



**Mansoura University**  
**Faculty of Computers and Information**  
**Department of Computer Science**  
**First Semester: 2020-2021**



**[MED121] Bioinformatics: An Introduction**  
**Grade: Third Year (Medical Informatics Program)**

**Sara El-Metwally, Ph.D.**  
**Faculty of Computers and Information,**  
**Mansoura University,**  
**Egypt.**

# AGENDA

- What is Bioinformatics.
- Bioinformatics Vs. Medical Informatics.
- Central Dogma of Life.
- DNA, Genomes.
- RNA, Transcriptomes.
- Proteins, Proteomics .
- Transcriptions and Translation.
- Omics era.

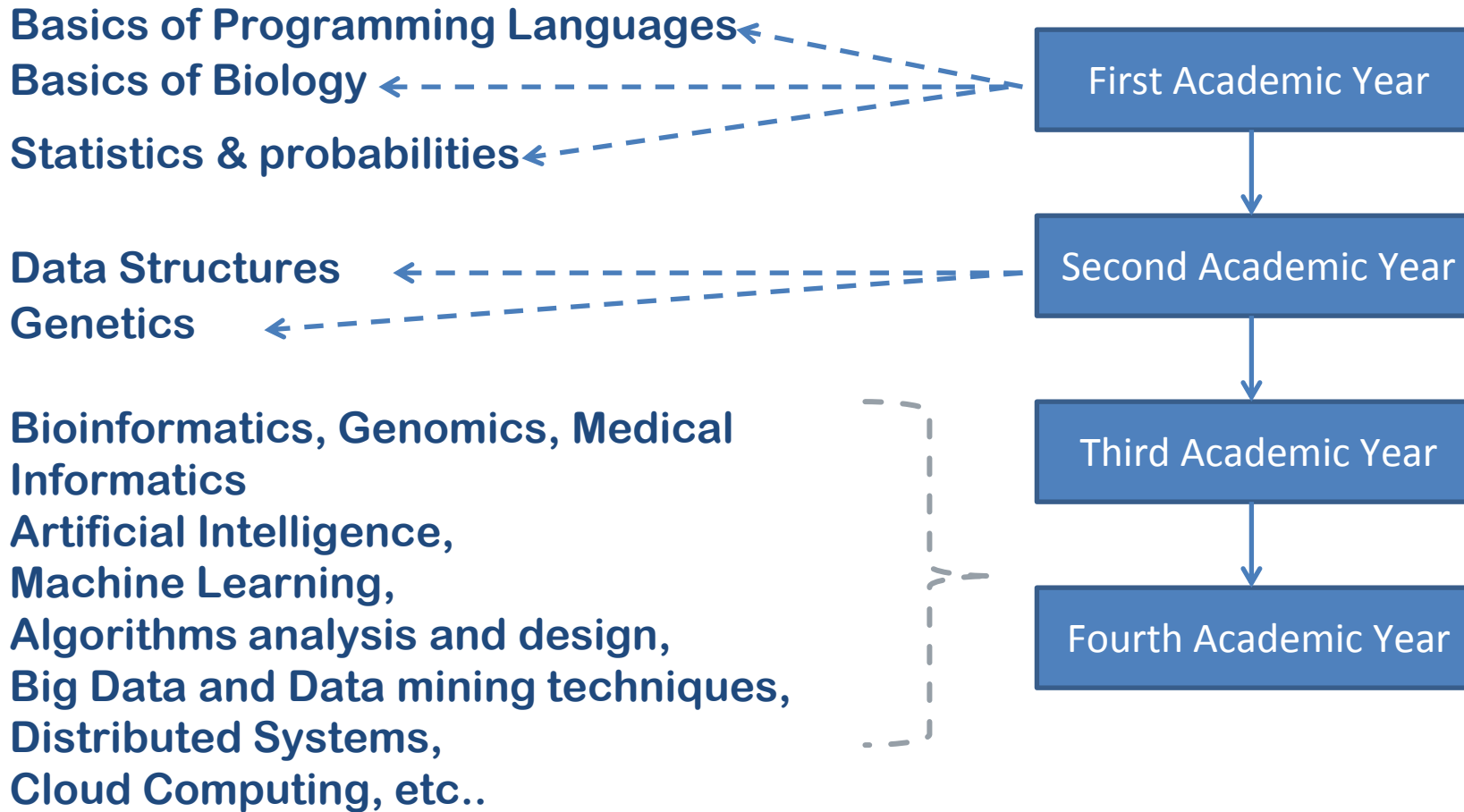
# COURSE OUTLINES

- **Course Meeting Time:** Wednesday, 2 :15 pm – 4.00 pm
- **Course Instructor:** Sara El-Metwally, PhD
- **Course TAs:** Eng. Aya Ayad, Eng. Mohamed Ashraf,
- **Course Labs:** Python
- **Course Grading:**
  - Midterm: 10%
  - Oral: 10%
  - Practical: 20%
  - Final: 60%
  - Any Projects are welcomed ! ( up to 40%)
  - Quizzes, attendance , etc..

# COURSE OBJECTIVES

- Bioinformatics course works as an interface between the computer science and biology and opens new opportunity to deal with different algorithms and tools for sequence comparisons, assembly, pattern matching, and efficient data structure to process a biological text.
- Understand the basic computer science terminologies regarding algorithm analysis and design, data analysis, and computer programming.
- Define the challenges that are facing the large scale biological data and think how to solve them through developing algorithms and bioinformatics tools.

# MANSOURA FCIS COURSE DEPENDENCY



# COURSE OUTPUT



<http://phdcomics.com/>

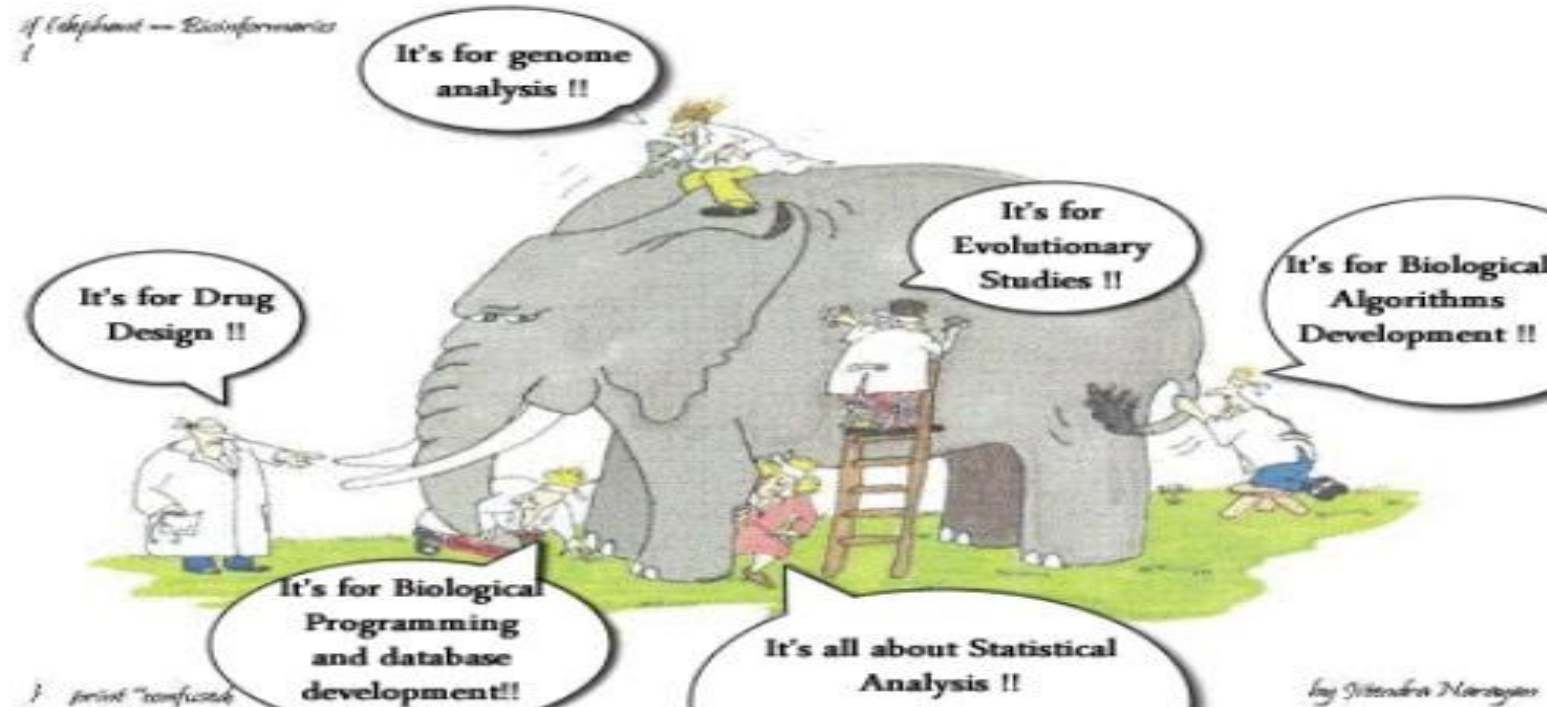
# LET'S START!

- What do you think “ **Bioinformatics** ” is? and why?



# LET'S START!

## What is Bioinformatics?



a professional in the pharmaceutical industry



a policeman worrying about DNA testing



a computer scientist developing bio-databases



a consumer concerned about GMOs (Genetically Modified Organisms)



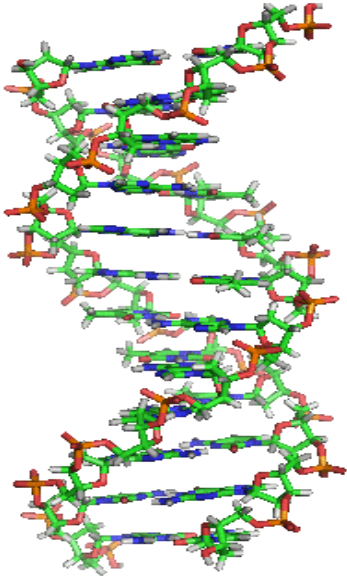
... ..



# WHAT IS BIOINFORMATICS?

- **Bioinformatics** is conceptualizing biology in terms of molecules and applying "informatics techniques" (derived from disciplines such as applied maths, computer science and statistics) to understand and organize the information associated with these molecules, on a large scale.
- **Bioinformatics** is computer aided biology! Also called computational biology.

# BIOINFORMATICS: A SIMPLE VIEW



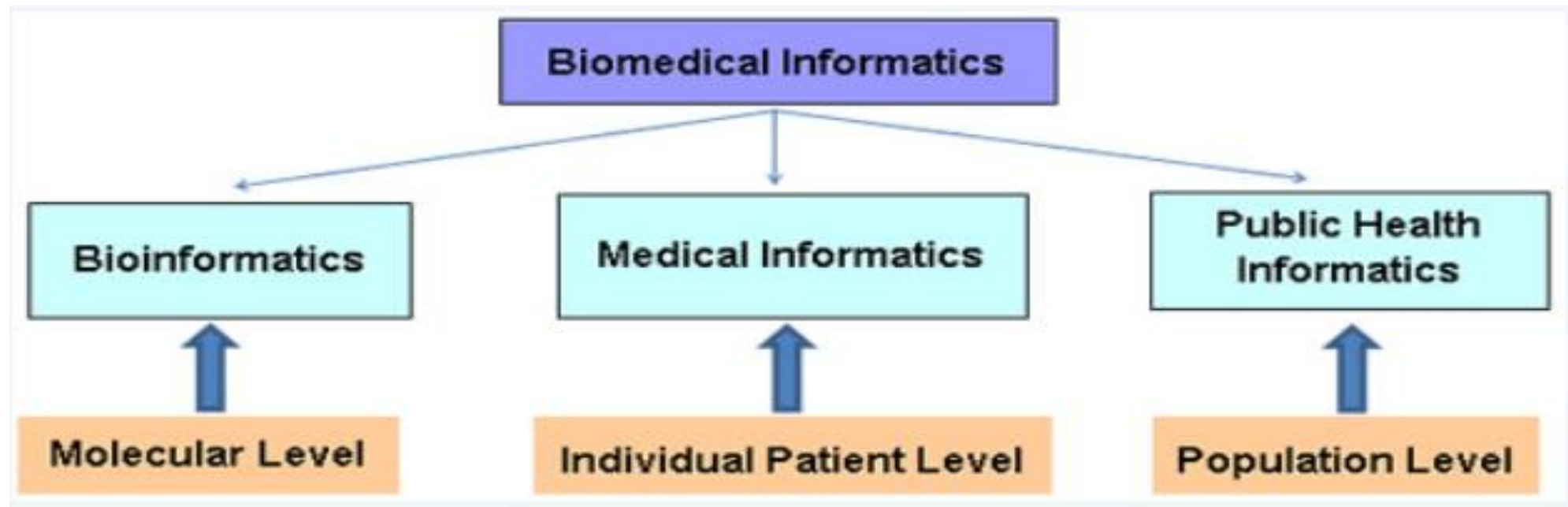
**Biological Data**

+



**Computer Calculations**

# BIOINFORMATICS, MEDICAL INFORMATICS & BIOMEDICAL INFORMATICS



<https://osteopathic.nova.edu/msbi/evolution.html>

# BIOINFORMATICS IS FUN!



# CENTRAL DOGMA OF LIFE

Central dogma of molecular biology



Central dogma of genomics

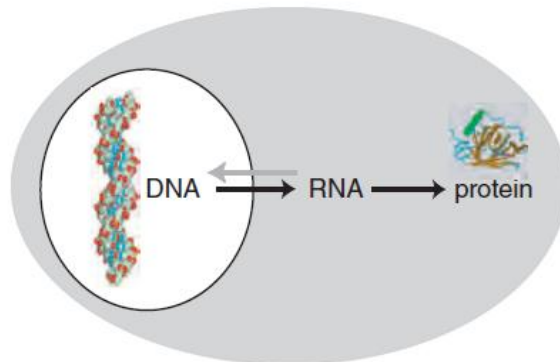


Photo credit :Functional Genomics Book by Jonathan Pevsner



Photo credit : <https://edu.t-bio.info/course/transcriptomics-1/>

# DNA

- **DNA stands** for Deoxyribo Nucleic Acid that carries the genetic instructions required for the development, functioning and reproduction of all known living organisms.
- In **eukaryotic organisms** (like animals, plants, and fungi), DNA is present in the nucleus of each cell.
- In **prokaryotic organisms** (single-celled organisms like bacteria and mitochondria), DNA is present in the cell's cytoplasm.

# DNA

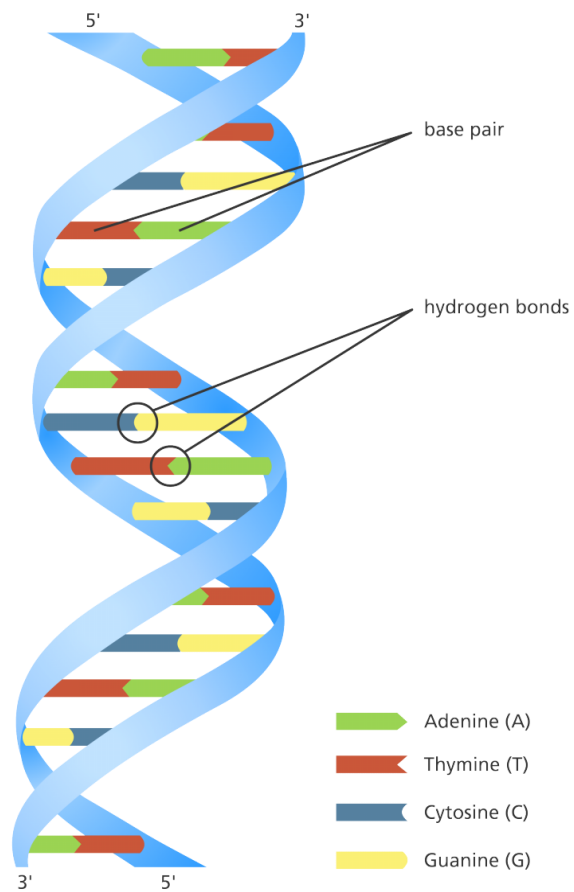


Photo Credit: <https://www.yourgenome.org/facts/what-is-dna>

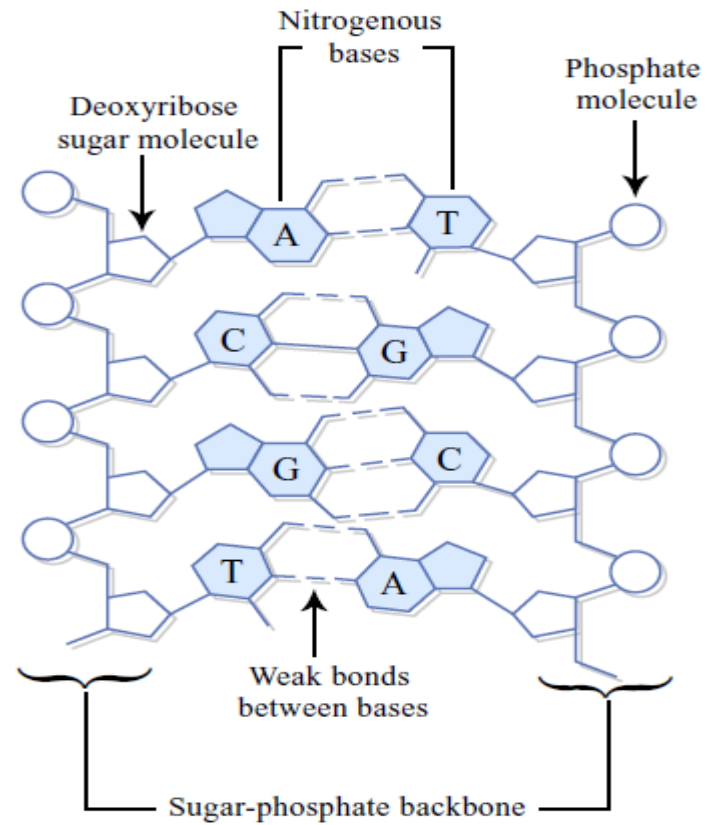


Photo Credit: <https://www.yourgenome.org/facts/what-is-dna>

# DNA

**TACTGTGACACTGT**

Y	Pyrimidine (C or T)
R	Purine (A or G)
W	Weak (A or T)
S	Strong (G or C)
K	Keto (T or G)
M	Amino (C or A)
D	A, G, T (not C - remember as after C)
V	A, C, G (not T - remember as after T/U - We'll get to "U" soon)
H	A, C, T (not G - remember as after G)
B	C, G, T (not A - remember as after A)
N	Any base
-	Gap

Photo credit : BIOSTAR Handbook



# DNA

sense/forward/+ /watson/code strand

5'

ATGACACTGTGACA

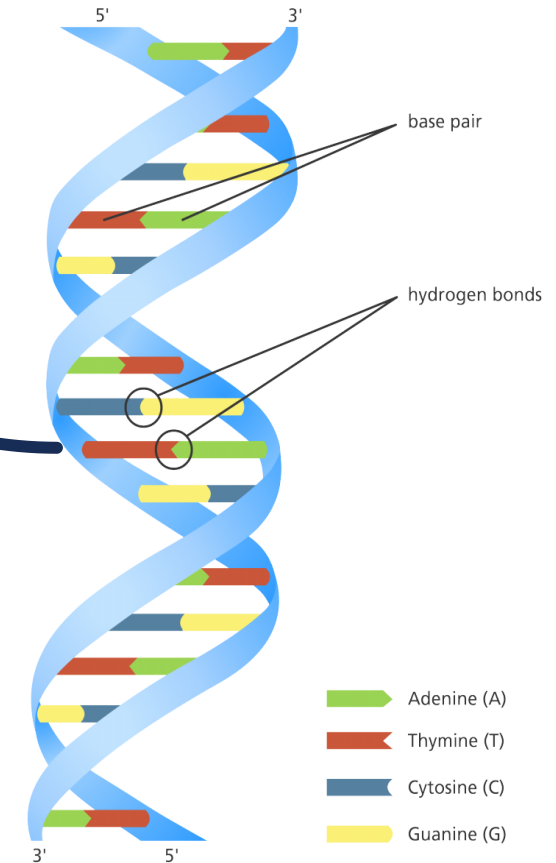
TACTGTGACACTGT

3'

anti-sense/reverse/- /crick/template strand

3'

5'



# GENOMES

- Genome refers to the complete set of genetic information in an organism, it is your DNA!
- Human genome contains approximately 3 billion base pairs, which reside in the 23 chromosomes.
- Mosquitoes have 3 pairs, i.e. tiger mosquito has 1,967, 000, 000 bp = 1,967 Mbp.

# GENOME SIZES

**Comparison of Genome Size in Different Organisms**






Species	<i>T2 phage</i>	<i>Escherichia coli</i>	<i>Drosophila melanogaster</i>	<i>Homo sapiens</i>	<i>Paris japonica</i>
Genome Size	170,000 bp	4.6 million bp	130 million bp	3.2 billion bp	150 billion bp
Common Name	 Virus	 Bacteria	 Fruit fly	 Human	 Canopy Plant

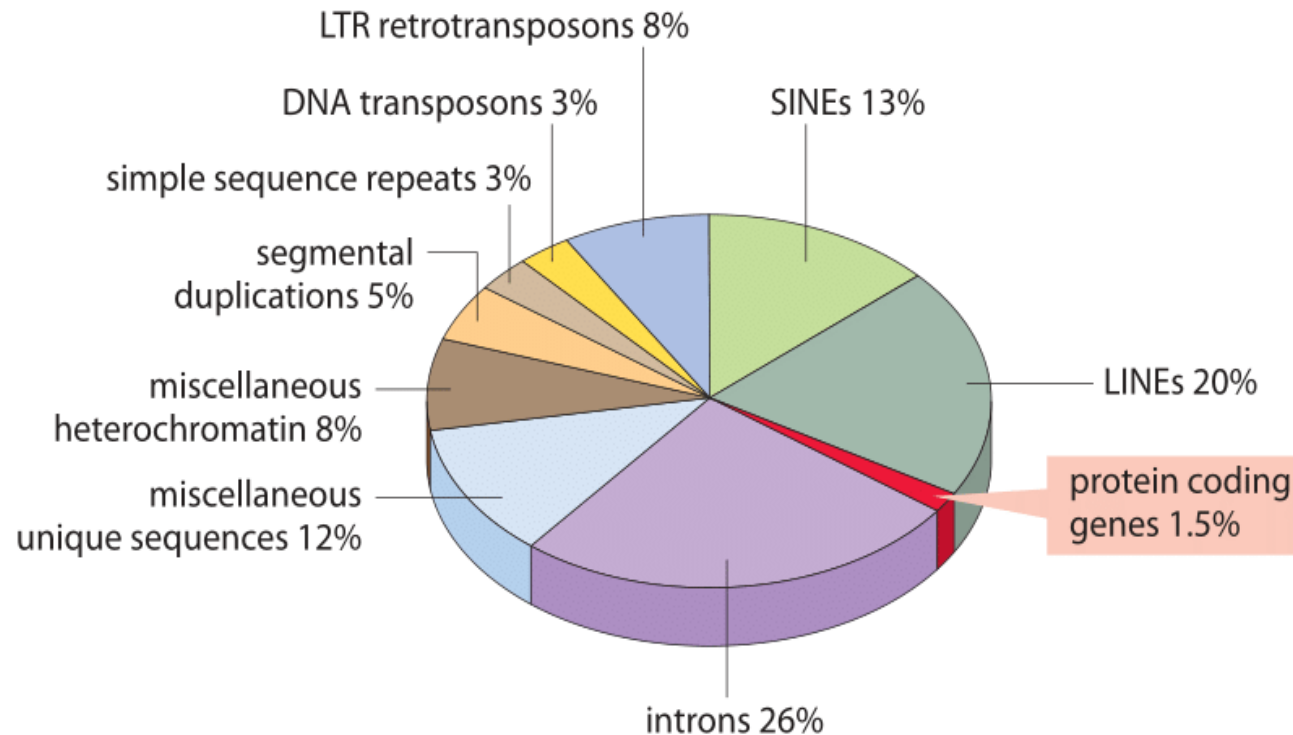
Photo Credit: <https://ib.bioninja.com.au/standard-level/topic-3-genetics/32-chromosomes/genome-size.html>

Organism	Completion date	Size	Description
phage phiX174	1978	5,368 bp	1st viral genome
human mtDNA	1980	16,571 bp	1st organelle genome
lambda phage	1982	48,502 bp	important virus model
HIV	1985	9,193 bp	AIDS retrovirus
<i>H. influenzae</i>	1995	1,830 Kb	1st bacterial genome
<i>M. genitalium</i>	1995	580 Kb	smallest bacterial genome
<i>S. cerevisiae</i>	1996	12.5 Mb	1st eukaryotic genome
<i>E. coli</i> K12	1997	4.6 Mb	bacterial model organism
<i>C. trachomatis</i>	1998	1,042 Kb	internal parasite of eukaryotes
<i>D. melanogaster</i>	2000	180 Mb	fruit fly, model insect
<i>A. thaliana</i>	2000	125 Mb	thale cress, model plant
<i>H. sapiens</i>	2001	3,000 Mb	human
SARS	2003	29,751 bp	coronavirus

Photo Credit: <http://book.bionumbers.org/how-many-genes-are-in-a-genome/>

# GENES

## main components of the human genome



	Organism	# of protein-coding genes	# of genes naïve estimate: (genome size / 1000)	BNID
viruses	HIV 1	9	10	105769
	Influenza A virus	10-11	14	105767
	Bacteriophage λ	66	49	105770
	Epstein Barr virus	80	170	103246
prokaryotes	<i>Buchnera sp.</i>	610	640	105757
	<i>T. maritima</i>	1,900	1,900	105766
	<i>S. aureus</i>	2,700	2,900	105500
	<i>V. cholerae</i>	3,900	4,000	105760
	<i>B. subtilis</i>	4,400	4,200	111448
	<i>E. coli</i>	4,300	4,600	105443
	<i>S. cerevisiae</i>	6,600	12,000	105444
eukaryotes	<i>C. elegans</i>	20,000	100,000	101364
	<i>A. thaliana</i>	27,000	140,000	111380
	<i>D. melanogaster</i>	14,000	140,000	111379
	<i>F. rubripes</i>	19,000	400,000	111375
	<i>Z. mays</i>	33,000	2,300,000	110565
	<i>M. musculus</i>	20,000	2,800,000	100308
	<i>H. sapiens</i>	21,000	3,200,000	100399, 111378
	<i>T. aestivum</i> (hexaploid)	95,000	16,800,000	105448, 102713

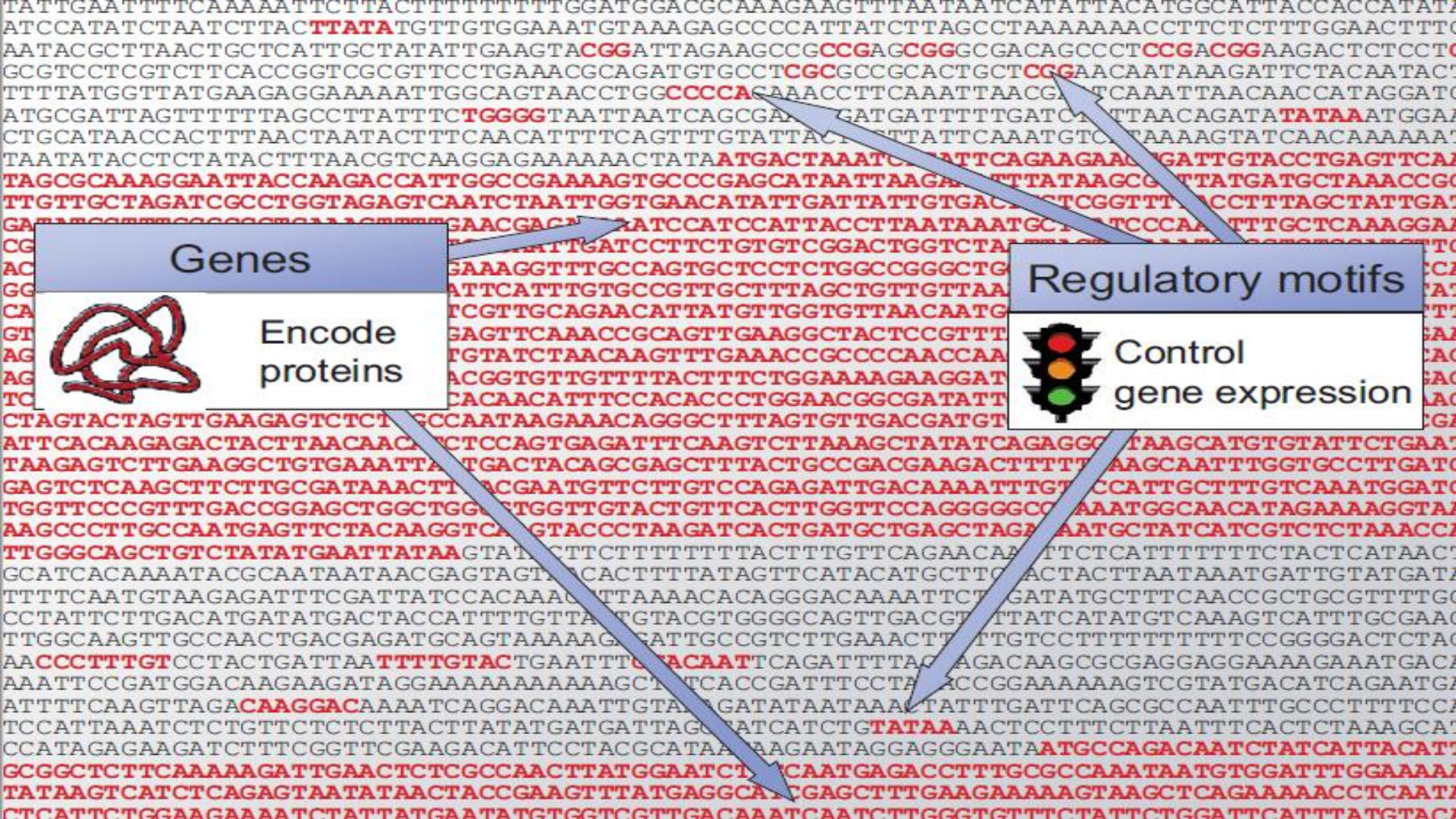
# GENES

- Gene is a **basic unit of heredity**, it is a sequence of bases that carries the instructions to make a particular protein.
- There are an estimated **19,000 - 20,000** human protein-coding genes.
- Not all of the DNA in a genome encodes protein, i.e. **~ 1.5 %** of human genome are protein coding regions.



ATCCATATCTAATCTTACTTATATGTTGTGGAAATGTAAAGAGCCCCATTATCTTAGCCTAAAAAACCTTCTCTTTGGAACTTTCT  
AATACGCTTAACCTGCTCATTGCTATATTGAAGTACGGATTAGAAGCCGCCGAGCGGGCGACAGCCCTCCGACGGAAGACTCTCCTC  
GCGTCCTCGTCTTCACCGGTCGCGTTCCTGAAACGCAGATGTGCCTCGCGCCGCACTGCTCCGAACAATAAAGATTCTACAATAC  
TTTTATGGTTATGAAGAGGAAAAATTGGCAGTAACCTGGCCCCACAAACCTTCAAATTAACGAATCAAATTAACAACCATAGGATC  
ATGCGATTAGTTTTTTTAGCCTTATTTCTGGGGTAATTAATCAGCGAAGCGATGATTTTTTGATCTATTAACAGATATATAAATGGA  
CTGCATAACCACCTTTAACTAATACTTTCAACATTTTCAGTTTGTATTACTTCTTATTCAAATGTCATAAAAAGTATCAACAAAAAA  
TAATAATACCTCTATACTTTAACGTCAAGGAGAAAAAACTATAATGACTAAATCTCATTCAGAAGAAGTGATTGTACCTGAGTTCAA  
TAGCGCAAAGGAATTACCAAGACCATTGGCCGAAAAGTGCCCGAGCATAATTAAGAAATTTATAAGCGCTTATGATGCTAAACCGC  
TTGTTGCTAGATCGCCTGGTAGAGTCAATCTAATTGGTGAACATATTGATTATTGTGACTTCTCGGTTTTACCTTTAGCTATTGAT  
GATATGCTTTTGCGCCGTCAAAGTTTTGAACGAGAAAAATCCATCCATTACCTTAATAAATGCTGATCCCAAATTTGCTCAAAGGAA  
CGATTTGCCGTTGGACGGTCTTATGTCACAATTGATCCTTCTGTGTCGGACTGGTCTAATTACTTTAAATGTGGTCTCCATGTTCT  
ACTCTTTTCTAAAGAAACTTGCACCGGAAAGGTTTGCCAGTGCTCCTCTGGCCGGGGCTGCAAGTCTTCTGTGAGGGTGATGTACCA  
GGCAGTGGATTGTCTTCTTCGGCCGCATTCATTTGTGCGGTTGCTTTAGCTGTTGTTAAAGCGAATATGGGGCCCTGGTTATCATAT  
CAAGCAAAATTTAATGCGTATTACGGTCGTTGCAGAACATTATGTTGGTGTAAACAATGGCGGTATGGATCAGGCTGCCTCTGT  
GTGAGGAAGATCATGCTCTATACGTTGAGTTCAAACCGCAGTTGAAGGCTACTCCGTTTTAAATTTCCGCAATTAAAAAACCATGA  
AGCTTTGTATTATTGCGAACACCCCTTGTTGTATCTAACAAAGTTTGAAACCGCCCCAACCAACTATAATTTAAGAGTGGTAGAAGTCA  
AGCTGCAAATGTTTTAGCTGCCACGTACGGTGTTGTTTTACTTTCTGGAAAAGAAGGATCGAGCACGAATAAAGGTAATCTAAGAC  
TCATGAACGTTTTATTATGCCAGATATCACAAACATTTCCACACCCTGGAACGGCGATATTGAATCCGGGCATCGAACGGTTAACAAAC  
CTAGTACTAGTTGAAGAGTCTCTCGCCAATAAGAAACAGGGCTTTAGTGTTGACGATGTCGCACAATCCTTGAAATTGTTCTCGCGA  
ATTCACAAGAGACTACTTAACAACATCTCCAGTGAGATTTCAAGTCTTAAAGCTATATCAGAGGGGCTAAGCATGTGTATTTCTGAAT  
TAAGAGTCTTGAAGGCTGTGAAATTAATGACTACAGCGAGCTTTACTGCCGACGAAGACTTTTTCAAGCAATTTGGTGCCTTGATC  
GAGTCTCAAGCTTCTTGCGATAAACTTTACGAATGTTCTTGTTCCAGAGATTGACAAAATTTGTTCCATTGCTTTGTCAAATGGATC  
TGGTTCCTCGTTTGACCGGAGCTGGCTGGGGTGGTTGTACTGTTCACTTGGTTCAGGGGGGCCCAAATGGCAACATAGAAAAGGTAA  
AAGCCCTTGCCAATGAGTTCTACAAGGTCAAGTACCCTAAGATCACTGATGCTGAGCTAGAAAATGCTATCATCGTCTCTAAACCA  
TTGGGCAGCTGTCTATATGAATTATAAGTATACTTCTTTTTTTTTACTTTGTTCAGAACAACCTTCTCATTTTTTTTTCTACTCATAACT  
GCATCACAAAATACGCAATAATAACGAGTAGTAACACTTTTATAGTTCATACATGCTTCAACTACTTAATAAATGATTGTATGATA  
TTTTCAATGTAAGAGATTTTCGATTATCCACAAACTTTAAAACACAGGGACAAAATTTCTTGATATGCTTTCAACCGCTGCGTTTTG  
CCTATTCTTGACATGATATGACTACCATTTTTGTTATTGTACGTGGGGCAGTTGACGTCTTATCATATGTCAAAGTCATTTGCGAAC  
TTGGCAAGTTGCCAACTGACGAGATGCAGTAAAAAGAGATTGCCGTCCTTGAAACTTTTTTGTCCTTTTTTTTTTTTCCGGGGACTCTAC  
AACCCCTTTGTCTACTGATTAAATTTTGTTACTGAATTTGGACAATTCAGATTTTAGTAGACAAGCGCGAGGAGGAAAAGAAATGACA





Genes



Encode  
proteins

Regulatory motifs



Control  
gene expression



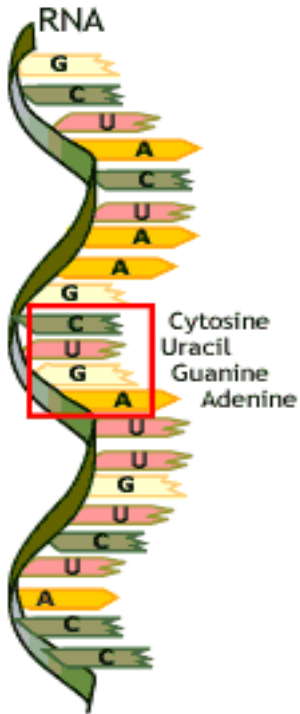
# GENOTYPE **VS.** PHENOTYPE

- **Genotype** is the set of genes in an organism.
- **Phenotype** is the set of all observable characteristics — which are influenced both by its genotype and by the environment.





# RNA



- RNA is like DNA except that:
  1. Backbone is a little different.
  2. Often single stranded.
  3. The base Uracil (U) is used in place of Thymine (T).
- RNA as a string of four alphabets:

**AUGACACUGUGACA**

# DNA vs. RNA

## DNA vs. RNA

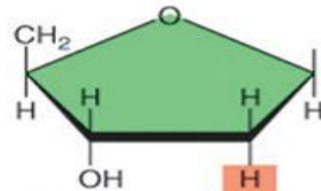


Double-stranded

b.

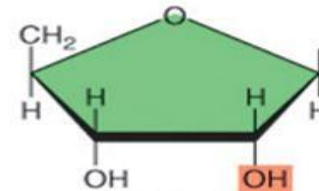


Generally single-stranded



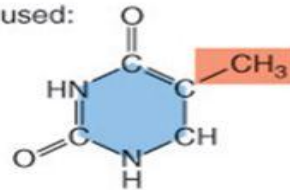
Deoxyribose as the sugar

c.



Ribose as the sugar

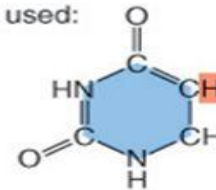
Bases used:



Thymine (T)  
Cytosine (C)  
Adenine (A)  
Guanine (G)

d.

Bases used:

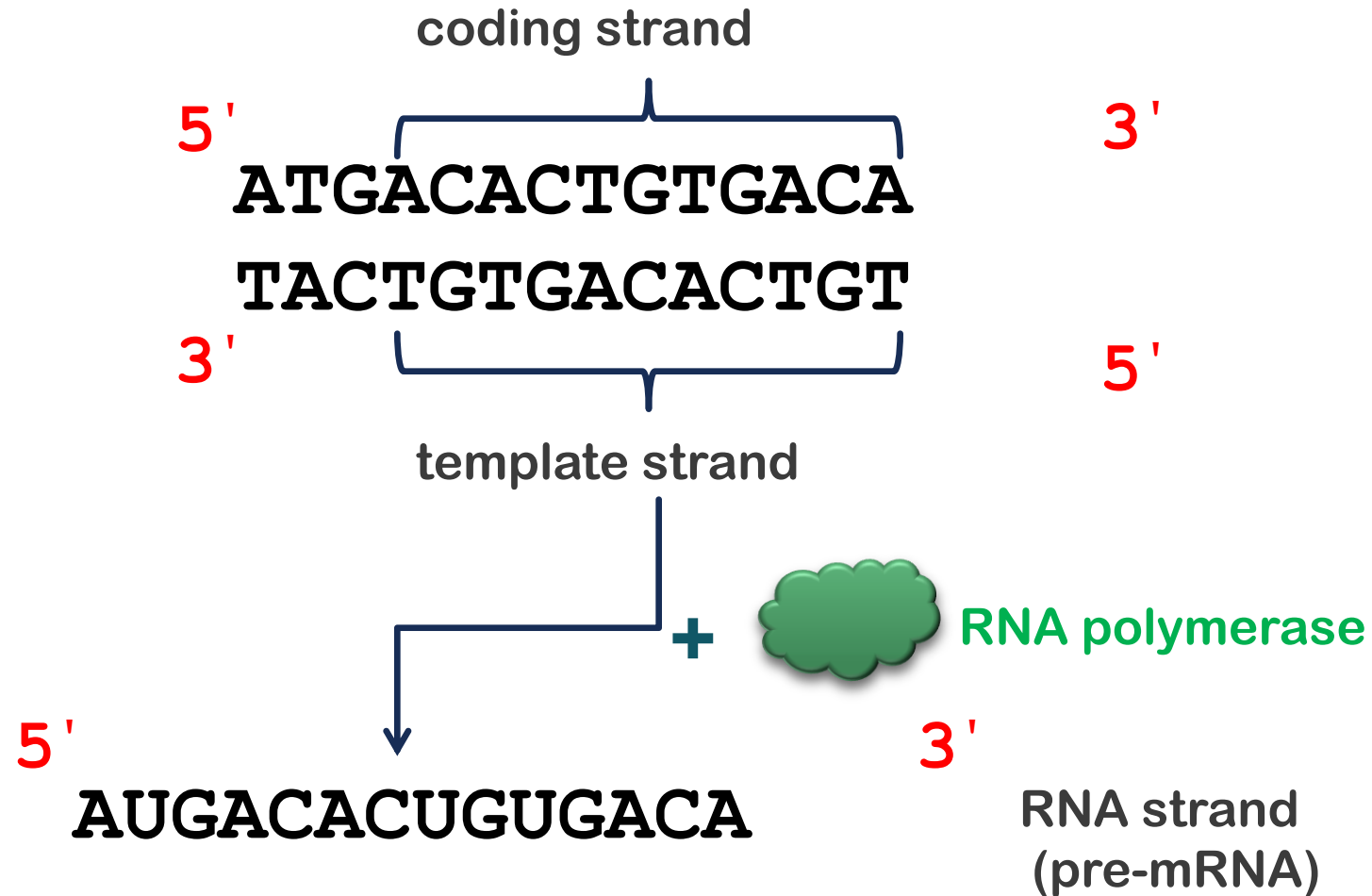


Uracil (U)  
Cytosine (C)  
Adenine (A)  
Guanine (G)

# TRANSCRIPTION

- Transcription is the process of transforming information in a gene in a DNA strand to an RNA strand, i.e. making RNA from a DNA template.
- Genetic information is carried on only one of the two strands of the DNA, which known as a coding strand.

# TRANSCRIPTION



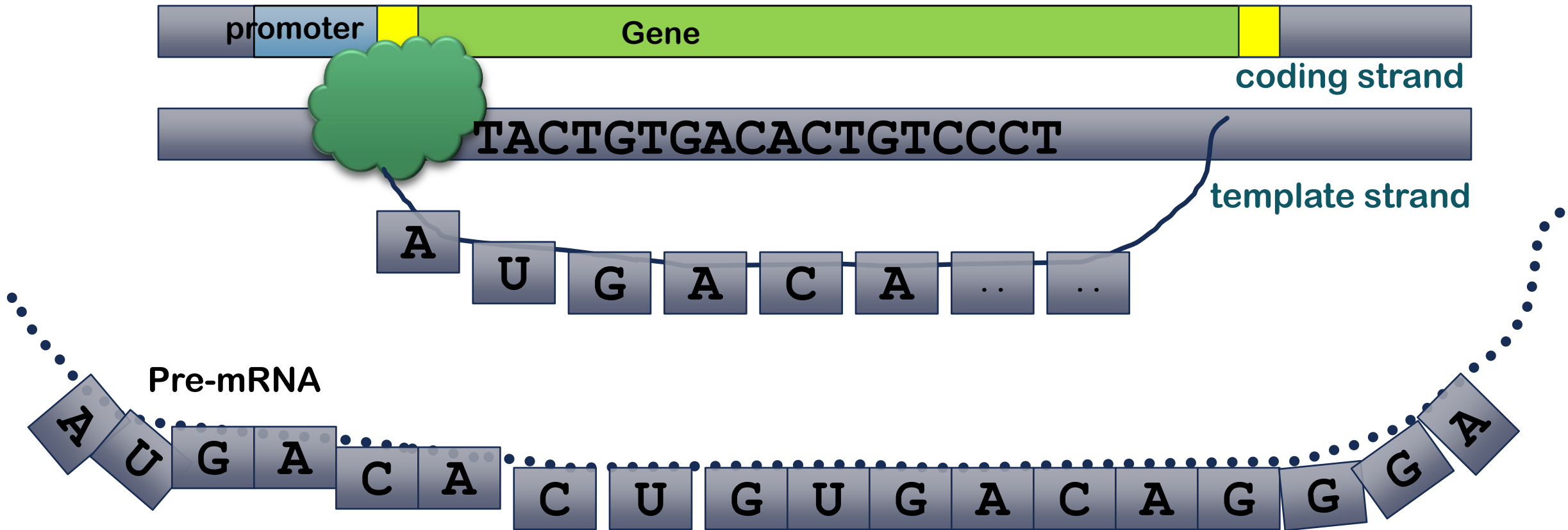
# TRANSCRIPTION

- RNA polymerase recognizes short sequence of bases called promoter sequence and bind to it.

RNA polymerase



# TRANSCRIPTION



# RNA PROCESSING

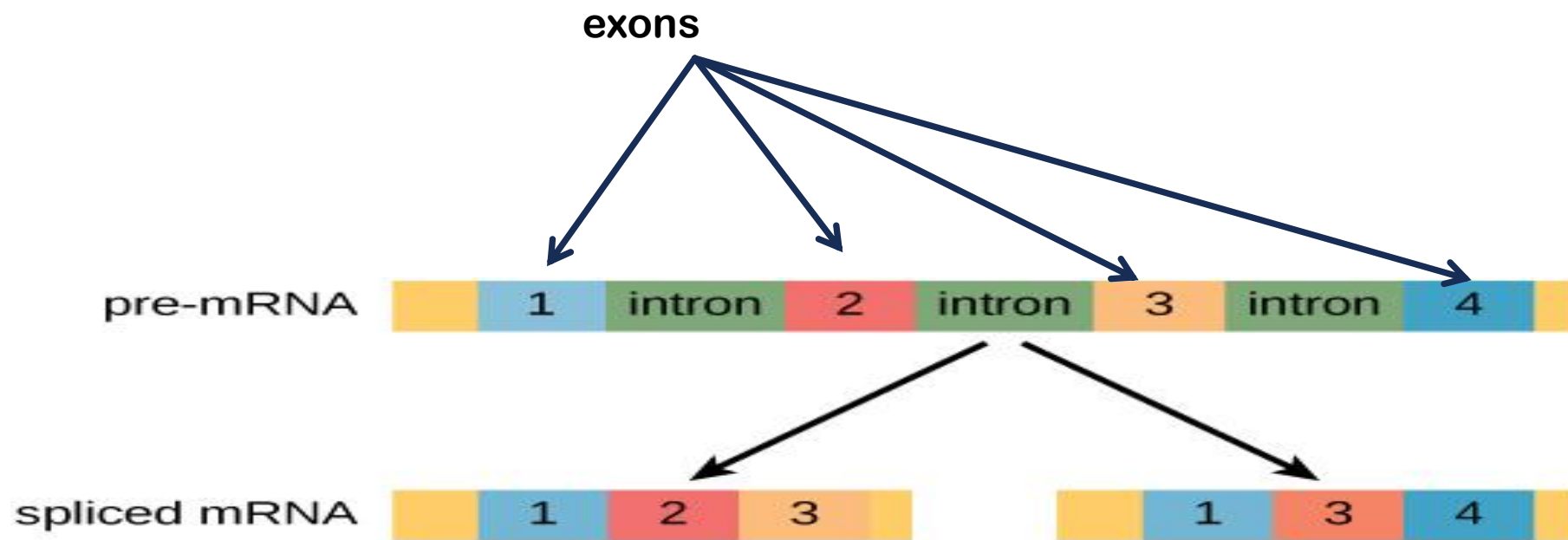


Photo credit: [https://en.wikipedia.org/wiki/Alternative\\_splicing](https://en.wikipedia.org/wiki/Alternative_splicing)

# ALTERNATIVE SPLICING

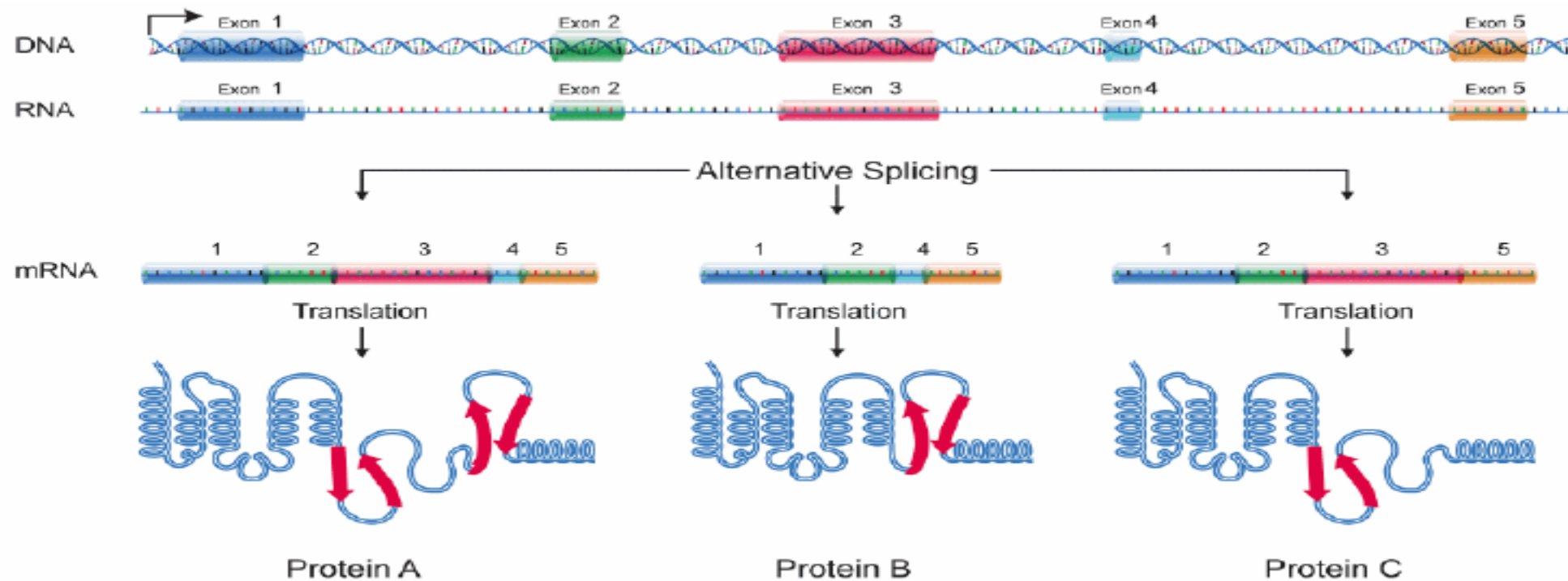


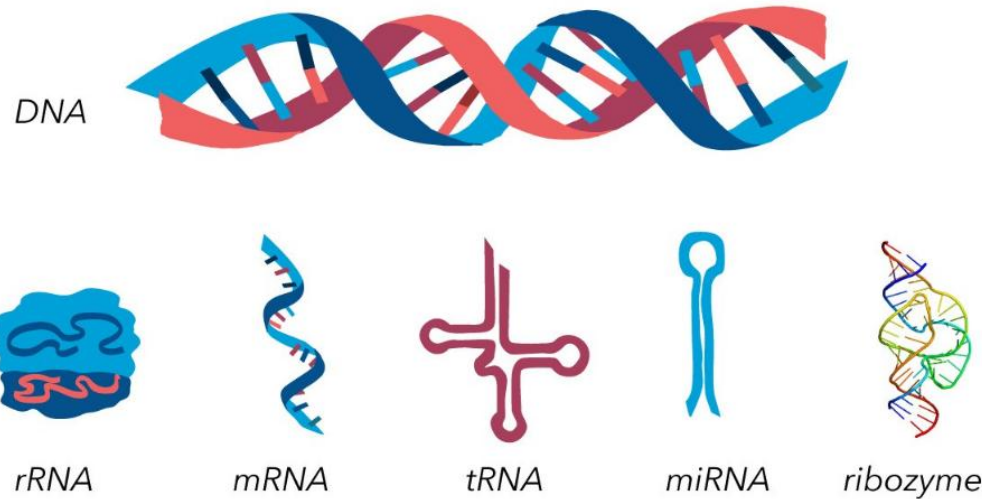
Photo credit: [https://en.wikipedia.org/wiki/Alternative\\_splicing](https://en.wikipedia.org/wiki/Alternative_splicing)



# NOTES

## ■ Other DNA products:

1. ribosomal RNA (**rRNA**), which includes major constituents of ribosomes.
2. transfer RNAs (**tRNAs**), which carry amino acids to ribosomes.
3. micro RNAs (**miRNAs**), which play an important regulatory role in various plants and animals.



# TRANSLATION

- **Ribosome** is the machine that synthesizes proteins from mRNA.
- Proteins are molecules composed of one or more **polypeptides**.
- a **polypeptide** is a polymer composed of **amino acids**.
- Cells build their proteins from **20 different amino acids**.
- A polypeptide can be thought of as a string composed from a 20-character alphabet.

# TRANSLATION

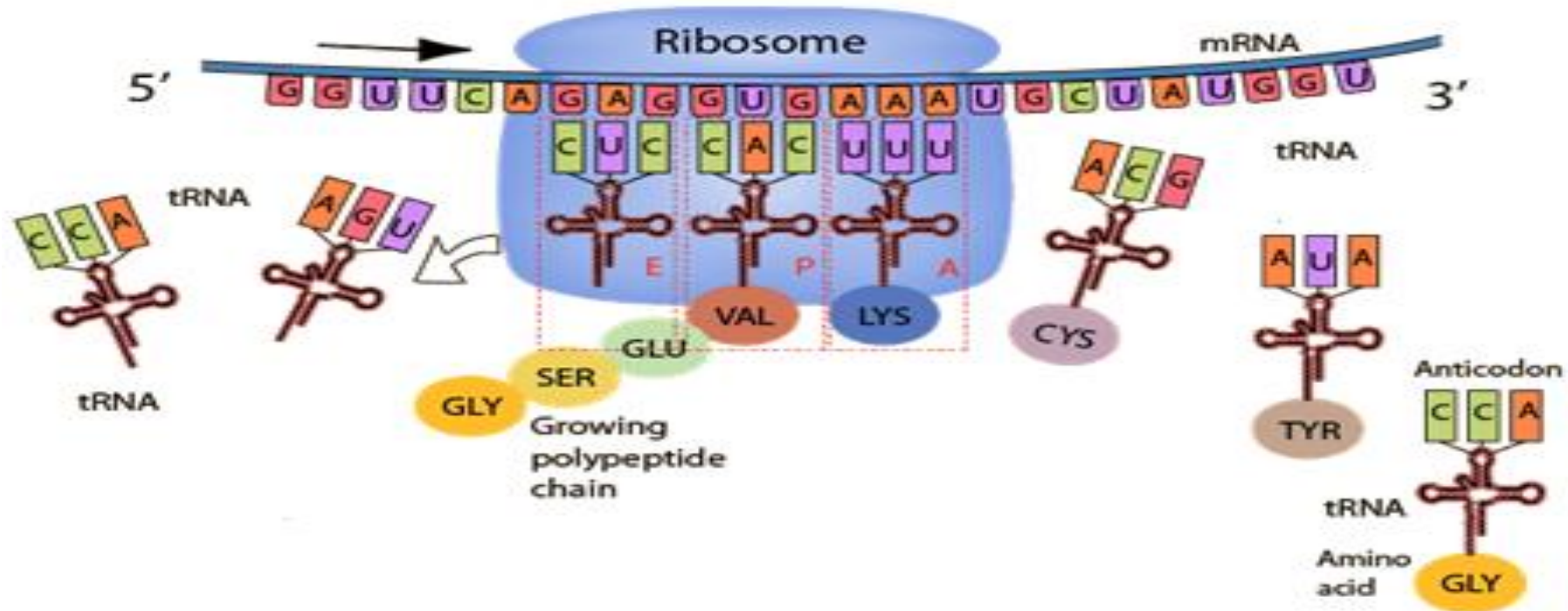


Photo credit: <http://hyperphysics.phy-astr.gsu.edu/hbase/Organic/translation.html>

# AMINO ACIDS

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

# PROTEIN SEQUENCE

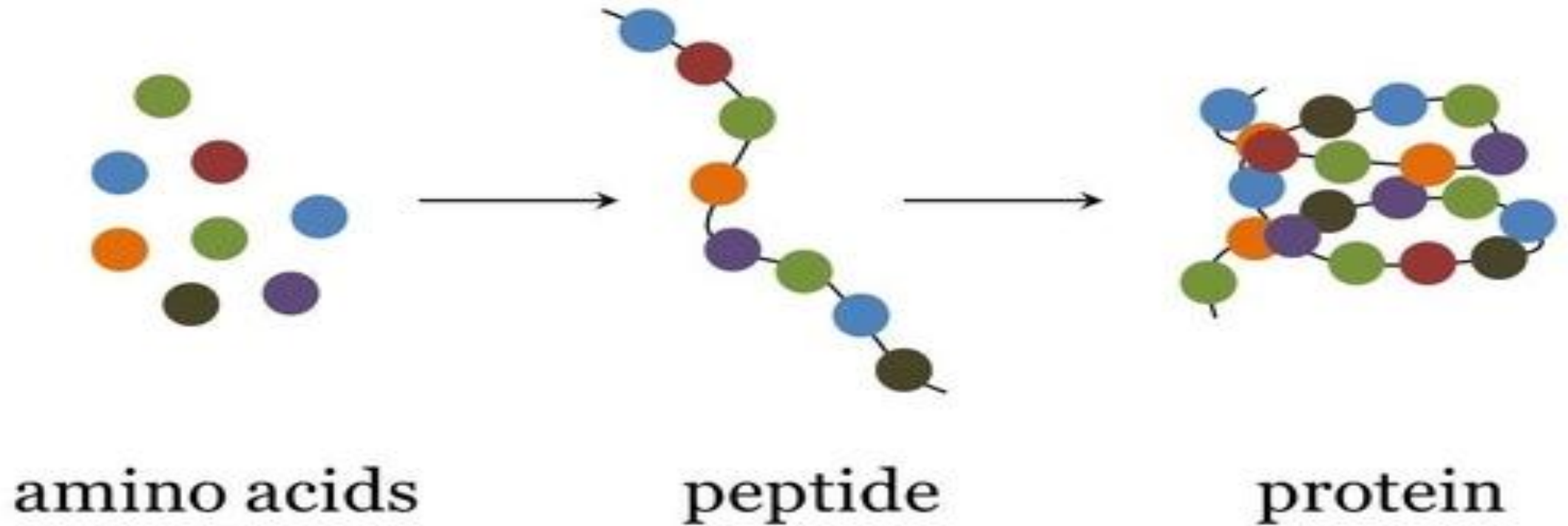


Photo Credit: <https://www.peptidesciences.com/information/peptides-vs-proteins/>

# PROTEIN SEQUENCE



**FLIVQNATPSVIE**

# VIDEO

- <https://www.youtube.com/watch?v=gG7uCskUOrA>

# OMICS

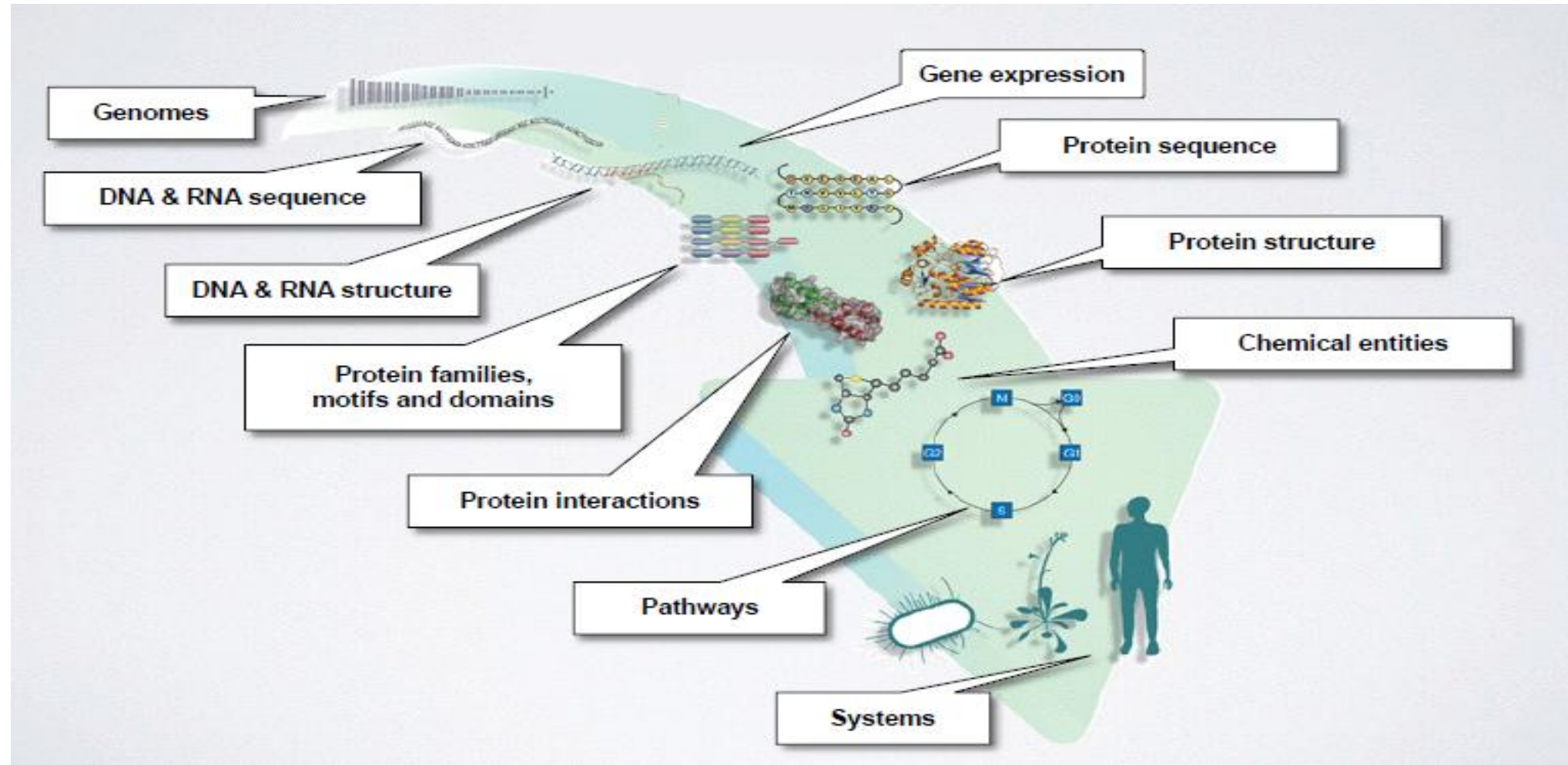
- Genomics: DNA.
- Transcriptomics: RNA.
- Proteomics: Proteins.
- Metagenomics: is a new research area focused on the analysis of mixture of DNA sequences extracted from different organisms i.e. viral, bacterial, or eukaryotic that are living together in a symbiotic community.
- Epigenomics is the study of the complete set of epigenetic modifications on the genetic material of a cell, known as the epigenome. The epigenome is made up of chemical compounds and proteins that can attach to DNA and direct such actions as turning genes on or off, controlling the production of proteins in particular cells.



# OMICS ERA (LARGE SCALE BIOLOGY)

Name	Study of
Genomics	entire genome of an organism
Transcriptomics	expressed genes
Exomics	coding sequences
Proteomics	proteins within an organism
Metagenomics	mixture of DNA sequences from different organisms
Metabolomics	metabolites within an organism
Interactomics	interactions between nucleotides, proteins and metabolites
Connectomics	neural pathways in the brain
Pharmacogenomics	application of genomics to pharmacology
Phenomics	observable phenotypes
Physiomics	functional behaviour of an organism
Exposomics	organism's environment

# TYPES OF BIOINFORMATICS DATA (**BIOLOGICAL MOLECULES**)





**Thank you!**