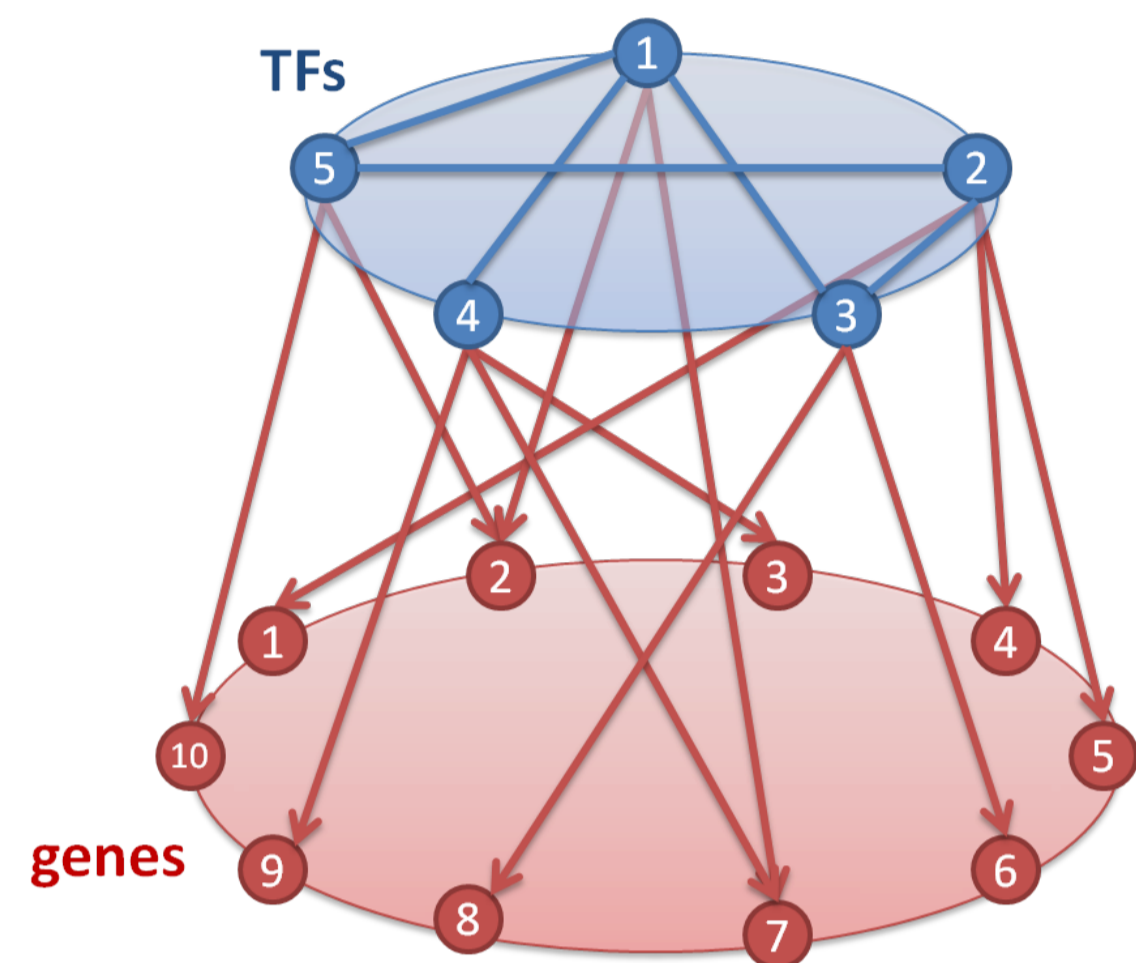


## Introduction

- ▶ TF-gene interactions can be used to learn TF-TF interactions.
- ▶ **Assumption:** if TF ‘A’ regulates TF ‘B’, then it is more likely that TF ‘A’ and TF ‘B’ will regulate similar sets of non-TF genes.
- ▶ **Idea:** Represent each TF as a vector of TF-gene interactions.



	TF <sub>1</sub>	TF <sub>2</sub>	TF <sub>3</sub>	TF <sub>4</sub>	TF <sub>5</sub>
g <sub>1</sub>	0	1	0	0	0
g <sub>2</sub>	1	0	0	0	1
g <sub>3</sub>	0	0	0	1	0
g <sub>4</sub>	0	1	0	0	0
g <sub>5</sub>	0	1	0	0	0
g <sub>6</sub>	0	0	1	0	0
g <sub>7</sub>	1	0	0	1	0
g <sub>8</sub>	0	0	1	0	0
g <sub>9</sub>	0	0	0	1	0
g <sub>10</sub>	0	0	0	0	1

- ▶ When TF-gene edges are not known in advance, these vectors can be *estimated* by regressing gene expression on TF expression.

## Learning TF-TF Interactions from TF-Gene Interactions

- ▶ Many algorithms for learning gene regulatory networks (such as relevance networks [1], ARACNE [2] and CLR [3]) compute pairwise measurements of similarity between random variables (typically, mutual information).

	$d$ variables (TFs)			
	$x_1$	$x_2$	...	$x_d$
$n$ samples	$x_1^{(1)}$	$x_2^{(1)}$	...	$x_d^{(1)}$
	$x_1^{(2)}$	$x_2^{(2)}$	...	$x_d^{(2)}$
	...	...	...	...
	$x_1^{(n)}$	$x_2^{(n)}$	...	$x_d^{(n)}$

- ▶ The usual approach consists in working with a matrix of microarray data where columns correspond to TFs and rows correspond to different samples.

$$\mathcal{L} = \{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n)}\}, \forall i, \mathbf{x}^{(i)} \in \mathbb{R}^d$$

	$d$ variables (TFs)			
	$y_1$	$y_2$	...	$y_d$
$m$ target genes	$y_1^{(1)}$	$y_2^{(1)}$	...	$y_d^{(1)}$
	$y_1^{(2)}$	$y_2^{(2)}$	...	$y_d^{(2)}$
	...	...	...	...
	$y_1^{(m)}$	$y_2^{(m)}$	...	$y_d^{(m)}$

- ▶ When TF-gene interactions are known, we propose to use an alternative data matrix where  $\mathbf{Y}_i^{(j)} = 1$  if TF  $i$  regulates gene  $j$ , and  $\mathbf{Y}_i^{(j)} = 0$  otherwise.

$$\mathcal{L}^* = \{\mathbf{y}^{(1)}, \dots, \mathbf{y}^{(m)}\}, \forall i, \mathbf{y}^{(i)} \in \{0, 1\}^d$$

- ▶ When TF-gene interactions are unknown, they can be estimated from microarray data. For each gene target  $t$ , solve

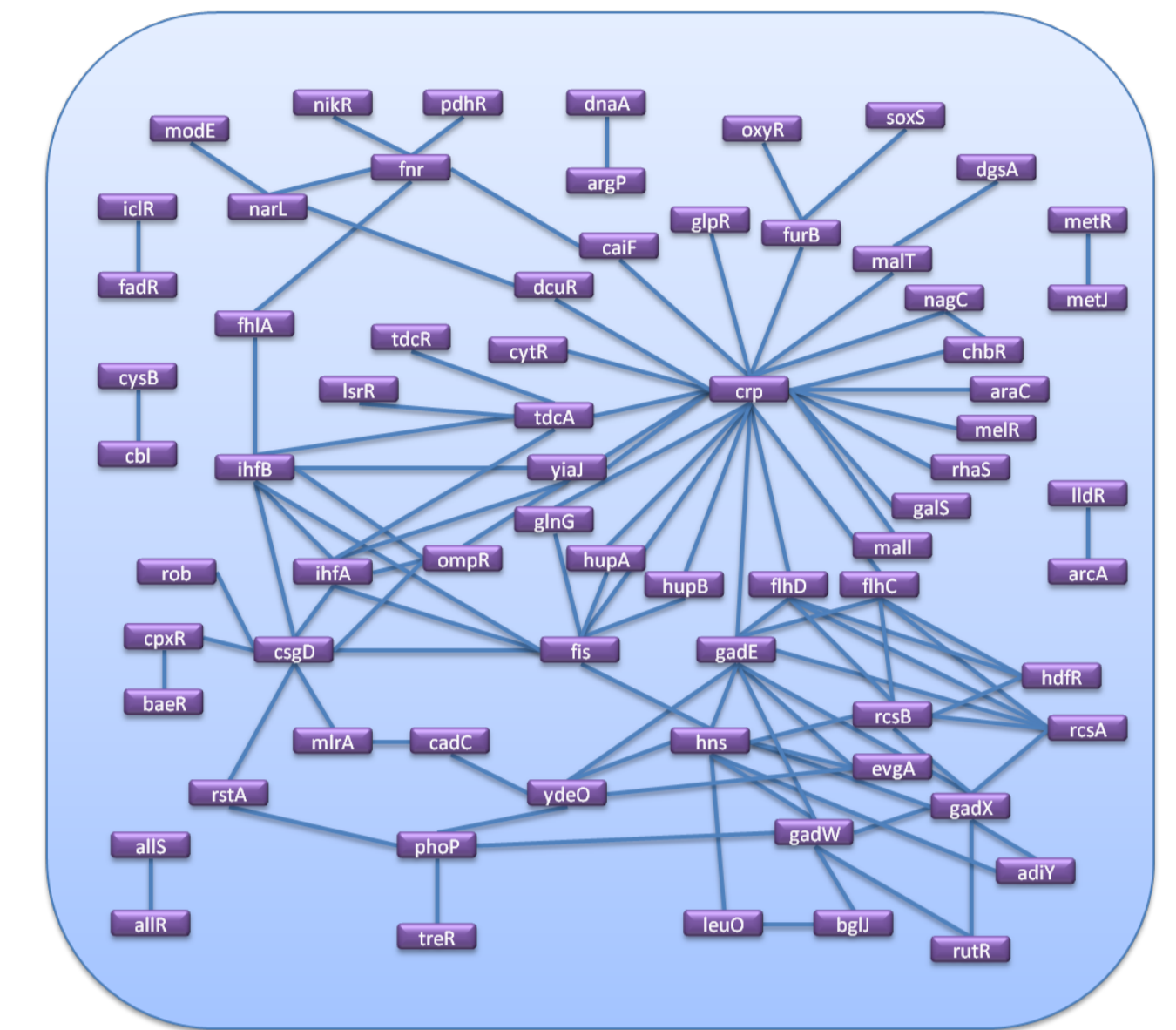
$$\min_{\beta^{(t)}} \|\mathcal{X} \cdot \beta^{(t)} - \mathbf{x}^{(t)}\|_2$$

where  $\mathbf{x}^{(t)} \in \mathbb{R}^n$  is the expression data for gene  $t$ ,  $\mathcal{X} \in \mathbb{M}_{n \times d}$  is the matrix of TF expression and  $\beta^{(t)} \in \mathbb{R}^d$ .

	$d$ variables (TFs)			
	$\hat{y}_1$	$\hat{y}_2$	...	$\hat{y}_d$
$m$ target genes	$ \beta_1^{(1)} $	$ \beta_2^{(1)} $	...	$ \beta_d^{(1)} $
	$ \beta_1^{(2)} $	$ \beta_2^{(2)} $	...	$ \beta_d^{(2)} $
	...	...	...	...
	$ \beta_1^{(m)} $	$ \beta_2^{(m)} $	...	$ \beta_d^{(m)} $

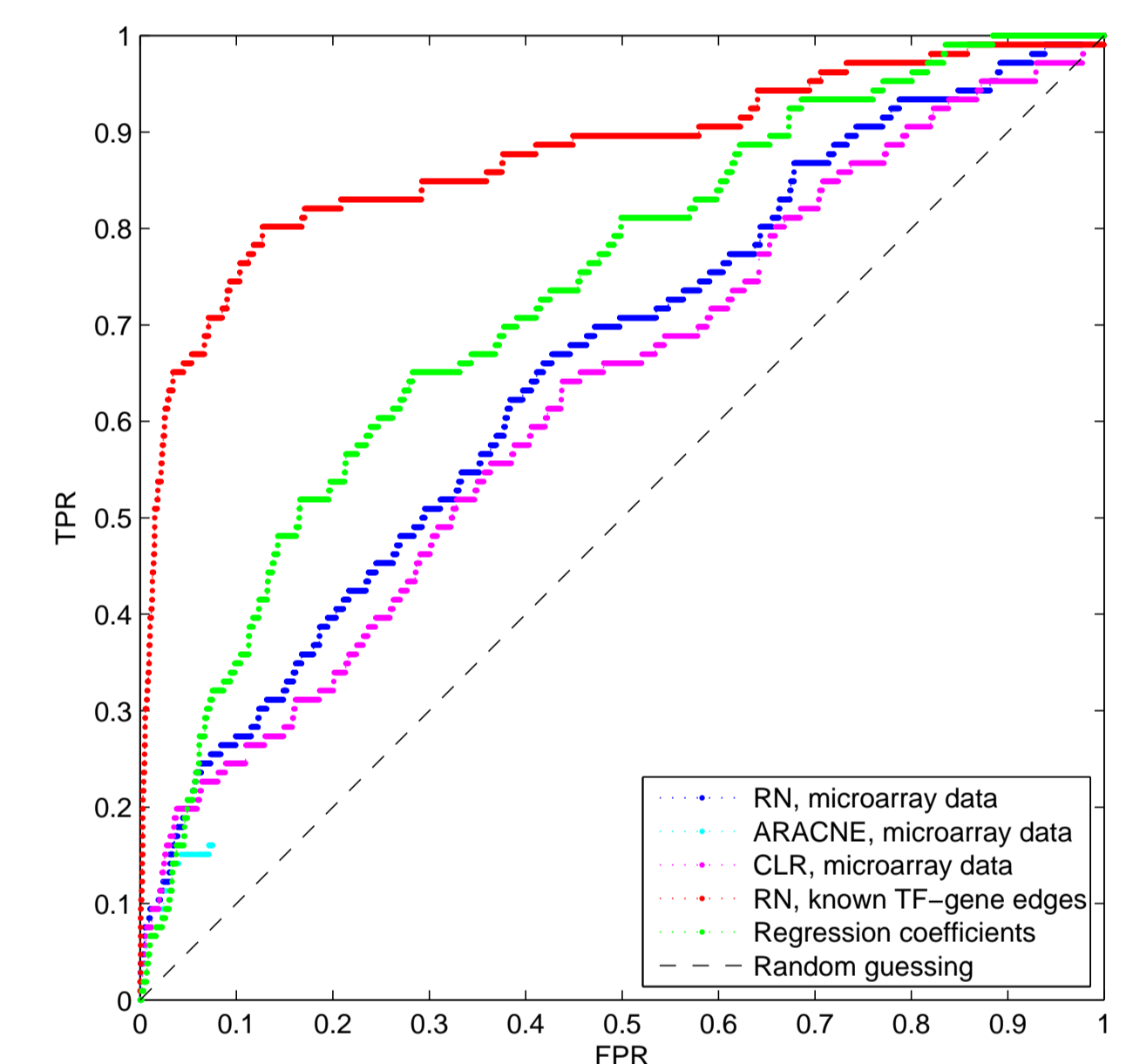
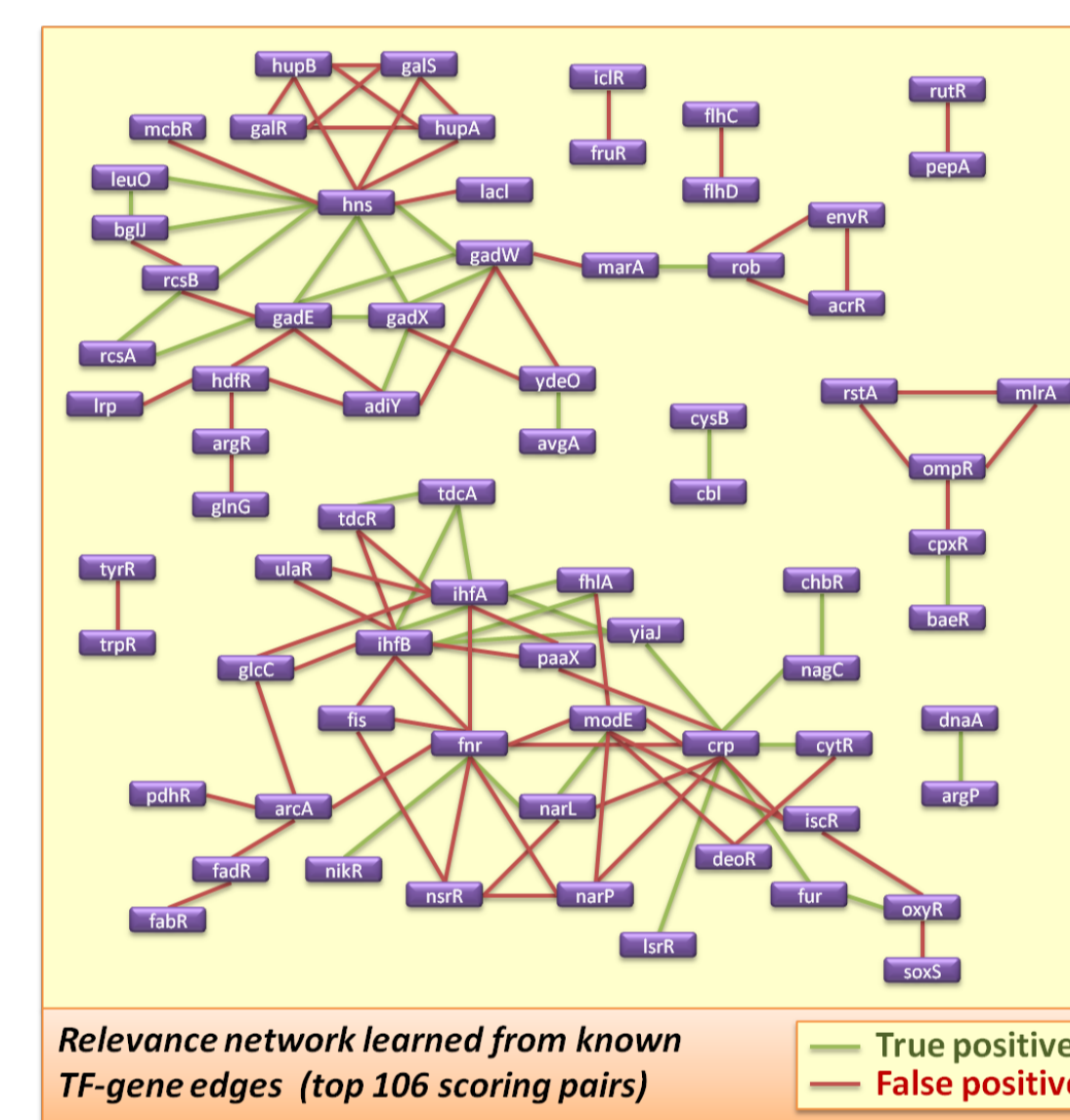
## Transcriptional Networks in *E. coli*

- ▶ Ground truth network from RegulonDB database [4], which contains 106 TF-TF interactions and 2,109 TF-gene interactions involving  $d = 126$  TFs and  $m = 984$  genes.



- ▶ Expression data from the Many Microbe Microarrays Database (M3D) [5], which contains  $n = 466$  microarray samples.

- ▶ We measured ROC performance for edgewise network reconstruction accuracy using the three types of data matrix representation.
- ▶ For the regression coefficients case, we simply ranked all pairwise Euclidean distances between columns.



## Discussion

- ▶ New approach intended to improve (not replace) existing TF-TF network reconstruction algorithms.
- ▶ By estimating TF-gene edges, we look jointly at microarray data for TFs and non-TF target genes (as opposed to alternatives that learn TF-TF networks using TF expression alone).
- ▶ Non-penalized linear regression was used only for illustration purposes. Sparse regression techniques may lead to better results, possibly closing the gap between the green and red ROC curves.
- ▶ Basis for two-phase network learning strategy: first, learn TF-gene edges using regression and then learn TF-TF edges as graphical models using the vectors of regression coefficients.

## References

- [1] Butte, A. J. and Kohane, I. S.: "Mutual information relevance networks: functional genomic clustering using pairwise entropy measurements," *Pac. Symp. Biocomputing*, 2000.
- [2] Margolin, A. et al.: "ARACNE: An algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context," *BMC Bioinformatics*, vol. 7, 2006.
- [3] Faith, J. J. et al.: "Large-Scale mapping and validation of *Escherichia coli* transcriptional regulation from a compendium of expression profiles," *PLoS Biol.*, vol. 5, n. 1, 2007.
- [4] Gama-Castro S. et al.: "RegulonDB (version 7.0): Transcriptional regulation of *Escherichia coli* K-12 integrated within genetic sensory response units (gensor units)," *Nucleic Acids Research*, 2010.
- [5] Faith, J. J. et al.: "Many Microbe Microarrays Database: uniformly normalized Affymetrix compendia with structured experimental metadata," *Nucleic Acids Research*, 2008.