

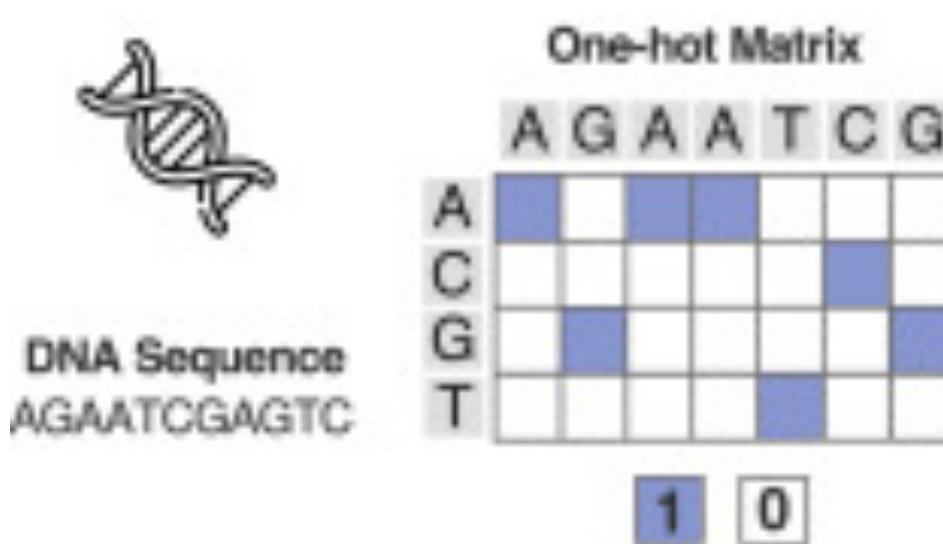
Bioinformatics for Beginners

Introduction to Bioinformatics: An Overview

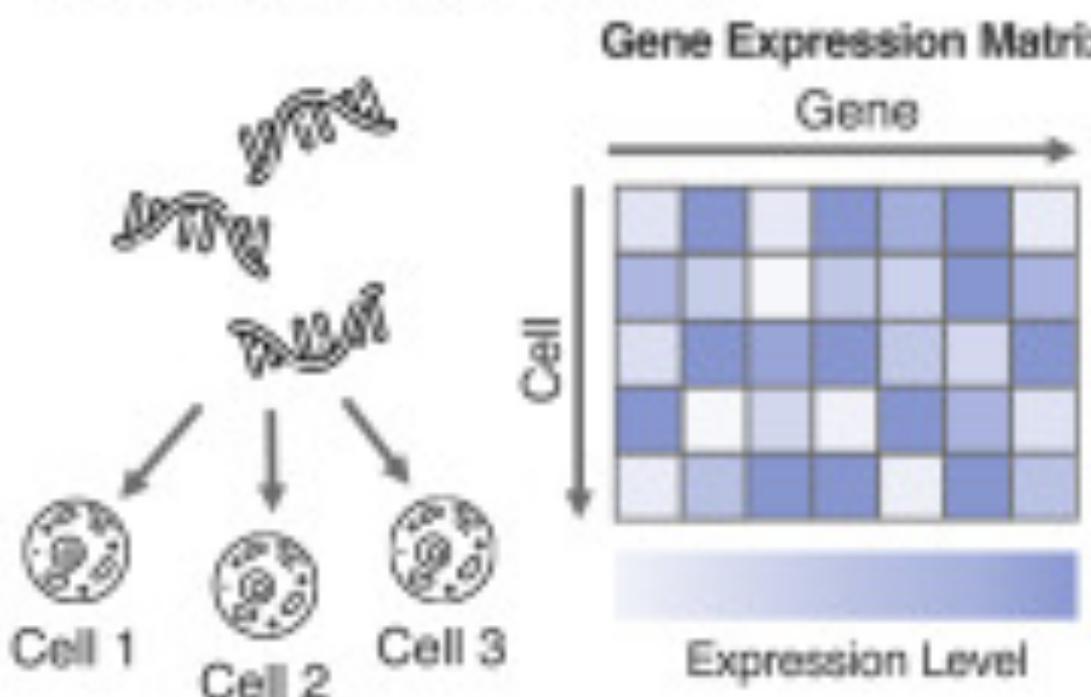
Tugce Bilgin Sonay, ZHAW, May 2023

Bioinformatics

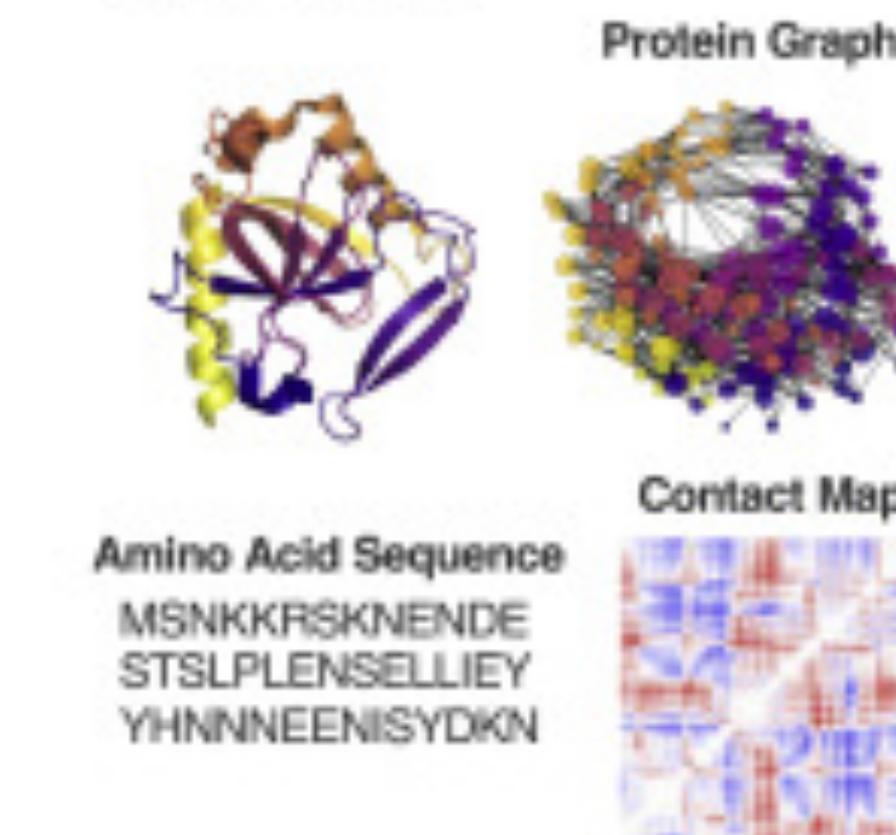
A DNA



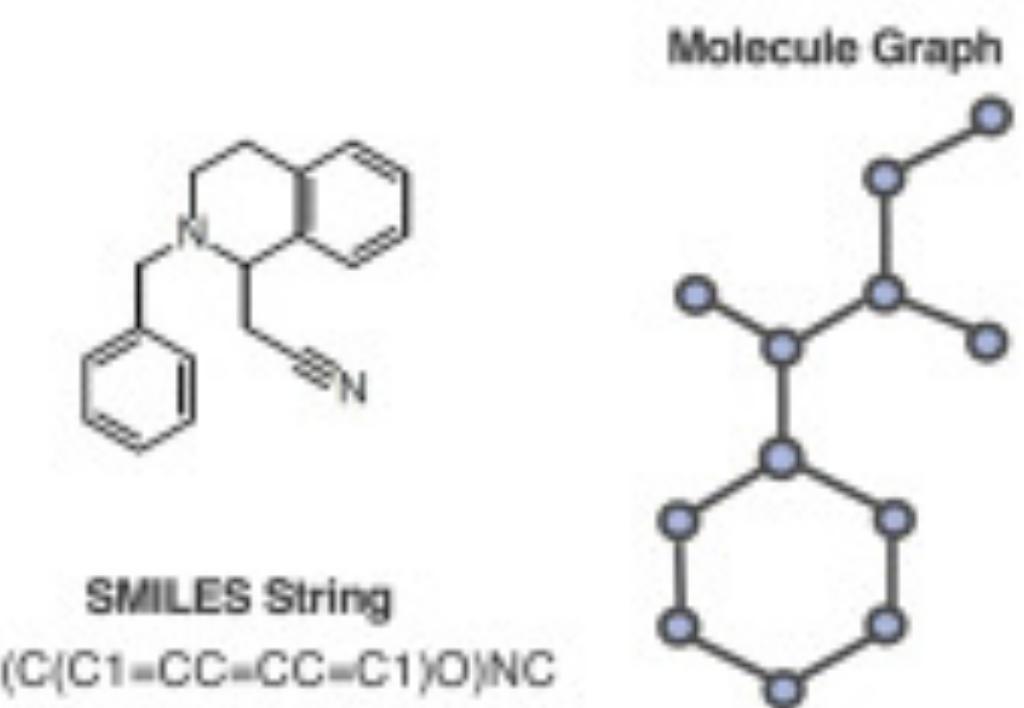
B Gene Expression



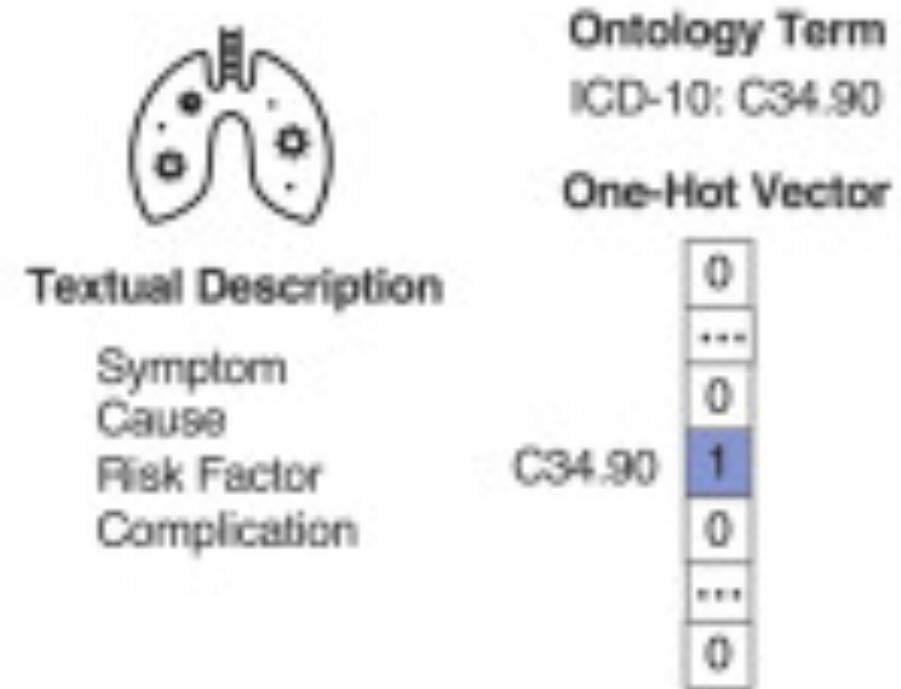
C Protein



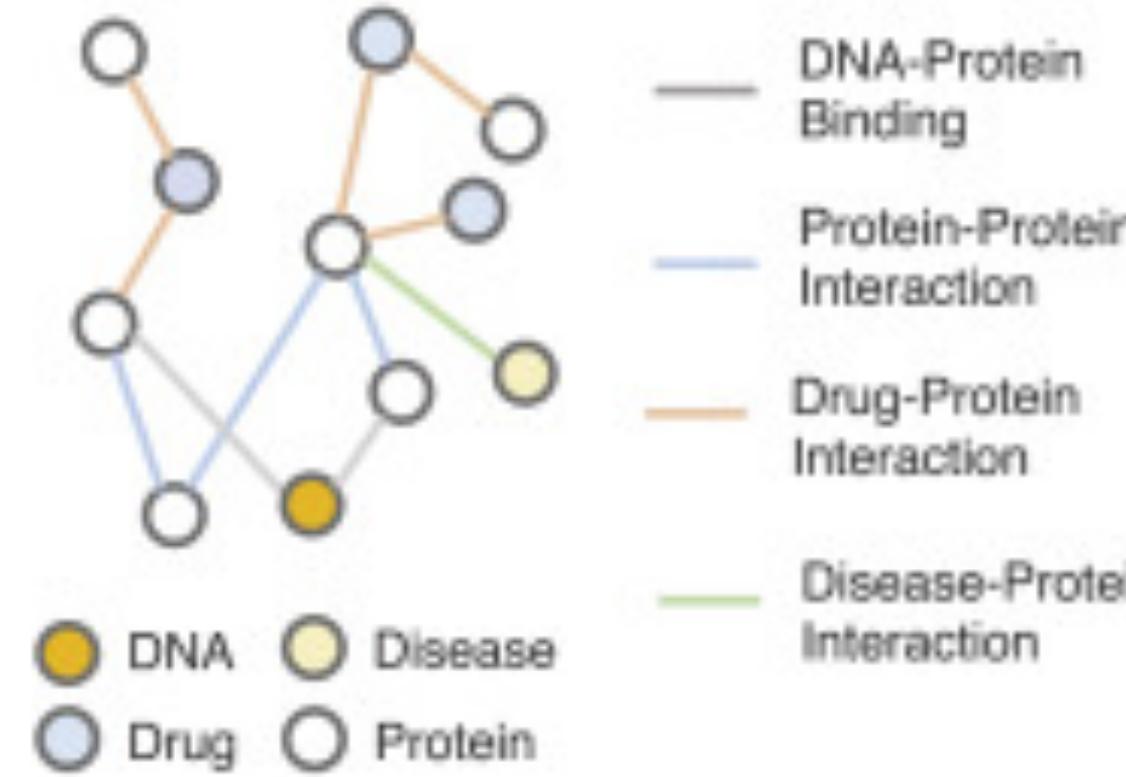
D Compound



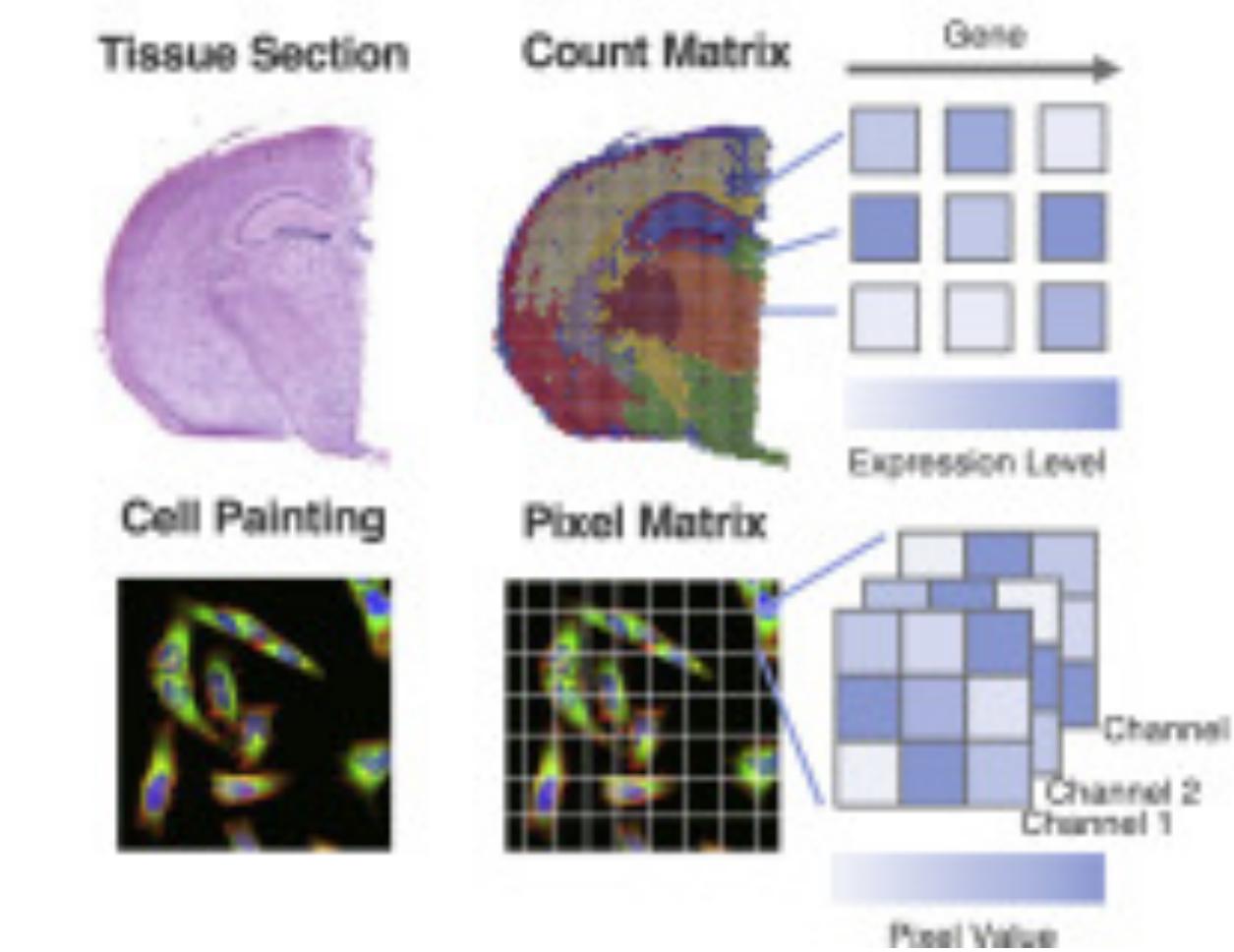
E Disease



F Biomedical Networks



G Spatial Data



H Text



Phenotypic Elasticity



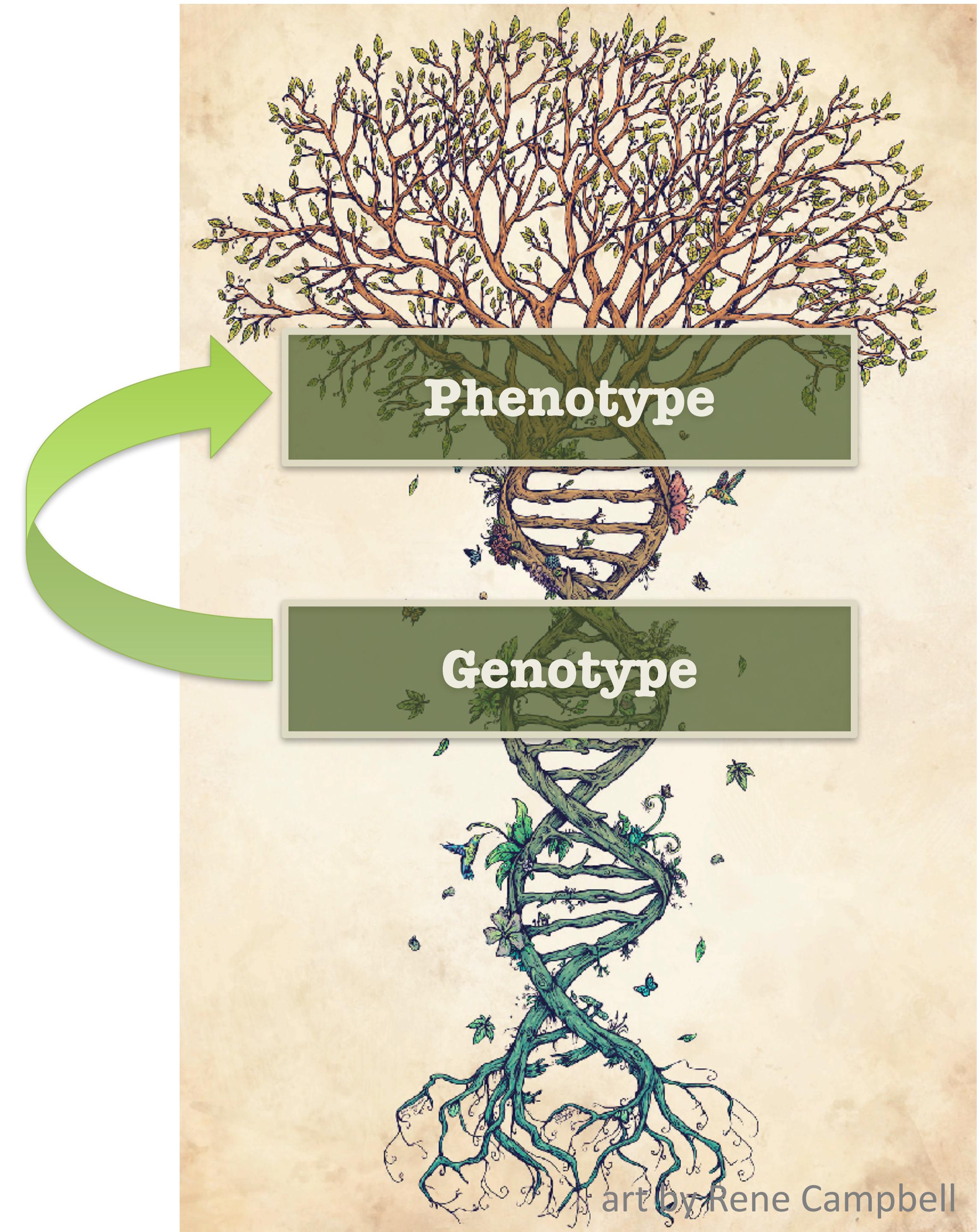
17° C



25° C

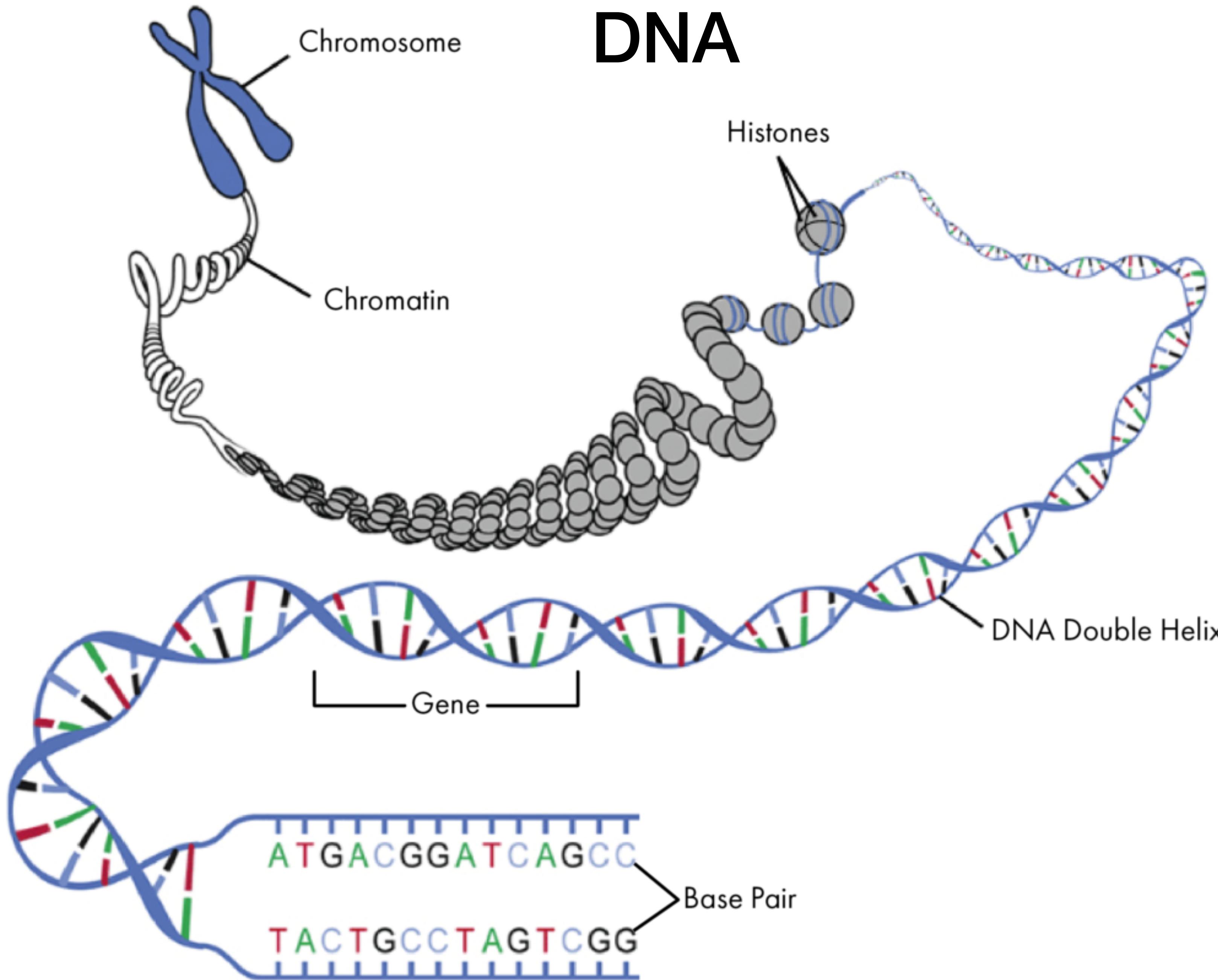


28° C

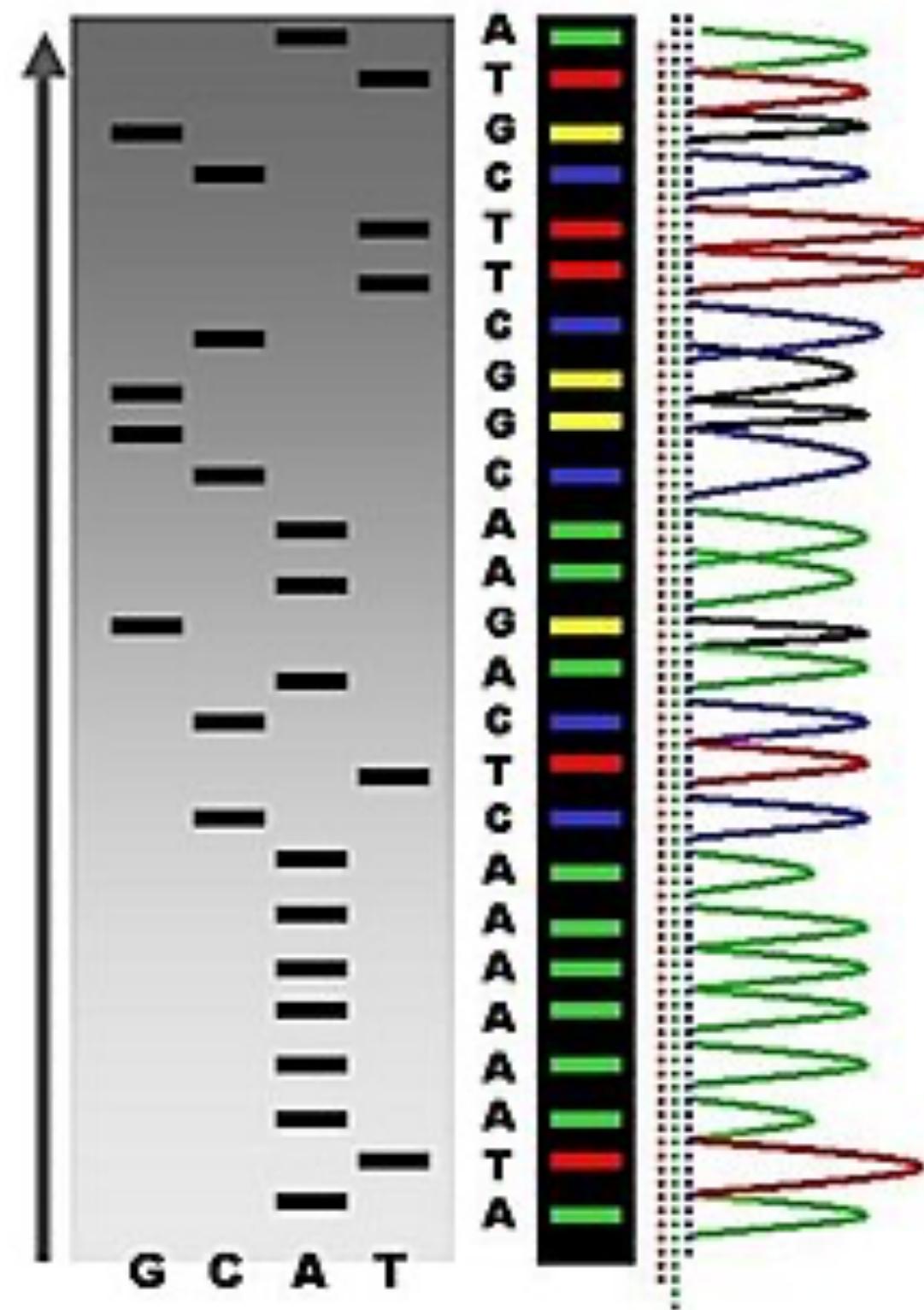


art by Rene Campbell

DNA



Sanger Sequencing



- 1977: First bacterial genome
- 1986: Genes for color blindness
- 1987: Sequence of the CRISPR
- 1989: Gene involved in cystic fibrosis
- 1994: Genes involved cancer
- 1995: Genome of *H.influenza*
- 1996: First eukaryotic genome
- 1998: First multicellular genome

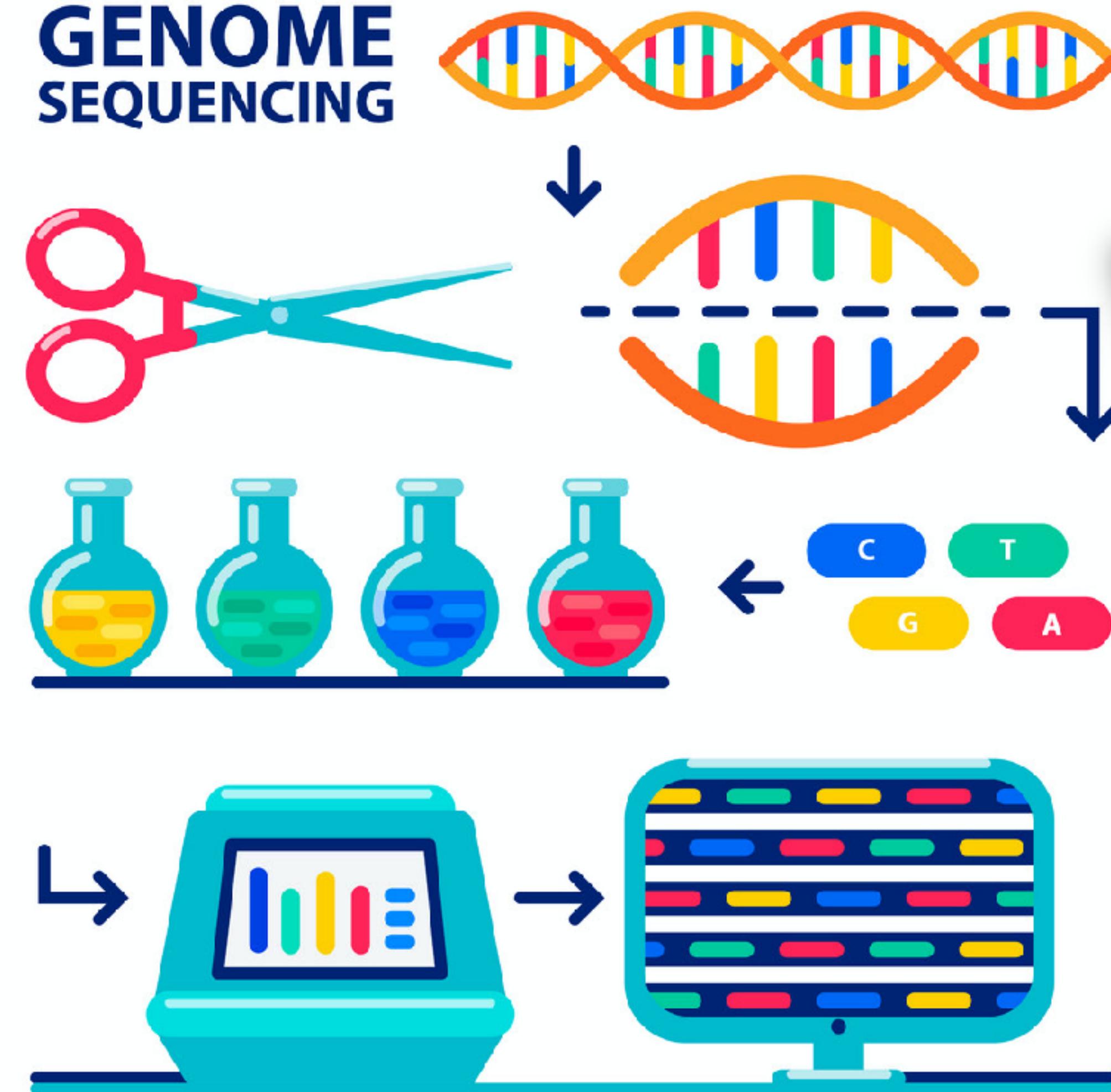


Human Genome Project

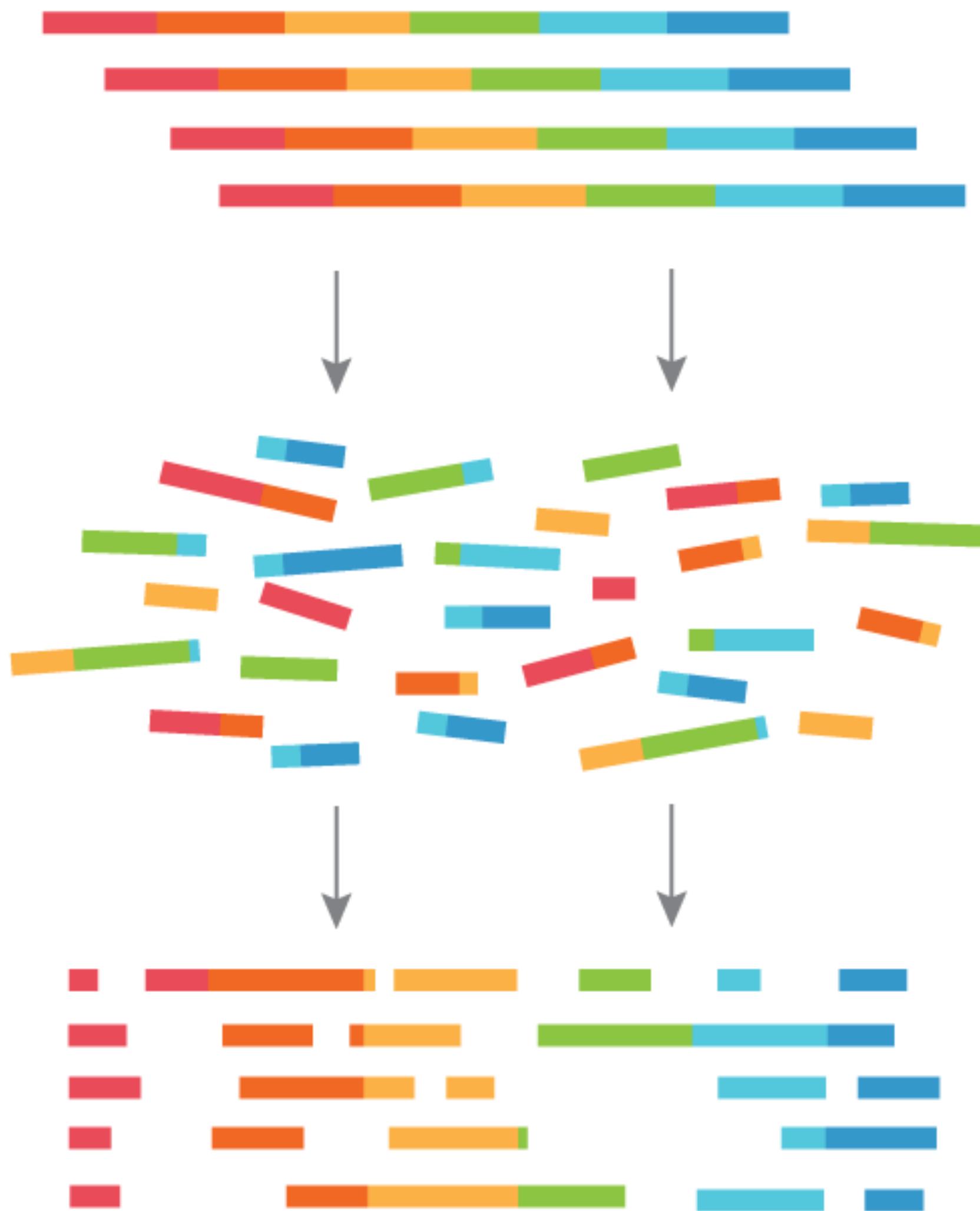
1990–2003

13 years
\$3 billion
20 countries

GENOME SEQUENCING



Next Generation Sequencing



ATGTTCCGATTAGGAAACCTATCTGTAACGTGTTCAATTCAAGTAAAAGGGAGGAAA

A few important findings of the HGP

- ✓ Ethical breakthrough: instant data sharing
- ✓ data analysis and storage technologies
- ✓ We have ~20,000 genes
- ✓ Our chromosomes have different lengths
 - ~ Most of the human genome is not coding proteins
 - ~ sequence is NOT equal to function

DNA Sequencing Techniques

Next Generation Sequencing

\$3,000,000,000 | 2003 Human Genome Project



\$20,000,000 | 2006 1st individual genome



\$2,000,000 | 2007 1st NGS Genome



\$200,000 | 2008 1st 30x genome



\$10,000 | 2010 1st sub-10K genome



\$1,000 | 2014 1st \$1,000 genome



DNA Sequencing Techniques

Next Generation Sequencing

\$3,000,000,000



\$20,000,000



\$2,000,000



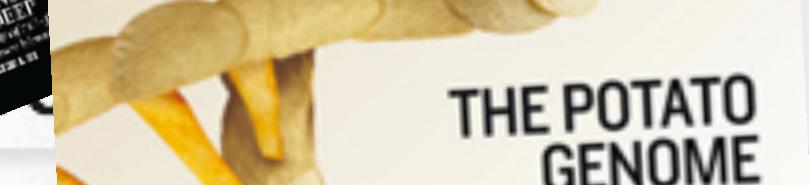
\$200,000



\$10,000



\$1,000



If you printed all your DNA
(6.4 billion letters), it would fill
4,200 books*

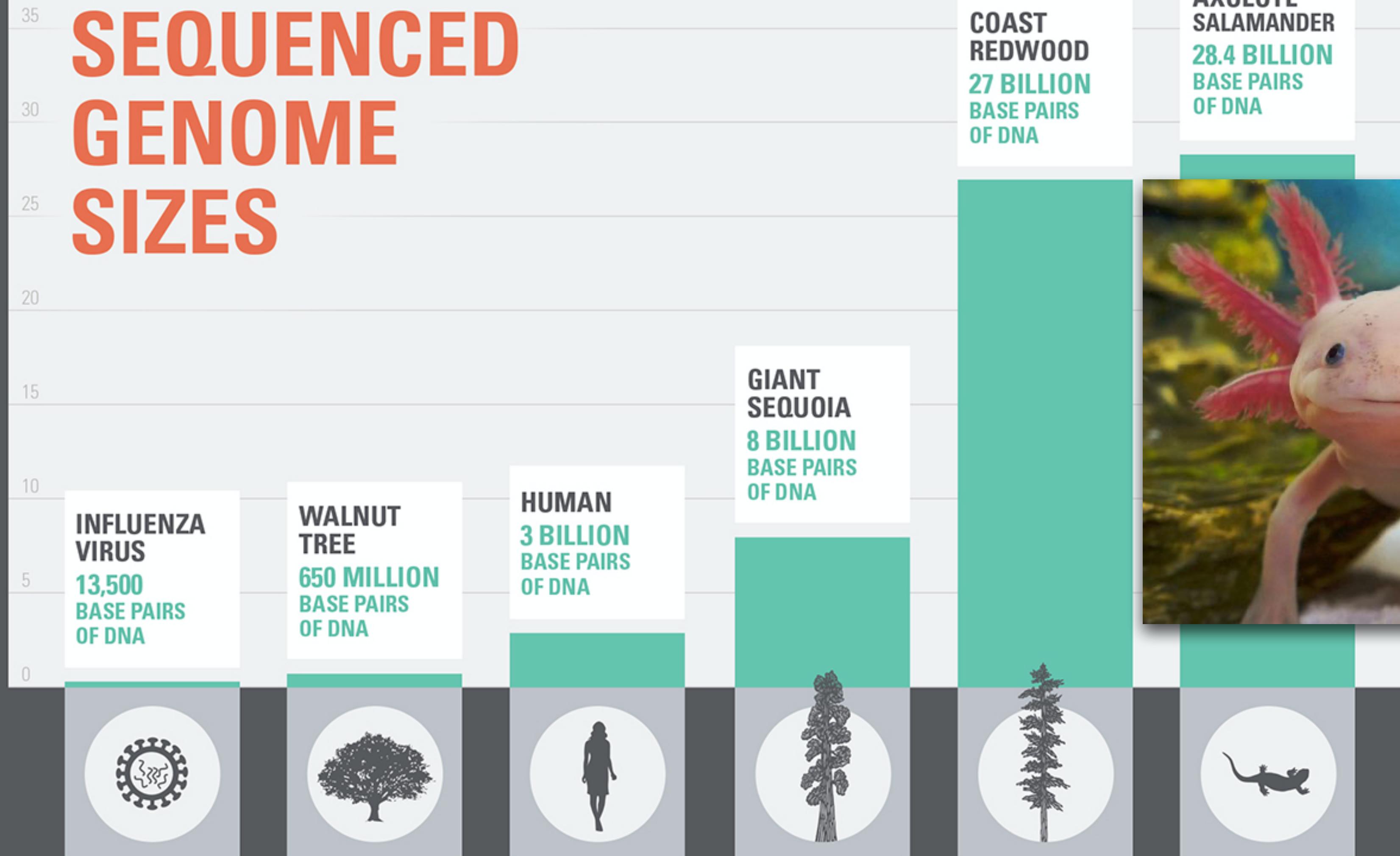


Veritas
sequences your
whole genome

(equivalent to
4,200 books)

* Assuming that a book, like Darwin's *Origin of Species*, has 500 pages.

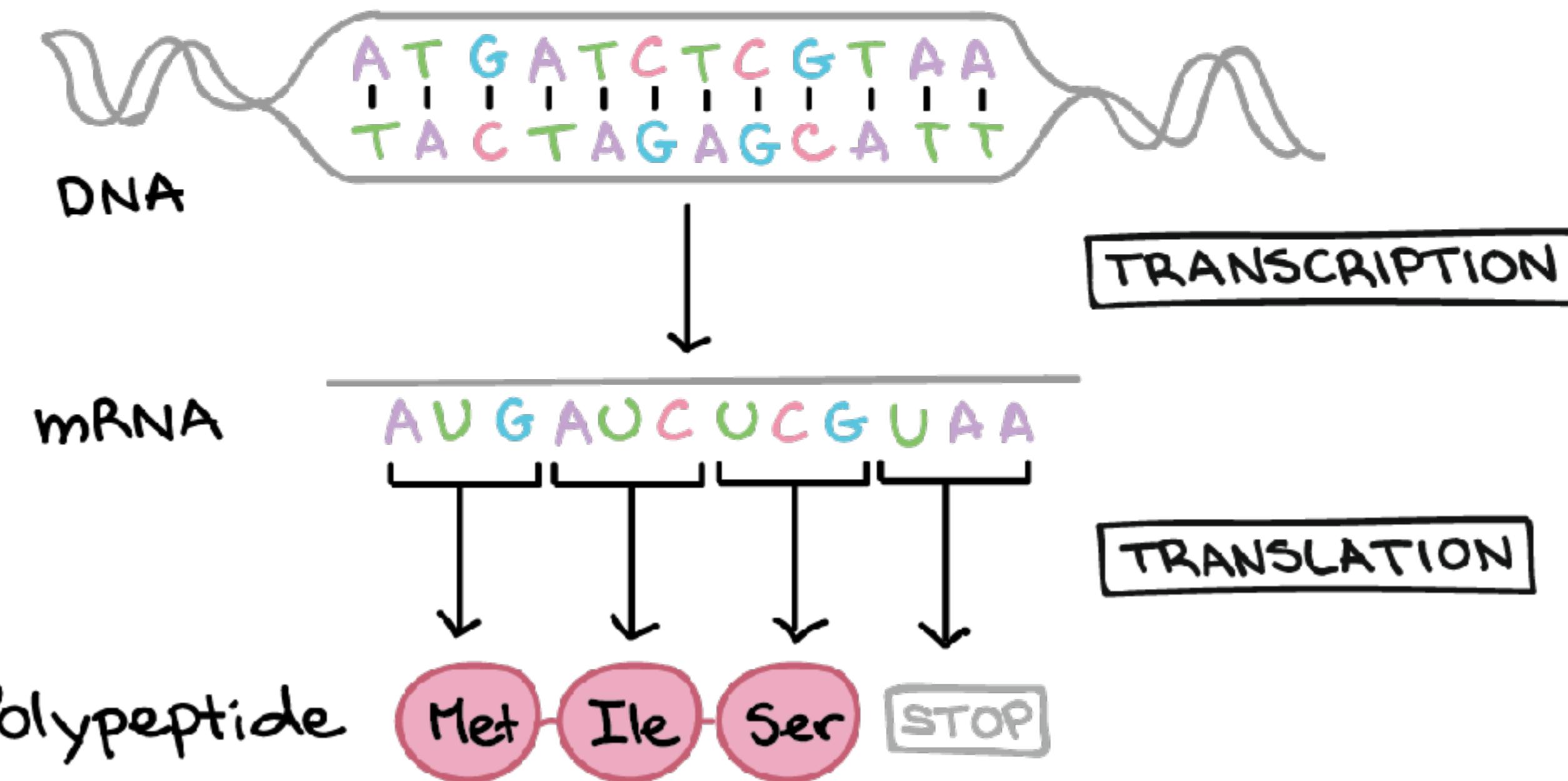
SEQUENCED GENOME SIZES



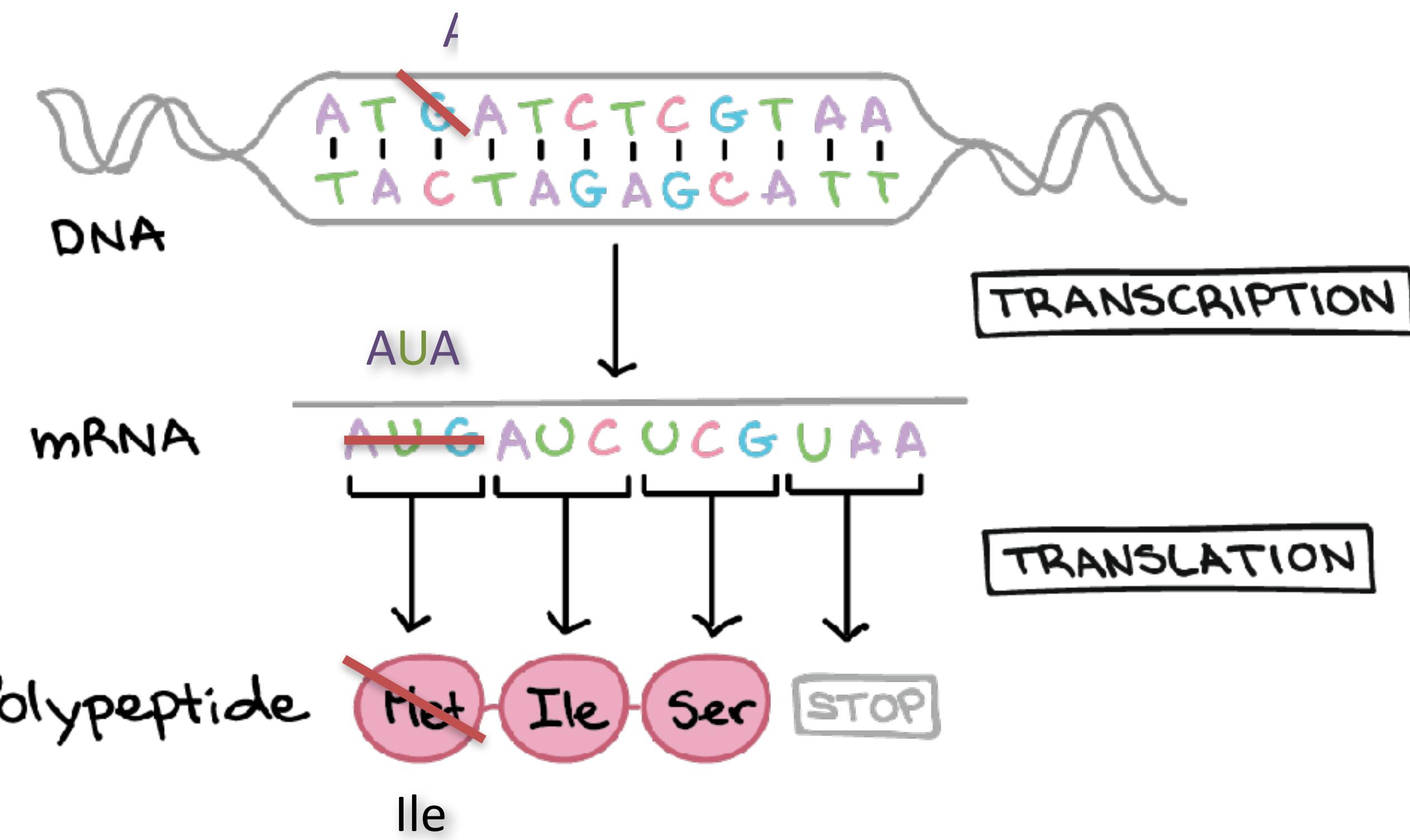
The coast redwood and giant sequoia genomes were sequenced by Save the Redwoods League, University of California, Davis, and Johns Hopkins University. Learn more at SaveTheRedwoods.org/RedwoodGenome.



THE CENTRAL DOGMA



THE CENTRAL DOGMA



11. AYKUT KENCE EVRİM KONFERANSI
18-19 ŞUBAT 2017
WWW.ODTÜAKEK.ORG



Dr. Tuğçe Bilgin Sonay

Sunum Konusu: DNA'nın en oynak elemanı: ardışık kısa tekrar düzleri ve fenotipik evrim

Yeditepe Üniversitesi Moleküler Biyoloji ve Genetik bölümünden mezun olan Dr. Tuğçe Bilgin Sonay, Zürich Üniversitesi'nde Evrimsel Biyoloji üzerine doktorasını bitirdi, yine aynı yerde araştırma görevlisi olarak çalışıyor. Evrim dilden psikoloji ve ikisinin kesişim noktaları olan bilincin ve empatinin ortaya çöküp gibi konulara ilgi duyuyor.

ODTÜ KKM - Kemal Kurdaş Konferans Salonu

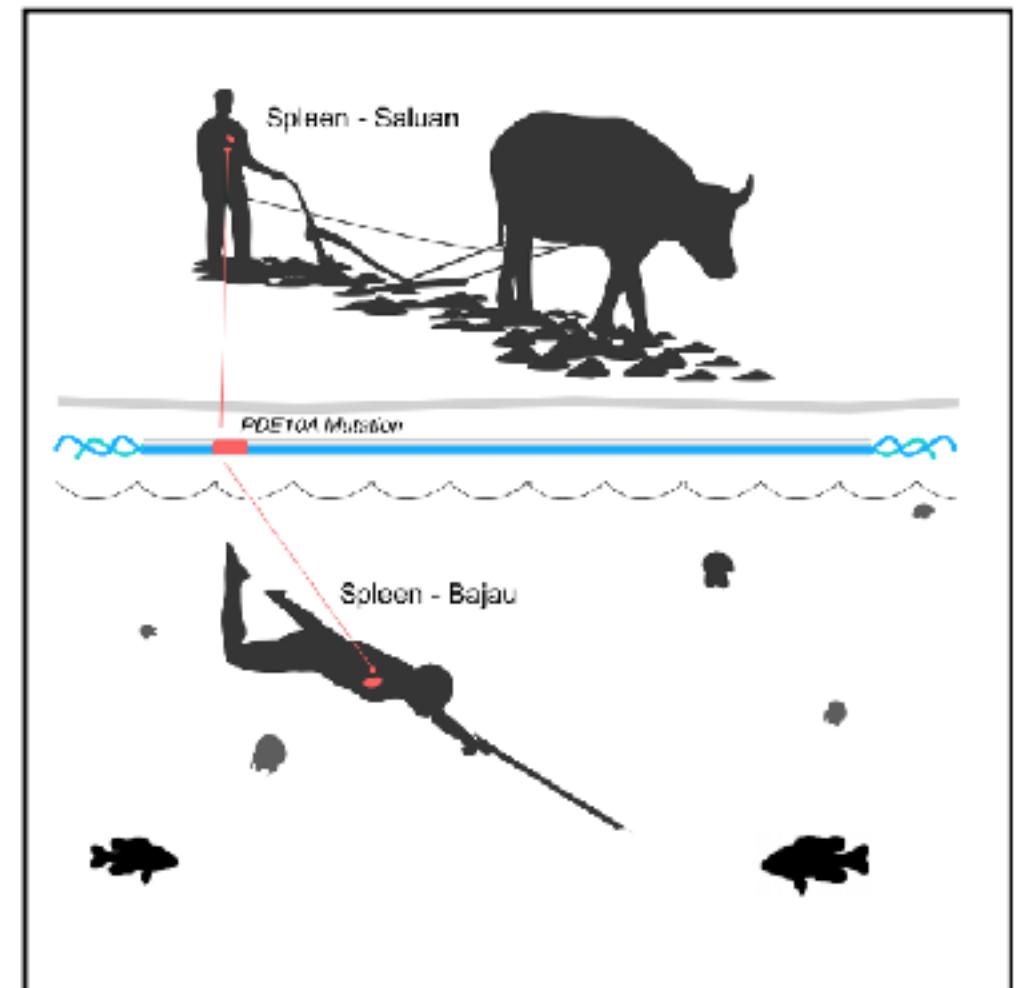




Cell

Physiological and Genetic Adaptations to Diving in Sea Nomads

Graphical Abstract



Authors

Melissa A. Ilardo, Ida Moltke,
Thorfinn S. Korneliussen, ...,
Suhartini Salingkat, Rasmus Nielsen,
Eske Willerslev

Correspondence

rasmus_nielsen@berkeley.edu (R.N.),
ewillerslev@snm.ku.dk (E.W.)

In Brief

Genetic and physiological adaptations enable the remarkable breath-holding ability of marine nomads.

Highlights

- The Bajau, or "Sea Nomads," have engaged in breath-hold diving for thousands of years
- Selection has increased Bajau spleen size, providing an oxygen reservoir for diving

Article





deniz
diyeti



sut
sindirim



malaryaya
dayanıklılık



koleraya
dayanıklılık



arsenikli
ortam



soguk
hava



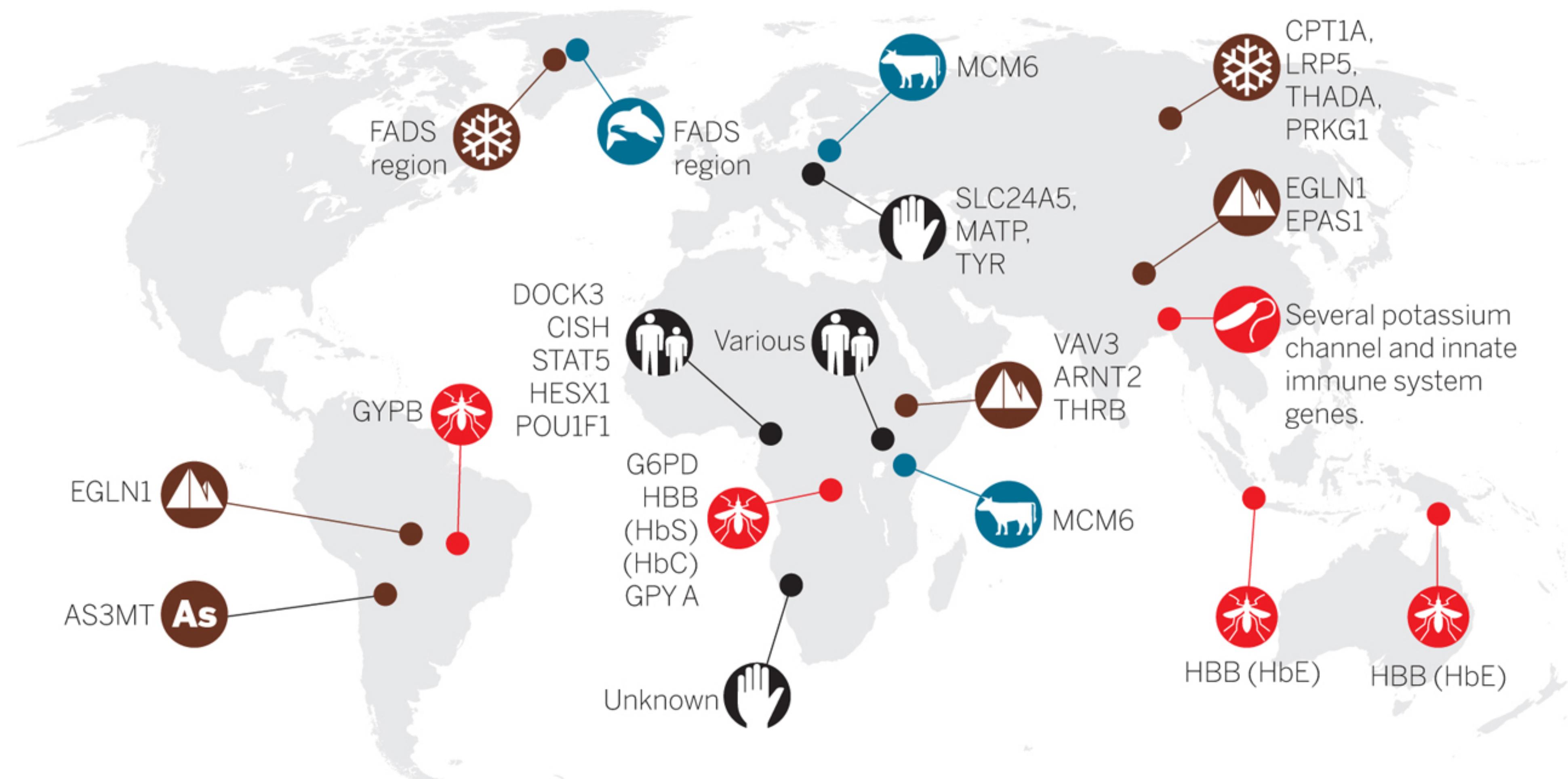
irtifa



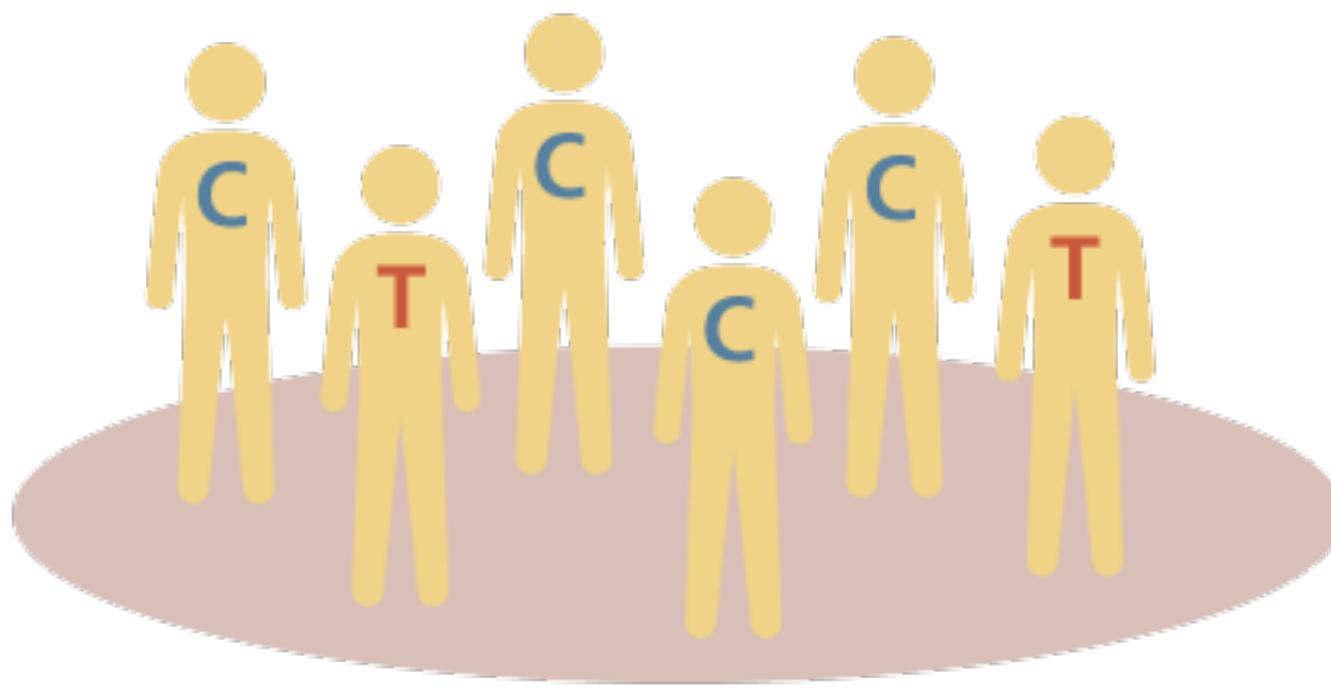
acik renk
ten



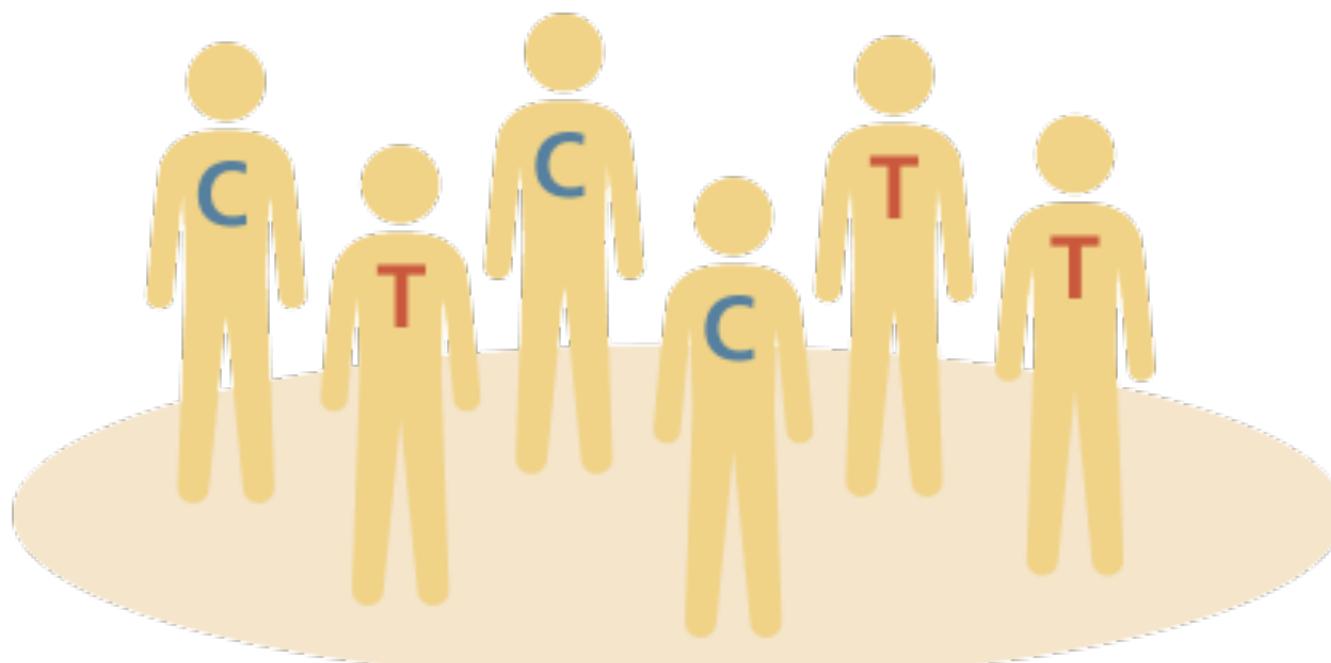
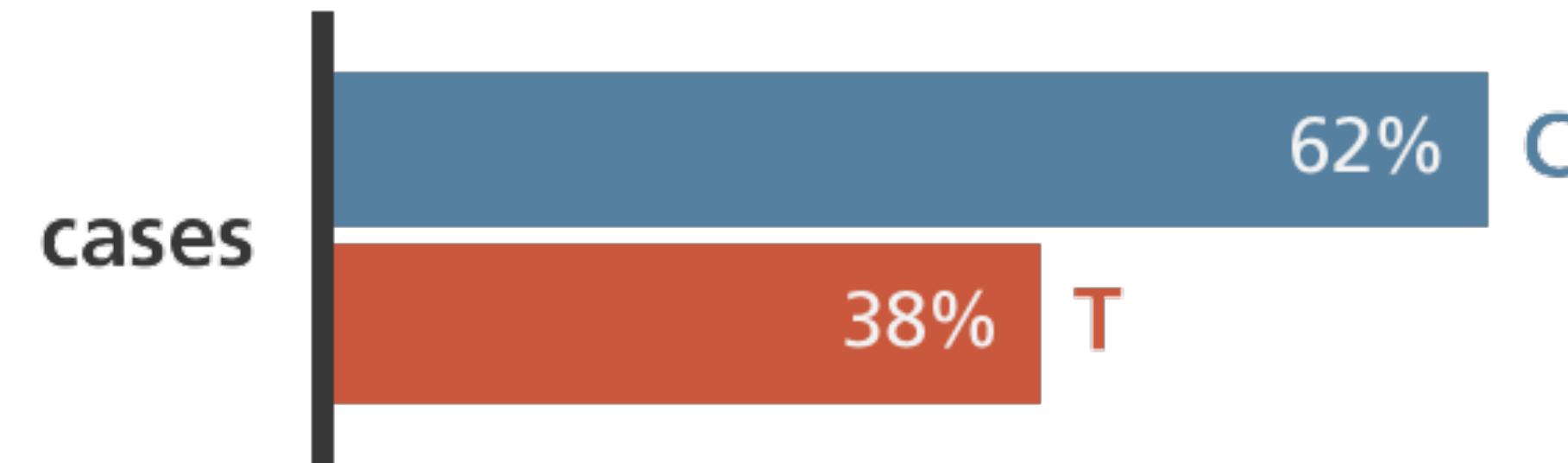
kisa
boy



Genome Wide Association Studies



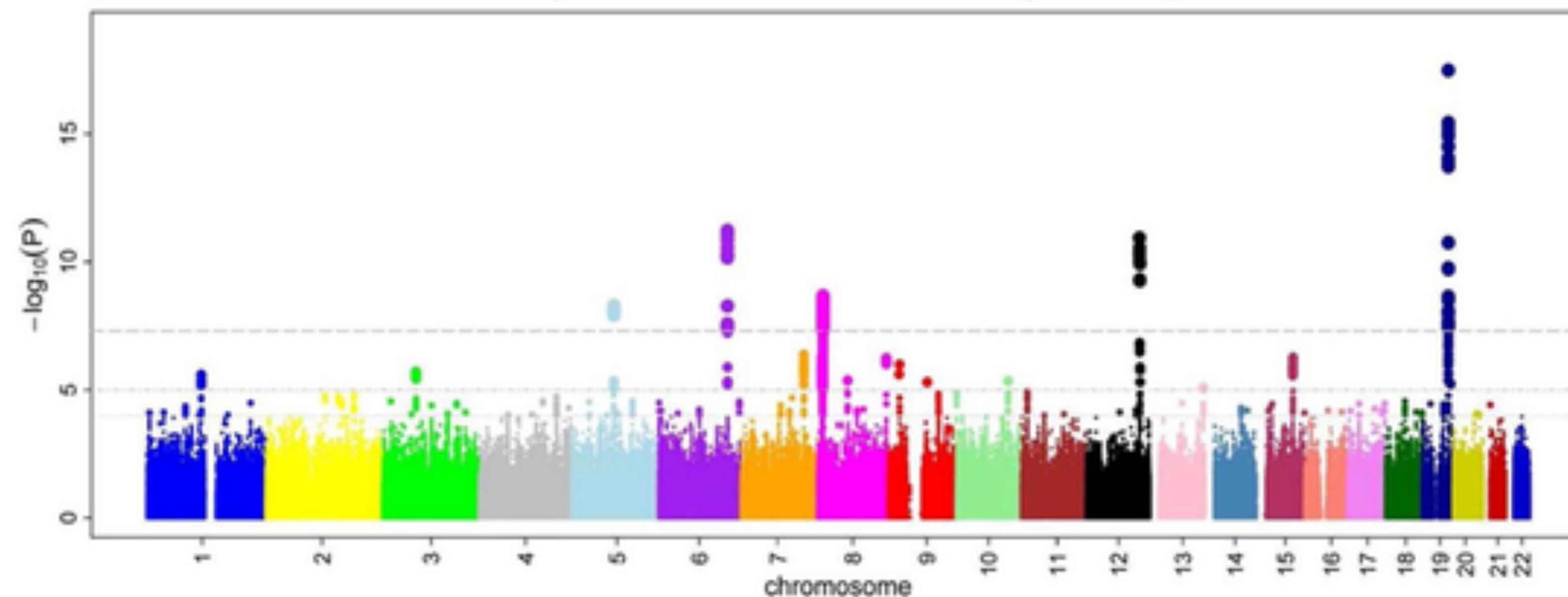
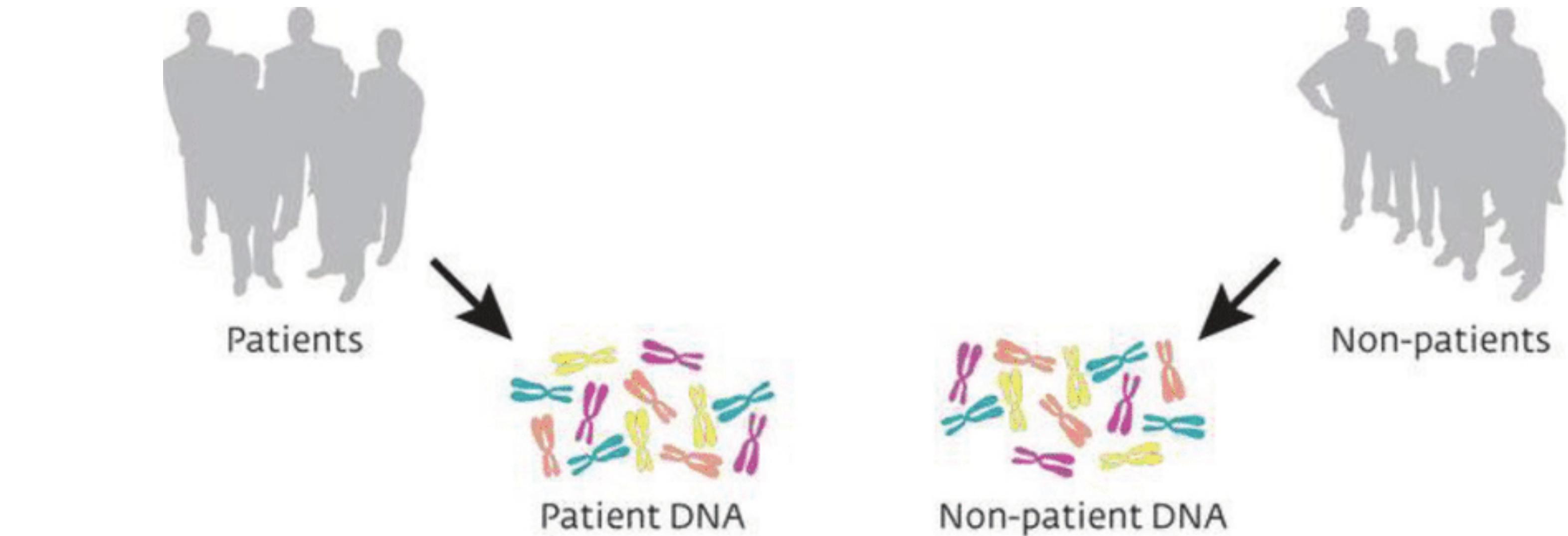
cases (n=1,000)
people with heart disease



controls (n=1,000)
people without heart disease



Genome Wide Association Studies



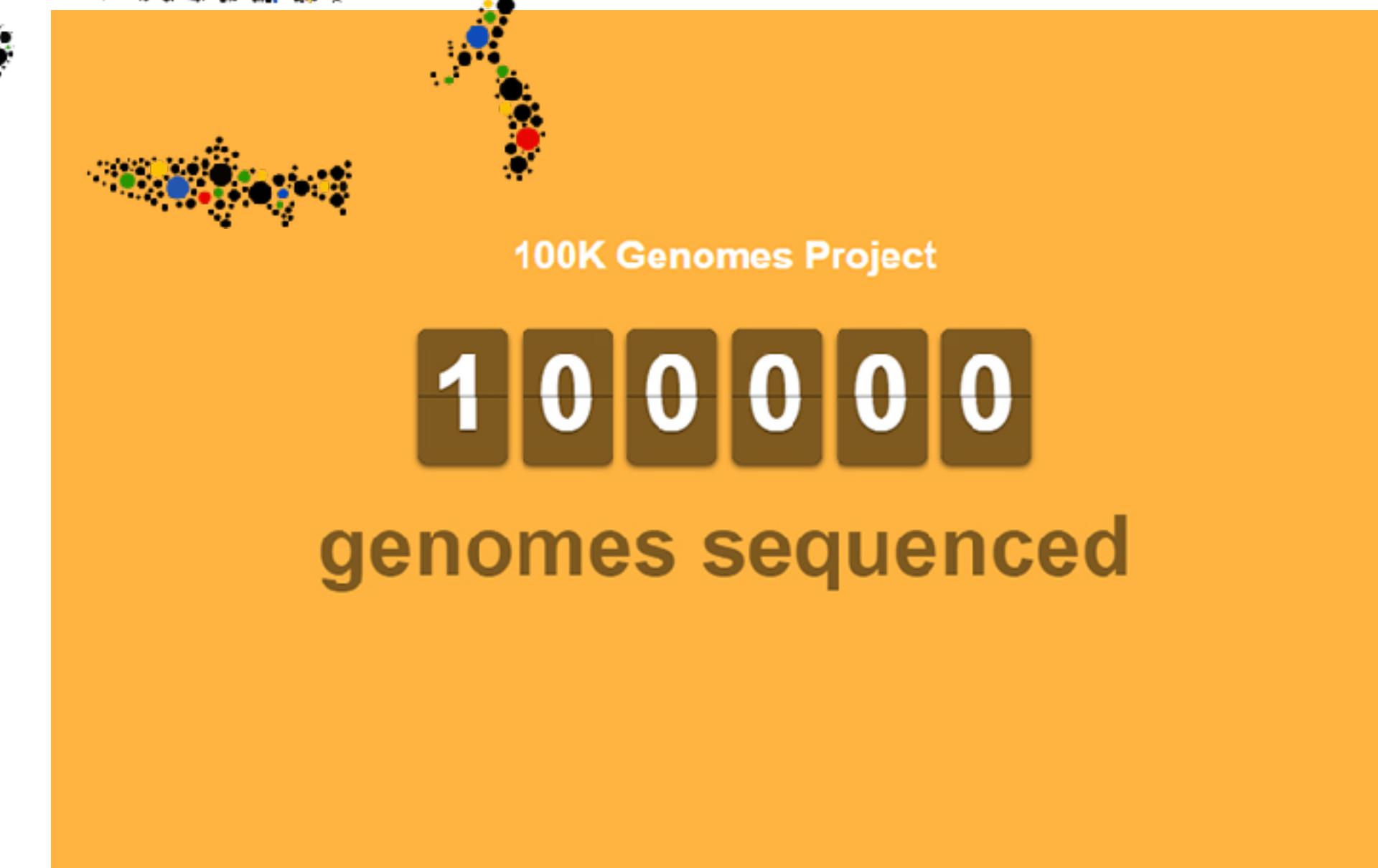
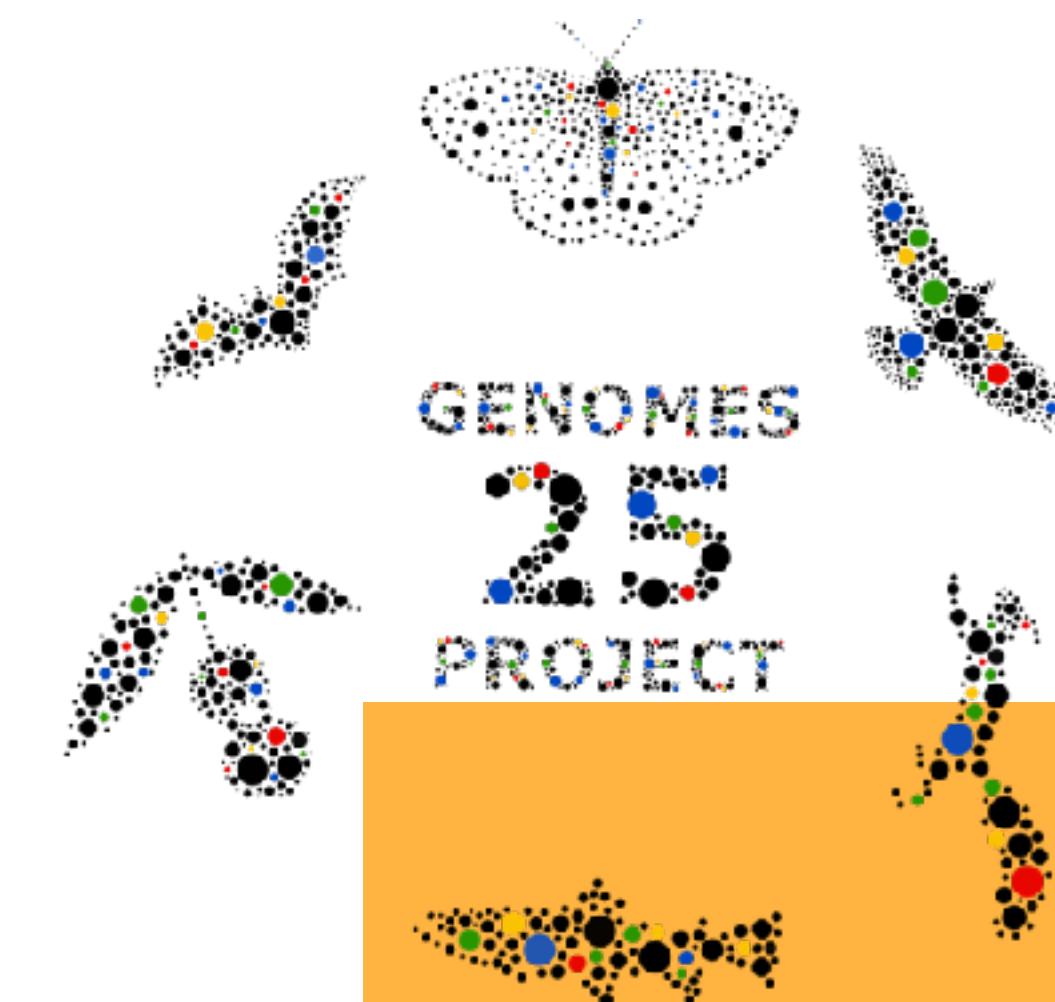
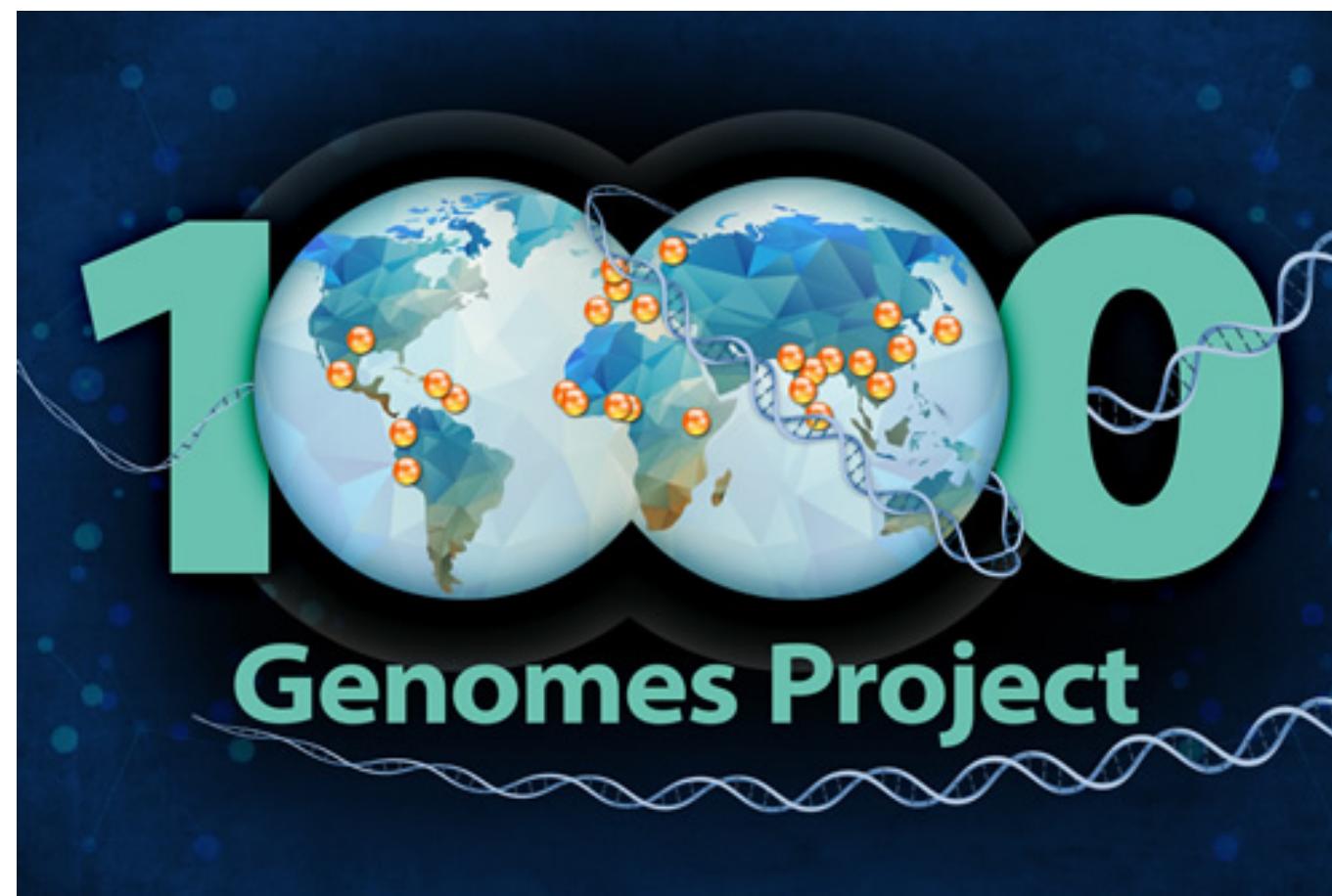
Genome Wide Association Studies



From 6734 people:
DNA sequence
Diseases, BMI, habits, ethnicity



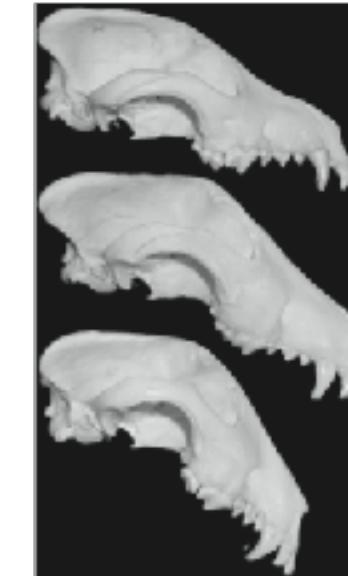
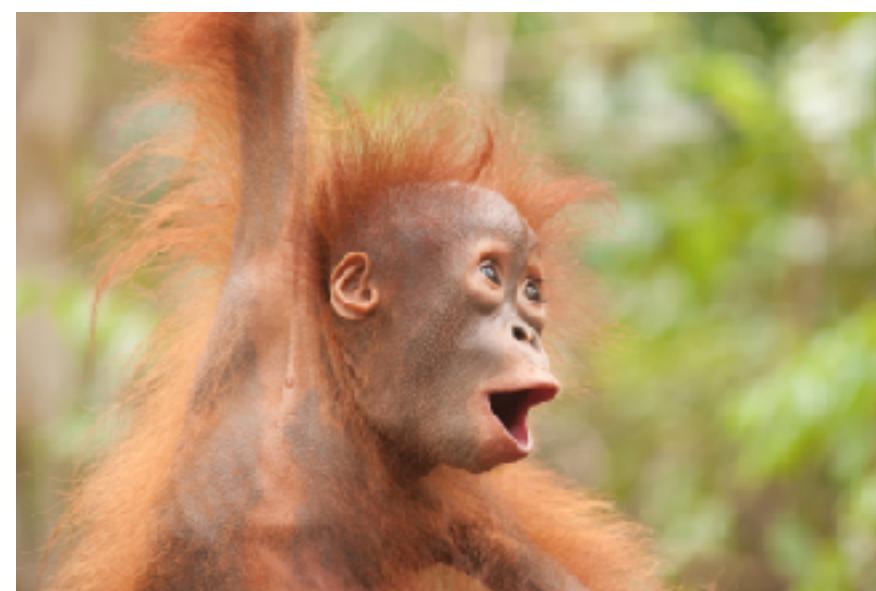
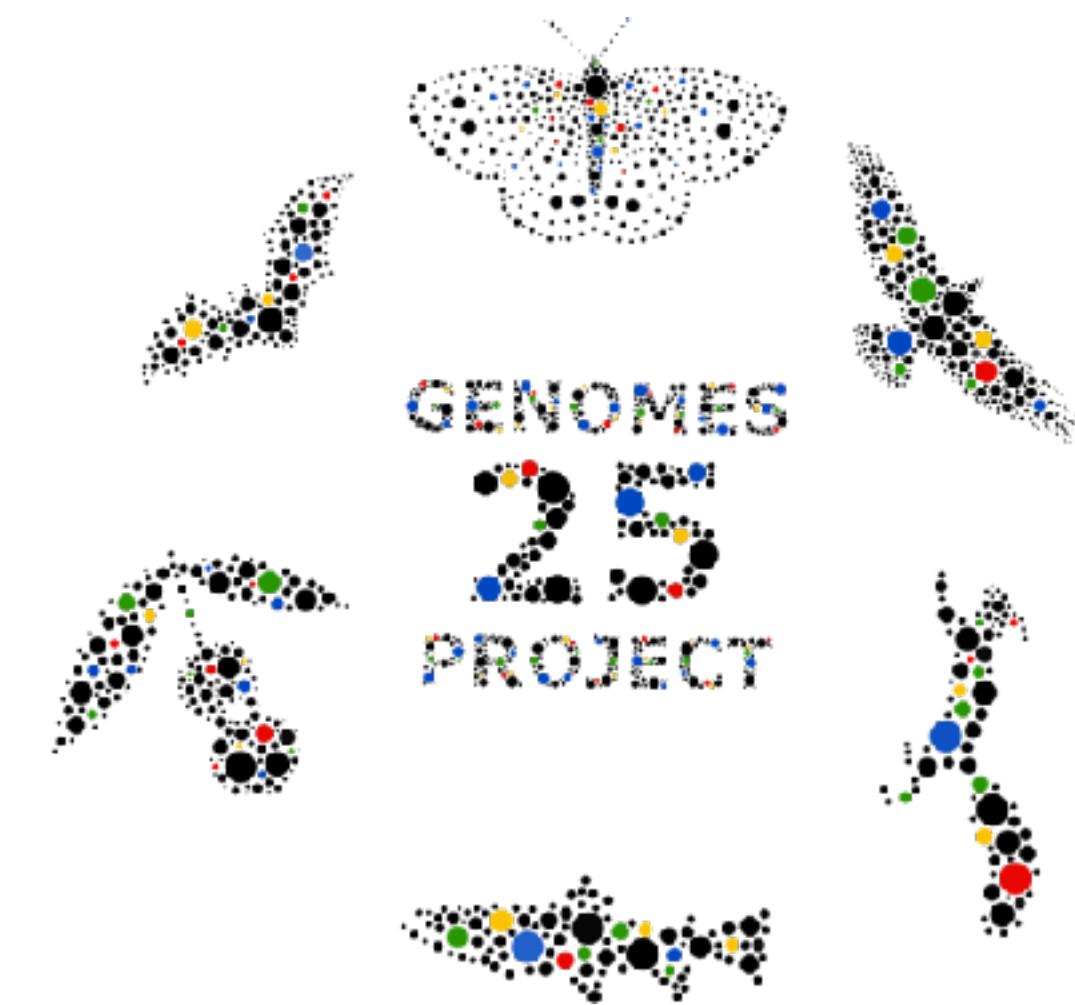
Genome Projects



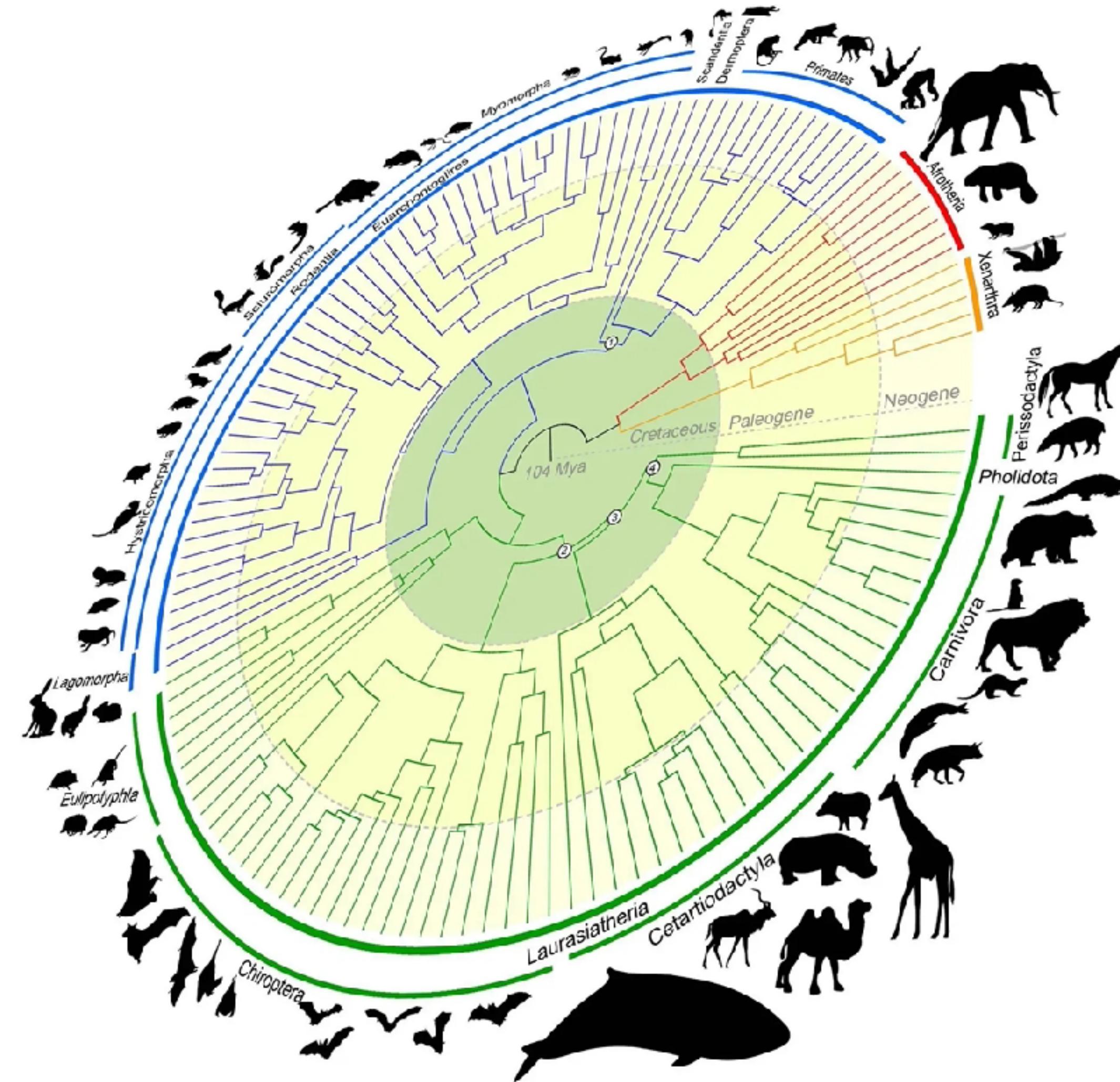
NIH HUMAN
MICROBIOME
PROJECT

Genomics

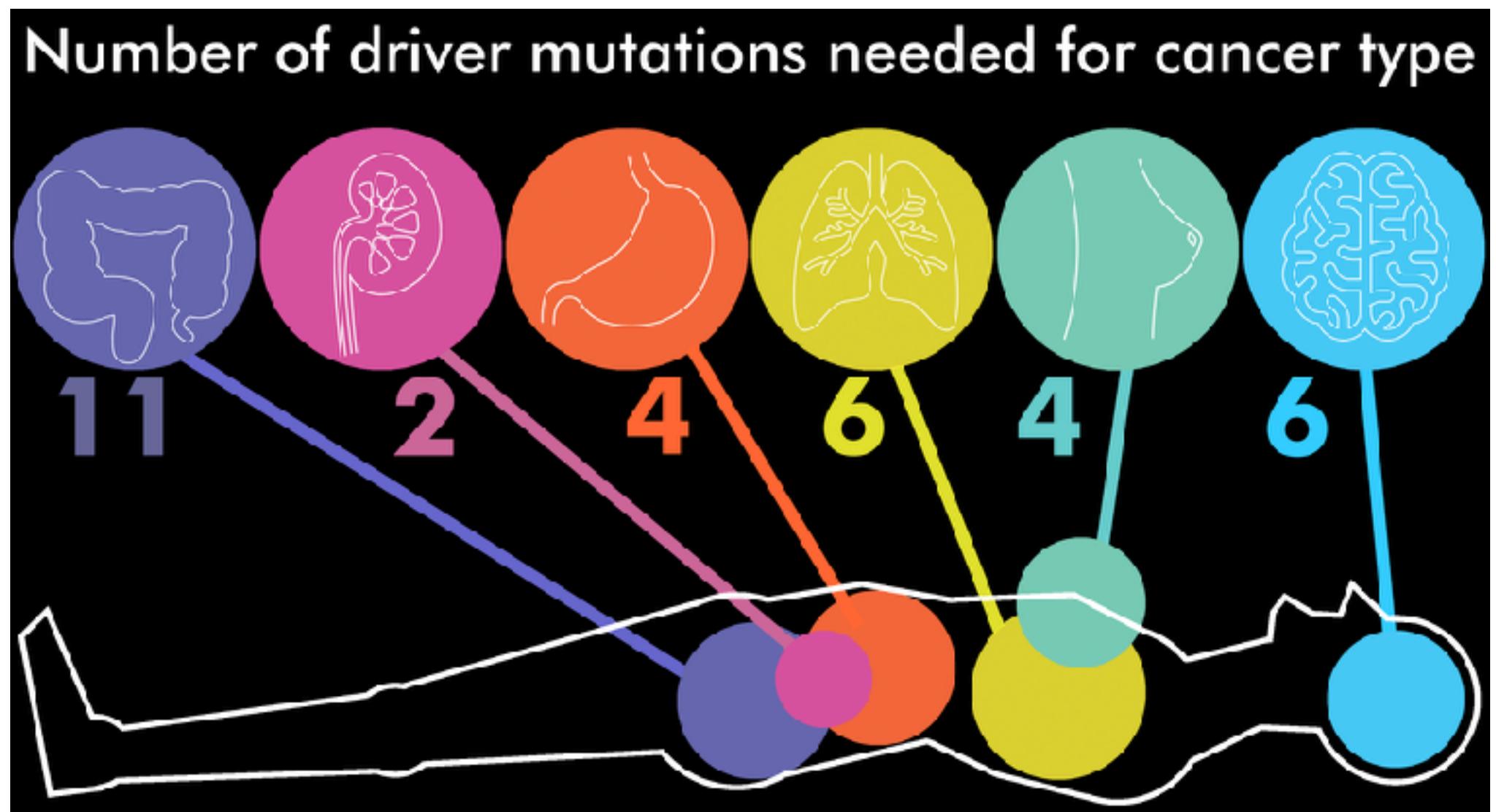
Biological Diversity and Conservation



Phylogenetic Trees

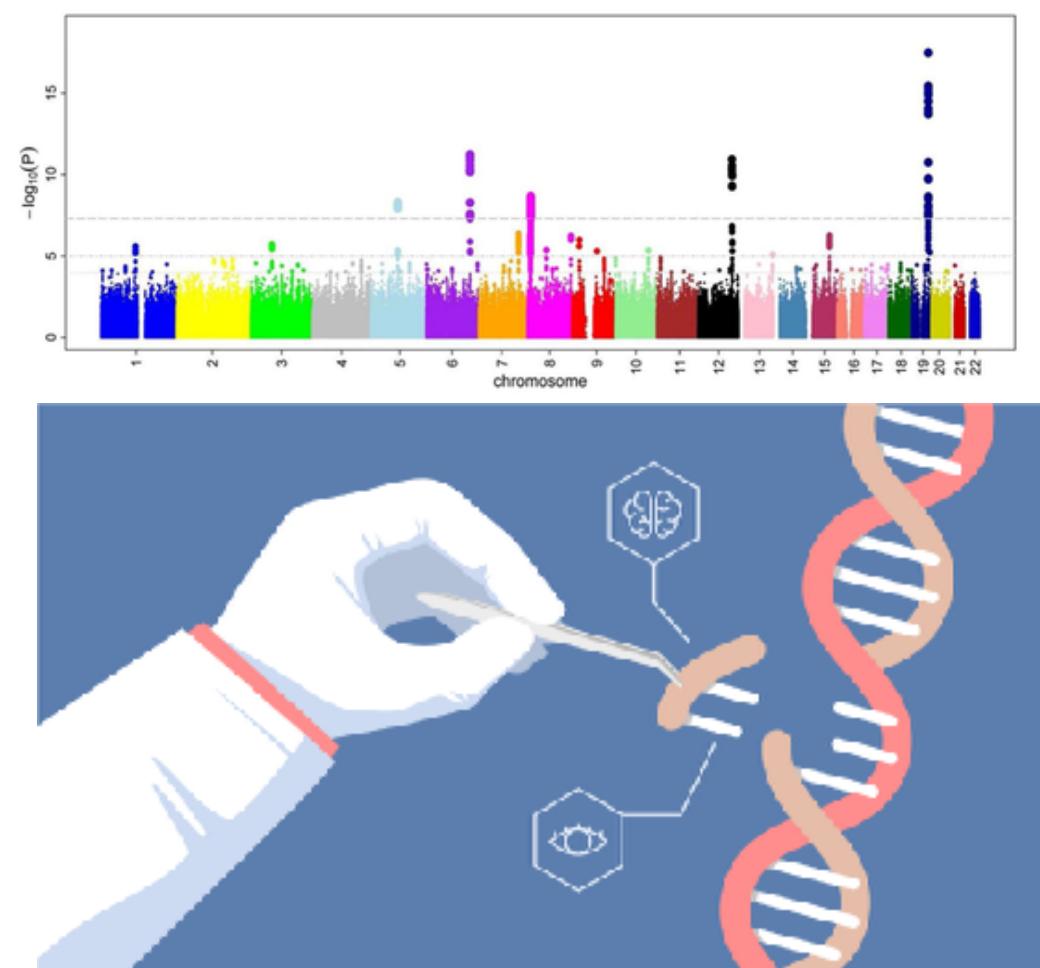


Most mutations are neutral



Which mutations?

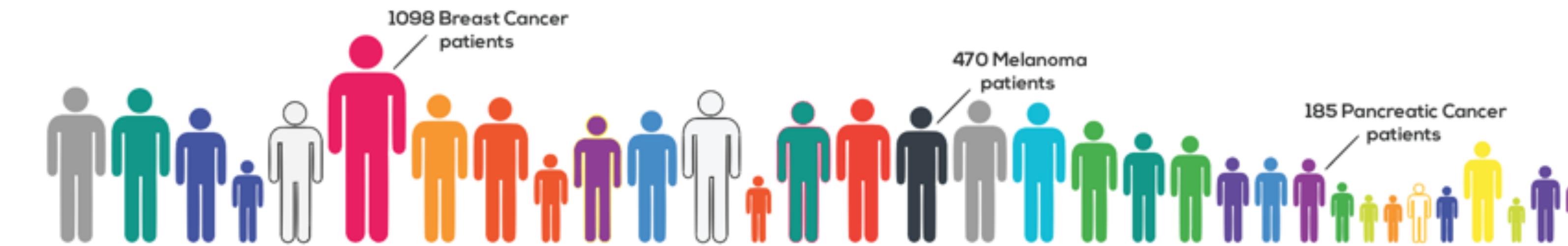
- In silico candidates - statistical prediction
- In vitro mutagenesis - experimentally mutate



Genomics

Precision Medicine

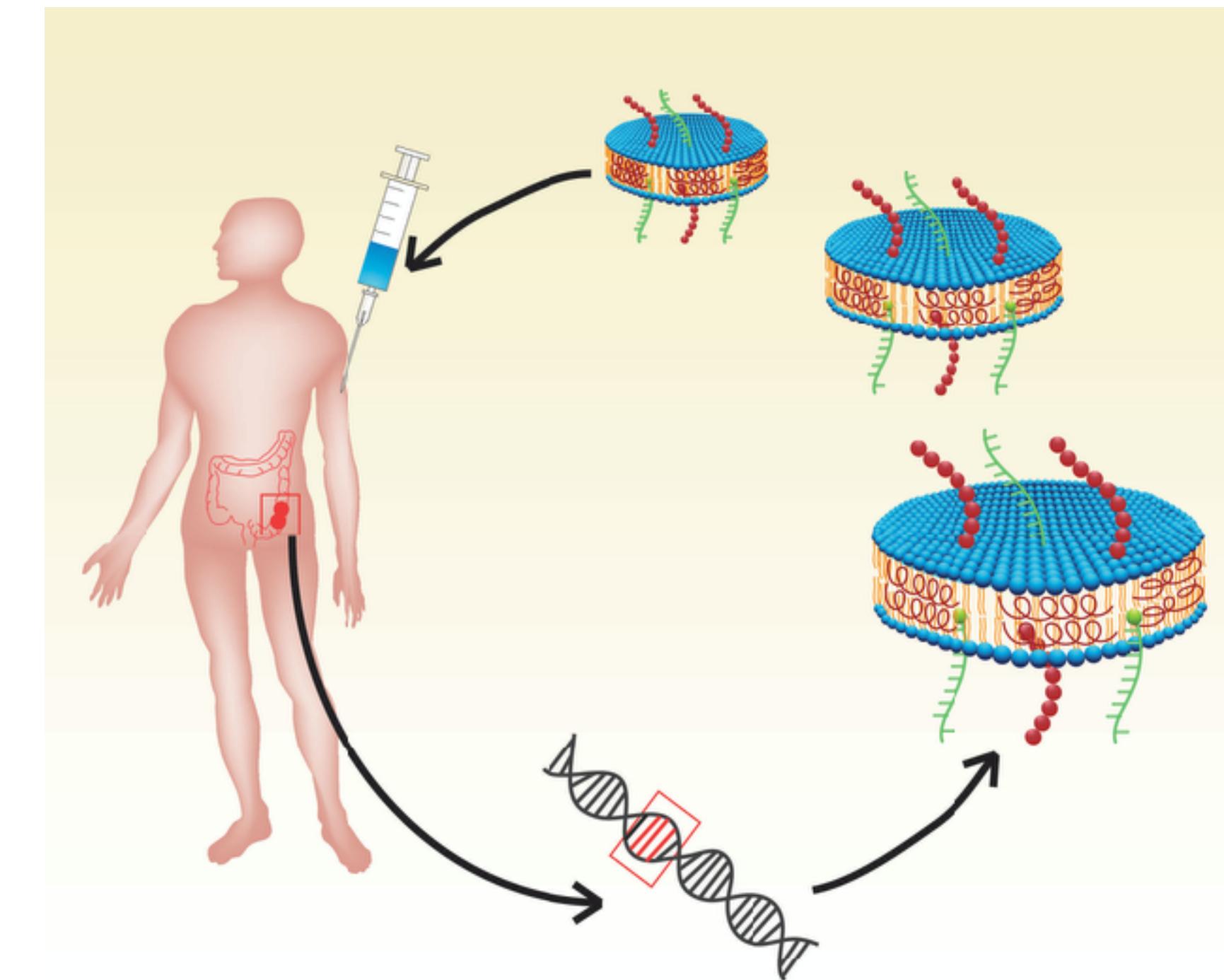
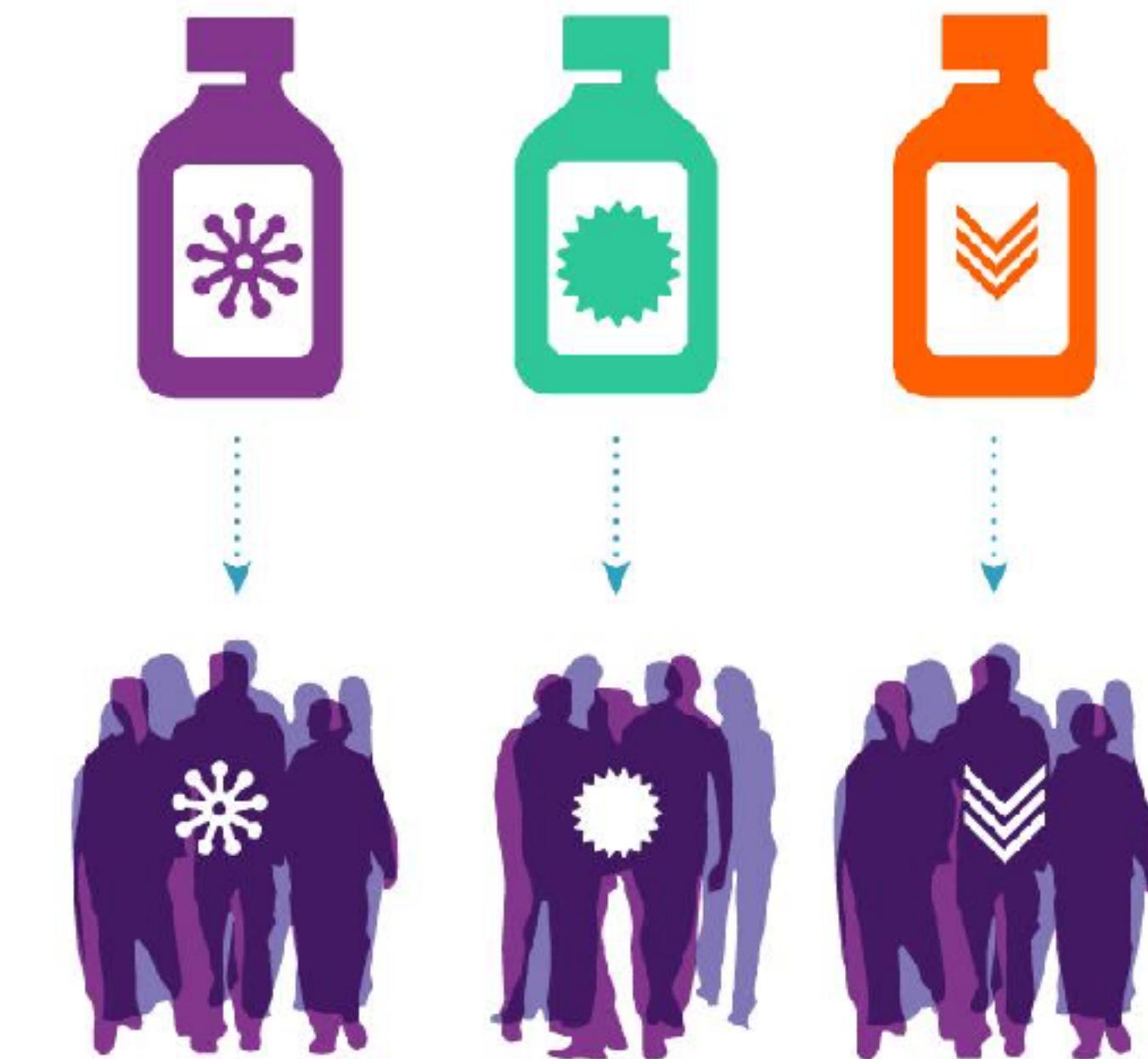
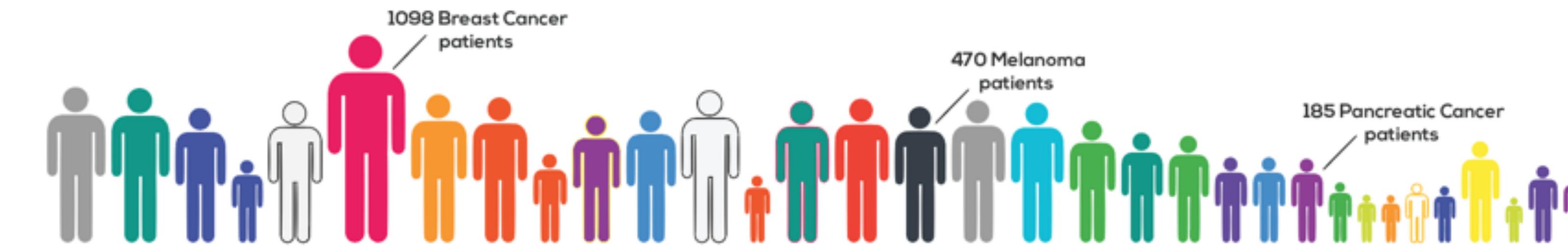
Matched tumor & normal tissues from more than **11,000** patients, representing **33** cancer types.



Genomics

Precision Medicine

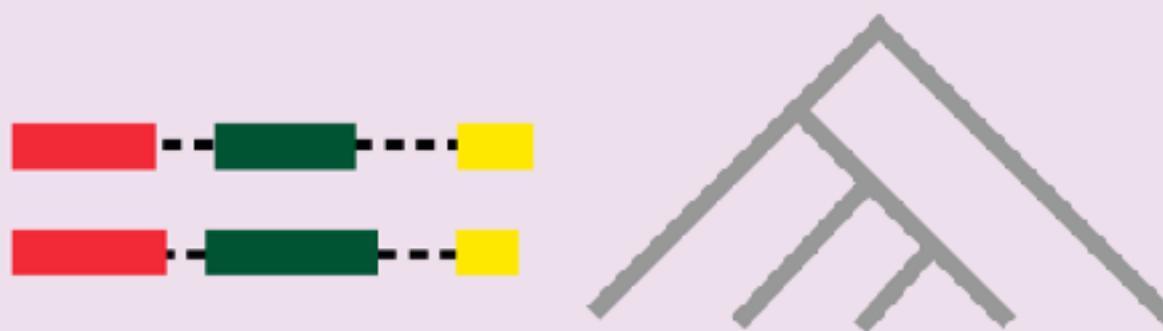
Matched tumor & normal tissues from more than **11,000** patients, representing **33** cancer types.



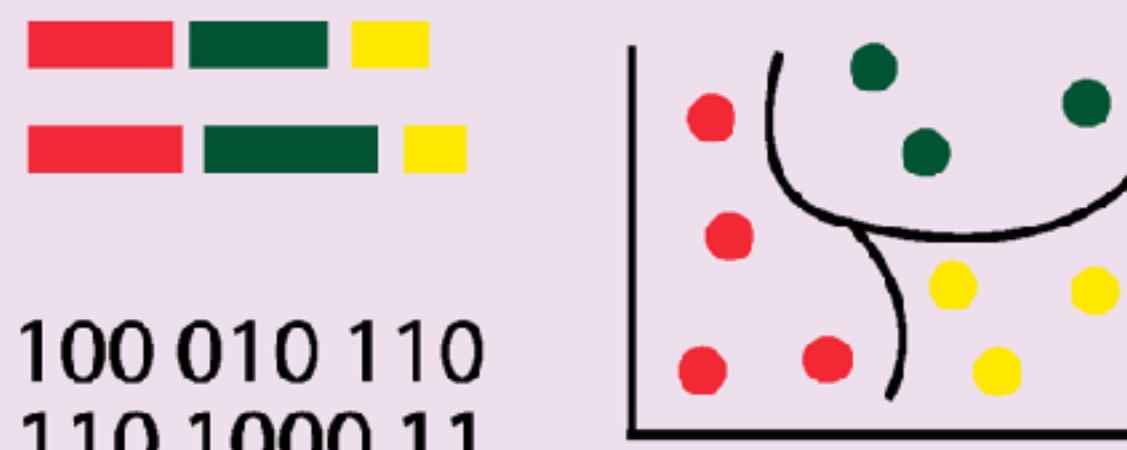
Machine and deep learning integration with bioinformatics

Molecular evolution

Phylogenetic inference

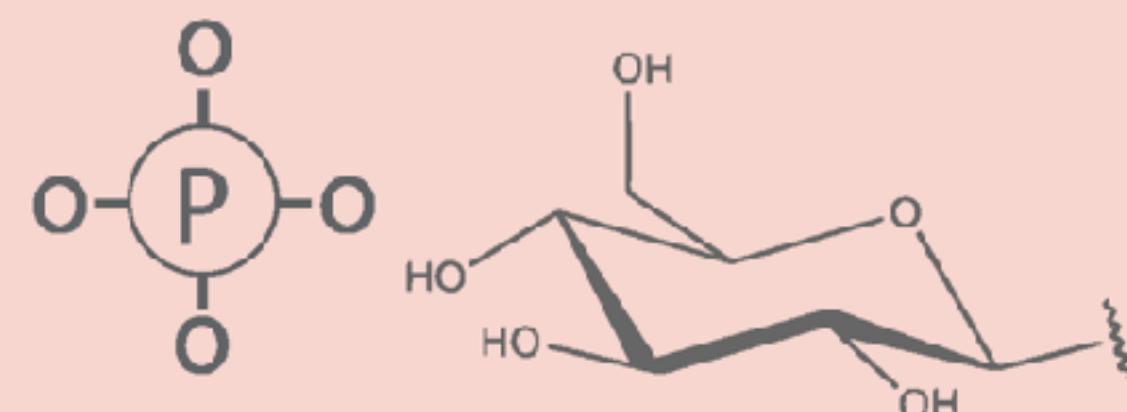


Alignment-free sequence classification

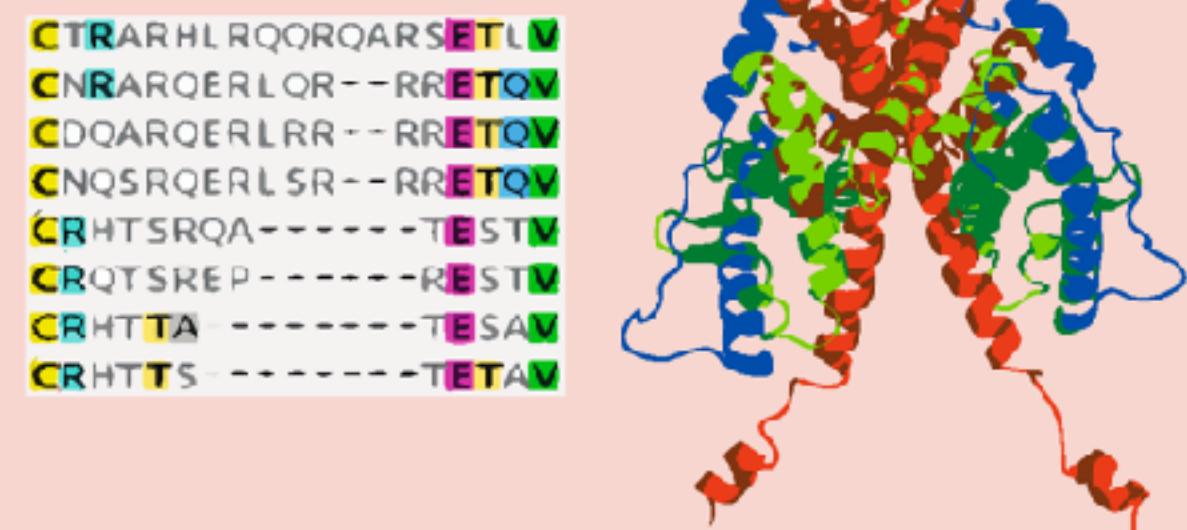


Protein structure Analysis

Post translational modification

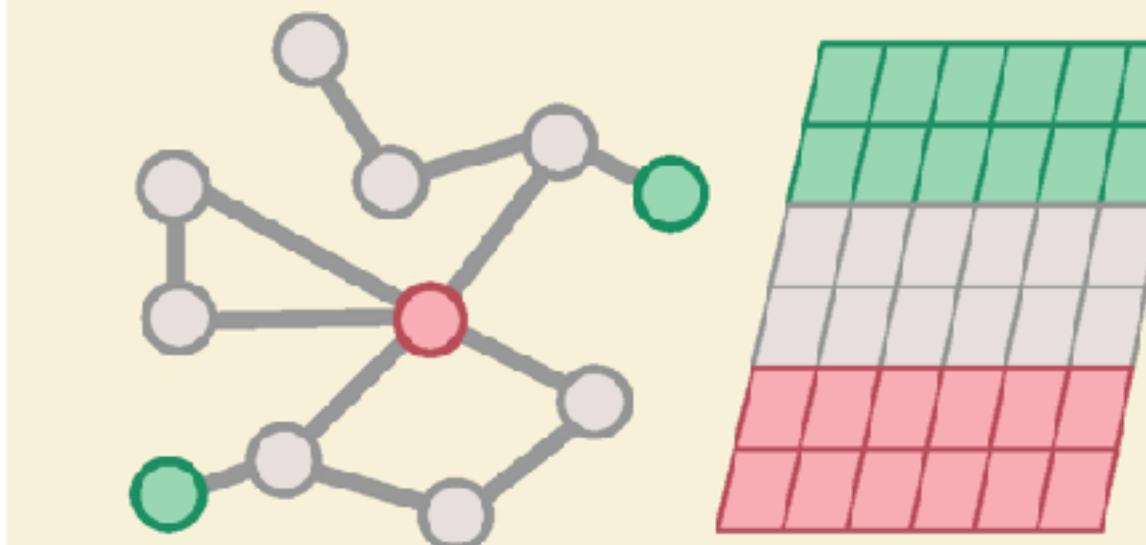


Folding and structure

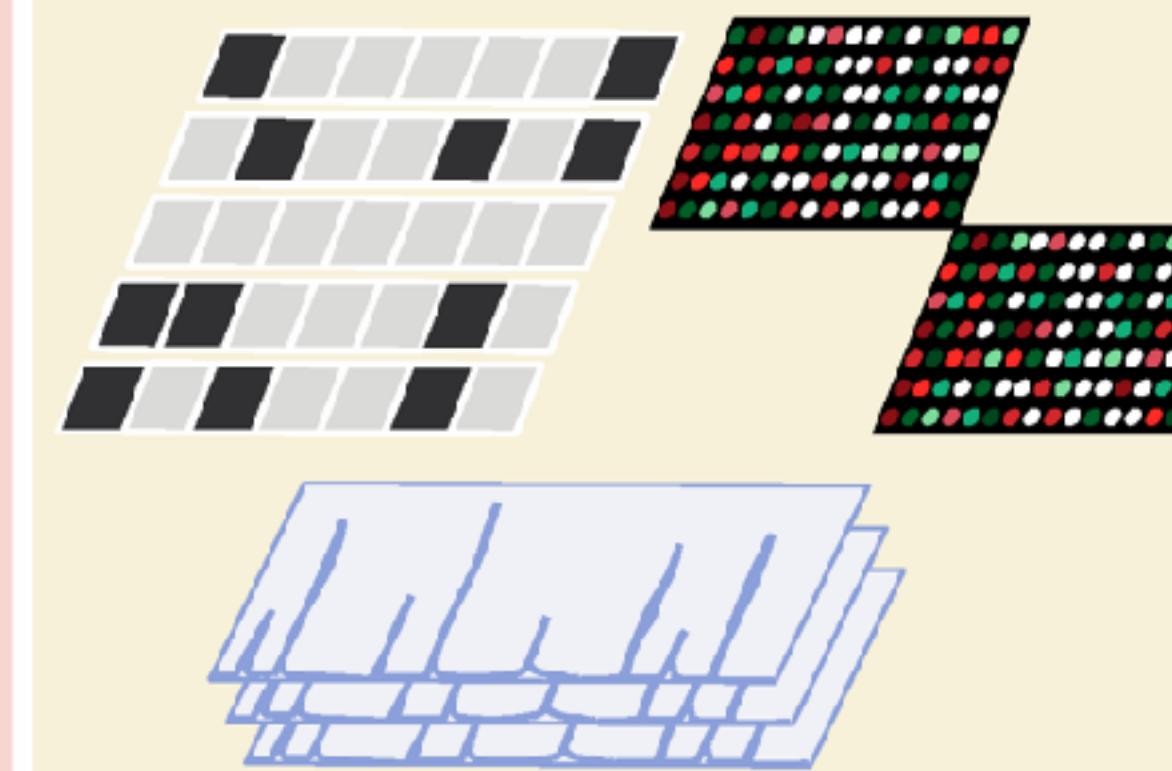


Systems biology

Biological Networks

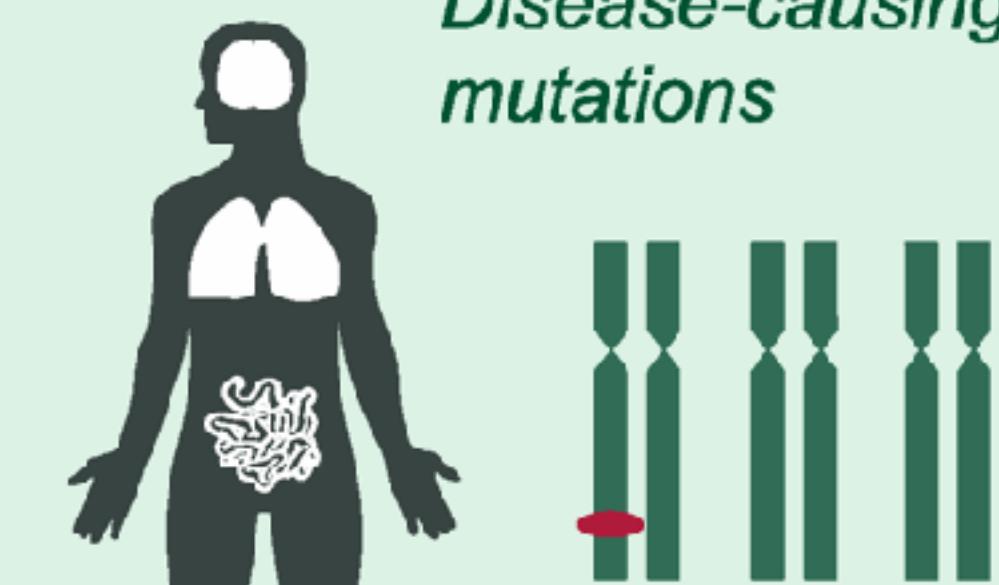


Multi-Omics integration

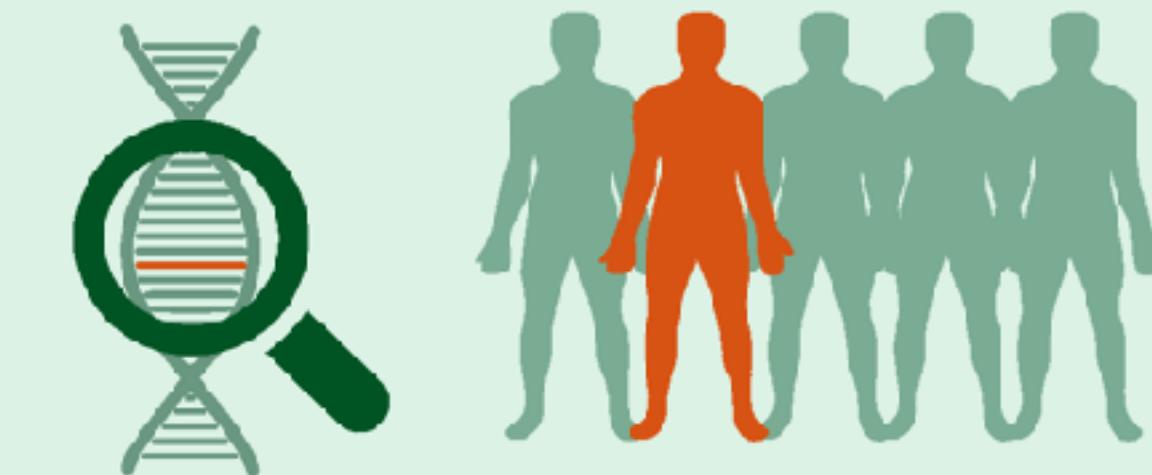


Genomics for Disease Research

Disease-causing mutations



Biomarkers discovery



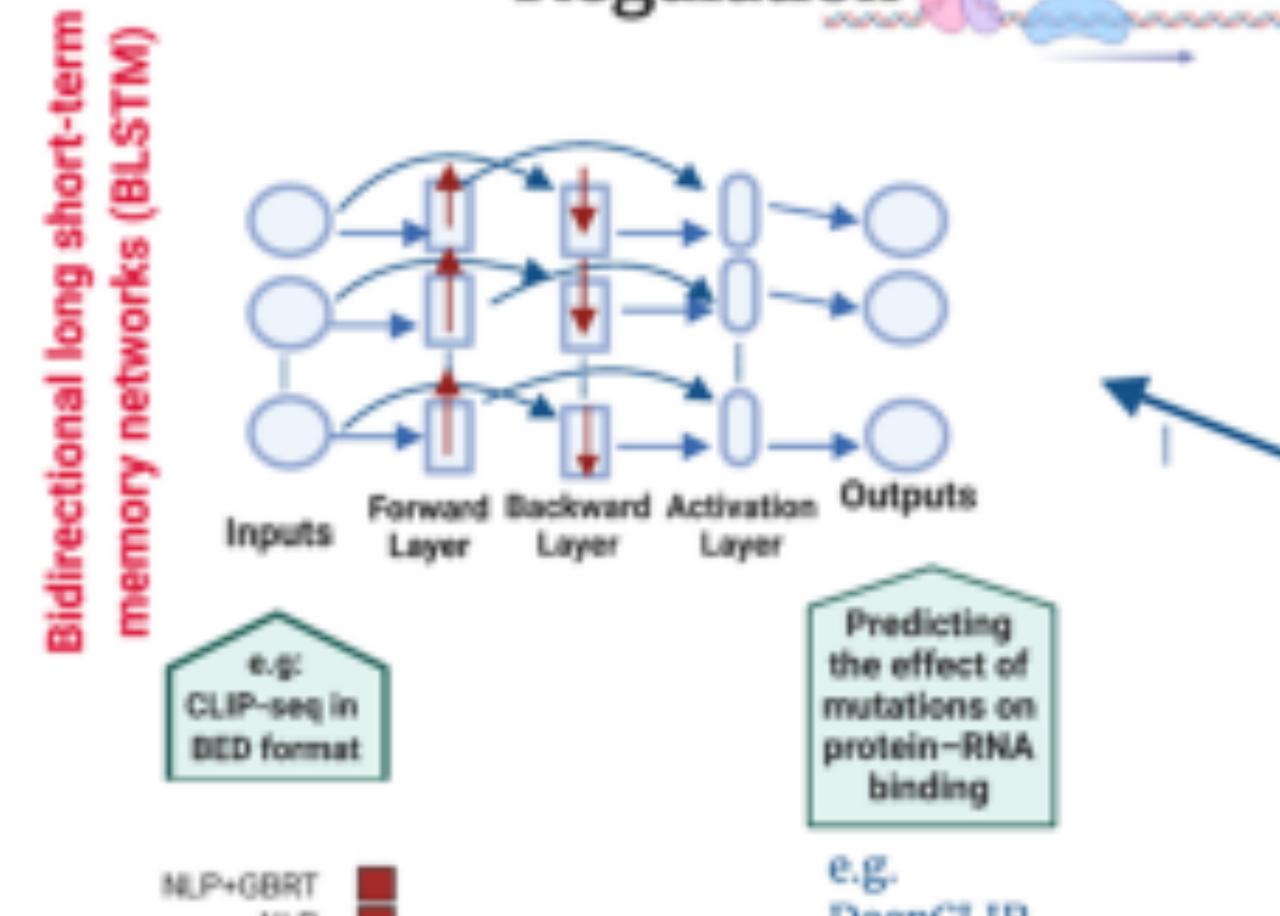
- Inference of tree topology
- Sequence classification
- Viral sequence identification
- Functional annotation

- Phosphorylation site prediction
- Protein glycosylation prediction
- Protein contact maps prediction
- Structural homology prediction

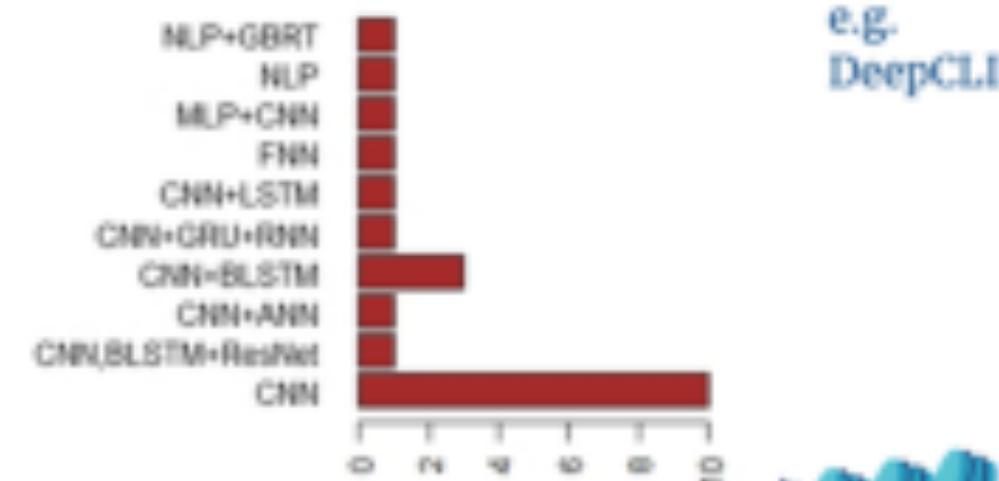
- Biological networks construction
- Biological interactions prediction
- Pathway dynamics prediction
- Platform integration frameworks

- Disease associated genes and mutations
- Biomarkers
- Precision medicine applications

Gene Expression & Regulation



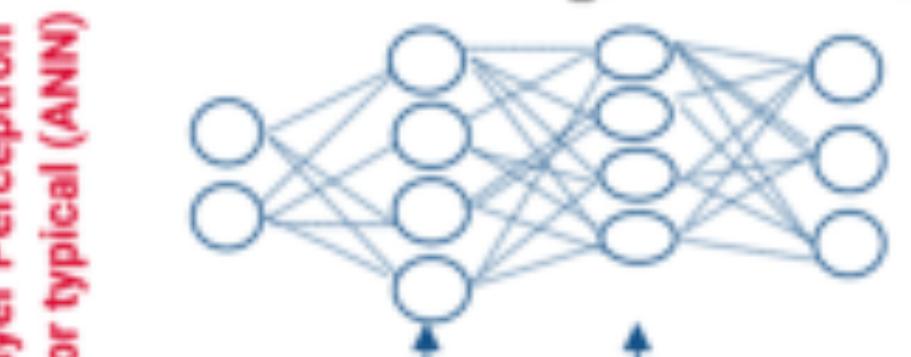
e.g:
CLIP-seq in BED format



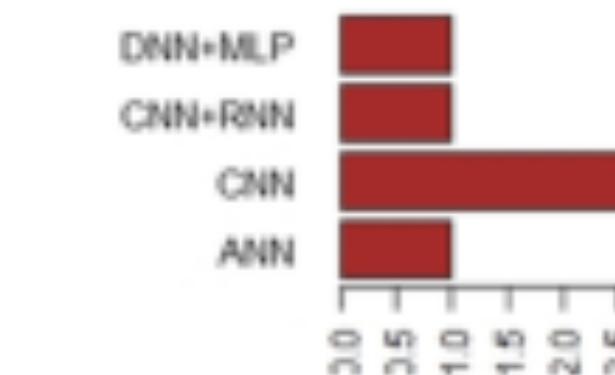
Predicting the effect of mutations on protein-RNA binding
e.g. DeepCLIP

Genomics

Variants calling & Annotation

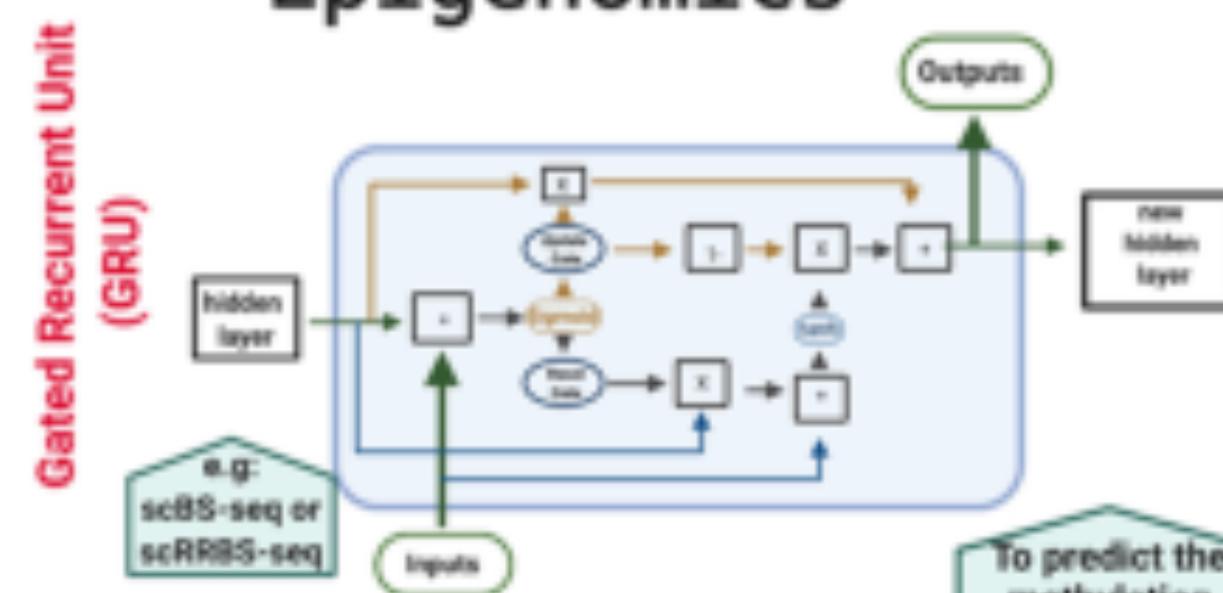


e.g:
WES data in VCF

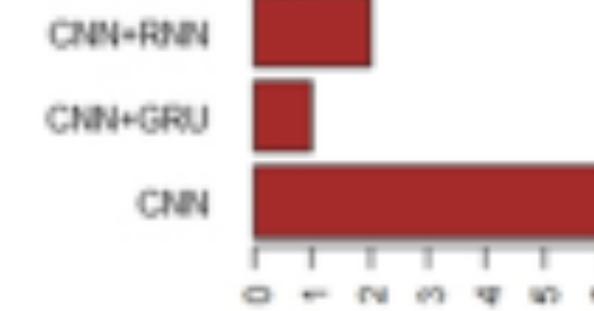


Classify true and false variants from WES data
e.g. GARFIELD - NGS

Epigenomics



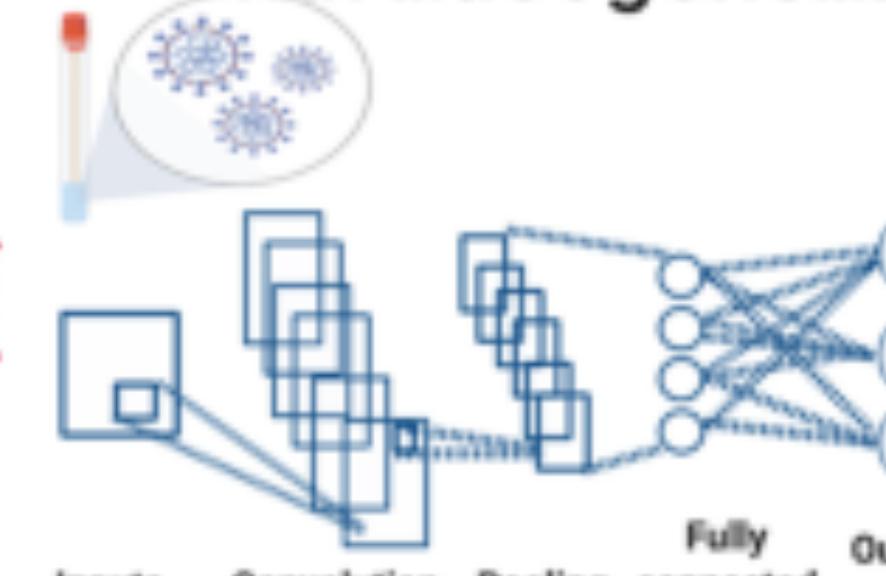
e.g:
scBS-seq or scRRBS-seq



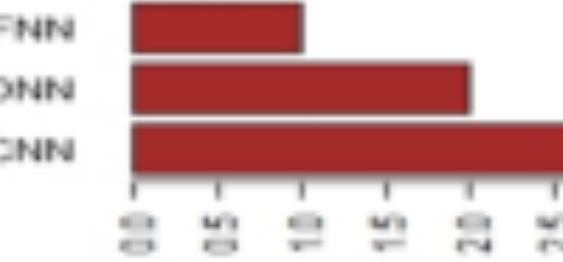
To predict the methylation states from single-cell data
e.g. DeepCpG

Pharmacogenomic

Convolutional Neural Network (CNN)



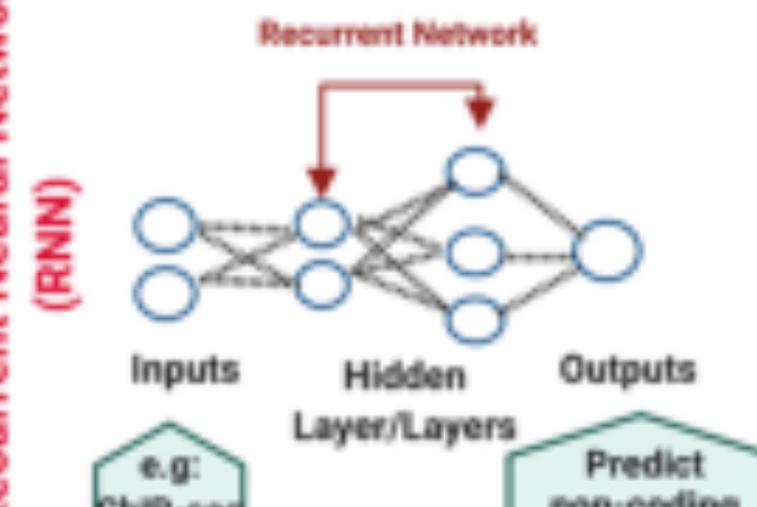
e.g:
RatSeq Data



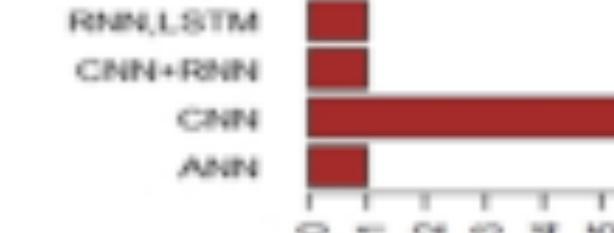
Predict the beta-lactamase (BLs) using protein or genome sequence datasets
e.g. DeepBL

Disease Variants

Recurrent Neural Network (RNN)



e.g:
ChIP-seq



Predict non-coding seq. variants impact
e.g. DeepMLO

Questions?

Introduction to Bioinformatics: An Overview

Tugce Bilgin Sonay, ZHAW, May 2023

Bioinformatics for Beginners

Mutations

Tugce Bilgin Sonay, ZHAW, May 2023

Mutations:

- **Single Nucleotide Polymorphisms – SNPs**

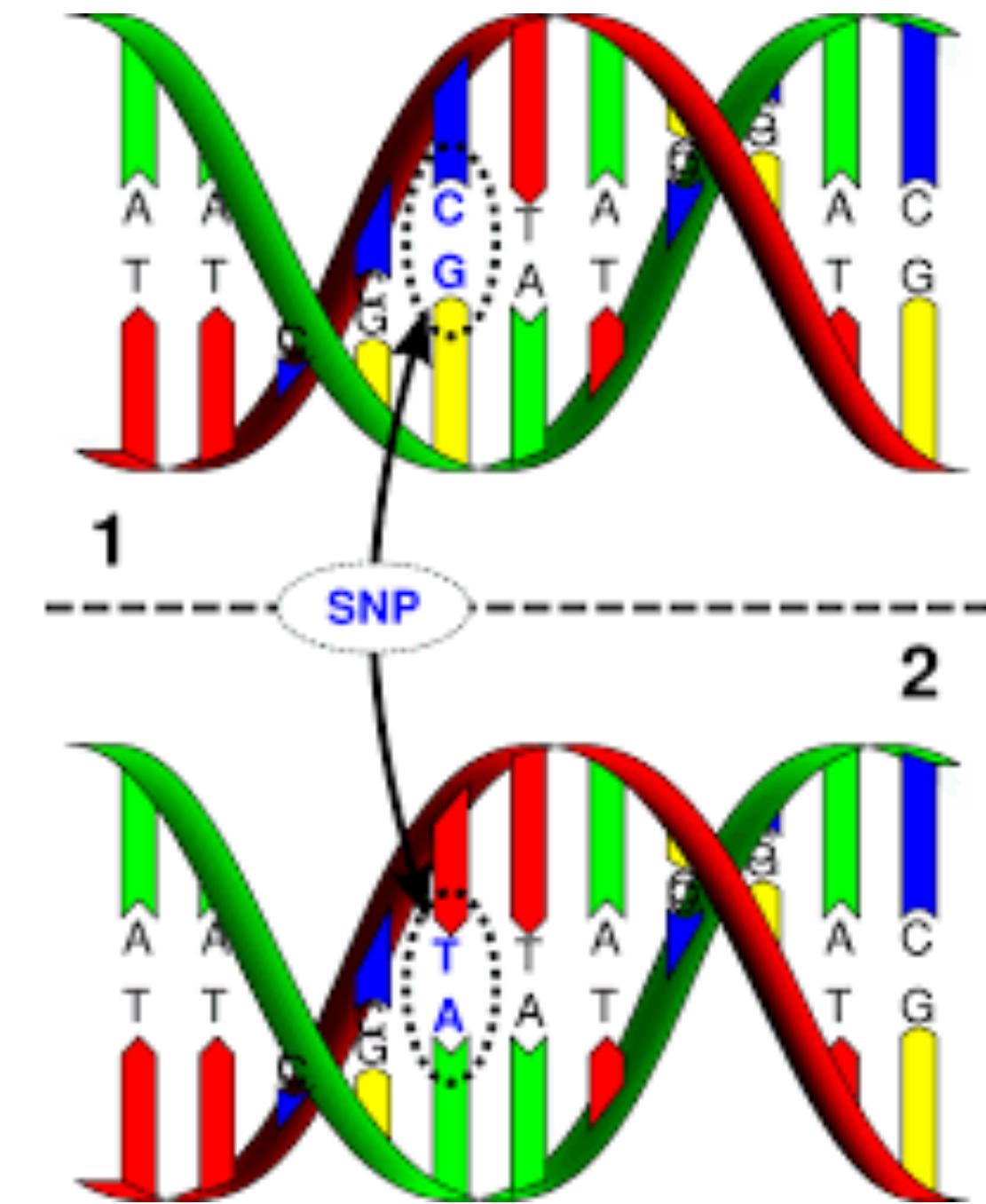
Mutation in 10-100 million bases per generation per cell

- **Indels**

Insertions and deletions of sequences, less frequent

- **Copy Number Mutations**

Long or short sequences



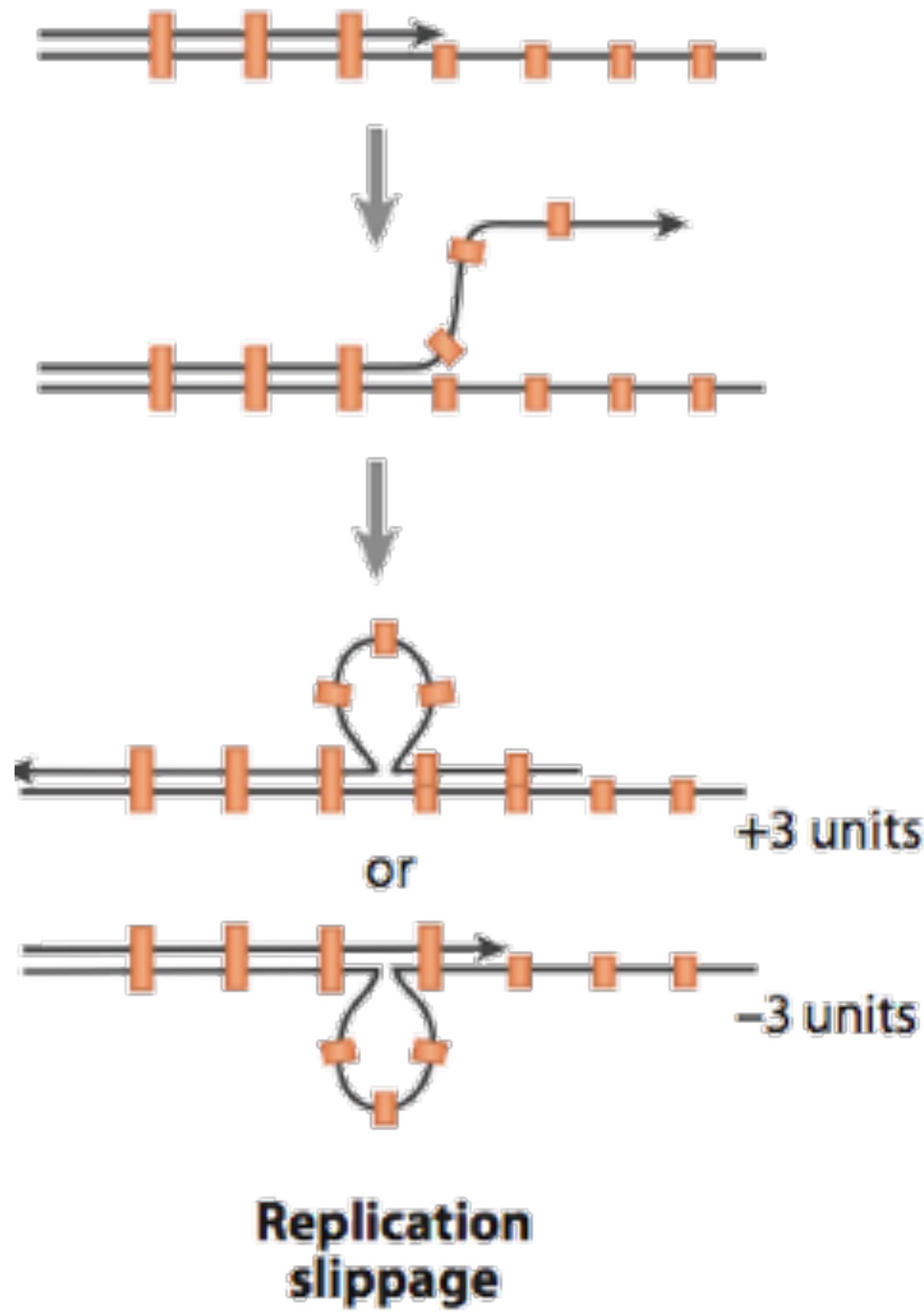
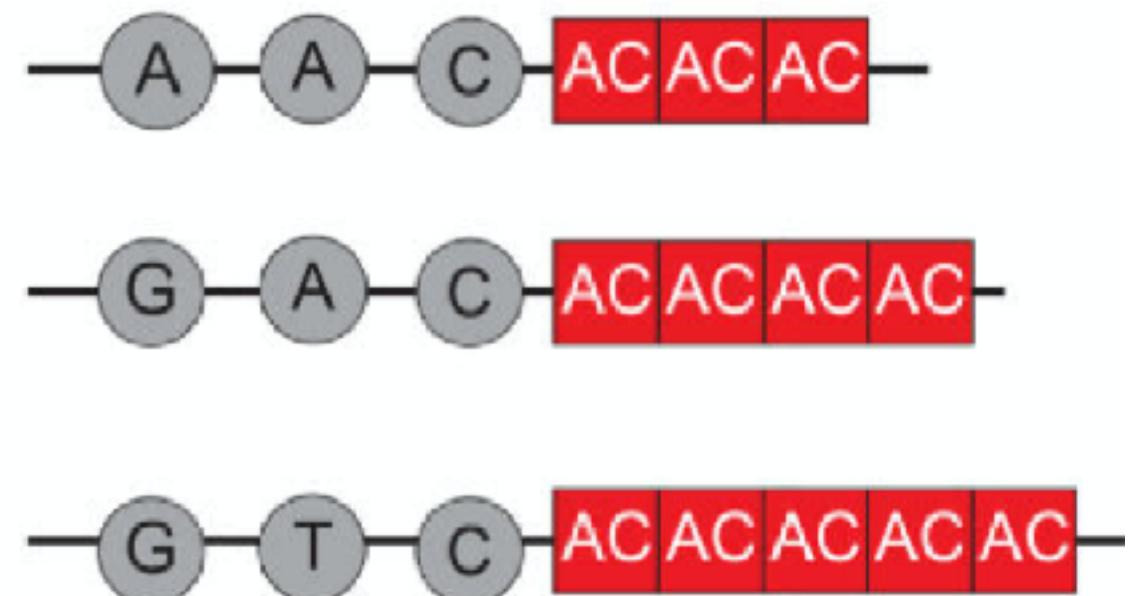
Short Tandem Repeats:

-Tandem repeat instability
20 % of functional variants

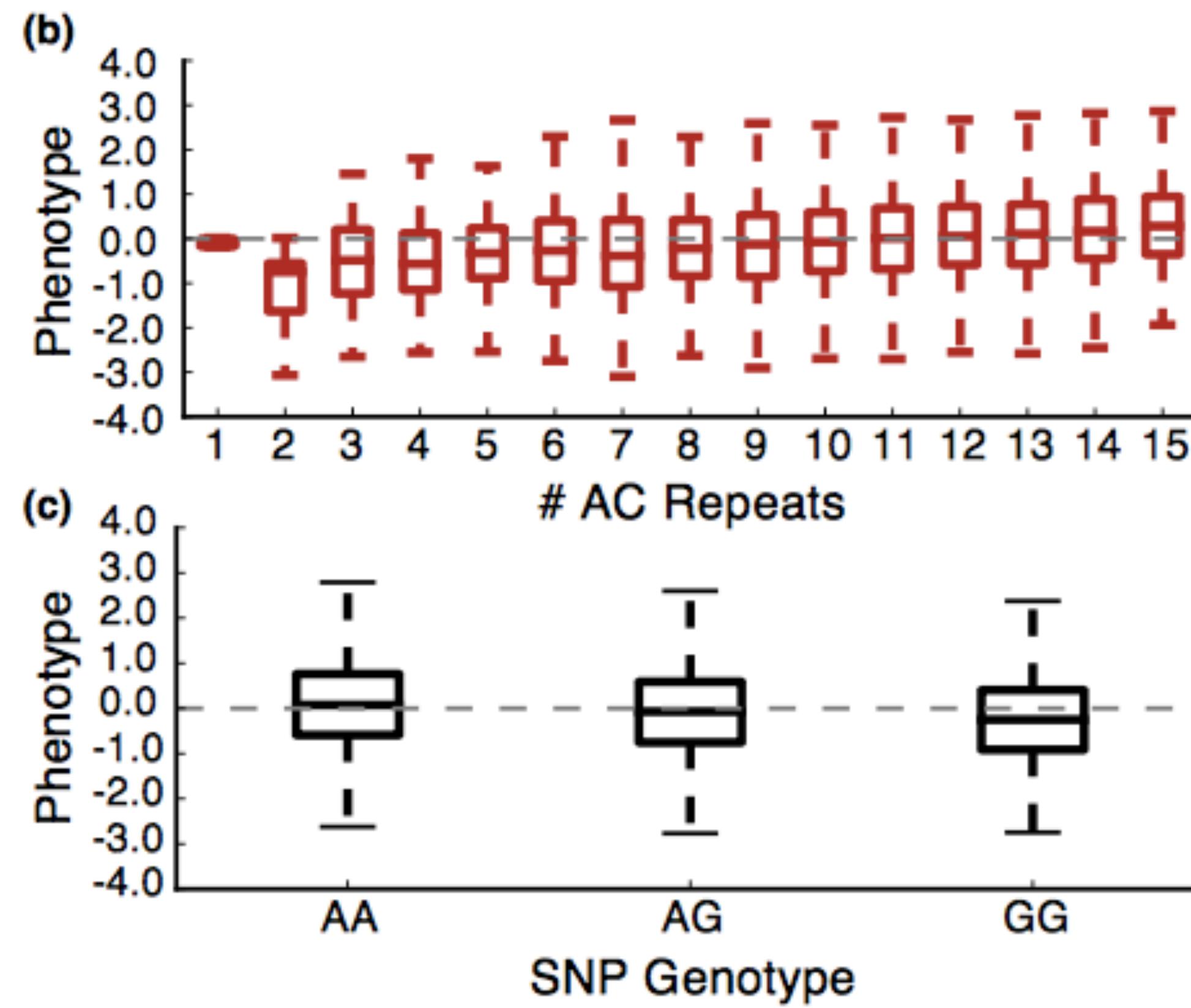
...ATTTAATTAAATTAAATTAA...

extremely high mutation rates: 10^{-2} - 10^{-5} per cell division

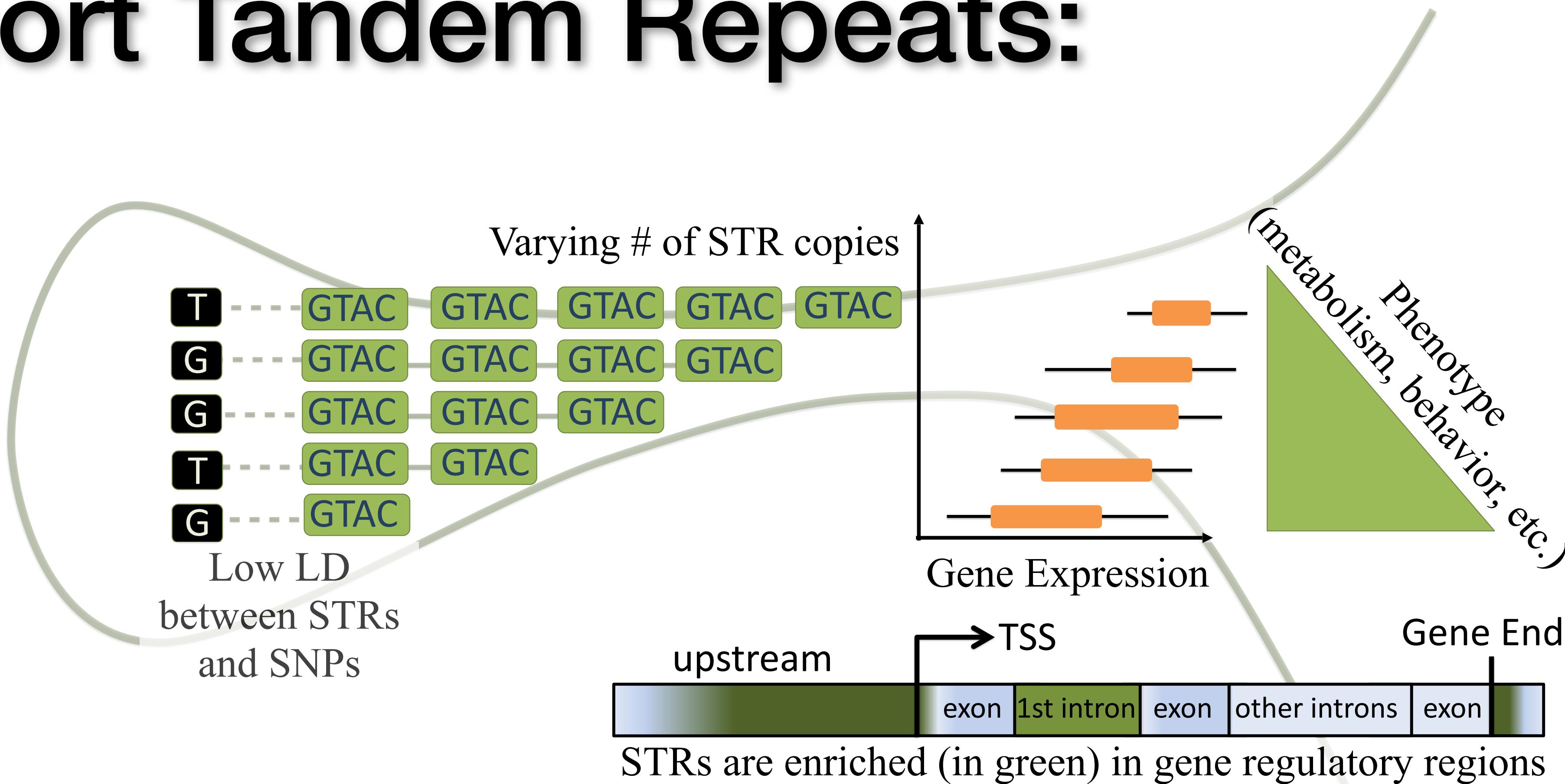
Microsatellites: Short Tandem Repeats in DNA



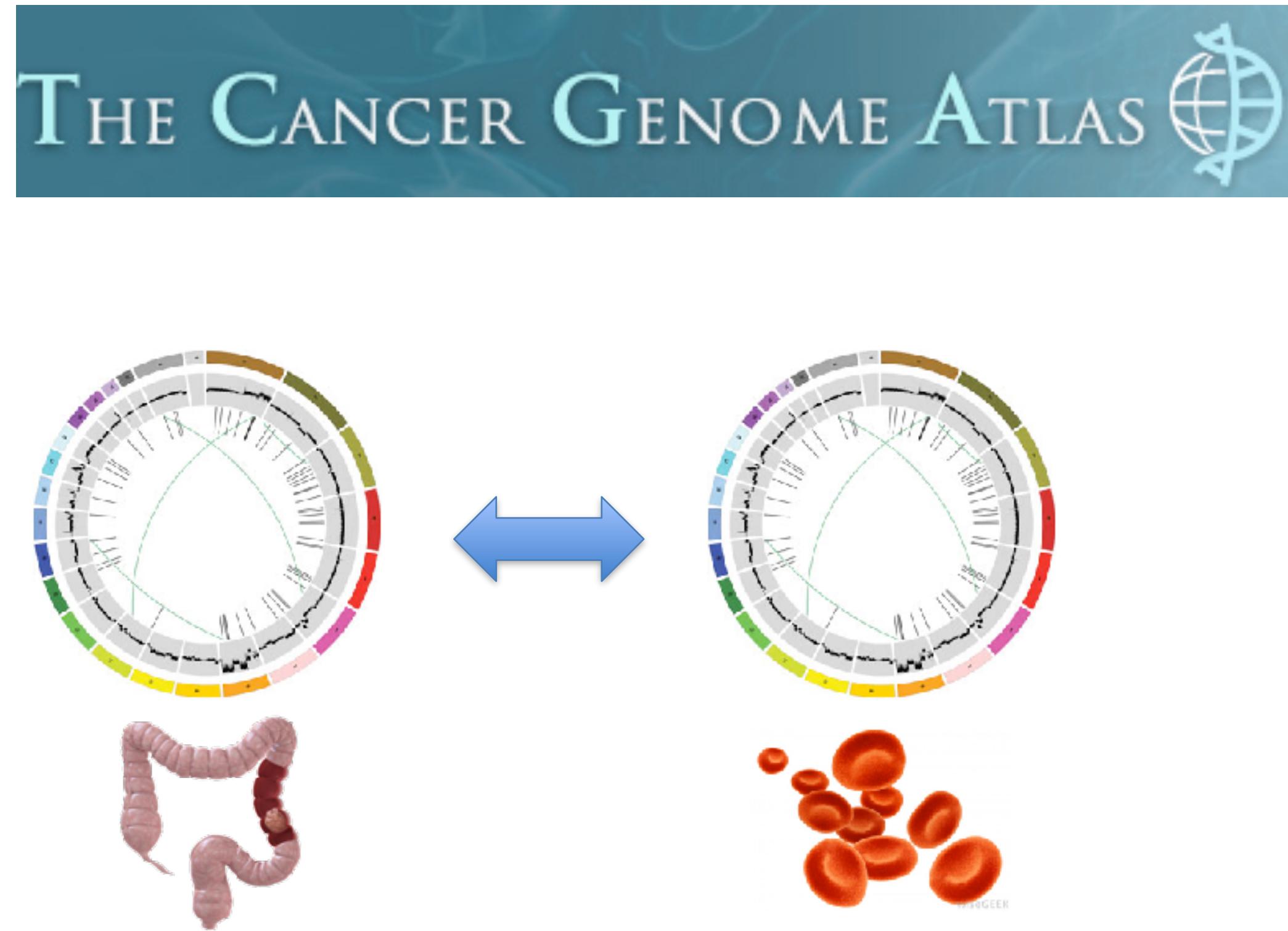
Short Tandem Repeats offer much more genotypic variation than SNPs



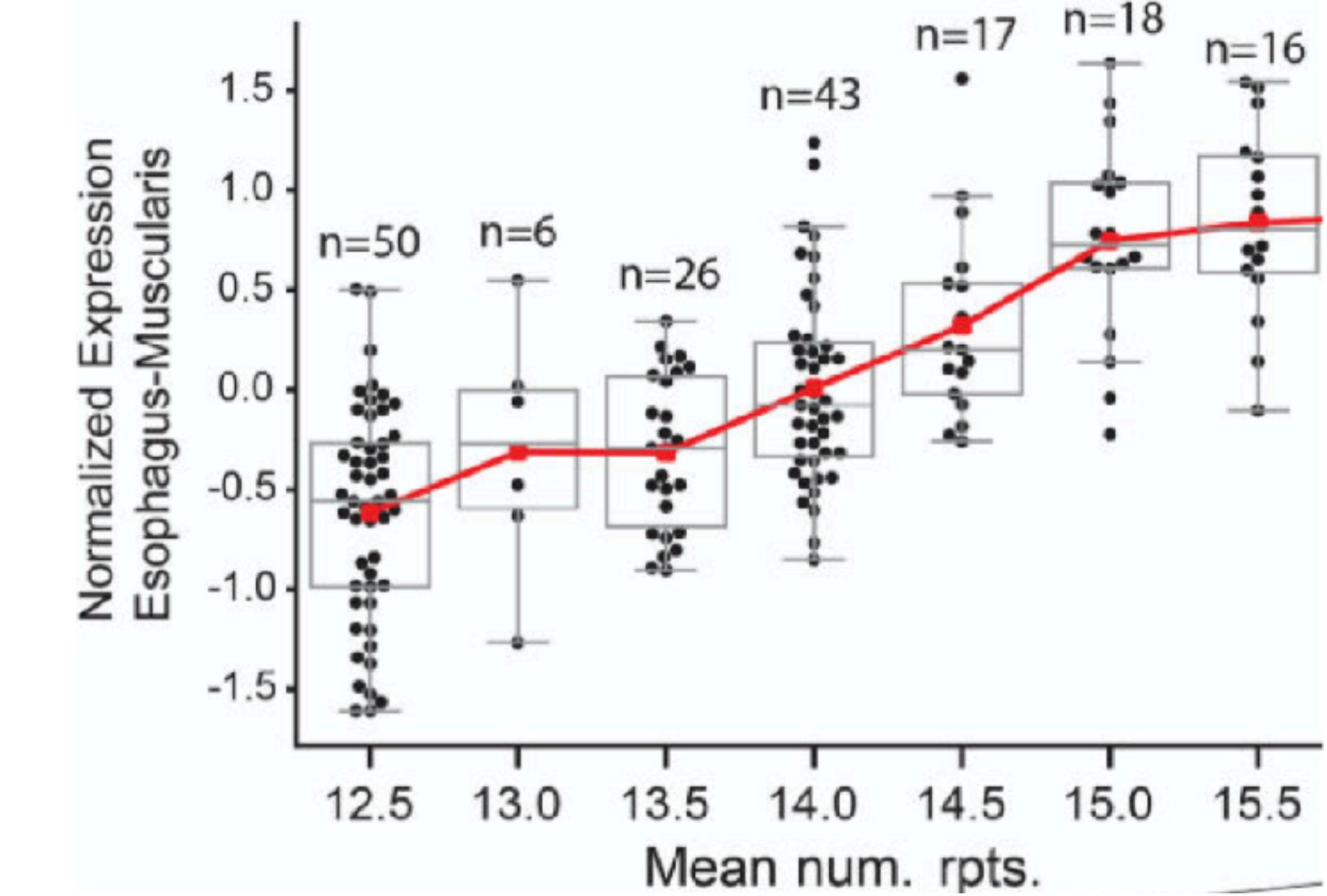
Short Tandem Repeats:



Short Tandem Repeats and Cancer:



Bilgin Sonay et al., BMC Genomics, 2015



Fosting et al., Nat Genet, 2020

Short Tandem Repeats and Cancer:



Developing the First Precision Immunotherapy

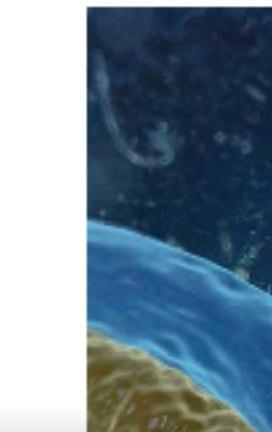
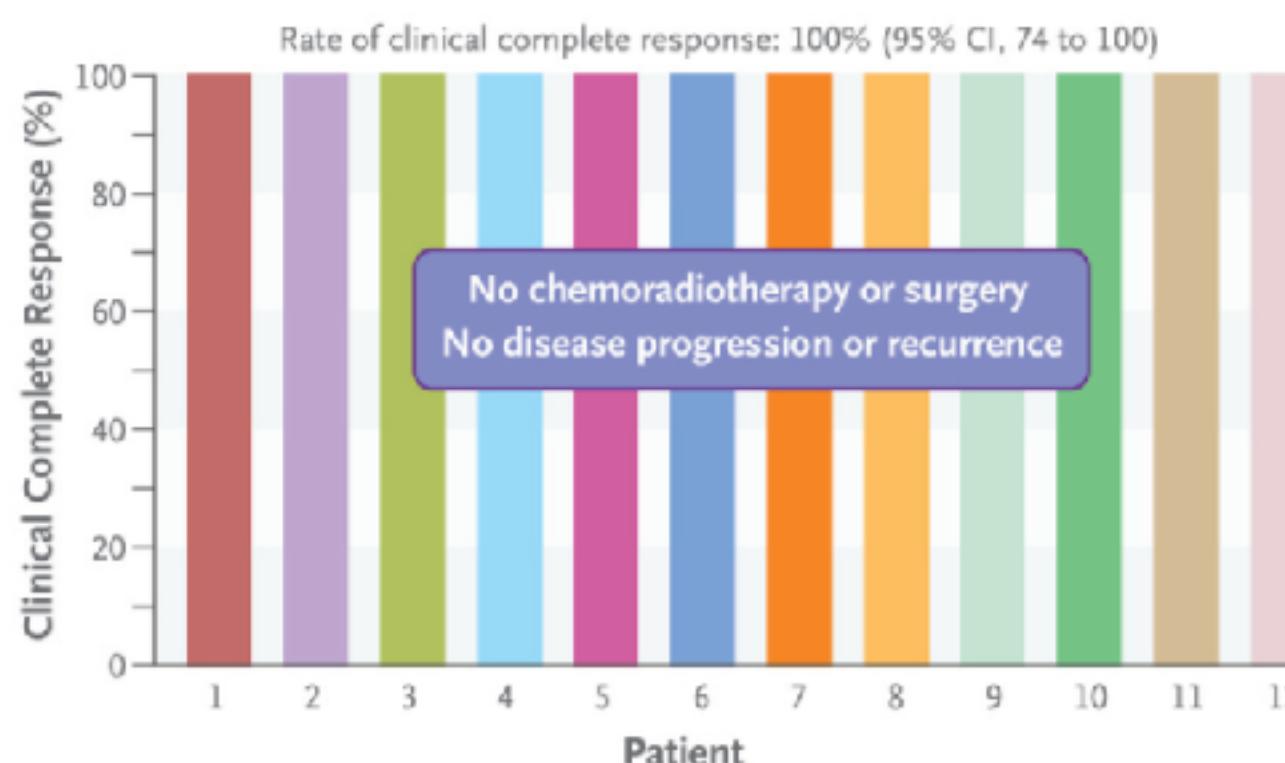
In May 2017, the Food and Drug Administration (FDA) approved the first drug to treat tumors based on their genetic characteristics, regardless of where in the body the cancer originated. Until now, drugs have been approved based on their cell or tissue of origin, such as the breast or the lung. But pembrolizumab (Keytruda®) doesn't target the genetic abnormalities of cancer cells specifically—it targets the immune system.



PD-1 Blockade in Mismatch Repair–Deficient, Locally Advanced Rectal Cancer

Cerck A et al. DOI: 10.1056/NEJMoa2201445

Overall Response to Dostarlimab in 12 Patients



Article | OPEN | Published: 06 June 2017

A molecular portrait of microsatellite instability across multiple cancers



Article | Published: 11 September 2017

Analysis of somatic microsatellite indels identifies driver events in human tumors



The Landscape of Microsatellite Instability in Colorectal and Endometrial Cancer Genomes



Deep learning can predict microsatellite instability directly from histology in gastrointestinal cancer

Cite as: D. T. Le et al., *Science* 10.1126/science.aan6733 (2017).

Mismatch-repair deficiency predicts response of solid tumors to PD-1 blockade

Short Tandem Repeats and Cancer:

Zürcher Hochschule
für Angewandte Wissenschaften



Maria
Anisimova

**Intrinsic Disorder as a potential link
between tandem repeats and cancer**

Verbiest, Delucchi, Bilgin Sonay
and Anisimova, 2021

**Mutation and selection processes
regulating short tandem repeats**

Verbiest, ..., Anisimova, Gymrek,
Bilgin Sonay, invited review in
JEB, 2022



Questions?

Mutations

Tugce Bilgin Sonay, ZHAW, May 2023

Bioinformatics for Beginners

Bioinformatics Resources

Tugce Bilgin Sonay, ZHAW, May 2023

Protein Databases

- Protein-protein interaction: <https://string-db.org/>
- Protein databases: <https://www.uniprot.org/> <https://www.rcsb.org/>

Bioinformatics Forums

Notebooks

- Jupyter Notebook, bioinformatician's diary: <https://jupyter.org>

Downstream Genomic

Omics Databases

- Bedtools : <https://bedtools.readthedocs.io/en/latest/>
- USCS Table Browser: <https://genome.ucsc.edu/cgi-bin/hgTables>
- Biomart Ensembl: <https://useast.ensembl.org/index.html>

Cancer Databases

Gene Databases

- ncbi entrez: <https://www.ncbi.nlm.nih.gov/>
- Pan Cancer, TCGA: <https://cancergenome.nih.gov>
- Gene Cards: <https://www.genecards.org>
- DAVID for Gene Ontology terms: <https://david.ncifcrf.gov/summary.jsp>

Genome Databases

- 1000 genomes new name : <https://www.internationalgenome.org>

Coding

- R Circos - <https://cran.r-project.org/web/packages/BioCircos/vignettes/BioCircos.html>
- R packages: <https://www.bioconductor.org>
- python packages: Biopython



CONDA[®]



MINI CONDA[®]

= conda
+ python
+ base packages



ANACONDA[®]

= miniconda
+ 150 high quality packages

Questions?

Bioinformatics Resources

Tugce Bilgin Sonay, ZHAW, May 2023