

京都府
日本



Annotating Gene Function by Combining Expression Data with a Modular Gene Network

Motoki Shiga, Ichigaku Takigawa,
Hiroshi Mamitsuka



Bioinformatics Center, ICR, Kyoto University, Japan

Table of Contents

1. Motivation

Gene expression analysis, Gene clustering

2. Proposed method

Probabilistic Model, Clustering Algorithm

3. Experiments

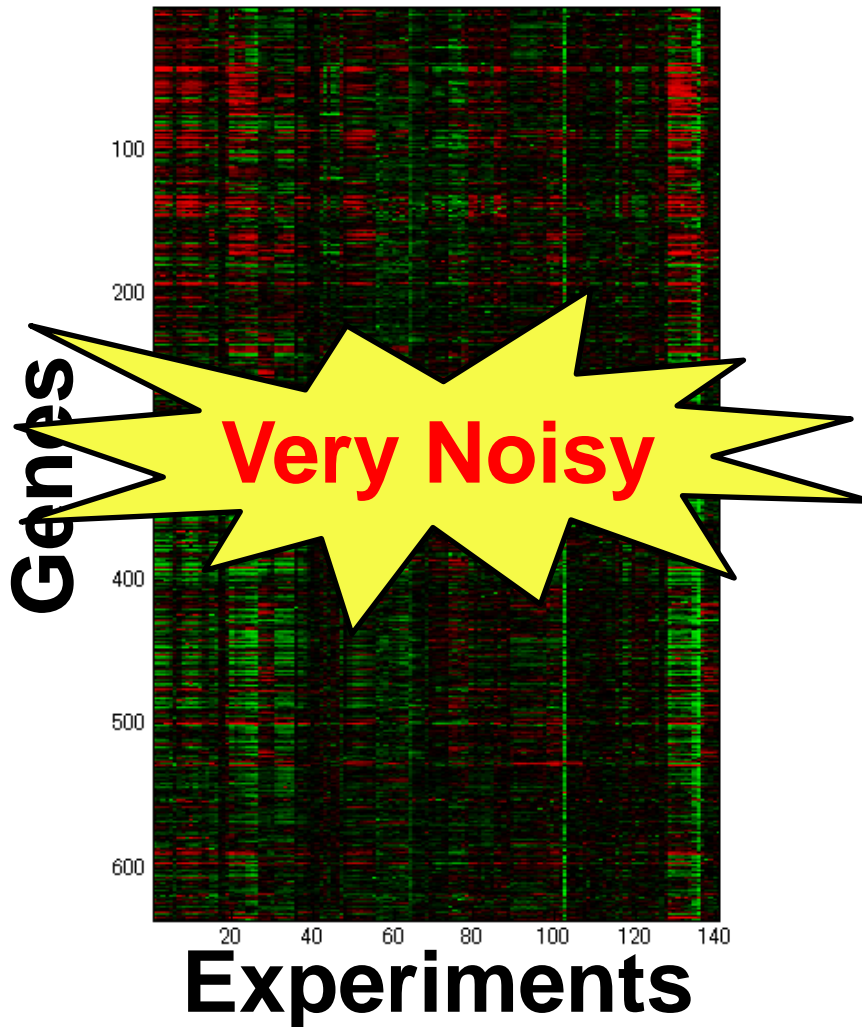
Datasets, Evaluation measure, Standard cluster

4. Summary

Gene Clustering using Expression Data

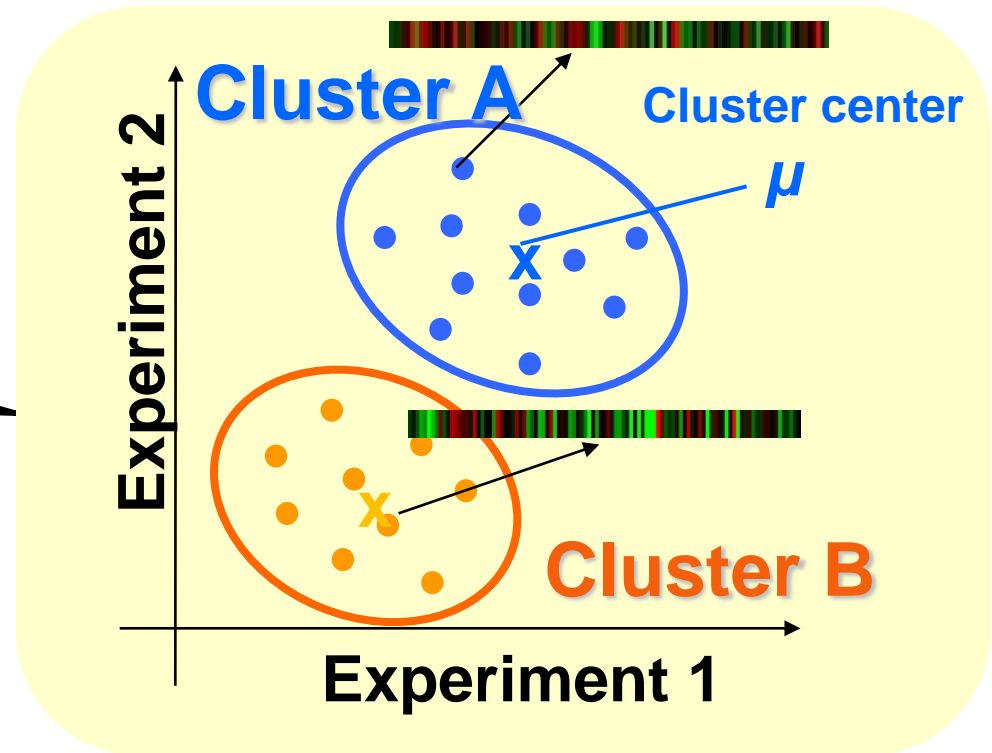
Expression Data

(Gasch, et al., Mol. Biol. Cell, 2000)



Clustering analysis

Ex. Hierarchical clustering,
K-means, Self-Organized Map



To improve accuracy,
use additional information

Gene Expression Analysis Using Additional Information

▶ Related work

Gene Ontology (GO) \Rightarrow a priori probability distribution
(Pan, Bioinformatics, 2006)

GO \Rightarrow Corrected distance
(Huang and Pan, Bioinformatics, 2006)

Metabolic pathway \Rightarrow Length of the shortest path
(Hanisch, et al., Bioinformatics, 2002)

▶ Proposed method

Metabolic pathway

\Rightarrow **Hidden modular random fields**
(Network modularity)

2. Proposed Method

Probabilistic model

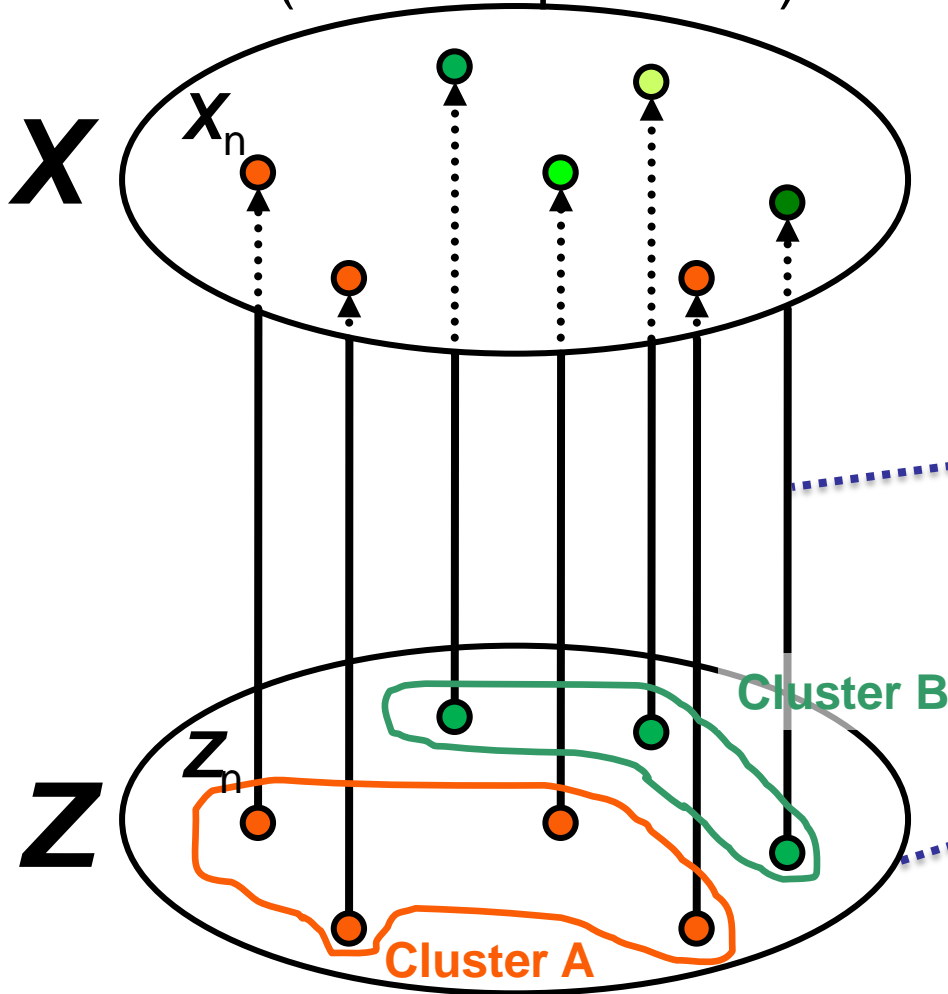
- Hidden modular random field using network modularity

Clustering algorithm

- EM algorithm + ICM

Existing Probabilistic Model (k-means)

Observable
(Gene expression)



Probability distribution
 $p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})$

$$p(\mathbf{X}|\mathbf{Z}) = \prod_{n=1}^N p(X_n|Z_n)$$

$$p(\mathbf{Z}) = \prod_{n=1}^N p(Z_n)$$

Independent

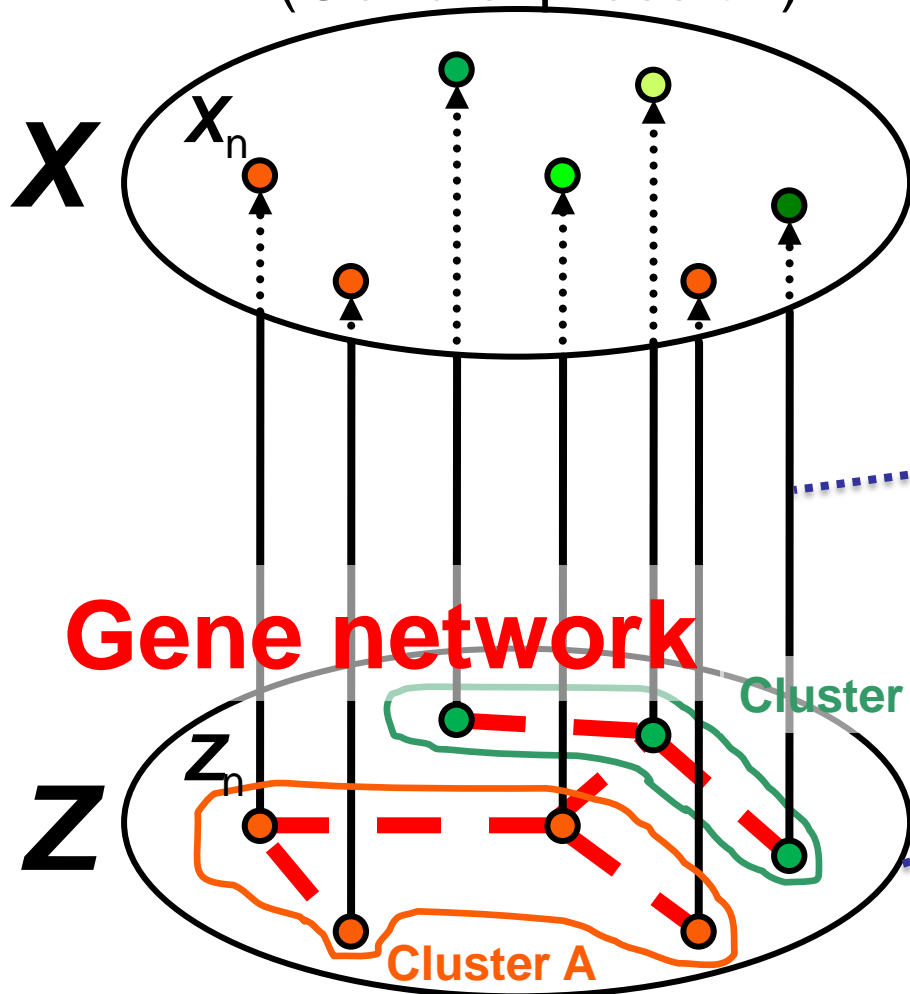
Hidden (Cluster labels)

Our Probabilistic Model

Observable
(Gene expression)

Hidden random field

Probability distribution
 $p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})$



$$p(\mathbf{X}|\mathbf{Z}) = \prod_{n=1}^N p(x_n|z_n)$$

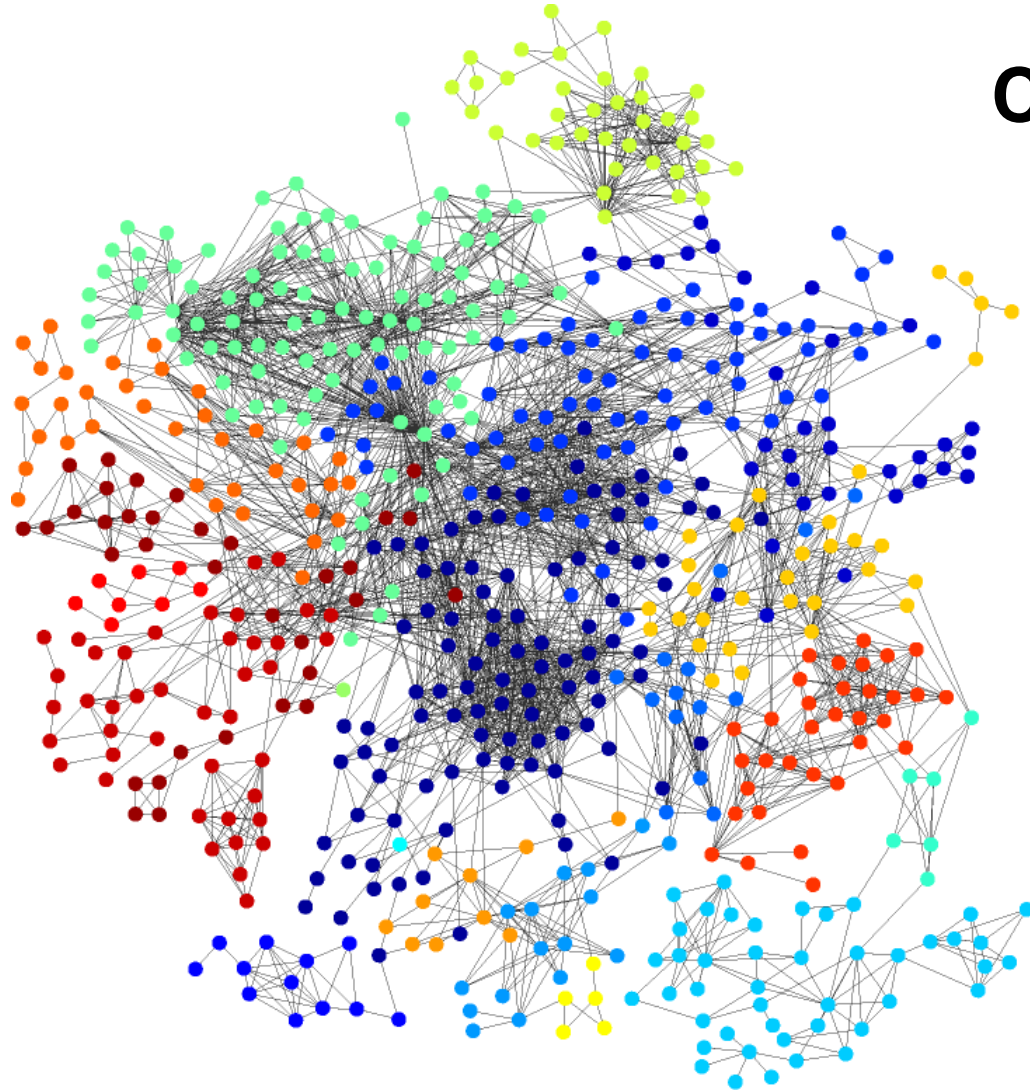
$$p(\mathbf{Z}) =$$

*Random
Field*

dependent

Hidden (Cluster labels)

Gene Networks



Complex networks

Ex. Gene network, WWW,
Social network, ..., etc.

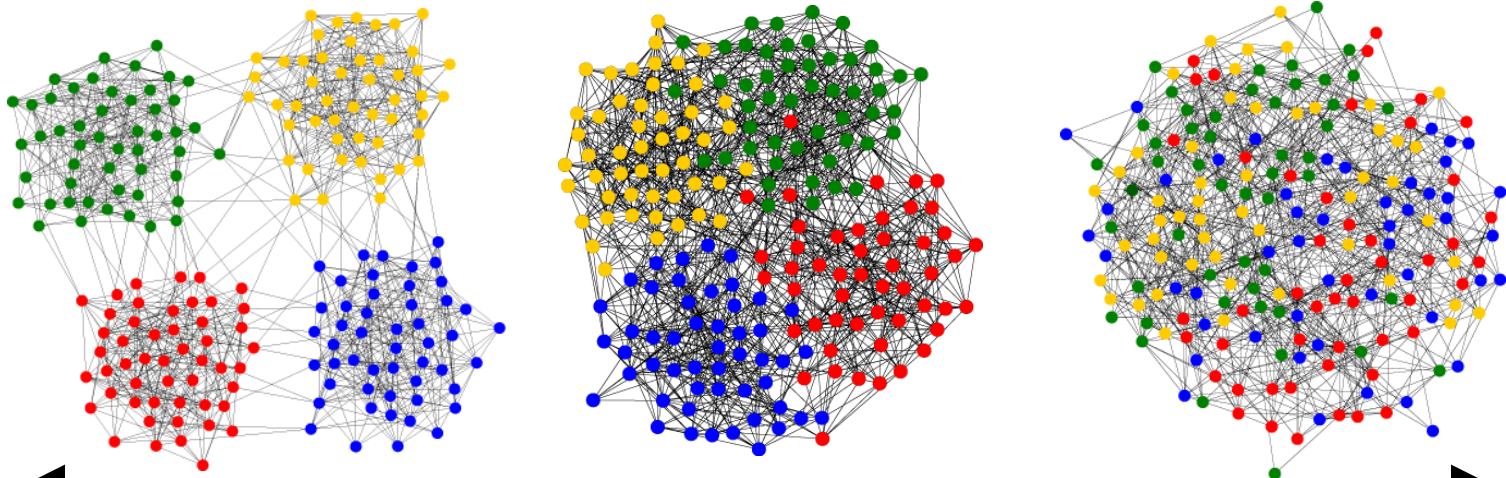
Properties

- Small world phenomena
- Scale-free
- Hierarchical structure
- Network modularity

Ravasz, et al., Science, 2002.
Guimera, et al., Nature, 2005.

Network Modularity

= density of intra-cluster edges



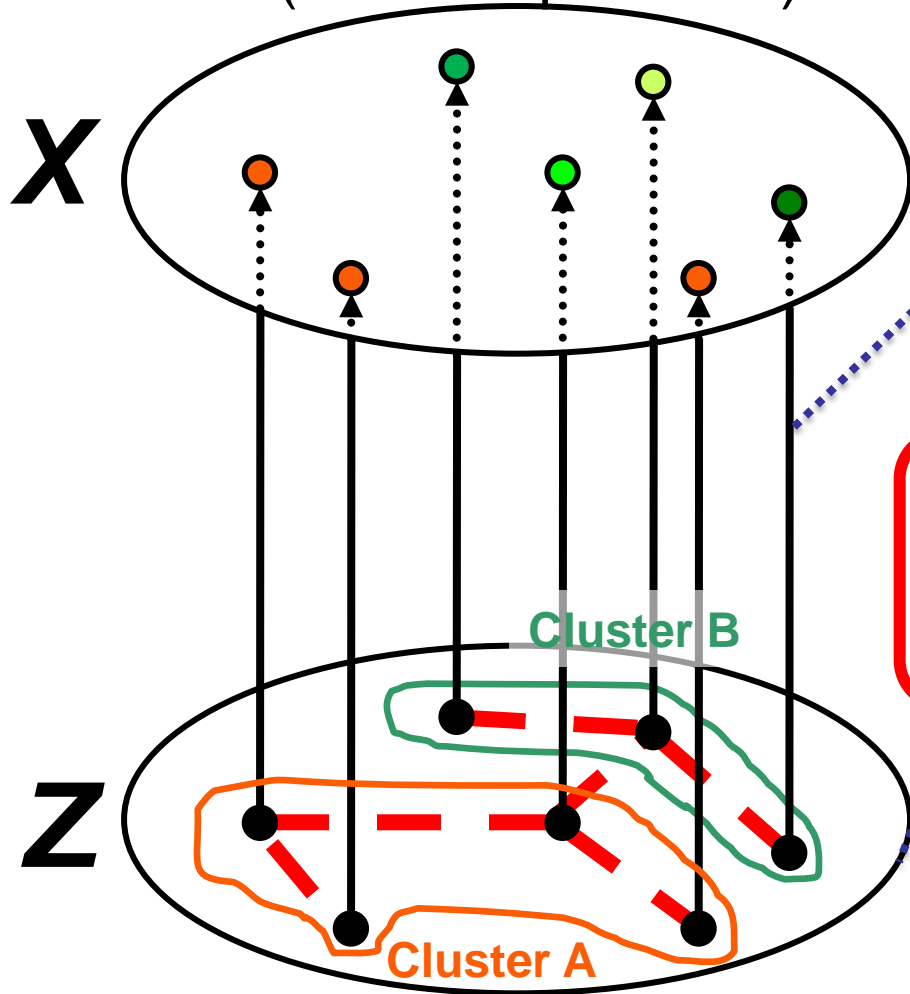
High ←————→ Low

$$V_G(\mathbf{Z}) = N \sum_{k=1}^K \left\{ \frac{l_k(\mathbf{Z})}{L} - \left(\frac{d_k(\mathbf{Z})}{2L} \right)^2 \right\}$$

intra-edges # total edges

Hidden Modular Random Field

Observable
(Gene expression)



Hidden (Cluster labels)

Probability distribution

$$p(\mathbf{X}, \mathbf{Z}) = p(\mathbf{X}|\mathbf{Z})p(\mathbf{Z})$$

$$p(\mathbf{X}|\mathbf{Z}) = \prod_{n=1}^N p(\mathbf{X}_n|Z_n)$$

Modular random field

$$p(\mathbf{Z}) = \text{Random field}$$

Network modularity

$$V_G(\mathbf{Z}) = N \sum_{k=1}^K \left\{ \frac{l_k(\mathbf{Z})}{L} - \left(\frac{d_k(\mathbf{Z})}{2L} \right)^2 \right\}$$

2. Proposed Method

Probabilistic model

- Hidden modular random field using network modularity

Clustering algorithm

- EM algorithm + ICM

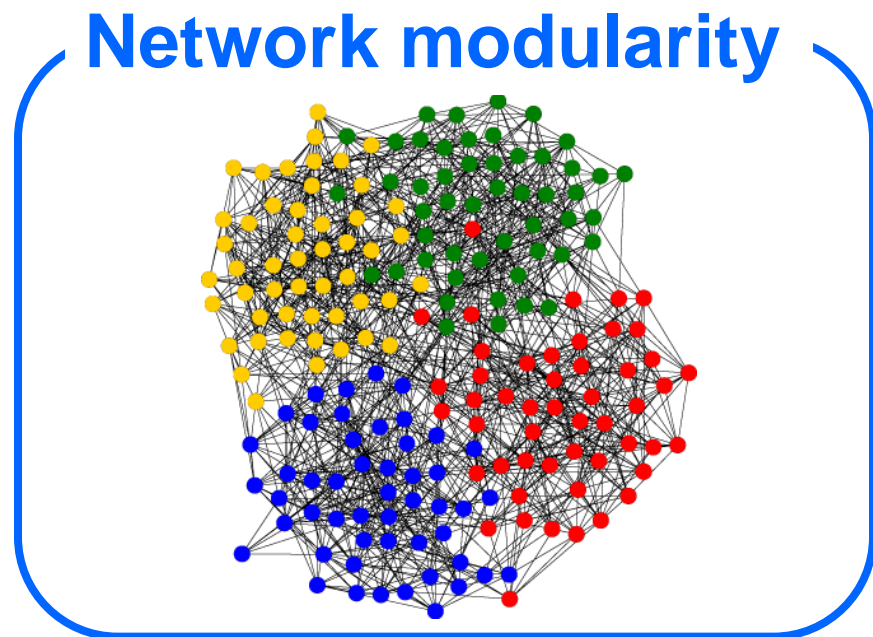
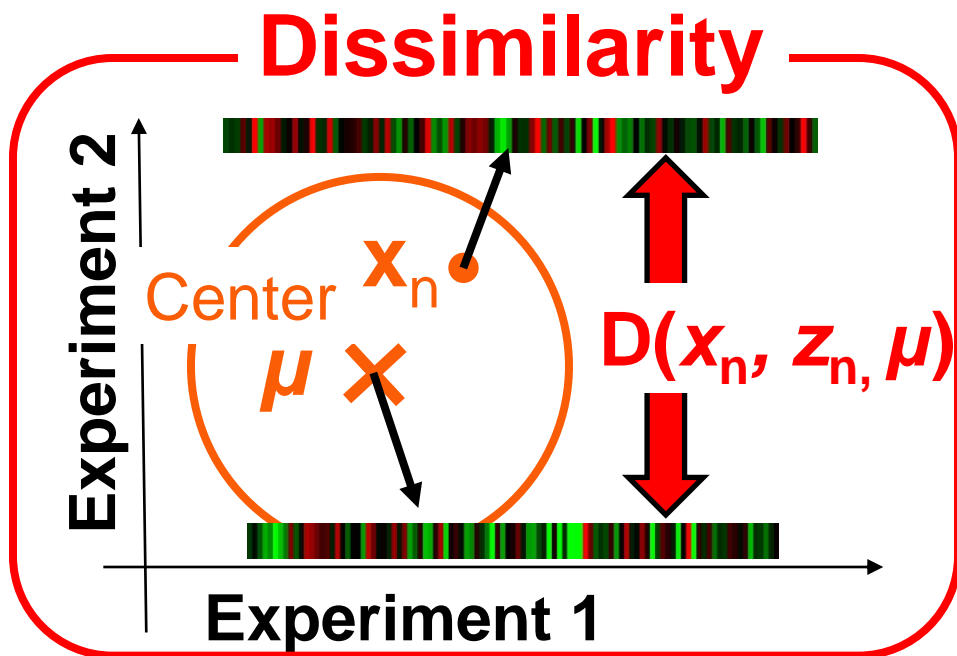
Cost function

is the negative pseudo-likelihood

$$-\log P(\mathbf{X}, \mathbf{Z}) = - \sum_{n=1}^N \log p(\mathbf{x}_n | z_n) - \log P(\mathbf{Z})$$

substitute **our probabilistic model** into $\log P(\mathbf{X}, \mathbf{Z})$

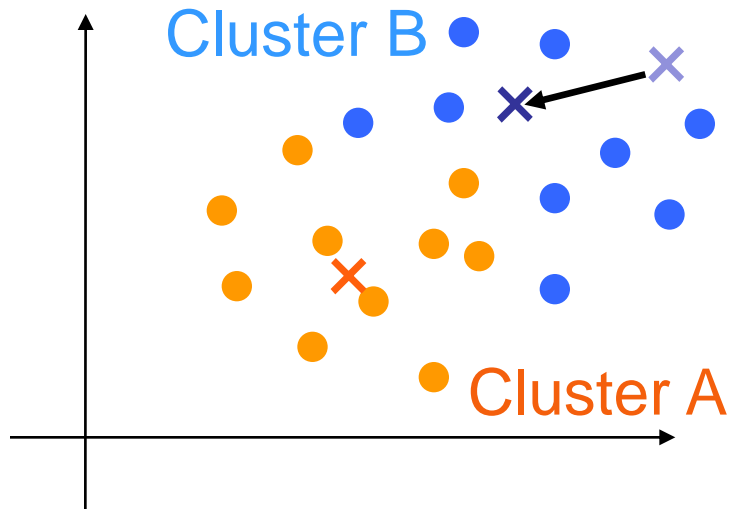
$$J = (1 - \omega) D(\mathbf{X}, \mathbf{Z}, \mu) - \omega V_G(\mathbf{Z})$$



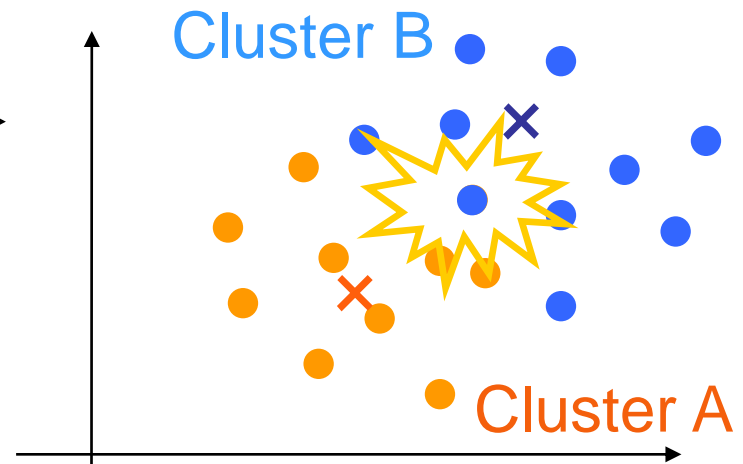
Clustering Algorithm

1. Set the number of cluster K
2. Initialize cluster labels Z and the cluster centers μ .
3. Iterate (i) and (ii) alternately until convergence.

(i) Update the cluster centers
by minimizing the cost J



(ii) Update cluster labels
by minimizing J



Iteration

Iterated Conditional Modes

(J. Besag, J. Roy. Stat. Soc. B, 1986) 13

3. Experiments

- Datasets
- Standard cluster
- Evaluation measure
- Results

Evaluation Measure

Normalized Mutual Information (NMI)

between estimated cluster and the standard cluster

$$NMI = \frac{H(C) + H(G) - H(C, G)}{\sqrt{H(C)}\sqrt{H(G)}}$$

$H(\mathbf{C})$: Entropy of probability variable \mathbf{C} ,

\mathbf{C} : Estimated clusters, \mathbf{G} : Standard clusters

The more **similar** clusters

\mathbf{C} and \mathbf{G} are, **the larger the NMI.**

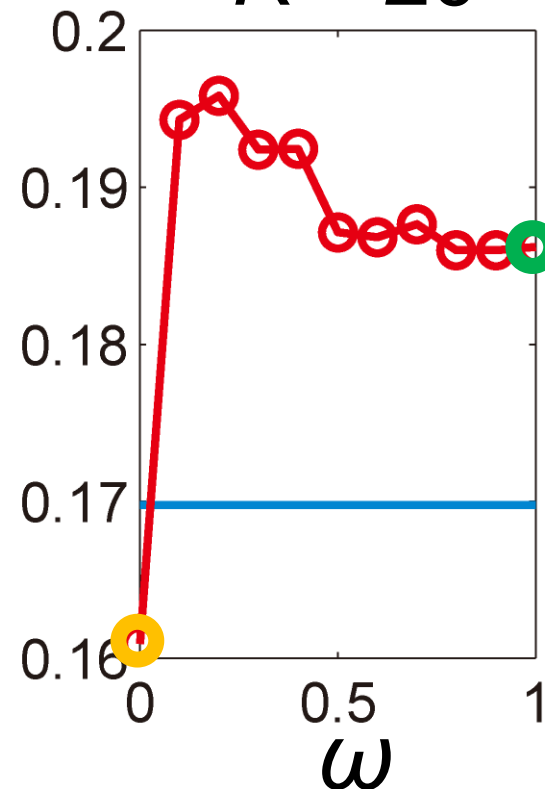
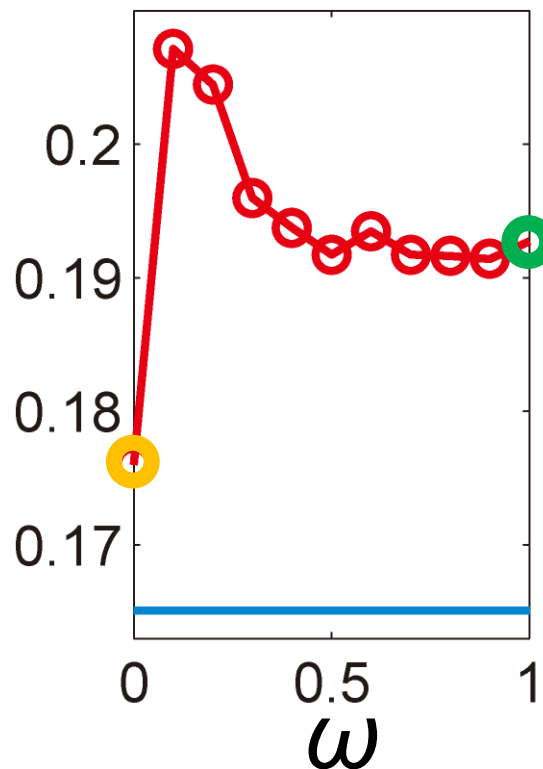
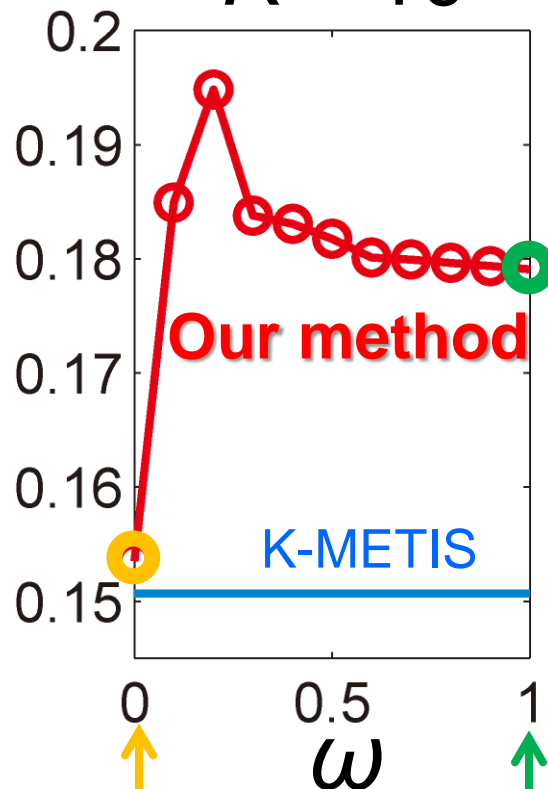
Performance results

$K = 10$

$K = 15$

$K = 20$

NMI



Gene expression only
(k-means)

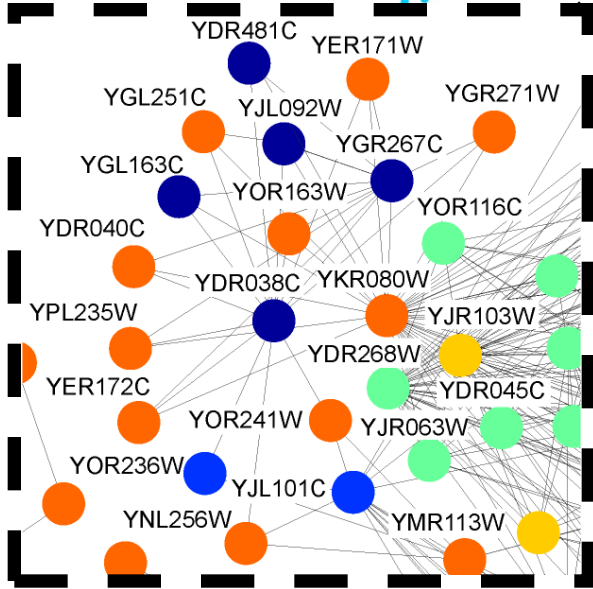
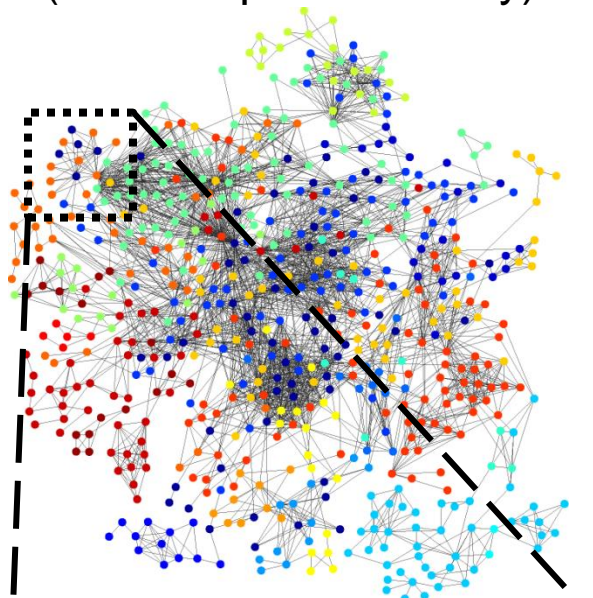
Gene network only
(maximum modularity)

*** Best NMI is in $0 < \omega < 1$.**
*** Better than three others.**
(i) k-means, (ii) max. Mod., (iii) k-METIS

Examples of Resultant Cluster (#cluster = 20)

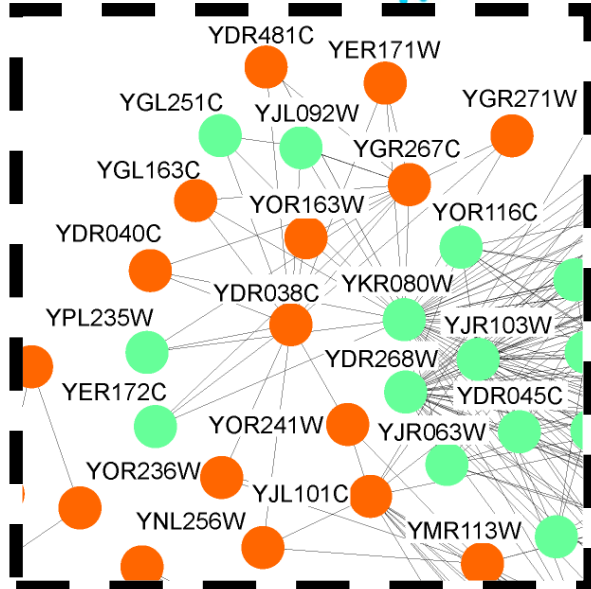
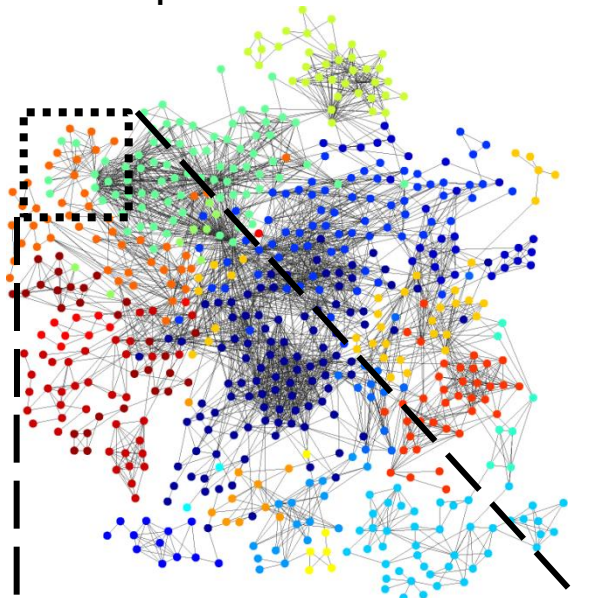
$\omega = 0$

(Gene expression only)



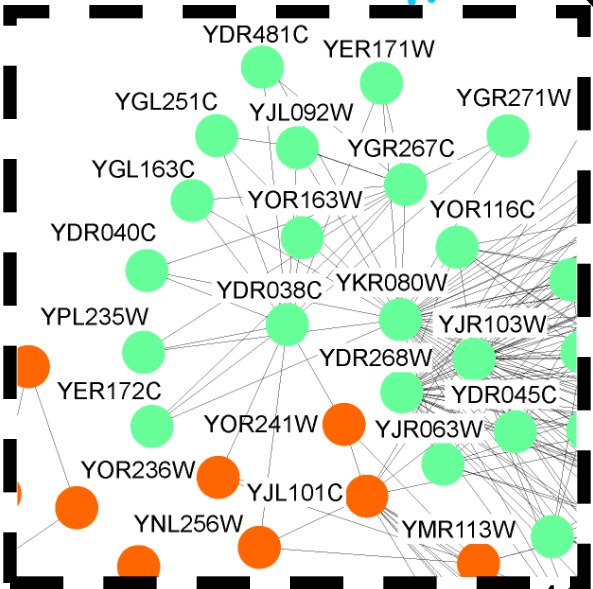
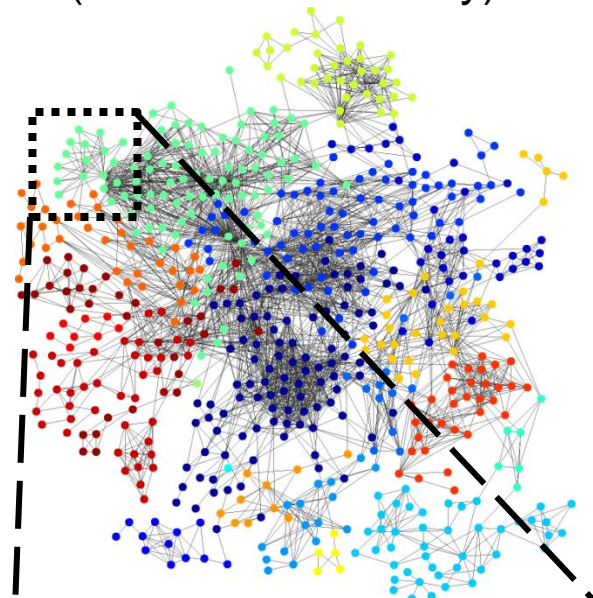
$\omega = 0.2$

(Both expression and network)



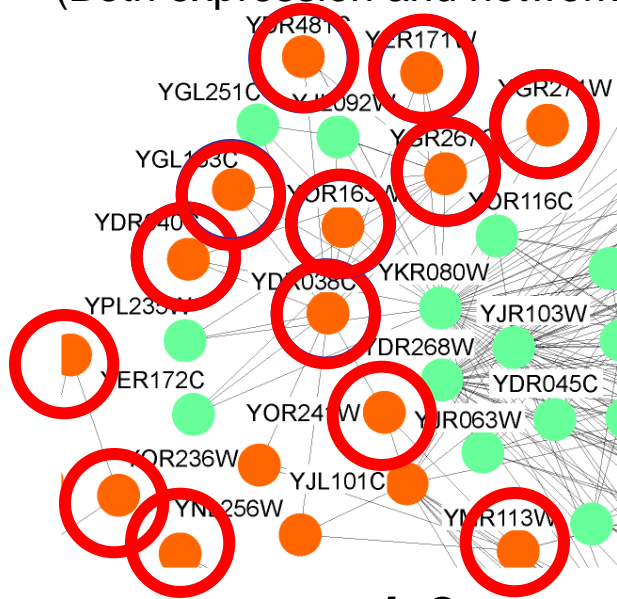
$\omega = 1$

(Gene network only)



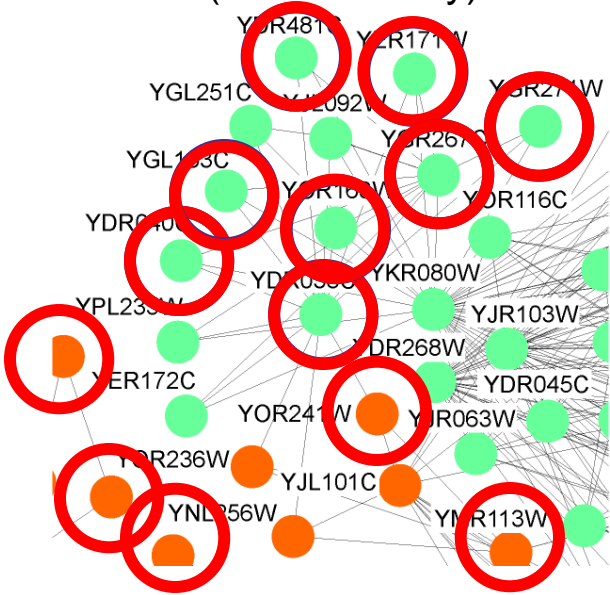
$\omega = 0.2$

(Both expression and network)

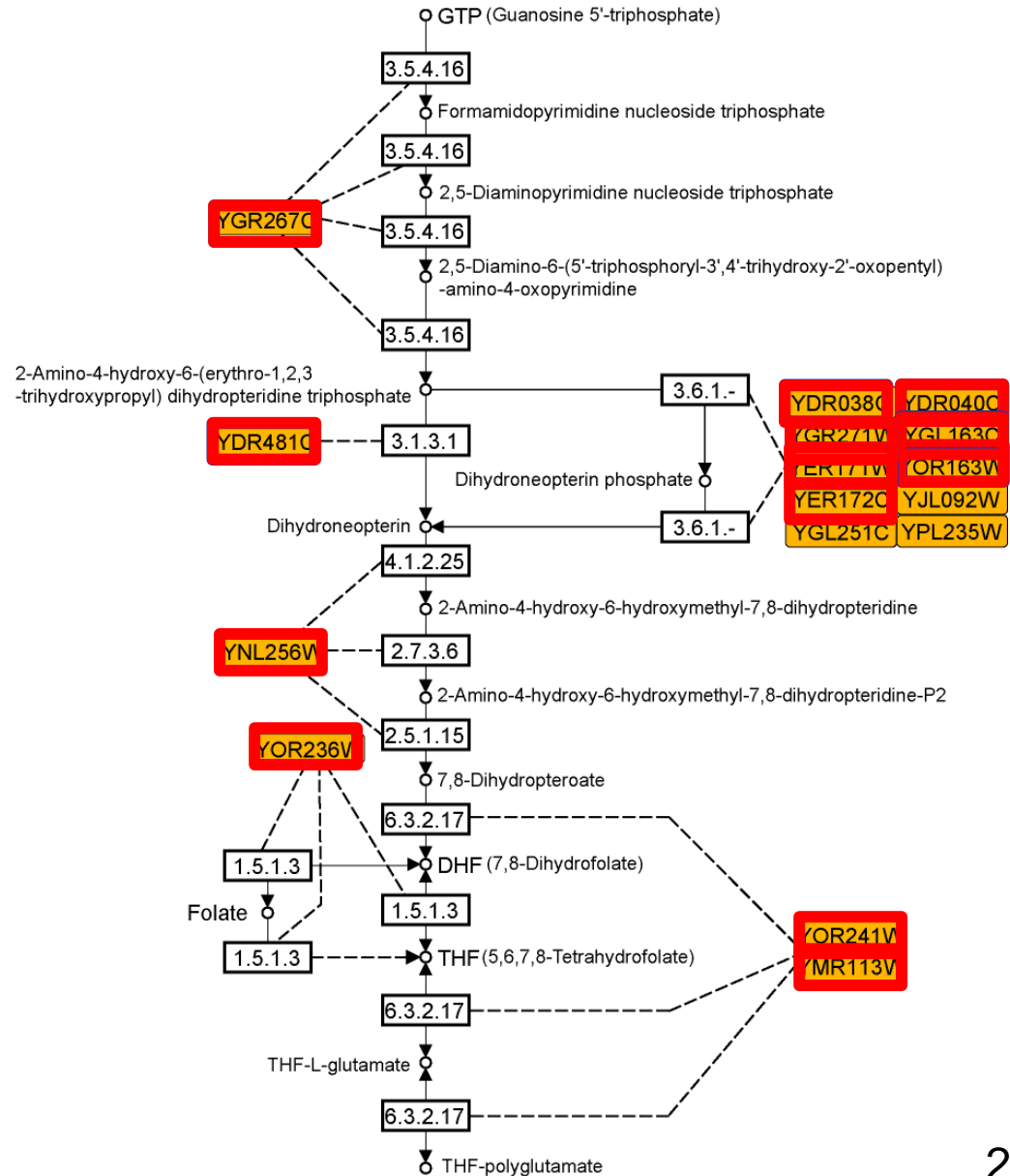


$\omega = 1.0$

(Network only)



Folate biosynthesis pathway



Summary

1. Proposed a new clustering method combining expression data and a network

- Probabilistic model
Hidden modular random fields (gene expression + network)
- Clustering algorithm
EM algorithm + ICM

2. Experimental results

- Best **NMI** by $0 < \omega < 1$ (**expression + network**).
- Better than three others.
 - $\omega=0$ (k-means)
 - $\omega=1$ (Maximizing network modularity)
 - k-METIS (Graph partitioning)

Thank you for your attention!