



Michigan State University
<http://binchenlab.org>

Big Data in Pharmacology

Bin Chen

Associate Professor

Dept. of Pediatrics and Human Development

Dept. of Pharmacology and Toxicology

College of Human Medicine

Michigan State University

Bin.Chen@hc.msu.edu @DrBinChen

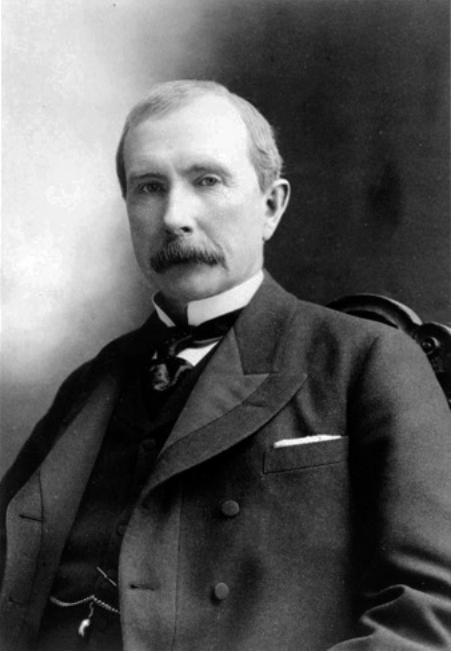
<http://binchenlab.org>

Learning objectives

- Understand pharmacology-related big data technologies and resources

Beyond Excel, GraphPad

Data = Oil



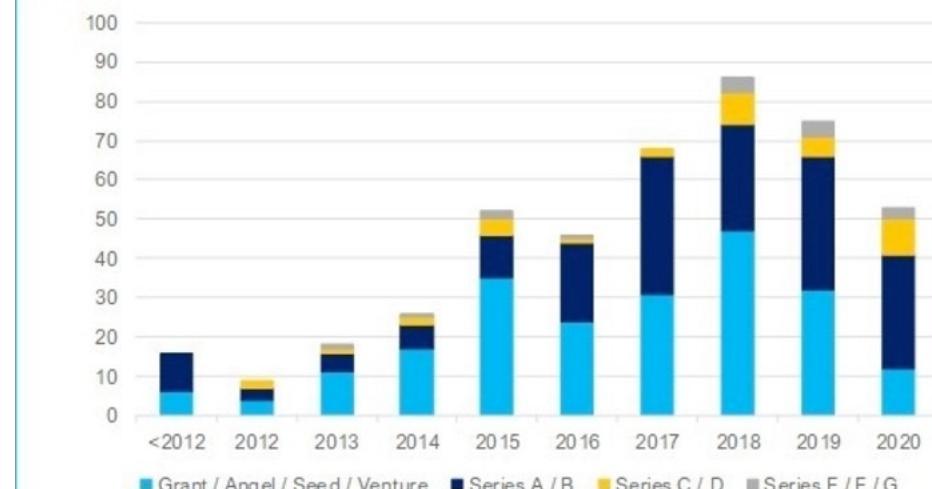
John D. Rockefeller in 1885



Company	Founding year	Investment (by 2021)
Insitro	2018	\$734M
XtalPi	2014	\$386M
Exscientia	2012	\$148M

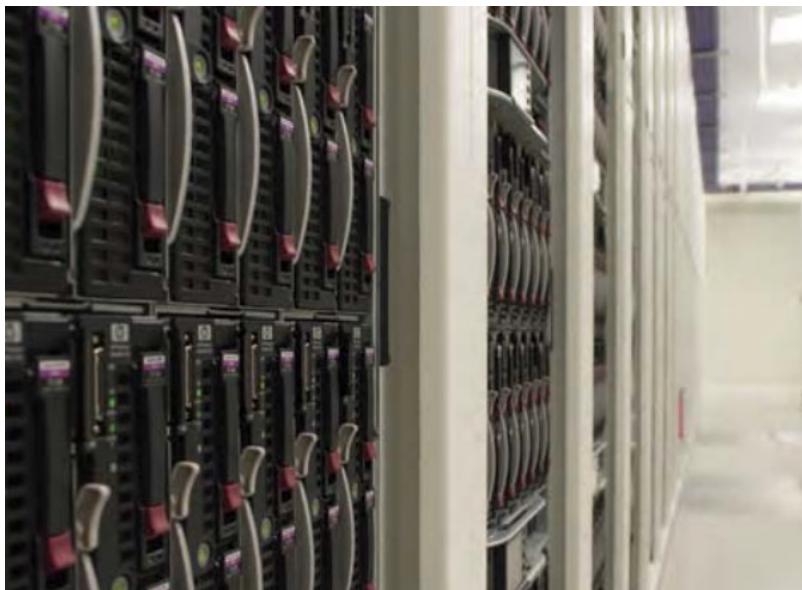
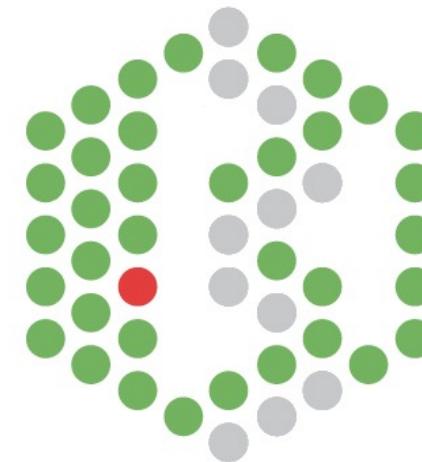
<https://www.nanalyze.com/2021/04/companies-ai-drug-discovery/>

Figure 2
AI in Drug Development Market - Fundings per Year
Number of Completed Funding Rounds



Source: Emersion Insights

EMBL-EBI



25 petabytes 2014

307 petabytes 2019

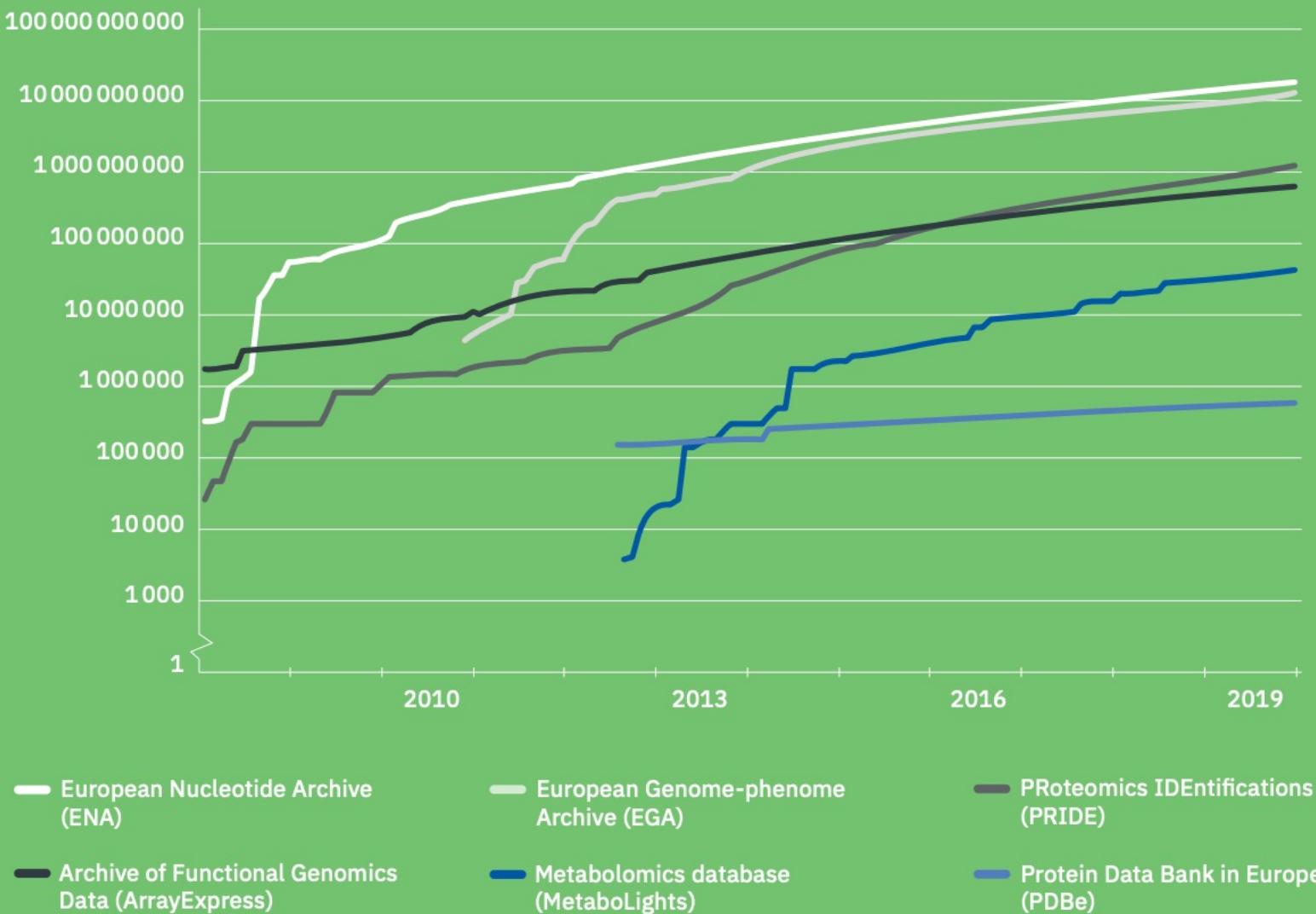
307 petabytes =



x 150,000

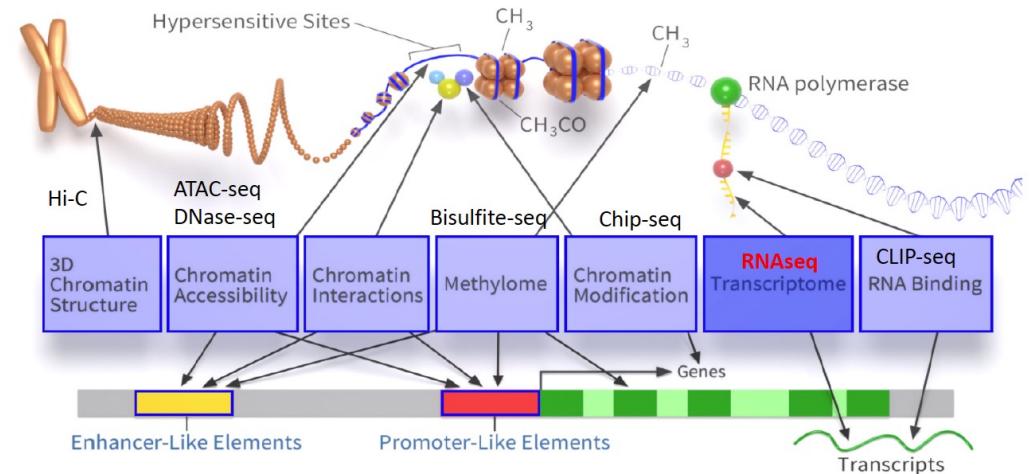
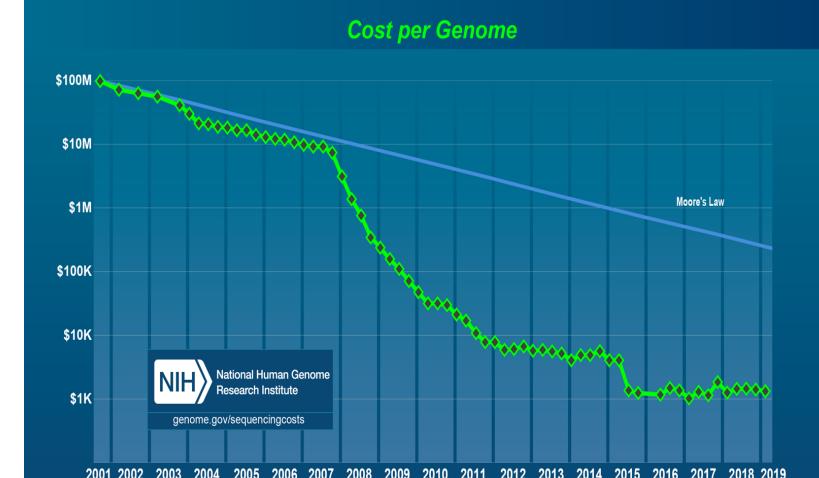
DATA GROWTH BY EMBL-EBI DATA RESOURCE

Volume of data (megabytes) per year (2008–2019)

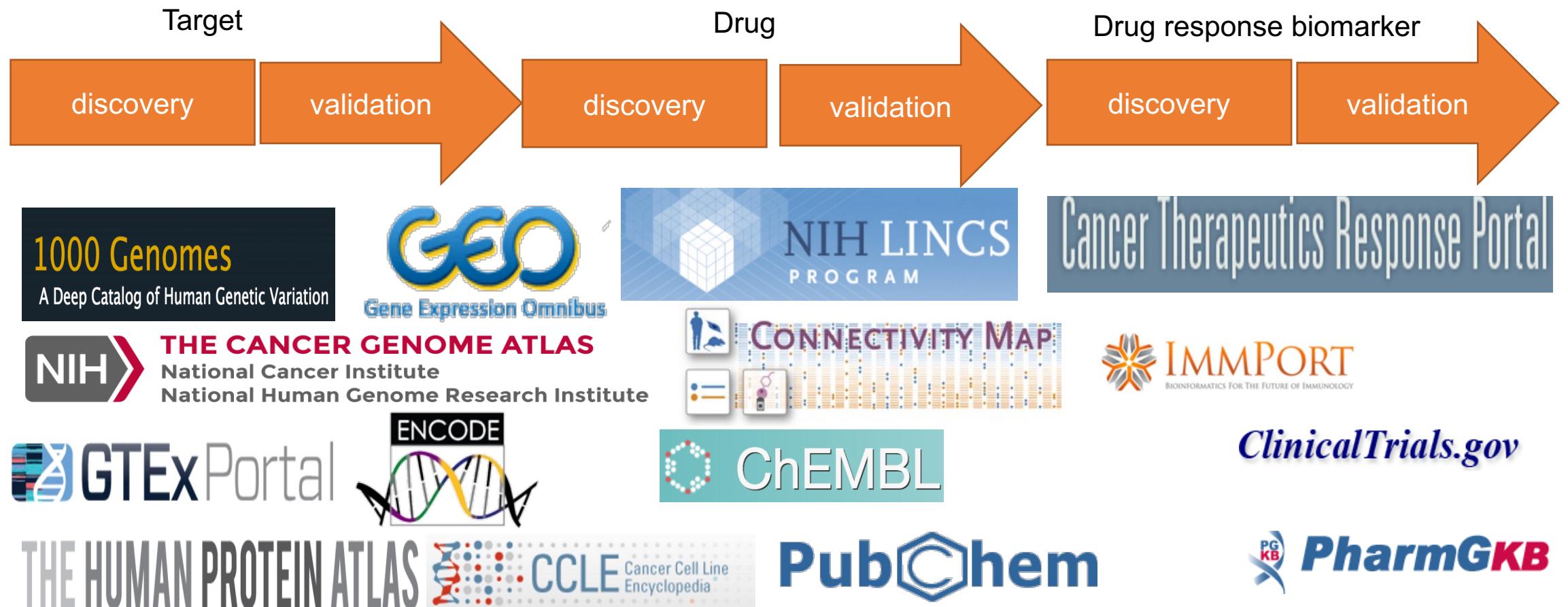


Why so much data available today?

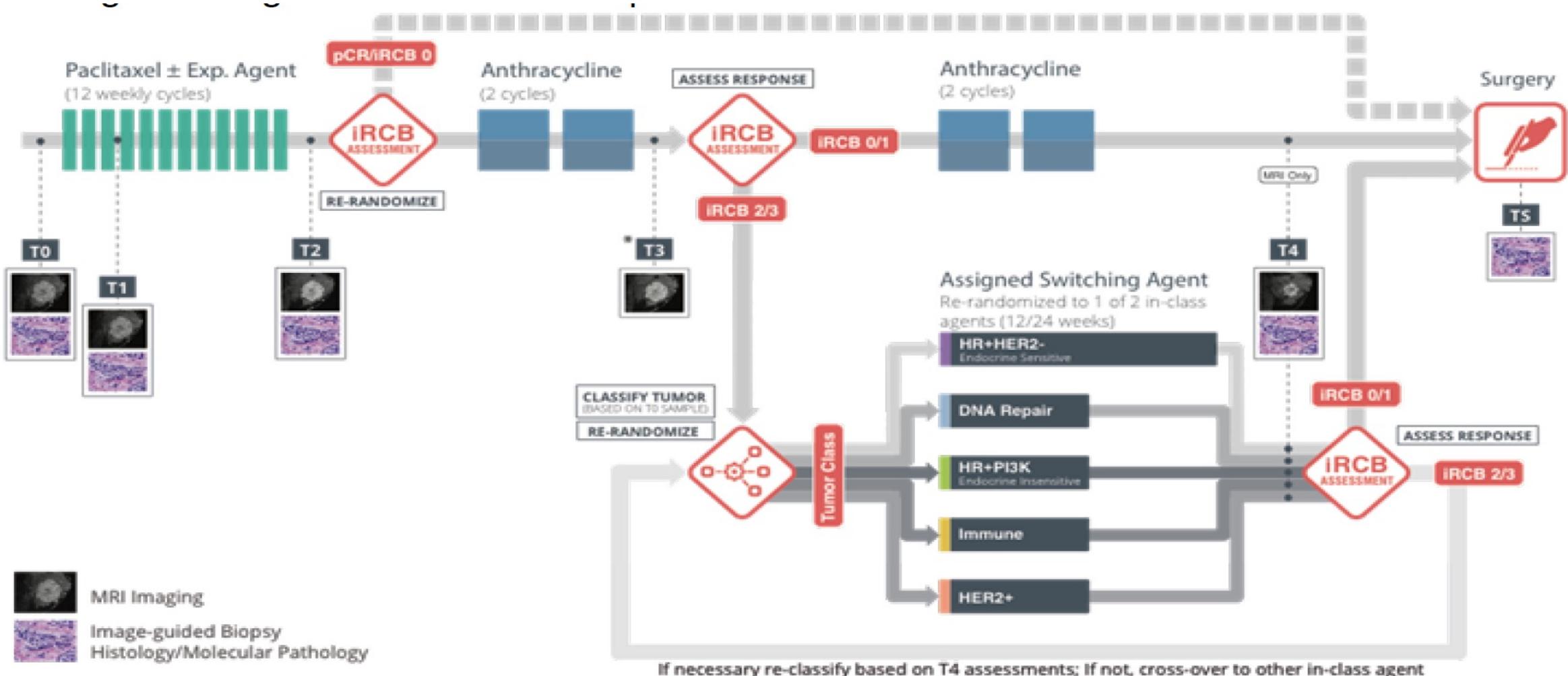
- Low sequencing costs
- High throughput technologies
- High performance computing



Big Data in drug discovery

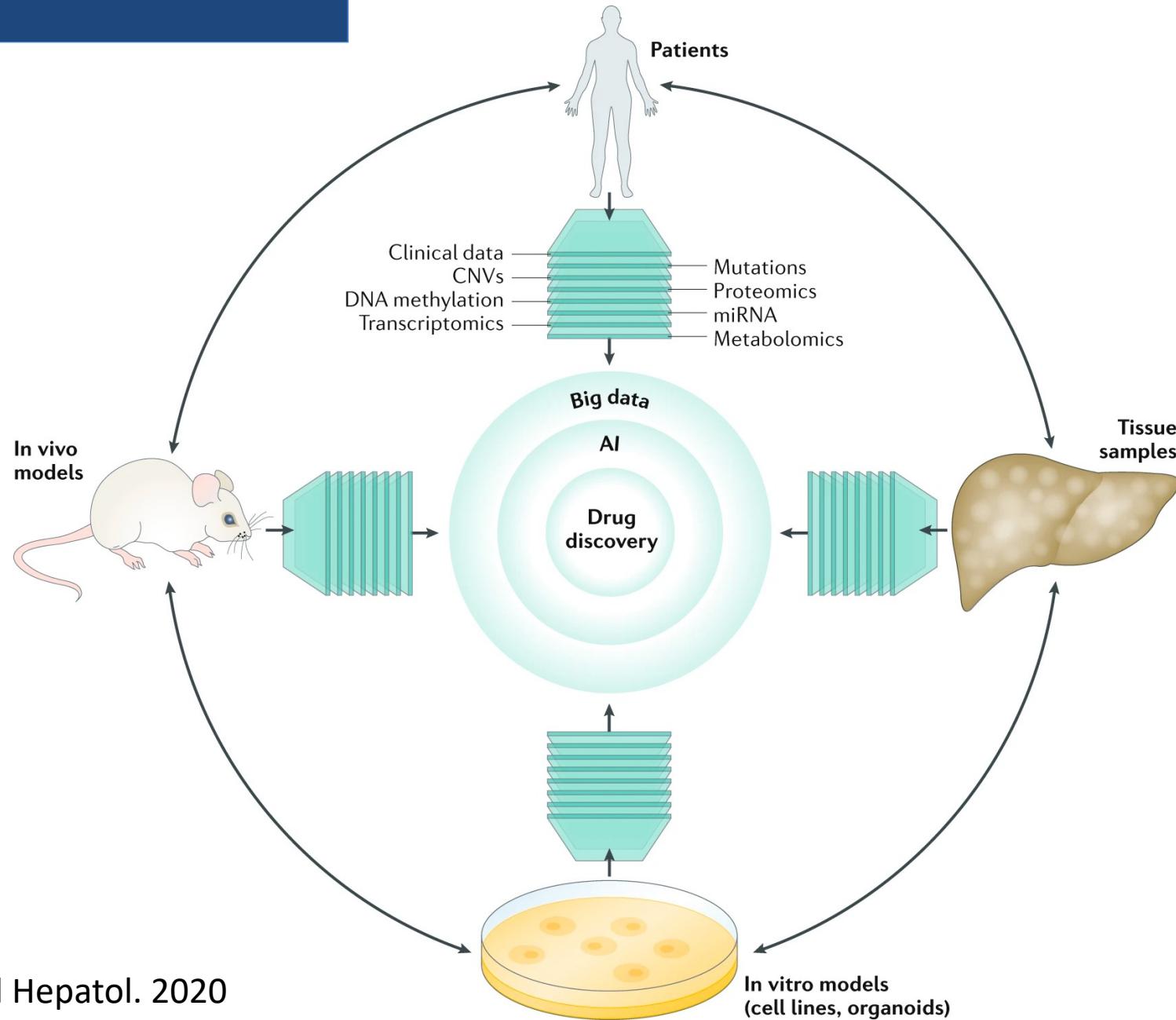


Big Data in the clinic



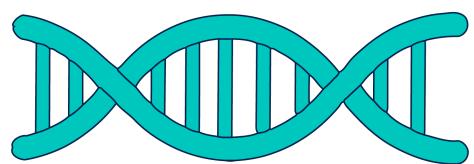
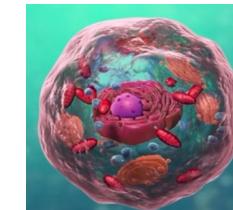
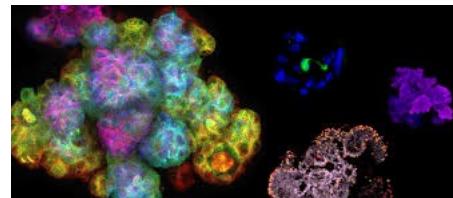
I-SPY 2 TRIAL: Neoadjuvant and Personalized Adaptive Novel Agents to Treat Breast Cancer (I-SPY 2)

Big Data in translational research

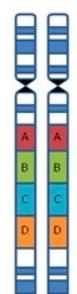


Big data you often encounter

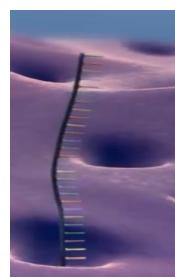
Molecular Profiling of biological systems with/without perturbation



DNA



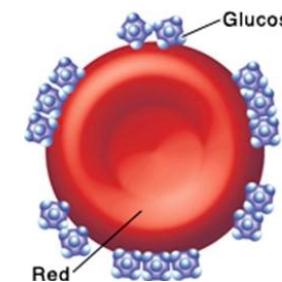
Chromosome
(copy number)



mRNA (gene)



Protein

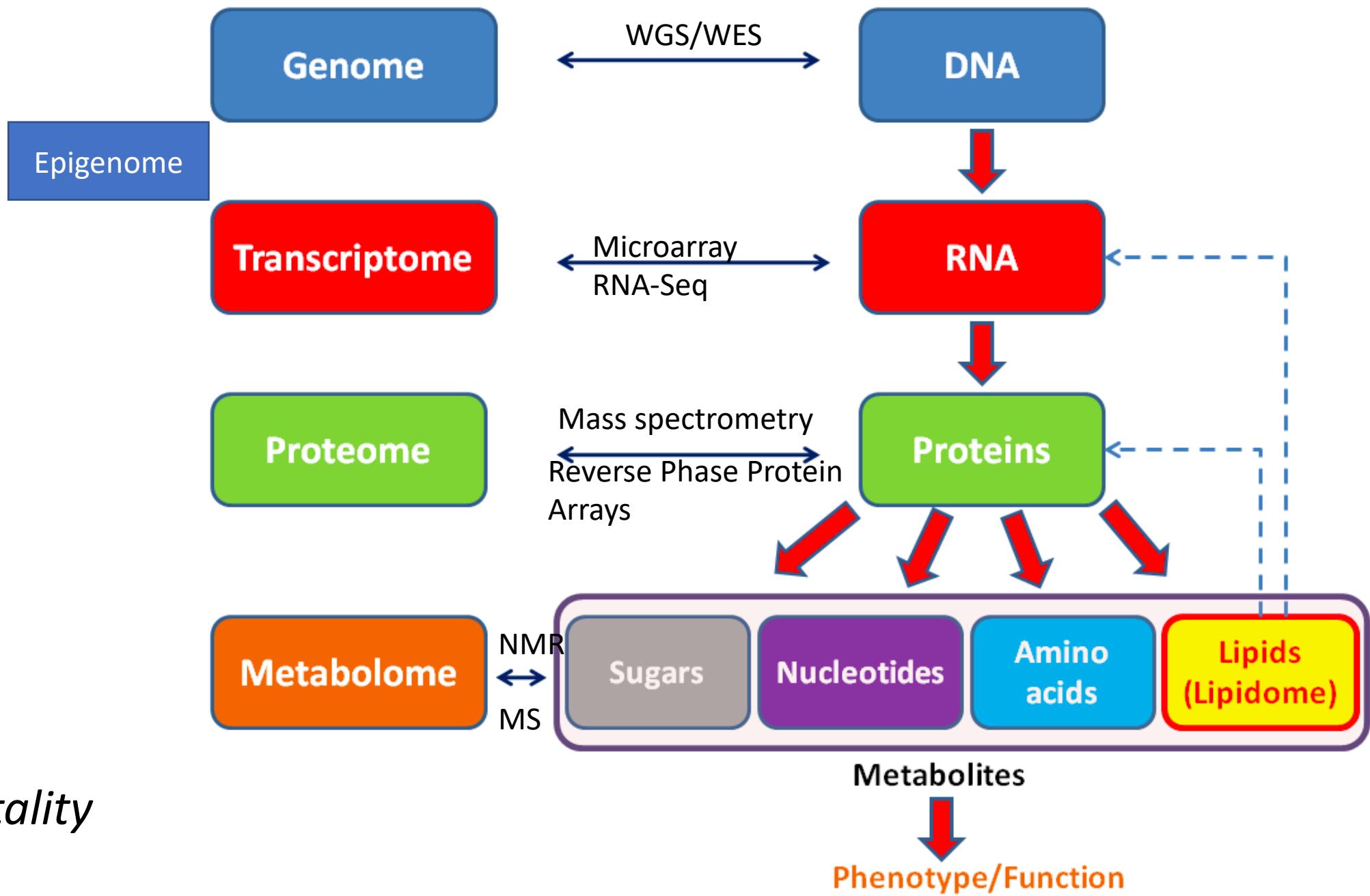


metabolite



With/without perturbation
(e.g., drug treatment, gene knock-down)

• • • •

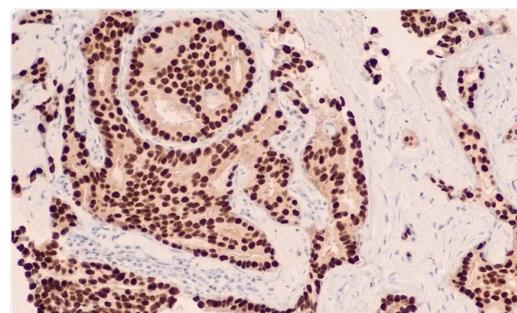
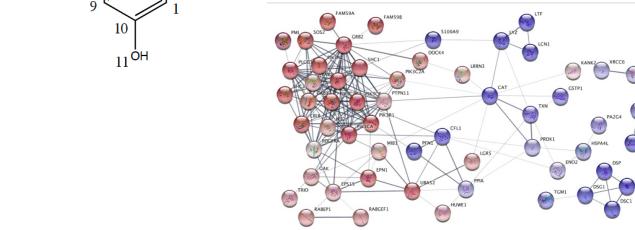
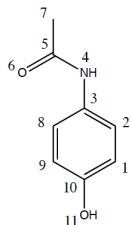
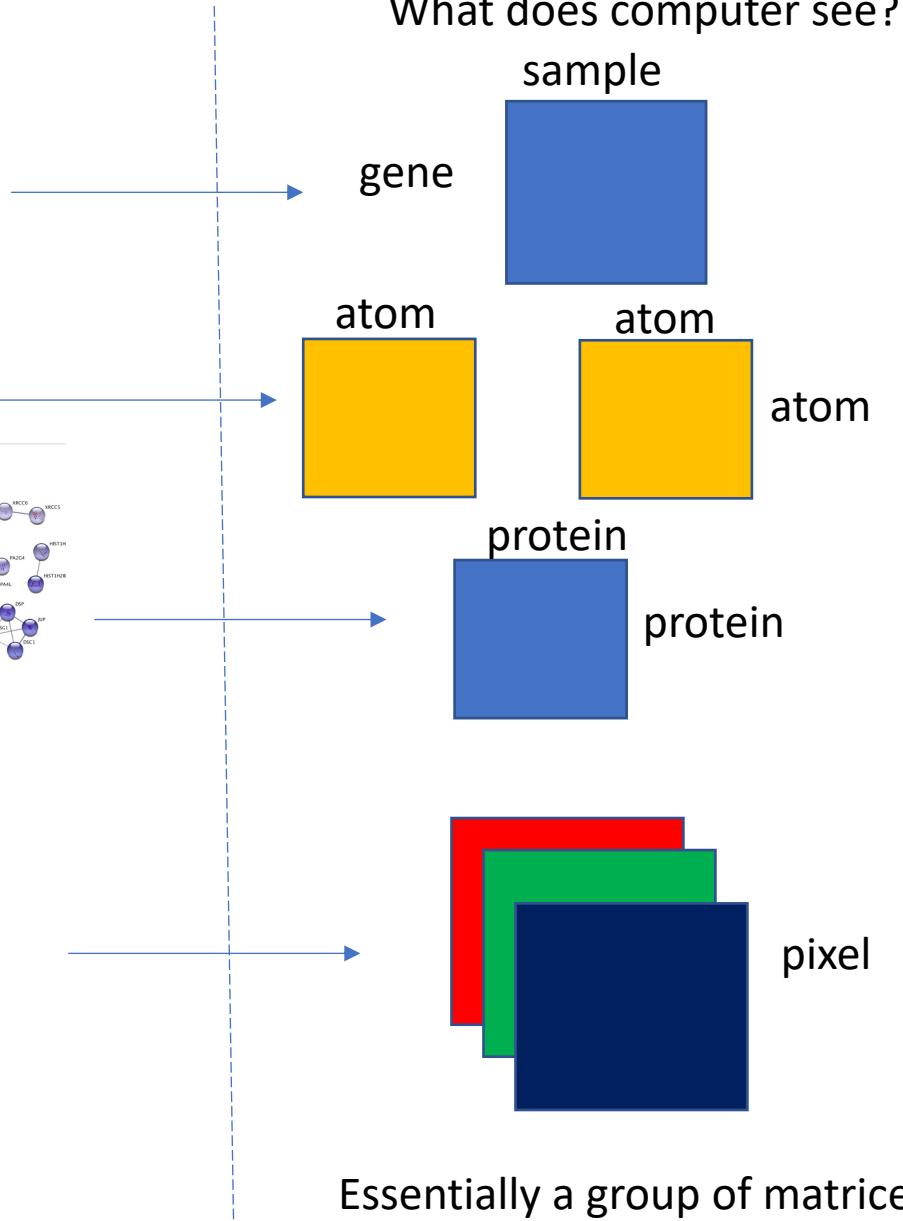
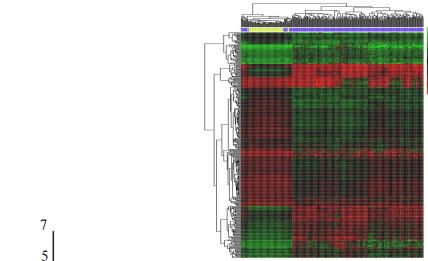


-Ome: *totality*

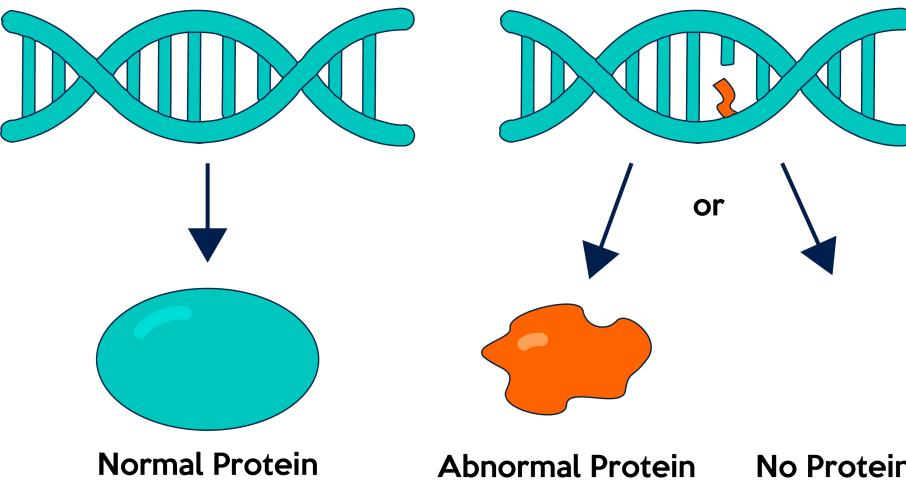
Commonly used big data types

- OMICS
 - Transcriptomics
 - Proteomics
- Graph
 - Chemical structures
 - Protein-Protein Interactions
- Image
 - MRI
 - IHC

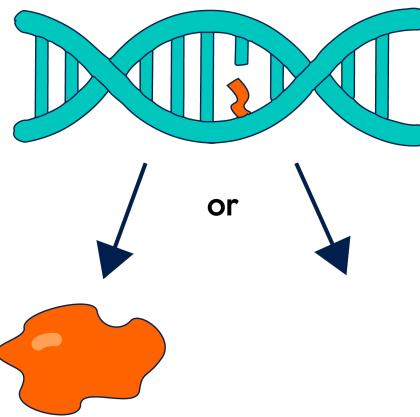
What do you see?



Normal Gene



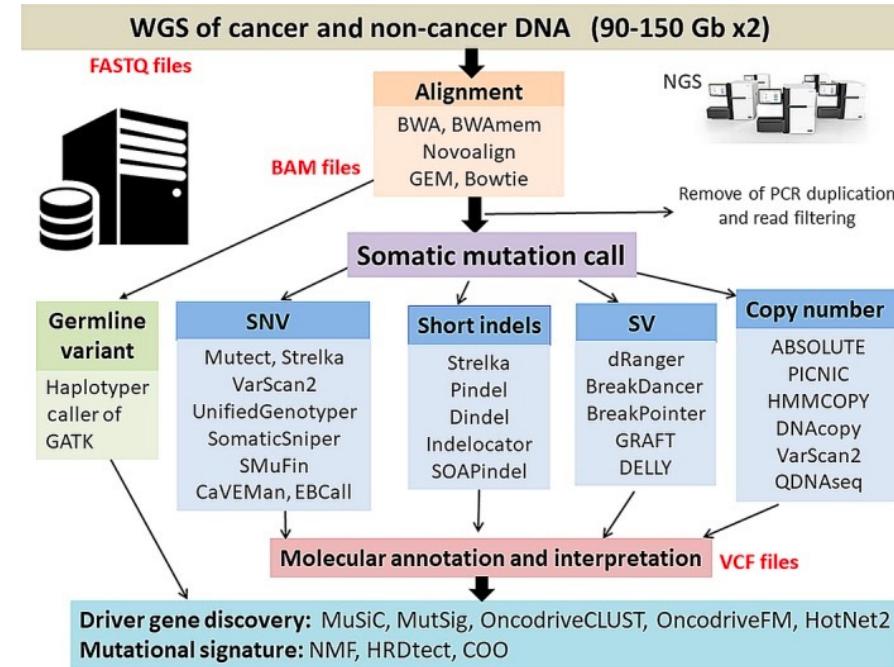
Mutated Gene



Normal Protein

Abnormal Protein

No Protein



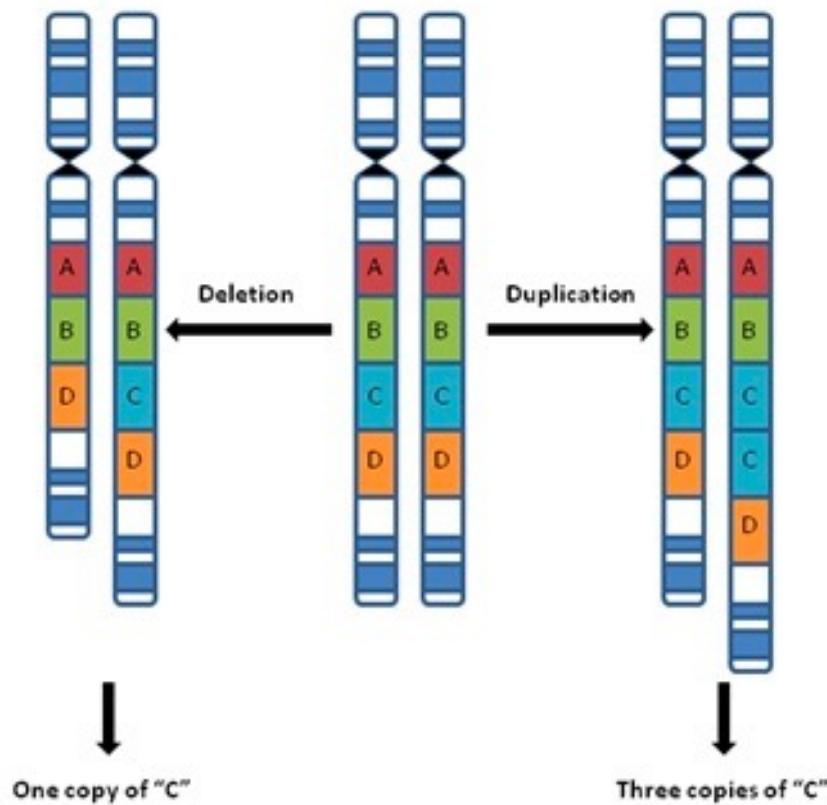
1656 cell lines

19K genes

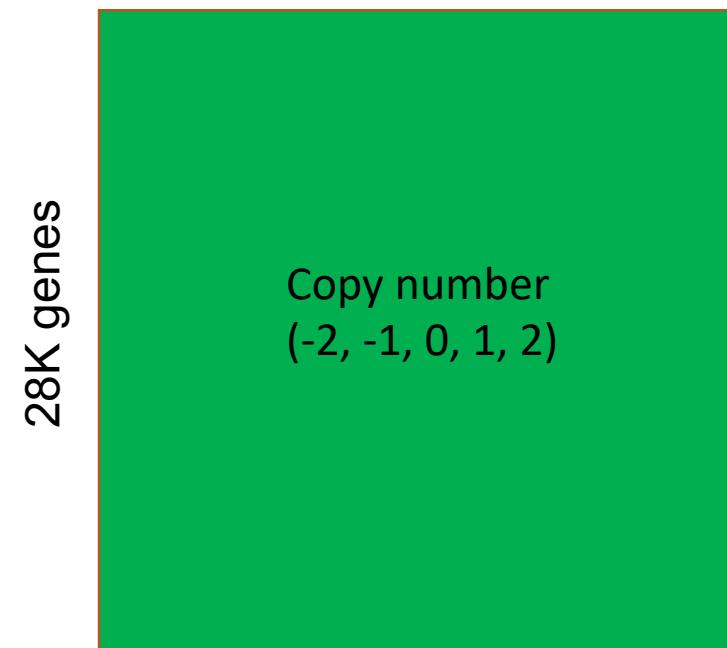
A	B	C	D	E	F	G	H	I	J	K	L	M	N	C P	Q	R	S
Hugo_Symbol	Entrez_Gene	NCBI_Build	Chromosome	Start_positio	End_positio	Strand	Variant_Classifi	Variant_Type	Reference_A	Tumor_Seq	dbSNP_RS	dbSNP_Val	Genome_Change	cDNA_Change	Codon_Change	Protein_Change	
VPS13D	55187	37	1	12359347	12359347	+	Nonsense_Mut	SNP	C	A	NA	NA	g.chr1:12359347C>A	A c.6122>A	c.(6121-6123)T>a aa	p.S2041*	
AADACL4	343066	37	1	12762308	12762322	+	In_Frame_Del	DEL	CTGGCGTG	-	rs58218424	byFrequency	g.chr1:12762308_12726	A c.786_800delCTGGCGTG	(784-801)tcctggcgta p.WRDA123del		
IFNL1R	163702	37	1	24484172	24484172	+	Silent	SNP	G	A			g.chr1:24484172G>A	A c.1011C>T	c.(1009-1011)ggC>ggT	p.G337G	
TMEM57	55219	37	1	25785018	25785019	+	Frame_Shift_Ins	INS	-	A			g.chr1:25785018_25785	A c.789_790insA	c.(790-792)aaafs	p.K264fs	
ZSCAN20	7579	37	1	33954141	33954141	+	Missense_Mut	SNP	T	G	NA	NA	g.chr1:33954141T>G	A c.4947>T	c.(493-495)gtggG ggG	p.V165G	
POU3F1	5453	37	1	38512139	38512139	+	Missense_Mut	SNP	C	G	NA	NA	g.chr1:38512139C>G	A c.277G>C	c.(277-279)Gcc>Gc	p.G93R	
MAST2	23139	37	1	46498028	46498028	+	Silent	SNP	C	T			g.chr1:46498028C>T	A c.3366C>T	c.(3364-3366)cgC>cgt	p.R1122R	
GBP4	115361	37	1	89657103	89657103	+	Silent	SNP	G	T	rs146643676		g.chr1:89657103G>T	A c.757C>A	c.(757-759)Cgg>Agg	p.R253R	
VAV3	10451	37	1	1.08E+08	1.08E+08	+	Splice_Site	SNP	A	G			g.chr1:108427170A>G	A CH-000001	c.e17+1		
NBPFL20	1E+08	37	1	1.48E+08	1.48E+08	+	Missense_Mut	SNP	T	G	rs36922345	NA	g.chr1:148346689G>G	A c.68A>C	c.(67-69)aAa>aCa	p.K23T	
FLG2	388698	37	1	1.52E+08	1.52E+08	+	Missense_Mut	SNP	C	T			g.chr1:152331291C>T	A c.70G>A	c.(70-72)Gag>Aag	p.E24K	
RHBG	57127	37	1	1.56E+08	1.56E+08	+	Frame_Shift_De	DEL	C	-	NA	NA	g.chr1:156339132delC	A c.92delC	c.(91-93)gcfts	p.A31fs	
FCRL5	83416	37	1	1.57E+08	1.57E+08	+	Missense_Mut	SNP	A	T			g.chr1:157497658A>T	A c.17097A>T	c.(1708-1710)Tcc>Ac	p.L570H	
SPTA1	6708	37	1	1.59E+08	1.59E+08	+	Missense_Mut	SNP	C	A	NA	NA	g.chr1:158639228C>A	A c.1803G>T	c.(1801-1803)aaG>aaT	p.K601N	
C1orf192	257177	37	1	1.61E+08	1.61E+08	+	Missense_Mut	SNP	A	G			g.chr1:161334861A>G	A c.428T>C	c.(427-429)atA>aCa	p.I143T	
DUSP12	11266	37	1	1.62E+08	1.62E+08	+	Missense_Mut	SNP	T	C			g.chr1:161719830T>C	A c.239T>C	c.(238-240)gtggC ggC	p.V80A	
QSOX1	5768	37	1	1.8E+08	1.8E+08	+	Silent	SNP	C	T			g.chr1:180153105C>T	A c.807C>T	c.(805-807)cTc>cT	p.L269L	
EPRS	2058	37	1	2.2E+08	2.2E+08	+	Silent	SNP	G	A			g.chr1:220156537G>A	A c.3294C>T	c.(3292-3294)gcC>cgt	p.A1098A	
RGS7	6000	37	1	2.41E+08	2.41E+08	+	Missense_Mut	SNP	G	C	NA	NA	g.chr1:241031950G>C	A c.546C>G	c.(544-546)gaC>gaG	p.D182E	
ZBTB18	10472	37	1	2.44E+08	2.44E+08	+	Missense_Mut	SNP	A	G			g.chr1:244217482A>G	A c.406A>G	c.(406-408)Agc>Gc	p.S136G	
PCBD1	5092	37	10	72645567	72645567	+	Silent	SNP	T	C			g.chr10:72645567T>C	A c.123A>G	c.(121-123)aaA>aaG	p.K41K	
ECF1	11319	37	10	74899185	74899186	+	Frame_Shift_Ins	INS	-	T	NA	NA	g.chr10:74899185_74899186	A c.1302_1303insA	c.(1300-1305)aaaafs	p.E435fs	
ZNF503	48858	37	10	77160064	77160064	+	Silent	SNP	G	A	NA	NA	g.chr10:77160064G>A	A c.384C>T	c.(382-384)ccC>cT	p.P128P	
PDZD7	79955	37	10	1.03E+08	1.03E+08	+	Missense_Mut	SNP	C	T			g.chr10:102781635C>T	A c.787G>A	c.(787-789)Ggt>Agt	p.G263S	
SORCS1	114815	37	10	1.08E+08	1.08E+08	+	Missense_Mut	SNP	A	C	NA	NA	g.chr10:108434902A>C	A c.1845T>G	c.(1843-1845)gtg>gag	p.D615E	
GFRAL1	2674	37	10	1.18E+08	1.18E+08	+	Missense_Mut	SNP	C	G	rs200978034		g.chr10:118030598C>G	A c.70G>C	c.(70-72)Gga>G	p.G24R	
CPXM2	119587	37	10	1.26E+08	1.26E+08	+	Missense_Mut	SNP	T	C	rs561902091		g.chr10:125516817TC	A c.1829A>G	c.(1828-1830)tcC>Gc	p.Y610C	
AL355531.2	0	37	10	1.31E+08	1.31E+08	+	Silent	SNP	C	T	rs57778684	byFrequency	g.chr10:131309064C>T	A c.462G>A	c.(460-462)agG>agA	p.R154R	
TCEGR1	256536	37	10	1.33E+08	1.33E+08	+	Missense_Mut	SNP	G	A	NA	NA	g.chr10:132915133G>A	A c.1324C>T	c.(1324-1326)Ccg>tcg	p.P442S	
MIL55B	777807	37	11	12717274	12717274	+	Silent	SNP	G	C	rs77008084		g.chr11:12717274G>C	A c.1316A>C	c.(1316-1316)A>Ac	p.T388T	

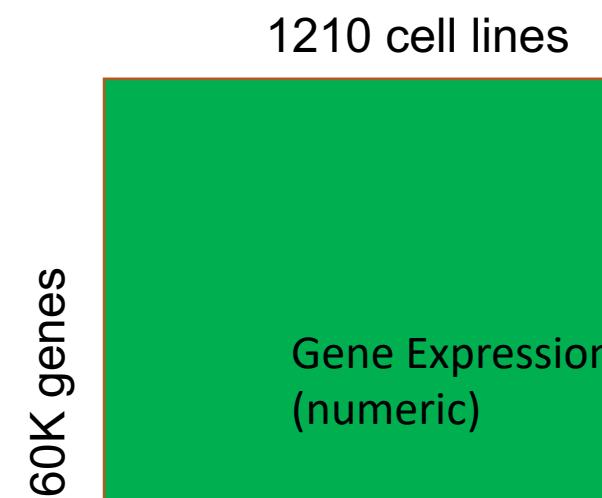
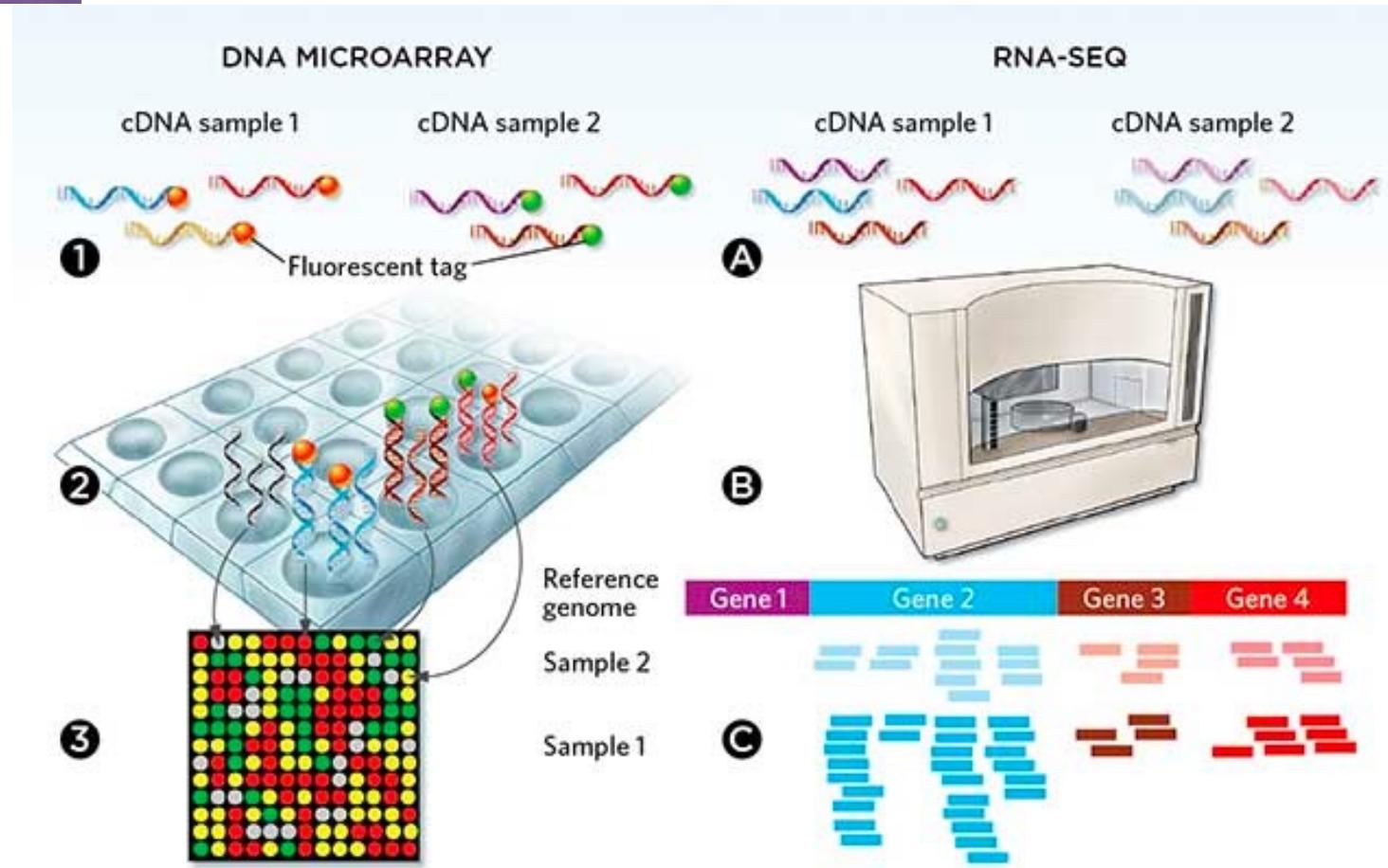
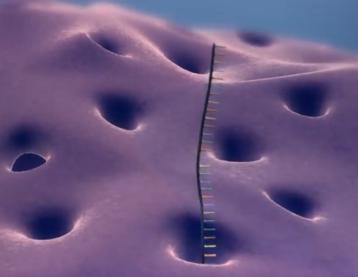
Mutation status (1 or 0)

Genome/copy number variation



1657 cell lines





[Modify Query](#)

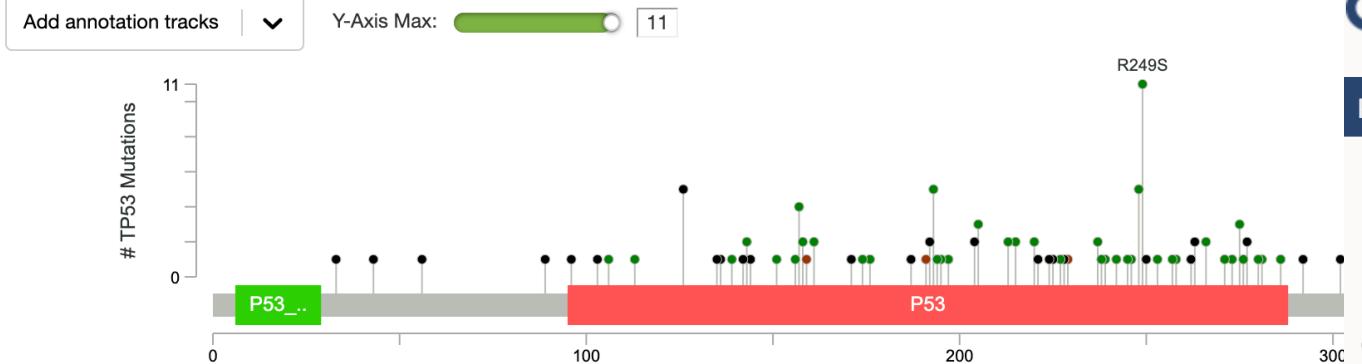

Liver Hepatocellular Carcinoma (TCGA, Firehose Legacy)

Samples with mutation and CNA data (366 patients/samples) - TP53

Queried gene is

[Oncoprint](#)
[Cancer Types Summary](#)
[Plots](#)
[Mutations](#)
[Co-expression](#)
[Comparison/Survival](#)
[CN Segments](#)
[Pathways](#)
[Download](#)

TP53


[Projects](#) [Data](#) [Tools](#) [News](#) [Help](#) [About](#) [Genome Version](#)

COSMIC v91, released 07-APR-20

COSMIC, the Catalogue Of Somatic Mutations In Cancer, is the world's largest and most comprehensive resource for exploring the impact of somatic mutations in human cancer.

ClinVar

Genomic variation as it relates to human health

 [Search ClinVar](#)
[Advanced search](#)
[About](#)
[Access](#)
[Submit](#)
[Stats](#)
[FTP](#)
[Help](#)

Was this helpful?

[Follow](#)

[Print](#)
[Download](#)

NM_004958.4(MTOR):c.7500T>G (p.Ile2500Met)

[Cite this record](#)

Interpretation:

Likely pathogenic

Review status:

criteria provided, single submitter

Submissions:

5 (Most recent: Jan 21, 2020)

Last evaluated:

Apr 11, 2019

Accession:

VCV000376455.2

Variation ID:

376455

Description:

single nucleotide variant

Species



Search

Metadata

Signature

Enrichment

Metadata Search

Human examples:

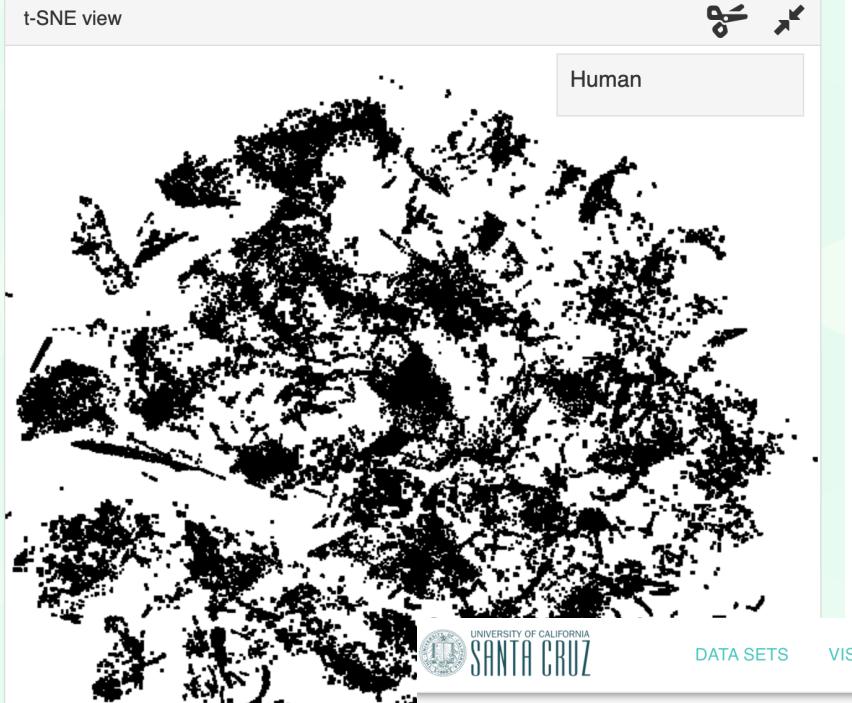
GSE81547, GSM2679484, Macrophage, Brain

Search for samples / GSE / GSM

Search

Tissue Types

- Cardiovascular System
- Connective Tissue
- Digestive System
- Immune System
- Integumentary System
- Muscular System
- Nervous System



Help with transcripts

Add Gene (e.g. KRAS)

KRAS

Study A GTEX Adipose Tissue

Expression Unit TPM

Gene Page

Copy

CSV

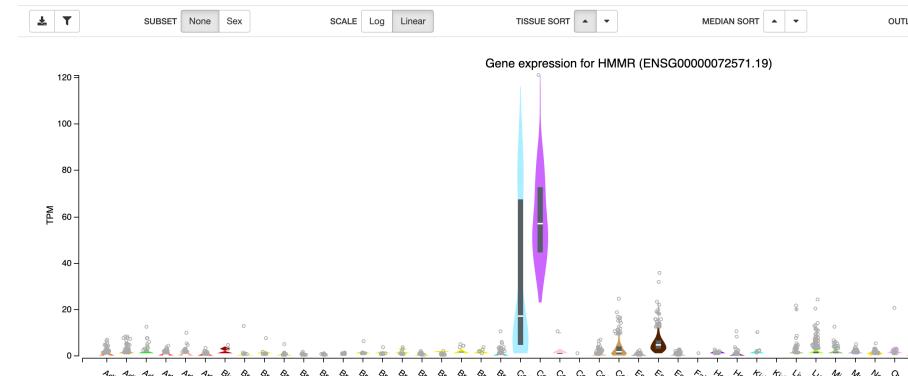
Gene Symbol	Gencode ID	Entrez Gene ID	Location	Gene Description
HMMR	ENSG0000072571.19	3161	chr5:163460203-163491945:+	hyaluronan mediated motility receptor [Source:HGNC Symbol;Acc:H
HMMR-AS1	ENSG00000251018.2	101927813	chr5:163483065-163494058:-	HMMR antisense RNA 1 [Source:HGNC Symbol;Acc:HGNC:49149]

Showing 1 to 2 of 2 entries

Gene expression for HMMR (ENSG0000072571.19)

Data Source: GTEx Analysis Release V6 (dbGaP Accession phs00424.v8.p2)

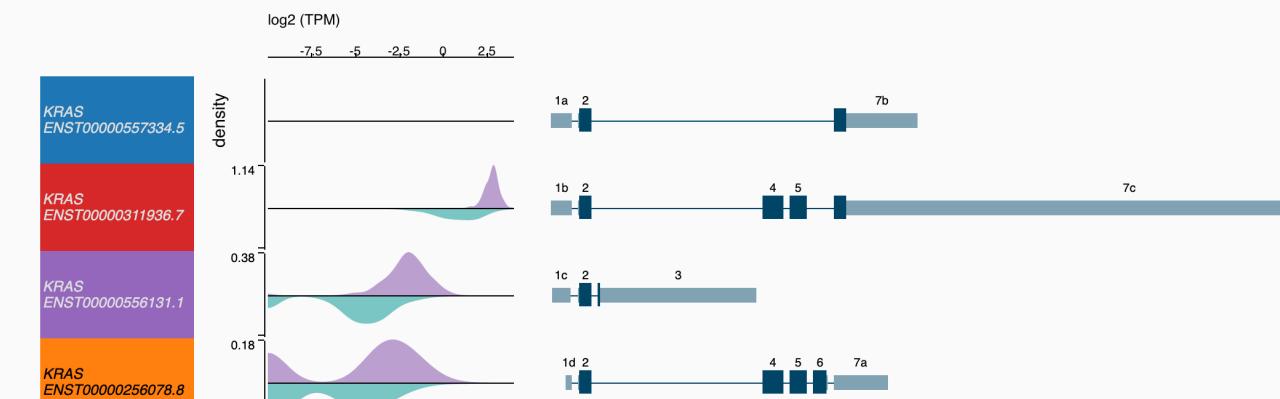
Data processing and normalization



Study A GTEX Adipose Tissue

Study B TCGA Uveal Melanoma

Expression Unit TPM





Keyword or GEO Accession

Search

Browse Content

Repository Browser

DataSets: 4348

Series: 134024

Platforms: 21247

Samples: 3740858

ArrayExpress – function

ArrayExpress Archive of Functional Genomics Data stores data from functional genomics experiments, and provides these data for reuse to the research community.



Browse ArrayExpress

SRA

SRA

Advanced

COVID-19 is an emerging, rapidly evolving situation.

Get the latest public health information from CDC: <https://www.coronavirus.gov>.

Get the latest research from NIH: <https://www.nih.gov/coronavirus>.

Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>

Full ▾

Send to: ▾

Design: HiSeq X Ten 150bp paired end sequencing of sample 90888_N5. Library was created with index CGCTCATT.

Submitted by: Cancer Research UK Cambridge Institute

Study: Liver Cancer Evolution - Lesion segregation

[PRJEB37808](#) • [ERP121138](#) • [All experiments](#) • [All runs](#)

[show Abstract](#)

European Genome-phenome Archive

All

Examples: EGAS000000000001, Cancer

Search

EGA home

About

Studies

Datasets

Data access committees

Data providers

Submit to EGA

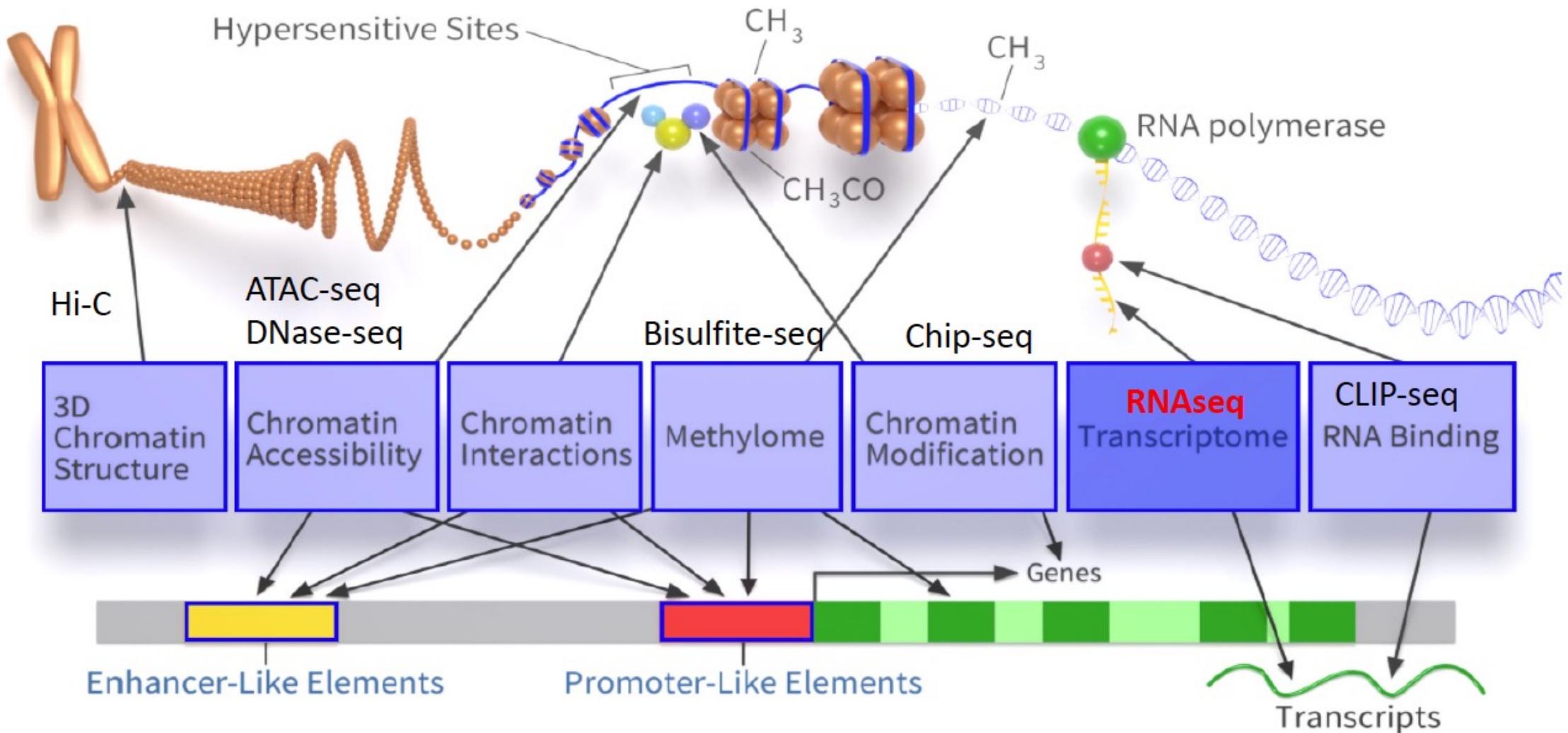
Contact Us

Login

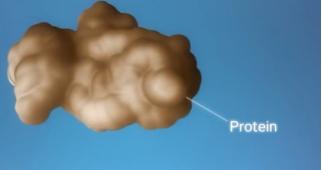
Help

- [Users FAQ](#)
- [Submitters FAQ](#)
- [Using your EGA account](#)
- [Contact Us](#)
- [EGA mailing list](#)

- 73320 experiments
- 2461459 assays
- 57.81 TB of archived data

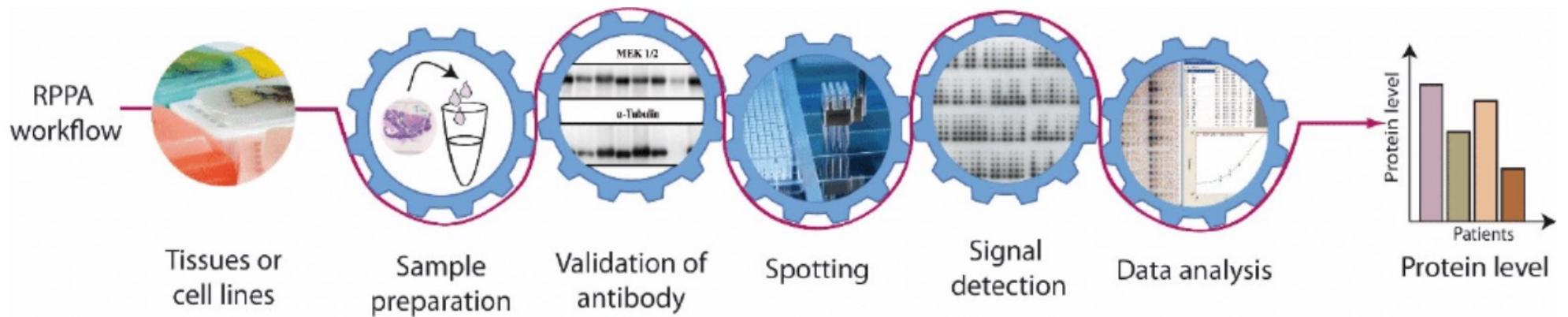


Based on an image by Darryl Leja (NHGRI), Ian Dunham (EBI), Michael Pazin (NHGRI)

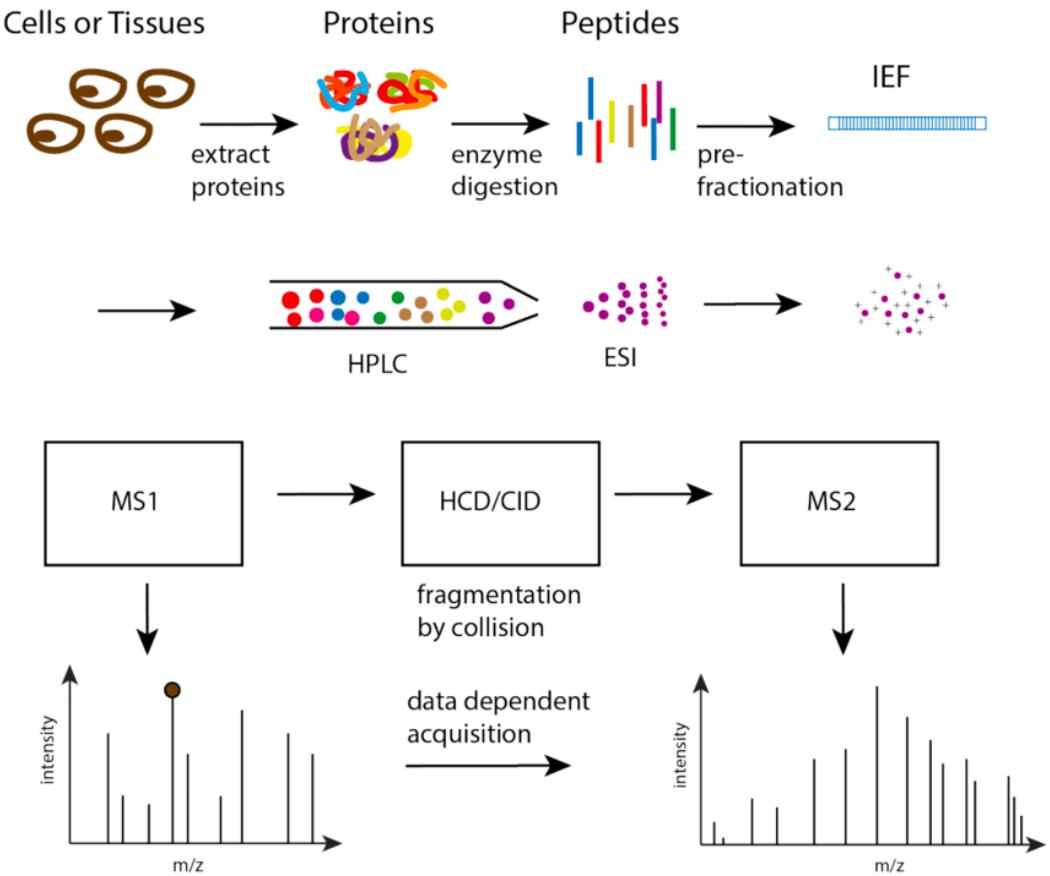


Proteome
Phosphoproteome

RPPA



MS



899 cell lines

214 proteins

Protein Expression
(numeric)

THE HUMAN PROTEIN ATLAS

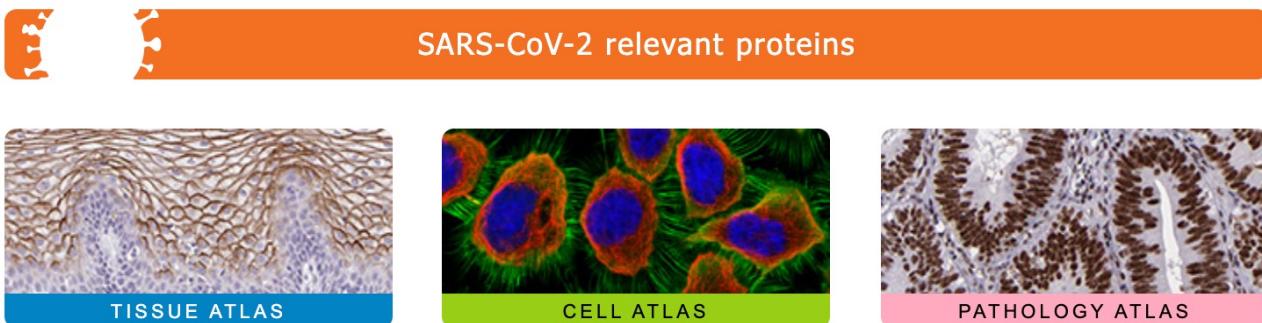


≡ MENU HELP NEWS

SEARCHⁱ

e.g. RBM3, insulin, CD36

Search Fields »



PRIDE Archive

PRoteomics IDEntifications Database

Home Resources Tools Docs About



NATIONAL CANCER INSTITUTE
Office of Cancer Clinical
Proteomics Research

Center for Strategic Scientific Initiatives



DATA PORTAL HOME



ASSAY PORTAL



Data Portal

CPTAC 3

(2016-present)

CPTAC 2 (2011-2016)

CPTC (2006-2011)

Latest Data Release and Publications:

July 2020

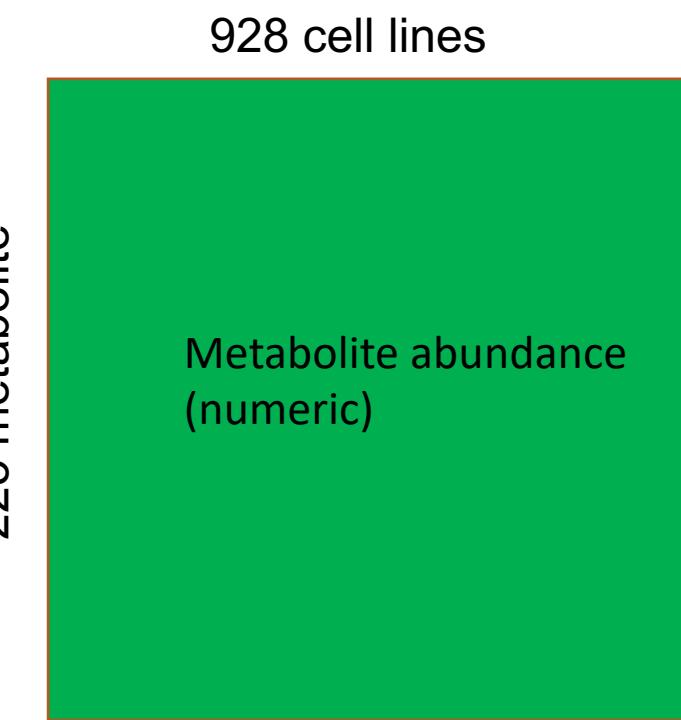
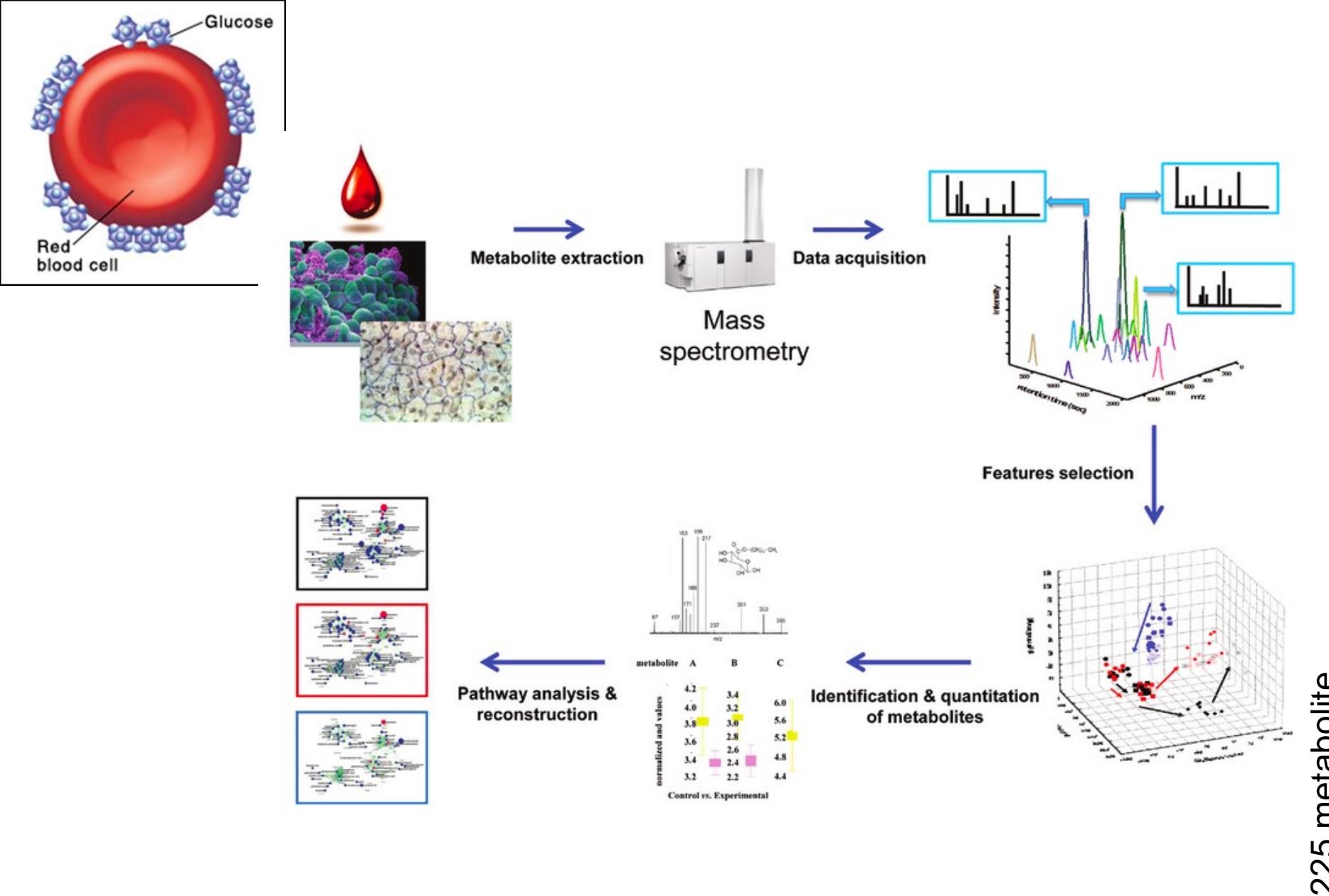
Proteogenomic Characterization Reveals Therap

Michael A. Gillette, Shankha Satpathy, Song Cao,

May 2020

CPTAC Head and Neck Squamous Cell Carcinom

Search



HMDB ID ↑

CAS Number

Name ↓↑

Structure

Formula

Average Mass ↓↑

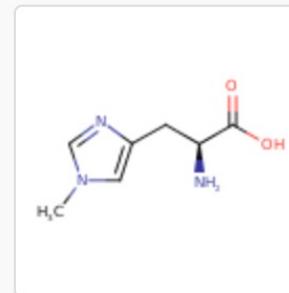
Monoisotopic Mass ↓↑

Biospecimen

HMDB0000001

332-80-9

1-Methylhistidine

C₇H₁₁N₃O₂

169.1811

169.085126611

Blood

Cerebrospinal

Feces

Saliva

Urine

HMDB0000002

109-76-2

1,3-Diaminopropane

MetaboLights

Examples: Alanine, Glucose

Home Browse Studies Browse Compounds Browse Species Download Help Give us feedback About

MetaboLights / Search

Filter your results

Type

- study
- compound

Technology

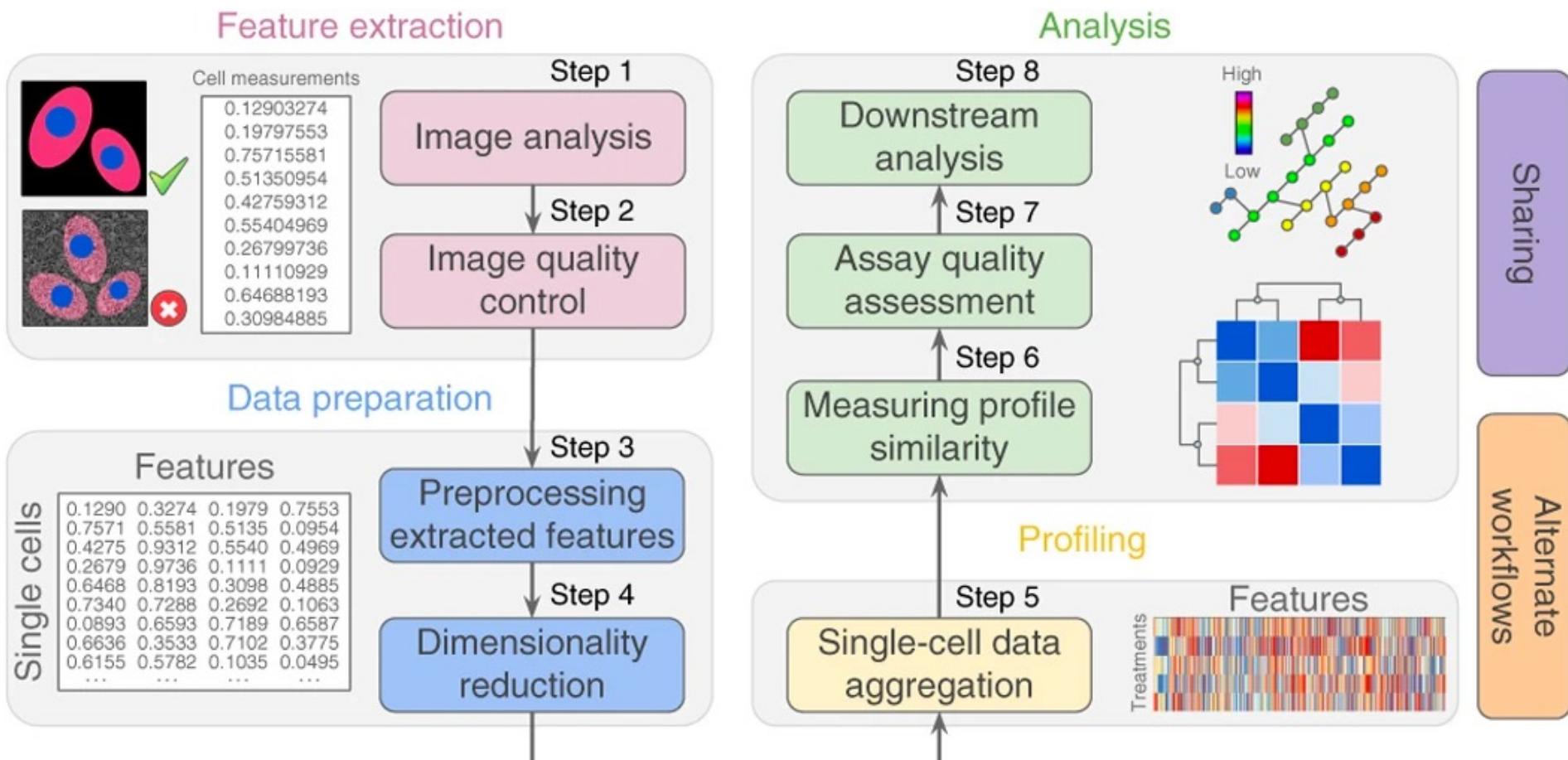
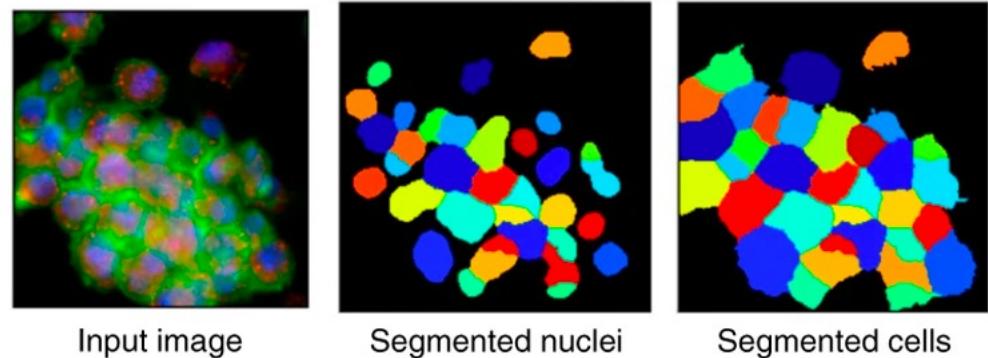
Organism

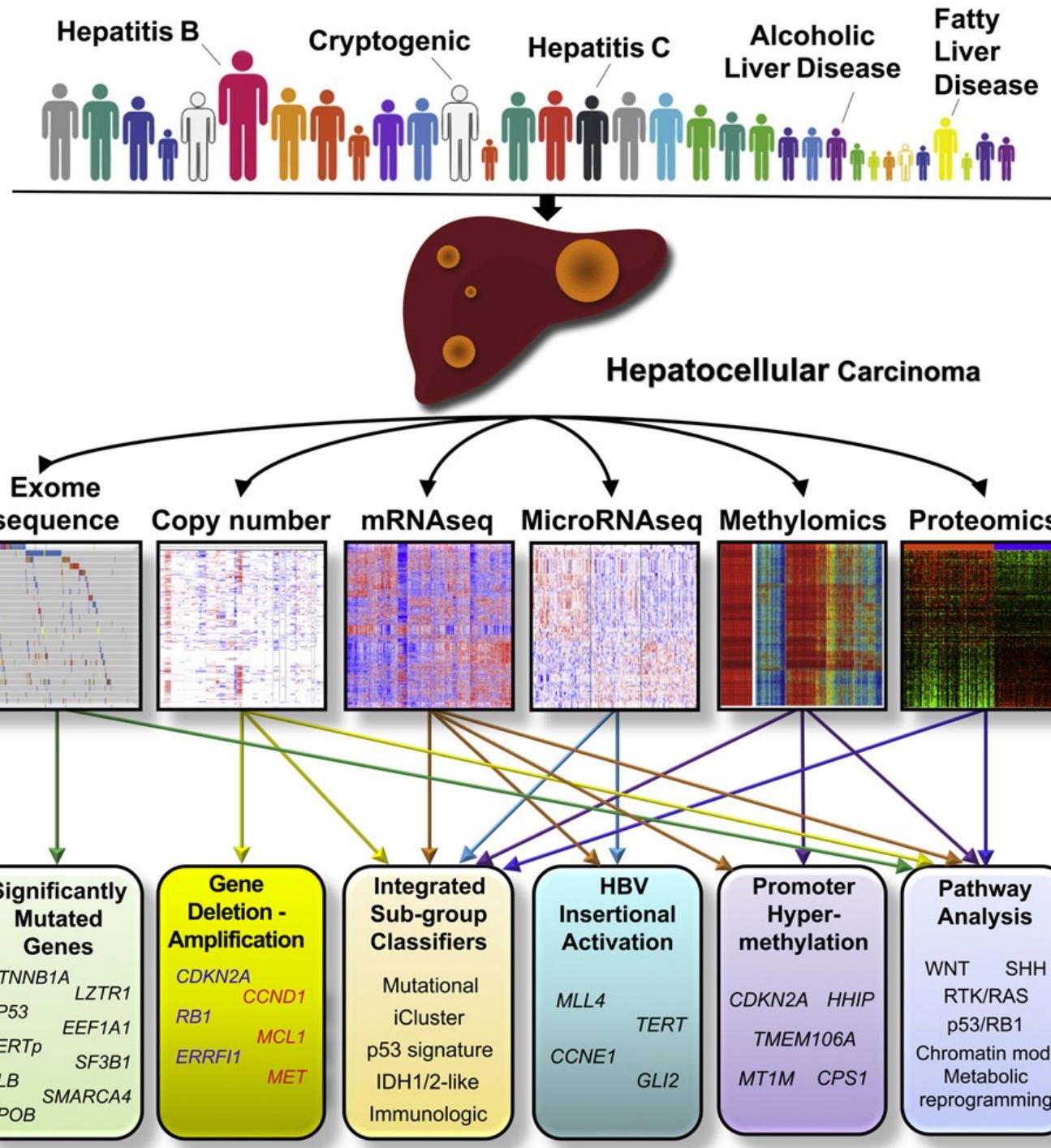
Organism Part

676 results , showing 1 to 10

Stable Isotope-Assisted Plant Metabolomics: Combi-Tracer-Based Labeling for Enhanced Untargeted Pro Annotation

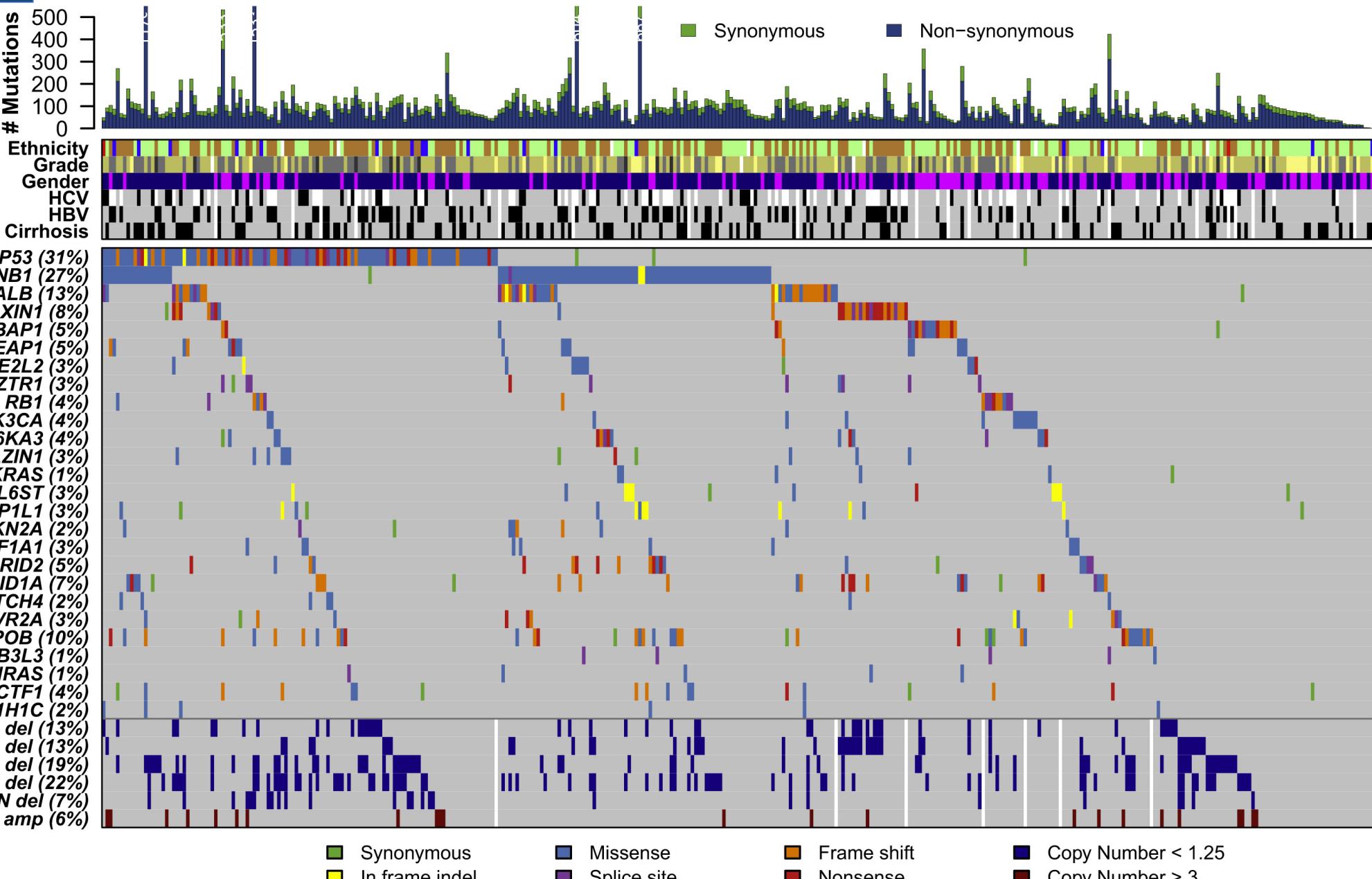
Study Identifier	MTBLS1217	Organism
Study Size	939.11MB	Study I
Submitted by	Christoph Bueschl	E-mail

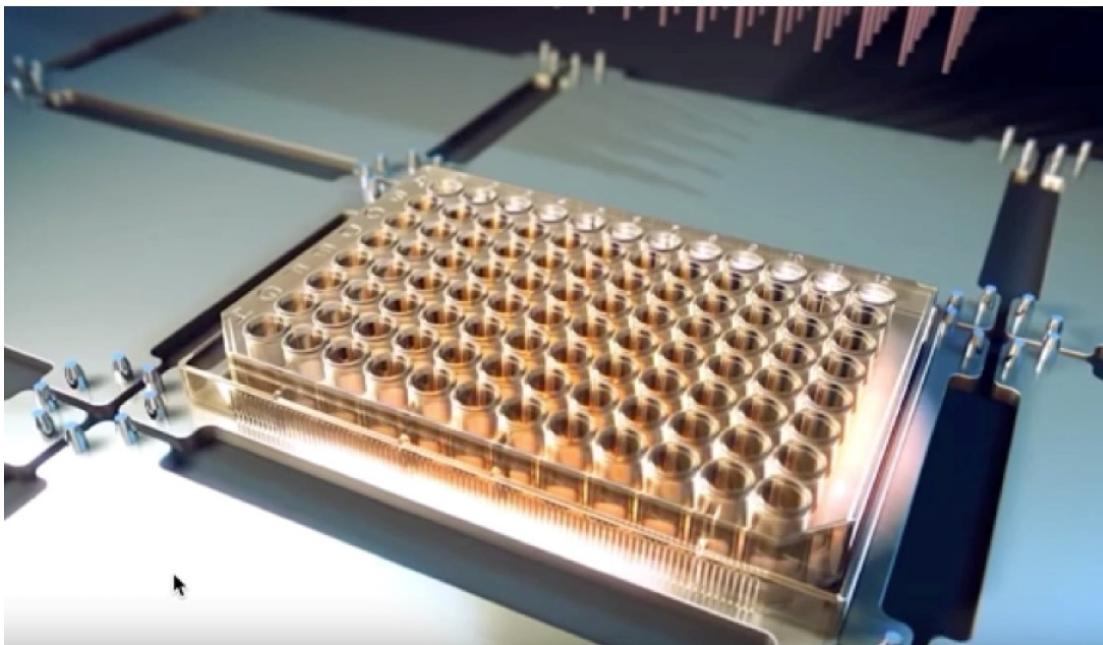
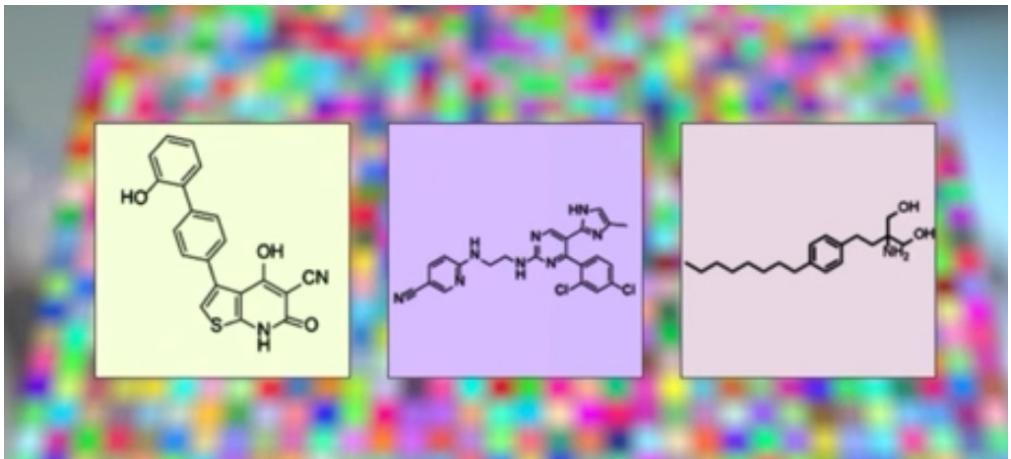




HCC Genetic Landscape

- African American
- American native
- Asian
- Caucasian
- Female
- Male
- G1
- G2
- G3
- G4
- NA

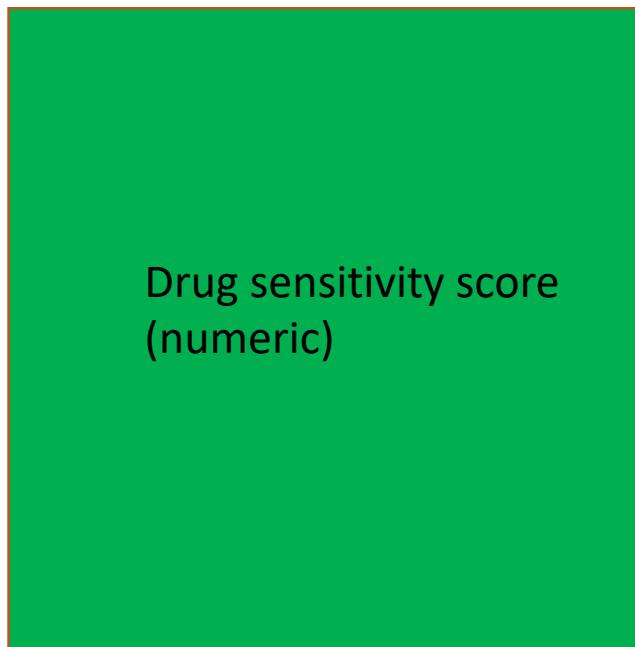


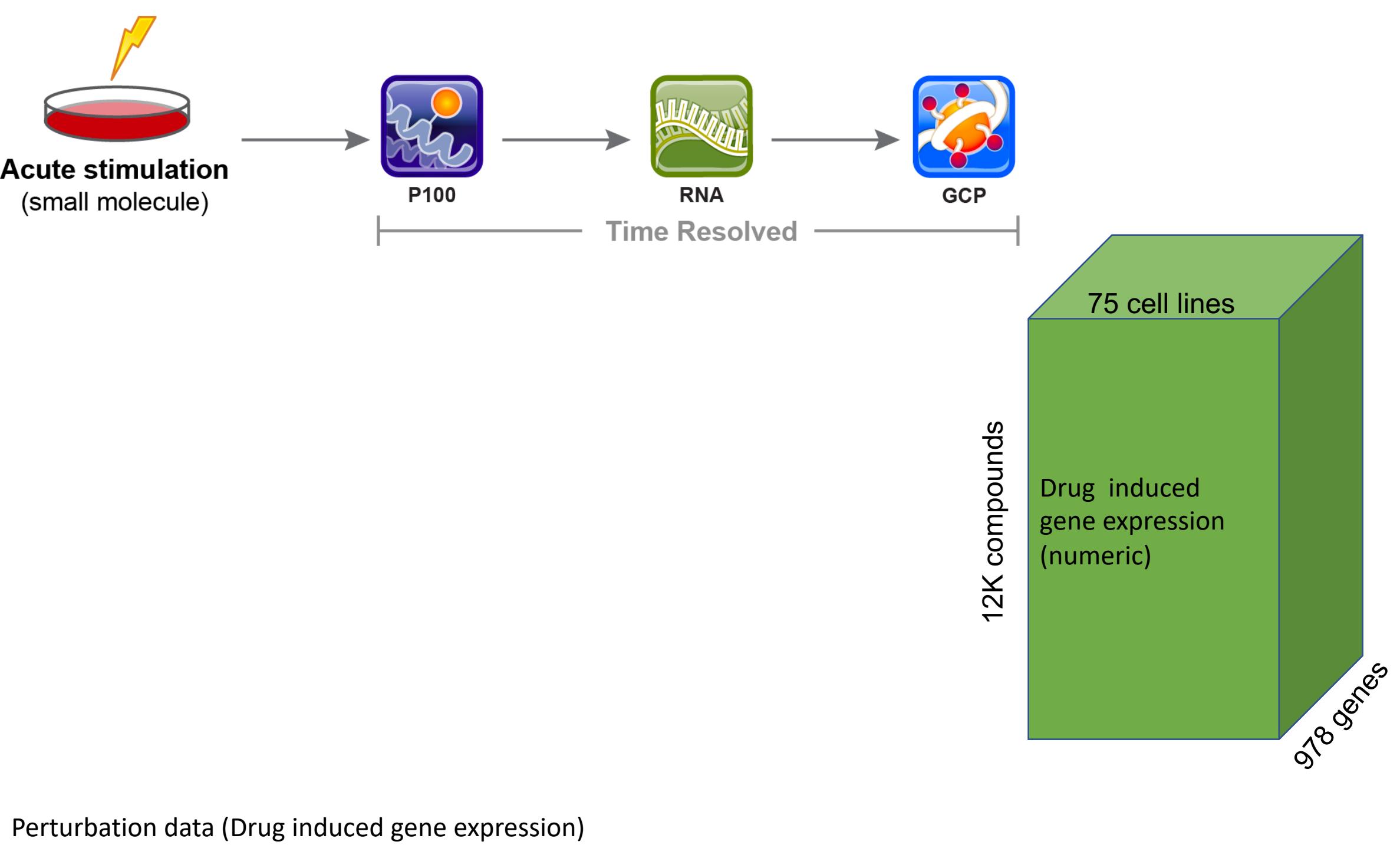


Perturbation data (Drug sensitivity)

5K compounds

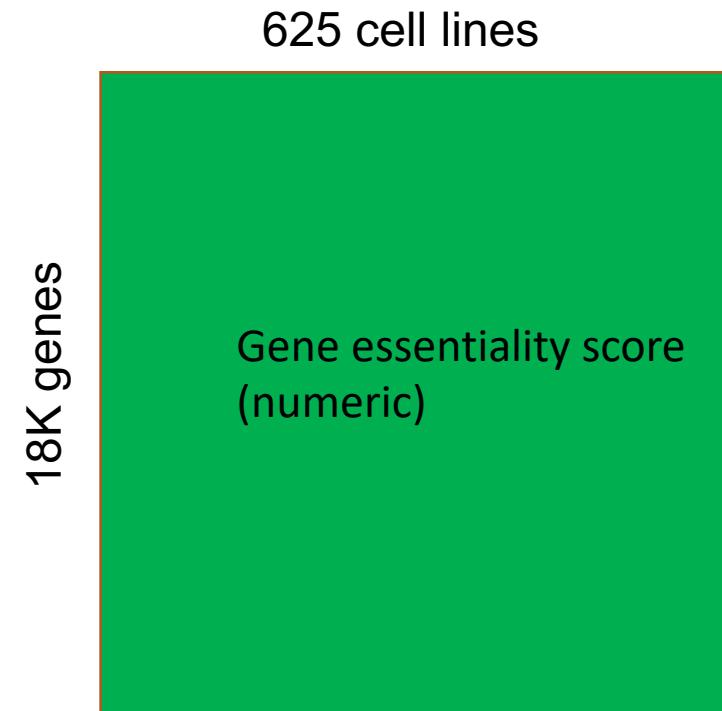
578 cell lines





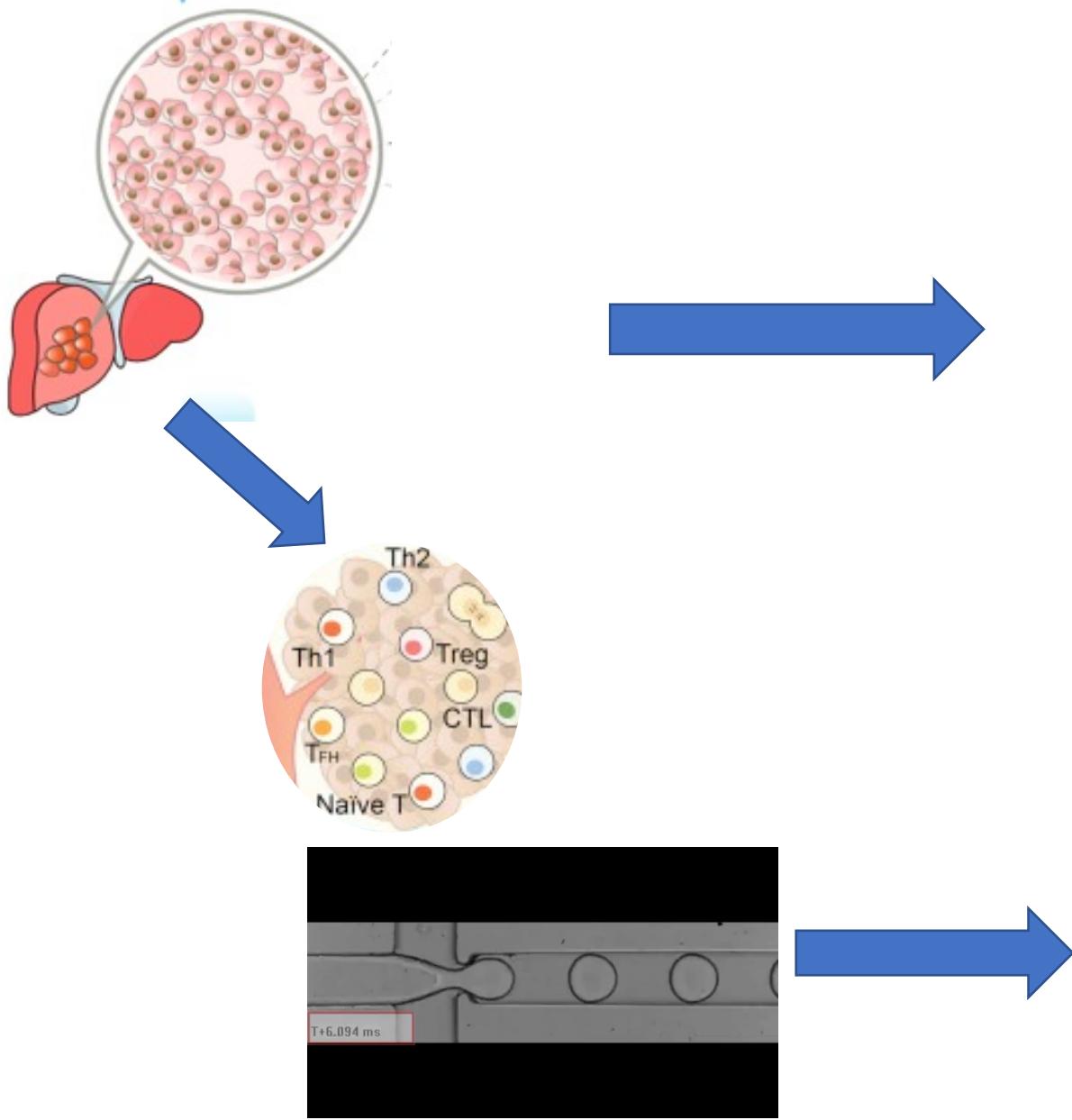


Perturbation data (Genome wide knock-out)



Demo

- <https://depmap.org/portal/>



20K patient tissues

60K genes

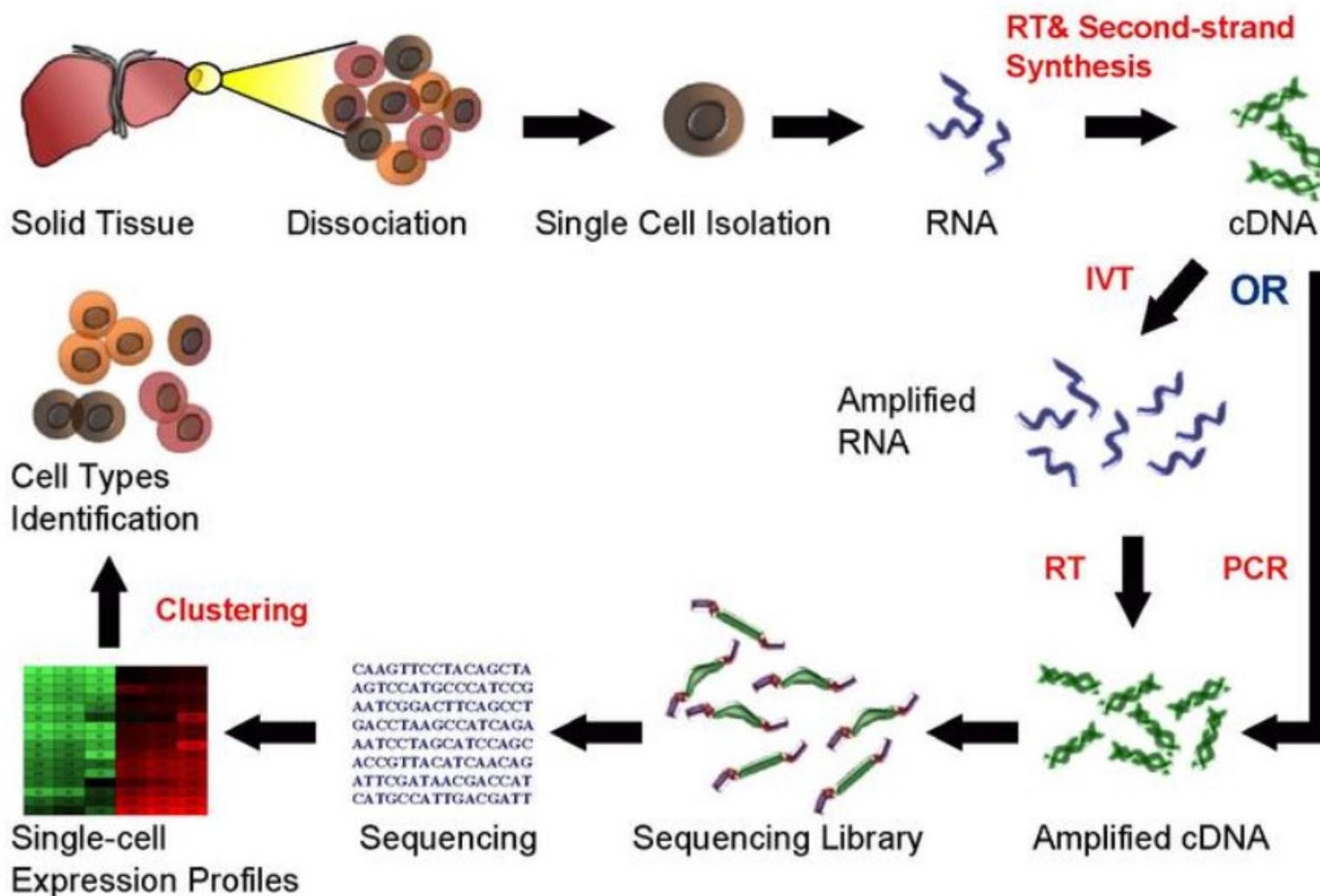
Gene Expression
(numeric)

Millions of single cells

60K
genes

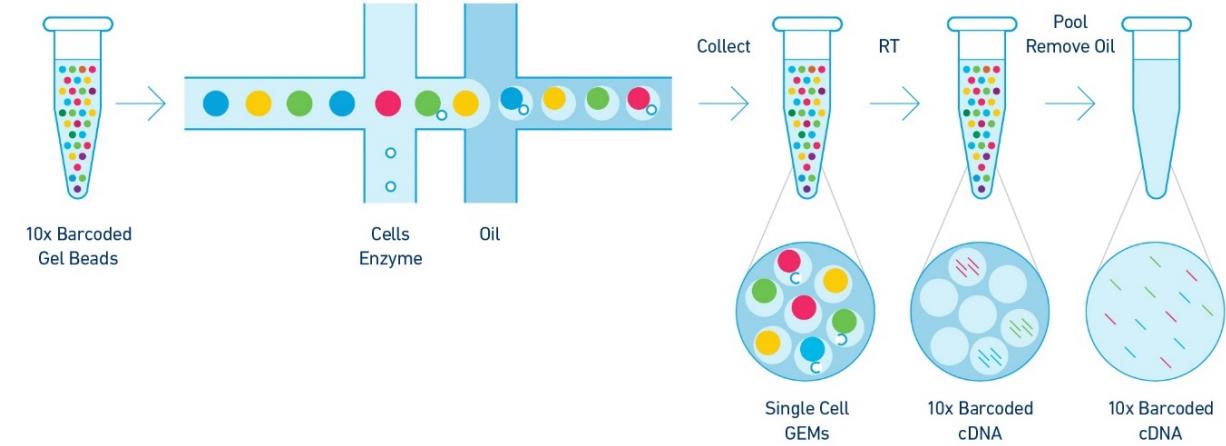
Gene Expression
(numeric)

Work flow of single cell RNA-seq



Overview about the scRNA-seq methods

	Micro-manipulation / Automated Pipetting	FACS	Microwell encapsulation	Droplet encapsulation
Cell Stress	Low	Moderate	Moderate	Moderate
Selection	Yes	Yes	No* / Yes ⁺⁺	No*
Doublet	Low	Low	Low-High	Moderate
Throughput	Low	Moderate	Moderate	High
Capture efficiency	Low	Moderate	Moderate	Low-Moderate
Academic / Commercial scRNA workflow	- CellenONE (Cellenion) [†] - Smart-Seq2 (42)	- MARS-Seq (39) - Smart-Seq2 (42)	- C1 (Fluidigm) - ddSeq (Biorad / Illumina) - ICell8 (Clontech) ⁺⁺ - Rhapsody (BD)	- InDrop (1CellBio) - DropSeq (Dolomite-bio) - 10X (Chromium)
Example of use	Fragile rare cells	Rare cells based on phenotype or marking	Large cell numbers	Large cell numbers



	FACS		Microwell encapsulation				Droplet encapsulation		
	Smart-Seq2	MARS-Seq	C1	ddSeq	ICell8	Rhapsody	InDrop	DropSeq	10X
Singlet Capture efficiency	82%	92%	39%	2.6%	37% ⁺⁺	Not reported	7%	Not reported	50%
Doublet rate	Not reported	2.27%	3-30%	5.8%	1.3-4%	0.6%	4%	0.36-11.3	1.6-3%
Reference	42	39	37 FWP	PB	PB	PB	36	37	26

[†]Automated pipetting system

^{*}Preselection or enrichment can be performed prior

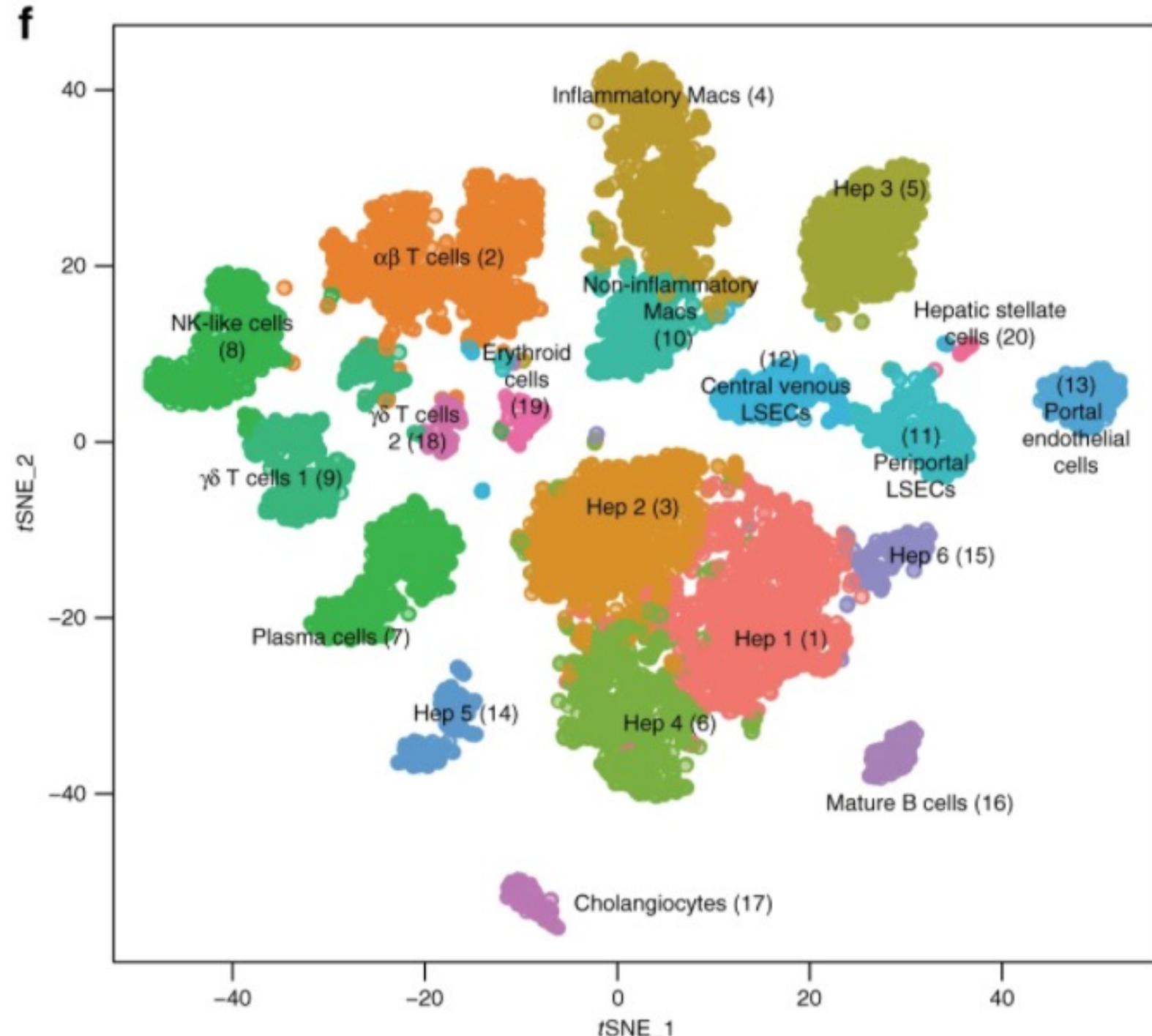
⁺⁺Only reagents added to wells containing singlets, determined by system

FWP: Fluidigm white paper

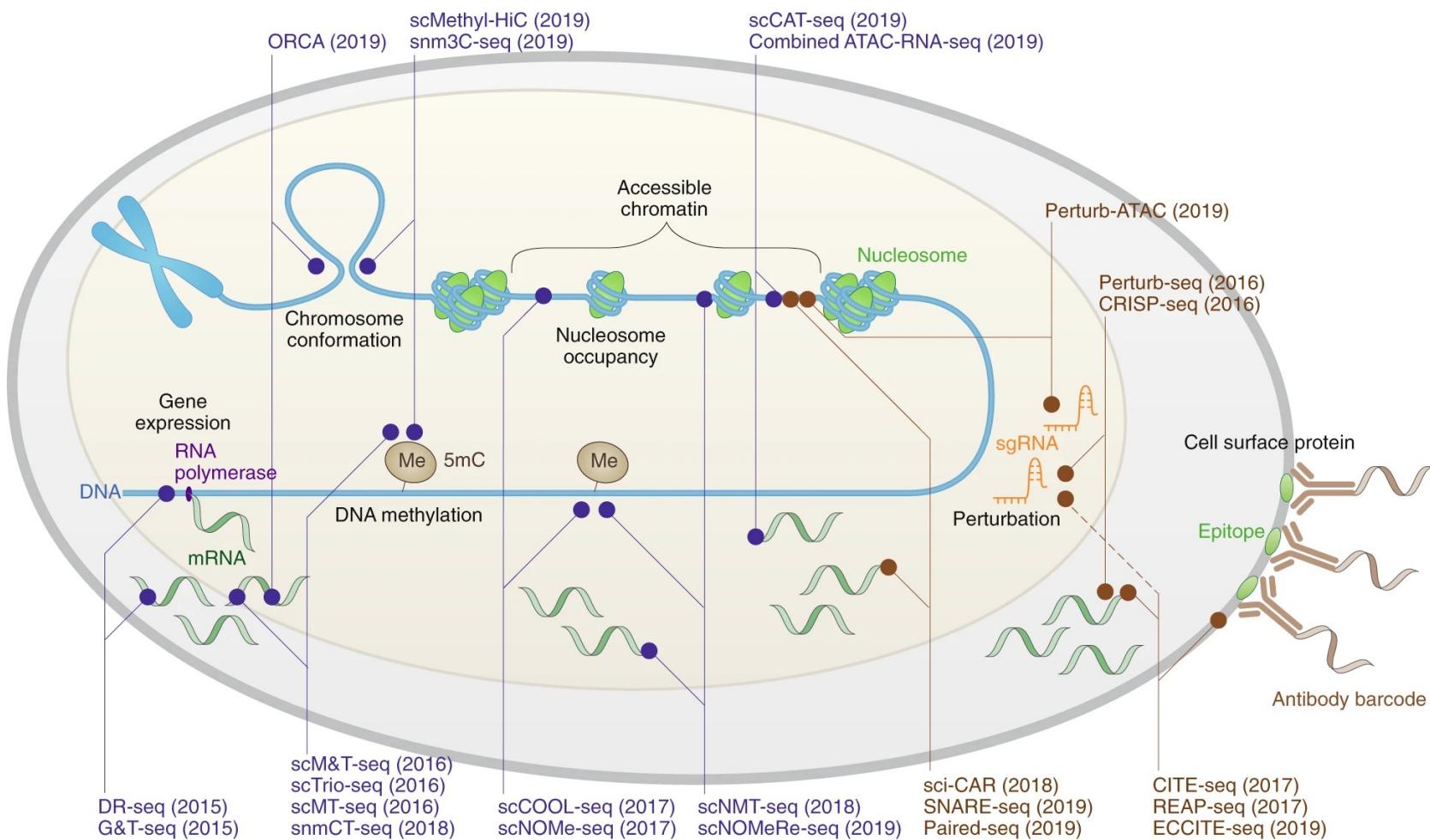
PB: Product brochure / manual

Main topics in single cell RNASeq

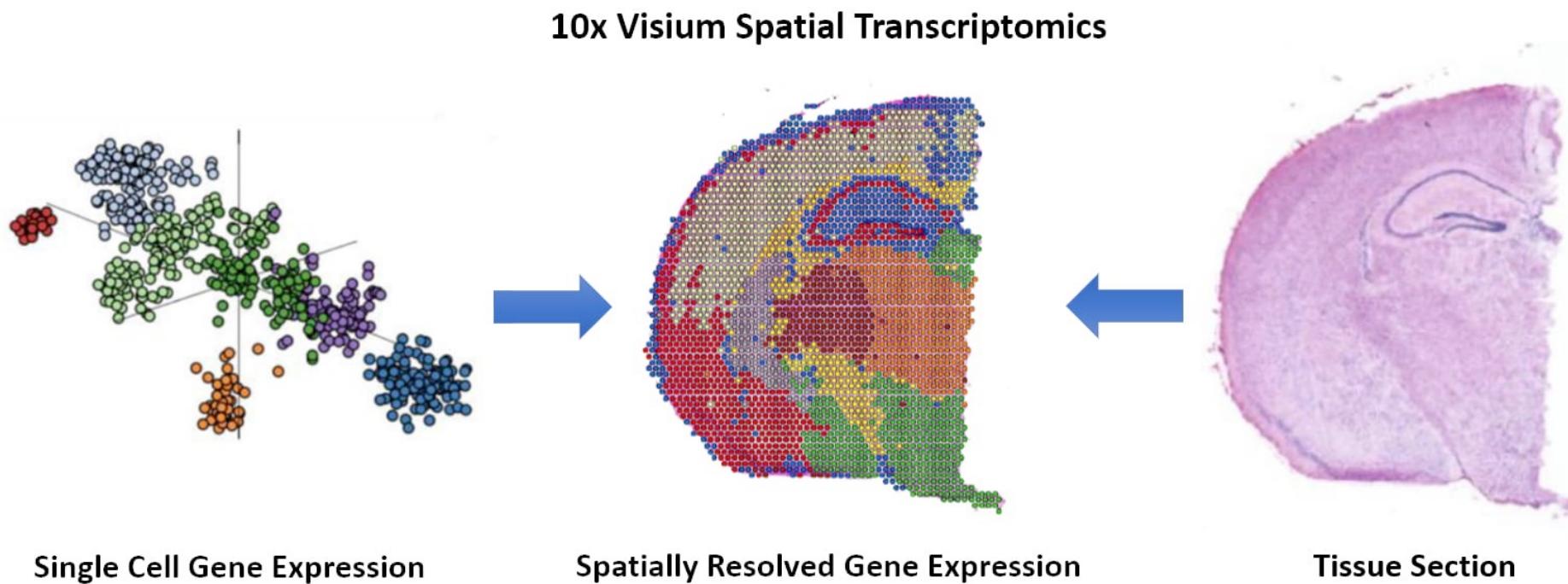
- Cell type identification
- Cell development
- Biomarker discovery



Multi-modal single cell



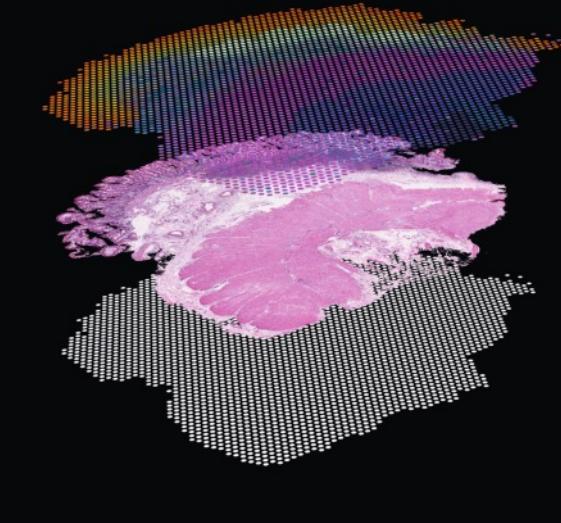
Spatial biology



Adapted from 10x Genomics



Method of the Year 2020:
Spatially resolved transcriptomics



Demo

- https://singlecell.broadinstitute.org/single_cell

Demo (Chemical Biology)

- <https://pubchem.ncbi.nlm.nih.gov/>
- <https://www.ebi.ac.uk/chembl/>
- <https://www.alphafold.ebi.ac.uk/>

AlphaFold Protein Structure Database

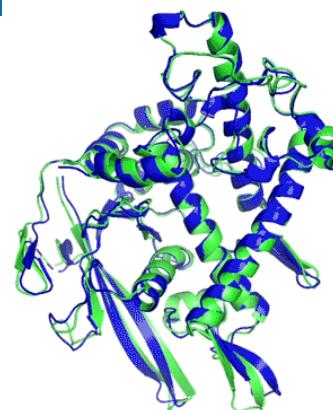
Developed by DeepMind and EMBL-EBI

Search for protein, gene, UniProt accession or organism

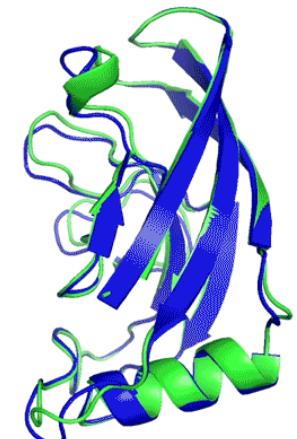
BETA

Search

Examples: Free fatty acid receptor 2 At1g58602 Q5VSL9 E. coli Help: AlphaFold DB search help



T1O37 / 6vr4
90.7 GDT
(RNA polymerase domain)



T1O49 / 6y4f
93.3 GDT
(adhesin tip)

● Experimental result

● Computational prediction

Data properties

- Matrix (small-samples, big features)
- Sparse
- Noisy (High-throughput)
- Batch effect
- Processed data



Bioinformatics

Big Data + AI ---> discovery

- Search
 - Search expression of MYC in lung cancer
- Compare
 - Is MYC expression higher in lung cancer than liver cancer?
- Predict
 - Predict drug sensitivity using pharmacogenomics data

Resources good for your publication

- GitHub for code sharing
- Figshare/synapse for data sharing
- Zenodo for code and data sharing

Resources

- Tutorial
 - http://manuals.bioinformatics.ucr.edu/home/R_BioCondManual/
 - <https://rstudio.com/resources/cheatsheets/>
 - <https://www.r-bloggers.com/>
 - <http://rafalab.github.io/pages/harvardx.html>
 - <https://liulab-dfci.github.io/bioinfo-combio>
- Troubleshooting
 - <https://stackoverflow.com/>
 - <https://www.biostars.org/>
- GitHub code repository
 - https://github.com/Bin-Chen-Lab/Awesome_BigData_AI_DrugDiscovery

ChatGPT would be the best resource for bioinformatics learning and coding.