

CpG Islands

- Regions in DNA sequences with increased occurrences of substring "CG"
- Rare: typically C gets methylated and then mutated into a T.
- Often around promoter or "start" regions of genes
- Few hundred to a few thousand bases long

Problem 1:

- **Input:** Small sequence **S**
- **Output:** Is **S** from a CpG island?
 - Build Markov models: M_+ and M_-
 - Then compare

Markov Models

+	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182

-	A	C	G	T
A	0.300	0.205	0.285	0.210
C	0.322	0.298	0.078	0.302
G	0.248	0.246	0.298	0.208
T	0.177	0.239	0.292	0.292

How to distinguish?

- Compute

$$S(x) = \log\left(\frac{P(x | M+)}{P(x | M-)}\right) = \sum_{i=1}^L \log\left(\frac{p_{x(i-1)xi}}{m_{x(i-1)xi}}\right) = \sum_{i=1}^L r_{x(i-1)xi}$$

r=p/m	A	C	G	T
A	-0.740	0.419	0.580	-0.803
C	-0.913	0.302	1.812	-0.685
G	-0.624	0.461	0.331	-0.730
T	-1.169	0.573	0.393	-0.679

Score(GCAC)

$$= .461 - .913 + .419 < 0.$$

GCAC not from CpG island.

Score(GCTC)

$$= .461 - .685 + .573 > 0.$$

GCTC from CpG island.

Problem 1:

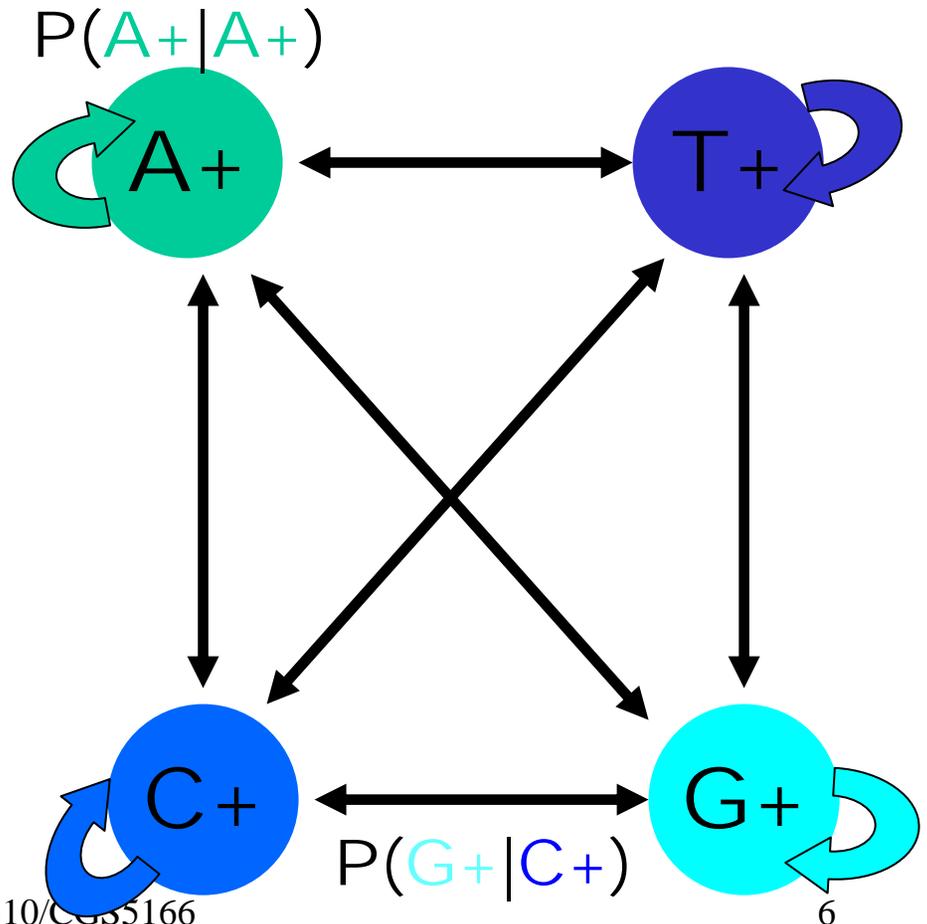
- **Input:** Small sequence **S**
- **Output:** Is **S** from a CpG island?
 - Build Markov Models: M_+ & M_-
 - Then compare

Problem 2:

- **Input:** Long sequence **S**
- **Output:** Identify the CpG islands in **S**.
 - Markov models are inadequate.
 - Need Hidden Markov Models.

Markov Models

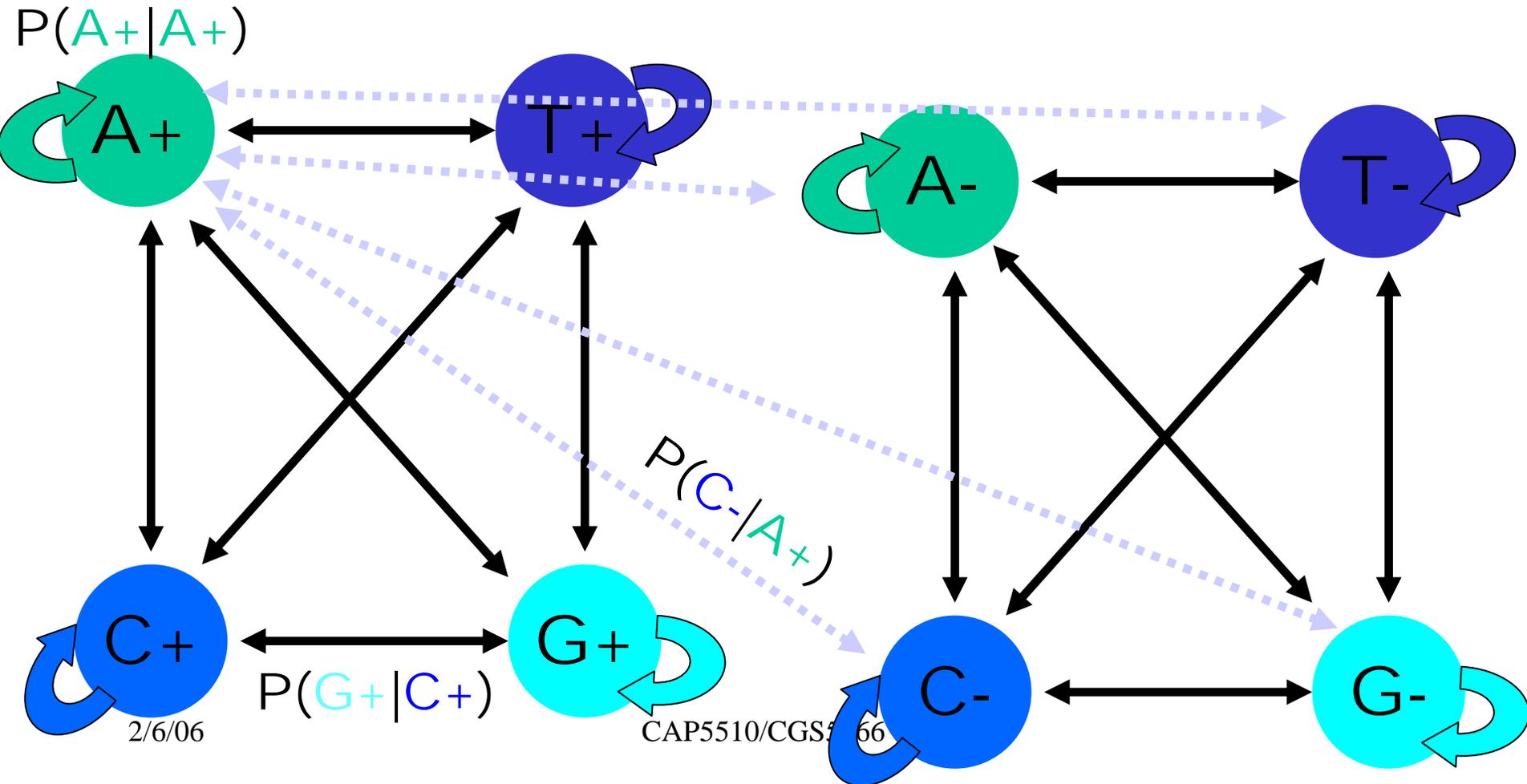
+	A	C	G	T
A	0.180	0.274	0.426	0.120
C	0.171	0.368	0.274	0.188
G	0.161	0.339	0.375	0.125
T	0.079	0.355	0.384	0.182



CpG Island + in an ocean of -

First order Hidden Markov Model

MM=16, HMM= 64 transition probabilities (adjacent bp)

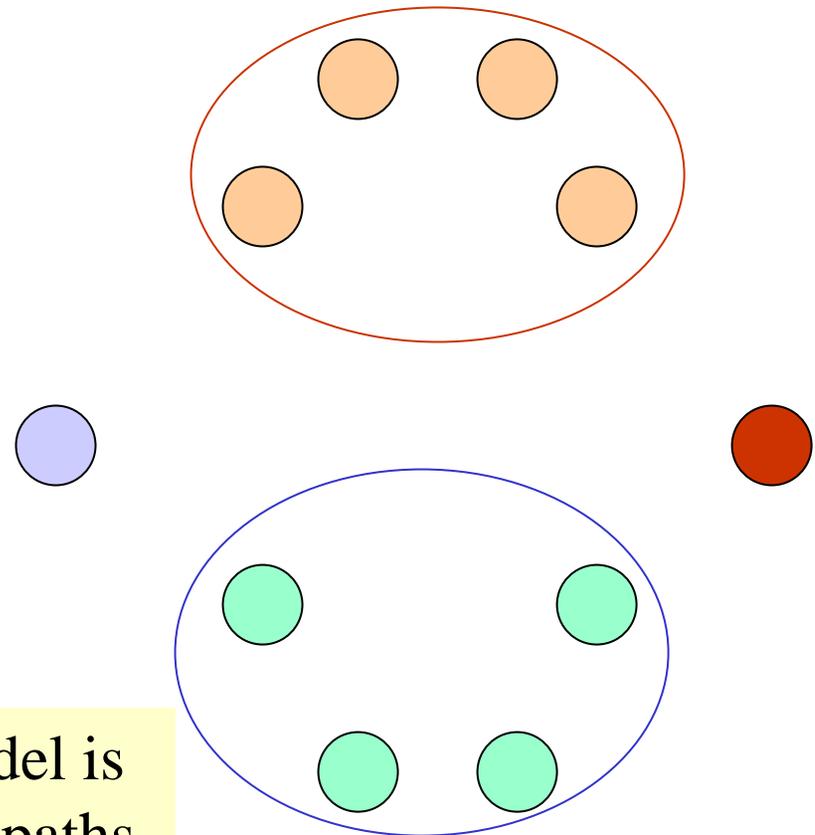


Hidden Markov Model (HMM)

- States
- Transitions
- Transition Probabilities
- Emissions
- Emission Probabilities

- What is hidden about HMMs?

Answer: The path through the model is hidden since there are many valid paths.



How to Solve Problem 2?

- Solve the following problem:

Input: Hidden Markov Model M ,
parameters Θ , emitted sequence S

Output: Most Probable Path Π

How: Viterbi's Algorithm (Dynamic Programming)

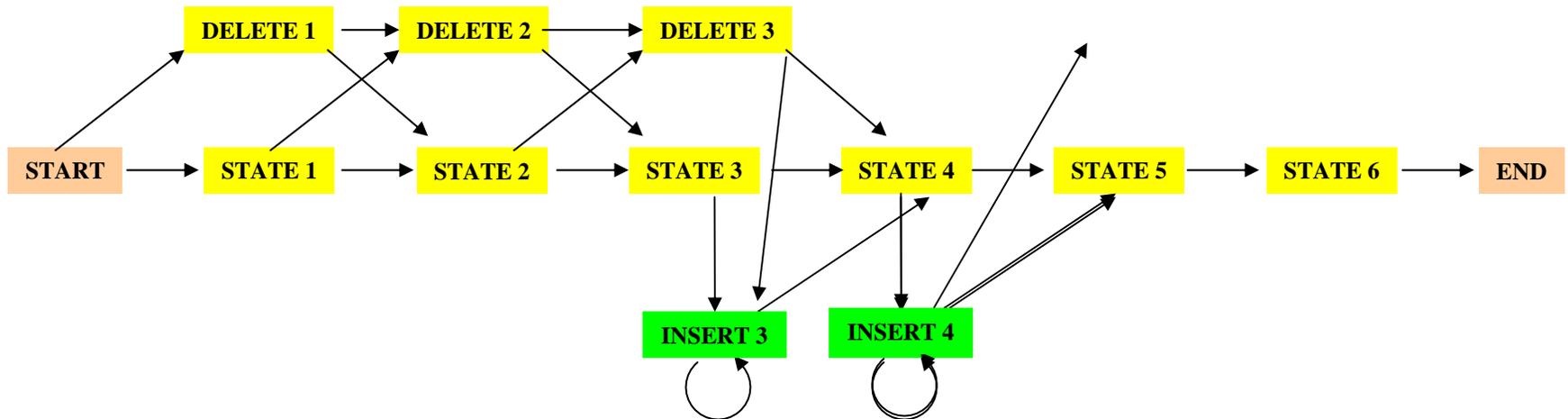
Define $\Pi[i,j]$ = MPP for first j characters of S ending in state i

Define $P[i,j]$ = Probability of $\Pi[i,j]$

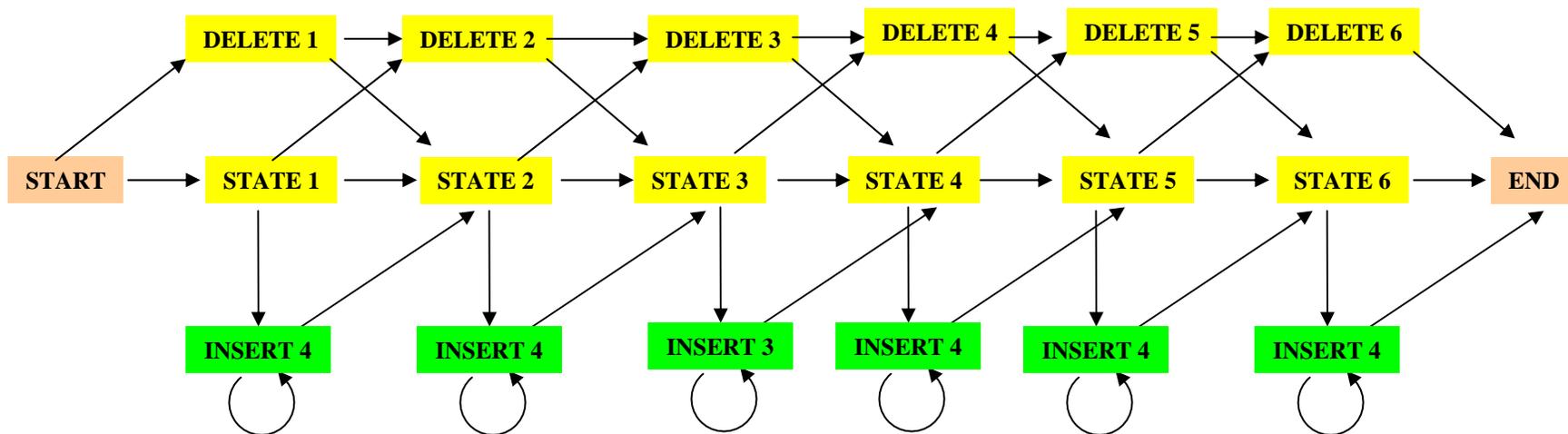
- Compute state i with largest $P[i,j]$.

Profile HMMs with InDels

- Insertions
- Deletions
- Insertions & Deletions



Profile HMMs with InDels

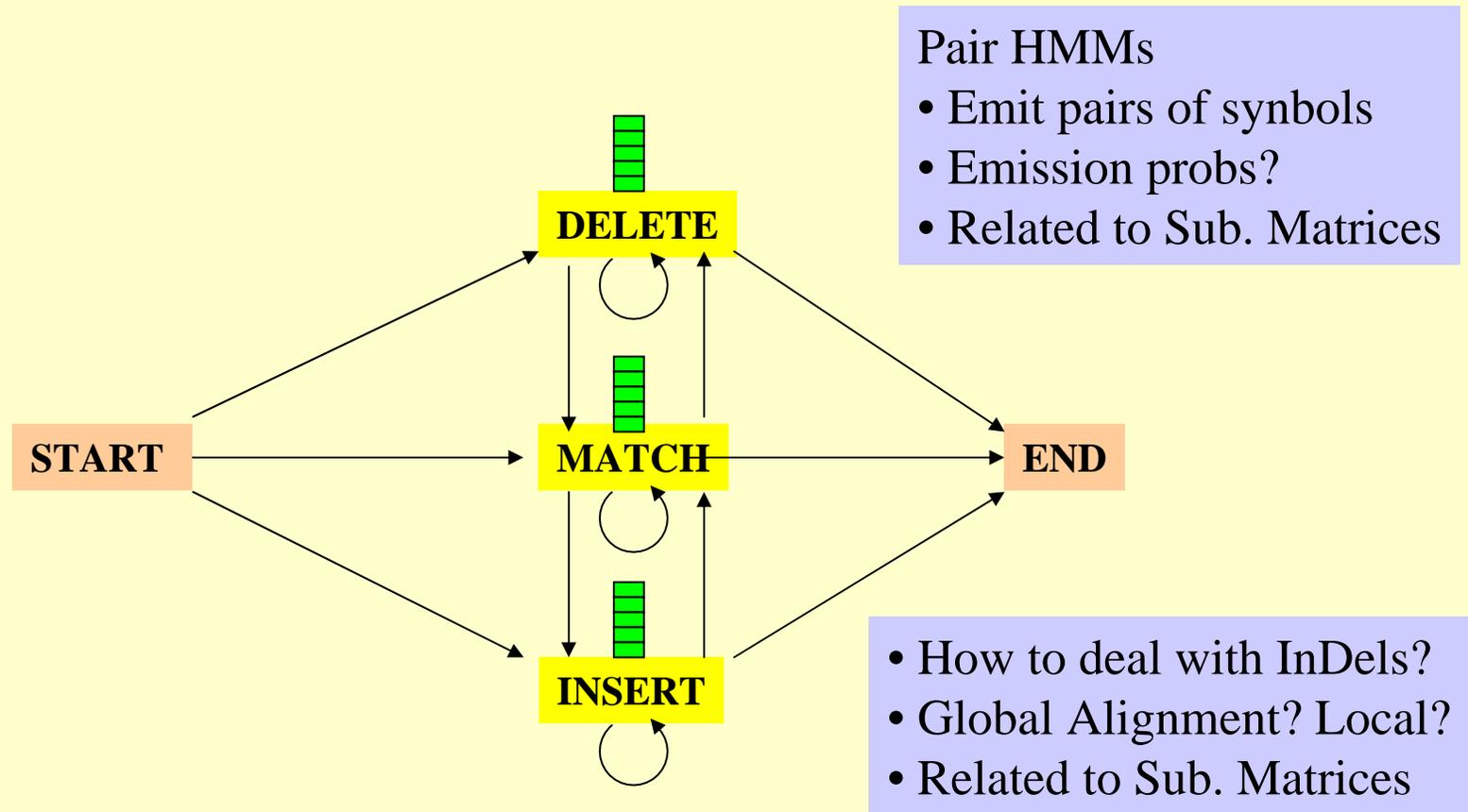


Missing transitions from **DELETE j** to **INSERT j** and
from **INSERT j** to **DELETE $j+1$** .

How to model Pairwise Sequence Alignment

LEAPVE

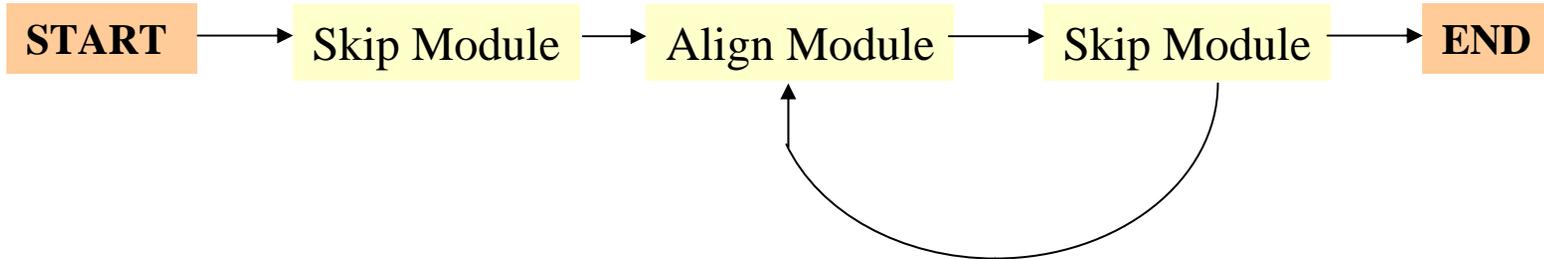
LAPVIE



How to model Pairwise Local Alignments?

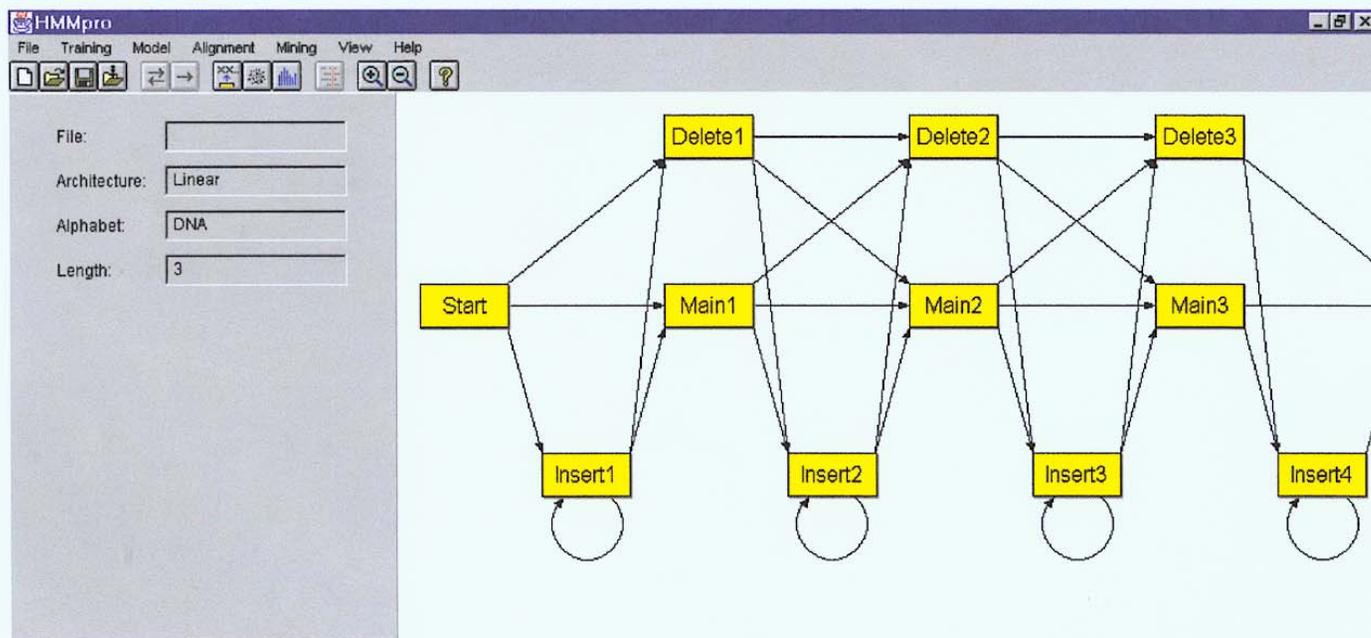


How to model Pairwise Local Alignments with gaps?



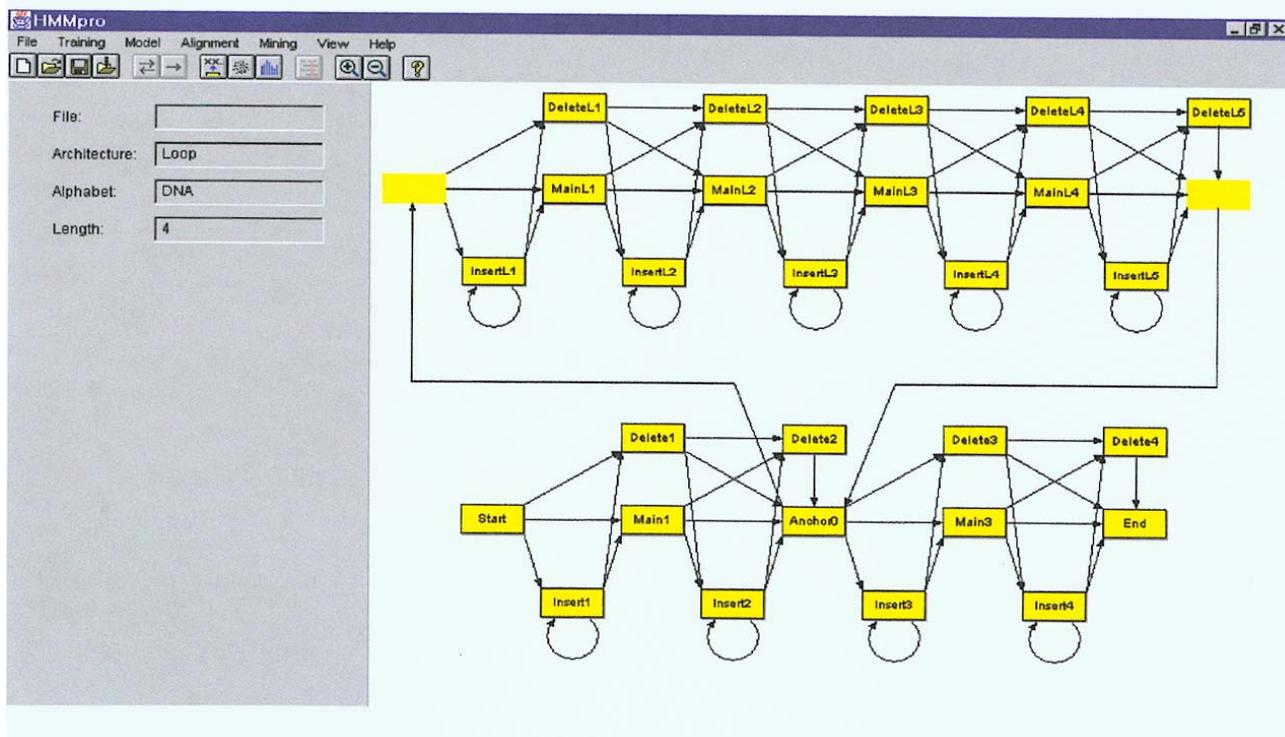
Standard HMM architectures

Linear Architecture



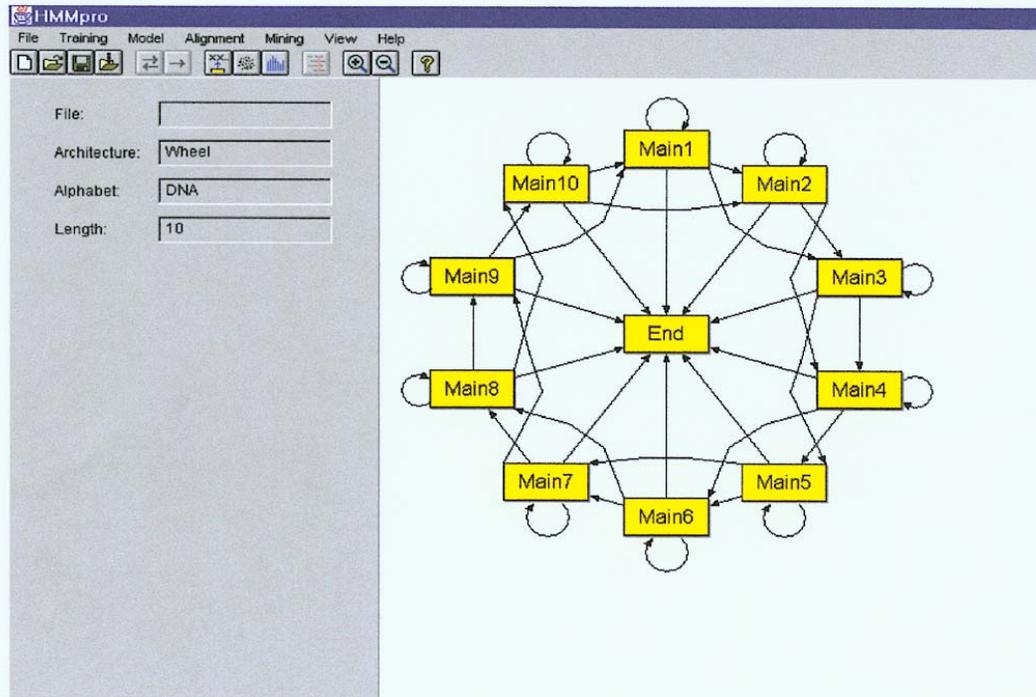
Standard HMM architectures

Loop Architecture



Standard HMM architectures

Wheel Architecture



Profile HMMs from Multiple Alignments

HBA_HUMAN	VGA--HAGEY
HBB_HUMAN	V-----NVDEV
MYG_PHYCA	VEA--DVAGH
GLB3_CHITP	VKG-----D
GLB5_PETMA	VYS--TYETS
LGB2_LUPLU	FNA--NIPKH
GLB1_GLYDI	IAGADNGAGV

Construct Profile HMM from above multiple alignment.

HMM for Sequence Alignment

A. Sequence alignment

N	•	F	L	S
N	•	F	L	S
N	K	Y	L	T
Q	•	W	-	T

RED POSITION REPRESENTS ALIGNMENT IN COLUMN

GREEN POSITION REPRESENTS INSERT IN COLUMN

PURPLE POSITION REPRESENTS DELETE IN COLUMN

B. Hidden Markov model for sequence alignment

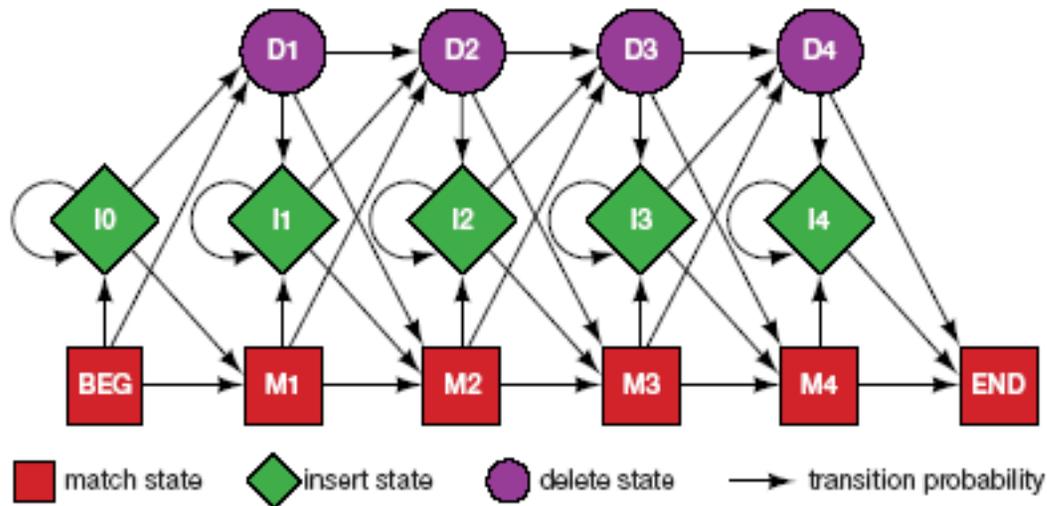


FIGURE 5.16. Relationship between the sequence alignment and the hidden Markov model of the alignment (Krogh et al. 1994). This particular form for the HMM was chosen to represent the sequence, structural, and functional variation expected in proteins. The model accommodates the identities, mismatches, insertions, and deletions expected in a group of related proteins. (A) A section of an msa. The illustration shows the columns generated in an msa. Each column may include matches and mismatches (*red* positions), insertions (*green* positions), and deletions (*purple* positions). (B) The HMM. Each column in the model represents the possibility of a match, insert, or delete in each column of the alignment in A. The HMM is a probabilistic representation of a section of the msa. Sequences can be generated from the HMM by starting at the beginning state labeled BEG and then by following

Problem 3: LIKELIHOOD QUESTION

- **Input:** Sequence **S**, model **M**, state **i**
- **Output:** Compute the probability of reaching state **i** with sequence **S** using model **M**
 - **Backward Algorithm (DP)**

Problem 4: LIKELIHOOD QUESTION

- **Input:** Sequence **S**, model **M**
- **Output:** Compute the probability that **S** was emitted by model **M**
 - **Forward Algorithm (DP)**

Problem 5: LEARNING QUESTION

- **Input:** model structure M , Training Sequence S
- **Output:** Compute the parameters Θ
- **Criteria:** ML criterion
 - maximize $P(S | M, \Theta)$ HOW???

Problem 6: DESIGN QUESTION

- **Input:** Training Sequence S
- **Output:** Choose model structure M , and compute the parameters Θ
 - No reasonable solution
 - Standard models to pick from

Iterative Solution to the **LEARNING QUESTION** (Problem 5)

- Pick initial values for parameters Θ_0
- Repeat
 - Run training set S on model M
 - Count # of times transition $i \Rightarrow j$ is made
 - Count # of times letter x is emitted from state i
 - Update parameters Θ
- Until (some stopping condition)

Entropy

- **Entropy** measures the variability observed in given data.

$$E = -\sum_c p_c \log p_c$$

- Entropy is useful in multiple alignments & profiles.
- Entropy is max when uncertainty is max.

G-Protein Couple Receptors

- Transmembrane proteins with 7 α -helices and 6 loops; many subfamilies
- Highly variable: 200-1200 aa in length, some have only 20% identity.
- [Baldi & Chauvin, '94] HMM for GPCRs
- HMM constructed with 430 match states (avg length of sequences); Training: with 142 sequences, 12 iterations

GPCR - Analysis

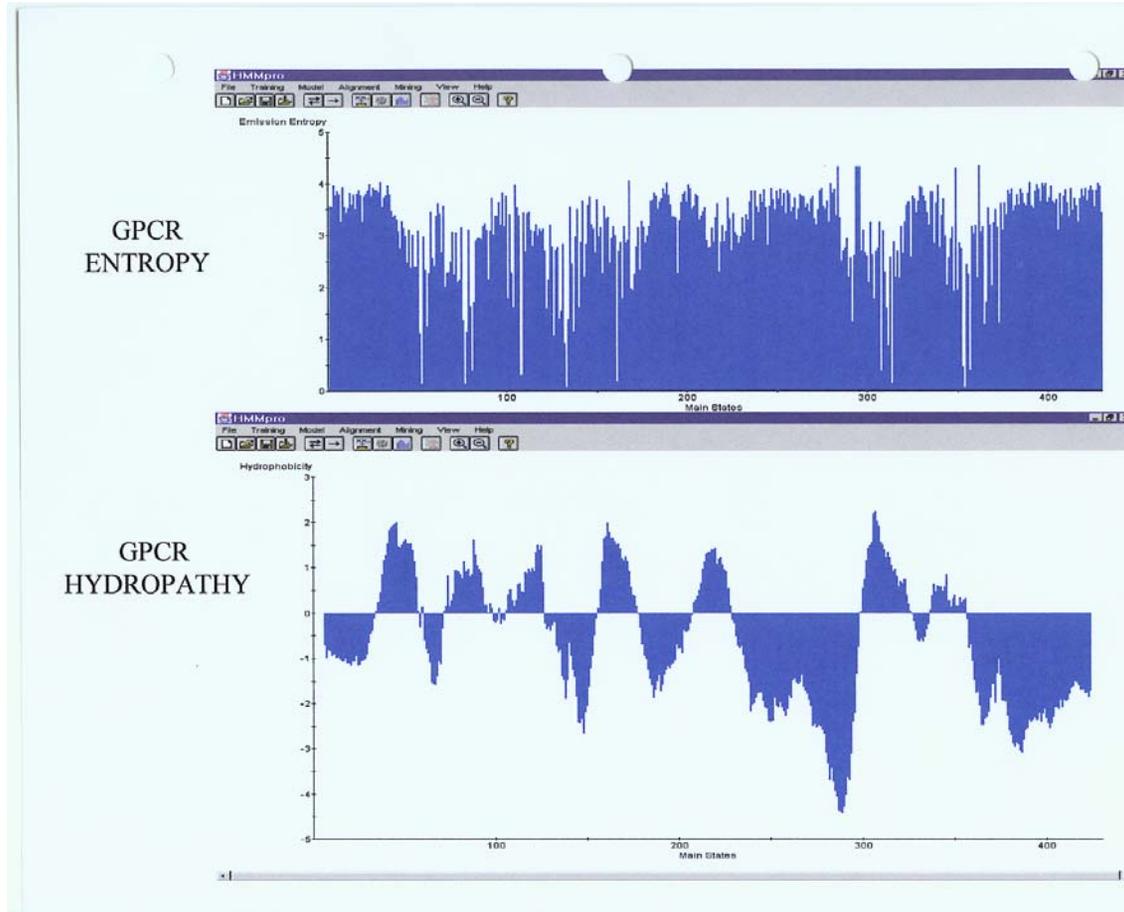
- Compute main state entropy values

$$H_i = -\sum_a e_{ia} \log e_{ia}$$

- For every sequence from test set (142) & random set (1600) & all SWISS-PROT proteins
 - Compute the negative log of probability of the most probable path π

$$\text{Score}(S) = -\log(P(\pi | S, M))$$

GPCR Analysis



Entropy

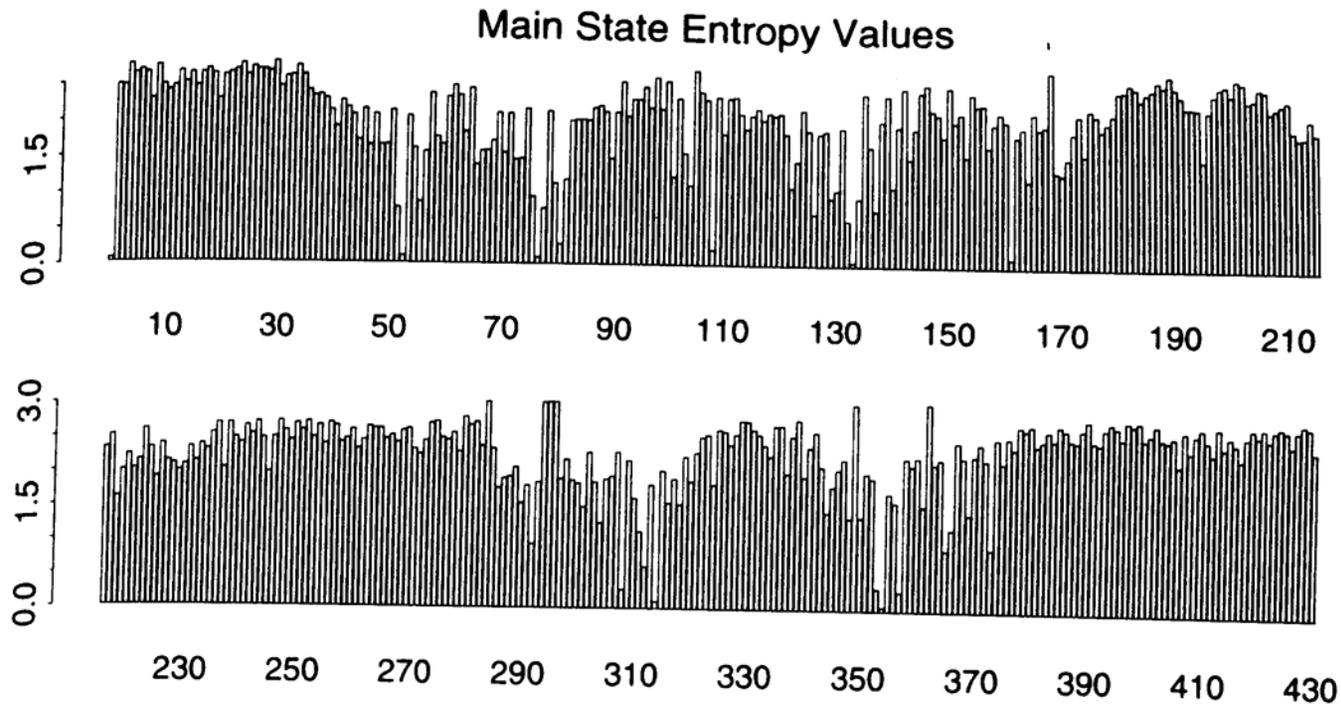


Figure 8.1: Entropy Profile of the Emission Probability Distributions Associated with the Main States of the HMM After 12 Cycles of Training.

GPCR Analysis (Cont'd)

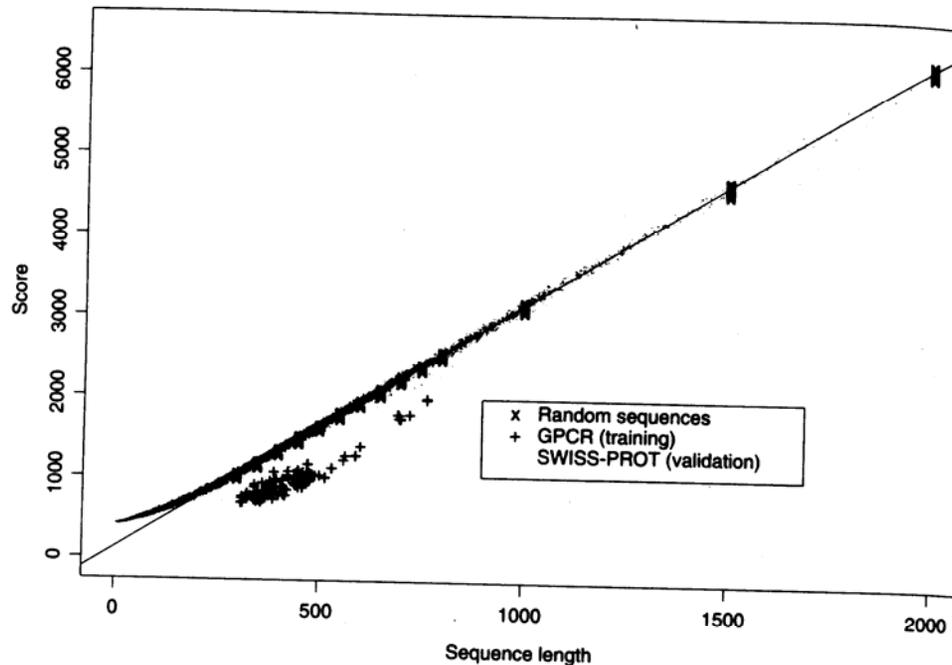


Figure 8.2: Scores (Negative Log-likelihoods of Optimal Viterbi Paths). Represented sequences consist of 142 GPCR training sequences, all sequences from the SWISS-PROT database of length less than or equal to 2000, and 220 randomly generated sequences with same average composition as the GPCRs of length 300, 350, 400, 450, 500, 550, 600, 650, 700, 750, 800 (20 at each length). The regression line was obtained from the 220 random sequences. The horizontal distances in the histogram correspond to normalized scores (6).

Applications of HMM for GPCR

- Bacteriorhodopsin
 - Transmembrane protein with 7 domains
 - But it is not a GPCR
 - Compute score and discover that it is close to the regression line. **Hence not a GPCR.**
- Thyrotropin receptor precursors
 - All have long initial loop on **INSERT STATE 20.**
 - Also clustering possible based on distance to regression line.

HMMs – Advantages

- Sound statistical foundations
- Efficient learning algorithms
- Consistent treatment for insert/delete penalties for alignments in the form of locally learnable probabilities
- Capable of handling inputs of variable length
- Can be built in a modular & hierarchical fashion; can be combined into libraries.
- Wide variety of applications: **Multiple Alignment, Data mining & classification, Structural Analysis, Pattern discovery, Gene prediction.**

HMMs – Disadvantages

- Large # of parameters.
- Cannot express dependencies & correlations between hidden states.

Perl: Practical Extraction & Report Language

- Created by Larry Wall, early 90s
- Portable, "glue" language for interfacing C/Fortran code, WWW/CGI, graphics, numerical analysis and much more
- Easy to use and extensible
- OOP support, simple databases, simple data structures.
- From interpreted to compiled
- high-level features, and relieves you from manual memory management, segmentation faults, bus errors, most portability problems, etc, etc.
- Competitors: Python, Tcl, Java

Perl Features

- Perl - many features
 - Bit Operations
 - Pattern Matching
 - Subroutines
 - Packages & Modules
 - Objects
 - Interprocess Communication
 - Threads , Process control
 - Compiling

BioPerl

- Routines for handling biosequence and alignment data.
- Why? Human Genome Project: Same project, same data. **different data formats!** Different input formats. Different output formats for comparable utility programs.
- BioPerl was useful to interchange data and meaningfully exchange results. "Perl Saved the Human Genome Project"
- Many routine tasks automated using BioPerl.
- String manipulations (string operations: substring, match, etc.; handling string data: names, annotations, comments, bibliographical references; regular expression operations)
- Modular: modules in any language

Managing a Large Project

- Devise a common data exchange format.
- Use modules that have already been developed.
- Write Perl scripts to convert to and from common data exchange format.
- Write Perl scripts to “glue” it all together.

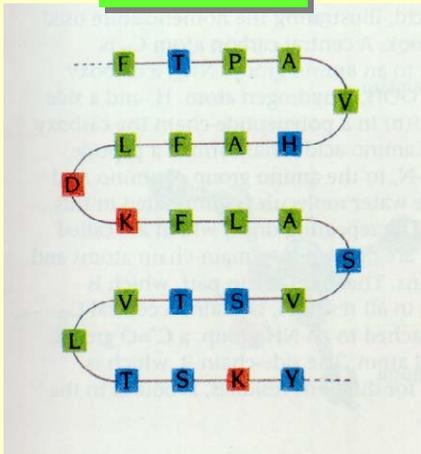
Miscellaneous

- pTk - to enable building Perl-driven GUIs for X-Window systems.
- BioJava
- BioPython
- The BioCORBA Project provides an object-oriented, language neutral, platform-independent method for describing and solving bioinformatics problems.

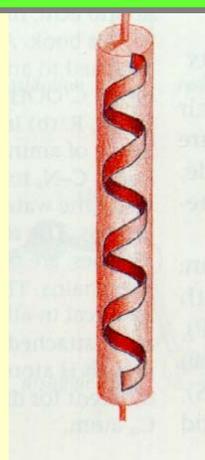
Protein Structures

- Sequences of amino acid residues
- 20 different amino acids

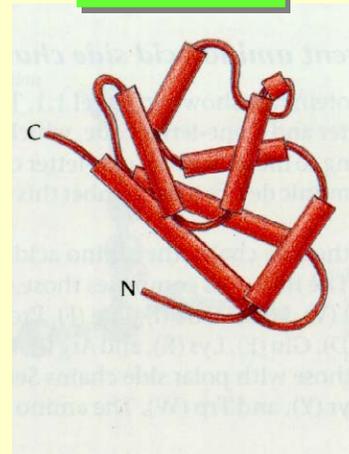
Primary



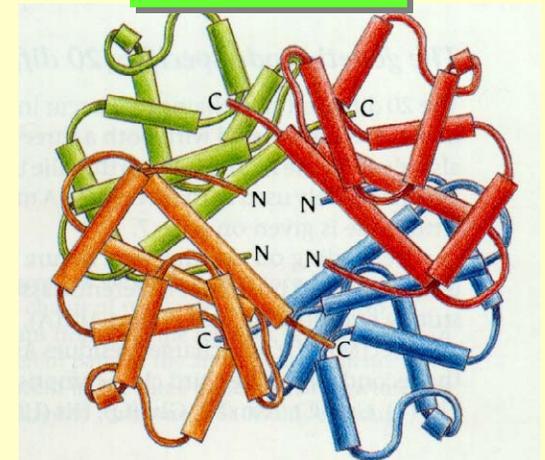
Secondary



Tertiary



Quaternary



Amino Acid Types

- **Hydrophobic** **I, L, M, V, A, F, P**
- **Charged**
 - **Basic** **K, H, R**
 - **Acidic** **E, D**
- **Polar** **S, T, Y, H, C, N, Q, W**
- **Small** **A, S, T**
- **Very Small** **A, G**
- **Aromatic** **F, Y, W**

All 3 figures are cartoons of an amino acid residue.

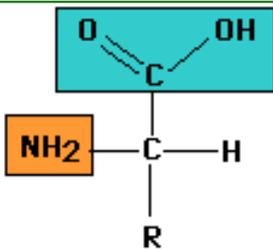
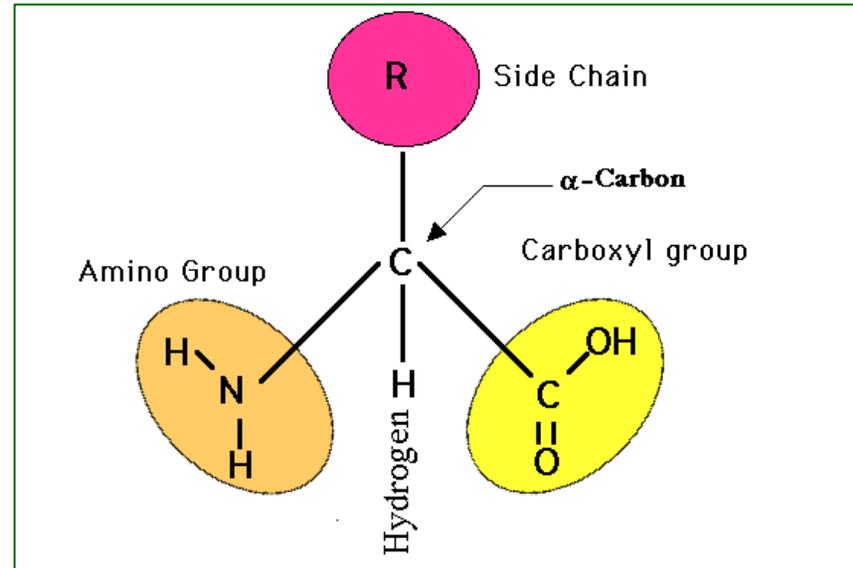
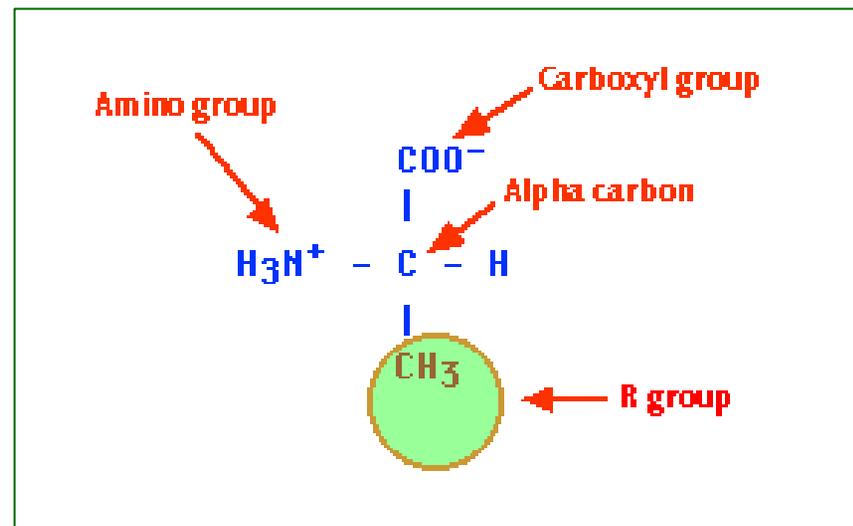


Fig. General formula for an amino acid molecule. "R" represents the variable groups that are attached to this basic molecule to make up the 20 common amino acids



Angles ϕ and ψ in the polypeptide chain

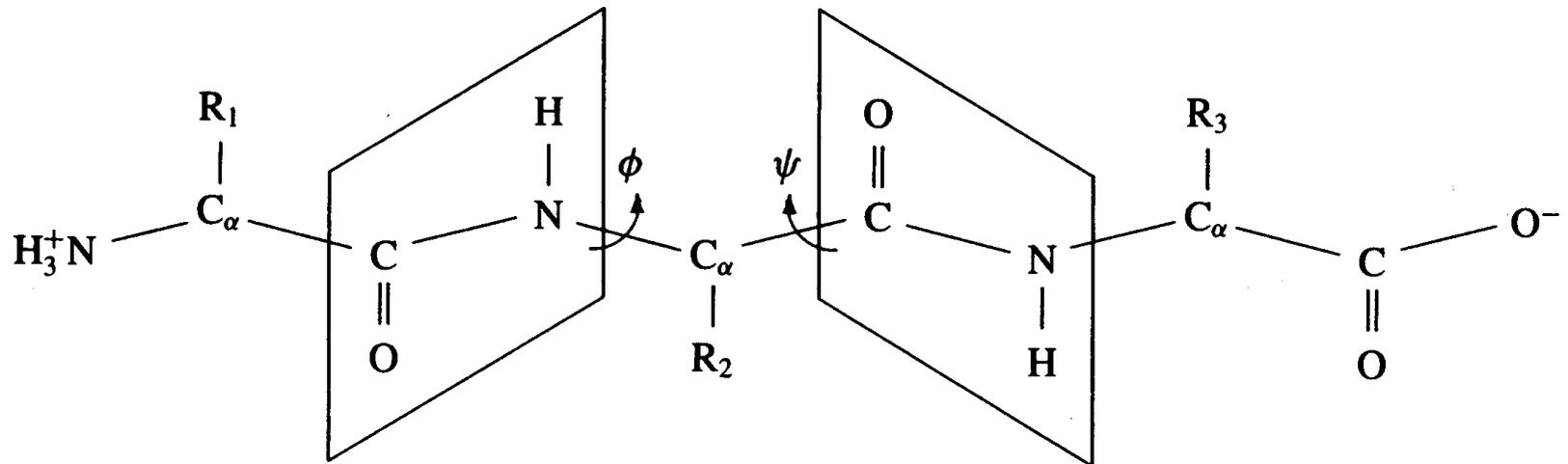
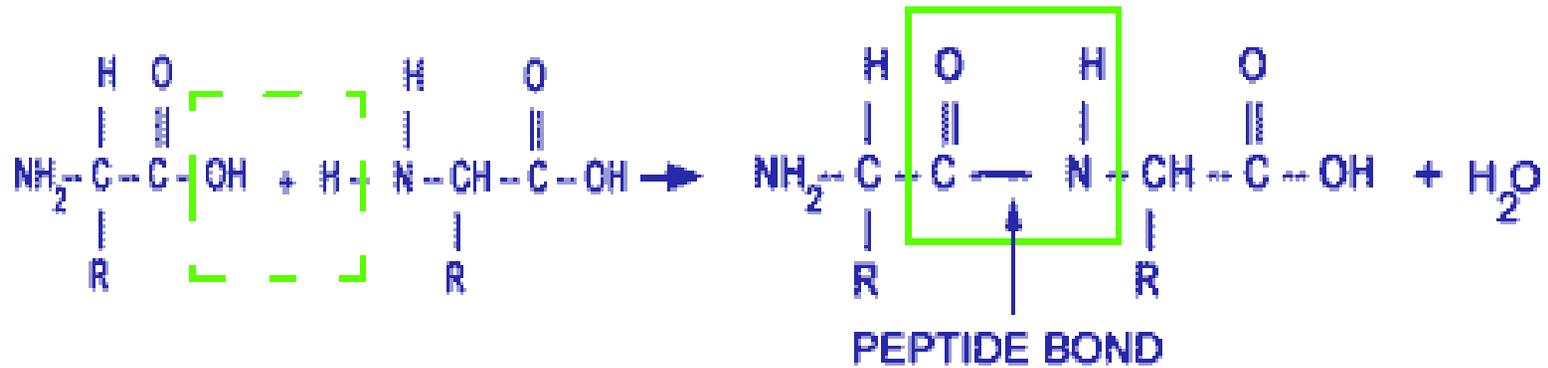
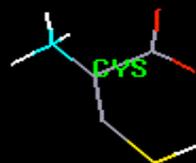
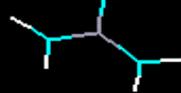
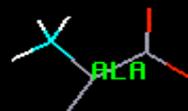
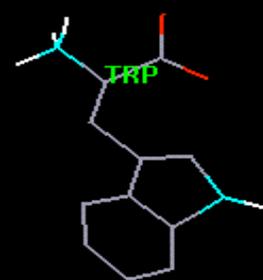
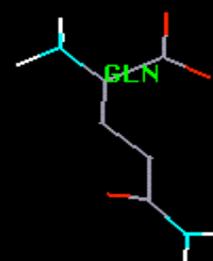
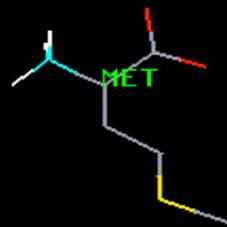
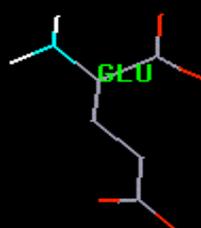
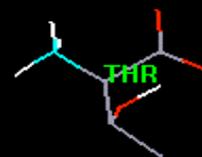
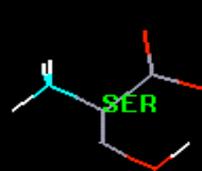
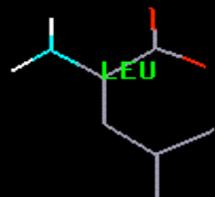
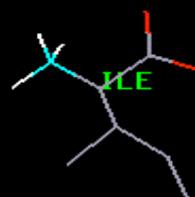
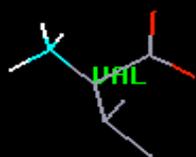
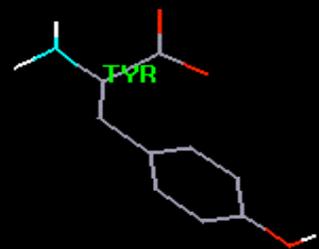


FIGURE 1.2

A polypeptide chain. The R_i side chains identify the component amino acids. Atoms inside each quadrilateral are on the same plane, which can rotate according to angles ϕ and ψ .

Peptide bonds in chains of residues





Proteins

- **Primary structure** is the sequence of amino acid residues of the protein, e.g.,

Flavodoxin:

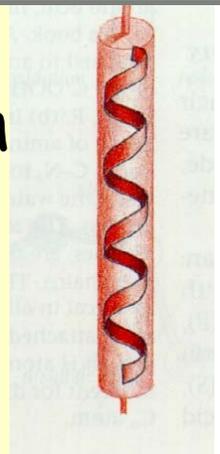
AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADA...

Secondary

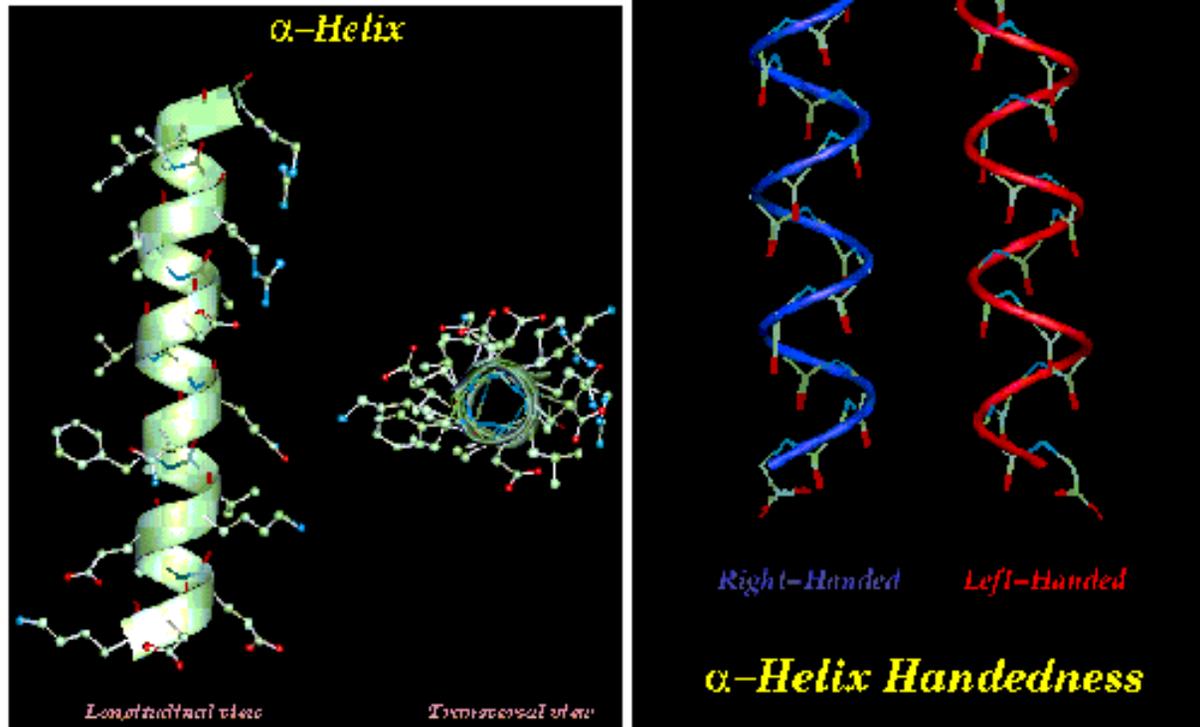
- Different regions of the sequence form local regular **secondary structures**, such

- **Alpha helix**, **beta strands**, etc.

AKIGLFYGTQTGVTQTIAESIQQEFGGESIVDLNDIANADA...



Alpha helices



(c) David Gilbert, Aik Choon Tan, Gilleain Torrance and Mallika Veeramalai 2002

16

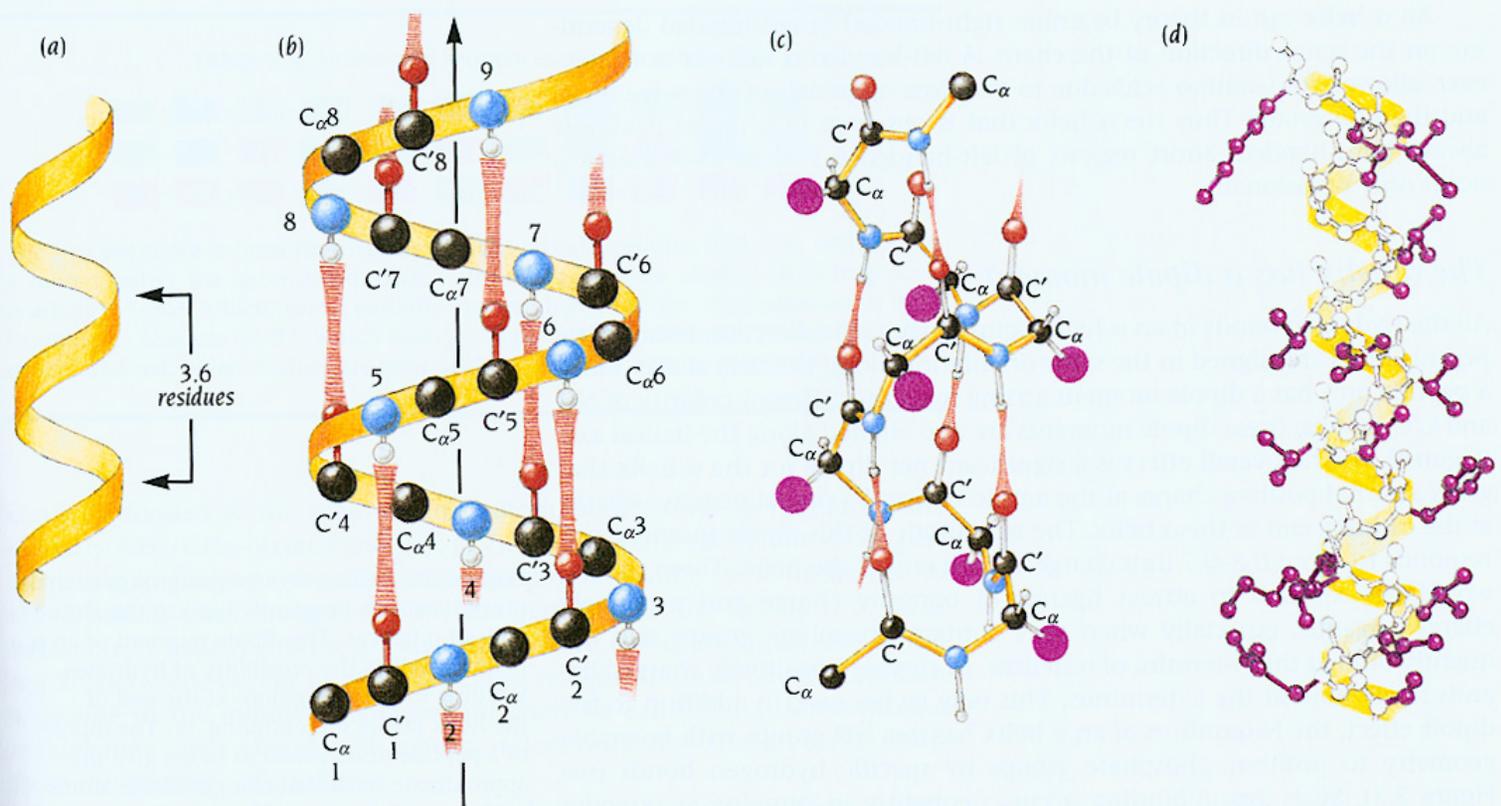
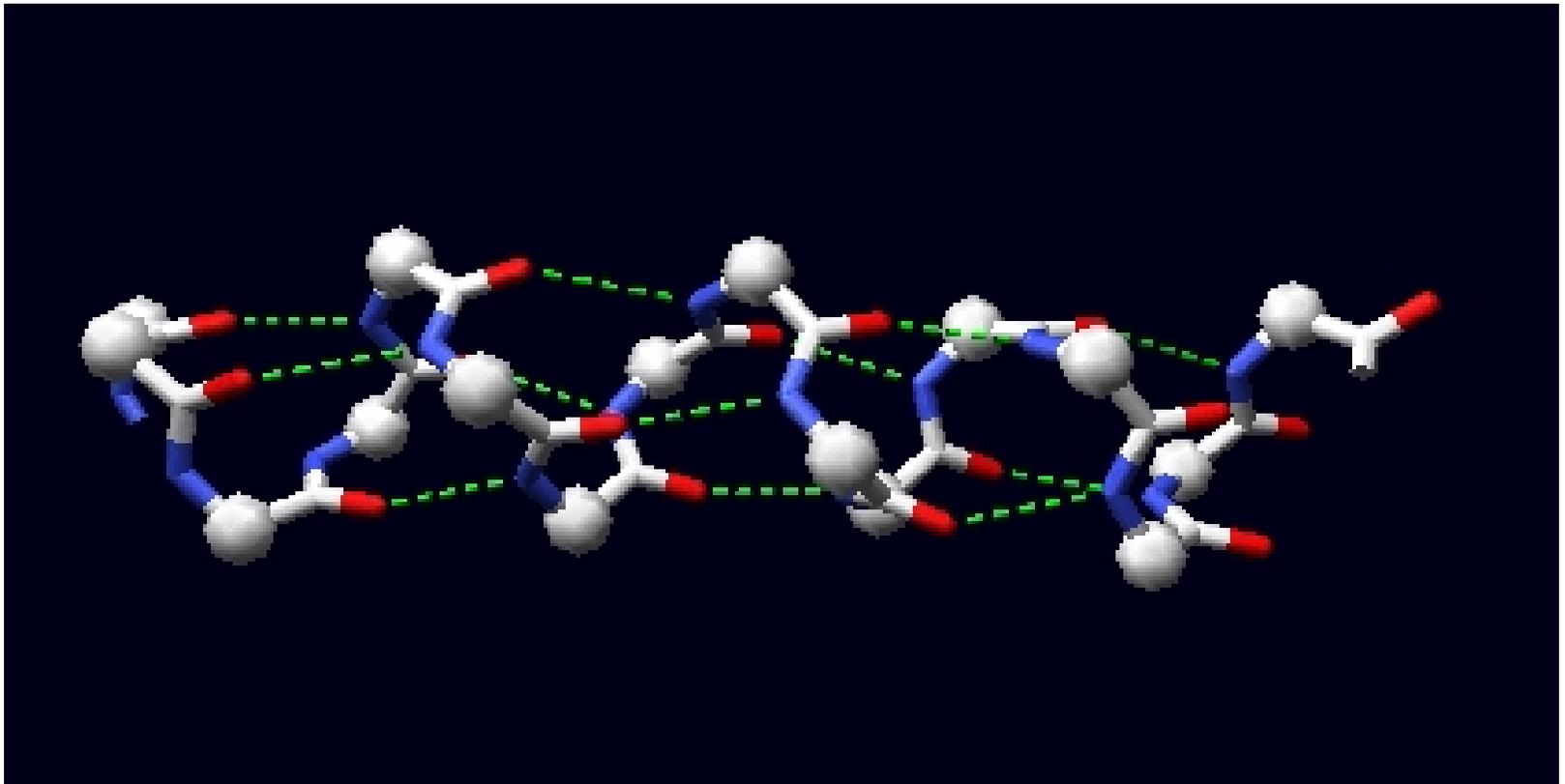


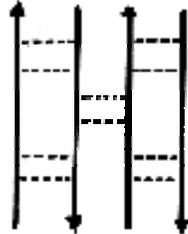
Figure 2.2 The α helix is one of the major elements of secondary structure in proteins. Main-chain N and O atoms are hydrogen-bonded to each other within α helices. (a) Idealized diagram of the path of the main chain in an α helix. Alpha helices are frequently illustrated in this way. There are 3.6 residues per turn in an α helix, which corresponds to 5.4 Å (1.5 Å per residue). (b) The same as (a) but with approximate positions for main-chain atoms and hydrogen bonds included. The arrow denotes the direction from the N-terminus to the C-terminus. (c) Schematic diagram of an α helix. Oxygen atoms are red, and N atoms are blue. Hydrogen bonds between O and N are red and striated. The side chains are represented as purple circles. (d) A ball-and-stick model of one α helix in myoglobin. The path of the main chain is outlined in yellow; side chains are purple. Main-chain atoms are not colored. (e) One turn of an α helix viewed down the helical axis. The purple side chains project out from the α helix.

Alpha Helix

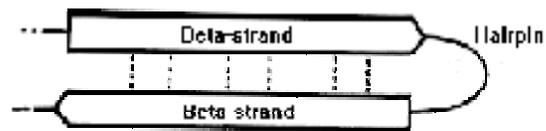


Beta sheet

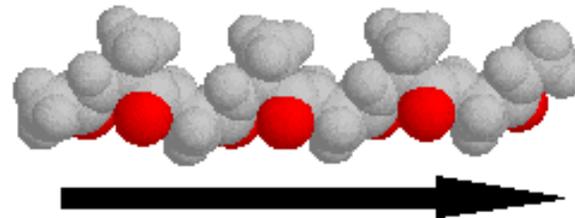
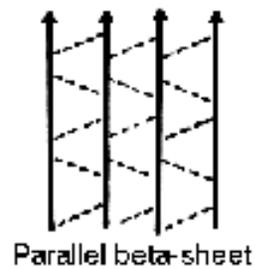
Antiparallel beta-sheet



The beta-hairpin turn.



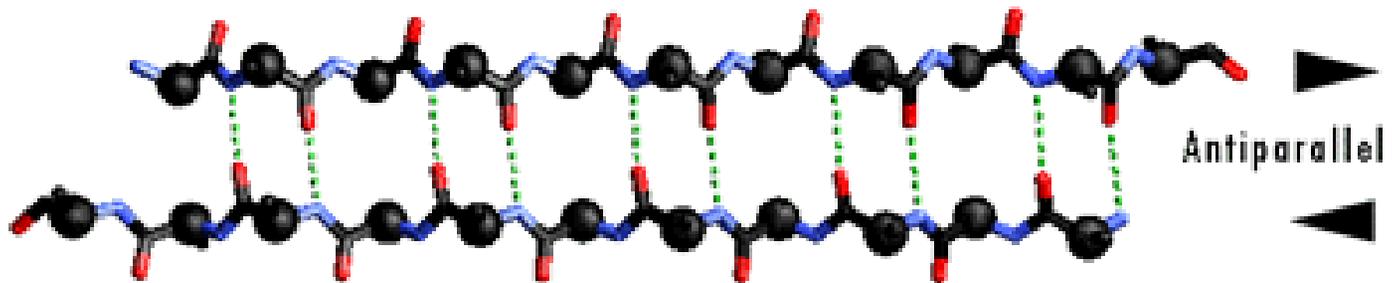
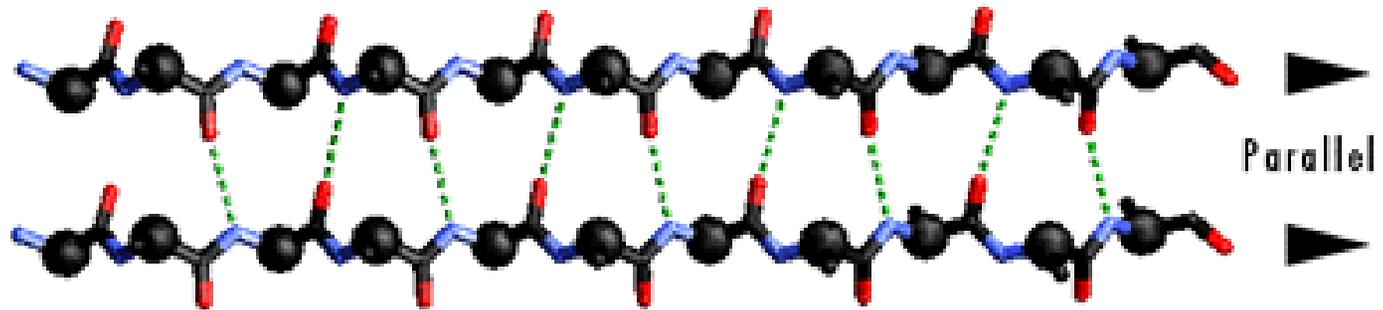
The dashed lines indicate main chain hydrogen bonds.



(c) David Gilbert, Aik Choon Tan, Gillesain Torrance and Mallika Veeramalai 2002

17

Beta Strand

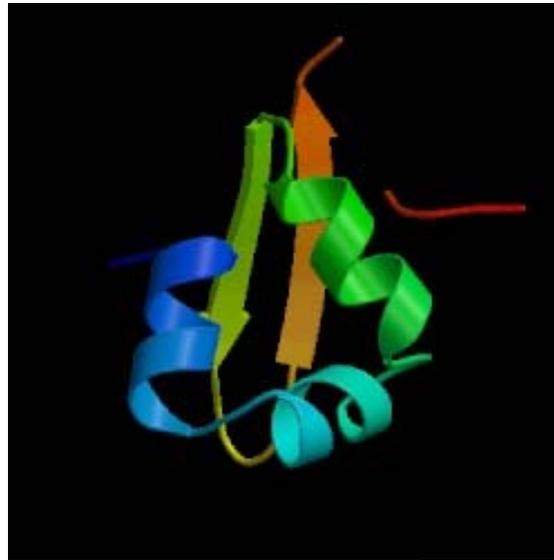


Proteins

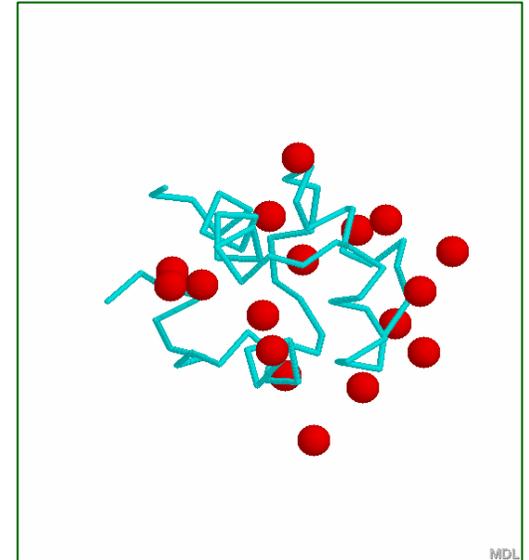
- **Tertiary structures** are formed by packing secondary structural elements into a globular structure.



Myoglobin



Lambda Cro

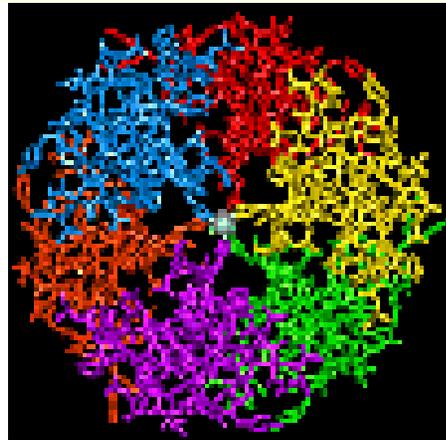
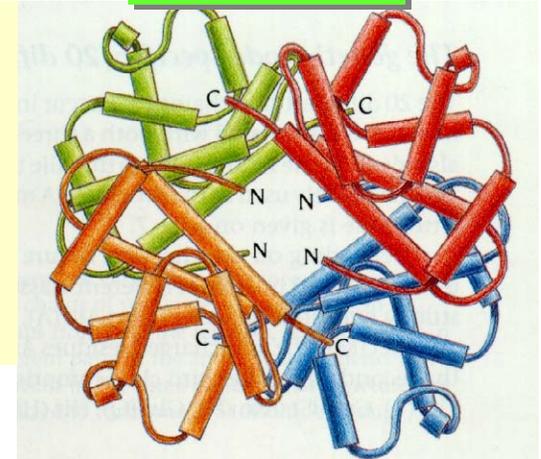


MDL

Quaternary Structures in Proteins

- The final structure may contain more than one “chain” arranged in a **quaternary structure**.

Quaternary



Insulin Hexamer

More on Secondary Structures

- **α -helix**

- Main chain with peptide bonds
- Side chains project outward from helix
- Stability provided by H-bonds between CO and NH groups of residues 4 locations away.

- **β -strand**

- Stability provided by H-bonds with one or more β -strands, forming β -sheets. Needs a β -turn.

Secondary Structure Prediction Software

254



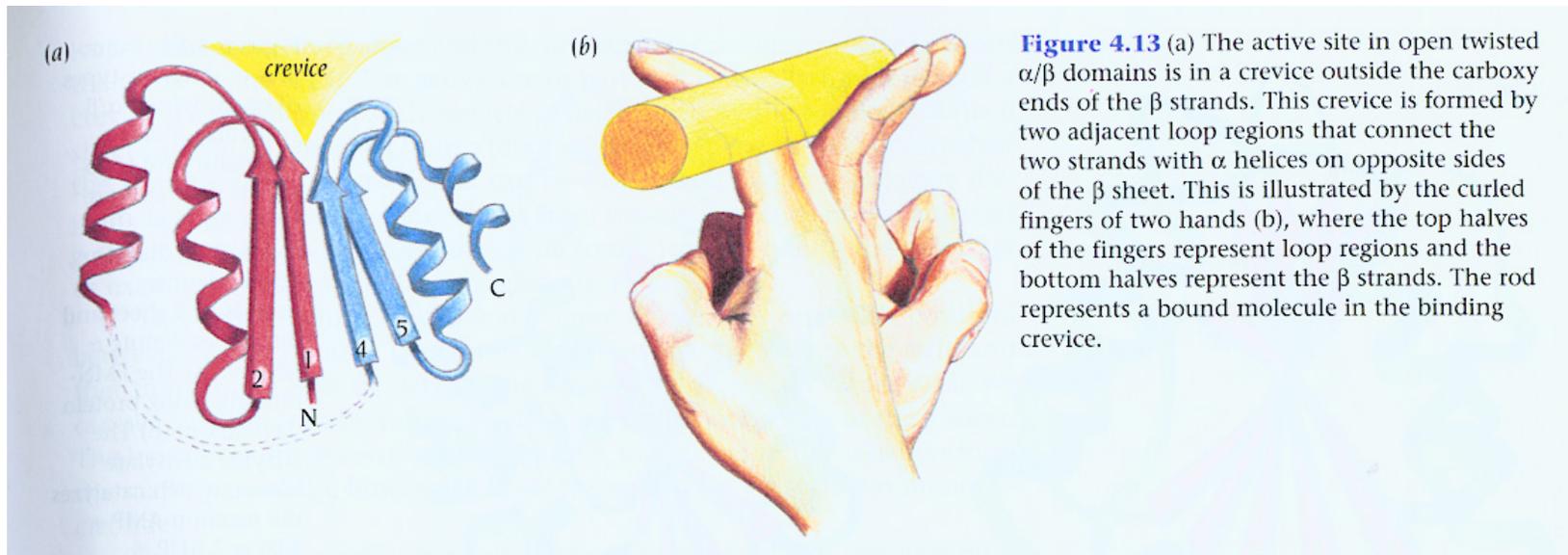
Figure 11.3 Comparison of secondary structure predictions by various methods. The sequence of flavodoxin, an α/β protein, was used as the query and is shown on the first line of the alignment. For each prediction, H denotes an α helix, E a β strand, T a β turn; all other positions are assumed to be random coil. Correctly assigned residues are shown in inverse type. The methods used are listed along the left side of the alignment and are described in the text. At the bottom of the figure is the secondary structure assignment given in the PDB file for flavodoxin (1OFV, Smith et al., 1983).

PDB: Protein Data Bank

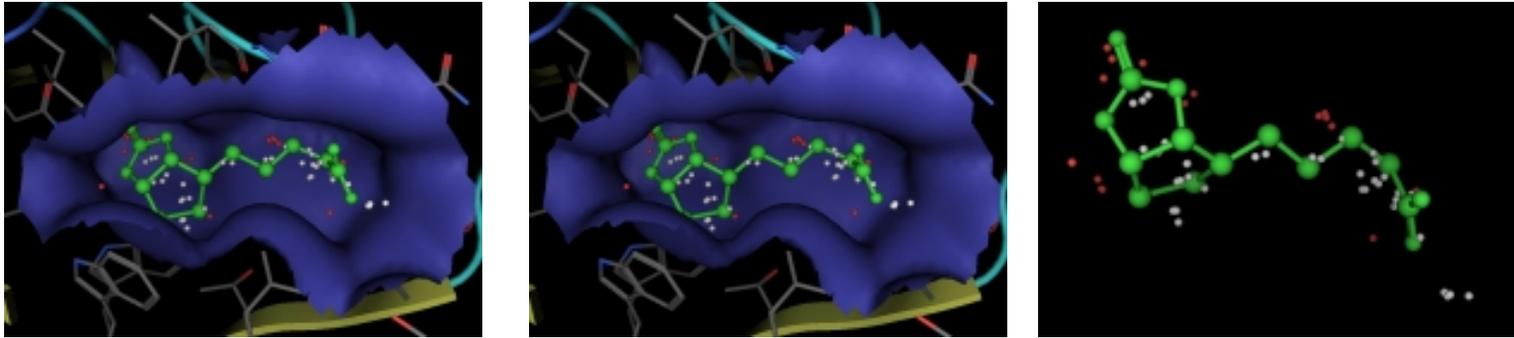
- Database of protein tertiary and quaternary structures and protein complexes.
<http://www.rcsb.org/pdb/>
- Over 29,000 structures as of Feb 1, 2005.
- Structures determined by
 - NMR Spectroscopy
 - X-ray crystallography
 - Computational prediction methods
- Sample PDB file: [Click here \[.\]](#)

Active Sites

Active sites in proteins are usually hydrophobic pockets/crevices/troughs that involve sidechain atoms.



Active Sites



Left PDB 3RTD (streptavidin) and the first site located by the MOE Site Finder. **Middle** 3RTD with complexed ligand (biotin). **Right** Biotin ligand overlaid with calculated alpha spheres of the first site.