

CS273B: Deep learning for Genomics and Biomedicine

Lecture 2: Genomics 101

09/27/2017

Anshul Kundaje, James Zou

Outline

1

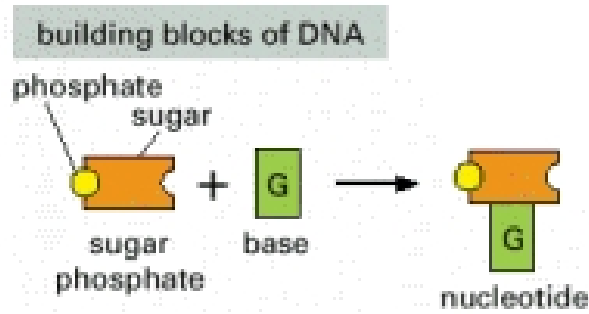
Anatomy of the human
genome

2

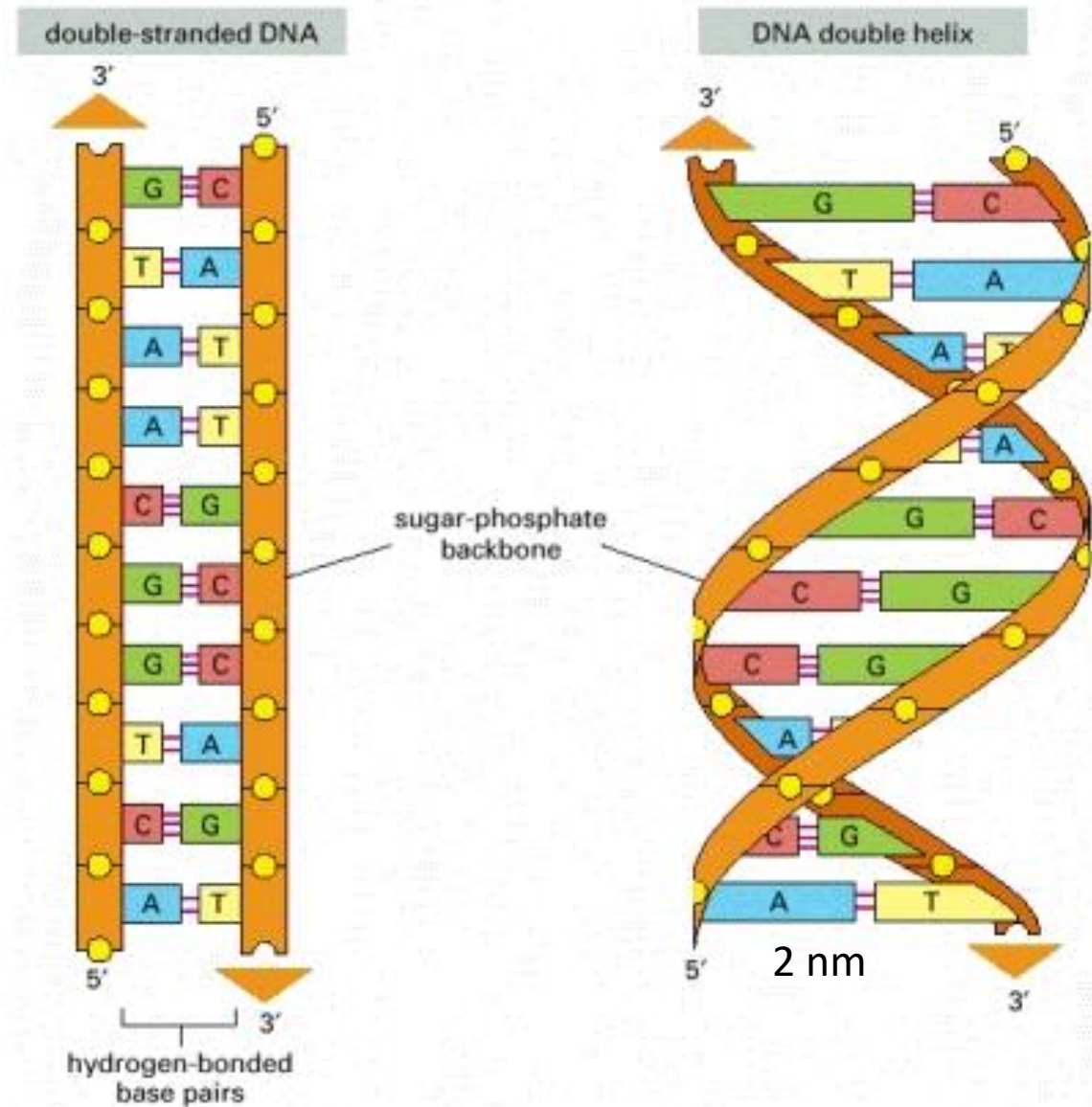
High throughput sequencing

Anatomy of the human genome

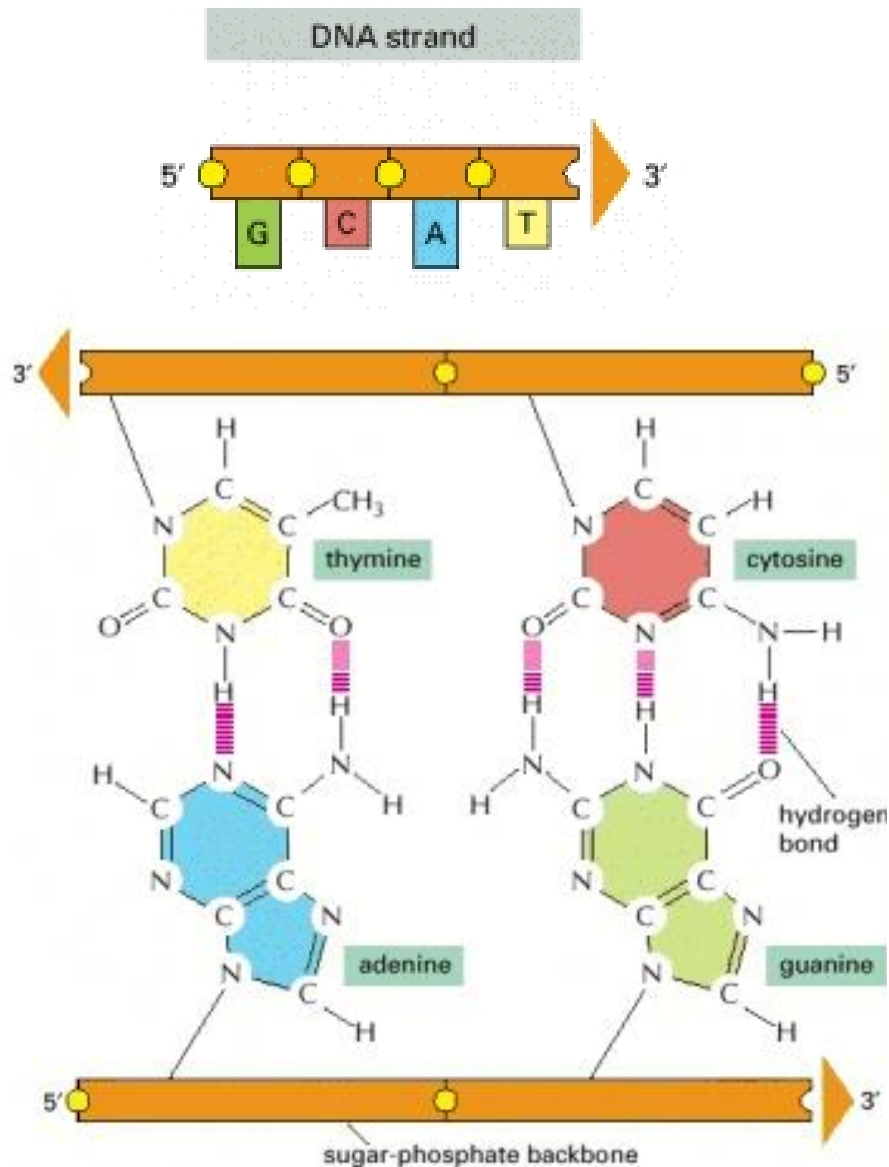
DNA: the molecule of heredity



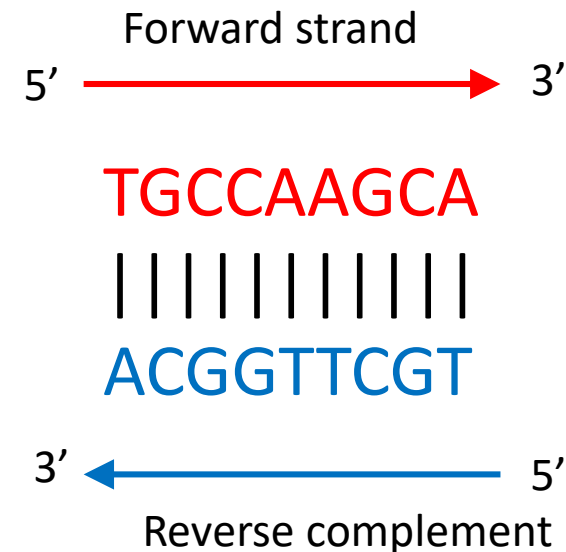
- DNA: Deoxyribose nucleic acid
- Double stranded biopolymer
- Canonical form: right handed double helix
- Phosphate backbone outside
- 4 Bases / nucleotides hidden on the inside (A, C, G, T)
- A pairs with T
- C pairs with G
- 1 helical turn is 10.4 nucleotides



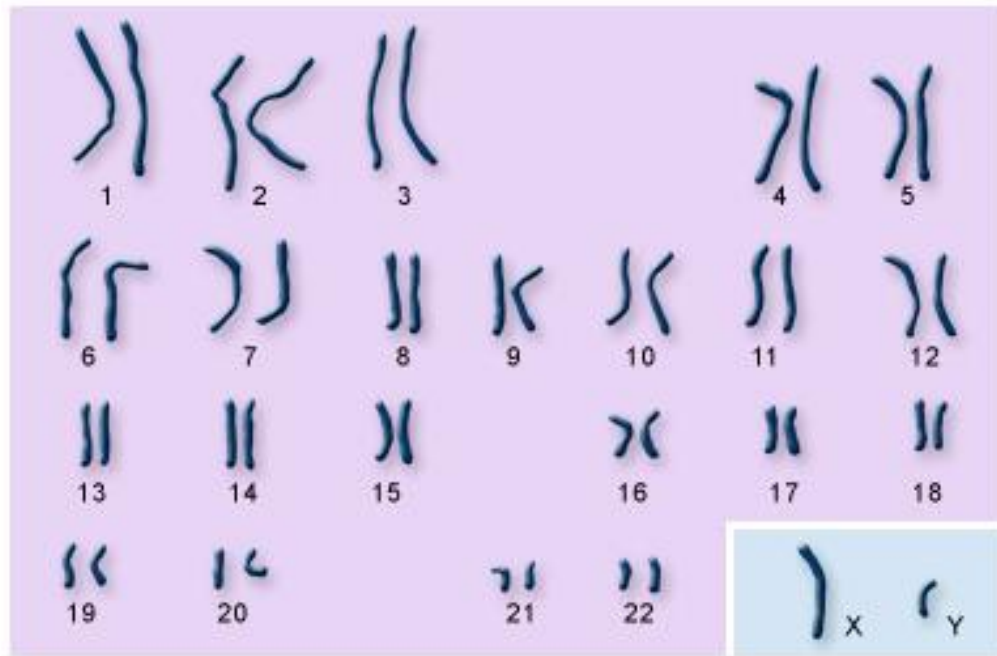
DNA is directional and has two complementary strands



- Weak hydrogen bonds hold the two strands together
- This allows low-energy opening and re-closing of two strands
- **Chemical Polarity:** Extension 5' → 3' tri-phosphate coming from newly added nucleotide



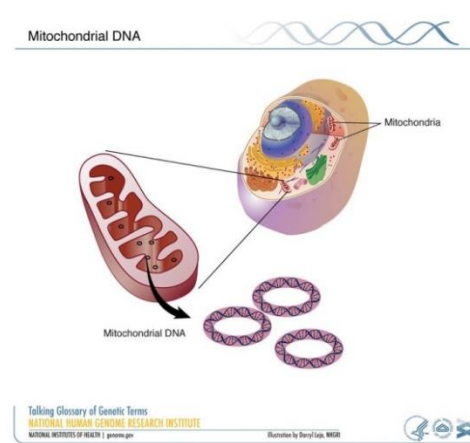
Chromosomes in humans



autosomes

sex chromosomes

U.S. National Library of Medicine



TGCCAAGCA

|||||
ACGGTTCGT

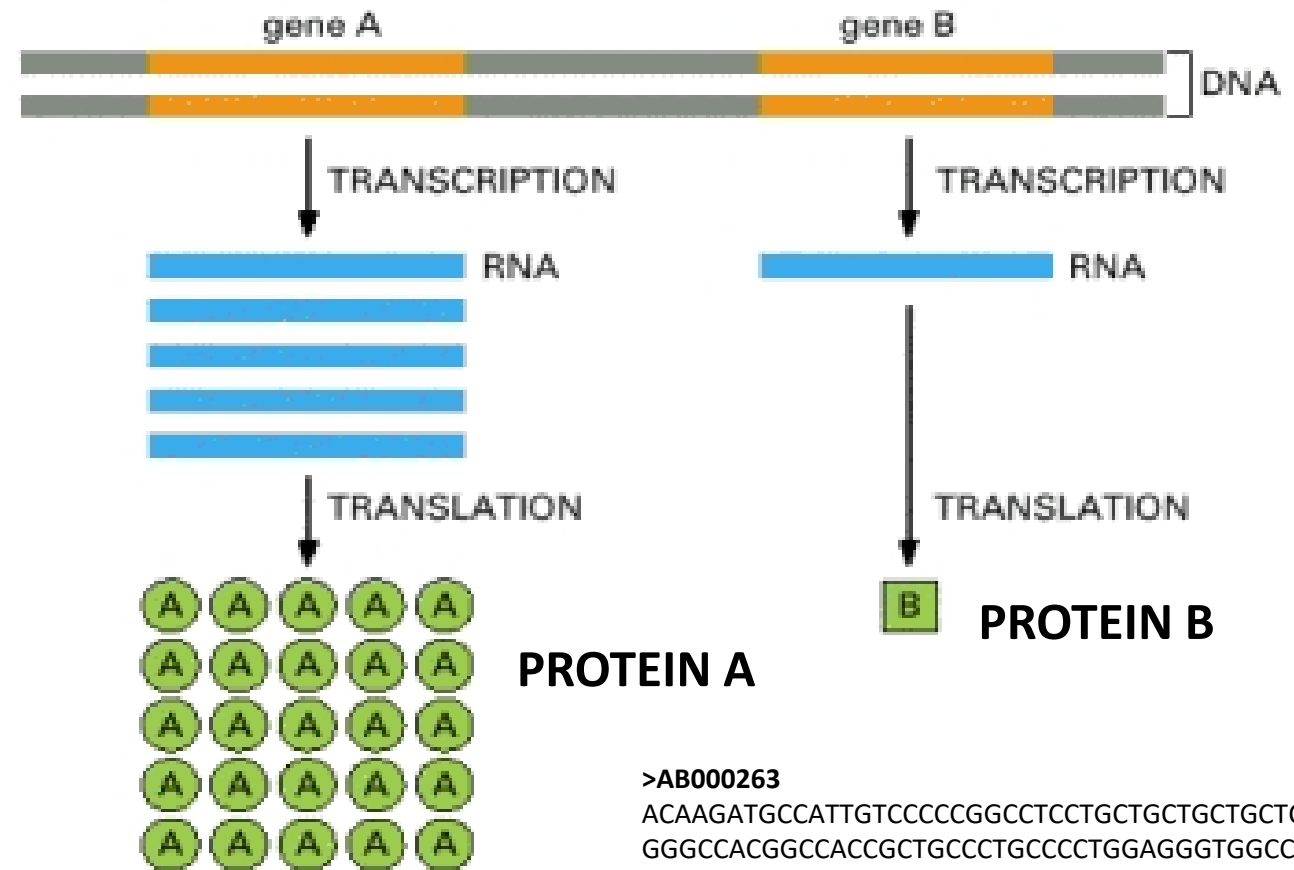
TGCCAAGCA

|||||
ACGGTTCGT

- Humans are diploid (2 copies of each chromosome)
- 22 pairs of autosomes
- Sex chromosomes: female (X,X) , male (X,Y)
- Mitochondrial DNA (circular, many copies per cell)
- Diploid Human genome = ~3 billion **base-pairs** X 2

Genes (DNA -> RNA -> protein)

- **The gene** is a fundamental functional unit of DNA whose sequence codes specific recipes to make other biomolecules (RNA and proteins)
 - In humans each DNA molecule \Leftrightarrow 25,000 protein-coding genes
 - Many other gene-like units encode short and long non-coding RNAs



Where to get genome and gene sequences

GENBANK: <https://www.ncbi.nlm.nih.gov/genbank/>

ENSEMBL: <https://www.ensembl.org/index.html>

```
>AB000263
ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTGCTCTCCG
GGGCCACGGCCACCGCTGCCCTGCCCCTGGAGGGTGGCCCCACC
GGCCGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAAGG
AAAAGCAGCCTCCTGACTTTCTCGTTGGTGGTTTGAGTGGACC
TCCCAGGCCAGTGCCGGGCCCCTCATAGGAGAGGAAGCTCGGG
AGGTGGCCAGGCGGCAGGAAGGCGCACCCCCCAGCAATCCGC
GCGCCGGGACAGAATGCCCTGCAGGAACCTTCTTCTGGAAGACCT
TCTCCTCTGCAAATAAACCTCACCCATGAATGCTCACGCAAGT
TTAATTACAGACCTGAA
```

FASTA format

mRNA: The messenger

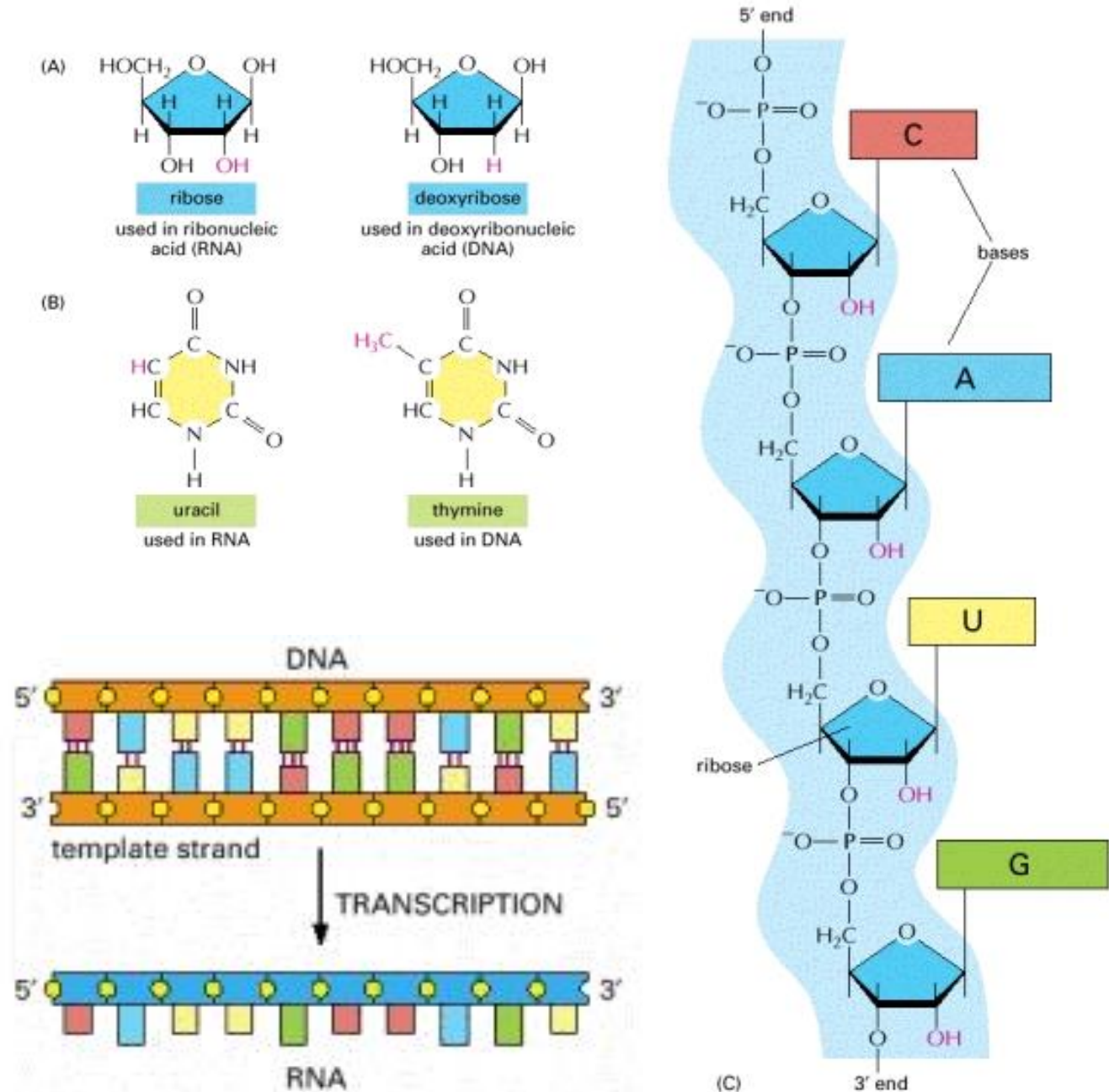
- Information changes medium
 - single strand vs. double strand
 - ribose vs. deoxyribose sugar
 - Uracil (U) instead of Thymine (T)

DNA sequence

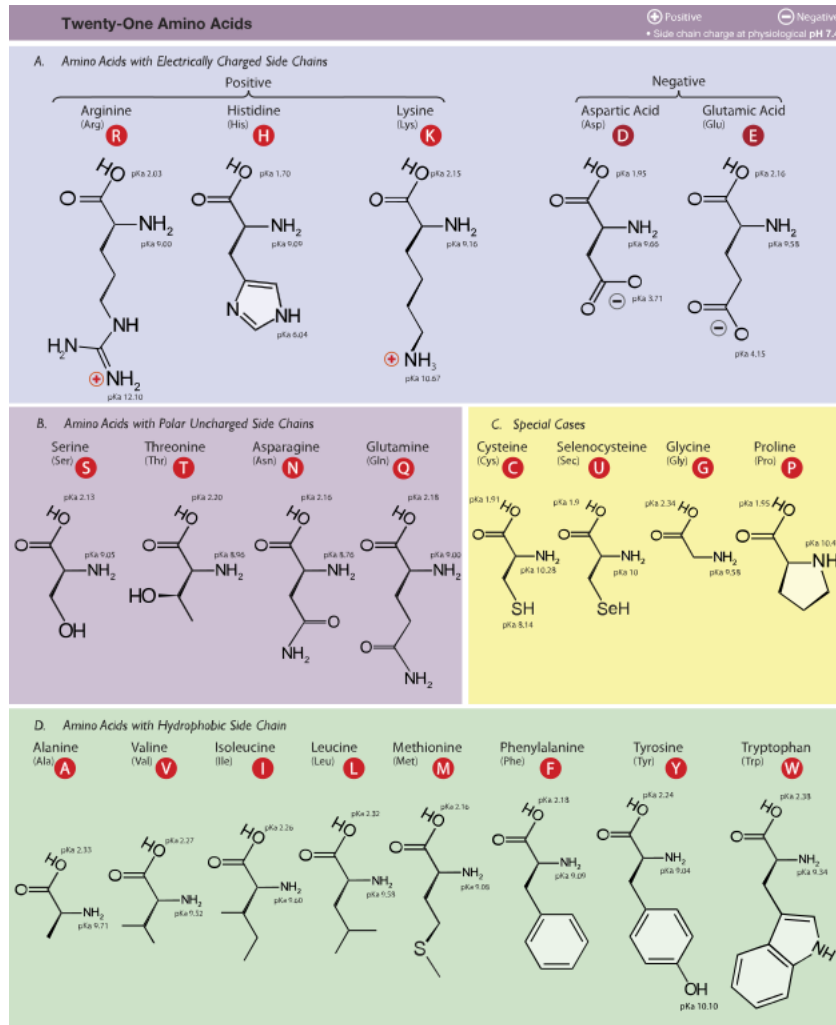
A T T A C G G T A C C G T

U A A U G C C A U G G C A

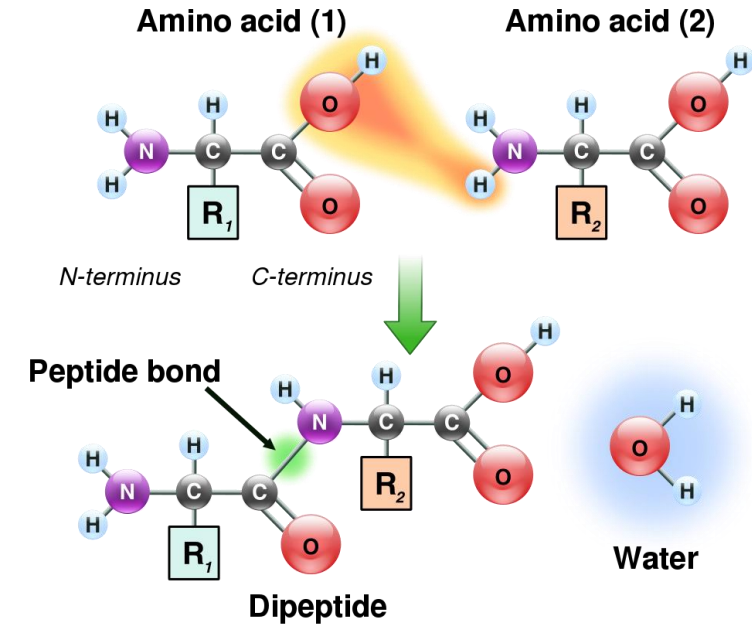
Corresponding RNA sequence



Proteins – chain of 20 (+2) amino acids



Amino acid	3-letter ^[132]	1-letter ^[132]
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



FASTA format

>DROME_HH_Q02936

MRHIAHTQRCLSRSLTSLVALLLIVLPMVFSPAHCSPGRGLGRHRARNLYPLVL
KQTIPNLSEYNTSASGPLEGVIRRDSPKFKDLVPNYNRDILFRDEEGTGADRLM
SKRCKEKLNVLAYSVSMNEWPGIRLLVTESWDEYHHGQESLHEGRAVTIATS
DRDQSKYGMLARLAVEAGFDWVSYSRRHIYCSVKSDSSISSHVHGCFTPES
TALLESQVVRKPLGELSIGDRVLSMTANGQAVYSEVILFMDRNLEQMNFVQLH
TDGGAVLTVTPAHLVSVWQPESQKLTFVFADRIEKNQVLVRDVTGELRPQR
VVKVGSVRSKGVVAPLTREGTIVVNSVAASCYAVINSQSLAHWGLAPMRLST
LEAWLPAKEQLHSSPKVVSSAQQQNGIH WYANALYKVKDYVLPQSWRHD

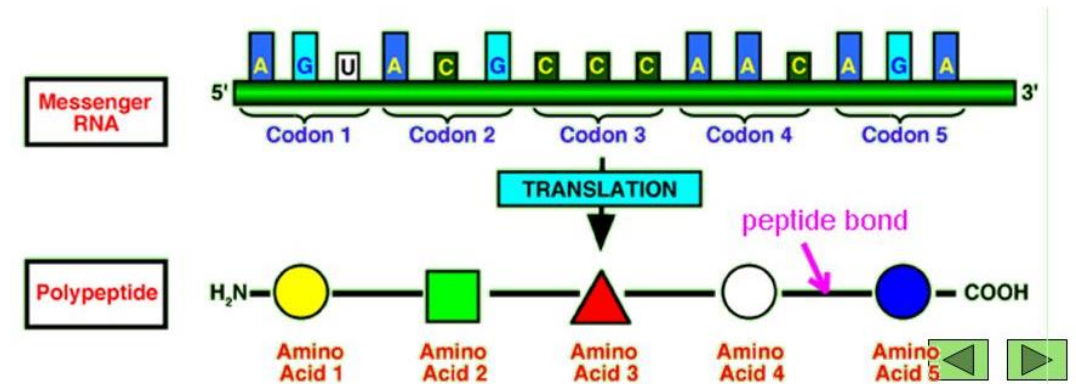
Where to get protein sequences

UNIPROT: <http://www.uniprot.org/>

Translation (RNA->protein): The Genetic Code

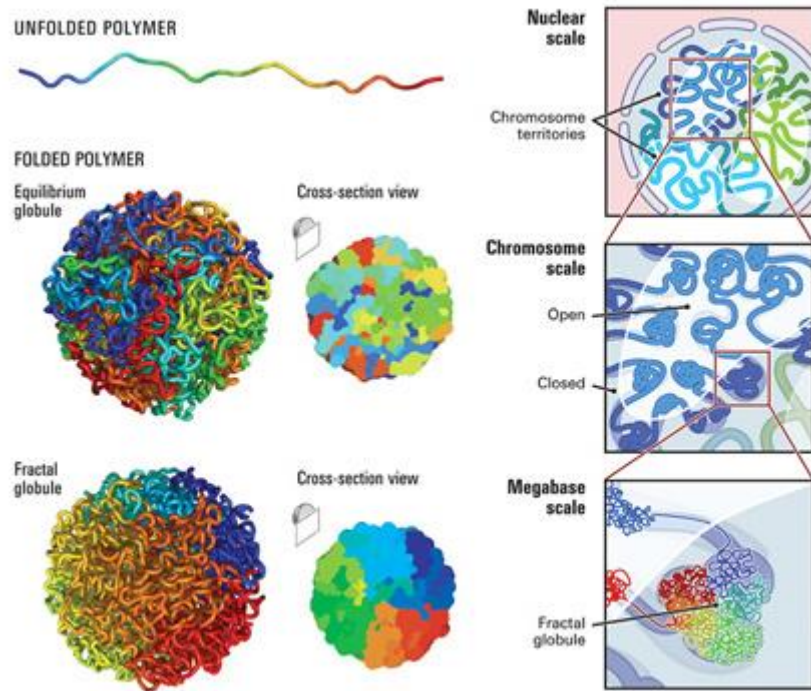
- A *triplet* of nucleic acids (*codon*) codes for one amino acid
- The code is redundant. E.g., both GGC and GGA code for Gly (Glycine)

		Second Base					
		U	C	A	G		
First Base	U	UUU } Phenylalanine F UUC } UUA } Leucine L UUG }	UCU } UCC } Serine S UCA } UCG }	UAU } Tyrosine Y UUA } Stop codon UAG } Stop codon	UGU } Cysteine C UGC } UGA } Stop codon UGG } Tryptophan W	Third Base	U C A G
	C	CUU } Leucine L CUC } CUA } CUG }	CCU } CCC } Proline P CCA } CCG }	CAU } Histidine H CAC } CAA } Glutamine Q CAG }	CGU } CGC } Arginine R CGA } CGG }		U C A G
	A	AUU } Isoleucine I AUC } AUA } Methionine start codon M AUG }	ACU } ACC } Threonine T ACA } ACG }	AAU } Asparagine N AAC } AAA } Lysine K AAG }	AGU } Serine S AGC } AGA } Arginine R AGG }		U C A G
	G	GUU } Valine V GUC } GUA } GUG }	GCU } GCC } Alanine A GCA } GCG }	GAU } Aspartic acid D GAC } GAA } Glutamic acid E GAG }	GGU } GGC } Glycine G GGA } GGG }		U C A G



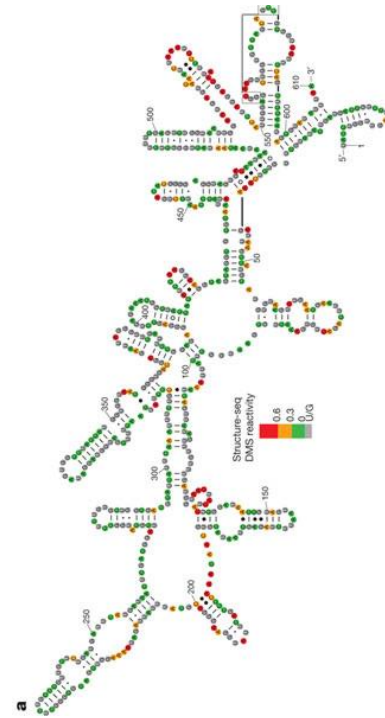
Picture: https://media1.shmoop.com/images/biology/biobook_dna_graphik_22.png

DNA, RNA and proteins form complex secondary and tertiary structures

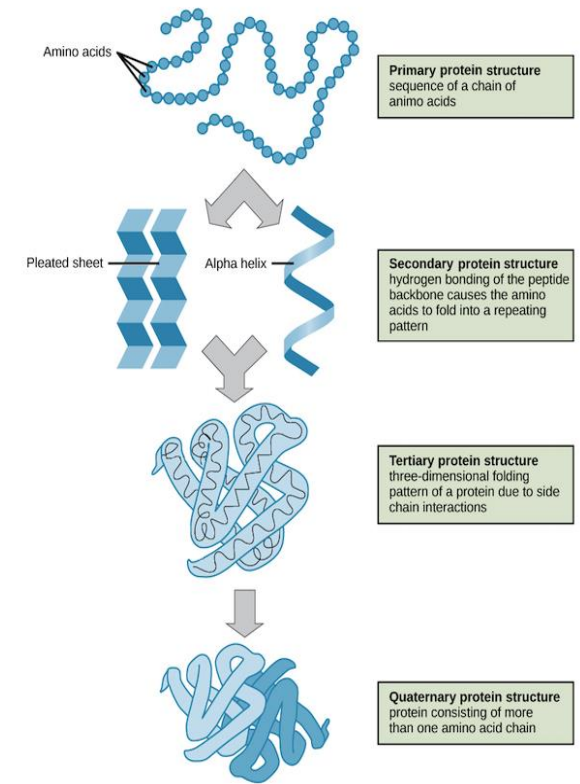


3D genome architecture
(DNA folding)

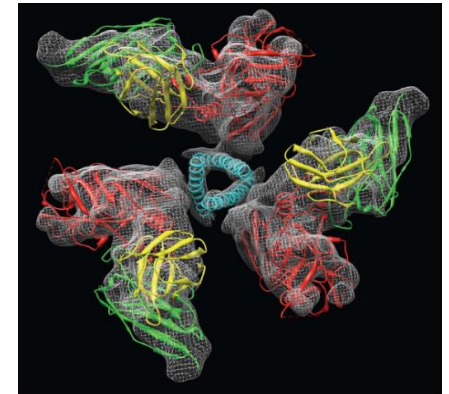
Major ML challenge: Predicting structure



RNA secondary structure



<https://ka-perseus-images.s3.amazonaws.com/71225d815cafcc09102504abdf4e10927283be98.png>



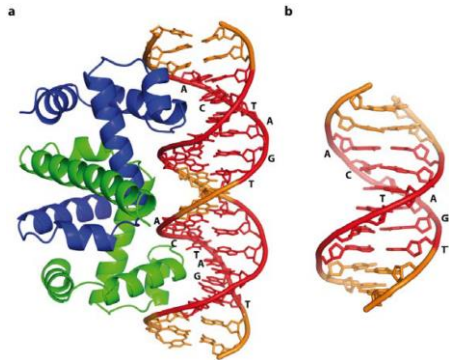
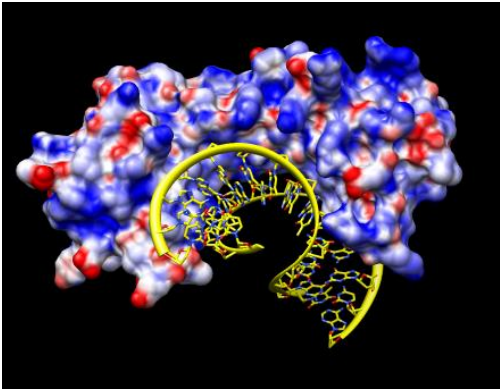
https://www.fei.com/uploadedImages/FEISite/Pages/Markets/Life_Sciences/Structural_Biology/Sub-Cellular_Imaging/Sriram_001_465x.png

Protein structure

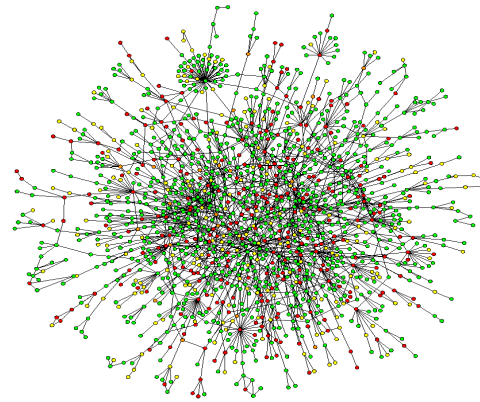
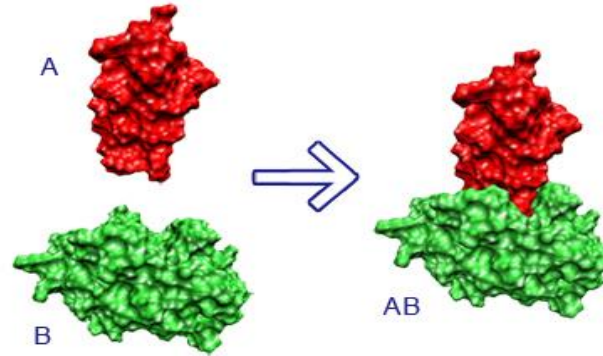
Interactions between proteins, DNA, RNA and small molecules (drugs)

Major ML challenge: Predicting interactions between biomolecules

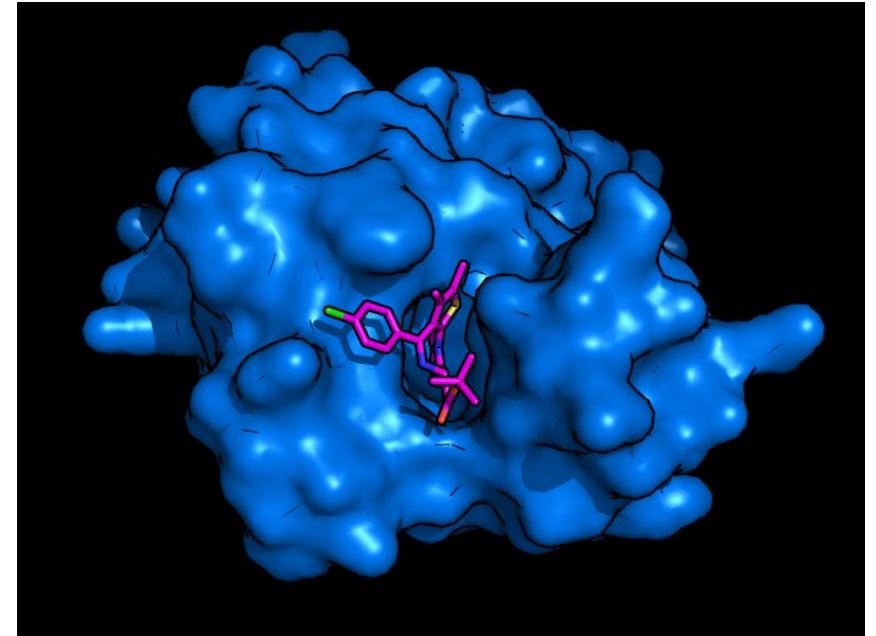
Protein-DNA interactions



Protein-Protein interactions



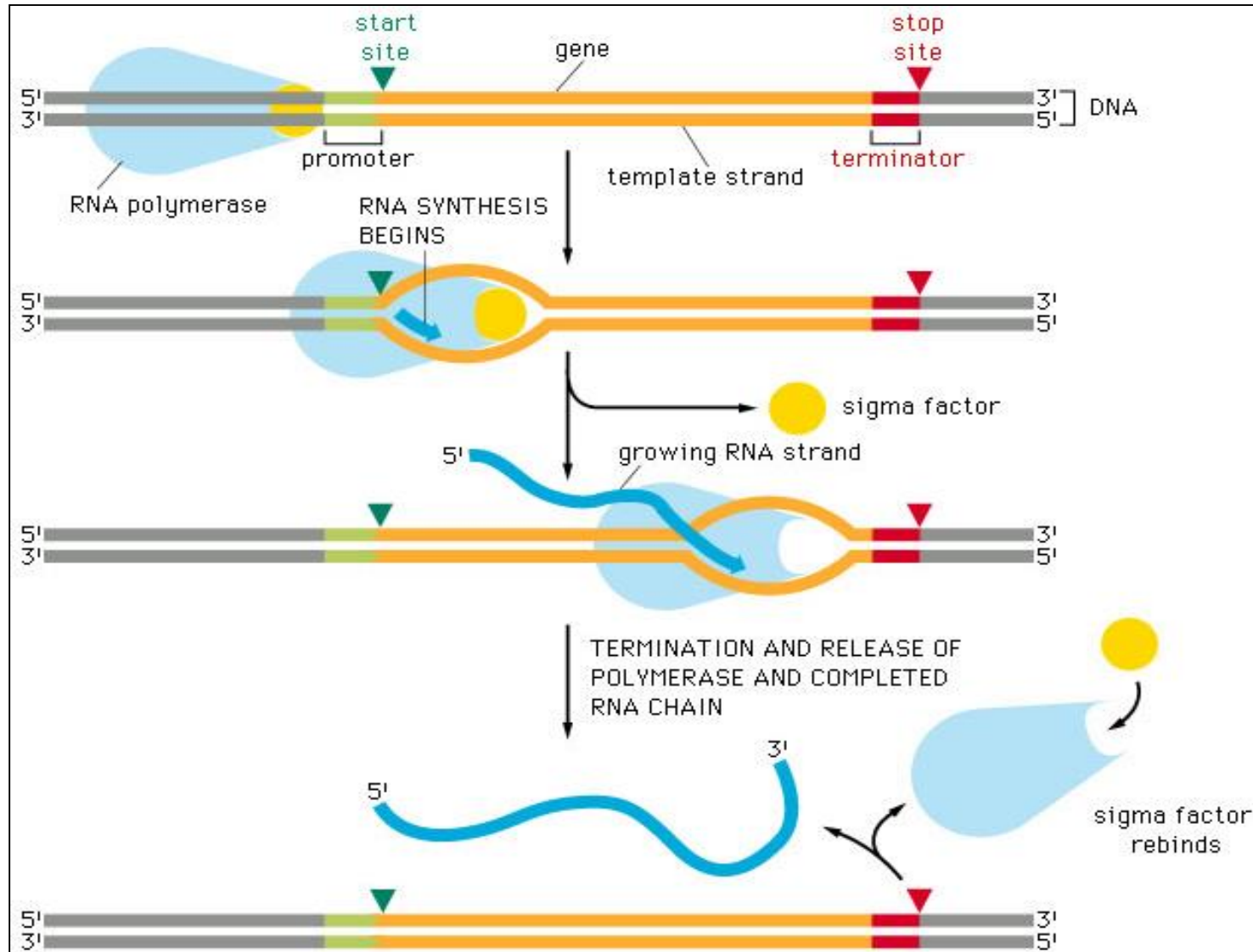
PPI networks



Protein-small molecule interactions

(Drug discovery, cheminformatics,
Quantitative structure activity
relationships (QSAR))

From DNA to RNA: Transcription

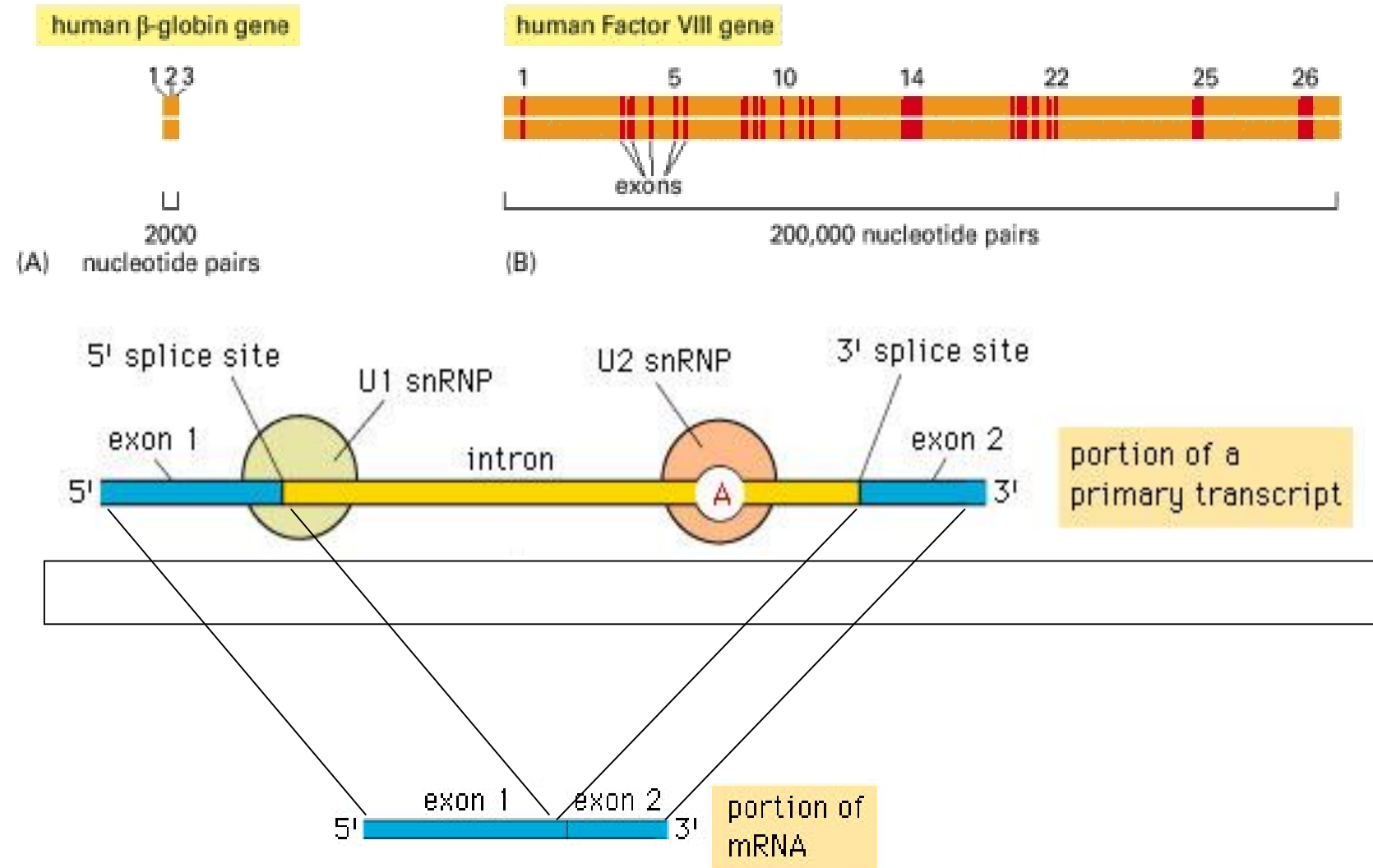


Machine learning problems

1. Predict genes in the genome
2. Predict gene transcription start sites (TSS) and termination sites (TTS)

Parts of a gene: Exons and Introns (Splicing)

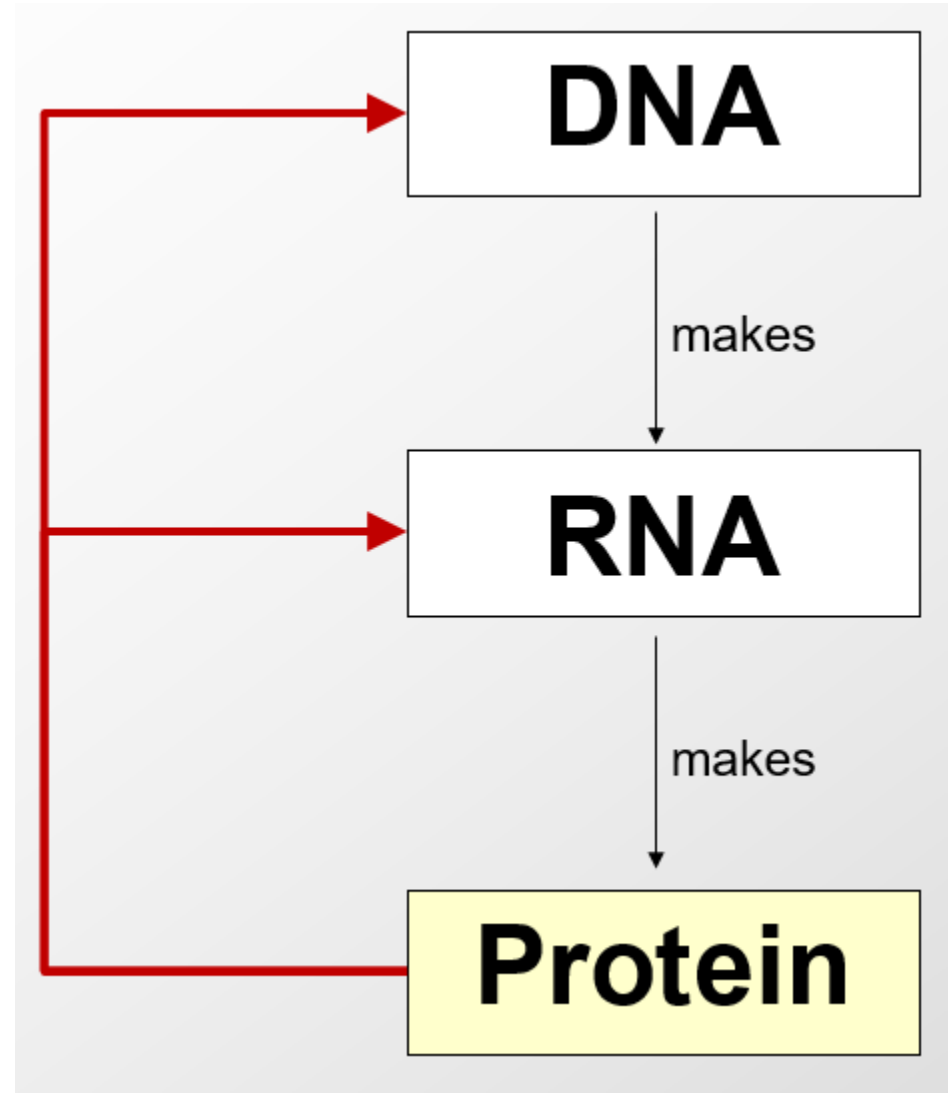
- Not every part of a gene is coding
- Exons (coding) interrupted by non-translated introns
- Introns are **spliced** out
- **Alternative splicing:** different exon subsets for the same gene => different protein products



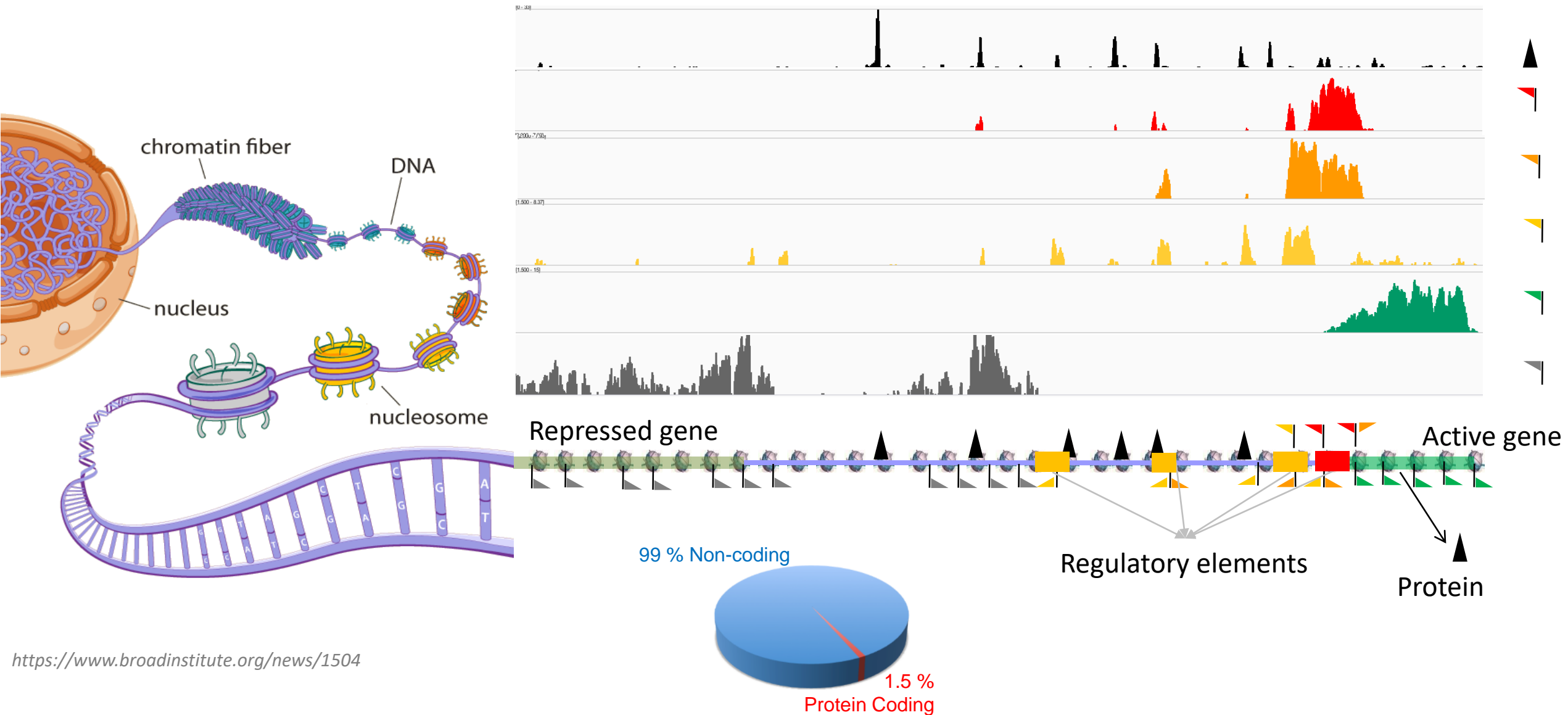
Major ML challenge: Predict splice sites, gene isoforms and alternative splicing events

What is gene regulation?

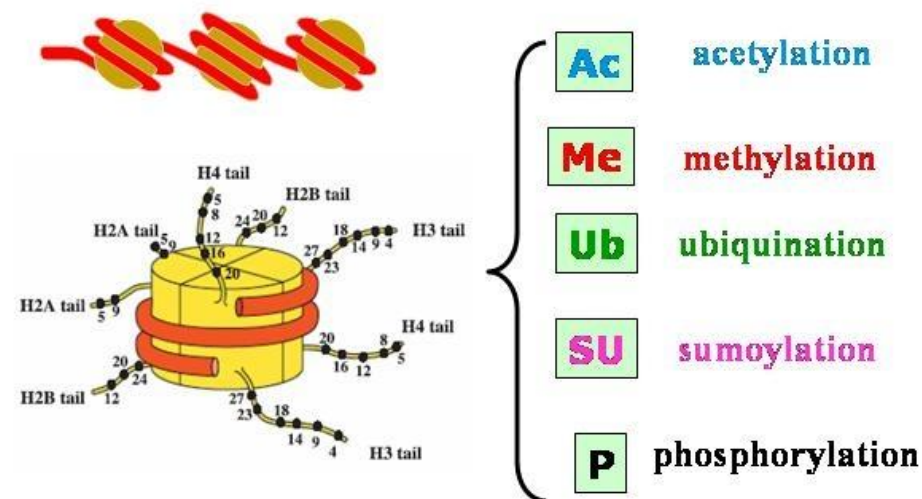
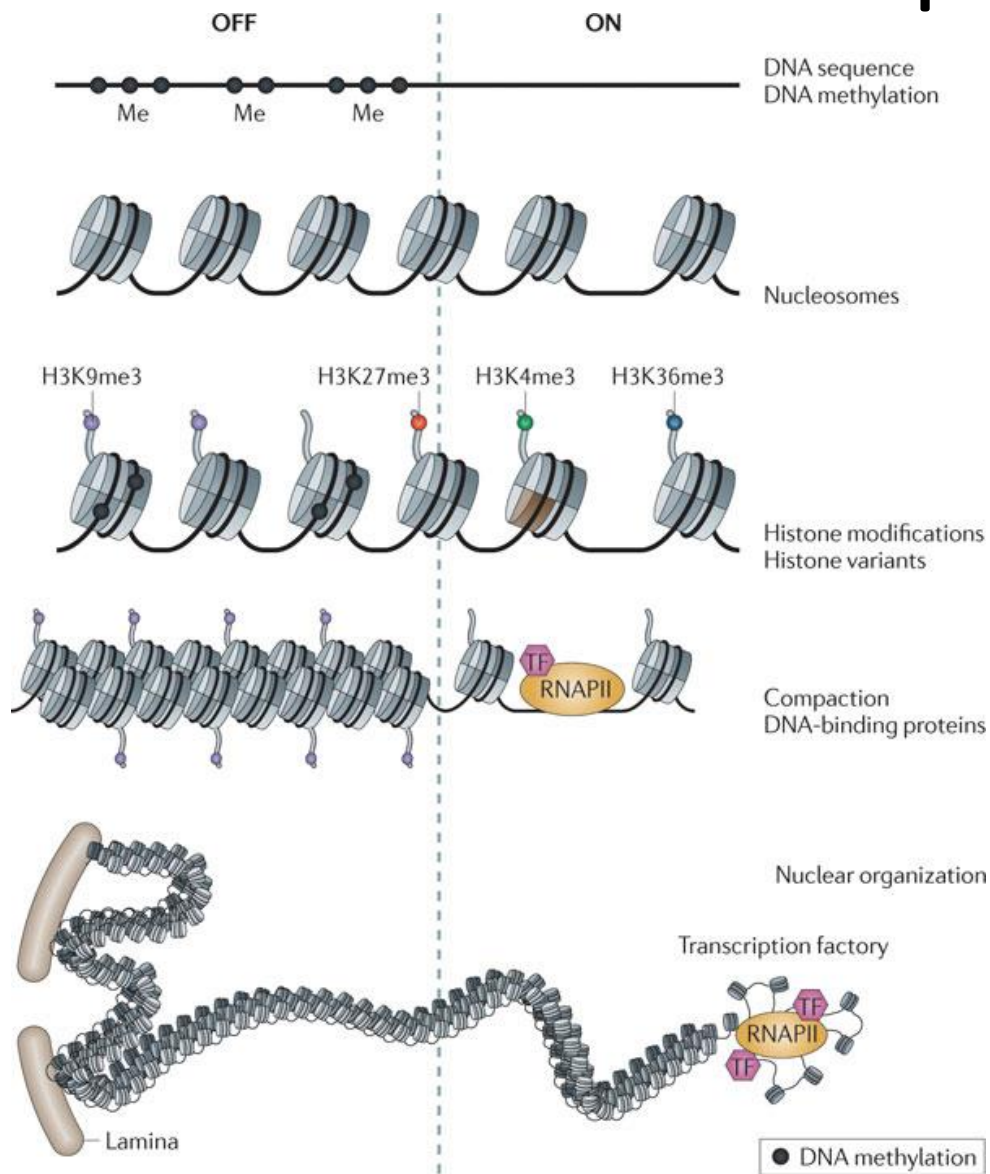
Gene regulation
(When? Where?
How much?)



Regulatory control elements and epigenomic marks



Chromatin and epigenomic modifications

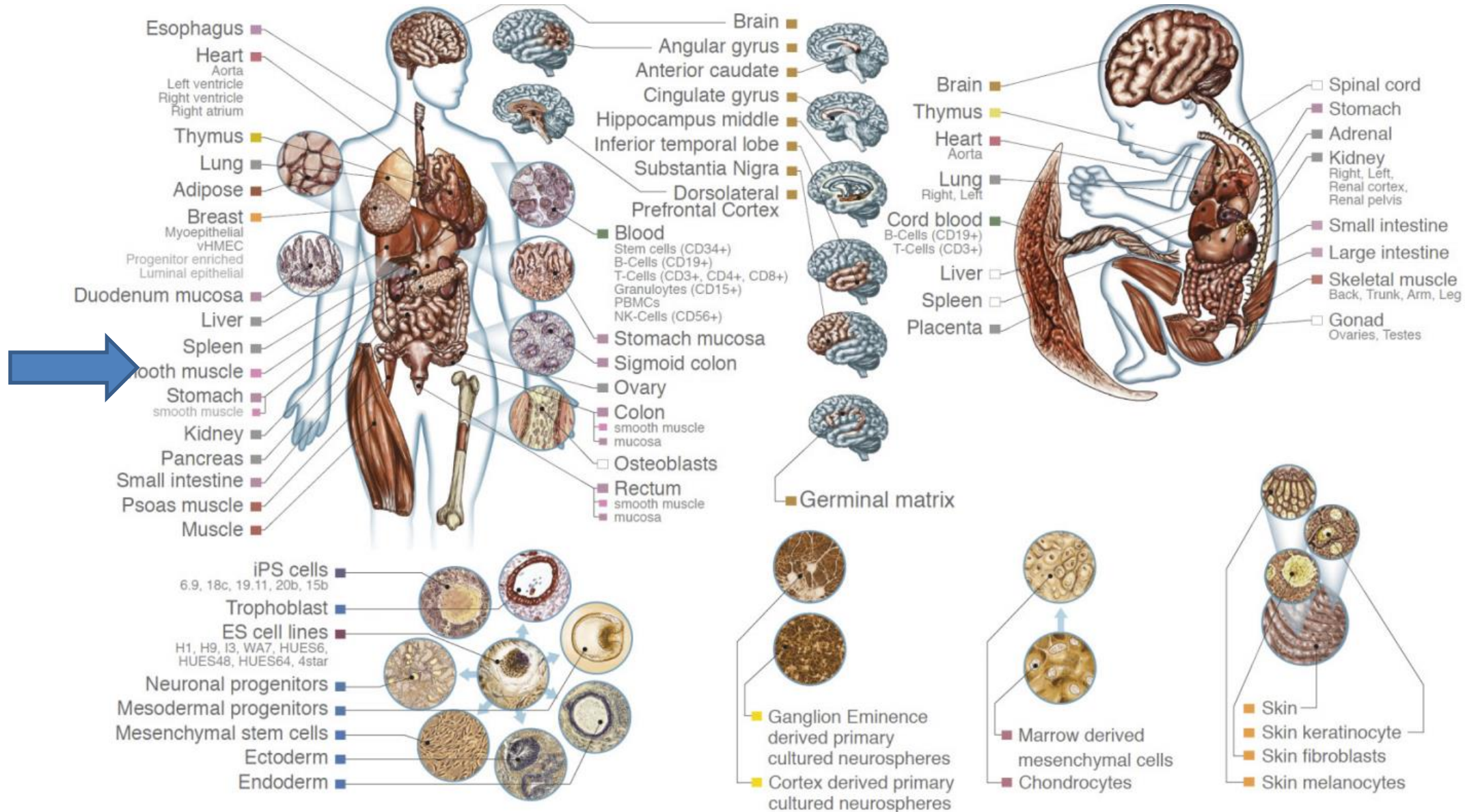


The figure illustrates nucleosome models and major posttranslational modifications which play essential roles in gene expression regulation and disease processes

Major ML challenge: Predict combinatorial epigenomic code

One genome \Leftrightarrow Many cell types

ACCAGTTACGACGG
 TCAGGGTACTGATA
 CCCCAAACCGTTGA
 CCGCATTTACAGAC
 GGGGTTTGGGTTTT
 GCCCCACACAGGTA
 CGTTAGCTACTGGT
 TTAGCAATTTACCG
 TTACAACGTTTACA
 GGGTTACGGTTGGG
 ATTTGAAAAAAGT
 TTGAGTTGGTTTTT
 TCACGGTAGAACGT
 ACCTTACAAA.....



Differential activation/repression of control elements and genes defines cell type identity and state



- Active control elements
- Active control elements
- Active genes
- Repressed elements

- ~25,000 genes
- ~2 million novel putative control elements!
- cell-type specific usage of elements

Major ML challenges:

1. Predict regulatory elements (REs) and their cell type specific activity
2. Predict which REs regulate which genes in which cell types
3. Predict gene expression from regulatory element activation

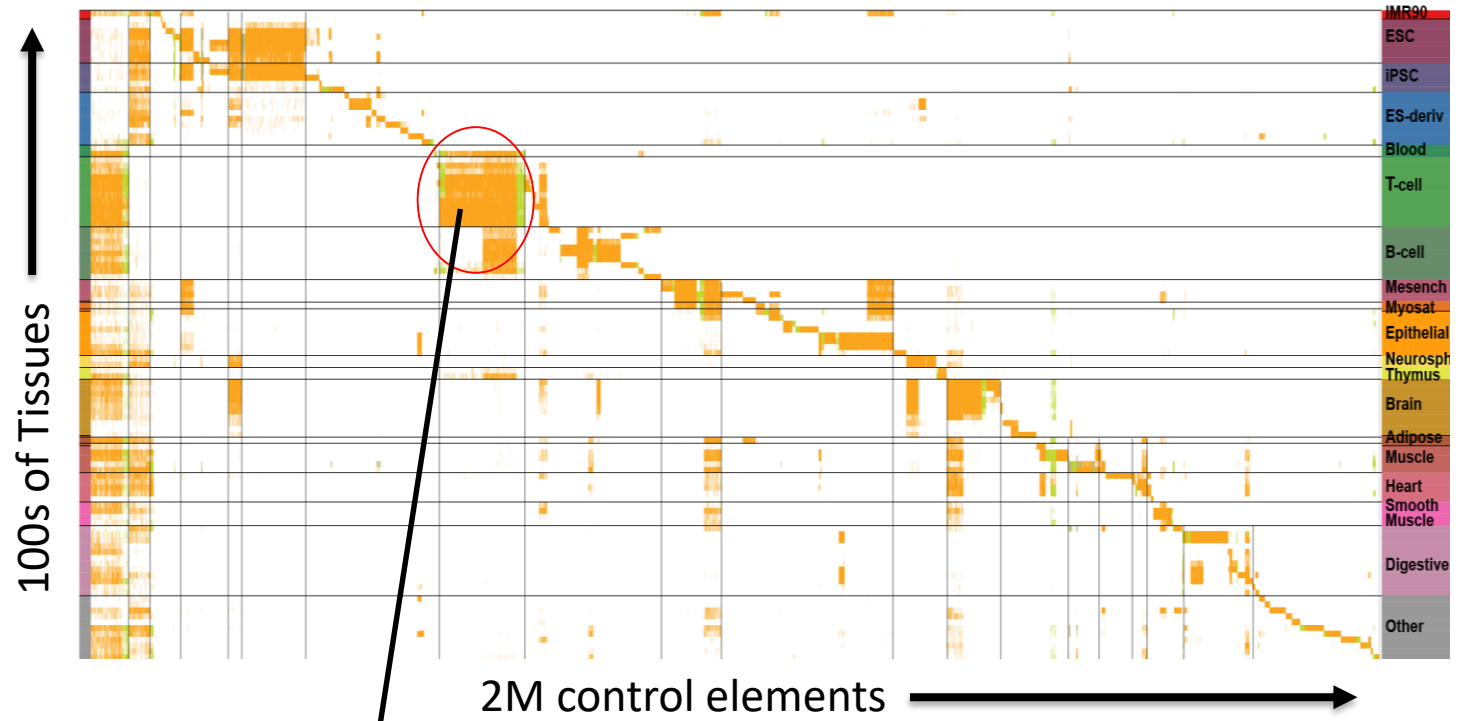
What code in the non-coding DNA controls cell-type specific activation of regulatory elements?



Regulatory proteins bind DNA words (landing pads) in control elements!



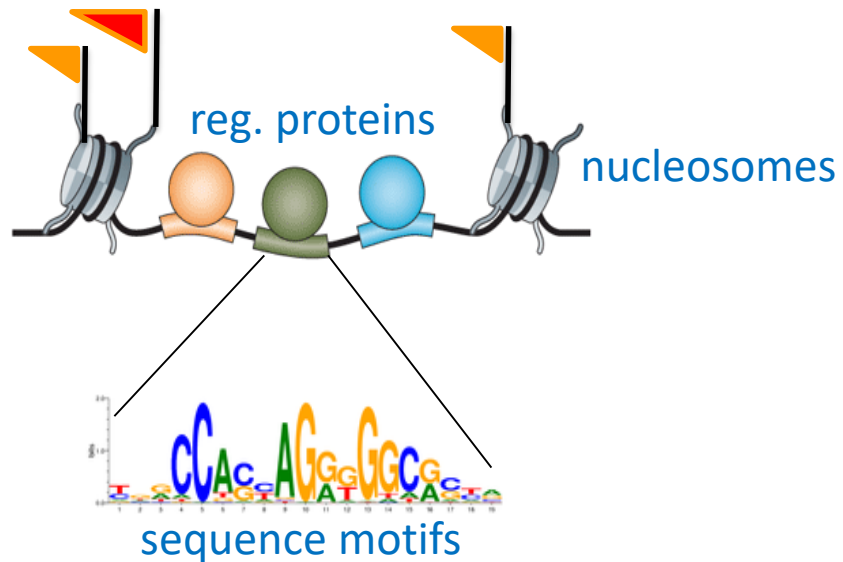
‘Motif Discovery’



AATCAGTTATCTGTTGTATACCCGGAGTCC
AGGTCGAATGCAAAACGGTTCTTGACGTA
GAGATAACCGCTTGATATGACTCATTTGCC
ATATTCCGGACGCTGTGACGATCCGGTTTG
GAACGCAAACAGTTCAGTGCTTATCATGAA

The multiple facets of genomic regulatory code

epigenomic marks



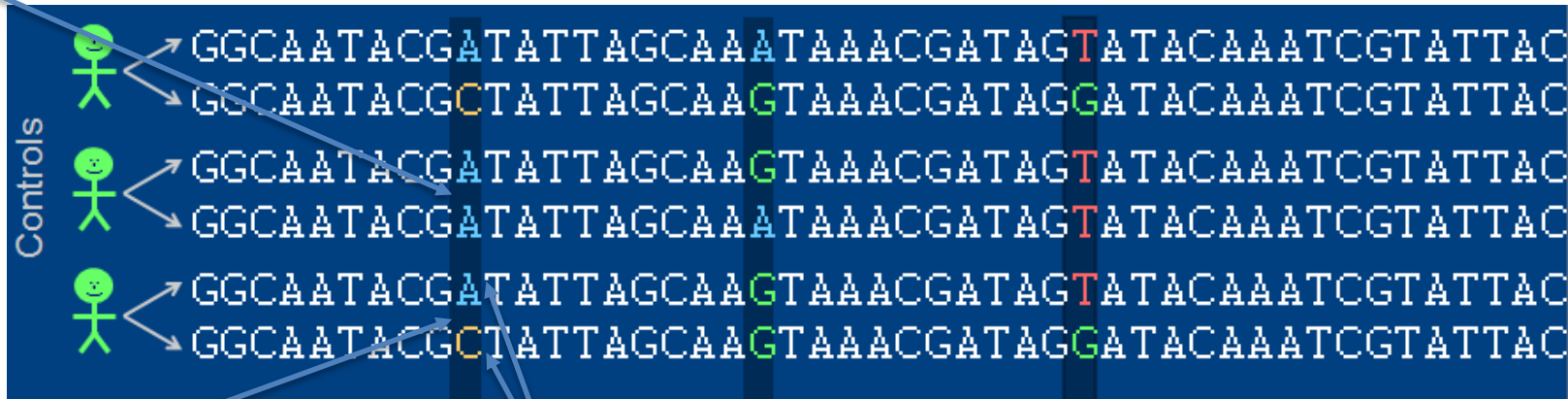
Adapted from Shlyueva et al. (2014) Nature Reviews Genetics.

Major ML challenges:

1. Predict DNA sequence affinity of individual proteins
2. Predict binding landscape of proteins in different cell types and states
3. Predict combinatorial binding patterns of proteins
4. Predict combinatorial regulatory grammars encoded in non-coding regulatory elements

Genetic variation across individuals

Homozygous
(identical alleles)



Heterozygous
(different alleles)

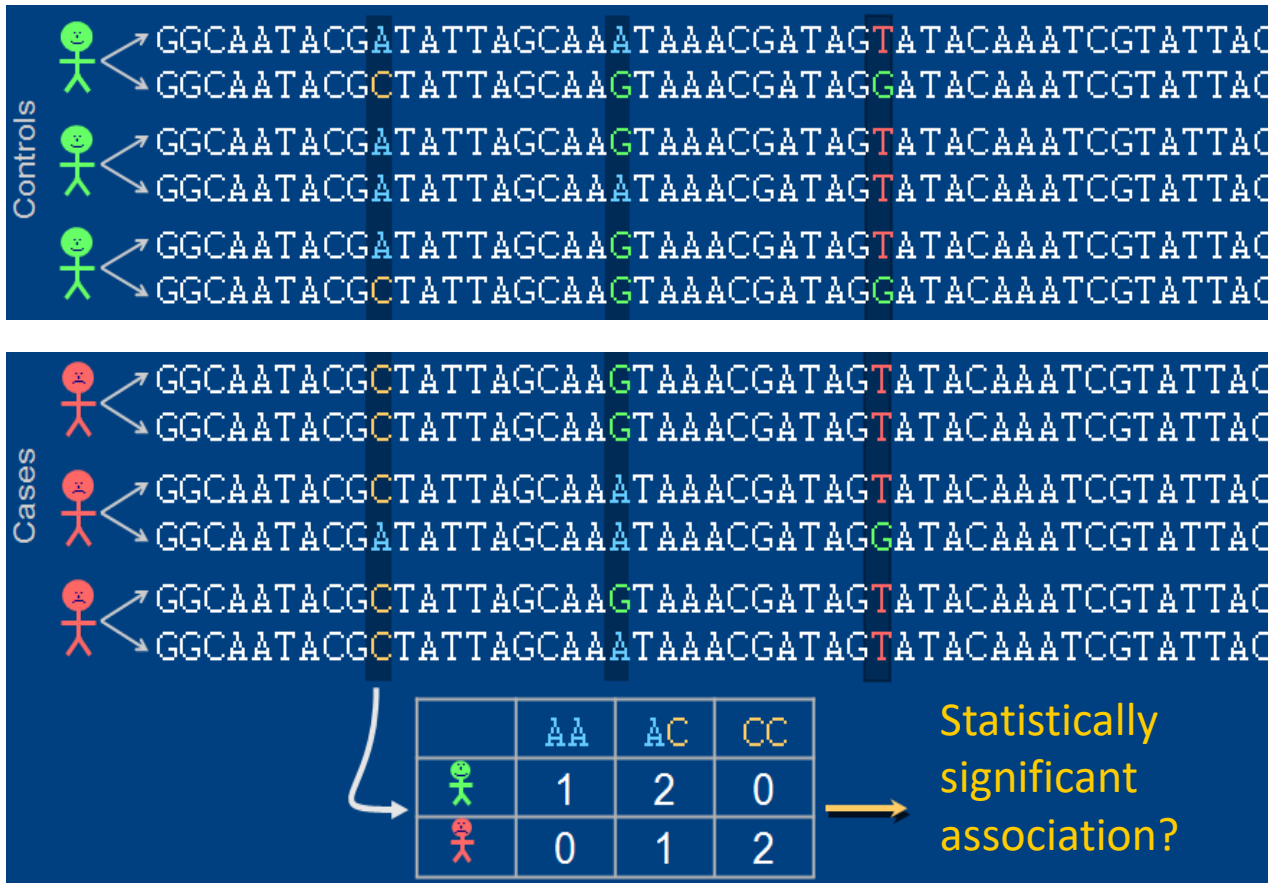
Human genome: ~3 billion bp X 2 copies

Types of genetic Variants:

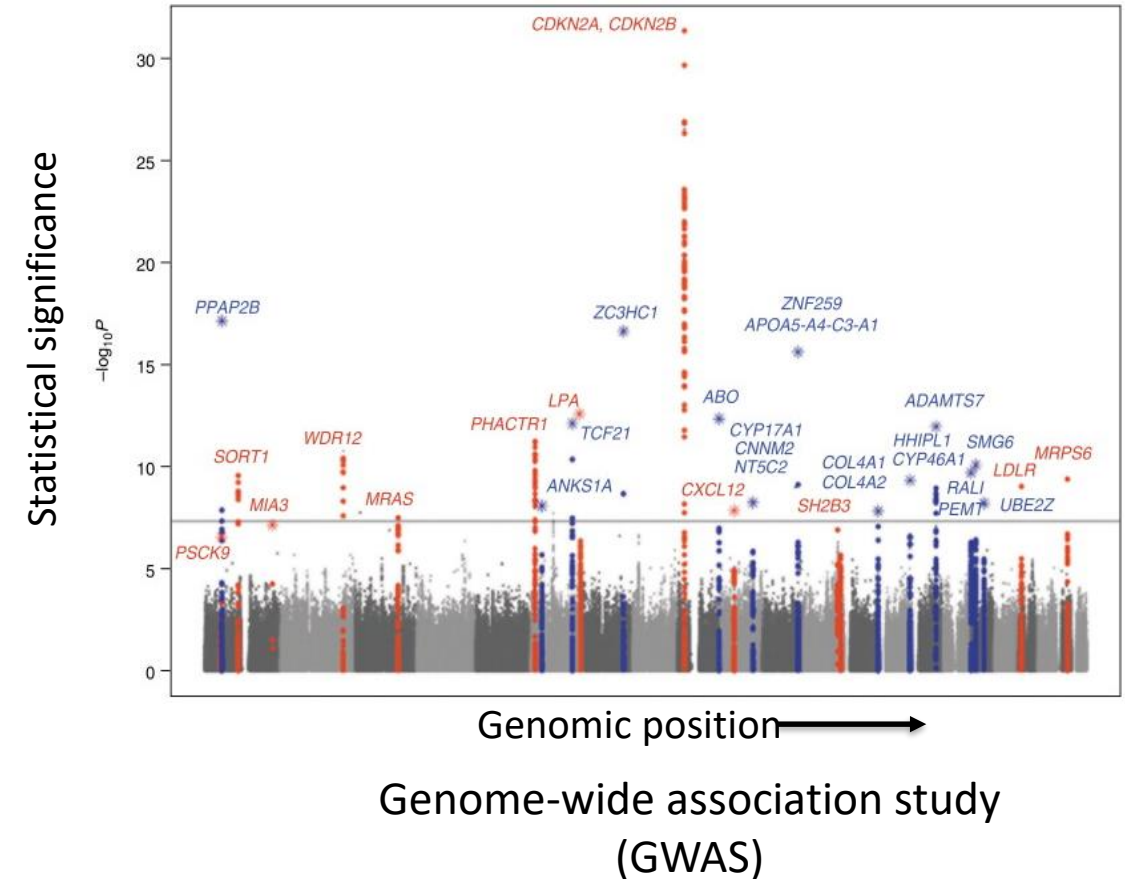
- Single nucleotide polymorphisms (SNPs): substitutions
 - ~3M **common** SNPs (> 2-5% minor allele frequency (MAF) in human population)
 - Rare (< 2-5% MAF frequency) and private SNPs (in a single individual)
- Short and large insertions, deletions, inversions, translocations

Germline vs. Somatic mutations: Genetic alteration acquired by a cell that can be passed to the progeny of the **mutated** cell in the course of cell division. **Somatic mutations** differ from **germ line mutations**, which are inherited genetic alterations that occur in the germ cells (i.e., sperm and eggs)

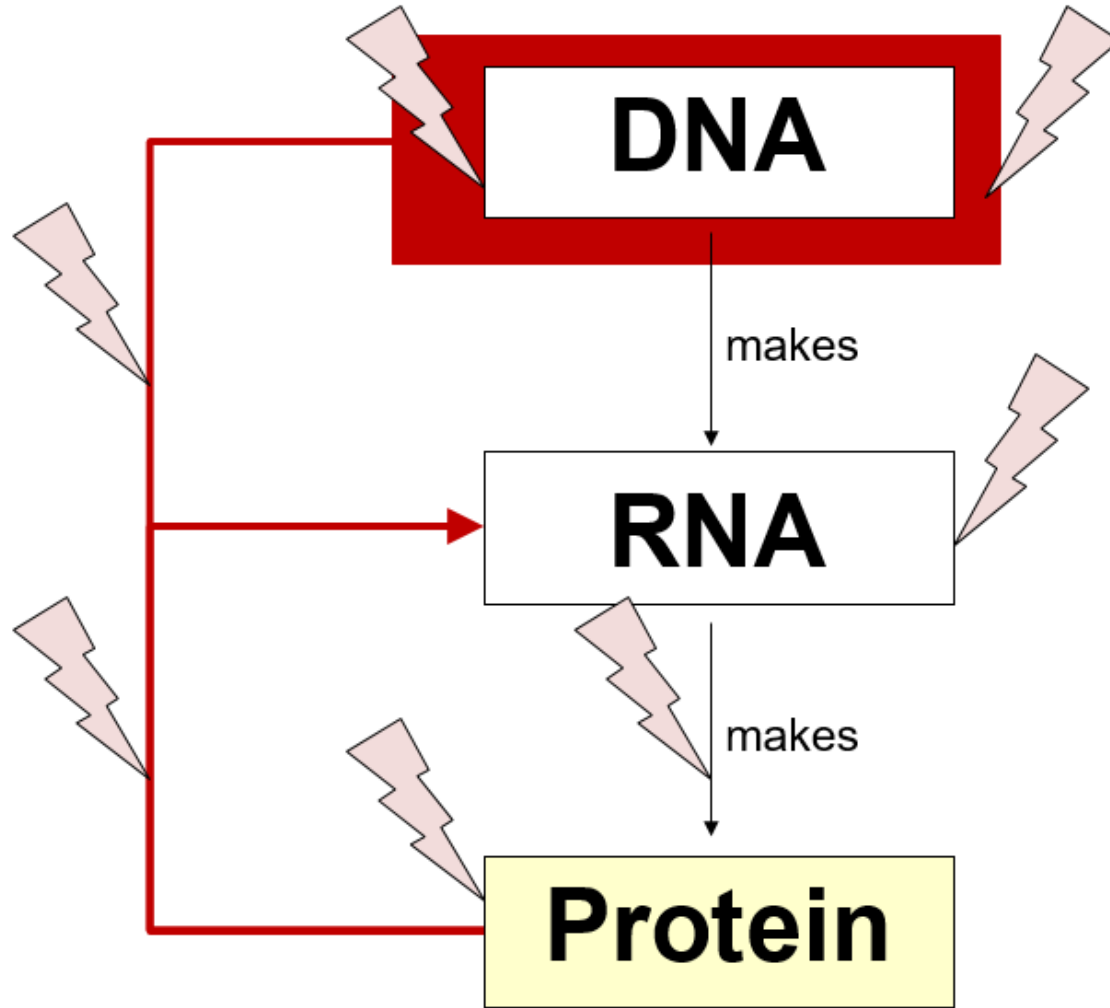
Case-control studies to identify disease-associated genetic variants



Correlation structure between variants makes it difficult to identify causal variant



The role of genetic variation in regulation

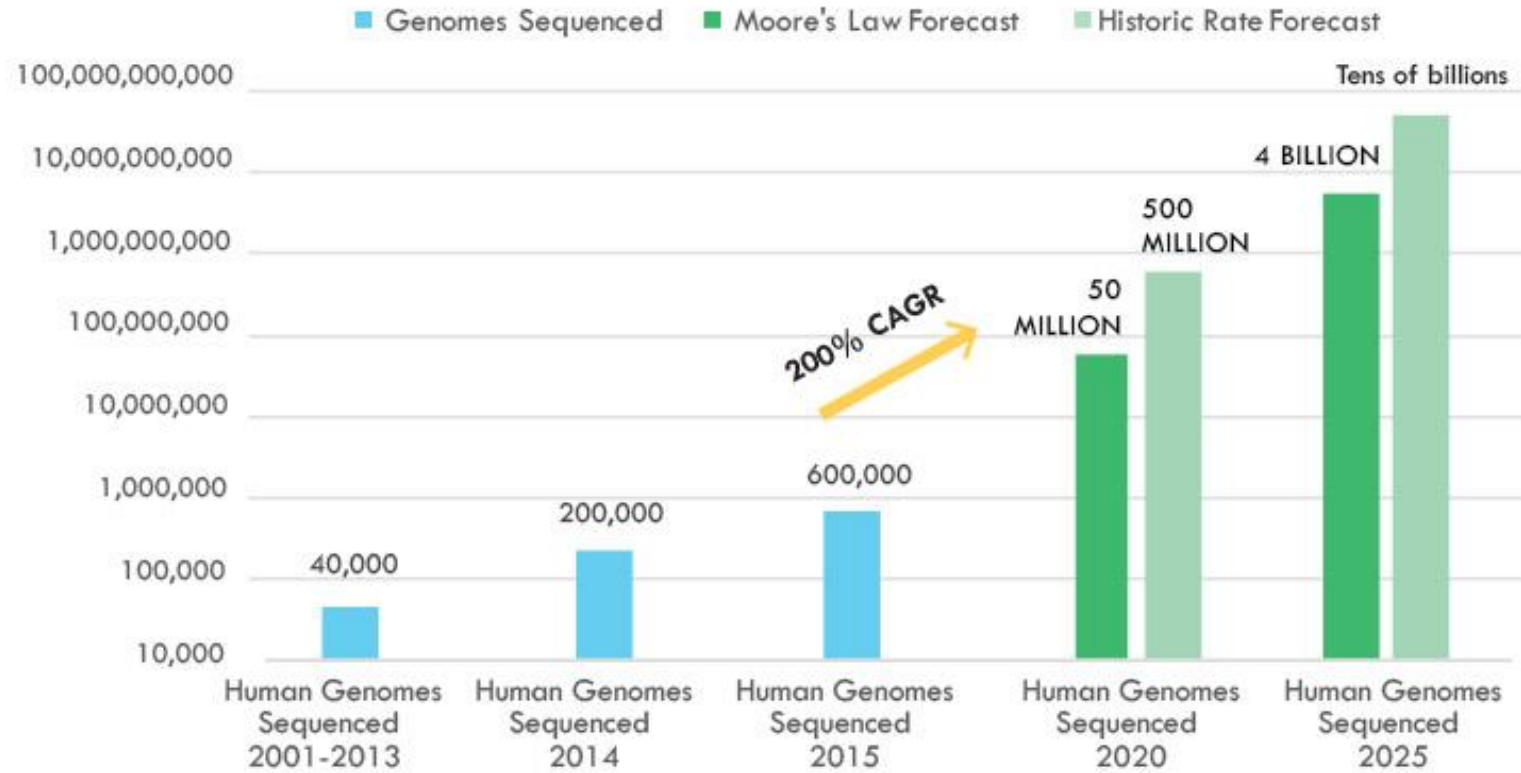


Major computational challenges:

1. Predict genetic variants from genome sequencing data
2. Which genetic variants are benign vs. harmful i.e. which genetic variants are associated with different phenotypes (diseases/traits)
3. What is the genetic architecture of a disease/trait (relationship between variants to phenotype)
4. Predict the molecular effect of a genetic variant i.e. how it affects cellular function

Introduction to high-throughput sequencing

The Number of Human Genomes Sequenced (log scale)



Source: National Human Genome Research Institute (NHGRI), ARK Investment Management LLC



GA II
1.6 billion bp per day
(2008)



GA IIX
5 billion bp per day
(2009)



HiSeq 2500
60 billion bp per day
(2012)

Images: www.illumina.com/systems

Numbers: www.politigenomics.com/next-generation-sequencing-informatics

Dates: Illumina press releases



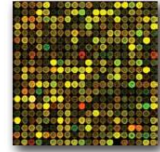
Oxford Nanopore
technology

Sequencing technologies



Sanger DNA sequencing

1977-1990s



DNA Microarrays

Since mid-1990s



2nd-generation DNA sequencing

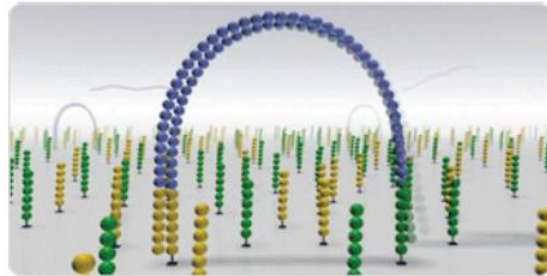
Since ~2007



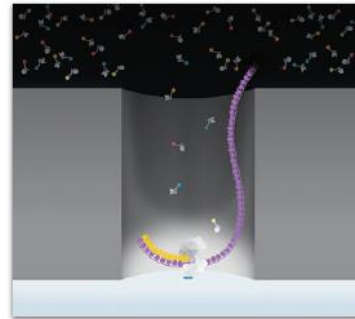
3rd-generation & single-molecule DNA sequencing

Since ~2010

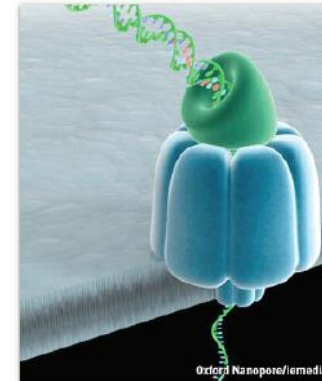
Since 2005, many DNA sequencing instruments have been described and released. They are based on a few different principles



Synthesis / ligation



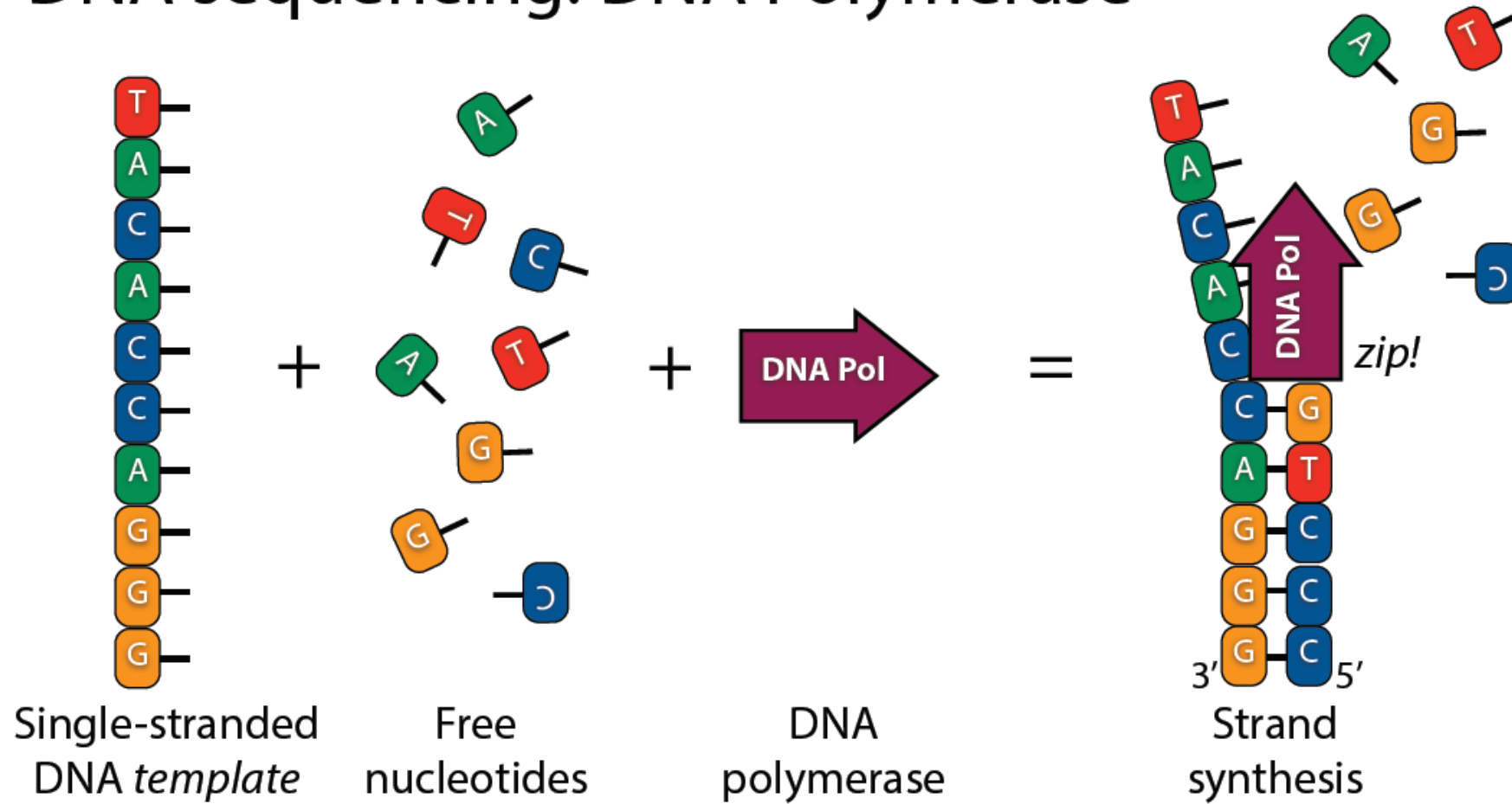
SMRT cell



Nanopore

Sequencing by synthesis (“massively parallel sequencing”) provides greatest throughput, and is the most prevalent today

DNA sequencing: DNA Polymerase

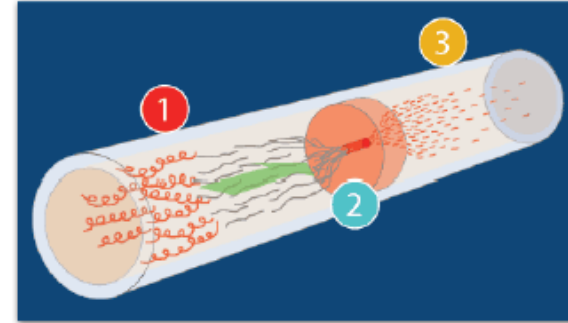


DNA polymerase moves along the template in one direction, integrating complementary nucleotides as it goes

Sequencing by synthesis

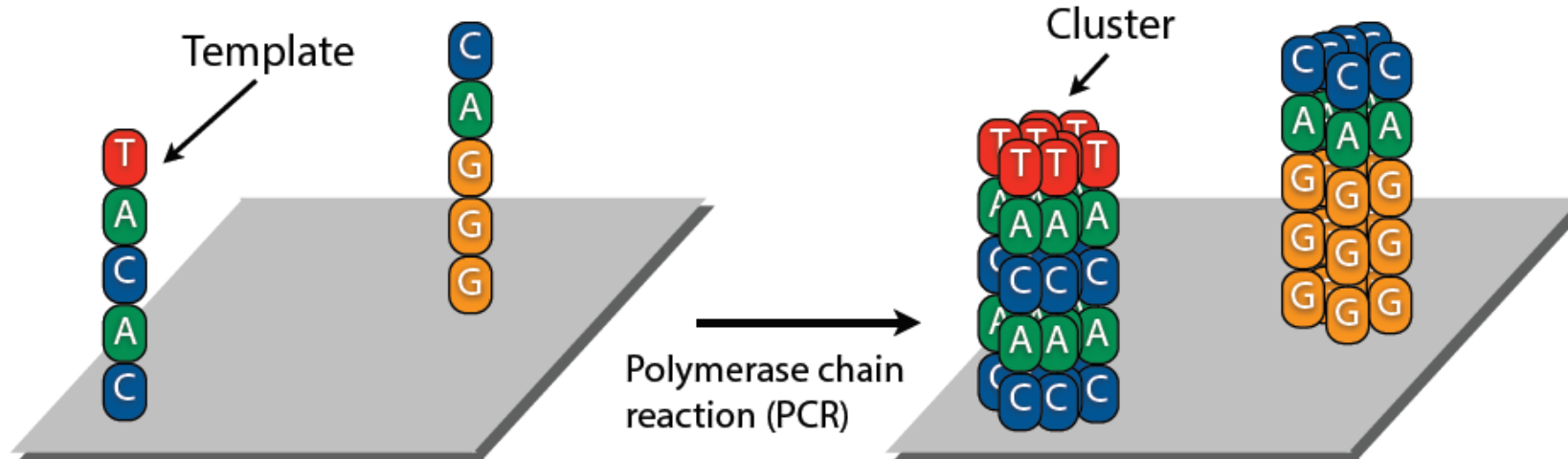
1. Take DNA sample, which includes many copies of the genome, and chop it into single-stranded fragments (“templates”)

E.g. with ultrasound waves,
water-jet shearing (pictured),
divalent cations



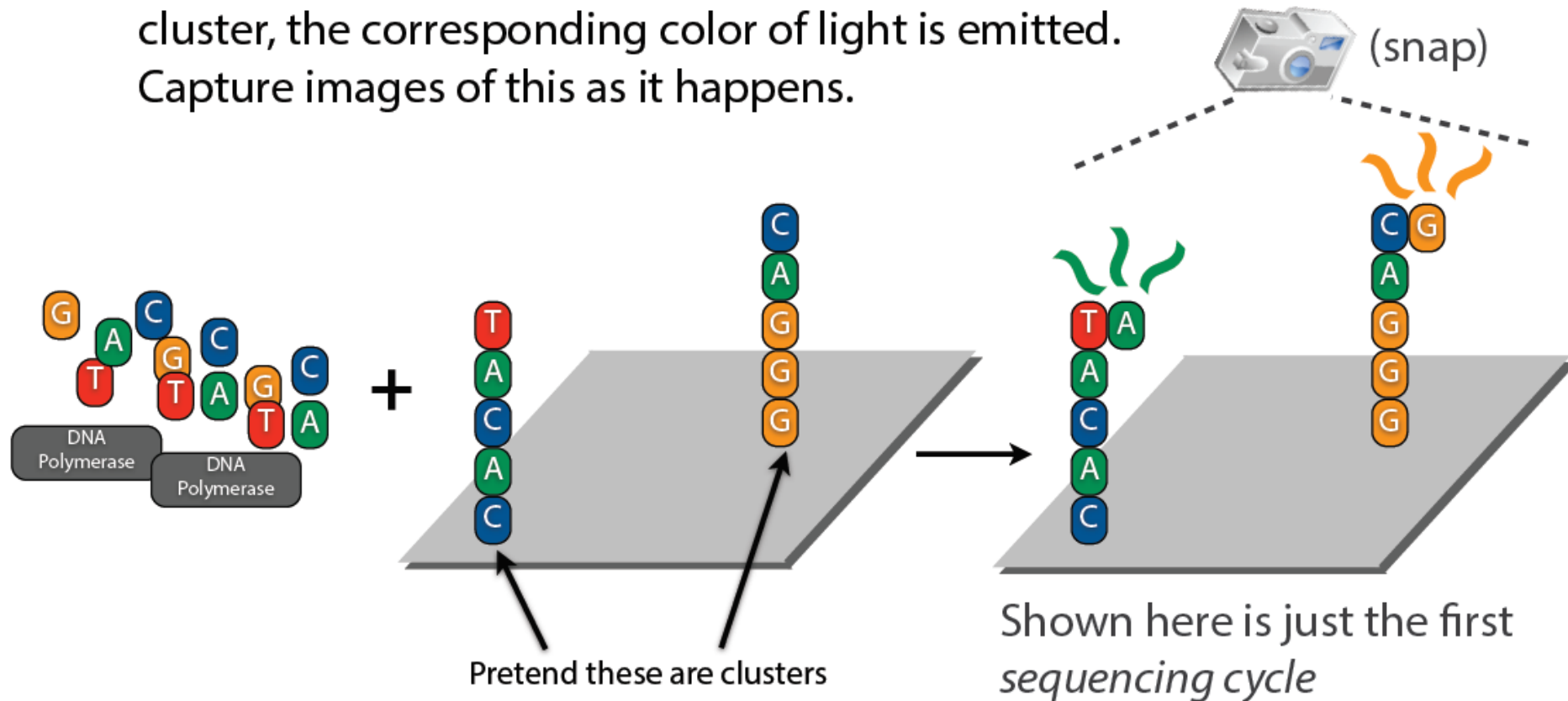
Picture: http://www.jgi.doe.gov/sequencing/education/how/how_1.html

2. Attach templates to a surface
3. Make copies so that each template becomes a “cluster” of clones



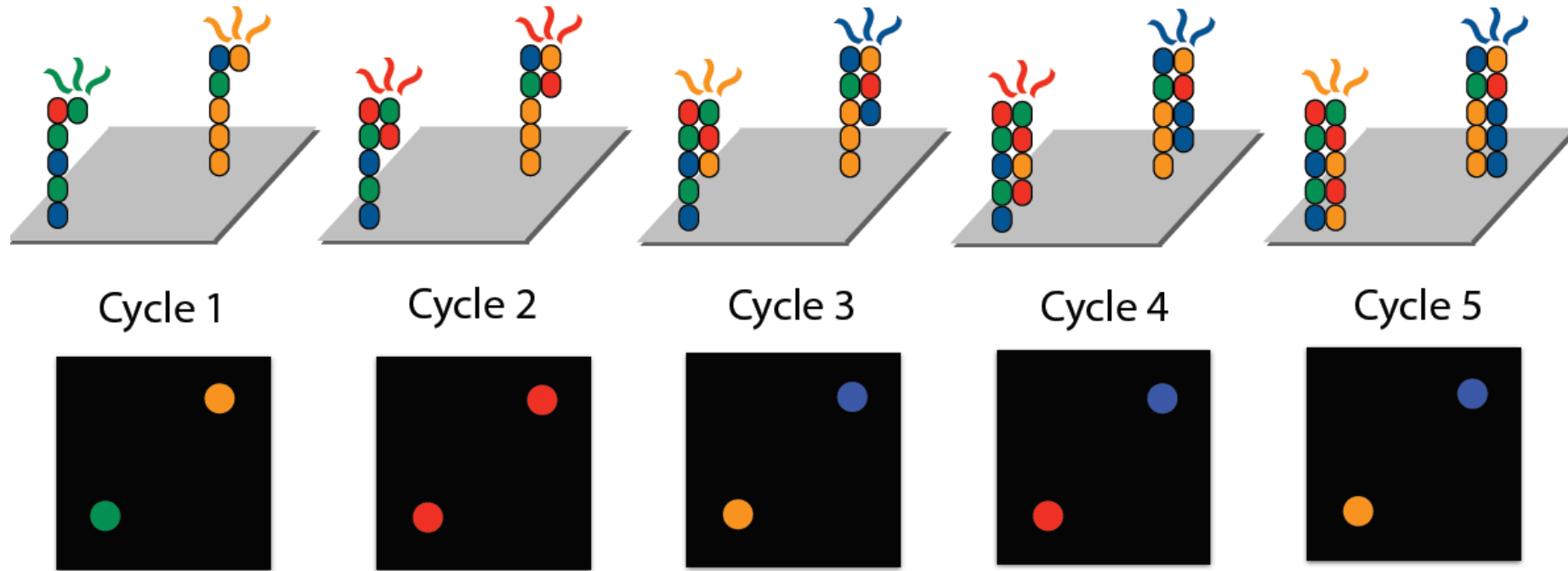
Sequencing by synthesis

4. Repeatedly inject mixture of *color-labeled* nucleotides (A, C, G and T) and DNA polymerase. When a complementary nucleotide is added to a cluster, the corresponding color of light is emitted. Capture images of this as it happens.



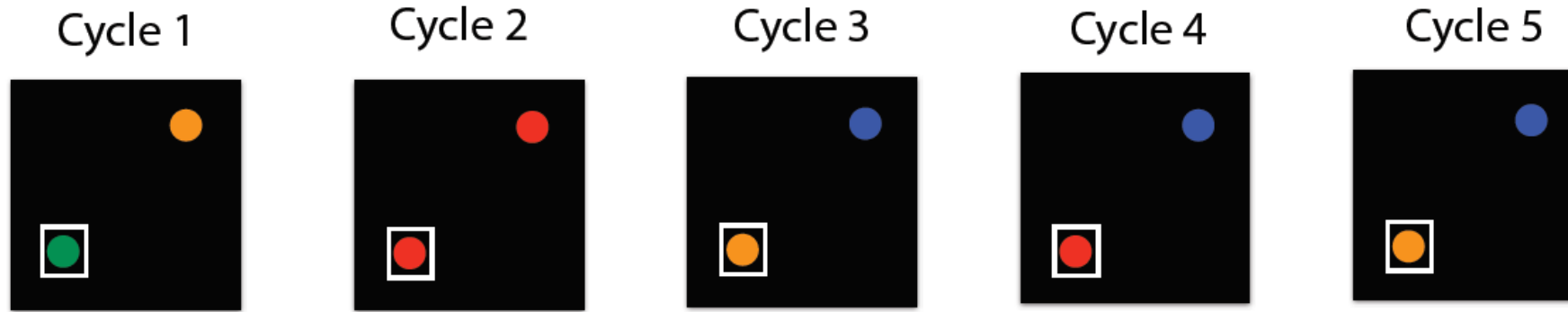
Sequencing by synthesis

5. Line up images and, for each cluster, turn the series of light signals into corresponding series of nucleotides

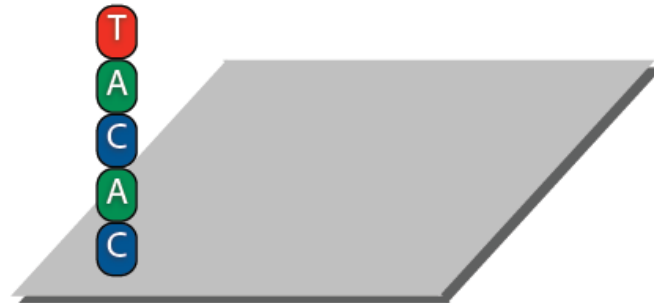


Sequencing by synthesis

5. Line up images and, for each cluster, turn the series of light signals into corresponding series of nucleotides

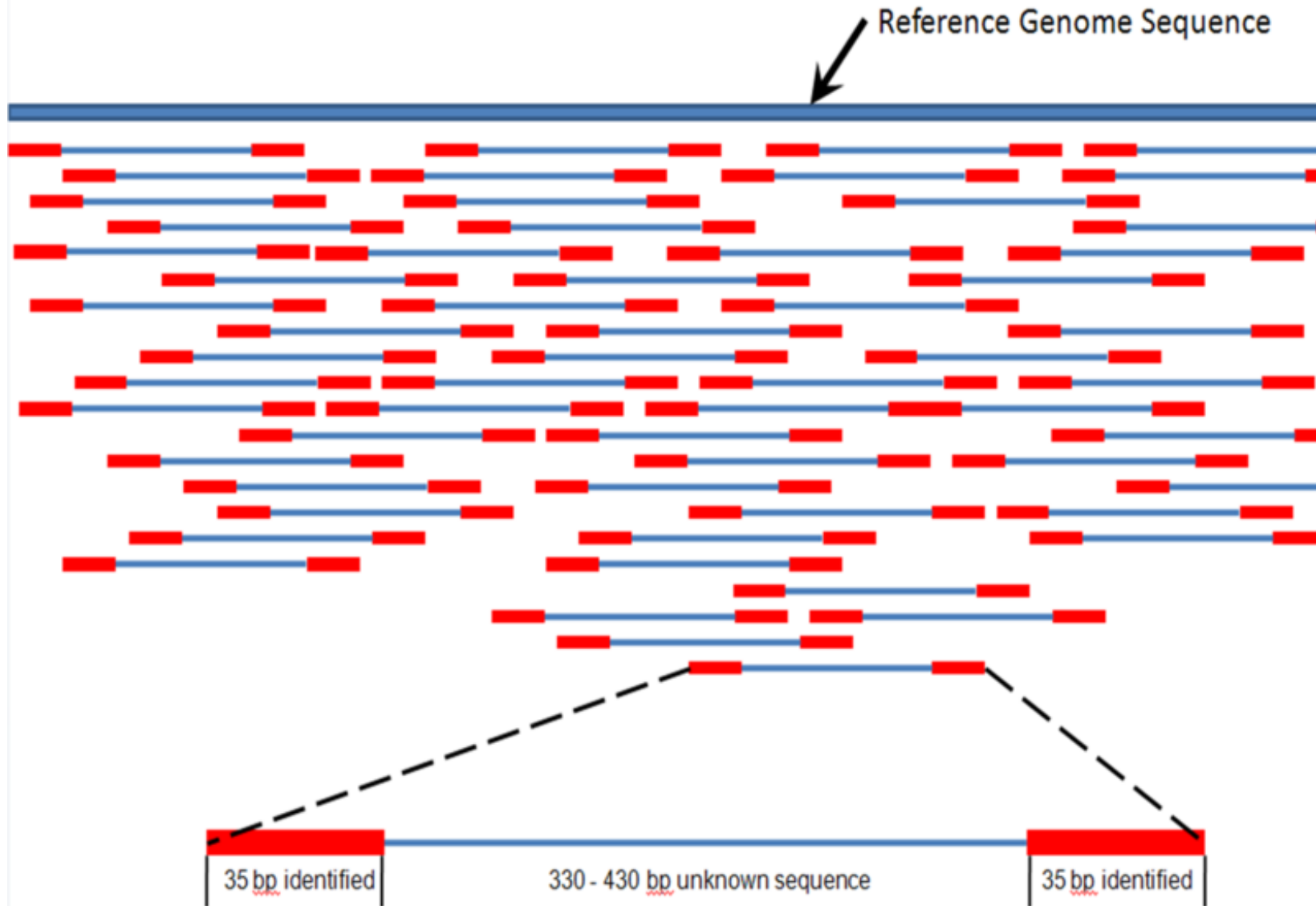


“Base caller” software looks at this cluster across all images and “calls” the complementary nucleotides: **TACAC**, corresponding to the template sequence



TACAC is a “sequence read,” or “read.”
Actual reads are usually 100 or more nucleotides long.

Genome assembly



Mapping reads to reference genome

Naïve method

- Scan whole genome with every read
- Problem: Too slow

Indexing + Alignment approach

- Create a compressed reference 'genome index'
 - a map of where each short subsequence of length 'k' hits the genome
- Map reads using index via smart alignment algorithms and data structures (e.g suffix array)
- Allow for errors: insertions, deletions, mismatches in alignments

Run times for indexing alignment

- Indexing human genome ~ 3 hours
- Alignment speed: 2 million 35 bp reads on 1 processor ~20 mins
- Alignment speed depends on error rate

ACGTTACCGAATCGATCAAGTCGA
TAC



Nature Reviews | Genetics

http://www.nature.com/nrg/journal/v14/n5/box/nrg3433_BX2.html

