



# Molecular features based cancer classification and stratification

Dongfang Wang

Tsinghua University

2016-09-02



# Cancer heterogeneity and stratification

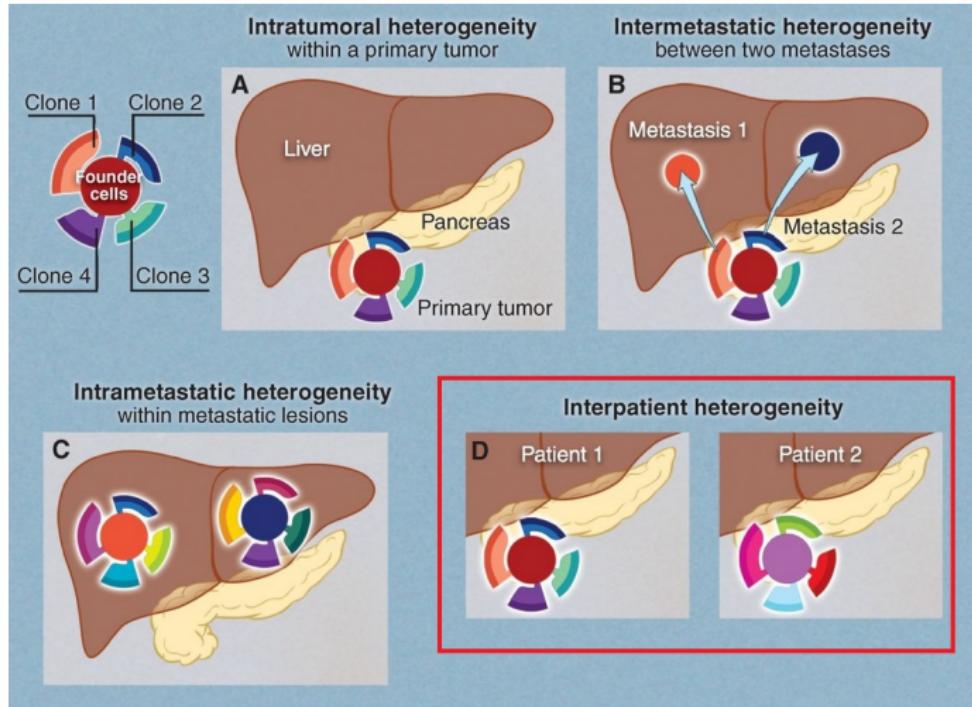
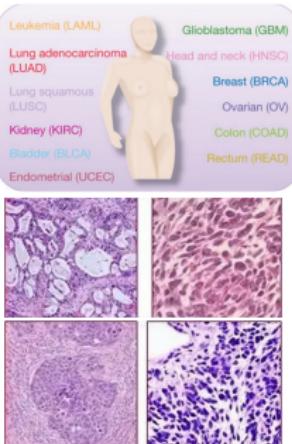


Figure: Four types of genetic heterogeneity in tumors. [Vogelstein et al., Science, 2013](#)



# Cancer stratification:

- Better undersatnding of tumor biology
- Subtype specific drug and therapy
- Patient tailored treatment
- ...



1.0 Primary tumor apparatus

2.0 Clinical Pathology

Fig: four types of lung cancer

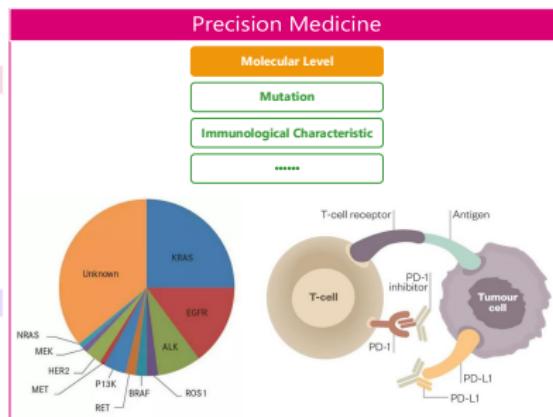


Figure: Development of cancer stratification. Figs from Internet.



# Molecular features based stratification

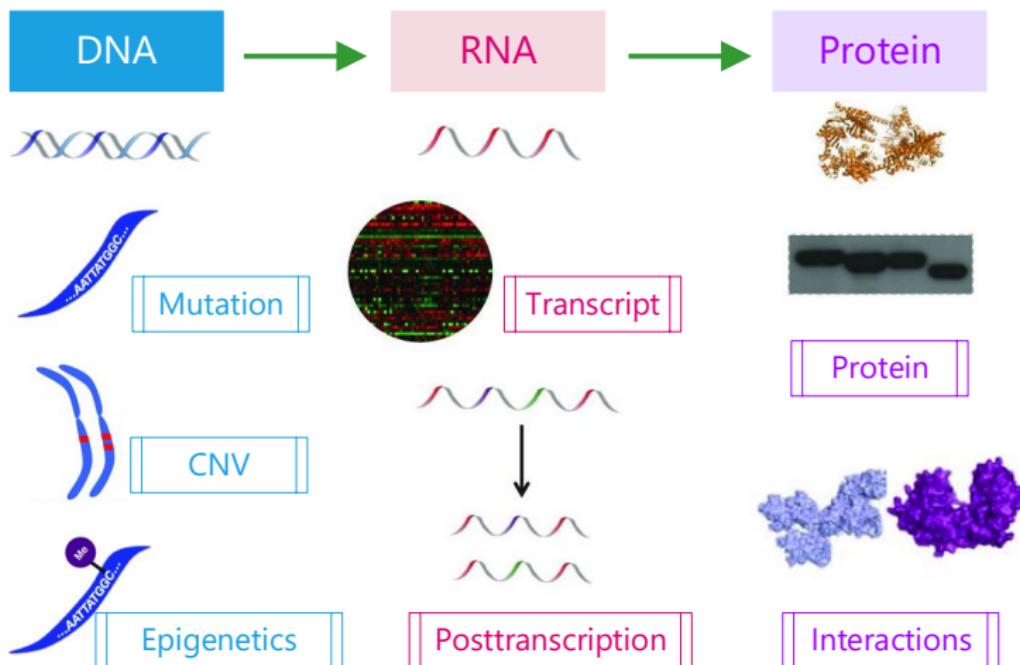


Figure: Various molecular features available for tumor patients. Figs from Internet.



## Different molecular levels indicates various kinds of information for subtyping:

- Different driver mutations ([Ciriello et al., \*Nat.Genetics\*, 2013, etc.](#))
- CpG island hypermethylation defines CIMP(CpG island methylator phenotype) subtypes ([Issa et al., \*Nat.Rev.Cancer\*, 2004, etc.](#))
- ER/HER2/PR define breast cancer subtypes ([Cancer Genome Atlas, \*Nature\*, 2012, etc.](#) )
- ...



# Computational Concerns

## ① Collections of available datasets:

### 1. Public available cancer datasets

- Multi-omics cancer samples: TCGA, ICGC ( $\sim 10^2 - 10^3$ )
- Gene Expression data: GEO, ArrayExpress ( $\sim 10^3 - 10^5$ )
- More omics data in cancer cell lines: CCLE
- Biological knowledge: COSMIC, MSigDB, DAVID

### 2. Own collected datasets: specific and accurate information ( $\sim 10^2$ )

## ② Characteristics and difficulties:

- High-dimensional features with limited samples
- Multi-source information
- Heterogenous datasets



# Content

- ① Cancer subtypes across multiple datasets
- ② Cancer subtypes indicated by multi-omics
- ③ Gene mutation and Sorafenib drug response



# Cancer subtypes across multiple datasets

- More samples ( $\sim 2000$  HCCs) to support possible discovery.

| Group                          | mRNA type   | miRNA type | Total       |
|--------------------------------|-------------|------------|-------------|
| HCC                            | 2031        | 625        | 2109        |
| Adjacent                       | 1199        | 224        | 1277        |
| Cirrhotic                      | 40          | 0          | 40          |
| Healthy                        | 8           | 10         | 18          |
| <b>Total number of samples</b> | <b>3286</b> | <b>859</b> | <b>3444</b> |

Table I: Number of samples in curated HCC datasets

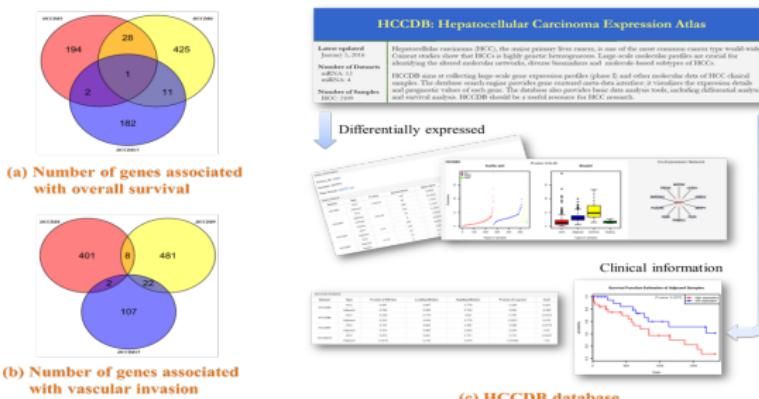


Figure: HCC Database made by Zhang Guchao



# Summary of current methods

- **Strategy 1:**

Stack all datasets into one big matrix after "removing" batch effects.

(J Cquinney et al., *Nature Medicine*, 2015; Johnson WE et al., *Biostatistics*, 2007; etc. )

- **Strategy 2:**

Compare "similarity" of clusters from different datasets.

(Hoshida Y, et al. *Cancer Research*, 2009; CR Planey, et al., *Genome Medicine*, 2016. )

**Similarity computation is based on signatures of each dataset.**



- **Key questions:** How to represent shared information across different datasets/studies/cohorts?
- **Hypothesis:** driver signature modules and their degrees of perturbation are shared among datasets.

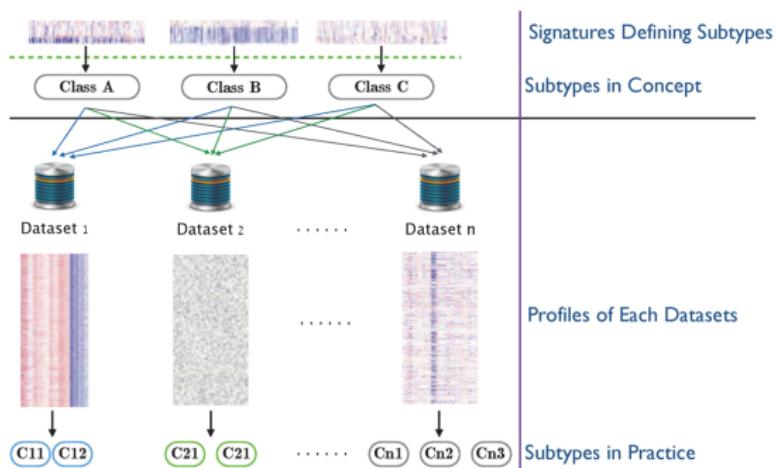


Figure: Basic framework of our model.



## - Definition of signature modules:

- ① Firstly, we divide every dataset into appropriate clusters  $c_{ij}$ , and obtain representative gene lists  $g_{ij}$ , where  $i$  indexes dataset and  $j$  cluster.
- ② Then integrating with literature genes and PPI network, we define signature modules  $M_k$ .

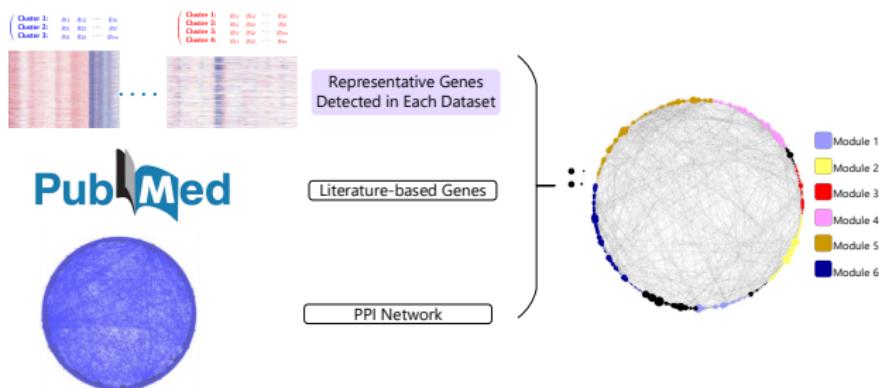


Figure: Discovery of signature modules across multiple datasets.



## - Functional annotation of signature modules:

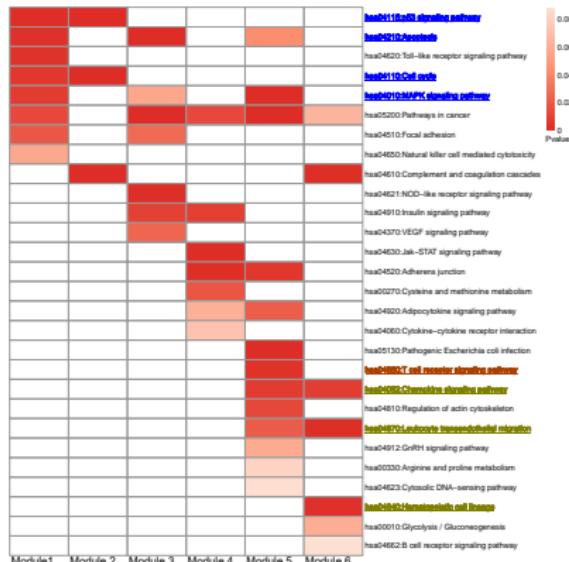


Figure: Function annotation of signature modules by DAVID.(Huang DW, et al., *Nature Protoc.*, 2009)



## - Perturbation of signature modules in each cluster $c_{ij}$ :

- ▶ Compute median expression as the template of every cluster in each dataset.
- ▶ Compute z-score of each template's enrichment in all 6 modules.



Figure: Perturbation vectors of every cluster in each dataset.



- Consensus perturbation vectors across datasets:

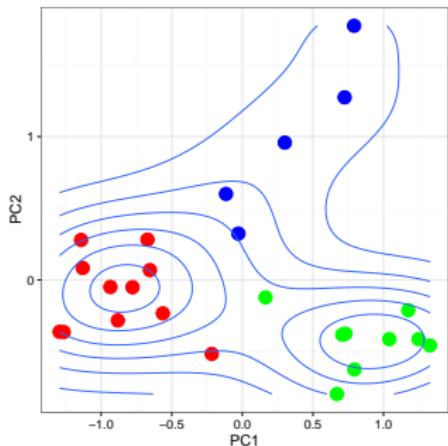


Figure: PCs of perturbation vectors for each cluster.

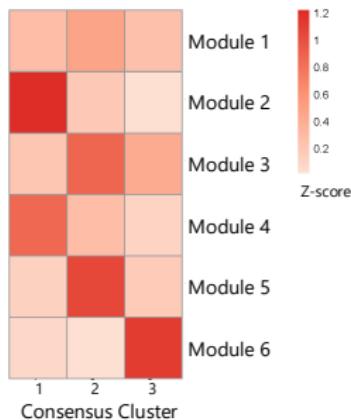


Figure: Heatmap for three consensus perturbation centers.



- Survival analysis results compared with clustering by single dataset:

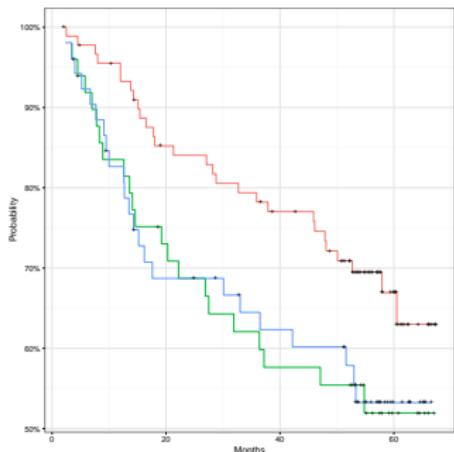


Figure: Survival analysis by single clustering of HCCDB6. cox: 0.05

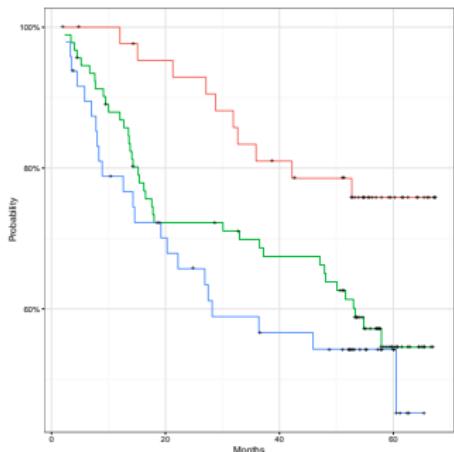


Figure: Survival analysis by consensus clustering of HCCDB6. cox: 0.008



# Summary

- Try to perform collective clustering across different datasets from one type of cancer,
  - Based on their shared driver signature modules.
  - Preliminary results support our hypothesis.
- 
- The current model is still very simple.
  - Further improvements and some details need to be re-considered.



# Cancer subtypes indicated by multi-omics

An example:

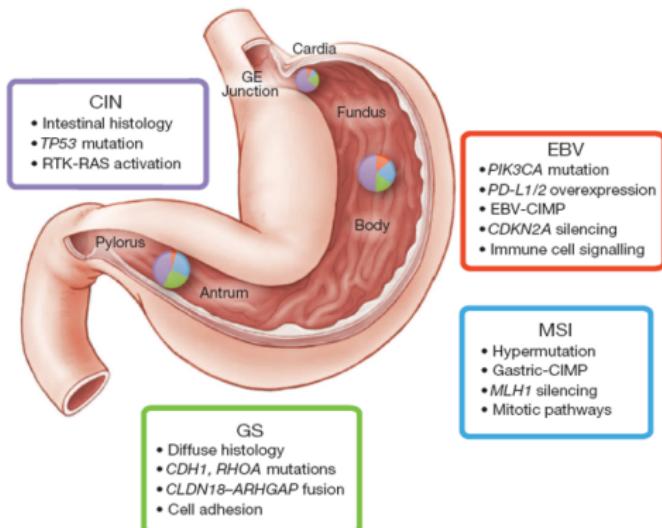


Figure: Subtypes of gastric adenocarcinoma, (TCGA, *Nature*, 2014)



# Summary of current methods

Table 1. Selected methods for unsupervised clustering of multi-omics data.

| Strategy                          | Description  | Methods   | Basic model  | Refs.  |
|-----------------------------------|--|---|--|--|
| Direct integrative clustering     | Put multi-omics datasets into a stacked matrix, and take this stacked matrix as input for the following clustering analysis                | Super $k$ -means<br>iCluster+<br>JIVE<br><b>LRAcluster</b><br>jNMF<br>Pathifier | Direct clustering with BIC<br>Latent factor analysis<br>Low rank-based approximation<br>Low rank-based approximation<br>Non-negative matrix factorization<br>Pathway-based integration | [12]<br>[13]<br>[14]<br>[15]<br>[16]<br>[17] |
| Clustering of clusters            | Perform clustering analysis on every single omics dataset and then integrate the primary clustering results into final cluster assignments | COCA<br>MDI<br>BCC<br>SNF   | Clustering of intermediated clusters<br>Latent <i>Dirichlet</i> allocation<br>Latent <i>Dirichlet</i> allocation<br>Similarity network fusion  | [3]<br>[18]<br>[19]<br>[20]                  |
| Regulatory integrative clustering | Focus on driver variations by considering the regulatory structures between different molecular layers                                     | PARADIGM  | Pathway-based integration  | [21]   |

Figure: Summary of current integrative clustering methods. [Wang, et al., Quantitative Biology, 2016.](#)



# LRAcluster: overview

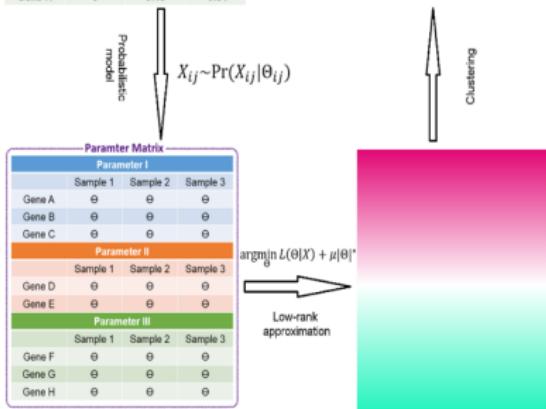
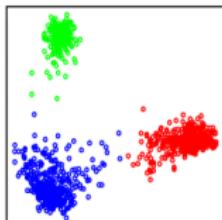
| Somatic Mutation |          |          |          |
|------------------|----------|----------|----------|
|                  | Sample 1 | Sample 2 | Sample 3 |
| Gene A           | 1        | 0        | 0        |
| Gene B           | 1        | 0        | 1        |
| Gene C           | 0        | 1        | 1        |

| RNA sequencing |          |          |          |
|----------------|----------|----------|----------|
|                | Sample 1 | Sample 2 | Sample 3 |
| Gene D         | 103      | 96       | 132      |
| Gene E         | 27       | 42       | 35       |

| Copy Number Variation |          |          |          |
|-----------------------|----------|----------|----------|
|                       | Sample 1 | Sample 2 | Sample 3 |
| Gene F                | -0.87    | 1.02     | 0        |
| Gene G                | -0.87    | 0        | -0.34    |
| Gene H                | 0        | 0.45     | -0.34    |



$$\min_{\Theta} -\log \mathcal{P}(X|\Theta) + \mu|\Theta|_*$$

- $X$ : omics data matrix;
- $\Theta$ : parameter matrix;
- $\mathcal{P}$ : probability distribution, Gaussian for real number, Binomial for binary, etc.
- $||_*$ : nuclear norm.

**Note:** cooperate with Wu Dingming (Wu et al., *BMC Genomics*, 2016). I helped design this algorithm and did experiments.



# LRAcluster: results

Simulation: mixture of samples from three kinds of cancers

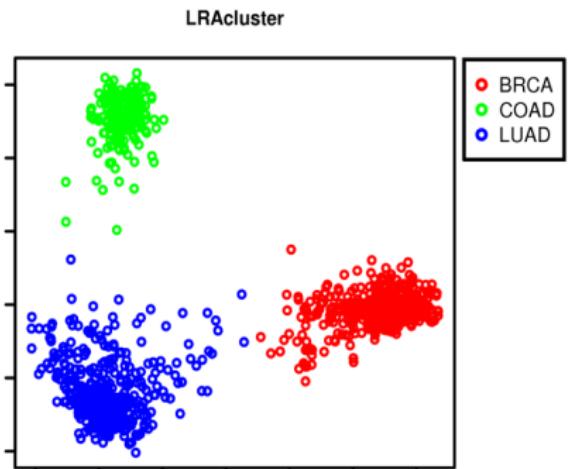


Figure: Clustering results of LRAcluster in simulation datasets.

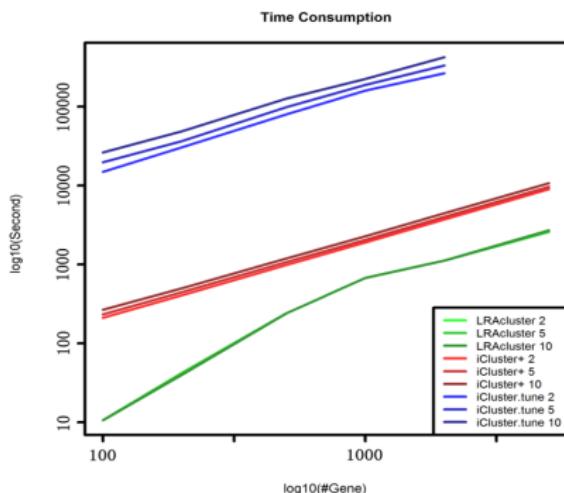


Figure: Time consumption of LRAcluster compared with similar methods.  $x$ -axis: number of dimensions.



# LRAcluster: summary

- A unified probabilistic framework for analyzing high dimensional multi-omics data.
- A fast tool to analyze multi-omics data and visualize them in a low dimensional space.
- Applied to pan-cancer and single cancer datasets.

Table 1. The results of pan-cancer analysis.<sup>1)</sup>

|                         | BRCA <sup>2)</sup>    | COAD <sup>2)</sup>    | GBM <sup>2)</sup>      | HNSC <sup>2)</sup>      | KIRC <sup>2)</sup>    | LGG <sup>2)</sup>       | LUAD <sup>2)</sup>    | LUSC <sup>2)</sup>     | PRAD <sup>2)</sup>    | STAD <sup>2)</sup>    | THCA <sup>2)</sup>    | Total <sup>2)</sup>     |
|-------------------------|-----------------------|-----------------------|------------------------|-------------------------|-----------------------|-------------------------|-----------------------|------------------------|-----------------------|-----------------------|-----------------------|-------------------------|
| C1 <sup>2)</sup>        | 1 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 286 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 6 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 293 <sup>2)</sup>       |
| C2 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 1 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 411 <sup>2)</sup>     | 412 <sup>2)</sup>       |
| <b>C3<sup>2)</sup></b>  | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>41<sup>2)</sup></b> | <b>0<sup>2)</sup></b>   | <b>0<sup>2)</sup></b> | <b>451<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b>  | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>492<sup>2)</sup></b> |
| C4 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>       | 231 <sup>2)</sup>     | 0 <sup>2)</sup>       | 231 <sup>2)</sup>       |
| C5 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 293 <sup>2)</sup>     | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 293 <sup>2)</sup>       |
| C6 <sup>2)</sup>        | 0 <sup>2)</sup>       | 190 <sup>2)</sup>     | 0 <sup>2)</sup>        | 1 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>         | 2 <sup>2)</sup>       | 0 <sup>2)</sup>        | 1 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 194 <sup>2)</sup>       |
| C7 <sup>2)</sup>        | 3 <sup>2)</sup>       | 17 <sup>2)</sup>      | 0 <sup>2)</sup>        | 0 <sup>2)</sup>         | 1 <sup>2)</sup>       | 0 <sup>2)</sup>         | 406 <sup>2)</sup>     | 7 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 3 <sup>2)</sup>       | 437 <sup>2)</sup>       |
| C8 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>         | 240 <sup>2)</sup>     | 0 <sup>2)</sup>         | 0 <sup>2)</sup>       | 0 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 240 <sup>2)</sup>       |
| C9 <sup>2)</sup>        | 448 <sup>2)</sup>     | 0 <sup>2)</sup>       | 1 <sup>2)</sup>        | 2 <sup>2)</sup>         | 1 <sup>2)</sup>       | 0 <sup>2)</sup>         | 4 <sup>2)</sup>       | 1 <sup>2)</sup>        | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 0 <sup>2)</sup>       | 457 <sup>2)</sup>       |
| <b>C10<sup>2)</sup></b> | <b>8<sup>2)</sup></b> | <b>1<sup>2)</sup></b> | <b>0<sup>2)</sup></b>  | <b>195<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b>   | <b>6<sup>2)</sup></b> | <b>60<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>0<sup>2)</sup></b> | <b>270<sup>2)</sup></b> |
| Total <sup>2)</sup>     | 460 <sup>2)</sup>     | 208 <sup>2)</sup>     | 42 <sup>2)</sup>       | 484 <sup>2)</sup>       | 242 <sup>2)</sup>     | 452 <sup>2)</sup>       | 418 <sup>2)</sup>     | 74 <sup>2)</sup>       | 294 <sup>2)</sup>     | 231 <sup>2)</sup>     | 414 <sup>2)</sup>     | 3319 <sup>2)</sup>      |

Figure: Apply LRAcluster on pan-cancer analysis.



# Gene mutation and Sorafenib drug response

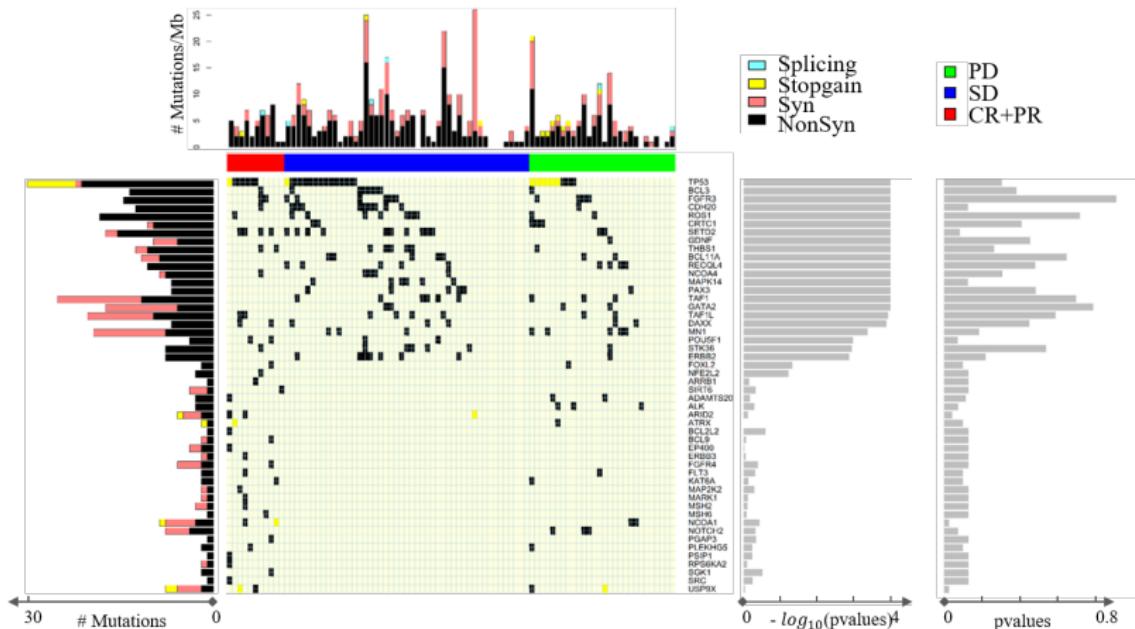


Figure: Overview of sorafenib reponse associated mutations.



# Conclusion

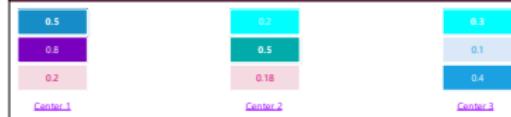
- Molecular features-based cancer stratification and classification have profound impacts on understanding cancer biology and cancer therapy.
- To overcome the computational difficulties, such as high-dimension, limited samples, multi-source information, we've built some models, including:
  - Cancer subtypes across multiple datasets
  - Cancer subtypes indicated by multi-omics
- It's still arguable about how to validate cancer stratification results.



## **Supplement**



### Consensus Clustering Signatures



Clustering of perturbation vectors

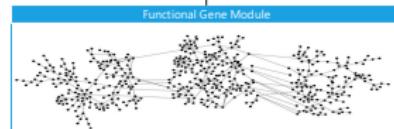
Obtain representative vector of each cluster

Perturbation vector: the extent of perturbation of every module in each cluster

| Dataset 1 |           |           |
|-----------|-----------|-----------|
| Cluster 1 | Cluster 2 | Cluster 3 |
| 0.4       | 0.3       | 0.3       |
| 0.8       | 0.5       | 0.1       |
| 0.2       | 0.18      | 0.2       |

| Dataset 2 |           |           |
|-----------|-----------|-----------|
| Cluster 1 | Cluster 2 | Cluster 3 |
| 0.05      | 0.9       |           |
| 0.6       | 0.07      |           |
| 0.2       | 0.13      |           |

| Dataset 3 |           |           |           |
|-----------|-----------|-----------|-----------|
| Cluster 1 | Cluster 2 | Cluster 3 | Cluster 4 |
| 0.43      | 0.07      | 0.95      | 0.34      |
| 0.87      | 0.69      | 0.9       | 0.11      |
| 0.25      | 0.24      | 0.1       | 0.27      |

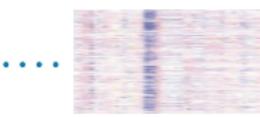


Clustering and signature detection in each dataset

(Cluster 1:  $g_{11} \ g_{12} \ \dots \ g_{1n}$   
Cluster 2:  $g_{21} \ g_{22} \ \dots \ g_{2n}$   
Cluster 3:  $g_{31} \ g_{32} \ \dots \ g_{3n}$ )

(Cluster 1:  $g_{11} \ g_{12} \ \dots \ g_{1k}$   
Cluster 2:  $g_{21} \ g_{22} \ \dots \ g_{2l}$ )

(Cluster 1:  $g_{11} \ g_{12} \ \dots \ g_{1k}$   
Cluster 2:  $g_{21} \ g_{22} \ \dots \ g_{2l}$   
Cluster 3:  $g_{31} \ g_{32} \ \dots \ g_{3m}$   
Cluster 4:  $g_{41} \ g_{42} \ \dots \ g_{4n}$ )





# Clustering of single dataset

## Outliers detection

- PCA: use the first two PCs, and compute the distance from the origin.
- KS test: check whether the distribution of one sample's gene expression profile was significantly different from the overall distribution.

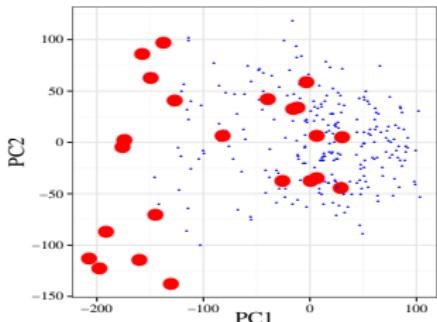


Figure: PCA of HCCDB3 expression profile with outliers labeled.



## Review of NMF methods

$$X \approx WH \text{ s.t. } W, H \geq 0$$

- Learn **parts** of an object.
- Applied widely in biological subtyping since (Brunet et al, *PNAS*, 2004).
- The obtained coefficients matrix could be used for clustering:

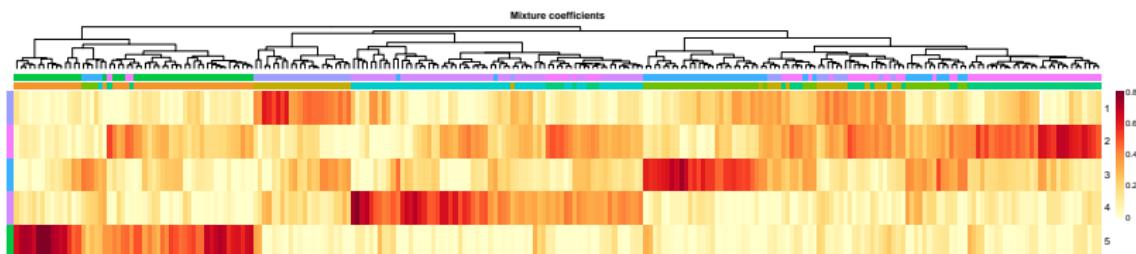


Figure: An example of NMF coef matrix. (HCCDB3)



- **NMF consensus clustering:** bootstrap samples, perform clustering separately and construct a consensus matrix. (See next page for an example)
- Use **cophenetic correlation** to determine the factorization rank. (how faithfully a dendrogram preserves the pairwise distances between the original unmodeled data points)

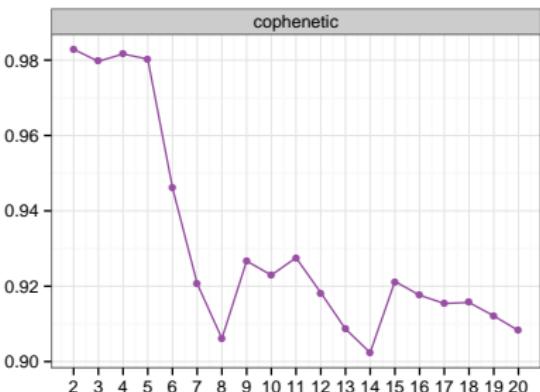


Figure: Cophenetic correlation for different  $K$ . (HCCDB3)

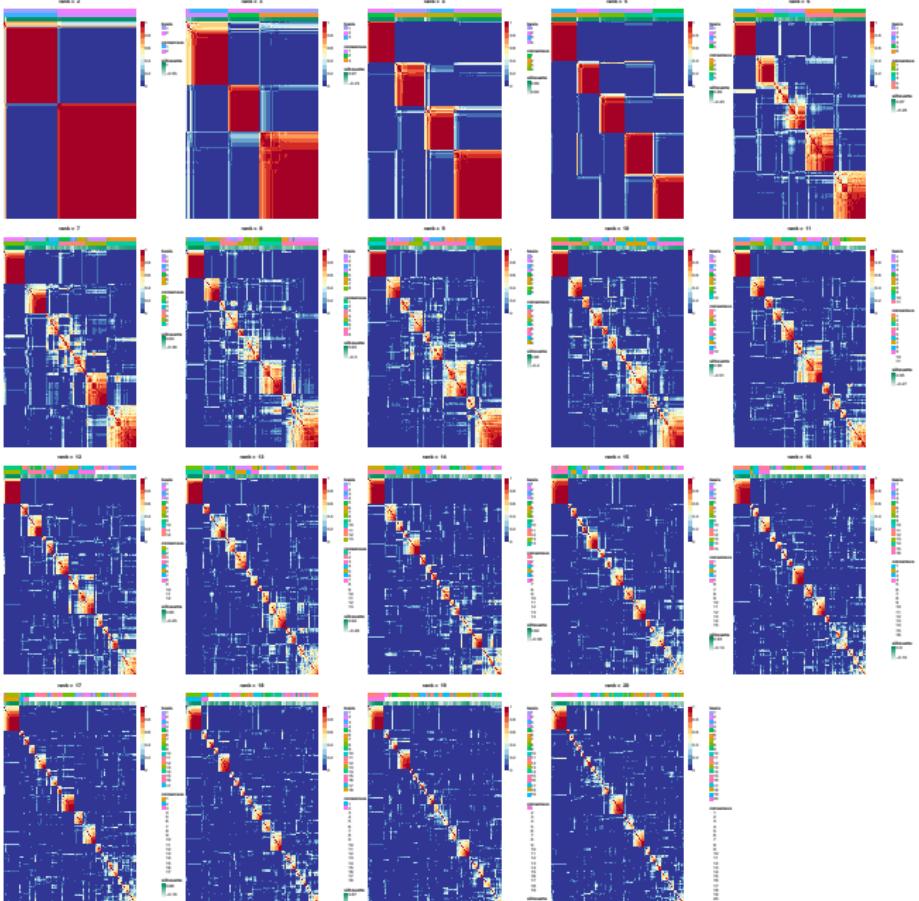


Figure: An example of NMF consensus matrix for different rank  $K$  (HCCDB3)



# Discovery of signature modules

- ① For each cluster in every single dataset, define its signature gene list by:
  - Differentially expressed (`t.test` and `wilcoxon` test,  $fdr < 1e - 5$  )
  - Higher mean expression in this cluster than others.
- ② Map literature genes on PPI to extract PPI sub-network.
- ③ For every gene at the PPI sub-network, count its number of appearance in signature gene lists of step 1 as its initial value.
- ④ Perform random walk with restart on sub-network of step 2 with initial values defined at step 2.



# Definition of perturbation vectors

- Now we have **6 signature modules** (6 signature gene lists), and for each cluster in every single dataset, we could define its representative "sample" vector  $g$  which is median of all samples in this cluster.
- Compute enrichment score  $s$  of  $g$  in every signature module.
- Then permute the gene lists, and generate a null distribution for  $s$ .
- Compute the  $z$ -score  $z$  for  $s$ , denoted by  $z_i, i = 1, \dots, 6$
- The dim 6 vector  $z$  is the perturbation vector for every cluster.



# Re-classification of single samples

- For every sample, we could also define a perturbation vector  $z$ .
- Then we compute the distance from it to all perturbation vectors defined for this dataset.



# Survival analysis of HCCDB15

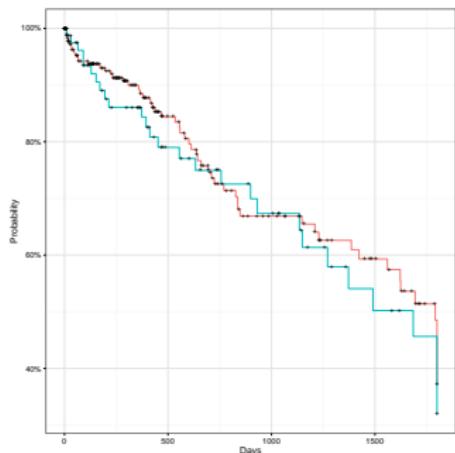


Figure: Survival analysis by single clustering of HCCDB15. cox: 0.56

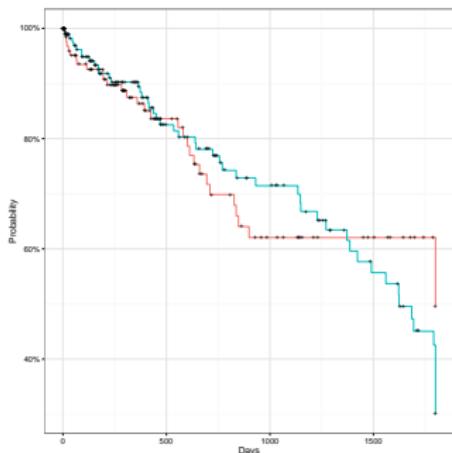


Figure: Survival analysis by consensus clustering of HCCDB15. cox: 0.62



## 基于低秩逼近的方法



### ➤ 迭代求解框架

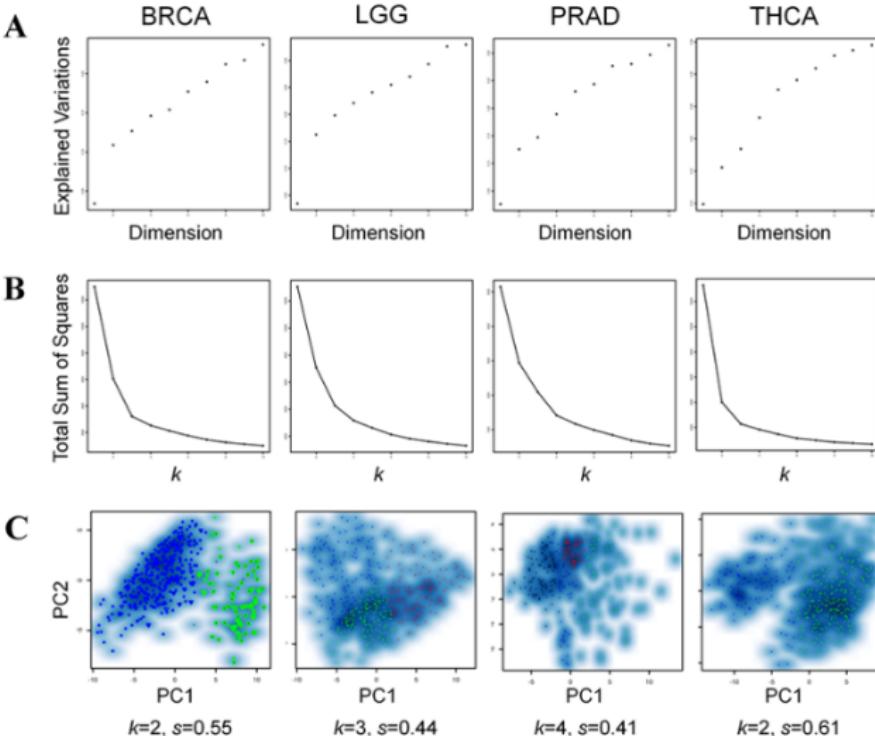
$$(1) \theta^{(2n+1)} = \theta^{(2n)} + \delta_n \partial \ln(\Pr(X; \theta^{(2n)}))$$

$$(2) \theta^{(2n+2)} = D_\mu(\theta^{(2n+1)})$$

$$D_\mu(X) = U \begin{bmatrix} \{s_1 - \mu\}^+ & 0 & \cdots & 0 \\ 0 & \{s_2 - \mu\}^+ & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \{s_R - \mu\}^+ \end{bmatrix} V^T$$

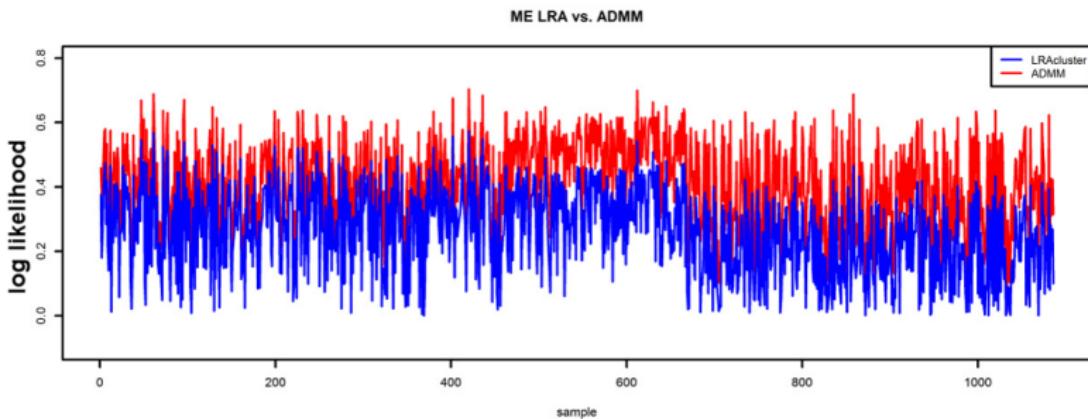


# More LRAcluster results



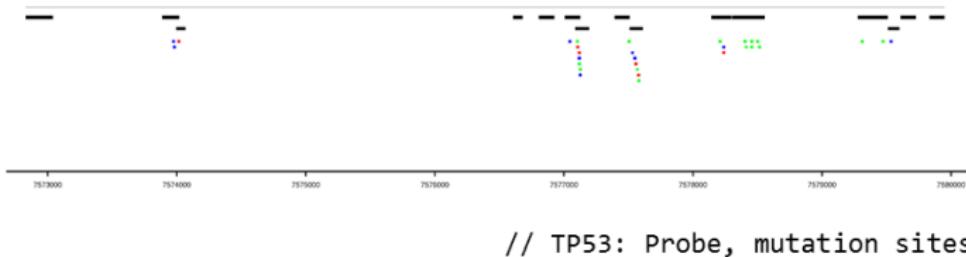


# More LRAcluster results





## Introduction



- Targeted exome sequencing 440 designed genes of advanced HCC patients (recurrent in 2 years) with Sorafenib treatment.
- 101 patients with PR+CR, SD, PD respectively
- We want to detect significant signatures that can predict the drug response and then validate it.