


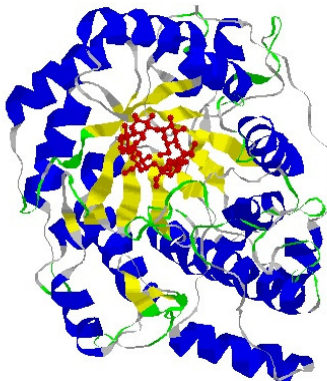
Introduction to Bioinformatics

Sami Khuri
Department of Computer Science
San José State University
San José, California, USA
khuri@cs.sjsu.edu
www.cs.sjsu.edu/faculty/khuri

@2010 Sami Khuri



Biology Review



- DNA
- RNA
- Proteins
- Central Dogma
- Transcription
- Translation

Last revised: July 18, 2010

@2010 Sami Khuri

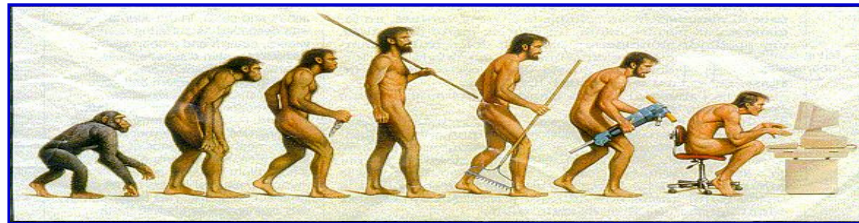
Understanding Biology I

Nothing in biology makes sense, except in the light of **evolution**.

Dobzhansky, Russian geneticist (1900-75)



"I've only just bought this bronze stuff and you're telling me I ought to upgrade to iron?"



©2010 Sami Khuri

Understanding Biology II



- All organisms are (probably) **evolutionarily** related to each other; i.e., descended from a single common ancestor.
- **Living organisms** are “imperfect replication machines”.
- Biology is not an exact science.

©2010 Sami Khuri

“We are our Proteins” Doolittle

Source: George Poste

Limitless Diversity From Combinatorial Assemblies of Limited Building Blocks

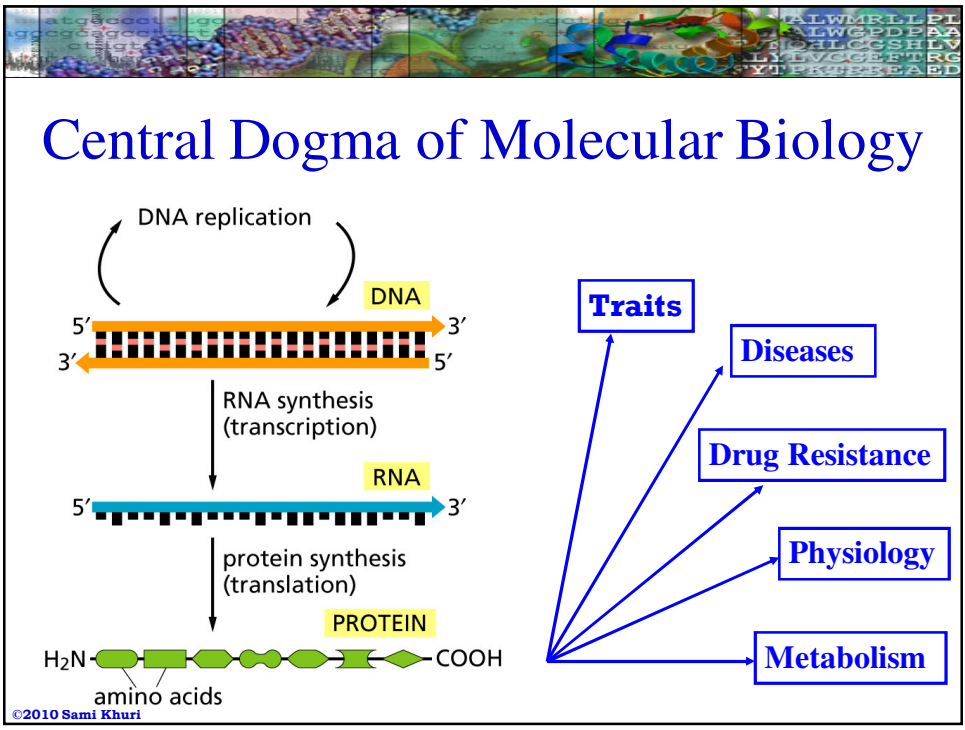
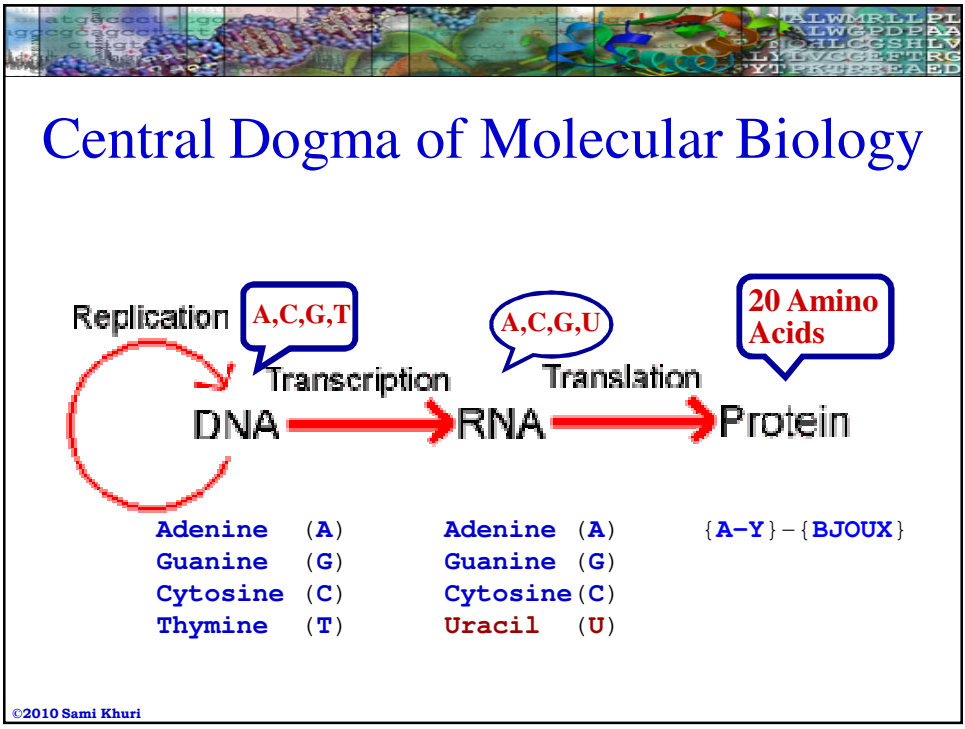
©2010 Sami Khuri

Protein Factory

Proteins: basis of how biology gets things done.

A typical **protein** is 300-500 amino acids long and folds into a 3-dimensional structure which determines its properties.

©2010 Sami Khuri





Prokaryotes and Eukaryotes

A **cell** is the fundamental working unit of every living organism.

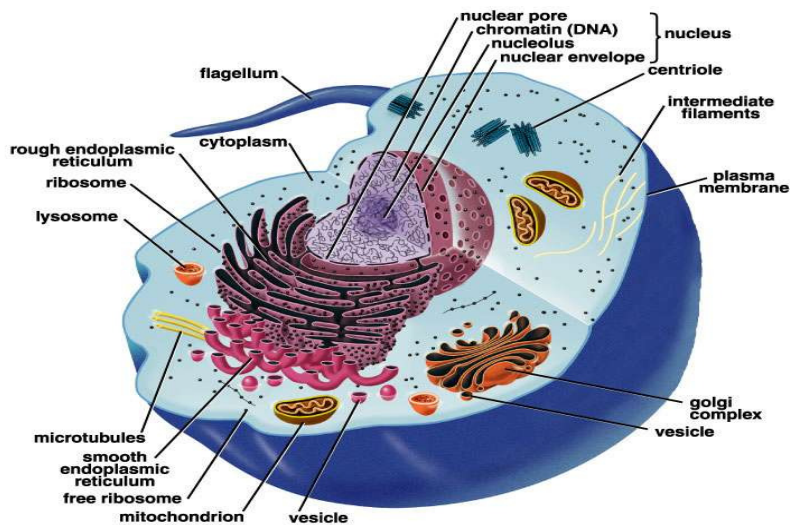
There are two kinds of cells:

- **prokaryotes**, which are mostly single-celled organisms with **no cell nucleus**: archaea and bacteria.
- **eukaryotes**, which are higher level organisms, and their cells have **nuclei**: animals and plants.

©2010 Sami Khuri



Generalized Animal Cell



©2010 Sami Khuri



Proteins and Nucleic Acids

All living organisms have a similar molecular chemistry (biochemistry). The main actors in the chemistry of life are molecules called:

- **proteins**: which are responsible for what a living being is and does in a physical sense.
“We **are** our proteins” R. Doolittle.
- **nucleic acids**: which encode the information necessary to produce proteins and are responsible for passing the “recipe” to subsequent generations.

©2010 Sami Khuri



DNA and RNA

- Living organisms contain two kinds of nucleic acids:
 - **Ribonucleic acid (RNA)**
 - **Deoxyribonucleic acid (DNA)**
- The **central dogma** states that information flows from **DNA** to **RNA** to **protein**.
- The function of a **protein** is determined by its unique three-dimensional structure.

©2010 Sami Khuri



DNA and Chromosomes

- The **human genome**: a complete set of instructions for making an organism, consists of tightly coiled threads of **DNA** and associated protein molecules, organized into structures called **chromosomes**.
- Besides the reproductive cell and red blood cell, every single **cell** in the human body contains the **human genome**.

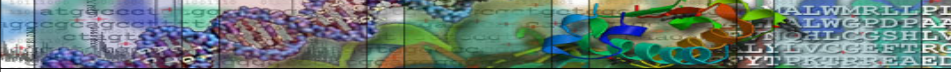
©2010 Sami Khuri



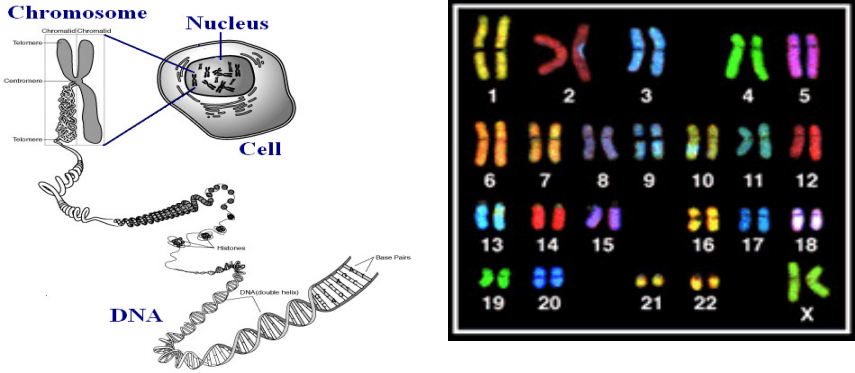
Autosomal and Sex Chromosomes

- The **human genome** is distributed along 23 pairs of chromosomes
 - 22 autosomal pairs
 - the sex chromosome pair, XX for females and XY for males.
- In each pair, one chromosome is **paternally** inherited, the other **maternally** inherited.

©2010 Sami Khuri



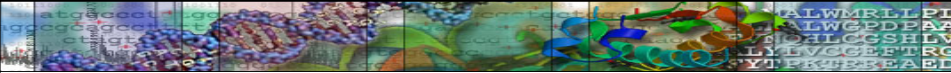
Chromosomes and Genome



The diagram illustrates the hierarchical structure of DNA. At the top, a **Chromosome** is shown as a condensed structure with labels for **Centromere**, **Chromatid**, and **Chromatid**. This is shown within a **Nucleus** and a **Cell**. Below, the DNA is shown as a **DNA (double helix)** with labels for **Histones** and **Base Pairs**. To the right, a human karyotype shows 22 numbered pairs of autosomes and the X and Y sex chromosomes.

Number of chromosomes in a genome is characteristic of a **species**.
The human **DNA** contains about three billion **base pairs** (A-T or C-G).

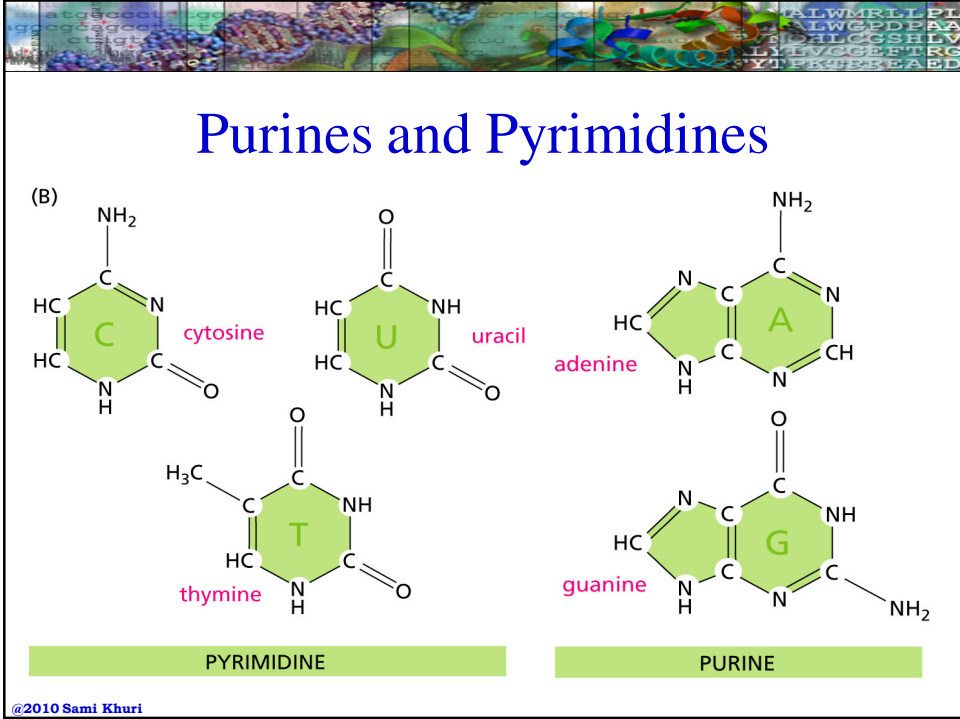
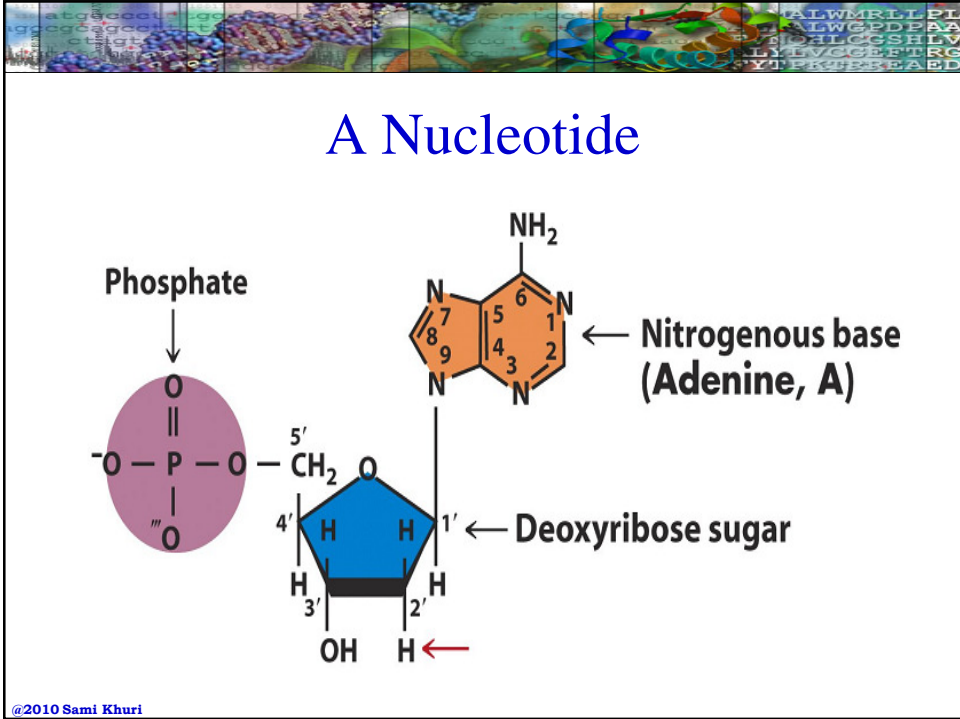
©2010 Sami Khuri



DNA Structure

- A **deoxyribonucleic acid** or **DNA** molecule is a double-stranded polymer composed of four basic molecular units called nucleotides.
- Each nucleotide comprises
 - a phosphate group
 - a deoxyribose sugar
 - one of four nitrogen bases:
 - purines: **adenine** (A) and **guanine** (G)
 - pyrimidines: **cytosine** (C) and **thymine** (T).

©2010 Sami Khuri





Double Helix

- The binding of two nucleotides forms a base pair.
- The double helix is formed by connecting complementary nucleotides A-T and C-G on two strands with hydrogen bonds.
- Knowledge of the sequence on one strand allows us to infer the sequence of the other strand.
- The bases are arranged along the sugar phosphate backbone in a particular order, known as the DNA sequence, encoding all genetic instructions for an organism.

©2010 Sami Khuri



DNA Phosphodiester Backbone

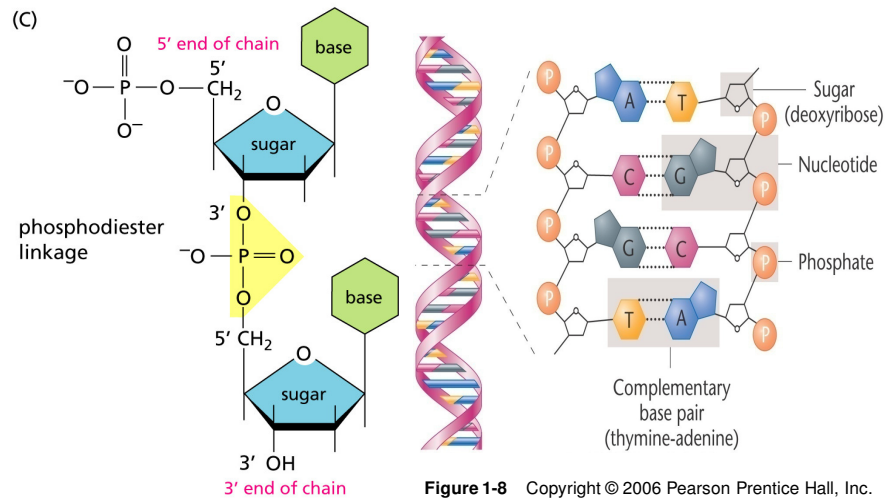
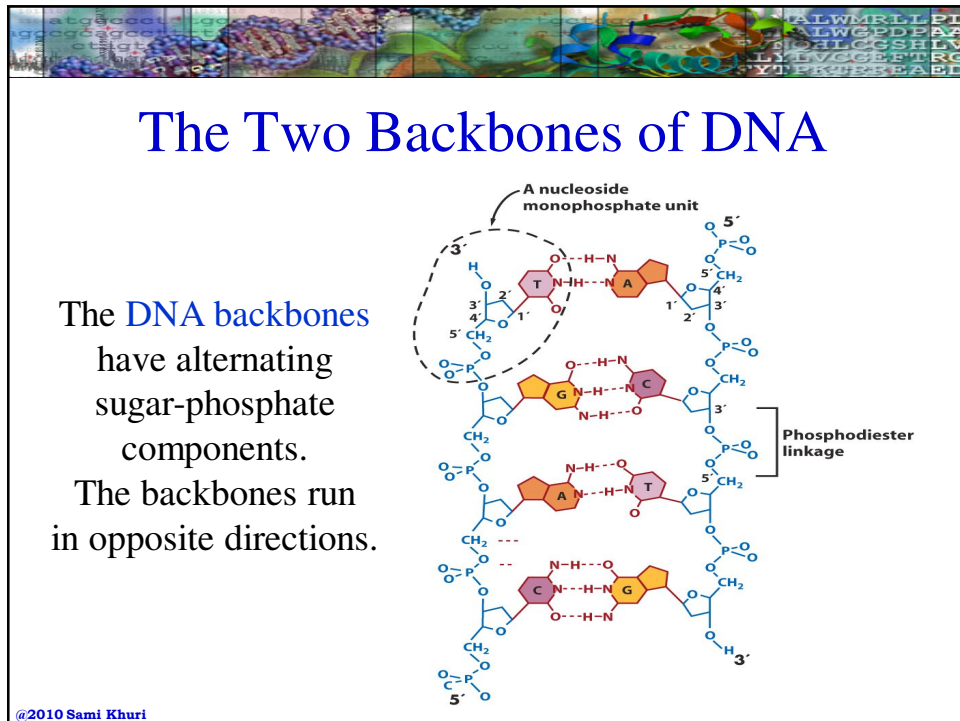
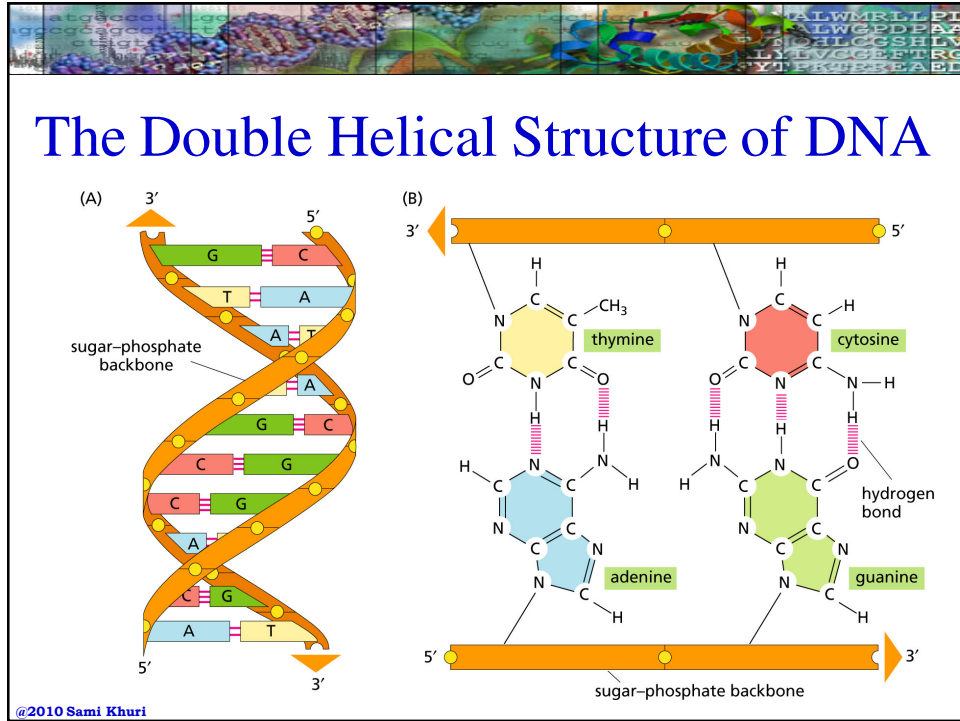


Figure 1-8 Copyright © 2006 Pearson Prentice Hall, Inc.

©2010 Sami Khuri





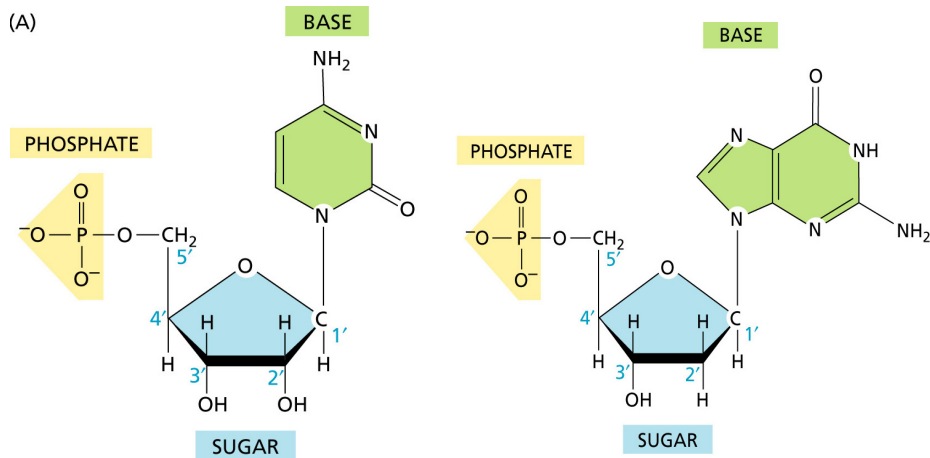
DNA and Chromosomes

- The **genome** is a complete set of instructions for making an organism, consists of tightly coiled threads of **DNA** organized into structures called **chromosomes**.
- Besides the reproductive cell and red blood cell, every single **cell** in the human body contains the **human genome**.

@2010 Sami Khuri

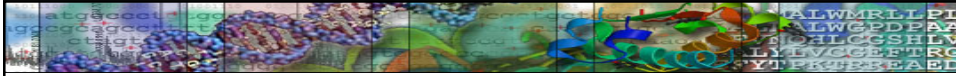


Sugars Present in Nucleic Acids



Pentose sugar present in **RNA** Pentose sugar present in **DNA**

@2010 Sami Khuri

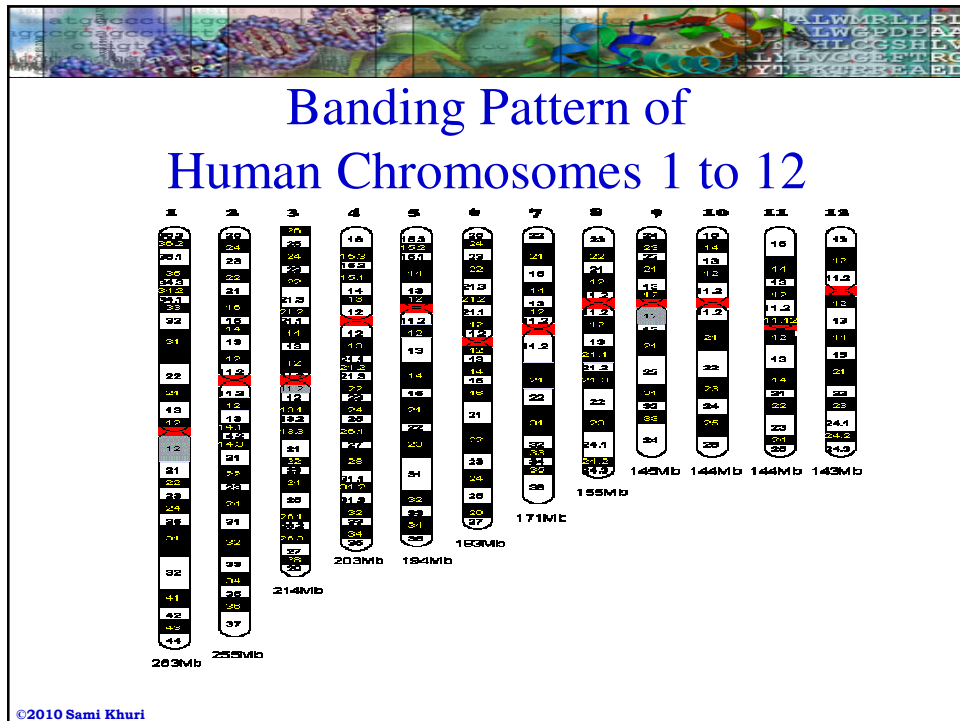


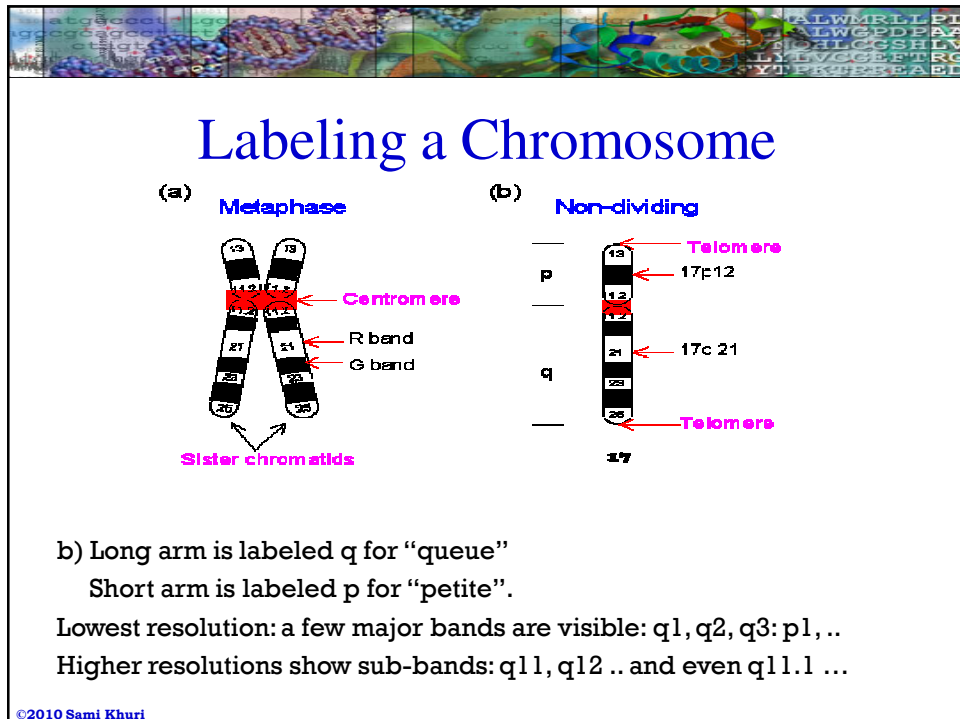
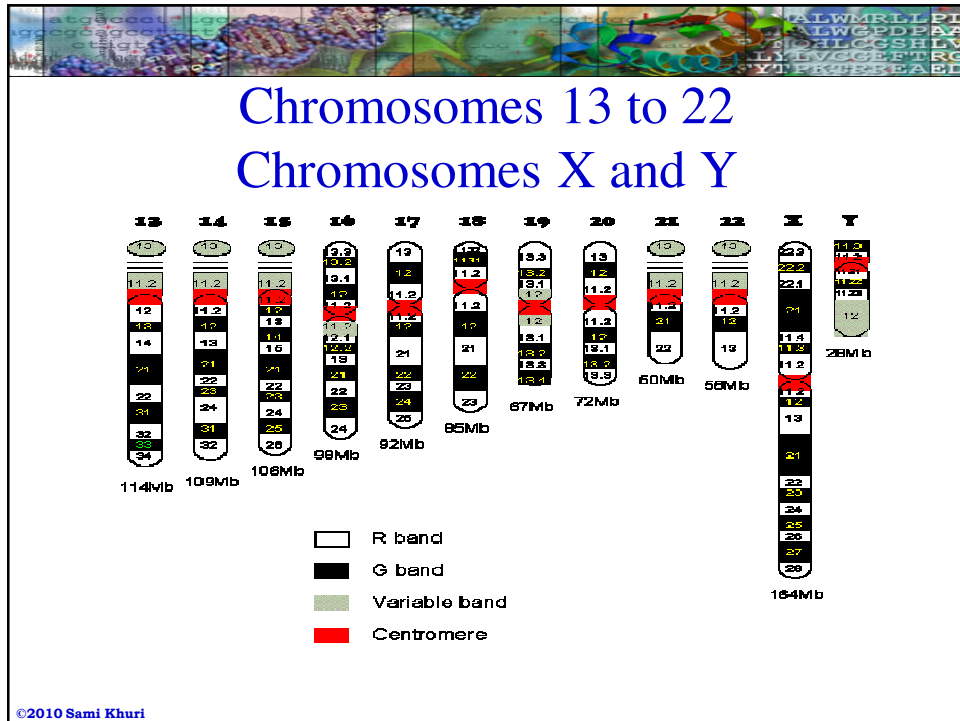
Pairs of Chromosomes in Species

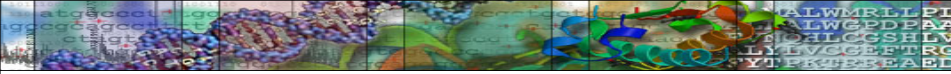
Table 3-2 Numbers of Pairs of Chromosomes in Different Species of Plants and Animals

Common name	Scientific name	Number of chromosome pairs	Common name	Scientific name	Number of chromosome pairs
Mosquito	<i>Culex pipiens</i>	3	Wheat	<i>Triticum aestivum</i>	21
Housefly	<i>Musca domestica</i>	6	Human	<i>Homo sapiens</i>	23
Garden onion	<i>Allium cepa</i>	8	Potato	<i>Solanum tuberosum</i>	24
Toad	<i>Bufo americanus</i>	11	Cattle	<i>Bos taurus</i>	30
Rice	<i>Oryza sativa</i>	12	Donkey	<i>Equus asinus</i>	31
Frog	<i>Rana pipiens</i>	13	Horse	<i>Equus caballus</i>	32
Alligator	<i>Alligator mississippiensis</i>	16	Dog	<i>Canis familiaris</i>	39
Cat	<i>Felis domesticus</i>	19	Chicken	<i>Gallus domesticus</i>	39
House mouse	<i>Mus musculus</i>	20	Carp	<i>Cyprinus carpio</i>	52
Rhesus monkey	<i>Macaca mulatta</i>	21			

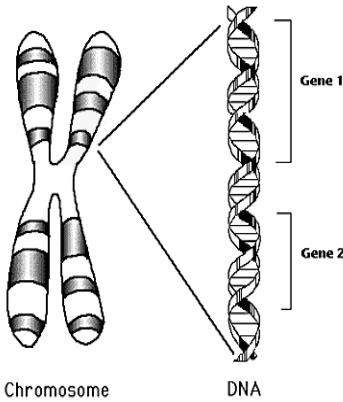
©2010 Sami Khuri








Genes

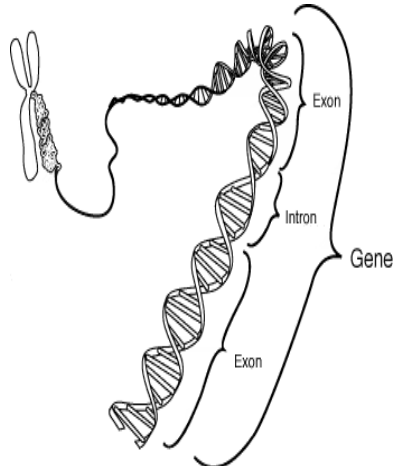


- A **gene** is a specific sequence of nucleotide bases along a chromosome carrying information for constructing a protein.
- **Genes** are part of the chromosomes.
- The distance between **genes** is often much larger than the genes themselves.

©2010 Sami Khuri



Exons and Introns



In eukaryotes, genes consist of:

- **exons**
protein-coding regions
- **introns**
noncoding regions.

Approximately 5-10% of the gene is made up of exons while the rest are introns.

www.accessexcellence.org/AB/GG/gene.html

©2010 Sami Khuri



Ribonucleic Acid - RNA

- **RNA** is found in the cell and can also carry genetic information.
- While DNA is located primarily in the nucleus, **RNA** can also be found in the **cytoplasm**: the cellular liquid outside the nucleus.
- **RNA** is built from the nucleotides **cytosine**, **guanine**, **adenine** and **uracil (U)** (instead of thymine).
- **RNA** has its sugar phosphate backbone containing **ribose**.
- **RNA** forms a **single strand**.


©2010 Sami Khuri



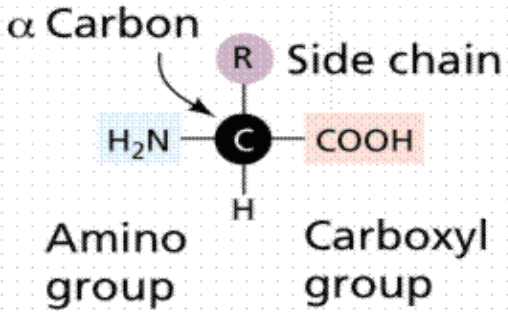
Proteins

- 20 different **amino acids** are used to synthesize **proteins**.
- The shape and other properties of each **protein** is dictated by the precise sequence of **amino acids** in it.
- The function of a **protein** is determined by its unique three-dimensional structure.

©2010 Sami Khuri




Structure of the Amino Acid



It is the structure of the R group that determines which of the 20 amino acids it is and its special properties.

©2010 Sami Khuri



The Twenty Amino Acids

$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_3 \end{matrix}$ Alanine A	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH} \\ \\ \text{H}_3\text{C}-\text{CH}_3 \end{matrix}$ Valine V	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{CH} \\ \\ \text{H}_3\text{C}-\text{CH}_3 \end{matrix}$ Leucine L	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH} \\ \\ \text{CH}_2 \\ \\ \text{CH}_3 \end{matrix}$ Isoleucine I	$\begin{matrix} \text{COO}^- \\ \\ \text{HN}-\text{C}-\text{H} \\ \\ 2\text{HC}-\text{CH}_2 \end{matrix}$ Proline P
$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{S} \\ \\ \text{CH}_3 \end{matrix}$ Methionine M	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C}_6\text{H}_5 \end{matrix}$ Phenylalanine F	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{H} \end{matrix}$ Tryptophan W	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{H} \end{matrix}$ Glycine G	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{OH} \\ \\ \text{S} \end{matrix}$ Serine S
$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{HC}-\text{OH} \\ \\ \text{CH}_3 \end{matrix}$ Threonine T	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{SH} \end{matrix}$ Cysteine C	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{NH}_2 \end{matrix}$ Asparagine N	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{NH}_2 \end{matrix}$ Glutamine Q	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{OH} \end{matrix}$ Tyrosine Y
$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{O}^- \end{matrix}$ Aspartic Acid D	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{C} \\ \\ \text{O}^- \end{matrix}$ Glutamic Acid E	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{NH}_3^+ \end{matrix}$ Lysine K	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{CH}_2 \\ \\ \text{CH}_2 \\ \\ \text{NH} \\ \\ \text{NH}_2 \end{matrix}$ Arginine R	$\begin{matrix} \text{COO}^- \\ \\ \text{H}_3\text{N}^+-\text{C}-\text{H} \\ \\ \text{HC}=\text{C} \\ \\ \text{NH} \\ \\ \text{H} \end{matrix}$ Histidine H


©2010 Sami Khuri

Orange:
nonpolar and hydrophobic.

The other amino acids are:
polar and hydrophilic - "water loving".

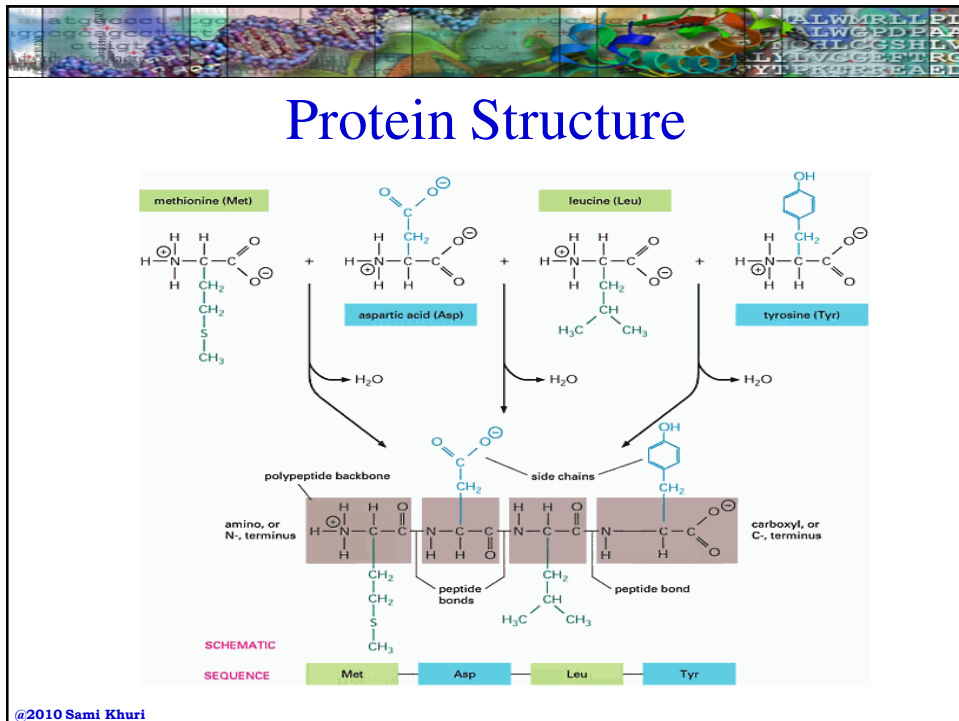
Magenta:
acidic - "carboxy" group in the side chain.

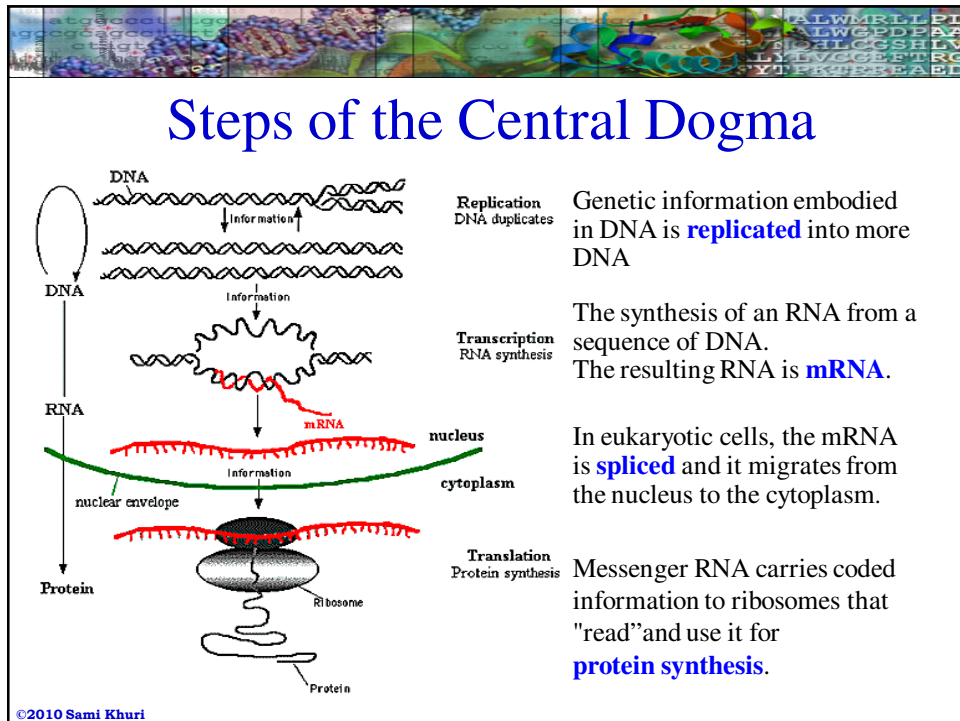
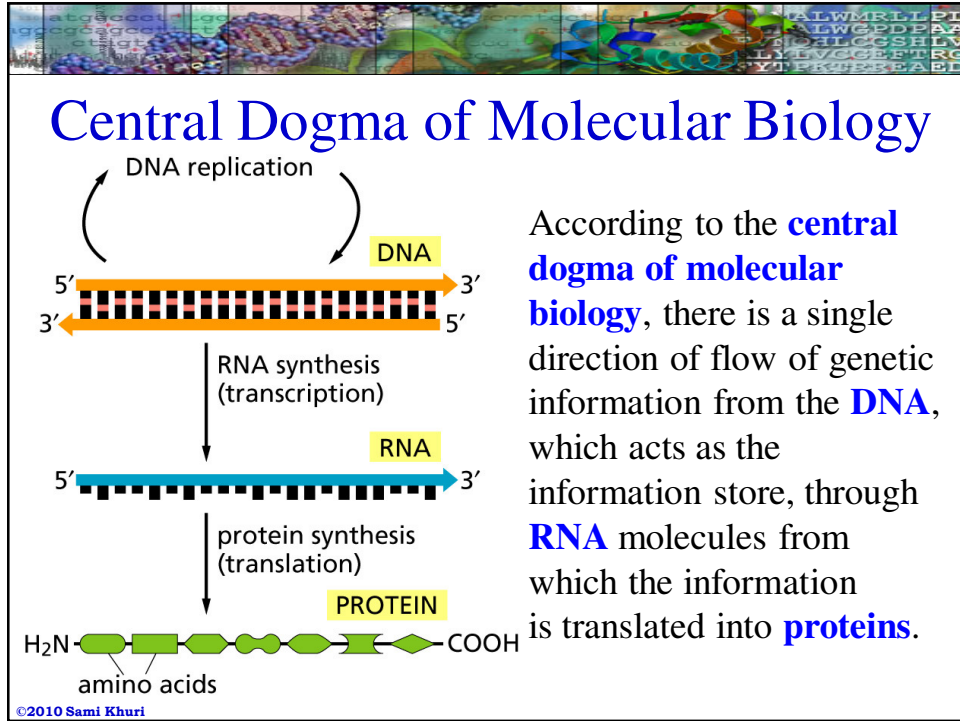
Light blue:
basic - "amine" group in the side chain.

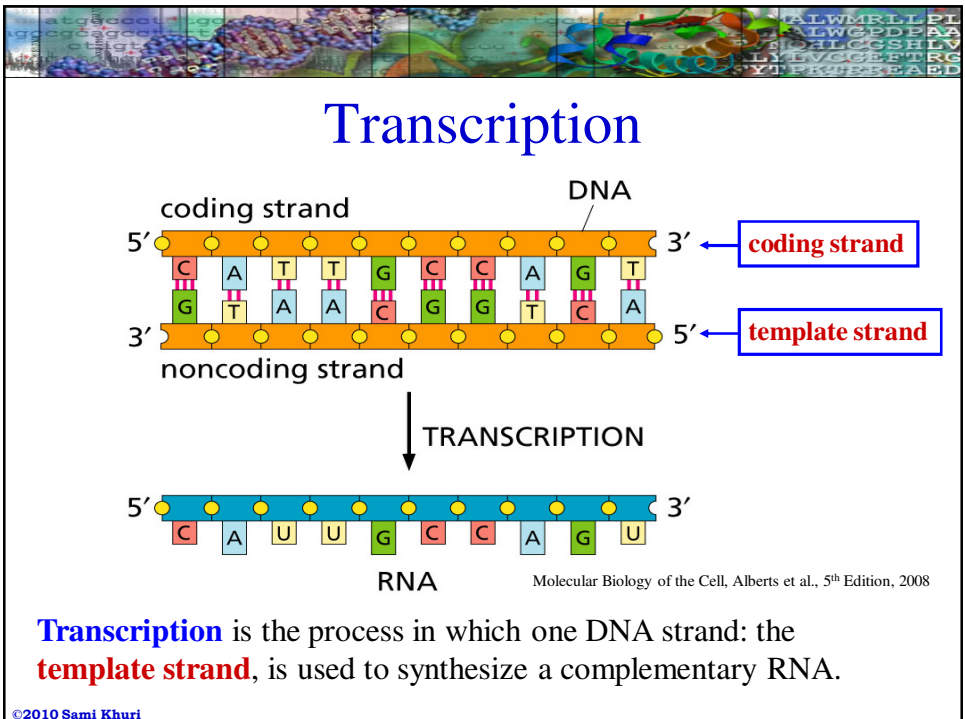
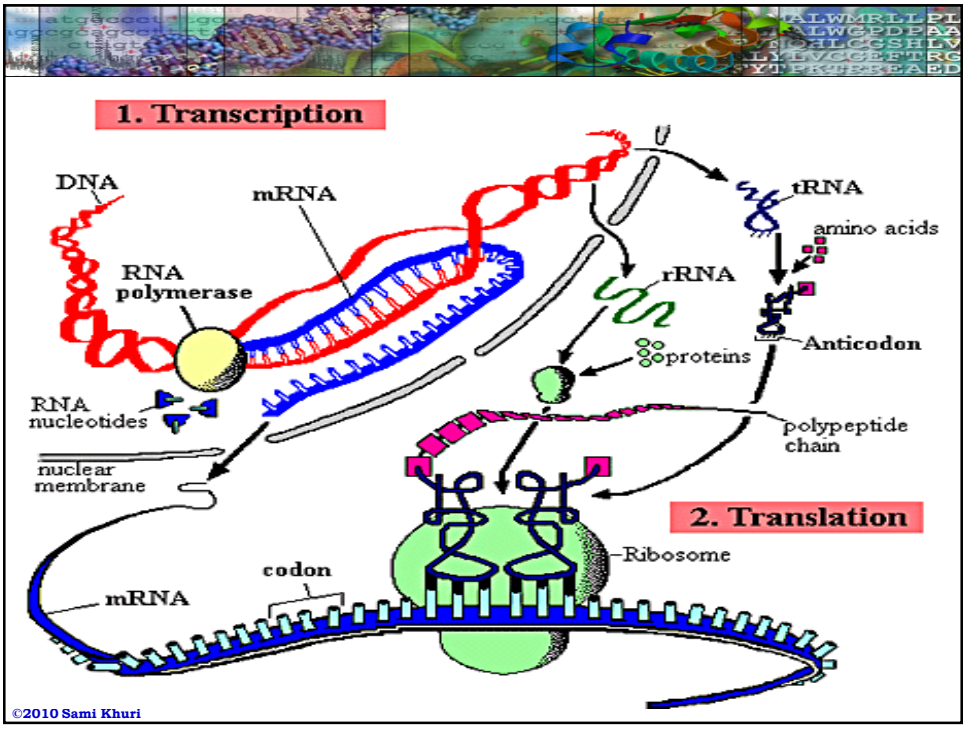


	One-letter code	Three-letter-code	Name
1	A	Ala	Alanine
2	C	Cys	Cysteine
3	D	Asp	Aspartic Acid
4	E	Glu	Glutamic Acid
5	F	Phe	Phenylalanine
6	G	Gly	Glycine
7	H	His	Histidine
8	I	Ile	Isoleucine
9	K	Lys	Lysine
10	L	Leu	Leucine
11	M	Met	Methionine
12	N	Asn	Asparagine
13	P	Pro	Proline
14	Q	Gln	Glutamine
15	R	Arg	Arginine
16	S	Ser	Serine
17	T	Thr	Threonine
18	V	Val	Valine
19	W	Trp	Tryptophan
20	Y	Tyr	Tyrosine

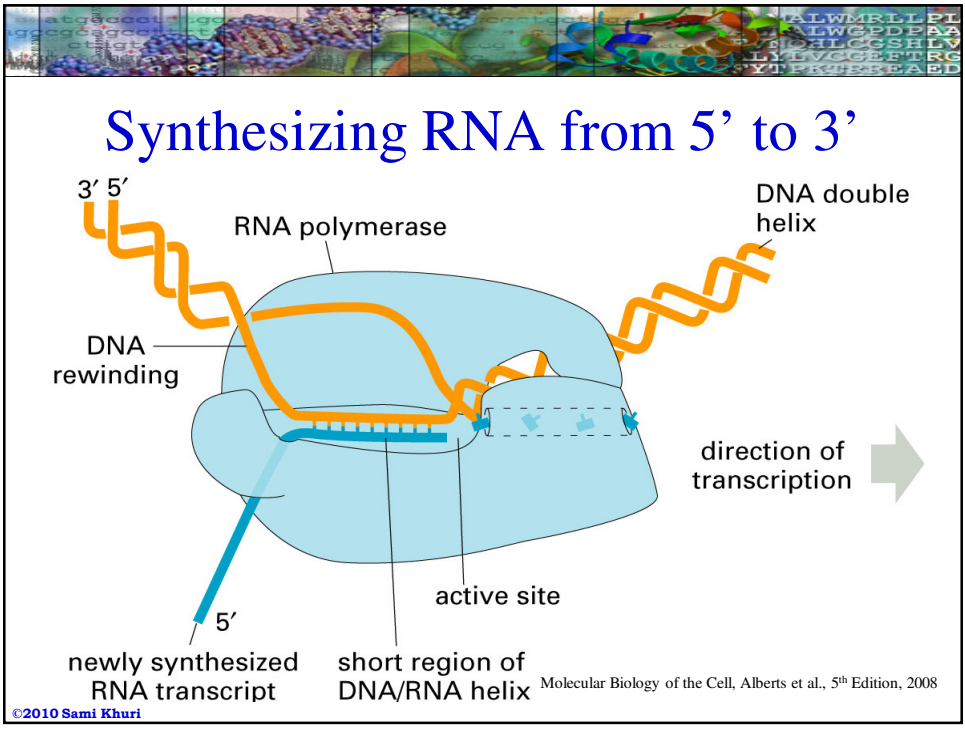
@2010 Sami Khuri David Gilbert







Transcription is the process in which one DNA strand: the **template strand**, is used to synthesize a complementary RNA.



The Genetic Code

		SECOND BASE				
		U	C	A	G	
FIRST BASE	U	UUU } Phe	UCU } Ser	UAU } Tyr	UGU } Cys	THIRD BASE
	UUC } Phe	UCC } Ser	UAC } Tyr	UGC } Cys		
	UUA } Leu	UCA } Ser	UAA } Stop	UGA } Stop		
	UUG } Leu	UCG } Ser	UAG } Stop	UGG } Trp		
C	CUU } Leu	CCU } Pro	CAU } His	CGU } Arg	U	
	CUC } Leu	CCC } Pro	CAC } His	CGC } Arg		
	CUA } Leu	CCA } Pro	CAA } Gln	CGA } Arg		
	CUG } Leu	CCG } Pro	CAG } Gln	CGG } Arg		
A	AUU } Ile	ACU } Thr	AAU } Asn	AGU } Ser	C	
	AUC } Ile	ACC } Thr	AAC } Asn	AGC } Ser		
	AUA } Met	ACA } Thr	AAA } Lys	AGA } Arg		
	AUG } Met	ACG } Thr	AAG } Lys	AGG } Arg		
G	GUU } Val	GCU } Ala	GAU } Asp	GGU } Gly	A	
	GUC } Val	GCC } Ala	GAC } Asp	GGC } Gly		
	GUA } Val	GCA } Ala	GAA } Glu	GGA } Gly		
	GUG } Val	GCG } Ala	GAG } Glu	GGG } Gly		

©2010 Sami Khuri



Transfer RNA and Translation

- The translation from nucleotides to amino acid is done by means of **transfer RNA (tRNA)** molecules, each specific for one amino acid and for a particular **triplet** of nucleotides in mRNA called a **codon**.
- The family of tRNA molecules enables the codons in a mRNA molecule to be **translated** into the sequence of amino acids in the protein.

©2010 Sami Khuri



Codons and Anticodons

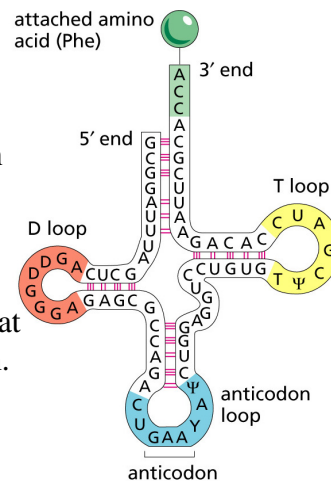
At least one kind of **tRNA** is present for each of the 20 amino acids used in protein synthesis.

Each kind of **tRNA** has a sequence of 3 unpaired nucleotides - the **anticodon** - which can bind to the complementary triplet of nucleotides - the **codon** - in an **mRNA** molecule.

The reading of codons in mRNA requires that the anticodons bind in the opposite direction.

Anticodon: 3' AAG 5'

Codon: 5' UUC 3'



©2010 Sami Khuri



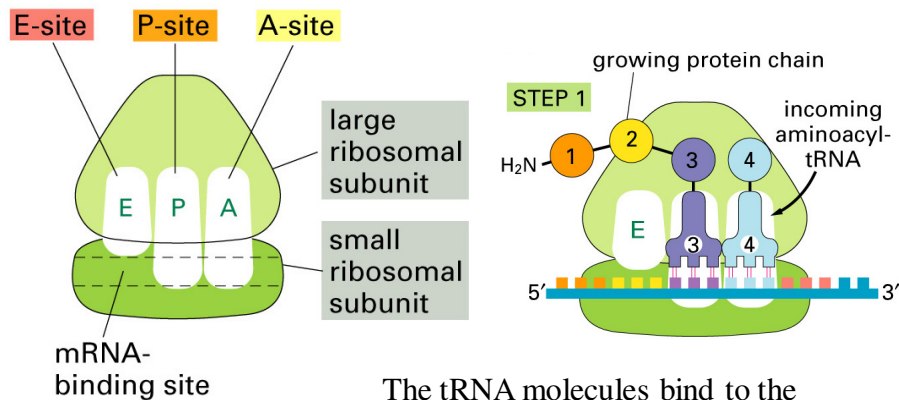
Start and Stop Codons

- The codon AUG serves two related functions
 - It begins most messages; that is, it signals the start of translation placing the amino acid methionine at the amino terminal of the polypeptide to be synthesized.
 - When it occurs within the message, it guides the incorporation of methionine.
- Three **codons**, UAA, UAG, and UGA, act as signals to terminate translation. They are called **STOP codons**.

©2010 Sami Khuri



Translation



Binding site of ribosome for the mRNA and the three tRNA binding sites.

The tRNA molecules bind to the ribosome and are the physical link between the mRNA and the growing protein chain.

©2010 Sami Khuri



Steps of Translation: Initiation

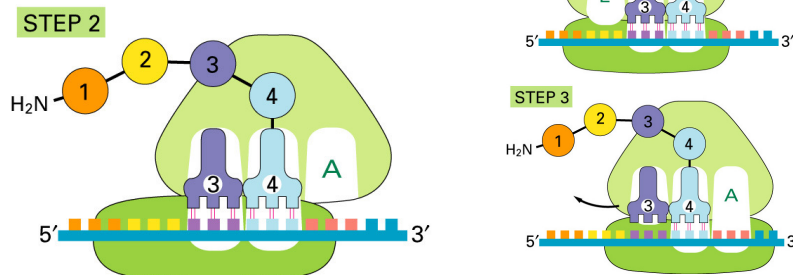
- The small subunit of the ribosome binds to a site “upstream” of the start of the message.
- It proceeds downstream until it encounters the **start codon** AUG.
- It is then joined by the large subunit and a special **initiator tRNA**. The initiator tRNA binds to the **P site** on the ribosome.
- In eukaryotes, **initiator tRNA** generally carries methionine (Met).

©2010 Sami Khuri



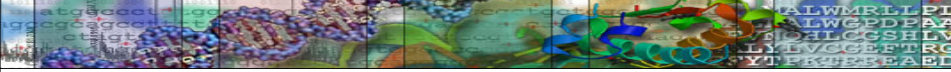
Steps of Translation: Elongation

An **aminoacyl-tRNA** able to base pair with the next codon on the mRNA arrives at the **A site**.



The preceding amino acid is linked to the incoming amino acid with a **peptide bond**.

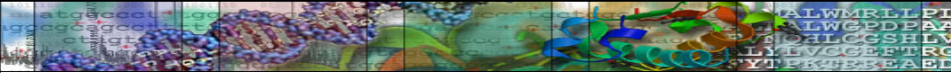
©2010 Sami Khuri



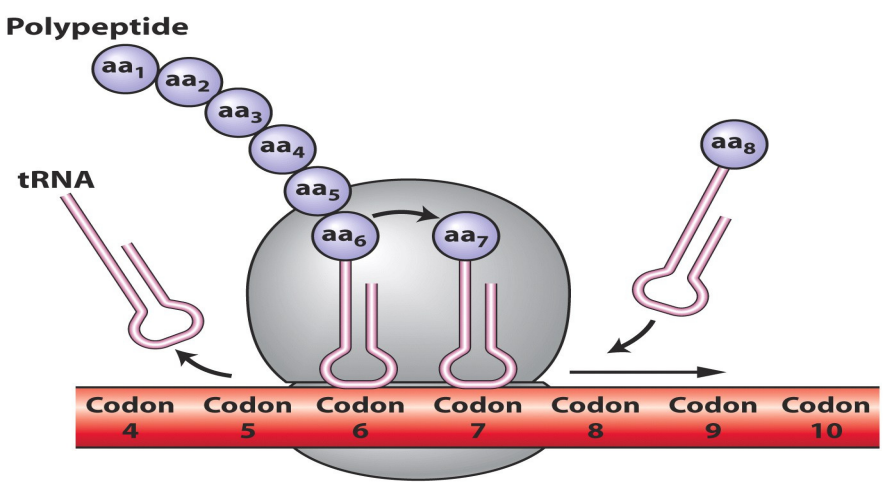
Steps of Translation: Termination

- The end of the message is marked by a **STOP codon**: **UAA**, **UAG**, **UGG**.
- No **tRNA** molecules have anticodons for **STOP codons**. A protein release factor recognizes these codons when they arrive at the **A site**.
- Binding of this protein releases the **polypeptide** from the ribosome.
- The **ribosome** splits into its subunits, which can later be reassembled for another round of **protein synthesis**.

©2010 Sami Khuri



Chain of Amino Acids



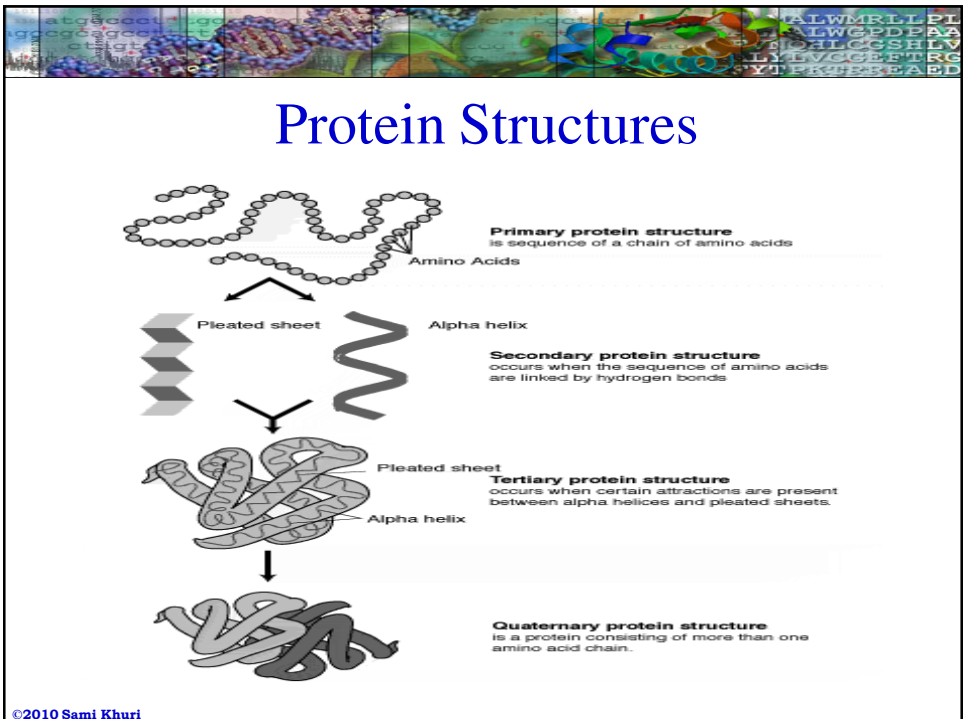
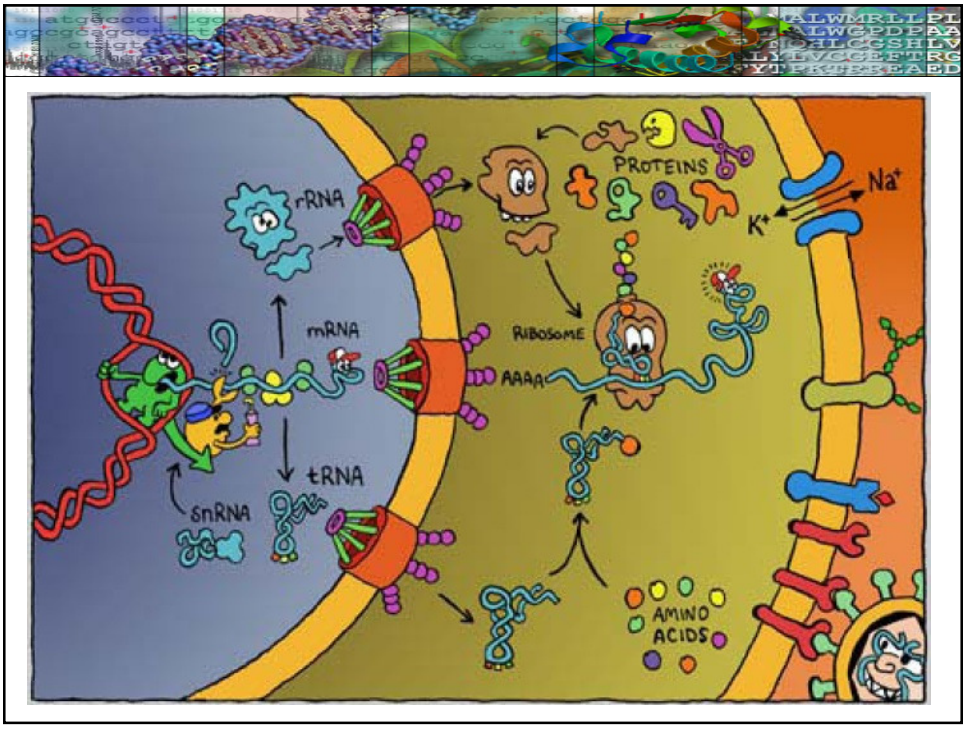
Polypeptide

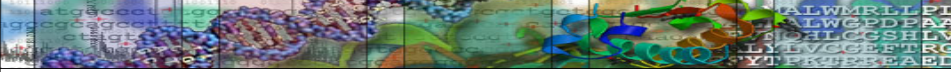
aa₁ aa₂ aa₃ aa₄ aa₅ aa₆ aa₇ aa₈

tRNA

Codon 4 Codon 5 Codon 6 Codon 7 Codon 8 Codon 9 Codon 10

©2010 Sami Khuri





3 Reading Frames of mRNA

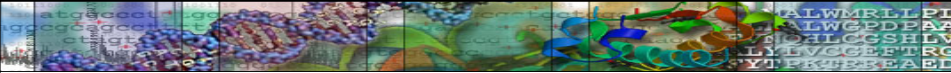
5' AGC GUU ACC AU 3'

1 CUC AGC GUU ACC AU
— Leu — Ser — Val — Thr —

2 C UCA GCG UUA CCA U
— Ser — Ala — Leu — Pro —

3 CU CAG CGU UAC CAU
— Gln — Arg — Tyr — His —

©2010 Sami Khuri



Six Reading Frames

GCT ACG CTT CGG AGC
CGA TGC CTC GAA GCT CG

©2010 Sami Khuri



Sequencing SARS



in vivo → **in vitro** → **in silico** <http://www.bcgsc.ca/bioinfo/SARS>

©2010 Sami Khuri