

# Network-based Stratification of Tumor Mutations

Dongfang Wang

December 28, 2013

Introduction  
oooooooo

Method  
oooooooo

Performance Analysis  
oooooooooooooooooooo

Discussion  
oooo

# Content

## 1 Introduction

## 2 Method

## 3 Performance Analysis

## 4 Discussion

Introduction  
oooooooo

Method  
oooooooo

Performance Analysis  
oooooooooooooooooooo

Discussion  
oooo

# Content

## 1 Introduction

## 2 Method

## 3 Performance Analysis

## 4 Discussion

# Introduction

## Cancer

- ① complex : driven by a combination of genes
- ② **heterogeneous:** gene combinations vary greatly between patients

# Introduction

## Cancer

- ① complex : driven by a combination of genes
- ② **heterogeneous:** gene combinations vary greatly between patients

## TCGA:The Cancer Genomes Atlas

profile thousands of tumors at multiple layers of genome-scale information

- ① +20 Cancer cohorts with 50-800 individuals
- ② samples of different measurement types:
  - mRNA expression, Protein expression
  - Copy number variations, Single nucleotide polymorphisms, Methylation, miRNA
  - Patient genomes (**somatic mutations**)

# Introduction

## Cancer

- ① complex : driven by a combination of genes
- ② **heterogeneous:** gene combinations vary greatly between patients

## TCGA:The Cancer Genomes Atlas

profile thousands of tumors at multiple layers of genome-scale information

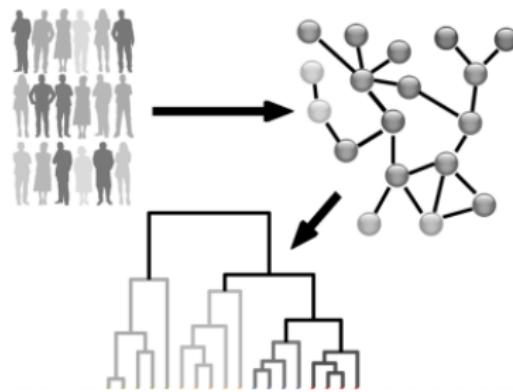
- ① +20 Cancer cohorts with 50-800 individuals
- ② samples of different measurement types:
  - mRNA expression, Protein expression
  - Copy number variations, Single nucleotide polymorphisms, Methylation, miRNA
  - Patient genomes (**somatic mutations**)

## Informatic methods: integrate and interpret genome-scale information

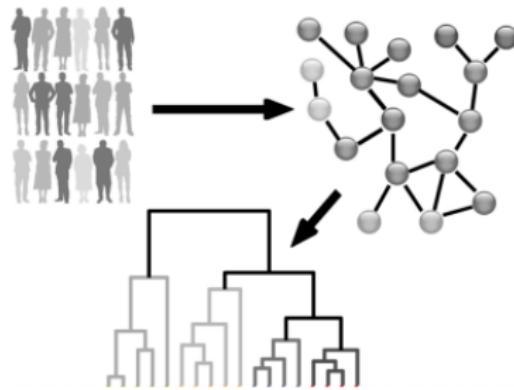
- ① molecular process driving tumor progression
- ② need in the clinic (genome-scale tumor profiles are unable to derive clinically relevant data)



## Stratification: driving cancer into subtypes



## Stratification: driving cancer into subtypes

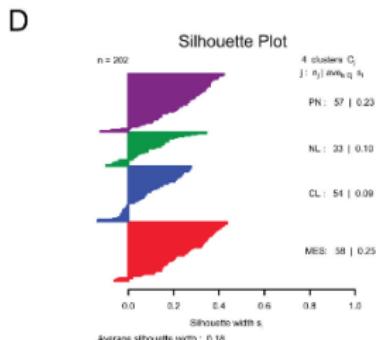
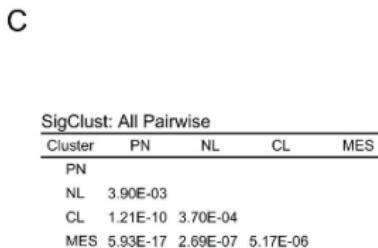
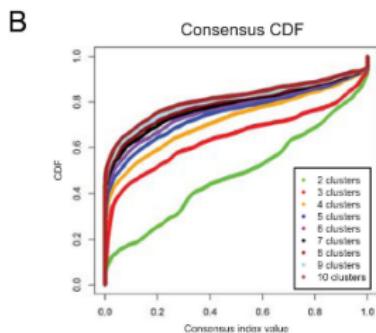
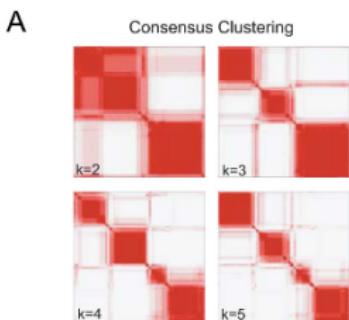


## Why stratify?

subtype: clinical and biological meaning

- ① Better patient prognostics
  - ② Better undersatnding of tumor biology
  - ③ subtype specific drug targets
  - ④ Patient tailored treatment

## Efforts to Stratify Using Gene Expression

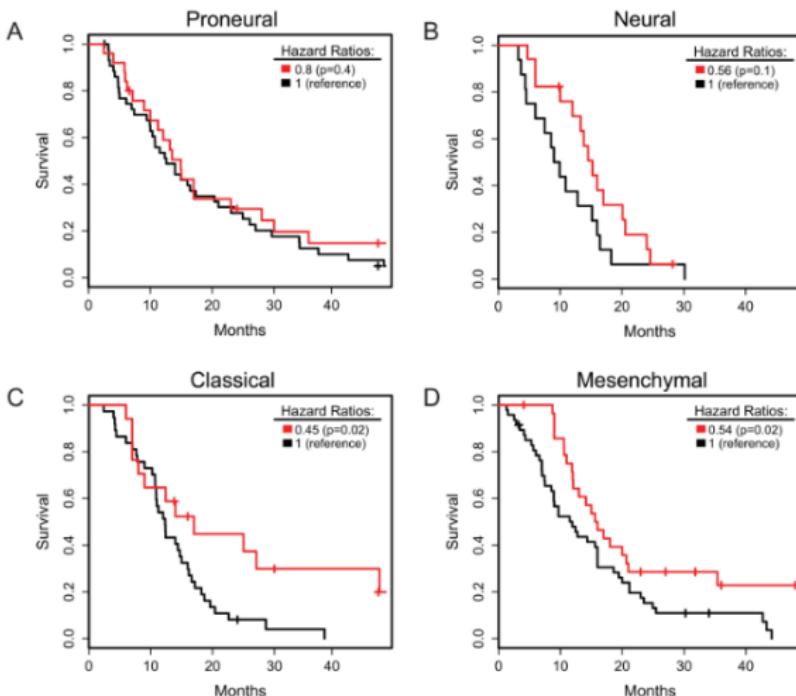


## factor analysis

(202 samples&1740 genes)

## consensus clustering

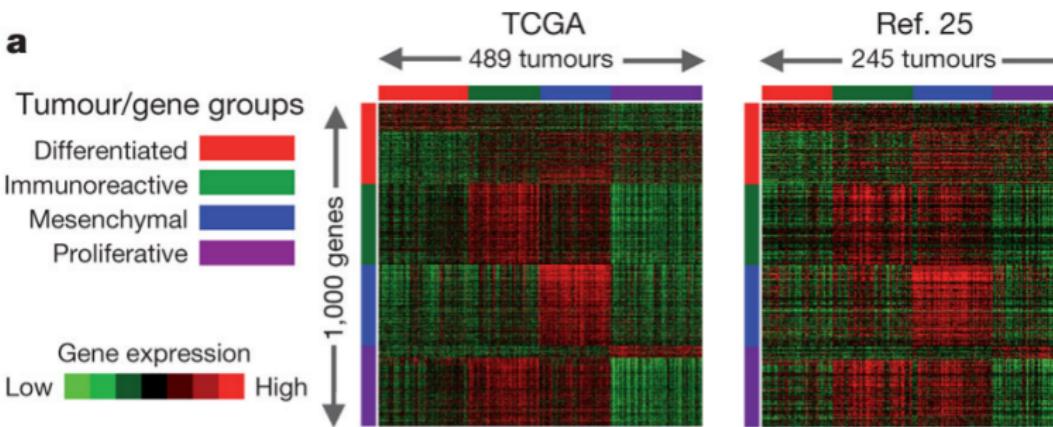
Verhaak, R.G. et al. Integrated genomic analysis identifies clinically relevant subtypes of glioblastoma characterized by abnormalities in PDGFRA, IDH1, EGFR, and NF1. *Cancer Cell* 17, 98-110 (2010)



- More intensive therapy: concurrent chemotherapy/radiation and/or >3 cycles of chemotherapy
- Less intensive therapy: non-concurrent chemotherapy/radiation or <4 cycles of chemotherapy

greatest benefit in Classical and no benefit in Proneural

# Problems



No association to a clinical phenotype was reported (for these subgroups).

T.C.G.A.R.N. (TCGA), Integrated genomic analyses of ovarian carcinoma. Nature 474, 609-15 (2011).

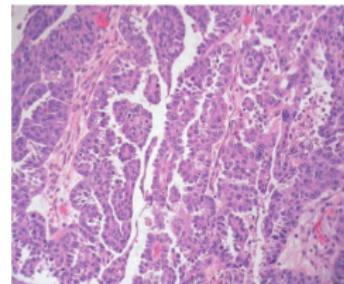
## Problems

- ① RNA sample quality
- ② lack of reproducibility
- ③ ample opportunities of overfitting of data

# Somatic Mutations

## Somatic mutations:

- high-throughput sequencing
- compare the genome or exome of a patient's tumor to that of the germ line to identify mutations that have become enriched in the tumor cell population
- contain causal drivers of tumor progression



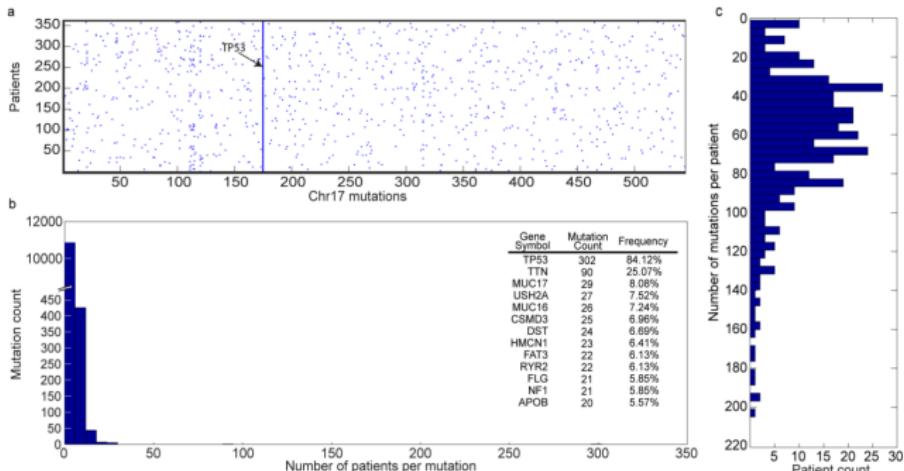
The Cancer Genome Atlas 

Somatic mutations in high grade serous ovarian cancer

(359 matched patient/tumor exome sequenced with Illumina GAILx; 11,231 somatic mutations )

# Challenges

- ➊ extremely sparse
- ➋ heterogeneous, clinically identical patients share no more than a single mutation



# Hypothesis

- ① Cancer is a disease not of individual mutations, nor of genes, but of combinations of genes acting in molecular networks corresponding to hallmark process
- ② although two tumors may not have any mutations in common, they may **share the networks affected by these mutations**

NATURE METHODS | ARTICLE OPEN



## Network-based stratification of tumor mutations

Matan Hofree, John P Shen, Hannah Carter, Andrew Gross & Trey Ideker

Affiliations | Contributions | Corresponding author

*Nature Methods* 10, 1108–1115 (2013) | doi:10.1038/nmeth.2651

Received 14 February 2013 | Accepted 12 August 2013 | Published online 15 September 2013

Introduction  
oooooooo

Method  
oooooooo

Performance Analysis  
oooooooooooooooooooo

Discussion  
oooo

# Content

## 1 Introduction

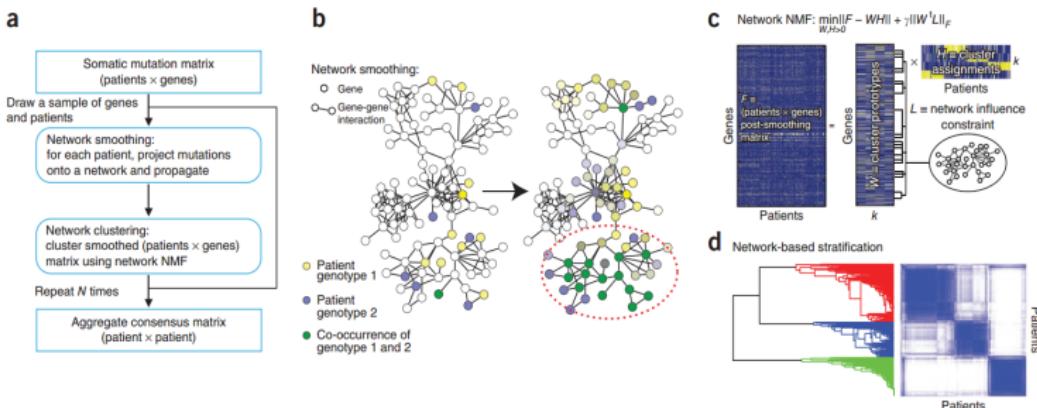
## 2 Method

## 3 Performance Analysis

## 4 Discussion

# Overview of NBS Methods

- ① Somatic mutation matrix
- ② Network smoothing
- ③ Network cluster
- ④ Consensus clustering



# Data Source

## Patient Mutation Profile

- TCGA data portal
- ovarian cancer, uterine endometrial carcinoma, lung adenocarcinoma somatic mutation data
- binary matrix; (0,1) states on genes

## Molecular Network Data

Gene interaction networks:

- STRING v.9
- HumanNet v.1
- PathwayCommons

# Network Smoothing

- ① Mapping a patient mutation profile onto a molecular network
- ② "Network propagation": smooth the mutation signal across the network

# Network Smoothing

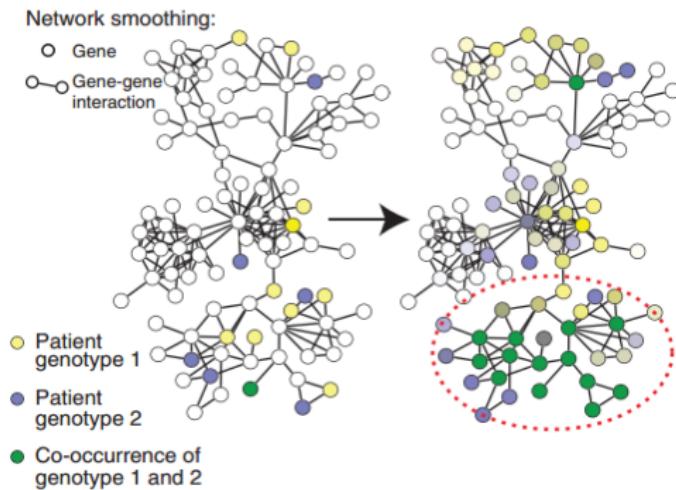
- ① Mapping a patient mutation profile onto a molecular network
- ② "Network propagation": smooth the mutation signal across the network

## Random walk with restarts on a network

$$F_{t+1} = \alpha F_t A + (1 - \alpha) F_0$$

- $F_0$ : a patient-by-gene matrix
  - $A$ : degree-normalized adjacency matrix of the gene interaction network
  - $\alpha$ : tuning parameter governing the distance that a mutation signal is allowed to diffuse through the network during propagation
- 
- ①  $\alpha$ : only a minor effect on the results of the NBS over a sizeable range (0.5-0.8)
  - ② rows are **quantile normalized** to ensure that the smoothed mutation profile for each patient follows the same distribution

# An Intuition for Network Smoothing

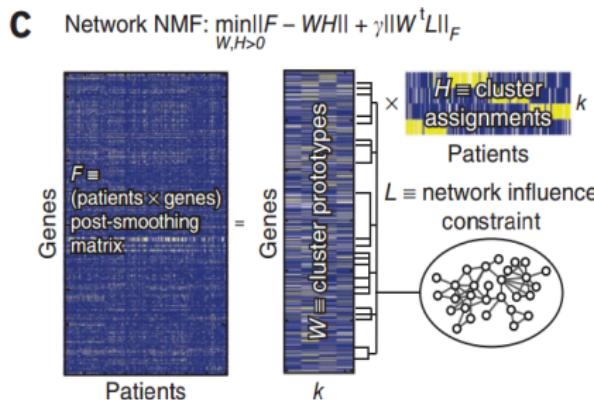


Genotype A: [...,0,0.3,1,0.4,0.4,1,0.2,0,0,0.3,0,0,0,0.4,0,1,0.1,0,...]  
Genotype B: [...,0,0.2,1,0.2,0.4,0.2,0,0,0.4,1,0,0,0,0,0.2,0.4,0.5,0,...]

# Network Clustering

## Network-regularized NMF

$$\min_{W,H>0} ||F - WH||^2 + \gamma \text{trace}(W^T KW)$$

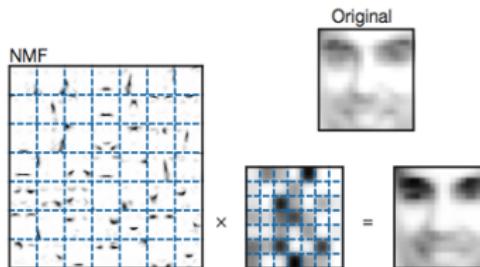


- ① NMF: learn the parts of the object
- ② Network-regularized: respect the structure of an underlying gene interaction network

# Non-negative matrix factorization

$$\mathbf{V} \approx \mathbf{WH}$$

- ①  $\mathbf{V}$  is an  $n \times m$  matrix, and  $\mathbf{M} \in R^{n \times r}$ ,  $\mathbf{H} \in R^{r \times m}$ .



- ② all elements of the three matrices are **non-negative**  
③  $\mathbf{W}$  : every column is a **basis image**  
 $\mathbf{H}$  : **encoding**; coefficients of the linear combination of basis images

$$\vec{V}_i = \mathbf{W}\vec{H}_i$$

## Network-regularized

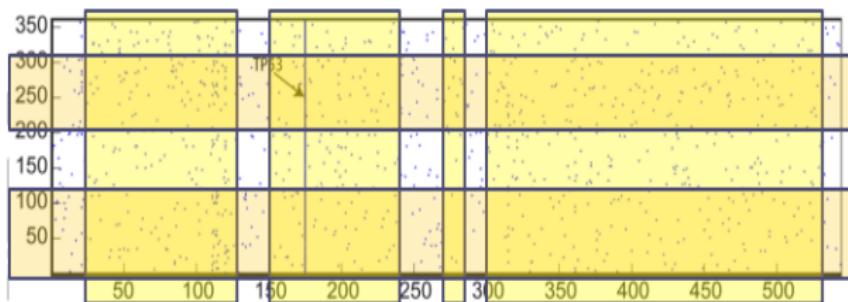
$$\min_{W,H>0} \quad ||F - WH||^2 + \gamma trace(W^T KW)$$

$$trace(W^T KW) = \sum_{ij} k_{ij}(w_i)^T (w_j)$$

constrain the basis vectors in  $W$  to respect local network neighbourhoods

# Consensus Clustering

- ① Draw a bootstrap sample of genes from  $patients \times genes$  matrix



- ② re-sampling 1000 times (80% of the patients and 80% of the mutants genes at random without replacement )
- ③ A similarity matrix of patients (how many times of a pair of patients appears in the same subtype)
- ④ stratify the patients

# Content

## 1 Introduction

## 2 Method

## 3 Performance Analysis

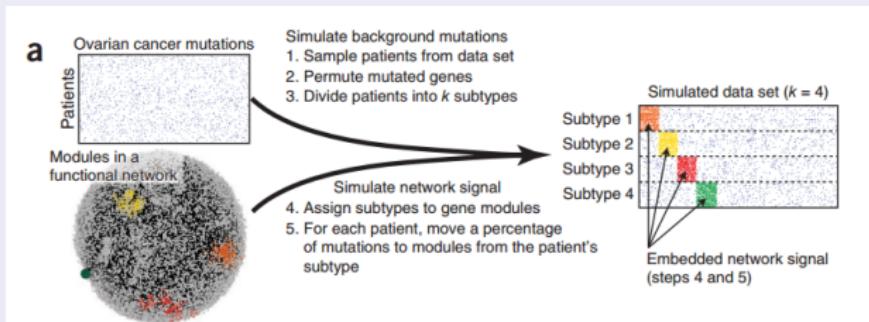
## 4 Discussion

# Simulations

determine the ability of NBS to recover subtypes from somatic mutation profiles

## How to generate 'ground truth' ?

### ① generate background mutation matrix (no subtype signal)



### ② add network-based signal

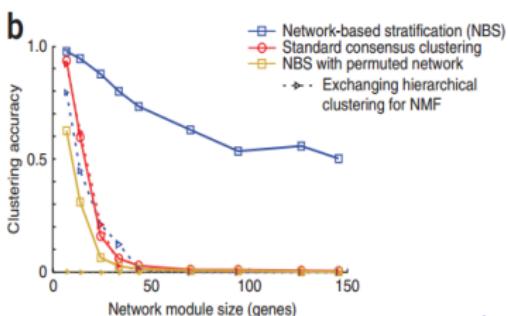
- establish a set of network communities
- assign a small number of network modules to each subtype → 'driver' subnetwork

### ③ reassign a fraction of the patient's mutations to genes covered by the driver modules for that patient's subtype

## A reasonable model of a pathway-based genetic disease

- ① driven by genetic circuits corresponding to a molecular network
- ② altered by mutations at multiple genes
- ③ many additional mutations that are noncausal 'passengers'

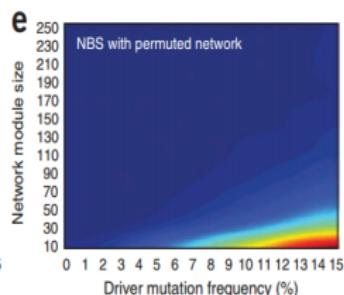
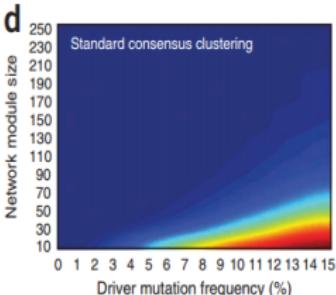
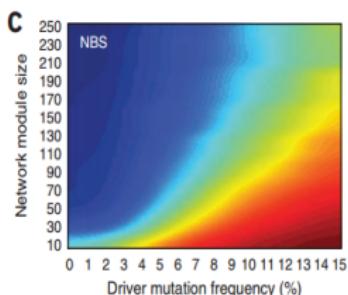
# Simulation Results



network-based  
Non-negative Matrix Factorization

large module size & low mutation frequency

no subtypes (permuted network)



| which NBS clusters recover simulated subtype assignments, evaluated with and without

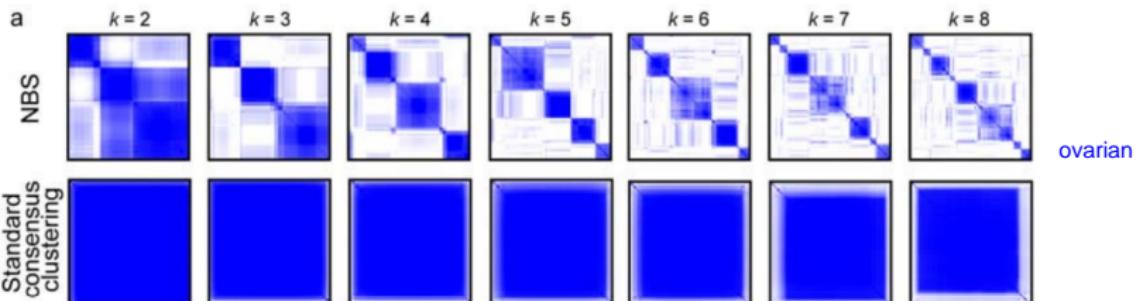
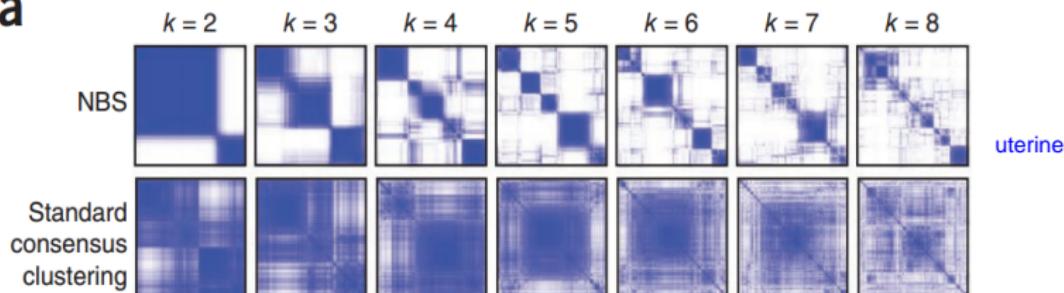
# Performance Analysis

## Overview

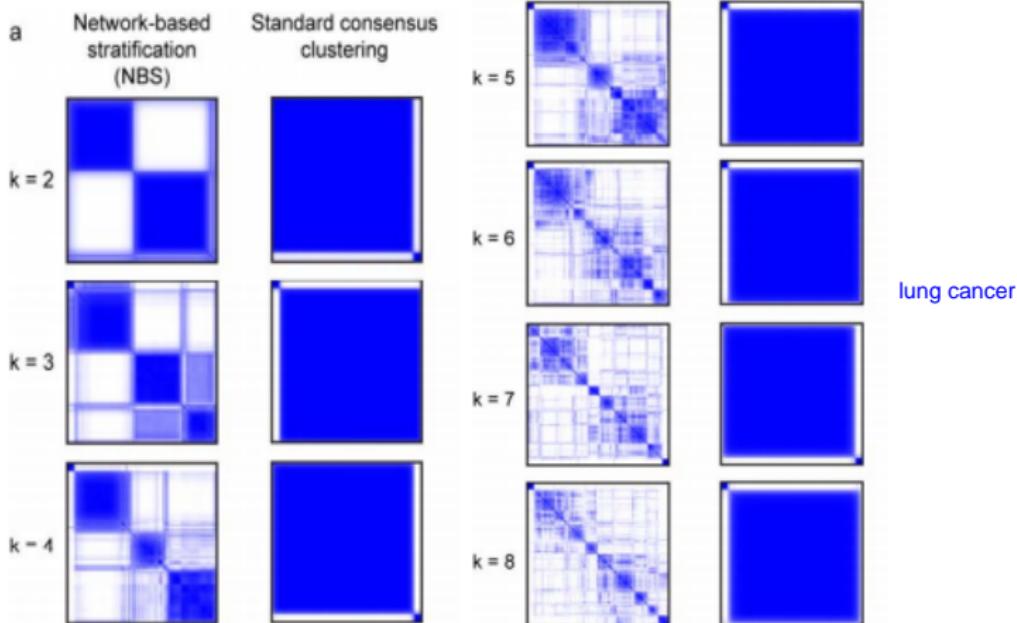
- ① result of clustering
- ② clinical and biological meaning
- ③ prediction

## Diverse Cancer types

uterine ; ovarian; lung cancer (co-clustering matrix)

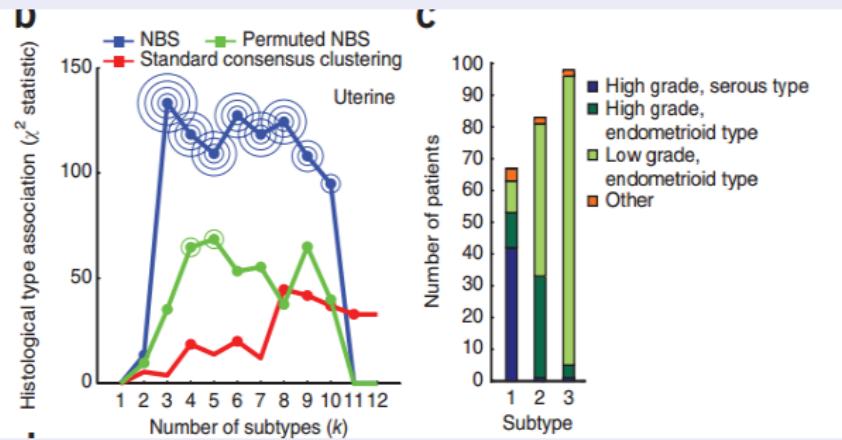
**a**

# Diverse Cancer types



# Biological Importance of Identified Subtypes

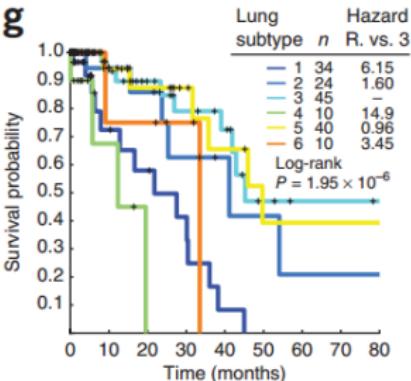
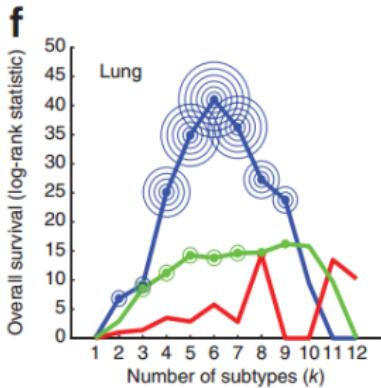
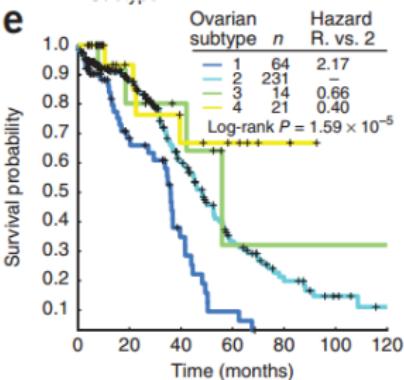
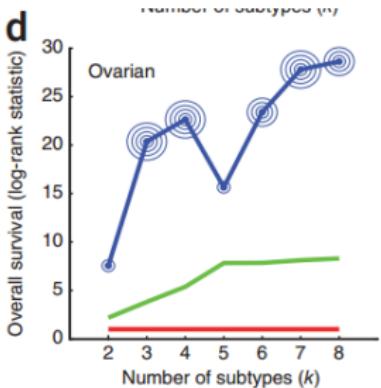
## uterine: histological basis



p-value of significance (concentric circles surrounding a data point)

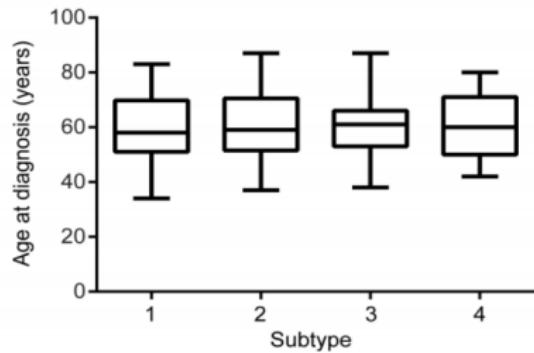
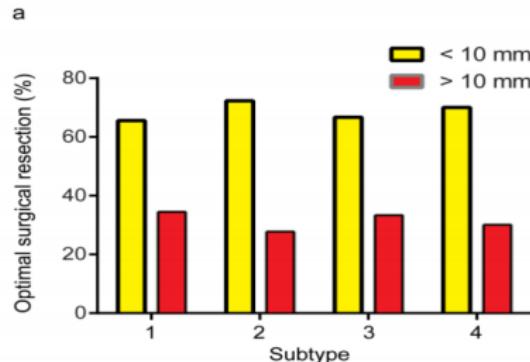
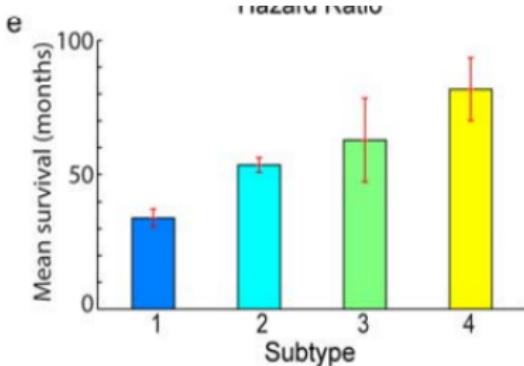
chi-square statistic: determine if a distribution of observed frequencies differs from the theoretical expected frequencies

# Survival Analysis

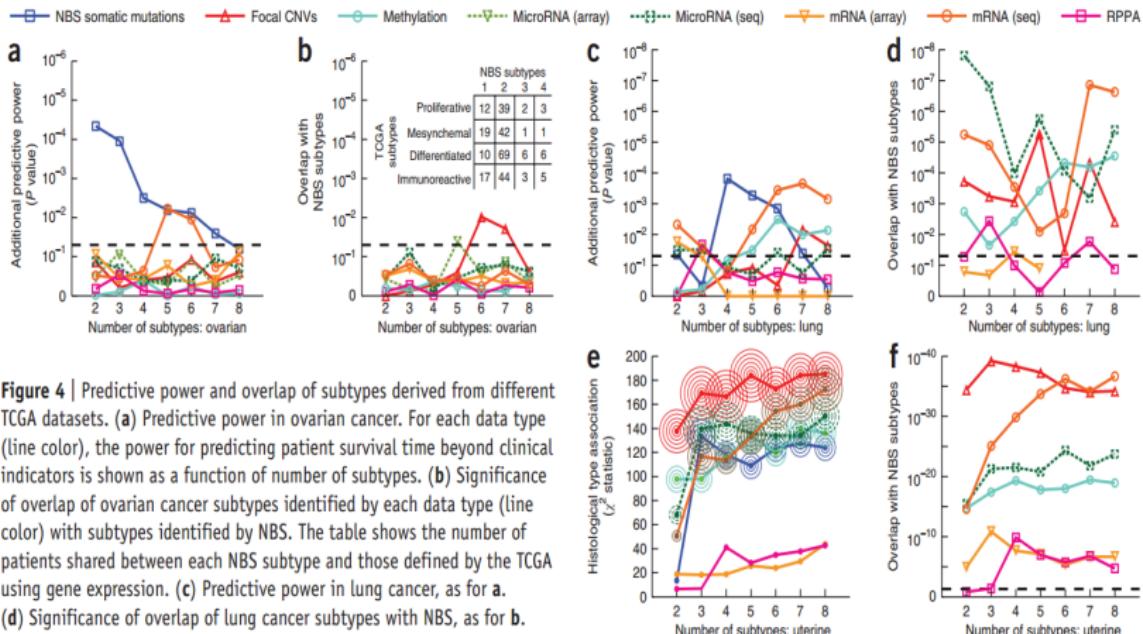


log-rank statistic: When are two KM curves statistically equivalent?

## Independent of other clinical marker



# Derived from other data type

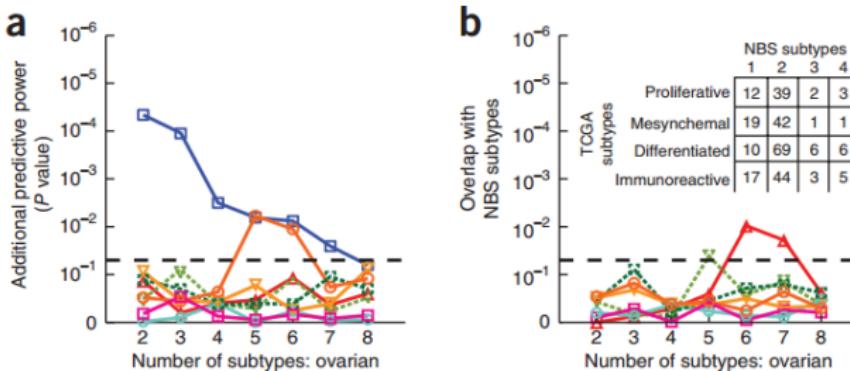


**Figure 4 |** Predictive power and overlap of subtypes derived from different TCGA datasets. **(a)** Predictive power in ovarian cancer. For each data type (line color), the power for predicting patient survival time beyond clinical indicators is shown as a function of number of subtypes. **(b)** Significance of overlap of ovarian cancer subtypes identified by each data type (line color) with subtypes identified by NBS. The table shows the number of patients shared between each NBS subtype and those defined by the TCGA using gene expression. **(c)** Predictive power in lung cancer, as for **a**. **(d)** Significance of overlap of lung cancer subtypes with NBS, as for **b**. **(e)** Association between uterine cancer subtype and tumor histology (y axis) as a function of the number of subtypes.  $P$  value of significance is indicated by concentric circles as in **Figure 3**. Colors are as in other panels, symbols have been omitted for clarity. **(f)** Significance of overlap of uterine cancer subtypes with NBS, as for **b**. Dashed horizontal lines indicate the  $P = 0.05$  threshold of significance.

## Compare with Other Data Type

CNV, methylation,mRNA expression,miRNA expression and protein profiles

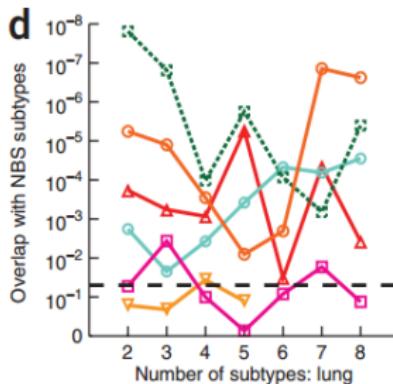
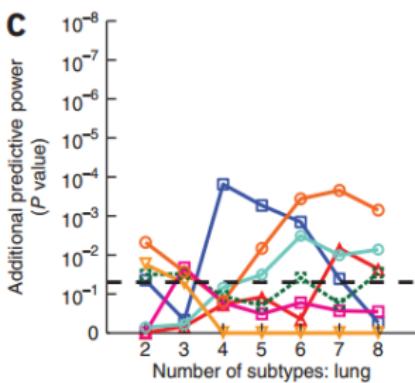
- ① **ovarian cancer:** all other data types had inferior ability to predict survival beyond what would be predicted from clinical covariate and led to different subtype assignments



## Compare with Other Data Type

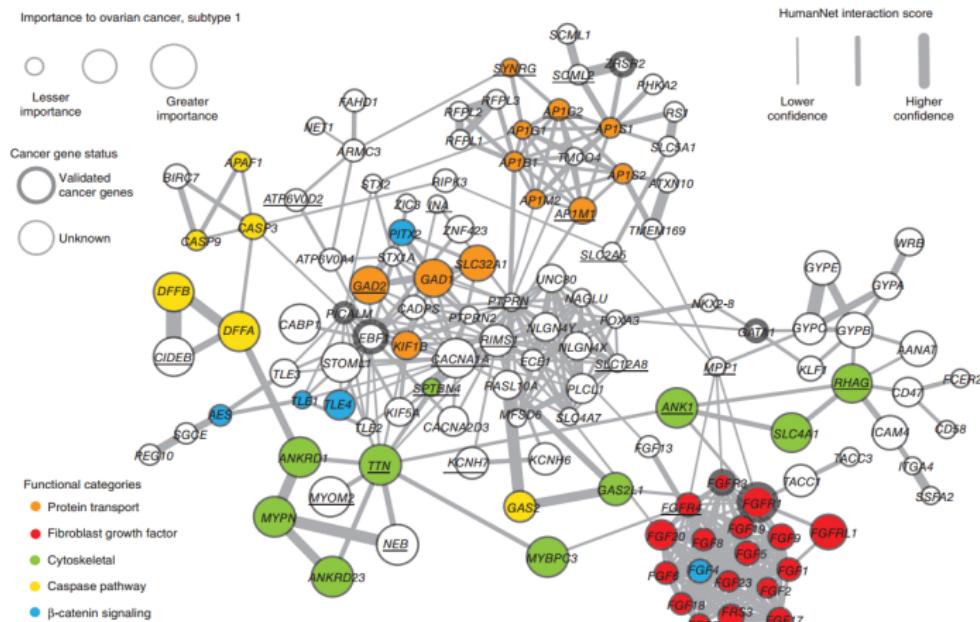
CNV, methylation,mRNA expression,miRNA expression and protein profiles

- ① **lung cancer**: both NBS subtypes and those based-on RNA-seq had good predictive power and had some overlaps



## Distinct Network Modules

Identify the regions of the network that are most responsible for discriminating the somatic mutation profiles of tumors of different subtypes



## Wilcoxon rank-based test

# Distinct Network Modules

## Ovarian Subtype 1

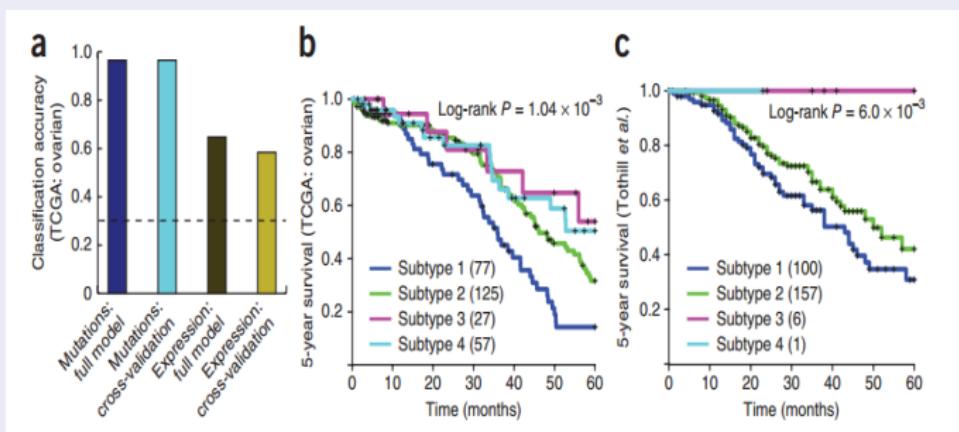
- ① worst overall survival time and shortest platinum-free interval
- ② contain over 20 genes in the FGF signaling pathway which has been previously implicated as a driver of tumor progression and associated with resistance to platinum

NBS not only can stratify patients into clinically informative subtypes but may help identify the molecular network regions commonly mutated in each subtype

# Predictive Signatures

- Applicable to new patients not in the TCGA ?
- assign a patient to one of the existing NBS subtypes

## Nearest shrunken centroid approach



## Nearest shrunken centroid approach

- Compute a standardized centroid for each class. This is the average gene expression for each gene in each class divided by the within-class standard deviation for that gene.

genes :  $i = 1, 2, \dots, p$  samples :  $j = 1, 2, \dots, n$  class :  $k = 1, 2, \dots, C$

$$\bar{x}_{ik} = \sum_{j \in C_k} \frac{x_{ij}}{n_k} \quad \bar{x}_i = \sum_{j=1}^n \frac{x_{ij}}{n}$$

$$d_{ik} = \frac{\bar{x}_{ik} - \bar{x}_i}{s_i} \quad s_i^2 = \frac{1}{n-K} \sum_k \sum_{j \in C_k} (x_{ij} - \bar{x}_{ik})^2$$

- shrinkage each  $d_{ik}$  towards 0 by soft-thresholding



$$d'_{ik} = \text{sign}(d_{ik})(|d_{ik}| - \Delta)_+$$

- shrunken centroid

$$\bar{x}'_{ik} = \bar{x}_i + s_i d'_{ik}$$

- Nearest centroid classification takes the gene expression profile of a new sample, and compares it to each of these class centroids. The class whose centroid that it is closest to, in squared distance, is the predicted class for that new sample.

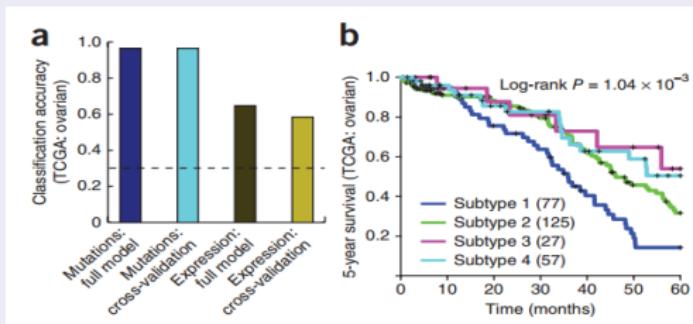
### Why shrinkage ?

- ① it can make the classifier more accurate by reducing the effect of noisy genes
- ② it does automatic gene selection. In particular, if a gene is shrunk to zero for all classes, then it is eliminated from the prediction rule.

Tibshirani, R., Hastie, T., Narasimhan, B. & Chu, G. Diagnosis of multiple cancer types by shrunken centroids of gene expression. Proc. Natl. Acad. Sci. USA 99, 6567-6572 (2002).

## How about use mRNA expression data to assign a new patient?

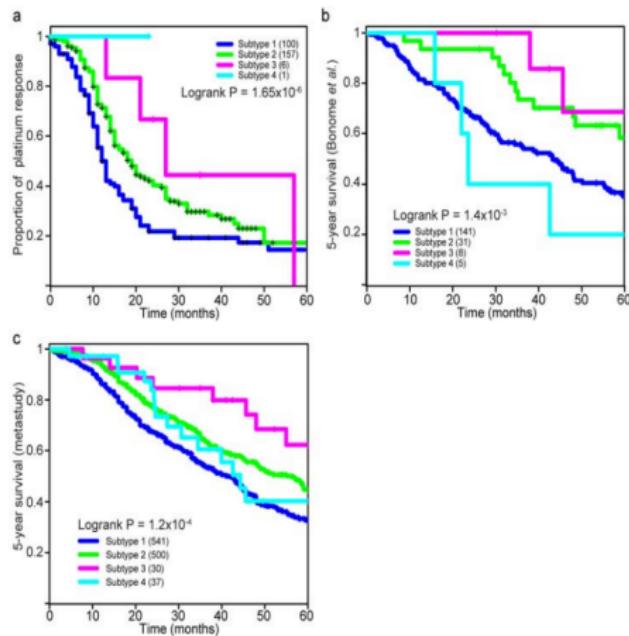
- mRNA expression data are presently much more widely available than are full genome or exome sequences



- classifier defined by NBS
- shrunken centroids: mRNA expression
- patients assigned by mRNA expression

# Gene Expression as a Surrogate Biomarker

TCGA and two independent studies (include gene expression profiles but lack somatic mutation profiles)



# Content

## ① Introduction

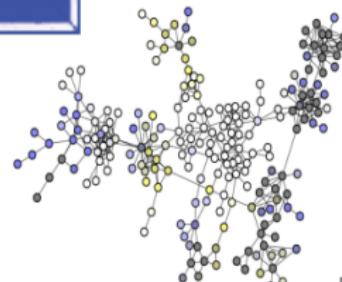
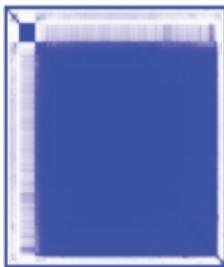
## ② Method

## ③ Performance Analysis

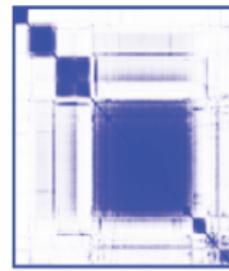
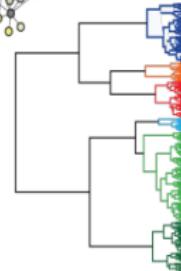
## ④ Discussion

# Summary

Regular Consensus Clustering NMF



Network based stratification



# Conclusion

- ① Network-based stratification recovers biologically relevant subtypes
- ② Somatic mutation subtypes are different from those recovered from other molecular profiles.
- ③ These subtypes can be recapitulated using gene expression.
- ④ Each subtype seems to have specific effected subnetworks.

# Discussion

## Somatic Mutation

- ① differential measurements between tumor and normal tissues, whereas expression profiles are absolute measurements  
"base-line" : germ line ; leave only tumor-specific
- ② Causal genetic events underlying tumor progression

## Method

- ① Use network as prior knowledge
- ② NMF clustering

## A comprehensive research

Simulation analysis;clinically and biologically meaningful;predictive signature

# Discussion

How to perform effective pattern recognition from massive biological data?

**① Generative model:**

- mostly used;
- data reduction;
- low dimension

**Assume latent variables**

**② Why discriminative model?**

- sample space;
- separability;
- unsupervised;
- feature selection;**
- class label learning**