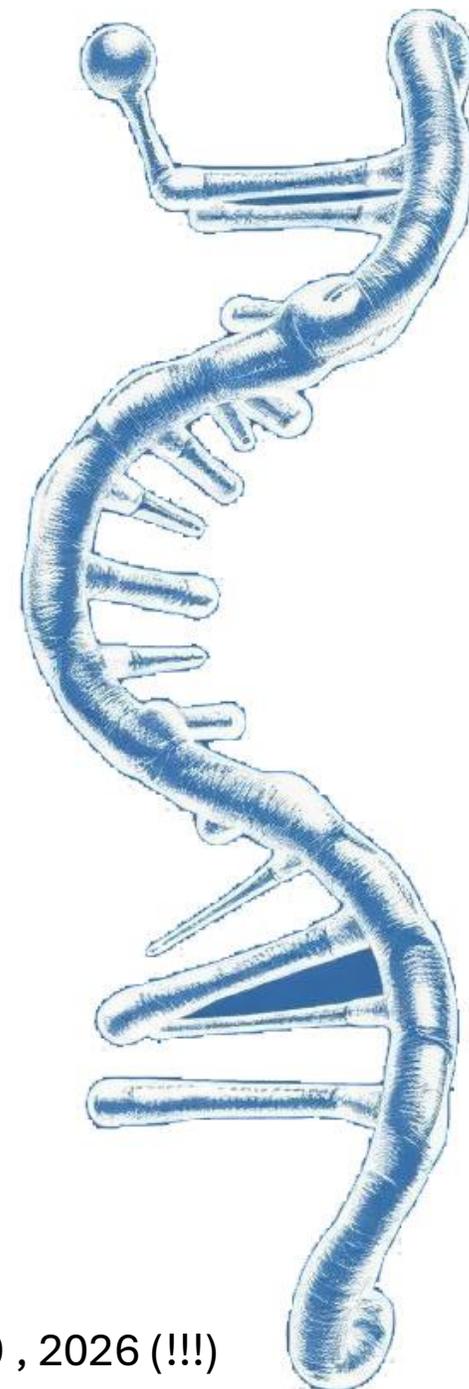


Watson and Crick, Nature 1953

AI and Genetic Part 1

Jean Martin Beaulieu
January 29, 2026



GPT 5.0 , 2026 (!!!)

Advancing regulatory variant effect prediction with AlphaGenome

<https://deepmind.google.com/science/alphagenome/>

<https://doi.org/10.1038/s41586-025-10014-0>

Received: 16 May 2025

Accepted: 4 December 2025

Published online: 28 January 2026

Open access

 Check for updates

Žiga Avsec^{1,2}, Natasha Latysheva^{1,2}, Jun Cheng^{1,2}, Guido Novati^{1,2}, Kyle R. Taylor^{1,2}, Tom Ward^{1,2}, Clare Bycroft^{1,2}, Lauren Nicolaisen^{1,2}, Eirini Arvaniti^{1,2}, Joshua Pan^{1,2}, Raina Thomas¹, Vincent Dutordoir¹, Matteo Perino¹, Soham De¹, Alexander Karollus¹, Adam Gayoso¹, Toby Sargeant¹, Anne Mottram¹, Lai Hong Wong¹, Pavol Drotár¹, Adam Kosiorek¹, Andrew Senior¹, Richard Tanburn¹, Taylor Applebaum¹, Souradeep Basu¹, Demis Hassabis¹ & Pushmeet Kohli¹

Deep learning models that predict functional genomic measurements from DNA sequences are powerful tools for deciphering the genetic regulatory code. Existing methods involve a trade-off between input sequence length and prediction resolution, thereby limiting their modality scope and performance^{1–5}. We present AlphaGenome, a unified DNA sequence model, which takes as input 1 Mb of DNA sequence and predicts thousands of functional genomic tracks up to single-base-pair resolution across diverse modalities. The modalities include gene expression, transcription initiation, chromatin accessibility, histone modifications, transcription factor binding, chromatin contact maps, splice site usage and splice junction coordinates and strength. Trained on human and mouse genomes, AlphaGenome matches or exceeds the strongest available external models in 25 of 26 evaluations of variant effect prediction. The ability of AlphaGenome to simultaneously score variant effects across all modalities accurately recapitulates the mechanisms of clinically relevant variants near the *TALI* oncogene⁶. To facilitate broader use, we provide tools for making genome track and variant effect predictions from sequence.

Data availability

All primary experimental datasets used for the training and evaluation of AlphaGenome in this study were obtained from publicly accessible sources. A comprehensive manifest detailing these data sources, including specific repositories (such as ENCODE portal, GTEx portal, 4D Nucleome portal, ClinVar and gnomAD), individual accession numbers, relevant version information and direct URLs where applicable, is provided in Supplementary Table 2. This study did not generate new primary experimental data requiring deposition.

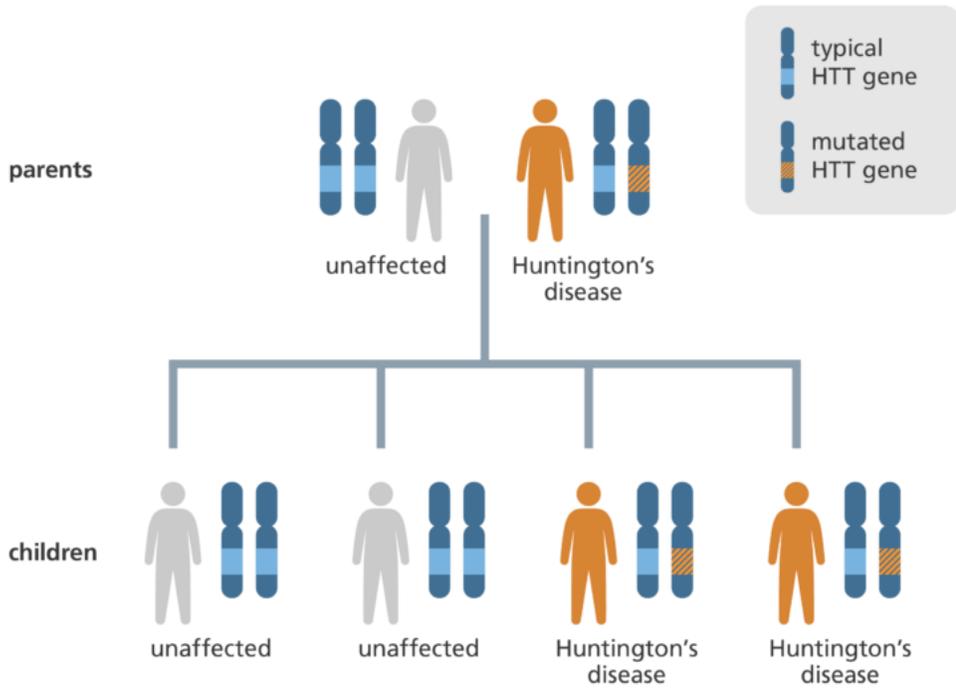
Code availability

AlphaGenome is available for non-commercial use through an online API at <http://deepmind.google.com/science/alphagenome>, with an accompanying Python software development kit provided to interact with the model. We also provide a genome interpretation suite to facilitate the exploration and interpretation of AlphaGenome. This offers a range of functionalities, such as streamlined variant scoring with quantile calibration and identification of critical sequence regions through contribution scores from ISM-based experiments. The model source code, weights, variant scoring implementations and a selection of variant evaluation datasets and predictions are available at https://github.com/google-deepmind/alphagenome_research.

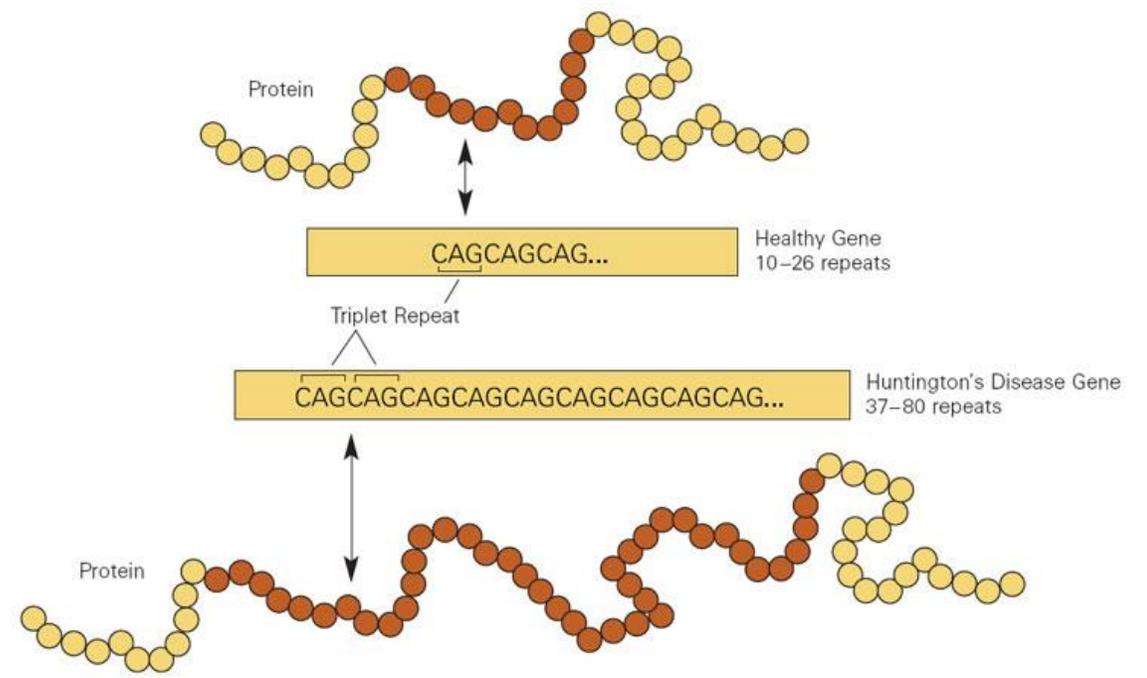


Part 1: Genomes and Genes

Genetic variants and disease causation



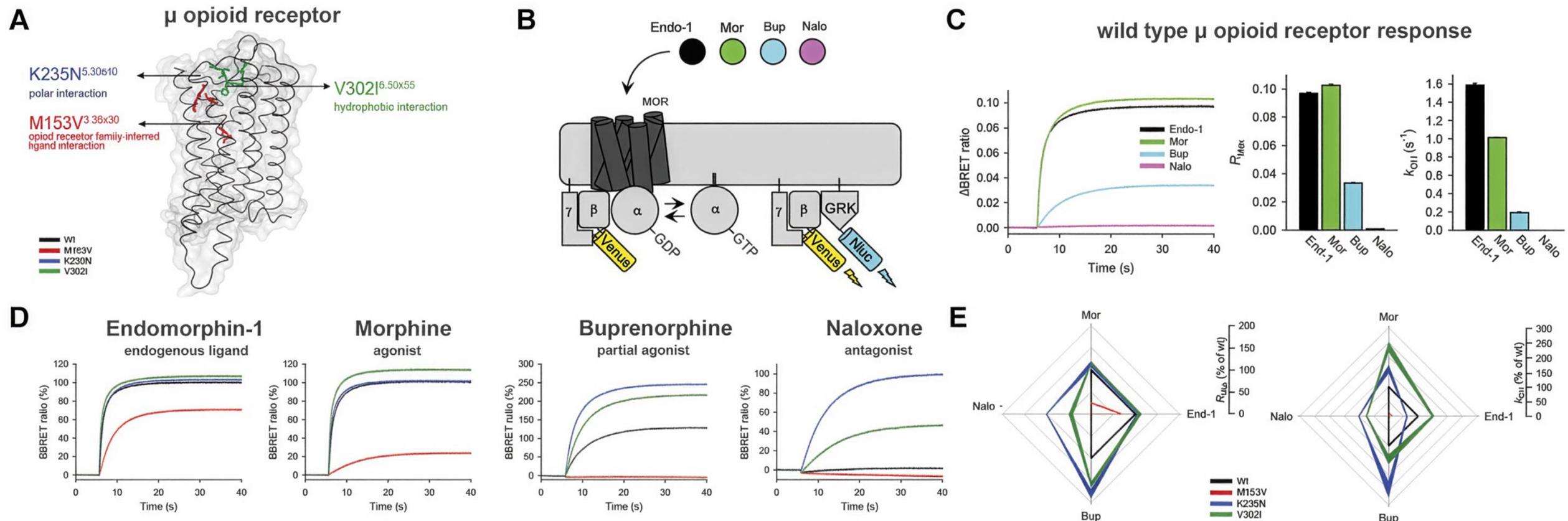
Family study



Gain of Function (GOF)
or
Loss of Function (LOF)

Impact of genetic variants on drug response (Pharmacogenetics)

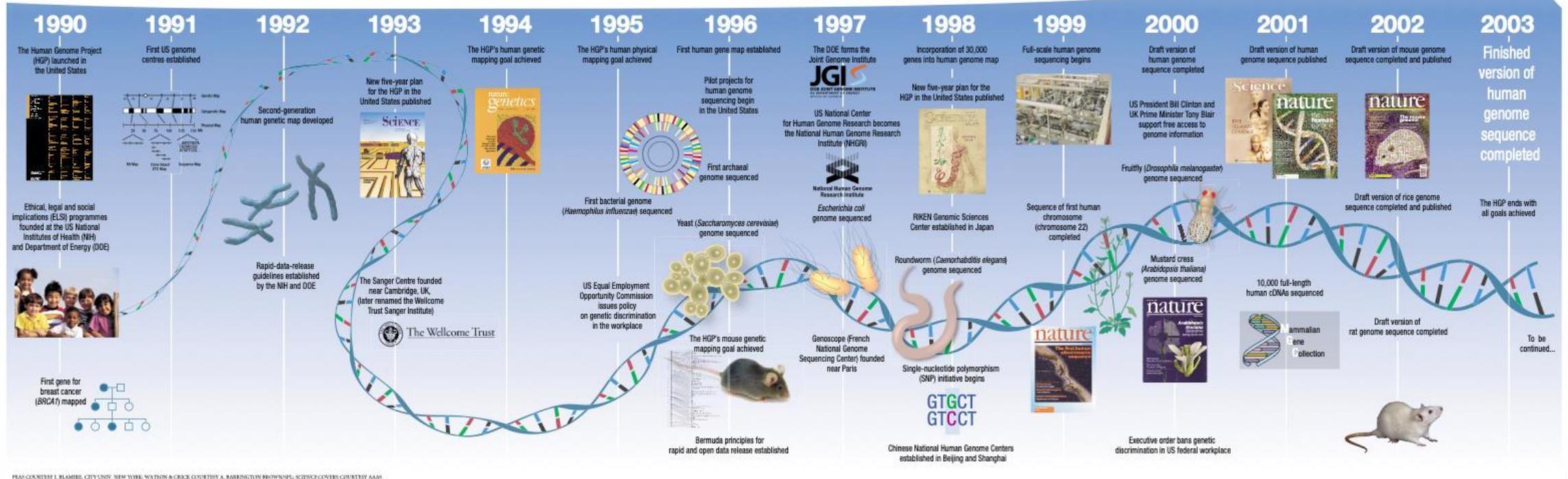
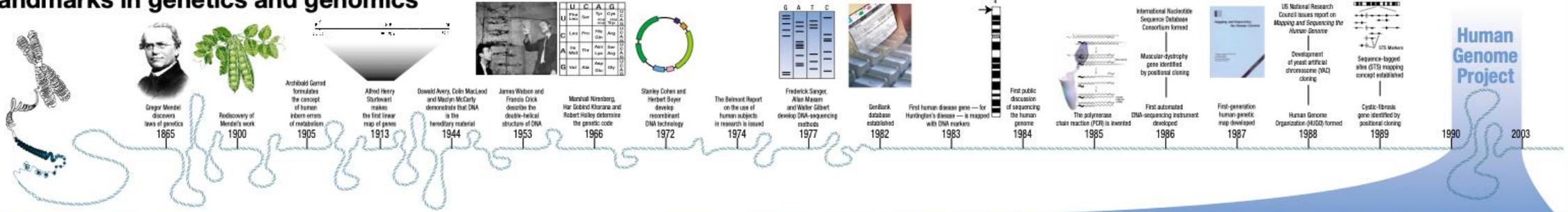
Impact of variants on functional response



Genetic variants can also affect compound metabolism, clearance etc...

Human Genome Project, Sequencing it all.

Landmarks in genetics and genomics



HEAD COURTESY J. BLAMIRE, CITY UNIV., NEW YORK; WATSON & CRICK COURTESY A. BARRINGTON BROWN/SCIENCE; SCIENCE COVERS COURTESY AAAS

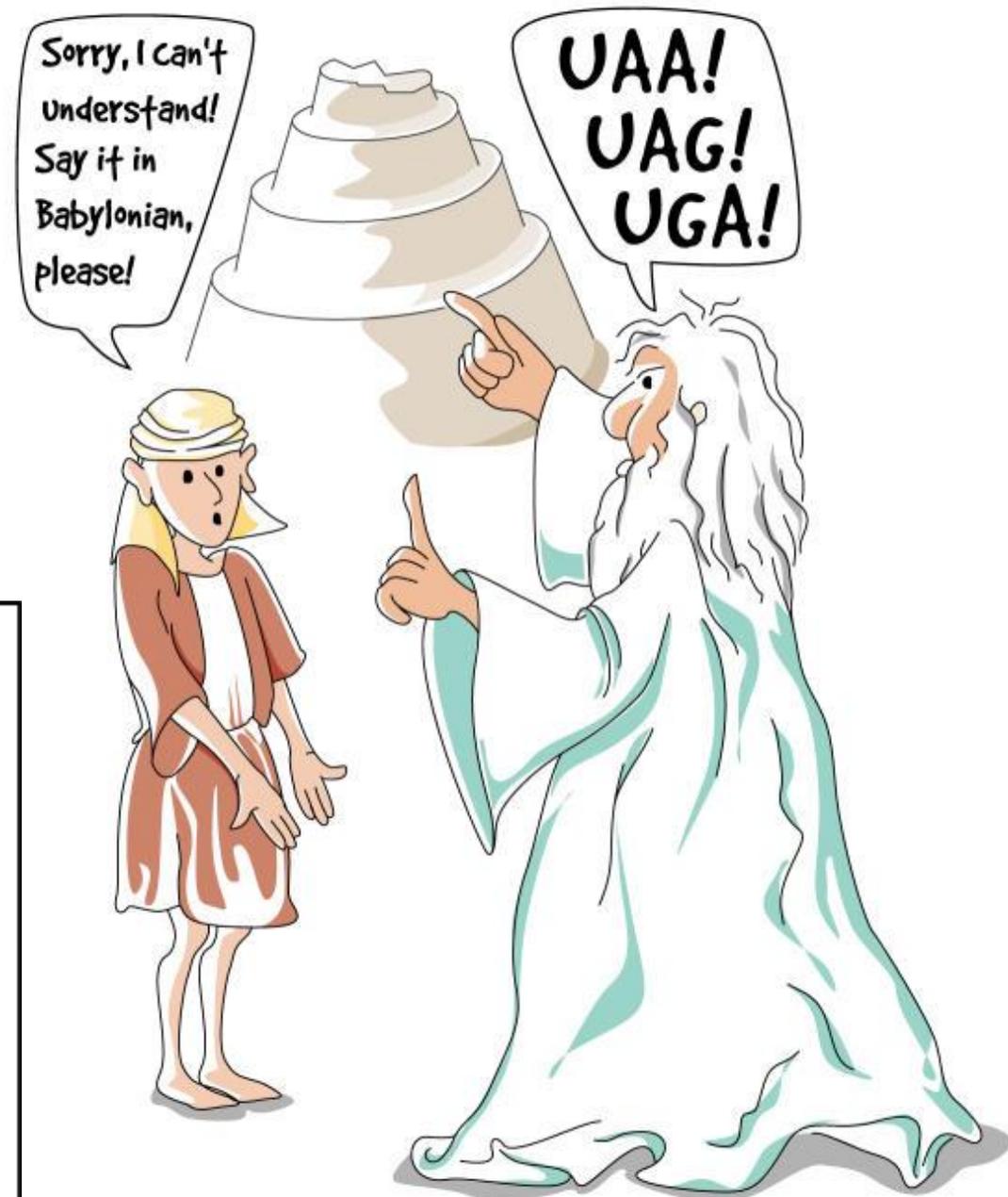
"It is humbling for me and awe inspiring to realise that we have caught the first glimpse of our own instruction book, previously known only to God." (Francis Collins, Project Director)

"We now have the possibility of achieving all we ever hoped for from medicine." (UK Science Minister Lord Sainsbury)

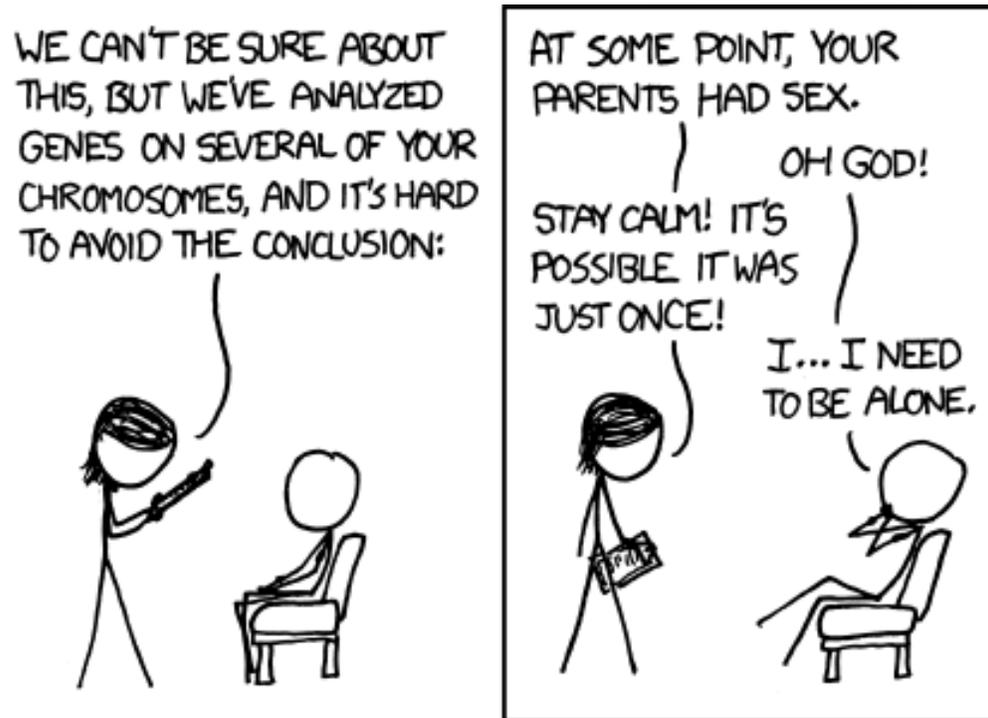
Or the limits of linear determinisms

•"As was predicted... getting the sequence will be the easy part... The hard part will be finding out what it means." — [Sydney Brenner](#).

(Nobel Prize of Medicine 2002, for their discoveries concerning genetic regulation of organ development and programmed cell death)



Weigman 2004, EMBO Reports



Genotype is one determinant of phenotype, it is not all ... or is it?



**Original
(Rainbow)**



**Clone
(Copy-Cat)**



What is a genome?

Viruses, Prokaryotes and Eukaryotes

Virus

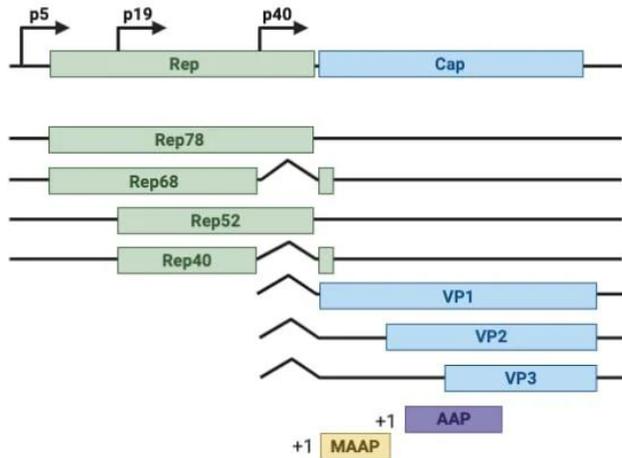
DNA or RNA

Single or double stranded

Data Compaction

-Small regulatory sequences

-Genes can sometimes overlap



Addgene

Prokaryotes (no nucleus, bacteria)

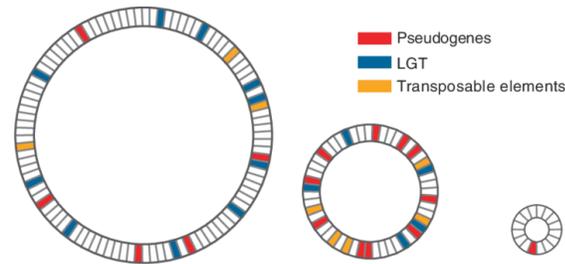
Double stranded DNA

Data Compaction

Mostly Naked

Circular

Trans-specific information exchange
(Plasmids)



	Free-living	Recent or facultative pathogen	Obligate symbiont or pathogen
Genome size	Large (5-10 MB)	Intermediate (2-5 MB)	Small (0.5-1.5 MB)
Number of pseudogenes	Few	Many	Rare
Incidence of LGT	Frequent	Frequent to rare	Rare to none
Selfish genetic elements	Few	Common	Rare
Genome organization	Stable or unstable	Unstable	Stable
Effective population size	Large	Small	Small

Ochman & Davalos, 2006 Science

Eukaryotes (cells with nucleus)

Double stranded DNA

Organized on chromosomes

Attached to proteins (chromatin)

Lots of non-coding sequences

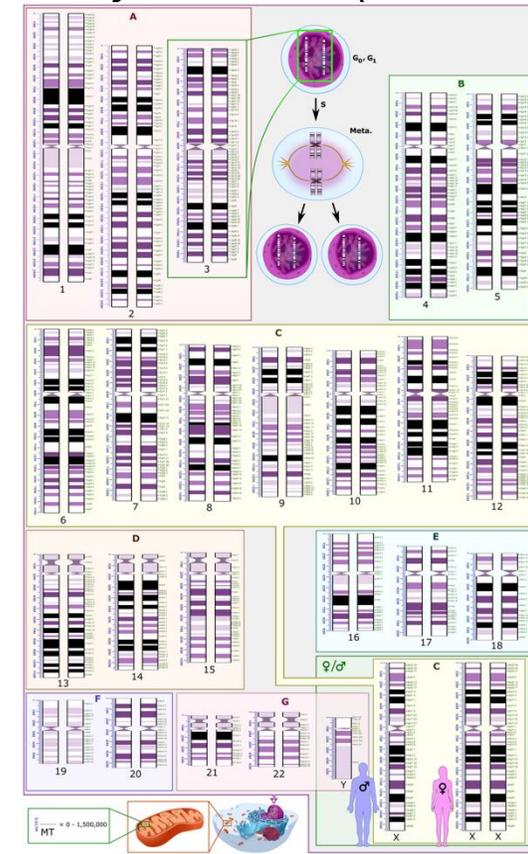
Vestigial genes

Gene duplication

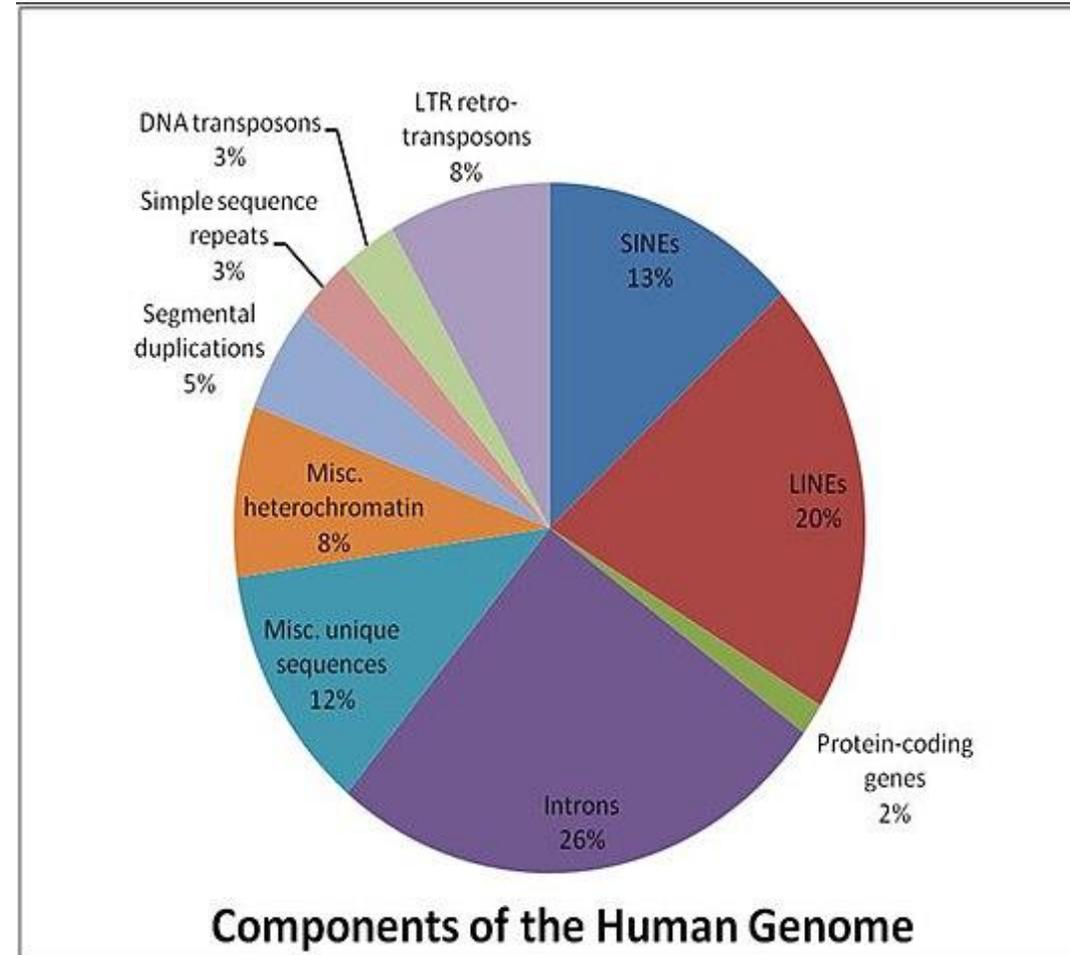
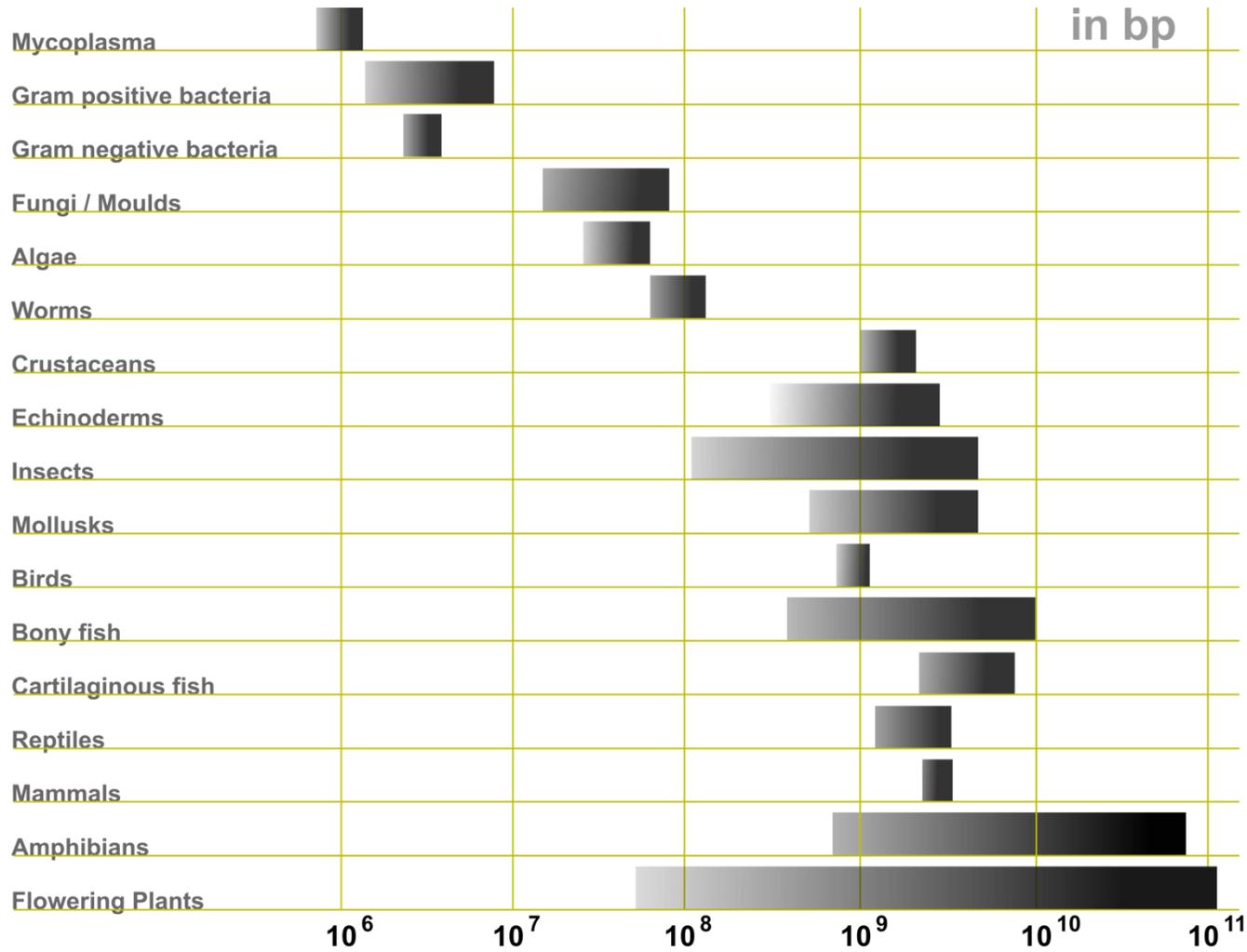
Heteroploidies

(more than one copies of each gene)

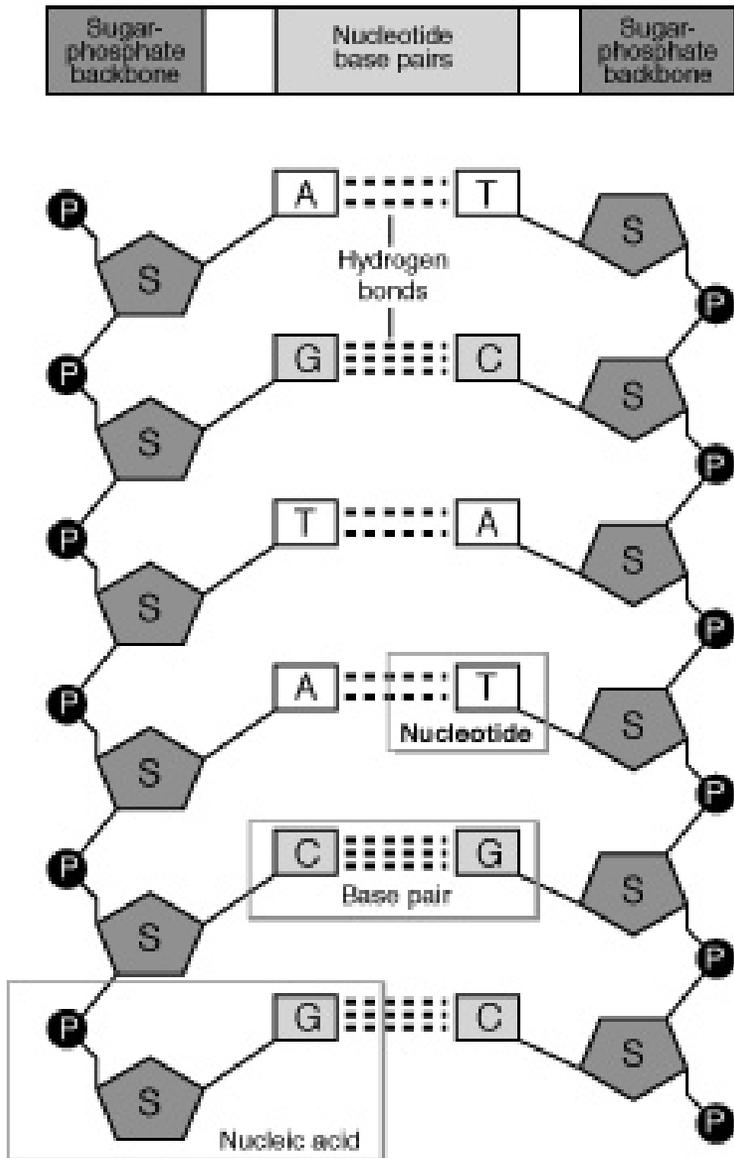
Secondary Genomes (mitochondria)



Genome size

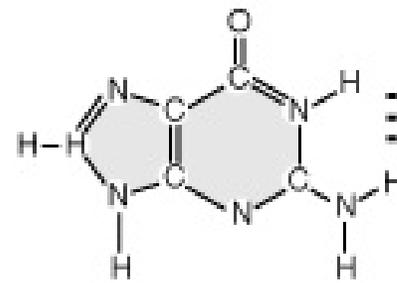


Deoxyribonucleic Acid (DNA)

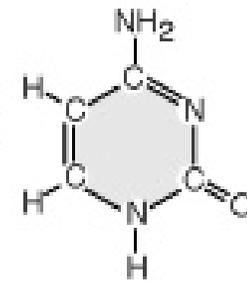


Nucleotides

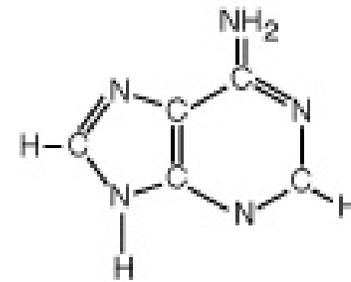
G Guanine



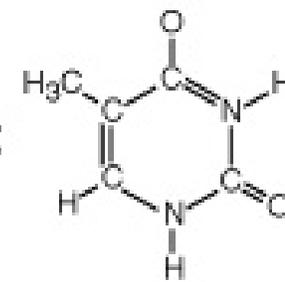
C Cytosine



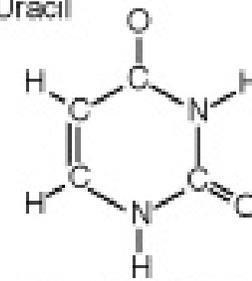
A Adenine



T Thymine



U Uracil



replaces Thymine in RNA

Why is nucleotide complementarity important ?

Biological Processes

Genome Replication

RNA synthesis

Protein Synthesis

Regulation of Gene expression

3D structure of the genome

RNA regulation

RNA transport

Methods

Gene Sequencing

PCR

In situ detection of nucleic acids (FISH)

Genes

Messengers

Modulation of Gene expression

Genome Engineering

Synthetic Biology.

Etc...

The complementarity of the bases is the mean to read and write information in nucleic acid.

DNA sequences encodes:

Genes

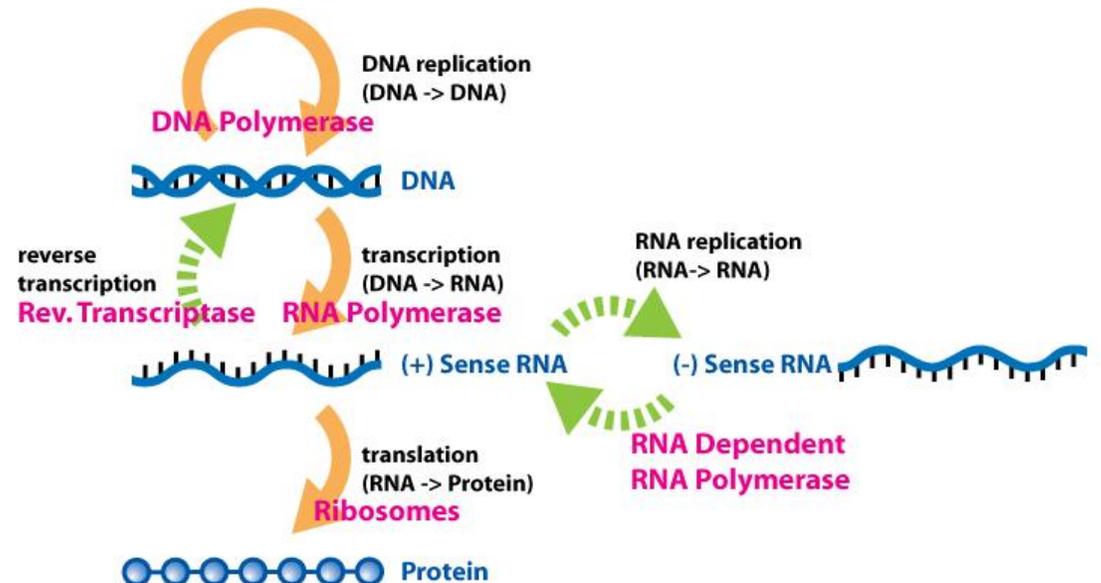
- The nature of the gene product : **Protein** or **RNA**
- The Expression level of genes: **Yes/No, How Much, When** and **Where**.

Genome

- Genome maintenance: Structure, Replication
- Junk, parasites or drivers of evolution ?

This information is encoded in multiple languages.
The genetic code use for protein synthesis is just one of these languages.

The Dogma:



General Genome Databases NCBI from NIH

www.ncbi.nlm.nih.gov

An official website of the United States government [Here's how you know](#)

NIH National Library of Medicine
National Center for Biotechnology Information

Log in

All Databases Search

- NCBI Home
- Resource List (A-Z)
- All Resources
- Chemicals & Bioassays
- Data & Software
- DNA & RNA
- Domains & Structures
- Genes & Expression
- Genetics & Medicine
- Genomes & Maps
- Homology
- Literature
- Proteins
- Sequence Analysis
- Taxonomy
- Training & Tutorials
- Variation

Welcome to NCBI

The National Center for Biotechnology Information advances science and health by providing access to biomedical and genomic information.

[About the NCBI](#) | [Mission](#) | [Organization](#) | [NCBI News & Blog](#)

Submit

Deposit data or manuscripts into NCBI databases



Download

Transfer NCBI data to your computer



Learn

Find help documents, attend a class or watch a tutorial



Develop

Use NCBI APIs and code libraries to build applications



Analyze

Identify an NCBI tool for your data analysis task



Research

Explore NCBI research and collaborative projects



Popular Resources

- PubMed
- Bookshelf
- PubMed Central
- BLAST
- Nucleotide
- Genome
- SNP
- Gene
- Protein
- PubChem

NCBI News & Blog

- BankIt Submitters: Upcoming Changes to How You Submit to GenBank
27 Jan 2026
- Are you a GenBank submitter? Do you use BankIt or the GenBank app in the...
14 Jan 2026
- GenBank Now Supports EGAPx-Based Annotation
14 Jan 2026
- With the latest release of EGAPx we're excited to announce...
13 Jan 2026
- An Updated Bacterial and Archaeal Reference Genome Collection is Available!
13 Jan 2026
- Download the updated bacterial and...

General Genome Databases Ensembl EMBL

www.ensembl.org/index.html

Ensembl BLAST/BLAT | VEP | Tools | BioMart | Downloads | Help & Docs | Blog

Search all species...

Tools **BioMart >** **BLAST/BLAT >** **Variation Effect Predictor >**

[All tools](#) Export custom datasets from Ensembl with this data-mining tool Search our genomes for your DNA or protein sequence Analyse your own variants and predict the functional consequences of known and unknown variants

Search

All species for

Go

e.g. **BRCA2** or **rat 5:62797383-63627669** or **rs699** or **coronary heart disease**

All genomes

-- Select a species --

Pig breeds
Pig reference genome and 20 additional breeds

[View full list of all species](#)

Favourite genomes

Human
GRCh38.p14
[Still using GRCh37?](#)

Mouse
GRCm39

Zebrafish
GRCz11

Ensembl is a public and open project providing access to genomes, annotations, tools and methods. Its goal is to enable genomic science by providing high-quality, integrated and consistent annotation on all cellular genomes within a harmonious, scalable and accessible infrastructure.

Ensembl Release 115 (September 2025)

- ~121,000 new protein-coding transcripts have been added to the GRCh38 human reference gene set.
- Two new breeds of cattle have been added: UOA_Tuli_1 and UOA_Wagyu_1
- The sheep reference has been updated to ARS-UI_Ramb_v3.0
- Two new export modes are now available for Newick trees

[More release news](#) on our blog

Ensembl Rapid Release

New genome assemblies are now being released to the [Ensembl Beta site](#).

All Rapid Release data, including release 65, has been uploaded into the new Ensembl Beta site.

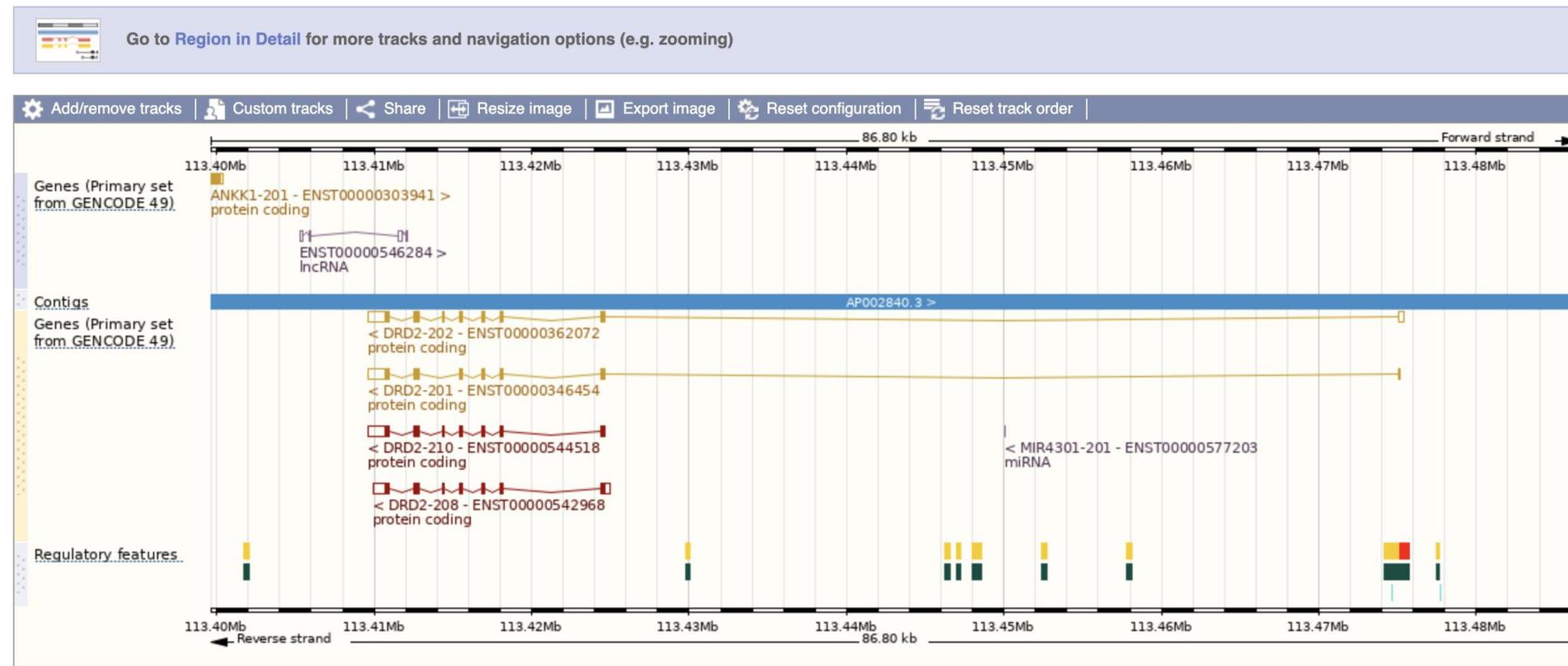
The Ensembl Rapid Release website will remain active for the foreseeable future, however, the data and species set will no longer be updated.

[Find out more on our blog](#)

Compare genes across species Find SNPs and other variants for my gene Gene expression in different tissues Retrieve gene sequence Find a Data Display Use my own data in Ensembl

Example of gene entry (human interface)

Name	DRD2 (HGNC Symbol)
MANE	This gene contains MANE Select ENST00000362072 , ENSP00000354859
UniProtKB	This gene has proteins that correspond to the following UniProtKB identifiers: P14416
RefSeq	This Ensembl/Gencode gene contains transcript(s) for which we have selected identical RefSeq transcript(s) . If there are other RefSeq transcripts available they will be in the External references table
CCDS	This gene is a member of the Human CCDS set: CCDS8361.1 , CCDS8362.1
Ensembl version	ENSG00000149295.15
Other assemblies	This gene maps to 113,280,327-113,347,124 in GRCh37 coordinates. View this locus in the GRCh37 archive: ENSG00000149295
Gene type	Protein coding
Annotation method	Annotation for this gene includes both automatic annotation from Ensembl and Havana manual curation, see article .



There are several data bases one can use (an AI generated list).

Core nucleotide/genome sequence databases

- GenBank (NCBI) – primary public DNA sequence archive, part of the INSDC with daily exchange with ENA and DDBJ.
- European Nucleotide Archive (ENA, EMBL-EBI) – Europe’s partner archive for nucleotide sequences.
- DNA Data Bank of Japan (DDBJ) – Japan’s partner archive within INSDC.
- NCBI Genome resources – genome assemblies, annotations, variation viewers, BioProject/BioSample.

Human gene, variant, and GWAS resources

- GeneCards – integrated human gene database aggregating genomic, transcriptomic, proteomic, genetic and clinical info.
- dbSNP (NCBI) – catalog of short genetic variants.
- dbVar (NCBI) – structural variation database (CNVs, insertions, deletions, inversions, etc.).
- dbGaP (NCBI) – genotype–phenotype database for controlled-access human studies.

NHGRI–EBI GWAS Catalog – curated catalog of published human GWAS

General genome resource lists

- NHGRI “Online Research Resources” – curated list of major genome-related databases and tools.
- NCBI “All Resources” guide – consolidated entry point to NCBI genomic and related databases.
- U of T Genome Databases guide – field-specific links including model organism genome resources.
(<https://guides.library.utoronto.ca/c.php?g=251933&p=1675867>)

Functional genomics and signatures

- GEO (Gene Expression Omnibus) – public repository for array- and sequence-based functional genomics data.
- GEO Datasets – curated GEO datasets with analysis tools.
- MSigDB (Molecular Signatures Database) – curated gene set collections for GSEA.

Population genomics reference resources

- 1000 Genomes / International Genome Sample Resource (IGSR) – deep catalog of common human genetic variation, ongoing maintenance and mapping to current assemblies.

Psychiatric Genomics Consortium (PGC) – Central hub for GWAS of schizophrenia, bipolar disorder, MDD, ADHD, ASD, PTSD, Tourette, OCD, etc.; provides extensive summary-statistics downloads per disorder and cross-disorder analyses.

Ageing and longevity genomics

- Human Ageing Genomic Resources (HAGR) – umbrella for ageing-related databases (GenAge, AnAge, GenDR, LongevityMap).
- LongevityMap – human variants associated with longevity.

Biobanks and population cohorts (with genomic data)

- “All of Us” Research Program (NIH) – large population biobank with genomic and environmental data.
- Biobank Graz, China Kadoorie Biobank, Estonian Biobank, EuroBioBank, UKBiobank etc. – large international biobanks with biospecimens and associated data.

Genome Geography, vocabulary

Locus: a location on the genome

Chromosome: A large nuclear structure that contains a long molecule of DNA associated to proteins (chromatin). Chromosomes encode subsets of the genome.

Gene: two meanings. The Mendelian gene is a basic unit of heredity. The molecular gene is a sequence in DNA that is transcribed to produce RNA.

Single Nucleotide Polymorphisms (SNPs): Variations between multiple genomes within a species

Haplotypes: a set of closely linked DNA variations (like SNPs or alleles) on a single chromosome that are inherited together from a single parent

Haplogroup: haplotypes inherited from a single parent. (e.g.: y chromosome-father, mtDNA-mother) This is used to define ancestry.

HapMap/1000 Genomes. These are databases of human genetic variations based of haplotype frequency.

The International Genome Sample Resource

The 1000 Genomes Project created a catalogue of common human genetic variation, using openly consented samples from people who declared themselves to be healthy. The reference data resources generated by the project remain heavily used by the biomedical science community.

The International Genome Sample Resource (IGSR) maintains and shares the human genetic variation resources built by the 1000 Genomes Project. We also update the resources to the current reference assembly, add new data sets generated from the 1000 Genomes Project samples and add data from projects working with other openly consented samples.

Explore the data sets in IGSR through our data portal

Access HGSCV data

View variants in genomic context in Ensembl

Sample	Populations
HGI01002	Individuals in West Africa: The Gambia - Mandinka
HGI01003	Mixed in Sierra Leone
HGI01017	East in Nigeria
NA19100	Colombian in Medellin, Colombia
GM19020	Peruvian in Lima, Peru
HGI00999	Sri Lankan Tamil in the UK
HGI00998	White British in the UK
NA19424	Mixed in Korea
GM19129	Yoruba in Ibadan, Nigeria
HGI00999	Japanese in Tokyo, Japan

Population	A	G	AIA	AIG
ASW	0.230 (28)	0.770 (94)	0.066 (4)	0.328 (20)
ESN	0.066 (13)	0.934 (185)	0.010 (1)	0.111 (11)
GWD	0.066 (15)	0.934 (211)	0.009 (1)	0.115 (13)
YRI	0.079 (17)	0.921 (199)	0.157 (17)	0.843 (91)
AMR	0.365 (253)	0.635 (441)	0.147 (51)	0.435 (151)
CLM	0.452 (85)	0.548 (103)	0.223 (21)	0.457 (43)

Latest Announcements

Wednesday July 23, 2025

New publications describing variants from the 1000 Genomes Project samples from HGSCV and collaborators

[Complex genetic variation in nearly complete human genomes](#) by Logsdon, G.A., et al describes an extensive catalog of variation from the near complete assemblies of 65 human genomes from diverse 1000 Genomes Project samples. These high quality assemblies enable a more comprehensive insight into all variant types, including those within complex regions.

Data is available via the [HGSCV3](#) collection.

[Structural variation in 1,019 diverse humans based on long-read sequencing](#) by Schloissnig, S et al describes the characterisation of structural variants in 1019 samples from 26 different the 1000 Genomes Project populations. This study used intermediate-coverage long read sequencing and a novel integration of linear and graph genome-based analyses.

Data is available via the [1KG_ONT_VIENNA](#) collection.

[All announcements](#)

Genetic Vocabulary

Mutation: New variant (loaded term, use with care)

Penetrance: the proportion of individuals with a specific genotype (genetic variant) who express the corresponding, observable trait or phenotype

Expressivity: degree to which a specific genotype manifests as a particular phenotype (observable trait or disease) in an individual,

Types of genetic studies:

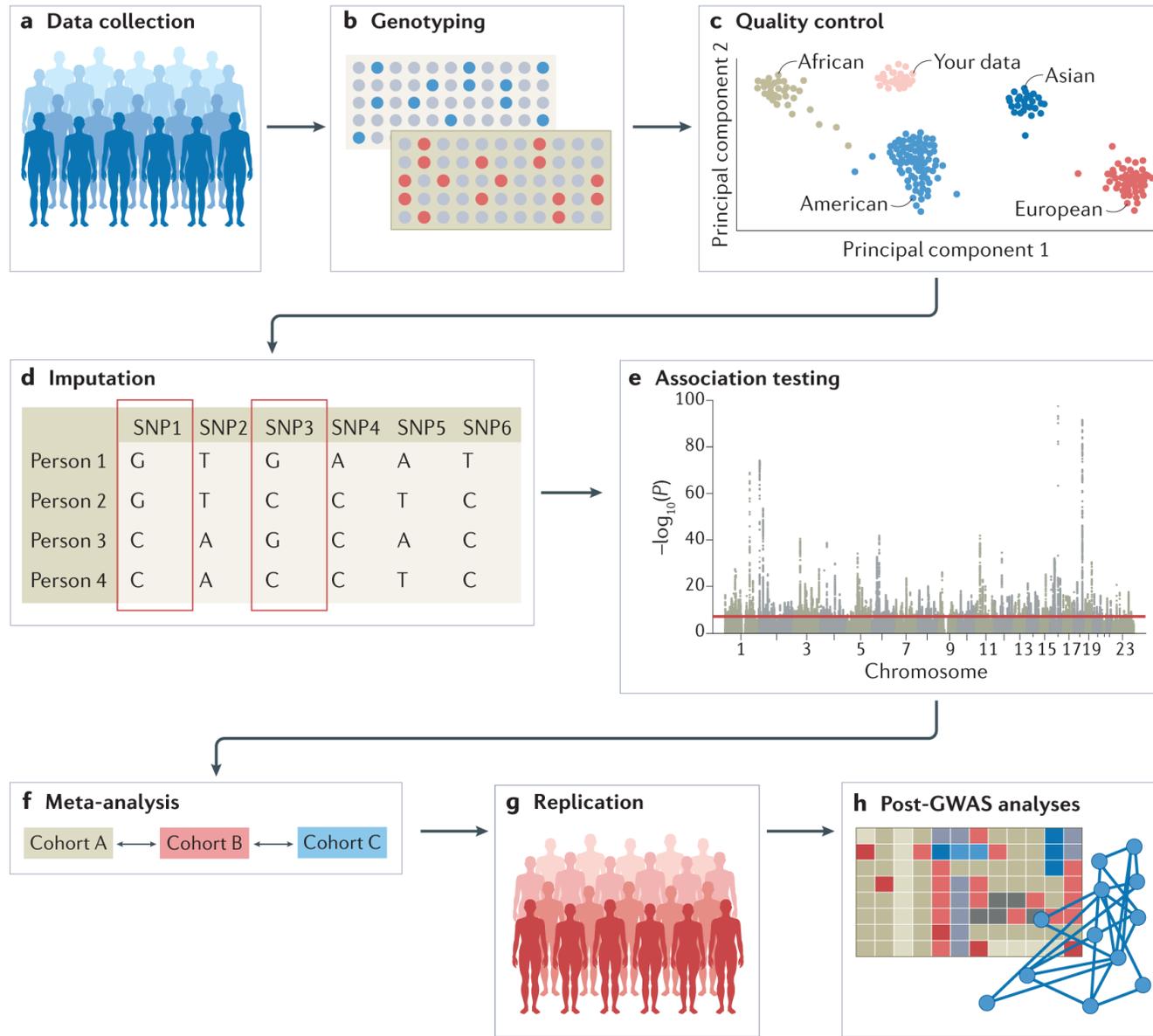
-**Family study:** Study members of a family to identify variants that correlate with a phenotype.

-**Rare variants:** Identify a variant (less than 1% of population (often new) that correlates with the apparition of a phenotype. Generally, these variants have a strong effect on the phenotype.

-**Genome Wide Association Study (GWAS):** Study of the association of **multiple SNPs** covering the whole genome with a given continuous (e.g. BMI) or non-continuous (disease diagnostic) outcome .

-**Transcript Wide Association Study (TWAS):** Study of the association of multiple **RNA Transcripts** covering the whole genome with a given continuous (e.g. BMI) or non- continuous (disease diagnostic) outcome .

Genome Wide Association Studies (GWAS)



Overall ,a GWAS is a collection of regressions

In GWAS, one typically fit a separate regression model per variant. A common basic formula for a quantitative trait is a linear regression:

$$Y_i = \beta_0 + \beta_1 G_i + \beta_2 C_{i1} + \dots + \beta_k C_{ik} + \varepsilon_i$$

- Y_i : phenotype (e.g., height) of individual i .
- G_i : genotype dosage for the tested SNP in individual i (coded 0, 1, 2 for copies of the effect allele).
- C_{ij} : covariates such as age, sex, principal components of ancestry, batch, etc.
- β_1 : SNP effect size (per-allele change in phenotype).
- ε_i : residual error term.

A case–control GWAS compares genetic variants between people with a disease (cases) and people without it (controls) to find variants associated with disease risk. This is not a continuous variable

Case–control GWAS, use logistic regression, for example:

$$\text{logit}(P(Y_i = 1)) = \beta_0 + \beta_1 G_i + \beta_2 C_{i1} + \dots + \beta_k C_{ik}$$

Here Y_i is disease status (0 = control, 1 = case),

$\exp(\beta_1)$ is the odds ratio per effect allele.

Can this be improved by ML in place of running ML on top of it?

For a good read on GWAS go see:

<https://doi.org/10.1038/s43586-021-00056-9>

Limitations of GWAS:

Sample composition: Population bias

Quality of the measures: What is the diagnostic

Pooling or not pooling: Heterogeneity inside groups

Most SNPs are not coding: No obvious functional outcome

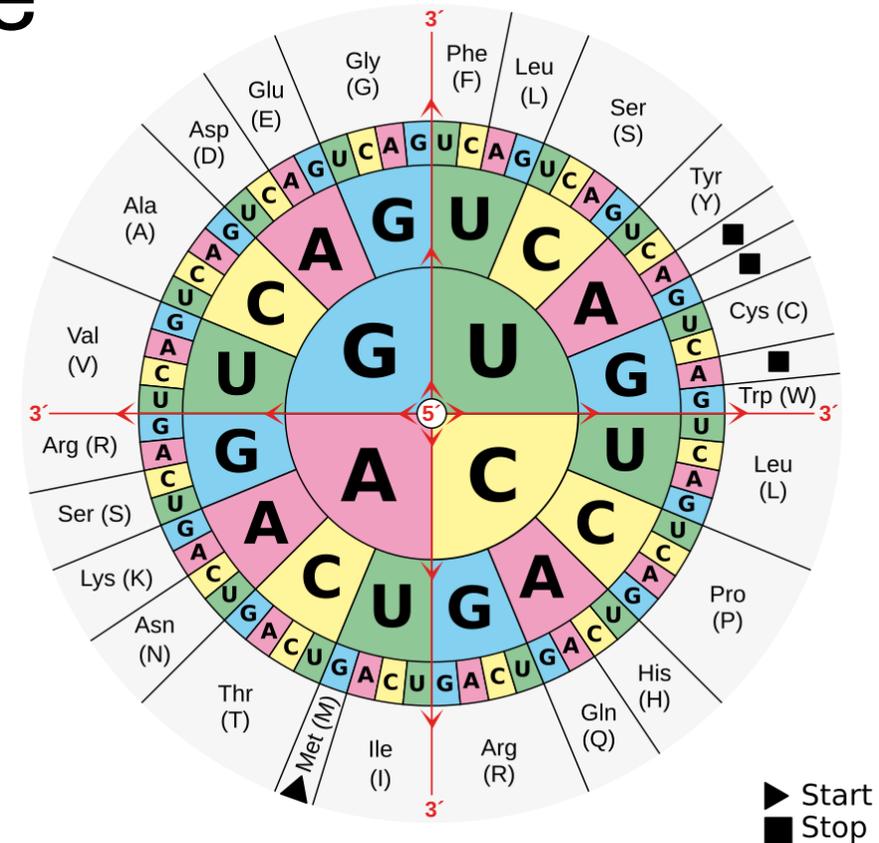
Identification of the affected gene is not always clear.

Polygenic Risk ...

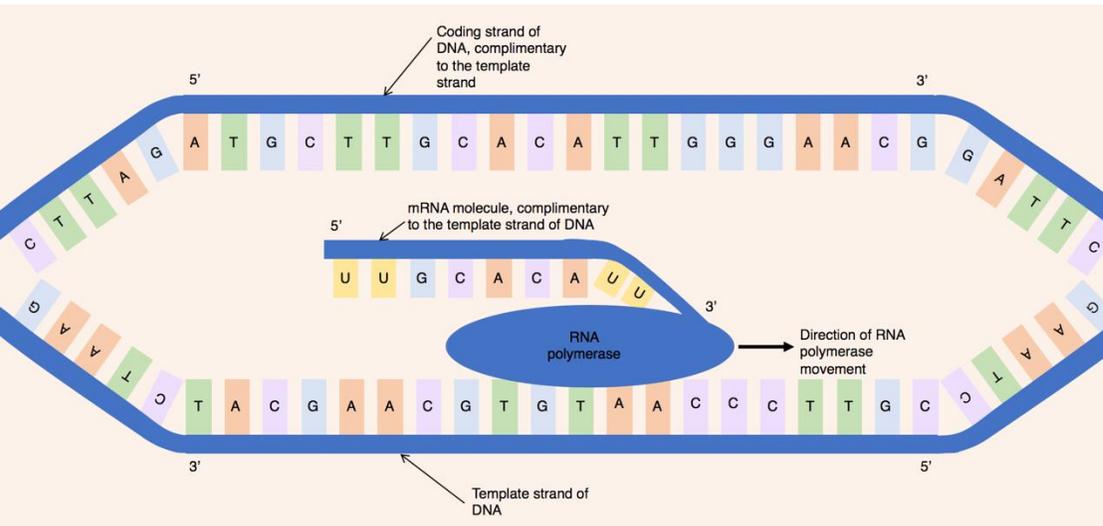
Part 2: Languages of the Genome

Some languages of the genome.

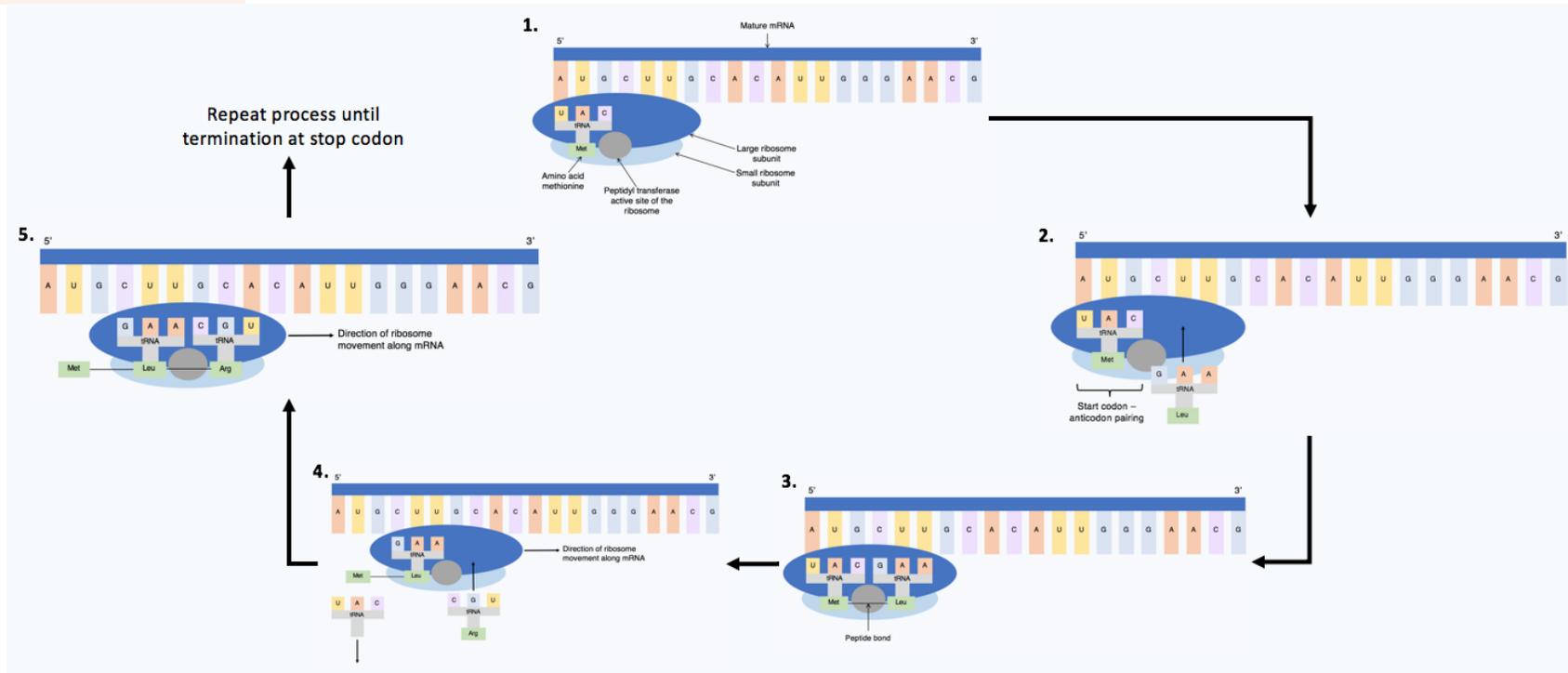
- Protein Coding
- mRNA splicing
- Regulation of Gene expression
- Regulation of Chromatin Structure
- Regulation mRNA translation



Transcription (nucleus)

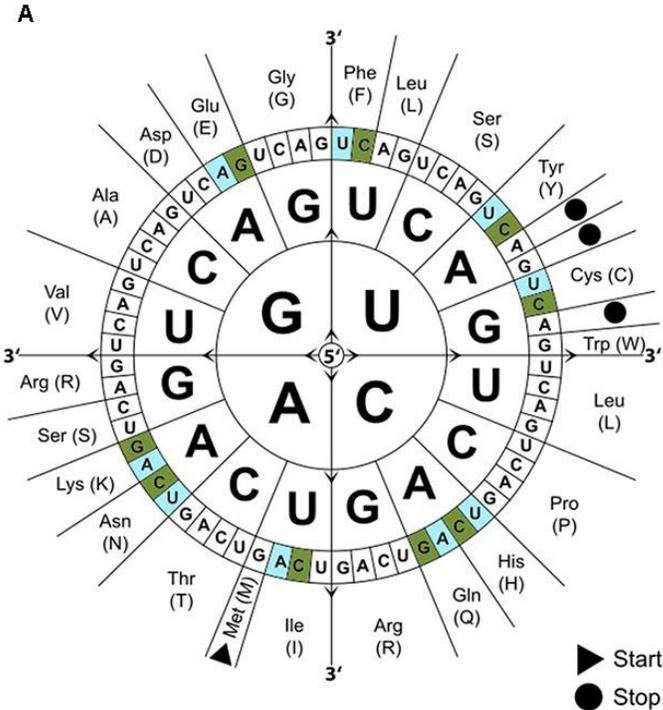


Translation (cytoplasm)



The codon language

Frequency
Can be regulatory

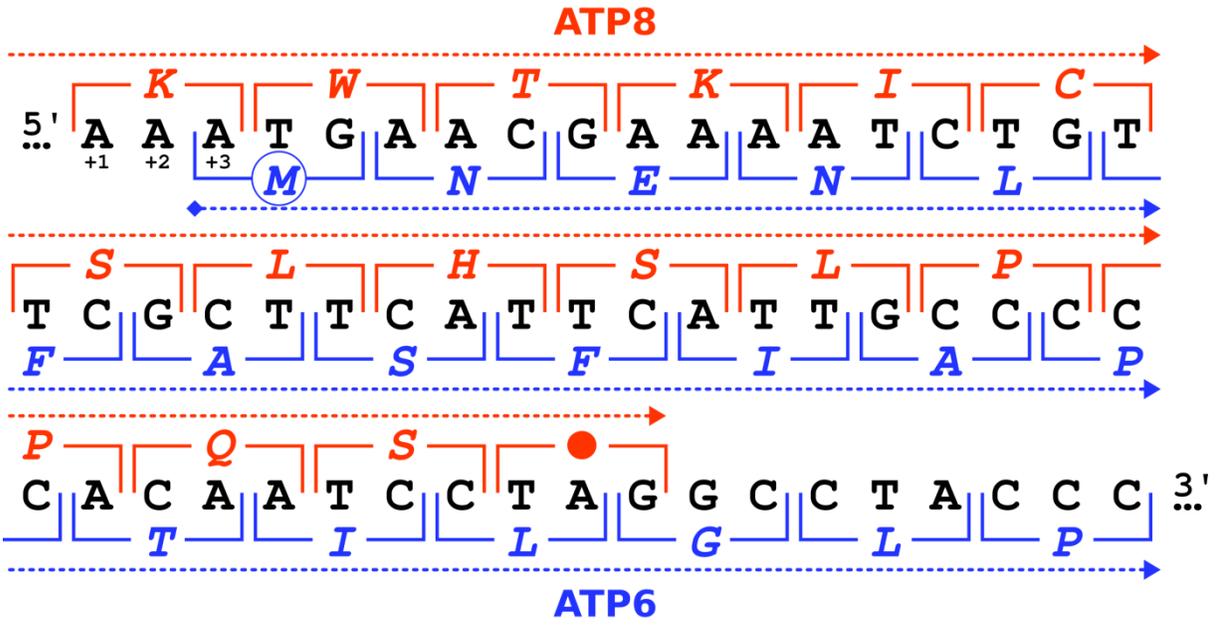


B

Amino acid	Codon	<i>P. patens</i>	<i>A. thaliana</i>	<i>S. cerevisiae</i>	<i>S. pombe</i>	<i>H. Sapiens</i>
Cys	UGC	+	-	+	+	+
	UGU	-	-	+	-	-
Glu	GAA	-	-	-	-	-
	GAG	+	+	-	+	+
Phe	UUC	+	+	+	+	+
	UUU	-	-	-	-	-
His	CAC	+	+	+	+	+
	CAU	-	-	-	-	-
Ile	AUA	-	-	/	-	-
	AUC	+	+	/	-	+
	AUU	-	-	-	-	-
Lys	AAA	-	-	-	-	-
	AAG	+	+	+	+	+
Asn	AAC	+	+	+	+	+
	AAU	-	-	-	-	-
Gln	CAA	-	-	-	+	-
	CAG	+	+	+	-	+
Tyr	UAC	+	+	+	+	+
	UAU	-	-	-	-	-

Reading Frame

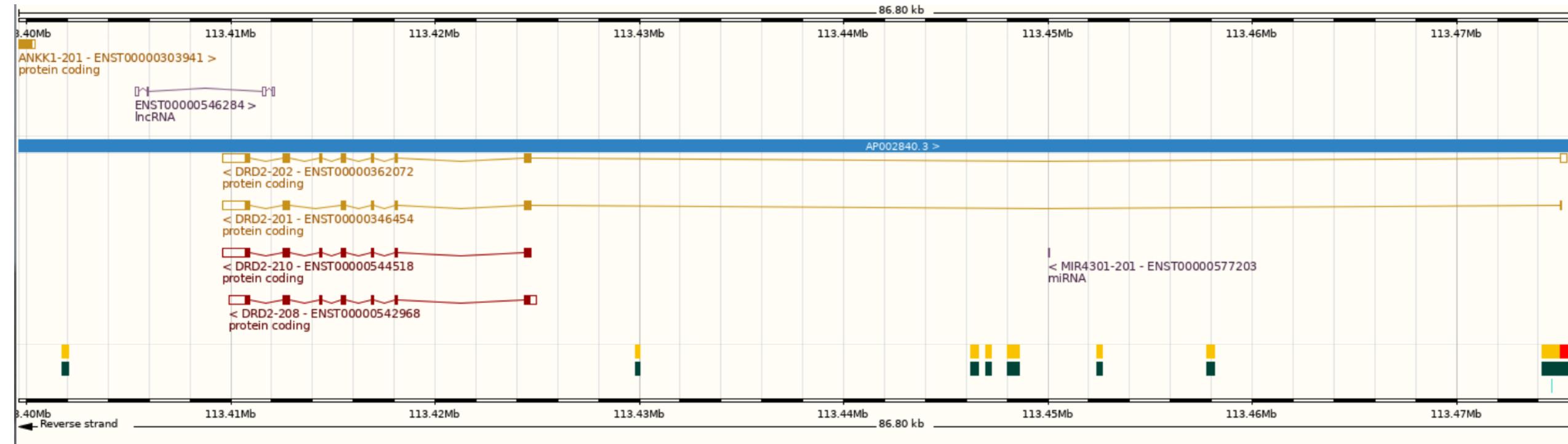
e.g. Mitochondrial ATP synthase



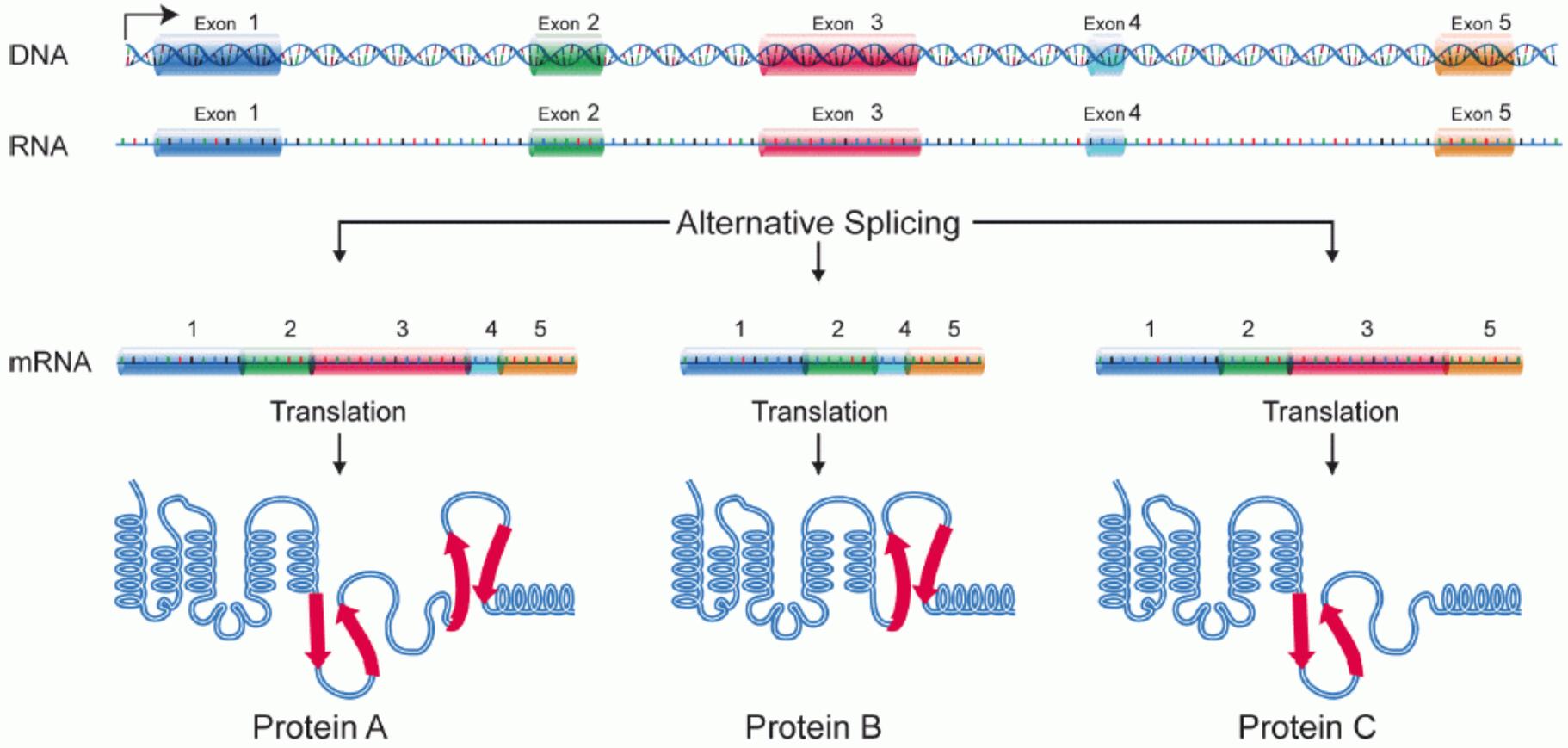
RNA Splicing

Exon: encodes protein

Intron: does not encode protein

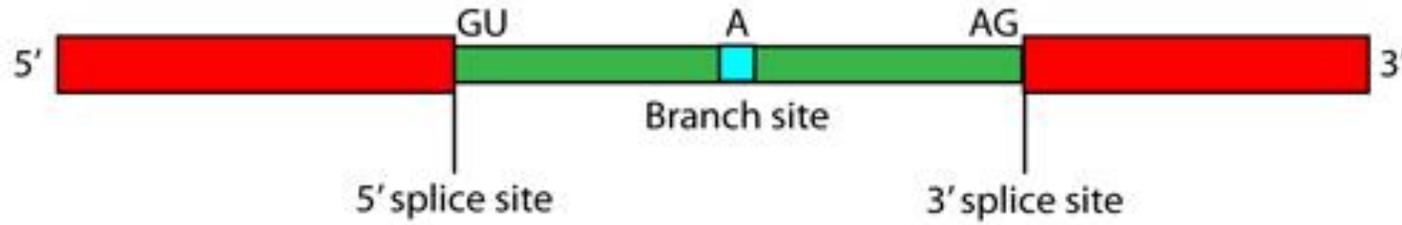


Alternative splicing

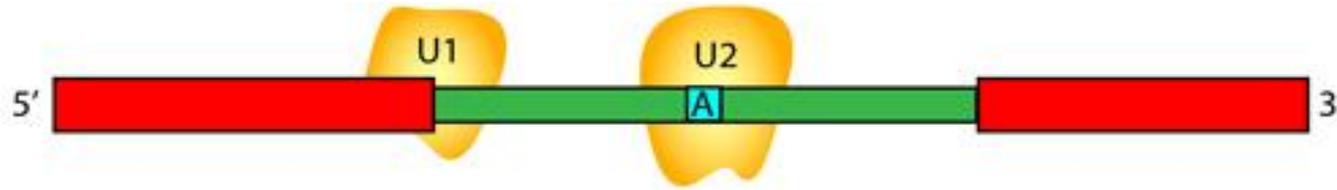


Splicing code

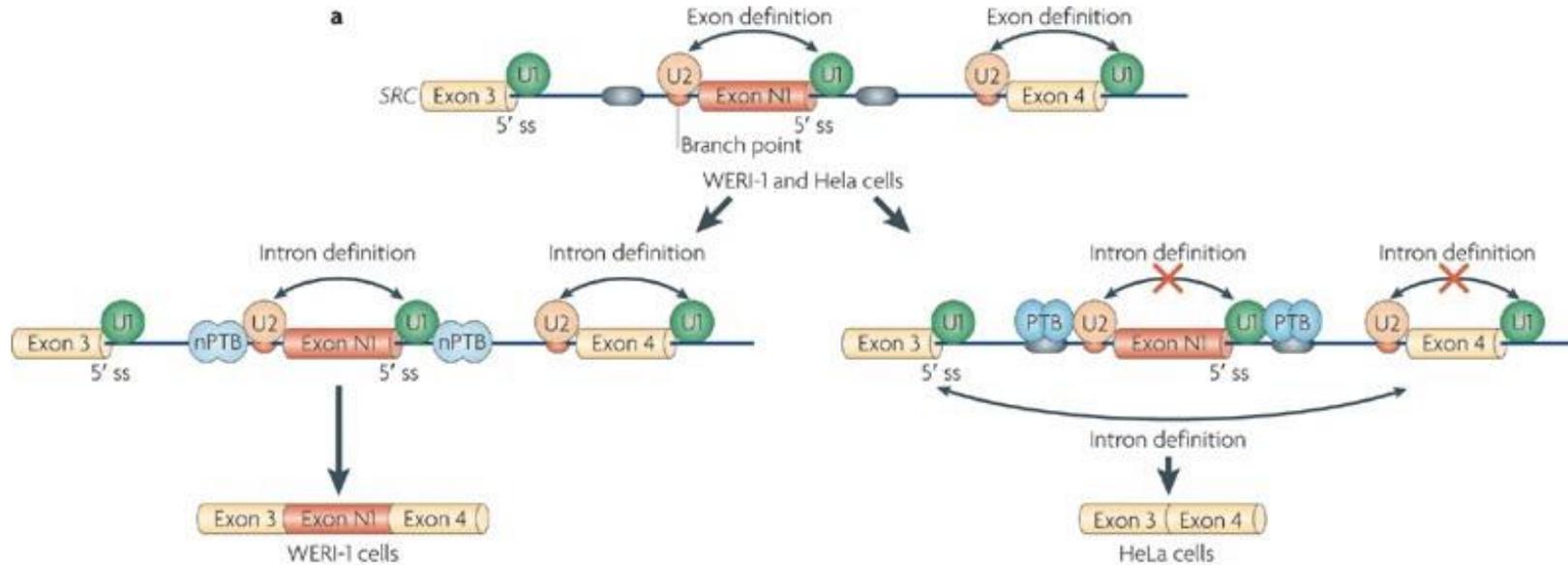
Located 18–40 nucleotides upstream of the 3' splice site. In mammals $yUnAy$. $y = C/U$ $n = \text{any base}$



The Code



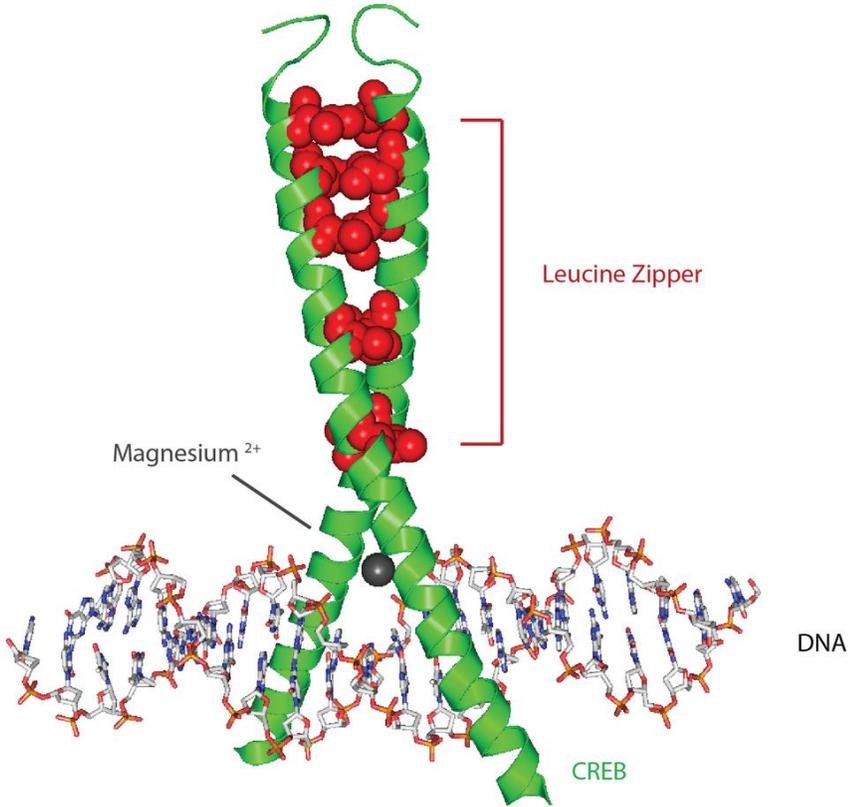
The readers (ribonucleoproteins)



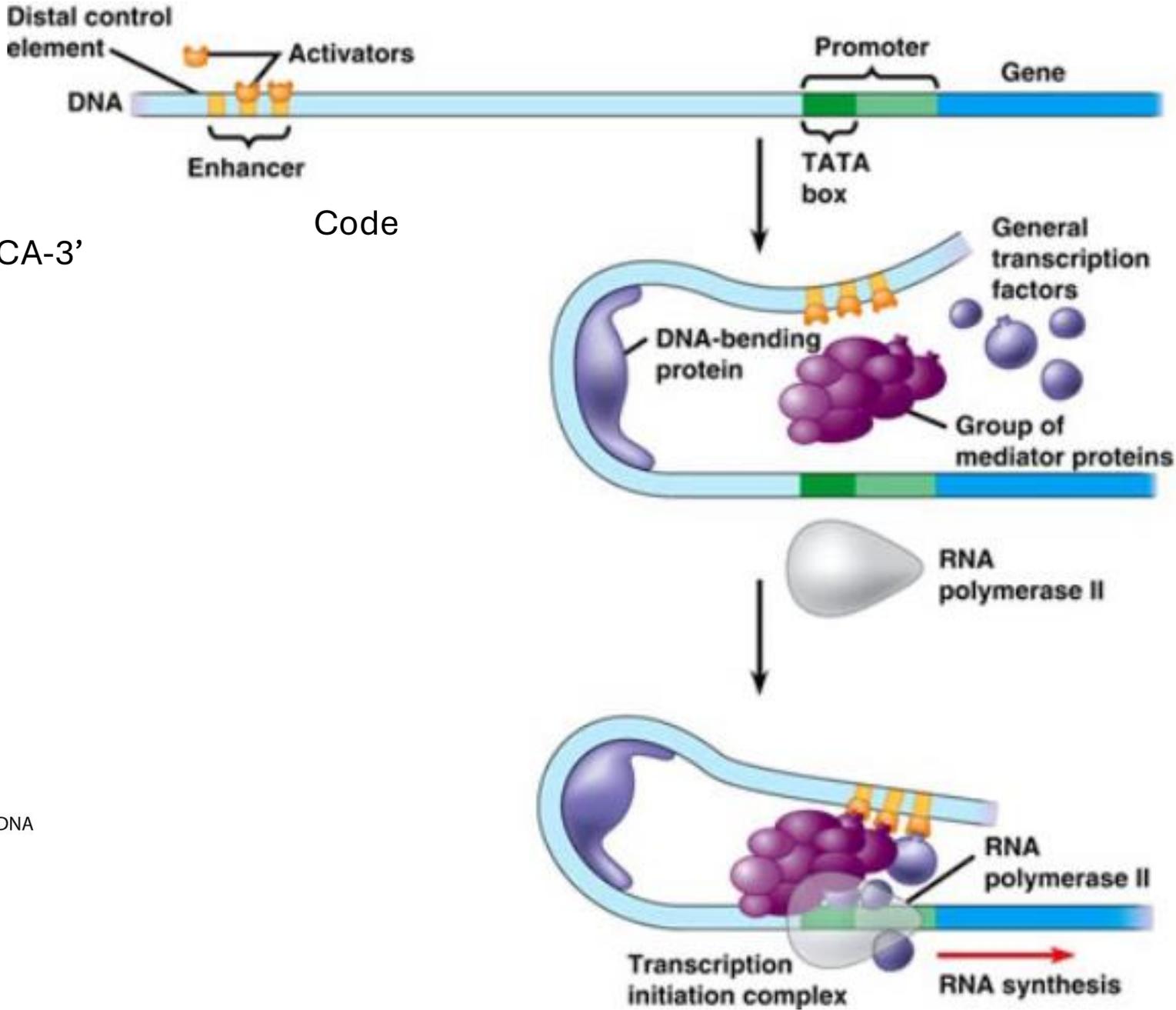
Alternative regulation

Regulating expression Languages of transcription

e.g. cAMP Responsive Element: 5'-TGACGTCA-3'



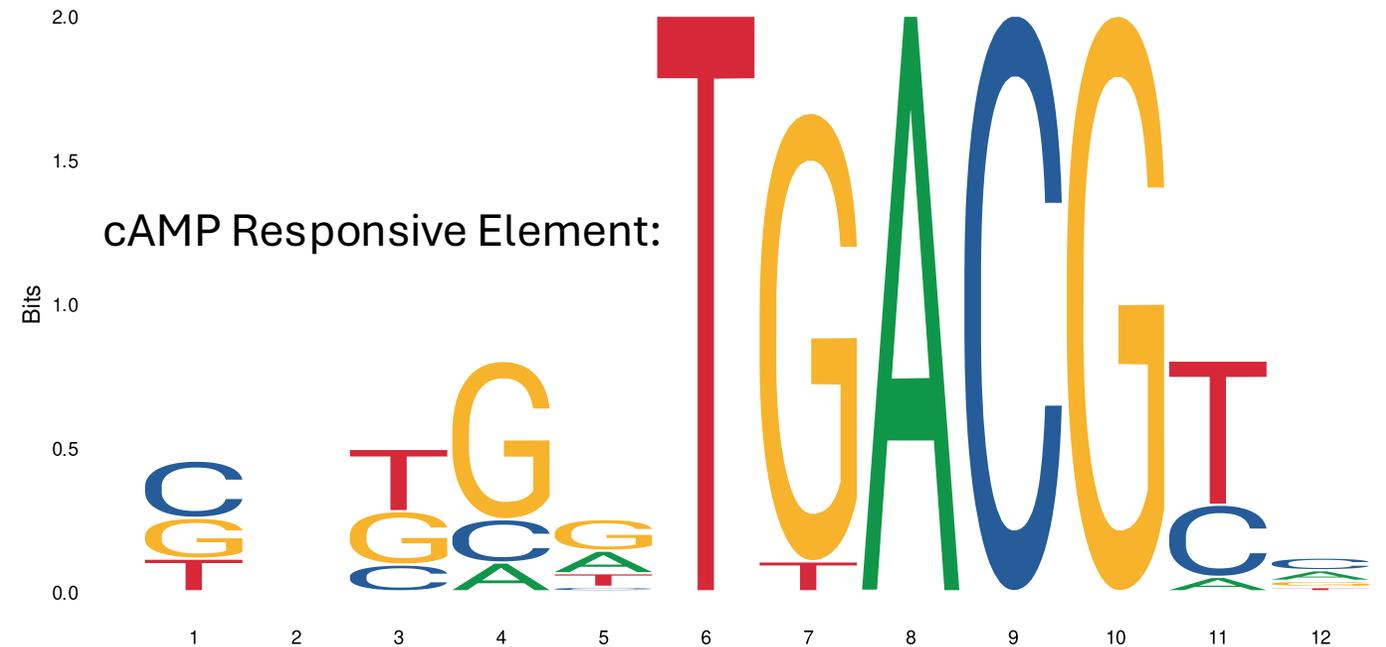
Reader



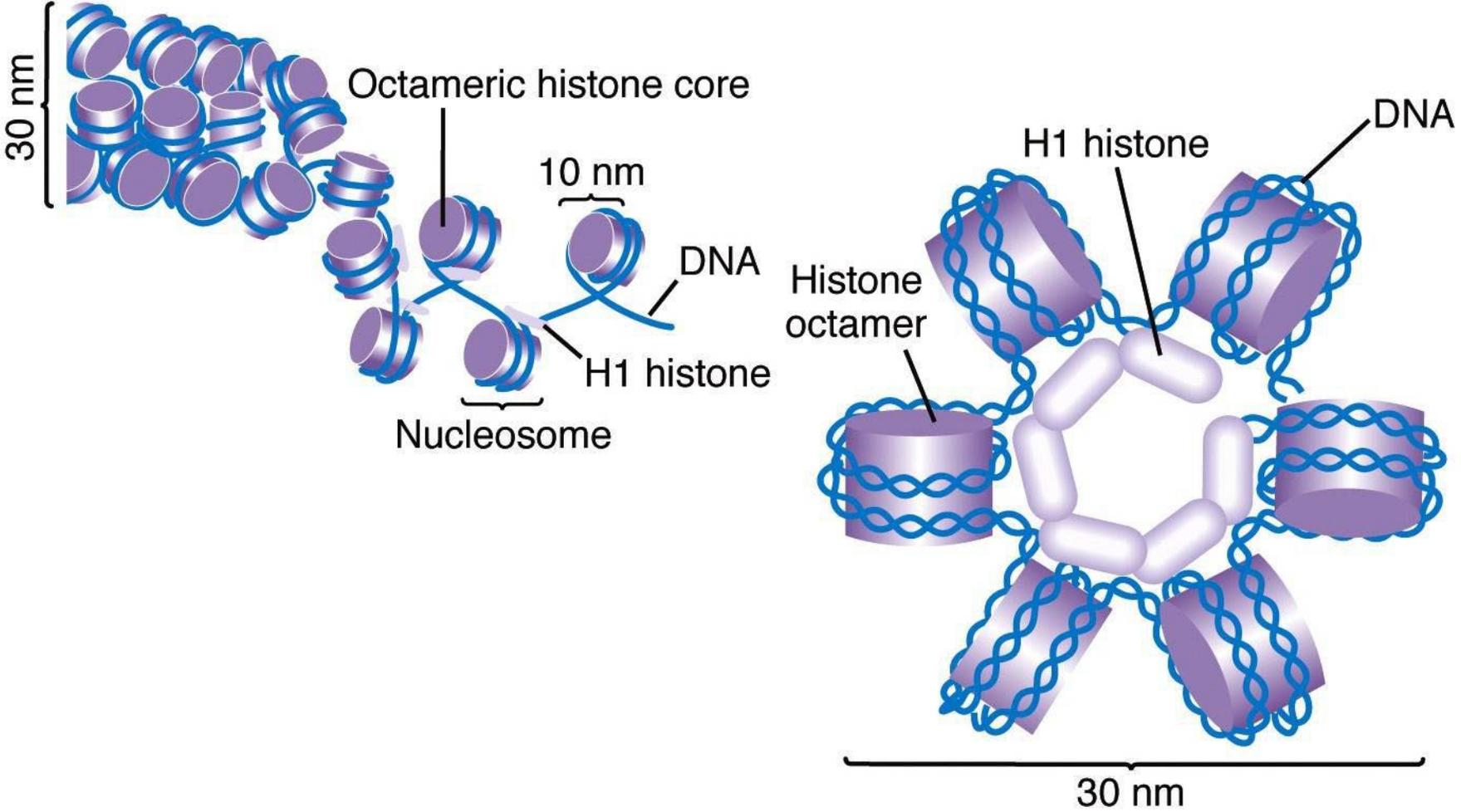
Database of binding sites

The screenshot shows the JASPAR 2026 website. The header includes the JASPAR 2026 logo, a search bar with the text "Search JASPAR database...", and navigation links for "Home", "About", "Search", "Browse JASPAR CORE", "Unvalidated Profiles", "Deep Learning Collection", "Browse Collections", "Tools", "RESTful API", "Download Data", "Matrix Clusters", "Genome Tracks", "Enrichment Analysis", "TFBS extraction", "isMOTIFin", and "LLM-extracted TF Targets". The main content area features a search bar, a section for "Browse JASPAR CORE for 6 different taxonomic groups" with images for Fungi, Insecta, Nematoda, Plantae, Urochordata, and Vertebrata, and a large banner titled "The high-quality transcription factor binding profile database" with a "JASPAR interactive tour" button. Below the banner is a "Citing JASPAR 2026" section with a citation and a "Download" button. At the bottom, there are buttons for "Q&A Forum", "RESTful API", and "pyJASPAR".

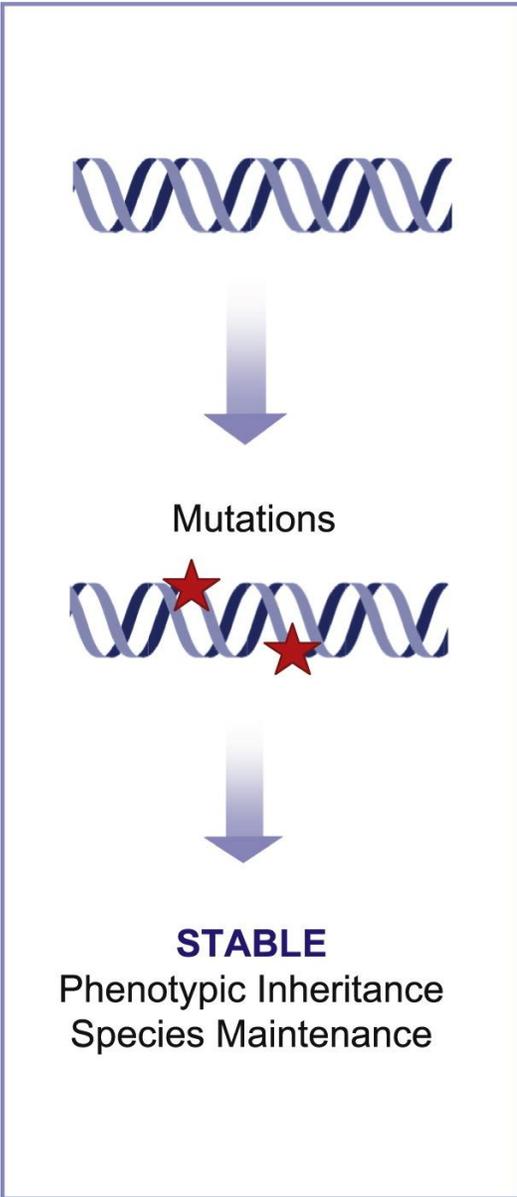
These binding site are often degenerated (receptor ligand analogy)



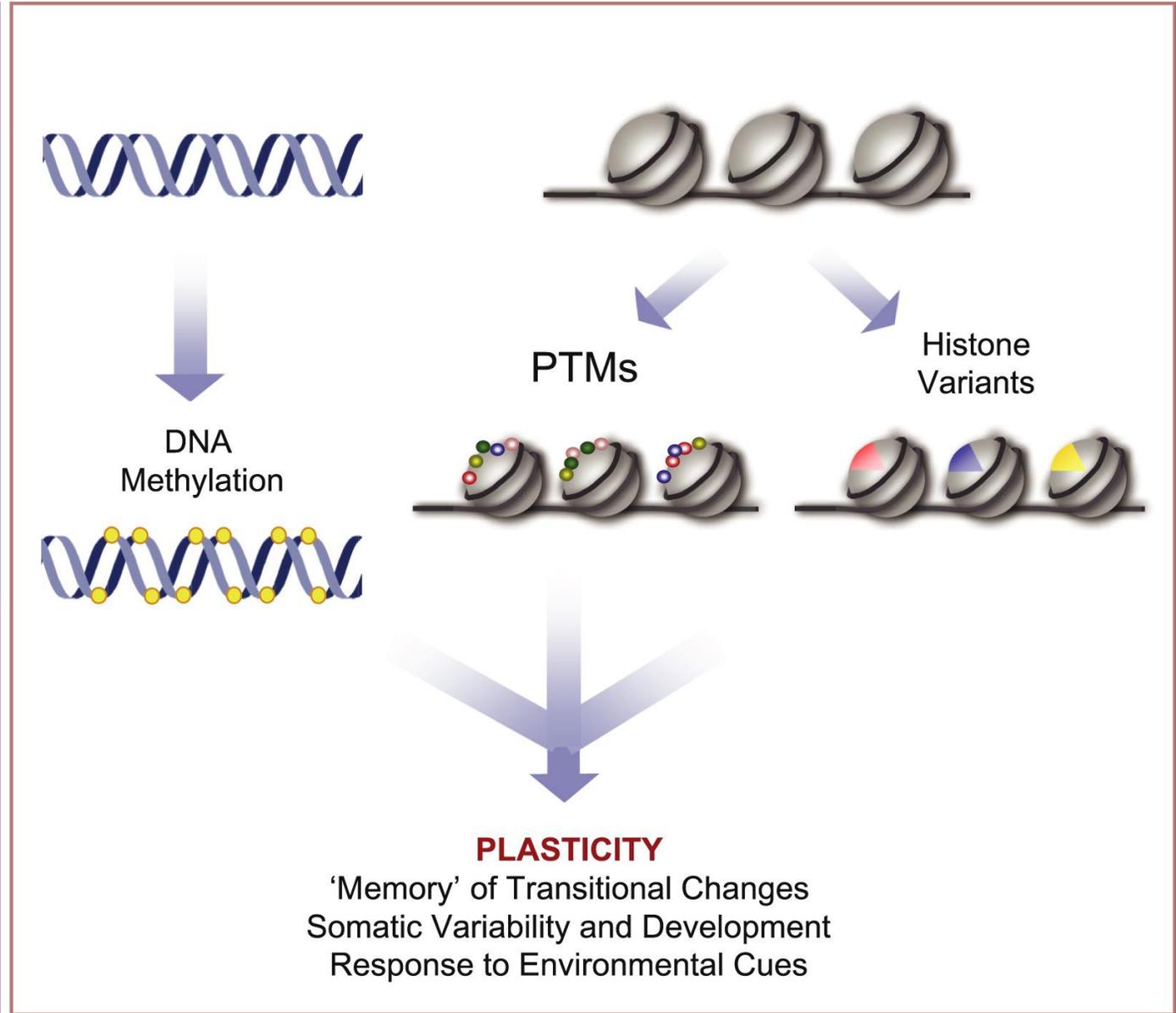
Chromosomes Eukaryotes



GENETICS

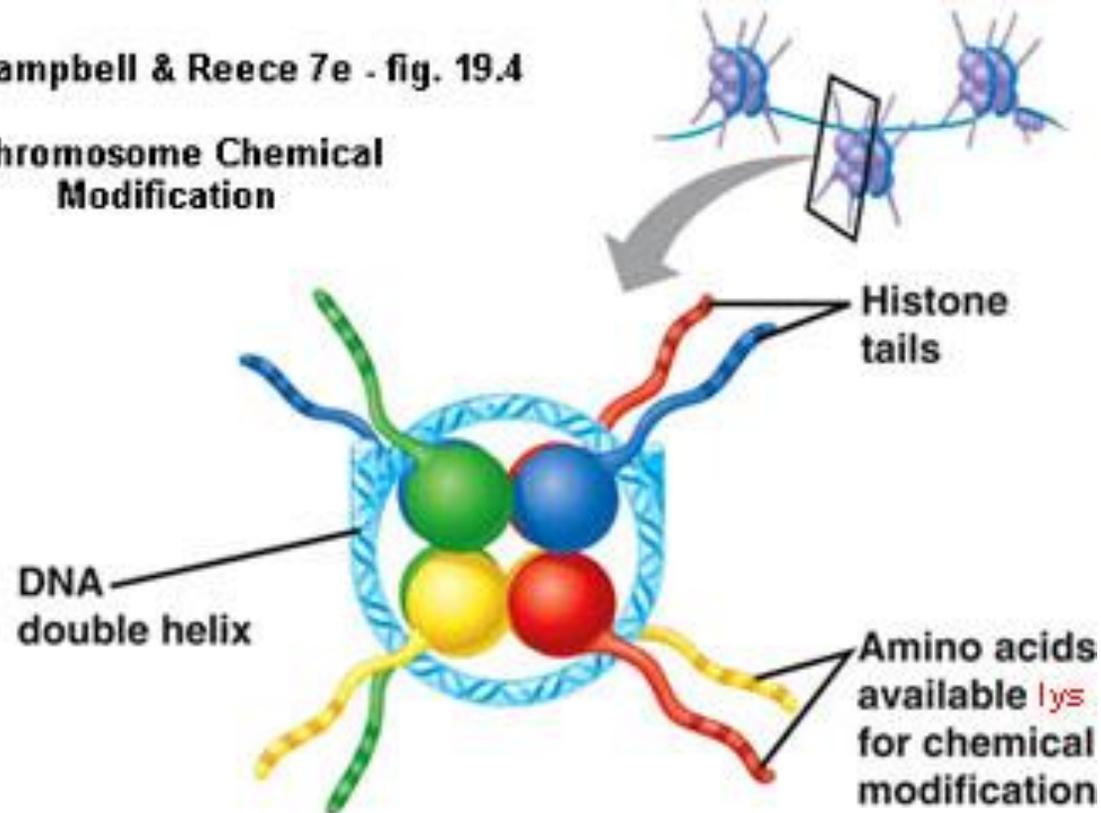


EPIGENETICS

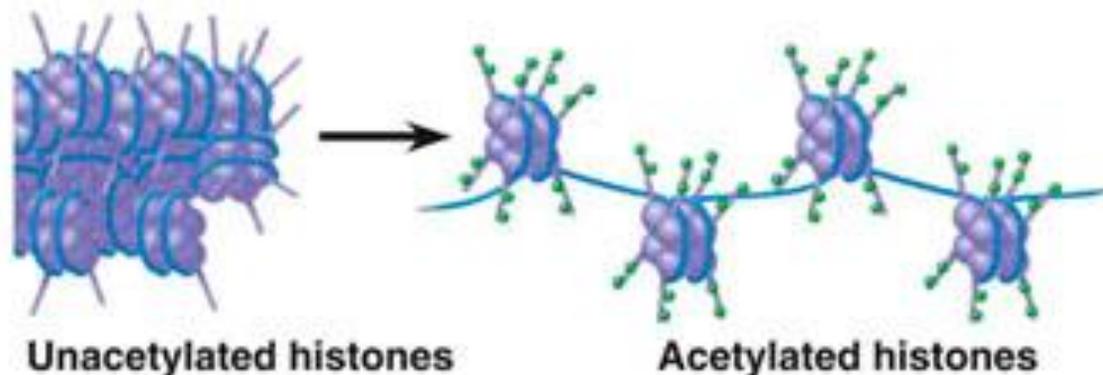


Campbell & Reece 7e - fig. 19.4

**Chromosome Chemical
Modification**



(a) Histone tails protrude outward from a nucleosome



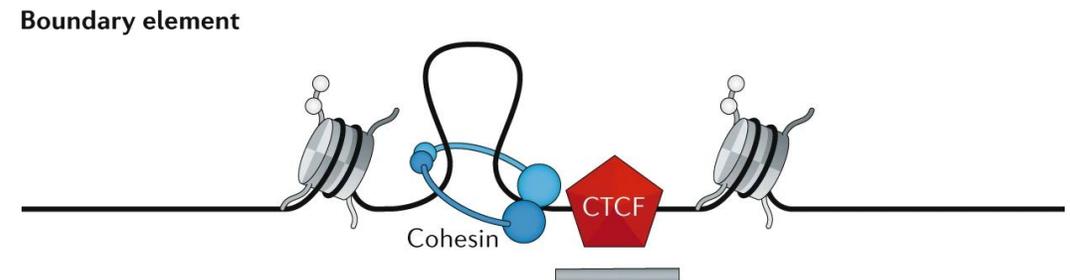
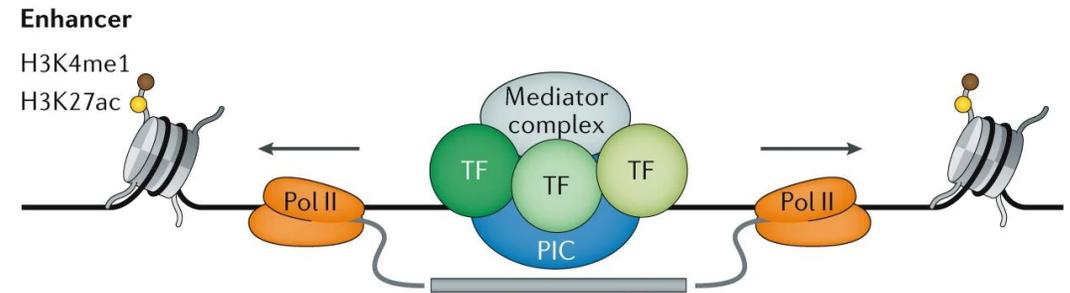
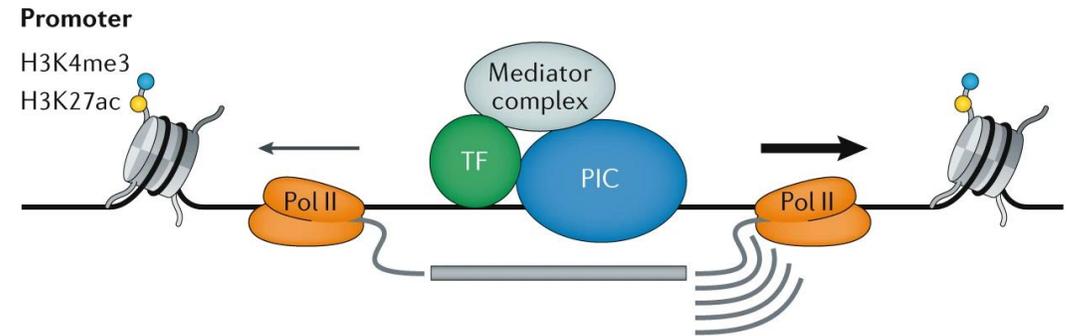
(b) Acetylation of histone tails promotes loose chromatin structure that permits transcription

Boundary elements

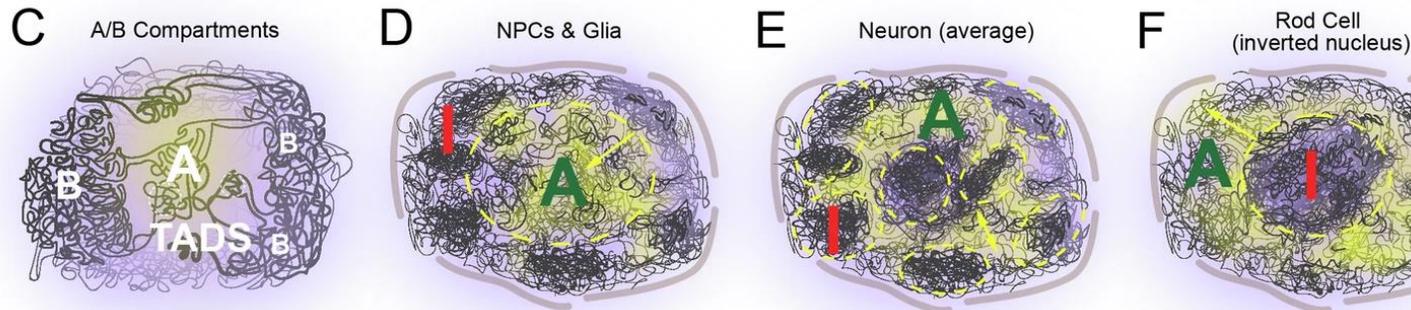
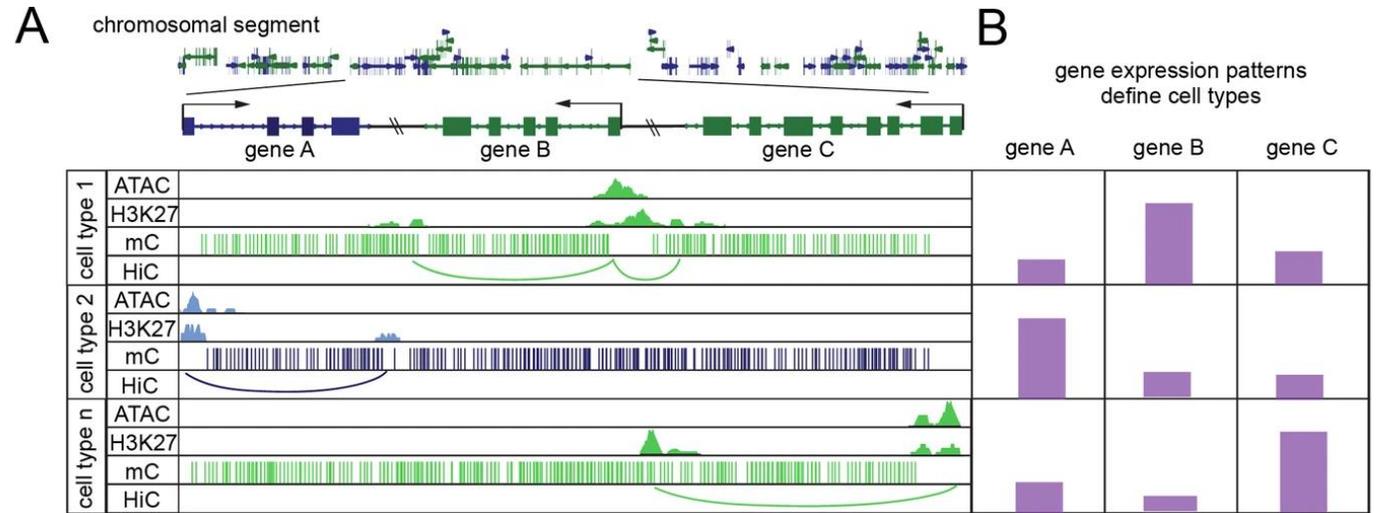
Sequences that isolate chromatin domains so that effects in one domain does not propagate

Overview of the three fundamental classes of regulatory elements in the genome: promoters, enhancers and boundary elements. The elements are indicated with a grey bar, and their bound factors and chromatin modifications when in an active state are shown. The transcription factors (TFs) in dark green represent ubiquitous TFs, whereas the TFs in light green represent tissue-specific and developmental stage-specific TFs. CTCF, CCCTC-binding factor*; H3K4me1, histone H3 monomethylated at Lys4; H3K4me3, histone H3 trimethylated at Lys4; H3K27ac, histone H3 acetylated at Lys27; PIC, pre-initiation complex; Pol II, RNA polymerase II. (Oudelaar & Higgs, Nature Rev Gen. 2020)

*CTCF is a DNA binding protein binding CCCTC repeats



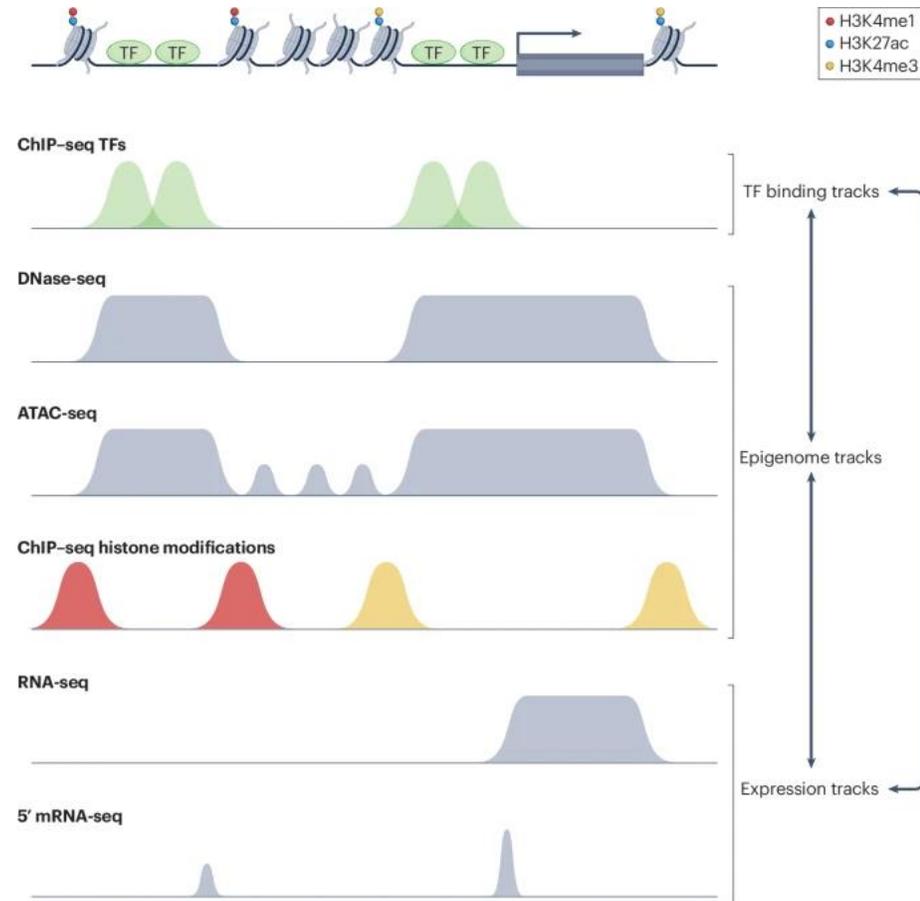
Differential regulation of gene expression defines cell type



A, B A combination of 4DN-type techniques is comparable in resolution with transcriptomics to define cell types. **A** Genome browser tracks for 3 cell types depicting ATAC, H3K27ac, DNA methylation, and HiC loops. **B** A cell type gene expression matrix for 3 cell types and 3 genes highlighting how cell types are defined by 4DN characteristics. **C–F** Neural precursors and glia (and most cells in the body) predominantly show inactive (I, red) chromatin localized at the nuclear exterior while active (A, green) chromatin is located centrally. Neurons do not seem to follow this rule as strictly, with the extreme example of nocturnal photoreceptors being complexly inverted. (Logeman et al., 2025 Mol Psy).

What data can be used to train deep learning models to predict gene transcription ?

Binding of transcription factors (TFs) influences gene expression and can be measured with chromatin immunoprecipitation followed by sequencing (ChIP-seq)¹²⁹. Regulatory elements often reside in open chromatin, which can be measured with DNase I hypersensitive sites sequencing (DNase-seq) or assay for transposase-accessible chromatin with sequencing (ATAC-seq)¹³⁰. Histone modifications that correlate with gene expression patterns can also be measured with ChIP-seq¹³¹. RNA sequencing (RNA-seq) or 5' mRNA sequencing (5' mRNA-seq) provides measurements of gene expression. 5' mRNA-seq techniques include cap analysis of gene expression¹³², global run-on with cap capture¹³³ or transcription start site sequencing¹³⁴, which measure the 5'-capped ends of nascent RNAs (reviewed elsewhere¹³⁵). Unlike RNA-seq, these techniques do not measure transcription termination, transcript stability or splicing events. All of these different data types can be used to train sequence-to-expression models. There are extensive mutual causality links between these features (indicated by arrows): DNA sequence and chromatin accessibility dictate TF binding; some TFs can open up chromatin or recruit histone-modifying enzymes, and transcription is controlled by all of these features but conversely also affects chromatin states.



AlphaGenome Model

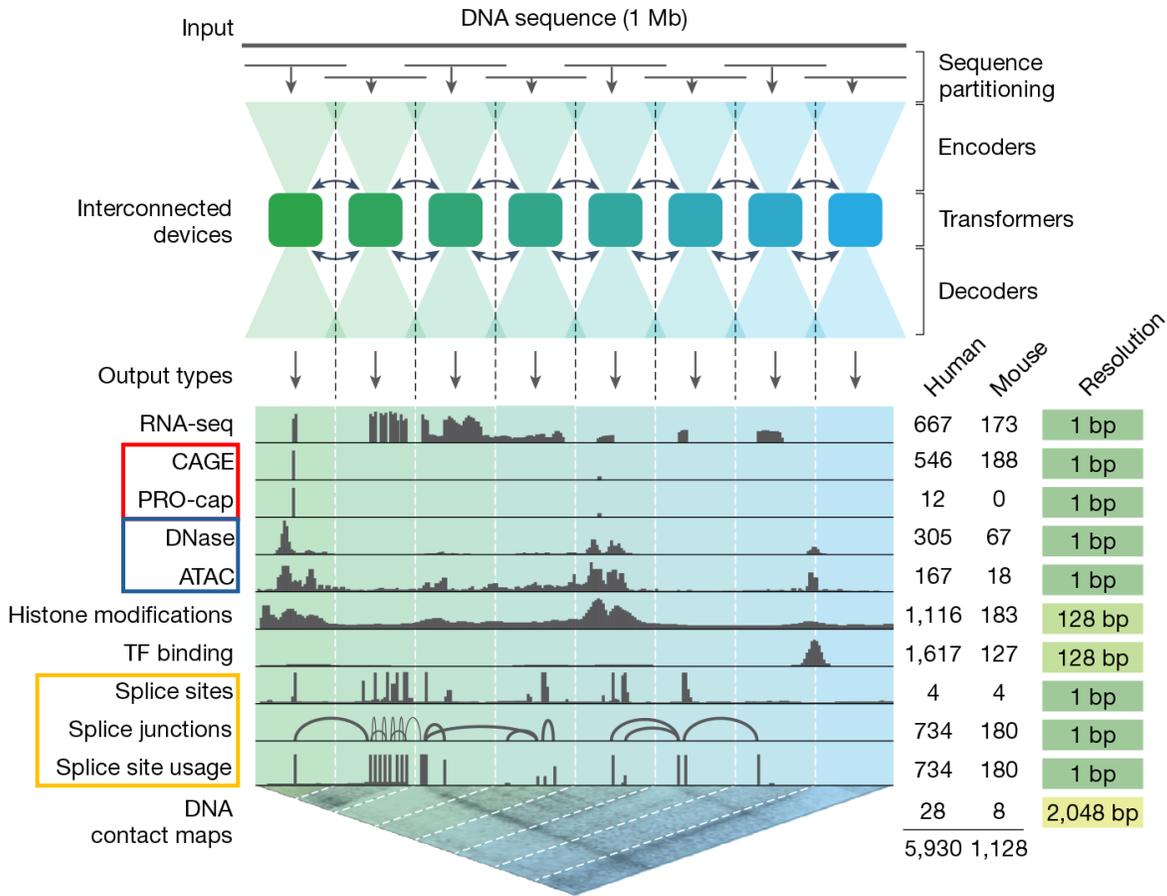
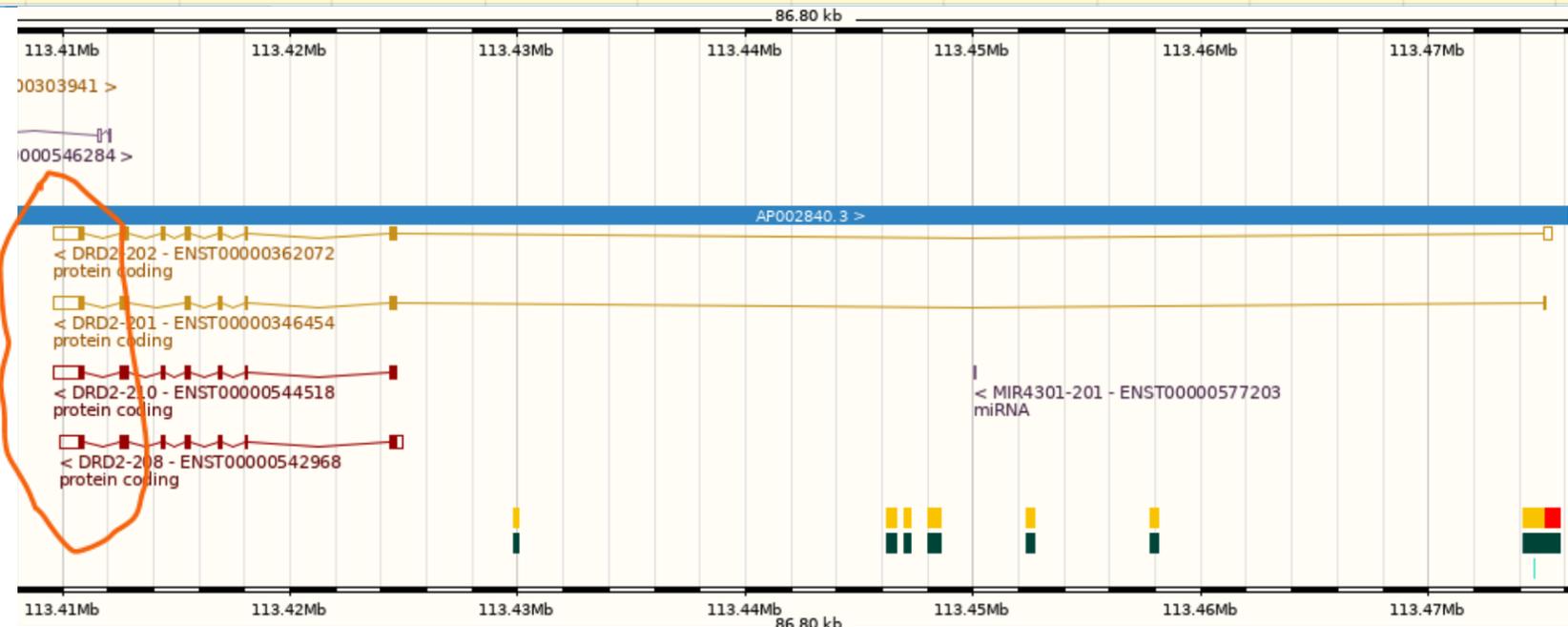
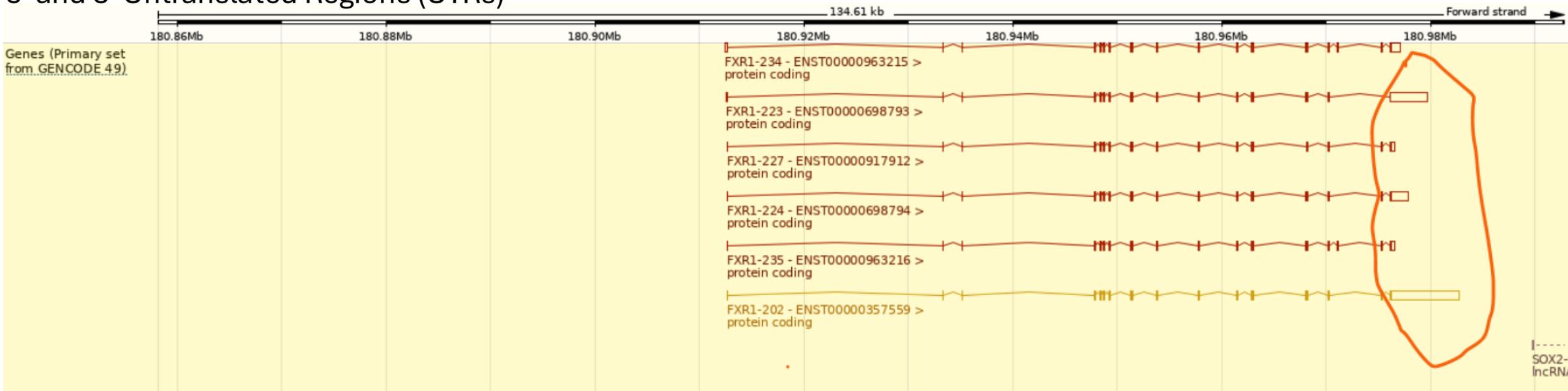


Fig. 1 | AlphaGenome model architecture, training regimes and comprehensive evaluation performance. a, Model overview. AlphaGenome processes 1 Mb of DNA sequences and species identity (human/mouse) to predict 5,930 human or 1,128 mouse genome tracks across diverse cell types and 11 output types at specific resolutions (far right). Computation leverages sequence parallelism, breaking the 1 Mb of DNA sequence into 131-kb chunks processed across devices. The core architecture features a U-Net-style design comprising an encoder (downsampling the sequence), transformers with inter-device communication and a decoder (upsampling), which feed into task-specific output heads at their respective resolutions (detailed in Extended Data Fig. 1). Nature | Vol 649 | 29 January 2026

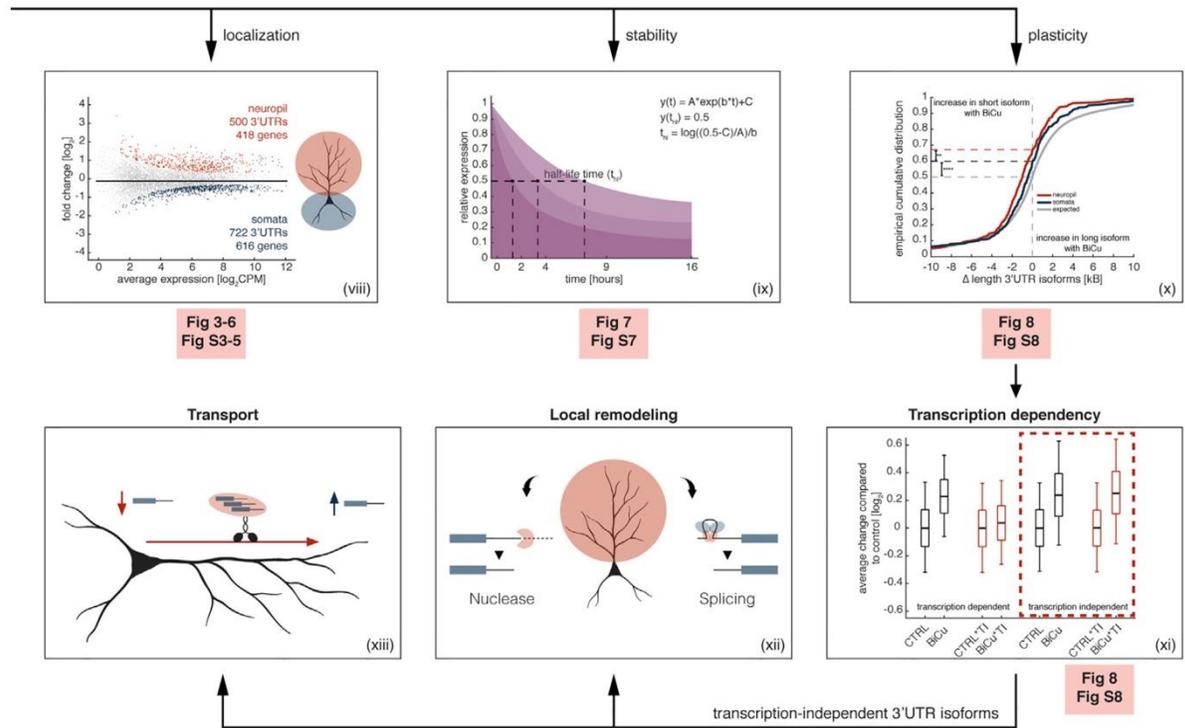
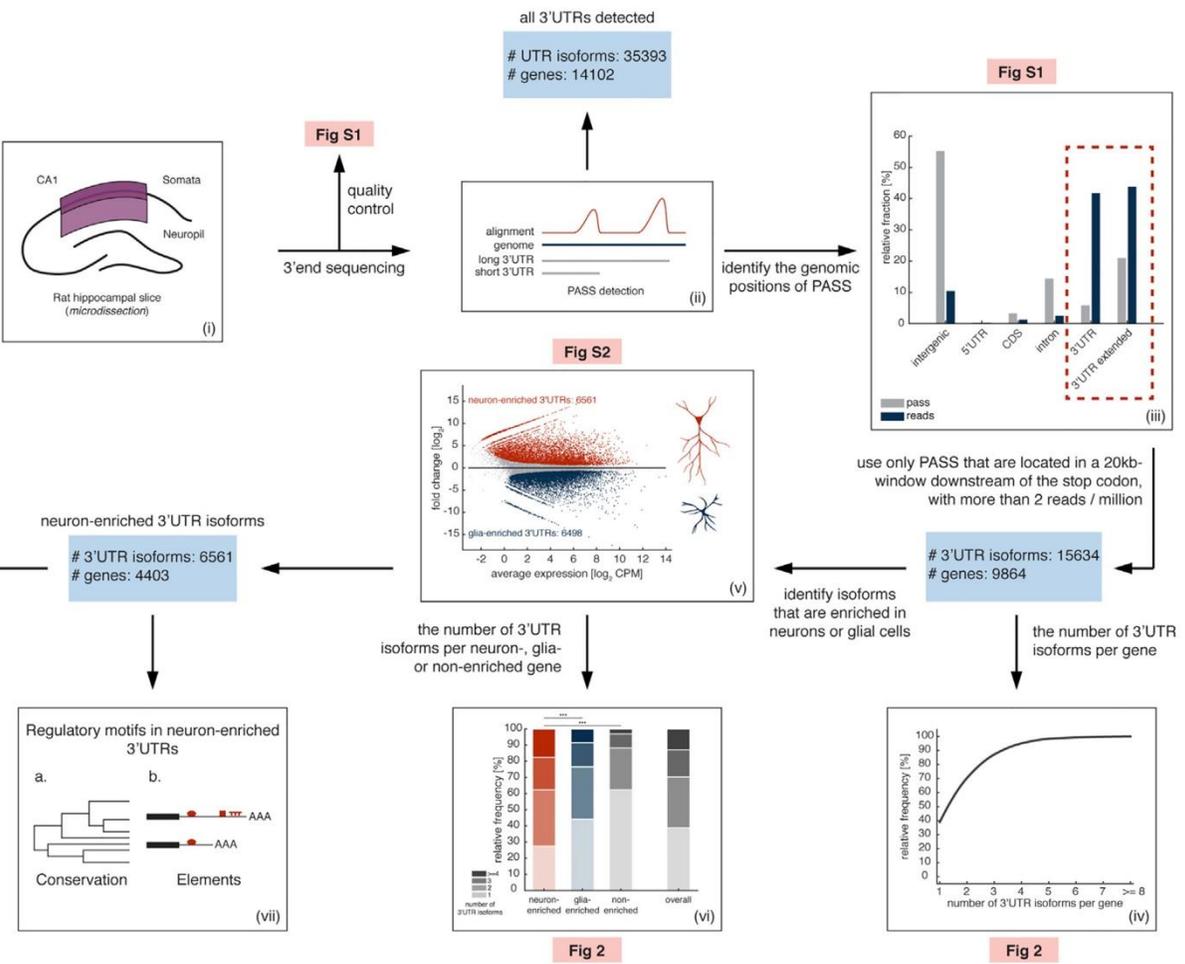
The other languages of RNA

5' and 3' Untranslated Regions (UTRs)

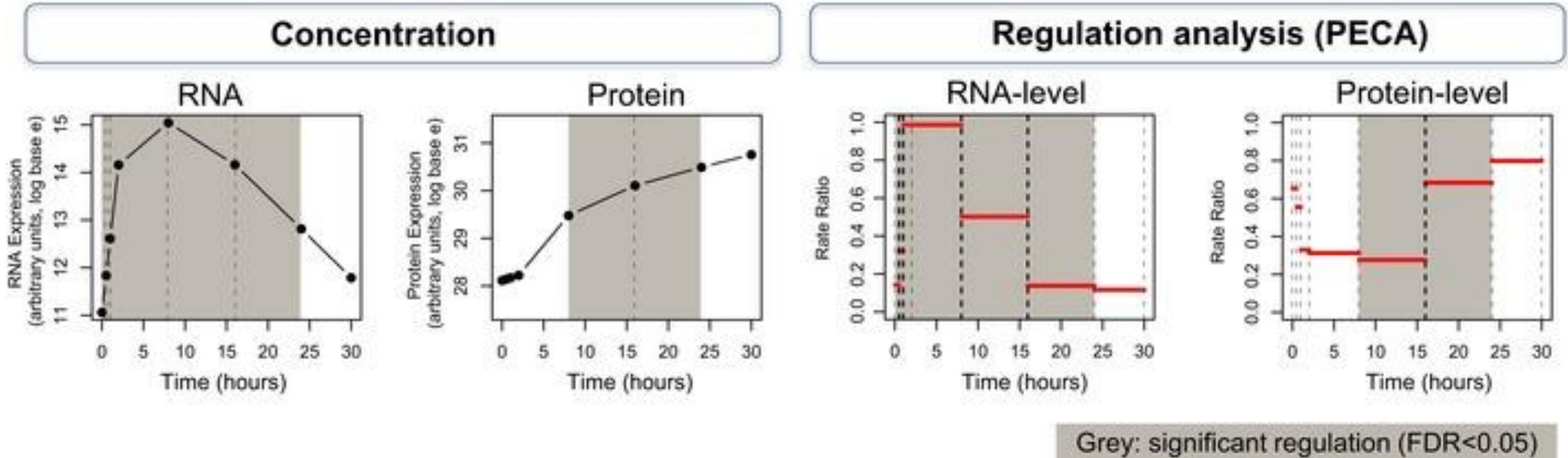


UTRs affect the localization Translation and stability of their respective mRNAs.

Alternative 3'UTRs Modify the Localization, Regulatory Potential, Stability, and Plasticity of mRNAs in Neuronal Compartments



mRNA Levels are not directly related to protein levels



The example shows the chaperone GRP78, a key ER stress protein. mRNA and protein concentrations are shown on the left; PECA results are shown on the right for RNA and protein level, respectively. Intervals with significant regulation as determined by PECA are gray shaded (FDR < 0.05). The value of PECA is illustrated at the 16-h time point at which mRNA concentration decreases, while protein concentration still rises. PECA highlights that there is a significant RNA- and protein-level regulation around this time point—a signal that would otherwise likely have been overlooked.

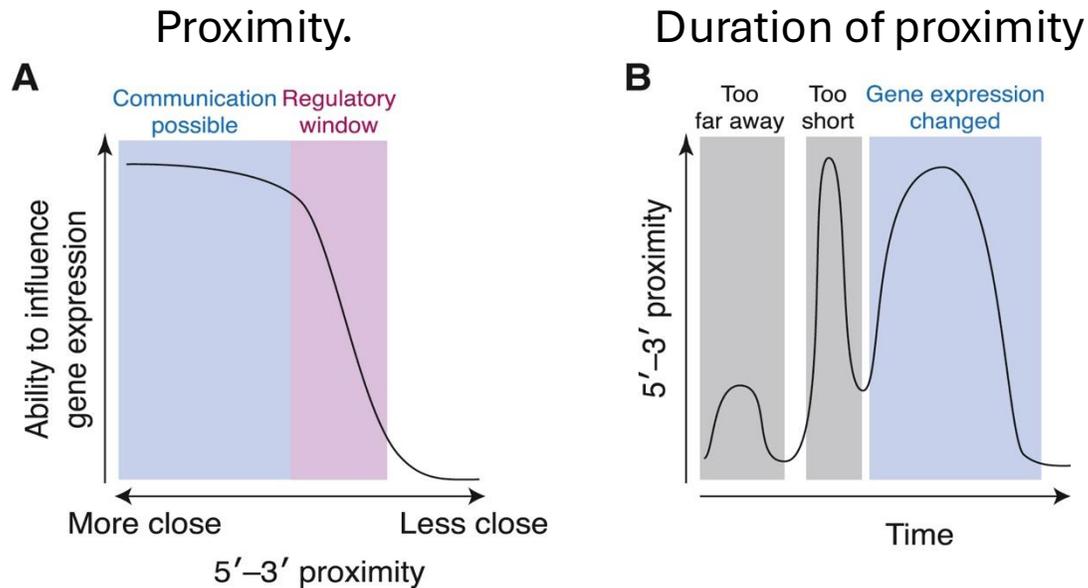
. 2016 Jan 20;12(1):855. doi: [10.15252/msb.20156423](https://doi.org/10.15252/msb.20156423)

For details on Protein Expression Control Analysis (PECA):

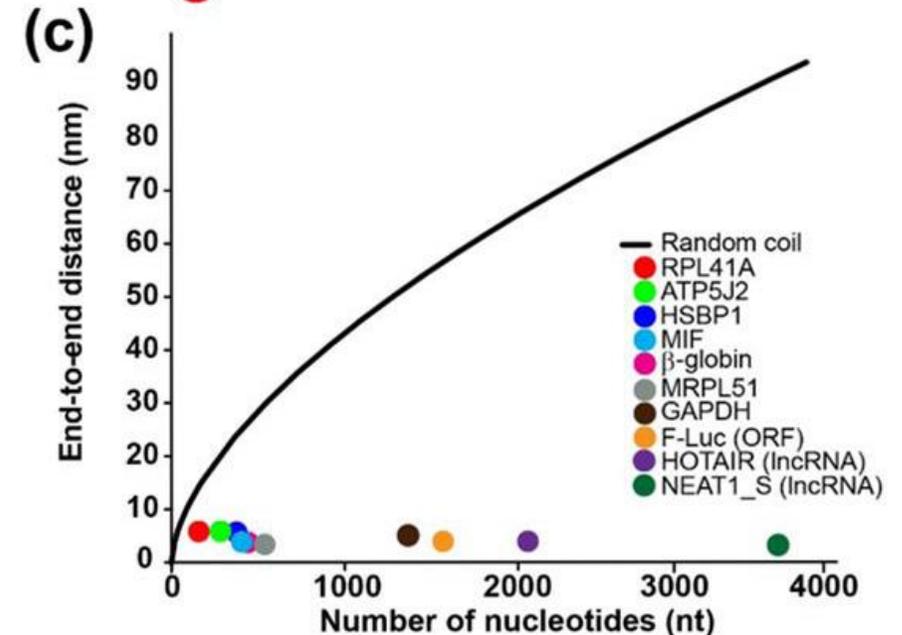
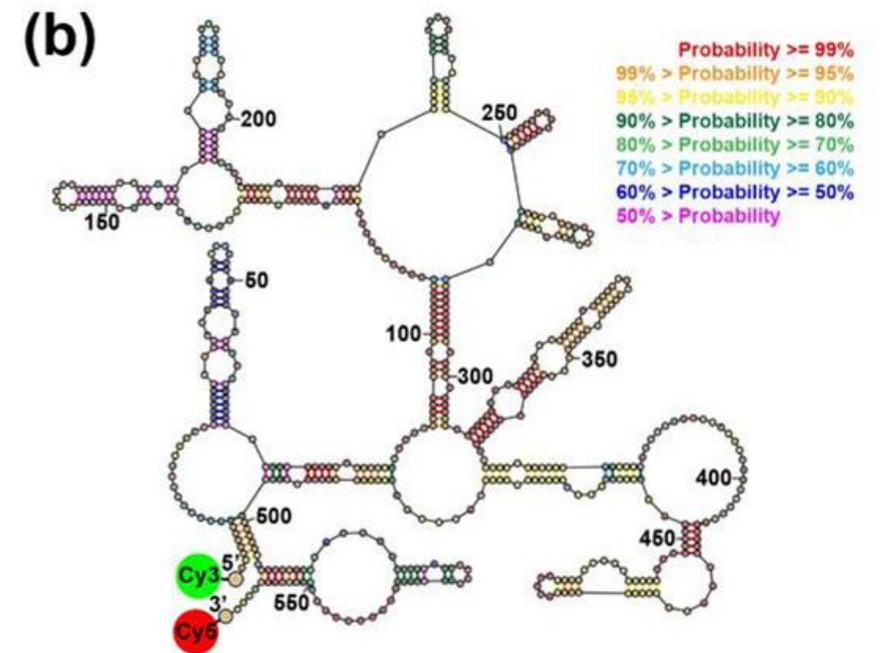
J Proteome Res. 2013 Nov 26;13(1):29–37. doi: [10.1021/pr400855q](https://doi.org/10.1021/pr400855q)

Not restricted to UTRs mRNA form complex structures guided by Nucleotide interactions

The RNA sequence is not only a set of instruction to make proteins
It is also a mode of regulation



<https://doi.org/10.1016/j.molcel.2018.10.047>



<https://www.biorxiv.org/content/10.1101/2020.04.29.069203v1.full>

Three main mode of regulations

- 1- Non-specific regulators of protein synthesis.
- 2- Specific RNA binding proteins (RBPs); regulate stability, translation, transport.
- 3- Small regulating RNAs (e.g. miRNAs.)

Regulation of mRNAs by RBP and miRNA is critically determined by the formation of specific secondary structures in the mRNAs

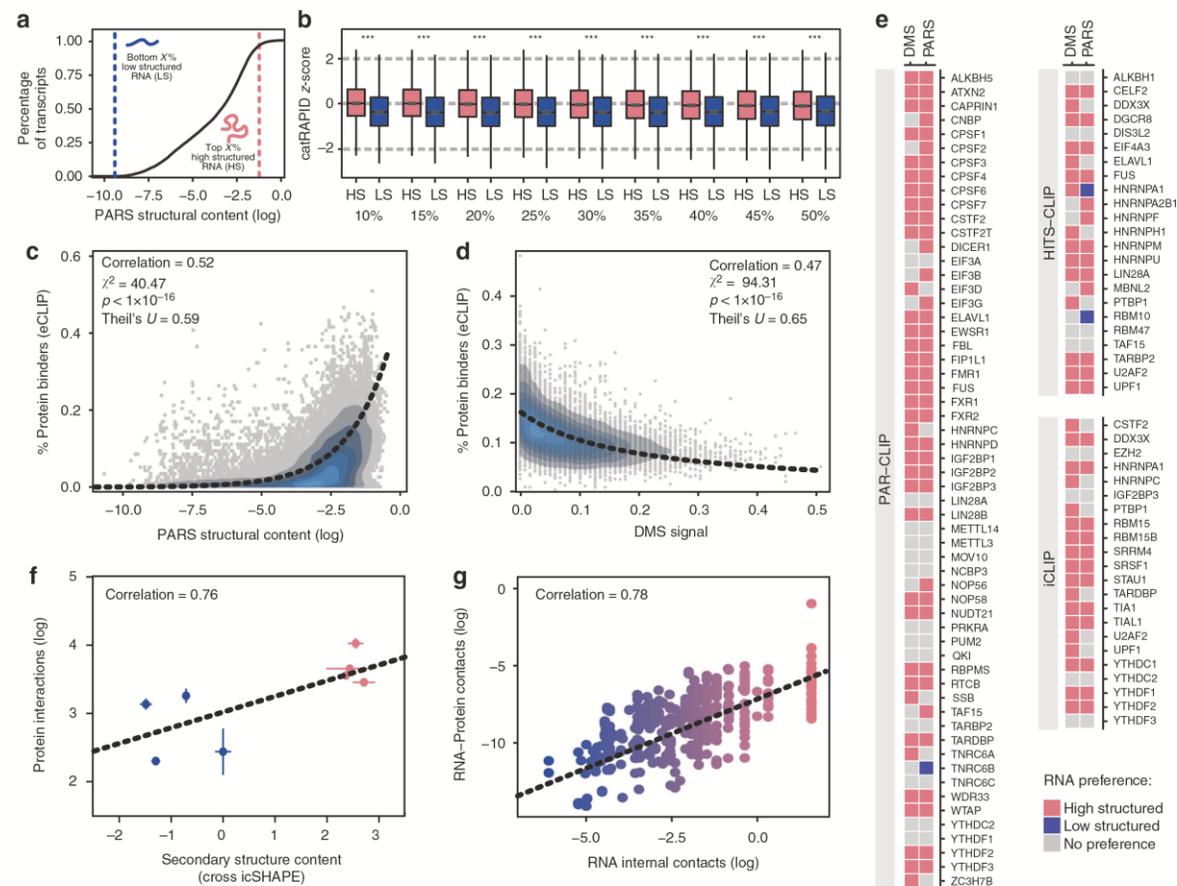
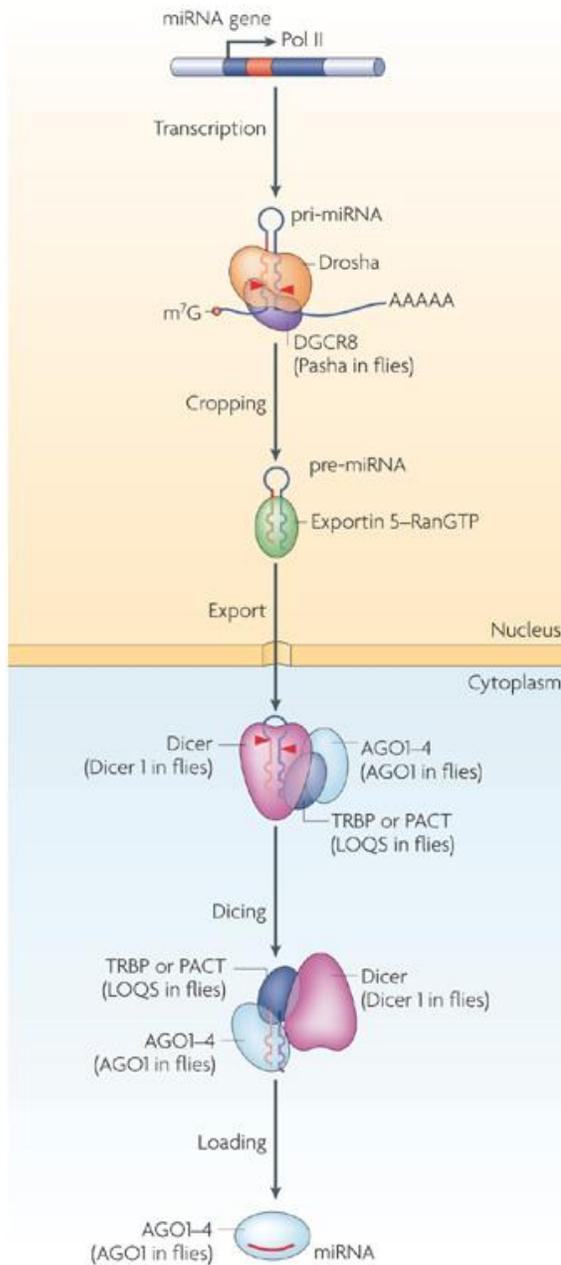


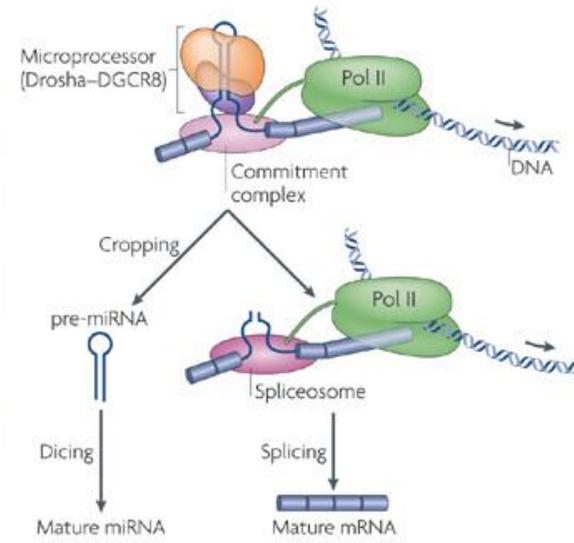
Fig. 1 The amount of protein structure correlates with the number of interactions. **a** Cumulative distribution function (CDF) for the secondary structure content of all human RNAs measured by parallel analysis of RNA structure (PARS)^{8,69}. Vertical lines indicate a certain fraction (X%) of RNAs with the lowest secondary content (LS; blue) and the same fraction with the highest secondary content (HS; pink). **b** catRAPID predictions of protein interactions with human RNAs ranked by structural content measured by PARS (118 RNA-binding proteins (RBPs) for which enhanced crosslinking and immunoprecipitation (eCLIP) information is also available)³¹. The fractions 10%, 15%, ..., 50% refer to the comparison between equal-size HS and LS sets. The results indicate that catRAPID is able to distinguish HS and LS groups significantly and consistently through the different fractions (p value $< 10^{-16}$; Kolmogorov-Smirnov (KS) test). The boxes show the interquartile range (IQR), the central line represents the median, the whiskers add 1.5 times the IQR to the 75 percentile (box upper limit) and subtract 1.5 times the IQR from the 25 percentile (box lower limit). s.d. is shown. **c** Relationship between number of protein interactions (eCLIP) and structural content measured by PARS³⁰. The fitting line corresponds to the formula $y = \exp(\alpha + \beta x)$, where $\alpha = -0.75$; $\beta = 0.67$; p value estimated with KS test. **d** Relationship between number of protein interactions and structural content measured by dimethyl sulfate modification (DMS)⁹. The fitting line corresponds to the formula $y = 1/(\alpha + \beta x)$, where $\alpha = 2.60$; $\beta = 87.36$; p value estimated with KS test. **e** Structural preferences of RBPs measured with three different CLIP techniques (photoactivatable ribonucleoside-enhanced CLIP (PAR-CLIP), high-throughput sequencing-CLIP (HITS-CLIP) and individual nucleotide resolution CLIP (iCLIP)). The colour indicates the RNA-binding preference of each protein: pink, high structured; blue, low structured; grey, no preference. **f** Correlation between structural content (CROSS predictions of icSHAPE experiments) and protein interactions of eight transcripts revealed by protein microarrays (Pearson's correlation). s.d. is shown. **g** Analysis of Protein Data Bank (PDB) structures containing protein-RNA complexes reveals a trend between protein (inter) and RNA (intra) contacts (196 different pairs; Pearson's correlation)

Micro RNAs

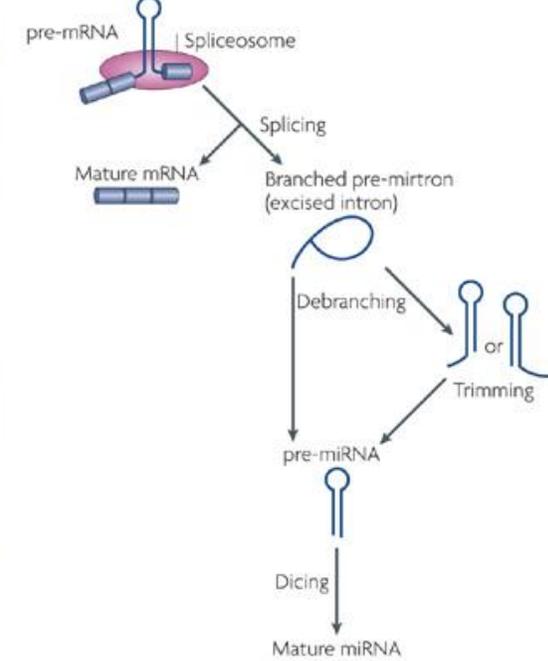
a Biogenesis of canonical miRNA



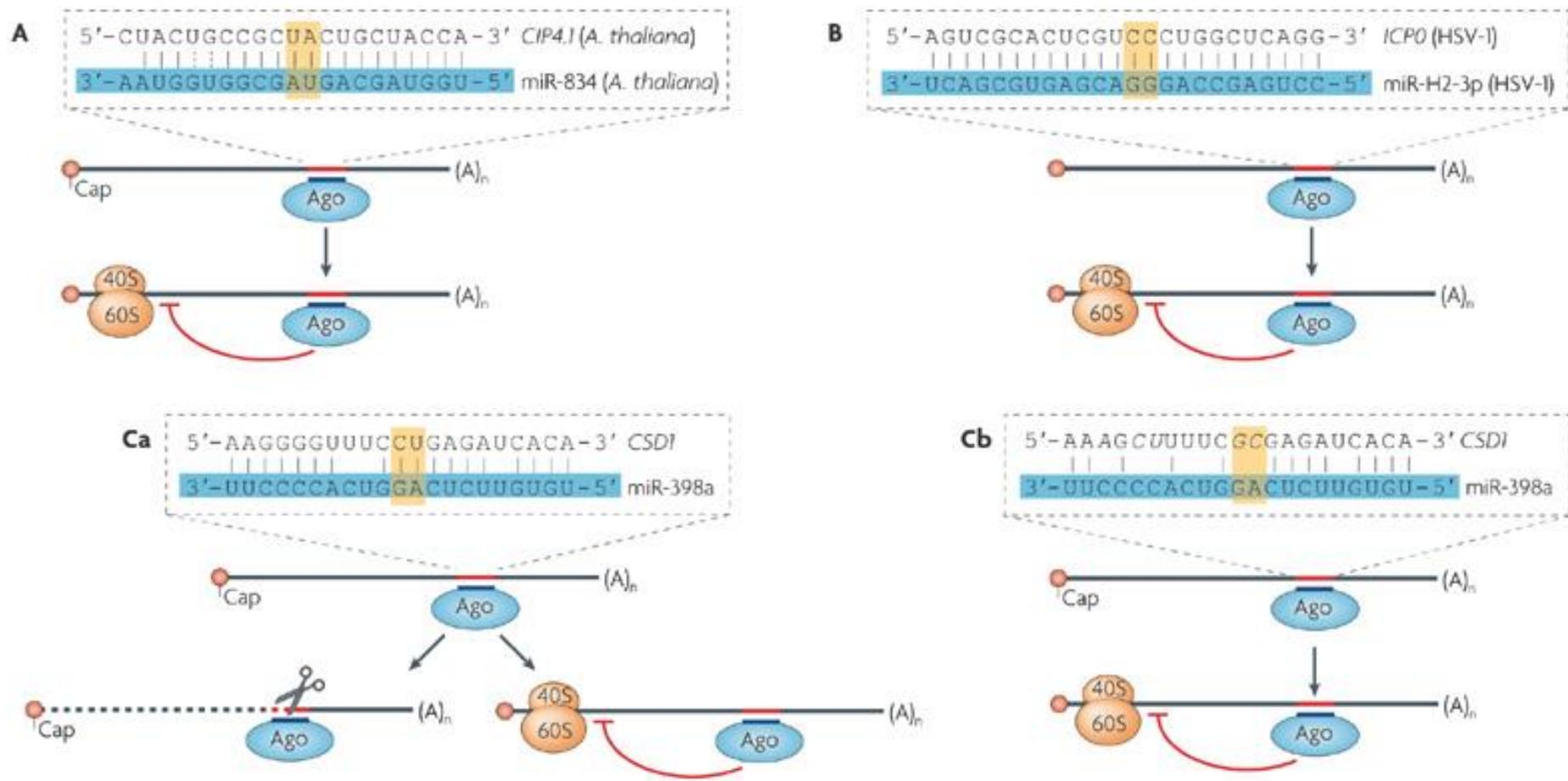
b Canonical intronic miRNA



c Non-canonical intronic small RNA (mirtron)



miR Generally stall the translation of mRNA or promote their degradation. However, there are evidence for other roles



The regulation of mRNA translation and stability is an additional layer of regulation.

-It can allow for subcellular differences in gene expression.

-It can allow to prime a system by stocking RNA on stalled ribosomes for fast synthesis of proteins upon Changes of biological needs.

RNA levels are a function of:

-Gene transcription

-RNA stability

Presence of mRNA does not mean
the corresponding protein is made

Protein levels are a function of

-mRNA translation

-protein degradation

Presence of protein does not mean
It is active.

Neither mRNA level nor protein level is a sure predictor of functional response.

They are just indications that the system response status has changed.

Determination of a functional state requires a combination of datasets.

However, single datasets are a good shortcut sometimes.

It all depends on what the question is.

How do you think this can affect ML model predictions?

Are non-ML models possible?