 2^{54}



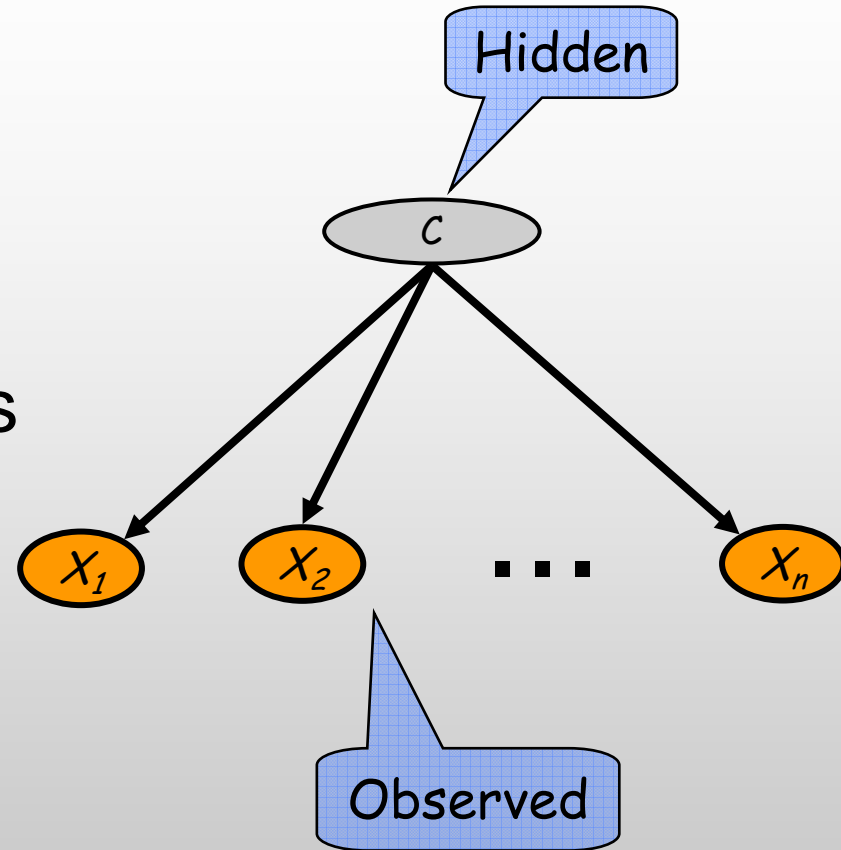
Hidden Variable(s)

A simple model of clustering

- ◆ C - gene's cluster
- ◆ X_1, \dots, X_n - expression of the gene in different experiments

Independence assumption:

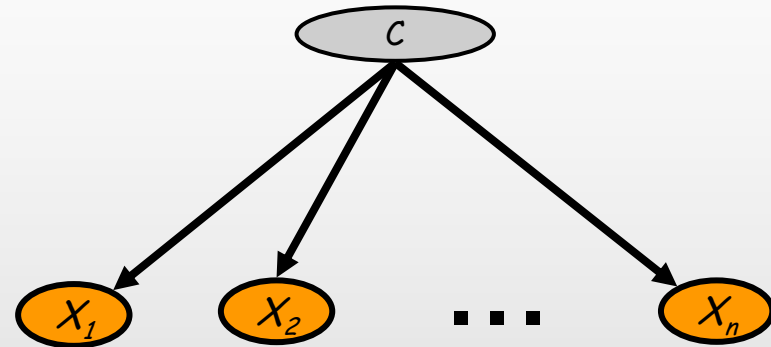
- ◆ $I(X_i; X_j | C)$



Hidden Variable(s)

Marginal distribution:

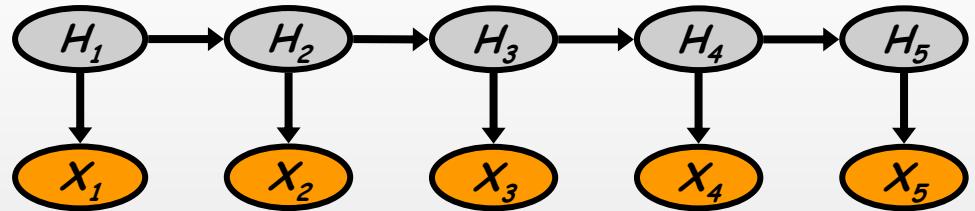
$$P(X_1, \dots, X_n) = \sum_c P(X_1 | c) \cdots P(X_n | c) P(c)$$



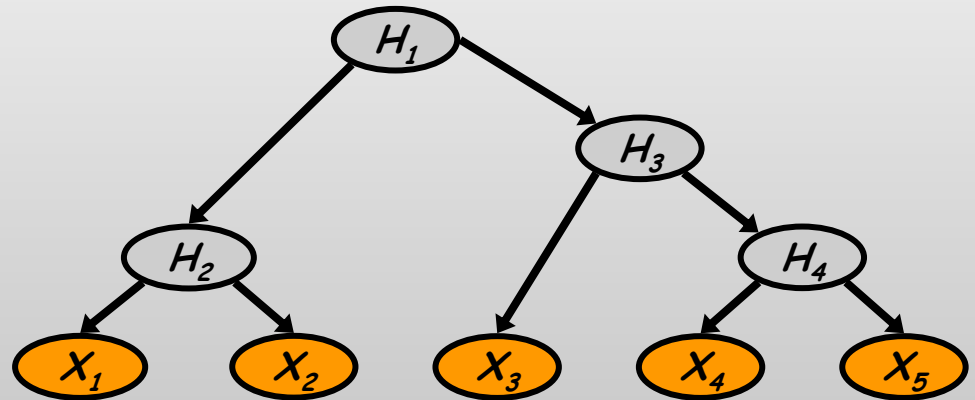
- ◆ Compact representation
 $(n+1)k$ vs. 2^n params
- ◆ No conditional independencies in the **marginal** distribution
- ◆ The variable C “channels” the dependencies between observed variables

Hidden Variables

Hidden Markov Model



Phylogenetic Trees



- ◆ The topology of hidden variables poses different constraints on the marginal distribution

Inference - Queries

◆ Posterior probabilities

- Probability of any event given any evidence

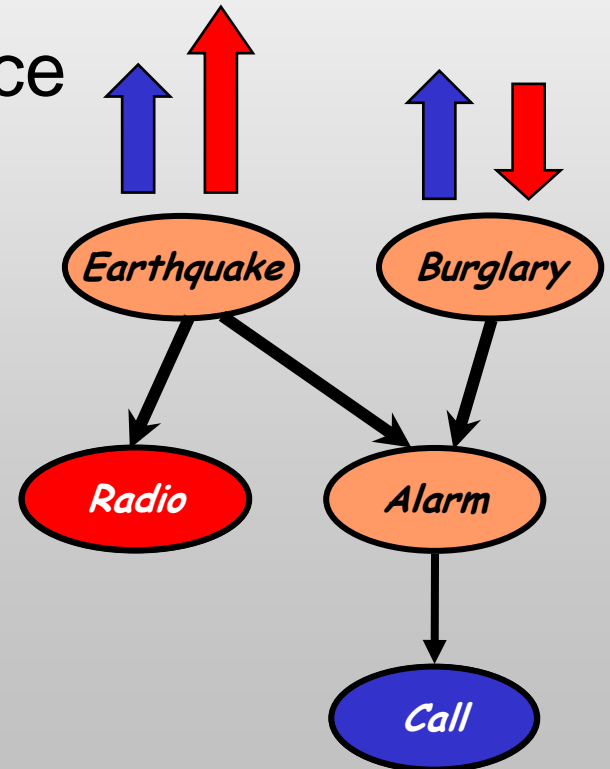
◆ Most likely explanation

- Scenario that explains evidence

◆ Rational decision making

- Maximize expected utility
- Value of Information

◆ Effect of intervention



Inference - Algorithms

Complexity:

- Worst case - exponential cost

Yet,

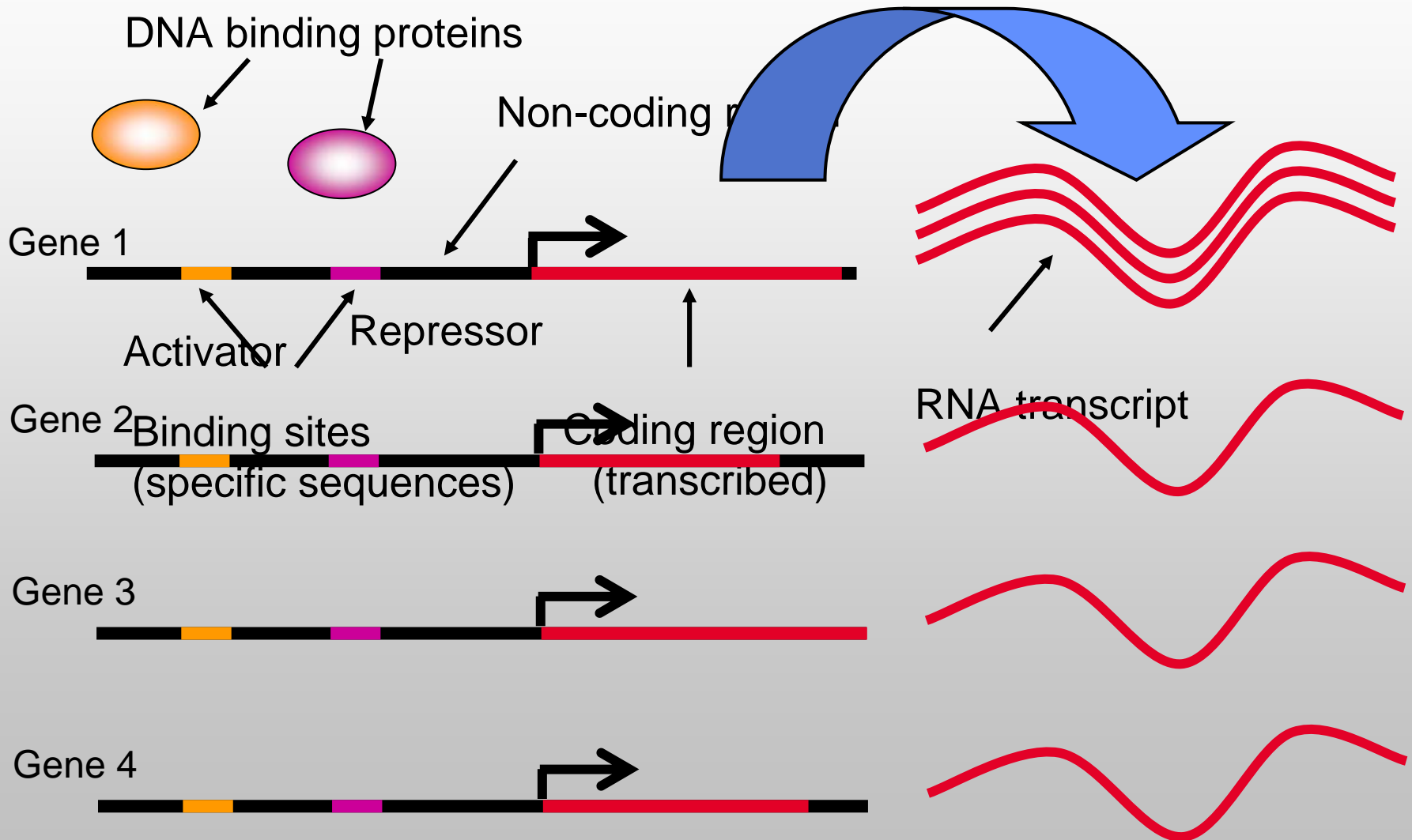
- ◆ Generic exact inference algorithms based on dynamic programming
 - Efficient in some network topologies
- ◆ Approximate inference algorithms
 - With the appropriate “dark art” they perform well

For the purposes of this tutorial we assume we can solve queries in networks.

Outline

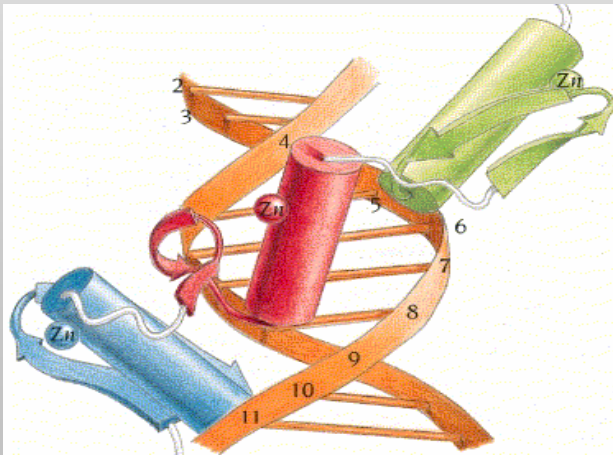
- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

Transcriptional Regulation

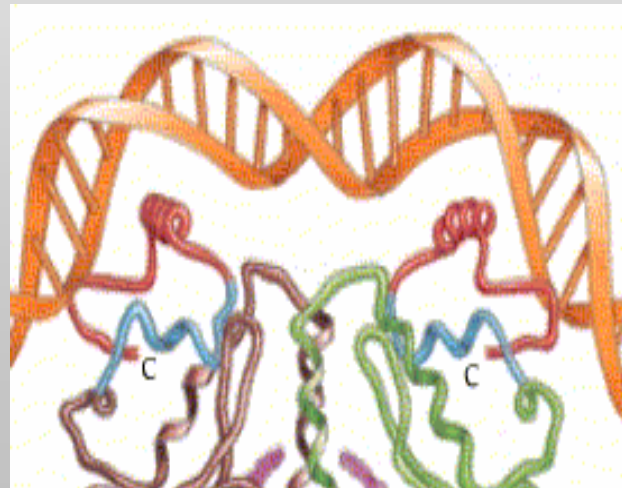


Transcription Factor Binding Sites

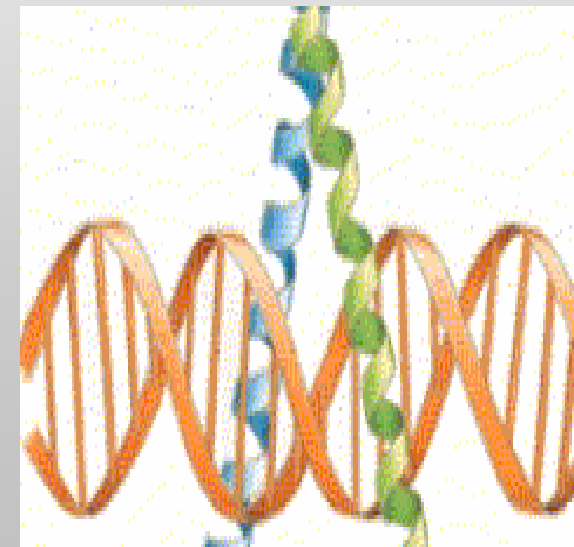
- ◆ Gene regulatory proteins contain structural elements that can “read” DNA sequence “motifs”
- ◆ The amino acid – DNA recognition is not straightforward
- ◆ Experiments can pinpoint binding sites on DNA



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)

N**N****G****G****G****G****C****N****G****G****G****C**

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



GCGGGGCCGGGC
 TGGGGGCGGGGT
 AGGGGGCGGGGG
 TAGGGGCGGGGC
 TGGGGGCGGGGT
 TGGGGGCGGGGC
 ATGGGGCGGGGC
 GTGGGGCGGGGC
 AAAGGGCCGGGC
 GGGAGGCCGGGA
 GGGGGGGGGGGG

Is this sufficient?

0 1 2 3 4 5 6 7 8 9 10 11 12

How to model binding sites ?

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



Profile: Independence model

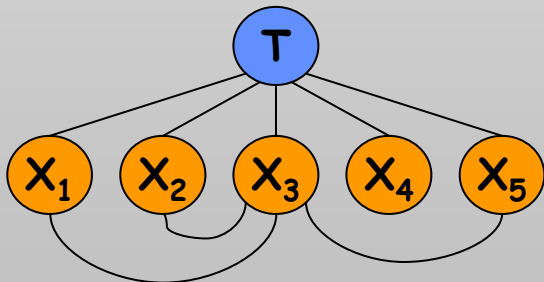
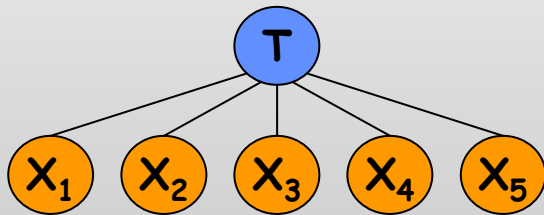
Tree: Direct dependencies

Mixture of Profiles:

Global dependencies

Mixture of Trees:

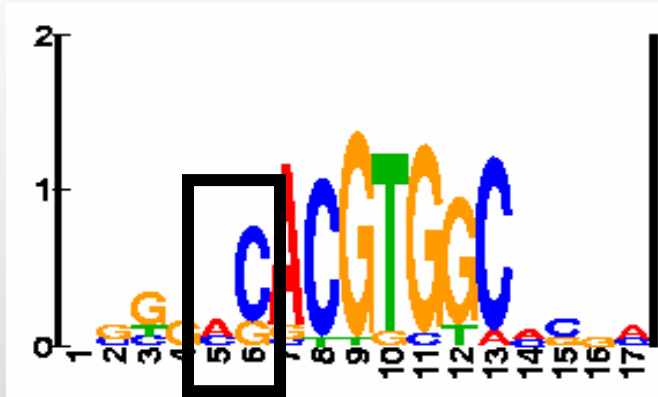
Both types of dependencies



$$P(X_1 \dots X_5) = \sum_T P(T) \prod_{i=1}^5 P(X_i | T) \prod_{(i,j) \in \text{edges}} P(X_i, X_j | T)$$

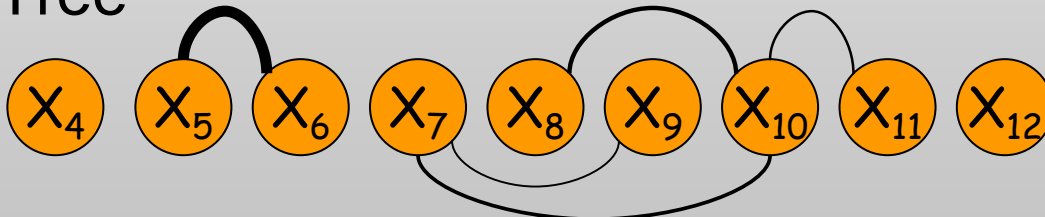
Arabidopsis ABA binding factor 1

Profile



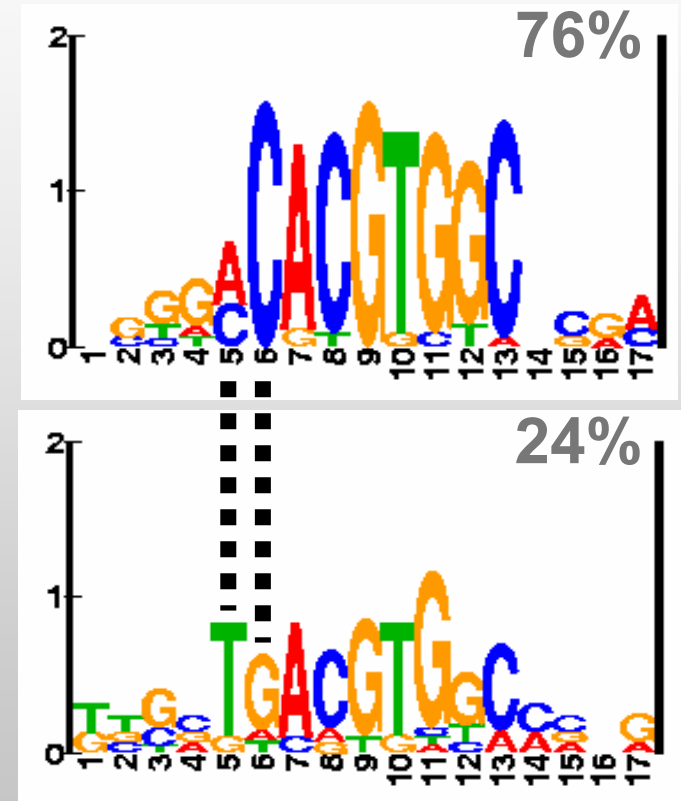
Test LL per instance **-19.93**

Tree



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

Mixture of Profiles



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

The Knowledge Acquisition Bottleneck

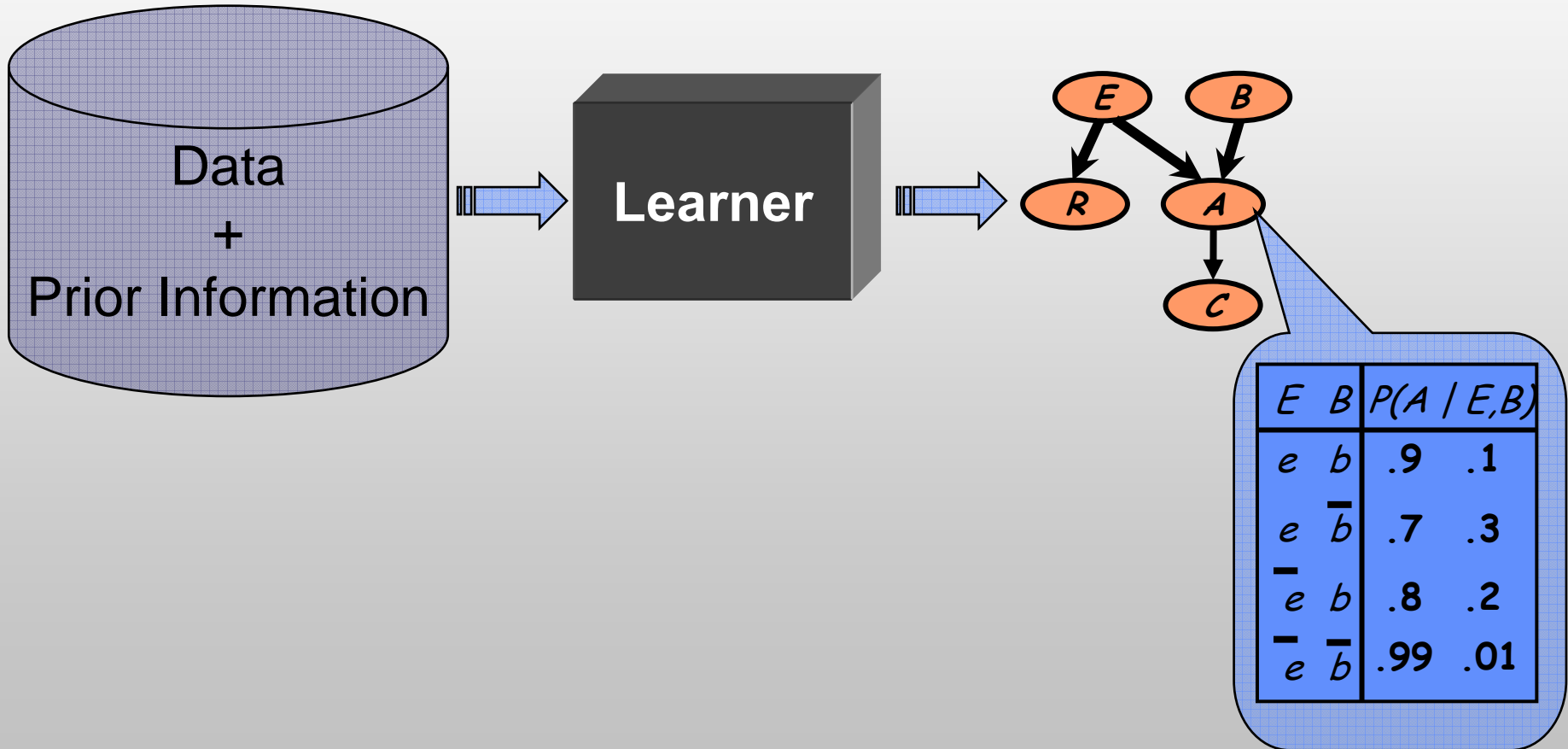
How do we construct these models?

- ◆ Knowledge acquisition is an expensive process
- ◆ Often we don't have an expert

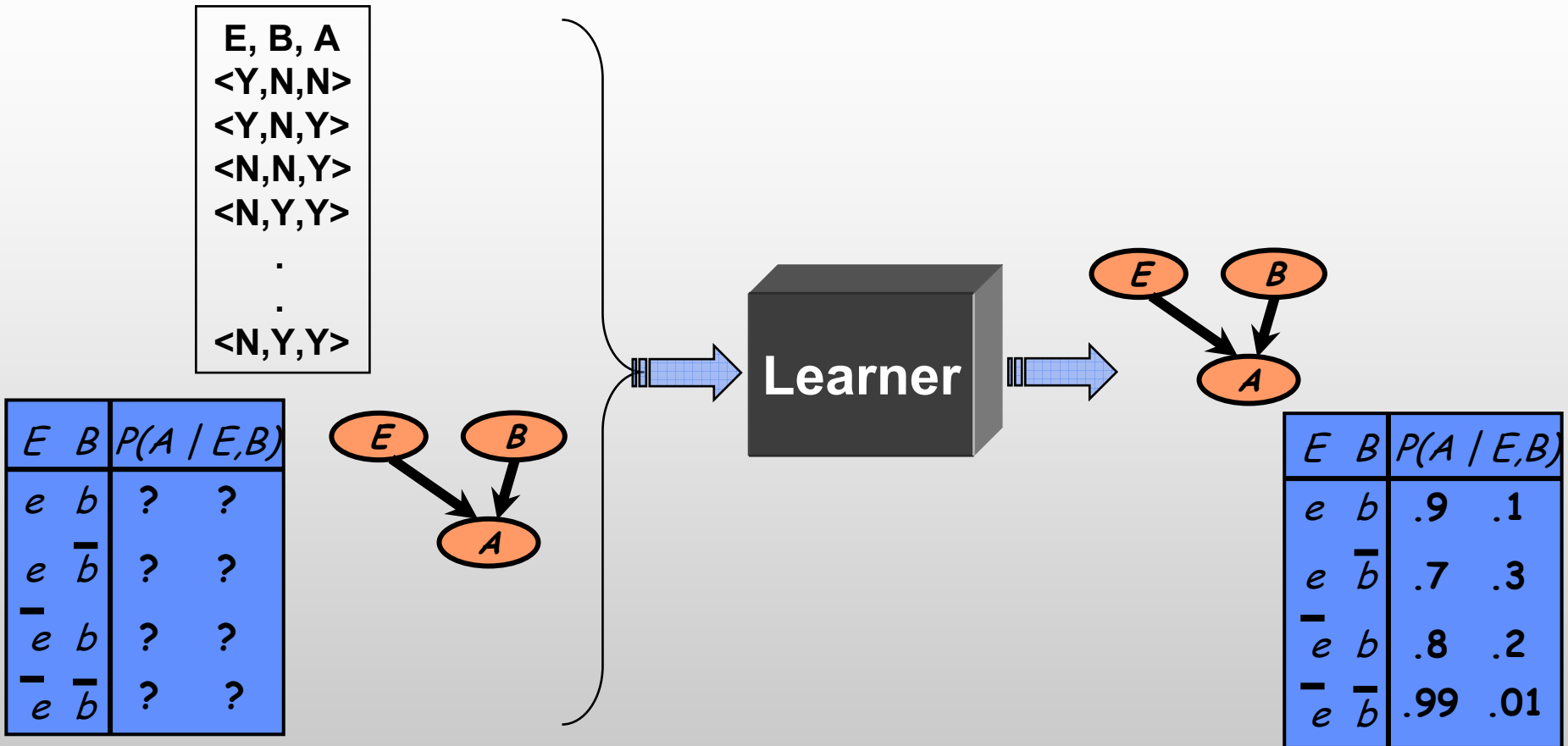
Harnessing Data

- ◆ Amount of available information growing rapidly
- ◆ Learning allows us to construct models from raw data
- ◆ The the details of learned models provide insights about the data

Learning Bayesian networks

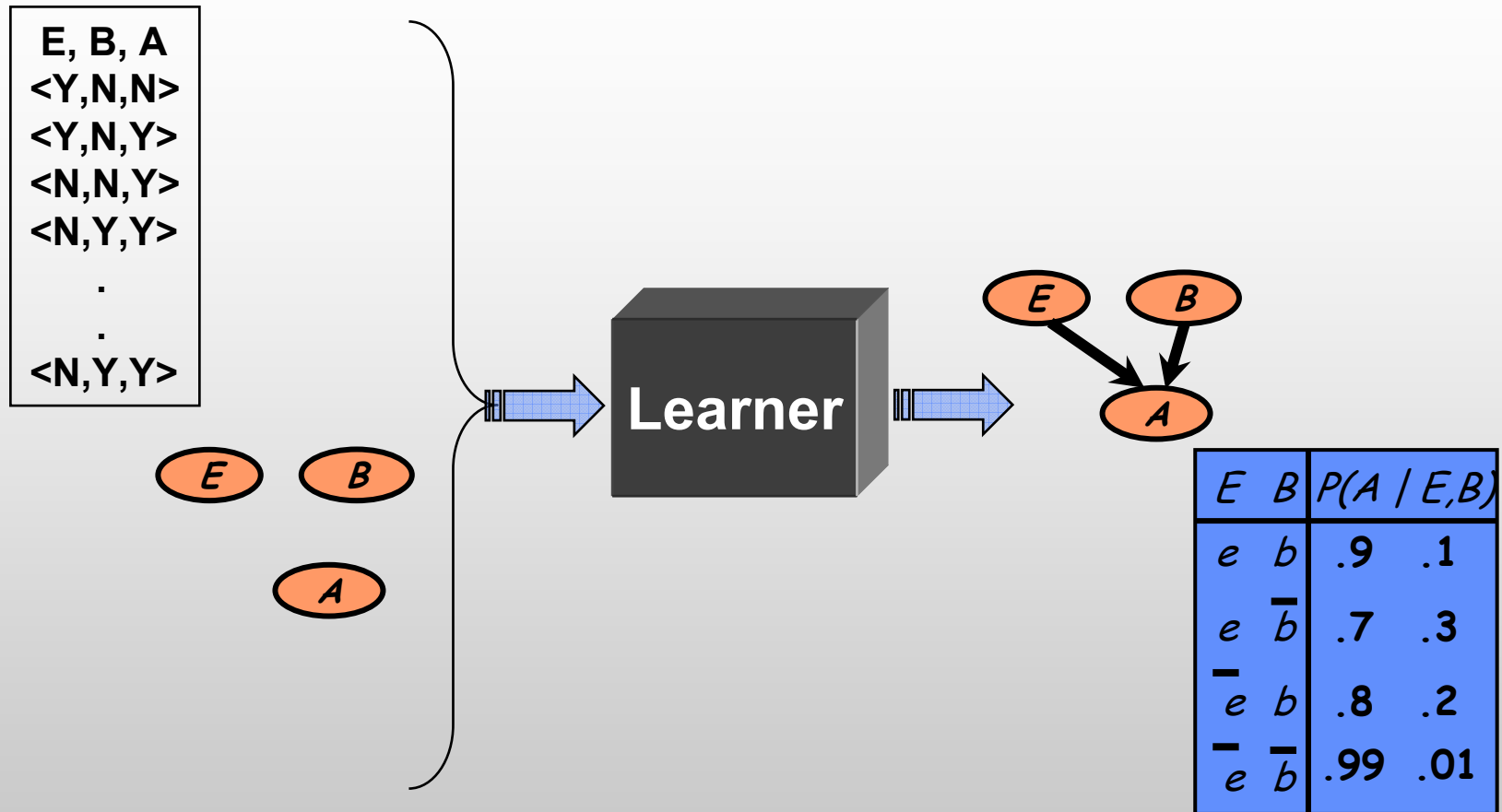


Known Structure, Complete Data



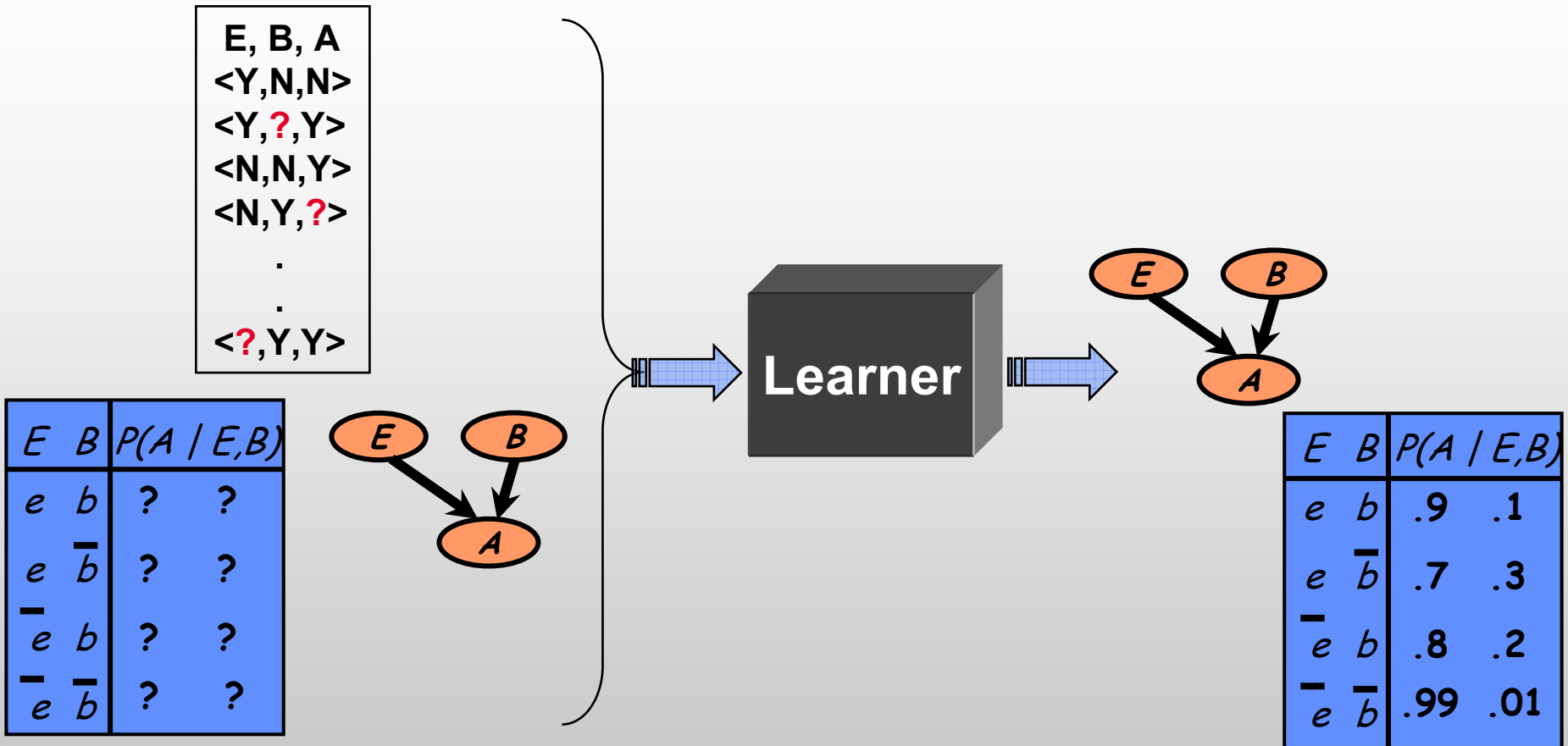
- ◆ Network structure is specified
 - Learner needs to estimate parameters
- ◆ Data does not contain missing values

Unknown Structure, Complete Data



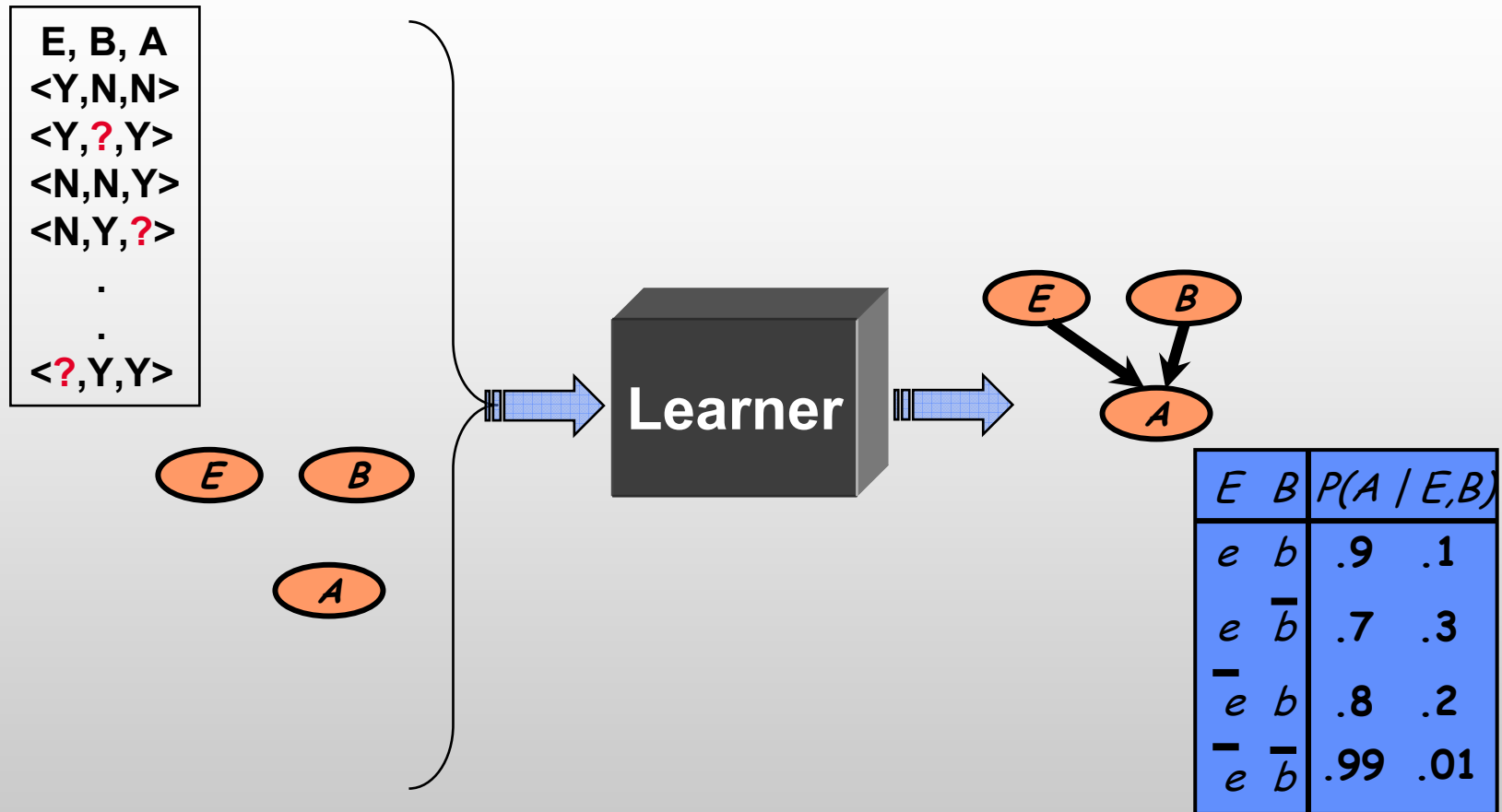
- ◆ Network structure is not specified
 - Inducer needs to select arcs & estimate parameters
- ◆ Data does not contain missing values

Known Structure, Incomplete Data



- ◆ Network structure is specified
- ◆ Data contains missing values
 - Need to consider assignments to missing values

Unknown Structure, Incomplete Data



- ◆ Network structure is not specified
- ◆ Data contains missing values
 - Need to consider assignments to missing values

The Learning Problem

	Known Structure	Unknown Structure
Complete Data	Statistical parametric estimation (closed-form eq.)	Discrete optimization over structures (discrete search)
Incomplete Data	Parametric optimization (EM, gradient descent...)	Combined (Structural EM, mixture models...)

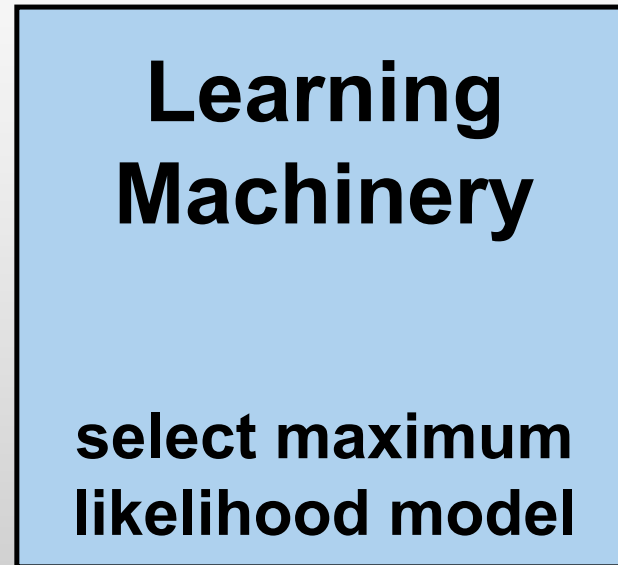
Outline

- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ **Transcriptional regulation**
- ◆ Gene expression
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

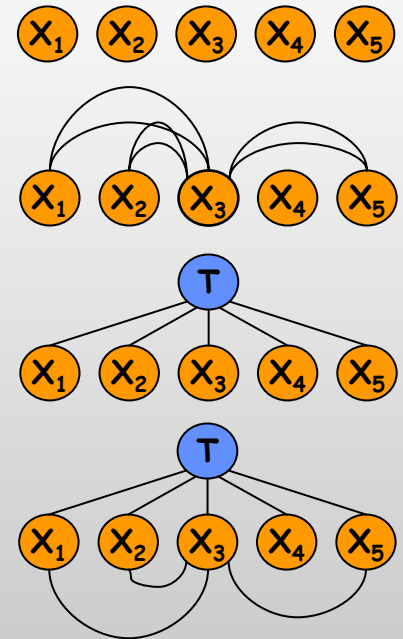
Learning models: Aligned binding sites

Aligned binding sites

```
GCGGGGCCGGGC  
TGGGGGCGGGGT  
AGGGGGCGGGGG  
TAGGGGCCGGGC  
TGGGGGCGGGGT  
AAAGGGCCGGGC  
GGGAGGCCGGGA  
GCGGGGCGGGGC  
GAGGGGACGAGT  
CCGGGGCGGTCC  
ATGGGGCGGGGC
```



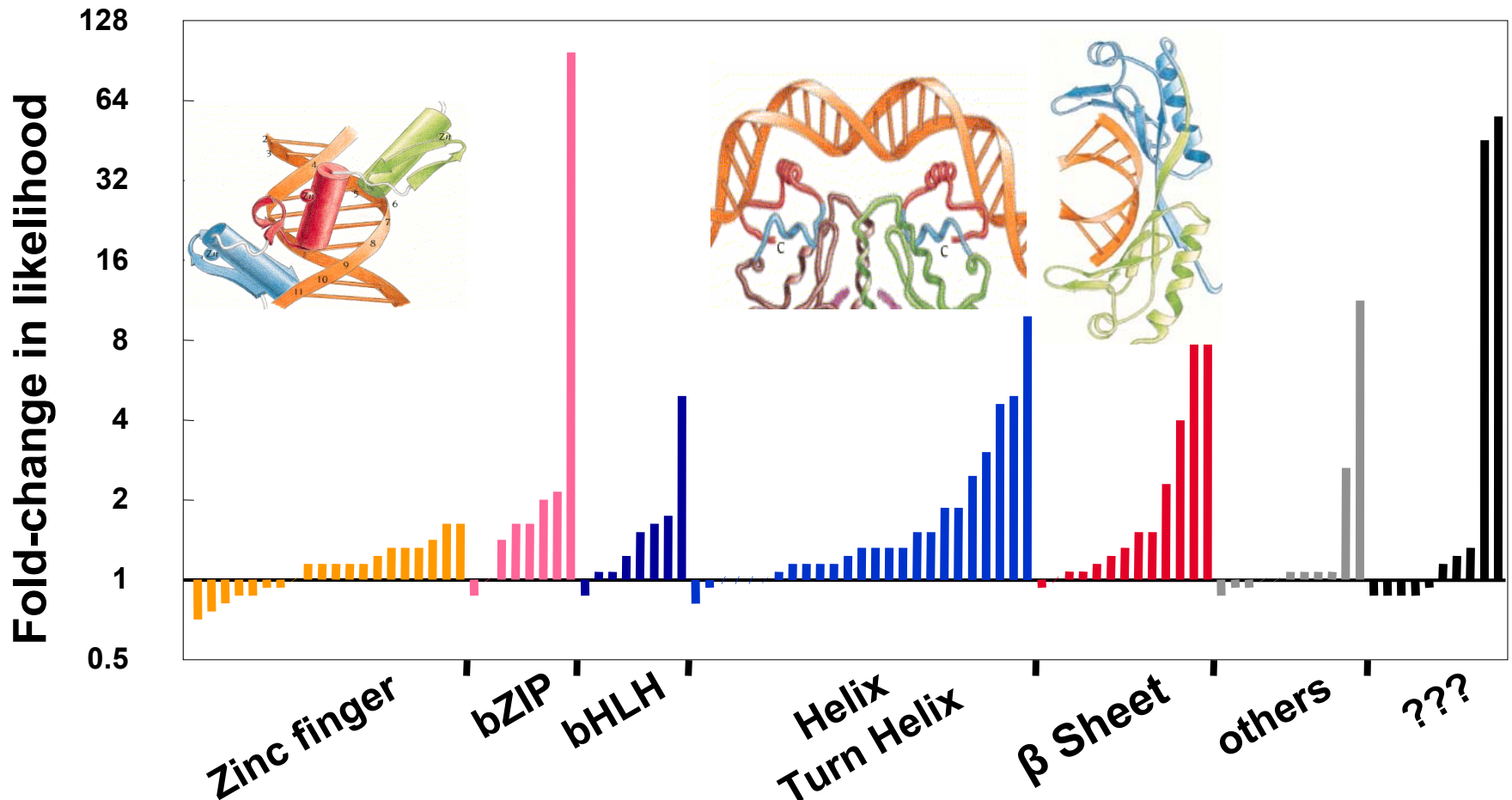
Models



Learning based on methods for probabilistic graphical models (*Bayesian networks*)

Likelihood improvement over profiles

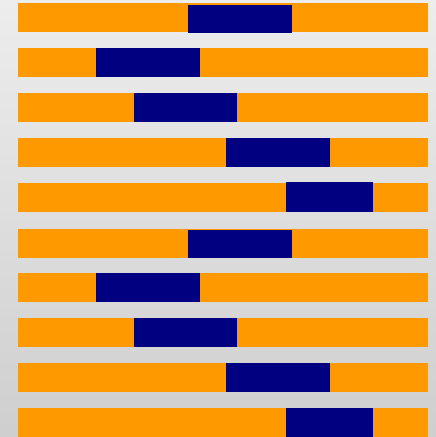
TRANSFAC: 95 aligned data sets



Motif finding problem

Input: A set of potentially co-regulated genes

Output: A common motif in their promoters

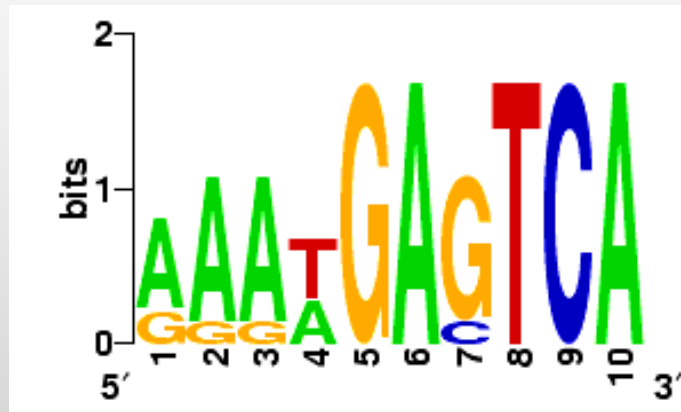


Sources of data:

- ◆ 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)

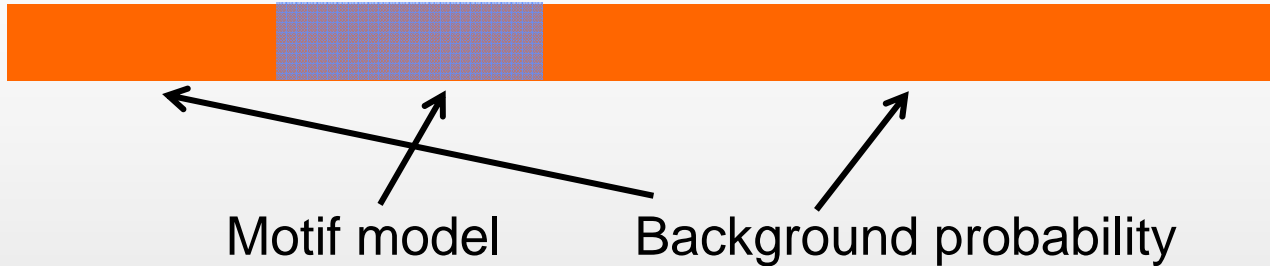
Example

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



5' - TCTCTCTCCACGGCTAATTAGGTGATCATGAAAAATGAAAAATTCATGAG**AAAAGAGT**CAGACATCGAAACATACAT ...*HIS7*
 5' - ATGGCAGAATCACTTTAAACGTGGCCCCACCCGCTGCACCCTGTGCATTTTGTACGTTACTGCG**AAATGACT**CAACG ...*ARO4*
 5' - CACATCCAACGAATCACCTCACCGTTATCG**TGACTCACTT**TCTTTTCGCATCGCCGAAGTGCCATAAAAAATATTTTTT ...*ILV6*
 5' - TGCGAAC**AAAAGAGT**CA**T**TACAACGAGGAAATAGAAGAAAATGAAAAATTTTCGACAAAATGTATAGTCATTTCTATC ...*THR4*
 5' - ACAAAGGTACCTTCCTGGCCAATCTCACAGATTTAATATAGTAAATTGTCATGCATAT**TGACTCAT**C**CC**GAACATGAAA ...*ARO1*
 5' - ATTGAT**TGACTCAT**TTTCTCTGACTACTACCAGTTCAAAATGTTAGAGAAAAATAGAAAAGCAGAAAAAATAAATAA ...*HOM2*
 5' - GCGCCACAGTCCGCGTTTGGTTATCCGGC**TGACTCATTCTGACTCTTTT**TTGGAAAGTGTGGCATGTGCTTCACACA ...*PRO3*

Probabilistic Model



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

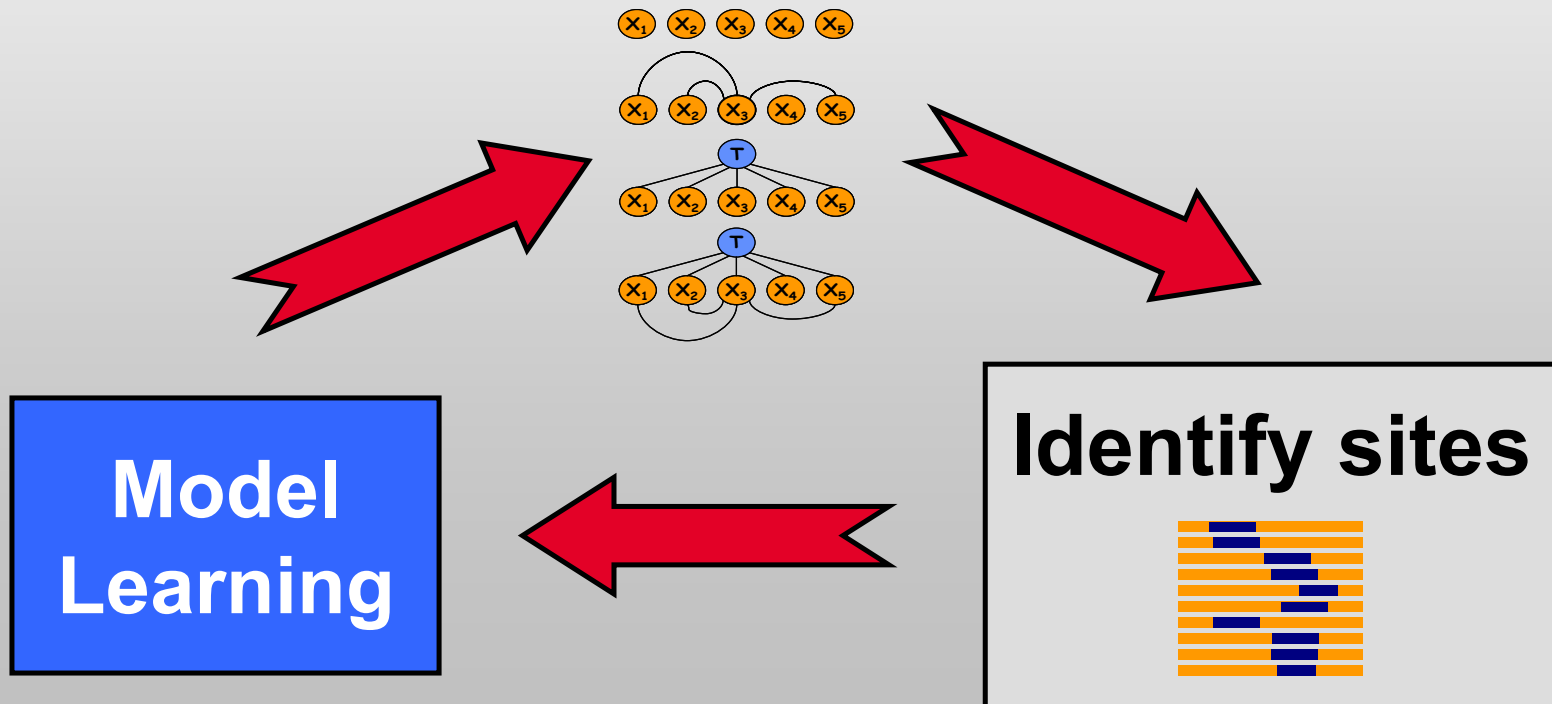
Hidden variable:

- ◆ Location of motif within each sequence

Learning models: unaligned data

EM (MEME-like)

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



ChIP location analysis

Yeast genome-wide location analysis

Target genes annotation for 106 TFs

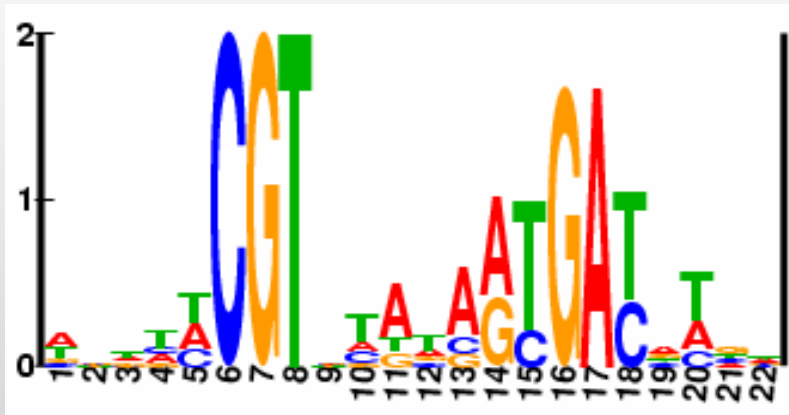
	Gene	ABF1 Targets	ZAP1 Targets
↑ # genes ~ 6000 ↓	YAL001C	+		-
	YAL002W	-		+
	YAL003W	+		-
	YAL005C	-		-
	.	.		.
	.	.		.
	.	.		.
	YAL010C	+		-
	YAL012C	-		+
	YAL013W	-		+
YPR201W	-		-	

Example: Models learned for ABF1 (YPD)

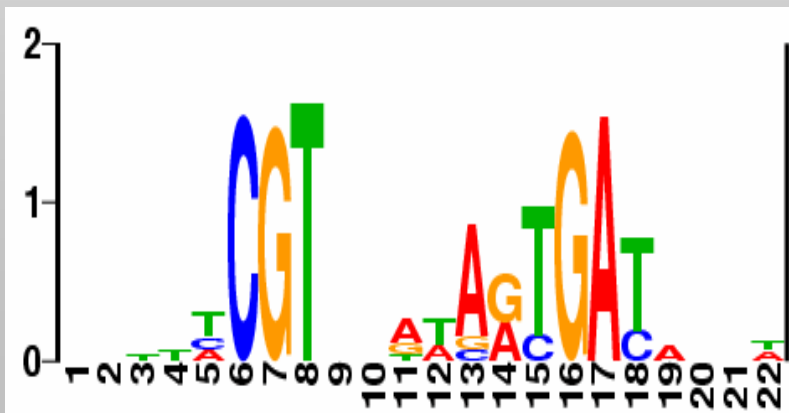
Autonomously replicating sequence-binding factor 1

Known profile

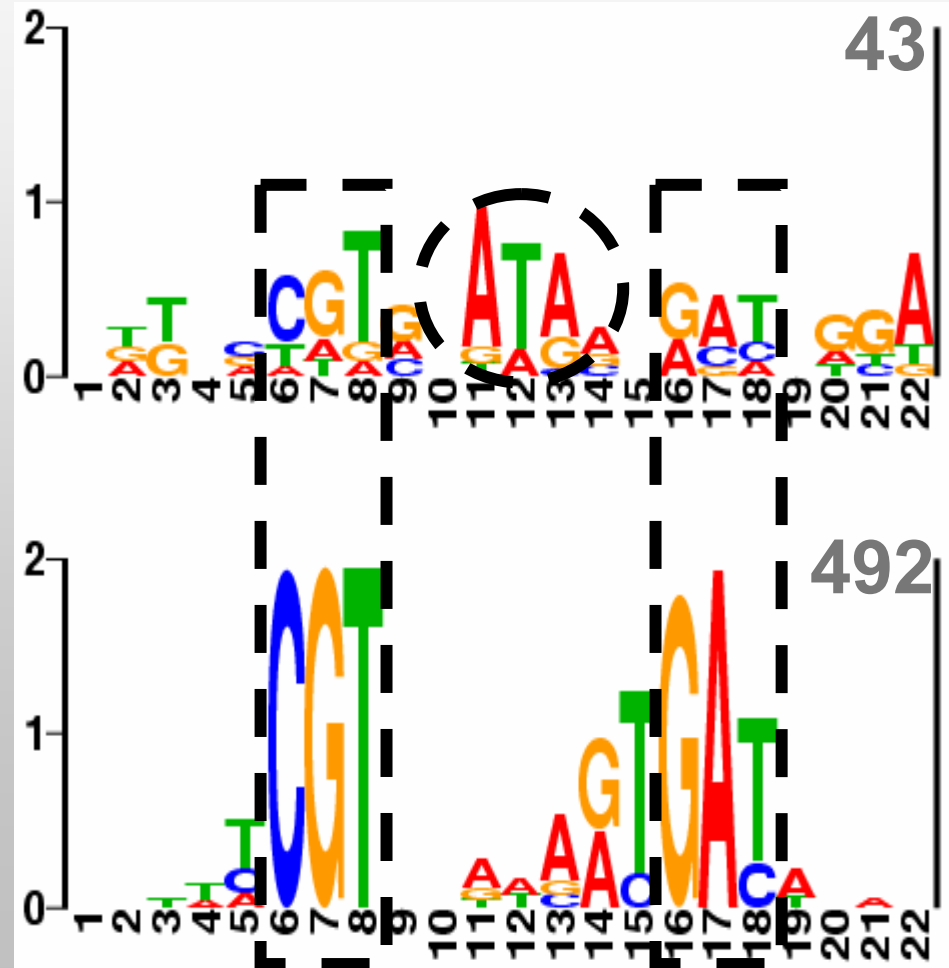
(from TRANSFAC)



Learned profile



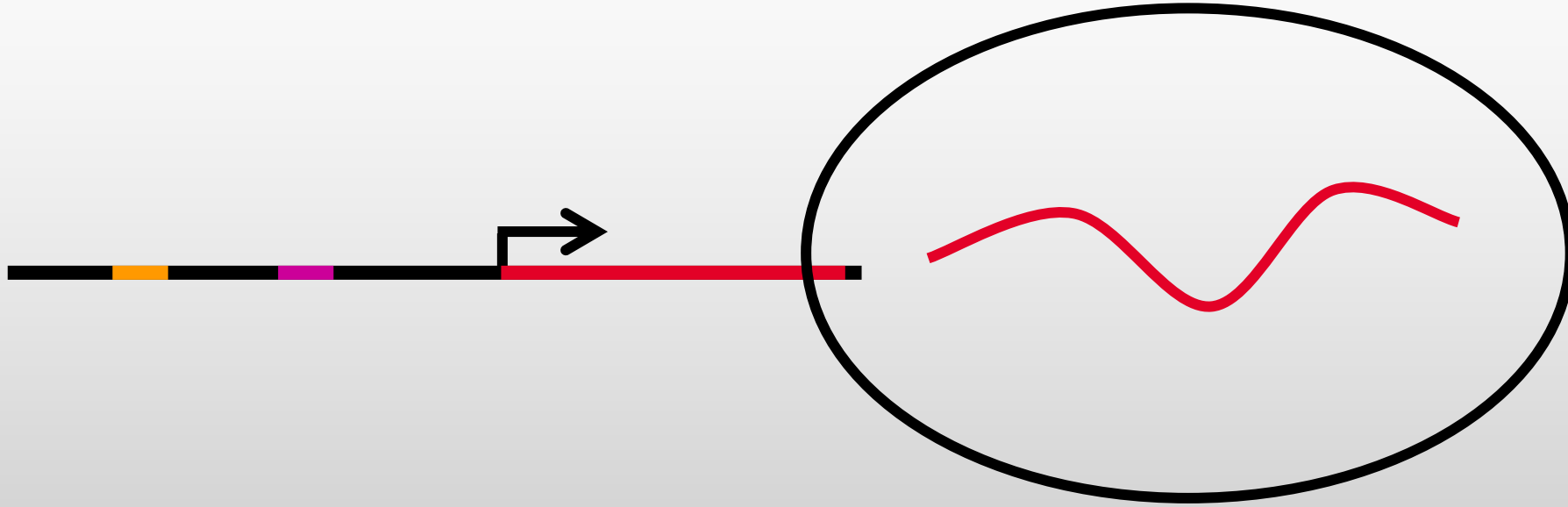
Learned Mixture of Profiles



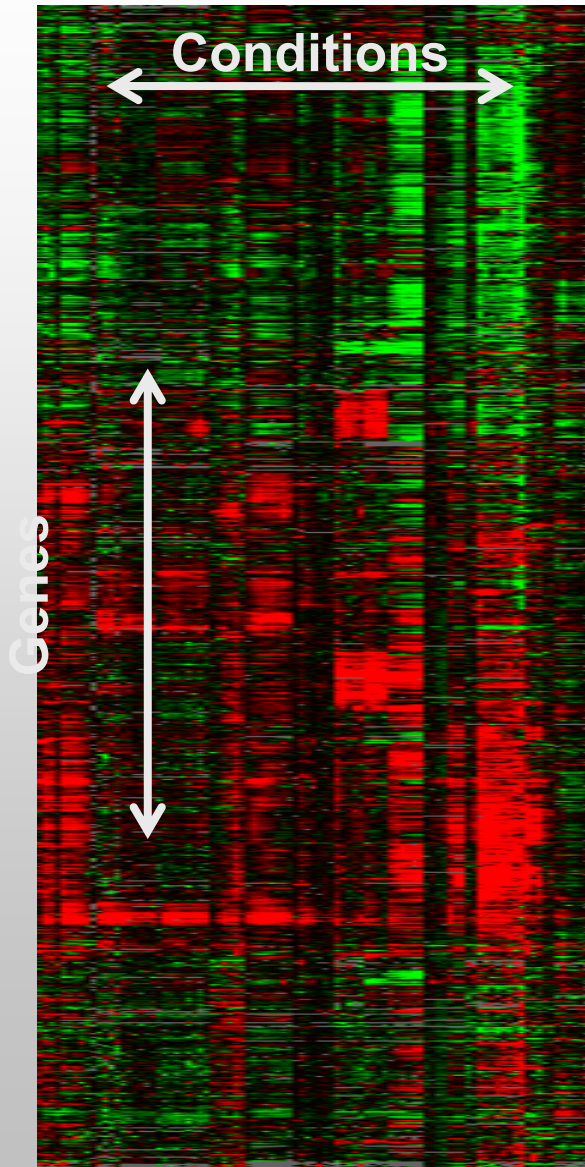
Outline

- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ **Gene expression**
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

Transcriptional Regulation



Expression Data



- ◆ 1000s of genes

- ◆ 10-100s of arrays

Possible designs

- ◆ Biopsies from different patient populations

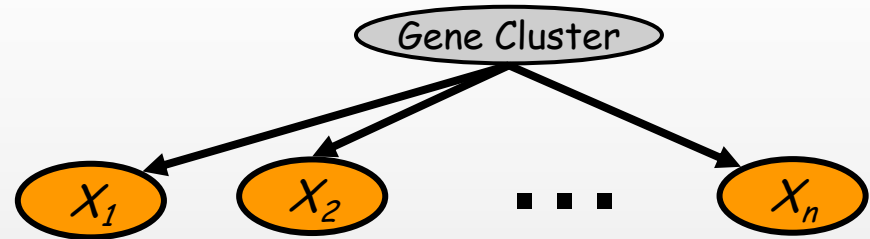
- ◆ Time course

- ◆ Different perturbations

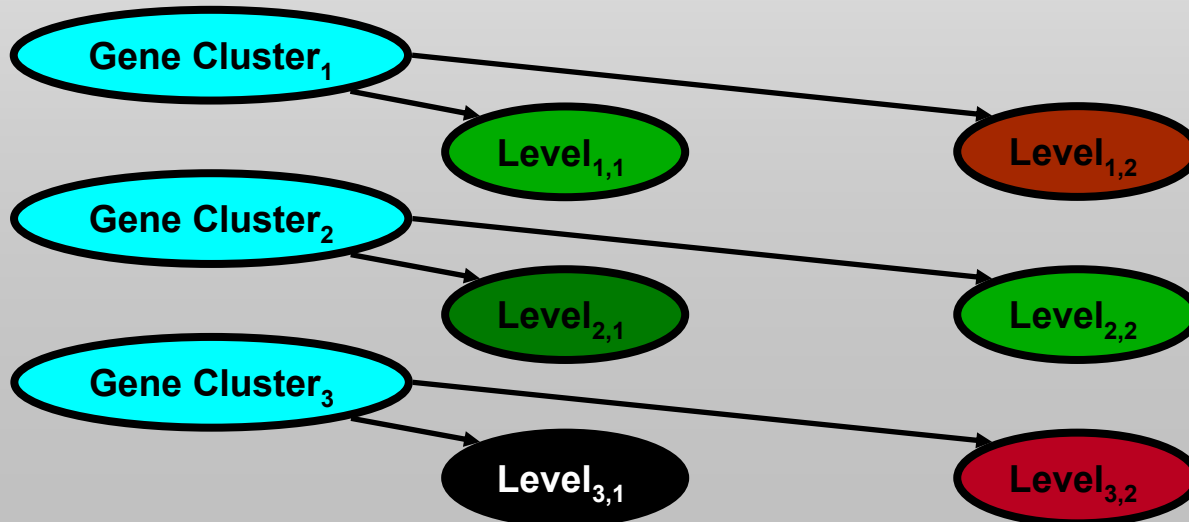
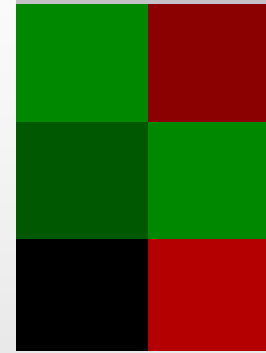
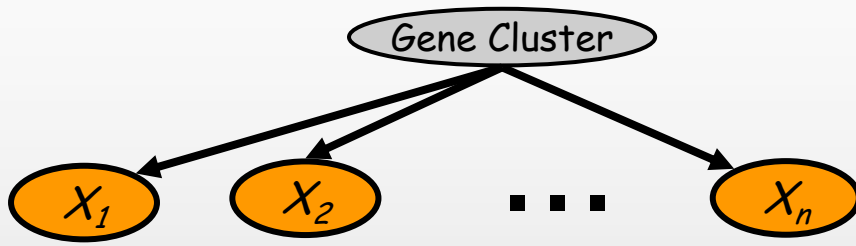
Clustering Gene Expression Profiles

Clustering model

- ◆ “Cluster” hidden variable explains dependencies among measurement of a gene in different conditions
- ◆ Each gene is viewed as a sample from the same distribution



Clustering Genes



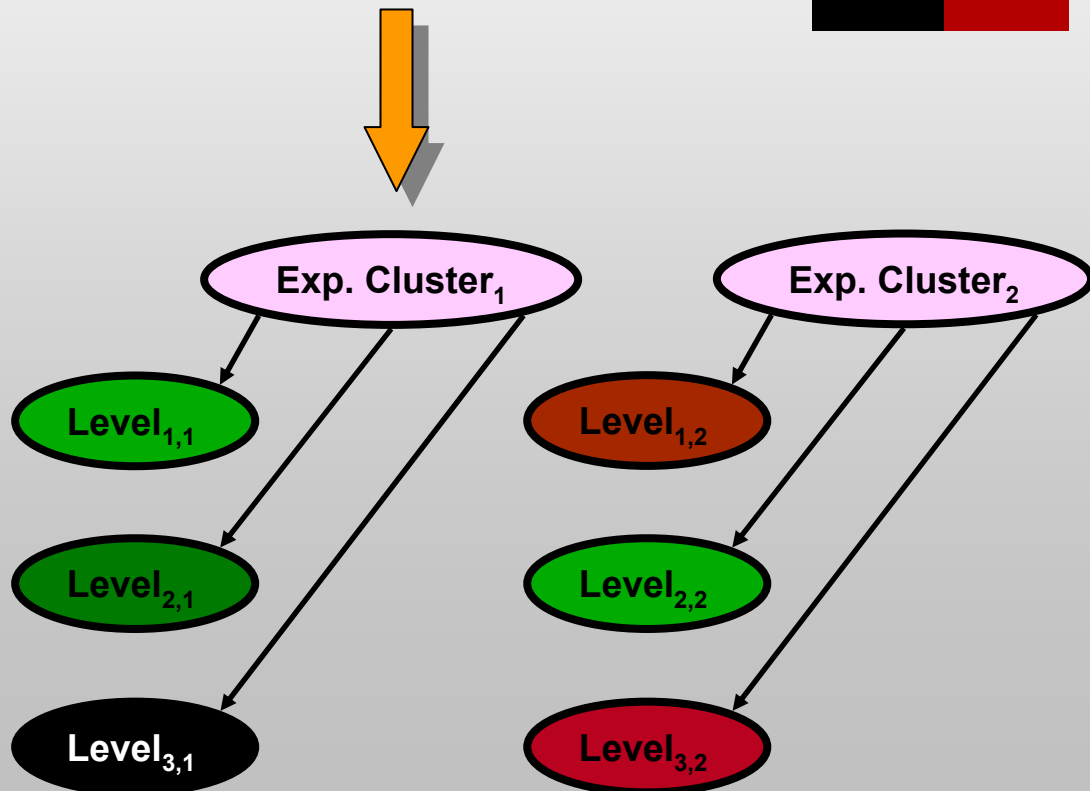
Parameters are shared by network copies

Clustering Conditions

Can we cluster both genes and conditions?

Now each condition is an instance

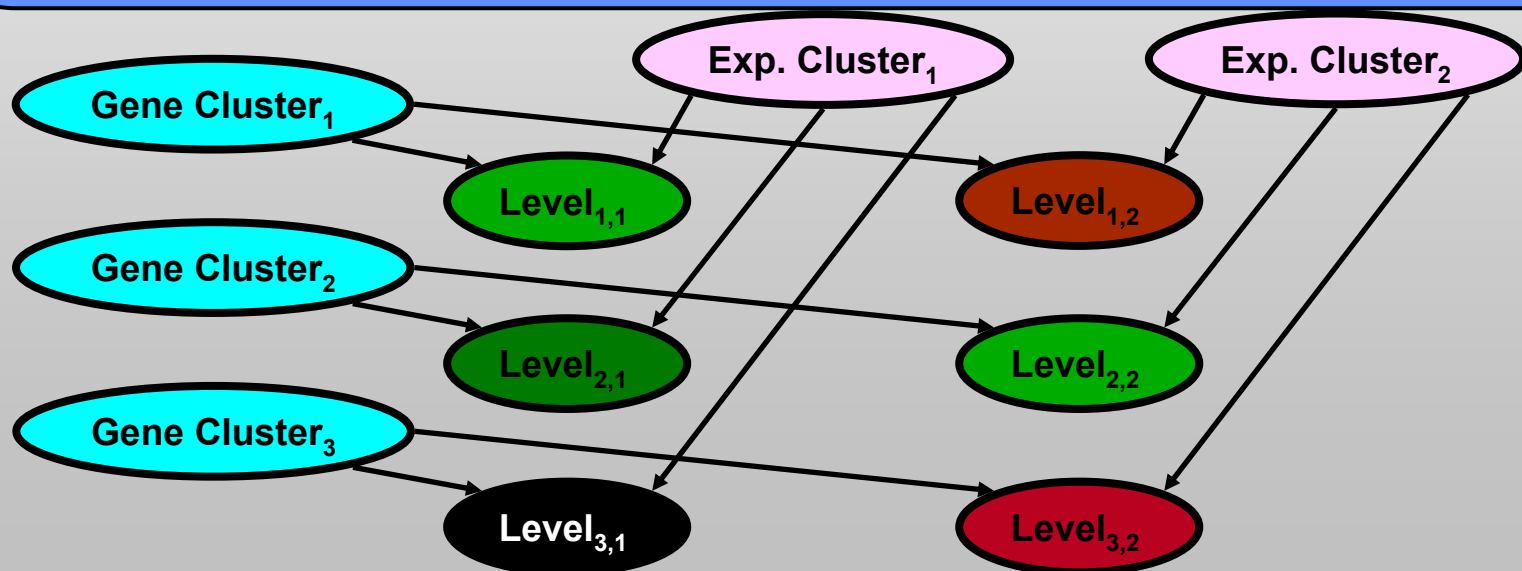
Parameters are shared by network copies



Joint Clustering?

A single network that spans the whole data

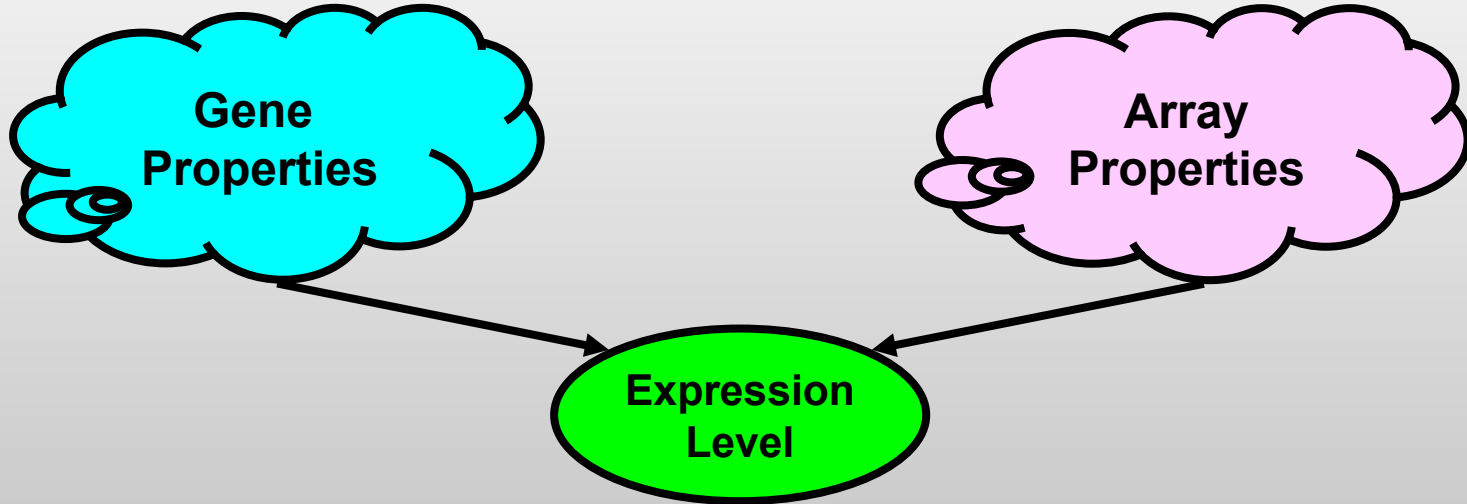
- Each expression variable has its own parameters
- # parameters \gg # observations



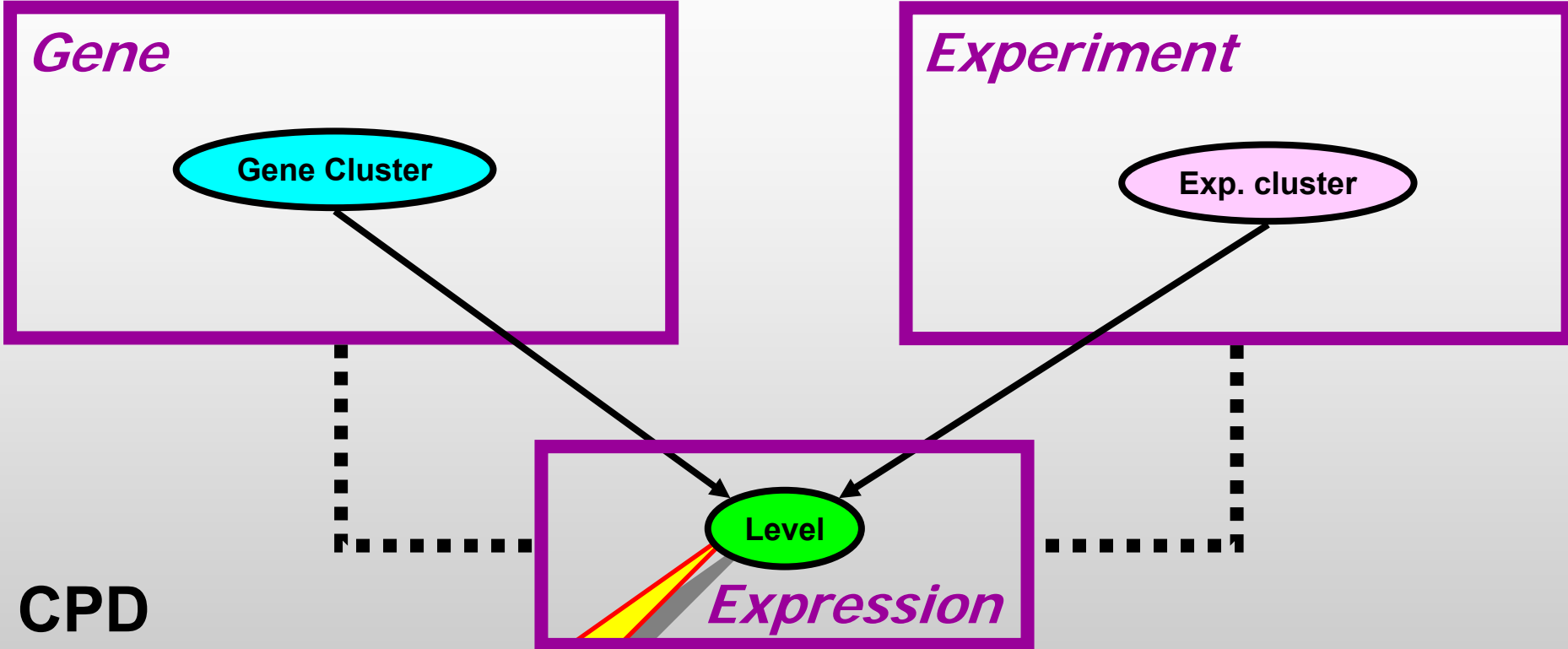
Relational Approach

Key Idea:

Expression level is “explained” by properties of gene and properties of experiment

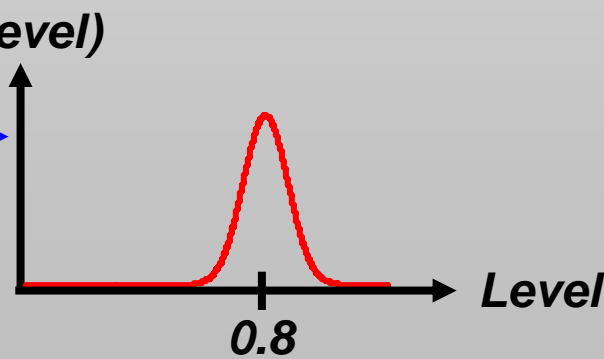
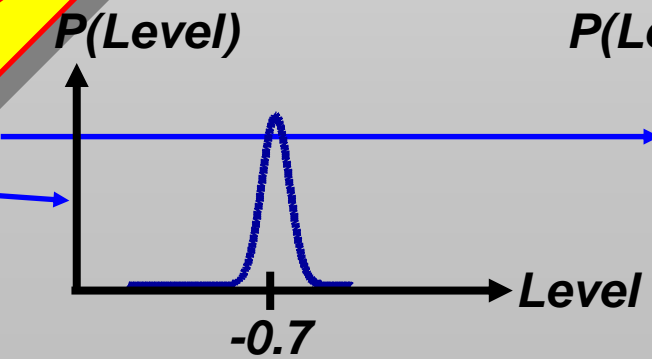


Probabilistic Relational Models

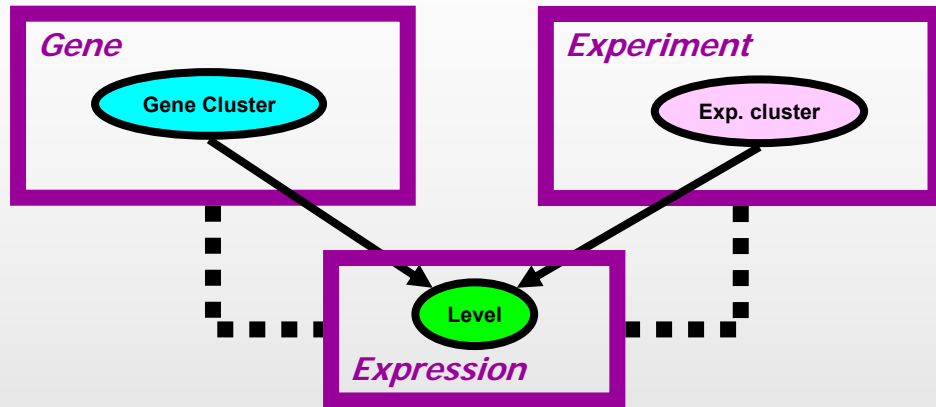


CPD

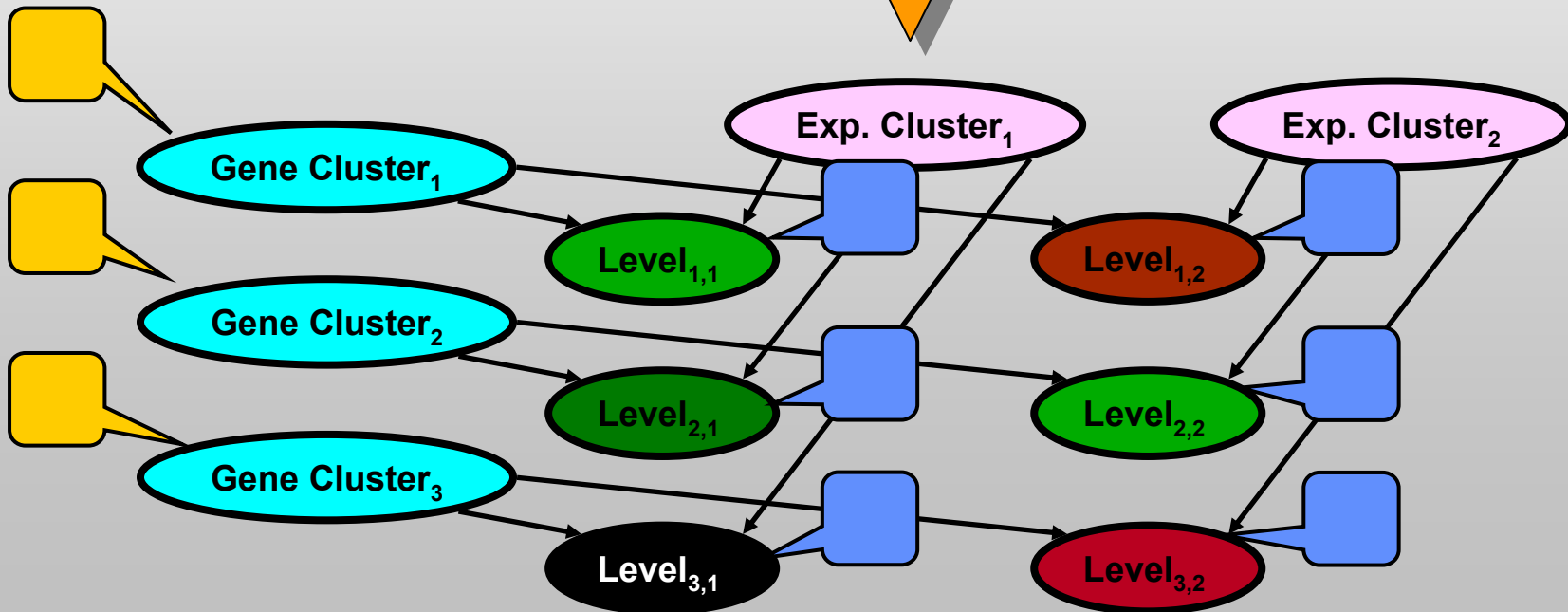
$G_{Cluster}$	$E_{Cluster}$	μ	σ
1	1	0.8	1.2
1	2	-0.7	0.6
...



Unrolling a Relational Network



Parameters of variables from the same template are shared



Expression + Annotations

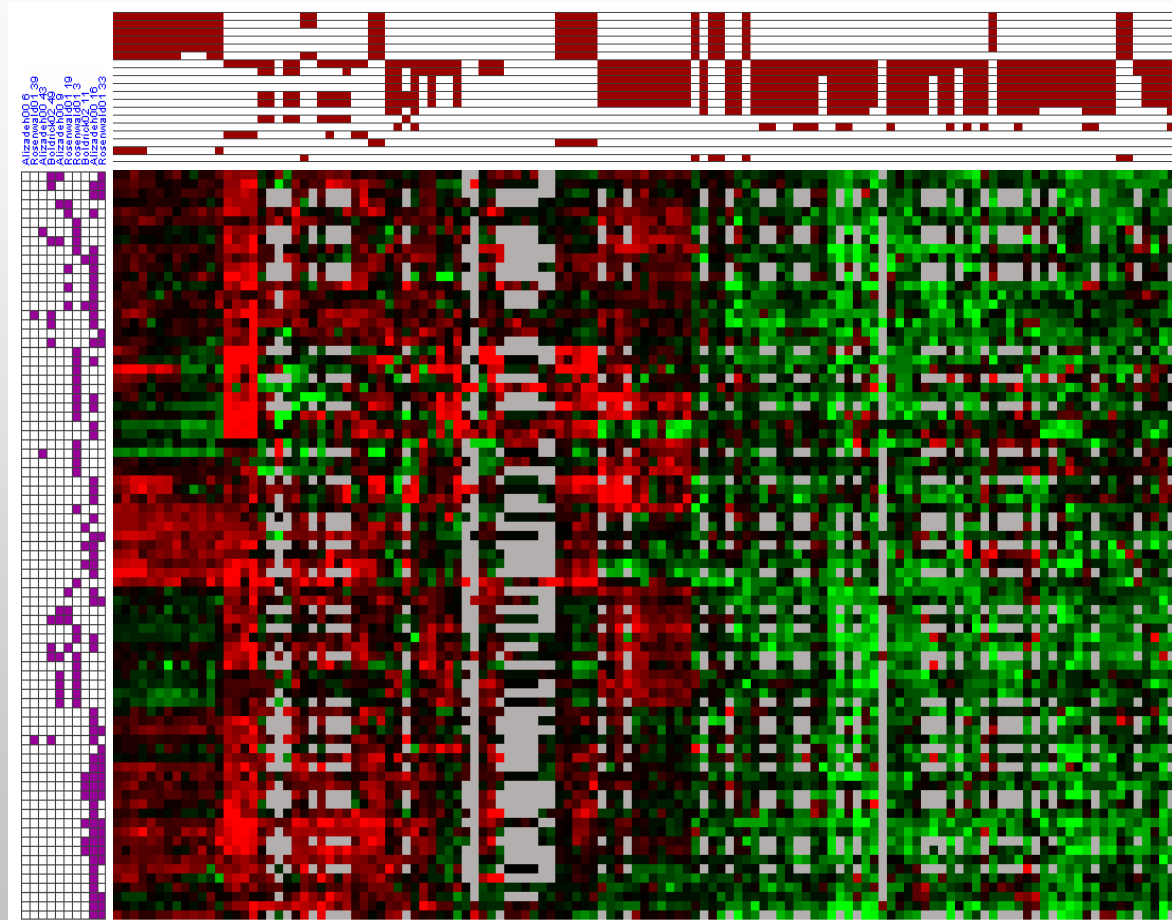
Array annotations:

Tissue type, Clinical conditions, Treatments, Prognosis

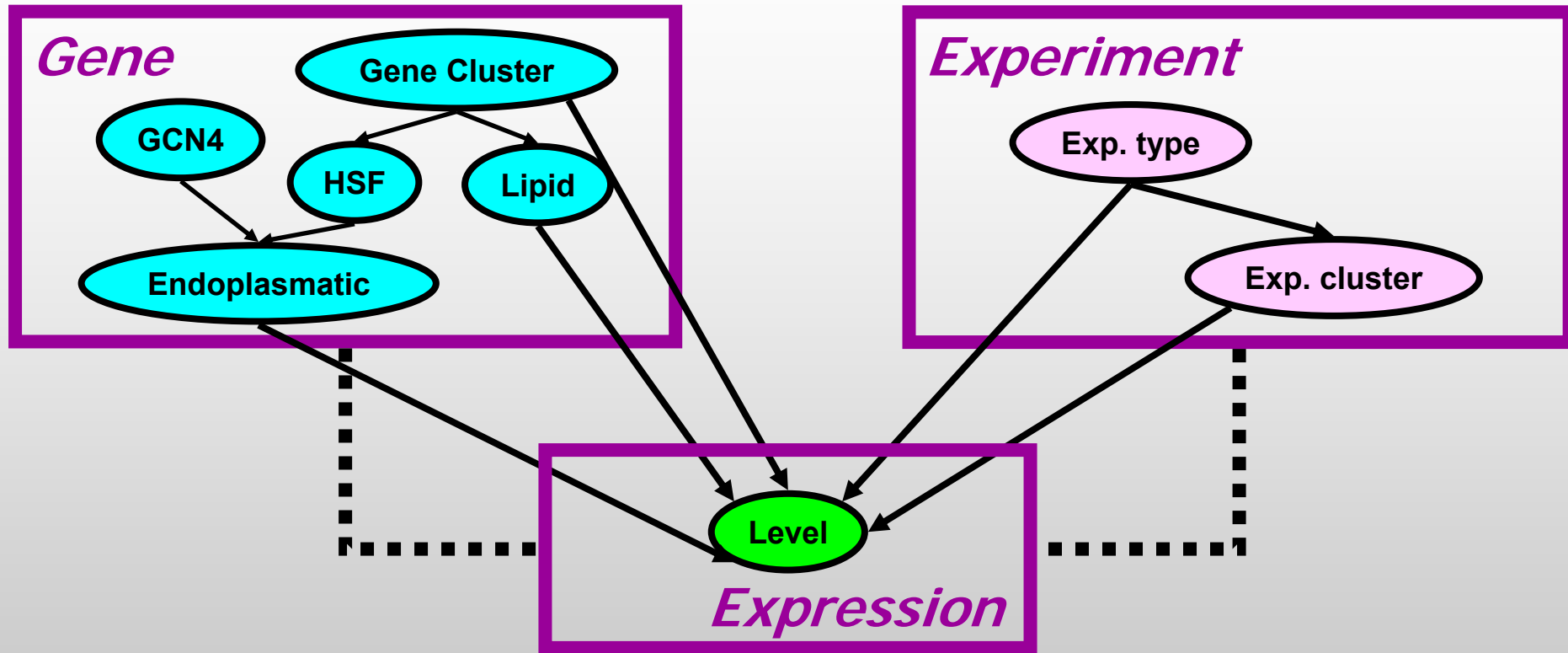
Gene annotations:

Function, Process, Regulatory regions, Cellular location, protein family

Relational models!

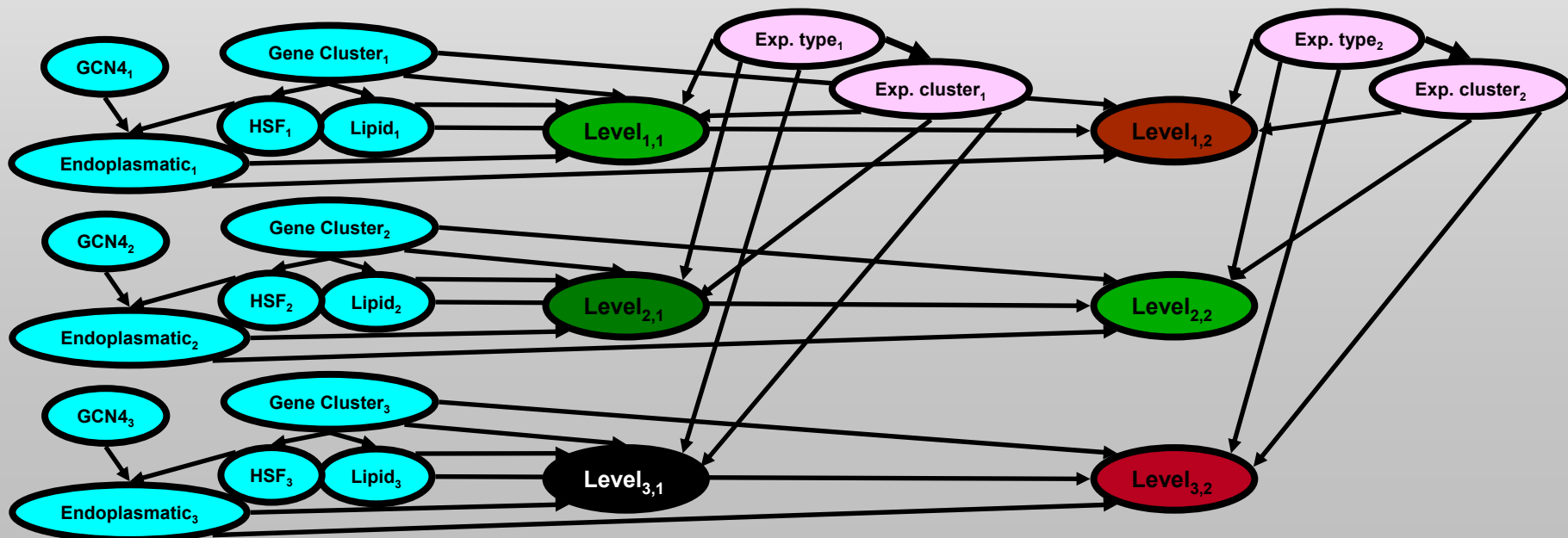
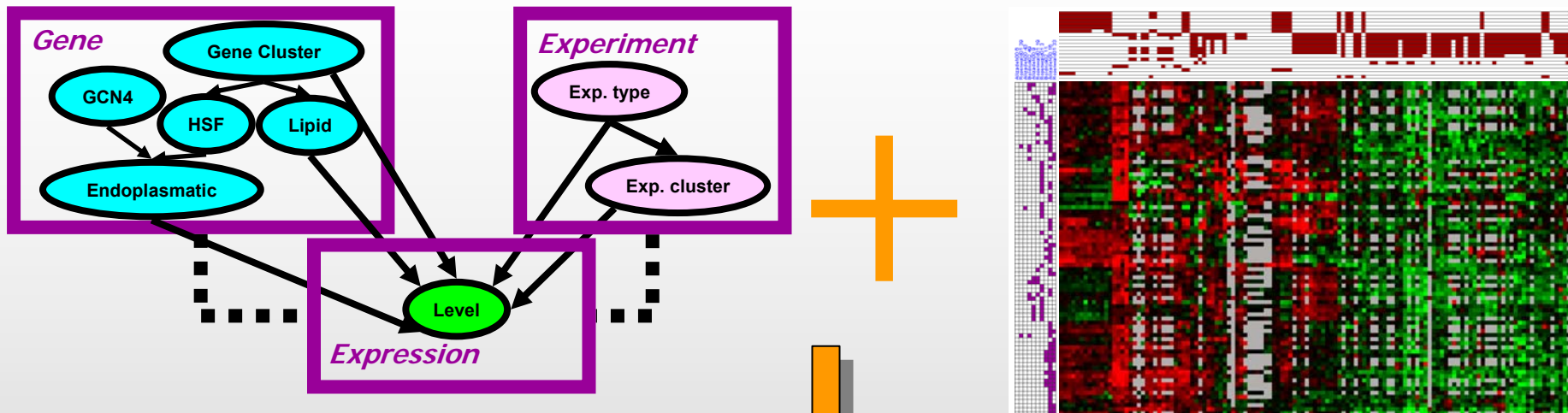


Adding Additional Data



- ◆ Annotations
- ◆ Binding sites
- ◆ Experimental details

Semantics



TF to Expression

Key Question:

- ◆ Can we explain changes in expression?

General model:

- ◆ Transcription factor binding sites in promoter region should “explain” changes in transcription



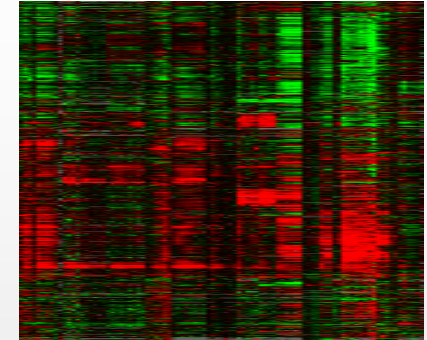
Goal

ACTAGTGCTGA

CTATTATTGCA

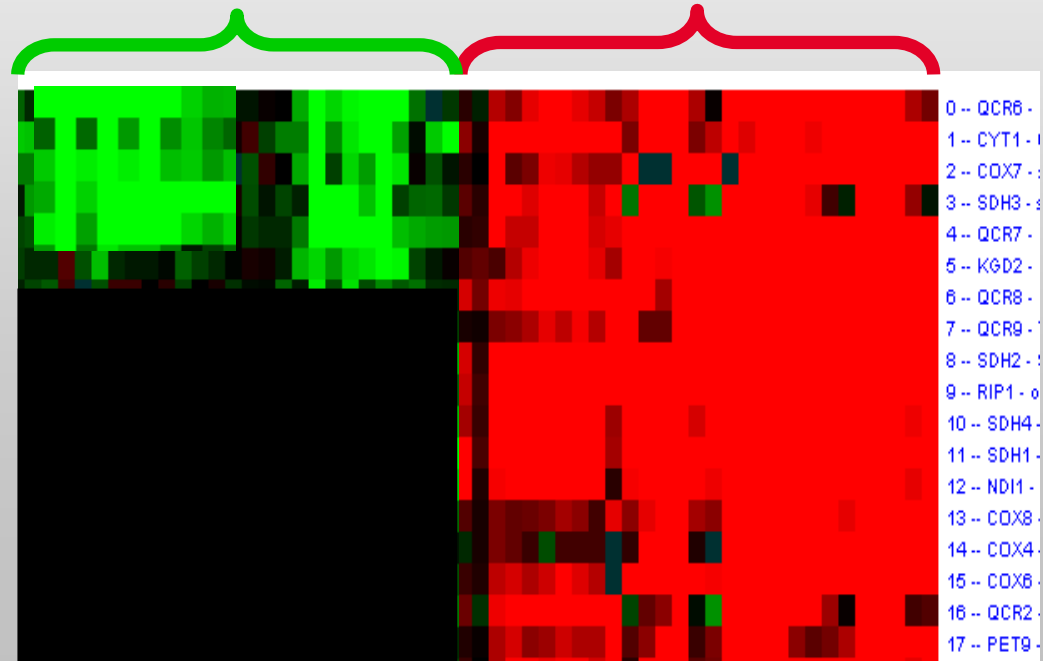
CTGATGCTAGC

+



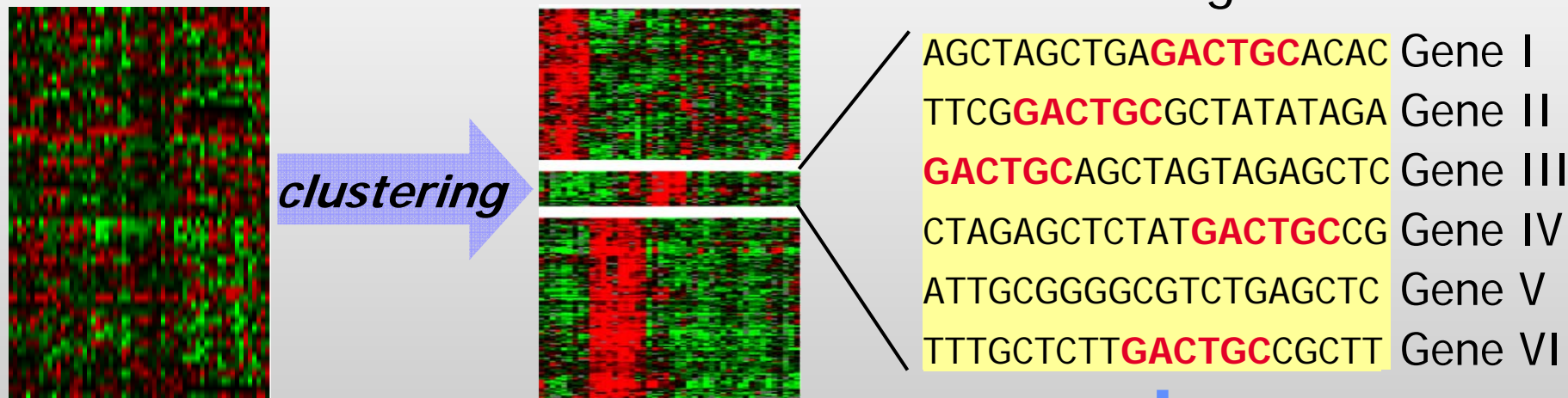
t_2 Motif t_1 Motif

$R(t_2)$ {
 ACCTAGCTTAGACTGCACTGATCGAG
 CCCGACATAGCTTGGACTGCGCTATA
 TAGACTGCAGCTAGTAGAGCTCTGCTAG
 AGCTCTATGACTGCCGATTGCGGGCGT
 CTGAGCTCTTTGCTCTTGACTGCCGCTT
 TTGATATTATCTCTCTGCTCGTGACTGC
 TTTATTGTGGGGGACTGCTGATTATGC
 TGCTCATAGGAGAGACTGCGAGAGTCGT
 CGTAGACTGCGTTCGTCTGATGATGCT
 GCTGATCGATCGACTGCCTAGCTAGTA
 GATCGATGTGACTGCAGAAGAGAGAGGG
 TTTTTTCGCGCCCGCCCGCGGACTGCT
 CGAGAGGAAGTATATATGACTGCGCGCG
 CCGCGCGCACGGACTGCAGCTGATGCAT
 GCATGCTAGTAGACTGCCTAGTCAGCTG
 CGATCGACTCGTAGCATGCATCGACTGC
 AGTCGATCGATGCTAGTTATTGACTGC
 GTAGTAGTGCGACTGCTCGTAGCTGTAG



“Classical” Approach

- ◆ Cluster gene expression profiles
- ◆ Search for motifs in control regions of clustered genes



Procedural

- Apply separate method to each type of data
- Use output of one method as input to the next
- Unidirectional information flow

Motif →

GACTGC

Flow of Information

Sequence

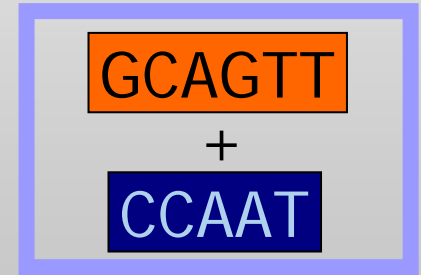
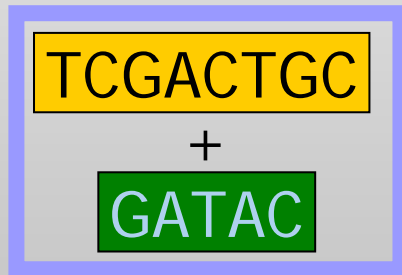
```
ACGATGCTAGTGTAGCTGATGCTGATCGATCGTACGTGCTAGCTAGCTAGCTAGCCCAATAGCTAGC
AGCTAGCTCGACTGGATACGGGTCGACTGCTCAAACACACACAACACCAAATGCCAATGTGG
TACTGATGATCGTAGTAACCACTGTCGATGATGCTGTGGGGGTATCGATGCCAATCACCCCCGCT
CGATCGATCGACTGCAGCTAGCTAGCTGATCAAAAACACCATACGCCCCCAATCTGCTCGTAGCAT
GCTAGCTAGCTGATCGATCAGCTACTCGACTGCGGATACGCTAGCTACTTTTTTTTTTTTGTAGCA
CCCAACTGACTGATCGTAGTCAGTACGTACGATCGTGACTGATCCAATGTCGTCGAGCAGTTTACG
TATCGACTGCGCATGCTAGCTGCTCGAAAAAAAACGTGCAAGTTCCAATGCTCGCCCCCCCCC
CCCGACTGATCGATACGACTAGTCGACTGCTGATCGATCGTAGCTCCAATATATATAGCAGTTACGGCG
```

Genes

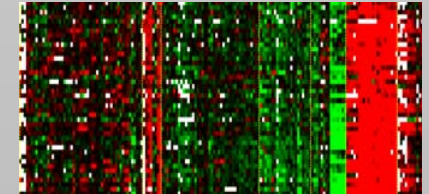
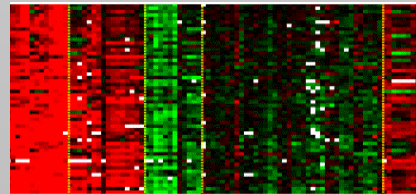
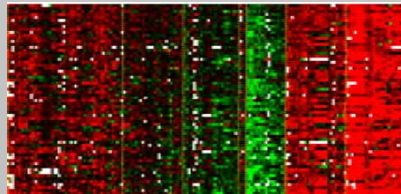
Motifs



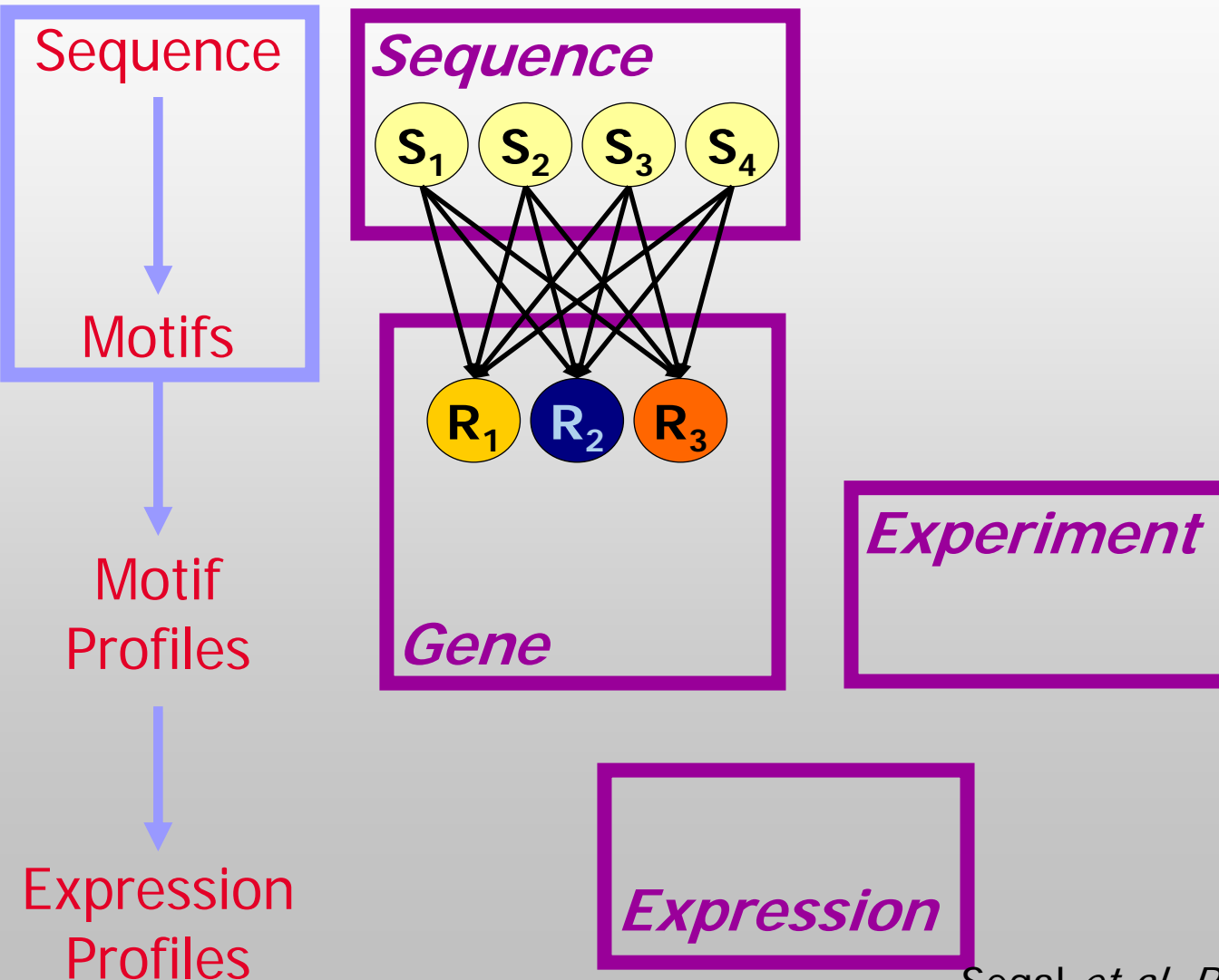
Motif Profiles



Expression Profiles

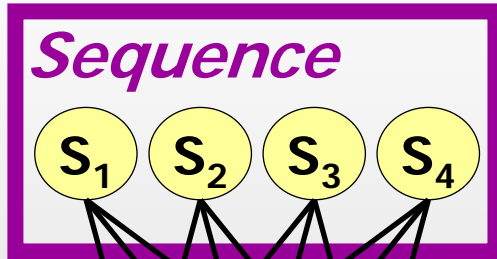


Unified Probabilistic Model

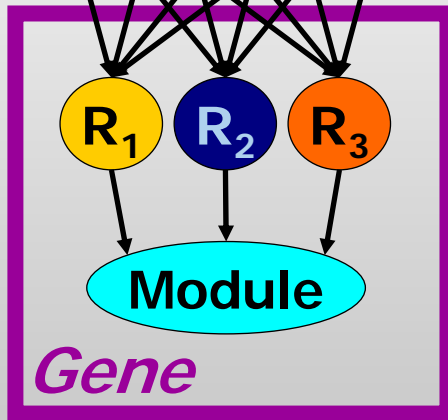


Unified Probabilistic Model

Sequence



Motifs



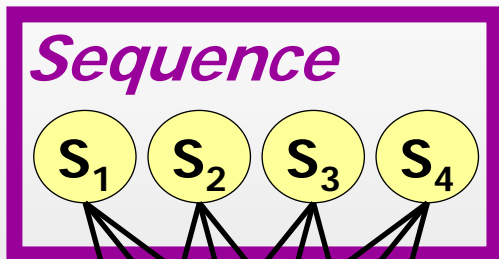
Motif Profiles



Expression Profiles

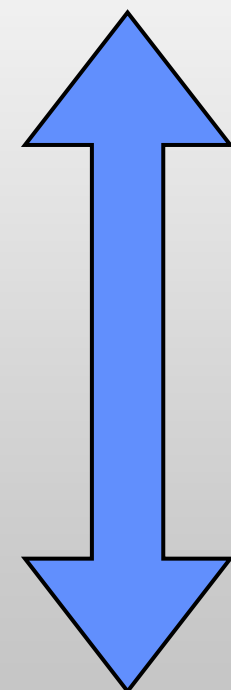
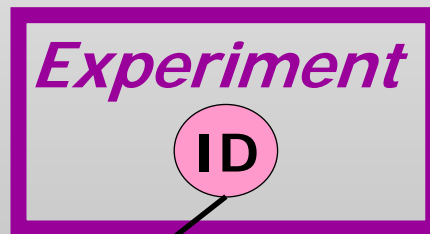
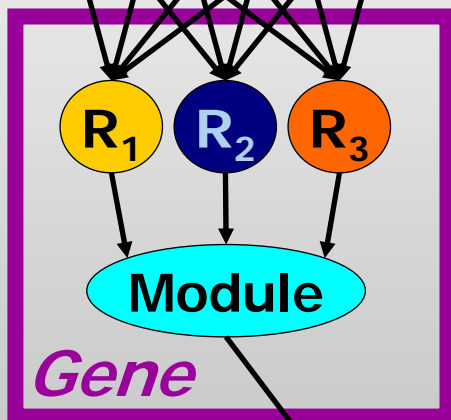
Unified Probabilistic Model

Sequence



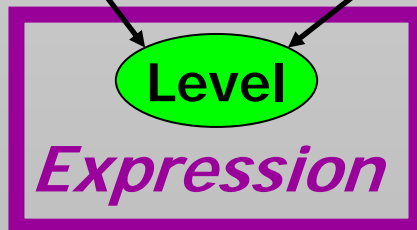
← Observed

Motifs



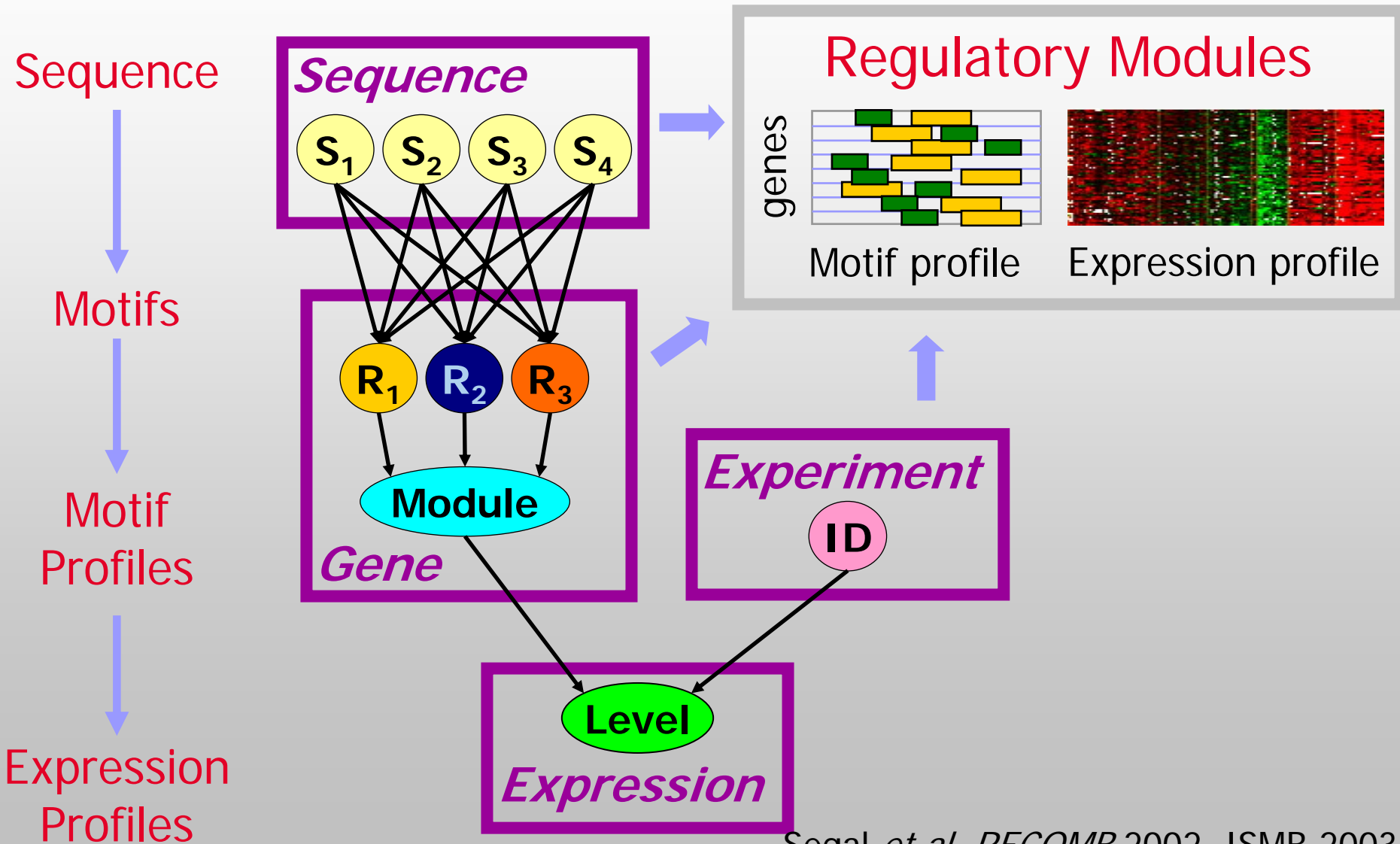
Motif Profiles

Expression Profiles

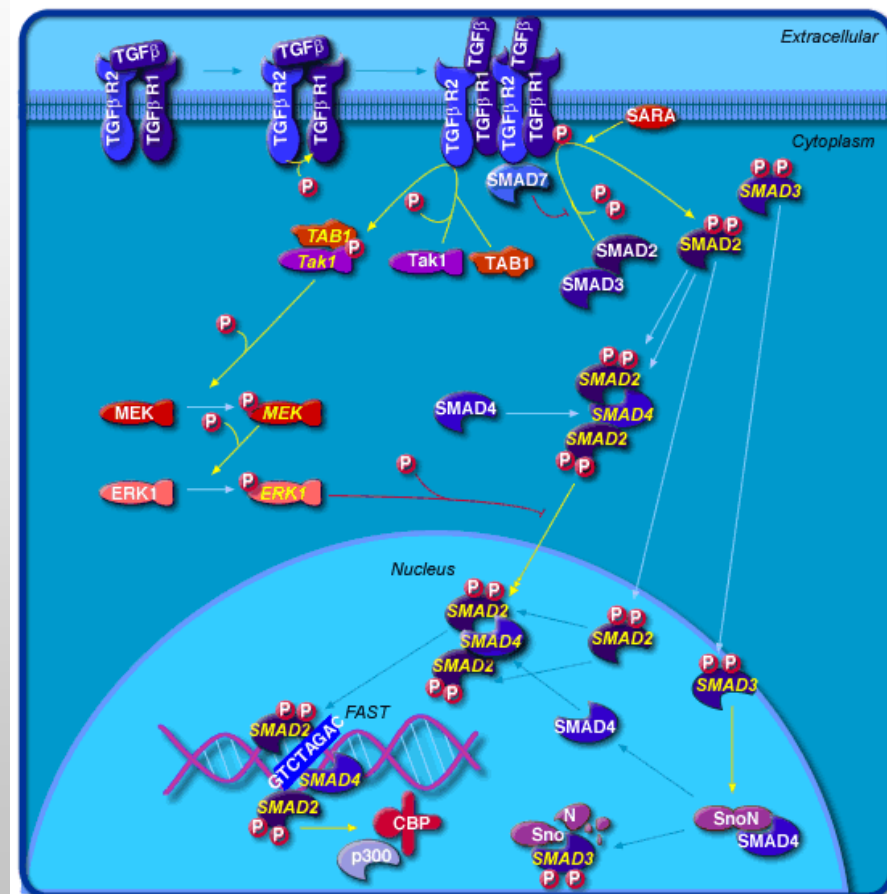
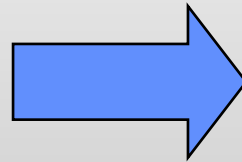
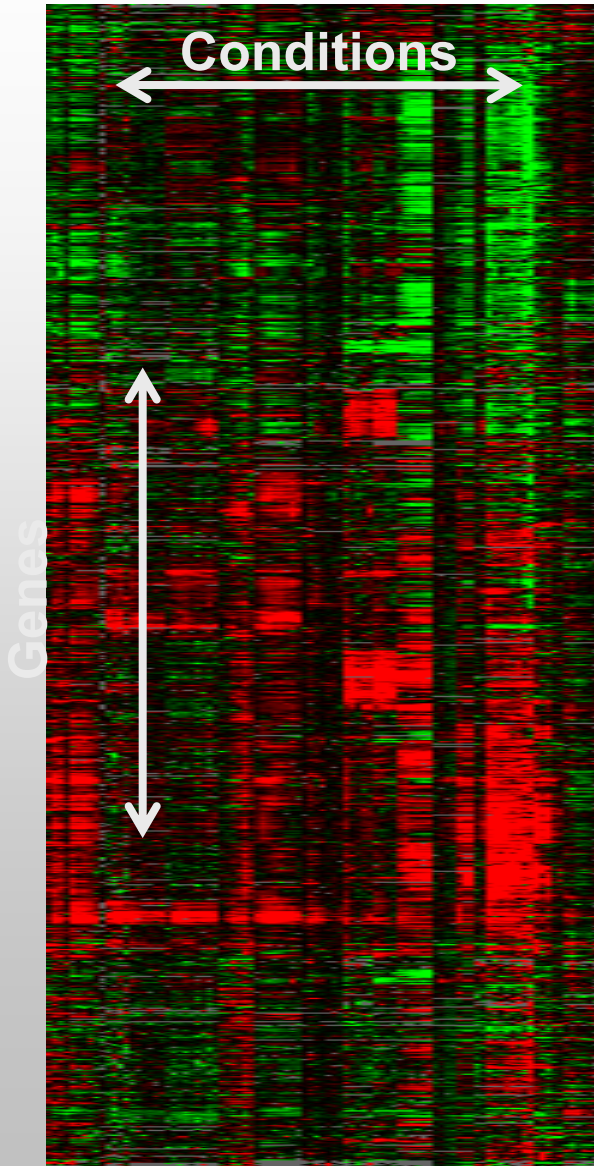


← Observed

Probabilistic Model

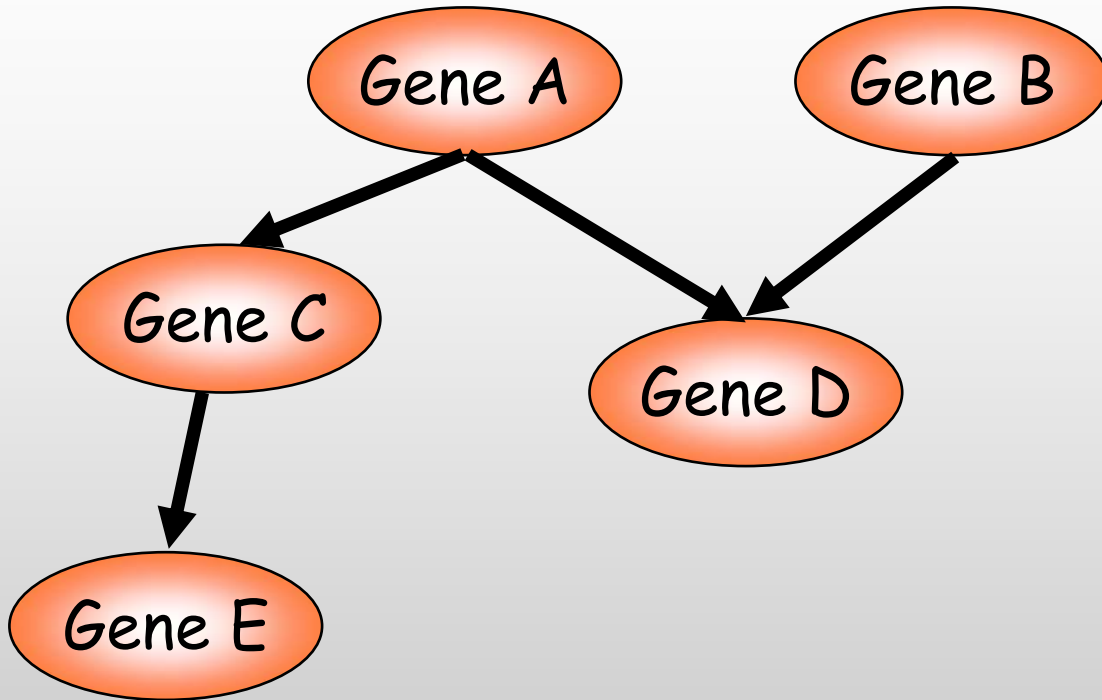


Goal: Reconstruct Cellular Networks

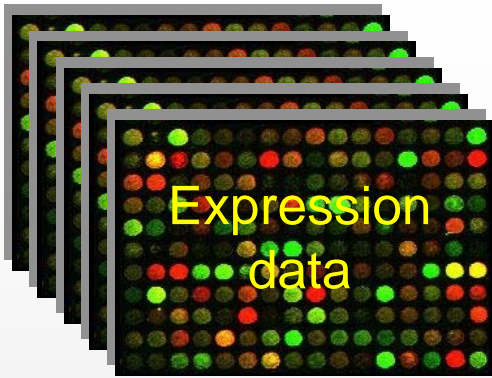


Biocarta. <http://www.biocarta.com/>

Causal Reconstruction for Gene Expression

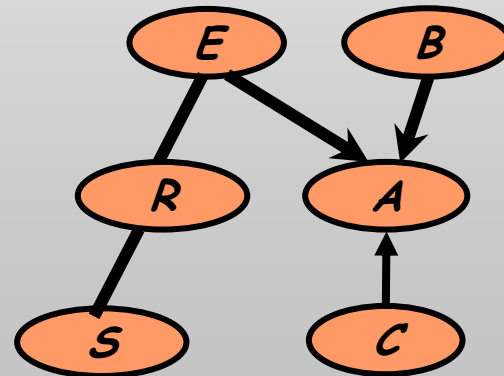


- ◆ Use language of Bayesian networks to reconstruct causal connections



Guided K-means
Discretization

Bayesian Network
Learning Algorithm

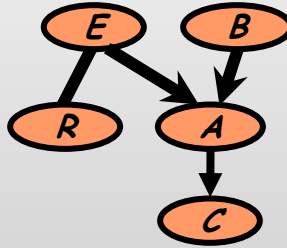
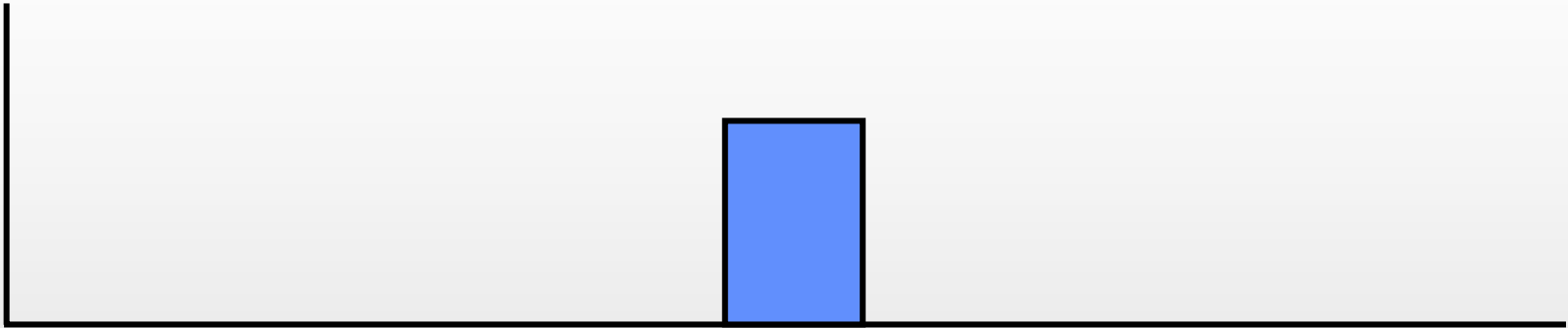


Properties of interactions
among genes

Critical question: do we believe the structure?

Discovering Structure

$P(G|D)$

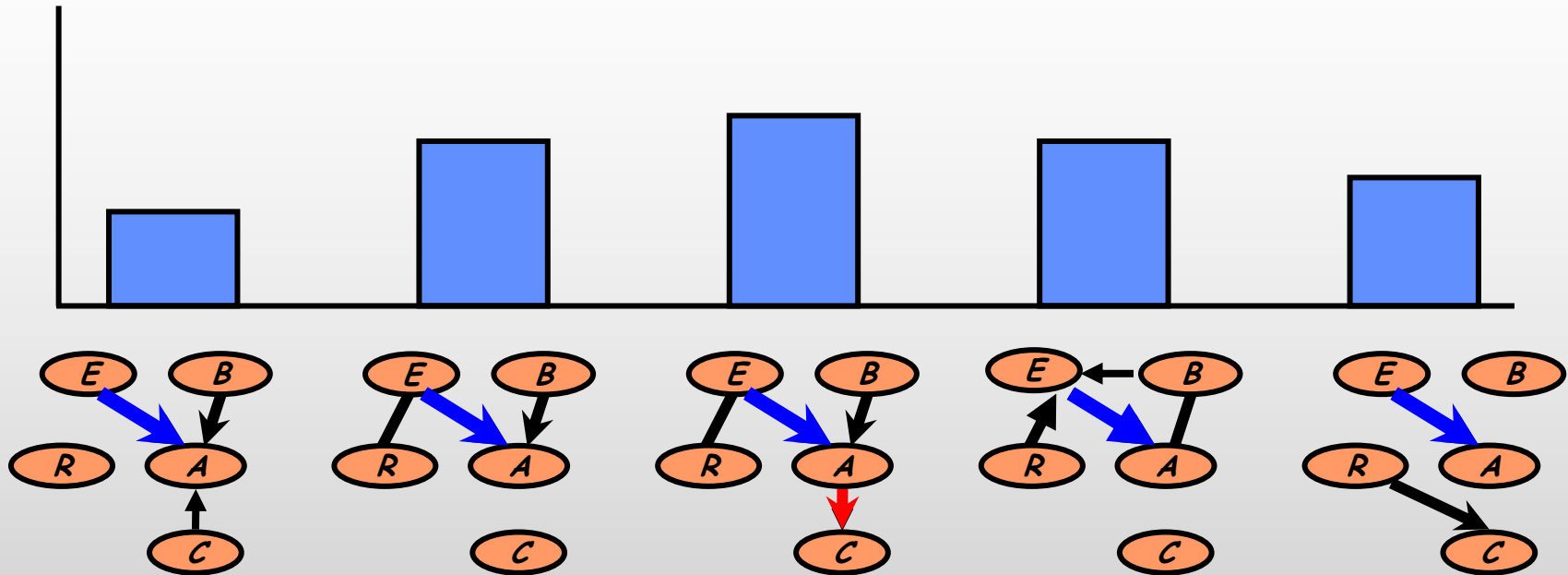


◆ Model selection

- Pick a single high-scoring model
- Use that model to infer domain structure

Discovering Structure

$P(G|D)$



Problem

- Small sample size \Rightarrow many high scoring models
- Answer based on one model often useless
- Want features common to many models

Bayesian Approach

- ◆ Posterior distribution over structures
- ◆ Estimate probability of **features**
 - Edge $X \rightarrow Y$
 - Path $X \rightarrow \dots \rightarrow Y$
 - ...

$$P(f \mid D) = \sum_G f(G) P(G \mid D)$$

Two estimation methods:

Feature of G ,
e.g., $X \rightarrow Y$

• Bootstrap

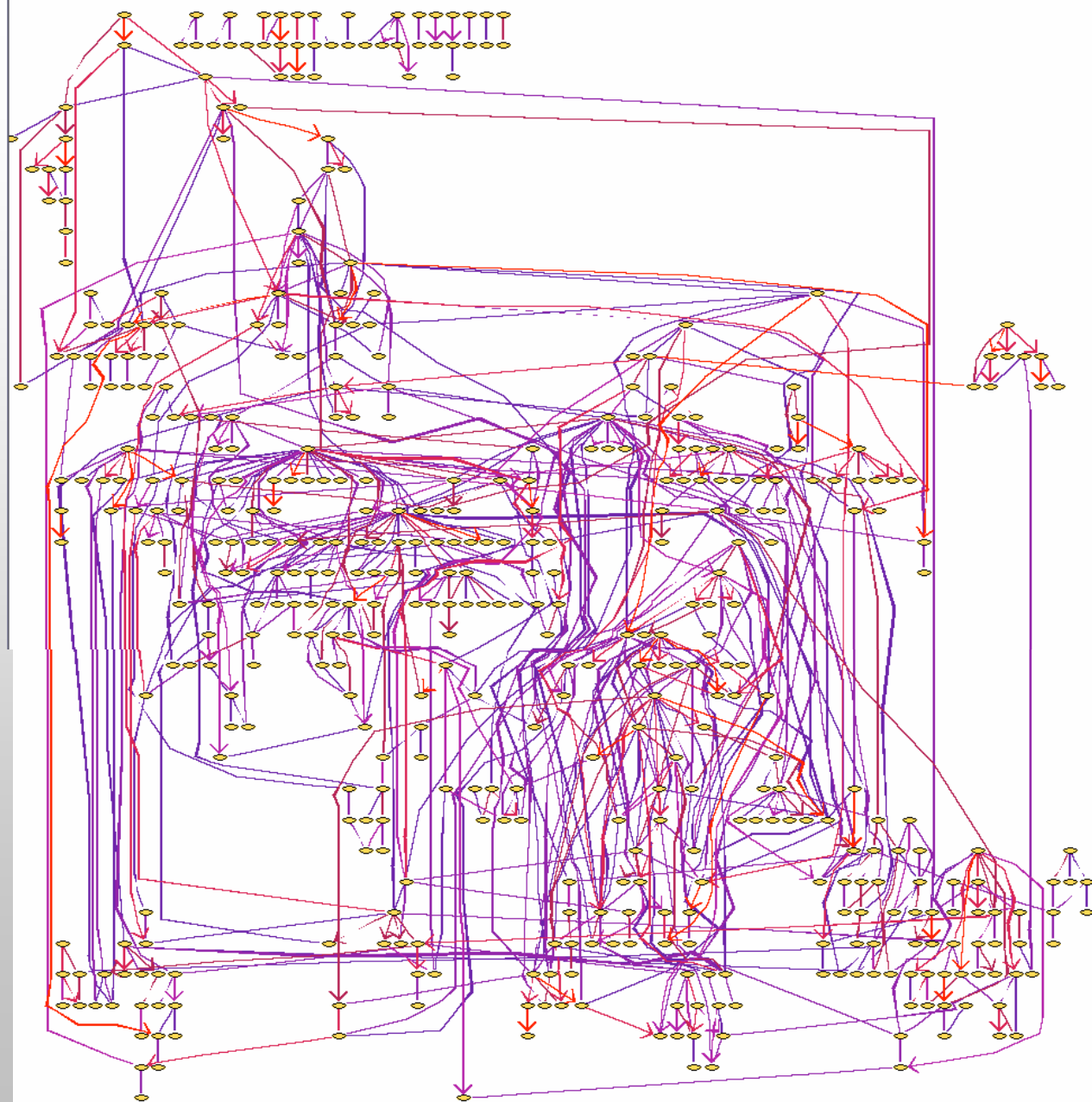
• Markov Chain Monte Carlo

Bayesian score
for G

Indicator function
for feature f

Experiment

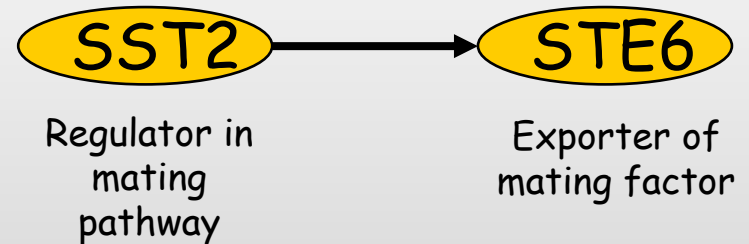
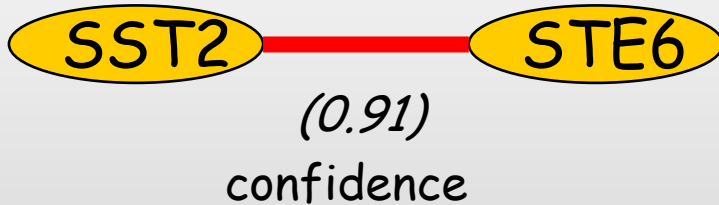
- ◆ 300 deletion knockout in yeast [Hughes et al 2000]
- ◆ 600 genes
- ◆ Color code showing confidence on edges



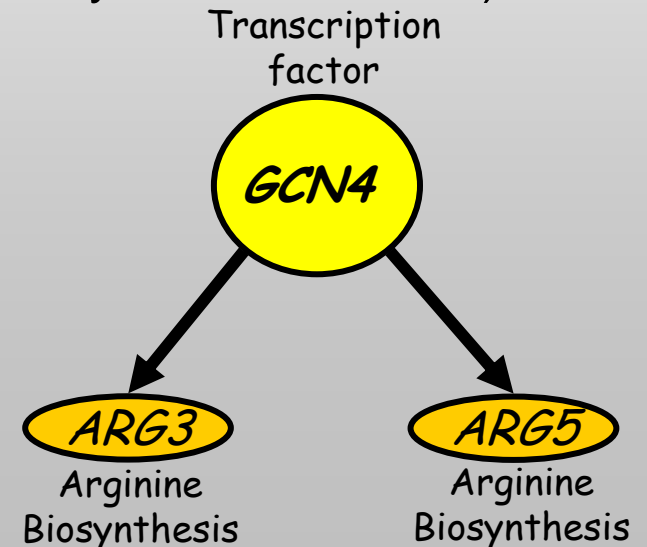
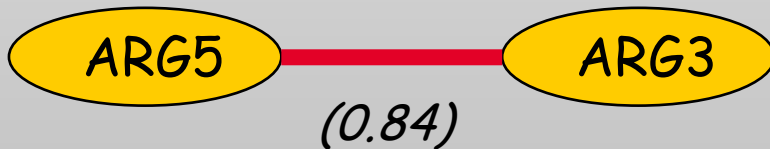
Markov Relations

Question: Do X and Y directly interact?

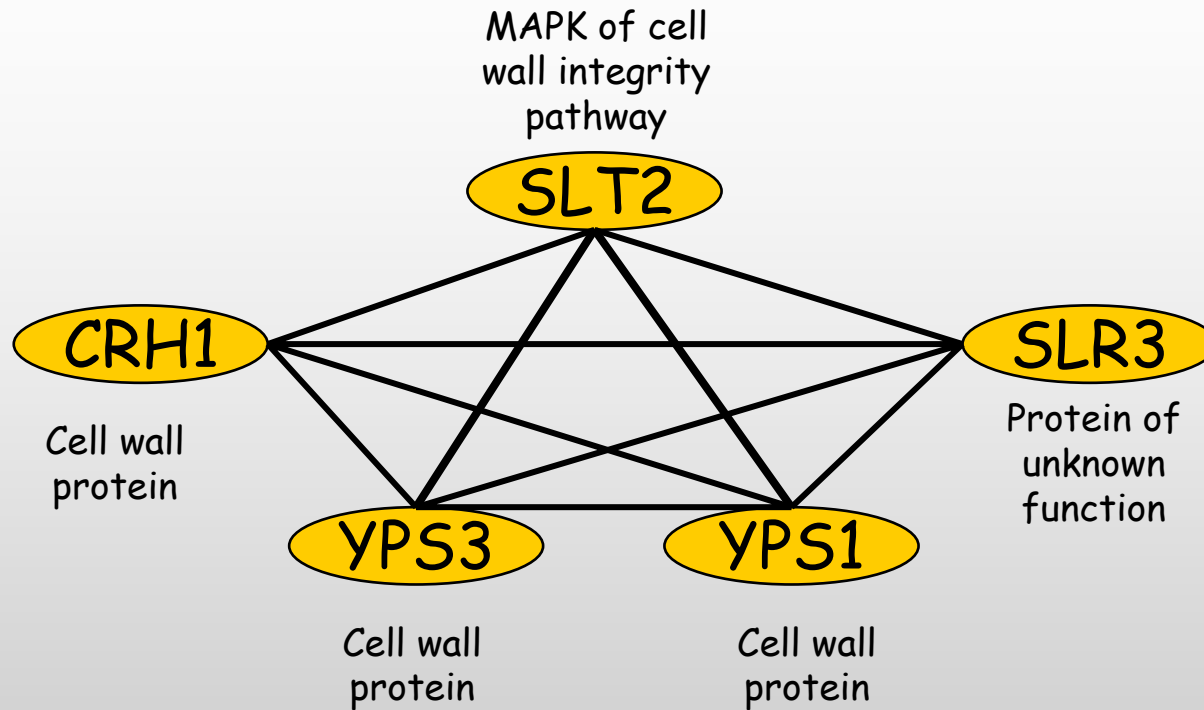
Parent-Child (one gene regulating the other)



Hidden Parent (two genes co-regulated by a hidden factor)

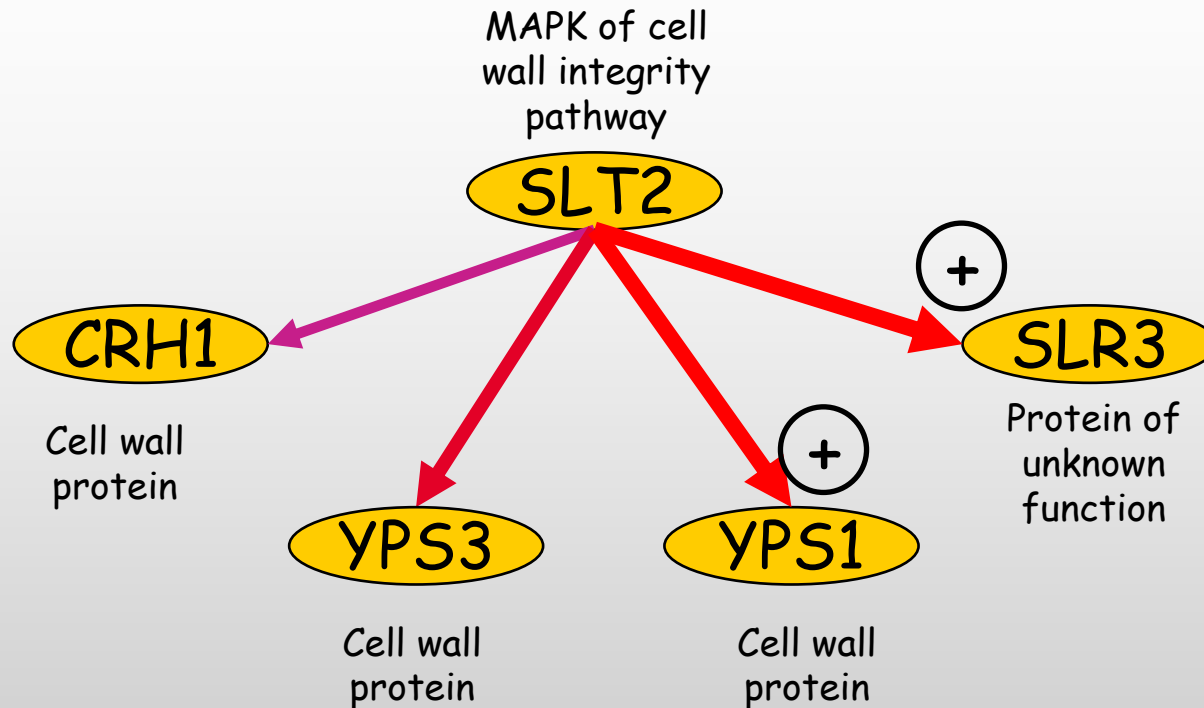


Separators: Intra-cluster Context

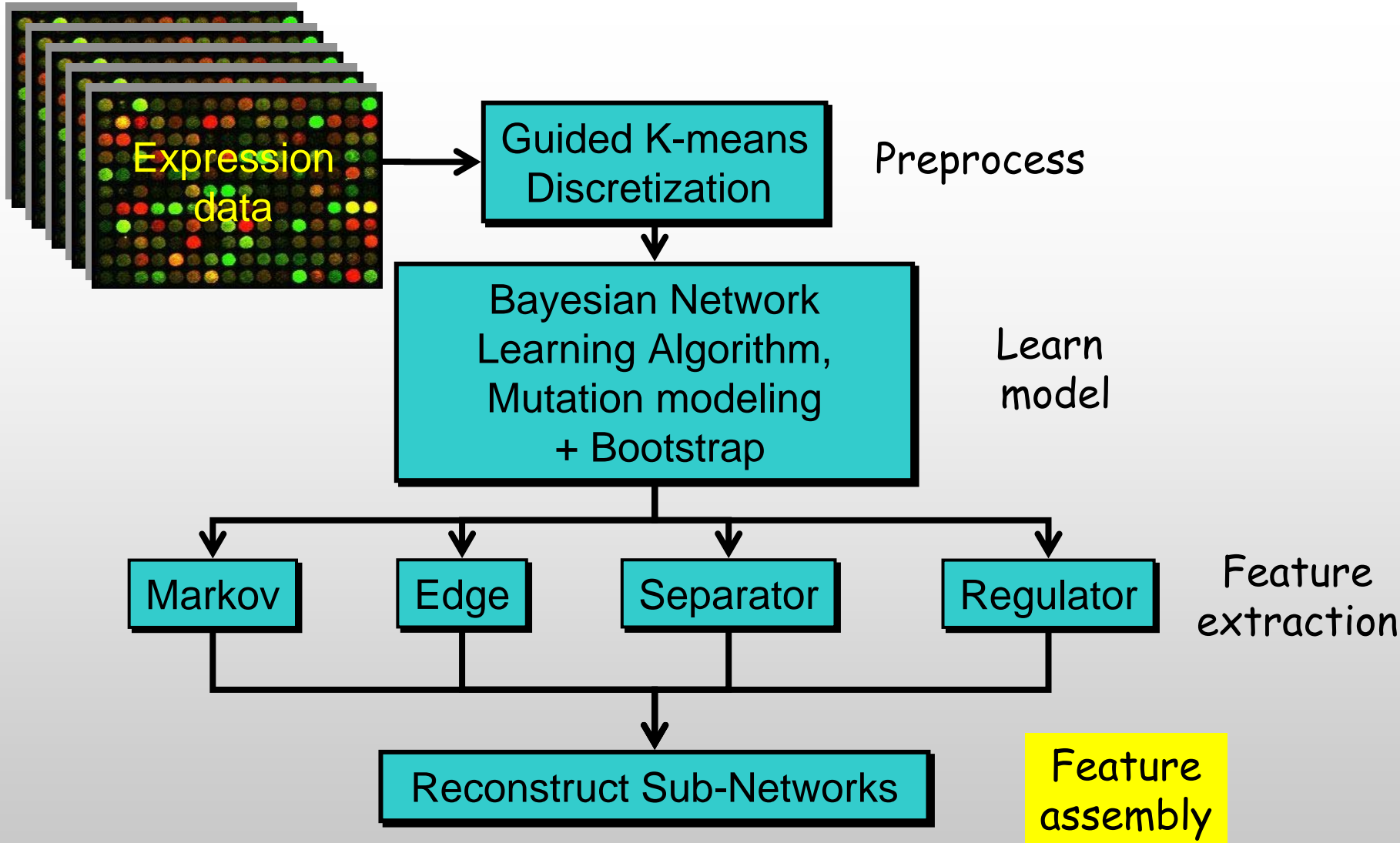


- ◆ All pairs have high correlation
- ◆ Clustered together

Separators: Intra-cluster Context

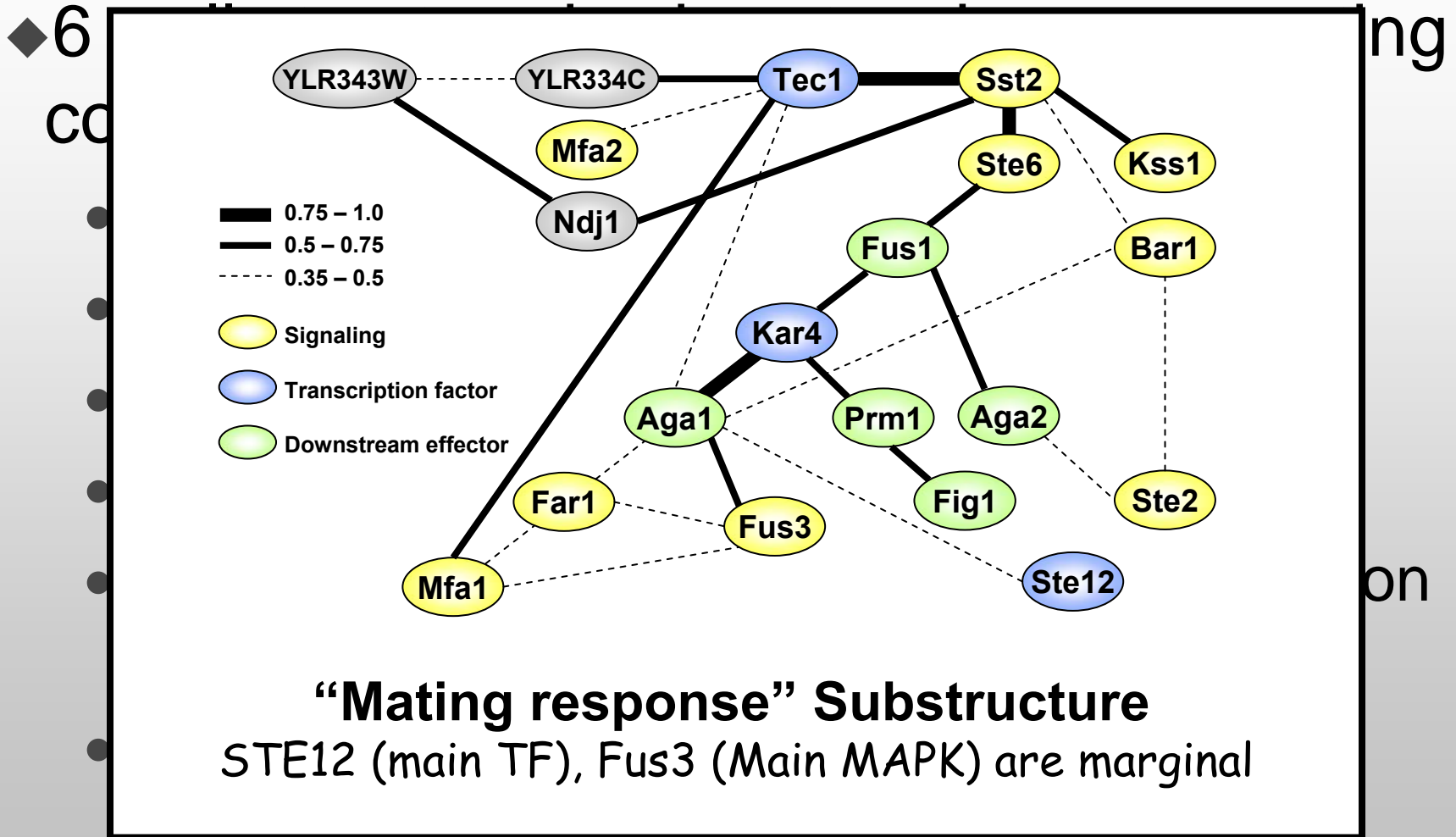


- ◆ **SLT2**: Pathway regulator, explains the dependence
- ◆ Many signaling and regulatory proteins identified as direct and indirect separators

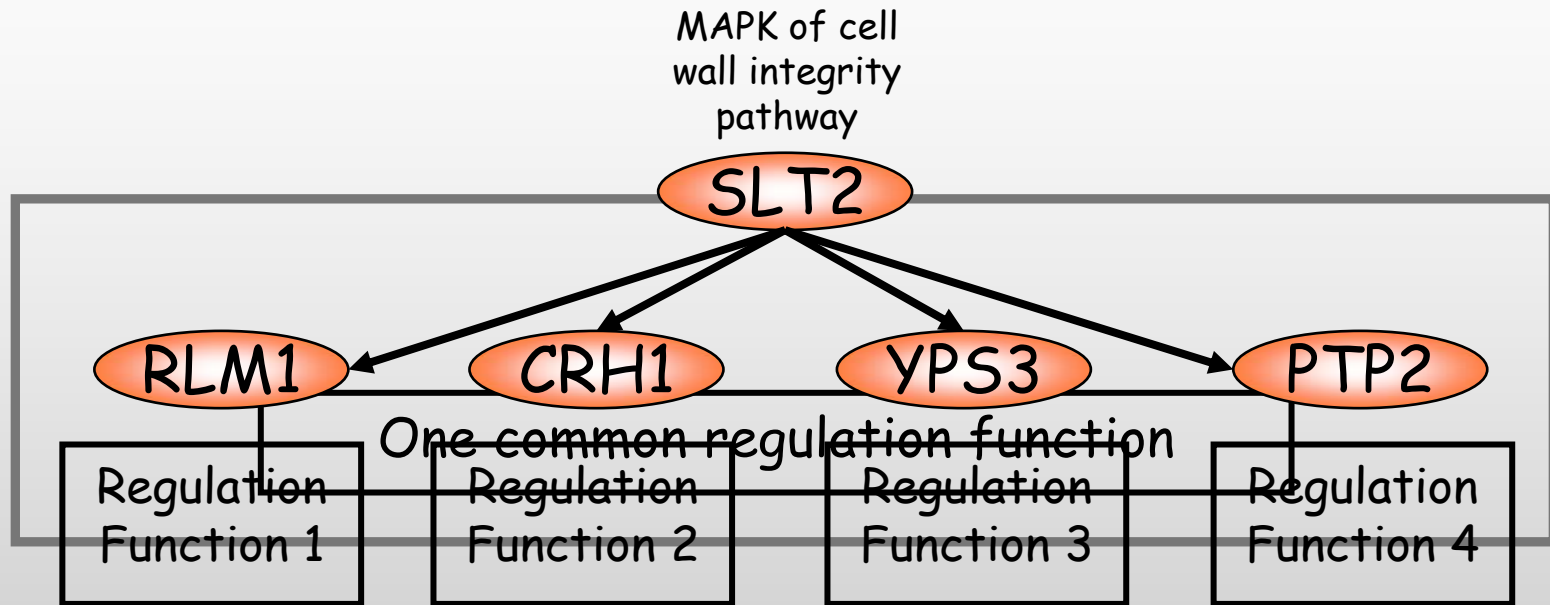


Global network → Local features → Sub-network

Subnetworks in Compendium Dataset

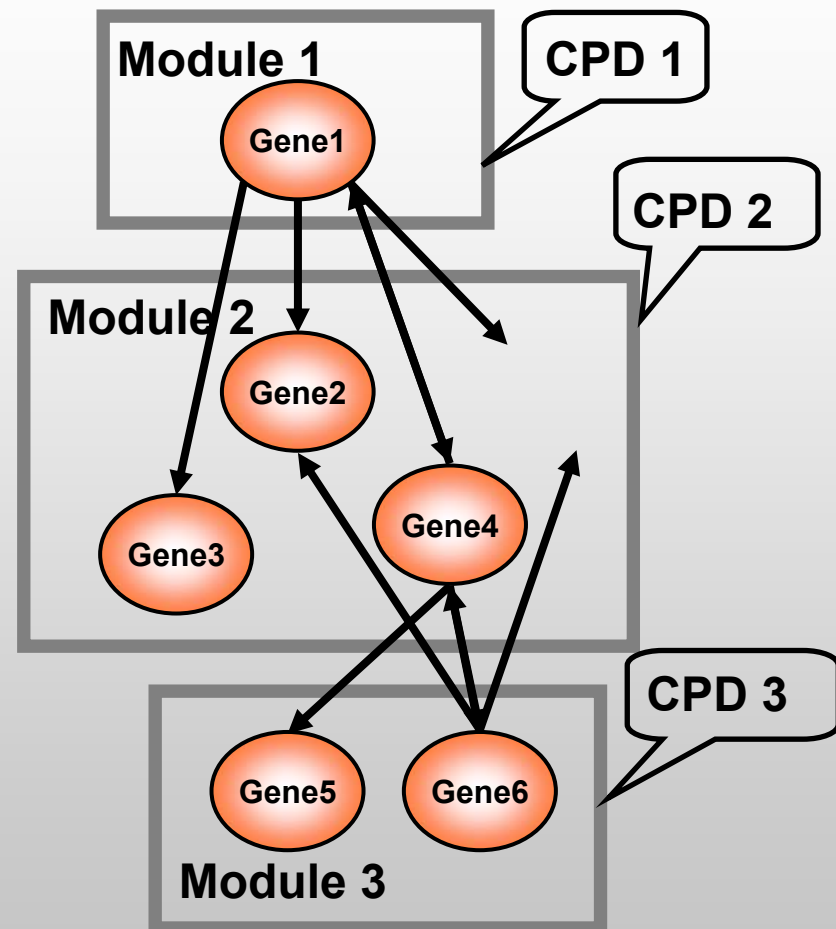
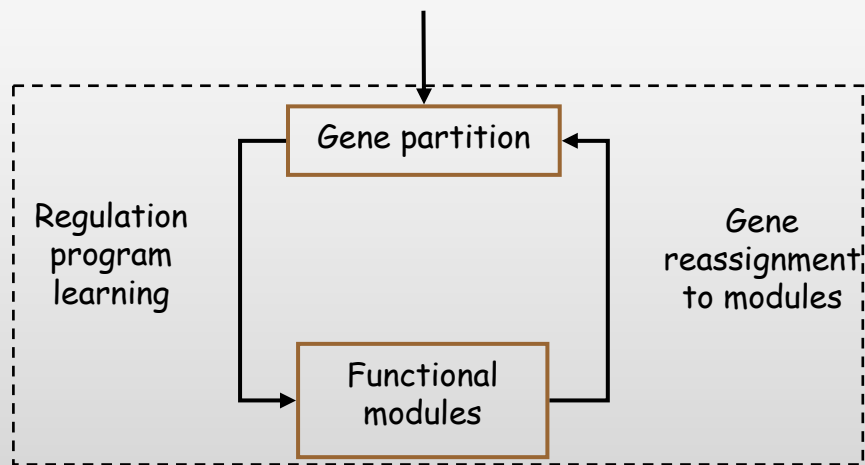


From Networks to Modules

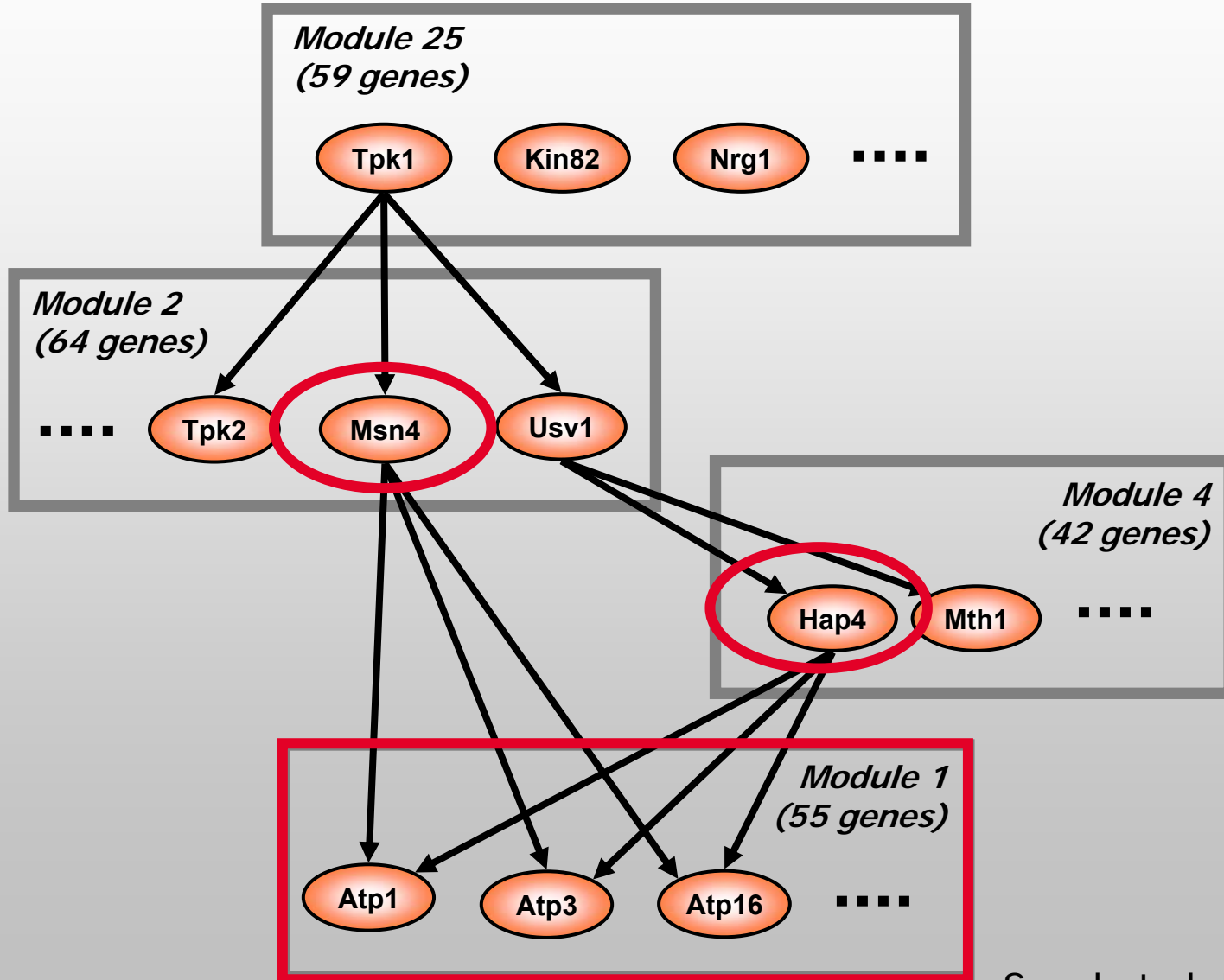


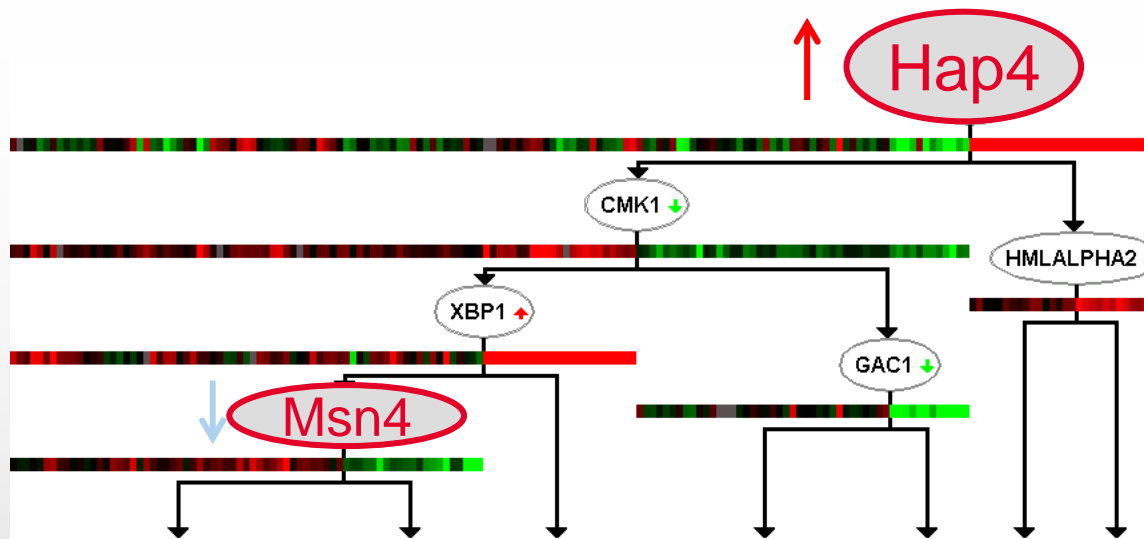
Idea: enforce common regulatory program

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



Learned Network (fragment)



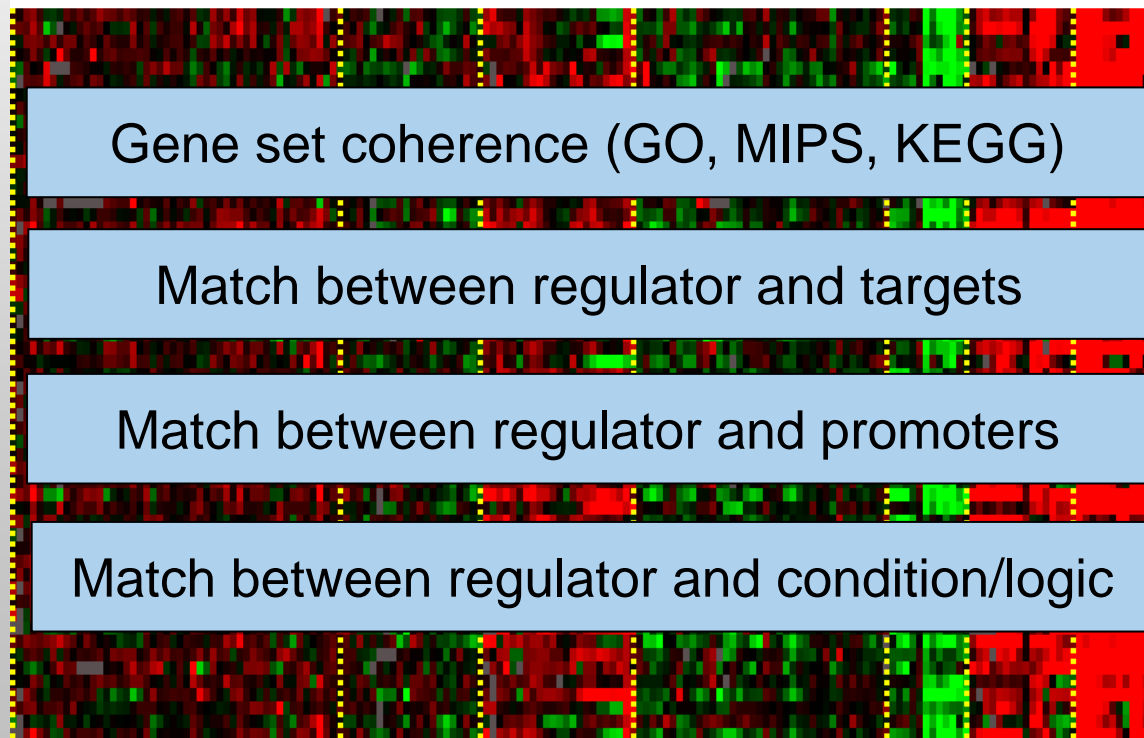
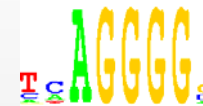


Oxid. Phosphorylation (26, 5×10^{-35})
 Mitochondrion (31, 7×10^{-32})
 Aerobic Respiration (12, 2×10^{-13})

■ HAP4 Motif



■ STRE (Msn2/4)

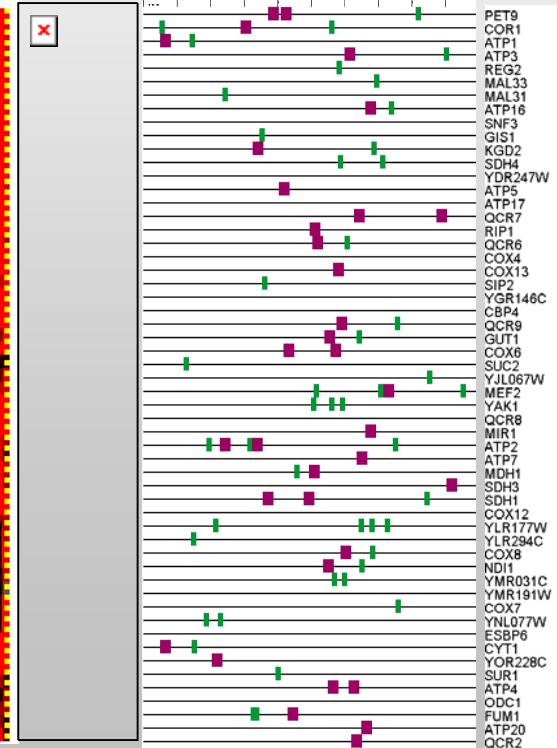


Gene set coherence (GO, MIPS, KEGG)

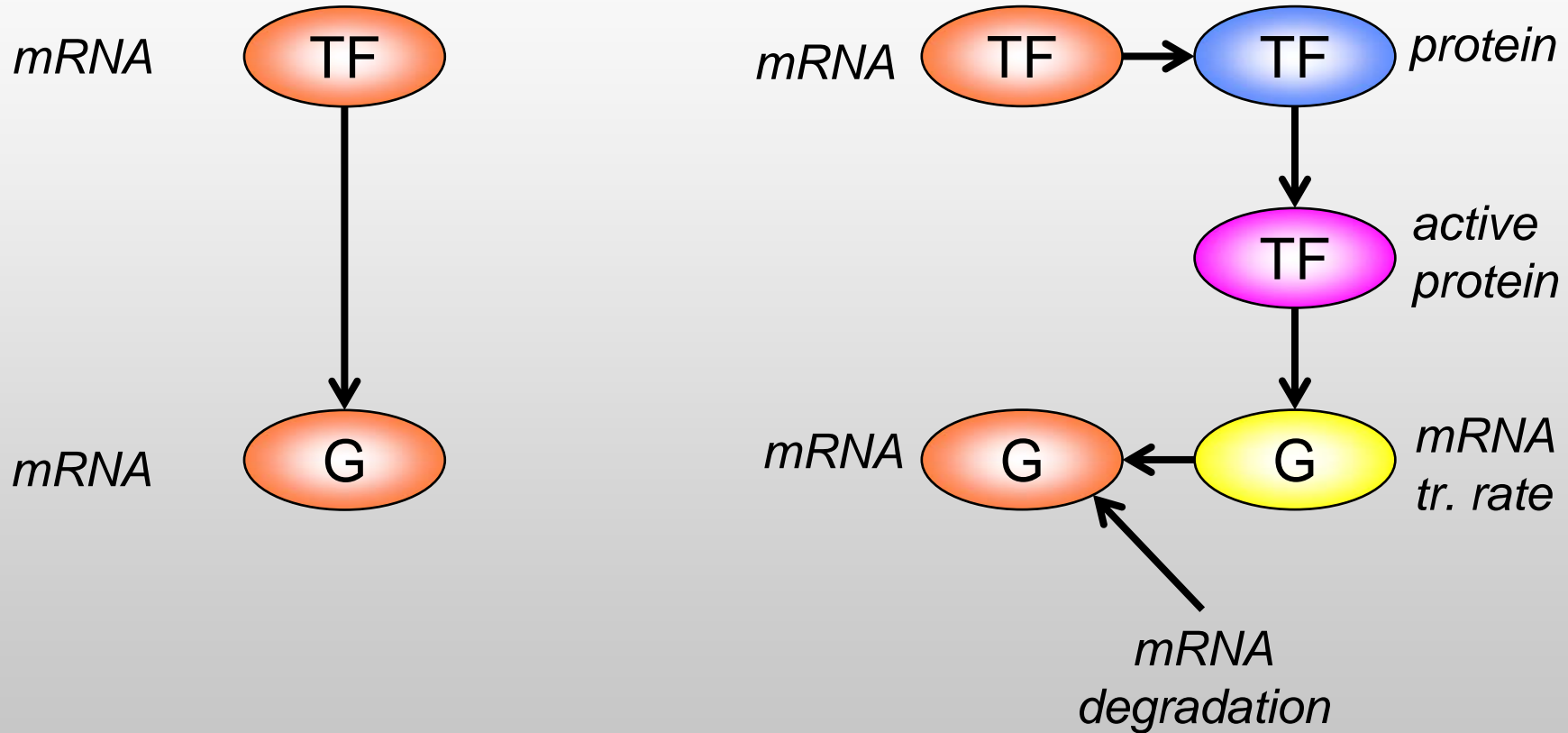
Match between regulator and targets

Match between regulator and promoters

Match between regulator and condition/logic

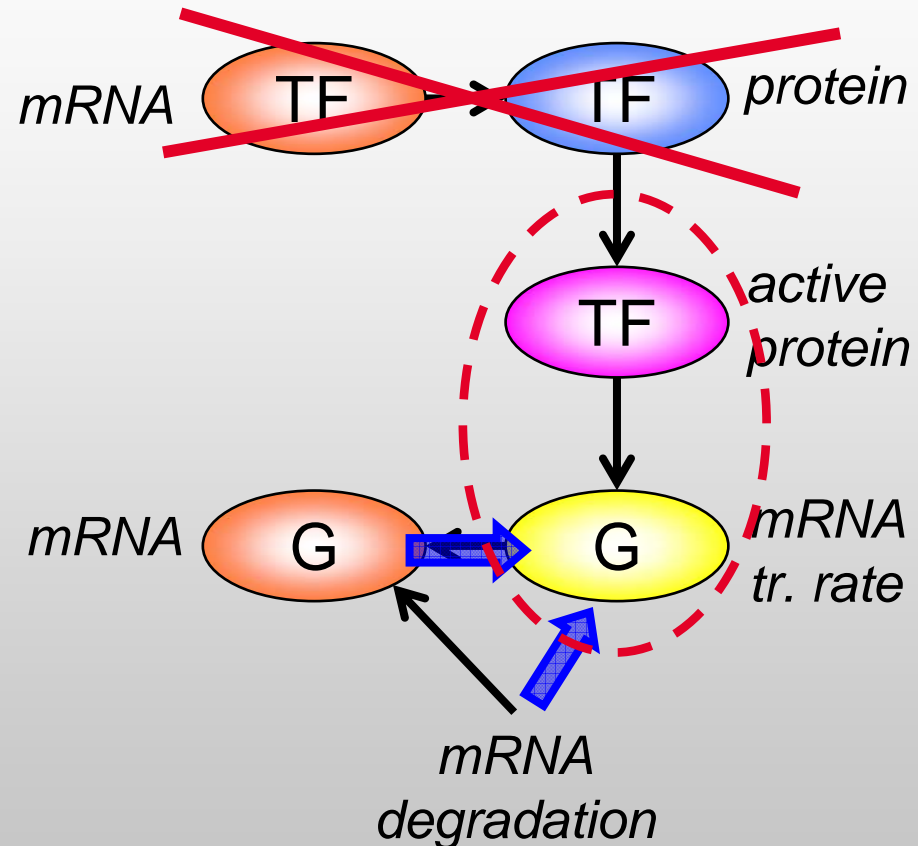


A Major Assumption

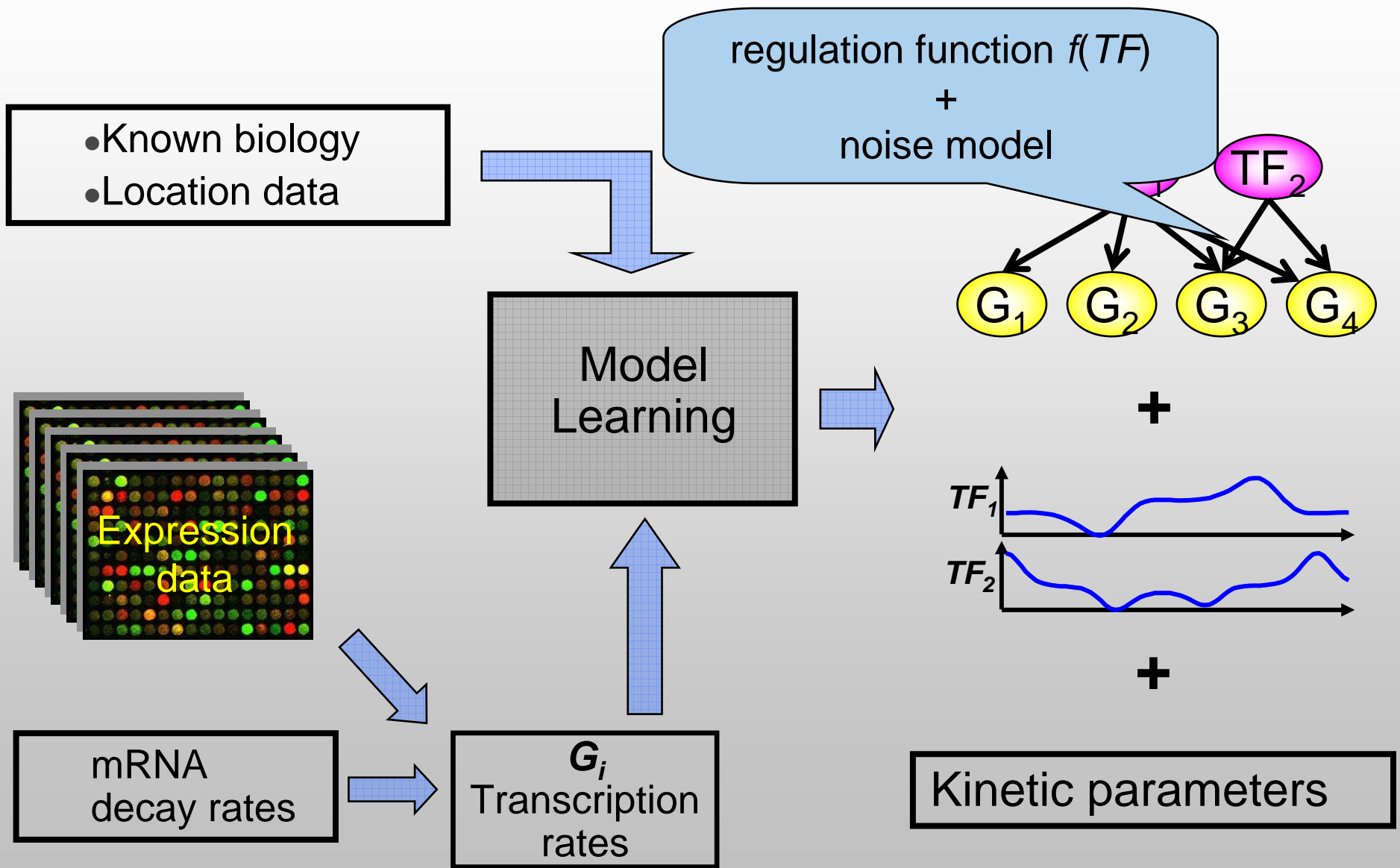


Realistic Regulation Modeling

- ◆ Model the closest connection
- ◆ Active protein levels are not measured
- ◆ Transcript rates are computed from expression data and mRNA decay rates



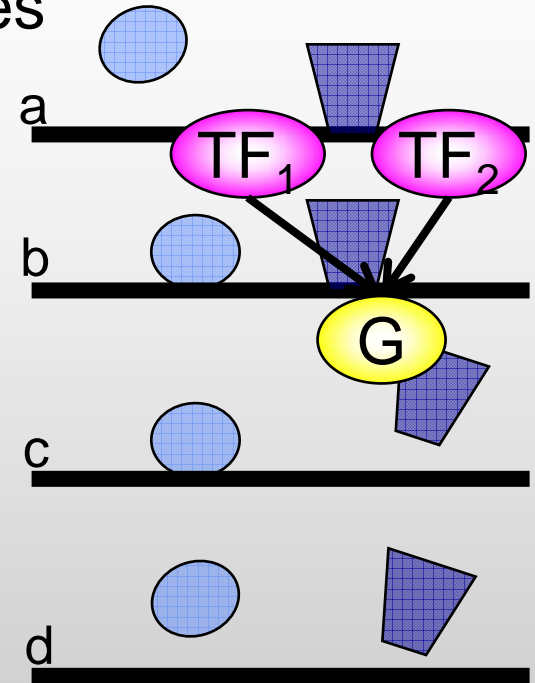
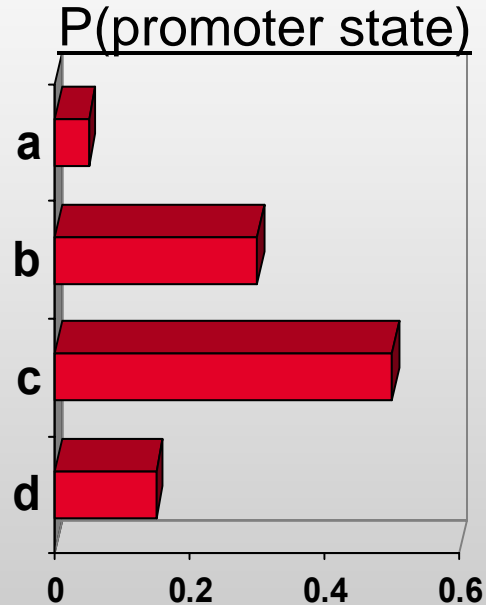
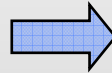
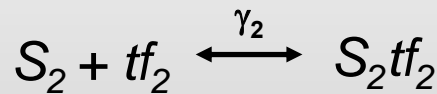
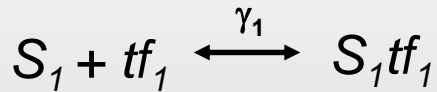
New Proposed Scheme



General Two Regulator Function

I. Compute distribution of promoter states

State Equations

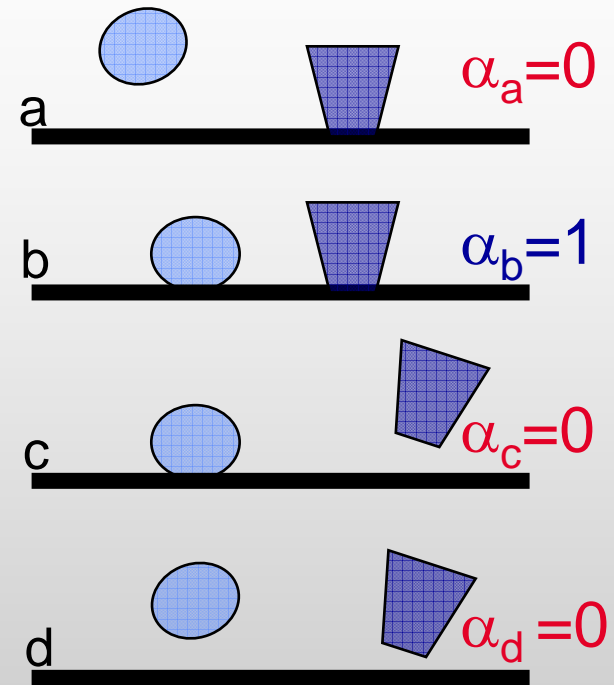
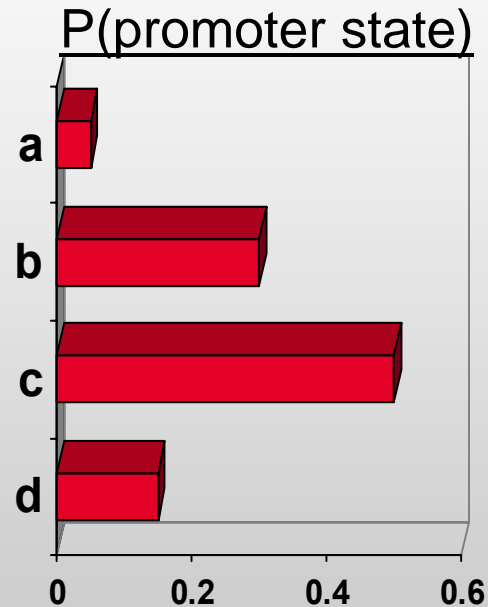
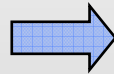
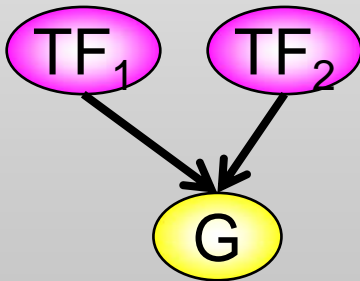
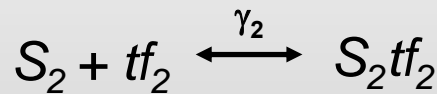
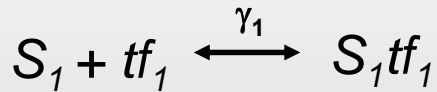


$f(TF_1, TF_2)$ should describe mean transcription rate of G

General Two Regulator Function

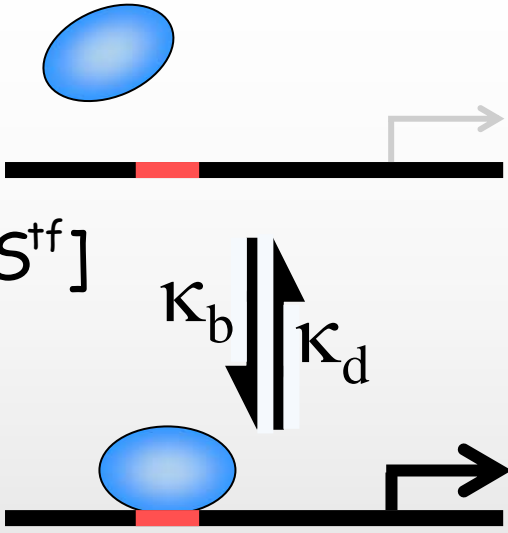
II. Assign activation level to each state

State Equations



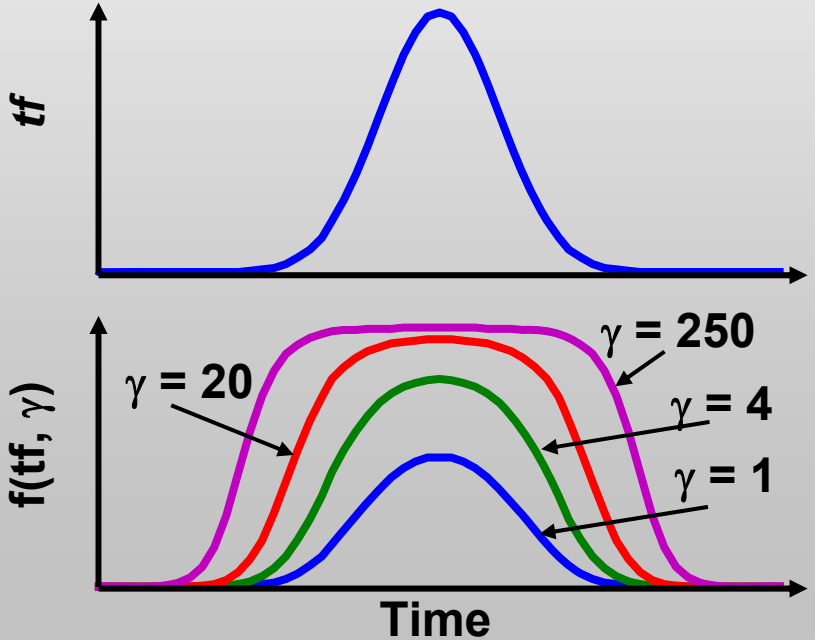
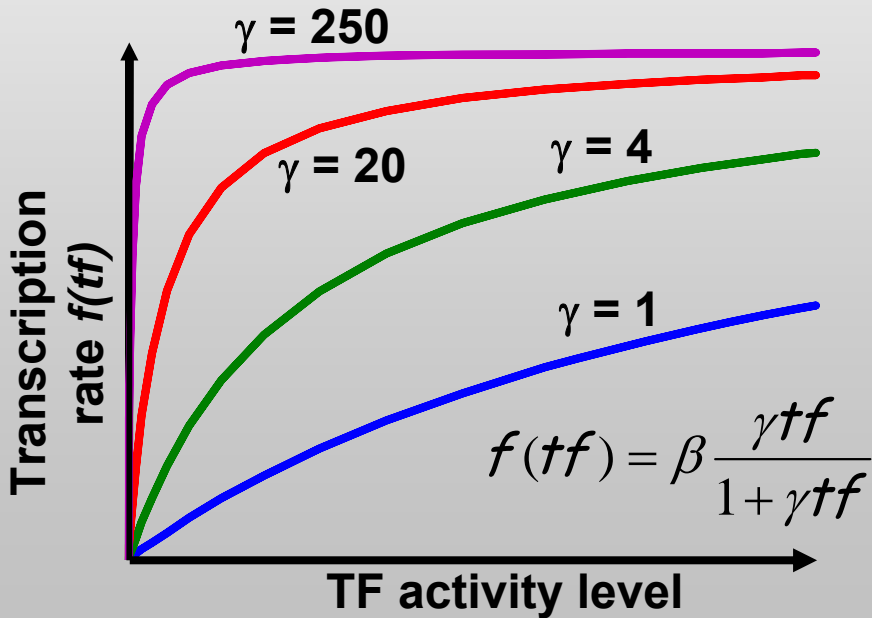
$$f = (\alpha_a X \text{ [bar a] } + \alpha_b X \text{ [bar b] } + \alpha_c X \text{ [bar c] } + \alpha_d X \text{ [bar d] }) \times \beta$$

Example: One Activator Function



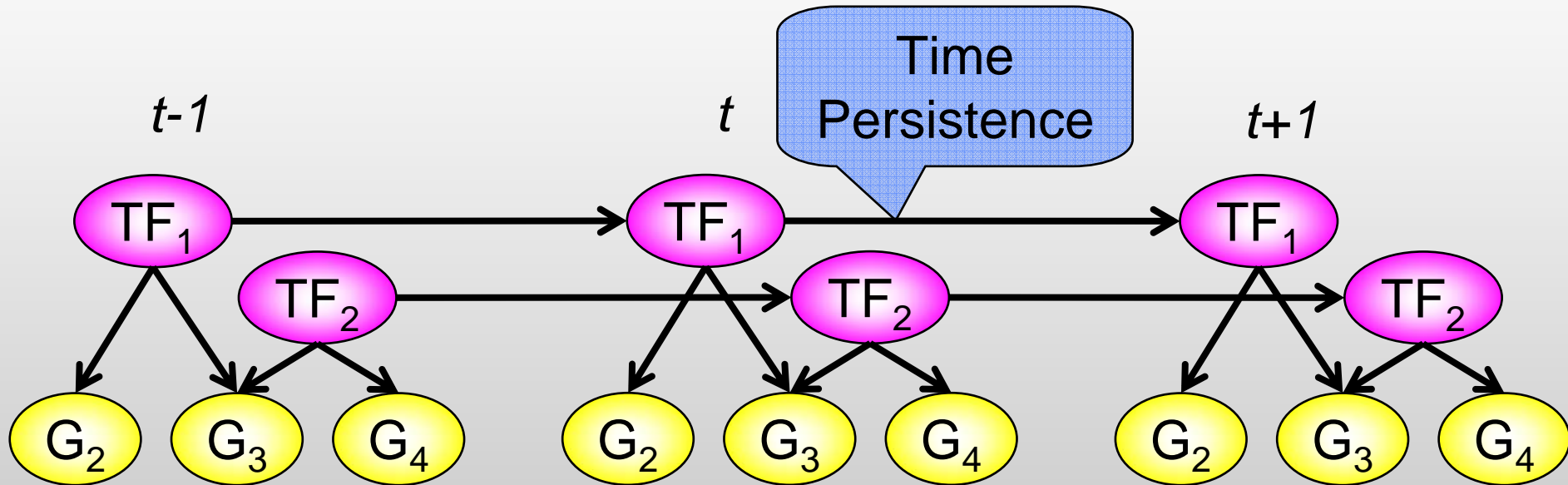
$$k_b[S^-][tf] = k_d[S^{+f}]$$

$$[S^-] + [S^{+f}] = 1$$

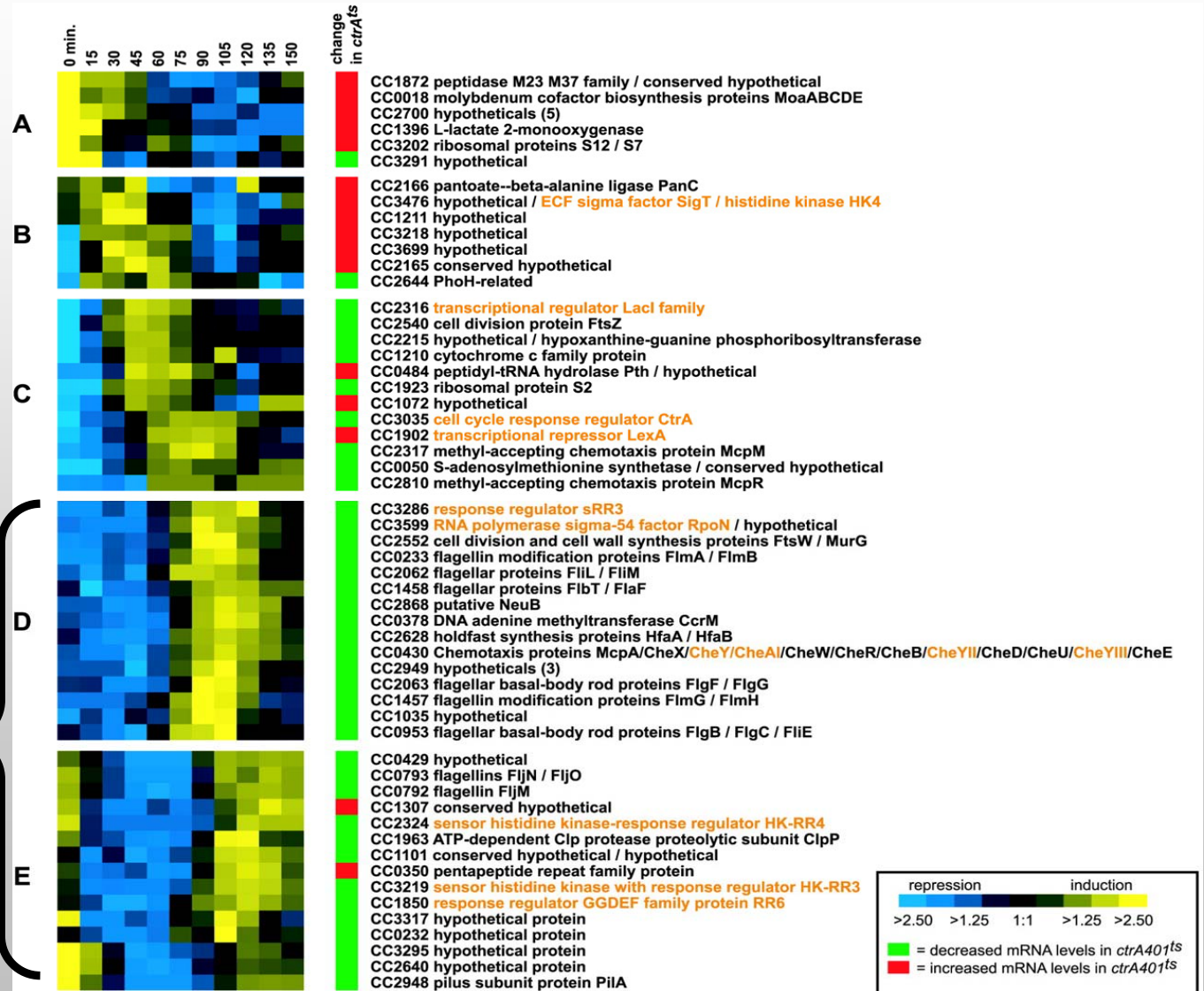


Adding a Temporal Aspect

For time series – add explicit time modeling



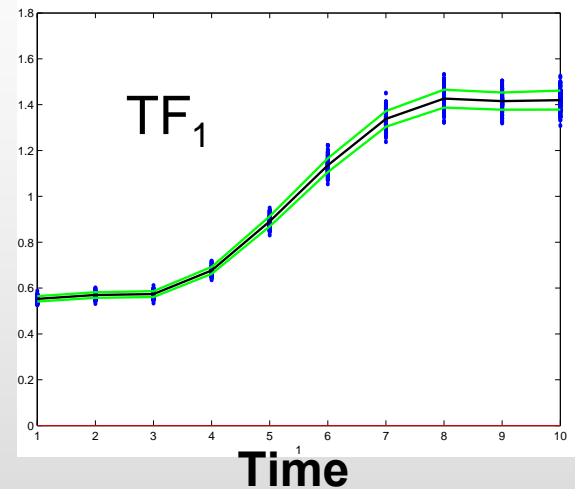
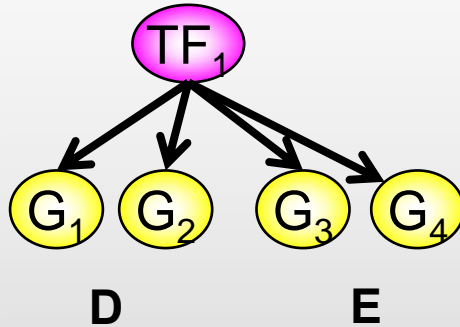
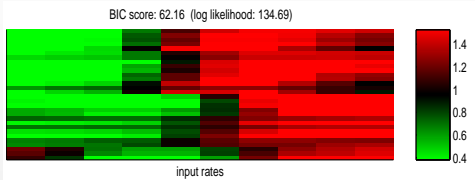
Caulobacter CtrA regulon



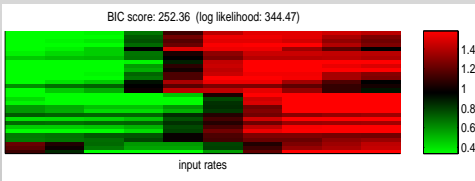
Is There a Second Regulator?

Caulobacter CtrA regulon

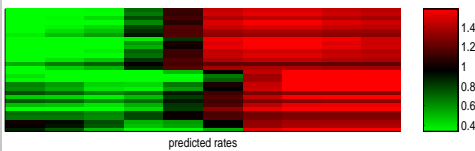
Input
r



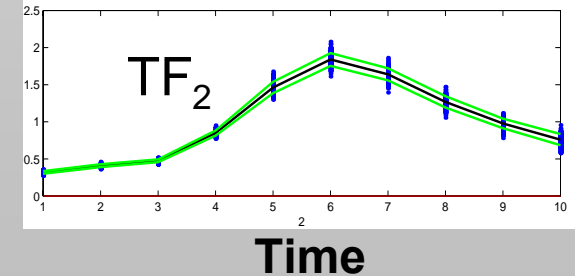
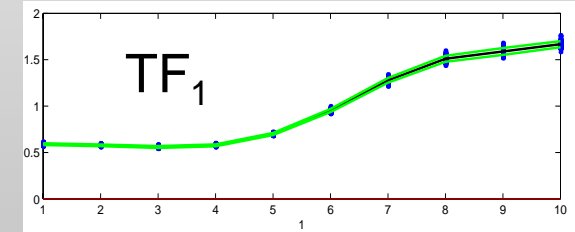
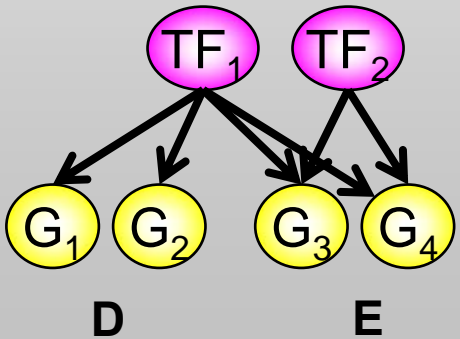
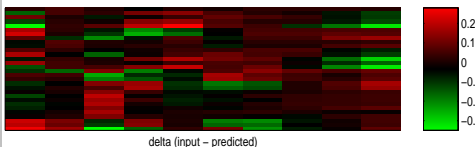
Input
r



Pred.
r



error



i1

show in animation:

input r -> model -> predicted h -> pred r -> error

remove p-val and say in words

stress "realistic"

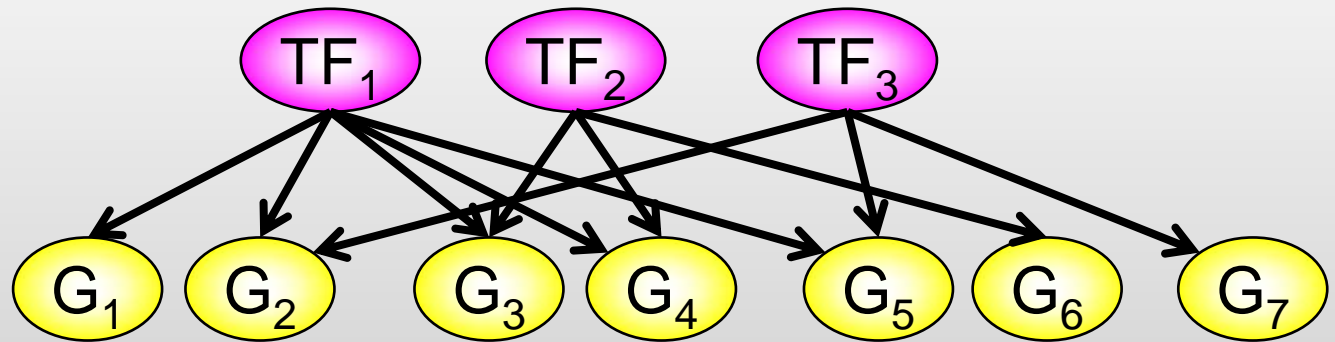
iftach, 3/31/2004

Multiple Regulon Experiments

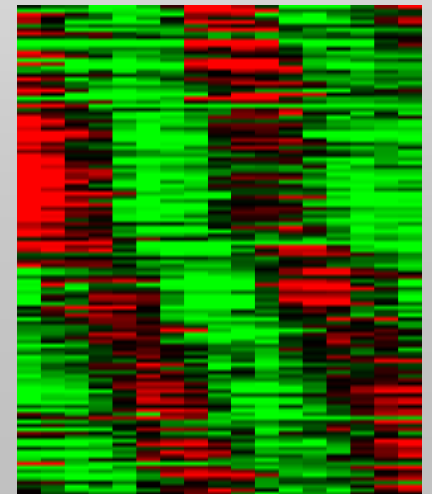
Can we describe the cell transcriptome using a small number of hidden regulators?

few TFs

many genes



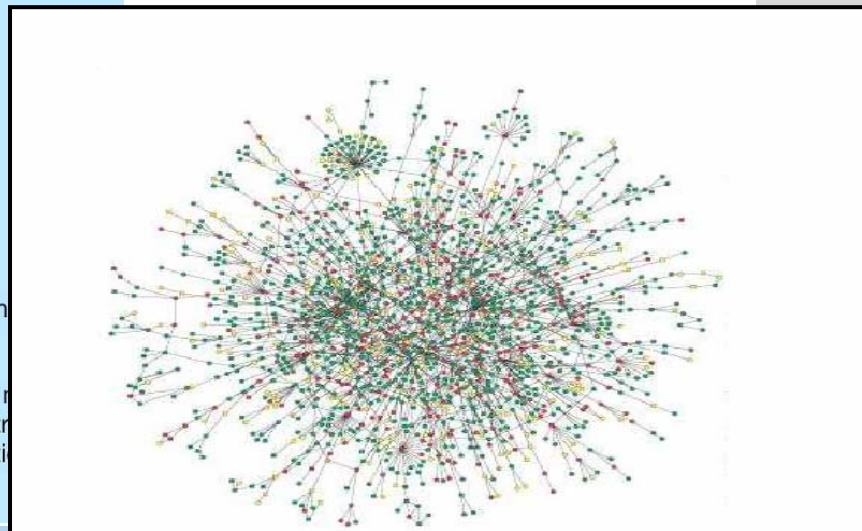
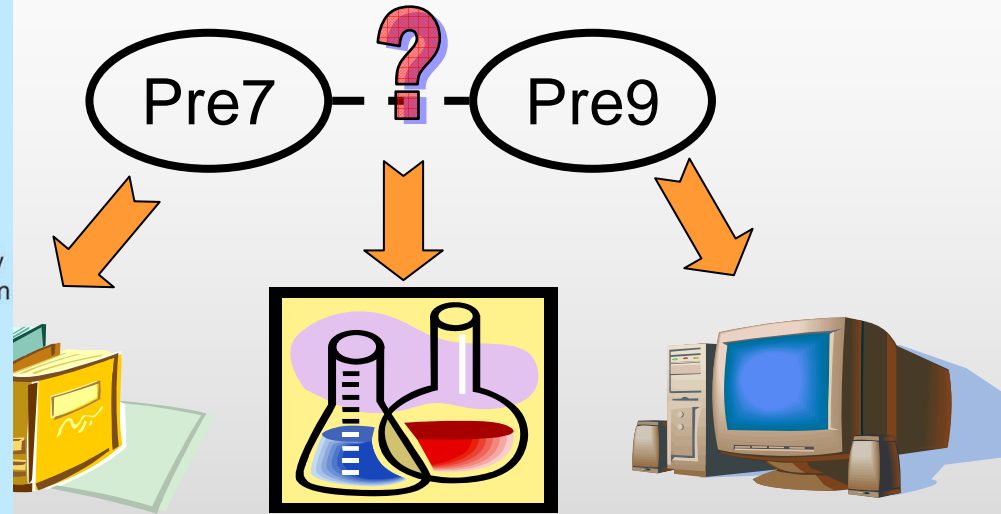
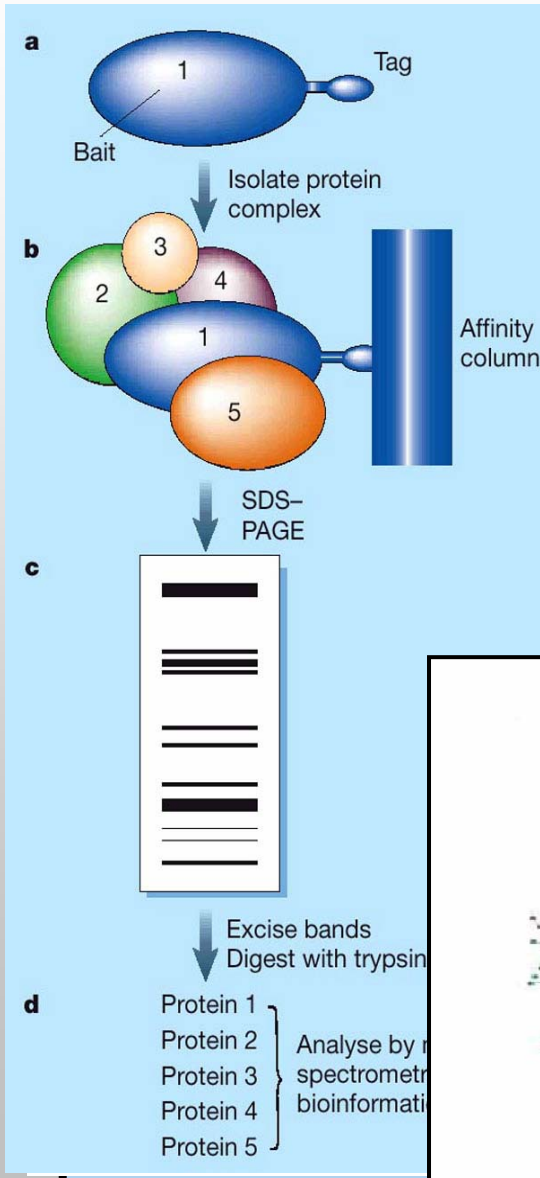
- “Realistic” dimensionality reduction
- Allows prediction of target gene dynamics



Outline

- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ **Markov Networks**
- ◆ Protein-Protein Interactions
- ◆ Discussion

Protein-Protein Interactions



Yeast two hybrid
 et al, 2000
 et al, 2001
 et al, 2002, 2000

Using Protein-Protein Interactions

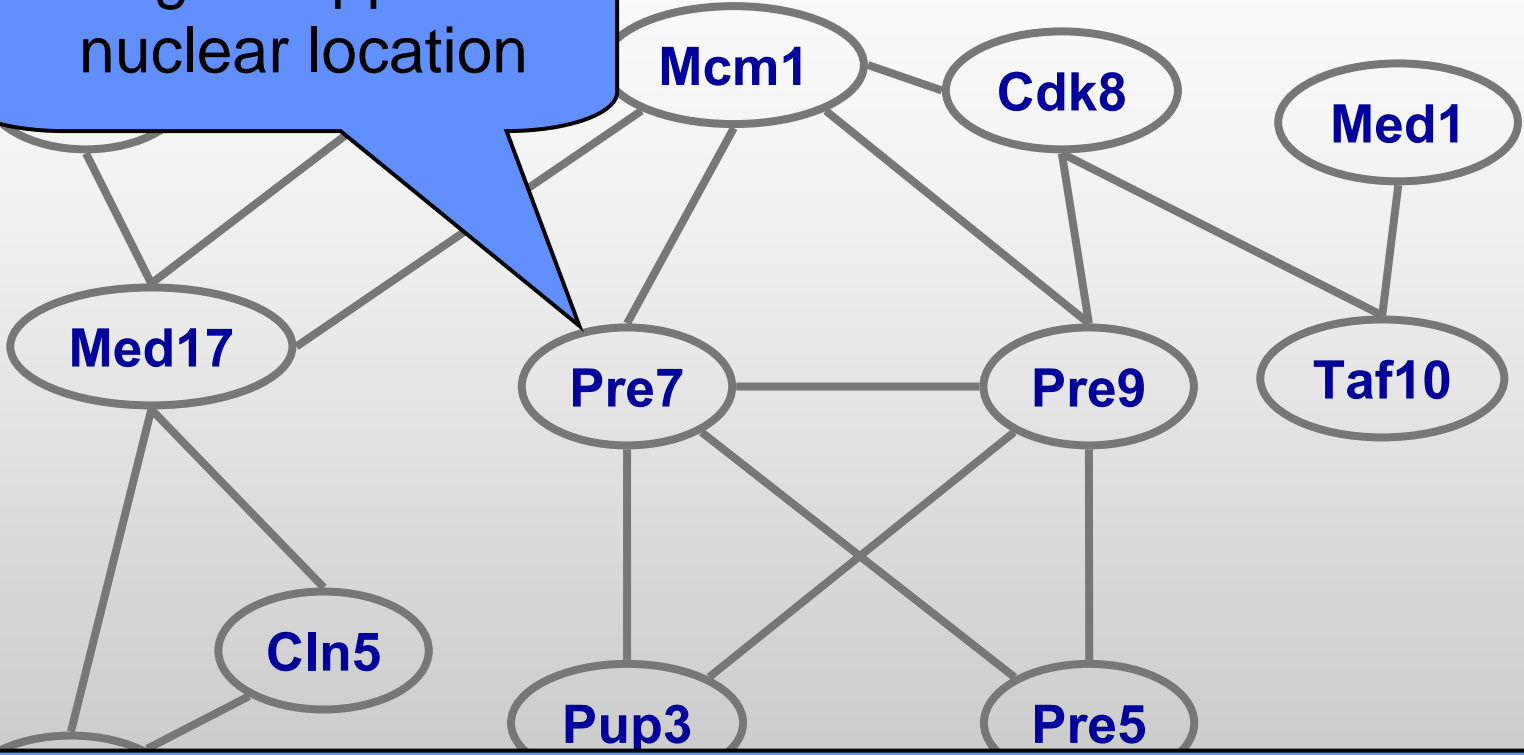
- ◆ Can we use interactions to better understand protein attributes?

Intuition: Interacting proteins tend to be similar

- In the same cellular compartment
- Involved in the same function
- Have similar expression patterns
- ...

Motivation

Stronger support for nuclear location

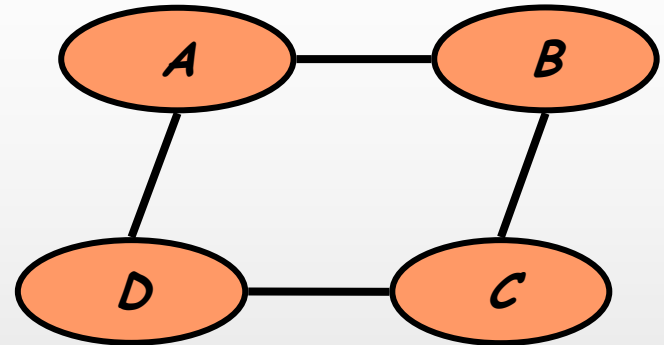


How do we formulate this type of reasoning?

Med1

Markov Networks

A	B	f_1
0	0	-1
0	1	0
1	0	1
1	1	1



Define joint probability distribution

$$P(A, B, C, D) = \frac{1}{Z} \text{Exp}(f_1(A, B) + f_2(B, C) + f_3(C, D) + f_4(D, E))$$

Normalization constant

Potential function

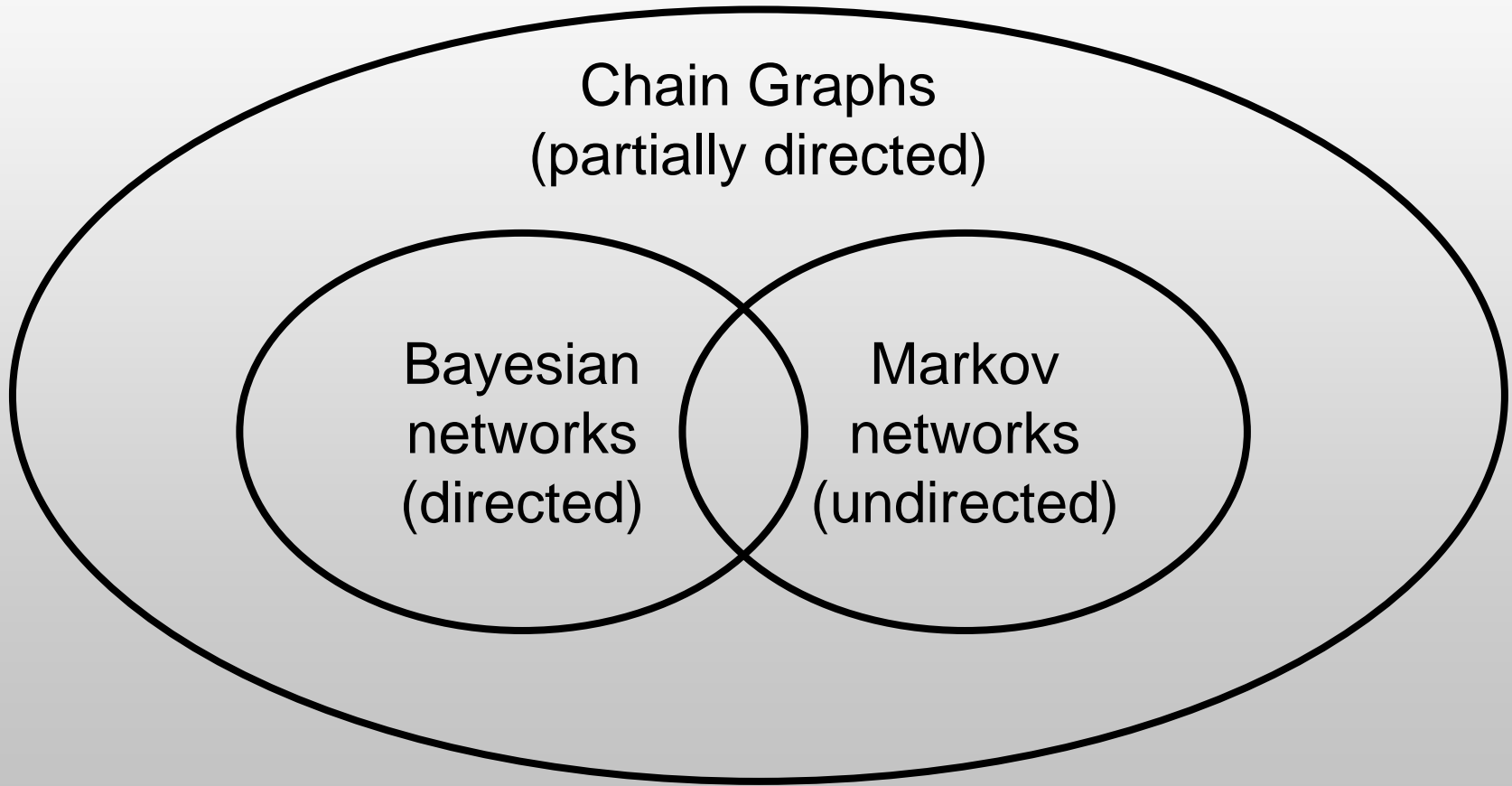
Undirected graph:

- Edge $X - Y$ if there is a factor that includes both X and Y in the same scope

Markov Networks vs Bayesian Network

- ◆ Undirected graph
 - no acyclicity constraints
- ◆ Potential functions
 - less natural and interpretable than conditional distributions
- ◆ Inference is similar to that of Bayesian networks
- ◆ Learning is computationally harder

Relationship between Directed & Undirected Models



Chain Graphs
(partially directed)

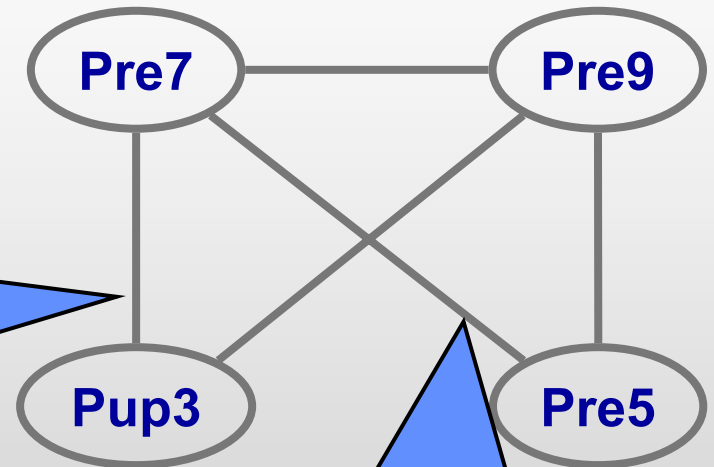
Bayesian
networks
(directed)

Markov
networks
(undirected)

Relational Markov Networks

- ◆ Similar to Relational Bayesian Networks
- ◆ Duplicate potentials

Pre7	Pup3	f
0	0	0
0	1	-1
1	0	-1
1	1	4



Pre7	Pre5	f
0	0	0
0	1	-1
1	0	-1
1	1	4

Outline

- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

Relational Markov Networks for Protein-Protein Interaction

◆ Random variable for each attribute of protein

- *Pre7.nucleus*
 - *Pre7.mitochondria*
 - *Pre7.cytoplasam*
 - ...
 - *Pre7.ribosomal*
 - *Pre7.DNA-binding*
 - ...
- } Cellular compartment
- } Functional category (GO)

◆ Introduce potential between interacting pairs

$$\prod_{p \text{ interacts with } q} f_{\text{nucleus}}(p.\text{nucleus}, q.\text{nucleus})$$

Relational Markov Networks for Protein-Protein Interaction

Three phase process

◆ Model construction

- Use interaction network to construct model

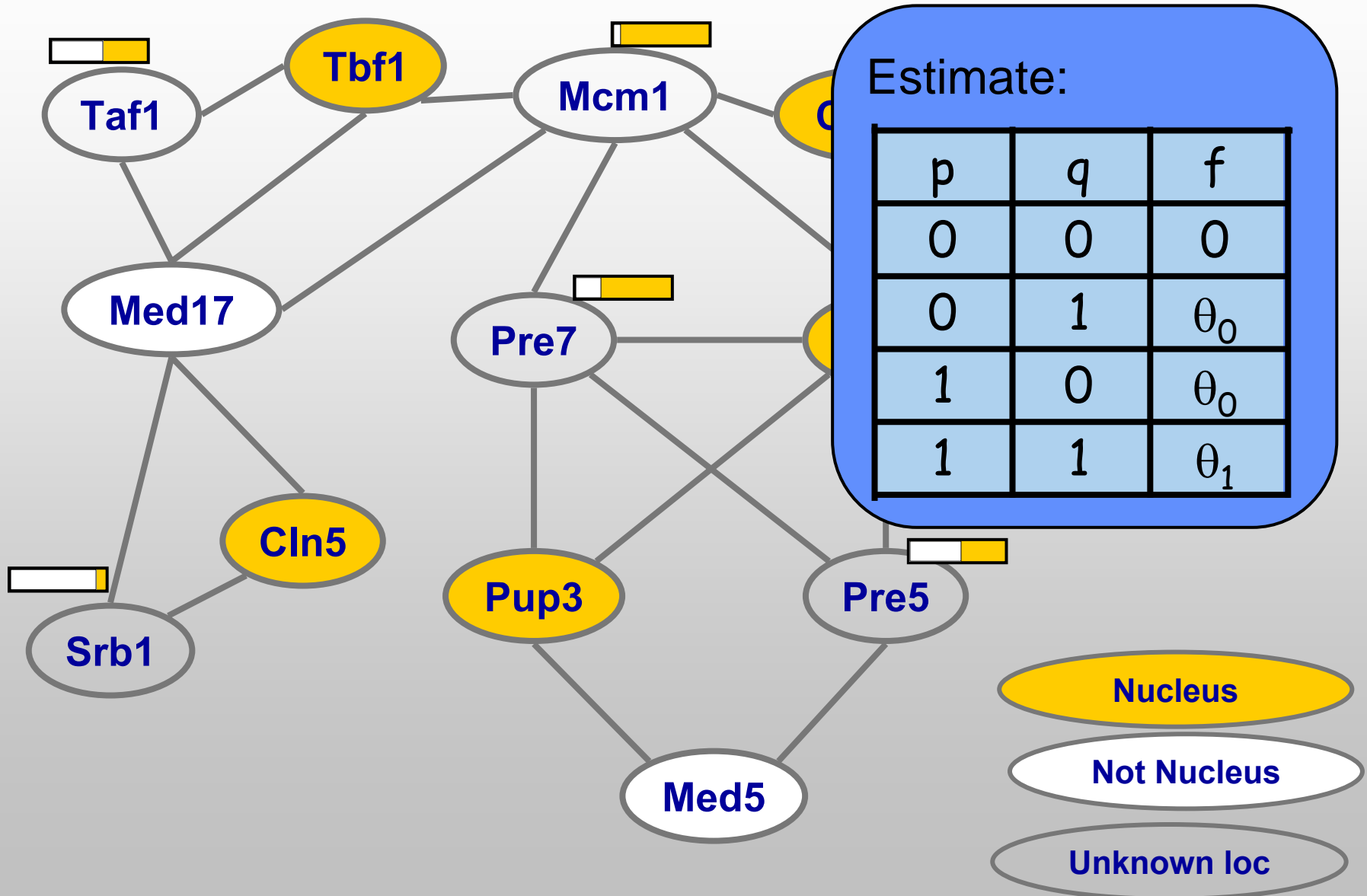
◆ Learning phase

- Use known proteins attributes to estimate potentials for each type of attribute

◆ Prediction phase

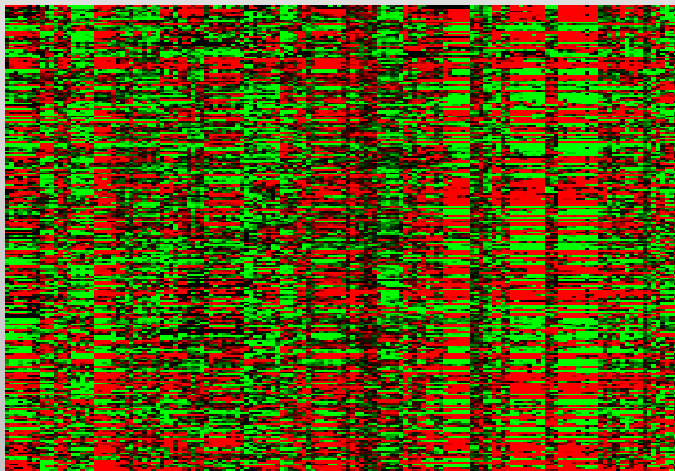
- Use inference to predict attributes for all proteins given evidence
- Simultaneous predictions for all the proteins in the network

Relational Markov Networks for Protein-Protein Interaction



Inferring “Pathways”

- ◆ **Assumption:** pathways exhibit two properties
 - Have similar expression profiles
 - Protein products more likely to interact
- ◆ Use both types of data to find pathways

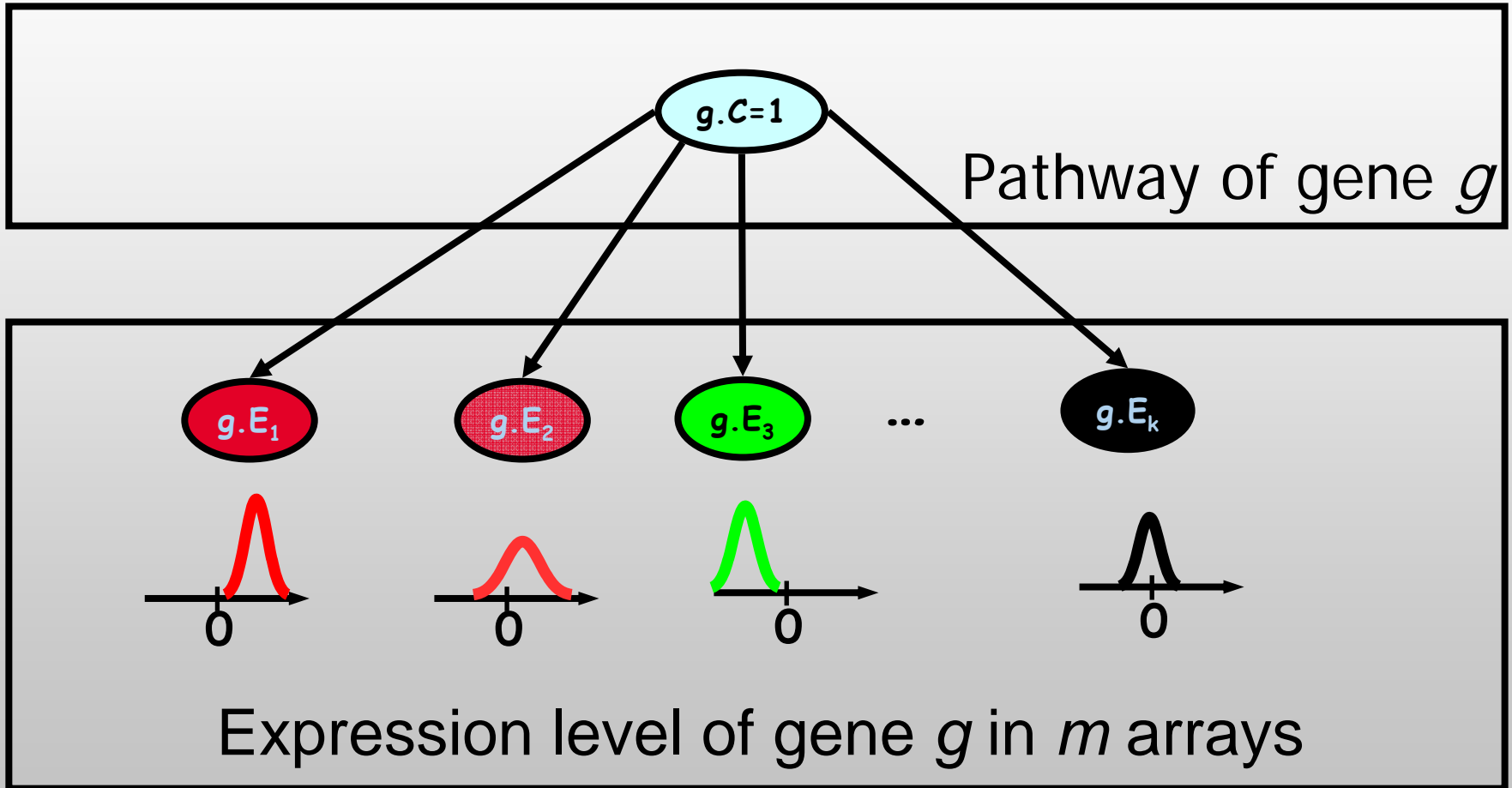


Probabilistic Model

- ◆ Genes are partitioned into “pathways”:
 - Every gene is assigned to one of ‘k’ pathways
 - Random variable for each gene with domain $\{1, \dots, k\}$
- ◆ Expression component:
 - Model likelihood is higher when genes in the same pathway have similar expression profiles
- ◆ Interaction component:
 - Model likelihood is higher when genes in the same pathway interact

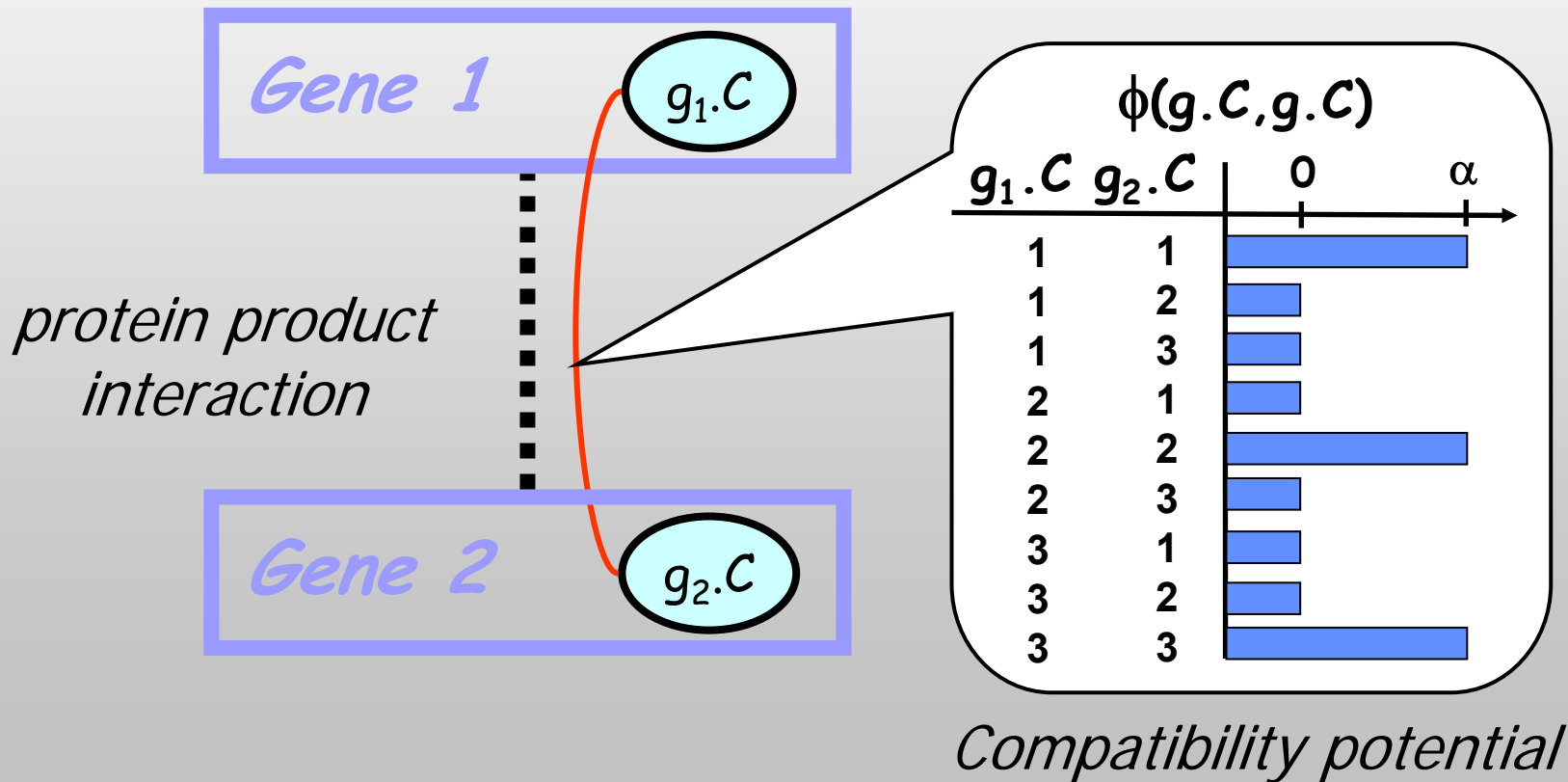
Expression Component

Naïve Bayes

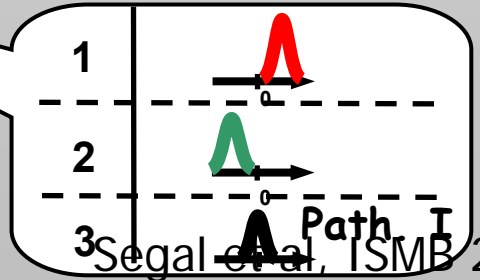
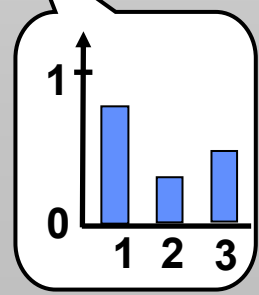
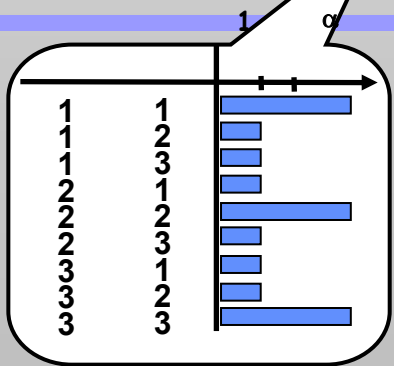
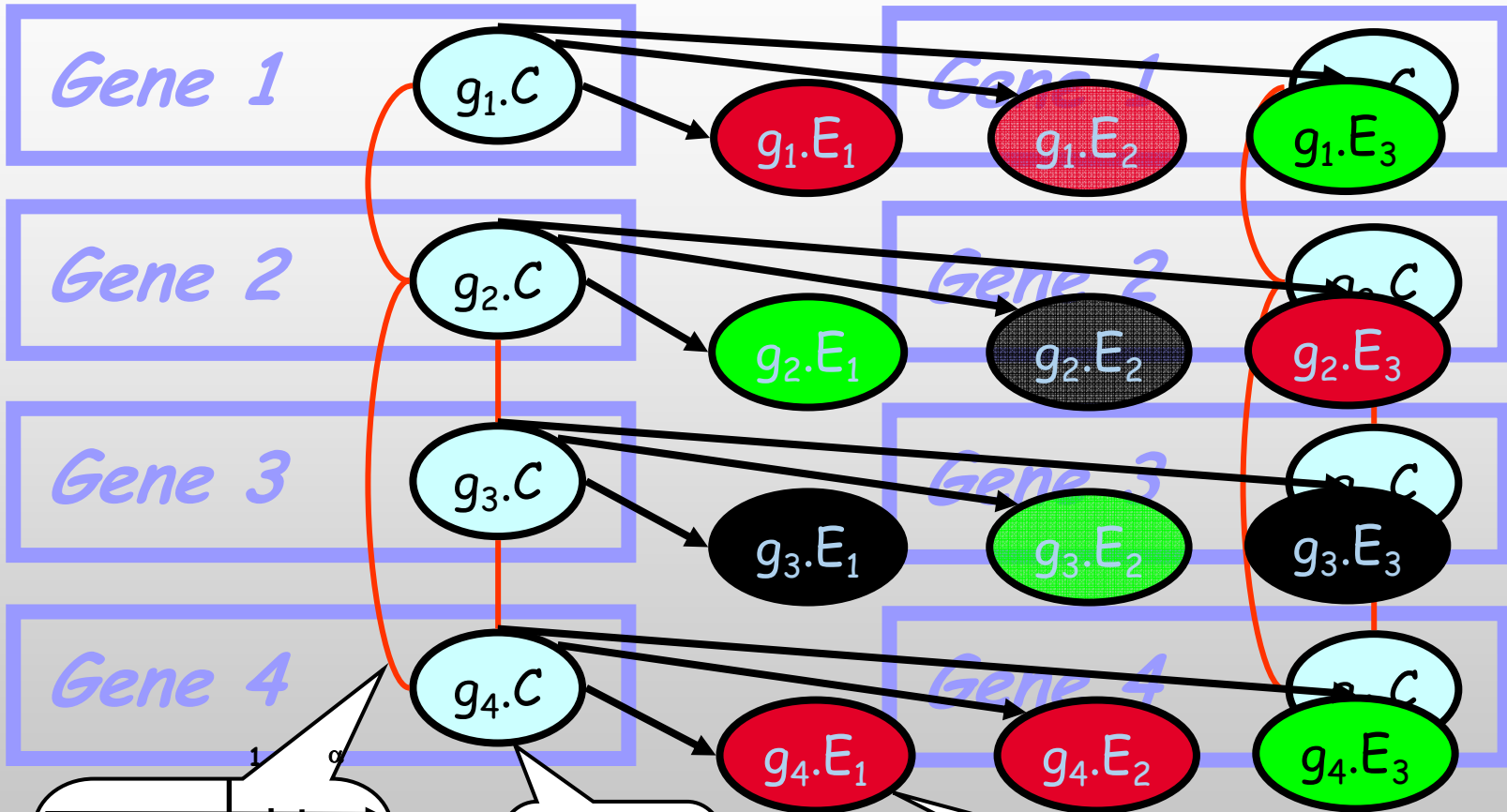


Protein Interaction Component

- ◆ Interacting genes are more likely to be in the same pathway

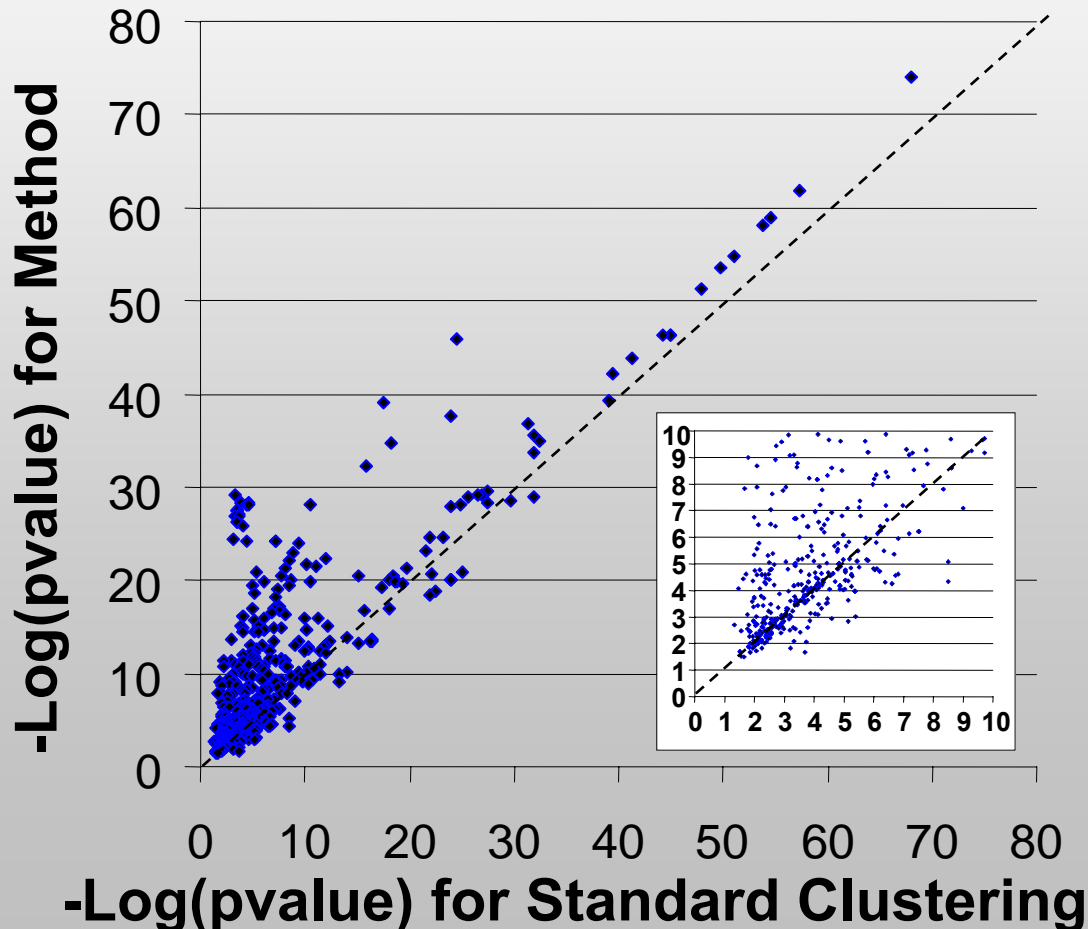


Joint Probabilistic Model

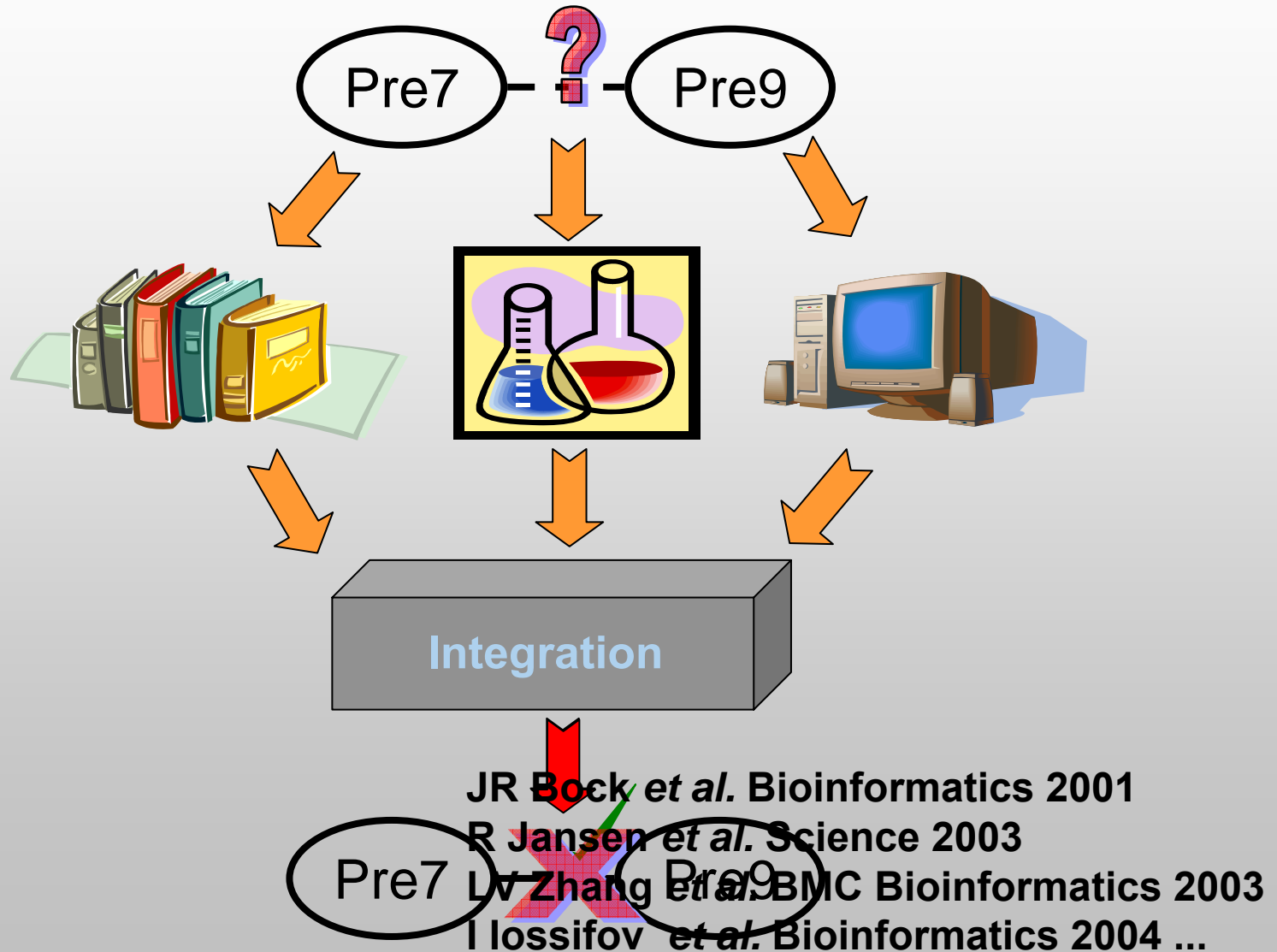


Comparison to Clustering

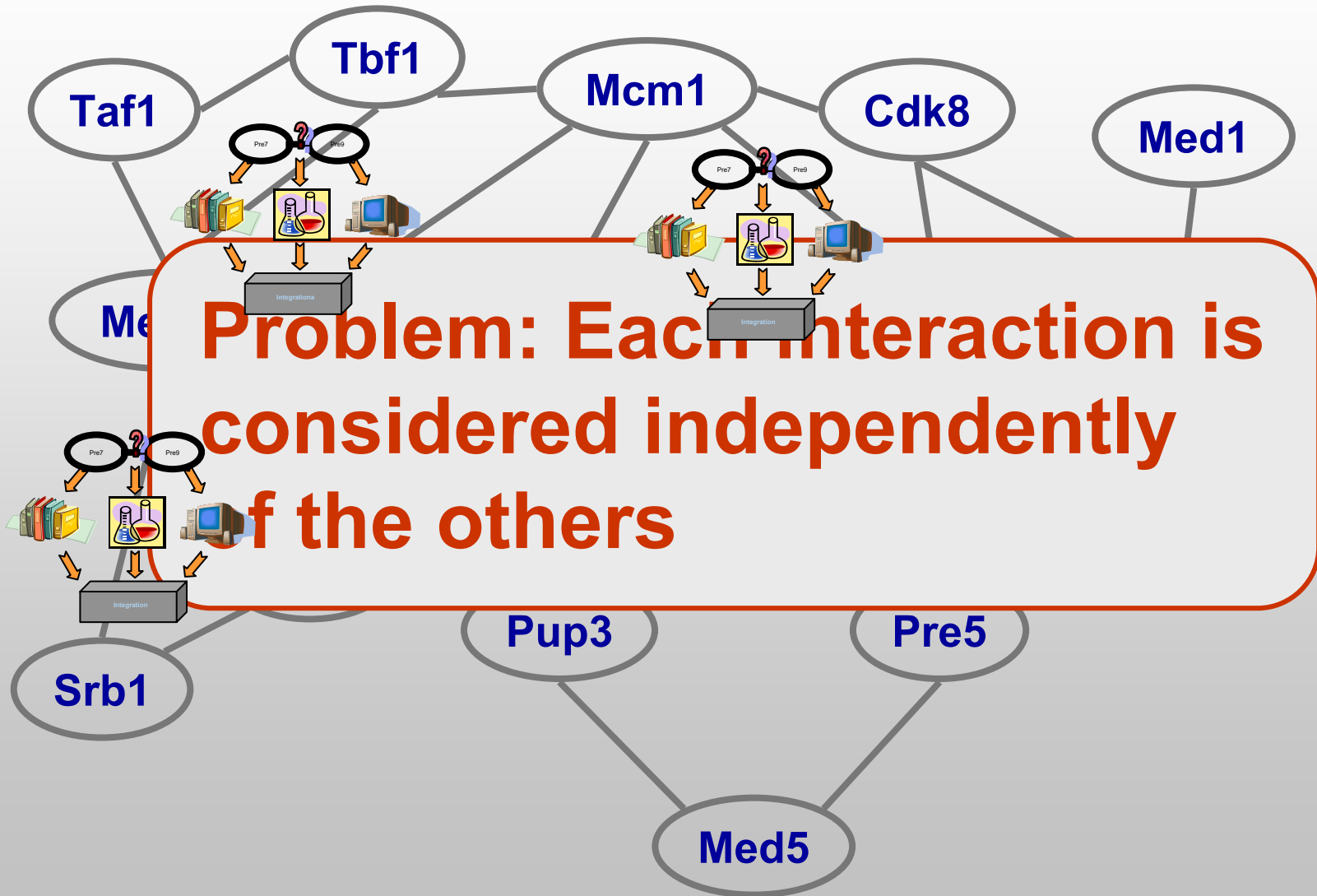
- ◆ Check enrichment of known gene annotations in pathways
- ◆ Calculate significance (negative log p-value)



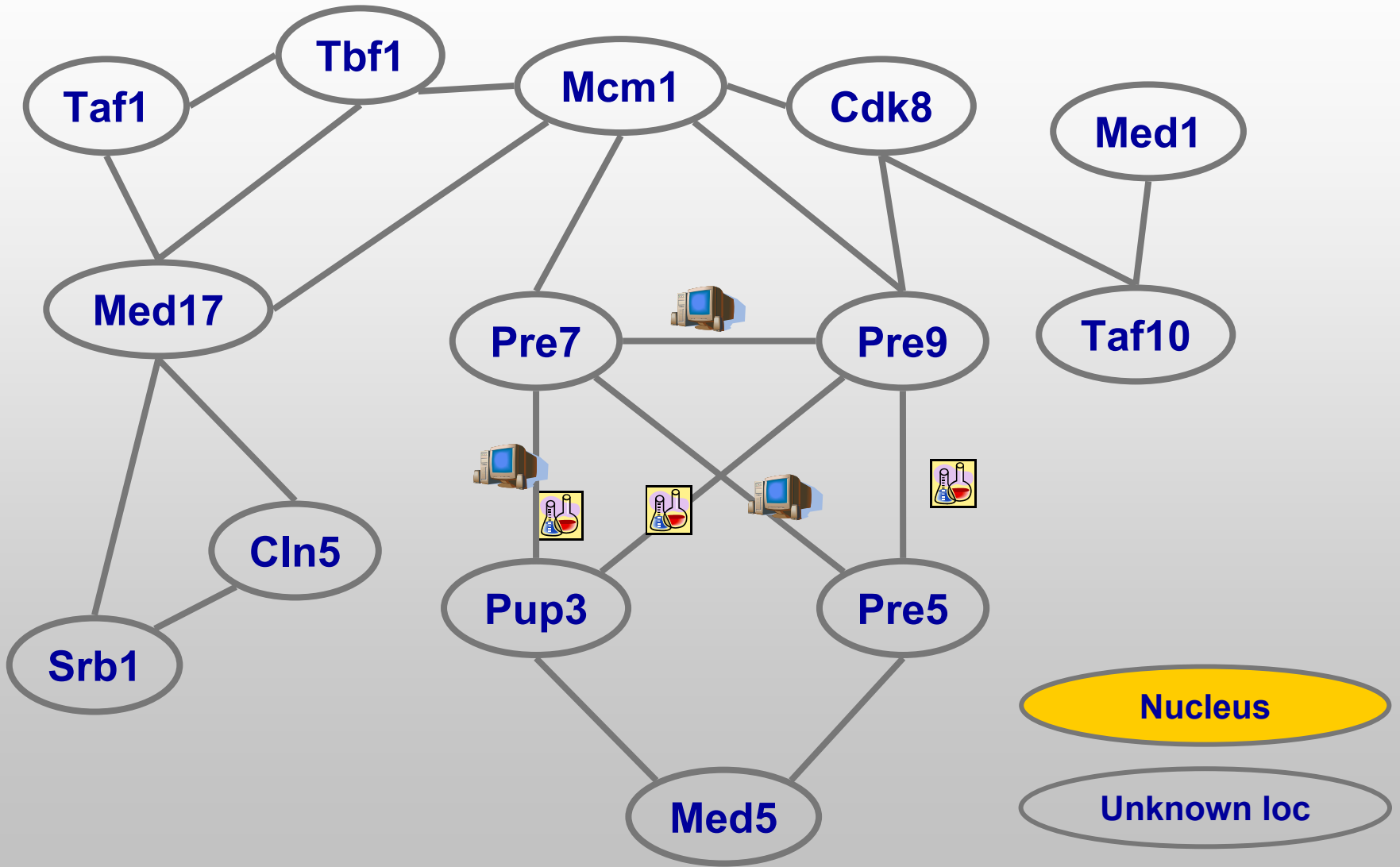
Predicting Protein-Protein Interactions



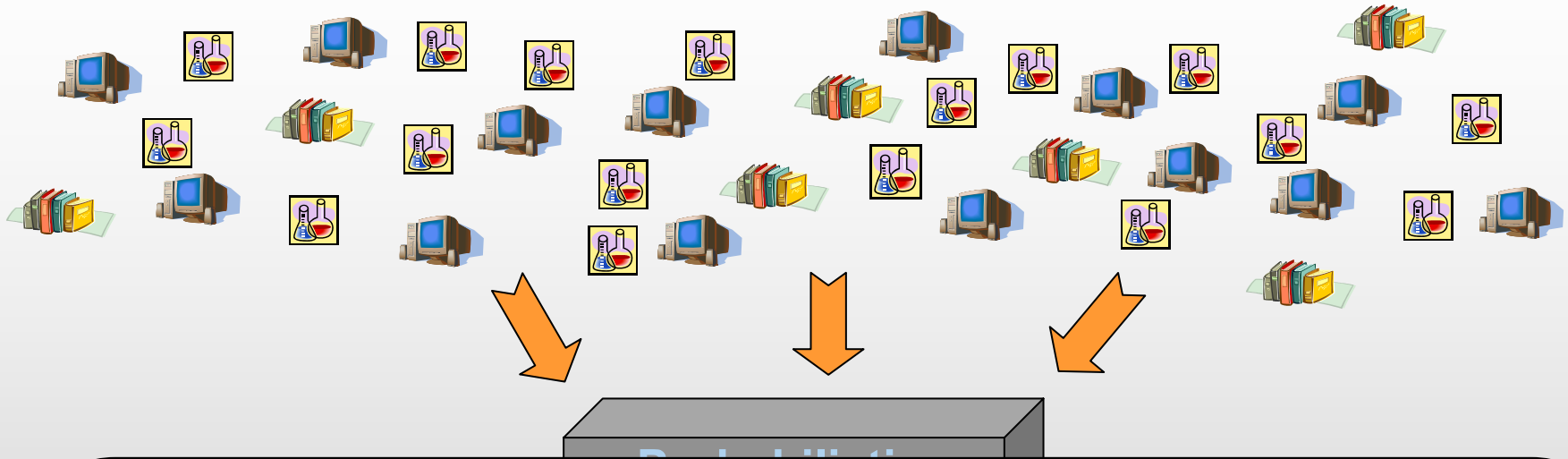
Predicting Interactions



Motivation



Design Plan



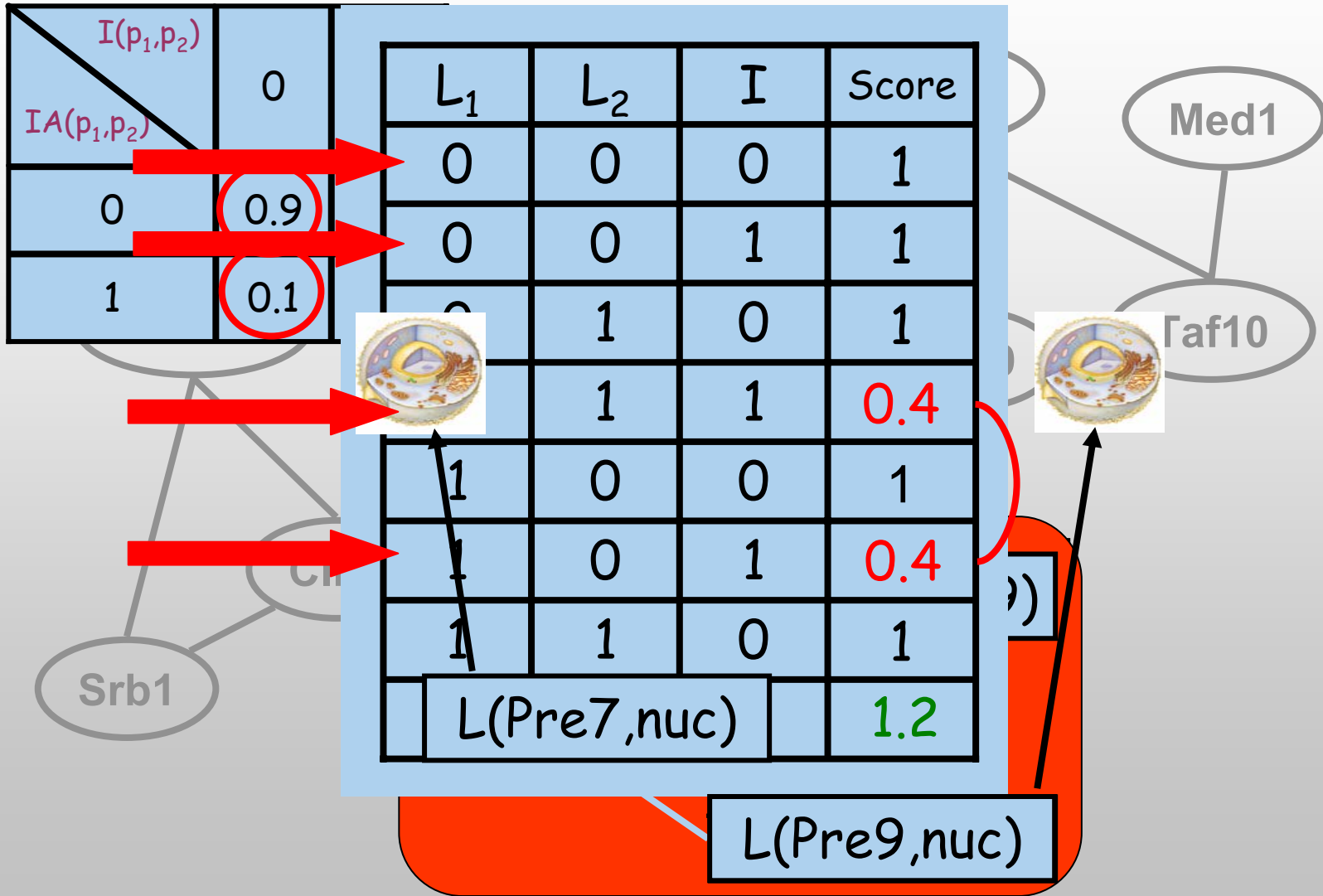
New variables denoting

- Whether two proteins interact
- Experimental observations about each interaction

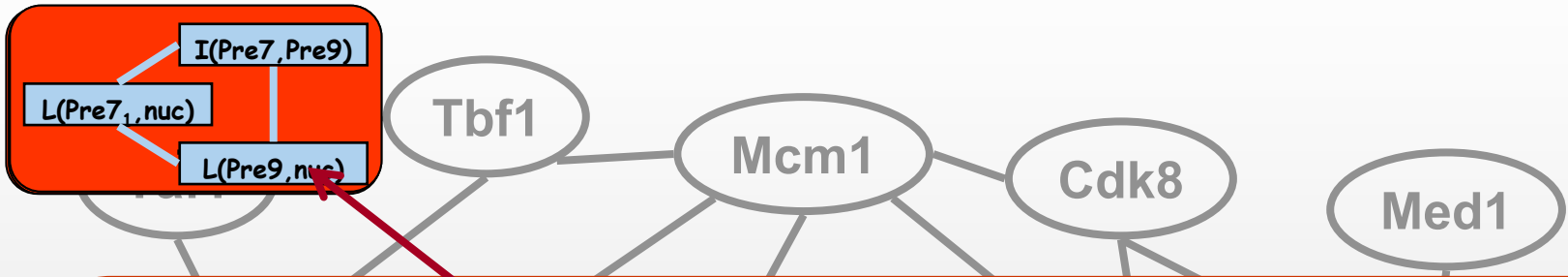
Main difficulty:

High connectivity between these variables

Building the Model

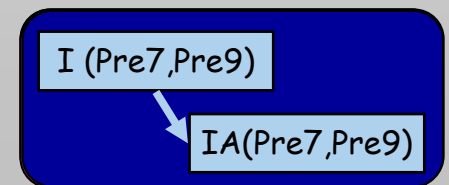
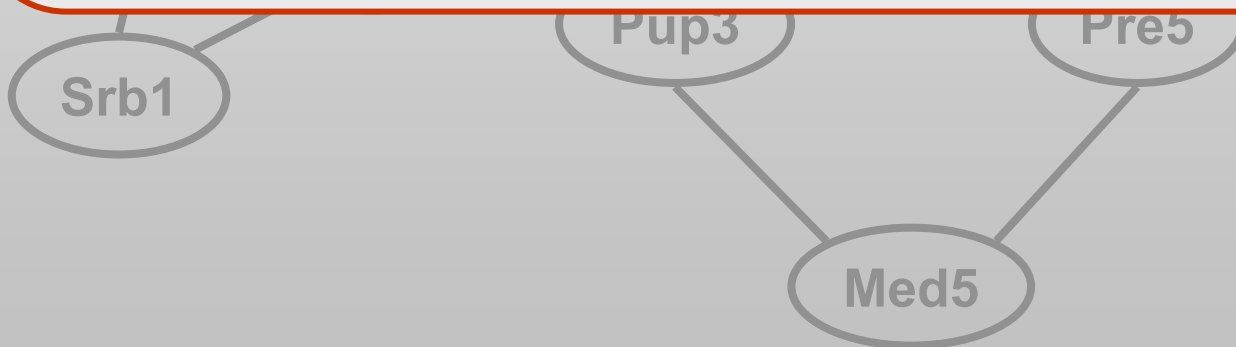


Using a Relational Model

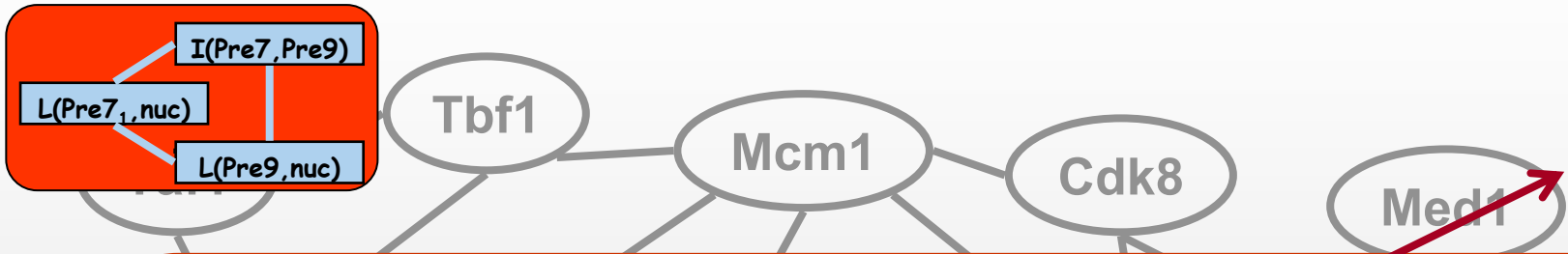


So far, equivalent to integrated prediction of each interaction independently

$$p(x) = \frac{1}{Z} \prod_{i \neq j} \left(\prod_l \psi(I_{i,j}, L_i^l, L_j^l) \prod_a p(I A_{i,j}^a | I_{i,j}) \right)$$

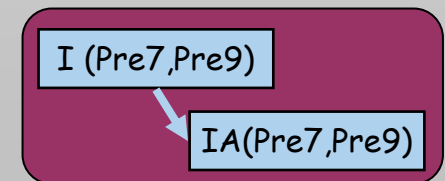
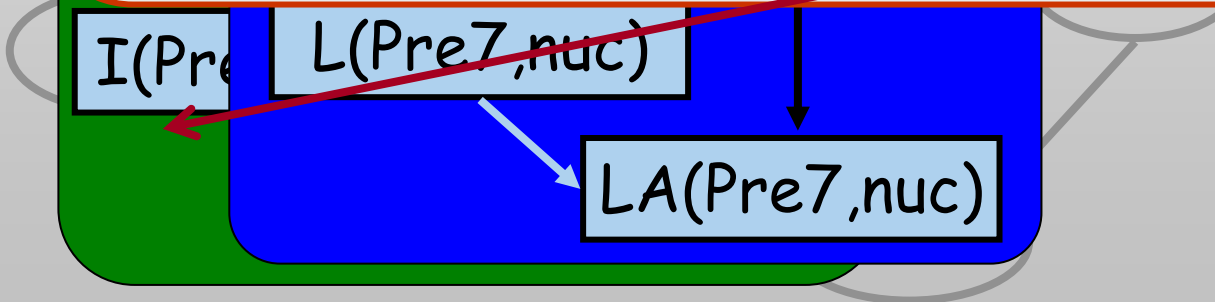


Building the Model

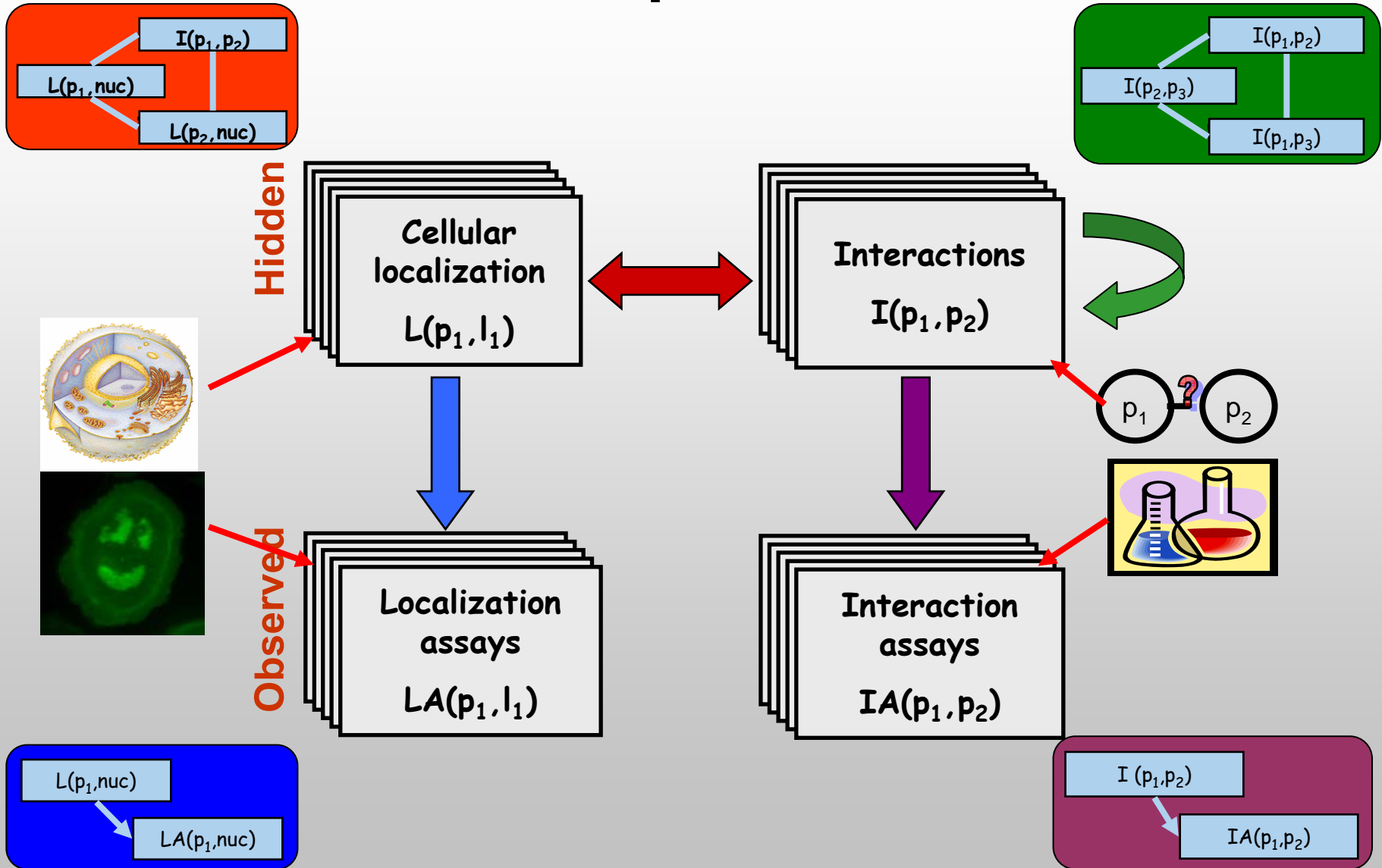


$$p(x) = \frac{1}{Z} \prod_{i \neq j} \prod_l \psi(I_{i,j}, L_i^l, L_j^l) \prod_{i \neq j} \prod_a p(IA_{i,j}^a | I_{i,j}^a)$$

$$\prod_{i \neq j \neq k} \psi(I_{i,j}, I_{j,k}, I_{i,k}) \prod_i \prod_l p(LA_i^l | L_i^l)$$



The Complete Model



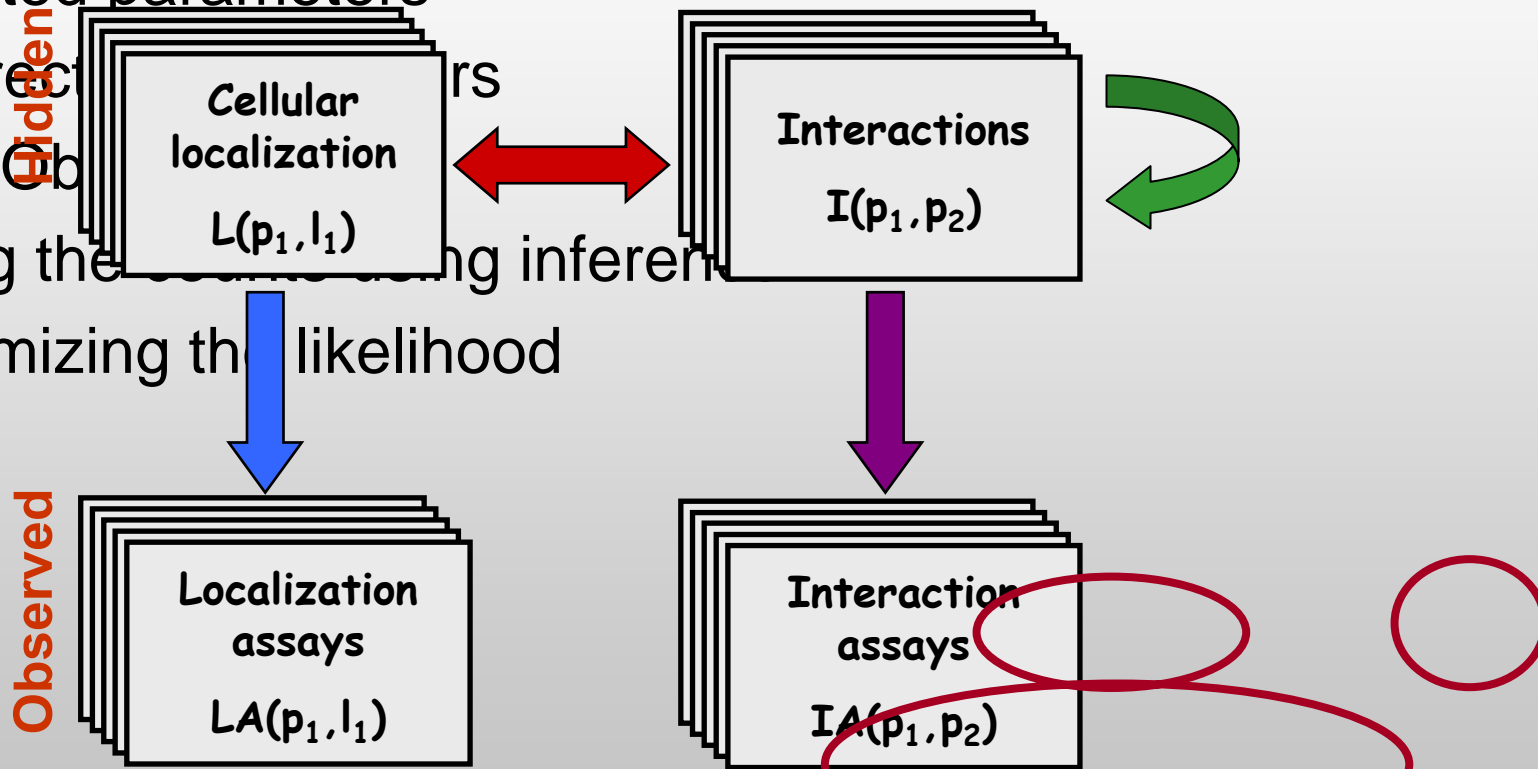
Learning the Parameters

◆ Maximizing the likelihood (fully observed case)

- Directed parameters
- Undirected parameters

◆ Partially Observed

- Filling the missing information
- Maximizing the likelihood



$$p(x) = \frac{1}{Z_x} \prod_c \psi_c(x_c) \prod_i p(y_i | Pa_i)$$

Model Evaluation: *S.cerevisiae*

Large scale data:

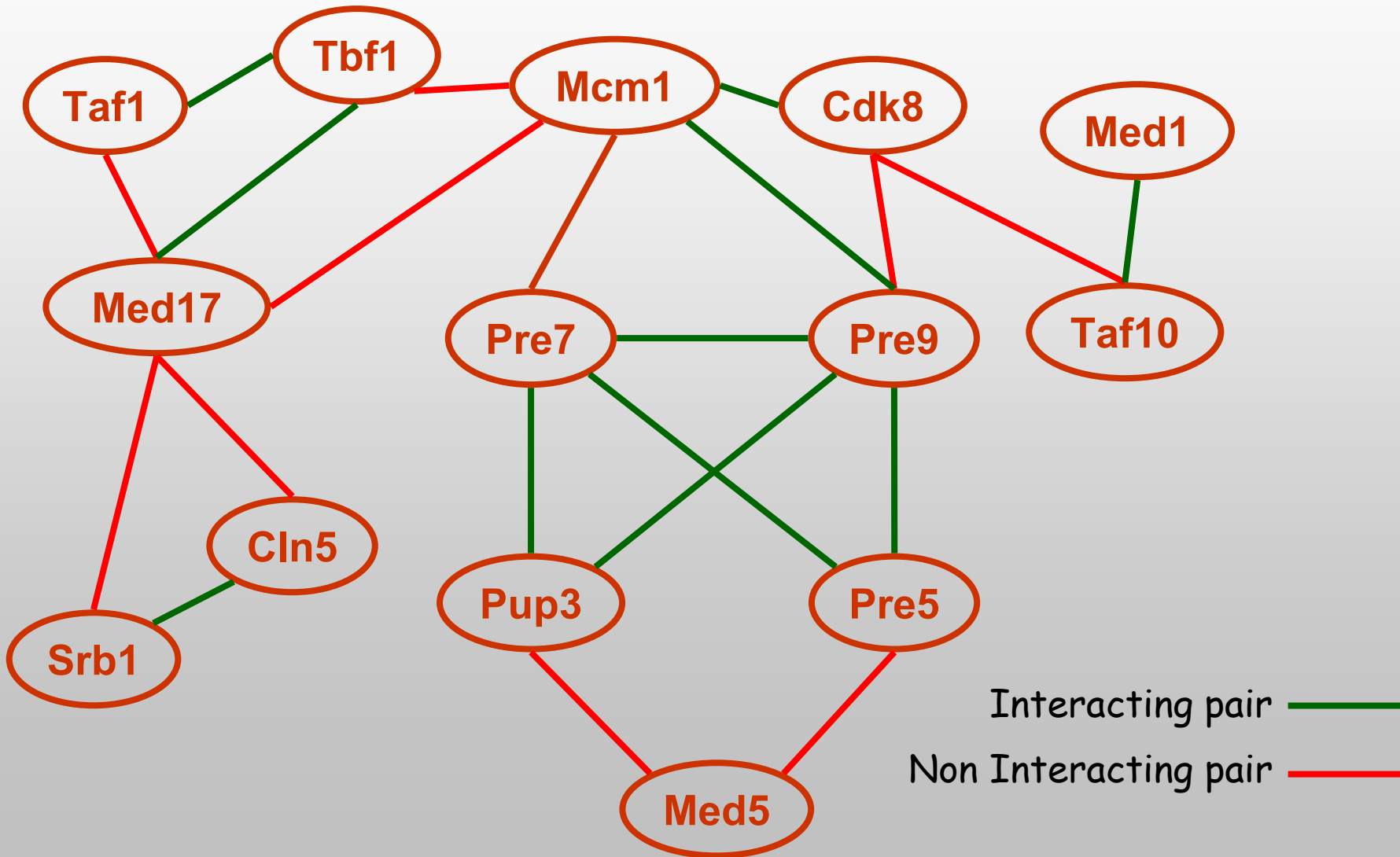
- Yeast two hybrid (Ito *et al.* + Uetz *et al.*)
- Complexes (MIPS)
- Correlated domain signatures (Sprinzak *et al.*)
- Protein localization (Huh *et al.*)

38,000 potentials

37 free parameters

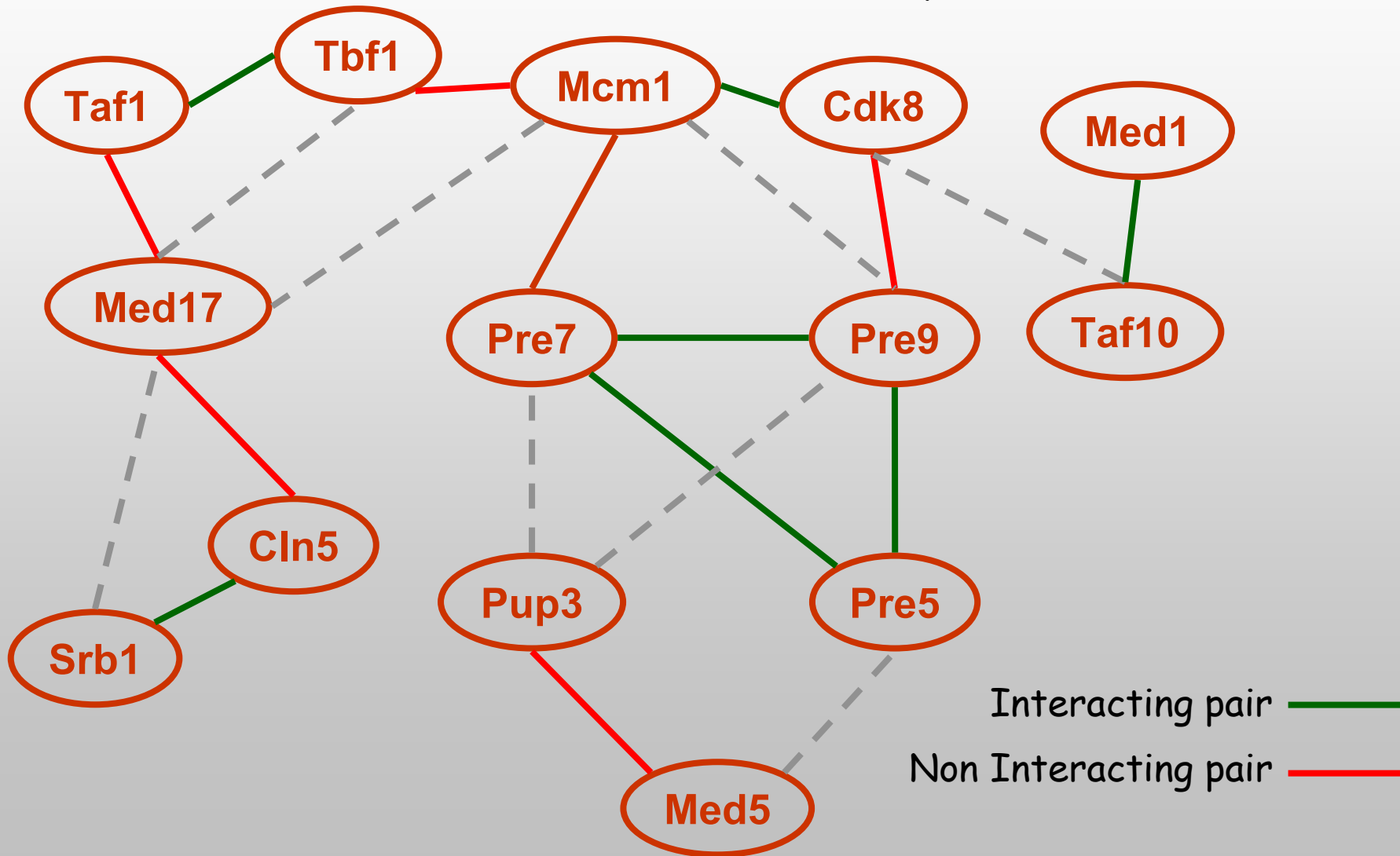
Med5

Evaluation: Cross Validation



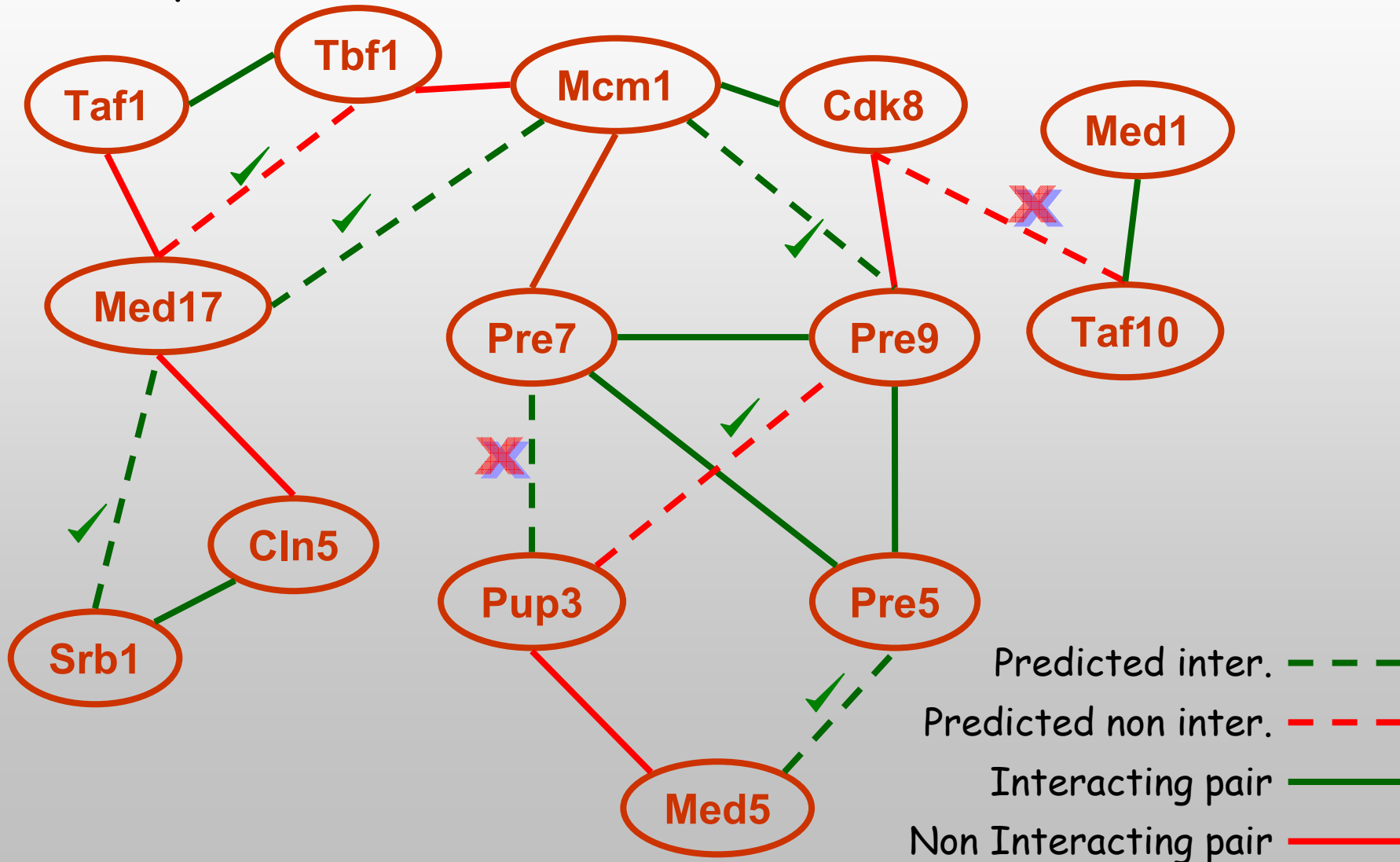
Evaluation: Parameter Estimation

Hide a set of test interactions and learn parameters

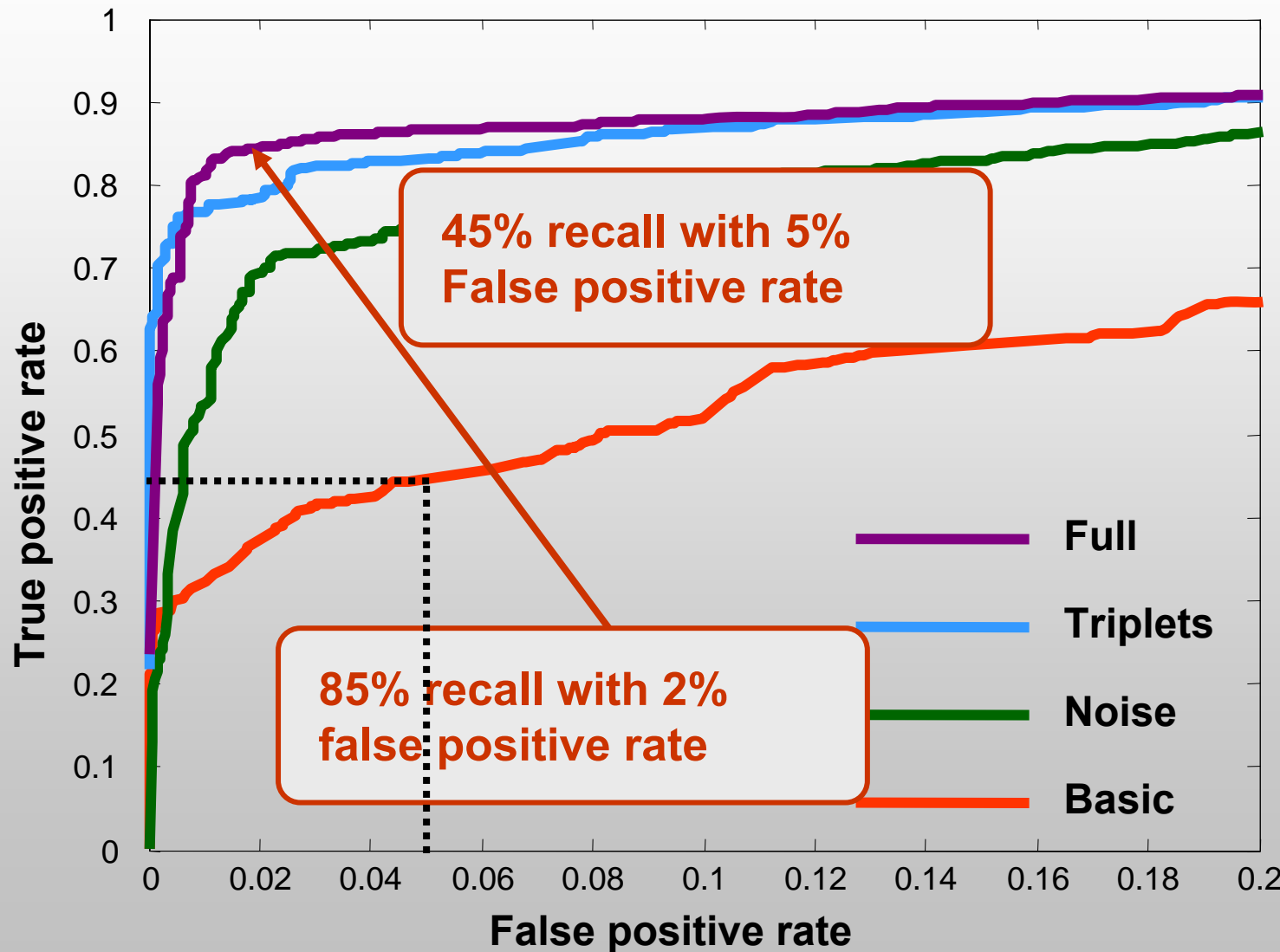
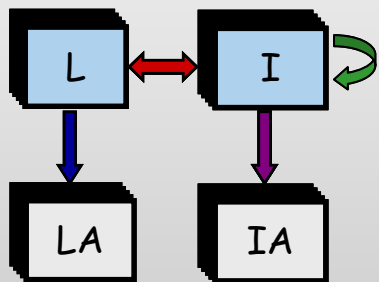


Evaluation: Validate Predictions

Use known prediction parameters to predict hidden interactions



Evaluation: ROC curve



Outline

- ◆ Introduction
- ◆ Bayesian Networks
- ◆ Learning Bayesian Networks
- ◆ Transcriptional regulation
- ◆ Gene expression
- ◆ Markov Networks
- ◆ Protein-Protein Interactions
- ◆ Discussion

Philosophy

Identify biological problem



Formulate model



Learning/Inference procedures



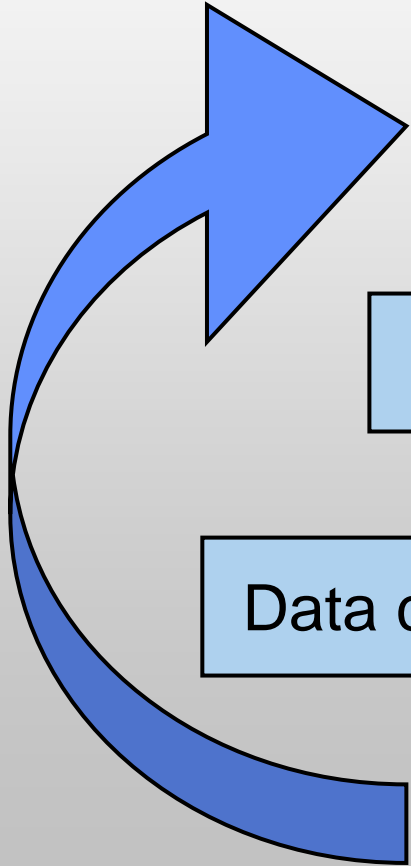
Data curation, apply procedure, get results



Biological interpretation



New discoveries



Recap

- ◆ Models of evolution
 - Pedigree analysis
 - Sequence evolution
- ◆ Transcription Factors
 - Binding sites
- ◆ Gene Expression
 - Clustering, interaction networks
- ◆ Protein-Protein interaction networks
- ◆ Combination of subsets of these

Additional Areas

◆ Gene finding

- Extended HMMs + evolutionary models

◆ Analysis of genetic variation

- SNPs, haplotypes, and recombination

◆ 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

The END