

Giới thiệu phương pháp phân tích dữ liệu phiên mã gene theo cấu trúc không gian và theo từng tế bào đơn

Phuc-Loi Luu, PhD
Bioinformatics Team Lead, Pacific Informatics
Loi.lp@pacificinformatics.com.vn
HCMC, 12/08/2023

Contents

0. Self-introduction and Acknowledgement
1. Human Genome and Next-Generation Sequencing
2. Bulk RNA Sequencing (RNA-seq)
3. Single cell RNA Sequencing (scRNA-seq)
4. Bulk ATAC Sequencing (ATAC-seq)
5. Single-cell ATAC Sequencing (scATAC-seq)
6. Discussion: optimization and debugging for **scRNA/ATAC-seq**

Self-introduction and Acknowledgement

Acknowledgement



All patients who donate samples

All member of Prof. Susan Clark lab

- Prof. Susan Clark
- Dr Clare
- Dr Ruth
- Jenny
- Wenija
- Dilys
- Dr Qian
- Dr Amanda
- Dr Jo
- Dr Braydon
- Dr Shalima



Acknowledgement Collaborators and Students in Viet Nam

- NGUYEN CANH HIEP, MD, PhD
- LUONG THI MY HANH, MD, PhD
- TRẦN THỊ THANH KHƯƠNG, PhD
- CAO THỊ TÀI NGUYÊN, PhD
- TRỊNH VĂN NGŨ, PhD
- LE MINH THONG, PhD
- Bsc. LÊ NHẤT THÔNG
- Bsc. PHẠM MAI TÂM
- Bsc. NGUYỄN TẤN THANH GIANG
- Bsc. NGUYỄN ANH XUÂN
- Bsc. LE VAN HIEU
- Bsc. NGUYỄN MINH HOÀNG
- Bsc. TRẦN BÁ THIỀN
- ONG PHUC THINH, MD
- NGUYEN HUY THINH, MD
- DAO NGOC BAC, MD, Msc
- NGUYEN PHAN XUAN TRUONG, MD
- TRAN NGUYEN TRONG PHU, MD

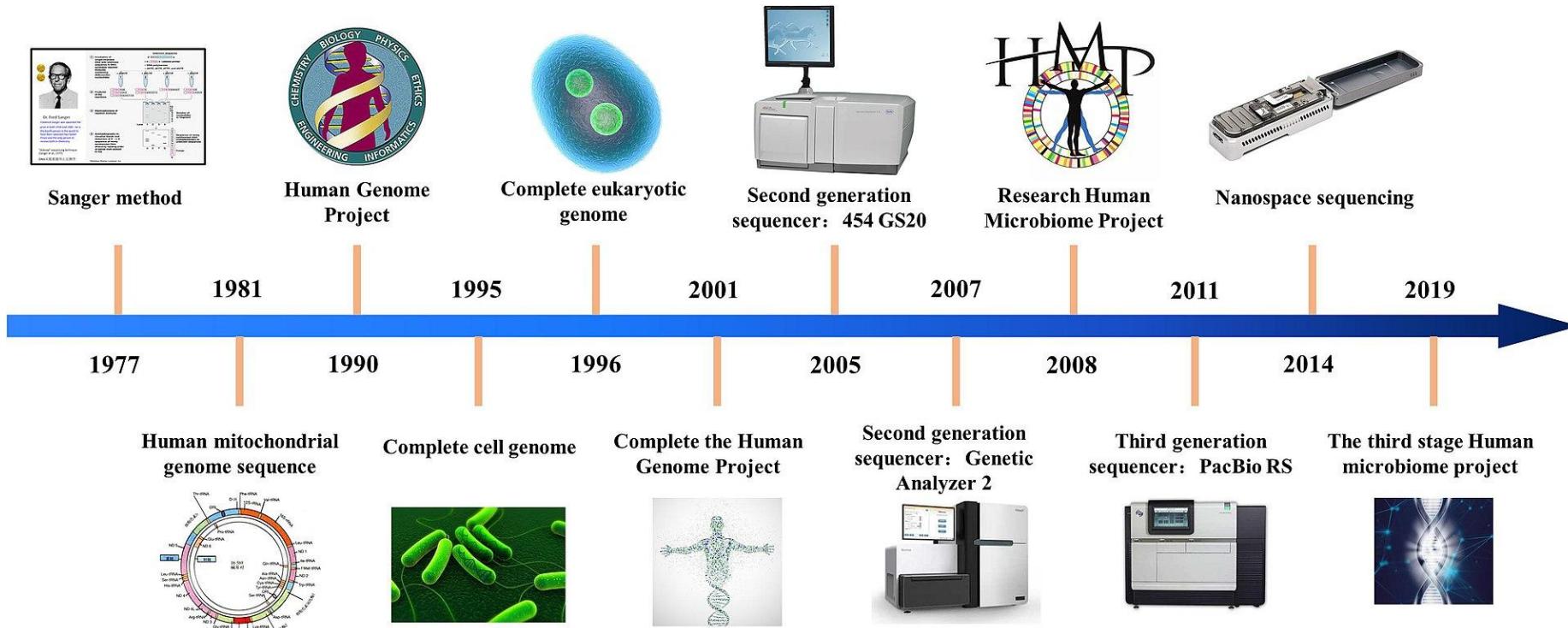
**All member of
VnPathoinformatics group**

The Bioinformatics Team

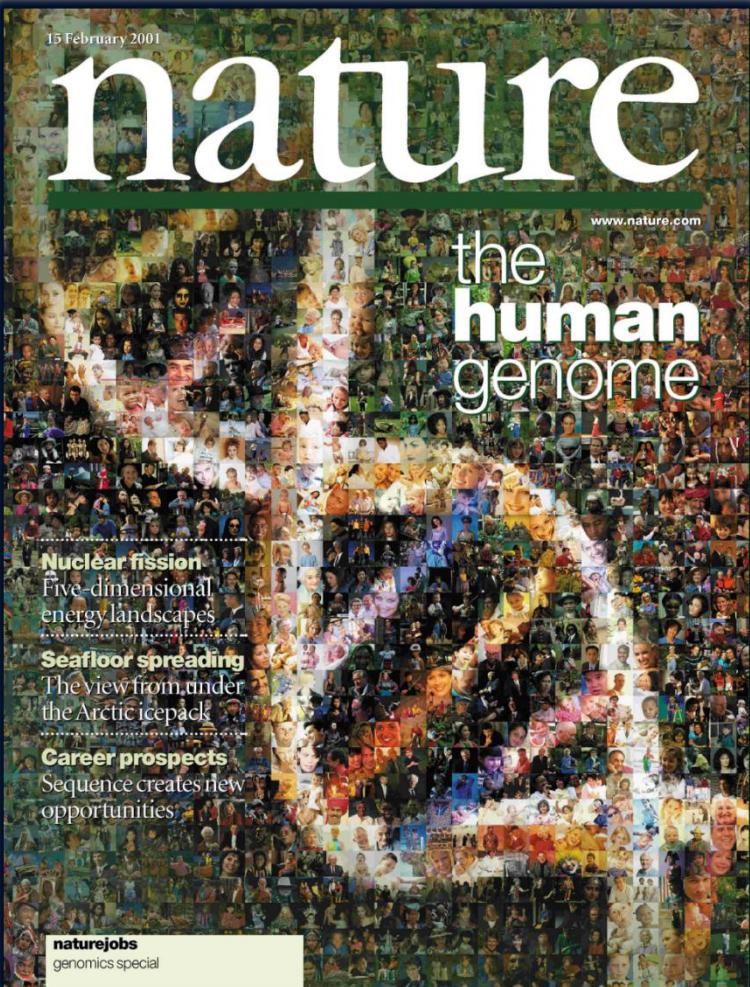
Name and Role	Profile Picture	Profile Picture	Name and Role
Luu Phuc-Loi , PhD Team Lead			Nguyen Viet Tuan , PhD Consultant
Du Hoang Tien , Msc. Associate			Phan Vo Thu Nga , Bsc. Associate
Nguyen Le Duc Minh , Dr. Genetics Consultant			Tran Thien Tan , Bsc. Technician
Dao Khuong Duy , Bsc. Technician Nguyen Minh Hoang , Bsc. Technician	 		Le Van Hieu , Bsc. Technician

Human Genome and Next-Generation Sequencing

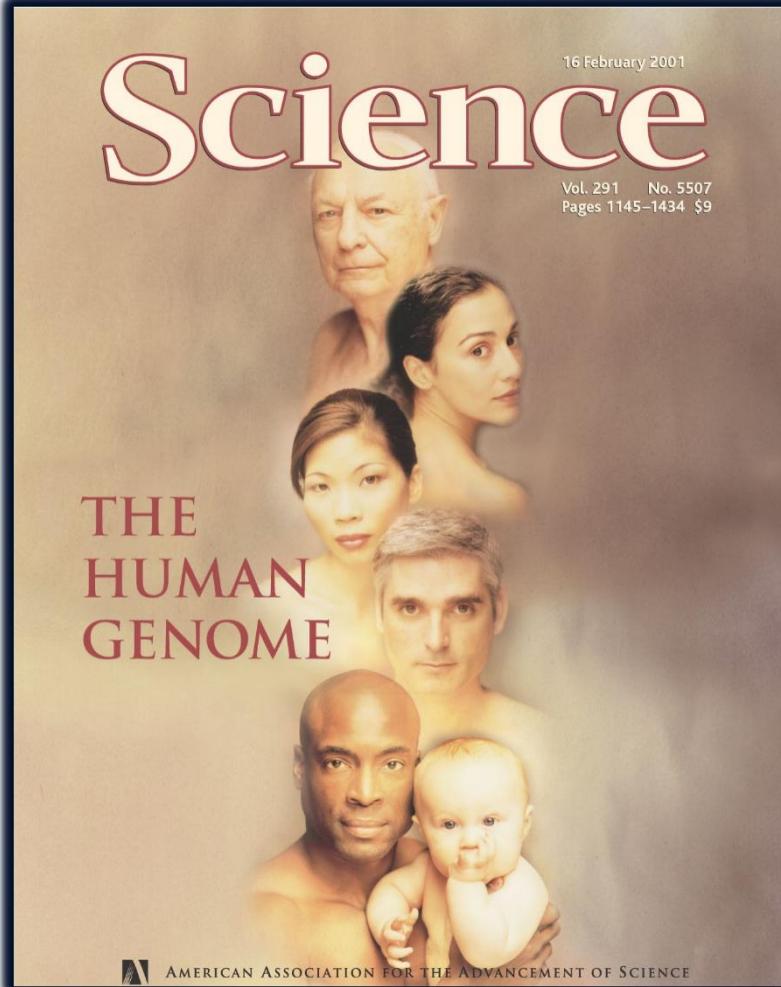
High-throughput sequencing (HTS) methods



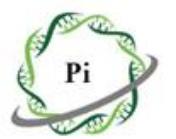
February 2001: Papers Reporting Draft Sequence of Human Genome



HGP Paper



Venter/Celera Paper



PACIFIC INFORMATICS
Your Trusted Partner in Bioinformatics

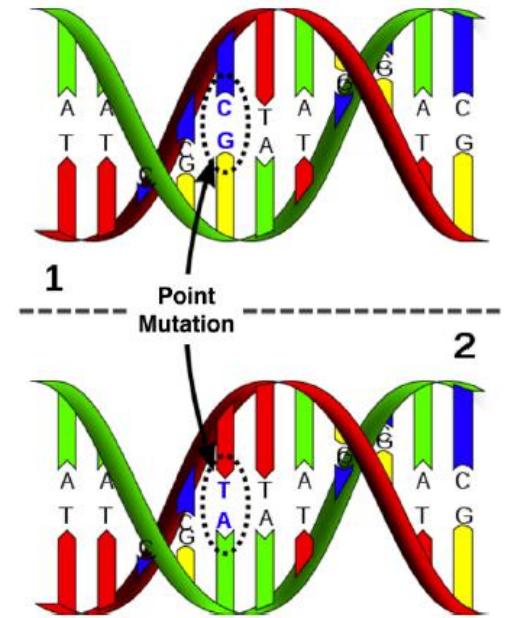
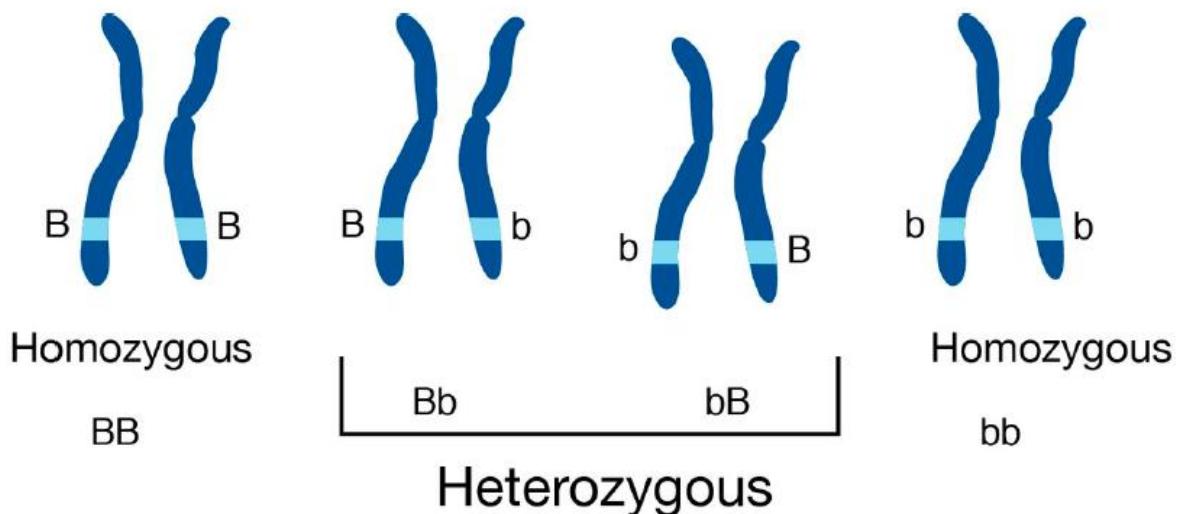


48541 agcccttcaa agaaatgttc tcagcaggca tggagccag gacttgctcc ctttggtag
48601 agagccgggt tgaaggtgac tgaagtgaaa tggacagta gaggcgaaaa gggtgtgag
48661 ttccctggagg tgggggtgt gggAACCTgc tttgtactga gatgcacccc tgccagttct
48721 gcctgaagat ttgaggcggg gggcaggggg gcggagtgaa gtcattttac tggtaagtaa
48781 ttttaaacct tttaatatta aagcaaacgt ggatatgtaa tgaatgaaat tcattctgga
48841 atgaaaaatt cacgtatgt taaaaataaa cacgggctt cagagaggac tttctggctg
48901 gcagcagact ccagattccc agggccctg caccctcctc tgcccacagg gcacctaatt
48961 ggagaaggtg tggaggaga gccaggccgg agtcagagca cactggtgac tccacatttg
49021 cagcgtgccccc tgcctctc ctgaggctt gcaacgtgca atatgctaag caaactcccc
49081 ctgtccccgt ccagttctg aggacaagag ccaccacctg tagcaaataa agaccaggca
49141 accctttgac tcatcttgc gagtctctgg aatcagaggg tagccacatc gctgagaggt
49201 ggagtgaagc actcgggtga aaaggtacaa ggaagtcaagg gacaggagtg tggggacatc
49261 acctagacaa tgacagagaa gaggggcaca gccgagttag gggagagggg ccggcagtcc
49321 tacatccccct ggcctgaagc acgctccagg gcagaaggaa aaacactgtc tttgggttcc
49381 aagagacctg agttcaaatt ctggctccac cactgaccac ctgtgttaacc ttgaactgct
49441 gtcgcctgaa cctcagggtt cccttctaaa aatagaggaa aaaaggatgc atttctcctt
49501 gcccctgtga gaacgaaatg gtgcaagcac caaggagcct cagcaaagggt cggcctgcc
49561 cccgcctggc caaaccttc ctcttcaggaa gccacggca accgtagttt gacagaagag
49621 cagcacctt attaatgtc tccctcaggat tgccttgcg caagtccacct aacctctgt
49681 ggctgcttcc tcattggaa aatatggctt ccagtaaaac ctgcctgttc cacctcttgg
49741 ggcacttggc aaacagcaaa agagtccaaa tgtgcaggct gggccaggcg cagtggctca
49801 tgcctgtaat cccagcaatt taggaagcca aggtggcggt atcacctgag gtcaggagtt
49861 tgagaccagc ctggccaaca tggtaaaacc ttgtctctac aaaaatacaa aaattagccg
49921 gcatgtatgg cgggtgcctg taatcccagt tactcggag gctgaggcaaa gagaatcgct
49981 tgaacccgga agggaaagggt tgcagttagc caagattgtt ccactgcact ccagctgggg
50041 caacagagcg agactctgtc taaaaaaaaaaaaaaa aaaaaaaaaaa aaacaatgca gagctggctg
50101 tgtaaaaaaac ctgttccact gcagggccca gtgtccacca ggctgggggtt caggccatag
50161 ggggtgggggc ccagcatcag cctctcaggaa gccctgggg gggggcgca tccctgtcccc
50221 ctcgtggctt ggatgtgttcc tagcccaagt cctagttac acctgcccgtc gcctggcctc
50281 tcaggagagg cccagggta ggaggagcat ggttaaagggtt aagctgattt ggaagtccgc
50341 tggggaaa gcaactcctt gcacattggaa ggaacccgaga aagactgacc ccgaggacag
50401 cagccagcat ggccttcctt gggagcccat gttggggat tcctgctgca gccaaggctc
50461 agcccttggc gtcgcagggtt ctggctctgg cctctccccc tcccatgcag gggcaggg
50521 gagatggctt ctgaggacct tggcagctt tggccctggg aatagattt ccagggagct
50581 ttaaaggcgc tgagtgtgtc atccagctaa gcctggggaa ggagcttggc tcaggtcctg
50641 acaggtgtga cagggatggg gactggaaag taagagatga aaccctggct ggaggctgtg
50701 agcttccaca gccagcgctt gacaggagggat acccactgtt gcccacca

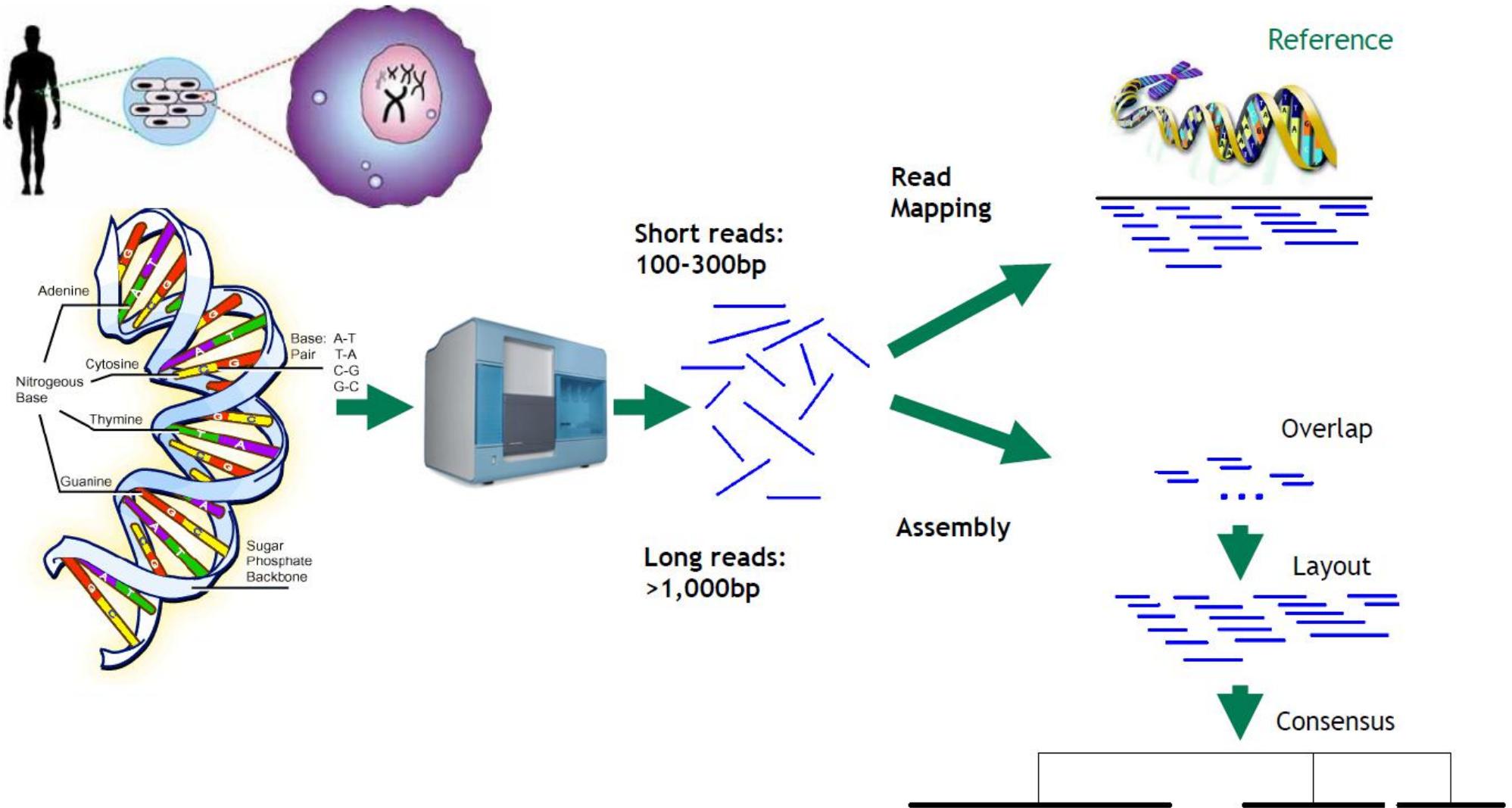
<https://www.ncbi.nlm.nih.gov/nuccore/806904736>

Human Genome Variation

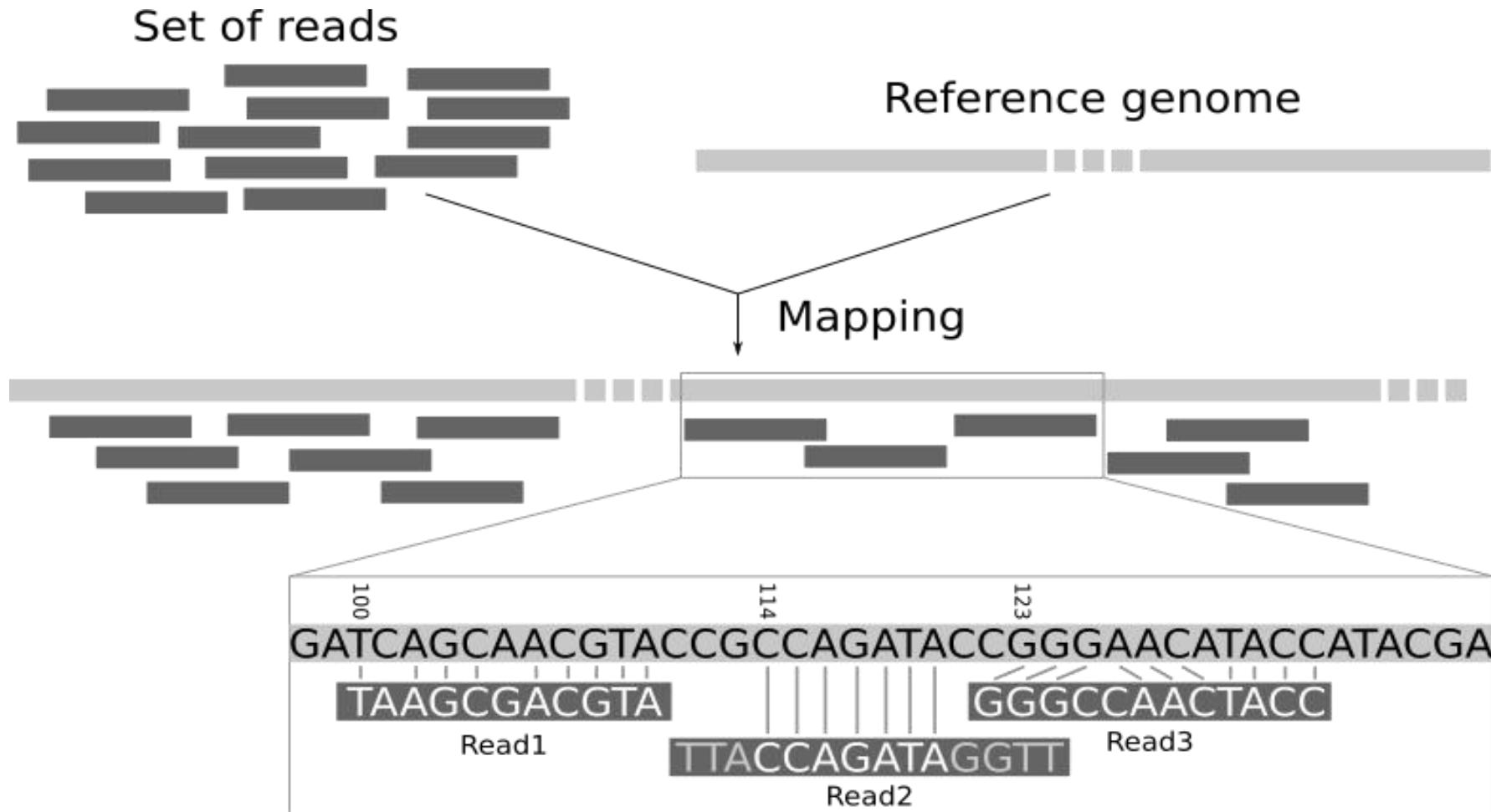
- Humans genomes are >99% similar by sequence
- A typical human genome has ~5 million variants with 3-4 million single-nucleotide variants
- Humans are diploid



Sequencing DNA



Read Alignment



Read Alignment

Reference sequence

	0	10	20
Reference	AGATTTCGATTGAGACTGT A - CTGATCAGGT		
read1	► AGATTTCGA		
read2	◀ TTTCGATT		
read3	◀ ATTGAGACTGT A - CT - ATC		
read4	► TGAG - CTGCATCTGATCA		
read7	◀ GAGACTGT A - CT		
read5	◀ AG - CTGC A - CTGAAACAG		
read8	► GACTGT A - CTGA		
read6	► G - CTGC A - CTGATCAGGT		

Variant Calling - Data Transformation

Alignment

Reference	0	10	20
	AGA	TTCGATTGAGACTGTA	-CTGATCAGGT
read1	>	AGATTCA	
read2	<	TTCGATT	
read3	<	ATTGAGACTGTA	-CT-ATC
read4	>	TGAG	-CTGCATCTGATCA
read7	<	GAGACTGTA	-CT
read5	<	AG	-CTGCA-CTGAACAG
read8	>	GACTGTA	-CTGA
read6	>	G	-CTGCA-CTGATCAGGT



List of Variants

CHR	POS	ID	REF	ALT	GT
chr1	12	.	GA	G	0/1
chr1	17	rs123	T	C	0/1

[No Title]

Genotype (GT):

0/0: hom. reference

0/1: heterozygous

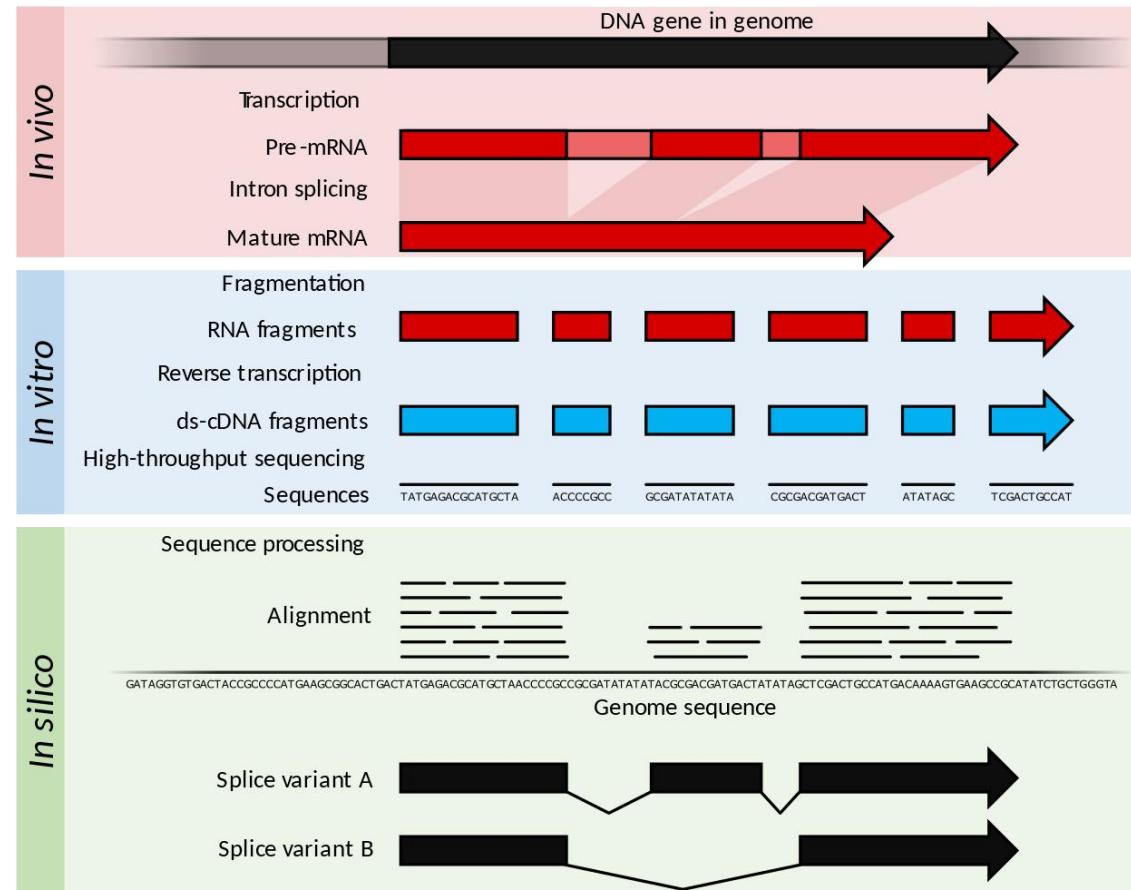
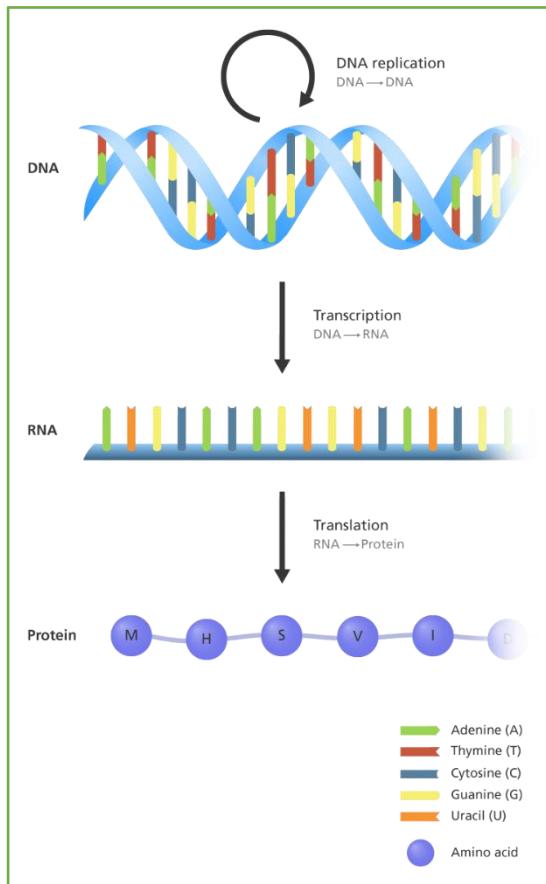
1/1: hom. alternative

SAM/BAM file

VCF/BCF file

Bulk RNA Sequencing (RNA-seq)

What is RNA sequencing (RNA-seq)?



What is RNA sequencing (RNA-seq)?

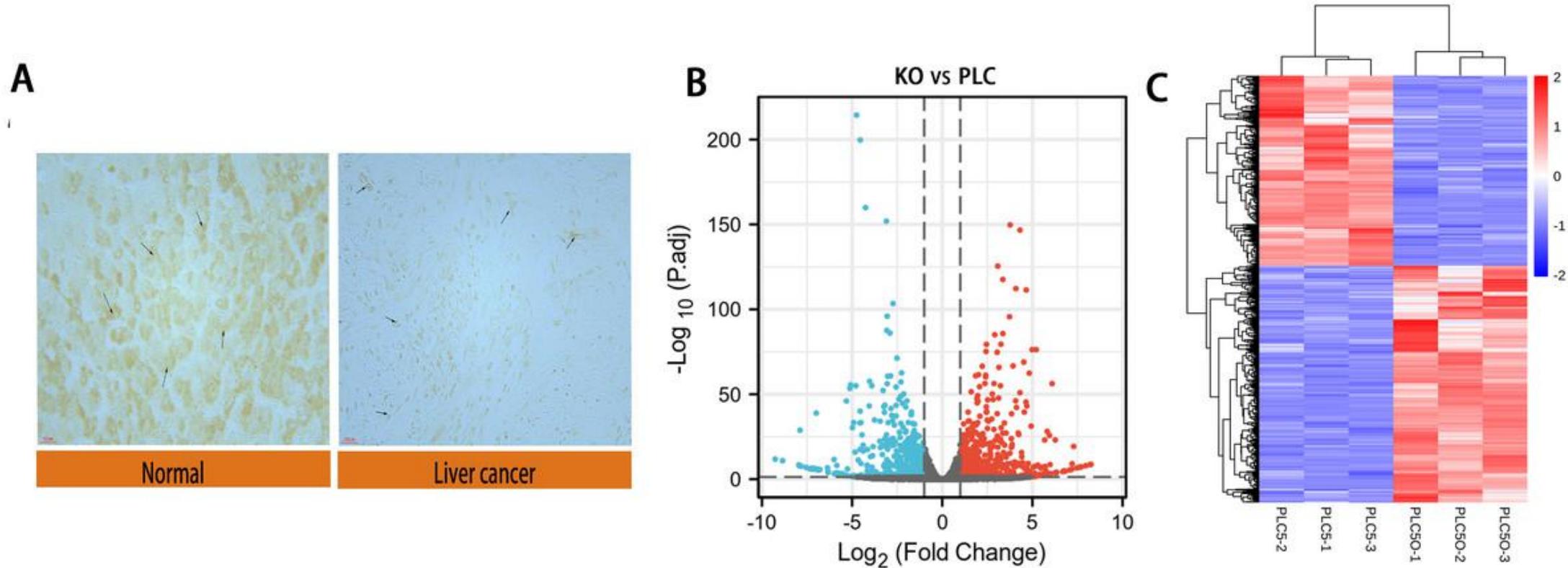


Figure: Differential expression of DAPK1 and volcano plot and heat maps of differential genes.

(A) Relative expression of DAPK1 compared to normal liver tissue samples. (B, C) Volcano plot and heat map showing the 732 differentially expressed genes. The color red indicates upregulated genes, and blue indicates downregulated genes. PLC, PLC/PRF/5 cells; KO, DAPK1-knockout PLC/PRF/5 cells. padj < 0.5, logFC |1|.

What is RNA sequencing (RNA-seq)?

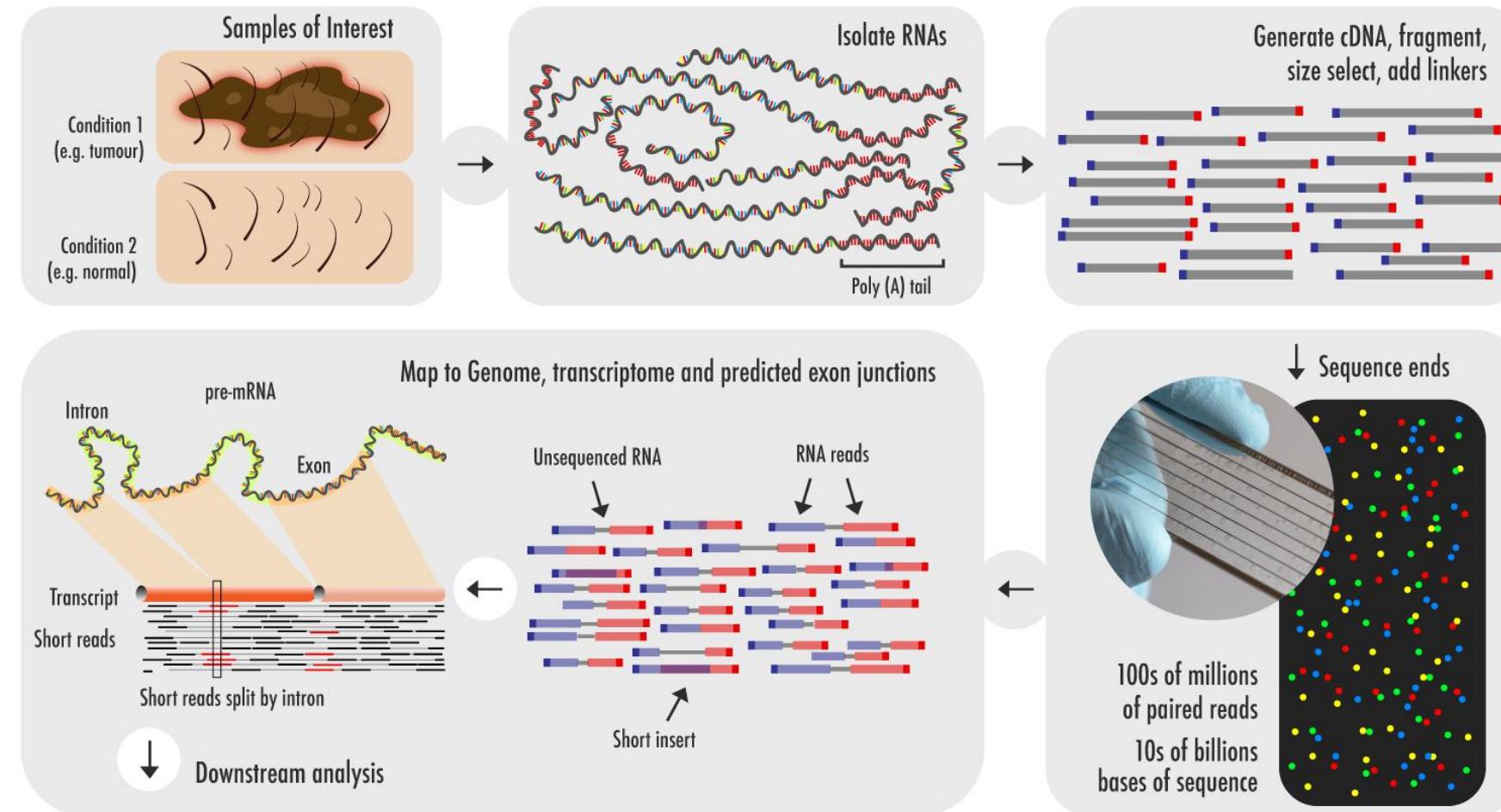
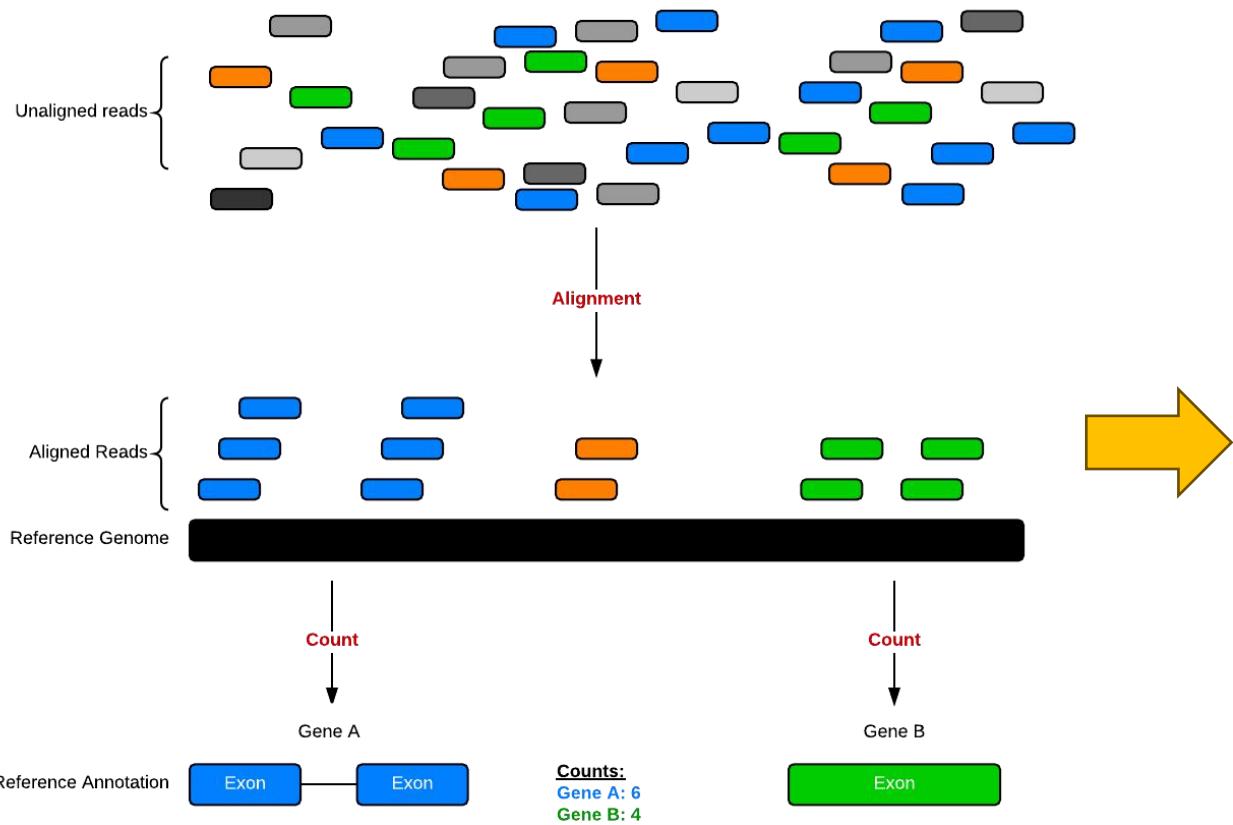


Figure. A workflow for RNA-seq. Credit: Technology Networks.

How can RNA sequencing quantify gene expression?



countData

	ctrl_1	ctrl_2	exp_1	exp_1
geneA	10	11	56	45
geneB	0	0	128	54
geneC	42	41	59	41
geneD	103	122	1	23
geneE	10	23	14	56
geneF	0	1	2	0
...
...
...

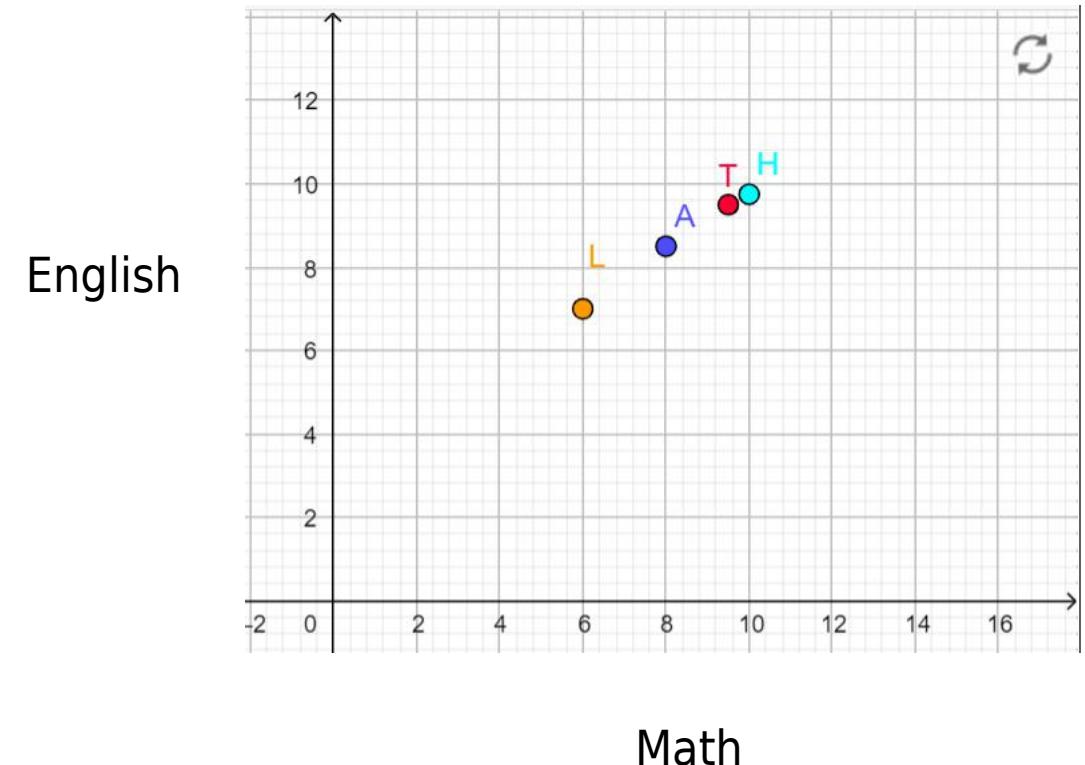
colData

	treatment	sex
ctrl_1	control	male
ctrl_2	control	female
exp_1	treatment	male
exp_2	treatment	female

Sample names:
ctrl_1, ctrl_2, exp_1, exp_2

Data of life

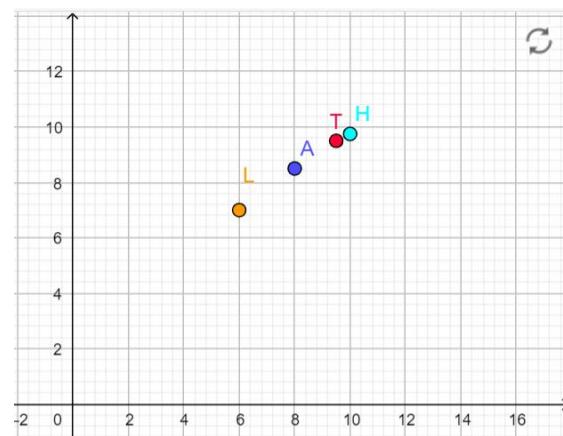
	Name	Math	English
1	Kim Anh	8	8.5
2	Hiep	10	9.75
3	Hanh	9.75	10 (French)
4	Loi	6	7
5	Truong	9.5	9.5



Data of life

	Name	Gender	Age	Math	English	History	isDoc
1	Kim Anh	f	16	8	8.5	10	No
2	Hiep	m	39	10	9.75	6	yes
3	Hanh	f	18	9.75	10 (French)	9.5	yes
4	Loi	m	41	6	7	8	No
5	Truong	m	34	9.5	9.5	9	yes

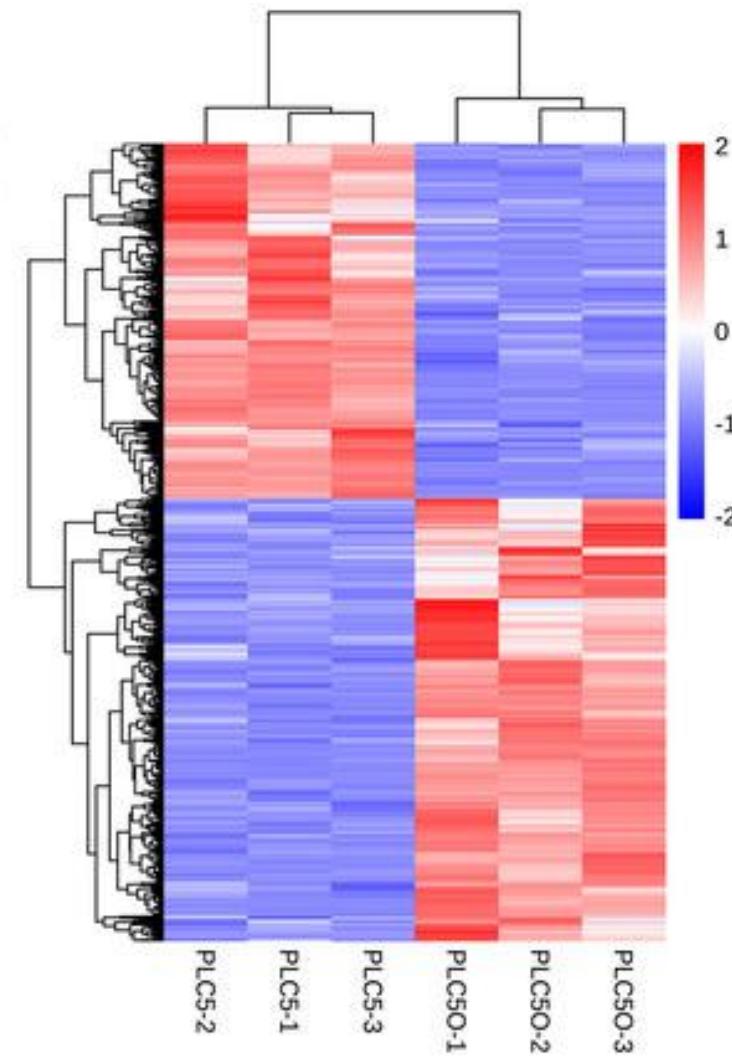
PC2 (14%)



PC1 (62%)

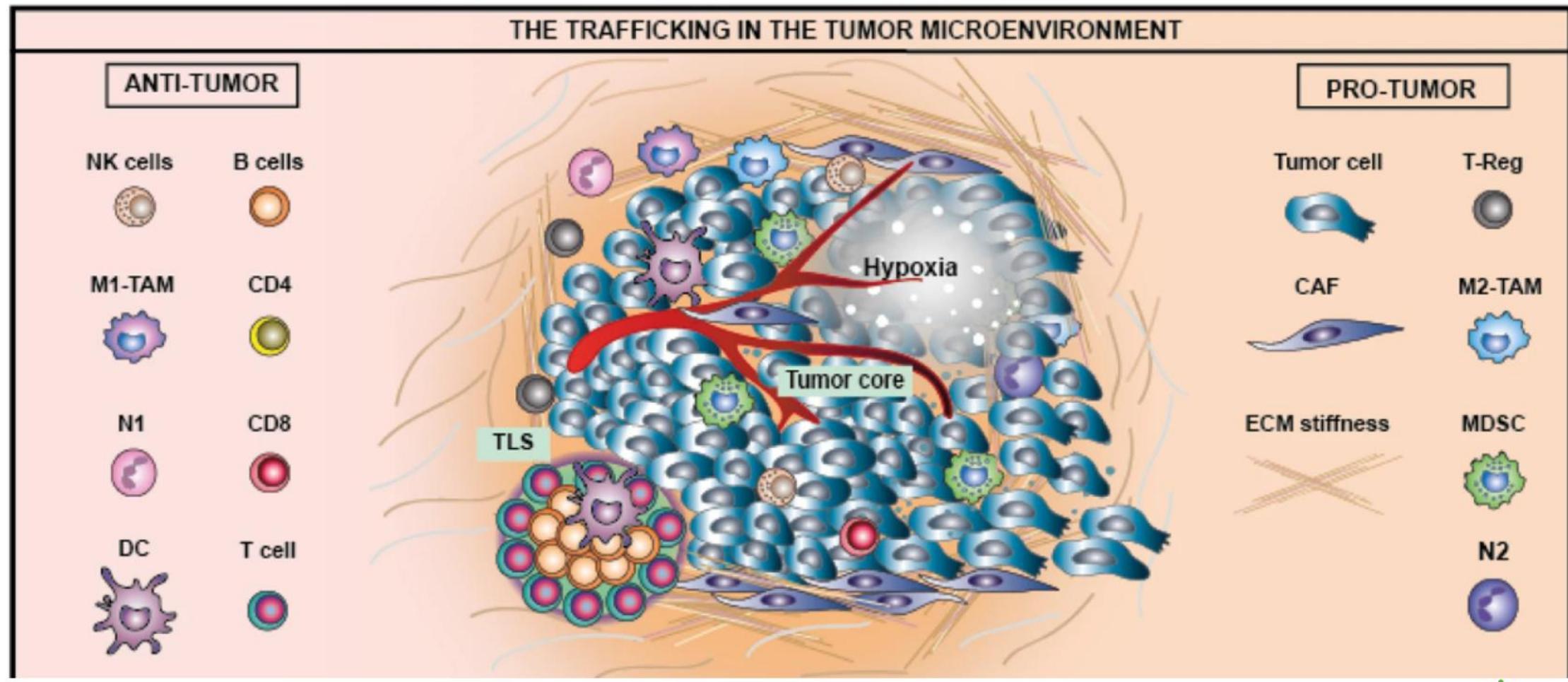
Heatmap

	A	B	C	D
1		2014	2015	2016
2	January	600	708	594
3	February	607	984	749
4	March	901	886	908
5	April	608	615	835
6	May	715	833	734
7	June	520	663	618
8	July	731	521	950
9	August	709	663	987
10	September	633	863	979
11	October	533	651	841
12	November	996	958	749
13	December	792	717	875



Single cell RNA Sequencing (scRNA-seq)

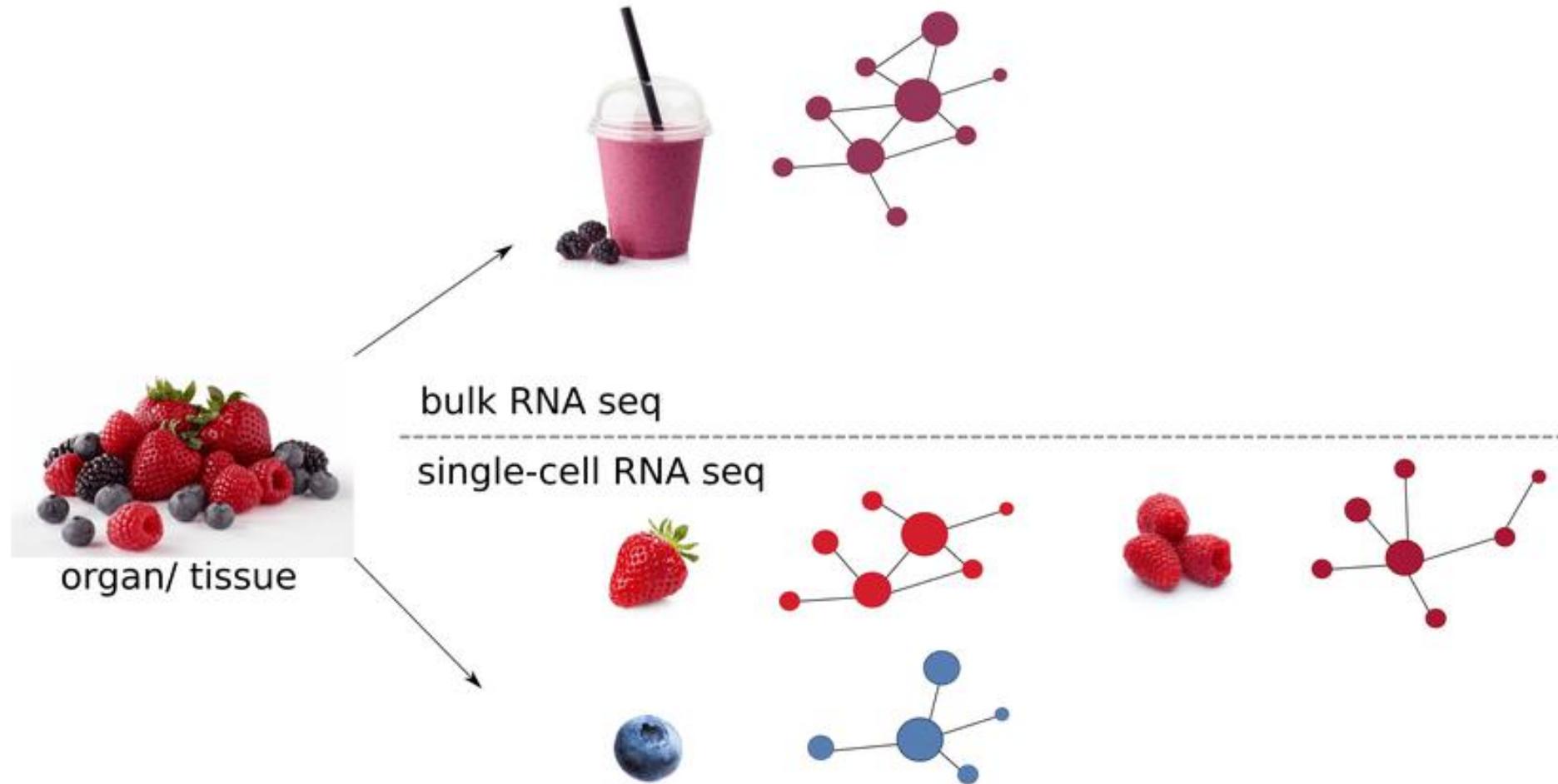
Single-cell applications: Tumor microenvironment



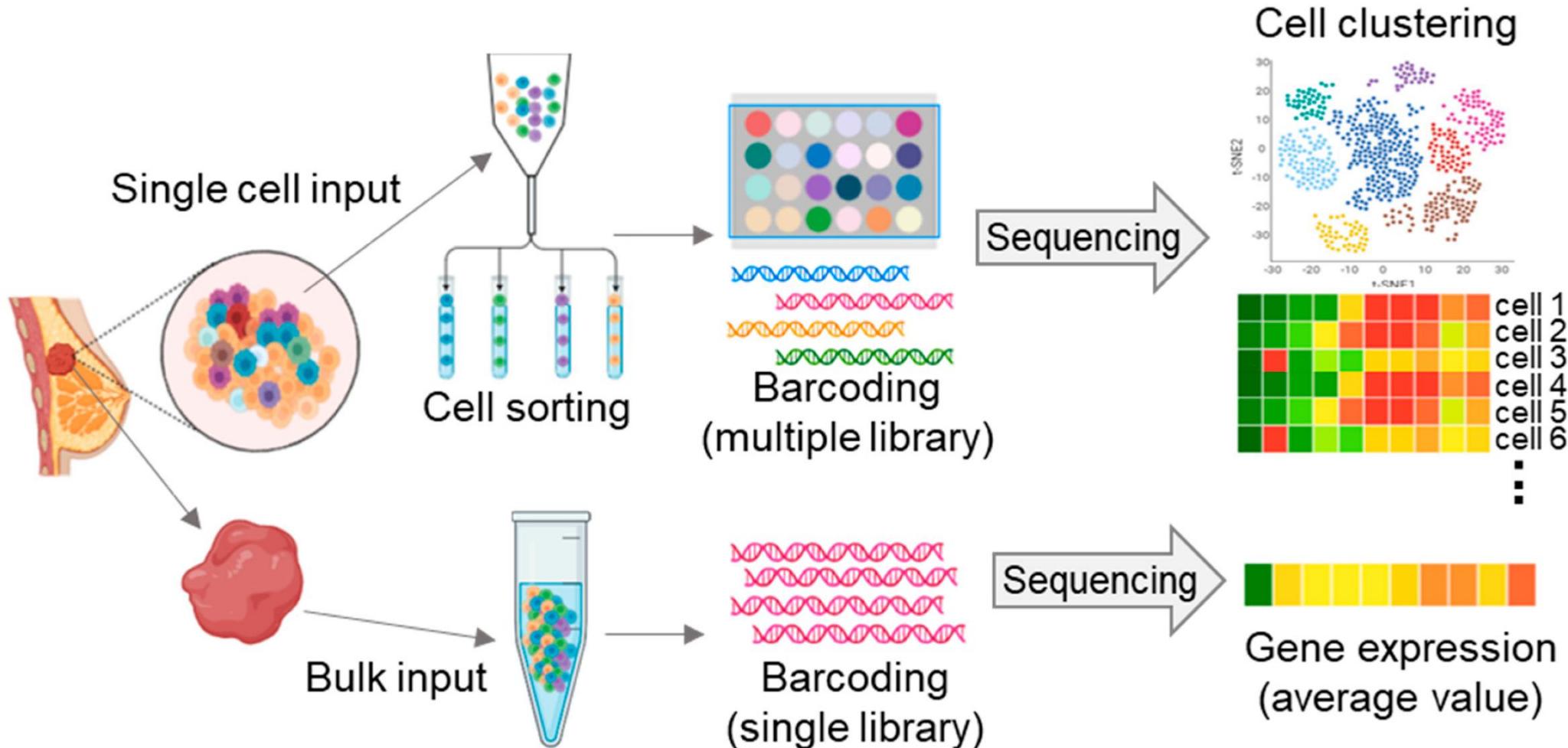
Source: <https://doi.org/10.1186/s13046-020-01586-y>

Tobias Rausch

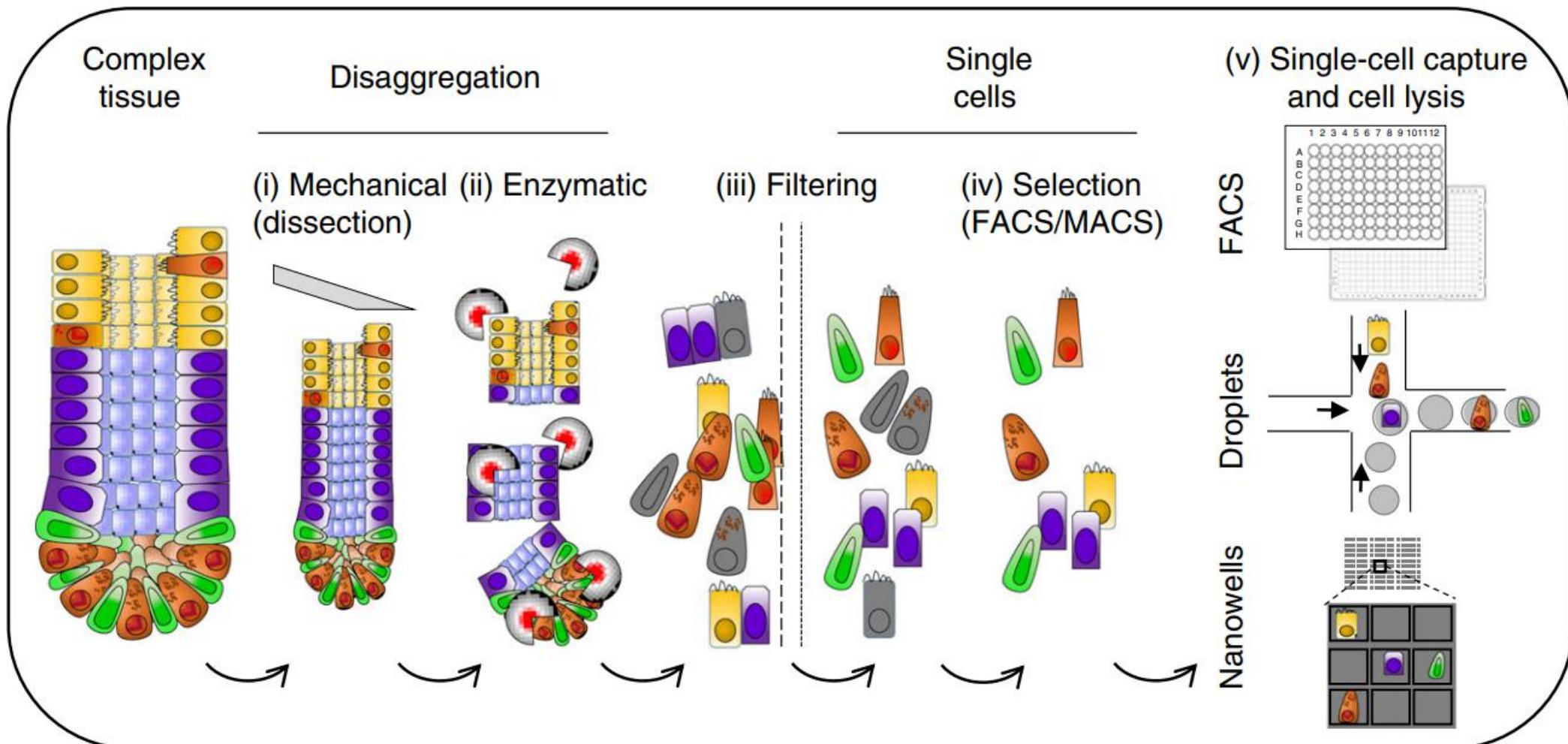
Bulk RNA-seq vs Single-cell RNA-seq



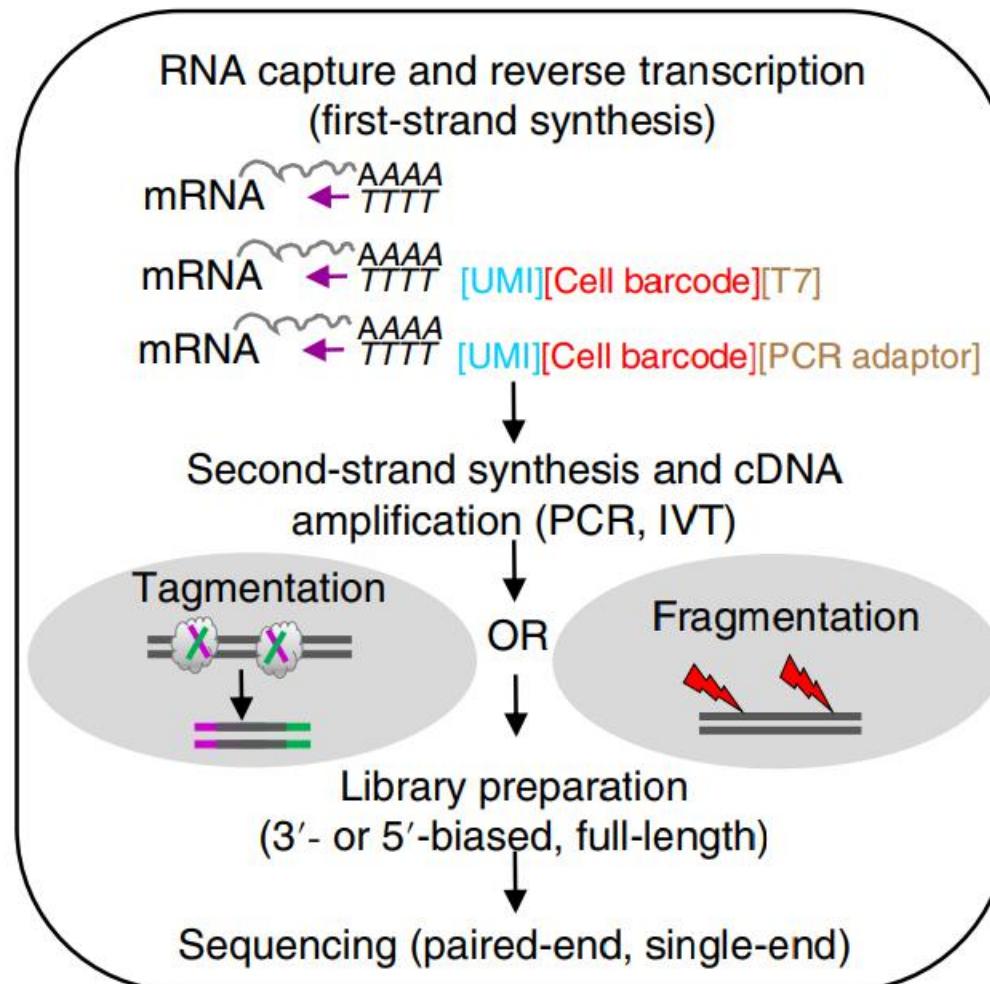
Bulk RNA-seq vs Single-cell RNA-seq



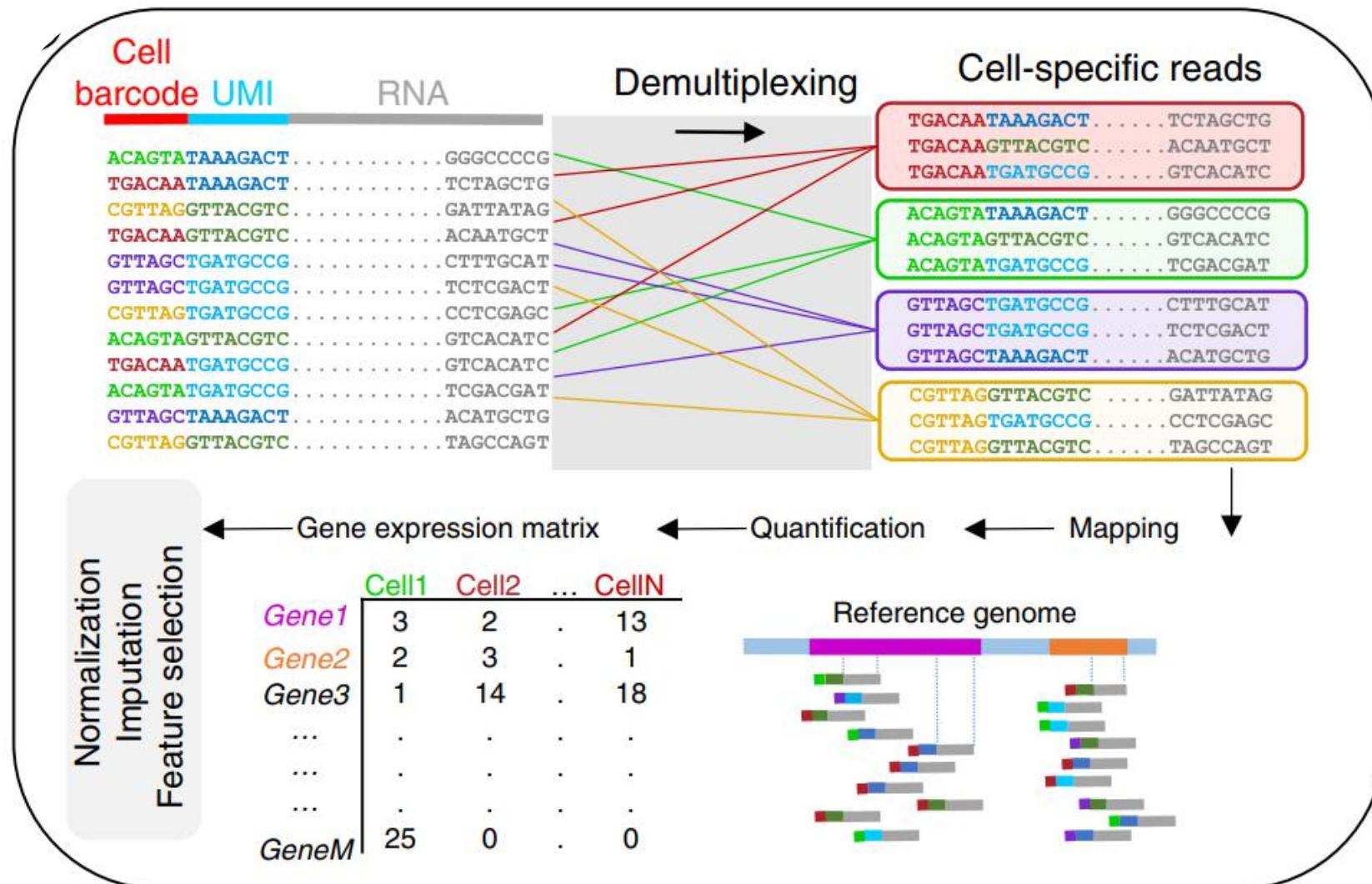
Single-cell RNA-seq: sample preparation



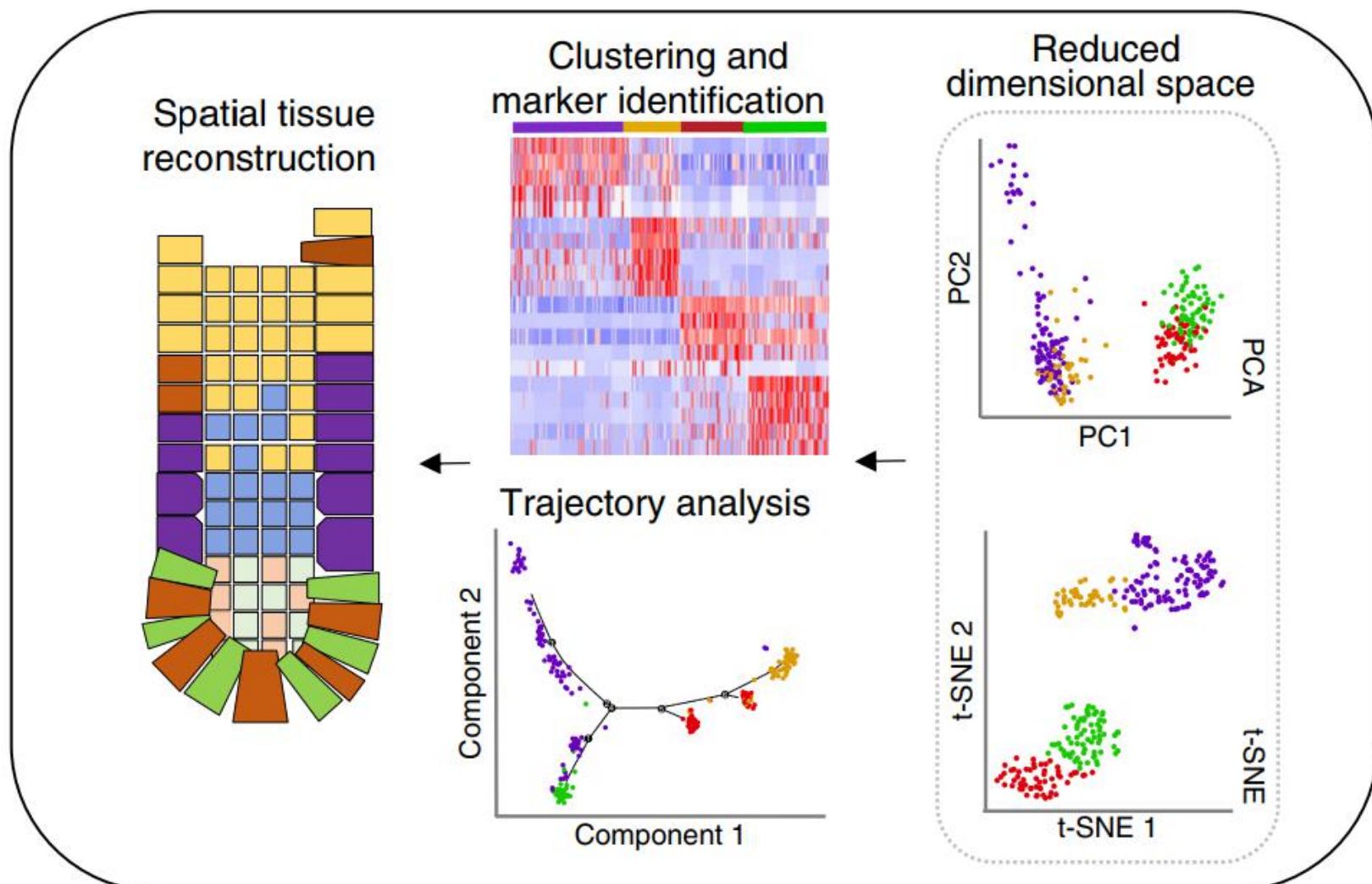
Single-cell RNA-seq: sequencing



Single-cell RNA-seq: data processing

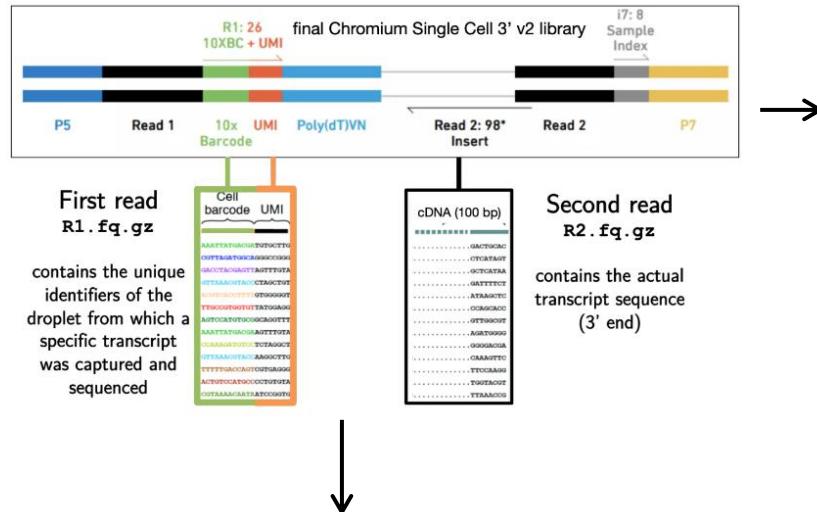


Single-cell RNA-seq: data processing

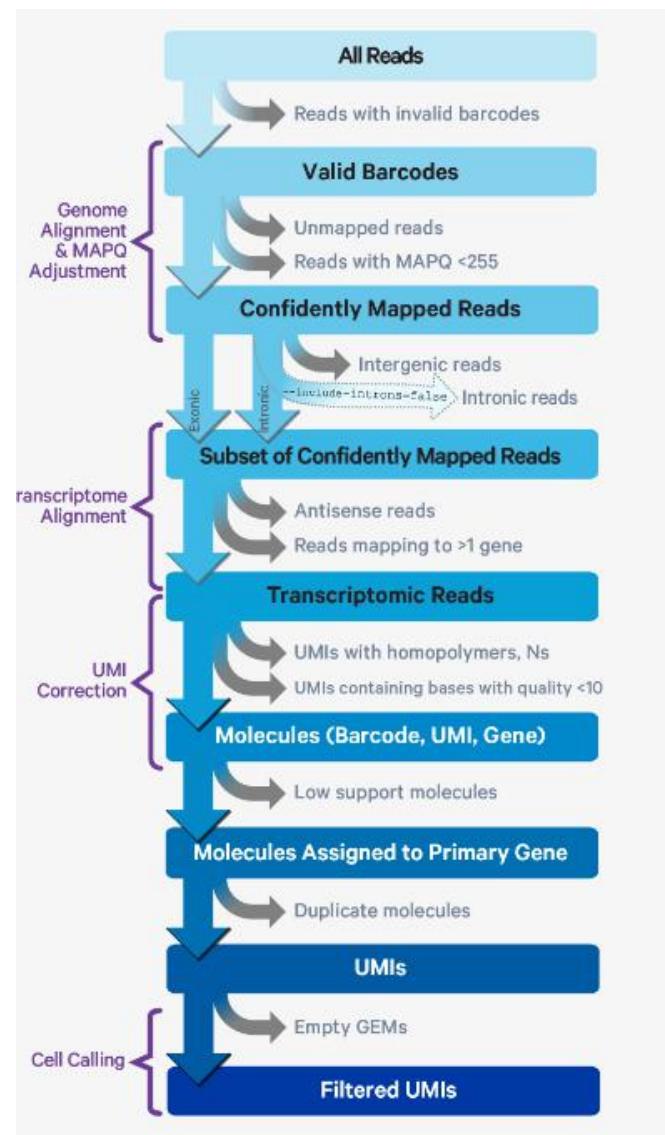


Upstream Analysis Process

Input: Fastq



Sequencing multi-QC



Cellranger Pipeline

Individual sample
CellRanger

CellRanger
Multi-QC report

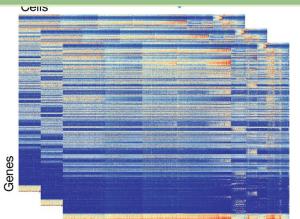
DATA PRE-PROCESSING



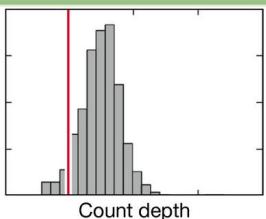
Raw
data
process



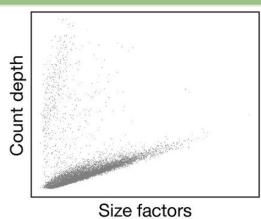
Alignment
to
transcripto-
me



Quantification into
Raw Counts Matrix

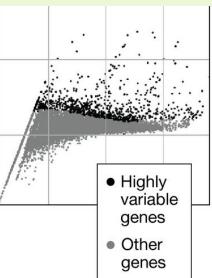


Quality control



Normalization

Feature selection

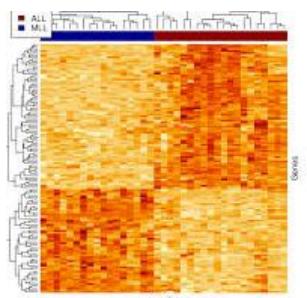


Visualization



Differentially gene
expressed features
(cluster biomarkers)

Gene
Ontology
analysis

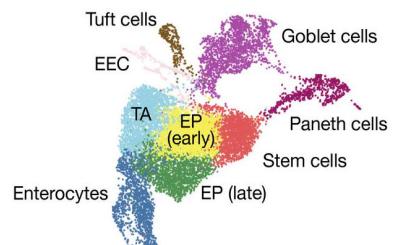


Pseudo-bulk
differential gene
expression analysis

SUBSETTING

Cluster
annotation

Clustering



DOWNTREAM ANALYSIS

Bulk ATAC-seq

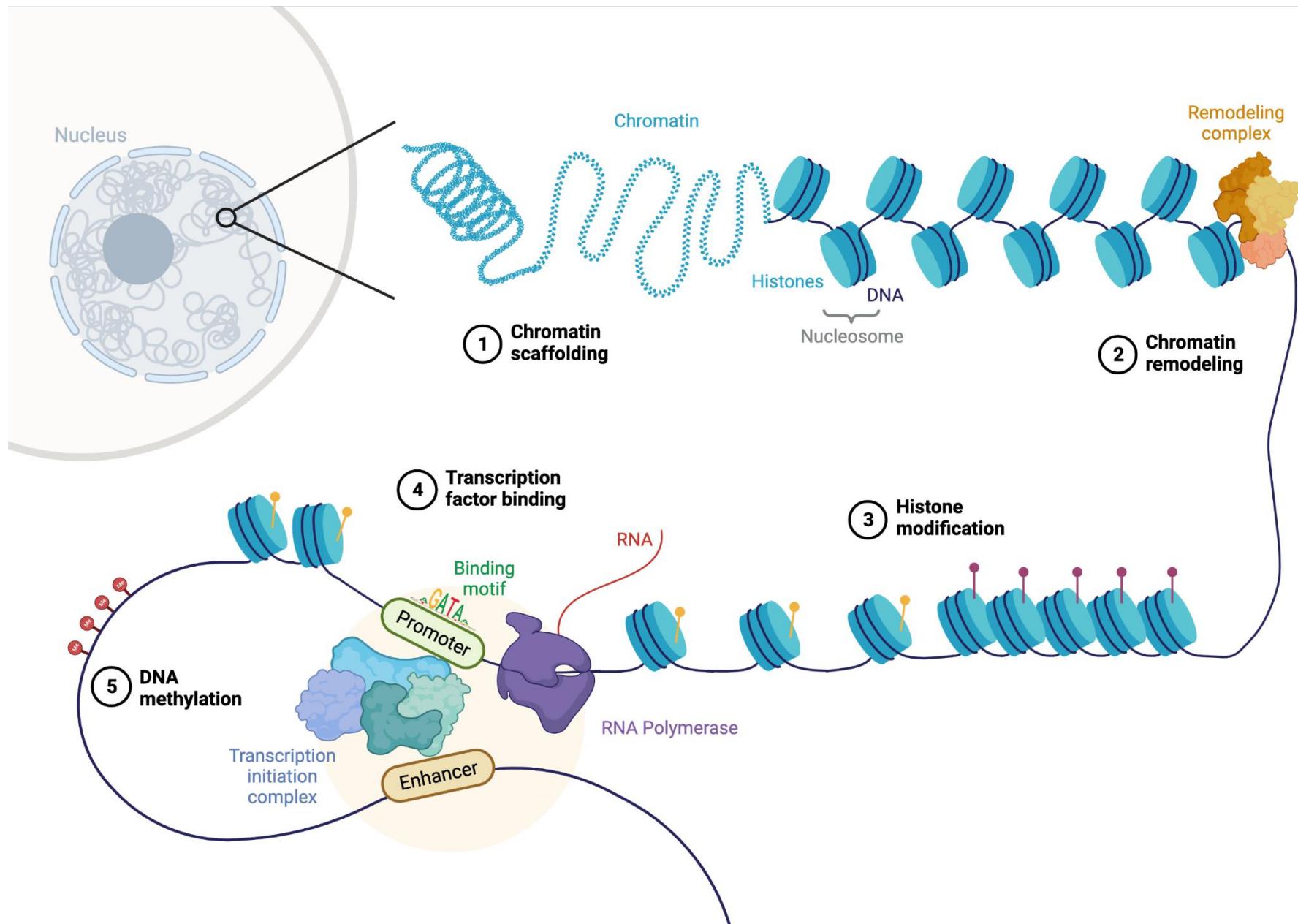


Figure. Overview of mechanisms influencing chromatin accessibility. Created with [BioRender.com](https://www.biorender.com).

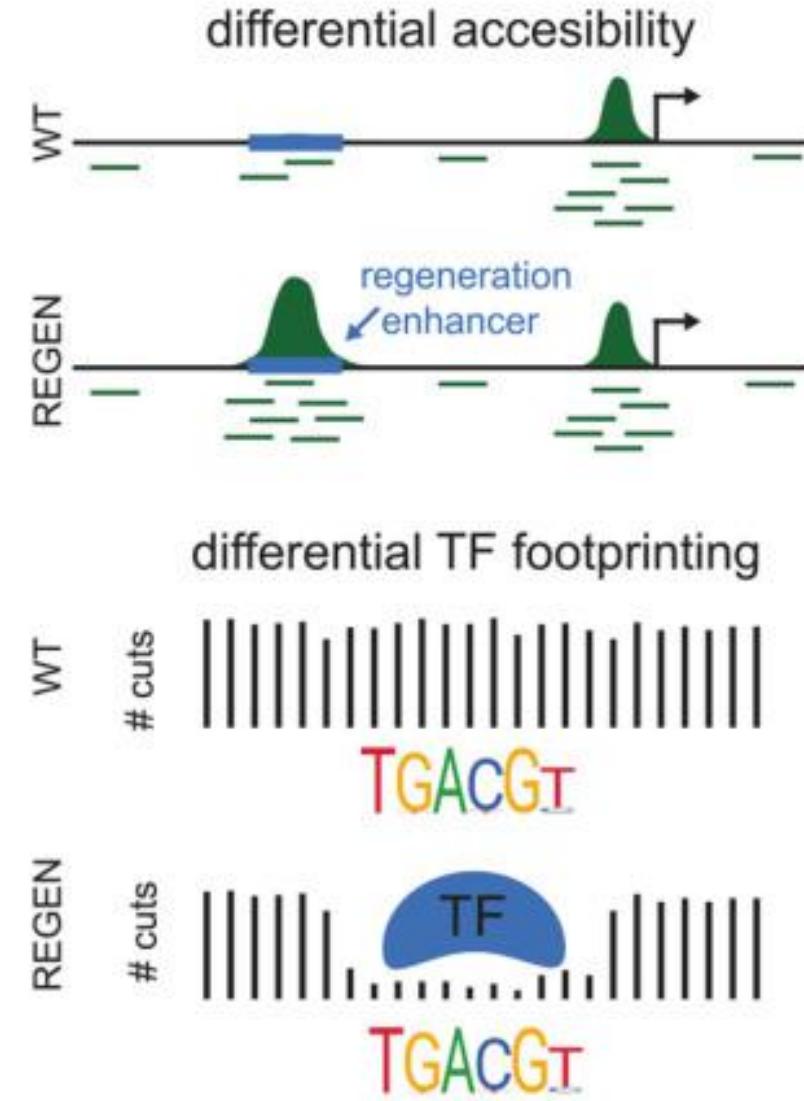
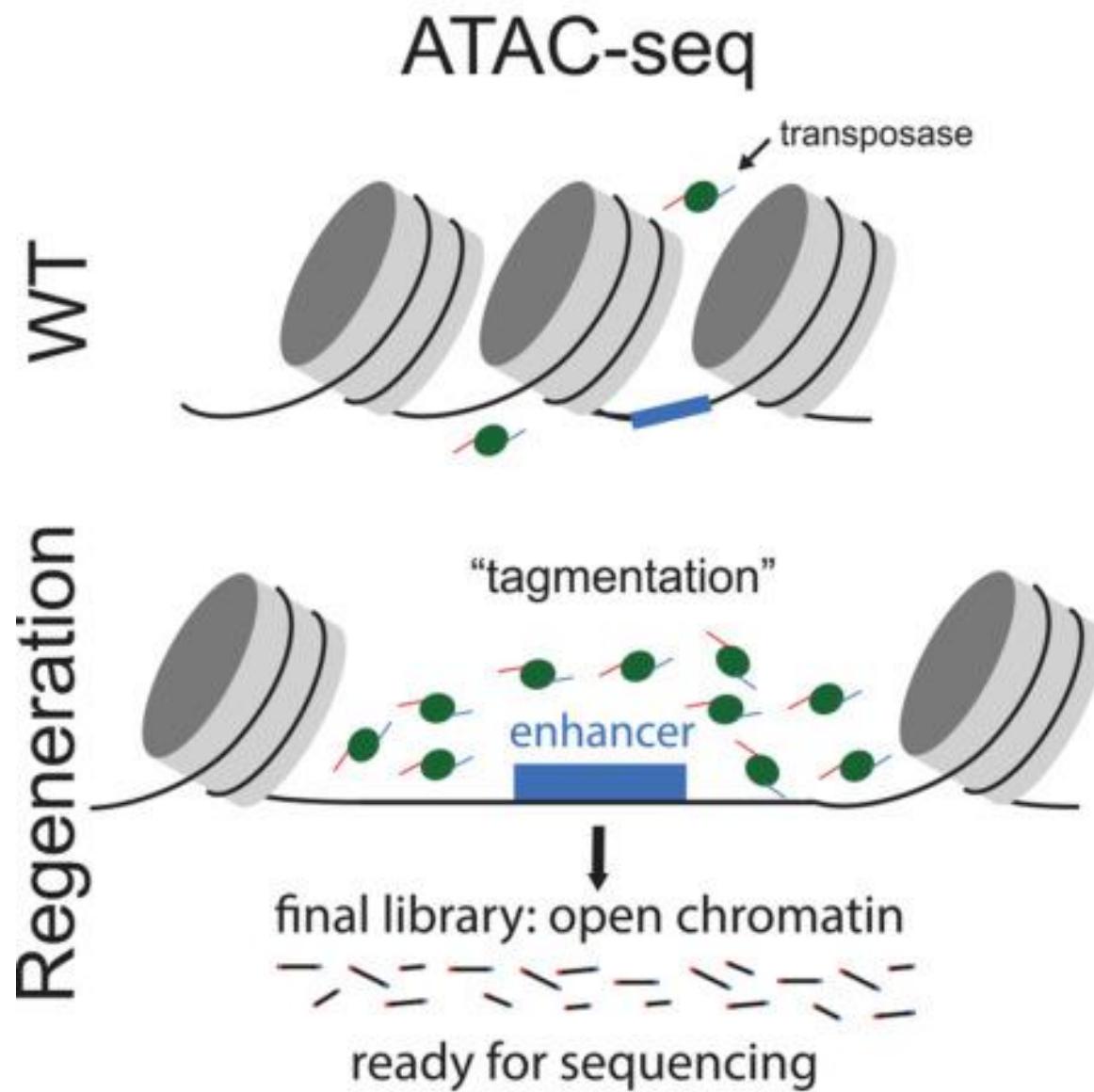


Figure. Overview of an ATAC-seq seq experiment to assay regeneration-responsive chromatin.

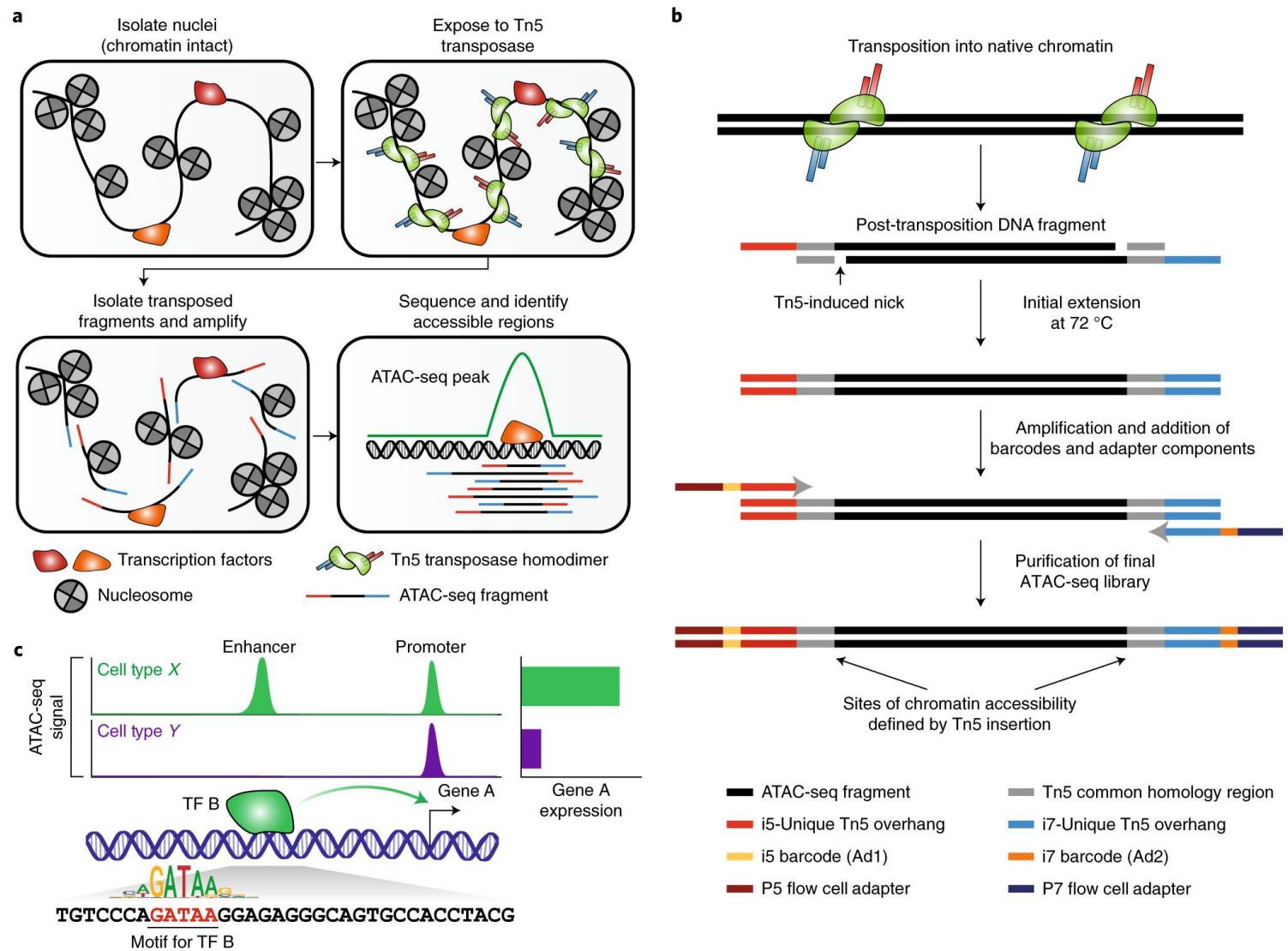
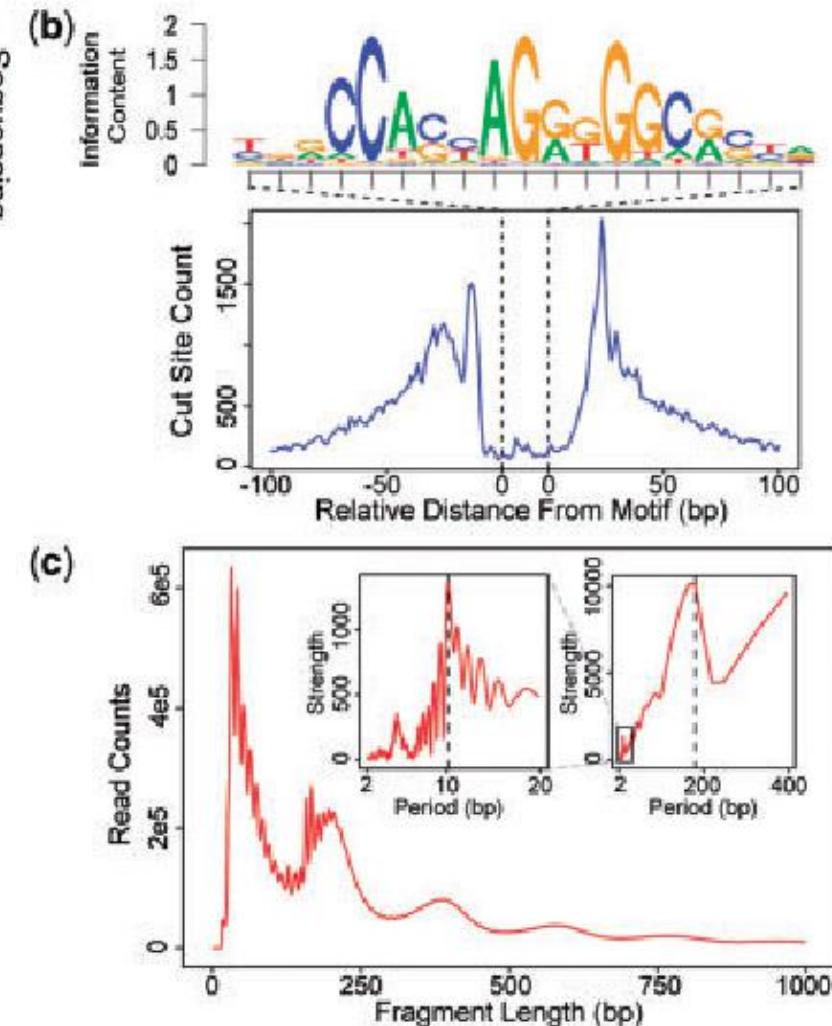
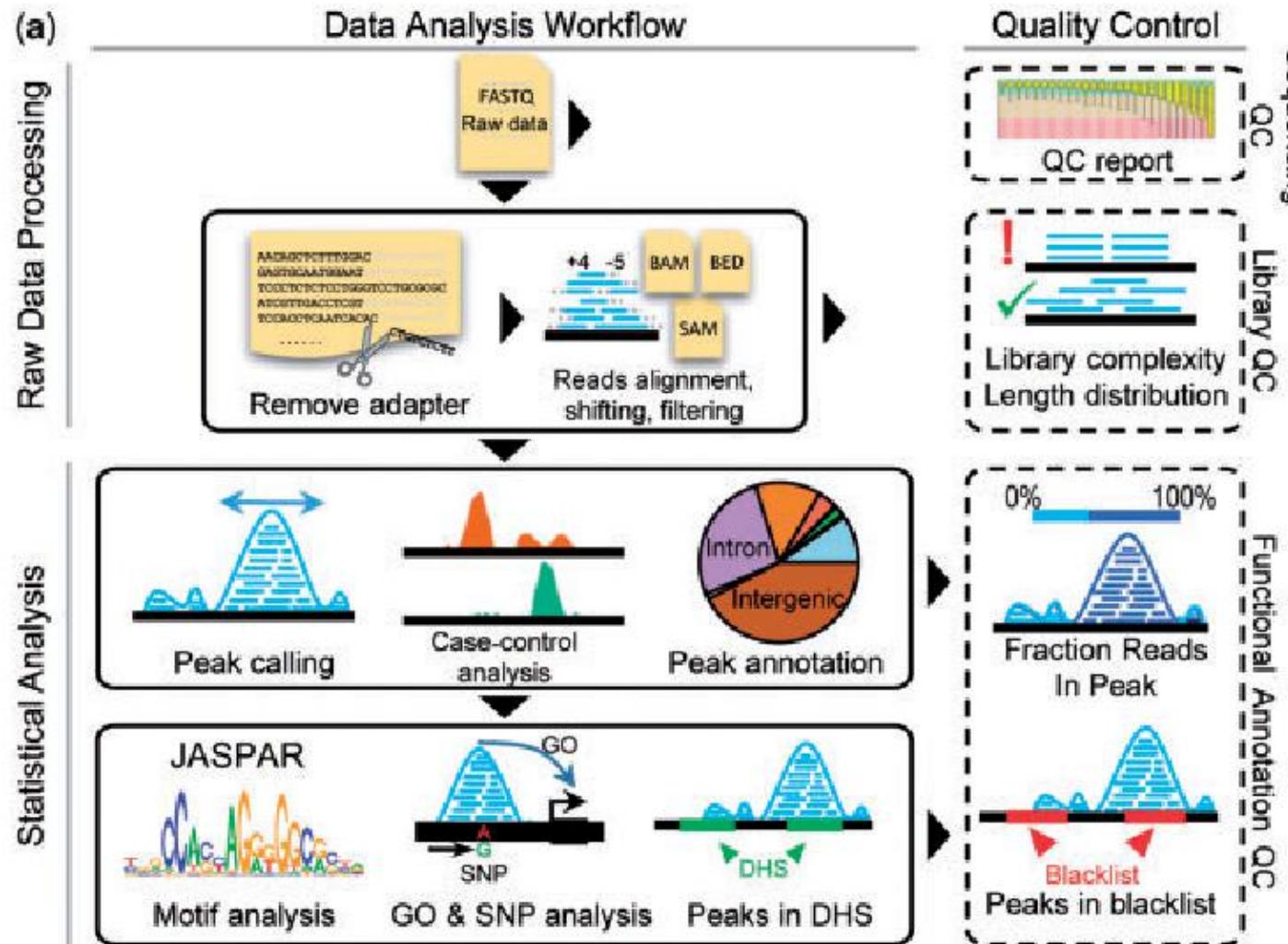


Figure. Schematic of the ATAC-seq transposition reaction and library preparation.

Bulk ATAC-seq analysis



Single-cell ATAC-seq

scATAC-seq analysis

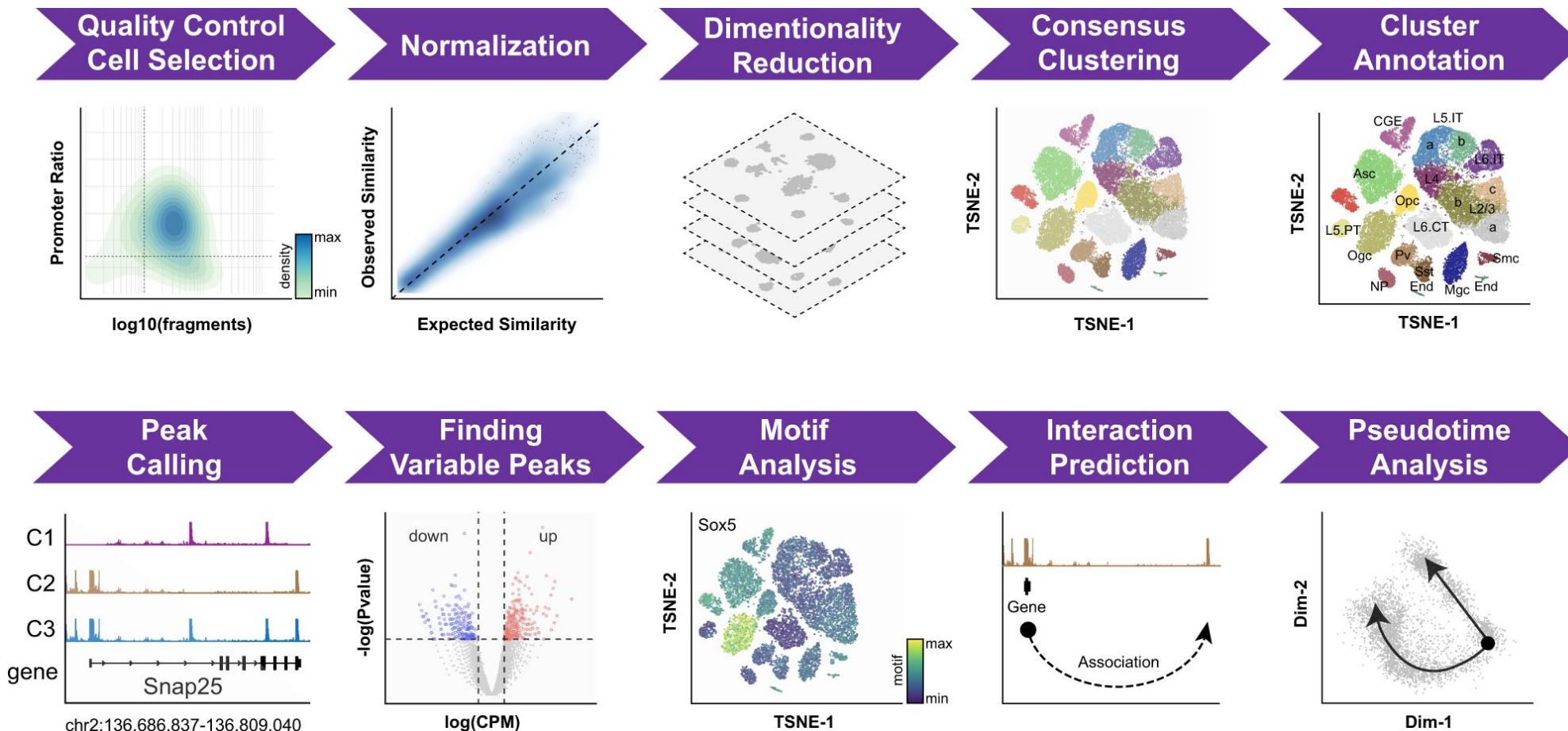
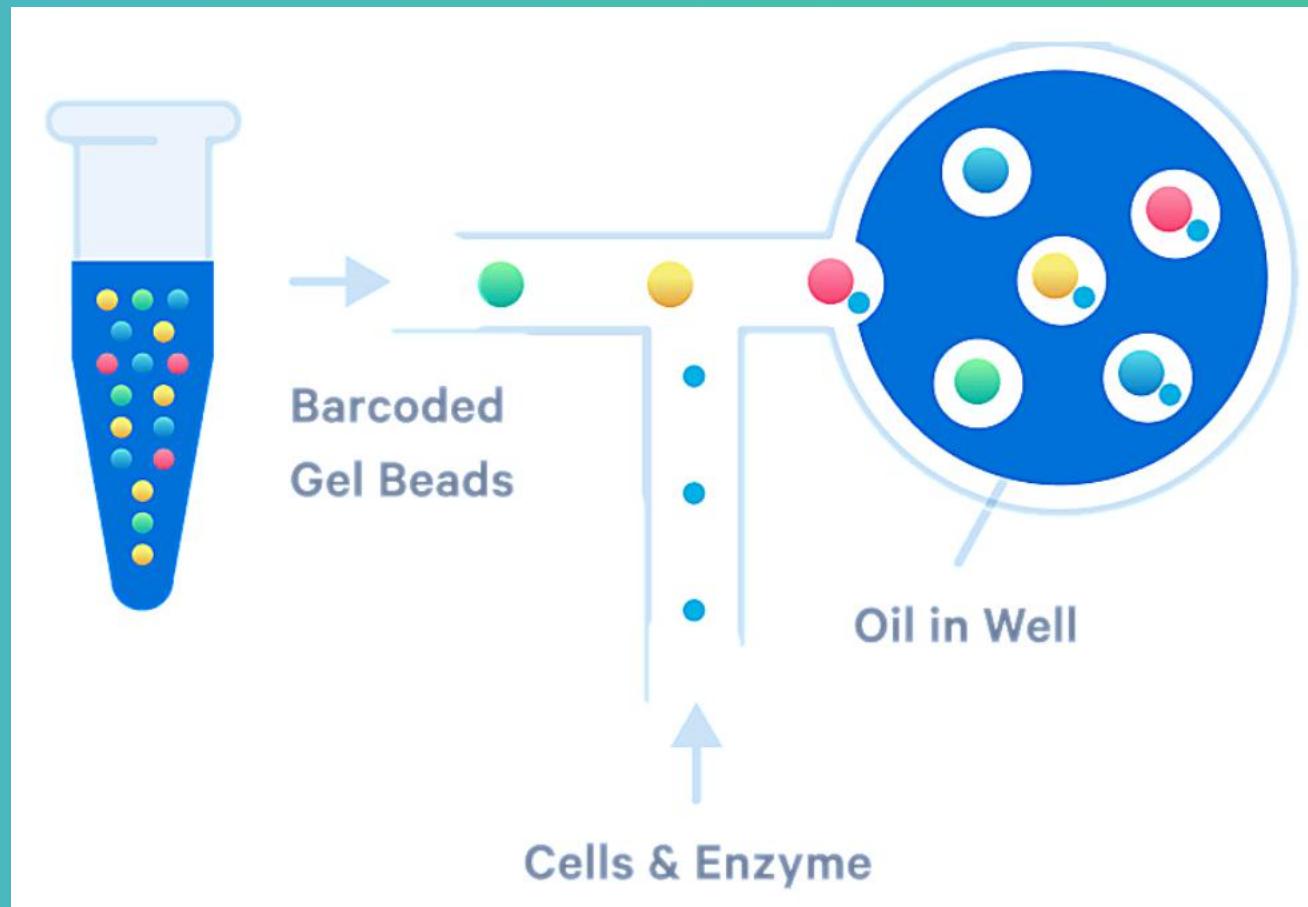
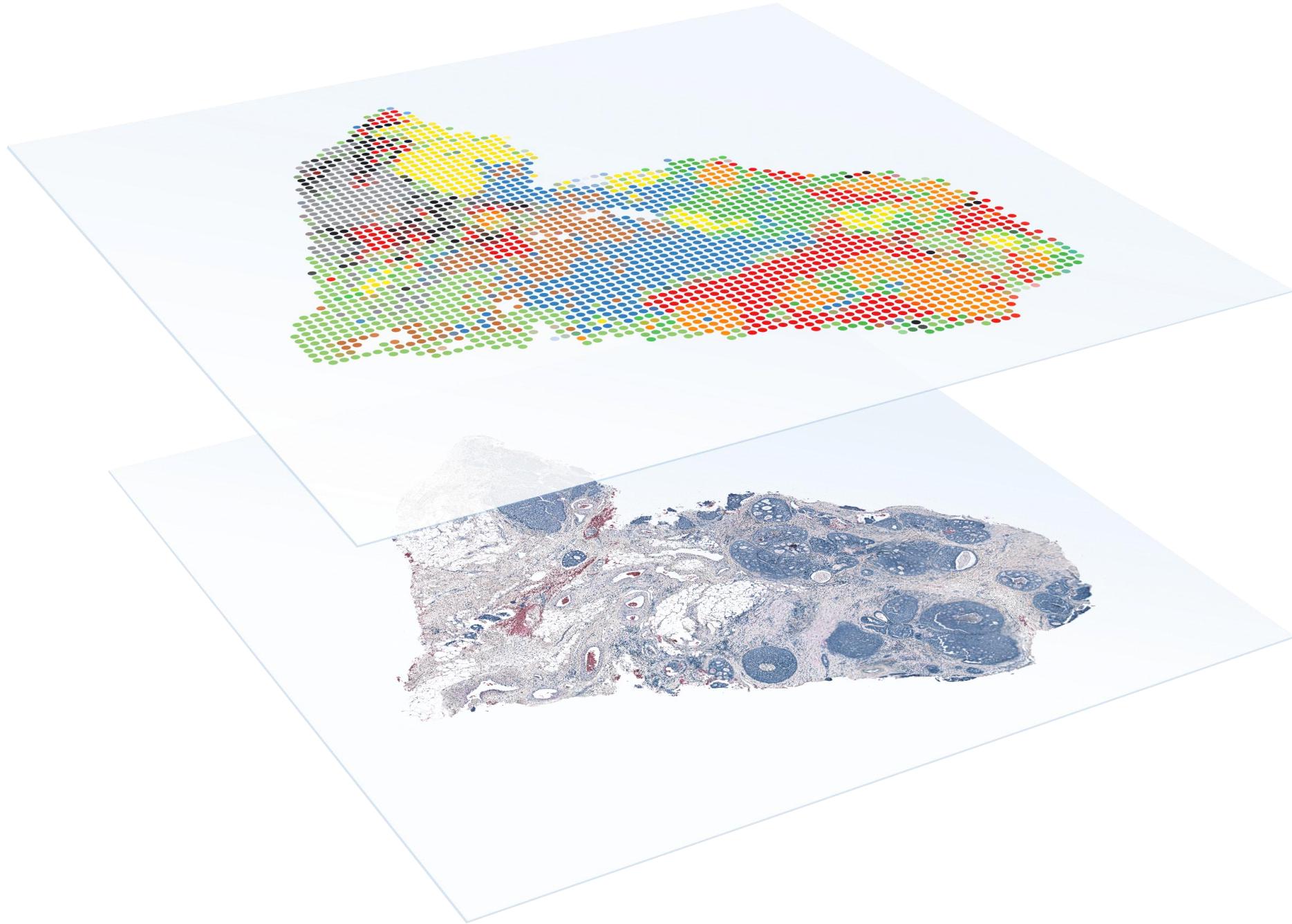


Figure: Schematic overview of SnapATAC analysis workflow.

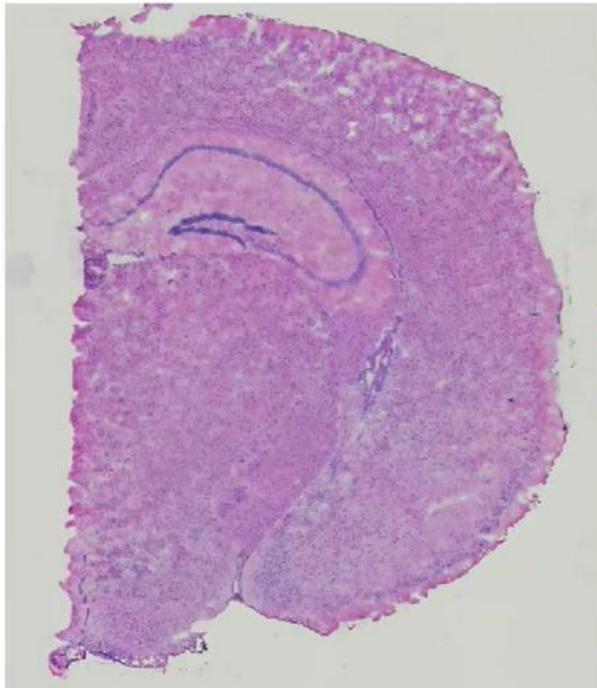
Chromium Single Cell



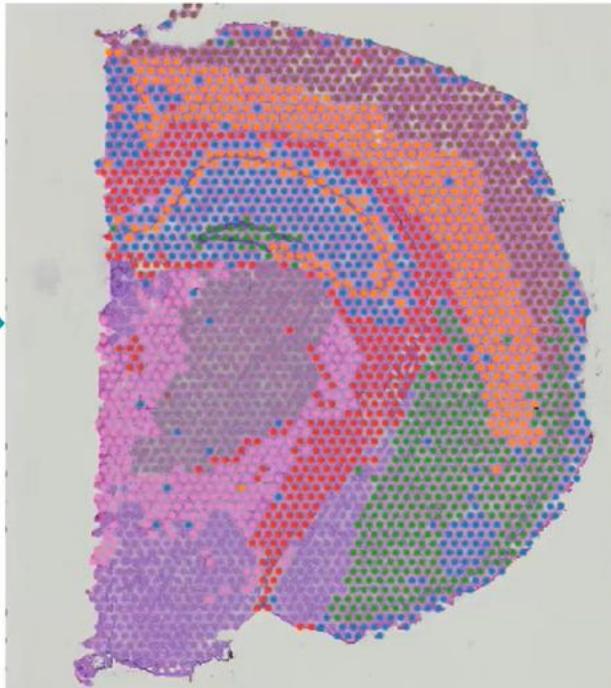


Spatial transcriptomics with Visium

Histology image



Histology + Gene expression



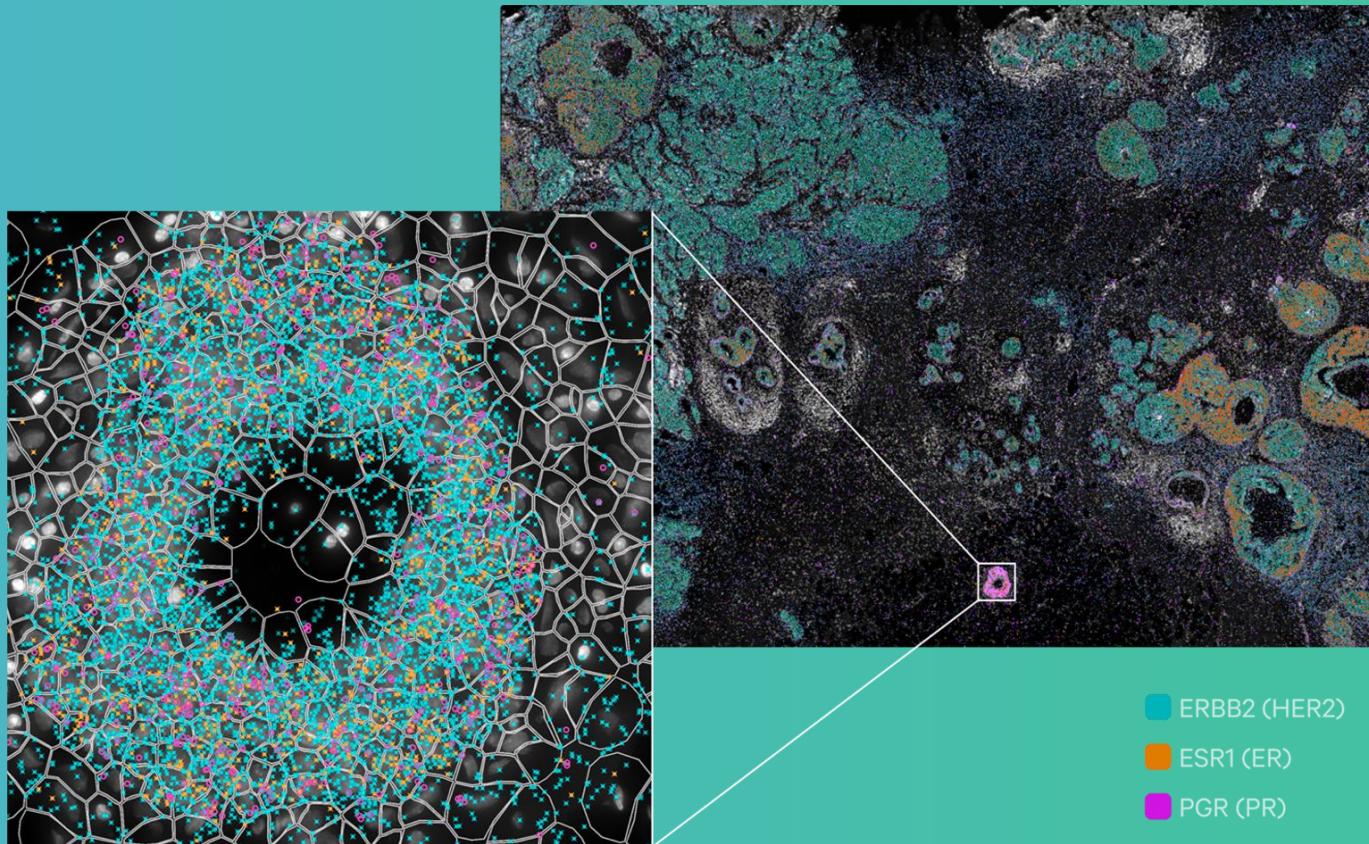
Gene expression



Each dot represents gene expression from one capture spot

Where do **active T cells** Locate?

Xenium In Situ



In this breast cancer sample
(Stage II-B, ER+/PR-/HER2+), Xenium identifies a
previously unknown triple-positive region.



**THANK YOU FOR YOUR
ATTENTION and Q&A!**

Please contact:loi.ip@pacificinformatics.com.vn
For further information!