# Introduction to Molecular Biology
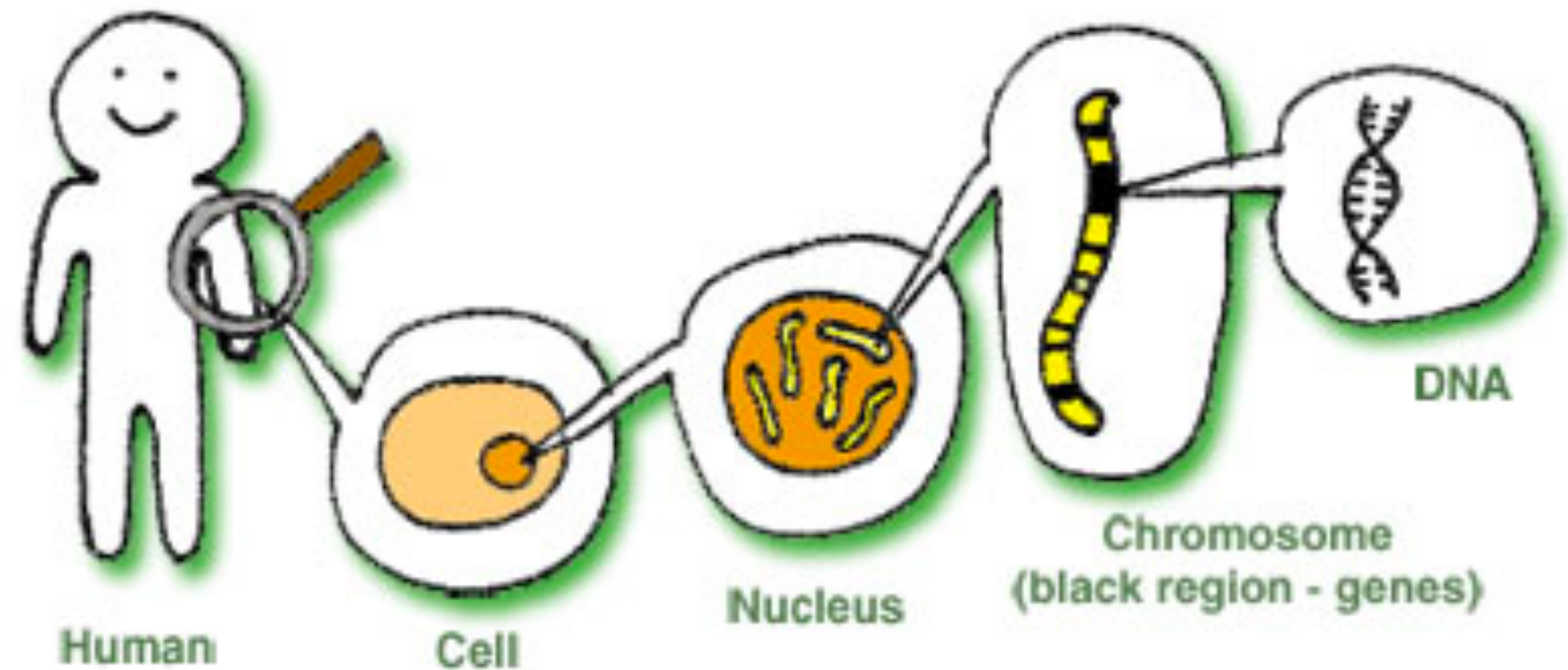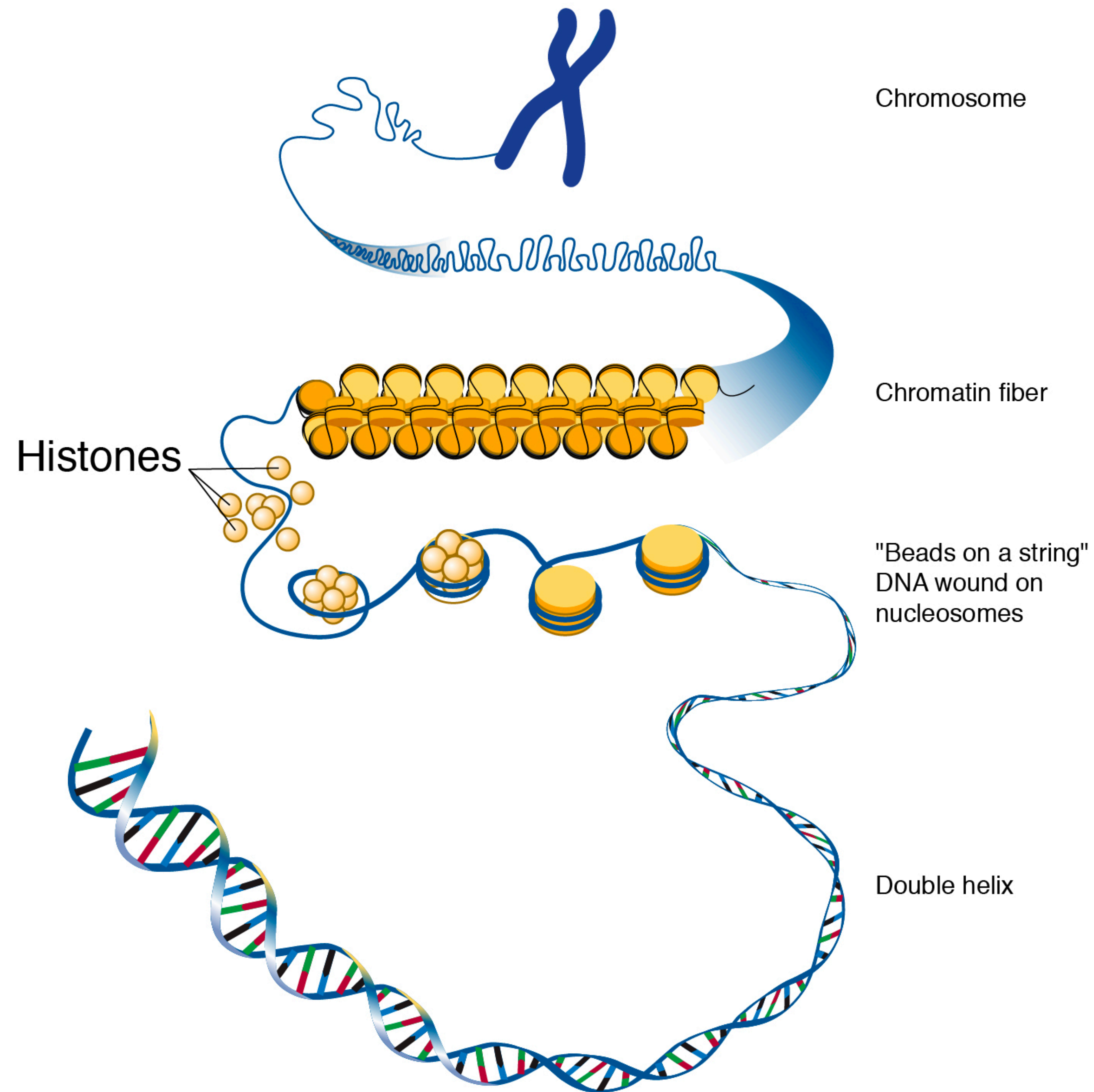
CS 4364 & 5364

# At the highest level

Organism are made up of one or multiple cells

inside the cell is the nucleus, which contains the DNA

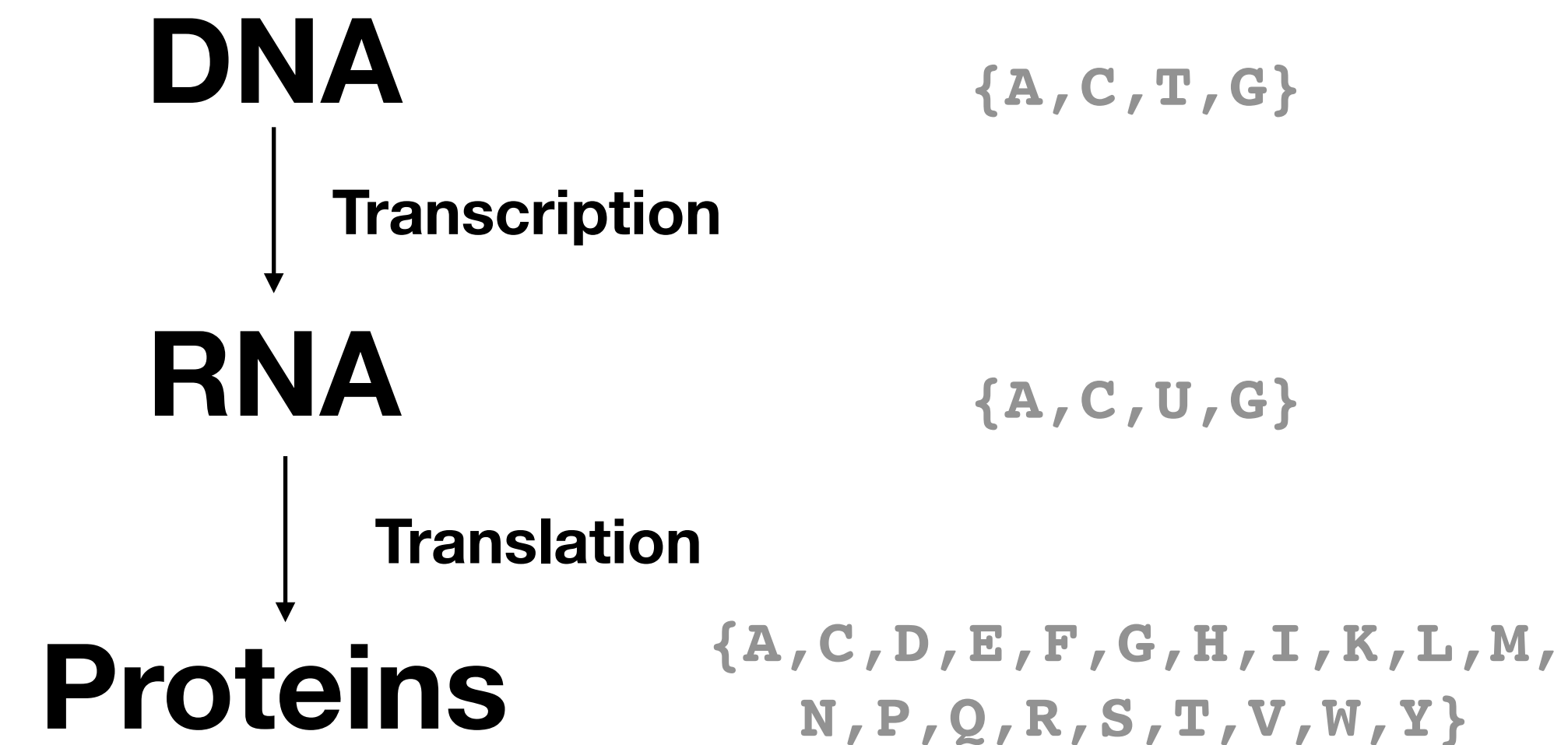humans are *diploid* meaning we have 2 copies of each chromosome (one from each parent)



Human    Cell    Nucleus    Chromosome (black region - genes)    DNA

Chromosome

Chromatin fiber

Histones

"Beads on a string"
DNA wound on
nucleosomes

Double helix

# The Central Dogma

**DNA**

- double stranded
- contains all of the information for "you"
- only about 1.5% of the human genome encodes proteins

**DNA**                    {A,C,T,G}

$\downarrow$ **Transcription**

**RNA**                    {A,C,U,G}

$\downarrow$ **Translation**

**Proteins**        {A,C,D,E,F,G,H,I,K,L,M,
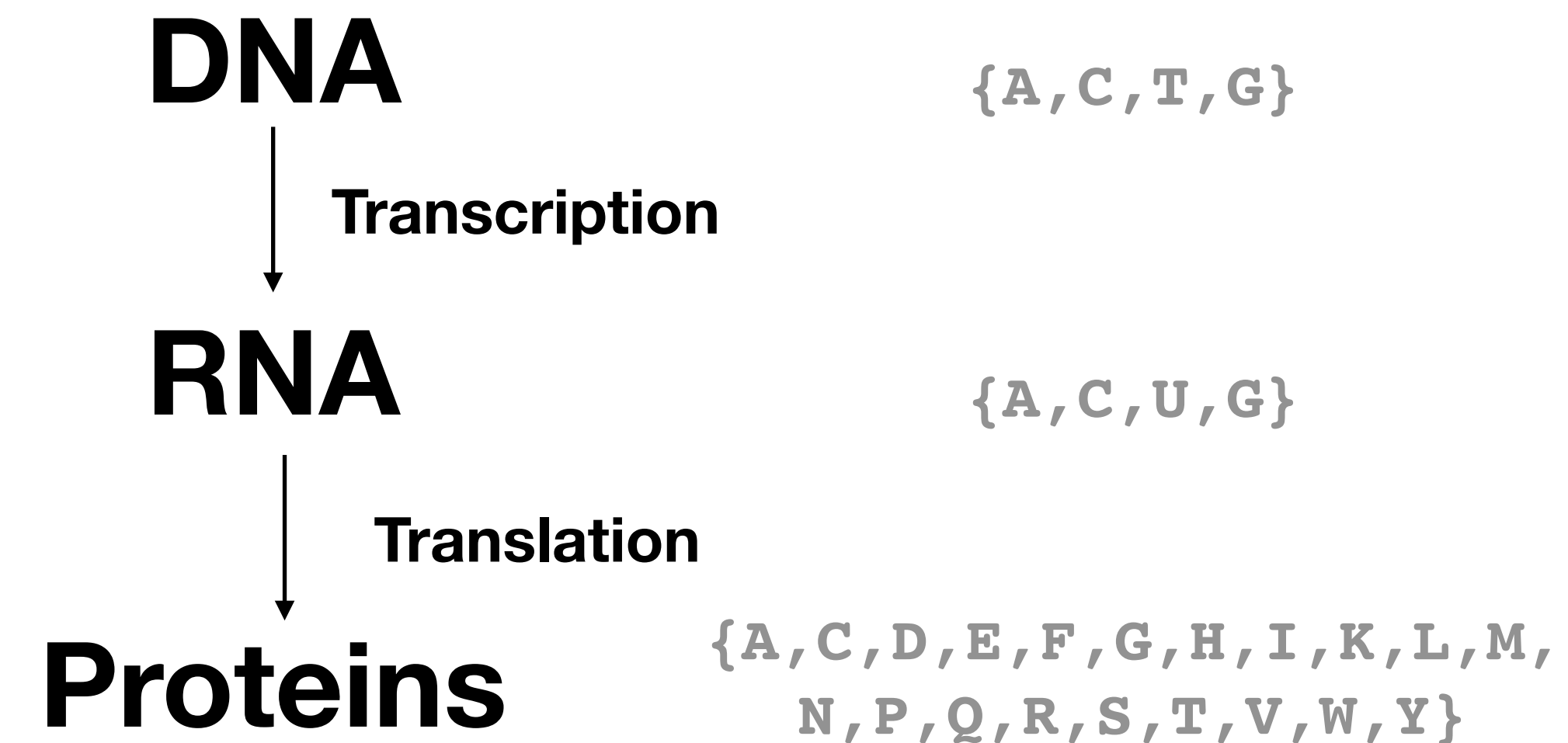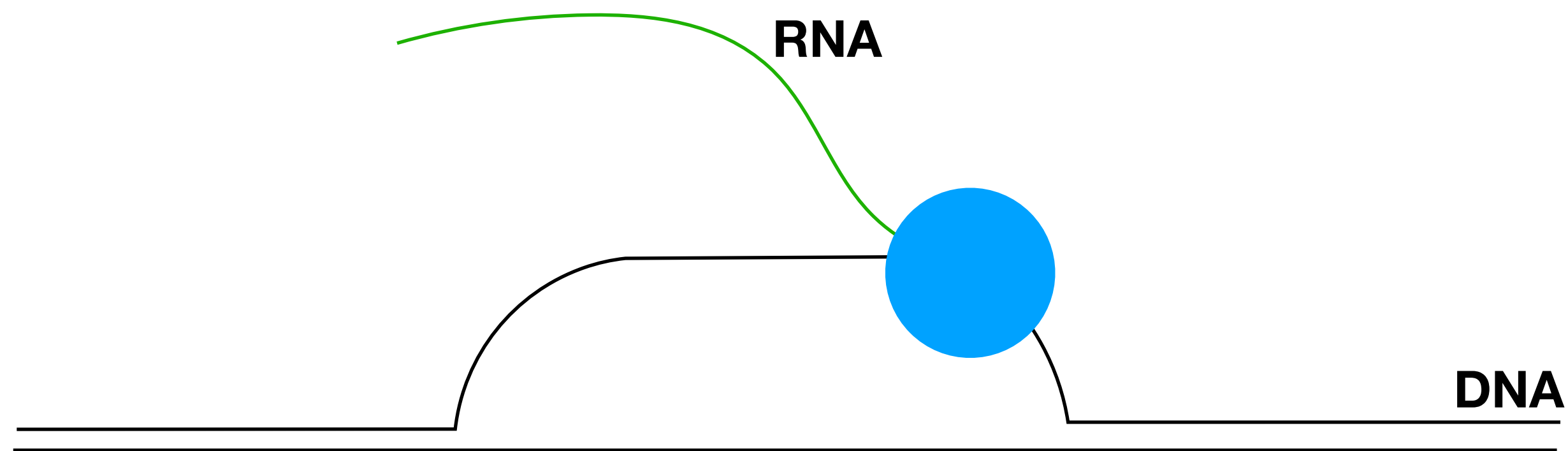                      N,P,Q,R,S,T,V,W,Y}

**DNA**

# The Central Dogma

**Transcription**
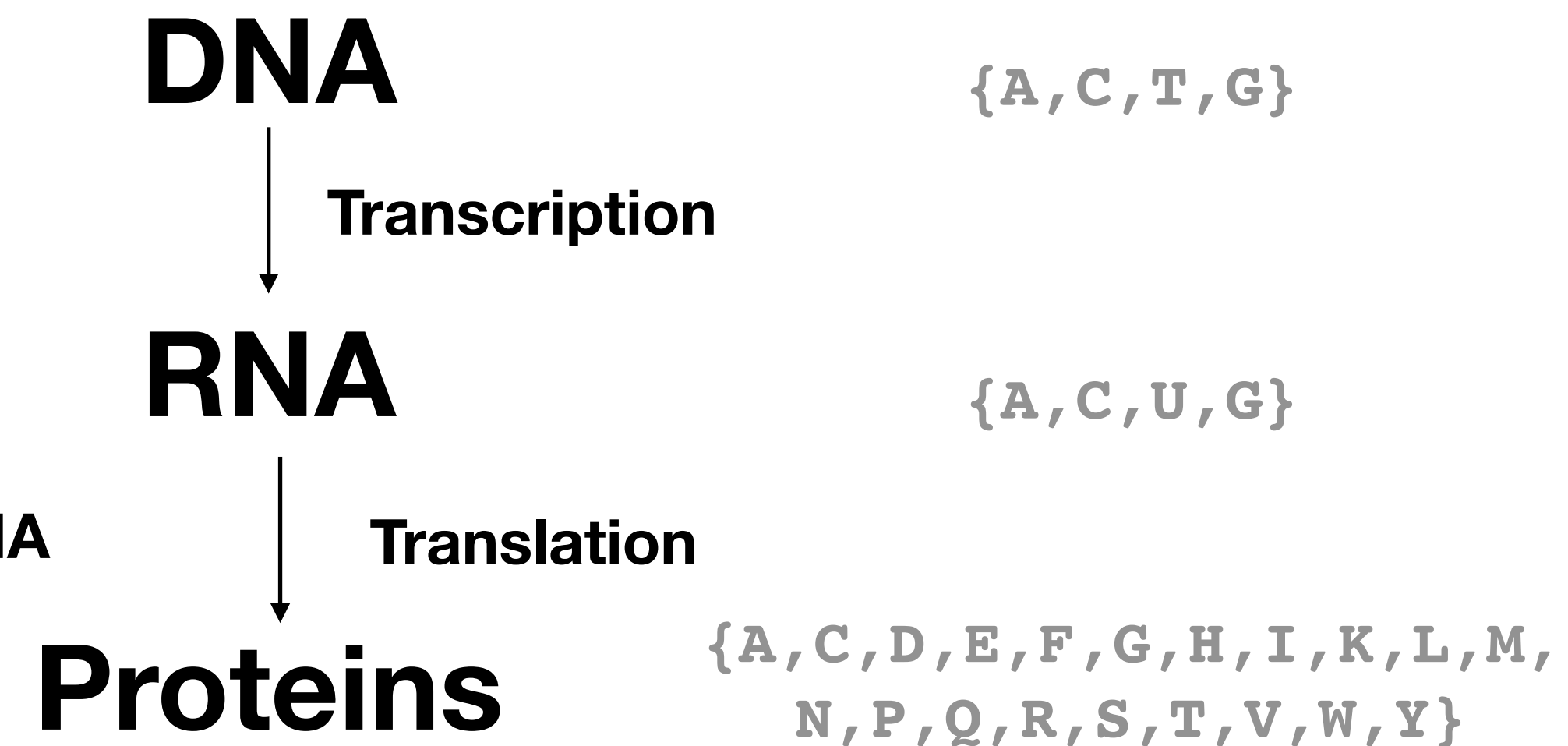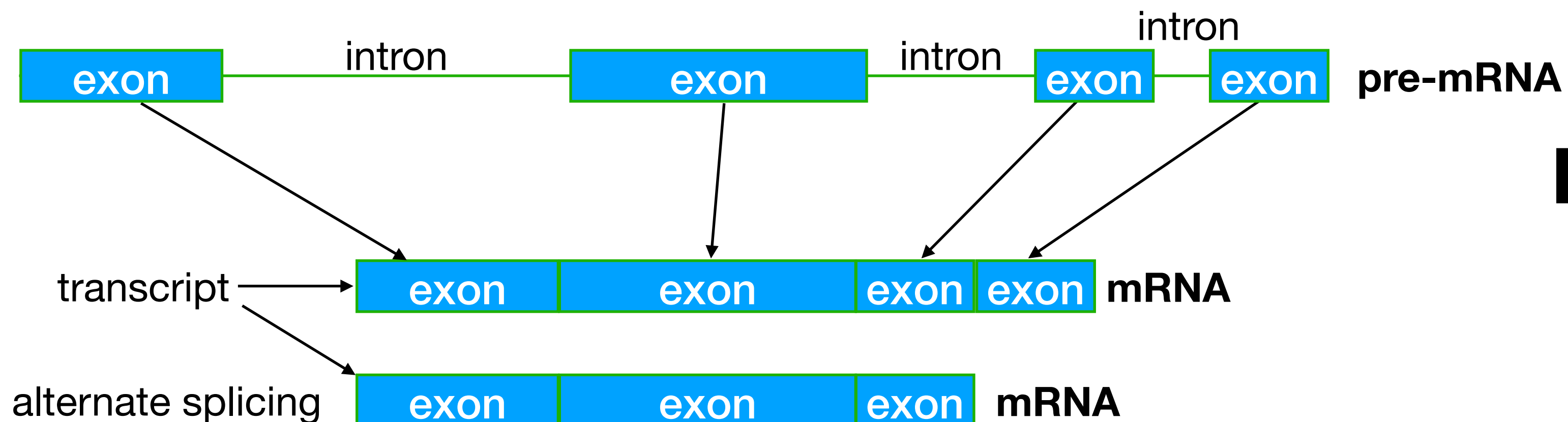
- process of uncoiling, seperating, and copying DNA into RNA
- first stage is called "pre-mRNA" in the case of protein coding genes

RNA

DNA

DNA {A,C,T,G}

Transcription

RNA {A,C,U,G}

Translation

Proteins {A,C,D,E,F,G,H,I,K,L,M, N,P,Q,R,S,T,V,W,Y}

# The Central Dogma

**RNA**

- pre-mRNA undergo splicing to remove the *introns* and leave only (some) *exons*
- some RNA perform functions on their own and are not spliced, called ncRNA (non-coding RNA)

**DNA**          $\{A,C,T,G\}$

Transcription

**RNA**          $\{A,C,U,G\}$

Translation

**Proteins**

$\{A,C,D,E,F,G,H,I,K,L,M, N,P,Q,R,S,T,V,W,Y\}$

exon — intron — exon — intron — exon — intron — exon    **pre-mRNA**

transcript → exon | exon | exon | exon    **mRNA**

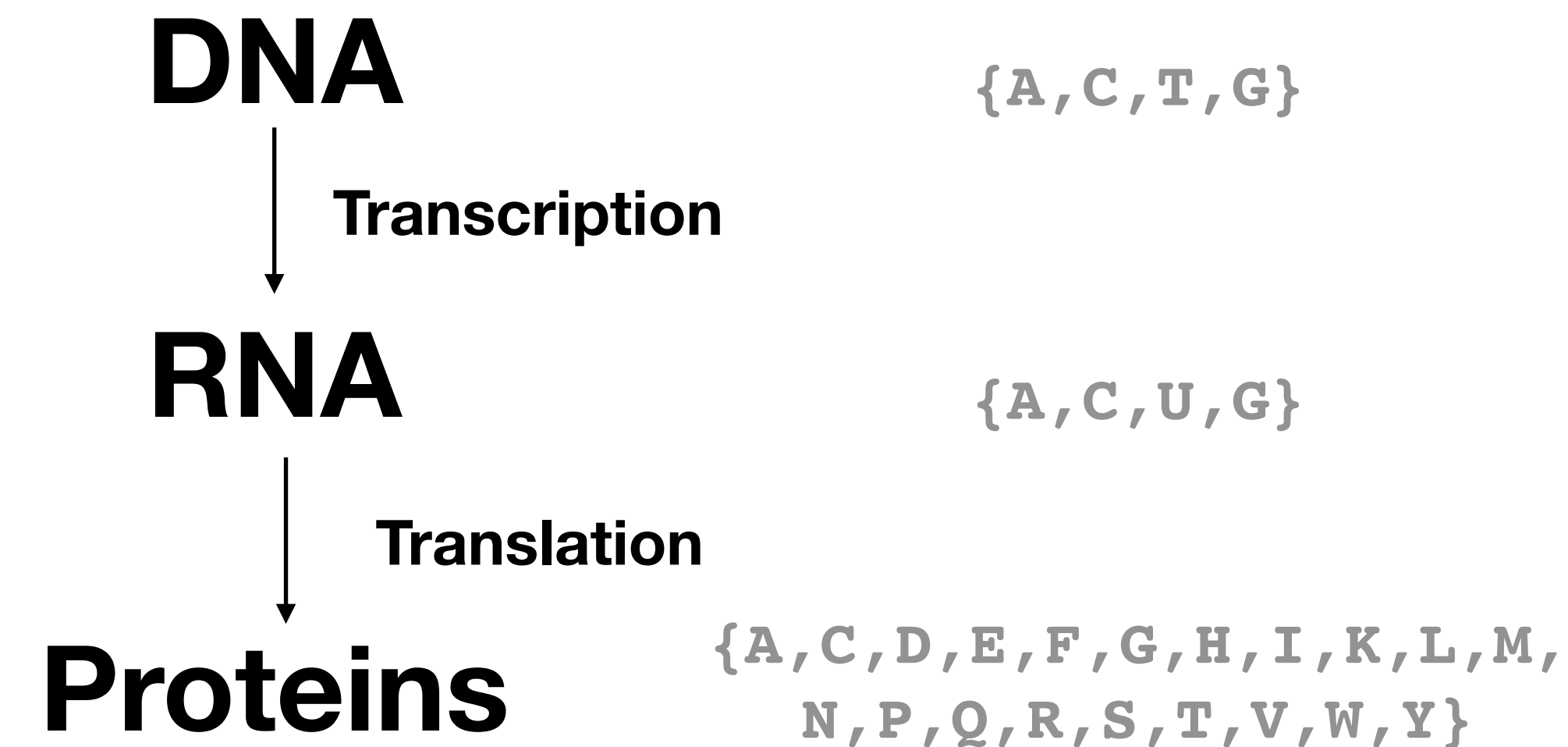alternate splicing → exon | exon | exon    **mRNA**

# The Central Dogma

**Translation**

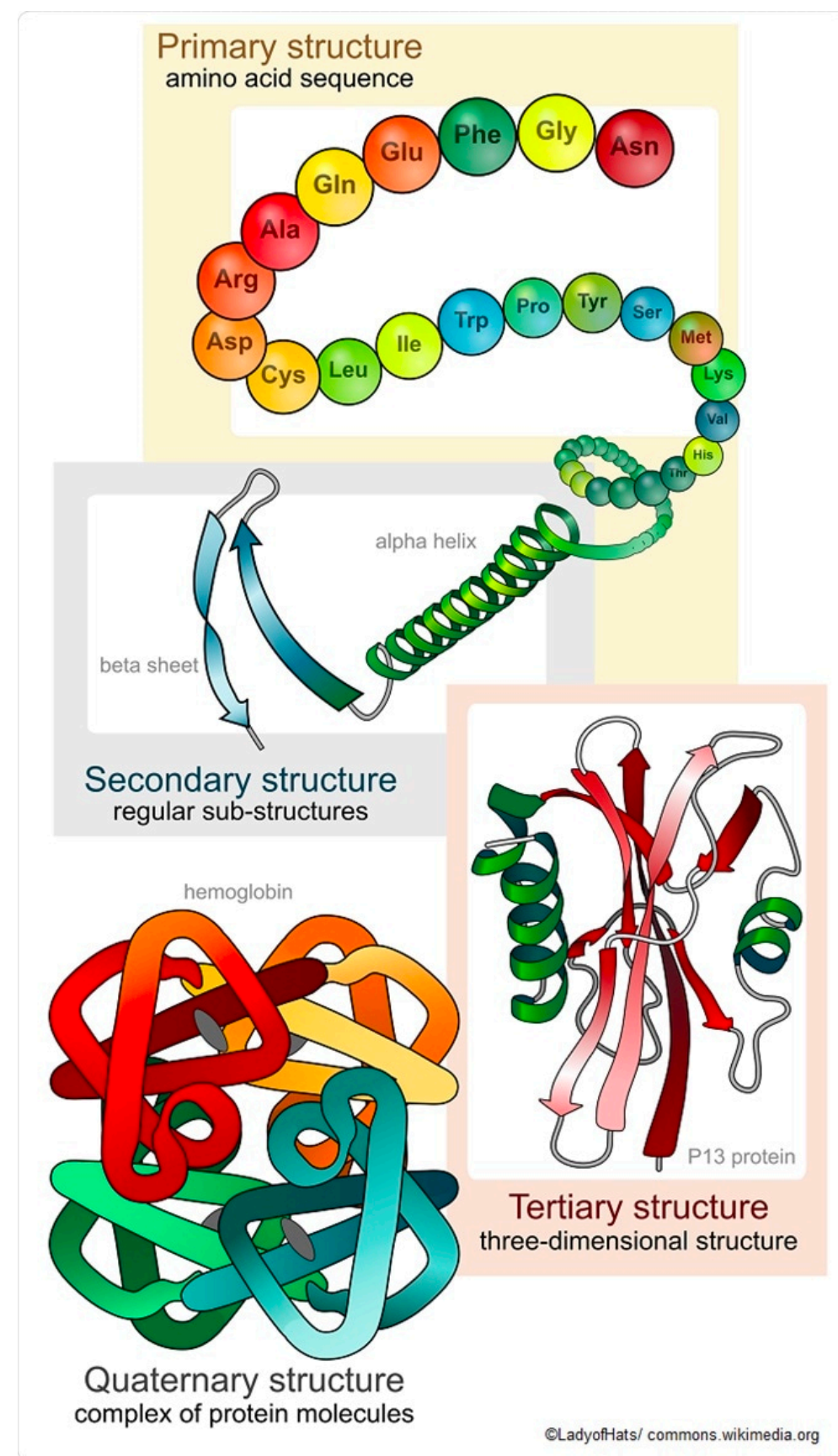- 3-letter groups of RNA characters, *codons,* are converted to amino acids, the building blocks for proteins

Second Character

| First Char. | | A | | C | | U | | G | | Third .Char |
|---|---|---|---|---|---|---|---|---|---|---|
| A | AAC | N | ACC | | AUC | | AGC | S | | C |
| | AAU | | ACU | T | AUU | I | AGU | | | U |
| | AAA | | ACA | | AUA | | AGA | | | A |
| | AAG | K | ACG | | AUG | M/start | AGG | R | | G |
| C | CAC | H | CCC | | CUC | | CGC | | | C |
| | CAU | | CCU | P | CUU | L | CGU | | | U |
| | CAA | | CCA | | CUA | | CGA | R | | A |
| | CAG | Q | CCG | | CUG | | CGG | | | G |
| U | UAC | Y | UCC | | UUC | F | UGC | C | | C |
| | UAU | | UCU | S | UUU | | UGU | | | U |
| | UAA | stop | UCA | | UUA | | UGA | stop | | A |
| | UAG | | UCG | | UUG | L | UGG | W | | G |
| G | GAC | D | GCC | | GUC | | GGC | | | C |
| | GAU | | GCU | A | GUU | V | GGU | | | U |
| | GAA | | GCA | | GUA | | GGA | G | | A |
| | GAG | E | GCG | | GUG | | GGG | | | G |

**DNA**   {A,C,T,G}

**Transcription**

**RNA**   {A,C,U,G}

**Translation**

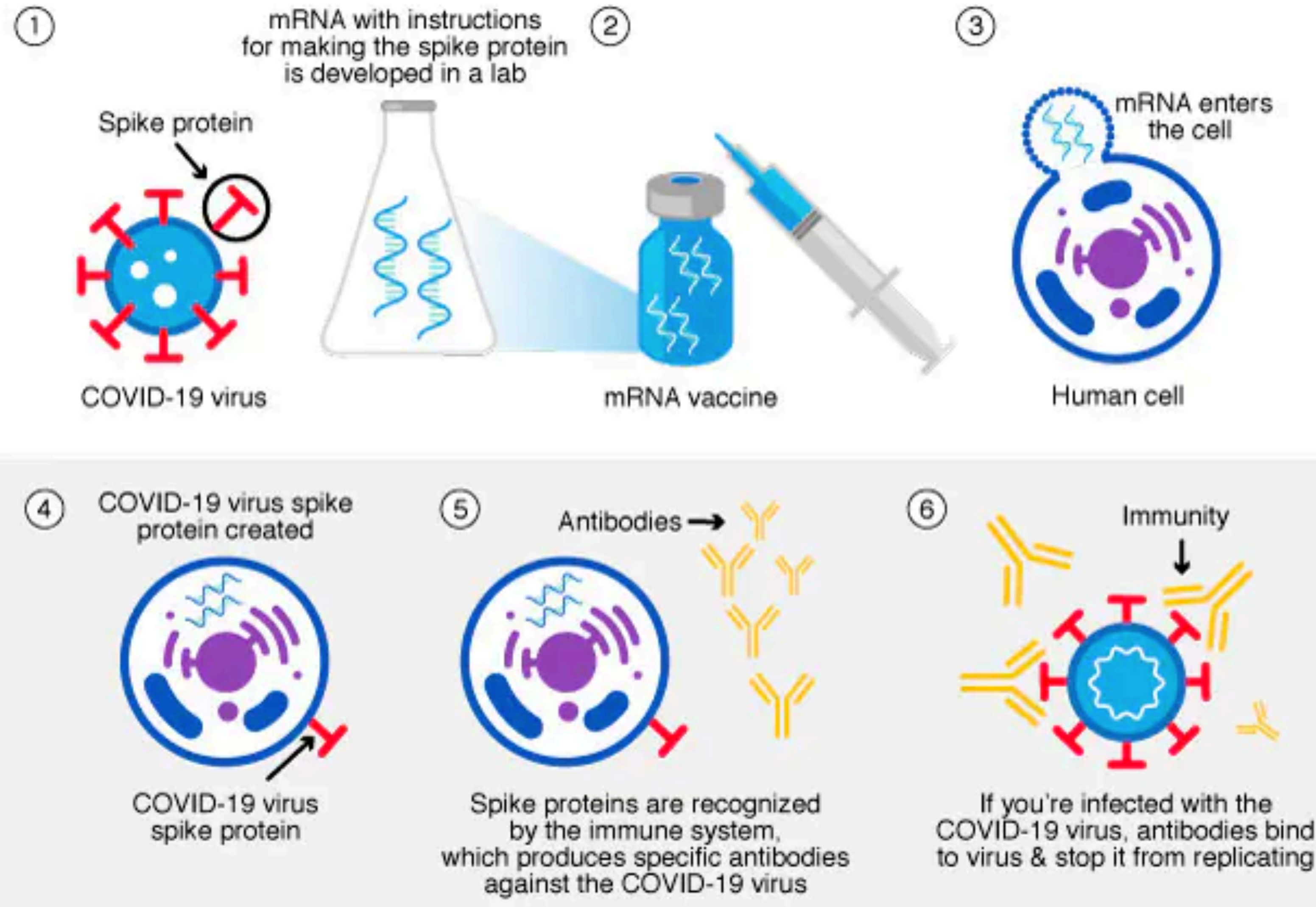**Proteins**   {A,C,D,E,F,G,H,I,K,L,M, N,P,Q,R,S,T,V,W,Y}

# The Central Dogma

**Proteins**

- Do stuff in the cell, including help with translation and transcription



**DNA** {A,C,T,G}

Transcription

**RNA** {A,C,U,G}

Translation

**Proteins** {A,C,D,E,F,G,H,I,K,L,M, N,P,Q,R,S,T,V,W,Y}

# Quick Diversion



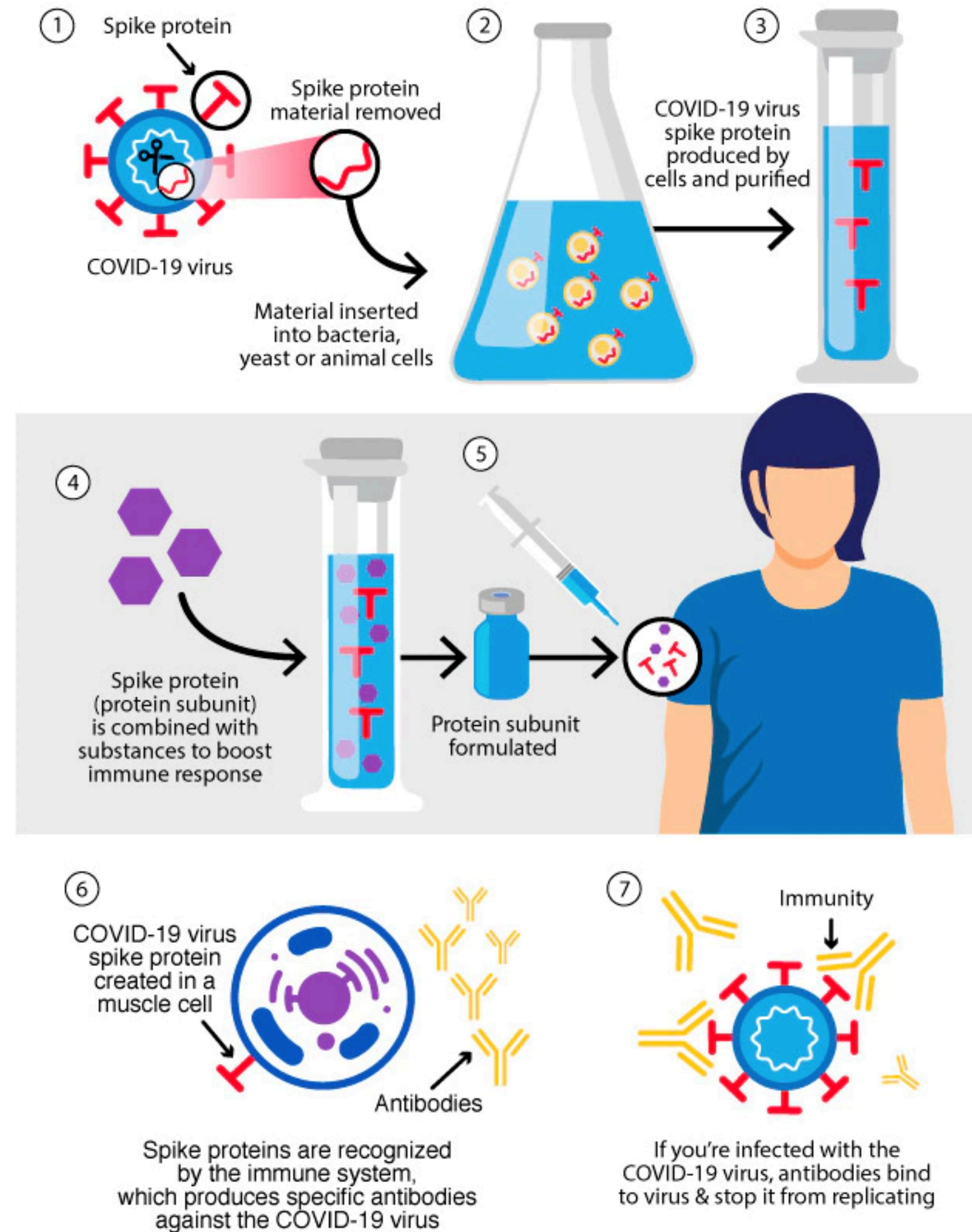① mRNA with instructions for making the spike protein is developed in a lab

Spike protein

COVID-19 virus

mRNA vaccine

② 

③ mRNA enters the cell

Human cell

④ COVID-19 virus spike protein created

COVID-19 virus spike protein

⑤ Antibodies →

Spike proteins are recognized by the immune system, which produces specific antibodies against the COVID-19 virus

⑥ Immunity

If you're infected with the COVID-19 virus, antibodies bind to virus & stop it from replicating

# Quick Diversion



① Spike protein → Spike protein
COVID-19 virus
Spike protein genetic material extracted
Genetic material inserted into inactive (harmless) virus
Unrelated, harmless (vector) virus

② Viral vector vaccine

③ Harmless virus enters cell
Human cell

④ COVID-19 virus spike protein created
COVID-19 virus spike protein

⑤ Antibodies →
Spike proteins are recognized by the immune system, which produces specific antibodies against the COVID-19 virus

⑥ Immunity
If you're infected with the COVID-19 virus, antibodies bind to virus & stop it from replicating

# Quick Diversion



① Spike protein

Spike protein
Spike protein material removed

COVID-19 virus

Material inserted into bacteria, yeast or animal cells

② COVID-19 virus spike protein produced by cells and purified

③

④ Spike protein (protein subunit) is combined with substances to boost immune response

⑤ Protein subunit formulated

⑥ COVID-19 virus spike protein created in a muscle cell

Antibodies

Spike proteins are recognized by the immune system, which produces specific antibodies against the COVID-19 virus

⑦ Immunity

If you're infected with the COVID-19 virus, antibodies bind to virus & stop it from replicating

# Genetic Variants

When copying a genome "errors" may occur, these changes are what make people different

- 99.99% of our genomes are identical
- **Single Nucleotide Polymorphism (SNP)** -- a change at a single base
- **Structural Variants (SV)** -- large scale changes

TACACCGTA**C**GATCG
copying
TACACCGTA**A**GATCG
**SNP**

**SVs**

TACA**CCGTAC**GATCG
copying
TACA**CATGCC**GATCG
inversion

TACA**CCGTAC**GATCG
copying
TACAGAT**CCGTAC**CG
translocation

TACA**CCGTAC**GATCG
copying
TACA**CCGTAC**GAT**CCGTAC**CG
duplication

TACA**CCGTAC**GATCG
copying
TACAGATCG
deletion

# Genetic Variants

- **Deleterious Mutations** -- changes that are harmful (lethal) to a cell
- **Germline Mutations** -- changes passed to offspring
- **Somatic Mutations** -- those not passed down
- **Heterozygous** -- different beween copies
- **Homozygous** -- same on both copies
- **Allele** -- specific position on a chromosome



```
...TACACCGTACGATCG...
...ATGTGGCATGCTAGC...   Copy 1


...TACACCGTAAGATCG...
...ATGTGGCATTCTAGC...   Copy 2
```

homozygous allele    heterozygous allele

# Sanger Sequencing

The basis of all modern sequencing.

**ATGTGGCATGCTAGCTAGCCCTACGTATTGCAGGAT**

**TACACCGTACGATCG ATCG**<span style="color:magenta">**G**</span>

**extend one base
at a time**

***primer* sequence
(matches exactly)**

**end with a
"special" base**

# Sanger Sequencing

. . .
TACACCGTACGATCGATCG**G**
TACACCGTACGATCGATC**G**
TACACCGTACGATCGAT**C**
TACACCGTACGATCGA**T**
TACACCGTACGATCG**A**

The basis of all modern sequencing.



**longer sequences move though
the gel more slowly**

# Sanger Sequencing

. . .
TACACCGTACGATCGATCG**G**
TACACCGTACGATCGATC**G**
TACACCGTACGATCGAT**C**
TACACCGTACGATCGA**T**
TACACCGTACGATCG**A**

The basis of all modern sequencing.

A  T  C  G  G  A  C  C

**Reading Window**

**longer sequences move though the gel more slowly**

**time**

# Second Generation Sequencing

Also called next generation sequencing

Based on the same principles, but at a much larger scale

Improvements were made in the amplification and reading with better microscopes

With this came shorter sequences
- Sanger could do >1,000 bases (characters) at once but all done by hand, so 10s of sequences, very accurate
- Illumina (current standard) ~250 base reads, 1,000,000s of sequences, some errors

# Second Generation Sequencing



from Compeau and Pevzner 2018

# Second Generation Sequencing

NextGen sequencing also introduced paired-end reads
- take a long piece of sequence (much longer than the read size, but predictable size)
- sequence both ends but keep them together
- gives two reads that you know are a certain distance from each other



genome

fixed length genome fragment

read 1          mate-pair distance          read 2

# Third Generation Sequencing

Recently Pacific Biosciences and Oxford Nanopore have introduced new technologies that:
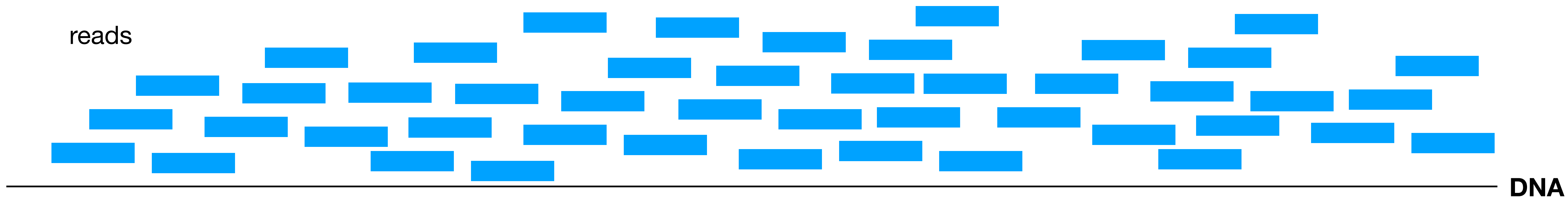- have long reads
- with high(er) error rates

| | Sanger | Next-Generation | Third-Generation |
|---|---|---|---|
| **Launched** | 1977 Basic chemistry<br>1998 Modern form | 2005 with significant improvements since | 2010 with significant improvements since |
| **Estimated Error Rate** | 0.001% - 1% | 0.46% - 2.4% | 11% - 14%<br>(but decreasing) |
| **Cost** | 🪙🪙🪙 | 🪙 | 🪙🪙 |
| **Throughput** | 🧬 | 🧬🧬🧬 | 🧬🧬 |
| **Currently Available Platforms** | Applied Biosystems* | Illumina<br>Ion Torrent*<br>Qiagen (Europe)<br>Complete Genomics (China)** | Pacific Biosciences<br>Oxford Nanopore |
| **Clinical Uses** | Many (but dwindling) | Many (and growing) | Niche uses (today) |

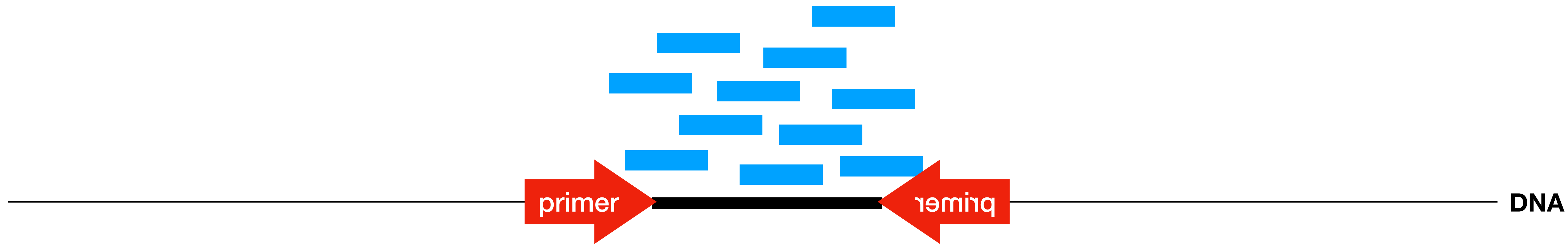*Part of Thermo Fisher      **Part of BGI

# Sequencing Applications

reads

**DNA**

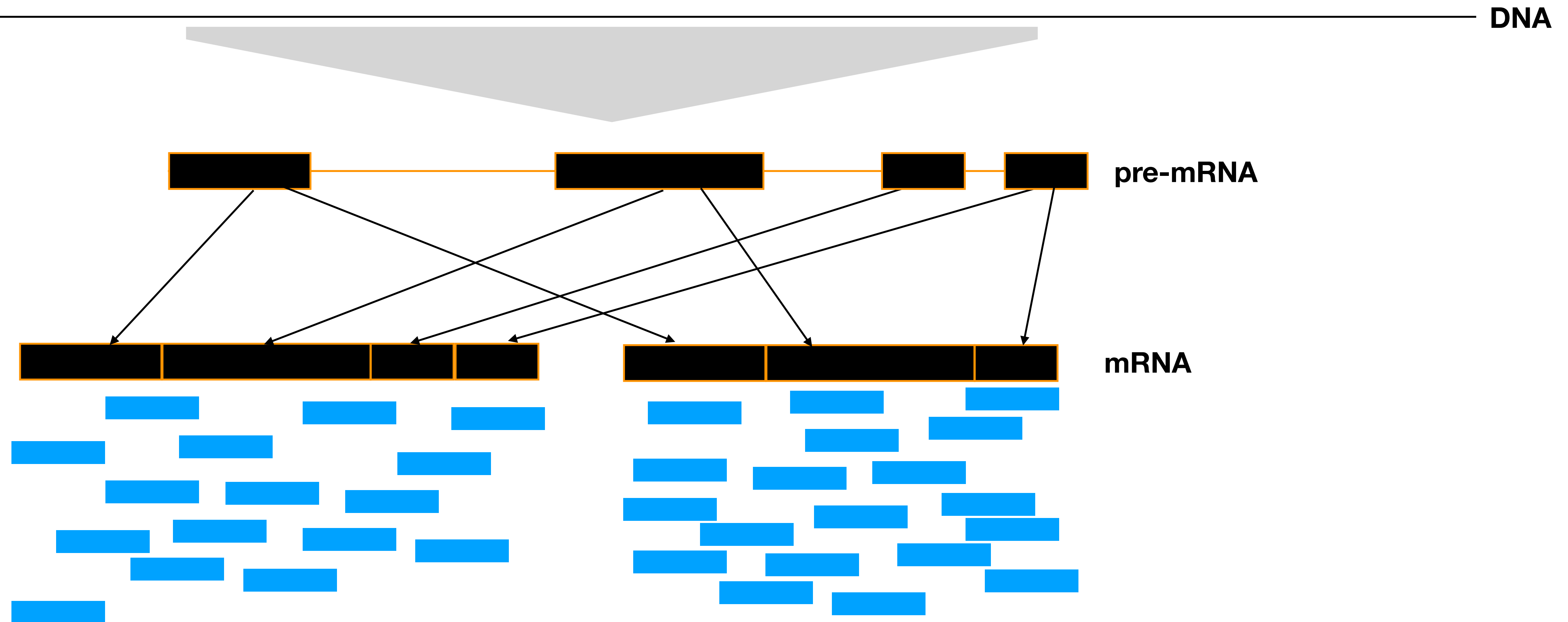whole genome sequencing

# Sequencing Applications



methylation

bisulphite sequencing

DNA

# Sequencing Applications



primer      primer      DNA

targeted sequencing

# Sequencing Applications



DNA

pre-mRNA

mRNA

RNA sequencing

# Sequencing Applications



DNA

binding

chromatin immunoprecipitation (ChIP) sequencing

# History

**1866** -- Gregor Mendel discovers genetics using pea plants

**1869** -- DNA was discovered

**1944** -- Avert and McCarty show DNA carried genetic information

**1953** -- Watson and Crick discovered the 3D structure of DNA

**1961** -- Nirenberg maps DNA to proteins

**1968** -- Discovery of restriction enzymes

**1970s** -- Development of the first sequencing techniques

**1985** -- Development of PCR

**1986** -- Discovery of RNA splicing

**1980-1990** -- Complete sequencing of genomes of small organisms

**1990** -- Launch of the Human Genome Project

**1998** -- Discovery of post-transciption RNA interference

**2000** -- Announcement of the draft human genome

# Major Ongoing Projects

ENCODE (The **Enc**yclopedia **o**f **D**NA **E**lements)
- Effort to identify all functions elements in the human genome

1000 Genomes Project
- Large sample size will hopefully show all (most) of the variation within the population

UK BioBank
- 500,000 UK genomes in great details

SRA (Sequence Read Archive)
- Public repository of all types of sequencing data

GWAS Catalog (Genome Wide Association Studies)
- Multiple studies for many possible purposes (i.e. cancer, disorders, etc.)