

Introduction to Bioinformatics and Molecular Biology

Jianlin Cheng, PhD

School of Electrical Engineering and Computer Science
University of Central Florida



2006

Free for academic use. Copyright @ Jianlin Cheng & original sources for some materials

Goals

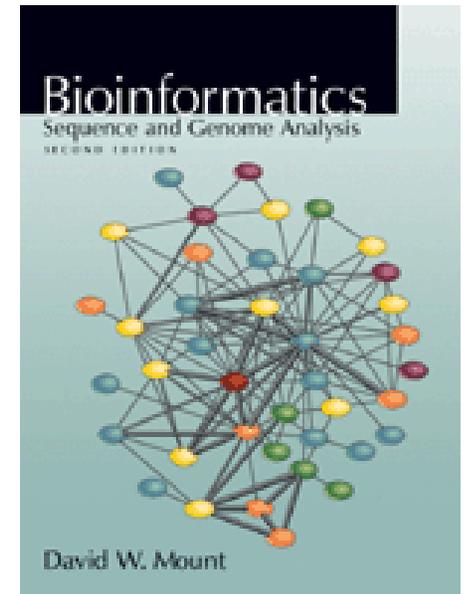
- Introduce fundamental problems, concepts, methods, and applications in Bioinformatics
- Emphasize both the methods and the practical use of bioinformatics tools and databases.
- Audience: computer scientists, engineers, biologists, statisticians, ...
- Prerequisite: some background in programming OR molecular biology.

Ten Topics

1. Introduction to Molecular Biology and Bioinformatics
2. Pairwise Sequence Alignment Using Dynamic Programming
3. Practical Sequence/Profile Alignment Using Fast Heuristic Methods (BLAST and PSI-BLAST)
4. Multiple Sequence Alignment
5. Gene and Motif Identification
6. Phylogenetic Analysis
7. Protein Structure Analysis and Prediction
8. RNA Secondary Structure Prediction
9. Clustering and Classification of Gene Expression Data
10. Search and Mining of Biological Databases, Databanks, and Literature

Resources

- Slides (**main materials**, http://www.eecs.ucf.edu/~jcheng/cheng_courses.html)
- A general introduction book (David Mount. *Bioinformatics: Sequence and Genome Analysis*, 2004. additional reading materials)
- Other online Bioinformatics courses (additional reading materials)
- Reference books (complementary)



Grading

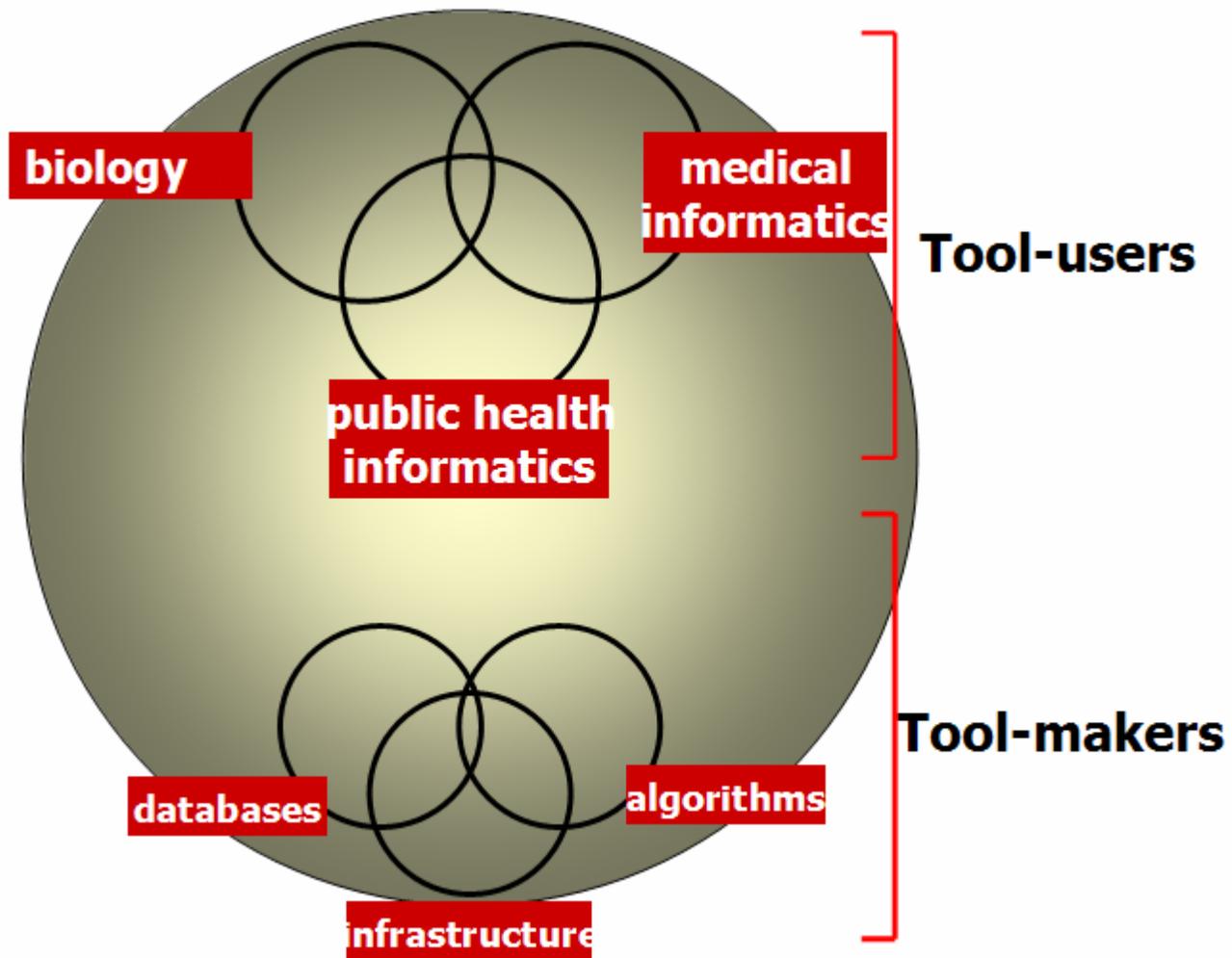
- Homework (20%)
- Project (20%)
- Midterm (30%)
- Final (30%)

What's Bioinformatics?

An interdisciplinary science of developing and applying computational techniques to address problems in molecular biology (or more broadly biology: computational biology).

Two main fields:

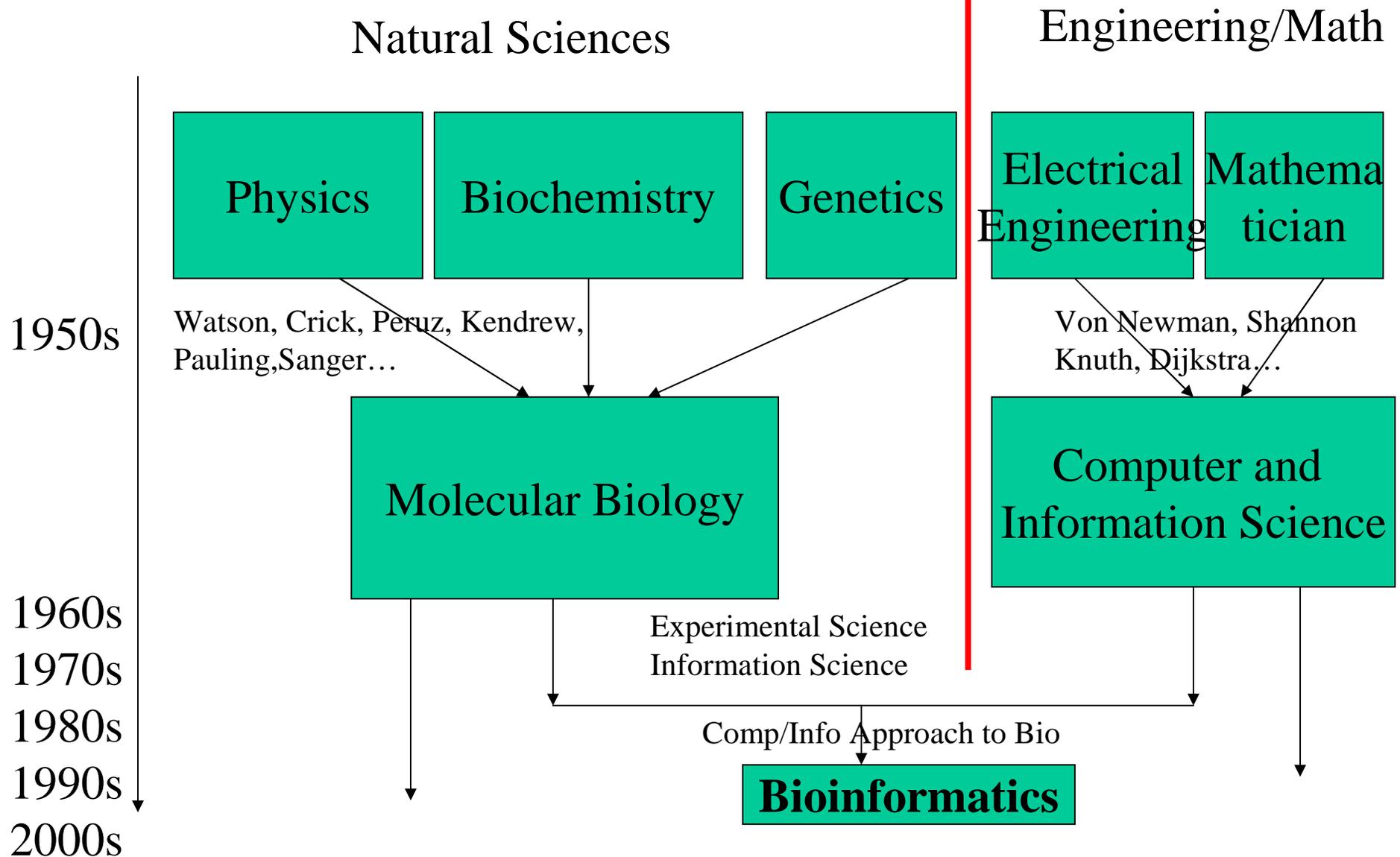
- Develop bioinformatics algorithms and tools
- Apply bioinformatics tools to address biological problems



Adapted from J. Pevsner, 2005

History of Bioinformatics

How does a new interdisciplinary science emerge?



Major Events in Bioinformatics History

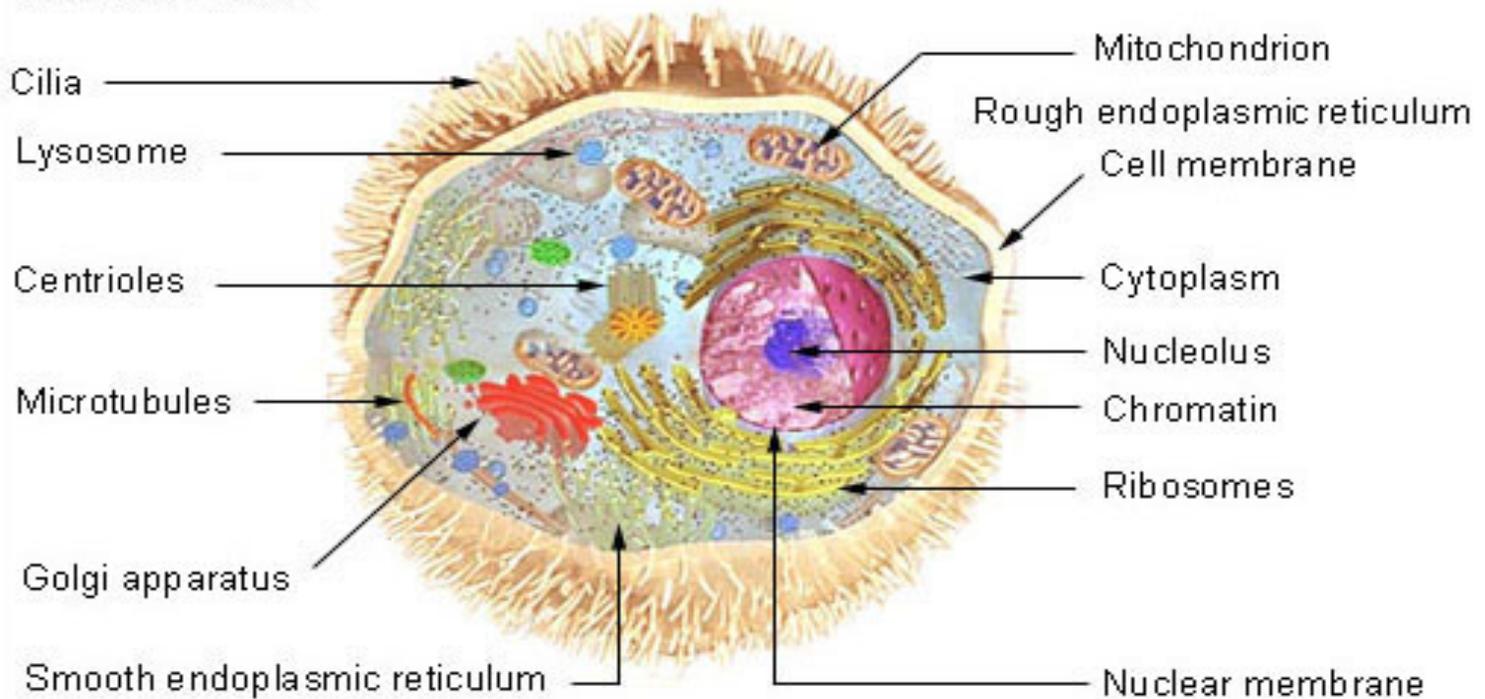
<http://www.ncbi.nlm.nih.gov/Education/BLASTinfo/milestones.html>

- 1962** Pauling's theory of molecular evolution
- 1965** Margaret Dayhoff's Atlas of Protein Sequences
- 1970** Needleman-Wunsch algorithm
- 1977** DNA sequencing and software to analyze it (Sateen)
- 1981** Smith-Waterman algorithm developed
- 1981** The concept of a sequence motif (Doolittle)
- 1982** GenBank Release 3 made public
- 1982** Phage lambda genome sequenced
- 1983** Sequence database searching algorithm (Wilbur-Lipman)
- 1985** FASTP/FASTN: fast sequence similarity searching
- 1988** National Center for Biotechnology Information (NCBI) created at NIH/NLM
- 1988** EMBnet network for database distribution
- 1990** BLAST: fast sequence similarity searching
- 1991** EST: expressed sequence tag sequencing
- 1993** Sanger Centre, Hinxton, UK
- 1994** EMBL European Bioinformatics Institute, Hinxton, UK
- 1995** First bacterial genomes completely sequenced
- 1996** Yeast genome completely sequenced
- 1997** PSI-BLAST
- 1998** Worm (multicellular) genome completely sequenced
- 1999** Fly genome completely sequenced
- 2004** Human genome sequenced

Introduction to Molecular Biology

- Cell is the unit of structure and function of all living things.

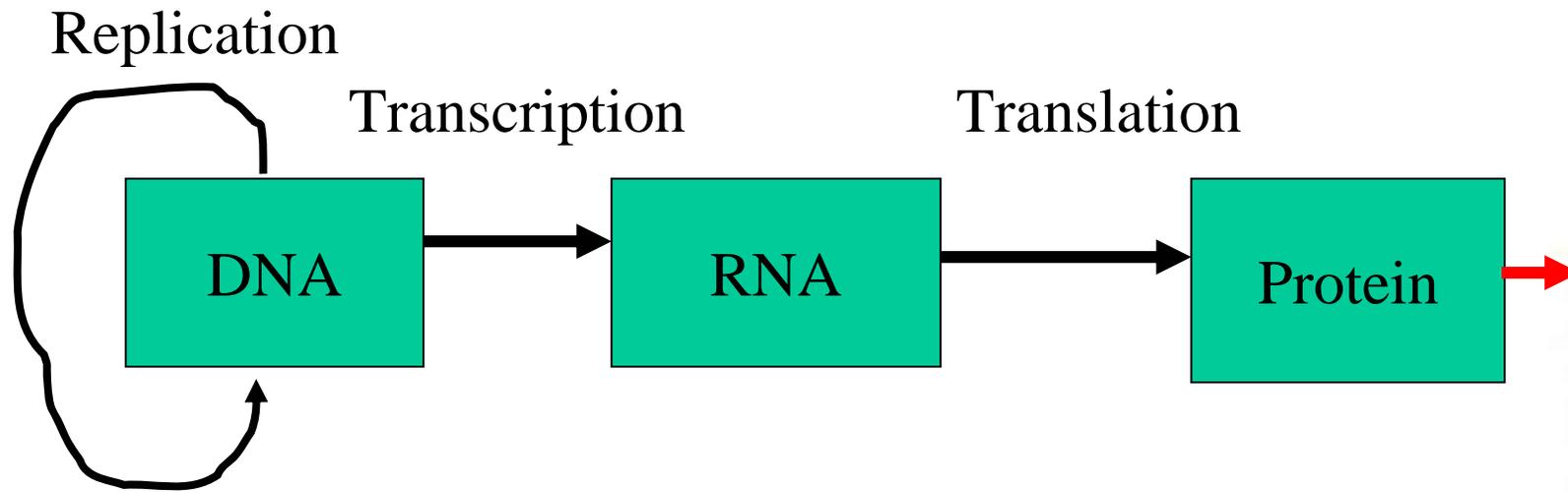
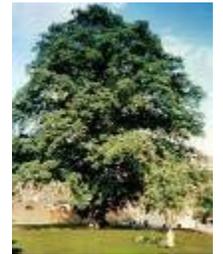
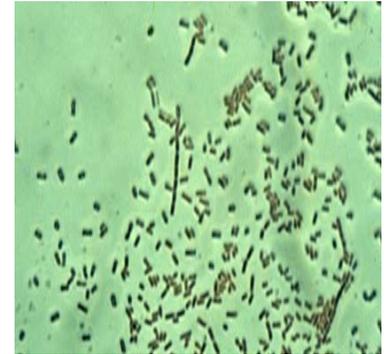
Cell Structure



Two types of cells: eukaryote (higher organisms) and prokaryote (lower organisms)

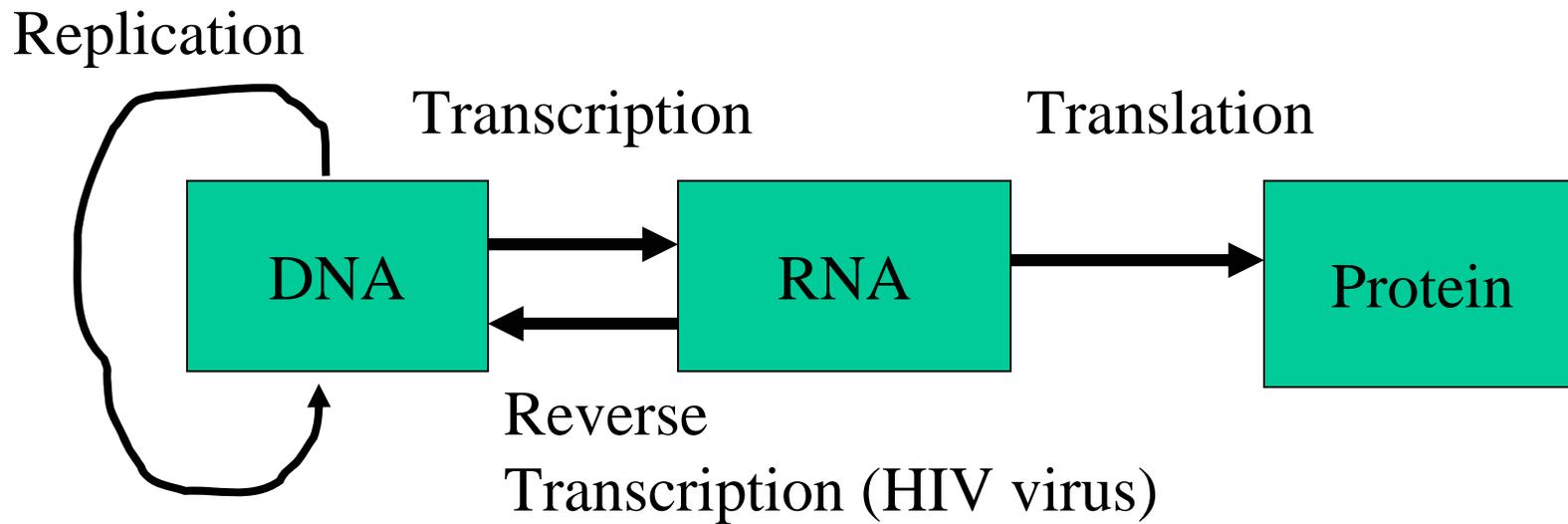
Central Dogma of Molecular Biology

Phenotype

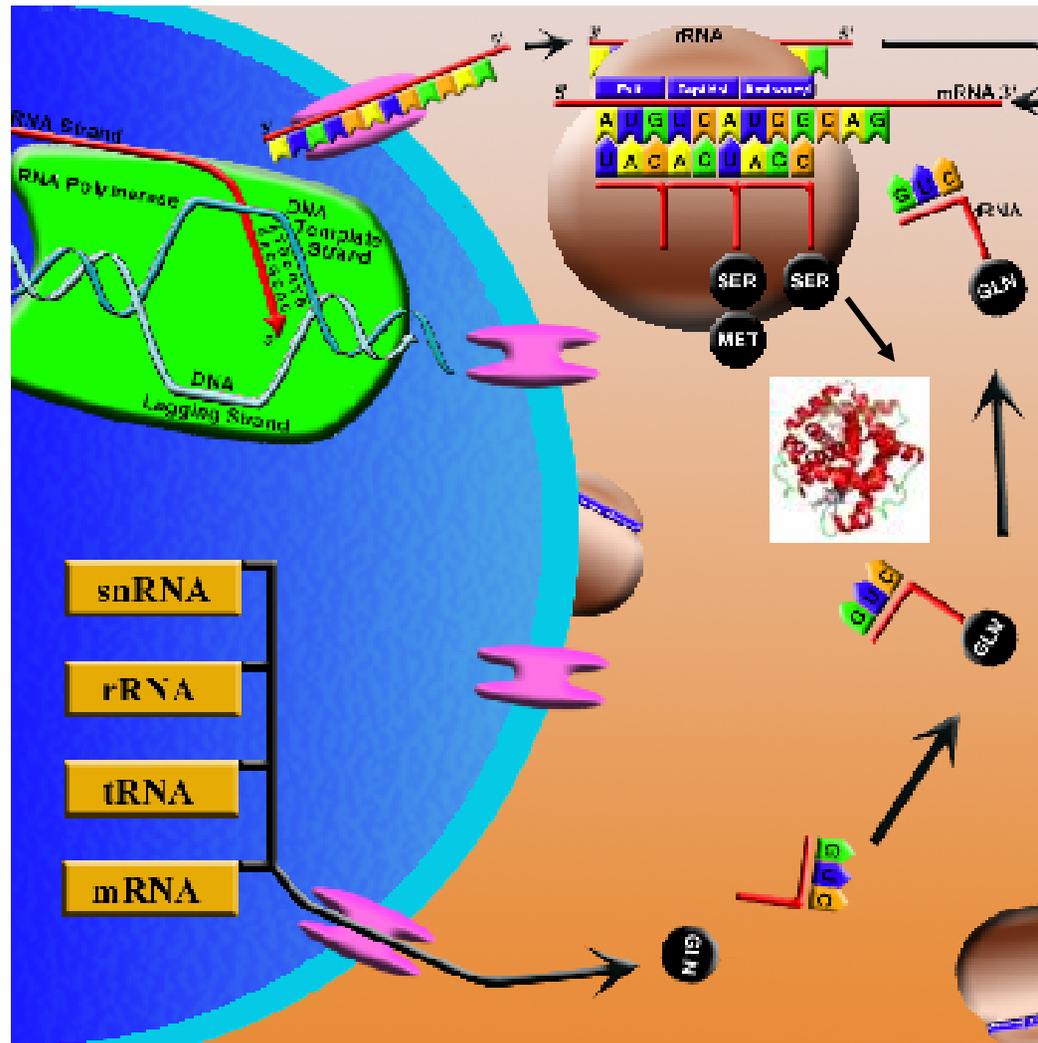


Genotype

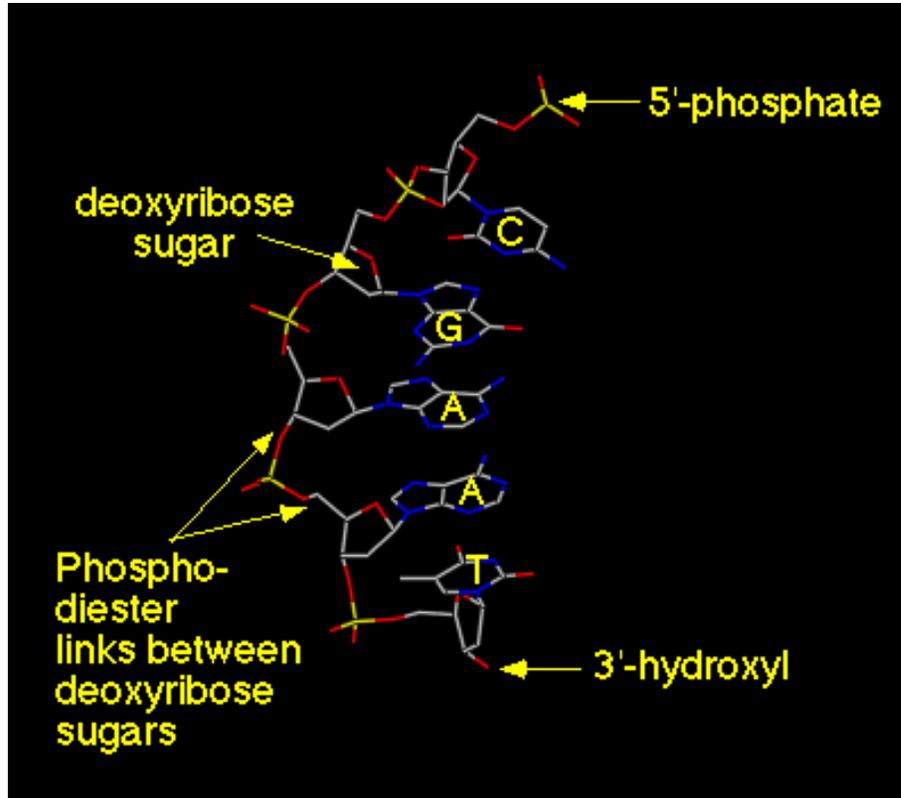
Central Dogma of Molecular Biology



Information flow



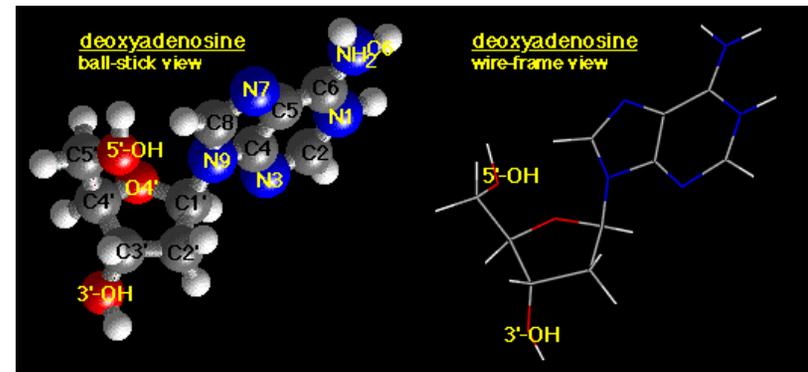
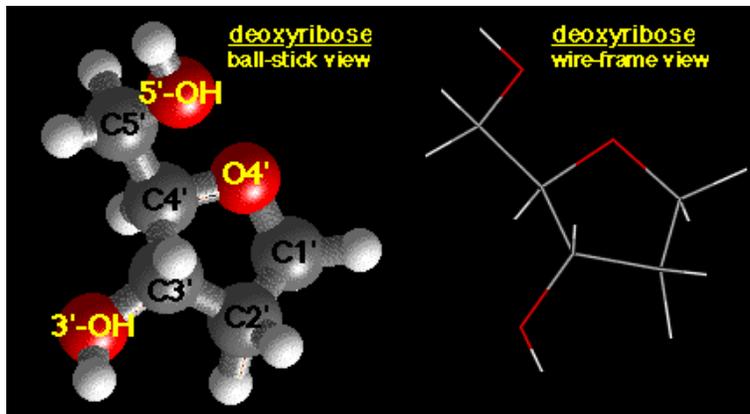
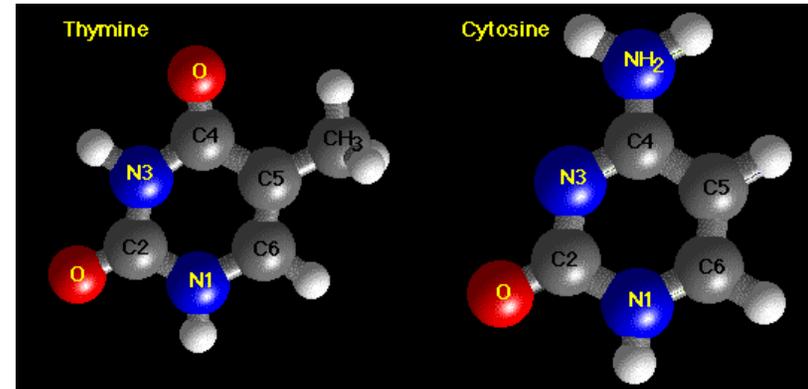
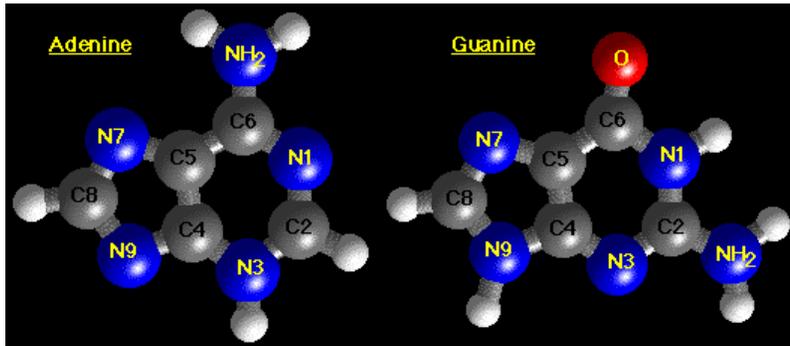
DNA (Deoxyribose Nucleotide Acids)

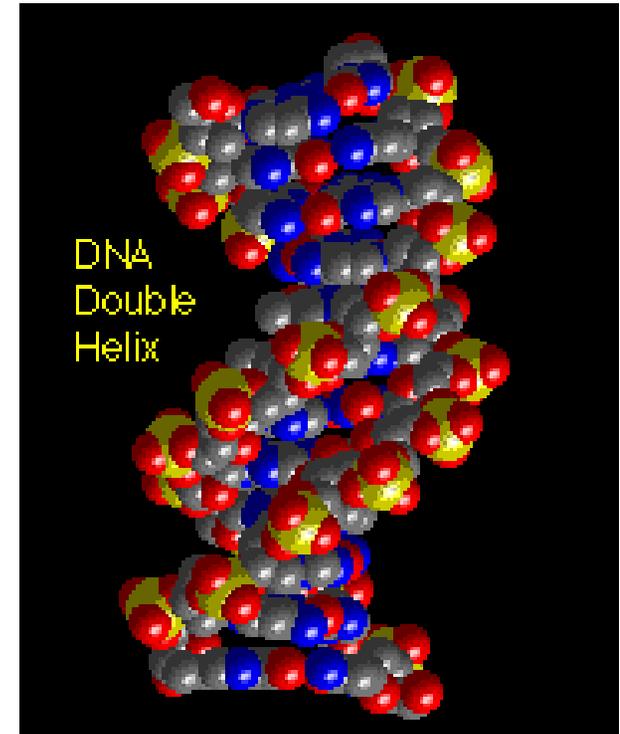
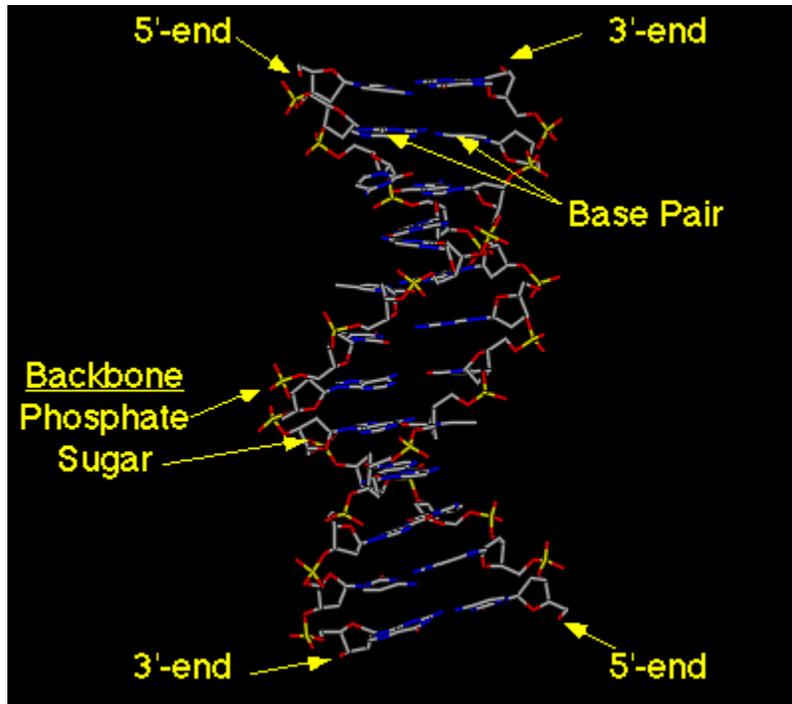


DNA is a polymer. The monomer units of DNA are nucleotides, and the polymer is known as a "polynucleotide." Each nucleotide consists of a 5-carbon sugar (deoxyribose), a nitrogen containing base attached to the sugar, and a phosphate group.

A is for adenine
G is for guanine
C is for cytosine
T is for thymine

CGAATGGGAAA.....



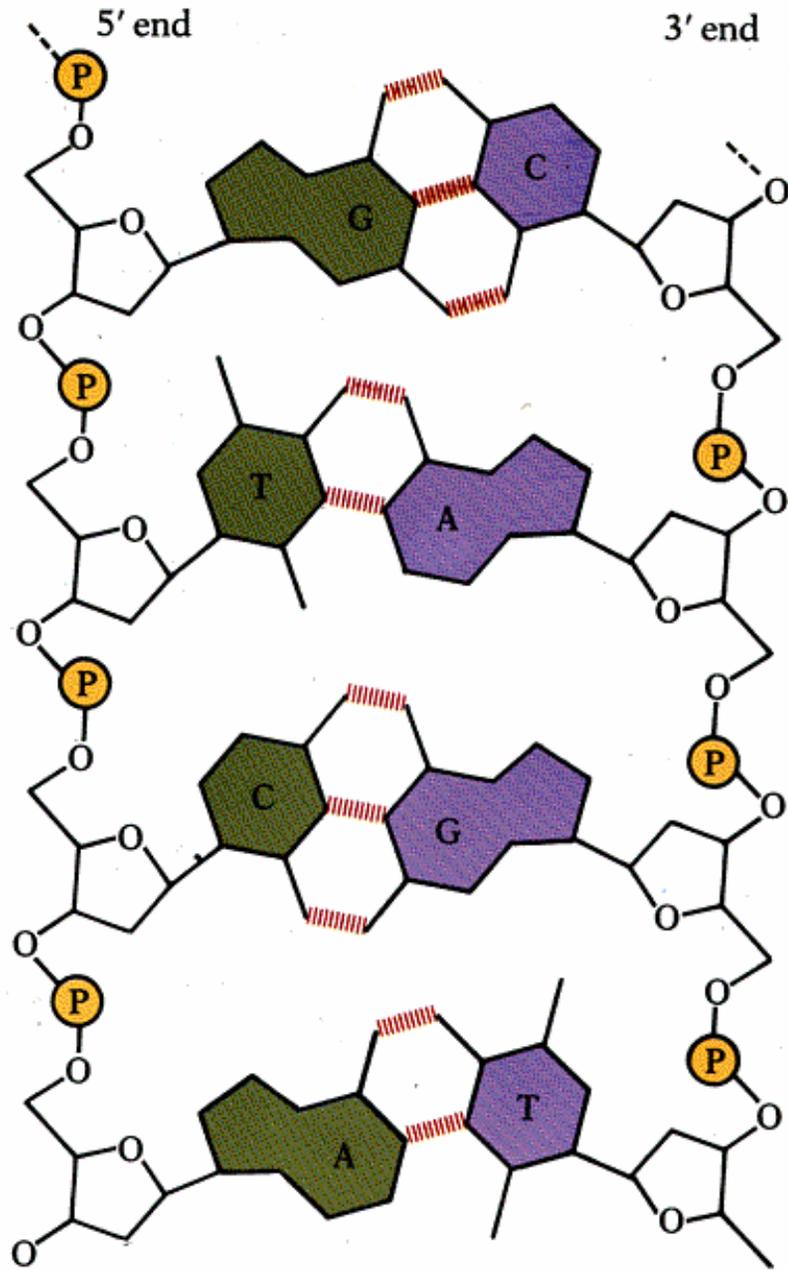


Base Pairs:

A-T (2 H-bonds)

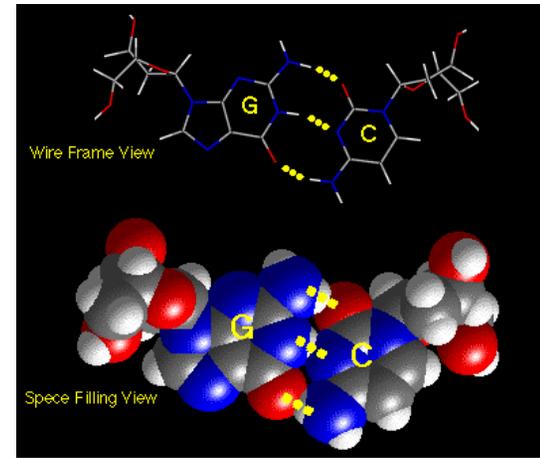
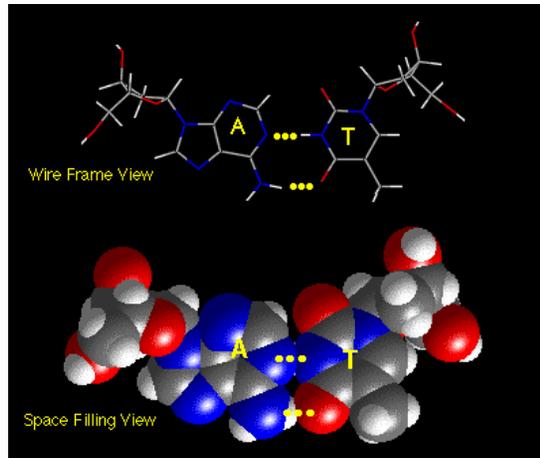
C-G (3 H-bonds)

Hydrogen bonds: non-covalent bonds mediated by hydrogen atoms



Uncoiled DNA Molecule

Source: Dr. Gary Stormo, 2002



James Watson & Francis Crick



Maurice
Wilkins



Rosalind
Franklin



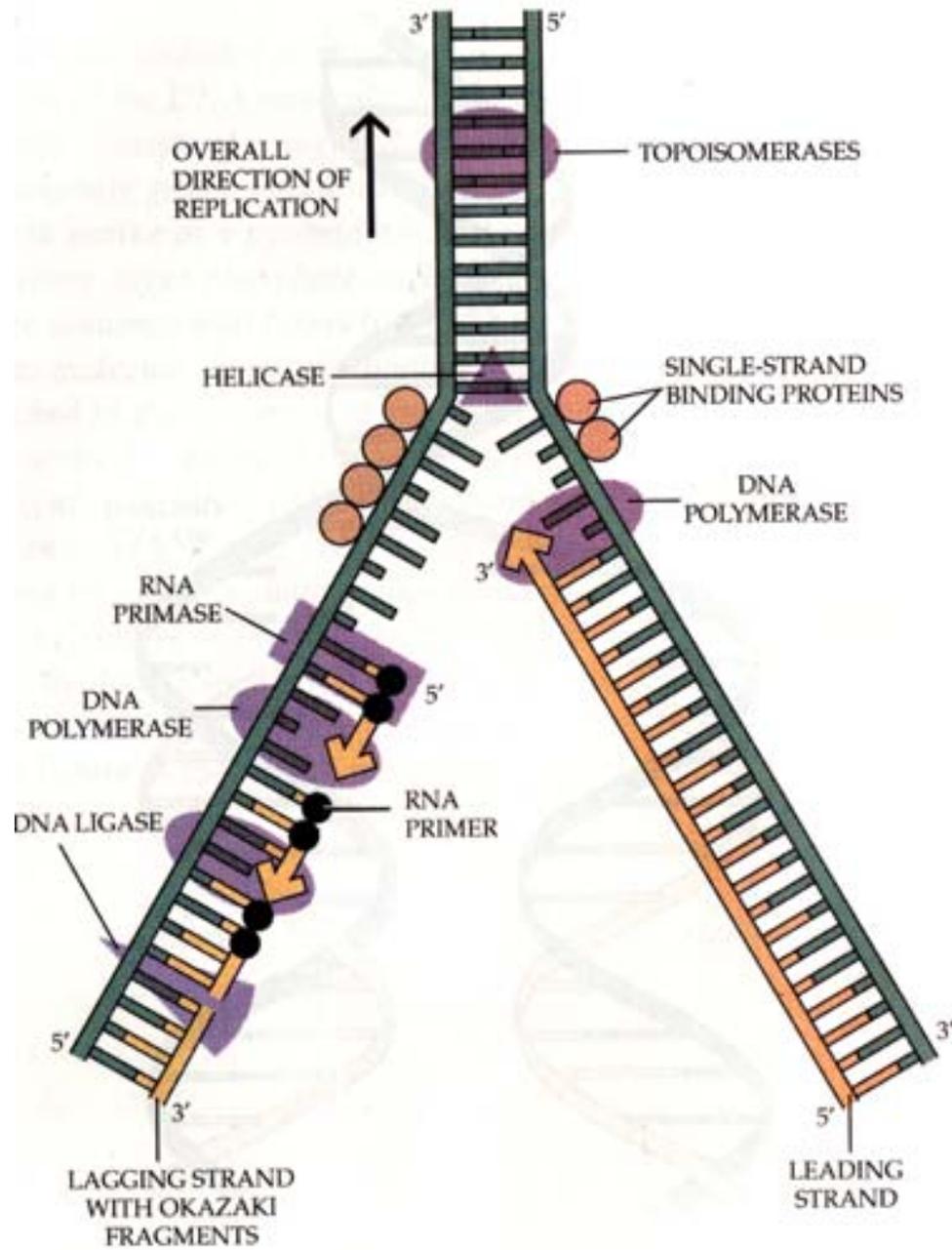
Linus
Pauling



Erwin
Chargaff

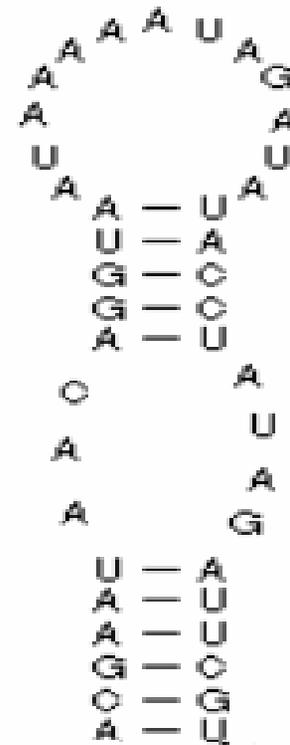
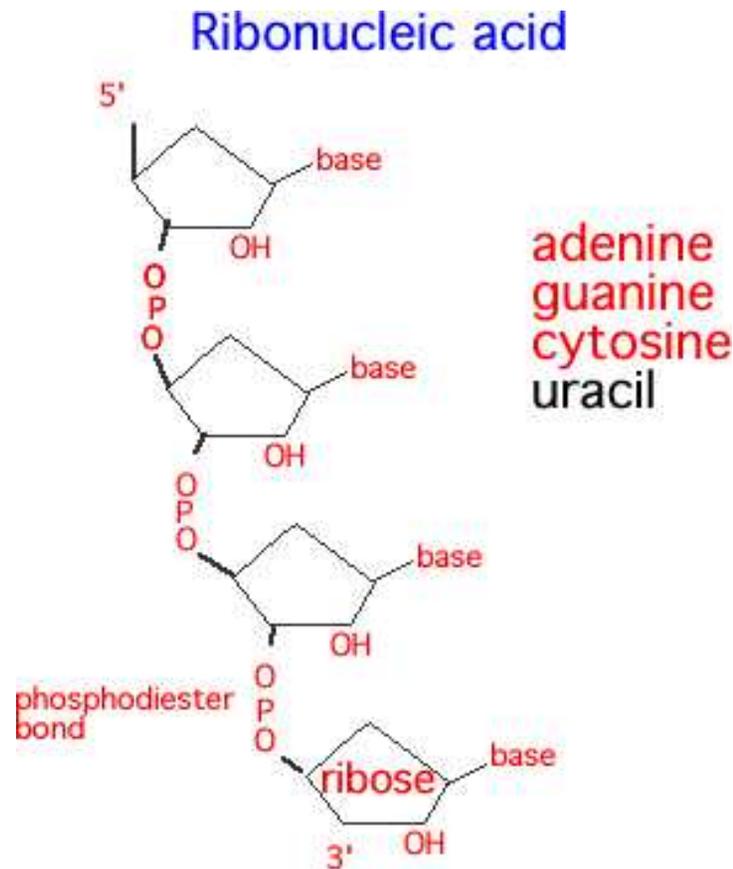
Fundamental Problems: How genetic information pass from one cell to another and from one generation to next generation

DNA Replication



DNA
Polymerase

RNA (Ribose Nucleotide Acids)

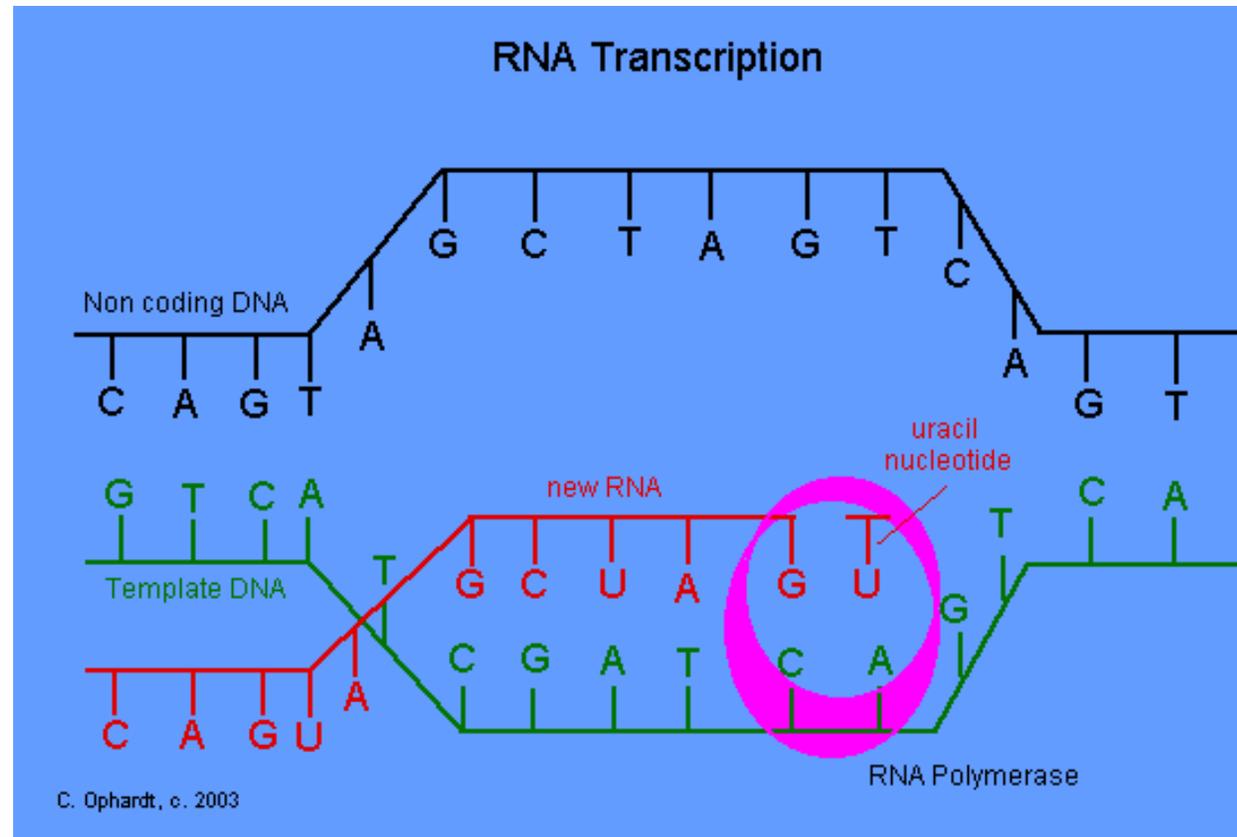


ACGAAUAACAGGUAUAUAAAAUAGAUUACCUAUAGAUUCGU

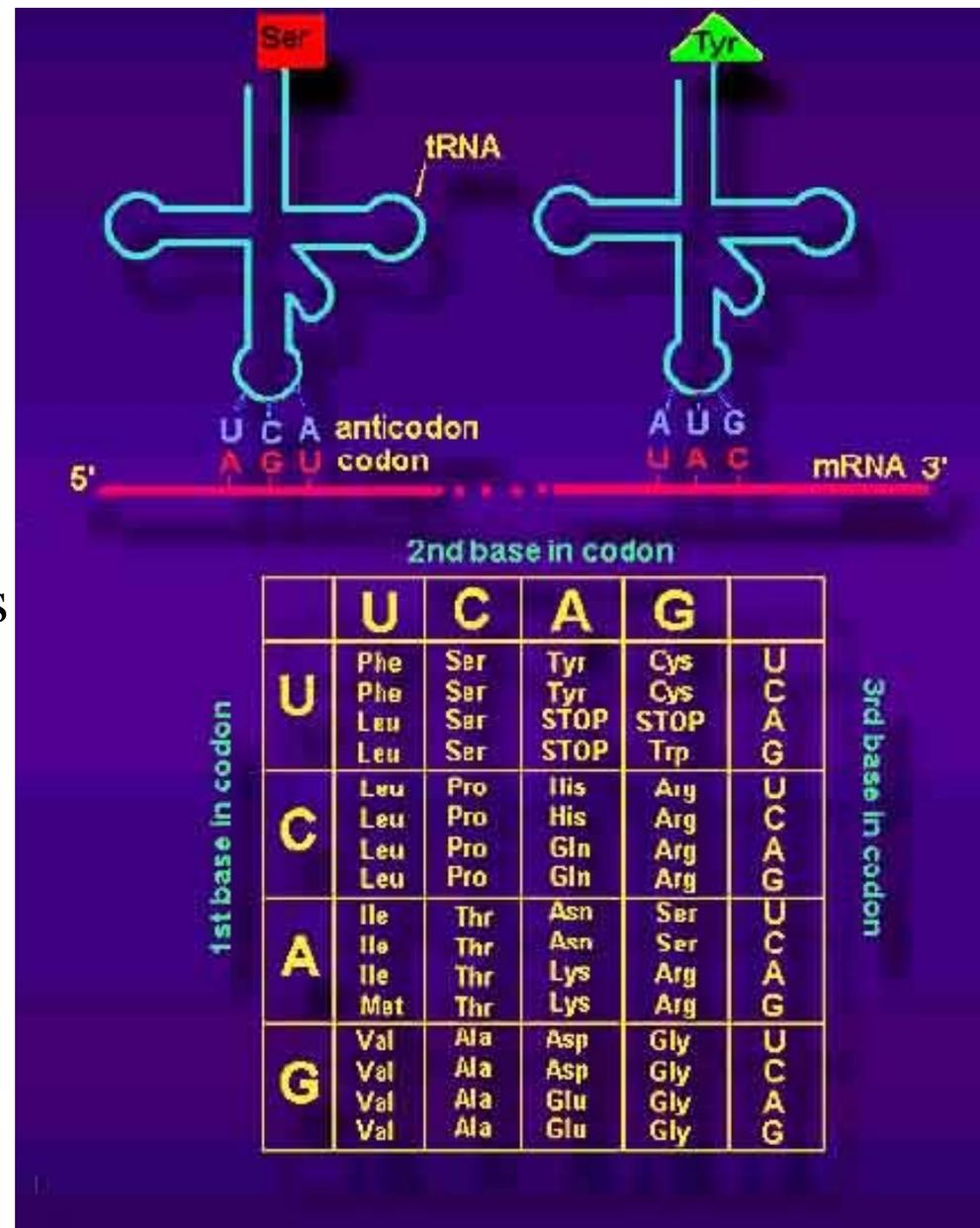
Different Kinds of RNA

- mRNA: messenger RNA
carry genetic information out of nucleus for protein synthesis (transcription process: RNA polymerase)
- rRNA: ribosomal RNA
constitute 50% of ribosome, which is a molecular assembly for protein synthesis
- tRNA: transfer RNA
decode information (map 3 nucleotides to amino acid); transfer amino acid
- snRNA: small RNA molecules found in nucleus
involve RNA splicing

Transcription of Gene into RNA



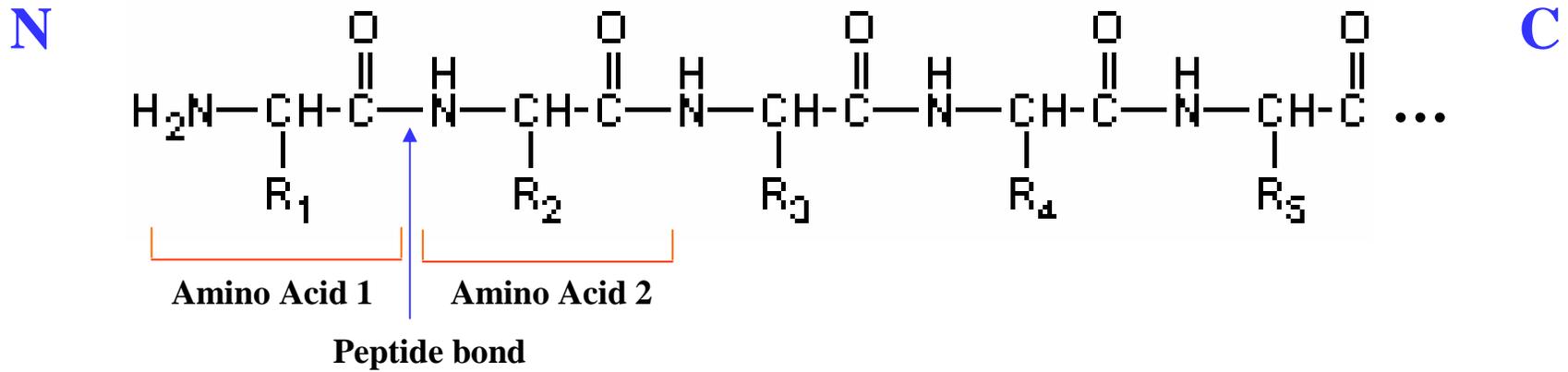
Genetic Code and Translation



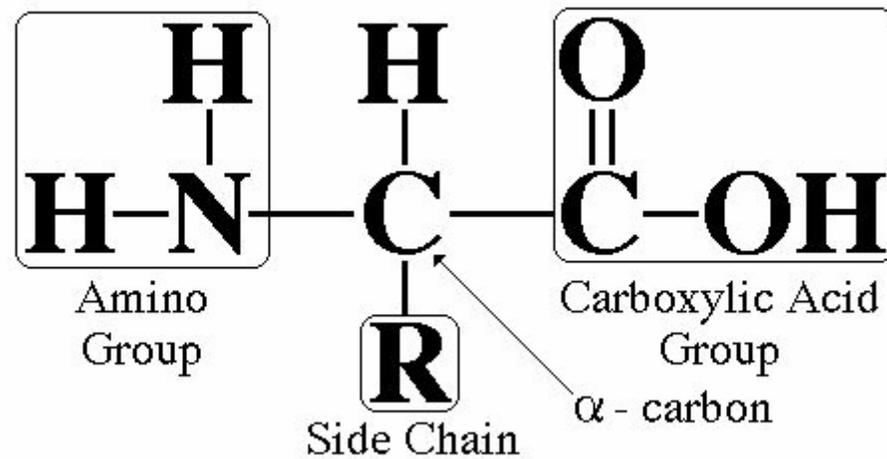
Three Nucleotides is called a codon.

Protein Sequence

A directional sequence of amino acids/residues



Amino Acid Structure



Amino Acids

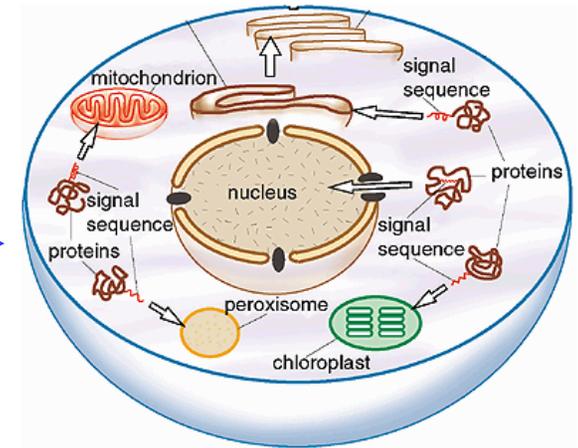
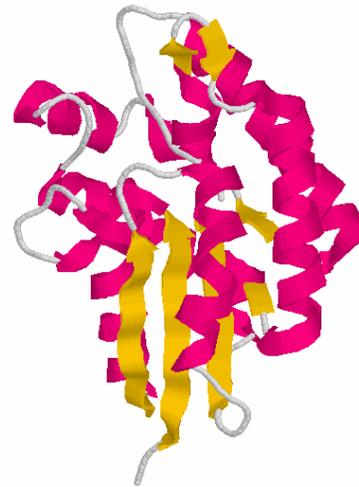
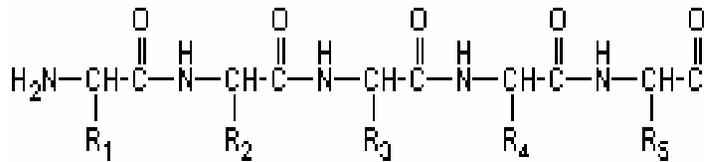
Amino acid	Abbrev.	Side chain	Hydrophobic	Polar	Charged	Small	Tiny	Aromatic or Aliphatic	van der Waals volume	Codon	Occurrence in proteins (%)
Alanine	Ala, A	-CH ₃	X	-	-	X	X	-	67	GCU, GCC, GCA, GCG	7.8
Cysteine	Cys, C	-CH ₂ SH	X	-	-	X	-	-	86	UGU, UGC	1.9
Aspartate	Asp, D	-CH ₂ COOH	-	X	negative	X	-	-	91	GAU, GAC	5.3
Glutamate	Glu, E	-CH ₂ CH ₂ COOH	-	X	negative	-	-	-	109	GAA, GAG	6.3
Phenylalanine	Phe, F	-CH ₂ C ₆ H ₅	X	-	-	-	-	Aromatic	135	UUU, UUC	3.9
Glycine	Gly, G	-H	X	-	-	X	X	-	48	GGU, GGC, GGA, GGG	7.2
Histidine	His, H	-CH ₂ -C ₃ H ₃ N ₂	-	X	positive	-	-	Aromatic	118	CAU, CAC	2.3
Isoleucine	Ile, I	-CH(CH ₃)CH ₂ CH ₃	X	-	-	-	-	Aliphatic	124	AUU, AUC, AUA	5.3
Lysine	Lys, K	-(CH ₂) ₄ NH ₂	-	X	positive	-	-	-	135	AAA, AAG	5.9
Leucine	Leu, L	-CH ₂ CH(CH ₃) ₂	X	-	-	-	-	Aliphatic	124	UUA, UUG, CUU, CUC, CUA, CUG	9.1
Methionine	Met, M	-CH ₂ CH ₂ SCH ₃	X	-	-	-	-	-	124	AUG	2.3
Asparagine	Asn, N	-CH ₂ CONH ₂	-	X	-	X	-	-	96	AAU, AAC	4.3
Proline	Pro, P	-CH ₂ CH ₂ CH ₂ -	X	-	-	X	-	-	90	CCU, CCC, CCA, CCG	5.2
Glutamine	Gln, Q	-CH ₂ CH ₂ CONH ₂	-	X	-	-	-	-	114	CAA, CAG	4.2
Arginine	Arg, R	-(CH ₂) ₃ NH-C(NH) NH ₂	-	X	positive	-	-	-	148	CGU, CGC, CGA, CCG, AGA, AGG	5.1
Serine	Ser, S	-CH ₂ OH	-	X	-	X	X	-	73	UCU, UCC, UCA, UCG, AGU, AGC	6.8
Threonine	Thr, T	-CH(OH)CH ₃	X	X	-	X	-	-	93	ACU, ACC, ACA, ACG	5.9
Valine	Val, V	-CH(CH ₃) ₂	X	-	-	X	-	Aliphatic	105	GUU, GUC, GUA, GUG	6.6
Tryptophan	Trp, W	-CH ₂ C ₈ H ₆ N	X	-	-	-	-	Aromatic	163	UGG	1.4
Tyrosine	Tyr, Y	-CH ₂ -C ₆ H ₄ OH	X	X	-	-	-	Aromatic	141	UAU, UAC	3.2

Hydrophilic



Central Dogma of Proteomics

AGCWY.....



Cell

Sequence

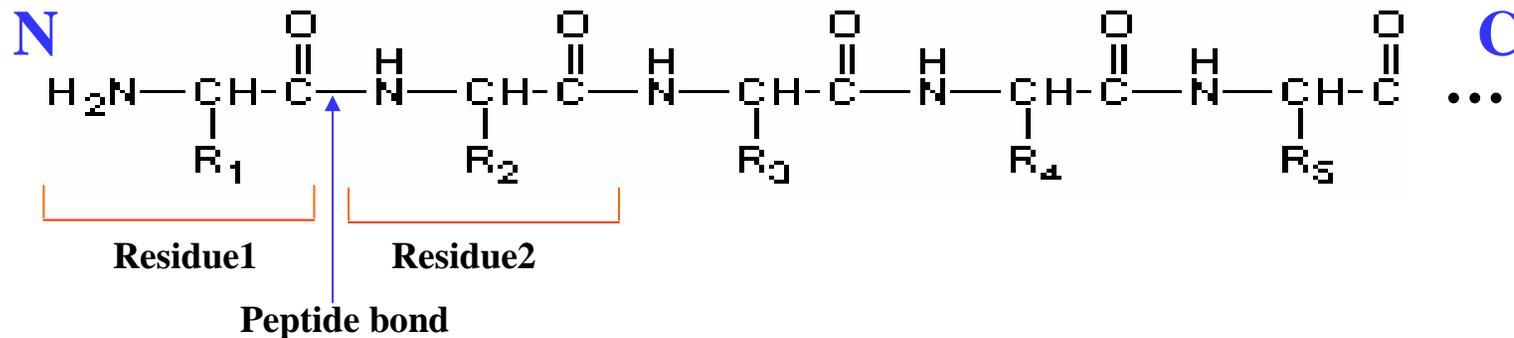
Structure

Function

From Sequence to Consequence

Four Levels of Protein Structure

Primary Structure (a directional sequence of amino acids/residues)



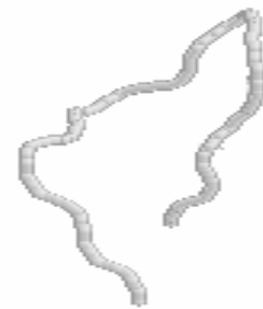
Secondary Structure (helix, strand, coil)



Alpha Helix



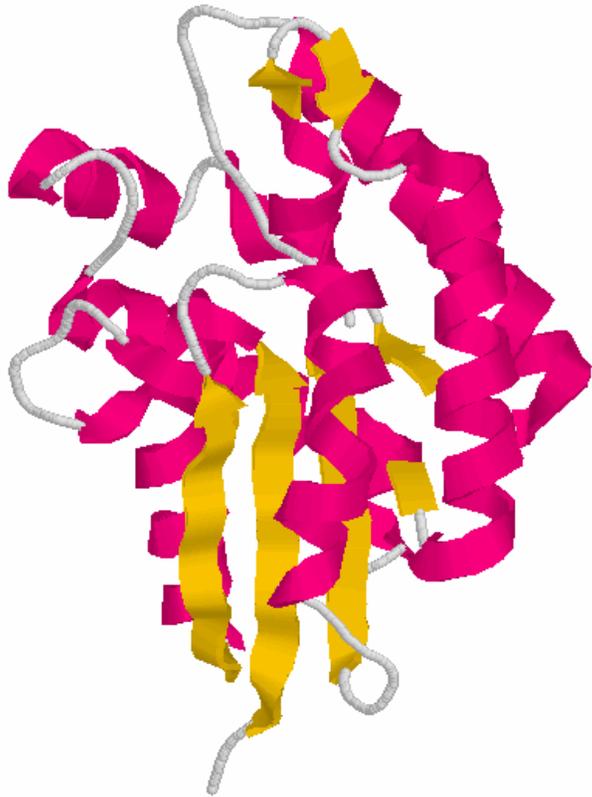
Beta Strand / Sheet



Coil

Four Levels of Protein Structure

Tertiary Structure



Quaternary Structure (complex)



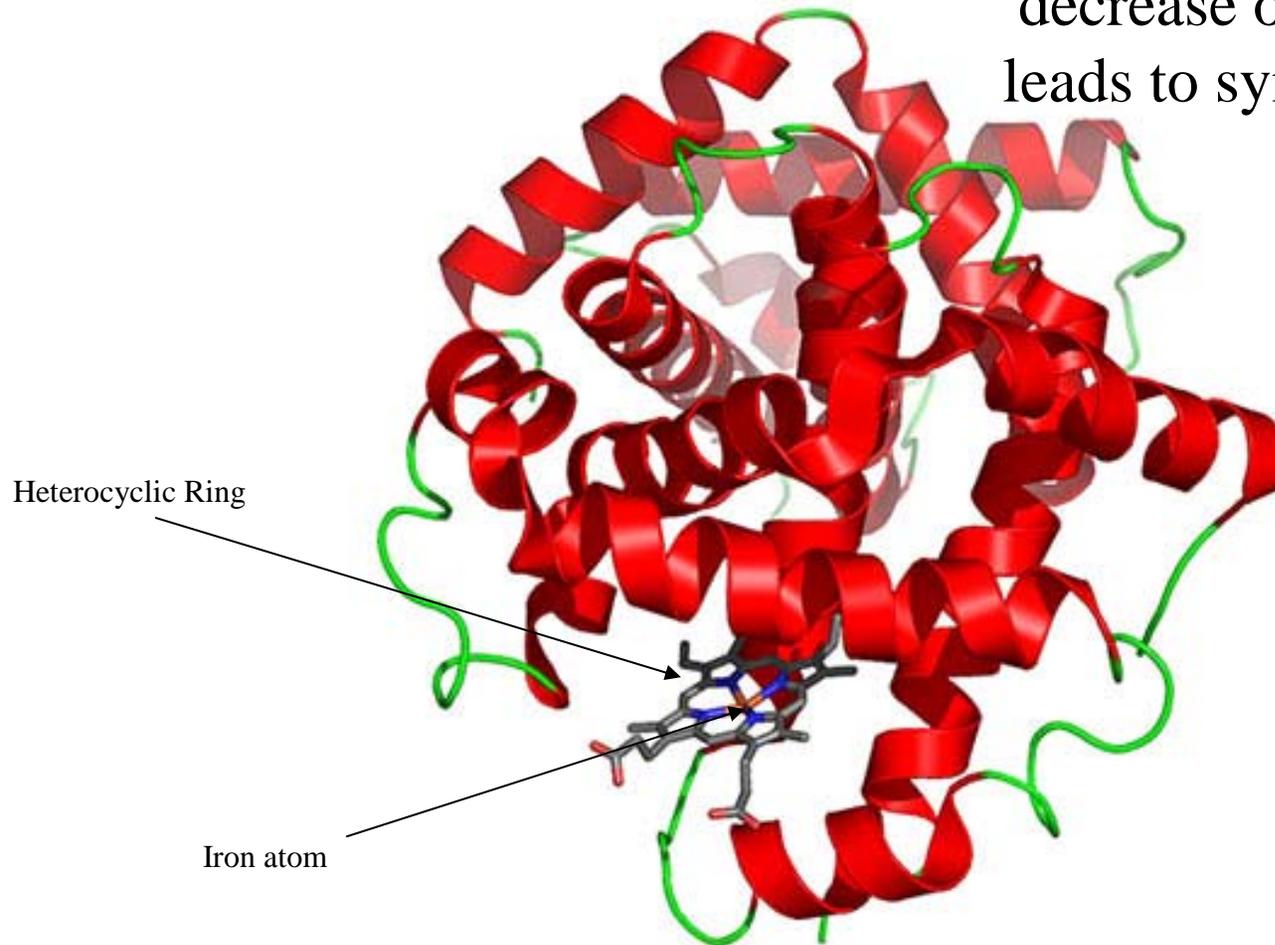
G Protein Complex

Protein Function

- Enzymatic catalysis
- Transport and storage
- Coordinated motion
- Immune protection
- Generation and transmission of nerve impulses
- Control of growth and differentiation
- Structural materials

Example: Hemoglobin Transports Oxygen

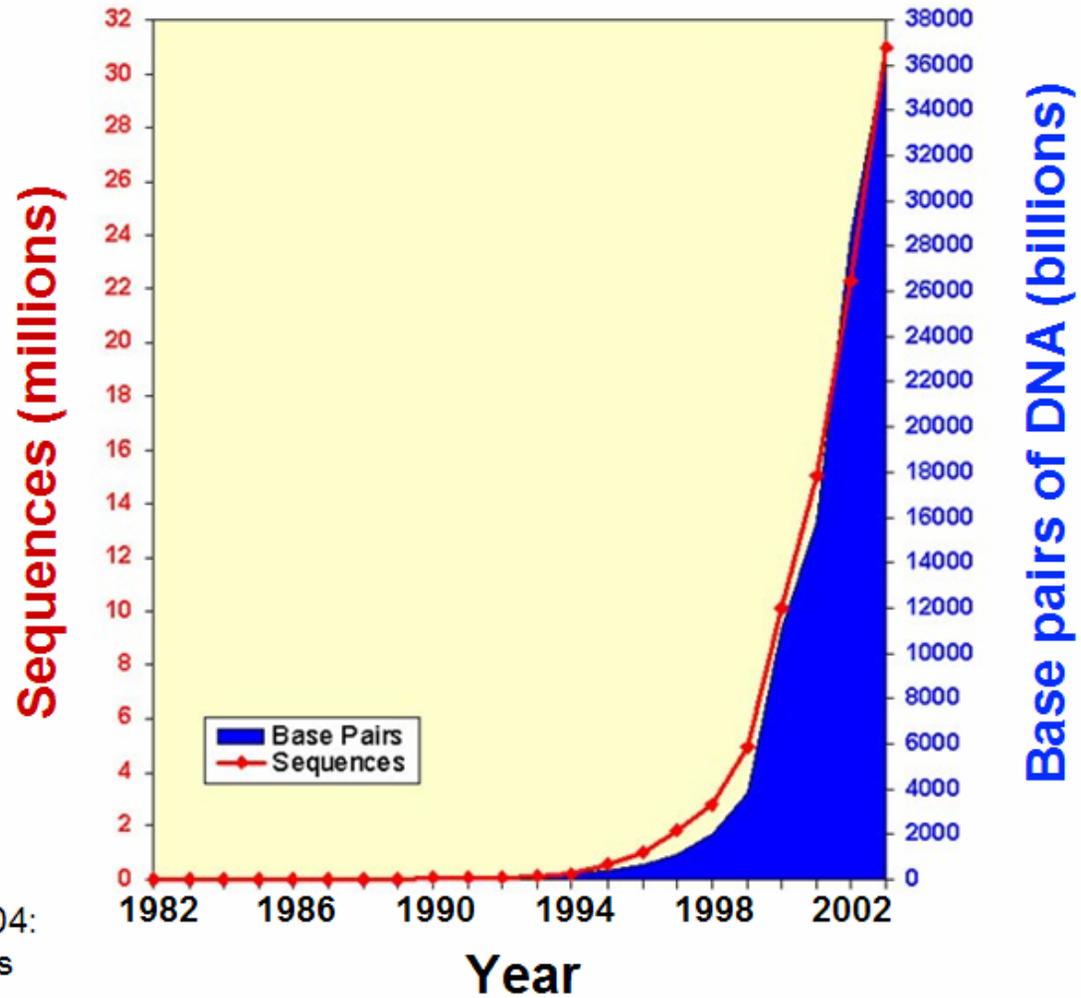
Decreased levels of hemoglobin, with or without an absolute decrease of red blood cells, leads to symptoms of anemia.



Huge Amount of Biological Data

- DNA/RNA Sequence Data (GeneBank)
(human genome project)
- Protein Sequence Data (SwissProt and PIR)
- Protein Structure Data (Protein Data Bank)
- Gene Expression Data (Gene expression omnibus)
- Protein Function (Gene Ontology)
- And many more.....

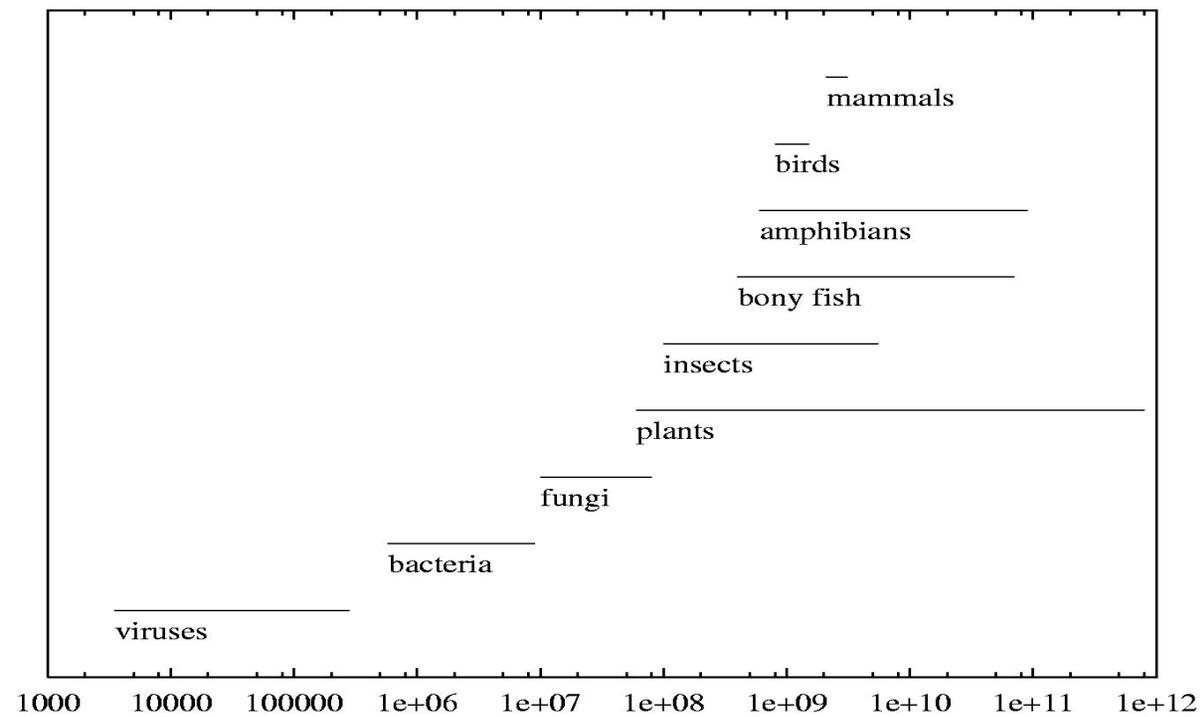
Growth of GenBank



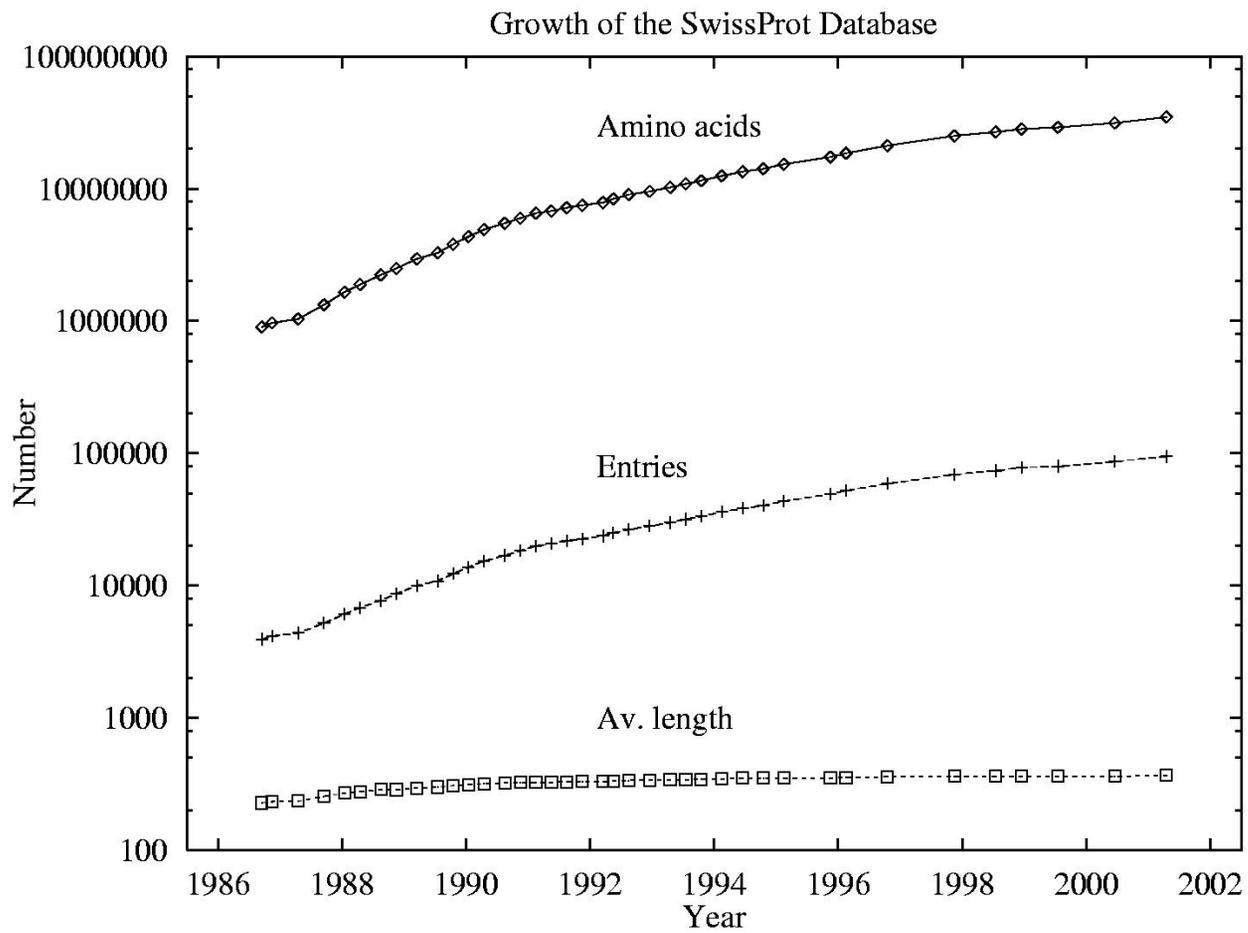
Updated 8-12-04:
>40b base pairs

Source: J. Pevsner, 2005

Genome Size



Source: P. Baldi, 2003



Source: P. Baldi, 2003

What can we do with these huge amount of data?



Find buried treasure - Doug Brutlag, 1999.

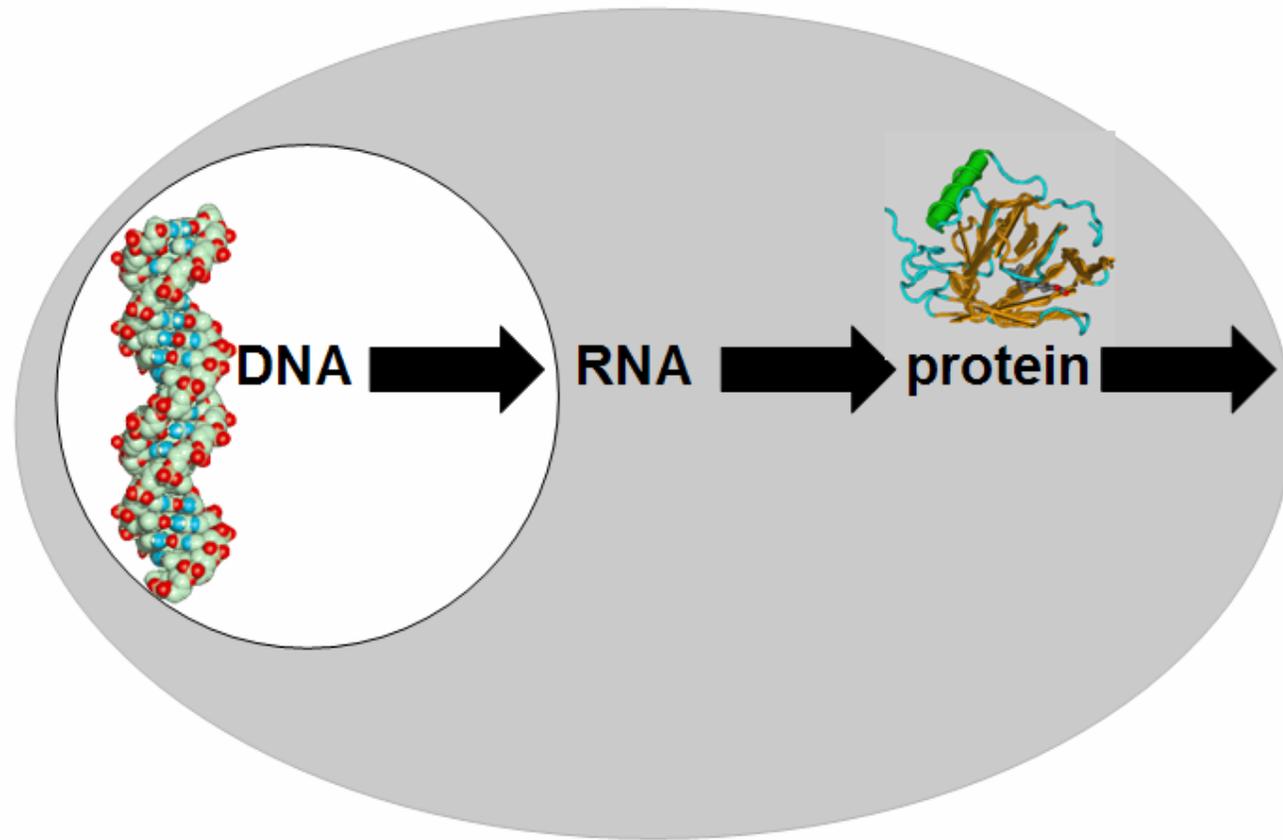
What can we do with these huge amount of data?

- Store (databases)
- Search (PSI-BLAST)
- Analyze / Annotate (visualization, interpretation, classification, pattern recognition)
- Generate new biological knowledge and build biological models (prediction)

Typical Problems

- Where is a gene? (about 1.5% human DNA codes gene)
- Are two protein sequences similar? Implication?
- Find protein sequences similar to my protein
- What does a protein's structure look like?
- Does this gene mutation relate to a disease (breast cancer)? How?

Central dogma of molecular biology



genome **→** transcriptome **→** proteome

Central dogma of bioinformatics and genomics

Genomics, transcriptomics, and proteomics

Source: J. Pevsner, 2005

Ten Topics

- 1. Introduction to Molecular Biology and Bioinformatics
- 2. Pairwise Sequence Alignment Using Dynamic Programming
- 3. Practical Sequence/Profile Alignment Using Fast Heuristic Methods (BLAST and PSI-BLAST)
- 4. Multiple Sequence Alignment
- 5. Gene and Motif Identification
- 6. Phylogenetic Analysis
- 7. Protein Structure Analysis and Prediction
- 8. RNA Secondary Structure Prediction
- 9. Clustering and Classification of Gene Expression Data
- 10. Search and Mining of Biological Databases, Databanks, and Literature