

# Introduction to Bioinformatics

Winter 2024

- Lecturer: Dr. Dvir Aran

- Teaching Assistance:

**Almog Angel**

**Ziv Cohen**

- Course web site:

<http://webcourse.cs.technion.ac.il/236523>

# gal iznko

סטודנט למדעי המחשב  
ומדעי החיים, עם התמחות  
בביו-אינפורמטיקה  
אוניברסיטת תל-אביב



# Introduction – Dr. Dvir Aran



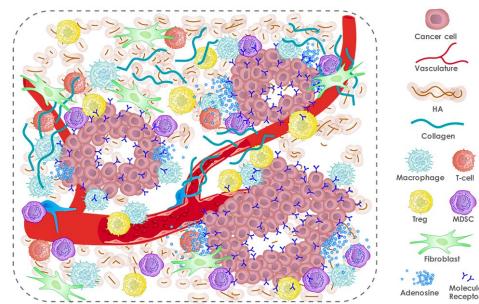
Computer Science  
Computational Biology  
Mathematics  
Faculty of Medicine



Asaf Hellman   Shai Shalev-Shwartz  
Epigenetics   Machine learning



Atul Butte  
Translational Bioinformatics  
Clinical Informatics



Data scientist  
Clinical Informatics  
Digital Medicine



Biomedical Data Science  
Aran Lab @ Technion

# Course Information

- Lectures

Wednesday, 10:30 – 12:30

- Recitations: Almog Angel

  - Wednesday, 14:30-15:30

  - Thursday, 13:30-14:30

- Course Website

<https://webcourse.cs.technion.ac.il/236523/>

# Course Structure and Requirements

- Class Structure
  - 2 hours lecture
  - 1 hour tutorial
- Home-work
  - Homework assignments will be given every second week.
  - The homework will be done in pairs.
  - 4/4 homework assignments must be submitted.
  - A final project will be conducted in pairs.
- Final projects will be analyzing genomic data and presenting your research and results.

# Course Regulations

- Staff
  - Dr. Dvir Aran
  - TA: Almog Angel
  - HW: Ziv Cohen
- Reception hour
  - Wednesday, 12:30
  - By appointment
- Please avoid emails regarding the course not through representatives.

# Grading

- 40% homework assignments (10 points each)
- 60% final project

# Course Objectives

- Learn underlying ideas of common algorithms in bioinformatics.
- Learn to translate a biological problem into a computational problem.
- Learn to read scientific papers, propose and conduct independent research.

As part of the course we will introduce the R programming language and learn to analyze DNA and RNA sequencing data and conduct basic statistical and machine-learning analyses.



# Class structure

- Lectures:
  1. Problem in biology.
  2. Computer science concept.
  3. Applying the computer science concept to solve the biology problem.
- Recitations:
  - Hands-on analysis of biological data
  - R, statistics, algorithms.

# Course subjects

1. Comparative genomics (2 classes)
2. Gene expression (3 classes)
3. Genetics (2 class)
4. Frontiers in bioinformatics:
  - Cancer genomics
  - Single-cell genomics
  - Genome engineering



Dynamic programming, linear algorithms hash functions, clustering algorithms, differential expression analysis, deconvolution, burrows-wheeler transformation, machine learning algorithms, energy functions optimization, deep learning, statistical hypothesis testing, dimensionality reduction, population structure algorithms, genome wide association studies, polygenic risk scores and more...

# Aligning expectations

Evolution

Sequencing technologies

CRISPR

immunology

Cancer

Research



Statistics

Machine learning

Genetic diseases

Gene expression

Algorithms

Molecular biology

What is Computational Biology/Bioinformatics?

# What is Computational Biology/Bioinformatics?

**Computational biology** and **bioinformatics** is an interdisciplinary field that develops and applies **computational methods** to analyze large collections of biological data, such as genetic sequences, cell populations or protein samples, to make new predictions or **discover new biology**.

<https://www.nature.com/subjects/computational-biology-and-bioinformatics>

# Why Study Computational Biology?

Interdisciplinary

Biology

Computer Science

Data Science

Mathematics

Statistics

## Best Jobs

1. Data scientist
2. Statistician
3. University Professor
4. Occupational Therapist
5. Genetic Counselor
6. Medical Services Manager
7. Information Security Analyst
8. Mathematician
9. Operations Research Analyst
10. Actuary
11. Software Developer

## Worst Jobs

224. Taxi Driver
223. Logging Worker
222. Newspaper Reporter
221. Retail Salesperson
220. Enlisted Military Personnel
219. Correctional Officer
218. Disc Jockey
217. Nuclear Decontamination Technician
216. Advertising Salesperson
215. Broadcaster

### Data Scientist

Overall Rating: 1/200

Median Salary: \$114,520

Work Environment

Stress

Projected Growth

Very Good

Low

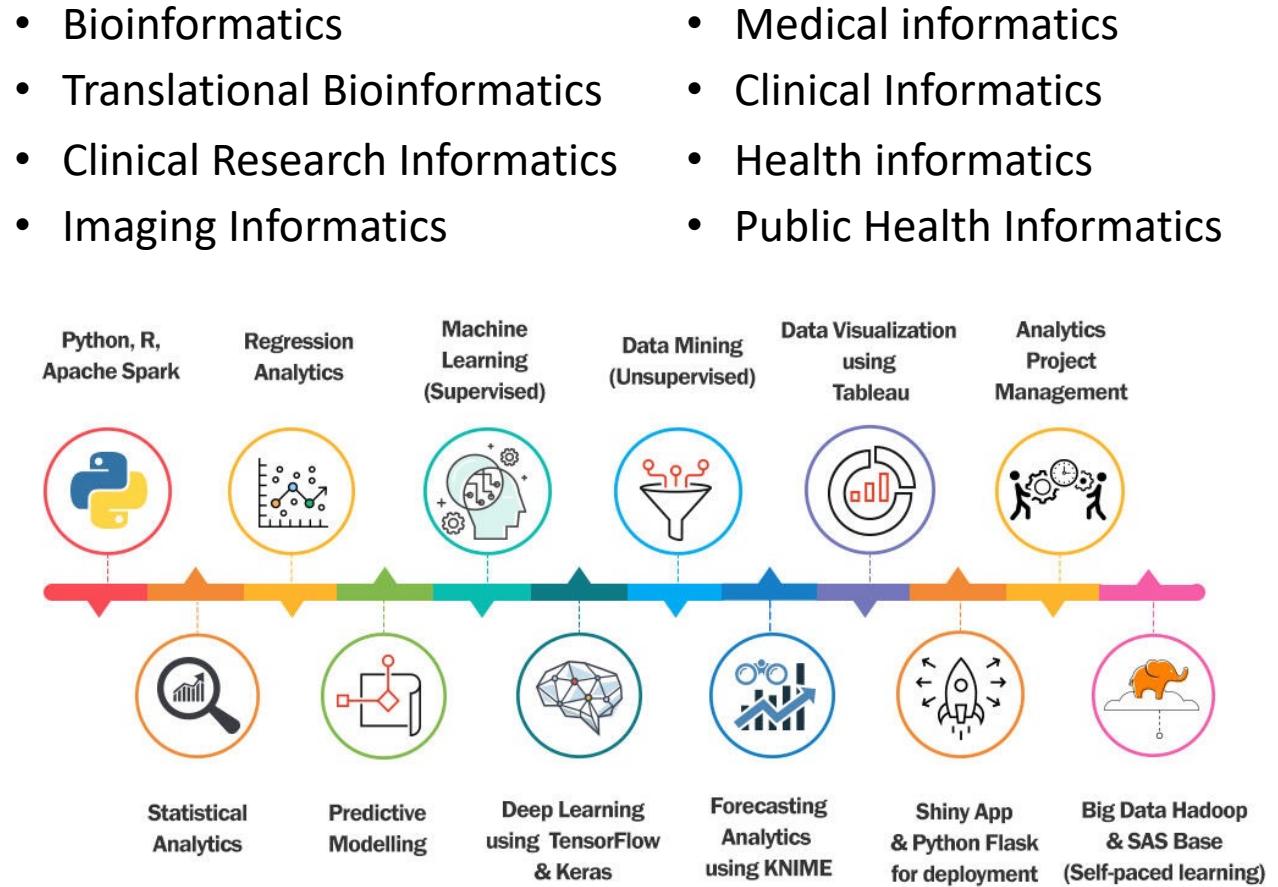
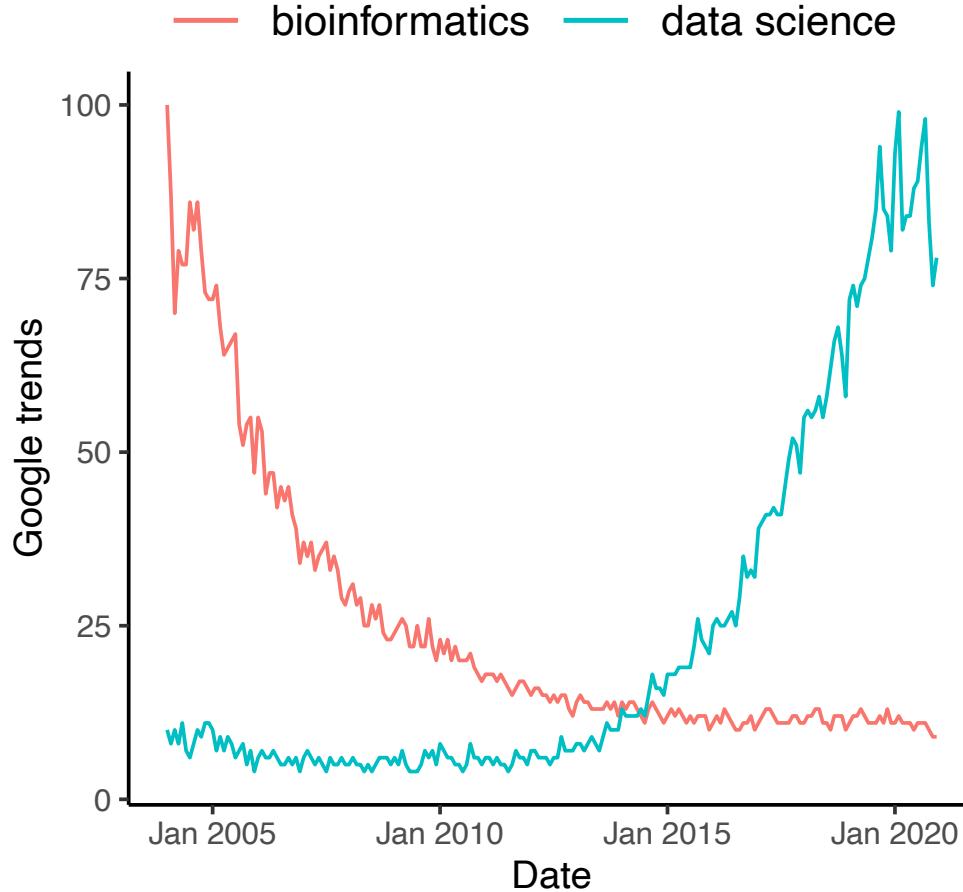
Very Good

4/200

42/200

37/200

# Bioinformatics or Biomedical Data Science?

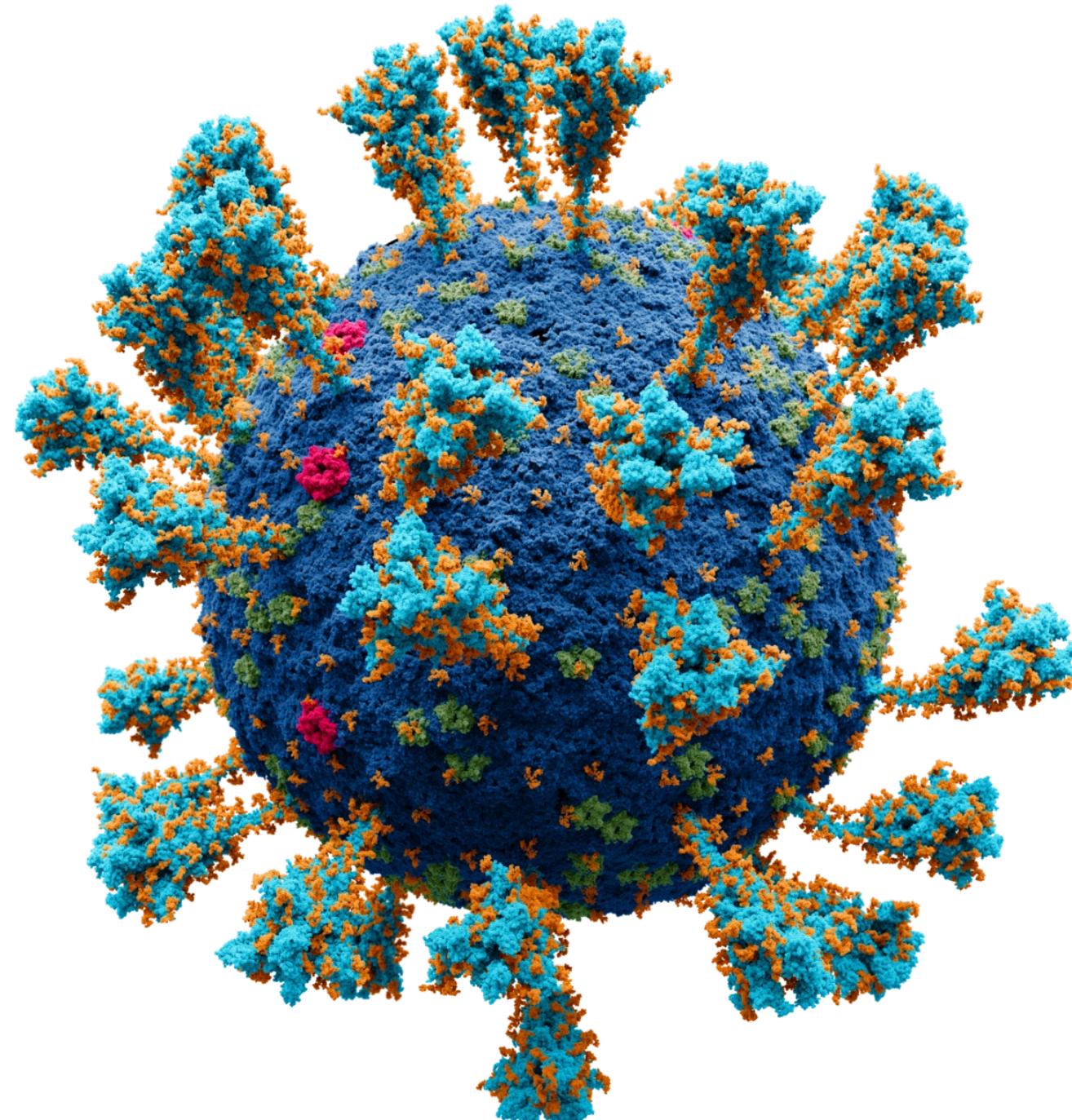


This image contains a single line of binary code, which is extremely long and consists of approximately 999,999,999 bits (approximately 124,999,999 bytes or about 124.999999 MB). The binary sequence is composed of a series of 1s and 0s, with no visible text, symbols, or other markings.

# SARS-CoV-2

**S - Spike glycoprotein precursor**

**Spike Glycoprotein**

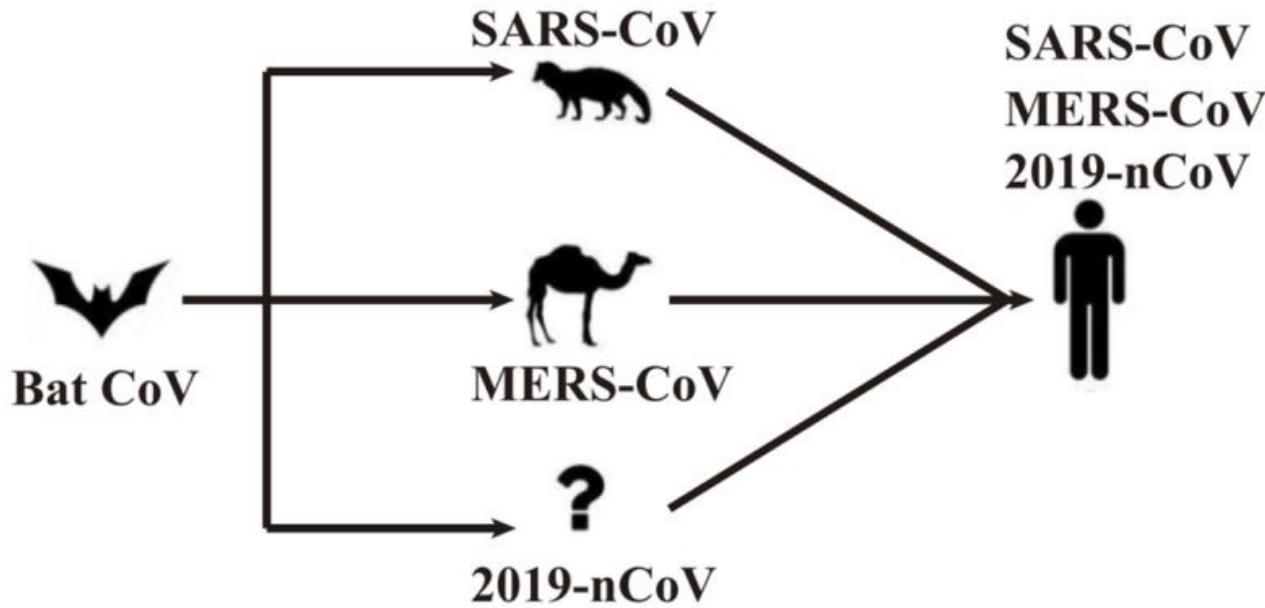
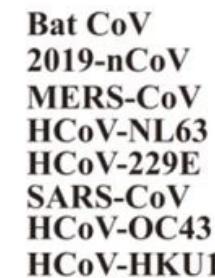
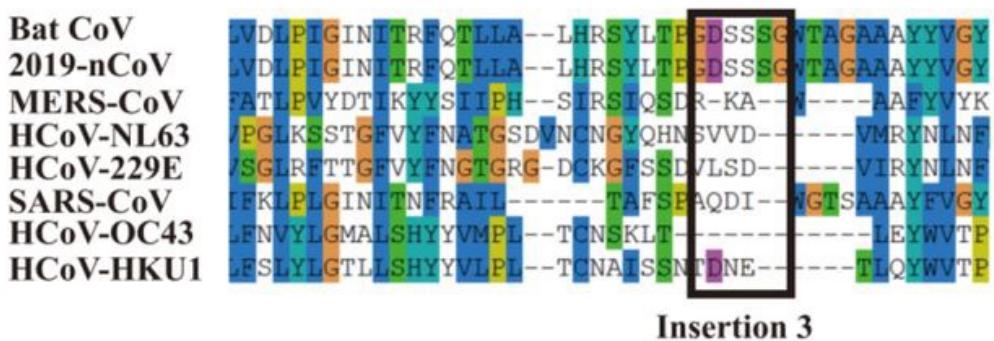
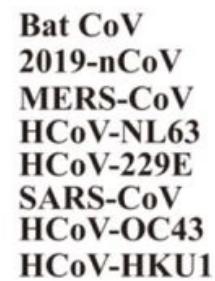
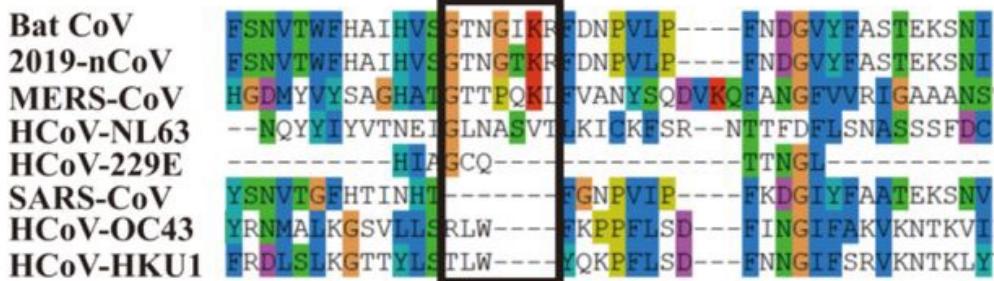


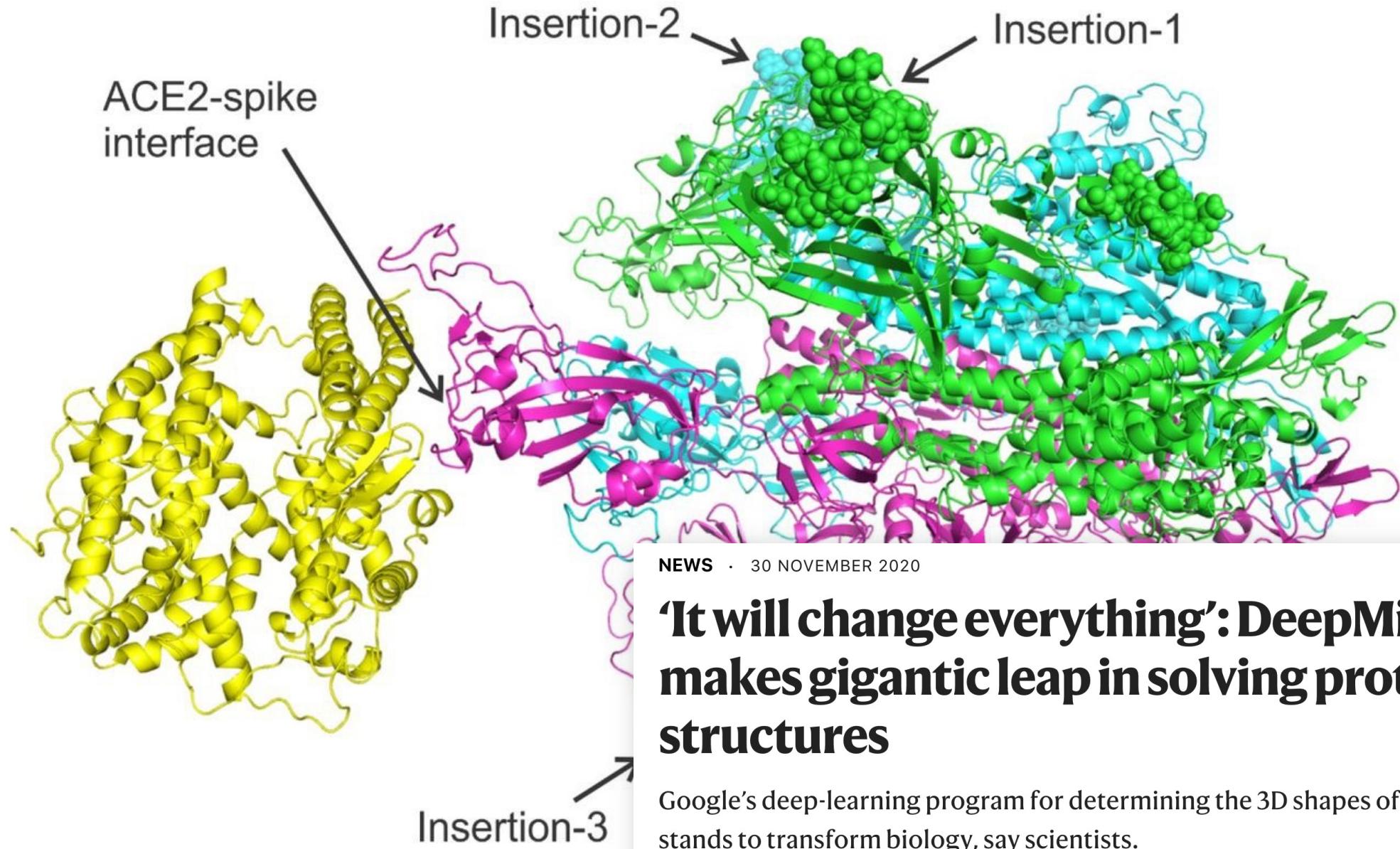
CUUUAAAAGAAGGUAAAUCAAUGAUUAUCUUUCUAGUAAAGGUAGACUUAAAAGAGAAAACAACAGAGUUGUUUUUCUAGGUGAUGUUCUUGUU  
AACACUAAACGAACA**AUG**UUUGUUUUUCUUGUUUAUUGCACUAGCUCUAGCAGUGUUAAUCUACACCAGAACUAAUACCCCCUGCAUACACUAAU  
CUUUCACACGUGGUUUUUACCCUGACAAAGUUUUACAGAUCCUCAGUUUACAUCAACUCAGGACUUGUUCUACCUUUCUUUCCAUGUUACUUGGUUCCAU  
GCUAUACAUGUCUCUGGGACCAAUGGUACUAAGAGGUUJGAUAAACCCUGCUACCAUUAAAUGAUGGUGUUUUUGCUUCCACUGAGAACGUCAACAUAAAG  
AGGCUGGAUUUUUGGUACUACUUAGAUUCGAAGACCCAGUCCUACUUUUGUUAACCGCUACUAUGUUGUAAAAGUCUGUGAAUUCAUUUUGUAAUG  
AUCCAUUUUUGGUUUUUACCCACAAAAACACAAAAGUUGGAUGGAAAGUGAGUUCAGAGUUUUCUAGUGCGAAUAAJUGCACUUUUGAAUAGUCUCAG  
CCUUUCUUAUGGACCUUGAAGGAAAACAGGGUAAAUCUAGGGAAUUUGGUUAAAAGAAUAAUGAUGGUAAAAGUAAUUCUAAGCACACGCC  
UAUUAUUUAGUGCGUGAUCUCCCUCAGGGUUUUUCGGCUUAGAACCAUUGGUAGAUUUGCCAAUAGGUAAAACAUCACUAGGUUCAACUUACUUGCUUAC  
AUAGAAGUUUUUGACUCCUGGUGAUUCUUCAGGUJUGGACAGCUGGUUGCGAGCUUAAAUGUGGUUAUCUCAACCUAGGACUUUCUAAUAAAUAU  
GAAAUGGAACCAUUACAGAUGCUGAGACUGACUUGACCCUCUCAGAAACAAAGUGUACGUUGAAAUCUUCACUGUAGAAAAGGAAUCUAACACUUC  
UAACUUUAGAGUCCAACCAACAGAAUCUAUUGUUAGAUUUCCUAAAUAACAAACUUGUGCCUUUUGGUAGGUUUAAAACGCCACCAUGCAUCUGUUUAUG  
CUUGGAACAGGAAGAGAAUCAGAACUGUGUUGCUGAUUAUCUGGUCAUAAAUCGCAUCAUUUCCACUUUAGGUUAUGGAGUGUCUCCUACAAAUA  
AAUGAUCUCUGCUUACUAAUGCUAUGCAGAUUJGUAAAAGAGGUGAAGUGACAGACAAUCGCUCCAGGGCAAACUGGUAGAUUGCUGAUUAU  
AAAAUUACCAGAUGUUUACAGGCUGGUUAAGCUUACAAUCUUGAUUCUAAAGGUJUGGUAAAUAUACUGUUAAGAUUGUUUAGGAAGU  
CUAAUCUAAACCUUUUGAGAGAUUUCAACUGAAAUCUACAGGCCGUAGCACACCUUGUAAUGGUUGAAGGUUUAAAUGUUACUUUCCUUACAAUCA  
UAUGGUUUCCAACCCACUAAUGGUUGGUACCAACAGAGUAGUACUUUCUUUGAACUUCUACUGCACCAGCAACUGUUJUGGUAGGUAAAAGUC  
UACUAAUUUGGUAAAACAAAGUGUCAUUUCAACUCAUUGGUUACAGGCACAGGUUCUACUGAGUACAAAAGUUUCGCCUUUCCAACAAUUG  
GCAGAGACAUUGCUGACACUACUGAUGCUGGUCCACAGACACUUGAGAUUCUUGACAUUACCCAUGUUCUUUUGGUUGGUAGUUAACACCAGGA  
ACAAAACUUCUAAACCAGGUUGCUGUUCUUUACAGGAUGUUACUGCACAGAAGGUCCUGUUGCUAUCAUGCAGAUCAACUUACCUACUUGGCUGUUUAUC  
UACAGGUUCUAAUGUUUUCAAACACUGCAGGCUGGUAAAAGGGGUAGUACAACUCAUAUGAGUGUACAUACCCAUUGGUAGGUUAUGCGCUA  
GUUAUCAGACUACAGUAAUUCUCCUCGGCGGACGUAGGUAGCUACAUUCCACUACUGUACUUGGUAGGUAAAUCAGUUGCUUACU  
AAUACUCUAAUGCCAUACCCACAAAUUUACUAAUAGGUUACAGACAGGUUACAGUGUCAUGACCAAGACAUCAUGGUAGAUUGUACAAUUGU  
UGAUUCAACUGAAUGCAGCAUCUUUGUUGCAAUAAUGGCAGGUUUUGUACACAAUAAAACCGUGCUUACUGGUAGUACAAGACAAAACACCAAG  
AAGUUUUUGCACAAGUCAACAAUAAAACACCACCAUUAAAAGAUUUUGGUUUUAAAUCACAAUAAAACCAGAACUCAAAACCAAGCAAGAGG  
UCAUUUAUUGAAGAUCUACUUUCAACAAAGUGACACUUGCAGAUGCUGGUUCAUCAACAAUAAAUGGUUACAGGUUGGUUACAGGUUG  
UGCACAAAAGUUUACGCCUUACUGUUUUGCCACUUGCUCACAGAUGAAAUGAUUGCUCAUACUUCUGCACGUAGCAGGUACAAUCACUUCUGGUUG  
CCUUUUGGUAGGUGCUGCAUUACCAUUUGCUAUGCAGGUUAAUGGUUAAAUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG  
AUUGCCAACCAUUUAAAGUGCUAUUGGCCAAAAGUACACUUCACUUCUCCACAGCAAGUGGCACUUGGUAAAUCUCAAGAUGGUUACCA  
UUUAAACACGCUUGUAAAACAUCUAGCUCCAAUUUUGGUCAUUUCAAGGUUUAAAUGAUUACCUUACGUUGACAAAGUUGAGGCUGAAGUG  
AUAGGUUGAUCACAGGCAGACUCAAGGUUGCAGACAUAGUGACUACAAUAAAAGAGGACUGGUAGGUAAAUCAGAGGUUCUGCUAAUCU  
UCAGAGUGGUACUUGGACAUAAAAGAGGUUGAUUUUGGUUAGGGCUAUCAUCUUAUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG  
GACUUAUGGUCCUGCACAAGAAAAGAACUUCACAUCUGGUCCAUUUGUCAUGAUGGUAAAAGCACACUUCUCCUGUGAAGGUUGGUU  
ACUGGUUUGUACACAAAGGAAUUUUAUGAACCAACAGACACACAUUUGUGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG  
GUUUUAUGAUCCUUUGCAACCUGAAUUAAGACUCAUUCAGGAGGUAGAUAAAUAUAAAAGAAUCAUACACAGGUUGGUUACAGGUUG  
UAUAGGUUCAGUUGUAAAACAUCUAAAAGAAAAGGUAGGCCUCAAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG  
AGUUAUAAAAGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG  
CUCAAGGGCUGUUGUUCUUGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUGGUUACAGGUUG



CUUUAAAAGAAGGC **AAAUCAAUGUAUGAUUUUACUCUUCUUAGUAAAGGUAGACUUAAUUAGAGAAAACAACAGAGUUGUUUUUCUAGUGAUGUUUCUUGU**  
**AACAACUAAACGAACA** AUGUUUGUUUUUCUUGUUUUUUGCCACUAGUCUAGCAGUGUUAUCUACACCAGAACUAAUACCCCCUGCAUACACUAAU  
CUUUCACACGUGGUGUUUAUACCUGACAAGUUUUUCAGAUCCUCAGUUUACAUCAACUCAGGACUUGUUCUACCUCUUUCCAUGUUACUUGGUUCCAU  
GCUAUACAUGUCUCUGGGACCAAUGGUACUAAGAGGUUUGAUACCCUGGUCCUACCAUUUAUGAUGGUUUAUUGCUUCCACUGAGAGUCUAACAUAAAG  
AGGCUGGAUUUUUGGUACUACUUUAGAUUCGAAGACCCAGUCCUACUUUUGUUAACCGUACUAAGUUGUUAUUAAGUCUGUGAUUUCAAUUUGUAAUG  
AUCCAUUUUUGGUUUUAUACCACAAAAACAACAAAAGUUGGAUGGAAAGUGAGUUCAGAGUUUAUCUAGUGCAGAUAAUUGCACUUUUGAAUAGUCUCAG  
CCUUUCUUUAUGGACCUUGAAGGAAAACAGGGUAAUUUCAAAAAUCUAGGGAAUUUGGUUUAAGAAUAUUGAUGGUUUAUUAACAGCACACGCC  
UAUUAUUUAUGCGUGAUCUCCCUCAGGGUUUUUCGGCUUUAAGACCAUUGGUAGAUUUGCCAAUAGGUUAUACACUAGGUUCAACCUUACUUGCUUAC  
AUAGAAGUUUUUGACUCCUGGUGAUUCUUUCAGGUUGGACAGCUGGUUGCCAGCUUAAAUGGGGUUAUCUCAACCUAGGACUUUCUAAUAAAUAU  
GAAAUGGAACCAUUACAGAUGCUGUAGACUGACCCUCUCAGAAACAGUGUACGUUGAAUCCUUCAGUAGAAAAGGAAUCUAACCUUC  
UAACUUUAGAGGUCAACCAACAGAAUCUAUUGUUAGAUU **CCUAAUAAUACAAACUUGUGCCUUUUGGUGAAGUUUUUACGCCACCAAGAUUUGCAUCUGUUUAUG**  
**CUUGGAACAGGAAGAGAAUCAGCAACUGUGUUGCUGAUUAUUCUGUCCUAAUAAUUCGCAUCAUUUCCACUUUUAAGGUUAUGGAGUGUCUCCUACUAAUUA**  
**AAUGAUCUCUGCUUUAACUAAUGCUAUGCAGAUUCAUUGUAAUAGAGGUGAUGAAGUCAGACAAUCGCUCCAGGGCAAACUGGAAAGAUUGCUGAUUAUAAUUA**  
**AAAAUUACCAGAUGAUUUUACAGGCUGCGUUUAAGCUUAGGAAUUCUAAACAUUCAAGGUUGGUGGUAAUUAUACUGUUAUAGUUUUAGGAAGU**  
**CUAAUCUAAACCUUUUUGAGAGAGAUUUCAACUGAAAUCUAAUCAGGCCGUAGCACACCUUUGUAUUGGUUAGGUUUUAUUGUUAUACUUCUUACAAUCA**  
**UAUGGUUCCAACCCACUAUUGGUUGGUUACCAACCAUACAGAGUAGUACUUUCUUJUGAACUUCUACUGCACCGAGCAACUGUUUGGGACCUUAAAAGUC**  
**UACUAAUUUGGUAAAAACAAAUGGUUCAACUUCUAAUGGUUACAGGCACAGGUUUCUACUGAGUACAAACAAAGUUUCUGCUUUCCAACAAUUG**  
**GCAGAGACAUUGCUGACACUACUGAUGCUGUCCGUGAUCCACAGACACUUGAGAUUCUUGACAUUACACCAUGUUCUUUUGGGUGGUAGGUUAACACCAGGA**  
ACAAAACUUCUAAACCAGGUUGCUGUUCUUUAUCAGGAUGUUAACUGCACAGAAGUCCUGUUGCUAUCAUGCAGAUCAACUUACUCCUACUUGGCUGUUUAUUC  
UACAGGUUCUAAUGUUUUUCAAACACGUGCAGGCUGUUUAUAGGGCUGAACAGUCAACAUCUAAUAGAGUGUGACAUACCAUUGGUGCAGGUUAUGCGCUA  
GUUAUCAGACUACUAAUUCUCCUCGGCGGGACGUAGUGUAGCUAGUCAUCCAUUGGUACACUACUAAUAGUGUGACAUACCAUUGGUGCAGGUUAUGCGCUA  
AAUACUCUAAUUGCCAUACCCACAAAAUUUACUAAUAGGUUACCAAGAAAUCUACAGUGUACUAGUAGACAUACAGUAGAUUGUACAAUUGUACAUUUGUGG  
UGAUUCAACUGAAUGCAGCAUCUUUGUUGCAAAUAGGCAGUUUUGUACACAAUUAACCGUGCUUUAACUGGAAUAGCUGUUGACAAAGACAAAACACCAAG  
AAGUUUUUGCACAAGUCAACAAACUAAACACCACCAUUUAAGAUUUUGGGUUUUUUUUUACAAAUAUACAGAACUCAACAGCAAGGAG  
UCAUUUAUUGAAGAUCUACUUUCAACAAAGUGACACUUGCAGAUGCUGGUUCAUCAACAAACAAUAGGUUACAGAACUCAACAGCAAGG  
UGCACAAAAGUUUACGCCUUACUGUUUUGCCACCUUUGCUCACAGAUGAAUUGCUAAUACACUUCUGCACGUUAGCAGGGUACAAUCACUUCUGGUUGGA  
CCUUUUGGUGCAGGUGCUGCAUUACAAUACUUGCUAAUAGGUUUAUGGUUACAGAACUCAACAGCAAGGAG  
AUUGCCAACCAUUUAUAGGUUACUUGGCCAAUUCAGACUACUUCUCCACAGCAAGUGCAGUCAUUGGAAAACUCAAGAUGGGUACCAACAAAGCACAAG  
UUUAAACACGCUUGUUAAACAACUUAAGCUCCAAUUUUGGUCAUUUCAAGGUUUUAAAUGAUUACCUUACGUUACAGUAGGAGGUACAGUUG  
AUAGGUUGAUACAGGCAGACUACAAAGUUUGCAGACAAUAGUGACUACAAUUAUAGAGGCUGCAGAAAUCAGAGCUUCUGCUAAUCUUCUGCUACUAAAUG  
UCAGAGUGUGUACUUGGACAAUCAAAAAGAGUUGAUUUUGGGAAAGGGCUAUCAUCUUAUGGUUACAGAACUCAACAGCAAGGAG  
GACUUAUGGUCCUGCACAAGAAAAGAACUUCACAACUGCUCCUGCCAUUUGUCAUGAUGGAAAAGCACACUUCUCCUGUGAAGGUGUUCUUGUUCAAAUGGCACAC  
ACUGGUUUGUAACAAAGGAAUUUUUUAUGAACCAACAGACAACACAAUUGUGUCUGGUACAGUGUUGUUAUAGGAAUUGUCAACACACA  
GUUUUAUGAUCCUUUGCAACCUGAAUUAGACUACAUUCAAGGAGGUAGAUAAAUAUUAAGAAUCAUACACCAGAUGUUGUUAUAGGAGACAUUCUGGU  
UAAUGCUUCAGUUGUAAACAUUCAAAAAGAAAAGGUACCGCCUCAAUGAGGUUGCCAAGAAUUAAGAAUUCUCAUCGAUCUCCAAAGAACUUGGAAAGUAUGAGC  
AGUUAUAAAAGGGCAUGGUACAUUUGGUAGGUUGGUUAGCUGGUUGGUUAGGUUAGGUUGGUACAGUUGCUGUAGUUGU  
CUCAAGGGCUGUUGUUCUUGGUAGGUCCUGCUGCAAUUUUGAUGAAGACGACUUCUGAGCCAGUGCUAAAGGGAGUCAAUUACACAC **UAA** ACGAACUUUAGGU

MFVFLVLLPLVSSQCVNLTRTQLPPAYTNSFRGVYYPDKVFRSSVLHSTQDLFLPFFSNVTWFHAIHSGTNGTKRFDNP  
VLPFNDGVYFASTEKSNIIRGWIFGTTLDSKTQSLLIVNNATNVVIKVCEFQFCNDPFLGVYYHKNNSWMESEFRVYSSA  
NNCTFEYVSQPFLMDLEGKQGNFKNLREFVFKNIDGYFKIYSKHTPINLVRDLPQGFSALEPLVDLPIGINITRFQTLLALHRS  
YLTPGDSSSGWTAGAAAYVGYLQPRTFLLKYNENGTTDAVDCALDPLSETKCTLKSFTVEKGIYQTSNFRVQPTESIVRFP  
NITNLCPFGEVFNATRFASVYAWNRKRISNCVADYSVLYNSASFSTFKCYGVSPKLNDLCFTNVYADSFVIRGDEVRQIAPG  
QTGKIADYNYKLPDDFTGCVIAWNSNNLDSKVGGNYYRLFRKSNLKPFERDISTEIYQAGSTPCNGVEGFNCYFPLQSY  
GFQPTNGVGYQPYRVVVLSELHAPATVCGPKNSTNLVKNKCVNFNFNGLTGTGVLTESNKKFLPFQQFGRDIADTTDAV  
RDPQTLEILDITPCSFGGVSVITPGNTSNQAVLYQDVNCTEVPAIHADQLPTWRVYSTGSNVFQTRAGCLIGAEHVN  
NSYECDIPIGAGICASYQTQTNSPRRARSVASQSIIAYTMSLGAENSVAYSNNSIAIPTNFTISVTTEILPVSMKTSVDCTMYI  
CGDSTECSNLLQYGSFCTQLNRALTGIAVEQDKNTQEVAQVKQIYKTPPIKDFGGFNFSQILPDPSKPSKRSFIEDLLFNKV  
TLADAGFIKQYGDCLGRIAARDLICAQKFNGLTVLPPLLTDEMIAQYTSALLAGTITSGWTFGAGAALQIPFAMQMAYRFN  
GIGVTQNVLYENQKLIANQFNSAIGKIQDSLSSSTASALGKLQDVVNQNAQALNTLVQLSSNFGAISSVLNDILSRLDKVEA  
EVQIDRLITGRLQLTYVTQQLIRAAEIRASANLAATKMSECVLGQSKRVDFCGKGYHLMSFPQSAPHGVVFLHVTYVPA  
QEKNFTTAPAICHDGKAHFREGVFSNGTHWFVTQRNFYEPQIITDNTFVSGNCDVVIGIVNNNTVYDPLQPELDSFKEE  
LDKYFKNHTSPDVDLGDISGINASVVNIQKEIDRLNEVAKNLNESLIDLQELGKYEQYIKWPWYIWLGFIAGLIAIVMVTIML  
CCMTSCSCLKGCCSCGSCCKFDEDDSEPVLKGVKLHYT

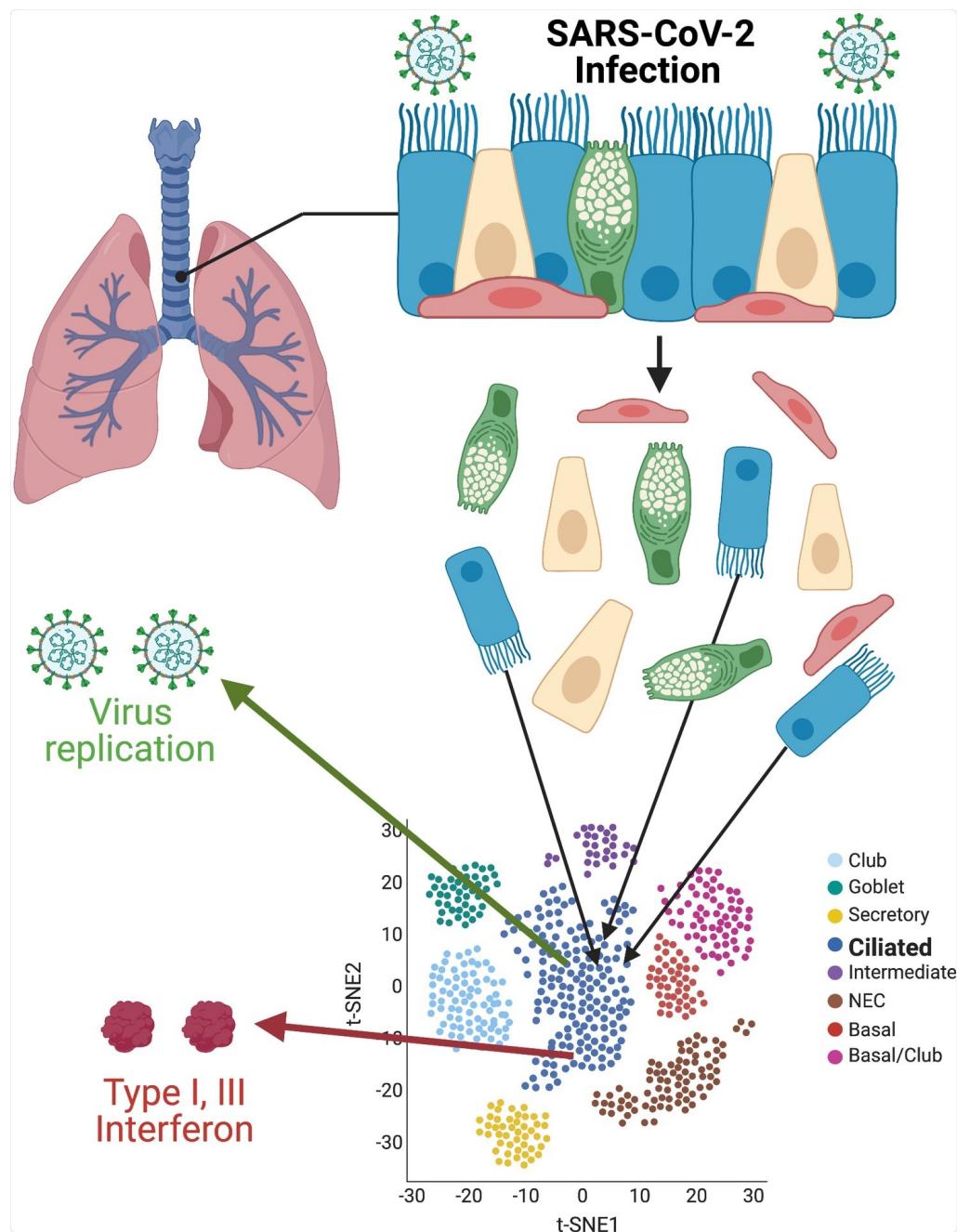
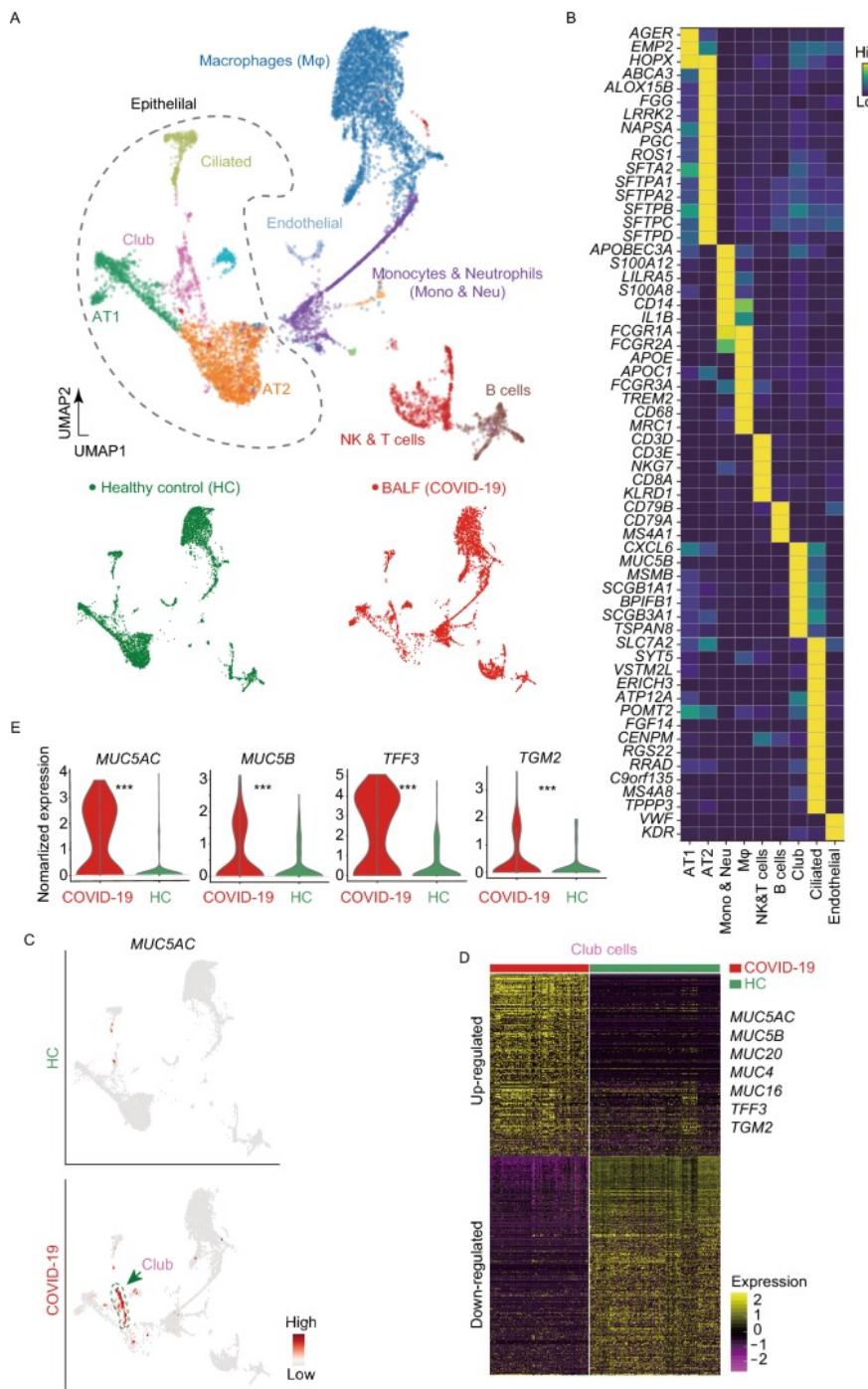
**A****B**



NEWS · 30 NOVEMBER 2020

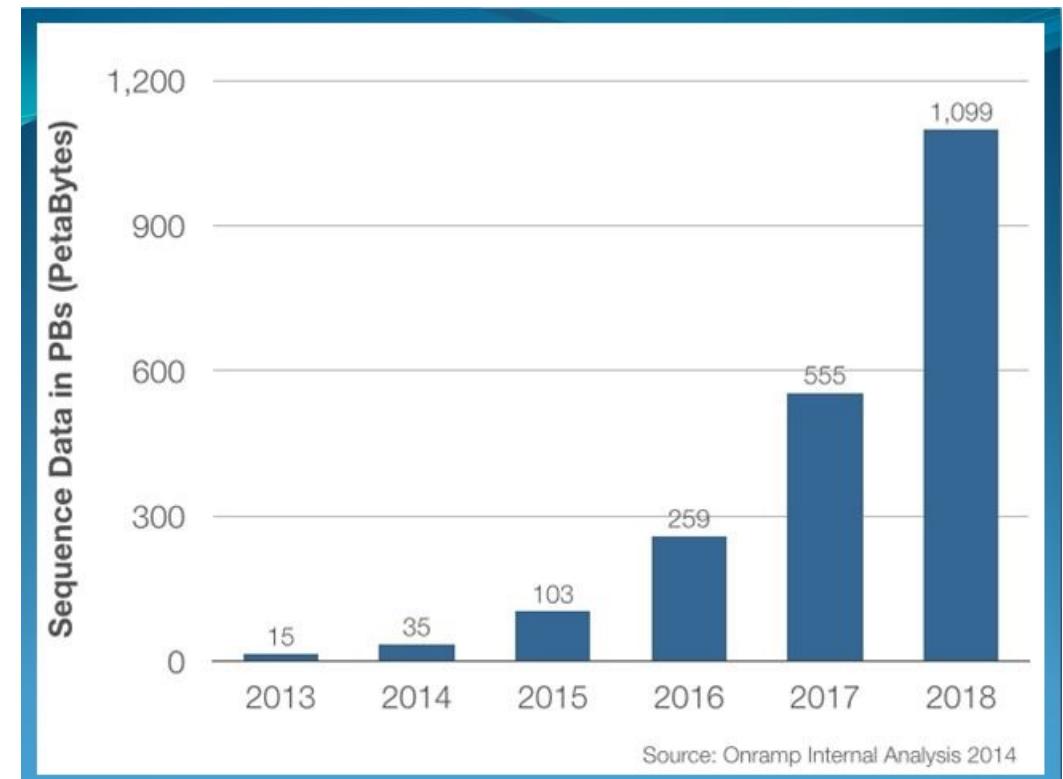
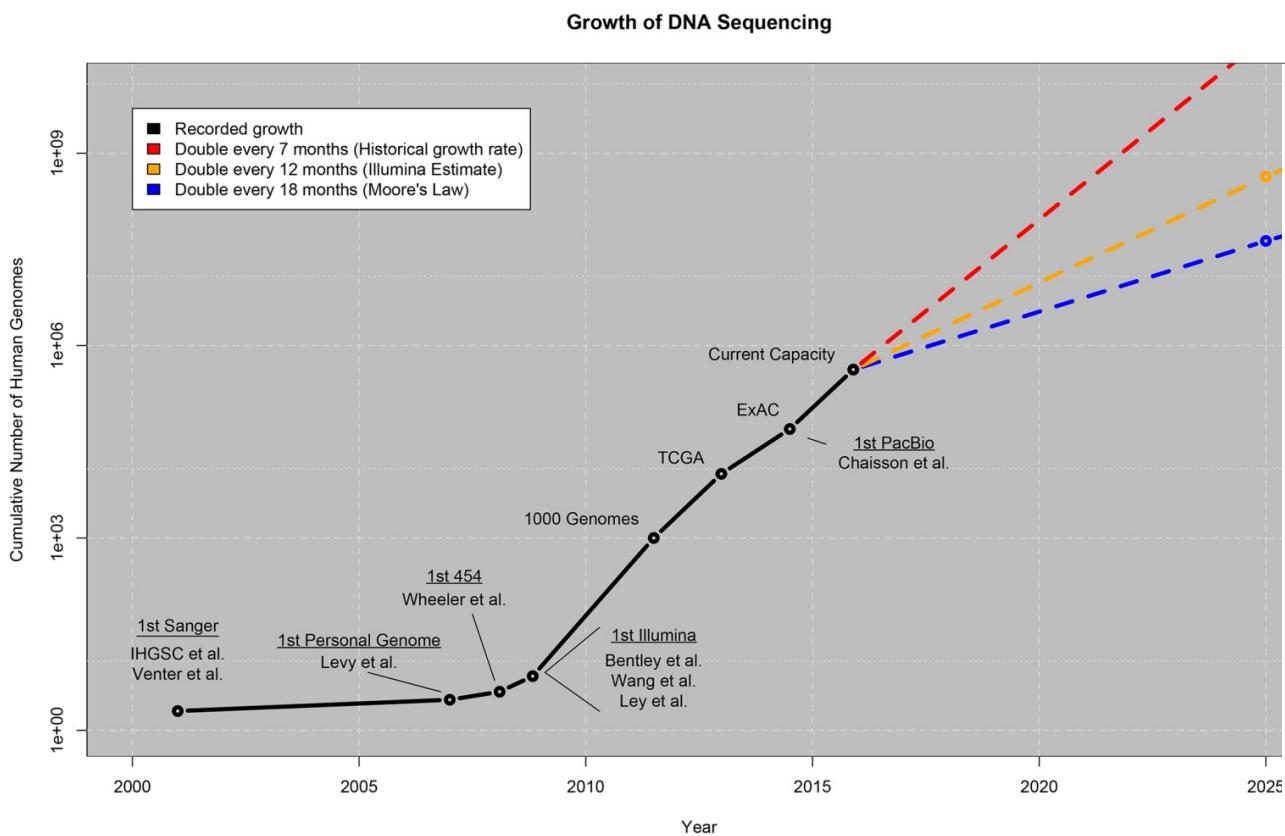
## 'It will change everything': DeepMind's AI makes gigantic leap in solving protein structures

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

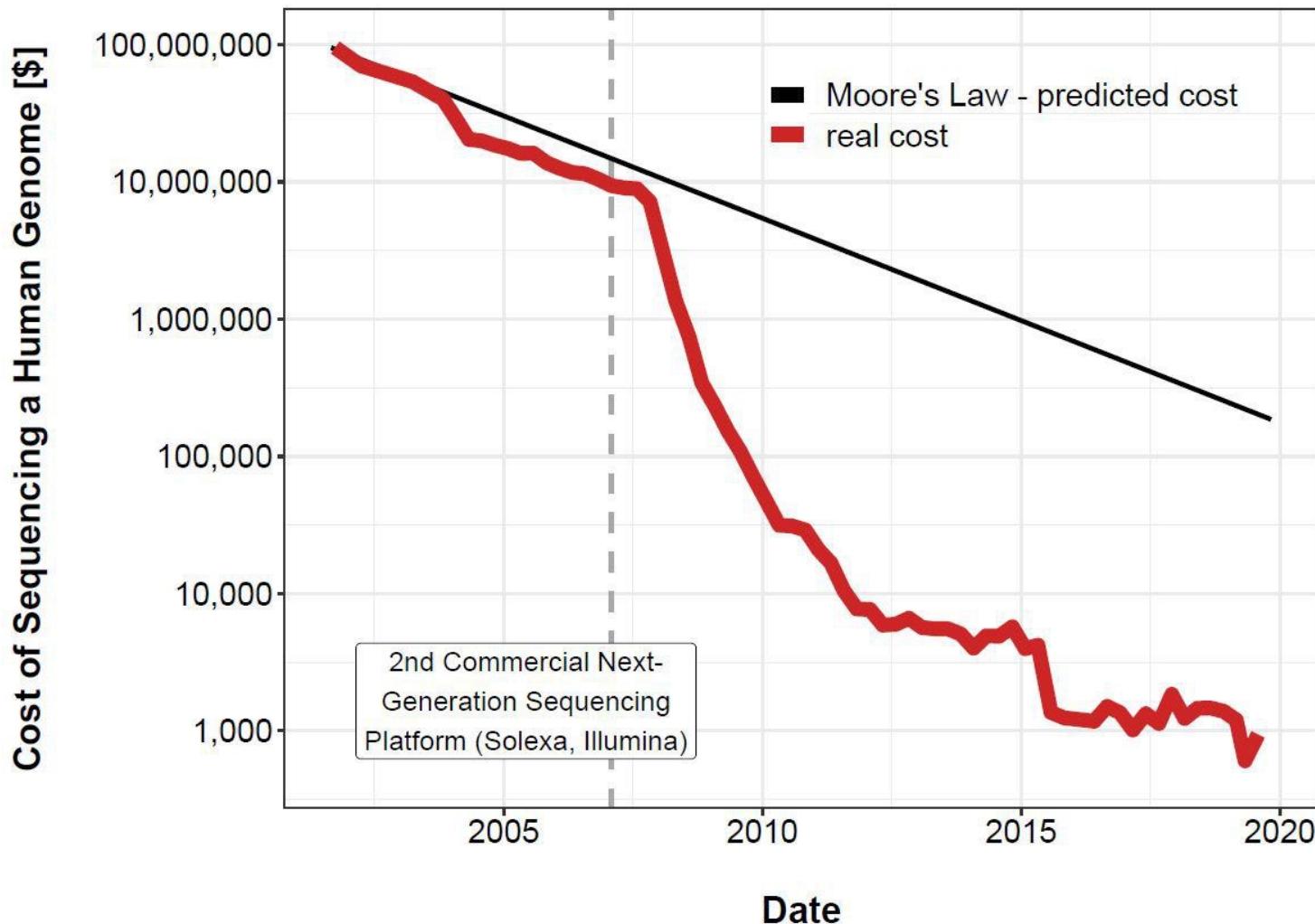




# Information explosion in biology



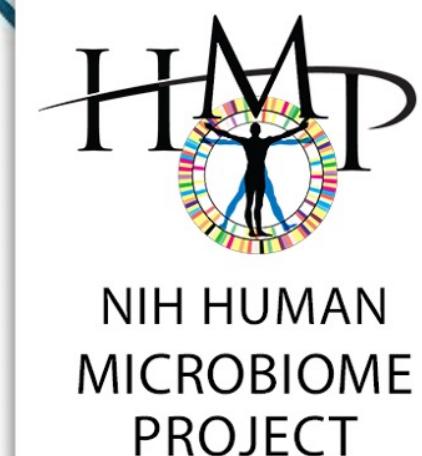
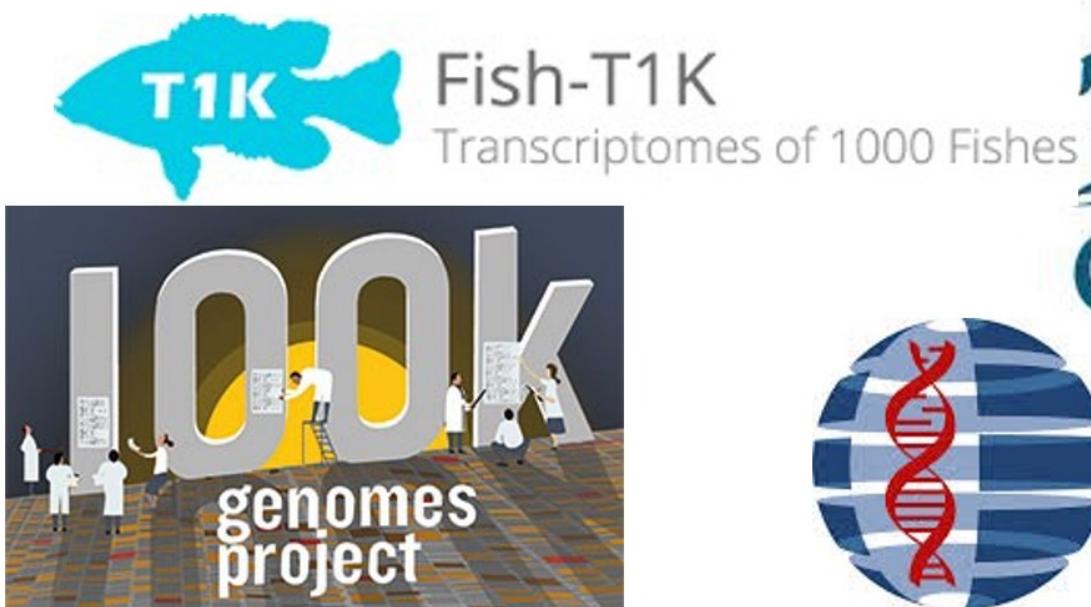
# Cost of sequencing the human genome



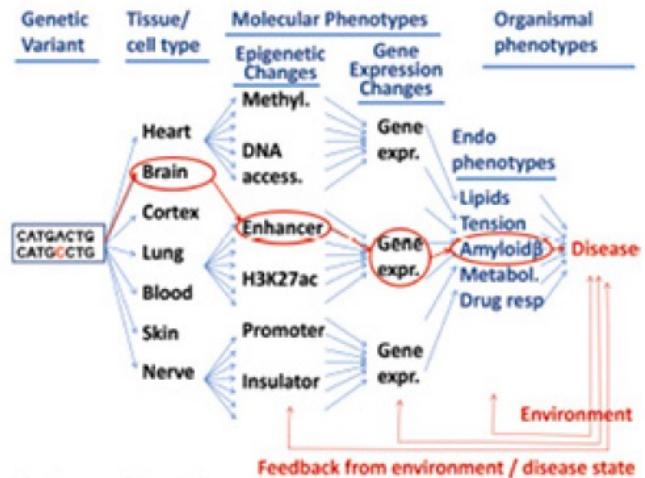
# A Deluge of Data



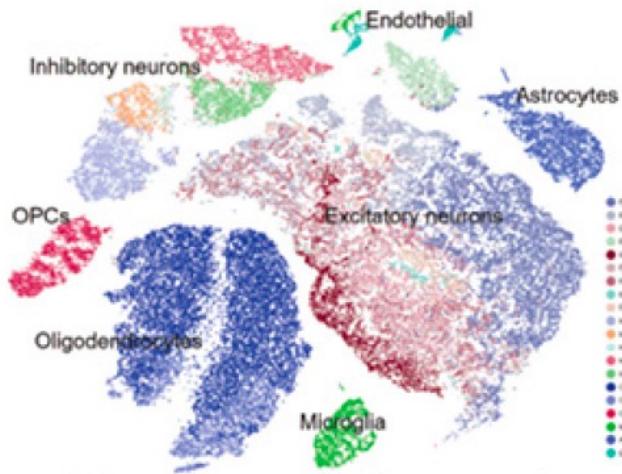
1000 Plant  
Genomes



# Many different applications



Mediation analysis/QTLs



Single-cell dissection

cAGxTGCCc  
CTCF

Gc AGAG G  
GAF

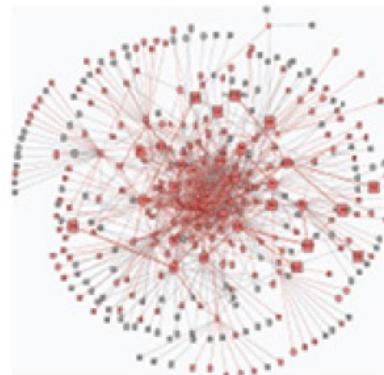
xT GCA TACTT Tx Gca  
Su(Hw)

G GATG cGTG  
BEAF-32

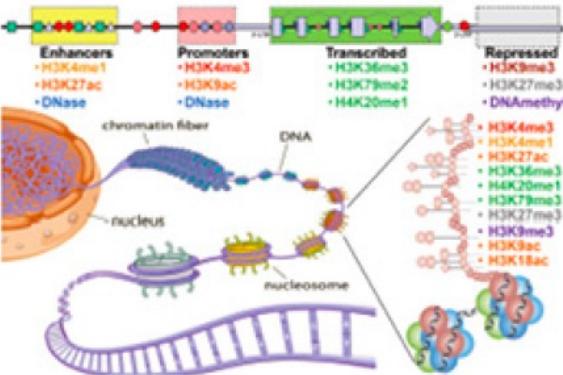
GGT CACAC TG  
CP190

GTATAAACAGT  
Mod(mdg4)

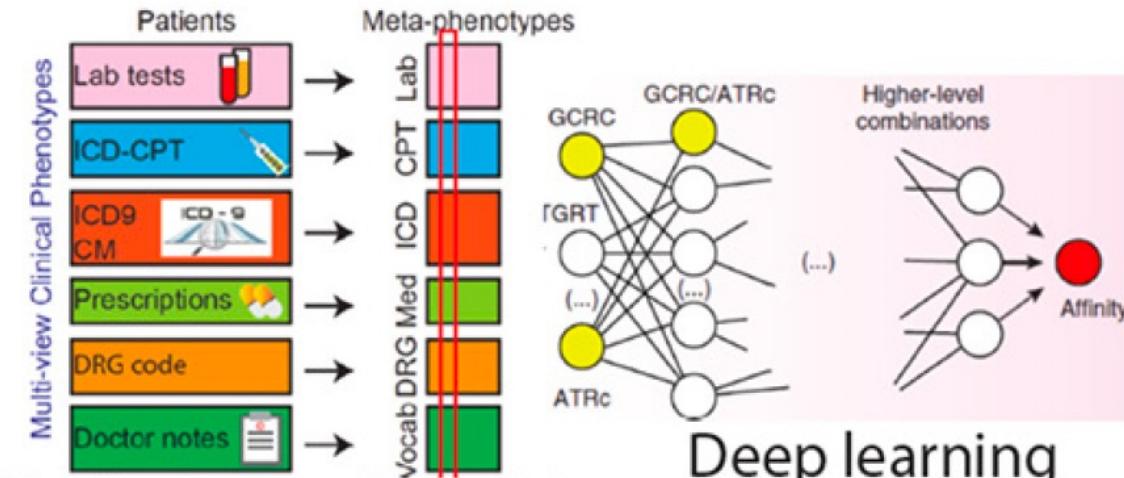
DNA motifs



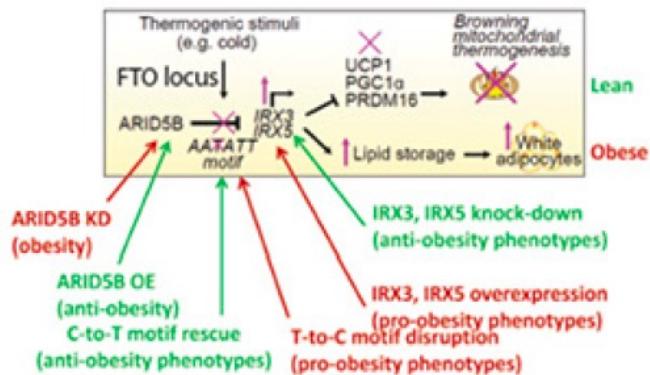
Gene networks



Epigenomics



Medical record models



Manipulate disease circuitry



Multi-phenotype

# mRNA vaccines

- The full genetic sequence of SARS-CoV-2 was deposited on January 5, 2020 (and made public on January 10)

Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome

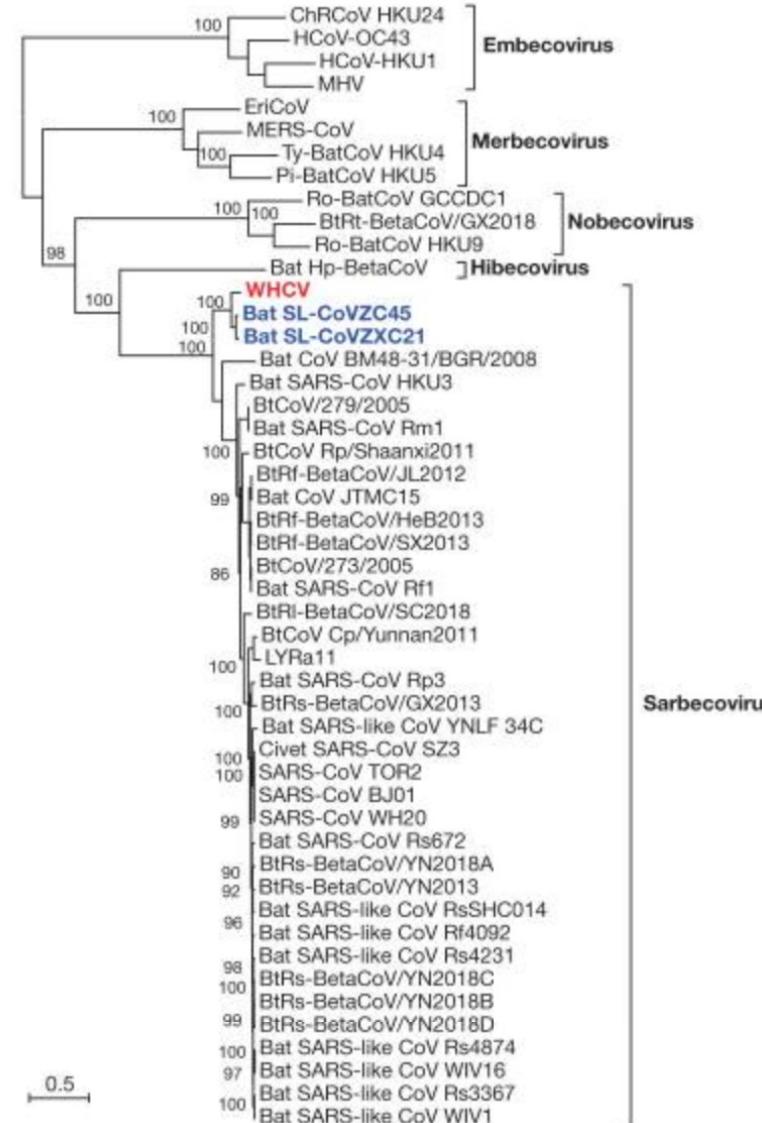
GenBank: MN908947.3

[FASTA](#) [Graphics](#)

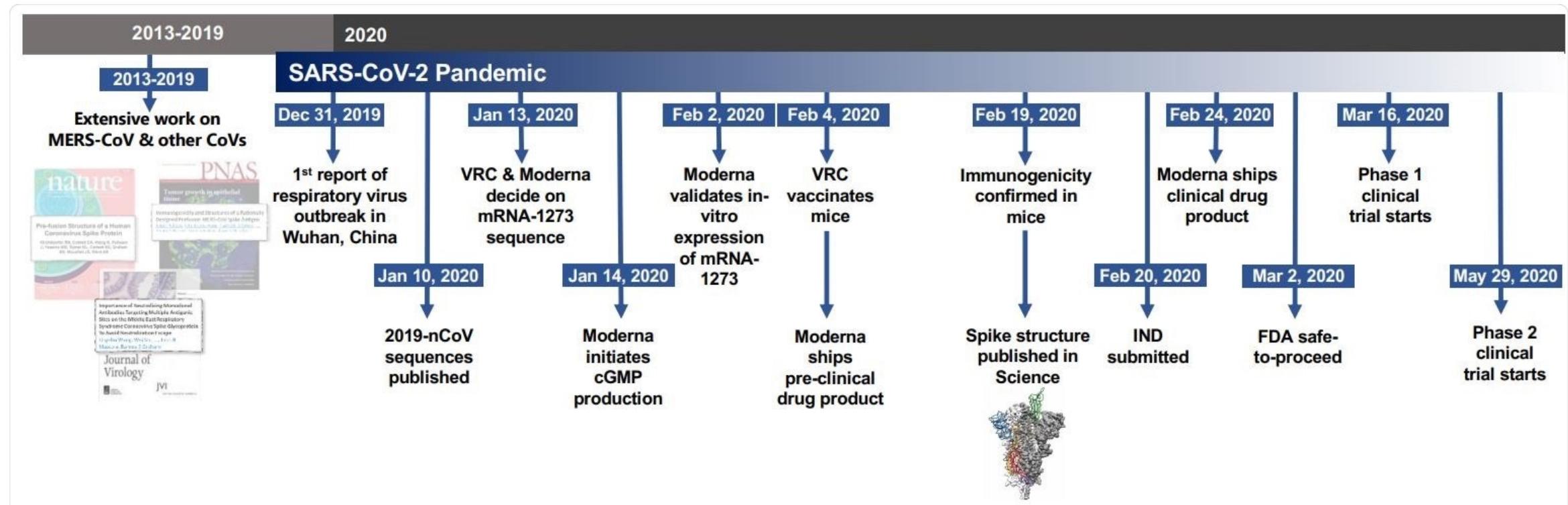
Go to:

LOCUS MN908947 29903 bp ss-RNA linear VRL 18-MAR-2020  
DEFINITION Severe acute respiratory syndrome coronavirus 2 isolate Wuhan-Hu-1, complete genome.  
ACCESSION MN908947  
VERSION MN908947.3  
KEYWORDS .  
SOURCE Severe acute respiratory syndrome coronavirus 2 (SARS-CoV-2)  
ORGANISM [Severe acute respiratory syndrome coronavirus 2](#)  
Viruses; Riboviria; Orthornavirae; Pisuviricota; Pisoniviricetes;  
Nidovirales; Cornidovirinae; Coronaviridae; Orthocoronavirinae;  
Betacoronavirus; Sarbecovirus.  
REFERENCE 1 (bases 1 to 29903)  
AUTHORS Wu,F., Zhao,S., Yu,B., Chen,Y.M., Wang,W., Song,Z.G., Hu,Y.,  
Tao,Z.W., Tian,J.H., Pei,Y.Y., Yuan,M.L., Zhang,Y.L., Dai,F.H.,  
Liu,Y., Wang,Q.M., Zheng,J.J., Xu,L., Holmes,E.C. and Zhang,Y.Z.  
TITLE A new coronavirus associated with human respiratory disease in  
China  
JOURNAL Nature 579 (7798), 265-269 (2020)  
PUBMED 32015508  
REFERENCE 2 (bases 1 to 29903)  
AUTHORS Wu,F., Zhao,S., Yu,B., Chen,Y.-M., Wang,W., Hu,Y., Song,Z.-G.,  
Tao,Z.-W., Tian,J.-H., Pei,Y.-Y., Yuan,M.L., Zhang,Y.-L.,

<https://www.ncbi.nlm.nih.gov/nuccore/MN908947>

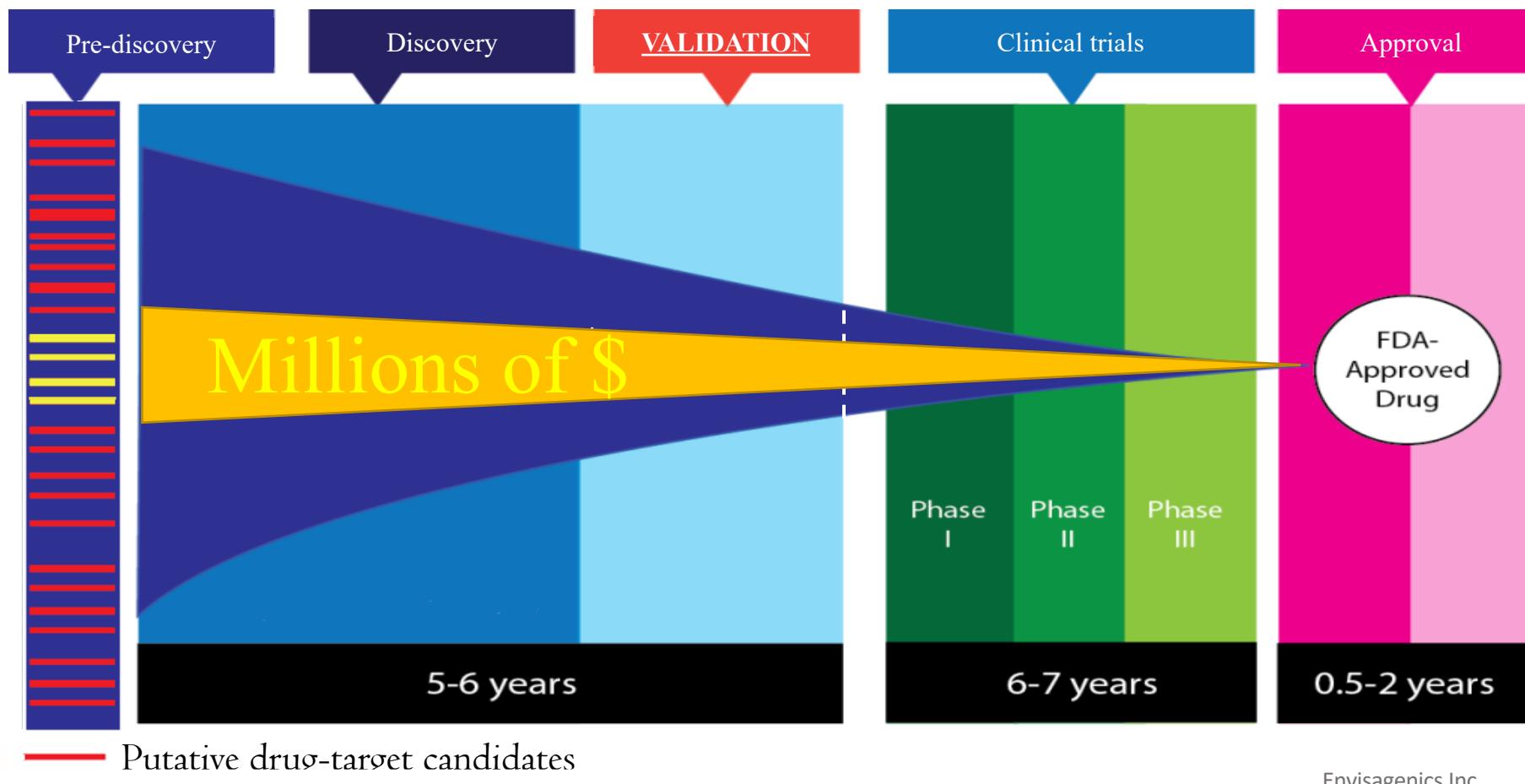


# Moderna vaccine development timeline



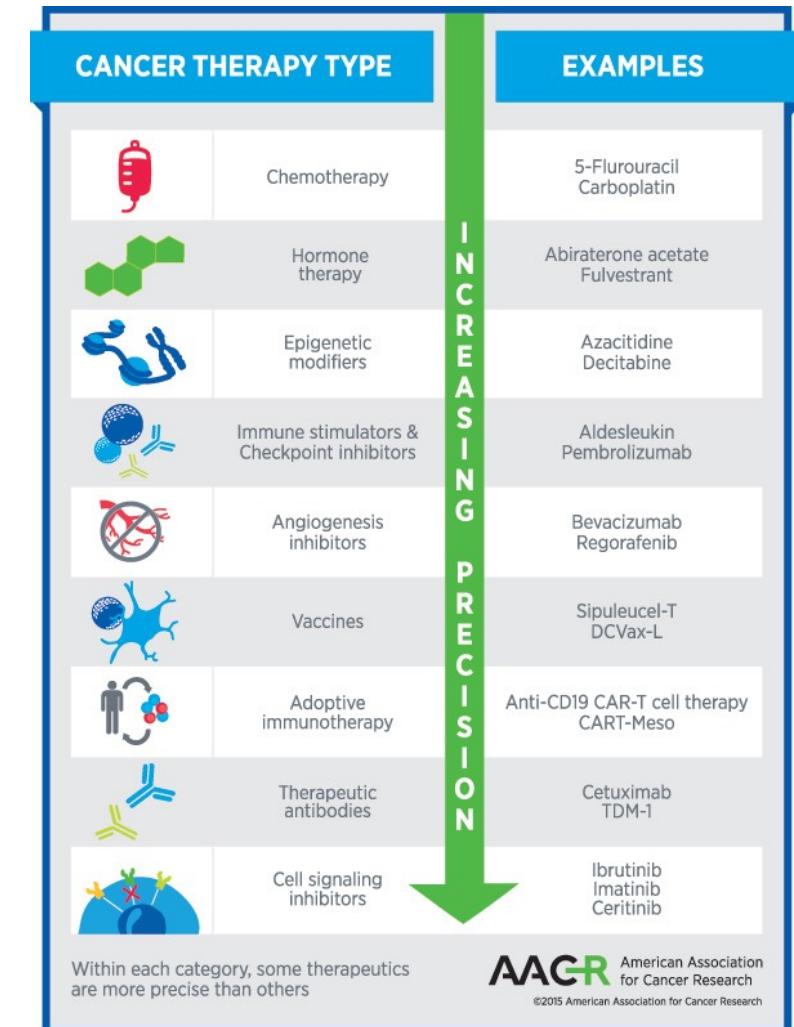
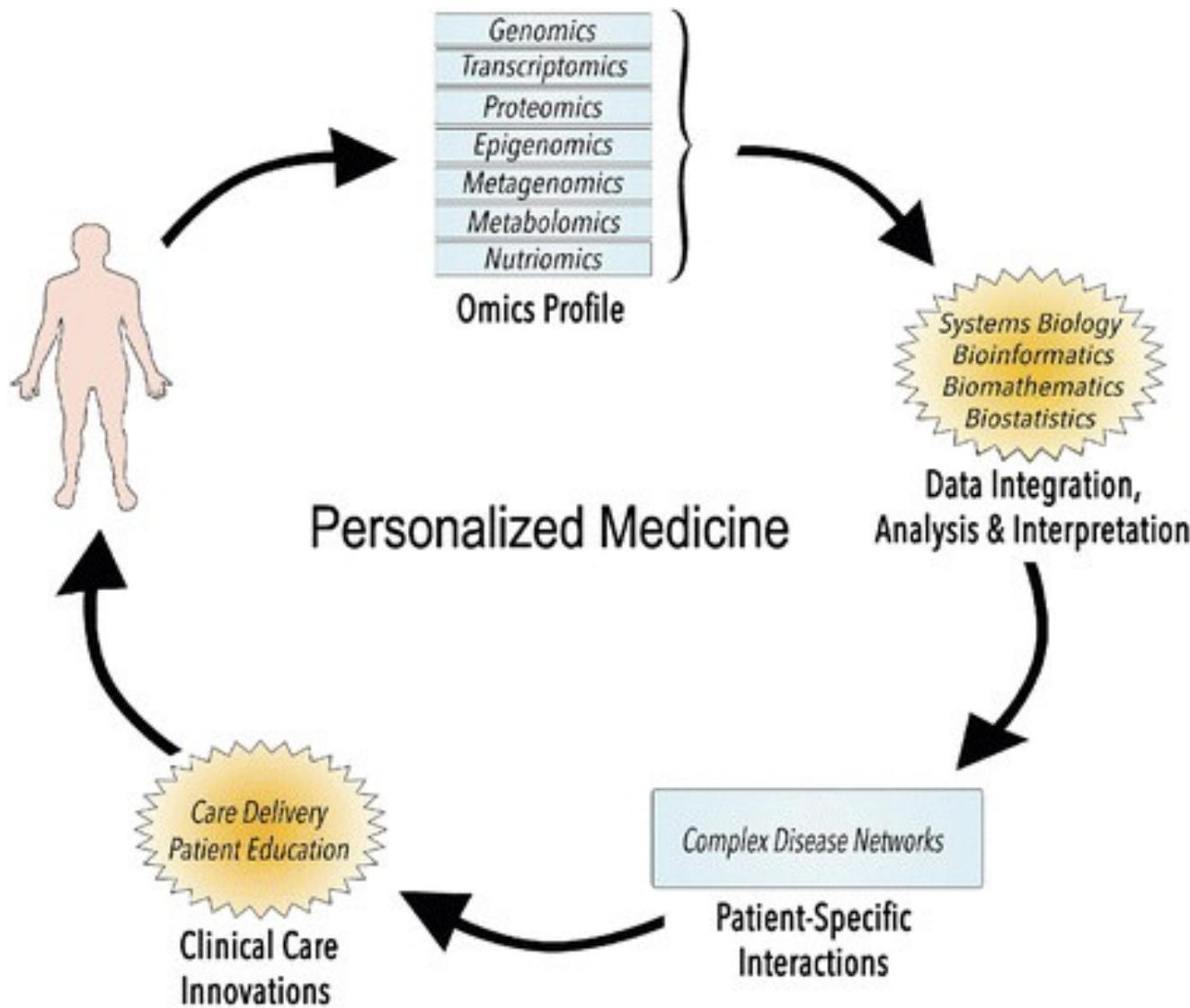
2 days from sequence to vaccine candidate!

# Bioinformatics can dramatically reduce the cost and time for developing a new drug



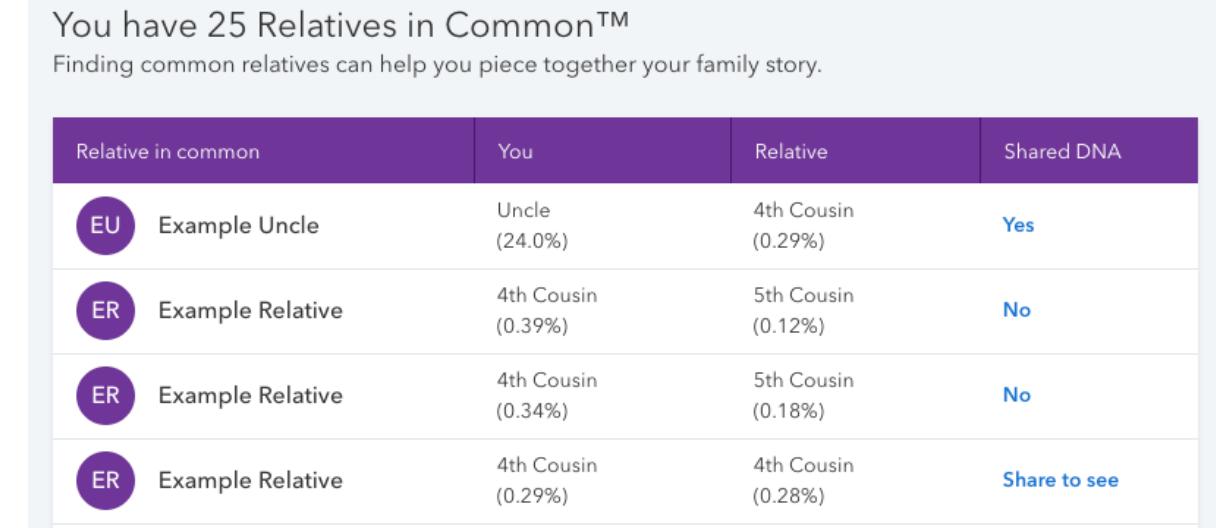
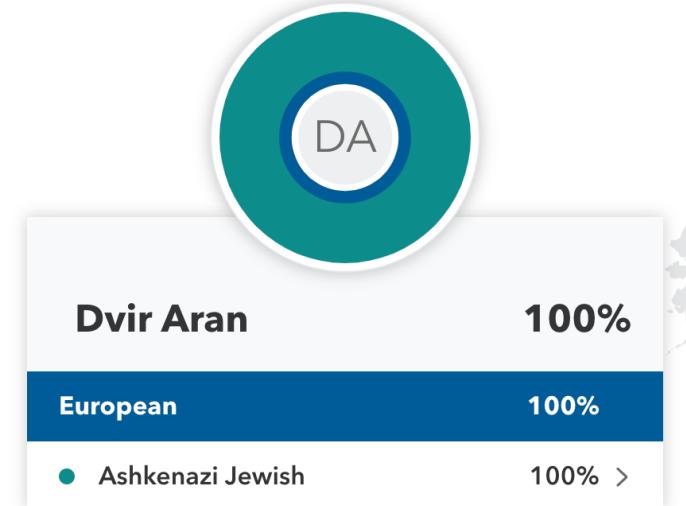
From drug discovery through FDA approval, developing a new medicine takes at least 10 years on average and costs an average of \$2.6 billion.\* Less than 12% of the candidate medicines that make it into Phase I clinical trials will be approved by the FDA.

# Personalized/precision medicine



# Ancestry composition and identifying far relatives

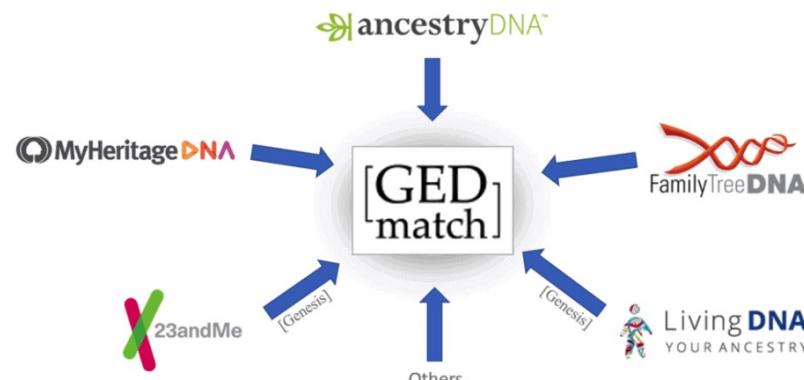
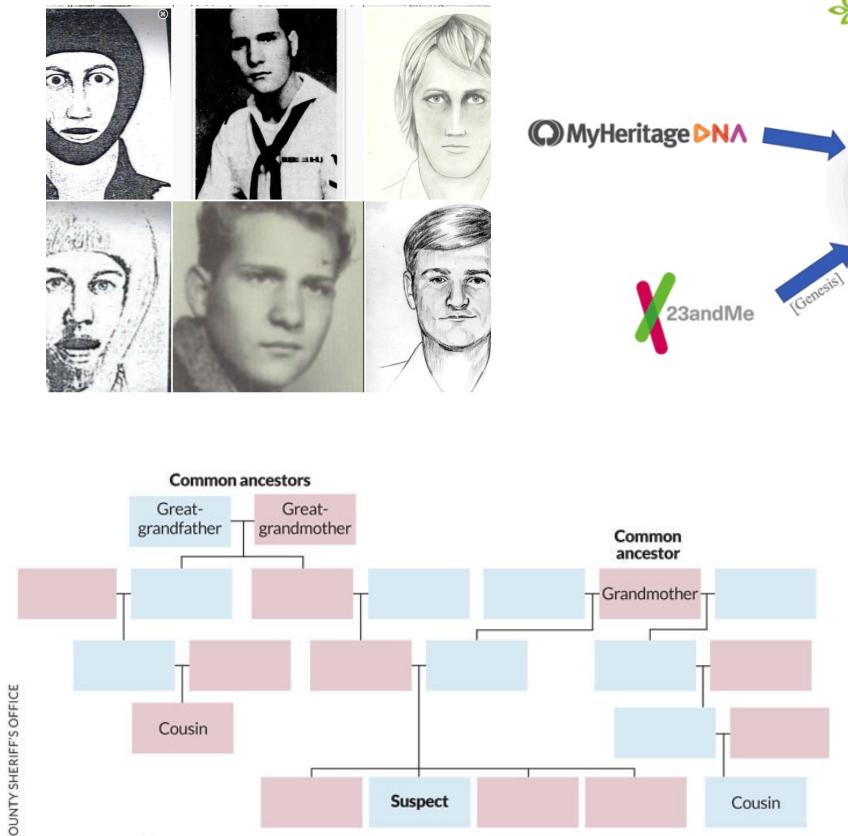
- ~50 million people in the world deposited their genetic information in public databases.
- Bioinformatics methods allow to reconstruct family trees.



# Identify far relatives... and criminals!

## Long-range familial search

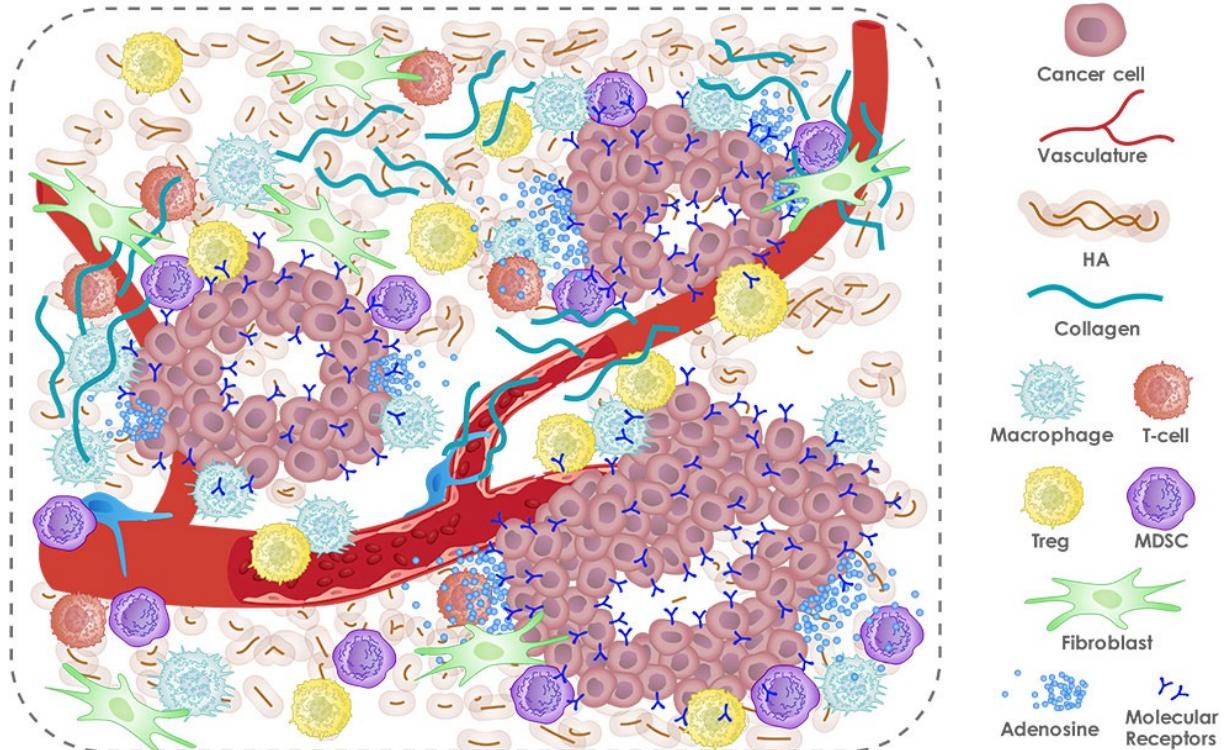
### Golden State Killer



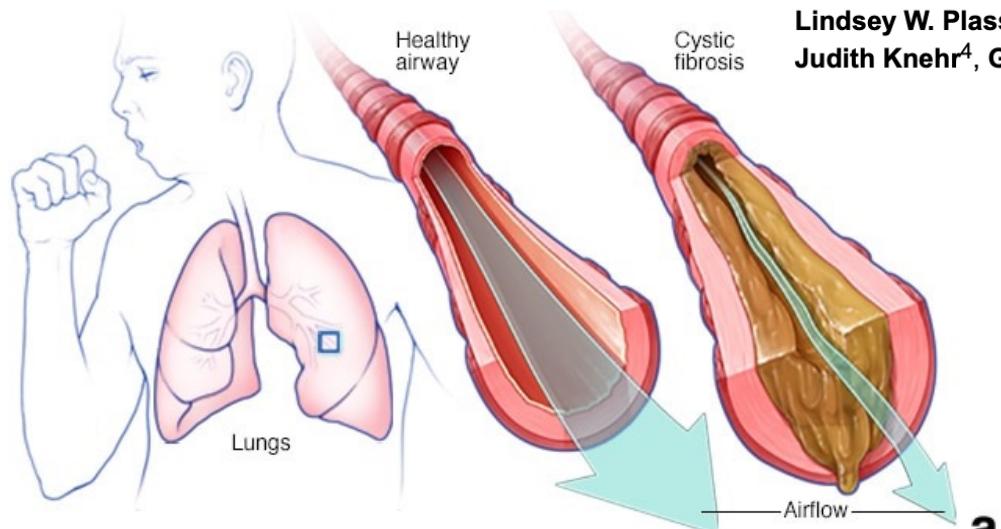
# Single-cell biology

- How a single cell develops into an adult animal with multiple organs and billions of cells?
- How a mutation causes a specific disease?
- How is the immune system shaped in the tumor microenvironment?

The Tumor Microenvironment



# Cystic Fibrosis



© MAYO FOUNDATION FOR MEDICAL EDUCATION AND RESEARCH. ALL RIGHTS RESERVED.

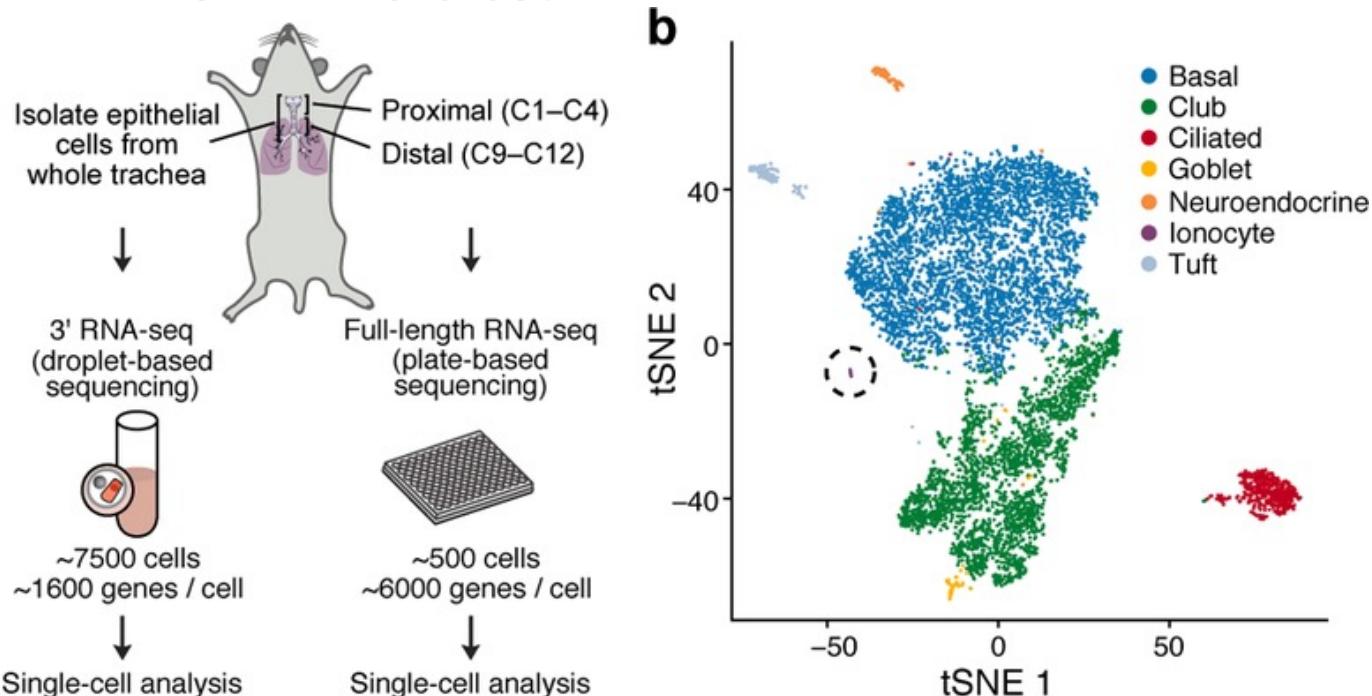
- Accumulation of thick, sticky mucus
- Coughing or shortness of breath
- Frequent chest infections
- Salty-tasting skin
- Poor growth and poor weight gain despite normal food intake
- **Mutation in the CFTR gene.**
- 1/2500 babies
- 1/25 every AZ is a carrier

## A single cell atlas of the tracheal epithelium reveals the CFTR-rich pulmonary ionocyte

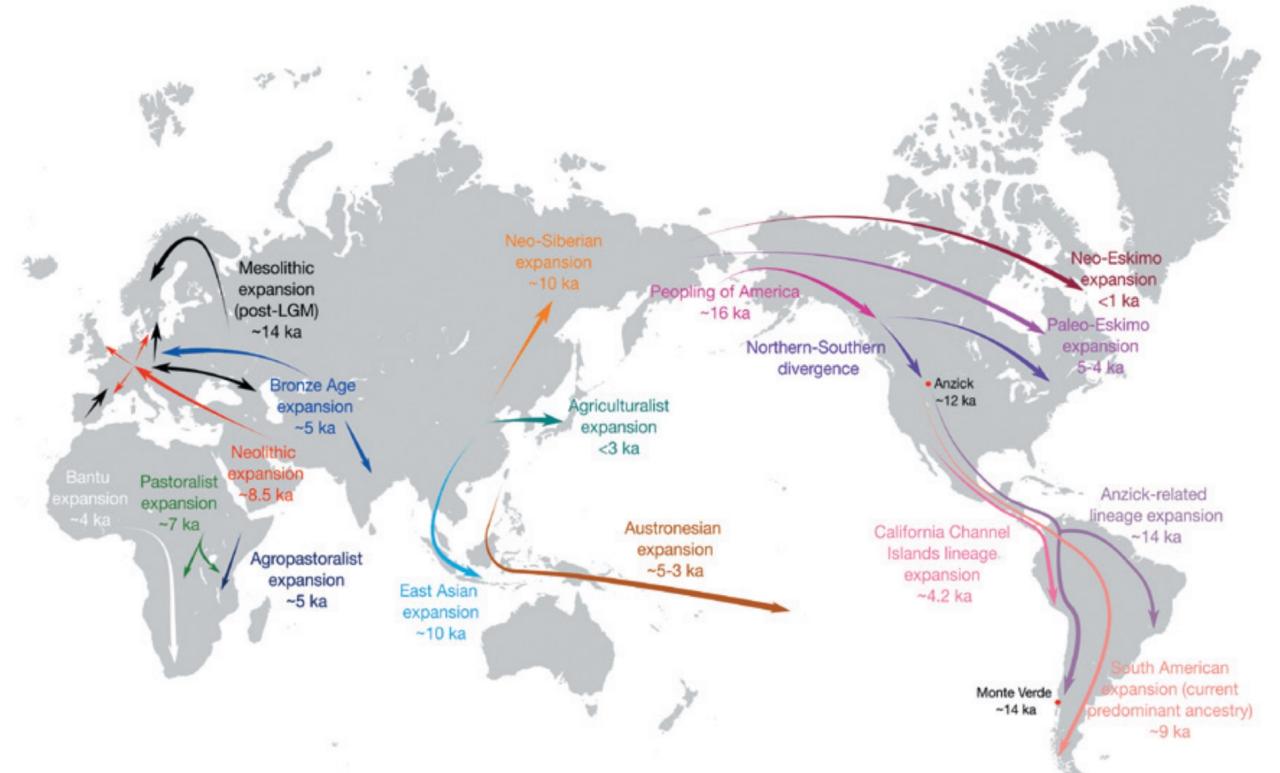
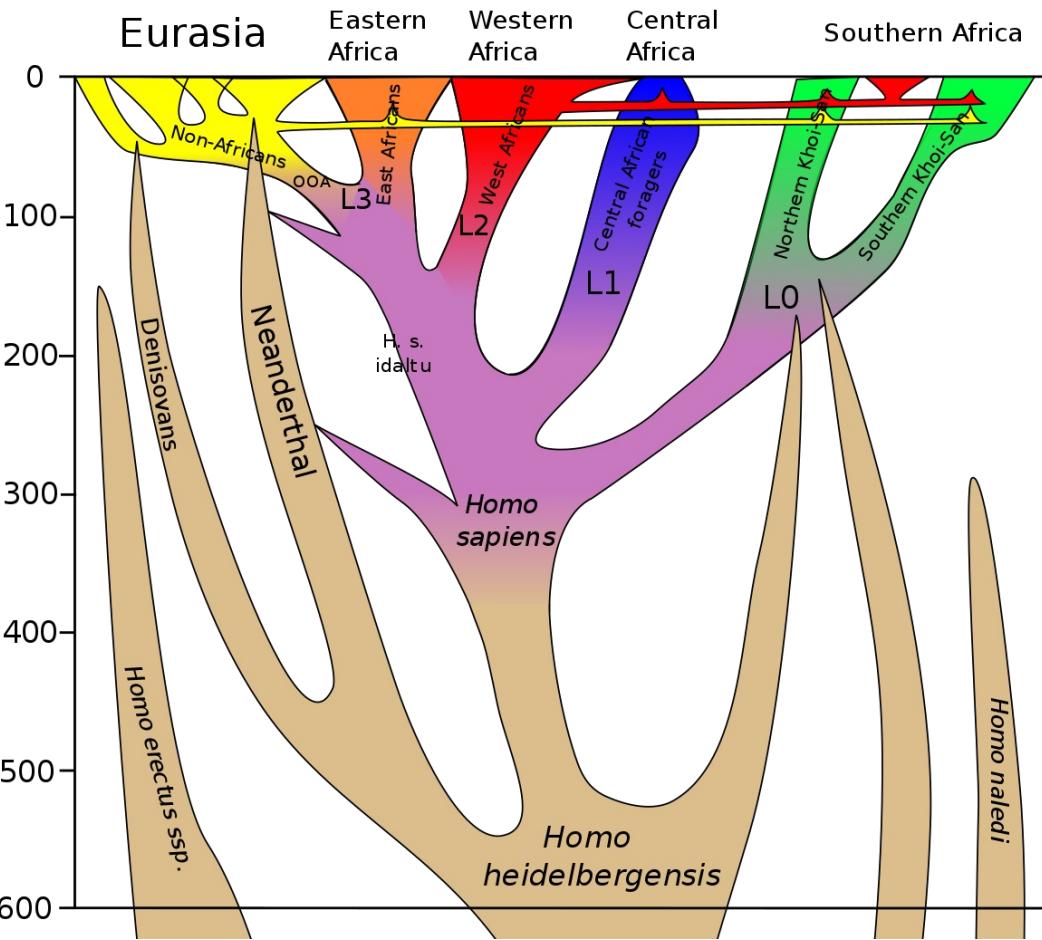
Lindsey W. Plasschaert<sup>#1</sup>, Rapolas Žilionis<sup>#2,3</sup>, Rayman Choo-Wing<sup>1</sup>, Virginia Savova<sup>2</sup>, Judith Knehr<sup>4</sup>, Guglielmo Roma<sup>4</sup>, Allon M. Klein<sup>2,†</sup>, and Aron B. Jaffe<sup>1,†</sup>

## A revised airway epithelial hierarchy includes CFTR-expressing ionocytes

Daniel T. Montoro<sup>#1,2,3</sup>, Adam L. Haber<sup>#4</sup>, Moshe Biton<sup>#4,5</sup>, Vladimir Vinarsky<sup>1,2,3</sup>, Brian Lin<sup>1,2,3</sup>, Susan Birket<sup>6,7</sup>, Feng Yuan<sup>8</sup>, Sijia Chen<sup>9</sup>, Hui Min Leung<sup>10,11</sup>, Jorge Villoria<sup>1,2,3</sup>, Noga Rogel<sup>4</sup>, Grace Burgin<sup>4</sup>, Alexander Tsankov<sup>4</sup>, Avinash Waghray<sup>1,2,3</sup>, Michal Slyper<sup>4</sup>, Julia Waldmann<sup>4</sup>, Lan Nguyen<sup>4</sup>, Danielle Dionne<sup>4</sup>, Orit Rozenblatt-Rosen<sup>4</sup>, Purushothama Rao Tata<sup>12,13,14,15</sup>, Hongmei Mou<sup>16,17</sup>, Manjunatha Shivaraju<sup>1,2,3</sup>, Hermann Bihler<sup>18</sup>, Martin Mense<sup>18</sup>, Guillermo J. Tearney<sup>10,11</sup>, Steven M. Rowe<sup>6,7</sup>, John F. Engelhardt<sup>8</sup>, Aviv Regev<sup>4,19,§</sup>, and Jayaraj Rajagopal<sup>1,2,3,§</sup>



# Human history

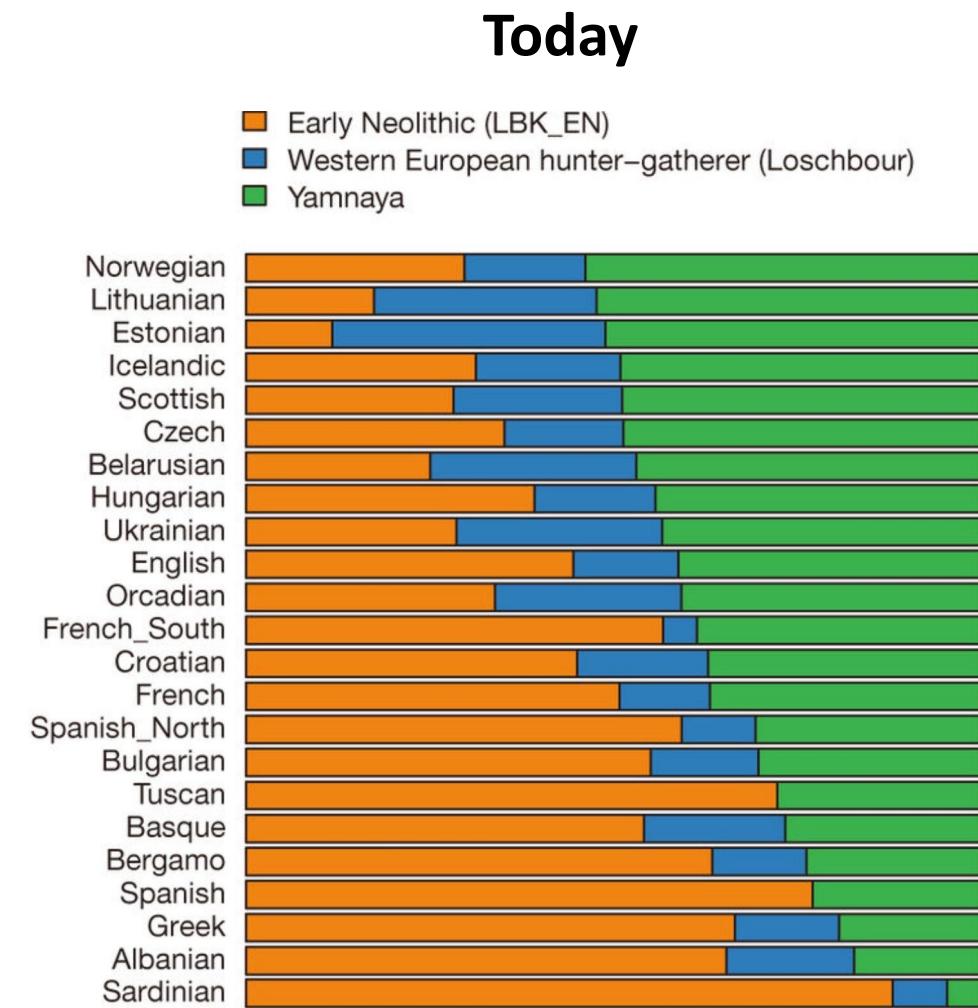
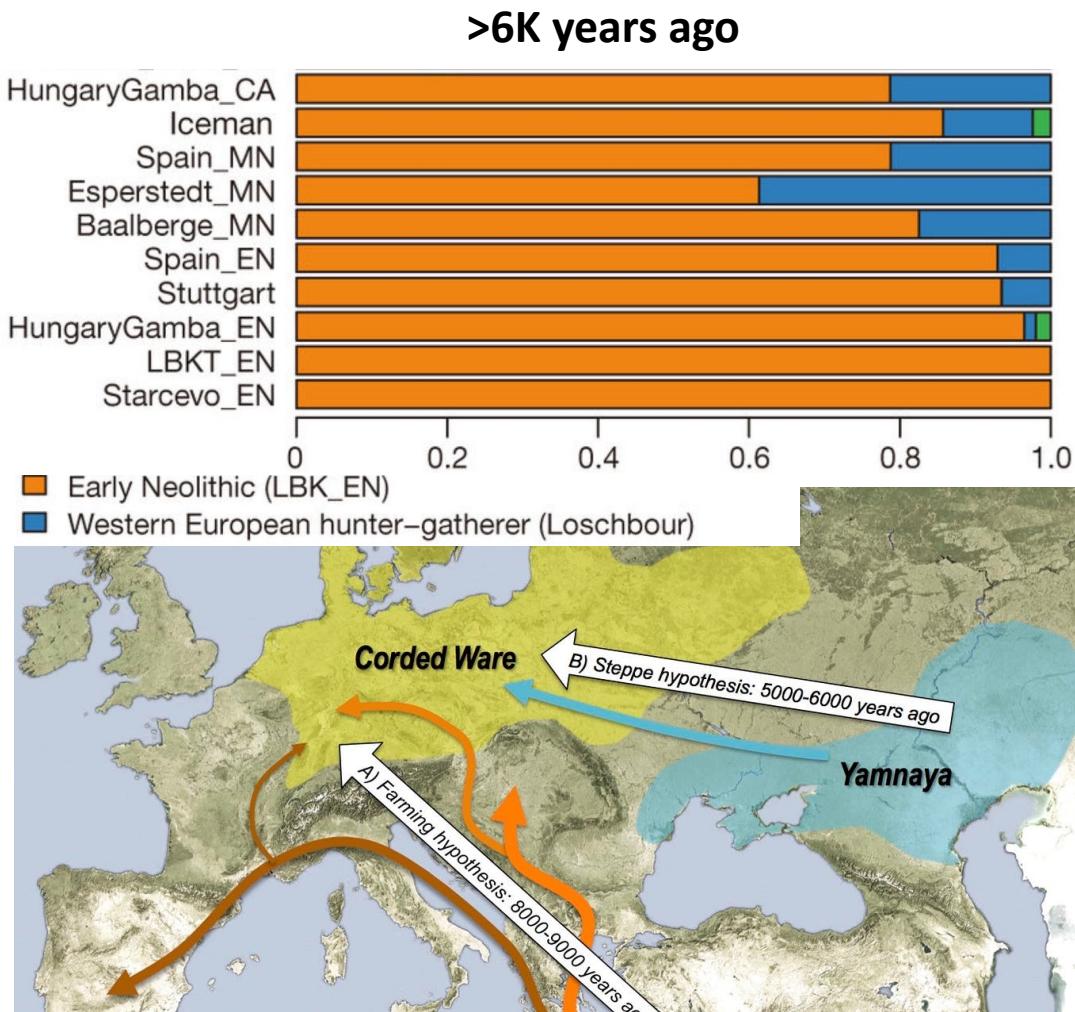


# Ancient DNA



Svante Paabo

# Ancient DNA



# Israeli companies in biomedical data science

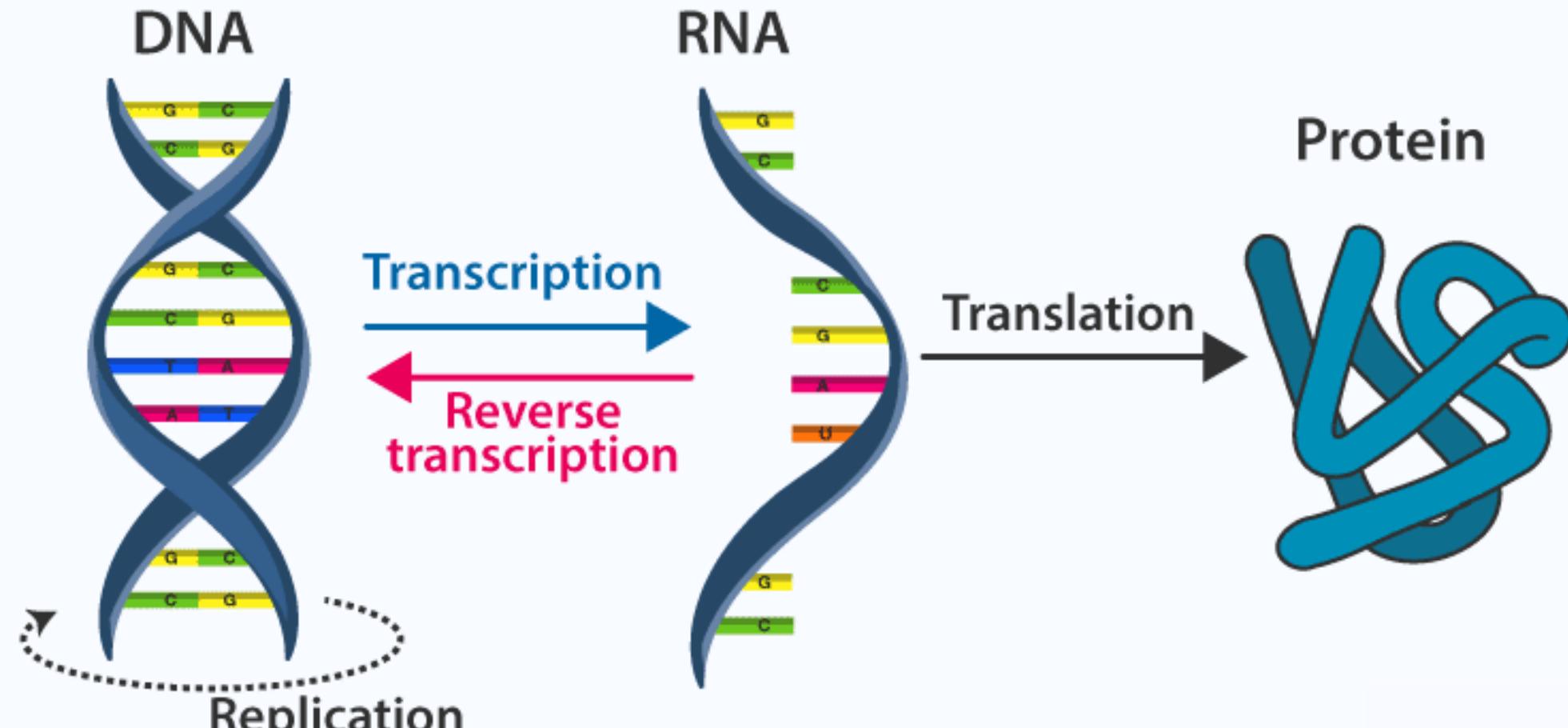




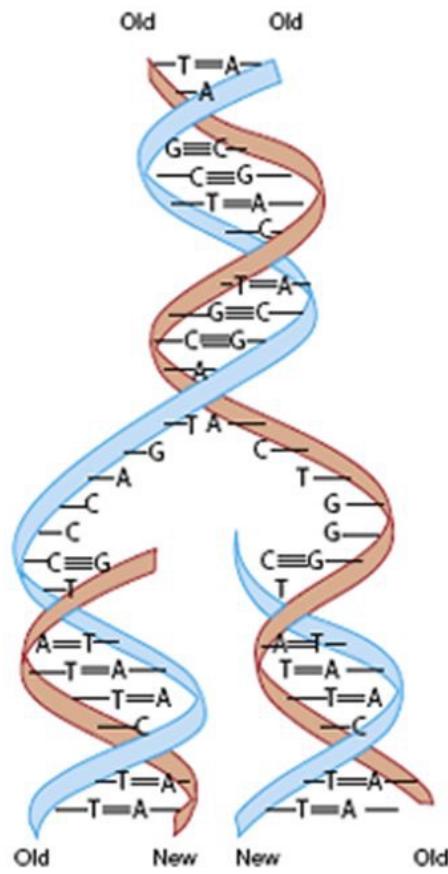
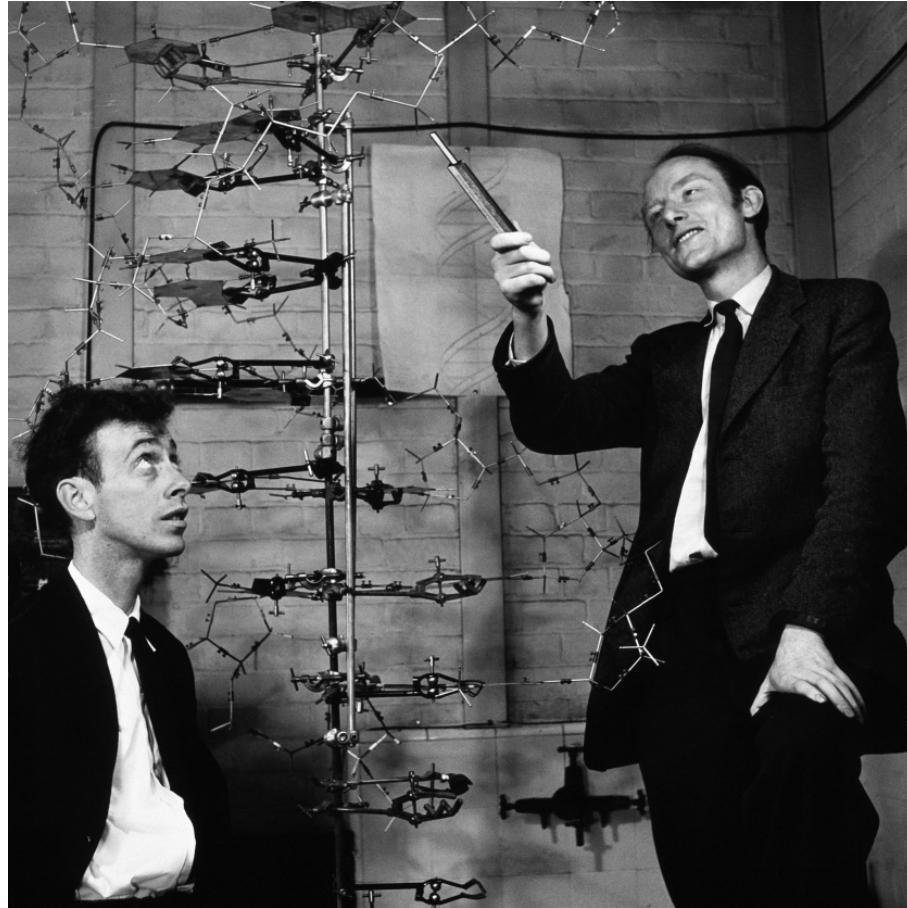
# Biology primer

Quick introduction to molecular biology  
and information transfer within the cell

## CENTRAL DOGMA : DNA TO RNA TO PROTEIN



# The Central Dogma of Molecular Biology



8

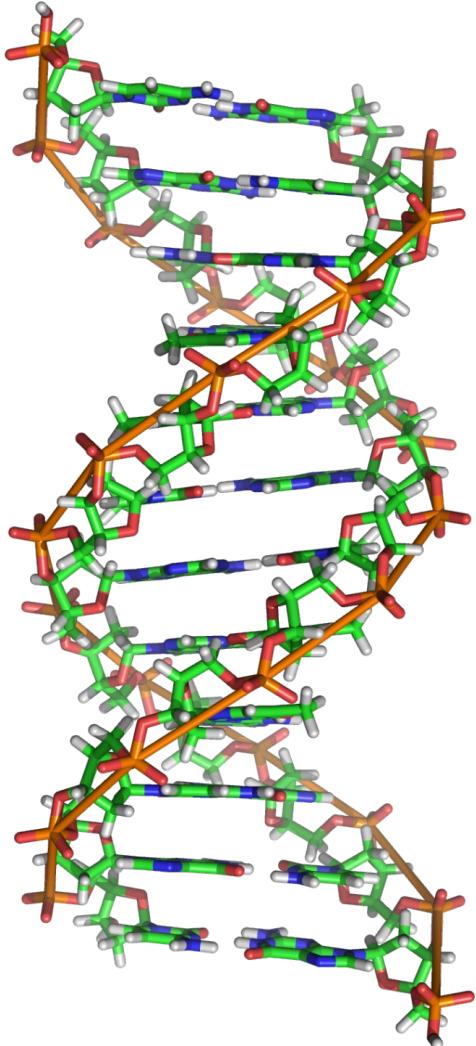
"It has not escaped our notice that the specific pairing we have postulated immediately suggests a possible copying mechanism for the genetic material."

-Watson and Crick

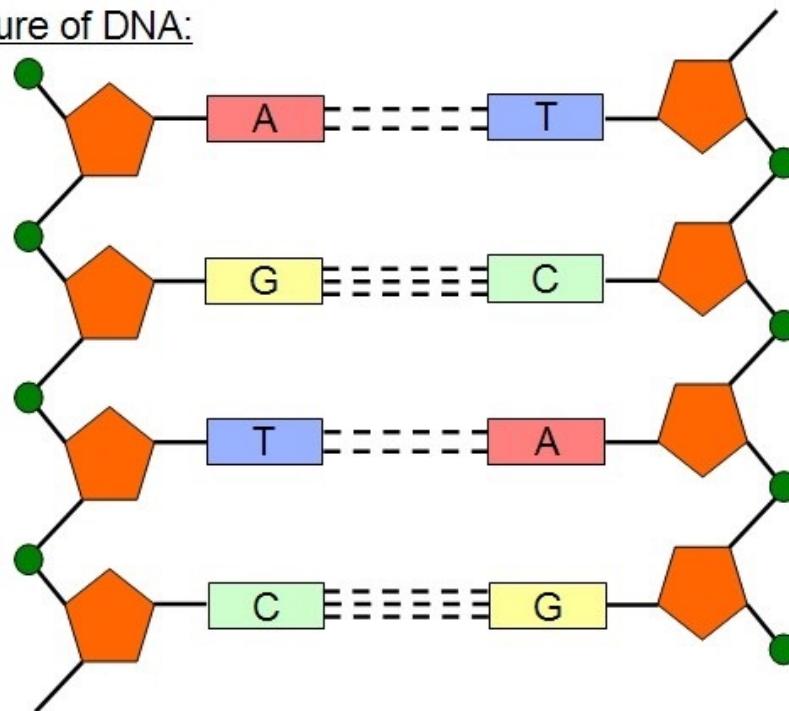
- Cellular molecules:**
1. DNA
  2. RNA
  3. Protein

# DNA

Each strand composed of sequence of covalently bonded **nucleotides (bases)**.



Structure of DNA:



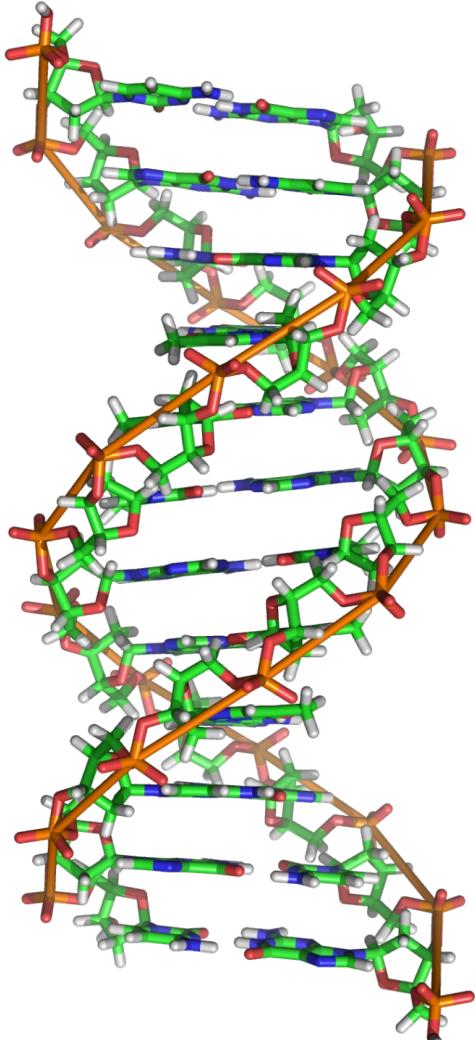
A $\leftarrow\rightarrow$ T, C $\leftarrow\rightarrow$ G Watson-Crick base-pairing

**Four nucleotides:**

- A (adenine)
- C(cytosine)
- T (thymine)
- G(guanine)

# DNA

Each strand composed of sequence of covalently bonded **nucleotides (bases)**.



5' ...ACGTGACTGAGGACCGTG...  
...|||...|||...|||...|||...|||...

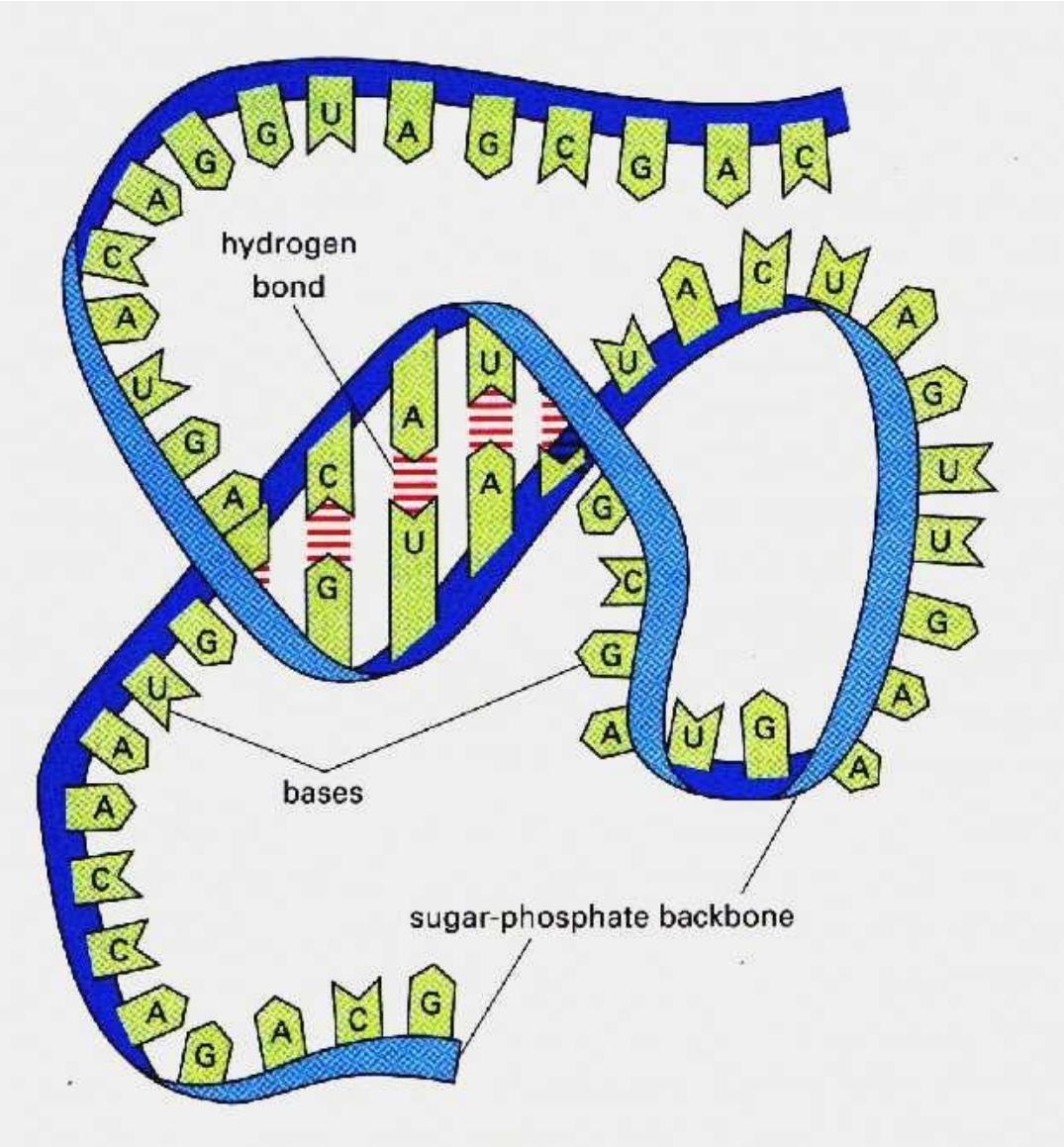
3' ...TGCAC TGACT CCTGGCAC...  
                                      5'

Pair of strings  
from 4-character  
alphabet

5' ...ACGTGACTGAGGACCGTG  
CGACTGAGACTGACTGGGT  
CTAGCTAGACTACGTTTA  
TATATATATACTCGTCGT  
ACTGATGACTAGATTACAG  
TGATTTAAAAAAATATT... 3'

Single string  
from 4-character  
alphabet

# RNA



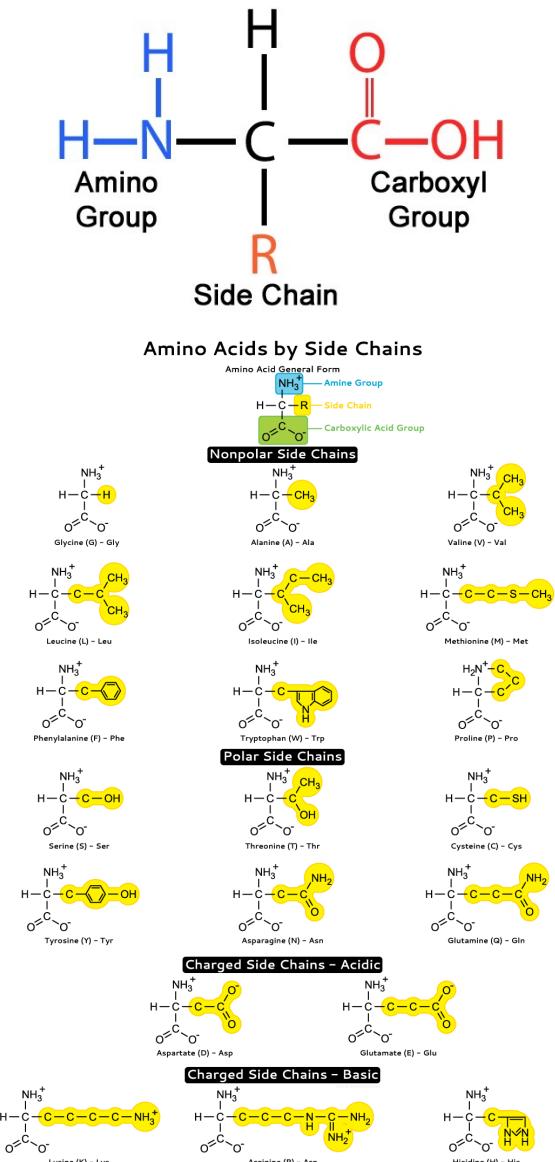
- **Single-stranded**
  - A(adenine)
  - C(cytosine)
  - U(uracil)
  - G(guanine)
- Can fold into **structures** due to base complementarity.  
 $A \leftarrow \rightarrow U, C \leftarrow \rightarrow G$
- Comes in many flavors:  
mRNA, rRNA, tRNA, tmRNA, snRNA, snoRNA, scaRNA, aRNA, asRNA, piwiRNA, etc.

# Protein

- String of amino acids: 20 letter alphabet

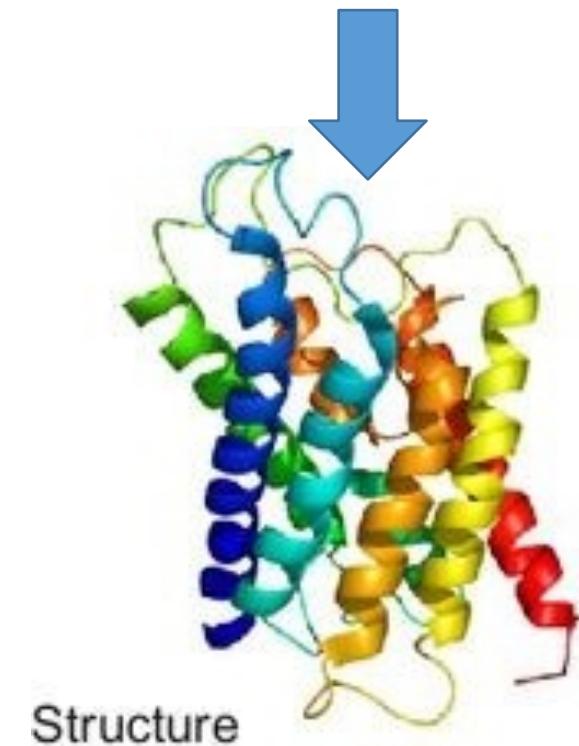
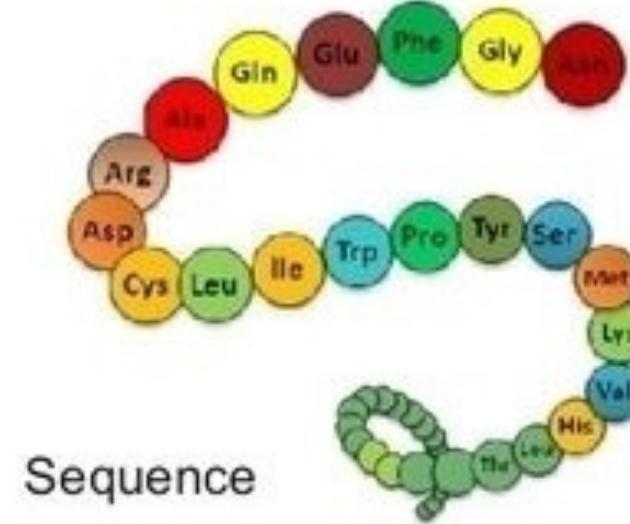
...DTIGDWNSPSFFGIQLVSSVHT  
 TLWYRENAFPVLGGFSWLSWFNW  
 HNMGYYPVYHIGYPMIRCGTHL  
 VPMQFAFQSIARSFALVHNAPM  
 VLKINPHERQDPVFWPCLYYSVD  
 IRSMHIGYPMIRCYQA...

Amino Acid	3-Letters	1-Letter
Alanine	Ala	A
Arginine	Arg	R
Asparagine	Asn	N
Aspartic acid	Asp	D
Cysteine	Cys	C
Glutamic acid	Glu	E
Glutamine	Gln	Q
Glycine	Gly	G
Histidine	His	H
Isoleucine	Ile	I
Leucine	Leu	L
Lysine	Lys	K
Methionine	Met	M
Phenylalanine	Phe	F
Proline	Pro	P
Serine	Ser	S
Threonine	Thr	T
Tryptophan	Trp	W
Tyrosine	Tyr	Y
Valine	Val	V

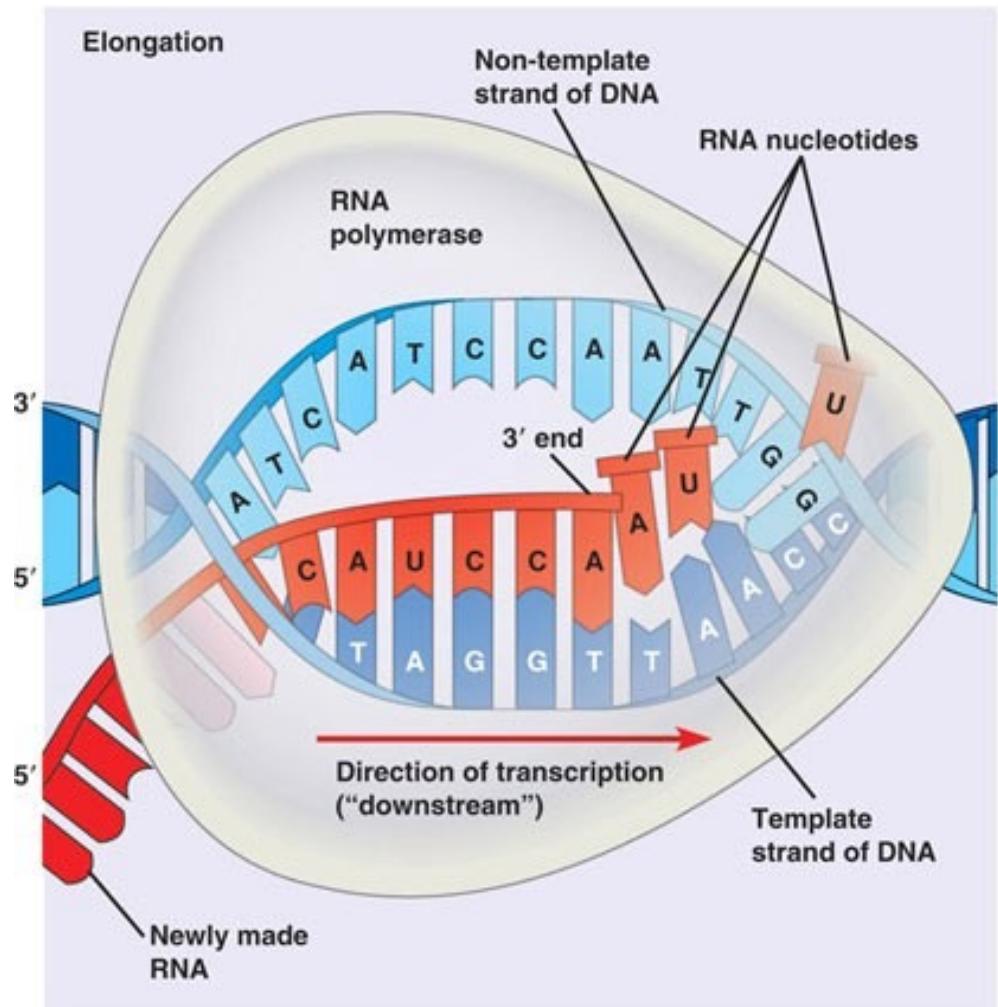


# Protein

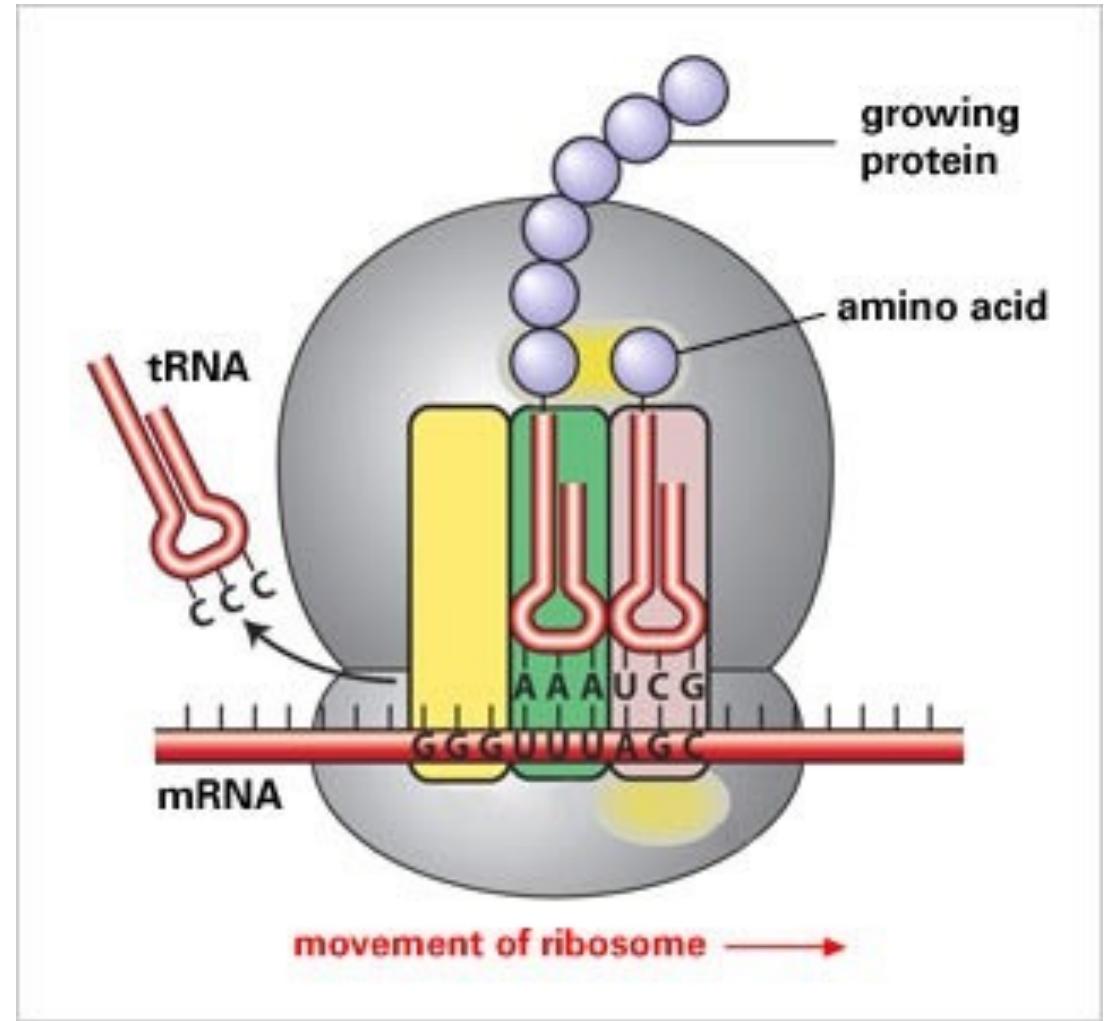
- String of amino acids: 20 letter alphabet
- Folds into 3D structures to perform various functions in cells



# Transcription and Translation

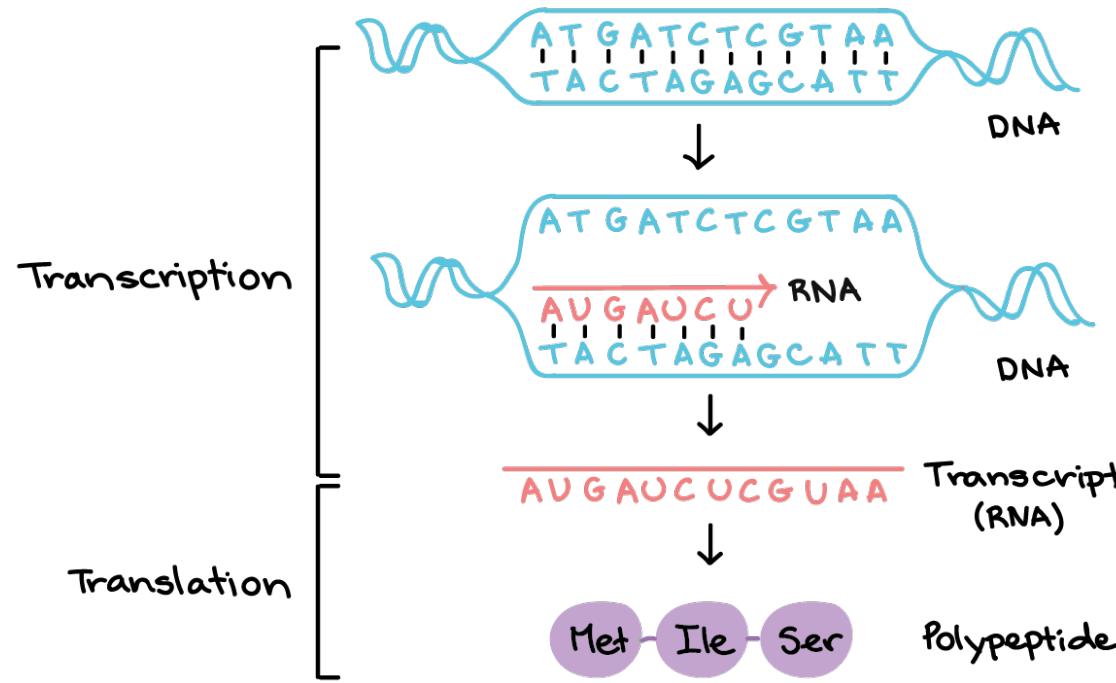


<http://dna-rna.net/wp-content/uploads/2011/08/rna-transcription2.jpg>



[http://www.frontiers-in-genetics.org/en/pictures/translation\\_1.jpg](http://www.frontiers-in-genetics.org/en/pictures/translation_1.jpg)

# Transcription and Translation



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

		Second base							
		U	C	A	G				
First base	U	UUU UUC UUA UUG	UCU UCC UCA UCG	UAU UAC	UGU UGC	Cysteine C			
	C	CUU CUC CUA CUG	CCU CCC CCA CCG	CAU CAC	CGU CGC	Stop codon Tryptophan W			
A	AUU AUC AUA AUG M	I	ACU ACC ACA ACG	CAA CAG	CAU CAC	Leucine L Proline P Glutamine Q Histidine H			
G	GUU GUC GUA GUG	V	GCU GCC GCA GCG	GAA GAG	GAU GAC	Serine S Arginine R Lysine K Asparagine N			
						Alanine A Aspartic acid D Glutamic acid E Glycine G			

<http://bioinfo.bisr.res.in/project/crat/pictures/codon.jpg>

# Primer on Molecular Biology

## Three fundamental molecules:

### 1. DNA

Information storage.

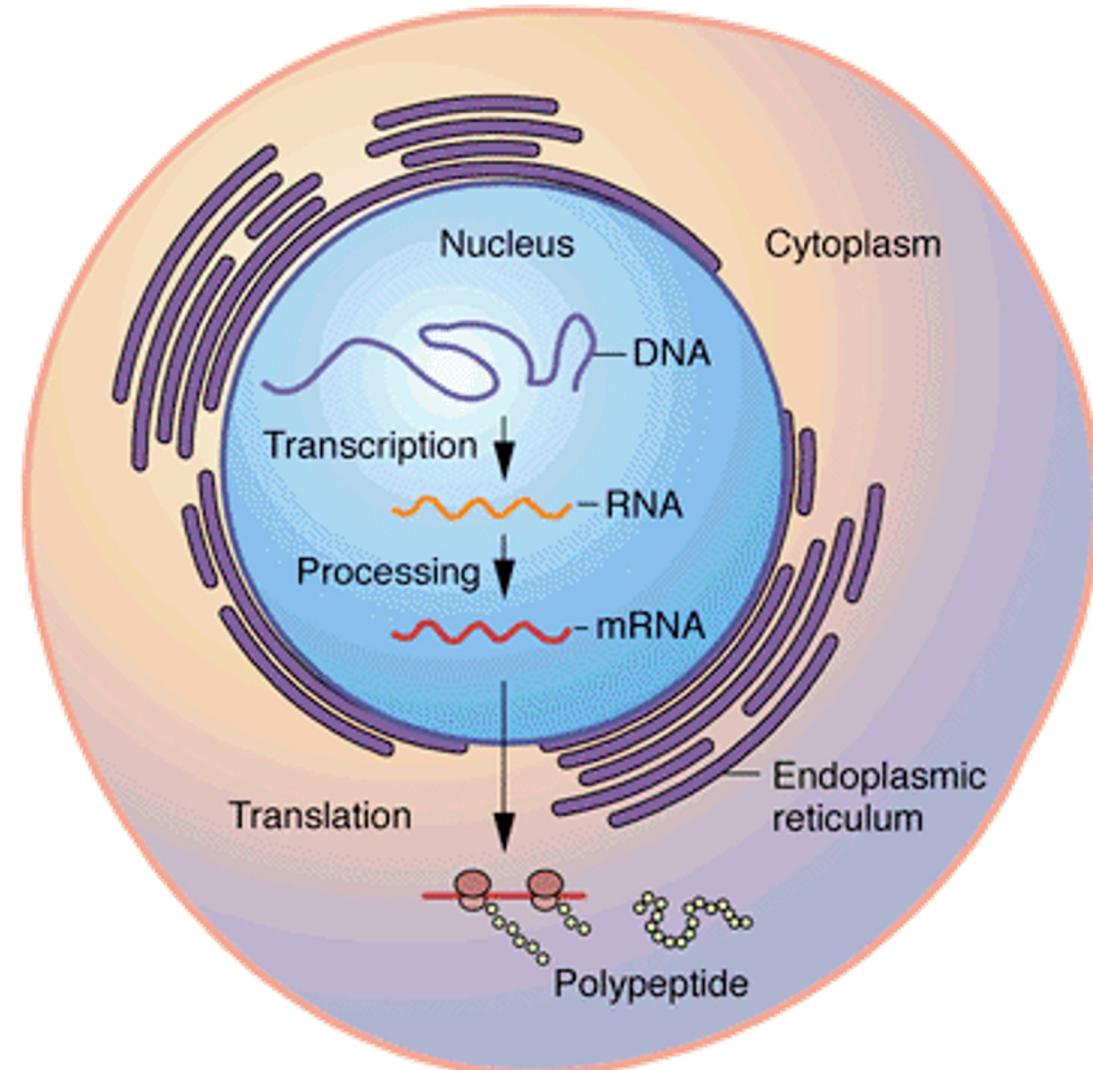
### 2 RNA

Old view: Mostly a “messenger”.

New view: Performs many important functions.

### 3. Protein

Perform most cellular functions  
(biochemistry, signaling, control, etc.)

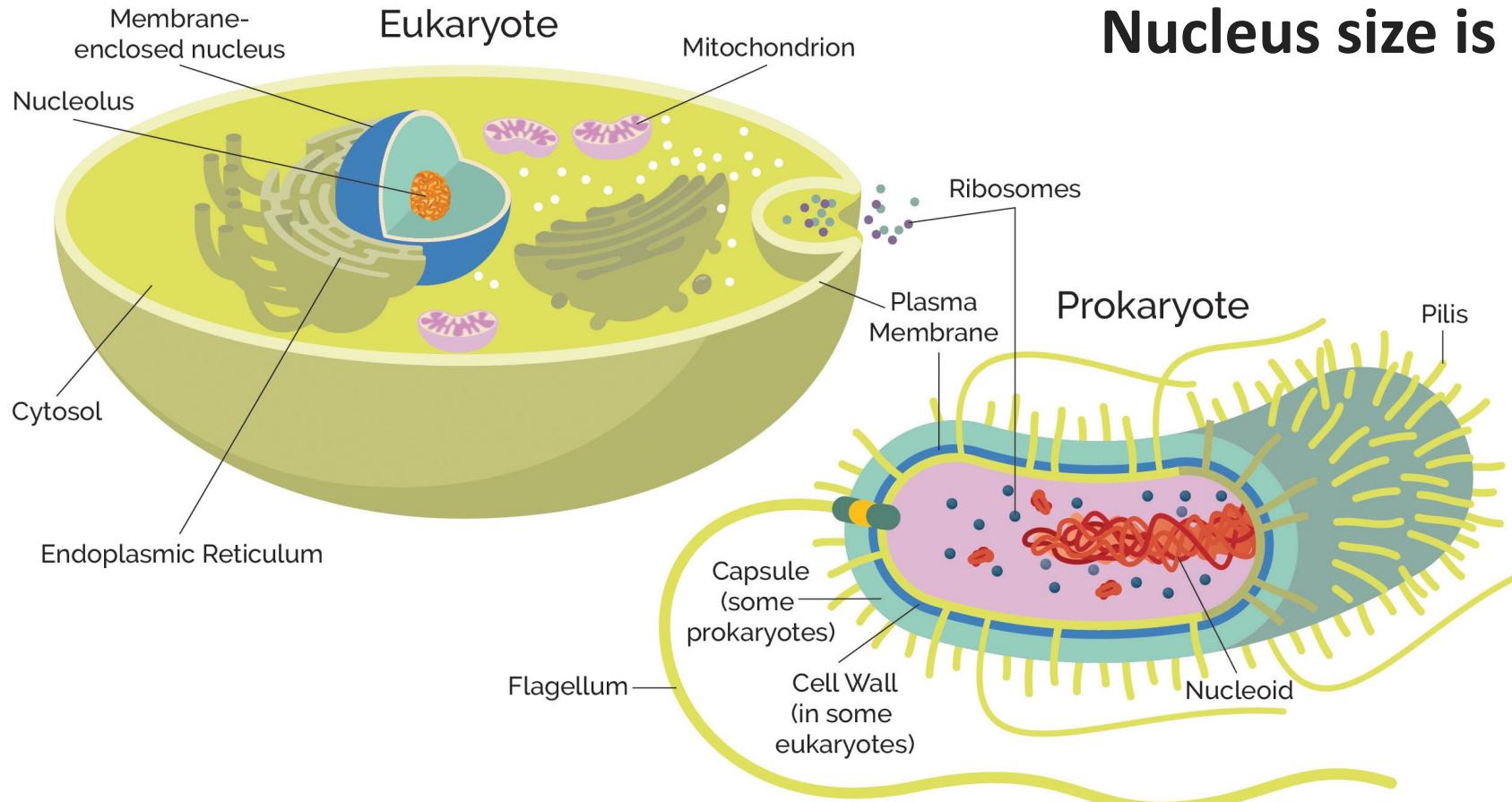


# Genomes size

Species	<i>T2 phage</i>	<i>Escherichia coli</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Paris japonica</i>
Genome Size	170,000 bp	4.6 million bp	130 million bp	3.2 billion bp	150 billion bp
Common Name	 Virus	 Bacteria	 Fruit fly	 Human	 Canopy Plant

Stretched out – human genome is 2.13 meters long

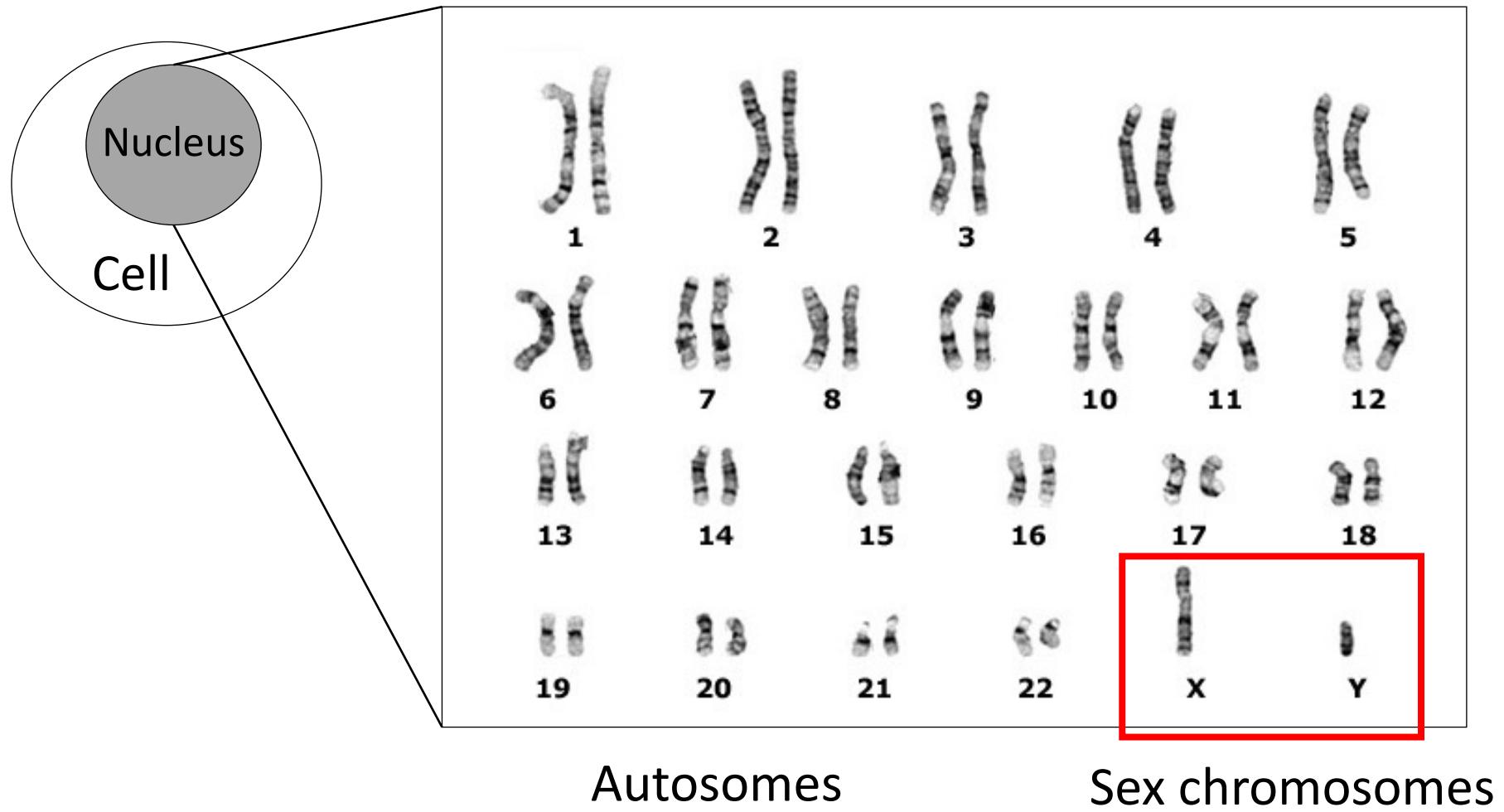
# Chromosomes inside the cell



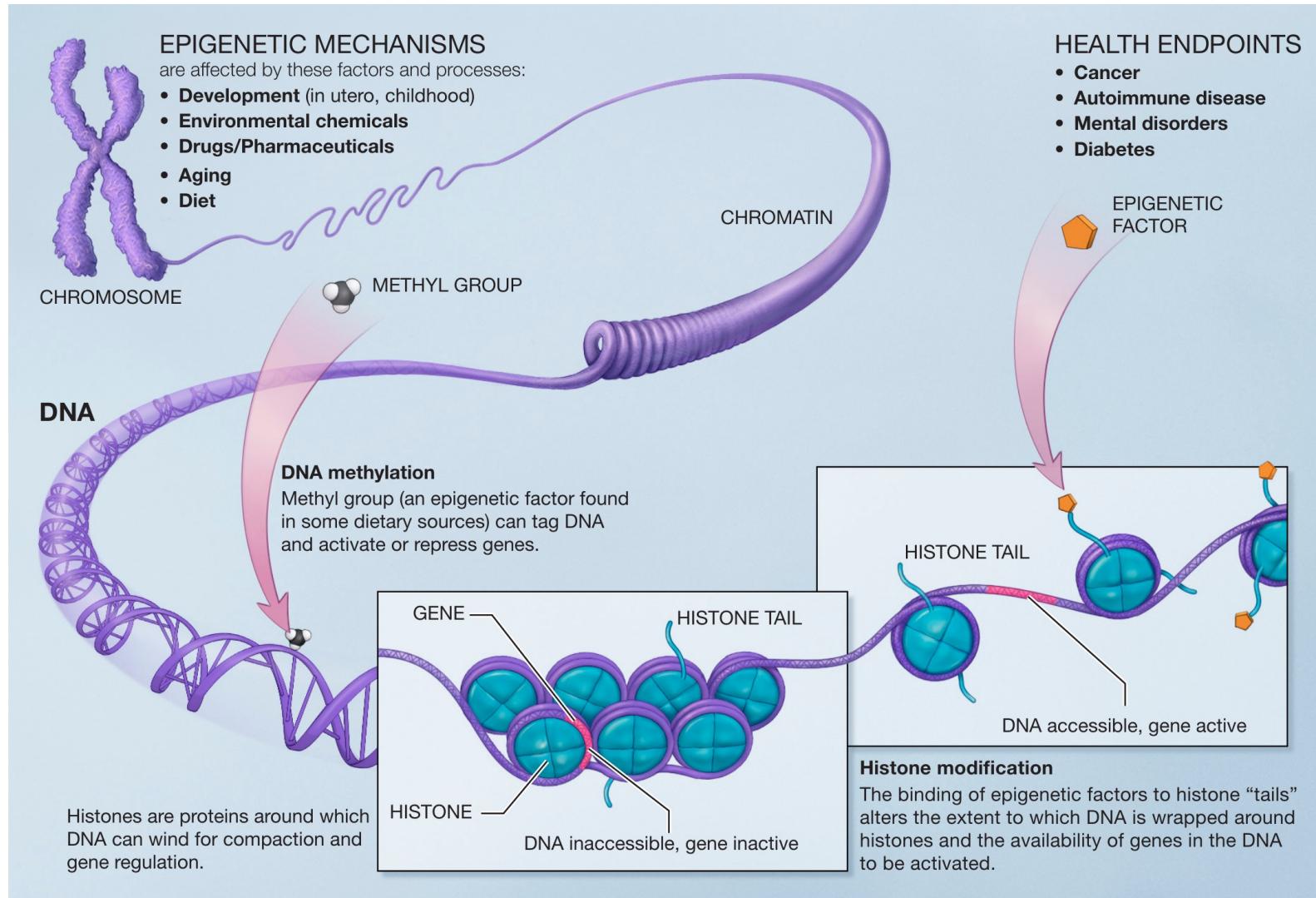
30,000,000,000,000 in human body

# Bird's eye view of the human genome

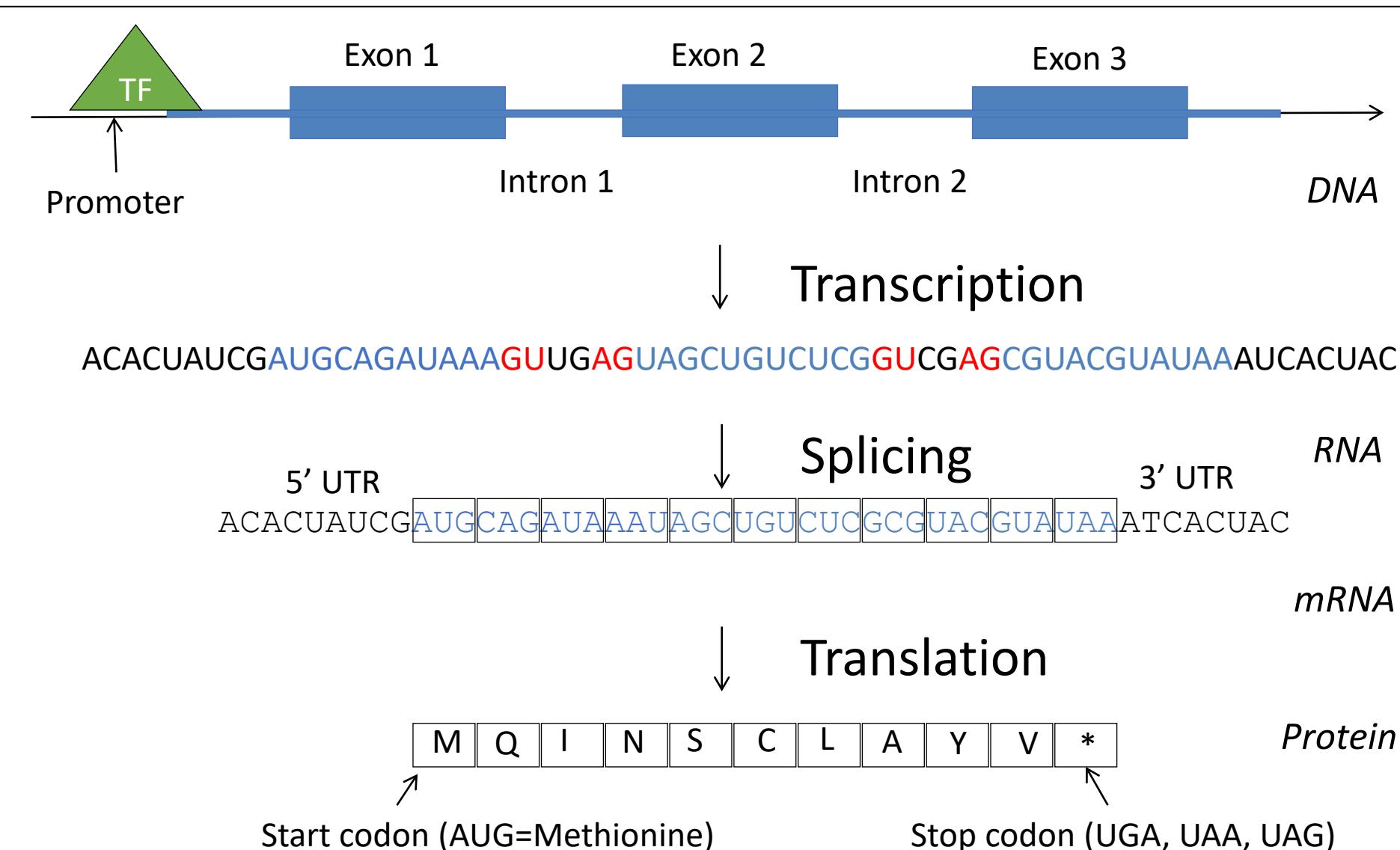
---



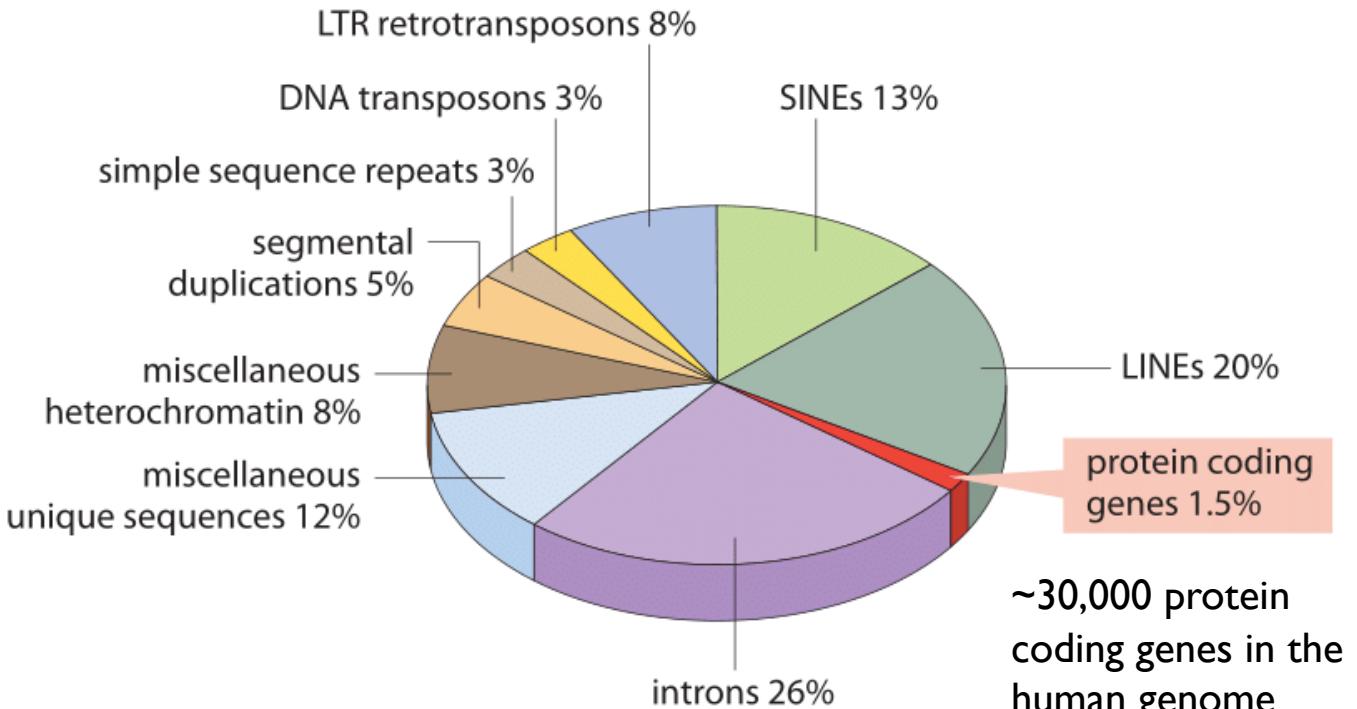
# Diverse epigenetic modifications



# The structure of a gene



# Organization of the human genome



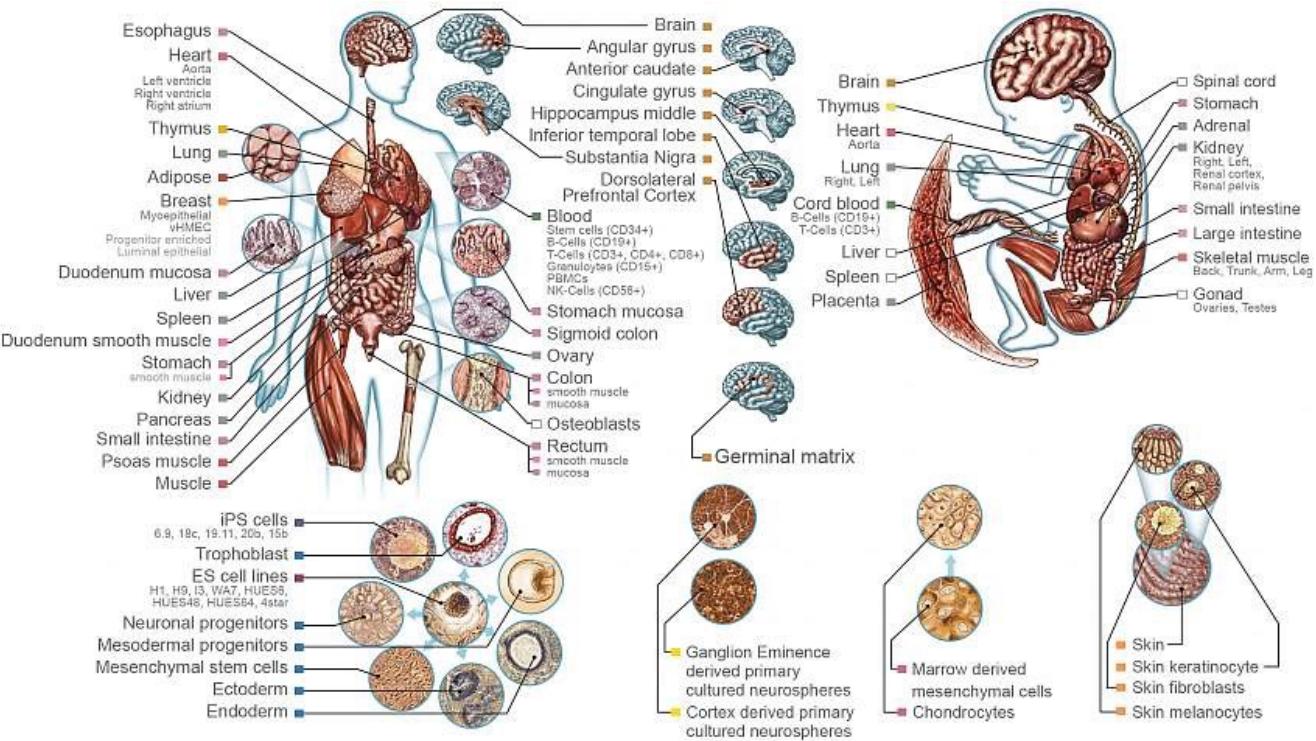
<http://book.bionumbers.org/how-many-genes-are-in-a-genome/>



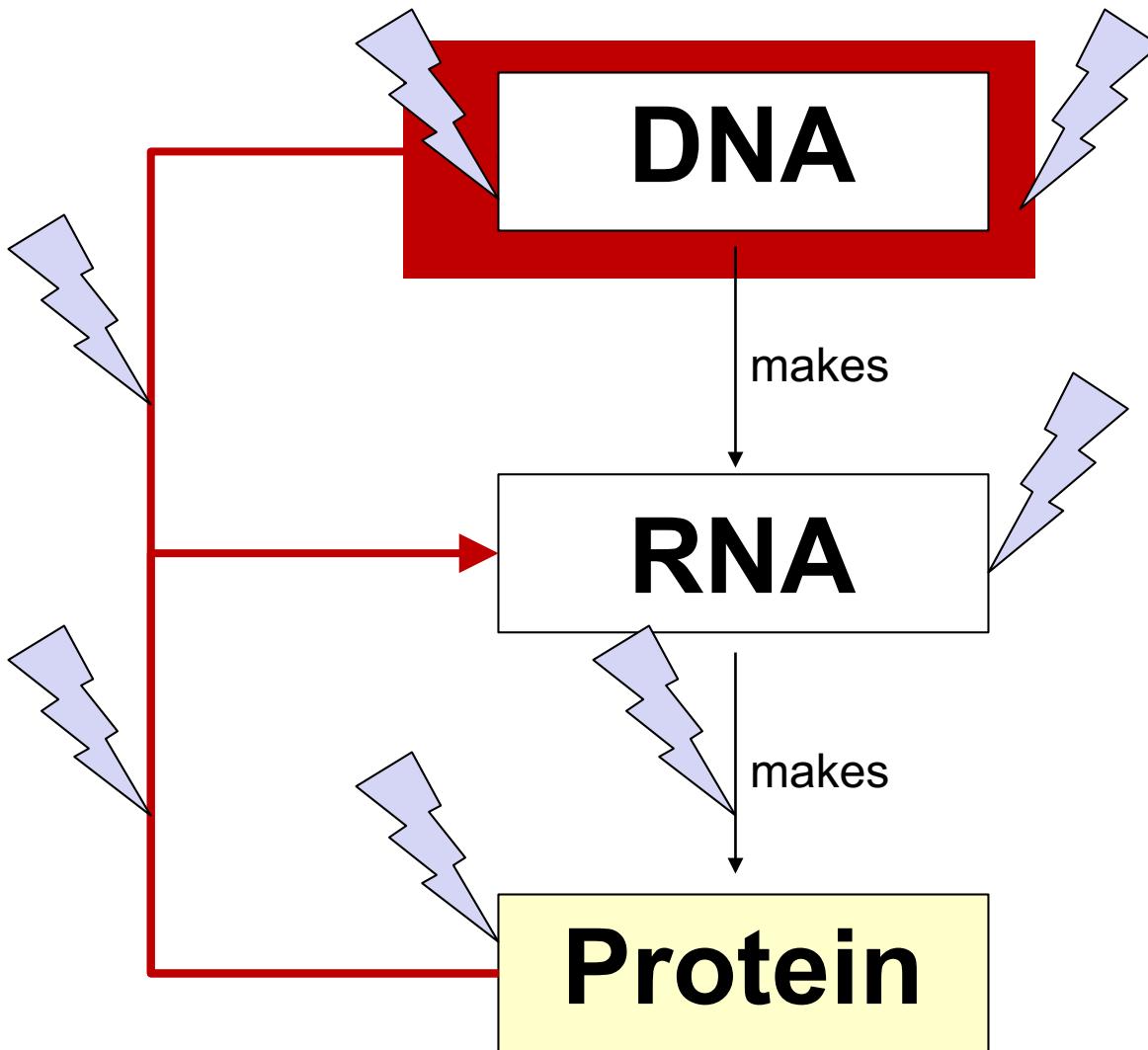
<http://www.biocomicals.com>

# One genome ↔ Many cell types

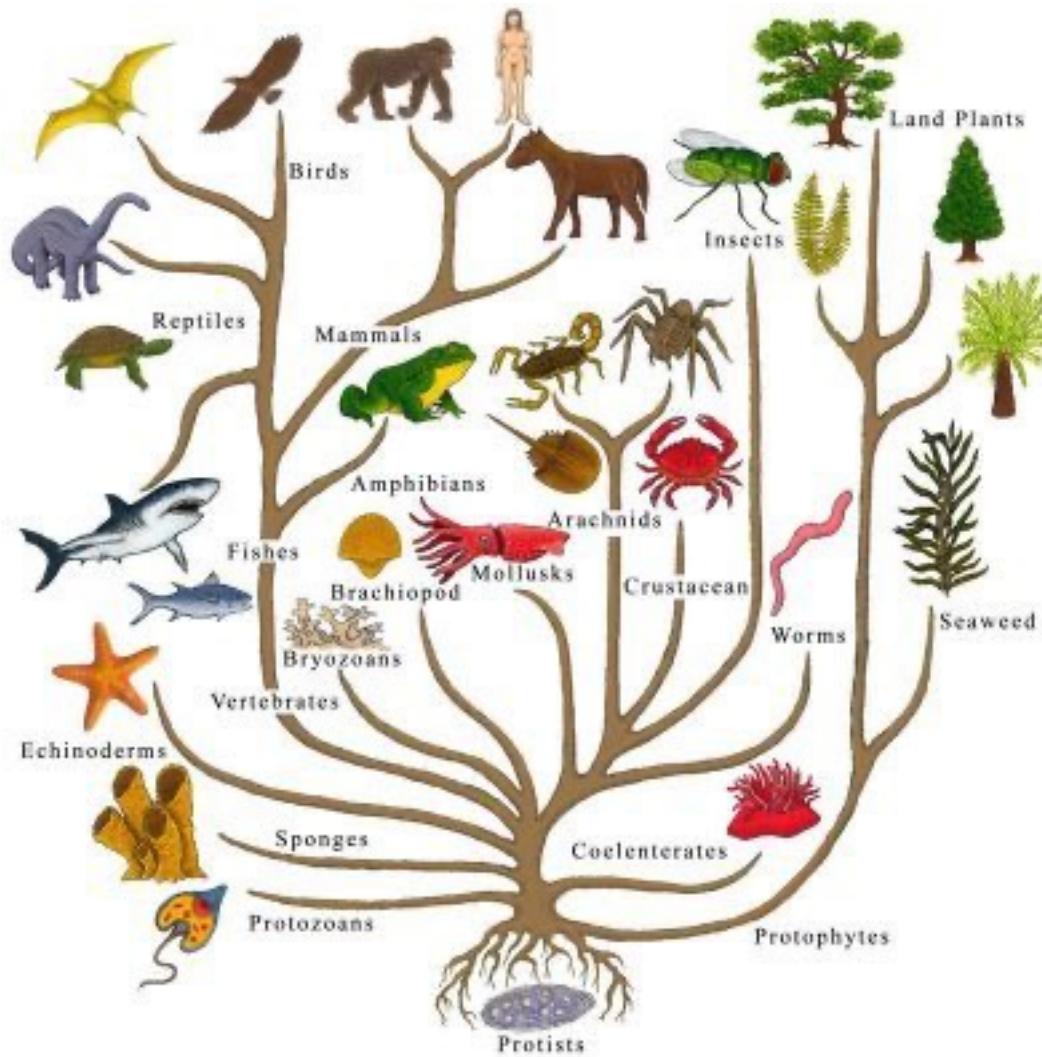
ACCAGTTACGACGGTCA  
 GGGTACTGATACCCCAA  
 ACCGTTGACCGCATTAA  
 CAGACGGGGTTGGGTT  
 TTGCCCCACACAGGTAC  
 GTTAGCTACTGGTTAG  
 CAATTACCGTTACAAC  
 GTTACAGGGTTACGGT  
 TGGGATTGAAAAAAAG  
 TTTGAGTTGGTTTTTC  
 ACGGTAGAACGTACCGT  
 TACCAAGTA



# The role of genetic alterations



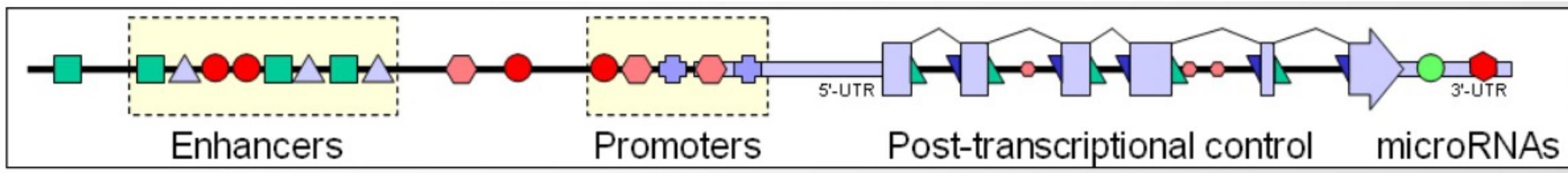
# Phylogenetic tree of life



(c)

	310	320	330	340	350
Human	GTACGTGAACGGGAAT	TGGGTGCCGGGGCAAGCCGAGCCG			
Monkey	GTACGTGAACGGGGAGT	GGGTGCCCGGGGGCAAGCCGAGCCG			
Fat-rumped sheep	GTACGTGAACGGGAGT	GGGTGCCCTTG	GGGGCAAGCCGAGCCG		
Sheep	GTACGTGAACGGGAGT	GGGTGCCCGGGGGCAAGCCGAGCCG			
Goat	GTACGTGAAT	GGGGAGTGGGTGCCCGGGGGCAAGCCGAGCCG			
Bovine	GTACGTGAAT	GGGGAGTGGGTGCCCGGGGGCAAGCCGAGCCG			
camel	GTACGTGAACGGCGAGT	GGGTGCCAGGGGGCAAGCCGAGCCG			T
Donkey	GTACGTGAACGGGAGT	GGGTGCCCGGGGGCAAGCCGAGCCG			
Horse	GTACGTGAACGGGAGT	GGGTGCCCGGGGGCAAGCCGAGCCG			
Dog	GTACGTGAACGGAGAGT	GGGTGCCAGGGGGCAAGCCGAGCCC			
Cat	GTATGTGAACGGAGAGT	GGGTGCCAGGGGGCAAGCCGAGCCC			
Mouse	ATATGTGAACGGGAGT	GGGTACCTGGGGCAAACCAAGAGCCT			
Rat	GTATGTGAAT	GGGGAGTGGGTACCTTG	GGGGCAAACCAAGAGCCT		
Pig	GTACGTGAACGGGAGT	GGGTGCCAGGGGGCAAGCCGAGCCG			
Chicken	GTACGTCAACGGCGAGT	GGGTGCCGGCGGCAAGCCGAGCCG			
Zebrafish	ATACGTGAACGGTGAAT	GGGTGCCCGGTGGGAAACCCGAACCC			

# The components of genomes and gene regulation



**Goal: A systems-level understanding of genomes and gene regulation:**

- The genome: Map reads, align genes/genomes, assembly strategies
- The genes: Protein-coding exons, introns, non-coding RNA, RNA folding
- The control regions: Promoters, enhancers, insulators, chromatin states
- The actual words: Regulatory motifs, high-resolution accessibility maps
- The regulators: Transcription factors, chromatin modifiers, nucleosomes
- The dynamics: Changing maps between cell types, across development
- The networks: regulator→enhancer→target, ChIP-seq, correlated activity
- The grammars: TF/motif/mark combinations, predictive models
- Human variation: Human diversity, population genomics, linkage maps
- Evolution: Phylogenetics, phylogenomics, coalescent, human ancestry
- GWAS/QTLs: Genome variation ↔organismal/molecular phenotypes
- Disease: Personal (epi)genomics, pharmacogenomics, synthetic biology



Donald Knuth

Professor emeritus of Computer Science  
Stanford University  
Turing Award winner  
“father of the analysis of algorithms.”

*“I can’t be as confident about computer science as I can about biology. Biology easily has 500 years of exciting problems to work on. It’s at that level.”*