

Lecture 1: Introduction to Genomic Data

What Genomic Data Looks Like

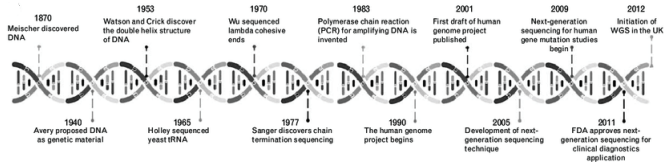
Yen-Yi Ho

Department of Statistics

- What is genomic data?
- Sequencing technologies and data generation
- Downstream analyses: RNA-seq, variant interpretation, population genomics, epigenomics

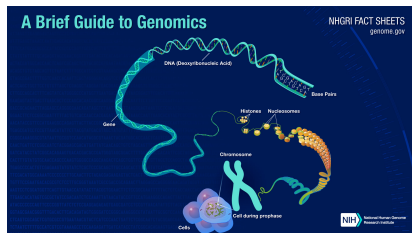


A long time ago in a galaxy far,
far away....



What is Genomic Data? — Big picture

- **Genomic data** = information derived from an organism's entire DNA sequence (the *genome*).
- Contains nucleotide sequences (A, C, G, T) and observed variation between individuals.
- Used to study *variation*, *disease*, *evolution*, and *biology*.



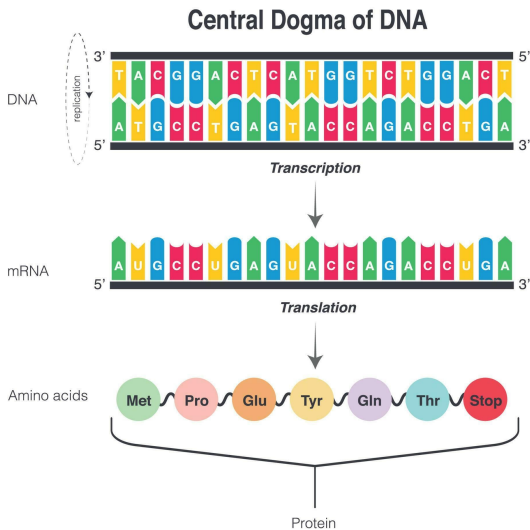
In human genome, every single cell in the body contains a complete copy of the approximately 3 billion (3×10^9) DNA base pairs (letters).

Structure of genomic data: DNA

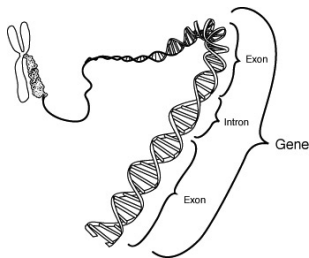
- Data is composed of nucleotides: A, C, G, T.
- Organized into **reads**, **contigs**, **chromosomes**, and **annotations**.
- Common file formats: FASTQ, FASTA, VCF, GTF/GFF.



Central Dogma of Molecular Biology

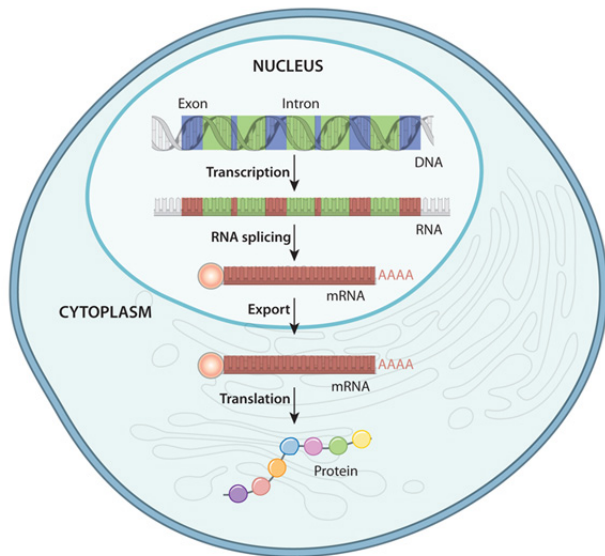


- A locatable region of genomic sequence, corresponding to a unit of inheritance, which is associated with regulatory regions, transcribed regions, and or other functional sequence regions.
- Or simply, a piece of “useful” DNA sequence.

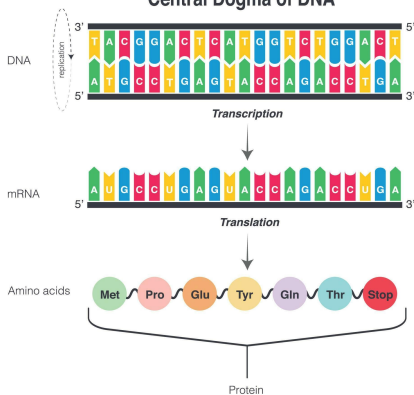


- **enhancer**: a region for enhancing gene expression. Not necessarily close to the gene.
- **promoter**: at the beginning of the gene, helps transcription.
- **exons**: the “useful” part of the gene, will appear in the mRNA product.
- **introns**: the “spacer” between exons, will NOT be in the mRNA product.
- **splicing**: the process to remove introns and join exons.
- **alternative splicing**: different splicing patterns for the same pre-mRNA. For example, mRNA could be from exons 1 and 2 or exons 1 and 3. Those are different **transcripts** of the same gene.

Gene structure and splicing



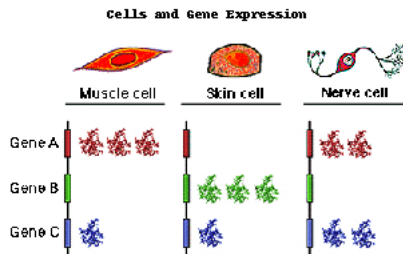
Central Dogma of DNA



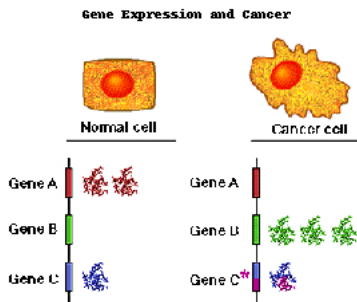
Gene Expression

Gene expression is a term that is used to describe the entire process of translation and transcription of a gene. Gene expression is a highly specific process. Only a small fraction of the genes are expressed, or turned "on," in any particular type of cell.

gene expression in different tissues

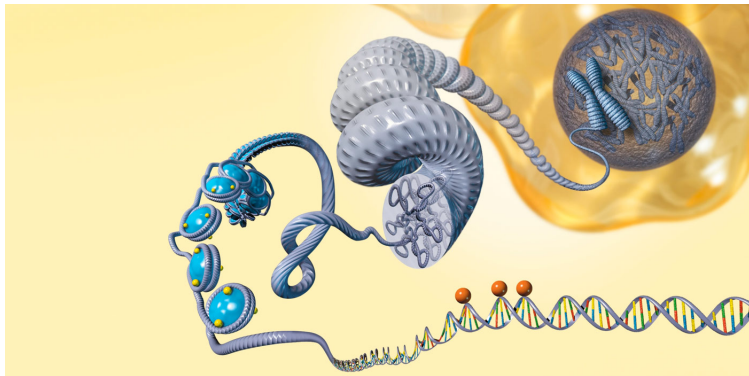


gene expression in the same tissue, but different points in time



Epigenetics

Non-DNA sequence related, heritable mechanisms to control gene expressions. Example: DNA methylation, histone modifications.



Source of Variations



C-Series



R-Series



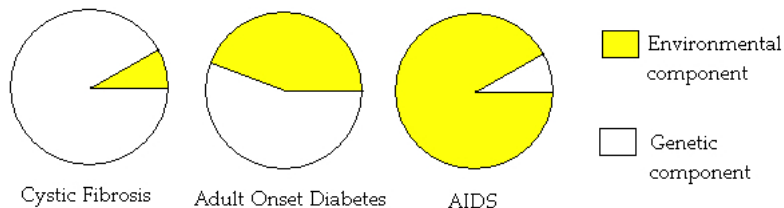
K-Series



BB-Series

Environment Vs. Gene

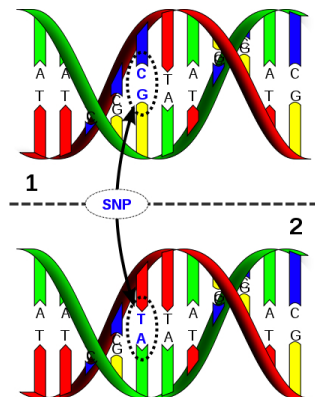
Any two individuals are 99.9% identical in their DNA



Genetic Variations (Polymorphisms)

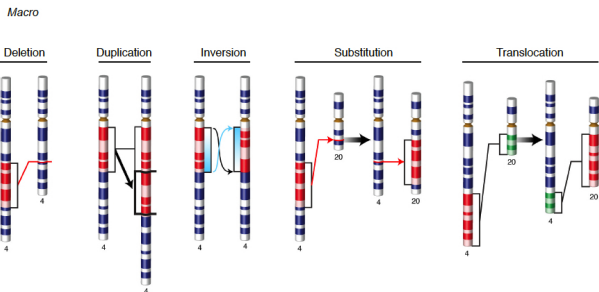
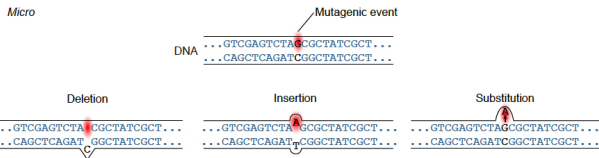
That 0.1 % is very important in defining our differences

- single nucleotide polymorphisms (SNPs, every 300 nucleotide on average)
- small-scale mutation, insertions, deletions
- copy number variations (AAGAAGAAGAAG)

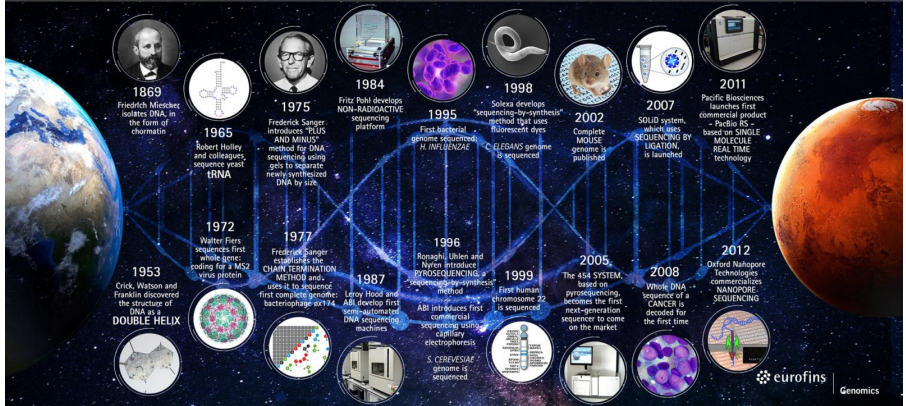


source: <http://ghr.nlm.nih.gov/handbook/genomicresearch/snp>

Mutations



A JOURNEY THROUGH THE HISTORY OF DNA SEQUENCING



How Sequencing Works?

- Determining the order of nucleotides (A, T, C, G)
- Technologies differ in chemistry, accuracy, read length, throughput
- Two major eras:
 - Next-Generation Sequencing (NGS)
 - Long-read (3rd generation) sequencing

- Massively parallel sequencing
- Short reads (50–300 bp)
- High throughput, low cost
- Dominant platforms: Illumina, Ion Torrent

Sequencing-by-Synthesis (SBS)

- Bridge amplification on flowcell
- Reversible fluorescent terminators
- Very high accuracy
- Applications: WGS, RNA-seq, ChIP-seq



Ion Torrent Sequencing

- Semiconductor detection of pH change
- No optical system needed
- Read length: 200–400 bp
- Limitations: homopolymer errors



Oxford Nanopore Technologies (ONT)

- Nanopores embedded in membranes measure ionic current
- Very long reads (10 kb–1 Mb)
- Portable devices (MinION)
- Applications: assembly, metagenomics, transcriptomics

Oxford Nanopore Sequencing



PacBio SMRT Sequencing (Long Reads)

- Single molecule real-time sequencing
- HiFi reads: 10–25 kb with very high accuracy (Q20+)
- Ideal for structural variants and de novo assembly



How Sequencing Data Is Generated

- **Step 1: Sample Preparation**

- Extract DNA or RNA
- Assess purity and concentration

- **Step 2: Library Preparation**

DNA is fragmented and adapters are added.

- **Step 3: Sequencing Run**

Depends on platform:

- Illumina: sequencing-by-synthesis with fluorescent nucleotides
- Nanopore: strands pass through nanopores
- PacBio: circular consensus sequencing

- **Step 4: Base Calling**

Raw signals → nucleotide sequences

- Illumina: fluorescence intensities
- ONT: electrical current patterns
- PacBio: light pulses



Data Output Formats

- FASTQ — sequences with quality scores
- BAM/CRAM — aligned reads (compressed)
- VCF — variants
- FASTA — reference sequences

```
@SEQ_ID
GATTTGGGGTC AAGAAAAGCA
+
!!"(((((((**...**)%}))}
```

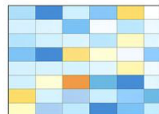
Raw Reads



Assembled Genome

CHROM	POS	REF	ALT
17	1234567	C	T
20	1110696	G	D
3	16512	T	AG

Variant Calls

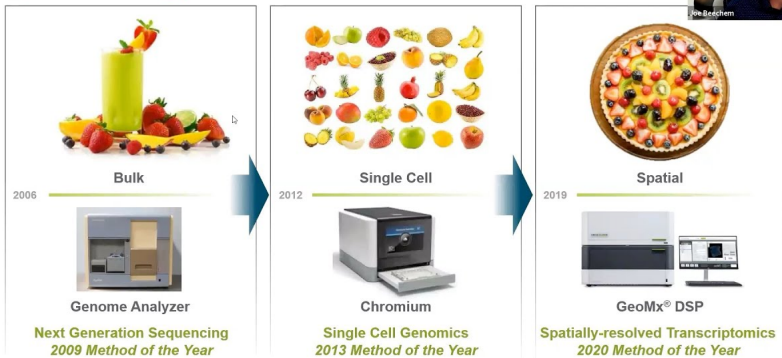


Functional Genomics

Applications in Research & Clinical Genomics

- Whole-genome/exome sequencing
- RNA-seq for transcriptomics
- Metagenomics for microbial community profiling
- Epigenomics (whole genome bisulfite sequencing, ATAC-seq)
- Structural variants and long-read phasing

Spatial Biology is the Next Life Sciences Revolution



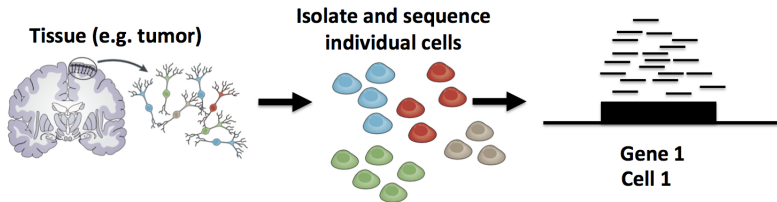
Adapted from concept and images by Dr. Aviv Regev of The Broad Institute; *Nature Methods*

nanoString

Single-Cell Sequencing: Overview

- Measures transcriptomic or epigenomic profiles at **single-cell resolution**.
- Enables discovery of:
 - Cell types and subtypes
 - Developmental trajectories
 - Cellular heterogeneity in disease
- Common platforms:
 - Droplet-based (10x Genomics)
 - Plate-based (SMART-seq2)
- Typical output: **cell** \times **gene** count matrix (usually sparse).

Single-cell RNA-Seq (scRNA-Seq)



Read Counts

	Cell 1	Cell 2	...
Gene 1	18	0	
Gene 2	1010	506	
Gene 3	0	49	
Gene 4	22	0	
...			

Compare gene expression profiles of single cells

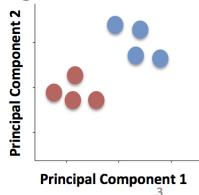
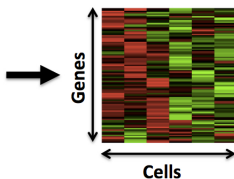


Figure adapted from <https://speakerdeck.com/stephaniehicks/welcome-to-the-world-of-single-cell-rna-sequencing>

Spatial Transcriptomics

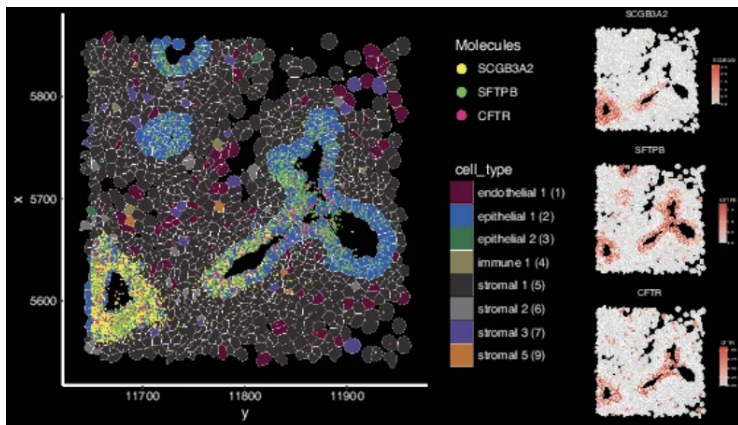


Figure adapted from Quach H, et al. Early human fetal lung atlas reveals the temporal dynamics of epithelial cell plasticity. Nat Commun 15: 5898 (2024). doi: 10.1038/s41467-024-50281-5

Spatial Transcriptomics: Overview

- Measures gene expression while preserving **spatial context** in tissues.
- Key question: *Where are gene programs active within the tissue microenvironment?*
- Major platform categories:
 - **Spot-based sequencing platforms**
 - 10x Visium
 - Slide-seq / Slide-seqV2
 - **High-resolution imaging-based platforms**
 - **10x Xenium** (RNA in situ imaging, subcellular resolution)
 - **NanoString CosMx SMI** (Spatial Molecular Imager; single-cell/subcellular resolution)
 - MERFISH, seqFISH
- Output varies by platform:
 - Spot-based: **spot** \times **gene** counts (multi-cell)
 - Imaging-based: **cell** \times **gene** or **subcellular transcripts with coordinates**

So ... what about analytical approaches?



Much to learn, you still have

Popular neural network architectures

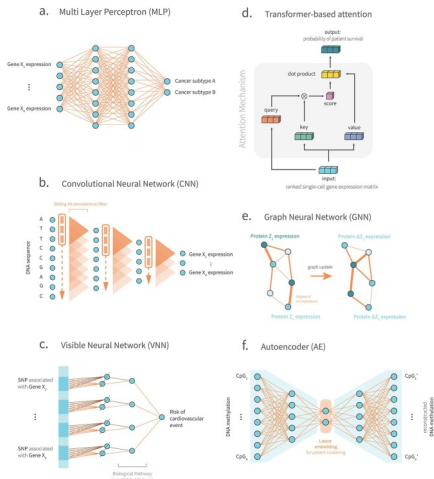
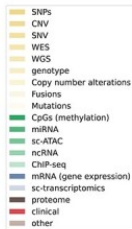


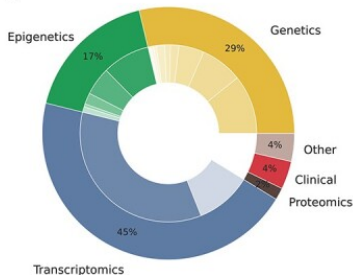
Figure adapted from Van Hilten, A., Katz, S., Saccenti, E., Niessen, W. J., & Roshchupkin, G. V. (2024). Designing interpretable deep learning applications for functional genomics: a quantitative analysis. *Briefings in Bioinformatics*, 25(5).

Data Types and Biological Fields

Characteristics of the input data



a.



b.

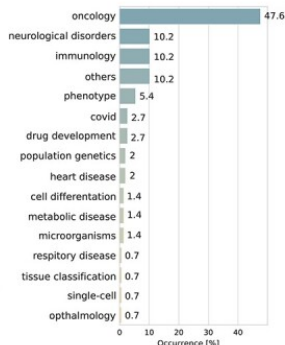


Figure adapted from Van Hilten, A., Katz, S., Saccenti, E., Niessen, W. J., & Roshchupkin, G. V. (2024). Designing interpretable deep learning applications for functional genomics: a quantitative analysis. *Briefings in Bioinformatics*, 25(5).

What makes genomic data unique?

- **Very large:** single-cell & individual genomes are gigabytes; cohorts are terabytes.
- **Hierarchical:** bases \rightarrow genes \rightarrow chromosomes \rightarrow genomes.
- **Complex formats:** multiple specialized file formats and metadata.
- **Sensitive:** potential to identify individuals — ethical and legal considerations.

