

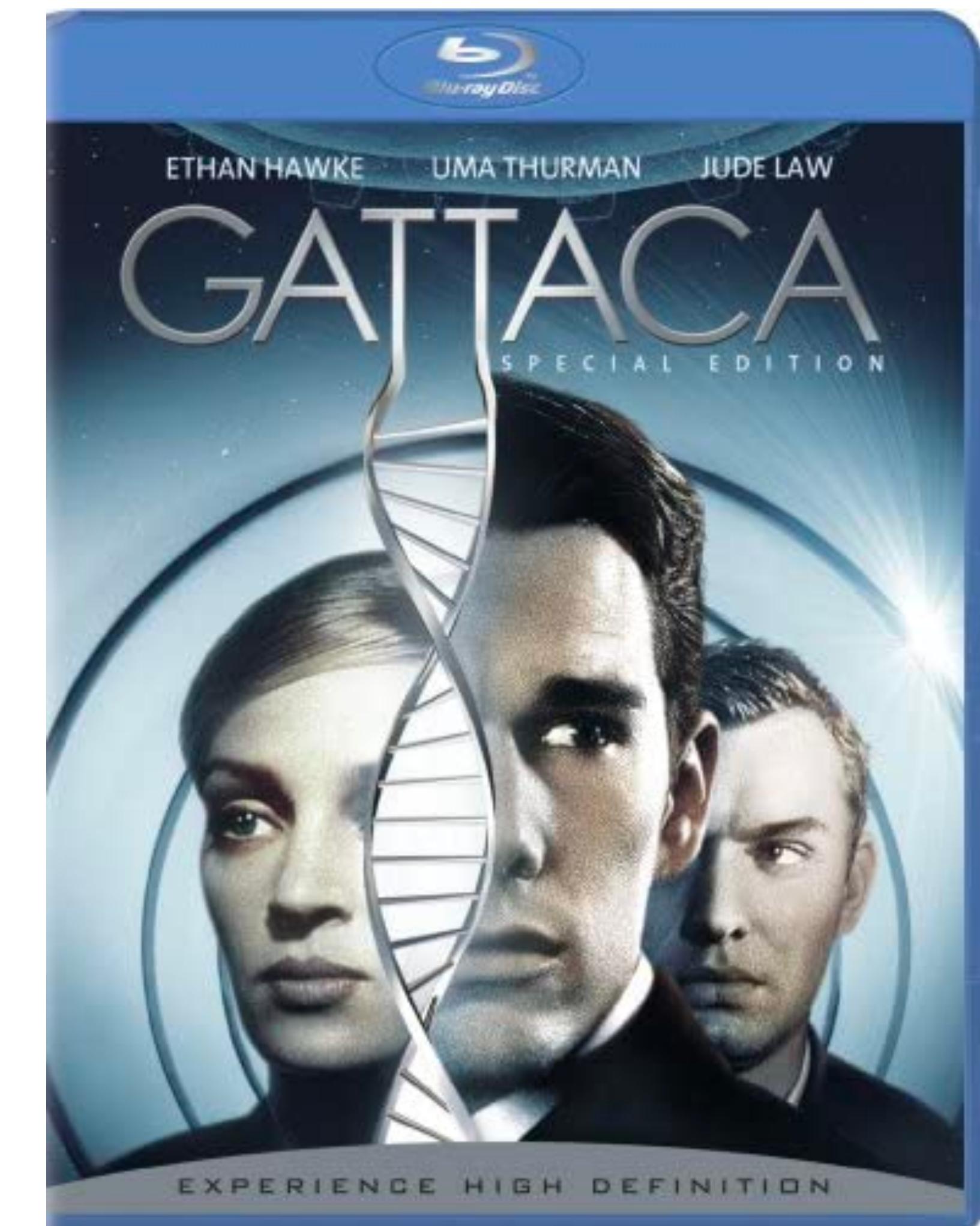
# Bioinformatics analysis for next generation sequencing

Hirak Sarkar  
Princeton University  
[www.hiraksarkar.com](http://www.hiraksarkar.com)

# Tools for measuring the “sight” and “sound” of biology (Genome and Genomics)

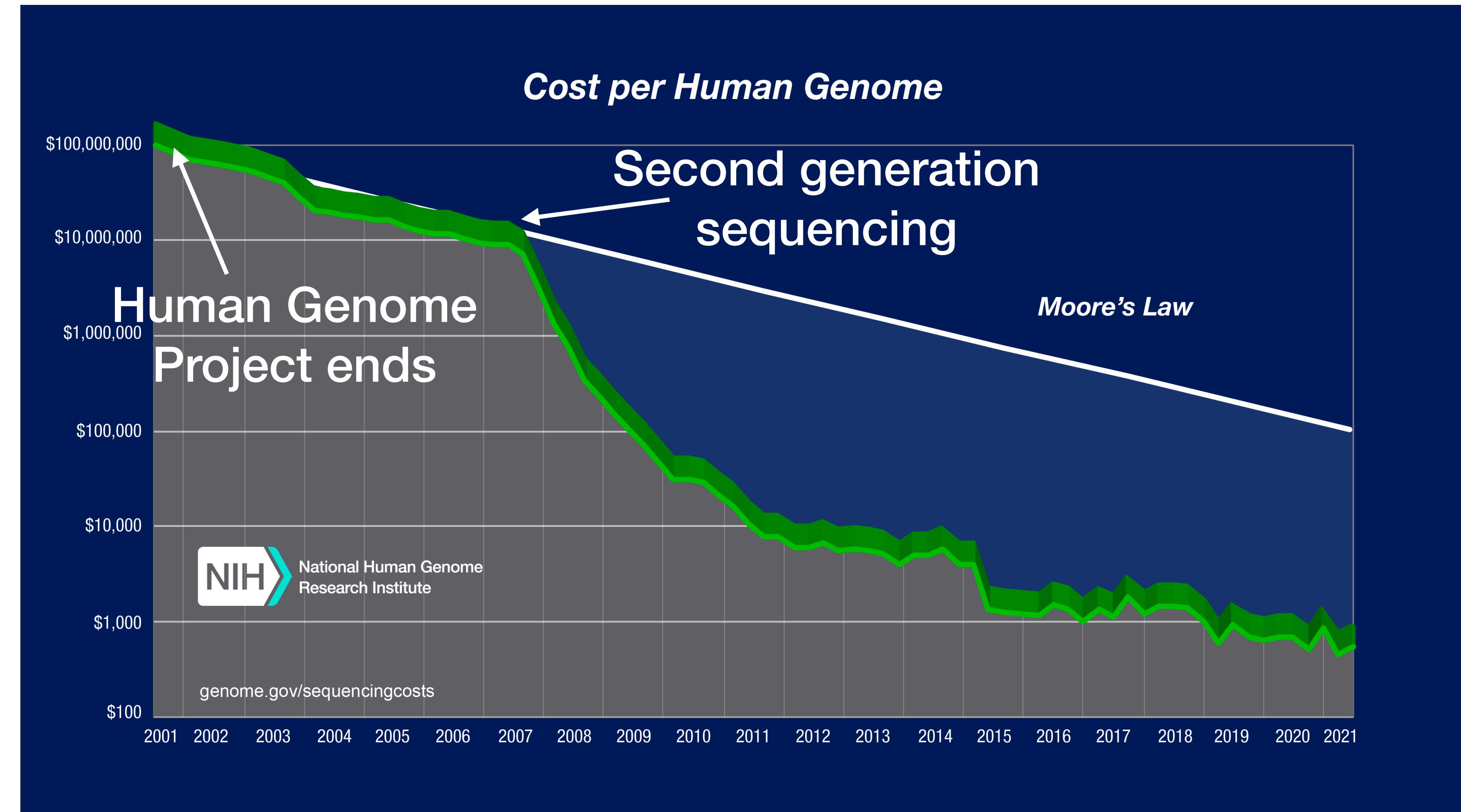
1997

1993

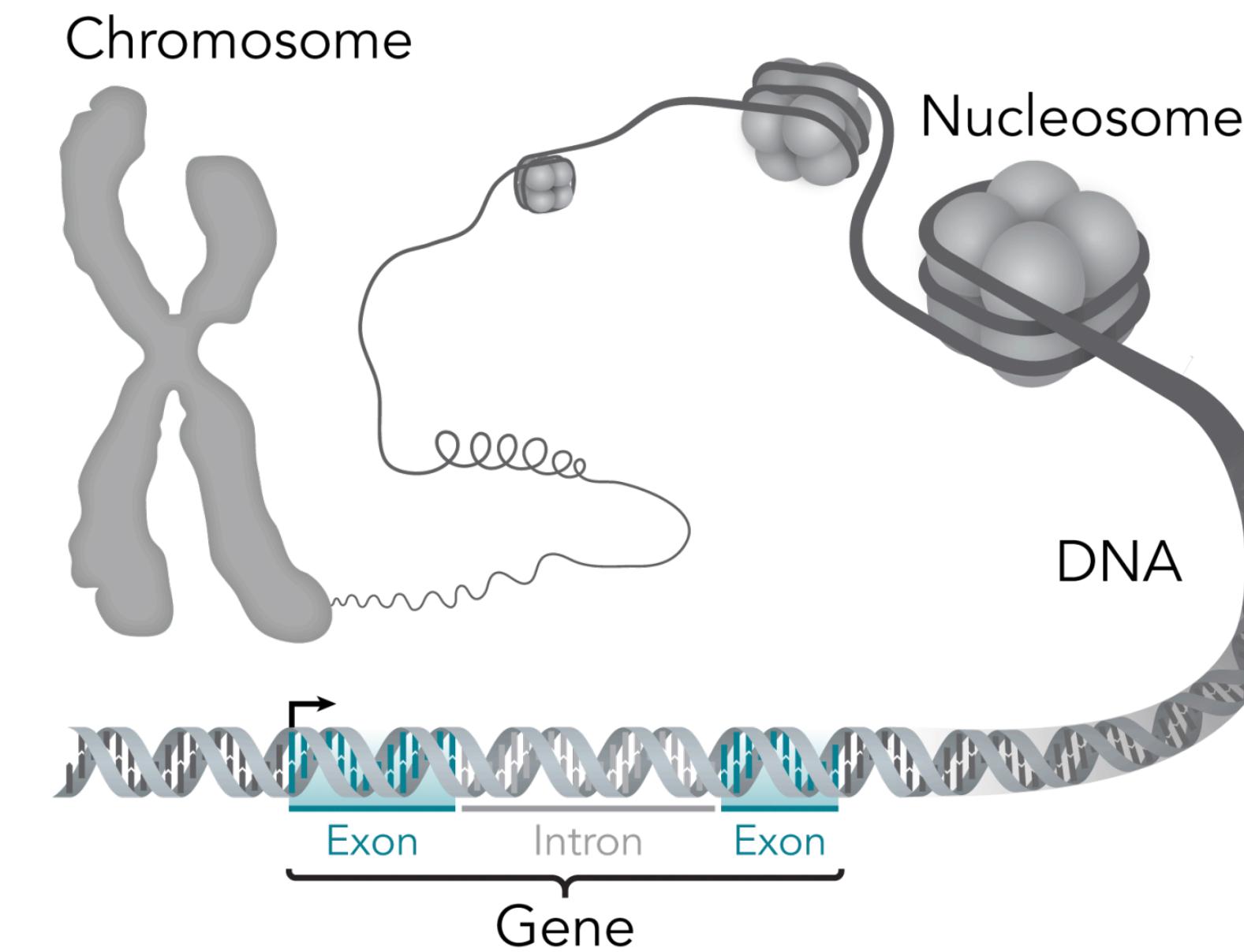
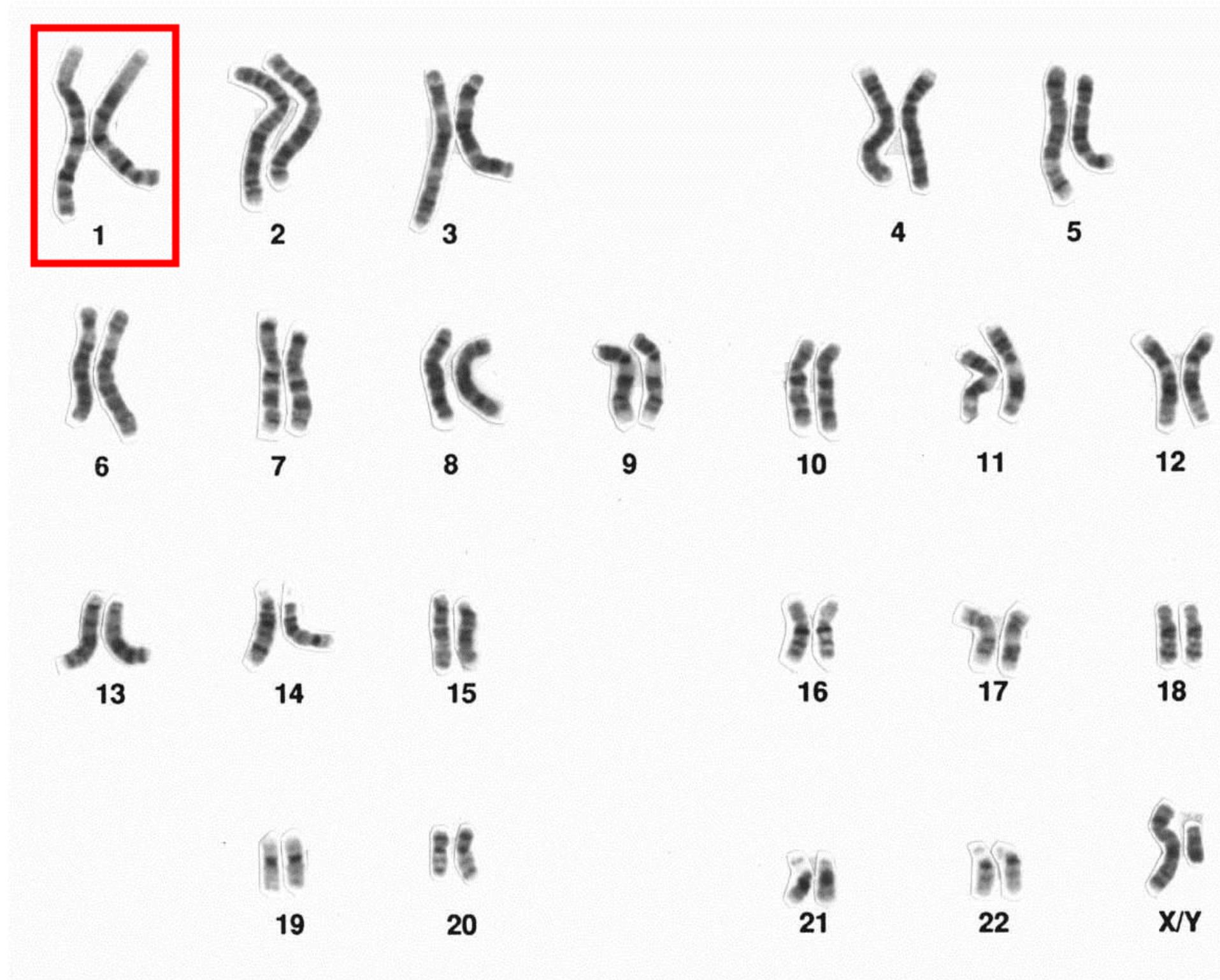


# Tools for measuring the “sight” and “sound” of biology (Genome and Genomics)

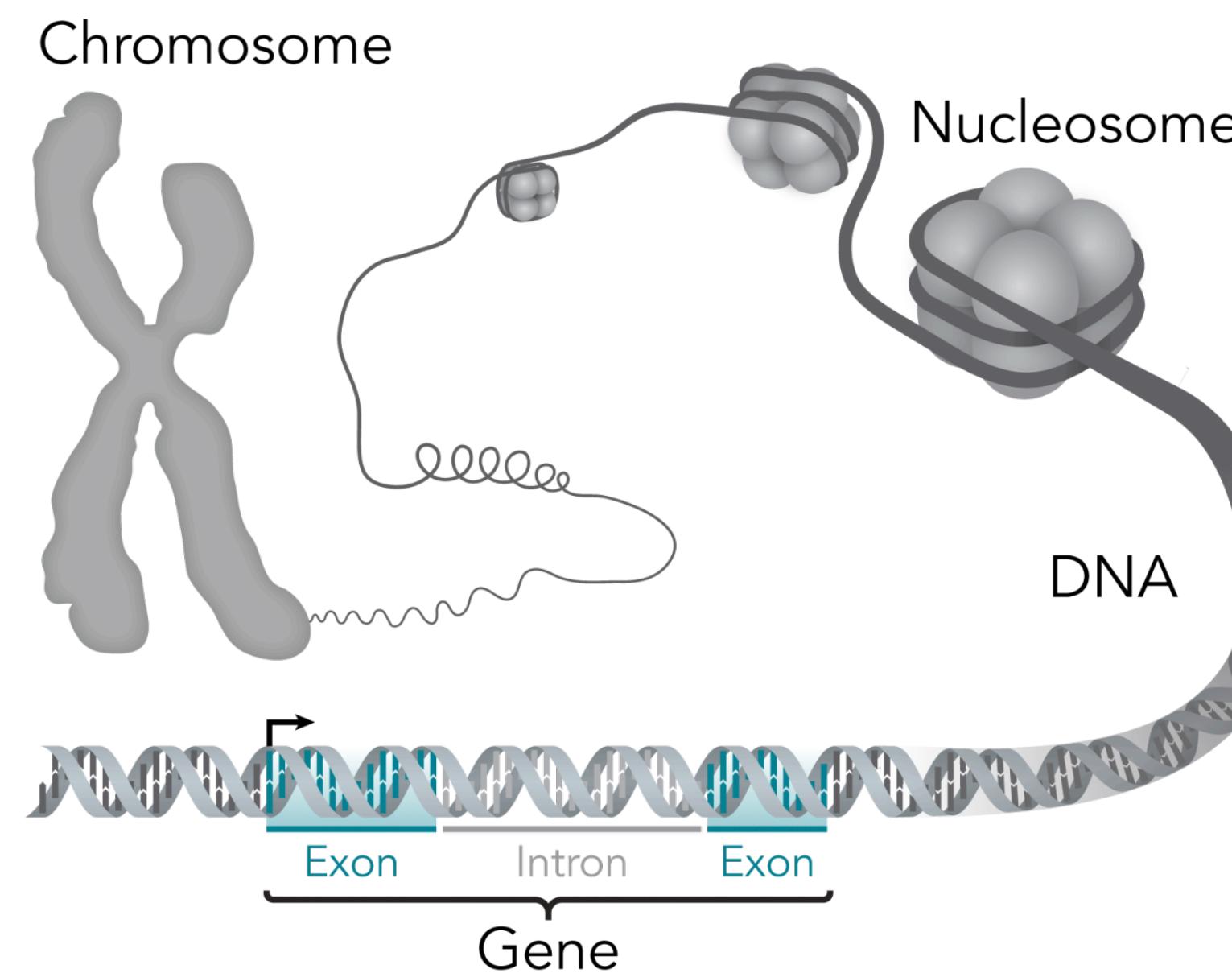
2001



# Tools for measuring the “sight” and “sound” of biology (Genome and Genomics)



# Tools for measuring the “sight” and “sound” of biology (Genome and Genomics)



## NovaSeq X series specifications

Output Range ~165 Gb - 16 Tb

Single reads per run 1.6 billion - 52 billion

Read length 2 × 150 bp

Run time ~13 hr - 48 hr

[View All NovaSeq X Specs](#)

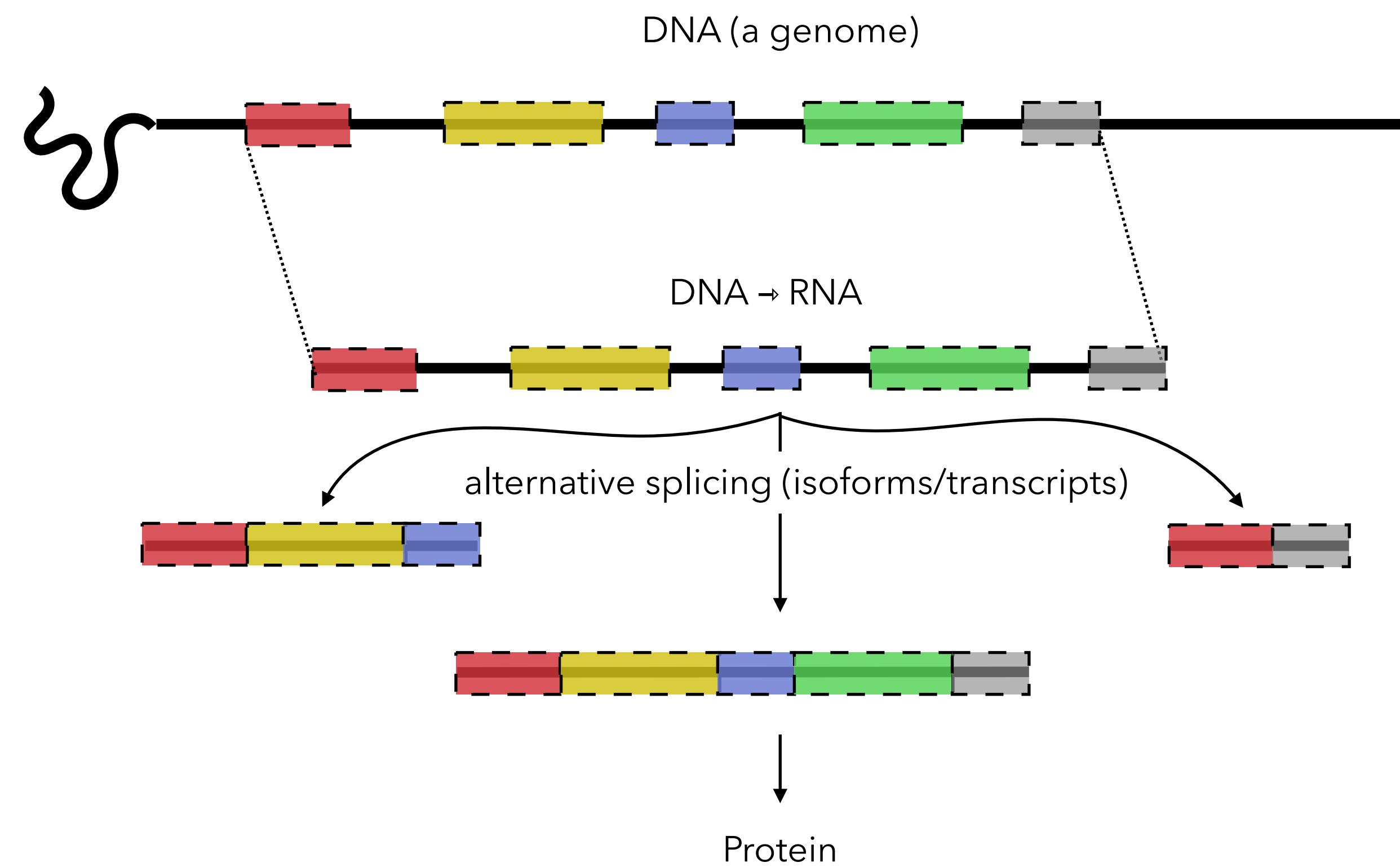
View AR

ATAGCAAGCTTT

# How to measure gene expression

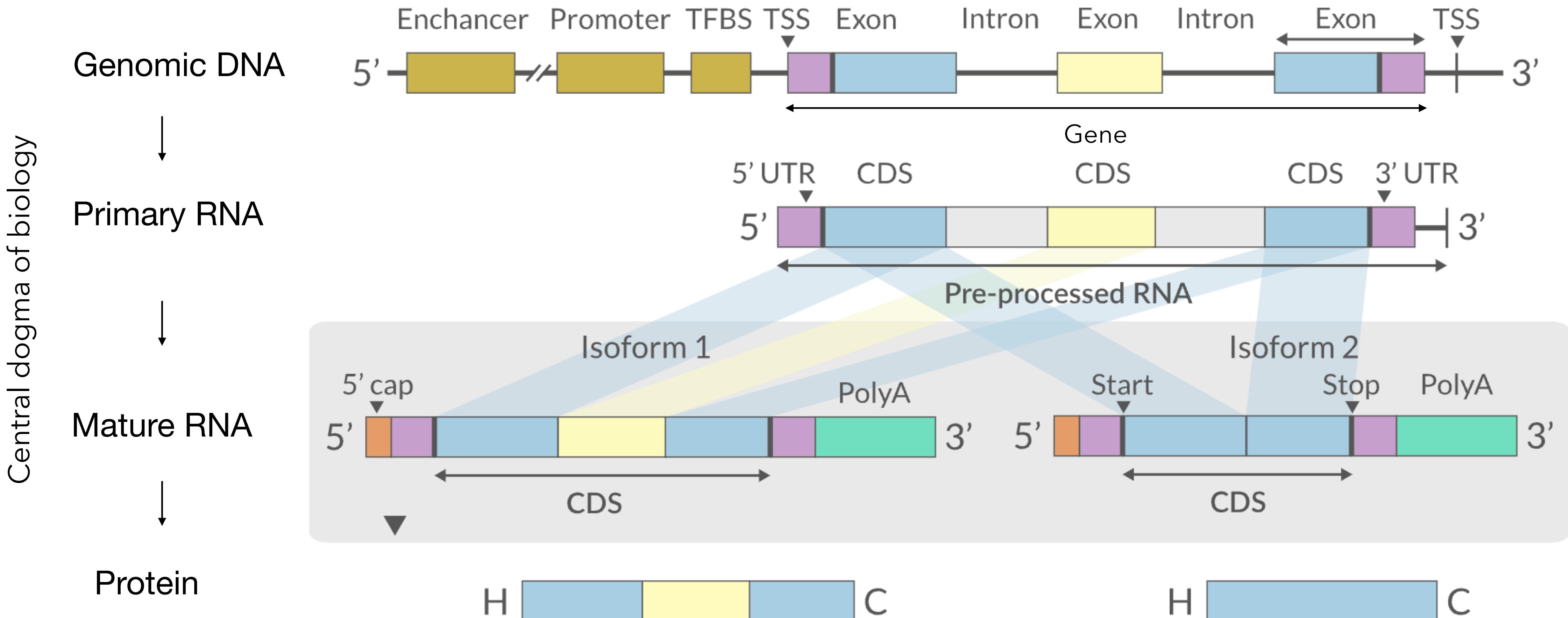
## RNA Sequencing (simplified)

Central dogma of biology



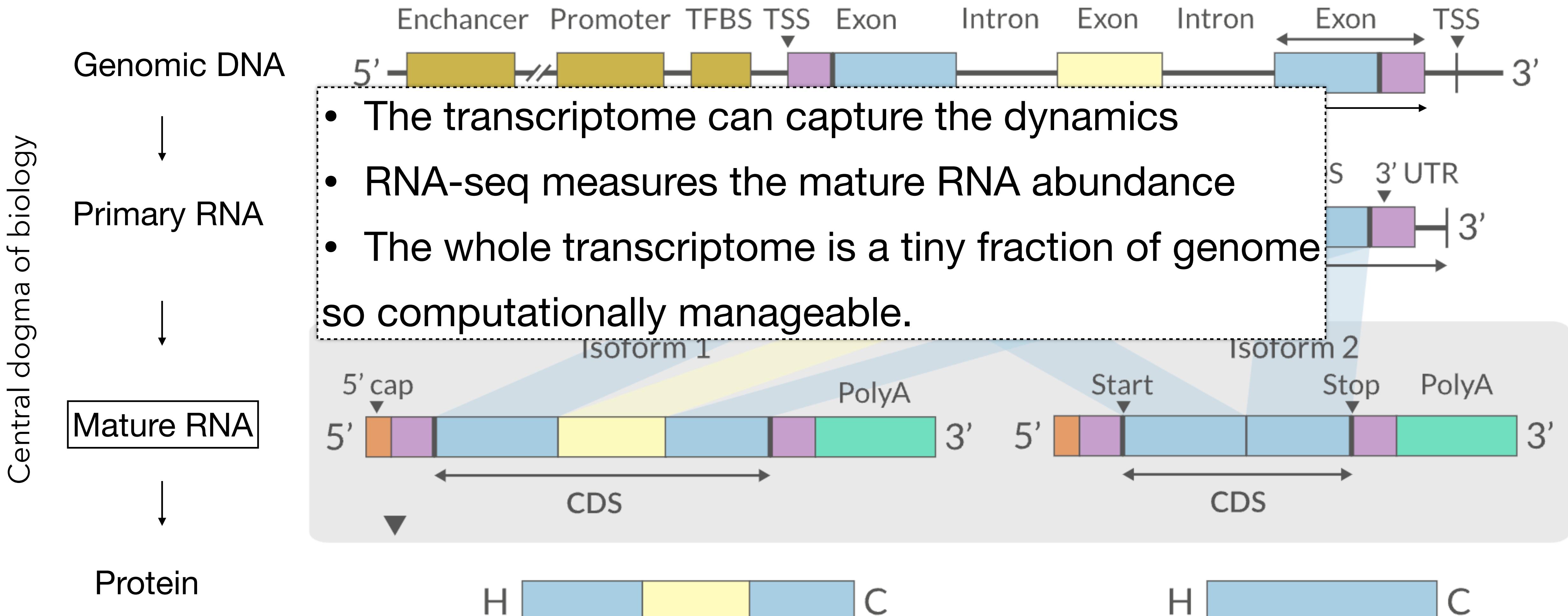
# How to measure gene expression

## RNA Sequencing (*not simplified*)



# How to measure gene expression

## RNA Sequencing (*not simplified*)

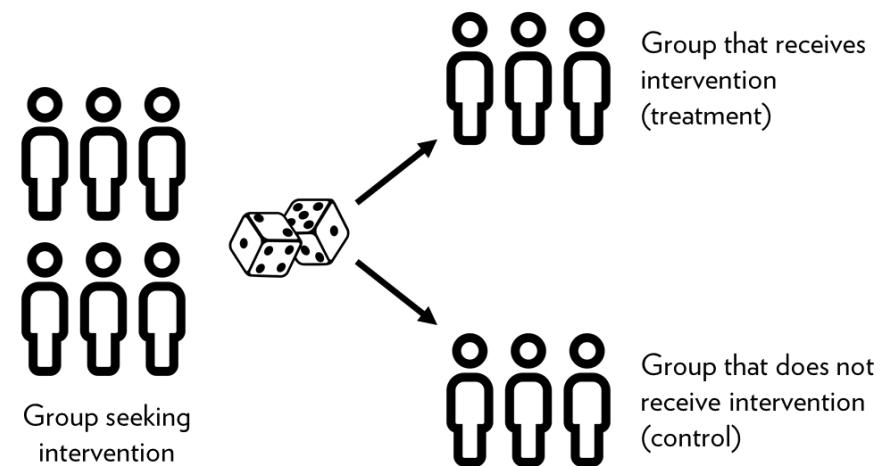
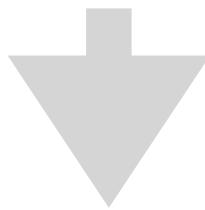


# What can we do with RNA-seq

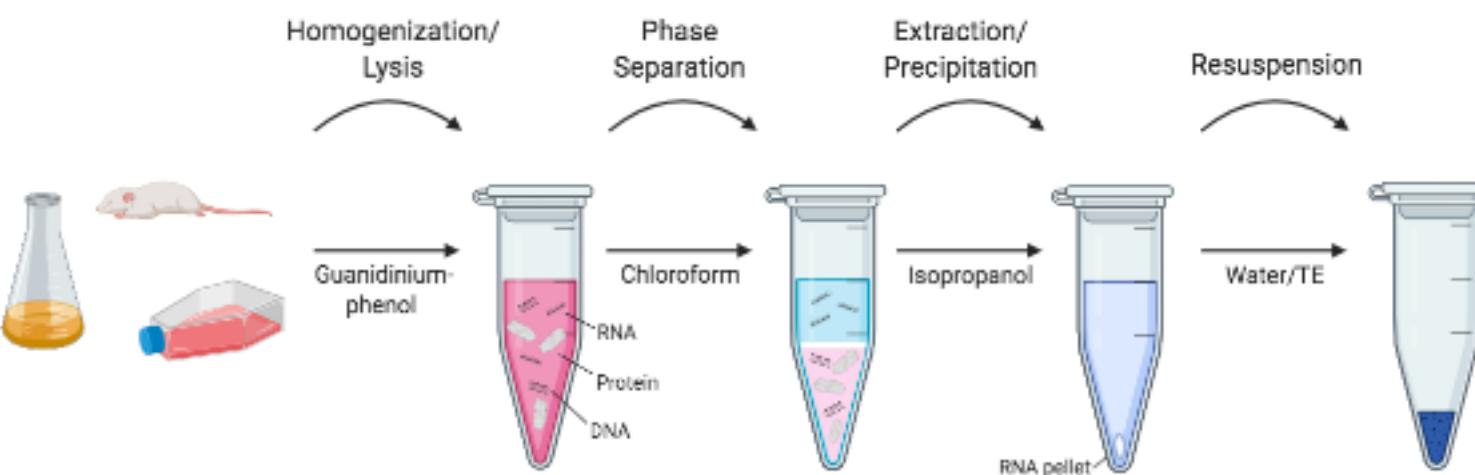
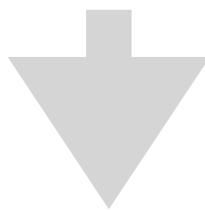
- Measure the expression of genes across different conditions.
- Understanding function of gene.
- Differential gene expression analysis.
  - Measure the degree of change in different conditions (e.g. a sample with special disease).
  - Discover pathways, gene networks, co-expression.
  - Fusion detection from genes
  - Much much more ...

# How can we do such exciting stuff

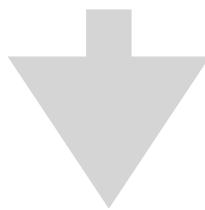
Experimental design



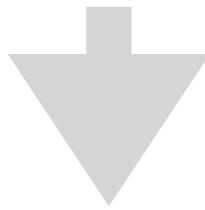
Experimental design



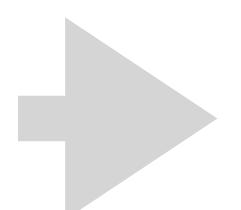
Library preparation



Sequencing

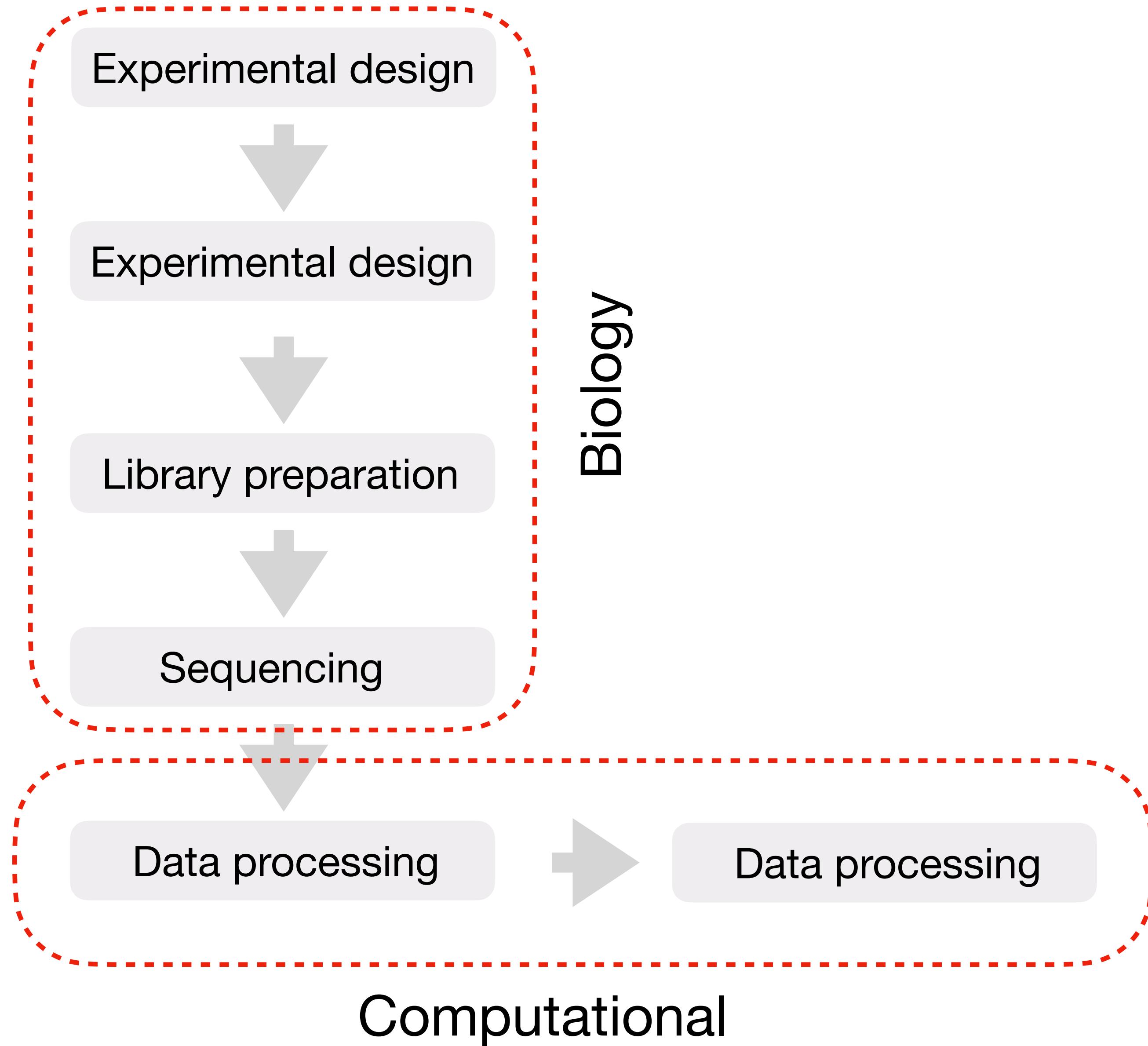


Data processing

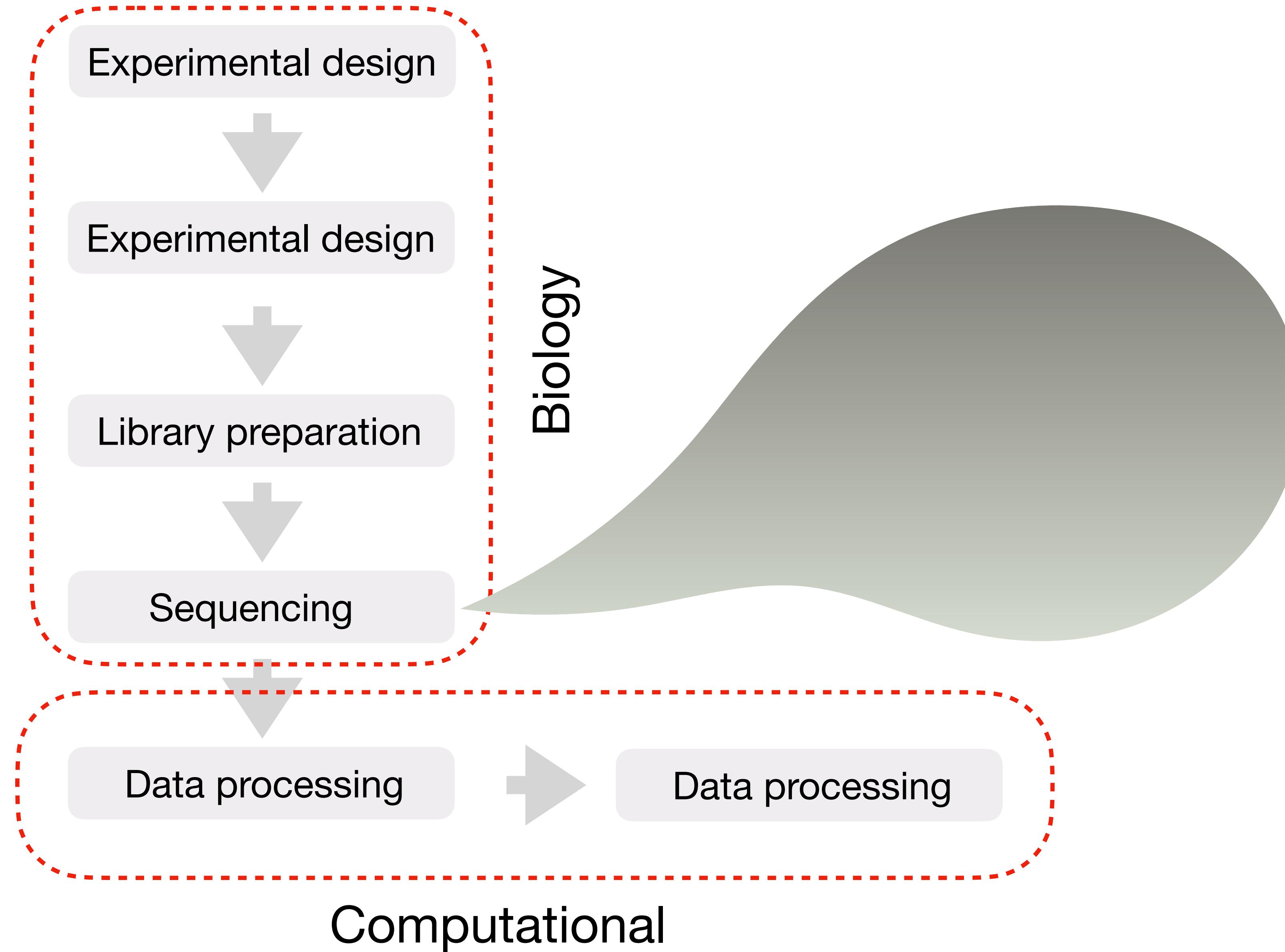


Data processing

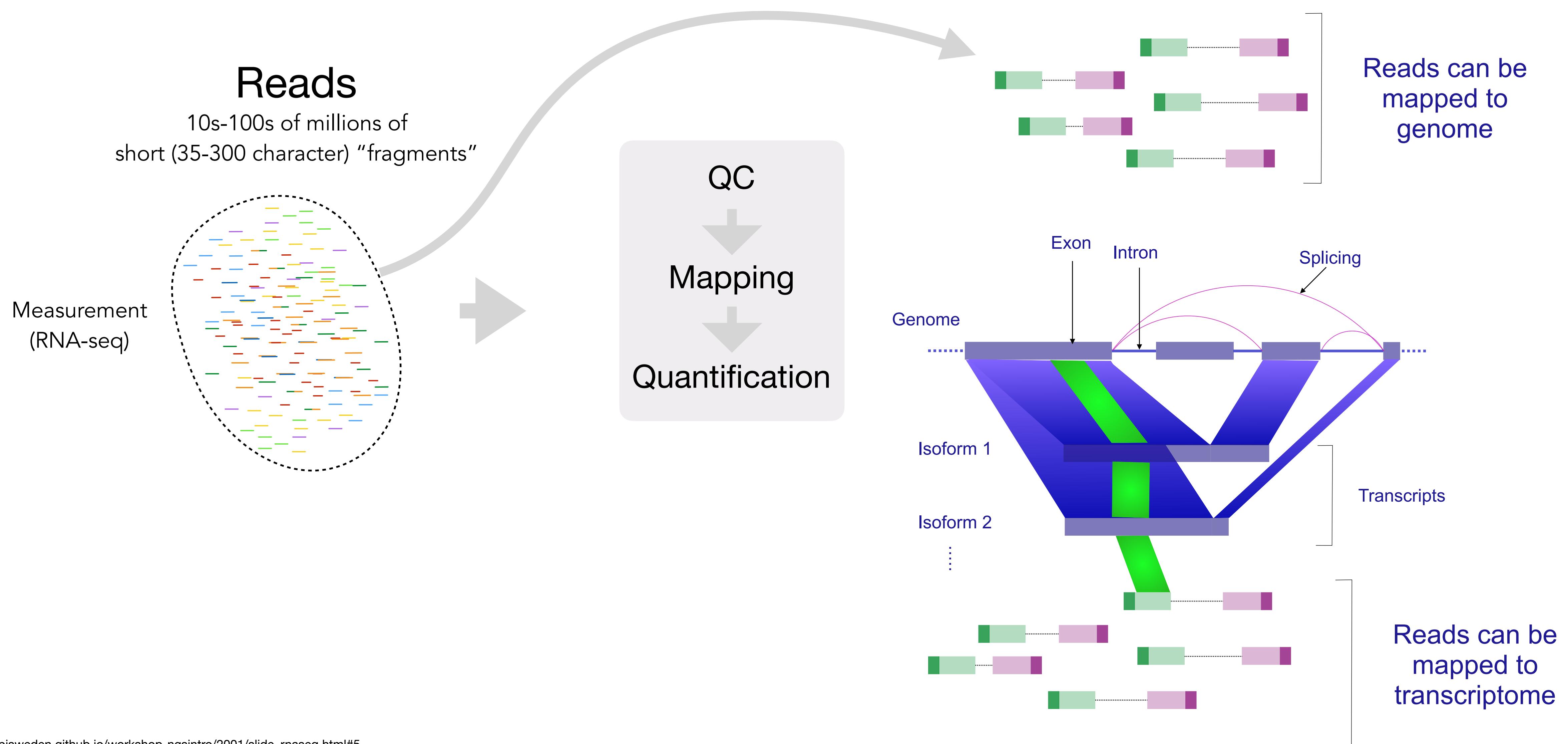
# How can we do such exciting stuff



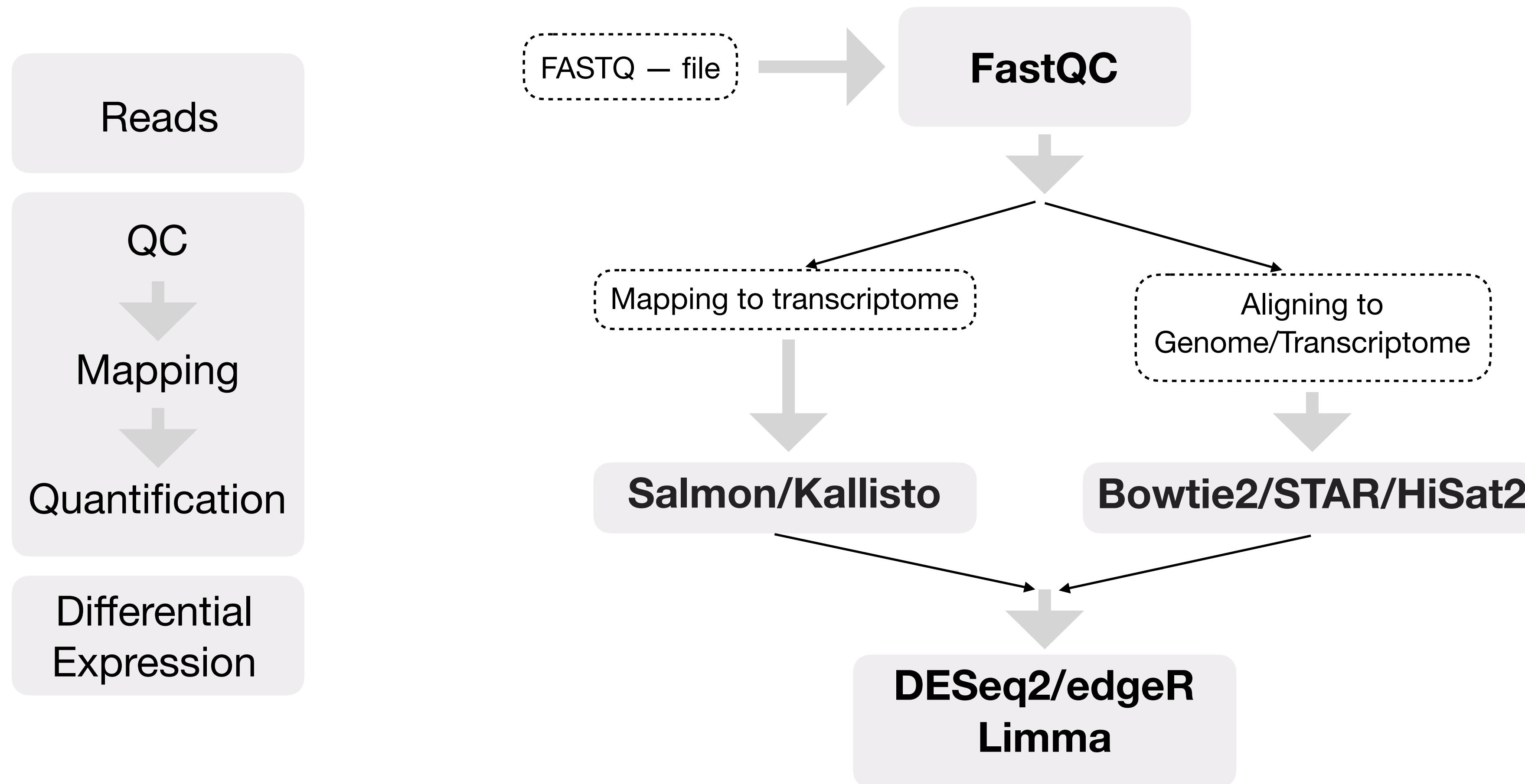
# How can we do such exciting stuff



# Basic workflow for RNA-seq (or any sequence based)

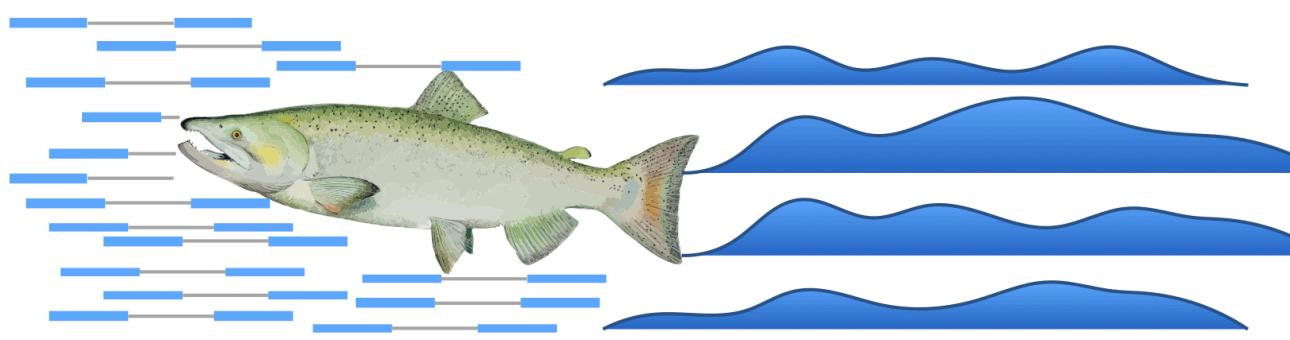


# Basic workflow for RNA-seq (or any sequence based)



# Tools we will use/learn today

Salmon



ostrokach/**gseapy**

Gene Set Enrichment Analysis in Python



# Acknowledgments

- [https://rpubs.com/qshenfeng/DESeq2 analysis tutorial](https://rpubs.com/qshenfeng/DESeq2_analysis_tutorial)
- [https://rpubs.com/qshenfeng/bulk RNAseq](https://rpubs.com/qshenfeng/bulk_RNAseq)
- [https://biocorecrg.github.io/PHINDaccess RNAseq 2020/salmon.html](https://biocorecrg.github.io/PHINDaccess_RNAseq_2020/salmon.html)
- <https://www.hadriengourle.com/tutorials/rna/>
- [https://github.com/mousepixels/sanbomics scripts/blob/main/  
PyDeseq2 DE tutorial.ipynb](https://github.com/mousepixels/sanbomics_scripts/blob/main/PyDeseq2_DE_tutorial.ipynb)
- <https://www.hadriengourle.com/tutorials/rna/>
- Documentation pages from softwares

# Real big data – (Song of Cells)\*

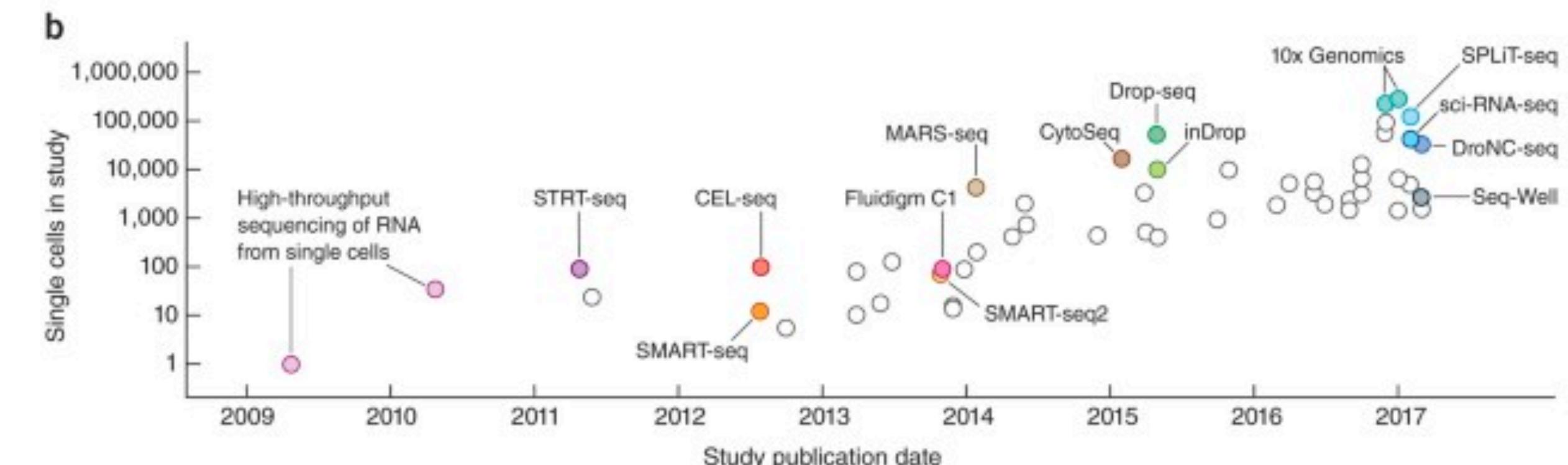
nature.com/nmeth

January 2020 Vol. 17 No. 1

## nature methods



Localization microscopy twice as precise  
A cryo-EM-based structural proteomics approach  
Time-resolved crystallography at the European XFEL  
Magnetic resonance at high speed



### U.S. Single-cell Analysis Market

size, by application, 2020 - 2030 (USD Billion)

