

# **Applications of Pattern Recognition in Computational Biology**

Pattern Recognition Course (2110597)  
Chulalongkorn University

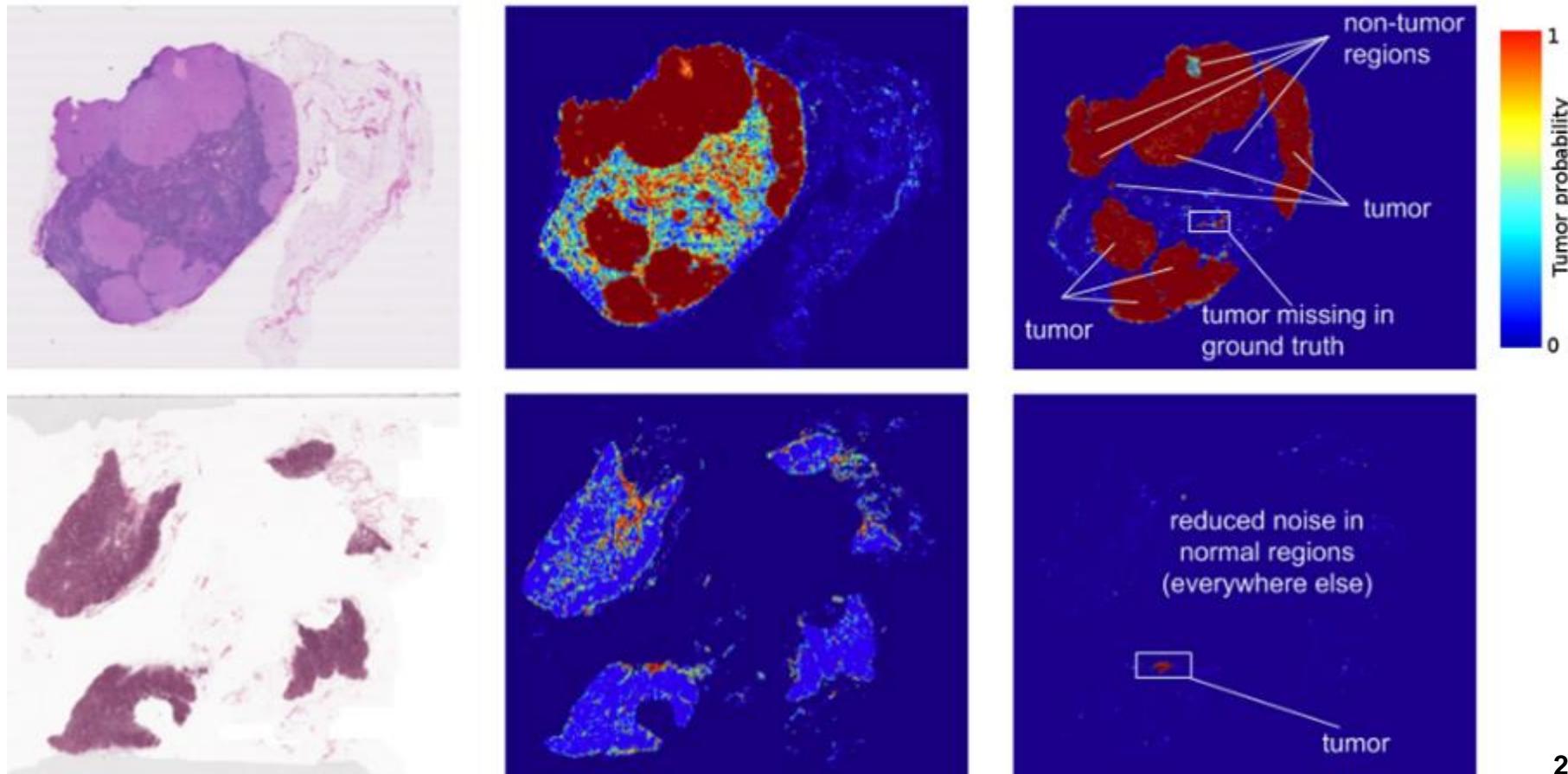
August 22<sup>nd</sup>, 2017

Instructor: Sira Sriswasdi (ສිරະ ສ්‍රීສ්වසදි)

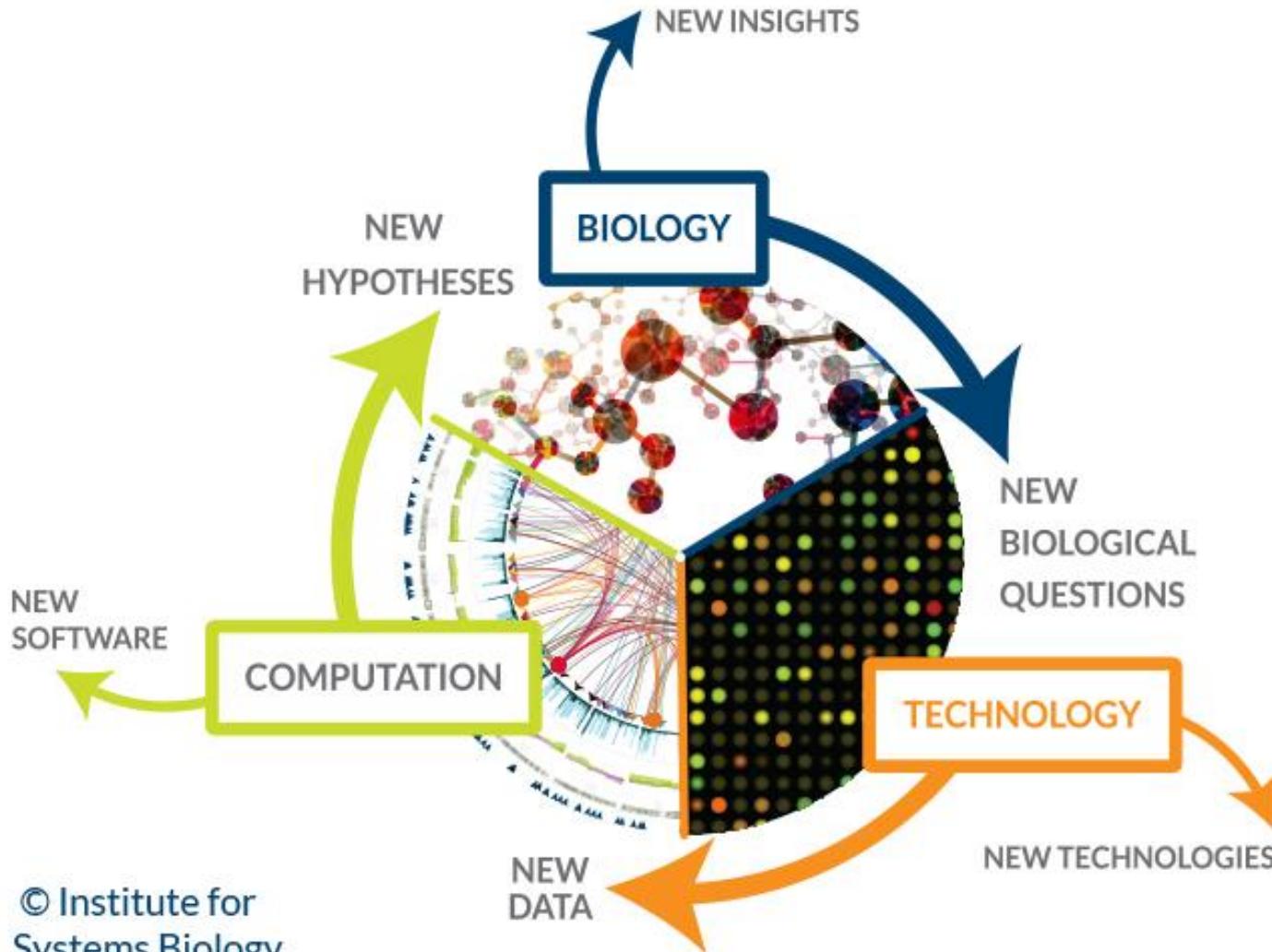
# Google Deep Learning system diagnoses cancer better than a pathologist with unlimited time

**89% accuracy vs 73%**

Ben Lovejoy - Mar. 3rd 2017 6:51 am PT  @benlovejoy



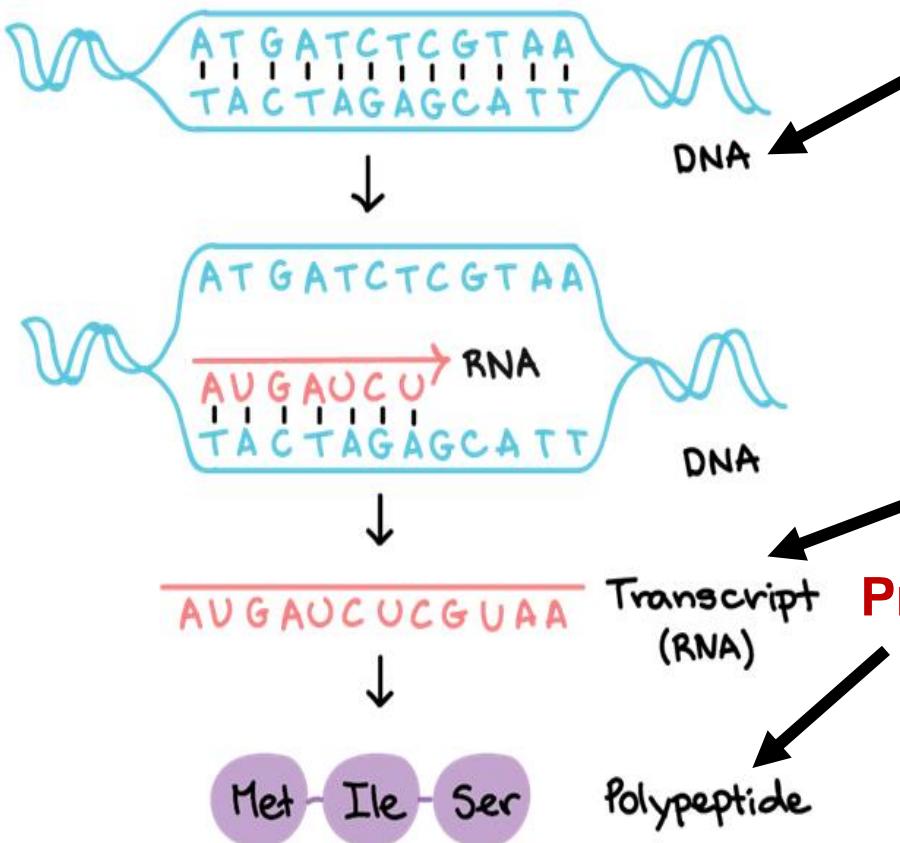
# Biology + Computation



# Data From High-Throughput Technology

## The Central Dogma

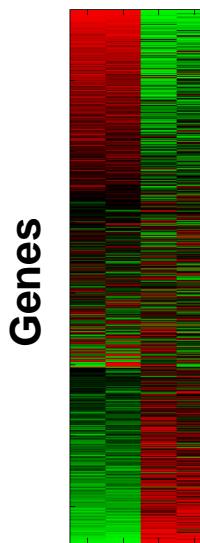
“Information Processing in Biology”



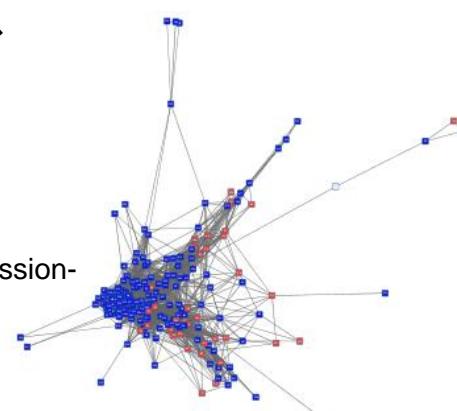
## Genome Sequence

ACCAGCGCGAAGCTCGGGCGGAGGGTTGA  
GCCACATGAGGCATGGCGACAATGAGGCGAG  
ACATGGCG **TGGCTGGC** TGTTACATTTGTTT  
GATGAAAAGCATAACCATGCGGATGATATT  
TATTATAGACTAGAGATGATTATTGAATAGAC  
**ATGCTCTAACCATTTAACCTCTAATTCCAC**

## RNA Expression

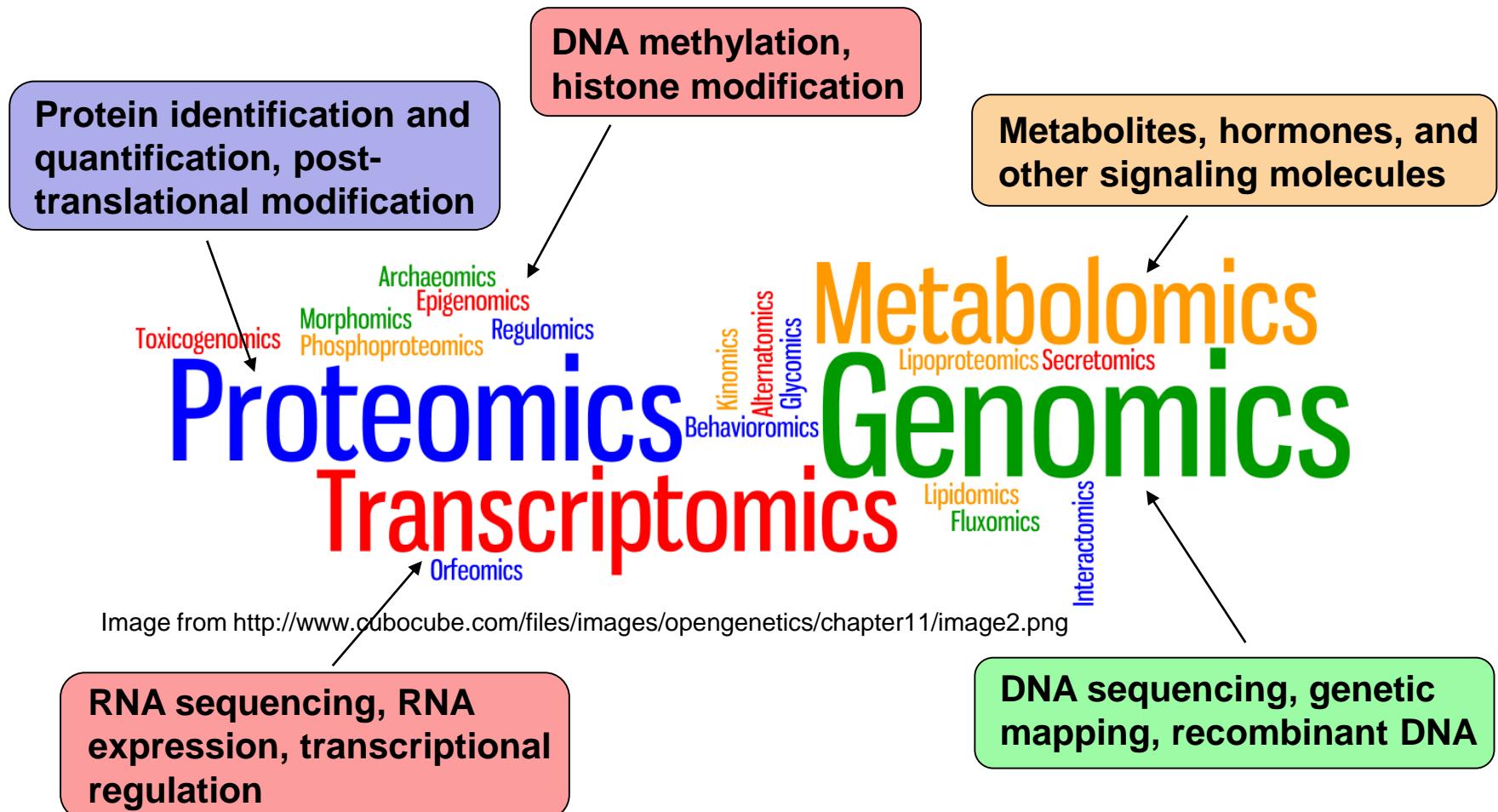


## Protein Interactions



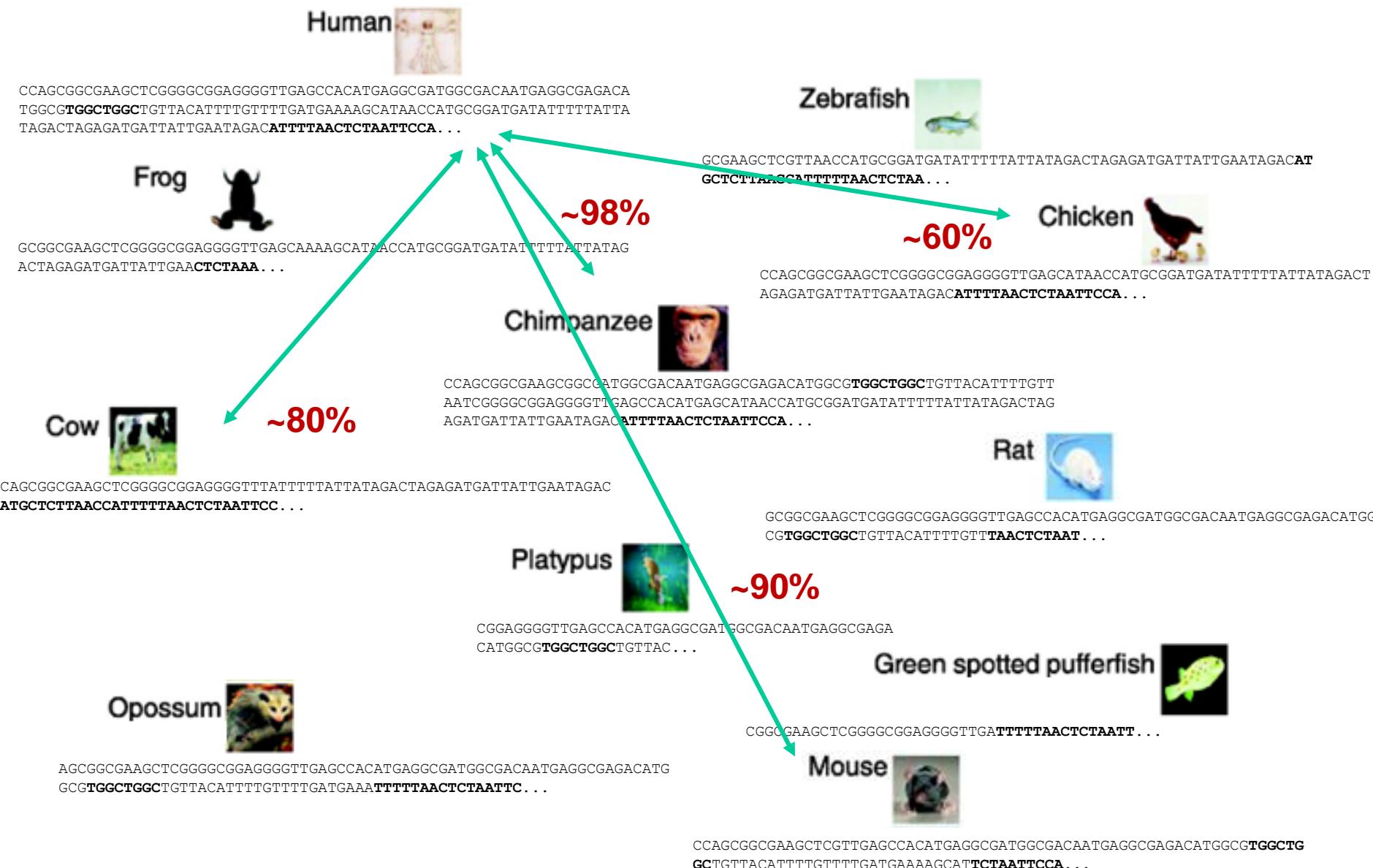
Conditions

# The Omics Era

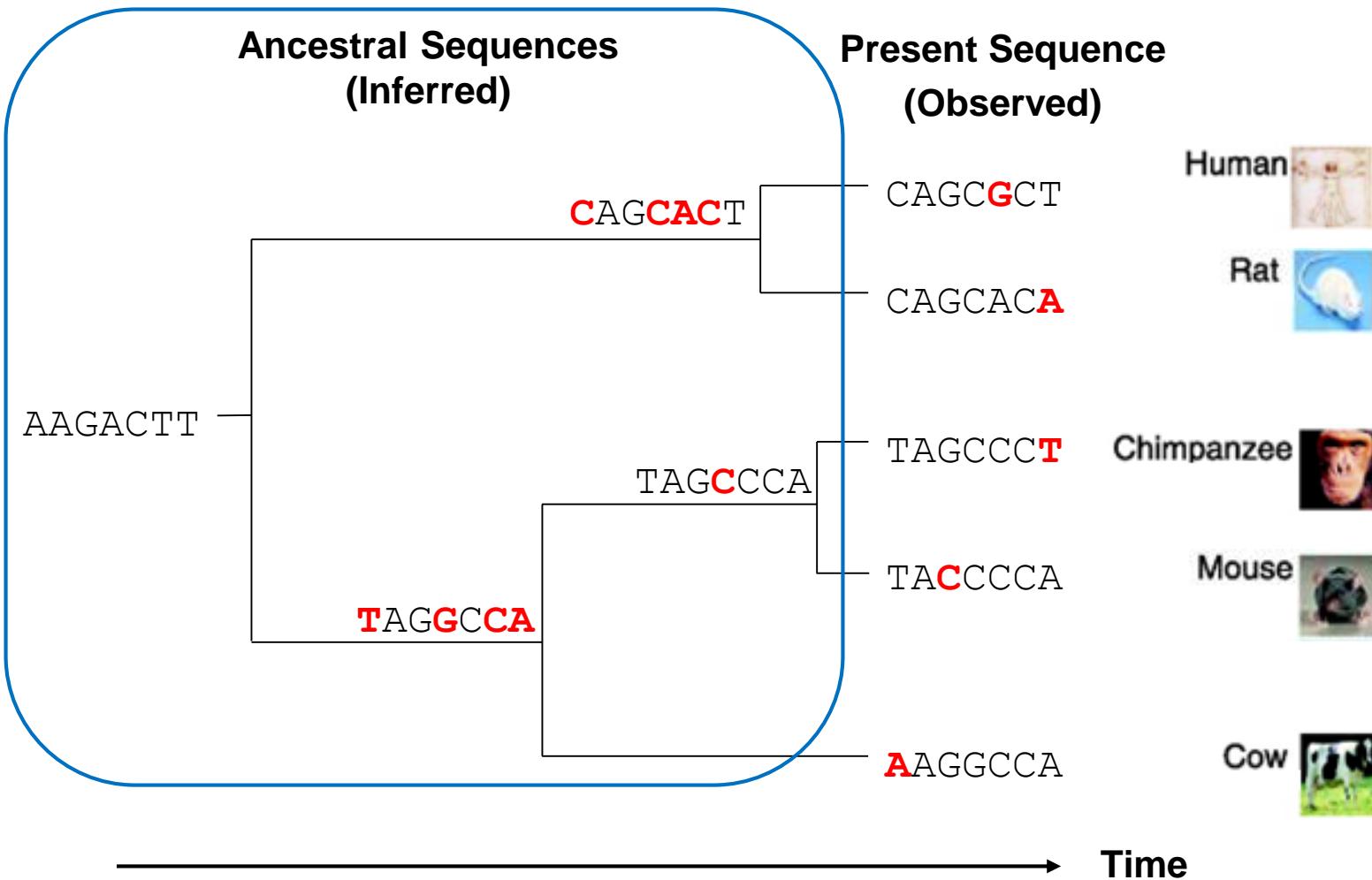


# **Application I: Evolutionary Genomics**

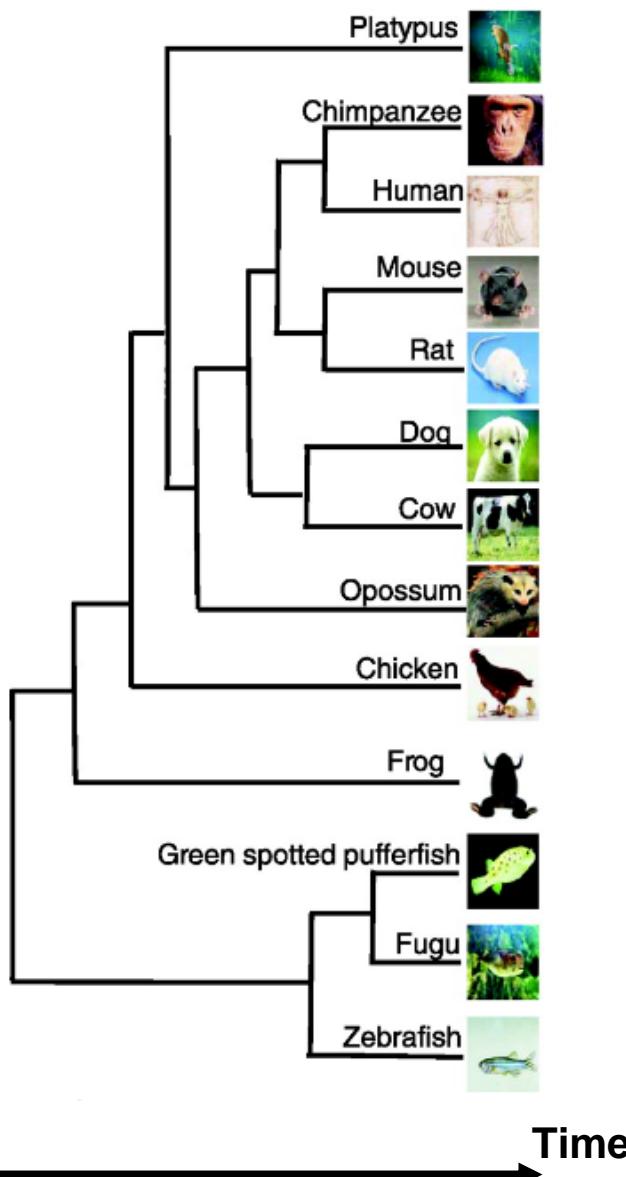
# Genome Sequence As Species Signature



# Evolution Of DNA Sequences

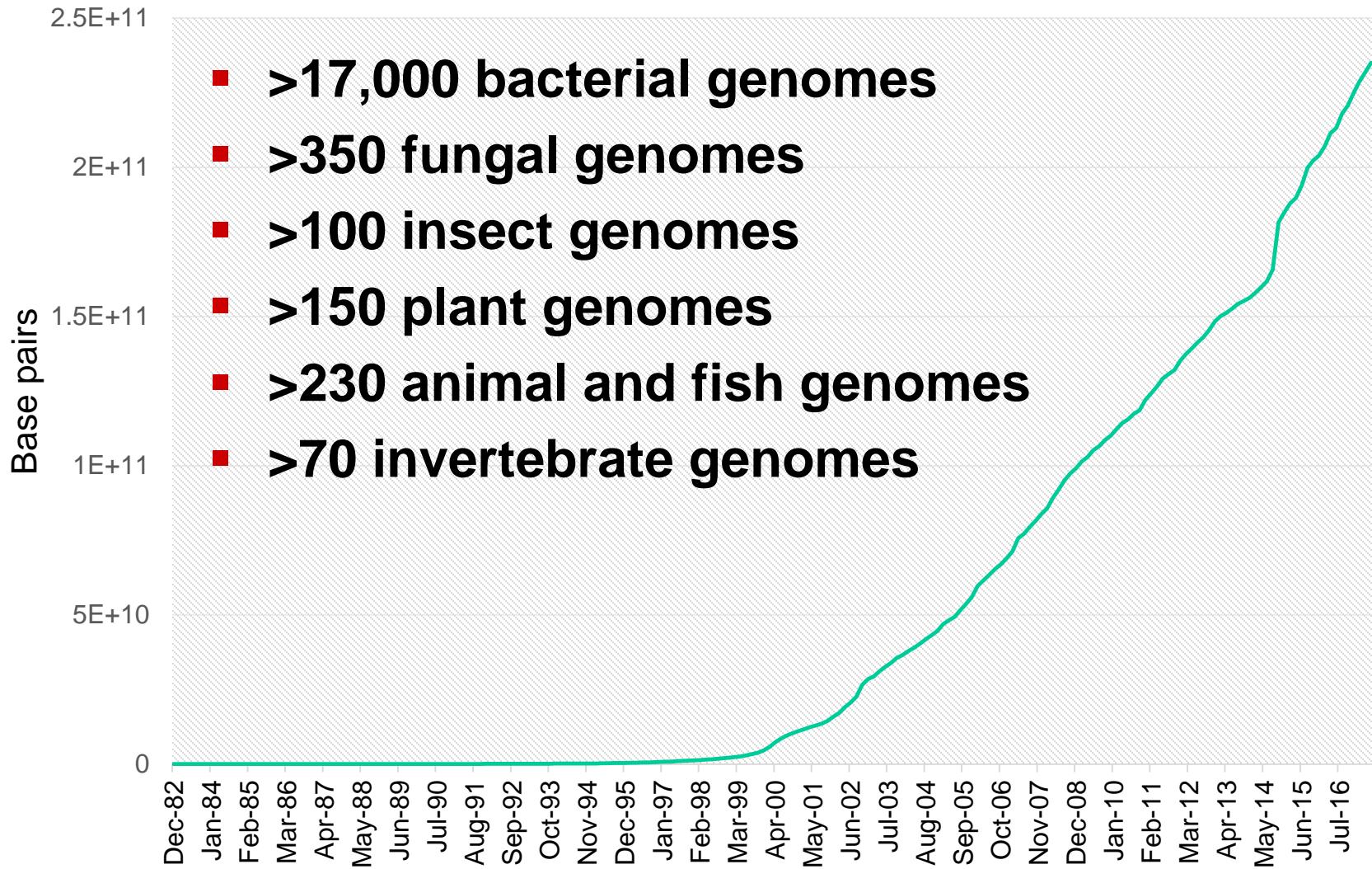


# Inferring Evolutionary History (Phylogenetics)



- Reconstruction of evolutionary events over millions of years
- Based on genome sequences of currently existing species
- Assume some models of evolution on DNA sequence, e.g.  $P(A \rightarrow T)$ ,  $P(G \rightarrow T)$
- Output the most likely tree topology and branch lengths
- Extremely large number of parameters, search spaces, and number of models to compare

# Growing Amount Of Genomic Data



# Forecasting And Regulating Evolution

- **Epidemiology**
  - Tracking the spread of disease outbreaks
  - Predict the next outbreaks and prepare vaccines in advance
- **Biotechnology**
  - Genetic engineering and breeding of new strains with desired characteristics and capabilities
- **Wildlife Conservation**
  - Pairing evolutionary history with climate/environmental changes can reveal the factors that drive animal evolution and extinction

# **Application II: Population Genetics**

# Tracing Population Structure Over Time

## Migration

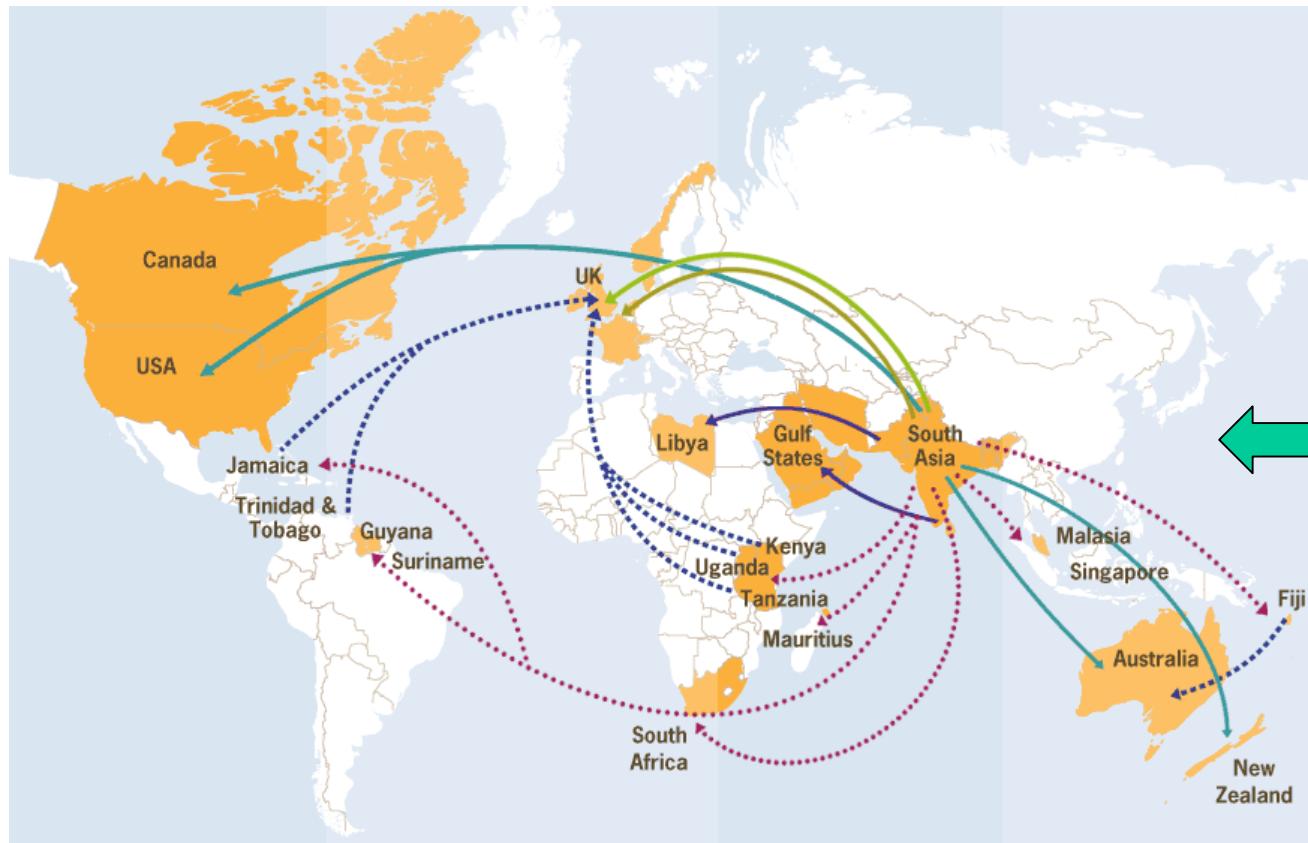


Image adapted from <https://www.theodysseyonline.com/why-people-migrate>

## Genetic Inheritance

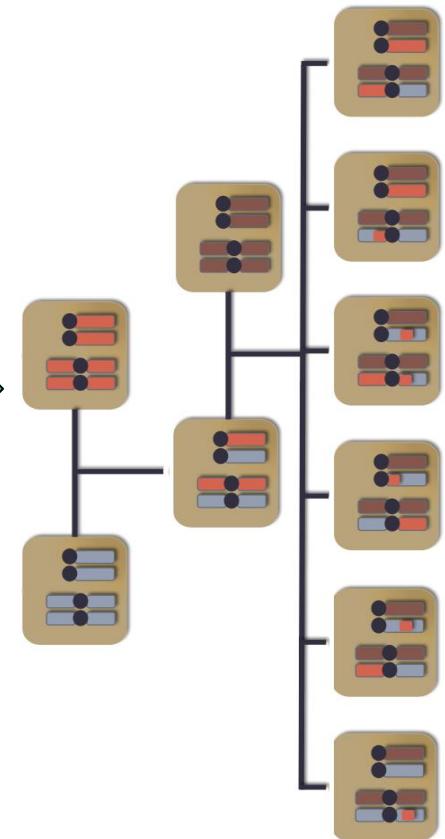


Image adapted from  
<https://wiki.uiowa.edu/display/2360159/Autosomal+Inheritance>

# Single Nucleotide Polymorphisms (SNPs) As Individual's Genetic Signature

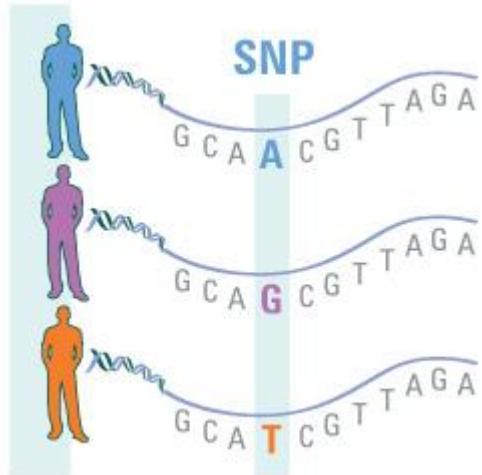
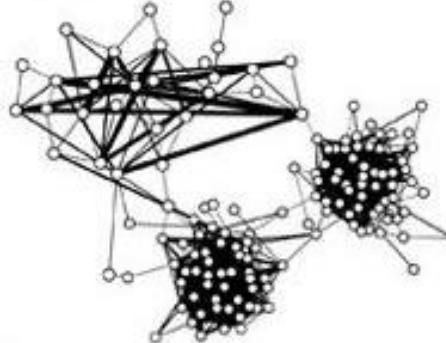


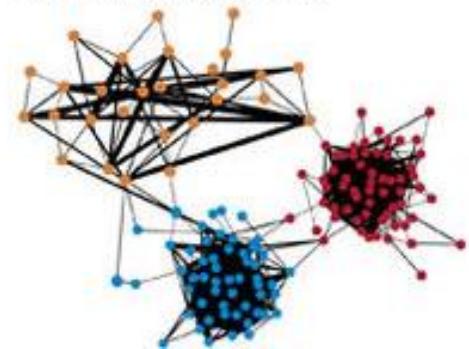
Image from <http://uvmgg.wikia.com/wiki/SNP>

## Identity-By-Descent (IBD) Network

Construct network from IBD.  
Join vertex pairs (genotyped samples) if IBD>12 cM.  
Edge weights are a function of total detected IBD.



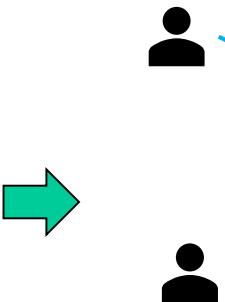
Detect network clusters.  
Recurisvely identify disjoint sets that maximize the modularity of the network. (Here one level of clustering hierarchy is shown.)



Han et al. Nature Communication 8, 14238 (2017)

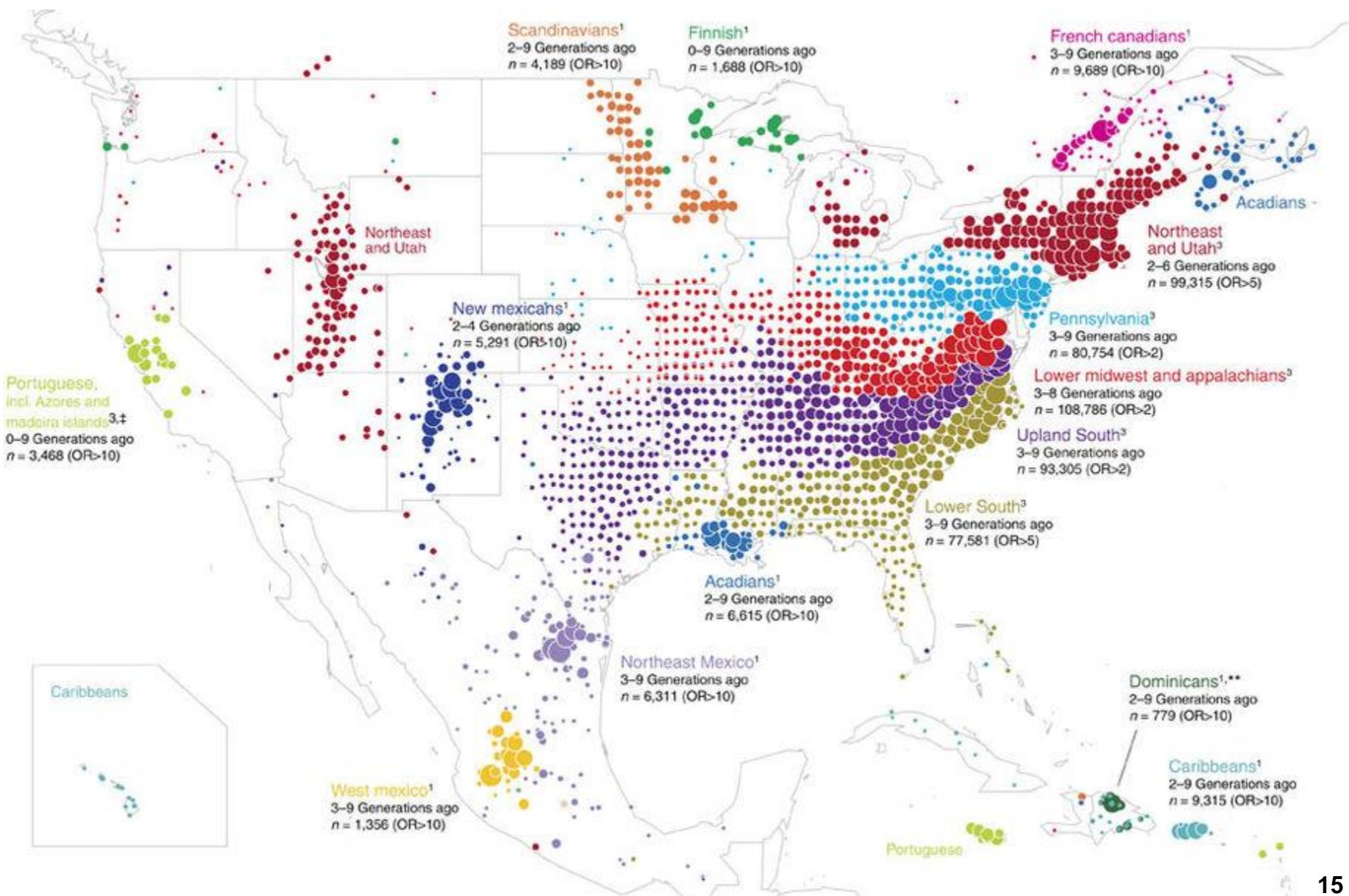
774,516 individuals

709,358 SNPs						
	A	T	T	C	...	G
	C	G	G	A	...	T
	C	G	G	A	...	G
	G	G	T	C	...	A



Connect individuals that share significant portion of consecutive SNPs

# Roots Of North American Population



# From Population To Personalized Medicine



## Carrier status

Find out if your children are at risk for inherited conditions, so you can plan for the health of your family.



## Health risks

Understand your genetic health risks. Change what you can, manage what you can't.



## Drug response

Arm your doctor with information on how you might respond to certain medications.



## Health tools

Document your family health history, track inherited conditions, and share the knowledge.



## Inherited traits

Explore your genetic traits for everything from lactose intolerance to male pattern baldness.



## Scientific advances

Keep receiving updates on your DNA as discoveries are made, so your knowledge grows as you do.

Image from <http://poshrx.com/23andme-is-back-on/>

## ▪ Social Sciences

- Tracking the dynamics of populations
- Understanding ethnic structures

## ▪ Medicine

- Identifying common genetic variations within a population that may be associated with drug targets
- Identifying disease risk factors

# **Application III: Gene Expression Analysis**

# Measuring Gene Expression



Image adapted from  
<https://www.khanacademy.org/science/biology/gene-expression-central-dogma/transcription-of-dna-into-rna/a/overview-of-transcription>

**Amount of RNA product of each gene**

gene symbol	target id	control 1	control 2	exp 1	exp 2	exp avera	log10 fold differen
RANBP2	NM_006267.4	6.48833	12.9884	0.176602	0.117152	0.143837678	1.80497163
CYTH1	XM_011525475.2	15.0945	11.3296	0.249429	0.173519	0.20804007	1.798369664
HNRNPK	NM_031263.3	117.923	184.554	2.67089	2.25173	2.452370922	1.779274833
PSMD2	NM_002808.4	47.5916	97.2485	0.855375	1.50337	1.1339952	1.778095392
NONO	NM_007363.4	91.2359	146.966	2.35517	1.58616	1.932789809	1.77750669
UGP2	NM_001001521.1	45.0717	70.7679	0.982185	0.906358	0.943510059	1.777123589
FDFT1	NM_004462.4	113.559	139.643	3.17652	1.55325	2.221245077	1.753523887
ATP6V1H	NM_015941.3	40.4483	34.6337	0.940035	0.463959	0.660407222	1.753387768
STK10	NM_005990.3	13.8212	25.0476	0.458813	0.246147	0.336058691	1.743240804
DDX3X	NM_001193416.2	59.55	105.791	1.66098	1.25129	1.441654488	1.740804056
KIF5B	NM_004521.2	10.4073	43.989	0.356018	0.437435	0.394632403	1.734148339
TRIB1	XR_428373.2	13.8993	10.355	0.211033	0.232964	0.221727517	1.733251922
SREBF2	NM_004599.3	51.8562	90.6209	1.50785	1.06592	1.267772642	1.732973154

## Sequencing Machine



Image from  
[https://support.illumina.com/sequencing/sequencing\\_instruments/hiseq-4000.html](https://support.illumina.com/sequencing/sequencing_instruments/hiseq-4000.html)



# RNA Sequencing (Counting) Biases

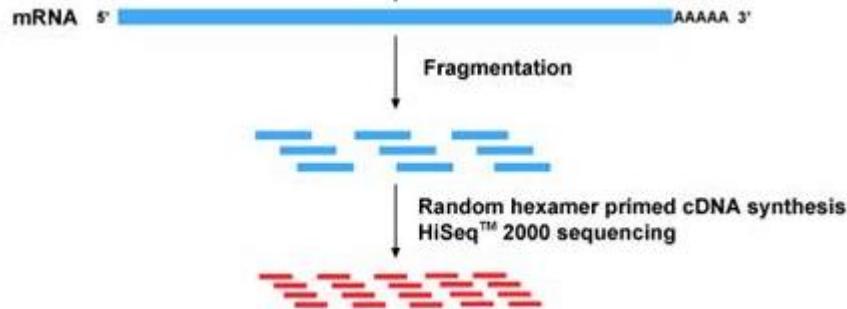


Image adapted from <http://bio.lundberg.gu.se/courses/vt13/rnaseq.html/>

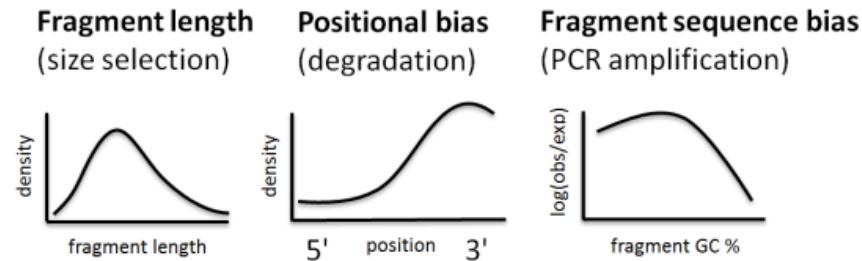


Image adapted from  
<https://mikelove.wordpress.com/2016/09/26/rna-seq-fragment-sequence-bias/>

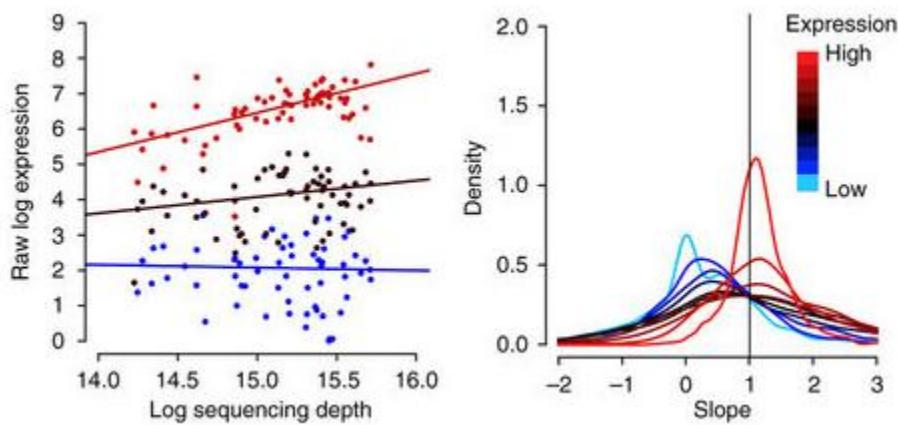
- Due to technological limitation, the entire length of RNA cannot be sequenced at once
  - Full-length RNA has to be fragmented
  - Bias in fragment length
- To increase sensitivity, fragmented RNA has to be amplified
  - Bias in signal amplification
- Sequencing is directional
  - Bias in head-to-tail read count



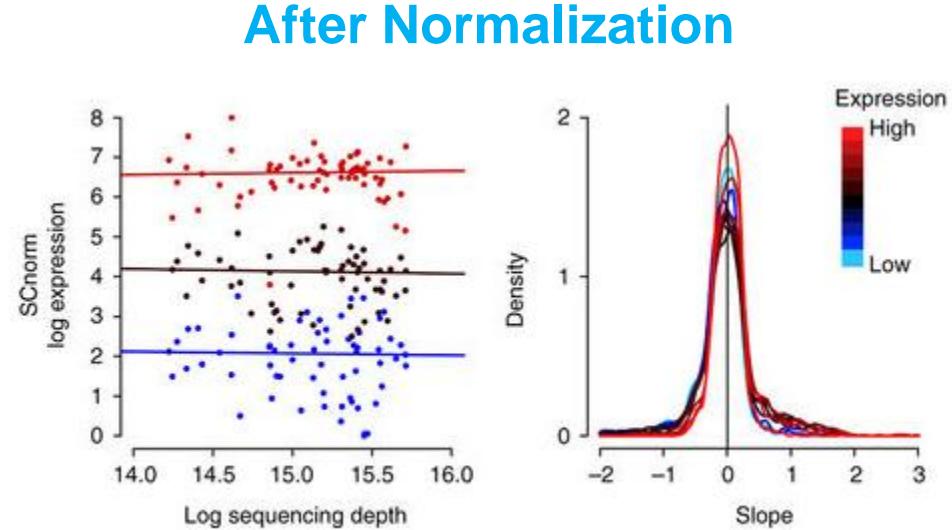
Bias correction

# Bias Normalization via Regression

Before Normalization



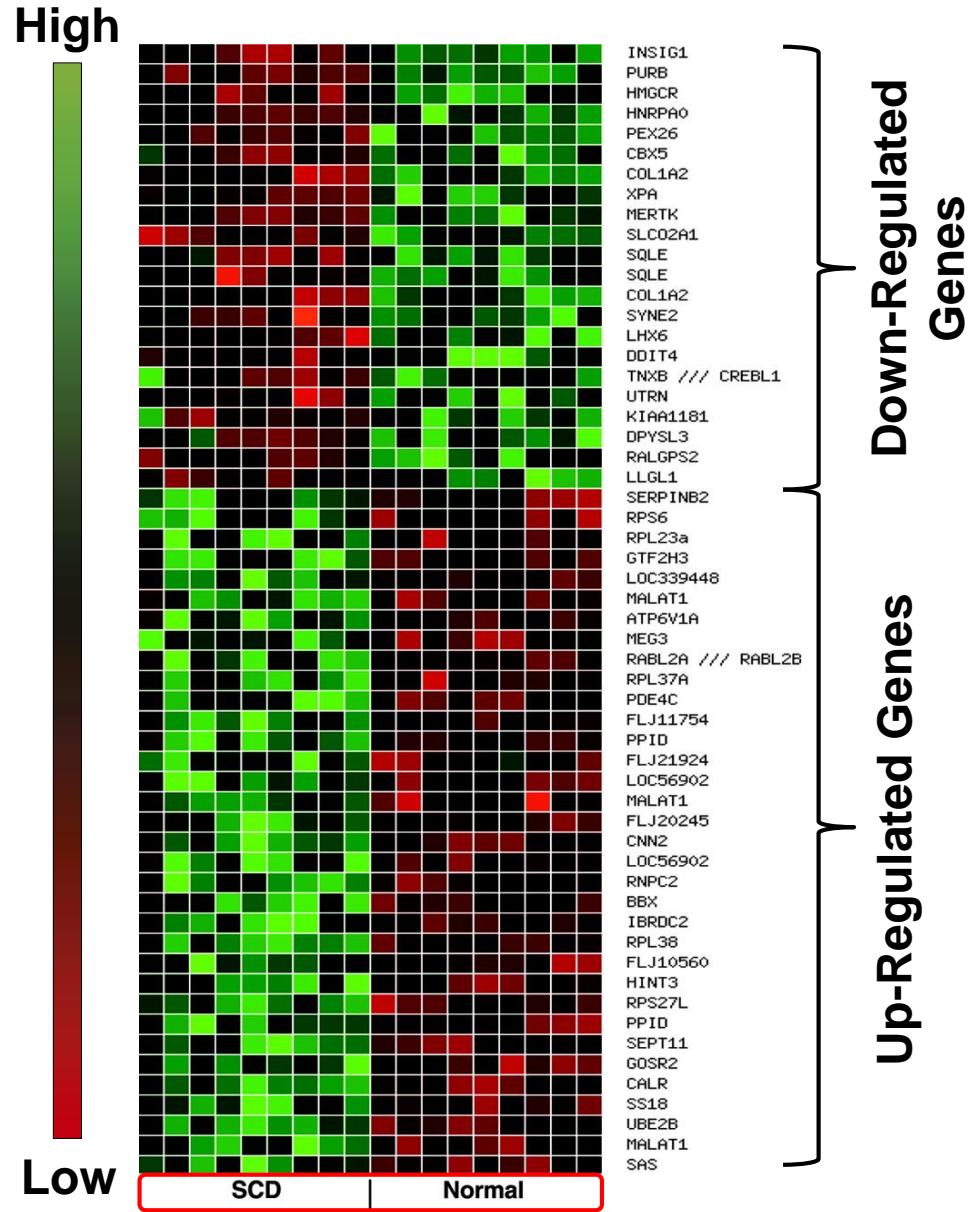
After Normalization



Adapted from Bacher *et al.* Nature Methods 14, 584-586 (2017)

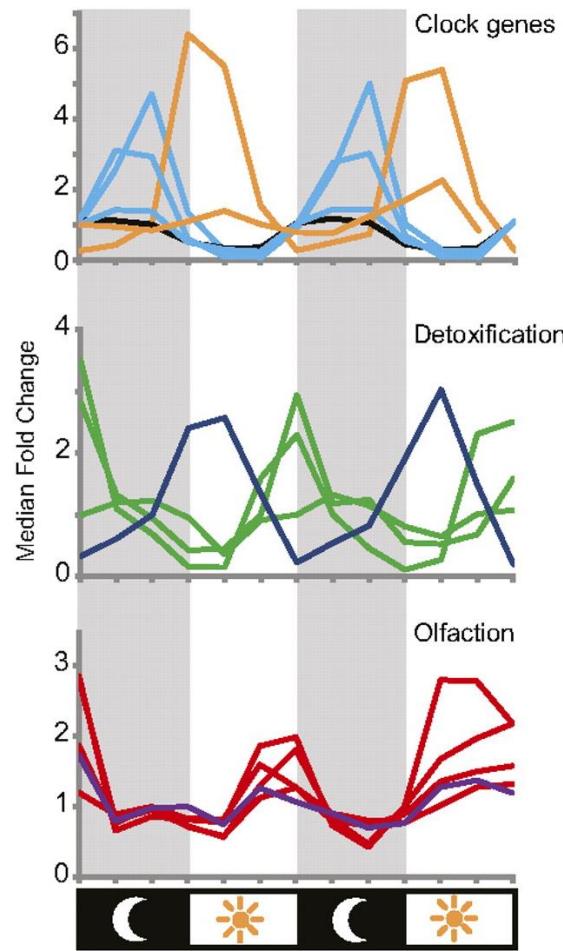
- Sequencing involves sampling of RNA transcripts
- Estimated expression levels of **low**, **medium**, and **high** expression genes are differently affected by the throughput of RNA sequencing experiment
- Normalization by regression corrected the biases

# What Can Gene Expression Tell Us?



Klings et al. Physiological Genomics 21, 293-298 (2005)

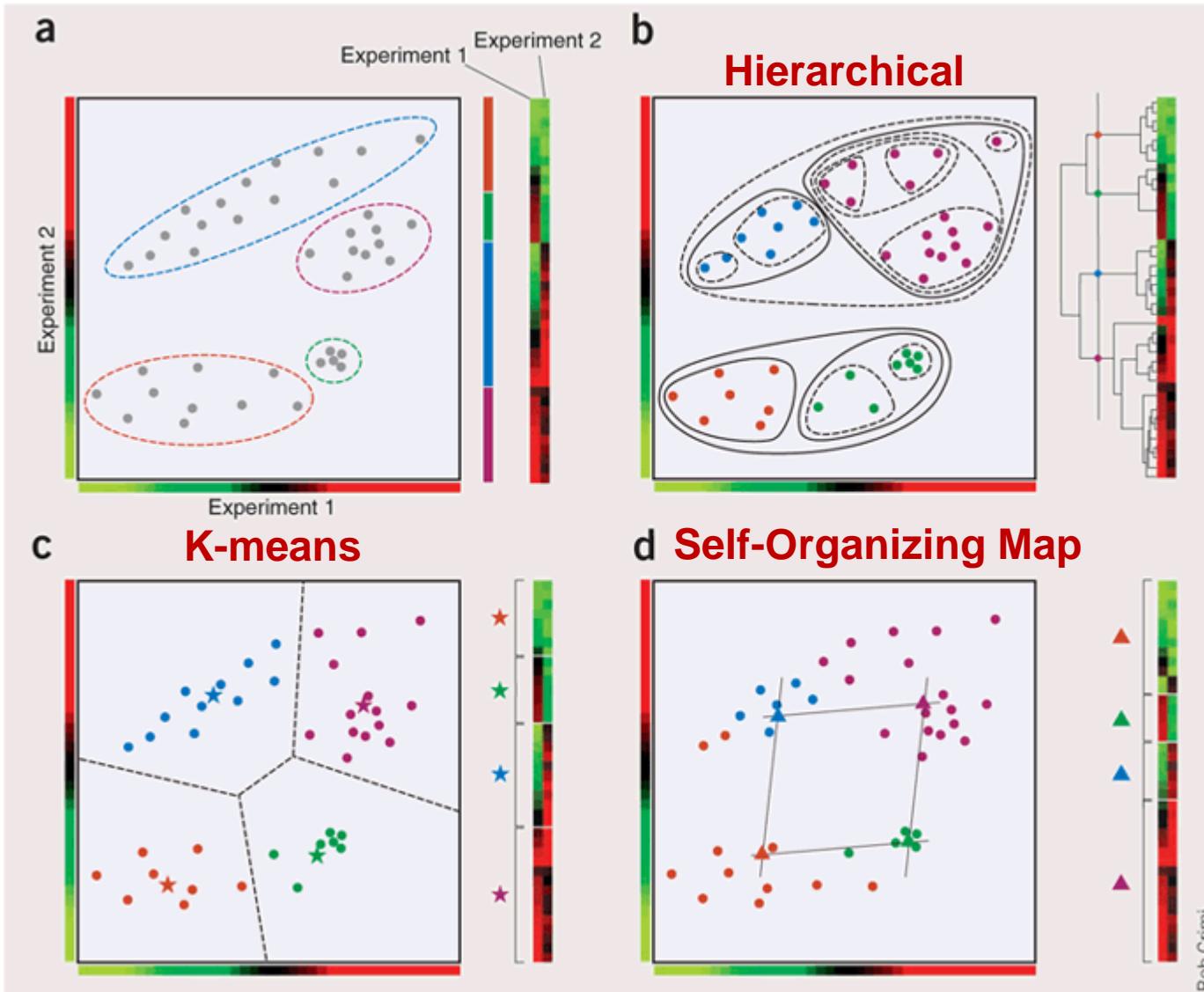
## Time Series Analysis



Adapted from Rund et al. PNAS 108, E421-430

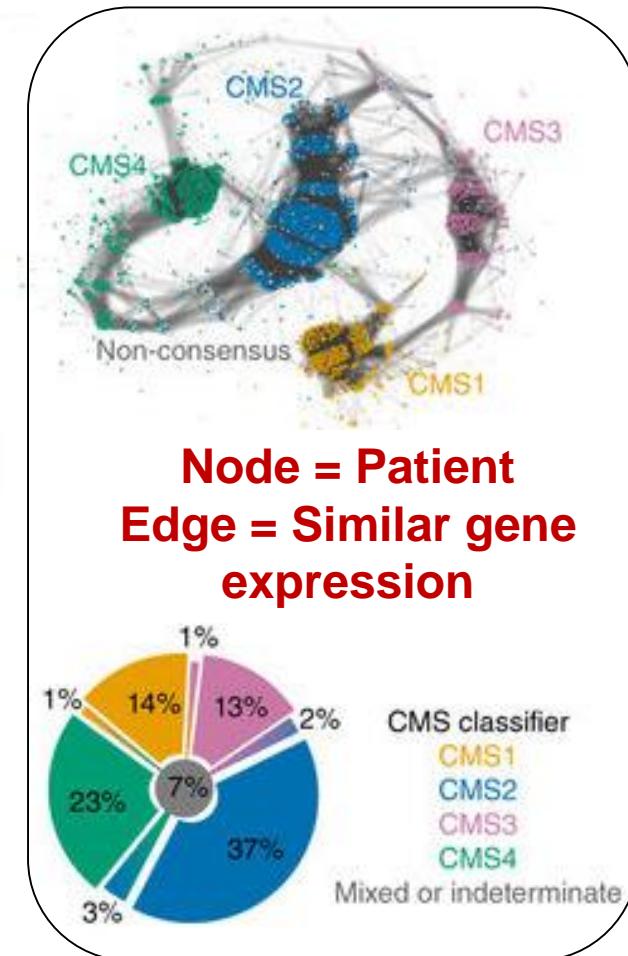
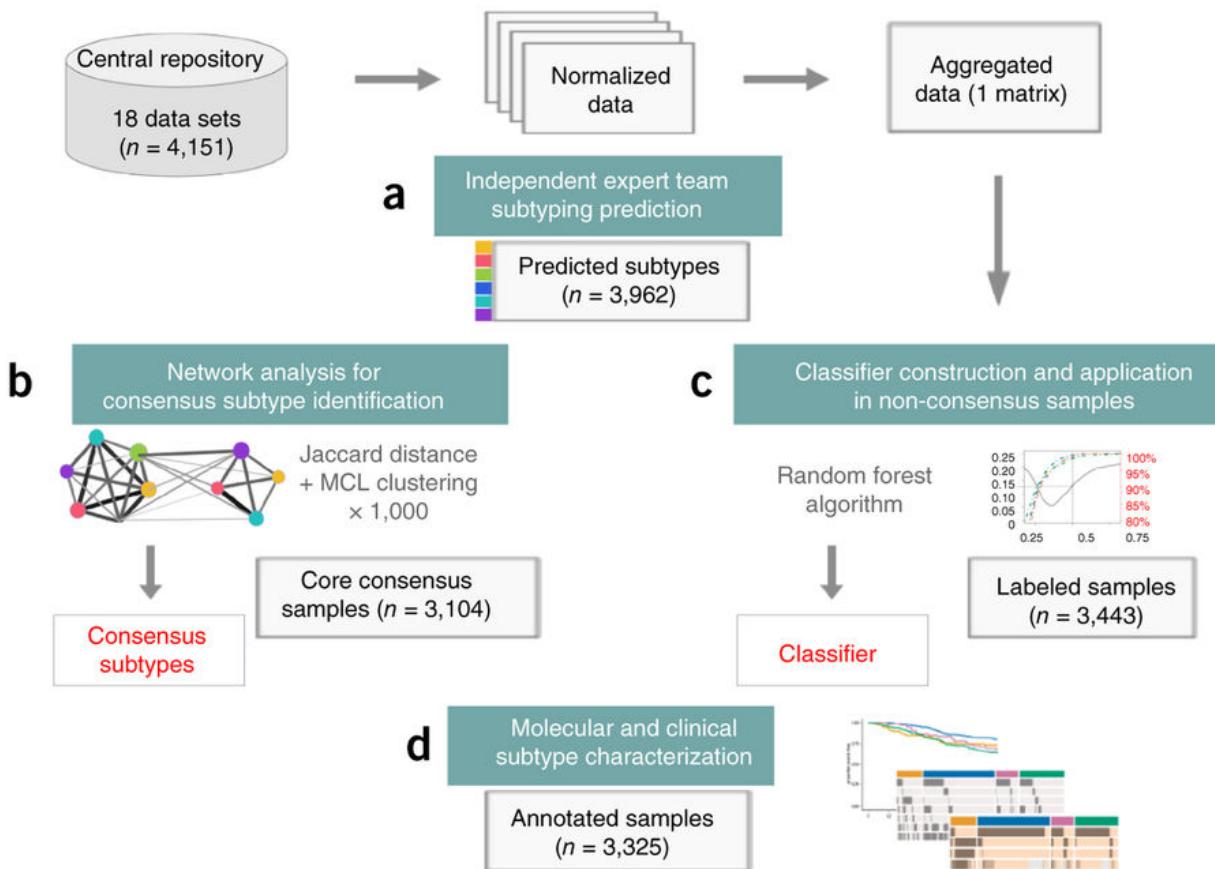
# Structure Behind Gene Expression Profiles

Each cluster represents a group of genes with similar functions



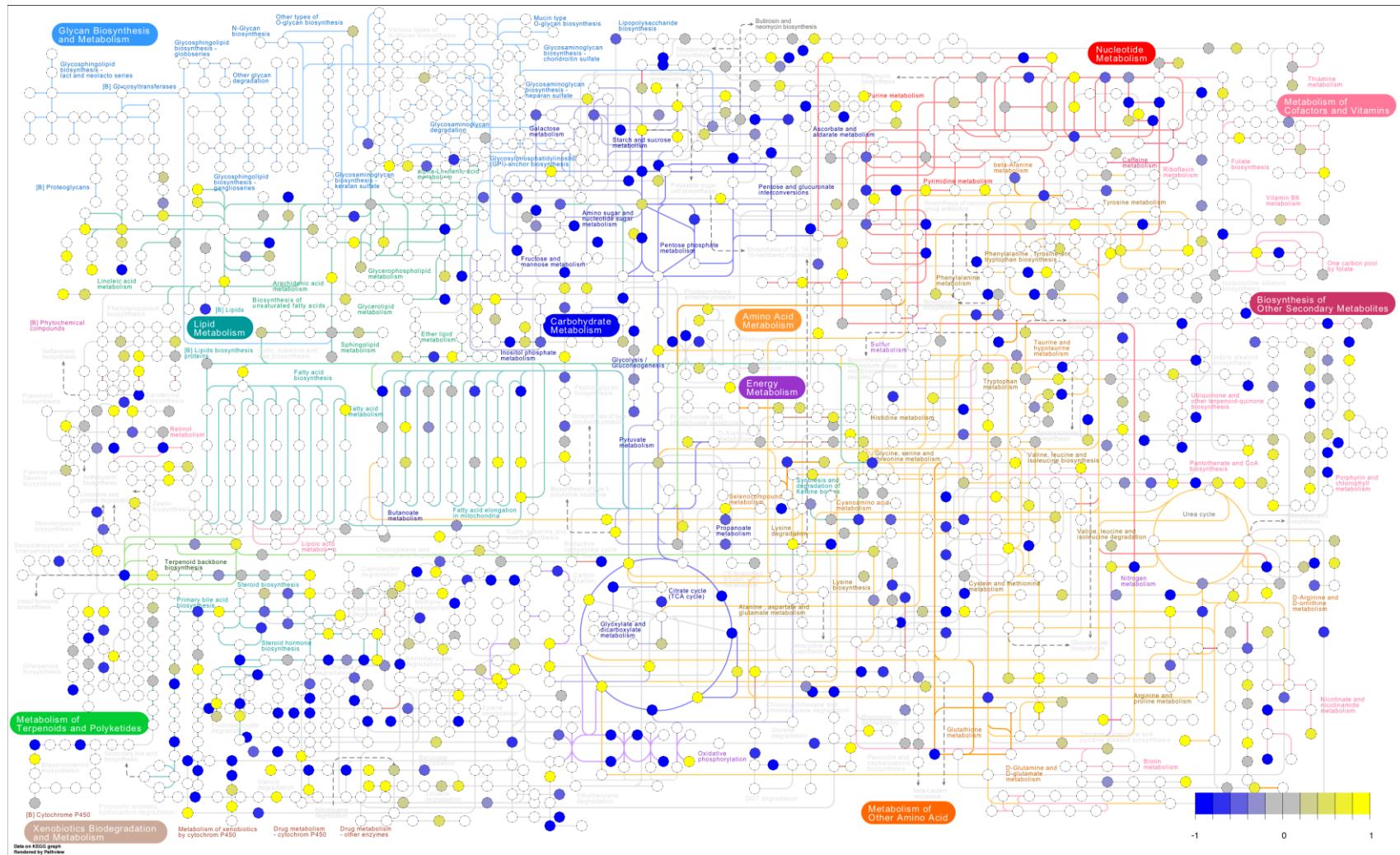
# Identifying Disease Subtypes

Gene expression data from >4,000 colorectal cancer patients



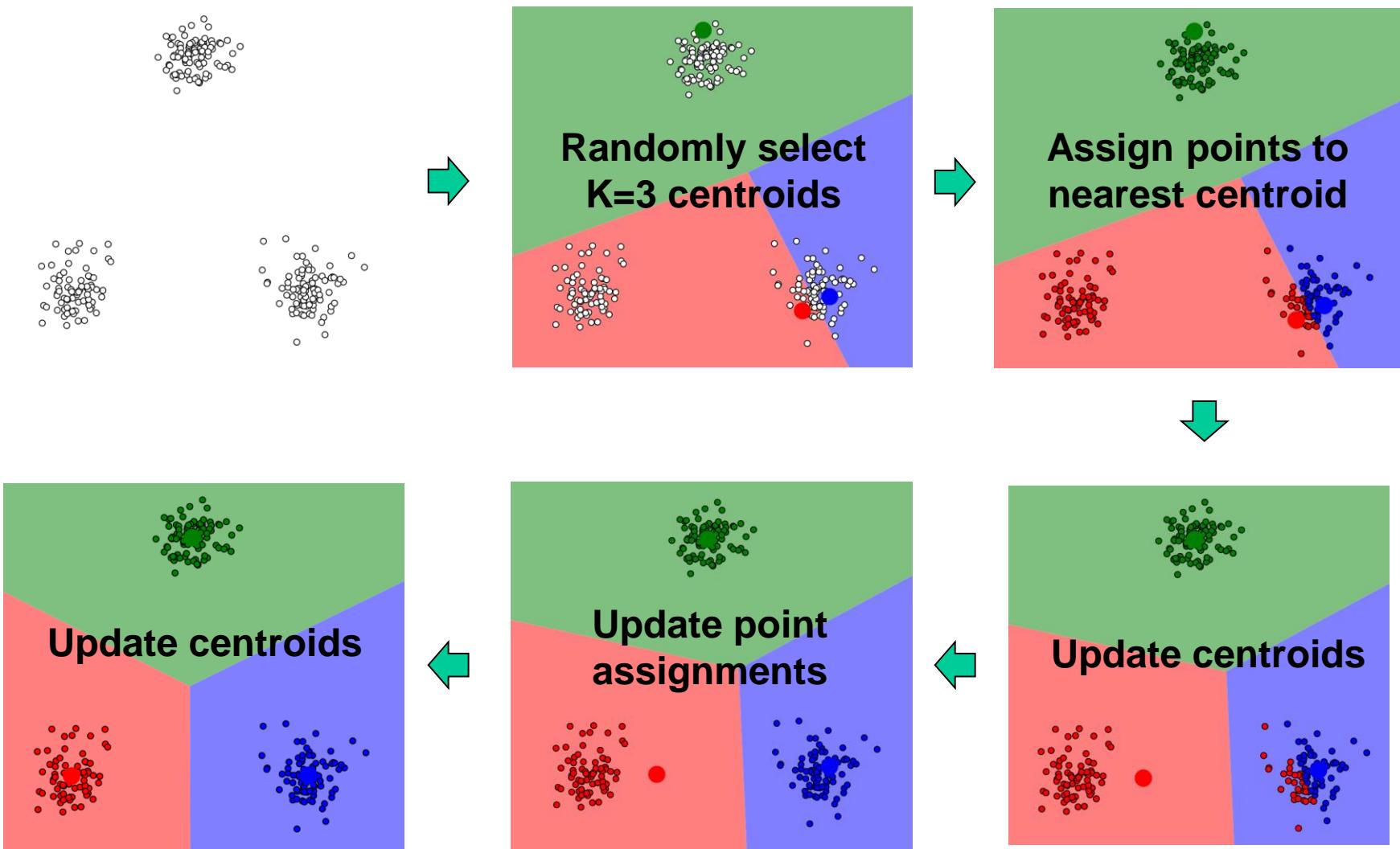
# Application Of Gene Expression Analysis

To understand the blueprint of biological systems and diseases



# A Break From Biology: Basic Clustering Techniques

# An Illustration Of K-Mean Clustering



# Characteristics Of K-Mean Clustering

- The number of clusters,  $K$ , is specified in advance.
- Euclidean distance
  - The nearest centroid minimizes the sum of squares,  $\|x-m\|^2$ .
- Always converge to a (local) minimum.
  - Poor starting centroid locations can lead to incorrect minima.

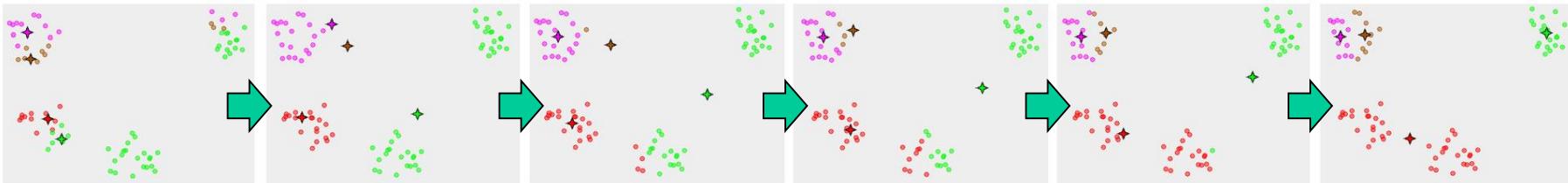
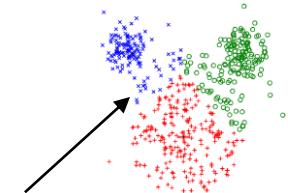
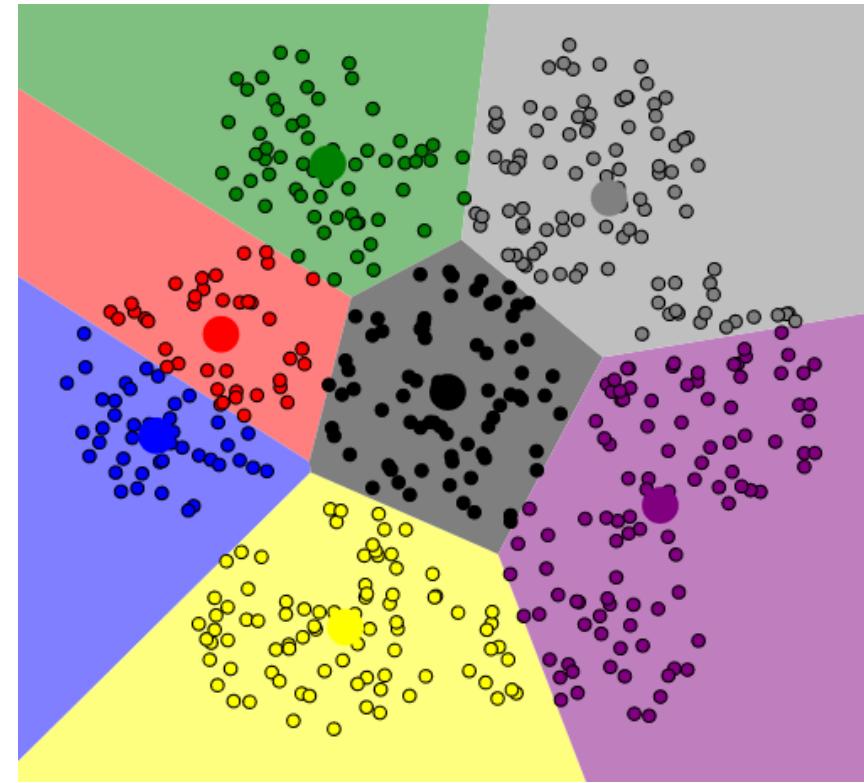
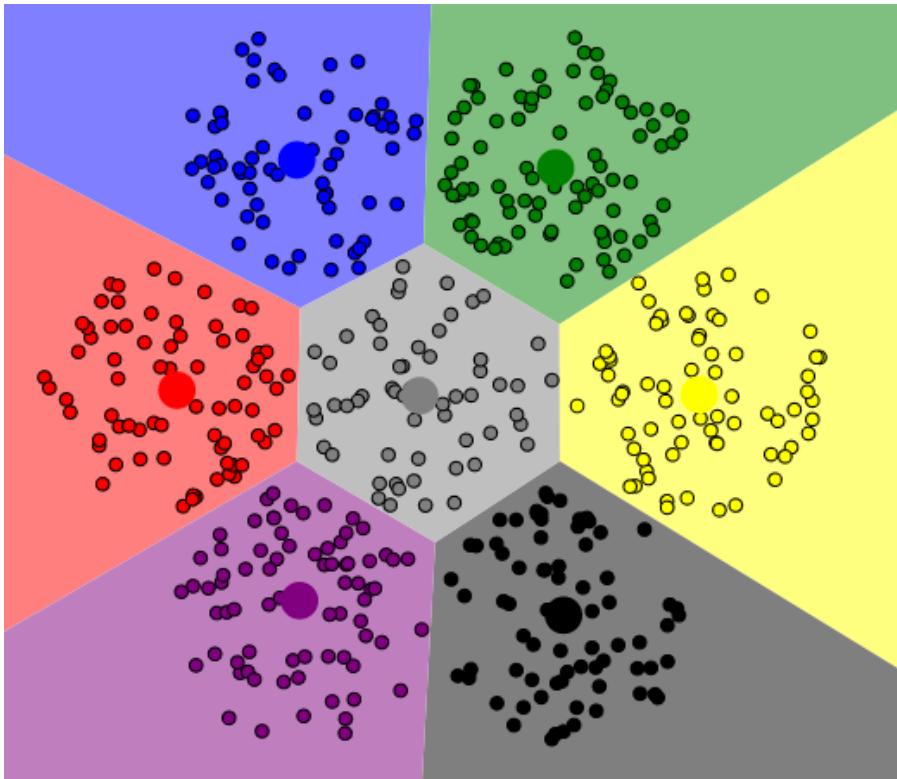


Image from [https://en.wikipedia.org/wiki/K-means\\_clustering](https://en.wikipedia.org/wiki/K-means_clustering)

- The model has several implicit assumptions:
  - Data points scatter around cluster's centers.
  - Boundary between adjacent clusters is always halfway between the cluster centroids.



# Effect Of Poor Initial Centroid Locations



# Distance Functions

- **Property of distance function**

- $d(x, y) \geq 0$  **non-negativity**
- $d(x, y) = 0 \Leftrightarrow x = y$  **identity**
- $d(x, y) = d(y, x)$  **symmetry**
- $d(x, z) \leq d(x, y) + d(y, z)$  **triangle inequality**

- **Example of distance functions**

- **Euclidean distance**
- **Squared Euclidean distance**
- **Manhattan distance**
- **Maximum distance**

$$d(x, y) = \sqrt{\sum (x_i - y_i)^2}$$

$$d(x, y) = \sum (x_i - y_i)^2$$

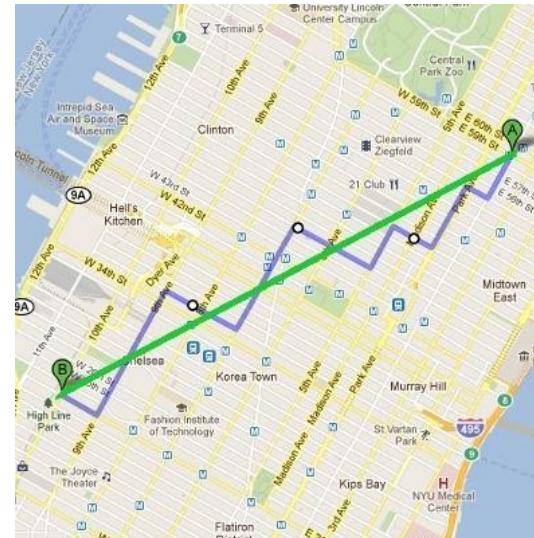
$$d(x, y) = \sum |x_i - y_i|$$

$$d(x, y) = \max |x_i - y_i|$$

- $\|\mathbf{x}\|_p = (\sum_i |\mathbf{x}_i|^p)^{1/p} \leftarrow p\text{-norm}$

# More About Distance Functions

Image from <https://www.quora.com/What-is-the-difference-between-Manhattan-and-Euclidean-distance-measures>



## ■ Manhattan distance

- $d(x, y) = \sum |x_i - y_i|$
- Can reflect driving distance

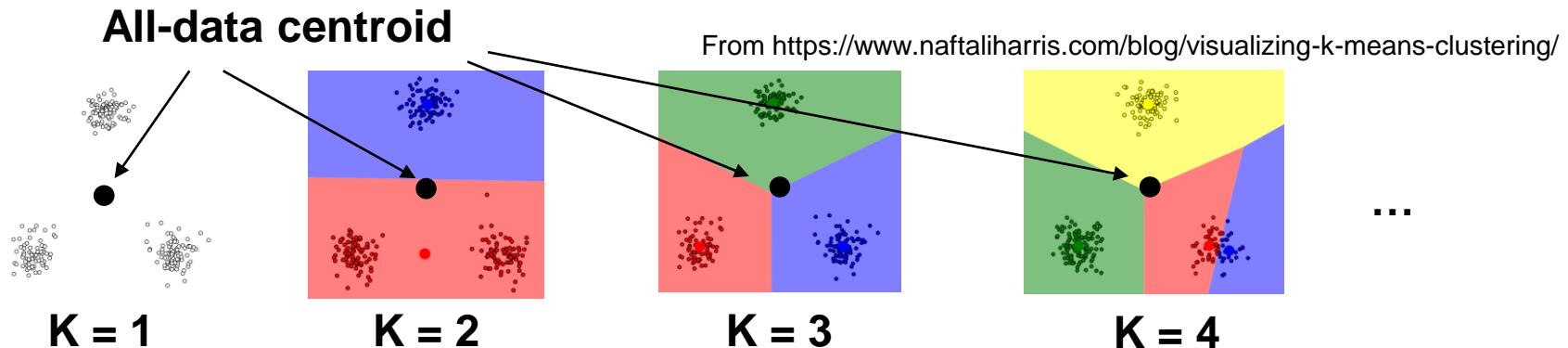
## ■ Hamming (Edit) distance

- For two string  $s$  and  $t$ ,  $d(s, t) = \# \text{ of mismatch positions between the two strings.}$
- Can reflect the extent of evolution between genes: more changes in sequence ~ more time has passed

ATGAGCATAACCATGC~~G~~GAT

ATGAGGATAACCCATGCC~~C~~GAT

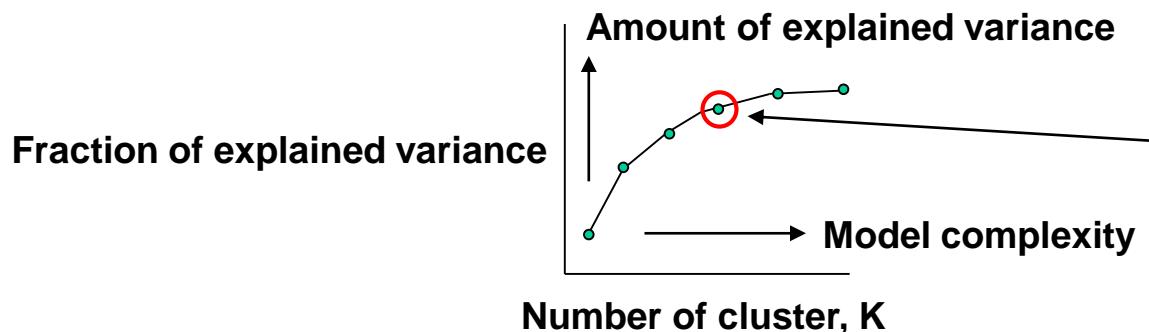
# Elbow Method For Selecting K In K-Mean



fraction of explained variance =  $\frac{\text{between-cluster variance}}{\text{all-data variance}}$

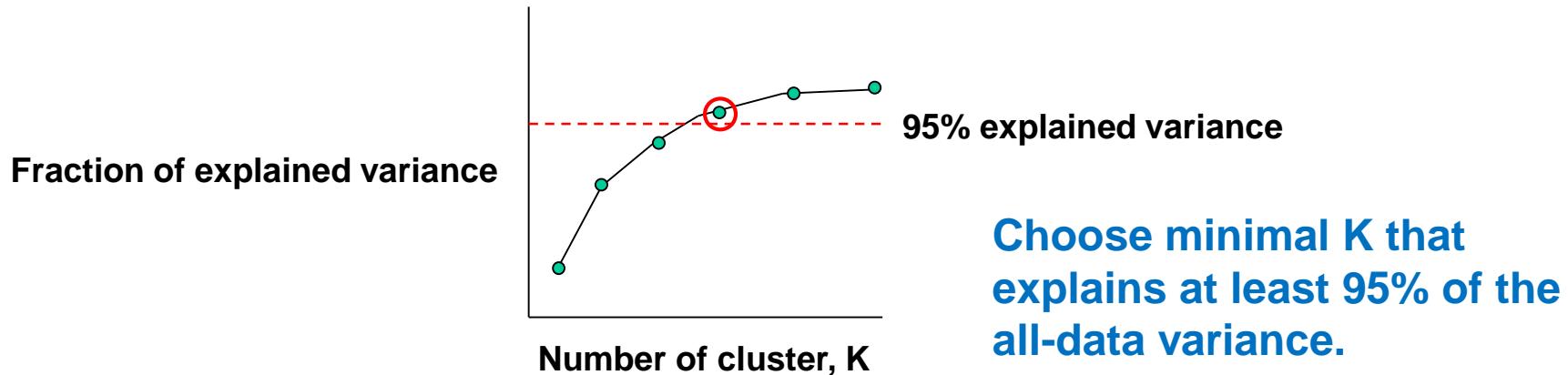
between-cluster variance =  $\sum_{i=1}^K \frac{n_i(M_i - M)^2}{K-1}$ , where  $n_i$  = size of  $i^{\text{th}}$  cluster,  
 $M_i$  = centroid of  $i^{\text{th}}$  cluster, and  
 $M$  = all-data centroid.

all-data variance =  $\sum_{i=1}^N \frac{(x_i - M)^2}{N-1}$ , where  $x_i$  =  $i^{\text{th}}$  data point and  $N$  = # of data.



The elbow method chooses K where increasing complexity doesn't yield much in return.

# Other Ways For Selecting K In K-Mean



$K = 2$   
 $K = 3$   
 $K = 4$   
⋮



Training  
K-mean  
Clustering  
Model



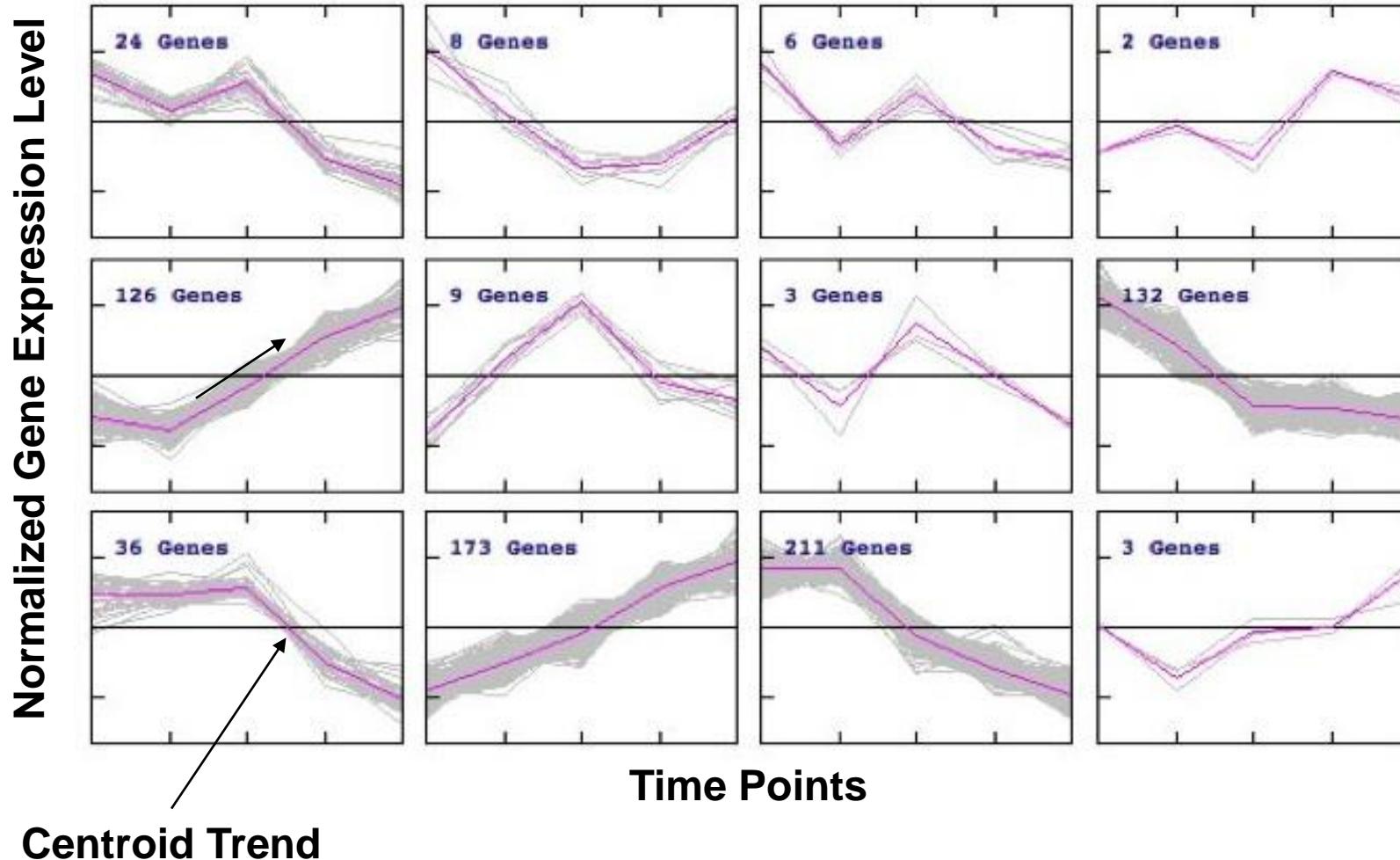
Testing /  
Cross-  
validation



K	Accuracy
2	50%
3	68%
4	83%
⋮	⋮

Choose K that maximizes certain objective (e.g. accuracy on testing data)

# K-Mean Clustering Of Gene Expression



# Hierarchical Clustering

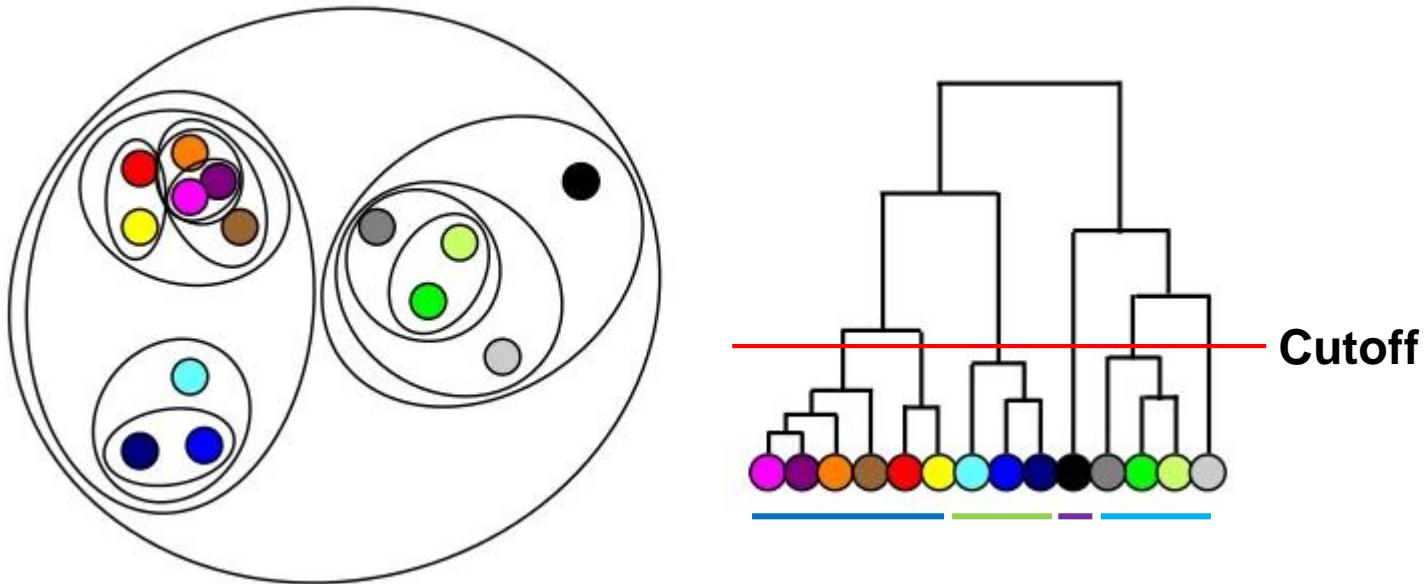
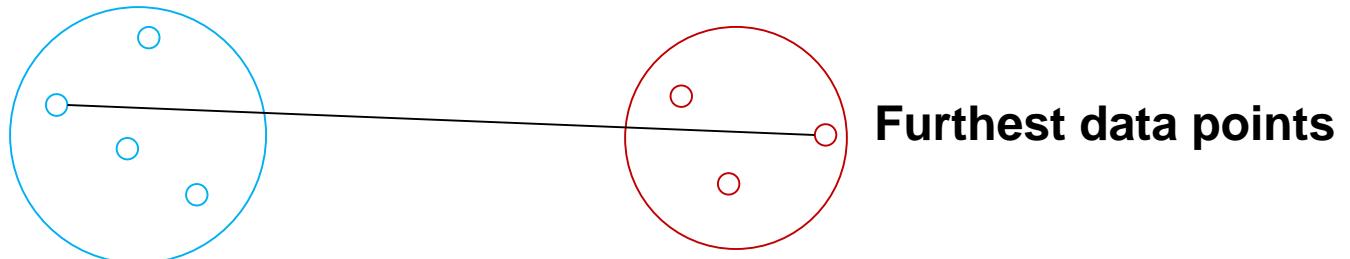


Image from <https://www.slideshare.net/ElenaSgis/data-preprocessing-and-unsupervised-learning-methods-in-bioinformatics>

- **Each step finds two data points or existing clusters that are closest to each other and group them together.**
- **The number of clusters is defined afterward by setting a cutoff on the distance.**
- **Choices of distance function.**
- **Choices of how to measure distance between clusters (e.g. using cluster centroids or closest members or all members).**

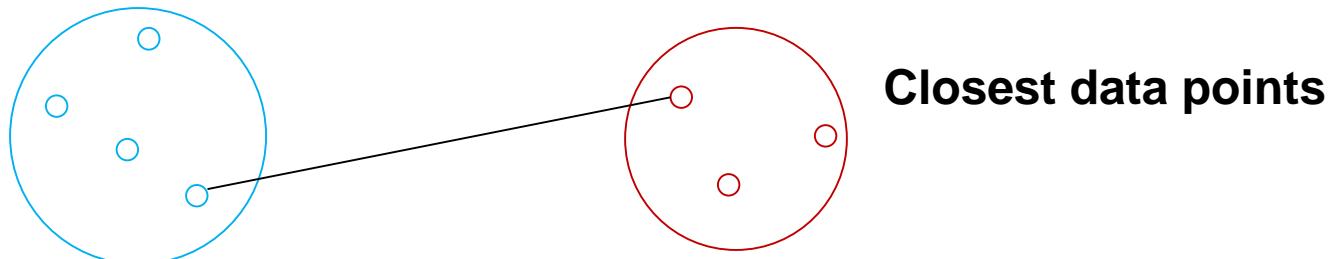
# Linkage Criteria: Distance Between Clusters

- Maximum or complete linkage



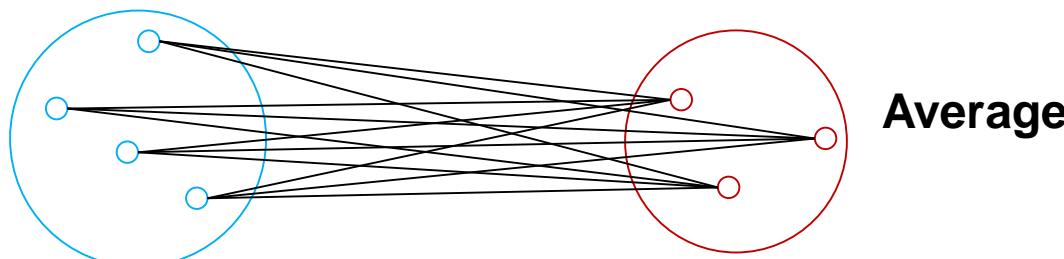
Furthest data points

- Minimum or single linkage



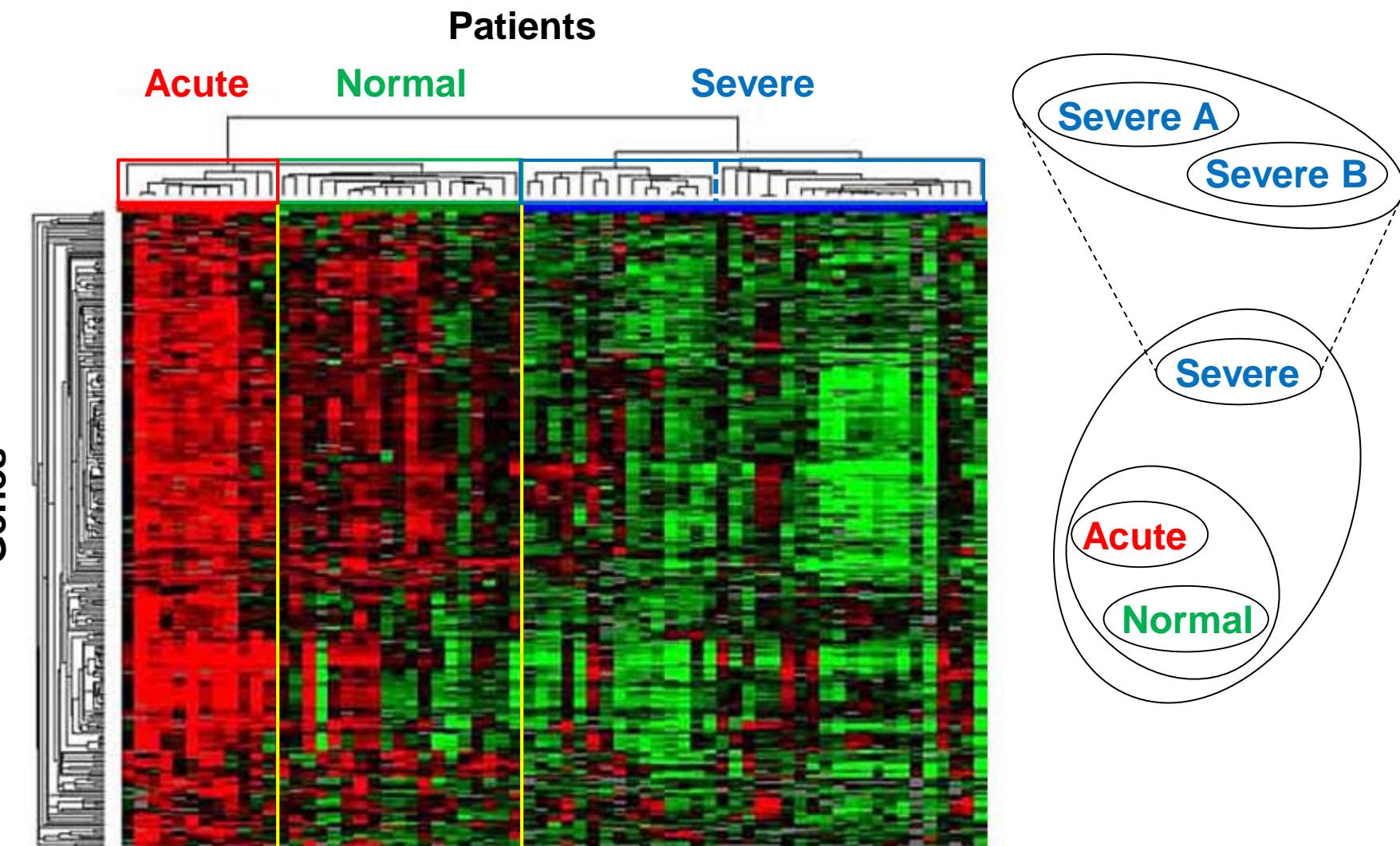
Closest data points

- Mean or average linkage



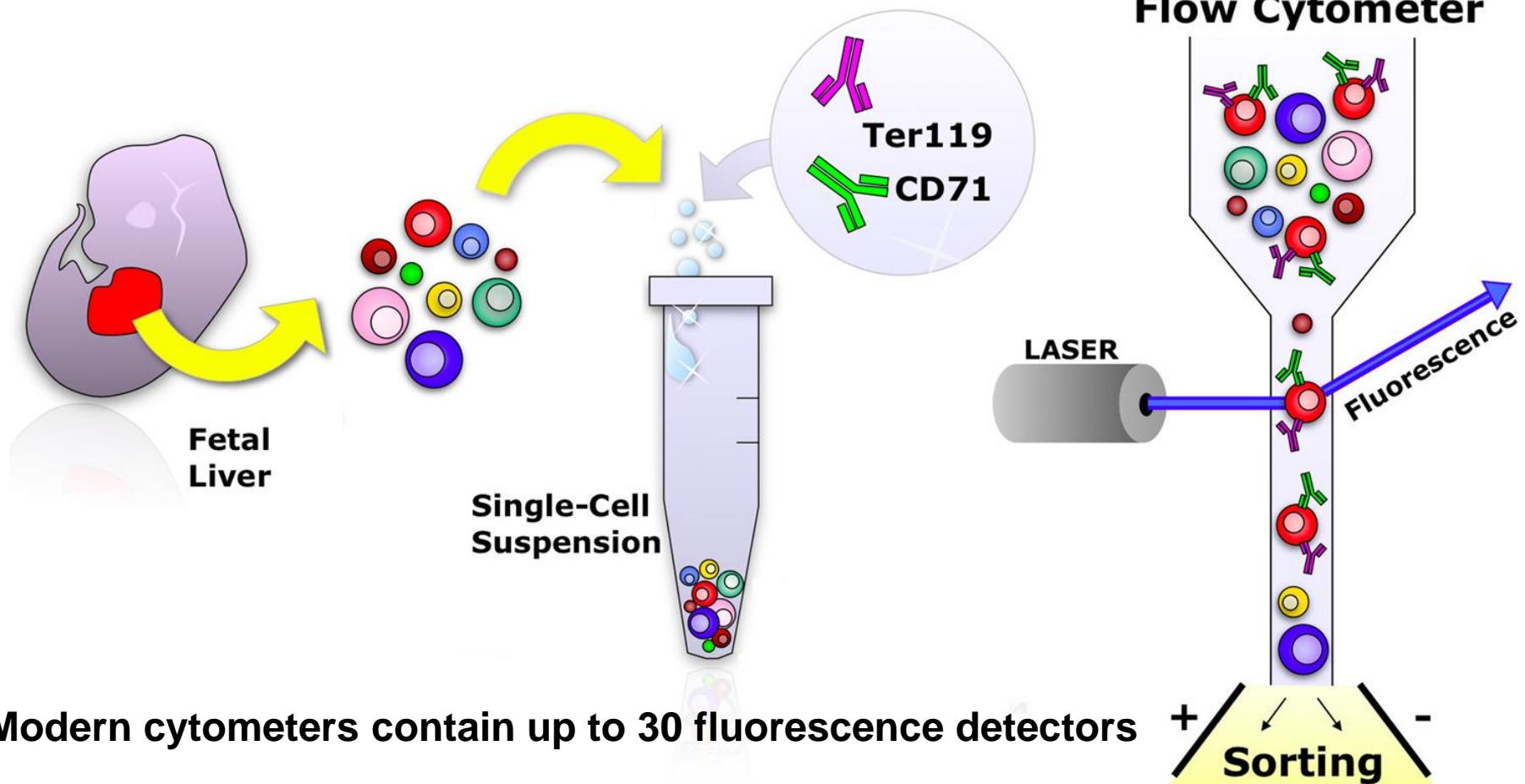
Average

# Hierarchical Clustering Of Gene Expression



# **Application IV: Annotation of Different Cell Types**

# Sorting Different Cell Types



Modern cytometers contain up to 30 fluorescence detectors

Image from [http://labs.umassmed.edu/socolovskylab/research~flow\\_cytometry.html](http://labs.umassmed.edu/socolovskylab/research~flow_cytometry.html)

# Example Of Flow Cytometry Data

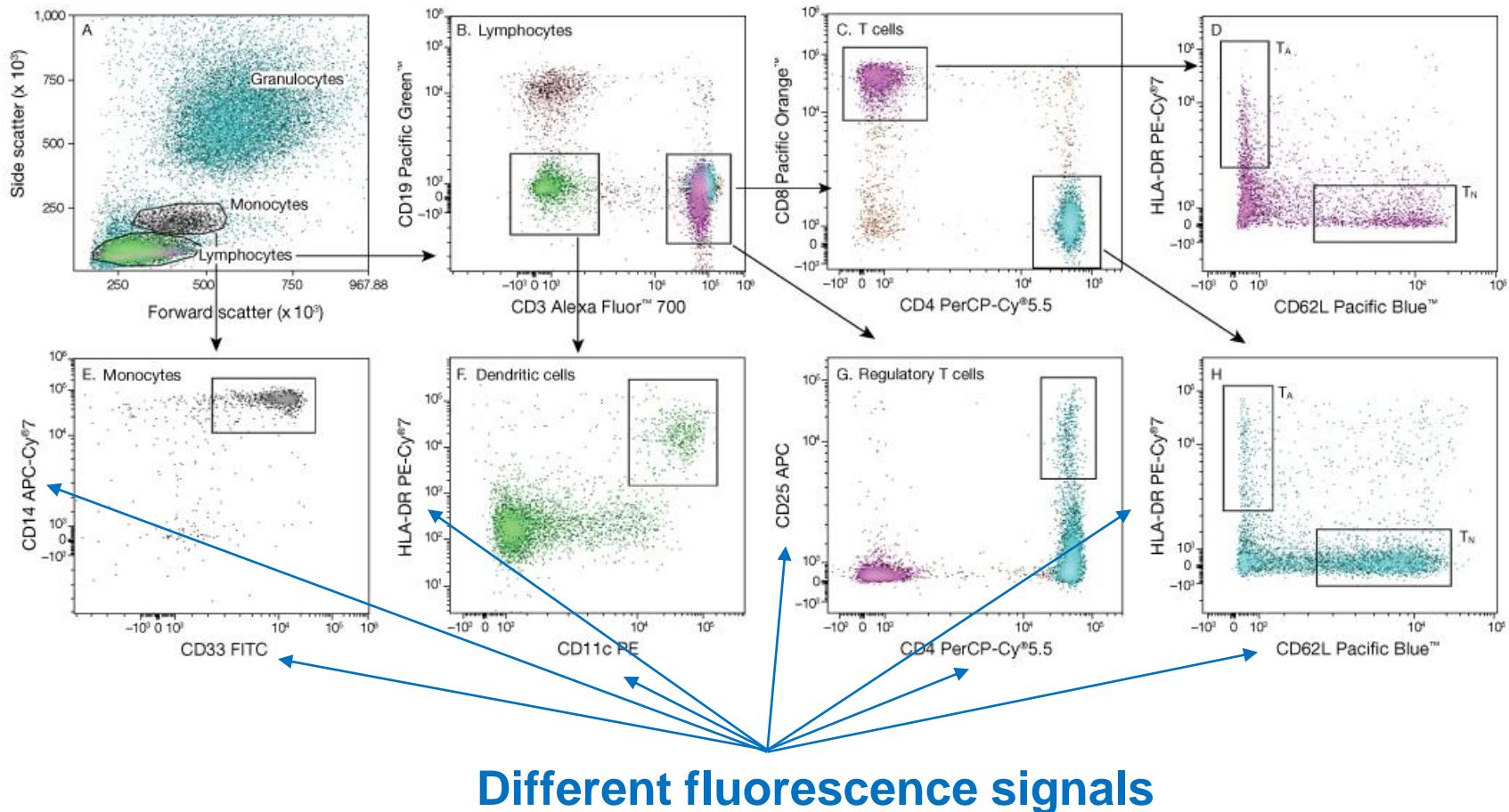
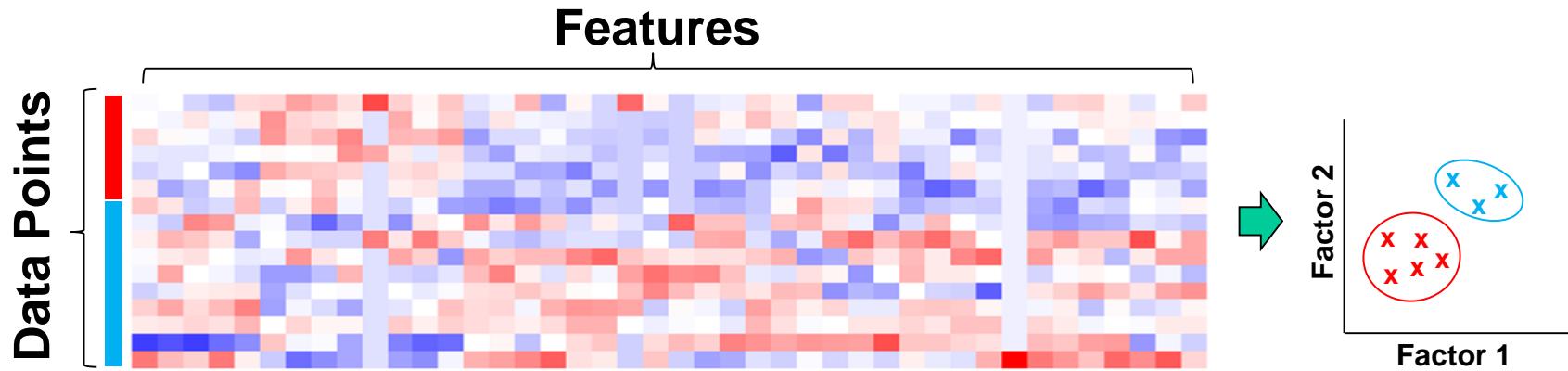


Image from <https://www.thermofisher.com/th/en/home/references/newsletters-and-journals/bioprobe...>

# Simplifying High-Dimensional Data



## Dimensionality Reduction

PCA, t-SNE

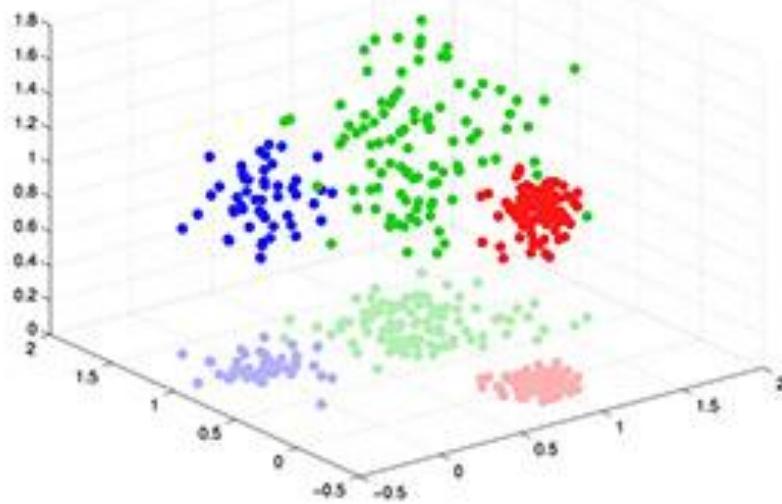


Image from <http://bigdata.csail.mit.edu/node/277>

## Classification

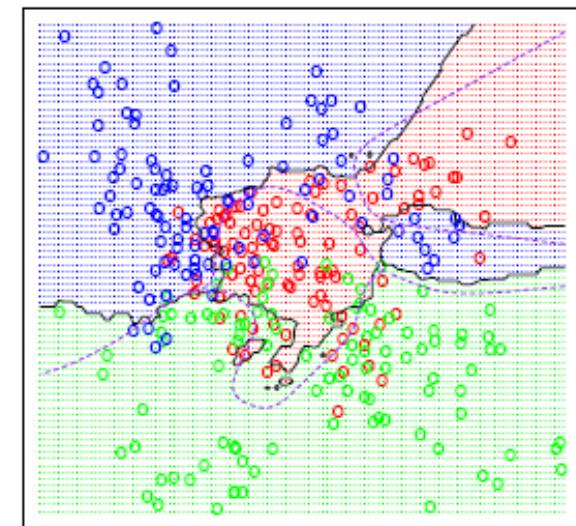
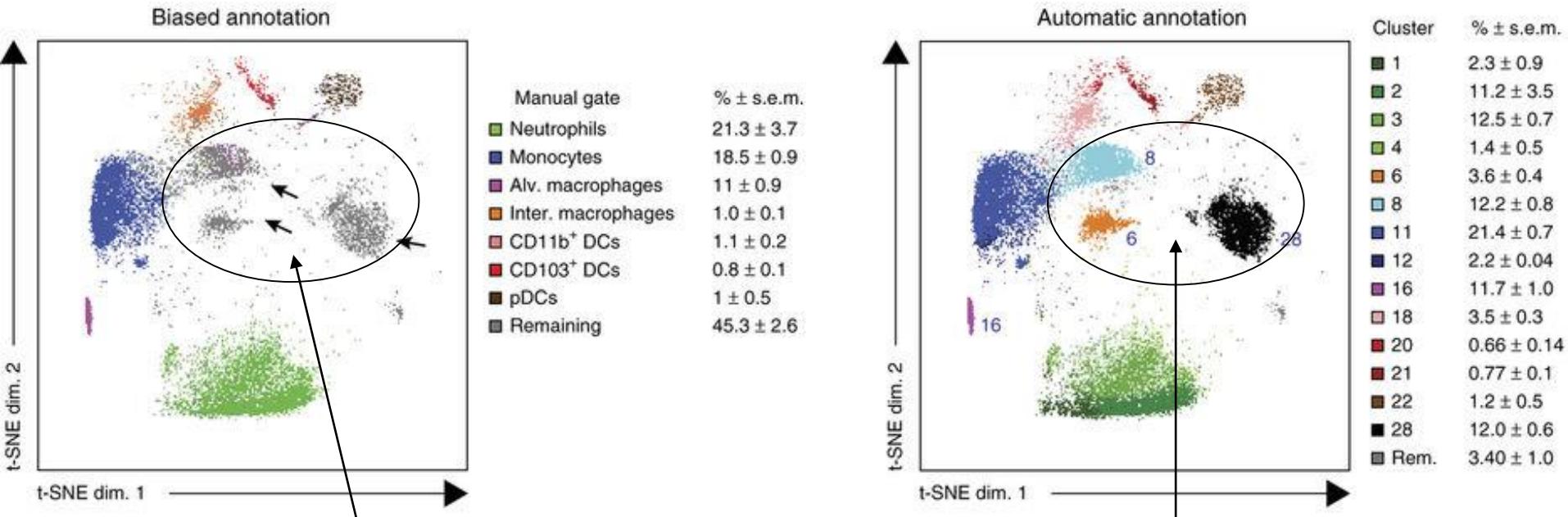


Image from <http://www.stat.ucla.edu/~ybzhao/teaching/stat101c/>

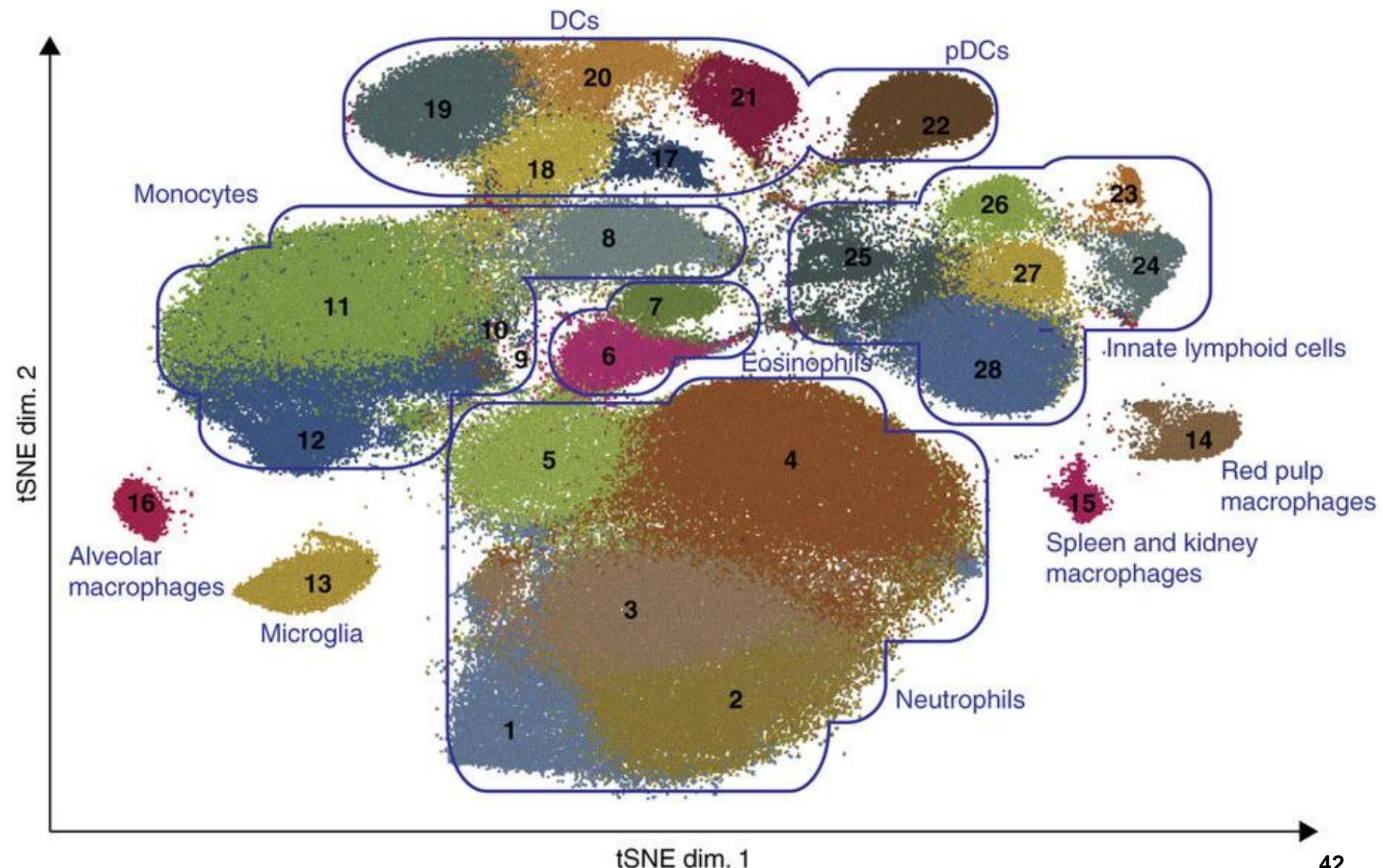
# Identification Of New Cell Types



**Unannotated cell types according to existing knowledge**

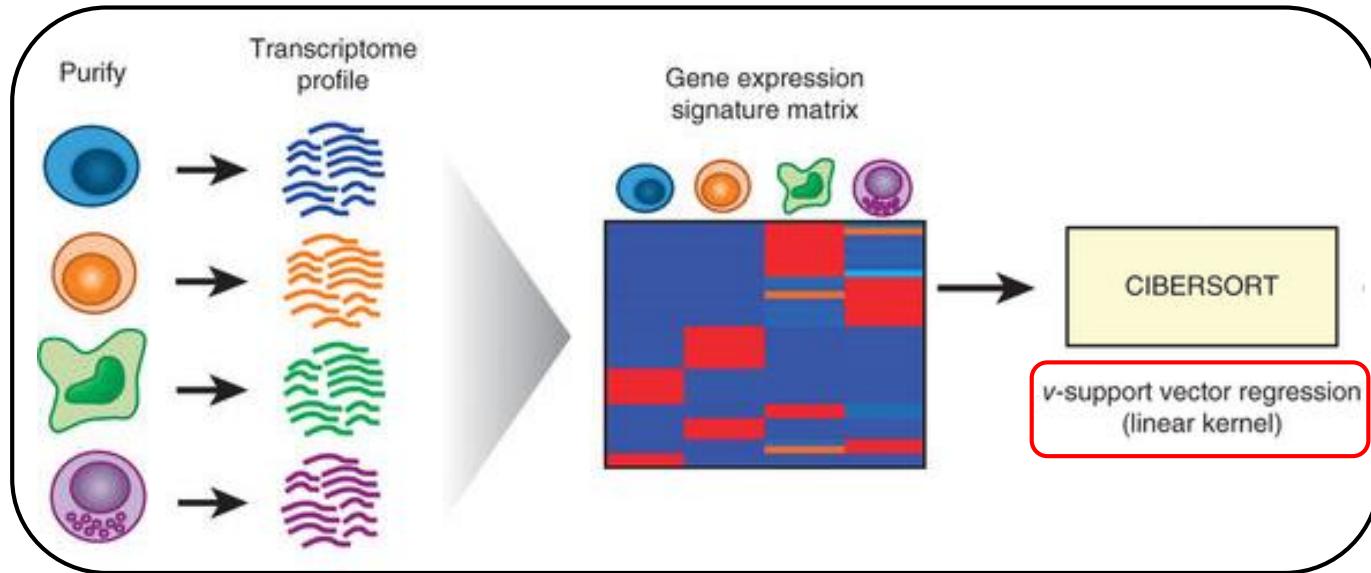
**Clustering process defines these as new cell types**

# Automatic Cell Type Classification

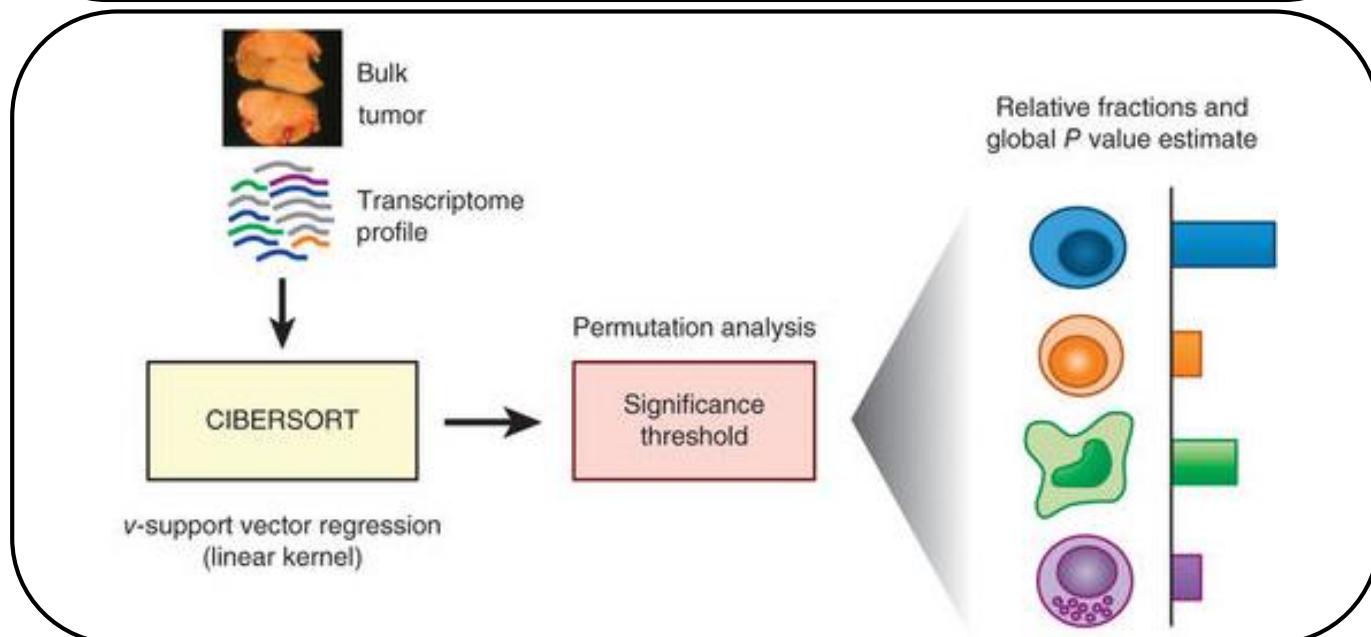


# Deconvolution Of Cell Type Composition

## Training

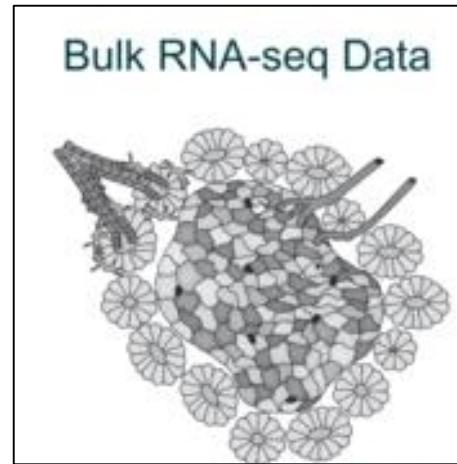


## Testing

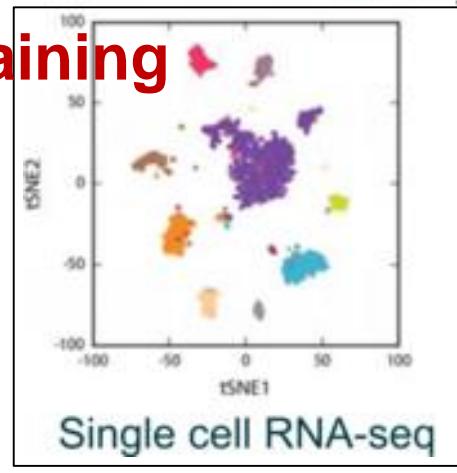


# Cell Composition Reflects Disease State

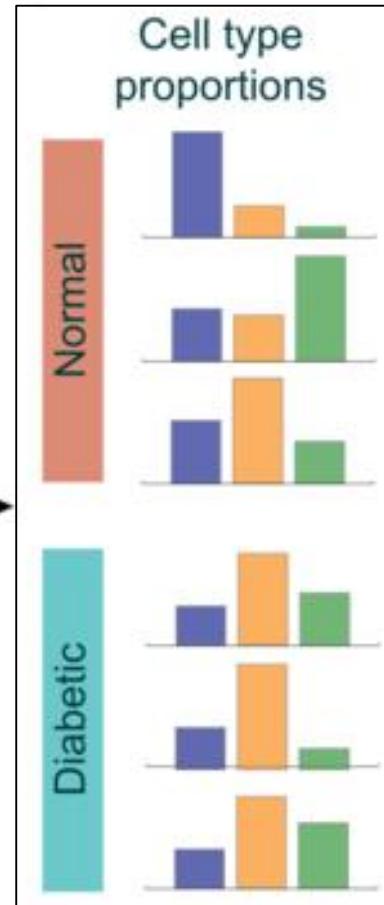
## New Data



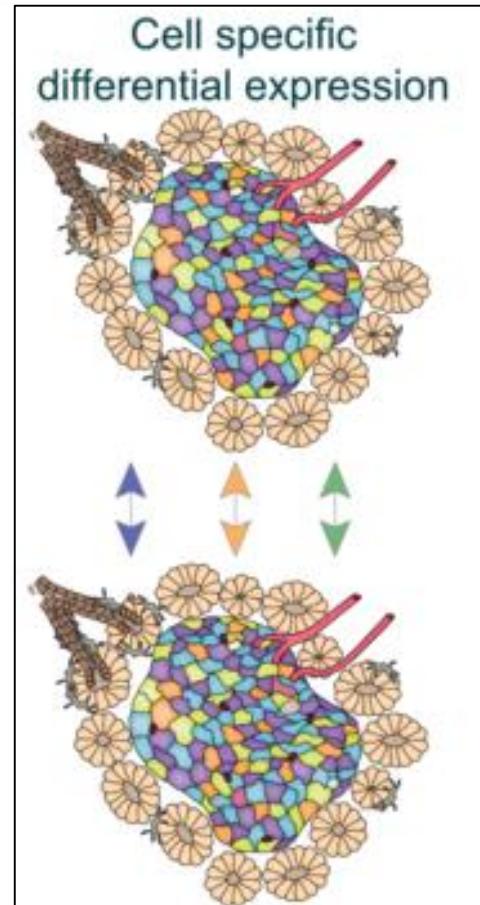
Deconvolution



*Bseq-SC*



Deconvolution



Diagnosis

# **Application V: Protein-To-DNA Binding Prediction**

# Transcription Factor (TF) Binding

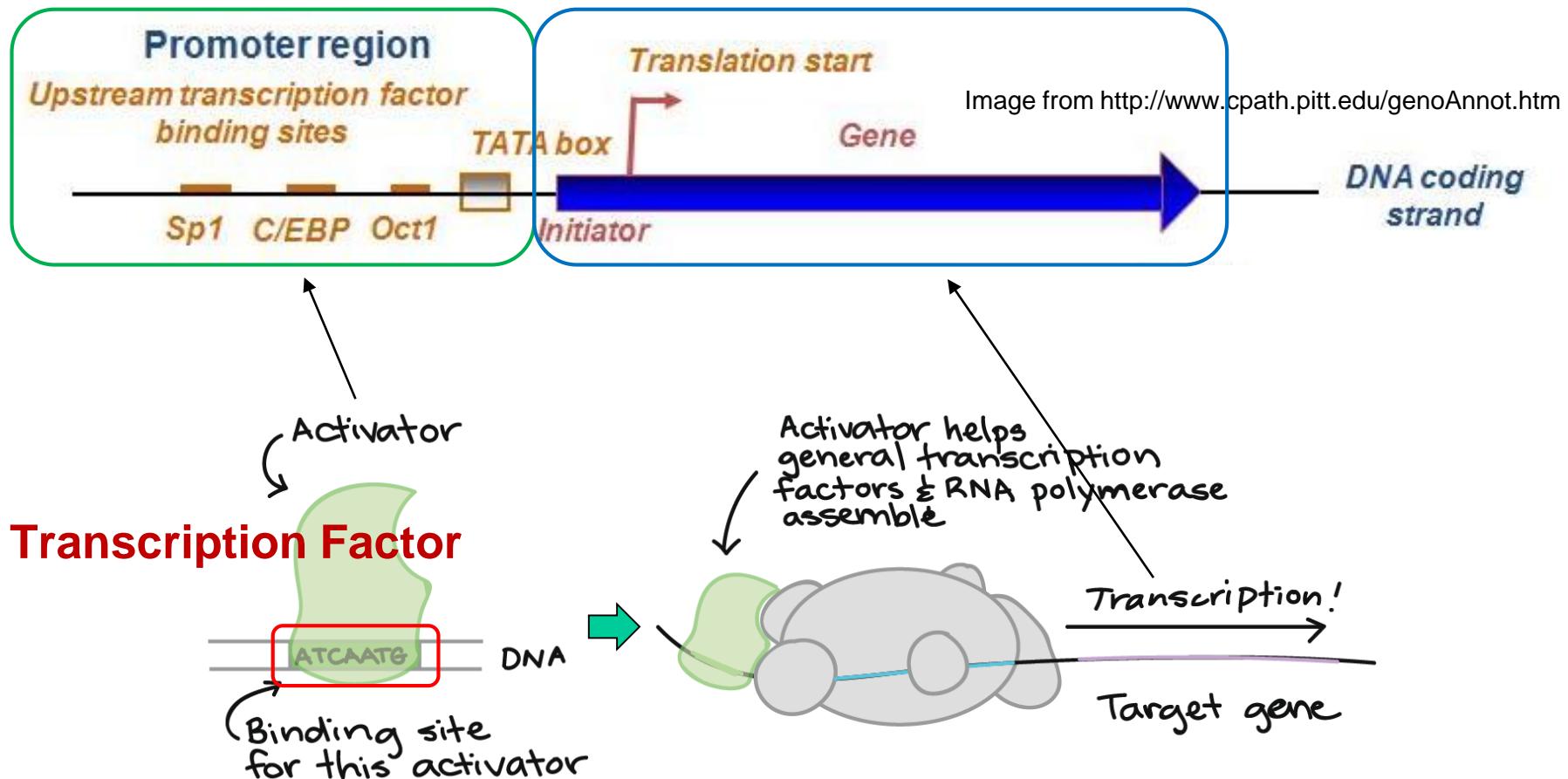


Image from <https://www.khanacademy.org/science/biology/gene-regulation/gene-regulation-in-eukaryotes/a/eukaryotic-transcription-factors>

# Sequence-Based Prediction Of TF Binding

## Expected Motif



**High chance of false positives due to flexibility of binding motifs.**

## Genome Sequence

... ACCAGCGCGAAGCTCGGGCGGAGGGT  
TGAGCCACATGAGGCATGGCGACA **TCCCATA**  
**TATGG**AGACATGGCGTGGCTGGCTGTTACATT  
TTGTTTGATGAAAAGCATAACCATGCGGATG  
ATATTTTATTATAGACTAGAGATGATTATTG  
AATAGAC **ATGCTCTAACCA**TTAAC**TCTA**  
**ATTCCAC**...  
... ACCAGCGCGAAGCTCGGGCGGAGGGT  
TGAGCCACATGAGGCATGGCGACAGGGACCT  
**CCGACCTTATAAGG**AGACATGGCGTGGCTGGC  
TGTTACATTTGTTTGATGAAAAGCATAACC  
ATGCGGATGATATTTTATTATAGACTAGAGA  
TGATTATTGAATAGGCCTACTTAC **ATGCTCT**  
**TAACCATTTAAC**TCTAATT**CCAC**...

- Predicted binding sites
- Downstream genes

# DNA Packaging Affects TF Binding

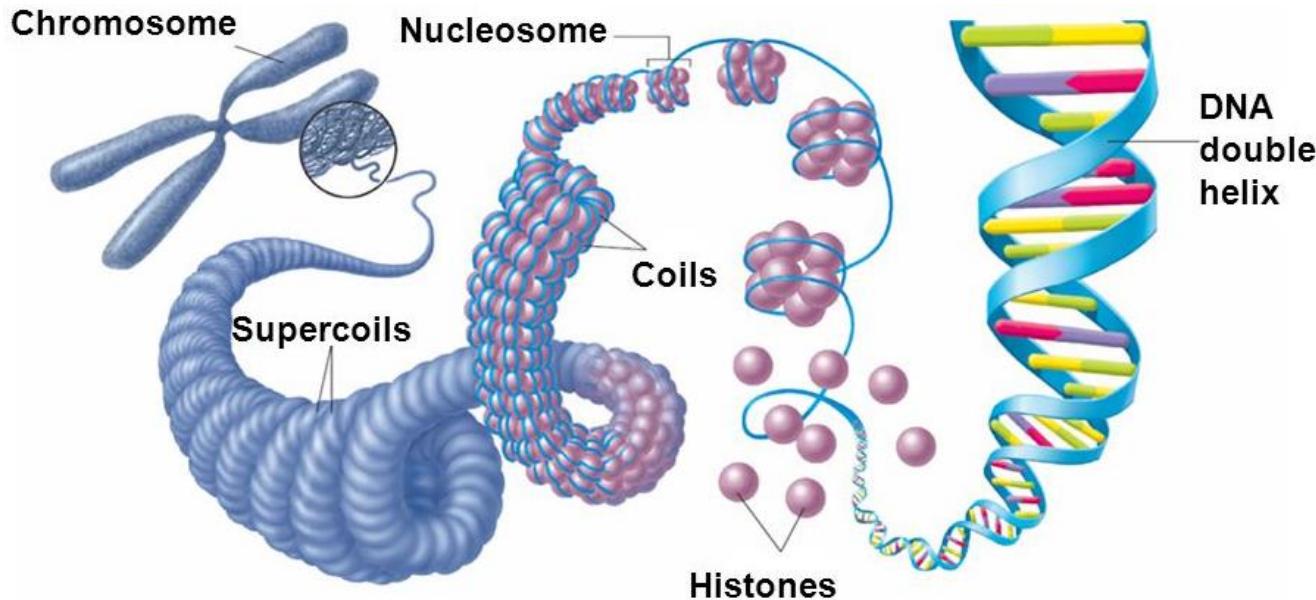
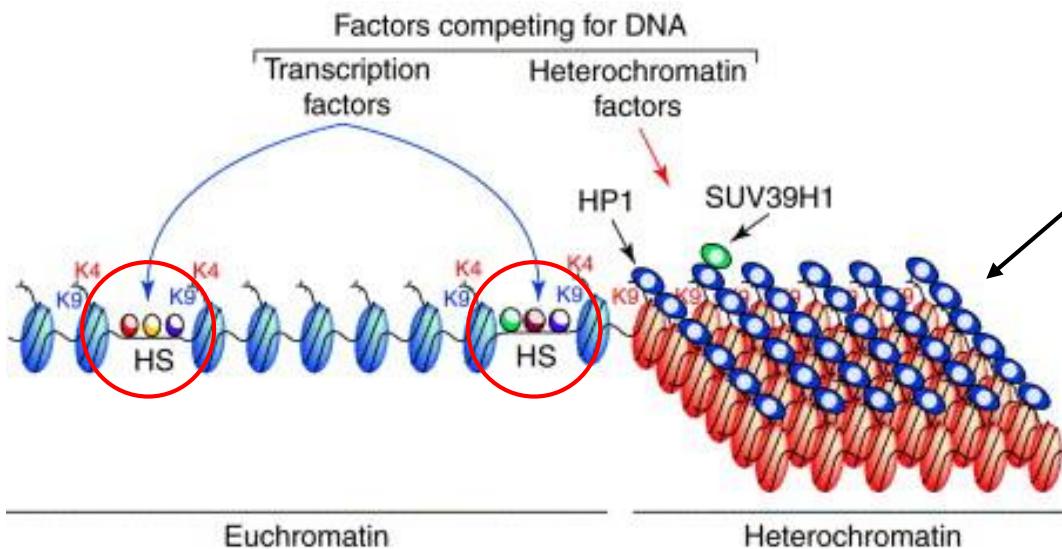


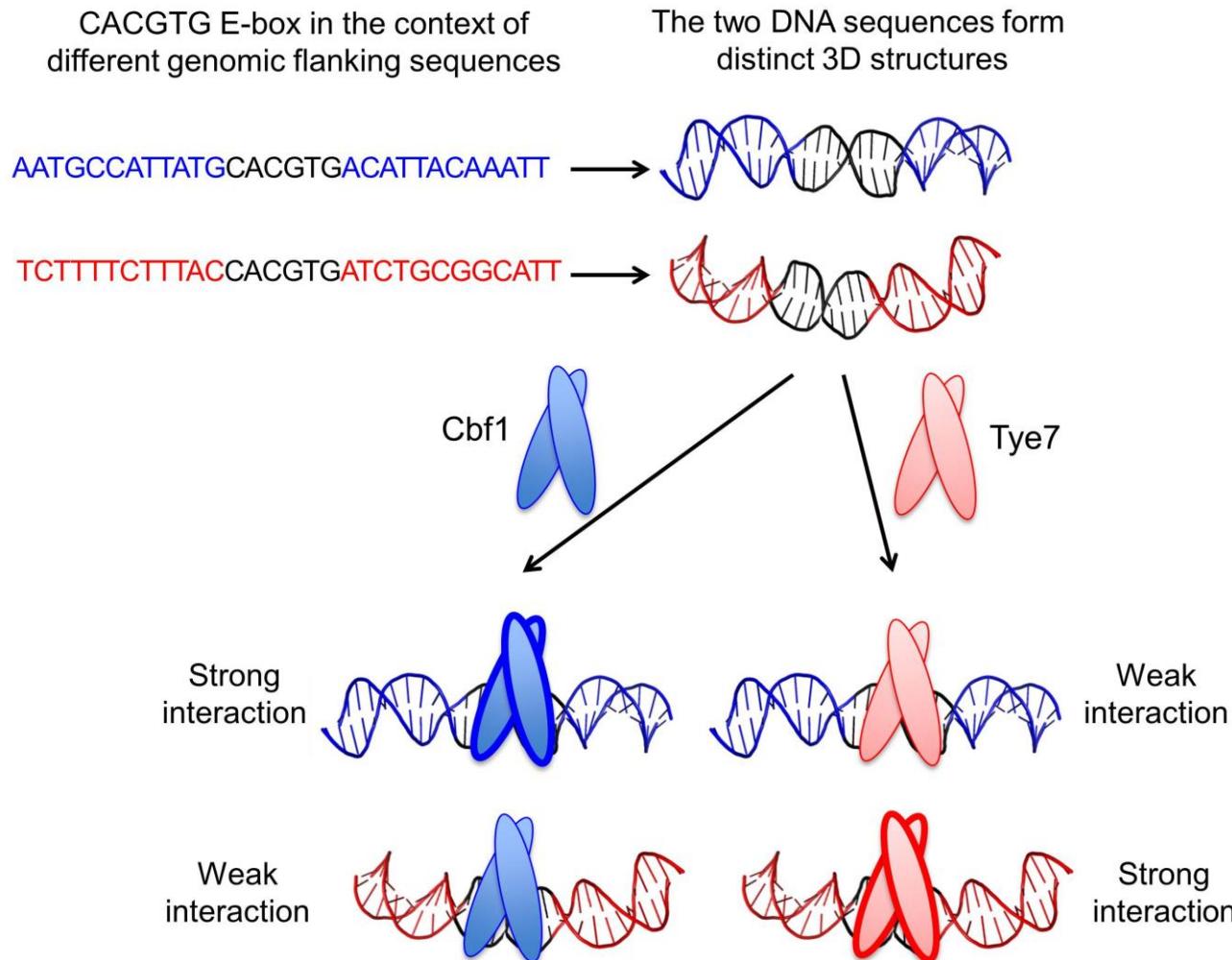
Image from [http://images.slideplayer.com/9/2508104/slides/slide\\_28.jpg](http://images.slideplayer.com/9/2508104/slides/slide_28.jpg)



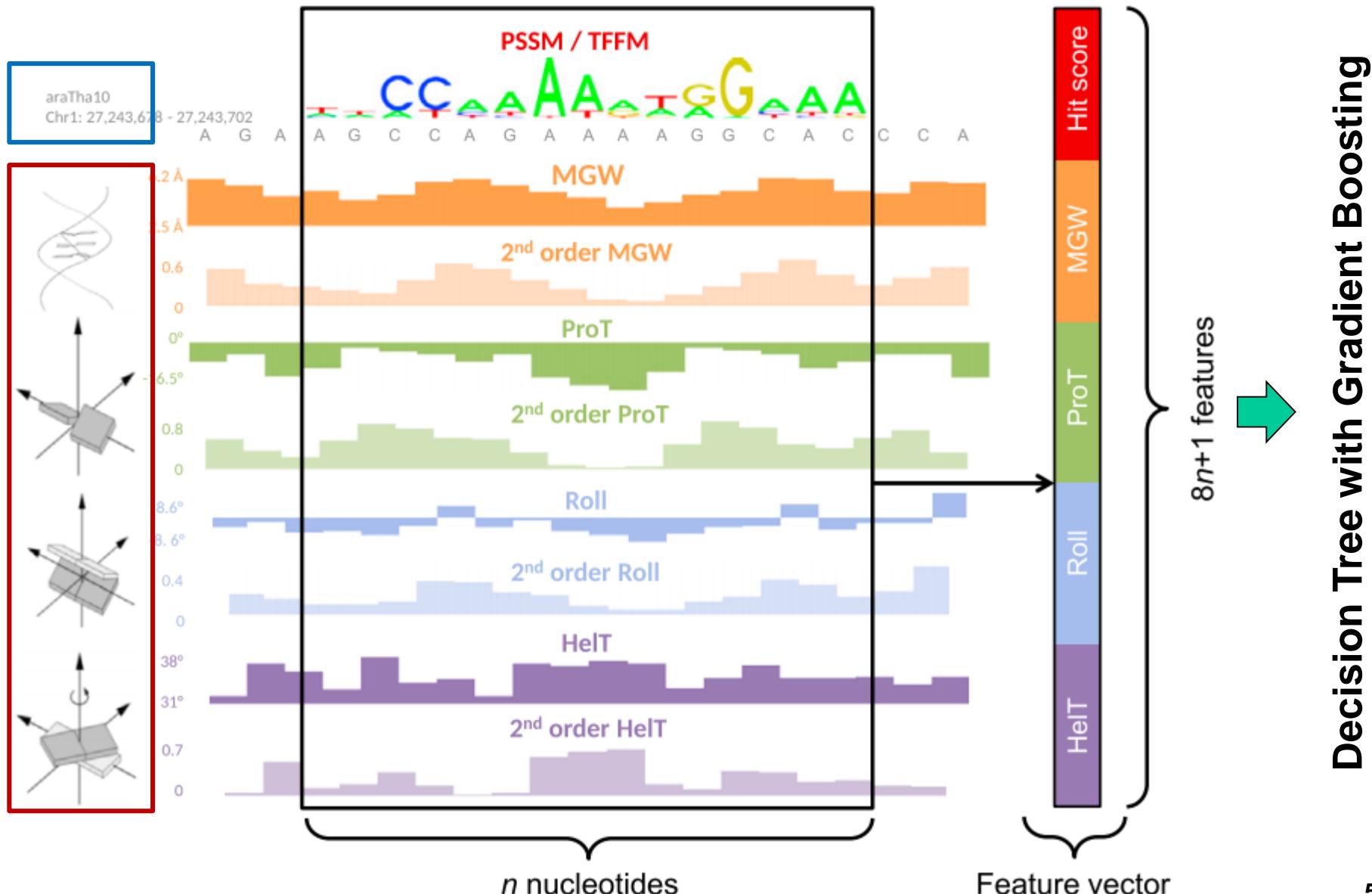
**Transcription factors  
cannot access tightly  
packed DNA region.**

Dillon *et al.* Trends in Genetics 18, 252-258 (2002)

# Local DNA 3D Structure Affects TF Binding



# Prediction Of TF Binding Using Physical Data



# **Application VI: Human Genome Structure**

# Human Genome

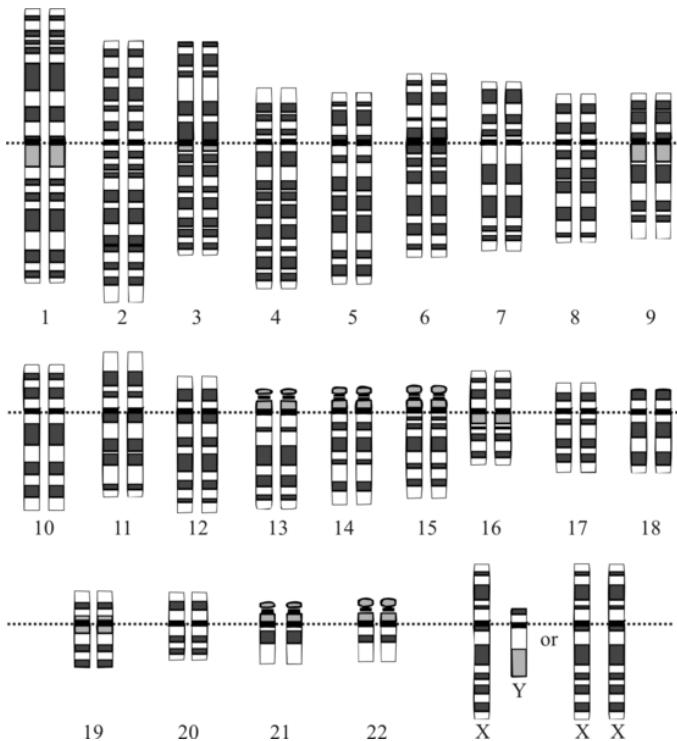


Image from [https://en.wikipedia.org/wiki/Human\\_genome](https://en.wikipedia.org/wiki/Human_genome)  
#Coding\_vs.\_noncoding\_DNA

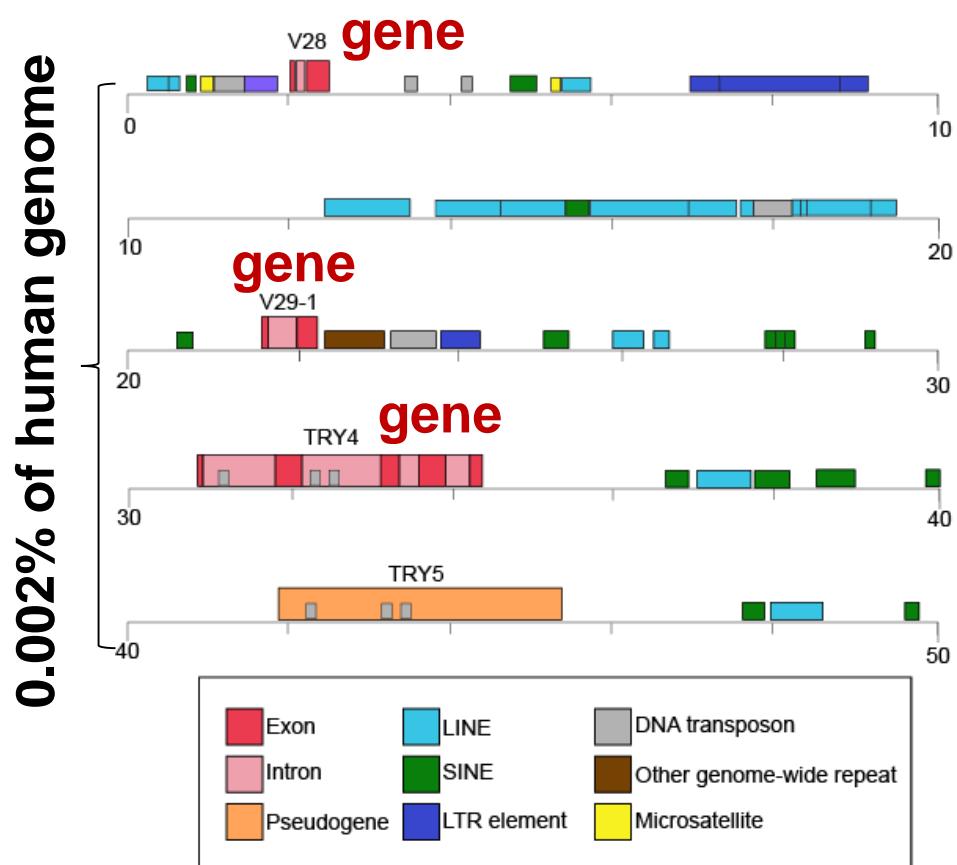
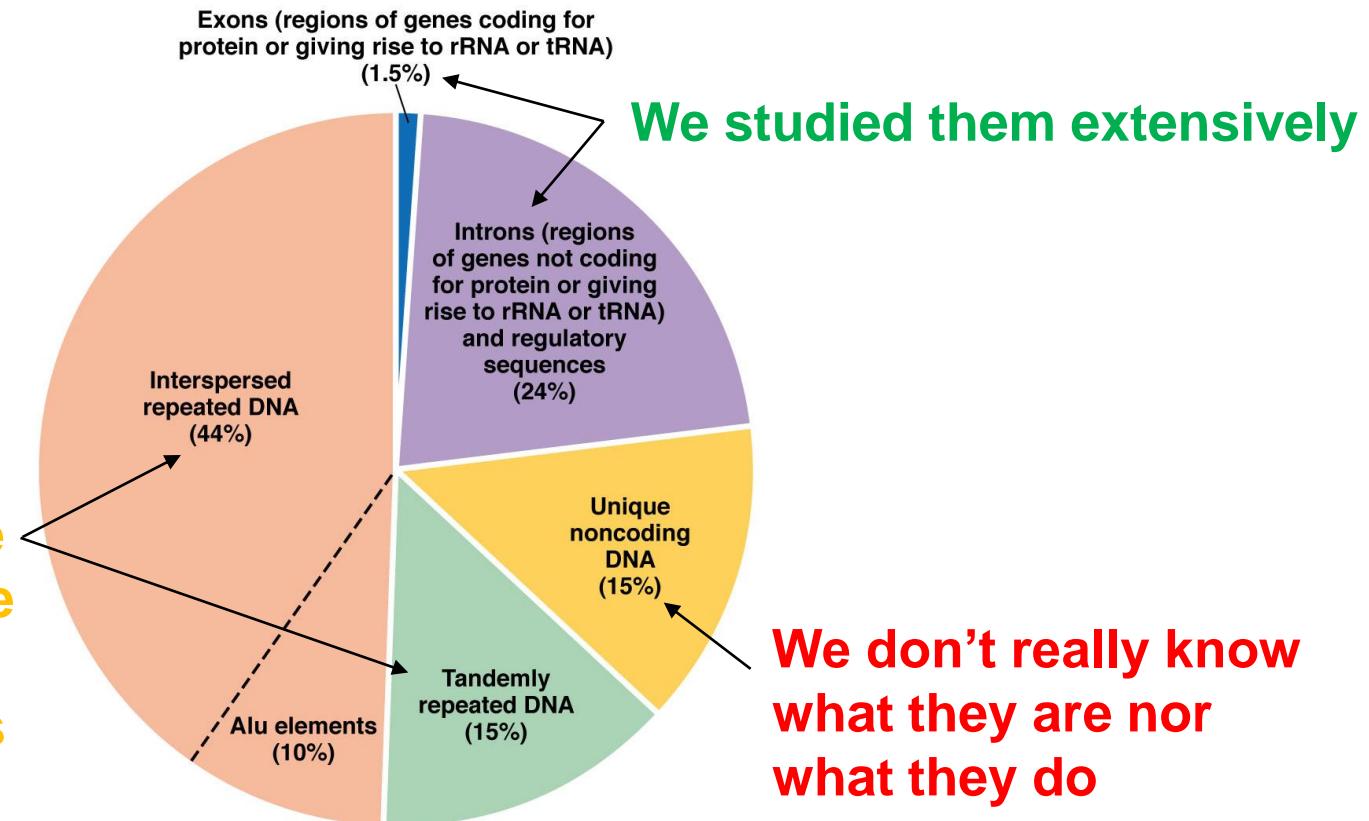


Image from [https://bio.libretexts.org/LibreTexts/University\\_of\\_California\\_Davis/BIS\\_2A%3A\\_Introductory\\_Biology\\_\(Easlon\)/Readings/26%3A\\_Genomes%3A\\_a\\_Brief\\_Introduction](https://bio.libretexts.org/LibreTexts/University_of_California_Davis/BIS_2A%3A_Introductory_Biology_(Easlon)/Readings/26%3A_Genomes%3A_a_Brief_Introduction)

- **3.2 Giga-basepairs x2 on 23 pairs of chromosomes.**
- **~0.1% variation between individuals.**
- **<1.5% code for proteins (**exons**).**

# Much Of Genome Remains Unknown



© 2012 Pearson Education, Inc.

Image from <http://www.mun.ca/biology/desmid/brian/BIOL2060/BIOL2060-18/CB18.html> /

- Using data from well-studied genomic regions to predict the structures and functions of other genomic regions in the same or in newly discovered species.

# Dynamic Bayesian Network (DBN)

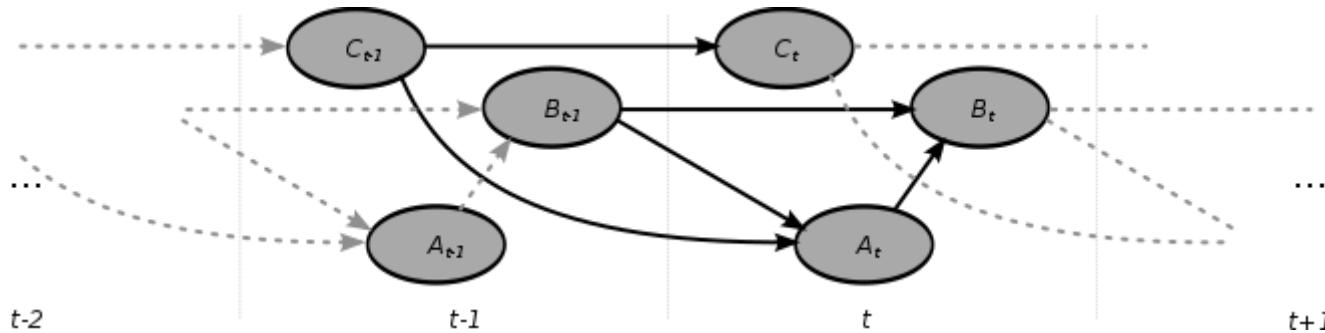


Image from [https://en.wikipedia.org/wiki/Dynamic\\_Bayesian\\_network](https://en.wikipedia.org/wiki/Dynamic_Bayesian_network)

- **The label of a genomic location can be predicted based on the labels and properties of nearby positions.**
- **For example,  $A_t$  may indicate label – such as exon or intron – at location  $t$  while  $B_t$  and  $C_t$  keep track of the properties of the genome – such as coding-vs-non-coding or loosely-packed-vs-tightly-packed – at location  $t$ .**
- **The probability of predicting a label depends on  $B_t$  and  $C_t$ .**

# Histone Modifications

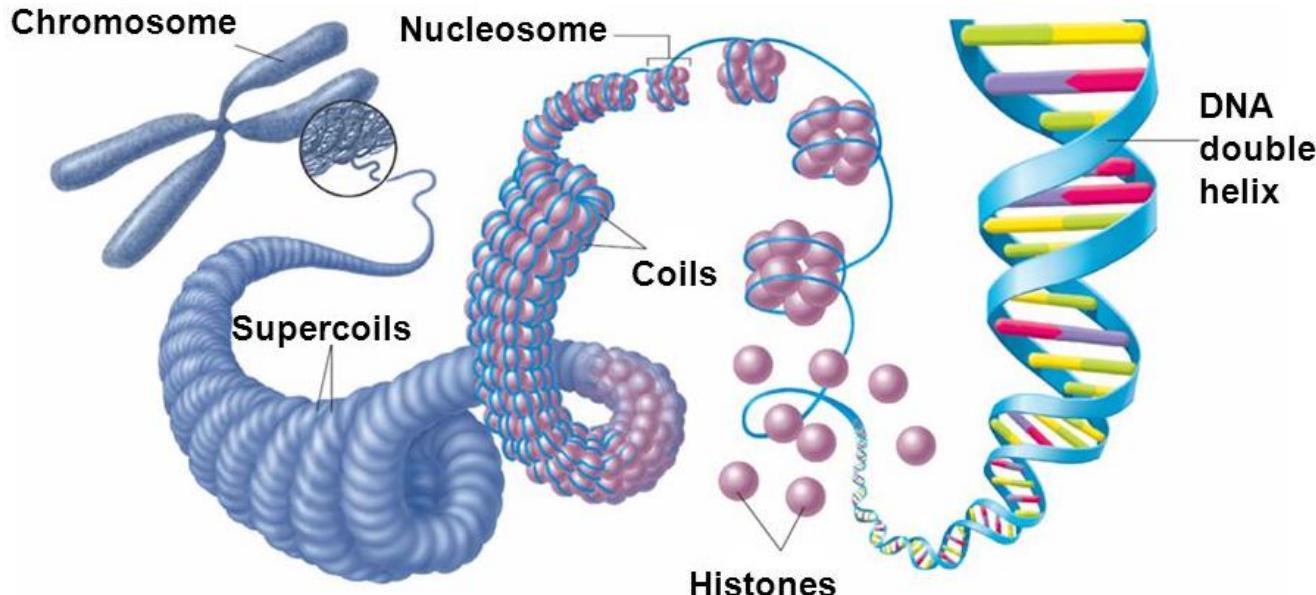


Image from [http://images.slideplayer.com/9/2508104/slides/slide\\_28.jpg](http://images.slideplayer.com/9/2508104/slides/slide_28.jpg)

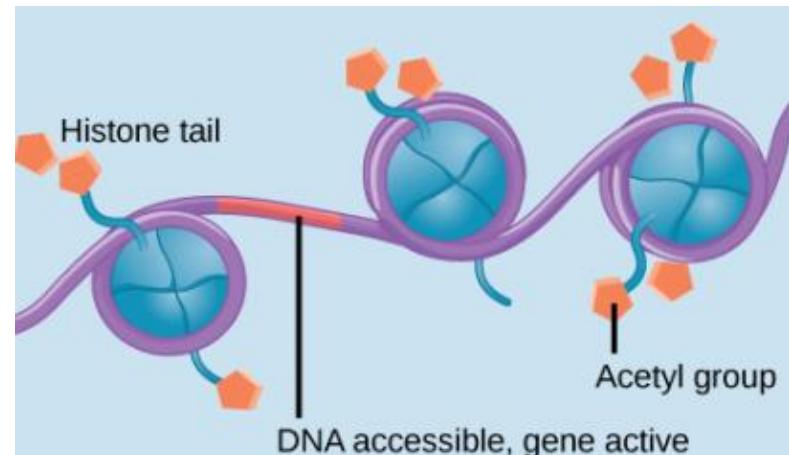
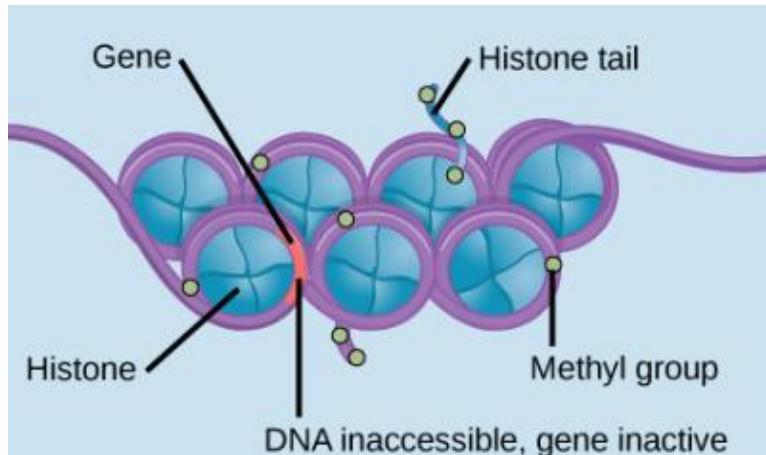
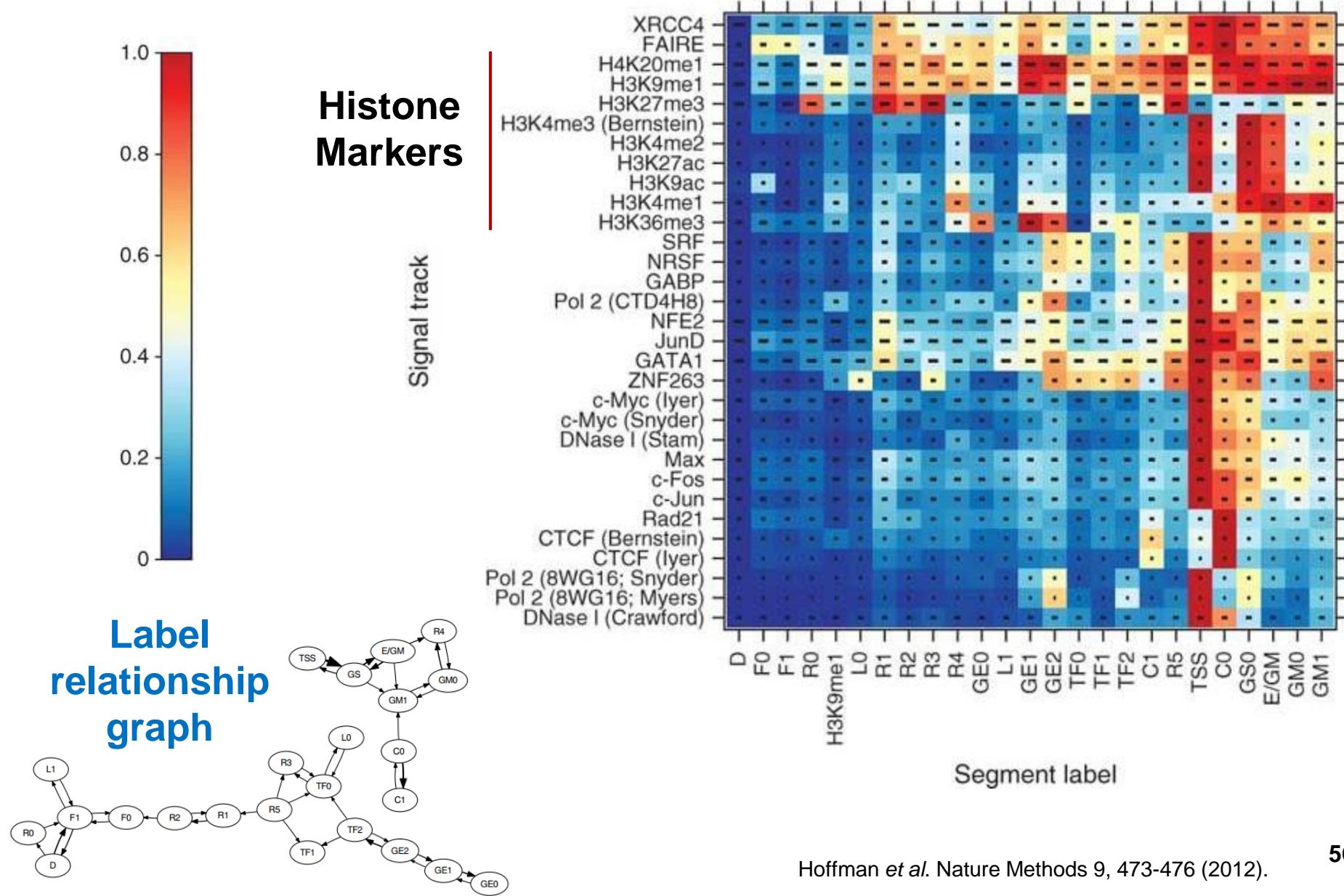
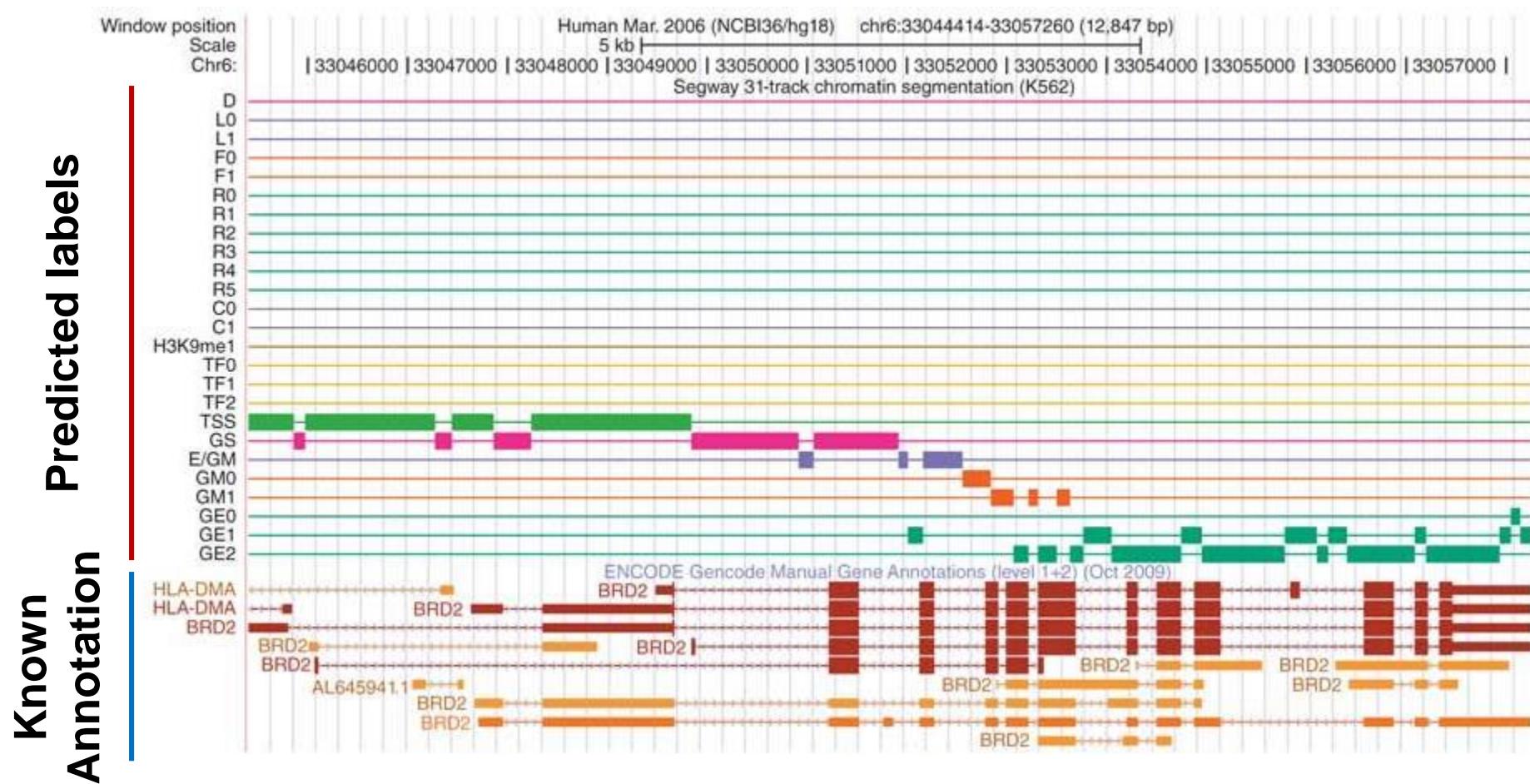


Image from <http://www.peirsoncenter.com/articles/dna-methylation-in-down-syndrome>

# Inferring Labels From Genomic Signals



## Decent Performance Without Using Sequence

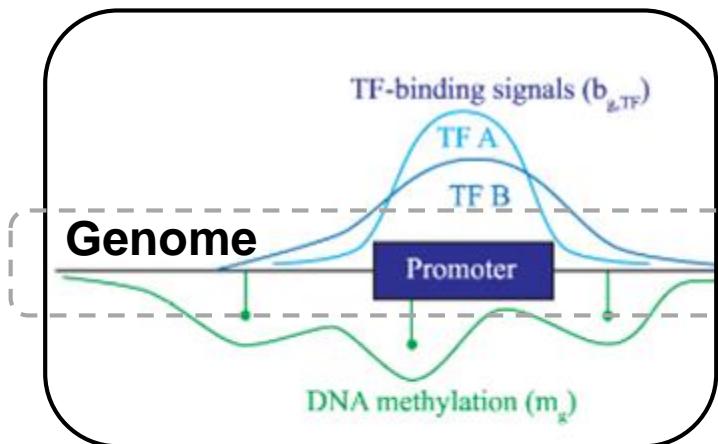


- Thick bars at the bottom are translated regions (exons).
  - There are multiple isoform of BRD2 (i.e. different start locations and exon structures)

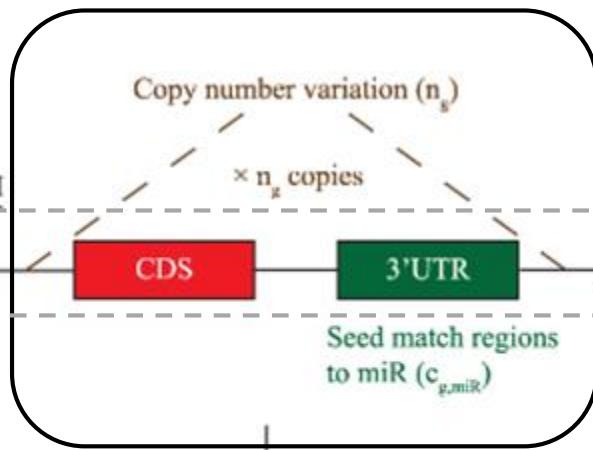
# **Application X: Regulation Of Gene Expression**

# Regulation Of Gene Expression

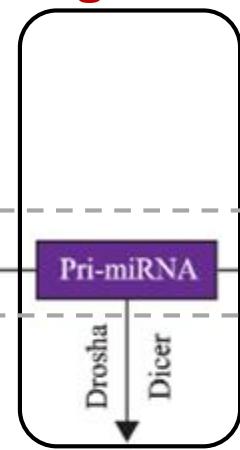
## ON/OFF Switch



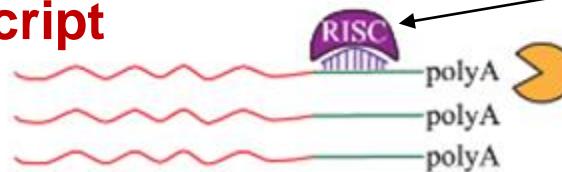
## Gene



## Long-Range Regulator

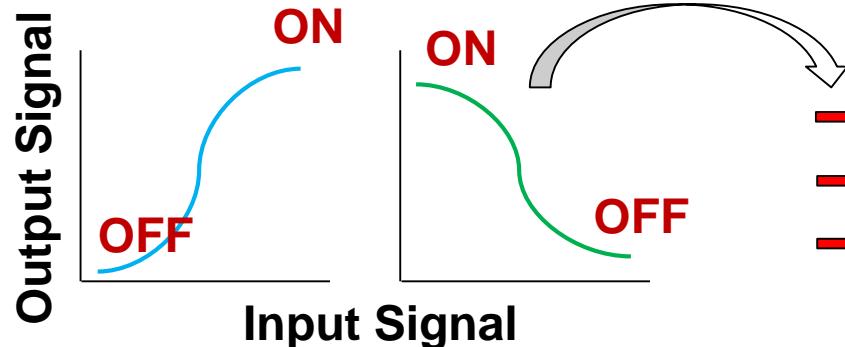


## RNA Transcript

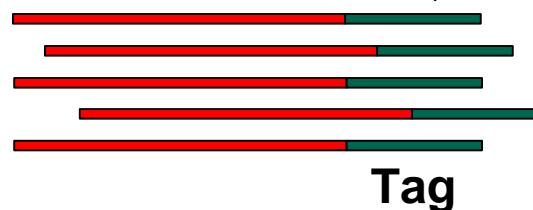


miRNA expression ( $z_{miR}$ )

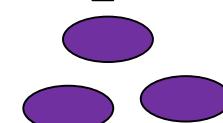
mRNA expression ( $y_g$ )



## Produce

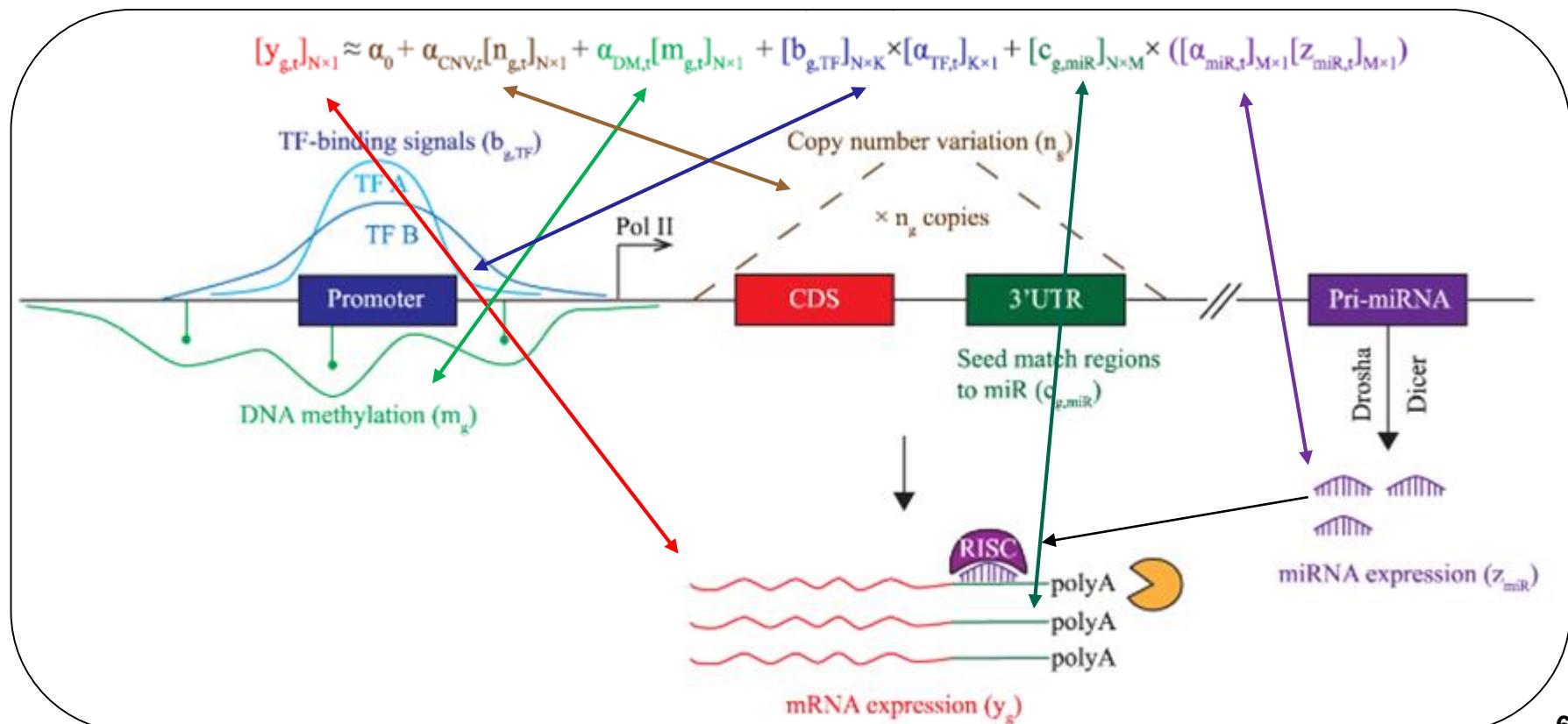
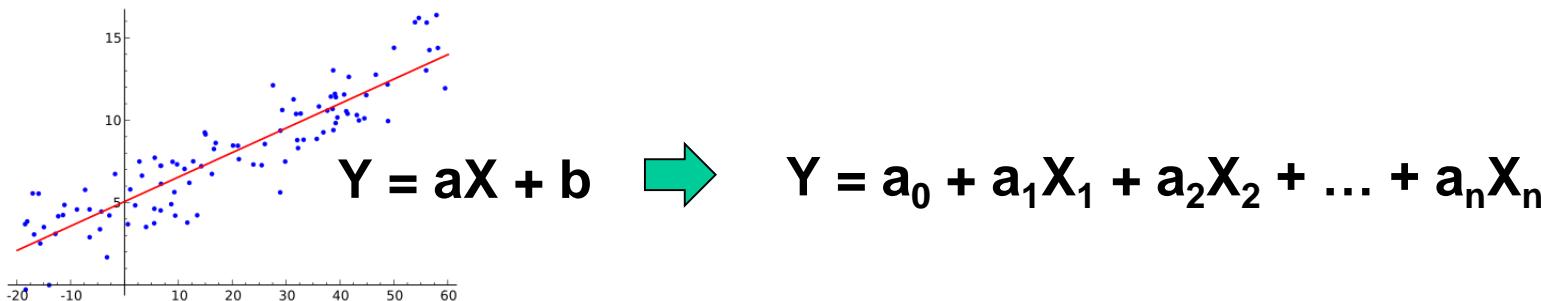


## Destroy



Tag

# Regression Model For Gene Expression I



# Regression Model For Gene Expression II

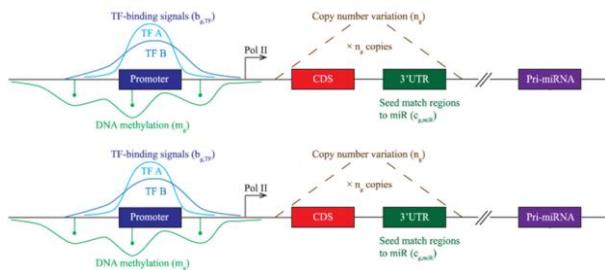
Sample #M

:

Sample #2

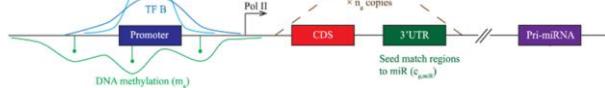
Sample #1

Gene #1:



$$Y_1 = a_{01} + a_{11}X_{11} + a_{21}X_{21} + \dots + a_{n1}X_{n1}$$

Gene #2:



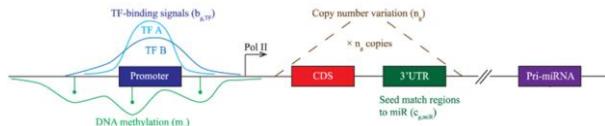
$$Y_2 = a_{02} + a_{12}X_{12} + a_{22}X_{22} + \dots + a_{n2}X_{n2}$$

Adapted from Li et al. PLoS Comp Biol 10, e1003908 (2014)

:

:

Gene #N:



$$Y_N = a_{0N} + a_{1N}X_{1N} + a_{2N}X_{2N} + \dots + a_{nN}X_{nN}$$

Same coefficients  $a_{i,j}$  across samples

# Regression Model For Gene Expression III

## Condition A

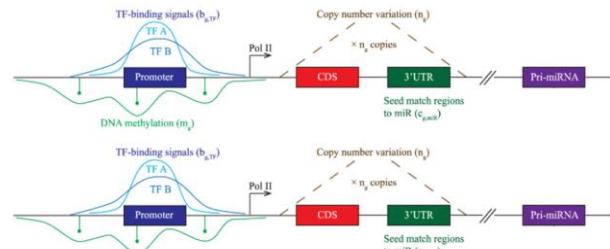
Sample  $A_M$

:

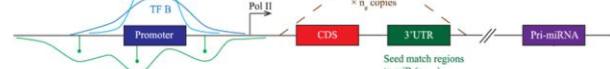
Sample  $A_2$

Sample  $A_1$

Gene #1:

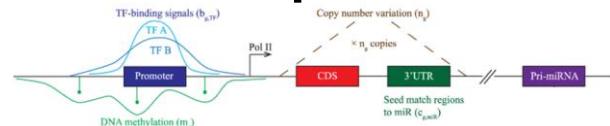


Gene #2:



Adapted from Li et al. PLoS Comp Biol 10, e1003908 (2014)

Gene #N:



## Condition B

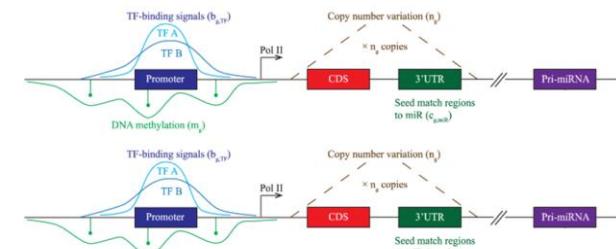
Sample  $B_P$

:

Sample  $B_2$

Sample  $B_1$

Gene #1:

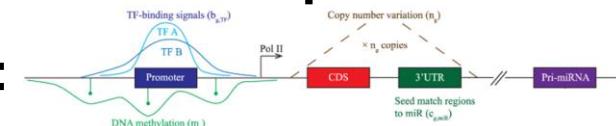


Gene #2:



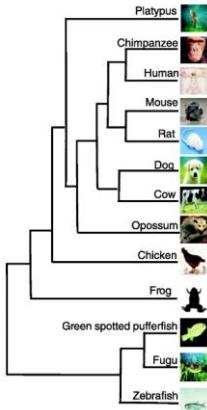
Adapted from Li et al. PLoS Comp Biol 10, e1003908 (2014)

Gene #N:



Same coefficients  $a_{i,j}$  across samples, different across conditions

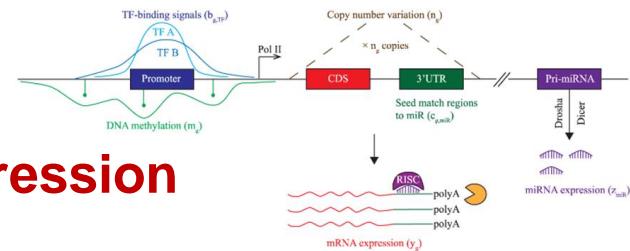
# Summary



# Generative Model

# Maximum Likelihood Estimator

Image from <http://physrev.physiology.org/content/89/3/921>



# Regression

Li et al. PLoS Comp Biol 10, e1003908 (2014)

# Dynamic Bayesian Network

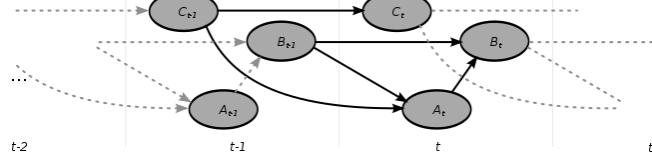
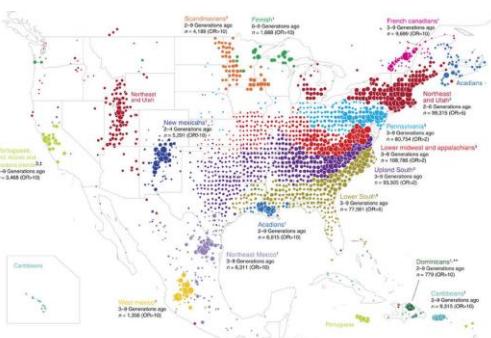
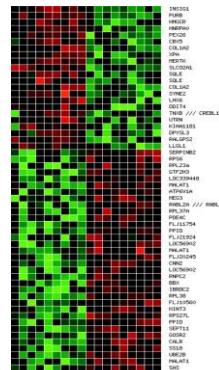


Image from [https://en.wikipedia.org/wiki/Dynamic\\_Bayesian\\_network](https://en.wikipedia.org/wiki/Dynamic_Bayesian_network)

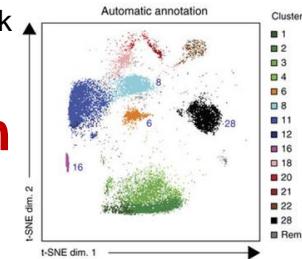
# Clustering



Han et al. Nature Communication 8, 14238 (2017)

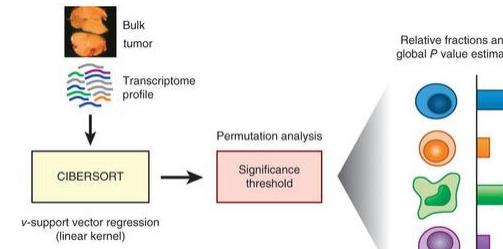


Klings et al. Physiological Genomics 21, 293-298 (2005)



Becher et al. Nature Immunology 15, 1181-1189 (2014)

# Support Vector Machine



Newman et al. Nature Methods 12, 453-457 (2015)