

---

# 3000788 Intro to Comp Molec Biol

Fall 2025



**Sira Sriswasdi, PhD**

- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

# Please introduce yourself

---

- Name (how should I call you)
- What do you want to learn from this course?
- Undergraduate background
- Graduate program and research interest
- Thesis advisor & topic (if already picked)

---

# 3000788 Intro to Comp Molec Biol

## Lecture 1: Computational thinking in biology

Fall 2025



**Sira Sriswasdi, PhD**

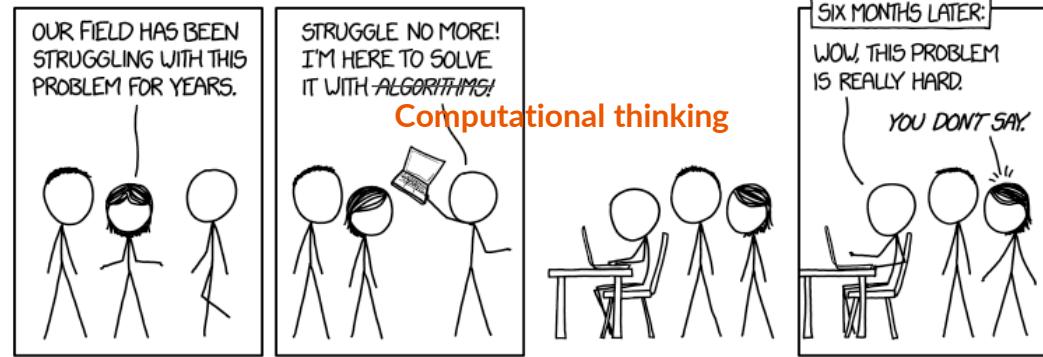
- Research Affairs
- Center of Excellence in Computational Molecular Biology (CMB)
- Center for Artificial Intelligence in Medicine (CU-AIM)

# About this course

---

## The Instructor

- Sira Sriswasdi (สิรະ ศรีสวัสดิ์)
- BSc in Mathematics
- PhD in Computational Biology
- Apply computational thinking and tools to answer biological and biomedical problems
- Machine learning and AI in medicine



Credit: "Here to Help" from [xkcd comic](#), reprinted under [Creative Commons License](#)

## Course Objectives

- Intended for new students
- Cover broad topics
- **Provide biological motivation (why) but focus on computational aspect (how)**
- Hands-on assignments with instruction

# Structure

---

- 28 sessions
- 6 assignments
  - 15% each
- Can work in group but write your own
- Don't abuse AI, and credit AI if used
- Can consult me

Week	Session	Date	Day	Time	Description	Activity
1	1	11-Aug-25	Mon		Holiday	
	2	13-Aug-25	Wed	10:00-11:30	Course introduction + history of computation in biology	
2	3	18-Aug-25	Mon	10:00-11:30	Sequence alignment	
	4	20-Aug-25	Wed	10:00-11:30	DNA sequencing techniques and applications	Assignment 1
3	5	25-Aug-25	Mon	10:00-11:30	Sequencing data processing workflows	
	6	27-Aug-25	Wed	10:00-11:30	Genome assembly and annotation	
4	7	1-Sep-25	Mon	10:00-11:30	Variant calling and genome-wide association study	
	8	3-Sep-25	Wed	10:00-11:30	Phylogenetics and evolutionary models	Assignment 2
5	9	8-Sep-25	Mon	10:00-11:30	Metagenomics	
	10	10-Sep-25	Wed	10:00-11:30	ChIP-seq and DNA motif detection	
6	11	15-Sep-25	Mon	10:00-11:30	Transcriptomics techniques	
	12	17-Sep-25	Wed	10:00-11:30	Differential expression analysis	Assignment 3
7	13	22-Sep-25	Mon	10:00-11:30	Functional enrichment analysis	
	14	24-Sep-25	Wed	10:00-11:30	Single-cell and spatial techniques	
8	15	29-Sep-25	Mon	10:00-11:30	Single-cell data processing	
	16	1-Oct-25	Wed	10:00-11:30	Proteomics and mass spectrometry	
9	17	6-Oct-25	Mon	10:00-11:30	Introduction to systems biology and dynamics	Assignment 4
	18	8-Oct-25	Wed	10:00-11:30	Multi-omics integration	
10	19	13-Oct-25	Mon		Holiday	
	20	15-Oct-25	Wed	10:00-11:30	Biological networks	
11	21	20-Oct-25	Mon	10:00-11:30	Chromatin conformation capture	Assignment 5
	22	22-Oct-25	Wed	10:00-11:30	RNA and protein structure models	
12	23	27-Oct-25	Mon	10:00-11:30	Online tools and resources	
	24	29-Oct-25	Wed	10:00-11:30	Microscopy data analysis	
13	25	3-Nov-25	Mon	10:00-11:30	Essential statistics for computational biologist	Assignment 6
	26	5-Nov-25	Wed	10:00-11:30	Omics data visualization	
14	27	10-Nov-25	Mon	10:00-11:30	Introduction to machine learning and AI	
	28	12-Nov-25	Wed	10:00-11:30	Machine learning applications	
15	29	17-Nov-25	Mon	10:00-11:30	Foundational AI models in biology	
	30	19-Nov-25	Wed	10:00-11:30	Biomarker discovery	



# **Self-study resources**

# Companion courses from MIT

## FOUNDATIONS OF COMPUTATIONAL AND SYSTEMS BIOLOGY

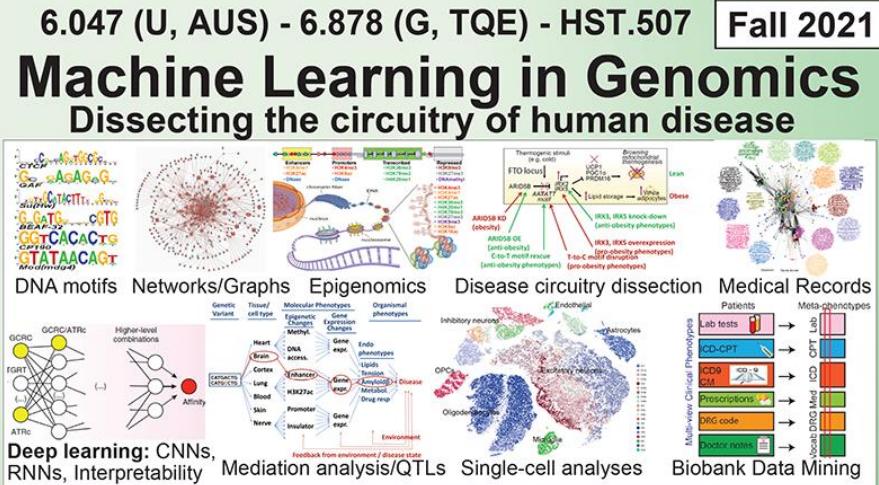
Lecture 1: Introduction to Computational and Systems Biology

Lecture 2: Local Alignment (BLAST) and Statistics

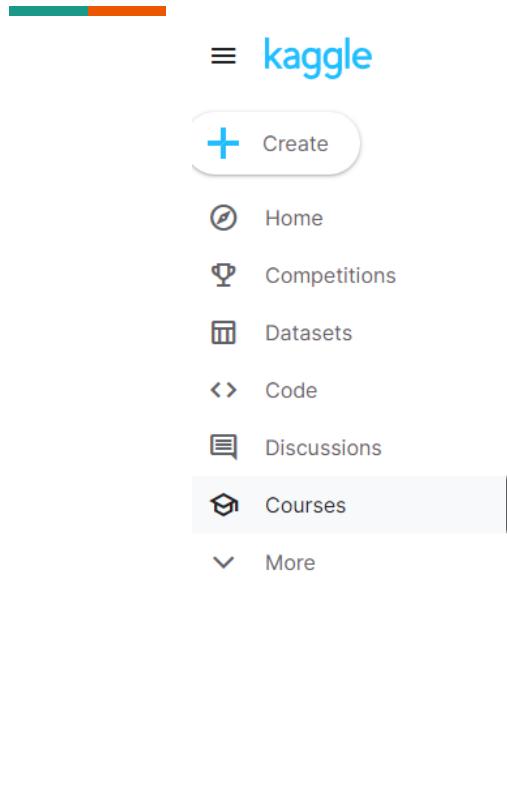
Lecture 3: Global Alignment of Protein Sequences (NW, SW, PAM, BLOSUM)

Lecture 4: Comparative Genomic Analysis of Gene Regulation

Visit <https://web.mit.edu/manoli/>



# Kaggle's programming courses



The image shows the left sidebar of the Kaggle website. At the top is a teal horizontal bar with a red progress bar underneath. Below it is the 'kaggle' logo with three horizontal bars to its left. A 'Create' button with a plus sign is highlighted with a light orange rounded rectangle. The sidebar menu includes: Home (with a compass icon), Competitions (with a trophy icon), Datasets (with a bar chart icon), Code (with a code editor icon), Discussions (with a comment icon), Courses (with a graduation cap icon), and More (with a downward arrow icon). The 'Courses' item is currently selected and has a light orange background.



A search bar with a magnifying glass icon and the placeholder text 'Search'.

## Explore Courses

-  **Intro to Programming**  
Get started with Python, if you have no coding experience.
-  **Python**  
Learn the most important language for data science.
-  **Intro to Machine Learning**  
Learn the core ideas in machine learning, and build your first models.
-  **Pandas**  
Solve short hands-on challenges to perfect your data manipulation skills.
-  **Intermediate Machine Learning**  
Handle missing values, non-numeric values, data leakage, and more.
-  **Data Visualization**  
Make great data visualizations. A great way to see the power of coding!

# MIT 6.0001/6.0002



## Introduction To Computer Science And Programming In Python

*6.0001 Introduction to Computer Science and Programming in Python* is intended for students with little or no programming experience. It aims to provide students with an understanding of the role computation can play in solving problems and to help students, regardless of their major, feel justifiably confident of their ...[Show more](#)

## Introduction To Computational Thinking And Data Science

6.0002 is the continuation of *6.0001 Introduction to Computer Science and Programming in Python* and is intended for students with little or no programming experience. It aims to provide students with an understanding of the role computation can play in solving problems and to help students, regardless of their major, ...



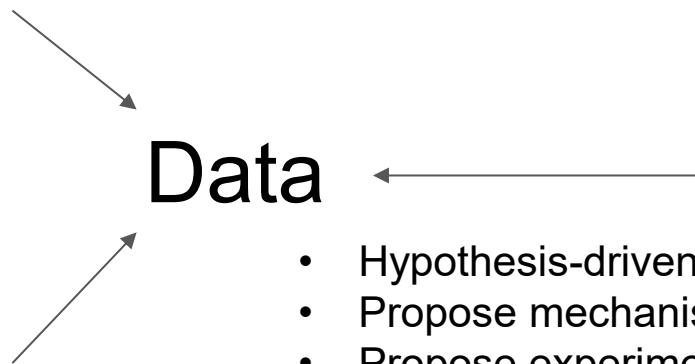
# Computational thinking

## A knowledge-centric view

---

Computational  
Skill

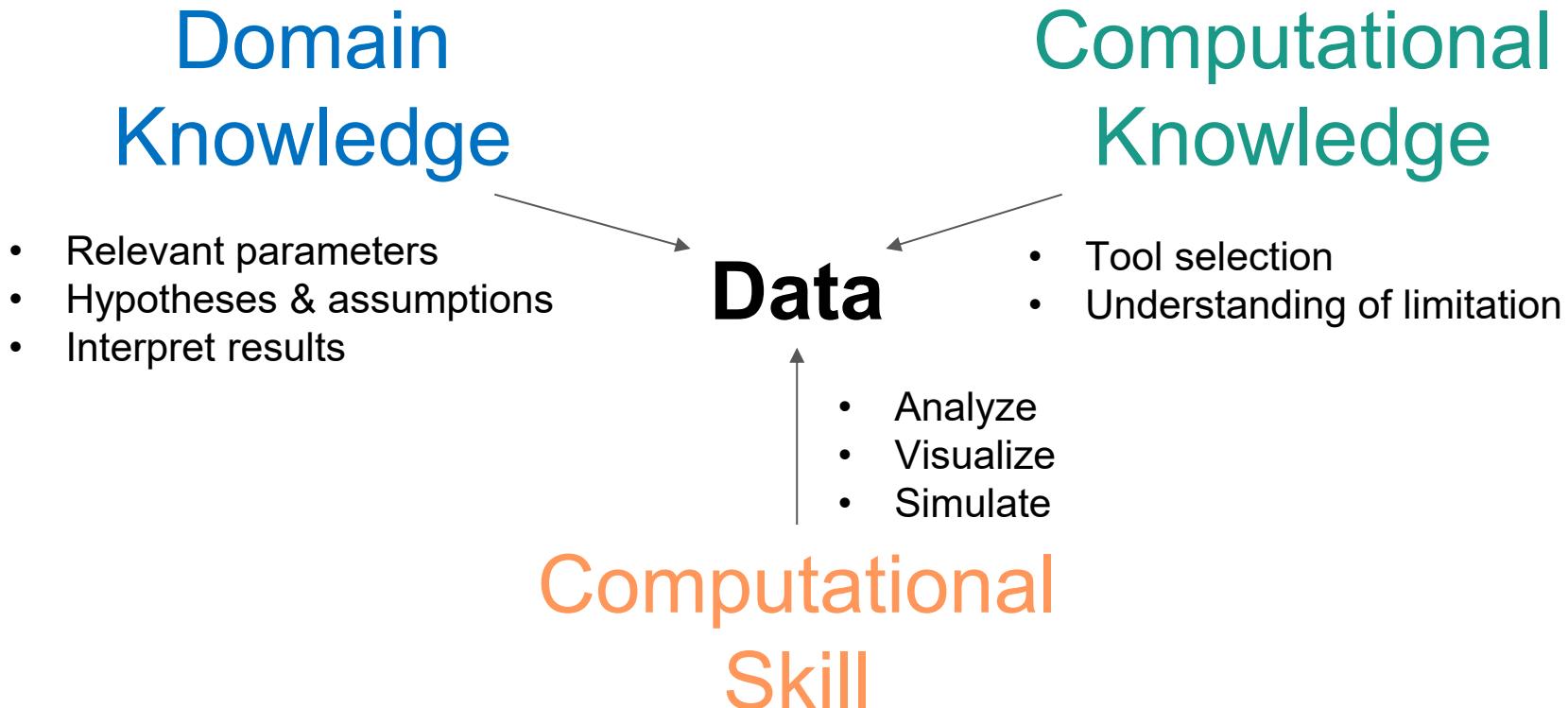
Computational  
Knowledge



Domain  
Knowledge

# A data-centric view

---



# Many tools for a task



- Computational knowledge is helpful for assessing tools and techniques
- Choose the right tools
- Know the limitation and assumption
- Also communicate with peers

## List of sequence alignment software

[Article](#) [Talk](#)

[Read](#) [Edit](#) [View history](#) [Tools](#)

 3 languages

From Wikipedia, the free encyclopedia

This list of sequence alignment software is a compilation of software tools and web portals used in pairwise [sequence alignment](#) and [multiple sequence alignment](#). See [structural alignment software](#) for structural alignment of proteins.

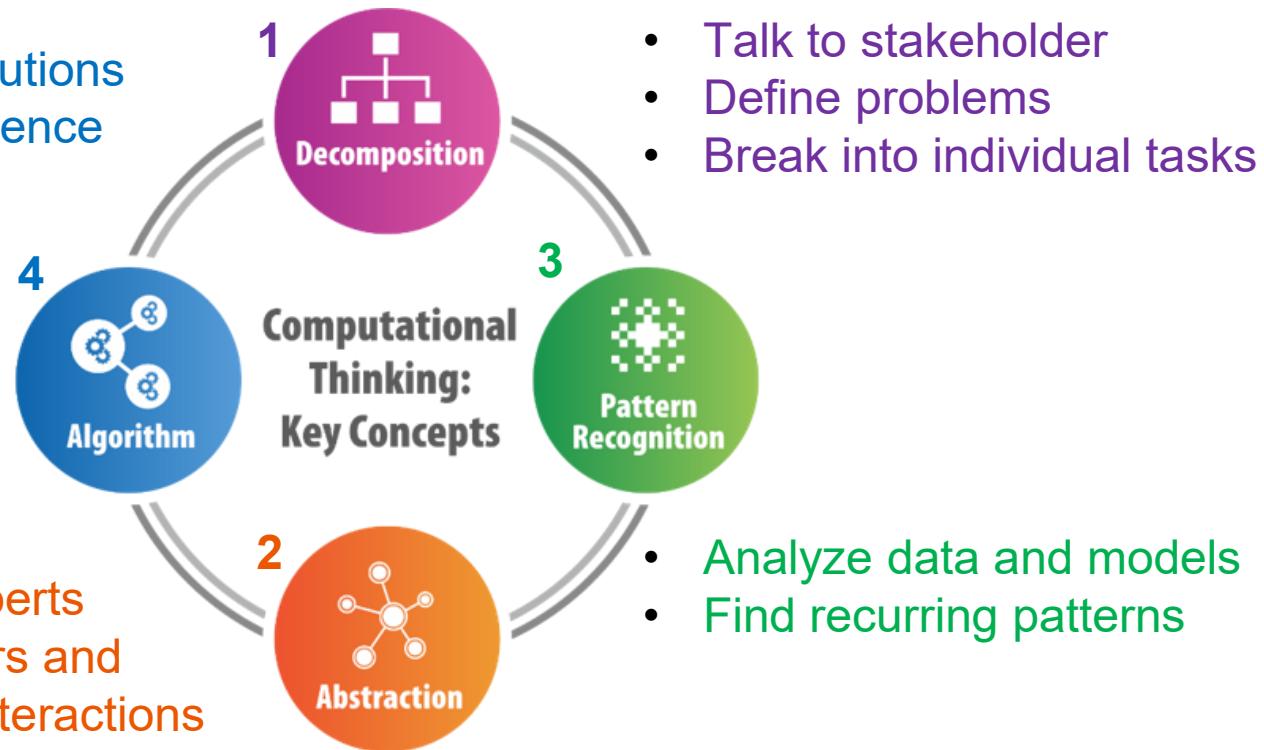
Database search only [\[edit\]](#)

Name	Description	Sequence type*	Authors	Year
BLAST	Local search with fast k-tuple heuristic (Basic Local Alignment Search Tool)	Both	Altschul SF, Gish W, Miller W, Myers EW, Lipman DJ. <sup>[1]</sup>	1990
HPC-BLAST	NCBI compliant multimode and multicore BLAST wrapper. Distributed with the latest version of BLAST, this wrapper facilitates parallelization of the algorithm on modern hybrid architectures with many nodes and many cores within each node. <sup>[2]</sup>	Protein	Burdyshaw CE, Sawyer S, Horton MD, Brook RG, Rekappalli B	2017
CS-BLAST	Sequence-context specific BLAST, more sensitive than BLAST, FASTA, and SSEARCH. Position-specific iterative version CSI-BLAST more sensitive than PSI-BLAST	Protein	Angermueller C, Biegert A, Soeding J. <sup>[3]</sup>	2013
CUDASW++	GPU accelerated Smith Waterman algorithm for multiple shared-host GPUs	Protein	Liu Y, Maskell DL and Schmidt B	2009/2010
DIAMOND	BLASTX and BLASTP aligner based on double indexing	Protein	Buchfink B, Xie C, Huson DH, Reuter K, Drost HG. <sup>[4][5]</sup>	2015/2021
FASTA	Local search with fast k-tuple heuristic, slower but more sensitive than BLAST	Both		
GGSEARCH, GLSEARCH	Global:Global (GG), Global:Local (GL) alignment with statistics	Protein		
Genome Magician	Software for ultra fast local DNA sequence motif search and pairwise alignment for NGS data (FASTA, FASTQ).	DNA	Hepperle D ( <a href="http://www.sequentix.de">www.sequentix.de</a> )	2020
	Genoogle uses indexing and parallel processing techniques for			

# Key procedures in computational thinking

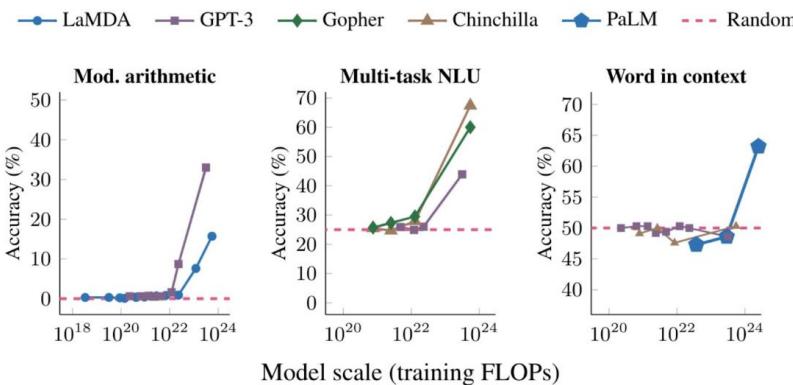
---

- Develop robust solutions
- Design user experience



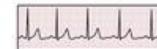
- Talk to domain experts
- Map out parameters and expected causal interactions

# Will large language model solve everything?



## Multimodal LLM

2. Interpret this ECG for indicators of myocardial infarction.



3. Provide a diagram showing the blocked coronary artery for this patient.

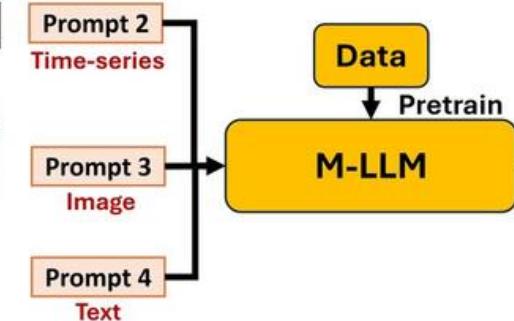


4. Create a rehabilitation plan video for a patient recovering from myocardial infarction.

Prompt 2  
Time-series

Prompt 3  
Image

Prompt 4  
Text



<https://hai.stanford.edu/news/examining-emergent-abilities-large-language-models>

Alsaad, R. et al. J of Medical Internet Research 2024

- LLM learns (mostly) correlation, not causation
- Require knowledge to verify
- Require mechanistic explanation to trust

# Danger of complexity

---



With four parameters I can fit an elephant, and with five I can make him wiggle his trunk.

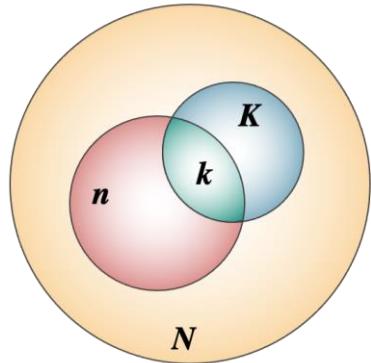
— *John von Neumann* —



# Statistical modeling

# Statistics are grounded in probability models

---



$$P(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$$

$\textcolor{brown}{N}$  = Background (e.g. Transcriptome  $\sim 40.000$  genes)

$\textcolor{red}{n}$  = Query list (e.g. upregulated genes)

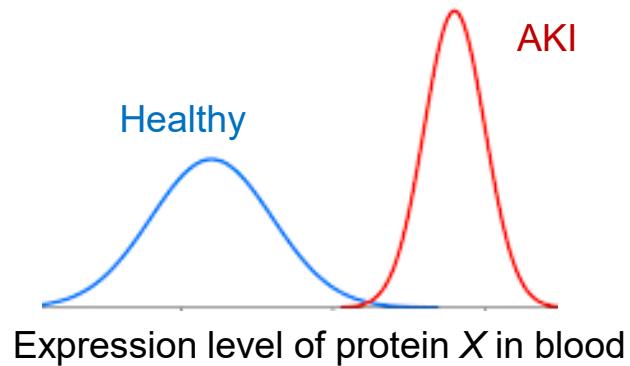
$\textcolor{teal}{K}$  = genes annotated in the pathway/set tested  
(e.g. Glycolysis)

$$\textcolor{teal}{k} = \textcolor{red}{n} \cap \textcolor{teal}{K}$$

- In a population  $N = 70$  million
  - There are estimated  $n = 1$  million with spontaneous nosebleed
  - There are estimated  $K = 20$  million exposed to high PM2.5
  - There are estimated  $k = 300,000$  with both
- How would you quantify the association between PM2.5 and nosebleed?
  - *Overrepresentation analysis with hypergeometric distribution*

# Unwrapping “canned” statistics

---



- How would you quantify the ability of  $X$  to distinguish healthy and AKI?
  - How about  $\text{Score}(X) = \text{Mean}_1 - \text{Mean}_2$  or  $\text{Abs}(\text{Mean}_1 - \text{Mean}_2)$  ?
  - How about  $\text{Score}(X) = \frac{\text{Mean}_1 - \text{Mean}_2}{\sqrt{\frac{1}{n}(\text{Variance}_1 + \text{Variance}_2)}} ?$
  - You have rediscovered *independent two-sample t-test* !!

# A peek behind popular statistical tests

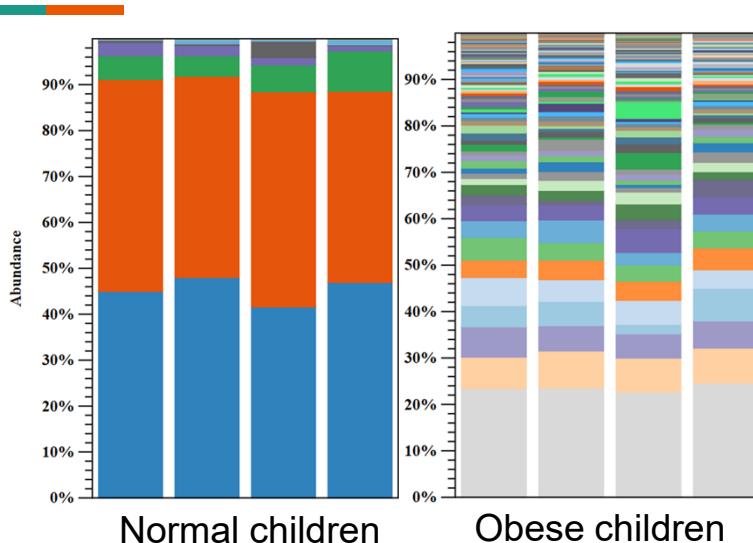
---

- **Mann-Whitney U test:**  
(aka Wilcoxon rank sum)

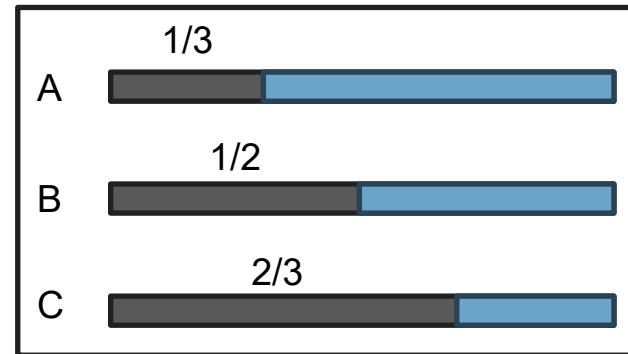
$$U = \sum_{i=1}^n \sum_{j=1}^m S(X_i, Y_j), \quad S(X, Y) = \begin{cases} 1, & \text{if } X > Y, \\ \frac{1}{2}, & \text{if } X = Y, \\ 0, & \text{if } X < Y. \end{cases}$$

- **Wilcoxon signed-rank test:**
  1. Compute  $|X_1|, \dots, |X_n|$ .
  2. Sort  $|X_1|, \dots, |X_n|$ , and use this sorted list to assign ranks  $R_1, \dots, R_n$
$$T = \sum_{i=1}^N \text{sgn}(X_i)R_i.$$
- **Sign test:** Each observation is equally likely to be + or -
  - $\Pr(k \text{ positive values from } N \text{ observations}) = \text{Binomial}(N, k, p = 0.5)$

# Turning verbal description into mathematical formula



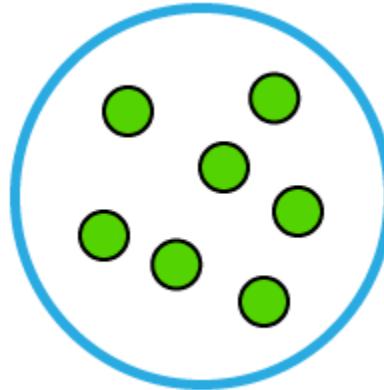
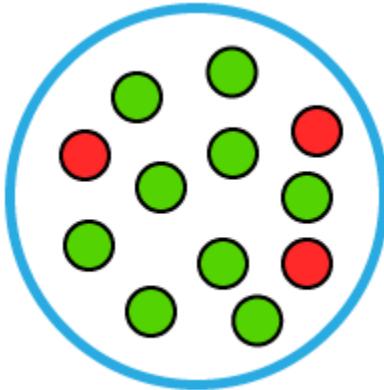
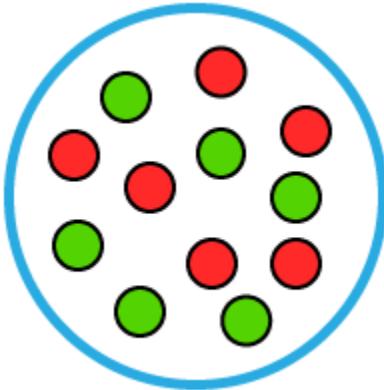
Colors indicate different microbial taxa in children gut microbiome



- How would you quantify the diversity of microbiome?
  - Number of different taxa =  $n$
  - Let  $p_1, \dots, p_n$  be the taxa frequency, how should diversity change with these parameters?
  - What should be the relationship between **score(A)**, **score(B)**, and **score(C)**?

# Entropy quantifies purity of a mixture

---



<https://www.javatpoint.com/entropy-in-machine-learning>

- **Entropy** =  $-p_1 \log_2(p_1) - p_2 \log_2(p_2) - \cdots - p_n \log_2(p_n)$

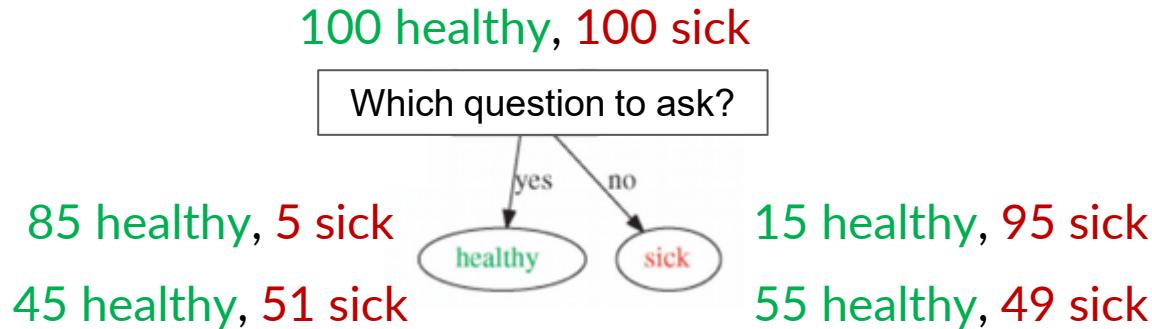
- Left:  $-\frac{6}{12} \log_2 \left(\frac{6}{12}\right) - \frac{6}{12} \log_2 \left(\frac{6}{12}\right) = \log_2(2) = 1$

- Middle:  $-\frac{3}{12} \log_2 \left(\frac{3}{12}\right) - \frac{9}{12} \log_2 \left(\frac{9}{12}\right) = 0.811278$

- Right:  $-\frac{0}{12} \log_2 \left(\frac{0}{12}\right) - \frac{12}{12} \log_2 \left(\frac{12}{12}\right) = 0$

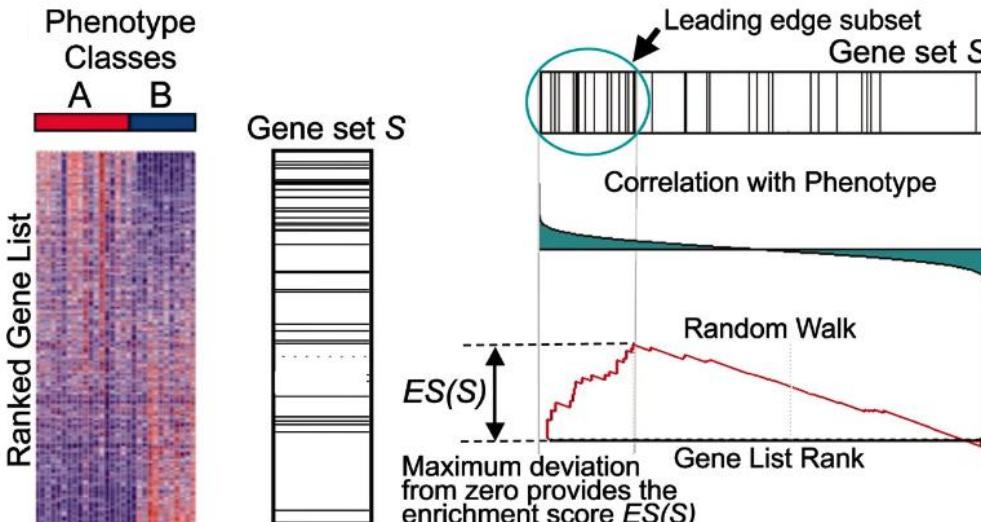
# Constructing a decision point

---



- Identify questions that produce the best separation of patients
  - **Gini impurity**:  $\sum p(1 - p)$
  - **Entropy**:  $-\sum p \ln(p)$
- This is the principle behind decision tree model

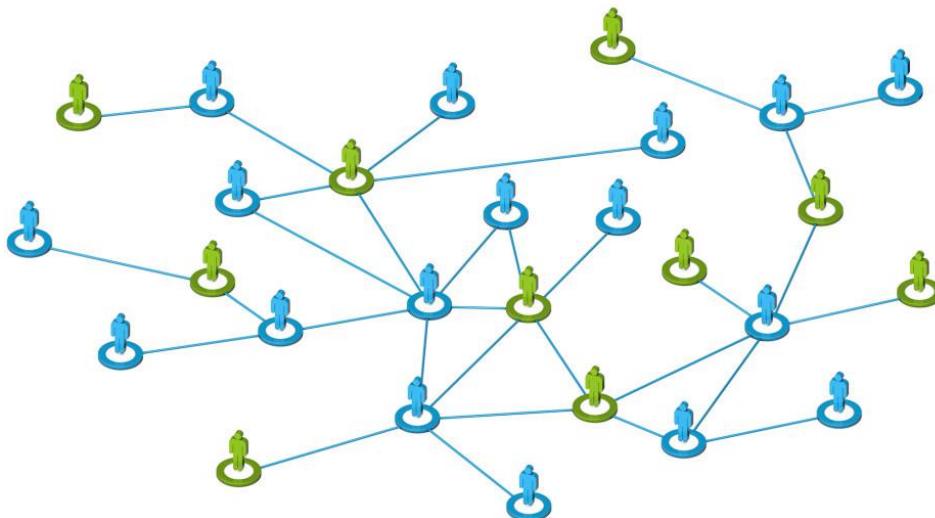
# Let's craft a statistical test



Subramanian et al. PNAS. 102:15545-15550 (2005)

- Sort 1,000 patients from high to low evaluated risks
- Found that female patients are concentrated at the top
- How will you test whether female patients have higher risks?

# Permutation test



- **Hypothesis:** FB friends tend to vote for the same political party
- How can you test this statement?

- Generate 1 million networks with the same number of people and friendship connections
- Count the number of connections between friends voting for the same party
- Compare to the original count



# Mechanistic modeling

# Importance of modeling

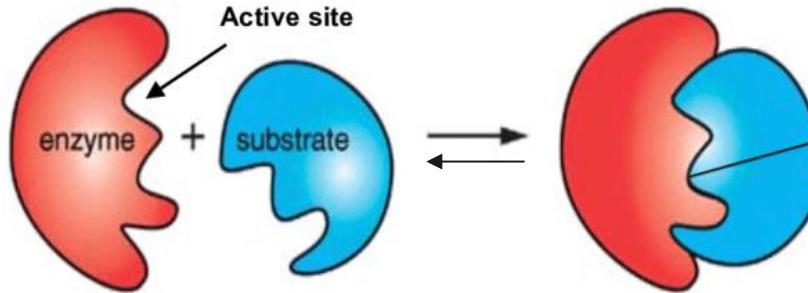
---



What I cannot create, I do not understand.

— *Richard P. Feynman* —

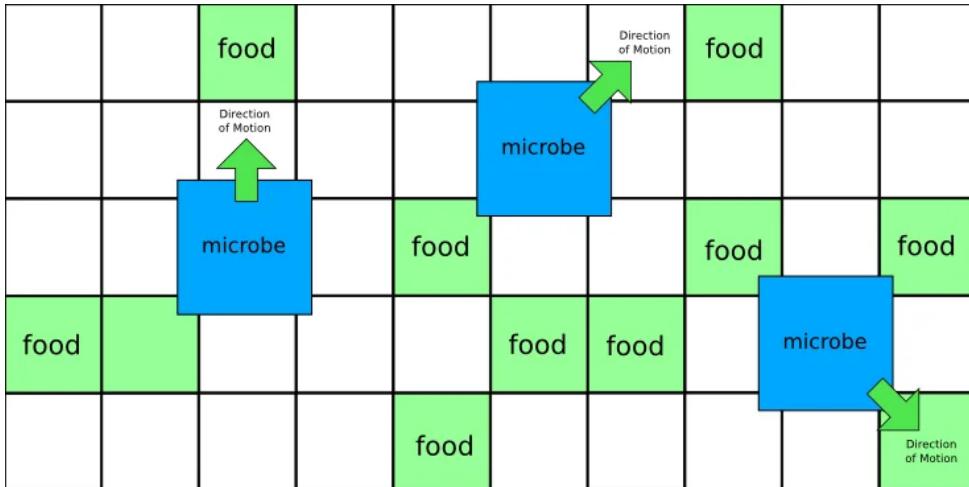
# Enzyme substrate binding model



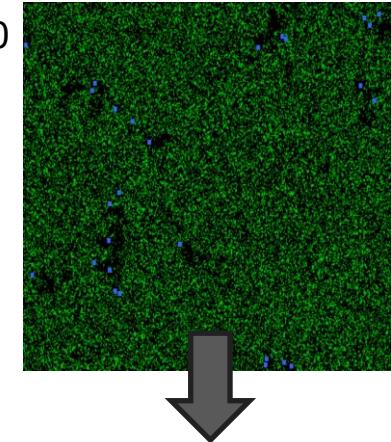
Graham Hutchings. "Development of new highly active nano gold catalysts for selective oxidation reactions" (2014)

- $K_{\text{dissociation}} = [E][S] / [E-S]$  at equilibrium, [ ] = concentration
- Association =  $P(E \text{ collides with } S \text{ in space}) \times Pr_{\text{binding}} \sim k_1 [E][S]$
- Dissociation = Number of E-S molecules  $\times Pr_{\text{dissociation}} \sim k_2 [E-S]$
- At equilibrium, Association = Dissociation
  - $k_1 [E][S] = k_2 [E-S] \rightarrow K_{\text{dissociation}} = k_2 / k_1 = [E][S] / [E-S]$

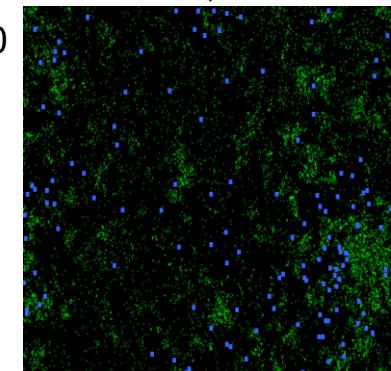
# Evolutionary simulation



Time = 0

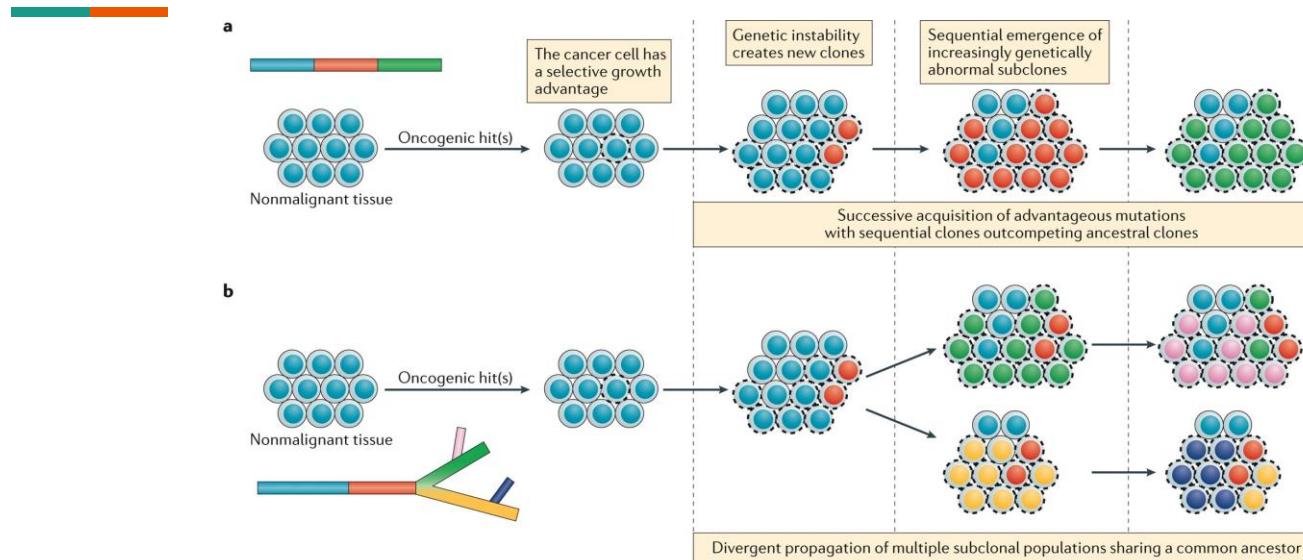


Time = 100



- Simulate bacteria movement behavior and food utilization
  - Random exploration versus guided chemotaxis
- [https://beltoforion.de/en/simulated\\_evolution/](https://beltoforion.de/en/simulated_evolution/)

# Spatial model for tumor growth



Dagogo-Jack and Shaw, Nat Rev Clin Oncol (2017)

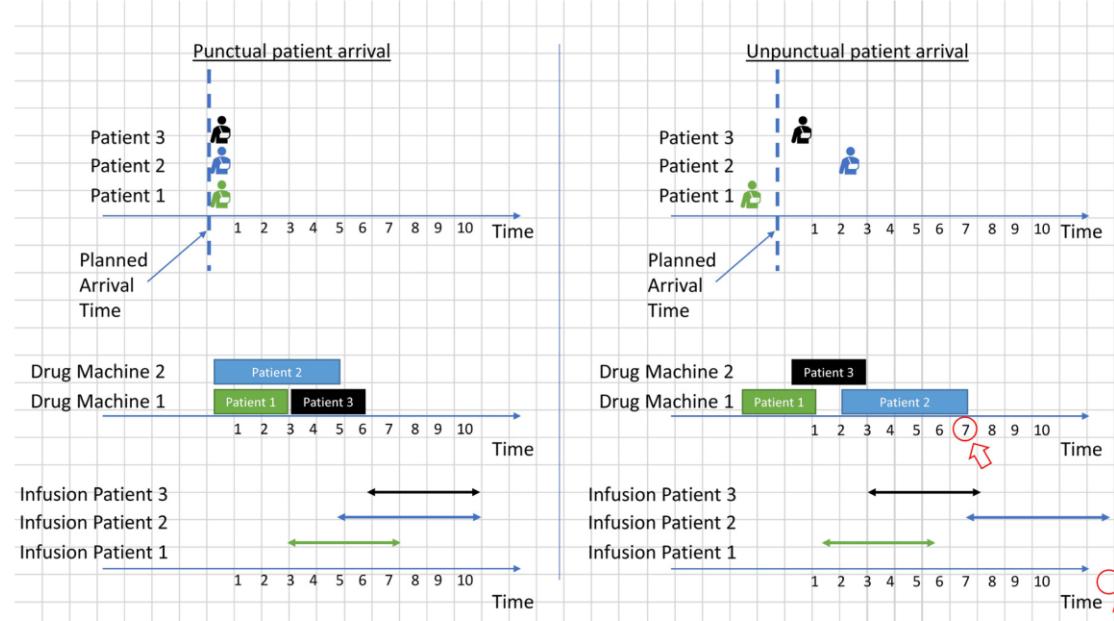
Nature Reviews | Clinical Oncology

- Possible actions: **Mutate, Divide, Die**
  - What are the parameters influencing these actions?
  - Competition for resources: tumor volume, relative division rates, ...



# Empirical, data-driven modeling

# Substituting theory with data

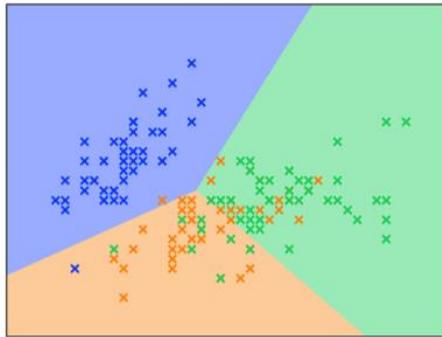


Hadid, M. et al. Int J Environ Res Public Health 2022

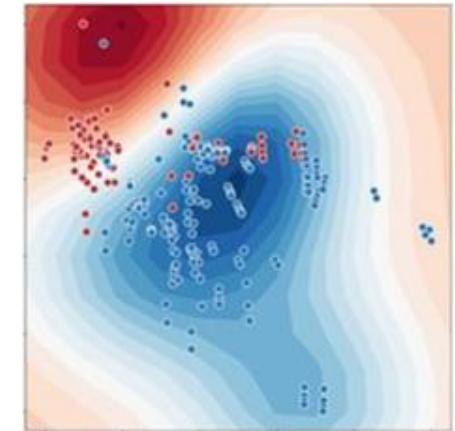
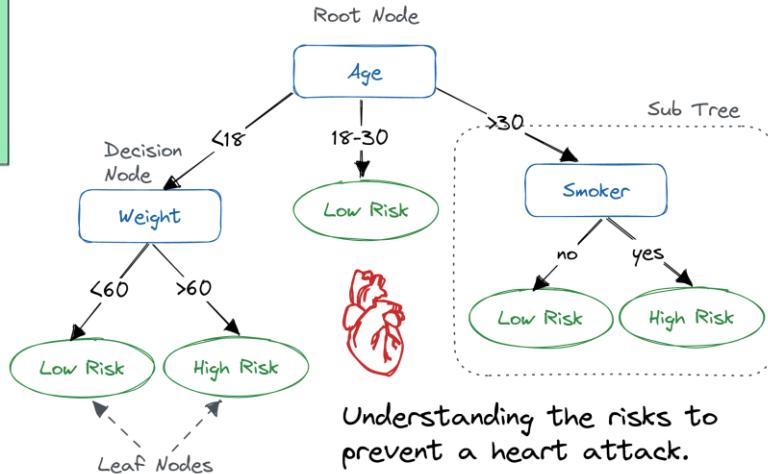
- Some systems cannot be captured by any theoretical distributions
- Build an empirical distribution for sampling using past observations

# Classical machine learning models

Linear: Score =  $(\text{input}_1 \times w_1) + \dots + (\text{input}_n \times w_n)$

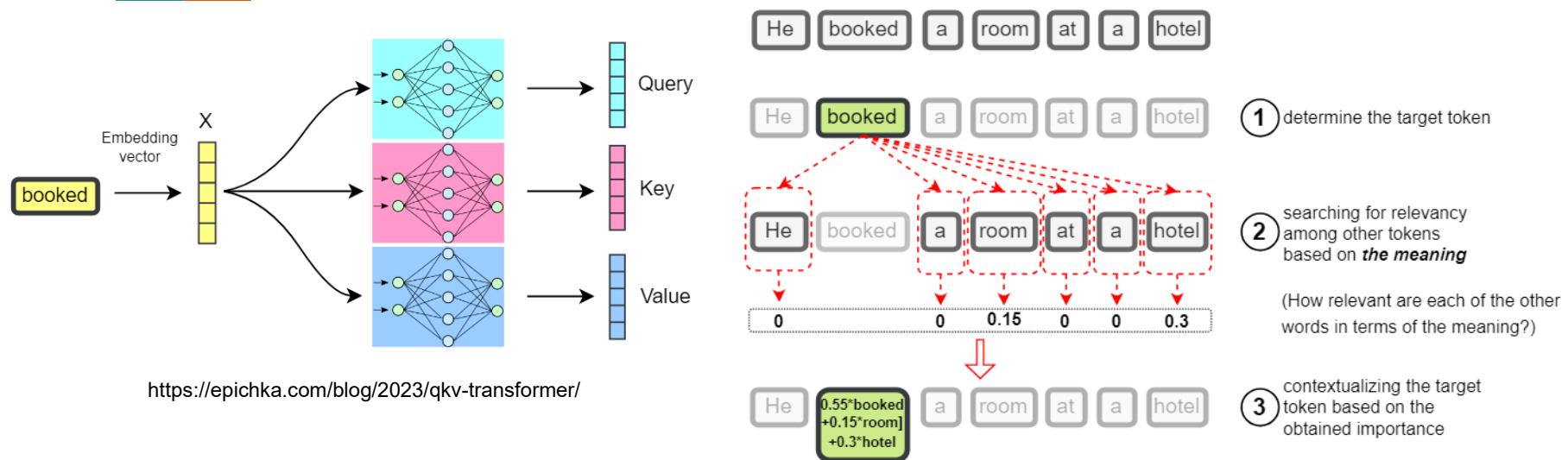


**Tree:** Collection of decisions



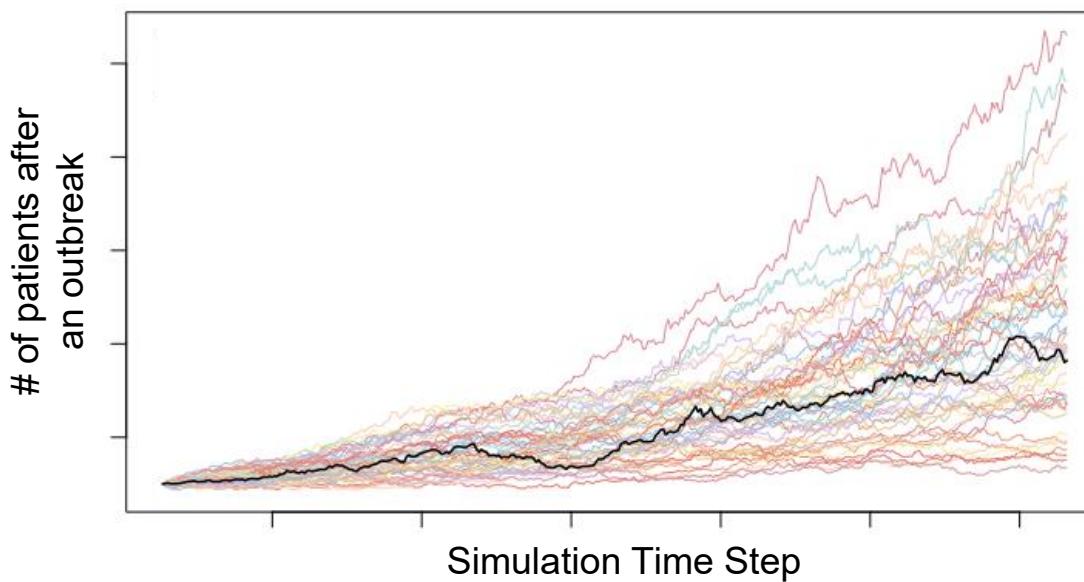
**Neighbor-based:**  
Predict using similarity  
to past observations

# Attention mechanism in transformer architecture



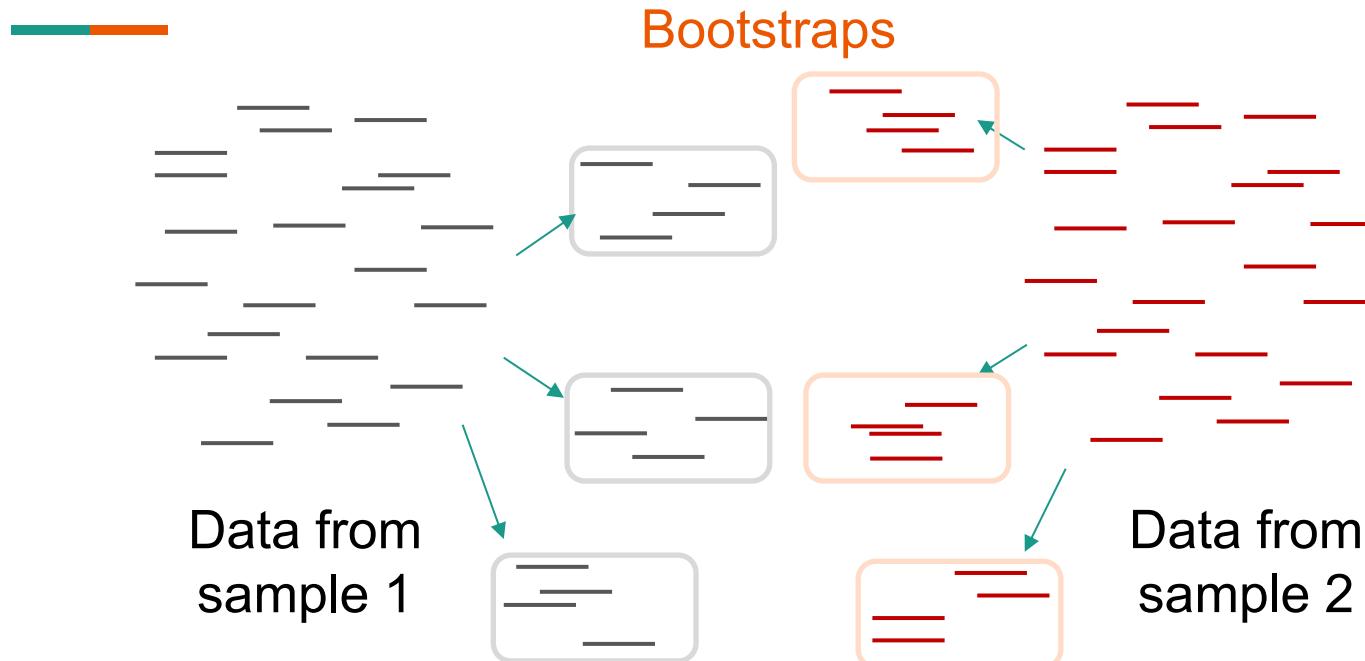
- Generate Query, Key, and Value for each input token
- Use Query-Key similarity to derive attention weights
- Apply weights to the corresponding Value to get the output

# Accounting for uncertainty: Monte Carlo simulation



- Different initial conditions can lead to different outcomes
- At each step, multiple outcomes have a probability to occur
- Simulate a population and analyze

# Variance estimate: Bootstrapping



- Instead of using the whole data to perform one calculation
- Sampling from data to perform multiple calculations → Estimate variance



# Optimization perspective

# Every algorithm involves optimization

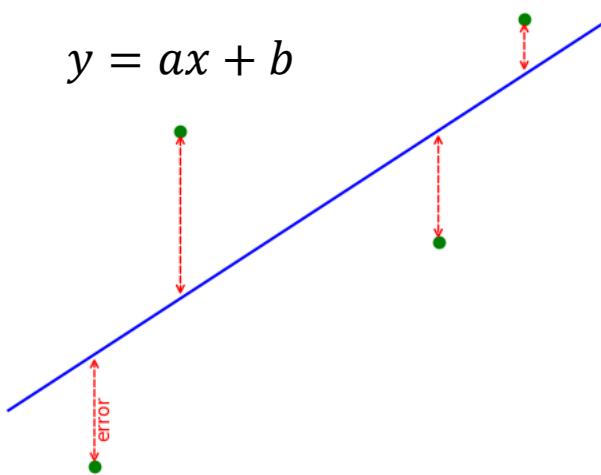
---

- **Curve fitting** = minimize square error between curve and data points
- **Drug-protein docking** = minimize energy of the drug-protein complex
- **Hospital queuing** = minimize average patient wait time
- **Cell type annotation** = maximize similarity to known cell types
- **Sequence alignment** = minimize number of mismatched position

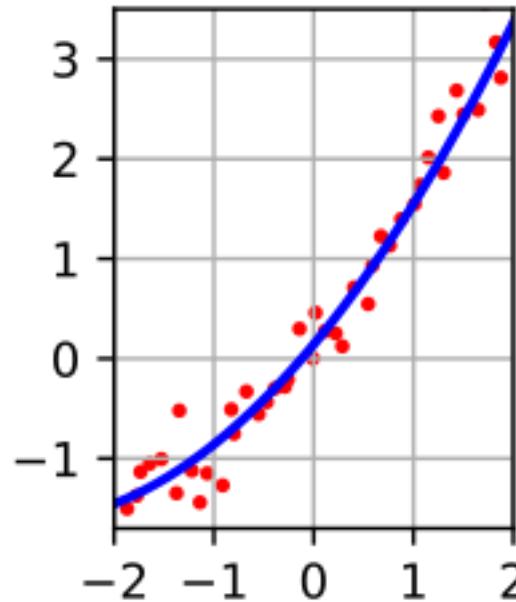
# Curve fitting with least square



$$y = ax + b$$



[https://en.wikipedia.org/wiki/Least\\_squares](https://en.wikipedia.org/wiki/Least_squares)



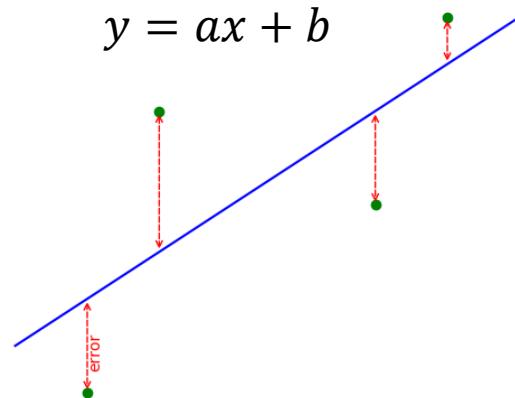
$$y = ax^2 + bx + c$$

- Finding the **best**  $a$ ,  $b$ , and  $c$  that make the curve fit the observations
- Minimize least square error  $\frac{1}{n}((\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_n - y_n)^2)$

# Some optimization approaches

---

- Minimize  $L(a, b) = \frac{1}{n}((\hat{y}_1 - y_1)^2 + \dots + (\hat{y}_n - y_n)^2)$
- Randomly **search** for  $(a, b)$ 
  - **Heuristic** guess:  $a = \frac{\sum y_i}{\sum x_i}$
- Use **calculus**
  - At optimal, slopes are zero
  - Solve  $\frac{\partial L(a,b)}{\partial a} = 0$  and  $\frac{\partial L(a,b)}{\partial b} = 0$



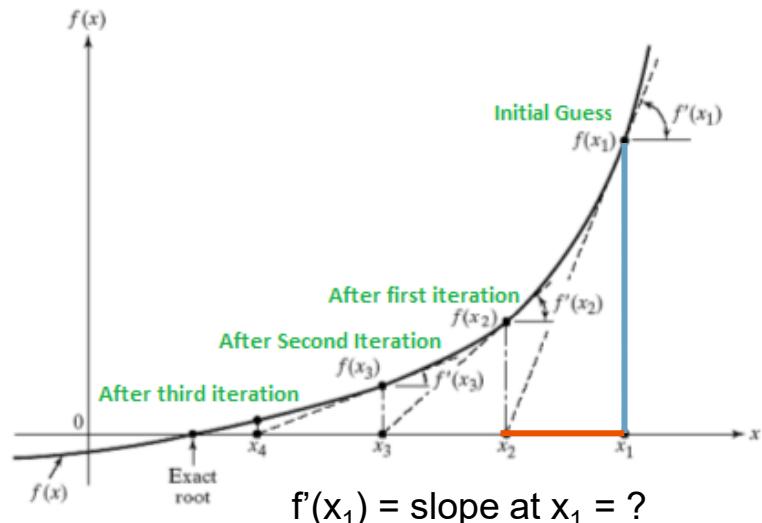
# Newton-Ralphson method

- Want to maximize  $L(x)$  by solving

$$f(x) = \frac{dL(x)}{dx} = 0$$

- Start with an initial guess  $x_1$
- Calculate  $f(x)$  and  $f'(x)$  at  $x_1$
- Define  $x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$
- Repeat the process for  $x_3, x_4, \dots$
- Stop when  $x_i$  converges to a value

$f(x)$  = the first derivative of the objective  $L(x)$



# Combinatorial decision: Knapsack problem

---



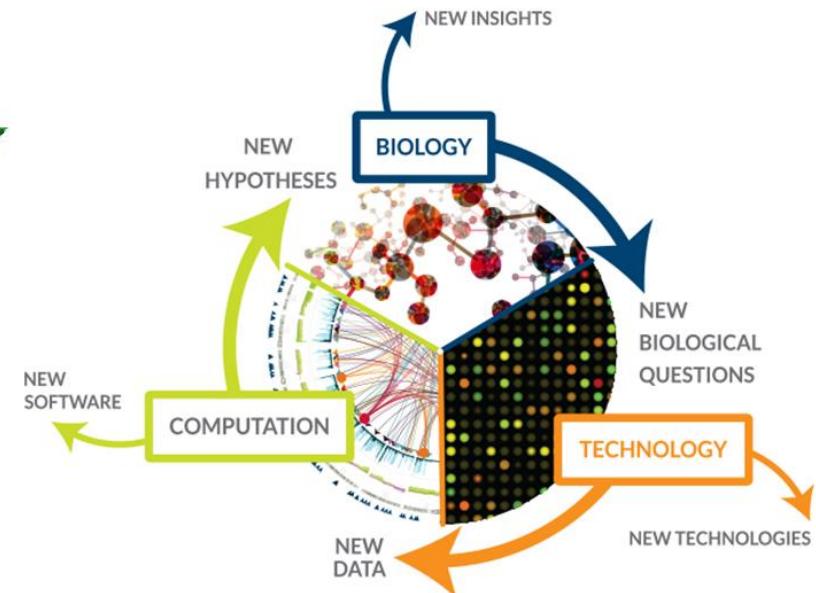
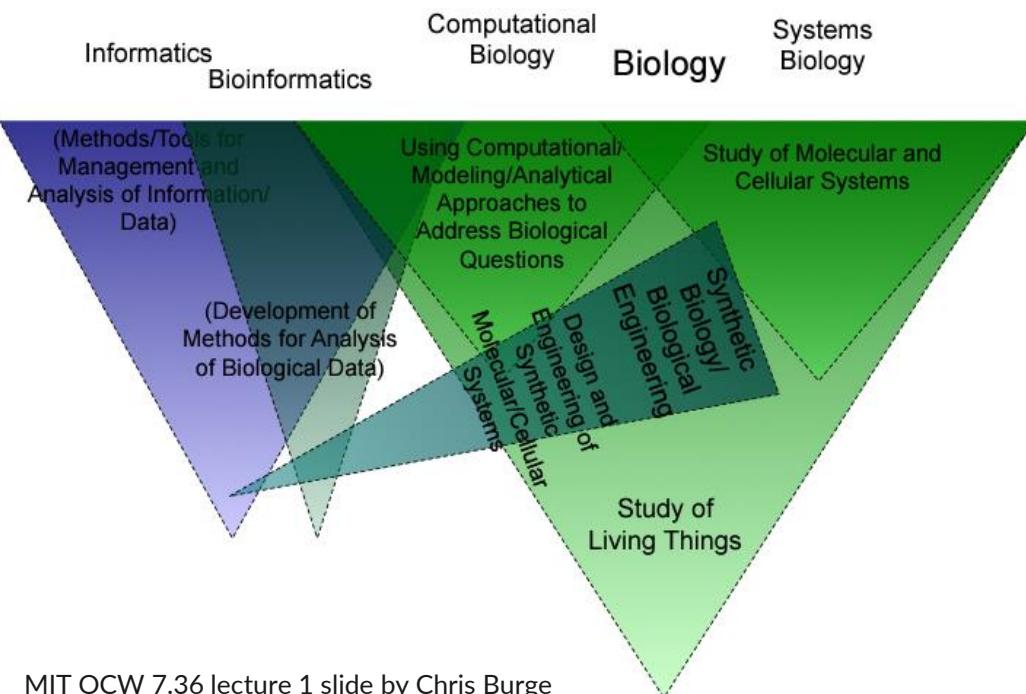
[https://en.wikipedia.org/wiki/Knapsack\\_problem](https://en.wikipedia.org/wiki/Knapsack_problem)

- Limited bag capacity
- Want to pick as much value as possible without exceeding the bag capacity
- General optimization setting
  - **Objective** = Total value
  - **Constraint** = Total weight  $\leq W$
- Evolution is a combinatorial optimization of the genome



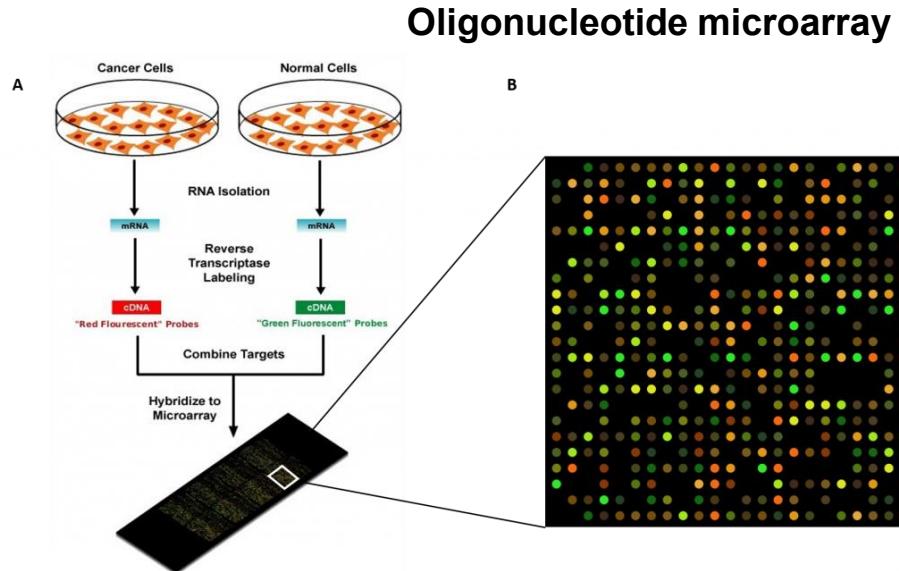
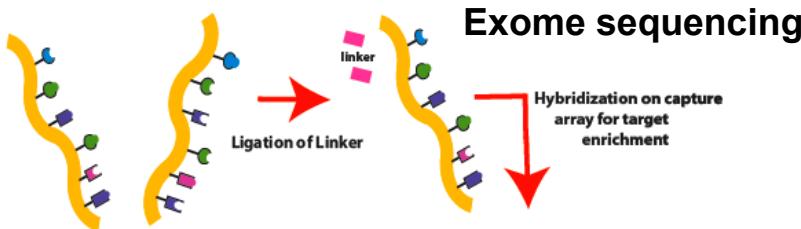
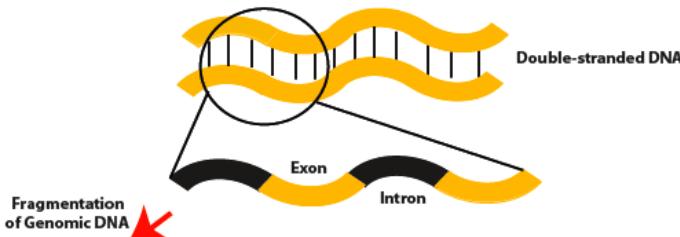
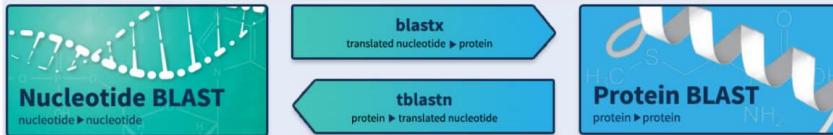
# Computational thinking in biology

# What is computational biology?



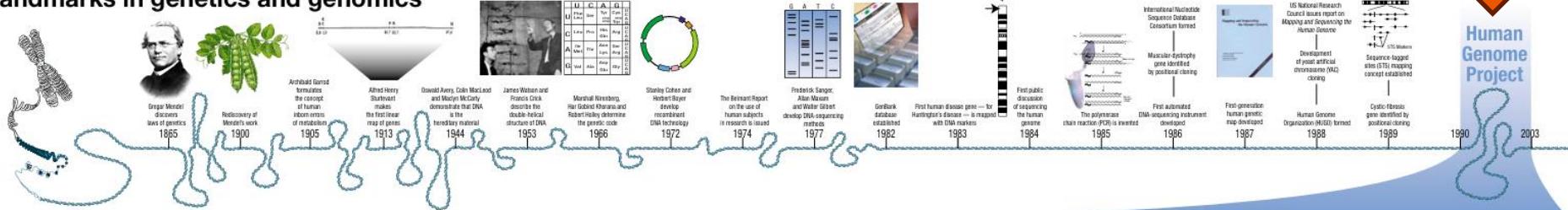
# The Big Bang moment of biology

## BLAST

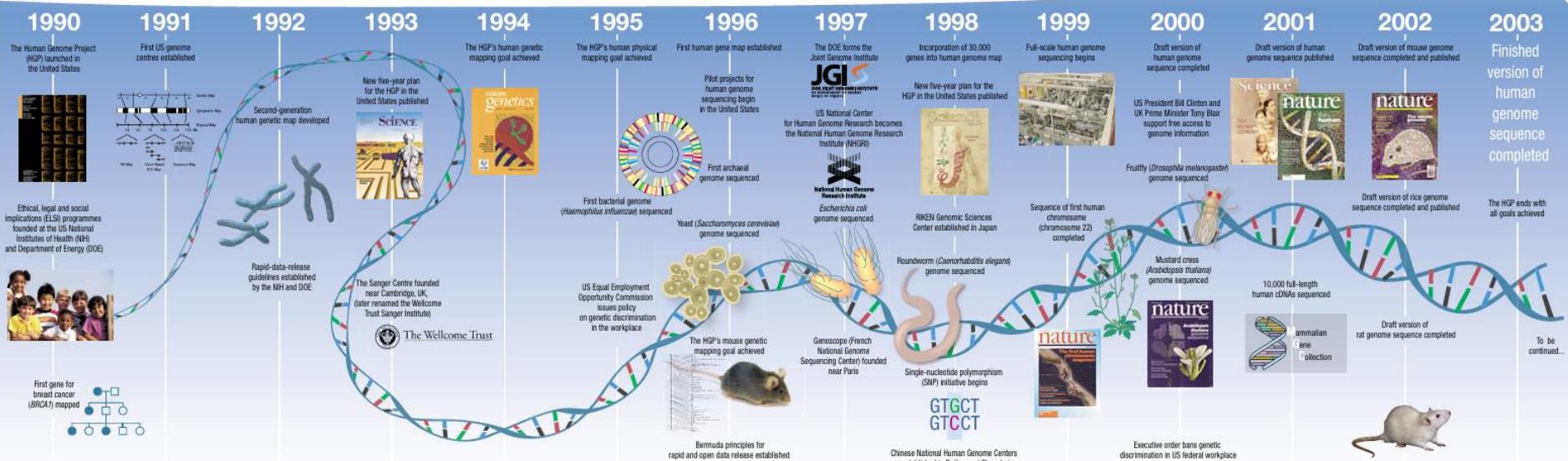


What does these techniques have in common?

## Landmarks in genetics and genomics



## Human Genome Project

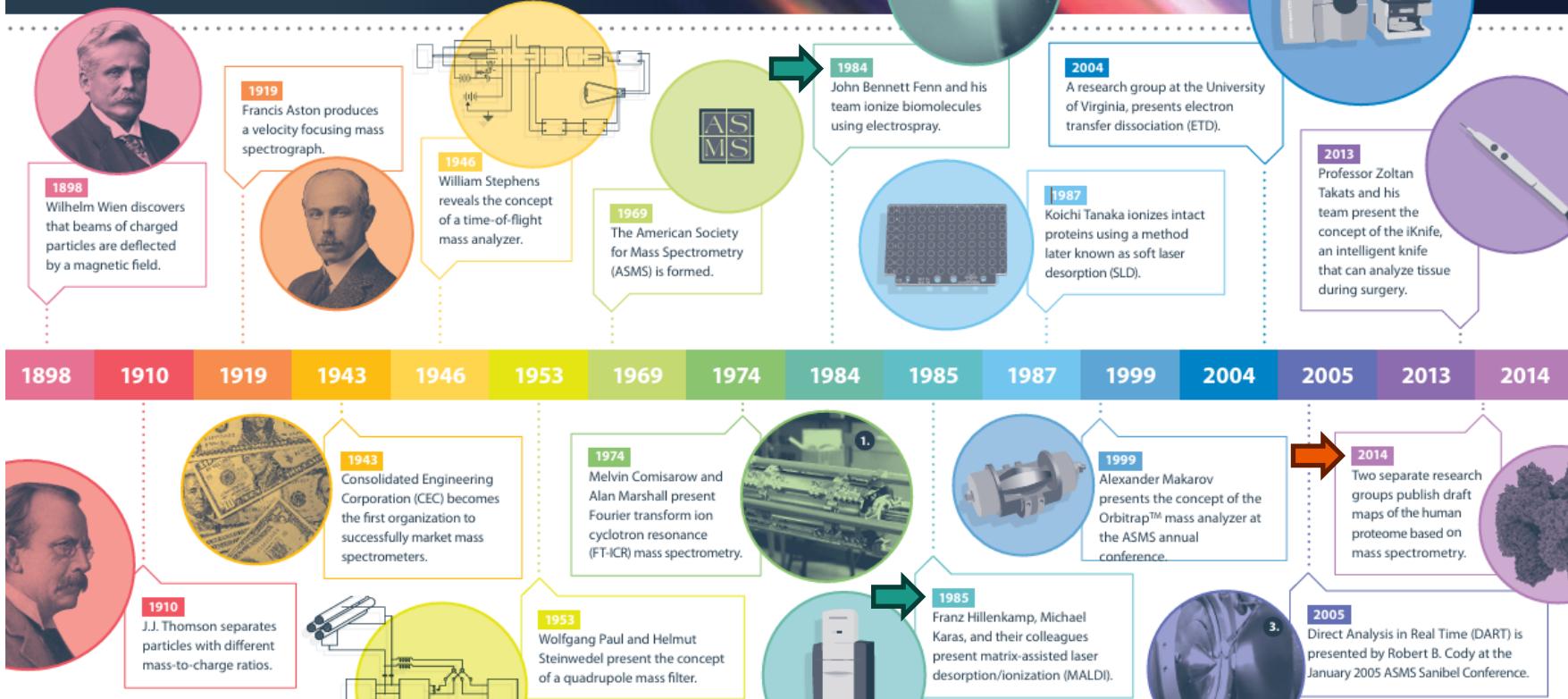


# The Evolution of Mass Spectrometry

A journey through time and technology

## References

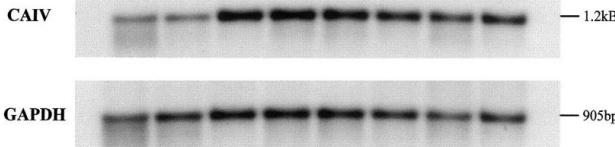
1. Wikimedia Commons, Jeff Dahl
2. Courtesy of Pacific Northwest National Laboratory
3. Wikimedia Commons, JEOL USA



# Digital transformation of biology

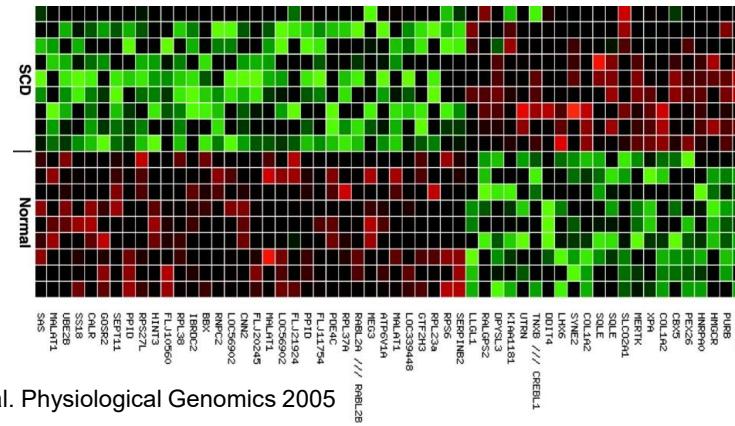


## Qualitative experiment



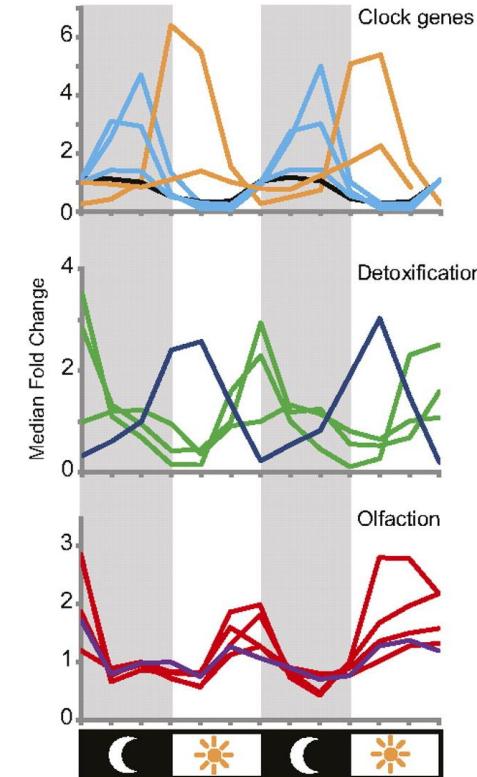
Rosen et al. Am J of Physiology 2001

## High-throughput data



Klings et al. Physiological Genomics 2005

## Time-series data

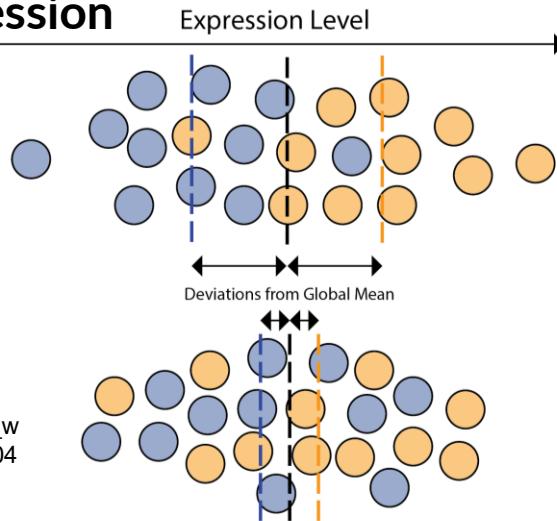
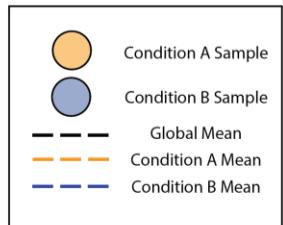


- Paradigm shift from data digitization and introduction of quantitative, high-throughput assays
- From hypothesis-driven to data-driven

Rund et al. PNAS 2011

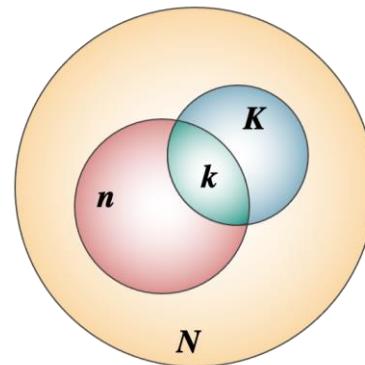
# Statistical approaches

## Differential Expression



[https://hbctraining.github.io/DGE\\_workshop\\_salmon\\_online/lessons/04\\_a\\_design\\_formulas.html](https://hbctraining.github.io/DGE_workshop_salmon_online/lessons/04_a_design_formulas.html)

## Overrepresentation



$$P(X = k) = \frac{\binom{K}{k} \binom{N-K}{n-k}}{\binom{N}{n}}$$

$\textcolor{brown}{N}$  = Background (e.g. Transcriptome  $\sim 40.000$  genes)

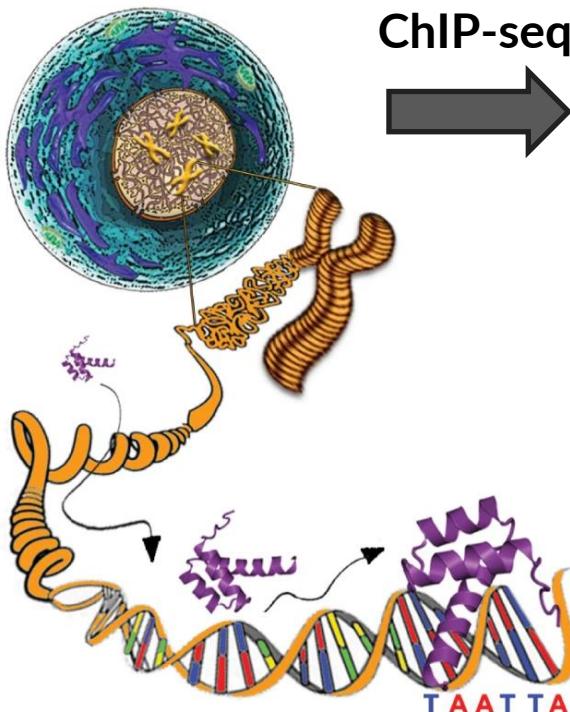
$\textcolor{red}{n}$  = Query list (e.g. upregulated genes)

$\textcolor{blue}{K}$  = genes annotated in the pathway/set tested  
(e.g. Glycolysis)

$\textcolor{teal}{k} = n \cap K$

- Theoretical models for data distribution
- Direct probability calculation to determine significance

# Bioinformatics = algorithms for biological data



ChIP-seq

## Gibbs Sampling

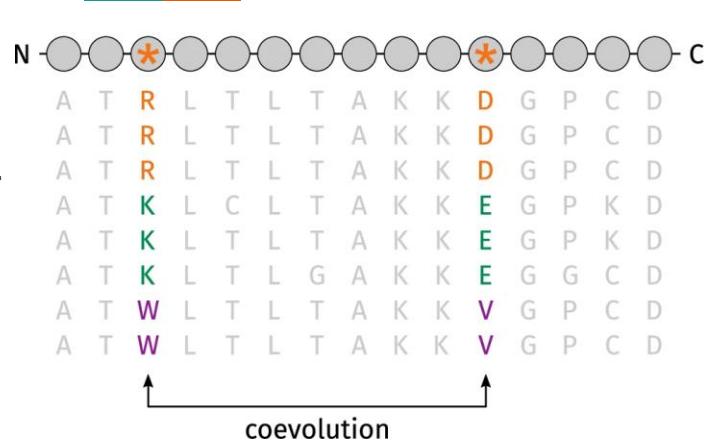
CGGGGCTATcCAgCTGGGTCGTACATTCCCCTT  
TTTGAGGGTGCCAATAAggGCAACTCCAAAGCGGACAAA  
GGATGgAtCTGATGCCGTTGACGACCTA  
AAGGAaGCAACcCCAGGAGGCCCTTGCTGG  
AGATTATAATGTCGGTCCtTGgAACTTC  
CAACTGAGATCATGCTGCATGCCcAtTTTCAAC  
TACATGATCTTGATGgcACTTGGATGAGGGAATGATGC

An Introduction to Bioinformatics Algorithm by Jones and Pevzner

- Use **algorithm** to find the optimal solutions with minimal resource requirement
- Use **statistics** to test whether the optimal solution is a *true signal*

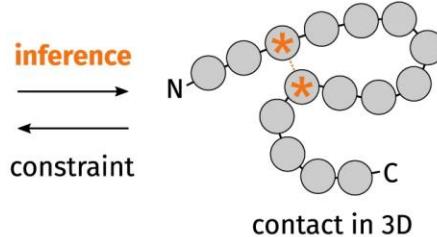
# Capturing co-evolution with mutual information

Different Species

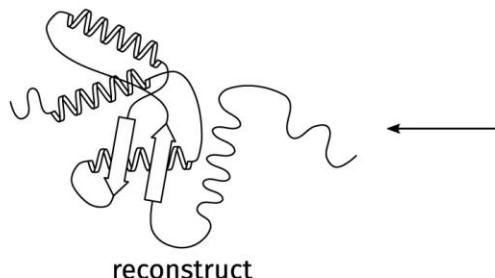


$$I(X;Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_{(X,Y)}(x,y) \log \left( \frac{P_{(X,Y)}(x,y)}{P_X(x) P_Y(y)} \right)$$

Mutual Information



Positions in protein(s) co-evolve because they interact in some fashions – usually physical bindings

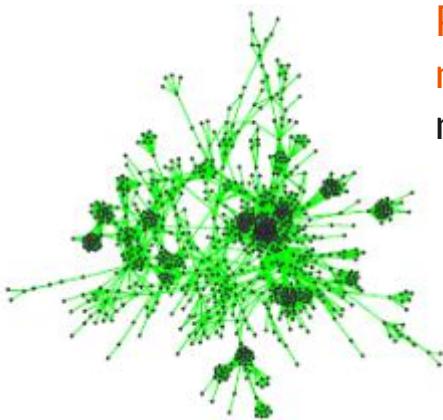


Co-evolution can be quantified and used to guide biological inferences

Bittrich et al. Sci Rep 2019

# Using graph theory to identify important genes

---

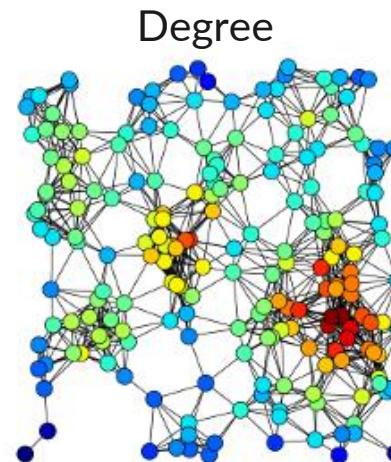


Proteins involved in many interactions might be important

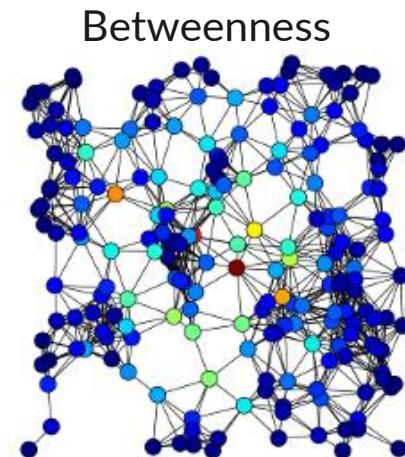
Proteins that connect other proteins might be important

**Node** = Protein

**Edge** = Protein-protein interaction



Degree



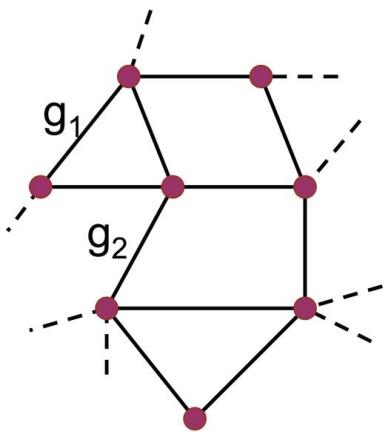
Betweenness

Images from wikipedia

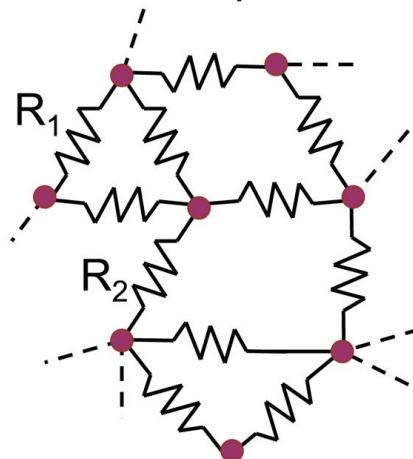
- Genes and proteins work together as a single interaction network

# Electrical circuit model for biological signal flows

Interactome network



Circuit representation



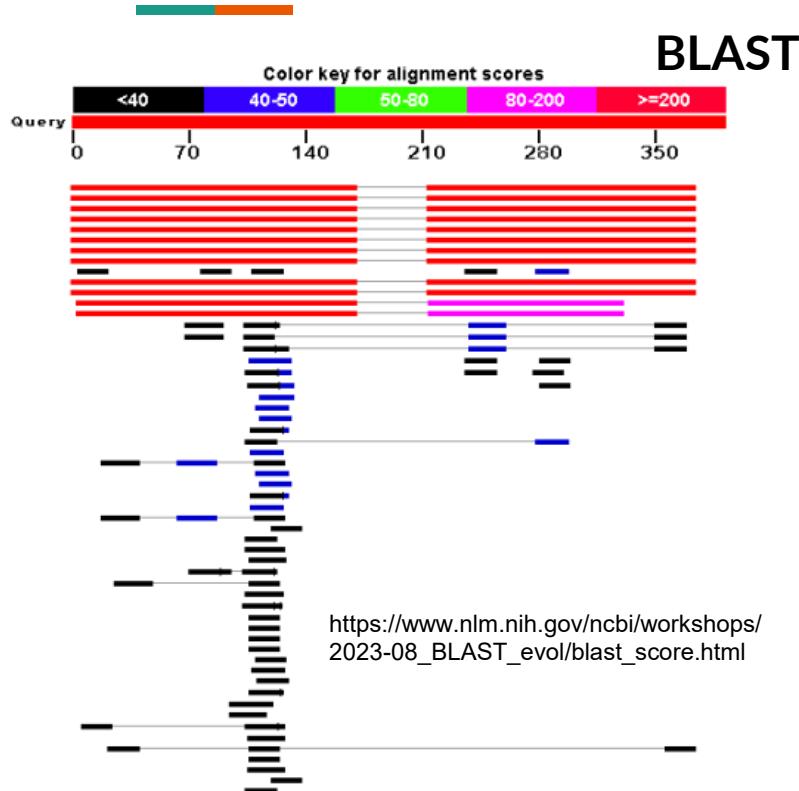
- protein
- zigzag line protein-protein interaction
- $g_n$  interaction confidence score
- $R_n$  resistance

Missiro et al. PLOS Comp Biol 2009

To determine how information flows between two nodes

Apply Kirchoff's laws to calculate the current through the circuit

# Collaboration between algorithm and statistics



[https://www.ncbi.nlm.nih.gov/ncbi/workshops/2023-08\\_BLAST\\_evol/blast\\_score.html](https://www.ncbi.nlm.nih.gov/ncbi/workshops/2023-08_BLAST_evol/blast_score.html)

- The best solution can be wrong if the true answer is not among the options
- Statistics can identify improbable answers



# Overview of this year's contents

# Module 1: DNA sequencing and applications

---

2	3	18-Aug-25	Mon	10:00-11:30	Sequence alignment	
	4	20-Aug-25	Wed	10:00-11:30	DNA sequencing techniques and applications	<b>Assignment 1</b>
3	5	25-Aug-25	Mon	10:00-11:30	Sequencing data processing workflows	
	6	27-Aug-25	Wed	10:00-11:30	Genome assembly and annotation	
4	7	1-Sep-25	Mon	10:00-11:30	Variant calling and genome-wide association study	
	8	3-Sep-25	Wed	10:00-11:30	Phylogenetics and evolutionary models	<b>Assignment 2</b>
5	9	8-Sep-25	Mon	10:00-11:30	Metagenomics	
	10	10-Sep-25	Wed	10:00-11:30	ChIP-seq and DNA motif detection	

- What can be inferred from sequence data?
- DNA sequencing: which platform to use?
- Evolutionary perspective
- Genomics and metagenomics
- Discovery of DNA binding and regulatory motifs

# Module 2: Transcriptomics

---

6	11	15-Sep-25	Mon	10:00-11:30	Transcriptomics techniques	
	12	17-Sep-25	Wed	10:00-11:30	Differential expression analysis	<b>Assignment 3</b>
7	13	22-Sep-25	Mon	10:00-11:30	Functional enrichment analysis	
	14	24-Sep-25	Wed	10:00-11:30	Single-cell and spatial techniques	
8	15	29-Sep-25	Mon	10:00-11:30	Single-cell data processing	

- Why is transcriptomics one of the most popular omics?
- How to identify differentially expressed genes?
- How to narrow down on affected pathways?
- Benefits of single-cell and spatial omics
- A quick look at single-cell transcriptomics data

# Module 3: Broad fields in computational biology

---

	16	1-Oct-25	Wed	10:00-11:30	Proteomics and mass spectrometry	
9	17	6-Oct-25	Mon	10:00-11:30	Introduction to systems biology and dynamics	<b>Assignment 4</b>
	18	8-Oct-25	Wed	10:00-11:30	Multi-omics integration	
10	19	13-Oct-25	Mon		Holiday	
	20	15-Oct-25	Wed	10:00-11:30	Biological networks	
11	21	20-Oct-25	Mon	10:00-11:30	Chromatin conformation capture	<b>Assignment 5</b>
	22	22-Oct-25	Wed	10:00-11:30	RNA and protein structure models	
12	23	27-Oct-25	Mon	10:00-11:30	Online tools and resources	
	24	29-Oct-25	Wed	10:00-11:30	Microscopy data analysis	

- Study of proteins
- Systems and dynamics modeling
- Biological networks
- Structural models
- Imaging

# Module 4: Computational tools



							<b>Assignment 6</b>
13	25	3-Nov-25	Mon	10:00-11:30	Essential statistics for computational biologist		
	26	5-Nov-25	Wed	10:00-11:30	Omics data visualization		
14	27	10-Nov-25	Mon	10:00-11:30	Introduction to machine learning and AI		
	28	12-Nov-25	Wed	10:00-11:30	Machine learning applications		
15	29	17-Nov-25	Mon	10:00-11:30	Foundational AI models in biology		
	30	19-Nov-25	Wed	10:00-11:30	Biomarker discovery		

- Statistics, revisited
- Data visualization
- Machine learning and AI
- Biomarker discovery

# Assignments: What to expect?

---

- Analyze data using online tools and pre-built software like BLAST, MEGA, MaxQuant, Web-Gestalt, CellProfiler, Cytoscape, etc.
- Identify key details in research articles
- Design experiments using omics techniques and bioinformatics tools
- Critique LLM's answers
- Edit R and Python scripts to analyze data

# **Any question or comment?**

---

See you next time