

Course Homepage

www.cs.fiu.edu/~giri/teach/Bioinf506.html

- Lecture notes, required reading material, homework, announcements, etc.
- Class: 2:00-4:45 PM, ECS 235.

1/17/06 CAP5510/CGSS166 1

Statistics Background

1/17/06 CAP5510/CGSS166 2

Basic Statistics

- Probability & Conditional Probability
- Probability space
- Random Variables
- Mean
- Standard Deviation
- Discrete vs Continuous
- Important results
 - Bayes theorem, Chebychev inequality, Markov's inequality, sum of expectations, Central Limit Theorem, ...

1/17/06 CAP5510/CGSS166 3

Discrete Distributions

- Bernoulli
 - Distribution of heads/tails
- Binomial
 - Distr of # of heads in n coin tosses
- Geometric
 - Distr of # of tails before a head appears
- Negative binomial
 - Distr of # of tosses before r heads appear
- Poisson
 - Distr of # of heads in a large # of tosses

1/17/06

CAP5510/CGSS166

4

Continuous Distributions

- Uniform
- Normal, Gaussian
- Exponential
- Gamma
- Beta
- Extreme-value

1/17/06

CAP5510/CGSS166

5

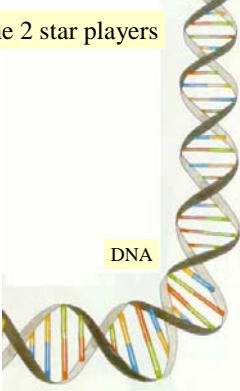
Molecular Biology Background

1/17/06

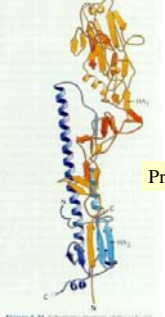
CAP5510/CGSS166

6

The 2 star players



DNA



Protein

Figure 4-24 Schematic diagram of the subunit monomers of hemoglobin from sulfonamide. The structure comprises about 500 amino acids arranged in two chains (H_β and H_γ). The H_β chain has a heme group in the diagram. The subunit is long (aligned with a long residual region) built up by residues from both chains and includes one of the longest α helices known in a globular structure, about 70 Å long. The globular head is formed by residues only from H_β. (L. Strydom et al. from Wilson, Harvard University.)

1/17/06 CAP5510/CGSS166 7

The Players

DNA

String with alphabet {A, C, G, T}
 Nucleotides/Bases

RNA

String with alphabet {A, C, G, U} Bases

Protein

String with 20-letter alphabet
 Amino acids/Residues

1/17/06 CAP5510/CGSS166 8

Typical DNA Sequence

```

1  gggagaacac  cggagaagg  agggagggc  gaagaaaagc  aacagaagcc  cagttgctgc
61  tccaggtccc  tggacagag  ctttttccat  gttgagactc  tctaaatgga  cgtgccccct
121  agtctctctt  agaagactg  cgtctctcta  aaggtcgacc  atggtgcccg  ggaccocgtg
181  tcttctagtg  tgcgtcttc  cccaggtctc  cctggggcgc  gggggccgpc  teattccaga
241  gctggggcgc  aagaatttc  cccggcctac  cagccgaccc  ttgtcccgcc  cttcggaaag
301  cgtctctcagc  gaatttgat  tgaggctgct  cagcatgttt  ggctgaagc  agagaccacc
361  ccccagcaag  gactgtgtg  tgccccctca  tatgttagat  ctgtaccgca  ggcactcagg
421  ccagccagga  gggcccgccc  cagaccaccg  gctggagagg  gcagccagcc  gcgccaaacac
481  cgtccgagc  tccatccag  aagaagcgt  ggggaactc  ccagagatga  gtgggaaaac
541  gggccggcgc  tctctctca  atttaagtc  tctcccagc  gacagtttc  tccatctcgc
601  agaactccag  attctccgg  aacagataca  ggaagctttg  ggaacagta  gttccagca
661  ccgaattaat  atttatgaa  ttataagcc  tgcagcagcc  aacttgaat  tctctgtac
721  cagactattg  gacaccaggt  tagtgaatca  gaacacaagt  cagttggaga  gcttcagct
781  cccccagct  gtgatcggt  ggaccaca  gggacacacc  aacctgggt  ttgtggtgga
841  agtggcccat  tttaggaga  acccaggtgt  ctccaagaga  catgtgagga  ttacaggtc
901  ttgpcacaa  gatgaacaa  gctggtcaca  gataaggcca  tgcctagta  ctttggaca
961  tgatggaaaa  ggacatccg  tccacaacg  agaaaagcgt  caagccaaac  acaaacagc
  
```

1/17/06 CAP5510/CGSS166 9

Typical protein sequence

```
/translation="MVAGTRCLLVLLLPQVLLGGAAGLIPELGRKKFAAASSRPLSRP  
SEDVLSSEFELRLLSMFGLKQRPTPSKDVVPPYMLDLYRRHSGQPAPADHRLERAA  
SRANTVRSFHHEEAVEELPEMSGKTARRFFNLSVPSDFLTSAELQIFREIQEAL  
GNSSFQHRINIYI I KPAANLKFVTRLLDTRLVNQNTSQWESFDVTPAVMRWTTQG  
HTNHGFVVEVAHLEENPGVSKRHVIRISRLHQDEHSWSQIRPLLVTFGHDGKGPLHK  
REKRQAKHKQRKRLKSSCKRHPLYVD FSDVGWNDWIVAPPGYHAFYCHGECFPPLADH  
LNSTNHAIVQTLVNSVNSKI PKACCVPELSAISMLYLDENEKVVLNQYQDMVVEGCG  
CR"
```

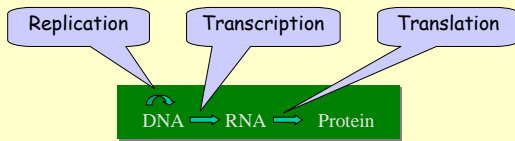
1/17/06

CAP5510/CGS5166

10

Central Dogma

- DNA acts as a template to replicate itself.
- DNA is transcribed into RNA.
- RNA is translated into **Protein**.

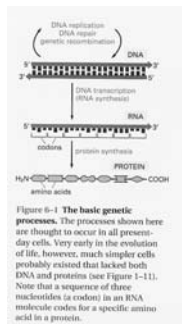


1/17/06

CAP5510/CGS5166

11

Basic Genetic Processes



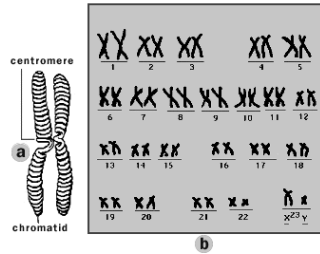
1/17/06

CAP5510/CGS5166

12

Chromosomes

Human chromosomes!



1/17/06

CAP5510/CGSS166

13

Chromosomes

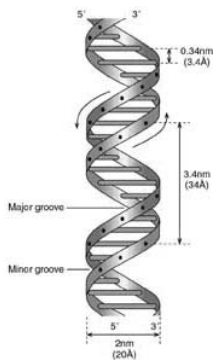


1/17/06

CAP5510/CGSS166

14

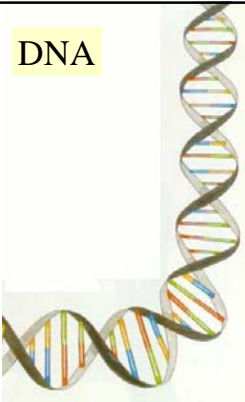
DNA Molecule



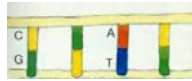
1/17/06

15

DNA



Complementary Bases



1/17/06

CAP5510/CGSS166

16

Proteins – Amino acids

amino acid	3 letter code	1 letter code
alanine	Ala	A
arginine	Arg	R
aspartic acid	Asp	D
asparagine	Asn	N
cysteine	Cys	C
glutamic acid	Glu	E
glutamine	Gln	Q
glycine	Gly	G
histidine	His	H
isoleucine	Ile	I
leucine	Leu	L
lysine	Lys	K
methionine	Met	M
phenylalanine	Phe	F
proline	Pro	P
serine	Ser	S
threonine	Thr	T
tryptophan	Trp	W
tyrosine	Tyr	Y
valine	Val	V

Table 1.1: Amino acid abbreviations

1/17/06

CAP5510/CGSS166

17

RNA

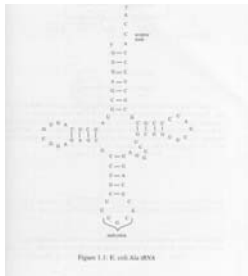



Figure 1.1: E. coli Ala-tRNA

1/17/06

CAP5510/CGSS166

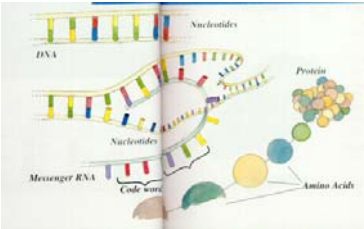
18

Genes



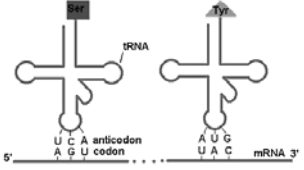
1/17/06
CAP5510/CGSS166
19

DNA → RNA → Protein



1/17/06
CAP5510/CGSS166
20

The Genetic Code



5' U C A anticodon Ser A G U codon A U G mRNA 3'

		2nd base in codon				
		U	C	A	G	
1st base in codon	U	Phe Phe Leu Leu	Ser Ser Ser STOP	Tyr Tyr STOP Trp	Cys Cys STOP Trp	U C A G
	C	Leu Leu Leu Leu	Pro Pro Pro Pro	His His Gln Gln	Arg Arg Arg Arg	U C A G
	A	Ile Ile Ile Met	Thr Thr Thr Thr	Asn Asn Lys Lys	Ser Ser Arg Arg	U C A G
	G	Val Val Val Val	Ala Ala Ala Ala	Asp Asp Glu Glu	Gly Gly Gly Gly	U C A G

The Genetic Code

1/17/06

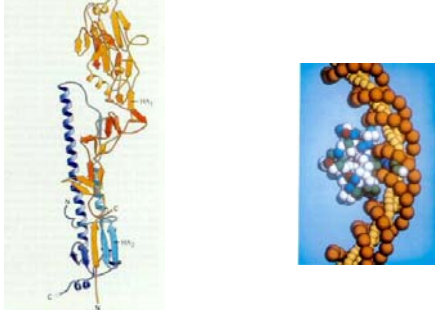
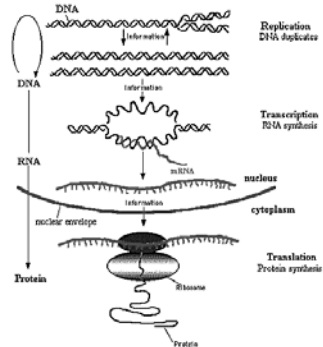


Figure 4.21 Schematic diagram of the tertiary structure of insulin. The structure comprises about 500 amino acids arranged in two chains, A1 and B1. The first half of each chain has a higher concentration of hydrophobic amino acids. The second half is charged with a long variable region that is the residue from both chains and includes one of the longest α helices known in a globular structure, about 70 Å long. The globular head is formed by residues only from the A1 chain. (Courtesy of Ben Wiley, Harvard University.)

1/17/06

AP5510/CG55166

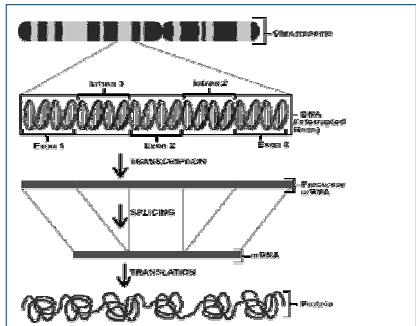
22



The Central Dogma of Molecular Biology

1/17/06

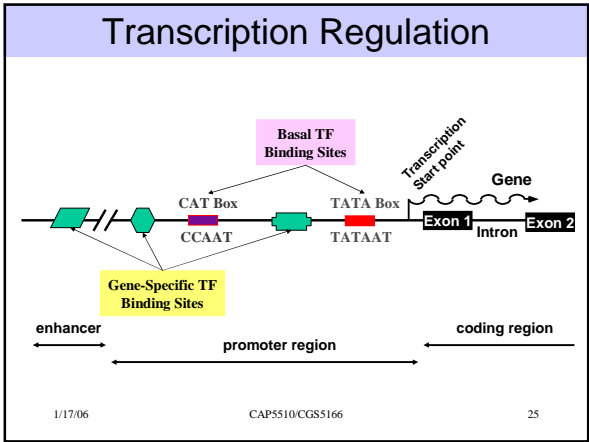
23

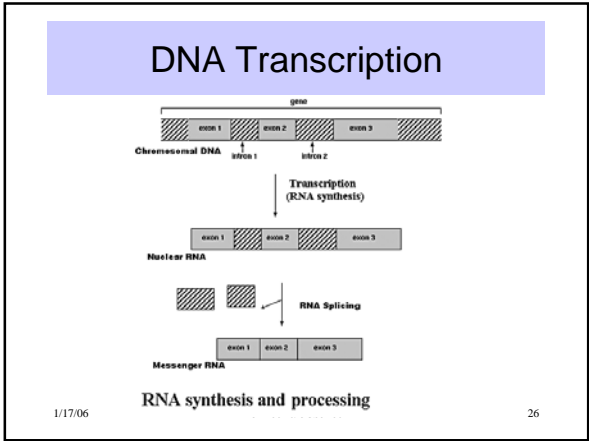


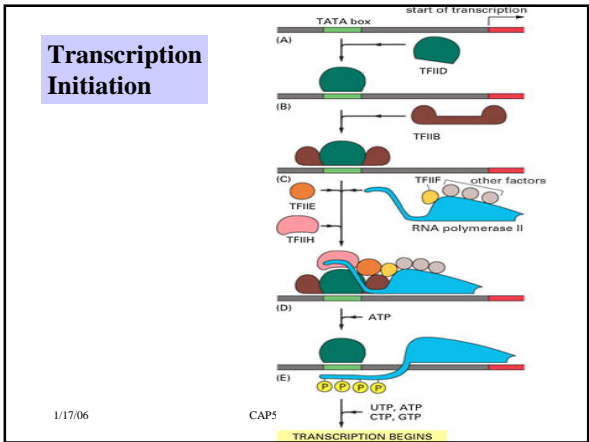
1/17/06

CAP5510/CG55166

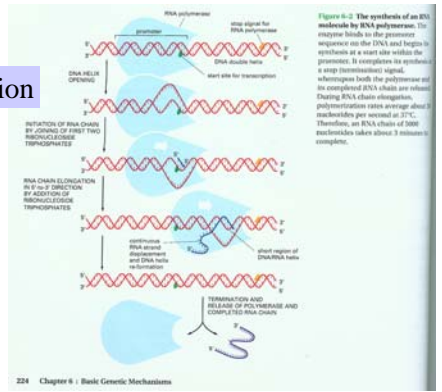
24







Transcription



1/17/06

CAP5510/CGSS166

28

Transcription Steps

- RNA polymerase needs many transcription factors (TFIIA,TFIIB, etc.)
- (A) The promoter sequence (TATA box) is located 25 nucleotides away from transcription initiation site.
 - (B) The TATA box is recognized and bound by transcription factor TFIID, which then enables the adjacent binding of TFIIB. DNA is somewhat distorted in the process.
 - (D) The rest of the general transcription factors as well as the RNA polymerase itself assemble at the promoter. What order?
 - (E) TFIIF then uses ATP to phosphorylate RNA polymerase II, changing its conformation so that the polymerase is released from the complex and is able to start transcribing. As shown, the site of phosphorylation is a long polypeptide tail that extends from the polymerase molecule.

1/17/06

CAP5510/CGSS166

29

Transcription Factors

- The general transcription factors have been highly conserved in evolution; some of those from human cells can be replaced in biochemical experiments by the corresponding factors from simple yeasts.

1/17/06

CAP5510/CGSS166

30

Protein Synthesis

1. Transcription

Protein synthesis

1/17/06 CAP5510/CGSS166 31

Protein Synthesis:

Incorporation of amino acid into protein

1/17/06 CAP5510/CGSS166 32

1/17/06 33
