

# Molecular biology cheat sheet

Nucleotides A, C, G, T

gene

DNA

5' ...TACCTACATCCACTCATC...AGCTACGTTCCCCGACTACGACATGGTGATT  
 ...ATGGATGTAGGTGAGTAG...TCGATGCAAGGGGCTGATGCTGTACCACTAA... 3'

exon                      intron                      exon

RNA

transcription

...AUGGAUGUAGAUGGGGCUGAUGCUGUACCACUAA

codon

translation

Protein

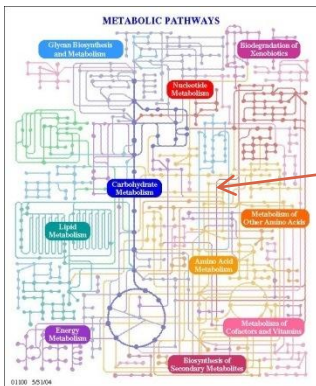
MDVDGLMLYH

Gene regulation

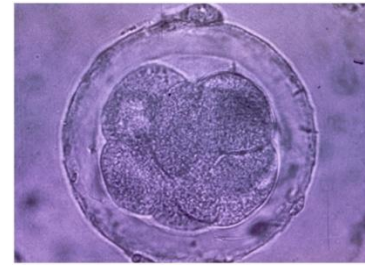
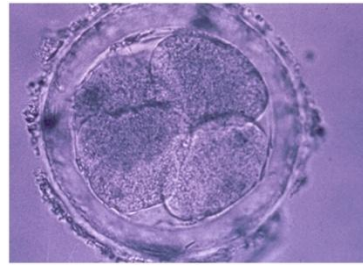
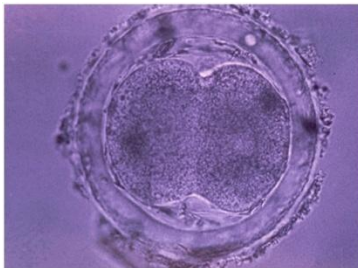
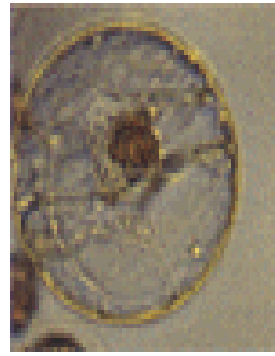
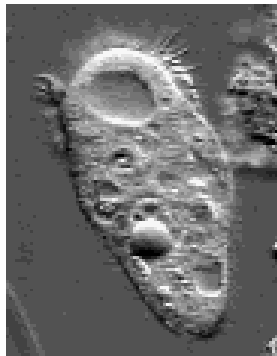
enzyme

recombination

Mother DNA — C G A A —  
 Father DNA — T C C T —  
 Daughter DNA — C C C A —

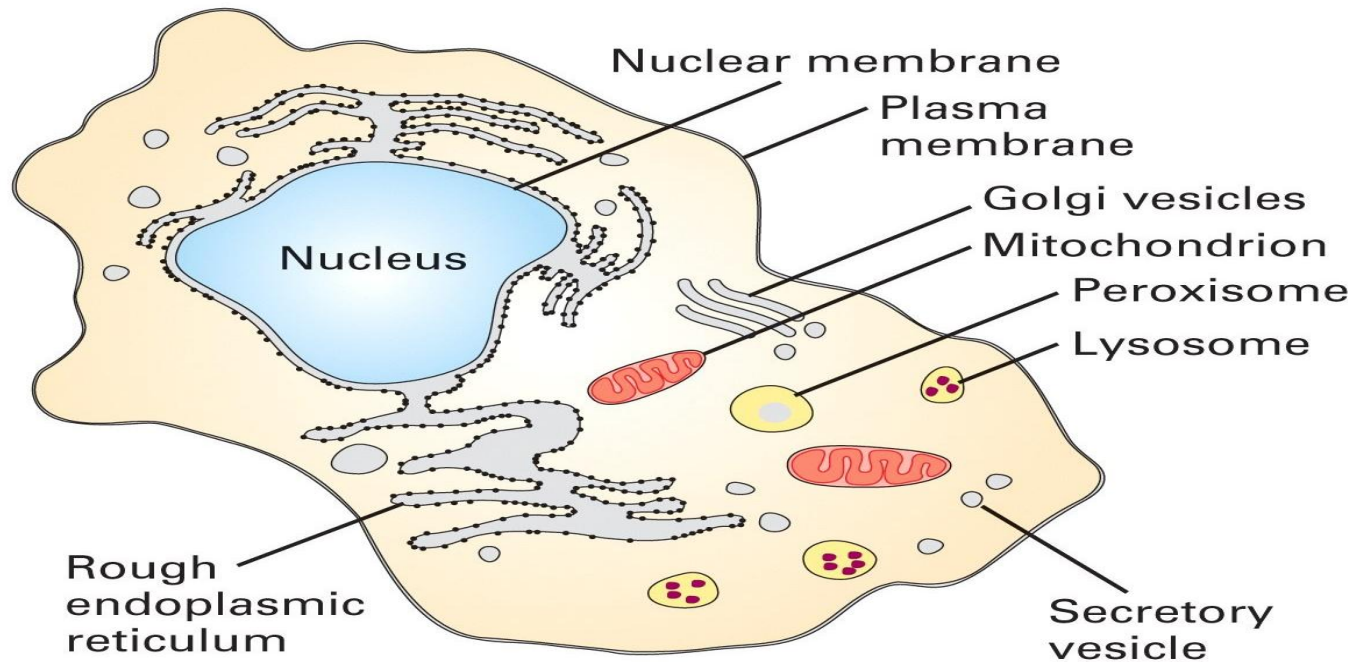


# Molecular biology primer



*Molecular Biology Primer by Angela Brooks, Raymond Brown, Calvin Chen, Mike Daly, Hoa Dinh, Erinn Hama, Robert Hinman, Julio Ng, Michael Sneddon, Hoa Troung, Jerry Wang, Che Fung Yung Edited for Introduction to Bioinformatics (Autumn 2007, Summer 2008, Autumn 2008) by Esa Pitkänen*

# Life begins with Cell



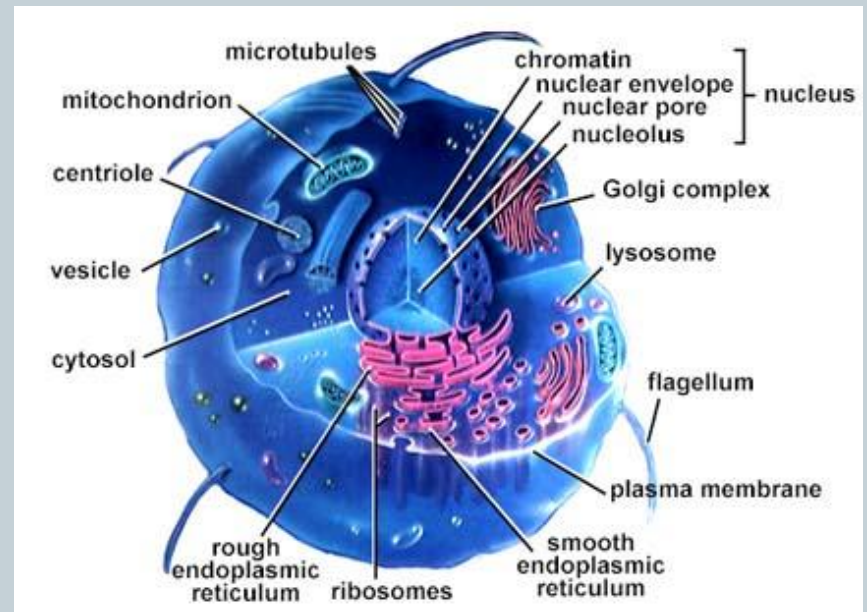
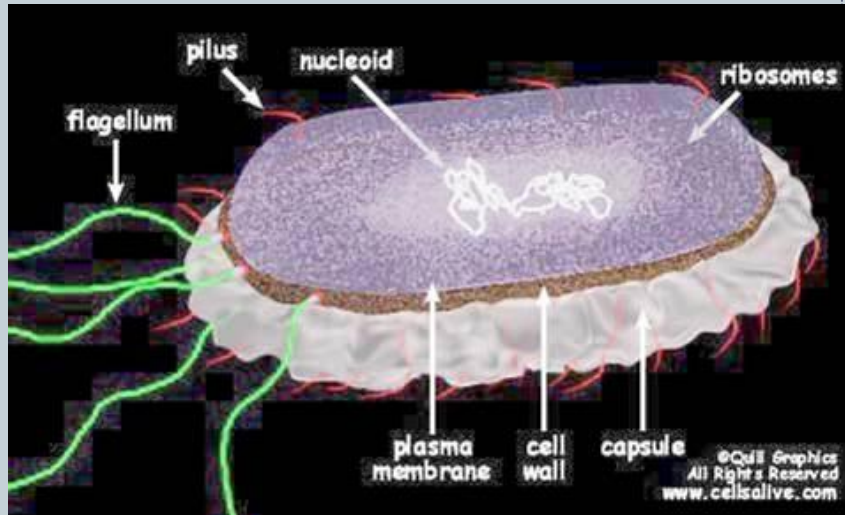
- A cell is a smallest structural unit of an organism that is capable of independent functioning
- All cells have some common features

# Cells



- **Fundamental working units** of every living system.
- Every organism is composed of one of two radically different types of cells:
  - **prokaryotic** cells or
  - **eukaryotic** cells.
- **Prokaryotes** and **Eukaryotes** are descended from the same primitive cell.
  - All prokaryotic and eukaryotic cells are the result of a total of 3.5 billion years of evolution.

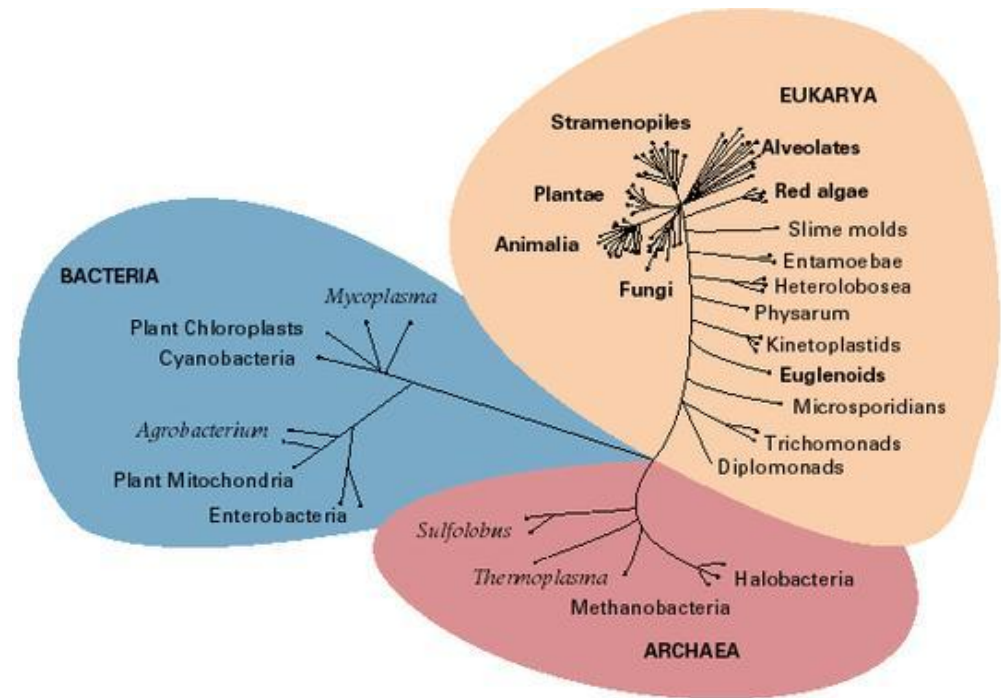
# Two types of cells: Prokaryotes and Eukaryotes



# Prokaryotes and Eukaryotes

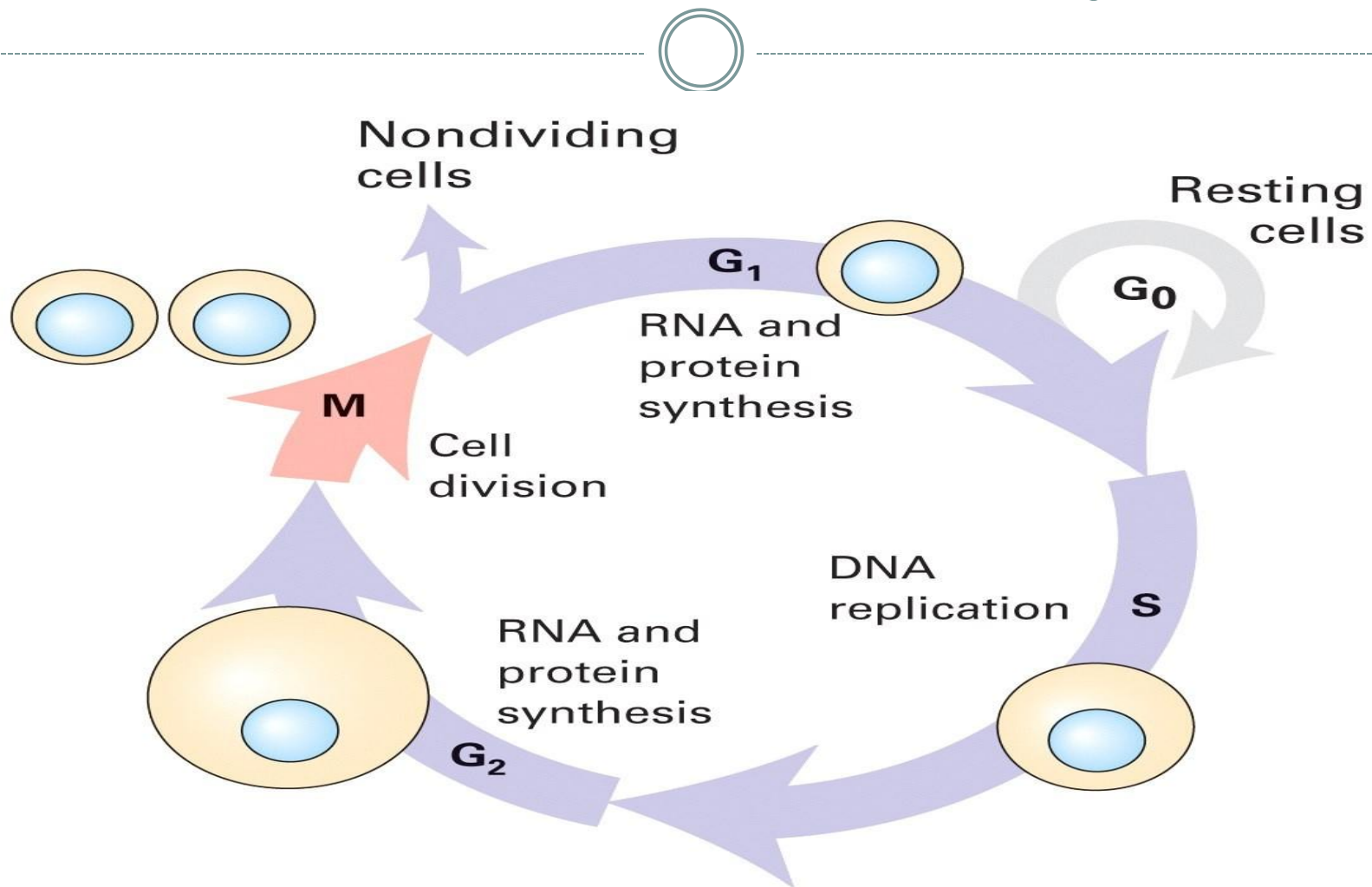


- According to the most recent evidence, there are three main branches to the tree of life
- Prokaryotes include Archaea (“ancient ones”) and bacteria
- Eukaryotes are kingdom Eukarya and includes plants, animals, fungi and certain algae





# All Cells have common Cycles



- Born, eat, replicate, and die

# Common features of organisms



- Chemical energy is stored in ATP
- Genetic information is encoded by DNA
- Information is transcribed into RNA
- There is a **common triplet genetic code**
- Translation into proteins involves ribosomes
- Shared metabolic pathways
- Similar proteins among diverse groups of organisms

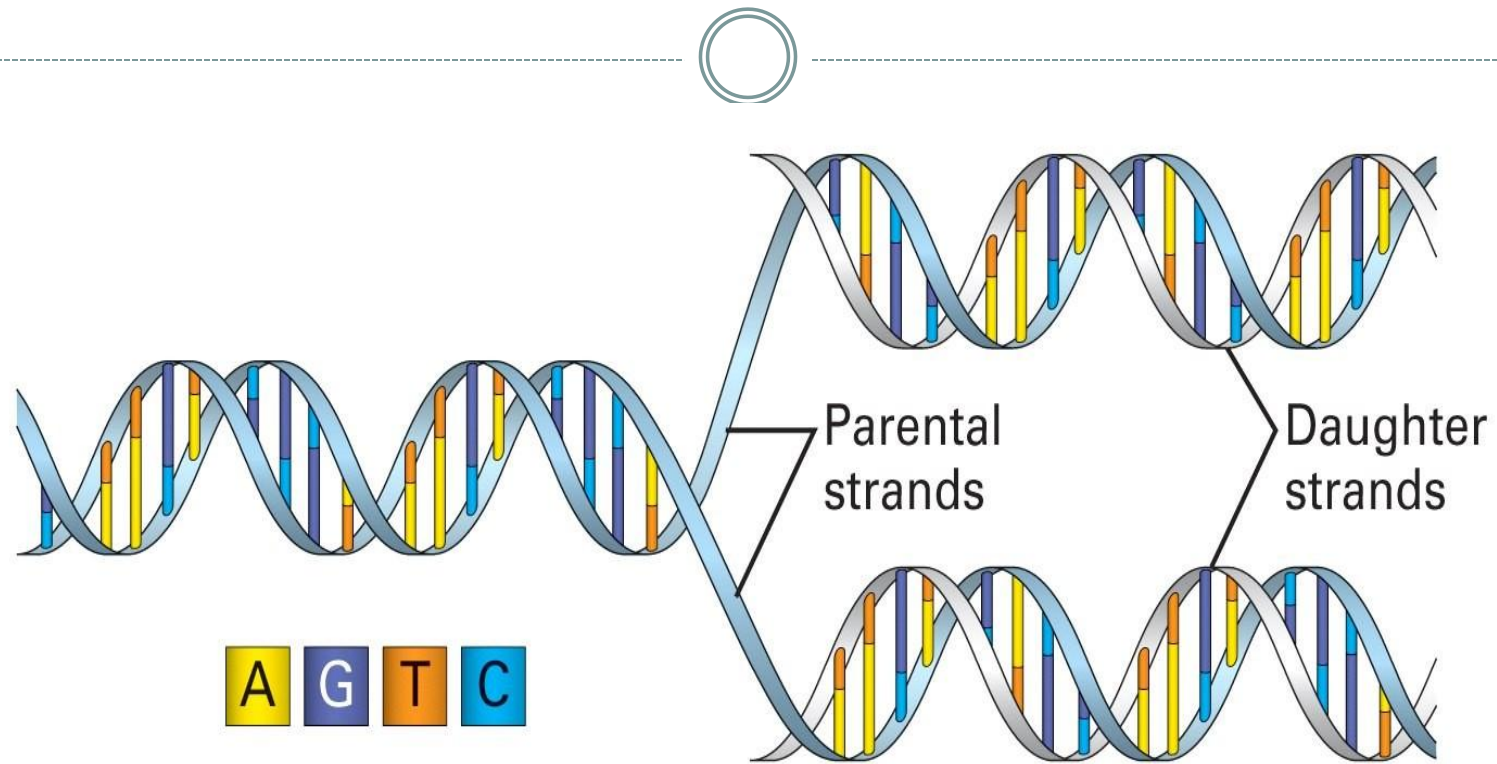


# All Life depends on 3 critical molecules



- **DNAs (Deoxyribonucleic acid)**
  - Hold information on how cell works
- **RNAs (Ribonucleic acid)**
  - Act to transfer short pieces of information to different parts of cell
  - Provide templates to synthesize into protein
- **Proteins**
  - Form enzymes that send signals to other cells and regulate gene activity
  - Form body's major components (e.g. hair, skin, etc.)
  - “Workhorses” of the cell

# DNA: The Code of Life

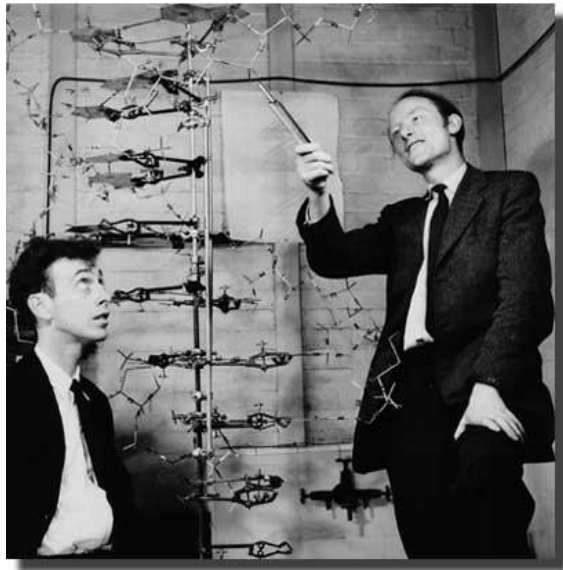


- The structure and the four genomic letters code for all living organisms
- Adenine, Guanine, Thymine, and Cytosine which pair A-T and C-G on complimentary strands.

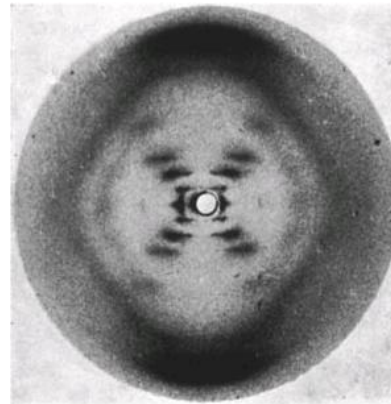
# Discovery of the structure of DNA



- **1952-1953** James D. Watson and Francis H. C. Crick deduced the double helical structure of DNA from X-ray diffraction images by Rosalind Franklin and data on amounts of nucleotides in DNA



James Watson and  
Francis Crick

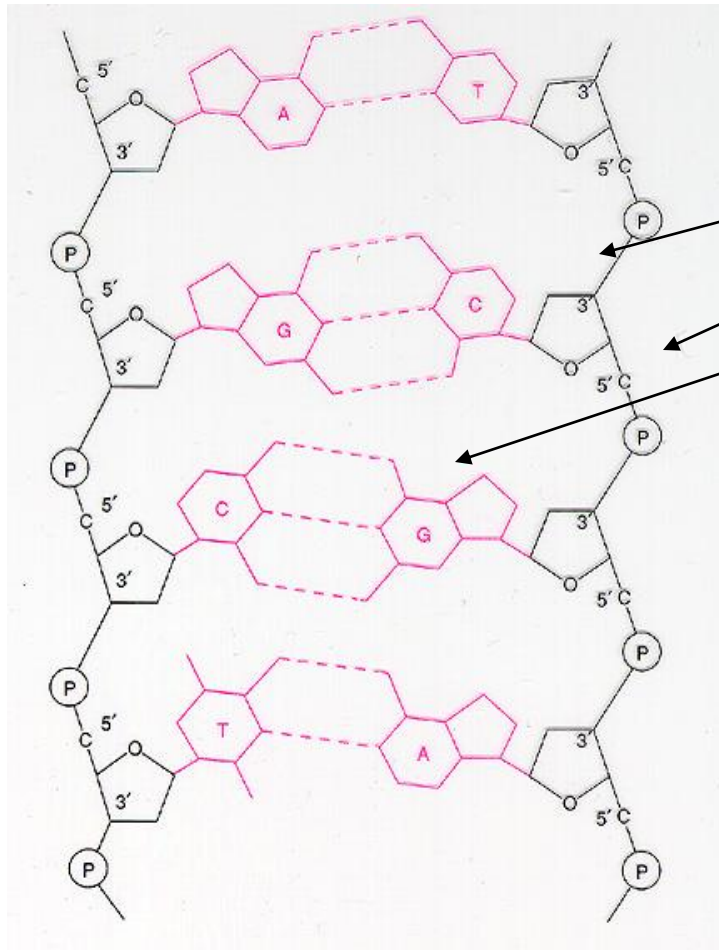


"Photo 51"



Rosalind  
Franklin

# DNA, continued



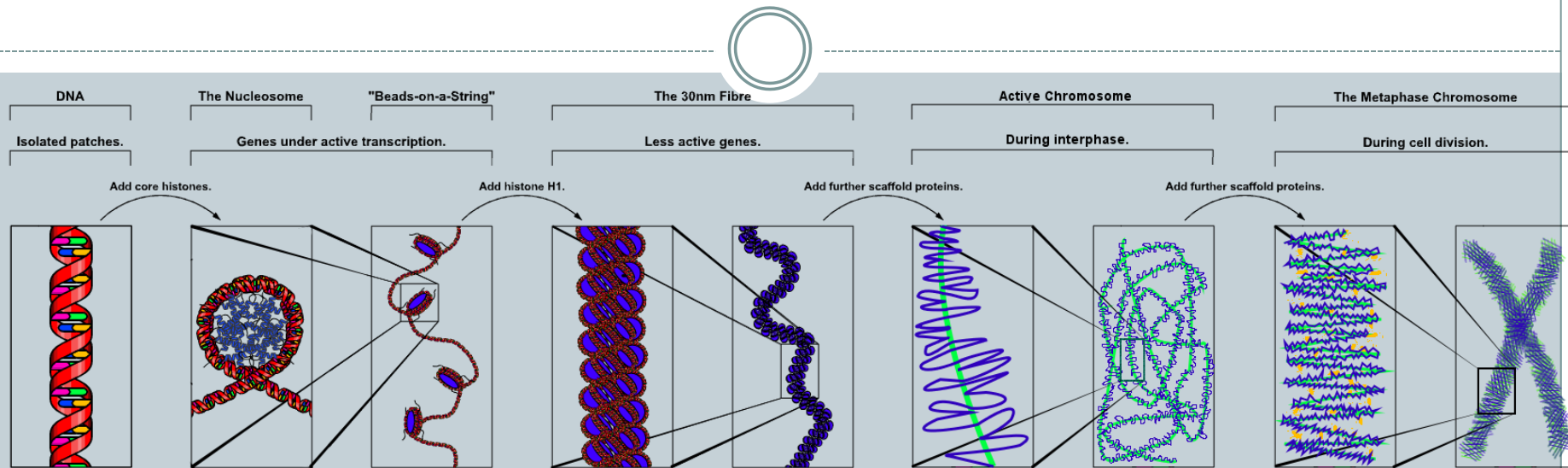
- DNA has a double helix structure which is composed of
  - sugar molecule
  - phosphate group
  - and a base (A,C,G,T)

- By convention, we read DNA strings in direction of transcription: from 5' end to 3' end

5' ATTAGGCC 3'

3' TAAATCCGG 5'

# DNA is contained in chromosomes



In eukaryotes, DNA is packed into *chromatids*

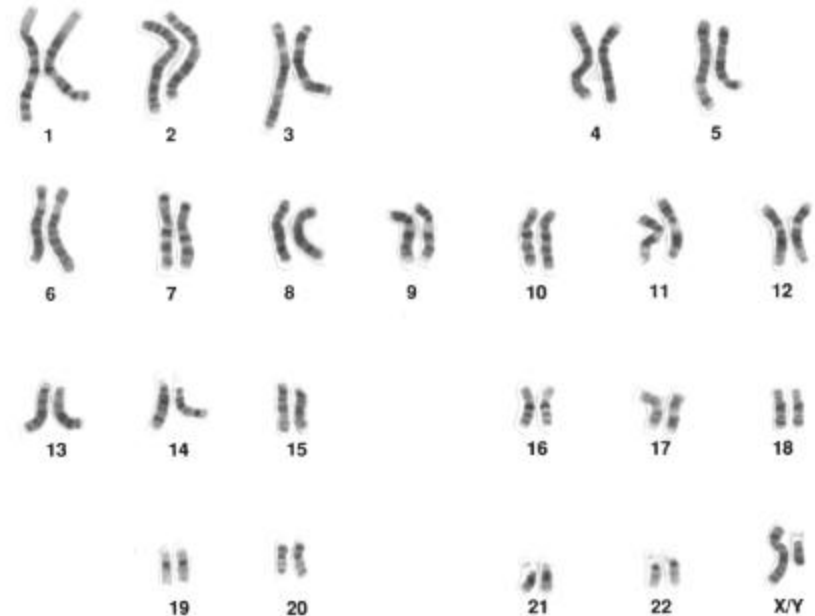
- In metaphase, the “X” structure consists of two identical chromatids

In prokaryotes, DNA is usually contained in a single, circular chromosome

# Human chromosomes



- Somatic cells in humans have 2 pairs of 22 chromosomes + XX (female) or XY (male) = total of 46 chromosomes
- Germline cells have 22 chromosomes + either X or Y = total of 23 chromosomes



Karyogram of human male using Giemsa staining (<http://en.wikipedia.org/wiki/Karyotype>)

# Length of DNA and number of chromosomes



Organism	#base pairs	#chromosomes (germline)
Prokaryotic		
Escherichia coli (bacterium)	$4 \times 10^6$	1
Eukaryotic		
Saccharomyces cerevisia (yeast)	$1.35 \times 10^7$	17
Drosophila melanogaster (insect)	$1.65 \times 10^8$	4
Homo sapiens (human)	$2.9 \times 10^9$	23
Zea mays (corn / maize)	$5.0 \times 10^9$	10

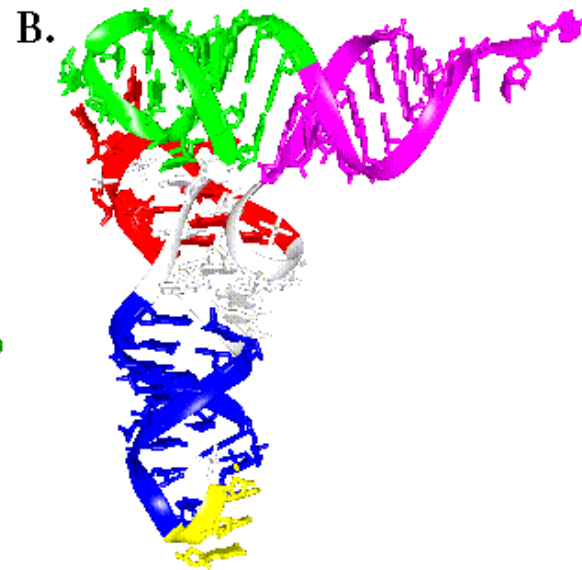
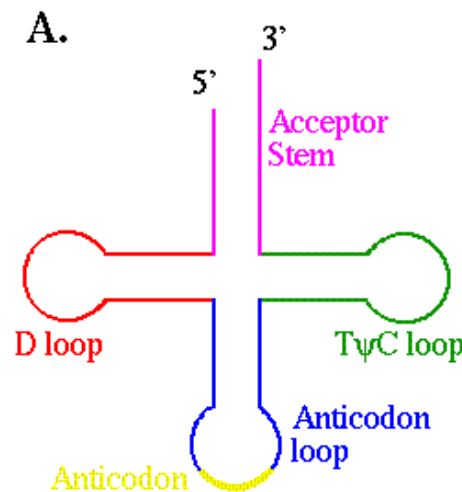


# Hepatitis delta virus, complete genome

1 atgagccaag ttccgaacaa ggattcgcgg ggaggataga tcagcgcccg agaggggtga  
61 gtcggtaaag agcattggaa cgtcggagat acaactccca agaaggaaaa aagagaaagc  
121 aagaagcggg tgaatttccc cataacgcca gtgaaactct aggaagggga aagaggggag  
181 gtggaagaga aggaggcggg cctcccgatc cgaggggccc ggcggccaag tttggaggac  
241 actccggccc gaagggttga gagtaccca gagggaggaa gccacacgga gtagaacaga  
301 gaaatcacct ccagaggacc ctttcagcga acagagagcg catcgcgaga gggagtagac  
361 catagcgata ggaggggatg ctaggagttg ggggagaccg aagcgaggag gaaagcaaag  
421 agagcagcgg ggctagcagg tgggtgttcc gccccccgag aggggacgag tgaggcttat  
481 cccggggaac tcgacttatc gtccccacat agcagactcc cggaccccct ttcaaagtga  
541 ccgagggggg tgactttgaa cattggggac cagtggagcc atgggatgct cctcccgatt  
601 ccgccaagc tccttcccc caagggtcgc ccaggaatgg cgggaccca ctctgcaggg  
661 tccgcgttcc atcctttctt acctgatggc cggcatggtc ccagcctcct cgctggcgcc  
721 ggctgggcaa cattccgagg ggaccgtccc ctcggtaatg gcgaatggga cccacaaatc  
781 tctctagctt cccagagaga agcgagagaa aagtggctct cccttagcca tccgagtgga  
841 cgtgcgtcct ccttcggatg cccaggtcgg accgcgagga ggtggagatg ccatgccgac  
901 ccgaagagga aagaaggacg cgagacgcaa acctgcgagt ggaaaccgc tttattcact  
961 ggggtcgaca actctgggga gaggaggag ggtcggctgg gaagagtata tcctatggga  
1021 atccctggct tccccttatg tccagtccct ccccggtccg agtaaagggg gactccggga  
1081 ctcttgcatt gctggggacg aagccgcccc cgggcgctcc cctcgttcca ctttcgaggg  
1141 ggttcacacc cccaacctgc gggccggeta ttcttcttcc ccttctctcg tcttctcgg  
1201 tcaacctcct aagttcctct tcctcctcct tgctgaggtt ctttcccccc gccgatagct  
1261 gctttctctt gttctcgagg gccttccttc gtcggtgatc ctgcctctcc ttgtcgggtga  
1321 atcctcccct ggaaggcctc ttectaggtc cggagtctac ttccatctgg tccgttcggg  
1381 ccctcttcgc cgggggagcc ccctctccat ccttatcttt ctttccgaga attcctttga  
1441 tgtttcccag ccagggatgt tcatectcaa gtttcttgat tttcttctta accttccgga  
1501 ggtctctctc gagttcctct aacttcttcc ttccgctcac cactgctcg agaacctctt  
1561 ctctcccccc gcggtttttc cttccttcgg gccggctcat cttcgactag aggcgacggt  
1621 cctcagtact cttactcttt tctgtaaaga ggagactgct ggccctgctg cccaagtctg  
1681 ag

# RNA

- RNA is similar to DNA chemically. It is usually only a single strand. T(hyamine) is replaced by U(racil)
- Several types of RNA exist for different functions in the cell.

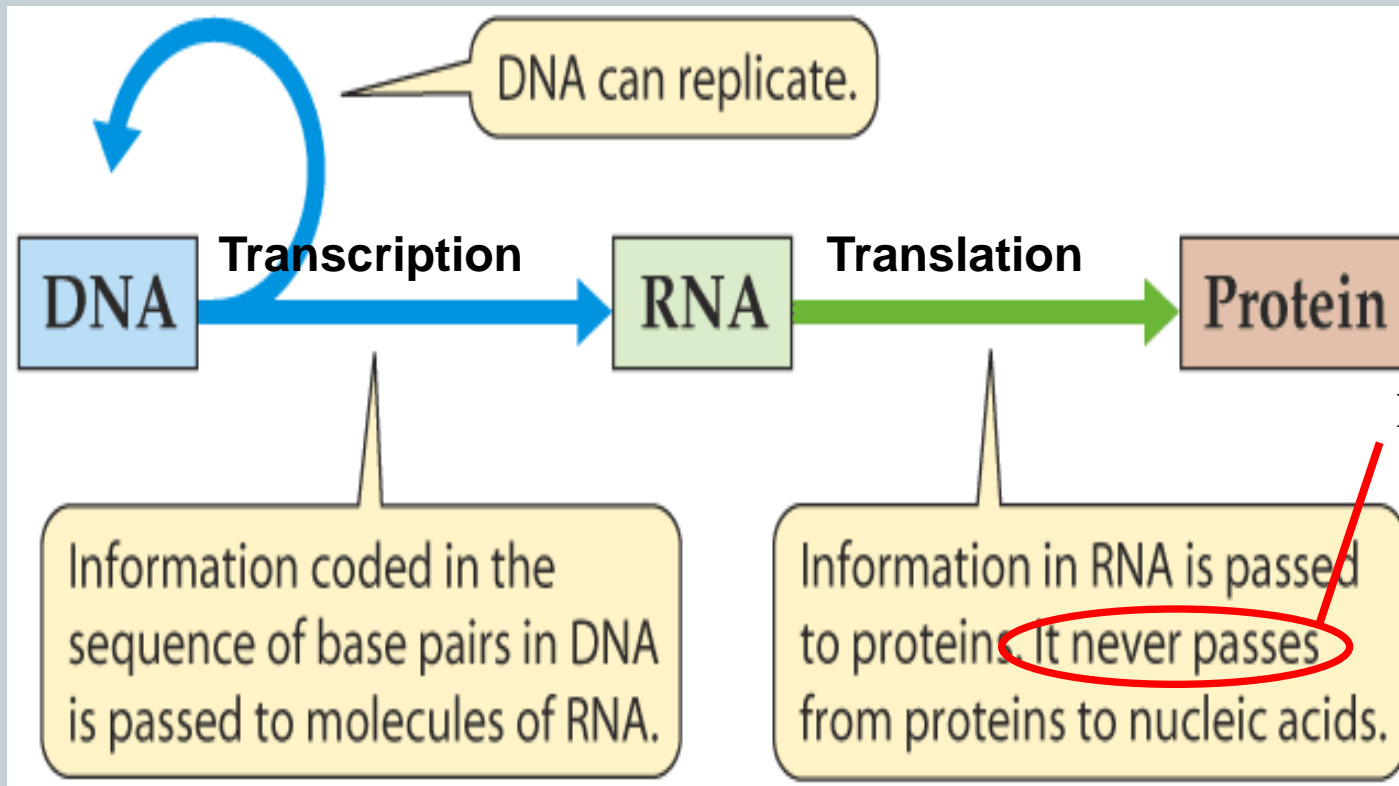


# DNA, RNA, and the Flow of Information



"The central dogma"

Replication



Is this true?

# Proteins



- Proteins are polypeptides (strings of amino acid residues)
- Represented using strings of letters from an alphabet of 20: AEGLV...WKKLAG
- Typical length 50...1000 residues



*Urease enzyme from Helicobacter pylori*

# Amino acids



$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ (\text{CH}_2)_3 \\   \\ \text{NH} \\   \\ \text{C}=\text{NH}_2 \\   \\ \text{NH}_2 \end{array}$ <p>Arginine (Arg / R)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{CH}_2 \\   \\ \text{C}=\text{O} \\   \\ \text{NH}_2 \end{array}$ <p>Glutamine (Gln / Q)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{C}_6\text{H}_5 \end{array}$ <p>Phenylalanine (Phe / F)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{C}_6\text{H}_4 \\   \\ \text{OH} \end{array}$ <p>Tyrosine (Tyr / Y)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{C}_8\text{H}_6\text{N}_2 \end{array}$ <p>Tryptophan (Trp, W)</p>
$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ (\text{CH}_2)_4 \\   \\ \text{NH}_2 \end{array}$ <p>Lysine (Lys / K)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{H} \end{array}$ <p>Glycine (Gly / G)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_3 \end{array}$ <p>Alanine (Ala / A)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{C}_3\text{H}_3\text{N}_2 \end{array}$ <p>Histidine (His / H)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{OH} \end{array}$ <p>Serine (Ser / S)</p>
$\begin{array}{c} \text{H}_2 \\   \\ \text{C} \\ / \quad \backslash \\ \text{H}_2\text{C} \quad \text{CH}_2 \\   \quad \quad   \\ \text{H}_2\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \end{array}$ <p>Proline (Pro / P)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{CH}_2 \\   \\ \text{COOH} \end{array}$ <p>Glutamic Acid (Glu / E)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{COOH} \end{array}$ <p>Aspartic Acid (Asp / D)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{H} - \text{C} - \text{OH} \\   \\ \text{CH}_3 \end{array}$ <p>Threonine (Thr / T)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{SH} \end{array}$ <p>Cysteine (Cys / C)</p>
$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{CH}_2 \\   \\ \text{S} \\   \\ \text{CH}_3 \end{array}$ <p>Methionine (Met / M)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array}$ <p>Leucine (Leu / L)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH}_2 \\   \\ \text{C}=\text{O} \\   \\ \text{NH}_2 \end{array}$ <p>Asparagine (Asn / N)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{HC} - \text{CH}_3 \\   \\ \text{CH}_2 \\   \\ \text{CH}_3 \end{array}$ <p>Isoleucine (Ile / I)</p>	$\begin{array}{c} \text{H} \\   \\ \text{H}_3\text{N}^+ - \alpha\text{C} - \text{C} \begin{array}{l} \diagup \text{O} \\ \diagdown \text{O} \end{array} \\   \\ \text{CH} \\ / \quad \backslash \\ \text{CH}_3 \quad \text{CH}_3 \end{array}$ <p>Valine (Val / V)</p>

# How DNA/RNA codes for protein?

- DNA alphabet contains four letters but must specify protein, or polypeptide sequence of 20 letters.
- Dinucleotides are not enough:  $4^2 = 16$  possible dinucleotides
- Trinucleotides (triplets) allow  $4^3 = 64$  possible trinucleotides
- Triplets are also called *codons*

		Second letter				
		U	C	A	G	
First letter	U	UUU Phenyl-alanine UUC UUA Leucine UUG	UCU UCC Serine UCA UCG	UAU Tyrosine UAC UAA Stop codon UAG Stop codon	UGU Cysteine UGC UGA Stop codon UGG Tryptophan	U C A G
	C	CUU CUC Leucine CUA CUG	CCU CCC Proline CCA CCG	CAU Histidine CAC CAA Glutamine CAG	CGU CGC Arginine CGA CGG	U C A G
	A	AUU Isoleucine AUC AUA AUG Methionine; start codon	ACU ACC Threonine ACA ACG	AAU Asparagine AAC AAA Lysine AAG	AGU Serine AGC AGA Arginine AGG	U C A G
	G	GUU Valine GUC GUA GUG	GCU GCC Alanine GCA GCG	GAU Aspartic acid GAC GAA Glutamic acid GAG	GGU GGC Glycine GGA GGG	U C A G

# How DNA/RNA codes for protein?



- Three of the possible triplets specify "stop translation"
- Translation usually starts at triplet AUG (this codes for methionine)
- Most amino acids may be specified by more than triplet
- How to find a gene? Look for start and stop codons (not that easy though)

		Second letter				
		U	C	A	G	
First letter	U	UUU Phenyl-alanine UUC UUA Leucine UUG	UCU UCC Serine UCA UCG	UAU Tyrosine UAC UAA Stop codon UAG Stop codon	UGU Cysteine UGC UGA Stop codon UGG Tryptophan	U C A G
	C	CUU CUC Leucine CUA CUG	CCU CCC Proline CCA CCG	CAU Histidine CAC CAA Glutamine CAG	CGU CGC Arginine CGA CGG	U C A G
	A	AUU Isoleucine AUC AUA AUG Methionine; start codon	ACU ACC Threonine ACA ACG	AAU Asparagine AAC AAA Lysine AAG	AGU Serine AGC AGA Arginine AGG	U C A G
	G	GUU Valine GUC GUA GUG	GCU GCC Alanine GCA GCG	GAU Aspartic acid GAC GAA Glutamic acid GAG	GGU GGC Glycine GGA GGG	U C A G



# Proteins: Workhorses of the Cell

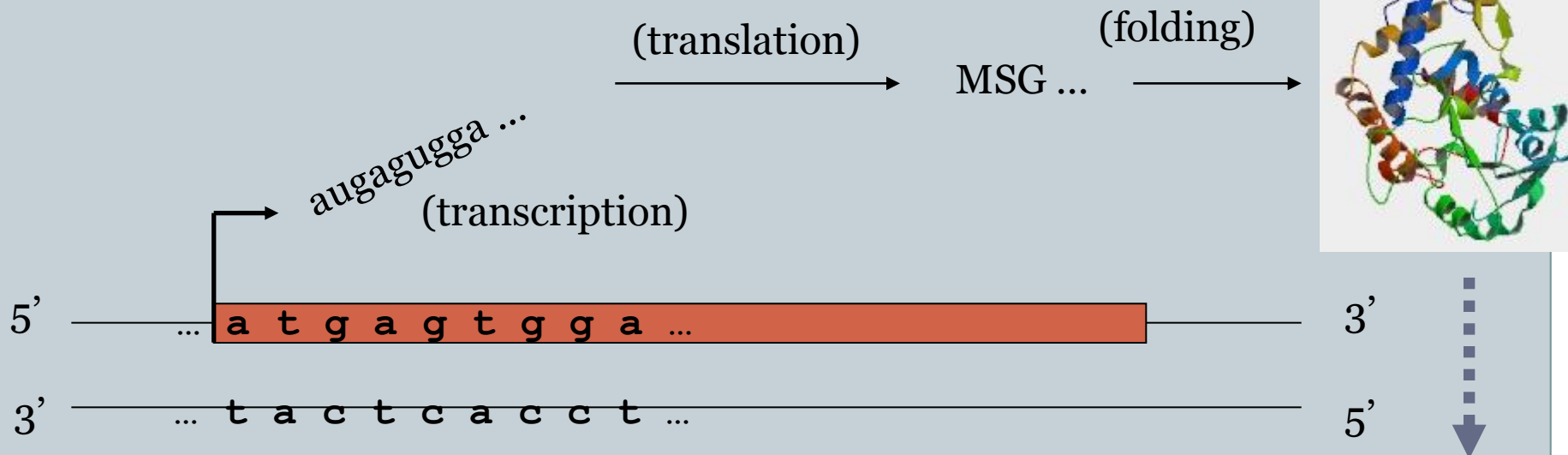


- 20 different **amino acids**
  - different chemical properties cause the protein chains to fold up into specific three-dimensional structures that define their particular functions in the cell.
- Proteins do all essential work for the cell
  - build cellular structures
  - digest nutrients
  - execute metabolic functions
  - mediate information flow within a cell and among cellular communities.
- Proteins work together with other proteins or nucleic acids as "molecular machines"
  - structures that fit together and function in highly specific, lock-and-key ways.

# Genes



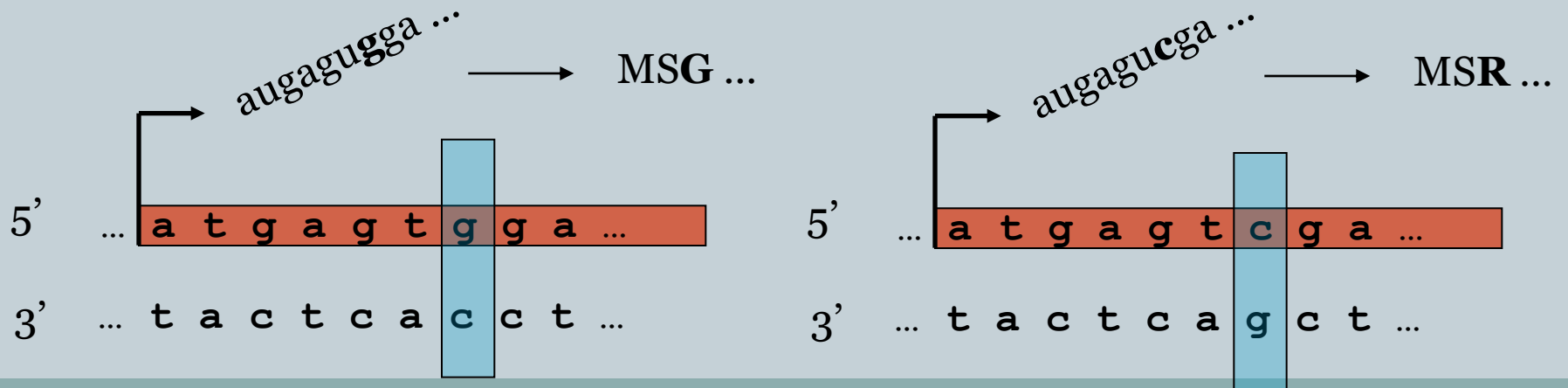
- “A gene is a union of genomic sequences encoding a coherent set of potentially overlapping functional products” --Gerstein et al.
- A DNA segment whose information is expressed either as an RNA molecule or protein



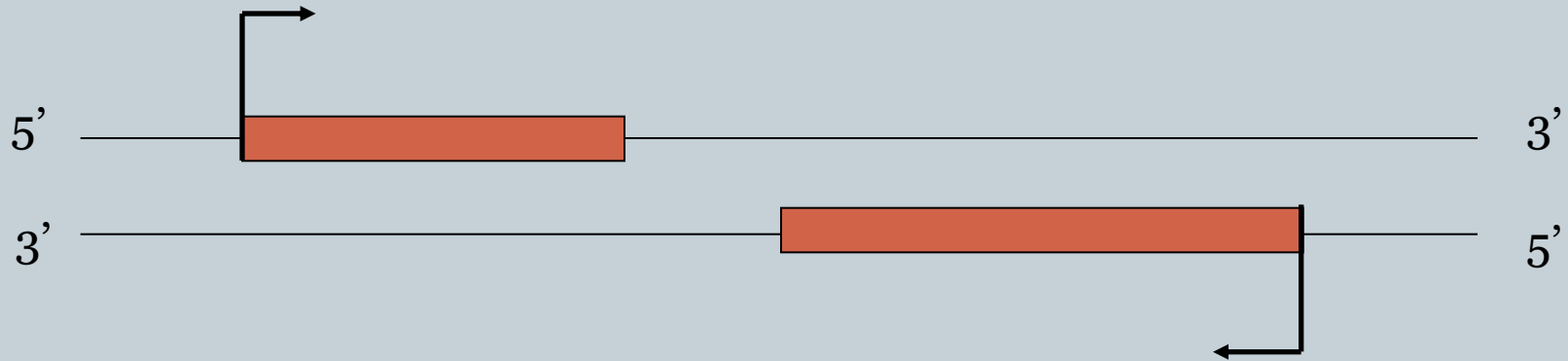
# Genes & alleles



- A gene can have different variants
- The variants of the same gene are called *alleles*



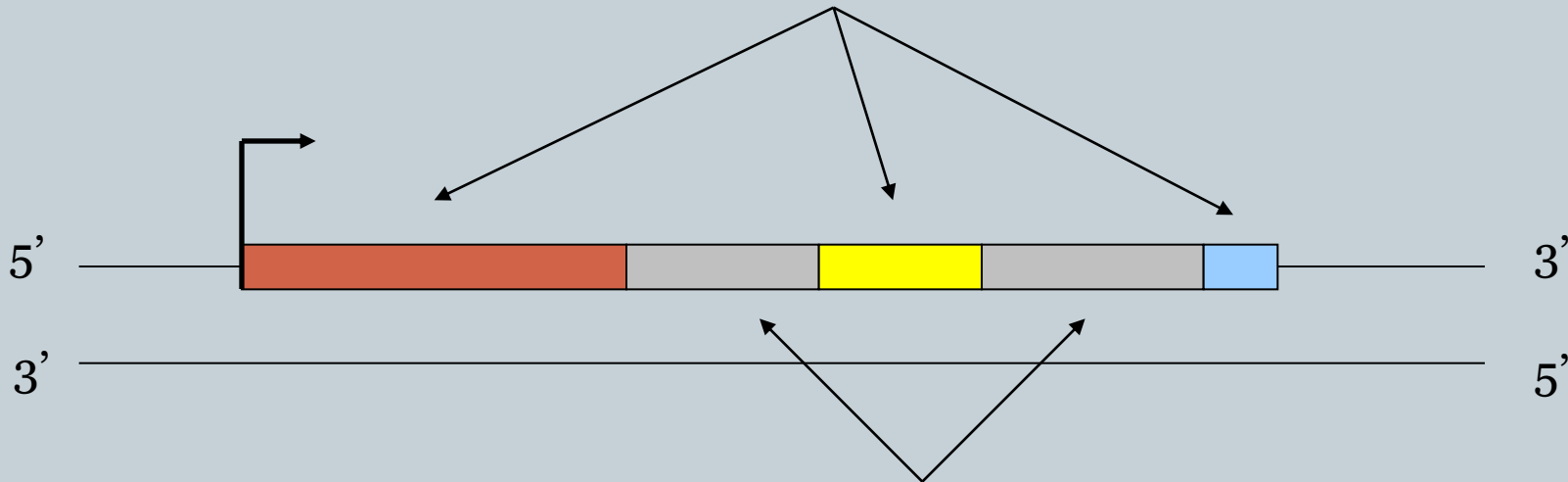
# Genes can be found on both strands



# Exons and introns & splicing



Exons



Introns are removed from RNA after transcription

Exons are joined:

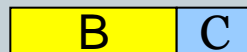
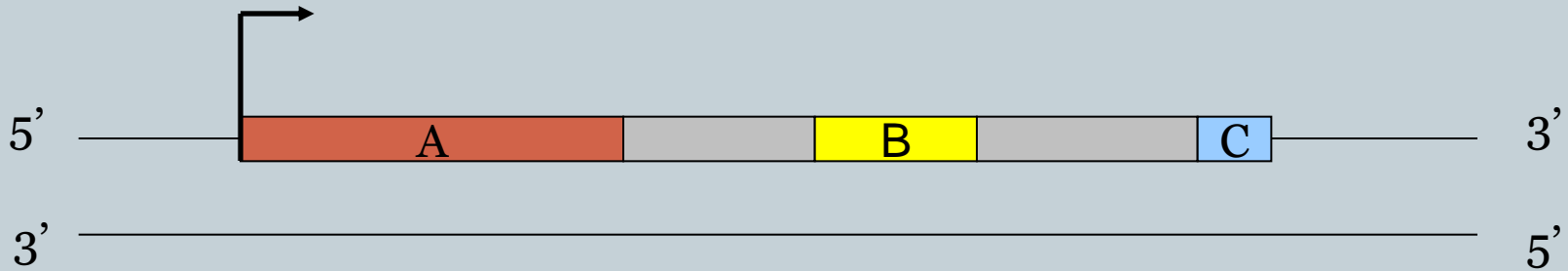


This process is called *splicing*

# Alternative splicing



Different *splice variants* may be generated

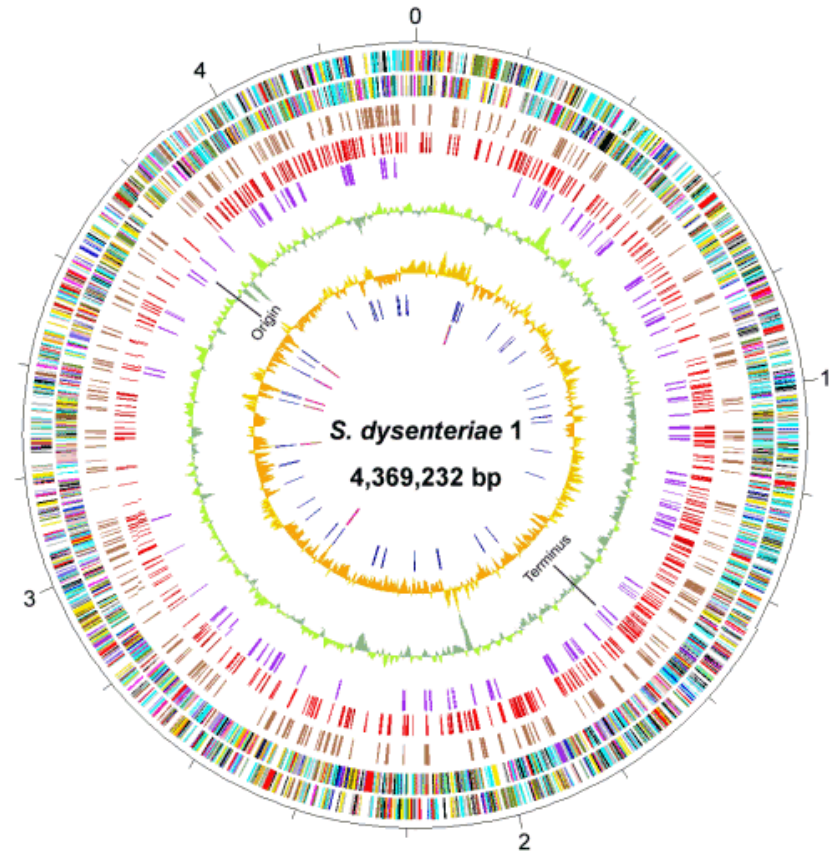


...

# Where does the variation in genomes come from?



- Prokaryotes are typically haploid: they have a single (circular) chromosome
- DNA is usually inherited vertically (parent to daughter)
- Inheritance is clonal
  - Descendants are faithful copies of an ancestral DNA
  - Variation is introduced via mutations, transposable elements, and horizontal transfer of DNA



Chromosome map of *S. dysenteriae*, the nine rings describe different properties of the genome  
[http://www.mgc.ac.cn/ShiBASE/circular\\_Sd197.htm](http://www.mgc.ac.cn/ShiBASE/circular_Sd197.htm)



# Causes of variation



- Mistakes in DNA replication
- Environmental agents (radiation, chemical agents)
- Transposable elements (transposons)
  - A part of DNA is moved or copied to another location in genome
- Horizontal transfer of DNA
  - Organism obtains genetic material from another organism that is not its parent
  - Utilized in genetic engineering

# Biological string manipulation

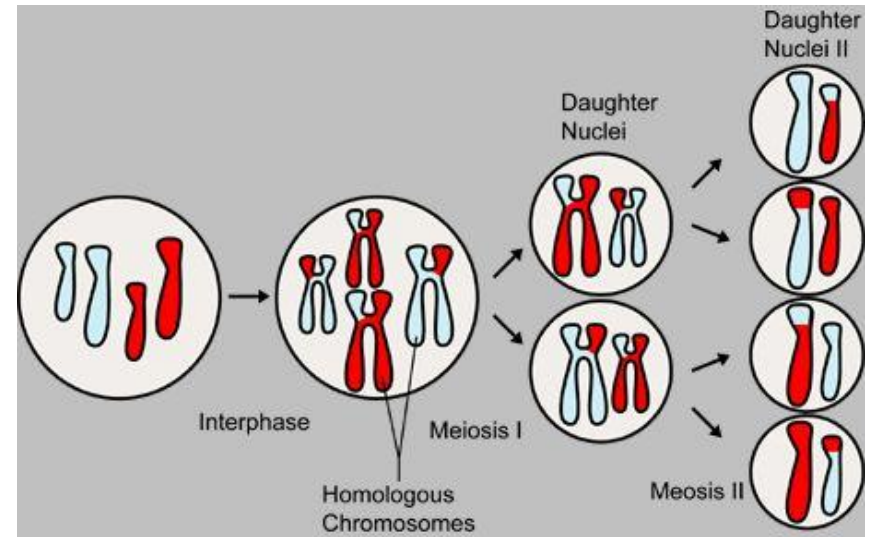


- Point mutation: substitution of a base
  - ...ACG**G**CT... => ...ACG**C**CT...
- Deletion: removal of one or more contiguous bases (substring)
  - ...TT**G**ATCA... => ...TTTCA...
- Insertion: insertion of a substring
  - ...GGCTAG... => ...GG**TCAACT**AG...

# Meiosis



- Sexual organisms are usually diploid
  - Germline cells (gametes) contain  $N$  chromosomes
  - Somatic (body) cells have  $2N$  chromosomes
- Meiosis: reduction of chromosome number from  $2N$  to  $N$  during reproductive cycle
  - One chromosome doubling is followed by two cell divisions



Major events in meiosis

<http://en.wikipedia.org/wiki/Meiosis>

<http://www.ncbi.nlm.nih.gov/About/Primer>

# Recombination and variation



- Recap: Allele is a viable DNA coding occupying a given locus (position in the genome)
- In recombination, alleles from parents become shuffled in offspring individuals via chromosomal crossover over
- Allele combinations in offspring are usually different from combinations found in parents
- Recombination errors lead into additional variations

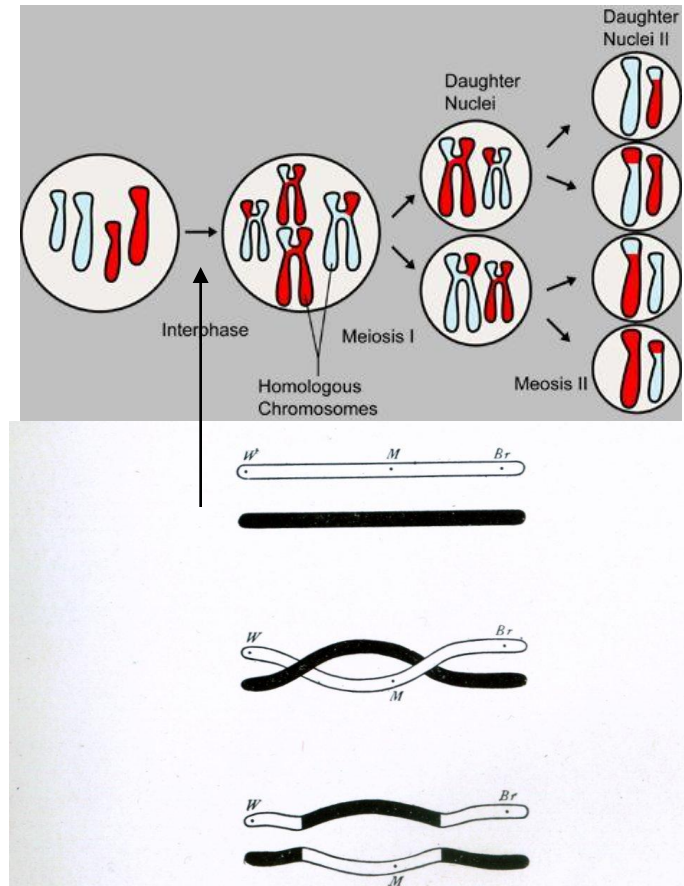
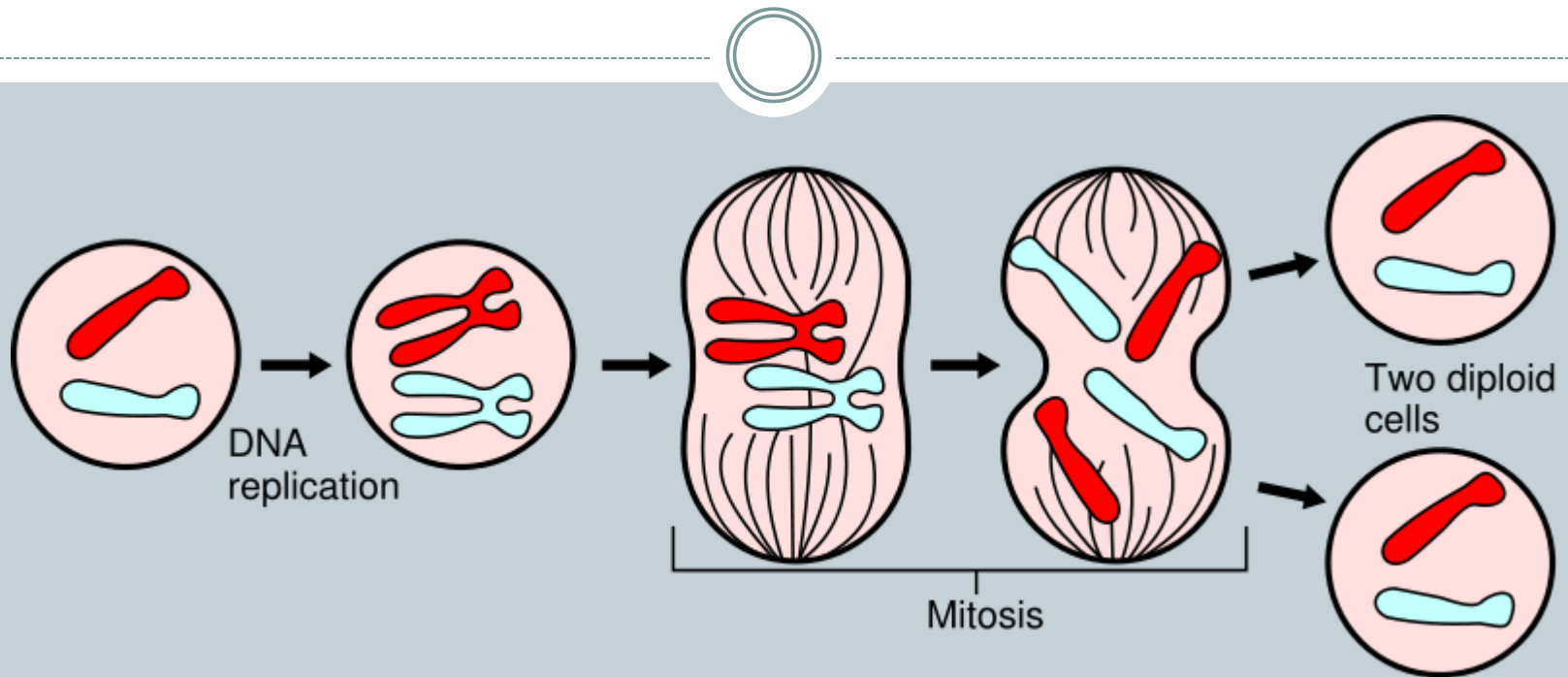


FIG. 65. Scheme to illustrate double crossing over.

Chromosomal crossover as described by  
T. H. Morgan in 1916

# Mitosis



- Mitosis: growth and development of the organism
  - One chromosome doubling is followed by one cell division