

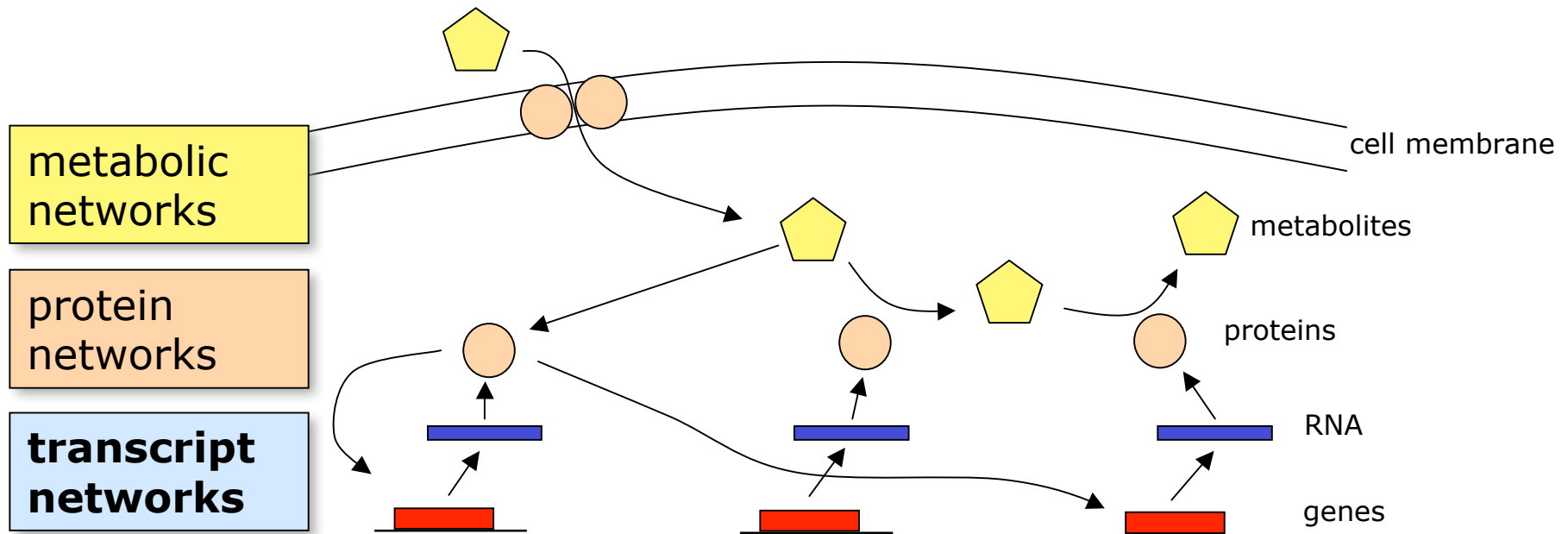
Gene network inference in diseases and drug discovery

Diego di Bernardo
Telethon Institute of Genetics and Medicine

19 Dicembre 2007, Napoli



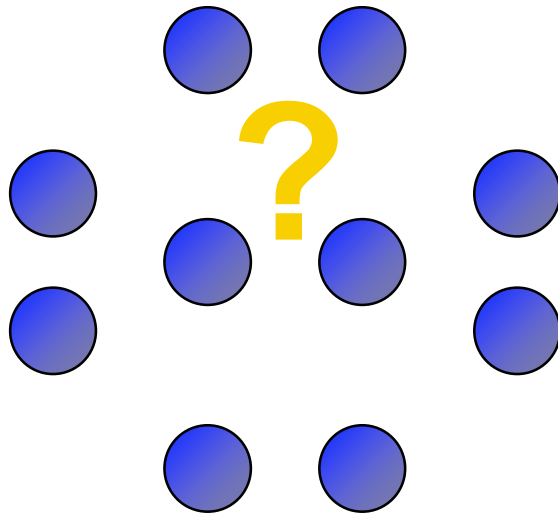
Gene Networks



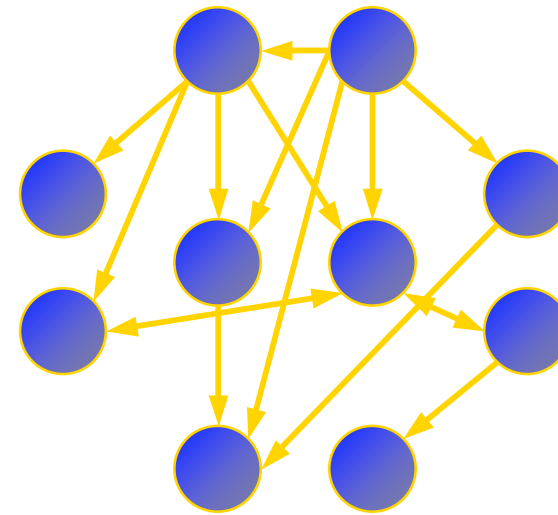
Our focus: methods to decode transcription regulation networks

Reverse engineering (or system id) gene networks:

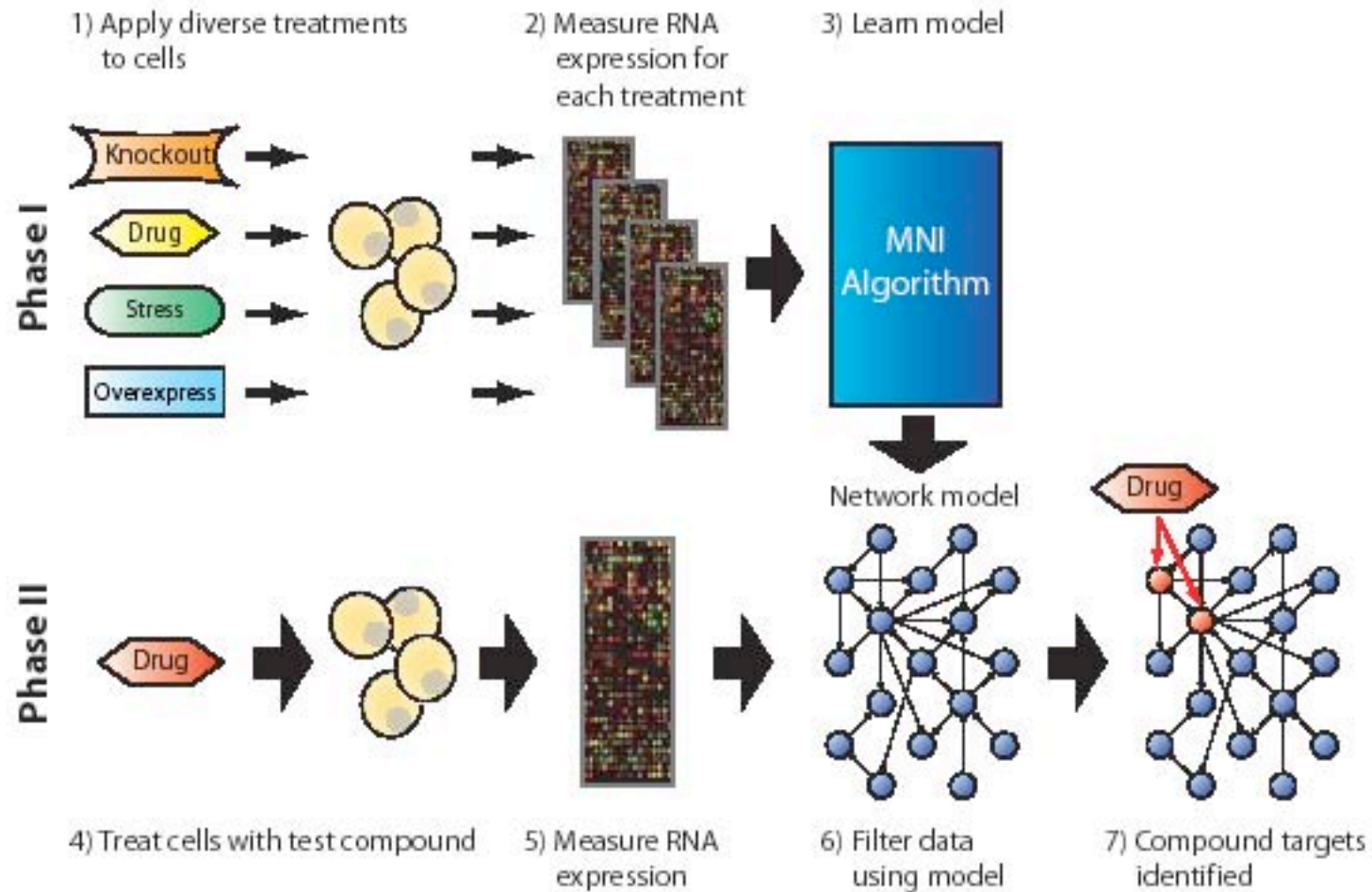
Unknown network



Inferred network

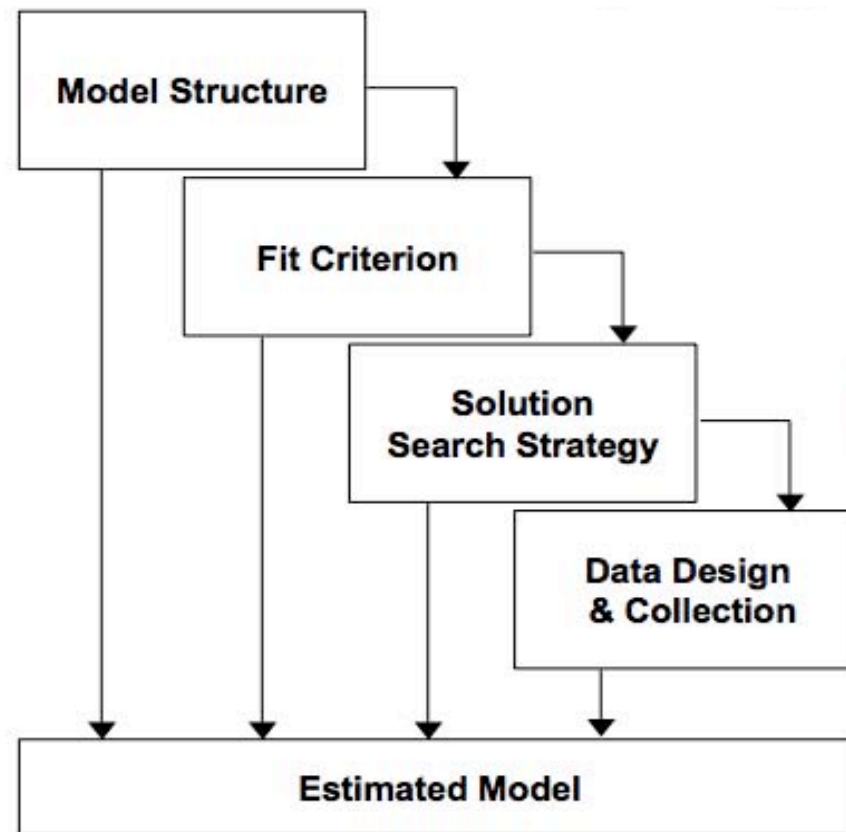


Aim of reverse engineering: gene function and drug MOA

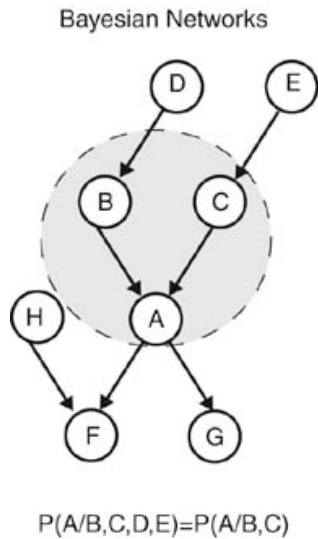


Reverse-engineering strategy:

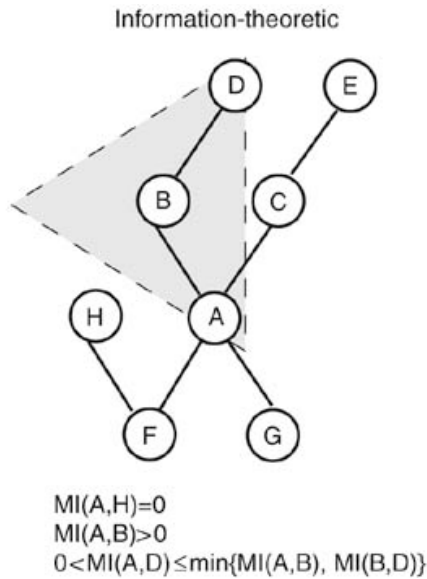
- Choose a **model**
- Choose a **fit criterion (cost function)** to measure the fit of the model to the data
- Define a **strategy** to find the parameters that best fit the data (i.e. that minimise cost function)
- Perform appropriate experiments to collect the experimental data:



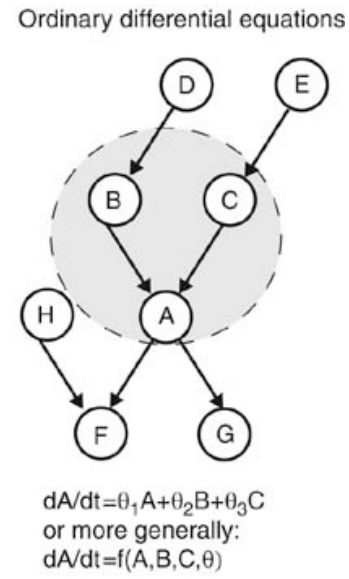
Comparison of different models



Bayesian network (Hartemink, A. (2005) Nature Biotechnology, 2005.)

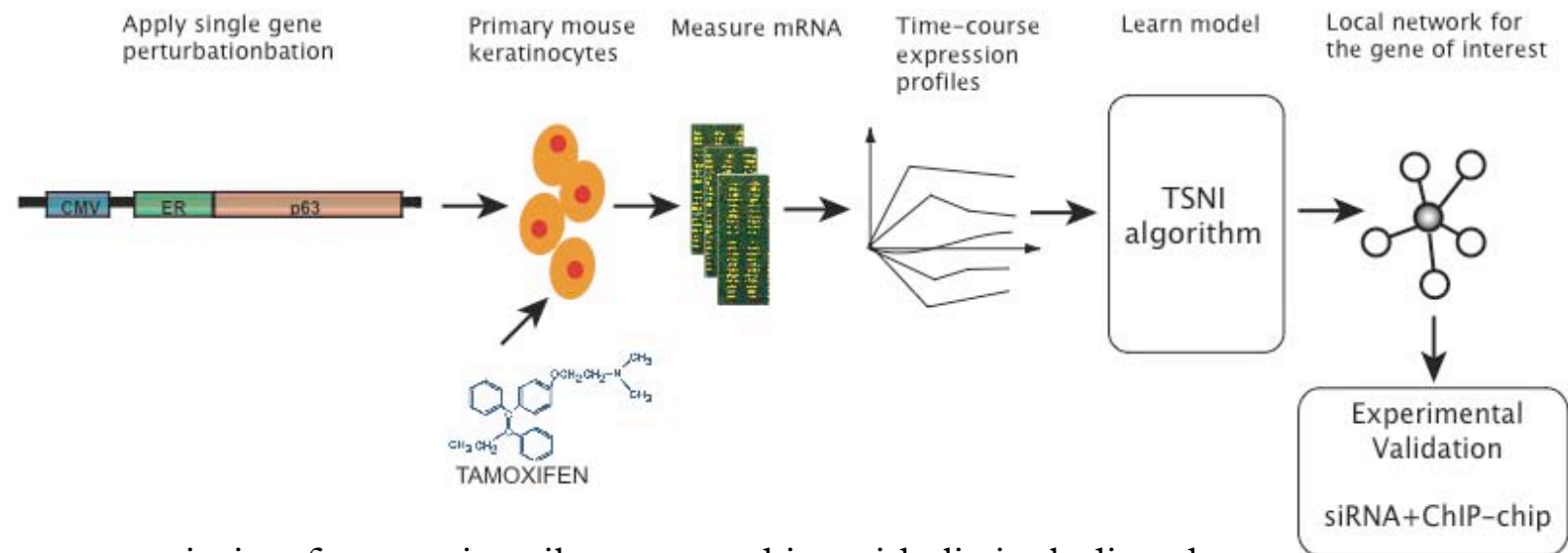


Information-theoretic: (Basso et al., Nature Genetics, 2006)



ODEs (Gardner, di Bernardo et al, Science, 2003; di Bernardo et al, Nature Biotechnology, 2005)

Understanding gene function in a genetic disease:

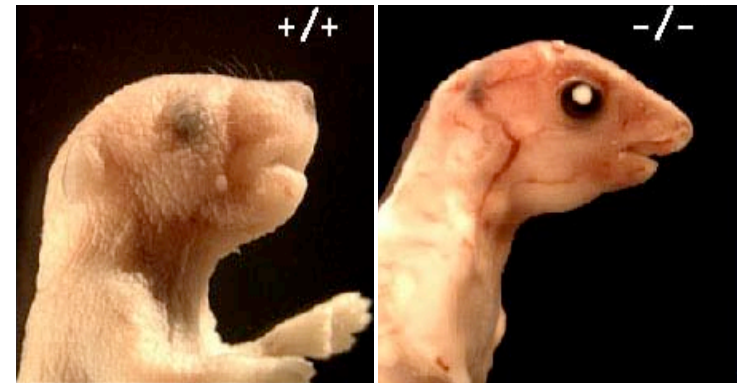
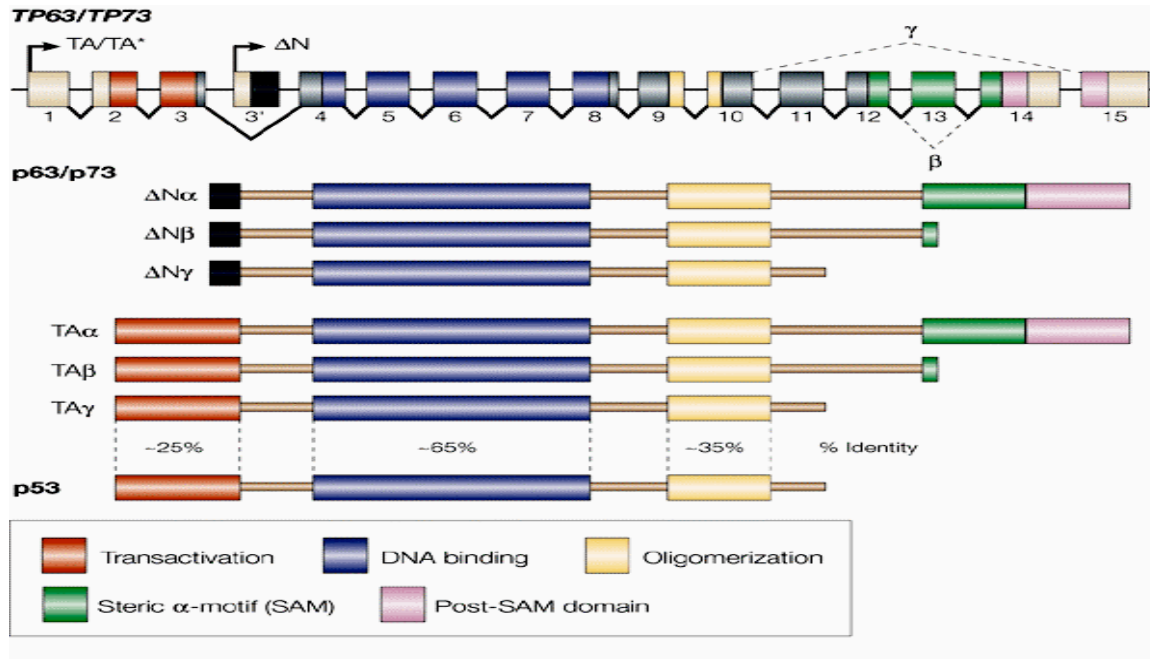


P63 is a transcription factor primarily expressed in epithelia including the proliferative compartments of the skin.

It plays an essential role in modulating cellular differentiation by unknown mechanisms.

Mutation in its DNA binding and SAM domains have been linked to five human malformation syndromes.

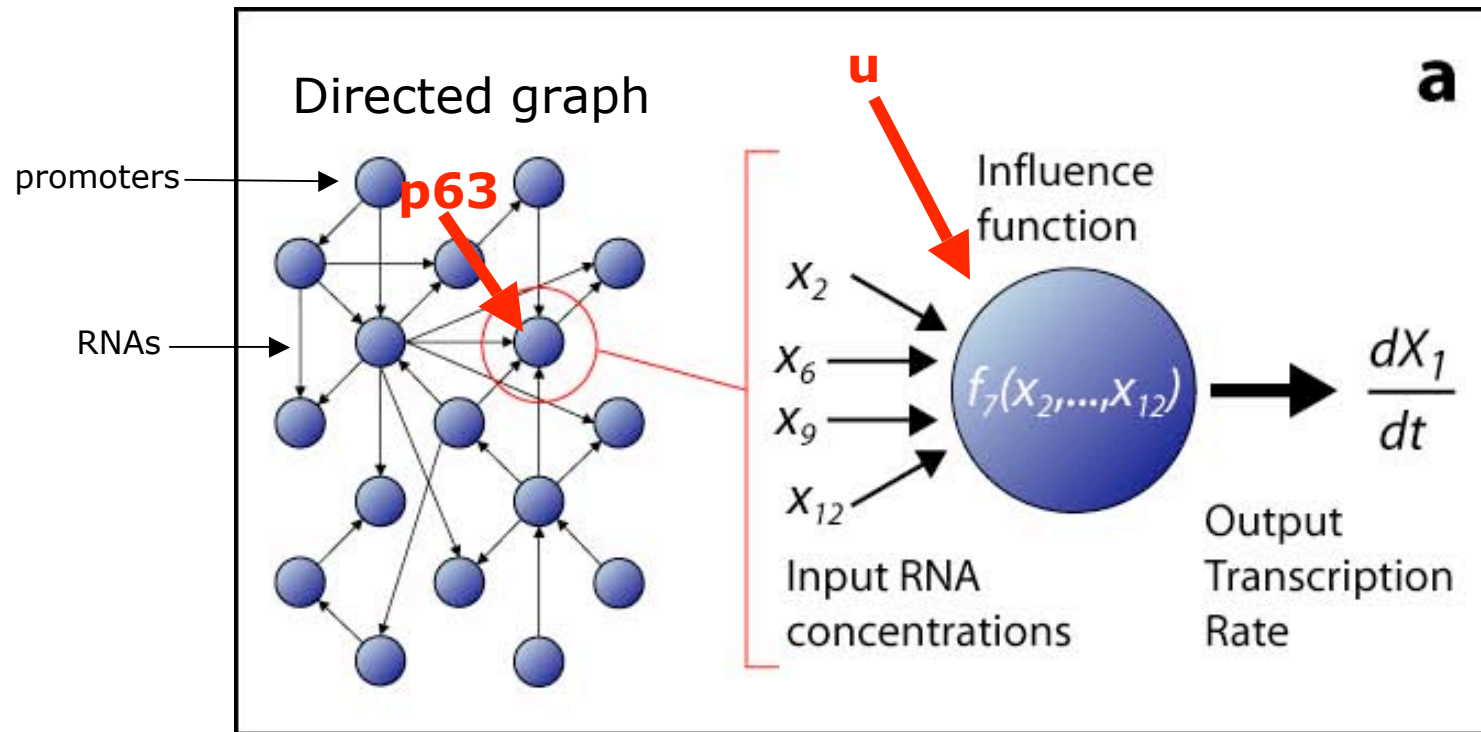
The transcription factor p63:



- P63 is a Transcription Factor (unknown targets)
- p63 gene has three different promoter that gives rise to two different transcripts (TA and ΔN);
- each transcript has three alternative splicing isoforms (α,β,γ)

- Lack of epidermis and hair follicles
- Defects in all squamous epithelia
- Limb truncations
- Craniofacial defects: lack of tooth primordia and eyelids
- Maxilla and mandible are truncated and secondary palate fails to close

Model structure: Ordinary differential equations



$$\frac{dX_1}{dt} = \mathbf{a}_2 X_2 + \mathbf{a}_6 X_6 + \mathbf{a}_9 X_9 + \mathbf{a}_{12} X_{12} + \mathbf{b}_1 u$$

Model structure:

$$\left\{ \begin{array}{l} x'_1(t) = a_{11}x_1(t) + a_{12}x_2(t) + \dots + a_{1n}x_n(t) + b_1u(t) \\ \dots \\ x'_n(t) = a_{n1}x_1(t) + a_{n2}x_2(t) + \dots + a_{nn}x_n(t) + b_nu(t) \end{array} \right.$$

Gene effect
unknown

Model structure - Discrete model:

$$x_1(t_{k+1}) = a^d_{11}x_1(t_k) + a^d_{12}x_2(t_k) + \dots + a^d_{1n}x_n(t_k) + b^d_1u(t_k)$$

.....

$$x_n(t_{k+1}) = a^d_{n1}x_1(t_k) + a^d_{n2}x_2(t_k) + \dots + a^d_{nn}x_n(t_k) + b^d_nu(t_k)$$

Or in matrix format:

$$\mathbf{x}(t_{k+1}) = \mathbf{A}^d \mathbf{x}(t_k) + \mathbf{b}^d u(t_k)$$

For $k=1..M$

$$\mathbf{X}(t_{k+1}) = \mathbf{A}^d \mathbf{X}(t_k) + \mathbf{b}^d u(t_k)$$

Where $(\mathbf{N} \times \mathbf{M}) = (\mathbf{N} \times \mathbf{N})(\mathbf{N} \times \mathbf{M}) + (\mathbf{N} \times 1)(1 \times \mathbf{M})$

Fit criterion (Least Squared Error):

$$\dot{X}(t_k) = AX(t_k) + BU(t_k) \quad k = 1 \dots M$$

$$X(t_{k+1}) = A_d * X(t_k) + B_d * U(t_k)$$

$$X(t_{k+1}) = \begin{bmatrix} A_d & B_d \end{bmatrix} * \begin{bmatrix} X(t_k) \\ U(t_k) \end{bmatrix}$$

$$N \times M = N \times (N+1) \quad (N+1) \times M$$

$$X(t_{k+1}) = H * Y(t_k)$$

$$B = (A_d + 1)^{-1} B_d$$

$$H = \begin{bmatrix} A_d & B_d \end{bmatrix}$$

$$Y(t_k) = \begin{bmatrix} X(t_k) \\ U(t_k) \end{bmatrix}$$

Solution search strategy:

- Usually $N \gg M$ (more unknown than data points; i.e. 1000 genes e 20 points)
- From the previous equation we know that a unique solution will exist only if $N \leq M-1$
- We need a trick to reduce number of unknowns by dimensional reduction:
 - Regression with Variable selection
 - Clustering
 - Regression with SVD

Solution search strategy: Regression using PCS/SVD

$$X(t_{k+1}) = H * Y(t_k)$$

$$H = [A_d \quad B_d]$$

$$Y(t_k) = \begin{bmatrix} X(t_k) \\ U(t_k) \end{bmatrix}$$

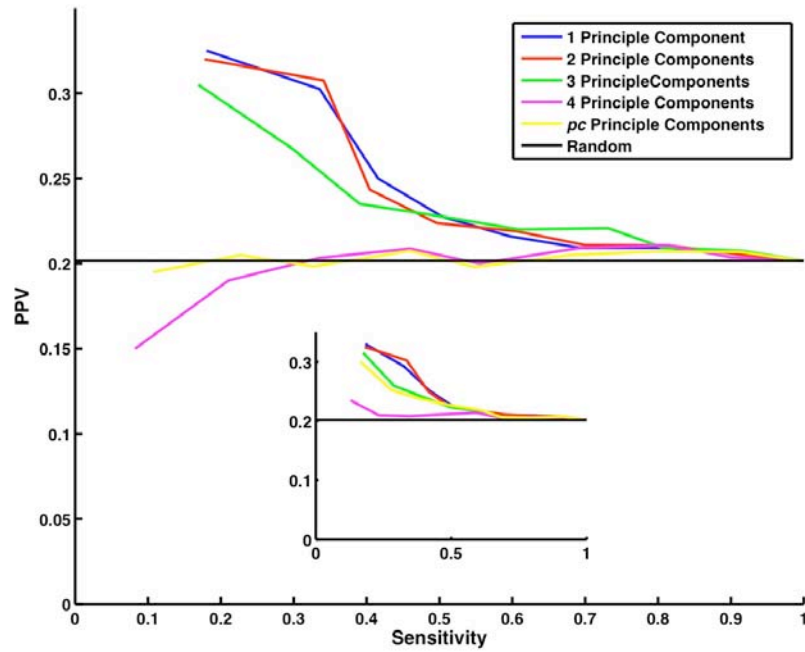
- Singular Value Decomposition:
 - $Y = UDV'$ $(N+1 \times M) = (N+1 \times M)(M \times M)(M \times M)$



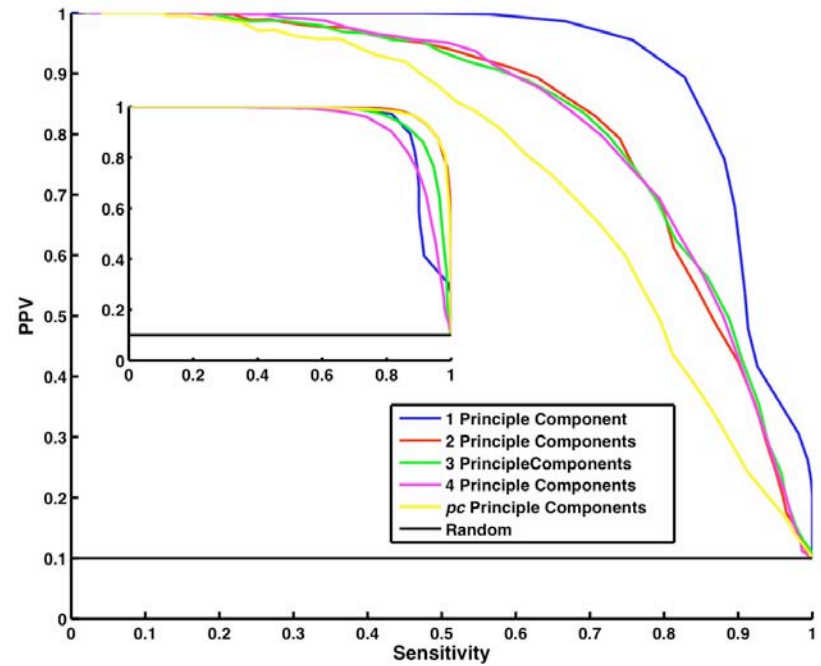
- $Y = U_r D_r V_r'$ $(N+1 \times M) = (N+1 \times K)(K \times K)(K \times M)$
- $X(t_{k+1}) = HY = (N \times N+1)(N+1 \times M) =$
- $X(t_{k+1}) = H U_r D_r V_r' = (H U_r) D_r V_r' = \mathbf{(N \times K)}(K \times K)(K \times M)$

Preliminary Results on noisy “*in silico*” data

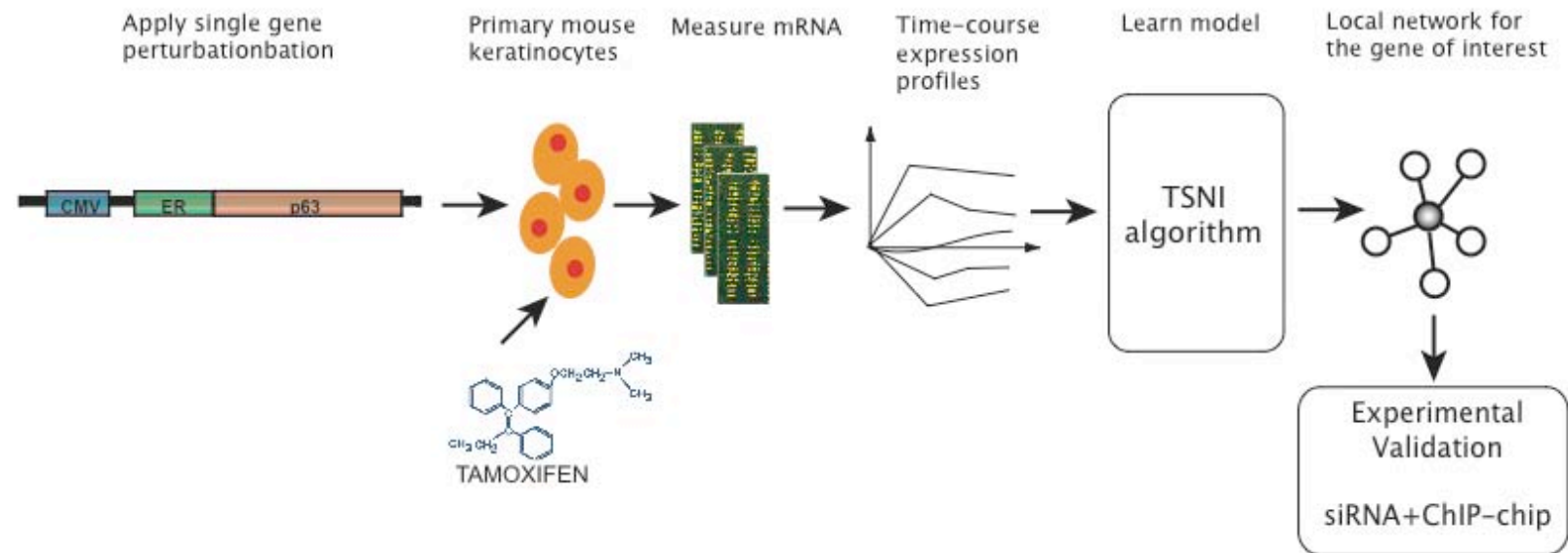
Network of 10 genes (A)



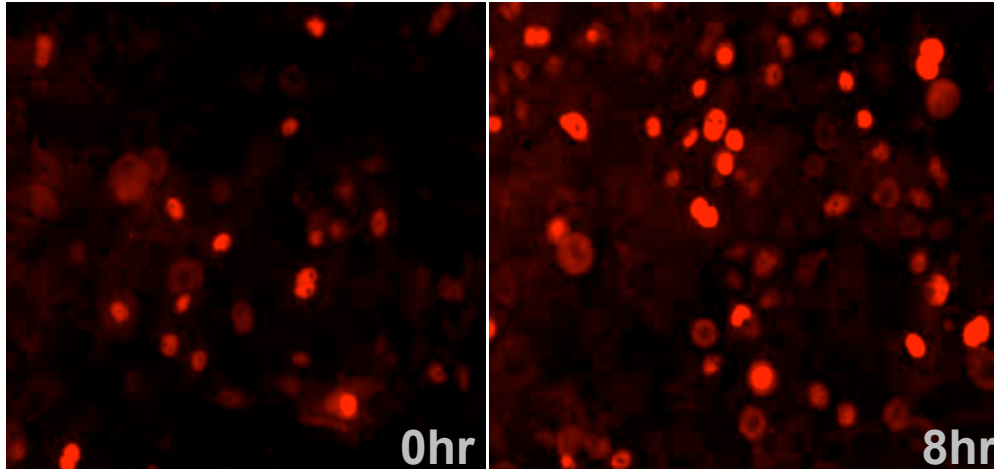
Target prediction in 1000 gene network (B)



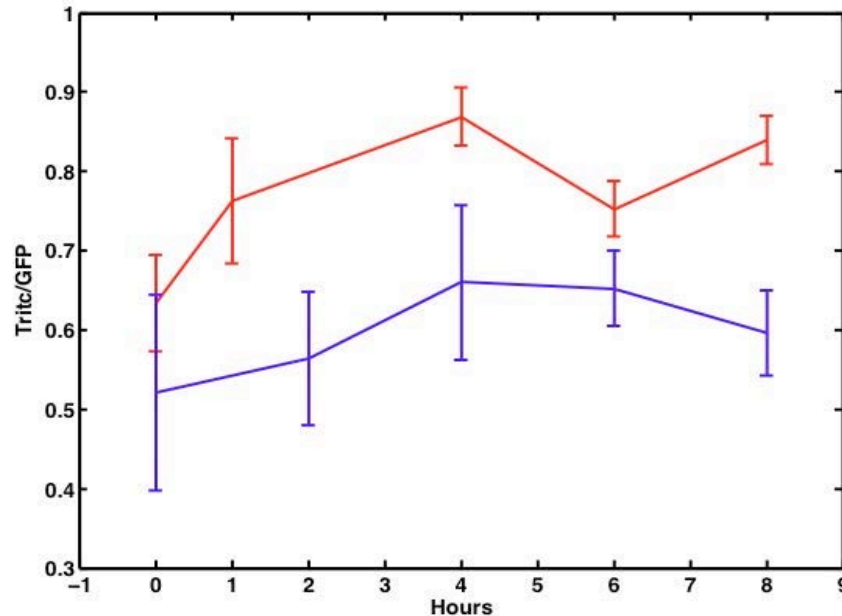
Application to p63:



Localization of $\Delta Np63$ upon tamoxifen treatment by immunofluorescence



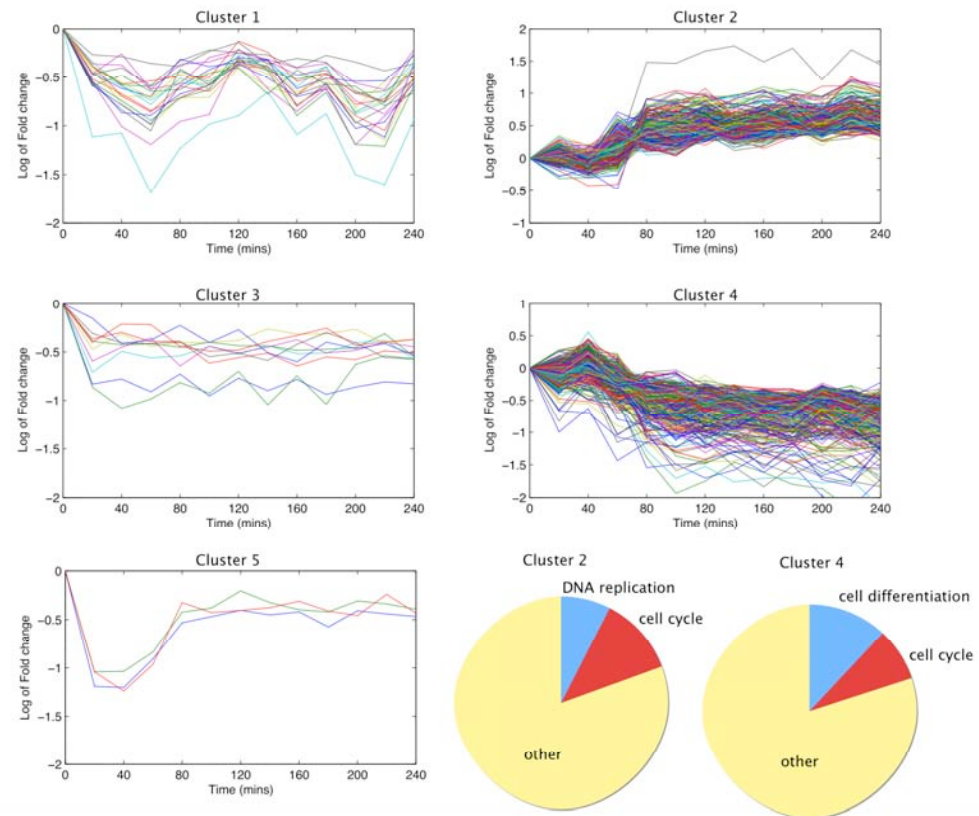
Murine keratinocytes are infected with Gingo ErDNp63 retrovirus. The cells are collected at the following time points from Tamoxifen induction: 1hr, 2hrs, 4hrs, 8hrs.



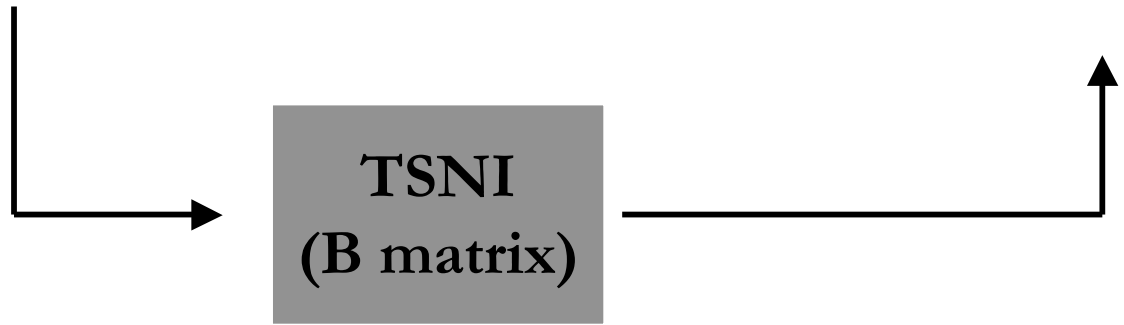
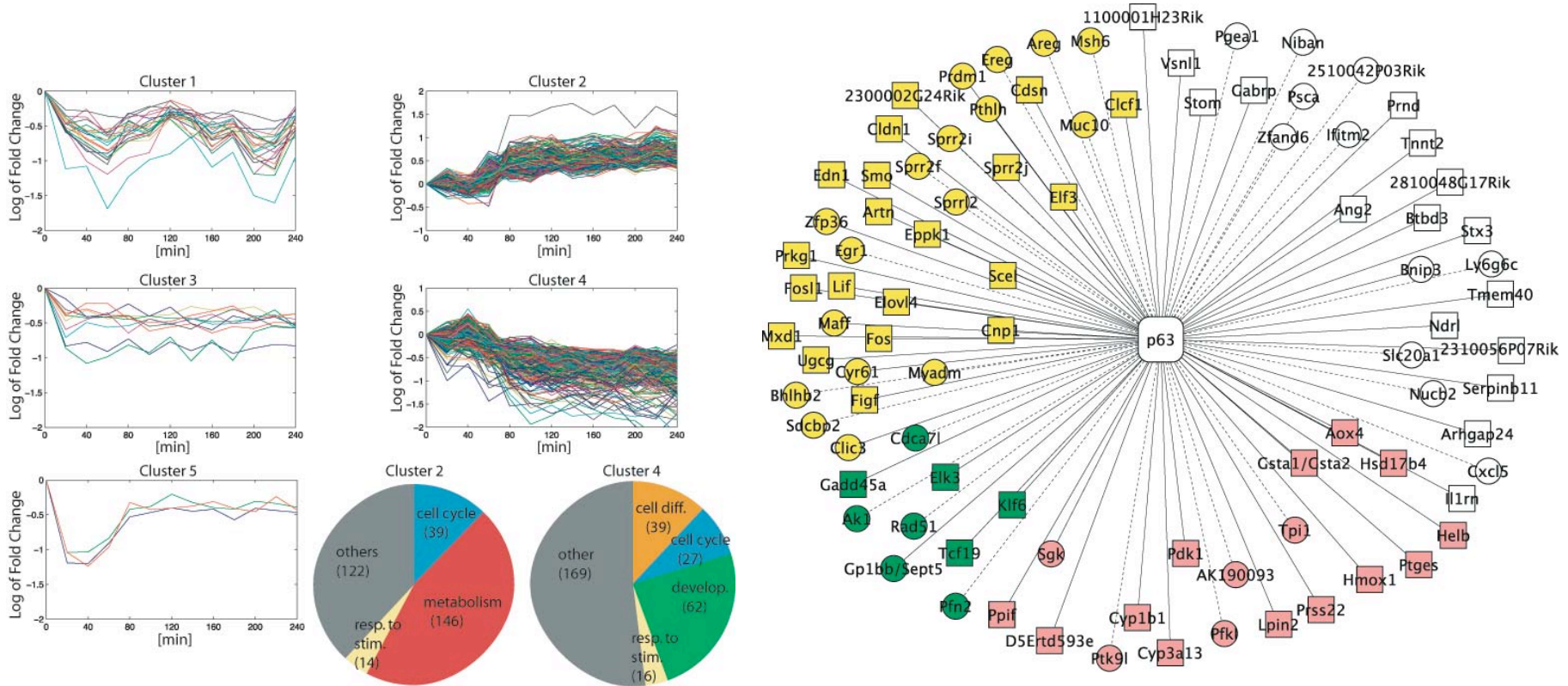
Immunofluorescence plot: red and blue lines represent two different immunofluorescence experiments.

$$\dot{x}(t) = A \cdot x(t) + Bu(t)$$

Hierarchical Clustering of 786 expression profiles:



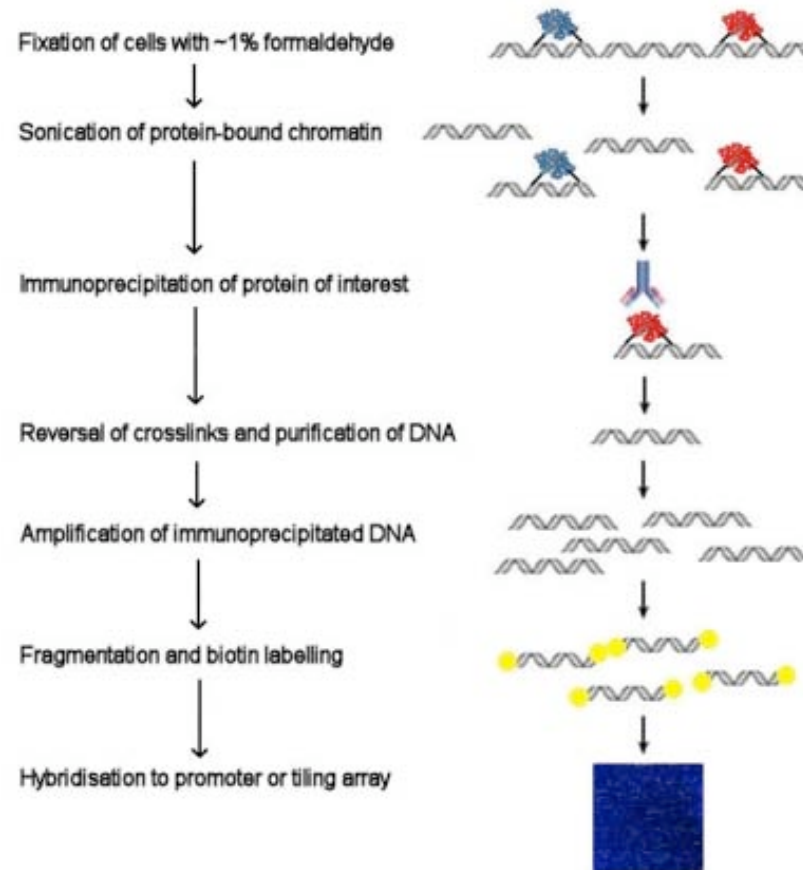
Application of TSNI on 786 expression profiles:



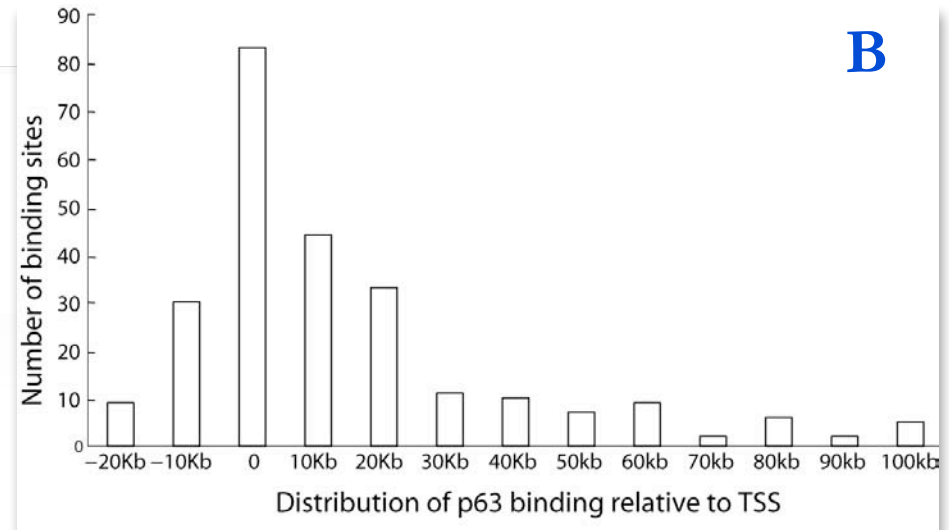
Validation of predicted targets: (1) ChIP-chip

- We selected the top 100 predicted targets and as control the bottom 200 predicted genes.
- For each gene we analysed 20 Kbp upstream+gene sequence via ChIP-chip (Agilent 200k custom-array)
- We found about 300 p63-binding sites bu ChIP-chip with a p-value<0.01

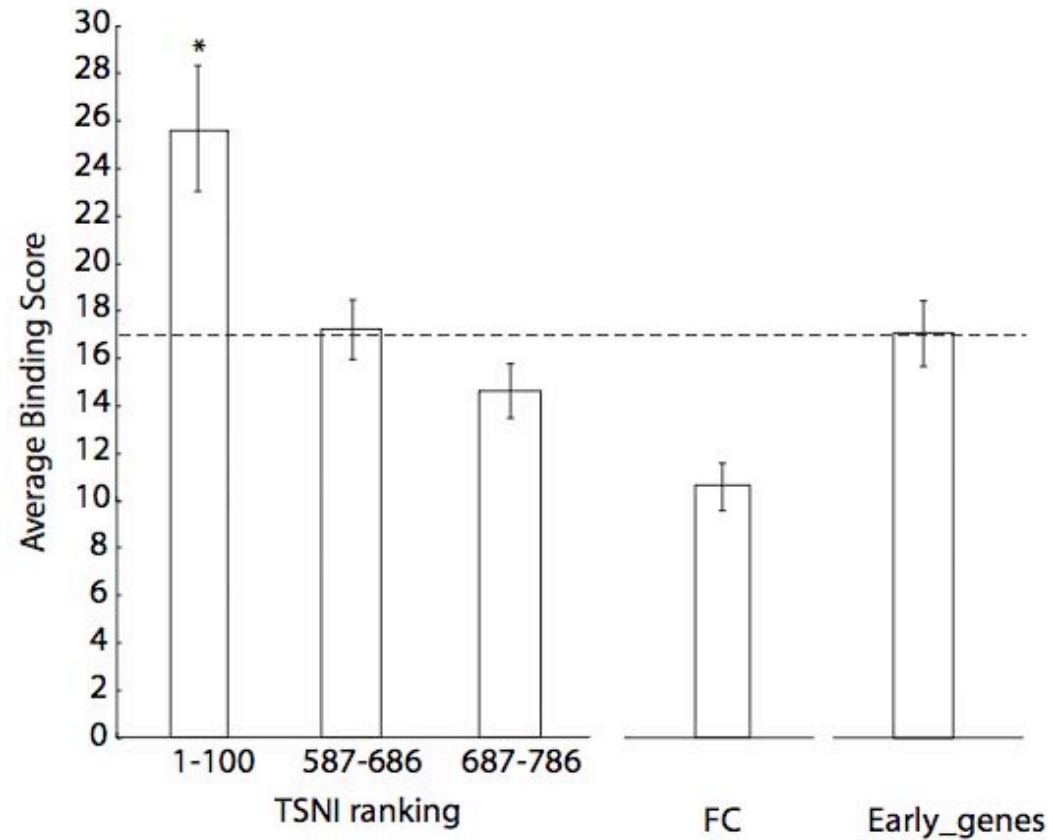
Validation of predicted targets: (1) ChIP-chip



Validation of predicted targets: (1) ChIP-chip

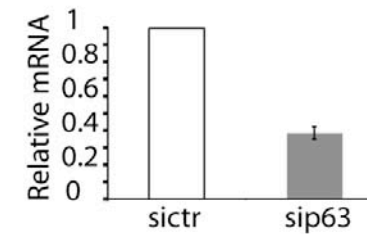
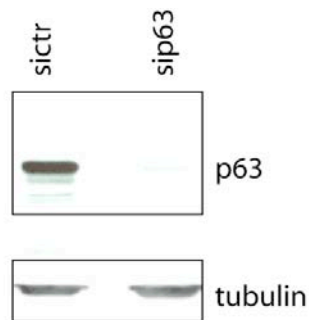
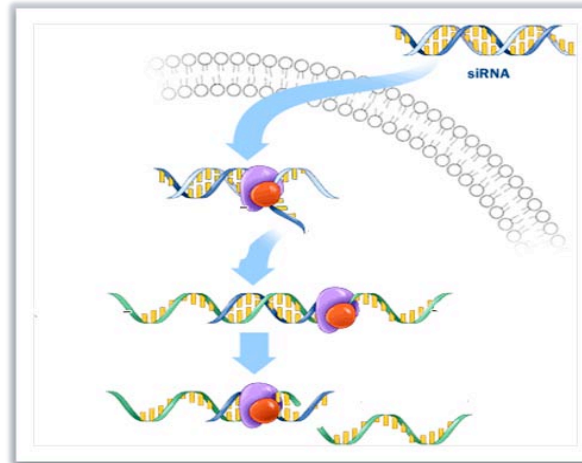


TSNI results are confirmed by ChIP-chip

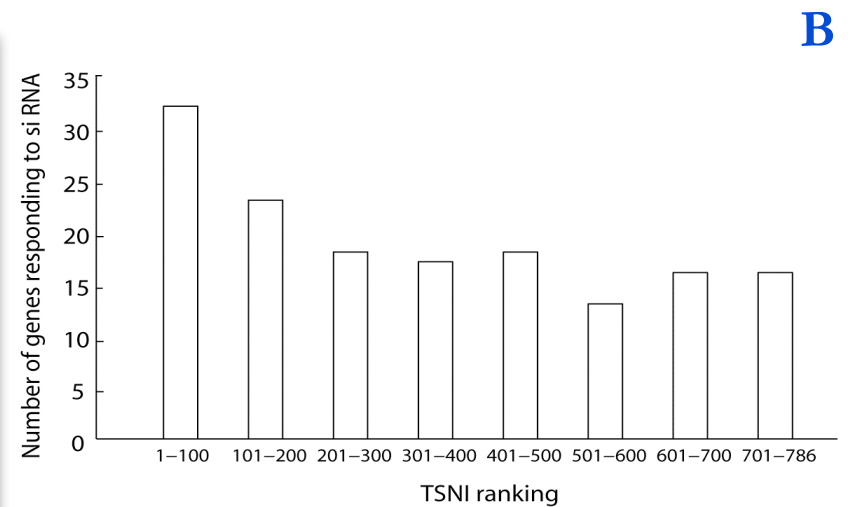
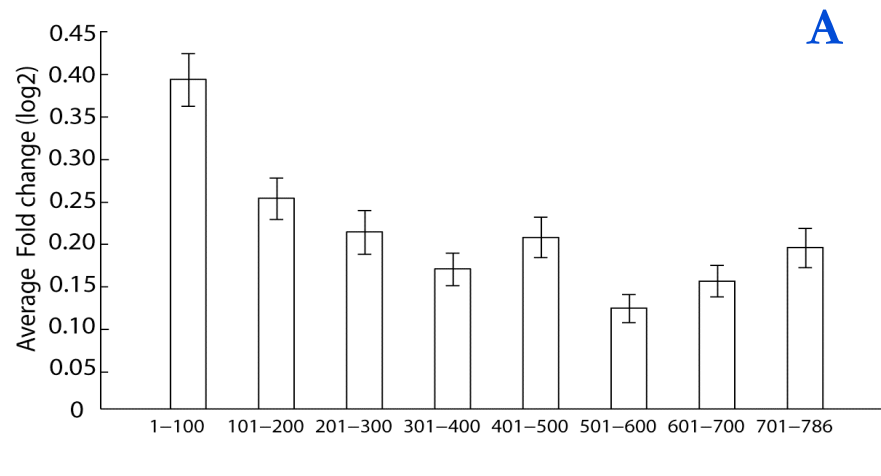


Knocking down p63 expression via siRNA:

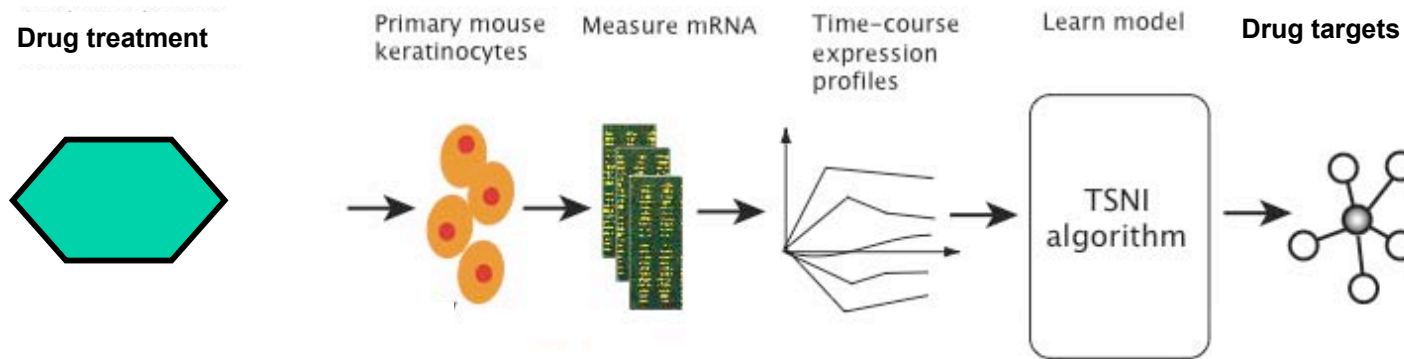
RNA interference against the p63 DNA binding domain



TSNI results are confirmed by siRNA

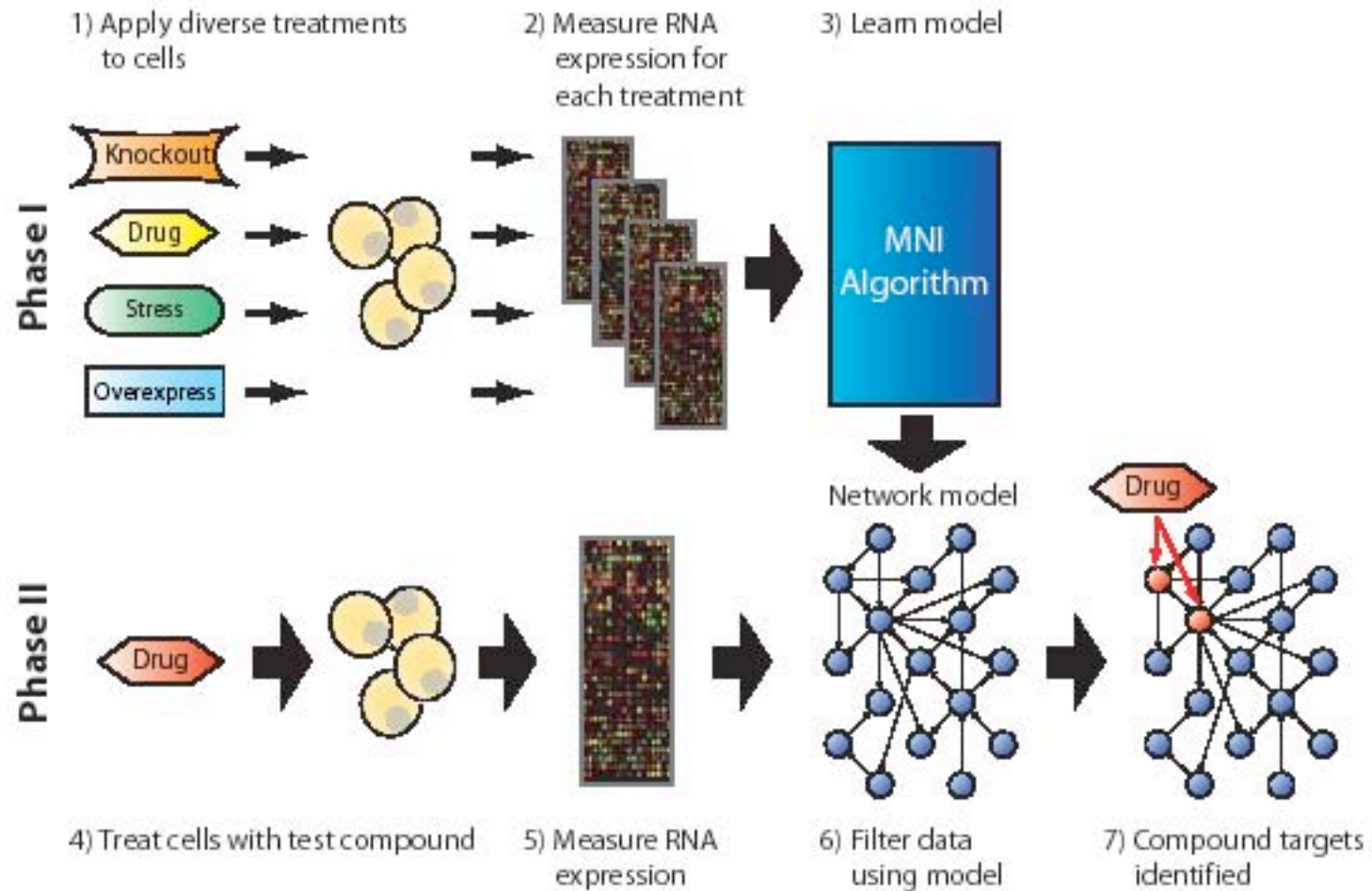


Application of TSNI to drug discovery :



A “steady-state” approach to drug target identification:

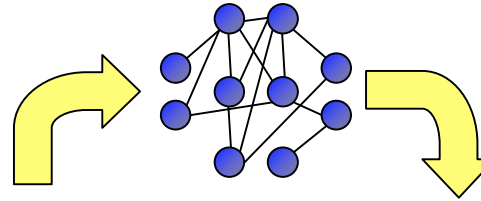
D di Bernardo et al
Nature Biotechnology, 2005



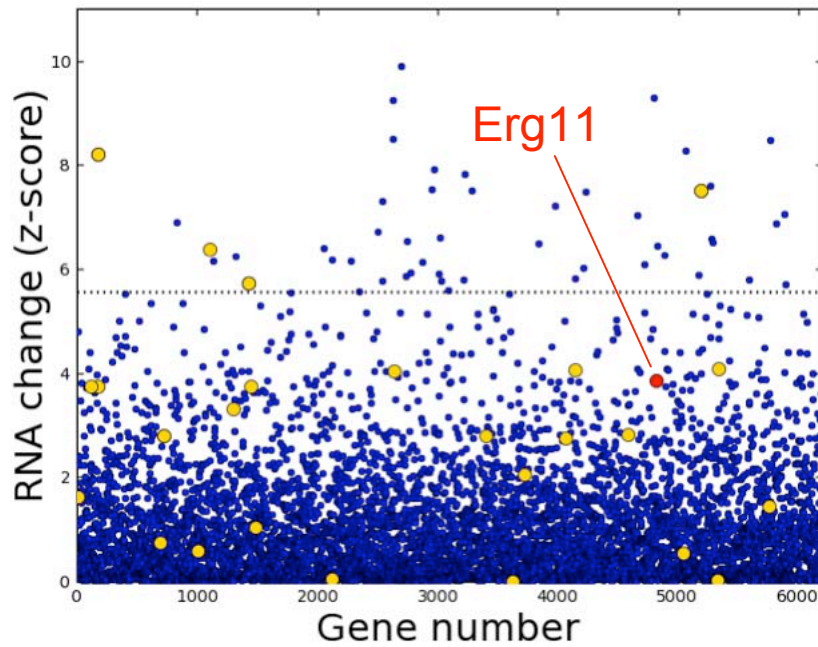
MNI identifies target of itraconazole in *S. cerevisiae*

Itraconazole treatment: a known target is ERG11

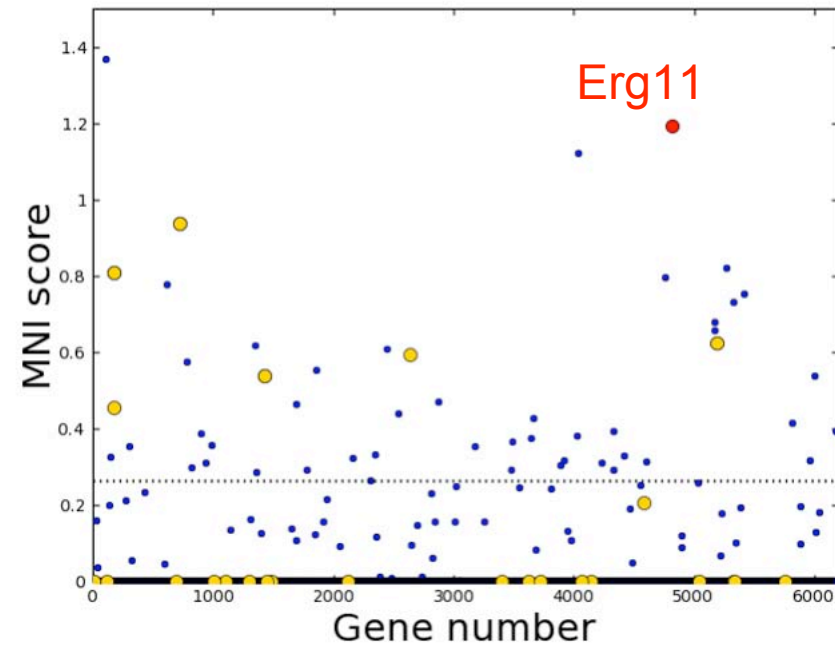
Filter through MNI-inferred network model



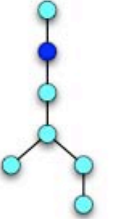

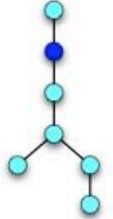





Expression Change



MNI Predictions



MNI results:

Compound	Known pathway	Known target	Predicted pathway	Ranked target genes rank
Itraconazole	ergosterol biosynthesis 	Erg11 	steroid metabolism 	ERG11 , ERG24, ERG1, ERG25, CYB5, ERG27, ATF2 
Hydroxyurea	DNA replication 	Rnr2, Rnr4 	DNA replication 	RNR4 , RNR2 , RNR1, RNR3 

Pathway data from GO database

MNI identifies MOA for 8 out of 9 drugs:

Table 2 Pathways and associated genes targeted by drug compounds

Drug	Known pathway	Known target	Significant GO ontology (rank, <i>P</i> -value)	Ranked pathway genes (rank)
Terbinafine	Ergosterol biosynthesis ⁴¹	Erg1	Steroid metabolism (1, 10 ⁻¹⁴)	<i>ERG7</i> (4), <i>ERG1</i> (5), <i>ERG8</i> (11), <i>ERG26</i> (13), <i>UPC2</i> (17), <i>ERG28</i> (18), <i>ERG11</i> (20), <i>DAP1</i> (33), <i>HES1</i> (34), <i>ATF2</i> (36), <i>ERG5</i> (49)
Lovastatin	Ergosterol biosynthesis ⁴²	Hmg2, Hmg1	Lipid metabolism (1, 10 ⁻⁴)	<i>BST1</i> (1), <i>ERG1</i> (18), <i>YSR3</i> (23), <i>HMG2</i> (30), <i>LCB5</i> (31), <i>ERG13</i> (36), <i>VRG4</i> (48)
Itraconazole	Ergosterol biosynthesis ⁴³	Erg11	Steroid metabolism (1, 10 ⁻⁸)	<i>ERG11</i> (2), <i>ERG24</i> (4), <i>ERG1</i> (6), <i>ERG25</i> (13), <i>CYB5</i> (16), <i>ERG27</i> (19), <i>ATF2</i> (23)
Hydroxyurea	DNA replication ⁴⁴	Rnr2, Rnr4	Heteroduplex formation (1, 10 ⁻⁴)	<i>RAD51</i> (15), <i>RAD54</i> (47)
			DNA replication (2, 10 ⁻²)	<i>RNR4</i> (2), <i>RNR2</i> (6), <i>RNR1</i> (14), <i>RNR3</i> (23)
Cycloheximide	Protein biosynthesis ⁴⁵	Ribosome	Nuclear mRNA splicing, via spliceosome (1, 10 ⁻⁴)	<i>SYF1</i> (3), <i>SMD3</i> (19), <i>HSH49</i> (42)
			–	<i>RPL26B</i> (32), <i>RPS29A</i> (34)
Tunicamycin	N-linked glycosylation ⁴⁶	Alg7	Protein-ER targeting (1, 10 ⁻³)	<i>SEC62</i> (1), <i>SIL1</i> (31), <i>SEC59^a</i> (43)
Nikkomycin	Cell wall chitin biosynthesis ⁴⁷	Chs3	Protein amino acid alkylation (1, 10 ⁻³)	<i>SWD2</i> (3), <i>RMT2</i> (6)
Drugs not in the original compendium data set				
3-aminotriazole	Histidine biosynthesis ⁴⁸	His3	Organic acid metabolism (1, 10 ⁻⁷)	<i>FRM2</i> (8), <i>BIO5</i> (9), <i>YAT2</i> (10), <i>ARO10</i> (18), <i>ARO9</i> (20), <i>CHAI</i> (21), <i>BIO3</i> (31), <i>ARG1</i> (33), <i>ARG4</i> (37), <i>HIS5[†]</i> (42), <i>LYS1</i> (47), <i>SAM2</i> (50)
	Oxygen and reactive oxygen species metabolism ³⁰	Cta1		
Dyclonine	Ergosterol biosynthesis ¹	Erg2	Sterol biosynthesis (1, 10 ⁻¹⁸)	<i>ERG3</i> (1), <i>ERG6</i> (2), <i>CYB5</i> (3), <i>ERG2</i> (4), <i>ERG11</i> (6), <i>ERG28</i> (10), <i>ERG1</i> (12), <i>ERG5</i> (13), <i>ERG27</i> (18), <i>MVD1</i> (23), <i>ERG24</i> (30), <i>ERG26</i> (37)

Acknowledgements (<http://dibernardo.tigem.it> for more info):

Dr Alberto Ambesi
(Sequence analysis)

Giusy Della Gatta
(Experimental work)

Mukesh Bansal
(Computational work)



Velia Siciliano
Lucia Marucci
Giulia Cuccato



Vincenzo Belcastro

Francesco Iorio

Mario Lauria

...and Dr Caterina Missero (CEINGE, Napoli) for supervision of experiments; Prof. Tim Gardner and Prof. Jim Collins (Boston University, USA) for MNI/NIR related work.



Selection of Significant Genes in Time Series experiment:

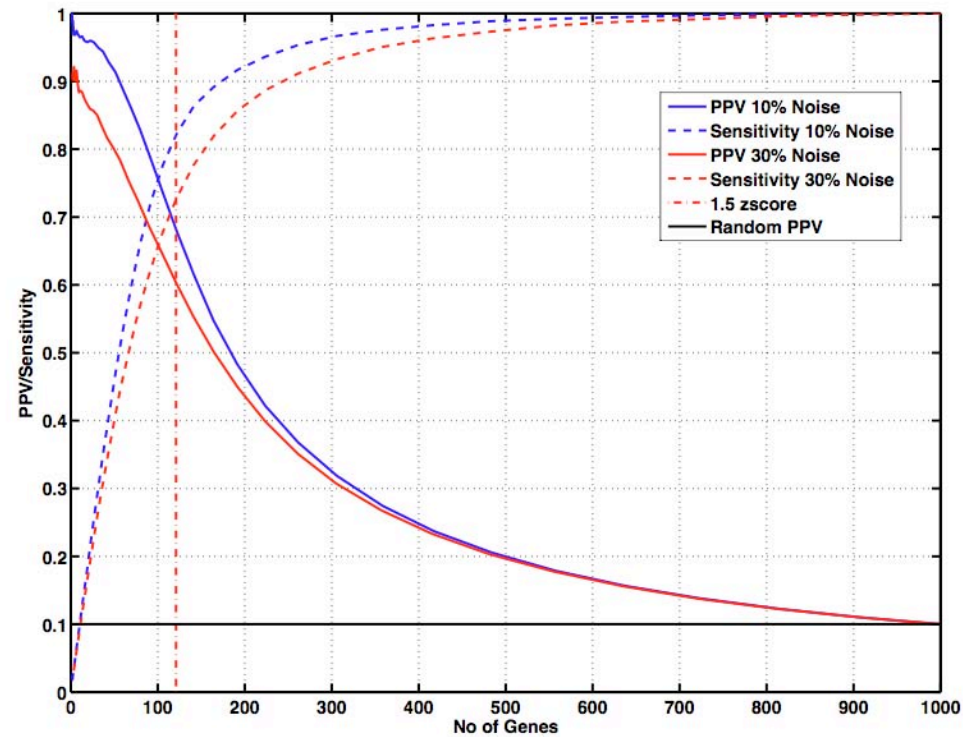
- First filter
 - Area under curve
- Second filter
 - Standard deviation is computed at each time point using smoothing and generalized cross validation algorithm
- Statistical test for significance (Chi-Square)

in silico with dual perturbation

$$dX_1/dt = a_2 X_2 + a_6 X_6 + a_9 X_9 + a_{12} X_{12} + b_{11} u_1 + b_{12} u_2$$

p63

Tamoxifen



Improvement of TSNI to identify the network (A)

- Integral Approach instead of differential approach

$$\int_0^{t_k} \dot{X}(t)dt = A \int_0^{t_k} X(t)dt + B \int_0^{t_k} U(t)dt \quad k = 1 \dots M$$

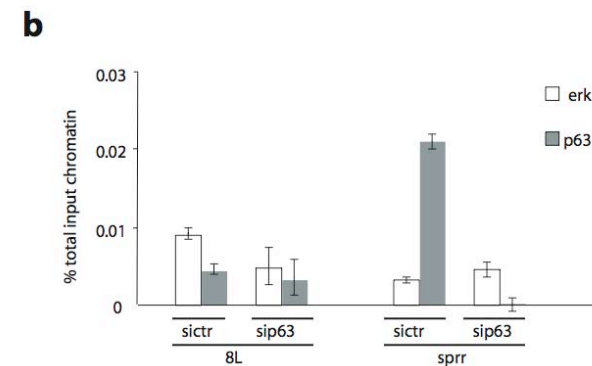
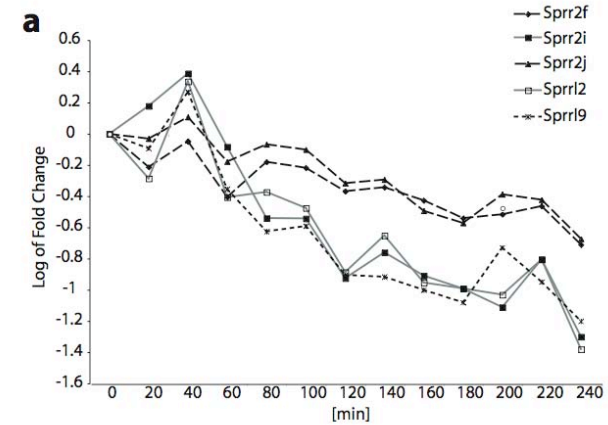
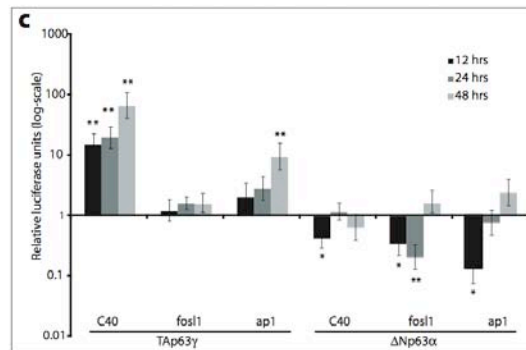
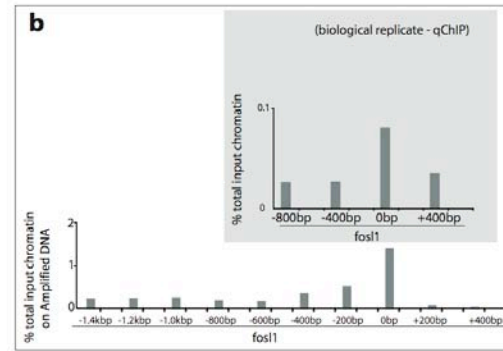
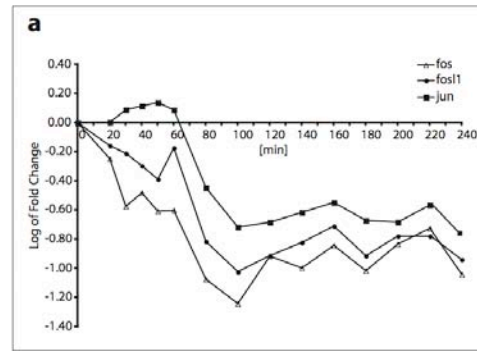
$$X(t_k) = A \int_0^{t_k} X(t)dt + B \int_0^{t_k} U(t)dt \quad k = 1 \dots M$$

- Bootstrapping

Mukesh Bansal, Diego di Bernardo

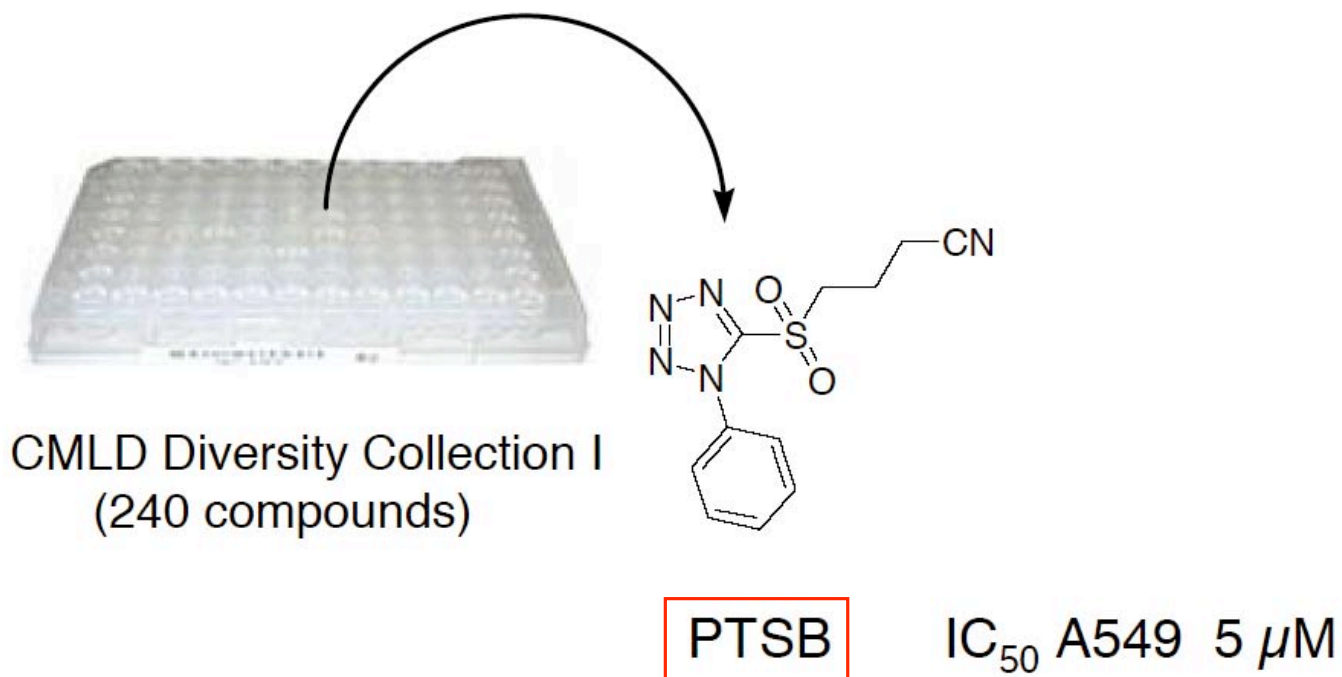
Inference of gene networks from temporal gene expression profiles. IET (Under Review)

(1a) Time Series Network Identification Algorithm (NEW) :



P63 directly regulates at early times AP1 complex and markers of terminal keratinocyte differentiation -*Manuscript submitted*-

Identified novel anticancer compound via chemical screen



- PTSB inhibits growth in yeast and tumor cell lines

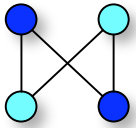
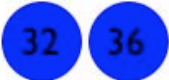
In collaboration with Schaus and Elliot laboratories

Dept. of Chemistry, Boston University

Center for Methodology and Library Development (CMLD), Boston U.

MNI identifies two enzyme targets of PTSB

Identifies thioredoxin (TRX2) and thioredoxin reductase (TRR1)

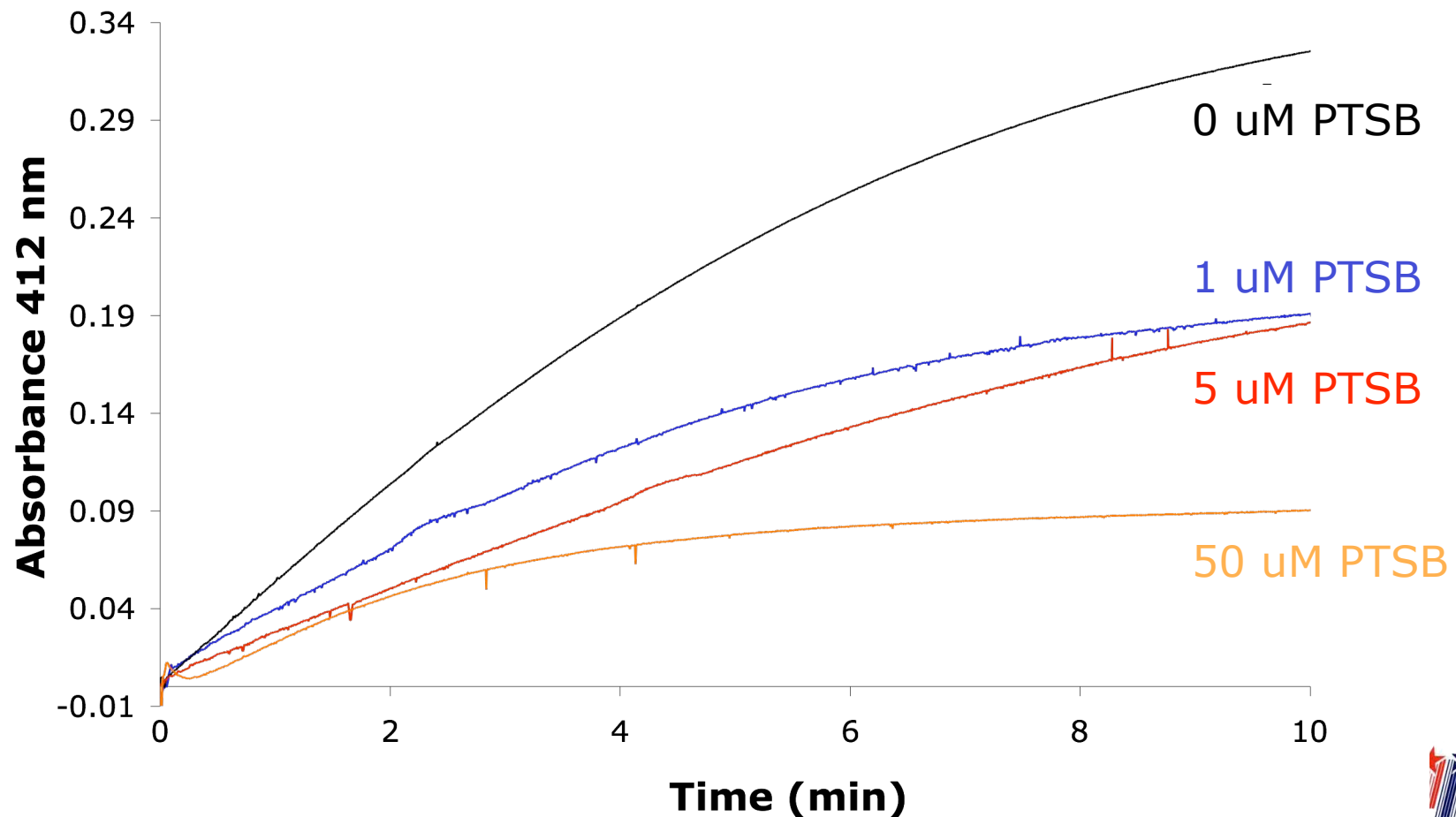
Compound	Known pathway	Known target	Predicted pathway	Ranked target genes rank
PTSB	unknown	unknown	cell redox homeostasis 	TRR1, TRX2 

TRR1/TRX2 activity inhibited in presence of PTSB

Assay:

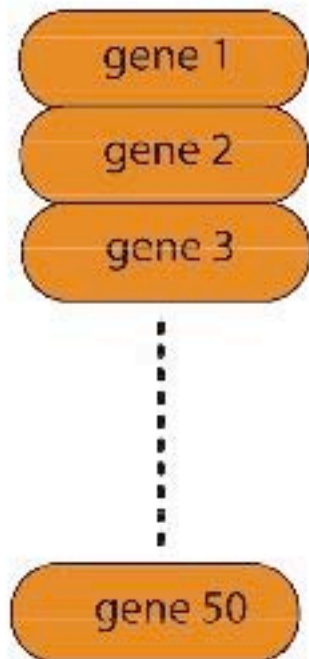
Thioredoxin reduction of dithio(bis)nitrobenzoic acid (DTNB)

- Product of reaction = thiolate anion, measured via A412

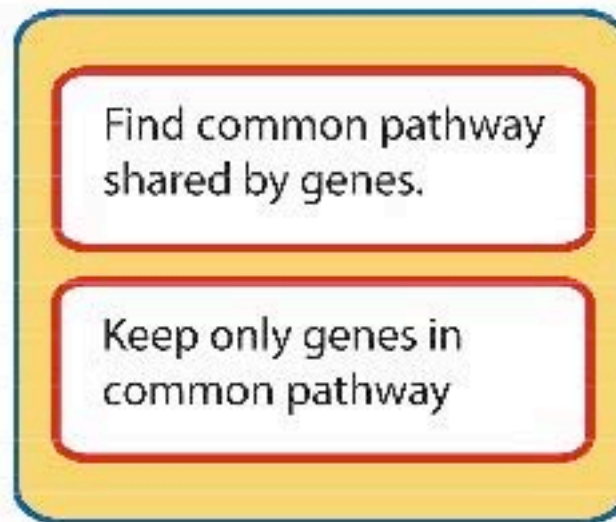


Geno Ontology Filtering

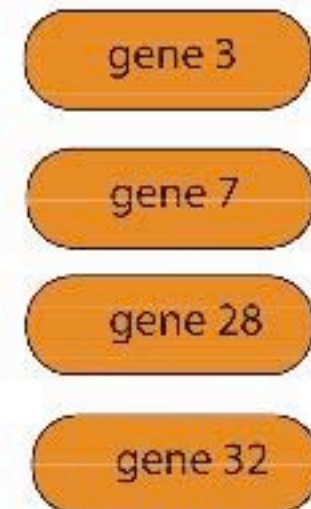
Top 50 genes ranked by MNI



Filter gene list using gene ontology



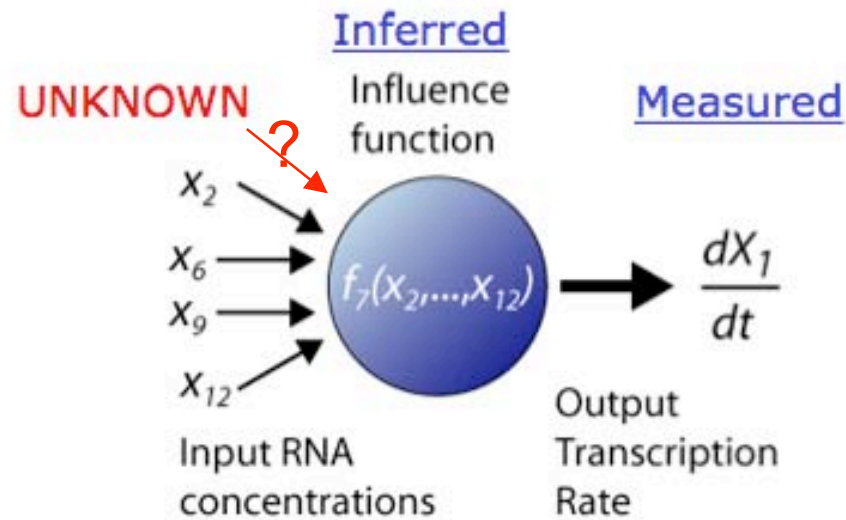
Genes selected after filtering



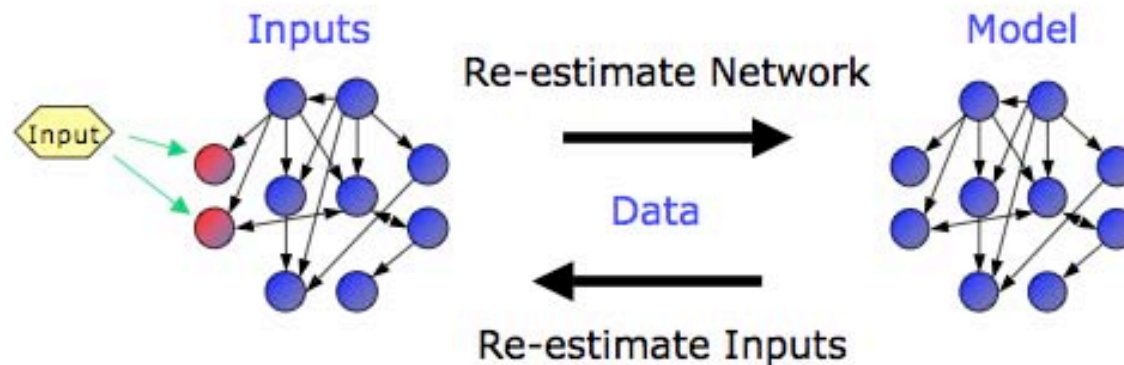
MNI: Microarray Network Identification

The MNI Algorithm: an unsupervised approach

MNI Method



MNI Algorithm: Recursively estimate inputs and model using data



Model structure:

$$\left\{ \begin{array}{l} \mathbf{x}'_1(t) = a_{11}\mathbf{x}_1 + a_{12}\mathbf{x}_2 + \dots + a_{1n}\mathbf{x}_n + \mathbf{b}_1\mathbf{u} \\ \dots \\ \mathbf{x}'_n(t) = a_{n1}\mathbf{x}_1 + a_{n2}\mathbf{x}_2 + \dots + a_{nn}\mathbf{x}_n + \mathbf{b}_n\mathbf{u} \end{array} \right. \quad \text{Drug effect unknown}$$

Or in matrix format:

$$\mathbf{x}' = \mathbf{A}\mathbf{x} + \mathbf{b}\mathbf{u}$$

Handwritten annotations:
- A red arrow points from the word "NETWORK" to the matrix \mathbf{A} .
- A red arrow points from the word "DRUG EFFECT" to the vector \mathbf{b} .
- A red question mark is placed below the matrix \mathbf{A} .
- A red question mark is placed below the vector \mathbf{b} .

Fit criterion and search solution strategy – phase 1:

$$0 = \underline{a_{11}}x_{11} + a_{12}x_{21} + \dots + a_{1n}x_{n1} + \cancel{b_1u}$$

.....

$$0 = \underline{a_{11}}x_{1h} + a_{12}x_{2h} + \dots + a_{1n}x_{nh} + \cancel{b_1u}$$

Choose only the h experiments where gene 1 has not been perturbed and solve the eqs. for gene 1 with $a_{11} \neq 0$.

Repeat for all the genes and obtain matrix **A**

...we still did not say how to find the h experiments where gene i has not been perturbed, this is done simply by choosing those experiments where the gene has changed less...

Fit criterion and search solution strategy – phase II:

Say \mathbf{x}_d the expression profile following the drug treatment, then the predicted targets of the drug are:

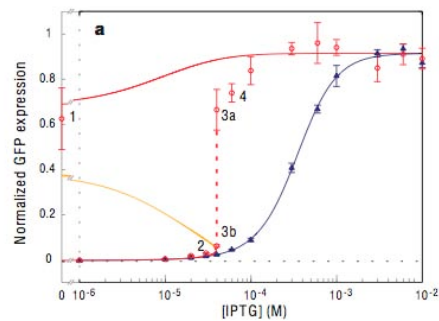
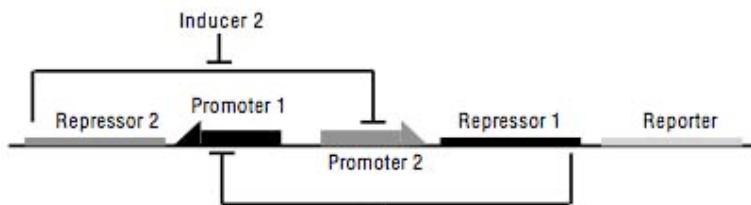
$$\mathbf{bu} = -\mathbf{Ax}_d$$

Other projects going on in our lab:

- REVERSE ENGINEERING BASED ON MI AND BAYESIAN APPROACH
- IDENTIFICATION OF DRUG TARGETS
- SYNTHETIC BIOLOGY - COBIOS (coordinator) 120kE/yr for 3 yrs for my lab

Construction of a genetic toggle switch in *Escherichia coli*

Timothy S. Gardner[†], Charles R. Cantor^{*} & James J. Collins^{††}



A synthetic oscillatory network of transcriptional regulators

Michael B. Elowitz & Stanislas Leibler

