



## Intro to R for Data Analysis and Visualization Summer 2022

# Connecting Patient Datasets to Preclinical Hypotheses in Cancer

Riyue Sunny Bao PhD  
Associate Professor of Medicine  
Co-Director of Cancer Bioinformatics (CBS)  
Co-Director of Translational Immuno-Informatics Lab (TIIL)

July 8<sup>th</sup>, 2022

# Outline

## Schedule

- 8:30 AM – 9:30 AM : Lecture
- 9:30 AM – 10:00 AM : Hands-on practice

## Key take home messages

- How to associate gene expression data with clinical outcome
- Hands on: Use (public) gene expression data to discover tumor subtypes and survival analysis

# Class materials

Download class materials from GitHub

- MSTP\_summer2022.lecture.pdf
- MSTP\_summer2022.handsOn.Rmd

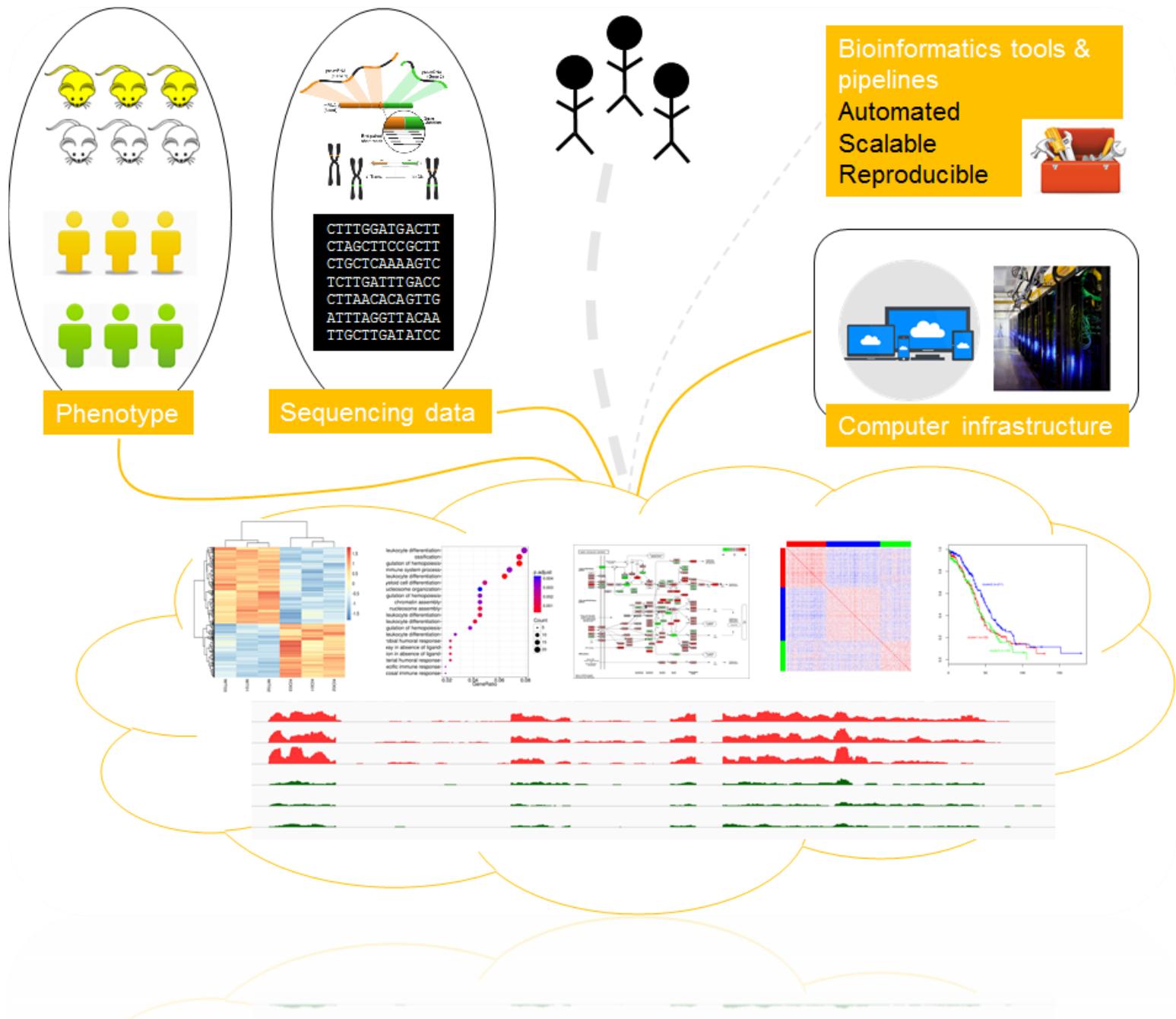
[https://github.com/riyuebao/MSTP\\_summer2022\\_class07](https://github.com/riyuebao/MSTP_summer2022_class07)

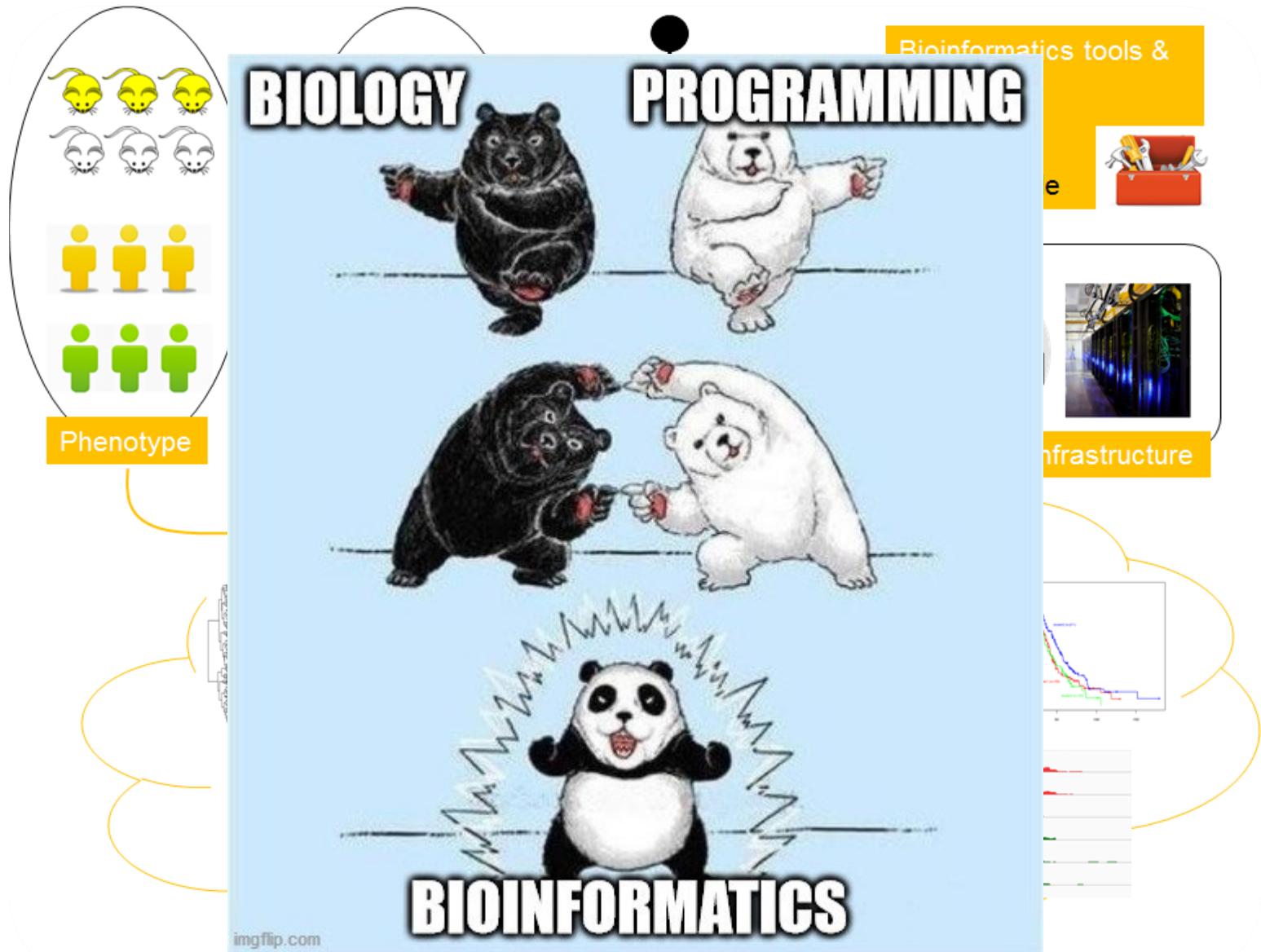
What you need

- Rstudio (or R console) on personal computers (hands on practice)

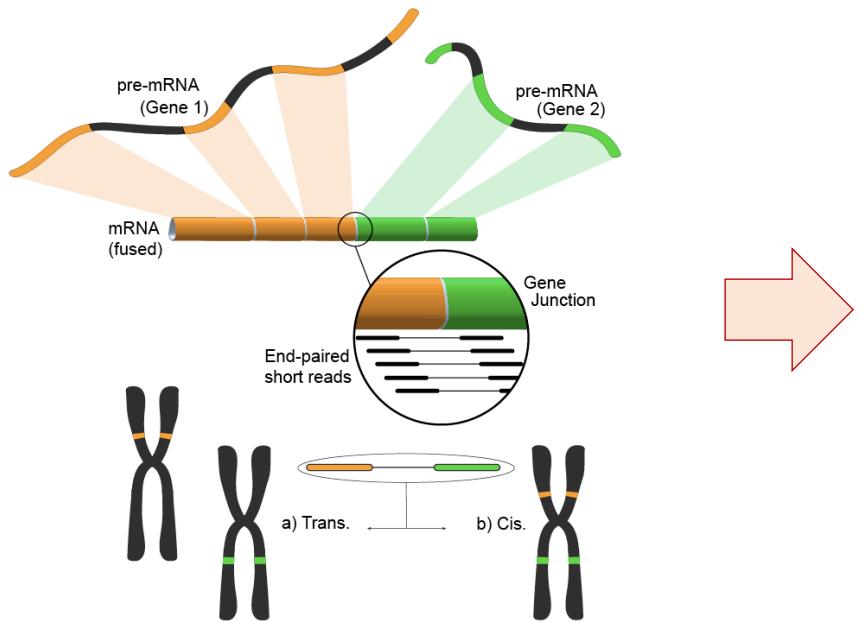
# Objective

- Learn the background and application of The Cancer Genome Atlas (TCGA)
- Learn the structure and access of Genomics Data Commons (GDC)
- Explore datasets hosted on GDC
- Practice how to connect gene expression with clinical data
  - Use gene expression to identify tumor subtype
  - Detect survival difference between subtypes
  - Produce high-quality plots for publication

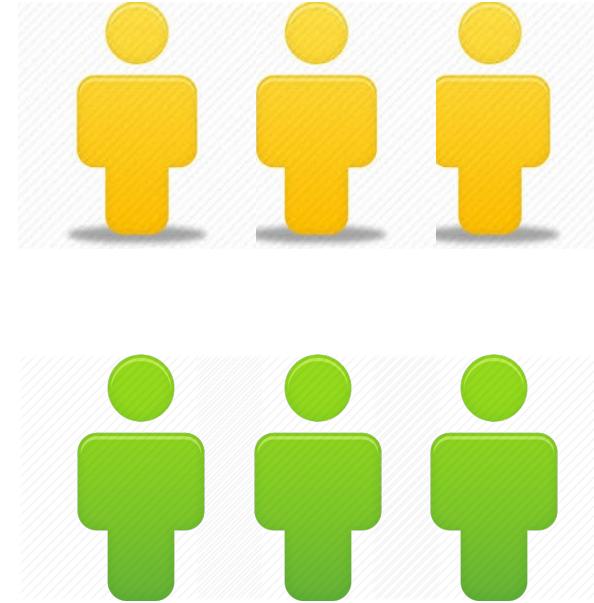




# Background

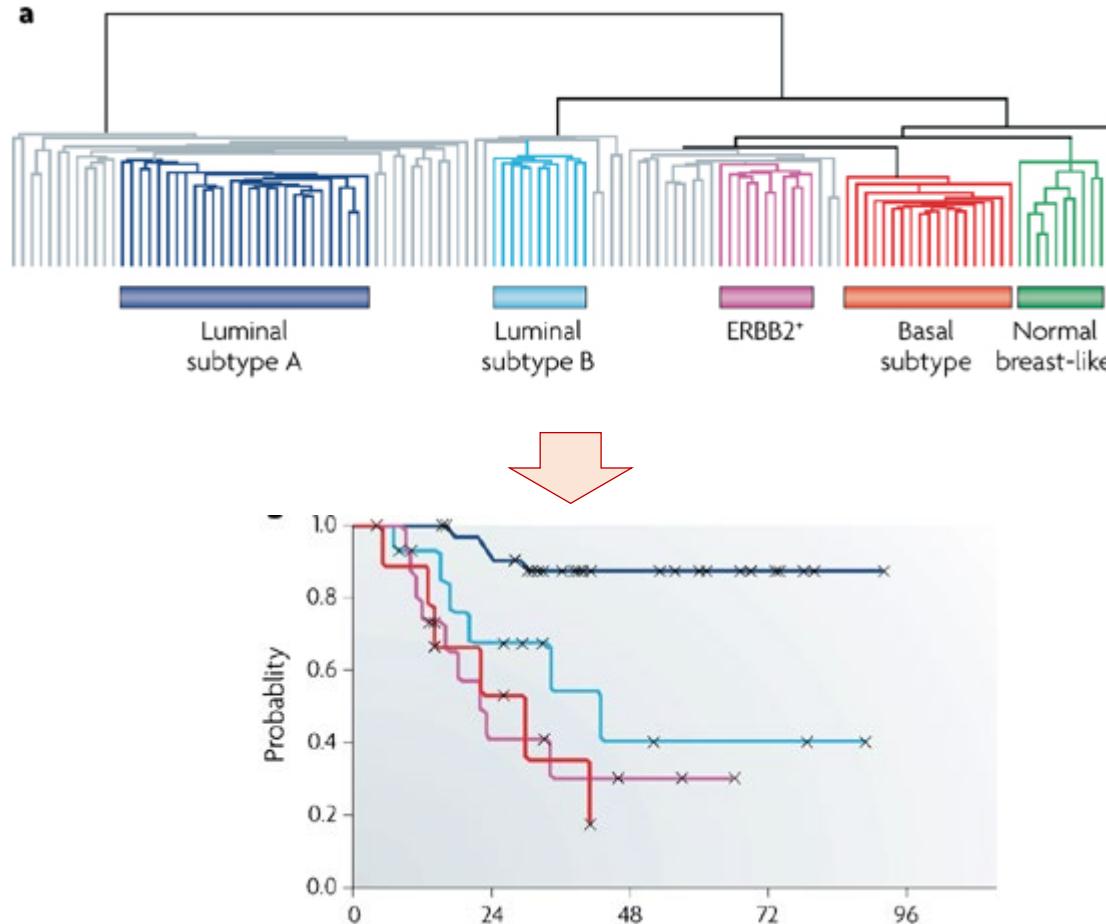


Gene Expression



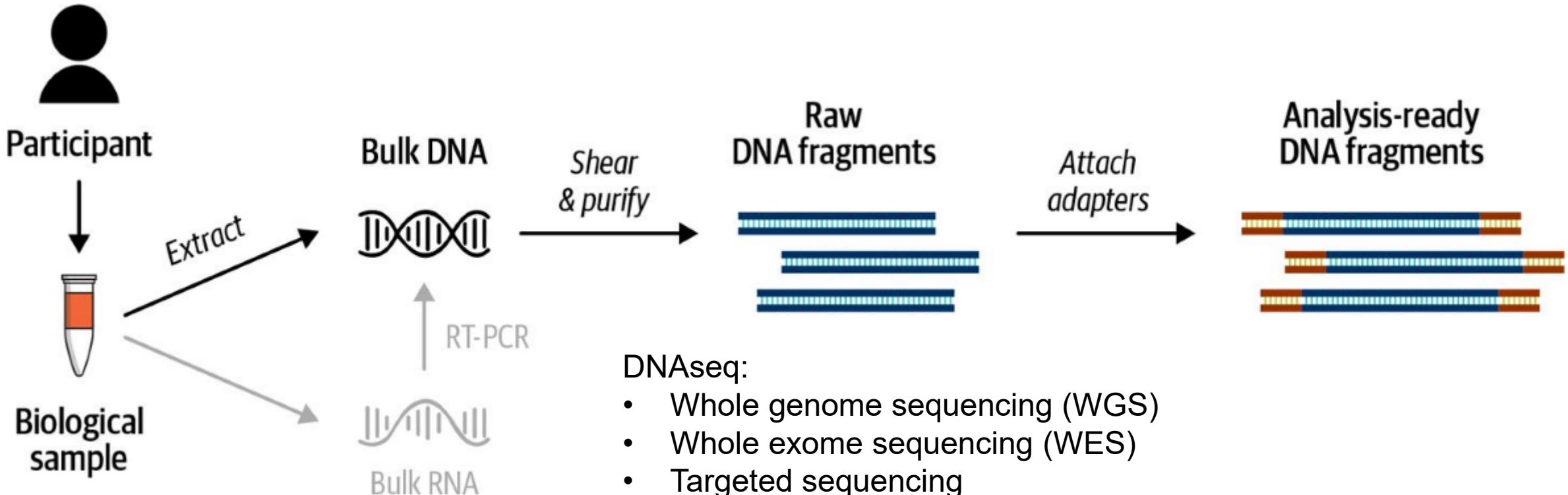
Patient's clinical data

# Background



- Correlate gene expression with clinical data (tumor stage, tumor grade, time to death, time to relapse, etc.)
- Identify tumor subtypes through sample clustering
- Detect survival difference between tumor subtypes
- Discover gene signatures to predict patient classes

# NGS experiments



## DNAseq:

- Whole genome sequencing (WGS)
- Whole exome sequencing (WES)
- Targeted sequencing

## RNAseq:

- Whole transcriptome sequencing (WTS)
- mRNA sequencing
- Small RNA sequencing (miRNA, snoRNA, etc.)

# Main types of public cancer omics databases

Sequencing based dbs

Protein based dbs

Imaging based dbs

Bulk tissue databases

Single cell databases

- 
- Mine data and form new hypothesis
  - Validate your existing discoveries using a much larger cohort
  - Identify novel targets for cancer therapy and bring it to clinic
  - Identify new patient subtypes and develop new biomarkers
  - Develop innovative computational tools
  - Meta-analysis across many cohorts
  - Perform machine learning and deep learning using cohorts of sufficient sample sizes

# Genomic Data Commons (GDC)

<https://portal.gdc.cancer.gov/>

## Harmonized Cancer Datasets Genomic Data Commons Data Portal

Get Started by Exploring:



e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

### Data Portal Summary

[Data Release 26.0 - September 08, 2020](#)

FILES

590,367

GENES

23,399

MUTATIONS

3,287,299



3.53 PB cancer genomics data!!  
Open to the scientific community!!  
Data download is very stable and fast!!

With the latest Gen3 technology

<https://stats.gen3.org/>

### Cases by Major Primary Site

#### Program

- TCGA
- TARGET
- GENIE
- BEATAML1.0
- CGCI
- CMI
- CPTAC
- CTSP
- FM
- HCMI
- MMRF
- NCICCR
- OHSU
- ORGANOID
- VAREPOP

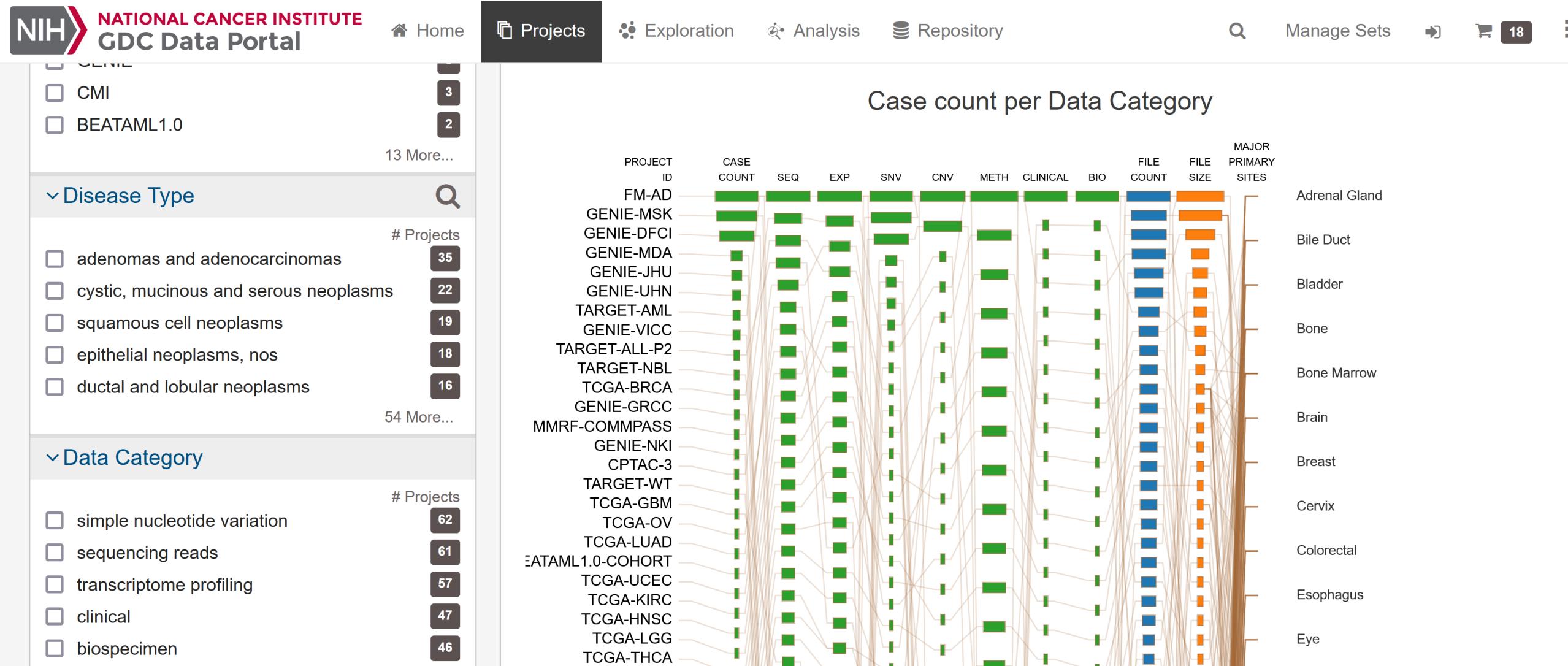
54 disease types

#### Disease Type

# Pro

- adenomas and adenocarcinomas
- cystic, mucinous and serous neoplasms
- squamous cell neoplasms
- epithelial neoplasms, nos
- ductal and lobular neoplasms
- complex epithelial neoplasms
- complex mixed and stromal neoplasms
- gliomas
- nevi and melanomas
- mesothelial neoplasms
- myeloid leukemias
- transitional cell papillomas and carcinomas
- acinar cell neoplasms
- adnexal and skin appendage neoplasms

# Genomic Data Commons (GDC)



# Genomic Data Commons (GDC): TCGA-LUAD



# Genomic Data Commons (GDC) : TCGA-LUAD

NATIONAL CANCER INSTITUTE  
GDC Data Portal

Home Projects Exploration Analysis Repository

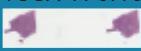
Manage Sets

Program Name IS TCGA AND Project Id IS TCGA-LUAD AND Data Type IS Annotated Somatic Mutation

Cases Slides Image

TCGA-44-6779 - TCGA-LUAD

TCGA-44-6779-01Z-00-DX1



TCGA-62-A472 - TCGA-LUAD

TCGA-44-6779-01A-01-TS1



TCGA-38-4628 - TCGA-LUAD

TCGA-69-7980 - TCGA-LUAD

TCGA-50-6595 - TCGA-LUAD

TCGA-78-7633 - TCGA-LUAD

TCGA-86-8055 - TCGA-LUAD

TCGA-55-A4DF - TCGA-LUAD

TCGA-49-6744 - TCGA-LUAD



## ARTICLE

OPEN

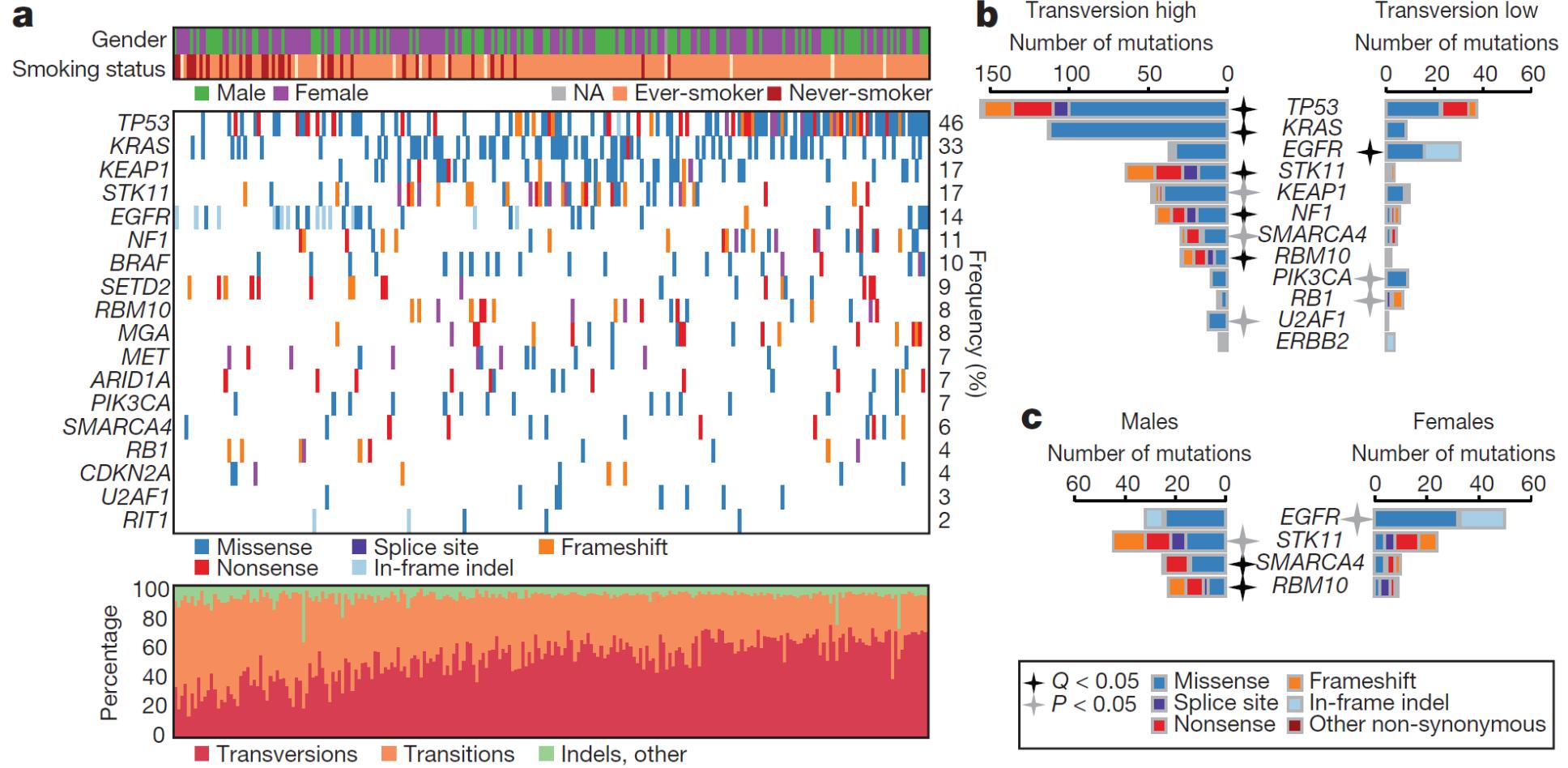
doi:10.1038/nature13385

# Comprehensive molecular profiling of lung adenocarcinoma

The Cancer Genome Atlas Research Network\*

Adenocarcinoma of the lung is the leading cause of cancer death worldwide. Here we report molecular profiling of 230 resected lung adenocarcinomas using messenger RNA, microRNA and DNA sequencing integrated with copy number, methylation and proteomic analyses. High rates of somatic mutation were seen (mean 8.9 mutations per megabase). Eighteen genes were statistically significantly mutated, including *RIT1* activating mutations and newly described loss-of-function MGA mutations which are mutually exclusive with focal MYC amplification. *EGFR* mutations were more frequent in female patients, whereas mutations in *RBM10* were more common in males. Aberrations in *NF1*, *MET*, *ERBB2* and *RIT1* occurred in 13% of cases and were enriched in samples otherwise lacking an activated oncogene, suggesting a driver role for these events in certain tumours. DNA and mRNA sequence from the same tumour highlighted splicing alterations driven by somatic genomic changes, including exon 14 skipping in *MET* mRNA in 4% of cases. MAPK and PI(3)K pathway activity, when measured at the protein level, was explained by known mutations in only a fraction of cases, suggesting additional, unexplained mechanisms of pathway activation. These data establish a foundation for classification and further investigations of lung adenocarcinoma molecular pathogenesis.

# Genomic Data Commons (GDC): TCGA-LUAD



# Genomic Data Commons (GDC): TCGA-OV

## ARTICLE

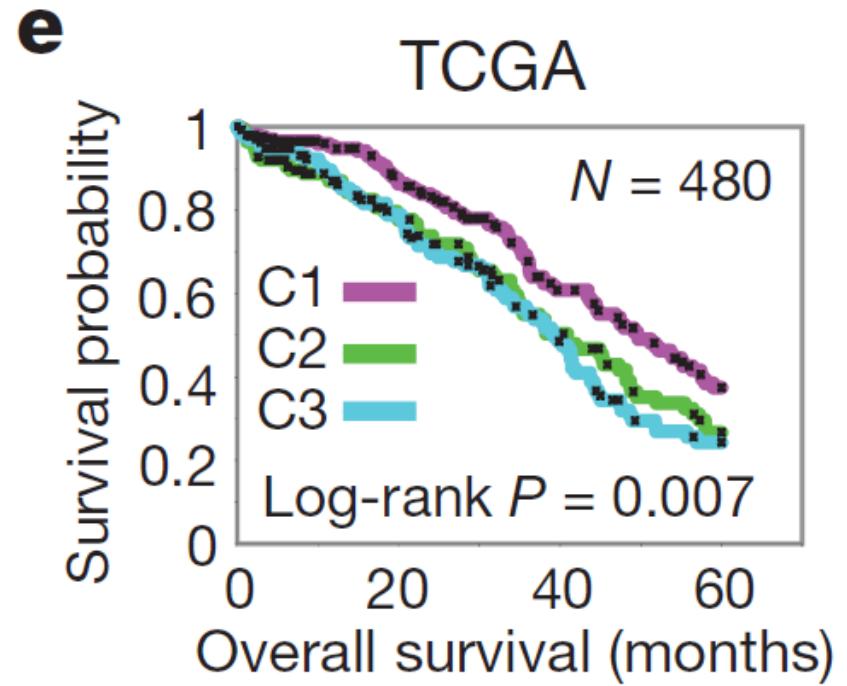
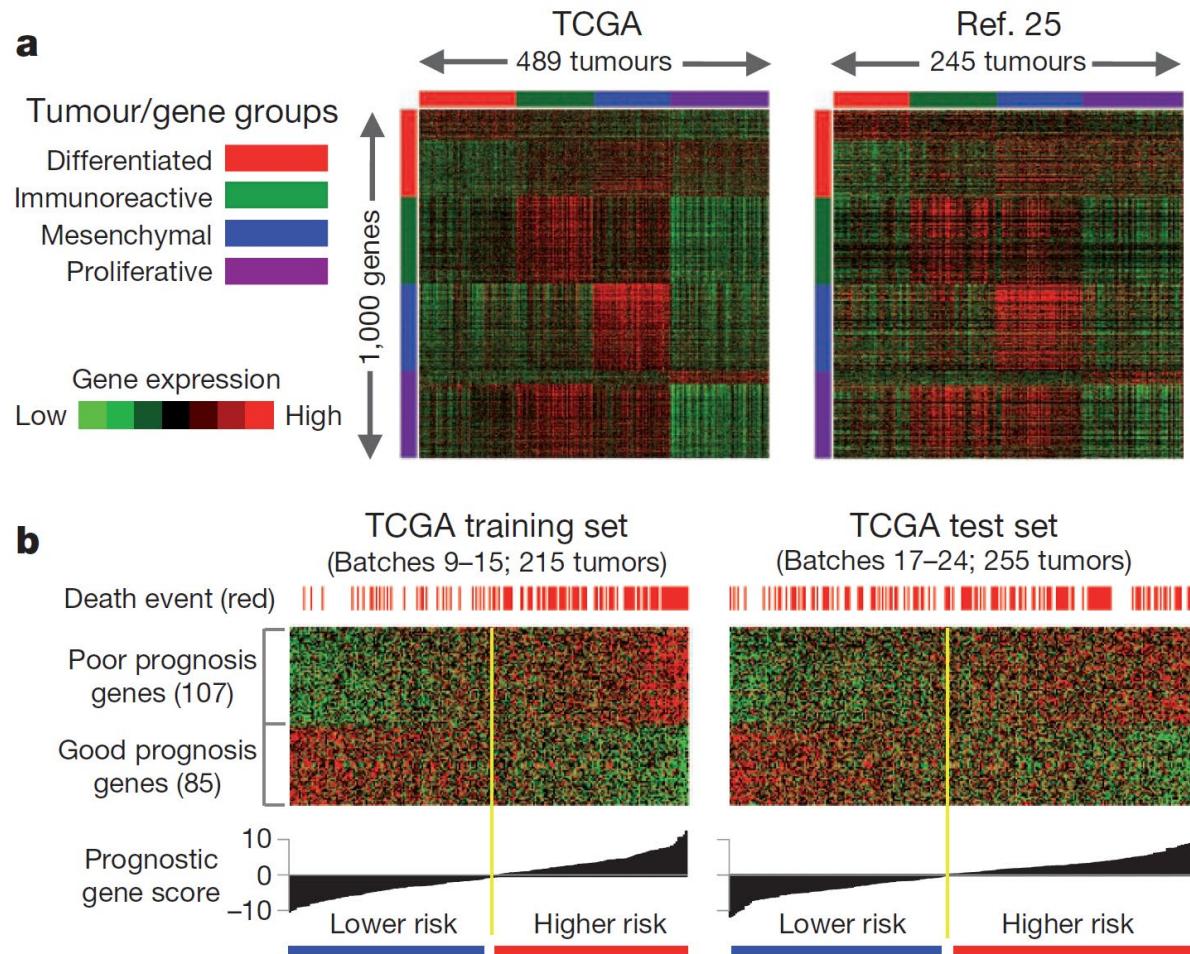
doi:10.1038/nature10166

# Integrated genomic analyses of ovarian carcinoma

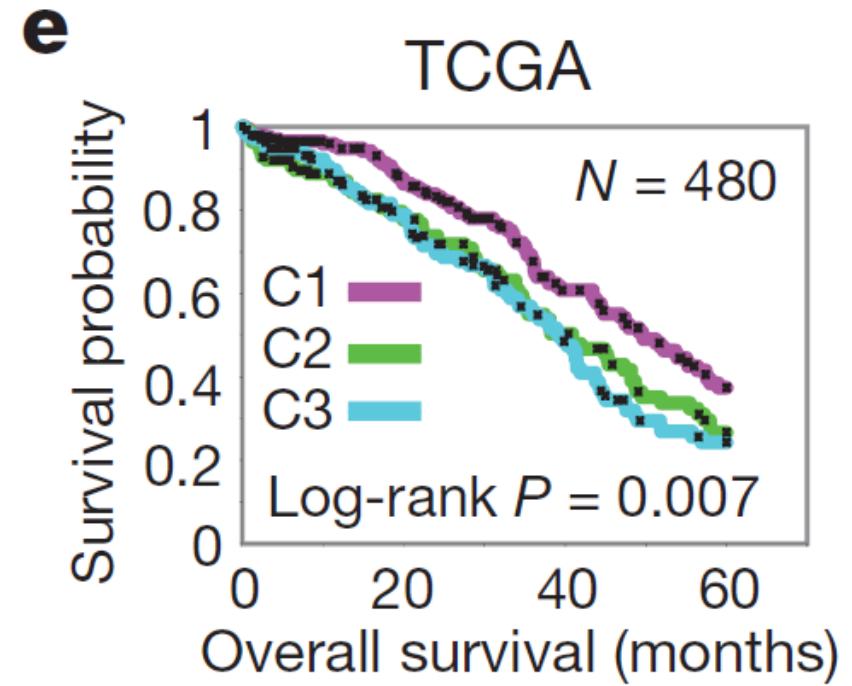
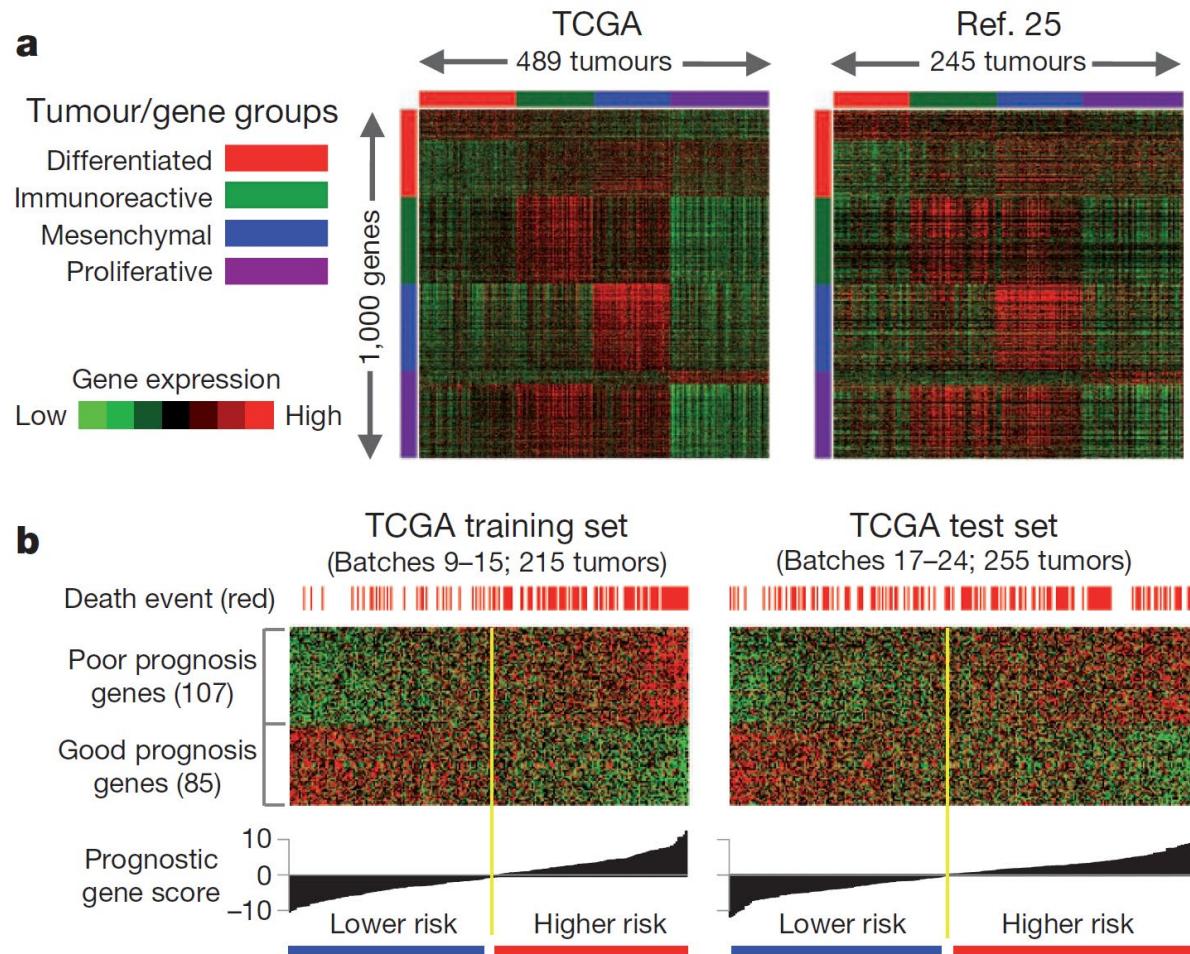
The Cancer Genome Atlas Research Network\*

A catalogue of molecular aberrations that cause ovarian cancer is critical for developing and deploying therapies that will improve patients' lives. The Cancer Genome Atlas project has analysed messenger RNA expression, microRNA expression, promoter methylation and DNA copy number in 489 high-grade serous ovarian adenocarcinomas and the DNA sequences of exons from coding genes in 316 of these tumours. Here we report that high-grade serous ovarian cancer is characterized by TP53 mutations in almost all tumours (96%); low prevalence but statistically recurrent somatic mutations in nine further genes including *NF1*, *BRCA1*, *BRCA2*, *RBI* and *CDK12*; 113 significant focal DNA copy number aberrations; and promoter methylation events involving 168 genes. Analyses delineated four ovarian cancer transcriptional subtypes, three microRNA subtypes, four promoter methylation subtypes and a transcriptional signature associated with survival duration, and shed new light on the impact that tumours with *BRCA1/2* (*BRCA1* or *BRCA2*) and *CCNE1* aberrations have on survival. Pathway analyses suggested that homologous recombination is defective in about half of the tumours analysed, and that NOTCH and FOXM1 signalling are involved in serous ovarian cancer pathophysiology.

# Genomic Data Commons (GDC): TCGA-OV



# Genomic Data Commons (GDC): TCGA-OV



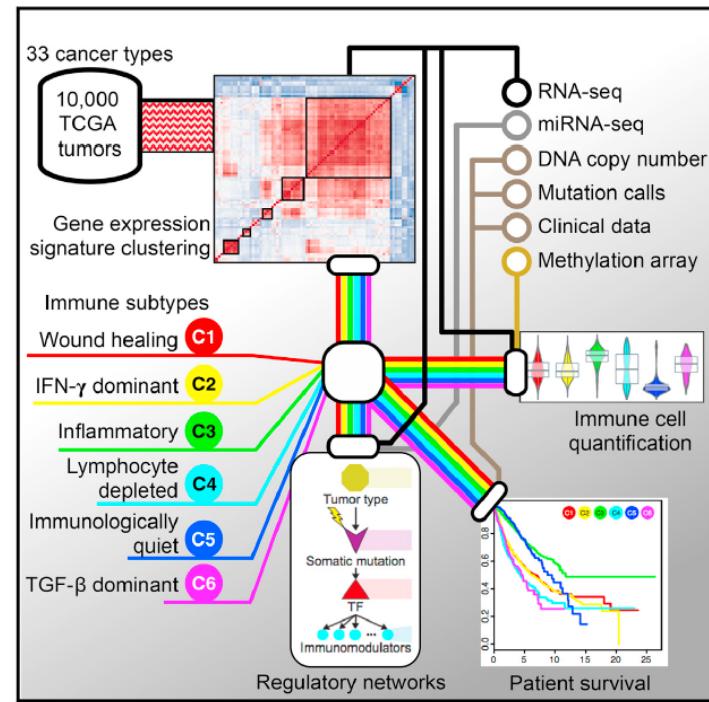
Today's hands on practice!

# Genomic Data Commons (GDC): TCGA-pan-cancer

## Immunity

### The Immune Landscape of Cancer

#### Graphical Abstract



#### Resource

#### Authors

Vésteinn Thorsson, David L. Gibbs,  
Scott D. Brown, ..., Mary L. Disis,  
Benjamin G. Vincent, Ilya Shmulevich

#### Correspondence

vesteinn.thorsson@systemsbiology.org  
(V.T.),  
benjamin.vincent@unchealth.unc.edu  
(B.G.V.),  
ilya.shmulevich@systemsbiology.org (I.S.)

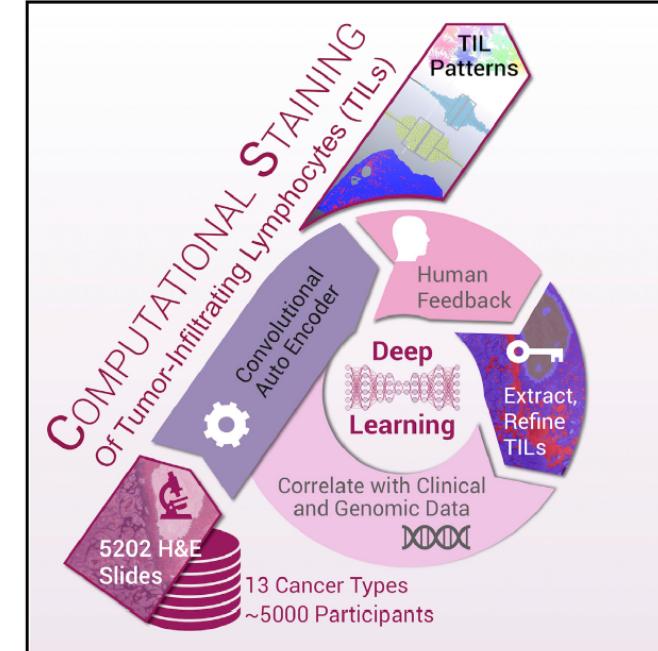
#### In Brief

Thorsson et al. present immunogenomics analyses of more than 10,000 tumors, identifying six immune subtypes that encompass multiple cancer types and are hypothesized to define immune response patterns impacting prognosis. This work provides a resource for understanding tumor-immune interactions, with implications for identifying ways to advance research on immunotherapy.

## Cell Reports

### Spatial Organization and Molecular Correlation of Tumor-Infiltrating Lymphocytes Using Deep Learning on Pathology Images

#### Graphical Abstract



#### Resource

#### Authors

Joel Saltz, Rajarsi Gupta, Le Hou, ..., Alexander J. Lazar, Ashish Sharma, Vésteinn Thorsson

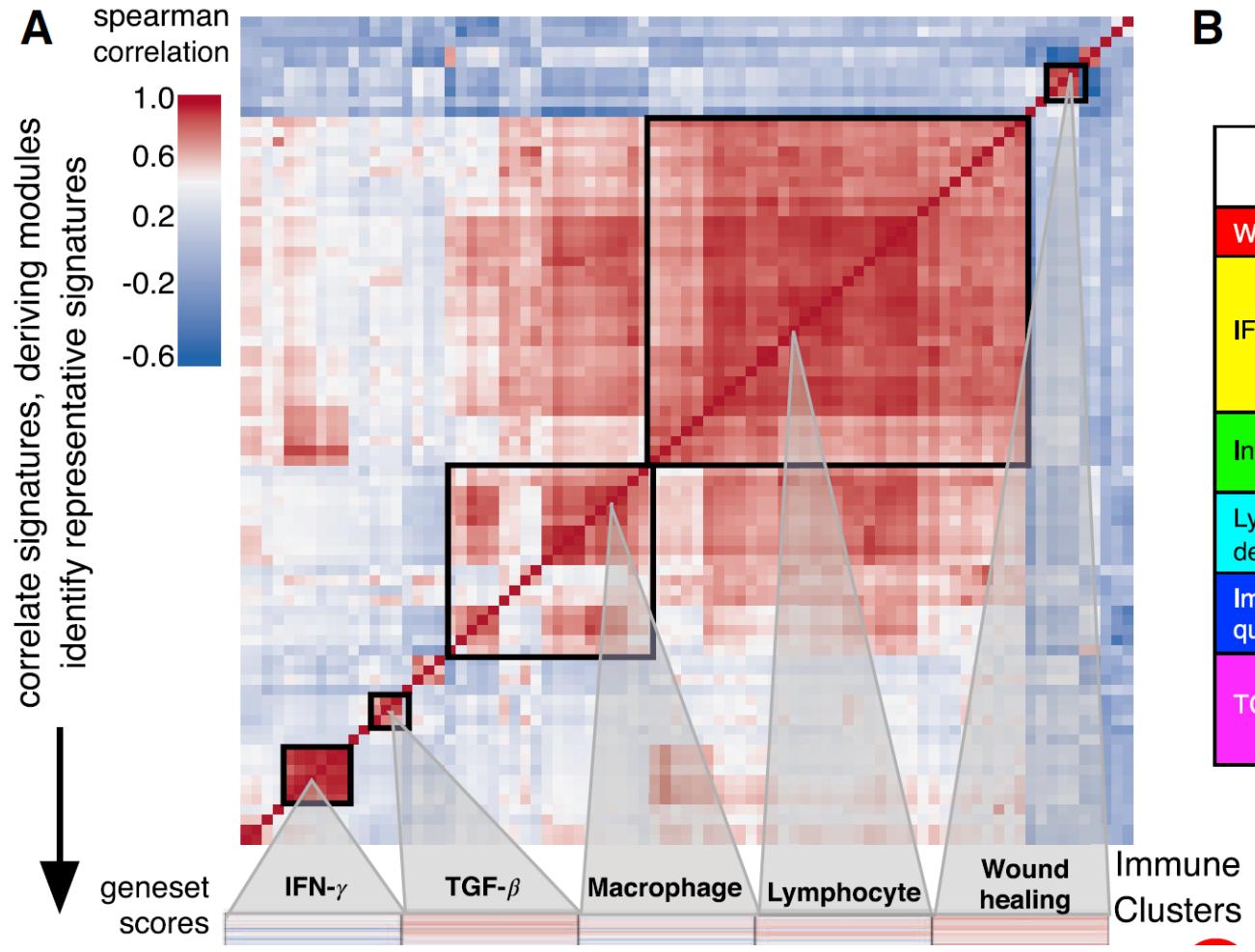
#### Correspondence

joel.saltz@stonybrookmedicine.edu (J.S.),  
vesteinn.thorsson@systemsbiology.org (V.T.)

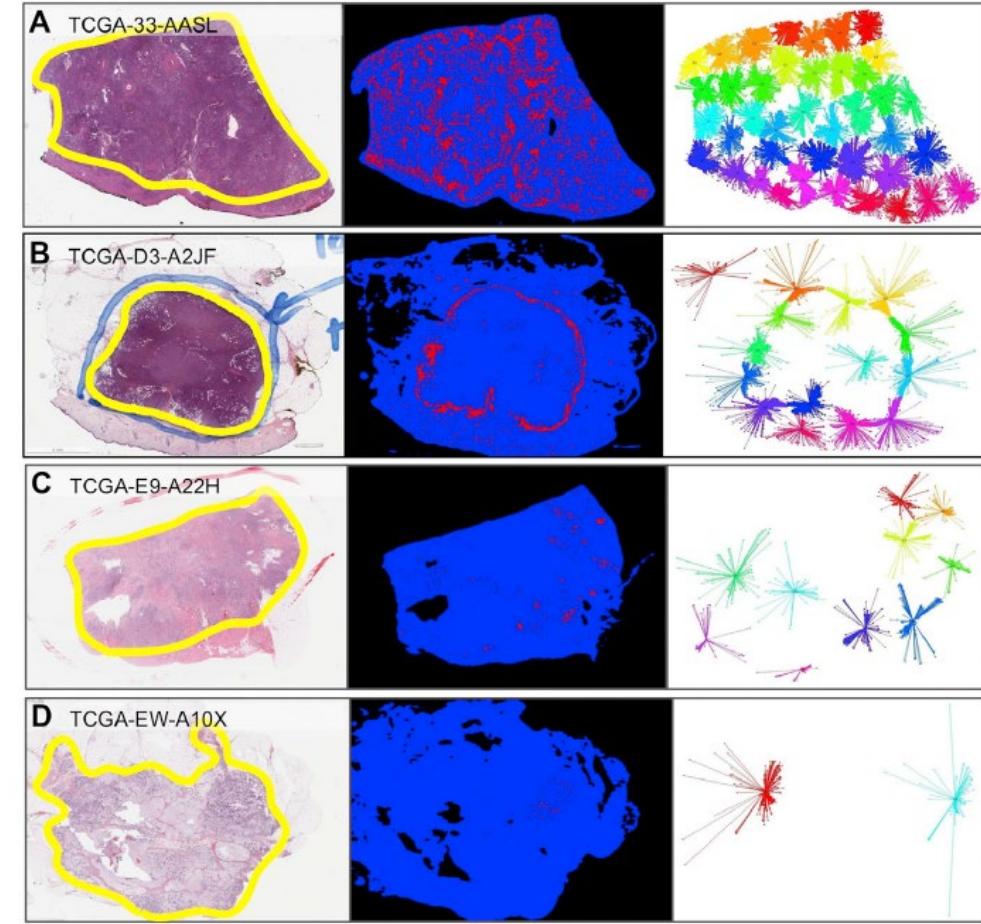
#### In Brief

Tumor-infiltrating lymphocytes (TILs) were identified from standard pathology cancer images by a deep-learning-derived “computational stain” developed by Saltz et al. They processed 5,202 digital images from 13 cancer types. Resulting TIL maps were correlated with TCGA molecular data, relating TIL content to survival, tumor subtypes, and immune profiles.

# Genomic Data Commons (GDC): TCGA-pan-cancer



**B**



**E**

Participant Barcode	Study	Number of TIL Patches	TIL fraction	Number of TIL Clusters	Cluster Size Mean	Within-Cluster Dispersion Mean	Cluster Extent Mean	Ball Hall Baffery Index	Banfield Baffery Index	C Index	Determinant Ratio Index	Global Pattern
TCGA-33-AASL	LUSC	26245	20.6	40	656.1	293456	41.0	447	159518	0.015	2065.4	Brisk Diffuse
TCGA-D3-A2JF	SKCM	6832	4.9	18	379.6	238600	82.1	771	43456	0.022	790.0	Brisk Band-like
TCGA-E9-A22H	BRCA	1000	1.5	10	100.0	54876	51.9	560	6174	0.025	343.0	Non-brisk Multifocal
TCGA-EW-A10X	BRCA	285	0.1	2	142.5	430332	223.0	3093	2283	0.000	29.6	Non-brisk Focal

# Genomic Data Commons (GDC)

<https://datacommons.cancer.gov/>

# Connecting Data to Accelerate Cancer Research

The NCI Cancer Research Data Commons (CRDC) is a cloud-based data science infrastructure that provides secure access to a large, comprehensive, and expanding collection of cancer research data. Users can explore and use analytical and visualization tools for data analysis in the cloud.

## Explore

### REPOSITORIES



#### Cancer Data Service (CDS)

Store and share NCI-funded data that are not hosted elsewhere to further advance scientific discovery across a broad range of research areas.



#### Clinical Trial Data Commons (CTDC)

Store and share data from NCI Clinical Trials. The resource is expected to launch in 2020.



**Genomic Data Commons (GDC)**  
Share, analyze, and visualize harmonized genomic data, including TCGA, TARGET, and CPTAC.



#### Imaging Data Commons (IDC)

Share, analyze, and visualize multi-modal imaging data from both clinical and basic cancer research studies.



#### Integrated Canine Data Commons (ICDC)

Share data from canine clinical trials, including the PRE-medical Cancer Immunotherapy Network Canine Trials (PRECINCT) and the Comparative Oncology Program.



**Proteomic Data Commons (PDC)**  
Share, analyze, and visualize proteomic data, such as CPTAC and The International Cancer Proteogenome Consortium (ICPC).

### INFRASTRUCTURE



#### Cancer Data Aggregator (CDA)

Enables users to query and connect data distributed across the CRDC for integrative analysis.



#### Center for Cancer Data Harmonization (CCDH)

Provides semantic services and tools that facilitate interoperability of data across CRDC.



**Data Commons Framework (DCF)**  
Provides secure user authentication and authorization and permanent digital object identifiers for data objects.

### ANALYTICAL RESOURCES



#### Broad FireCloud

Access elastic compute capacity of Google Cloud Platform to perform large-scale multi-omics analyses.



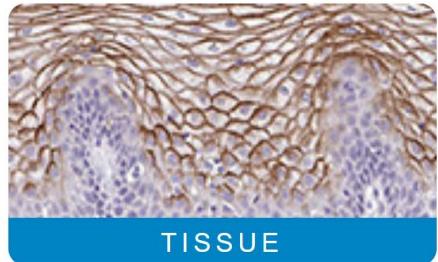
#### Institute for Systems Biology Cancer Genomics Cloud (ISB-CGC)

Access data sets using fully interactive web-based applications, including BigQuery, which is hosted on Google Cloud Platform.

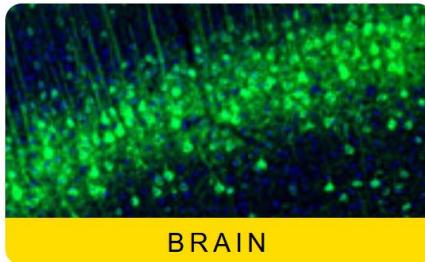


**Seven Bridges Cancer Genomics Cloud (SB-CGC)**  
Explore and analyze large datasets alongside secure and scalable analytical resources for large-scale computational research.

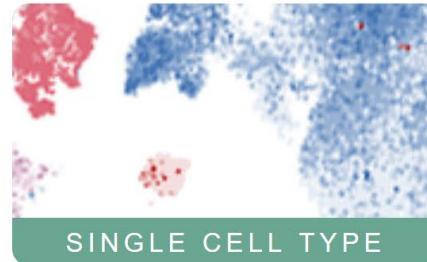
# The Human Protein Atlas



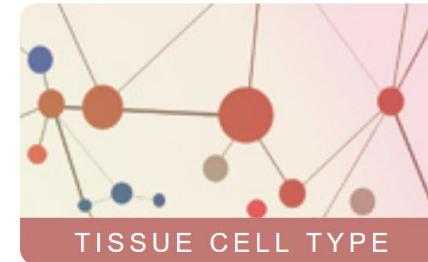
TISSUE



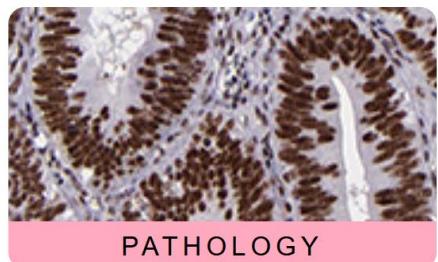
BRAIN



SINGLE CELL TYPE



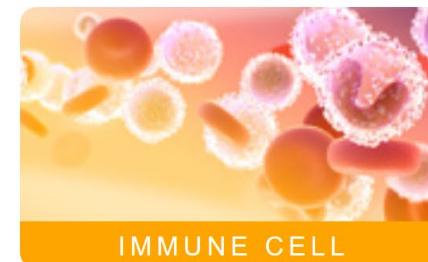
TISSUE CELL TYPE



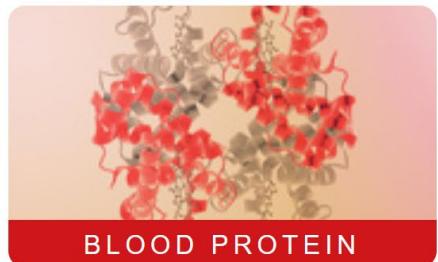
PATHOLOGY

The open access resource for human proteins

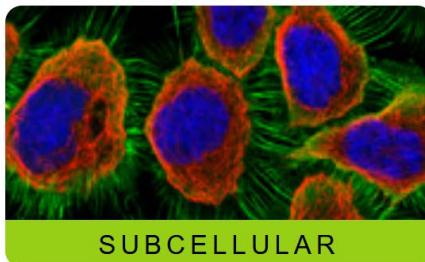
Search for specific genes/proteins or  
explore the 10 different sections



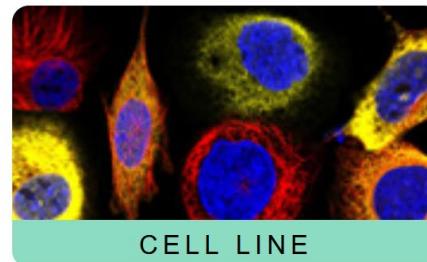
IMMUNE CELL



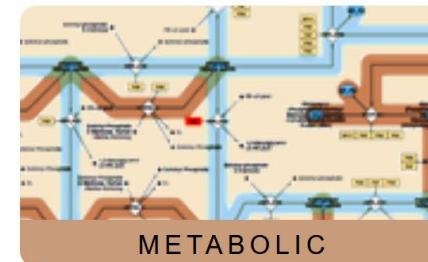
BLOOD PROTEIN



SUBCELLULAR



CELL LINE



METABOLIC

# The Human Protein Atlas

ACE2

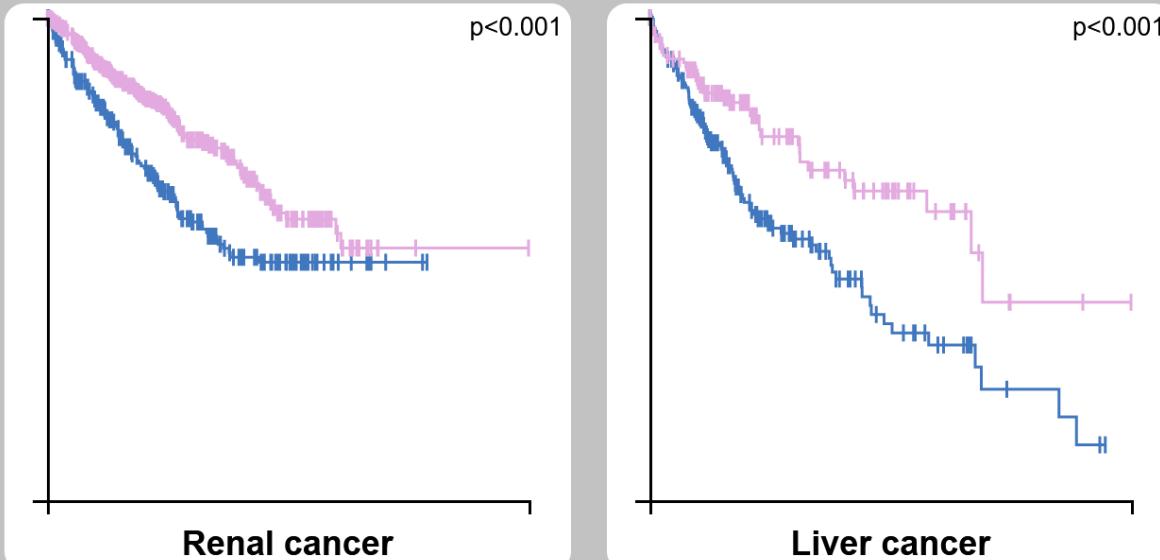


PROTEIN SUMMARY	
RNA DATA	
GENE/PROTEIN	
ANTIBODIES AND VALIDATION	
HUMAN PROTEIN ATLAS SUMMARY <sup>i</sup>	
<b>Protein<sup>i</sup></b>	Angiotensin I converting enzyme 2
<b>Gene name<sup>i</sup></b>	ACE2
<b>Tissue specificity<sup>i</sup></b>	Tissue enhanced (gallbladder, intestine, kidney)
<b>Tissue expression cluster<sup>i</sup></b>	Intestine & Liver - Lipid metabolism (mainly)
<b>Single cell type specificity<sup>i</sup></b>	Cell type enriched (Proximal enterocytes)
<b>Single cell type expression cluster<sup>i</sup></b>	Enterocytes - Digestion (mainly)
<b>Immune cell specificity<sup>i</sup></b>	Not detected in immune cells
<b>Brain specificity<sup>i</sup></b>	Not detected in human brain
<b>Cancer prognostic summary</b>	Prognostic marker in renal cancer (favorable) and liver cancer (favorable)
<b>Predicted location<sup>i</sup></b>	Membrane, Secreted (different isoforms)

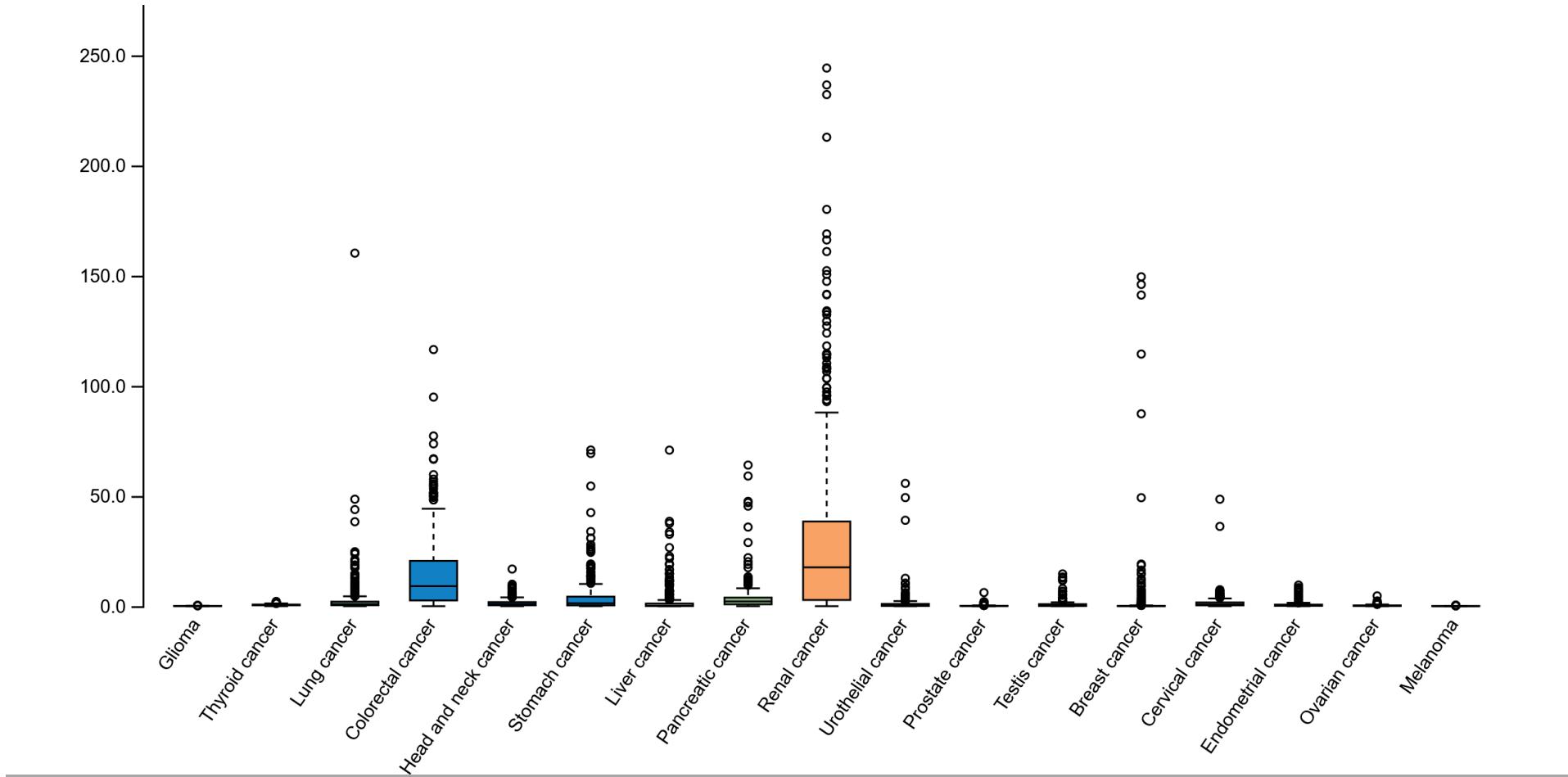
# The Human Protein Atlas

## PROGNOSTIC SUMMARY<sup>i</sup>

Prognostic marker in [renal cancer](#) (favorable) and [liver cancer](#) (favorable)

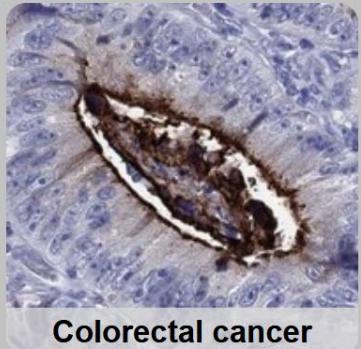


# The Human Protein Atlas

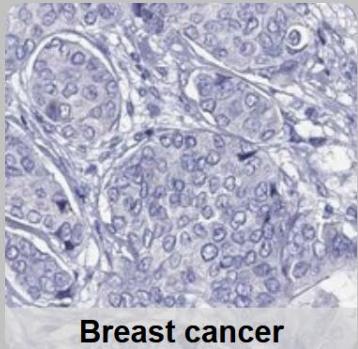


# The Human Protein Atlas

## PROTEIN EXPRESSION<sup>i</sup>



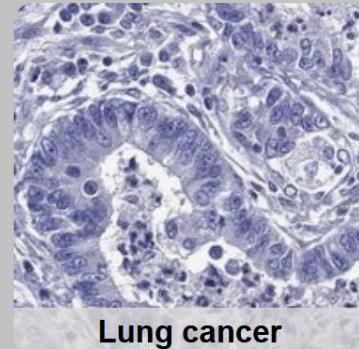
Colorectal cancer



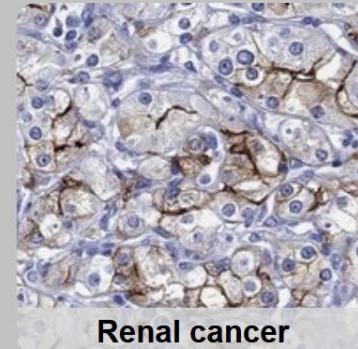
Breast cancer



Prostate cancer



Lung cancer



Renal cancer

## PROTEIN EXPRESSION SUMMARY<sup>i</sup>

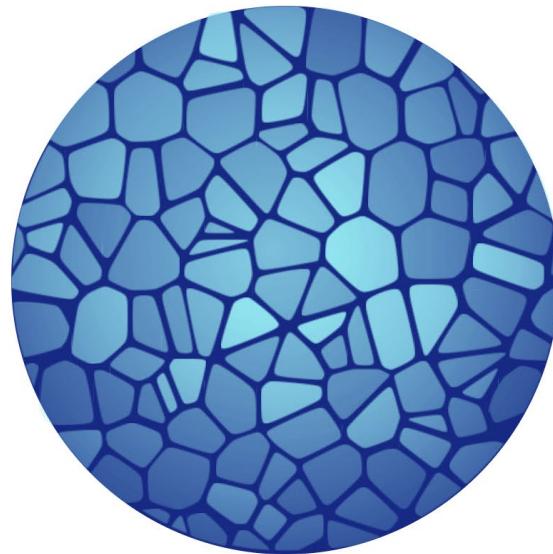
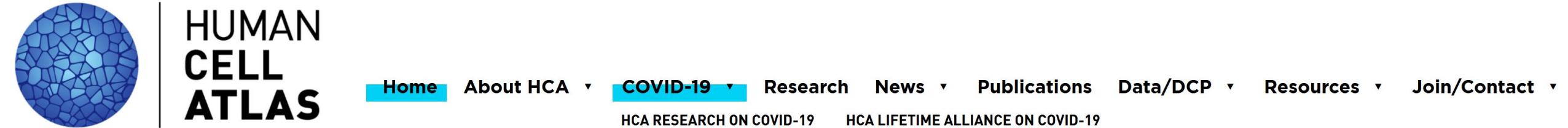
HPA000288 CAB026174 CAB026213 CAB080024 CAB080025 CAB080027 CAB080028

Organ

Expression

Alphabetical

# The Human Cell Atlas



HUMAN  
CELL  
ATLAS

# The Human Cell Atlas

Immune Cell Atlas > Study overview Help & resources Create a study Sign in

## Study: ICA: Blood Mononuclear Cells (2 donors, 2 sites) 13316 cells

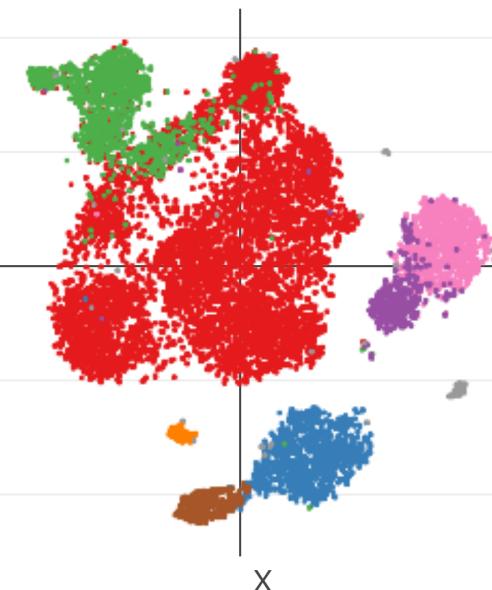
Summary Explore Download

Genes (e.g. "PTEN NF2")



Scatter

lineage\_kmeans



### AllCells - Primary Lineages

Show all		Hide all
1. T	8326	
2. CD14+ Monocyte	1264	
3. NK	1394	
4. Memory B cell	491	
5. DC	142	
6. CD16+ Monocyte	398	
7. Naive B cell	1169	

### OPTIONS

#### Clustering

lineage\_kmeans

#### Annotation

AllCells - Primary L...

+ Create annotation

#### Subsampling

All Cells

Reset view

# The Human Cell Atlas



ARTICLE

Check for updates

<https://doi.org/10.1038/s41467-021-24467-0>

OPEN

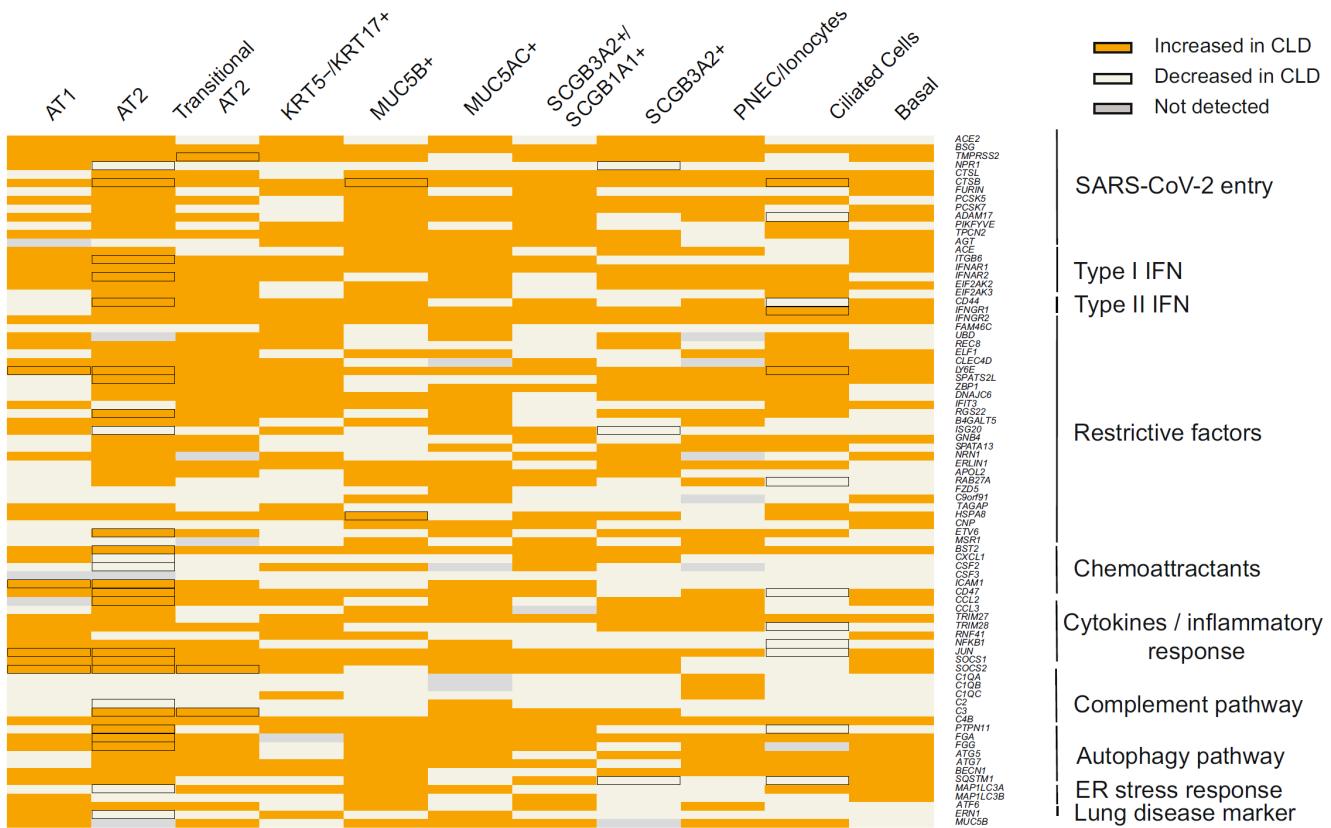
## Chronic lung diseases are associated with gene expression programs favoring SARS-CoV-2 entry and severity

Linh T. Bui<sup>1,87</sup>, Nichelle I. Winters<sup>2,87</sup>, Mei-I Chung<sup>1</sup>, Chitra Joseph<sup>1,3</sup>, Austin J. Gutierrez<sup>1,3</sup>, Arun C. Habermann<sup>2</sup>, Taylor S. Adams<sup>4</sup>, Jonas C. Schupp<sup>1,4</sup>, Sergio Poli<sup>1,5</sup>, Lance M. Peter<sup>1</sup>, Chase J. Taylor<sup>2</sup>, Jessica B. Blackburn<sup>2</sup>, Bradley W. Richmond<sup>1,2,6</sup>, Andrew G. Nicholson<sup>7,8</sup>, Doris Rassl<sup>9</sup>, William A. Wallace<sup>10,11</sup>, Ivan O. Rosas<sup>12</sup>, R. Gisli Jenkins<sup>1,3</sup>, Naftali Kaminski<sup>1,4</sup>, Jonathan A. Kropski<sup>2,6,13,88</sup>, Nicholas E. Banovich<sup>1,88</sup>, and the Human Cell Atlas Lung Biological Network\*

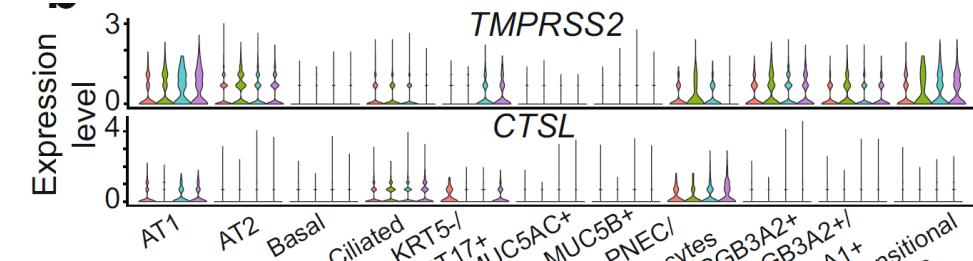
Patients with chronic lung disease (CLD) have an increased risk for severe coronavirus disease-19 (COVID-19) and poor outcomes. Here, we analyze the transcriptomes of 611,398 single cells isolated from healthy and CLD lungs to identify molecular characteristics of lung cells that may account for worse COVID-19 outcomes in patients with chronic lung diseases. We observe a similar cellular distribution and relative expression of SARS-CoV-2 entry factors in control and CLD lungs. CLD AT2 cells express higher levels of genes linked directly to the efficiency of viral replication and the innate immune response. Additionally, we identify basal differences in inflammatory gene expression programs that highlight how CLD alters the inflammatory microenvironment encountered upon viral exposure to the peripheral lung. Our study indicates that CLD is accompanied by changes in cell-type-specific gene expression programs that prime the lung epithelium for and influence the innate and adaptive immune responses to SARS-CoV-2 infection.

# The Human Cell Atlas

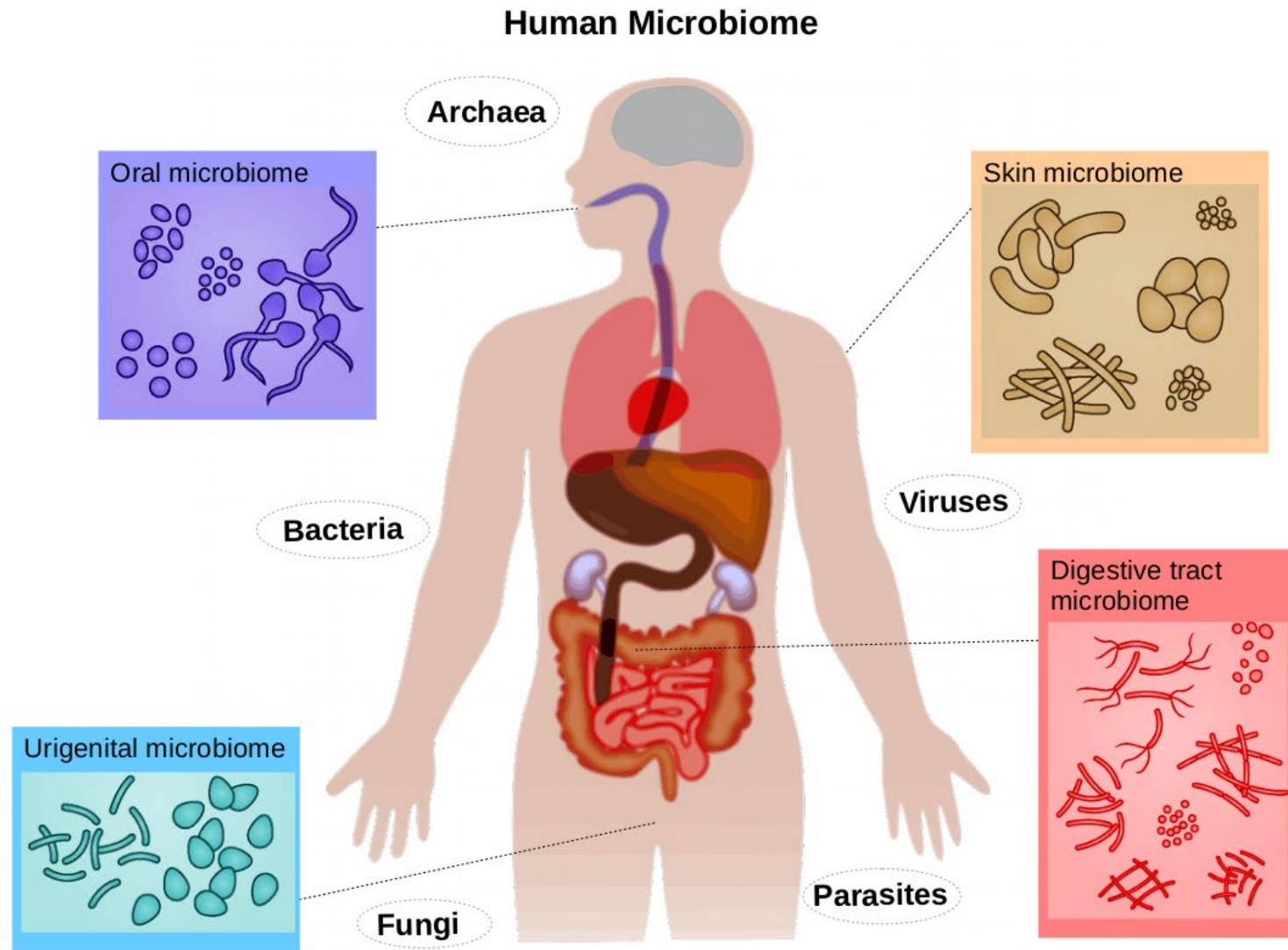
**a**



**b**



# The Human Microbiome



# The Human Microbiome – Cancer Microbiome



RESEARCH

## CANCER IMMUNOTHERAPY

# The commensal microbiome is associated with anti-metastatic melanoma

Vyara Matson,<sup>1\*</sup> Jessica Fessler,<sup>1\*</sup> Riyue Bao,<sup>2,3\*</sup>  
Maria-Luisa Alegre,<sup>4</sup> Jason J. Luke,<sup>4</sup> Thomas F.

Anti-PD-1-based immunotherapy has had a major impact on cancer treatment, but it has only benefited a subset of patients. Among the variables that contribute to this heterogeneity is differential composition of the gut microbiome. We analyzed baseline stool samples from 111 patients before immunotherapy treatment, through an integrated approach of metagenomic sequencing, metagenomic shotgun sequencing, and a targeted PCR assay for selected bacteria. A significant association was found between gut microbial composition and clinical response. Bacterial species associated with anti-PD-1 responders included *Bifidobacterium longum*, *Clostridium difficile*, and *Enterococcus faecium*. Reconstitution of germ-free mice with fecal microbiota from anti-PD-1 responders could lead to improved tumor control, augmented anti-PD-L1 therapy. Our results suggest that the gut microbiome may have a mechanistic impact on antitumor immunity in humans.

RESEARCH

## CANCER IMMUNOTHERAPY

# Gut microbiome modulates response to anti-PD-1 immunotherapy in melanoma patients

RESEARCH

## CANCER IMMUNOTHERAPY

# Gut microbiome influences efficacy of PD-1-based immunotherapy against epithelial tumors

Bertrand Routy,<sup>1,2,3</sup> Emmanuelle Le Chatelier,<sup>4</sup> Lisa Derosa,<sup>1,2,3</sup> Connie P. M. Duong,<sup>1,2,5</sup> Maryam Tidjani Alou,<sup>1,2,3</sup> Romain Daillère,<sup>1,2,3</sup> Aurélie Fluckiger,<sup>1,2,5</sup> Meriem Messaoudene,<sup>1,2</sup> Conrad Rauber,<sup>1,2,3</sup> Maria P. Roberti,<sup>1,2,5</sup> Marine Fidelle,<sup>1,3,5</sup> Caroline Flament,<sup>1,2,5</sup> Vichnou Poirier-Colame,<sup>1,2,5</sup> Paule Opolon,<sup>6</sup> Christophe Klein,<sup>7</sup> Kristina Iribarren,<sup>8,9,10,11,12</sup> Laura Mondragón,<sup>8,9,10,11,12</sup> Nicolas Jacquemet,<sup>1,2,3</sup> Bo Qu,<sup>1,2,3</sup> Gladys Ferrere,<sup>1,2,3</sup> Céline Clémenson,<sup>1,13</sup> Laura Mezquita,<sup>1,14</sup> Jordi Remon Masip,<sup>1,14</sup> Charles Nalbet,<sup>15</sup> Solenn Brosseau,<sup>15</sup> Coureche Kaderbhai,<sup>16</sup> Corentin Richard,<sup>16</sup> Hira Rizvi,<sup>17</sup> Florence Levenez,<sup>4</sup> Nathalie Galleron,<sup>4</sup> Benoit Quinquis,<sup>4</sup> Nicolas Pons,<sup>4</sup> Bernhard Ryffel,<sup>18</sup> Véronique Minard-Colin,<sup>1,19</sup> Patrick Gonin,<sup>1,20</sup> Jean-Charles Soria,<sup>1,14</sup> Eric Deutsch,<sup>1,13</sup> Yohann Loriot,<sup>1,3,14</sup> François Ghiringhelli,<sup>16</sup> Gérard Zaleman,<sup>15</sup> François Goldwasser,<sup>9,21,22</sup> Bernard Escudier,<sup>1,14,23</sup> Matthew D. Hellmann,<sup>24,25</sup> Alexander Eggertmont,<sup>1,2,14</sup> Didier Raoult,<sup>26</sup> Laurence Albiges,<sup>1,3,14</sup> Guido Kroemer,<sup>8,9,10,11,12,27,28\*</sup> Laurence Zitvogel<sup>1,2,3,5\*</sup>

# The Human Microbiome – Cancer Microbiome

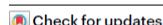
## Microbiome analyses of blood and tissues suggest cancer diagnostic approach

<https://doi.org/10.1038/s41586-020-2095-1>

Received: 7 June 2019

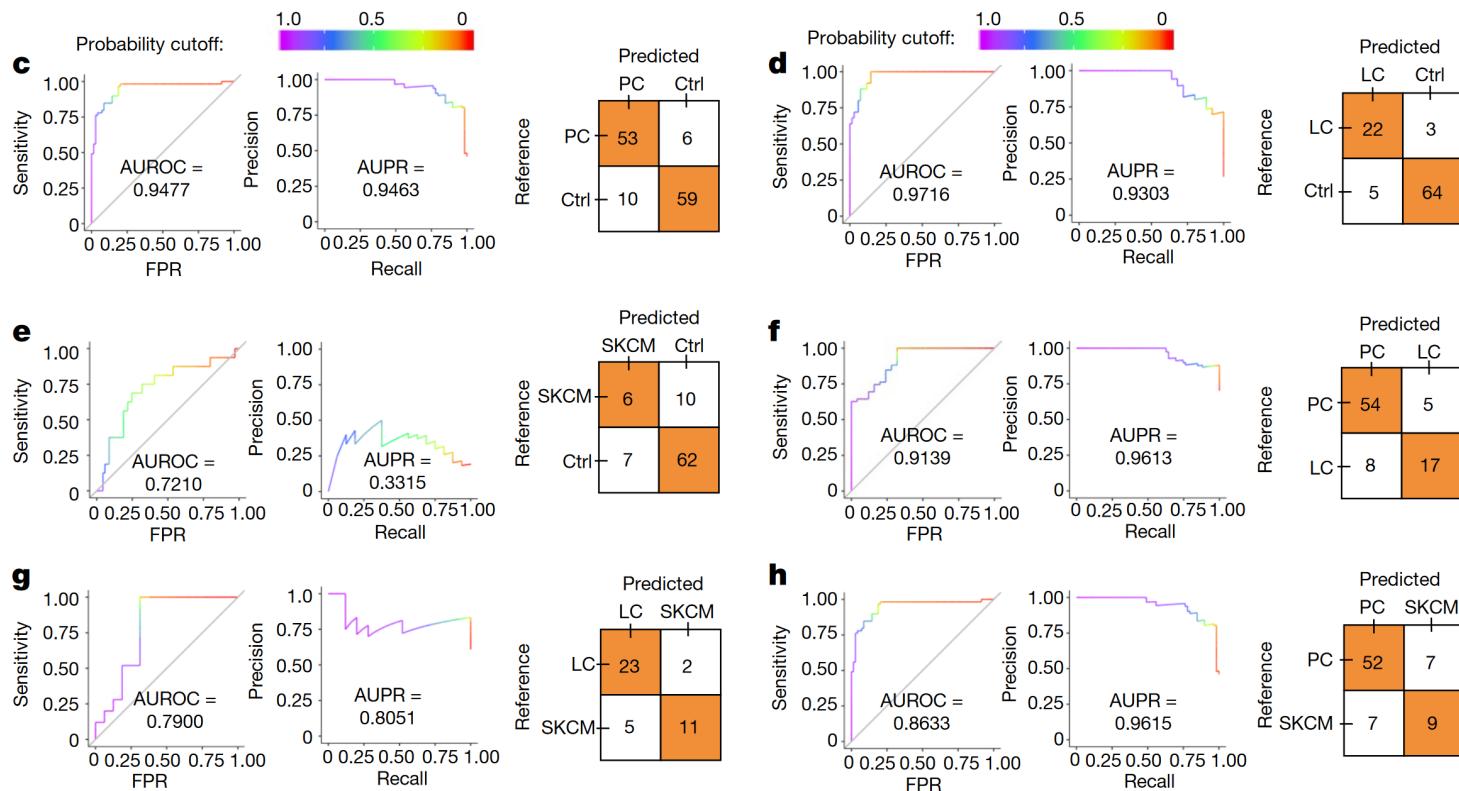
Accepted: 6 February 2020

Published online: 11 March 2020



Gregory D. Poore<sup>1,2</sup>, Evgenia Kopylova<sup>2,9,12</sup>, Qiyun Zhu<sup>2</sup>, Carolina Carpenter<sup>3</sup>, Serena Fraraccio<sup>3</sup>, Stephen Wandro<sup>3</sup>, Tomasz Kosciolek<sup>2,10</sup>, Stefan Janssen<sup>2,11</sup>, Jessica Metcalf<sup>4</sup>, Se Jin Song<sup>3</sup>, Jad Kanbar<sup>5</sup>, Sandrine Miller-Montgomery<sup>1,3</sup>, Robert Heaton<sup>6</sup>, Rana McKay<sup>7</sup>, Sandip Pravin Patel<sup>1,2,7</sup>, Austin D. Swafford<sup>3</sup> & Rob Knight<sup>1,2,3,8,✉</sup>

Systematic characterization of the cancer microbiome provides the opportunity to develop techniques that exploit non-human, microorganism-derived molecules in the diagnosis of a major human disease. Following recent demonstrations that some types of cancer show substantial microbial contributions<sup>1–10</sup>, we re-examined whole-genome and whole-transcriptome sequencing studies in The Cancer Genome Atlas<sup>11</sup> (TCGA) of 33 types of cancer from treatment-naïve patients (a total of 18,116 samples) for microbial reads, and found unique microbial signatures in tissue and blood within and between most major types of cancer. These TCGA blood signatures remained predictive when applied to patients with stage Ia–IIC cancer and cancers lacking any genomic alterations currently measured on two commercial-grade cell-free tumour DNA platforms, despite the use of very stringent decontamination analyses that discarded up to 92.3% of total sequence data. In addition, we could discriminate among samples from healthy, cancer-free individuals ( $n=69$ ) and those from patients with multiple types of cancer (prostate, lung, and melanoma; 100 samples in total) solely using plasma-derived, cell-free microbial nucleic acids. This potential microbiome-based oncology diagnostic tool warrants further exploration.



# Sample data

## ARTICLE

---

---

doi:10.1038/nature10166

# Integrated genomic analyses of ovarian carcinoma

The Cancer Genome Atlas Research Network\*

A catalogue of molecular aberrations that cause ovarian cancer is critical for developing and deploying therapies that will improve patients' lives. The Cancer Genome Atlas project has analysed messenger RNA expression, microRNA expression, promoter methylation and DNA copy number in 489 high-grade serous ovarian adenocarcinomas and the DNA sequences of exons from coding genes in 316 of these tumours. Here we report that high-grade serous ovarian cancer is characterized by *TP53* mutations in almost all tumours (96%); low prevalence but statistically recurrent somatic mutations in nine further genes including *NFI*, *BRCA1*, *BRCA2*, *RBI* and *CDK12*; 113 significant focal DNA copy number aberrations; and promoter methylation events involving 168 genes. Analyses delineated four ovarian cancer transcriptional subtypes, three microRNA subtypes, four promoter methylation subtypes and a transcriptional signature associated with survival duration, and shed new light on the impact that tumours with *BRCA1/2* (*BRCA1* or *BRCA2*) and *CCNE1* aberrations have on survival. Pathway analyses suggested that homologous recombination is defective in about half of the tumours analysed, and that NOTCH and FOXM1 signalling are involved in serous ovarian cancer pathophysiology.

# Sample data

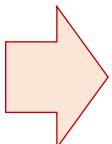
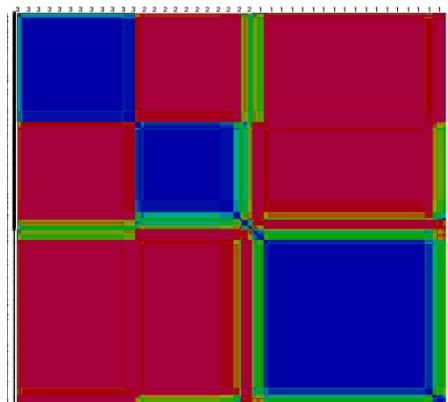
## ARTICLE

doi:10.1038/nature10166

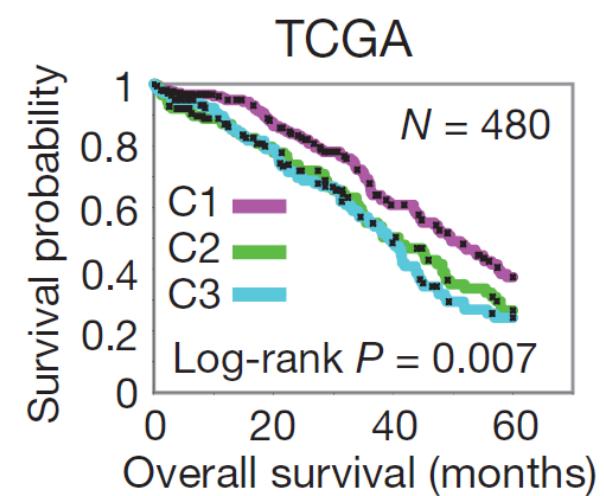
# Integrated genomic analyses of ovarian carcinoma

The Cancer Genome Atlas Research Network\*

150 most variable  
miRNAs for sample  
clustering



Cluster 1 has significantly  
better survival



NATIONAL CANCER INSTITUTE  
GDC Data Portal

Home Projects Exploration Analysis Repository

Manage Sets

Harmonized Cancer Datasets

# Genomic Data Commons Data Portal

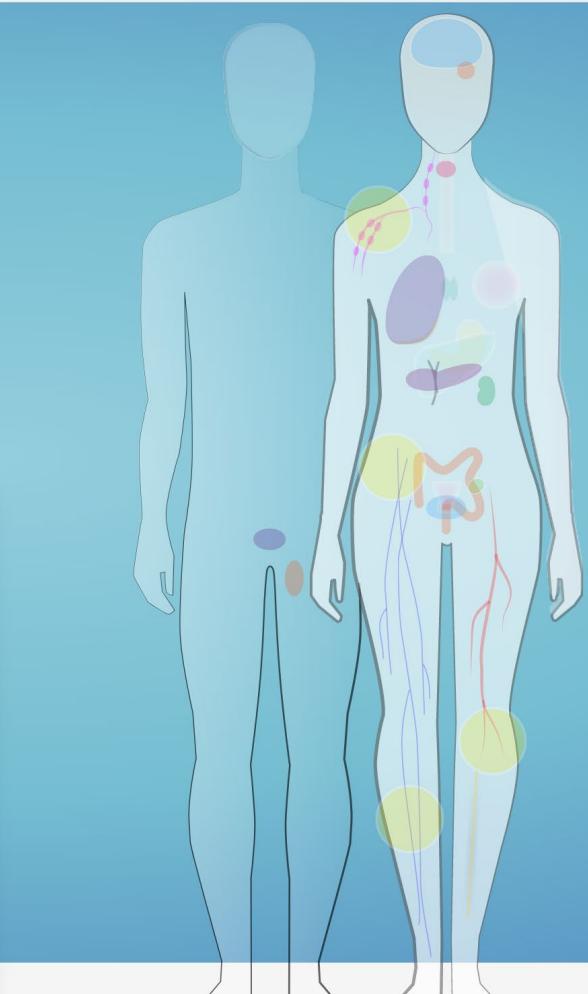
Get Started by Exploring:

Projects Exploration Analysis Repository

e.g. BRAF, Breast, TCGA-BLCA, TCGA-A5-A0G2

## Data Portal Summary Data Release 33.1 - May 31, 2022

PROJECTS	PRIMARY SITES	CASES
72	67	85,552
FILES	GENES	MUTATIONS
828,104	21,759	2,683,650



Cases by Major Primary Site

Primary Site	Cases
Adrenal Gland	1
Bile Duct	1
Bladder	1
Bone	1
Bone Marrow	9
Brain	1
Breast	9
Cervix	1
Colorectal	8
Esophagus	1
Eye	1
Head and Neck	3
Kidney	3
Liver	1
Lung	10
Lymph Nodes	1
Nervous System	4
Ovary	3
Pancreas	2
Pleura	1
Prostate	2
Skin	3
Soft Tissue	1
Stomach	1
Testis	1
Thymus	1
Thyroid	2
Uterus	2

# The Cancer Genome Atlas (TCGA)

## Ovarian cancer



<https://portal.gdc.cancer.gov>

Files

Cases



[Add a File Filter](#)

▼ Search Files



e.g. 142682.bam, 4f6e2e7a-b...



▼ Data Category

simple nucleotide variation

# Files

8,408

copy number variation

3,578

sequencing reads

2,684

biospecimen

2,601

dna methylation

1,869

4 More...



▼ Data Type

Annotated Somatic Mutation

# Files

4,900

Aligned Reads

2,684

RNA-seq

2,546

[Clear](#) Program Name IS TCGA AND Project Id IS TCGA-OV

Files

(23,453)

Cases

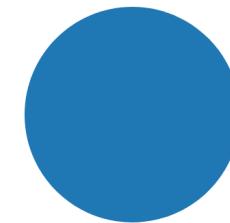
(608)

[Add All Files to Cart](#)

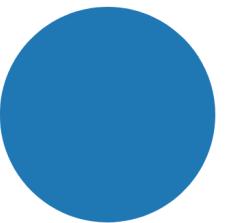
[Manifest](#)

[View Images](#)

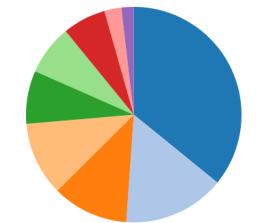
Primary Site



Project



Data Category



[Show More](#)

Showing 1 - 20 of 23,453 files 46.67 TB

Access

File Name

Cases Project

Data Category

Data Format

File Size

Annotations

[ec5ab2eb-37b0-4f3c-a35c-](#)

ID	Disease Type	Primary Site	Program	Cases
<a href="#">TCGA-BRCA</a>	Breast Invasive Carcinoma	Breast	TCGA	<a href="#">1,098</a>
<a href="#">TCGA-GBM</a>	Glioblastoma Multiforme	Brain	TCGA	<a href="#">617</a>
<a href="#">TCGA-OV</a>	Ovarian Serous Cystadenocarcinoma	Ovary	TCGA	<a href="#">608</a>
<a href="#">TCGA-LUAD</a>	Lung Adenocarcinoma	Lung	TCGA	<a href="#">585</a>
<a href="#">TCGA-UCEC</a>	Uterine Corpus Endometrial Carcinoma	Uterus	TCGA	<a href="#">560</a>
<a href="#">TCGA-KIRC</a>	Kidney Renal Clear Cell Carcinoma	Kidney	TCGA	<a href="#">537</a>
<a href="#">TCGA-HNSC</a>	Head and Neck Squamous Cell Carcinoma	Head and Neck	TCGA	<a href="#">528</a>
<a href="#">TCGA-LGG</a>	Brain Lower Grade Glioma	Brain	TCGA	<a href="#">516</a>
<a href="#">TCGA-THCA</a>	Thyroid Carcinoma	Thyroid	TCGA	<a href="#">507</a>
<a href="#">TCGA-LUSC</a>	Lung Squamous Cell Carcinoma	Lung	TCGA	<a href="#">504</a>
<a href="#">TCGA-PRAD</a>	Prostate Adenocarcinoma	Prostate	TCGA	<a href="#">500</a>
<a href="#">TCGA-SKCM</a>	Skin Cutaneous Melanoma	Skin	TCGA	<a href="#">470</a>
<a href="#">TCGA-COAD</a>	Colon Adenocarcinoma	Colorectal	TCGA	<a href="#">461</a>
<a href="#">TCGA-STAD</a>	Stomach Adenocarcinoma	Stomach	TCGA	<a href="#">443</a>
<a href="#">TCGA-BLCA</a>	Bladder Urothelial Carcinoma	Bladder	TCGA	<a href="#">412</a>
<a href="#">TCGA-LIHC</a>	Liver Hepatocellular Carcinoma	Liver	TCGA	<a href="#">377</a>
<a href="#">TCGA-CESC</a>	Cervical Squamous Cell Carcinoma	Cervix	TCGA	<a href="#">307</a>
<a href="#">TCGA-KIRP</a>	Kidney Renal Papillary Cell Carcinoma	Kidney	TCGA	<a href="#">291</a>
<a href="#">TCGA-SARC</a>	Sarcoma	Soft Tissue	TCGA	<a href="#">261</a>
<a href="#">TCGA-LAML</a>	Acute Myeloid Leukemia	Bone Marrow	TCGA	<a href="#">200</a>
<a href="#">TCGA-PAAD</a>	Pancreatic Adenocarcinoma	Pancreas	TCGA	<a href="#">185</a>
<a href="#">TCGA-ESCA</a>	Esophageal Carcinoma	Esophagus	TCGA	<a href="#">185</a>
<a href="#">TCGA-PCPG</a>	Pheochromocytoma and Paraganglioma	Adrenal Gland	TCGA	<a href="#">179</a>
<a href="#">TCGA-READ</a>	Rectum Adenocarcinoma	Colorectal	TCGA	<a href="#">172</a>
<a href="#">TCGA-TGCT</a>	Testicular Germ Cell Tumors	Testis	TCGA	<a href="#">150</a>
<a href="#">TCGA-THYM</a>	Thymoma	Thymus	TCGA	<a href="#">124</a>
<a href="#">TCGA-KICH</a>	Kidney Chromophobe	Kidney	TCGA	<a href="#">113</a>
<a href="#">TCGA-ACC</a>	Adrenocortical Carcinoma	Adrenal Gland	TCGA	<a href="#">92</a>
<a href="#">TCGA-MESO</a>	Mesothelioma	Pleura	TCGA	<a href="#">87</a>
<a href="#">TCGA-UVM</a>	Uveal Melanoma	Eye	TCGA	<a href="#">80</a>
<a href="#">TCGA-DLBC</a>	Lymphoid Neoplasm Diffuse Large B-Cell Lymphoma	Lymph Nodes	TCGA	<a href="#">58</a>
<a href="#">TCGA-UCS</a>	Uterine Carcinosarcoma	Uterus	TCGA	<a href="#">57</a>
<a href="#">TCGA-CHOL</a>	Cholangiocarcinoma	Bile Duct	TCGA	<a href="#">51</a>
<b>Total</b>				<b>11,315</b>

# TCGA cancer types (n=33)

- TCGA raw data were harmonized by NCI's GDC team
  - Release 6 June 2017
  - Release 33.1 (May 2022; current)
- Data types: from raw files to compiled results
- Result access: public or protected
- Apply for access: dbGap
- Download: GDC

# Access GDC data

```
[rbao@cri16in002 ~]$ gdc-client -h
usage: gdc-client [-h] [--version] {download,upload,interactive} ...

The Genomic Data Commons Command Line Client

optional arguments:
  -h, --help            show this help message and exit
  --version             show program's version number and exit

commands:
  {download,upload,interactive}
    download           for more information, specify -h after a command
    upload              download data from the GDC
    interactive         upload data to the GDC
                        run in interactive mode
```

## Bioconductor packages

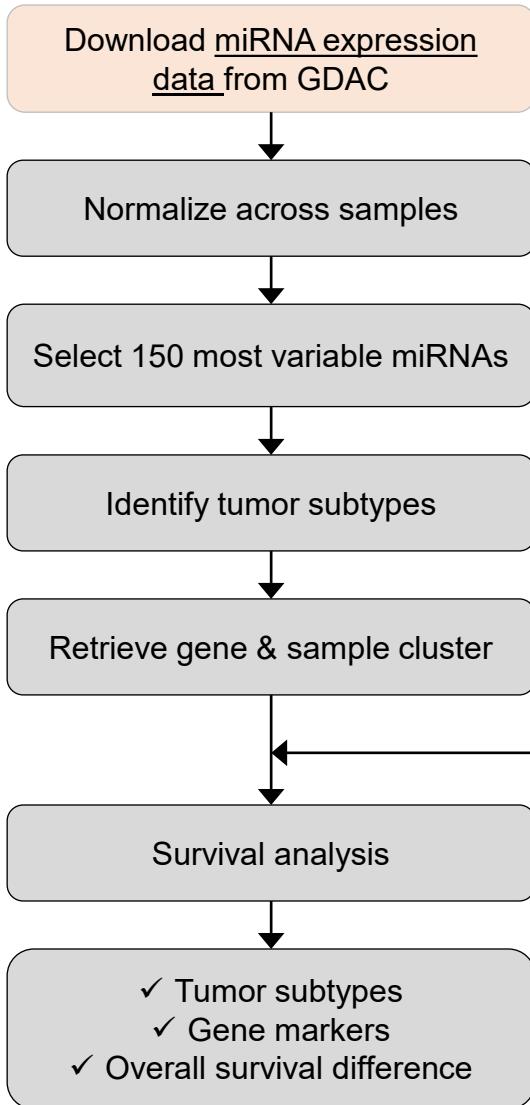
- library(TCGAbiolinks)
  - <https://bioconductor.org/packages/release/bioc/html/TCGAbiolinks.html>
- library(GenomicDataCommons)
  - <https://bioconductor.org/packages/release/bioc/html/GenomicDataCommons.html>

## The GDC Application Programming Interface (API): An Overview

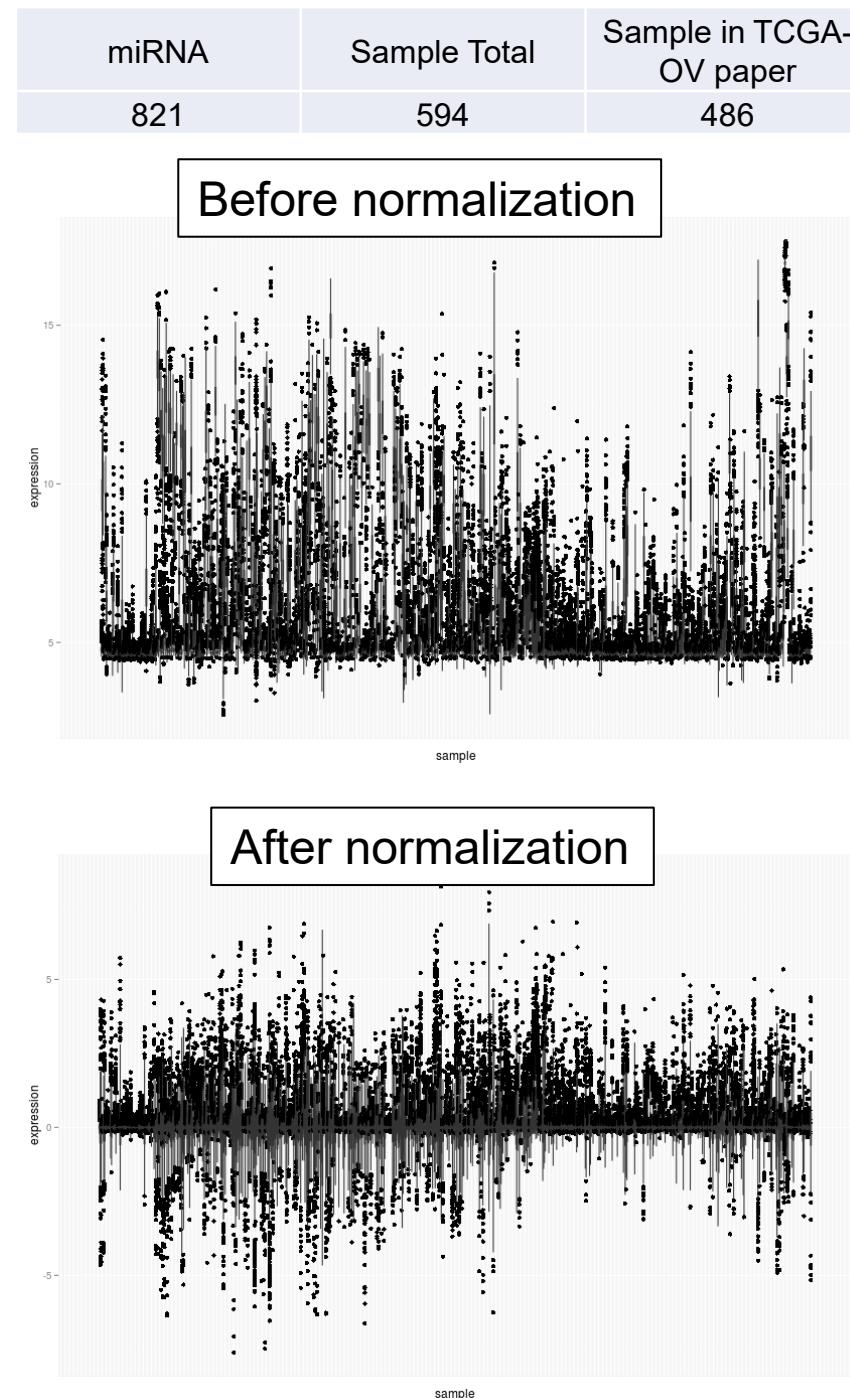
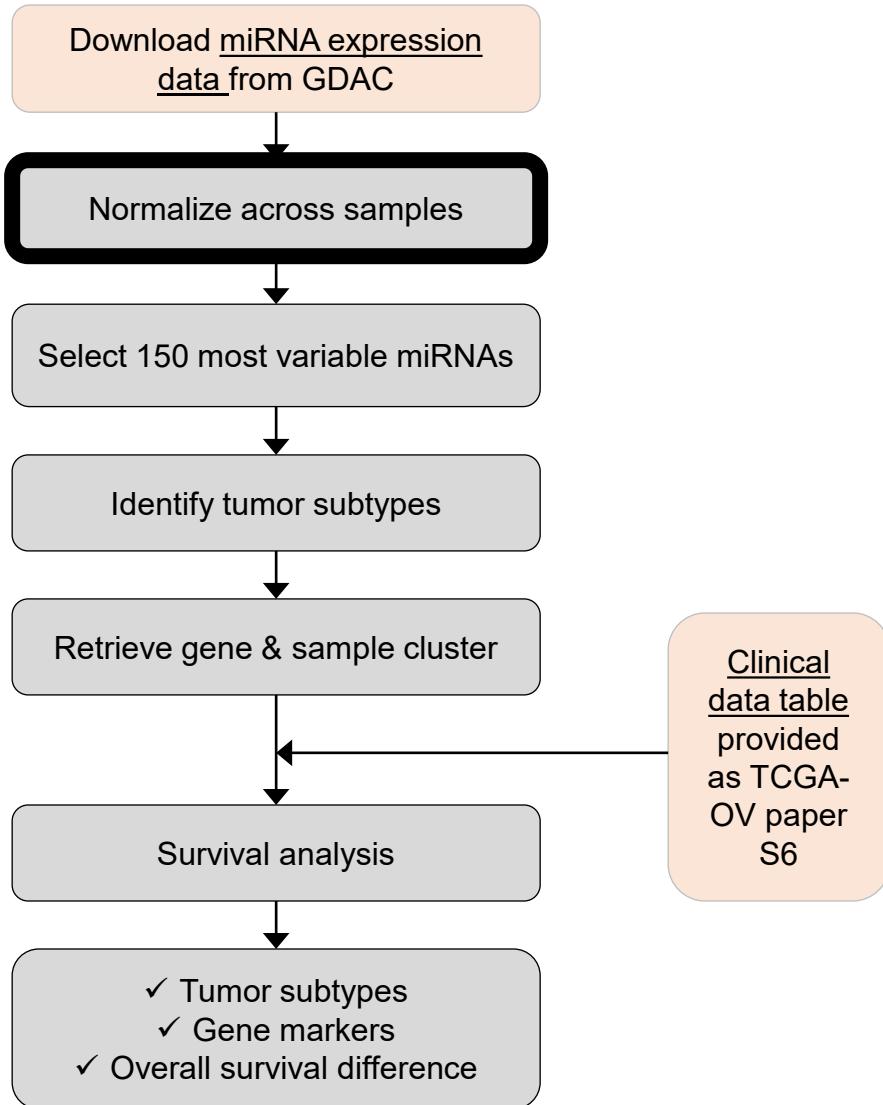
The GDC API drives the GDC Data and Submission Portals and provides programmatic access to GDC functionality. This includes searching for, downloading, and submitting data and metadata. The GDC API uses JSON as its communication format, and standard HTTP methods like `GET`, `PUT`, `POST` and `DELETE`.

```
curl https://api.gdc.cancer.gov/files/e6fe614f-7fb4-4d6c-85c1-311ec25af938?pretty=true
```

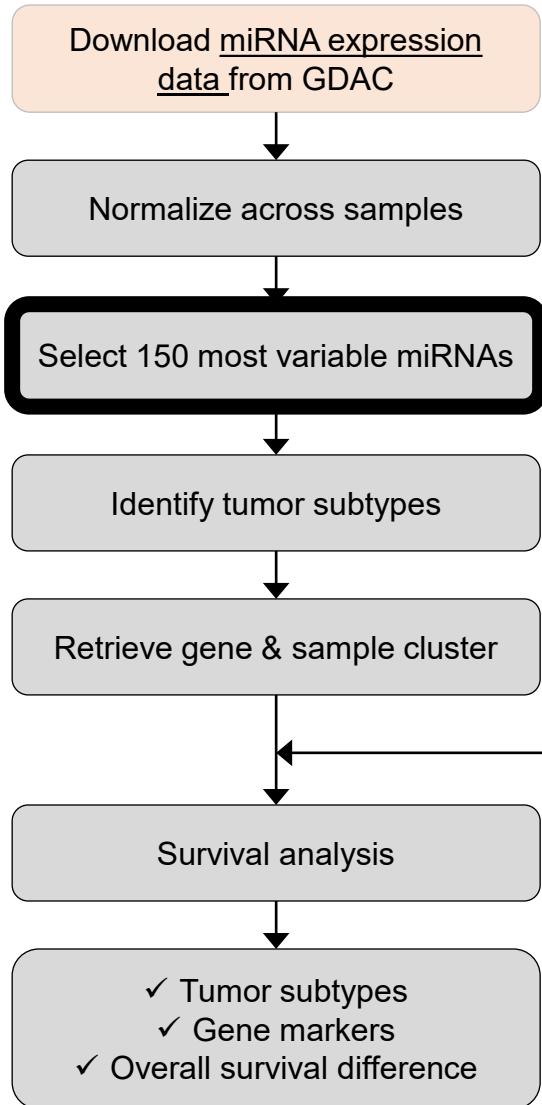
# Integrate Expression with Clinical Data



# Integrate Expression with Clinical Data

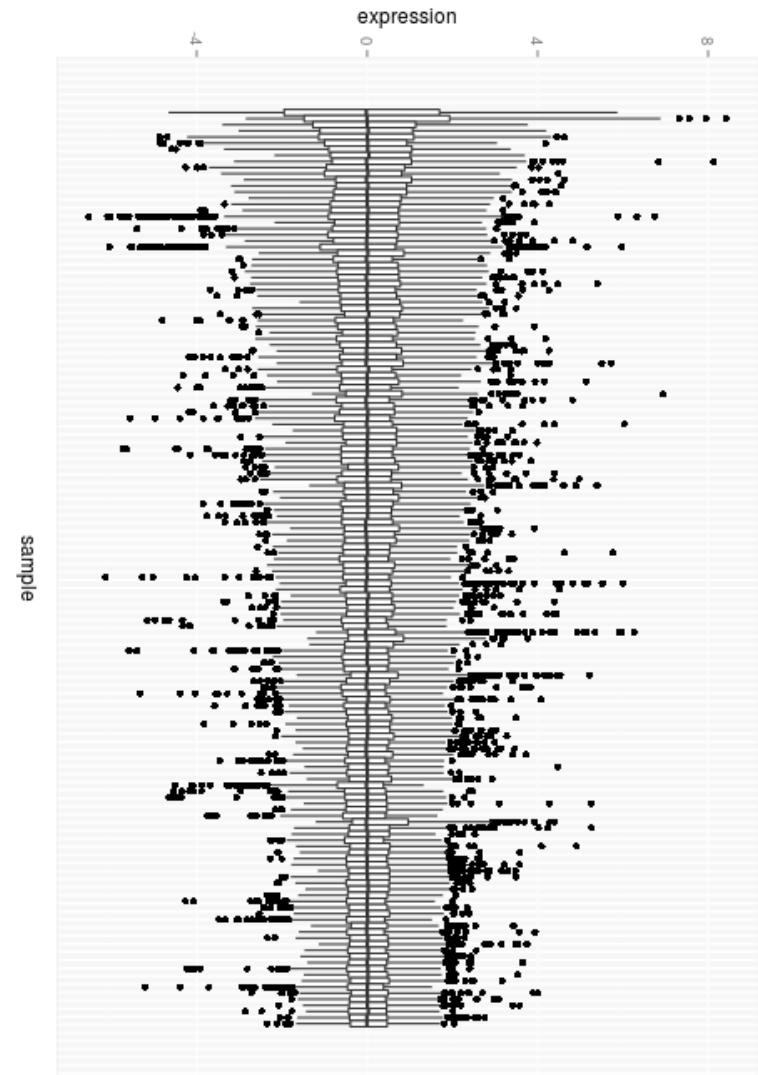


# Integrate Expression with Clinical Data



	mad
HSA-MIR-205	2.6824804
HSA-MIR-449A	2.3890068
HSA-MIR-31	1.8097683
HSA-MIR-224	1.6650244
HSA-MIR-451	1.6289215
HSA-MIR-10A	1.4811346
HSA-MIR-10B	1.4523870
HSA-MIR-31*	1.4370410
HSA-MIR-363	1.3751314
HSA-MIR-96	1.3709655
HSA-MIR-203	1.3489193
HSA-MIR-494	1.2932659

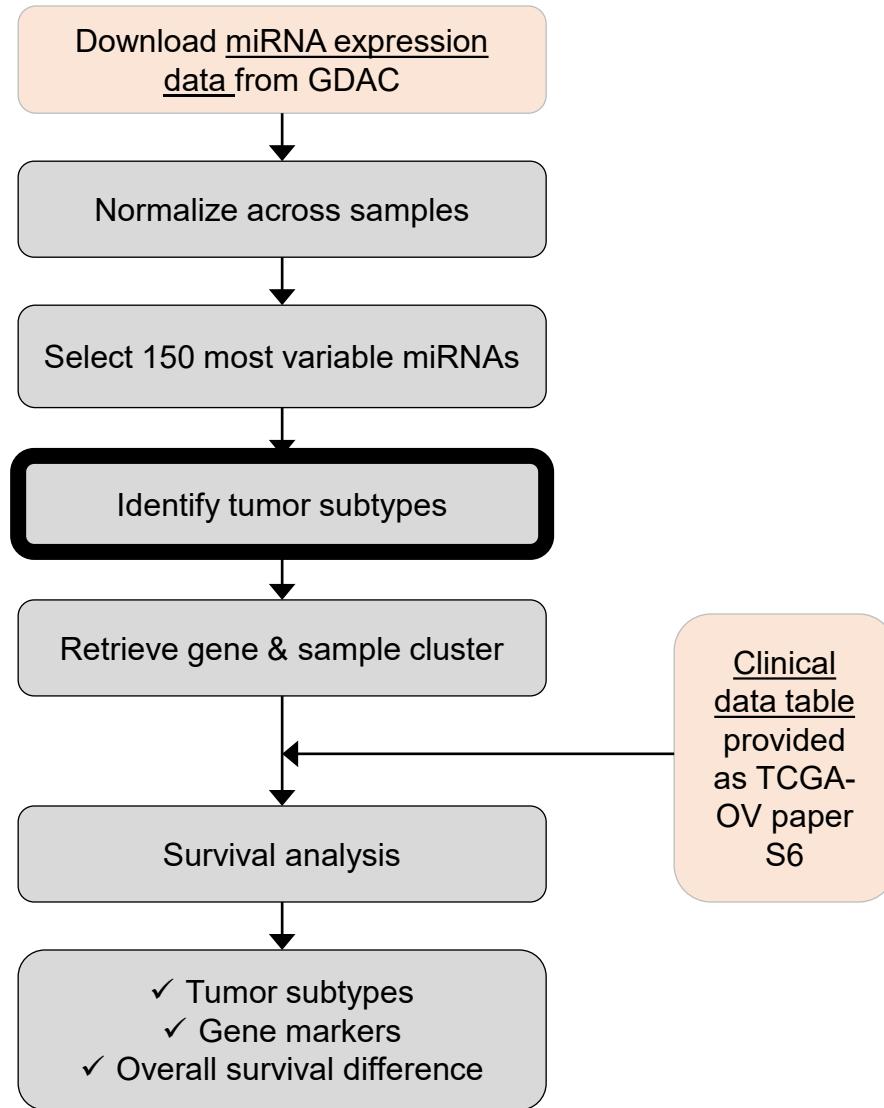
Showing 1 to 13 of 150 entries



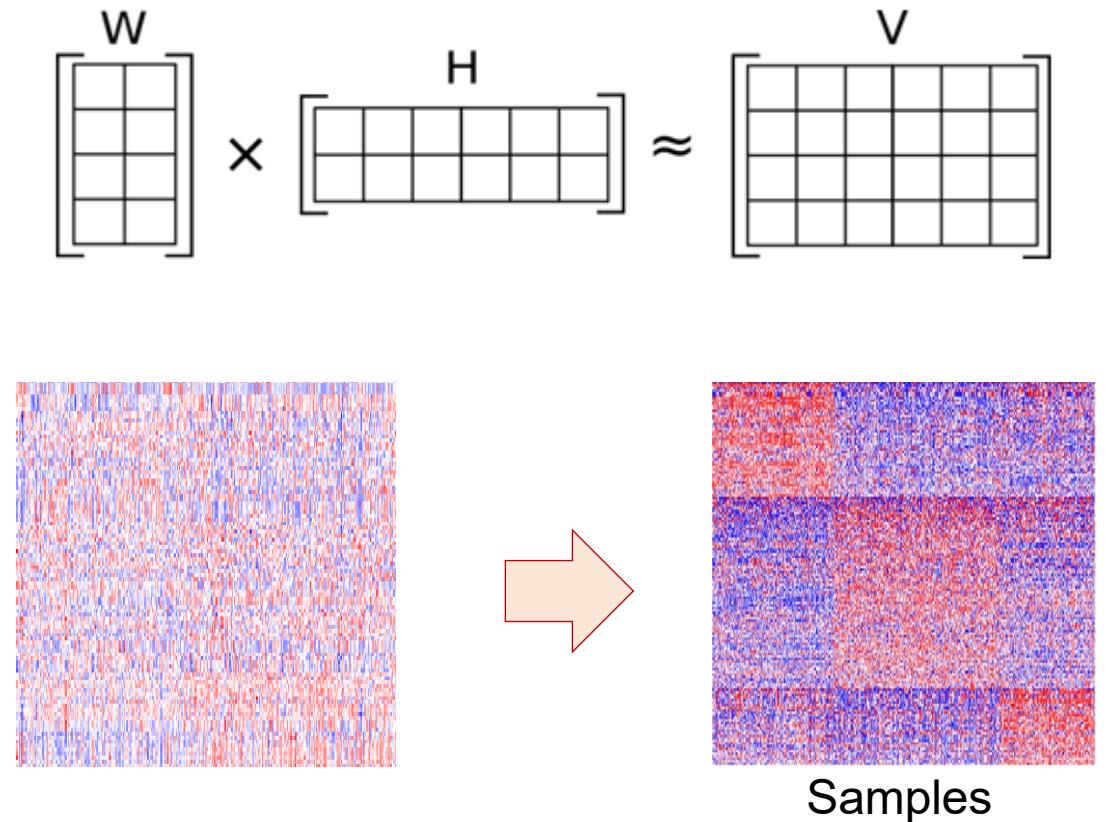
Median absolute deviation (MAD)

$$\text{MAD} = \text{median}_i ( |X_i - \text{median}_j(X_j)| )$$

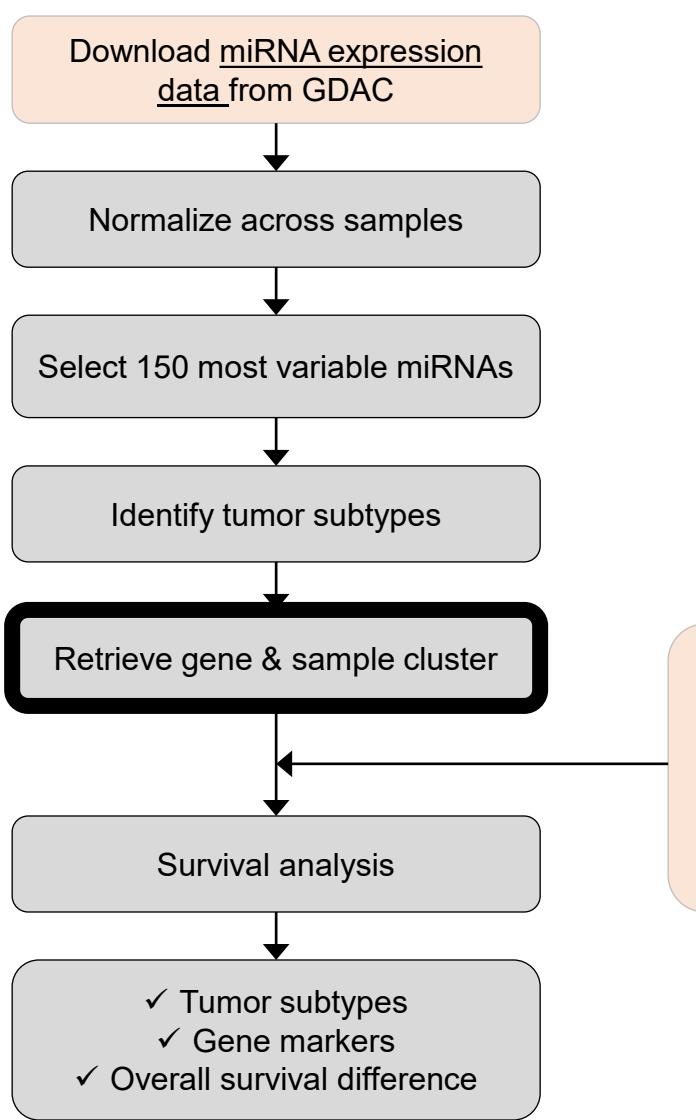
# Integrate Expression with Clinical Data



Non-negative matrix factorization (NMF), also non-negative matrix approximation is a group of algorithms in multivariate analysis and linear algebra where a matrix  $V$  is factorized into (usually) two matrices  $W$  and  $H$ , with the property that all three matrices have no negative elements.



# Integrate Expression with Clinical Data



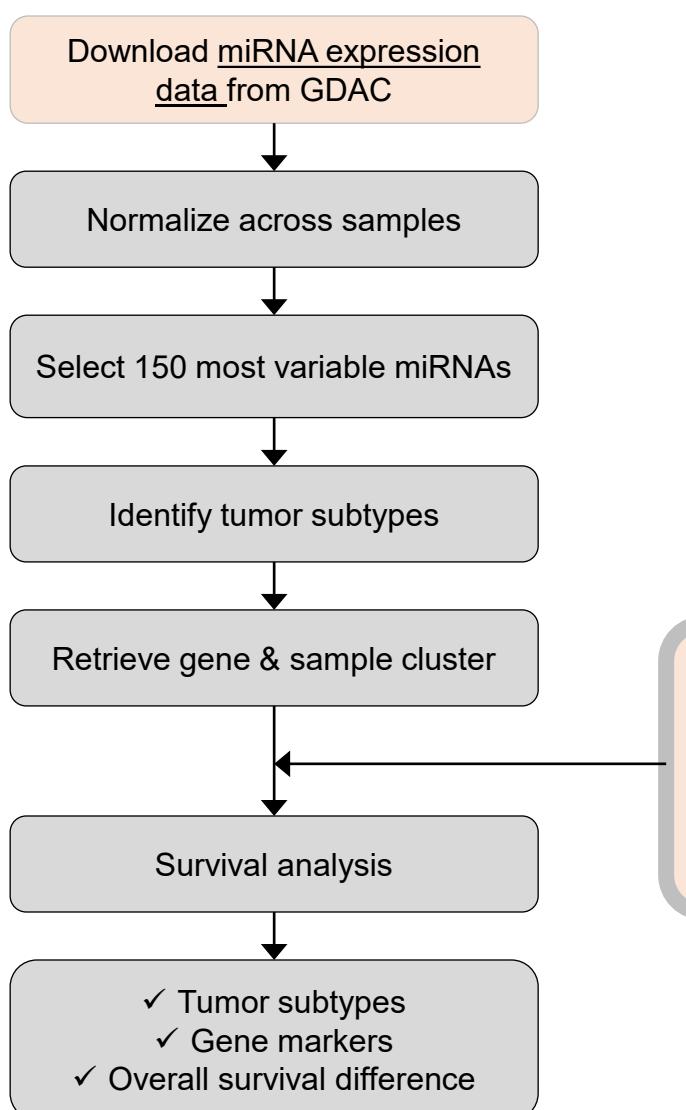
sample	cluster
TCGA.04.1331	1
TCGA.04.1338	1
TCGA.04.1341	1
TCGA.04.1343	1
TCGA.04.1348	1
TCGA.04.1362	1
TCGA.04.1365	1
TCGA.04.1530	1
TCGA.04.1542	1
TCGA.04.1648	1
TCGA.04.1649	1
TCGA.04.1651	1

Sample clusters

gene	cluster
HSA-MIR-205	1
HSA-MIR-494	1
HSA-MIR-144	1
HSA-MIR-142-5P	1
HSA-MIR-181A	1
HSA-MIR-151-3P	1
HSA-MIR-1225-5P	1
HSA-MIR-222	1
HSA-MIR-638	1
EBV-MIR-BART19-3P	1
HSA-MIR-21*	1
HSA-MIR-630	1

Gene clusters

# Integrate Expression with Clinical Data

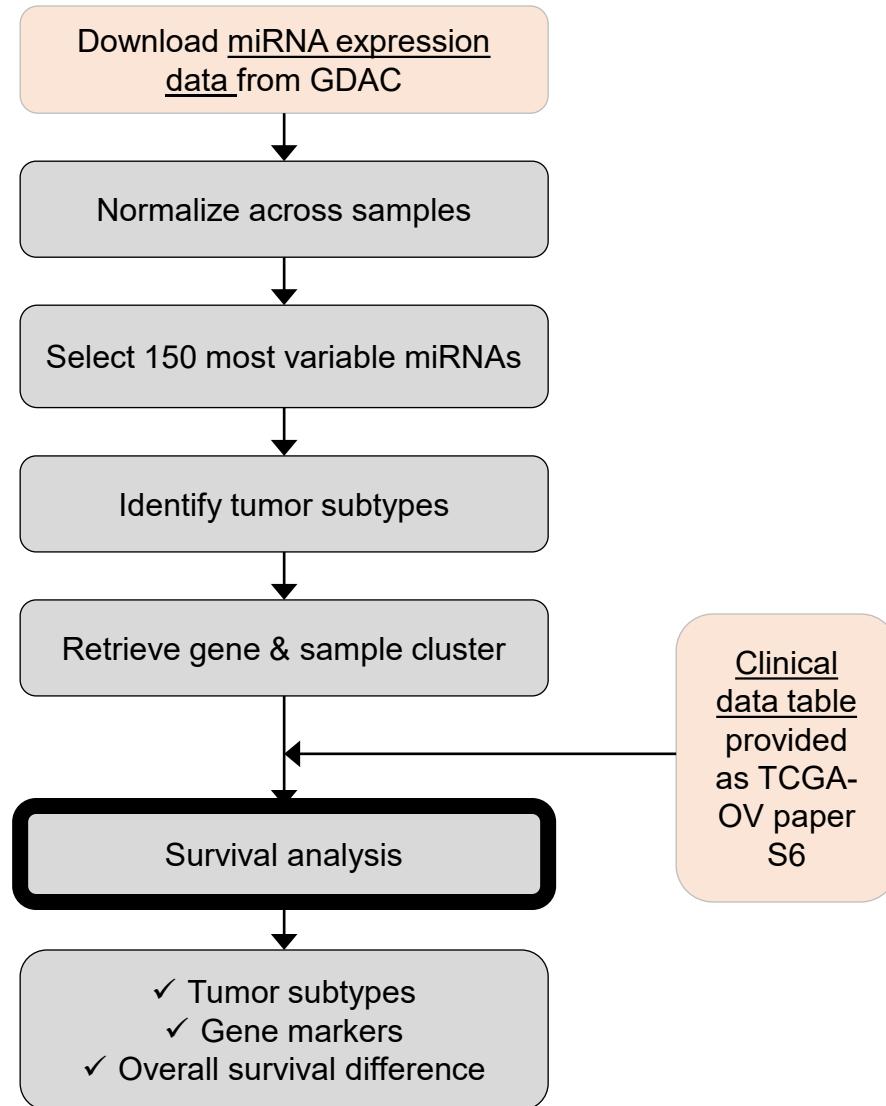


	sample	vital.status	overall.survival.month	age.at.diagnosis.year	tum.or.stage	tum.or.grade
1	TCGA.O4.1331	DECEASED	43.80	79.04	IIIC	G3
5	TCGA.O4.1338	LIVING	46.49	78.87	IIIC	G3
6	TCGA.O4.1341	LIVING	NA	85.52	NA	G3
8	TCGA.O4.1343	DECEASED	11.84	72.41	IV	G3
11	TCGA.O4.1348	DECEASED	48.62	44.48	IIIB	G3
17	TCGA.O4.1362	DECEASED	44.20	59.58	IIC	G3
19	TCGA.O4.1365	LIVING	76.33	87.47	IIIB	G3
25	TCGA.O4.1530	DECEASED	118.75	68.53	IIIC	G3
26	TCGA.O4.1542	DECEASED	83.97	52.78	IIIB	G2
29	TCGA.O4.1648	DECEASED	28.56	57.84	IIIC	G2
30	TCGA.O4.1649	LIVING	64.46	74.42	IIIC	G3
31	TCGA.O4.1651	DECEASED	36.07	53.78	IIIC	G3

Showing 1 to 13 of 486 entries

- Sample ID
- Vital status
- Overall survival
- Age at diagnosis
- Tumor stage
- Tumor grade
- **Sample cluster** (from the previous step)

# Integrate Expression with Clinical Data



- Survival analysis

- Study the time between entry to a study or an event (such as death)
- Calculate survival/risk difference and detect significance between groups
- Build models to predict prognosis

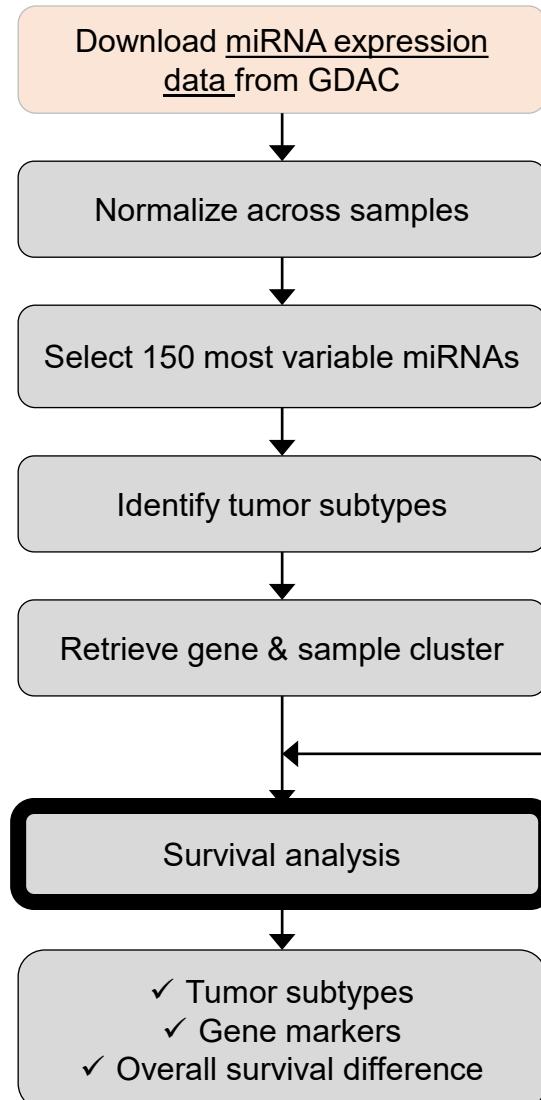
- Survival data

- Time to event (in year, month, etc.)
- Status (whether the event has happened?)
  - Censoring: only some individuals have experienced the event by the last follow up, while for others, the time is unknown
- “cumulative” survival time

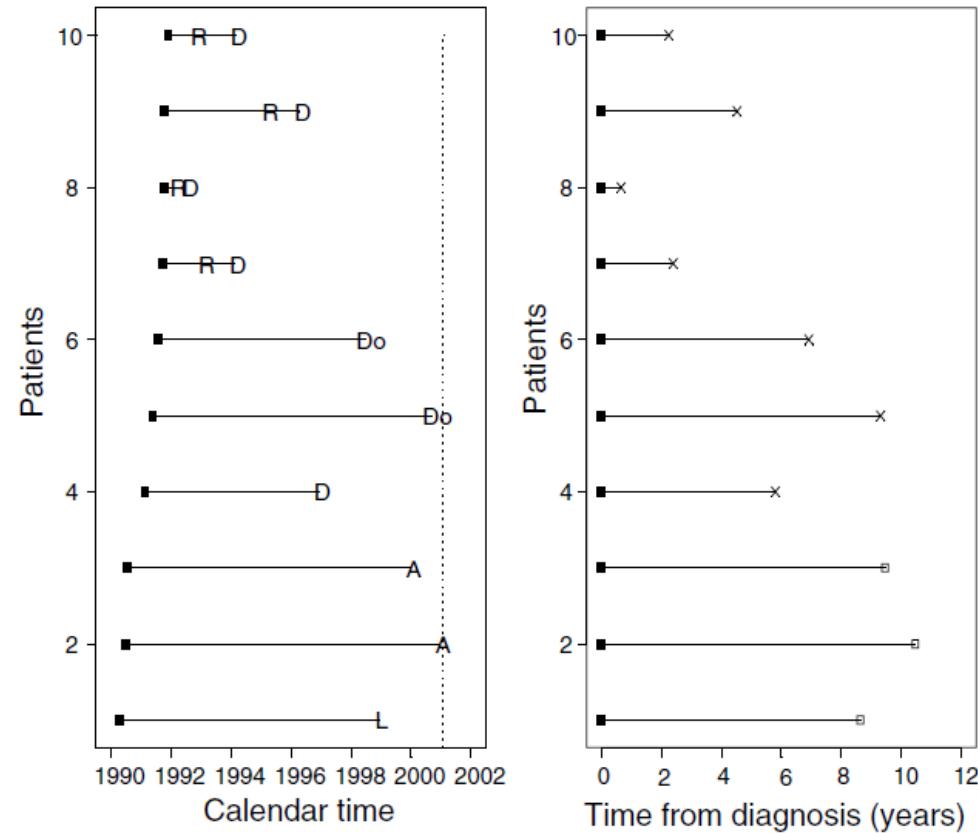
- Survival methods

- Kaplan-Meier estimator
- Log-rank test (Mantel-Haenzel test)
- Cox regression model (proportional hazard model)

# Integrate Expression with Clinical Data

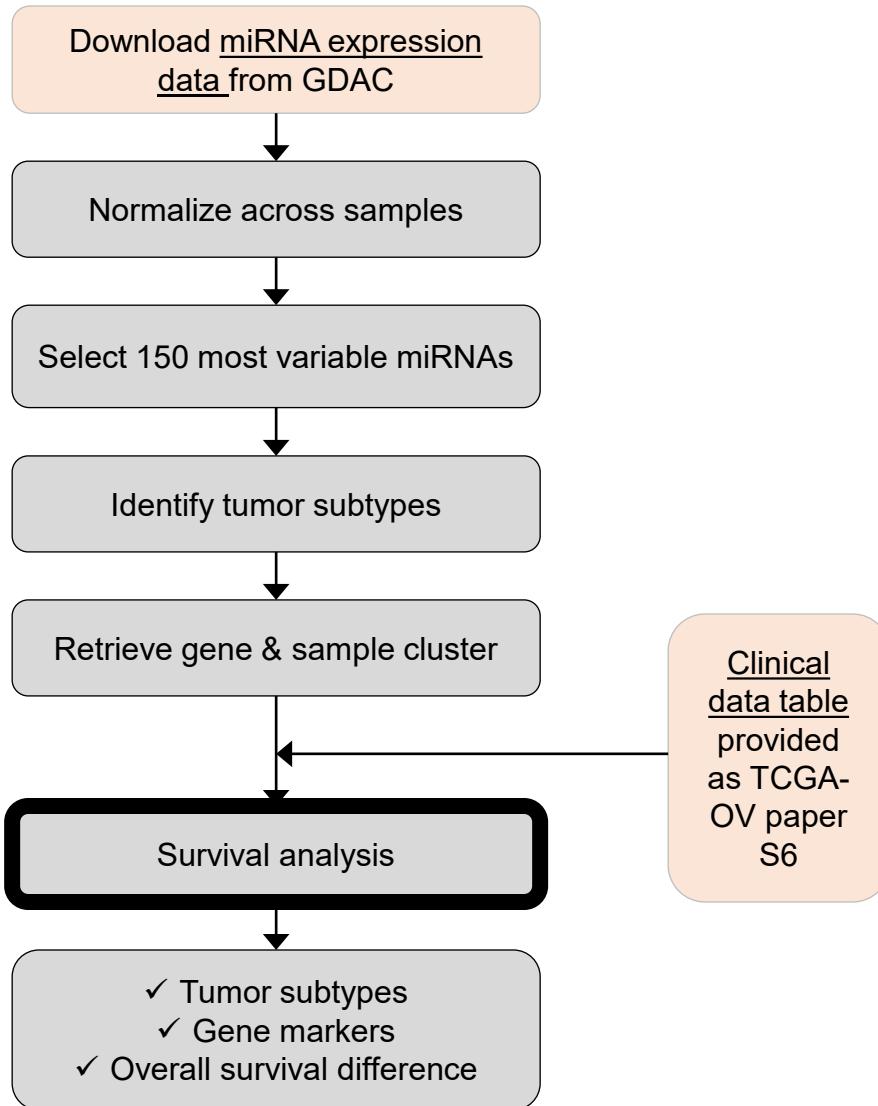


Right Censoring

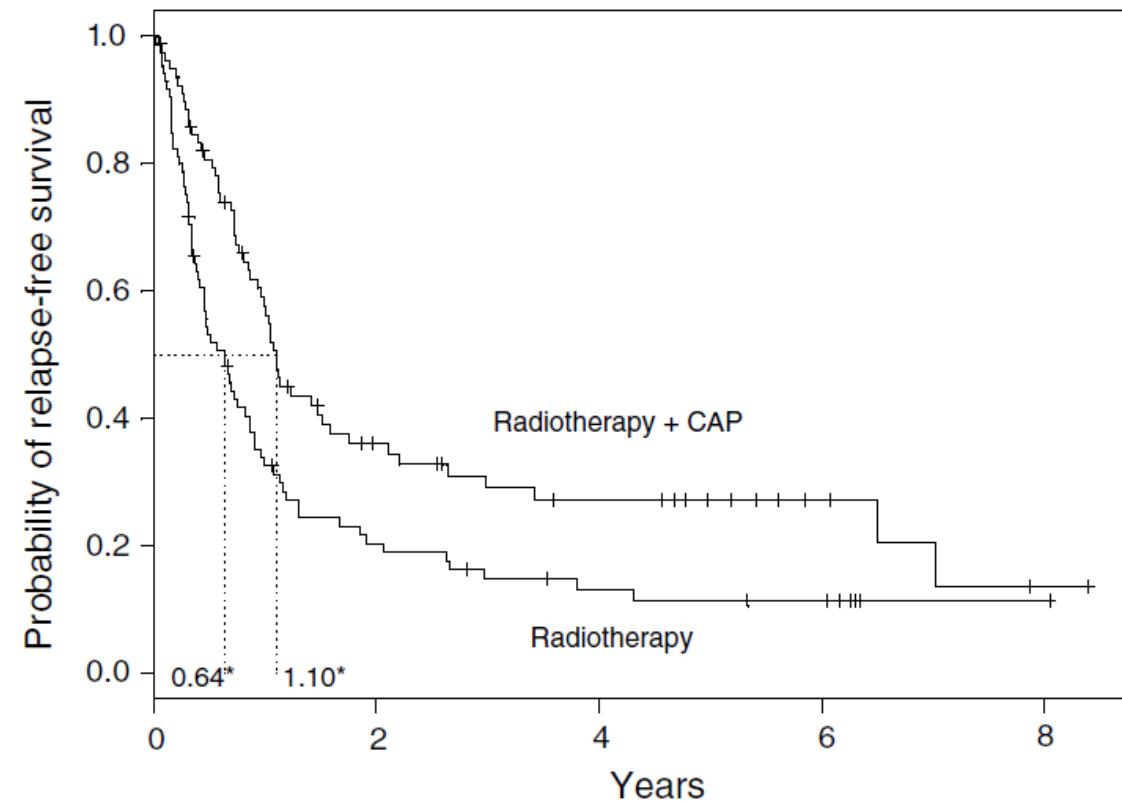


**Figure 1** Converting calendar time in the ovarian cancer study to a survival analysis format. Dashed vertical line is the date of the last follow-up, R = relapse, D = death from ovarian cancer, Do = death from other cause, A = attended last clinic visit (alive), L = loss to follow-up, X = death, □ = censored.

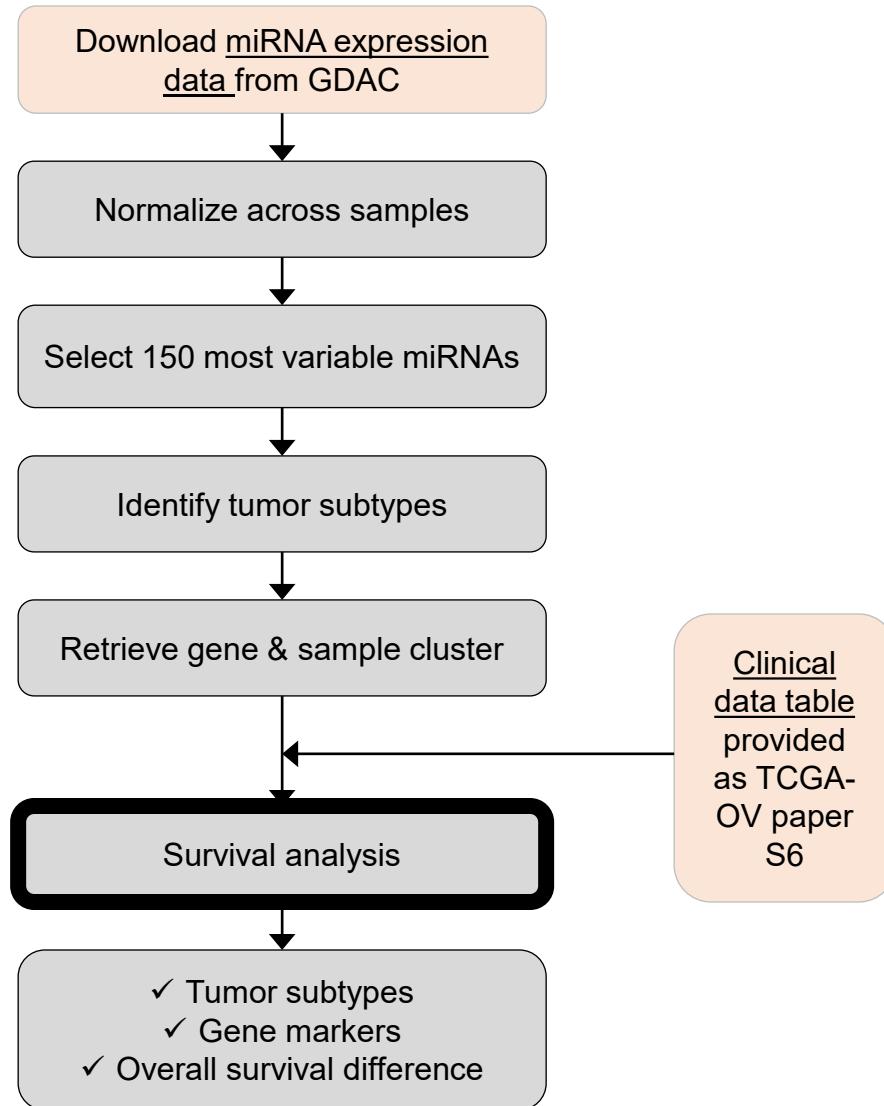
# Integrate Expression with Clinical Data



- **Kaplan-Meier estimator**
  - Stepwise function
  - Does not account for effect of other covariates (univariable test)



# Integrate Expression with Clinical Data



- Log-rank test
  - Chi-square test
  - Efficient in comparing groups differed by categorical variables, but not continuous ones
  - Univariable test

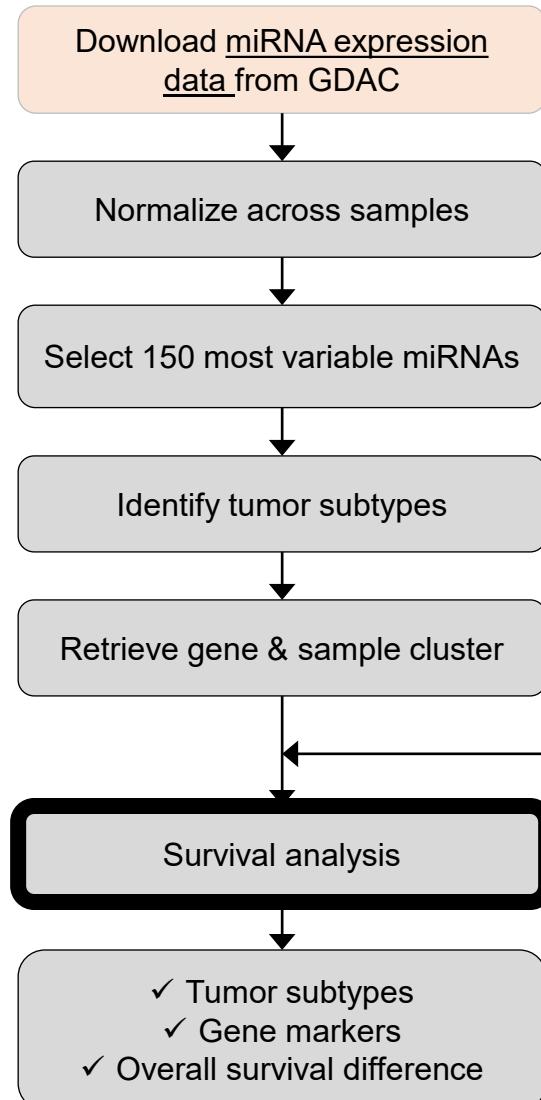
$$HR = \frac{O_1/E_1}{O_2/E_2}$$

The hazard ratio (HR) is a measure of the relative survival experience in the two groups

O: Observed

E: Expected (if no difference between group 1 and 2)

# Integrate Expression with Clinical Data



- Cox Proportional hazard model
  - Conveniently access the effect of continuous and categorical variables
  - Test the significance of factor of interest adjusting for other factors
  - Multivariable test!

```
S ~ sample.cluster +  
patient.age +  
tumor.grade
```

```
> summary(m3)  
Call:  
coxph(formula = surv ~ (clinical.sub$cluster + clinical.sub$age.at.diagnosis.year))  
  
n= 558, number of events= 292  
  
coef exp(coef) se(coef) z Pr(>|z|)  
clinical.sub$cluster2 0.451874 1.571254 0.145955 3.096 0.001962 **  
clinical.sub$cluster3 0.309421 1.362636 0.138750 2.230 0.025744 *  
clinical.sub$age.at.diagnosis.year 0.019574 1.019766 0.005354 3.656 0.000257 ***  
---  
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Cox Proportional hazard model

- $\beta$ : coefficient of explanatory variables or predictors
- $\exp(\beta)$ : the ratio of the hazards between two individuals whose values of  $x$  differ by one unit when all other covariates are held constant (**hazard ratio**, analogous to an **odds ratio** in the setting of multiple logistic regression analysis)
- Z: Wald statistics calculated by dividing  $\beta$  by its standard error
- P: P-value that corresponds to Z statistics. If  $P<0.05$ , then the null hypothesis of  $\beta$  equal to zero can be rejected at 95% confidence level

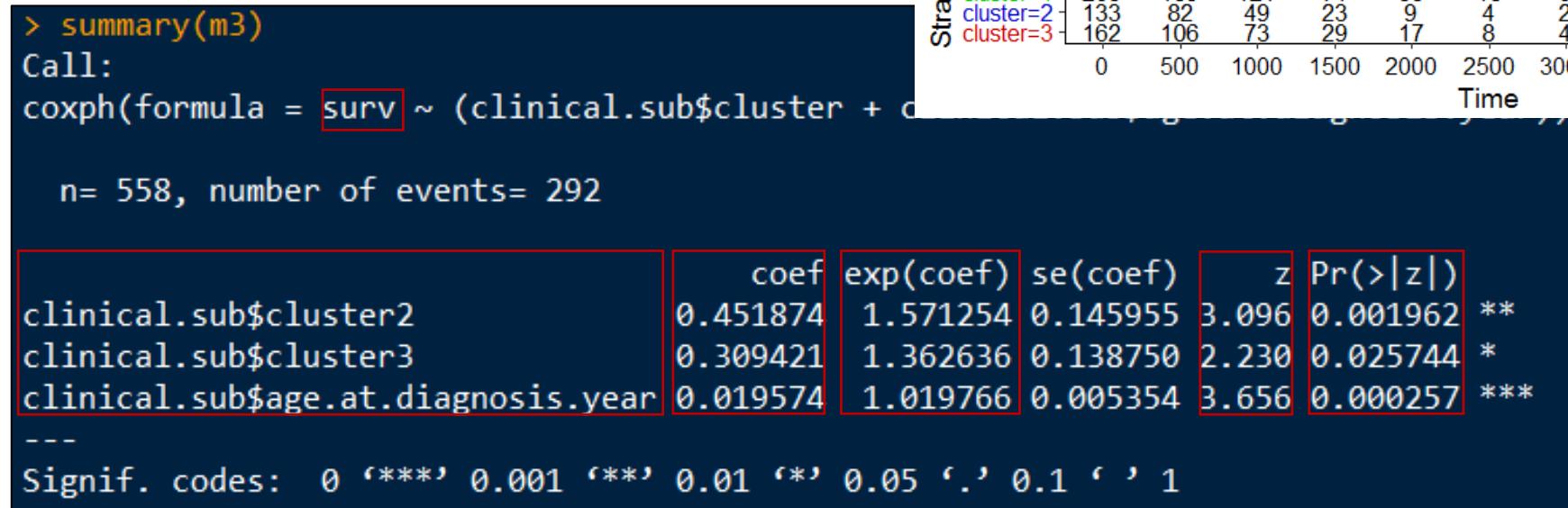
```
> summary(m3)
Call:
coxph(formula = surv ~ (clinical.sub$cluster + clinical.sub$age.at.diagnosis.year))

n= 558, number of events= 292

      coef  exp(coef) se(coef)   z Pr(>|z|)    
clinical.sub$cluster2 0.451874  1.571254 0.145955 3.096 0.001962 ***
clinical.sub$cluster3 0.309421  1.362636 0.138750 2.230 0.025744 *  
clinical.sub$age.at.diagnosis.year 0.019574  1.019766 0.005354 3.656 0.000257 ***
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

# Cox Proportional hazard model

- $\beta$ : coefficient of explanatory variables or predictors
- $\exp(\beta)$ : the ratio of the hazards between two individuals per unit when all other covariates are held constant (**hazard ratio** in the setting of multiple logistic regression analysis)
- Z: Wald statistics calculated by dividing  $\beta$  by its standard error
- P: P-value that corresponds to Z statistics. If P<0.05, the null hypothesis that the coefficient is equal to zero can be rejected at 95% confidence level



# Cox Proportional hazard model: visualize using forest plot

- $\beta$ : coefficient of explanatory variables
- $\exp(\beta)$ : the ratio of the hazards between a unit and a reference unit when all other covariates are held constant. This is called the **hazard ratio** in the setting of multiple logistic regression.
- Z: Wald statistics calculated by dividing  $\beta$  by its standard error.
- P: P-value that corresponds to Z statistic. A value that is less than or equal to zero can be rejected at 95% confidence level.

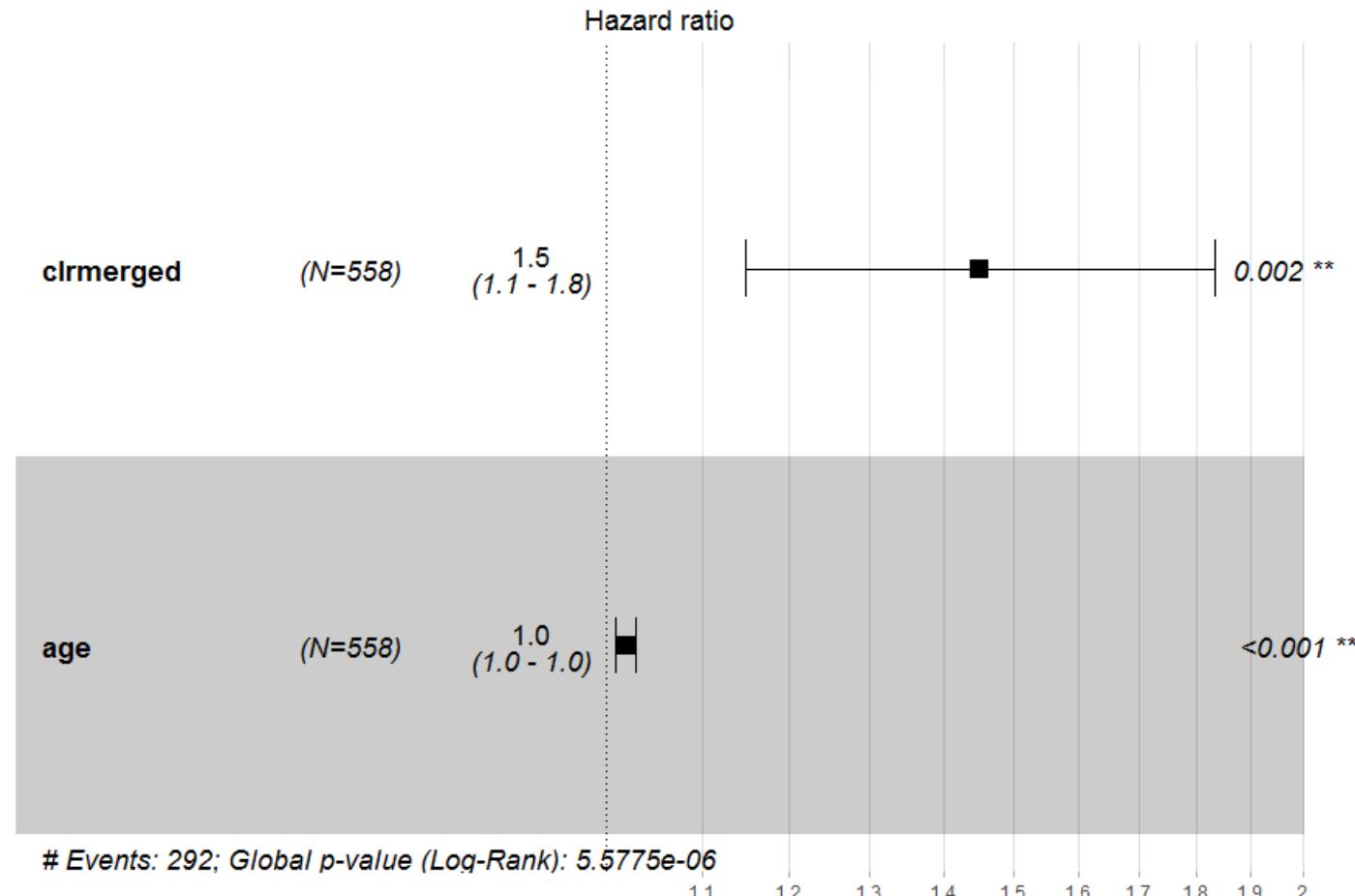
```
> summary(m3)
Call:
coxph(formula = surv ~ (clinical.sub$cluster2 + clinical.sub$cluster3 + clinical.sub$age.at.diagnosis.year))
n= 558, number of events= 292

# Events: 292; Global p-value (Log-Rank): 5.5775e-06

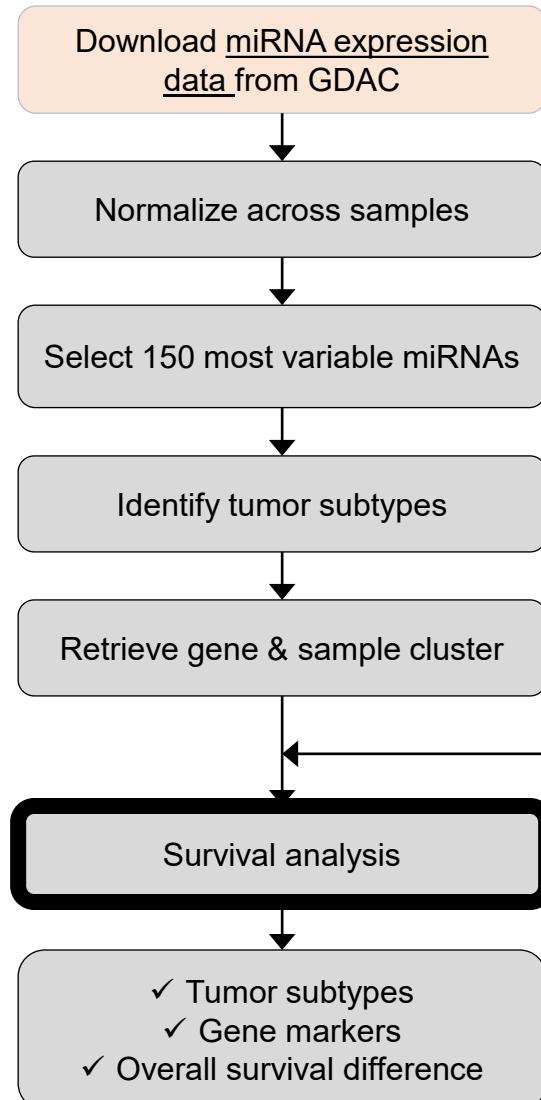
clirmerged      (N=558)      1.5
                           (1.1 - 1.8) 0.002 **

age              (N=558)      1.0
                           (1.0 - 1.0) <0.001 ***

Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```



# Integrate Expression with Clinical Data



- Survival methods

- Kaplan-Meier estimator
- Log-rank test (Mantel-Haenzel test)
- Cox regression model (proportional hazard model)

Library(survival)

survfit

survdiff

coxph

MSTP\_summer2022.hands on.Rmd

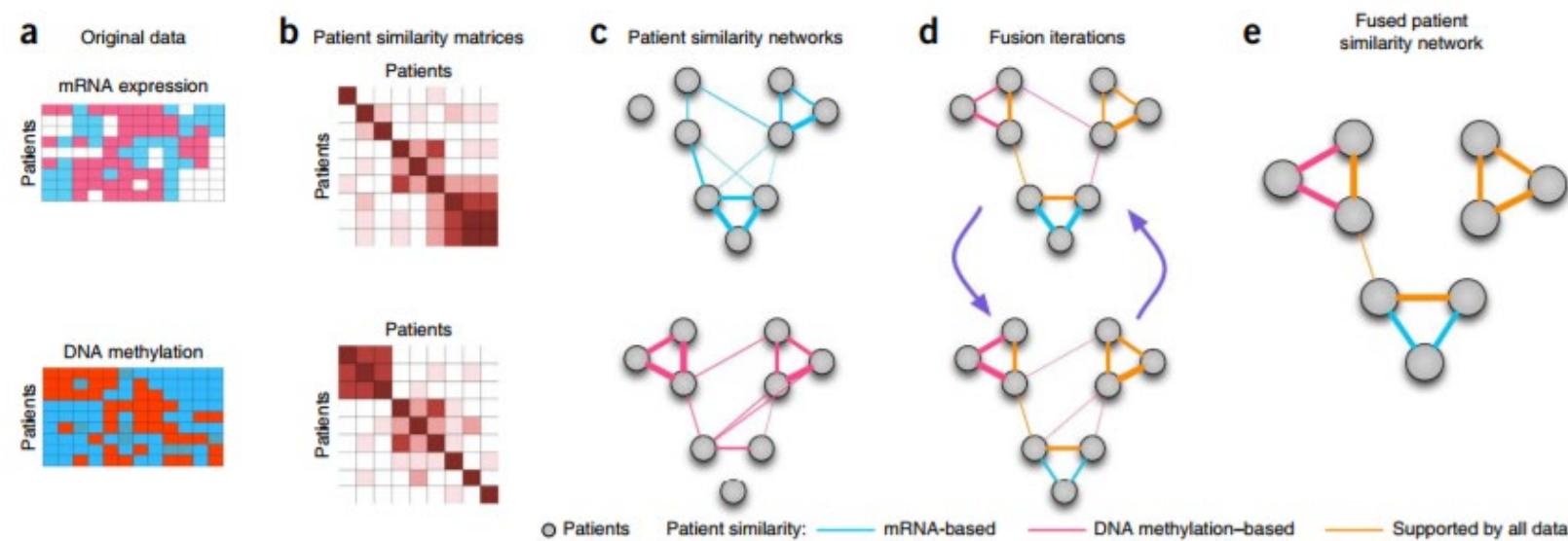
# Integration of multi-omics data

- Multiple types of genomic data, e.g. mRNA, miRNA, methylation
- SNF: similarity network fusion  
<http://compbio.cs.toronto.edu/SNF/SNF/Software.html>
- iClusterPlus: a joint latent variable model for integrative clustering  
<https://bioconductor.org/packages/release/bioc/html/iClusterPlus.html>

## Similarity network fusion for aggregating data types on a genomic scale

Bo Wang<sup>1,5</sup>, Aziz M Mezlini<sup>1,2</sup>, Feyyaz Demir<sup>1,2</sup>, Marc Fiume<sup>2</sup>, Zhuowen Tu<sup>3</sup>, Michael Brudno<sup>1,2</sup>, Benjamin Haibe-Kains<sup>4,5</sup> & Anna Goldenberg<sup>1,2</sup>

Recent technologies have made it cost-effective to collect diverse types of genome-wide data. Computational methods are needed to combine these data to create a comprehensive view of a given disease or a biological process. Similarity network fusion (SNF) solves this problem by constructing networks of samples (e.g., patients) for each available data type and then efficiently fusing these into one network that represents the full spectrum of underlying data. For example, to create a comprehensive view of a disease given a cohort of patients, SNF computes and fuses patient similarity networks obtained from each of their data types separately, taking advantage of the complementarity in the data. We used SNF to combine mRNA expression, DNA methylation and microRNA (miRNA) expression data for five cancer data sets. SNF substantially outperforms single data type analysis and established integrative approaches when identifying cancer subtypes and is effective for predicting survival.



# Thank you!



Questions



# Hands on practice!

