



# 计算系统生物学

王勇

#### 中国科学院数学与系统科学研究院



http://zhangroup.aporc.org Chinese Academy of Sciences





## Transcriptional Regulatory Network Inference

#### Yong Wang

http://zhangroup.aporc.org



http://zhangroup.aporc.org Chinese Academy of Sciences

## Outline

- Background: Definition of TRN inference
- Inferring TRN from sequence's perspective.
- Inferring TRN from gene expression's perspective (Method: Inferelator)
- Inferring TRN from transcription complexes' perspective (Method: TRNInfer)

#### **Transcription in higher eukaryotes**

ZHANGroup



🍩 🕞 🚳

A **transcription factor**(sometimes called a sequence-specific <u>DNA</u>-binding factor) is a <u>protein</u> that binds to specific <u>DNA sequences</u>, thereby controlling the flow (or <u>transcription</u>) of genetic information from DNA to <u>mRNA</u>

ZHANGroup





#### **Transcriptional Regulation**



ZHANGroup



#### **Transcriptional Factor**



A defining feature of transcription factors is that they contain one or more <u>DNA-binding domains</u>(DBDs)





Α









#### Structure



Schematic diagram of the amino acid sequence (amino terminus to the left and carboxylic acid terminus to the right) of a prototypical transcription factor that contains

- (1) a DNA-binding domain (**DBD**), which attach to specific sequences of DNA (<u>enhancer</u> or <u>promoter</u> sequences) adjacent to regulated genes.
- (2) signal sensing domain (**SSD**), which senses external signals and in response transmit these signals to the rest of the transcription complex, resulting in up or down regulation of gene expression. An optional **domain** (*e.g.*, a ligand binding domain).

(3) a transactivation domain (**TAD**), which contain binding sites for other proteins such as <u>transcription coregulators</u>. These binding sites are frequently referred to as **activation functions (AFs**).





#### **Trans-activating domain**

MADIS

	Annotated 9aaTA	D Peptide - KIX interaction (NMR data)
p53 TAD1	E TFSD LWKL	LSPEET <u>FSDLWK</u> LPE
p53 TAD2	D DIEQ WFTE	QAMDDLMLSPDDIEQWFTEDPGPD
MLL	S DIMD FVLK	DCGNI <u>LPSDIMDFVL</u> KNTP
E2A	D LLDF SMMF	PVGTDKELSDLLDFSMMFPLPVT
Rtg3	E TLDF SLVT	E2A homolog
CREB	R KILN DLSS	<u>RR</u> EILSRRP <u>SY</u> RK <u>IL</u> N <u>DL</u> SS <u>DAP</u>
CREBaB6	E AILA ELKK	CREB-mutant binding to KIX
Gli3	D DVVQ YLNS	TAD homology to CREB/KIX
Gal4	D DVYN YLFD	Pdrl and Oafl homolog
Oaf1	D LFDY DFLV	DLFDYDFLV
Pip2	D FFDY DLLF	Oafl homolog
Pdr1	E DLYS ILWS	EDLYSILWSDWY
Pdr3	T DLYH TLWN	Pdr1 homolog

Nine-amino-acid transactivation domain (9aaTAD)





#### **DNA-binding domain**

Family	<u>InterPro</u>	<u>Pfam</u>	<u>SCOP</u>
basic-helix-loop-helix <sup>[43]</sup>	<u>IPR001092</u>	<u> Pfam PF00010</u>	<u>SCOP 47460</u>
basic-leucine zipper ( <u>bZIP</u> ) <sup>[44]</sup>	<u>IPR004827</u>	<u> Pfam PF00170</u>	<u>SCOP 57959</u>
C-terminal effector domain of the bipartite response regulators	<u>IPR001789</u>	<u> Pfam PF00072</u>	<u>SCOP 46894</u>
GCC box			<u>SCOP 54175</u>
helix-turn-helix <sup>[45]</sup>			
homeodomain proteins - bind to homeobox			
DNA sequences, which in turn encode other transcription factors. Homeodomain proteins play critical roles in the regulation of development. <sup>[46]</sup>	<u>IPR009057</u>	<u> Pfam PF00046</u>	<u>SCOP</u> <u>46689</u>
lambda repressor-like	<u>IPR010982</u>		<u>SCOP 47413</u>
srf-like ( <u>serum response factor</u> )	<u>IPR002100</u>	<u> Pfam PF00319</u>	<u>SCOP 55455</u>
paired box <sup>[47]</sup>			
winged helix	<u>IPR013196</u>	<u> Pfam PF08279</u>	<u>SCOP 46785</u>
zinc fingers <sup>[48]</sup>			
* multi-domain Cys <sub>2</sub> His <sub>2</sub> zinc fingers <sup>[49]</sup>	<u>IPR007087</u>	<u> Pfam PF00096</u>	<u>SCOP 57667</u>
* Zn <sub>2</sub> /Cγs <sub>6</sub>			<u>SCOP 57701</u>
* Zn <sub>2</sub> /Cys <sub>8</sub> <u>nuclear receptor</u> zinc finger	<u>IPR001628</u>	<u> Pfam PF00105</u>	<u>SCOP 57716</u>





- There are approximately 2600 proteins in the <u>human</u> <u>genome</u> that contain DNA-binding domains, and most of these are presumed to function as transcription factors.
- 10% of genes in the genome code for transcription factors, which makes this family the single largest family of human proteins.
- the combinatorial use of a subset of the approximately 2000 human transcription factors easily accounts for the unique regulation of each gene in the human genome during <u>development</u>.

# Regulatory mechanism

- stabilize or block the binding of RNA polymerase to DNA
- catalyze the <u>acetylation</u> or deacetylation of <u>histone</u> proteins. The transcription factor can either do this directly or recruit other proteins with this catalytic activity. Many transcription factors use one or the other of two opposing mechanisms to regulate transcription:<sup>[13]</sup>
  - <u>histone acetyltransferase</u> (HAT) activity acetylates <u>histone</u> proteins, which weakens the association of DNA with <u>histones</u>, which make the DNA more accessible to transcription, thereby up-regulating transcription
  - <u>histone deacetylase</u> (HDAC) activity deacetylates <u>histone</u> proteins, which strengthens the association of DNA with histones, which make the DNA less accessible to transcription, thereby down-regulating transcription
- recruit <u>coactivator</u> or <u>corepressor</u> proteins to the transcription factor DNA complex



#### **Transcriptional Regulation**





ZHANGroup

#### Transcriptional Regulation: output

) 🖓 🚳



### Perspective I: Cis-regulatory elements

#### Learning problems:

 Understand which regulators control which target genes



ZHANGrou

### Perspective II: Target gene expression

## Learning problems:

 Understand which regulators control which target genes



### Perspective III: Transcriptional complex

### Learning problems:

 Understand which TF complex control which target genes

RNA polymerase

(transcription)

Nuclear membrane Ribosome (translation)

- Estimate the TF complex activity
- Correlate the expression of target genes with TF complex activity
- Select the TF complex to explain the data







 Gene regulatory networks (GRN): indirect gene-gene interactions (genetic interactions)





## GRN and TRN ?

 Transcription regulatory networks (TRN): direct interactions between TFs and genes (physical interactions)





## GRN and TRN ?

• GRN:

mRNA x(t)  $\rightarrow$  mRNA x(t): indirect interactions

ZHANGroup OB



## Outline

- Background: Definition of TRN inference)
- Inferring TRN from sequence's perspective.
- Inferring TRN from gene expression's perspective (Method: Inferelator)
- Inferring TRN from transcription complexes' perspective (Method: TRNInfer)



#### • **TF** binding sites discovery

- Cluster genes by expression profile, annotation, ... to find potentially coregulated genes
- Find overrepresented motifs in promoter sequences of similar genes (algorithms: MEME, Consensus, Gibbs sampler, AlignACE, ...)





#### **TFBS and PWM?**

- Transcription factor binding sites (TFBSs) are usually slightly variable in their sequences.
- A positional weight matrix (PWM) specifies the probability that you will see a given base at each index position of the motif.

Pos	<u>1</u>	<u>2</u>	<u>3</u>	<u>4</u>	<u>5</u>	<u>6</u>	7	<u>8</u>	9	<u>10</u>	<u>11</u>	<u>12</u>	<u>13</u>	<u>14</u>	<u>15</u>
Α	18	8	5	4	1	29	7	7	7	0	1	39	1	1	6
С	8	3	3	9	33	4	21	15	14	0	0	1	43	39	18
G	13	31	34	9	8	10	11	15	19	4	44	3	0	1	6
Т	7	4	4	24	4	3	7	9	6	42	1	3	2	5	16
Con	Ν	G	G	Т	С	Α	Ν	Ν	Ν	Т	G	Α	С	С	Ν











Isolation of cells of the immune response

- Use the correct number of cells:  $1\times10^6$  to  $10\times10^6$
- · Collect biological replicates of cells
- Choose an appropriate control for antibody specificity (knockout or RNAi knockdown)

Fragmentation by sonication or MNase treatment

- Shear chromatin to a size range of ~150–300 bp
- Sonicate chromatin extracts for non-histone proteins
   \* Sonication conditions should be determined empirically for each cell type
- Treat chromatin extracts with MNase for analysis of histone modifications

\* Do not overdigest chromatin

ChIP analysis of histone modifications, transcription factors or epigenetic regulators

- · Select antibody: monoclonal versus polyclonal
- · Choose reference control (Input or IgG)
- Perform ChIP with established protocols
- Purify DNA

Library construction

- Do end repair and adapter ligation
- Perform PCR using primers compatible with sequencing platform
  - \* Avoid overamplifying DNA

Sequencing

- Determine sequencing depth on the basis of the prevalence of binding throughout the genome: more sequencing tags may be needed for diffuse signals (such as H3K27me3)
- Perform single-end or paired-end sequencing

ChIP-Seq is a powerful tool with which to investigate protein-DNA interactions on a global scale. It is important that the appropriate controls for antibody specificity be determined before ChIP-Seq is begun. After isolation of the ideal number of cells, chromatin is sheared into an ideal size range by sonication or enzymatic means (micrococcal nuclease (MNase)). Next, highquality antibodies are used for ChIP to enrich for factor-occupied DNA sequences. After purification of ChIP-enriched DNA, a library is constructed to allow sequencing on nextgeneration sequencing (NGS) platforms. Library construction typically includes endrepair, the addition of single adenosine residues, adaptor ligation and PCR with primers compatible with the sequencing platform. After cluster generation, single- or paired-end sequencing is performed on nextgeneration sequencing platforms. RNAi, RNA-mediated interference; bp, base pairs.

ChIP-Seq: technical considerations for obtaining highquality data Benjamin L Kidder, Gangqing Hu & Keji Zhao, Nature Immunology 12, 918–922 (2011)

Sequencing on NGS platforms

#### **Calculation of PWM**

- 1. acggcagggTGACCc
- 2. aGGGCAtcgTGACCc
- 3. cGGTCGccaGGACCt
- 4. tGGTCAggcTGGTCt
- 5. aGGTGGcccTGACCc
- 6. cTGTCCctcTGACCc
- 7. aGGCTAcgaTGACGt
- 41. cagggagtgTGACCc
- 42. gagcatgggTGACCa
- 43. aGGTCAtaacgattt
- 44. gGAACAgttTGACCc
- 45. cGGTGAcctTGACCc
- 46. gGGGCAaagTGACTg

Position frequency matrix (PFM) (also known as *raw count matrix*)

Given N sequence fragments of fixed length, one can assemble a position frequency matrix (number of times a particular nucleotide appears at a given position). A normalized PFM, in which each column adds up to a total of one, is a matrix of probabilities for observing each nucleotide at each position.

## Position weight matrix (PWM) (also known as *position-specific scoring matrix*)

PFM should be converted to log-scale for efficient computational analysis. To eliminate null values before log-conversion, and to correct for small samples of binding sites, a sampling correction, known as *pseudocounts*, is added to each cell of the PFM.





#### **Position Weight Matrix**

					Со	nverti	ing a	PFM i	nto a	PWN	l					
Α	18	8	5	4	1	29	) 7	7	7	7	0	1	39	1	1	6
С	8	3	3	9	33	4	2	1 '	15	14	0	0	1	43	39	18
G	13	31	34	9	8	10	) 1	1 '	15	19	4	44	3	0	1	6
Т	7	4	4	24	4	3	7	7	9	6	42	1	3	2	5	16
For elei	each ment	matriz do:	х <i>w</i> (	(b,i) =	$\log_2 \frac{1}{2}$	$\frac{p(b,i)}{p(b)}$	$= \log_2$	$f_{b,}$	$\frac{1}{i} + \frac{\sqrt{4}}{4}$ $\frac{1}{p(b)}$	V		T	Ι			-
А	0.58	-0.44	-0.98	-1.21	-2.29	1.22	-0.60	-0.60	-0.60	-2.96	-2.29	1.62	-2.29	-2.29	-0.72	
С	-0.44	-1.49	-1.49	-0.30	1.39	-1.21	0.78	0.34	0.25	-2.96	-2.96	-2.29	1.76	1.62	0.46	
G	0.16	1.31	1.44	-0.30	-0.44	-0.17	-0.06	0.34	0.65	-1.21	1.79	-1.49	-2.96	-2.29	-0.64	
Т	-0.60	-1.21	-1.21	0.96	-1.21	-1.49	-0.60	-0.30	-0.78	1.73	-2.29	-1.49	-1.84	-0.98	0.23	

 $f_{b,i}$  – raw count (PFM matrix element) of nucleotide  $m{b}$  in column  $m{i}$ 

N – number of sequences used to create PFM (= column sum)

 $\frac{\sqrt{N}}{4}$  and  $\sqrt{N}$  - pseudocounts (correction for small sample size)

p(b) - background frequency of nucleotide b, this one usually defaults to 0.25

Hertz GZ, Stormo GD. Bioinformatics (1999)



**ZHANGroup** 

#### **TABLE 4.1.** Several Databases of TF Binding Sites

# Scoring putative transcriptional regulation by scanning the promoter with PWM

### <u>GGGTCAGCATGGCCA</u>

**ZHANGroup** 

Row

А	0.58	-0.44	-0.98	-1.21	-2.29	1.22	-0.60	-0.60	-0.60	-2.96	-2.29	1.62	-2.29	-2.29	-0.72
С	-0.44	-1.49	-1.49	-0.30	1.39	-1.21	0.78	0.34	0.25	-2.96	-2.96	-2.29	1.76	1.62	0.46
G	0.16	1.31	1.44	-0.30	-0.44	-0.17	-0.06	0.34	0.65	-1.21	1.79	-1.49	-2.96	-2.29	-0.64
Т	-0.60	-1.21	-1.21	0.96	-1.21	-1.49	-0.60	-0.30	-0.78	1.73	-2.29	-1.49	-1.84	-0.98	0.23

Absolute score of the site 
$$S = \sum_{i=1}^{m} w(b, i) = 11.57$$

 Max
 0.58
 1.31
 1.44
 0.96
 1.39
 1.22
 0.78
 0.34
 0.65
 1.73
 1.79
 1.62
 1.76
 1.62
 17.20

 Min
 -0.60
 -1.49
 -1.21
 -2.29
 -1.49
 -0.60
 -0.78
 -2.96
 -2.96
 -2.29
 -24.02

 $relative\_score = \frac{Absolute\_score-Minimum\_score}{Maximum\_score-Minimum\_score}$ 

 $=\frac{11.57 - (-24.02)}{17.20 - (-24.02)} = 0.86$ 





**ZHANGrou** 



• A consensus logo for the LexA-binding motif of several Gram-positive species.





# Binding sites database

Name	Organisms	Source	Access	URL
<u>RegTransBase</u>	Prokaryotes	Expert/literatur e curation	Public	[1]
<u>RegulonDB</u>	Escherichia coli	Expert curation	Public	[2]
PRODORIC	Prokaryotes	Expert curation	Public	[3]
TRANSFAC	Mammals	Expert/literatur e curation	Private	[4]
TRED	Human, Mouse, Rat	Computer predictions, manual curation	Public	[5]
DBSD	Drosophila species	Literature/Expe rt curation	Public	[6]
носомосо	Human	Literature/Expe rt curation	Public	[7],[8]



Program	Description
MatInspector	Utilizes a large library of matrix descriptions for TFBSs to locate matches in DNA sequences
MATCH	Uses a library of mononucleotide or dinucleotide weight matrixes from TRANSFAC 3.5 for searching potential TFBSs
YMF	Does an enumerative search to find the motifs with the highest $z$ scores
MotifSampler	Uses Gibbs sampling to find the PWM that represents the motif by modeling
	the background with a higher-order Markov model
PhyloScan	Uses evidence from matching sites found in cross-species to identify TFBSs
ANN-Spec	Uses an artificial neural network and a Gibbs sampling method to model the specificity of a DNA-binding protein
CONSENSUS	Searches for the PWM with the maximum information content
Weeder	Enumerates all the oligos of (or up to) a given length and determines their occurrences with possible substitutions in the input sequences
AlignACE	Uses Gibbs sampling algorithm to find a series of motifs as PWMs that are overrepresented in the input sequences
MEME	Uses EM algorithm to optimizes the $E$ value of a statistic related to the information content of the motif
GLAM	Uses a Gibbs sampling-based algorithm that optimizes the alignment width and obtains the best possible gapless multiple alignment

#### TABLE 4.2. Some Software for Searching TF Binding Sites

🍪 🕀 🊱



Databases	Websites
SCPD	http://rulai.cshl.edu/SCPD
CEPDB	http://rulai.cshl.edu/cgi-bin/CEPDB
LSPD	http://rulai.cshl.edu/LSPD
PlantProm DB	http://mendel.cs.rhul.ac.uk/mendel.php?topic=plantprom
EPD	http://www.epd.isb-sib.ch
CSHLmpd	http://rulai.cshl.edu/CSHLmpd2
MPromDb	http://bioinformatics.med.ohio-state.edu/MPromDb
OMGProm	http://bioinformatics.med.ohio-state.edu/OMGProm
HemoPDB	http://bioinformatics.med.ohio-state.edu/HemoPDB
OPD	http://www.opd.tau.ac.il/
HPD	http://zlab.bu.edu/mfrith/HPD.html
DCPD	http://www-biology.ucsd.edu/labs/Kadonaga/DCPD.htm
TiProD	http://tiprod.cbi.pku.edu.cn:8080/index.html
DBTSS	http://dbtss.hgc.jp/

ZHANGroup

#### **TABLE 4.3. Databases of Promoters and TSSs**





## Outline

- Background: Definition of TRN inference
- Inferring TRN from sequence's perspective.
- Inferring TRN from gene expression's perspective (Method: Inferelator)
- Inferring TRN from transcription complexes' perspective (Method: TRNInfer)



MADIS



ZHANGroup

Target gene expression





### Structure learning

- Learn structure of "regulatory network", "regulatory modules", etc.
- Fit interpretable model to training data
- Many computational and statistical challenges; often used for qualitative hypotheses rather than prediction



(Segal et al, 2003, 2004)
#### A list of relevant computational methods

🍪 🖓 🊳

Name	Description	Reference
GRAM	Searches for co-bound genes with a strict cutoff. Then relaxes cutoff for genes that co-express with the original set.	Bar-Joseph et al, 2003
SAMBA	Discretizes expression and binding data into gene properties. Algorithm then looks for genes with statistically significant common property sets.	Tanay et al, 2003
ReMoDiscovery	Stringent and relaxed two step procedure that combined motif, expression, and ChIP-chip data.	Lemmens et al, 2006
COGRIM	Uses a Bayesian network to model expression level as a function of transcription factor expression and binding.	Chen et al, 2007
Inferelator	Uses biclustering to group co-expressed genes and then machine learning to infer regulatory influence from RNA and protein expression levels.	Bonneau et al, 2006

**ZHANGroup** 

## **Differential Equation Models**

- Attempt to reconstruct the dynamical system that produced the gene expression data
  - Reduce dimensionality of the data
  - Approximate dynamics
    - Modeled using ordinary differential equations
  - Restrict model complexity
- Example system : The Inferelator



- Regulators (genes and environment)
  - Limited to transcription factors
  - Factors with correlated profiles are merged
- Genes
  - Clustered based on putative coregulation
  - Used cMonkey to form biclusters across genes and conditions [Bonneau, 2006]
    - Correlated expression
    - Shared regulatory sequence motifs





## Model Details

• Expression of *y* (*gene or bicluster mean*) is influenced by the expression of N regulators:

 $X = (X_1, X_2, ..., X_N)$ 

$$\tau \frac{dy}{dt} = -y + g(\beta \cdot Z)$$

 $Z = (z_1[X], z_2[X] \dots z_P[X])$ 



-



## Model Details

$$\tau \frac{dy}{dt} = -y + g\left(\beta \cdot Z\right) \qquad \qquad Z = (z_1[X], z_2[X] \dots z_P[X])$$

#### **Choice of Squashing Function**

压缩函数(Squashing Function)

$$g(\beta \bullet Z) = \frac{1}{1 + e^{-\beta \bullet Z}} \qquad g(\beta Z) = \begin{cases} \beta Z : & \text{if } \min(y) < \beta Z < \max(y) \\ \max(y) : & \text{if } \beta Z > \max(y) \\ \min(y) : & \text{if } \beta Z < \min(y) \end{cases}$$



### Model Details

 $\tau \frac{dy}{dt} = -y + g(\beta \cdot Z)$ 

 $Z = (z_1[X], z_2[X] \dots z_P[X])$ 

**Choice of Z:** 

MADIS

 $\beta \mathbf{Z} = \beta_1 x_1 + \beta_2 x_2 + \beta_3 \min(x_1, x_2)$ 







### Model Details

$$\tau \frac{dy}{dt} = -y + g(\beta \cdot Z)$$

Steady state

$$y = g(\beta \cdot Z_{SS})$$

Time course

$$\tau \frac{y_{m+1} - y_m}{\Delta t_m} + y_m = g(\sum_{j=1}^{P} \beta_j z_{mj}) \quad for \quad m = 1, 2, ..., T - 1$$

## Model Learning with LASSO

• LASSO, a.k.a. L1 shrinkage

$$\left(\hat{\alpha}, \hat{\beta}\right) = \underset{\alpha, \beta}{\operatorname{arg\,min}} \left\{ \sum_{i=1}^{N} \left( y_i - \alpha - \sum_{j=1}^{p} \beta_j z_{ij} \right)^2 \right\} \qquad \text{S.T.} \qquad \sum_{j=1}^{p} \left| \beta_j \right| \le t \left| \beta_{ols} \right|$$



(Bonneau, et al, Genome Biology, 2006)

**ZHANGroup** 

## Results



The inferred regulatory network of Halobacterium NRC-1

Regulators are indicated as circles

Target gene biclusters are indicated by rectangles



Journal club



## A Predictive Model for Transcriptional Control of Physiology in a Free Living Cell

Richard Bonneau et. al. Cell, Vol 131, 1354-1365, 28 Dec. 2007

Institute for Systems Biology, Seattle, WA 98103, USA Center for Genomics & Systems Biology, New York University, New York, NY 10003, USA On the cover: Brightly colored blooms of halophlic (喜盐的) organisms in the salt flats of the South San Francisco Bay (photograph by Michael Melford, courtesy Getty Images) serve as a vibrant backdrop(背景) for a segment of a predictive environmental and gene regulatory circuit determined for one of this ecosystem's principal inhabitants,

ZHANGroup

the archaeon(古代生物) Halobacterium salinarium NRC-1(一种嗜盐的古生菌, 一般 只生存在盐水池塘或是盐湖中). This organism possesses a number of fascinating adaptations for life in hypersaline (高盐) environments including the production of membrane pigments (细胞膜上产生色素) that mediate lightdriven energy production and flotation devices called gas vesicles for vertical mobility in search of favorable oxic regimes (氧载体). While the availability of unique adaptations is important, the integrated regulation of these and many other core physiological processes (生理学 过程) is vital for survival in this dynamic environment. In this issue, Bonneau et al. report a systems level regulatory circuit for the transcriptional control of 80% of all genes in this organism. This regulatory model accurately predicts the transcriptional changes that occur when Halobacterium is challenged with new environmental and genetic perturbations. Significantly, this study supports the claim that fundamental properties of biological systems and their environments should enable the rapid construction of highly accurate, predictive models of global gene regulation for both traditional model systems and for many more currently uncharacterized organisms.



# Faculty of 1000 Biology

 "This paper represents an exceptionally important milestone in the field..." Evaluated by Faculty of 1000 Biology member Charles Auffray (Centre National de la Recherche Scientifique (CNRS) - UMR 7091, France)

"Faculty of 1000 Biology"创办于2002年1月,根据全球资深科学家的意见,提供对近期发表的生物科学论文的快速评论,目的是帮助广大科研人员遴选和发现有价值的研究工作。





## Other comments

• Research Highlight by *Nature Reviews Microbiology* 6, 92 (February 2008)

• **Bio-IT World's Systems Biology newsletter.**(In the closing days to 2007, a really nice piece of systems biology work was published in the journal *Cell* )



# Why this paper

- Cell publish computational biology work
- From Institute for Systems Biology, Seattle
- The ISB founder, also the founder of systems biology Lee Hood is one of the co-authors.
- To taste the flavor of systems biology (network+perturbation+data integration)





# What they studied

- A largely uncharacterized organisms
- Easy to be cultured
- The environment significantly influences the dynamic expression



## Methodology

#### • Experiments:

- 1. Microarray data: Total 413 experiments (Time-course and steady state, 8 environment effects perturbation, combinatorial perturbation. 33 gene deletion and GTF overexpression)
  - 266 experiments in training set
  - 147 new experiments
- 2. ChIP-chip data

http://baliga.systemsbiology.net/egrin.php

### Computational prediction

- 1. Protein structure prediction
- 2. Function annotation algorithm
- 3. Biclustering algorithm (Data integration and dimensionality reduction)
- 4. Transcriptional regulatory network inference







### Outline

- Background: Definition of TRN inference
- Inferring TRN from sequence's perspective.
- Inferring TRN from gene expression's perspective (Method: Inferelator)
- Inferring TRN from transcription complexes' perspective (Method: TRNInfer)





### Motivation

- TF activity level cannot be measured directly by microarray due to post-translational modifications
- Most existing algorithms has an implicit assumption that TFAs are proportional to their mRNA levels (like the previous example)
- TF generally regulates a gene with many collaborators (Transcription complex)





#### Cellular post-translational modifications

This schematic figure shows the location and role of a selection of some of the most important of more than 200 types of posttranslational modification (PTM). PTMs are found on all types of protein, from nuclear transcription factors to metabolic enzymes, structural proteins and plasma membrane receptors. PTMs affect the physicochemical properties of proteins, which provides a mechanism for the dynamic regulation of molecular self-assembly and catalytic processes through the reversible molecular recognition of proteins, nucleic acids, metabolites, carbohydrates and phospholipids. Ac, acetyl group; GPI, glycosylphosphatidylinositol; Me, methyl group; P, phosphoryl group; Ub, ubiquitin.













- Use TF-TG relation benefit the regulatory network identification
- TF expression level is not a good measure of the TF activity. The activated protein level of a TF, rather than its expression level, is what controls gene expression.
- The activity of a transcription factor is regulated according to the cell's need, largely through signal transduction. It may not be directly observed, but can be reflected by the genes it regulates.



MADIS

ZHANGroup







### Framework for TRNinfer



Wang et al. Bioinformatics, 2007



Transcription Regulatory network

#### • The general form

The transcription processes can be represented by differential equations with gene expression and TFAs:

$$\dot{x}(t) = f(a(t)) - Kx(t) \tag{1}$$

where  $x(t) = (x_1(t), \dots, x_m(t))^T$  is gene expression level (RNA),  $a(t) = (a_1(t), \dots, a_c(t))$  denotes TF activity level (Protein).

#### • The linear form

the linear form of (1) is

$$\dot{x}(t) = Ja(t) + b(t) \tag{2}$$

where  $J = [J_{ij}]_{m \times c} = \partial f(a) / \partial a$  is an  $m \times c$  Jacobian matrix or connectivity matrix.



### Approximating TF activity

- TFs and many cooperative proteins regulate a gene by a transcription complex (TC).
- TF activity depends on TC.
- A TC is formed by a series of biochemical reactions:

$$A_0 + A_1 + A_2 \rightleftharpoons_{k_{-1}}^{k_1} A$$



### Approximating TF activity

• According to the law of mass action,

the governing equations of the above reactions are given by

ZHANGroup

$$\frac{da_i}{dt} = -k_1 a_0 a_1 a_2 + k_{-1} a \quad \text{for } i = 0, 1, 2,$$
$$\frac{da}{dt} = k_1 a_0 a_1 a_2 - k_{-1} a \quad .$$

• TF activity can be given

$$a = k_0 a_0 a_1 a_2 \approx k_1 x_0 x_1 x_2$$

#### a : TF activity x : gene expression





## LP model

For all L datasets, J should be as consistent as possible with all datasets, which can be achieved by

$$\min_{J} \sum_{k=1}^{L} |\dot{X}^{k} - JA^{k}| + \lambda |J|.$$
(10)

where the first term is to minimize the error between real data and the reconstructed model, whereas the second term is the sparsity term which forces J sparse by using  $L_1$  norm.



Experimental results

- In the budding yeast *S. cerevisiae*, ChIP-chip experiments have been utilized to elucidate the binding interactions between 6270 genes and 113 preselected TFs.
- By checking yeast protein complexes in MIPS, we found 26 TFs in transcriptional protein complexes.
- Among these 26 TFs, some are related to yeast cell cycle and some are related to polyphosphate metabolism in S. cerevisiae





### Yeast cell cycle data

- There are 11 TFs that are known to be related to cell-cycle regulation, among which 5 TFs are in 4 different TCs.
- Except these 5 TFs, we selected 8 genes that are closely related to cell cycle based on the information in YEASTRACT (<u>http://www.yeastract.com/index.php</u>).
- According to the gene expression data from Spellman et al., we generated 4 datasets with the number of time points as 18, 17, 24, and 14 respectively.

TFs	TCs	protein members
MBP1	510.190.70	MBP1 SWI6
MCM1	510.190.120	ARG82 ARG81 ARG80 MCM1
STB1	510.190.150	STB2 STB1 RPD3 SIN3
SWI4	510.190.60	SWI4 SWI6
SWI6	510.190.60	SWI4 SWI6

Table 3: TFs related to yeast cell cycle and their TCs.





### Yeast cell cycle data



The inferred yeast cell cycle transcriptional regulatory network. The red arrows in the figure indicate repression while the blue arrows indicate activation. The comparison results of LP method based on transcription complexes (LP TC), LP method based on only mRNA levels of TFs (LP mRNA) and SVD method based on mRNA levels of TFs (SVD mRNA). (a) on yeast cell cycle data set; (b) on yeast polyphosphate metabolism data set.

**ZHANGroup** 

MADIS



### Yeast cell cycle data

• We can check the periodicity of the activity levels of the TFs (or TCs) because it is believed that the activities of TFs related to cell cycle tend to be periodic. This fact can be confirmed by Fisher's g-test.

ZHANGrout

Table 3. The P-values of the periodicity for some TFs related with cell cycle

TFs	Experiment conditions	Expression	Activity
MBP1	alpha0min-alpha119min	0.525	0.003
SWI4	alpha0min-alpha119min	0.0064	0.00019
SWI6	alpha0min-alpha119min	0.367	0.00019
SWI4	cdc1510min-cdc15290min	0.132	0.01
SWI6	cdc1510min-cdc15290min	0.024	0.01

### Experimental results ---Polyphosphate metabolism data

- Among the TFs related to polyphosphate metabolism verified by the ChIP experiments, there are 14 TFs in 9 different TCs.
- Gene expression data: Ogawa N, DeRisi J, Brown PO (2000).
- Among the genes in this dataset, some genes of those with change of 2 fold up or down in at least two time points of the expression levels are believed to be closely related to polyphosphate metabolism.
- In such a way, totally 64 genes (including 14 TFs) form a test data

#### Polyphosphate metabolism data

MADIS

**ZHANGroup** 

Table 4: TFs related to polyphosphate metabolism and their TCs.

TFs	TCs	protein members
RTG1	510.190.130	RTG3 RTG1
RTG3	510.190.130	RTG3 RTG1
MET4	510.190.160.30	MET32 MET28 MET4
MET31	510.190.160.20	MET28 MET4 MET31
LEU3	510.190.210	LEU3
HAP5	510.160	HAP3 HAP2 HAP4 HAP5
HAP4	510.160	HAP3 HAP2 HAP4 HAP5
HAP3	510.160	HAP3 HAP2 HAP4 HAP5
GCR2	510.190.90	GCR2 GCR1
GCR1	510.190.90	GCR2 GCR1
GAL4	510.190.80	GAL3 GAL80 GAL4
CBF1	510.190.160.10	MET28 CBF1 MET4
ARG80	510.190.120	ARG81 ARG80 MCM1
ARG81	510.190.120	ARG81 ARG80 MCM1




## Polyphosphate metabolism data



Transcriptional regulatory network for polyphosphate metabolism. The red arrows in the figure indicate repression while the blue arrows indicate activation.

## ZHANGroup on

## Take-home messages

- Looking at the same transcriptional regulatory interactions from different perspectives.
- For inferring a TRN, one must first determine which genes or proteins are TFs.
- Furthermore, it is also very difficult to measure the protein concentration levels of TFs and determine their regulatory effects on gene transcription.
- The interactions or cooperations between multiple TFs and their coregulators is a big challenge
- We develop TRNinfer for inferring transcriptional networks by using transcription complexes. <u>http://zhangroup.aporc.org/ResourceBioinformatics</u>