

分子标志物  
揭示表达暗物质  
动态网络标志物

中国科学院系统生物学重点实验室  
陈洛南

分子标志物 → occurred disease

网络标志物 → occurred complex disease

动态网络标志物 → un-occurred disease

# Characterizing one sample

One sample →

- ✓ Sequence
- ✓ Gene expression
- ✓ Protein expression
- ✓ Metabolomics
- ✓ Methylation

Unstable feature  
不稳定

Multiple samples →

✗ **Network**

Stable feature  
稳定

**Clinic data or Single-cell data**

**One sample or a few samples for an individual**

***Large samples for population***

***Small samples for individuals***

**Single-cell data**

Question ?

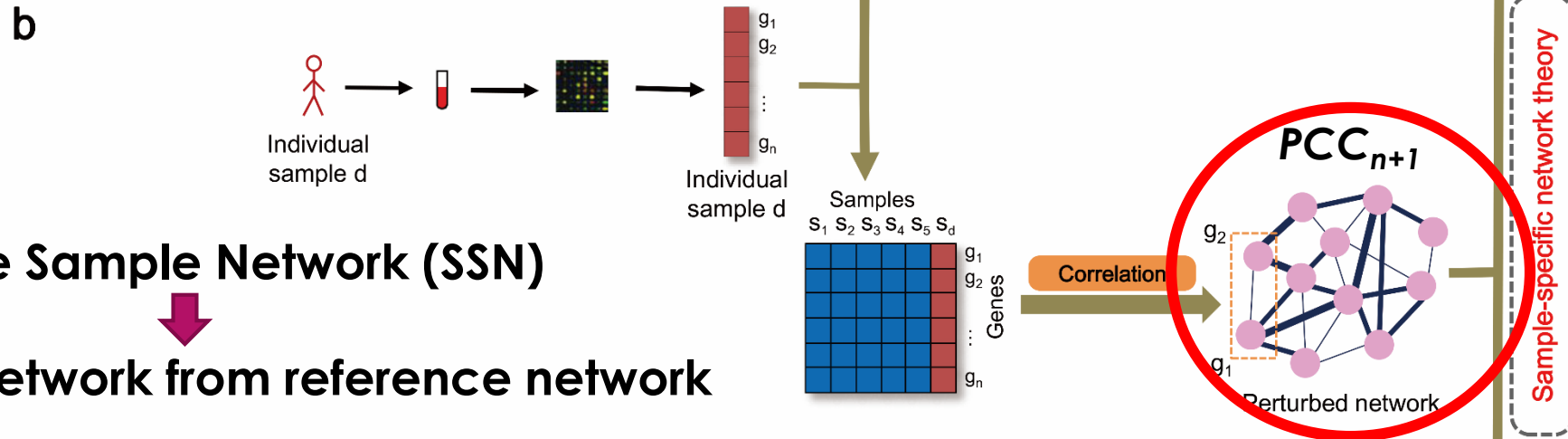
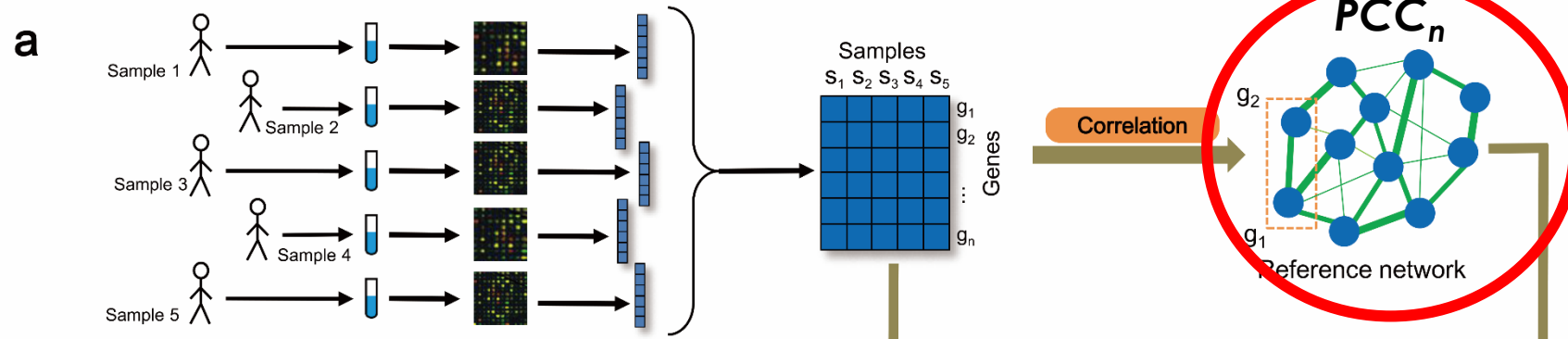
► Can we construct a network by a single sample data?

YES



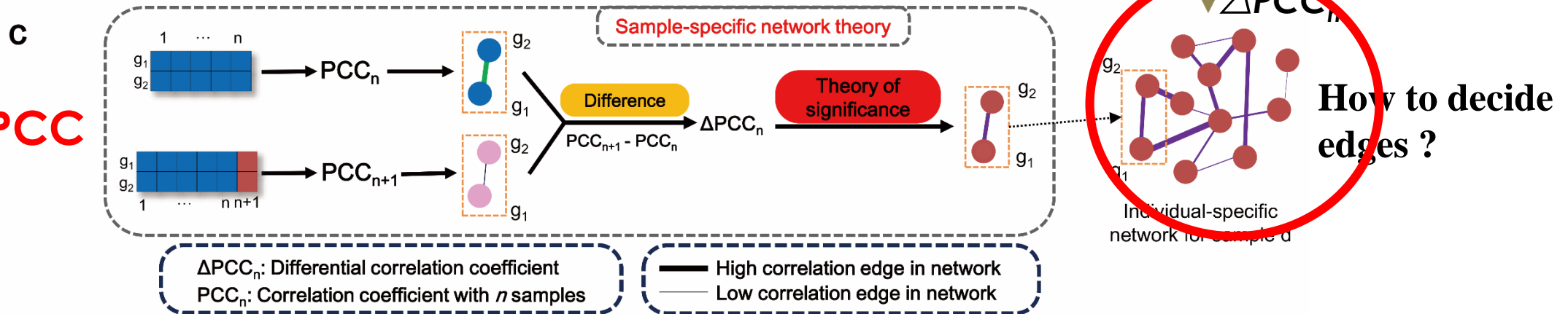
Single Sample Network (SSN)

PCC



Single Sample Network (SSN)

Differential network from reference network



# 单样本网络理论

一个样本  $\rightarrow$  一个网络

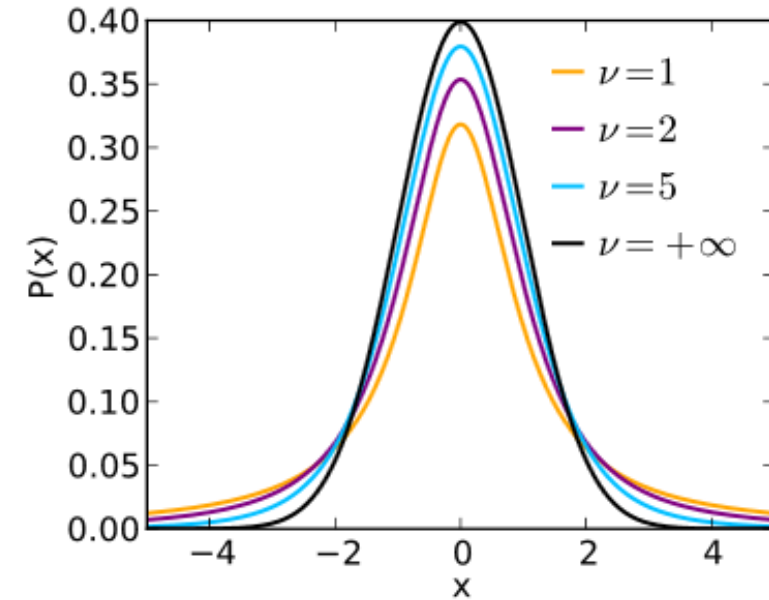
# Statistic of PCC

$$PCC_n = R_n = \frac{1}{n} \sum \left( \frac{x - \mu_x}{\sigma_x} \right) \left( \frac{y - \mu_y}{\sigma_y} \right)$$



$$\mu_{PCC} = E(R_n) = 0$$

$$\sigma_{PCC} = \sqrt{E(R_n^2)} = \sqrt{\frac{1 - R^2}{n - 2}}$$



**PCC distribution**

***PCC<sub>n</sub> follows Student T<sub>n-2</sub> distribution***

As  $n \rightarrow \infty$ ,  $T_{n-2} \rightarrow$  Gaussian Distribution



# Differential PCC

- Differential PCC distribution between  $n$  samples and  $n+1$  samples with  $n$  common samples.

$\Delta PCC_n = PCC_{n+1} - PCC_n$  follows ? distribution

$T_{n-1}$  distribution

$T_{n-2}$  distribution

# Distribution of differential PCC

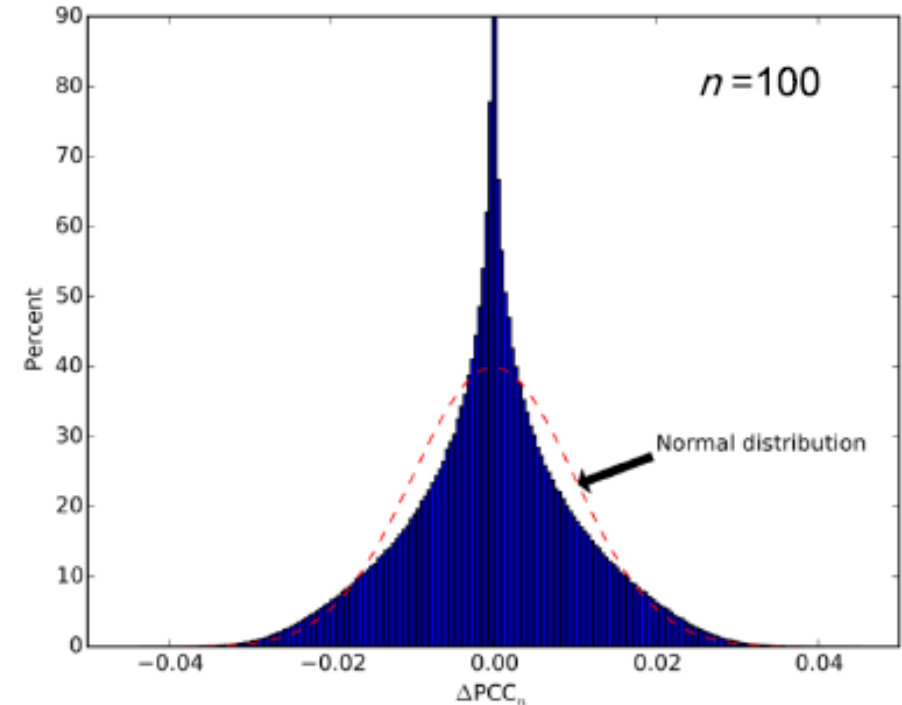
Thus, the mean and standard deviation of the differential  $PCC$

$$\mu_{\Delta PCC} = E(\Delta R_n) = 0$$

$$\sigma_{\Delta PCC} = \sqrt{E(\Delta R_n^2)} = \frac{1}{n-1} (1 - R_n^2)$$

We can use Z statistic to quantify differential PCC

$$v = \frac{\Delta R_n - \mu_{\Delta PCC}}{\sigma_{\Delta PCC} / m}$$



**Volcano distribution**

**Eiffel Tower Distribution**

埃菲尔铁塔 分布

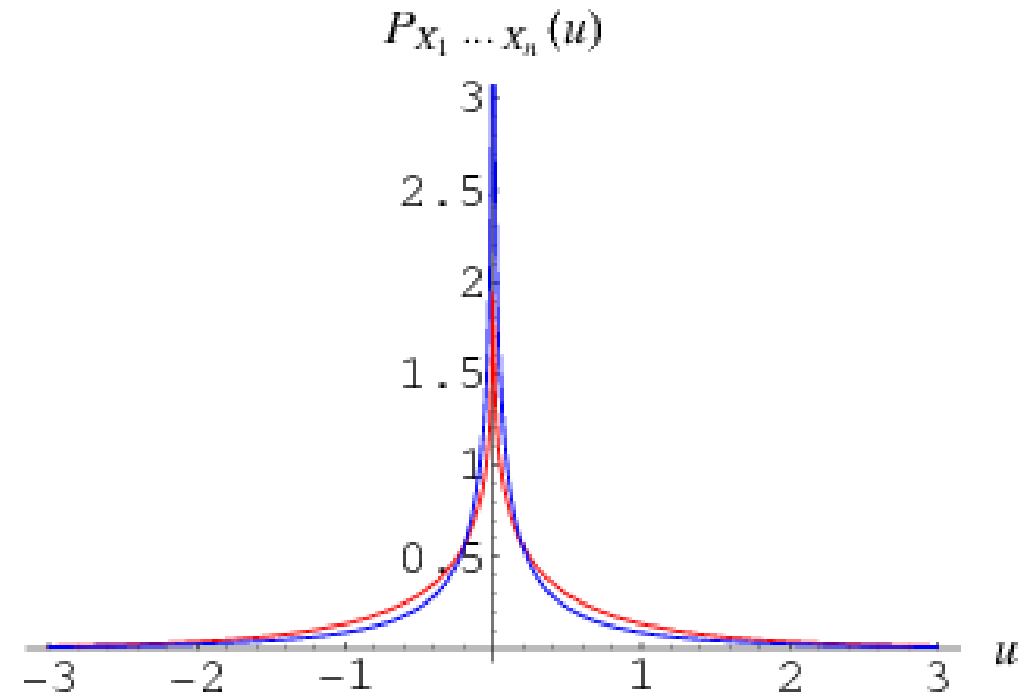
Theoretically, when  $n \rightarrow \text{infinite}$  (Gaussian assumption),  
Volcano distribution approaches normal product distribution with the correlation  $PCC_n$

# Normal Product Distribution

The distribution of a product  $u$  of two normally distributed variates  $x$  and  $y$  with zero means and variances  $\sigma_x^2$  and  $\sigma_y^2$  **without the correlation** is given by

$$P_{X \cdot Y}(u) = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} \frac{e^{-x^2/(2\sigma_x^2)}}{\sigma_x \sqrt{2\pi}} \frac{e^{-y^2/(2\sigma_y^2)}}{\sigma_y \sqrt{2\pi}} \delta(xy - u) dx dy$$

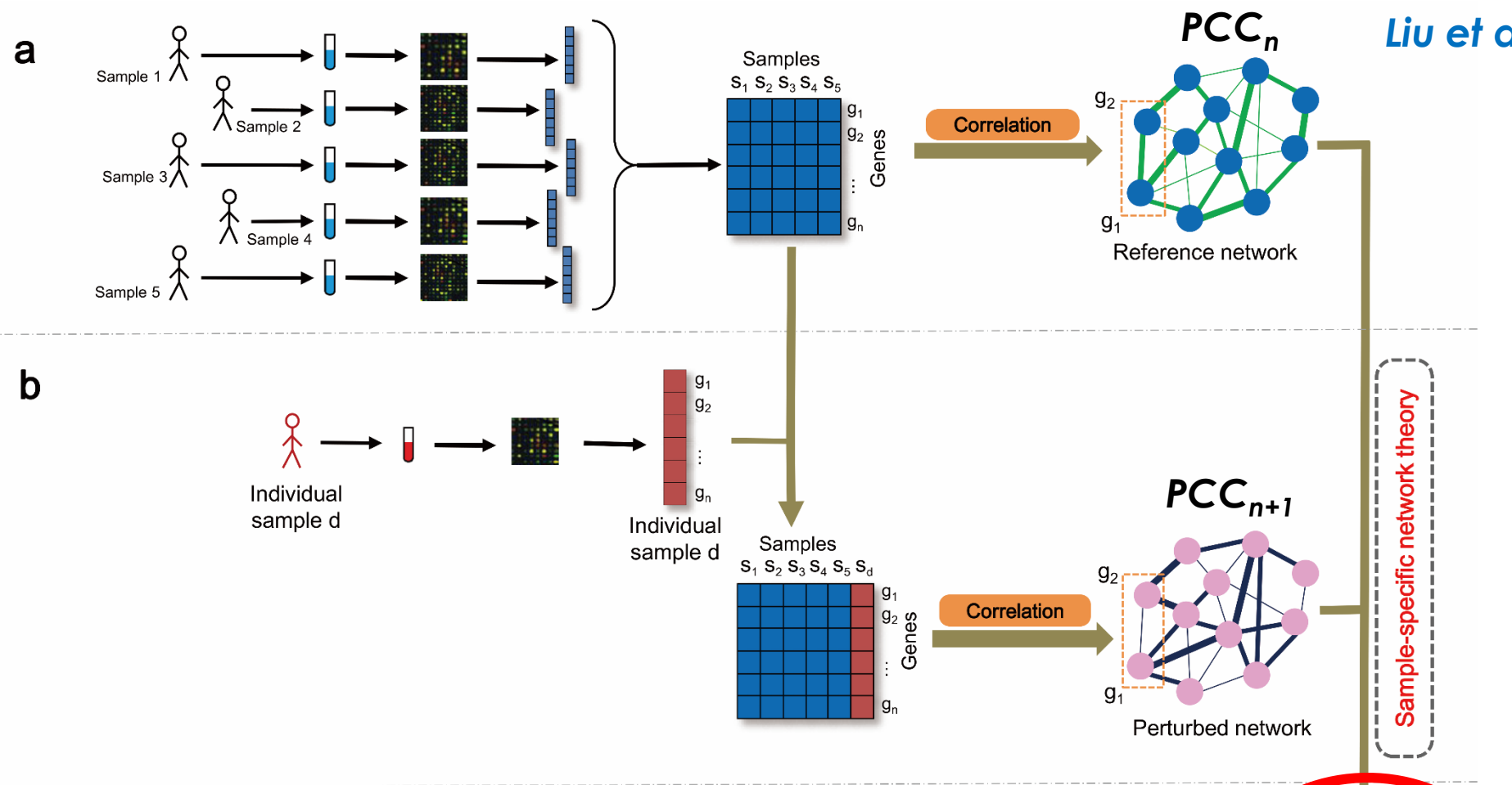
where  $\delta(x)$  is a delta function



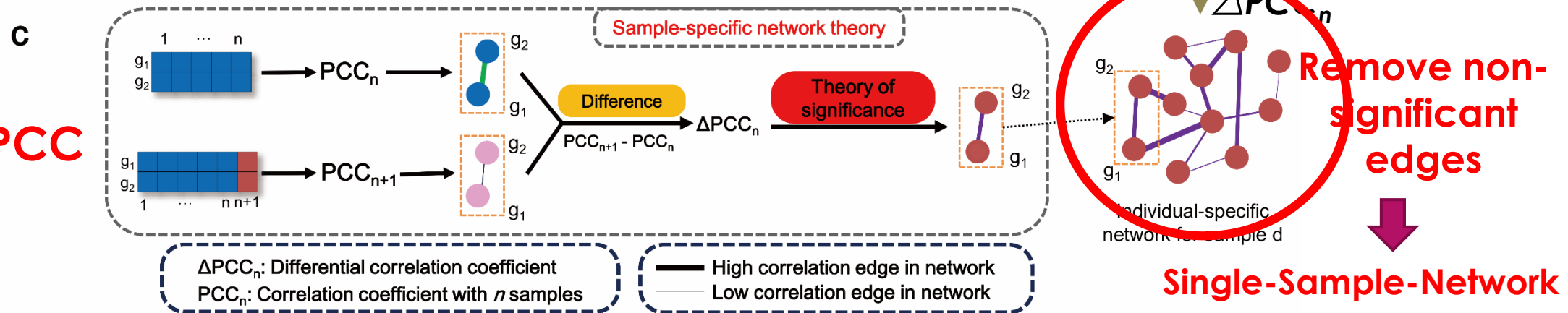
**Theoretically, when  $n \rightarrow \infty$ ,  $\frac{x - \mu_x}{\sigma_x} \frac{y - \mu_y}{\sigma_y} \rightarrow$**

- 1. With Gaussian assumption, Volcano distribution approaches normal product distribution with the correlation.**
- 2. Without Gaussian assumption, Volcano distribution approaches the product distribution of variables  $x$  and  $y$ .**

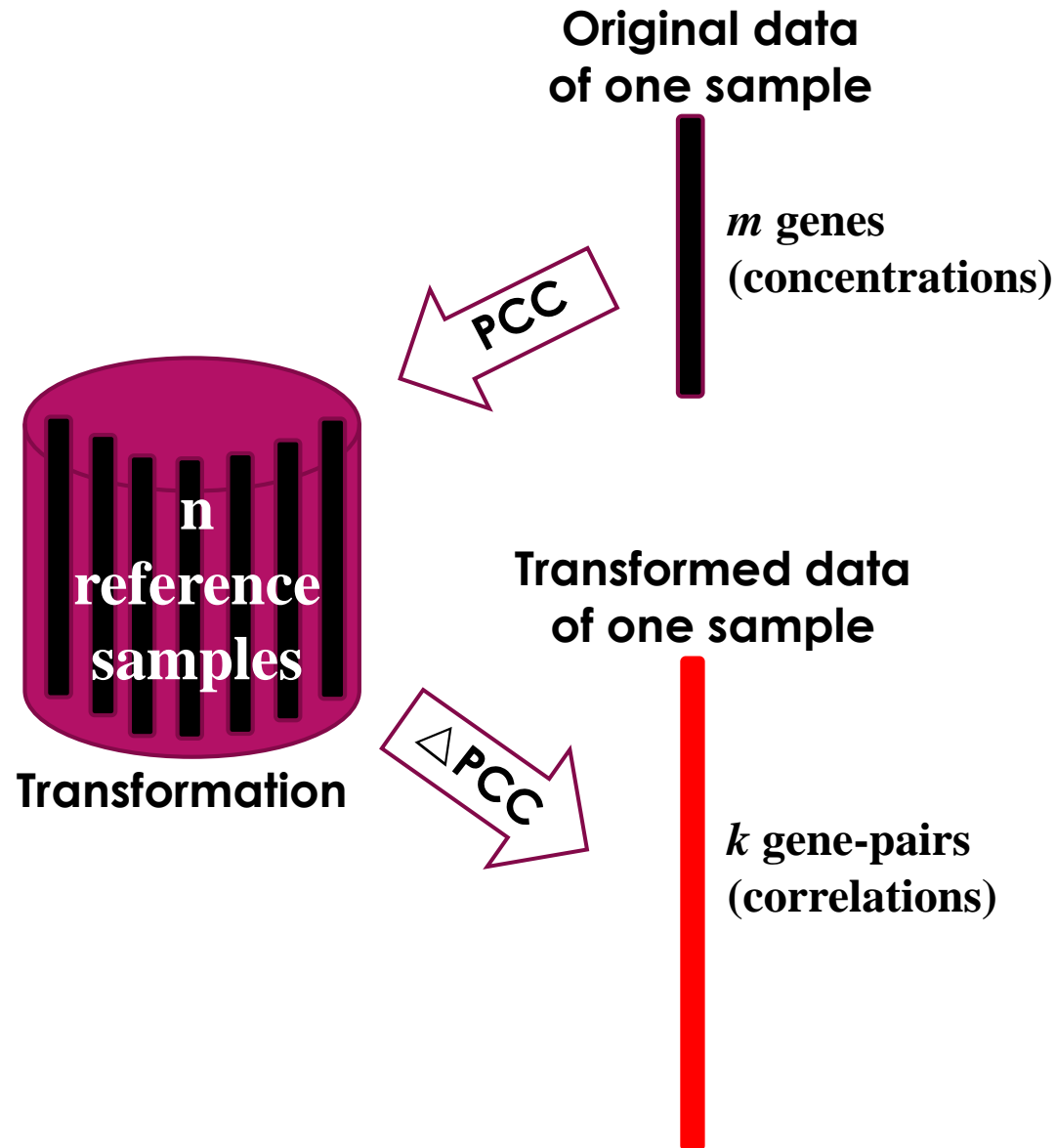
PCC



Differential PCC



# SSN → Data Transformation

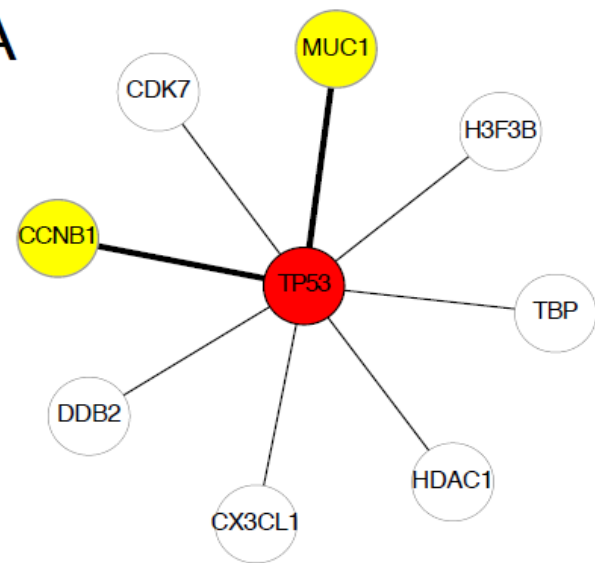


*Similar to  
Fourier or Wavelet transformation,  
changing one type data to another domain*

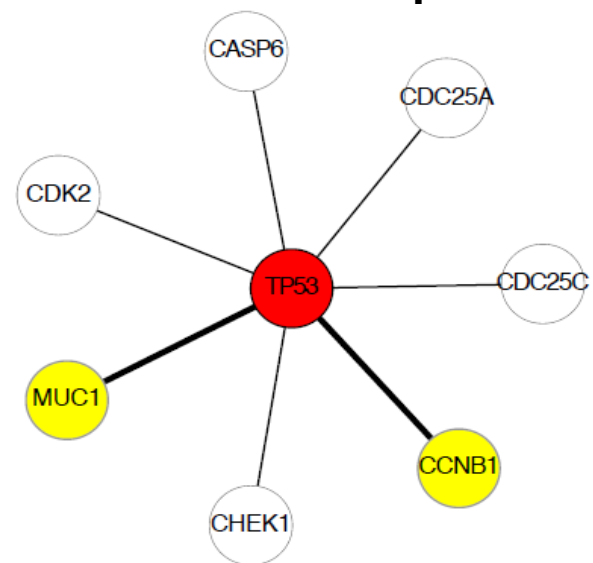
# Validation

- ▶ TCGA (expression data and sequence information)
  - ▶ Breast cancer
  - ▶ Glioblastoma multiforme
  - ▶ Kidney renal clear cell carcinoma
  - ▶ Lung adenocarcinoma
  - ▶ Lung squamous cell carcinoma
  - ▶ Ovarian serous cystadenocarcinoma
  - ▶ Stomach cancer
  - ▶ Thyroid carcinoma
  - ▶ Uterine Corpus Endometrial Carcinoma
- ▶ Constructing individual-network for every disease sample

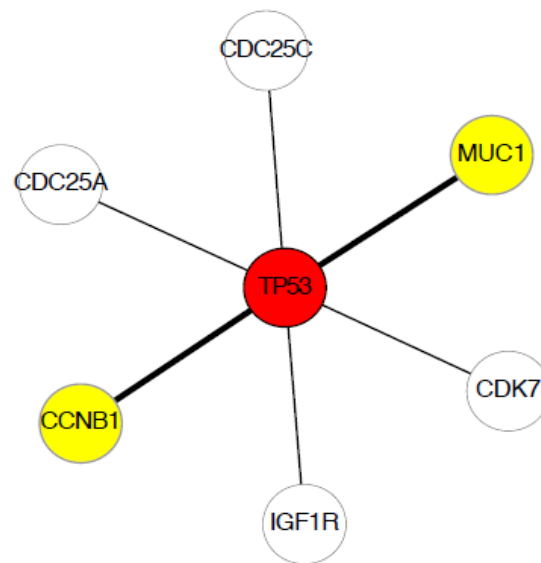
A



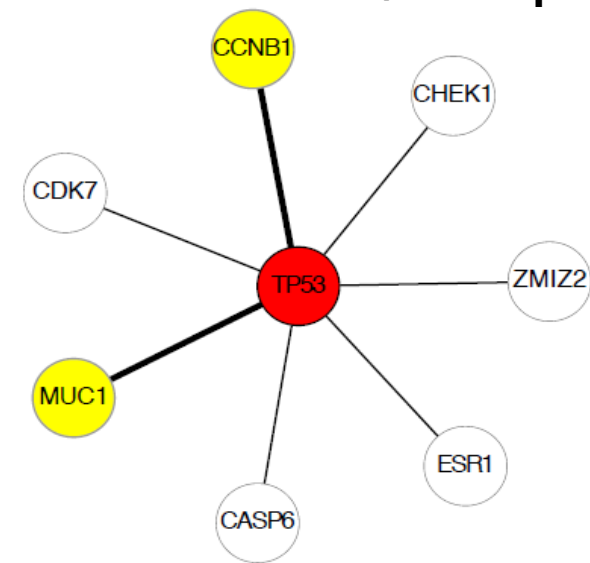
BRCA\_A0B0



BRCA\_A25A

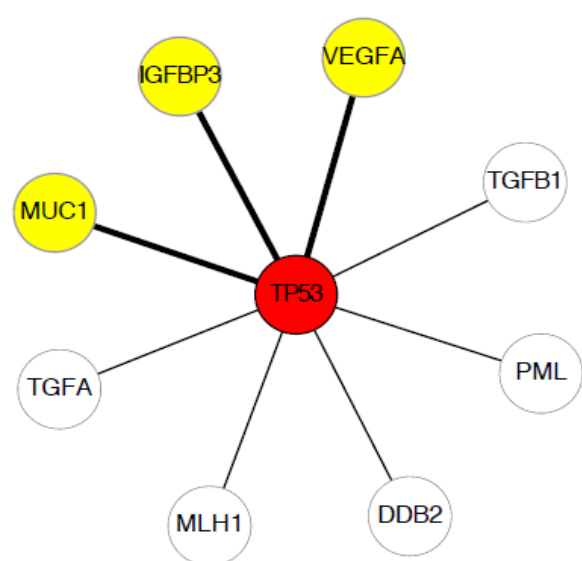


BRCA\_A0HO

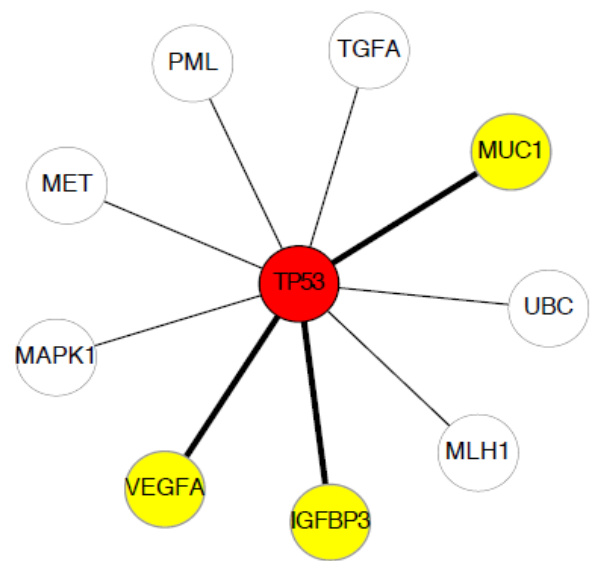


BRCA\_A0DP

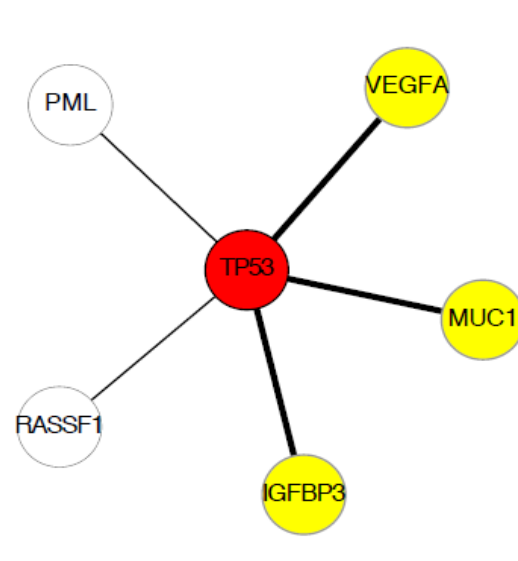
B



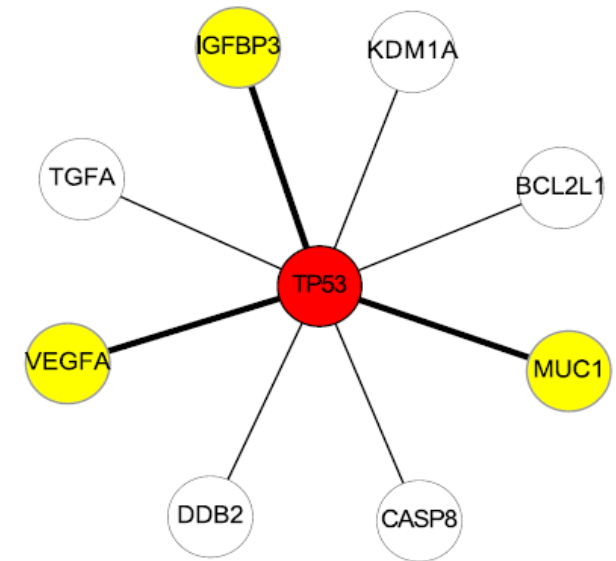
KIRC\_4807



KIRC\_3444



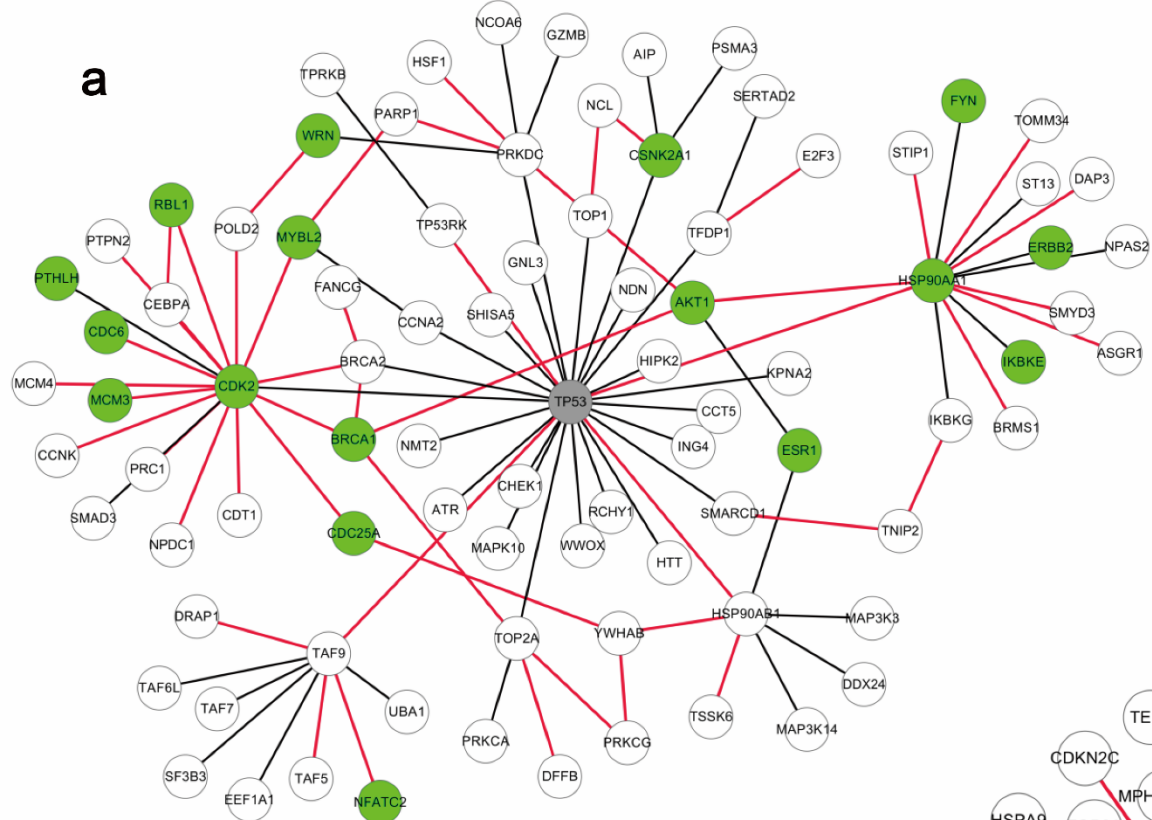
KIRC\_3362



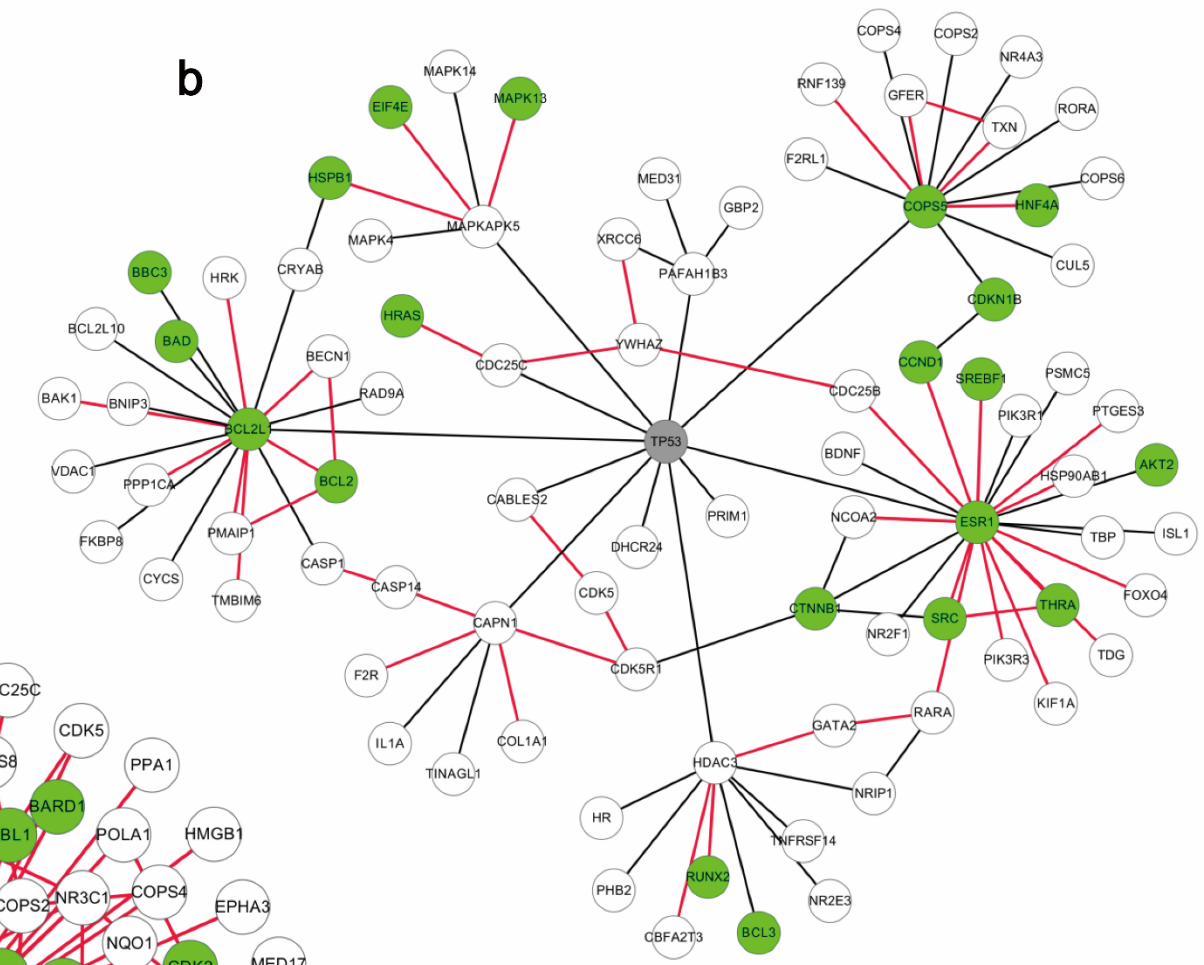
KIRC\_5457

## Individual-specific networks

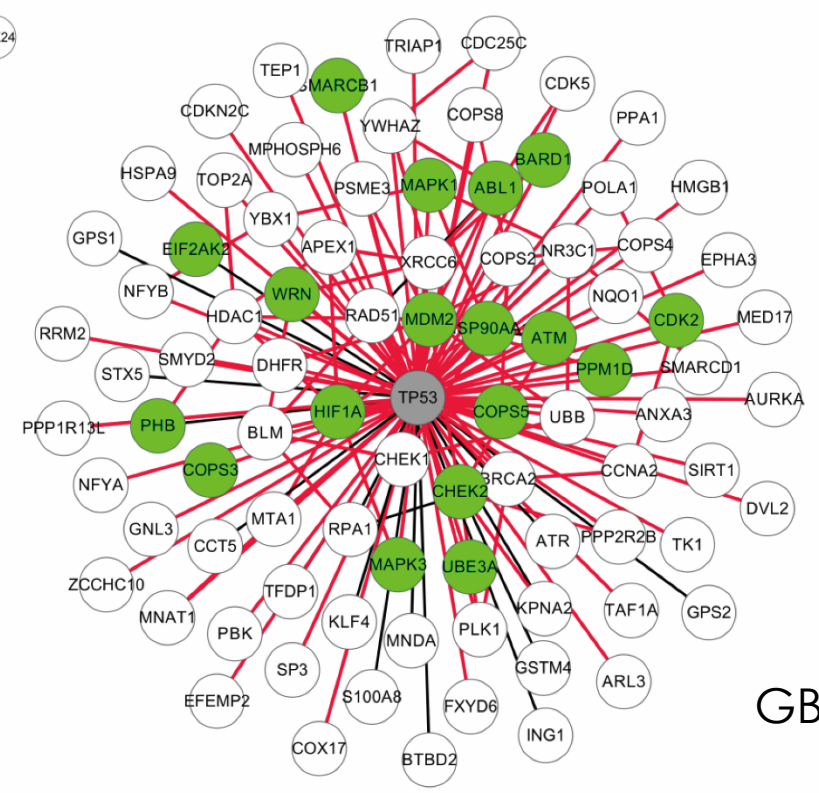
## Over 90% samples



STAD



BRCA

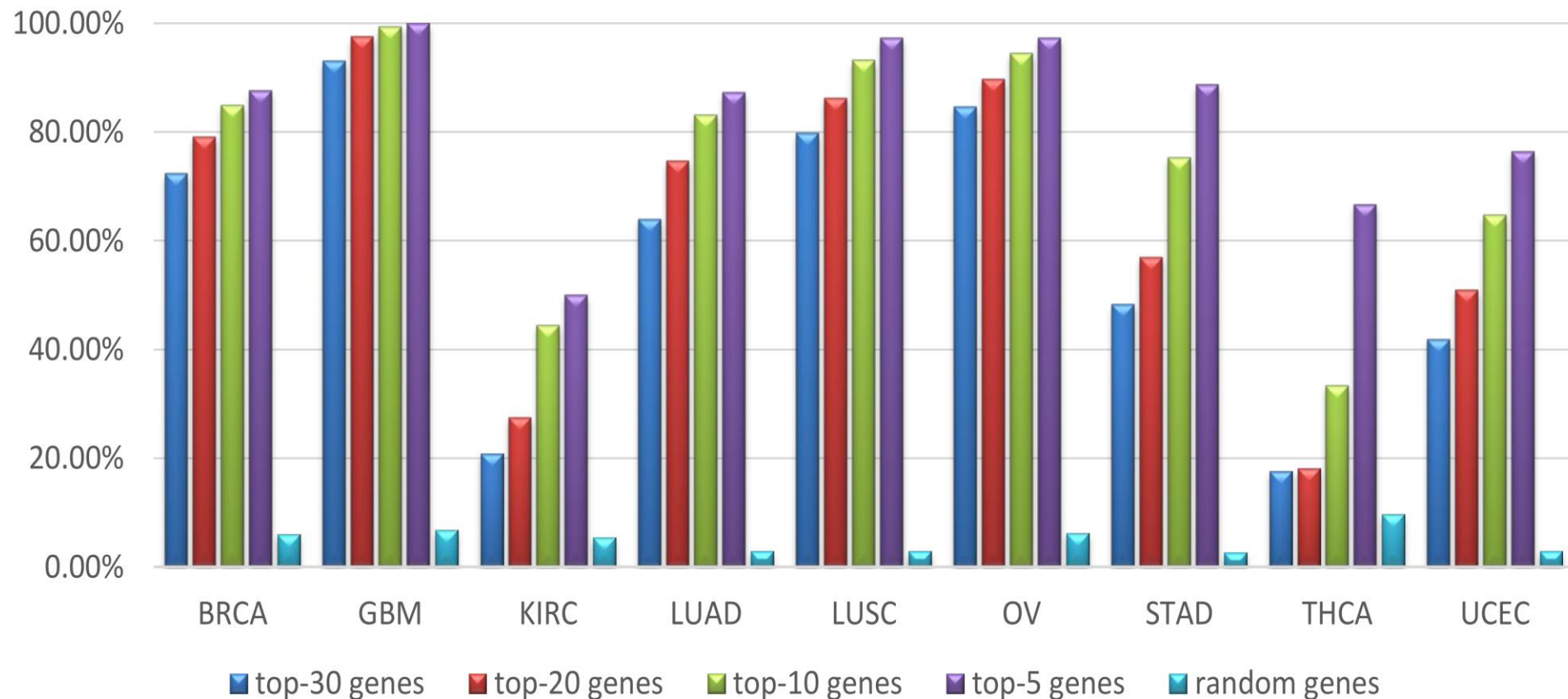


GBM

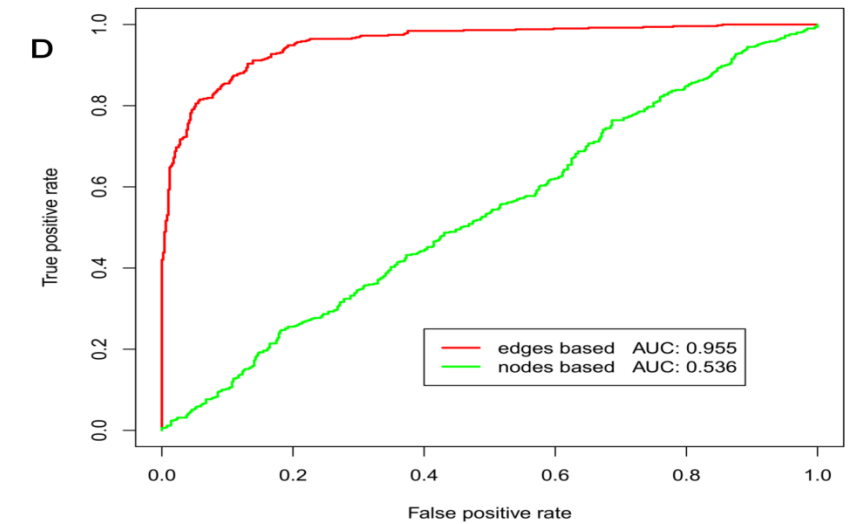
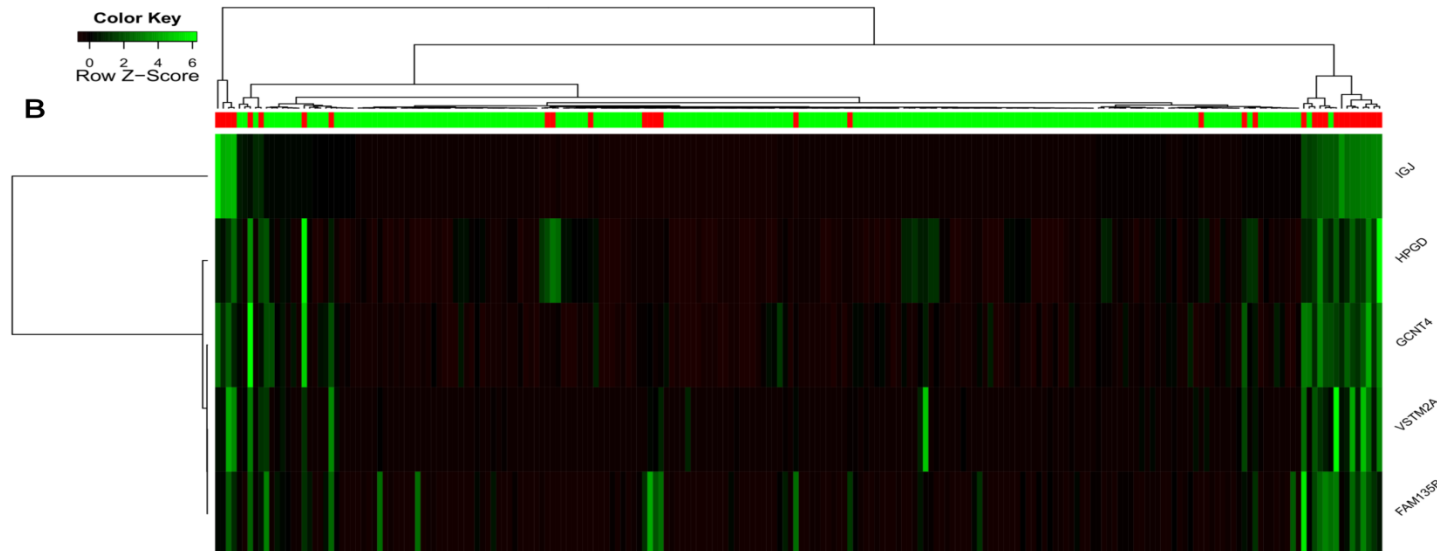
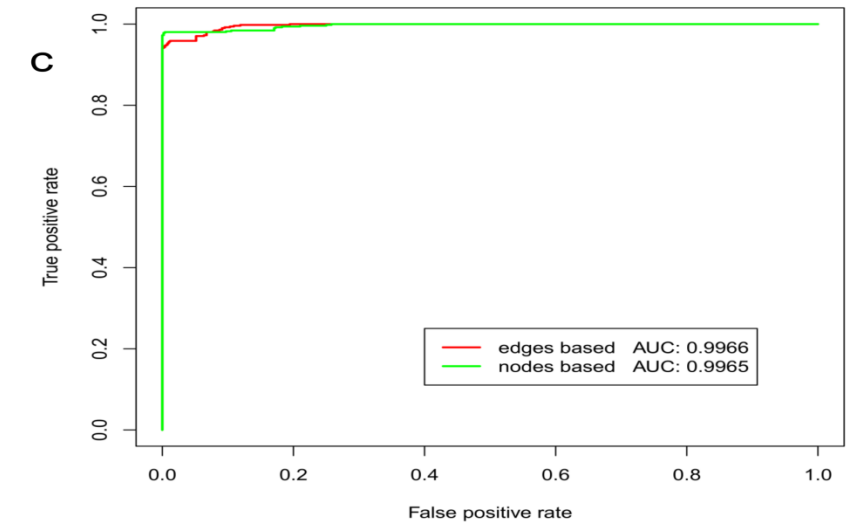
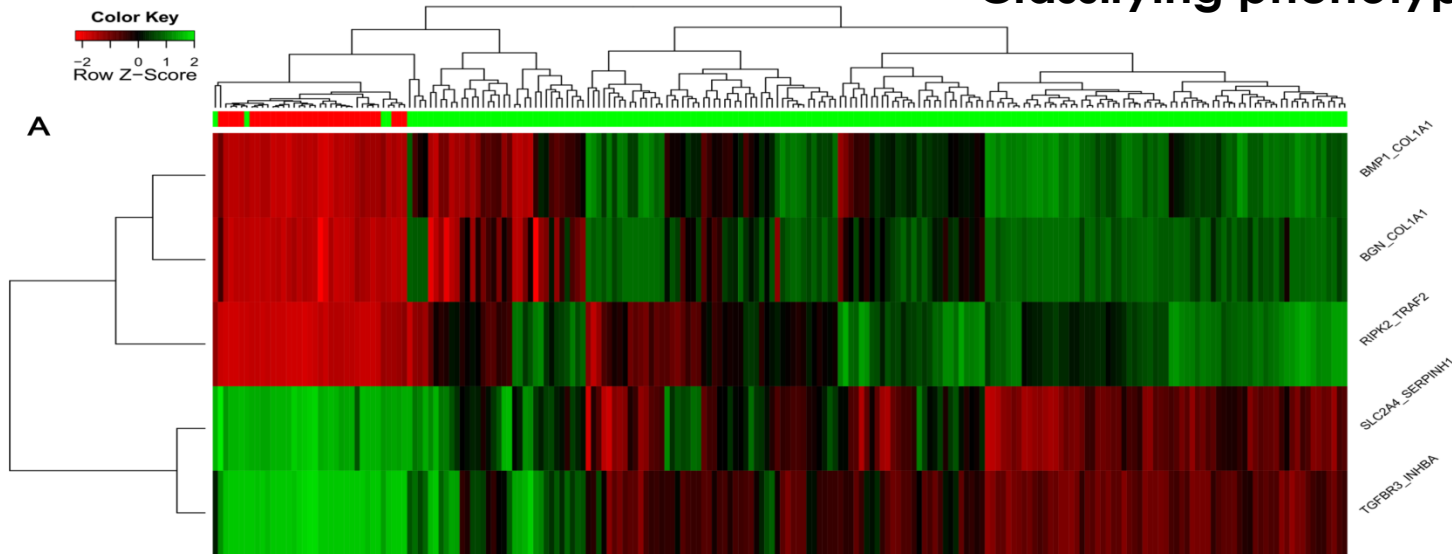


# Driver mutation prediction by expression data

Potential driver genes: Top  $n$  highly connected nodes in the individual network for each sample



# Classifying phenotypes and cancer subtypes



**E**

	BRCA	UCEC	KIRC	LUAD	OV	GBM	LUSC	THCA	STAD
Node markers	0.141	0.319	0.0026	0.0085	0.432	0.502	0.173	0.4	0.829
Edge markers	0.026	0.0037	<1E-6	0.0033	0.056	0.265	0.136	0.37	0.81

# Experiment-1 as network biomarkers

## Cholangiocarcinoma (胆管癌)

- ▶ GSE26566
- ▶ 163 samples
  - ▶ 104 cancerous
  - ▶ 59 non-cancerous

从不变 找到 变化  
**Find variant factors from  
invariant data**

Choose only non-differential genes

1312 non-differential genes were selected with the criterion of  $p\text{-value} > 0.6$

538 gene pairs are differentially correlated under the criterion of  $|\Delta PCC| > 0.9$

**For comparison purpose,**

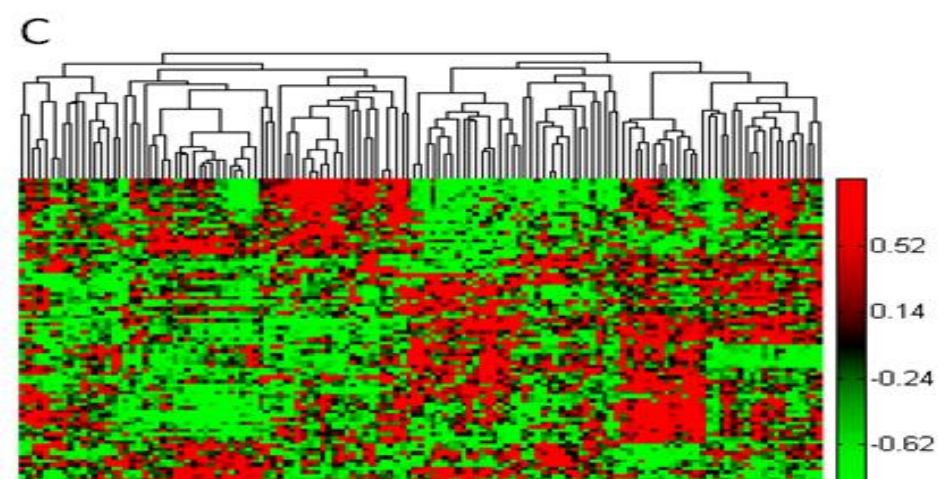
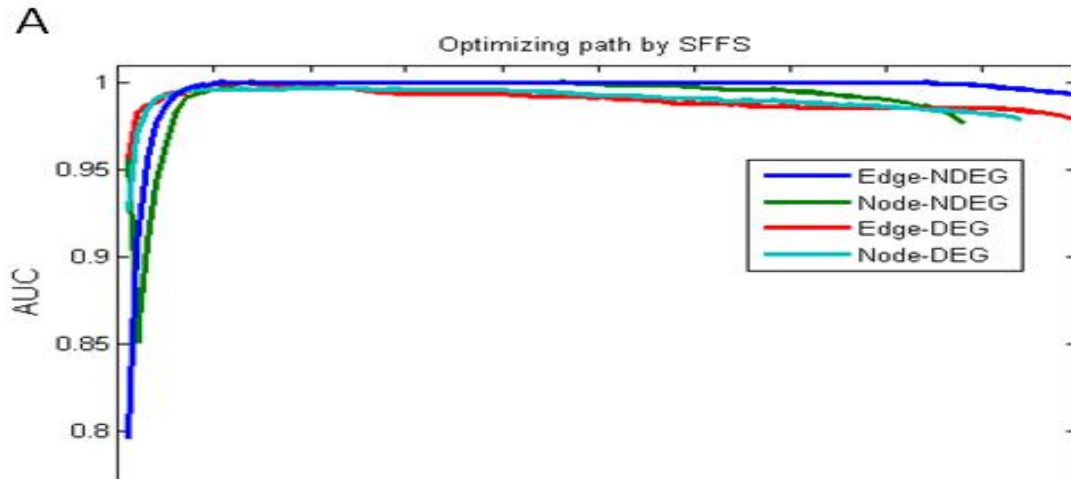
**delete all differential genes  
or  
choose only non-differential genes**

**Dark matter in gene expression 90%**

Dark matter in gene sequence 90%

Dark matter in Universe 90%

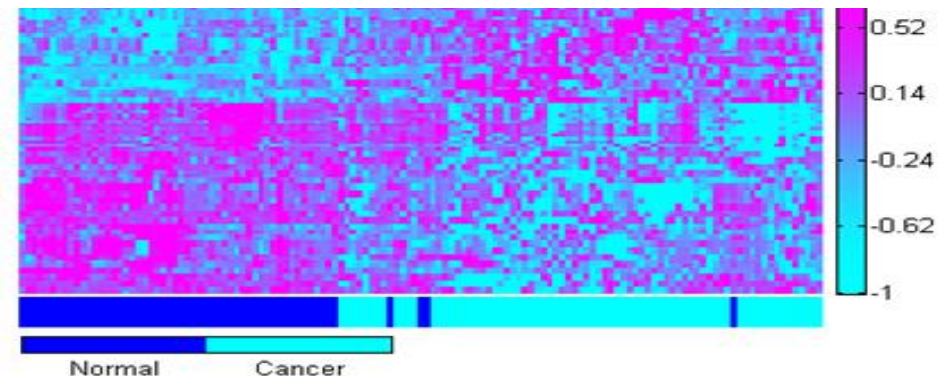
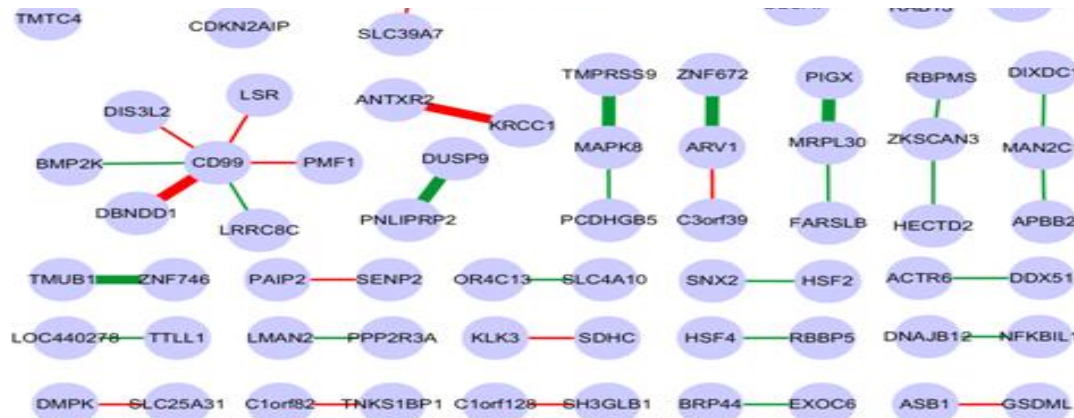
# Cholangiocarcinoma



(A) For NDEG and DEG, the curve by SFFS methods on the top 100 edge features and the involved genes are plotted. (B) The discriminative edge features are presented, where bold lines are

**Non-differential genes achieve the same or better effect than differential genes !!!**

**not in gene level but in network level**



Heatmap of the edge data by (2).

# Experiment-2 as network biomarkers

## Type I diabetes

- ▶ GSE9006
- ▶ 67 samples
  - ▶ 43 newly T1D samples
  - ▶ 24 normal samples

从不变 找到 变化

Choose only non-differential genes

1600 Non differential genes were selected with the criterion of  $P\text{-value} > 0.8$

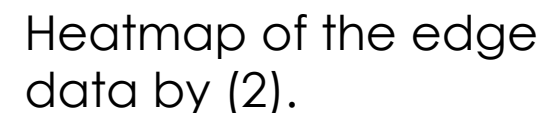
there are 365 Differential Correlation Pairs



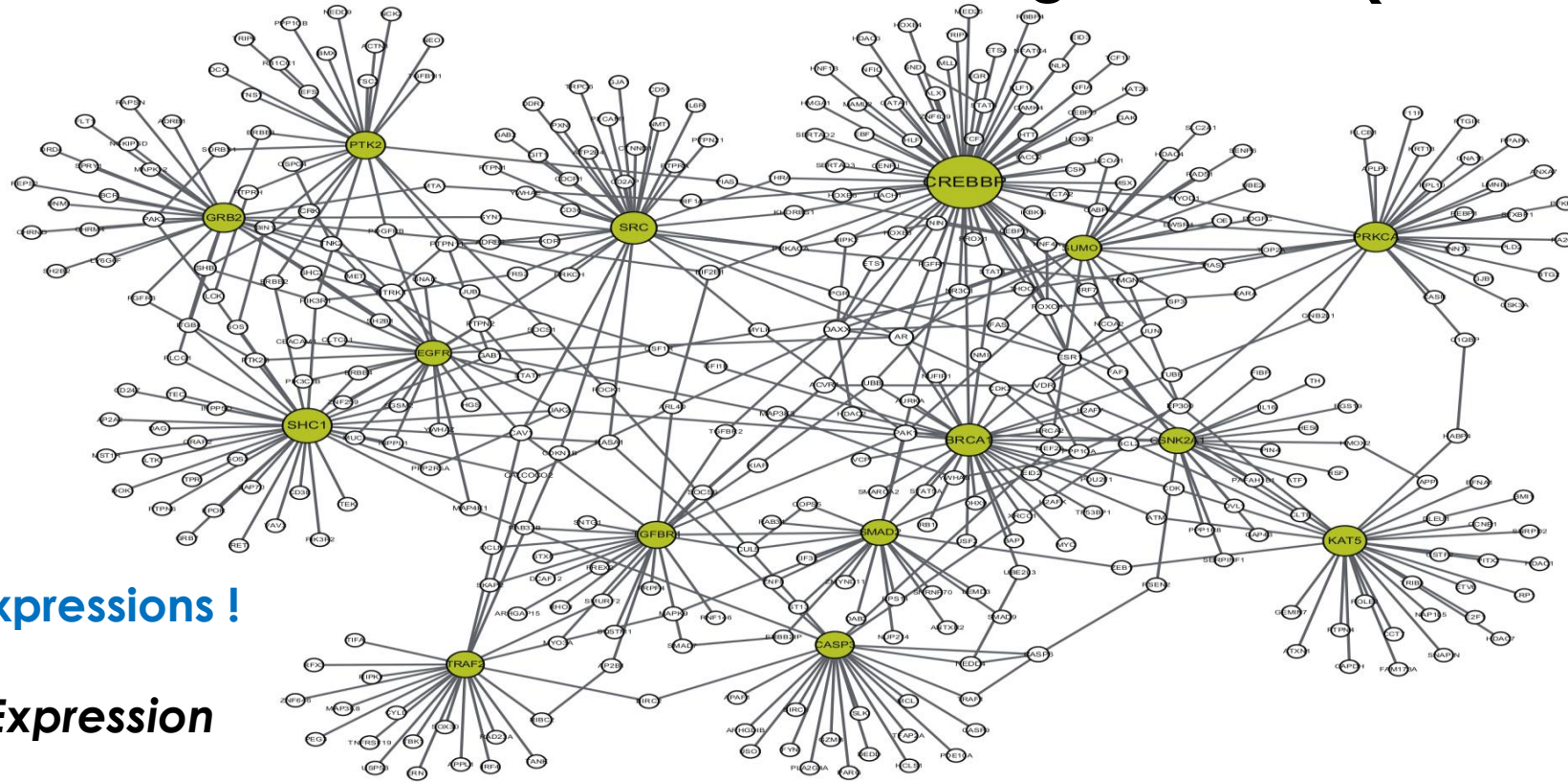
## A



F



# A Drug Resistance Gene Prediction for Lung Cancer (PC9 and PC9-DR)



No differential expressions !

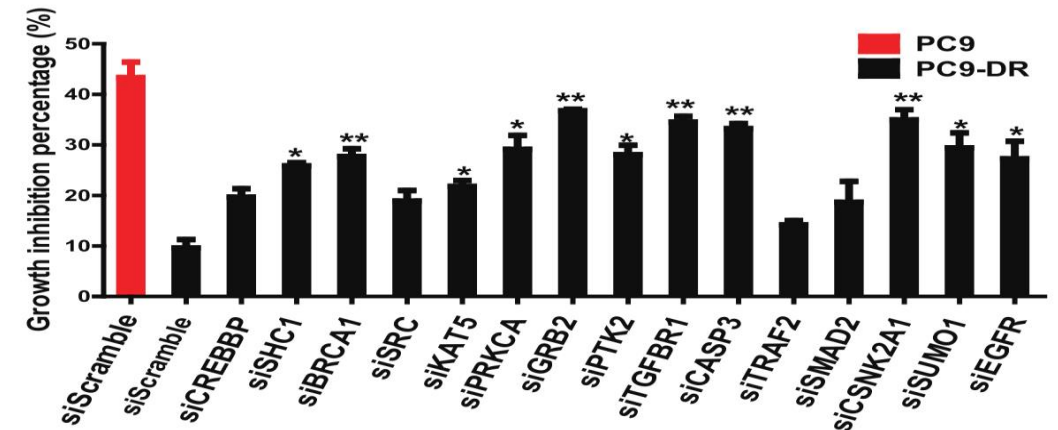


Dark Matter in Expression

B

Gene symbols	Degree	Fold change	Significance to drug resistance
CREBBP	68	0.754549298	No
SHC1	37	1.489521613	Yes*
BRCA1	33	0.838699498	Yes**
SRC	33	1.232327829	No
KAT5	30	1.273340765	Yes*
PRKCA	29	0.746732494	Yes*
GRB2	28	1.30560704	Yes**
PTK2	25	0.761979505	Yes*
TGFBFR1	24	0.712501428	Yes**
CASP3	24	0.876307508	Yes**
TRAF2	23	1.218907247	No
SMAD2	22	0.805846744	No
CSNK2A1	22	1.147045792	Yes**
SUMO1	21	0.694003697	Yes*
EGFR	21	0.897537327	Yes*

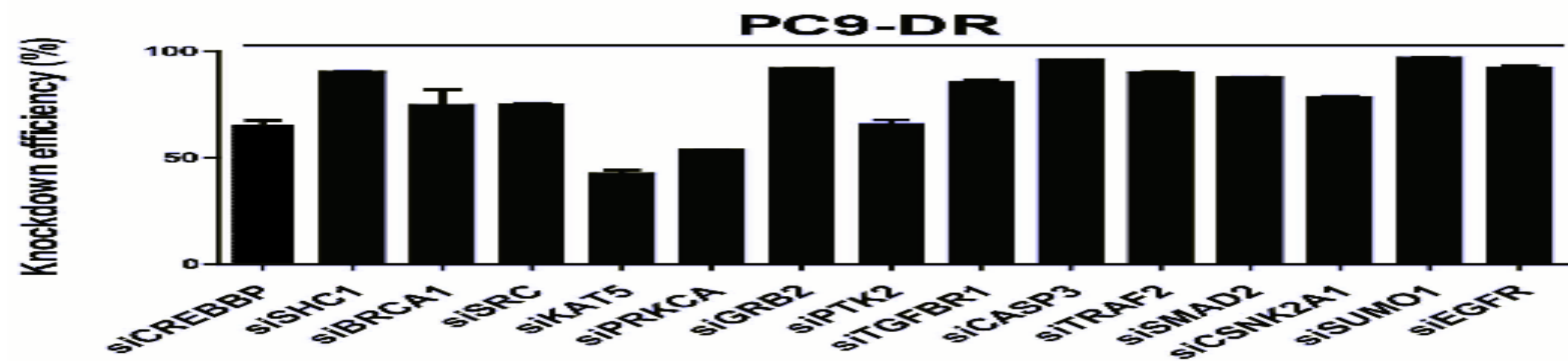
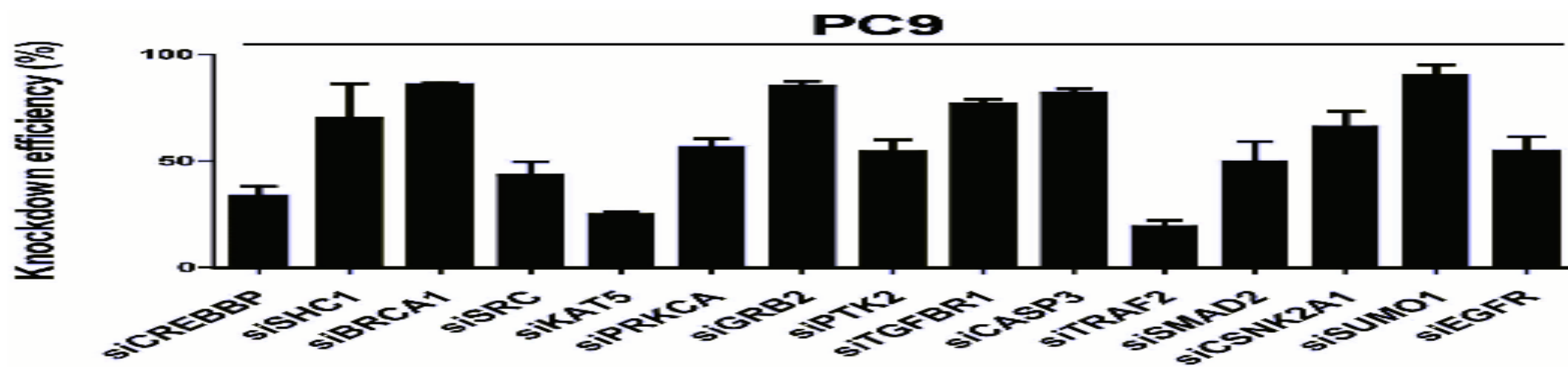
C

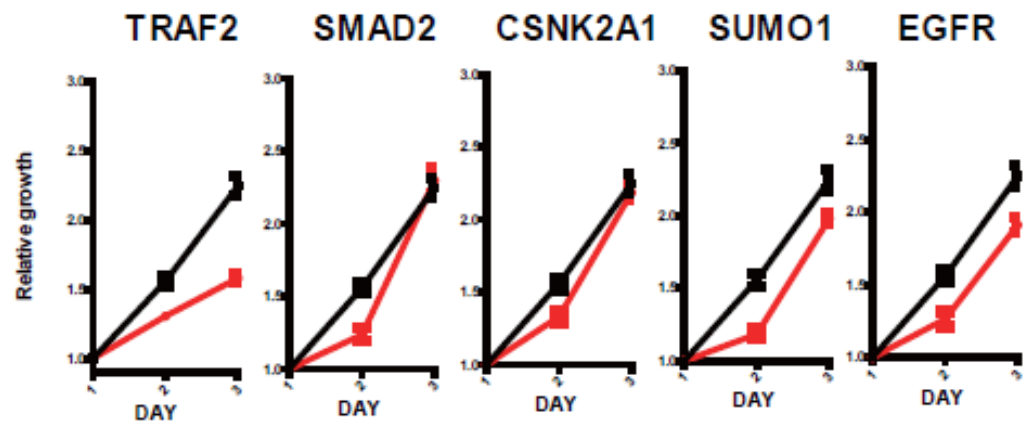
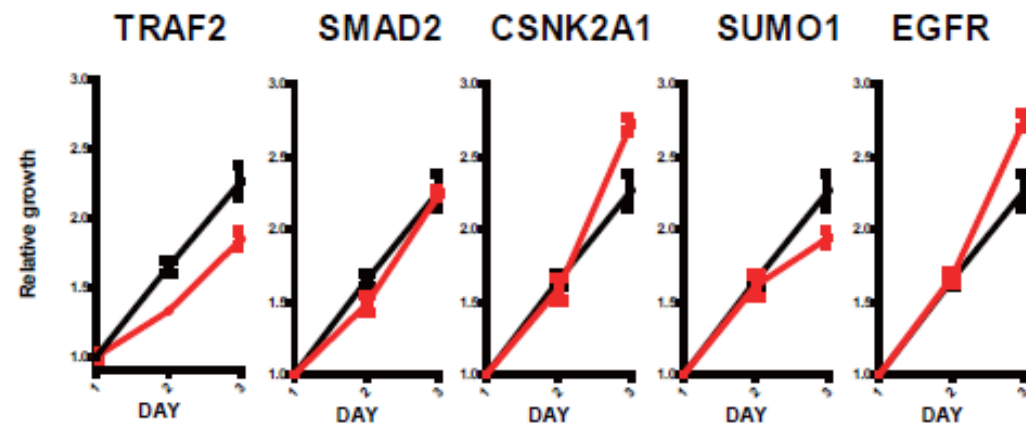
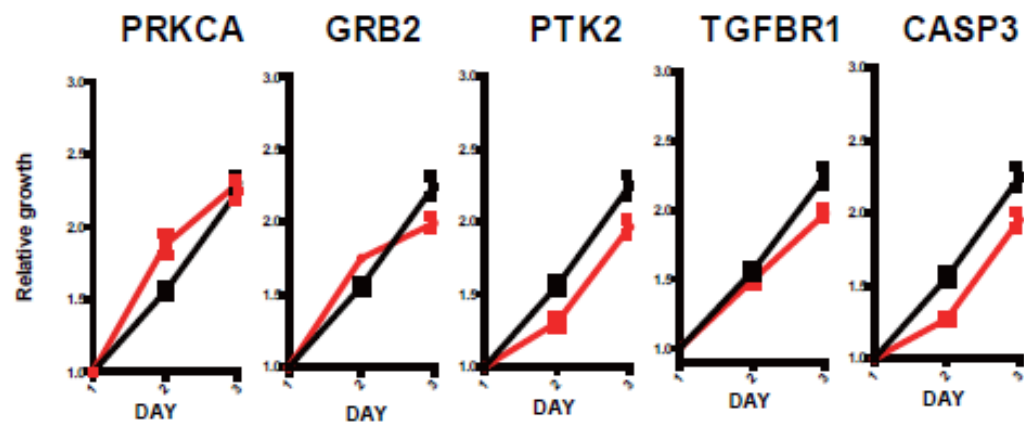
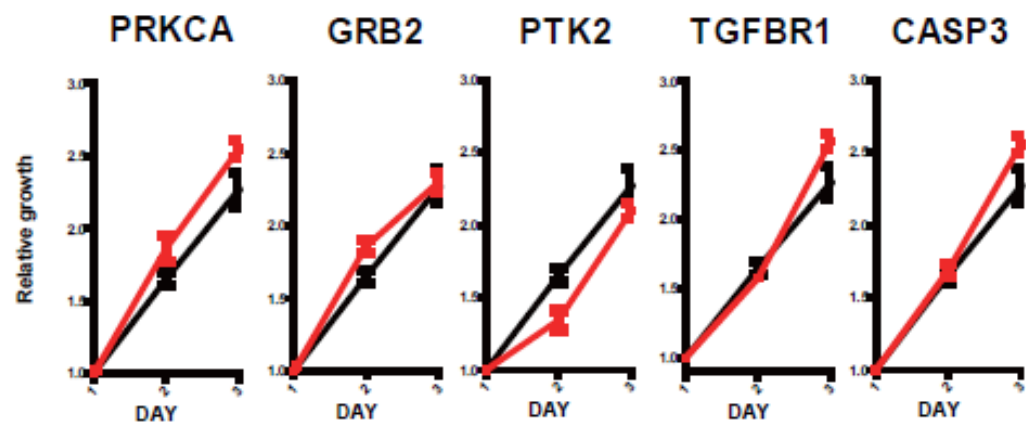
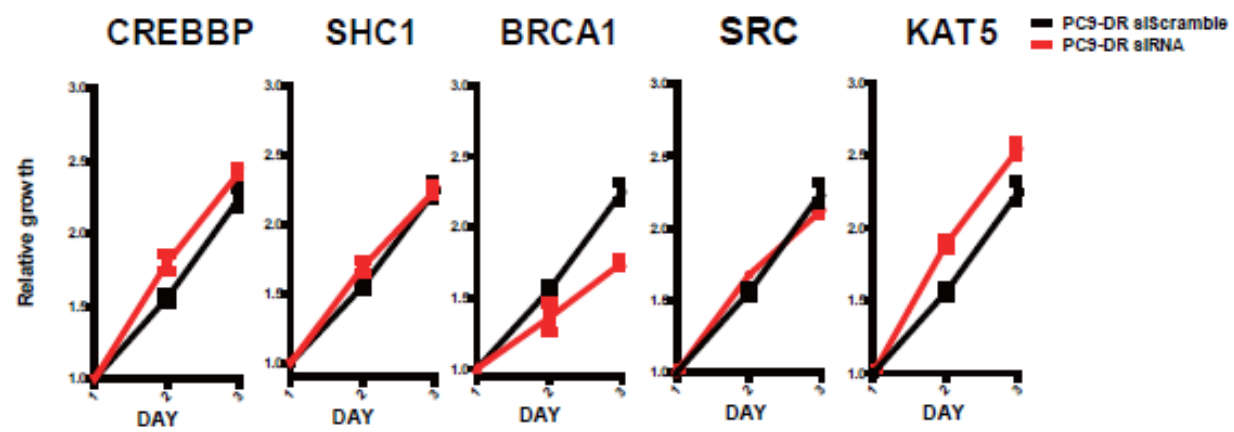
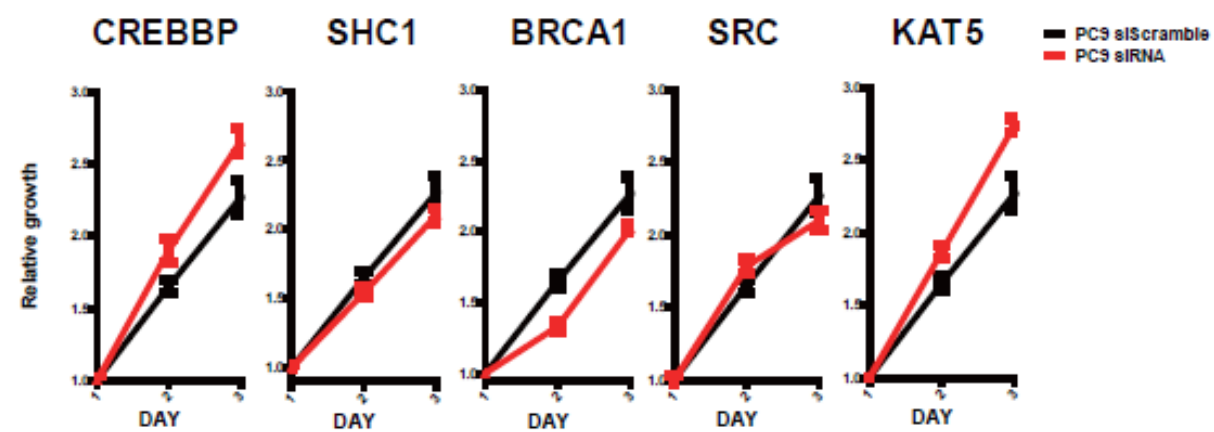




# Dark Matter

- ▶ Universe → 90% Dark Matter in Universe
- ▶ Non-coding RNA → 90% Dark Matter in sequence
- ▶ Non-differential expression → 90% Dark Matter in expression





# Conclusion

All SSNs are comparable, given a reference network

- ▶ Construct Single Sample Network
- ▶ Single-cell network is single-sample network
- ▶ Replace molecular biomarkers by network biomarkers

**Cross species, Cross tissues, Cross platforms**

- ▶ Predict key or driver genes even without differential expressions
- ▶ Achieve network biomarker for personalized medicine
- ▶ Single Sample DNB and DNB landscape

# Acknowledgments

- ▶ Xiaoping Liu,  
Chinese Academy of Sciences
- ▶ Kazuyuki, Aihara,  
The University of Tokyo
- ▶ Hongbin Ji, Yuetong Wang  
Chinese Academy of Sciences

**THANK YOU!**