

# Data, Knowledge, Modules and Models

Adam Arkin

Howard Hughes Medical Institute

Departments of Bioengineering and Chemistry

University of California

Physical Biosciences Division

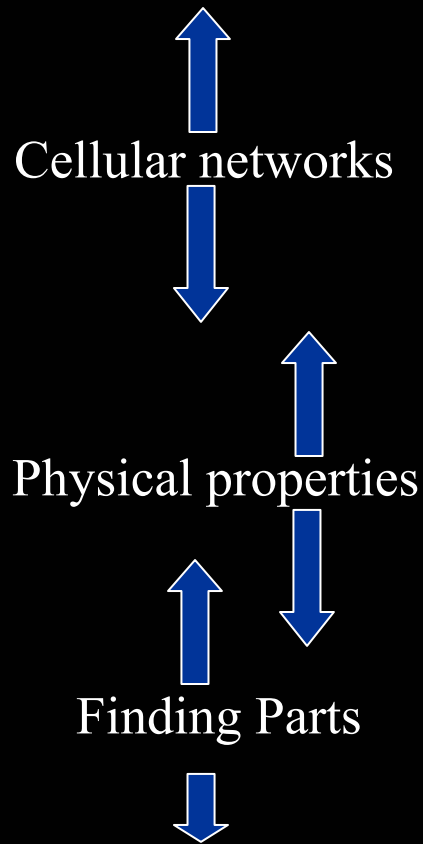
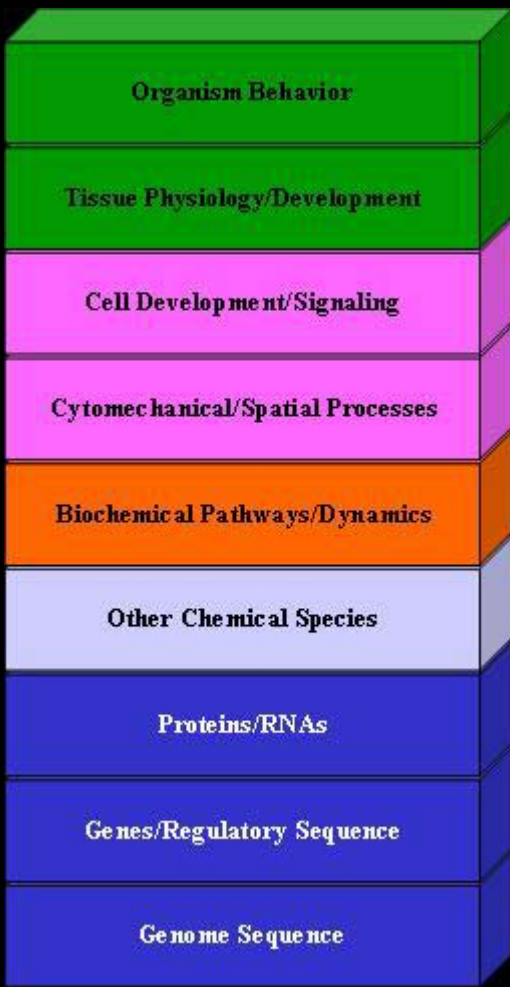
E.O. Lawrence Berkeley National Laboratory

Berkeley, CA 94720

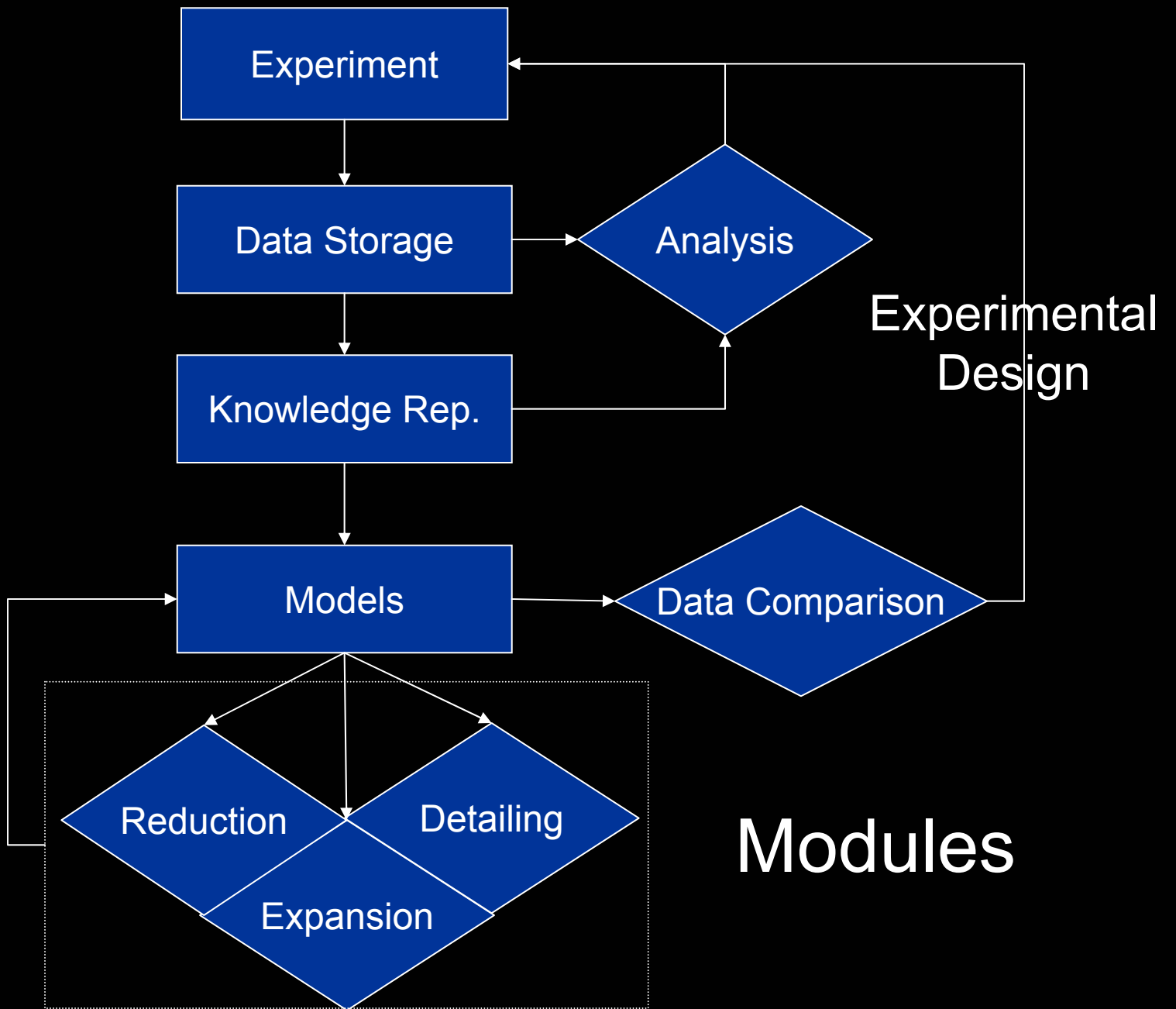
<http://genomics.lbl.gov>

<http://www.grc.uri.edu/programs/2001/bioinf.htm>

# Tools for “multilevel” analysis



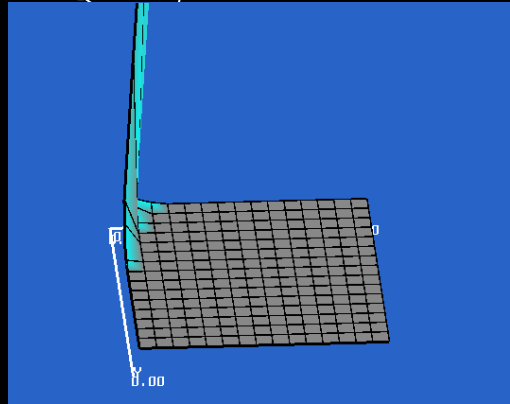
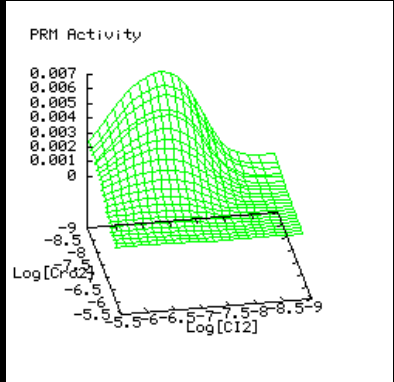
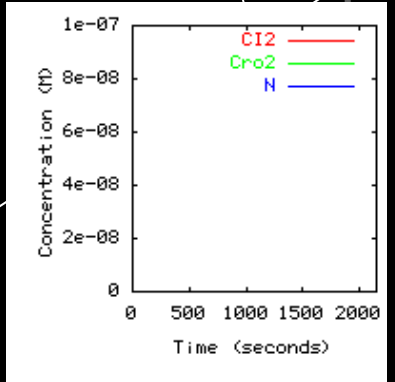
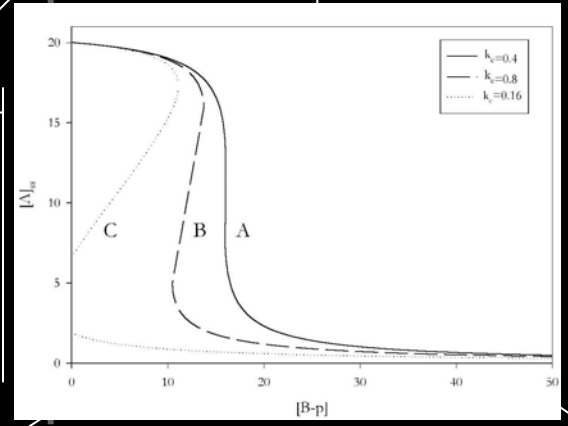
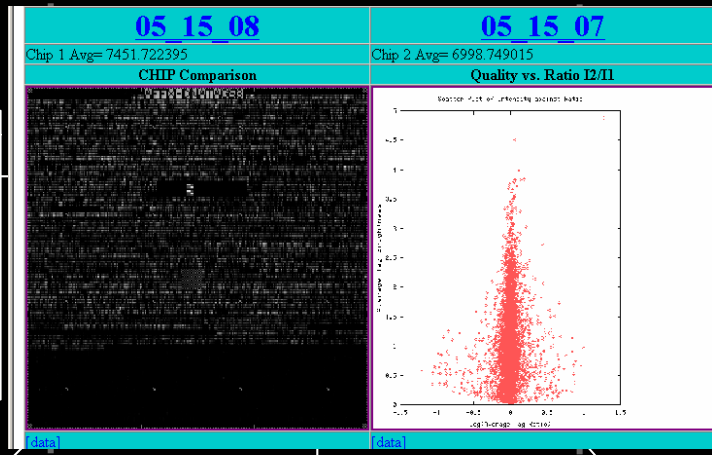
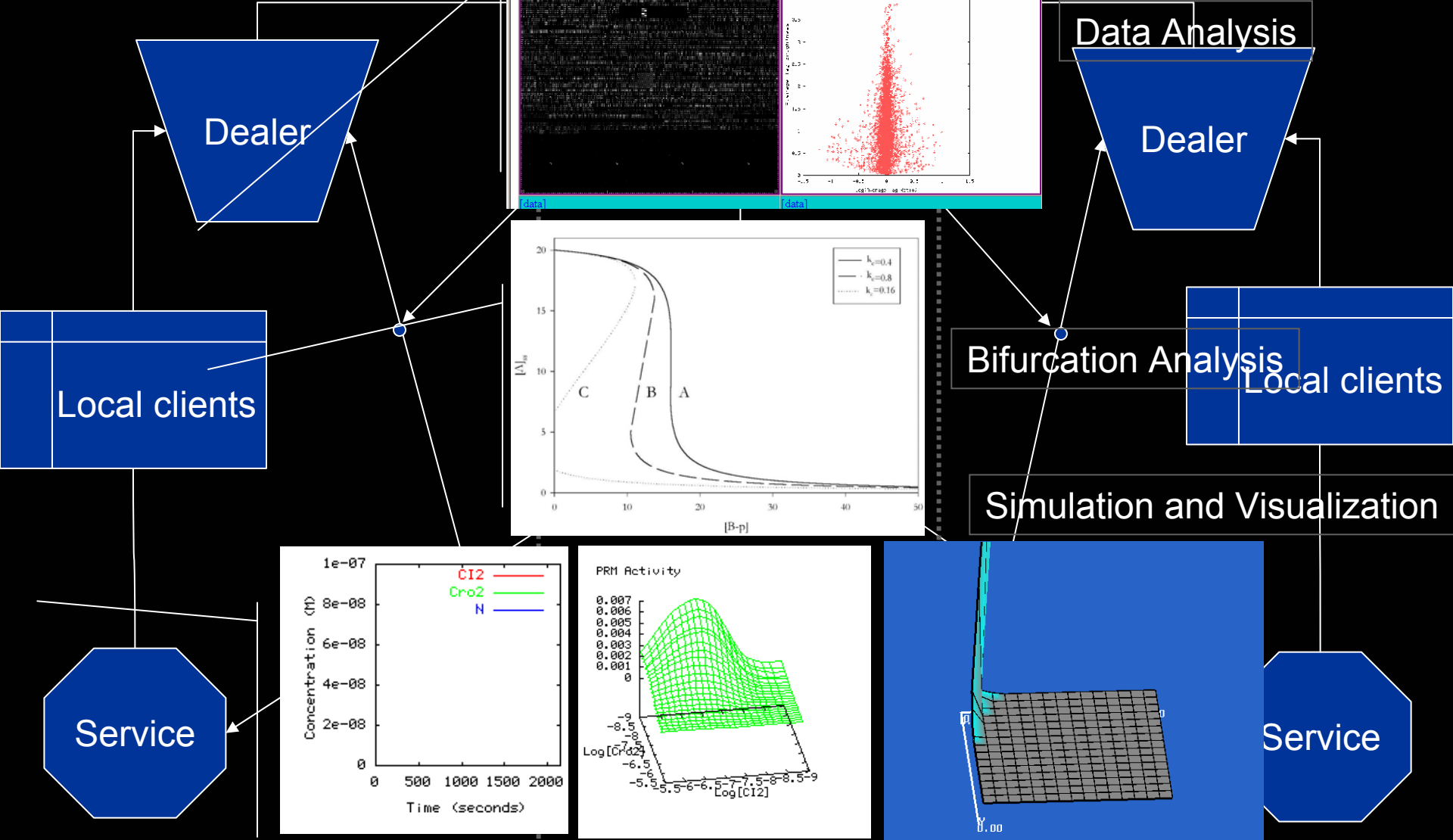
Epidemiological/Ecological Models		Cancer Dynamics	Multi-organism function: e.g. Infectious disease
Tissue Mechanics		Cell Behavior & Engineering	Organismal Behavior
Cytomechanical Analysis	Morphogenesis & Development	Homeostasis	Cell-Cell Interactions
Metabolic/Biosynthetic Analysis & Engineering	Signal Transduction Analysis	Gene expression/network Analysis	
Biochemical and Genetic Network Prediction			
Molecular Interaction Prediction	Chromatin Structure	Macromolecular Dynamics	
Protein 3° Struct	Protein Function ID	RNA Function ID	
Protein Sequence ID	Homology Modeling	RNA 3° Struct	
mRNA Regulation	mRNA Splicing	RNA 2° Struct	
ORF Identification	DNA Regulatory ID	RNA Gene ID	
Assembled Genomes		Polymorphisms	

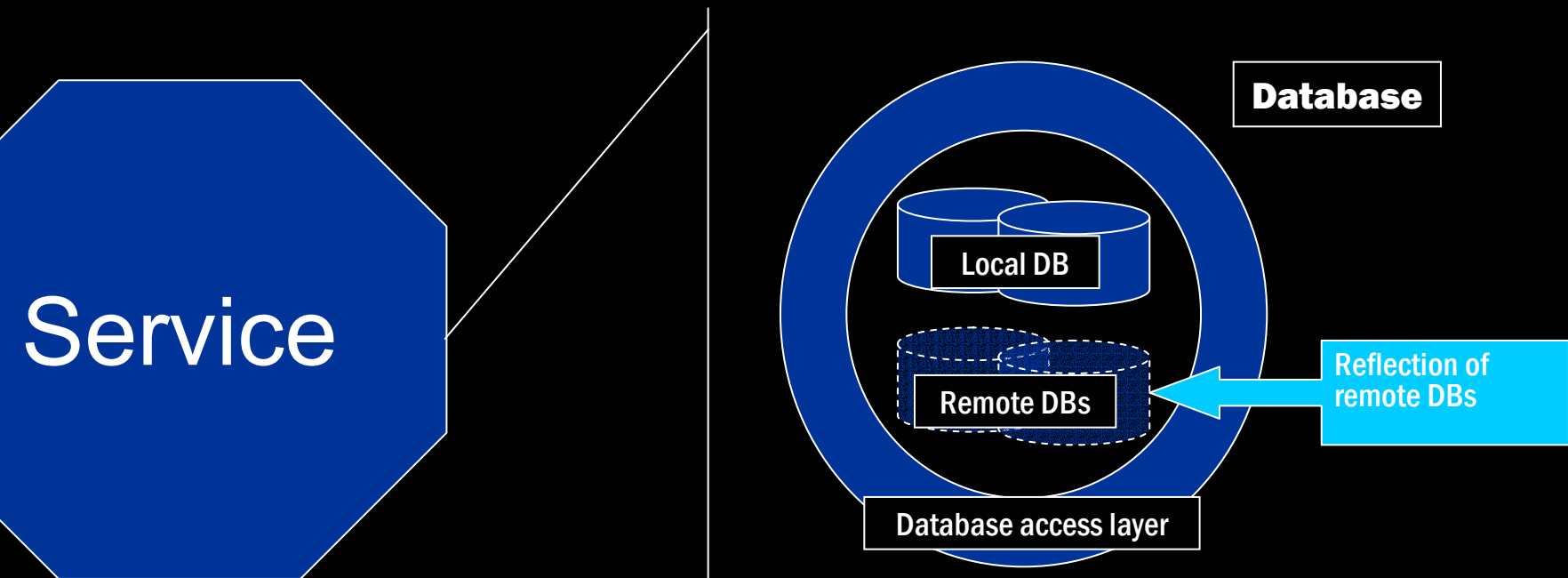


Building a problem solving environment



# Many client services



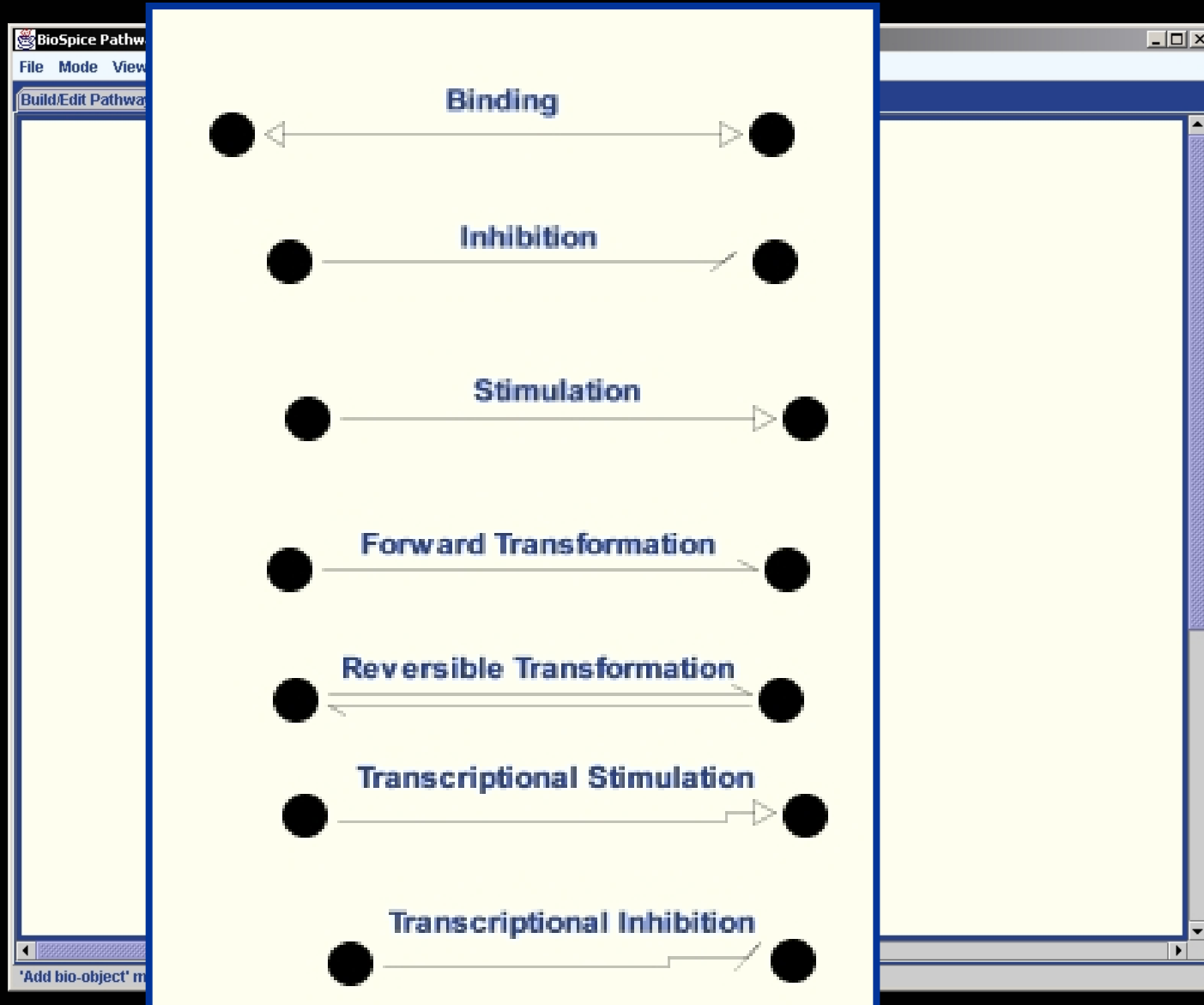


Schema Compliant with  
NCBI/BIND  
GATC/MGED  
GENBANK/PDB  
AND Glue for Models  
Not as nice as Shankar's





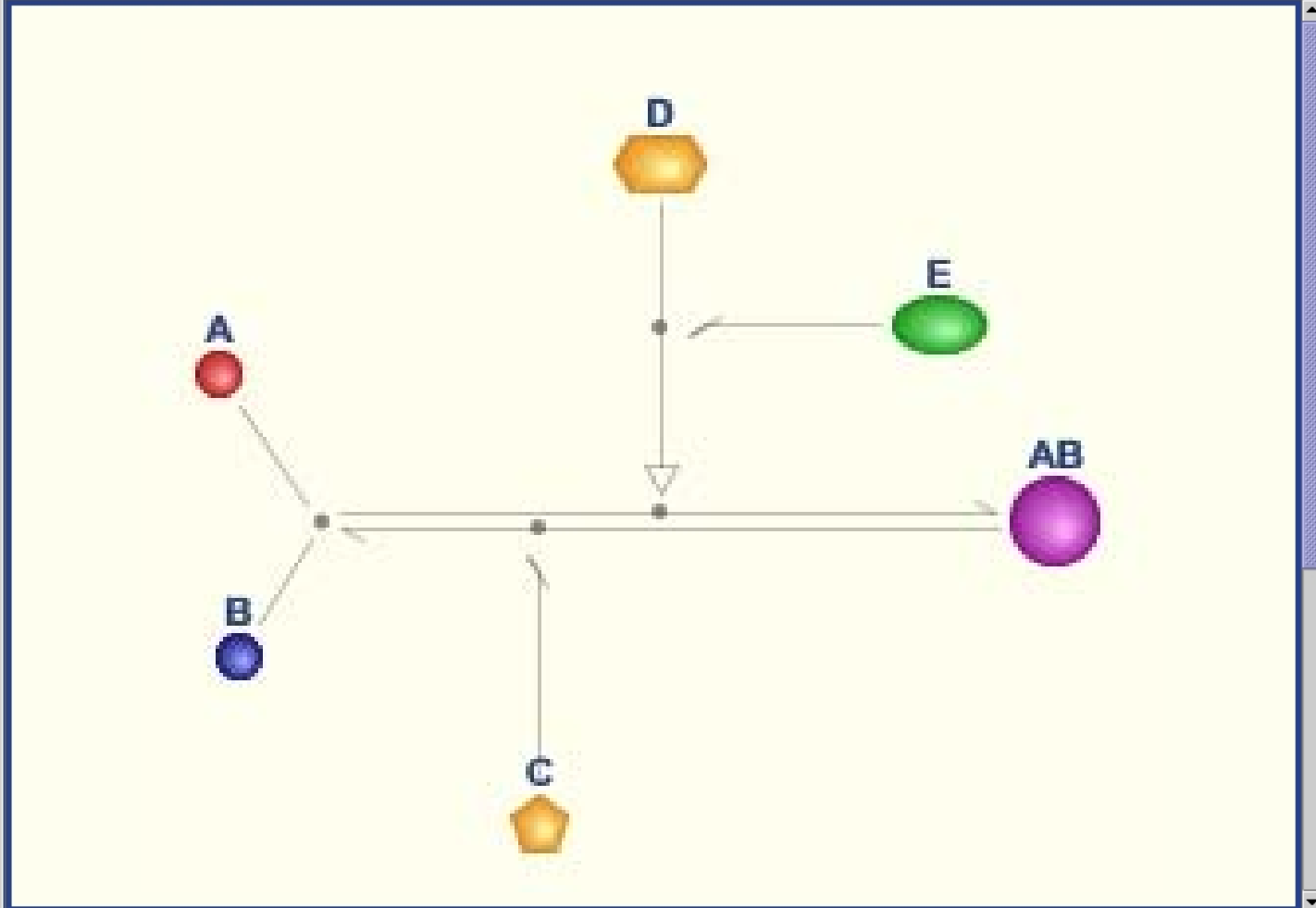
GUI must represent biological models at different levels of abstraction.



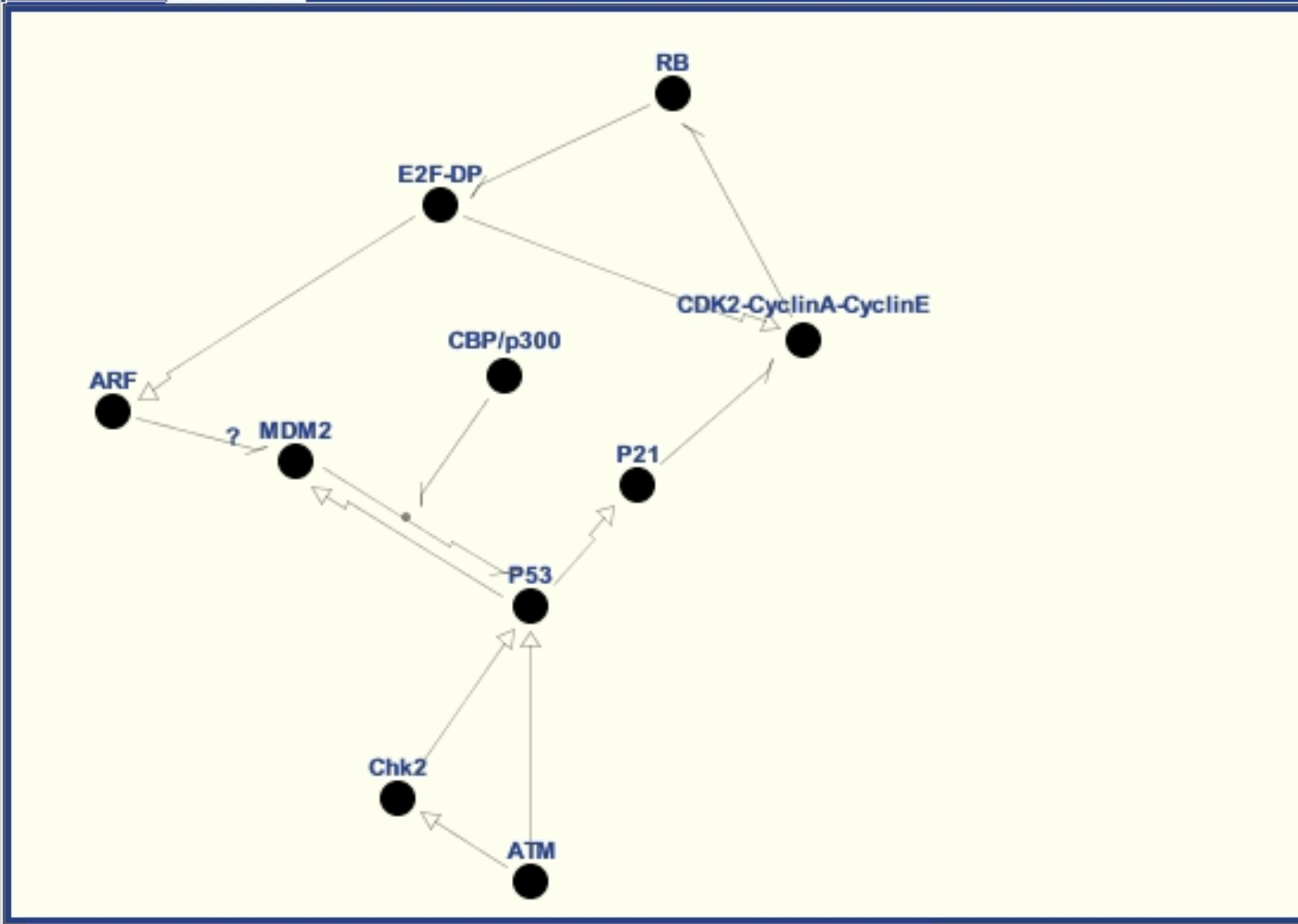
BioSpice Pathway Editor

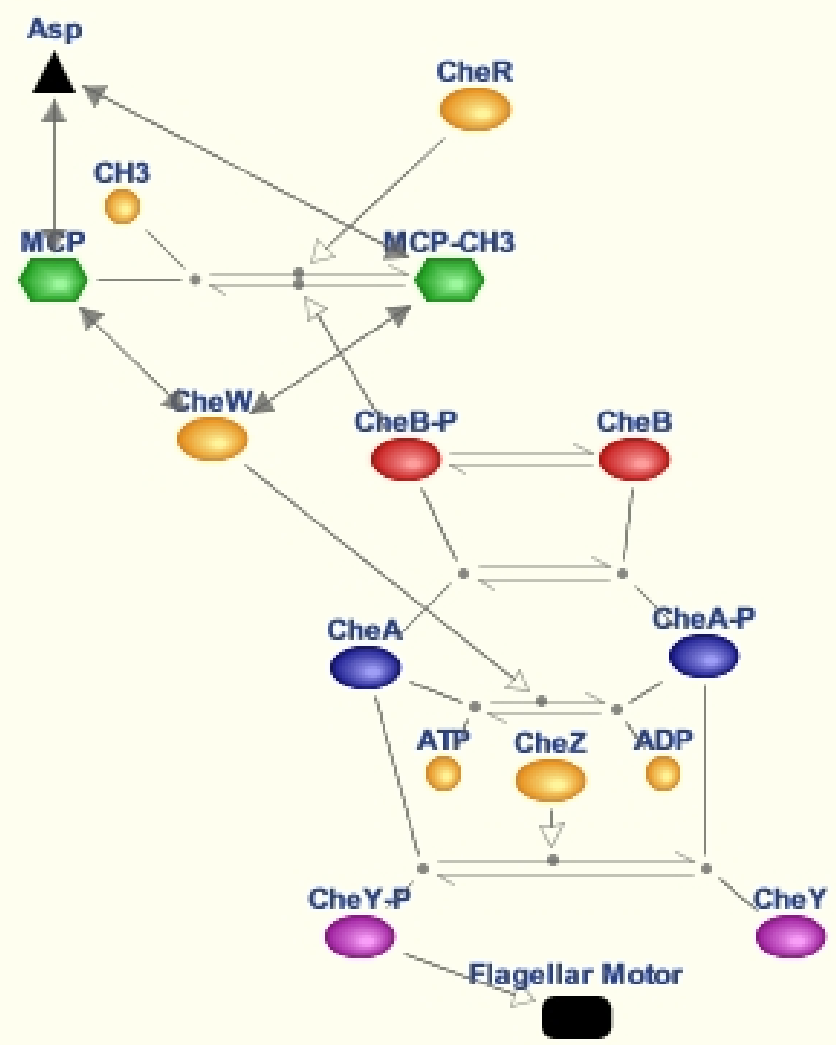
File Mode View Info

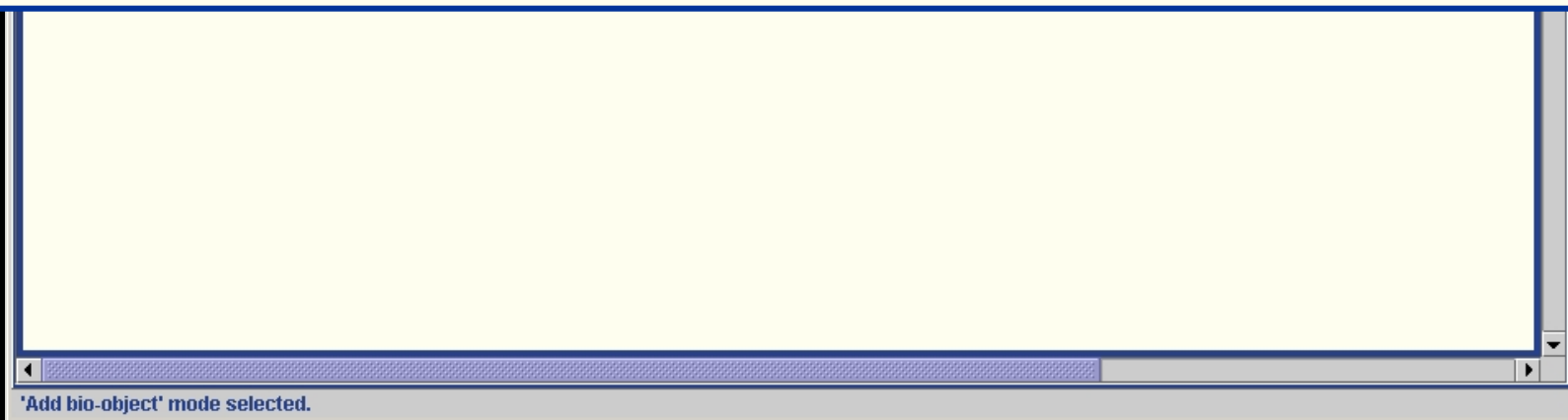
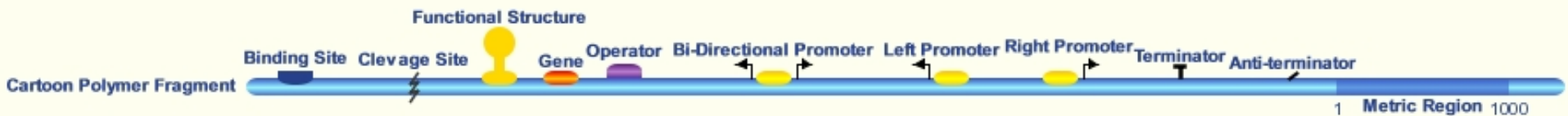
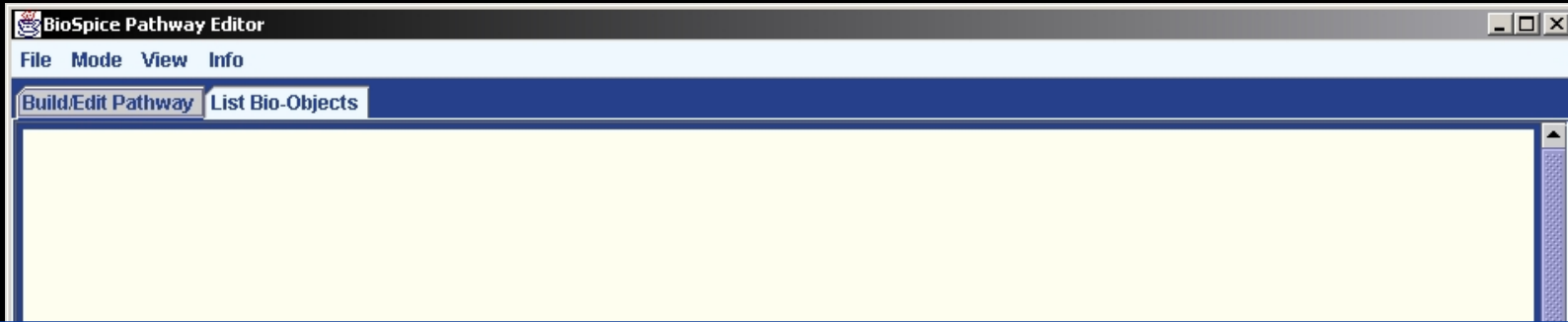
Build/Edit Pathway List Bio-Objects

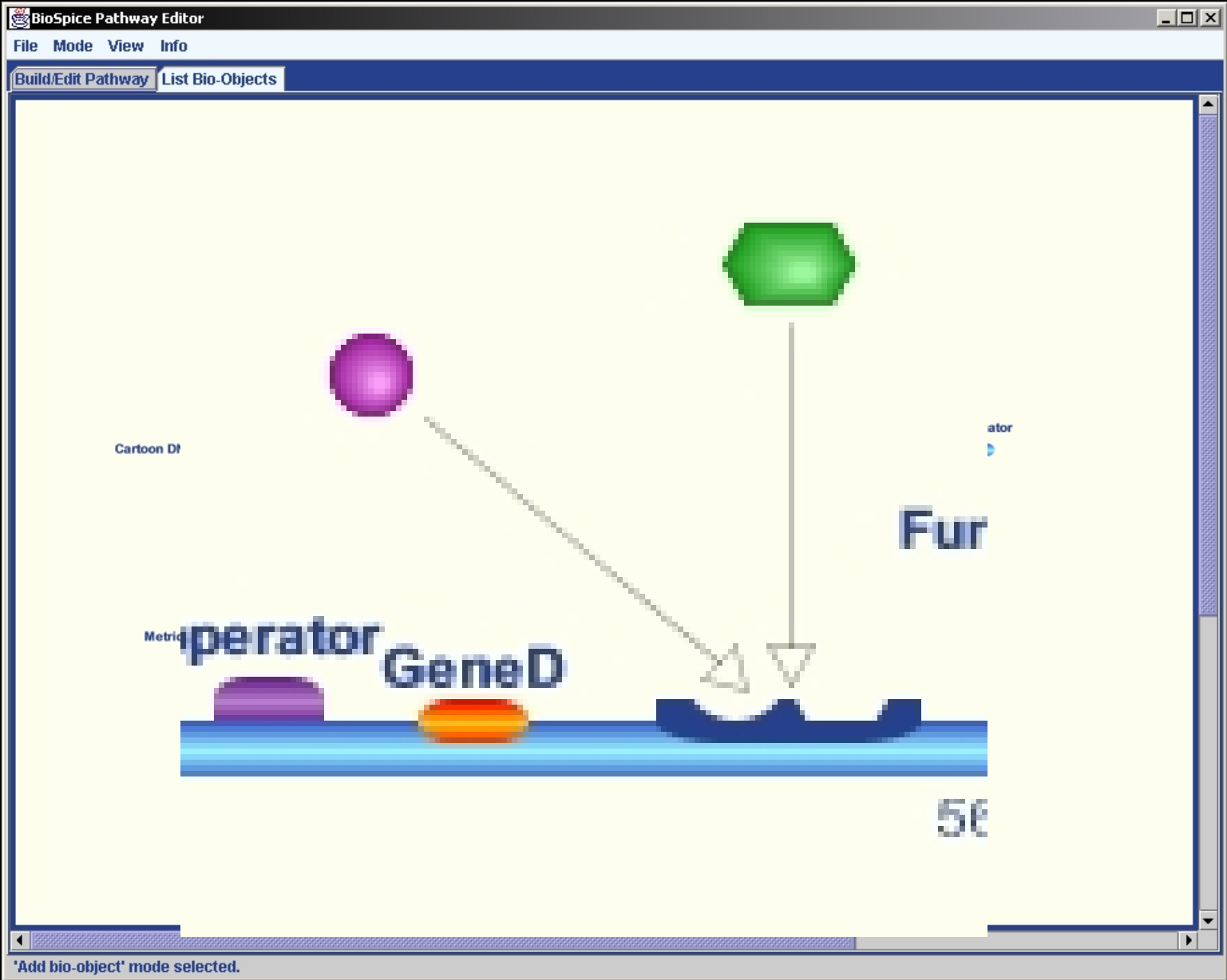


'Add bio-object' mode selected.







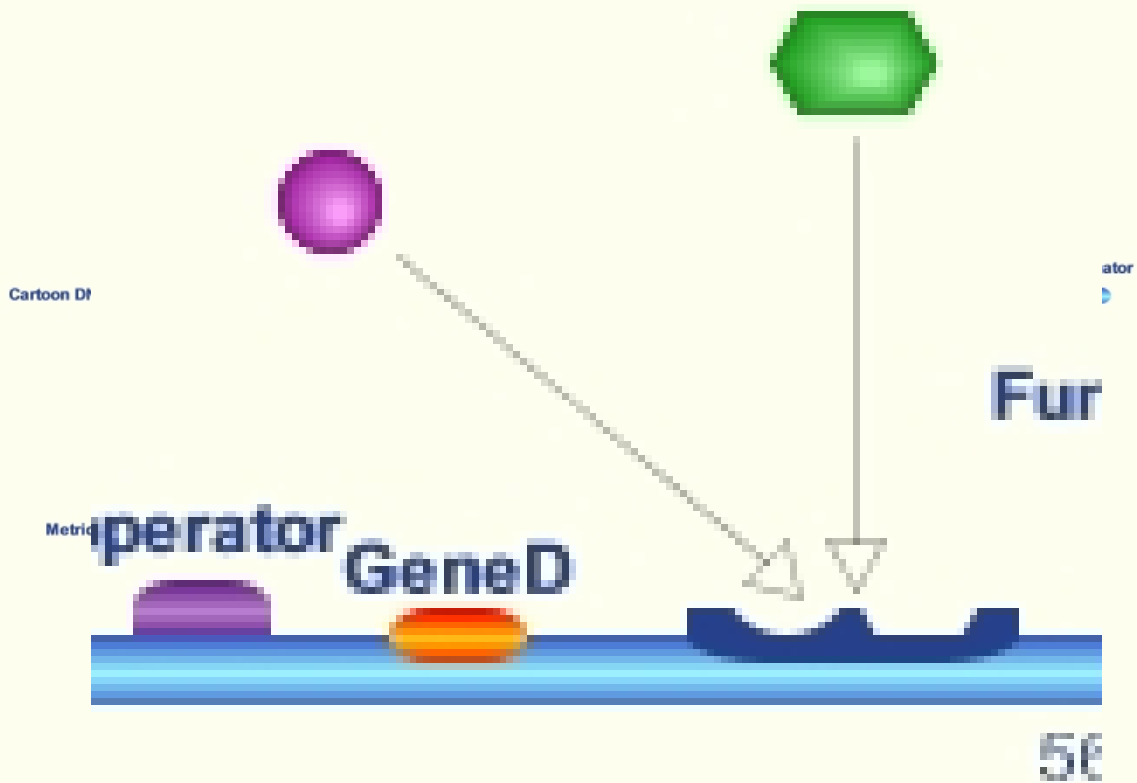


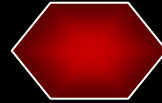
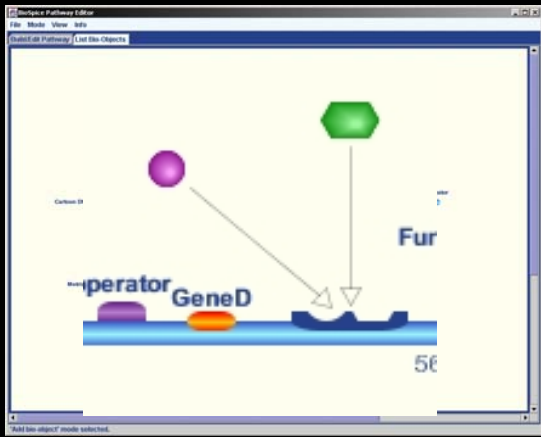
But how did we chose these icons?

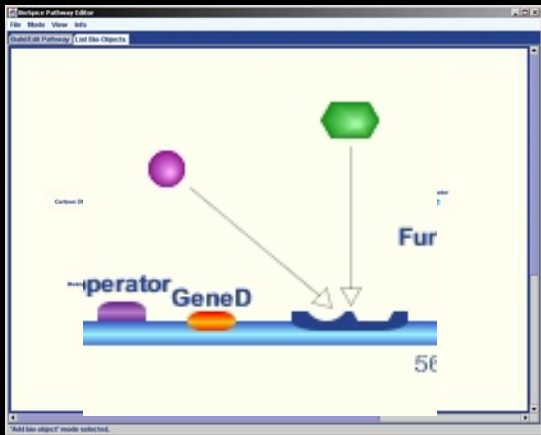
We didn't.

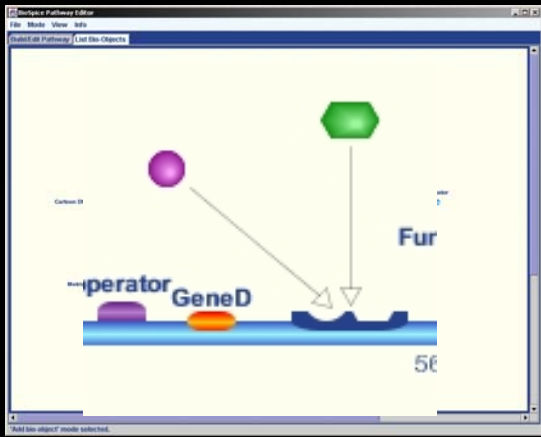
This is a big problem.



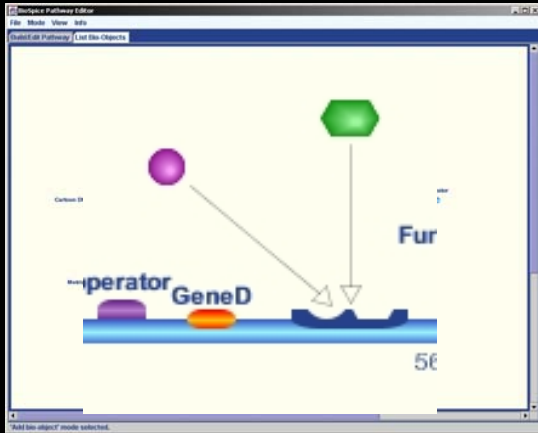












So how do we represent this information?

Depends on the data

Two/one hybrid data

Surface Plasmon Resonance

F.R.E.T.

Foot Printing

Depends on the model

Graphical

Thermodynamic

Kinetic

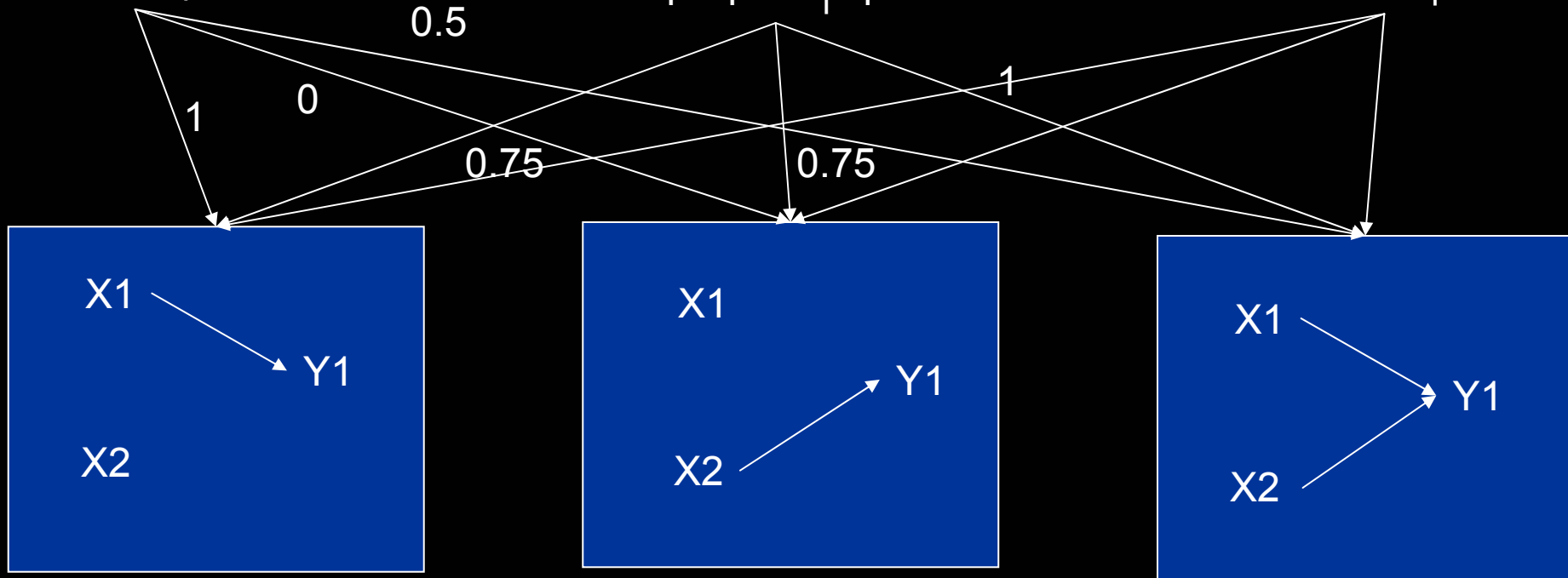
HYBRID

# Graphical Models: An incredibly stupid example!

X1	X2	Y1
0	0	1
0	1	1
1	0	0
1	1	0

X1	X2	Y1
0	0	0
0	1	1
1	0	1
1	1	1

X1	X2	Y1
0	0	0
0	1	0
1	0	0
1	1	1



Obviously, our data will be more complicated

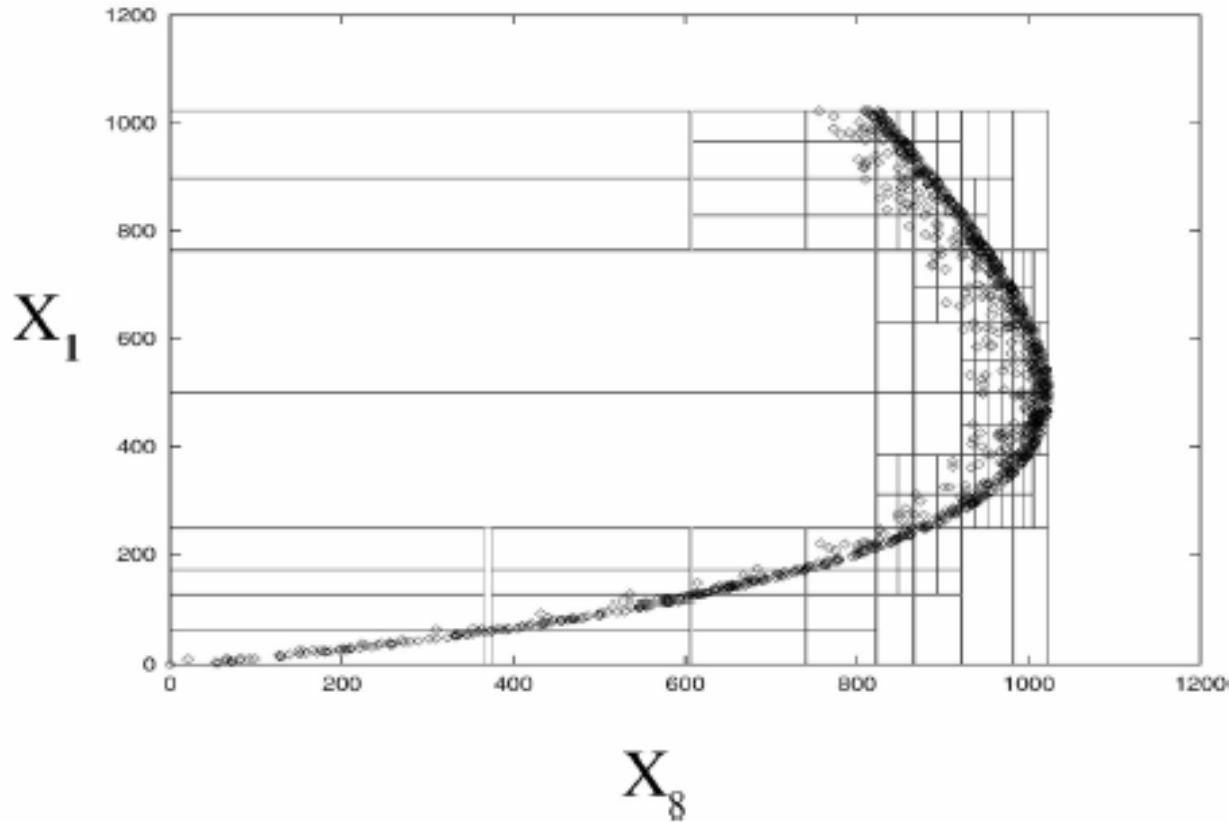


FIG. 5. The diamonds plot the values of  $X_1$  vs  $X_8$  obtained from the simulated time series. The rectangles are the result of a partitioning algorithm, see the text. From Ref. 32.

Samoilov, Arkin, Ross, *Chaos*, 11(1):108-114



This gets into models and modules but.....

For graphical models data representation can be VERY simple at first.  
For more complicated models we have to consider

1) What is a molecule?

Should we represent p53 as 1 molecule with  $2^{27}$  states?

2) What is an interaction?

Influence?

Direct binding?

How do we associate different data types with it?

3) How do we relate data at different “model levels” together?

# Knowledge representation for data classification and analysis

Aid to user in decision making.

Allows for data fusion.

For now:

Ontologies= Explicit specification of a conceptualization

Schema= A structure of tables in a database

Data Ontology

Analysis Ontology

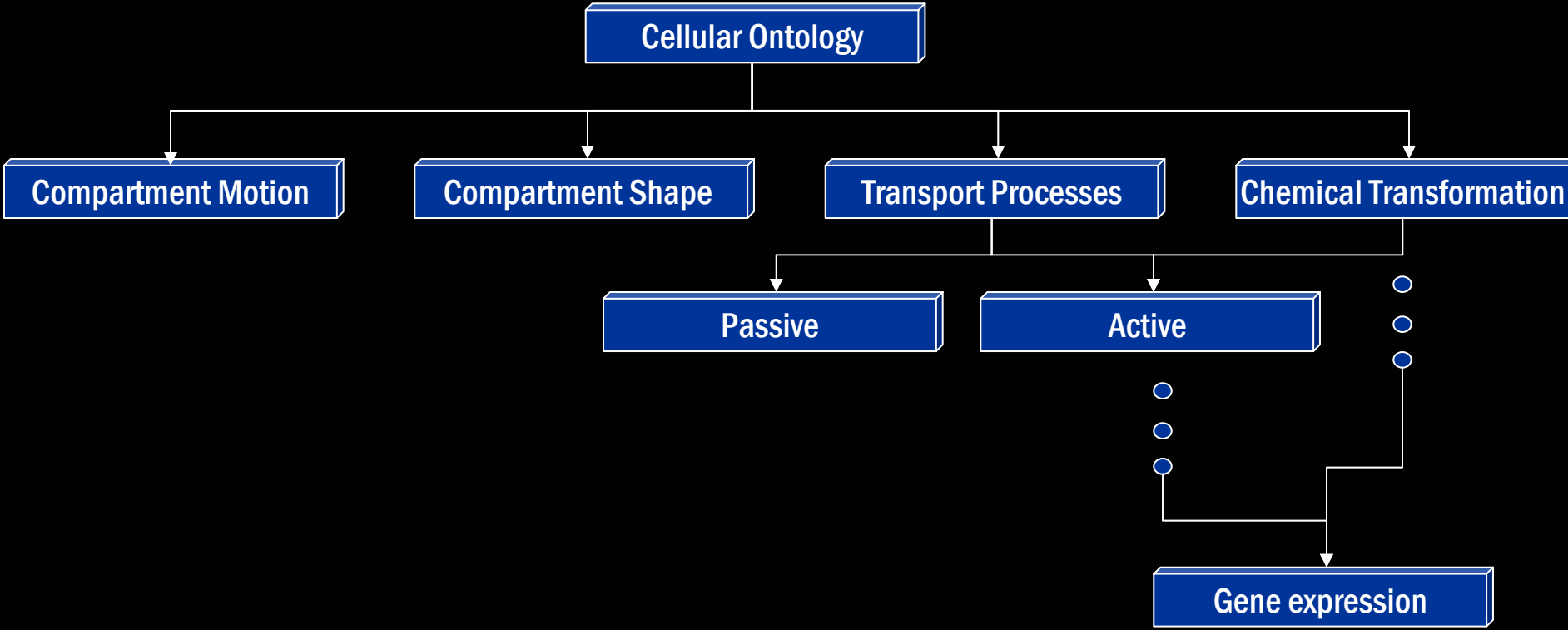
Mathematical Ontology

Differential, Algebraic, Stochastic

Cellular Ontology

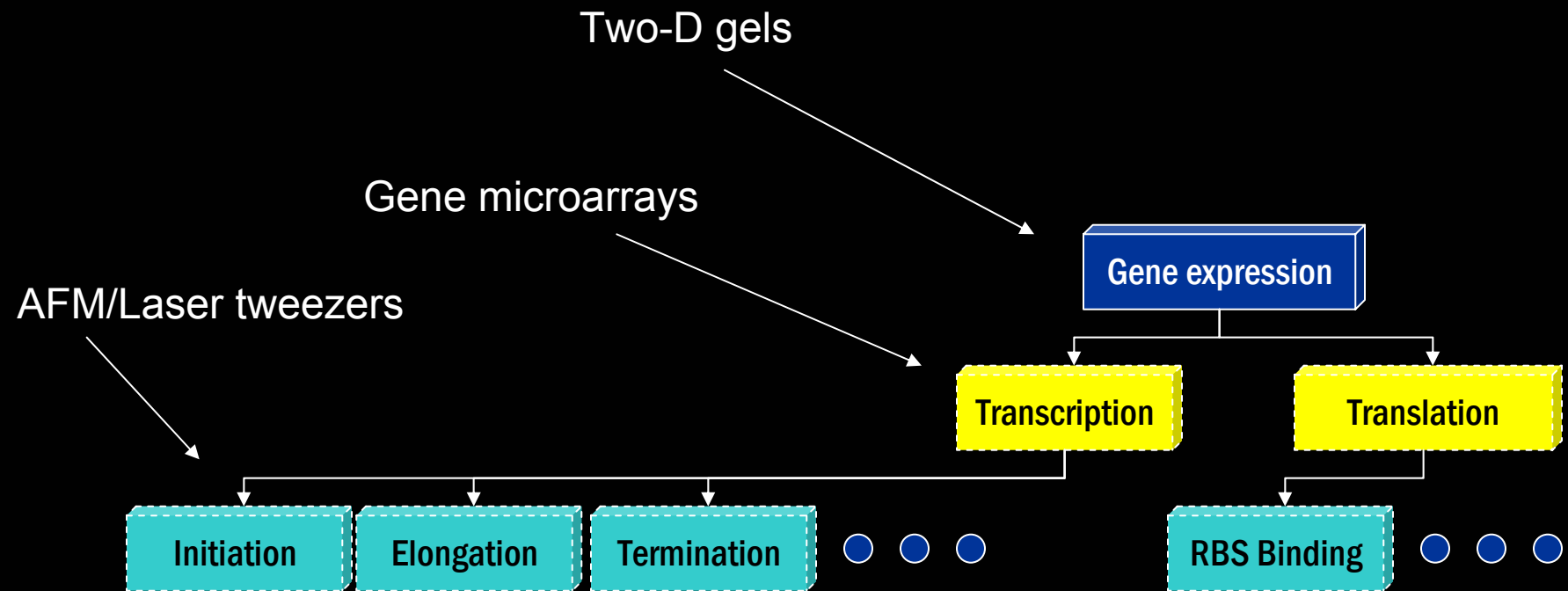
Motion, Shape Change, Transport, Transformation

# Knowledge representation for data classification and analysis



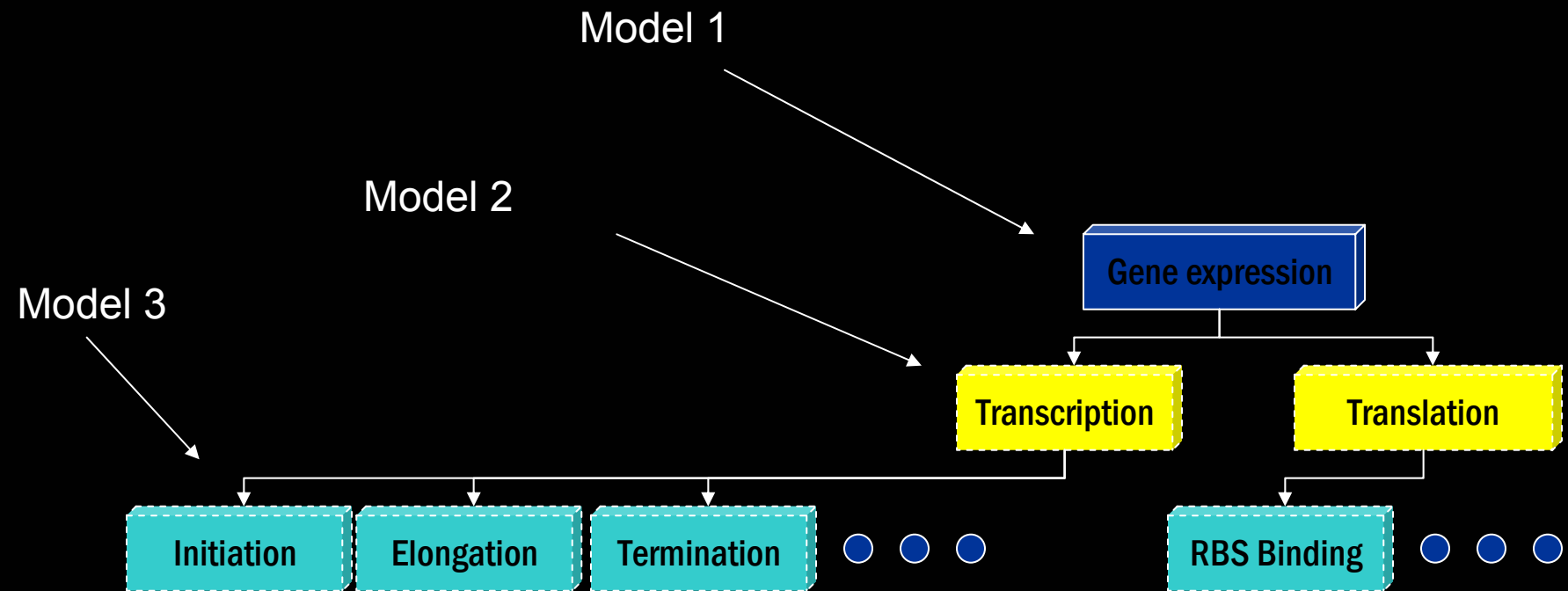
Bioontologies: <http://smi-web.stanford.edu/projects/bio-ontology/>

# Leaves of the ontologies: Cellular

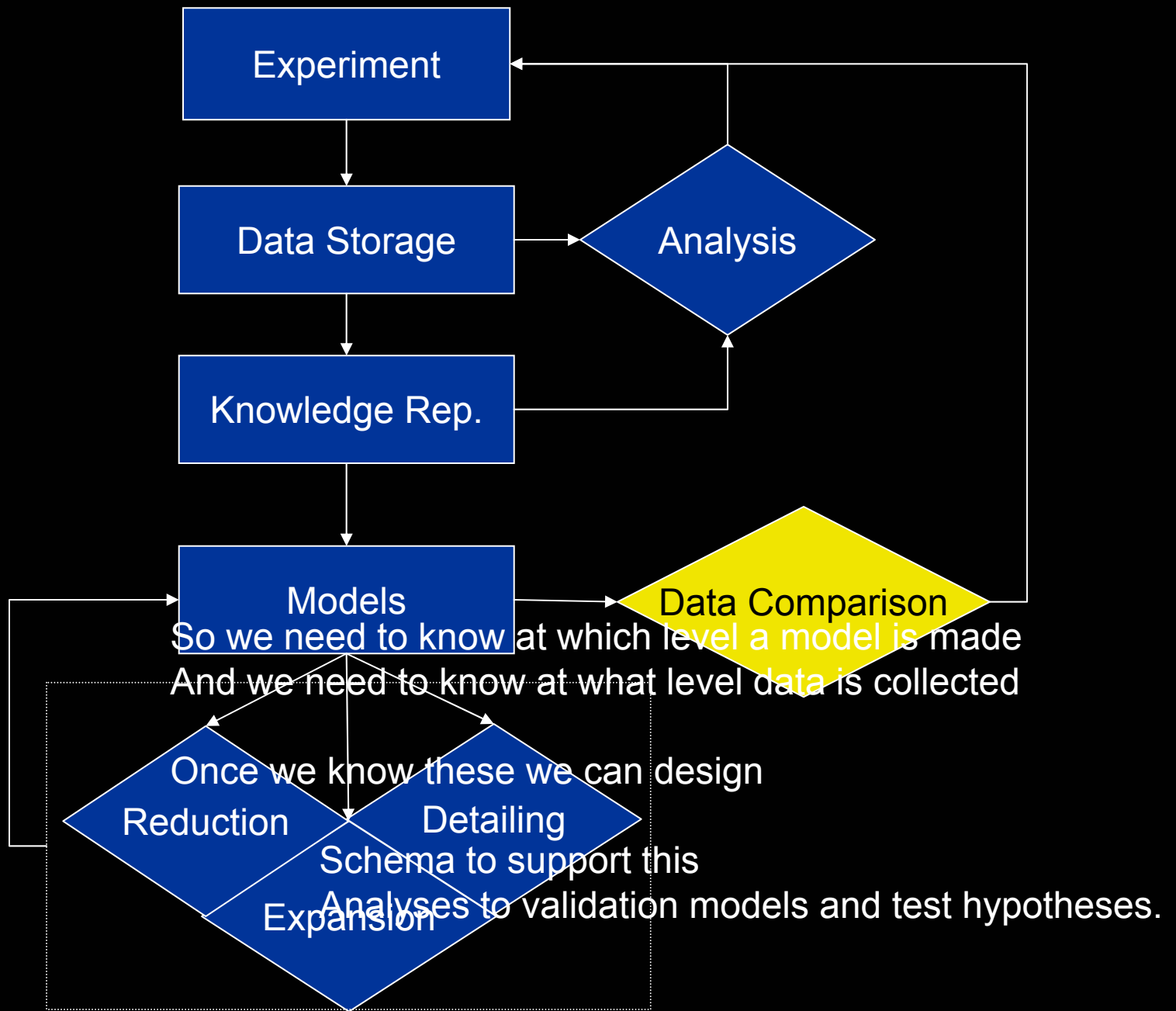


Forms a hierarchy for modeling and data

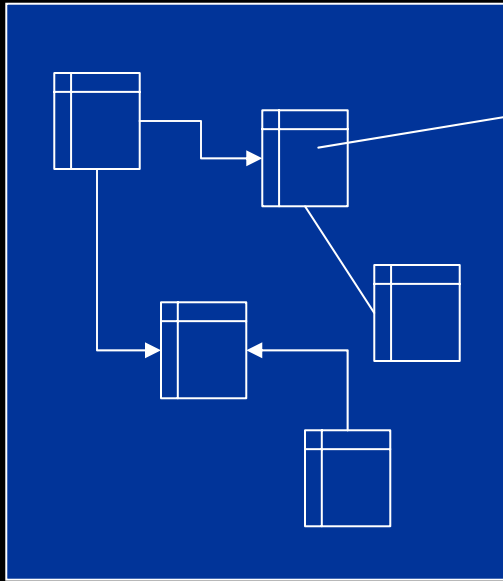
# Leaves of the ontologies: Cellular



Forms a hierarchy for modeling and data



# Technical fields for data comparison



For every piece of "data" in the data base

Data base schema  
Assumed to contain  
Organism  
Strain  
Protocols  
Etc.

## Comparison Tables

- Provenance
- Type
- Method
- Quality
- Revision
- Alternate

Comparison Tables

Provenance

Source1, Source2, ...

Authors



Comparison Tables

Type

Hypothetical

Calculated

Indirect

Direct

Comparison Tables

Method

HMM Predicted

Two-Hybrid

Microarray

Mass Spec

Comparison Tables

Quality

Format Correctness Checked

Relevance Scores

Interpretation Scores

Reproducibility

Current Accepted Data

Comparison Tables

Revisions

Author

Date

Why

Comparison Tables

Alternatives

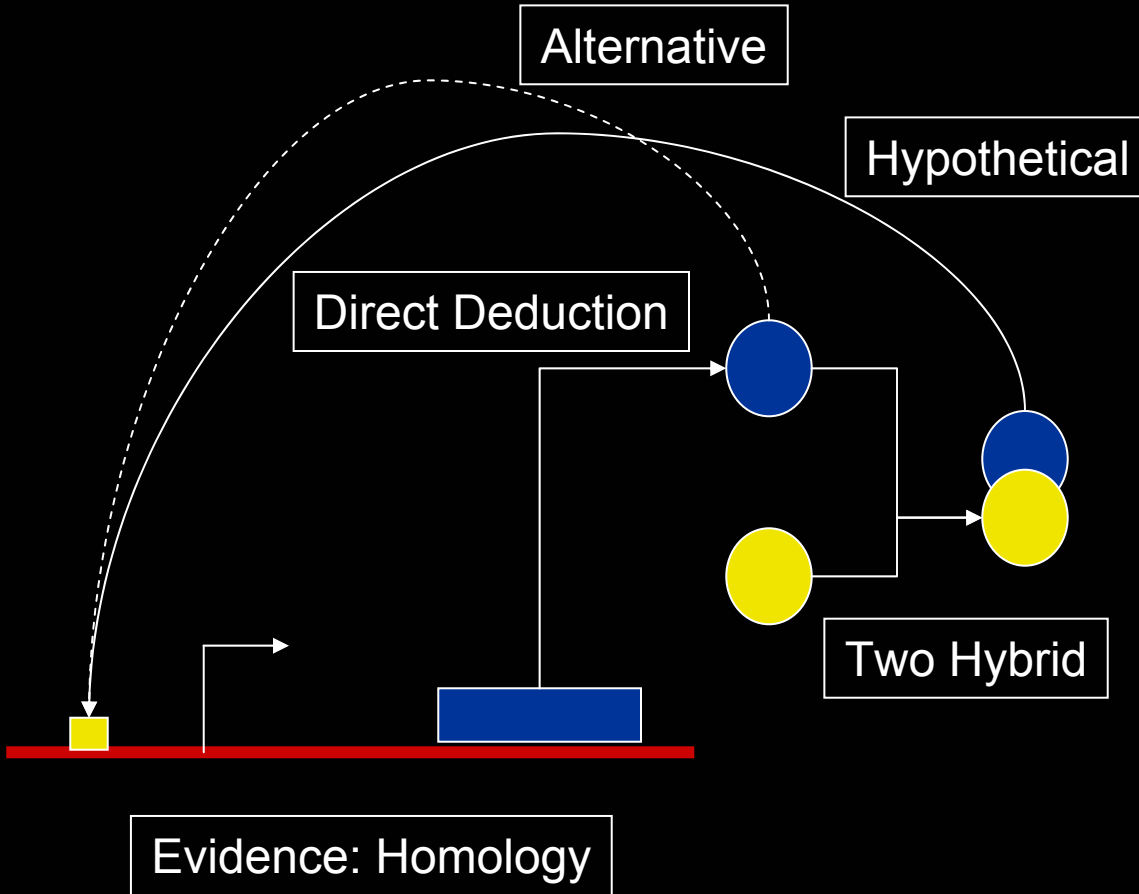
Author

Date

Why

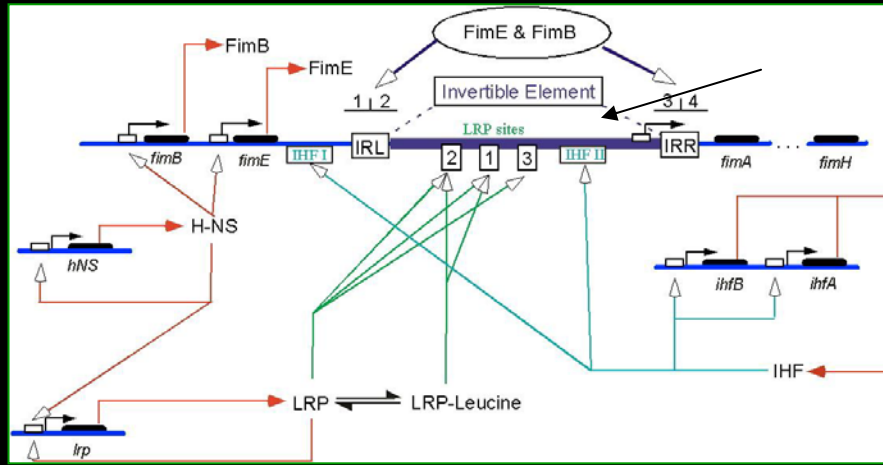
Alternatives are especially important for model data bases.

A model is a collection of data of a particular sort and hypotheses.





# The Modeling Process: Bottom Up!



This pathway is diagramed using a notation. Stochastic models are chosen for the DNA, and deterministic models for the proteins.

Equilibrium statistical thermodynamics for  $f, g$

$$\begin{aligned}
 \frac{dP_{on}}{dt} &= f * (1 - P_{on}) - g * P_{on} \\
 &= \frac{\sum_{s \in OFF} \alpha_s e^{-\Delta G_s / RT} [IHF] [FimE]^{j(s)} [FimB]^{k(s)} [Lrp^*]^{m(s)} [Lrp]^l(s)}{1 + \sum_{s \in OFF} e^{-\Delta G_s / RT} [IHF] [FimE]^{j(s)} [FimB]^{k(s)} [Lrp^*]^{m(s)} [Lrp]^l(s)} (1 - P_{on}) \\
 &\quad - \frac{\sum_{s \in ON} \alpha_s e^{-\Delta G_s / RT} [IHF] [FimE]^{j(s)} [FimB]^{k(s)} [Lrp^*]^{m(s)} [Lrp]^l(s)}{1 + \sum_{s \in ON} e^{-\Delta G_s / RT} [IHF] [FimE]^{j(s)} [FimB]^{k(s)} [Lrp^*]^{m(s)} [Lrp]^l(s)} P_{on} \quad (1)
 \end{aligned}$$

$$\frac{d[FimB]}{dt} = k_1 \frac{e^{-\Delta G_b / RT} [RNAP]}{1 + e^{-\Delta G_b / RT} [RNAP] + e^{-\Delta G_{hb} / RT} [H - NS]} - \lambda_1 [FimB] \quad (2)$$

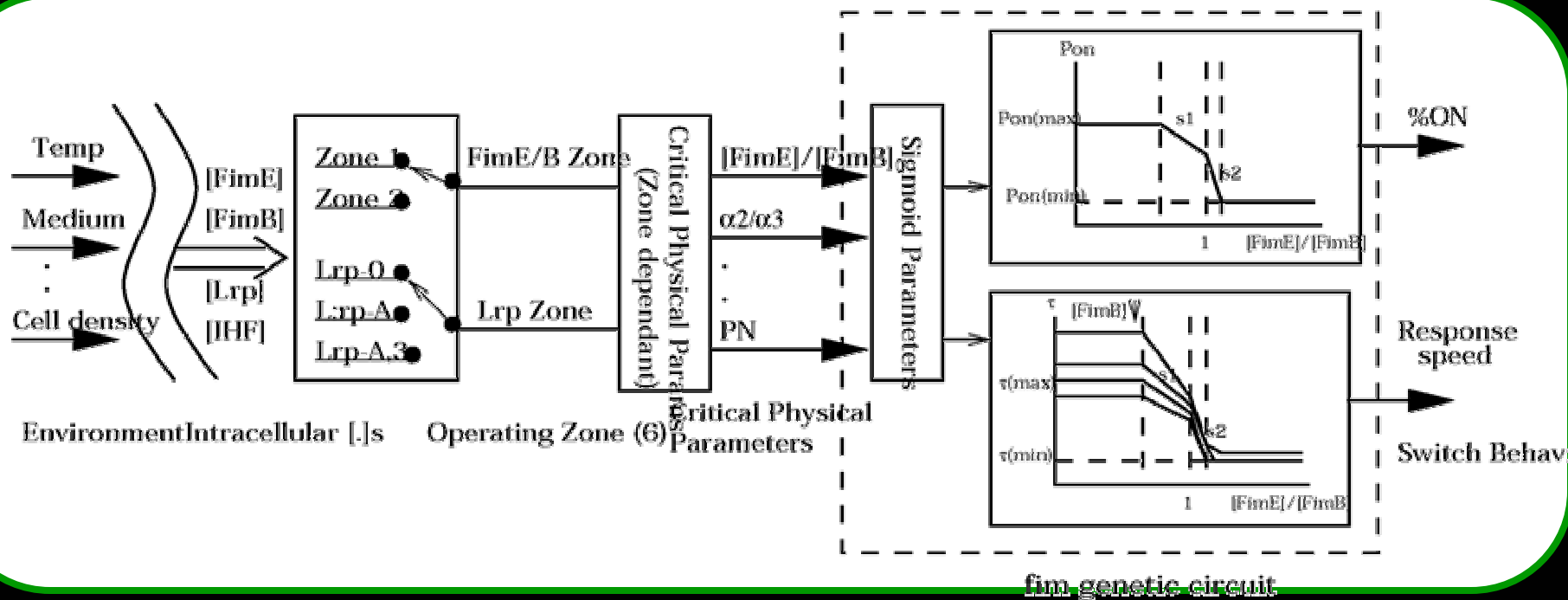
Protein degradation

$$\begin{aligned}
 \frac{d[FimE]}{dt} &= (k_2 \frac{e^{-\Delta G_e / RT} [RNAP]}{1 + e^{-\Delta G_e / RT} [RNAP] + e^{-\Delta G_{he} / RT} [H - NS]}) P_{on} \\
 &\quad + (k_3 \frac{e^{-\Delta G_e / RT} [RNAP]}{1 + e^{-\Delta G_e / RT} [RNAP] + e^{-\Delta G_{he} / RT} [H - NS]}) (1 - P_{on}) - \lambda [FimE] \quad (3)
 \end{aligned}$$

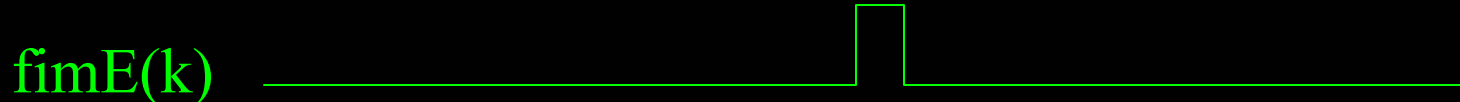
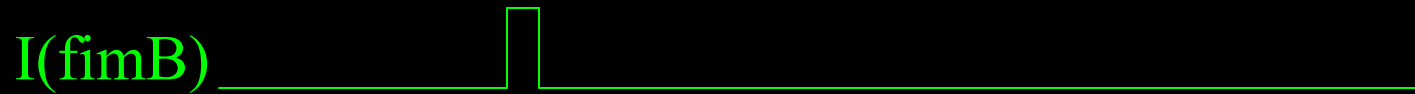
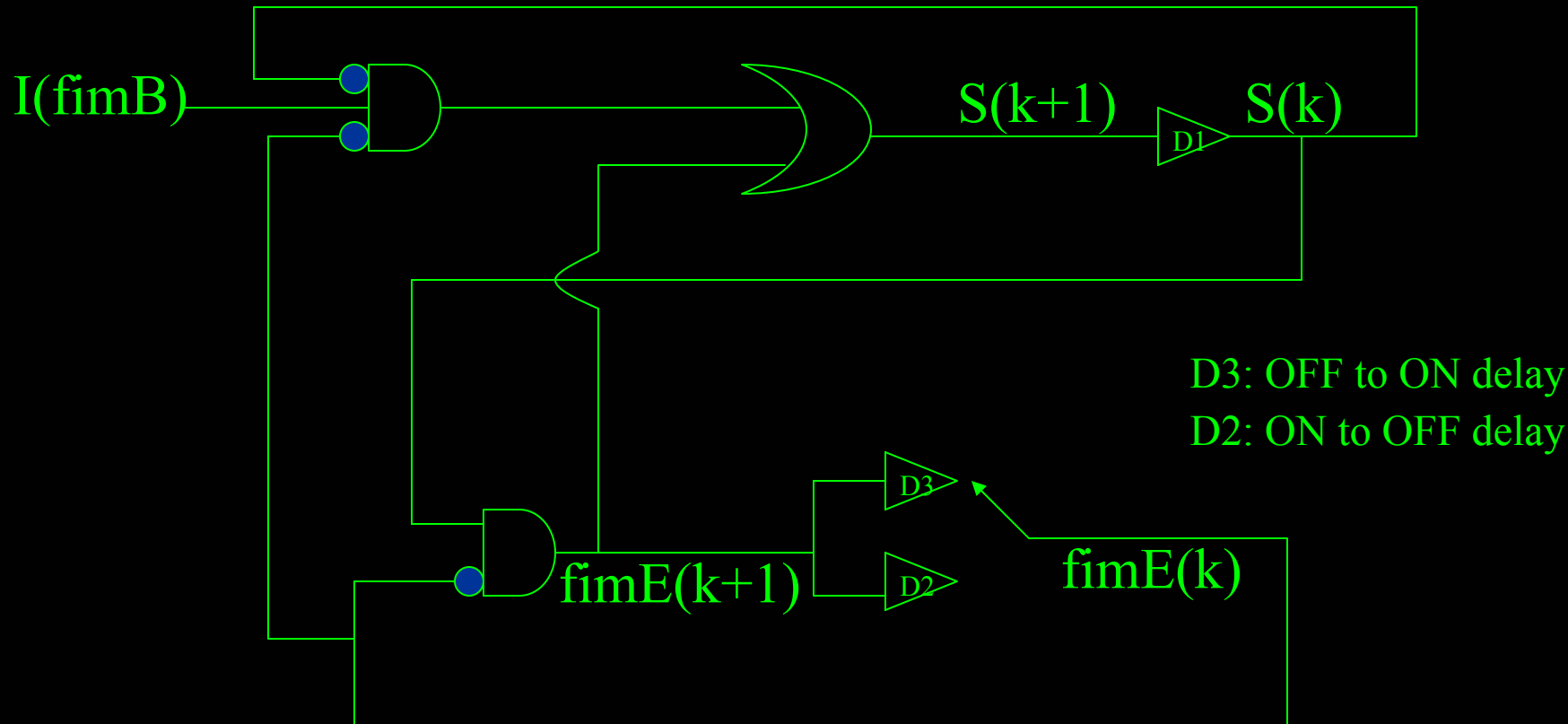
This is a hybrid model at the molecular level of abstraction. This now becomes a submodel of a larger infection model.



# Expert Logical Abstraction



Any expert examining these results can rapidly deduce a logical control diagram. This is not a formal abstraction but could be used as the basis for a simpler (lower-resolution) model of the switch to be used in the larger simulation.



Since biological models must represent models at these different levels of abstraction

How we represent molecules is VERY important.

Either we store models as unrelated directly to primary data in the database as a lump of variables and equations

Hard to modify

Hard to relate to data

Hard to deal with a family of models

Hard to deal with a linked set of abstractions (stay tuned)

Or we ensure that the objects that models describe are represented for modeling.



If a phosphorylatable protein is one molecule with internal state, how does a model specification refer to it?



If we represent all states the DB gets bloated?



If states group how do we group them?



## Adam's ignorant scheme

1. If data exists on a particular state of a molecule it is given its own record
2. This record is referenced back to the “parent” molecules
  1. Defined as the molecules from which this molecule may be created.
3. If data on a particular “possible” state does not exist do the same thing.
4. Define equivalence class specifications
  1. This is a model object!

Representation and interoperability: Models can be passed around

CellML- <http://www.cellml.org>

SBML- <http://www.cds.caltech.edu/erato/>

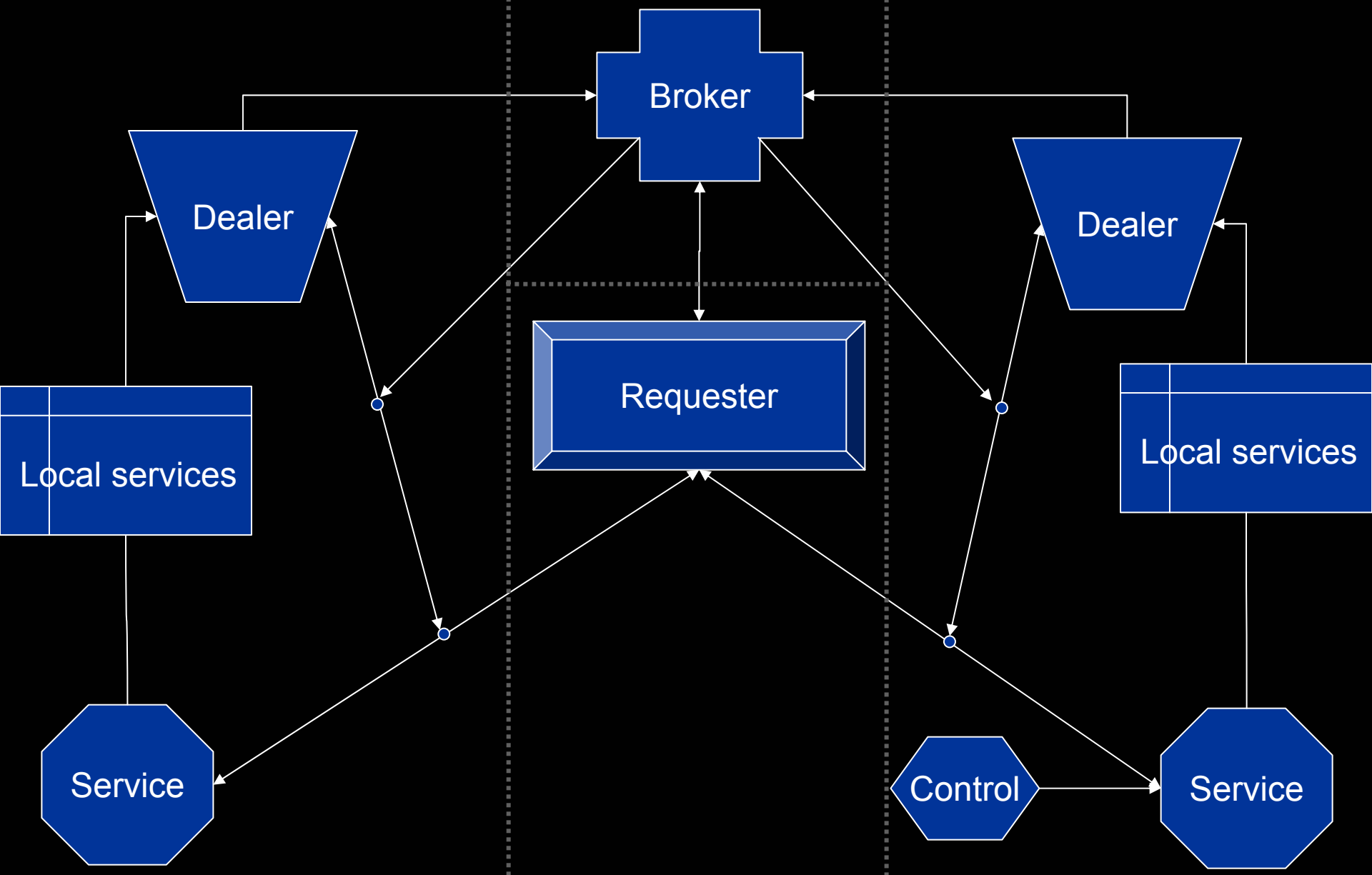
**JOIN THE DISCUSSION NOW!**

```
<model>
...
  <listOfUnitDefinitions>
    <unitDefinition name="volume">
      <listOfUnits>
        <unit kind="liters" scale="-3"/>
      </listOfUnits>
    </unitDefinition>
  </listOfUnitDefinitions>
...
</model>
```

These specifications are designed for sending models to simulators.  
But what if there is data comparison, etc.

However, this is “middleware” and can solve a plethora of problems

Distributed, loosely coupled, architecture



# Conclusions

- Data representation is the lowest level of representation of biological knowledge
- Models are particular “statements” of this knowledge.
- Databases and models must be linked for comparison
- How schema should be designed to facilitate this is research
  - Much can be fixed post-facto with “middleware”

# Data analysis, Modules and Models

Prospectus and Problems

Pessimistic optimism

Panglossian Pessimism



# Data analysis

An Example of Effect Detection

# Questions

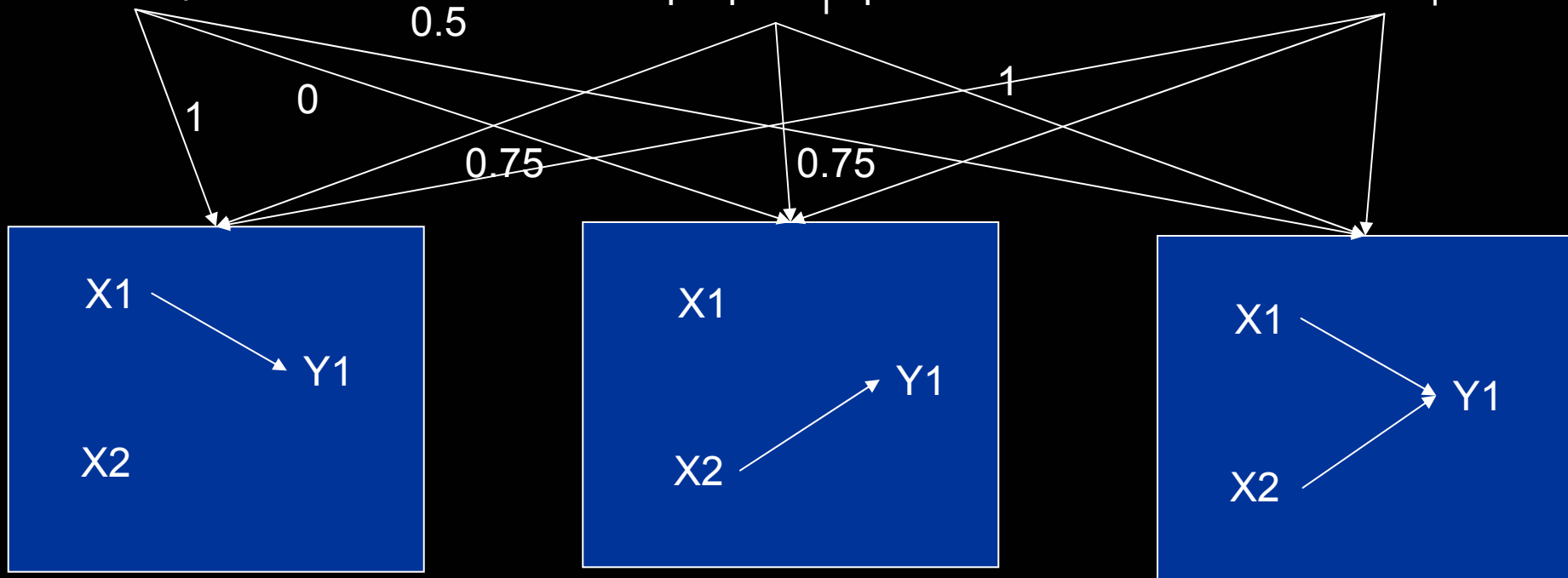
- What's my question?
  - Does this perturbation have an effect?
  - Is loss of this interaction responsible for  $X$ ?
  - Can the known network reproduce dynamics?
- What experimental design best answers my question?
- Do I need complicated statistics?

# Graphical Models: An incredibly stupid example!

X1	X2	Y1
0	0	1
0	1	1
1	0	0
1	1	0

X1	X2	Y1
0	0	0
0	1	1
1	0	1
1	1	1

X1	X2	Y1
0	0	0
0	1	0
1	0	0
1	1	1



## Graphical Models: An incredibly stupid example!

X1	X2	Y1
0	0	1
0	1	1
1	0	0
1	1	0

X1	X2	Y1
0	0	0
0	1	1
1	0	1
1	1	1

X1	X2	Y1
0	0	0
0	1	0
1	0	0
1	1	1

For simple discrete *combinational* data, we can enumerate all possible states and deduce directly.

For sequential data, large combinational data or continuous data things become more complicated.

We need roughly measures of the sort:

$$P(Y=y \mid X_1, X_2, \dots)$$

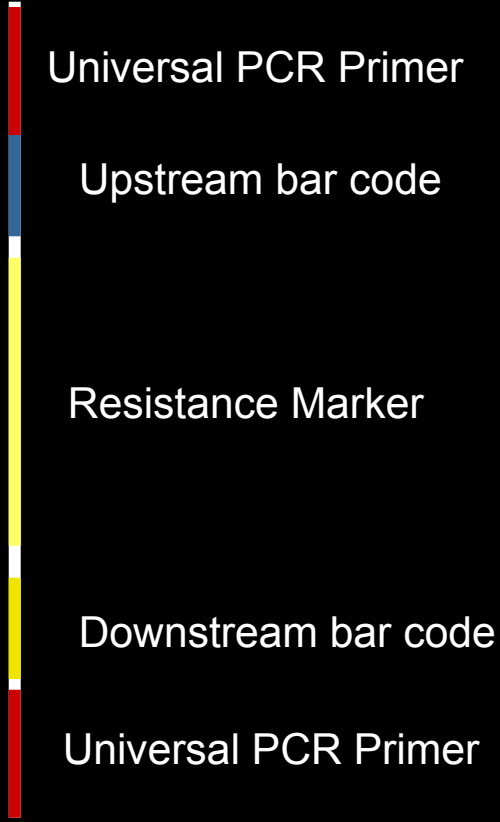
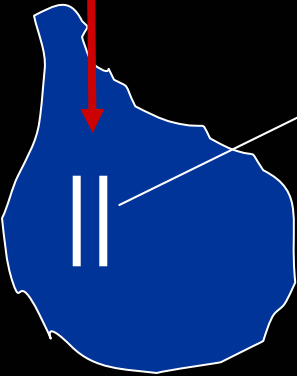
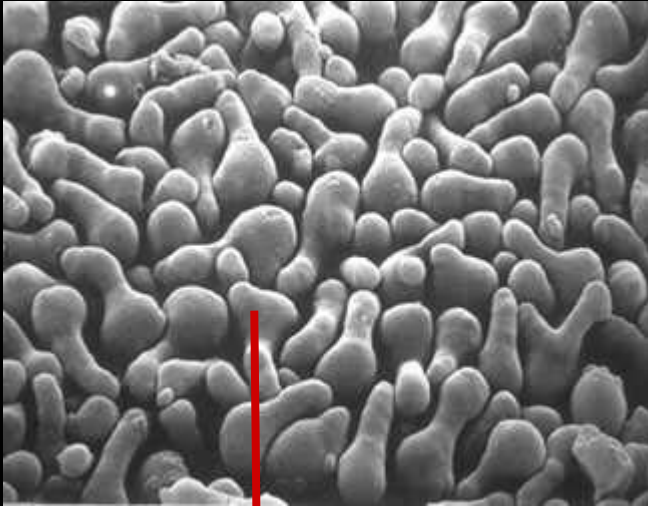
# Simpler Question

Effect Detection

Guri Giaever and Ron Davis

(confidential and VERY early)

# Yeast Haploinsufficiency Trials



Universal PCR Primer

So there are 8 bar codes:

Upstream bar code

Up  
Down

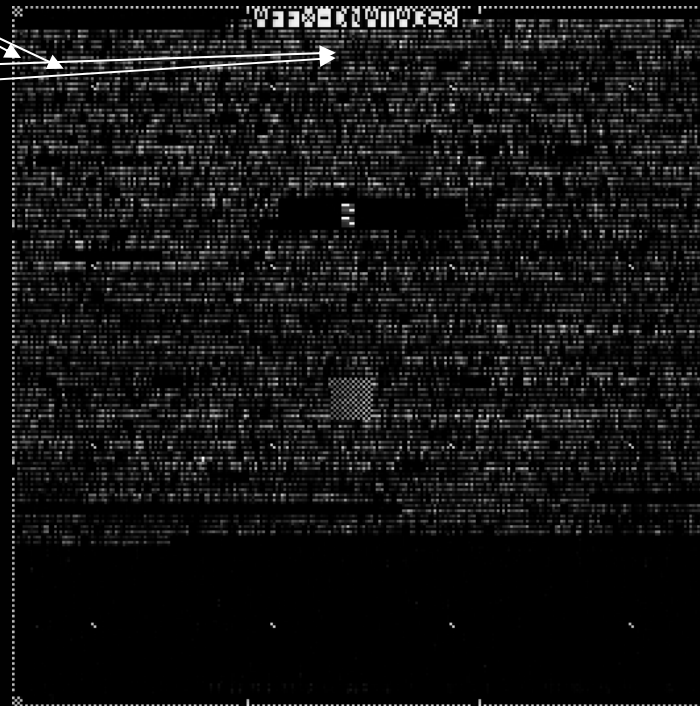
Resistance Marker

Antisense Up  
Antisense Down

Downstream bar code

Mismatches for all

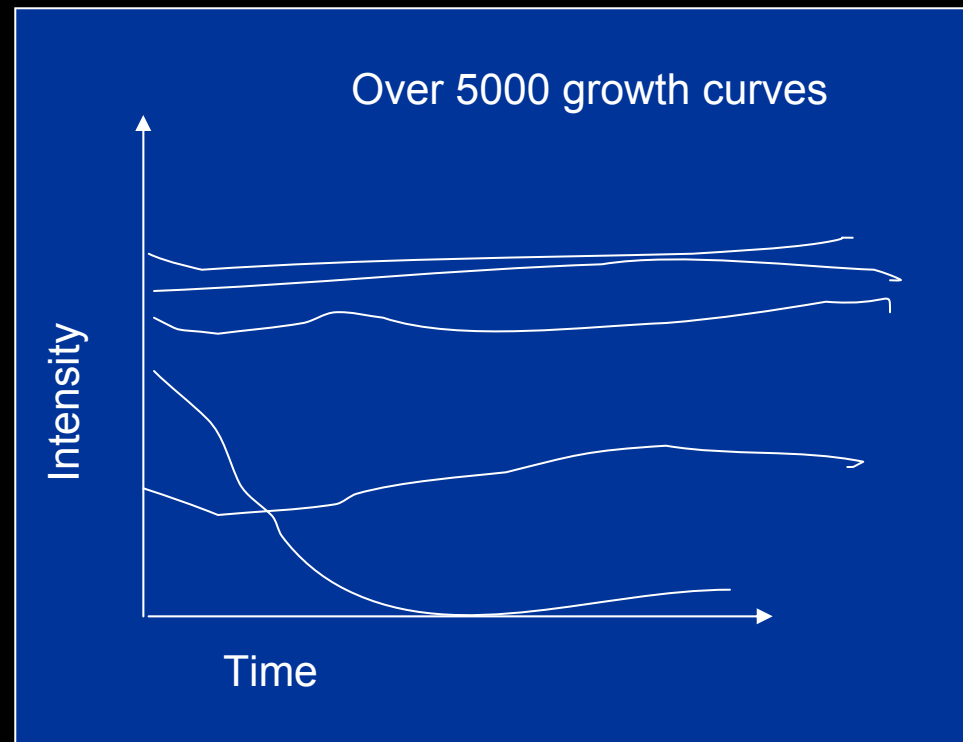
Universal PCR Primer



# Yeast Haploinsufficiency Trials



Sample every 0.5 population doubling.  
Dilute sample to standard OD.  
Examine on chip.



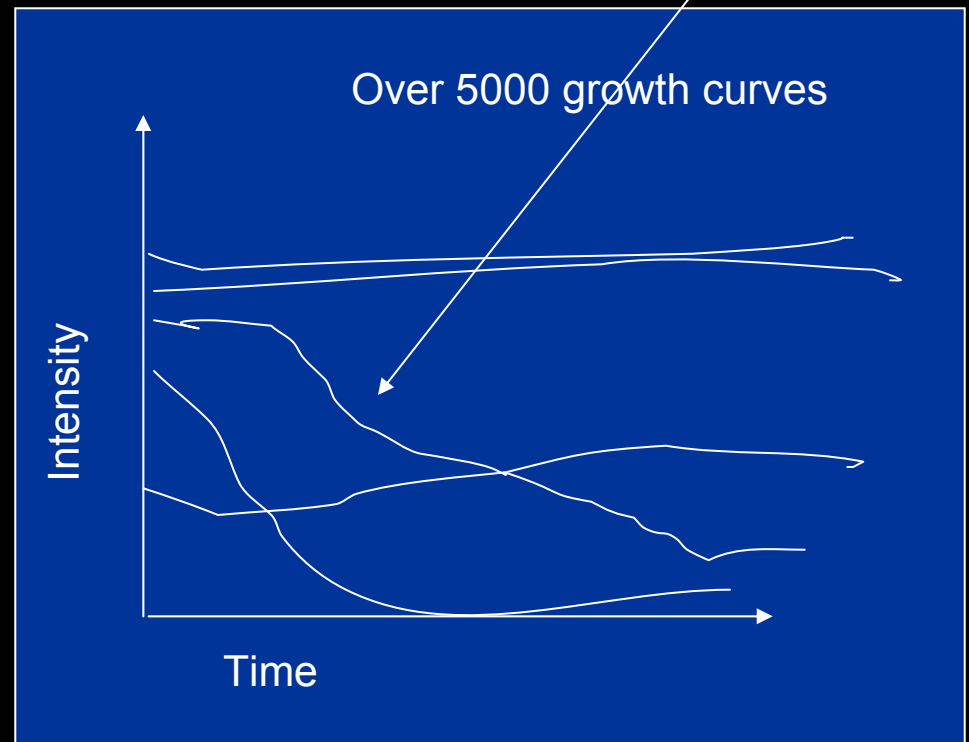


# Yeast Haploinsufficiency Trials



Drop drug in!

This guy now drops out!

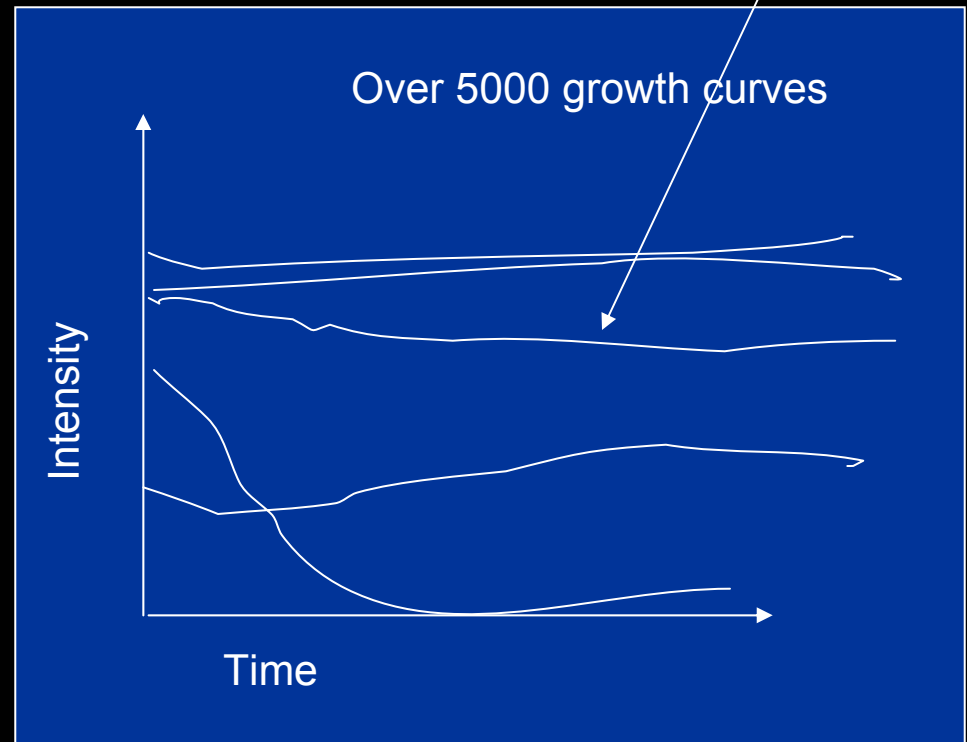


# Yeast Haploinsufficiency Trials



Drop drug in!

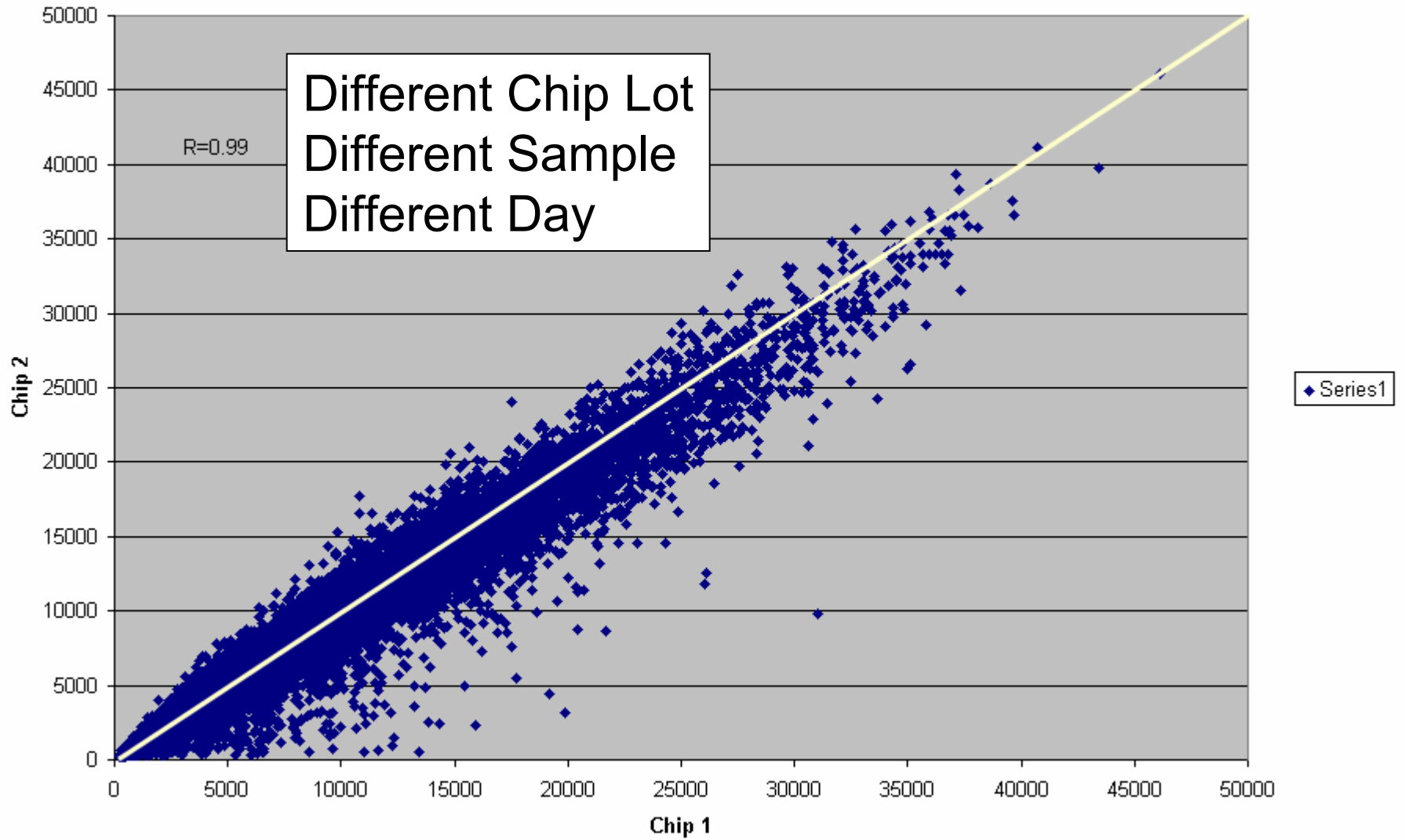
Is this guy dropping out?

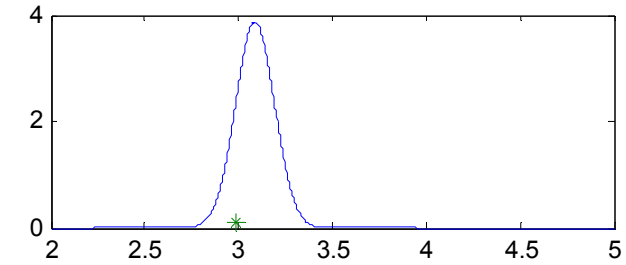
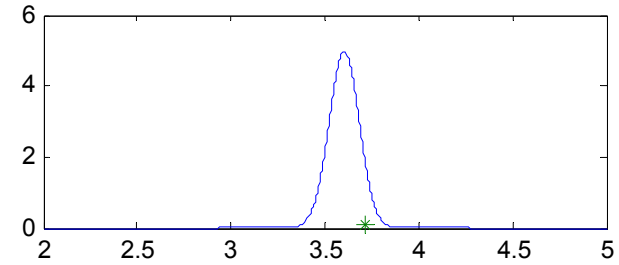
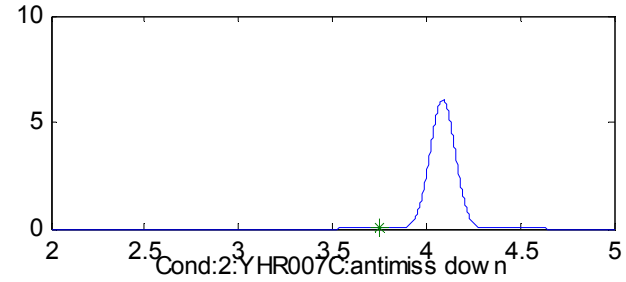
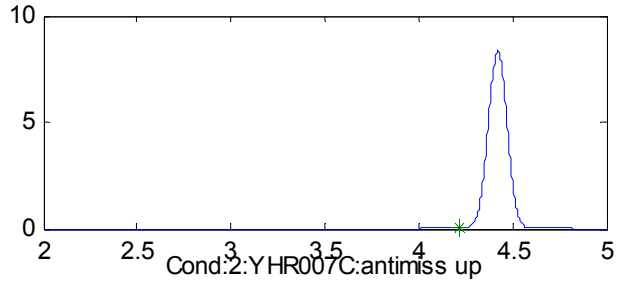
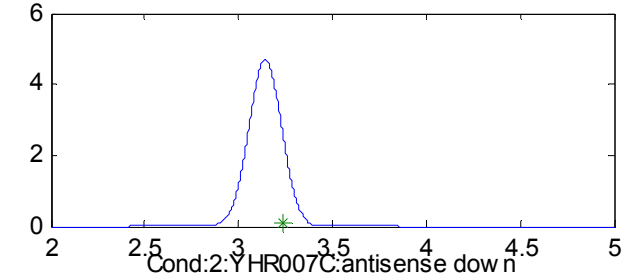
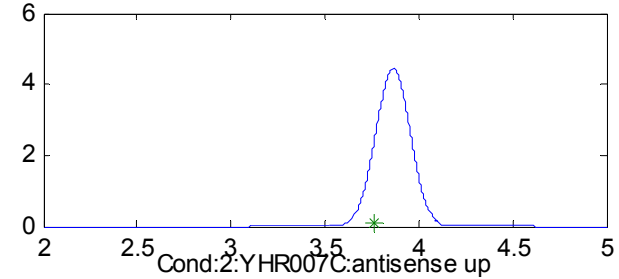
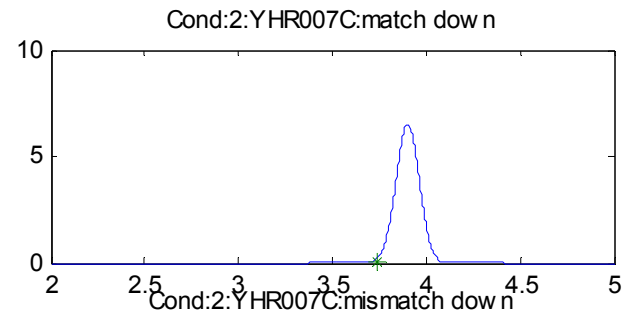
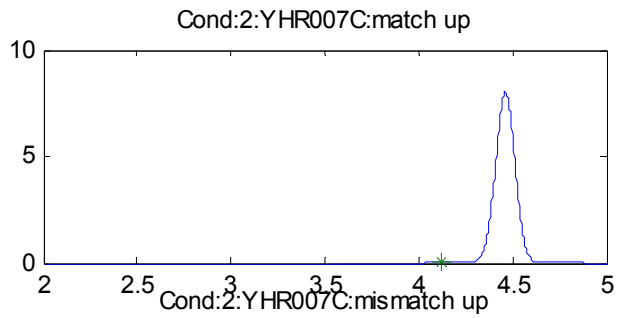


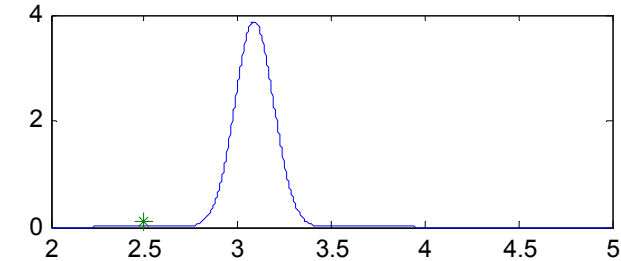
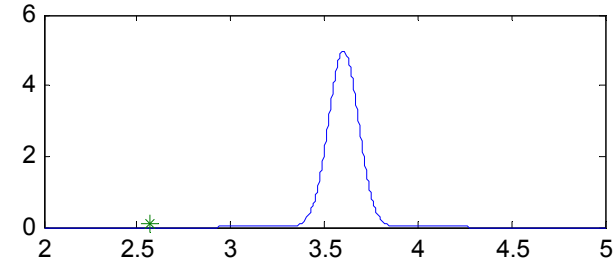
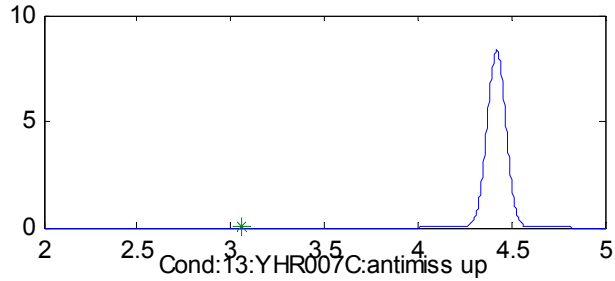
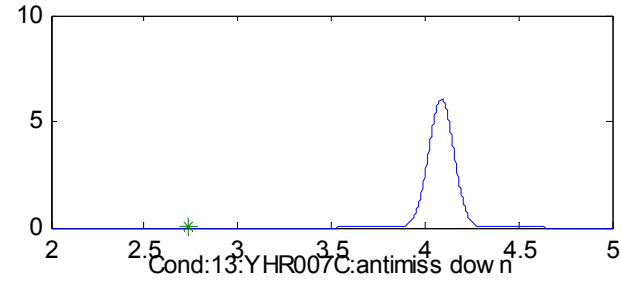
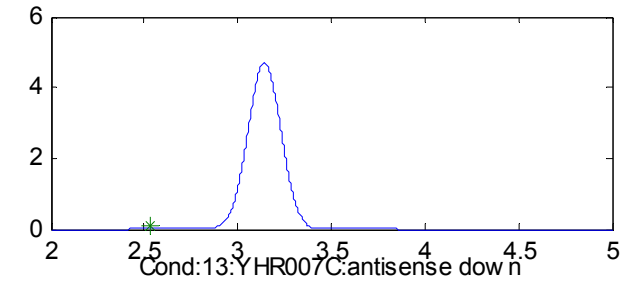
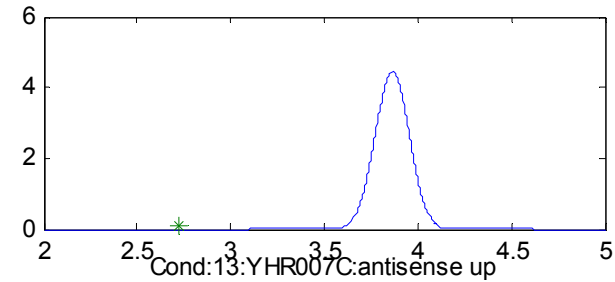
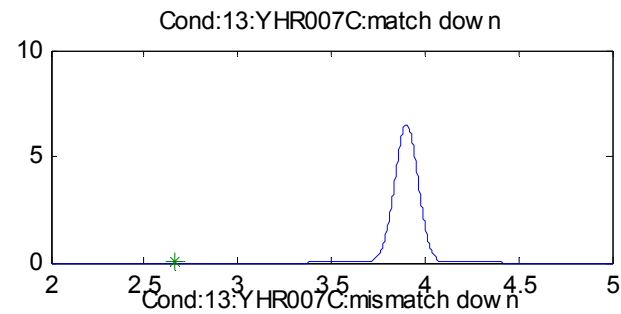
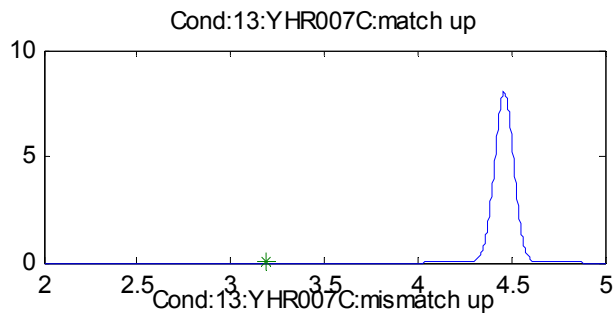
# Approach

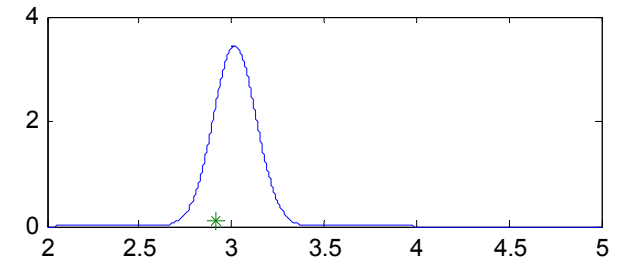
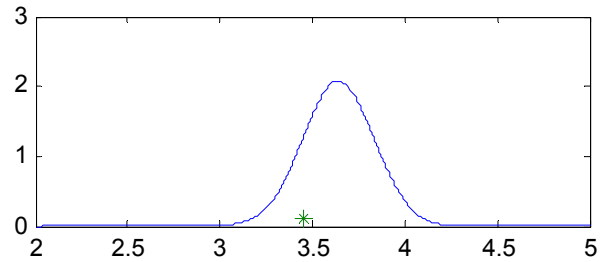
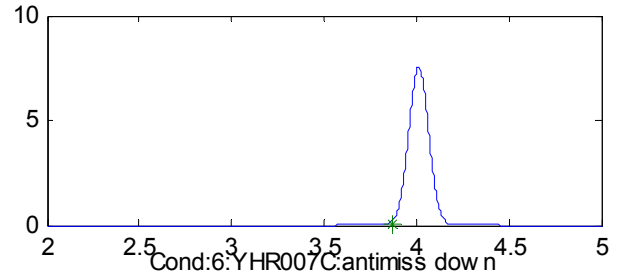
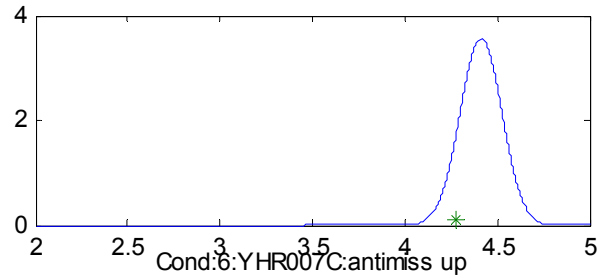
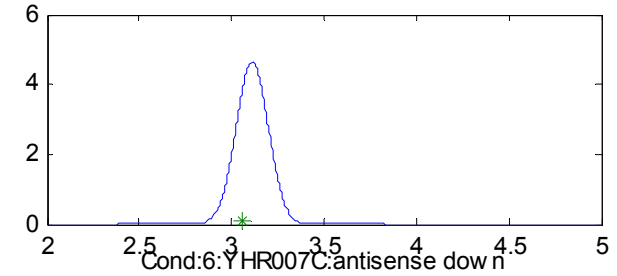
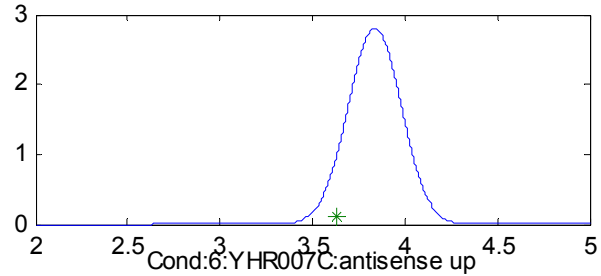
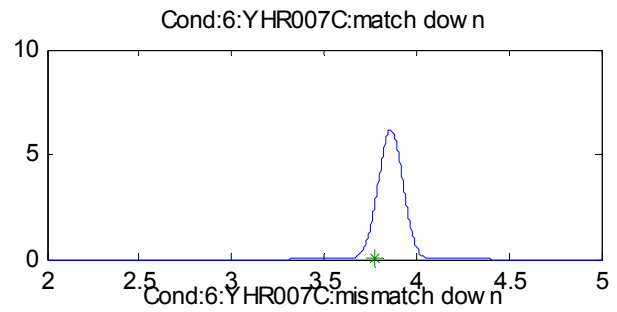
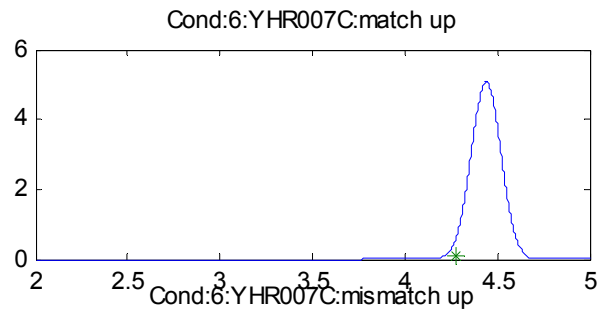
1. Look at time course, fit to growth curve, score by decay
2. Look at time zero, look at time= $T$ , score by ratio
3. Develop likelihood measure. (A stupid but effective one)
  1. Decide on a standard condition
  2. Do  $N$  replicates at each generation time
  3. Estimate the distribution of intensity for each tag at a given time.
  4. Score new experiments as the probability of being drawn from that distribution.

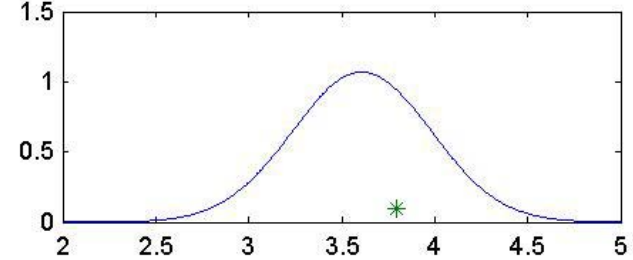
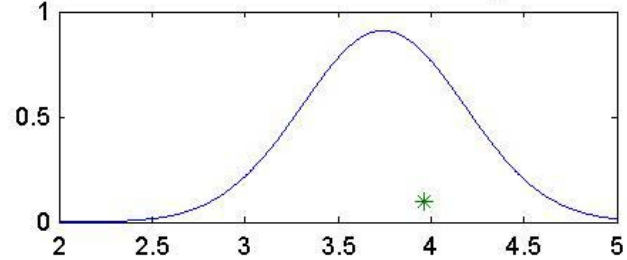
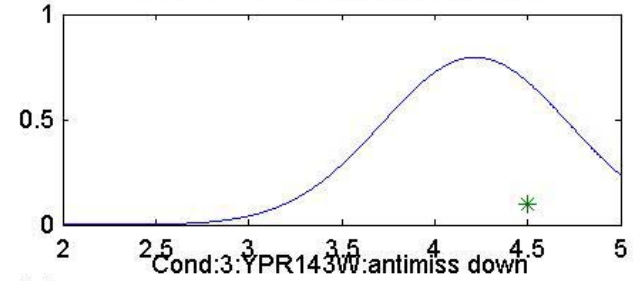
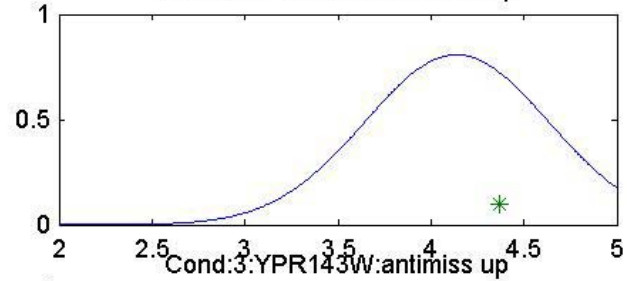
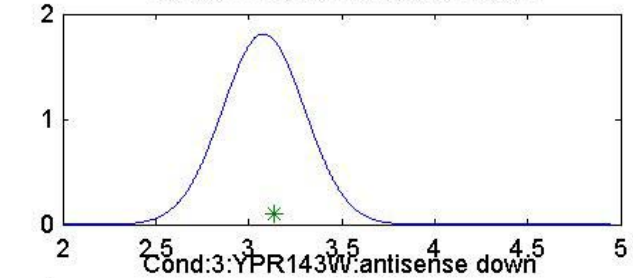
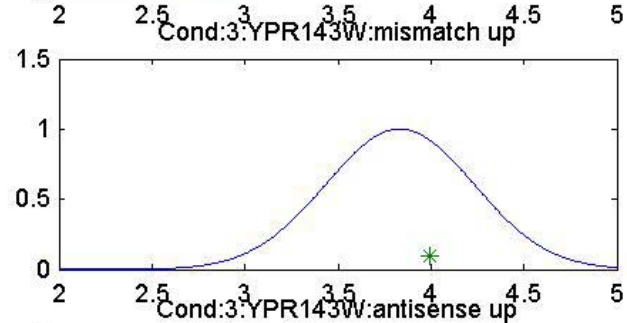
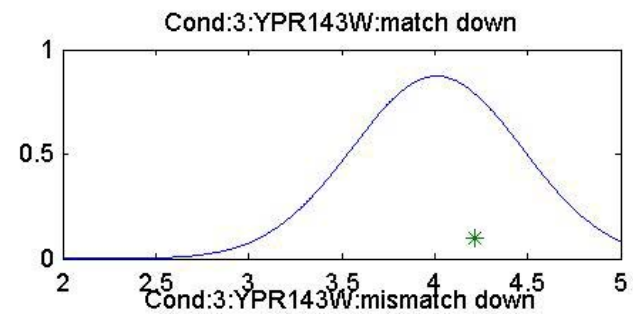
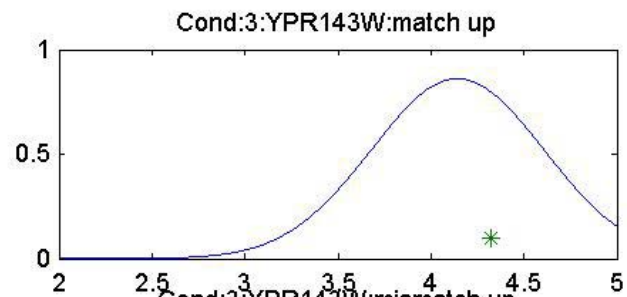
# Chip Reproducibility



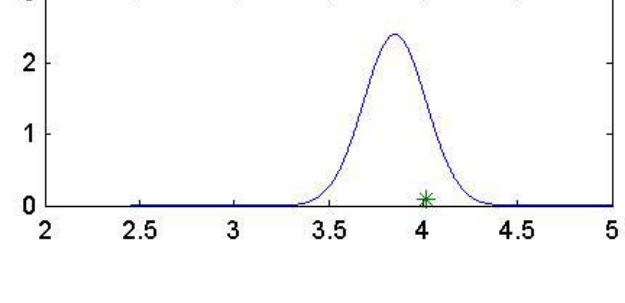
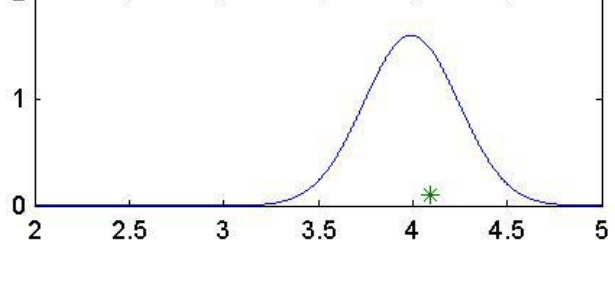
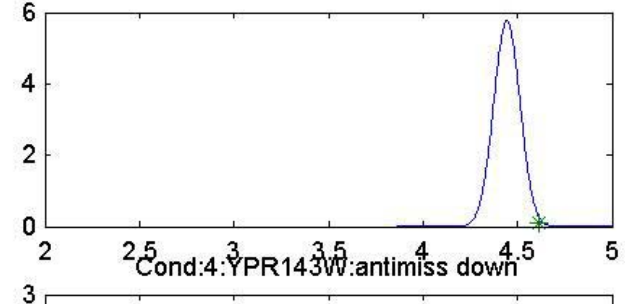
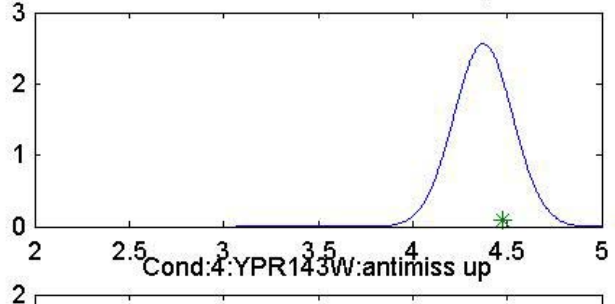
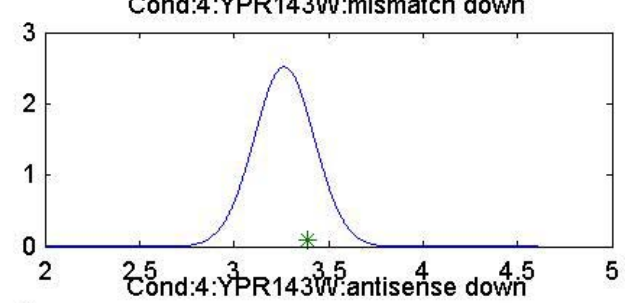
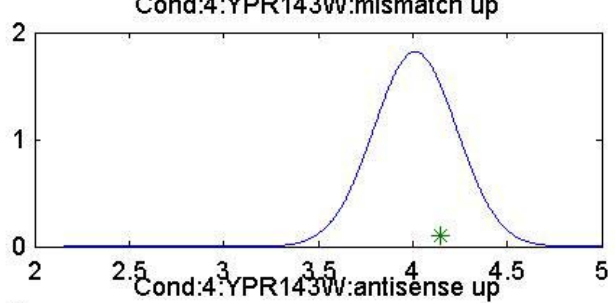
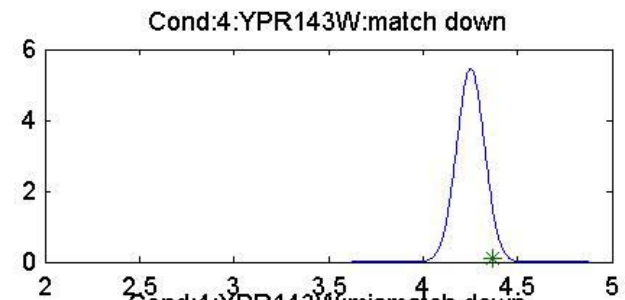
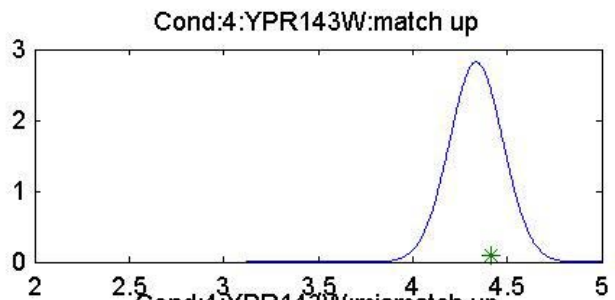




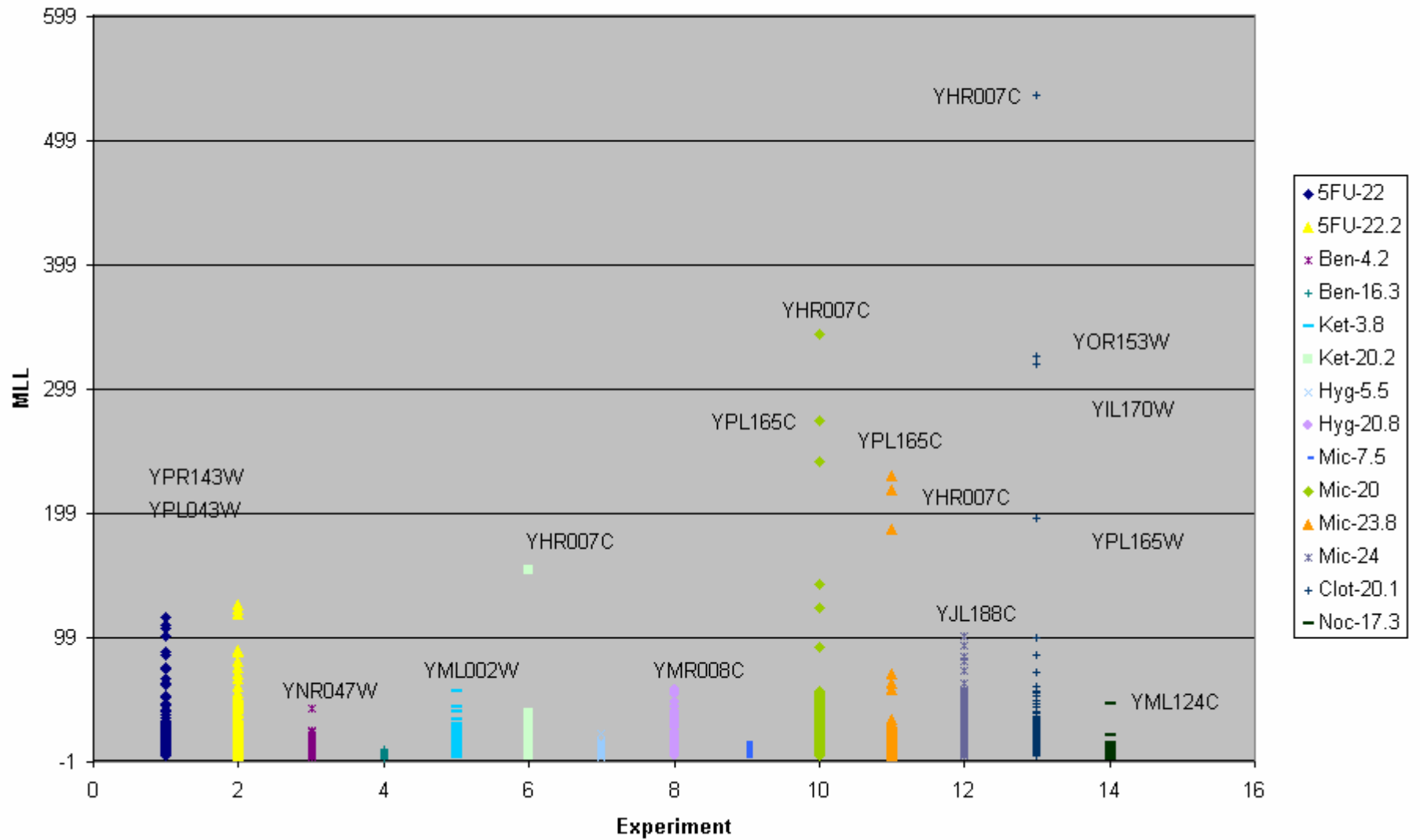








### Sensitivity Plots



## Experiment 6: Ket 20.2 hours

<b>YHR007C</b>	ERG11	cytochrome P450 lanosterol 14a-demethylase, ergosterol biosynthesis, lanosterol 14-alpha-demethylase, endoplasmic reticulum
<b>YML085C</b>	Tub1	alpha-tubulin, mitotic chromosome segregation*, structural protein of cytoskeleton, spindle pole body
<b>YJR097W</b>		

## Experiment 10: Mic 20 hours

YHR007C	ERG11	cytochrome P450 lanosterol 14a-demethylase, ergosterol biosynthesis, lanosterol 14-alpha-demethylase, endoplasmic reticulum
YPL165C	UNK	UNK
YIL170W	HXT12	Hexose permease transport
YOR153W	PDR5, LEM1, YDR1	multidrug resistance transporter
YLR208W	Sec13	cytoplasmic protein involved in release of transport vesicles from the ER, non-selective vesicle assembly, molecular_function unknown, cytoplasm
YML085C	Tub1	alpha-tubulin, mitotic chromosome segregation*, structural protein of cytoskeleton, spindle pole body

## Experiment 13: Clot 20.1 hours

<b>YHR007C</b>	ERG11	cytochrome P450 lanosterol 14a-demethylase, ergosterol biosynthesis, lanosterol 14-alpha-demethylase, endoplasmic reticulum
<b>YOR153W</b>	PDR5, LEM1, YDR1	multidrug resistance transporter
<b>YPL165C</b>	UNK	UNK
<b>YIL170W</b>	HXT12	Hexose permease transport
<b>YML085C</b>	Tub1	alpha-tubulin, mitotic chromosome segregation*, structural protein of cytoskeleton, spindle pole body

## Experiment 14: Noc 17.3 hours

YML124C	TUB3	alpha-tubulin, mitotic chromosome segregation*, structural protein of cytoskeleton, spindle pole bod
YJL014W	CCT3, BIN2, TCP3	Cytoplasmic chaperonin subunit gamma, protein folding*, chaperone, cytoplasm
YJR097W	RPL27A	Ribosomal protein L27A, protein biosynthesis, structural protein of ribosome, cytosolic large ribosomal (60S)-subunit

# Simplified Data Upload

QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites History Print Mail News RSS Feeds

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer NIST Eng. Stats GoogleBerkeley

Address http://gobi.lbl.gov/~aparkin/Projects/Giaever/DTrials/Experiments/QuickRat/ Go

---

Bio/Spice::QuickRat HHMI  
University of California  
Lawrence Berkeley National Laboratory

---

**Welcome to the temporary Chip Ratioing Portal.**

Here you can upload chips, ask for ratio analyses and view past analyses.

Click [here](#) to request analysis of current chips.  
Click [here](#) to view past analyses.

Otherwise add a new chip using the form below.

## Add Chip

Name:

Comment:

Pool:

Chip:

Exp:

Done Internet

# Easy access to chips/requests for analysis

QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Search Favorites History Print Mail News RSS

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer NIST Eng. Stats GoogleBerkeley

Address http://gobi.lbl.gov/~aparkin/Projects/Giaever/DTrials/Experiments/QuickRat/cgi-bin/ShowChips.cgi Go

---

HHMI  
University of California  
Lawrence Berkeley National Laboratory

## Bio/Spice::QuickRat

---

### Chip Pages

To look at a single chip analysis click on the chip name.  
To request a ratio analysis select the first and second chip and press Analyze

### Current Chips

Chip Name	Homo/Hetero	Comment	CEL Data	EXP Data	First	Second
<a href="#">032801</a>	<a href="#">Heterozygous</a>	03_01 het pool sampling 1X template	<a href="#">[032801.cel]</a> <a href="#">[annot]</a>	<a href="#">[032801.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">041701</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue t = -70C	<a href="#">[04_17_01.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_17_01.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">041702</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue time = 0 (on)	<a href="#">[04_17_02.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_17_02.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">041703</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue time = 0 (on) first automated fluidics washing	<a href="#">[04_17_03.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_17_03.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">042401</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue time = 24hrs no drug T1 G = 15.3	<a href="#">[04_24_01.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_24_01.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">042402</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue t = 24hrs G = 14.4 ben 0.2ug/ml	<a href="#">[04_24_02.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_24_02.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">042403</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue t = 24hrs G = 15.2 noc 0.2ug/ml	<a href="#">[04_24_03.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_24_03.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">042404</a>	<a href="#">Heterozygous</a>	04_01_1 het pool, DMSO/glycerol issue t = 24hrs G = 14.9 ben 0.4ug/ml	<a href="#">[04_24_04.cel]</a> <a href="#">[annot]</a>	<a href="#">[04_24_04.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">050101</a>	<a href="#">Heterozygous</a>	04_25_01 het -70C	<a href="#">[05_01_01.cel]</a> <a href="#">[annot]</a>	<a href="#">[05_01_01.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">050102</a>	<a href="#">Heterozygous</a>	04_25_01 t = 0 (on) G = 10	<a href="#">[05_01_02.cel]</a> <a href="#">[annot]</a>	<a href="#">[05_01_02.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">050103</a>	<a href="#">Heterozygous</a>	04_01 het noc 4ug/ml 24hrs	<a href="#">[05_01_03.cel]</a> <a href="#">[annot]</a>	<a href="#">[05_01_03.exp]</a>	<input type="radio"/>	<input type="radio"/>
<a href="#">050802</a>	<a href="#">Heterozygous</a>	het 04_25_01 noc 2ug/ml 24hr G = 17	<a href="#">[05_08_02.cel]</a> <a href="#">[annot]</a>	<a href="#">[05_08_02.exp]</a>	<input type="radio"/>	<input type="radio"/>

Internet



# Automated Statistical Analysis/QC

QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

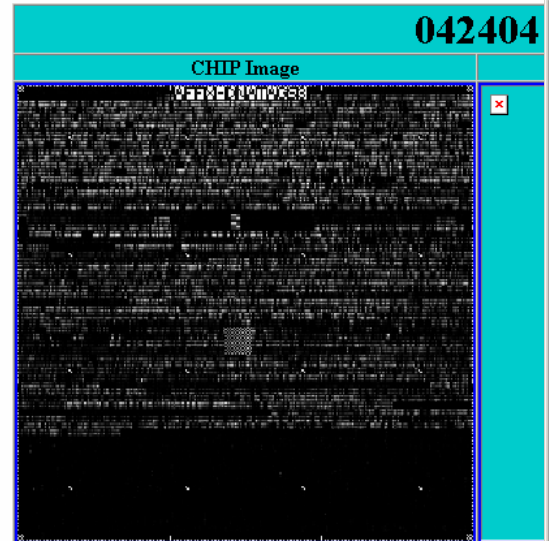
Back Forward Stop Home Search Favorites History Print Copy Paste

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer NIST Eng. Stats GoogleBerkeley

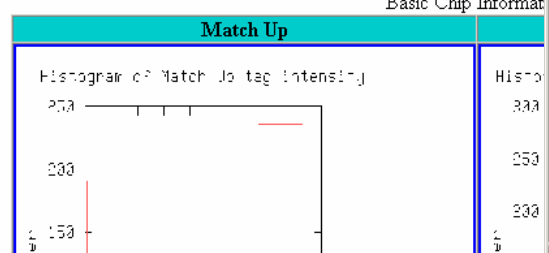
Address http://gobi.lbl.gov/~aparkin/Projects/Gaever/DTrials/Experiments/QuickRat/chips/042404/chip.html

## Bio/Spice::QuickRat Lawrence Berkeley

042404



04\_24\_04.cel 04\_24\_04.cel.annot  
[COMMENT](#) [ANN](#) [Heteroz](#)



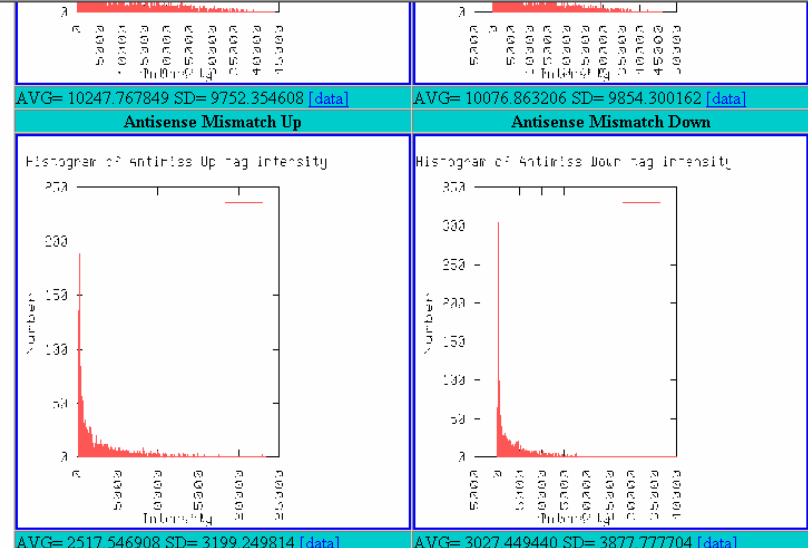
QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

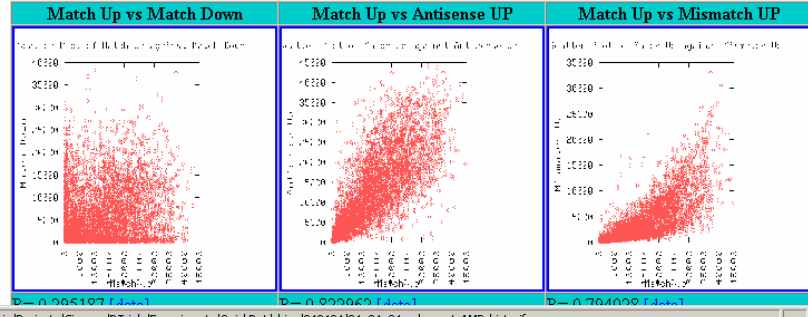
Back Forward Stop Home Search Favorites History Print Copy Paste

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer NIST Eng. Stats GoogleBerkeley

Address http://gobi.lbl.gov/~aparkin/Projects/Gaever/DTrials/Experiments/QuickRat/chips/042404/chip.html



Histograms for different tag types. Click on figure for full size version.



# Easy access to past analyses

QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Search Favorites History Print Mail News RSS Feeds

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer

Address http://gobi.lbl.gov/~aparkin/Projects/Giaever/DTrials/Experiments/QuickRat/cgi-bin/ShowAnalyses.cgi Go

---

Bio/Spice::QuickRat HHMI  
University of California  
Lawrence Berkeley National Laboratory

---

## Analysis Pages

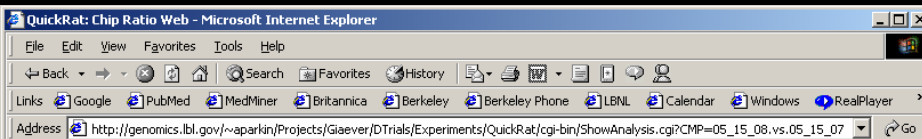
To look at a ratio analysis click on the analysis name.

### Current Analyses

Analysis	Chip1	Chip2
<a href="#">041702 vs 041703</a>	<a href="#">[04_17_02.cel]</a>	<a href="#">[04_17_02.exp]</a>
<a href="#">041702 vs 042401</a>	<a href="#">[04_17_02.cel]</a>	<a href="#">[04_17_02.exp]</a>
<a href="#">041702 vs 042403</a>	<a href="#">[04_17_02.cel]</a>	<a href="#">[04_17_02.exp]</a>
<a href="#">041702 vs 050103</a>	<a href="#">[04_17_02.cel]</a>	<a href="#">[04_17_02.exp]</a>
<a href="#">042401 vs 050103</a>	<a href="#">[04_24_01.cel]</a>	<a href="#">[04_24_01.exp]</a>
<a href="#">050101 vs 050102</a>	<a href="#">[05_01_01.cel]</a>	<a href="#">[05_01_01.exp]</a>
<a href="#">050102 vs 050802</a>	<a href="#">[05_01_02.cel]</a>	<a href="#">[05_01_02.exp]</a>
<a href="#">050902 vs 051105</a>	<a href="#">[05_09_02.cel]</a>	<a href="#">[05_09_02.exp]</a>
<a href="#">051101 vs 051102</a>	<a href="#">[05_11_01.cel]</a>	<a href="#">[05_11_01.exp]</a>
<a href="#">051101 vs 051103</a>	<a href="#">[05_11_01.cel]</a>	<a href="#">[05_11_01.exp]</a>
<a href="#">051603 vs 051604</a>	<a href="#">[05_16_03.cel]</a>	<a href="#">[05_16_03.exp]</a>
<a href="#">051603 vs 052202</a>	<a href="#">[05_16_03.cel]</a>	<a href="#">[05_16_03.exp]</a>
<a href="#">052203 vs 052204</a>	<a href="#">[05_22_03.cel]</a>	<a href="#">[05_22_03.exp]</a>
<a href="#">052205 vs 052206</a>	<a href="#">[05_22_05.cel]</a>	<a href="#">[05_22_05.exp]</a>
<a href="#">05_15_01 vs 05_15_02</a>	<a href="#">[05_15_01.CEL]</a>	<a href="#">[05_15_02.CEL]</a>
<a href="#">05_15_08 vs 05_15_07</a>	<a href="#">[05_15_08.CEL]</a>	<a href="#">[05_15_07.CEL]</a>
<a href="#">Adam1 vs Adam2</a>	<a href="#">[080101.CEL]</a>	<a href="#">[080102.CEL]</a>
<a href="#">Adam1 vs Adam3</a>	<a href="#">[080101.CEL]</a>	<a href="#">[080901.CEL]</a>
<a href="#">Adam1 vs Fake</a>	<a href="#">[080101.CEL]</a>	<a href="#">[080101_s.CEL]</a>
<a href="#">Adam1 vs Test</a>	<a href="#">[080101.CEL]</a>	<a href="#">[TPZ072600.CEL]</a>
<a href="#">Adam2 vs Fake</a>	<a href="#">[080102.CEL]</a>	<a href="#">[080101_s.CEL]</a>
<a href="#">Adam2 vs Test2</a>	<a href="#">[080102.CEL]</a>	<a href="#">[TPZ0815_3.CEL]</a>
<a href="#">Adam3 vs Adam4</a>	<a href="#">[080901.CEL]</a>	<a href="#">[081100.CEL]</a>
<a href="#">Guri1 vs Guri2</a>	<a href="#">[04_17_01.cel]</a>	<a href="#">[04_17_01.exp]</a>
<a href="#">Guri2 vs Guri3</a>	<a href="#">[04_17_02.cel]</a>	<a href="#">[04_17_02.exp]</a>

Internet

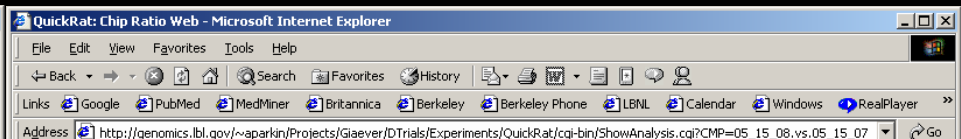
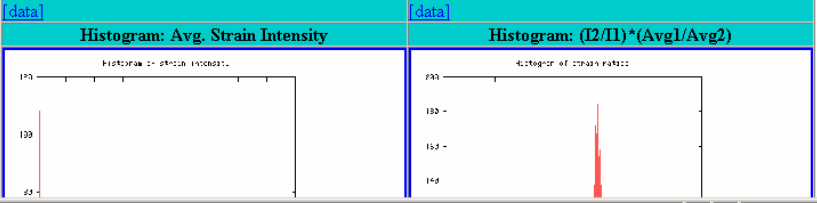
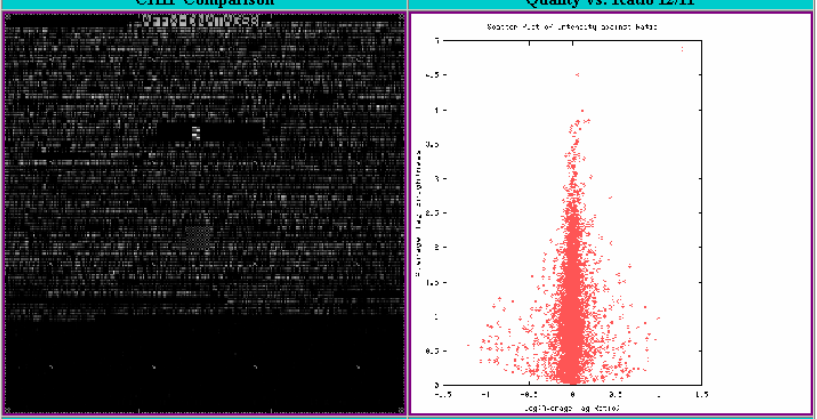
# Automated Analysis/Target Hyp.



Bio/Spice::QuickRat HHMI  
University of California  
Lawrence Berkeley National Laboratory

## Analysis Pages

<b>05 15 08</b>	<b>05 15 07</b>
Chip 1 Avg= 7451.722395	Chip 2 Avg= 6998.749015



QuickRat: Chip Ratio Web - Microsoft Internet Explorer

File Edit View Favorites Tools Help

Back Forward Stop Home Search Favorites History

Links Google PubMed MedMiner Britannica Berkeley Berkeley Phone LBNL Calendar Windows RealPlayer

Address http://genomics.lbl.gov/~aparkin/Projects/Gaever/DTrials/Experiments/QuickRat/cgi-bin/ShowAnalysis.cgi?CMP=05\_15\_08.vs.05\_15\_07

REGEXP:  Brightness Range:  L  H Ratio Range:  L  H Quality Range:  L  H

Sort By:  Ratio  Intensity  Quality  None  A  D Limit to results:

Query Results

## QUERY RESULTS

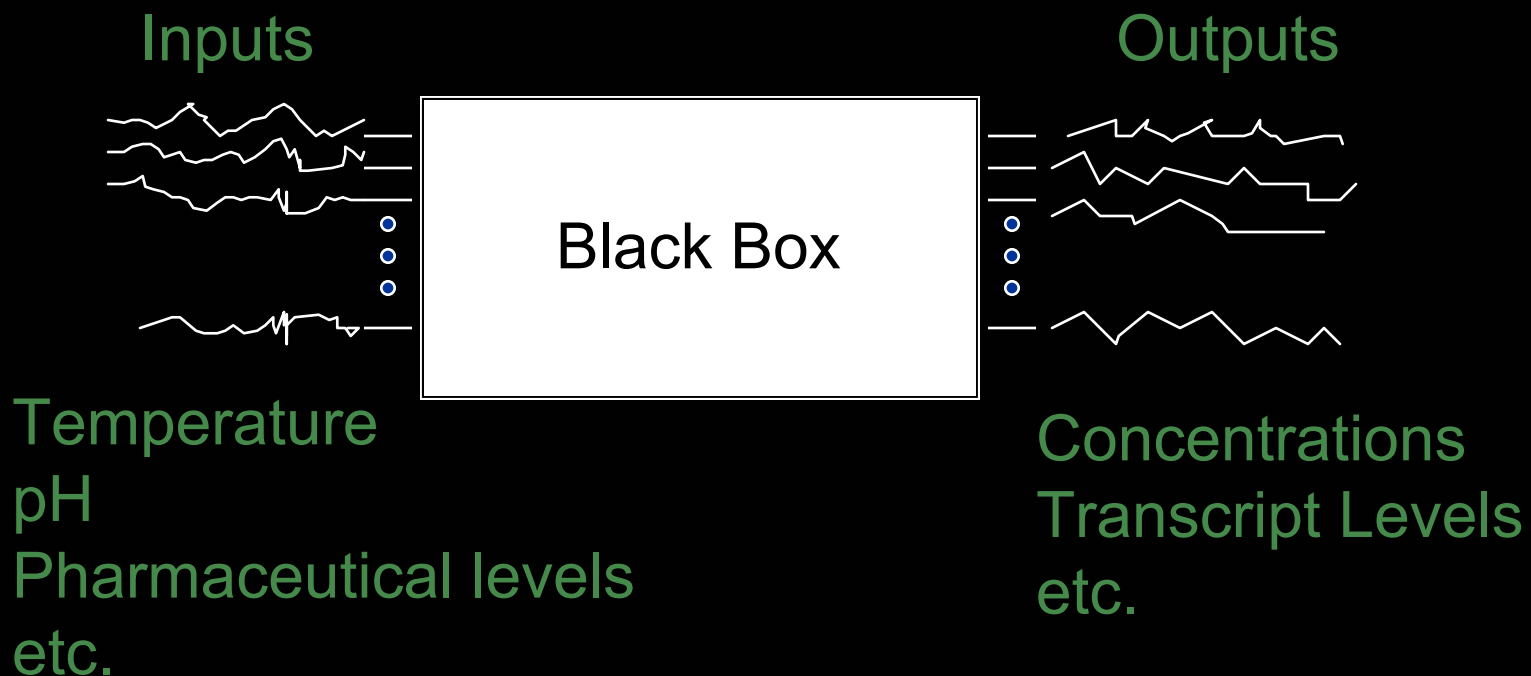
Strain	Ratio	Intensity	Gene	Quality	NTags	Viability
YML124C	-1.211994	0.569563	tub3	8.213675	4	INESSENTIAL
alpha-3 tubulin						
alpha-tubulin						
YCR036W	-1.067519	0.344817	REK1	7.271917	4	INESSENTIAL
rbokinase						
rbokinase						
YHR156C	-1.058958	0.538422	7.181411	4		INESSENTIAL
weak similarity to mouse kinesin KIF3B						
YCR028C	-1.032899	0.652064	FEN2	6.990636	4	INESSENTIAL
similarity to allantoin permease transporter						
Amino acid permease						
YCR085W	-1.031515	0.622995	6.984592	4		INESSENTIAL
hypothetical protein						
YGR029W	-0.985959	0.873954	ERV1	6.654740	4	ESSENTIAL
mitochondrial biogenesis and regulation of cell cycle						

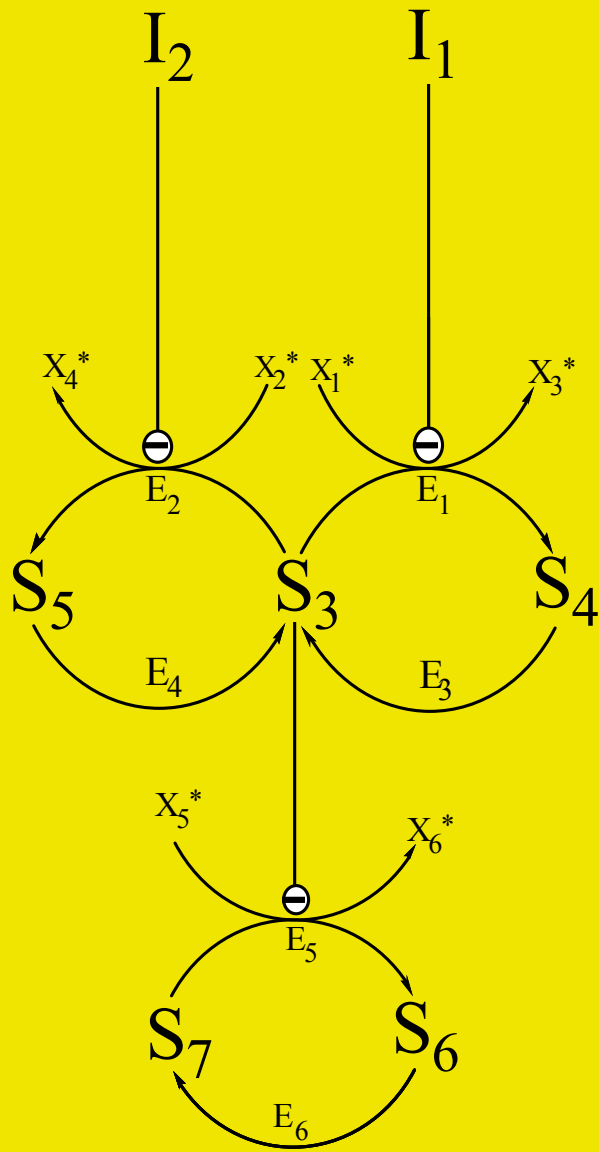
# Network Deduction and Modules

An incredibly naïve approach

# Correlation Metric Construction

A method to deduce reaction pathways directly from concentration time-series measurements

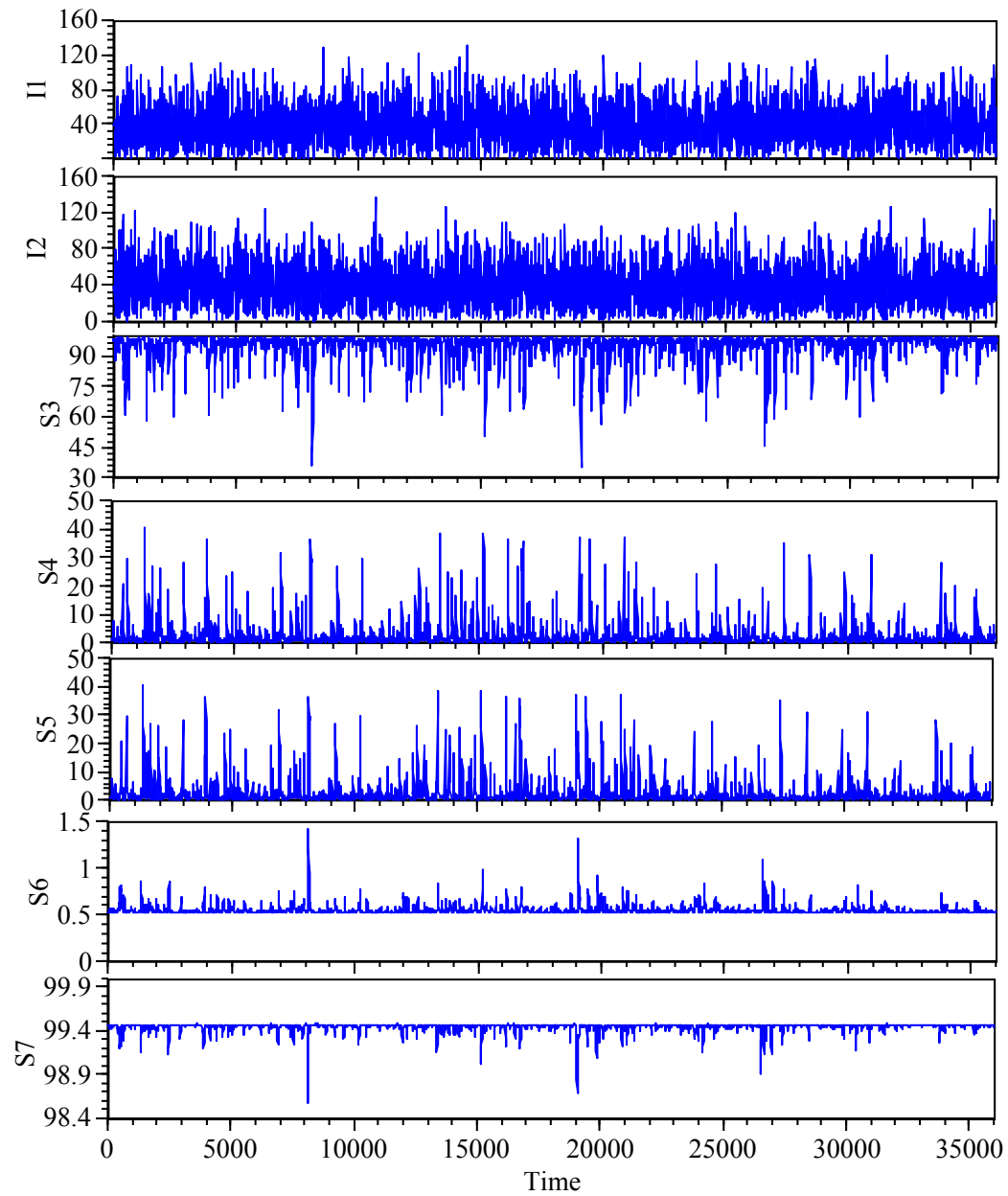




This is an abstract biochemical NAND gate

Based on a mechanism at the heart of switching between the glycolytic and gluconeogenic modes of the hexose phosphate pathway.

Can we deduce the network structure by perturbing it with inflows of the inputs and measuring the response of the concentrations?



# Measures of Dependency

## Linear Correlation

- o Measures the linear relationship between variables.
- o May be extended to multiple dependencies (i.e.  $y = f(x_1, x_2, x_3, \dots)$ ) by assumption of a linear regression model.
- o Very difficult to tell the significance of a given correlation since no distributional assumptions are made.

## Non-Parametric Rank Correlation (Spearman)

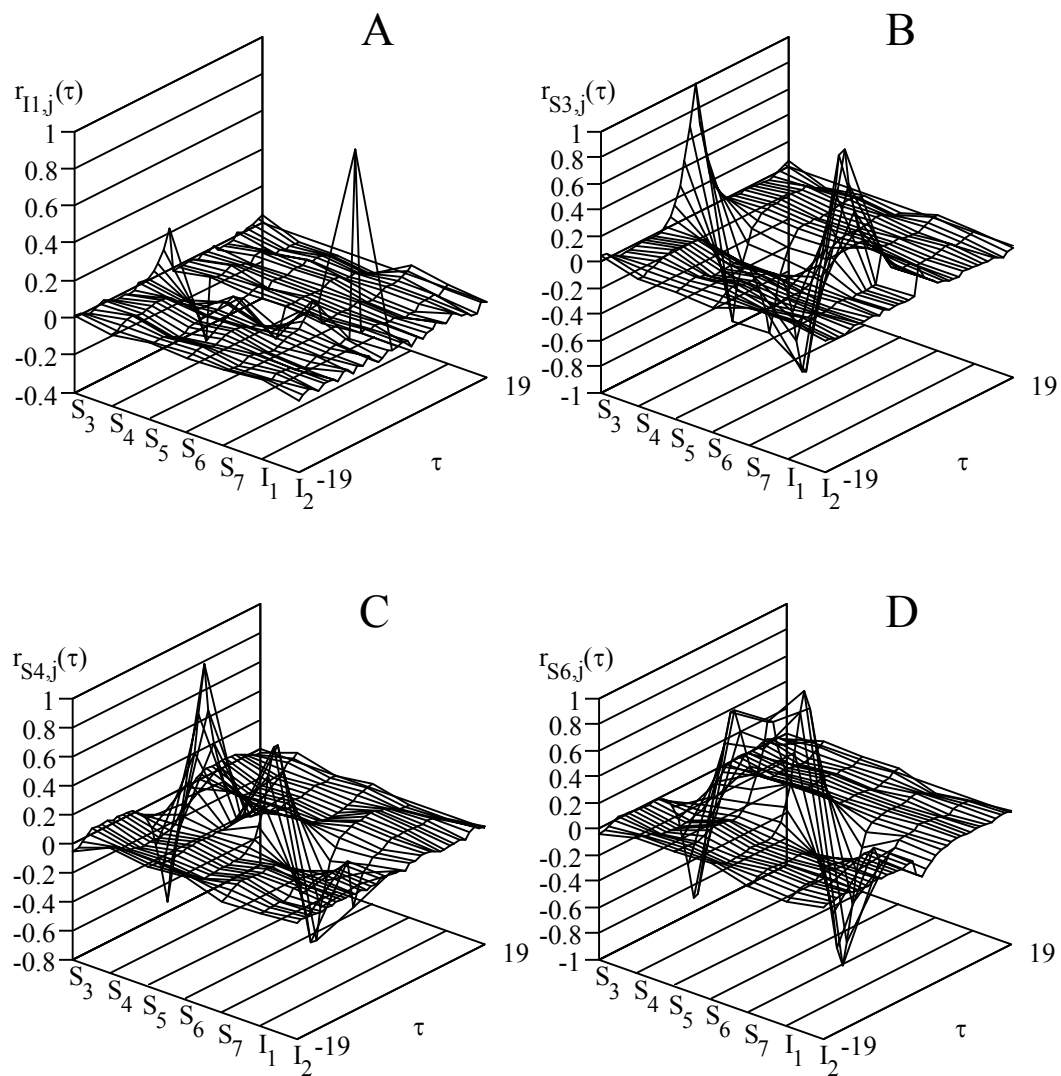
- o Measures monotonic relationships between variables.
- o Like Linear Correlation except that distribution of numbers is now known (uniform, exactly).
- o Robust to data defects.
- o Significances may be calculated...weak dependencies may be missed.

## Transinformation

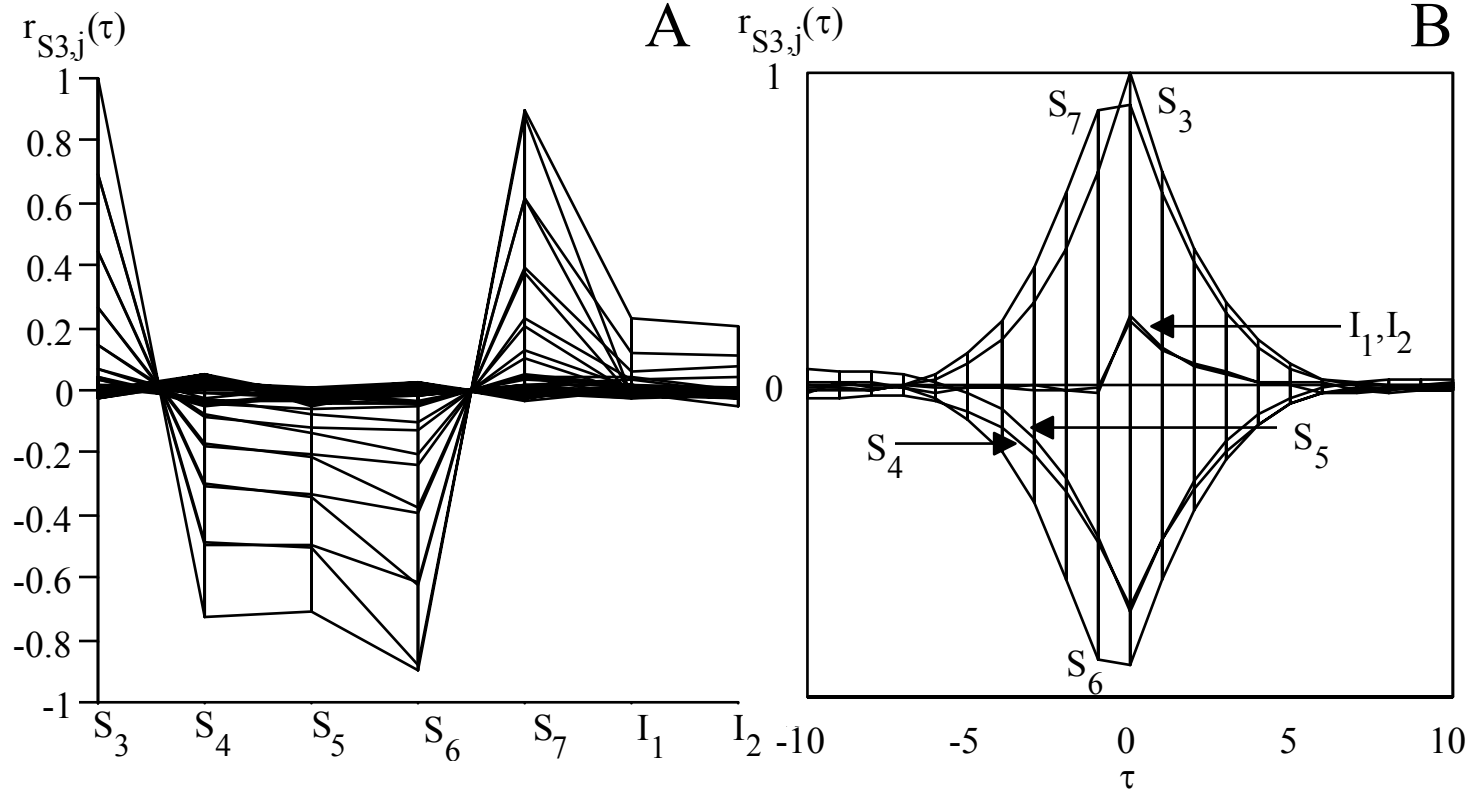
- o Measures the constraint on one variable given knowledge of another; i.e., requires that one variable merely be a function of the other.
- o The distribution is the quantity actually calculated.  $c_2$  then provides an accurate measure of significance.
- o Multiple dependencies easily incorporated by increasing the dimension of the distribution.



# Chemical NAND Gate Correlation Functions



# Correlation Function Projection



# Multidimension Scaling Solutions

A) Eigenvectors,  $z_k$ :

point/ $z_k$	1	2	3	4	5	6	7
1 ( $I_1$ )	6.68e-01	-5.84e-01	4.05e-01	5.51e-02	1.77e-02	-8.20e-09	-3.52e-10
2 ( $I_2$ )	7.00e-01	5.26e-01	-4.30e-01	4.93e-02	1.42e-02	-8.20e-09	-3.52e-10
3 ( $S_7$ )	-4.20e-01	7.29e-03	-8.16e-03	2.05e-01	1.90e-03	-6.82e-09	-7.67e-09
4 ( $S_6$ )	-4.20e-01	7.29e-03	-8.16e-03	2.05e-01	1.90e-03	-9.58e-09	6.97e-09
5 ( $S_4$ )	-1.44e-01	-5.51e-01	-4.02e-01	-1.60e-01	-7.55e-02	-8.20e-09	-3.52e-10
6 ( $S_5$ )	-7.15e-02	5.60e-01	4.30e-01	-1.38e-01	-7.65e-02	-8.20e-09	-3.52e-10
7 ( $S_3$ )	-3.14e-01	3.49e-02	1.27e-02	-2.16e-01	1.16e-01	-8.20e-09	-3.52e-10

B) Eigenvalues,  $l_k$ :

k	$l_k$	$a_{1,k}$	$a_{2,k}$
1	1.413496e+00	3.978311e-01	4.937104e-01
2	1.237497e+00	7.461268e-01	8.721279e-01
3	6.958617e-01	9.419783e-01	9.917824e-01
4	1.805556e-01	9.927960e-01	9.998381e-01
5	2.559576e-02	1.000000e+00	1.000000e+00
6	4.747935e-16	1.000000e+00	1.000000e+00
7	1.080736e-16	1.000000e+00	1.000000e+00

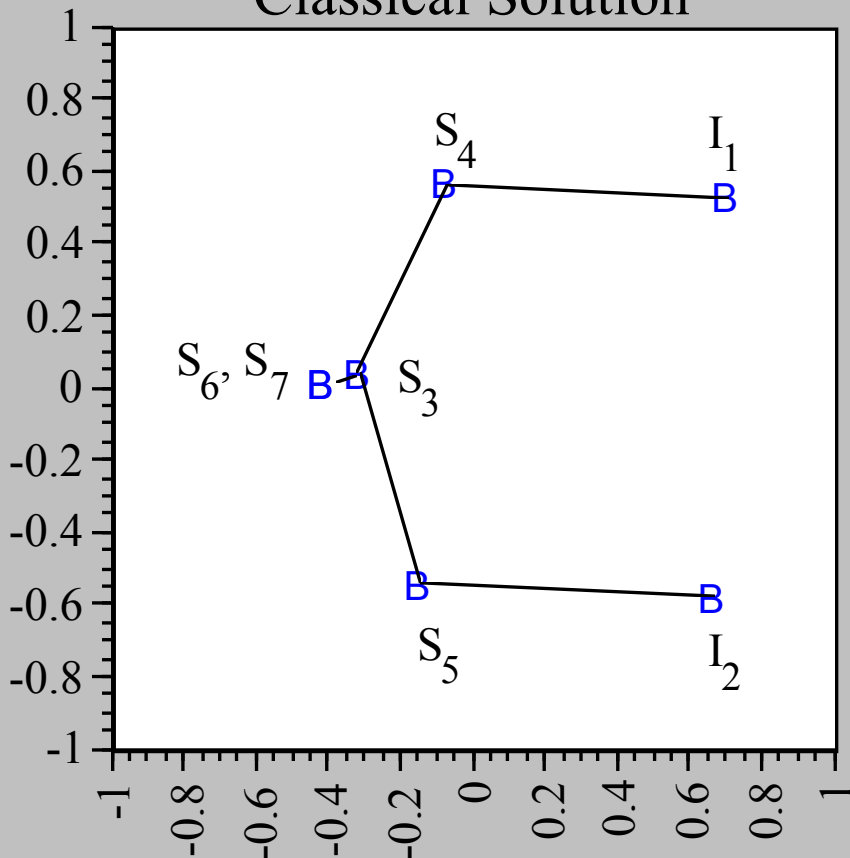
C) Significant Connections

	$S_4$	$S_5$	$S_6$	$S_7$
$I_1$	-0.31			
$I_2$		-0.32		
$S_3$	-0.72	-0.71	-0.90	0.90
$S_6$				-1.00

# MDS Solutions: Projection and Squash

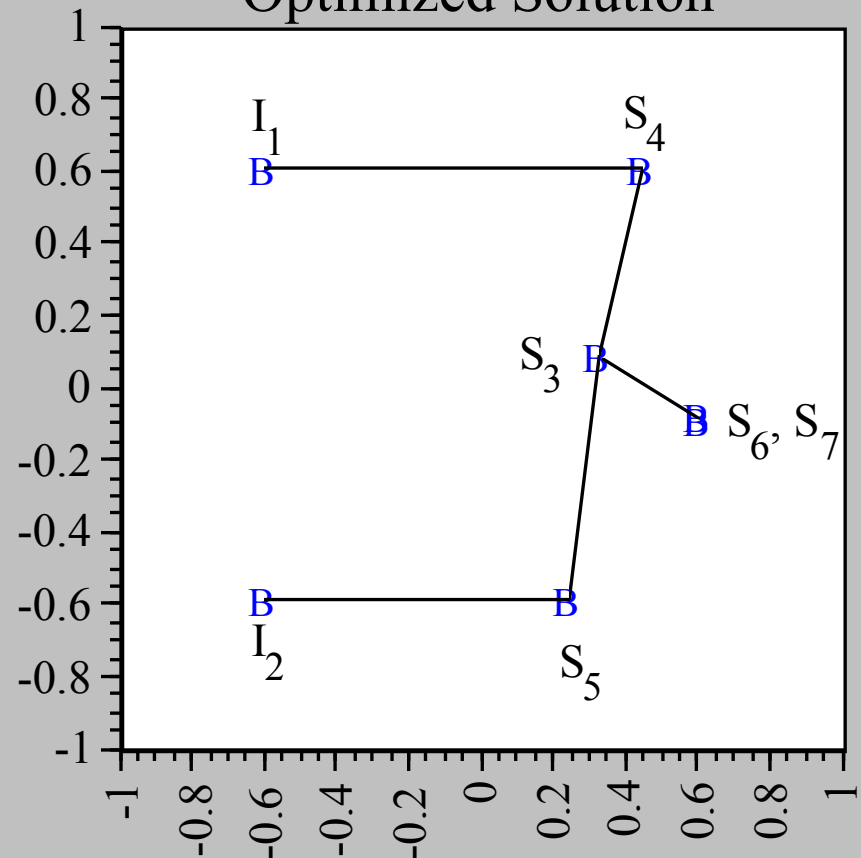
A

Classical Solution

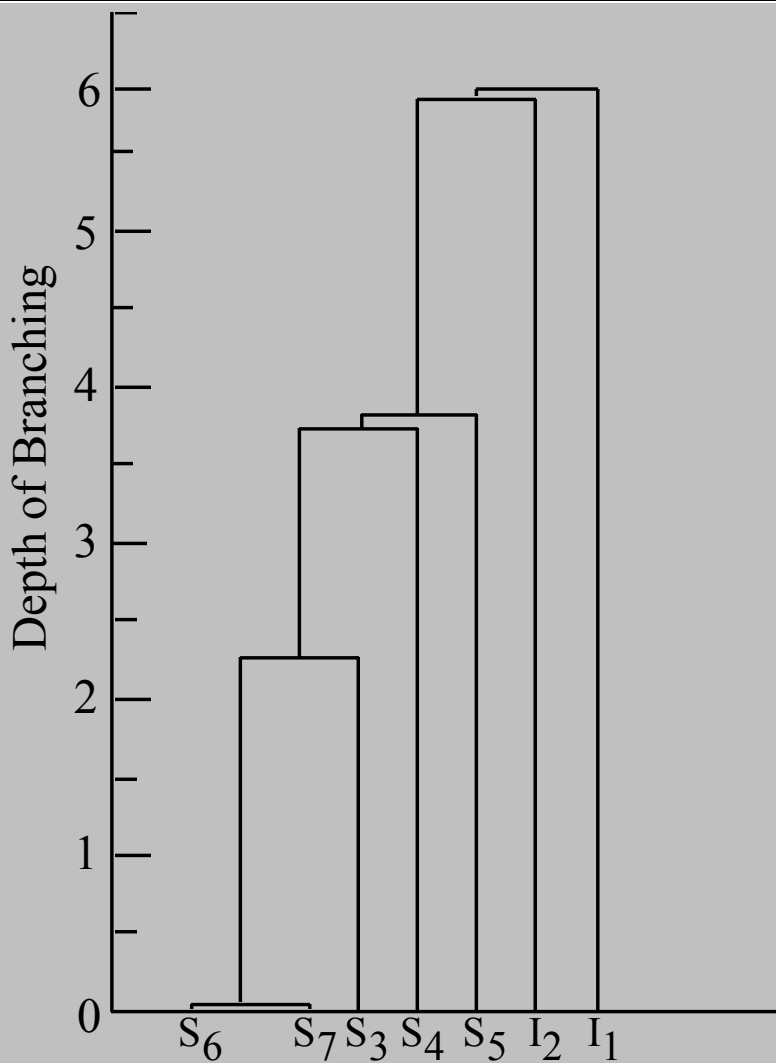


B

Optimized Solution



# Hierarchical Cluster Diagram



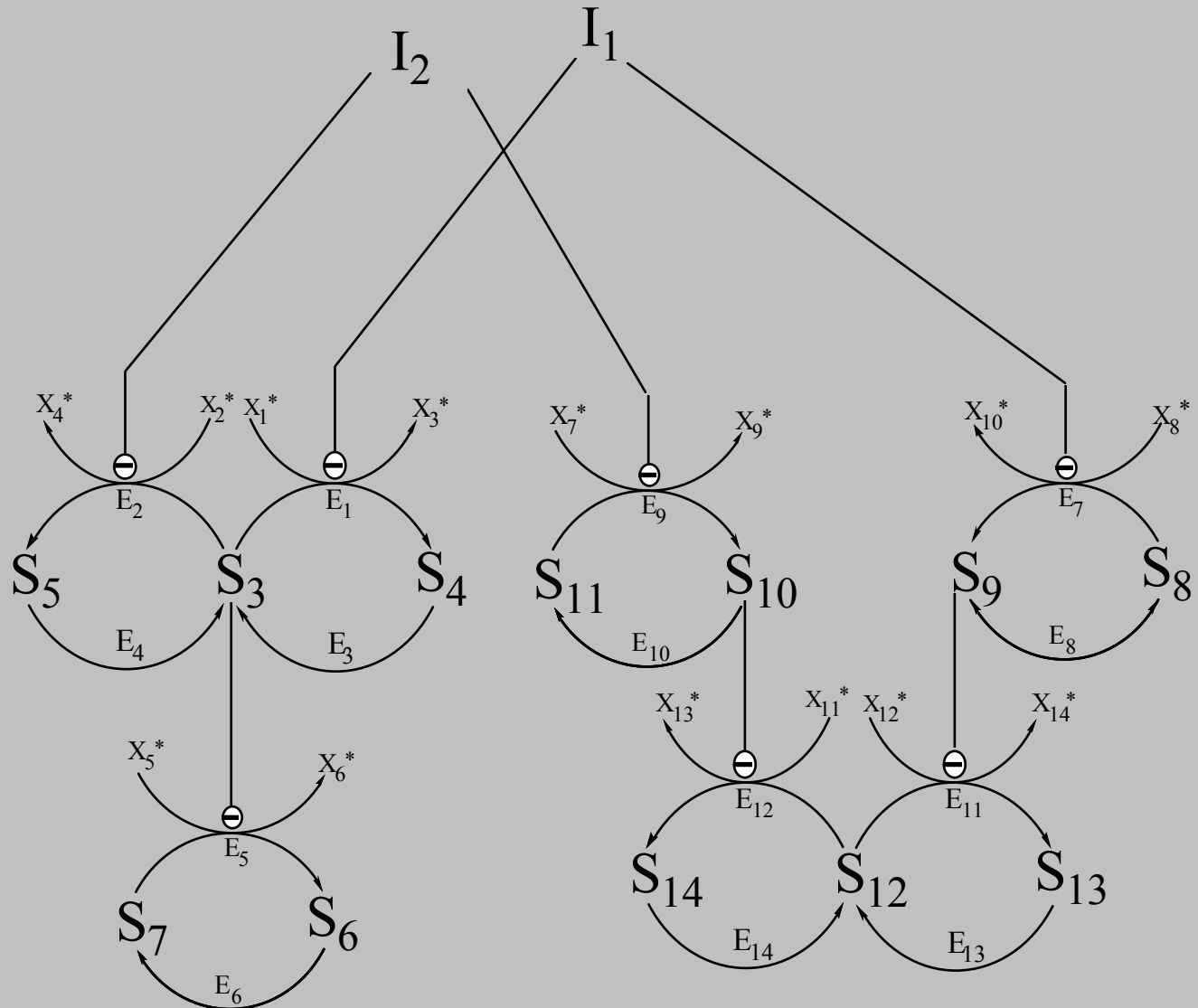
Tightly coupled species are grouped together.

Methodology for determining tightly coupled pathways?

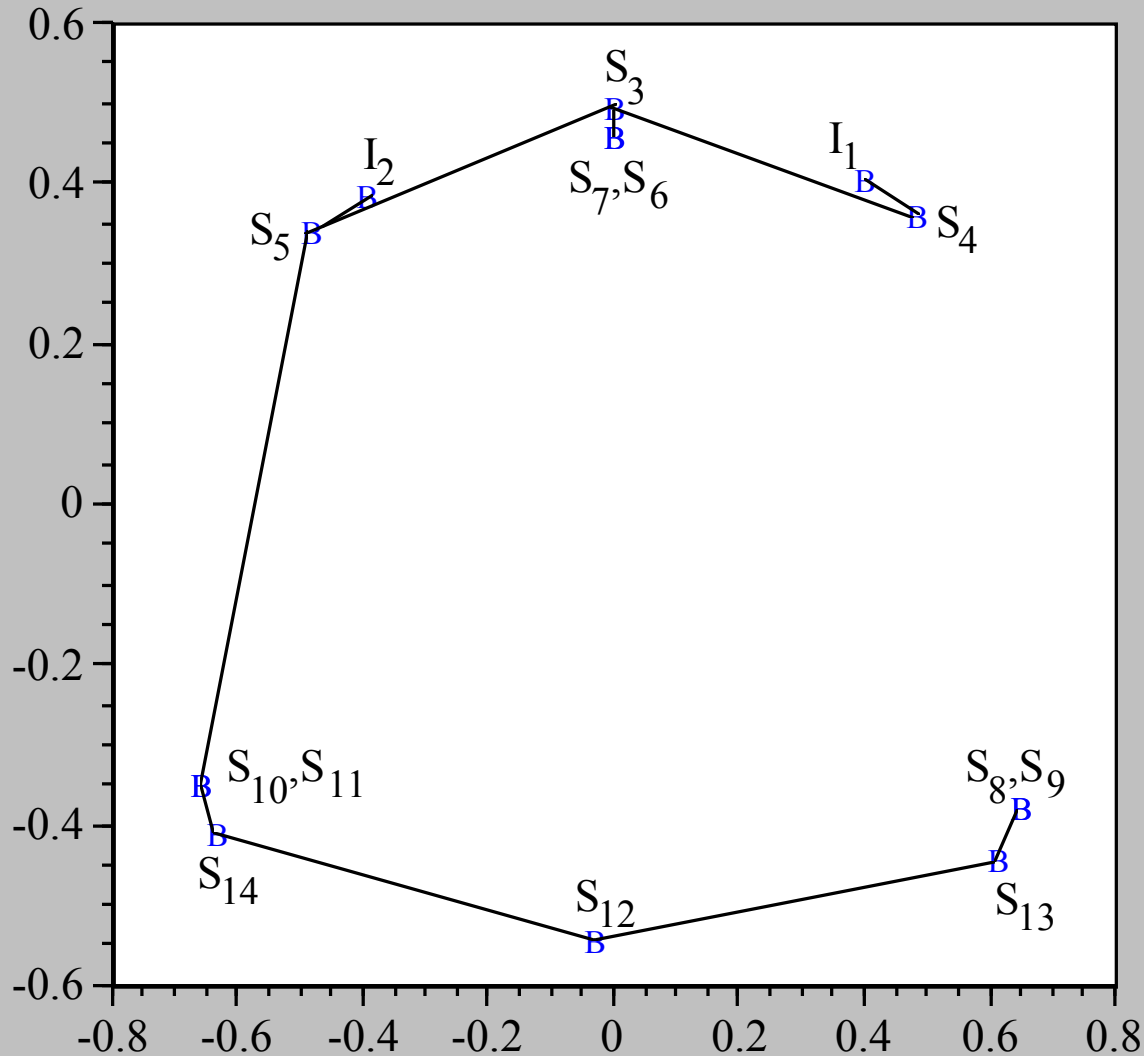
Or chemical subsystems that may be analyzed outside of the rest of the circuit?

**MORE ON MODULARITY:  
STAY TUNED!**

# A more complex case



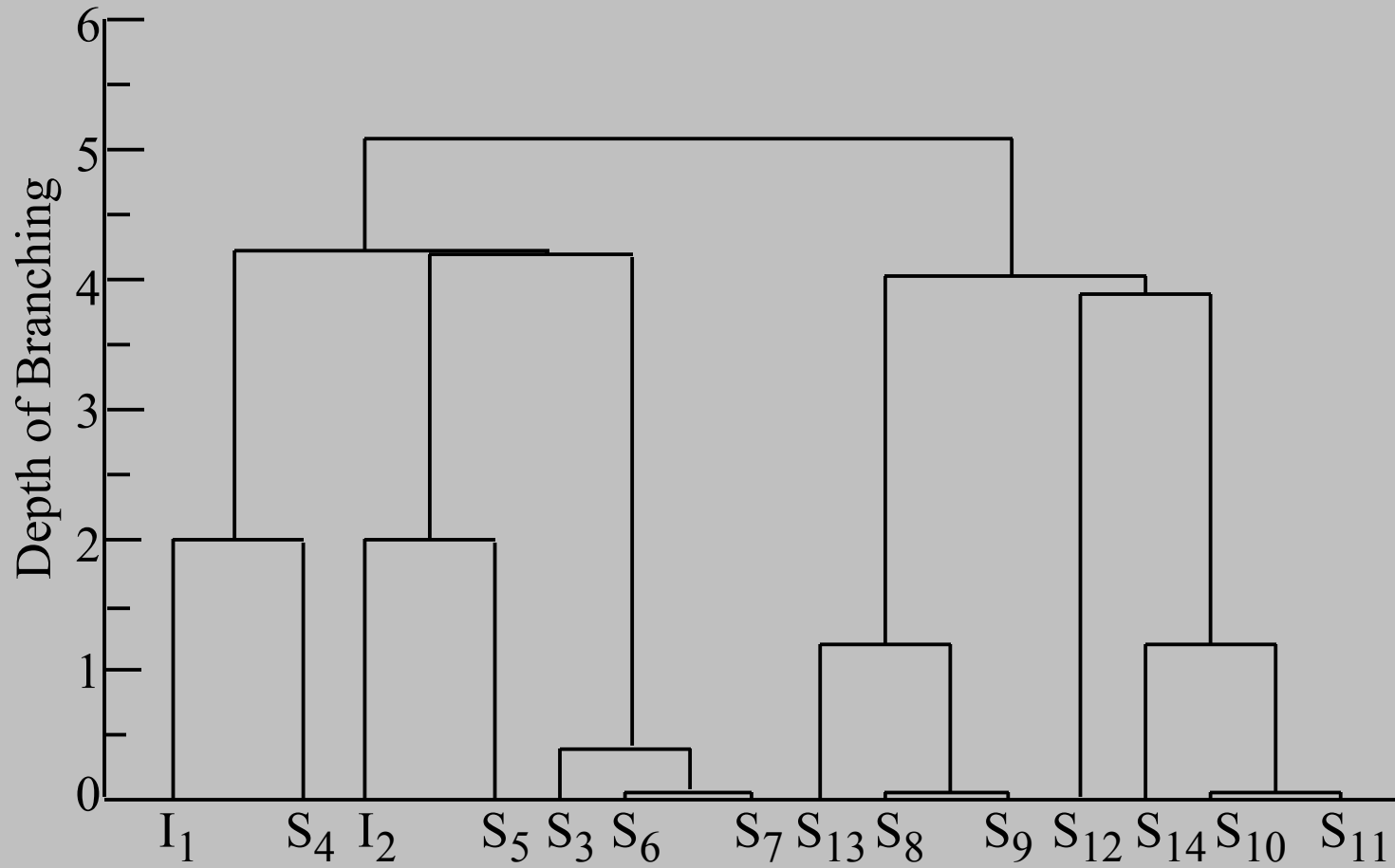
# The MDS Diagram



Note that the two parallel pathways groups on different sides of the diagram

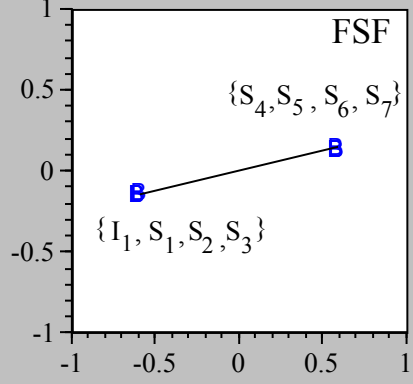
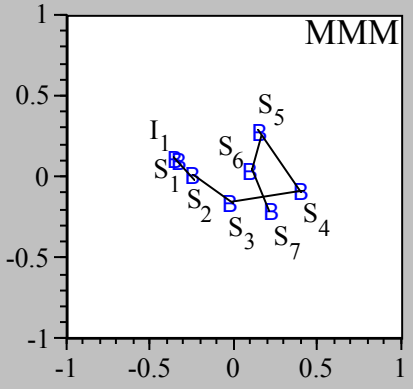
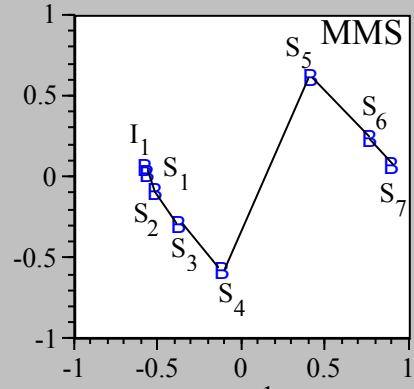
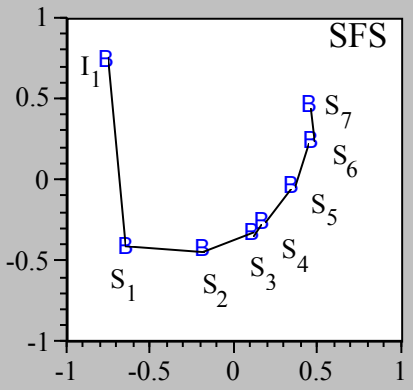
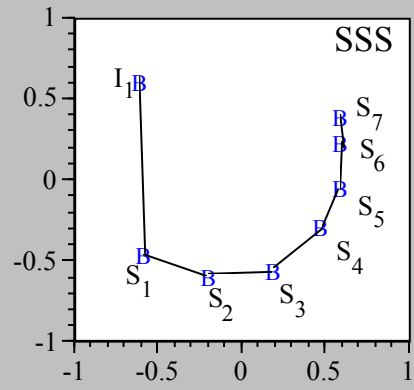
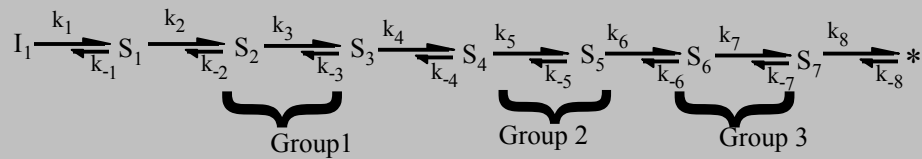
But also note the spurious connection between S5 and S10,11.

# Hierarchical Cluster

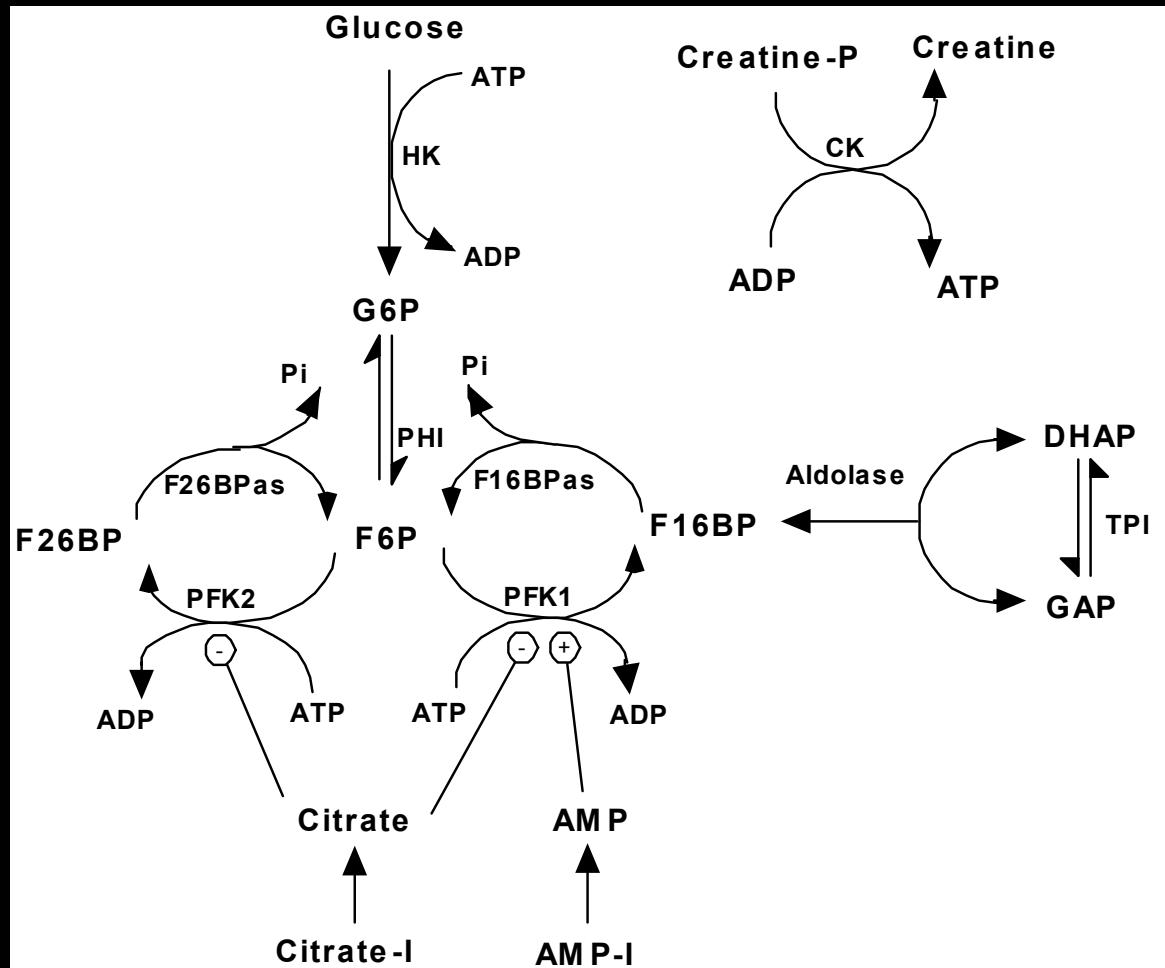




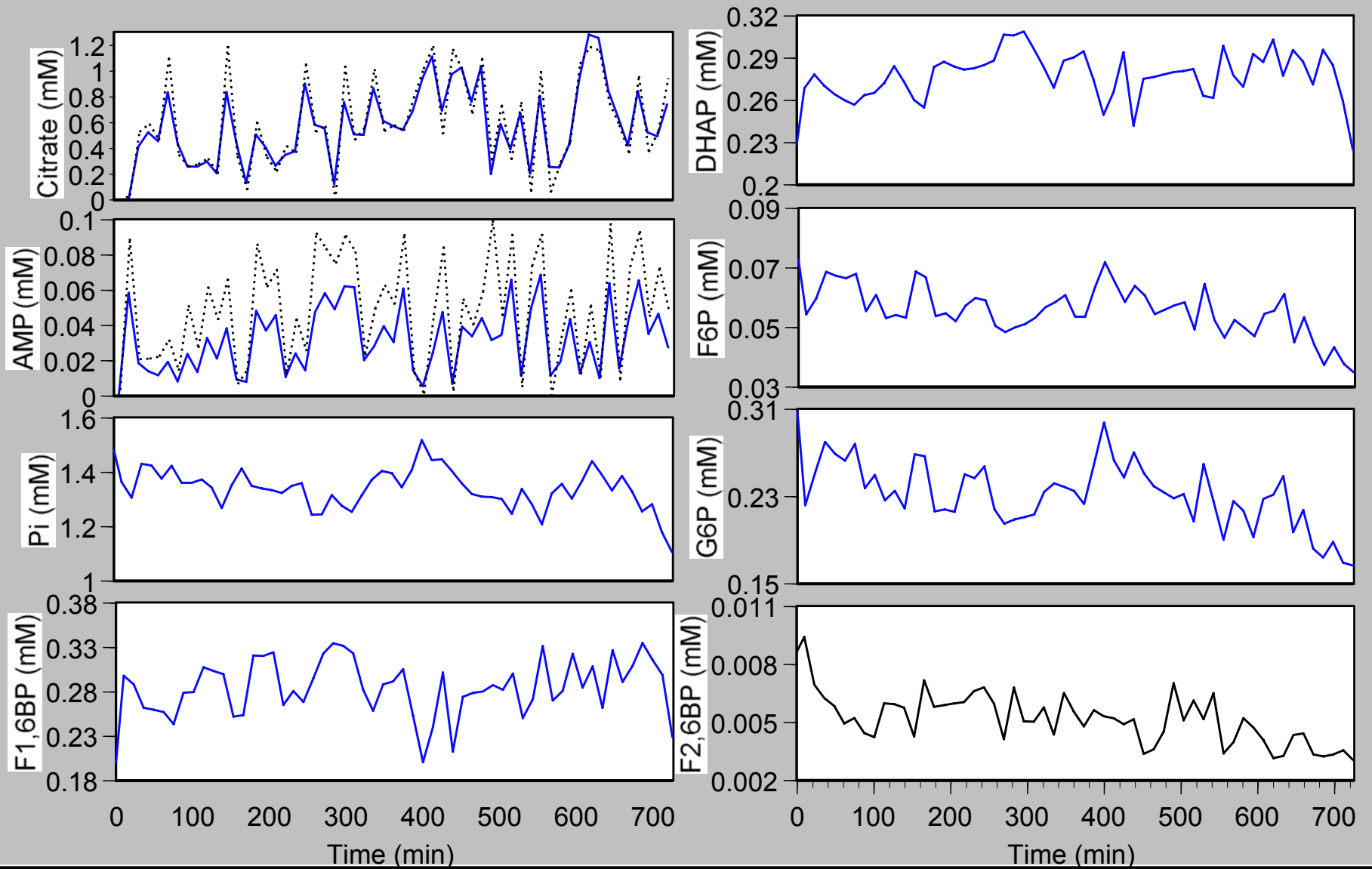
# Reaction Chain Cases

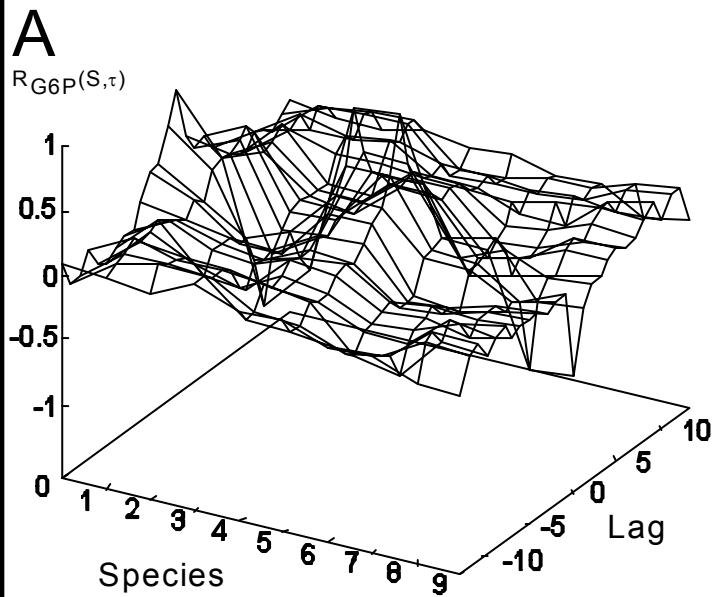


# Experimental Test System

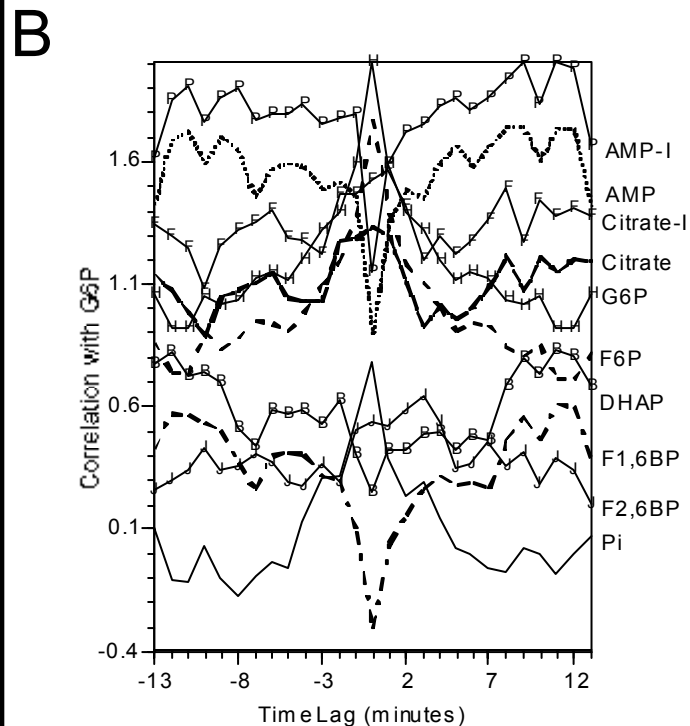


# Capillary Electrophoresis Data



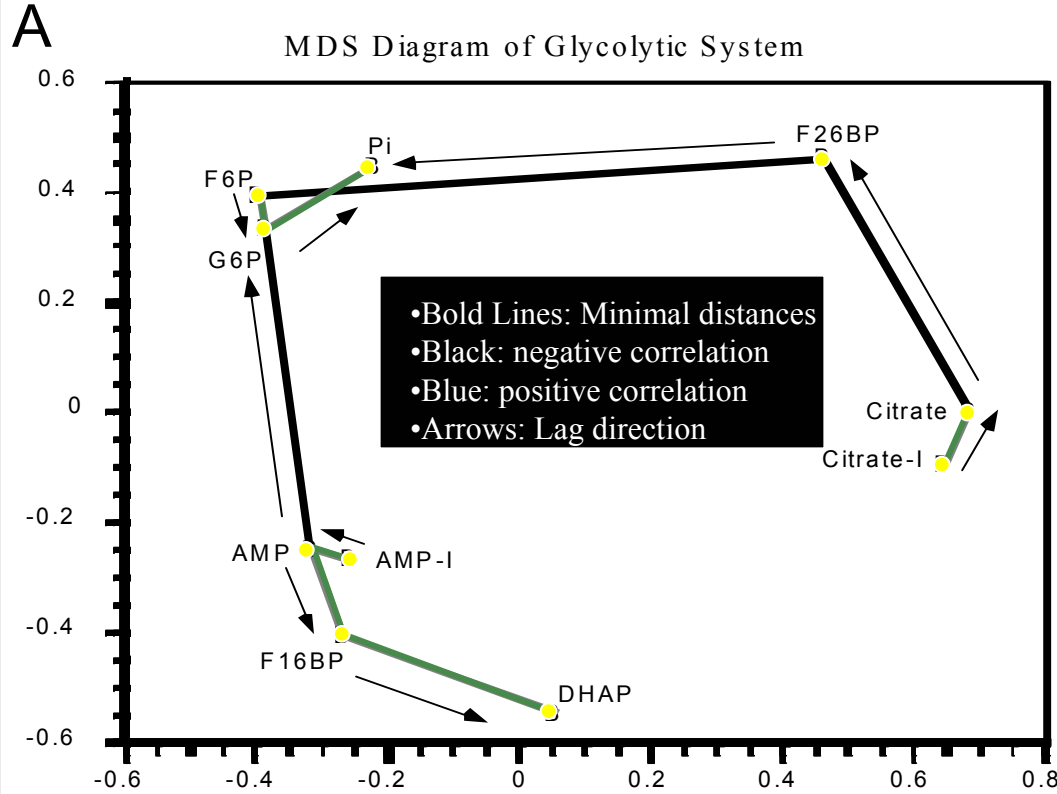


Experimental correlation of all 10 measured glycolytic intermediates with G6P.



Both strength of interaction and temporal ordering are implicit in time-dependent correlation function

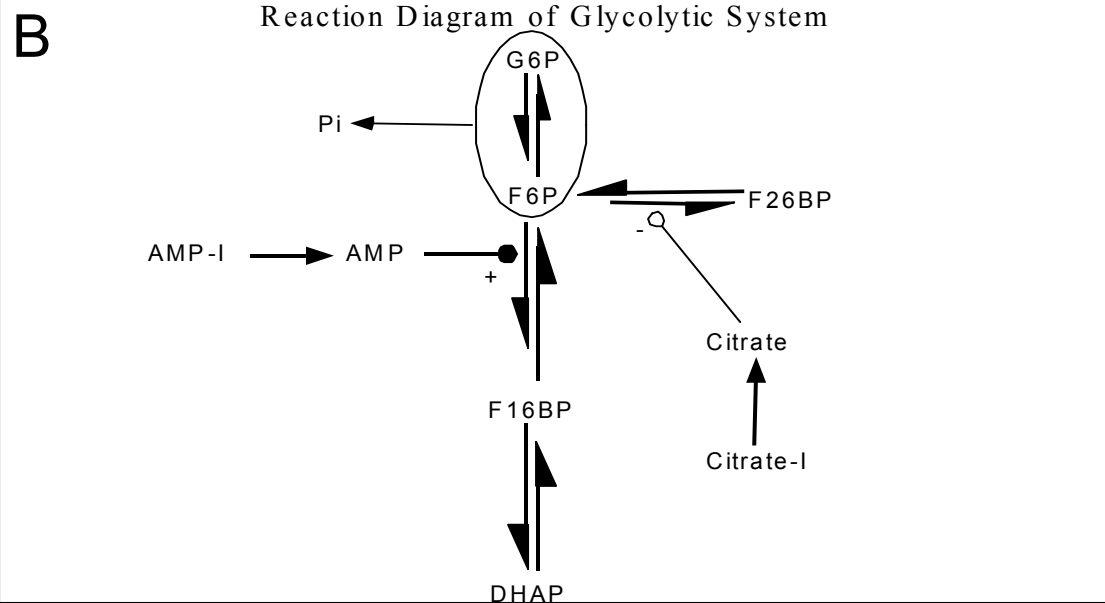
Projection of the surface down the time-lag axis. The results of each species are offset a bit from zero.



Multidimensional Scaling (MDS) diagrammatic summary of correlation functions.

Diagram summarizes:

- Interaction strengths
- Temporal ordering
- Probable network structure



“Expert system” prediction of reaction pathway from correlation matrix/MDS analysis

Correlation misses relationships such as this. Information is a better metric.

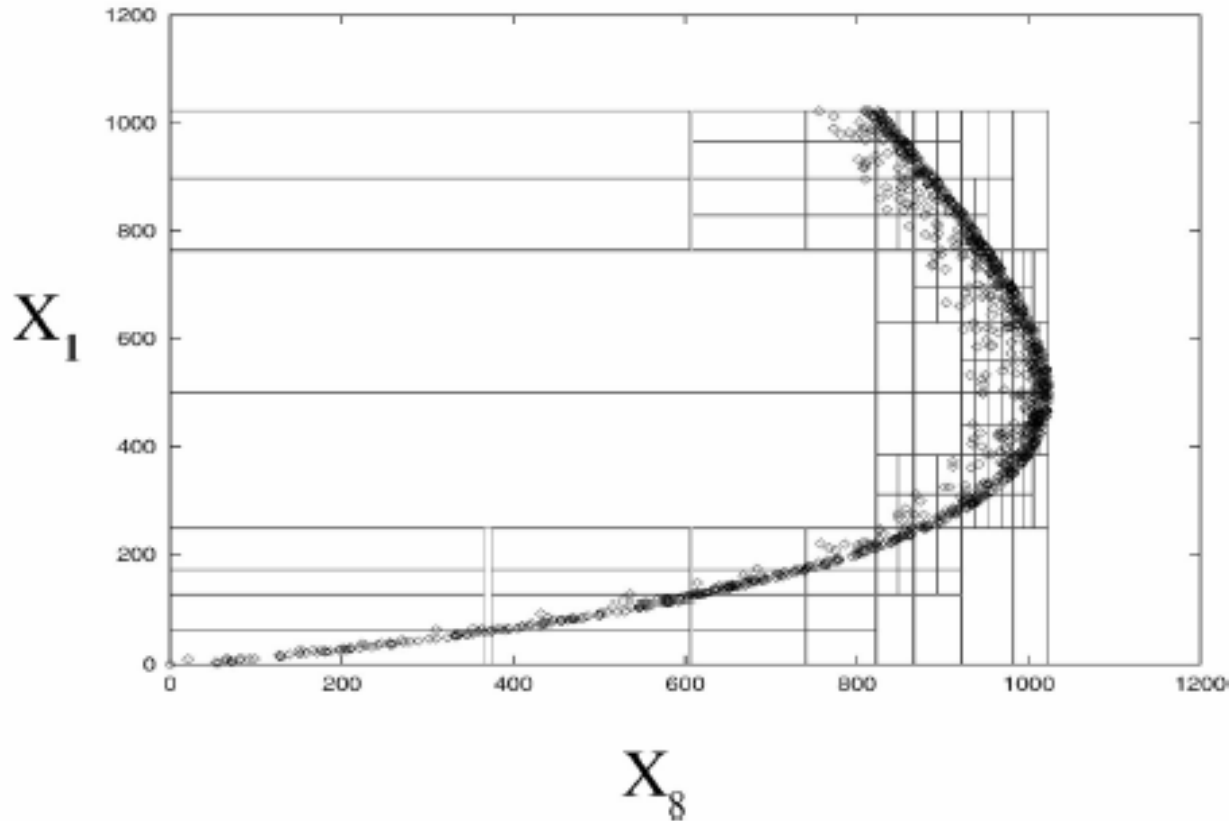


FIG. 5. The diamonds plot the values of  $X_1$  vs  $X_8$  obtained from the simulated time series. The rectangles are the result of a partitioning algorithm, see the text. From Ref. 32.

Samoilov, Arkin, Ross, *Chaos*, 11(1):108-114

Data problems: Theoretically  $N = 5 \cdot Q^M$

# Some short thoughts on modules

Considerations

# How do we define modules?

1. Repeated units
2. Evolutionary conservation
3. Time scale separation
  - a. fast and slow manifolds (Michaelis Menten)
  - b. on at different times
4. Spatial/Structural separation
  - a. no links between two subsystems
  - b. separated in space literally
  - c. weak coupling
5. Individually controllable
  - a. high impedances between subsystems?

Elliott will define more!



# Schematics for cells: Modules

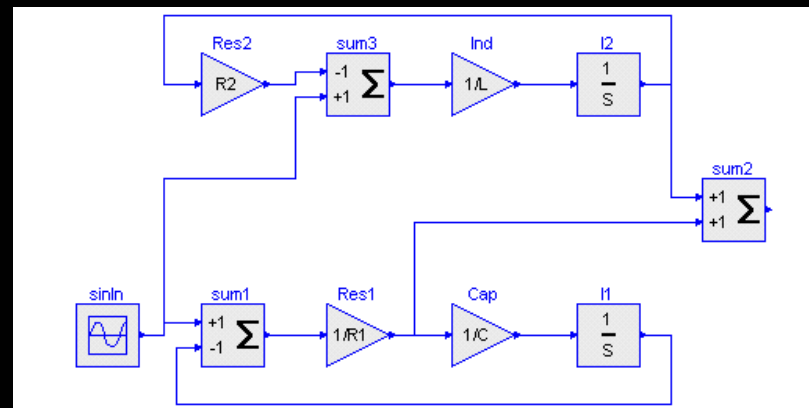
Motif= A repeated pattern of interactions among objects

Module

1. A motif

2. An elementary functional unit

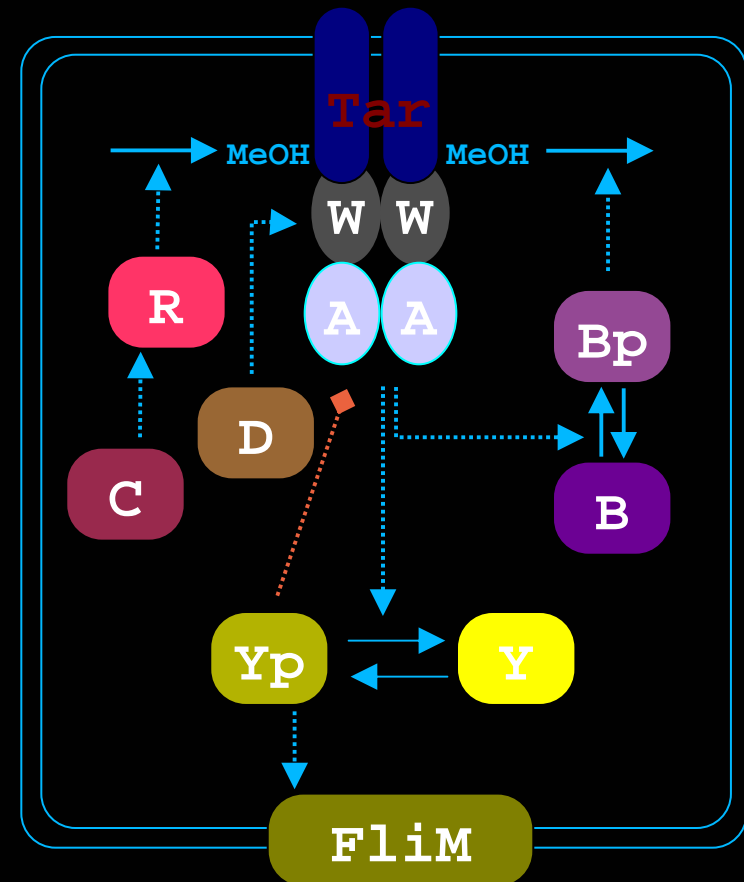
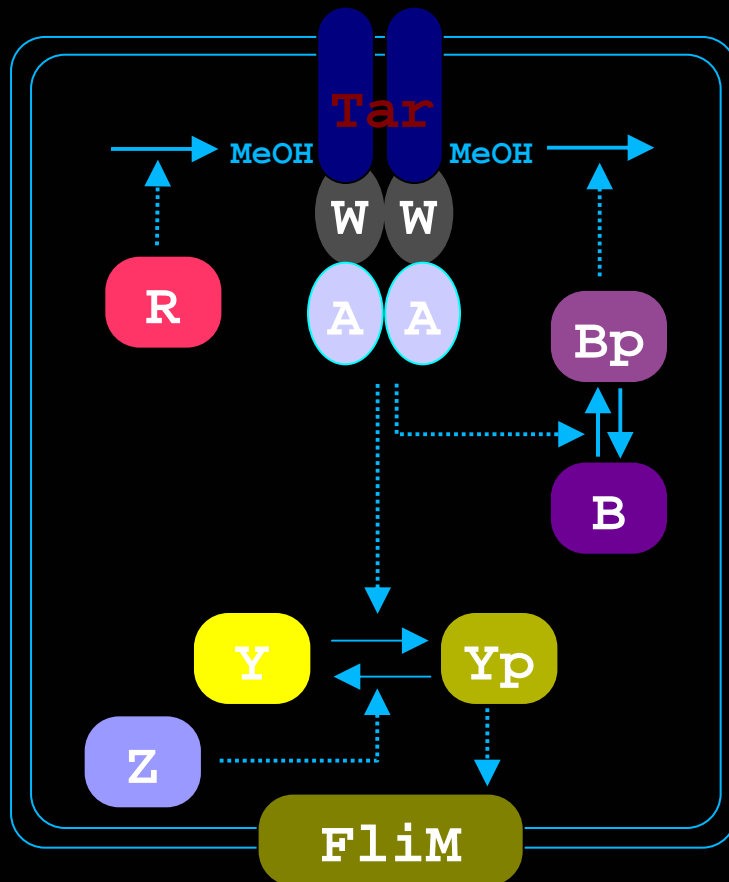
3. A compound functional unit that may be abstracted



# Chemotaxis Signal Transduction Pathways. A motif? A module?

*E. coli*

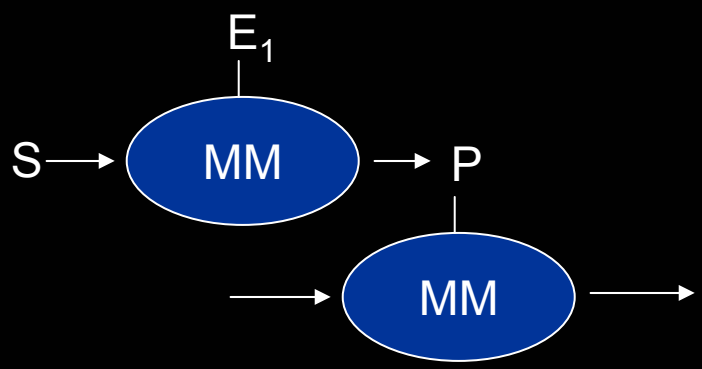
*B. subtilis*





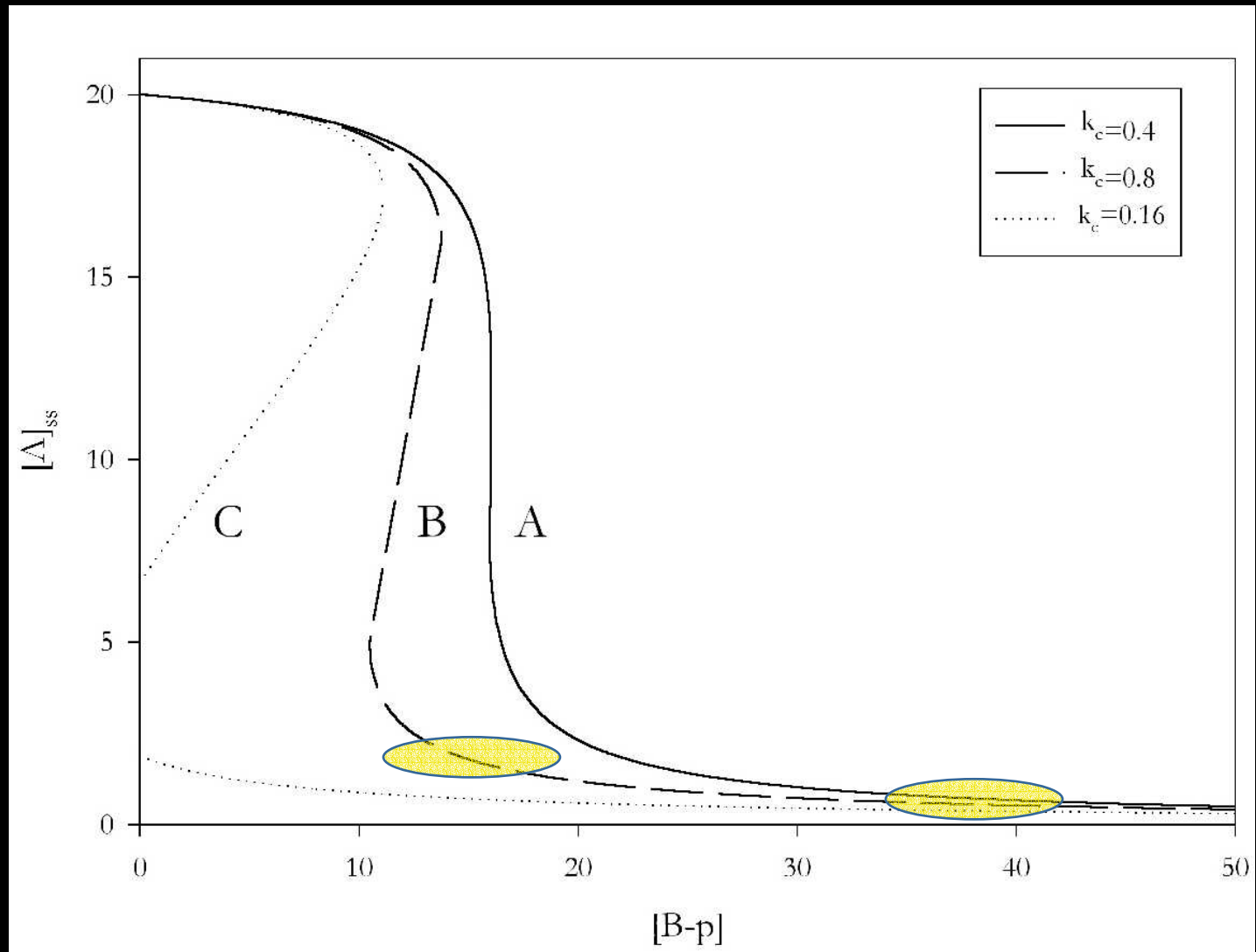
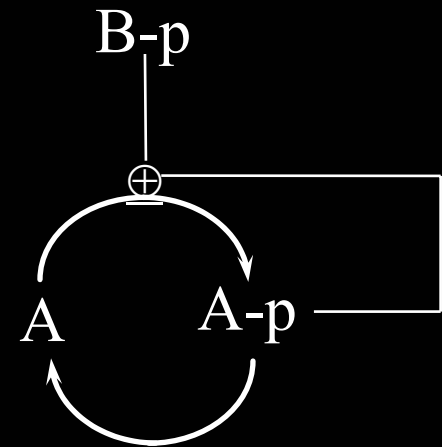


Problems of abstraction

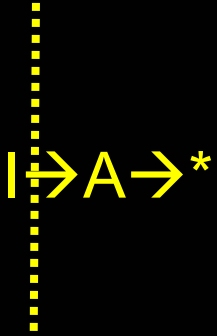


Problems of assumptions

# Problems of composition



# Brief Digression: Chemical Impedance



$$\frac{dA}{dt} = k_1[I] - k_2[A] \quad \Rightarrow \quad A_{t \rightarrow \infty}(I) = \frac{k_1}{k_2}[I]$$

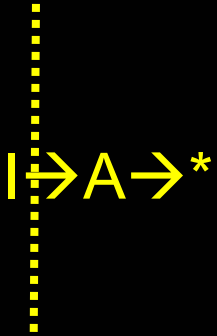
So A is the signal inside the cell that I is outside the cell.  
What if A signals to downstream targets by reacting with them?



$$\frac{dA}{dt} = k_1[I] - k_2[A] - k_3[A][B] \quad \Rightarrow \quad A_{t \rightarrow \infty}(I) = \frac{k_1[I]}{k_2 + k_3[B]}$$

The rates and concentrations of downstream processes degrade the signal from A.

# Brief Digression: Chemical Impedance



$$\frac{dA}{dt} = k_1[I] - k_2[A] \quad \Rightarrow \quad A_{t \rightarrow \infty}(I) = \frac{k_1}{k_2}[I]$$

But what if reaction is by reversible binding?



$$\frac{dA}{dt} = k_1[I] - k_2[A] - k_3[A][B] + k_4[C]$$

$$A_{t \rightarrow \infty}(I) = \frac{k_1[I]}{k_2}$$

The rates and concentrations of downstream processes don't affect the signal.