UNIVERSITY OF TORONTO

SickKids®

# Patient networks in cancer: a platform for data integration
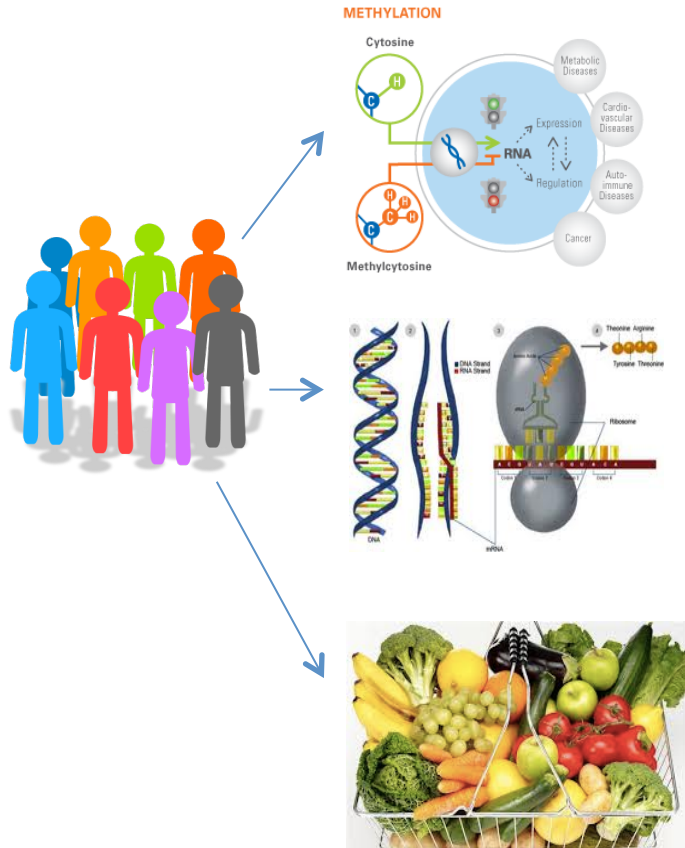
Anna Goldenberg
and
The Goldenberg Lab

# Outline

o Data integration – problem setup

o Patient network representation – why and how

o Similarity Network Fusion – novel integration method

o Network driven analysis:

  o Cancer heterogeneity

  o Differential feature selection

o Missing data

  o Random entries

  o Patients

o Taking networks further:

  o Survival analysis (novel formulation)
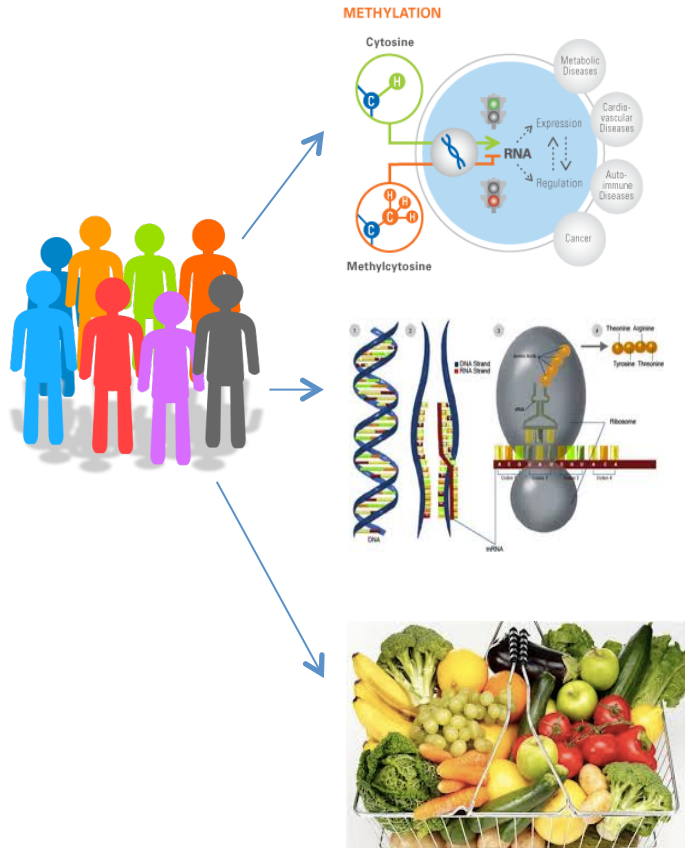
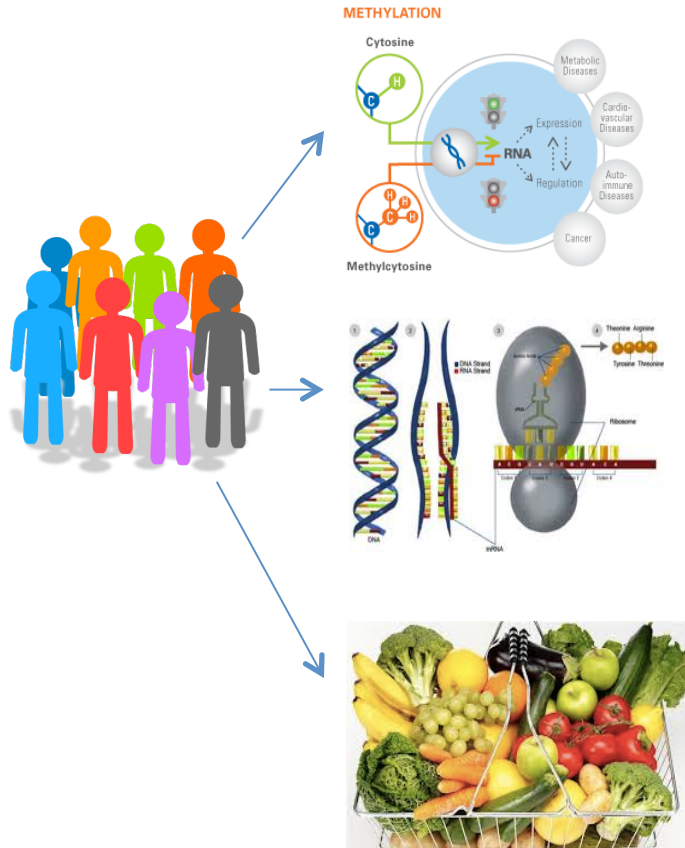    o Personalized medicine

# Problem setup

# Problem setup

# Problem setup

# Problem setup

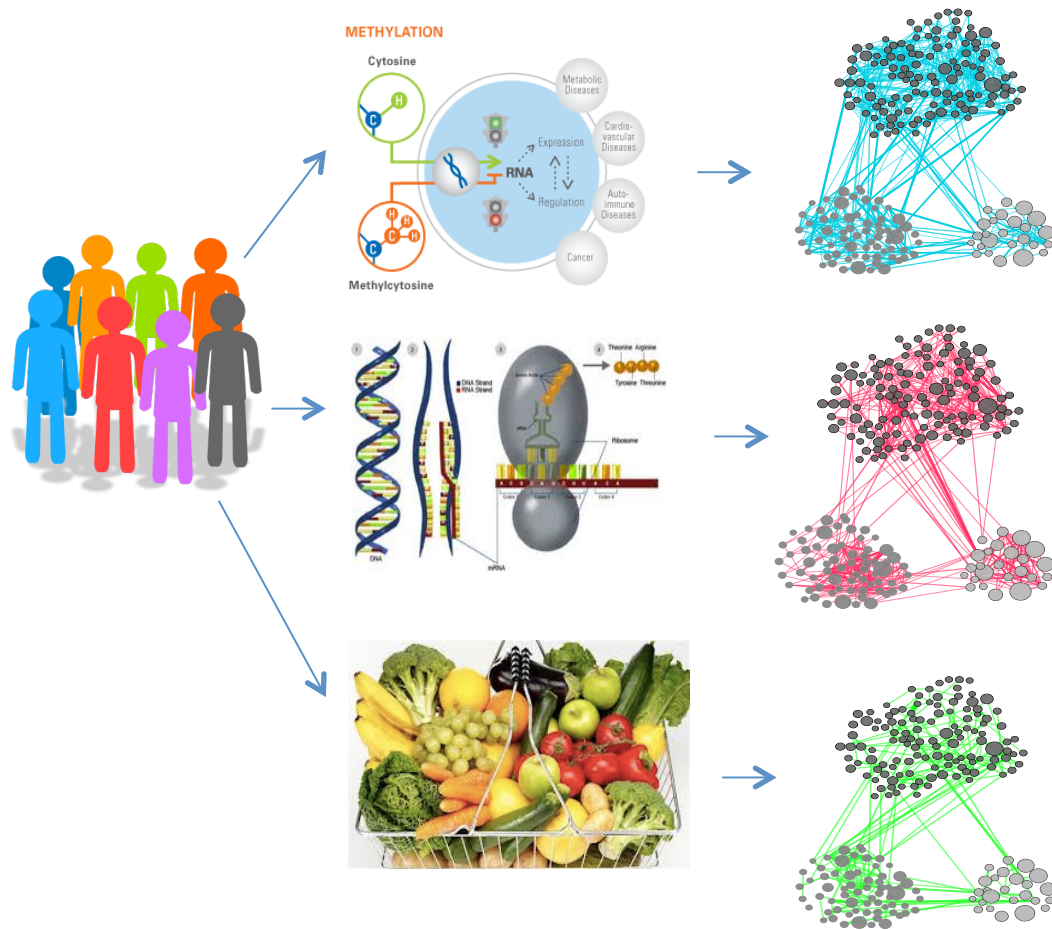## How to combine?



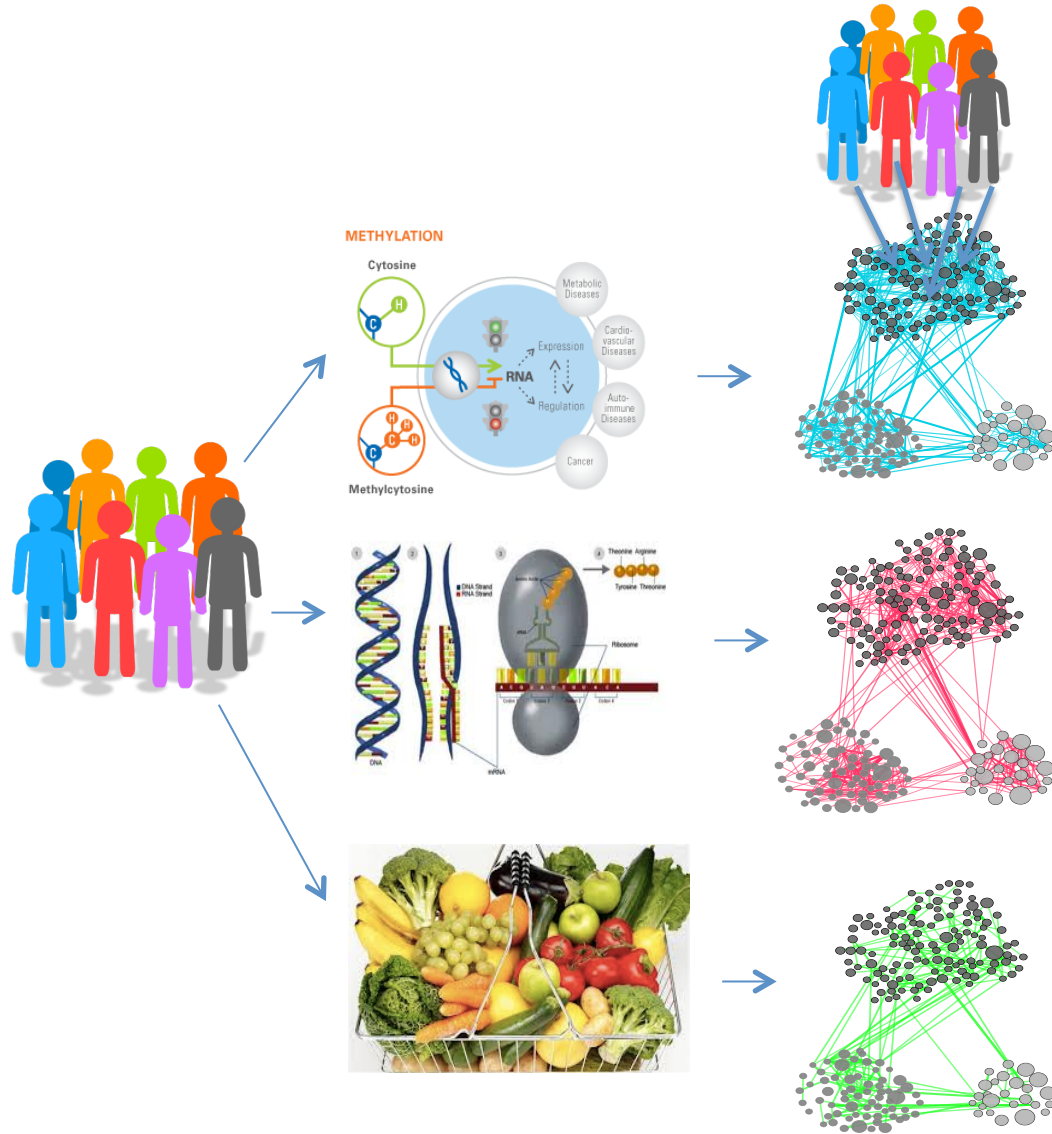## Issues

⊙ Large number of measurements, small sample sizes (p>>n)

⊙ Need to integrate common and complementary information

⊙ Not all measurements can be mapped to the same unit (gene)

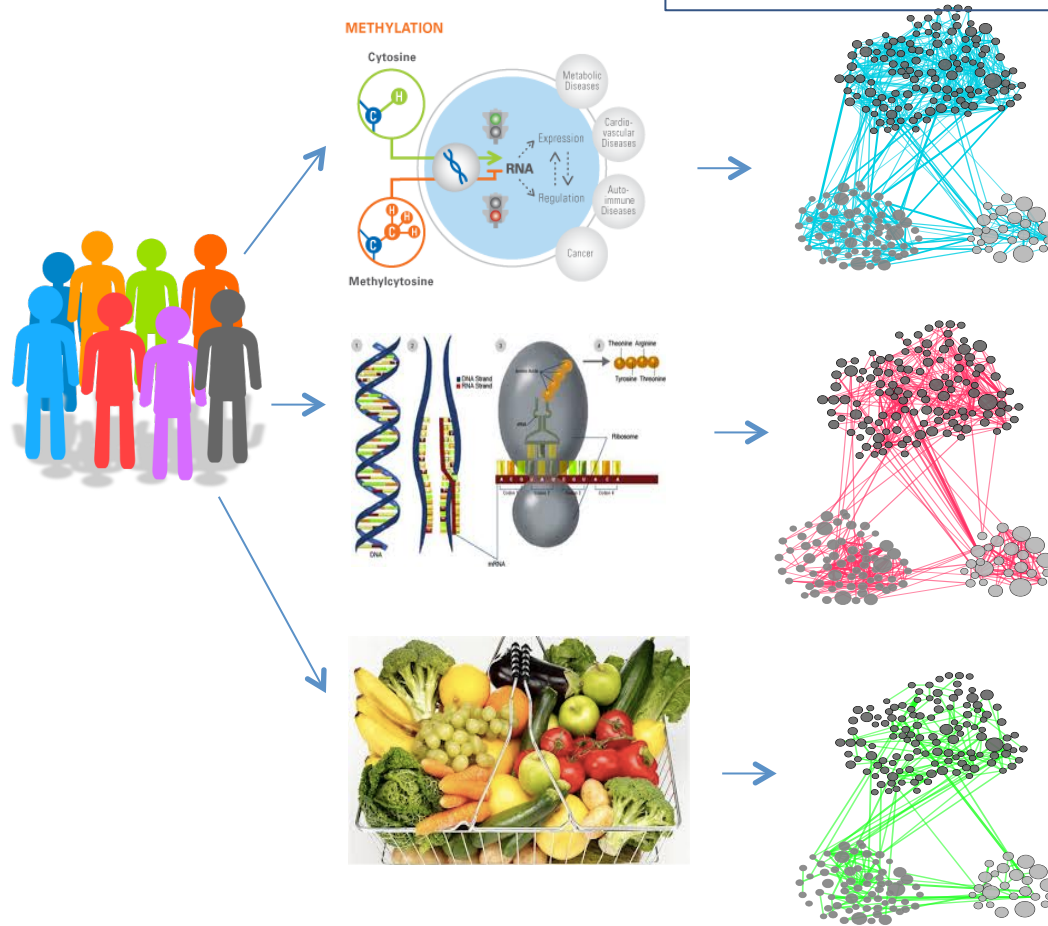# Goldenberg Lab: Similarity Network Fusion

# Goldenberg Lab: Similarity Network Fusion

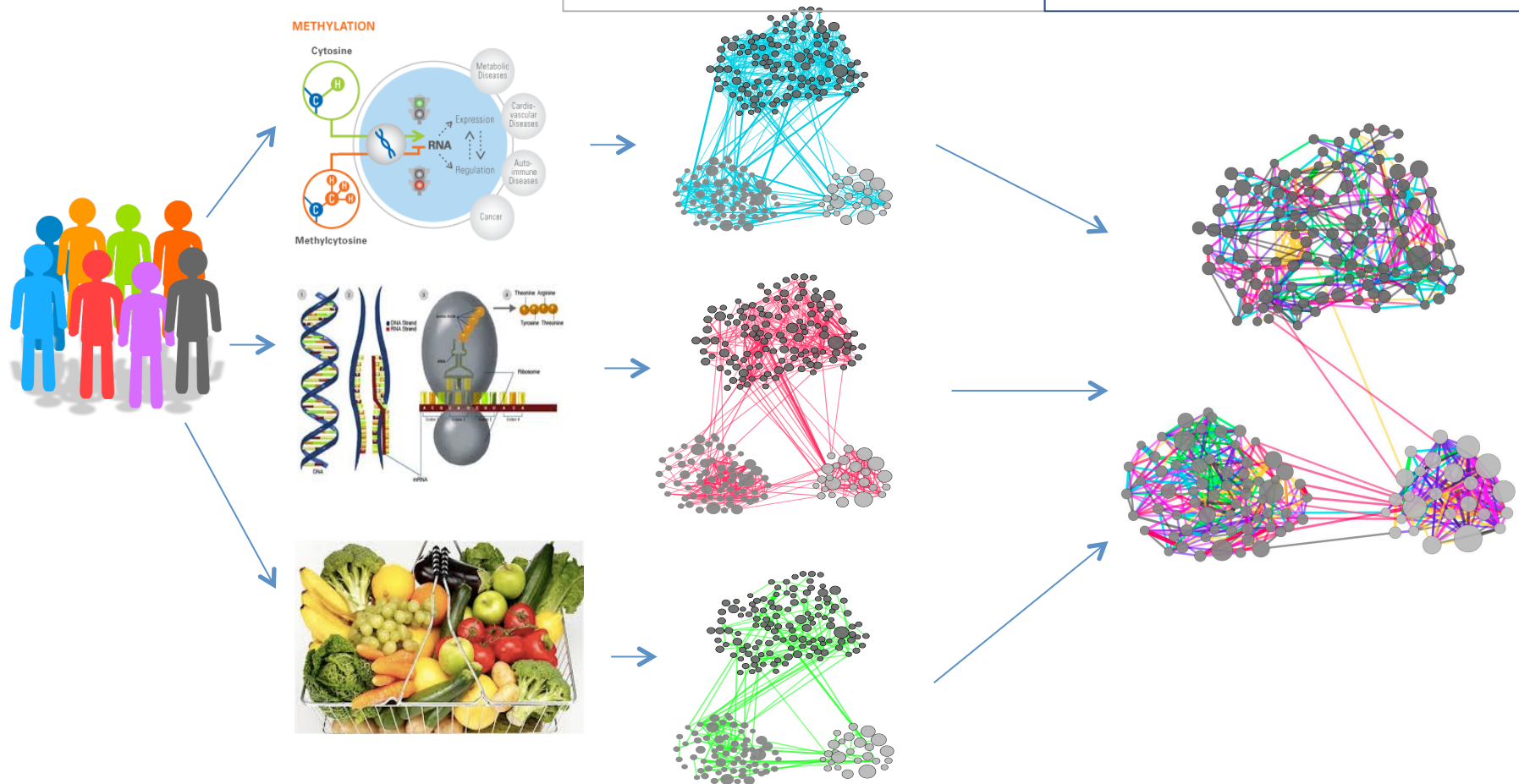# Goldenberg Lab: Similarity Network Fusion

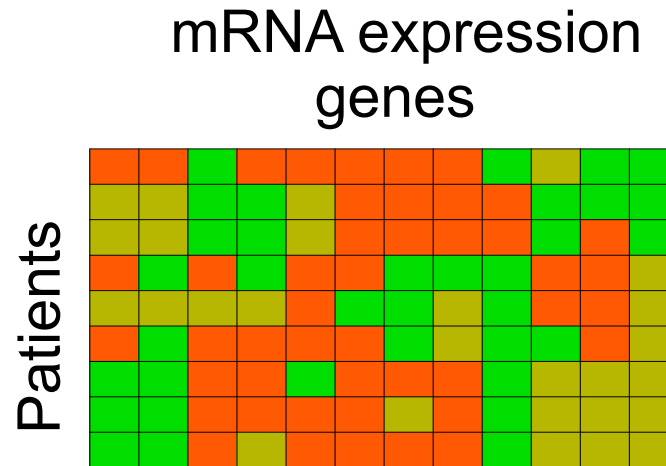**Step 1**. Construct a similarity network for each data source

# Goldenberg Lab: Similarity Network Fusion

# 1. Construct similarity networks

mRNA expression
genes
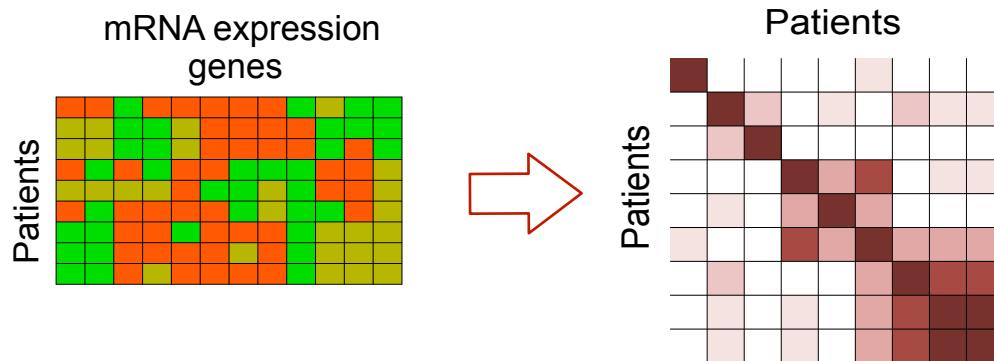
Patients

# 1. Construct similarity networks

Patient similarity:
$$W(i,j) = exp(\frac{\rho(x_i, x_j)^2}{\eta \xi_{ij}^2})$$

Adjacency matrix:
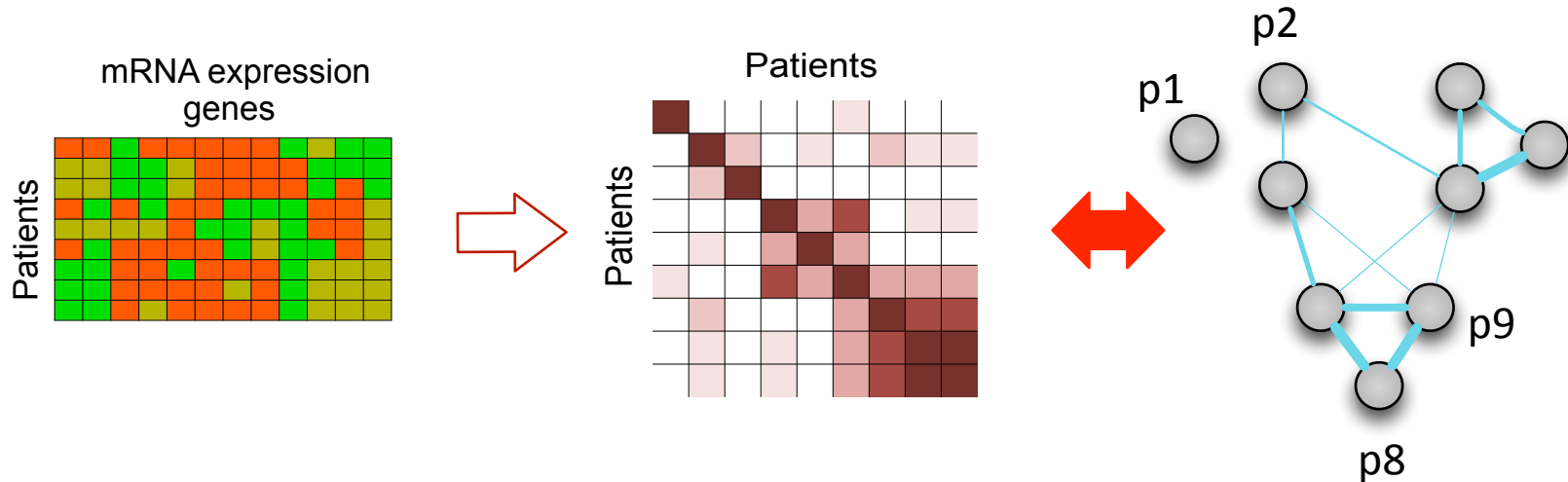$$P(i,j) = \frac{W(i,j)}{\sum_{k \in V} W(i,k)}$$



mRNA expression
genes

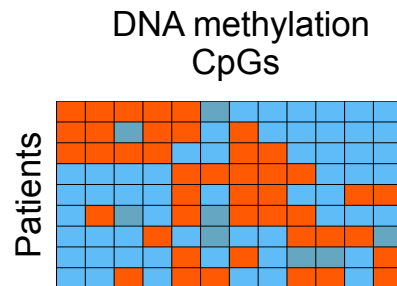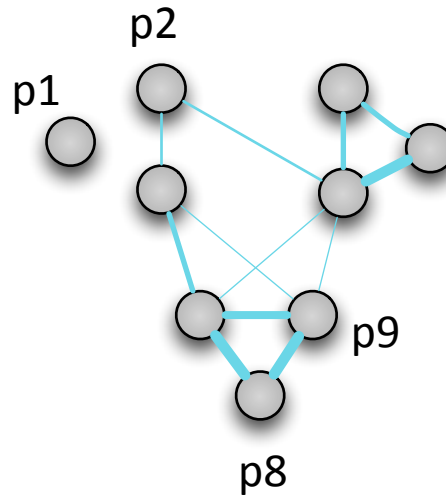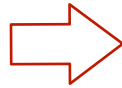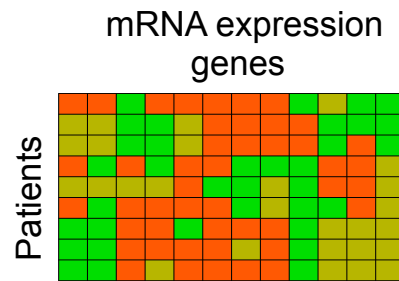Patients

Patients

Patients

# 1. Construct similarity networks

1) $$\mathcal{W}(i,j) = \begin{cases} W(i,j) \text{ if } x_j \in KNN(x_i) \\ 0 \text{ otherwise} \end{cases}$$

Sparsification

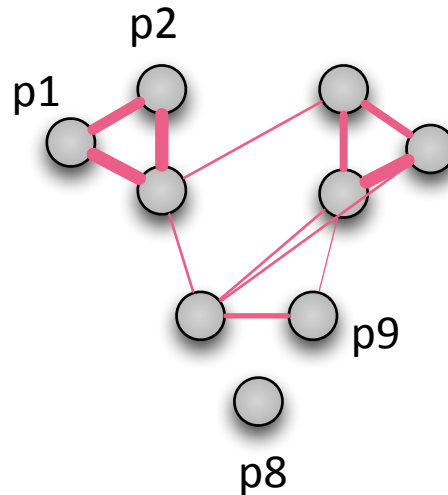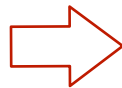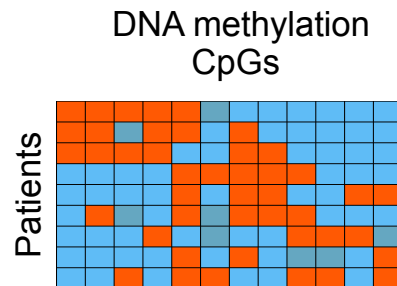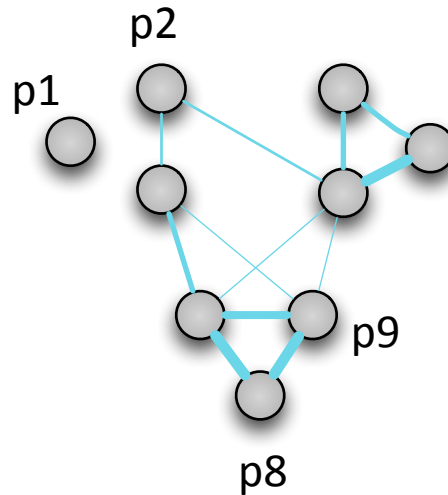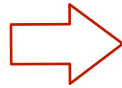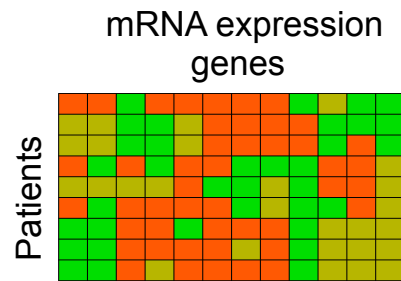2) $$\mathcal{P}(i,j) = \frac{\mathcal{W}(i,j)}{\sum_{x_k \in KNN(x_i)} \mathcal{W}(i,k)}$$

# 1. Construct similarity networks

mRNA expression genes

Patients

DNA methylation CpGs

Patients

p1
p2
p9
p8

# 1. Construct similarity networks



mRNA expression genes

Patients

p1  p2  p9  p8

DNA methylation CpGs

Patients

p1  p2  p9  p8

# 2. Combine networks

Sample Similarity Networks



| | | | | |
|---|---|---|---|---|
| ⬤ Patient | Patient similarity: | ━ mRNA-based | ━ DNA Methylation-based | ━ Supported by all data |

# Combine networks

Sample Similarity Networks

Fusion



$$P_{t+1}^{(1)} = \mathcal{P}^{(1)} \times (P_t^{(2)}) \times (\mathcal{P}^{(1)})'$$

$$P_{t+1}^{(2)} = \mathcal{P}^{(2)} \times (P_t^{(1)}) \times (\mathcal{P}^{(2)})'$$

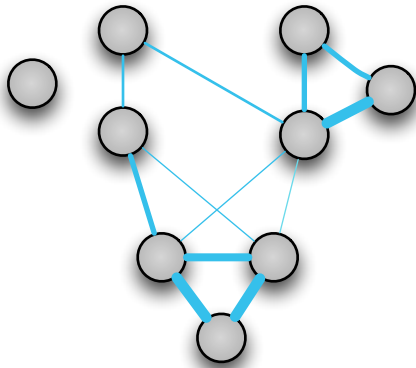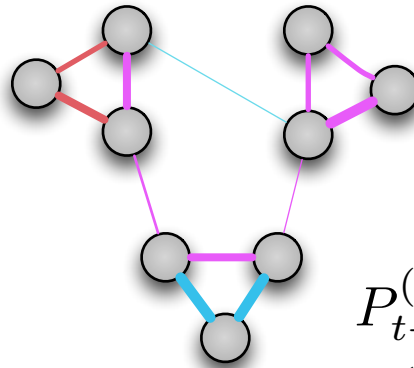○ Patient    Patient similarity:    —— mRNA-based    —— DNA Methylation-based    —— Supported by all data

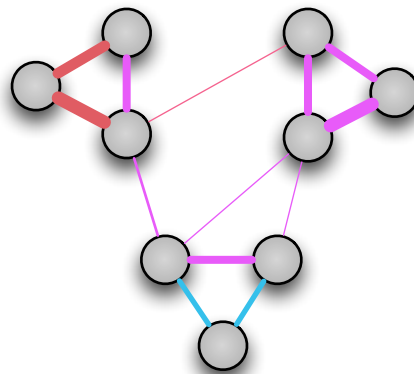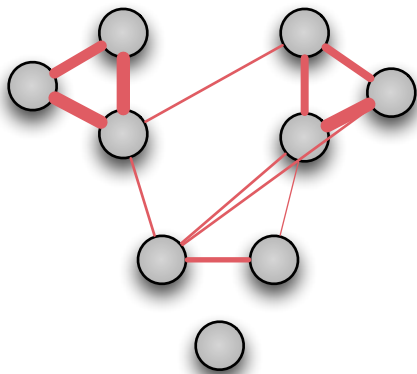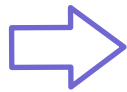# 2. Combine networks

Sample Similarity Networks

Fusion Iterations



$$P_{t+1}^{(1)} = \mathcal{P}^{(1)} \times (P_t^{(2)}) \times (\mathcal{P}^{(1)})'$$

$$P_{t+1}^{(2)} = \mathcal{P}^{(2)} \times (P_t^{(1)}) \times (\mathcal{P}^{(2)})'$$

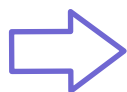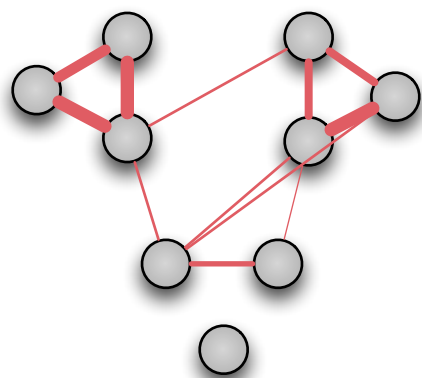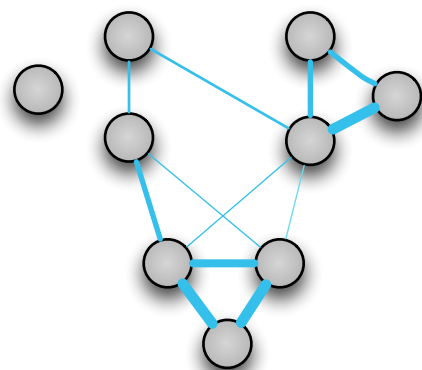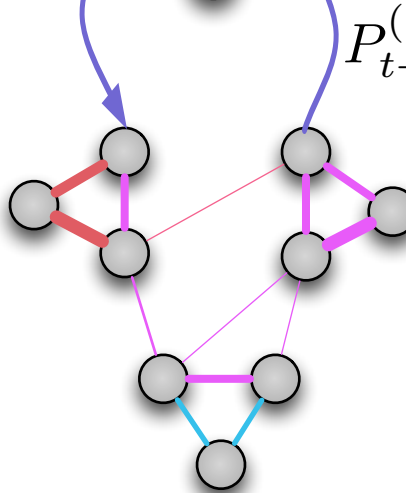Patient    Patient similarity:    —— mRNA-based    —— DNA Methylation-based    —— Supported by all data

# 2. Combine networks

Sample Similarity Networks

Fusion Iterations

Fused Similarity Network

$$\frac{\|W_{t+1} - W_t\|}{\|W_t\|} \leq 10^{-6}$$

Patient    Patient similarity:    —— mRNA-based    —— DNA Methylation-based    —— Supported by all data

# Network Fusion

Fusing 2 networks:

$$P_{t+1}^{(1)} = \mathcal{P}^{(1)} \times (P_t^{(2)}) \times (\mathcal{P}^{(1)})'$$

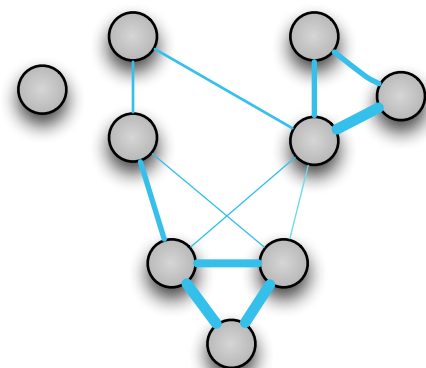$$P_{t+1}^{(2)} = \mathcal{P}^{(2)} \times (P_t^{(1)}) \times (\mathcal{P}^{(2)})'$$

Fusing m networks:

$$P_{t+1}^{(i)} = \mathcal{P}^{(i)} \times (\frac{1}{m-1} \sum_{j \neq i} P_t^{(j)}) \times (\mathcal{P}^{(i)})' + \eta I$$

# Experiments

**Data:**

    5 TCGA cancers

    METABRIC (Large

      Breast Cancer db)

**Comparative Methods:**

    Concatenation

    iCluster

    PDSB

    Multiple kernel learning

**Criteria:**

$-\log_{10}$(log rank pvalue)

Silhouette score  (cluster homogeneity)

Running time

# Simulation 1 – complementarity

# Simulation 1 convergence



20% of patients are mislabeled

# Simulation 2 - removing noise

# Simulation 2 - removing noise

# TCGA Data

| Cancer Type | Patients | mRNA | Methylation | miRNA | Controls | |
|---|---|---|---|---|---|---|
| | | | | | mRNA | Methylation |
| GBM | 215 | 12,042 | 1,491 | 534 | 10 | - |
| BIC | 105 | 17,814 | 23,094 | 1,046 | 63 | 27 |
| KRCCC | 124 | 20,532 | 24,976 | 1,046 | 68 | 199 |
| LSCC | 105 | 12,042 | 27,578 | 1,046 | - | 27 |
| COAD | 92 | 17814 | 27578 | 705 | 19 | 37 |

# Case study: Glioblastoma

DNA methylation data



mRNA expression



miRNA expression



Bo Wang

# Case study: Glioblastoma

DNA methylation data

FUSED

mRNA expression

miRNA expression

Similarity type

miRNA — DNA methylation

mRNA

Bo Wang

# Clinical properties of the subtypes



p-value = $2\times10^{-4}$

# Clinical properties of the subtypes



p-value = 2x10$^{-4}$

p-value = 3x10$^{-5}$

# Clinical properties of the subtypes

# Biological characterization of the subtypes

# Feature Selection

**Standard t-test**
**Differential analysis**



Bo Wang

# Feature Selection



**Standard t-test Differential analysis**

**Network-based NMI Differential analysis**

Bo Wang

# Gene Pre-selection in GBM



Genes are ordered by significance of the differential values between tumor and normal

Bo Wang

# Gene Pre-selection in GBM



Genes are ordered by significance of the differential values between tumor and normal

Bo Wang

# Gene Pre-selection in GBM



Genes are ordered by significance of the differential values between tumor and normal

Bo Wang

# Gene pre-selection across cancers



Bo Wang

# Clustering of the network



Bo Wang

# Patient networks framework advantages

- ✓ Creates a unified view of patients based on multiple heterogeneous sources

- ✓ Integrates gene and non-gene based data

- ✓ No need to do gene pre-selection

- ✓ Robust to different types of noise

- ✓ Scalable

# Patient networks

- ✓ Obtain superior results on regular tasks such as subtyping
- ✓ No fea~~ture~~ pre-se~~lection~~
- ✓ Imputation
- ✓ Can ea~~sily~~ ~~without~~ imput~~ing~~

Transformative power
of patient networks

# Breast Cancer (METABRIC example)

CNV and expression data
Discovery: 997 patients
Validation: 995 patients

Nature, 2012

| | PAM50 (5 clusters) | iCluster (10 clusters) | SNF (5 clusters) | SNF (10 clusters) |
|---|---|---|---|---|
| $P$ value discovery cohort | $3.0 \times 10^{-9}$ | $1.2 \times 10^{-14}$ | $6.10 \times 10^{-11}$ | $3.31 \times 10^{-12}$ |
| $P$ value validation cohort | $1.7 \times 10^{-9}$ | $2.9 \times 10^{-11}$ | $5.12 \times 10^{-13}$ | $7.86 \times 10^{-12}$ |
| CI discovery cohort | 0.560 | 0.621 | 0.638 | 0.638 |
| CI validation cohort | 0.551 | 0.605 | 0.633 | 0.633 |

established

# Breast Cancer (METABRIC example)

CNV and expression data
Discovery: 997 patients
Validation: 995 patients

Nature, 2012

| | PAM50 (5 clusters) | iCluster (10 clusters) | SNF (5 clusters) | SNF (10 clusters) |
|---|---|---|---|---|
| $P$ value discovery cohort | $3.0 \times 10^{-9}$ | $1.2 \times 10^{-14}$ | $6.10 \times 10^{-11}$ | $3.31 \times 10^{-12}$ |
| $P$ value validation cohort | $1.7 \times 10^{-9}$ | $2.9 \times 10^{-11}$ | $5.12 \times 10^{-13}$ | $7.86 \times 10^{-12}$ |
| CI discovery cohort | 0.560 | 0.621 | 0.638 | 0.638 |
| CI validation cohort | 0.551 | 0.605 | 0.633 | 0.633 |

established

So how many subtypes are there really in breast cancer?

# Predicting using the network

# Predicting using the network

Cox objective

$$lp(z) = \sum_{i=1}^{n} \delta_i \left( \mathbf{X}_i^T z - \log \left( \sum_{j \in R(t_i)} \exp\left( \mathbf{X}_j^T z \right) \right) \right)$$

# Predicting using the network

Cox objective
$$lp(z) = \sum_{i=1}^{n} \delta_i \left( \mathbf{X}_i^T z - \log \left( \sum_{j \in R(t_i)} \exp\left( \mathbf{X}_j^T z \right) \right) \right)$$

Our network-regularized objective

$$lp(\mathbf{z}) = \sum_{i=1}^{n} \delta_i \left( X_i^T \mathbf{z} - \log \left( \sum_{j \in R(t_i)} \exp(X_j^T \mathbf{z}) \right) \right) - \lambda \sum_i \sum_j (X_i^T z - X_j^T z)^2 w_{ij}$$

# Predicting using the network
# Breast Cancer (METABRIC example)

CNV and expression data
Discovery: 997 patients
Validation: 995 patients

Nature, 2012

| | PAM50 (5 clusters) | iCluster (10 clusters) | SNF (5 clusters) | SNF (10 clusters) | Network |
|---|---|---|---|---|---|
| $P$ value discovery cohort | $3.0 \times 10^{-9}$ | $1.2 \times 10^{-14}$ | $6.10 \times 10^{-11}$ | $3.31 \times 10^{-12}$ | – |
| $P$ value validation cohort | $1.7 \times 10^{-9}$ | $2.9 \times 10^{-11}$ | $5.12 \times 10^{-13}$ | $7.86 \times 10^{-12}$ | – |
| CI discovery cohort | 0.560 | 0.621 | 0.638 | 0.638 | 0.720 |
| CI validation cohort | 0.551 | 0.605 | 0.633 | 0.633 | 0.706 |

established