
3000788 Intro to Comp Molec Biol

Week 1: Course introduction and logistics

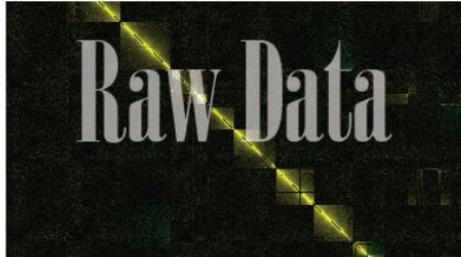
Fall 2024



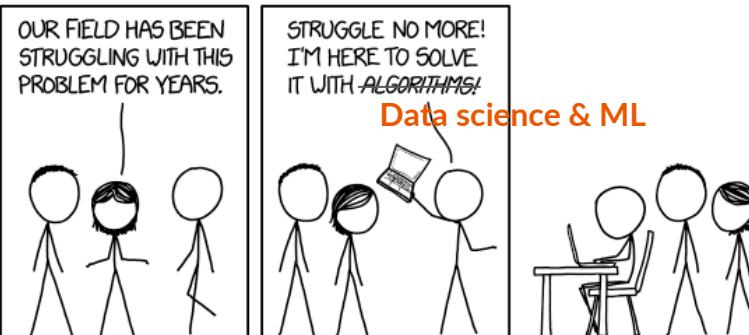
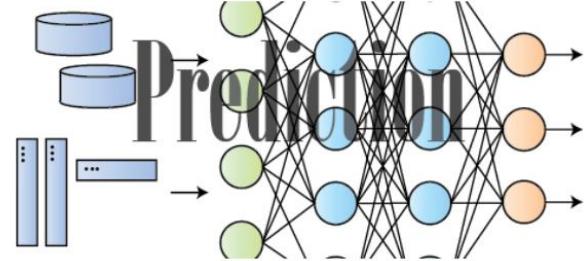
Sira Sriswasdi, PhD

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

About instructor



$$\begin{aligned} & + A \sum_{s=0}^{\infty} \sum_{t=0}^{\infty} \sum_{m=0}^{d-1} s \cdot B^t \cdot \mathbf{1}_{m=s \leq t} \cdot p_s^{t-s} \\ & + A \sum_{s=0}^{\infty} \sum_{t=0}^{\infty} \sum_{m=0}^{d-1} s \cdot B^t \cdot \mathbf{1}_{m < s \leq t} \cdot \mathbf{1}_{r \leq m} \cdot \binom{t-s+k-1}{k-1} \cdot p_s^{t-s} (1-p_s)^k \\ & + A \sum_{s=0}^{\infty} \sum_{t=0}^{\infty} \sum_{m=0}^{d-1} s \cdot B^t \cdot \mathbf{1}_{m < s \leq t} \cdot \mathbf{1}_{r > m} \cdot \binom{t-s}{k} \cdot p_s^{t-s} (1-p_s)^{k+1} \\ \sum_{\dots}^{\infty} \sum_{\dots}^{\infty} \sum_{\dots}^{d-1} s \cdot B^t \cdot \mathbf{1}_{m \geq s=t} & = \sum_{\dots}^{d-1} \sum_{\dots}^m t \cdot B^t \end{aligned}$$



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

- BS in Mathematics
- PhD in Computational Biology

Keywords

- Molecular Evolution
- Proteomics / Transcriptomics
- Machine Learning

About you

Please introduce yourself

- Name & nickname
- Graduate program & year
- Undergraduate background
- Research interest
- Thesis advisor & topic (if you already picked)

Fill out the pre-course questionnaire:

https://docs.google.com/forms/d/e/1FAIpQLSfXPgzsG5_Z7HvGTpsKw_mnwetwrfNkrzZ3-iC7b40R9XD2hw/viewform

About this course

- Survey broad topics in computational molecular biology
- Provide foundation for understanding deeper courses
 - Always ask “Why?”
- Mixing discussion and practice
- Exposure to basic bioinformatics workflows
- Some taste of
 - Python programming
 - Machine learning & data science
- Go to <https://github.com/cmb-chula/comp-biol-3000788/>

Course schedule



Module	Date	Topics	Assignment
1	August	DNA sequencing & applications	Problem set 1-2
2	September	Transcriptomics: bulk and single-cell	Problem set 3-5
3	October	Other omics - Proteomics - Chromatin organization - Biological networks	Problem set 6-8
4	November	Programming & machine learning	Problem set 9-10
		Post-course evaluation	

Course schedule (Thai section)

Session	Date	Time	Topics	Videos	Assignment
1	08/14/24	1-3pm	Course introduction	-	Pre-course evaluation
2	08/21/24	1-3pm	DNA sequencing and data processing	L4 / L5	PS1
3	08/28/24	1-3pm	Sequence alignment / Phylogenetics	L6 / L7	PS2
4	09/04/24	1-3pm	Metagenomics / Microarray	L8 / L9	PS3
5	09/11/24	1-3pm	RNA-seq / Expression data analysis	L11 / L13	PS4
6	09/18/24	1-3pm	RNA-seq analysis demo	-	PS5
7	09/25/24	1-3pm	Single-cell technique and demo	L14	-
8	10/02/24	1-3pm	Proteomics and demo	L16	PS6
9	10/09/24	1-3pm	Biological networks / Chromatin	L18 / L20	PS7
10	10/16/24	1-3pm	Dynamics simulation / Online resources	L21 / L22	PS8
-	10/23/24	-	-	-	-
11	10/30/24	1-3pm	Python programming session 1	L23 / L24	-
12	11/06/24	1-3pm	Python programming session 2	L25 / L26	PS9
13	11/13/24	1-3pm	Principles of machine learning	L27 / L28	-
14	11/20/24	1-3pm	Python programming session 3	L29	PS10
15	11/27/24	1-3pm	Deep learning in life sciences	L30	Post-course evaluation

Course schedule (English section)

Session	Date	Time	Topics	Videos	Assignment
1	08/16/24	1:30-3:30pm	Course introduction	-	Pre-course evaluation
2	08/23/24	1:30-3:30pm	DNA sequencing and data processing	L3 / L4	PS1
3	08/30/24	1:30-3:30pm	Sequence alignment / Phylogenetics	L5 / L6	PS2
4	09/06/24	1:30-3:30pm	Metagenomics / Microarray	L7 / L8	PS3
5	09/13/24	1:30-3:30pm	RNA-seq / Expression data analysis	L10 / L12	PS4
6	09/20/24	1:30-3:30pm	RNA-seq analysis demo	-	PS5
7	09/27/24	1:30-3:30pm	Single-cell technique and demo	L13	-
8	10/04/24	1:30-3:30pm	Proteomics and demo	L15	PS6
9	10/11/24	1:30-3:30pm	Biological networks / Chromatin	L17 / L18	PS7
10	10/18/24	1:30-3:30pm	Dynamics simulation / Online resources	L20	PS8
11	10/25/24	1:30-3:30pm	Python programming session 1	L21 / L22	-
12	11/01/24	1:30-3:30pm	Python programming session 2	L23 / L24	PS9
13	11/08/24	1:30-3:30pm	Principles of machine learning	L25 / L27	-
14	11/15/24	1:30-3:30pm	Python programming session 3	L26 / L28	PS10
15	11/22/24	1:30-3:30pm	Deep learning in life sciences	L29	Post-course evaluation

Learning strategy

- Before class
 - **Watch assigned videos** (1-2 per week, 1.5hr each)
 - **Submit your pre-class query:**
<https://docs.google.com/forms/d/e/1FAIpQLSdL8Kx9Z1L41ZJeLwu0jDXOOZ7wZ1Gc8NbVQH2nQ87ytV1Yw/viewform>
 - **Download data and install software**
- In class
 - **Lecture-style recap** of key topics
 - **Q & A** of your pre-class query
 - **Discussion**

Grading criteria

- Problem set [9% x 10 problem sets = 90%]
 - Can work with each other
 - **But write your own answer**
 - **You may use ChatGPT but also report how you used it**
- In-class activity [10%]
 - Pre-class query
 - Hands-on practice
 - Discussion
 - Attendance

Companion courses from MIT (7.36)

FOUNDATIONS OF COMPUTATIONAL AND SYSTEMS BIOLOGY



Lecture 1: Introduction to Computational and Systems Biology



Lecture 2: Local Alignment (**BLAST**) and Statistics

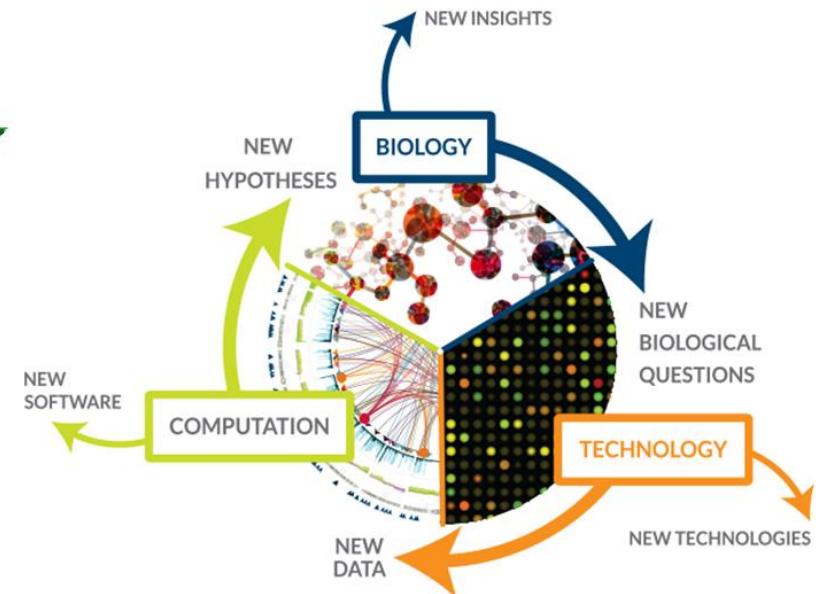
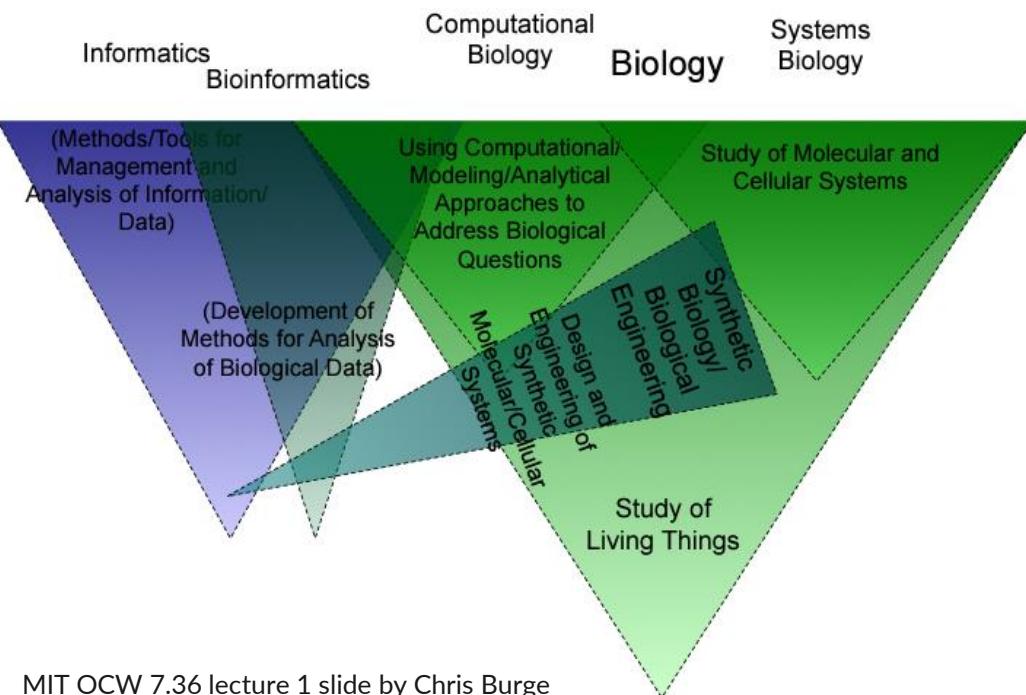


Lecture 3: Global Alignment of Protein Sequences (**NW, SW, PAM, BLOSUM**)

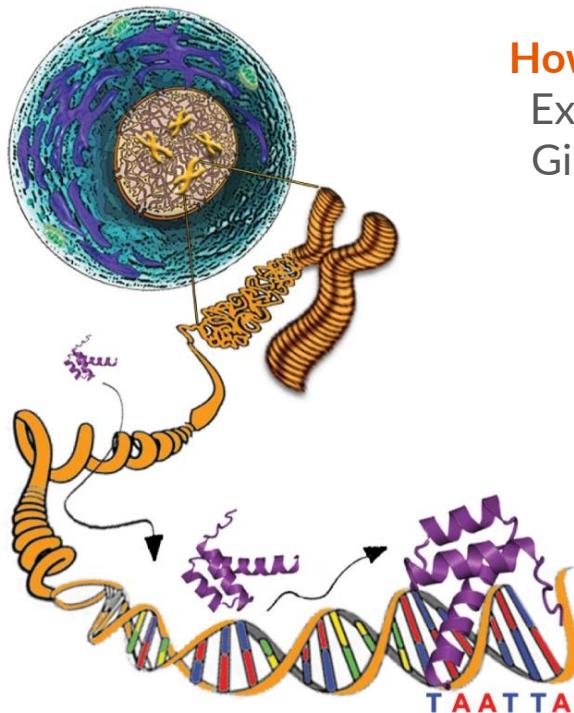


Lecture 4: Comparative Genomic Analysis of Gene Regulation

What is computational biology?



Biology-inspired motif analysis



How to identify the motif?

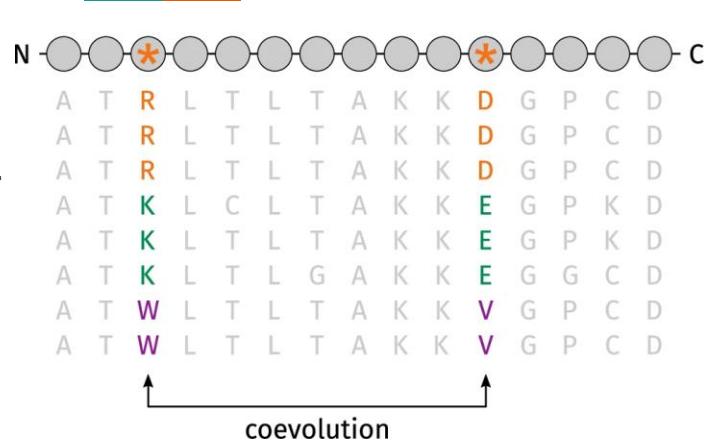
Experiment can isolate DNA sites with bound protein
Gibbs Sampling to identify similar patterns



CGGGGCTATcCAgTGGGTCGTACTCCCC
TTTGAGGGTGCCAATAAggGCAACTCCAAAGGCGACAA
GGATGAtTGATGCCGTTTGACGACCTA
AAGGAaGCAACCCCCAGGGAGGGCCTTTGCGG
AGATTATAATGTCGGTCCtTGgAACTC
CAACTGAGATCATGCTGCATTTCAAC
TACATGATTTTGATGgCACTTGGATGAGGGAATGATGC

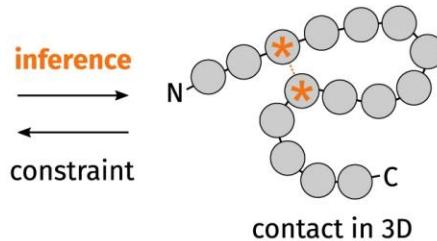
Capturing co-evolution with mutual information

Different Species

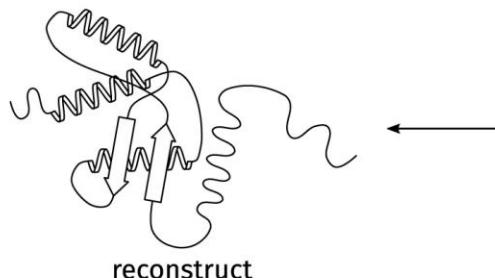


$$I(X;Y) = \sum_{y \in \mathcal{Y}} \sum_{x \in \mathcal{X}} P_{(X,Y)}(x,y) \log \left(\frac{P_{(X,Y)}(x,y)}{P_X(x) P_Y(y)} \right)$$

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 study biological networks

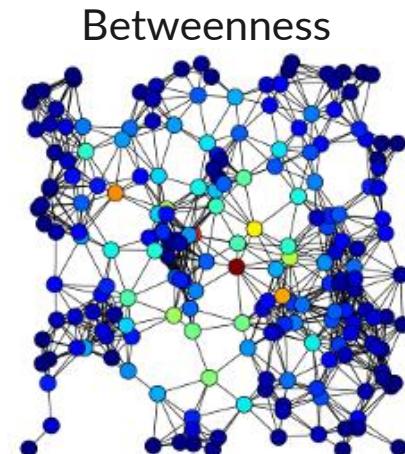
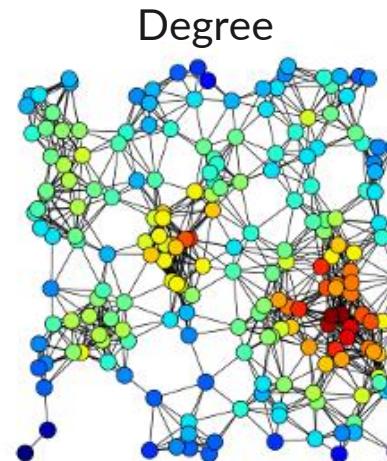


Proteins involved in
many interactions
might be important

Proteins that connect
other proteins might
be important

Node = Protein

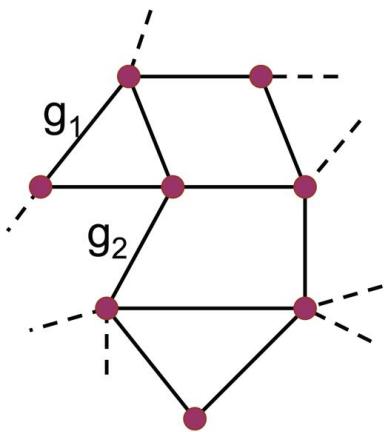
Edge = Protein-protein interaction



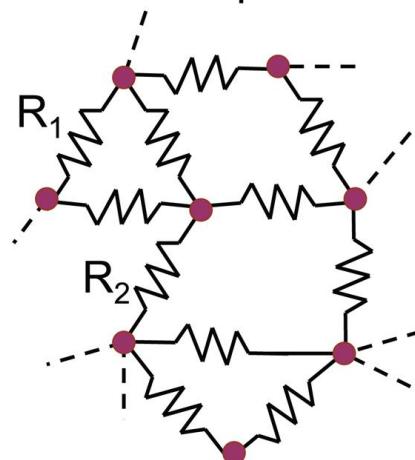
Images from wikipedia

Electrical circuit model for biological signal flows

Interactome network



Circuit representation



- protein
- w — 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

Course objectives

The diagram consists of three large, overlapping circles. The leftmost circle is light blue and contains the text 'Computational Knowledge' in bold blue and 'Mathematics Statistics' in blue. The middle circle is light orange and contains 'Domain Knowledge' in bold dark gray, followed by 'Biology', 'Chemistry', 'Immunology', and 'etc.' in a smaller dark gray font. The rightmost circle is light green and contains 'Data Skills' in bold blue, 'Programming' in blue, and 'Modeling' in blue.

**Computational
Knowledge**

Mathematics
Statistics

Domain Knowledge

Biology
Chemistry
Immunology
etc.

Data Skills

Programming
Modeling

Knowledge enables communication

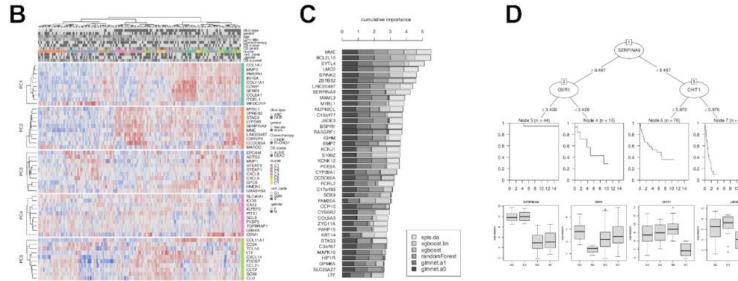


Module 0: Statistics and computational thinking

- P-values help distinguish biological pattern from random chance
 - Differential test of gene expression
 - Enrichment of biological function terms among differentially expressed genes
- How were they calculated?
 - Null hypothesis
 - Correction for multiple testing
- What can you do with computational thinking?
 - Data analysis
 - Modeling and simulation
 - Understand algorithm / workflow

What can you do with computational thinking?

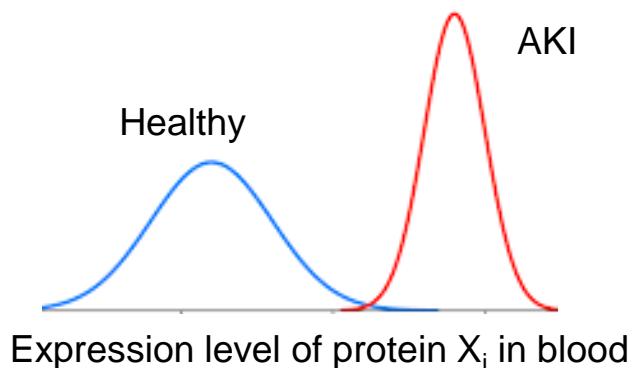
- Data analysis
 - Identify calculations that will support or disprove your hypothesis
 - Be aware of assumptions and limitations of each calculation
- Modeling and simulation
- Algorithm



Akhmedov, M. et al. NAR Genom and Bioinfor, 2(1):lqz019 (2019)

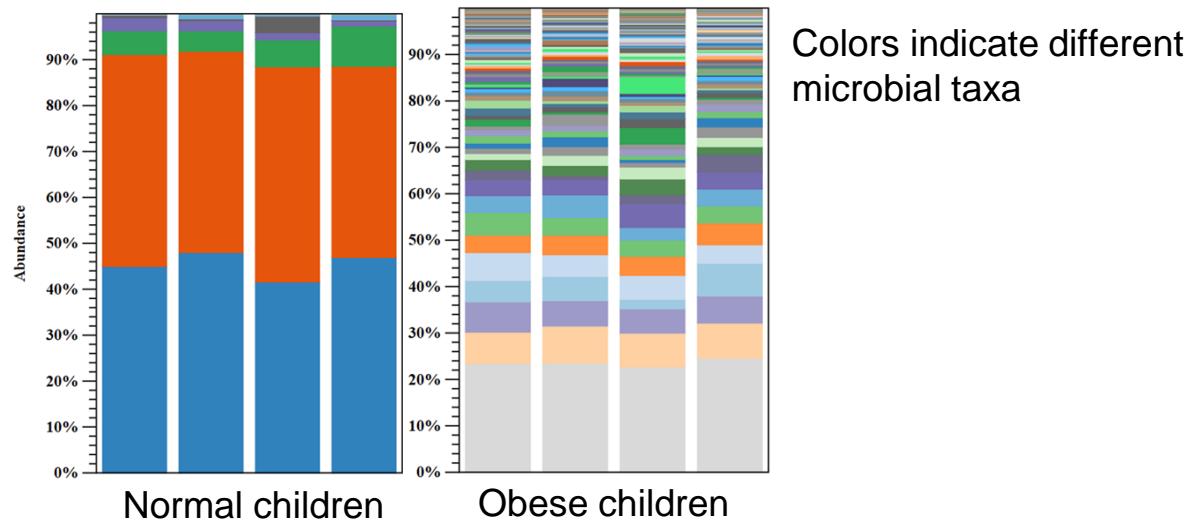
Gene ID	P61_2_C	P62_2_C	P63_2_C	P64_2_C	P68_2_C
ENSG00000000003.14	4.637576	6.183992	5.237635	2.372719	5.665966
ENSG00000000005.5	0	0	0	0	0
ENSG00000000419.12	11.22781	4.813792	2.99782	10.99452	10.7482
ENSG00000000457.13	7.656414	5.082675	7.710682	9.014404	8.488388
ENSG00000000460.16	3.172546	2.245954	5.974815	3.501081	4.162024
ENSG00000000938.12	0	0	0	0.042488	0
ENSG00000000971.15	6.626259	8.19511	5.904925	11.7748	2.050394
ENSG00000001036.13	1.790445	0.76823	3.670635	0.68115	1.894823
ENSG00000001084.11	19.53907	25.08378	11.04872	5.815902	20.23763
ENSG00000001167.14	15.34717	20.00867	17.10001	25.31168	27.41216
ENSG00000001460.17	0.889852	3.090642	0.744581	3.439525	2.417934
ENSG00000001461.16	3.771195	3.12468	1.385353	2.767444	2.973217
ENSG00000001497.16	16.75059	9.662455	15.4965	14.34071	10.62035
ENSG00000001617.11	2.998366	3.712208	3.885852	17.50663	3.019686

Quantifying the “goodness” of biomarkers



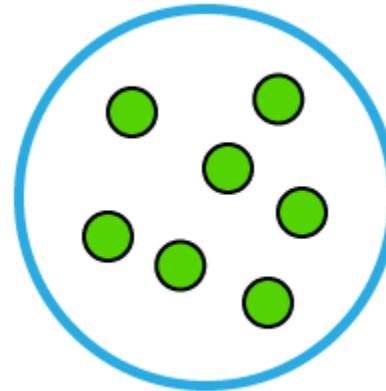
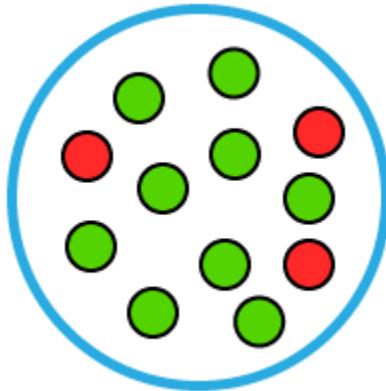
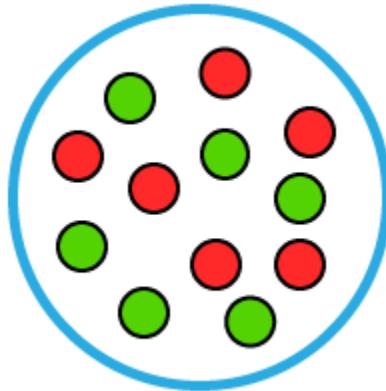
- Identify significantly differentially expressed proteins, X_1, X_2, \dots, X_{200} with t-test
- What if we want to identify the **best biomarkers**?
 - How to quantify the ability of a protein to distinguish healthy and AKI?
 - Independent two-sample t-test's statistics =
$$\frac{\text{Mean}_1 - \text{Mean}_2}{\sqrt{\frac{1}{n}(\text{Variance}_1 + \text{Variance}_2)}}$$

Quantifying diversity of gut microbiomes



- Visually, microbial taxa distributions in obese children are clearly more diverse
- Can we quantify this pattern?
 - Number of taxa
 - $\text{Entropy} = -\text{Frequency}_1 \log_2(\text{Frequency}_1) - \dots - \text{Frequency}_n \log_2(\text{Frequency}_n)$

Entropy quantifies purity of a mixture



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

- **Entropy** = $-\text{Freq}_1 \log_2(\text{Freq}_1) - \dots - \text{Freq}_n \log_2(\text{Freq}_n)$

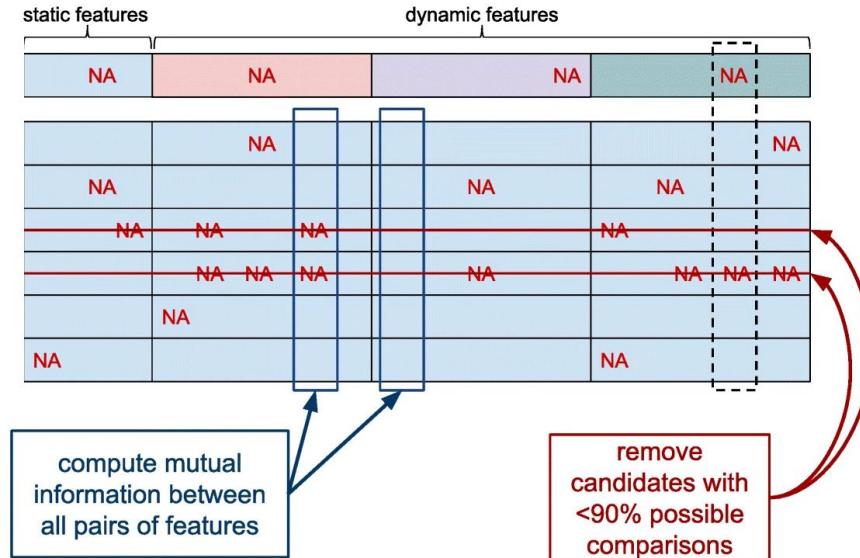
- Left sample: $-\frac{6}{12} \log_2\left(\frac{6}{12}\right) - \frac{6}{12} \log_2\left(\frac{6}{12}\right) = \log_2(2) = 1$

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

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

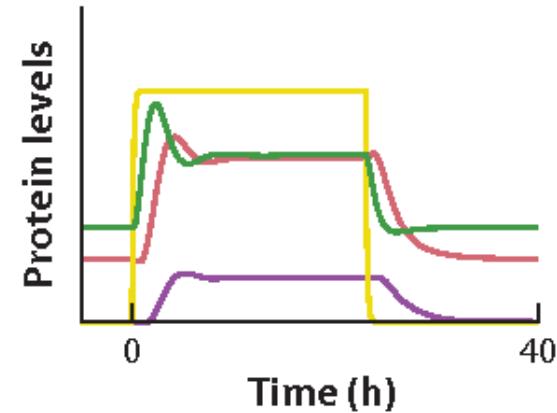
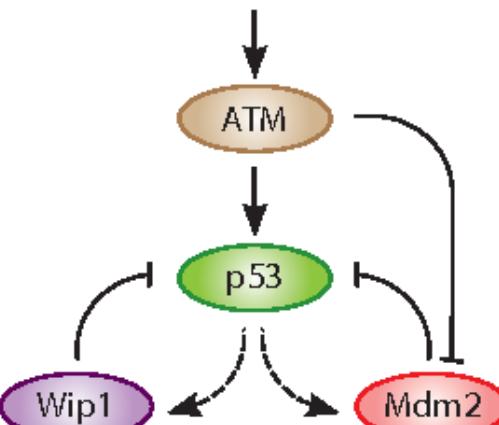
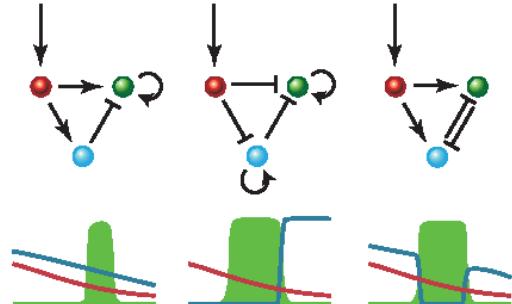
Mutual Information ~ correlation for categorical data

- MI = difference between $P(X, Y)$ and $P(X) P(Y)$
- If X and Y are statistically independent
 - $P(X, Y) = P(X) P(Y)$
 - $MI = 0$
- $MI(X; Y) = \sum_x \sum_y p(x, y) \log \left(\frac{p(x, y)}{p(x)p(y)} \right)$

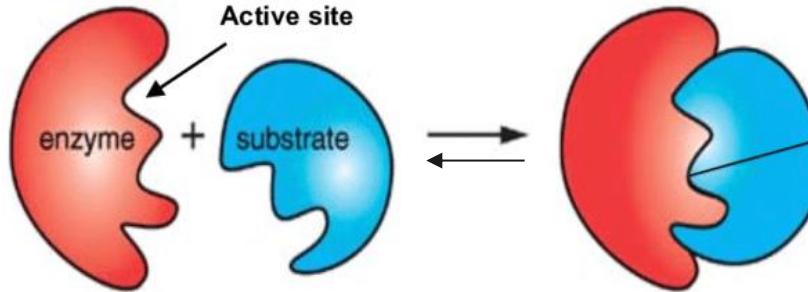


What can you do with computational thinking?

- Data analysis
- Modeling and simulation
 - Identify mechanisms that underlie the phenomenon or system of interest
 - Develop models
 - Study the behavior of the models / systems
 - Synthesize new data
- Algorithm



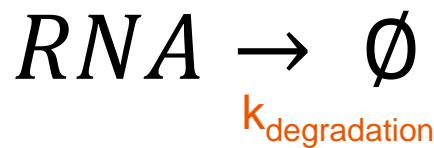
Enzyme substrate binding



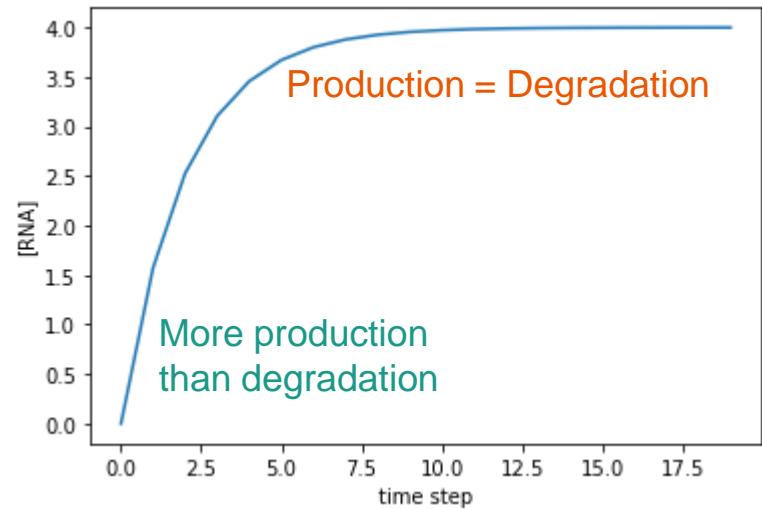
Graham Hutchings. "Development of new highly active nano gold catalysts for selective oxidation reactions" (2014)

- $K_{\text{dissociation}} = [E][S] / [E-S]$ at equilibrium
- Association = $P(E \text{ meeting } S) \times P_{\text{binding}} = k_1 [E][S]$
- Dissociation = Number of E-S molecules $\times P_{\text{dissociating}} = k_2 [E-S]$
- At equilibrium, Association = Dissociation
 - $k_1 [E][S] = k_2 [E-S] \rightarrow K_{\text{dissociation}} = k_2 / k_1$

A simple gene expression model

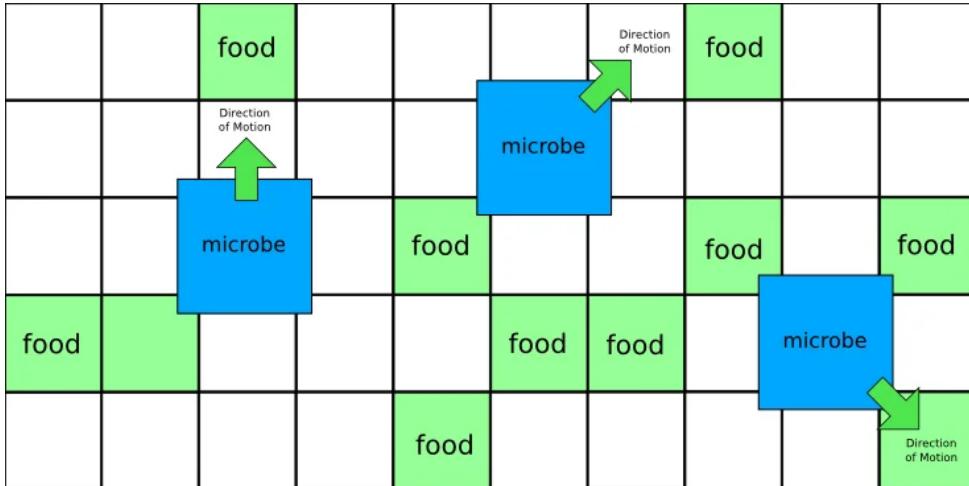


$k_{degradation}$

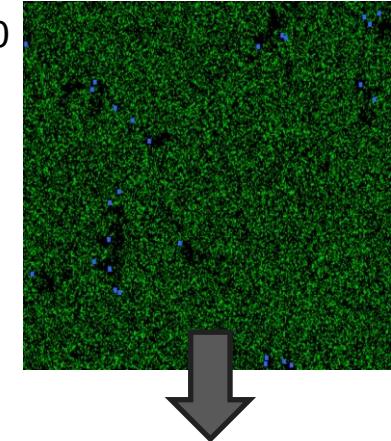


- $\frac{d[RNA]}{dt} = k_{transcription} - k_{degradation}[RNA]$
- This is called an **ordinary differential equation**

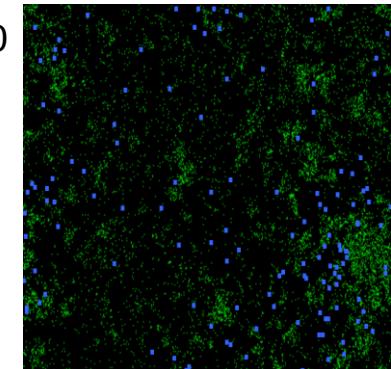
Microbial growth



Time = 0

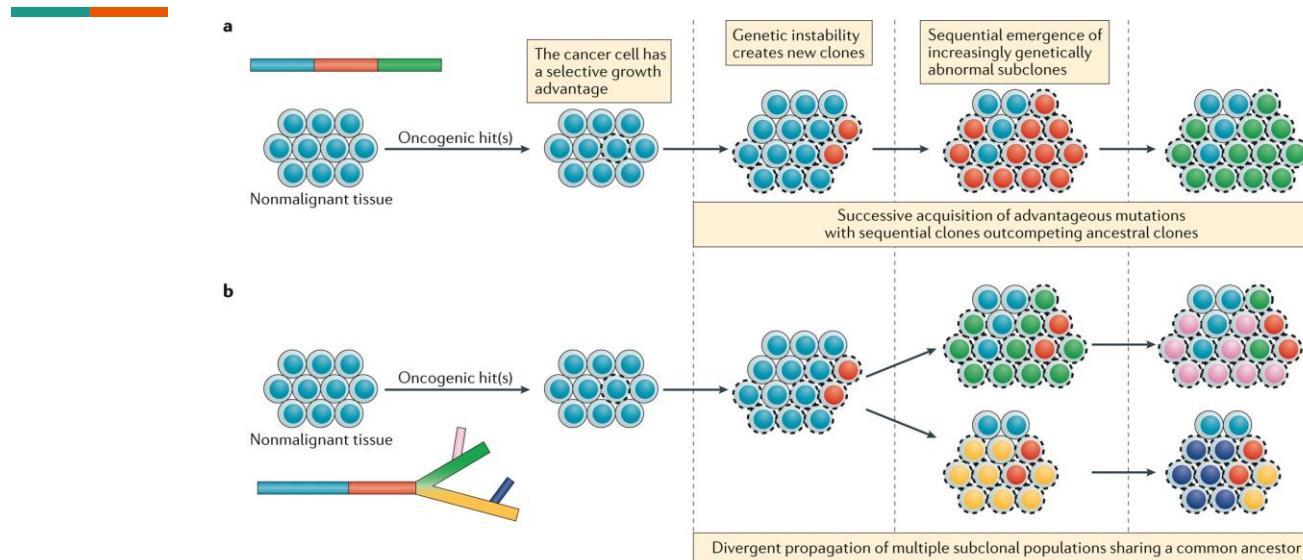


Time = 100



- Model different movement behaviors (probability of moving in a certain direction) with different genes
- See more at https://beltoforion.de/en/simulated_evolution/

Tumor growth

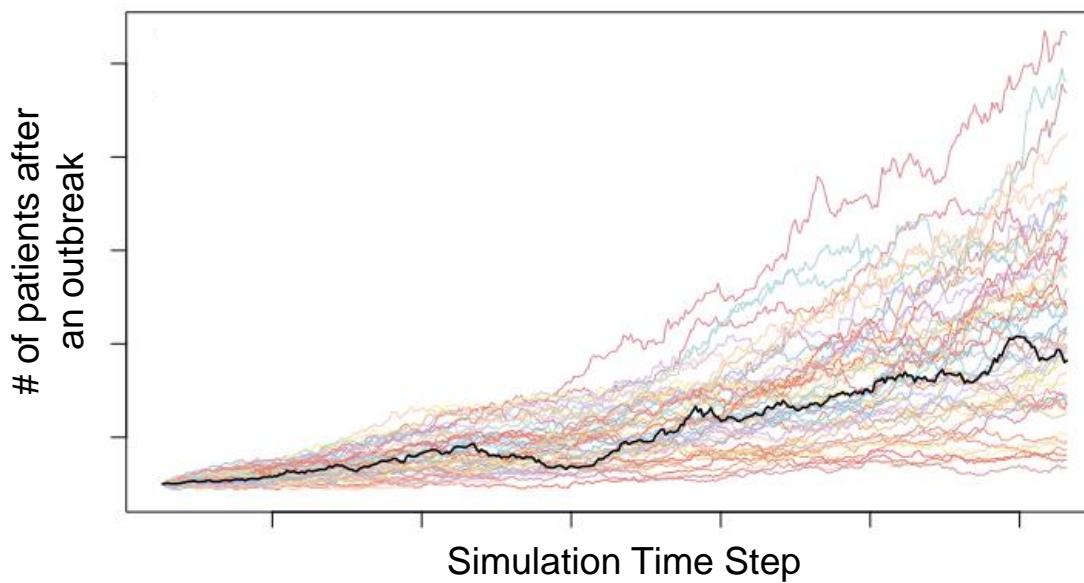


Dagogo-Jack and Shaw, Nat Rev Clin Oncol (2017)

Nature Reviews | Clinical Oncology

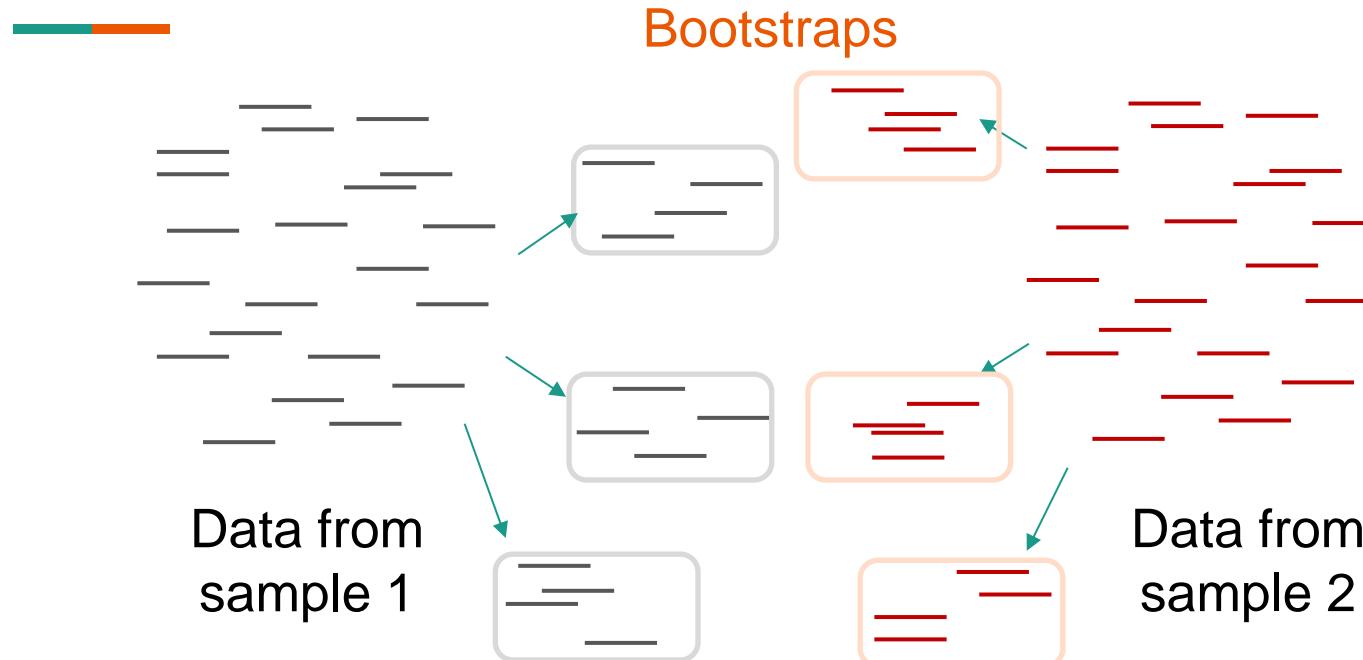
- Cell's actions: Gain mutation, Divide, Die
 - What are the parameters influencing these actions?

Monte Carlo technique



- Perform **repeated sampling** to explore **broad parameter scenarios**
- Provide **estimates for the probability** of different scenarios and outcomes
- May require **millions of repetition** (can utilize multiple CPUs)

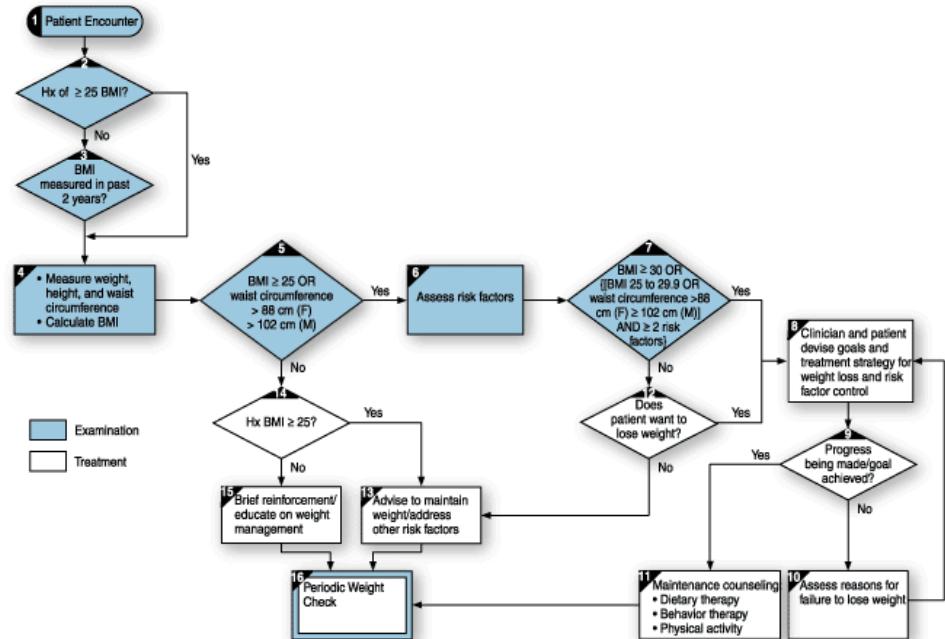
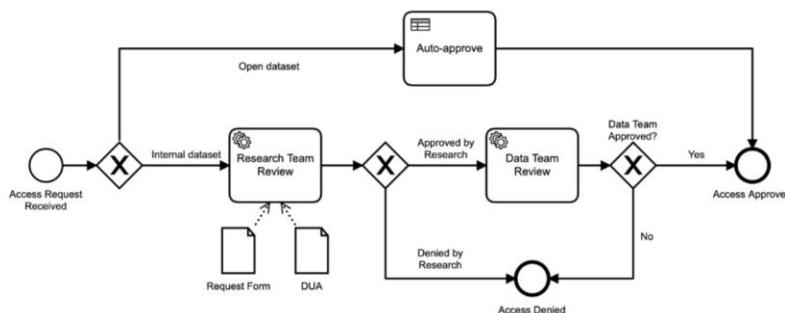
Bootstrapping



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

What can you do with computational thinking?

- Data analysis
- Modeling and simulation
- Algorithm
 - Formulating step-by-step instructions
 - Identifying weak points



* This algorithm applies only to the assessment for overweight and obesity and subsequent decisions based on that assessment. It does not include any initial overall assessment for cardiovascular risk factors or diseases that are indicated.

Knapsack problem

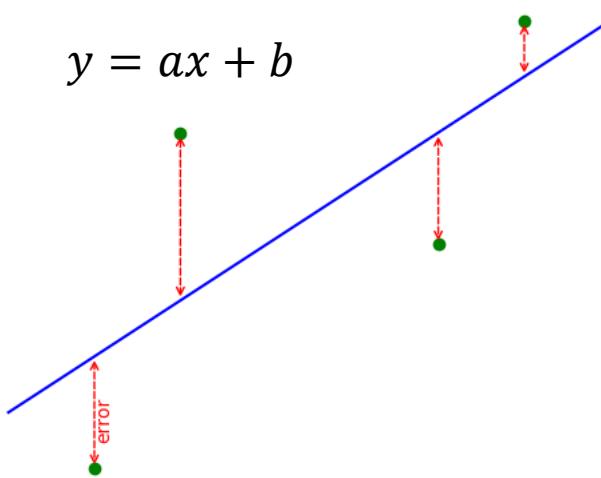


- Limited capacity
- Want to pick as much value as possible while not exceeding the capacity
- General optimization setting
 - **Objective** = Total value
 - **Constraint** = Total weight $\leq W$
- Paying **X** baht with the minimal number of coins and bank notes

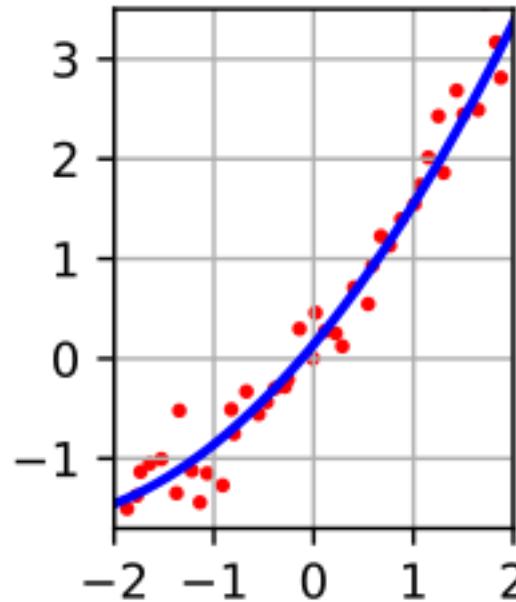
Curve fitting with least square

QUESTION

$$y = ax + b$$



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



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

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

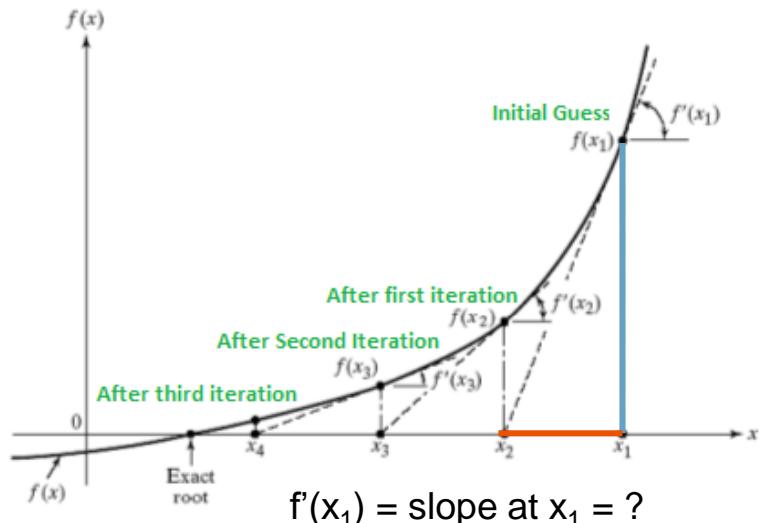
Newton-Ralphson method

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

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

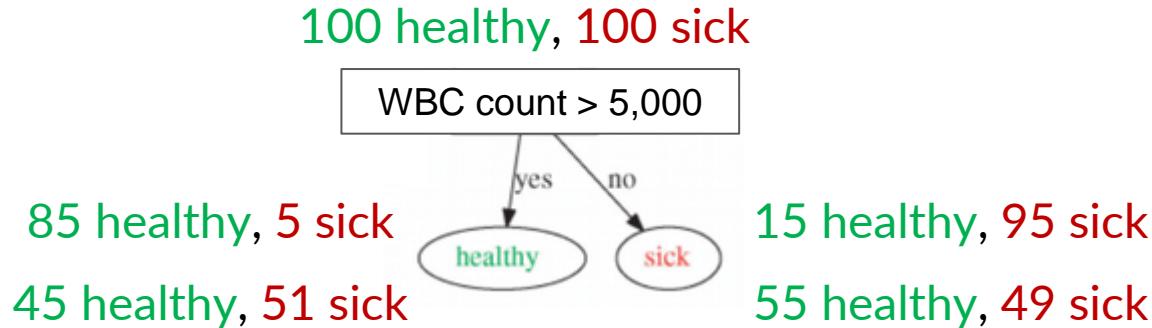
- Start with an initial guess x_1
- Calculate $f(x)$ and $f'(x)$ at x_1
- Define $x_2 = x_1 - \frac{f(x_1)}{f'(x_1)}$
- Repeat the process for x_3, x_4, \dots
- Stop when x_i converges to a value

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



$f'(x_1) = \text{slope at } x_1 = ?$

Constructing a decision point



- **Gini impurity:** $\sum p(1 - p)$
- **Entropy:** $-\sum p \ln(p)$
 - Minimal at $p = 0$ or $1 \rightarrow$ Perfect split
 - Maximal at $p = 0.5 \rightarrow$ 50-50 split
- **Search for variable and cutoff that strongly reduce impurity**

(Almost) all algorithms involve optimization

- **Curve fitting** = minimize square error
- **Drug-protein docking** = minimize energy
- **Hospital queuing** = minimize wait time
- **Diagnosis** = maximize accuracy
- **Triaging** = maximize number of severe patients at the top of the list

Module 1: DNA sequencing & applications

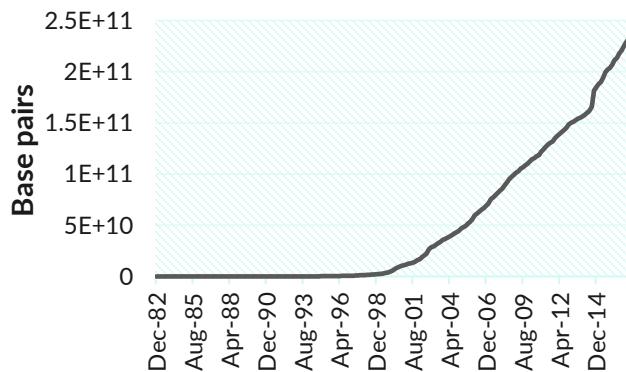
- What do you think kickstarted modern-day computational biology?
- First draft of human genome, 26 June 2000
 - BLAST sequence alignment
- Gene structure annotation
 - Exome sequencing
 - Oligo nucleotide microarray



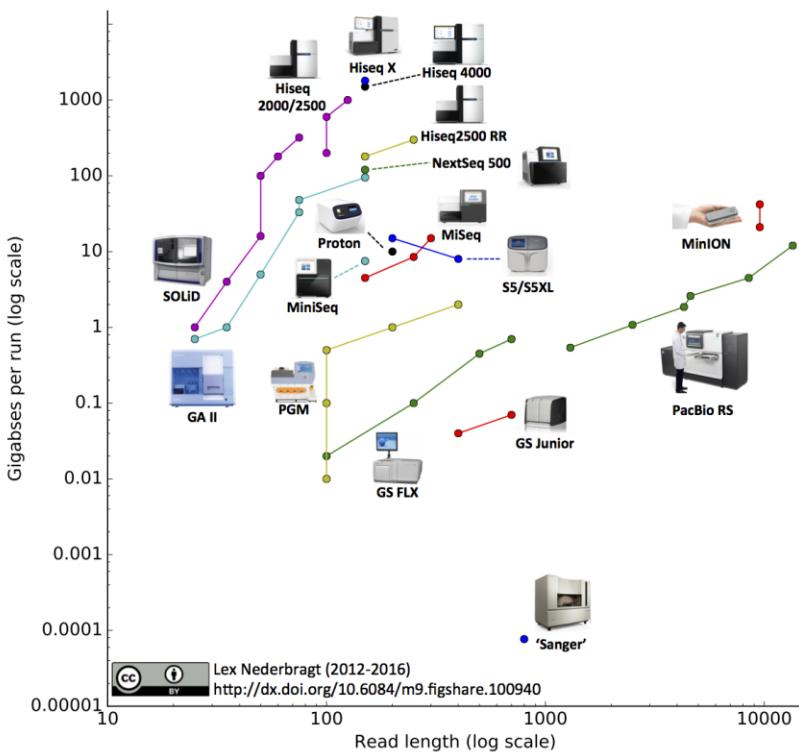
Improvements in DNA sequencing



Genome data on NCBI

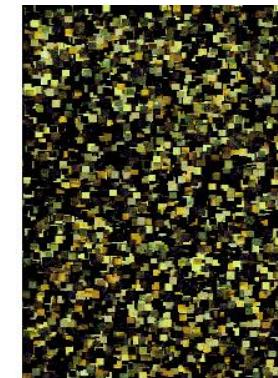
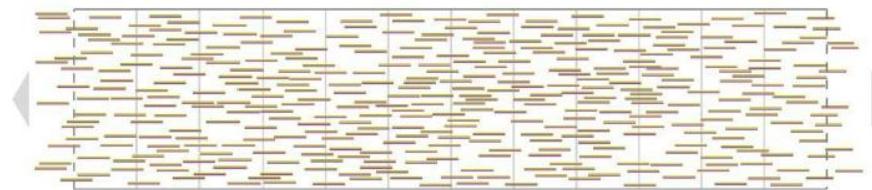


Plotted with data from
<https://www.ncbi.nlm.nih.gov/genbank/statistics/>

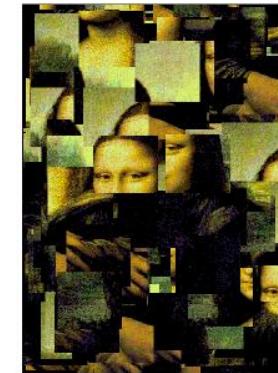
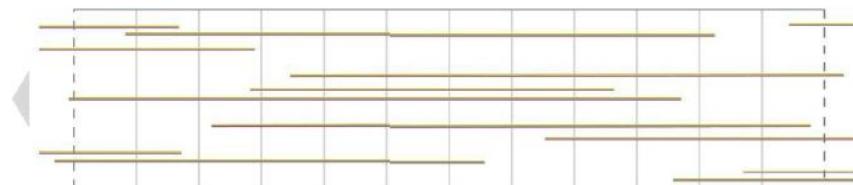


Combining short read and long read data

Short Reads

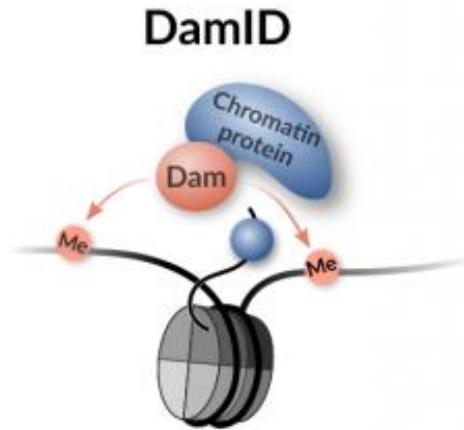


Long Reads

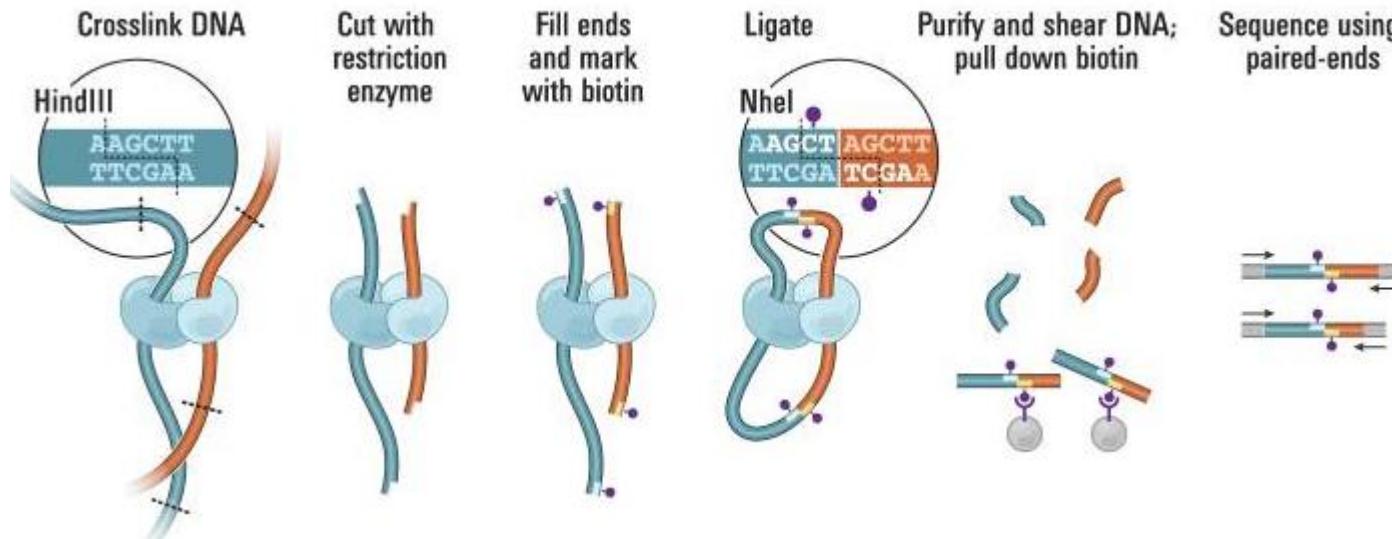


DNA sequencing applications

- Bisulfite-seq
 - DNA methylation
- ChIP-seq, DamID-seq
 - Histone modification, DNA-binding protein
- DNase-seq, MNase-seq, ATAC-seq
 - DNA accessibility
- 3C, 4C, 5C, Hi-C, ChIA-PET
 - Chromatin folding structure

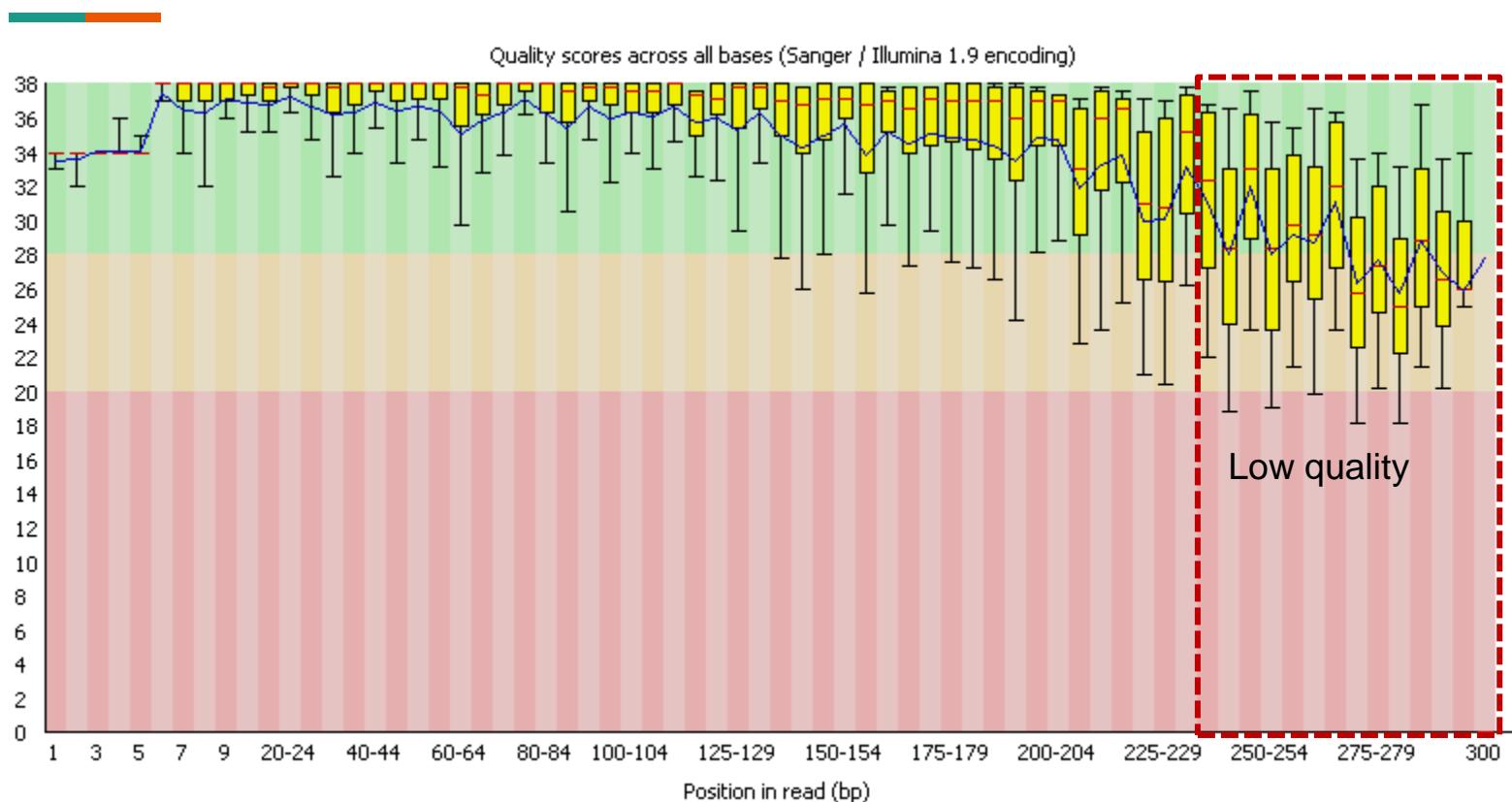


Hi-C: Chromatin folding structure



Lieberman-Aiden *et al.* Science 2009

DNA sequencing QC



Important file formats

```
Coor      12345678901234 5678901234567890123456789012345
ref       AGCATGTTAGATAA**GATAGCTGTGCTAGTAGGCAGTCAGGCCAT

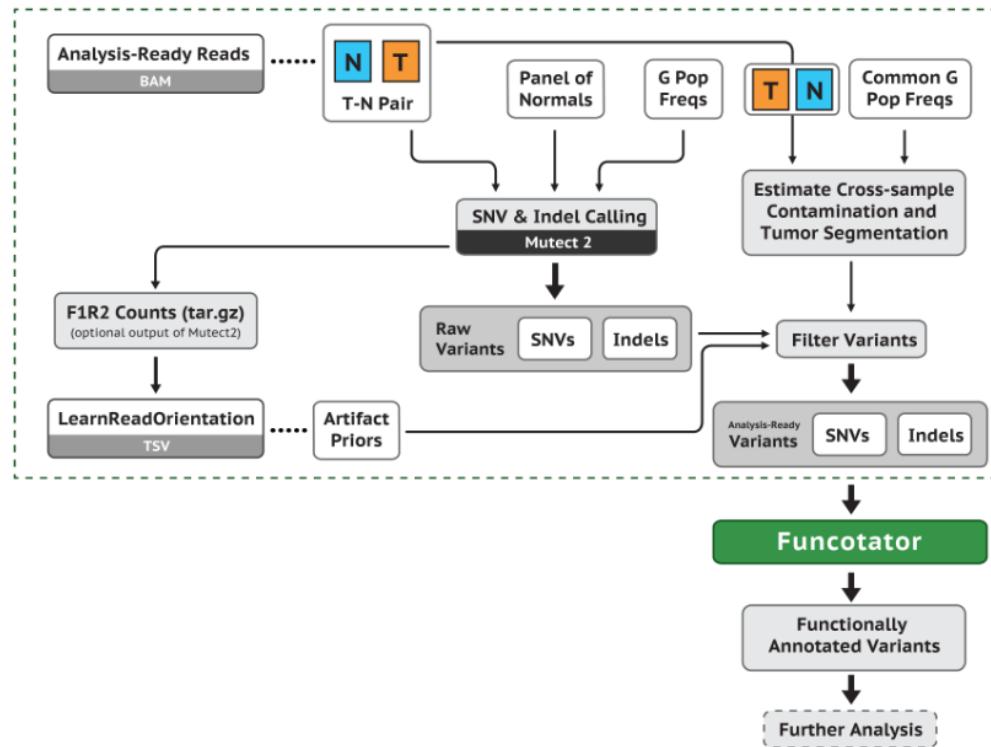
+r001/1    TTAGATAAAGGATA*CTG
+r002    aaaAGATAA*GGATA
+r003    gcctaAGCTAA
+r004    ATAGCT.....TCAGC
-r003    ttagctTAGGC
-r001/2    CAGCGGCAT
```

The corresponding SAM format is:¹

```
@HD VN:1.5 SO:coordinate
@SQ SN:ref LN:45
r001 99 ref 7 30 8M2I4M1D3M = 37 39 TTAGATAAAGGATACTG *
r002 0 ref 9 30 3S6M1P1I4M * 0 0 AAAAGATAAGGATA *
r003 0 ref 9 30 5S6M * 0 0 GCCTAAGCTAA * SA:Z:ref,29,-,6H5M,17,0;
r004 0 ref 16 30 6M14N5M * 0 0 ATAGCTTCAGC *
r003 2064 ref 29 17 6H5M * 0 0 TAGGC * SA:Z:ref,9,+,5S6M,30,1;
r001 147 ref 37 30 9M = 7 -39 CAGCGGCAT * NM:i:1
```

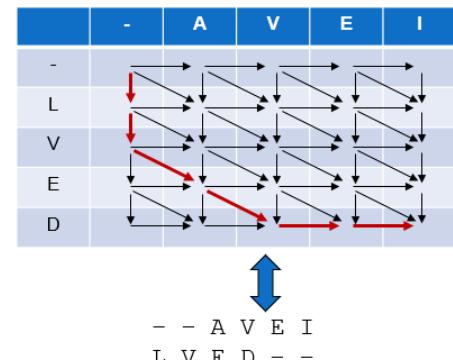
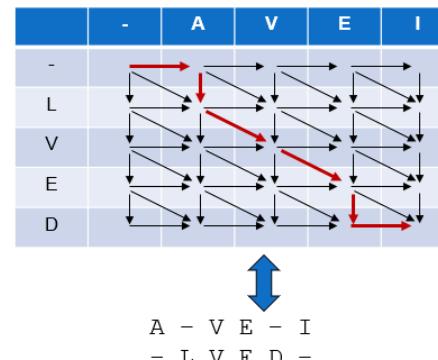
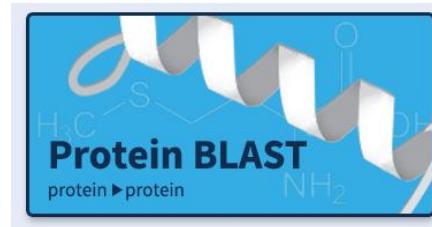
- SAM, BAM, BED, GTF, GFF, FASTA, FASTQ, VCF, MAF

GATK variant calling pipelines



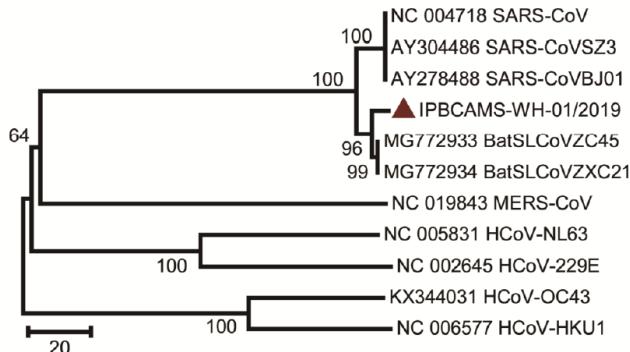
Basic Local Alignment Search Tool (BLAST)

- How does it work?
 - How should we adjust its parameters?
- What questions can it answer?
 - Functional annotation
 - Taxonomy annotation
 - A start for evolutionary analysis
- Variants of BLAST
 - MEGABLAST
 - PSI-BLAST



Phylogenetics

- Evolution is informative
- So many parameters and models
 - How to properly analyze?
 - Which tools to use?



Guo et al. Clin Infect Dis (2020)

MX: Analysis Preferences

Phylogeny Reconstruction

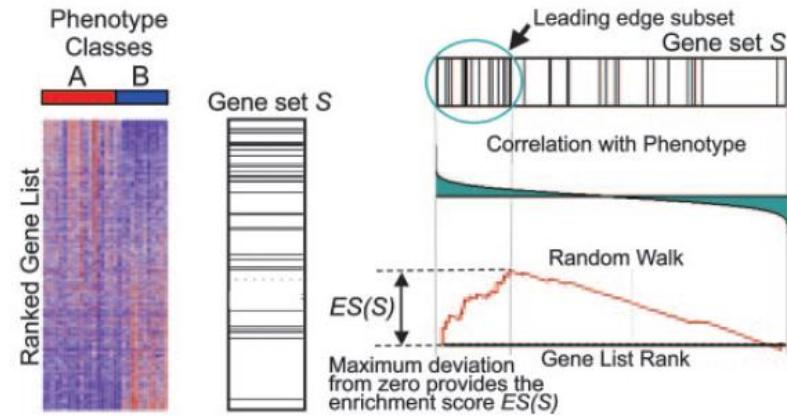
Option	Setting
ANALYSIS	
Statistical Method	→ Maximum Likelihood
PHYLOGENY TEST	
Test of Phylogeny	→ None
No. of Bootstrap Replications	→ None Bootstrap method
SUBSTITUTION MODEL	
Substitutions Type	→ Nucleotide
Model/Method	→ Tamura-Nei model
RATES AND PATTERNS	
Rates among Sites	→ Uniform Rates
No of Discrete Gamma Categories	→ Not Applicable
DATA SUBSET TO USE	
Gaps/Missing Data Treatment	→ Use all sites
Site Coverage Cutoff (%)	→ Not Applicable
TREE INFERENCE OPTIONS	
ML Heuristic Method	→ Nearest-Neighbor-Interchange (NNI)
Initial Tree for ML	→ Make initial tree automatically (Default - NJ/BioNJ)
Initial Tree File	→ Not Applicable
Branch Swap Filter	→ None
SYSTEM RESOURCE USAGE	
Number of Threads	→ 7

Module 2: Transcriptomics

- Microarray vs RNA-seq vs Nanostring
 - Different data models and computational analyses
- Pseudoalignment of RNA-seq
 - kallisto, salmon
- StringTie's hybrid alignment and *de novo* assembly pipeline
- Pairing RNA-seq processing and differential expression analysis tools
 - HTSeq-count with DESeq2
 - salmon with sleuth

Functional enrichment analysis

- Overrepresentation analysis
 - Frequent functional terms
 - Hypergeometric distribution
- Gene set enrichment analysis (GSEA)
 - Random walk model
- Gene-gene network topology
 - Frequent functional terms
 - Nearby gene on network



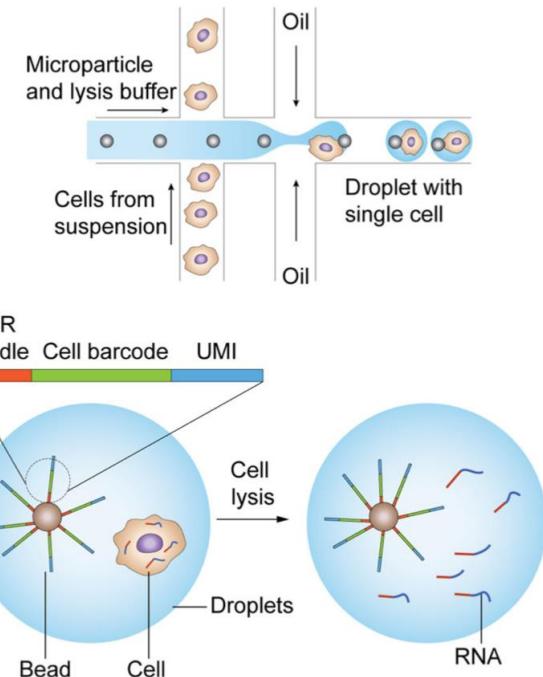
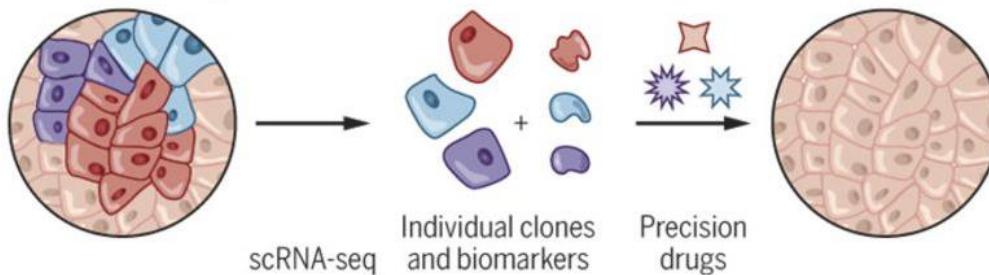
Subramanian *et al.* PNAS. 102:15545-15550 (2005)

Toward single-cell technology

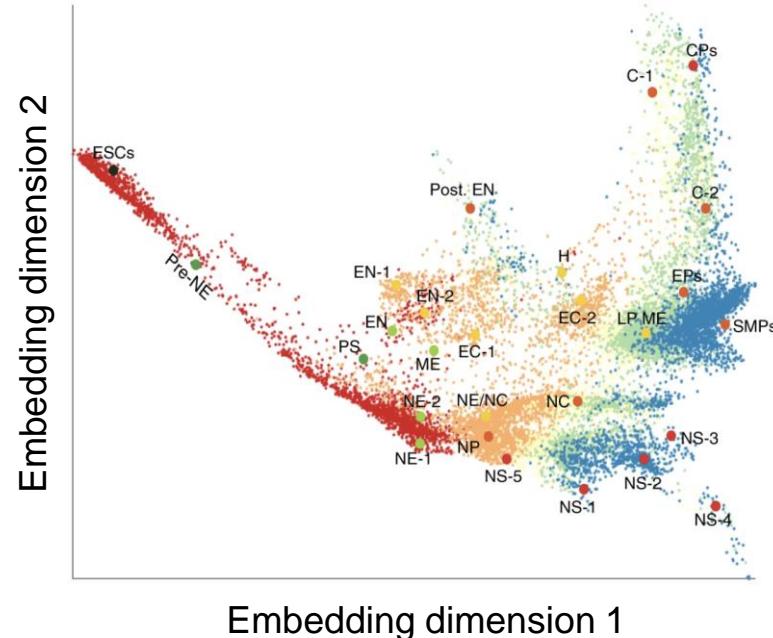
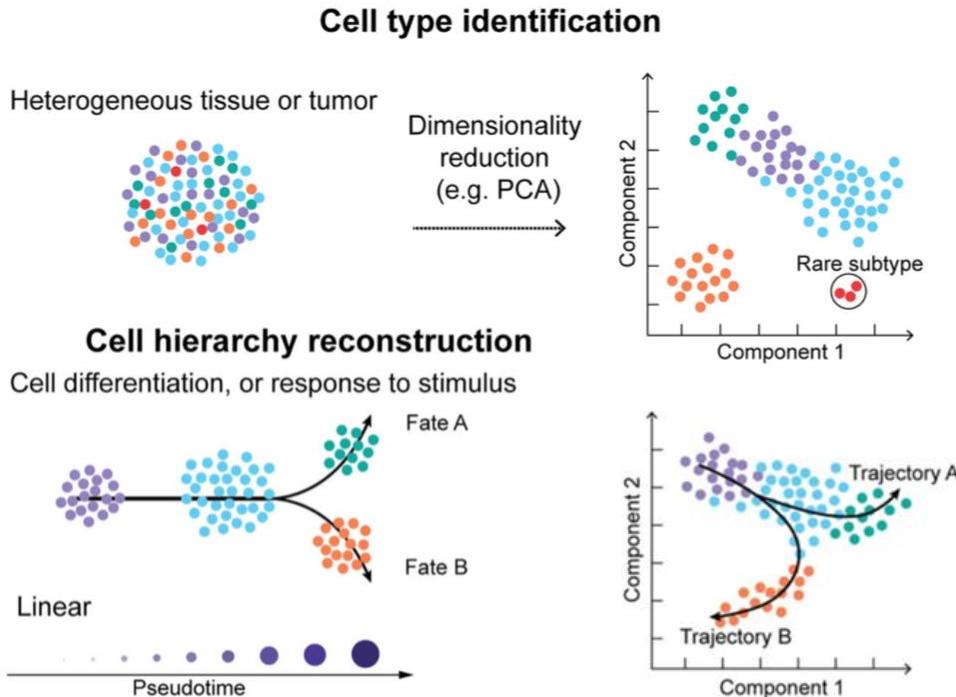
Bulk analysis



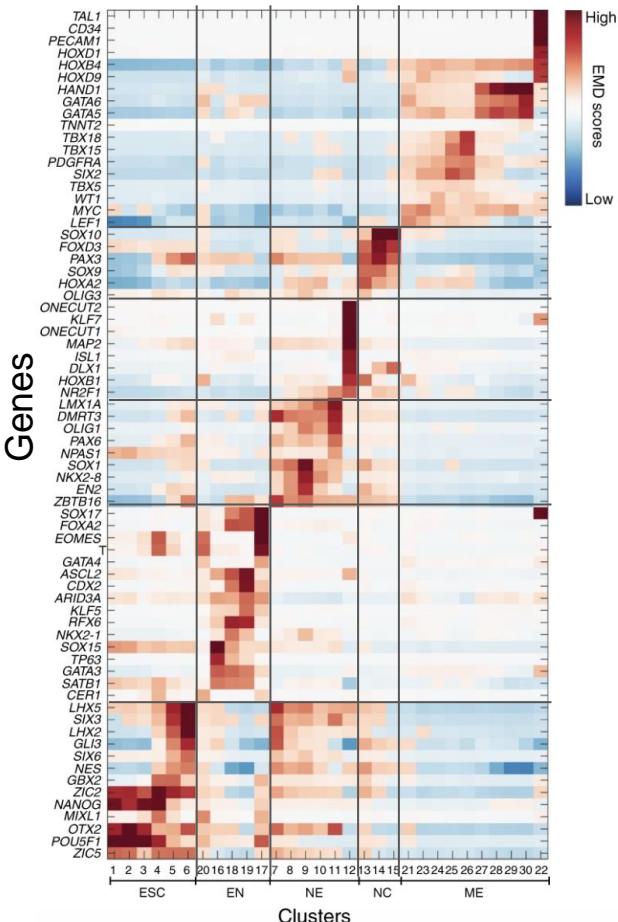
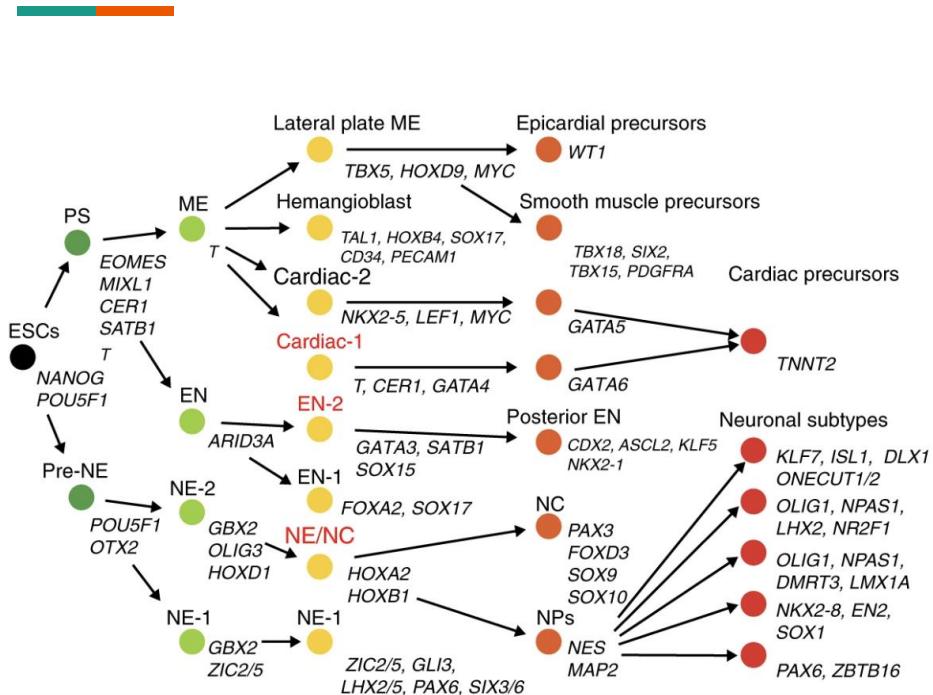
scRNA analysis



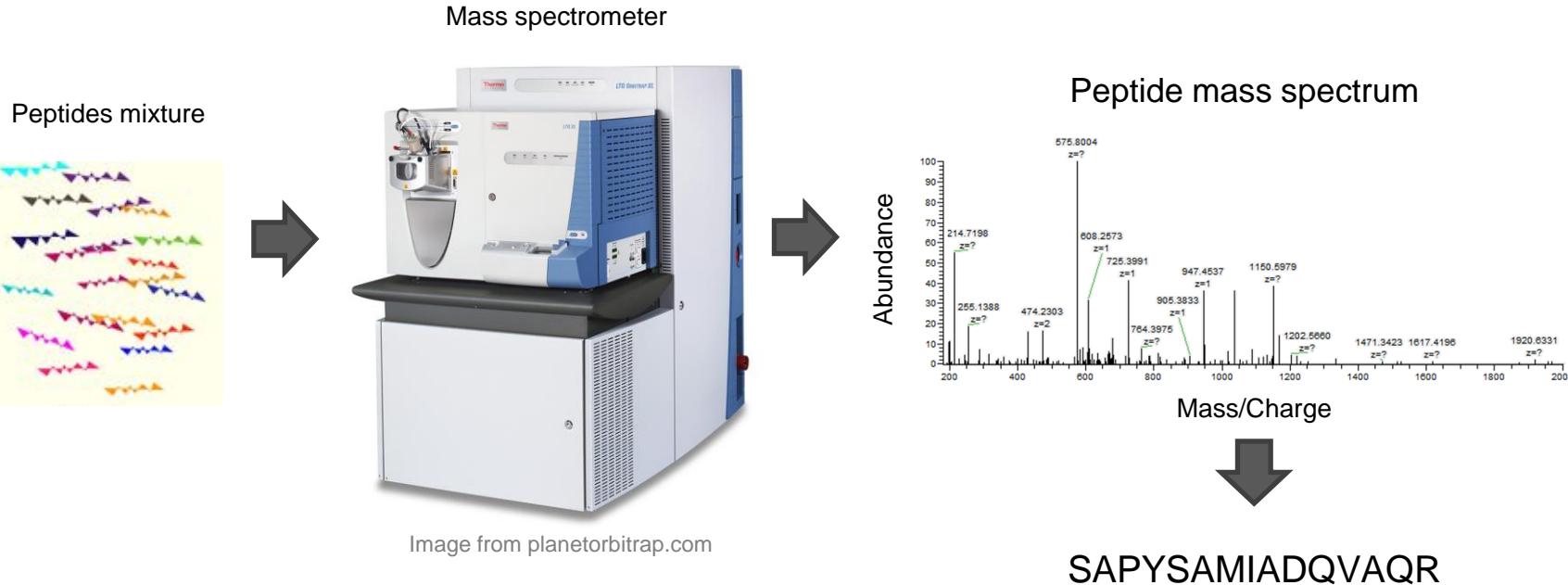
Power of single-cell data



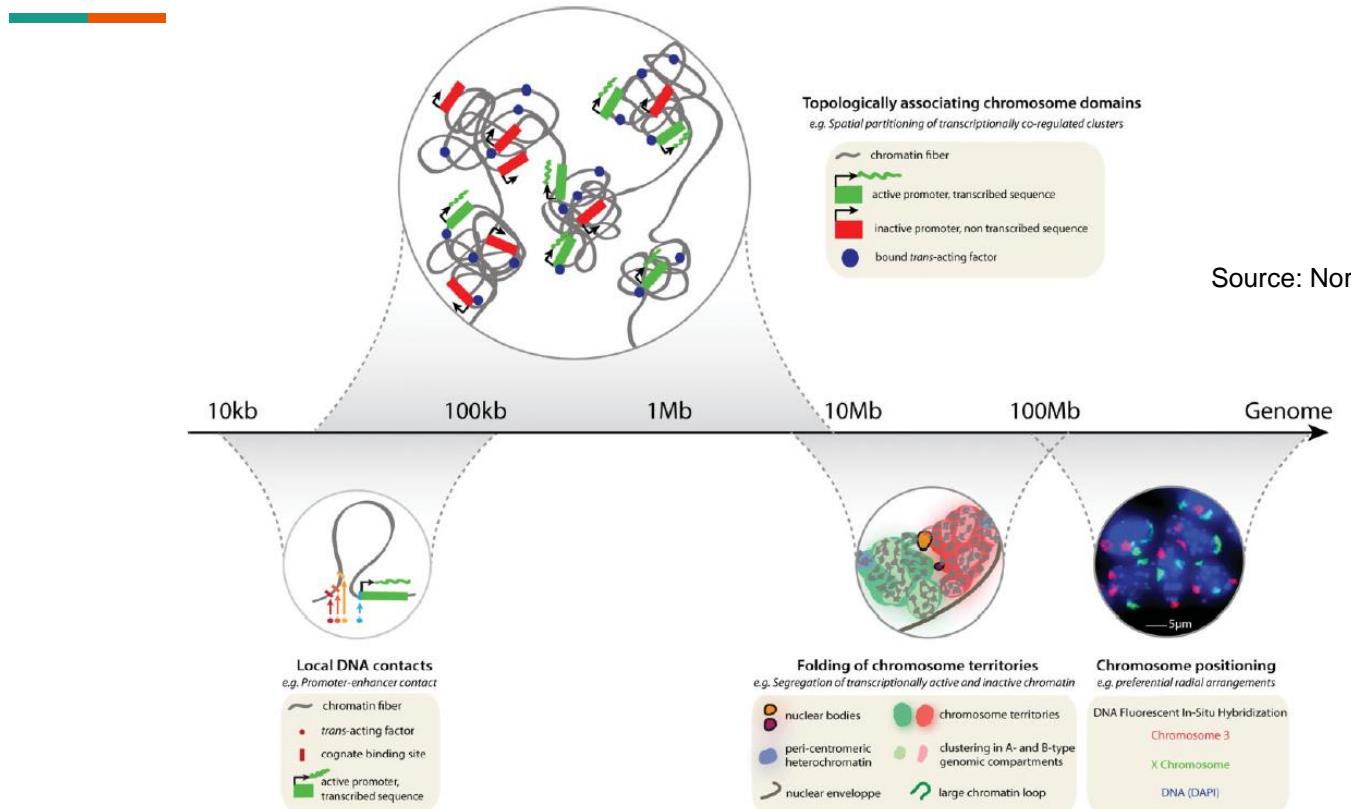
Power of single-cell data



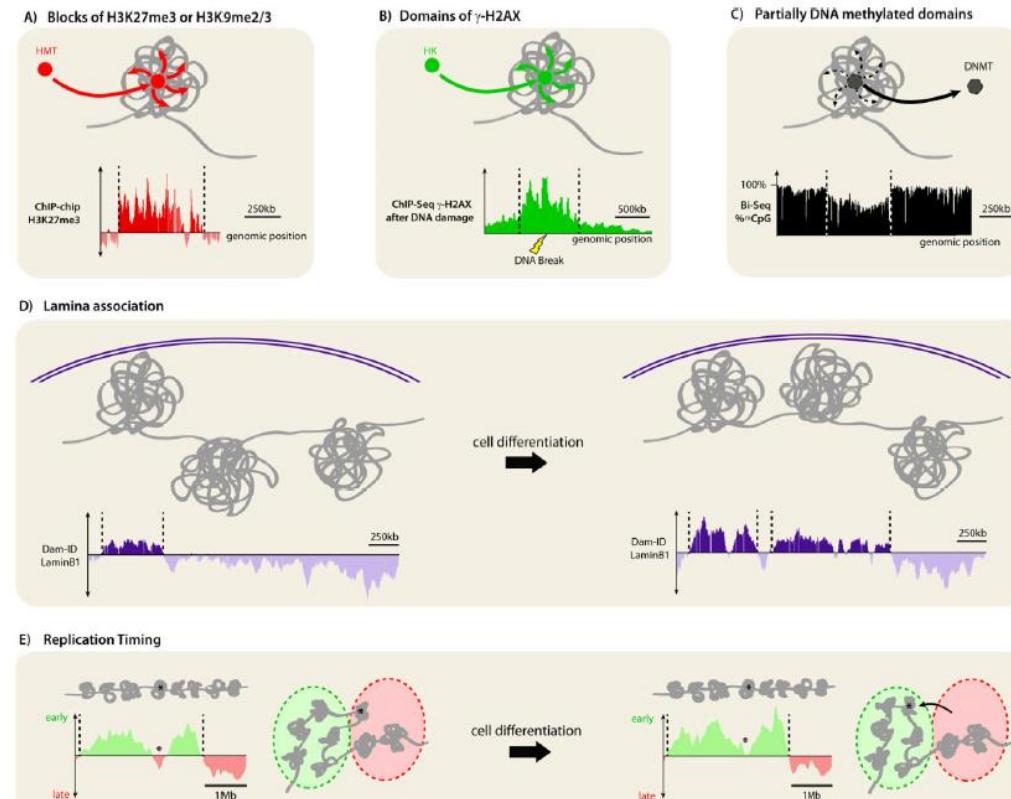
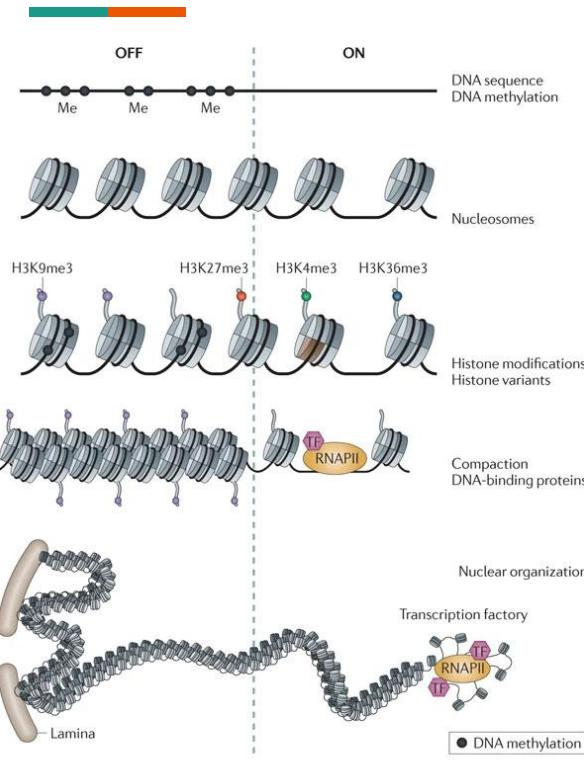
Module 3: Other omics



Chromatin organization

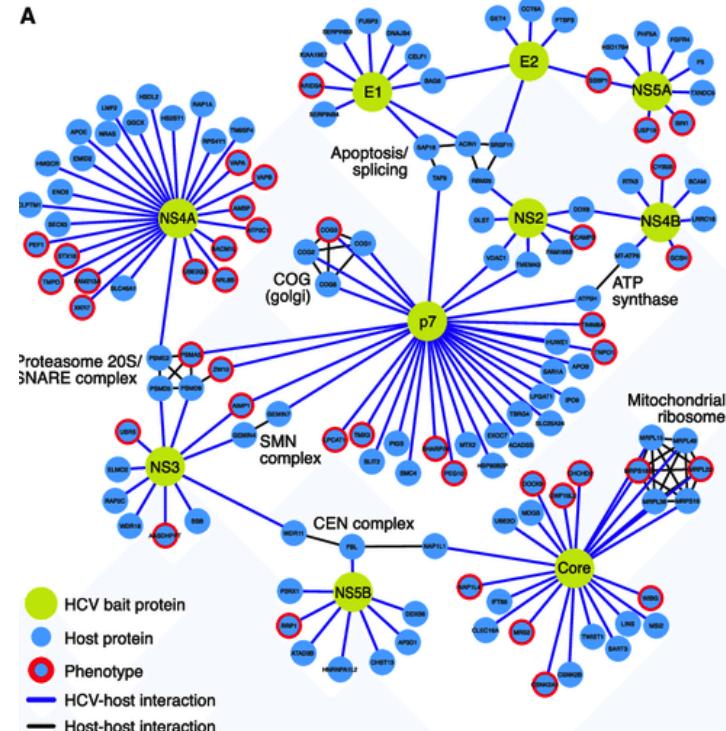
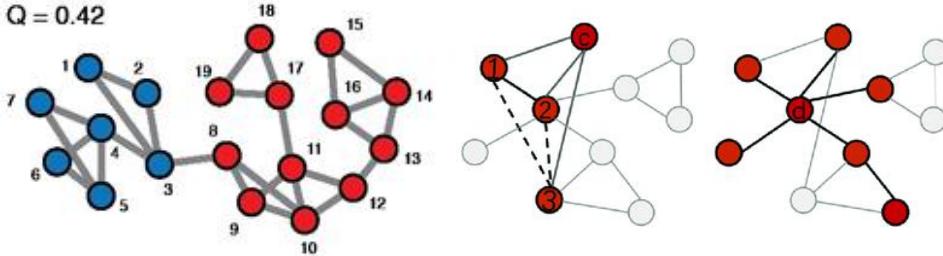


Epigenetics and gene expression regulation



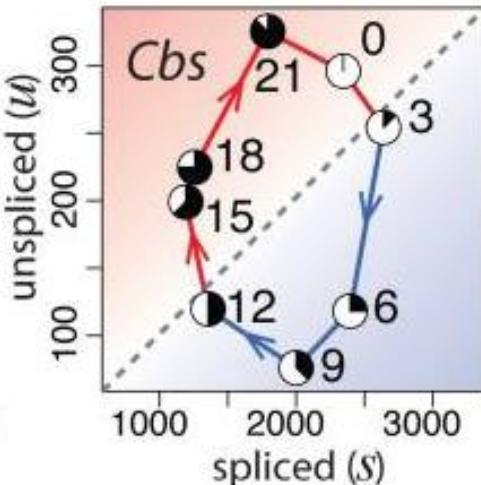
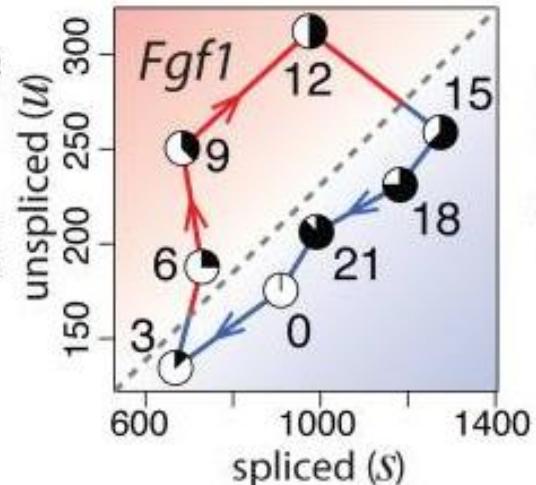
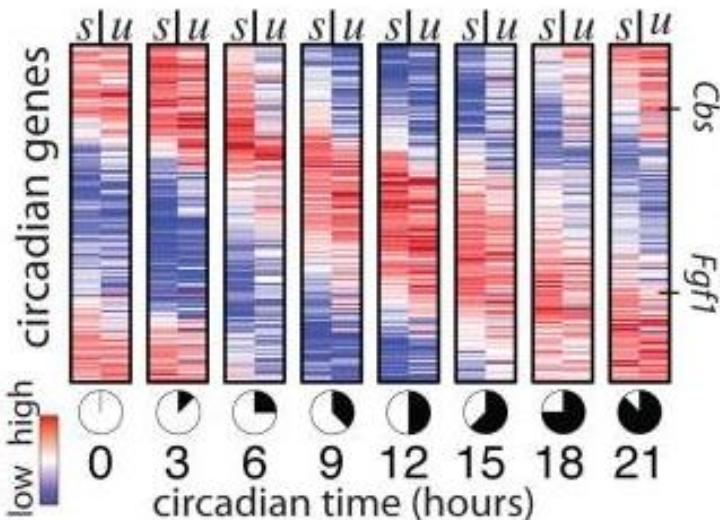
Biological networks

- Topological properties of networks
 - Relationship to biology
- Applications in biomedicine
- Visualization and analysis with CytoScape



Source: Ramage et al. Mol Cell (2015)

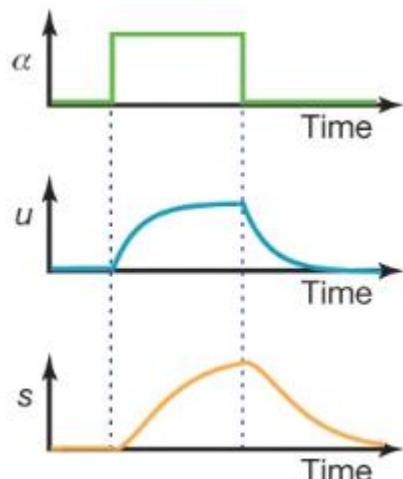
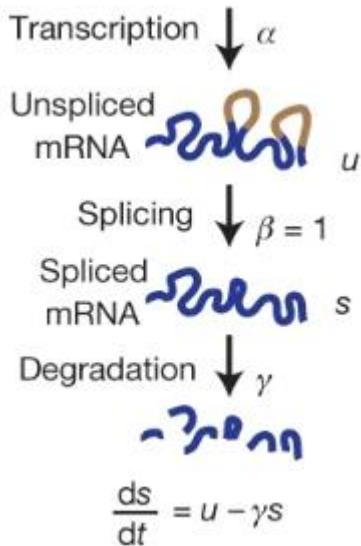
Dynamics modeling in Systems Biology



La Manno *et al.* Nature 2018

- Gene upregulation first produce unsспорced RNA
- Followed by processing into mature RNA and proteins

RNA velocity model



- Model RNA synthesis as differential equations
- Simulation can be performed to analyze the dynamics of the system
 - Try various parameter values

Human Protein Atlas: Protein localization

ACE2



RNA
TISSUE

RNA
BRAIN

RNA
SINGLE CELL

RNA
TISSUE CELL

RNA
PATHOLOGY

RNA
IMMUNE

MS
BLOOD

RNA
SUBCELL

RNA
CELL LINE



PROTEIN SUMMARY

RNA DATA

GENE/PROTEIN

ANTIBODIES AND VALIDATION



HUMAN PROTEIN ATLAS SUMMARYⁱ

Protein ⁱ	Angiotensin I converting enzyme 2
Gene name ⁱ	ACE2
Tissue specificity ⁱ	Tissue enhanced (gallbladder, intestine, kidney)
Tissue expression cluster ⁱ	Intestine & Kidney - Transmembrane transport (mainly)
Single cell type specificity ⁱ	Cell type enriched (Proximal enterocytes)
Single cell type expression cluster ⁱ	Enterocytes - Digestion (mainly)
Immune cell specificity ⁱ	Not detected in immune cells
Brain specificity ⁱ	Not detected in human brain
Cancer prognostic summary	Prognostic marker in renal cancer (favorable) and liver cancer (favorable)
Predicted location ⁱ	Membrane, Secreted (different isoforms)
Extracellular location ⁱ	Secreted to blood

MSigDB: Curated gene sets

Human MSigDB Collections



The 33196 gene sets in the Human Molecular Signatures Database (MSigDB) are divided into 9 major collections, and several sub-collections. See the table below for a brief description of each, and the [Human MSigDB Collections: Details and Acknowledgments](#) page for more detailed descriptions. See also the [MSigDB Release Notes](#).

Click on the "browse gene sets" links in the table below to view the gene sets in a collection. Or download the gene sets in a collection by clicking on the links below the "Download Files" headings. For a description of the [GMT file format](#) see the [Data Formats](#) in the [Documentation section](#). The gene sets can be downloaded as NCBI (Entrez) Gene Identifiers or HUGO (HGNC) Gene Symbols. There are also JSON bundles containing the HUGO (HGNC) Gene Symbols along with some useful metadata. An XML file containing all the Human MSigDB gene sets is available as well.

H: hallmark gene sets (browse 50 gene sets)

Hallmark gene sets summarize and represent specific well-defined biological states or processes and display coherent expression. These gene sets were generated by a computational methodology based on identifying overlaps between gene sets in other MSigDB collections and retaining genes that display coordinate expression. [details](#)

[Download GMT Files](#)
[Gene Symbols](#)
[NCBI \(Entrez\) Gene IDs](#)
[JSON bundle](#)

C1: positional gene sets (browse 299 gene sets)

Gene sets corresponding to human chromosome cytogenetic bands. [details](#)

[Download GMT Files](#)
[Gene Symbols](#)
[NCBI \(Entrez\) Gene IDs](#)
[JSON bundle](#)

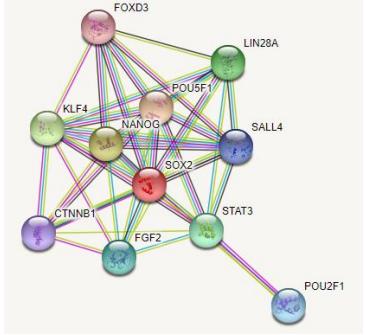
C2: curated gene sets (browse 6449 gene sets)

Gene sets in this collection are curated from various sources, including online pathway databases and the biomedical literature. Many sets are also contributed by individual domain experts. The gene set page for each gene set lists its source. The C2 collection is divided into the following two sub-collections: Chemical and genetic perturbations (CGP) and Canonical pathways (CP). [details](#)

[Download GMT Files](#)
[Gene Symbols](#)
[NCBI \(Entrez\) Gene IDs](#)
[JSON bundle](#)

- Function-based / Disease-based / Publication-based / Genomic location-based
- For gene panel, enrichment test, etc.

STRING: Protein-protein interaction



Viewers ▾ Legend > Settings > Analysis > Table > More Less

Network
Summary view: shows current interactions. Nodes can be moved; popups provide information on nodes & edges.

Experiments
Co-purification, co-crystallization, Yeast2Hybrid, Genetic Interactions, etc ... as imported from primary sources.

Databases
Known metabolic pathways, protein complexes, signal transduction pathways, etc ... from curated databases.

PubMed
Automated, unsupervised textmining - searching for proteins that are frequently mentioned together.

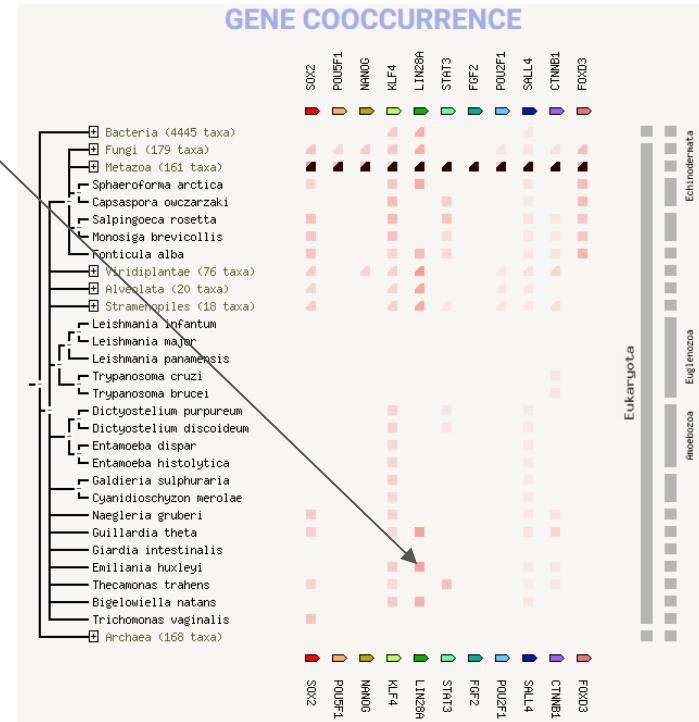
Cooccurrence currently showing Gene families whose occurrence patterns across genomes show similarities.

Coexpression
Proteins whose genes are observed to be correlated in expression, across a large number of experiments.

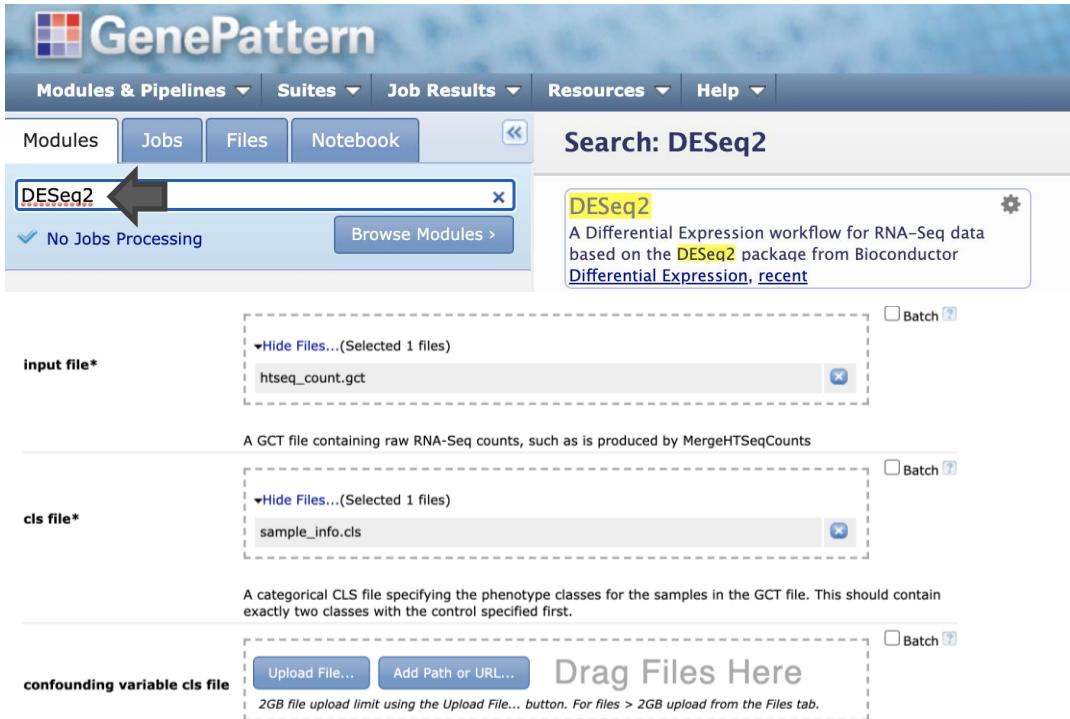
Neighborhood
Groups of genes that are frequently observed in each other's genomic neighborhood.

Fusion
Genes that are sometimes fused into single open reading frames.

LIN28A	25	PEDAARAAD-EPQLLHGAGICKWFNVRMGFGFLSMTARAGVALDPVVDV P A R D EP +G CKWF+V+ GFGF+ + + D+ PNTATRVDPEPA---PSGCKWFDFVKQKGFGFIDVE-----NQEQLD
EOD22827	19	
LIN28A	73	FVHQSKLHMEGFRSLKEGEAVEFTFKKSAGK--LESIRVTGPGGVFCIG FVHQ+ + +GFRSL EGEA+EF + AK L++I VTGGGG F G FVHQTDIKAGFRSLAEGLAEFKVSRDAKTNKLKAIEVTGPGGDVEG
EOD22827	58	
LIN28A	120	SERRP + R P
EOD22827	107	APREP



DESeq2 on GenePattern



The screenshot shows the GenePattern web interface. At the top, there's a navigation bar with links for Modules & Pipelines, Suites, Job Results, Resources, and Help. Below the navigation bar, there are tabs for Modules, Jobs, Files, and Notebook. A search bar contains the text "DESeq2". A message below the search bar says "No Jobs Processing" and "Browse Modules >". A large arrow points to the "DESeq2" entry in the search results. The main area has three input fields: "input file*" containing "htseq_count.gct", "cls file*" containing "sample_info.cls", and "confounding variable cls file" which is currently empty. There are also "Upload File..." and "Add Path or URL..." buttons, and a "Drag Files Here" area with a note about a 2GB file upload limit.

- Input expression data (RNA-seq read count) and sample label
- .gct and .cls are text files

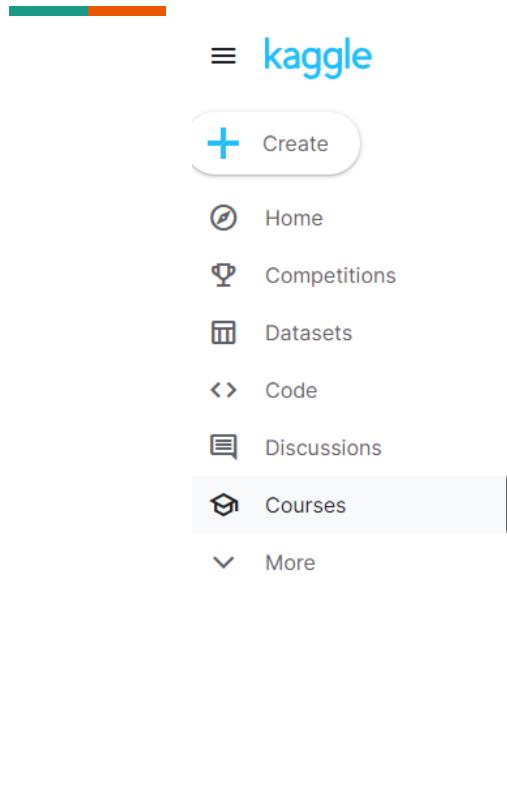
#	Sample ID	Sample Name	Description	sample1	sample2	sample3	sample4	sample5
1	60488	12		2	10	0	2	
2	Name		Description					
3	ENSG00000242268.2	ENSG00000242268.2		2	10	0	2	
4	ENSG00000270112.3	ENSG00000270112.3		8	6	0	0	
5	ENSG00000167578.15	ENSG00000167578.15		102	633	468	1200	
6	ENSG00000273842.1	ENSG00000273842.1		0	0	0	0	
7	ENSG0000078237.5	ENSG0000078237.5		370	364	1220	692	

```
12 2 1
# alive dead
1 1 0 0 1 0 0 0 1 0 1 1
```

Module 4: Programming and machine learning

- Kaggle's programming courses
- In-class practice & problem sets
 - Handling of tabular data
 - Statistical analyses
 - Data exploration
 - Visualization

Kaggle's programming



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 a three-line menu icon to its left. A 'Create' button with a plus sign is highlighted with a light orange rounded rectangle. The sidebar contains several links with icons: Home (compass), Competitions (trophy), Datasets (bar chart), Code (code editor), Discussions (comment), Courses (book), and More (down arrow). A vertical line separates the sidebar from the main content area.

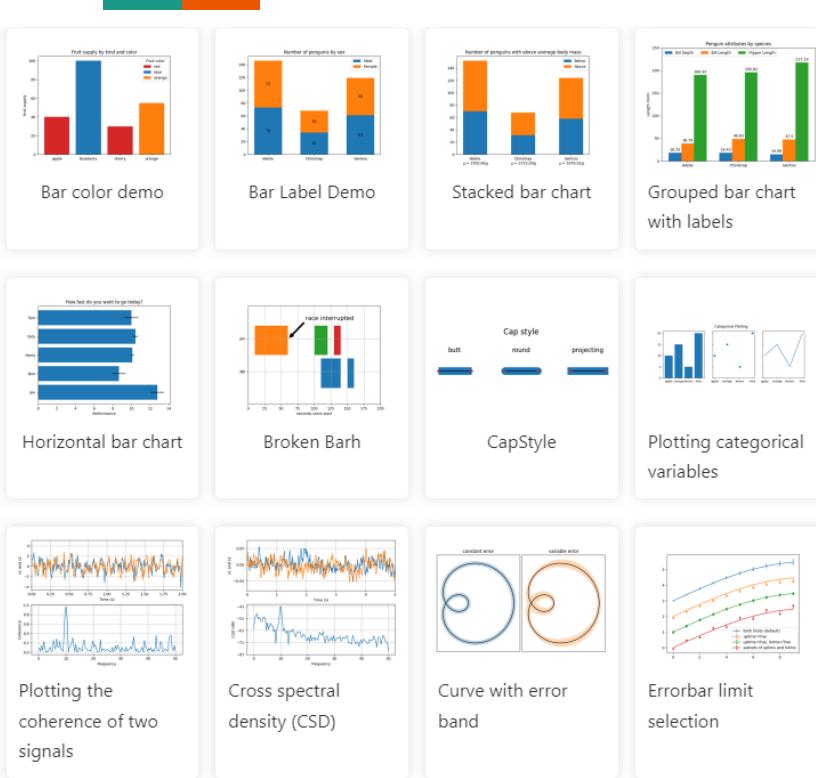


A search bar with a magnifying glass icon and the word 'Search'.

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!

Matplotlib (<https://matplotlib.org/stable/>)



```
import matplotlib.pyplot as plt
import numpy as np

# data from https://allisonhorst.github.io/palmerpenguins/

species = (
    "Adelie\n $\mu=$3700.66g",
    "Chinstrap\n $\mu=$3733.09g",
    "Gentoo\n $\mu=5076.02g$",
)
weight_counts = {
    "Below": np.array([70, 31, 58]),
    "Above": np.array([82, 37, 66]),
}
width = 0.5

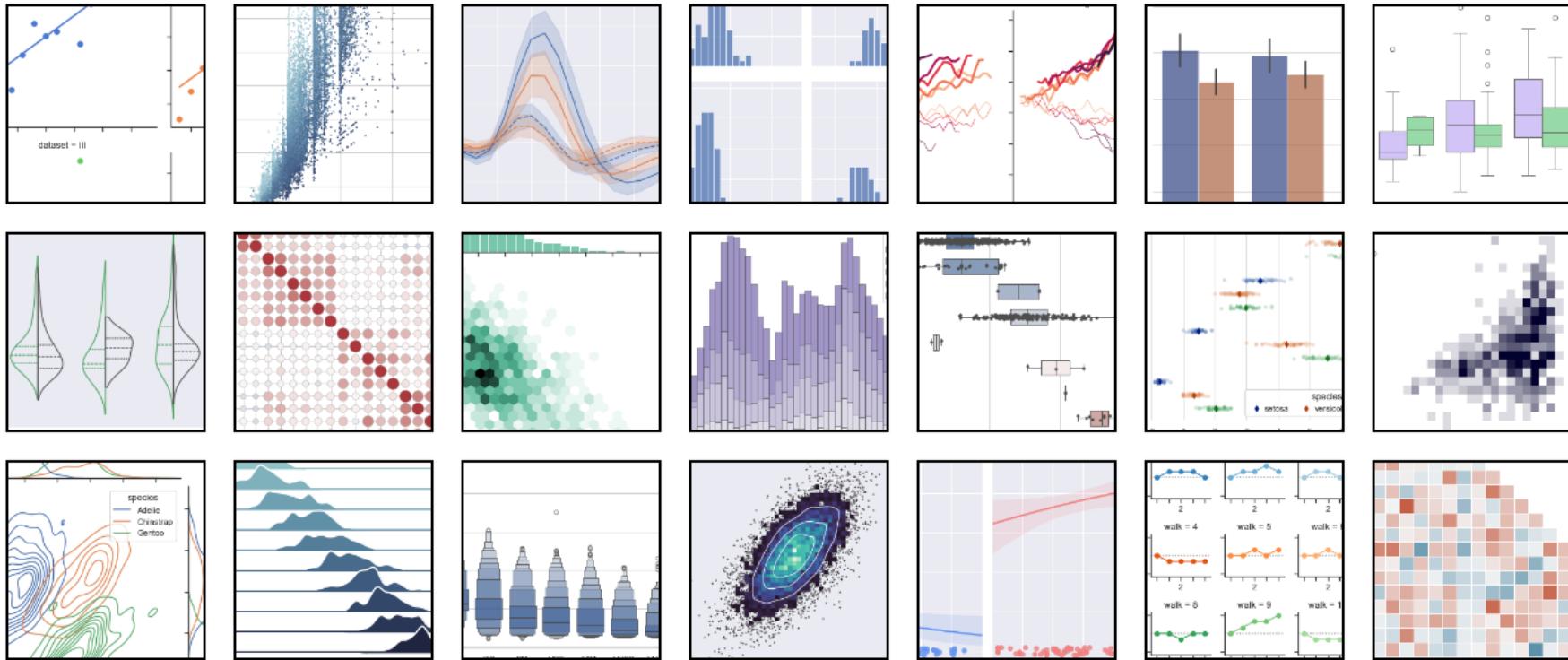
fig, ax = plt.subplots()
bottom = np.zeros(3)

for boolean, weight_count in weight_counts.items():
    p = ax.bar(species, weight_count, width, label=boolean, bottom=bottom)
    bottom += weight_count

ax.set_title("Number of penguins with above average body mass")
ax.legend(loc="upper right")

plt.show()
```

Seaborn (<https://seaborn.pydata.org>)

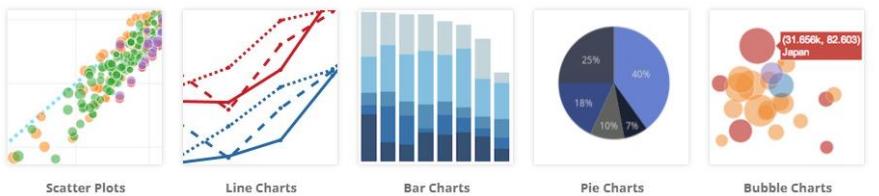


Plotly (<https://plotly.com/python/>)

Fundamentals



Basic Charts



Statistical Charts

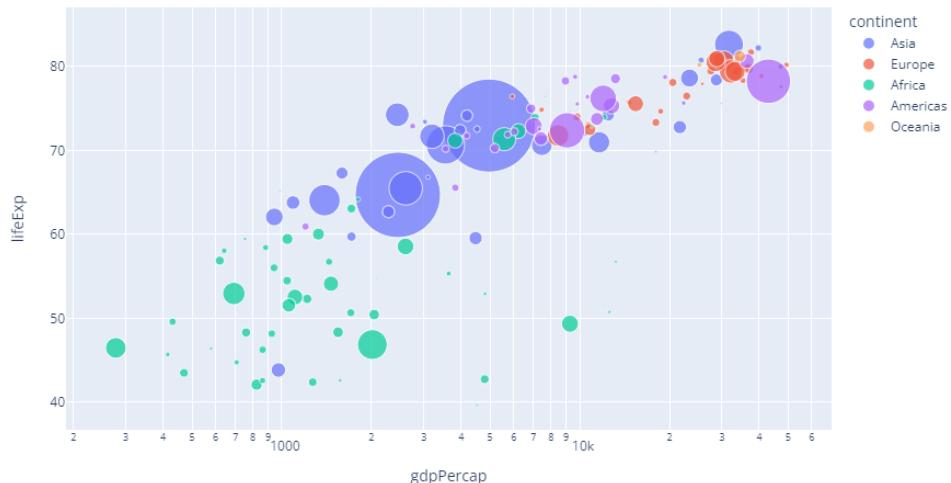


```
import plotly.express as px
```

```
df = px.data.gapminder()
```

```
fig = px.scatter(df.query("year==2007"), x="gdpPerCap", y="lifeExp",
                  size="pop", color="continent",
                  hover_name="country", log_x=True, size_max=60)
```

```
fig.show()
```

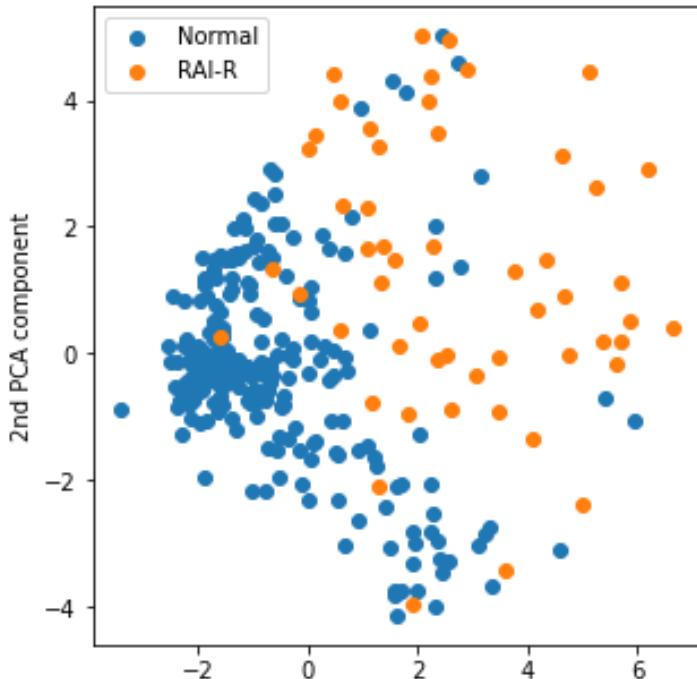


Machine learning

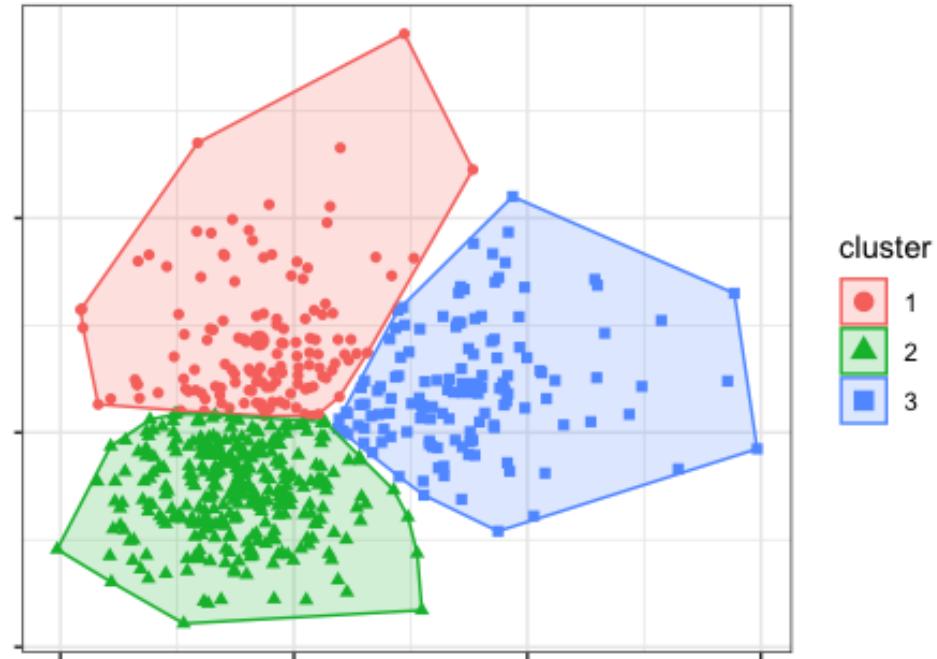
- Unsupervised learning
 - PCA, PCoA, t-SNE, UMAP
 - Clustering techniques
- Supervised learning
 - Predict cancer subtype
 - Identify potential biomarker genes
- (a touch of) Deep learning
 - AlphaFold

Two primary branches of unsupervised learning

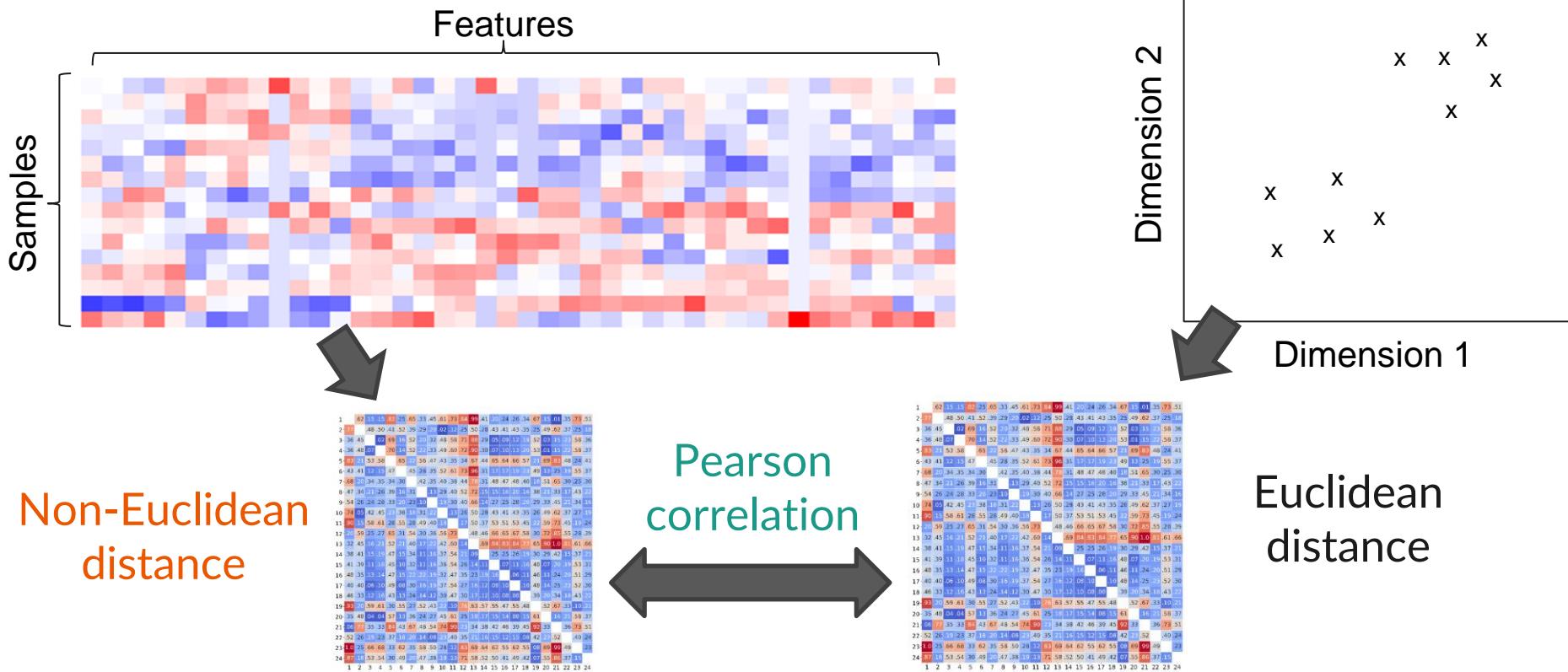
Dimensionality Reduction



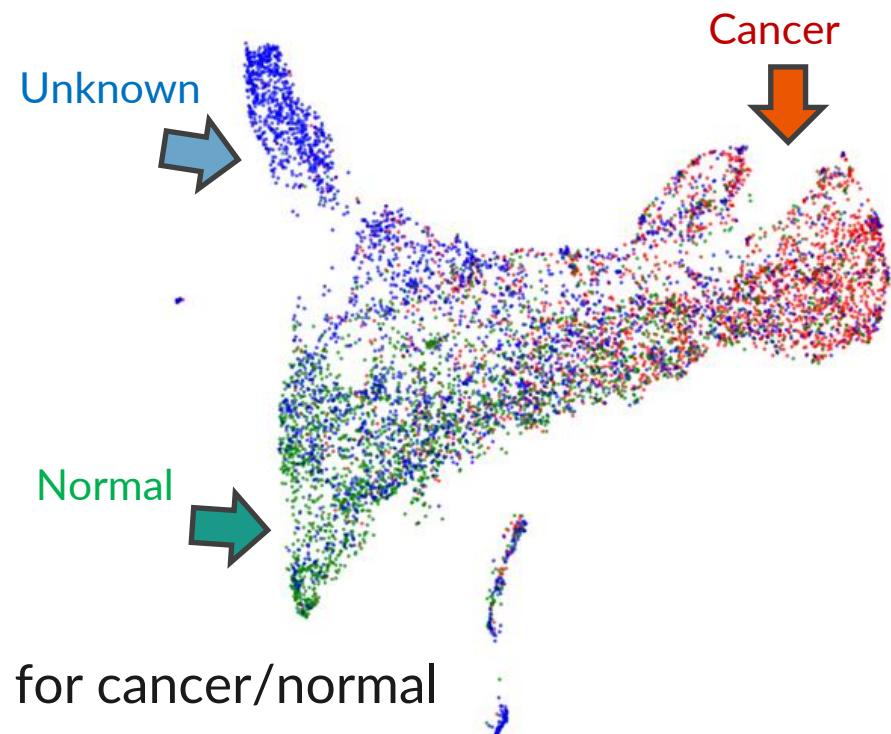
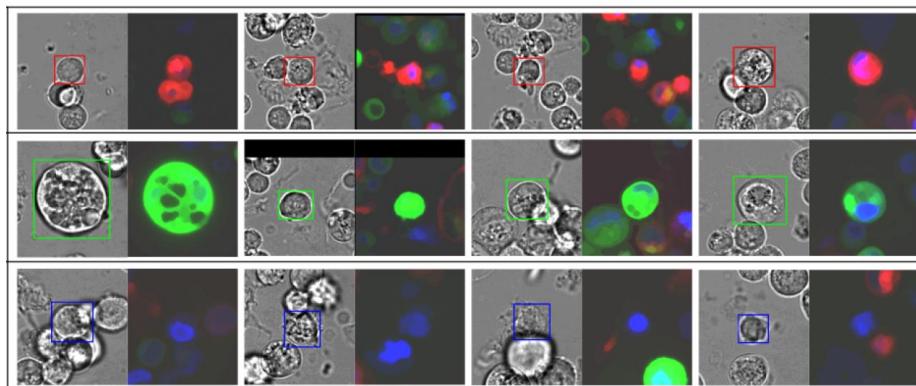
Clustering



Principal Coordinate Analysis (PCoA)

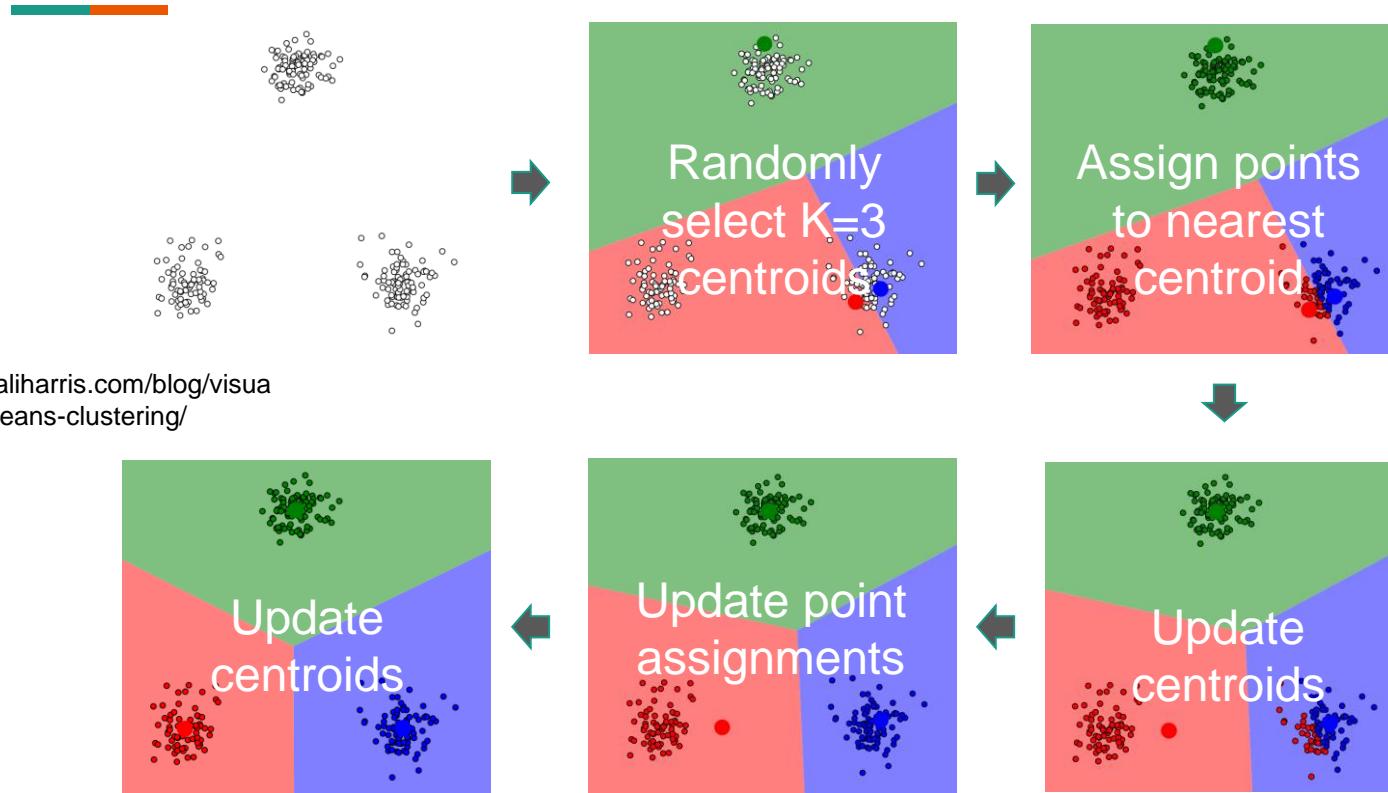


2D visualization for cell images

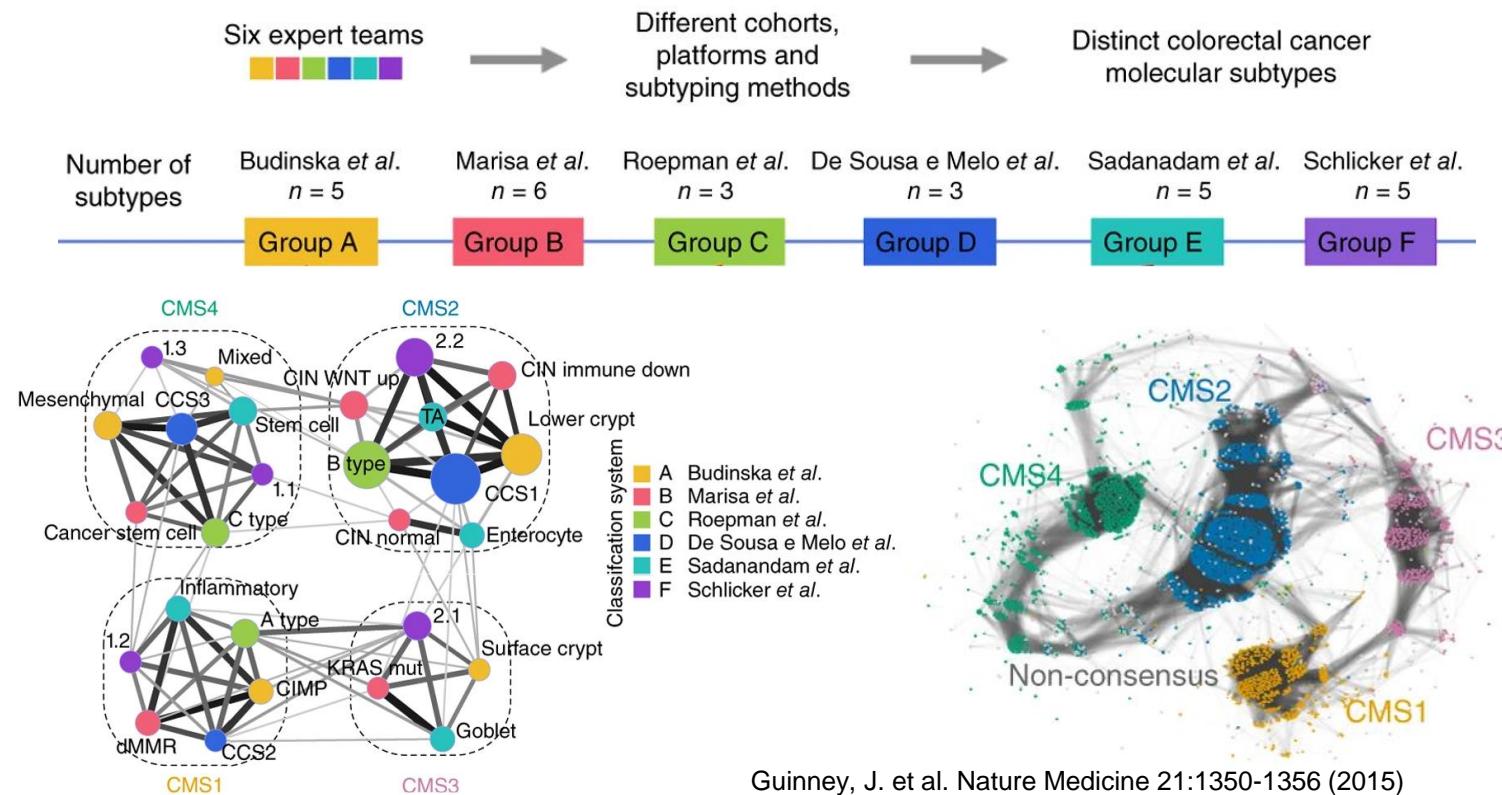


- Extracted from **deep learning model** for cancer/normal cell type classification

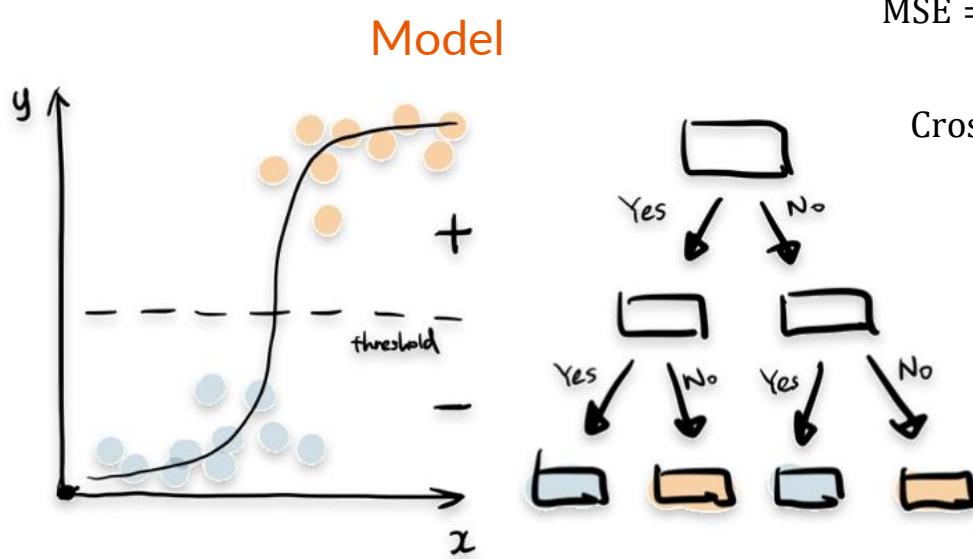
k-means algorithm



Consensus colorectal cancer subtyping



The cores of supervised learning

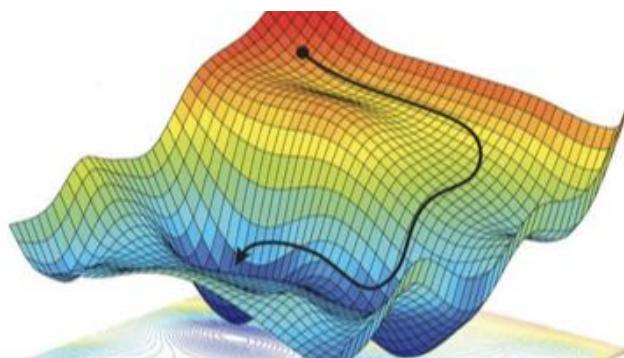


Objective / Loss Function

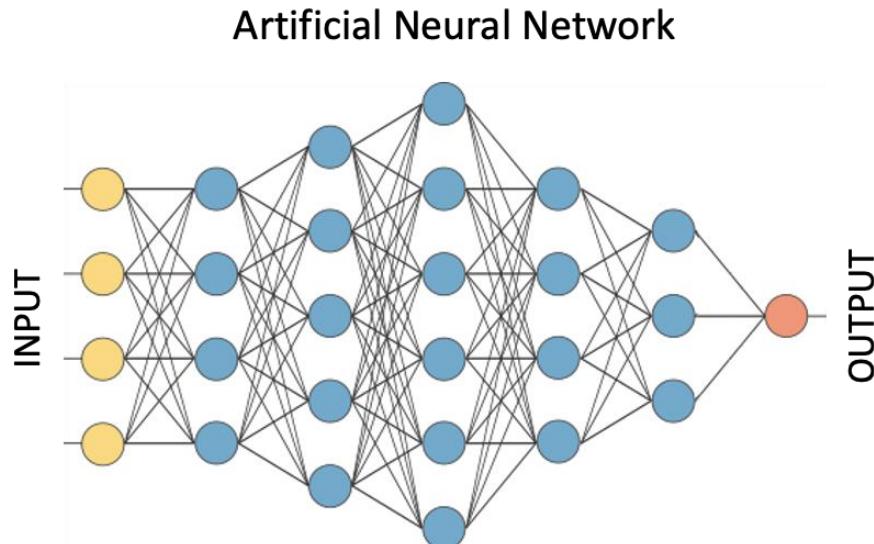
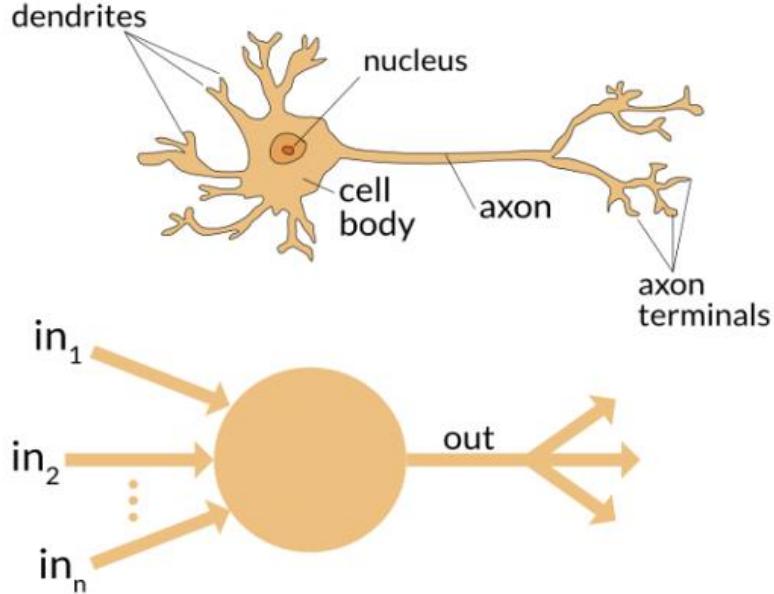
$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad \text{MAPE} = \frac{1}{n} \sum_{i=1}^n \frac{|y_i - \hat{y}_i|}{y_i} \times 100$$

$$\text{Crossentropy} = -\frac{1}{n} \sum_{i=1}^n y_i \ln(\hat{y}_i) + (1 - y_i) \ln(1 - \hat{y}_i)$$

Optimization Algorithm

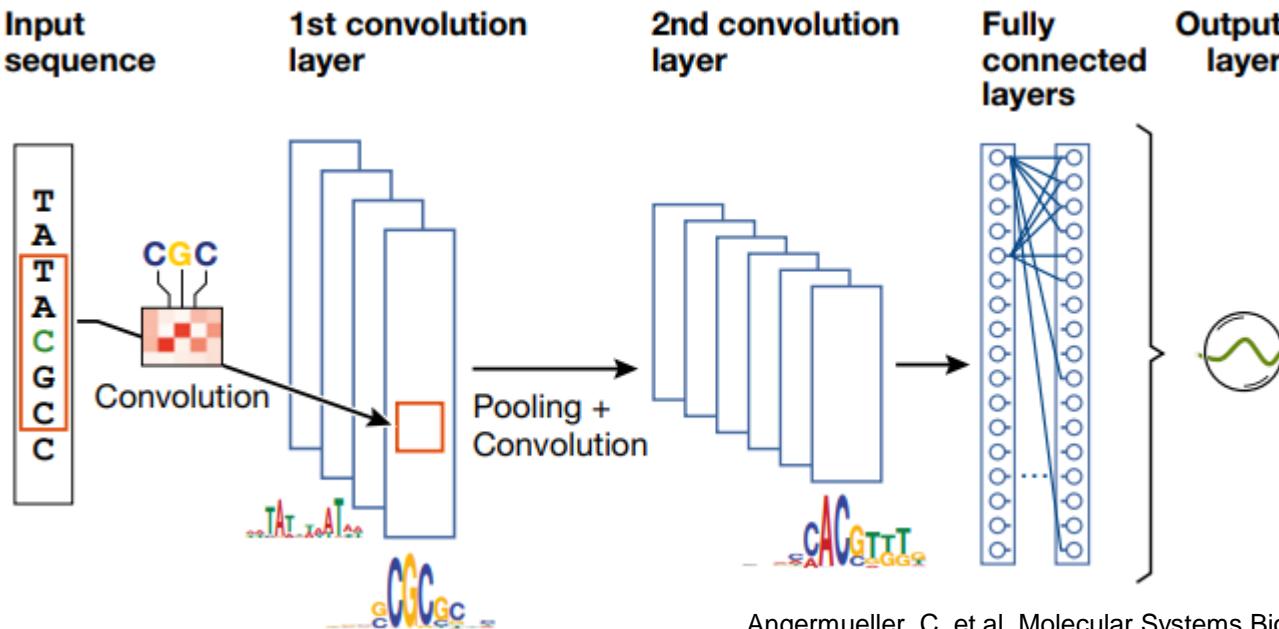


Artificial neural network



- Network of **simple** computation nodes: $out = f(w_1in_1 + w_2in_2 + \dots + w_nin_n)$

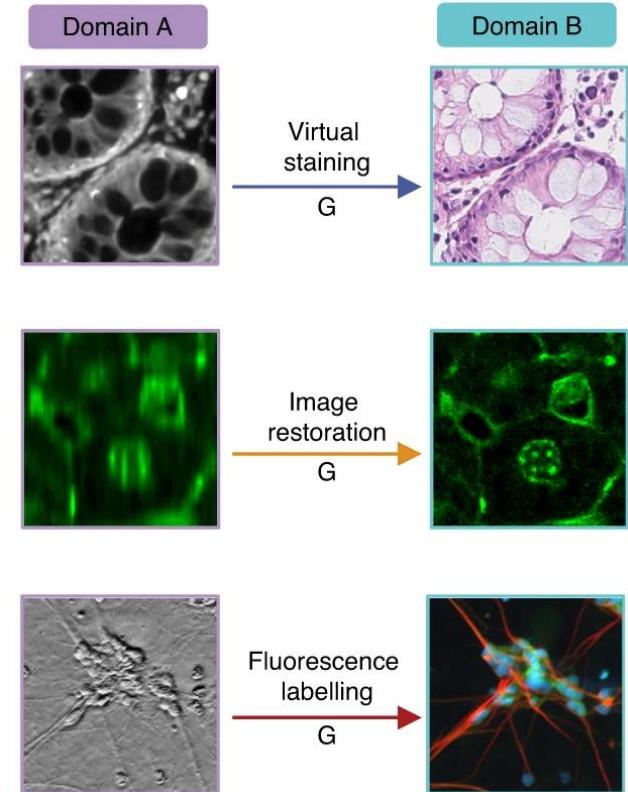
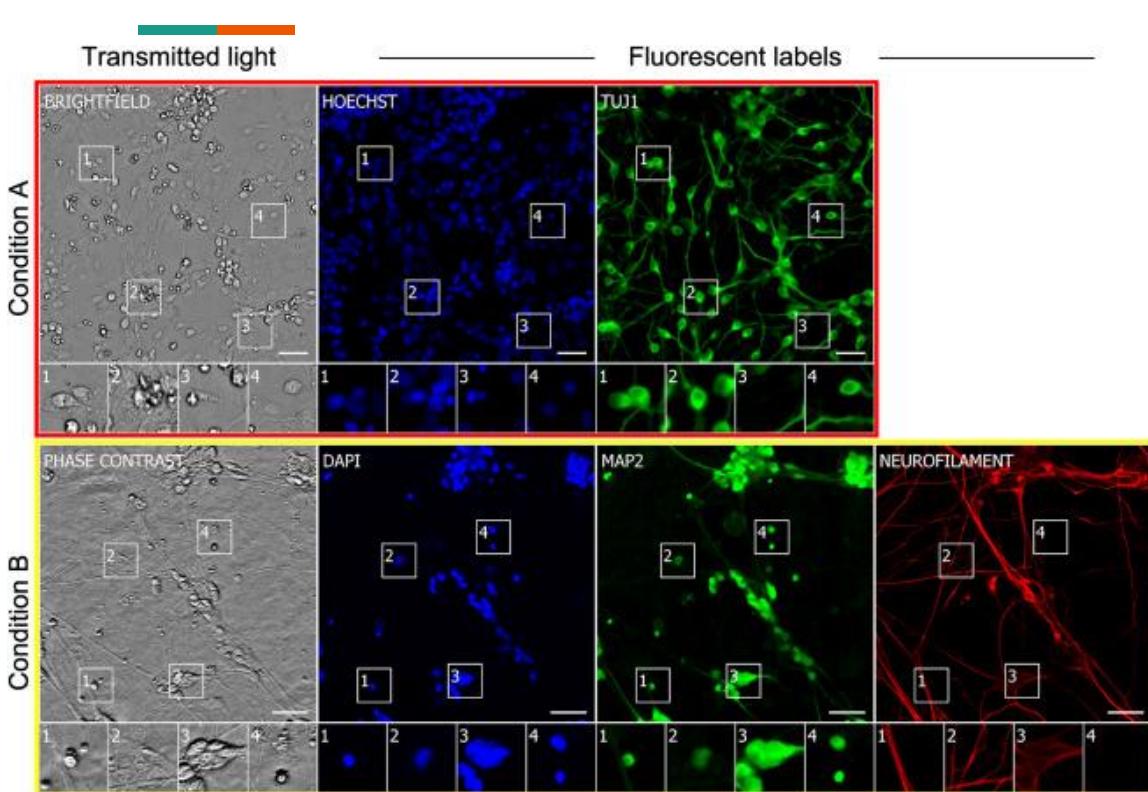
Convolution for DNA sequences



Angermueller, C. et al. Molecular Systems Biology 12:878 (2016)

- Motif = contextual pattern on DNA sequence

Virtual staining



From bioinformatics to deep learning

Published: 24 September 2018

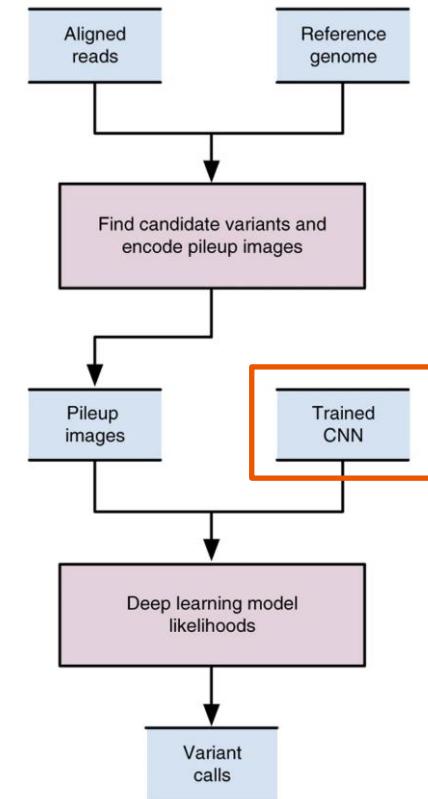
A universal SNP and small-indel variant caller using deep neural networks

Published: 27 July 2015

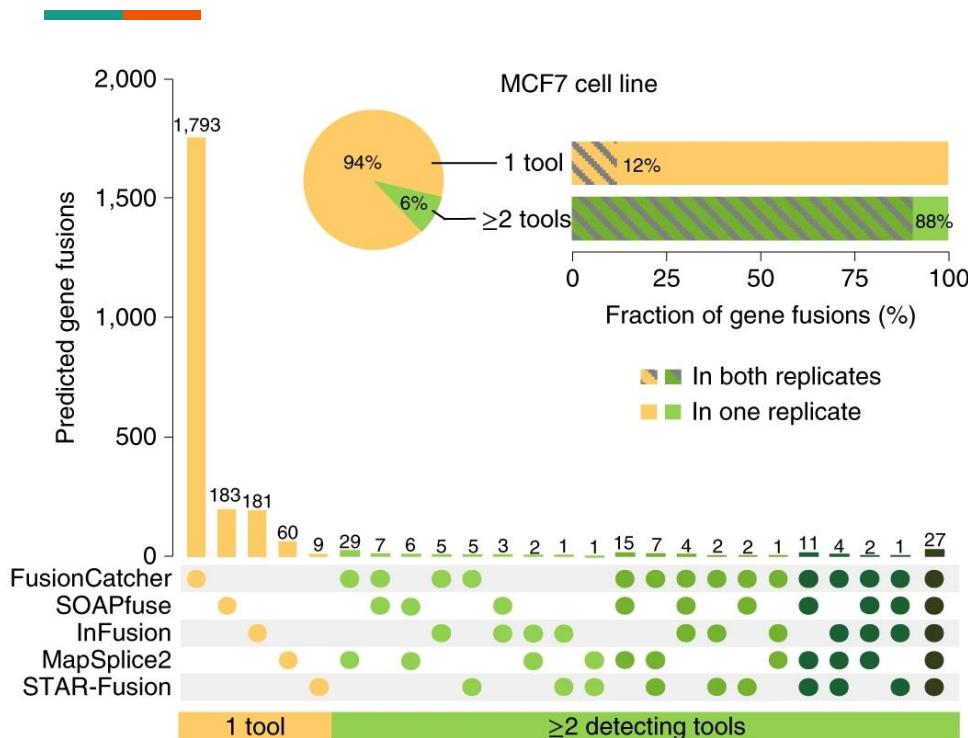
Predicting the sequence specificities of DNA- and RNA-binding proteins by deep learning

Article | Open Access | Published: 19 May 2022

Prediction of protein–protein interaction using graph neural networks



Aggregate scores from multiple tools



Any questions?
