

# Computational Modeling of Molecular Structure

Jianlin Cheng, PhD

Department of EECS

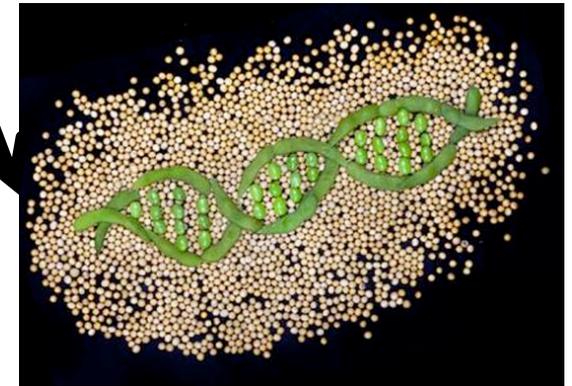
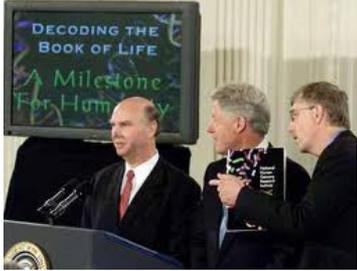
Informatics Institute

University of Missouri, Columbia

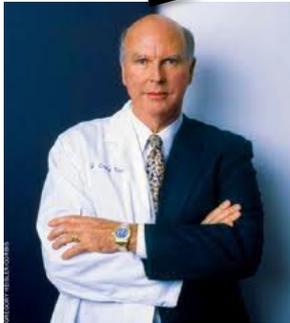
Spring, 2018

# The Genomic Era

Collins, Venter, Human Genome, 2000



# DNA Sequencing Revolution



Scientists



Government



Company

**\$1000  
Personal  
Genome**

# A Topic of Big Bio Data Analysis

## Science enters \$1,000 genome era

By Paul Rincon

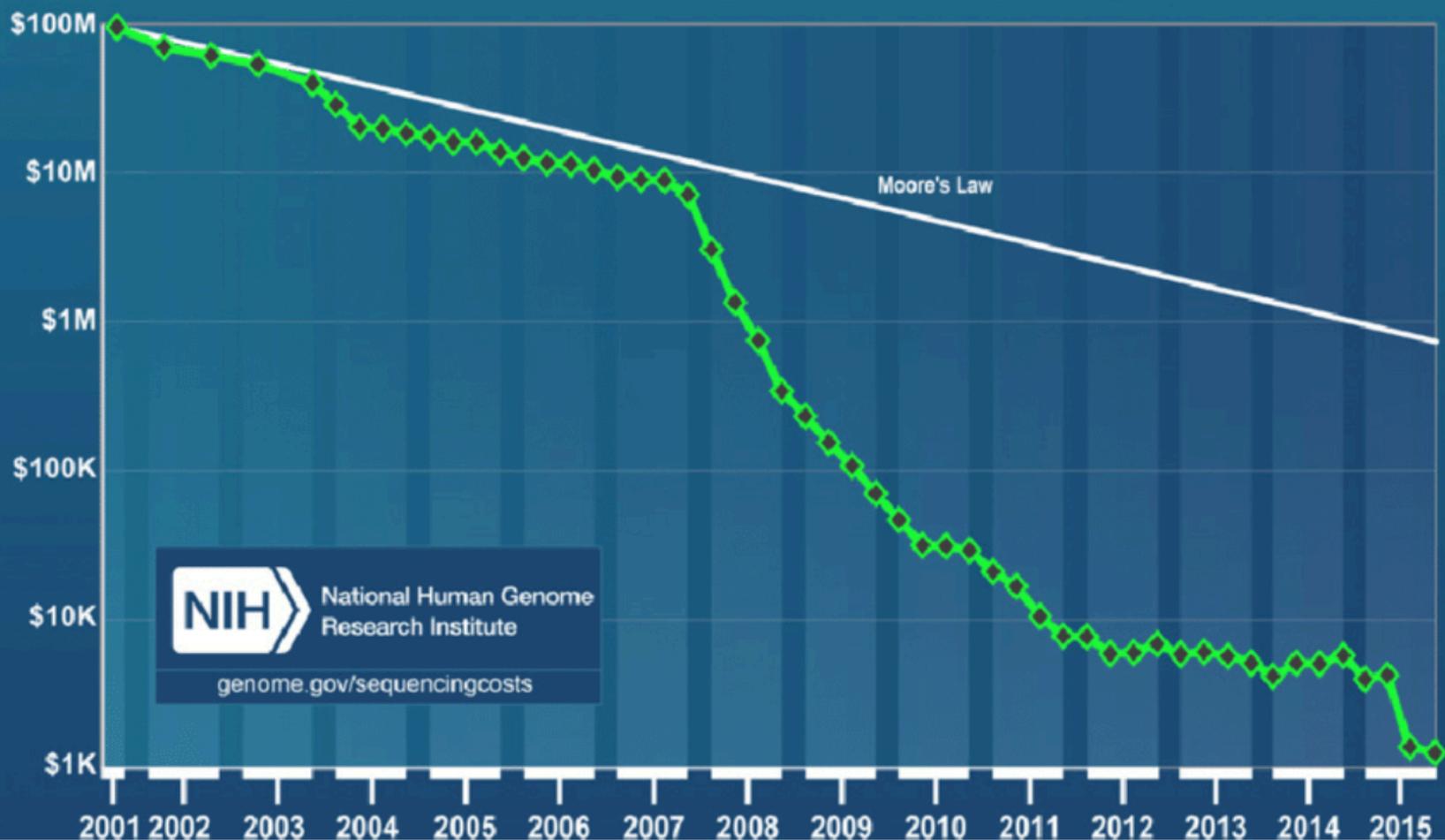
Science editor, BBC News website

---



The HiSeq X Ten is capable of sequencing five human genomes a day, Illumina claims

## Cost per Genome



# Objectives

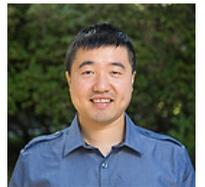
- Properties of molecular structures (proteins, RNA, genome / DNA)
- Computational representation of molecular structures
- Data-driven computational modeling of molecular structures
- Application of modeling of molecular structures such as drug design

# Significance of Studying Molecular Structures

- One foundation of life sciences
- Personal healthcare and medicine
- One major topic of bioinformatics and computational biology – an important field of computer science
- **A great application area of computer algorithms and data science**
- **A very interdisciplinary field (CS, data science, stats, math, biology, chemistry, physics)**

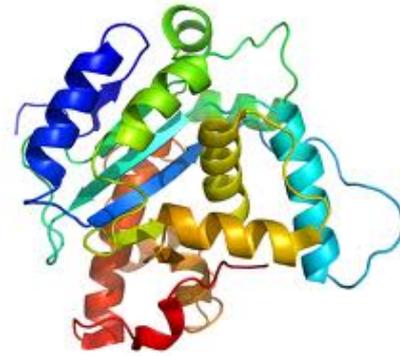
# A Good Career for CS Graduates

- Five PhD graduates are assistant / associate professors of bioinformatics in CS departments
- One PhD student secured a scientist position in a bioinformatics company
- One PhD student works for Microsoft
- Numerous other graduate students received good training and worked in data-intensive fields.

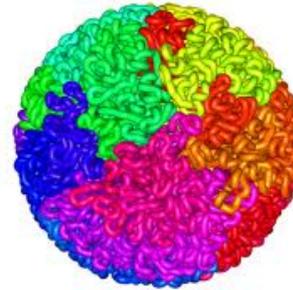


# Three Kinds of Structures

- **Protein Structure**



- **Genome Structure**

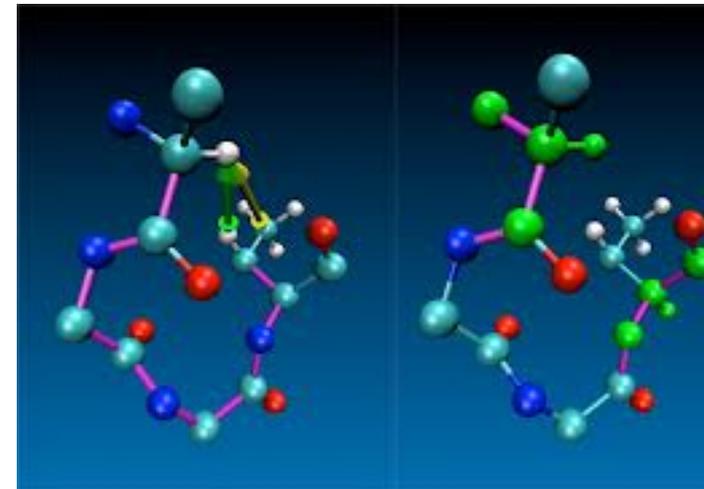
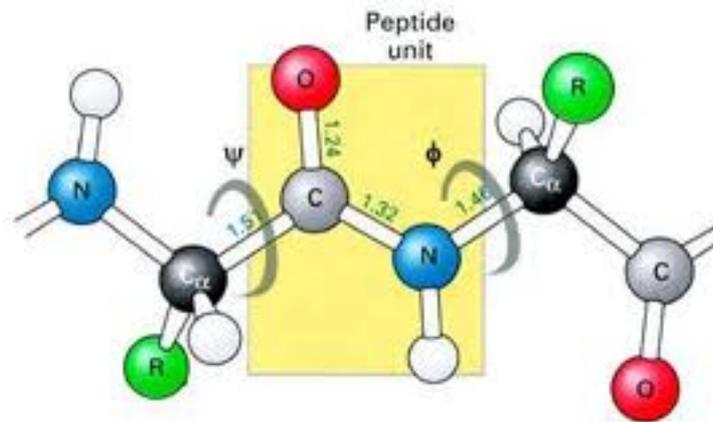
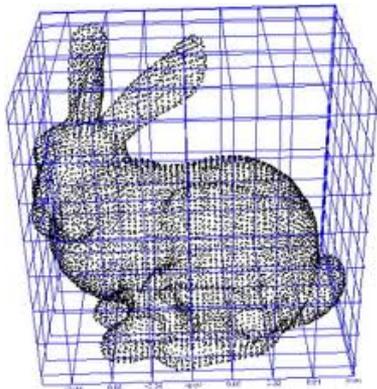
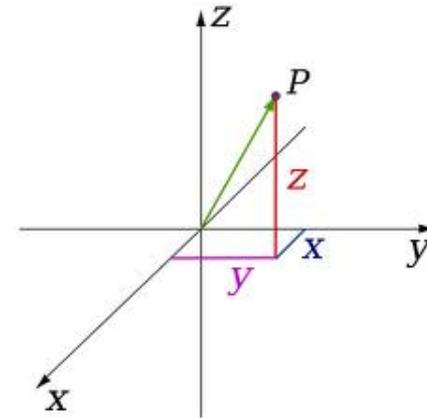


- **RNA Structure**



# Representation of Molecular Structures

- X, Y, Z coordinates
- Euclidean grid
- Vector and angles
- Computer graphics



# Algorithms

- Grid-based simulation (random walk)
- Vector-based simulation
- Angular-based simulation
- Gradient descent simulation and variants
- Simulated annealing
- Markov Chain Monte Carlo
- Probabilistic modeling
- Constraint-based optimization
- Machine learning (e.g., deep learning)

# Software Packages

- RasMol, Jmol, PyMol, Chimera
- Modeller, Rosetta, I-TASSER, MULTICOM, MTMG, CNS, CONFOLD, Zdock, MOGEN, LorDG, AutoDock
- Your own algorithm, implementation, and practice

# Course Format

- Course web site:  
[http://calla.rnet.missouri.edu/cheng\\_courses/cscmms2018/](http://calla.rnet.missouri.edu/cheng_courses/cscmms2018/)
- Problem solving
- Active learning by practicing
- Syllabus (see details)

# Teaching Format of Each Topic



**Group:**

**3-4 students per group**

**Rotate as topic coordinator**

**Each member participates  
in every topic**

**All members present  
the whole project**

# Grading

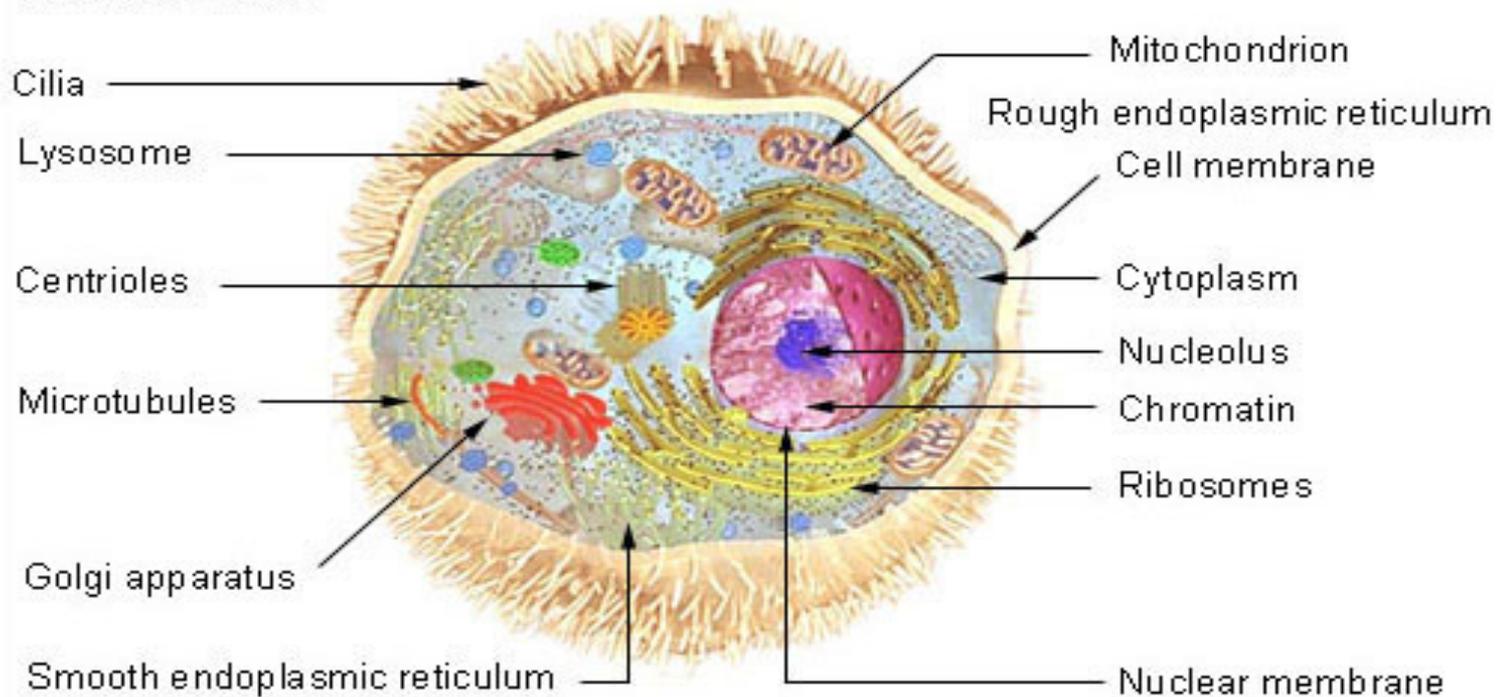
- Class discussions (15%)
- Literature reviews (10%)
- Topic plan presentation (20%, group)
- Topic implementation and report (45%, group)
- Final presentation (10%, group)
- Grade scale: A+, A, A-, B+, B, B-, C+, C, C-, and F.

# **Introduction to Molecular Biology for Computer Science and Engineering Students**

# Introduction to Molecular Biology

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

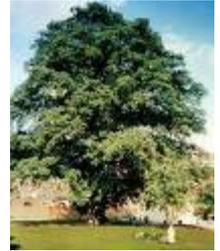
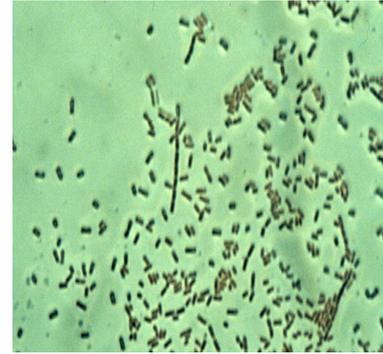
## Cell Structure



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

# Central Dogma of Molecular Biology

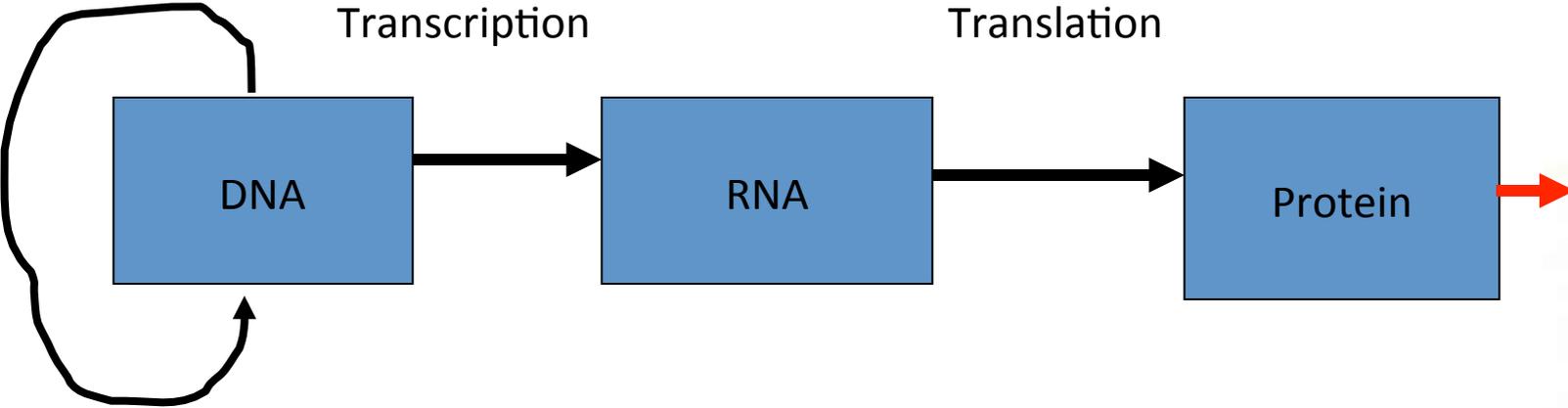
Phenotype



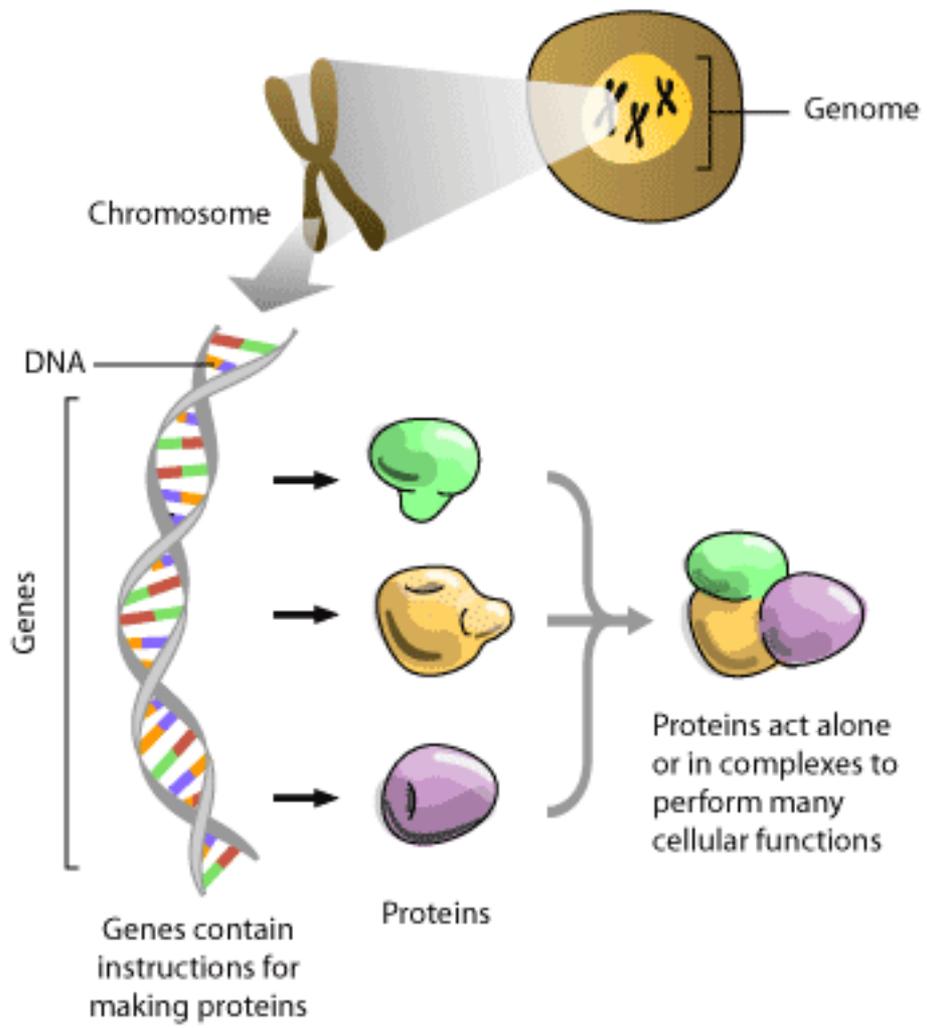
Replication

Transcription

Translation

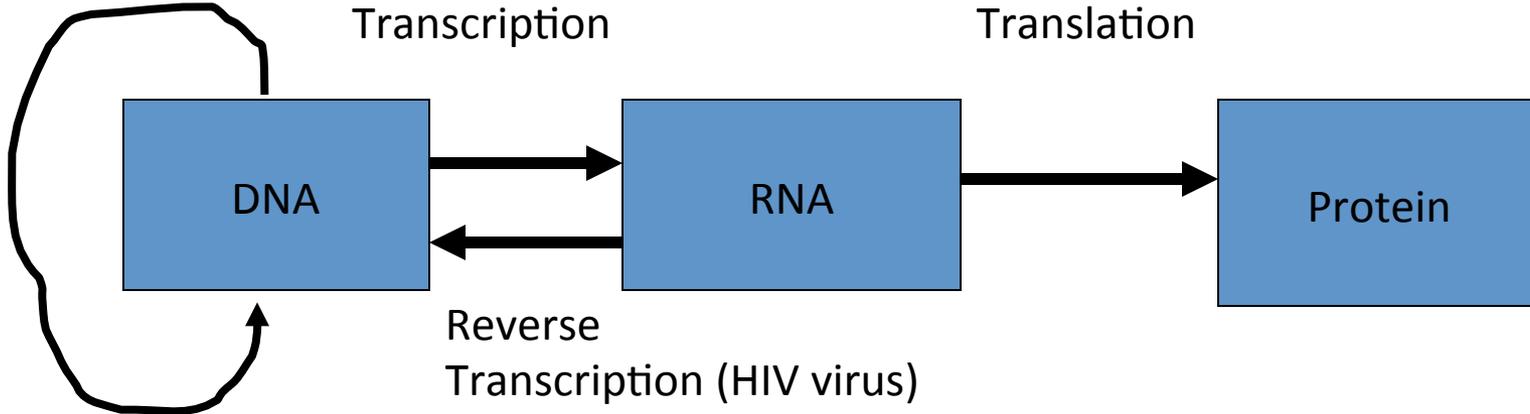


Genotype

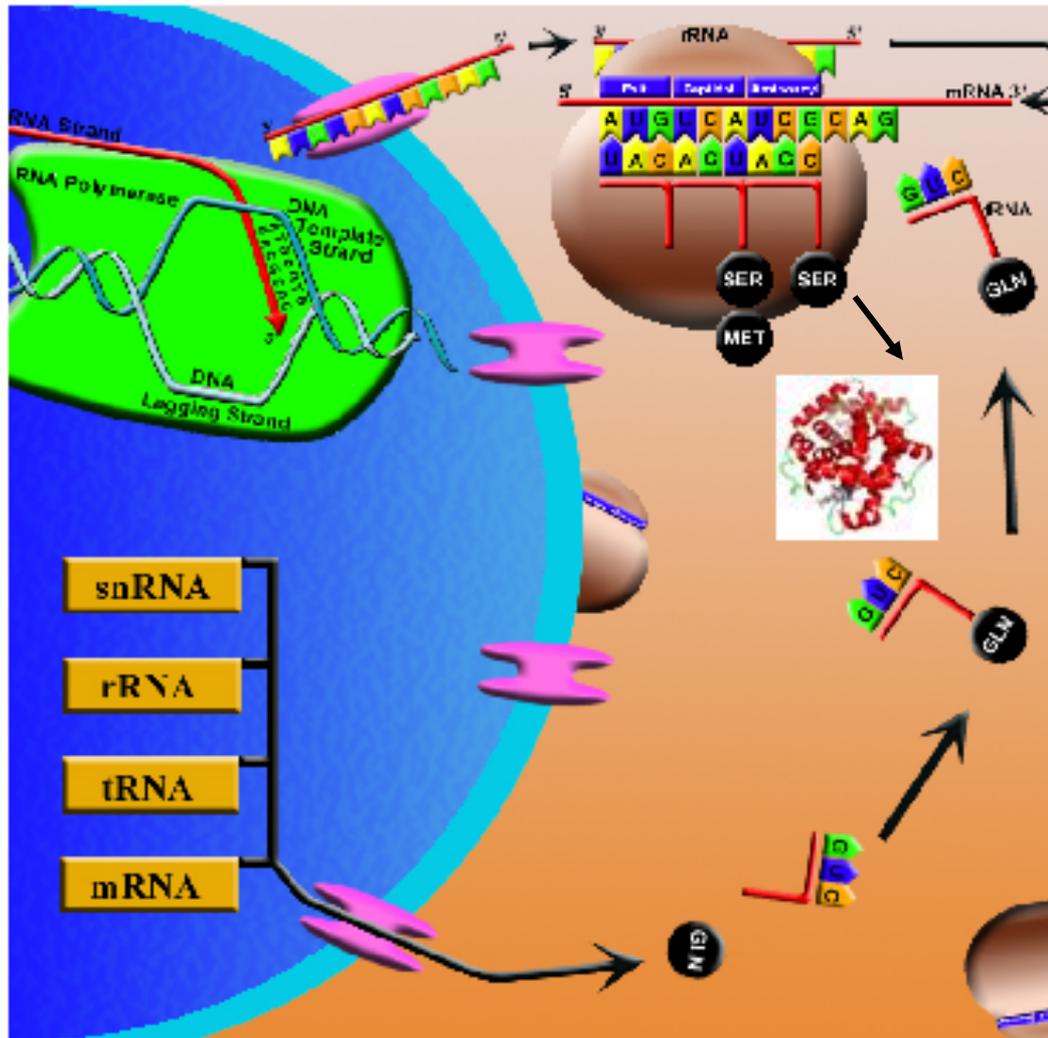


# Central Dogma of Molecular Biology

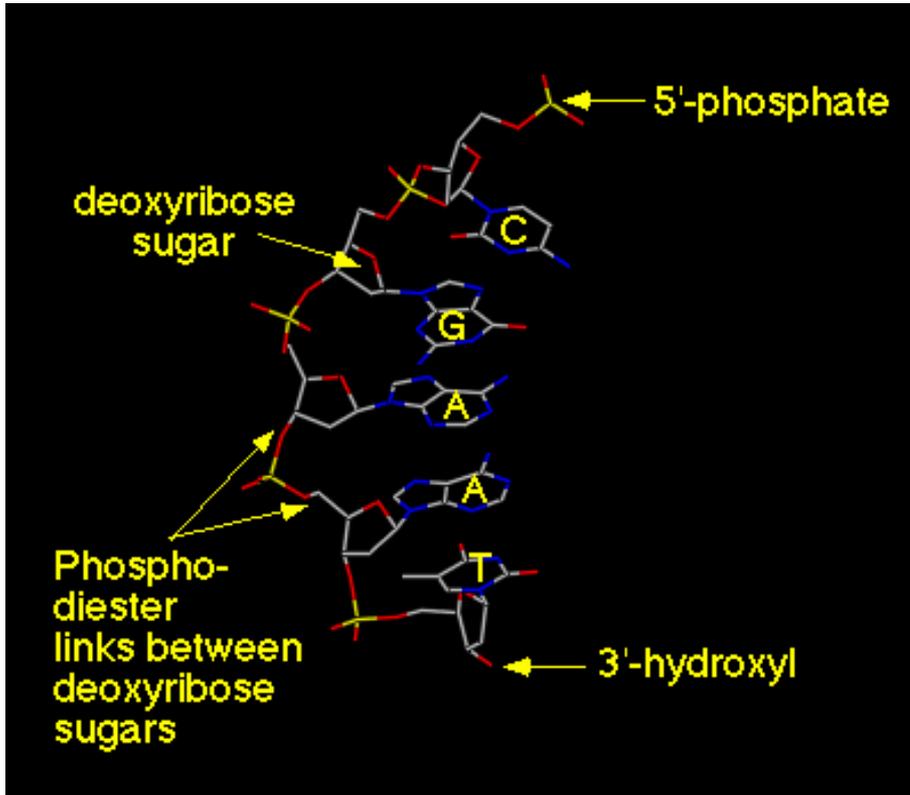
Replication



Information flow



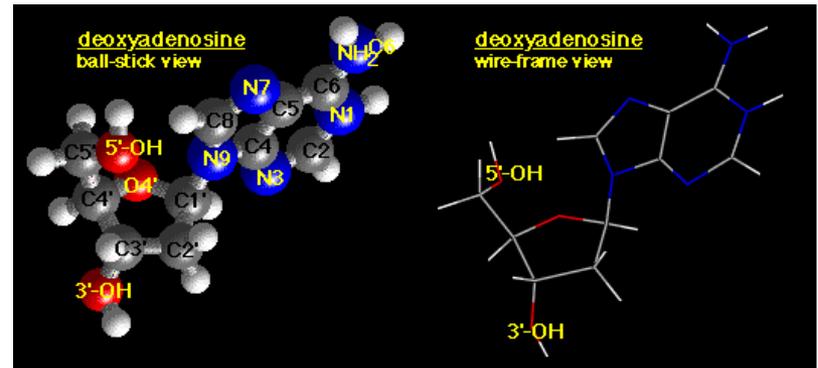
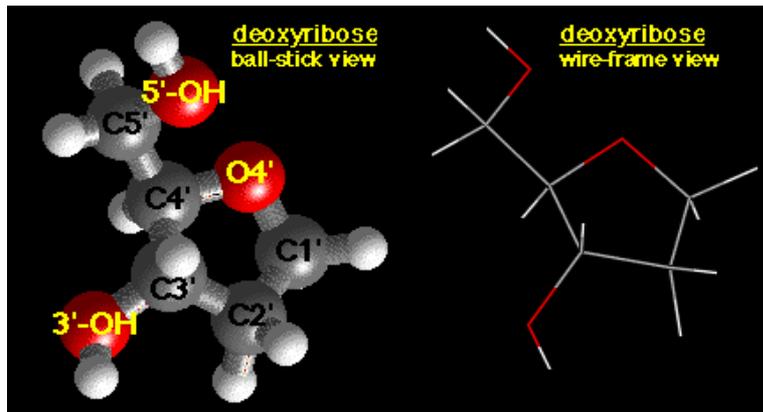
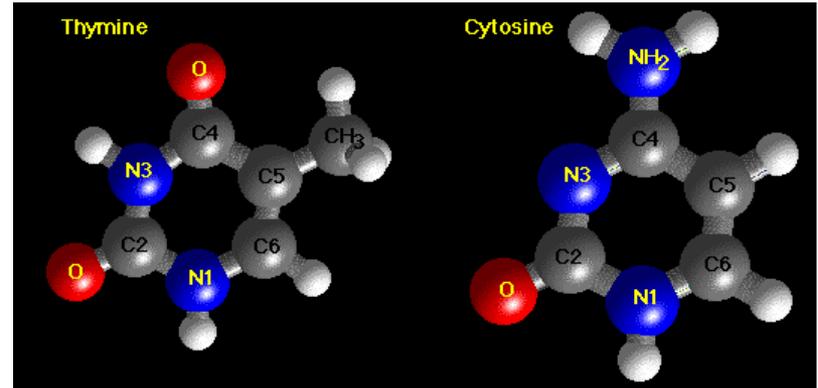
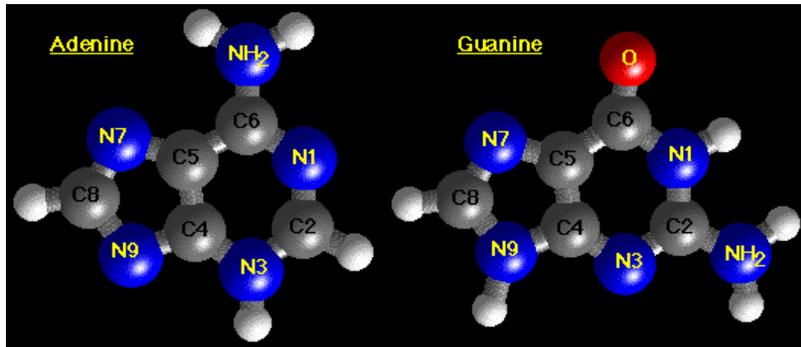
# DNA (Deoxyribose Nucleotide Acids)

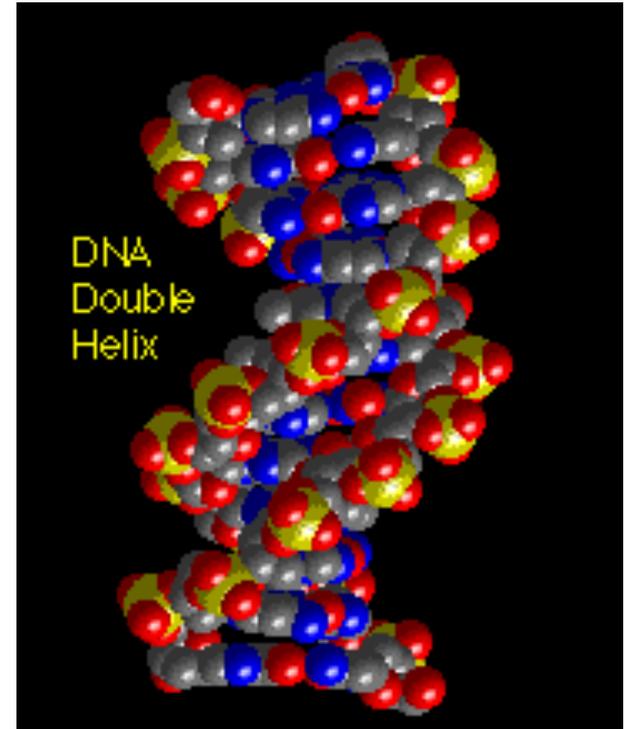
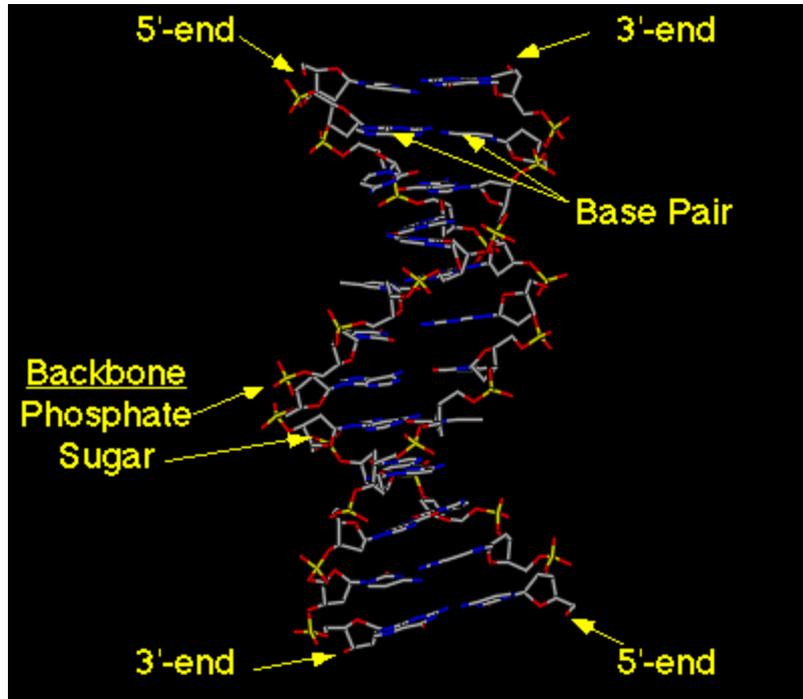


CGAATGGGAAA.....

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



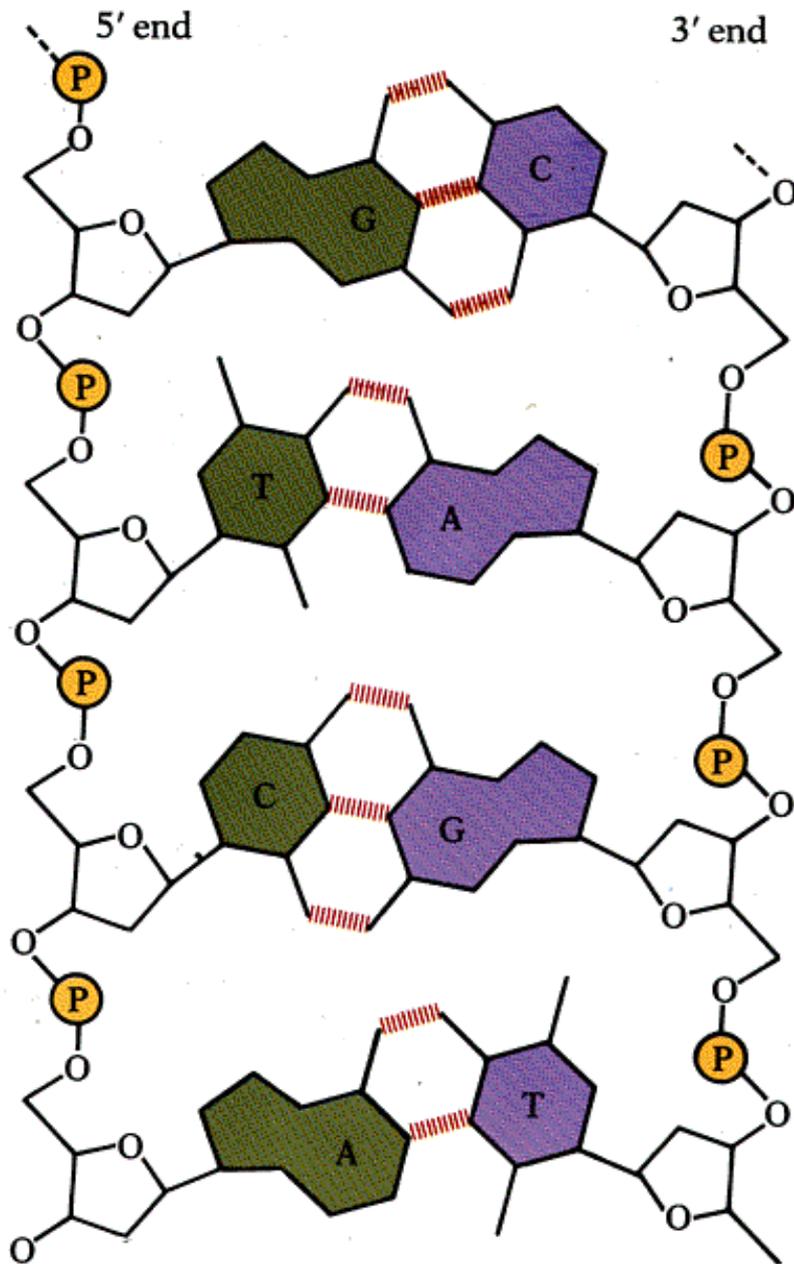


Base Pairs:

A-T (2 H-bonds)

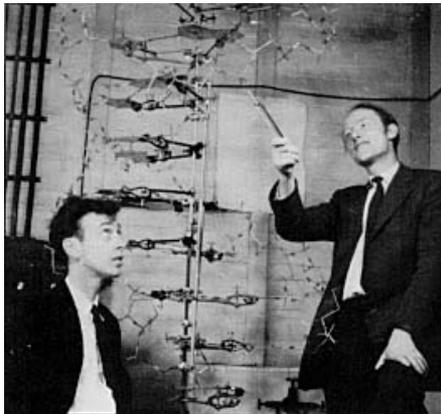
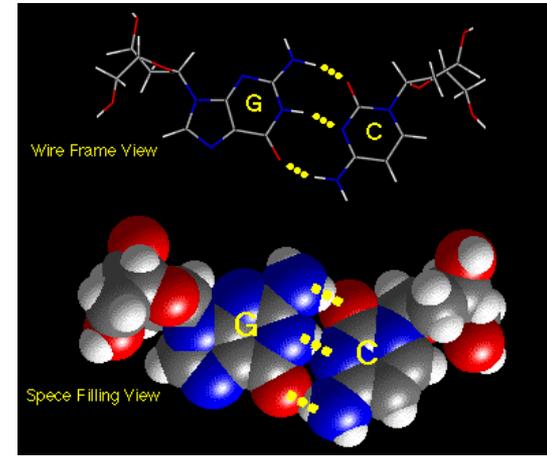
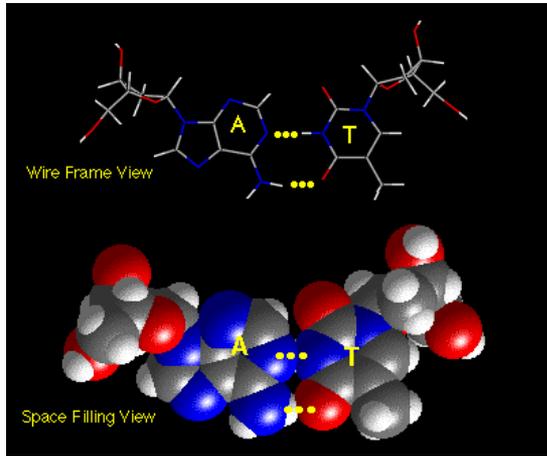
C-G (3 H-bonds)

Hydrogen bonds: non-covalent bonds mediated by hydrogen atoms



Uncoiled DNA Molecule

Source: Dr. Gary Stormo, 2002



James Watson & Francis Crick



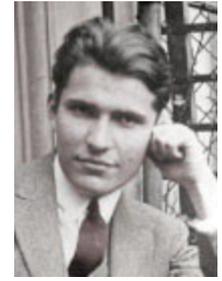
Maurice Wilkins



Rosalind Franklin



Linus Pauling

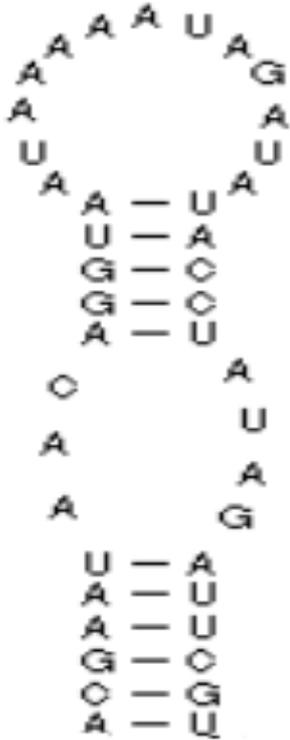
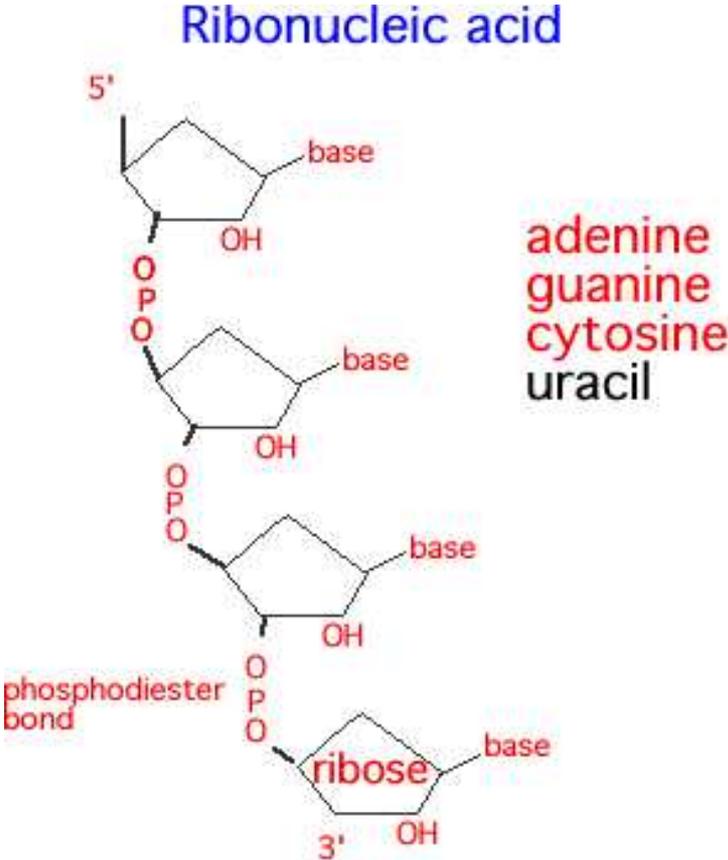


Erwin Chargaff

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



# RNA (Ribose Nucleotide Acids)

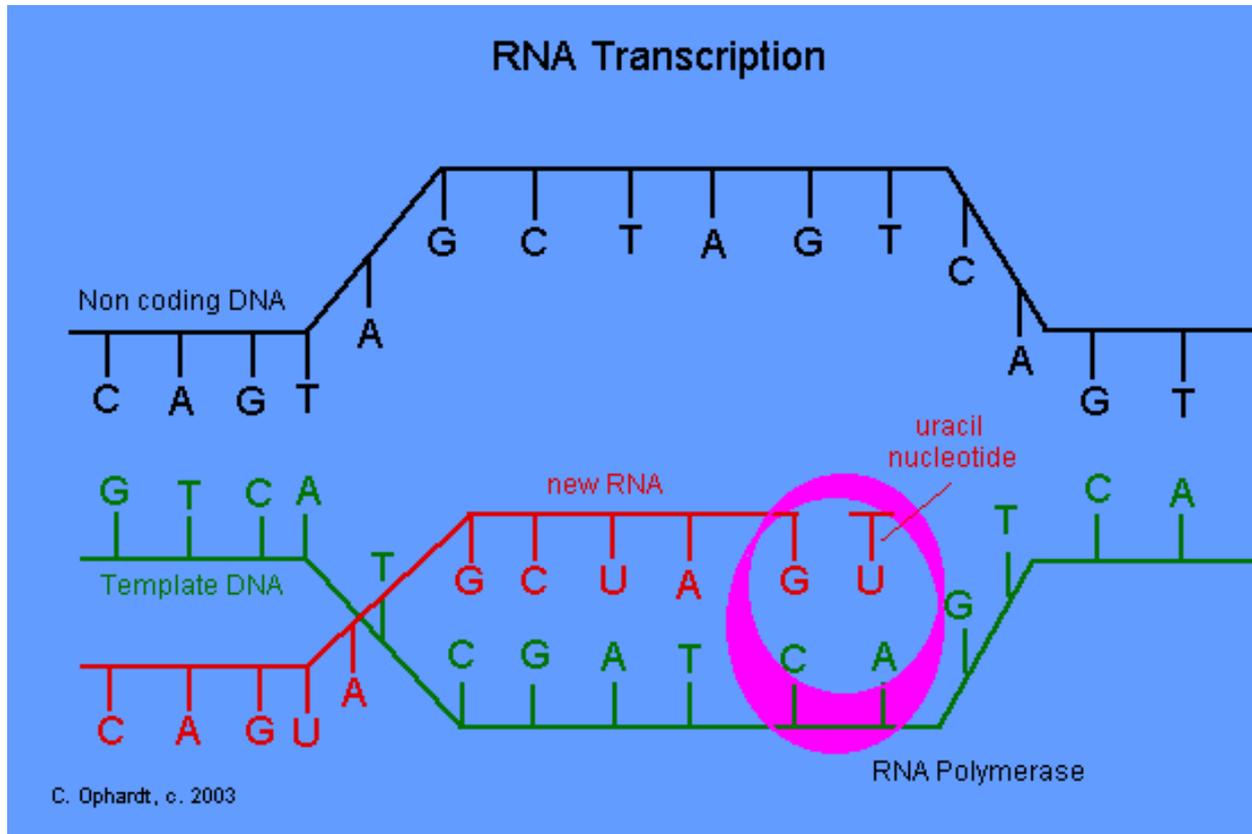


ACGAAUAACAGGUAUUAAAAUAGAUUACCUAUAGAUUCGU

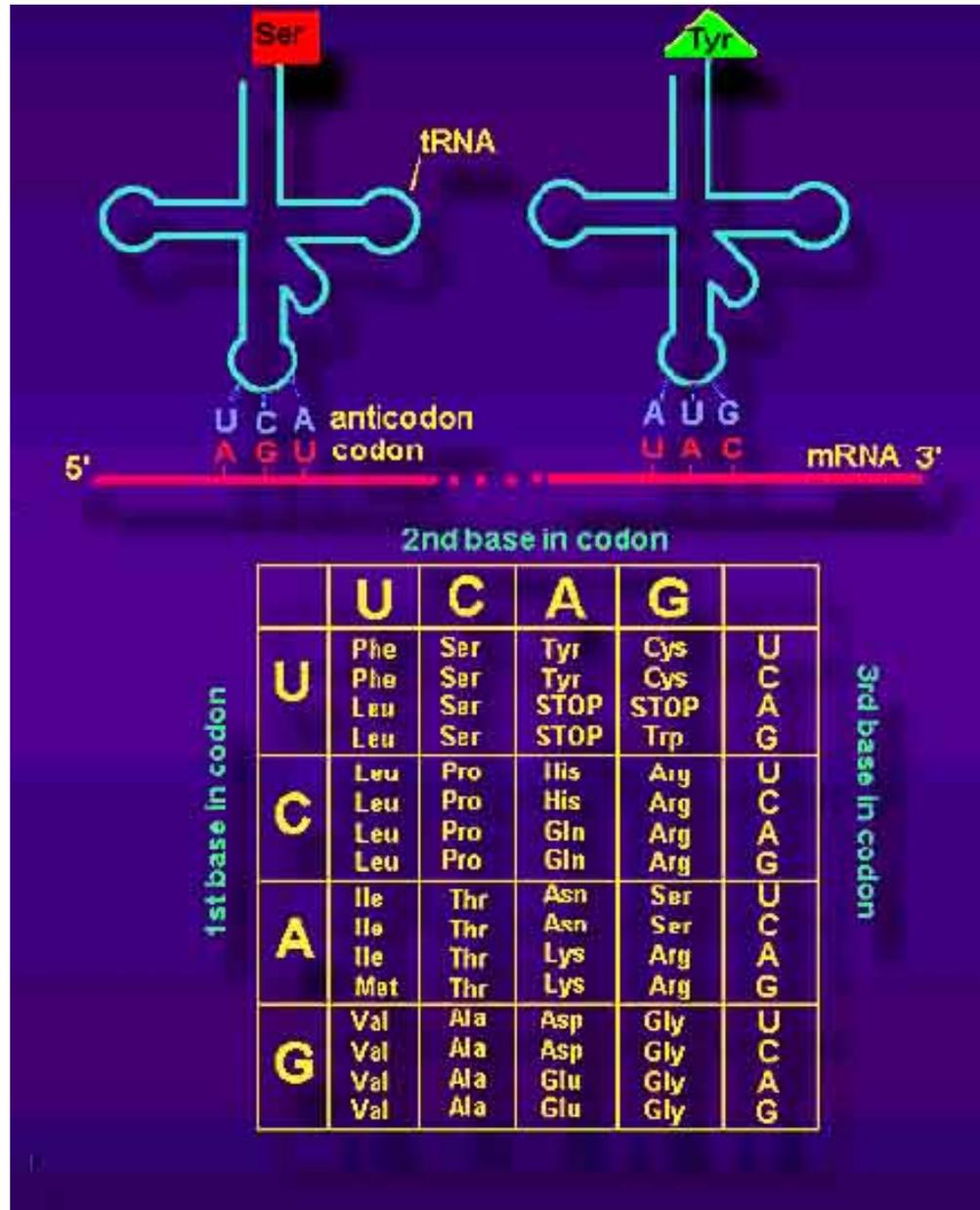
# Different Kinds of RNA

- 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
- Non-coding RNA

# Transcription of Gene into RNA



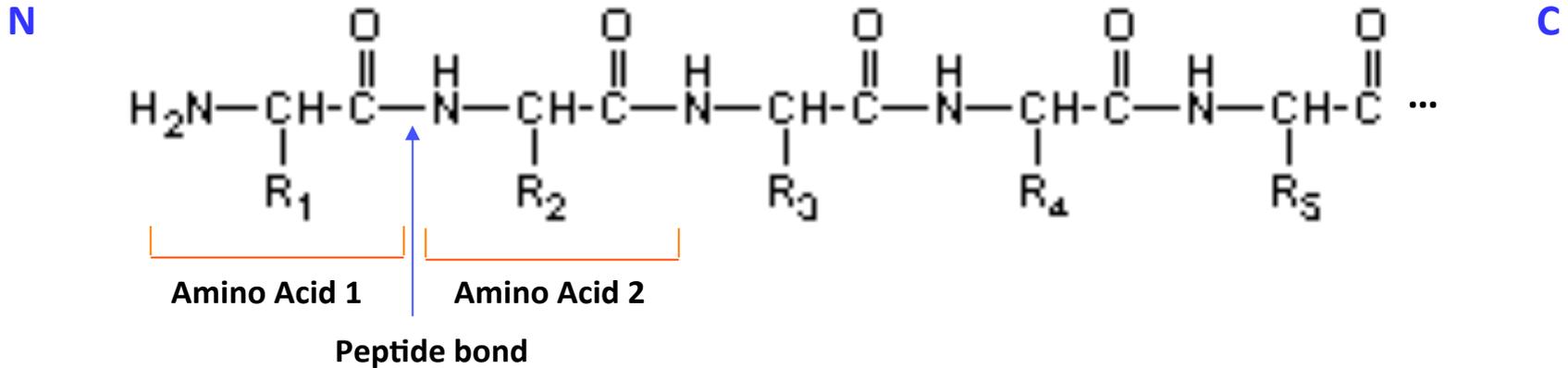
# Genetic Code and Translation



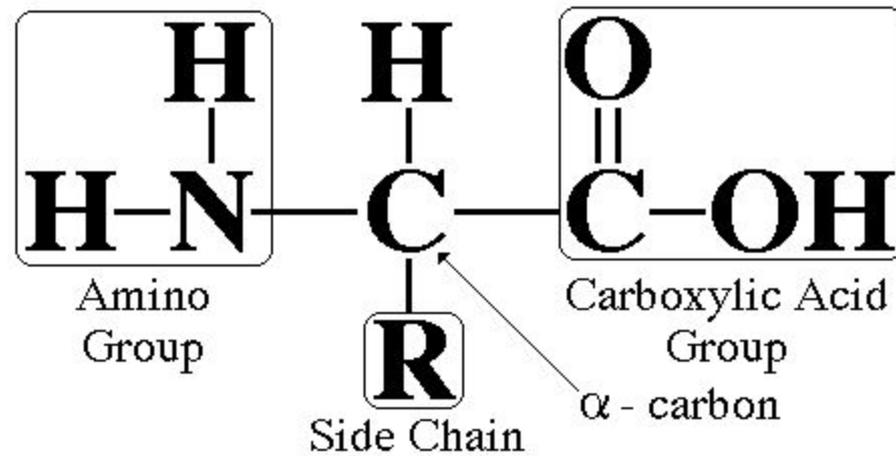
Three Nucleotides is called a codon.

# Protein Sequence

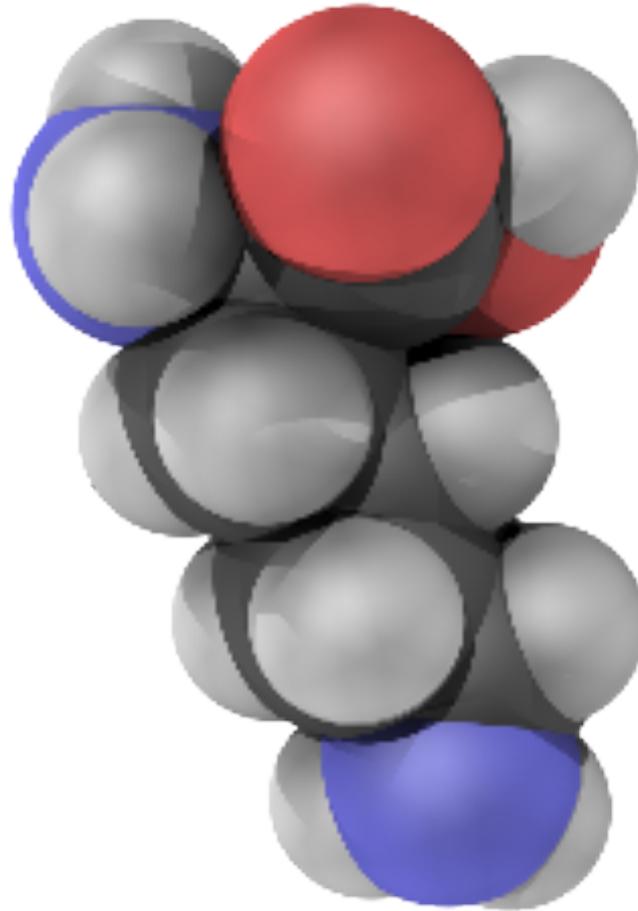
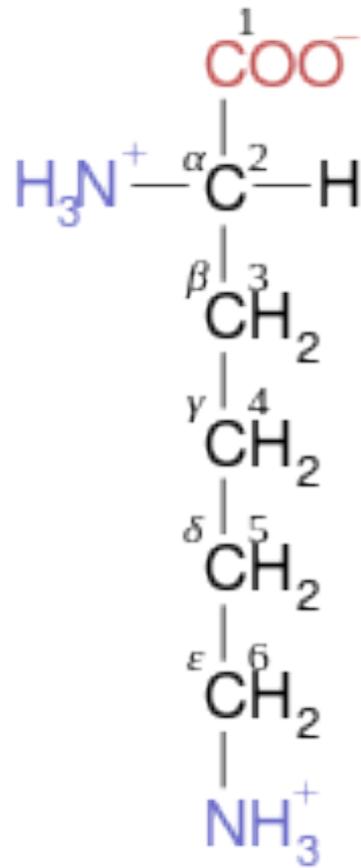
A directional sequence of amino acids/residues



# Amino Acid Structure



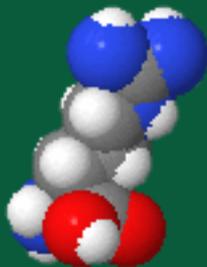
# Lysine



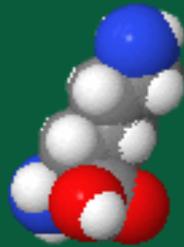
# Amino Acids

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 <sub>3</sub>	X	-	-	X	X	-	67	GCU, GCC, GCA, GCG	7.8
Cysteine	Cys, C	-CH <sub>2</sub> SH	X	-	-	X	-	-	86	UGU, UGC	1.9
Aspartate	Asp, D	-CH <sub>2</sub> COOH	-	X	negative	X	-	-	91	GAU, GAC	5.3
Glutamate	Glu, E	-CH <sub>2</sub> CH <sub>2</sub> COOH	-	X	negative	-	-	-	109	GAA, GAG	6.3
Phenylalanine	Phe, F	-CH <sub>2</sub> C <sub>6</sub> H <sub>5</sub>	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 <sub>2</sub> -C <sub>3</sub> H <sub>3</sub> N <sub>2</sub>	-	X	positive	-	-	Aromatic	118	CAU, CAC	2.3
Isoleucine	Ile, I	-CH(CH <sub>3</sub> )CH <sub>2</sub> CH <sub>3</sub>	X	-	-	-	-	Aliphatic	124	AUU, AUC, AUA	5.3
Lysine	Lys, K	-(CH <sub>2</sub> ) <sub>4</sub> NH <sub>2</sub>	-	X	positive	-	-	-	135	AAA, AAG	5.9
Leucine	Leu, L	-CH <sub>2</sub> CH(CH <sub>3</sub> ) <sub>2</sub>	X	-	-	-	-	Aliphatic	124	UUA, UUG, CUU, CUC, CUA, CUG	9.1
Methionine	Met, M	-CH <sub>2</sub> CH <sub>2</sub> SCH <sub>3</sub>	X	-	-	-	-	-	124	AUG	2.3
Asparagine	Asn, N	-CH <sub>2</sub> CONH <sub>2</sub>	-	X	-	X	-	-	96	AAU, AAC	4.3
Proline	Pro, P	-CH <sub>2</sub> CH <sub>2</sub> CH <sub>2</sub> -	X	-	-	X	-	-	90	CCU, CCC, CCA, CCG	5.2
Glutamine	Gln, Q	-CH <sub>2</sub> CH <sub>2</sub> CONH <sub>2</sub>	-	X	-	-	-	-	114	CAA, CAG	4.2
Arginine	Arg, R	-(CH <sub>2</sub> ) <sub>3</sub> NH-C(NH) NH <sub>2</sub>	-	X	positive	-	-	-	148	CGU, CGC, CGA, CGG, AGA, AGG	5.1
Serine	Ser, S	-CH <sub>2</sub> OH	-	X	-	X	X	-	73	UCU, UCC, UCA, UCG, AGU, AGC	6.8
Threonine	Thr, T	-CH(OH)CH <sub>3</sub>	X	X	-	X	-	-	93	ACU, ACC, ACA, ACG	5.9
Valine	Val, V	-CH(CH <sub>3</sub> ) <sub>2</sub>	X	-	-	X	-	Aliphatic	105	GUU, GUC, GUA, GUG	6.6
Tryptophan	Trp, W	-CH <sub>2</sub> C <sub>8</sub> H <sub>6</sub> N	X	-	-	-	-	Aromatic	163	UGG	1.4
Tyrosine	Tyr, Y	-CH <sub>2</sub> -C <sub>6</sub> H <sub>4</sub> OH	X	X	-	-	-	Aromatic	141	UAU, UAC	3.2

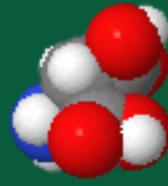
↑  
Hydrophilic



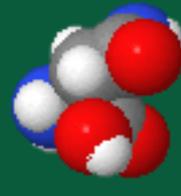
Arg



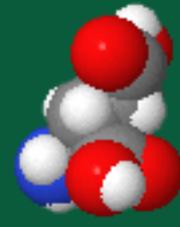
Lys



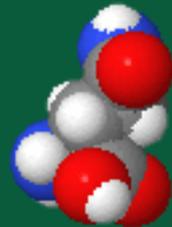
Asp



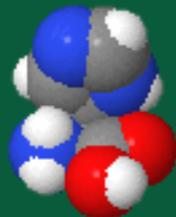
Asn



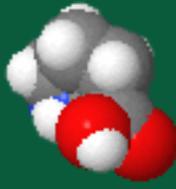
Glu



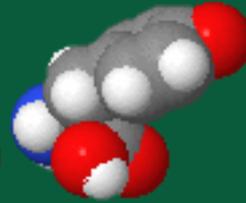
Gln



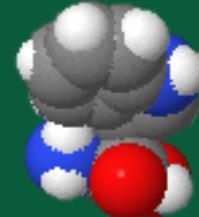
His



Pro



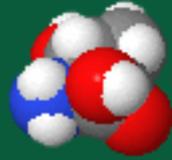
Tyr



Trp



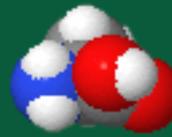
Ser



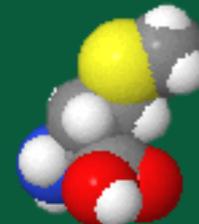
Thr



Gly



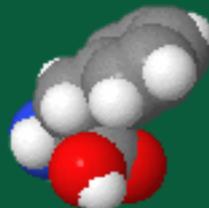
Ala



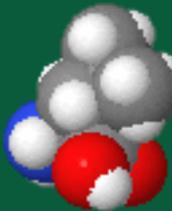
Met



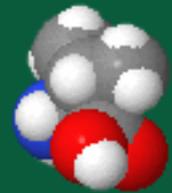
Cys



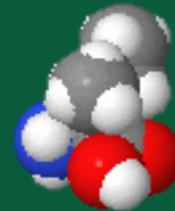
Phe



Leu



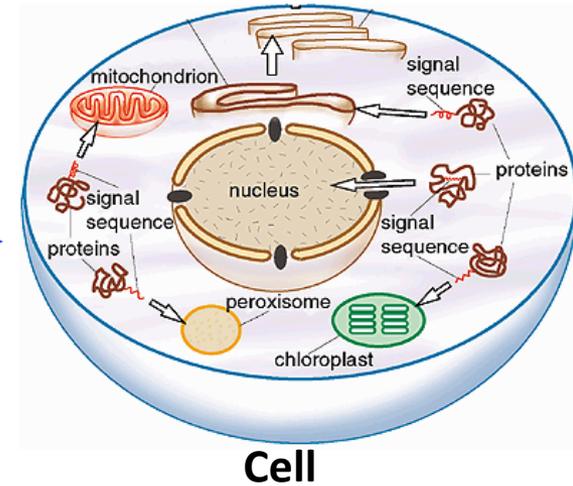
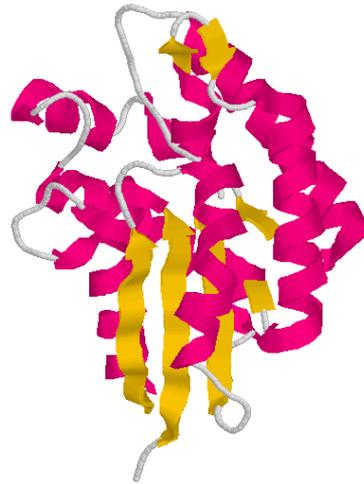
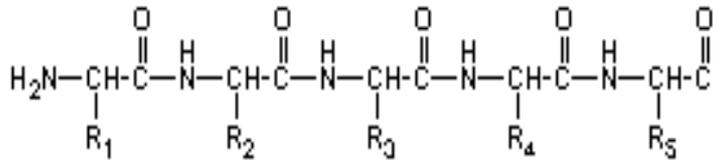
Val



Ile

# Central Dogma of Proteomics

AGCWY.....

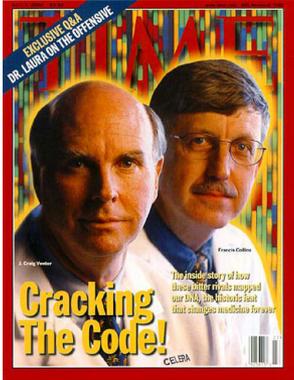


**Sequence**

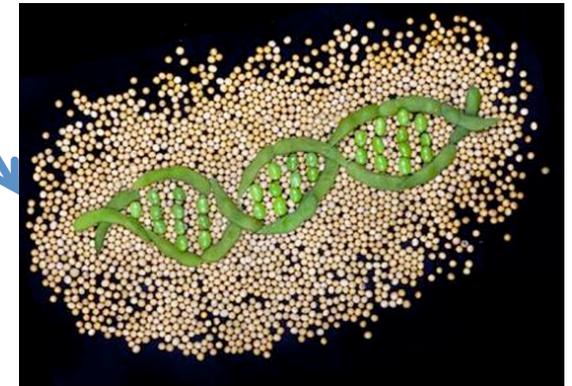
**Structure**

**Function**

# The Genomic Era



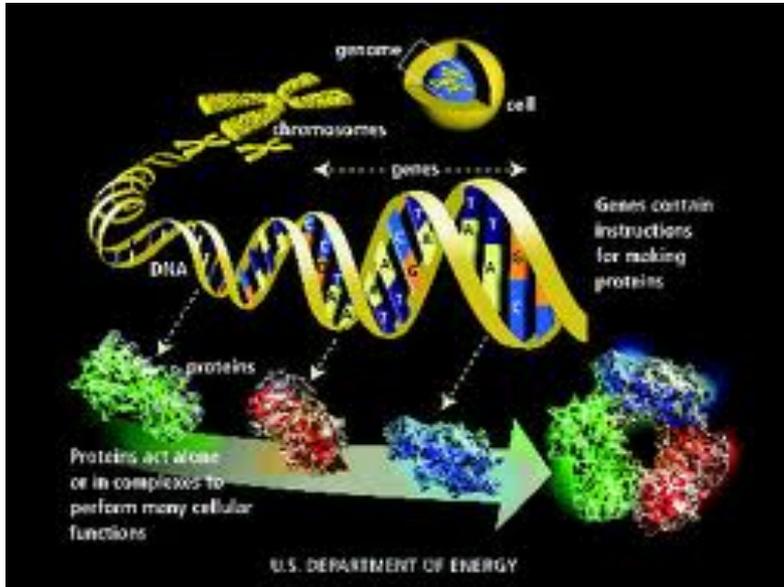
Collins, Venter, Human Genome, 2000



# Personal Genome's Implications

- **Personalized Disease Prevention**
- **Personalized Disease Diagnosis**
- **Personalized Medicine**
- **Personalized Health Care**
- **Precision Medicine**

# Genome Implications to Information Sciences and Life Sciences



## Elements and Systems

# Assignment One

Read one of the two articles and write a half page summary:

A. Sali. T. Blundell. Comparative Protein Modeling by Satisfaction of Spatial Restraints. JMB, 1993.

J. Li, J. Cheng. A Stochastic Point Cloud Sampling Method for Multi-Template Protein Comparative Modeling. Scientific Reports, 2016.

Submit your review summary (half page)

to [mumachinelearning@gmail.com](mailto:mumachinelearning@gmail.com) .

Due by Feb. 3 (Saturday).

Form your group (3 to 4 students per group)

# Acknowledgements

**images.google.com and all the authors  
providing valuable images**