

CSE8803/CX4803

Machine Learning in Computational Biology

Lecture 2: Primer on Molecular Biology

Yunan Luo

What is Computational Biology/Bioinformatics?

Computational biology and **bioinformatics** is an interdisciplinary field that develops and applies **computational methods** to analyze large collections of biological data, such as genetic sequences, cell populations or protein samples, to make new predictions or **discover new biology**.

<https://www.nature.com/subjects/computational-biology-and-bioinformatics>

Reading

Please read “[Molecular Biology for Computer Scientists](#)” by
Lawrence Hunter

Molecular biology

Molecular Biology is the field of **biology** that studies the composition, structure and interactions of cellular **molecules** – such as nucleic acids and proteins – that carry out the **biological** processes essential for the cell's functions and maintenance.

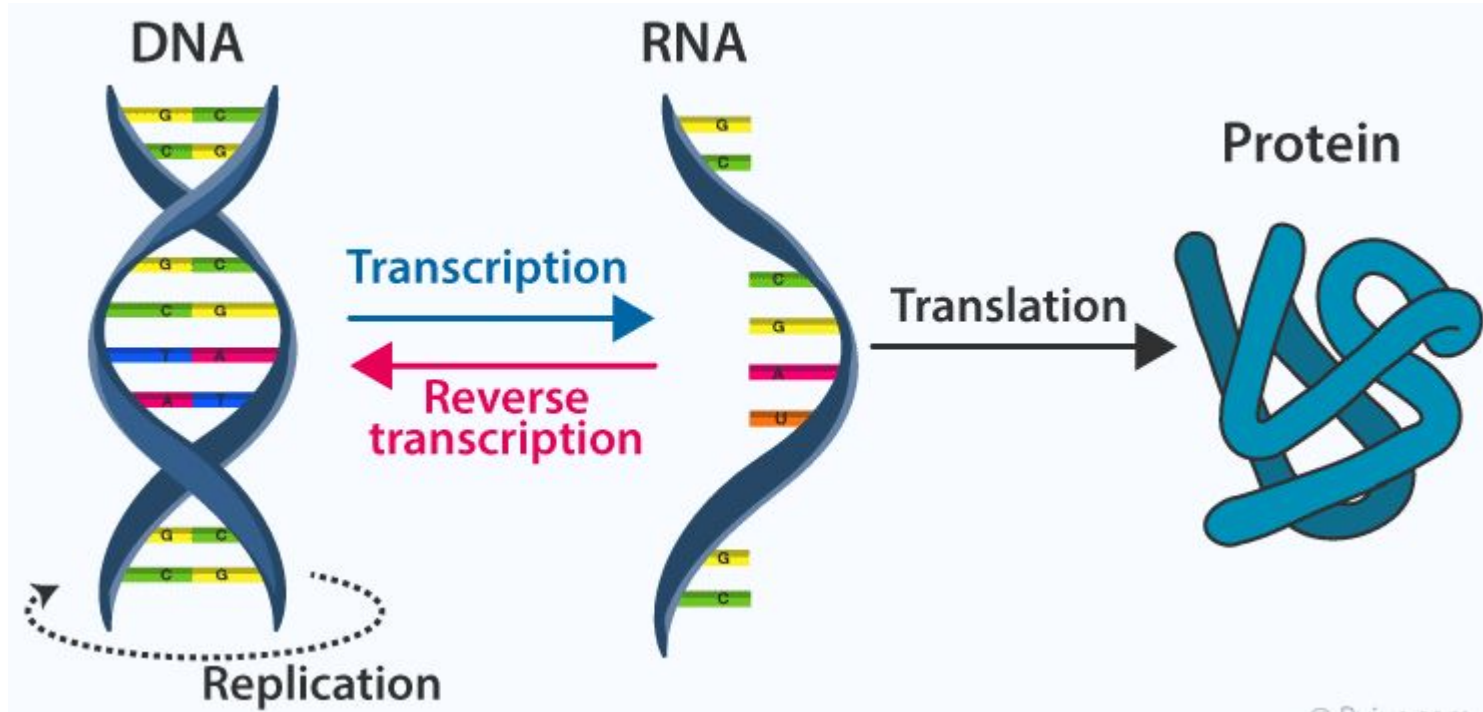
Cellular molecules:

1. DNA

2. RNA

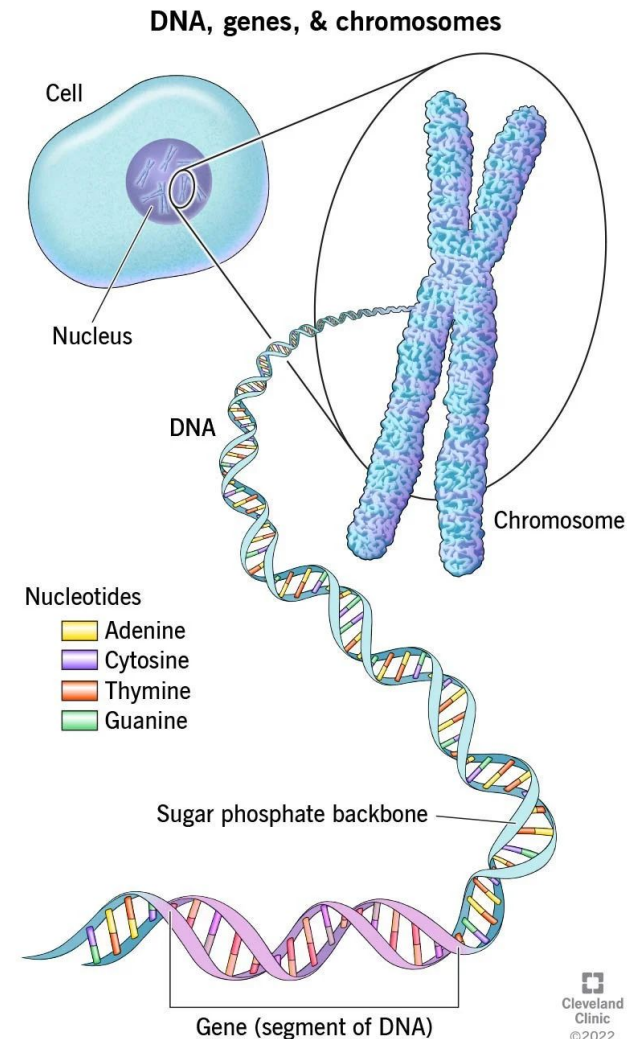
3. Protein

Central dogma: DNA → RNA → Protein

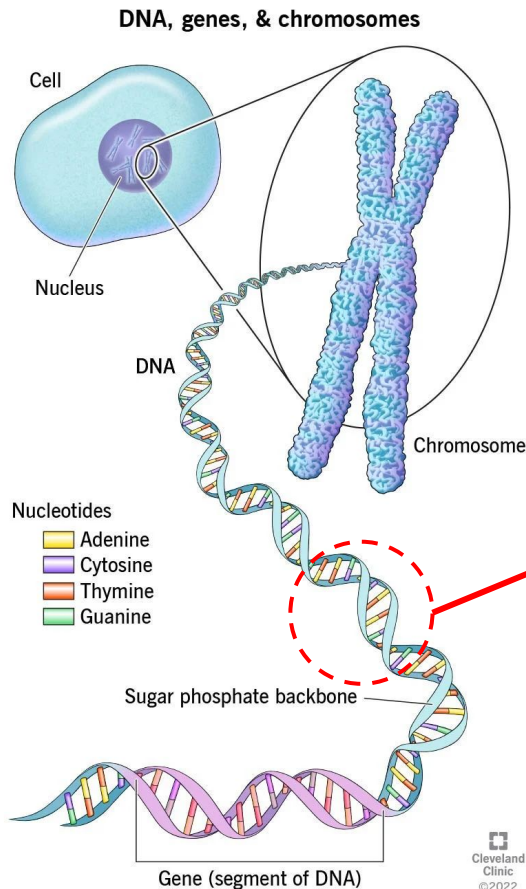


DNA

- DNA discovered as the physical (molecular) carrier of hereditary information
- DNA is a molecule: *deoxyribonucleic acid*
- DNA is a very “long” molecule
- DNA in human has 3 billion base-pairs
 - String of 3 billion characters! (about 6 feet long)
- DNA harbors “genes”
 - A gene is a substring of the DNA string
 - A gene “codes” for a protein



DNA



Cleveland
Clinic
©2022

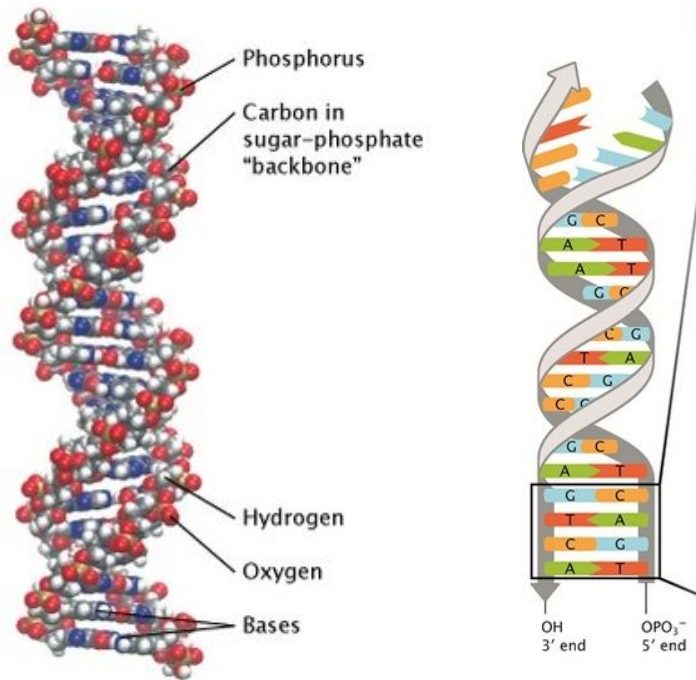


Watson & Crick, *Nature*, 1953

Double helical structure (discovered by Watson, Crick & Franklin)

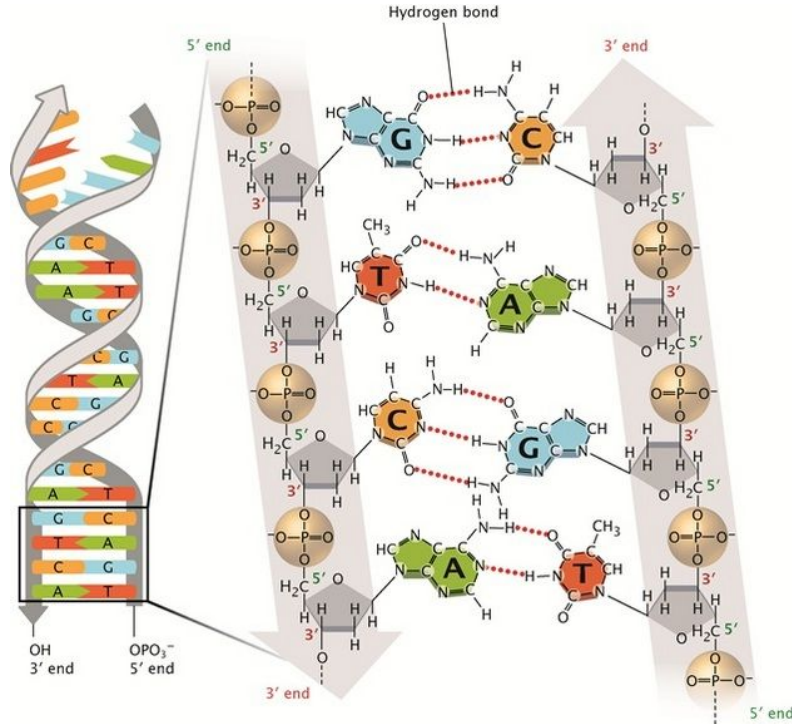
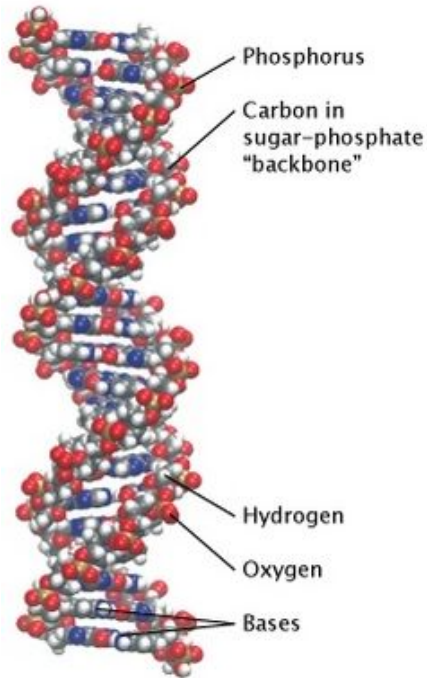
DNA

Each strand of the DNA is composed of sequence of covalently bonded **nucleotides (bases)**



DNA

Each strand of the DNA is composed of sequence of covalently bonded **nucleotides (bases)**



Four nucleotides:

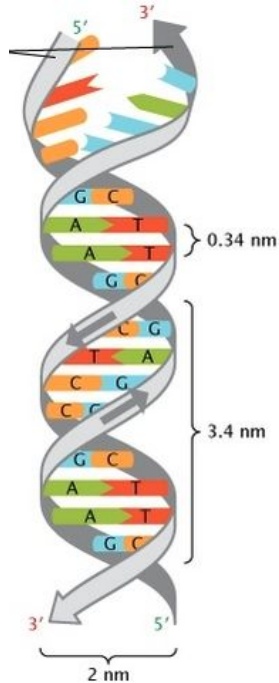
A (adenine)
C (cytosine)
T (thymine)
G (guanine)

Base pairing:

A <-> T
C <-> G

DNA

In the language of computer science, a DNA strand is a sequence s of over the alphabet of four characters $\{A, C, G, T\}$



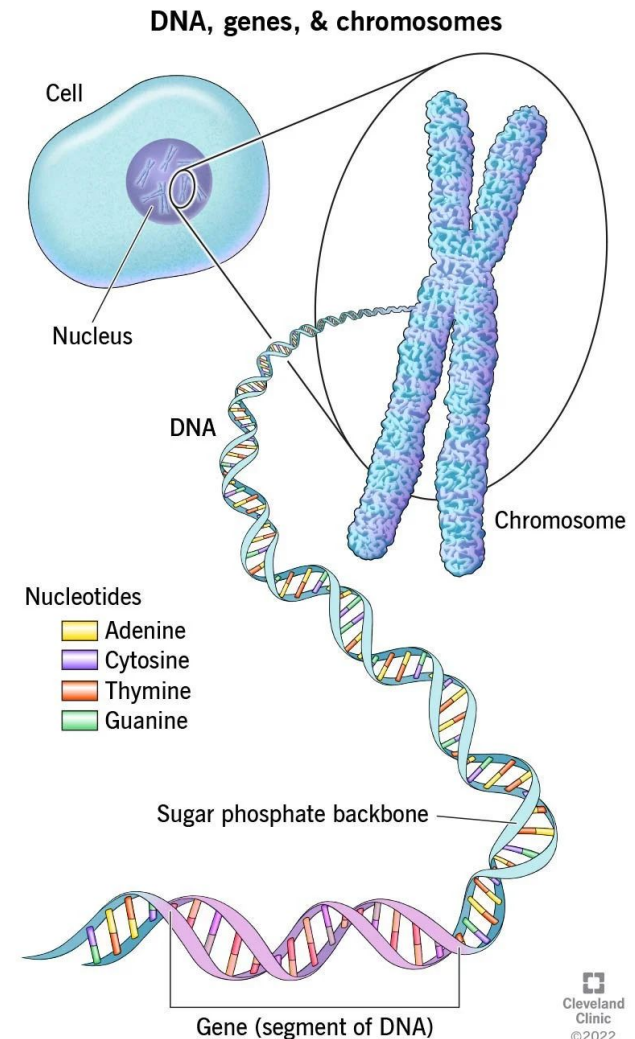
5' ...ACGTGACTGAGGACCGTG... 3'
... ||||| ||||| ||||| ||||| ...
3' ...TGCCTGACTCCTGGCAC... 5'

Or simply:

5' ...ACGTGACTGAGGACCGTG... 3'

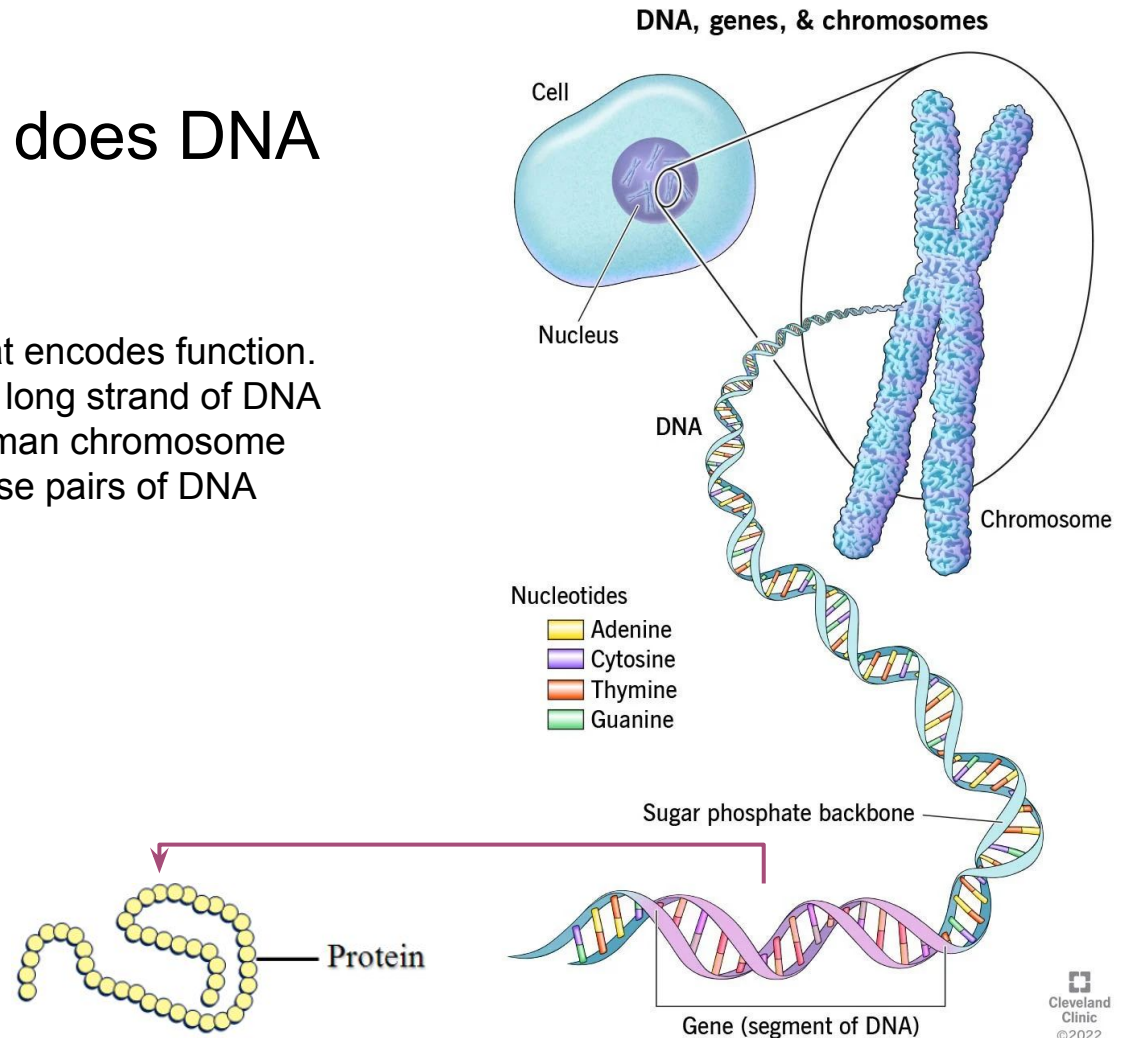
DNA to chromosome

A chromosome is a long DNA molecule with part or all of the genetic material of an organism.



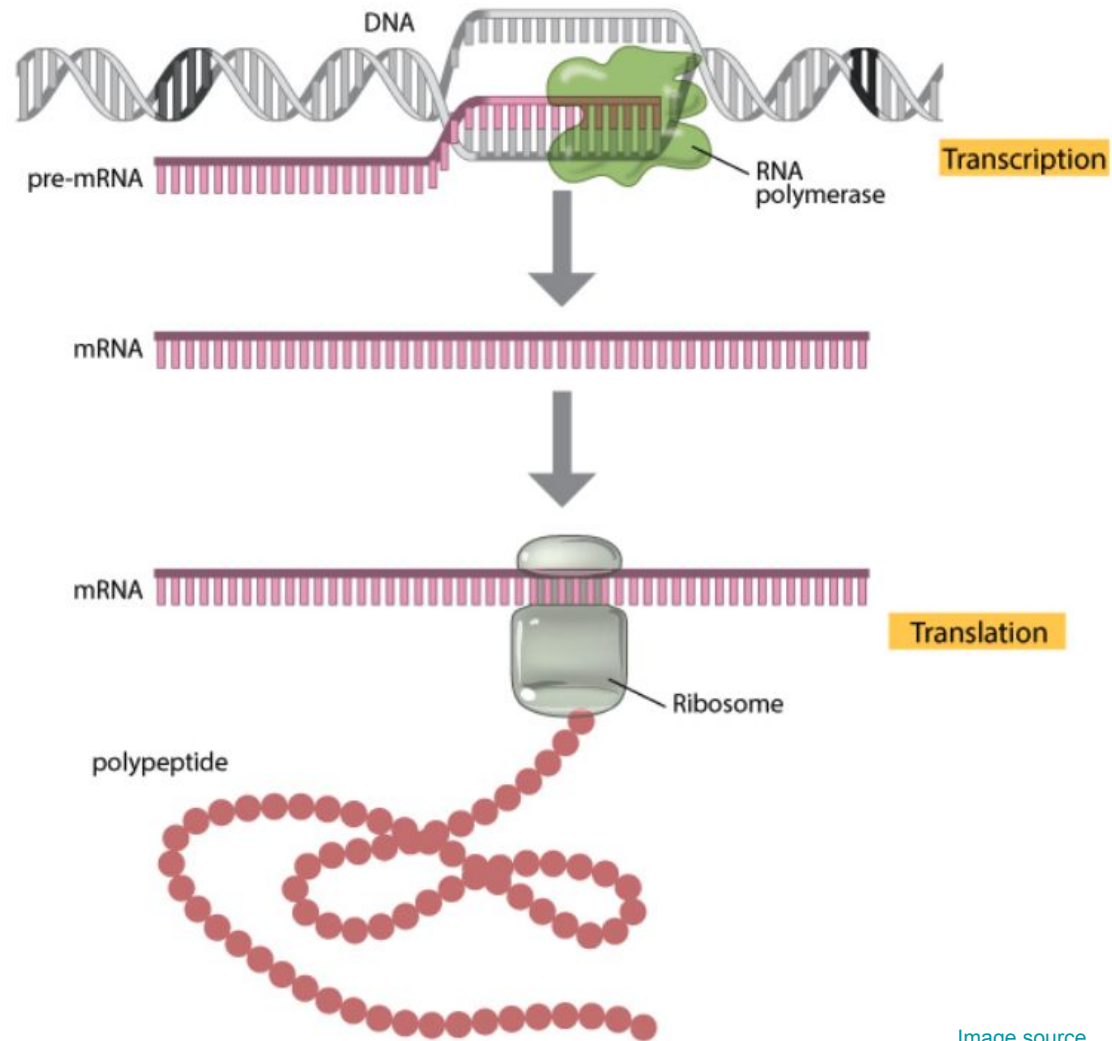
What information does DNA encode?

A **gene** is a region of DNA that encodes function. A **chromosome** consists of a long strand of DNA containing many genes. A human chromosome can have up to 500 million base pairs of DNA with thousands of genes.

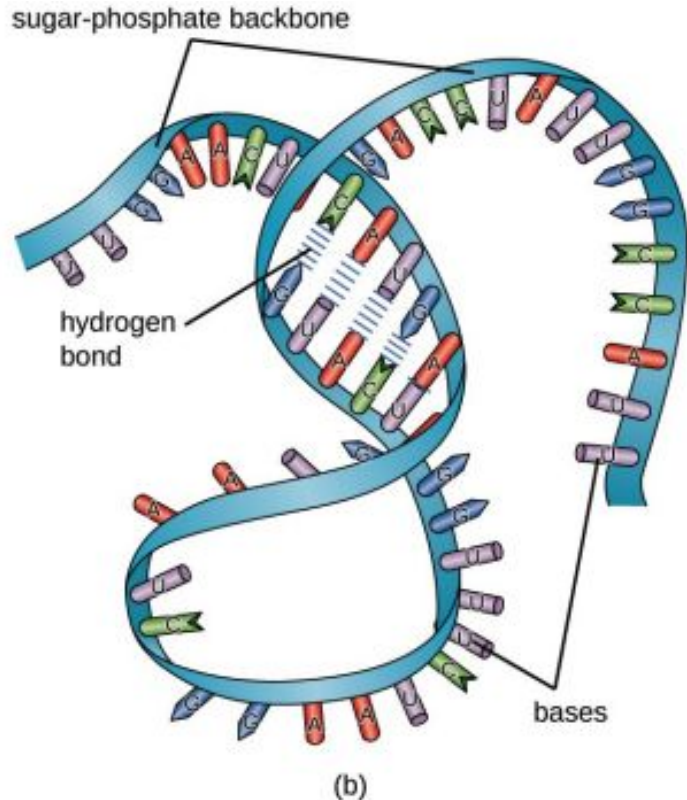


Central dogma: DNA → RNA → Protein

The process by which cells
“read” the genome



What is RNA?

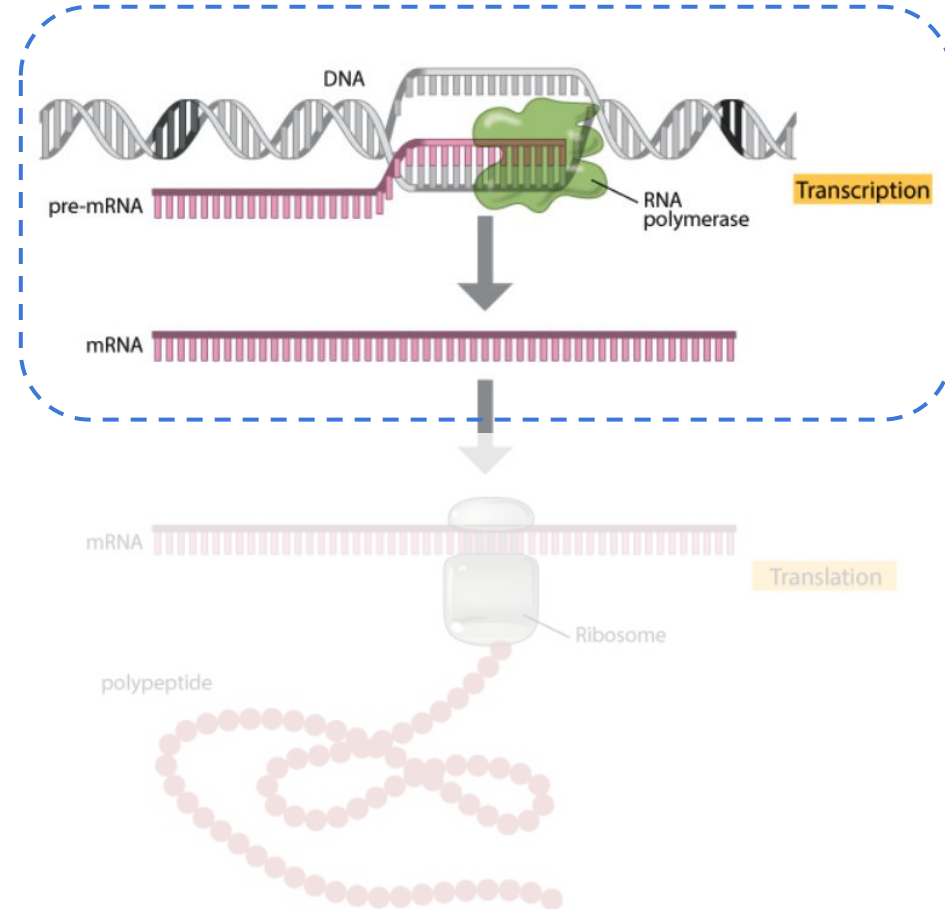


RNA = ribonucleic acid

- “U” instead of “T”
 - A (adenine)
 - C (cytosine)
 - U (uracil)
 - G (guanine)
- Usually **single stranded**
- Can fold into **structures** due to base complementarity
 - $A \leftrightarrow U$, $C \leftrightarrow G$
- Comes in many flavors:
 - mRNA, rRNA, tRNA, tmRNA, snRNA, snoRNA, scaRNA, aRNA, asRNA, piwiRNA, etc

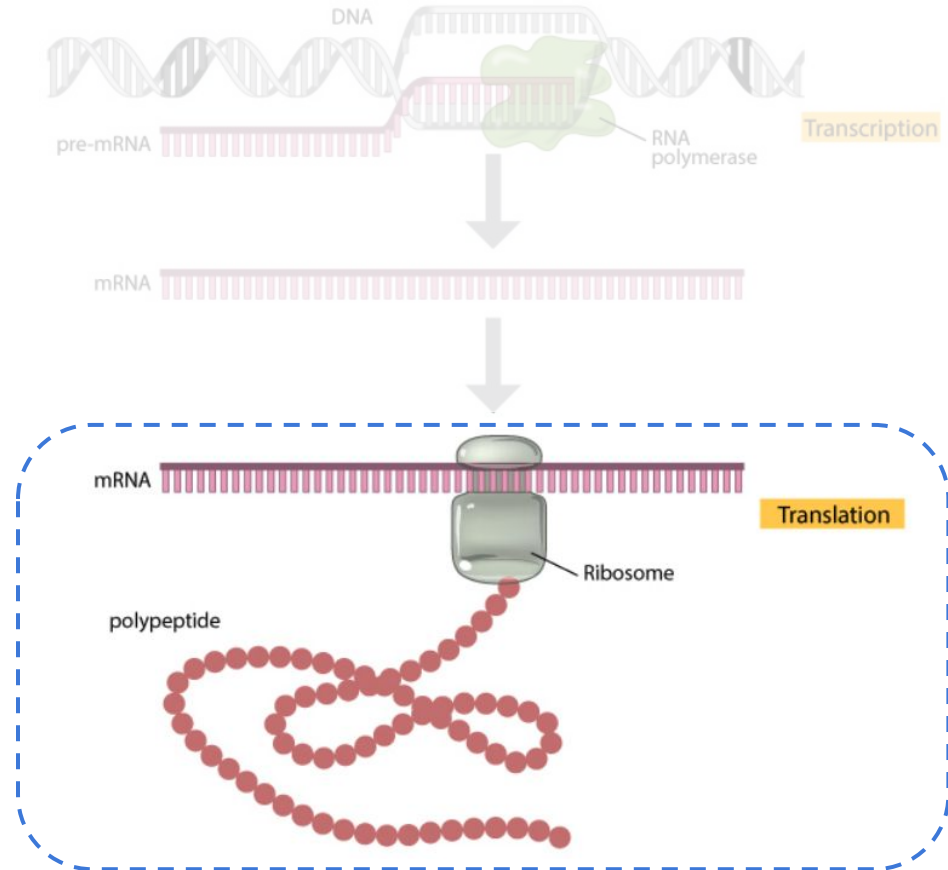
Transcription

- Process of making a single stranded mRNA using double stranded DNA as template
- Only genes are transcribed, not all DNA
- Gene has a transcription “start site” and a transcription “stop site”



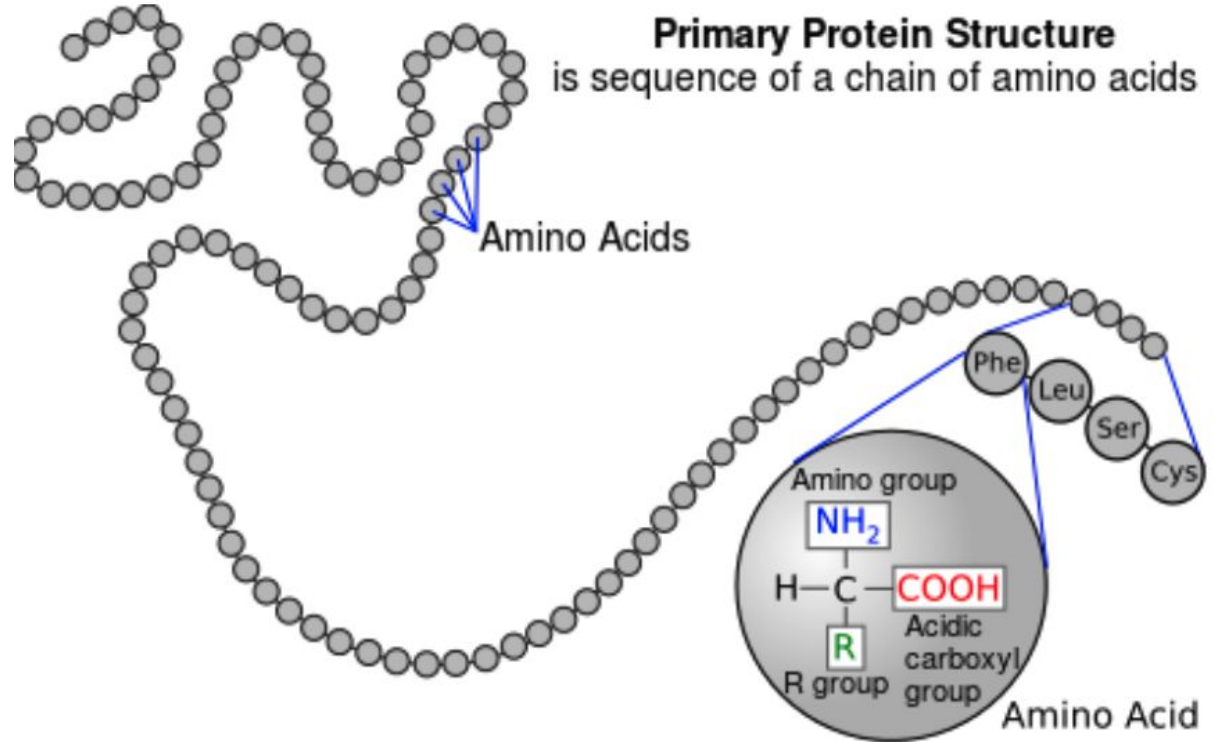
Translation

- Process of making an amino acid sequence from (single stranded) mRNA
- Each triplet of bases translates into one amino acid
- Each such triplet is called “codon”
- The translation is basically a table lookup

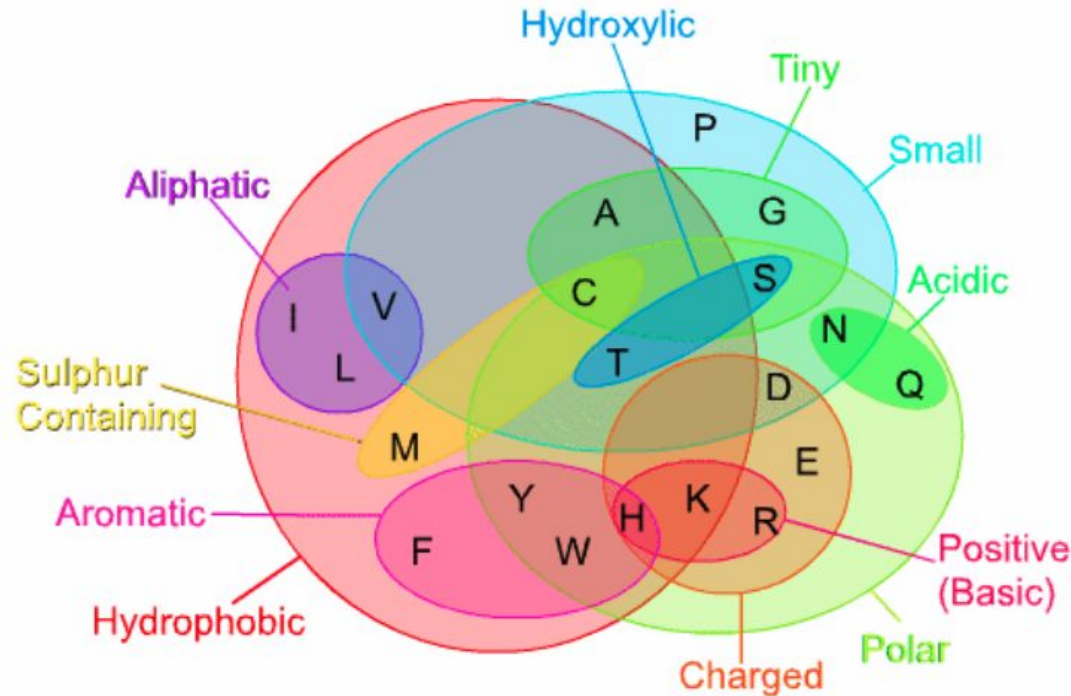


Protein sequence

In the language of CS, a protein sequence is a string **s** of over the alphabet of 20 characters



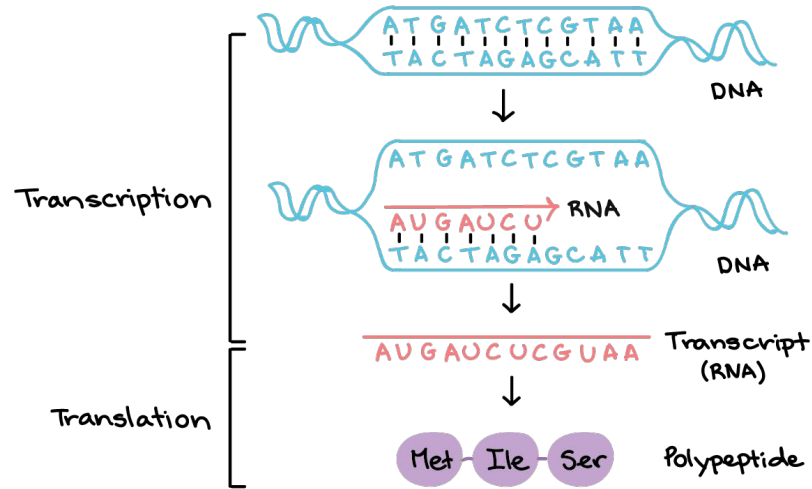
Alphabet: amino acids



Amino Acids

A alanine (ala)
R arginine (arg)
N asparagine (asn)
D aspartic acid (asp)
C cysteine (cys)
Q glutamine (gln)
E glutamic acid (glu)
G glycine (gly)
H histidine (his)
I isoleucine (ile)
L leucine (leu)
K lysine (lys)
M methionine (met)
F phenylalanine (phe)
P proline (pro)
S serine (ser)
T threonine (thr)
W tryptophan (trp)
Y tyrosine (tyr)

Genetic code: lookup table



[Image source](#)

		Second base				
		U	C	A	G	
First base	U	UUU } Phenyl-alanine F UUC } UUA } Leucine L UUG }	UCU } Serine S UCC } UCA } UCG }	UAU } Tyrosine Y UAC } UAA } Stop codon UAG } Stop codon	UGU } Cysteine C UGC } UGA } Stop codon UGG } Tryptophan W	U C A G
	C	CUU } Leucine L CUC } CUA } CUG }	CCU } Proline P CCC } CCA } CCG }	CAU } Histidine H CAC } CAA } Glutamine Q CAG }	CGU } Arginine R CGC } CGA } CGG }	U C A G
	A	AUU } Isoleucine I AUC } AUA } AUG } Methionine start codon M	ACU } Threonine T ACC } 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 } Alanine A GCC } GCA } GCG }	GAU } Aspartic acid D GAC } GAA } Glutamic acid E GAG }	GGU } Glycine G GGC } GGA } GGG }	U C A G

[Image source](#)

Primer on molecular biology

Three fundamental molecules:

1. DNA

Information storage.

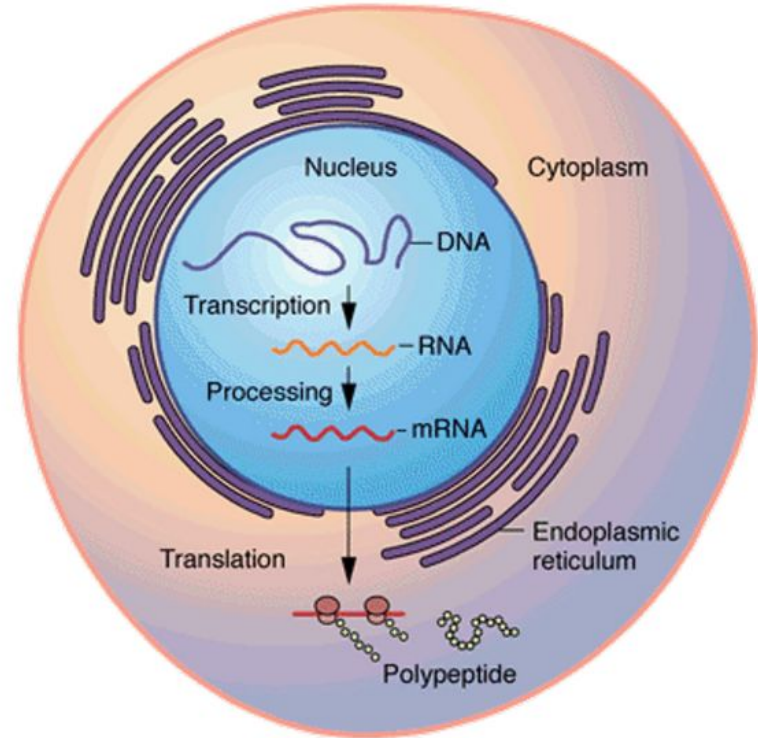
2. RNA

Old view: Mostly a “messenger”.

New view: Performs many important functions.

3. Protein

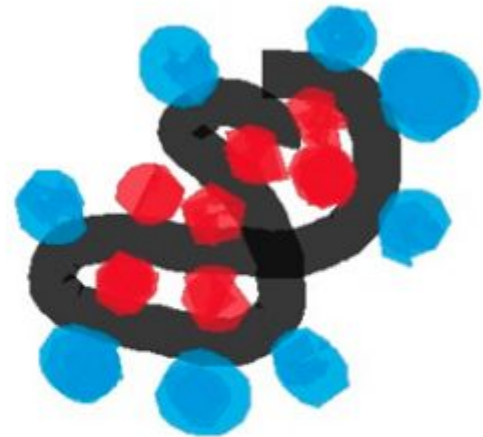
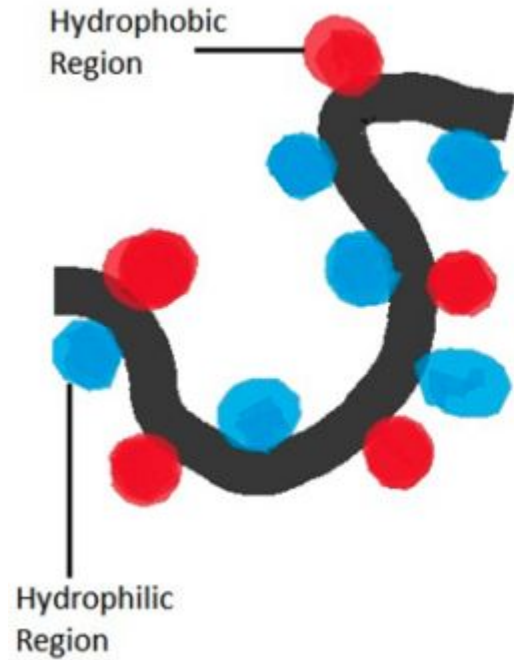
Perform most cellular functions
(biochemistry, signaling, control, etc.)



Summary: string transformation

- DNA = nucleotide sequence
 - Alphabet size = 4 (A,C,G,T)
- DNA to mRNA (single stranded)
 - Alphabet size = 4 (A,C,G,U)
- mRNA to amino acid sequence
 - Alphabet size = 20
- Amino acid sequence “folds” into 3-dimensional protein

Protein folding

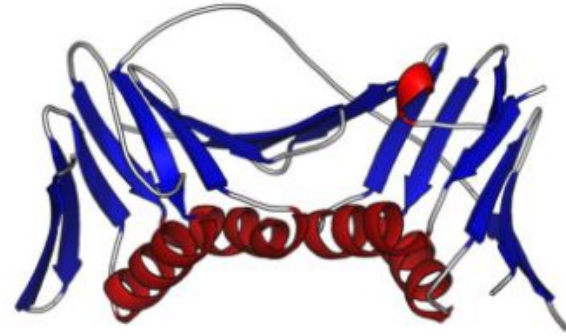


Protein in aqueous solution

Protein secondary structure

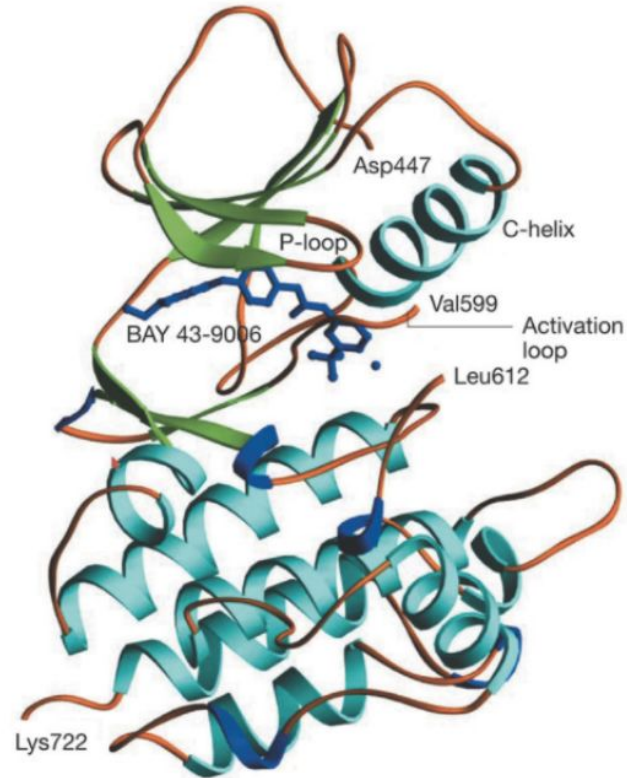


β -sheet

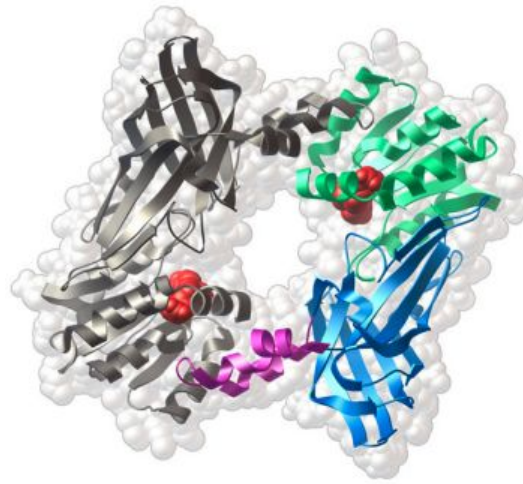


α -helix

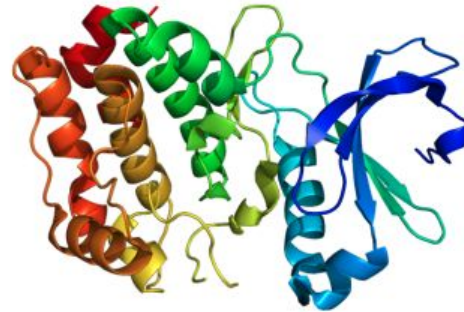
Tertiary structure



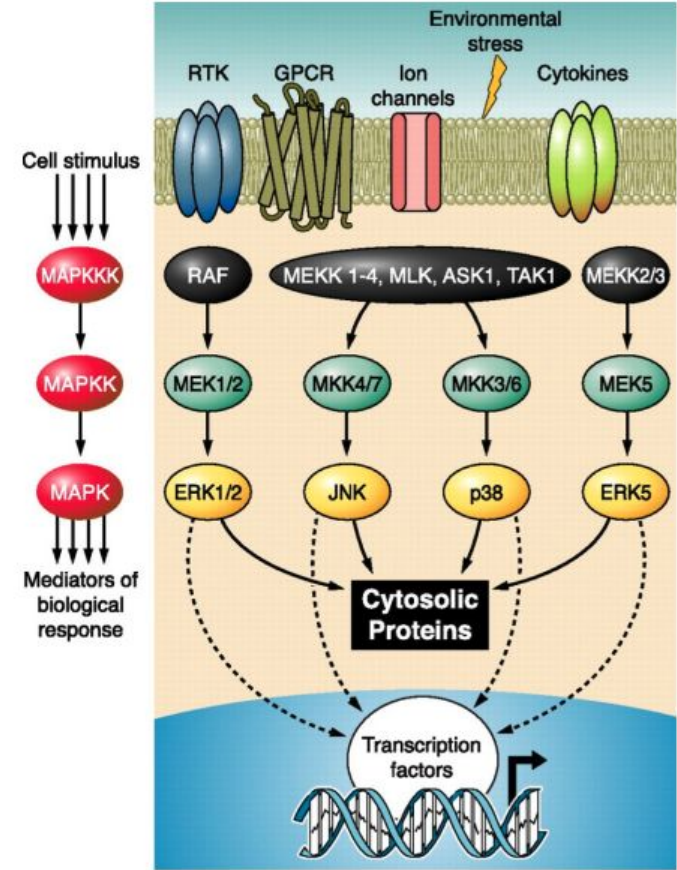
Protein function



Molecular switch

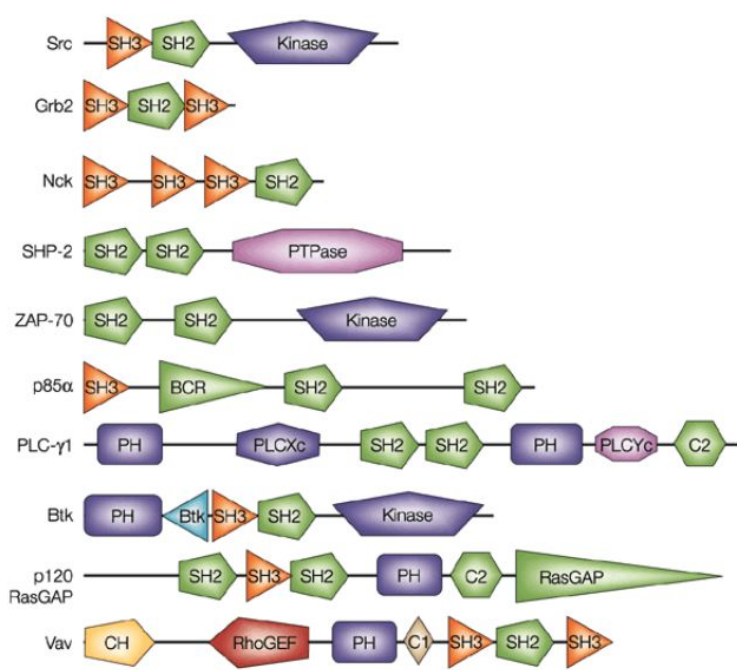


Enzyme

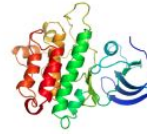


Signaling transduction

Protein domains



kinase



sh2



sh3

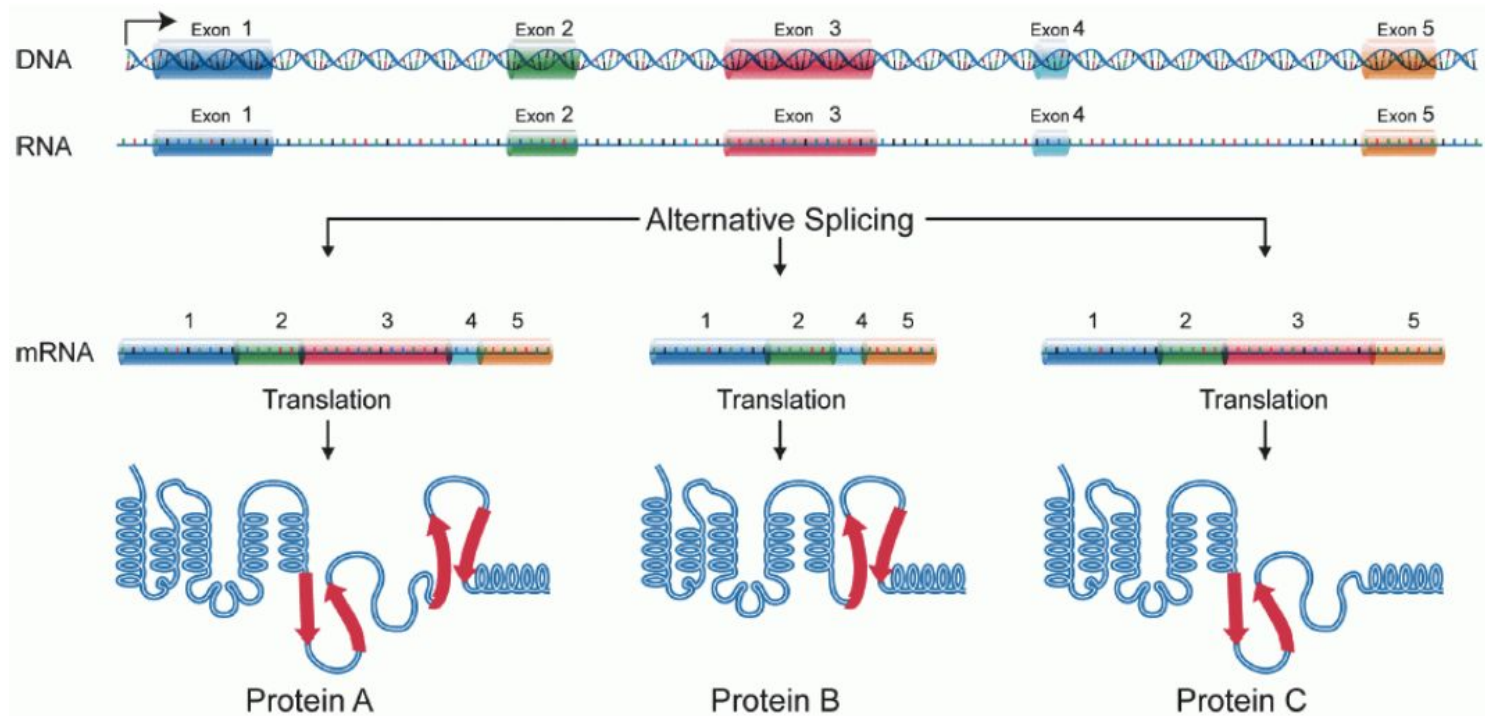


A protein domain is any identifiable longer contiguous **subsequence** of a protein that can **fold**, **function** and **exist independently** of the rest of the protein chain or structure.

TRUE or FALSE:

A gene is uniquely translated into a protein

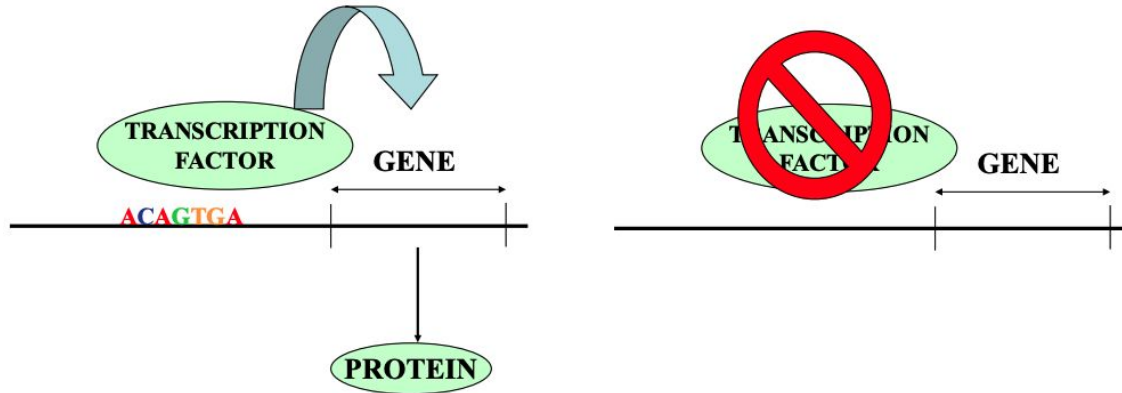
Gene structure



One gene can be translated into multiple different proteins

Gene expression

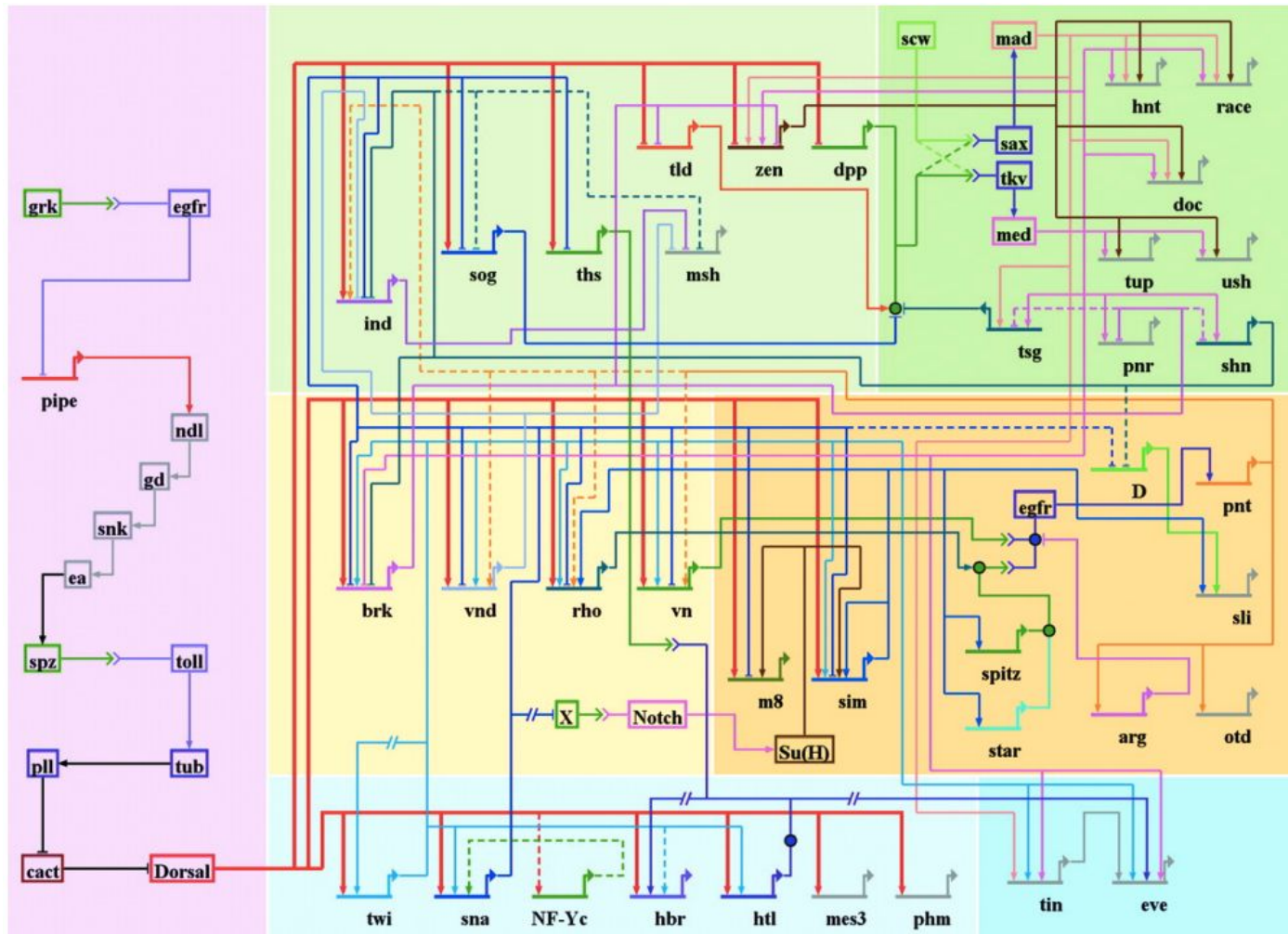
- Process of making a protein from a gene as template
- Transcription, then translation
- Can be regulated
 - The function of transcription factors is to regulate—turn on and off—genes in order to make sure that they are expressed in the desired cells at the right time and in the right amount.



Gene regulation

- Chromosomal activation/deactivation
- Transcriptional regulation
- Splicing regulation
- mRNA degradation
- mRNA transport regulation
- Control of translation initiation
- Post-translational modification

That is a “circuit” responsible for controlling gene expression



Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is ?% identical among individuals

Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is 99.9% identical among individuals

Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is 99.9% identical among individuals
- Human genome is 3 billion base-pairs (bp) long
- Genes and regulatory sequences make up ?% of human genome

Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is 99.9% identical among individuals
- Human genome is 3 billion base-pairs (bp) long
- Genes and regulatory sequences make up 5% of human genome

Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is 99.9% identical among individuals
- Human genome is 3 billion base-pairs (bp) long
- Genes and regulatory sequences make up 5% of human genome
- What's the rest doing?

Genome

- The entire sequence of DNA in a cell
- All cells have the same genome
 - All cells came from repeated duplications starting from initial cell (zygote)
- Human genome is 99.9% identical among individuals
- Human genome is 3 billion base-pairs (bp) long
- Genes and regulatory sequences make up 5% of human genome
- What's the rest doing?
 - We don't know for sure



Donald Knuth

Professor emeritus of Computer Science at Stanford University

Turing Award winner

“father of the analysis of algorithms.”

*“I can’t be as confident about computer science as I can about biology. **Biology easily has 500 years of exciting problems to work on.** It’s at that level.”*