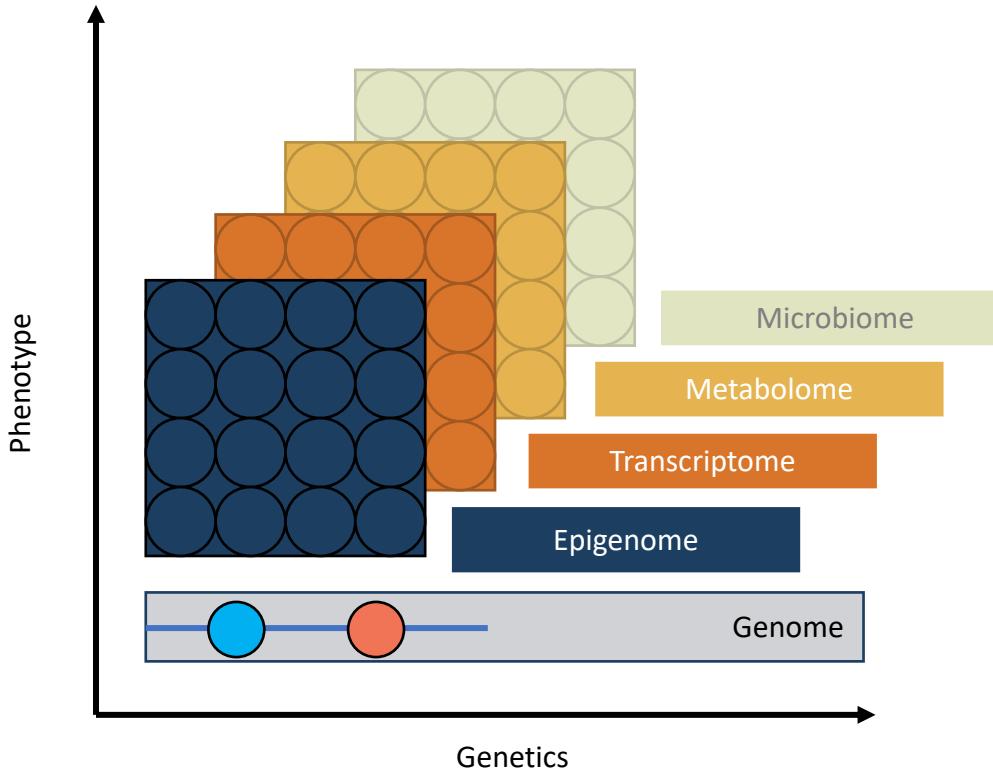
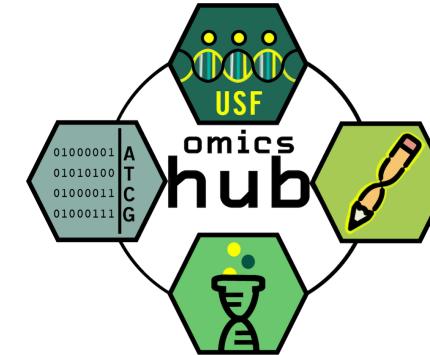


Applied Computational Genomics

(95555-PHC 6934-002, Fall 2020)

Time: 5:00 pm- 7:45 pm, Wednesday, Aug 24 – Dec 10, 2020



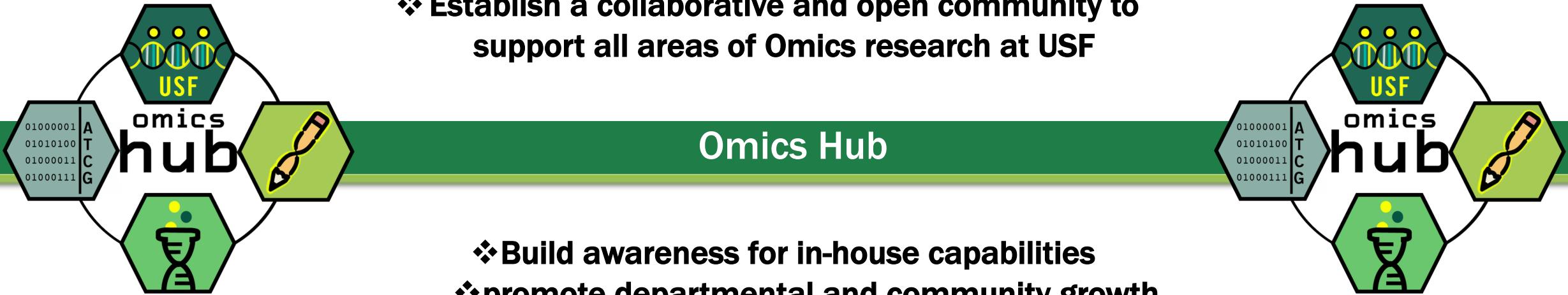
Instructor:

Chengqi(Charley) Wang, PhD
Jenna Oberstaller, PhD
Justin Gibbons, PhD
Thomas Keller, PhD
Anujit Sarkar, PhD
Min Zhang, MD

Teaching Assistant:
Jan Dahrendorff
Samira Jahangiri
Janelle Donglasan

Program Structure and mission overview

USF Genomics Program



- ❖ Build awareness for in-house capabilities
- ❖ promote departmental and community growth

- ❖ Computational and laboratory consulting in partnership with the Core
- ❖ Encourage grants submission as collaborators
- ❖ Provide training (bench, data analysis)

USF Genomics Core Facility

- ❖ Core equipment to which certified users have **SUPERVISED** access after training and project-consultation with our Core experts

Omics expertise, consultation, and community-building

Career biologists, bench scientists, bioinformaticians, programmers, and scientific writing experts here to help you design experiments, generate preliminary data, and collaborate with you to facilitate your omics research and strengthen the USF Genomics Program

Computational expertise

genomics@usf.edu



Jenna Oberstaller, PhD



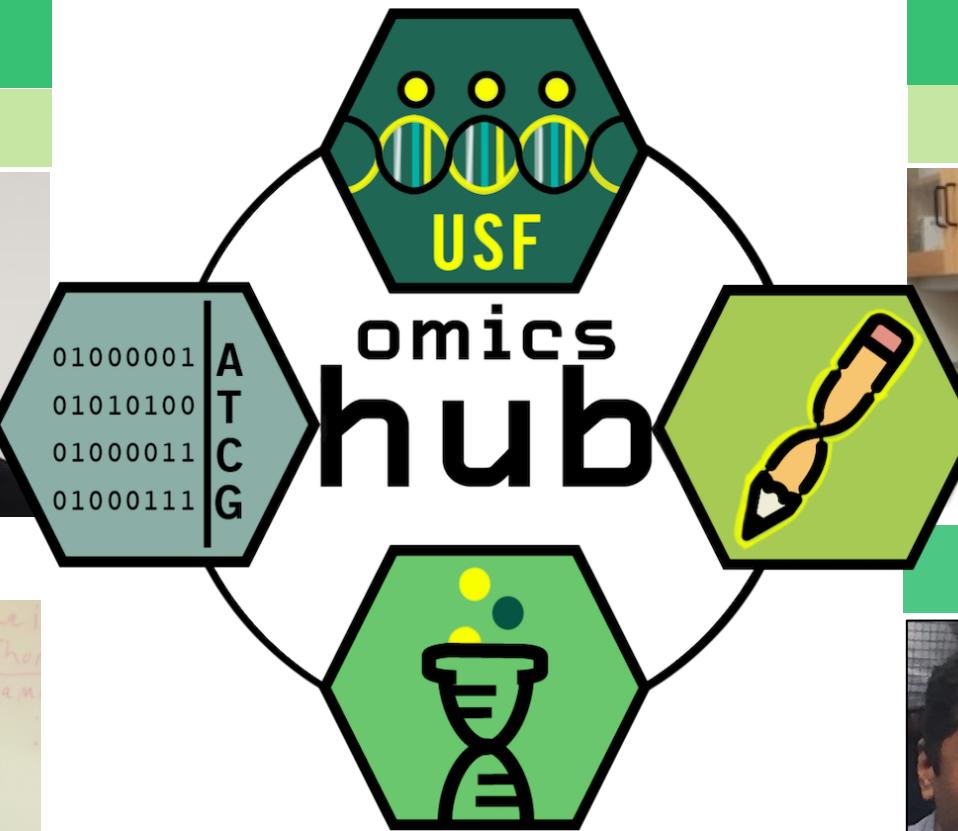
Thomas Keller, PhD



Chengqi(Charley)
Wang, PhD

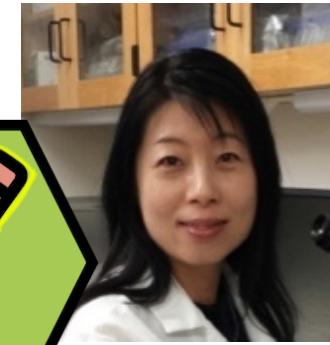


Justin Gibbons, PhD



Laboratory expertise

genomics@usf.edu



Min Zhang, MD
Core instrumentation
expert



Anujit Sarkar, PhD
All Things Microbiome

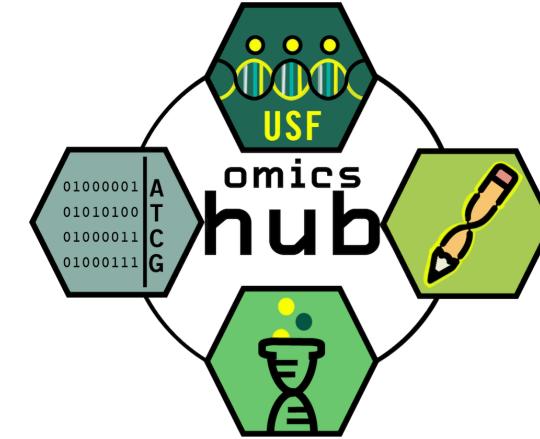


Swamy Rakesh
Adapa, MS
A founding Hub
contributor

genomics@usf.edu

Key Hub Course Concepts

- EXPERIMENTAL DESIGN is paramount
- BEST PRACTICES
 - at the bench
 - at the computer
- Data-analyses
 - QUALITY CONTROL
 - APPROPRIATE analyses to apply to answer YOUR biological question
 - many different methods, stats to be considered
 - REPRODUCIBILITY
- Training for biologists
 - we can't remove ALL the code
 - but focus = straightest path from data to biology/interpretation



USF GENOMICS CORE FACILITIES AND RESOURCES

- ❖ Supervised access to sequencing core
- ❖ Grant-application support
- ❖ Consultation for sequencing and data analysis

FREE 2 hours per USF PI!



TapeStation /
Bioanalyzer



M220 Focused-
Ultrasonicator



Qubit® 2.0
Fluorometer



Roche LightCycler 96
qPCR



Illumina - MiSeq
System



NextSeq 550
System



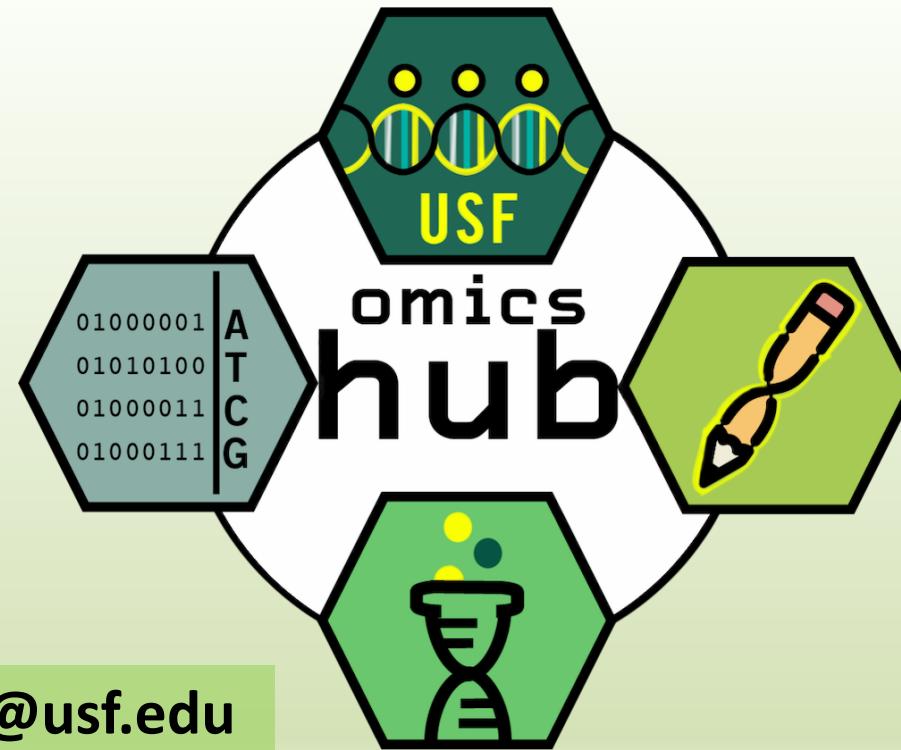
10x Chromium
SC

- ❖ Whole-Genome Sequencing
- ❖ RNA Sequencing
- ❖ Targeted Gene Sequencing (amplicon, gene panel)
- ❖ Gene Expression Profiling
- ❖ Targeted Gene Expression Profiling
- ❖ miRNA & Small RNA Analysis
- ❖ Methylation Sequencing
- ❖ 16S Methylation Sequencing
- ❖ Single-Cell Sequencing
- ❖ DNA, RNA Quantification
- ❖ Sequencing Consultation

USF Genomics Core Computational Facilities and Resources

- ❖ Omics Consultation
- ❖ Transcriptomics
- ❖ Epigenomics
- ❖ Phylogenomics
- ❖ Microbiome and Metagenomics
- ❖ Genome Assembly and Annotation
- ❖ Custom Scripts and Analysis Pipelines
- ❖ Data Visualization
- ❖ Genomics Training

genomics@usf.edu



Core Computational Resources:

- USF's High-Performance Computing (HPC) Cluster, managed by the USF Department of Research Computing**
- Sequence data are securely stored on HIPAA-compliant architecture of USF's HPC cluster**

USF Genomics home page



USF Genomics

Home

Researchers

Grants

Degree Programs >

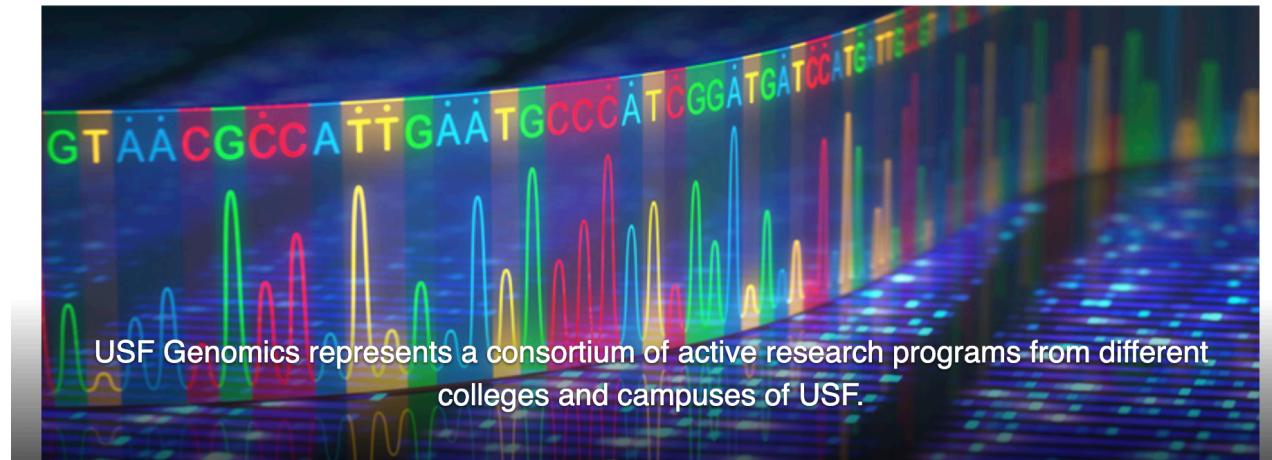
Undergraduate Opportunities >

Annual Symposium

College of Arts and Sciences
News

Genomics Training

Request Services



USF Genomics represents a consortium of active research programs from different colleges and campuses of USF.

About USF Genomics

Our projects integrate modern genomic approaches to understand and help develop solutions for some of the major challenges affecting global health and especially our Florida community. USF Genomics research projects are naturally interdisciplinary with studies that cross many boundaries of traditional scientific fields, integrating bioscience approaches in experimental and environmental research with computation biology. Opportunities exist for new graduate students to join well-funded research programs to study disease causes from underlying drug resistance and virulence genes in pathogenic organisms to changes of microbial communities within our bodies to the marine environment.

Register and Submit Abstracts for the USF Genomics Annual Symposium on November 15, 2019

Upcoming Events

USF Genomics Seminars: Tuesdays once a month, 4:00-5:00 pm - IDRB 302

USF Genomics Seminar - Dr. Maureen Groer - Microbiome State of the Science - Tuesday, September 17, 4:00-5:00 pm - IDRB 302

USF Genomics Forum Discussion Series - Tuesdays once a Month, 5:00-6:00 pm - IDRB 302

USF Genomics Annual Symposium: Personal Genomics - November 15, 2019

Genomics Laboratory and Data-Analysis Training Courses: February, May, September 2019

Introduction



Instructor:

- Chengqi(Charley) Wang, PhD
- My own research

Functional Genomics, Epigenetics, Machine Learning and Deep Learning algorithm



Instructor:

- Jenna Oberstaller, PhD
- My own research

Functional genomics, evolutionary genomics, systems biology, omics tools-development for biologists



Instructor:

- Justin Gibbons, PhD
- My own research

He is currently working on several projects ranging from drug resistance in malaria to how the microbiome effects the health of preterm infants

Introduction



Instructor:

- Thomas Keller, PhD
- My own research

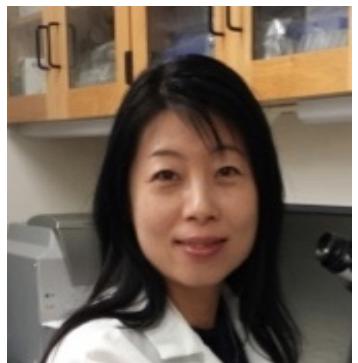
Epigenomics, microbiome, clinical applications to machine learning & deep learning



Invited lecturer:

- Anujit Sarkar, PhD
- My own research

Interrelationship between cardiovascular disease, T1D and oral microbiome; Role of gut microbiome in the long term growth and development of preterm infants



Invited Lecturer:

- Min Zhang, MD
- My own research

Molecular Genetics, Next Generation Sequencing, Single-cell Sequencing, Malaria parasite biology, Genetic analysis of critical metabolic process

Introduction



TA:

-- Jan Dahrendorff

-- My own research

Differential gene expression in longitudinal cohort of post traumatic stress disorder

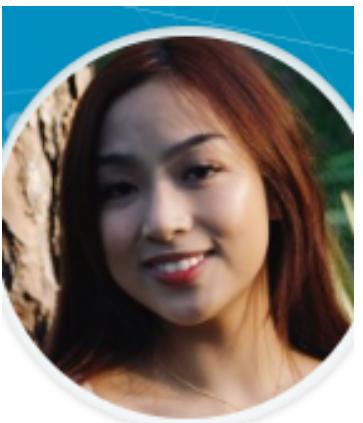


TA:

-- Samira Jahangiri

-- My own research

Variant analysis in Genome-wide associated studies, differential gene expression analysis and the functional study of disease causing mutations using high-throughput data.



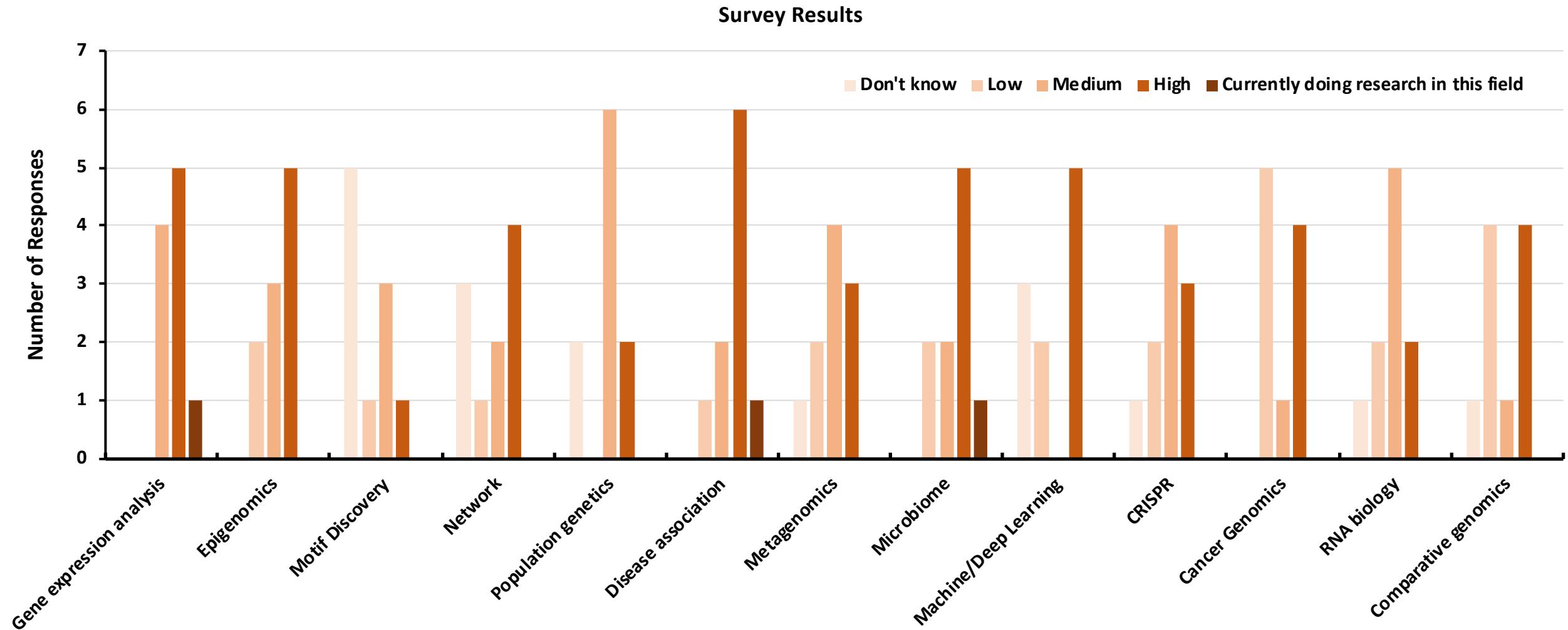
TA:

-- Janelle Donglasan

-- My own research

Identify genetic and epigenetic predictors of stress-related mental disorders, with a particular focus on depression and post-traumatic stress disorder

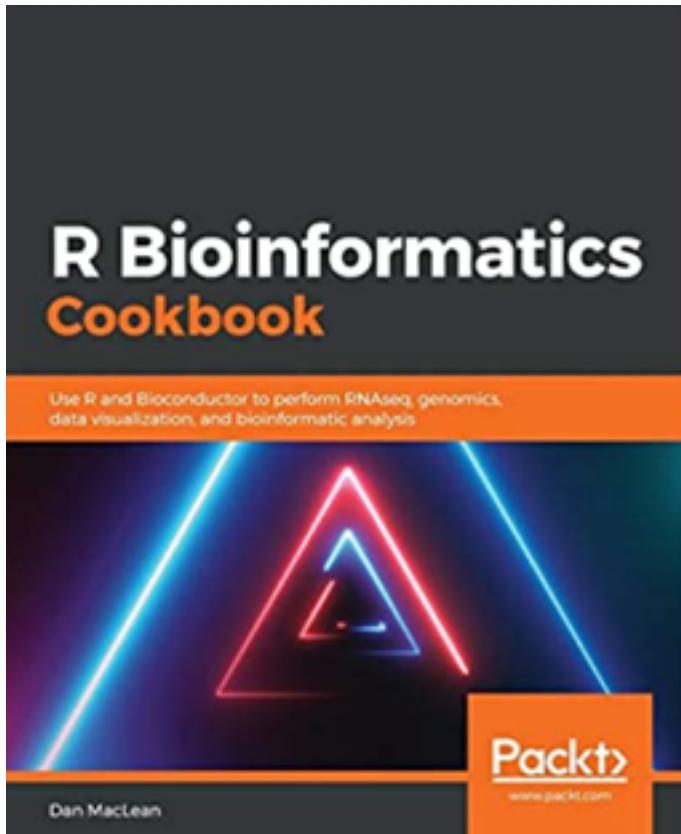
Students (From first-day survey)



Interest in specific topic

	Don't know	Low	Medium	High	Currently doing research in this field	Total
Gene expression analysis			3	5	1	9
Epigenomics		2	2	5		9
Motif Discovery	4	1	3	1		9
Network	2	1	2	4		9
Population genetics	2		5	2		9
Disease association		1	1	6	1	9
Metagenomics	1	1	4	3		9
Microbiome		2	1	5	1	9
Machine/Deep Learning	2	2		5		9
CRISPR	1	2	3	3		9
Cancer Genomics		5		4		9
RNA biology	1	2	4	2		9
Comparative genomics	1	3	1	4		9
Total	14	22	29	49	3	117

(Optional) Book for the Course



R Bioinformatics Cookbook: Use R and Bioconductor to perform RNAseq, genomics, data visualization, and bioinformatic analysis

Maclean Dan

Grading Criteria

Assessment	Percent of Final Grade
Project 1	30%
Project 2 (Final)	40%
In-Class Assignment	20%
Participation and Attendance	10%
Total	100%

In-class Assignment

Each class will have an In-Class Assignment to complete the code discussed in class (shell or R-script).

Each assignment must be submitted to the canvas within 72 hours.

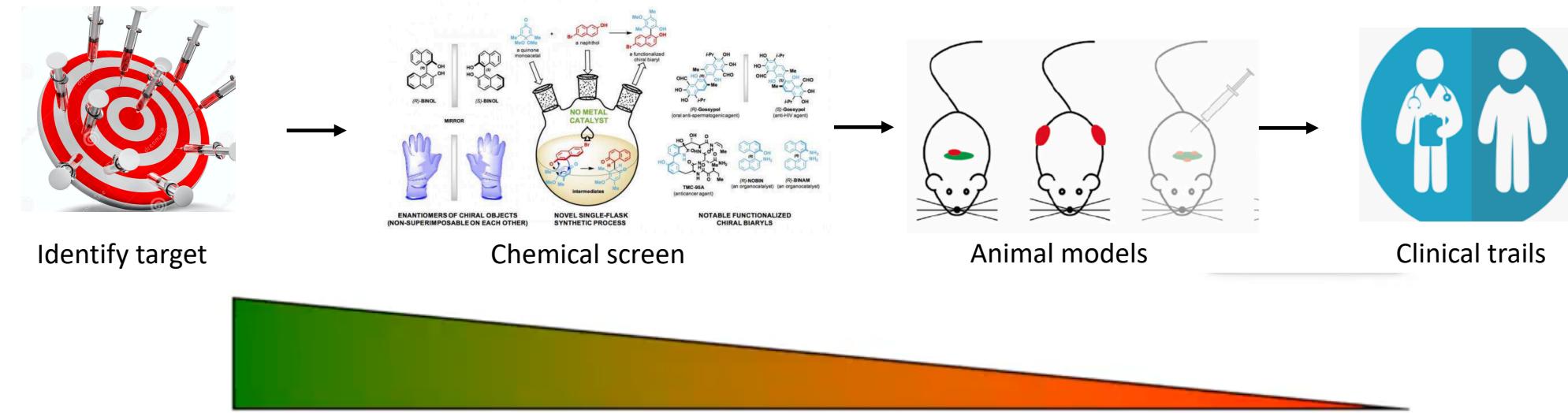
Lesson 2 [9/2/20]	Chengqi (Charley) Wang TA: Samira Jahangiri	<p>Part 1: Introduction to the Linux directory system and how to get around it</p> <p>Part 2: Text Editor (Vim) in Linux</p> <p>Part 3: Piping and Redirection in Linux</p> <p>Part 4: Slurm work load manager</p>	<p>Pre-class reading: Basic knowledge of Shell script (Learning the Shell) http://linuxcommand.org/lc3_learning_the_shell.php</p> <p>In-class assignment</p>	7:45 pm [9/5/20]
----------------------	--	--	--	---------------------

Q&A

Q&A (Every Thursday evening from 8:00pm -8:40 pm)

**How can genetics research help meet
outstanding medical challenges**

Informing therapeutic development

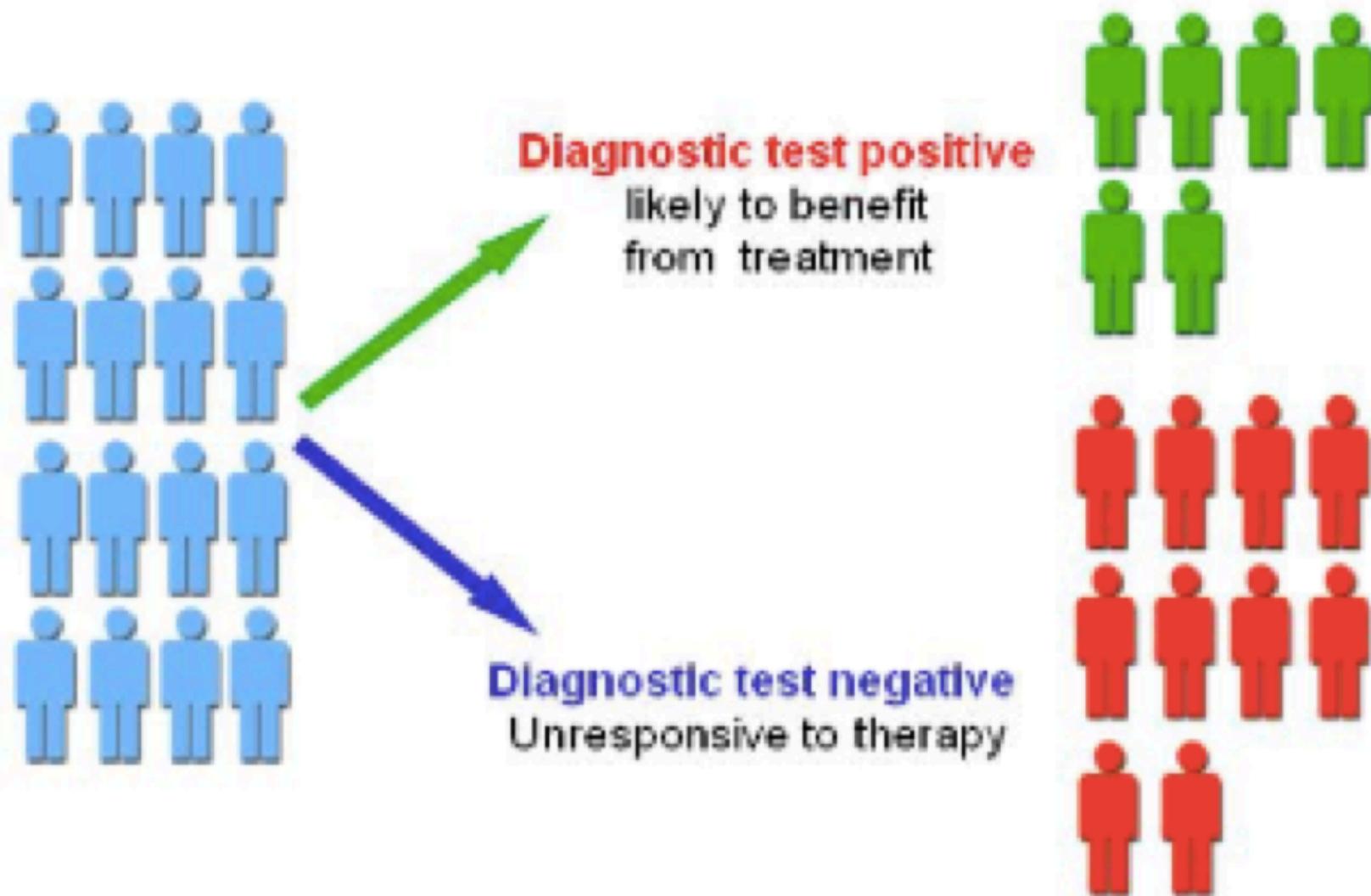


Most therapies developed through traditional approaches fail to have efficacy

Therapeutic development desperately needs into human biology

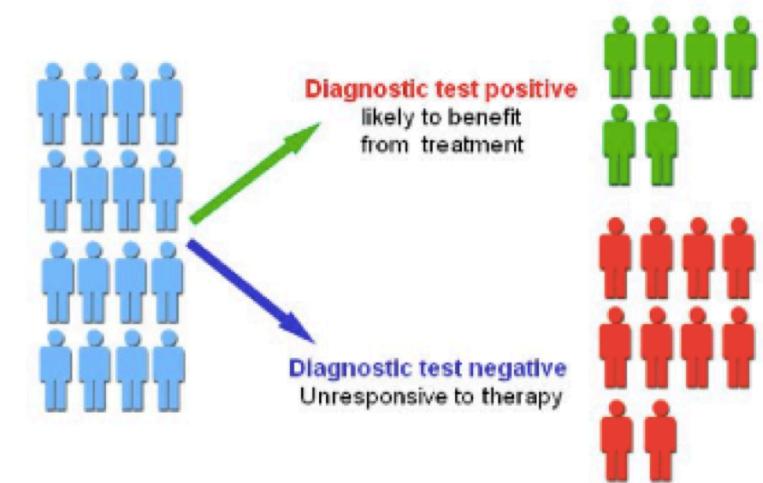
Genetics provides a validated means to these insights

Informing therapeutic development

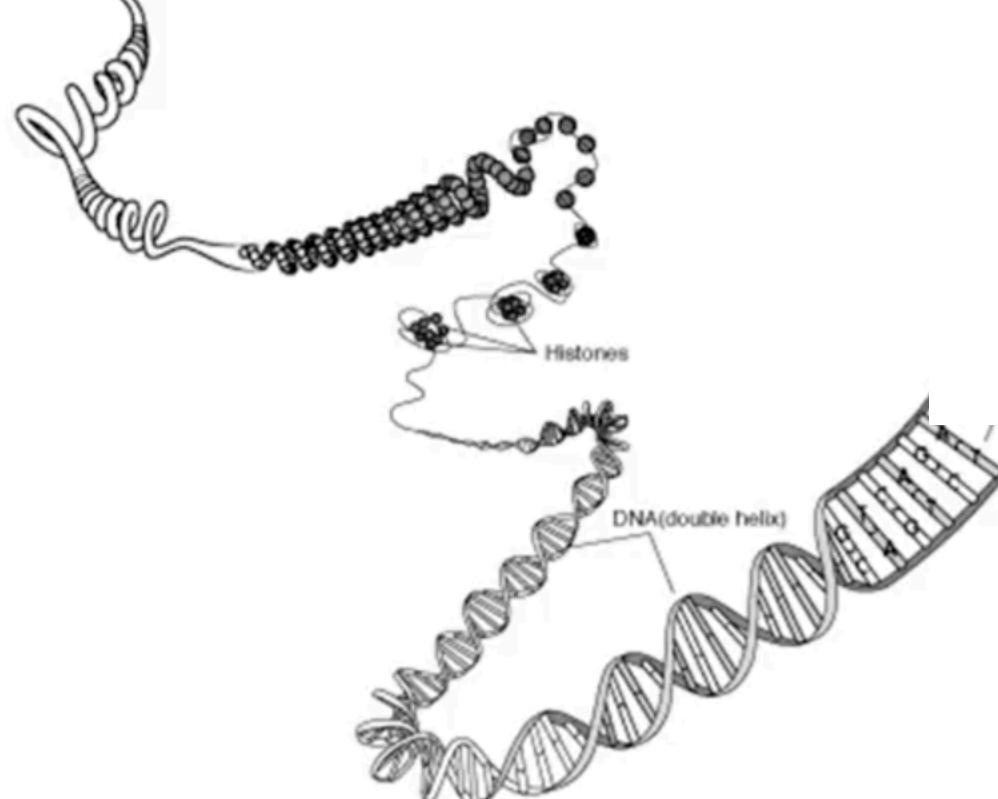
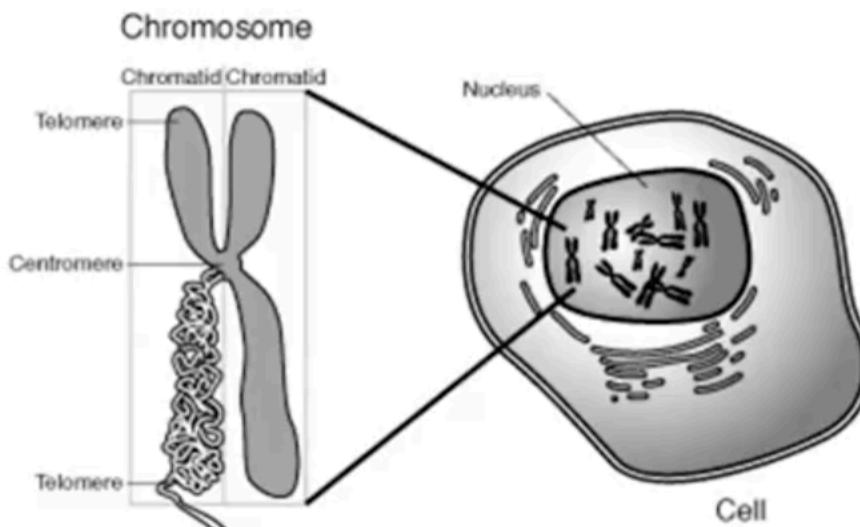


Genetics can in some cases provide individualized medical insights

- diagnostics in cases of severe genetic disorders
- identification of individuals more or less likely to benefit from specific therapeutic interventions
- prediction of individuals at risk for severe idiosyncratic adverse drug responses



Vastly oversold for years - technological advances are bringing what has long been a futuristic vision much closer to reality



The scope of the challenge

Within each cell:

2 copies of the genome

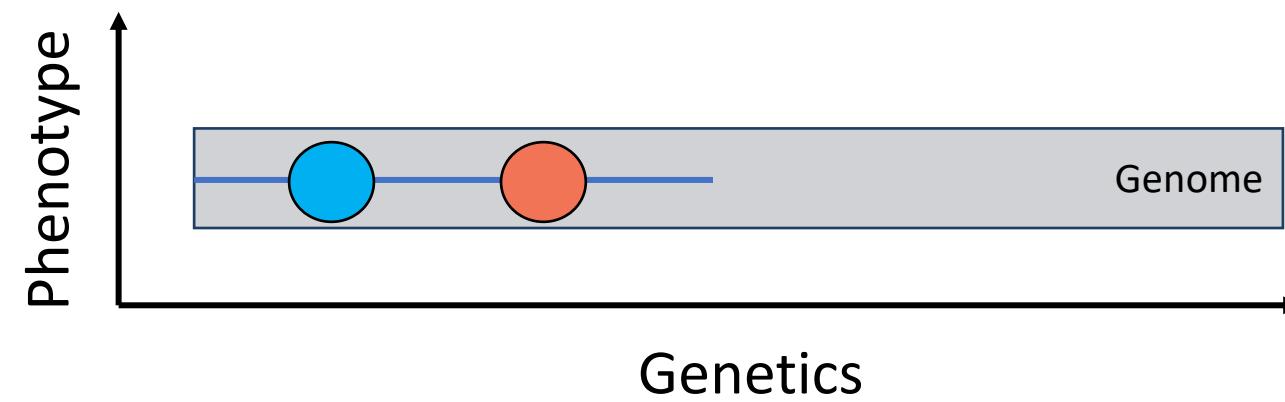
23 chromosomes

~20,000 genes

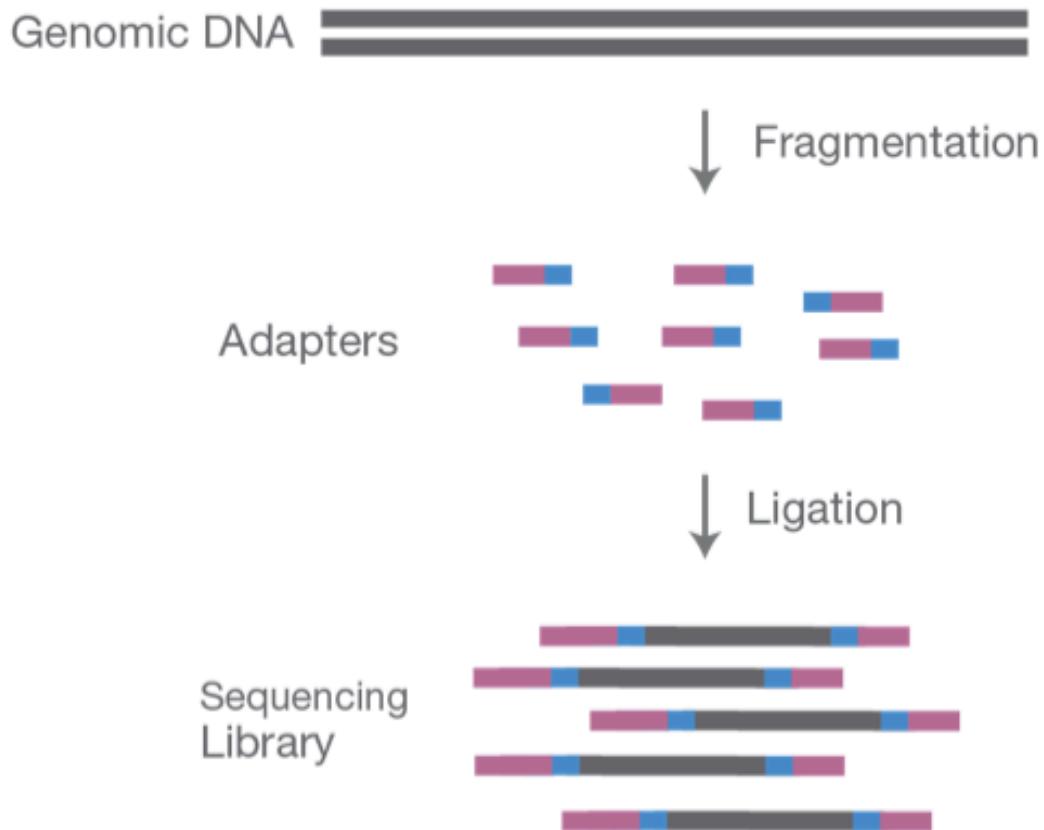
3.2B letters of DNA

Millions of polymorphic sites

Genomics

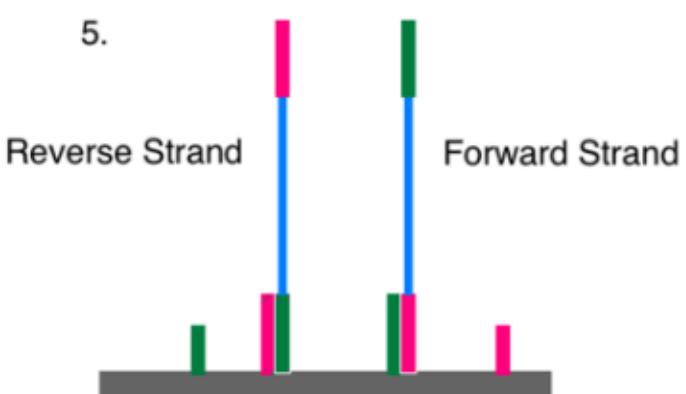
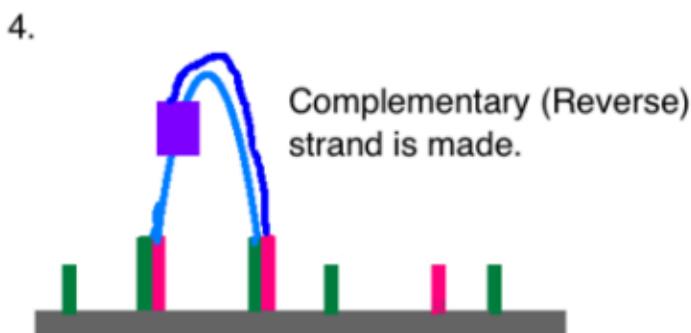
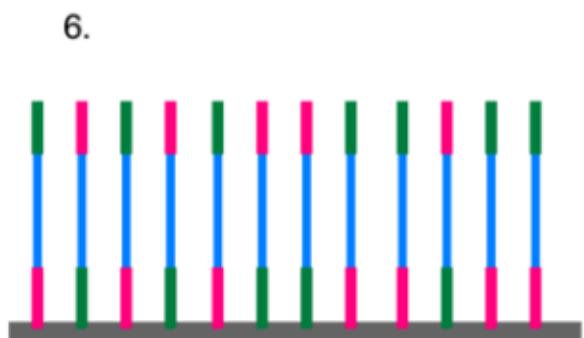
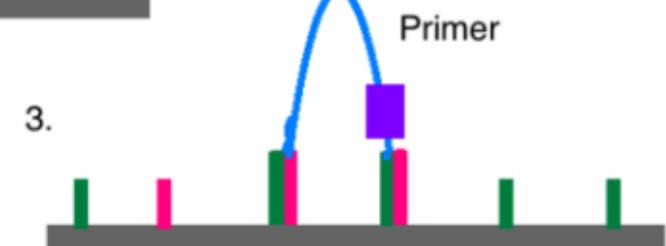
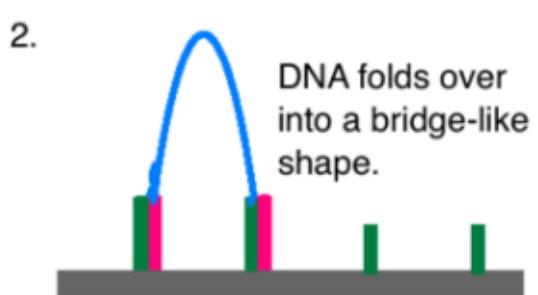
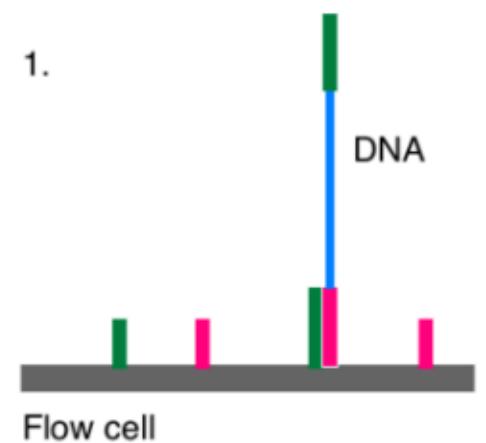


A. Library Preparation



NGS library is prepared by fragmenting a gDNA sample and ligating specialized adapters to both fragment ends.

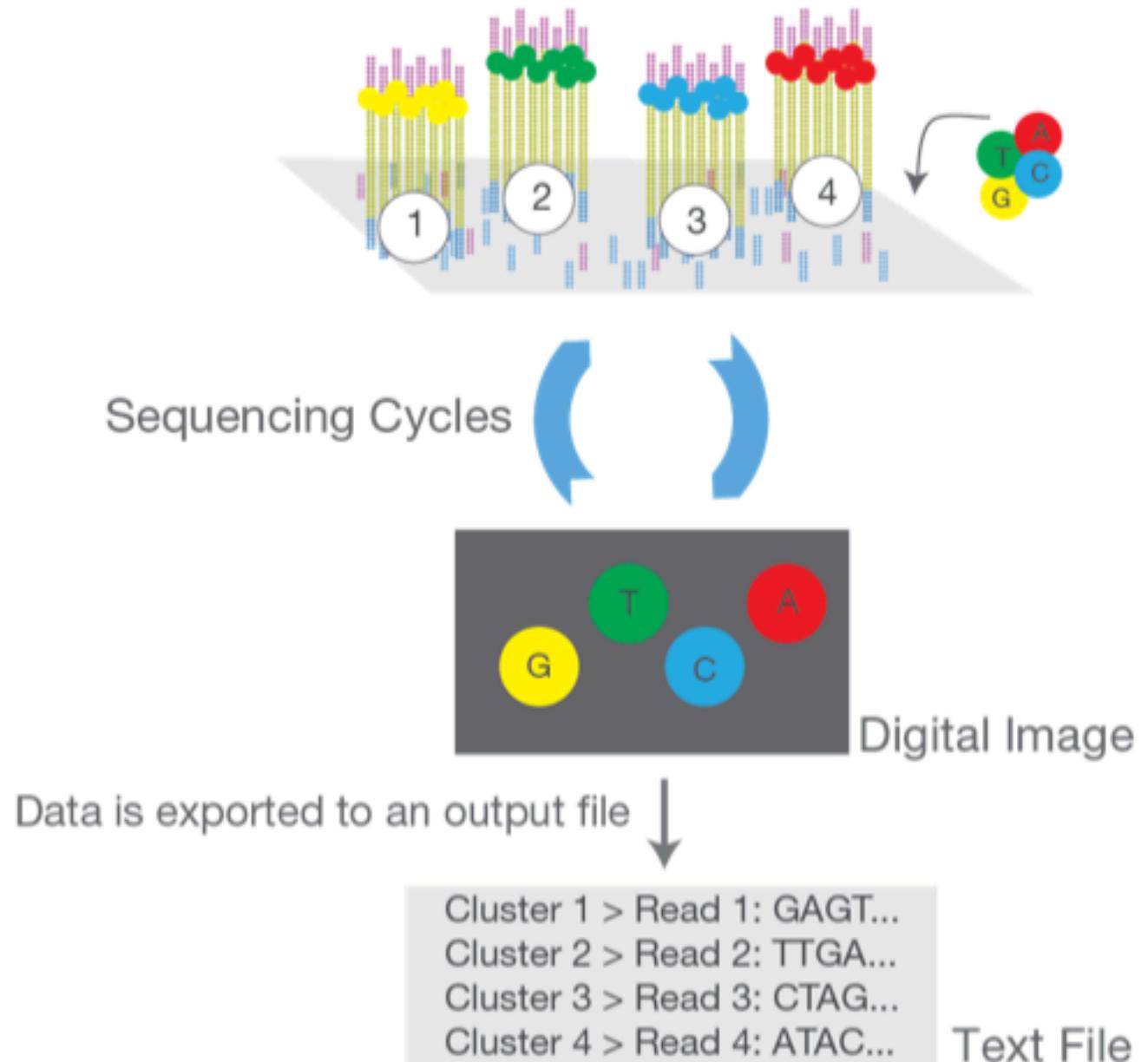
Illumina sequencing



B. Cluster Amplification

Illumina sequencing

C. Sequencing



Discovering genetic variation: Sequencing

Reference ACTGCACCTACGTGCCCTACGTACATACTCA

 CTGCACCTACGTGGC
 TGGCCCTACGTACATACTAC
 CCTACGTGGCCCT
 GTGGCCCTACGTC

Some mapping procedure ignores mismatches, and matched reads that were highly similar.

But they may represent true sequence variants

Statistical methods needed to distinguish true variants (Variant calling)

Whole genome variant callline: GATK HaplotypeCaller

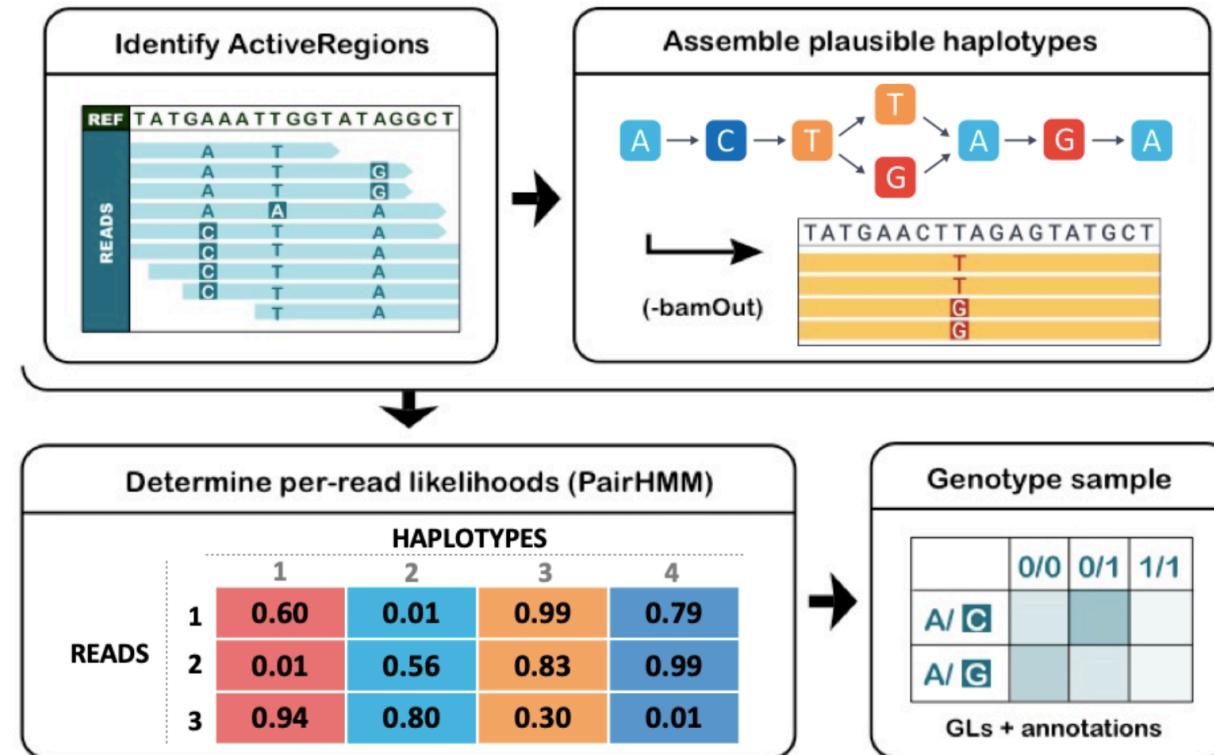
1. Use heuristic to find mismatches not explained by noise
2. Use assembly graph to identify possible haplotypes
3. For each haplotype, estimate

P(read | haplotype)

Using probabilistic sequence alignment

- Hidden Markov Model
 - States: insertion, deletion, substitution
 - Emission: pairs of aligned nucleotides/gaps
 - Transitions: equivalent to insertion/deletion/gap penalties from Smith-Waterman algorithm
 - Get **P(read | haplotype)** using forward-backward algorithm
4. Use Bayes rule to get $P(\text{Haplotype} | \text{read})$
 5. Assign genotypes to each sample based on the max a posteriori haplotypes

A haplotype is a group of genes or DNA sequences that are inherited together on the same chromosome



Type of genetic variation

- **99% of DNA is shared** between two individuals
- Variation in the remainder explains all our **predisposition** differences
- **Remaining** phenotypic variation: environmental/stochastic differences

Name	Example	Frequency in one genome
Single nucleotide polymorphisms (SNPs)	GAGGAGAACG [C/G] AACTCCGCCG	1 per 1,000 bp
Insertions/deletions (indels)	CACTATT C/[CTATGG] TGTCTAA	1 per 10,000 bp
Short tandem repeats (STRs)	ACGGCA GTCGTCGTCGTC ACCGTAT	1 per 10,000 bp
Structural variants (SVs) / Copy Number Variants (CNVs)	Large (median 5,000 bp) deletions, duplications, inversions	1 per 1,000,000 bp

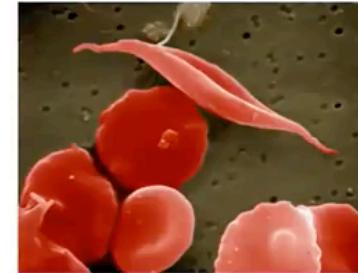
Type of genetic variation

CATGGTGCATCTGACTCCT**GAGGAGAAGTCTGCCGTTACTG**
CATGGTGCATCTGACTCCT**GT**GGAGAAGTCTGCCGTTACTG

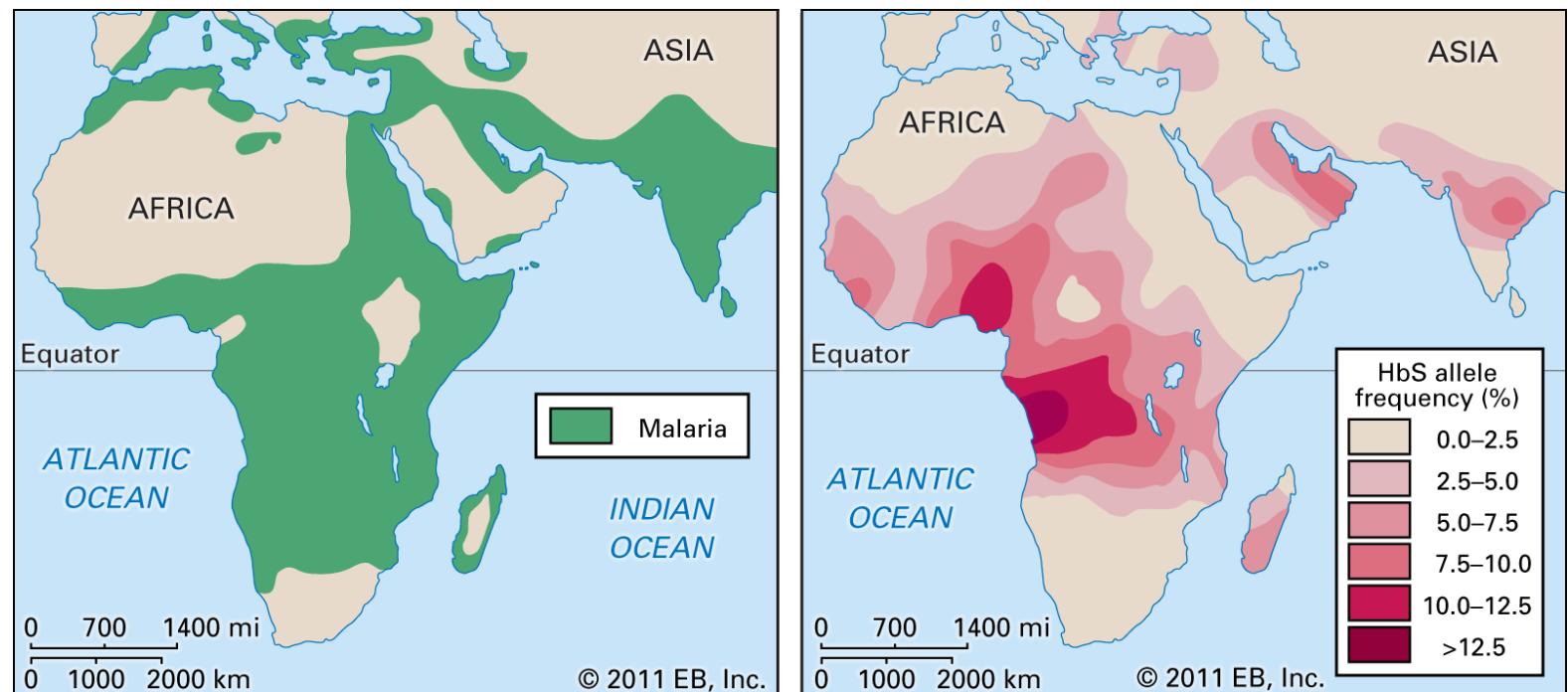
glutamic acid > valine

The disease is characterized by many of the symptoms of chronic anemia (fatigue, pale skin, and shortness of breath) as well as susceptibility to infection, jaundice and other eye problems, delayed growth, and episodic crises of severe pain in the abdomen, bones, or muscles.

Sickle Cell Anemia



The distribution of malaria and the distribution of sickle cell anemia overlap in areas of Africa, southern Asia, and the Mediterranean. The persistence of the *HbS* gene, which causes sickle cell anemia, has been explained by the fact that heterozygous persons are resistant to malaria.



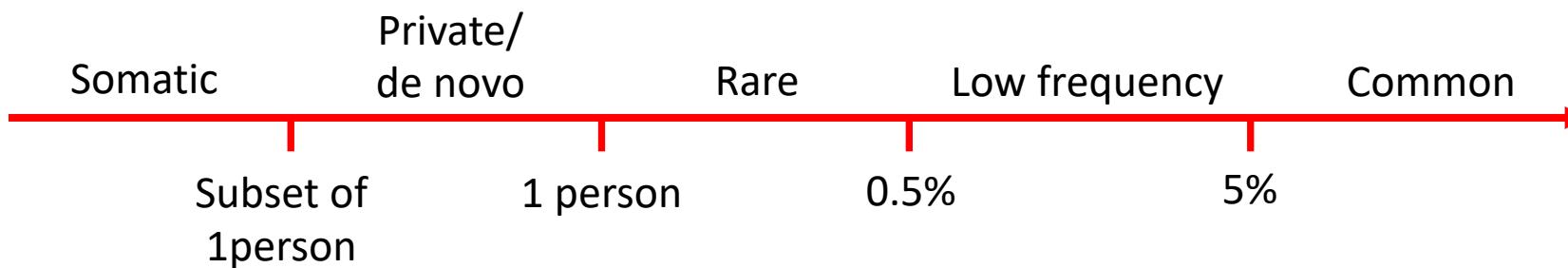
- Many modern analyses (GWAS, eQTL) focus on SNPs/indels
- Often have only two **alleles** (states)
- Identified as reference SNP clusters (**rsid**)
- Submitted sequences containing a variant are clustered to build a database (**dbSNP**)
- To date, >100 M known variants in dbSNP

rs189107123

GAGGAGAACG [C/G] AACTCCGCCG

Classifications of SNP alleles

Classifying variants by minor allele frequency:



Distinguishing the two alleles

- Matching the human reference sequence (Reference/Alternate)
- Being more frequent in the population (Major/Minor)
- Matching the most recent common ancestor between human and chimpanzee (ancestral/derived)
- Based on their disease association (risk/non-risk)

Example: rs189107123

GAGGAGAACG [C/G] AACTCCGCCG

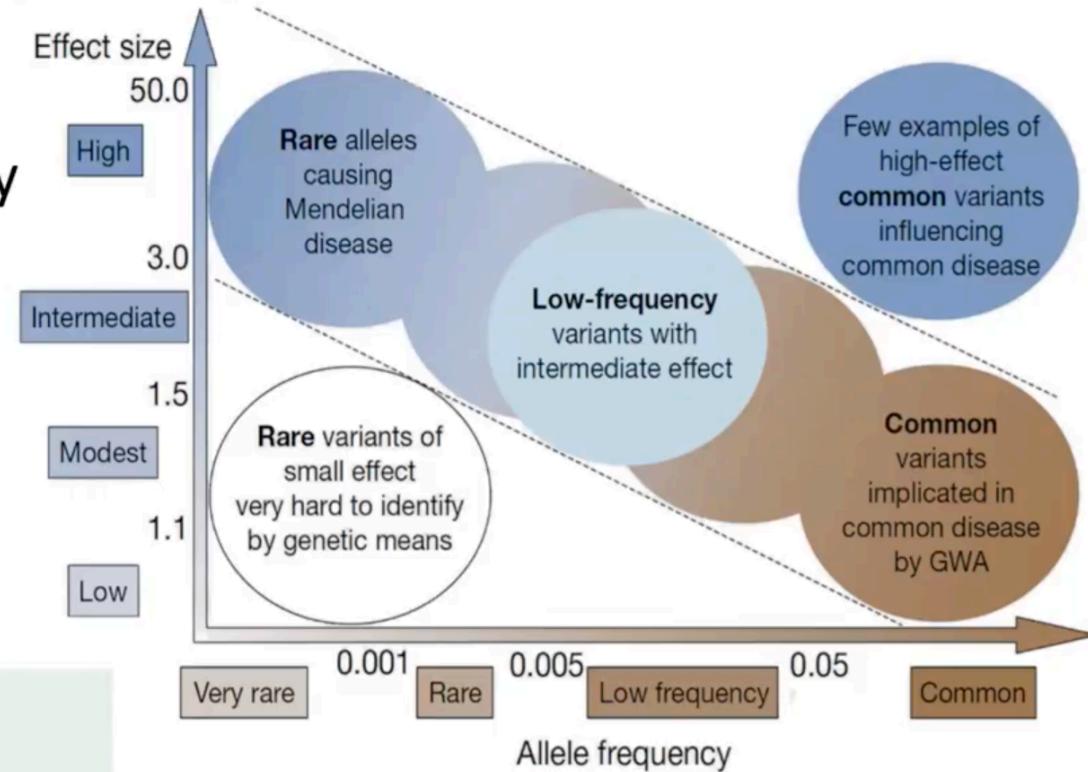
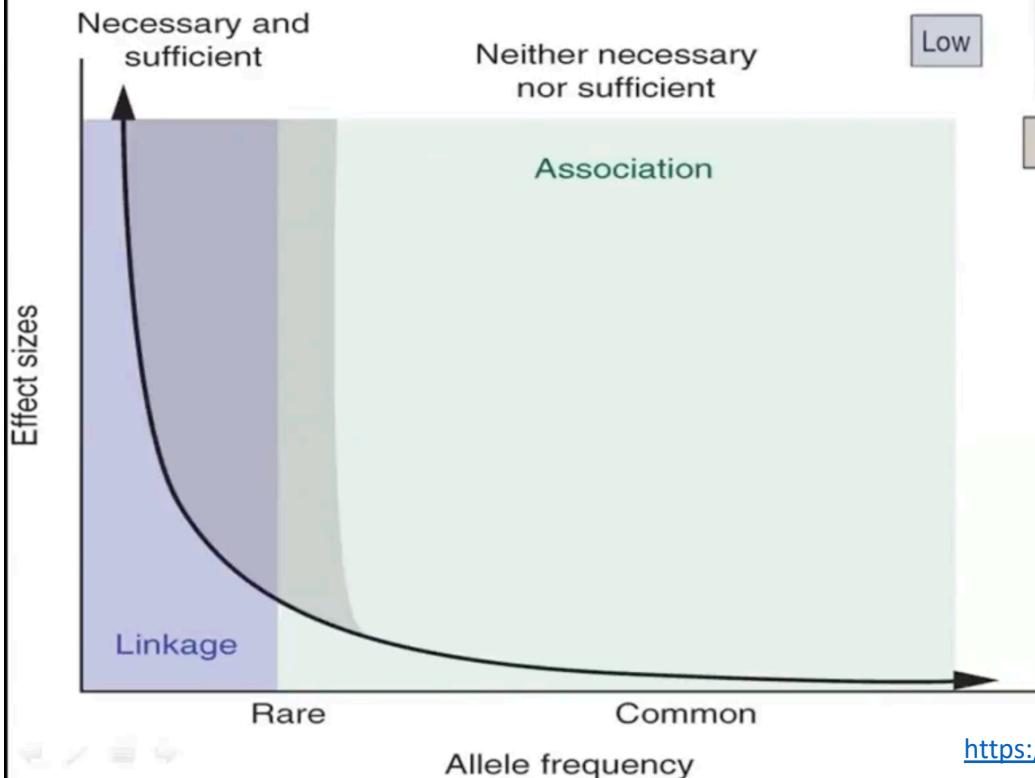
Reference allele: C

Minor allele: G (frequency 0.03 in Europeans)

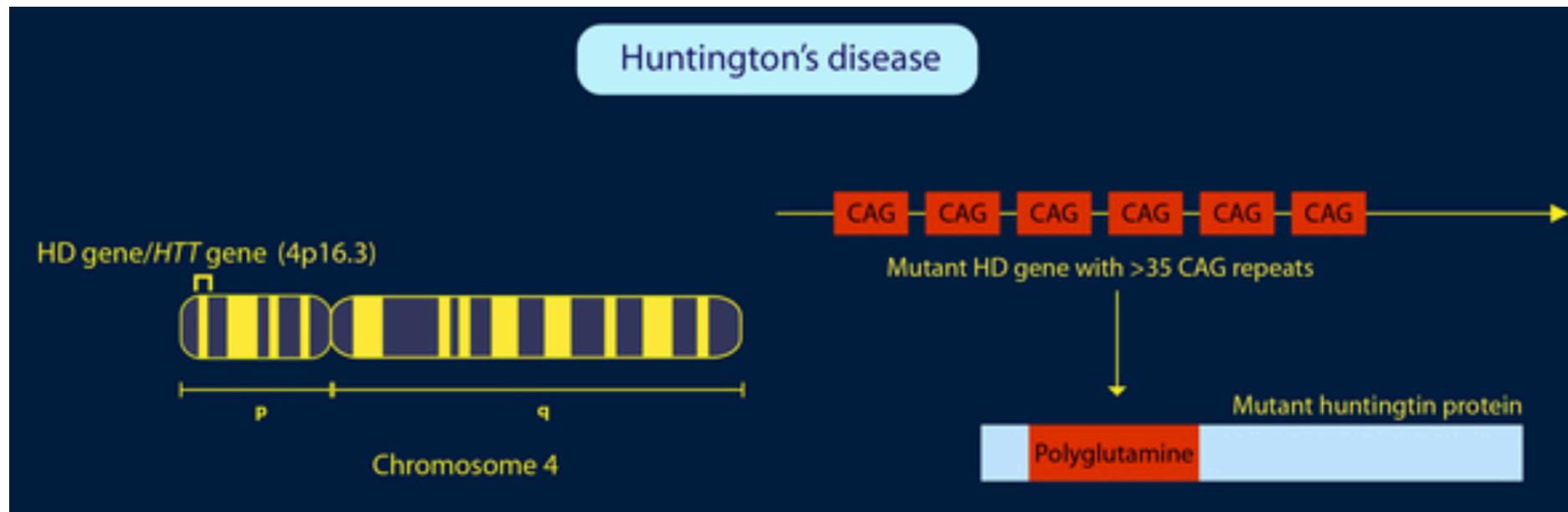
Ancestral allele: unknown (**why?**)

Common alleles typically have small effects

Discovery method tuned to variant effect size/frequency



Discovery method tuned to variant effect size/frequency



9

TCACAGCAGCAGCAGCAGCAGCAGCAGCAGCAGTTGCATTT

10

TCACAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGCAGTTGCATTT

12

TCACAGTTGCATTT

> 30 Huntington's Disease

Abnormal protein, damages neurons, brain cell death, mood, coordination, speaking, dementia, etc



Beyond SNPs: Insertion/Deletions

Cystic fibrosis transmembrane conductance regulator (CFTR) → Lung infections, Cystic fibrosis

CATTAAAGAAAATATCATCTTGGTGTTCCCTATGATGAAT.
CATTAAAGAAAATATCATTGGTGTTCCCTATGATGAATA

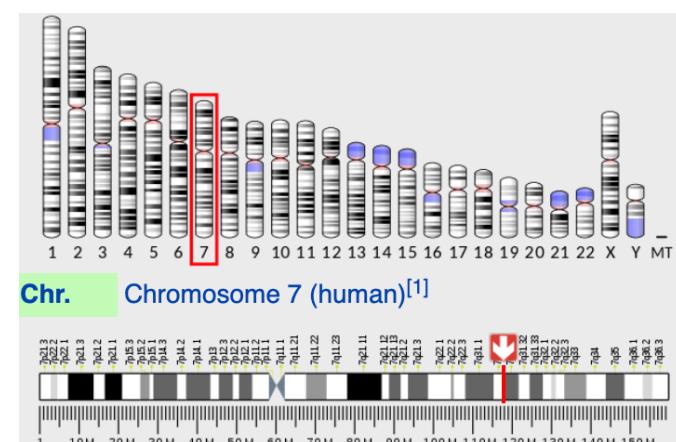
CFTR Sequence:

Nucleotide	ATC	ATC	C	T	T	GT	GT
Amino Acid	Ile	Ile	Phe		Gly	Val	
506	I	I	Phe I				510
Deleted in ΔF508							

ΔF508 CFTR Sequence:

Nucleotide	ATC	ATT	GGT	GT
Amino Acid	Ile	Ile	Gly	Val
506	I			

DeltaF508 ($\Delta F508$) results from a deletion (Δ) of three nucleotides which results in a loss of the amino acid phenylalanine (F) at the 508th position on the protein. As a result, the protein does not fold normally and is more quickly degraded. This protein does not escape the endoplasmic reticulum for further processing. Having two copies of this mutation (one inherited from each parent) is by far the most common cause of cystic fibrosis (CF), responsible for nearly two-thirds of cases worldwide



https://en.wikipedia.org/wiki/Cystic_fibrosis_transmembrane_conductance_regulator

“Unrelated” population of individuals measured for a phenotype

	SNPs										Phenotype	
Individual	A	C	G	A	A	C	G	G	T	A	A	1
	C	C	G	G	T	C	G	G	T	C	T	1
	C	C	T	A	T	G	A	A	A	A	A	0
	A	T	G	A	A	G	G	G	T	A	T	0

Population Association studies

- Find an association between a SNP and phenotype
- Really want to find a SNP that is causal.

Individual	SNPs										Phenotype
	A	C	G	A	A	C	G	G	T	A	
	A	C	G	G	T	C	G	G	T	C	1
	C	C	G	G	T	C	G	G	T	C	1
	C	C	T	A	T	G	A	A	A	A	0
	A	T	G	A	A	G	G	G	T	A	0

Testing for association

Individual	SNPs										Phenotype
	A	C	G	A	A	C	G	G	T	A	
	A	C	G	A	A	C	G	G	T	A	A
	C	C	G	G	T	C	G	G	T	C	T
	C	C	T	A	T	G	A	A	A	A	0
	A	T	G	A	A	G	G	G	T	A	0

R11209026	Allele A	Allele G
Case	22	976
Controls	68	932

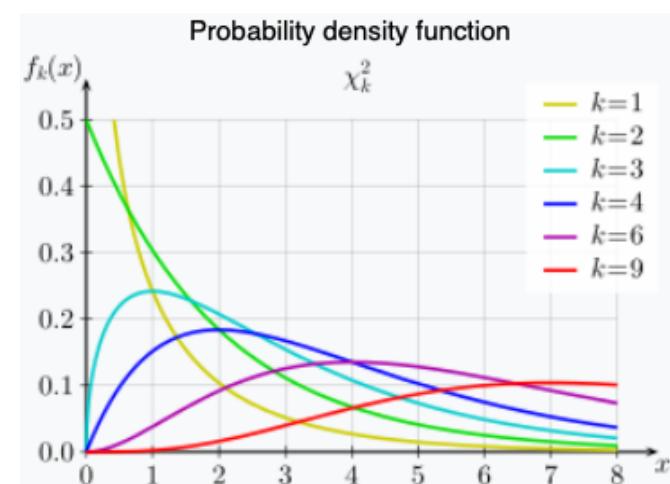
Expected

R11209026	Allele A	Allele G
Case	45	953
Controls	45	955

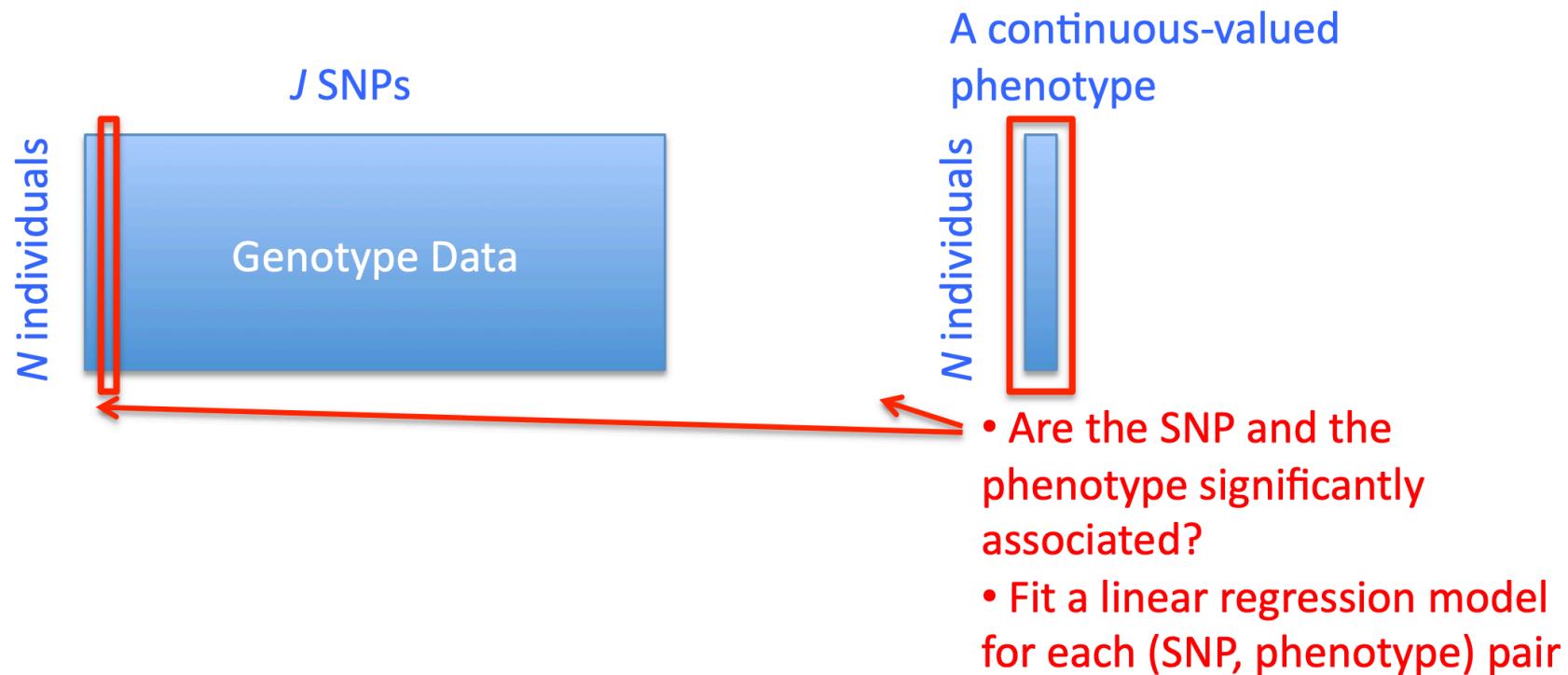
$$\chi^2 = \sum(O - E)^2/E$$

R11209026	Allele A	Allele G
Case	11.72	0.55
Controls	11.70	0.55

- Simplest tests (single marker regression, χ^2) rule the data
- Complex multi-marker models are often less reliable

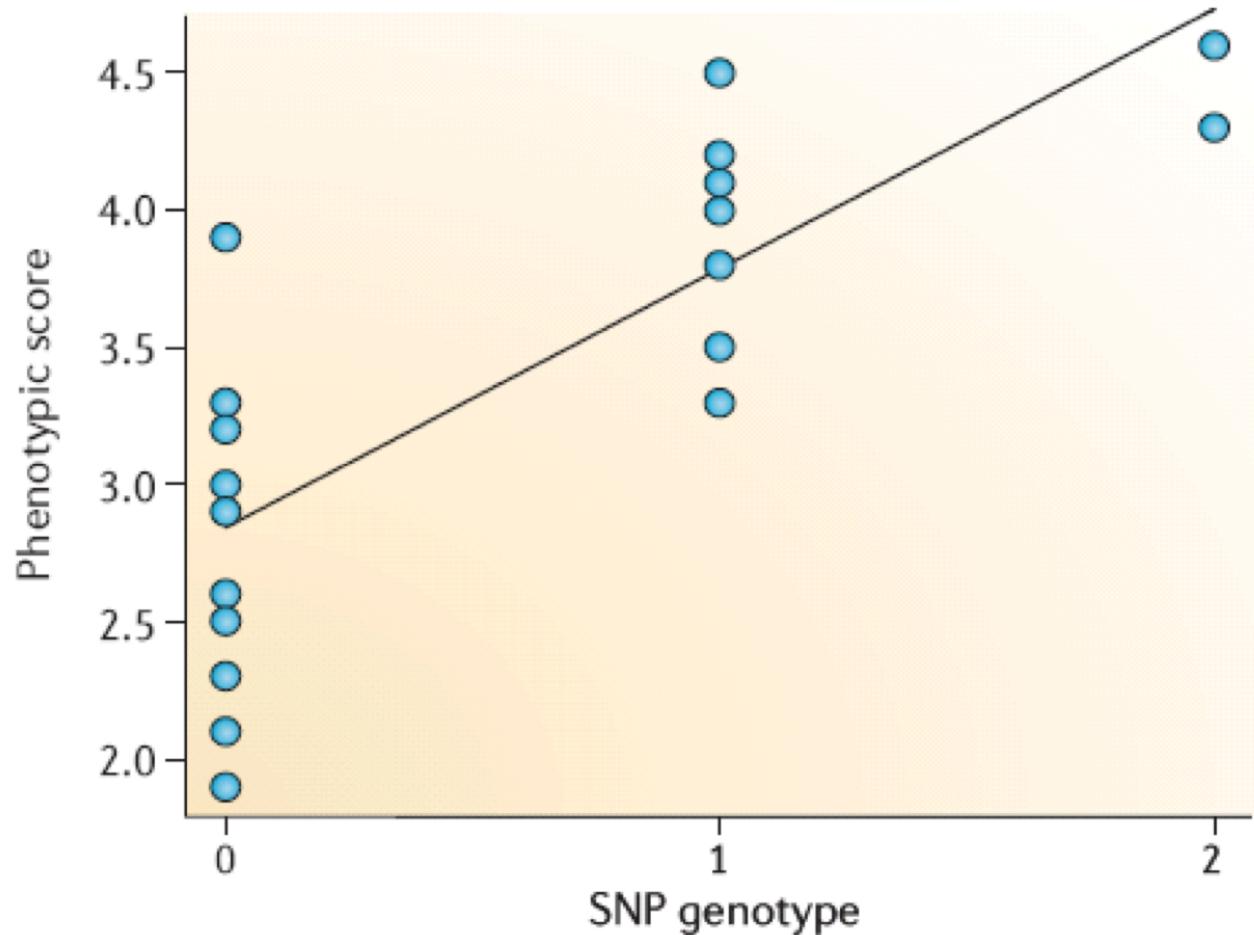


Using linear regression for SNP association testing



Using linear regression for SNP association testing

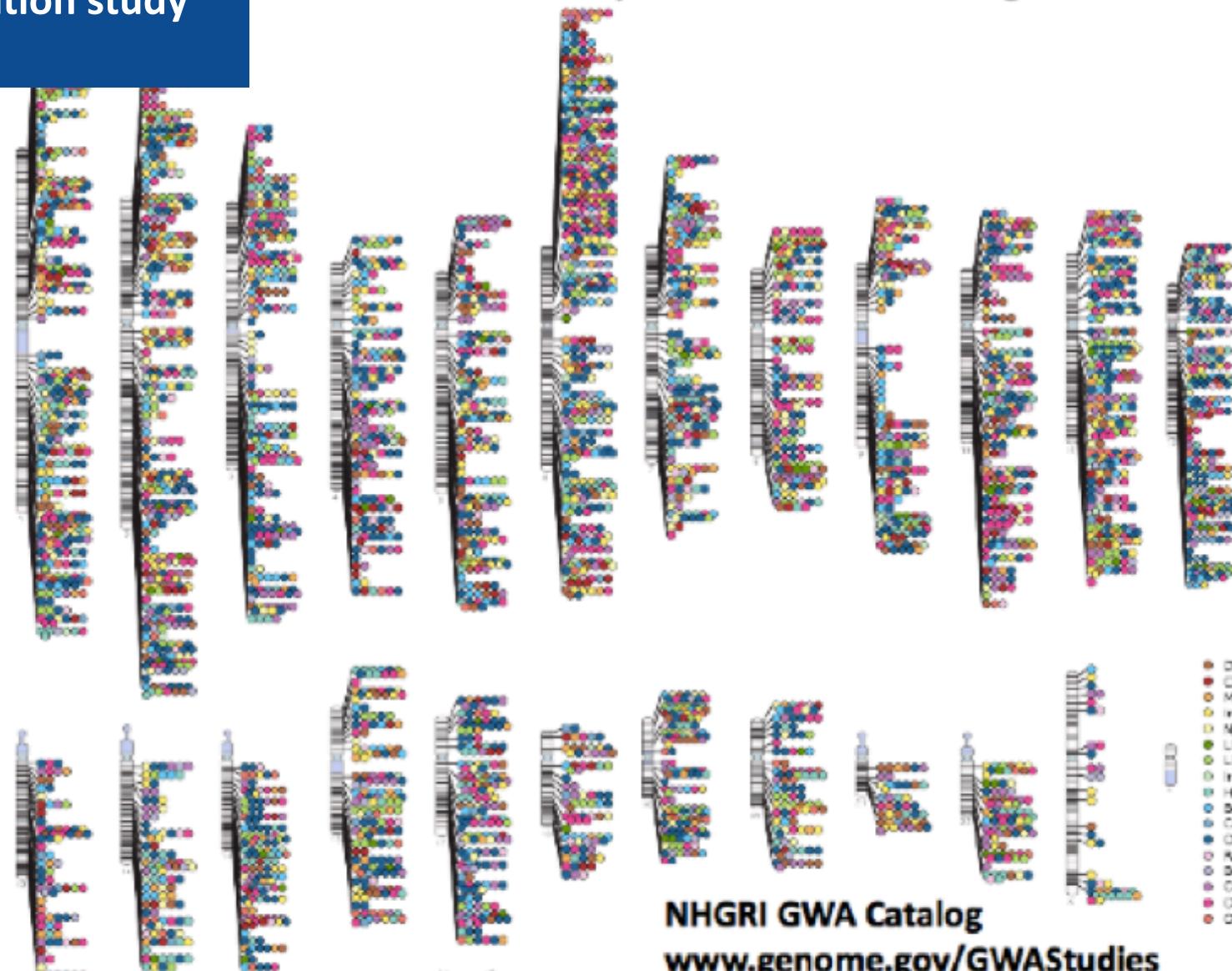
- Continuous-valued traits
 - Also called quantitative traits
 - Cholesterol level, blood pressure etc.
- One cannot create a contingency matrix as in case/control studies
- For each locus, fit a linear regression using the number of minor alleles at the given locus of the individual as covariate



Genome-wide association study (GWAS)

Perform a population association study across the genome

Published Genome-Wide Associations through 12/2013
Published GWA at $p \leq 5 \times 10^{-8}$ for 17 trait categories

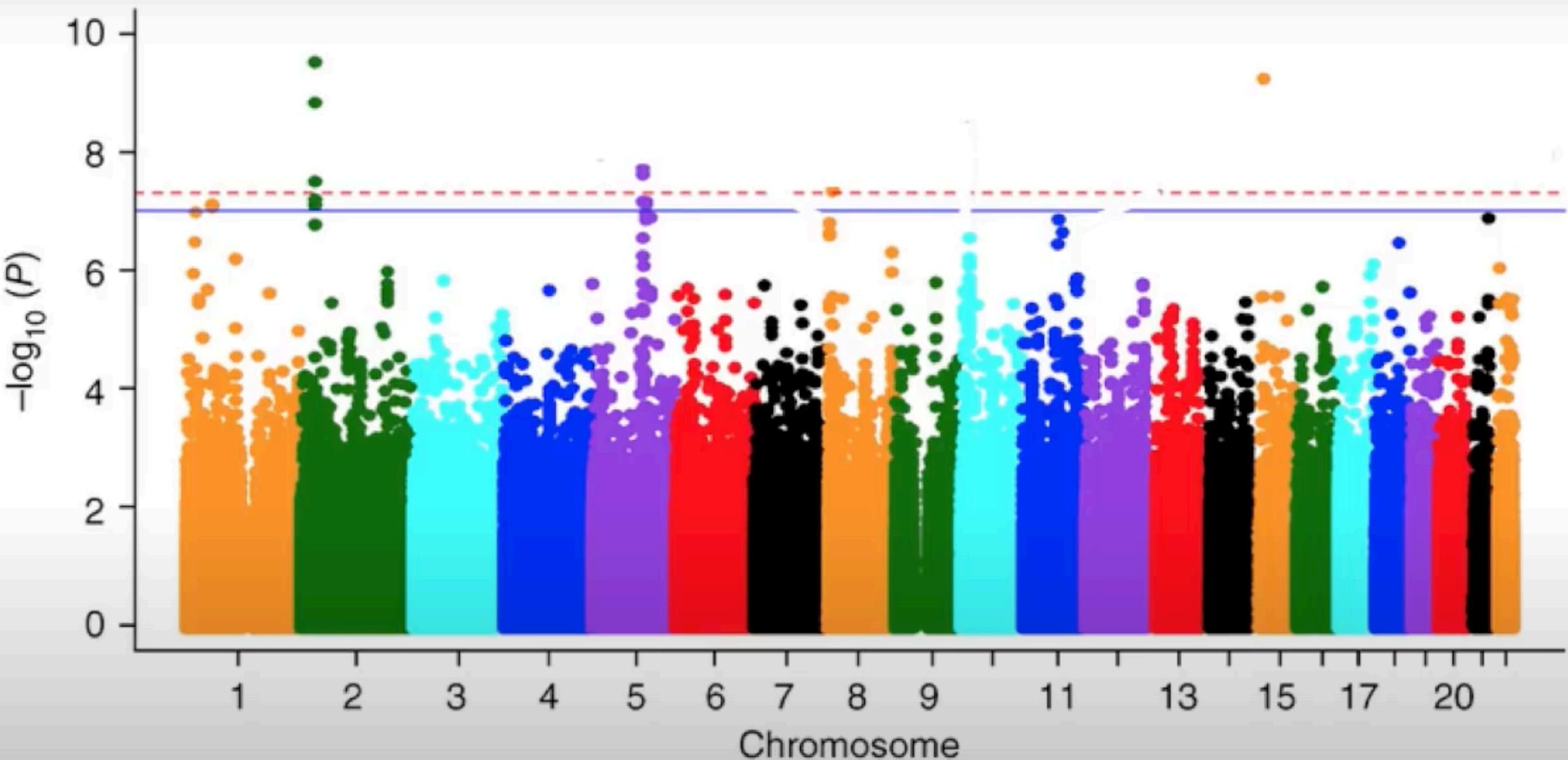


NHGRI GWA Catalog
www.genome.gov/GWASStudies
www.ebi.ac.uk/fgpt/gwas/

- In GWAS, we're performing $O(10^6)$ tests that are largely independent
 - Each study has hundreds of $p < .001$ purely by statistical chance (no real relationship to disease)
 - “Genome-wide significance” often set at $p = 5 \times 10^{-8}$ ($= .05 / 1 \text{ million tests}$)

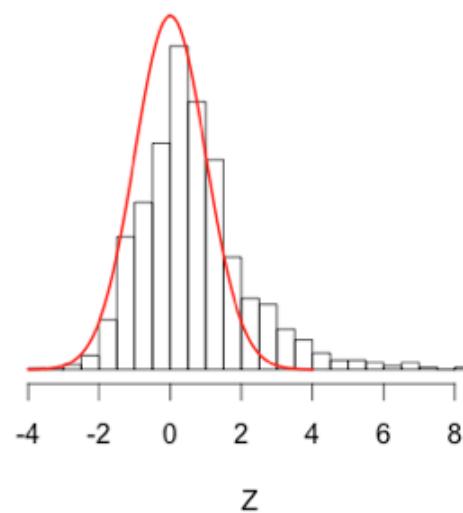
1996: Risch and Merikangas propose that a **p**-value of 5×10^{-8} (equivalent to a **p**-value of 0.05 after a Bonferroni correction for 1 million independent tests) is a conservative threshold for declaring significant association in a genome-wide study.

Understanding Manhattan Plots

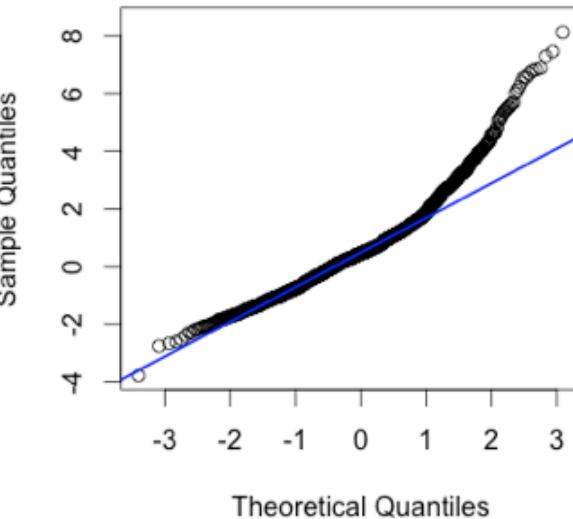


Understanding Q-Q Plots

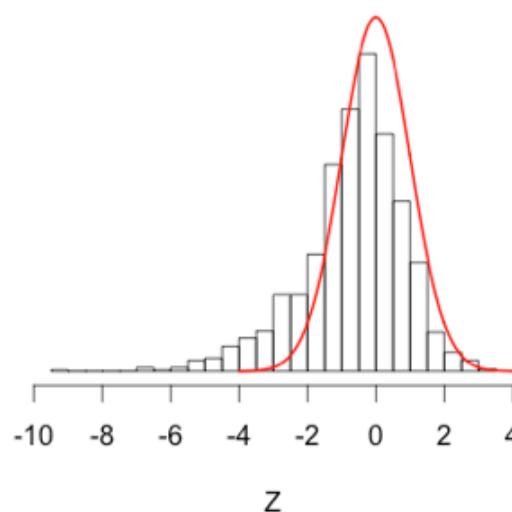
Skewed Right



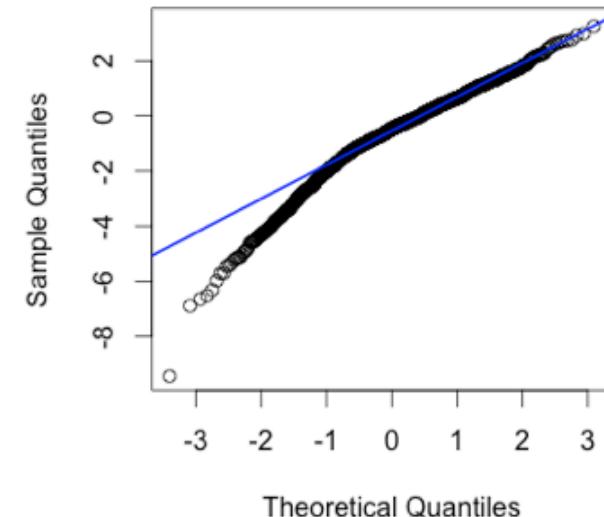
Normal Q-Q Plot



Skewed Left

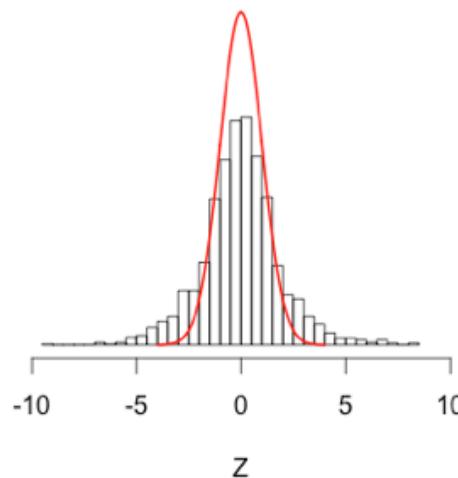


Normal Q-Q Plot

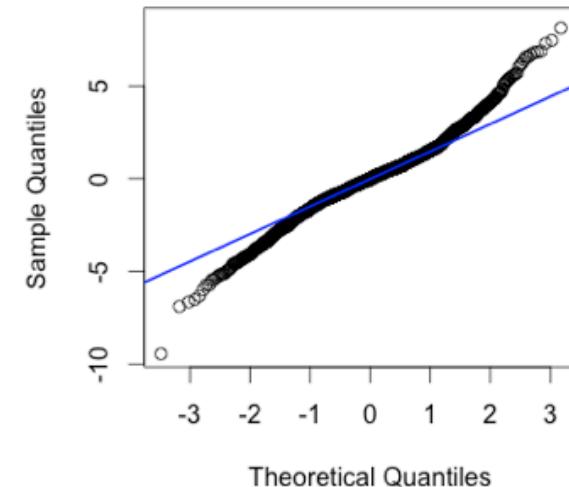


Understanding Q-Q Plots

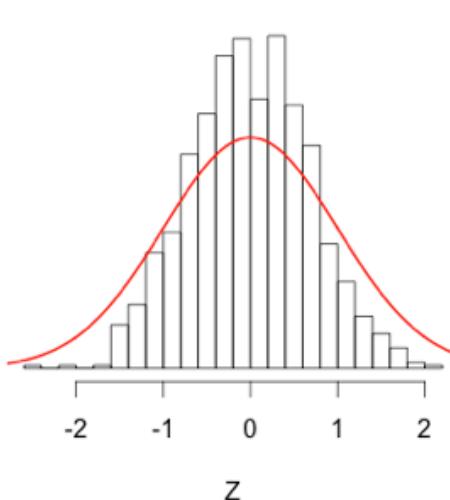
Fat Tails



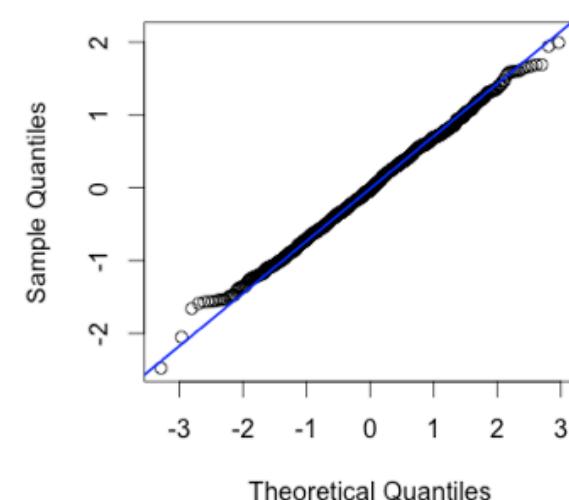
Normal Q-Q Plot



Thin Tails

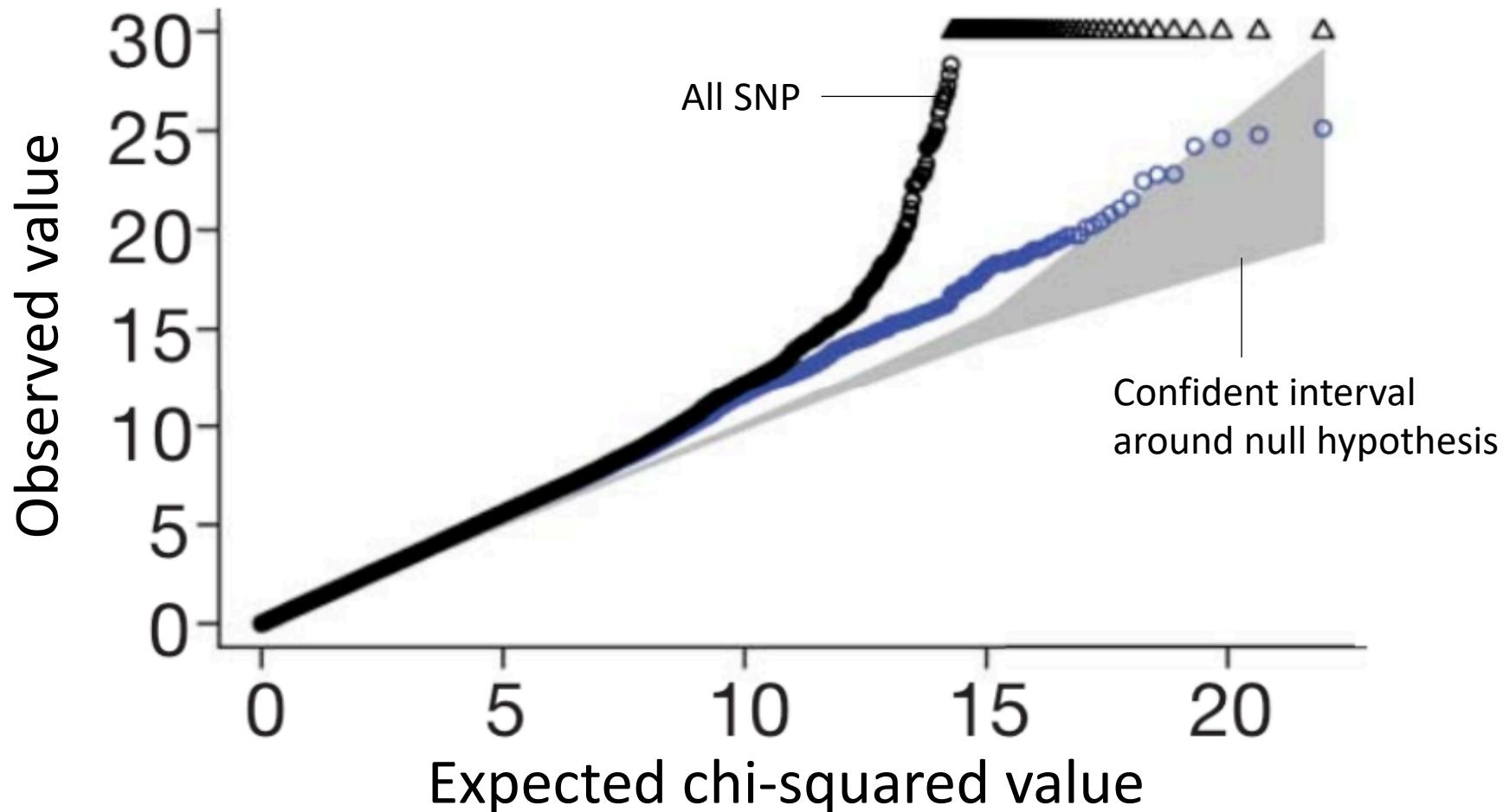


Normal Q-Q Plot



Understanding Q-Q Plots

A Q–Q plot is used to compare the shapes of distributions, providing a graphical view of how properties such as location, scale, and skewness are similar or different in the two distributions.



Genomics

Genomics is the most mature of the omics fields. In the realm of medical research, genomics focuses on identifying genetic variants associated with disease, response to treatment, or future patient prognosis.

GWAS is a successful approach that has been used to identify thousands of genetic variants associated with complex diseases (GWAS catalog <https://www.ebi.ac.uk/gwas/home>) in multiple human populations. In such studies, thousands of individuals are genotyped for more than a million genetic markers, and statistically significant differences in minor allele frequencies between cases and controls are considered evidence of association.

GWAS studies provide an invaluable contribution to our understanding of complex phenotypes. Associated technologies include genotype arrays, NGS for whole-genome sequencing, and exome sequencing

Rationale for the development of systems biology technologies that involve the integration of different omics data types to identify molecular patterns associated with disease

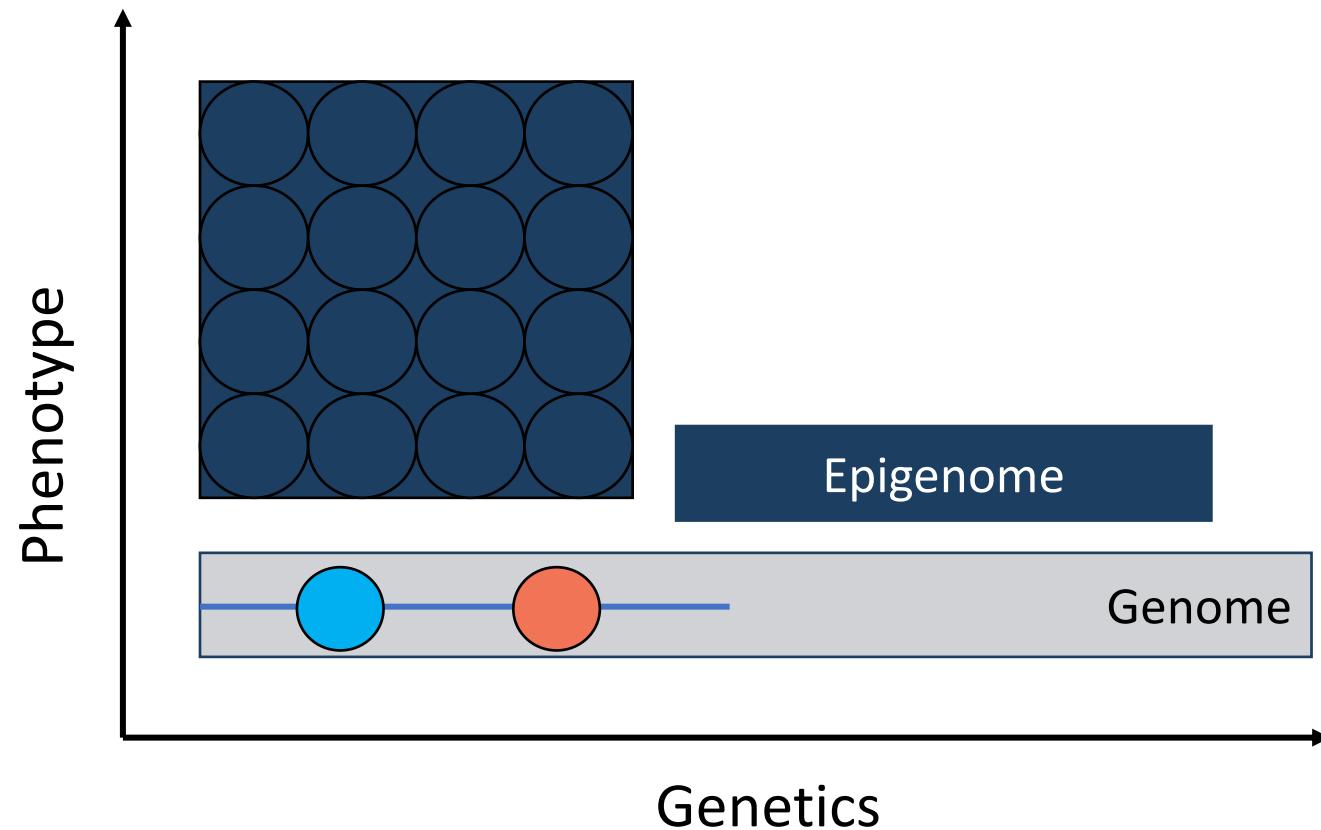
However, as our power to identify genetic variants associated with complex disease increased several realizations were reached that have shaped subsequent approaches to elucidating the causes of disease.

First, the loci that have been identified so far generally explain only a fraction of the heritable component for specific diseases.

Second, while Mendelian diseases generally result from changes in coding regions of genes, common diseases usually result from changes in gene regulation.

Third, the same genetic variants often contribute to different final outcomes, depending on the environment and genetic background.

Epigenomics



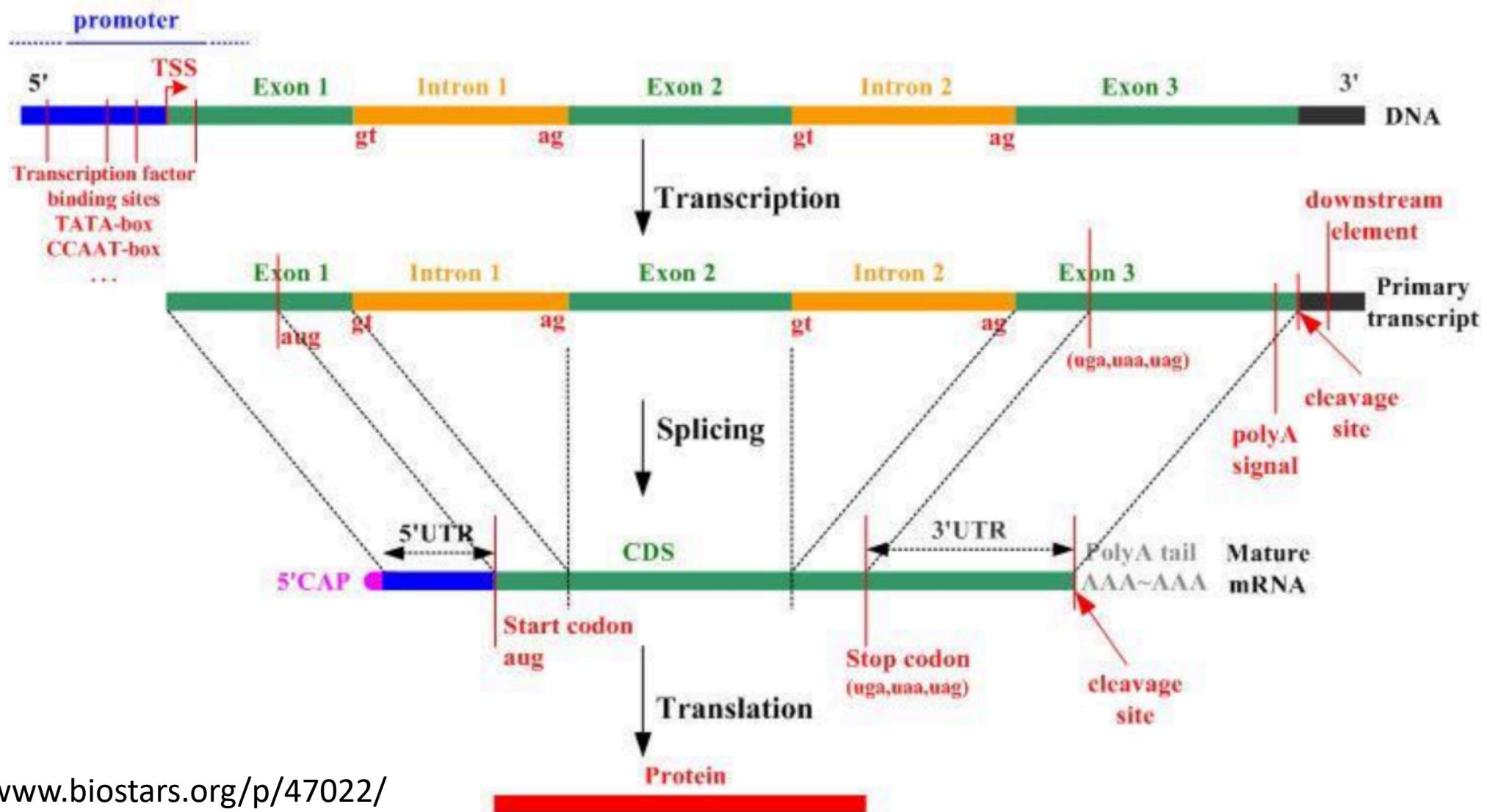
ORF and CDS

The region of the nucleotide sequences from the start codon (ATG) to the stop codon is called the **Open Reading frame**.

The **Coding Sequence (CDS)** is the actual region of DNA that is translated to form proteins. While the ORF may contain introns as well, the CDS refers to those nucleotides(concatenated exons) that can be divided into codons which are actually translated into amino acids by the ribosomal translation machinery.

Mainly: CDS means only that the sequence is known to be transcribed and, therefore, it is coding for something -- neither gene nor protein has to be known. Any full mRNA sequence (obtained from cDNA sequencing) will have a full coding sequence. ORF is usually predicted based on DNA sequence and not proven to be transcribed.

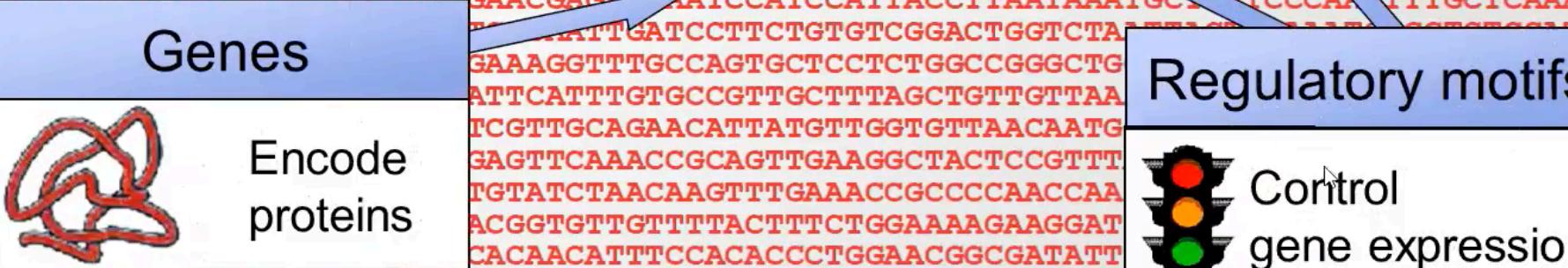
ORF and CDS



TATTGAATTTCAAAAATTCTTACTTTGGATGGACGCAAAGAAGTTAACATATTACATGGCATTACCACCATATA
ATCCATATCTAATCTTACTTATATGTTGGAATGTAAGAGGCCATTATCTTAGCCTAAAAAACCTCTTGGAACTTTC
AATACGCTTAAGTGCCTATTGCTATATTGAAGTACGGATTAGAAGCCGCCAGCGGGCGACAGCCCTCGACGGAAGACTCTCCTC
GCGTCCTCGTCTCACCGTCCGTTCTGAAACGCAGATGTGCCTCGGCCACTGCTCGAACAAATAAGATTCTACAATACI
TTTATGGTTATGAAGAGGAAAATTGGCAGTAACCTGGCCCCACAAACCTCAAAATTAAACGAATCAAATTAAACACCATAAGGATC
ATGCGATTAGTTTTAGCCTTATTCTGGGTAATTAAATCAGCGAAGCGATGATTGATCTATTAAACAGATATATAATGGAA
CTGCATAACCACCTTAACAACTTCAACATTTCAGTTGTATTACTTCTTATTCAAATGTCAAAAGTATCAACAAAAAI
TAATATACCTCTATACTTAACGTCAAGGAGAAAAACTATAATGACTAAATCTCATTAGAAGAAGTGATTGTACCTGAGTTCA
TAGCGCAAAGGAATTACCAAGACCATTGGCCAAAAGTGCCCGAGCATAATTAAAGAAATTATAAGCGTTATGCTAAACC
TTGTTGCTAGATGCCCTGGTAGAGTCATCTAATTGGTGAACATATTGATTATTGACTTCTCGGTTTACCTTAGCTATTGAT
GATATGCTTGCGCCGTCAAAGTTTGAACGAGAAAAATCCATTACCTTAATAAATGCTGATCCCATTGCTCAAAGGA
CGATTGCCGTTGGACGGTCTTATGTCACAATTGATCCTCTGTGCGACTGGCTAATTACTTAAATGTGGTCTCATGTT
ACTCTTTCTAAAGAAACTTGCACCGGAAAGGTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTGTGAGGGTGTACCA
GGCAGTGGATTGTCTTCCGCCGCATTCAATTGCGTTGCTTAGCTGTTAAAGCGAATATGGCCCTGGTATCATAT
CAAGCAAATTAAATGCGTATTACGGTCGTTGCAGAACATTATGTTGGTGTAAACAATGGCGGTATGGATCAGGCTGCCTGTT
GTGAGGAAGATCATGCTCTATCGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTAAATTCCGCAATTAAAAACCATGAA
AGCTTGTATTGCAACACCCCTGTTGATCTAACAGTTGAAACCGCCCAACCAACTATAATTAAAGAGTGGTAGAAGTCAC
AGCTGCAAATGTTAGCTGCCACGTACGGTGTGTTACTTCTGGAAAAGAAGGATCGAGCACGAATAAGGTATCTAAGAG
TCATGAACGTTATTATGCCAGATATCACAAACATTCCACACCCCTGGAACGGCGATATTGAATCCGGCATCGAACGGTAA
CTAGTACTAGTTGAAGAGTCTCGCCAATAAGAAACAGGGCTTACTGTTGACGATGTCGACAATCCTGAATTGTTCTCGCGA
ATTCAAGAGAGACTACTAACACATCTCCAGTGAGATTCAAGTCTTAAAGCTATATCAGAGGGCTAACGATGTGATTCTGAAT
TAAGAGTCTTGAAGGCTGTGAAATTAAATGACTACAGCGAGCTTACTGCCAGCGAACACTTTCAAGCAATTGGTGCCTGATC
GAGTCTCAAGCTTCTGCGATAAAACTTACGAATGTTCTTGCCAGAGATTGACAAAATTGTTCCATTGCTTGTCAAATGGATC
TGGTCCCGTTGACCGGAGCTGGCTGGGTGTTGACTGTTCACTGGTCCAGGGGCCAAATGGCAACATAGAAAAGGTA
AAGCCCTTGCCTAACGAGTTCTACAAGGTCAAGTACCCCTAACGATCACTGATGCTGAGCTAGAAAATGCTATCGTCTAA
TTGGGCAGCTGTCTATATGAATTATAAGTATACTTCTTTTACTTGTTCAGAACAACTCTCATTTTCTACTCATAACT
GCATCACAAAATACGCAATAAACGAGTAGTAACACTTTATAGTCATACATGCTTCAACTACTAACAAATGATTGTATGATA
TTTCAATGTAAGAGATTGATTATCCACAAACTTAAACACAGGGACAAAATTCTGATATGCTTCAACCGCTGCGTTGG
CCTATTCTGACATGATGACTACCATTTGTTATTGACGTGGGGCAGTGACGTCTTACATGTCAAAGTCATTGCGAAC
TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAGAGATTGCCGTTGAAACTTTGTCCTTTTCTTCCGGGACTCTAC
AACCTTTGTCCTACTGATTAATTGACTGAATTGGACAATTCAAGATTAGACAAAGCGCGAGGAGGAAAAGAAATGACA
AAATTCCGATGGACAAGAAGATAGGAAAAAAAGCTTCAACGATTCTAGACCGAAAAAGTCGTATGACATCAGAACATGA
ATTTCAAGTTAGACAAGGACAAAATCAGGACAAATTGTAAGATATAATAAAACTATTGATTGACGCCATTGCCCTTCC
TCCATTAAATCTCTGTTCTCTTACTTATATGATGATTAGGTATCTGATAAAACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGTTCGAAGACATTCTACGCATAATAAGAACAGGAGGAAATAATGCCAGACAACTATCATTACATT
GCGGCTCTCAAAAGATTGAACCTCGCCAACCTATGGAATTCTCAATGAGACCTTGCACCAATAATGTGGATTGGAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTATGAGGCATCGAGCTTGAAGAAAAAGTAAGCTCAGAAAAACCTCAATA
CTCATTCTGGAAGAAAATCTATTGAAATATGTTGCGTTGACAAATCAATCTTGGGTGTTCTATTCTGGATTCAATTGAC
AGGACTTGAAGCCCCGTCGAAAAAGAAAAGGCGGGTTGGTCTGGTACAATTATTGTTACTTCTGGCTGCTGAATGTTCAATATC
ACTTGGCAAATTGCGACTACAGGTCTACAACCTGGTCTAAATTGGTGGCAGTGTGATAACAATTGGATTGGTACGGTTCGT

TATTGAATTTCAAAAATTCTACTTTTTGGATGGACGCAAAGAAGTAAATAATCATATTACATGGCATTACCACCATATA
ATCCATATCTAATCTTAC**TTATA**TGTTGTGGAATGTAAGAGGCCATTATCTTAGCCTAAAAAACCTCTTGGAACTTT
AATACGCTTAAGTGCCTATTGCTATATTGAAGTA**CGG**ATTAGAAGCCG**CCGAG****CGGGCGACAGCCCT****CGCACGG**AAAGACTCTCCTC
GCGTCCTCGTCTCACCGTCGCTTGAACGCAGATGTGCCT**CGC**GCCGACTGCT**CCG**AACAATAAAGATTCTACAATAC
TTTATGGTTATGAAGAGGAAAATTGGCAGTAACCTGG**CCCCA**CAAACCTCAAAATTACAATCAAATTAAACACCAGGATC
ATGCGATTAGTTTTAGCCTATTTC**TGGGG**TAATTAAATCAGCGAAGCGATGATTGATCTATTAAACAGATA**TATAA**ATGGAA
CTGCATAACCACTTAACTAATACCAACATTTCAGTTGTATTACTCTTATTCAAATGTCATAAAAGTATCAACAAAAA
TAATATAACCTCTATACTTAACGTCAAGGAGAAAAACTATA**ATGACTAAATCTCATT**CAGAAGAAGTGATTGTACCTGAGTTCA
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAAGAAATTATAAGCGCTTATGATGCTAAACC
GG
TTGTTGCTAGATCGCTGGTAGAGTCAACTCTAATTGGTGAACATATTGATTATTGTGACTTCTCGGTTTACCTTAGCTATTGAT
GATATGCTTGCGCCGTCAAAGTTGAACGAGAAAAATCCATTACCTTAATAAATGCTGATCCAAATTGCTCAAAGGA
CGATTGCCGTTGGACGGTTCTTATGTCACAATTGATCCTCTGTGTCGGACTGGTCTAATTACTTAAATGTGGTCTCCATGTT
ACTCTTTCTAAAGAAACTTGCACCGGAAAGGTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTCTGTGAGGGTGTAC
GGCAGTGGATTGTCTCTCGGCCGCATTCAATTGTGCCGTTGCTTAGCTGTTAAAGCGAATATGGGCCCTGGTTATCATA
CAAGCAAATTAAATGCGTATTACGGCGTTGCAGAACATTATGTTGGTGTAAACAATGGCGGTATGGATCAGGCTGCCTGTT
GTGAGGAAGATCATGCTCTATACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTAAATTCCGCAATTAAAAACCATGA
AGCTTGTATTGCGAACACCCCTGTTGATCTAACAGTTGAAACCGCCCAACCAACTATAATTAAAGAGTGGTAGAAGTC
AGCTGCAAATGTTAGCTGCCACGTACGGTTGTTACTTCTGGAAAAGAAGGATCGAGCAGAACAGGTAATCTAAGAC
TCATGAACGTTATTATGCCAGATATCACAAACATTCCACACCCCTGGAACGGCGATATTGAATCCGGCATCGAACGGTTAAC
CTAGTACTAGTTGAAGAGTCTCGCCAATAAGAAACAGGGCTTACTGTTGACGATGTCGACAACTCTGAATTGTTCTCGCG
ATTCAACAGAGACTACTAACACATCTCCAGTGAGATTCAAGTCTTAAAGCTATATCAGAGGGCTAACGATGTGATTCTGA
TAAGAGTCTTGAAGGCTGTGAAATTAAATGACTACAGCGAGCTTACTGCCAGAACAGACTTTCAAGCAATTGGTGCCTGAT
GAGTCTCAAGCTTGCATAAAACTTACGAATGTTCTTCCAGAGATTGACAAAATTGTTCCATTGCTTGTCAAATGGATC
TGGTCCCCTTGAACGGAGCTGGCTGGGTGGTTGACTGTTCACTGGTCCAGGGGGCCAAATGGCAACATAGAAAAGGTA
AAGCCCTTGCCTAACAGGTTCTACAAGGTCAAGTACCCCTAACGATCACTGATGCTGAGCTAGAAAATGCTATCATGCTCTAA
TTGGGCAGCTGTCTATATGAATT**TATAA**GTATACTCTTTTACTTGTTCAGAACAACTCTCATTTTTCTACTCATAACT
GCATCACAAATACGCAATAATAACGAGTAGTAACACTTTATAGTTCATACATGCTCAACTACTAACAAATGATTGTATGATA
TTTCATGTAAGAGATTCGATTATCCACAAACTTAAACACAGGGACAAAATTCTGATATGCTTCAACCGCTGCCTTGC
CCTATTCTGACATGATATGACTACCATTGTTATTGTACGTGGGCAGTGACGTCTTATCATATGTCAAAGTCATTGCGAAC
TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAGAGATTGCCGTCTGAAACTTTGTCTTTTCCGGGGACTCTAC
AA**CCCTTGTC**CTACTGATTA**TTTGTA**TGAATT**GGACAAT**TCAGATTAGTAGACAAGCGCGAGGAGGAAAAGAAATGACA
AAATTCCGATGGACAAGAGATAGGAAAAAAAGCTTCACCGATTCTAGACCGGAAAAAGTCGTATGACATCAGAACATGA
ATTTCAAGTTAGA**CAAGGAC**AAAATCAGGACAAATTGTAAGATAATAAAACTATTGATTGACGCCATTGCCCTTCC
TCCATTAAATCTCTGTTCTCTTACTTATATGATGATTAGGTATCTG**TATAAA**ACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGTTCGAAGACATTCTACGCATAATAAGAACAGGAGGAAATA**ATGCCAGACAATCTATCATTAC**
GCAGCTCTCAAAAGATTGAACCTCTCGCCAACCTTGAATCTTCAATGAGACCTTGCAGCAAATAATGTTGGATTGGAAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTATGAGGCATCGAGCTTGAAGAAAAGTAAGCTCAGAAAAACCTCA
CTCATTCTGGAAAGAAAATCTATTGAATATGTTGCTGACAAATCAATCTTGGGTGTTCTATTCTGGATTCAATTGTA
AGGACTTGAAGGCCGTCGAAAAAGAAAGGCAGGTTGGTCTGGTACAATTATTGTTACTCTGGCTGCTGAATGTTCAATATC
ACTTGGCAAATTGCGACTACAGGTCTACAACTGGGTCTAAATTGGTGGCAGTGTGGATAACAAATTGGATTGGTACGGTTCTG
TGGTTTGTCTTTGCCCTGTCAGCTGGCTTGTGCTGAGTGTGGATAACAAATTGGATTGGTACGGTTCTG

TATTGAATTTCAAAAATTCTACTTTGGATGGACGCAAAGAAGTAAATAATCATATTACATGGCATTACCAACCATATA
ATCCATATCTAATCTTAC**TTATA**TGTGTGAAATGTAAGAGGCCATTATCTTAGCCTAAAAAACCTCTTGGAACTTTC
AATACGCTTAAC TGCTCATTGCTATATTGAAGTA**CGG**ATTAGAAGCCG**CCGAG****CGGGCGACAGCCCT****CCGACGG**AAGACTCTCCTC
GCGTCCTCGTCTCACCGTCGCCTGAAACGCAGATGTGCCT**CGC**GCCGACTGCT**CCG**AACAATAAGATTCTACAATACT
TTTATGGTTATGAAGAGGAAAATTGGCAGTAACCTGG**CCCCA**AAACCTCAAATTAAACGATTCAAATTAAACACCATAAGGATG
ATGCATTAGTTTTAGCCTTATTTC**TGGGG**TAATTAAATCAGCGAATTGATTGATTTGATGTTAACAGATA**TATAA**ATGGAA
CTGCATAACCACCTTAACAACTTCAACATTTCAGTTGTATTACATTATTCAAATGTCATAAAAGTATCAACAAAAAT
TAATATAACCTCTATACTTAAACGTCAAGGAGAAAAACTATA**ATGACTAAATU****TATT**CAGAAGAAATTGATTGTACCTGAGTTCAA
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAAGAATTATAAGCATTATGATGCTAAACCGG
TTGTTGCTAGATGCCCTGGTAGAGTCATTAATTGGTGAACATATTGATTATTGTGACTCGGTTTACCTTAGCTATTGAT
GATATTGCTTTGCTGGGAAATTGCTGGGAAATTGCTGGGAAATTGCTGGGAAATTGCTGGGAAATTGCTGGGAAATTGCTGGGAA
CG
AC
GG
CA
GI
AG
AG
TO
CTAGTACTAGTGAAGAGTCTCTGCCAATAAGAAACAGGGCTTGTGACGATGT
ATTCAACAAGAGACTACTTAACAAACCTCCAGTGAGATTCAAGTCTAAAGCTATATCAGAGGCTAAGCATGTATTCTGAAT
TAAGAGTCTTGAAGGCTGTGAAATTATGACTACAGCGAGCTTACTGCCGACGAAGACTTTTAAAGCAATTGGTGCCTTGATG
GAGTCTCAAGCTTCTTGCATAAAACTTACGAATGTTCTTGCCAGAGATTGACAAAATTGTCATTGCTTGTCAAATGGATC
TGGTTCCCCTTGACCGGAGCTGGCTGGTTGTTACTGTTCACTGGTCCAGGGGGCCTAAATGGCAACATAGAAAAGGTAA
AAGCCCTGCCAATGAGTTCTACAAGGTCAAGTACCTAAGATCACTGATGCTGAGCTAGAAATGCTATCATCGTCTAAACCA
TTGGGCAGCTGTCTATATGAATTATAAGTAACTTCTTTTTACTTGTTCAGAACATTCTCATTTTTCTACTCATAACT
GCATCACAAAATACGCAATAAAACGAGTAGTAACTTTATAGTTCATACATGCTTCAACTACTTAATAATGATTGTATGATA
TTTCAATGTAAGAGATTGCTGATTATCCACAAATTAAAACACAGGGACAAAATTCTGATATGCTTCAACCGCTGCCTTGG
CCTATTCTGACATGATATGACTACCATTGTTAAGTACGTGGGGCAGTTGACGTTATCATATGTCAAAGTCATTGCGAAG
TTGGCAAGTTGCCAAGTGACGAGATGCGAGTAAAAGGATTGCCGTCTGAAACTTGTGCTTTTTCCGGGGACTCTAC
AA**CCCTTG**CCTACTGATTAA**TTTTGTAC**TGAATT**CAACAA**TTCAGATTTCAGACAAGCGCAGGGAGGAAAAGAAATGACA
AAATTCCGATGGACAAGAAGATAGGAAAAAAAGCTACCCGATTTCCTAACCGAAAAAGTCGTATGACATCAGAATGA
ATTTCAGTTAGA**CAAGGAC**AAAATCAGGACAAATTGTAAGATATAATAAACTATTGATTCAAGCGCCAATTGCCCTTCCA
TCCATTAAATCTCTGTTCTCTTACTTATATGATGATTGATCATCTG**TATAAA**ACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGTTCGAAGACATTCTACGCATAAAAGAATAGGAGGGAATAATGCCAGACAATCTATCATTACATT
GGGGCTTTCAAAAAGATTGAACTCTGCCAACTTATGGAACTCTCAATGAGACCTTGCCTAAATAATGTGGATTGGAAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTATGAGGGCACTGAGCTTGAAGAAAAGTAAGCTCAGAAAACCTCAATA
CTCATTCTGGAAAGAAAATCTATTGAAATATGTTGCTGACAAATCAATCTGGGTGTTCTATTCTGGATTCAATTGTTACTTCTGGCTGCTGAATGTTCAATATC
AGGACTTGAAGCCCCTGCAAAAAGAAAGGCGGGTTGGTCCCTGGTACAATTGTTACTTCTGGCTGCTGAATGTTCAATATC
ACTTGGCAAATTGCAAGTACAGGTCTACAACCTGGGTCTAAATTGGTGGAGAACAAATTGGATTGGTACGGTTCGT



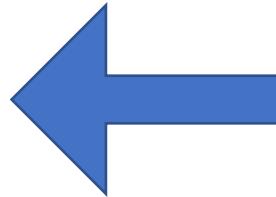
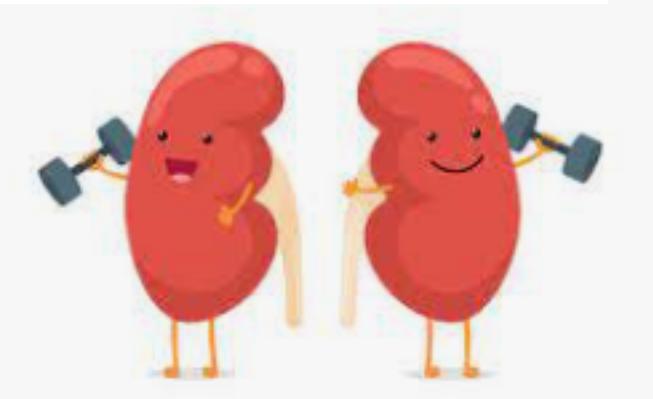
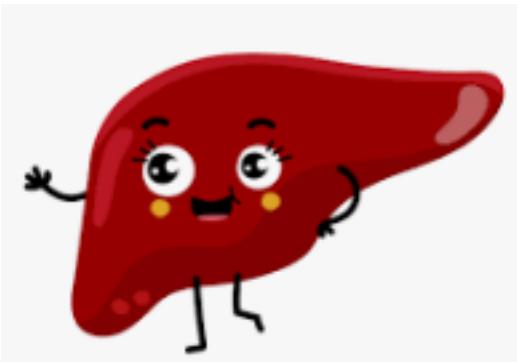
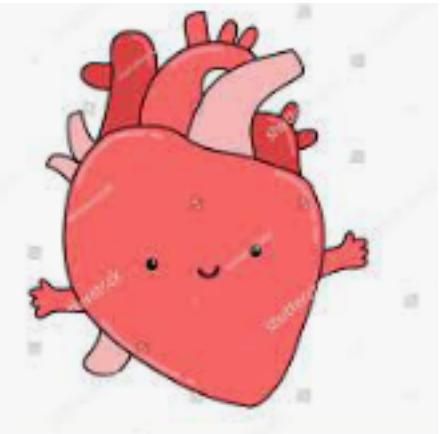
TATTGAATTTCAAAAATTCTTACTTTGGATGGACGCAAAGAAGTTAACATATTACATGGCATTACCACCATATA
ATCCATATCTAATCTTACTTATATGTTGGAATGTAAGAGGCCATTATCTTAGCCTAAAAAACCTCTTGGAACTTTC
AATACGCTTAAGTGCCTATTGCTATATTGAAGTACGGATTAGAAGCCGCCAGCGGGCGACAGCCCTCGACGGAAGACTCTCCTC
GCGTCCTCGTCTCACCGTCCGTTCTGAAACGCAGATGTGCCTCGGCCACTGCTCGAACAAATAAGATTCTACAATACI
TTTATGGTTATGAAGAGGAAAATTGGCAGTAACCTGGCCCCACAAACCTCAAAATTAAACGAATCAAATTAAACACCATAAGGATC
ATGCGATTAGTTTTAGCCTTATTCTGGGTAATTAAATCAGCGAAGCGATGATTGATCTATTAAACAGATATATAATGGAA
CTGCATAACCACCTTAACAACTTCAACATTTCAGTTGTATTACTTCTTATTCAAATGTCAAAAGTATCAACAAAAAI
TAATATACCTCTATACTTAACGTCAAGGAGAAAAACTATAATGACTAAATCTCATTAGAAGAAGTGATTGTACCTGAGTTCA
TAGCGCAAAGGAATTACCAAGACCATTGGCCAAAAGTGCCCGAGCATAATTAAAGAAATTATAAGCGCTATGCTAAACC
TTGTTGCTAGATGCCCTGGTAGAGTCATCTAATTGGTGAACATATTGATTATTGACTTCTCGGTTTACCTTAGCTATTGAT
GATATGCTTGCGCCGTCAAAGTTTGAACGAGAAAAATCCATTACCTTAATAAAATGCTGATCCCATTGCTCAAAGGA
CGATTGCCGTTGGACGGTCTTATGTCACAATTGATCCTCTGTGCGACTGGCTAATTACTTAAATGTGGTCTCATGTT
ACTCTTTCTAAAGAAACTTGCACCGGAAAGGTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTGTGAGGGTGTACCA
GGCAGTGGATTGTCTTCCGCCGCATTCAATTGCGTTGCTTAGCTGTTAAAGCGAATATGGCCCTGGTATCATAT
CAAGCAAATTAAATGCGTATTACGGTCGTTGCAGAACATTATGTTGGTGTAAACAATGGCGGTATGGATCAGGCTGCCTGTT
GTGAGGAAGATCATGCTCTATCGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTAAATTCCGCAATTAAAAACCATGAA
AGCTTGTATTGCAACACCCCTGTTGATCTAACAGTTGAAACCGCCCAACCAACTATAATTAAAGAGTGGTAGAAGTCAC
AGCTGCAAATGTTAGCTGCCACGTACGGTGTGTTACTTCTGGAAAAGAAGGATCGAGCACGAATAAGGTATCTAAGAG
TCATGAACGTTATTATGCCAGATATCACAAACATTCCACACCCCTGGAACGGCGATATTGAATCCGGCATCGAACGGTAA
CTAGTACTAGTTGAAGAGTCTCGCCAATAAGAAACAGGGCTTACTGTTGACGATGTCGACAATCCTGAATTGTTCTCGCGA
ATTCAAGAGAGACTACTAACACATCTCCAGTGAGATTCAAGTCTTAAAGCTATATCAGAGGGCTAACGATGTGATTCTGAAT
TAAGAGTCTTGAAGGCTGTGAAATTAAATGACTACAGCGAGCTTACTGCCAGCGAACACTTTCAAGCAATTGGTGCCTGATC
GAGTCTCAAGCTTCTGCGATAAAACTTACGAATGTTCTTGCCAGAGATTGACAAAATTGTTCCATTGCTTGTCAAATGGATC
TGGTCCCGTTGACCGGAGCTGGCTGGGTGTTGACTGTTACTGGTCCAGGGGCCAAATGGCAACATAGAAAAGGTA
AAGCCCTTGCCTAACGAGTTCTACAAGGTCAAGTACCTAACGATCACTGATGCTGAGCTAGAAAATGCTATCATGCTCTAA
TTGGGCAGCTGTCTATATGAATTATAAGTATACTTCTTTTACTTGTTCAGAACAACTCTCATTTTCTACTCATAACT
GCATCACAAATACGCAATAAACGAGTAGTAACACTTTATAGTCATACATGCTTCAACTACTAACAAATGATTGTATGATA
TTTCAATGTAAGAGATTGATTATCCACAAACTTAAACACAGGGACAAAATTCTGATATGCTTCAACCGCTGCTTTGG
CCTATTCTGACATGATGACTACCATTTGTTATTGACGTGGGGCAGTGACGTCTTACATGTCAAAGTCATTGCGAAC
TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAGAGATTGCCGTCTGAAACTTTGTCCTTTTCTCCGGGACTCTAC
AACCTTTGTCCTACTGATTAATTGACTGAATTGGACAATTCAAGATTAGACAGACGCGAGGAGGAAAAGAACATGACA
AAATTCCGATGGACAAGAAGATAGGAAAAAAAGCTTCAACGATTTCTAGACCGAAAAAGTCGTATGACATCAGAACATGA
ATTTCAAGTTAGACAAGGACAAATCAGGACAAATTGTAAGATATAATAAAACTATTGATTGACGCCATTGCTTCTTCC
TCCATTAAATCTCTGTTCTTACTTATATGATGATTAGGTATCTGATAAAACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGTTCGAAGACATTCTACGCATAATAAGAACAGGAGGAAATAATGCCAGACAACTATCATTACATT
GCGGCTCTCAAAAGATTGAACCTCGCCAACCTATGGAATTCTCAATGAGACCTTGCACCAATAATGTGGATTGGAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTTATGAGGCATCGAGCTTGAAGAAAAAGTAAGCTCAGAAAAACCTCAATA
CTCATTCTGGAAGAAAATCTATTGAATATGTTGCTGACAAATCAATCTGGGTGTTCTATTCTGGATTCAATTGAC
AGGACTTGAAGCCCCGTCGAAAAAGAAAAGGCGGGTTGGTCTGGTACAATTATTGTTACTTCTGGCTGCTGAATGTTCAATATC
ACTTGGCAAATTGCGACTACAGGTCTACAACCTGGTCTAAATTGGTGGCAGTGTGATAACAATTGGATTGGTACGGTTCGT

TATTGAATTTCAAAAATTCTACTTTTTGG**ATG**GACGCAAAGAAGTTAACATTAC**ATG**GATTACCAACCATA
ATCCATATCTAATCTTACTTAT**ATG**TGTGGAA**ATG**TAAAGAGCCCCATTATCTTAGCCTAAAAAACCTCTTGGAACTTCA
AATACGCTTAACGCTCATTGCTATATTGAAGTACGGATTAGAACGCCGAGCGGGCGACAGCCCTCGACGGAAGACTCTCCTC
GCGTCCTCGTCTCACCGTGCCTGAAACGCAG**ATG**TGCCTCGGCCGCAGTGCTCGAACAAATAAGATTCTACAATAC
TTTT**ATGGTT****ATG**AAGAGGAAAATTGGCAGTAACCTGGCCCCACAAACCTCAAAATTACGAATCAAATTACAAACCATTAGG**ATG**
ATGCGATTAGTTTTAGCCTATTCTGGGTAATTAATCAGCGAAGCG**ATG**ATTGGATCTATTAACAGATATATAA**ATG**GA
CTGCATAACCACCTTAACAACTTCAACATTTCAGTTGTATTACTCTTATTCAA**ATG**TCATAAAAGTATCAACAAAAAT
TAATATACCTCTACCTTAACGTCAAGGAGAAAAACTATA**ATG**ACTAAATCTCATTAGAAGAAGTGATTGTACCTGAGTTCAA
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAGAAATTATAAGCGCT**ATGATG**CTAAACC
TTGTTGCTAGATCGCTGGTAGAGTCATCTAATTGGTGAACATATTGATTATTGTGACTTCTCGGTTTACCTTAGCTATTGAT
GAT**ATG**CTTGCGCCGTCAAAGTTTGAAACGAGAAAATCCATTACCTTAATAA**ATG**GTGATCCCATTGCTCAAAGGA
CGATTGCGTTGGACGGTCTT**ATG**TCACAATTGATCCTCTGTGCGACTGGTCTAATTACTTAA**ATG**GGTCTCC**ATG**TTG
ACTCTTTCTAAAGAAACTTGACCGGAAAGGTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTGTGAGGGTG**ATG**TACCA
GGCAGTGGATTGTCTTCGGCCGCATTGTGCCGTGCTTAGCTGTTAAAGCGAAT**ATG**GGCCCTGGTTATCAT**AT**
CAAGAAAATTAA**ATG**CGTATTACGGTCGTTGCAGAACATT**ATG**TTGGTGTAAACA**ATG**GGCG**ATG**GATCAGGCTGCCTGT
GTGAGGAAGATC**ATG**CTCTATACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTAAATTCCGCAATTAAAAACC**ATG**AA
AGCTTGTATTGCGAACACCCTGTTGATCTAACAGTTGAAACCGCCCAACCAACTATAATTAAAGAGTGGTAGAAGTCAC
AGCTGCAA**ATG**TTTAGCTGCCACGTACGGTTGTTACTTCTGGAAAAGAAGGATCGAGCACGAATAAGGTAATCTAAGAG
TC**ATG**AACGTTATT**ATG**CCAGATATCACACATTCCACACCCCTGGAACGGCGATTGAATCCGGCATCGAACGGTAACAAAG
CTAGTACTAGTTGAAGAGTCTCGCCAATAAGAACAGGGCTTAGTGTGACG**ATG**TCGCACAATCCTGAATTGTTCTCGCGA
ATTCAAGAGACTACTTAACACATCTCAGTGAGATTCAAGTCTTAAAGCTATATCAGAGGGCTAACG**ATG**TGATTCTGAAT
TAAGAGTCTGAAGGCTGTGAAATT**ATG**ACTACAGCGAGCTTACTGCCGACGAAGACTTTCAAGCAATTGGTGCCTG**ATG**
GAGTCTCAAGCTTCTGCGATAAAACTTACGA**ATG**TCTTGCCAGAGATTGACAAAATTGTTCCATTGCTTGTCAA**ATG**GATC
TGGTTCCCGTTGACCGGAGCTGGCTGGGTTGTTACTGTTCACTGGTTCCAGGGGCCAA**ATG**GCAACATAGAAAAGGTA
AAGCCCTTGCCA**ATG**AGTTCTACAAGGTCAAGTACCTAAGATCACTG**ATG**GTGAGCTAGAAA**ATG**CTATCATCGCTCTAAACCA
TTGGGCAGCTGTCT**ATG**AATTATAAGTATACTTCTTTTTACTTGTCAAGAACAACTCTCATTTTTCTACTCATAACT
GCATCACAAATACGCAATAAAACGAGTAGTAACACTTTATAGTCATAC**ATG**CTTCAACTACTAACAA**ATG**ATTGT**ATG**ATA
TTTCA**ATG**TAAGAGATTGCGATTATCCACAAACTTAAACACAGGGACAAAATTCTGAT**ATG**CTTCAACCGCTCGTTGG
CCTATTCTTGAC**ATG**ATGACTACCAATTGTTATTGTACGTGGGCAGTGACGTCTTACAT**ATG**TCAAAGTCATTGCGAAG
TTGGCAAGTTGCCAACTGACGAG**ATG**CAGTAAAAGAGATTGCCGTCTGAAACTTTGTCTTTTTCCGGGGACTCTAC
AACCTTTGCTACTGATTAATTGTACTGAATTGGACAATTCAAGATTAGACAAAGCGCGAGGAGGAAAAGAA**ATG**ACA
AAATTCCG**ATG**GACAAGAAGATAGGAAAAAAAGCTTCACCGATTCTAGACCGAAAAAGTCGT**ATG**ACATCAGA**ATG**
ATTTCAGTTAGACAAGGACAAATCAGGACAAATTGTAAGATATAAAACTATTGATTCAAGGCCATTGCCCTTTCCA
TCCATTAAATCTCTGTTCTCTTACTTAT**ATG**ATGATTAGGTATCATCTGTATAAAACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGGTCGAAGACATCCTACGCATAATAAGAACAGGAGGAAATA**ATG**CCAGACAACTATCATTACATT
GCGGCTCTCAAAAGATTGAACCTCGCCA**ATG**GAATCTTCA**ATG**GAGACCTTGCACAAATA**ATG**TGGATTGGAAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTT**ATG**AGGCATCGAGCTTGAAGAAAAAGTAAGCTCAGAAAAACCTCAATA
CTCATTCTGGAAGAAAATCTATT**ATG**AA**ATG**TGGTCGTTGACAAATCAATTCTGGGTGTTCTATTCTGGATTCAATT**ATG**TAC
AGGACTTGAAGCCCCTCGAAAAAGAAAGGCAGGGTTGGTCTGGTACAATTATTGTTACTCTGGCTGCTGA**ATG**TTCAATATC
ACTTGGCAAATTGCAGCTACAGGTCTACAACTGGGTCTAAATTGGTGGCAGTGTGATAACAATTGGATTGGTACGGTTCGT

Extracting signal from noise

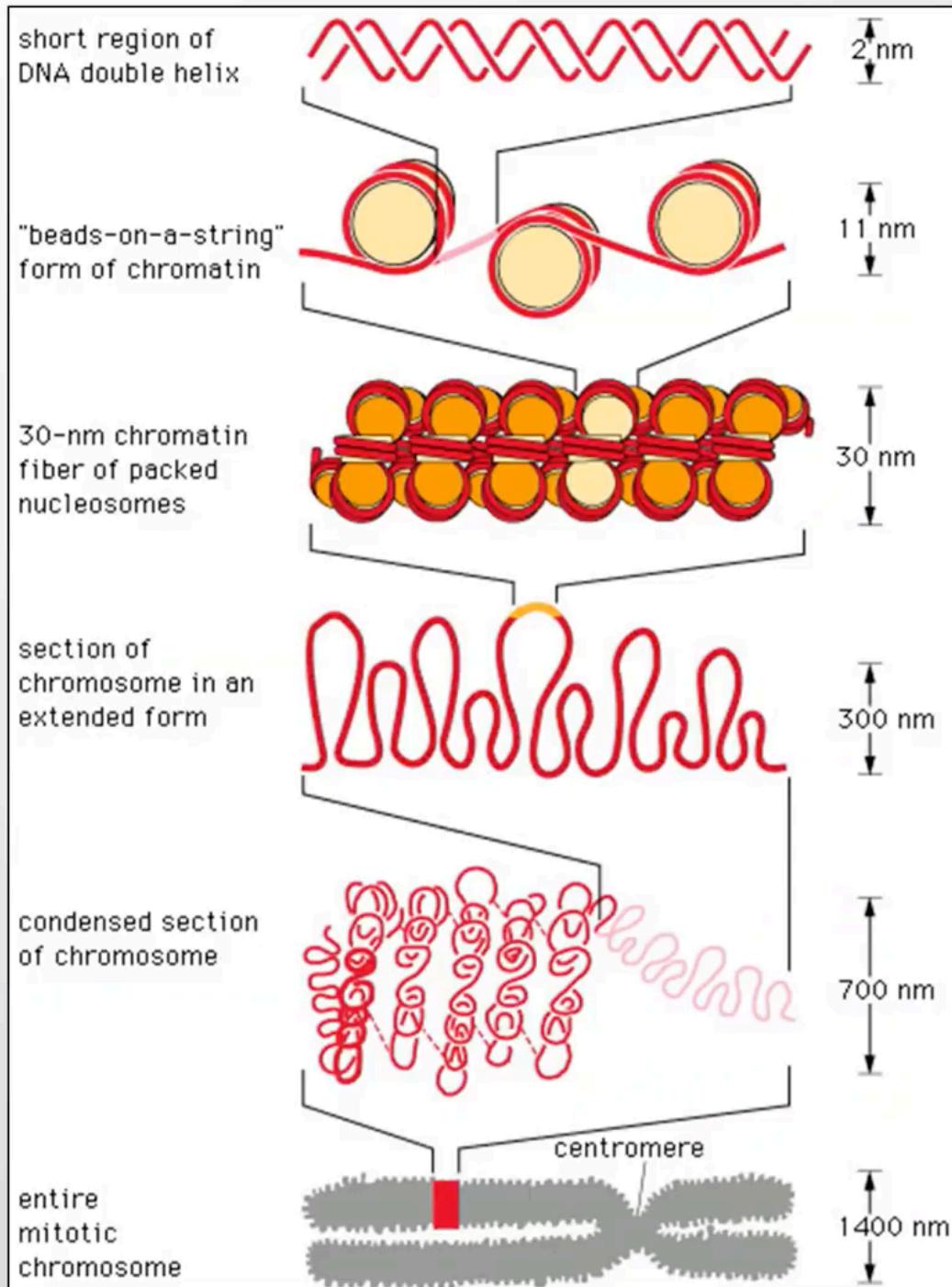
TATTGAATTTCAAAAATTCTACTT TTTGGATG GACGCCAAGAAGTTAATAATCATATTACATGCATTACACCATA
ATCCATATCTAATCTTA
AATACGCTTAACGCTC
GCGTCCTCGTCTCACCGGTCGGTCTCGAACAAACGAGATGCTCCGACTCGTCCGAACAAATAAGATTTCTACAAACT
TTTTATGGTTATG AAGAGGAAAATTGGCAGTAACCTGGCCCCACAAACCTCAAATTAAACGAATCAAATTAAACAACCATAAGGATG
ATGCGATTAGTTTTAGCCTTATTCTGGGTAAATTAACTCAGCGAAGCG**ATG**ATTTTGATCTATTAAACAGATATATAAATGGA
CTGCATAACCACCTTAACAACTTCAACATTTCAGTTGTATTACTCTTATTCAA**ATG**TCATAAAAGTATCAACAAAAAT
TAATATACCTCTACCTAACGTCAAGGAGAAAAACTATA**ATG**ACTAAATCTCATTCAAAGAAGTGATTGTACCTGAGTTCAA
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAAGAAATTATAAGCGCT**ATGATG**CTAAACC
TTGTTGCTAGATCGCTGGTAGAGTCATCTAATTGGTAAACATATTGATTATTGTGACTTCTCGGTTTACCTTAGCTATTGAT
GAT**ATG**CTTGCGCCGTCAAAGTTTGAAACGAGAAAAATCCATTACCTTAATAA**ATG**CTGATCCAAATTGCTCAAAGGA
CGATTGCCGTTGGACGGTTCTT**ATG**TCACAATTGATCCTCTGTGCGACTGGTCTAATTACTTAA**ATG**GGTCTCC**ATG**TTG
ACTCTTTCTAAAGAAACTTGACCCGAAAGGTTGCCAGTGCTCCTCTGGCCGGGCTGCAAGTCTGTGAGGGTG**ATG**TACCA
GGCAGTGGATTGTCTTCGCCGCATTGTGCCGTGCTTAGCTGTTAAAGCGAAT**ATG**GGCCCTGGTTATCAT**AT**
CAAGAAAATTAA**ATG**CGTATTACGGTCGTTGCAGAACATT**ATG**TTGGTGTAAACA**ATG**GCGGT**ATG**GATCAGGCTGCCTGT
GTGAGGAAGATC**ATG**CTCTACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTAAATTCCGCAATTAAAAACC**ATG**AA
AGCTTGTATTGCGAACACCCTGTTGATCTAACAGTTGAAACCGCCCAACCAACTATAATTAAAGAGTGGTAGAAGTCAC
AGCTGCAA**ATG**TTTAGCTGCCACGTACGGTTGTTACTTCTGGAAAAGAAGGATCGAGCACGAATAAGGTAATCTAAGAG
TC**ATG**AACGTTATT**ATG**CCAGATATCACACATTCCACACCCCTGGAACGGCGATTGAATCCGGCATCGAACGGTAACAAAG
CTAGTACTAGTTGAAGAGTCTCGCCAATAAGAACAGGGCTTACTGTTGACG**ATG**TCGCACAATTCTGAATTGTTCTCGCGA
ATTCAAGAGACTACTTAACACATCTCCAGTGAGATTCAAGTCTTAAAGCTATATCAGAGGGCTAACG**ATG**TGATTCTGAAT
TAAGAGTCTGAAGGCTGTGAAATT**ATG**ACTACAGCGAGCTTACTGCCGACGAAGACTTTCAAGCAATTGGTGCCTG**ATG**
GAGTCTCAAGCTTCTGCGATAAAACTTACGA**ATG**TCTTGCCAGAGATTGACAAAATTGTTCCATTGCTTGTCAA**ATG**GATC
TGGTTCCCCTTGAACGGAGCTGGCTGGGGTGGTTGACTGTTCACTGGTCCAGGGGGCCAA**ATG**GCAACATAGAAAAGGTA
AAGCCCTTGCCT**ATG**AGTTCTACAAGGTCAAGTACCTAAGATCACTG**ATG**CTGAGCTAGAAA**ATG**CTATCATCGCTCTAAACCA
TTGGGCAGCTGTCT**ATG**AATTATAAGTATACTTCTTTTTACTTGTCAAGAACAACTCTCATTTTTCTACTCATAACT
GCATCACAAACACGCAATAATAACGAGTAGTAACACTTTATAGTCATAC**ATG**CTTCAACTACTAAATAA**ATG**ATTGT**ATG**ATA
TTTCA**ATG**TAAGAGATTGCGATTATCCACAAACTTAAACACAGGGACAAAATTCTGAT**ATG**CTTCAACCGCTCGTTGG
CCTATTCTTGAC**ATG**ATGACTACCAATTGTTATTGTACGTGGGGCAGTGACGTCTTACAT**ATG**CTAAAGTCATTGCGAAG
TTGGCAAGTTGCCAACTGACGAG**ATG**CAGTAAAAGAGATTGCCGTCTGAAACTTTGTCCTTTTTCCGGGGACTCTAC
AACCCCTTGTCTACTGATTAATTGTACTGAATTGGACAATTCAAGATTCTGAT**ATG**CCAGACAACTCATCATTACATT
AAATTCCG**ATG**GACAAGAAGATAGGAAAAAAAGCTTCAACGATTTCCTAGACCGAAAAAGTCGT**ATG**ACATCAGA**ATG**
ATTTCAAGTTAGACAAGGACAAATCAGGACAAATTGTAAGATATAATAAAACTATTGATTCAAGGCCAATTGCCCTTTCCA
TCCATTAAATCTCTGTTCTCTTACTT**ATG**ATGATTAGGTATCATCTGTATAAAACTCCTTCTTAATTCACTCTAAAGCAT
CCATAGAGAAGATCTTCGGGTCGAAGACATCCTACGCATAATAAGAATAGGAGGGAATA**ATG**CCAGACAACTCATCATTACATT
GC GGCTCTTCAAAAGATTGAACCTCGCCA**ATG**GAATCTTCA**ATG**GAGACCTTGCCTAA**ATG**TGGATTGGAAAAA
TATAAGTCATCTCAGAGTAATATAACTACCGAAGTT**ATG**AGGCATCGAGCTTGAAGAAAAAGTAAGCTCAGAAAAACCTCAATA
CTCATTCTGGAAAGAAAATCTATT**ATG**AA**ATG**TGGTCGTTGACAAATCAATTCTGGGTGTTCTATTCTGGATTCAATT**ATG**TAC
AGGACTTGAAGCCCCTCGAAAAAGAAAGGCAGGGTTGGTCTGGTACAATTATTGTTACTCTGGCTGCTGA**ATG**TTCAATATC
ACTTGGCAAATTGCAGCTACAGGTCTACAACTGGGTCTAAATTGGTGGCAGTGGATAACAATTGGATTGGTACGGTTCTG

One genome – Many cells



```
ACCAGTTACGACGGTCA  
GGGTACTGATACCCCAA  
ACCGTTGACCGCATTAA  
CAGACGGGGTTGGGTT  
TTGCCCCACACAGGTAC  
GTTAGCTACTGGTTAG  
CAATTTACCGTTACAAC  
GTTTACAGGGTTACGGT  
TGGGATTGAAAAAAAG  
TTTGAGTTGGTTTTTC  
ACGGTAGAACGTACCGT  
TACCAAGTA
```

DNA-packing



Why packing

- DNA is very long
 - stretch out the entire human genome: 300 x distance between earth and sun

- Cell is very small

Compression

- Chromosome is 50,000 times shorter than extended DNA

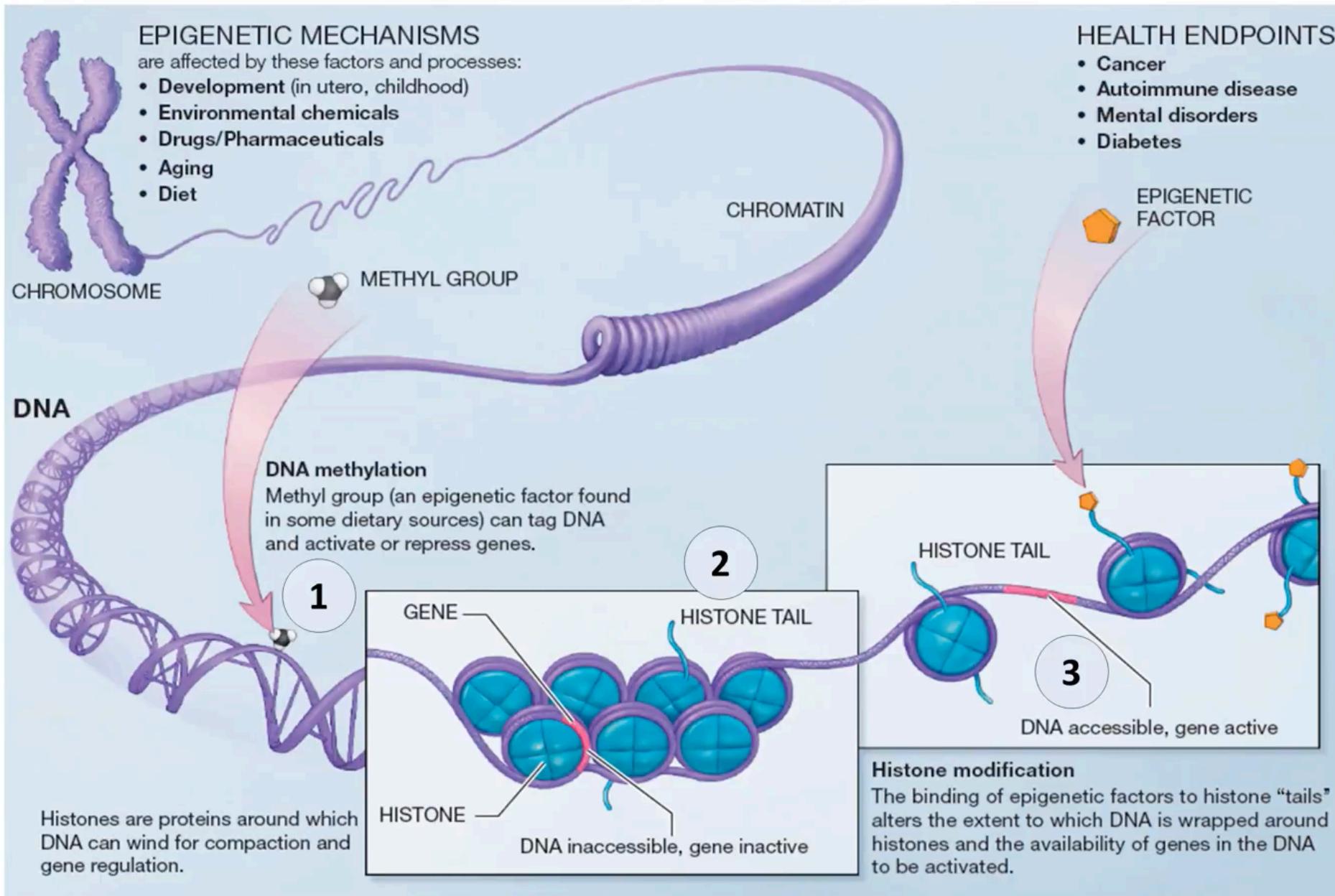
Using the DNA

- Before a piece of DNA is used for anything, this compact structure must open locally

Now emerging:

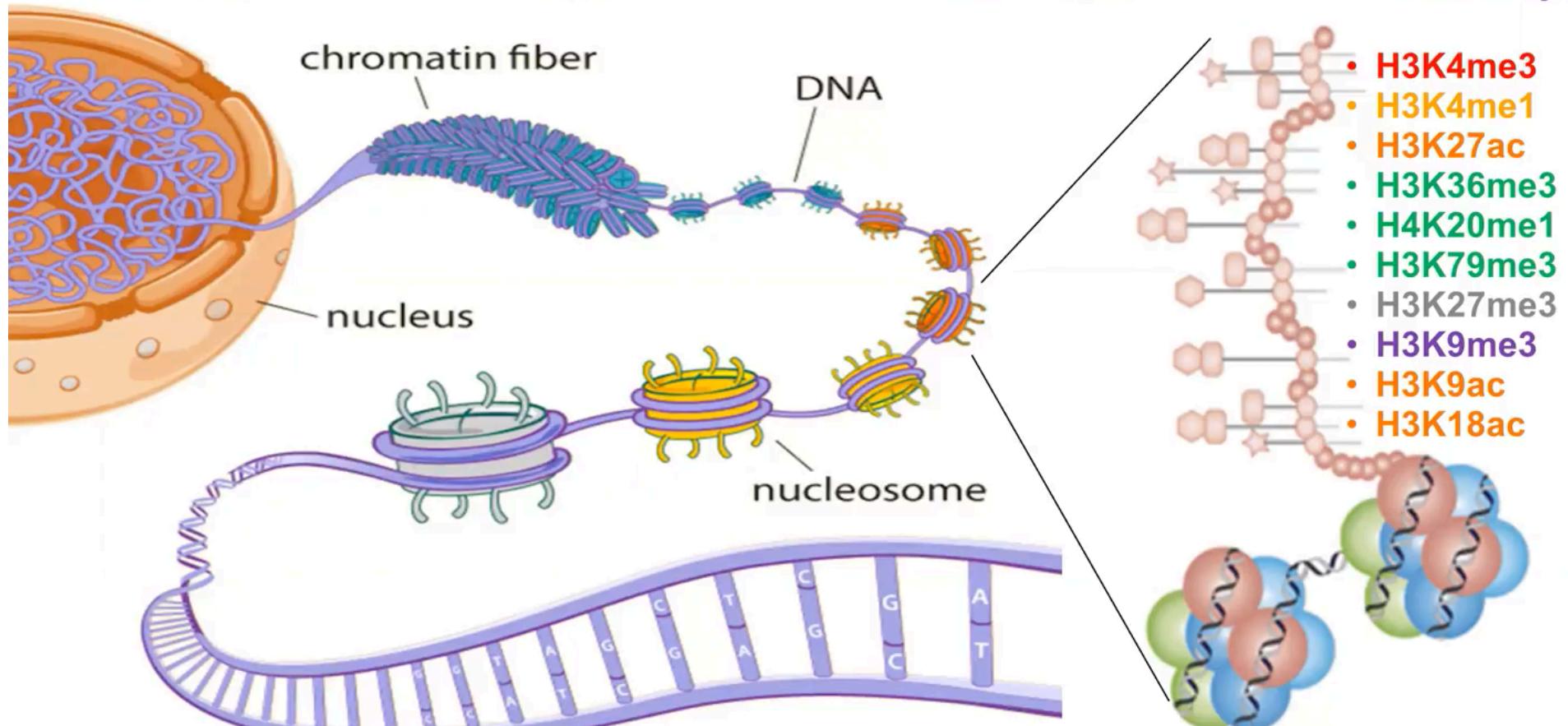
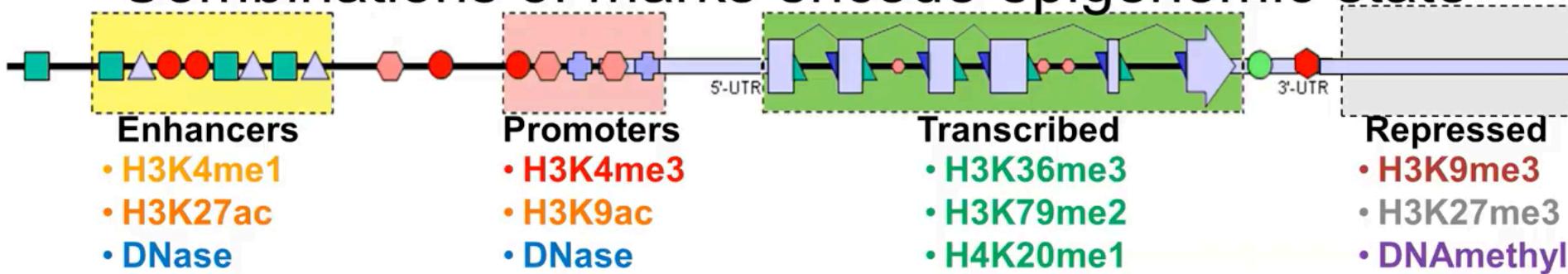
- Role of accessibility
- State in chromatin itself
- Role of 3D interaction

Three types of epigenetic modifications



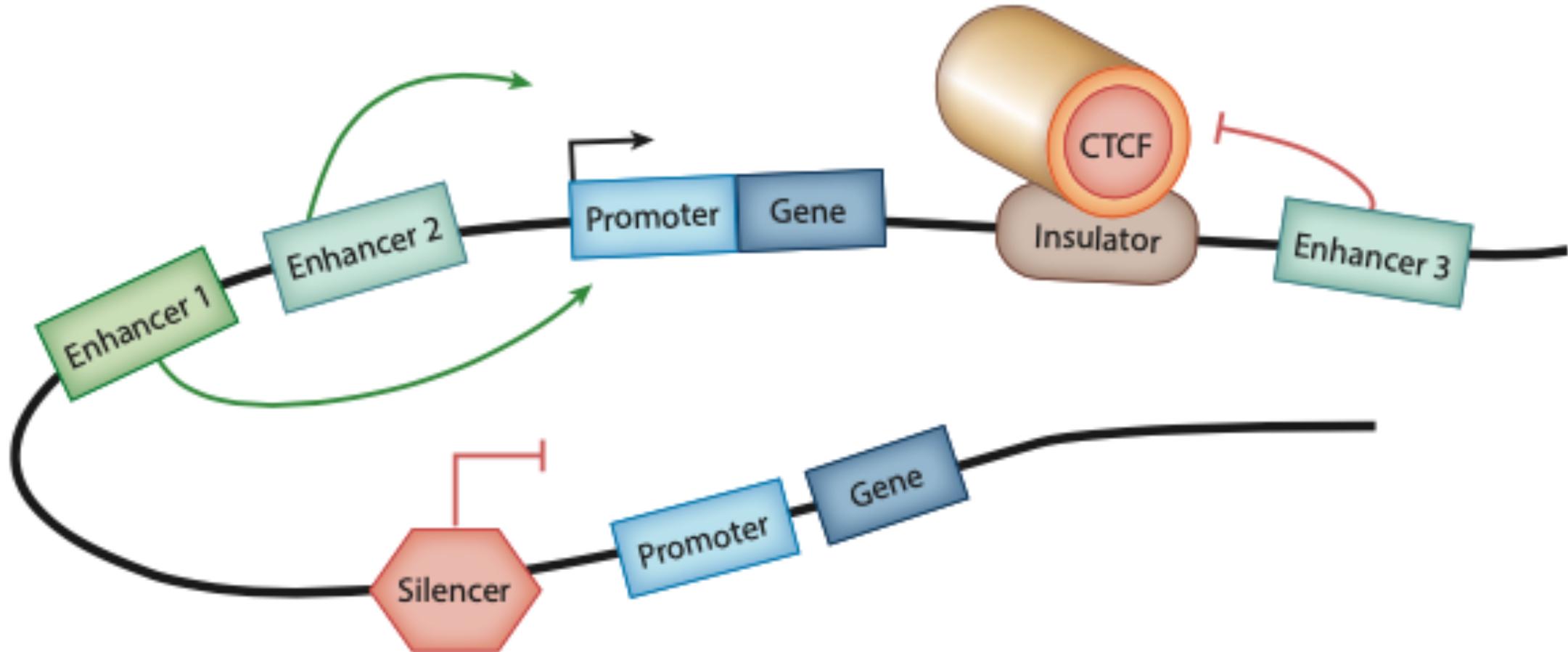
<https://en.wikipedia.org/wiki/Epigenetics>

Combinations of marks encode epigenomic state

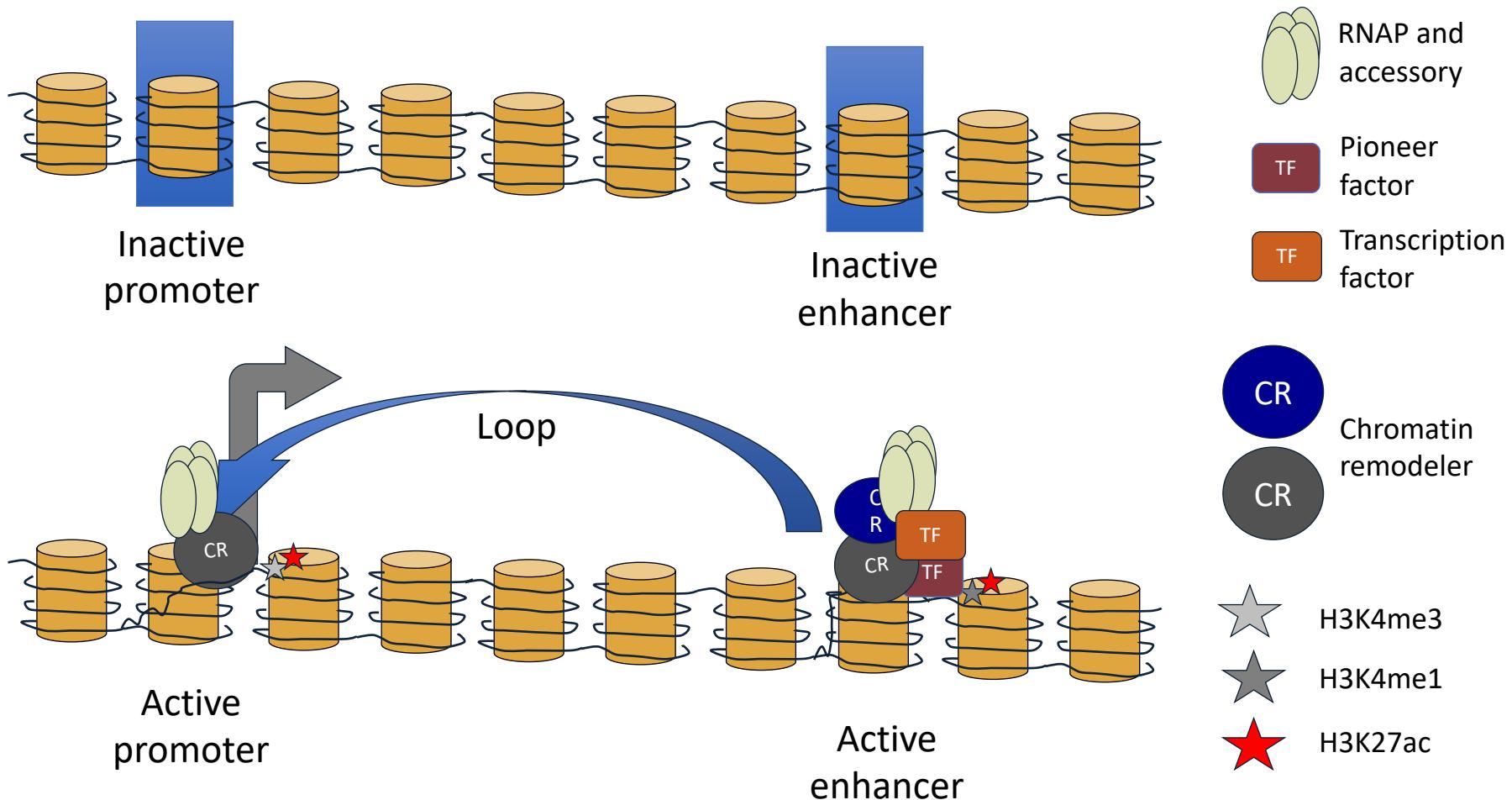


- 100s of known modifications, many new still emerging
- Systematic mapping using ChIP-, Bisulfite-, DNase-Seq

Types of gene regulatory elements



A model for step-wise activation of enhancers



Reviewed by Zaret & Carroll, Gene & Dev 2011;

DNA Methylation

Genetics

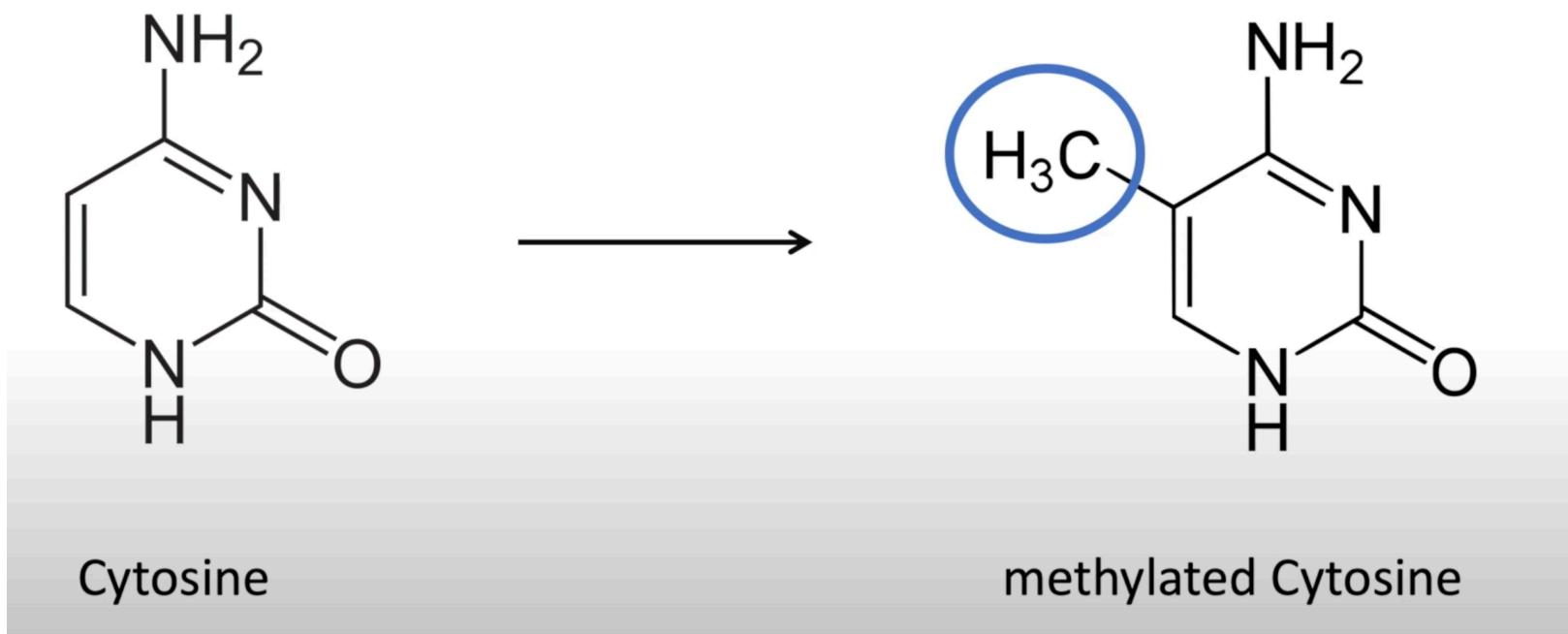


Changes in Phenotype

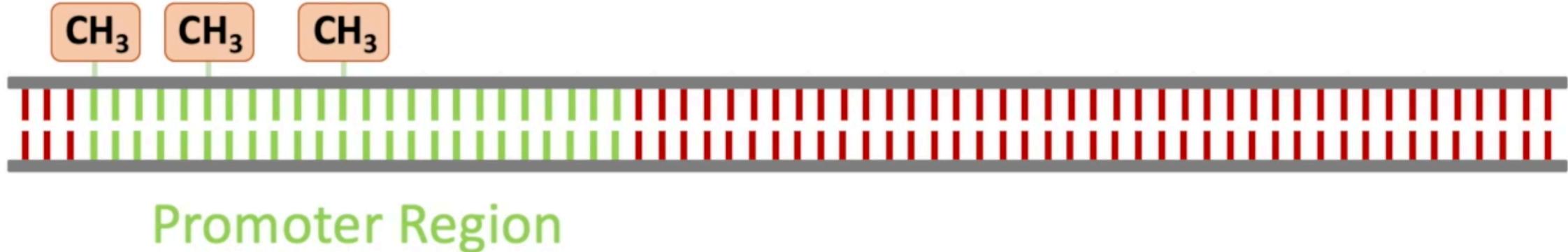
Epigenetics



DNA Methylation

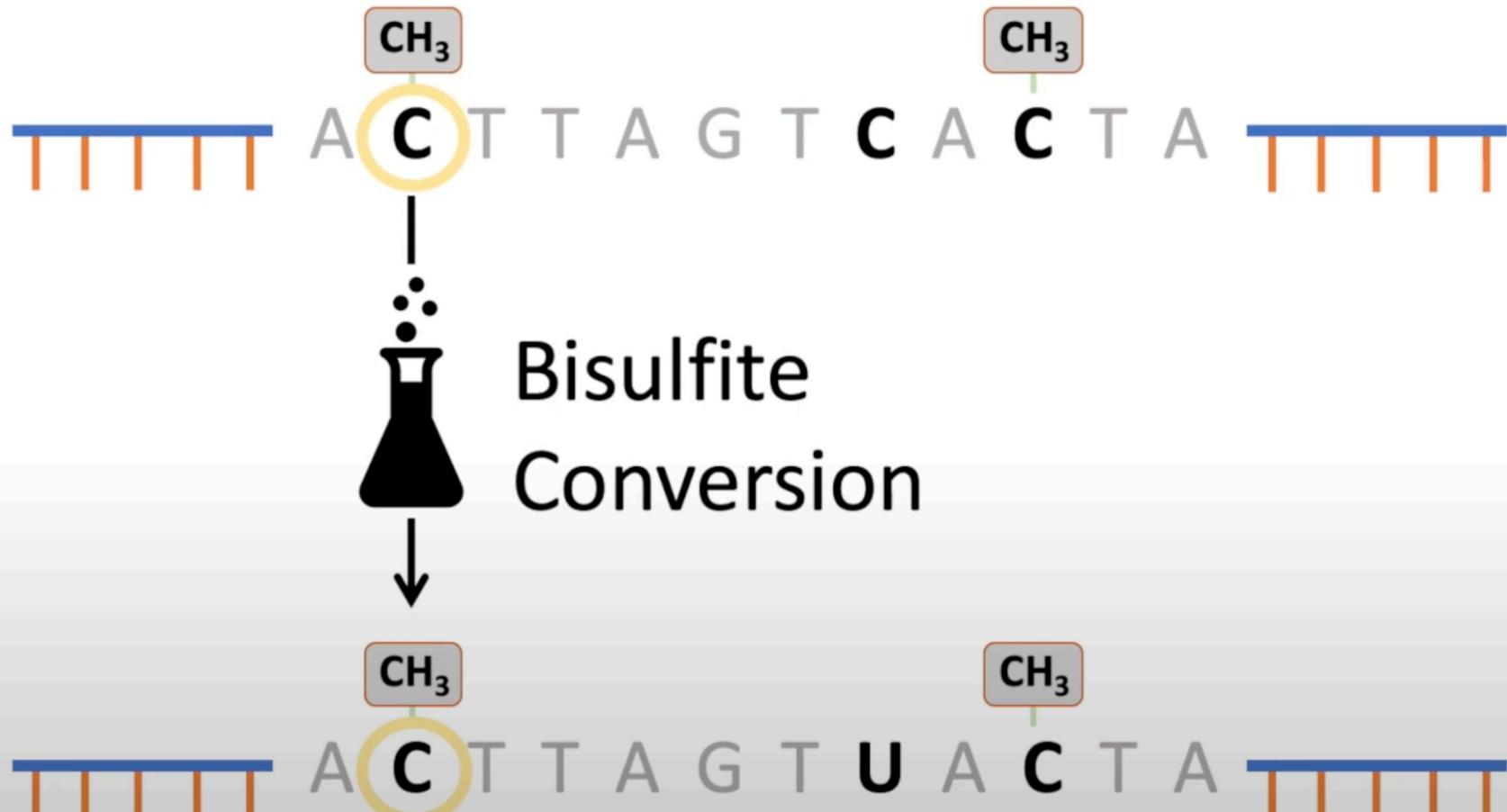


Function of DNA methylation

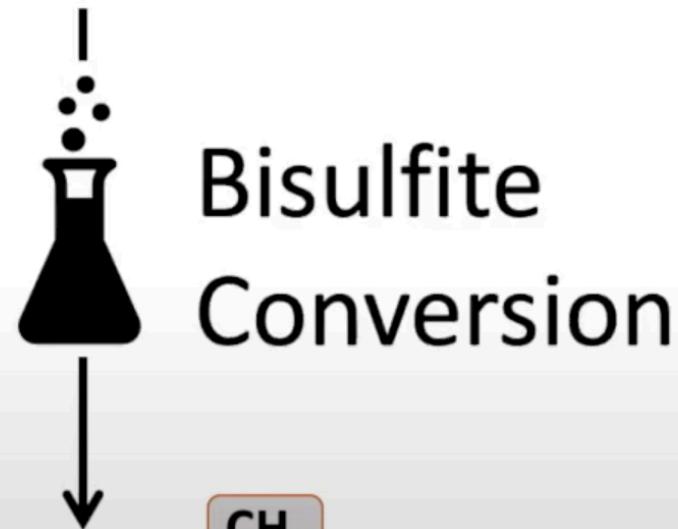


DNA Methylation → Inactivation of Genes

Bisulfite Conversion



Bisulfite Conversion



Cytosine to Uracil

Bisulfite Conversion

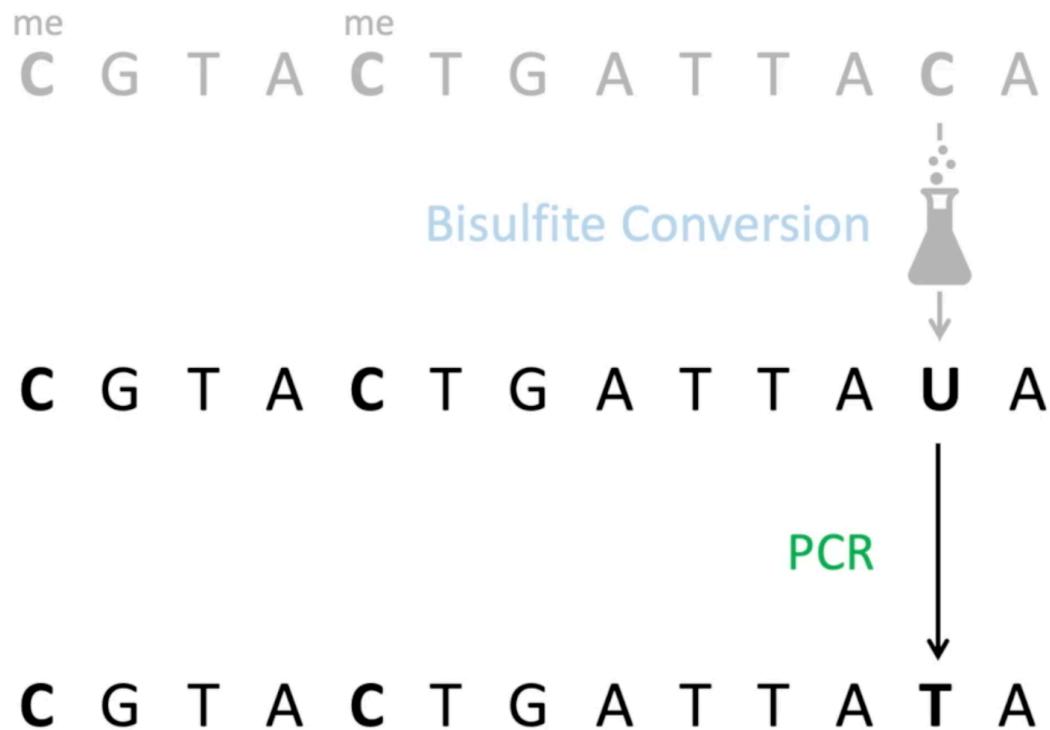
me
C G T A C T G A T T A C A

Bisulfite Conversion

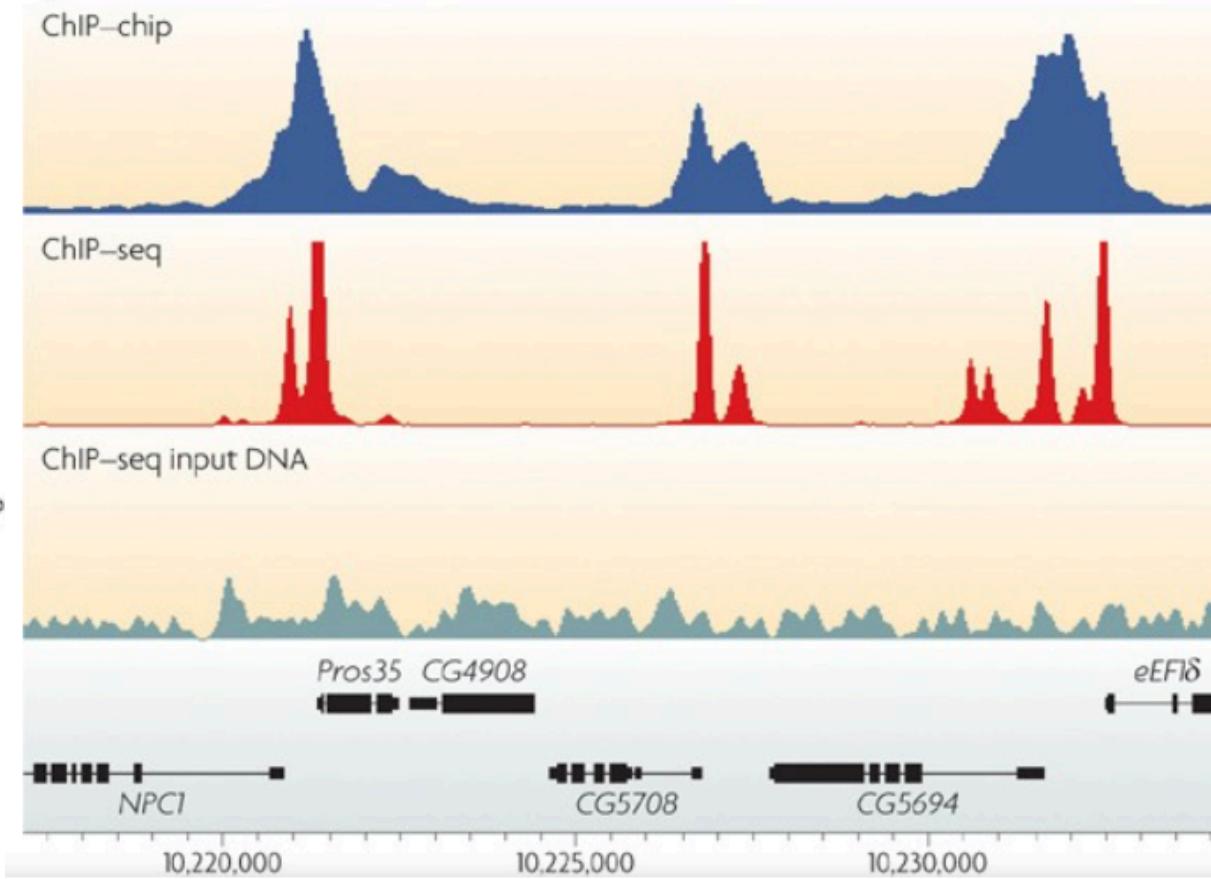
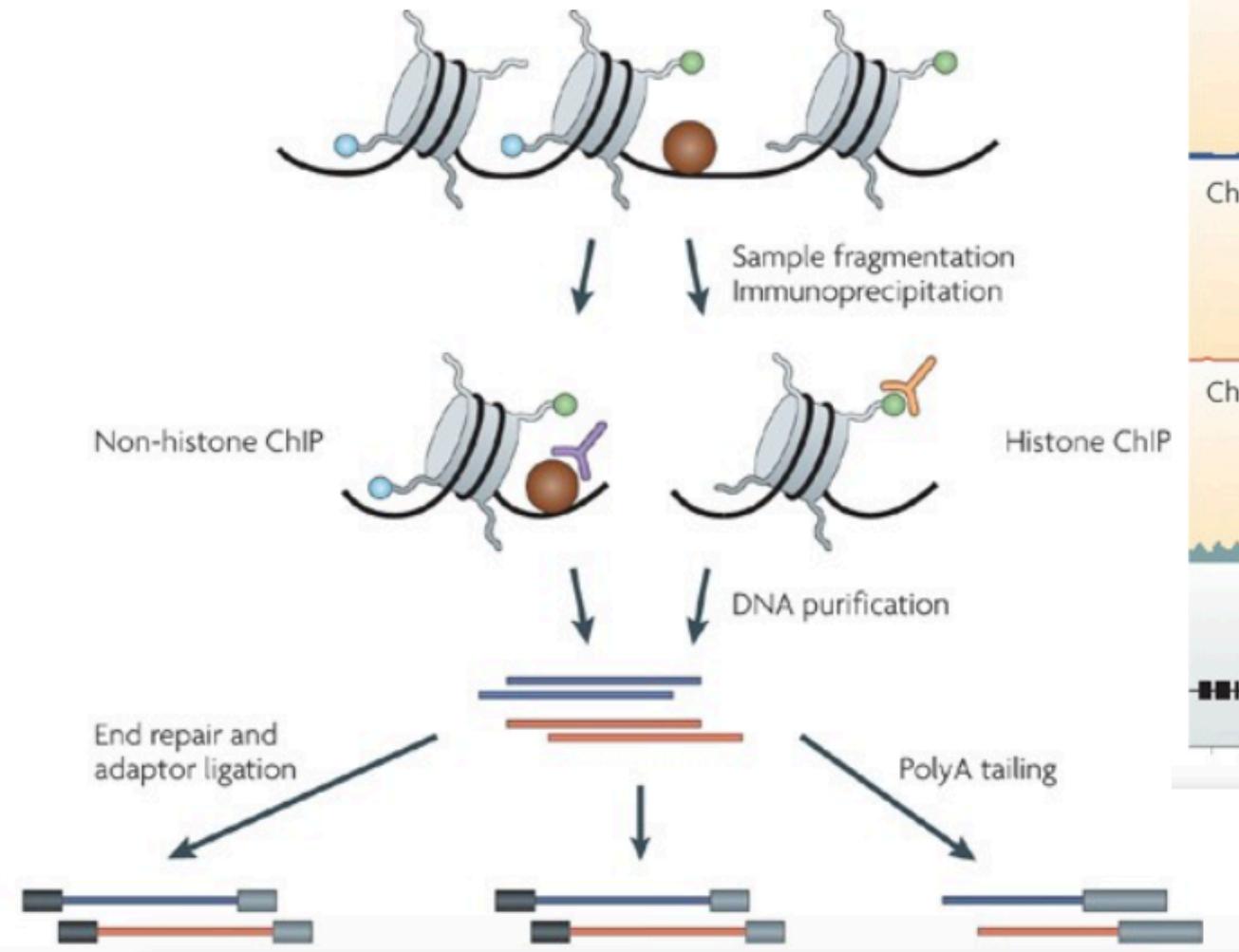
↓

C G T A C T G A T T A U A

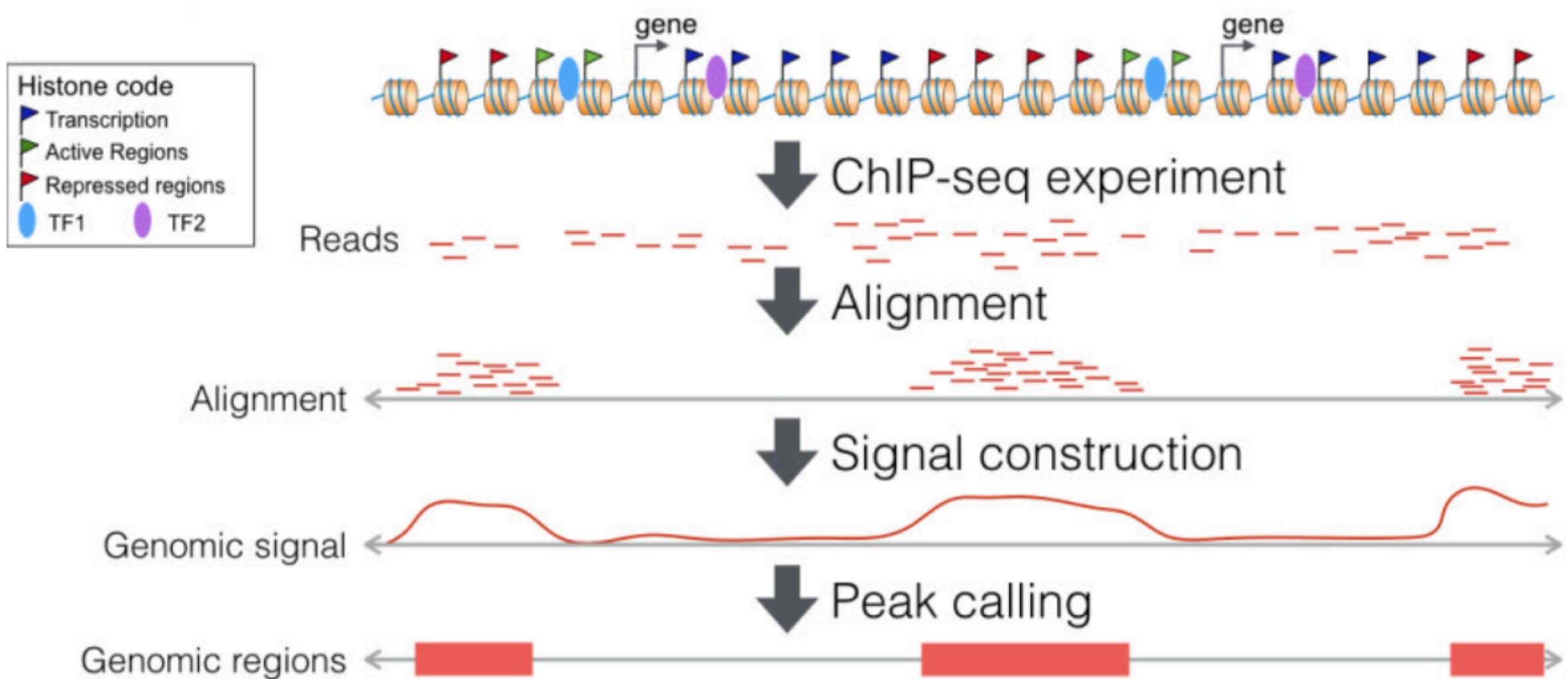
Bisulfite Conversion



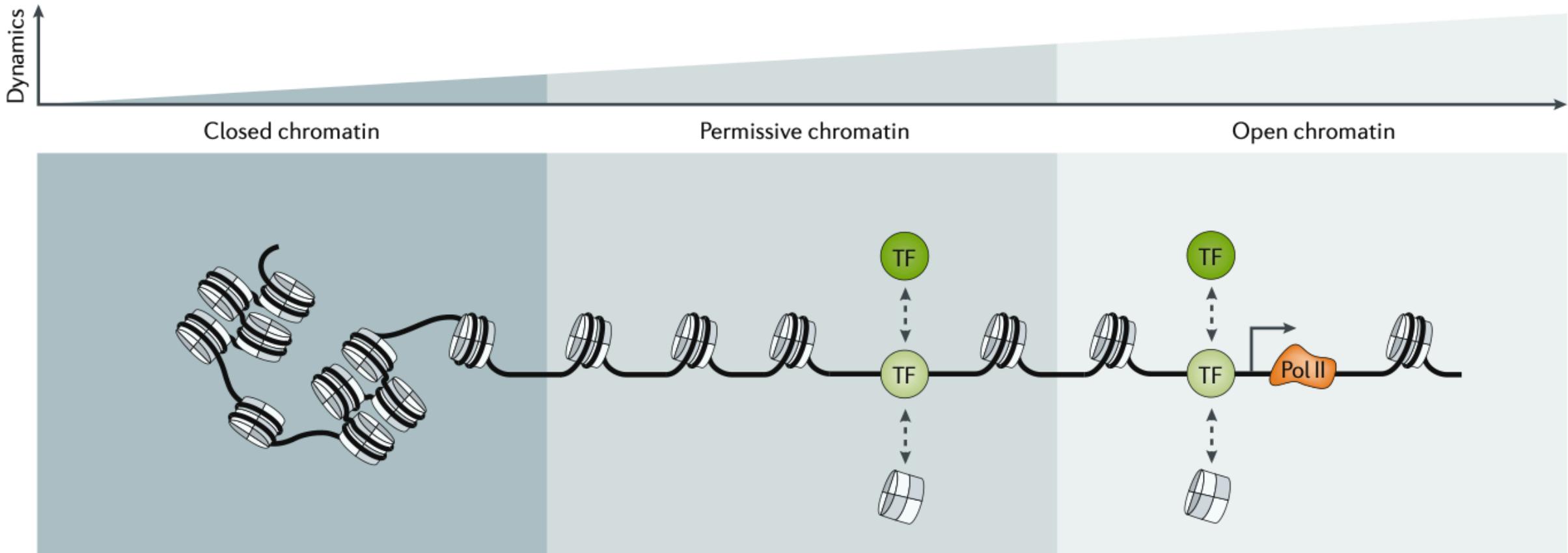
Chromatin Immunoprecipitation Sequencing (ChIP-Seq)



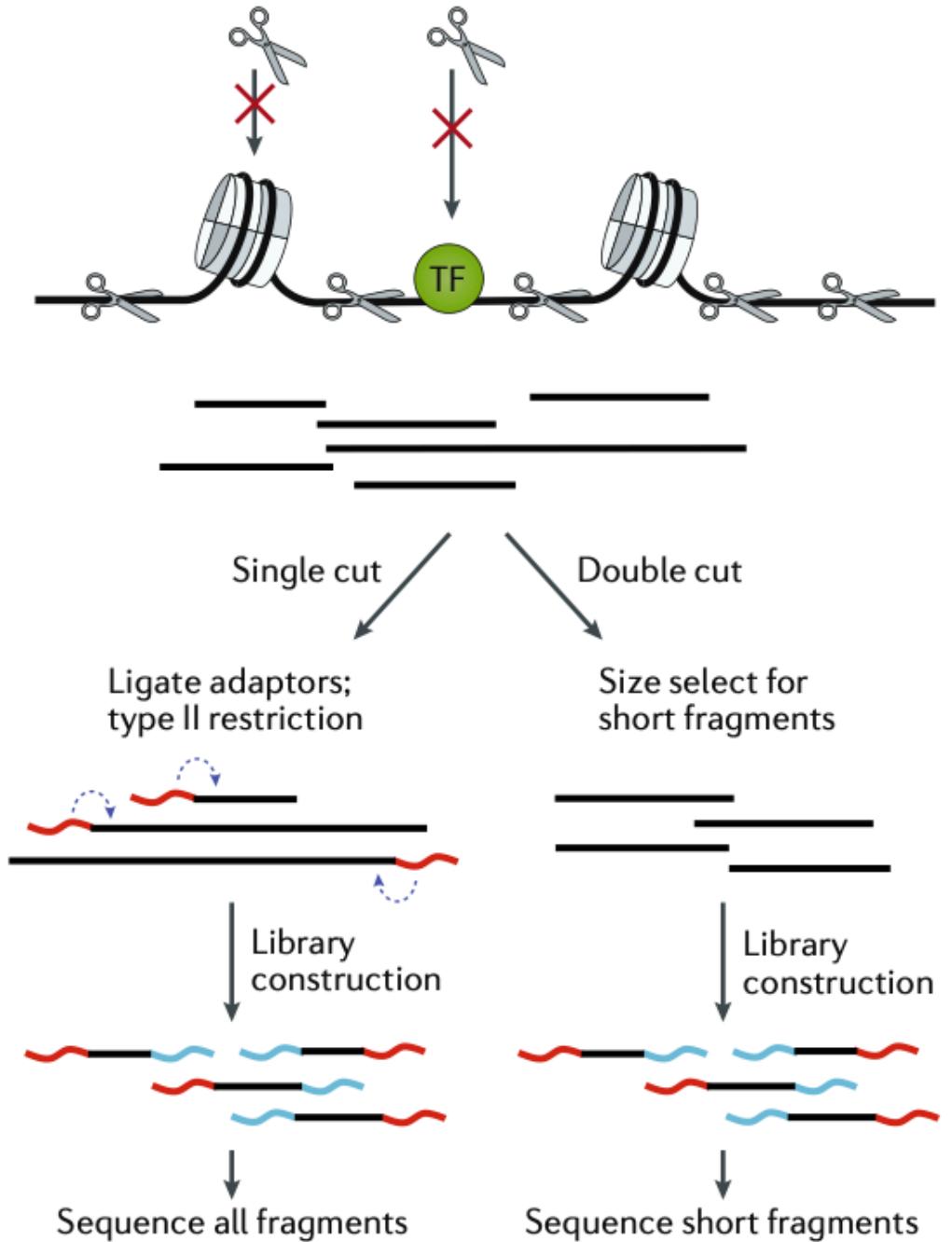
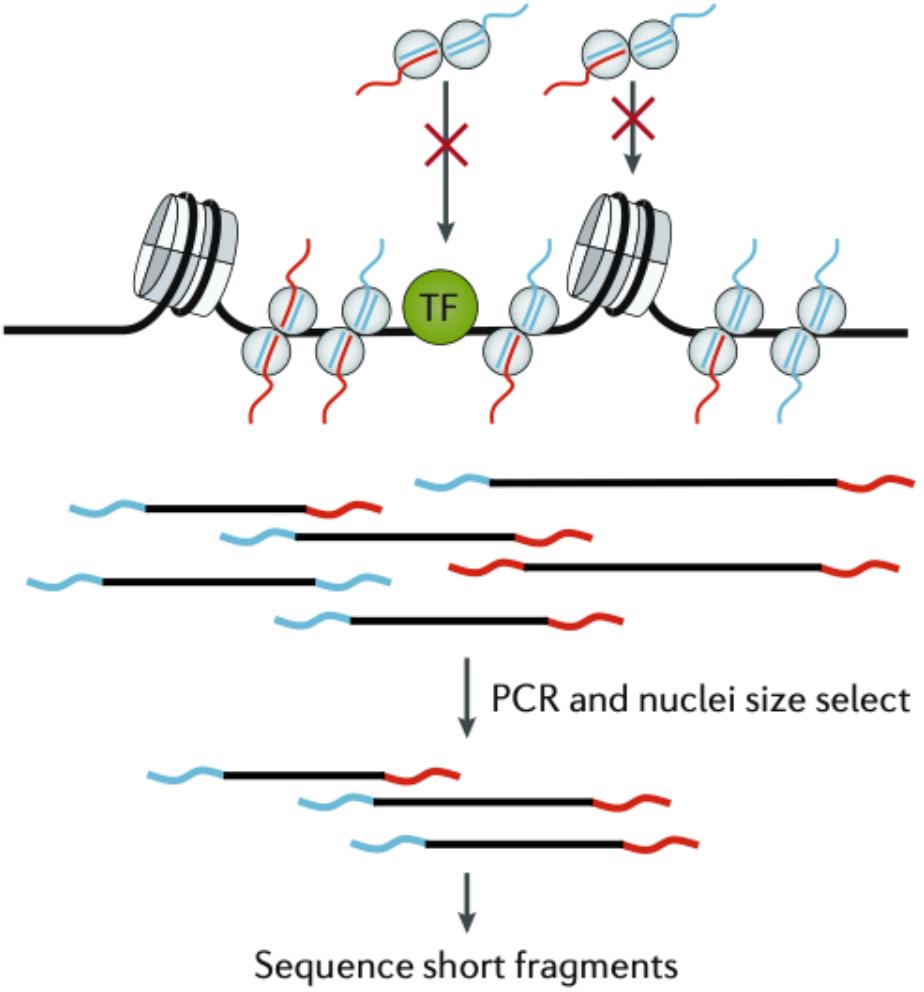
Typical pipeline for analysis of histone modification

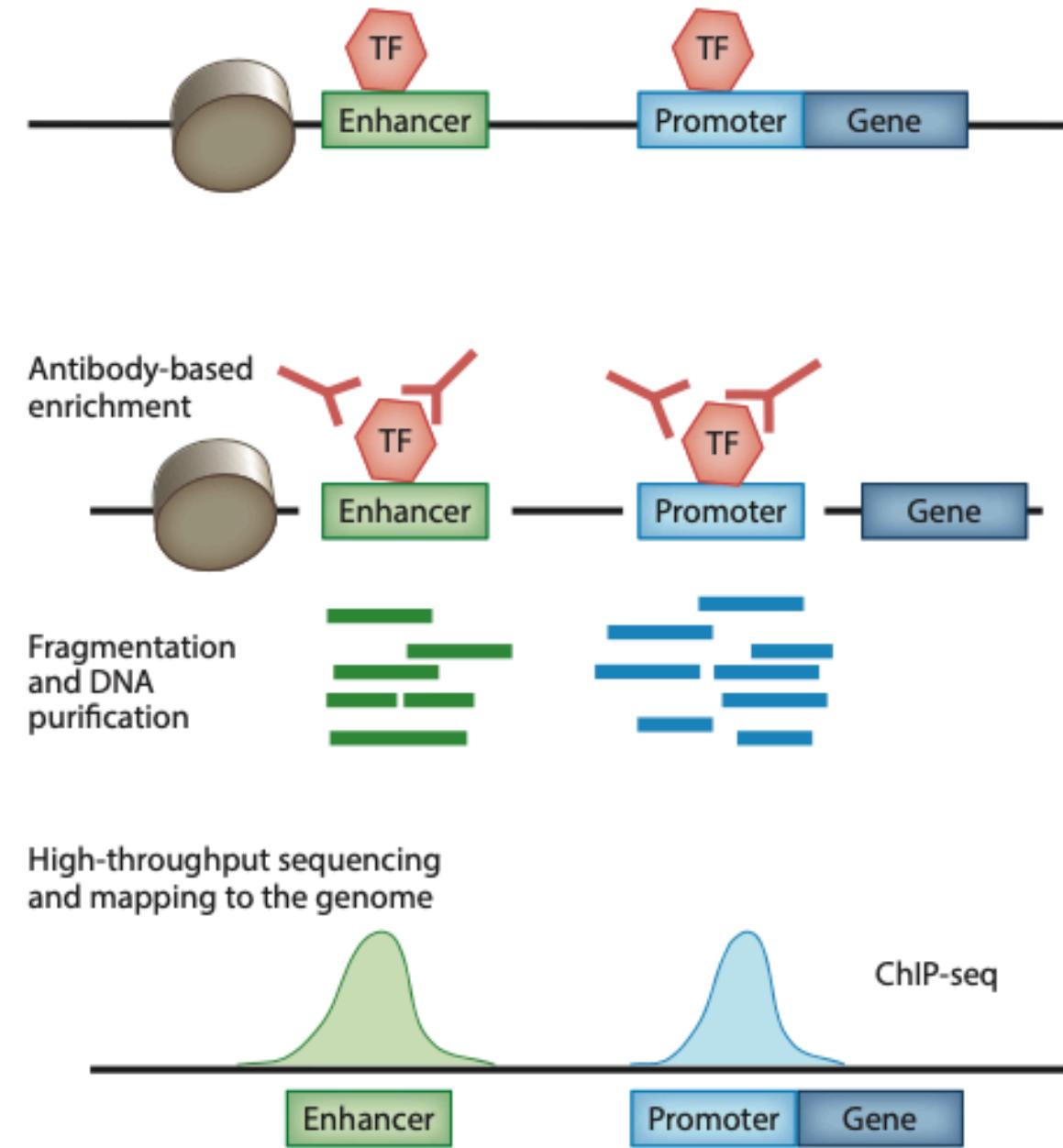
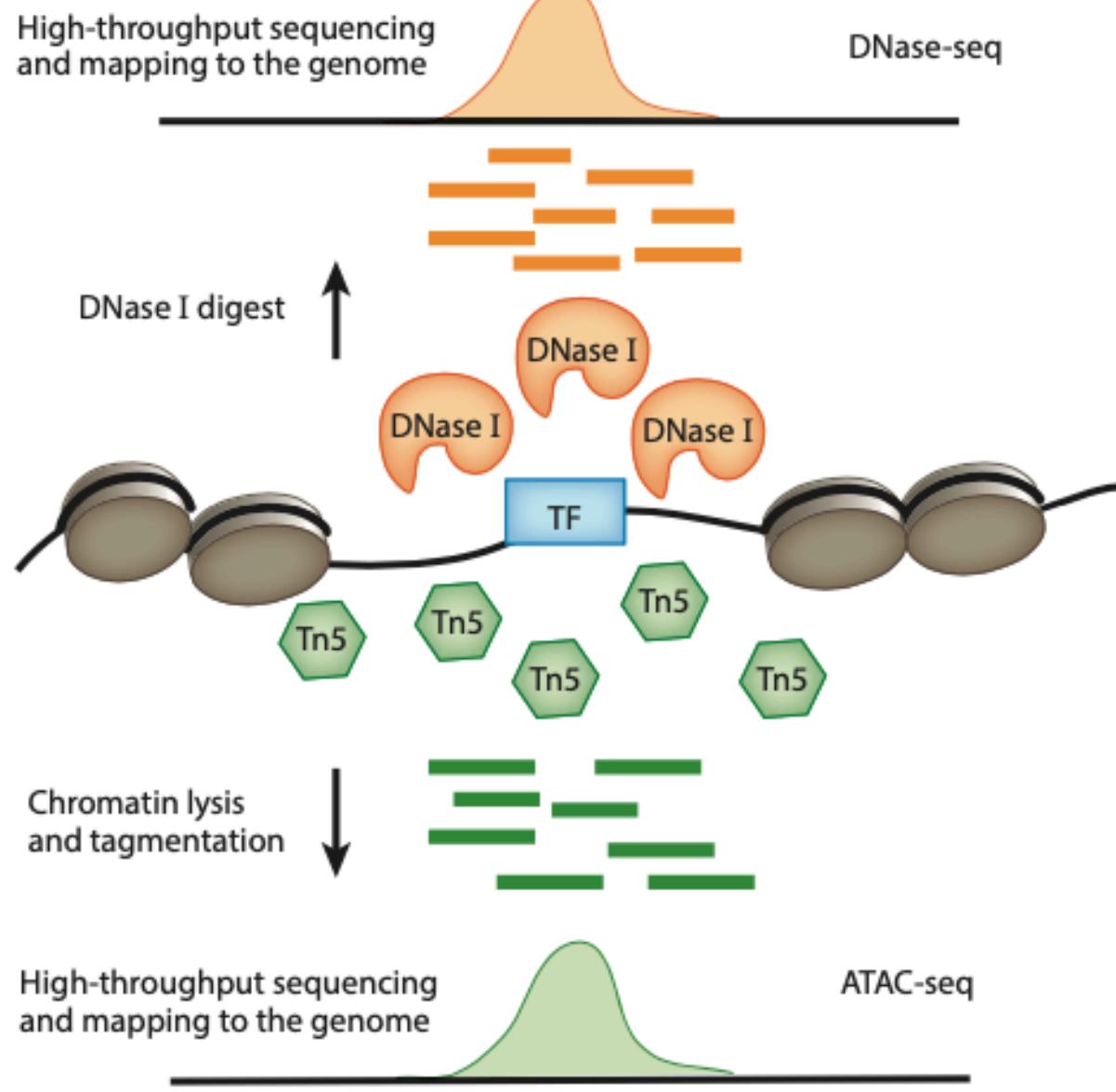


A continuum of accessibility states broadly reflects the distribution of chromatin dynamics across the genome

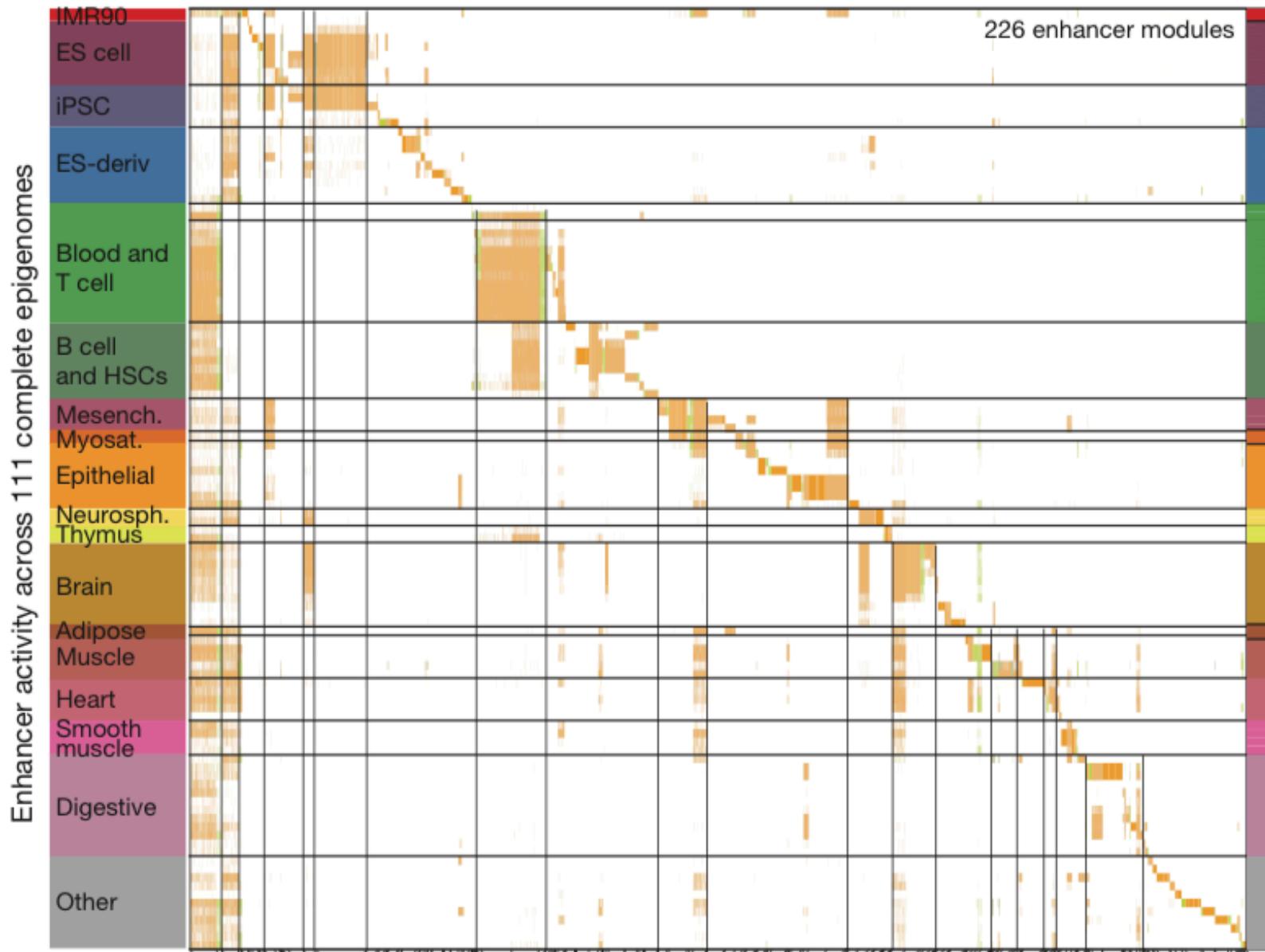


In contrast to closed chromatin, permissive chromatin is sufficiently dynamic for transcription factors to initiate sequence-specific accessibility remodeling and establish an open chromatin conformation (illustrated here for an active gene locus). Pol II, RNA polymerase II; TF, transcription factor.

a DNase-seq**b ATAC-seq (Assay for Transposase-Accessible Chromatin using sequencing)**

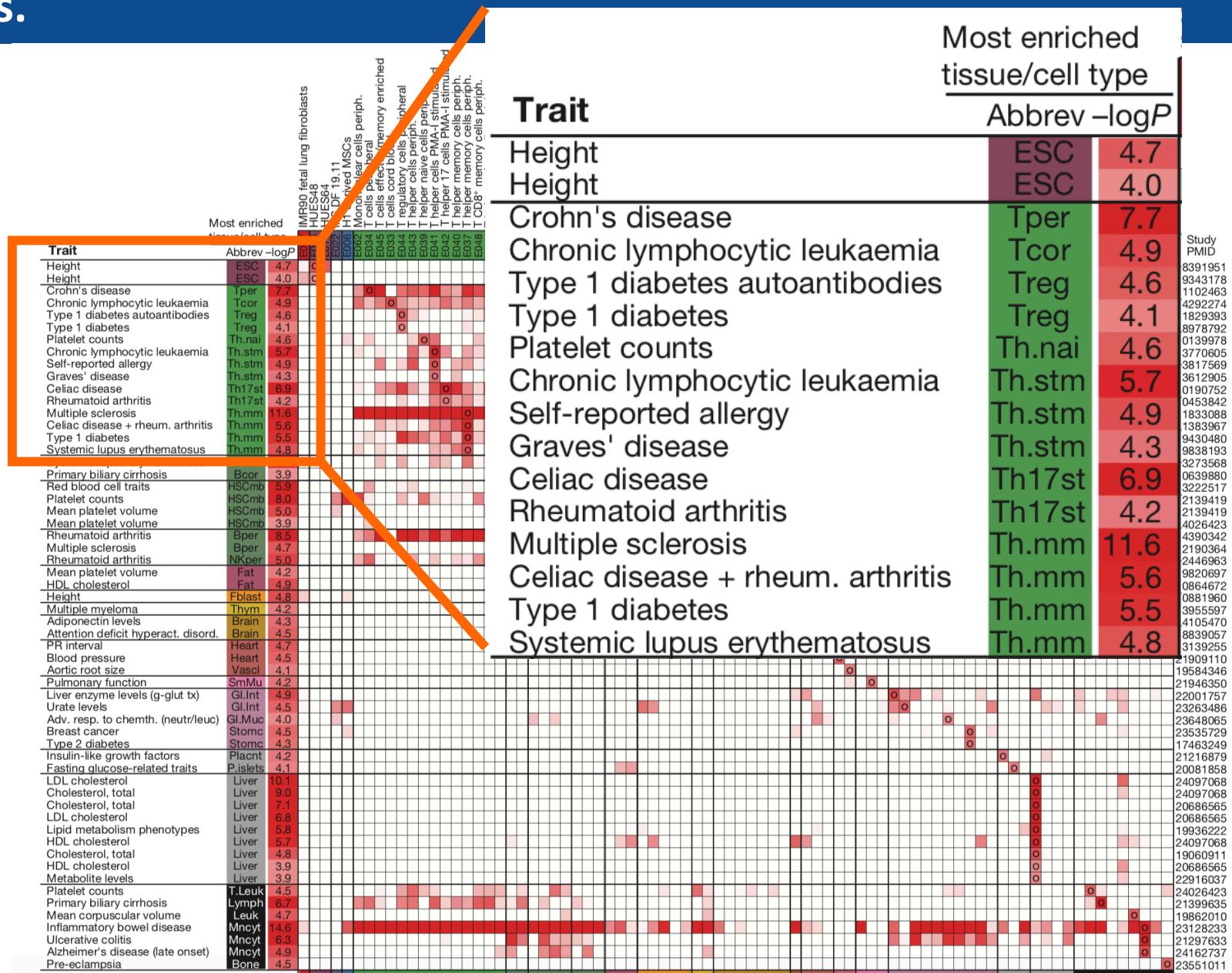


Regulatory modules from epigenome dynamics



Roadmap Epigenomics Consortium Nature 2015, 518: 317-330

Epigenomic enrichments of genetic variants associated with diverse traits.



Roadmap Epigenomics Consortium
Nature 2015, 518: 317-330

Epigenomics

Epigenomics focuses on genome-wide characterization of reversible modifications of DNA or DNA-associated proteins, such as DNA methylation, chromatin accessibility or histone acetylation.

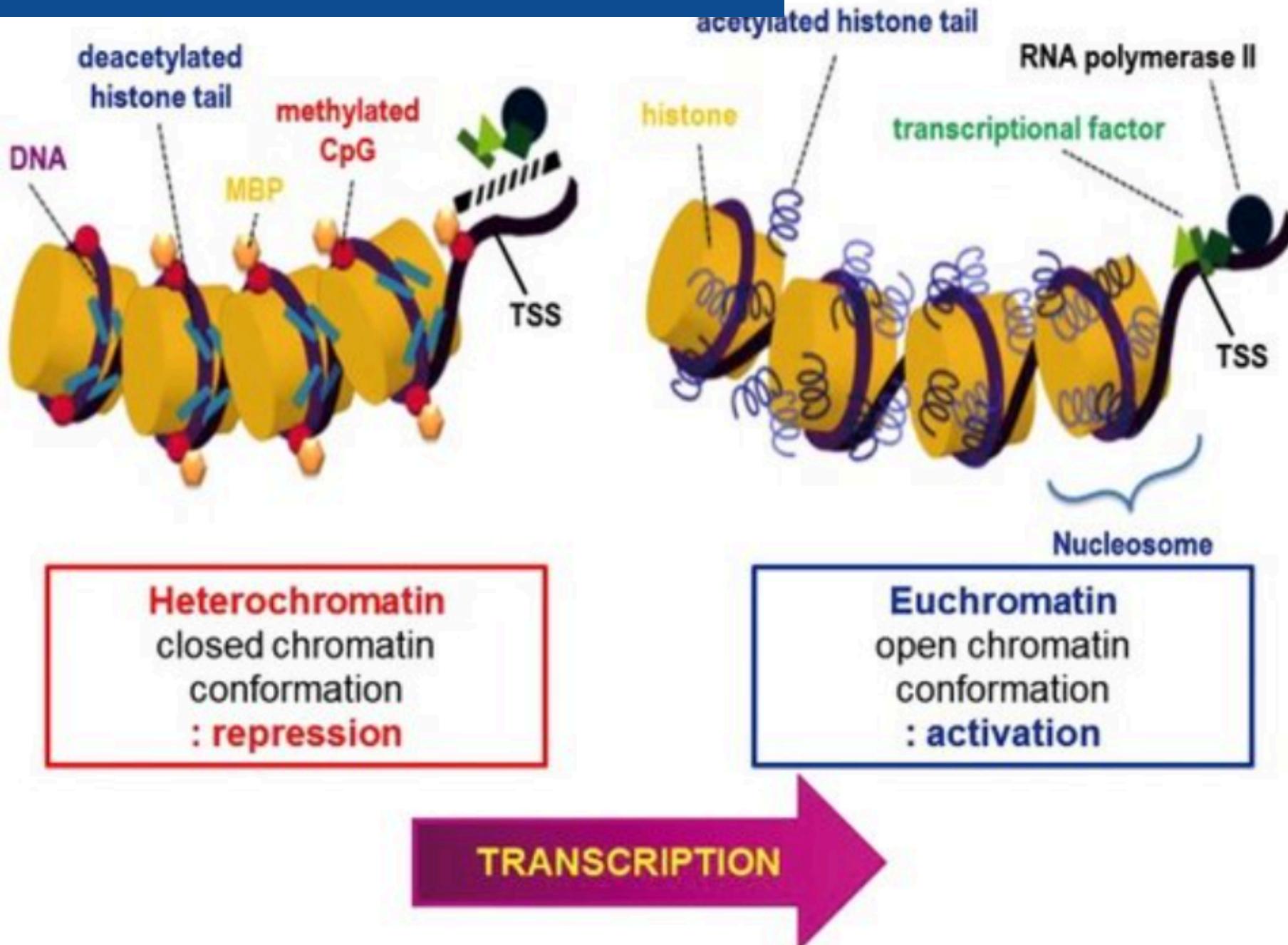
Those modifications can be influenced both by genetic and environmental factors, can be long lasting, and are sometimes heritable. While the role of epigenetic modifications as mediators of transgenerational environmental effects remains controversial

Their importance in biological processes and disease development is evident from many epigenome-wide association studies that have been reported. For example, differentially methylated regions of DNA can be used as indicators of disease status for metabolic syndrome (GenomeBiol. 2013;14:R115; Cell Metab. 2015;21:138–49), cardiovascular disease (PLoS One. 2010;5:e9692), cancer (Hum Mol Genet. 2001;10:687–92), and many other pathophysiologic states

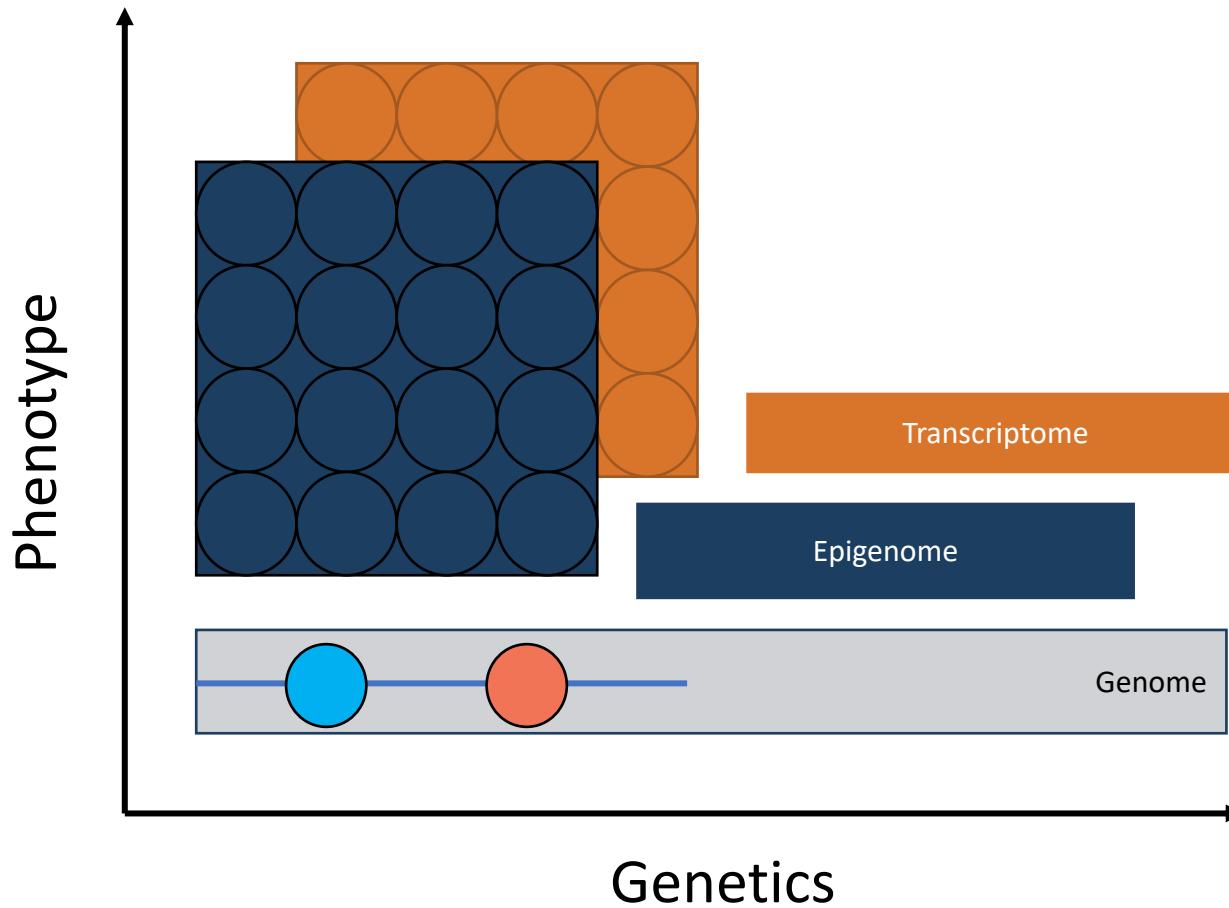
Epigenetic signatures are often tissue-specific (Cell. 2013;152:642–54), and several large consortia are focusing on establishing comprehensive epigenomic maps in multiple human tissues (Roadmap Epigenomics (<http://www.roadmapepigenomics.org/>) and International Human Epigenome Consortium (<http://ihec-epigenomes.org/>))

Data generated by these studies has great potential to enhance our functional interpretation of genetic variants residing in those regions or of epigenetic markers associated with disease independently of genetic variation

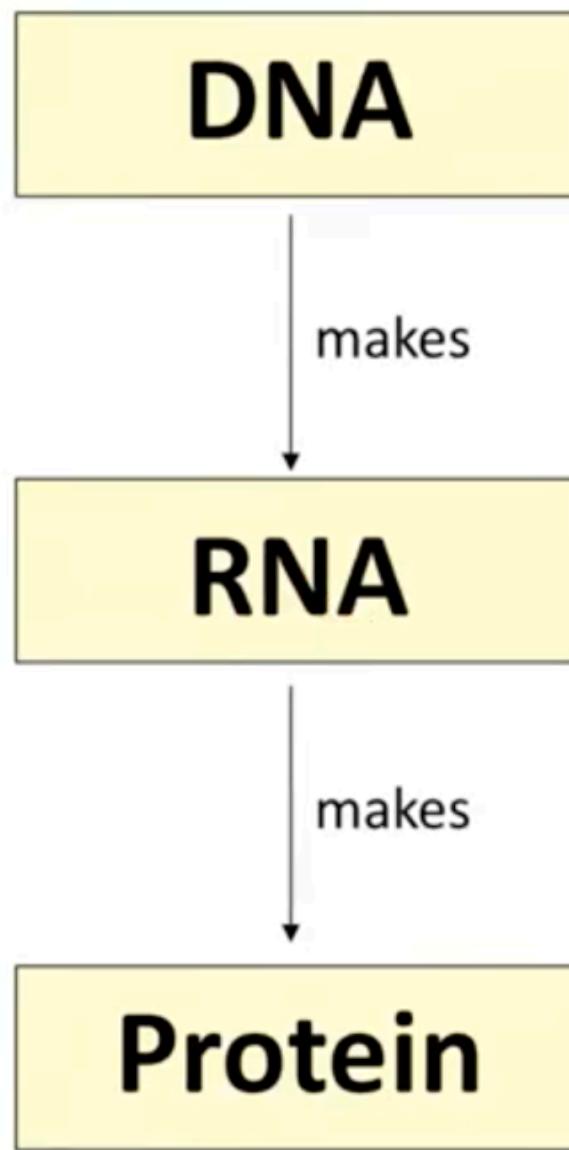
Relationship between epigenetics and gene transcription



Transcriptomics



Central Dogma



Information storage and
cross-generation transmission

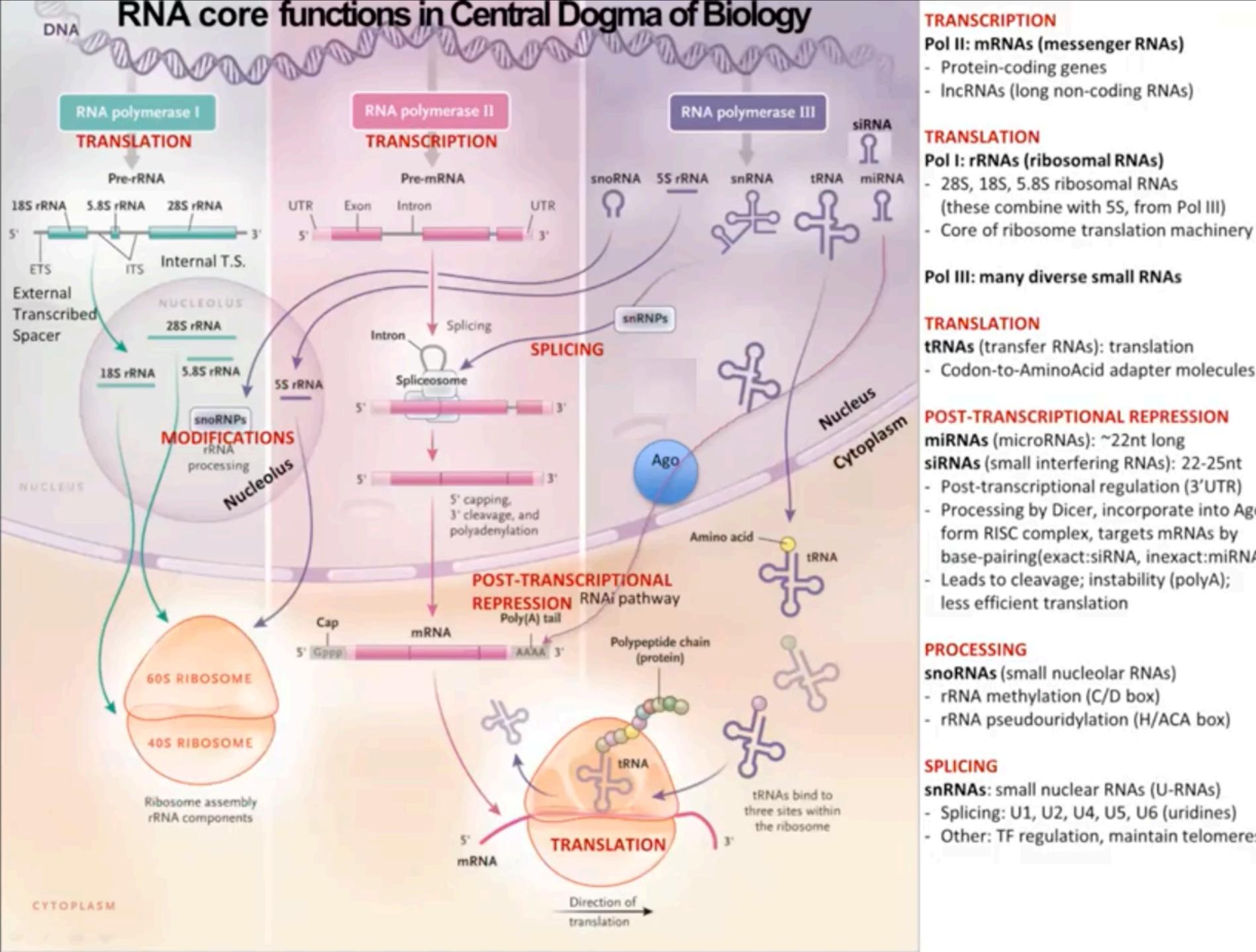
RNA can store information

Temporary copy (boring)

RNA can transmit information

Carry out all cellular functions

RNA can carry out enzymatic reactions



TRANSCRIPTION

Pol II: mRNAs (messenger RNAs)

- Protein-coding genes
- lncRNAs (long non-coding RNAs)

TRANSLATION

Pol I: rRNAs (ribosomal RNAs)

- 28S, 18S, 5.8S ribosomal RNAs
(these combine with 5S, from Pol III)
- Core of ribosome translation machinery

Pol III: many diverse small RNAs

TRANSLATION

tRNAs (transfer RNAs): translation

- Codon-to-Amino Acid adapter molecules

POST-TRANSCRIPTIONAL REPRESSION

miRNAs (microRNAs): ~22nt long

siRNAs (small interfering RNAs): 22-25nt

- Post-transcriptional regulation (3'UTR)
- Processing by Dicer, incorporate into Ago, form RISC complex, targets mRNAs by base-pairing(exact:siRNA, inexact:miRNA)
- Leads to cleavage; instability (polyA); less efficient translation

PROCESSING

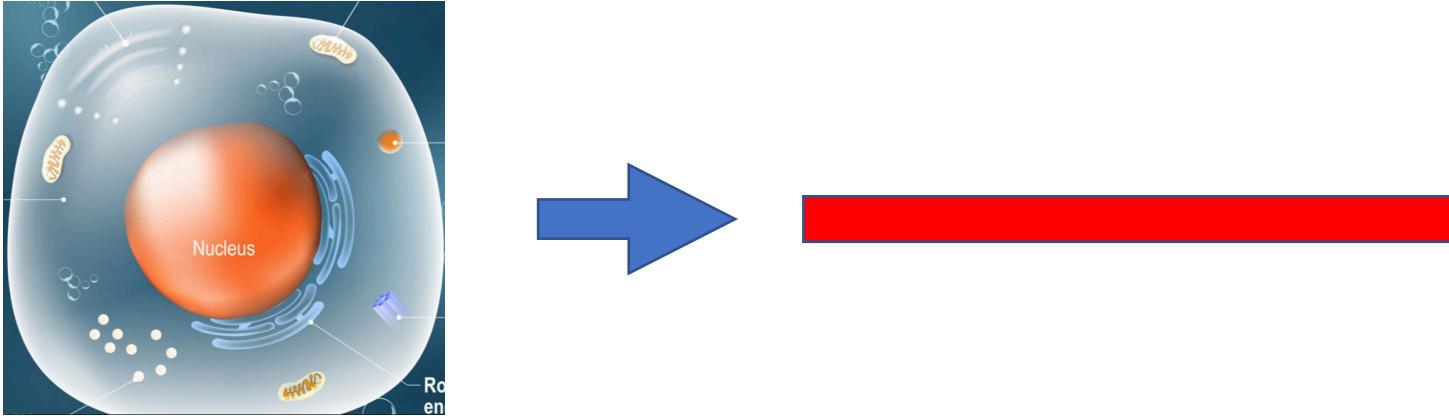
snoRNAs (small nucleolar RNAs)

- rRNA methylation (C/D box)
- rRNA pseudouridylation (H/ACA box)

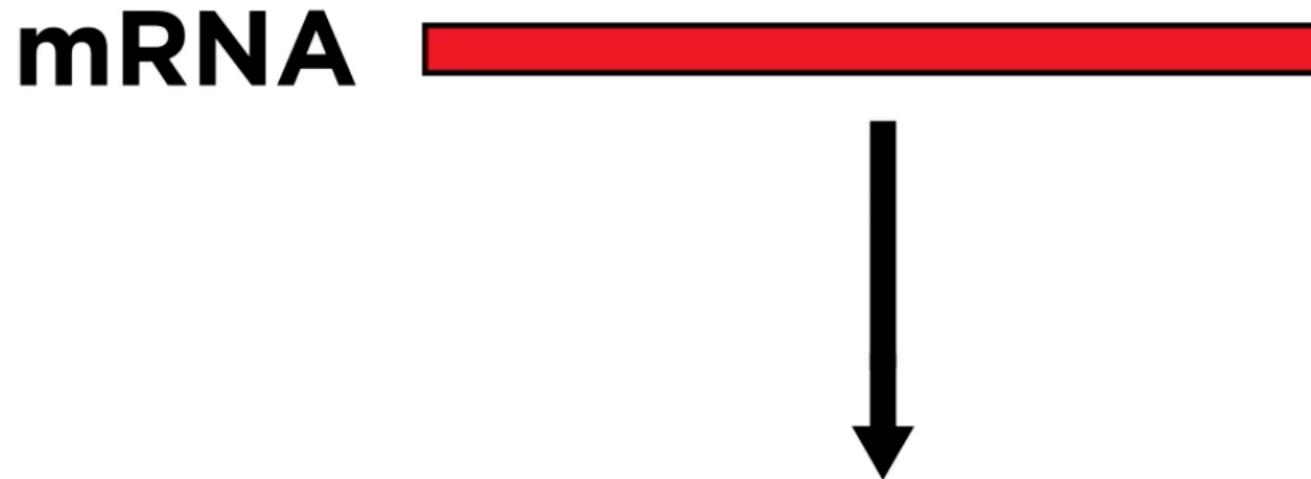
SPLICING

snRNAs: small nuclear RNAs (U-RNAs)

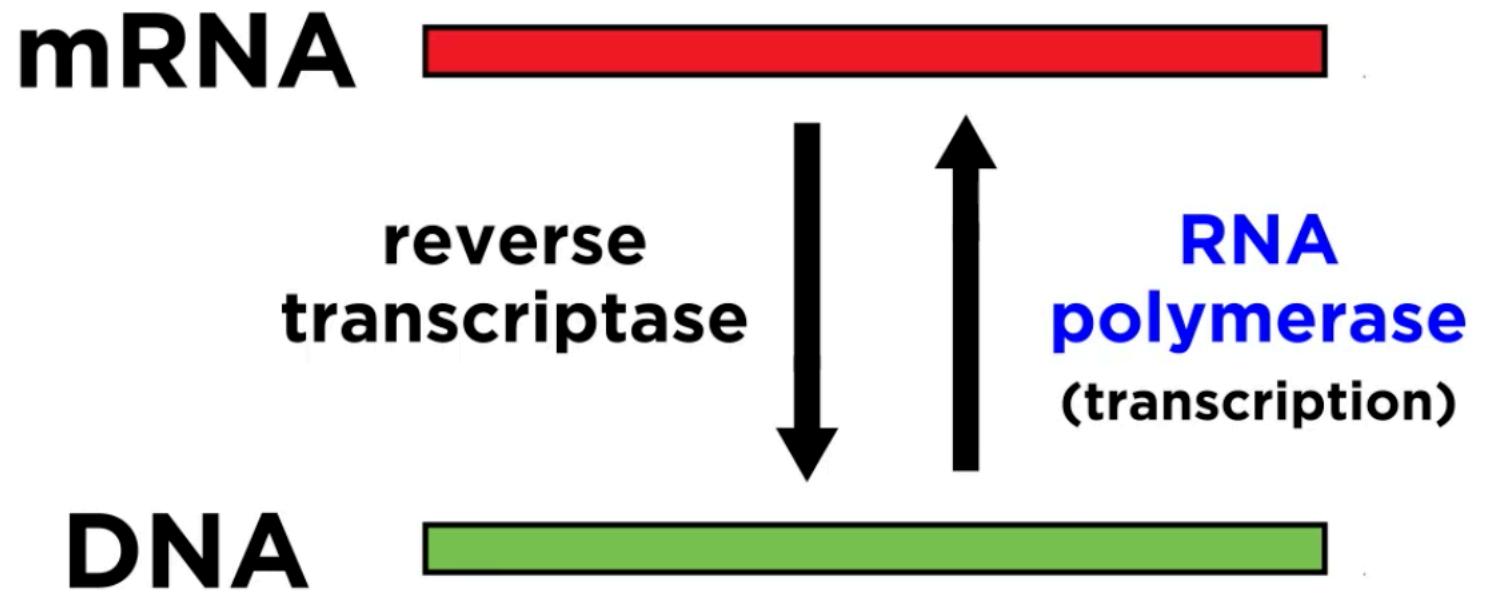
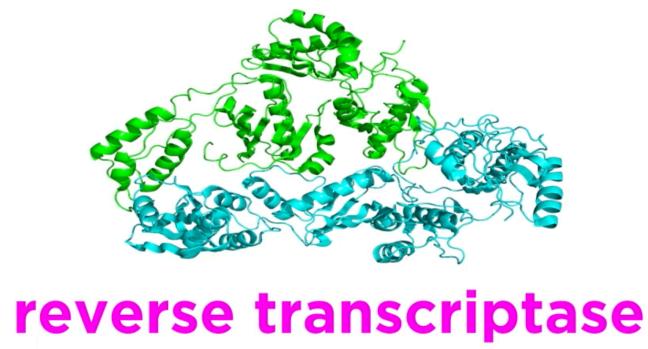
- Splicing: U1, U2, U4, U5, U6 (uridines)
- Other: TF regulation, maintain telomeres



- 1) Isolate the **mRNAs** from a cell
- 2) Performing the reverse transcriptase polymerase chain reaction (RT-PCR)



We are reverse engineering the **DNA template**
that would have generated the mRNA



mRNA AAA
 |
 |
 | template

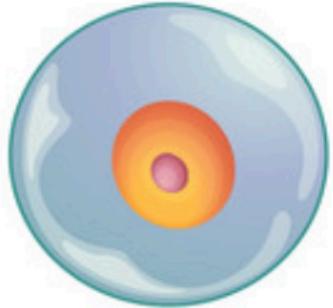
reverse
transcriptase
↓
 |
 | TTT
 | primer

DNA
primer |
 |
 |
 | →

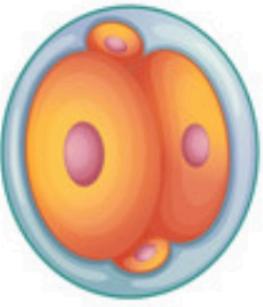
DNA
polymerase
↓

complementary DNA
(cDNA)

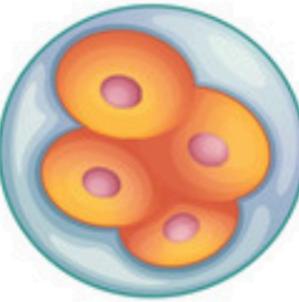
When a gene is being expressed during the embryonic development



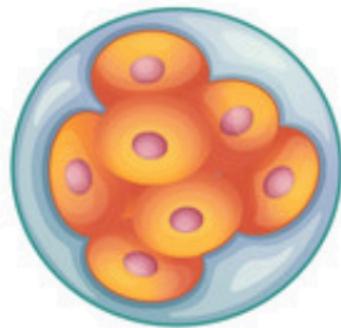
Fertilized egg



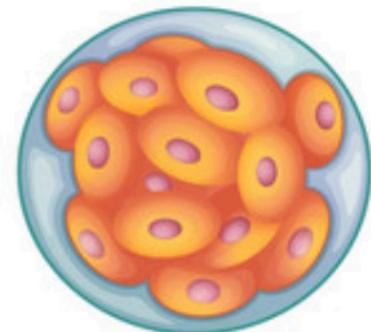
2-cell stage



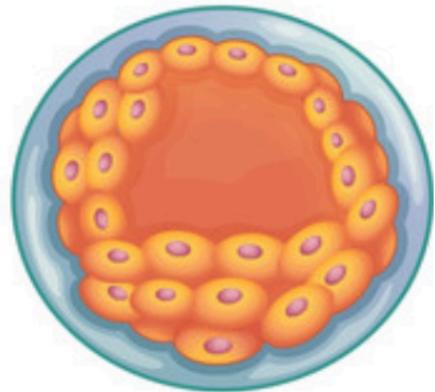
4-cell stage



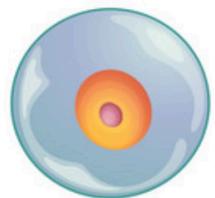
8-cell stage



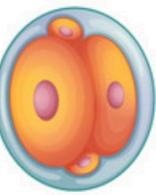
16-cell stage



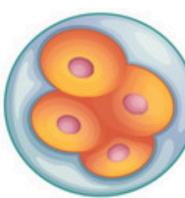
Blastocyst



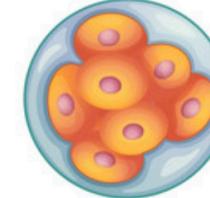
Fertilized egg



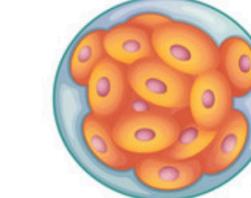
2-cell stage



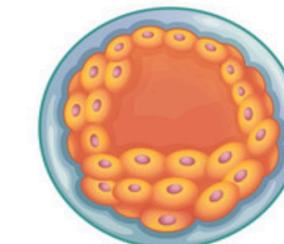
4-cell stage



8-cell stage



16-cell stage



Blastocyst



mRNA

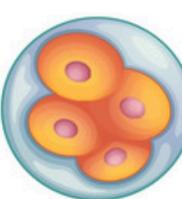




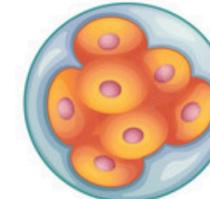
Fertilized egg



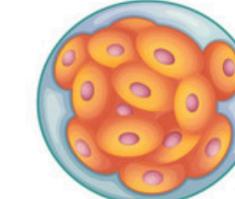
2-cell stage



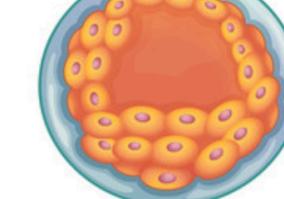
4-cell stage



8-cell stage



16-cell stage



Blastocyst



mRNA

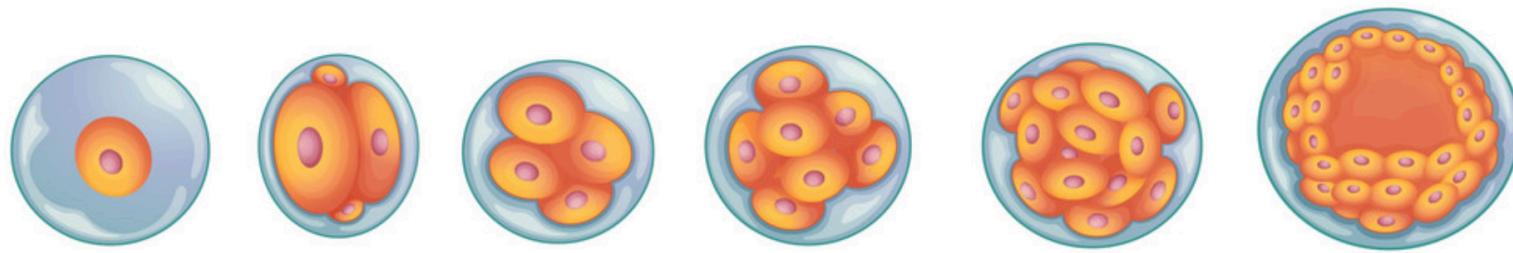


RT-PCR



cDNAs





Fertilized egg

2-cell stage

4-cell stage

8-cell stage

16-cell stage

Blastocyst

mRNA

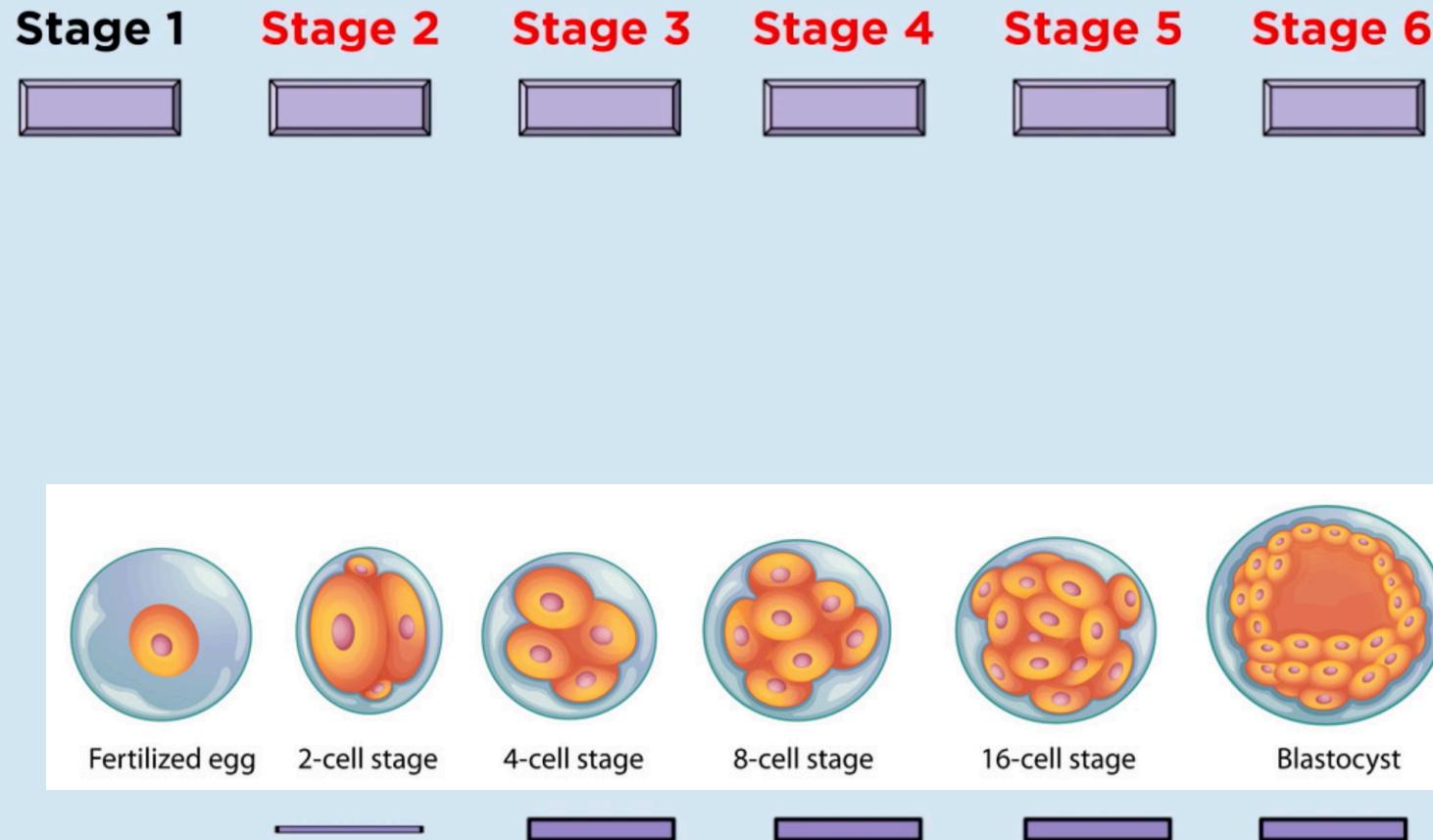
RT-PCR

cDNAs

The primer used for
amplification is
specific to this gene

Many copies of this cDNA molecule

Length of
DNA
fragment
(base pairs)



**this is how we know when
a gene is being expressed**

Gel Electrophoresis

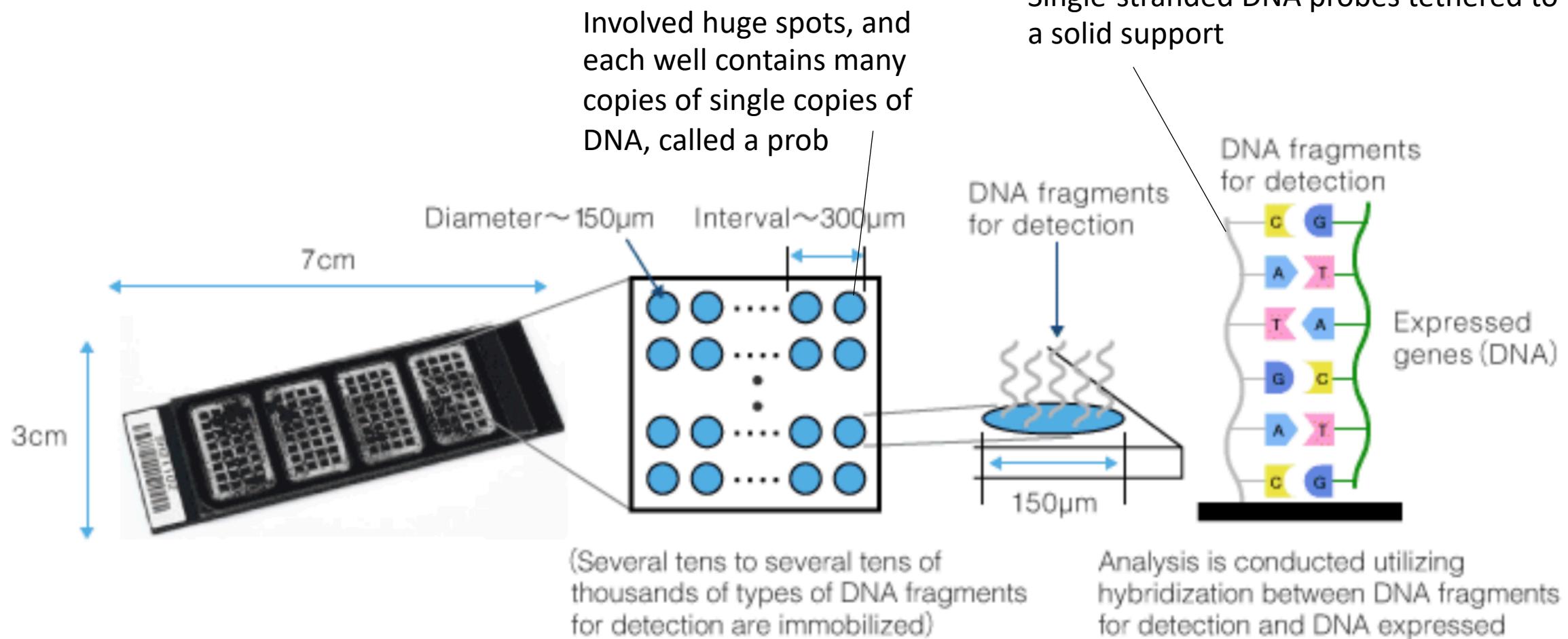


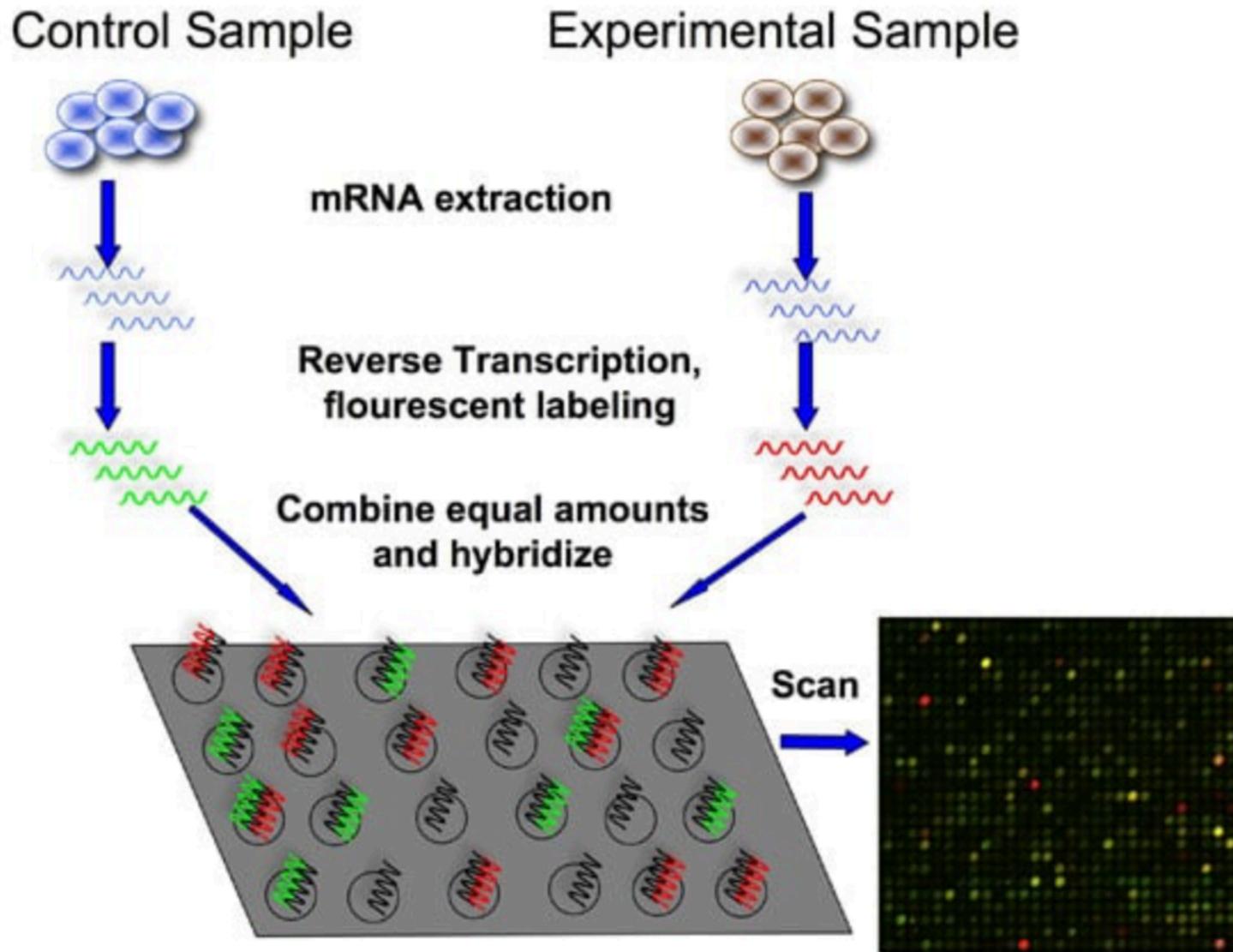
cDNA that
represents the
gene of interest

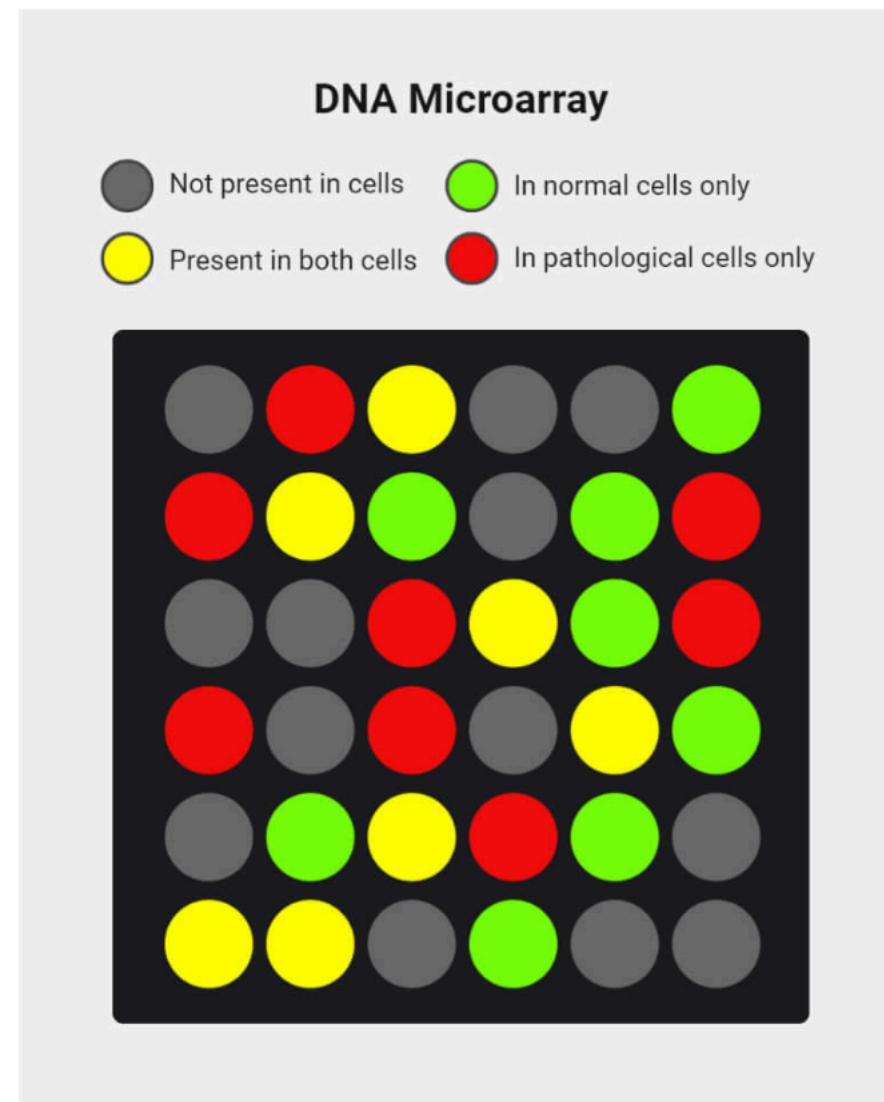
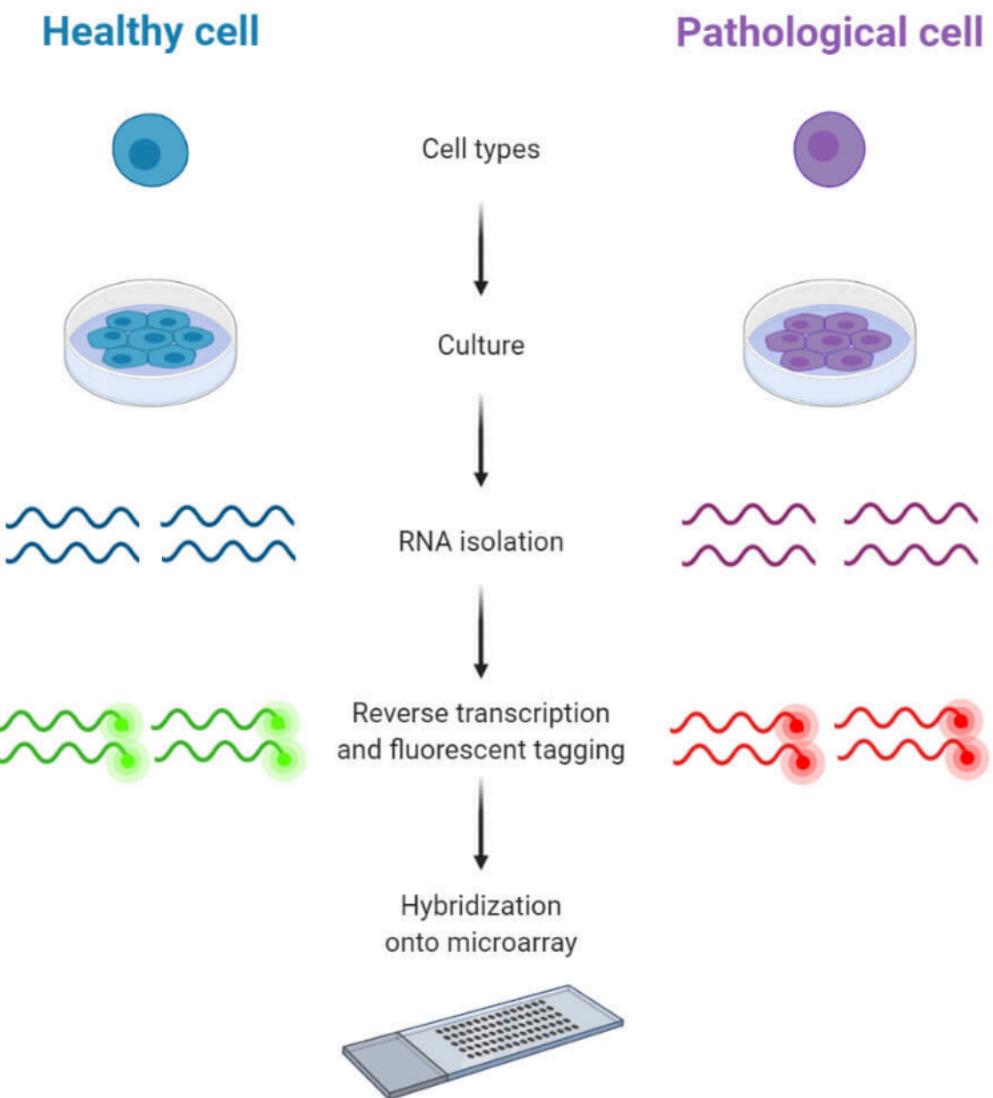


mRNA is produced
during that **stage**
of development

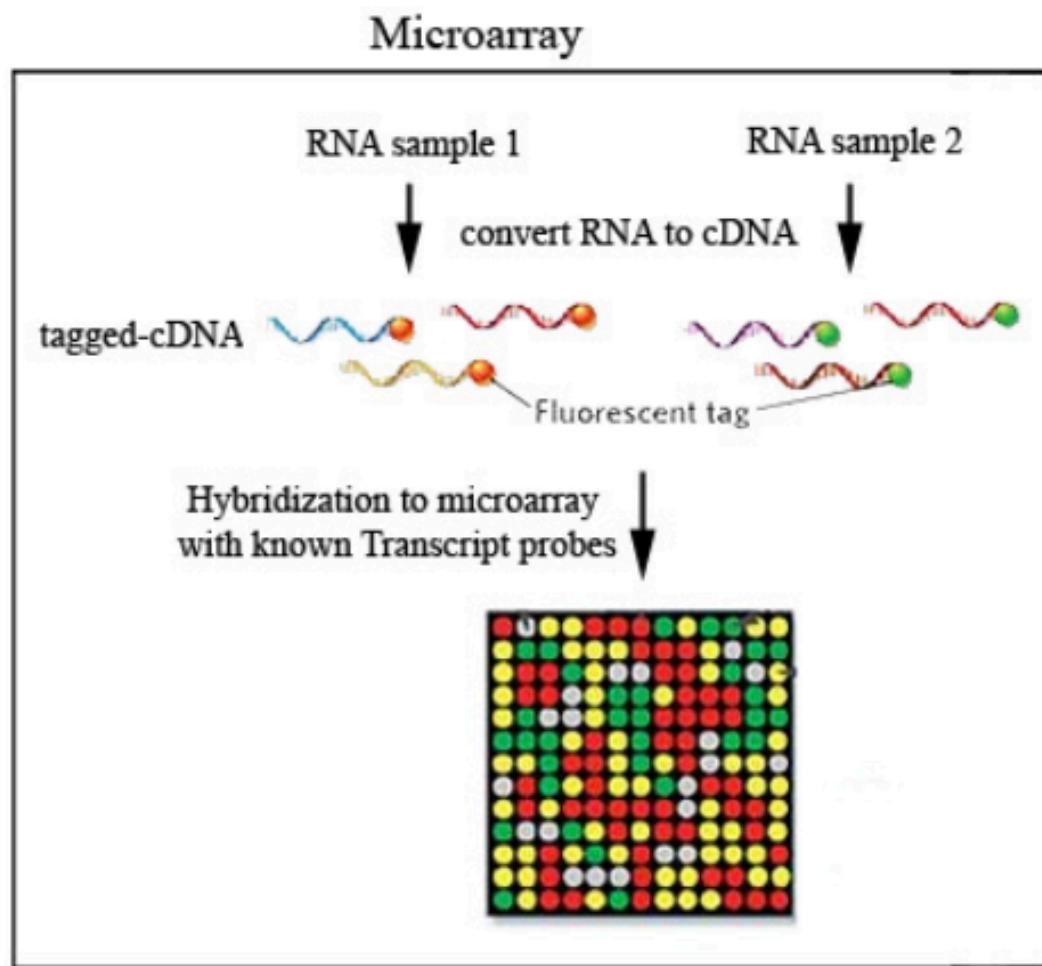
Microarray





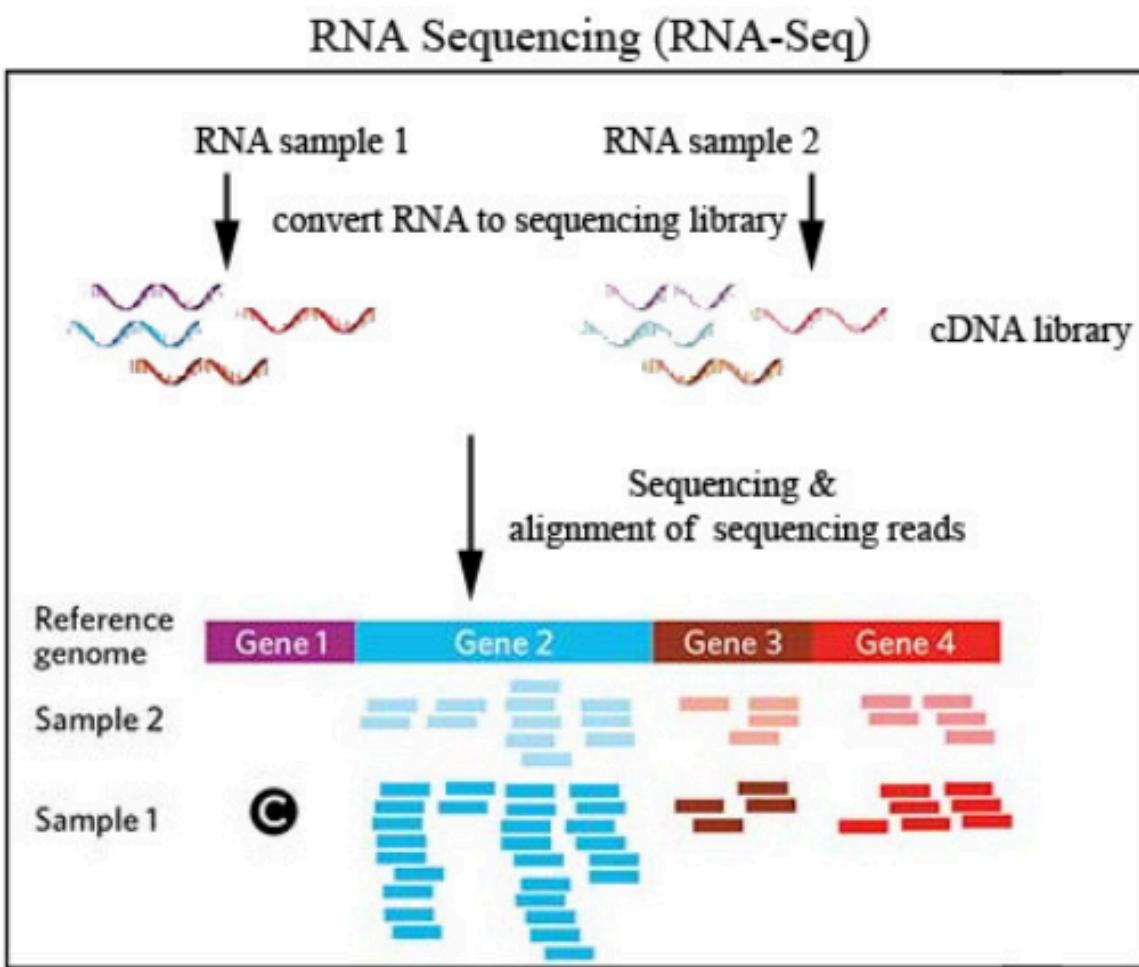


Microarray vs. RNA-Seq



relative intensity
= expression levels

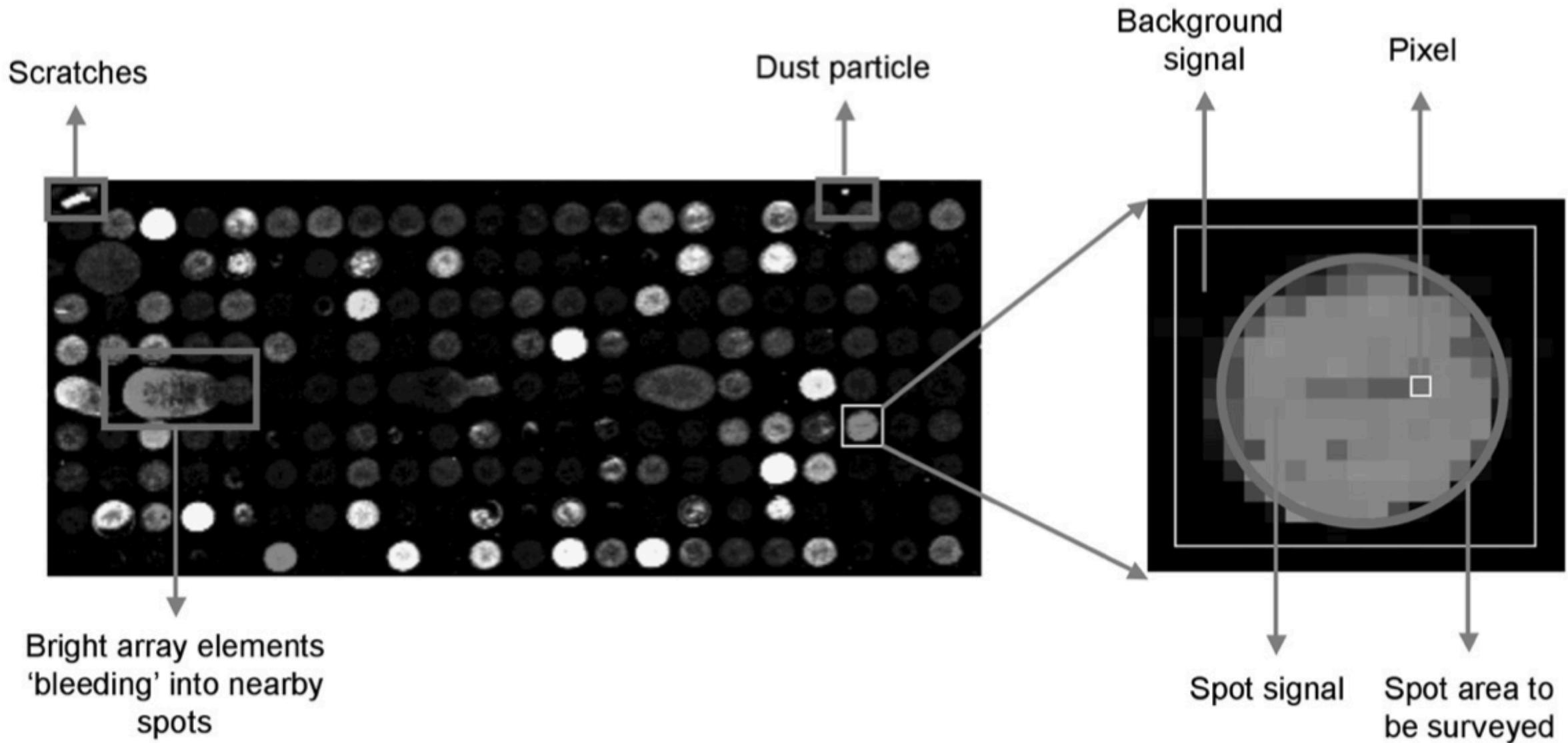
Low sensitivity
Low dynamic range
known transcript only
No alternative splicing information
lower cost



High sensitivity
High dynamic range
Novel transcripts sequences identified
structural variation & alternative splicing revealed
unlimited sample comparisons

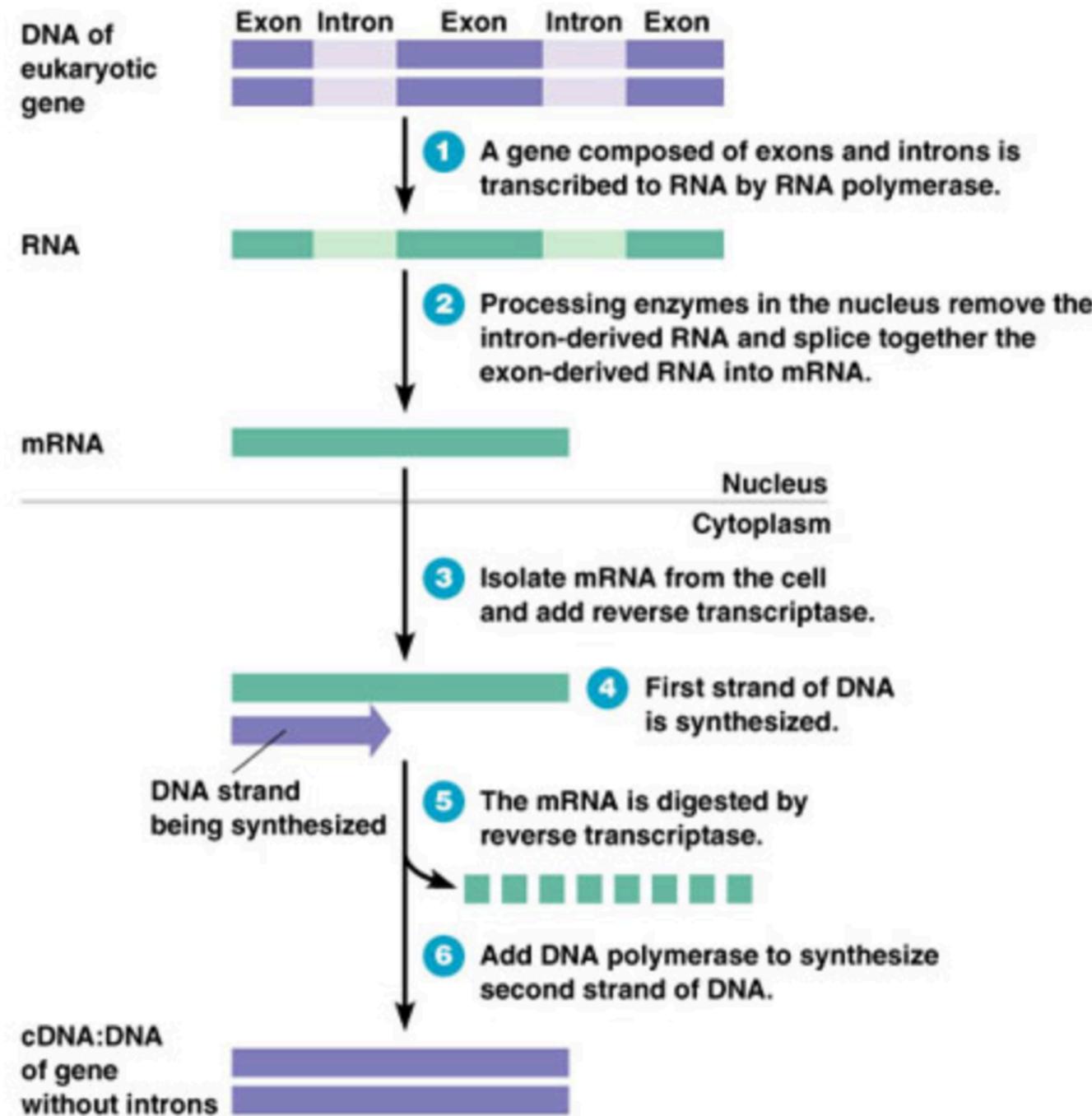
Sequencing Reads
= expression levels

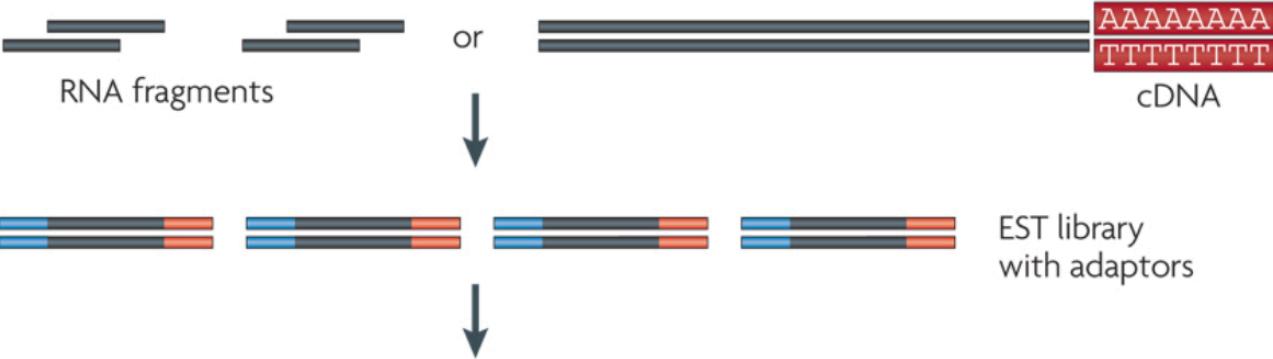
Microarray images are not perfect, as it is often the case, which leads to many problems with spurious signals from dust particles



Common RNA-Seq workflow

1. Isolate all mRNA
2. Convert to cDNA using reverse transcriptase



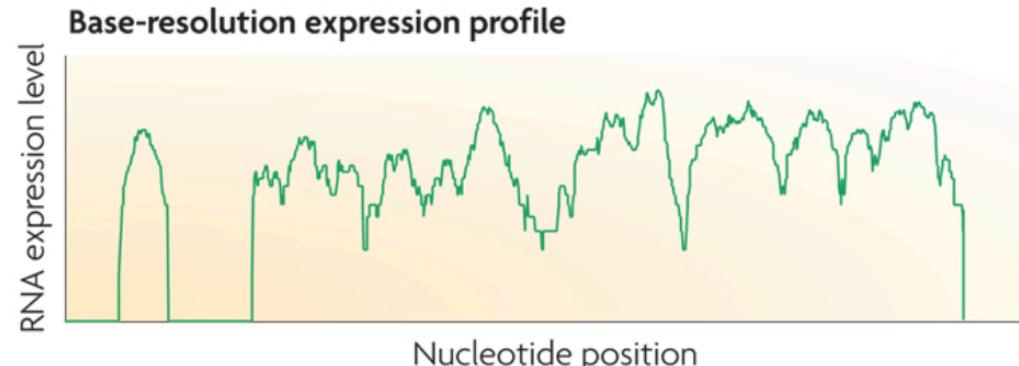
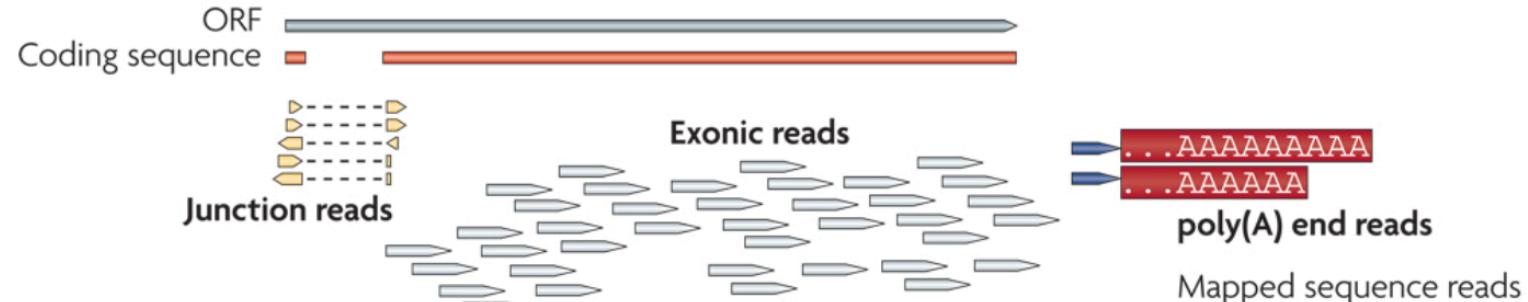


3. Sequence the cDNA

ATCACAGTGGGACTCCATAAATTTCT
CGAAGGACCAGCAGAACGAGAGAAAAA
GGACAGAGTCCCCAGCGGGCTGAAGGGG
ATGAAACATTAAAGTCAAACAATATGAA
.....

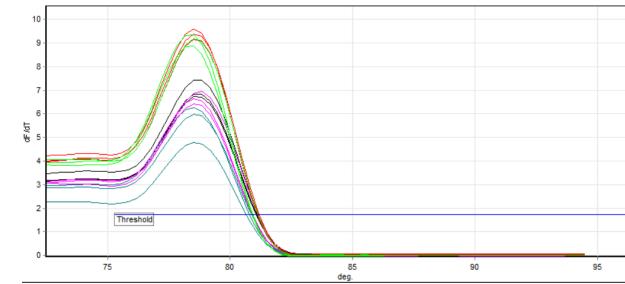
Short sequence reads

4. Map sequences to the genome



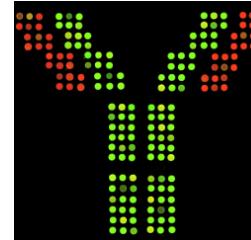
qPCR, Microarrays or RNA-sequencing: When To Choose One Over the Other?

1. qPCR



Choose qPCR for your experiment if you have to analyze the expression of a few genes with a known sequence

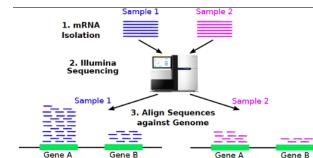
2. Microarrays



Consider going for microarrays if you are looking for an affordable and robust solution, especially if you do not know which genes you want to analyze – or for when you want to perform a whole transcriptome differentially expressed (DE) genes analysis and have a good reference sequence for your organism.

3. RNA-Seq

In comparison to microarrays, RNA-sequencing (or RNA-seq for short) enables you to look at differential expressions at a much broader dynamic range, to examine DNA variations (SNPs, insertions, deletions) and even discover new genes or alternative splice variations using just one dataset.



1. Exons are separated by introns – only <3% human genome codes exons.

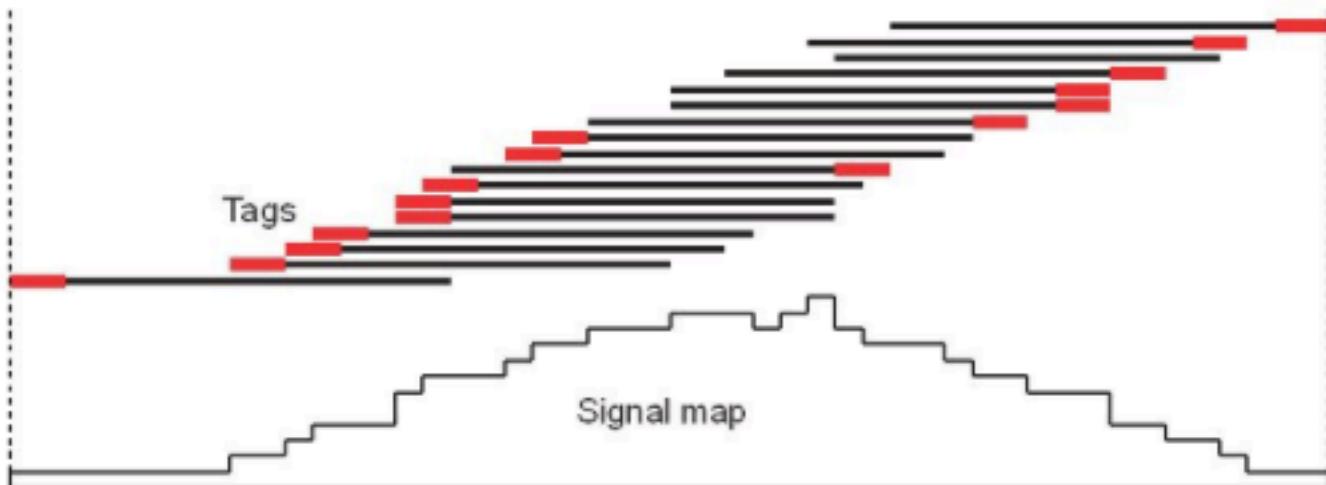
Mapping reads to the reference genome is challenging

2. The relative abundance of RNAs vary over a very large range

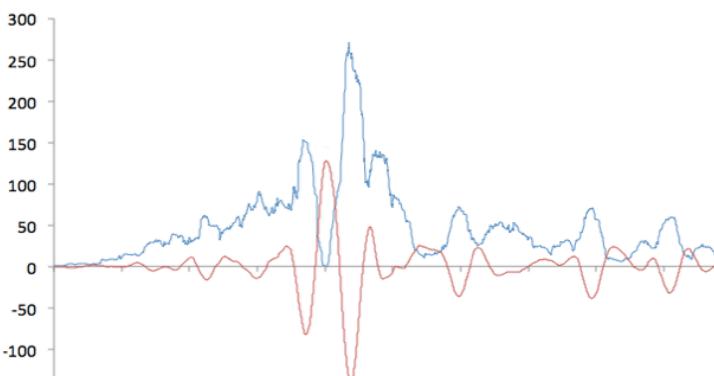
Since RNA sequencing works by random sampling, a small fraction of highly expressed genes may make up the majority of reads

3. RNAs come in a wide range of sizes

Small RNAs must be captured separately

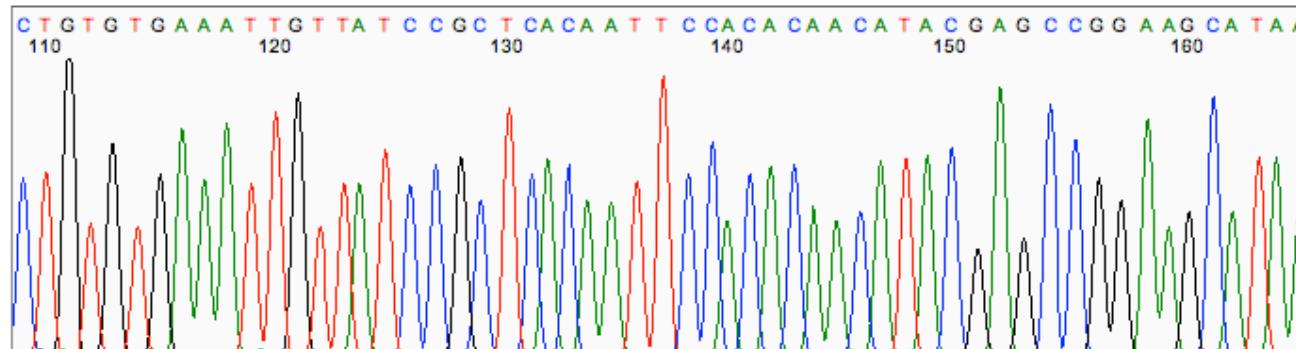


- ✓ Map reads (red) to the genome. Whole pieces of DNA are black.
- ✓ Count # of reads mapping to each DNA base → signal

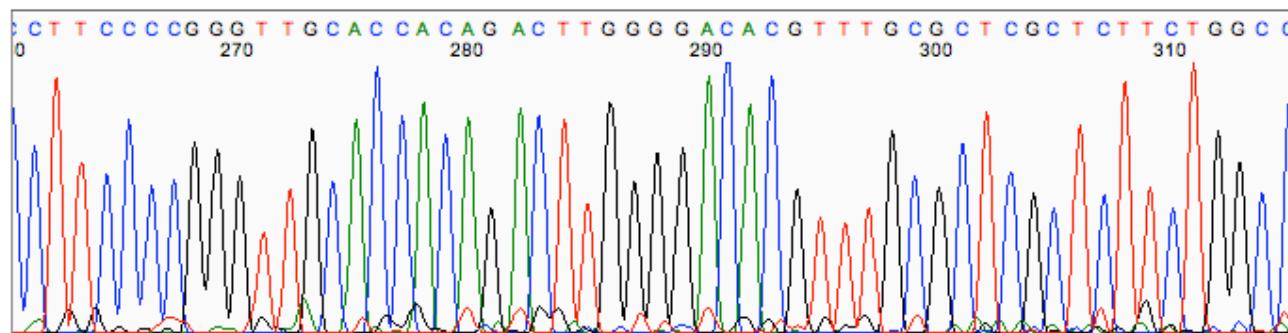


Noisy or Weak DNA Sequencing Lanes

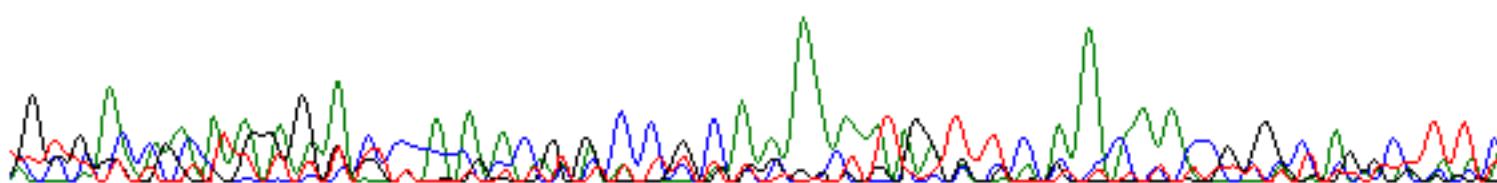
Here's an example of excellent sequence. Note the evenly-spaced peaks and the lack of baseline 'noise' (see further down for examples of higher baseline noise):



The next example has a little baseline noise, but the 'real' peaks are still easy to call, so there's no problem with this sample:



GTGA CAA AAG ACC AAA GG CCT CA AAA AAAT TC AA AA AC G G TA GC TT C
30 40 50 60 70



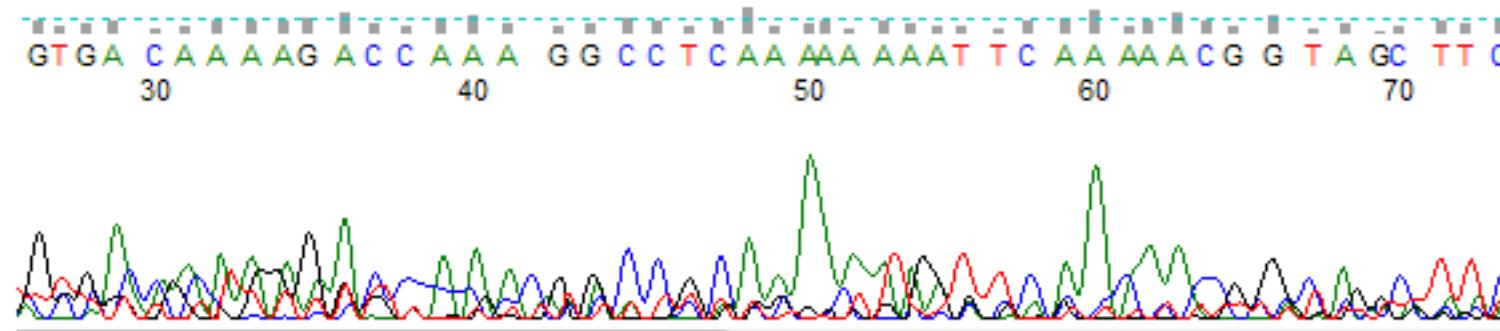
Noisy or Weak DNA Sequencing Lanes

Causes of weak signal lanes

Partially failed sequencing reaction.

Too much or too little DNA.

Partial loss of the sample during clean-up

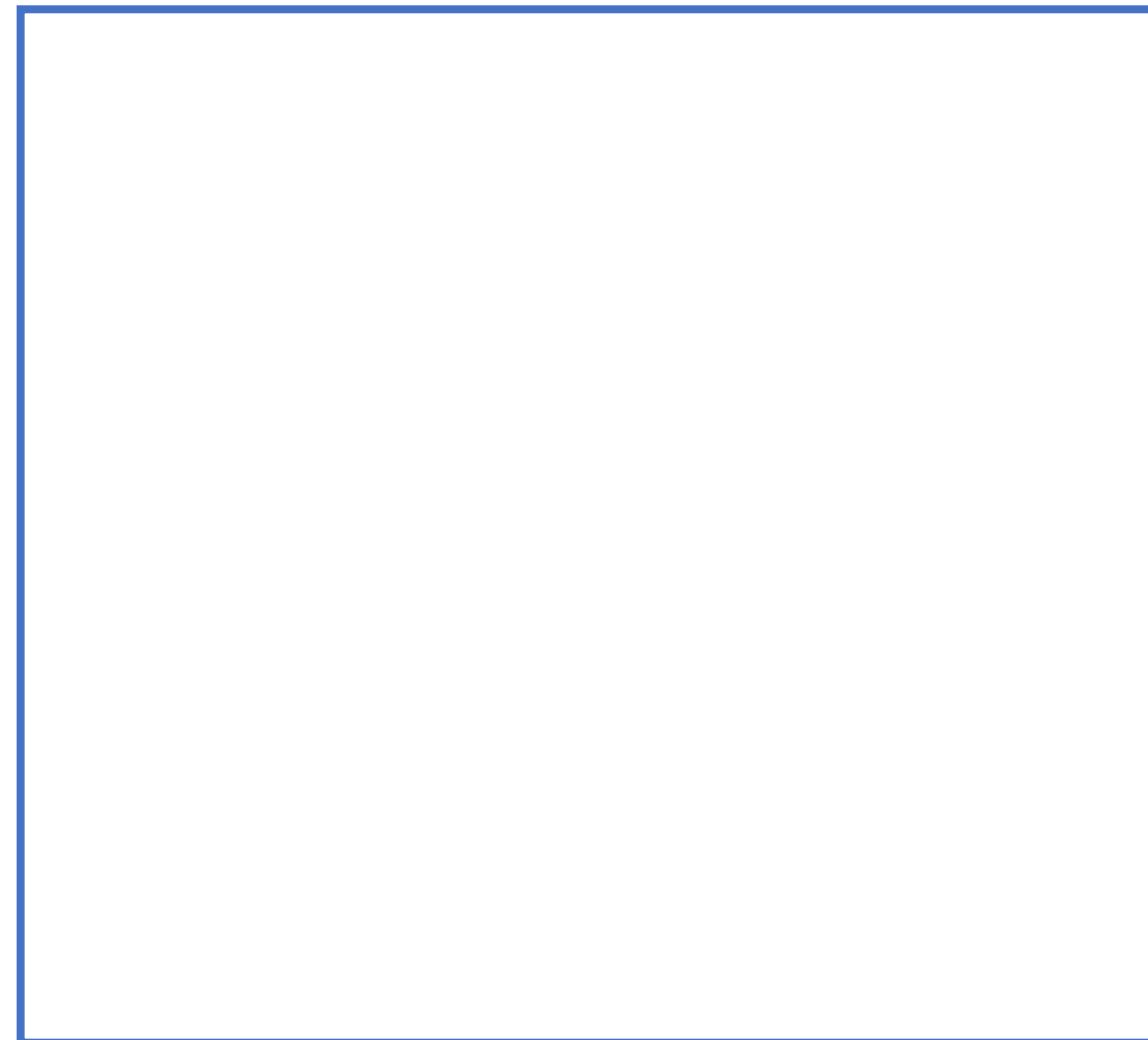


Expression analysis data matrix

Measure 20,000 genes in 100s of conditions

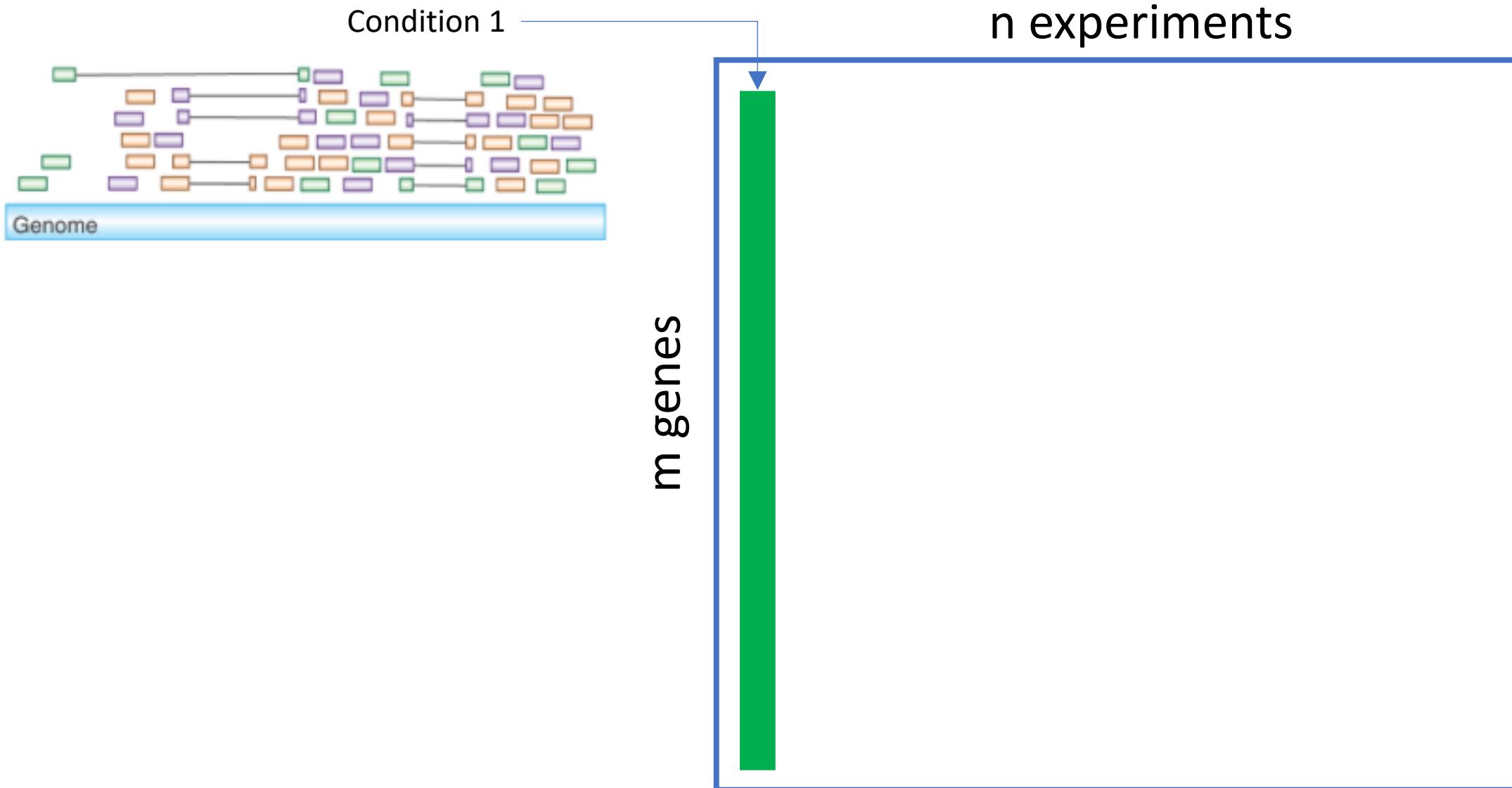
n experiments

m genes



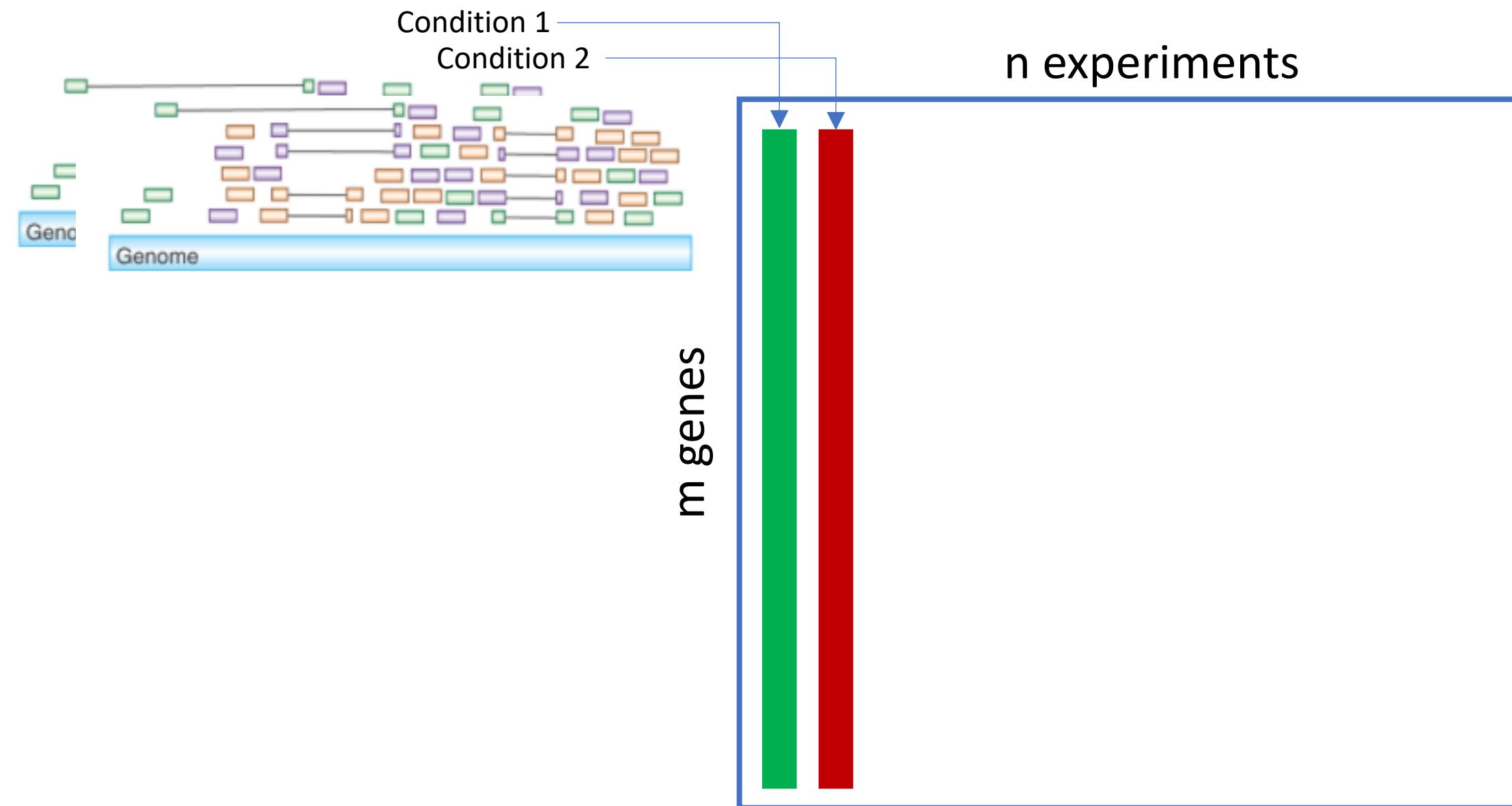
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions



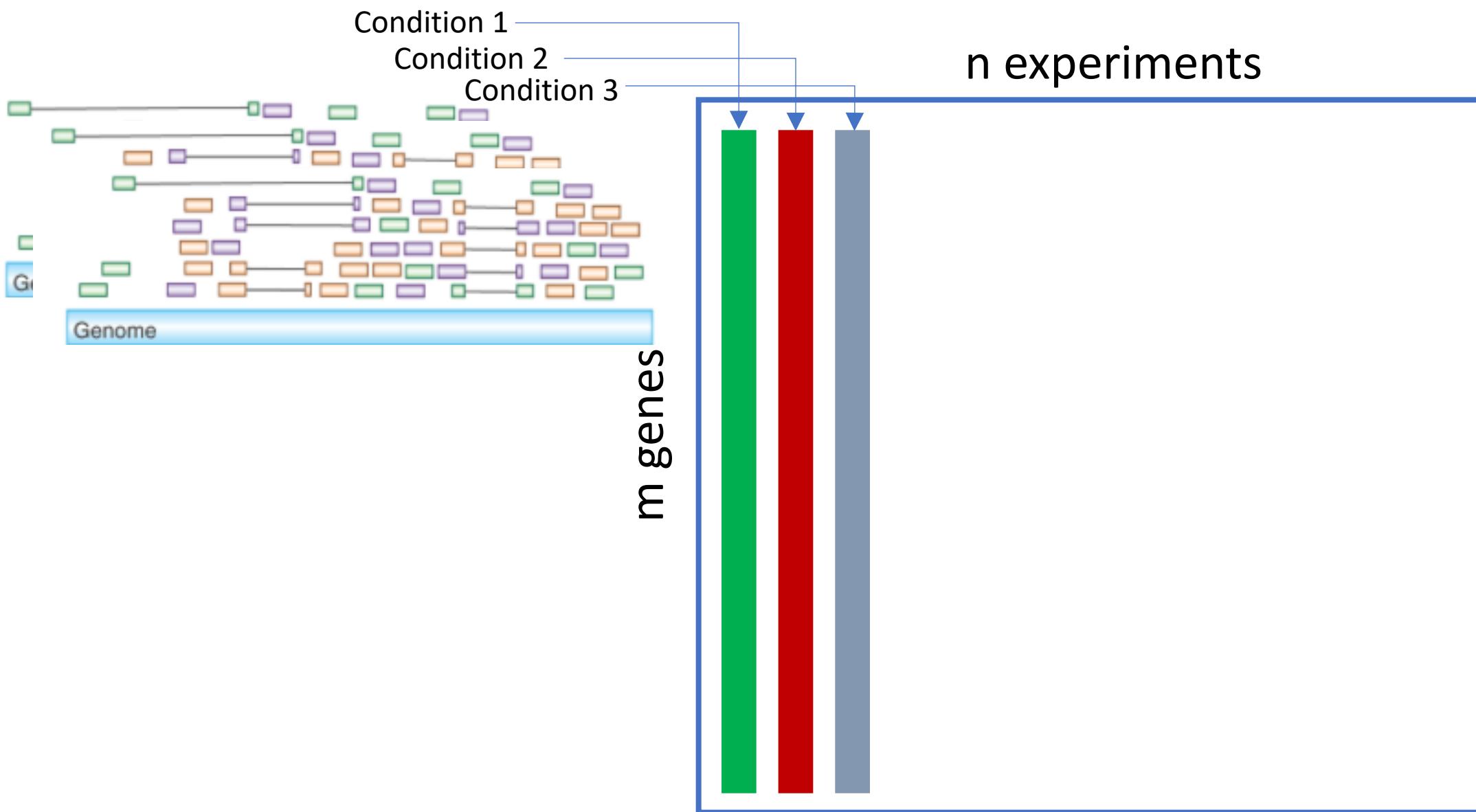
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions



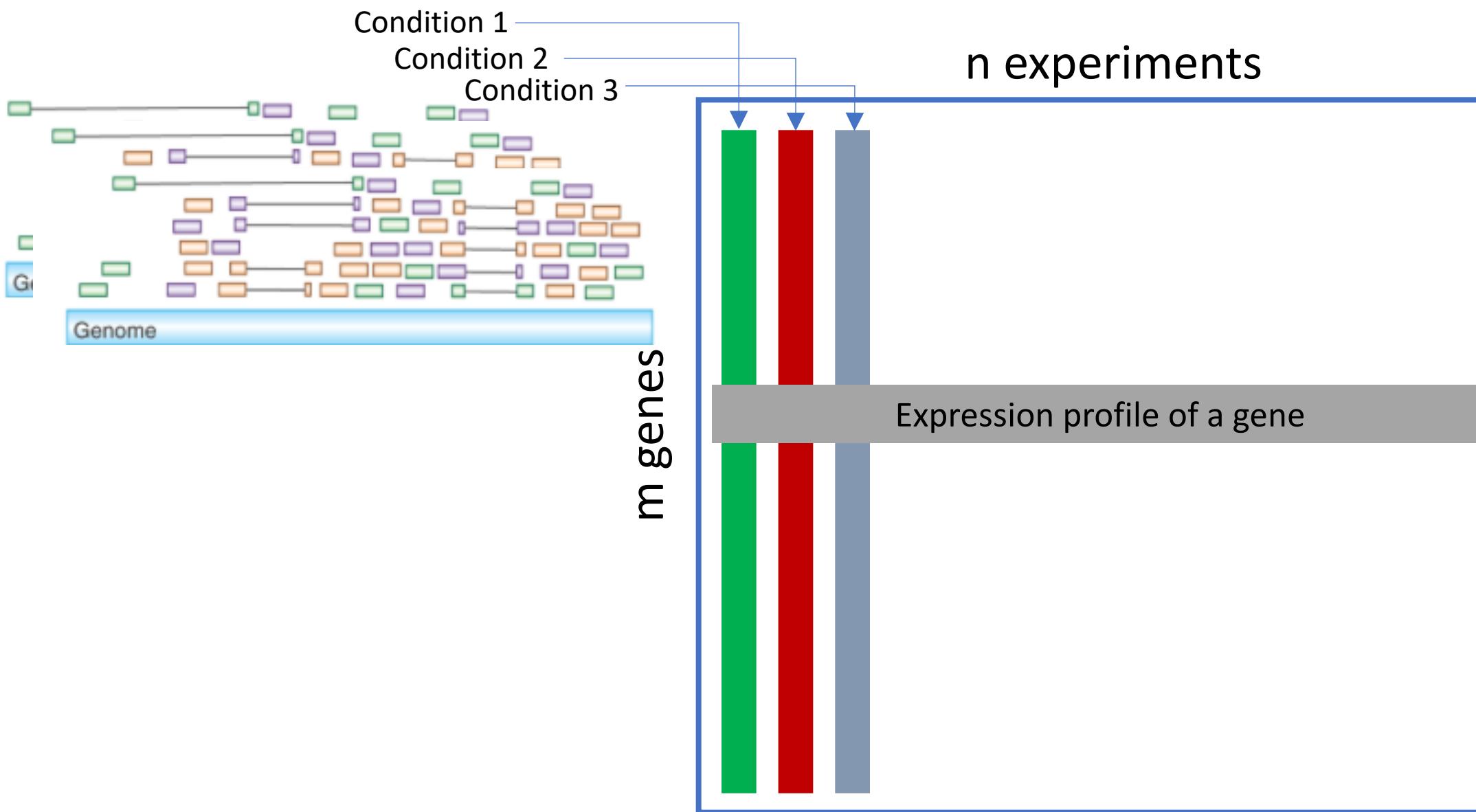
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions



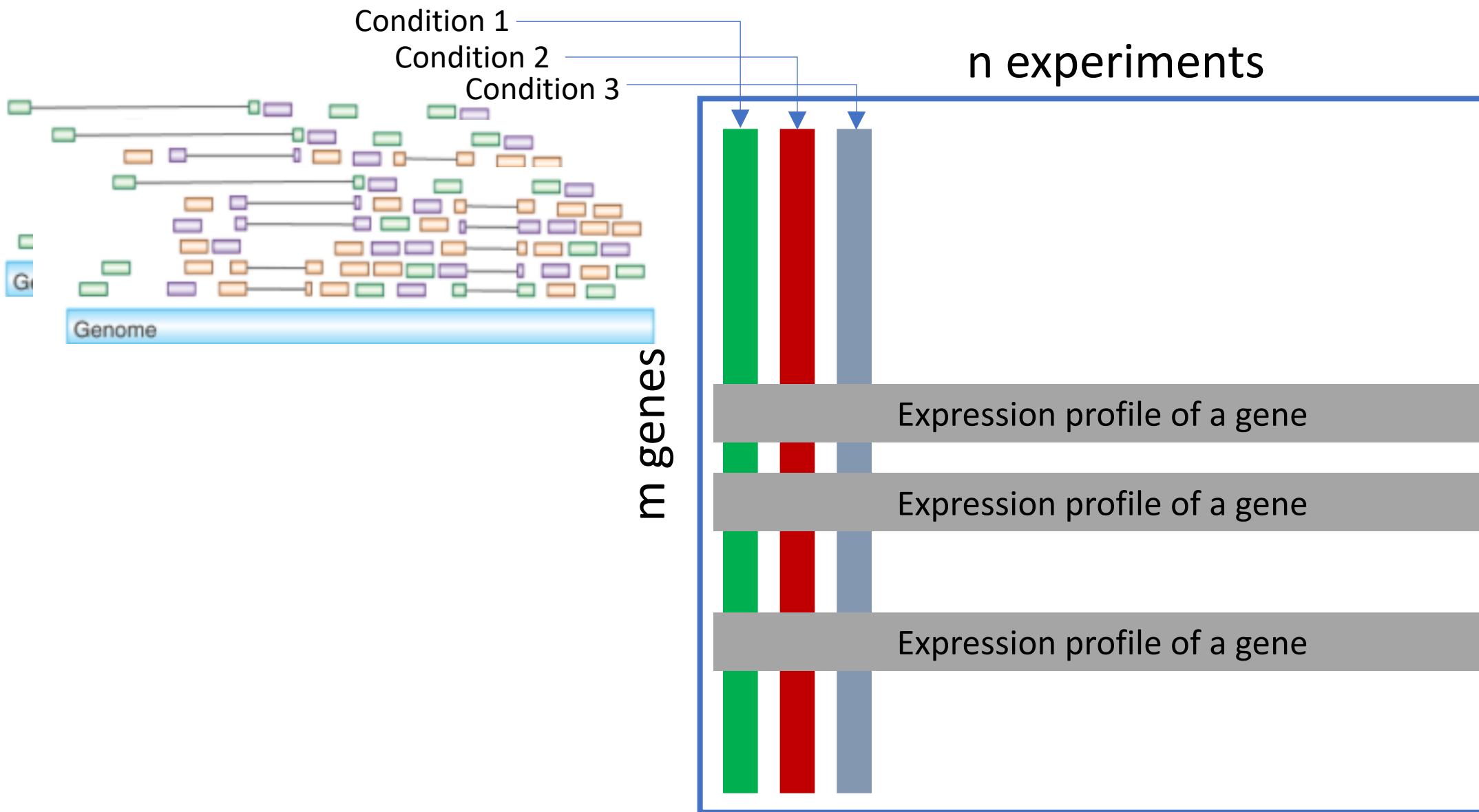
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions



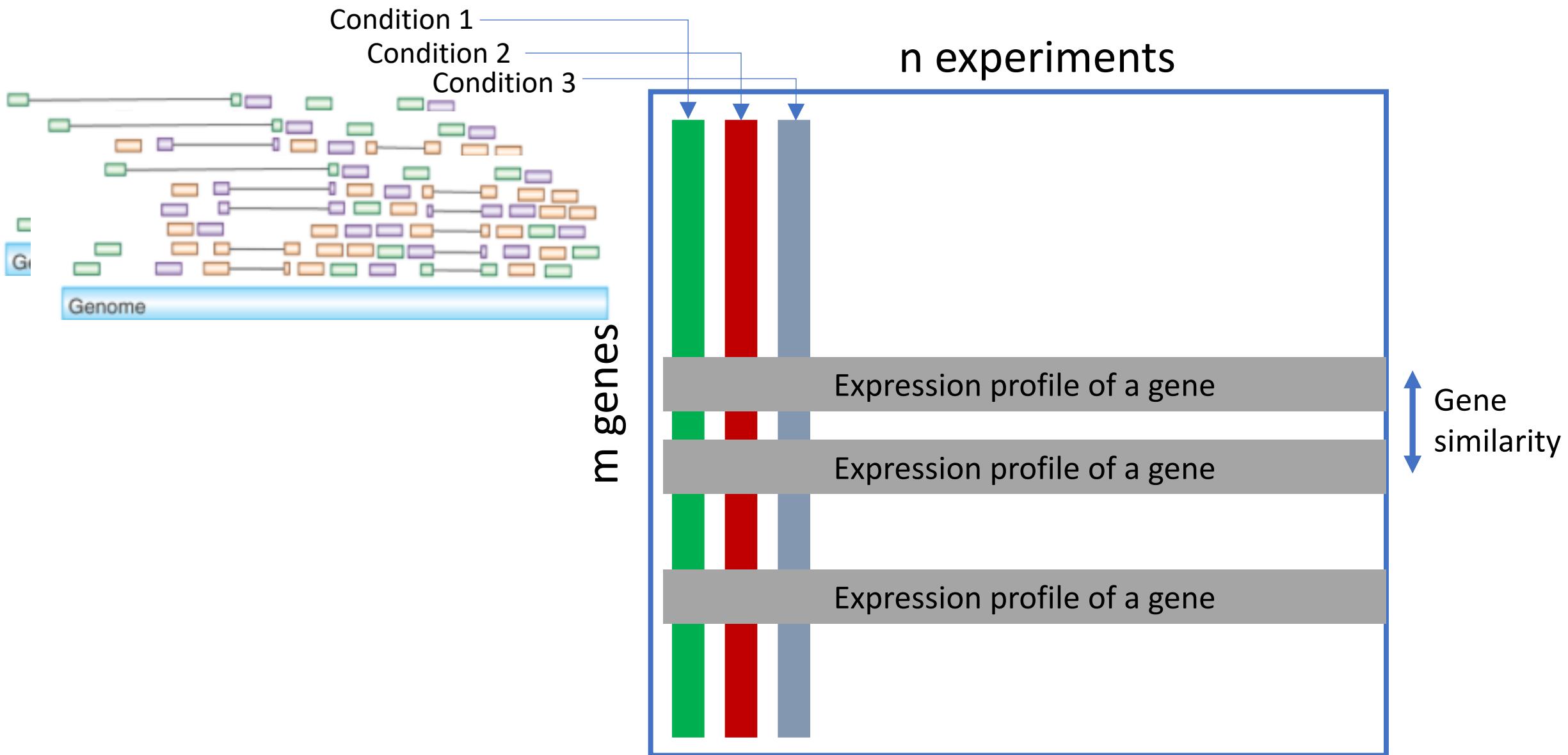
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions



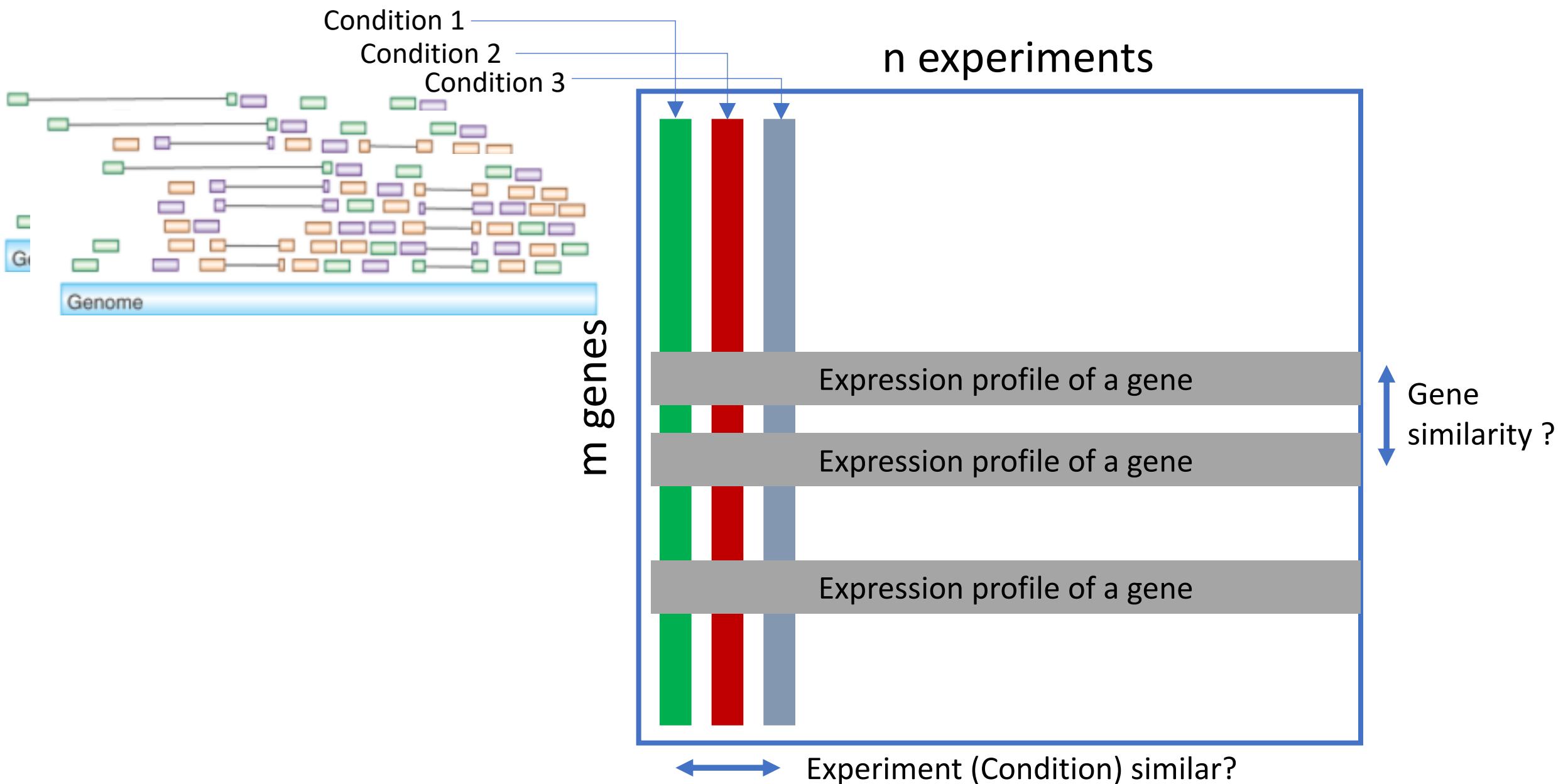
Expression analysis data matrix

Measure 20,000 genes in 100s of conditions

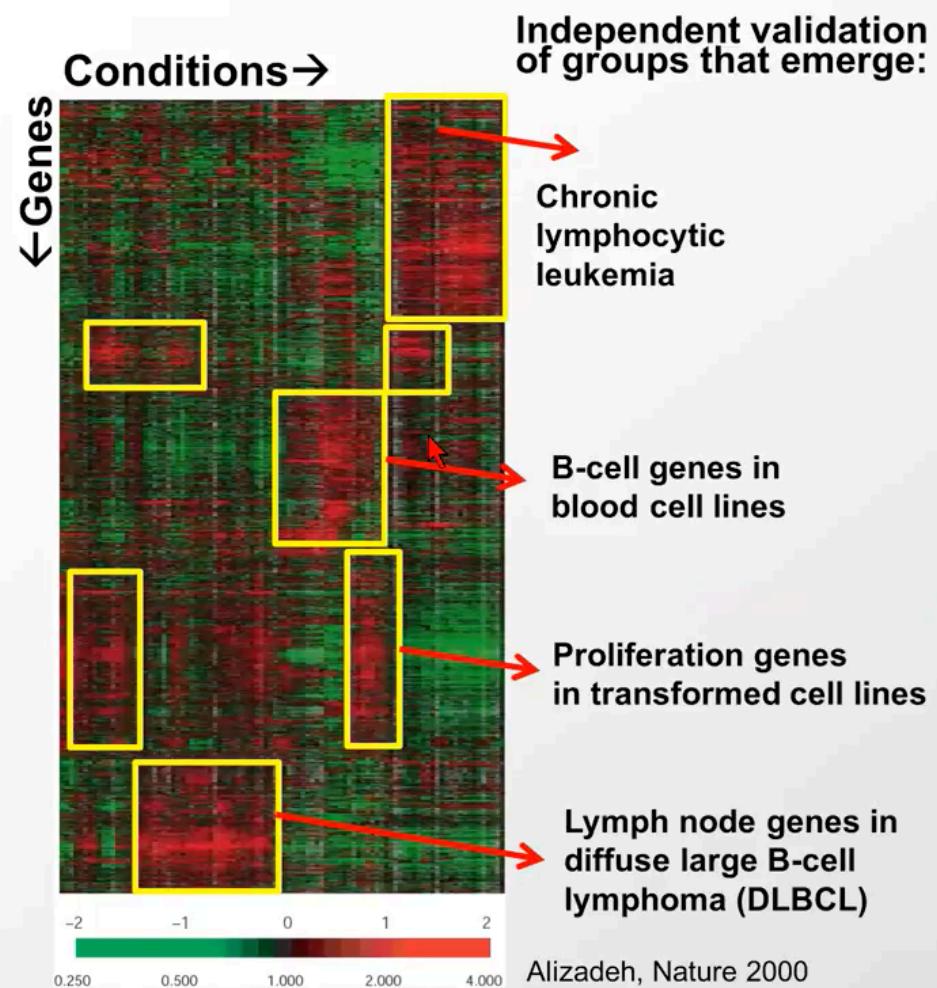


Expression analysis data matrix

Measure 20,000 genes in 100s of conditions

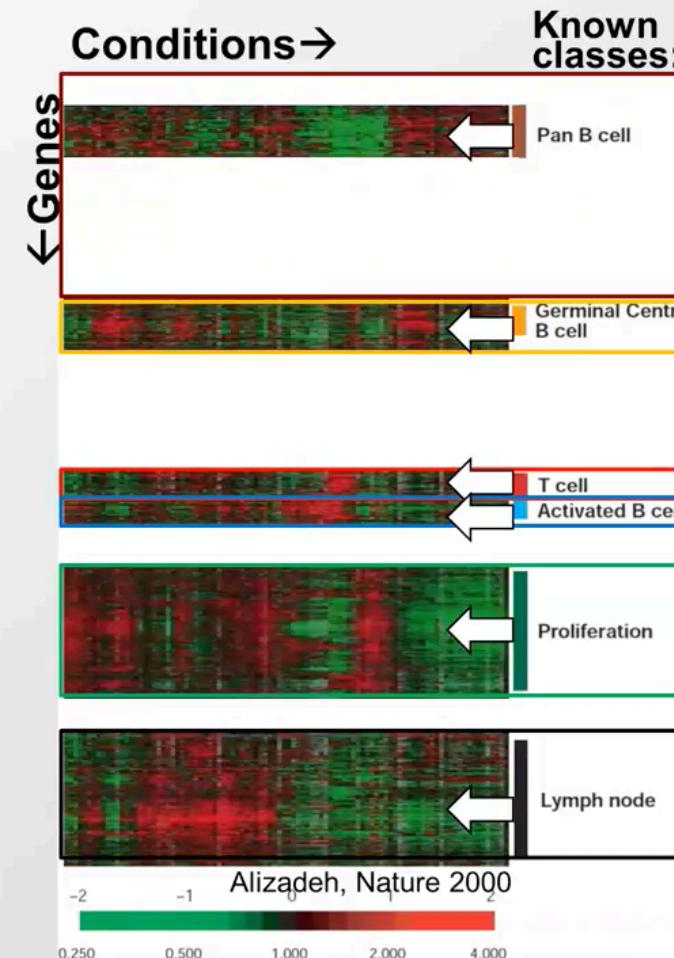


Clustering vs Classification



Goal of Clustering: Group similar items that likely come from the same category, and in doing so reveal hidden structure

- **Unsupervised learning**

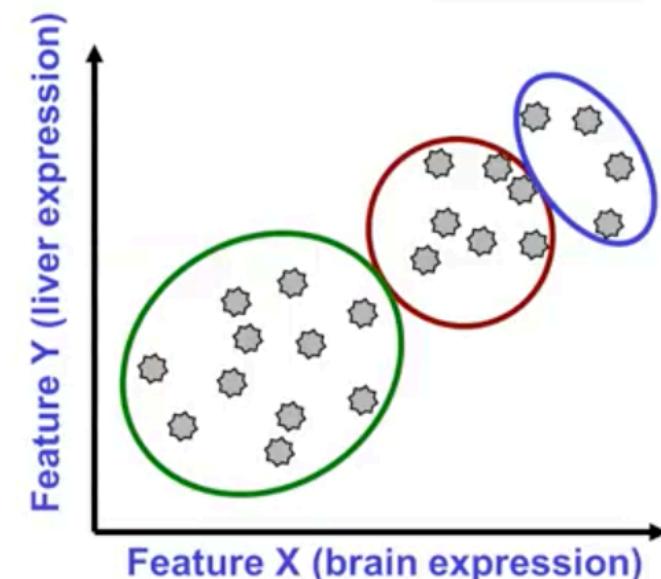
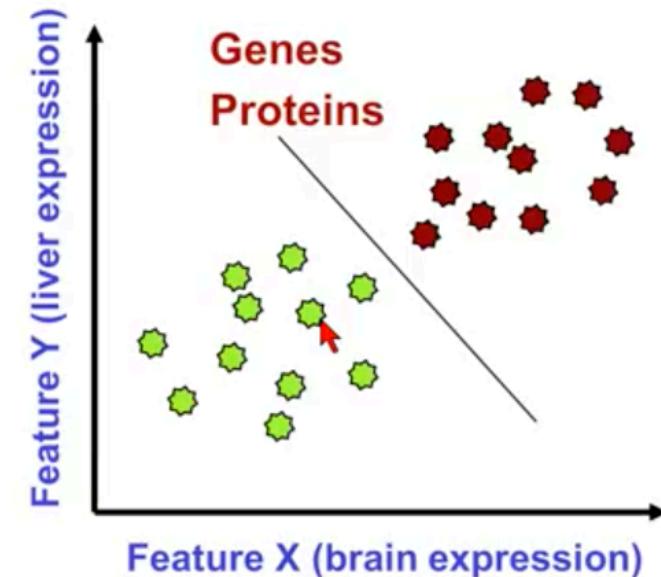


Goal of Classification: Extract features from the data that best assign new elements to ≥ 1 of well-defined classes

- **Supervised learning**

Clustering vs Classification

- Objects characterized by one or more features
- **Classification (supervised learning)**
 - Have labels for some points
 - Want a “rule” that will accurately assign labels to new points
 - Sub-problem: Feature selection
 - Metric: Classification accuracy
- **Clustering (unsupervised learning)**
 - No labels
 - Group points into clusters based on how “near” they are to one another
 - Identify structure in data
 - Metric: independent validation features



- Partitioning (e.g. k-means)
 - Divides objects into non-overlapping clusters such that each data object is in exactly one subset
- Agglomerative (e.g. hierarchical clustering)
 - A set of nested clusters organized as a hierarchy

The Basic Idea

- Assume a **fixed number** K of clusters
- Partition points into K compact clusters

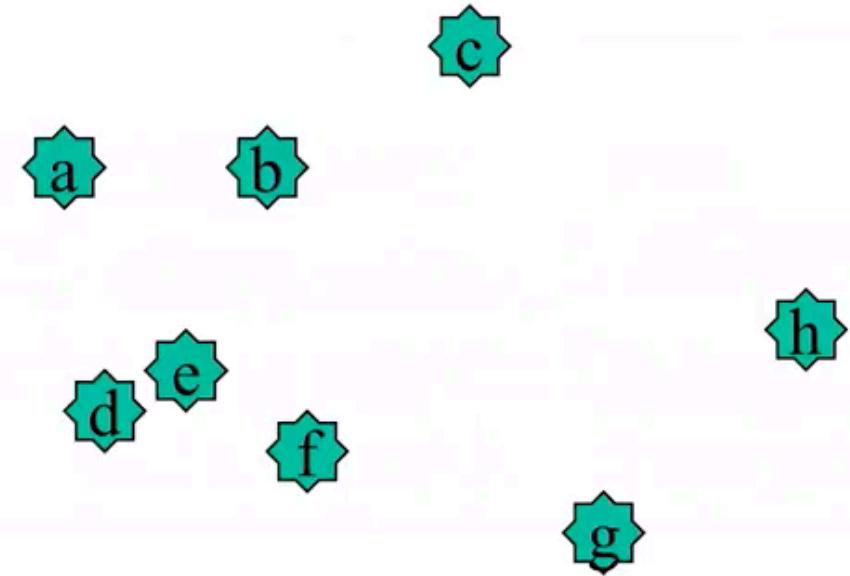
The Algorithm

- Initialize K cluster centers randomly
- Repeatedly:
 - Assign points to nearest center
 - Move centers to center of gravity of their points
- Stop at convergence (no more reassignments)

- How do we select K?
 - We can always make clusters “more compact” by increasing K
 - e.g. What happens if $K = \text{number of data points}$?
 - What is a meaningful improvement?
- Hierarchical clustering side-steps this issue

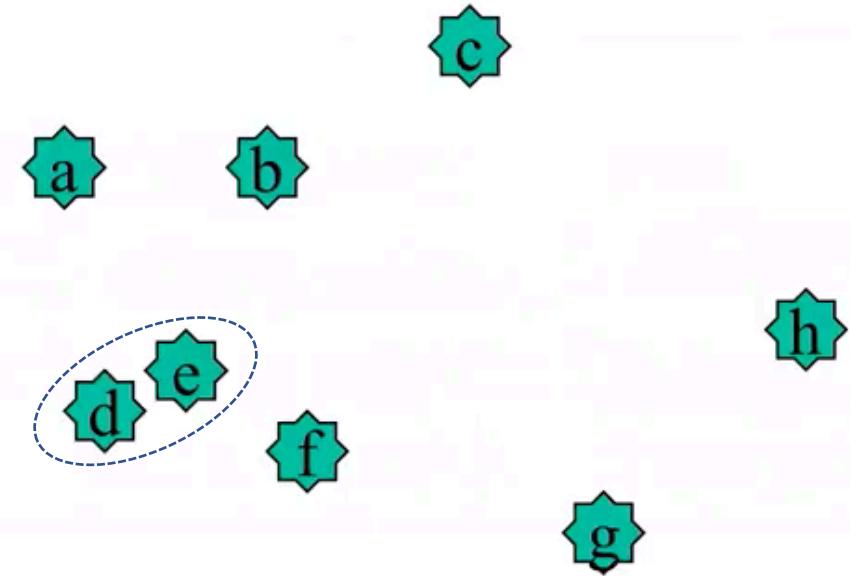
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



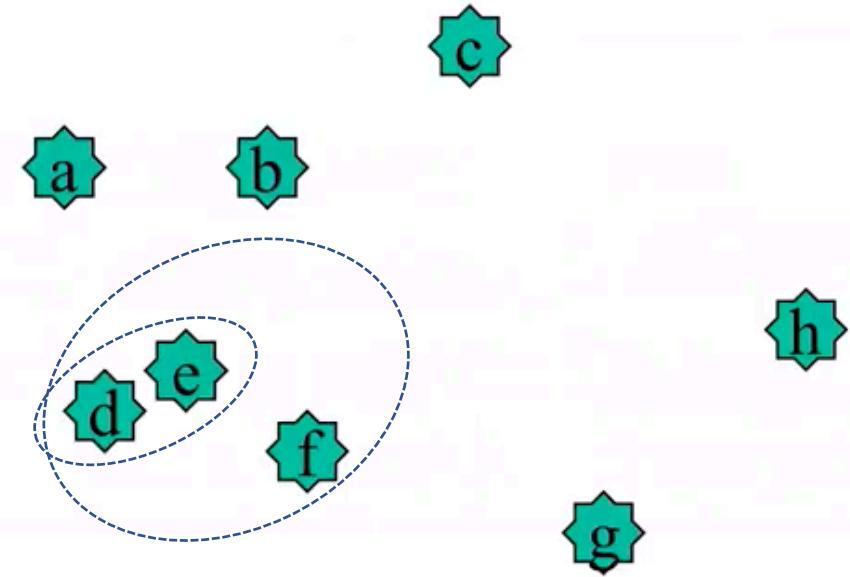
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



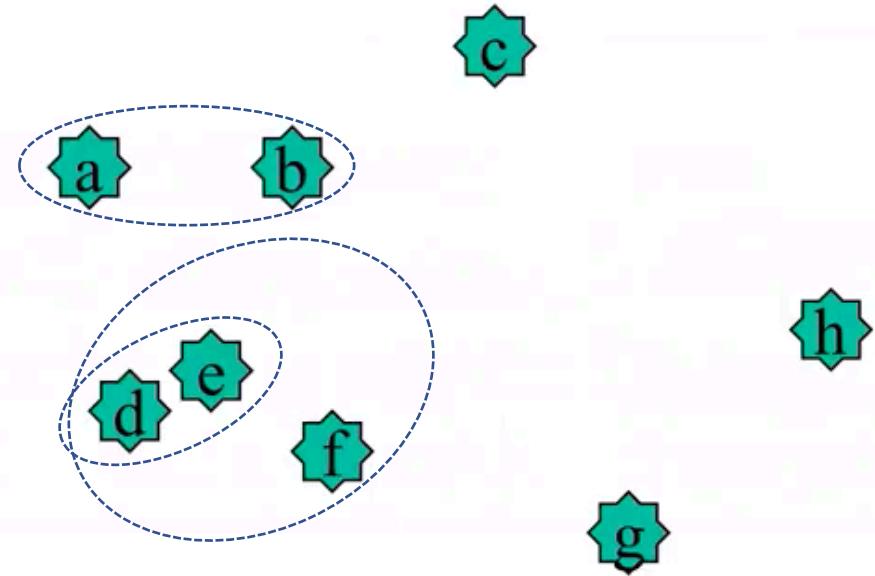
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



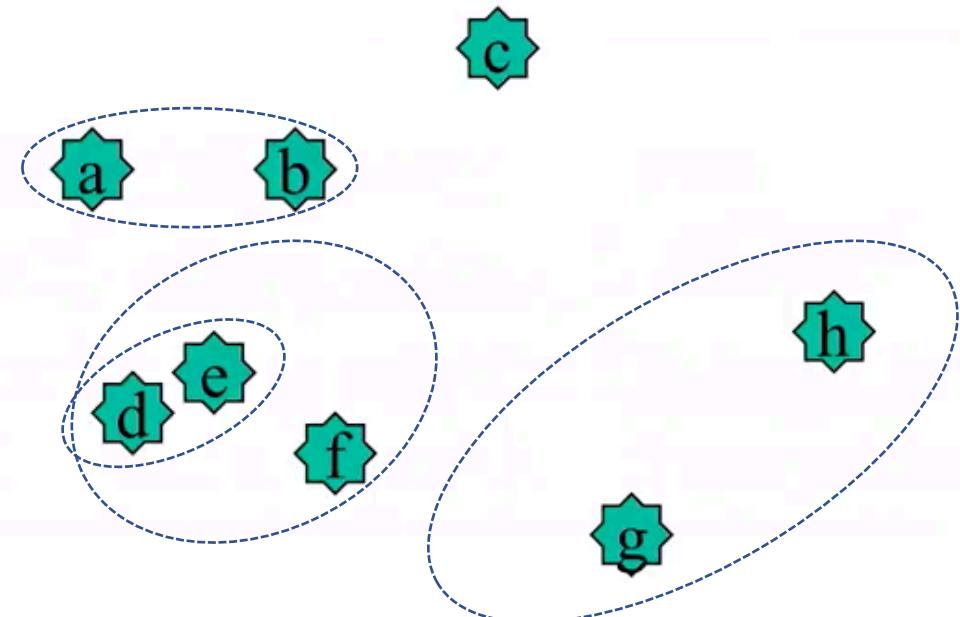
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



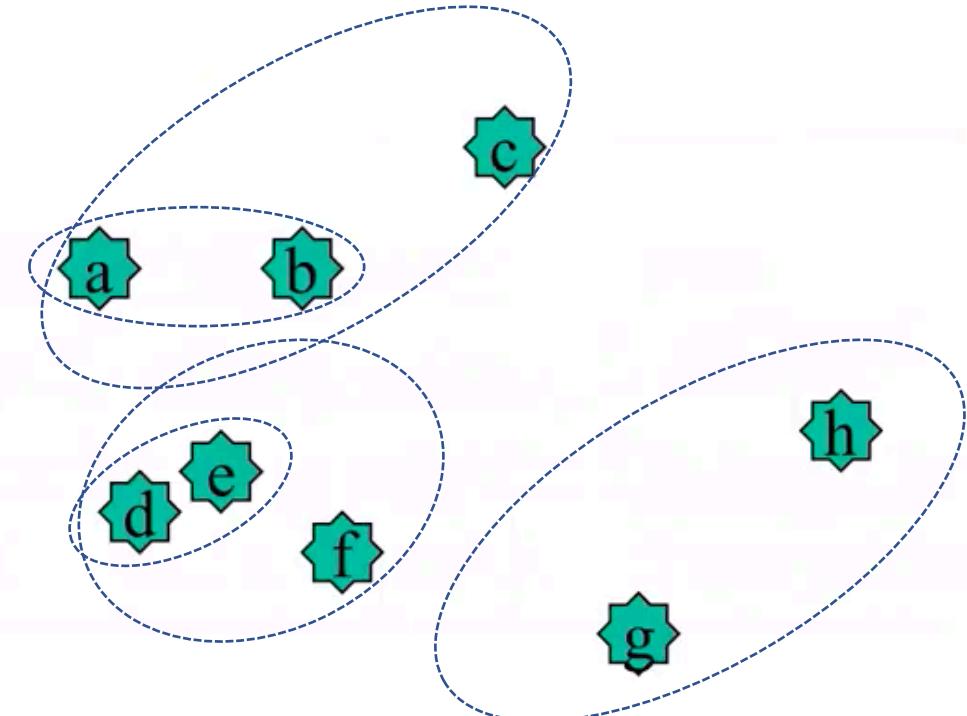
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



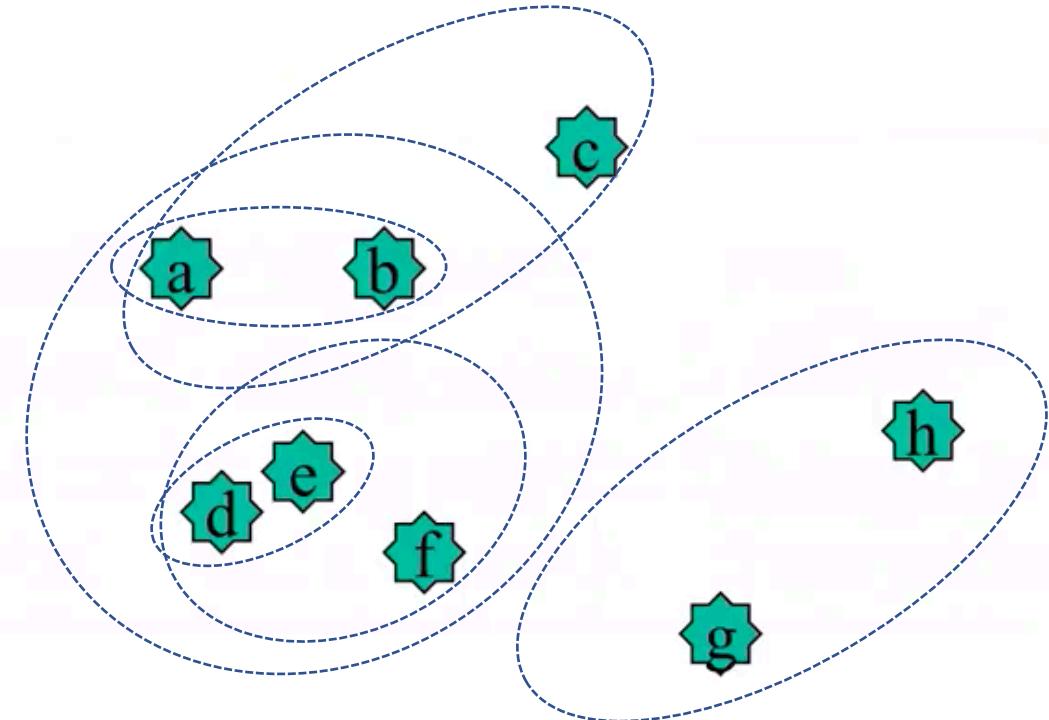
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



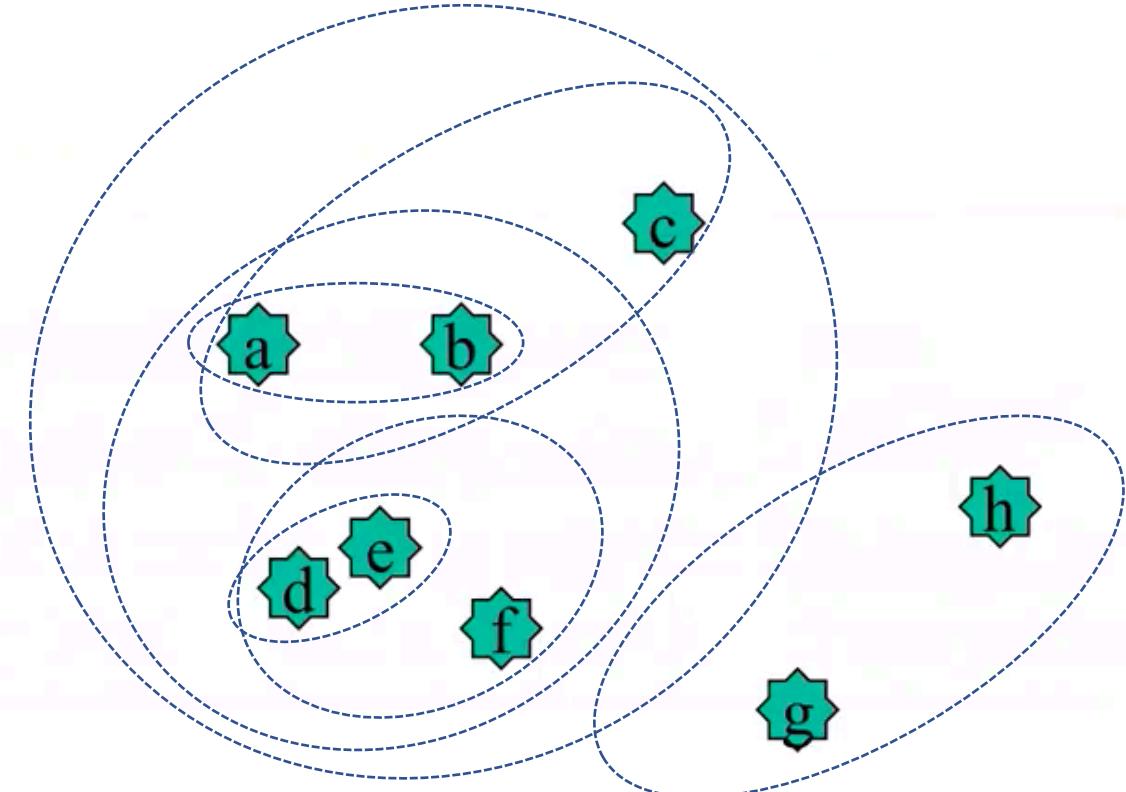
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



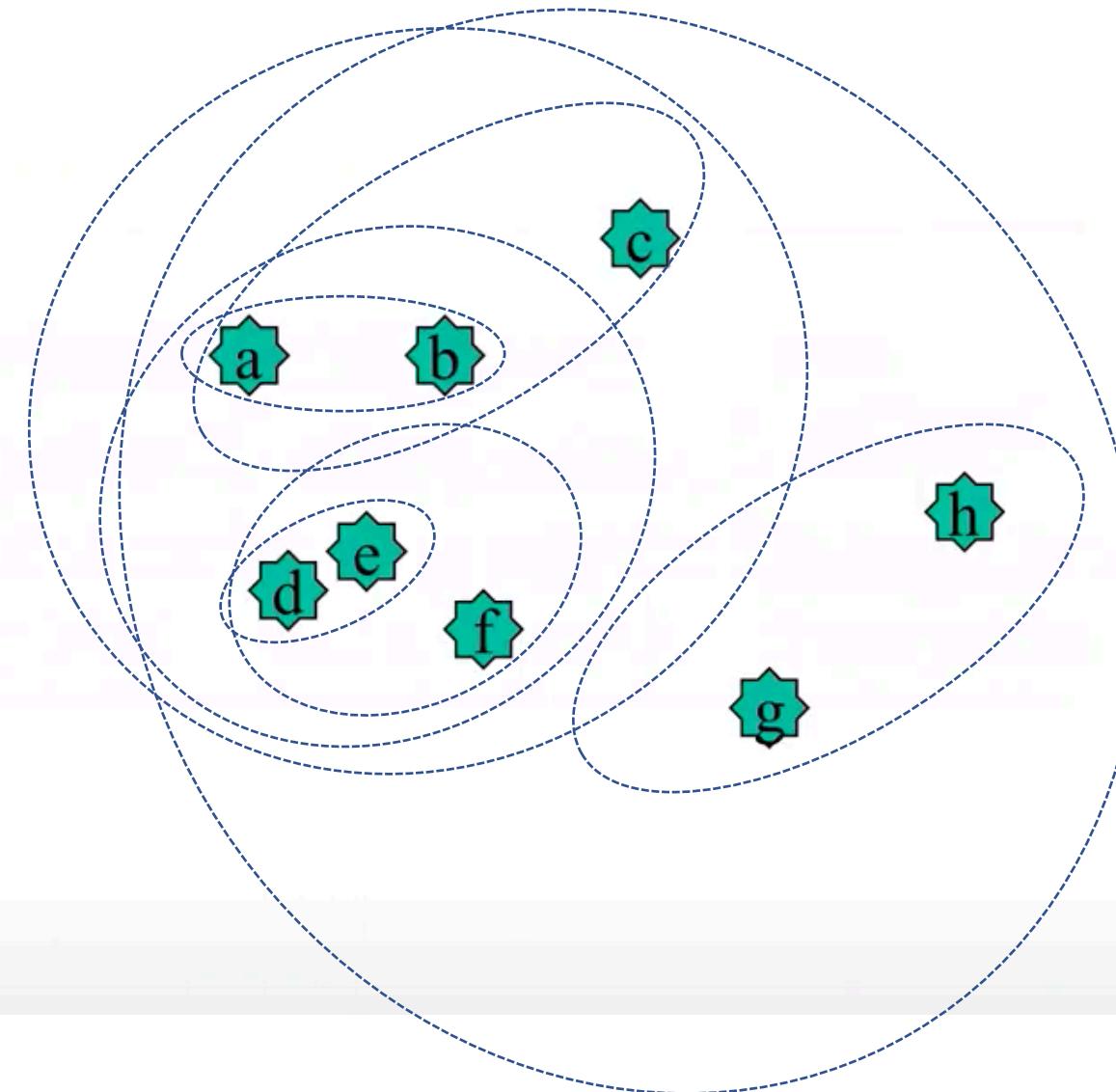
Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge

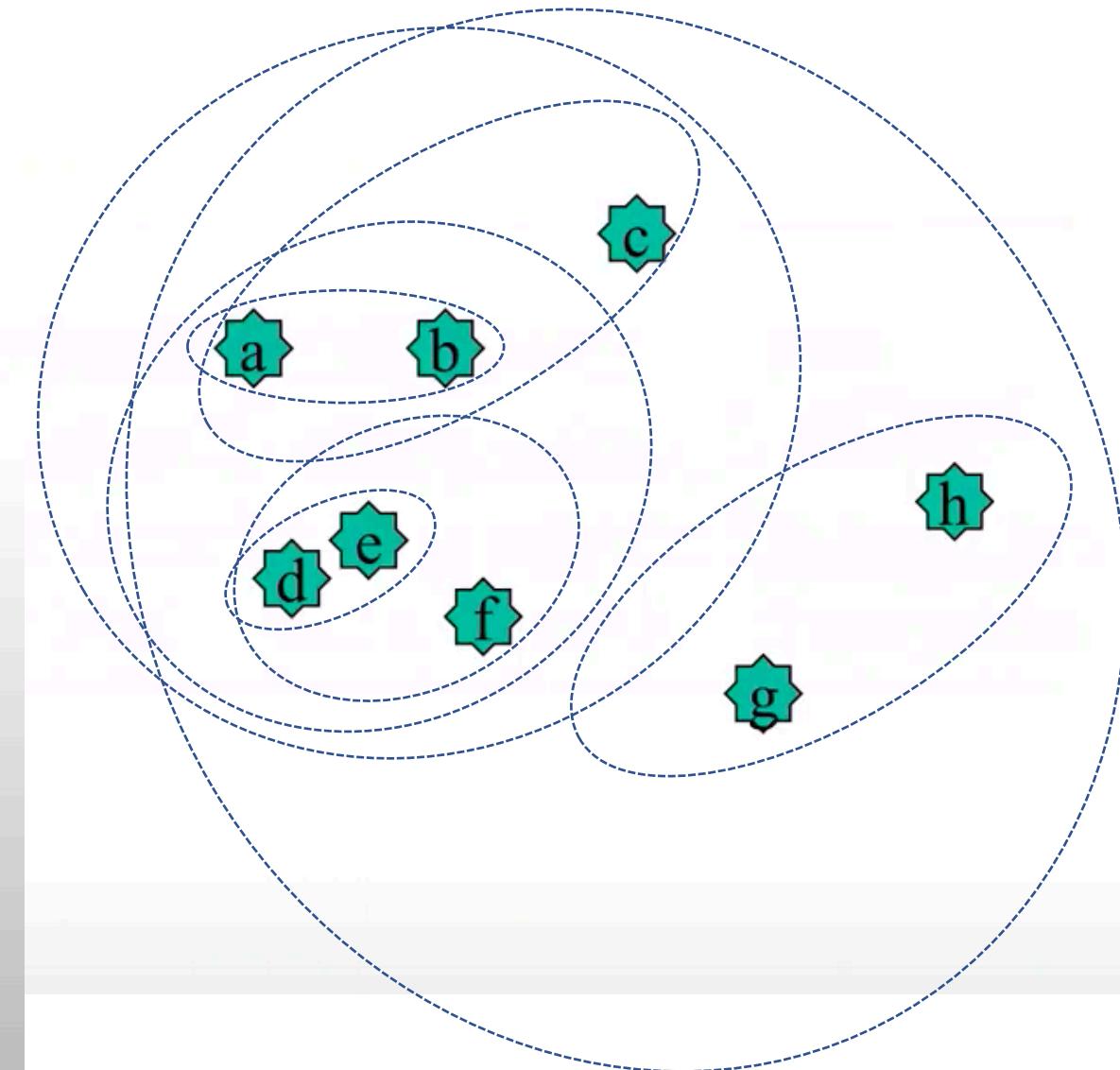
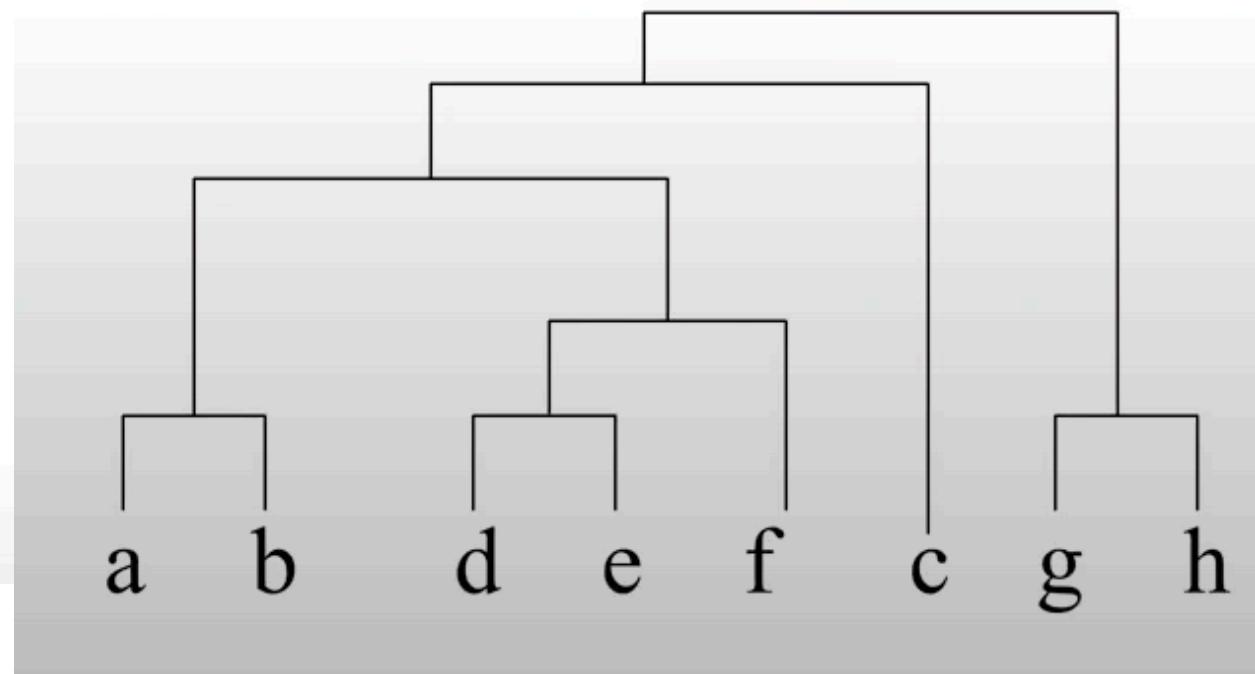


Most widely used algorithm for expression data

- Start with each point in a separate cluster
- At each step:
 - Choose the pair of **closest clusters**
 - Merge



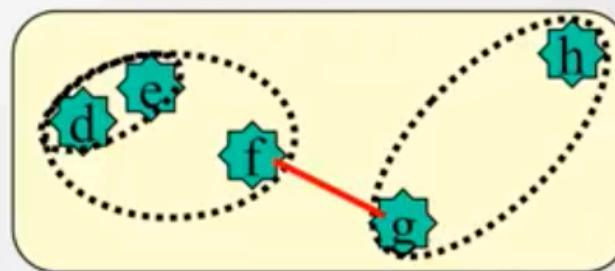
Most widely used algorithm for expression data



Distance between clusters

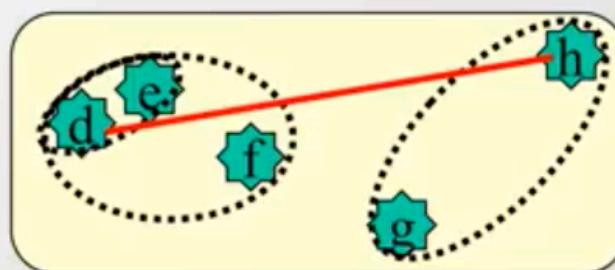
- $CD(X, Y) = \min_{x \in X, y \in Y} D(x, y)$

Single-link method



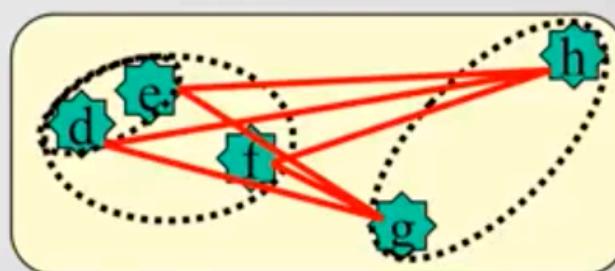
- $CD(X, Y) = \max_{x \in X, y \in Y} D(x, y)$

Complete-link method



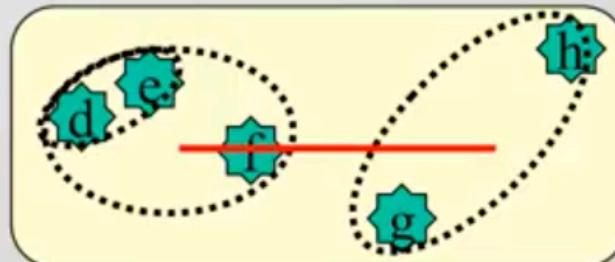
- $CD(X, Y) = \text{avg}_{x \in X, y \in Y} D(x, y)$

Average-link method



- $CD(X, Y) = D(\text{avg}(X), \text{avg}(Y))$

Centroid method



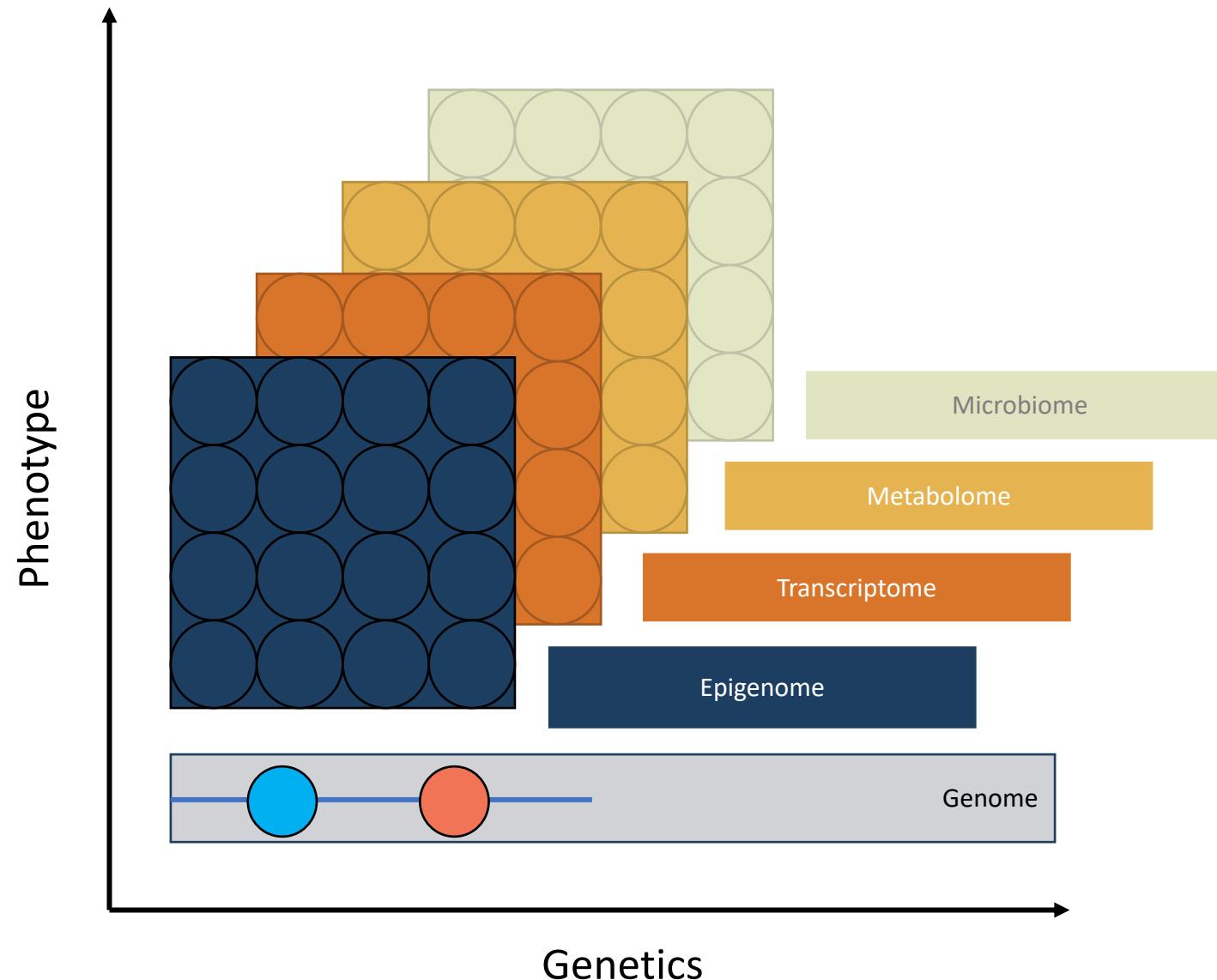
Transcriptomics

Transcriptomics examines RNA levels genome-wide, both qualitatively (which transcripts are present, identification of novel splice sites, RNA editing sites) and quantitatively (how much of each transcript is expressed). The central dogma of biology viewed RNA as a molecular intermediate between DNA and proteins, which are considered the primary functional read-out of DNA.

An even more significant contribution of these studies was the development of the non-coding RNA field. It is now clear that thousands of long non-coding RNAs transcribed in mammalian cells (<http://www.gencodegenes.org/>) play essential roles in many physiological processes, for example, brown adipose differentiation (Cell Metab. 2015;21:764–76), endocrine regulation (Nat Rev Endocrinol. 2015;11:151–60), and neuron development (Nat Rev Neurosci. 2016;17:537–49).

Dysregulation of long non-coding RNAs had been implicated in various diseases, such as myocardial infarction(J Hum Genet. 2006;51:1087–99), diabetes(Genes Dev. 2016;30:502–7), cancer (Nature. 2010;464:1071–6), and others (Cell Mol Life Sci. 2016;73:2491–509).

Microbiome





Microbiome

This term refers to the entire habitat, including the microorganisms (bacteria, archaea, lower and higher eukaryotes, and viruses), their genomes (i.e., genes), and the surrounding environmental conditions. This definition is based on that of “biome,” the biotic and abiotic factors of given environments. Others in the field limit the definition of microbiome to the collection of genes and genomes of members of a microbiota. It is argued that this is the definition of metagenome, which combined with the environment constitutes the microbiome. The microbiome is characterized by the application of one or combinations of metagenomics, metabonomics, metatranscriptomics, and metaproteomics combined with clinical or environmental metadata.

Metataxonomics

Metataxonomics is a term we propose and define as the high-throughput process used to characterize the entire microbiota and create a metataxonomic tree, which shows the relationships between all sequences obtained. While viruses are an integral part of the microbiota, no universal viral marker genes are available to perform such taxonomic assignments.

Microbiota. The assemblage of microorganisms present in a defined environment

EDITORIAL

Open Access



CrossMark

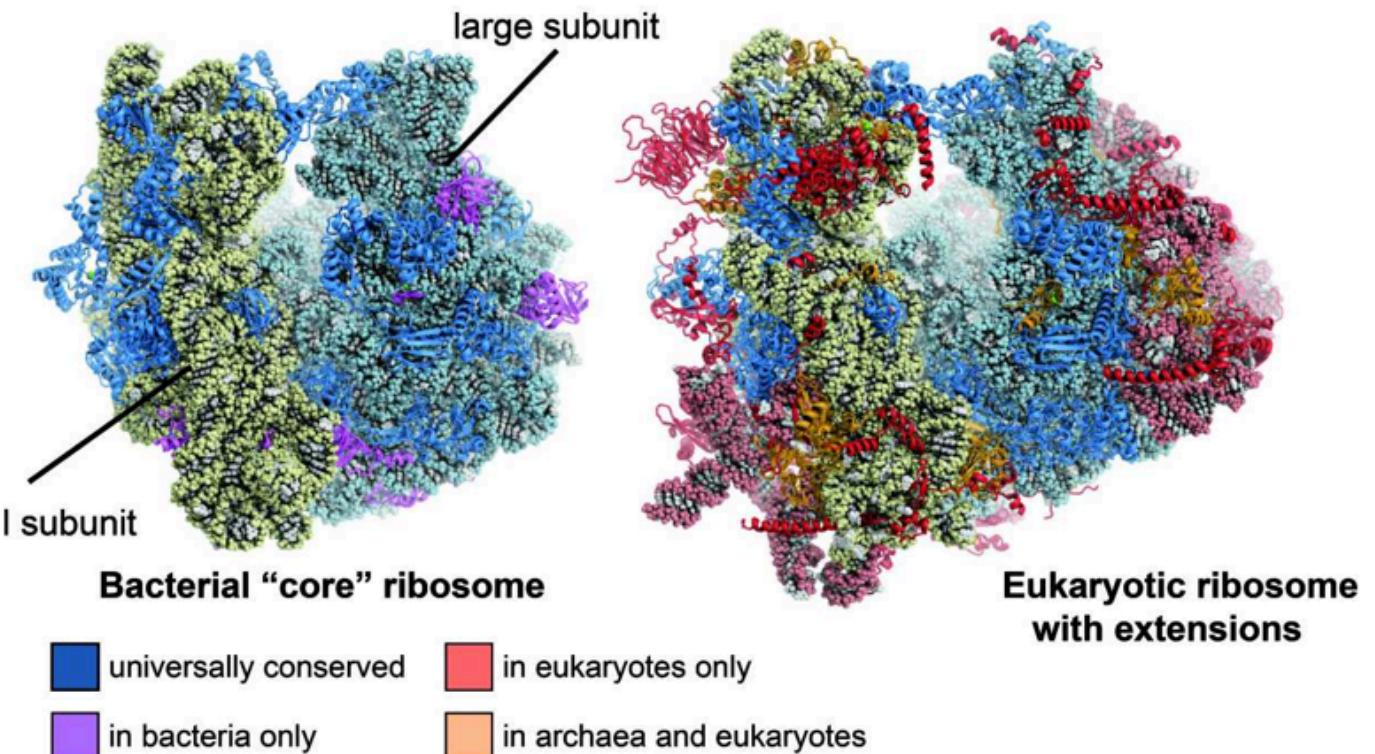
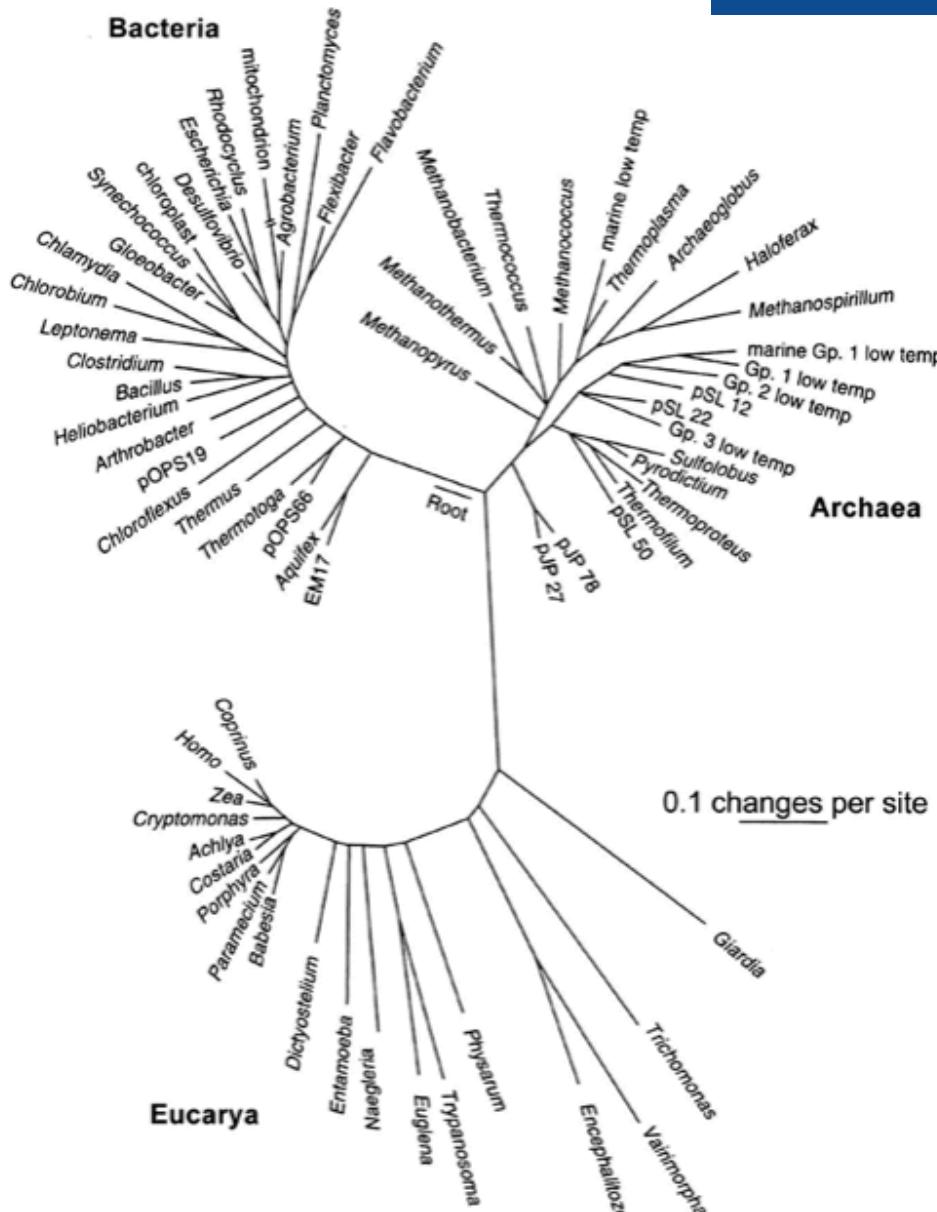
The vocabulary of microbiome research: a proposal

Julian R. Marchesi^{1,2} and Jacques Ravel^{3,4*}

Metagenome

The collection of genomes and genes from the members of a microbiota. This collection is obtained through shotgun sequencing of DNA extracted from a sample (metagenomics) followed by assembly or mapping to a reference database followed by annotation. Metataxonomic analysis, because it relies on the amplification and sequencing of taxonomic marker genes, is not metagenomics. Metagenomics is the process used to characterize the metagenome, from which information on the potential function of the microbiota can be gained.

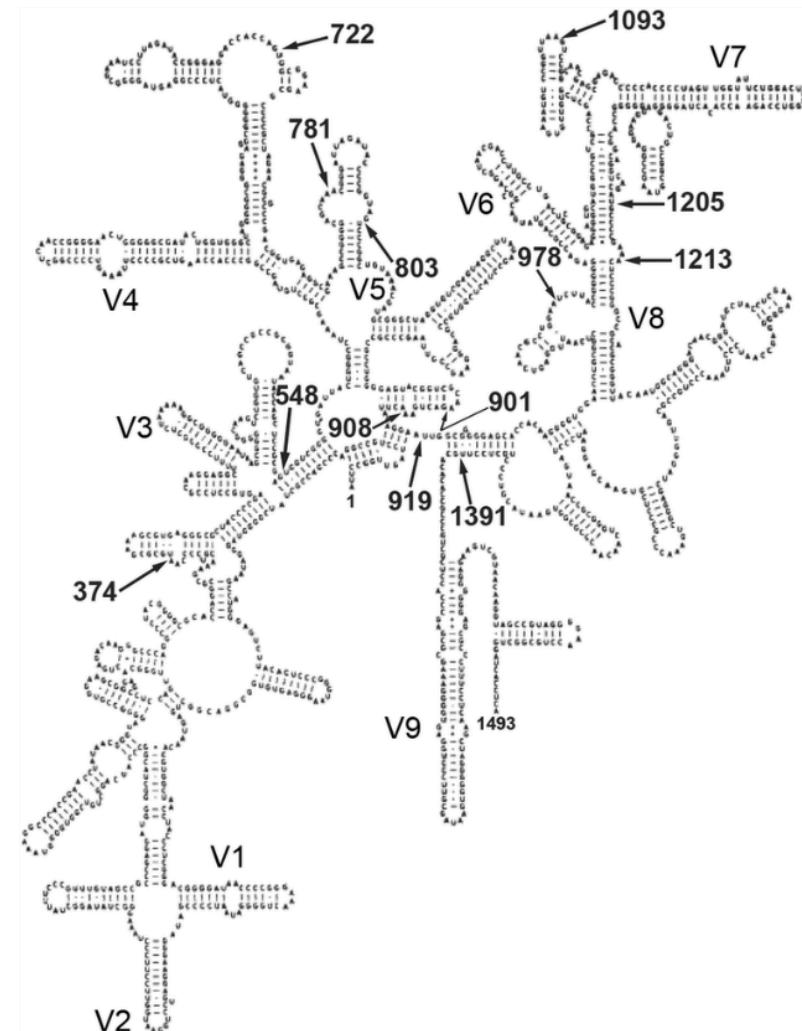
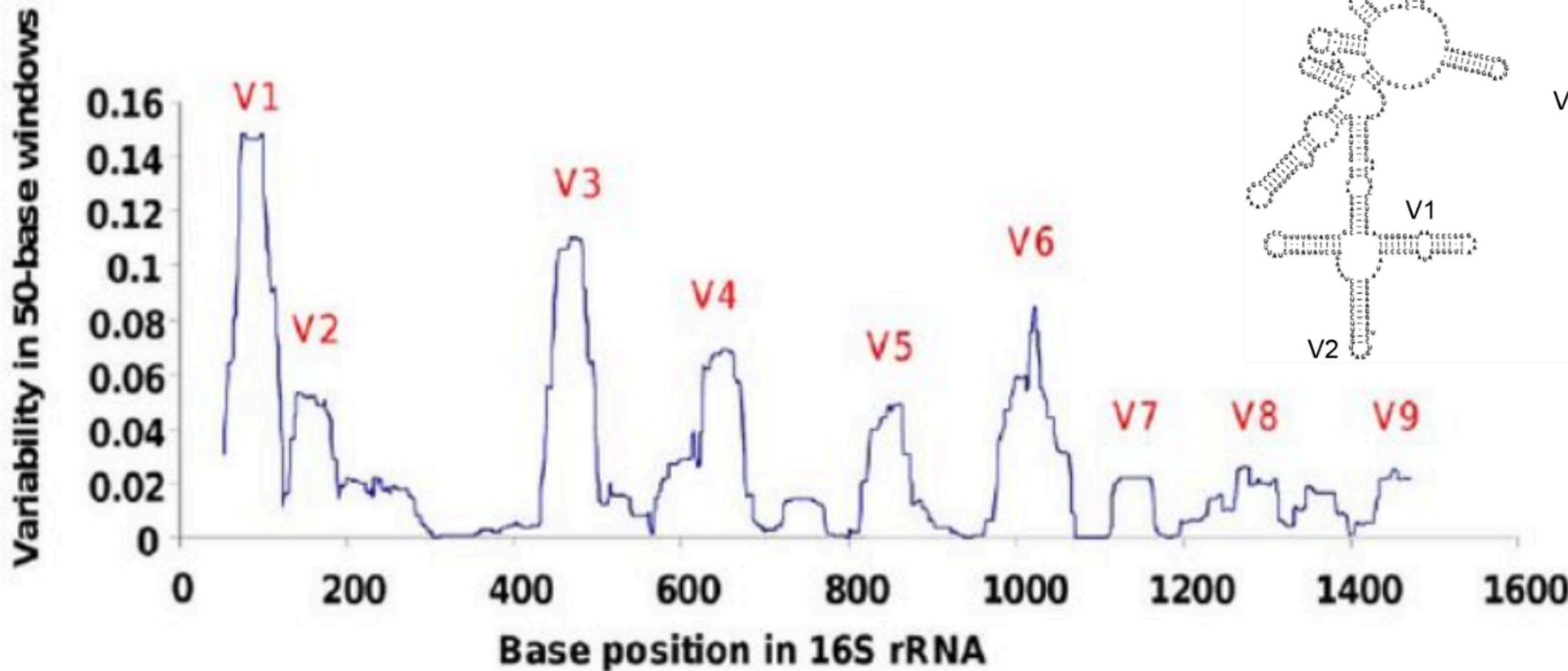
Core marker genes are specific taxonomic clades

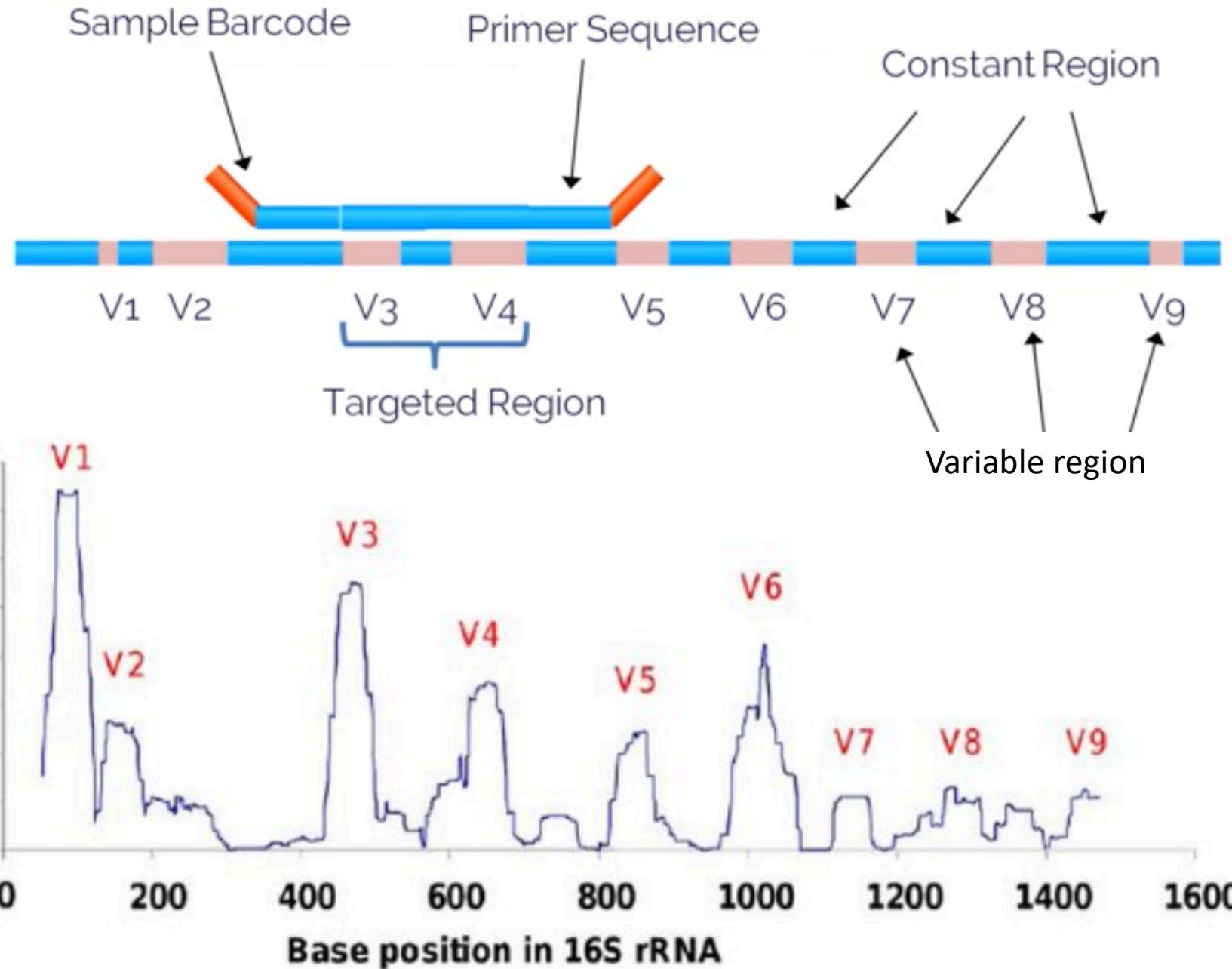


Phylogenetic tree based on the 16S rRNA gene sequence comparisons

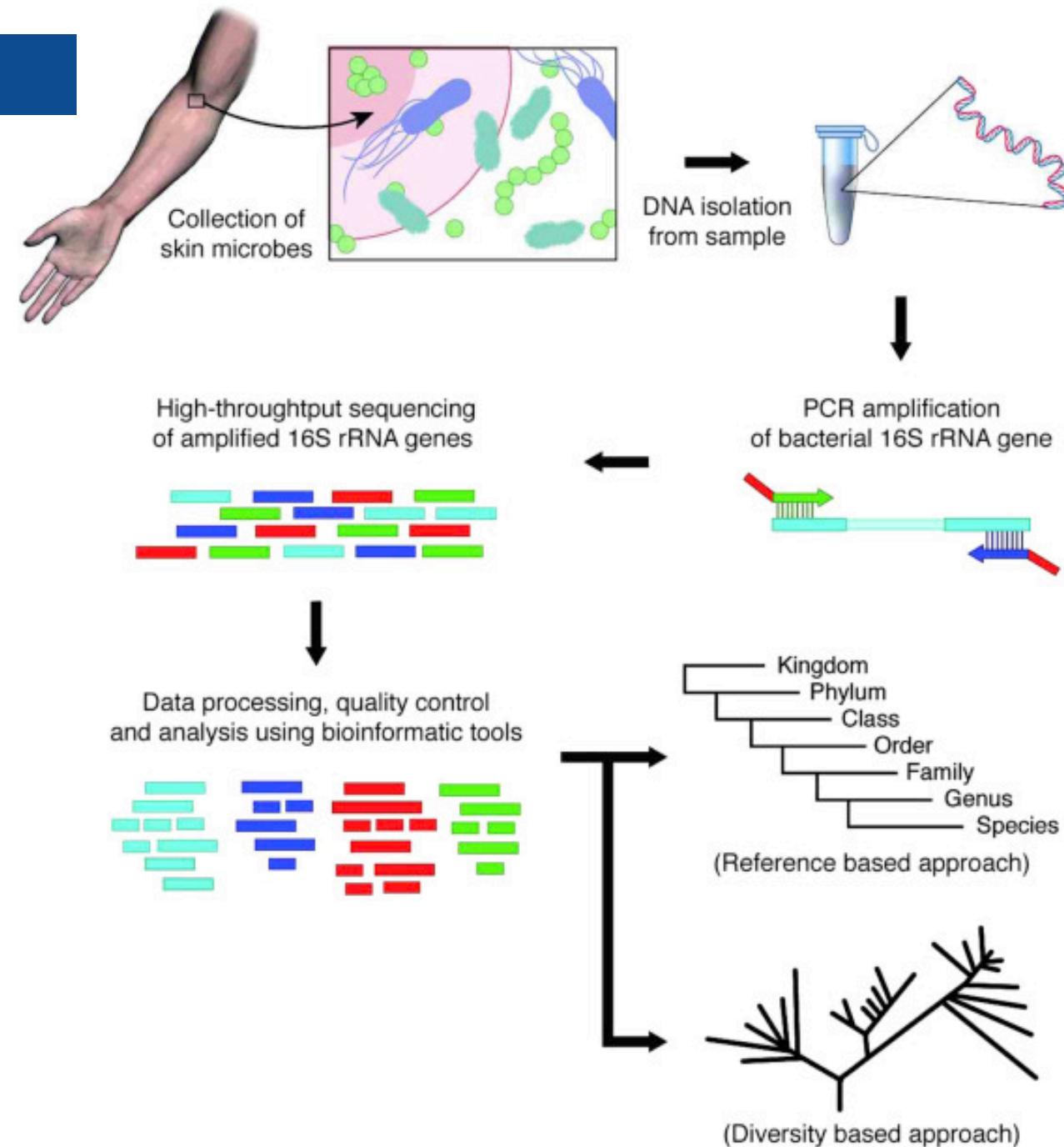
Bacterial diversity: 16S rRNA gene

Universal and variable regions of 16S rRNA used for PCR amplification and classification



B

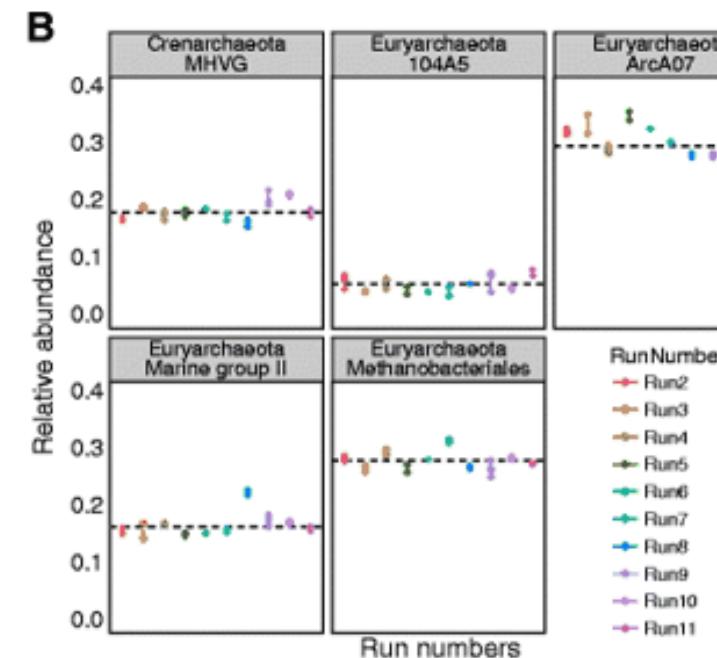
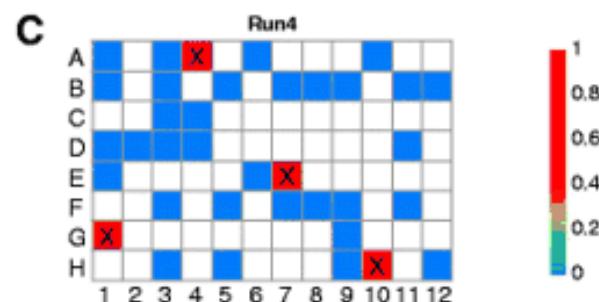
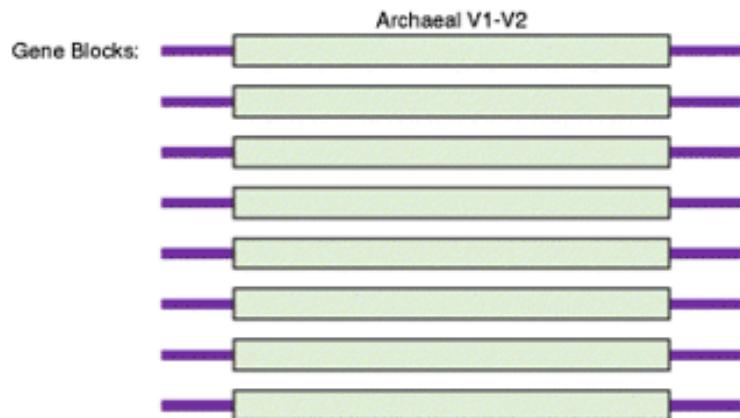
Basic workflow: 16S rRNA gene



Importance of Control

Negative control: e.g. air swab processed through entire DNA prep, sequencing analysis

Positive control: e.g. Mock community (purchased from BEI or made in house). Do not underestimate primer bias or chimeras.



Mock community contains 20 bacteria. Amplified and sequenced with various methods

Optimizing methods and dodging pitfalls in microbiome research

Dorothy Kim, Casey E. Hofstaedter, Chunyu Zhao, Lisa Mattei, Ceylan Tanes, Erik Clarke, Abigail Lauder, Scott Sherrill-Mix, Christel Chehoud, Judith Kelsen, Máire Conrad, Ronald G. Collman, Robert Baldassano, Frederic D. Bushman & Kyle Bittinger 

Microbiome 5, Article number: 52 (2017) | [Cite this article](#)

16k Accesses | 140 Citations | 69 Altmetric | [Metrics](#)

Suites of analysis tools



STUDY DESIGNS

Experimental and analytical tools for studying the human microbiome

Justin Kuczynski¹, Christian L. Lauber², William A. Walters¹, Laura Wegener Parfrey³, José C. Clemente³, Dirk Gevers⁴ and Rob Knight^{3,5}



[Genome Res.](#) 2009 Jul; 19(7): 1141–1152.

doi: [10.1101/gr.085464.108](https://doi.org/10.1101/gr.085464.108)

PMCID: PMC3776646

PMID: [19383763](https://pubmed.ncbi.nlm.nih.gov/19383763/)

Microbial community profiling for human microbiome projects: Tools, techniques, and challenges

[Micah Hamady](#)¹ and [Rob Knight](#)^{2,3}

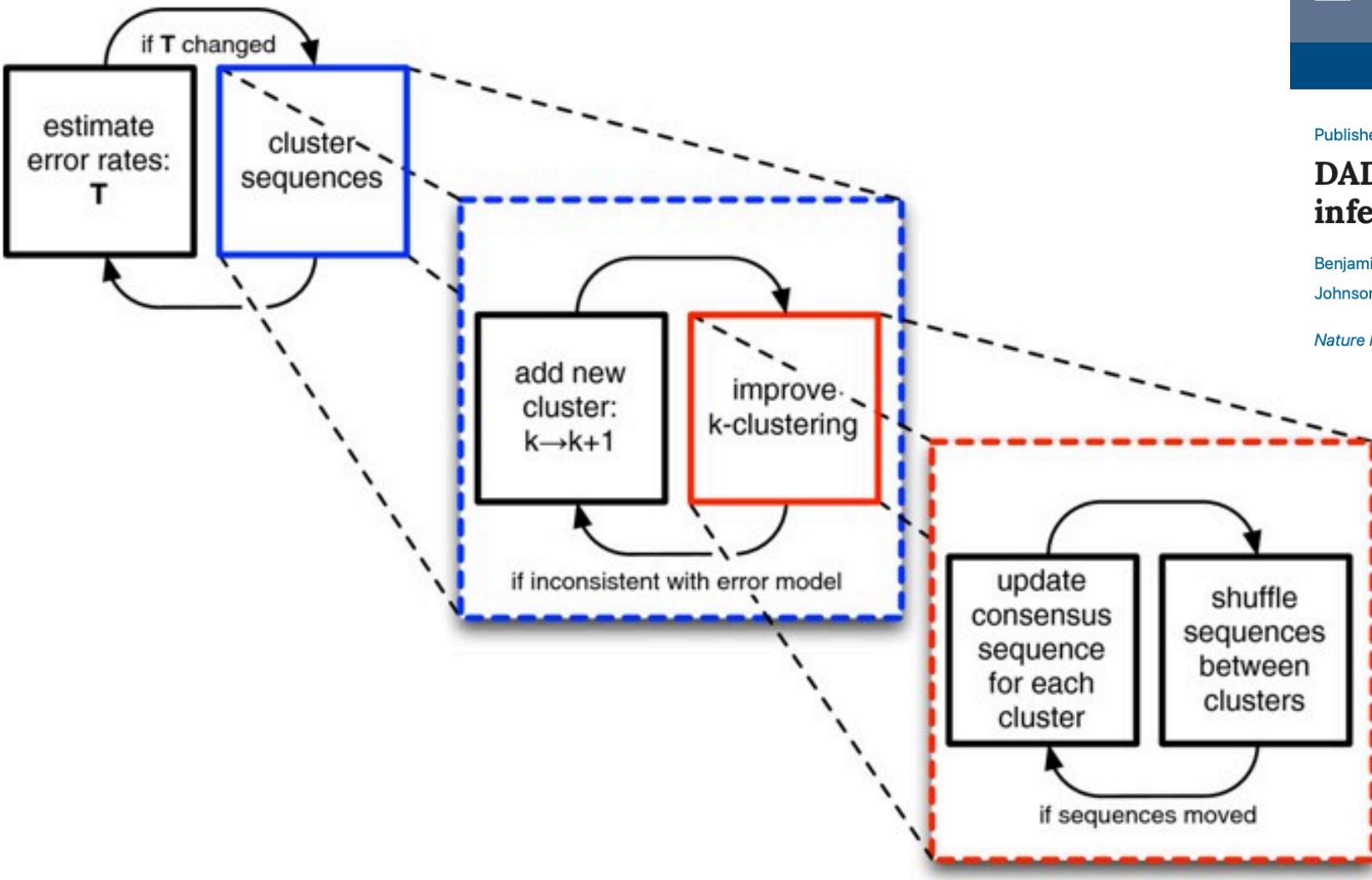
OTU table

OUT table is a sample-by-observation matrix

	Sam1	Sam2	Sam3	
Taxonomy	Otu1	0.34	0.32	0.29
	Otu2	0.12	0.17	0.10
	Otu3	0.07	0.03	0.11
	Otu4	0.06	0.02	0.09

OTU can have corresponding taxonomy information

OTUs were developed as a tool for classifying microbial species, but have also been repurposed to the task of correcting errors; the sequences within an OTU are typically interpreted as a taxonomic grouping without specifying whether the variation within an OTU represents errors or real diversity on a finer scale than that chosen to define the OTU. If the scale of the noise is smaller than that of the clusters, then the construction of OTUs will appropriately group error-containing sequences together with their true genotype.



DADA schematic. The basic structure of *DADA*, an algorithm to denoise amplicon sequence data. See Algorithm Algorithm 1 in the Methods section for the pseudocode and a more detailed description.

Published: 23 May 2016

DADA2: High-resolution sample inference from Illumina amplicon data

Benjamin J Callahan , Paul J McMurdie, Michael J Rosen, Andrew W Han, Amy Jo A Johnson & Susan P Holmes

Nature Methods 13, 581–583(2016) | [Cite this article](#)

BMC Bioinformatics

[Home](#) [About](#) [Articles](#) [In Review](#) [Submission Guidelines](#)

We'd like to understand how you use our website

Research article | [Open Access](#) | Published: 31 October 2012

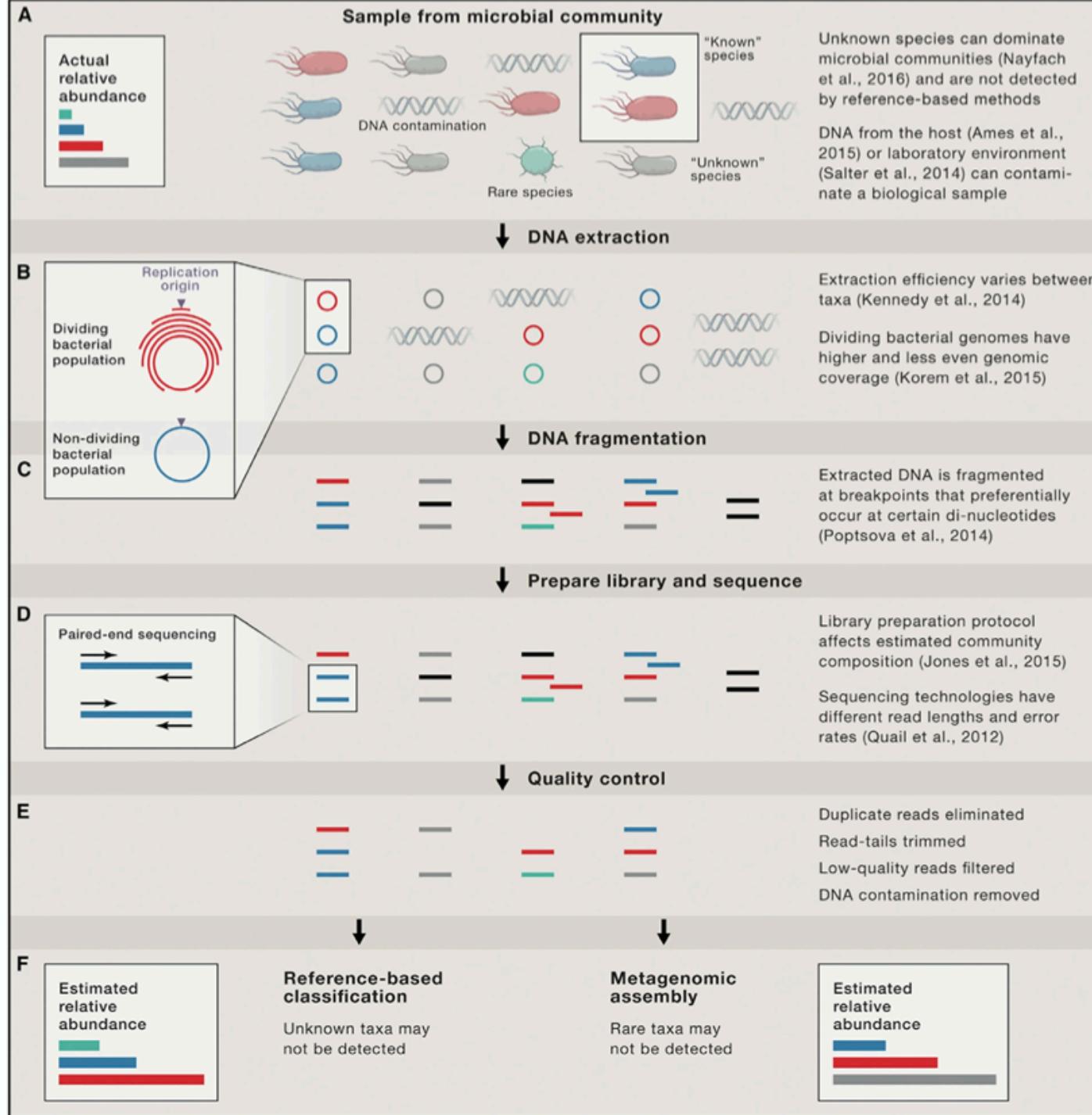
Denoising PCR-amplified metagenome data

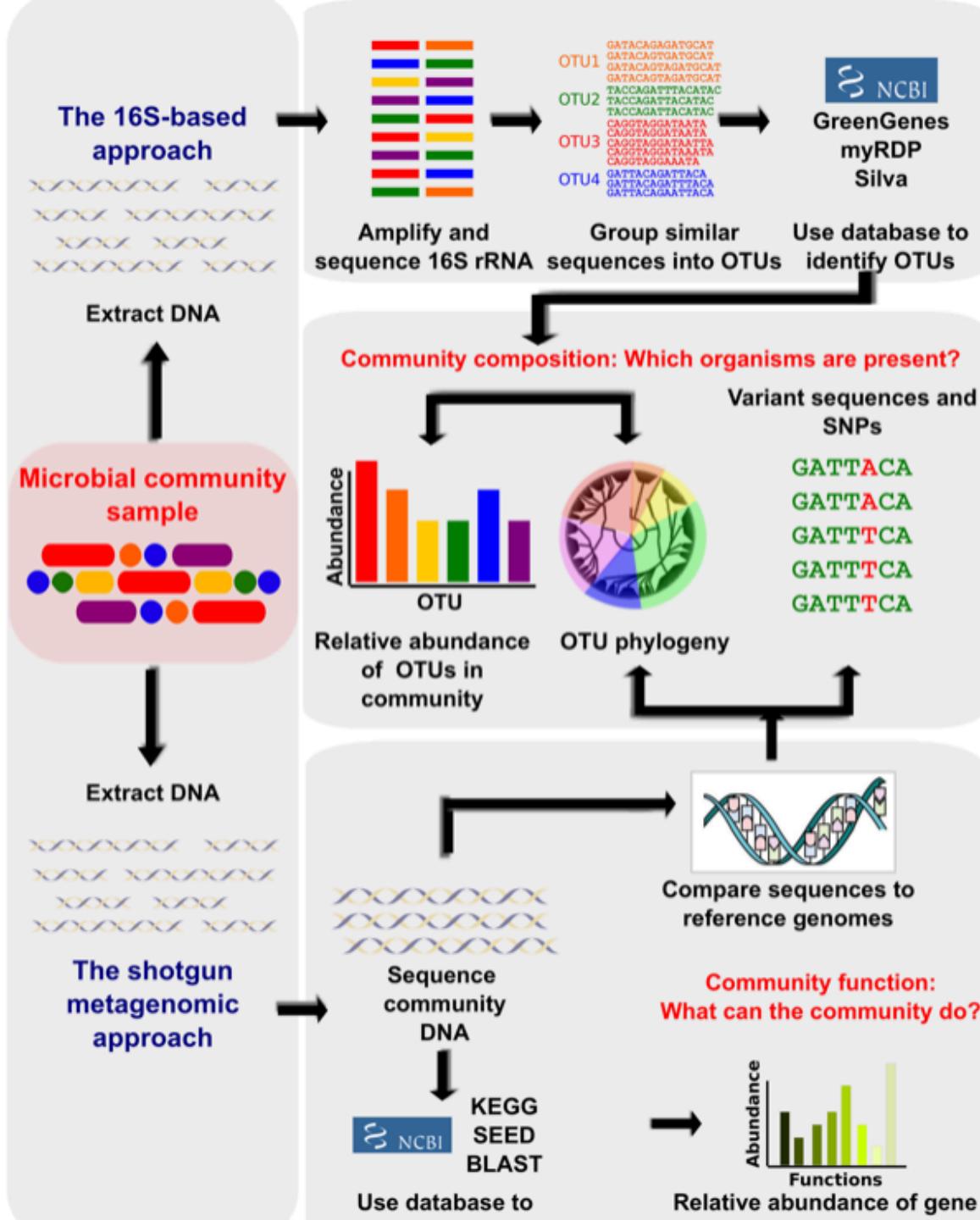
Michael J Rosen , Benjamin J Callahan, Daniel S Fisher & Susan P Holmes

BMC Bioinformatics 13, Article number: 283 (2012) | [Cite this article](#)

8770 Accesses | 30 Citations | 7 Altmetric | [Metrics](#)

Shotgun DNA Sequencing





Microbiome

Microbiomics is a fast-growing field in which all the microorganisms of a given community are investigated together.

Human skin, mucosal surfaces, and the gut are colonized by microorganisms, including bacteria, viruses, and fungi, collectively known as the microbiota (and their genes constituting the microbiome).

The human microbiome is enormously complex; for example, the gut contains roughly 100 trillion bacteria from 1000 different species. There are substantial variations in microbiota composition between individuals resulting from seed during birth and development, diet and other environmental factors, drugs, and age.

Many studies have implicated perturbations in gut bacteria in a variety of disorders, including diabetes, obesity, cancer, colitis, heart disease, and autism. The microbiome can be profiled by amplifying and then sequencing certain hypervariable regions of the bacterial 16S rRNA genes followed by clustering the sequences into operational taxonomic units.

Shotgun metagenomics sequencing, in which total DNA is sequenced, can provide additional resolution for distinguishing genetically close microbial species. Several analytic tools have been developed for analyzing NGS data from targeted 16S or metagenomics analysis, such as QIIME (quantitative insights into microbial ecology).

integrative analysis of multiple omics data

