

Introduction to Sequencing

CS 4390/5390

Fall 2019

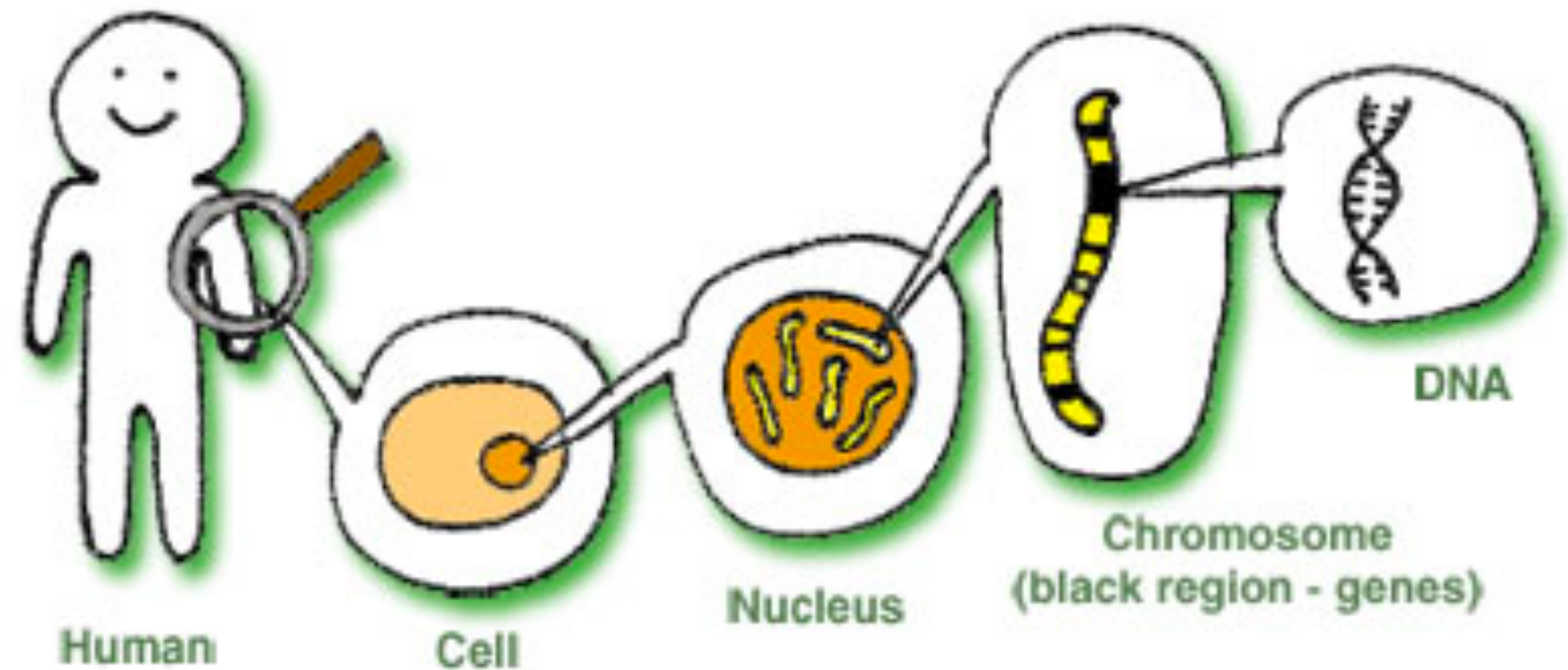
**Associated Reading: Mäkinen, et al. Chapter 1
Gibson and Muse, Chapter 2**

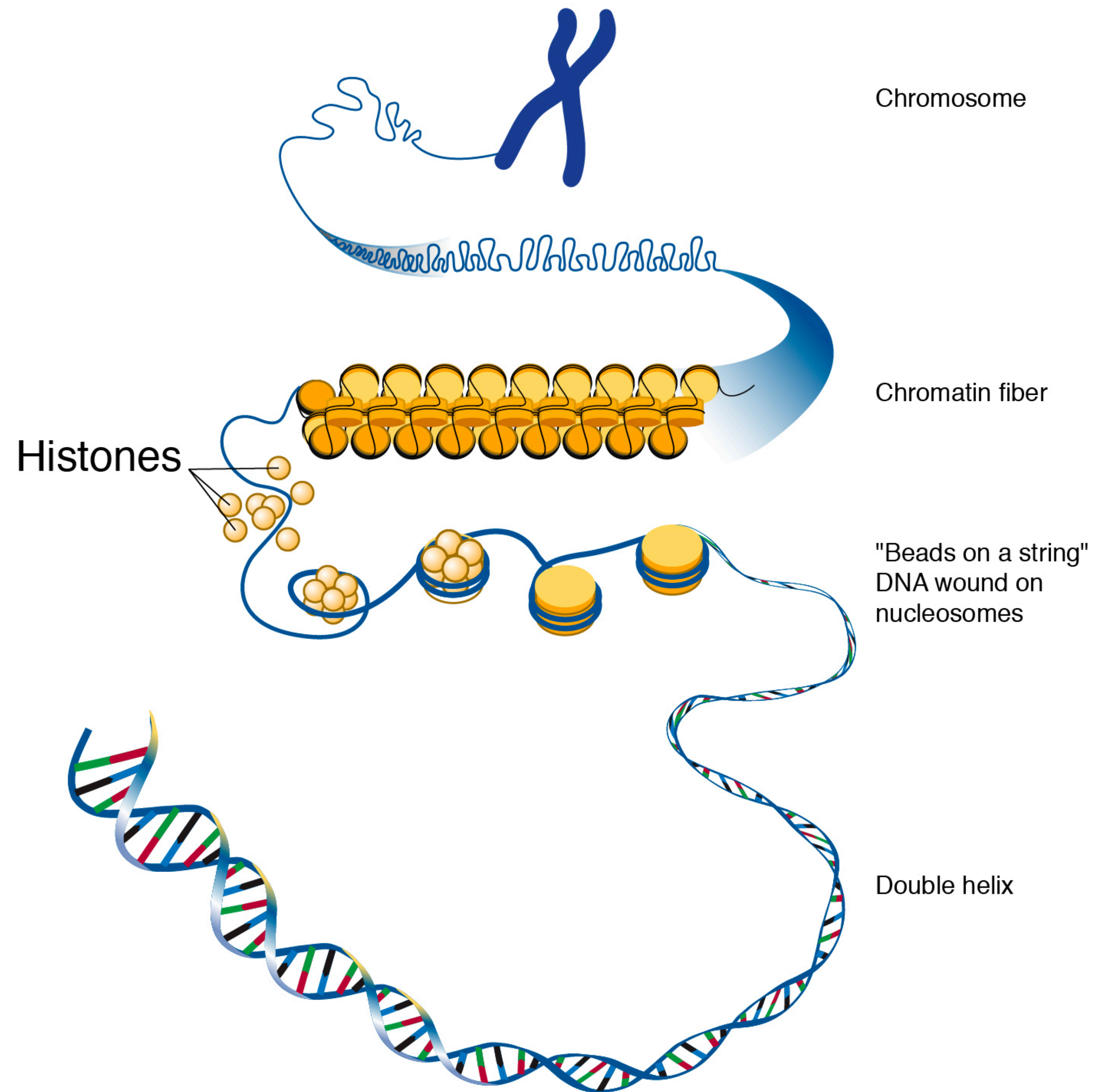
At the highest level

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)





The Central Dogma

DNA

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

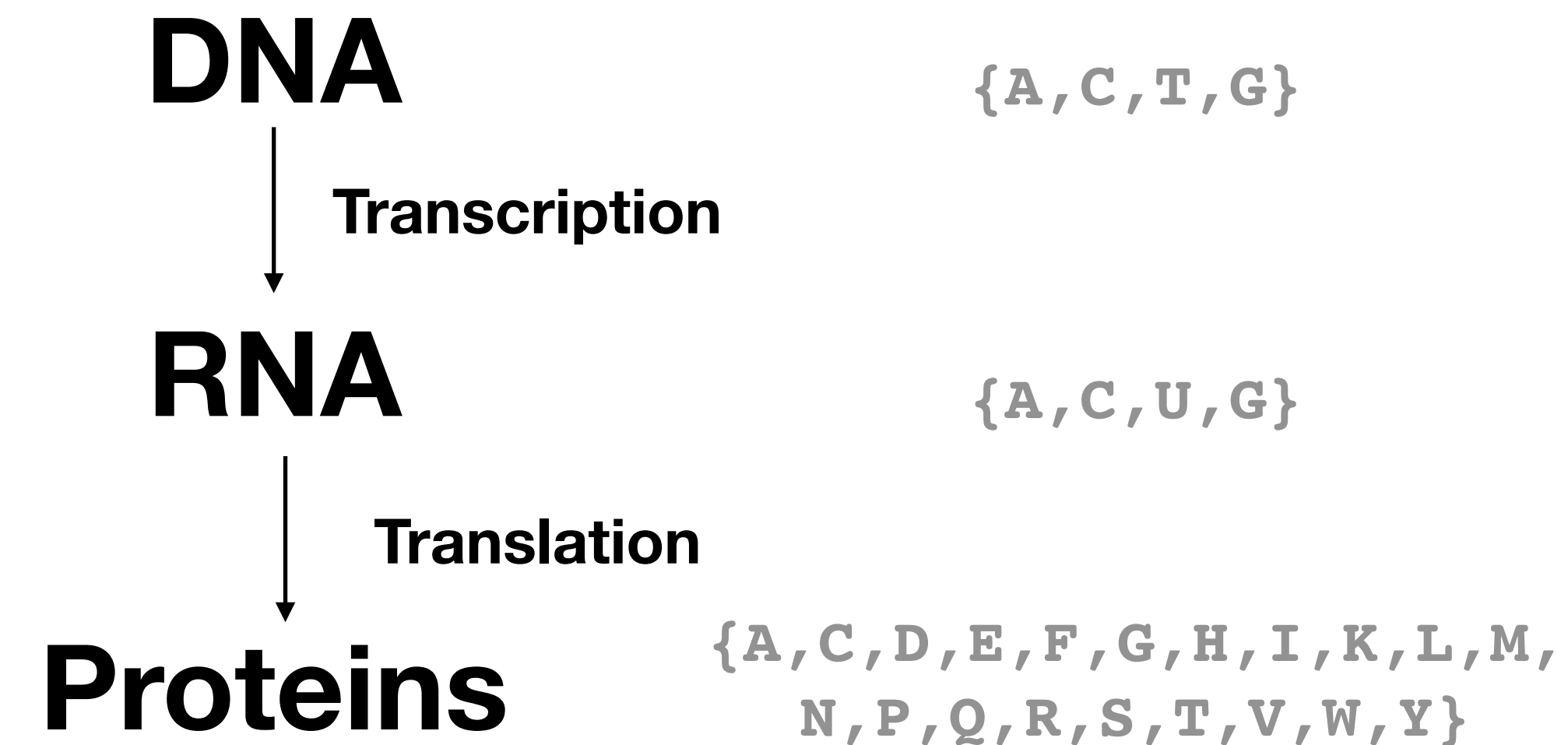
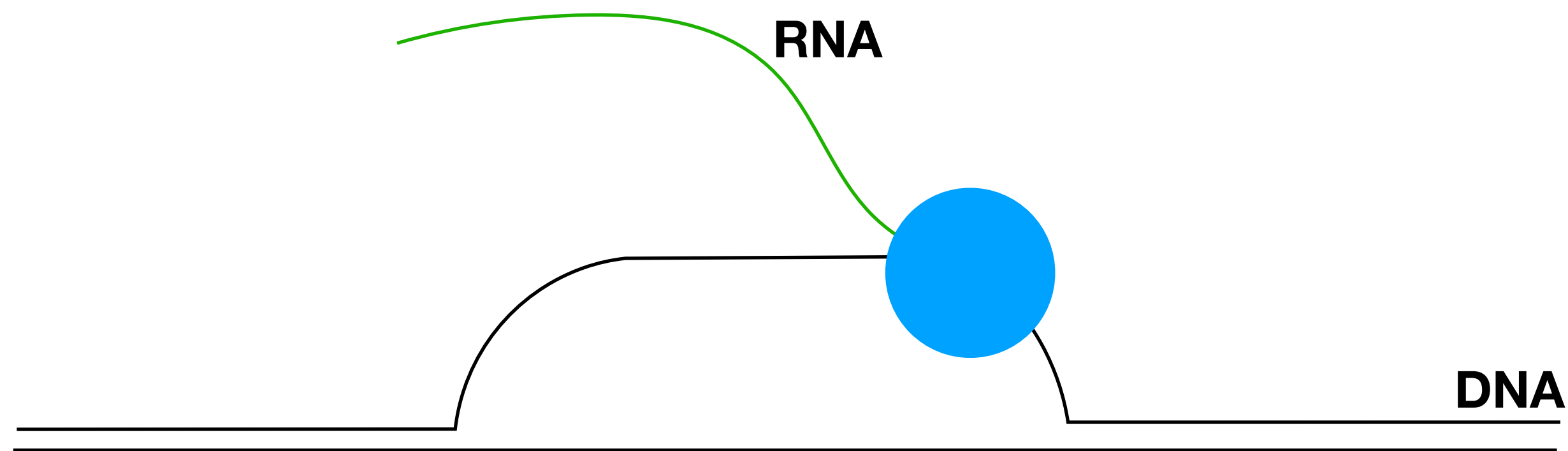


DNA

The Central Dogma

Transcription

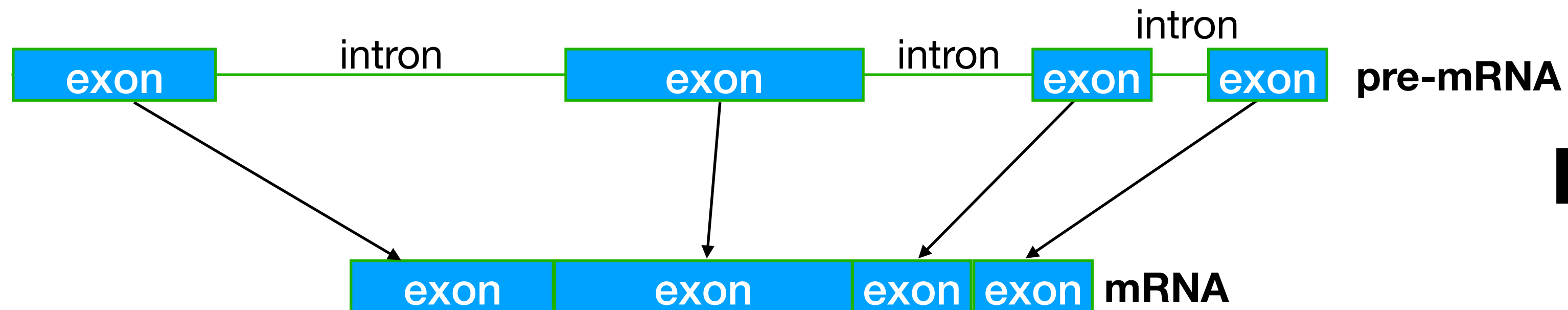
- process of uncoiling, separating, and copying DNA into RNA
- first stage is called "pre-mRNA" in the case of protein coding genes



The Central Dogma

RNA

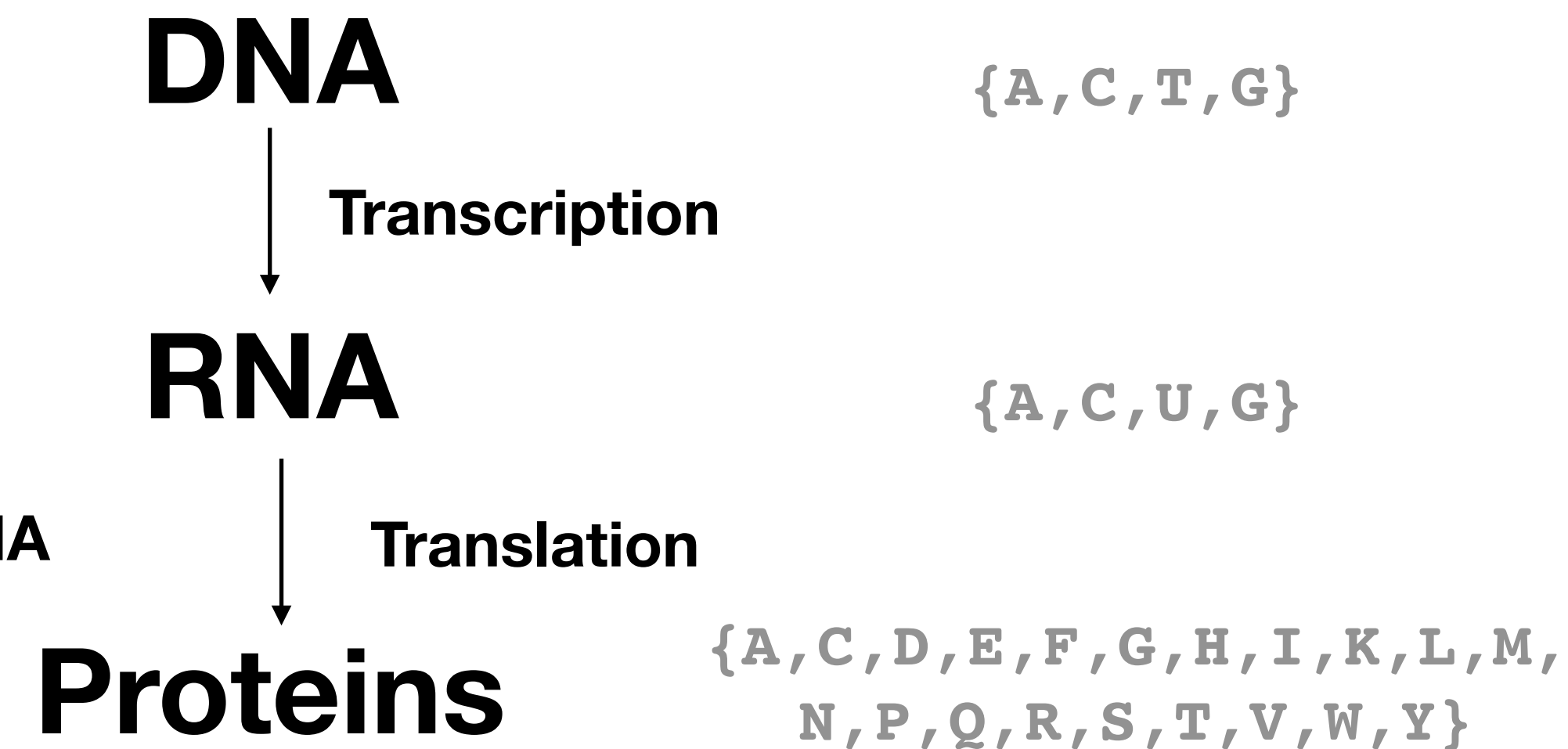
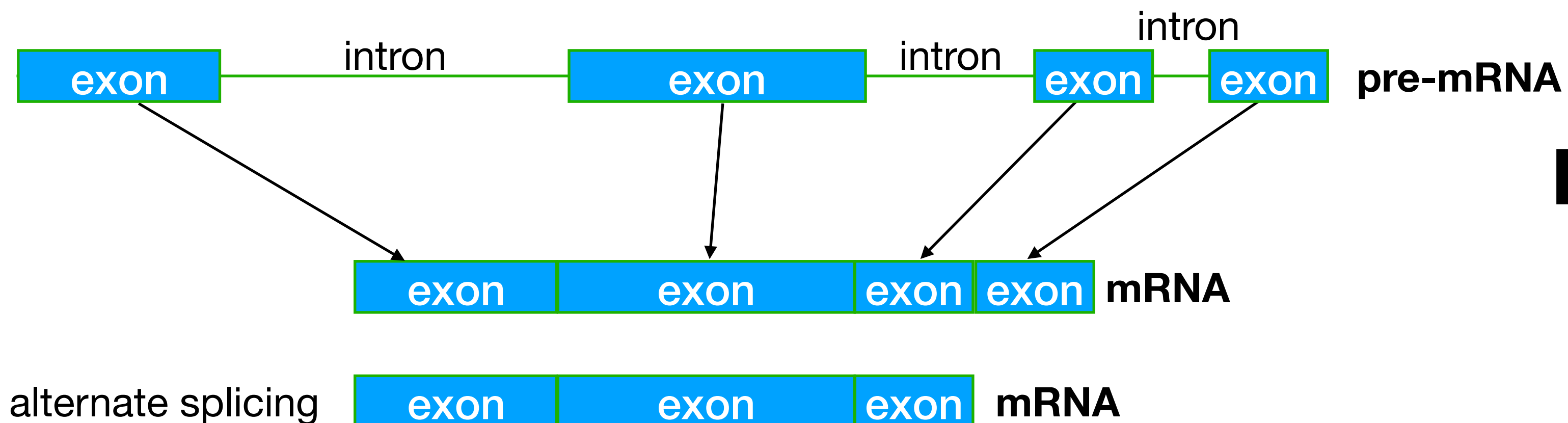
- pre-mRNA undergo splicing to remove the *introns* and leave only (some) *exons*
- some RNA perform functions on their own and are not spliced, called ncRNA (non-coding RNA)



The Central Dogma

RNA

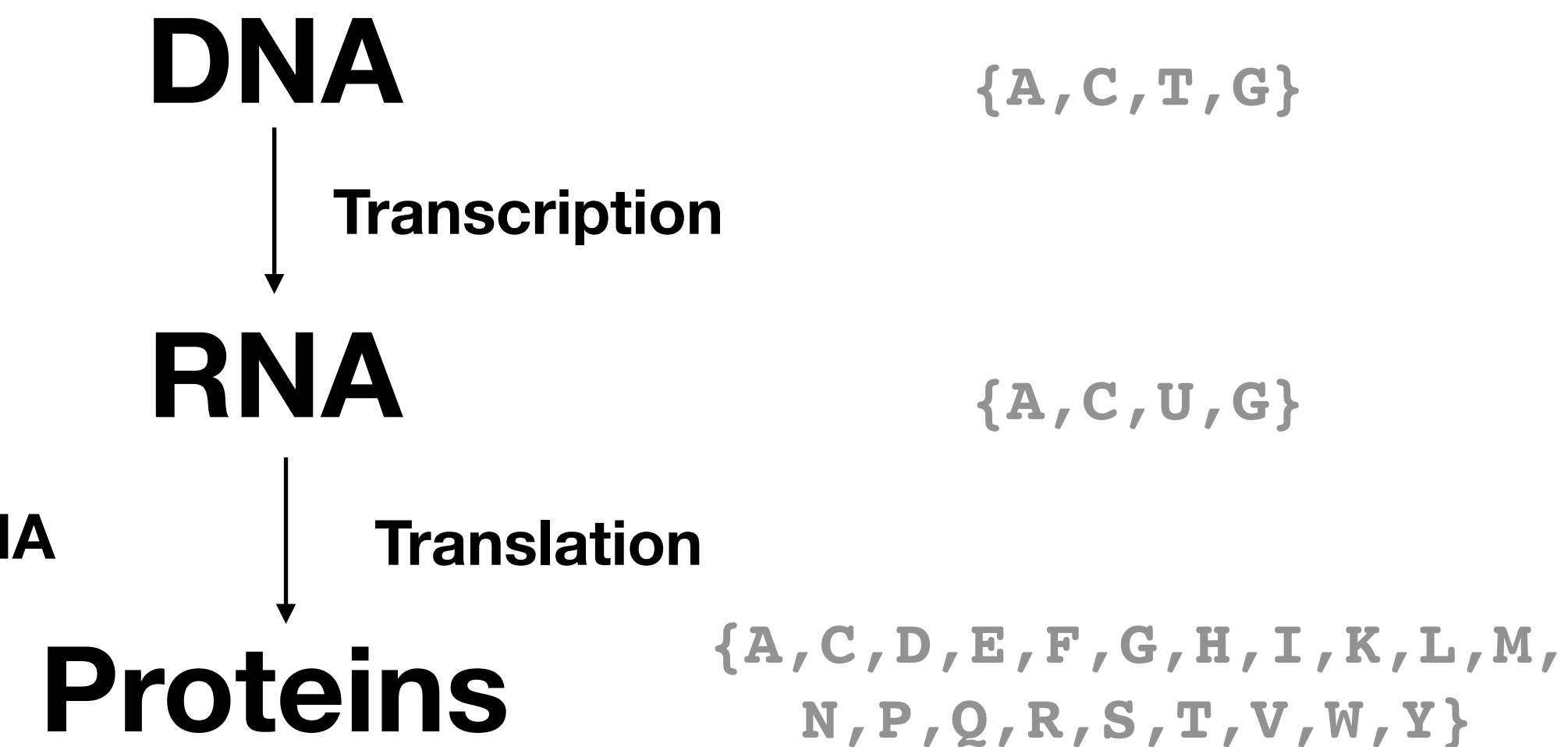
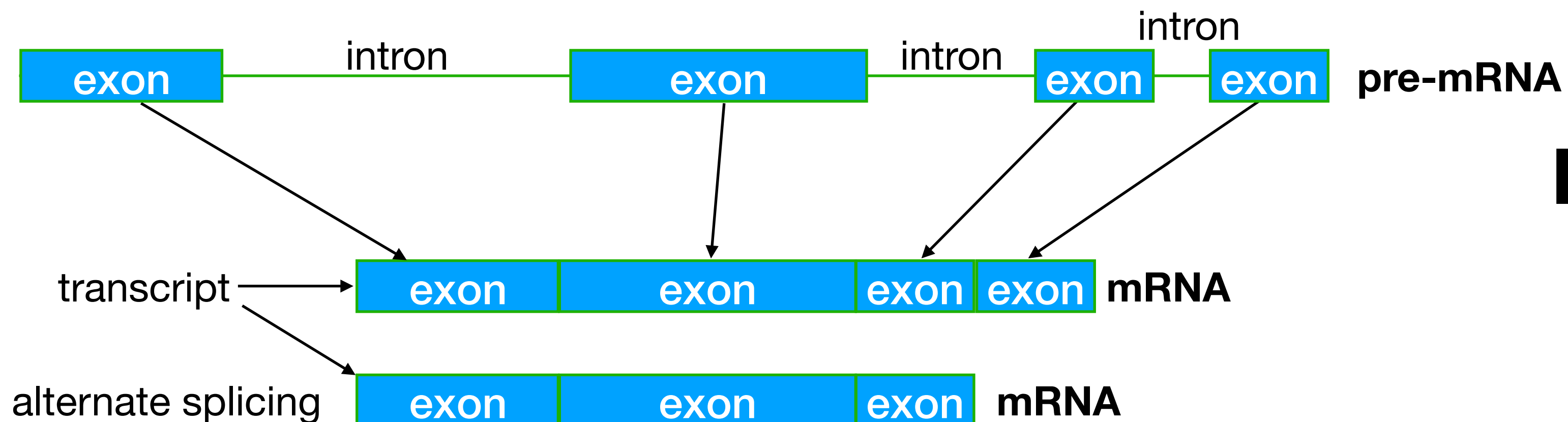
- pre-mRNA undergo splicing to remove the *introns* and leave only (some) *exons*
- some RNA perform functions on their own and are not spliced, called ncRNA (non-coding RNA)



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)

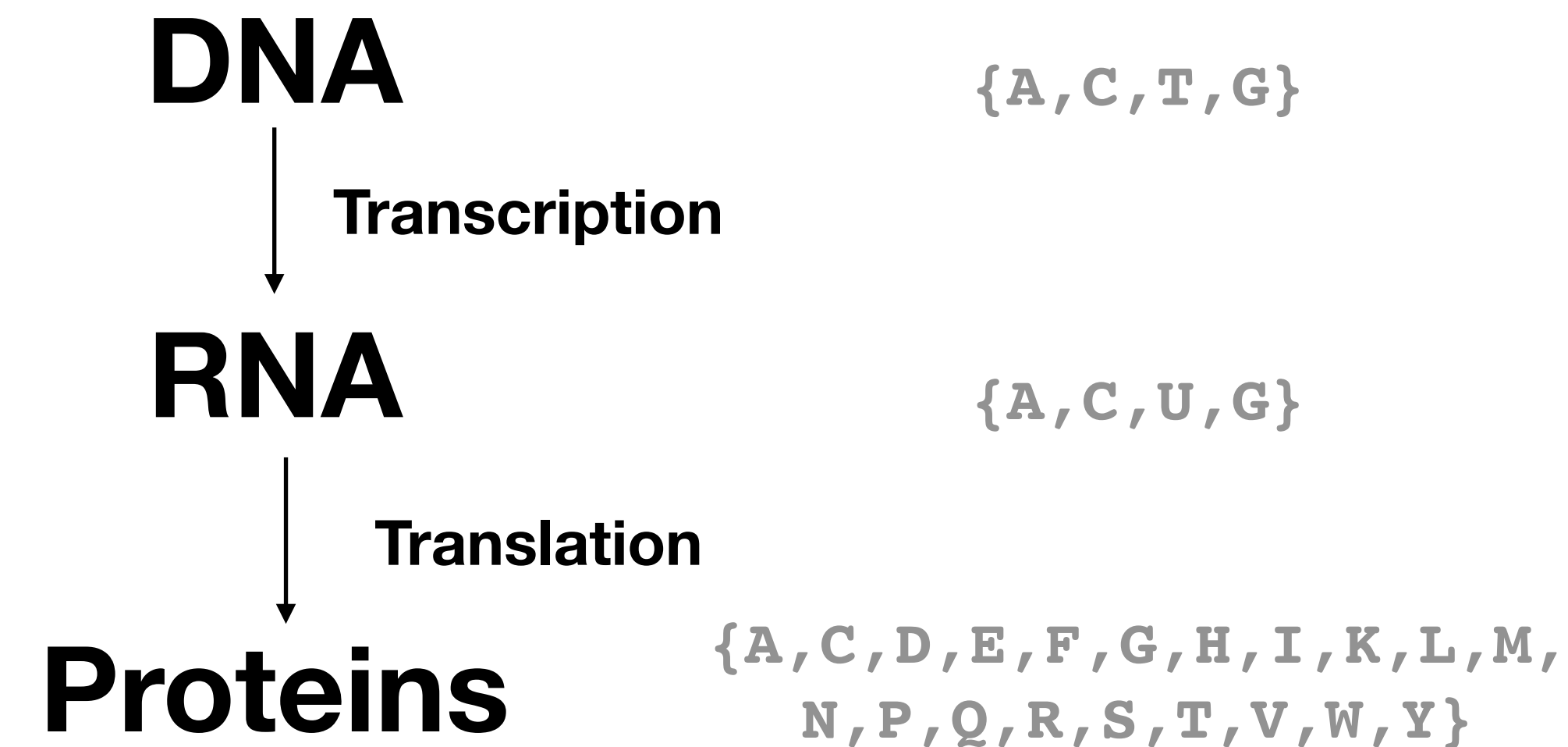


The Central Dogma

Translation

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

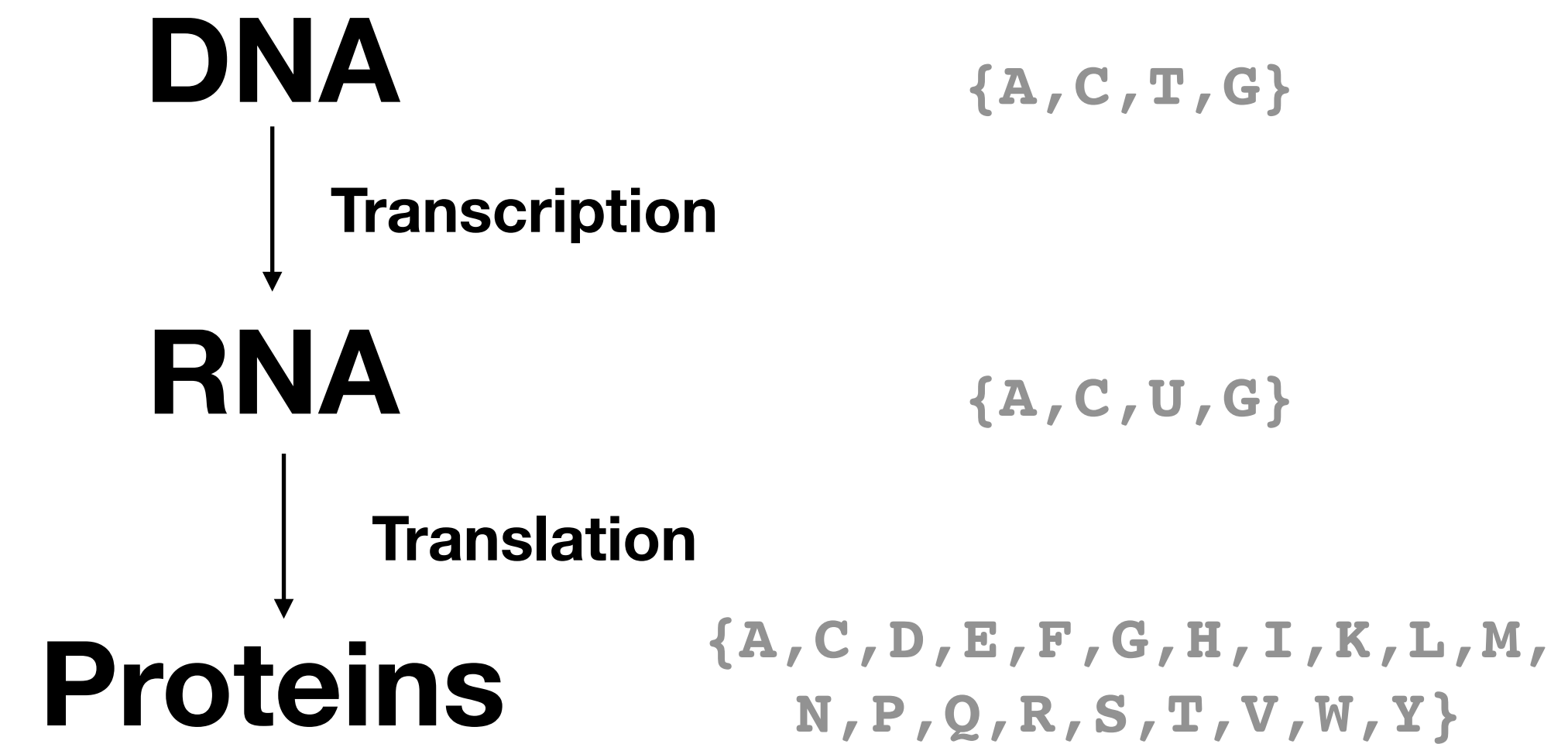
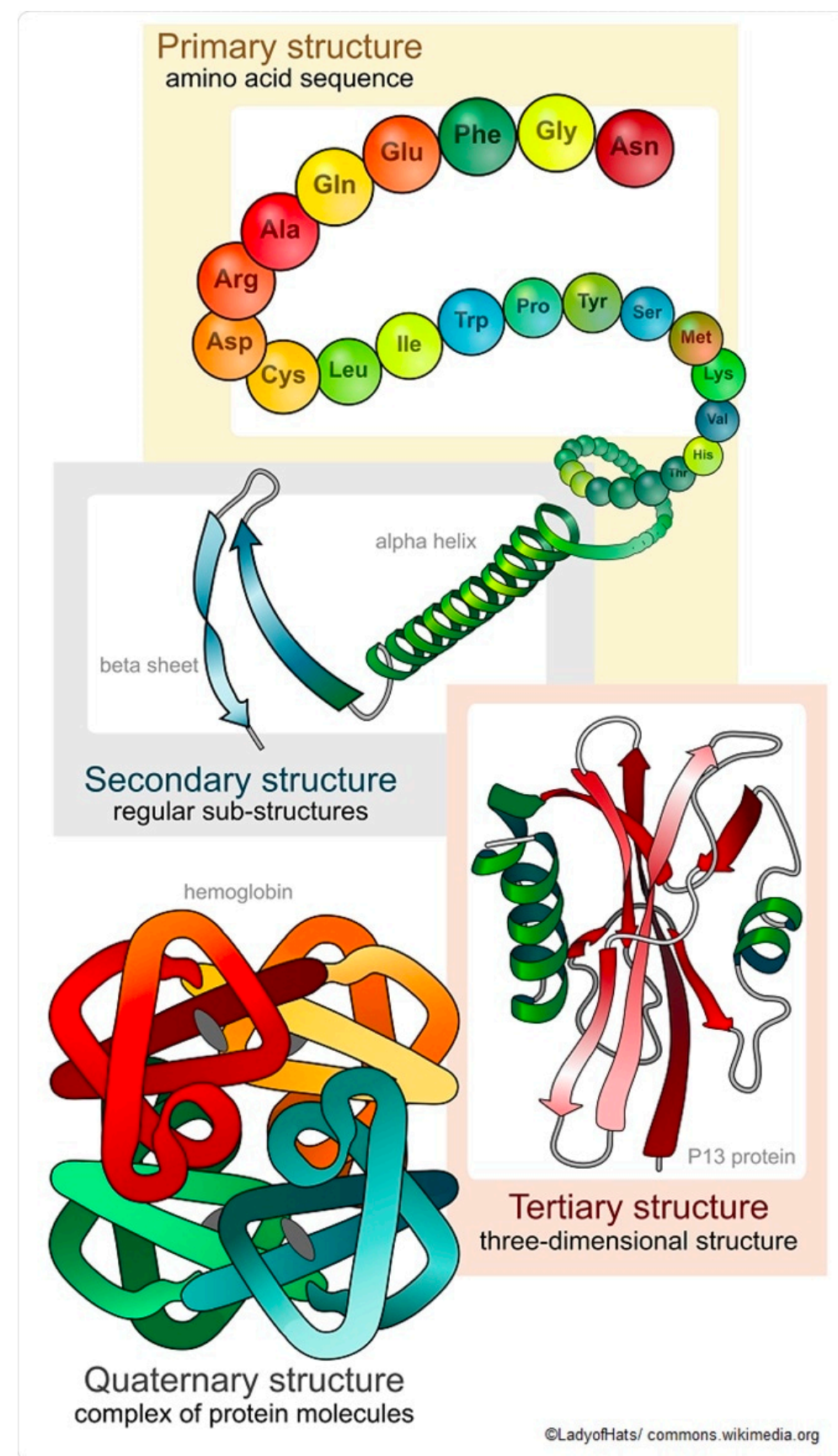
		Second Character								
		A		C		U		G		
First Char.	A	AAC	N	ACC	T	AUC	I	AGC	S	C
		AAU		ACU		AUU		AGU		U
		AAA	K	ACA		AUA	M/start	AGA	R	A
		AAG		ACG		AUG		AGG		G
	C	CAC	H	CCC	P	CUC	L	CGC	R	C
		CAU		CCU		CUU		CGU		U
		CAA	Q	CCA		CUA		CGA		A
		CAG		CCG		CUG		CGG		G
	U	UAC	Y	UCC	S	UUC	F	UGC	C	C
		UAU		UCU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



The Central Dogma

Proteins

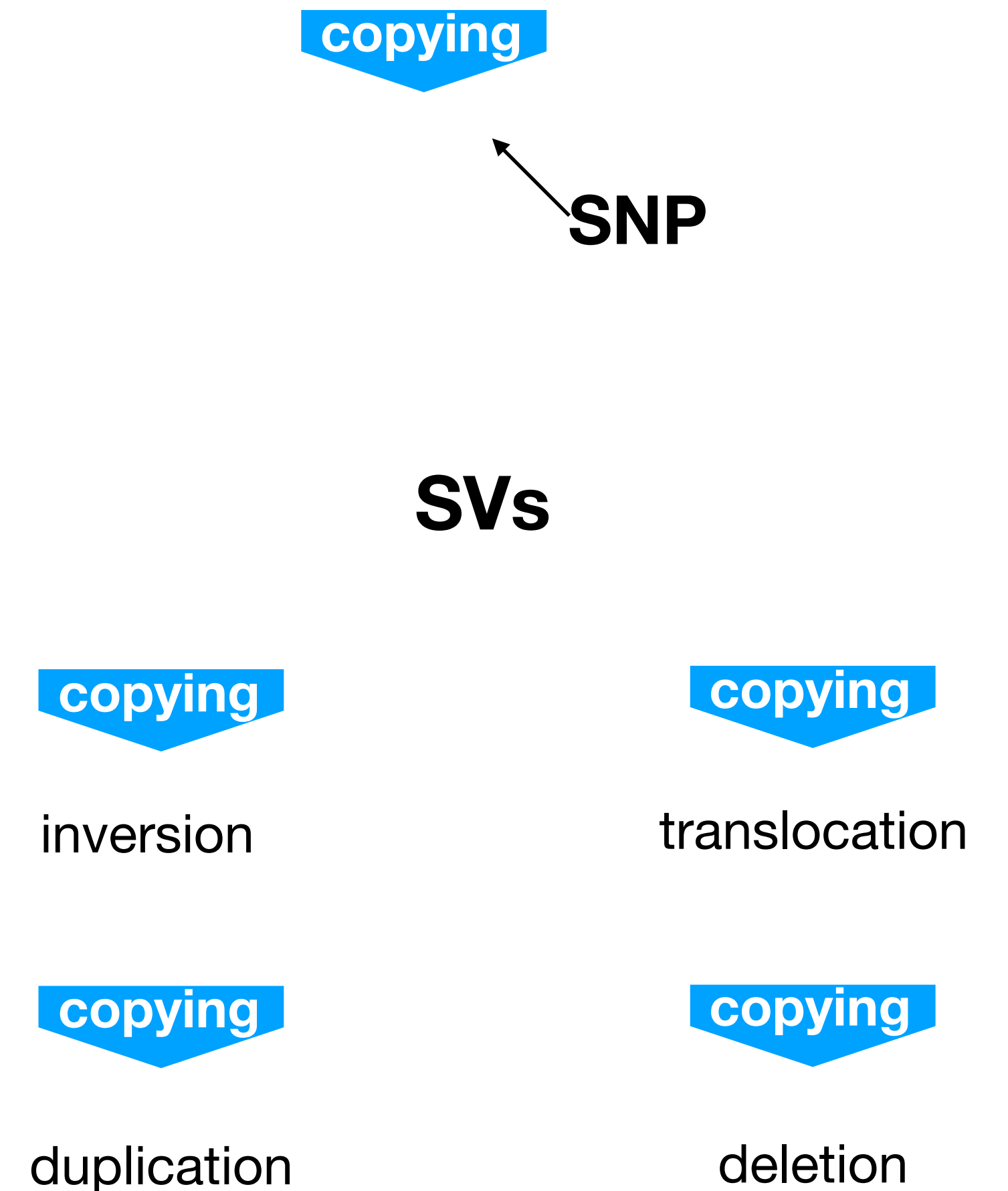
- Do stuff in the cell, including help with translation and transcription



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



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

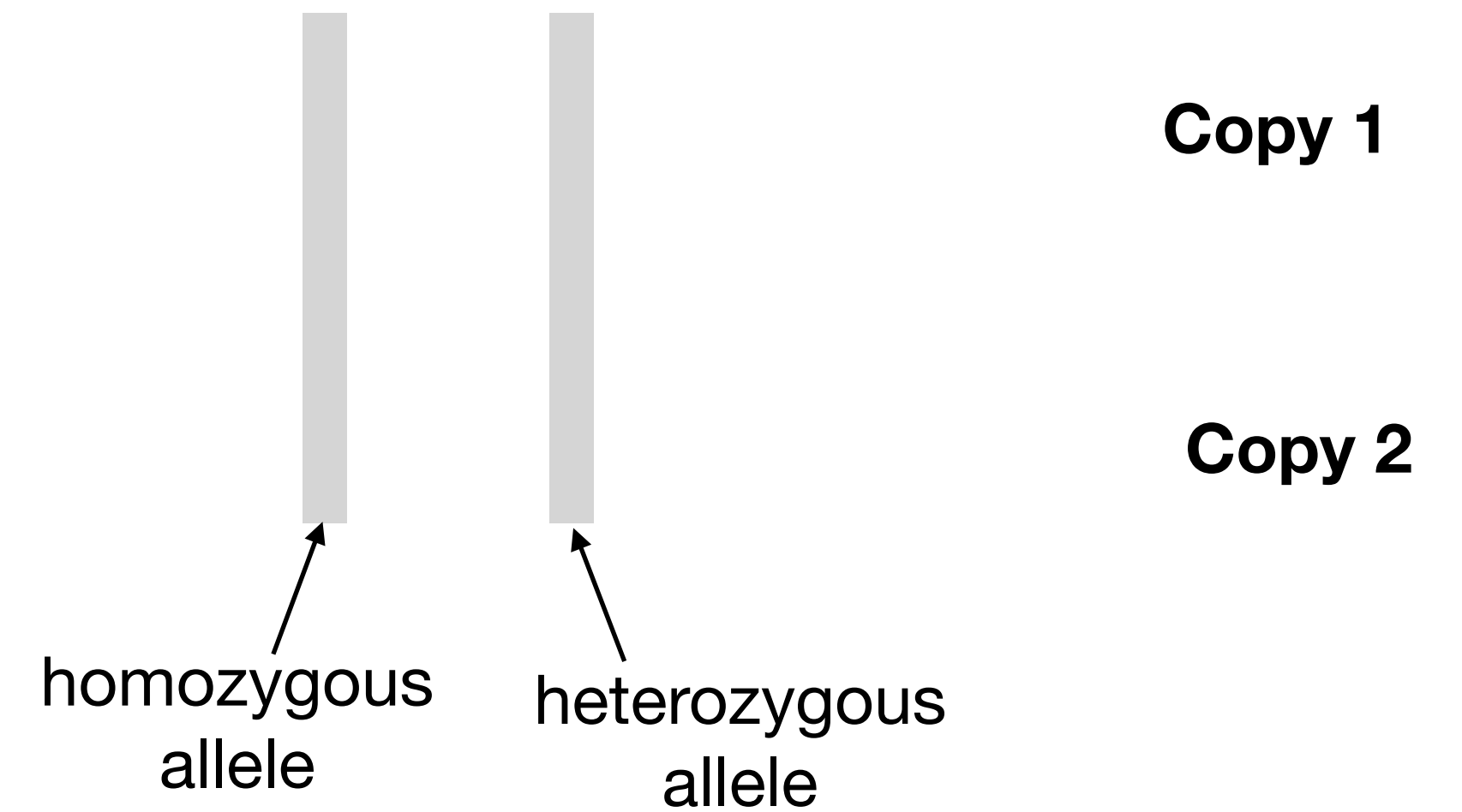
TACACCGTACGATCG
copying
TACACCGTAAGATCG
SNP

SVs

TACACCGTACGATCG copying TACACATGCCGATCG inversion	TACACCGTACGATCG copying TACAGATCCGTACCG translocation
TACACCGTACGATCG copying TACACCGTACGATCCGTACCG duplication	TACACCGTACGATCG copying TACAGATCG deletion

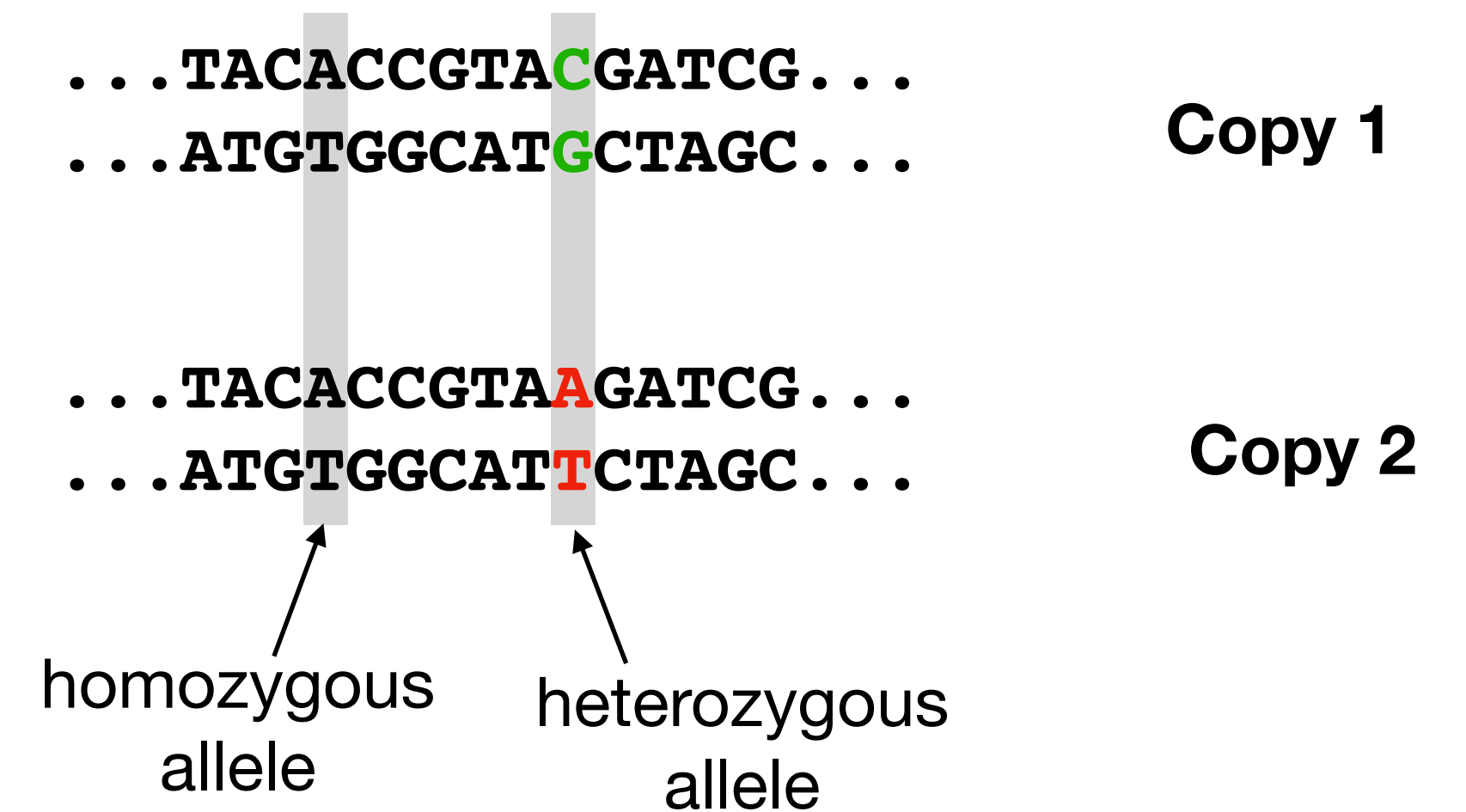
Genetic Variants

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

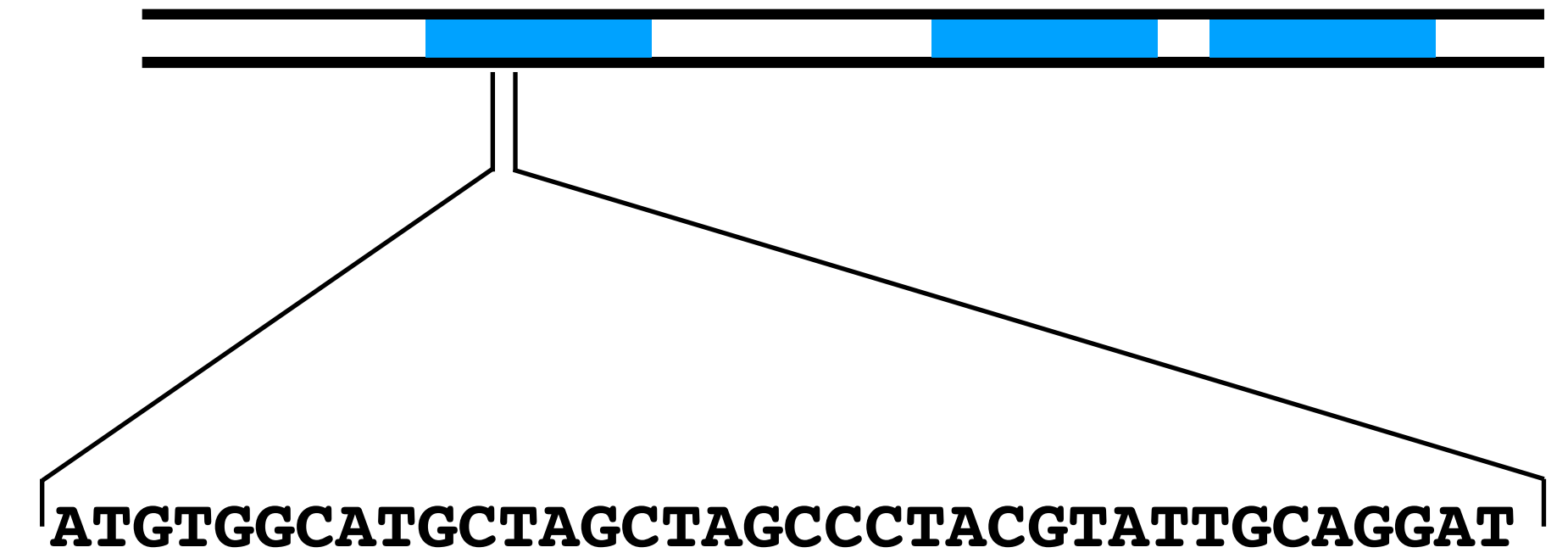


Genetic Variants

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

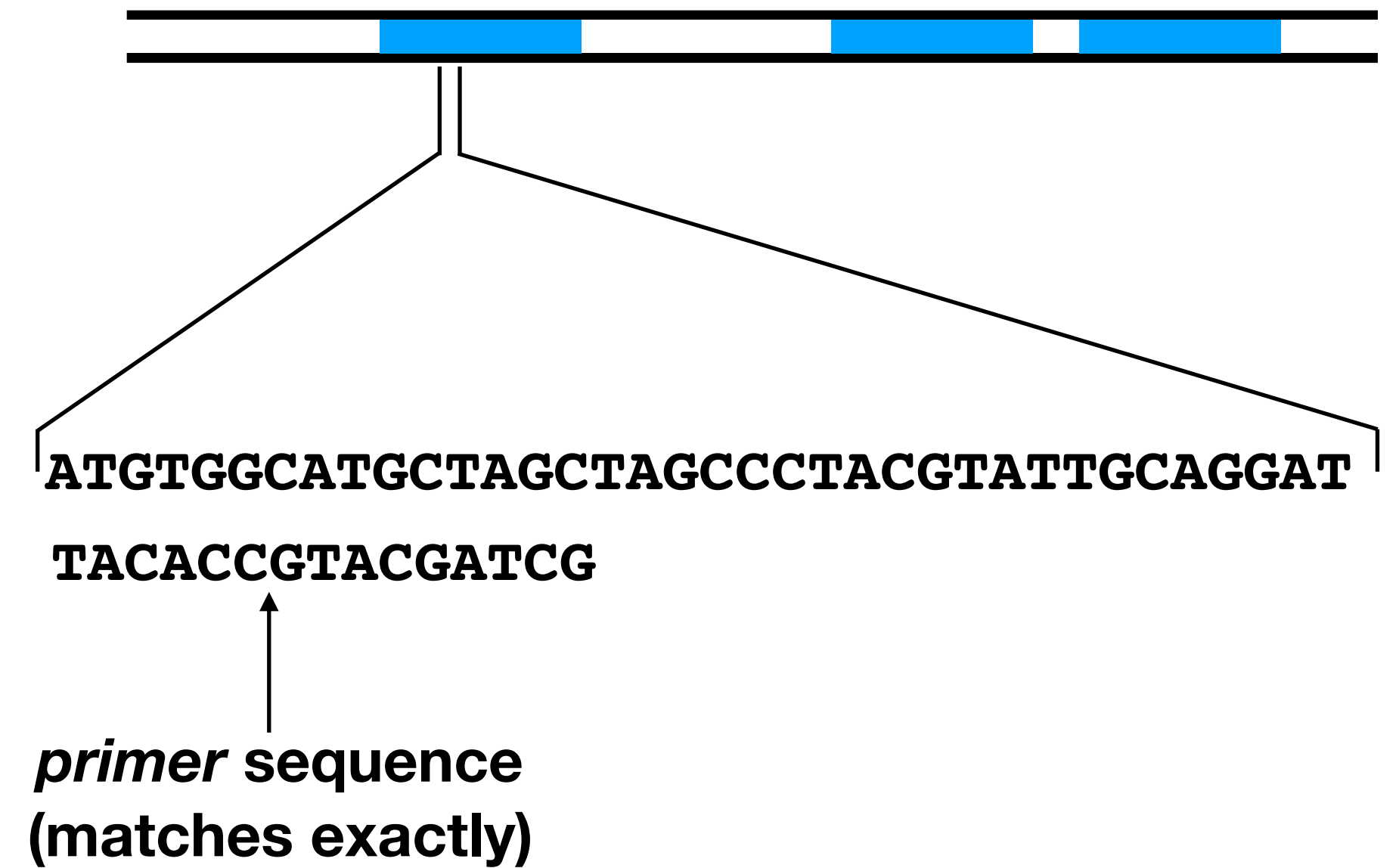


Sanger Sequencing



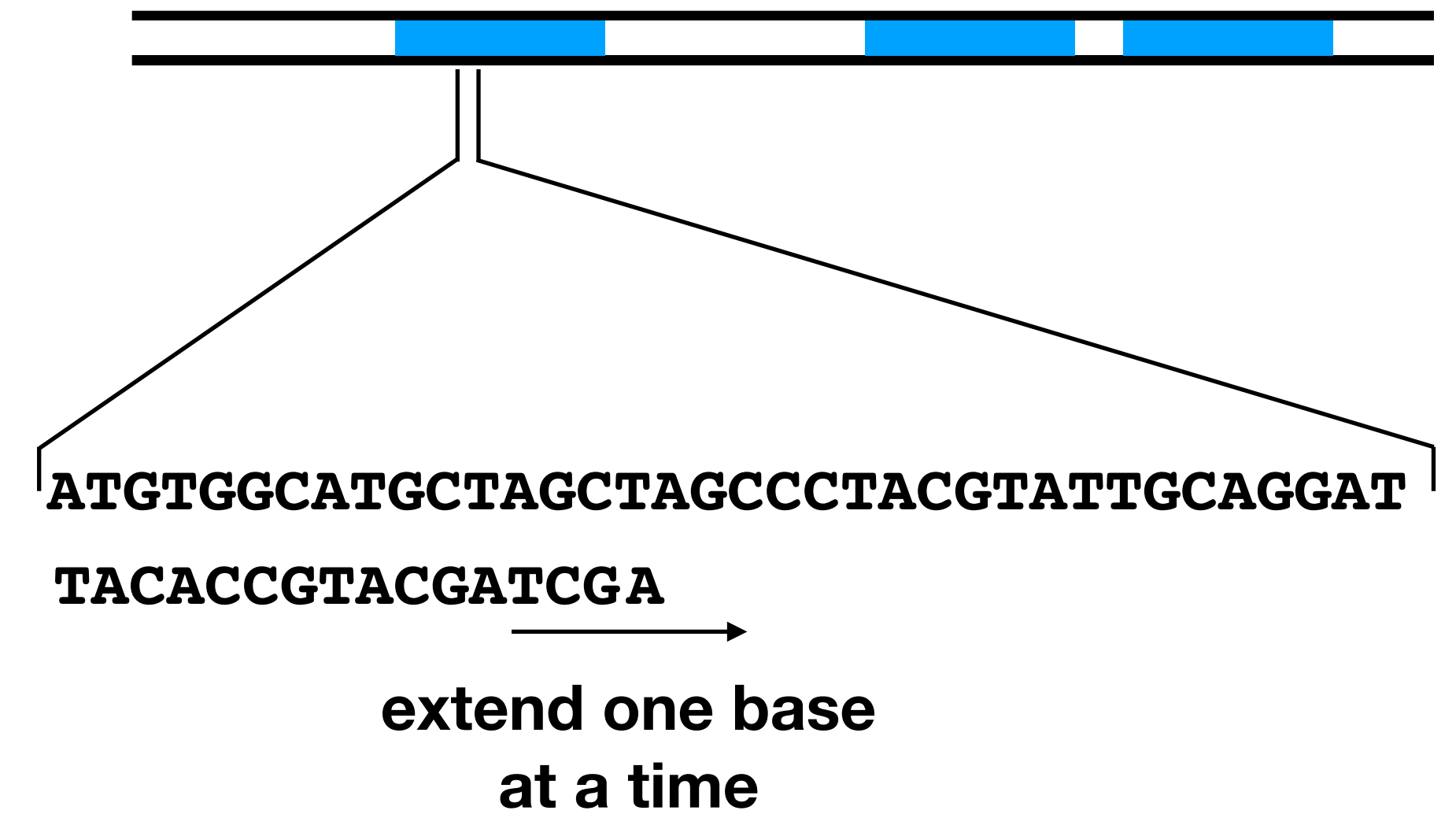
The basis of all modern sequencing.

Sanger Sequencing



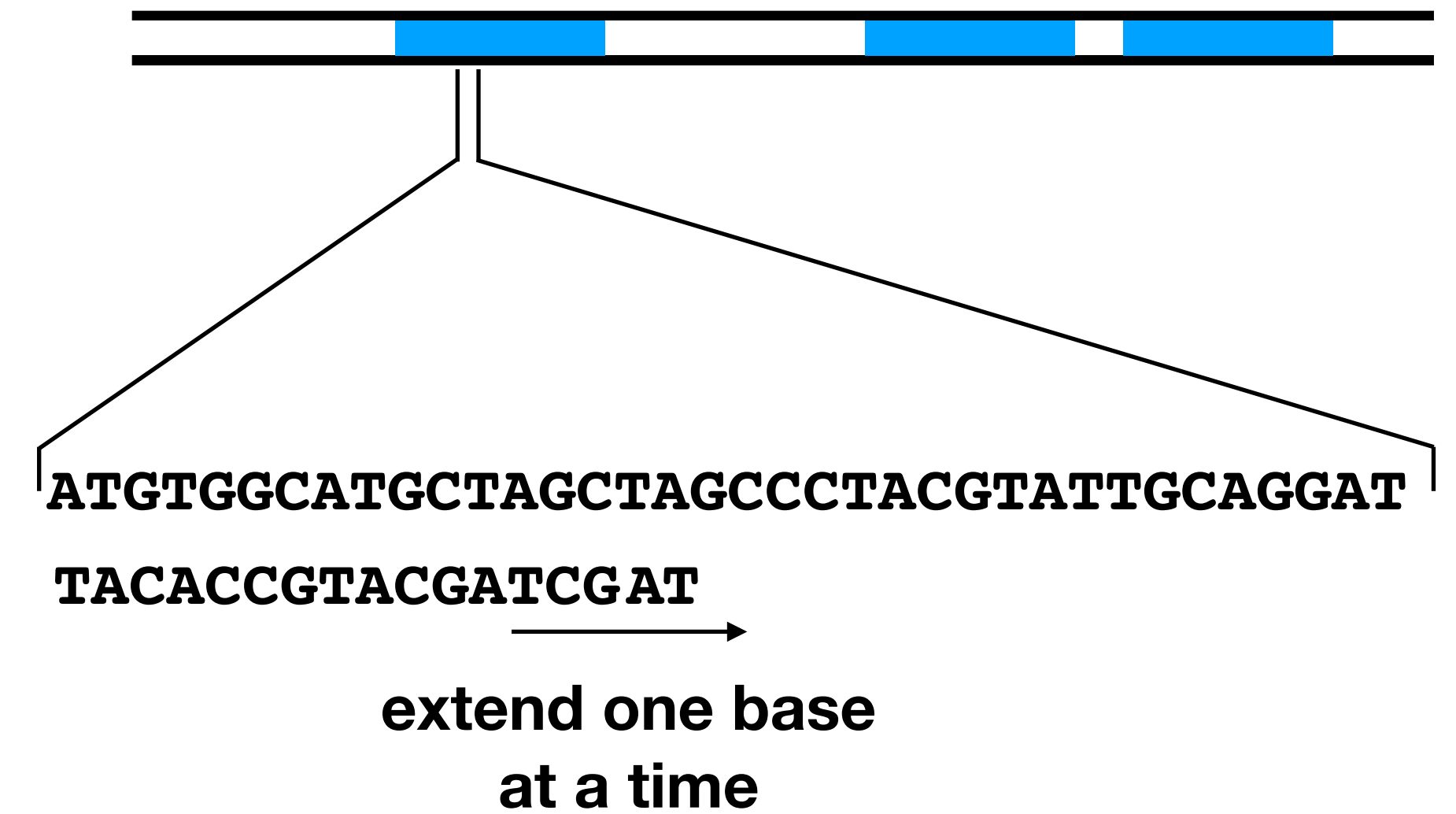
The basis of all modern sequencing.

Sanger Sequencing



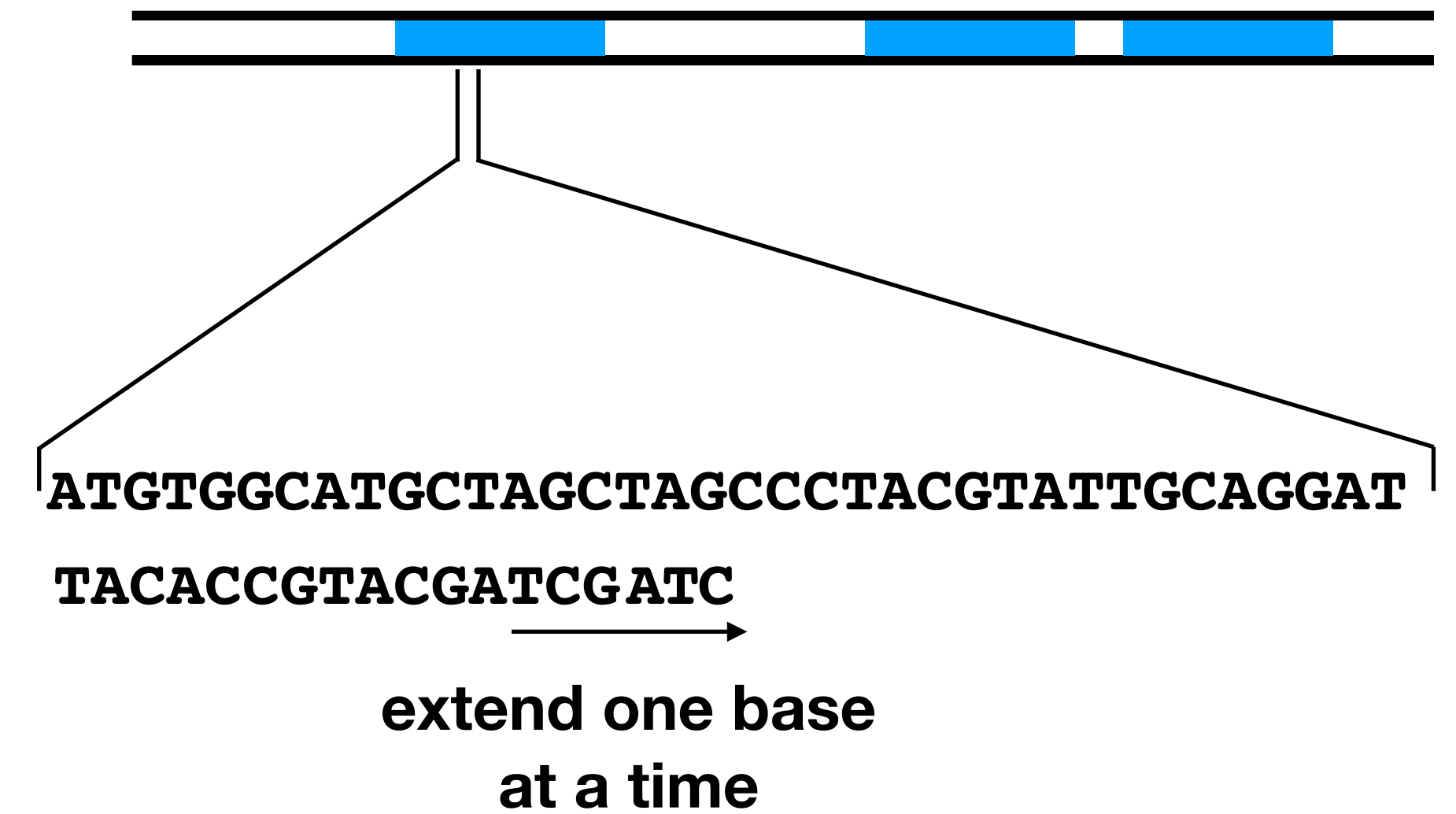
The basis of all modern sequencing.

Sanger Sequencing



The basis of all modern sequencing.

Sanger Sequencing



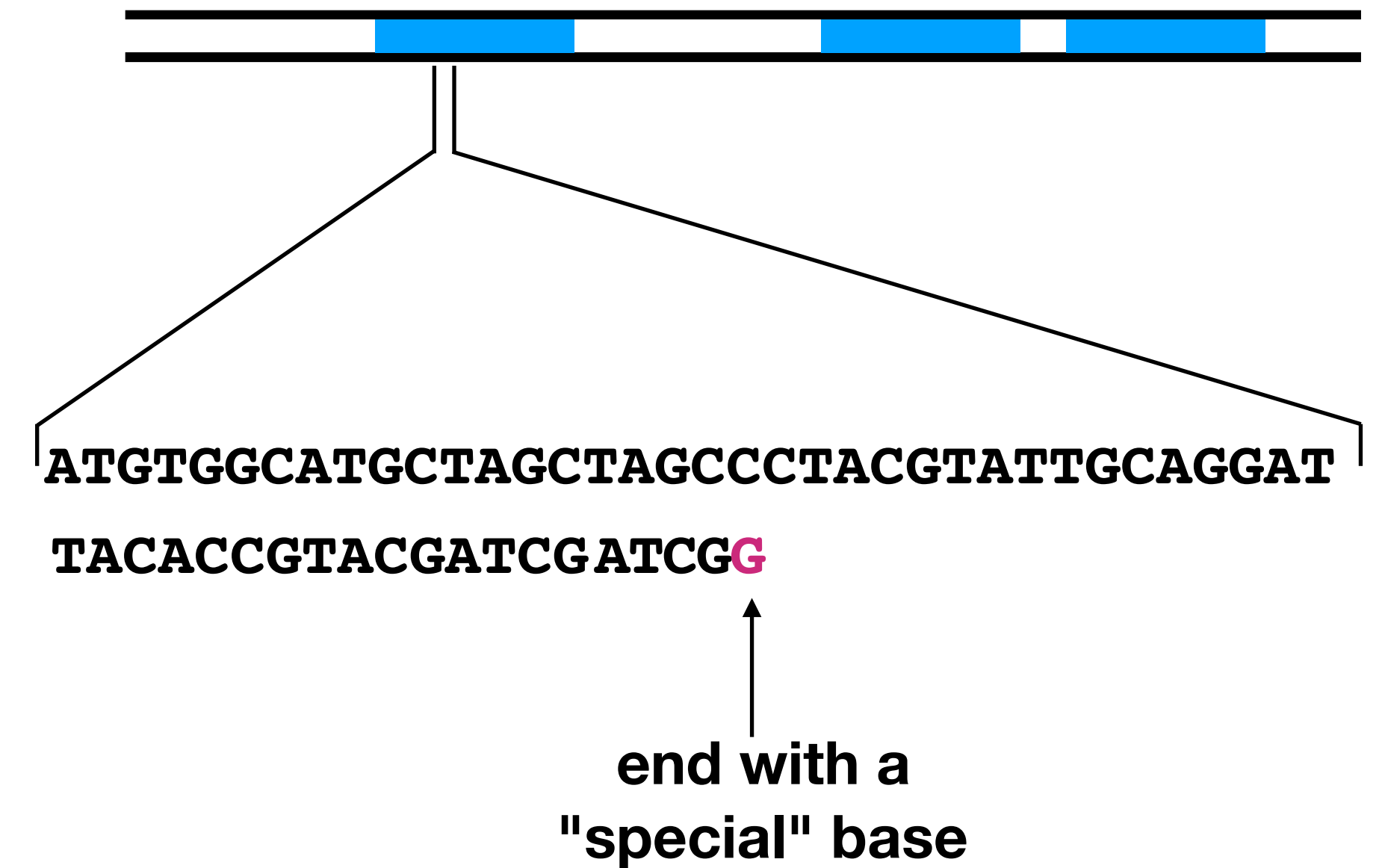
The basis of all modern sequencing.

Sanger Sequencing



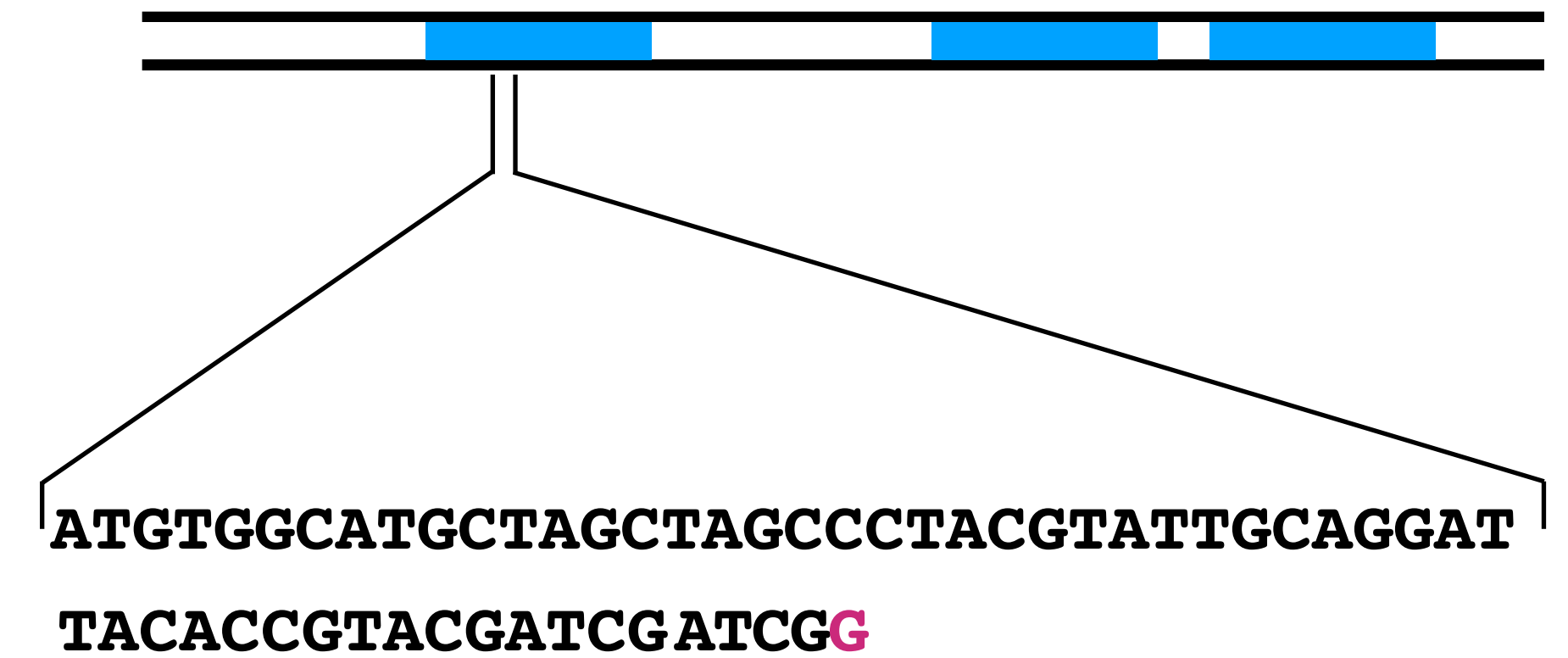
The basis of all modern sequencing.

Sanger Sequencing



The basis of all modern sequencing.

Sanger Sequencing



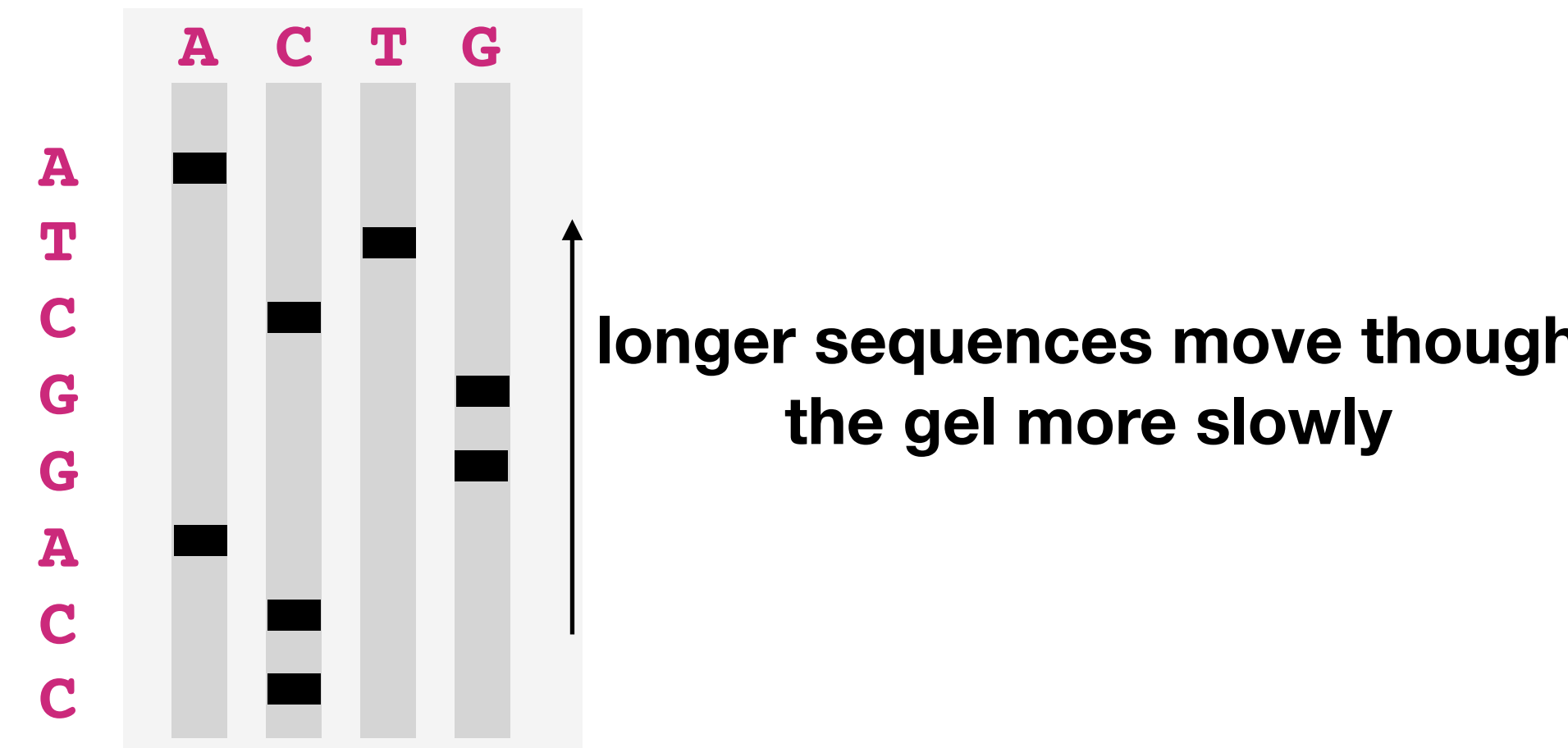
The basis of all modern sequencing.

Sanger Sequencing

...

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

The basis of all modern sequencing.

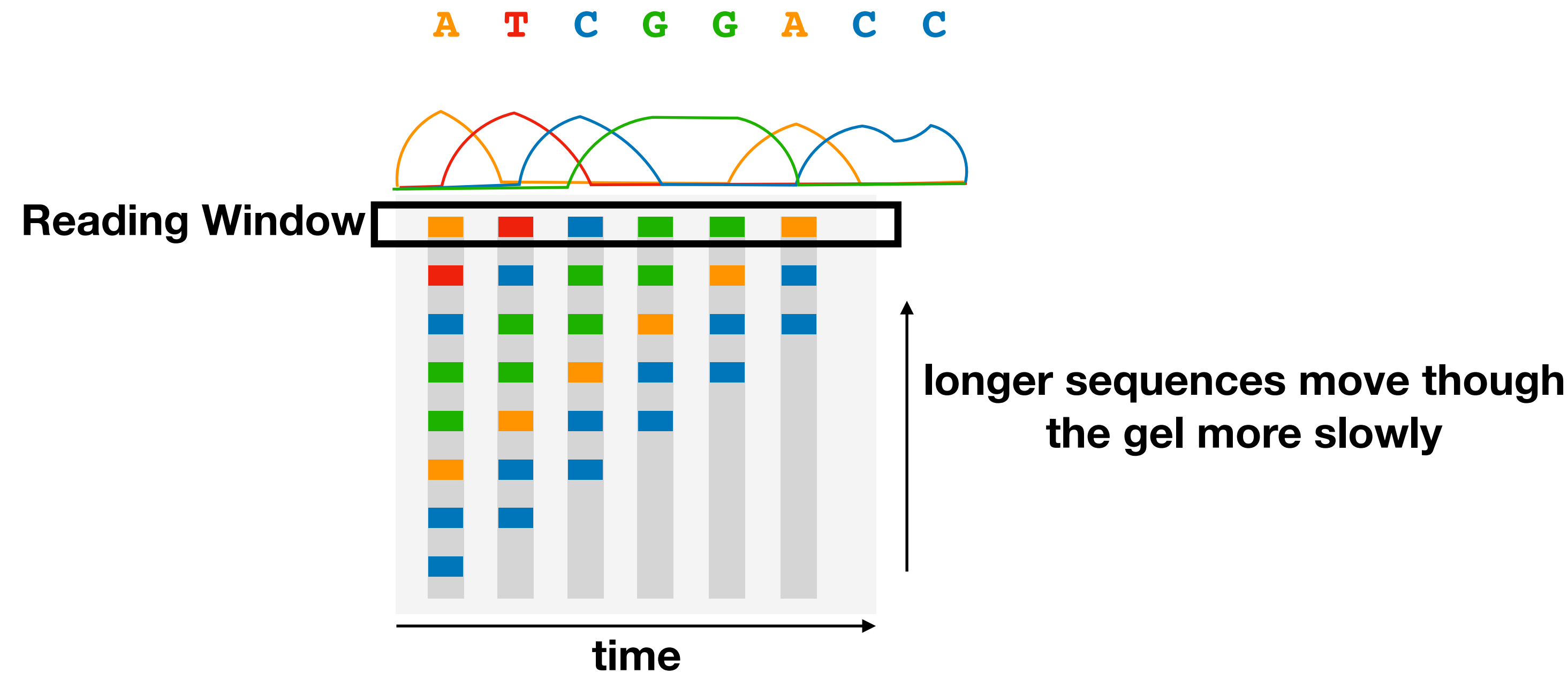


Sanger Sequencing

...

TACACCGTACGATCGATCGG
TACACCGTACGATCGATCG
TACACCGTACGATCGATC
TACACCGTACGATCGAT
TACACCGTACGATCGA

The basis of all modern sequencing.



Second Generation Sequencing

Also called next generation sequencing

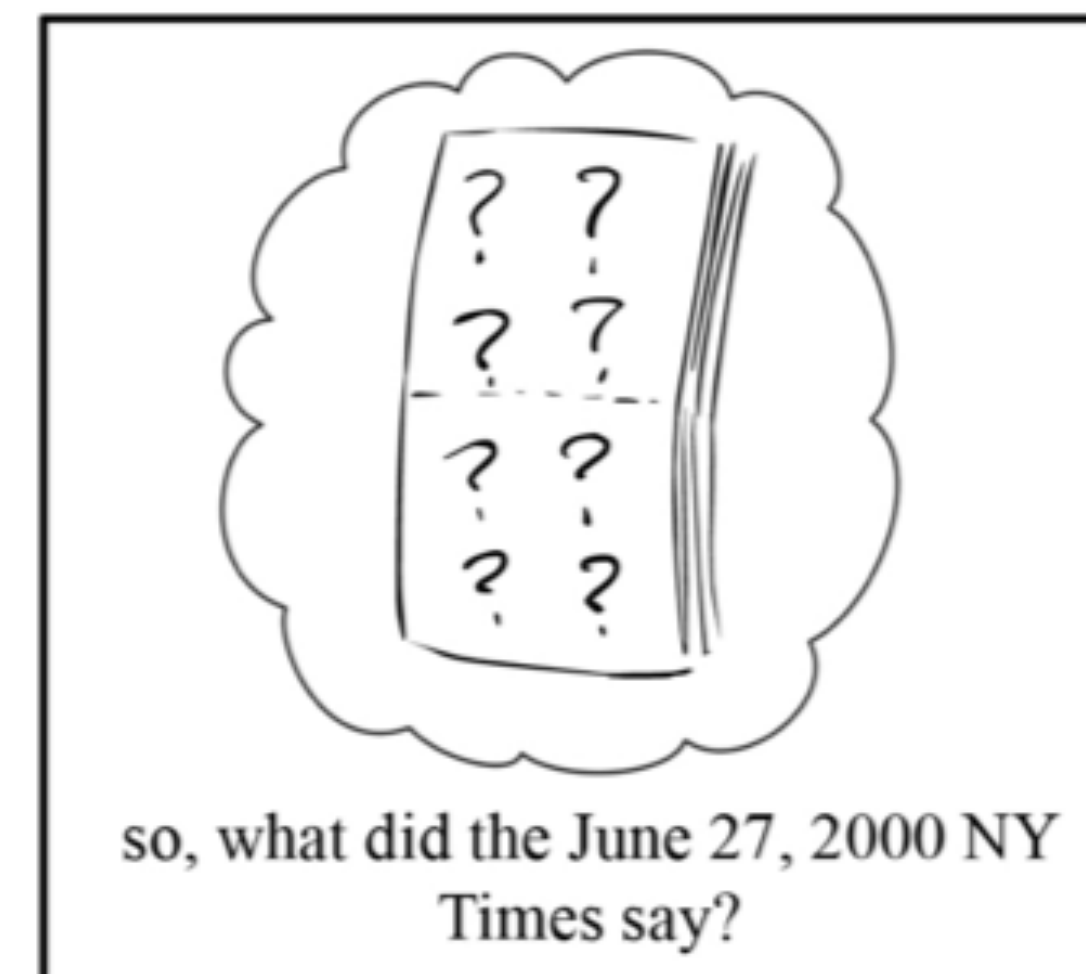
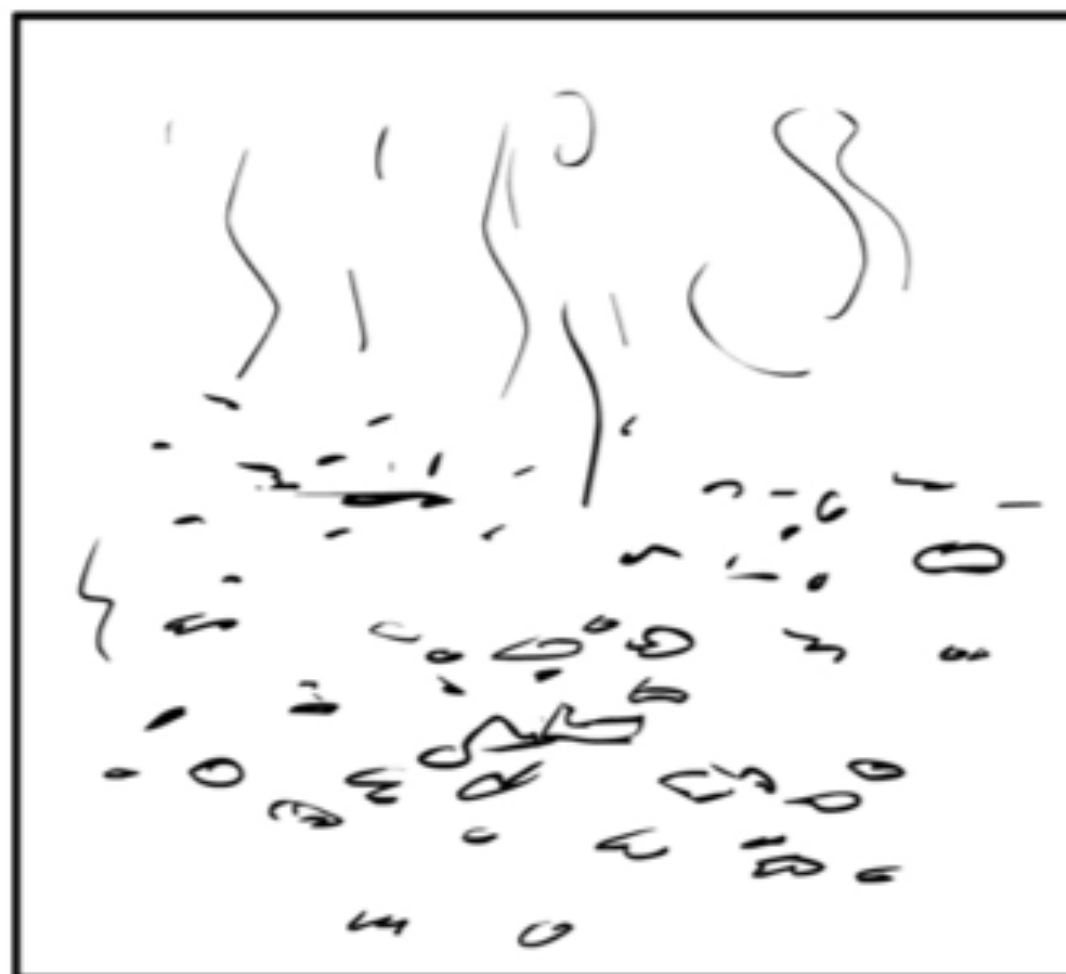
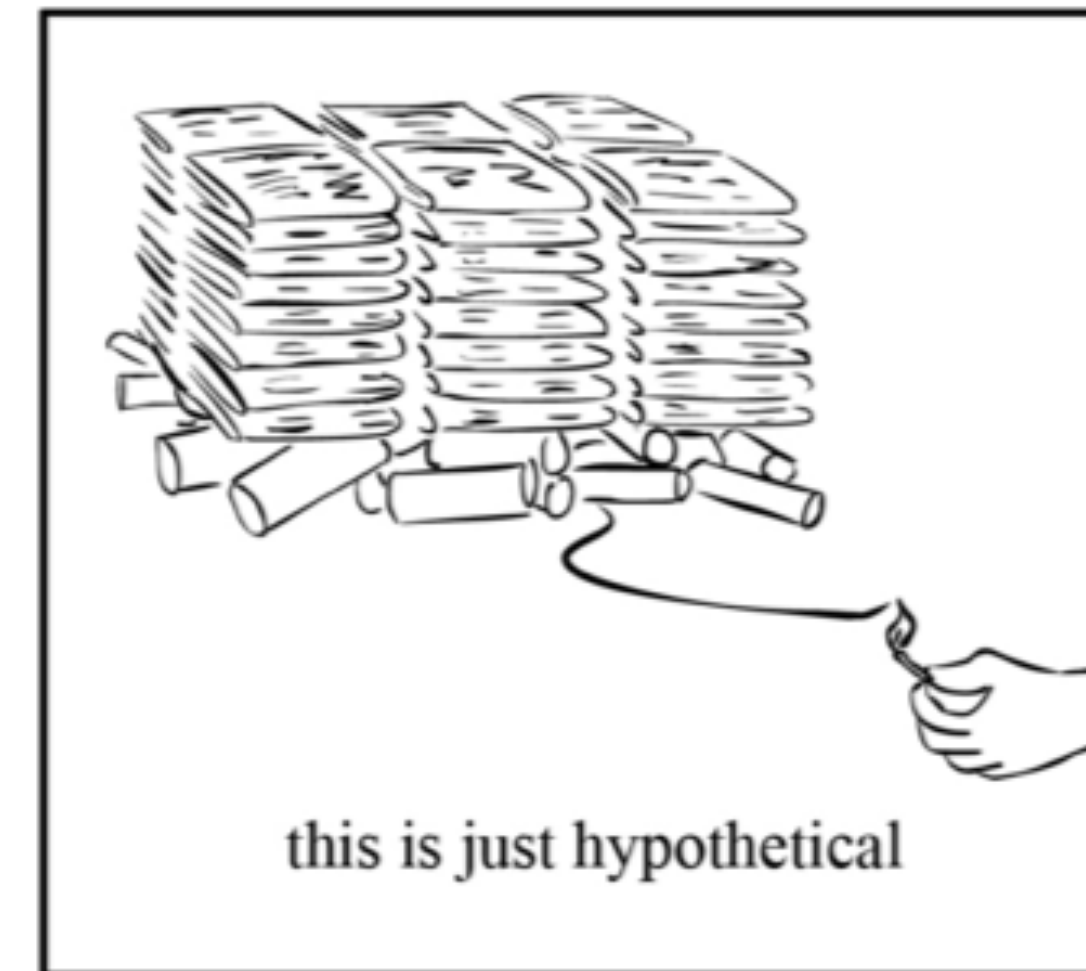
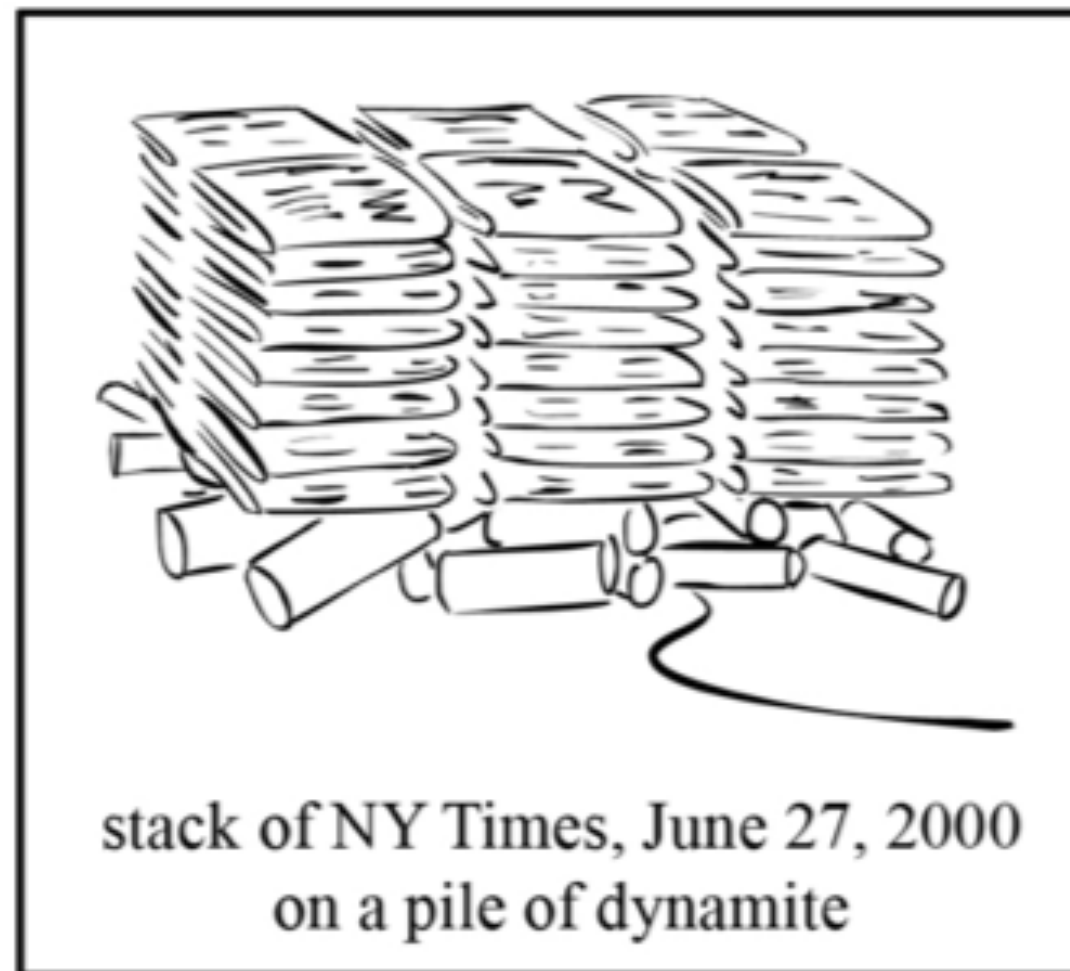
Based on the same principles, but at a much larger scale

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

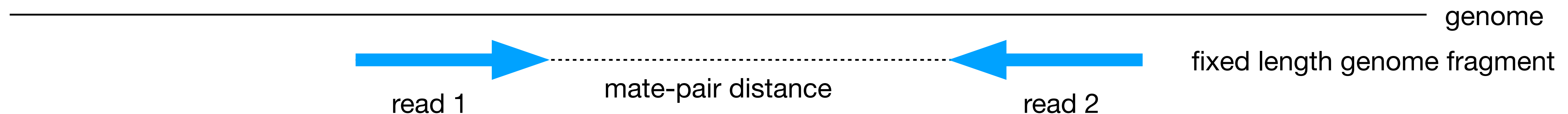
Second Generation Sequencing



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



Third Generation Sequencing

Recently Pacific Biosciences and Oxford Nanopore have introduced new technologies that:

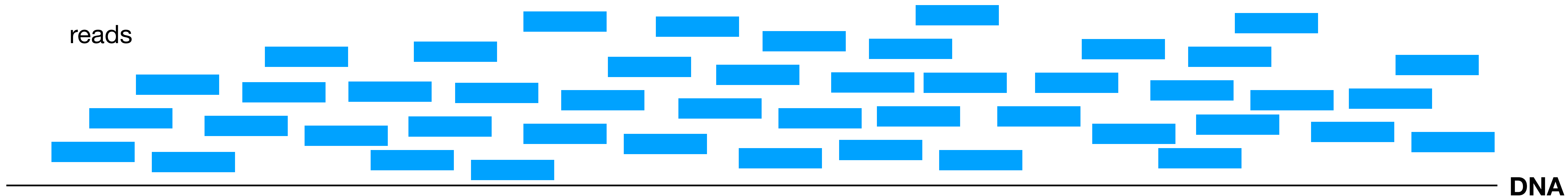
- have long reads
- with high(er) error rates

	Sanger	Next-Generation	Third-Generation
Launched	1977 Basic chemistry 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% (but decreasing)
Cost			
Throughput			
Currently Available Platforms	Applied Biosystems*	Illumina Ion Torrent* Qiagen (Europe) Complete Genomics (China)**	Pacific Biosciences Oxford Nanopore
Clinical Uses	Many (but dwindling)	Many (and growing)	Niche uses (today)

*Part of Thermo Fisher

**Part of BGI

Sequencing Applications



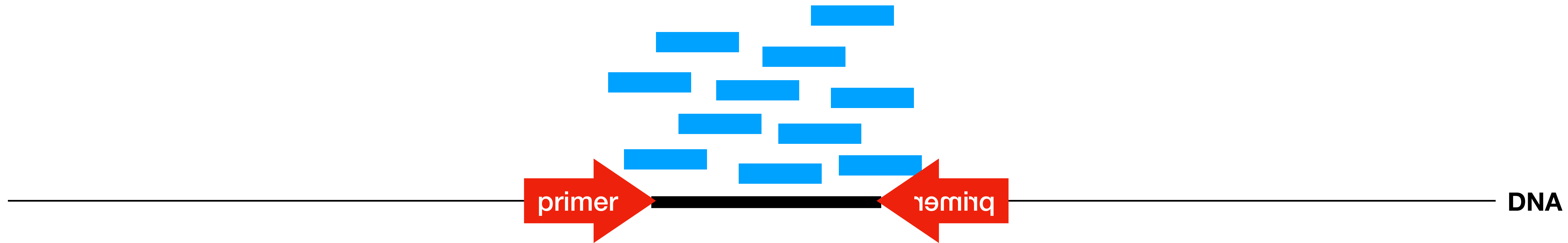
whole genome sequencing

Sequencing Applications



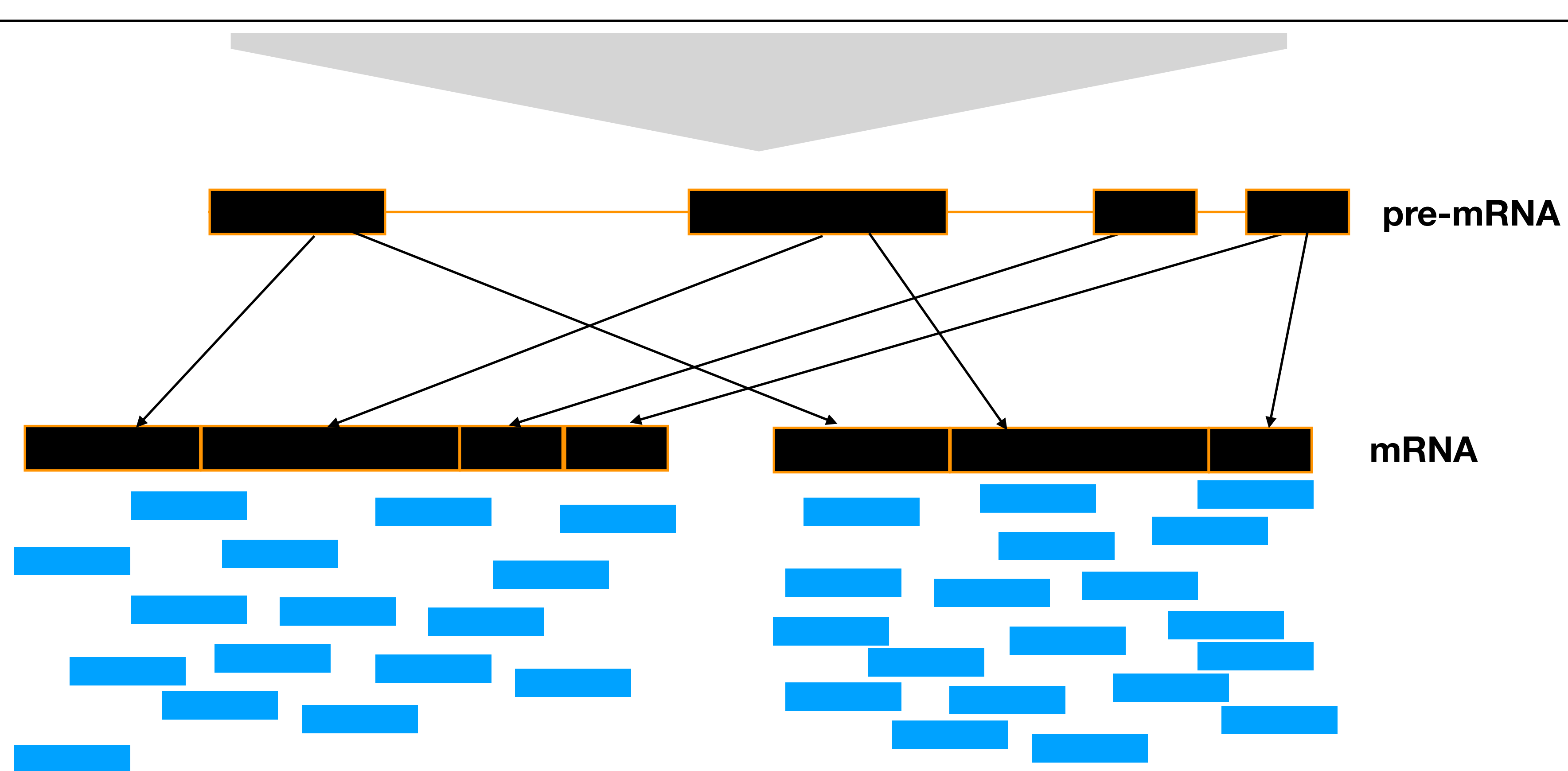
bisulphite sequencing

Sequencing Applications



targeted sequencing

Sequencing Applications



DNA

RNA sequencing

adapted from figure 1.2 in Mäkinen, *et al.* 2015

Sequencing Applications

