

## Quantitative Biology

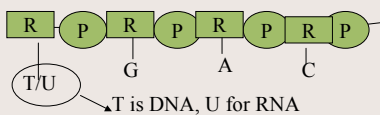
- Building mathematical models to understand and predict behavior of living systems
  - A wide range of hierarchical models of life:
    - Self contained replicating molecules (?)
    - Cells
    - Co-operating cells
    - Populations
    - Artificial life
- Extensive use of abstraction (good computer programming)

## CS: 426 Introduction to computational (molecular) biology

- A bottom to top approach
  - What is a “minimal” model of life?  
A self-replicating molecule that uses material and energy from the environment
  - The most basic biological molecules are DNA, RNA, proteins, and lipids. Which is a candidate for a minimal model?  
For a long time people try to prove that proteins are the best choice, but more recently RNA became the more natural choice.

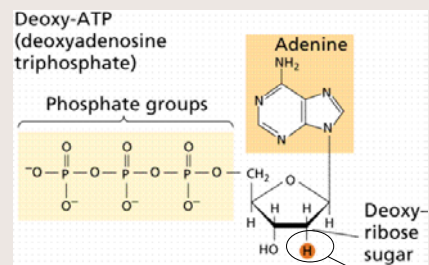
## Why RNA?

- RNA performs two functions coding information and synthesis
  - Like essentially all the important biological molecules RNA is a linear polymer that is made of a small number of monomers
  - The different monomers are attached to the same backbone (sugar rings (ribose) linked by phosphates)



## RNA/DNA (continue)

- Closer look at basic syntax

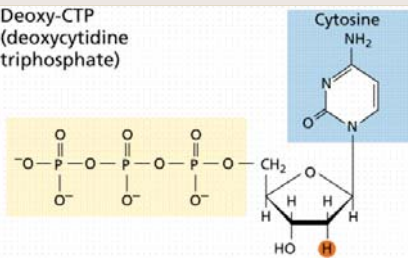


Called by “friends and family” : A

For RNA  
Replace H by OH

## More RNA/DNA

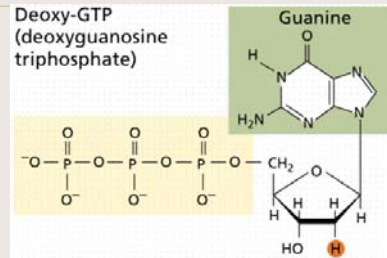
Deoxy-CTP  
(deoxycytidine  
triphosphate)



Called C

## RNA/DNA: basic

Deoxy-GTP  
(deoxyguanosine  
triphosphate)



Called G So far we have A,C,G for DNA we also have T for RNA U. Total of 4 "bases".

## The bases



Adenine



Guanine



Thymine  
*DNA only*

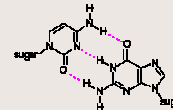


Uracil  
*RNA only*

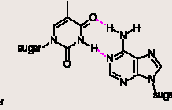


Cytosine

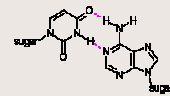
## The base pairs (hydrogen bonds)



CG

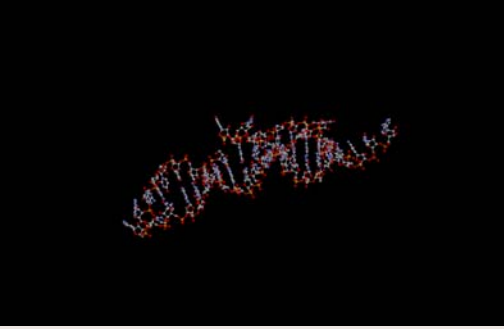


AT



AU

## RNA



Hydrogen bonding between complementing bases:  
G:C or A:T/U (DNA/RNA)

## RNA (the n-th time)

- The base pairing C:G and A:T determine structure of the RNA (must match!). One dimensional sequence determines shape e.g.
- GGGAGCUCAACUCUCCCCCCC  
UUUUCGAGGGUCAUCGGAAC  
CA
- Determine the shape from sequence using maximum pairing condition (RNA structure prediction problem).

## RNA (n+1) time



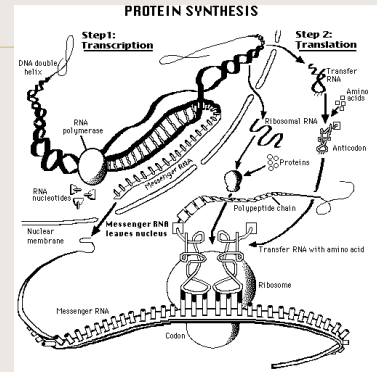
## RNA (n+2)

- RNA has globular shape
- Double molecular function:
  - Linear sequence stores information (information in biology is stored linearly)
  - Three dimensional shape determines enzymatic (machine-like) activity
- **Can be dangerous for the data to process itself!**

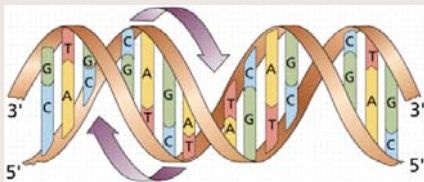
## Life cannot be made too simple

- As a cautionary step nature separated information from processing and cell function. Genetic information (in particular) must be protected exceptionally well.
- Instead of one molecule that does it all we now have two: DNA (information) Proteins (operations). RNA is kept in the middle probably for historical reasons and additional safety nets.

## From storage to machines



DNA is a “boring” molecule that “only” stores linear information



In contrast to RNA no complex three dimensional fold or enzymatic function (only hard disk no CPU). A **double** helix (A/B/Z) of complementing base pairs, very useful in information processing and corrections

DNA in the most common B form



Genetic Table From DNA To proteins

From 4 bases To 20 amino acids

Many errors Allowed.

Not all DNA codes proteins

	T	C	A	G
<b>T</b>	TTT Phe (F) TTC " TTA Leu (L) TTG "	TCT Ser (S) TCC " TCA " TCG "	TAT Tyr (Y) TAC TAA Ter TAG Ter	TGT Cys (C) TGC TGA Ter TGG Trp (W)
<b>C</b>	CTT Leu (L) CTC " CTA " CTG "	CCT Pro (P) CCC " CCA " CCG "	CAT His (H) CAC " CAA Gln (Q) CAG "	CGT Arg (R) CGC " CGA " CGG "
<b>A</b>	ATT Ile (I) ATC " ATA " ATG Met (M)	ACT Thr (T) ACC " ACA " ACG "	AAT Asn (N) AAC " AAA Lys (K) AAG "	AGT Ser (S) AGC " AGA Arg (R) AGG "
<b>G</b>	GTT Val (V) GTC " GTA " GTG "	GCT Ala (A) GCC " GCA " GCG "	GAT Asp (D) GAC " GAA Glu (E) GAG "	GGT Gly (G) GGC " GGA " GGG "

## More on DNA

- In higher organisms a lot of "junk" non coding DNA (up to 90 percent)
- Junk DNA may play a role in regulation controlling gene transcription/introns + exons
- Reading frames: Can shift a base to read different amino acids, e.g. (codon – 3 basepairs)
  - (TAA)(TCG)(AAT)(GGG)C == XSNG
  - T(AAT)(CGA)(ATG)(GGC) == NRMG
  - TA(ATC)(GAA)(TGG)GC == IEW
  - There are actually six ways of reading the DNA...
  - Open reading frame

## The monomers in a protein chain are the 20 amino acids

Name	Abbr.	Linear structure formula	AMINO ACIDS
<u>Alanine</u>	ala	a	CH <sub>3</sub> -CH(NH <sub>2</sub> )-COOH
<u>Arginine</u>	arg	r	HN=C(NH <sub>2</sub> )-NH-(CH <sub>2</sub> ) <sub>3</sub> -CH(NH <sub>2</sub> )-COOH
<u>Asparagine</u>	asn	n	H <sub>2</sub> N-CO-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Aspartic acid</u>	asp	d	HOOC-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Cysteine</u>	cys	c	HS-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Glutamine</u>	gln	q	H <sub>2</sub> N-CO-(CH <sub>2</sub> ) <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Glutamic acid</u>	glu	e	HOOC-(CH <sub>2</sub> ) <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Glycine</u>	gly	g	NH <sub>2</sub> -CH <sub>2</sub> -COOH
<u>Histidine</u>	his	h	NH-CH=N-C-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Isoleucine</u>	ile	i	CH <sub>3</sub> -CH <sub>2</sub> -CH(CH <sub>3</sub> )-CH(NH <sub>2</sub> )-COOH
<u>Leucine</u>	leu	l	(CH <sub>3</sub> ) <sub>2</sub> -CH-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Lysine</u>	lys	k	H <sub>2</sub> N-(CH <sub>2</sub> ) <sub>4</sub> -CH(NH <sub>2</sub> )-COOH
<u>Methionine</u>	met	m	CH <sub>3</sub> -S-(CH <sub>2</sub> ) <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Phenylalanine</u>	phe	f	Ph-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Proline</u>	pro	p	NH-(CH <sub>2</sub> ) <sub>3</sub> -CH-COOH
<u>Serine</u>	ser	s	HO-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Threonine</u>	thr	t	CH <sub>3</sub> -CH(OH)-CH(NH <sub>2</sub> )-COOH
<u>Tryptophan</u>	trp	w	Ph-NH-CH-C-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Tyrosine</u>	tyr	y	HO-p-Ph-CH <sub>2</sub> -CH(NH <sub>2</sub> )-COOH
<u>Valine</u>	val	v	(CH <sub>3</sub> ) <sub>2</sub> -CH-CH(NH <sub>2</sub> )-COOH

## The machines (proteins)

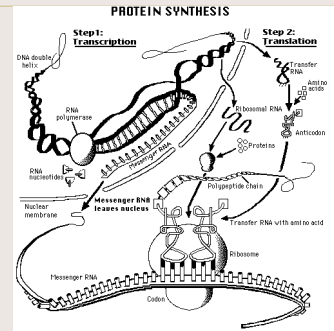


A signal RAS protein

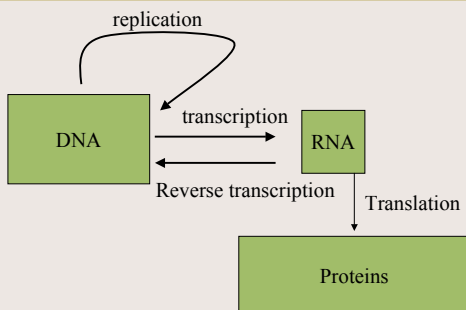


Cytochrome C  
Electron transport

## Another look on separation of information and function



## And in brief



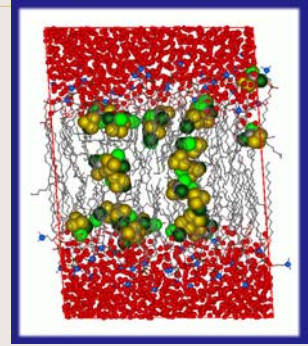
## More terminology

- We shall call a DNA segment that codes a protein – a gene
- The complete set of genetic material is called a genome
- Genomes are divided into chromosomes. Simple bacteria (*Escherichia Coli*) has 1 (5M basepairs). Humans have 46 chromosomes (3B basepairs)

## Experimental techniques of studying genomes (sequencing)

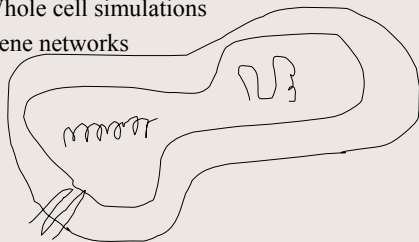
- Genetic linkage maps 1-10M bp
- Physical maps 0.1-1M bp
- Sequencing 1K-10K bp

## Lipids and cell walls



## What is life (again)?

- Divide inside from outside
  - Whole cell simulations
  - Gene networks



## What we do to understand biology?

- Use chemical physics principles
- Look for problems we already solved (homology)

## Computational tasks of molecular biology

---

- Assembly of DNA fragments
- Identifying Genes (DNA segments that code)
- Building evolutionary trees
- Gene networks & micro array analysis
- Whole cell simulations
- Identify gene function
  - Chemical physics
  - Homology