# Protein Structures
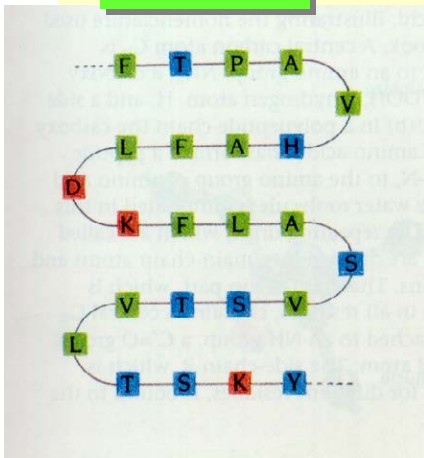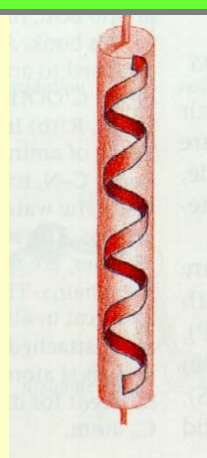
- Sequences of amino acid residues
- 20 different amino acids
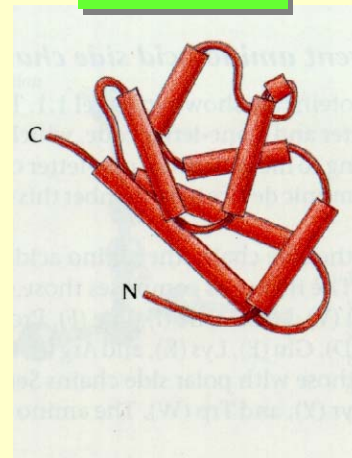
Primary

Secondary
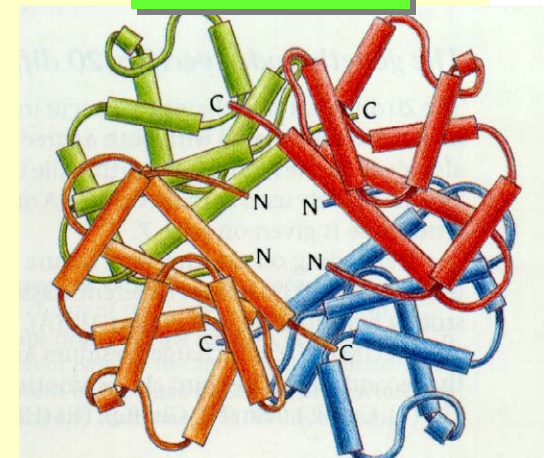
Tertiary

Quaternary

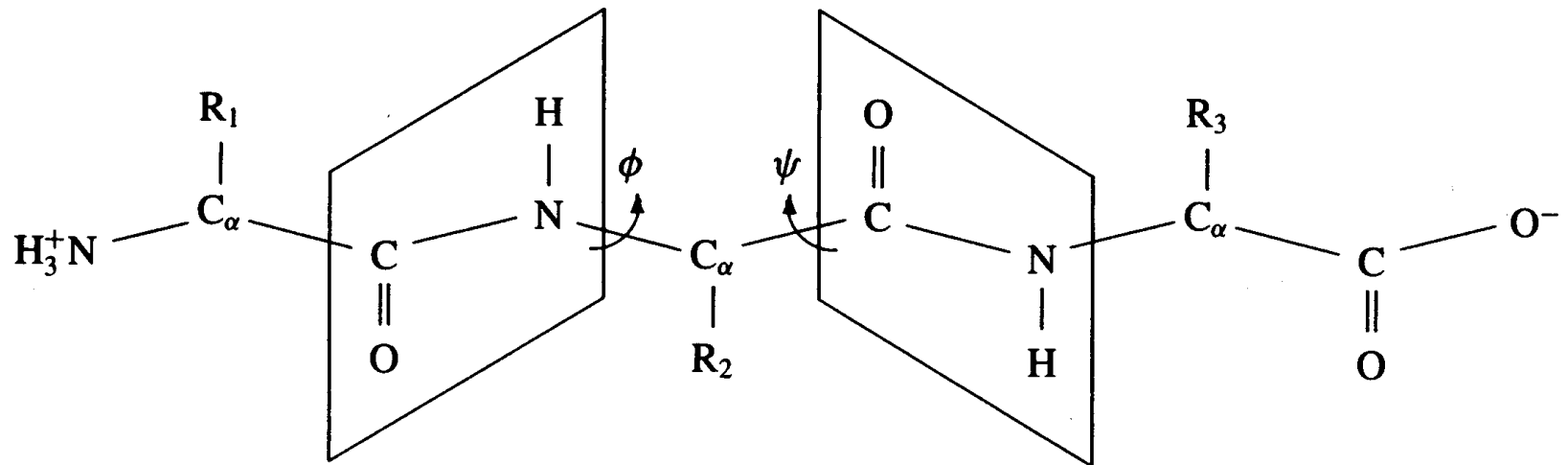# Angles $\phi$ and $\psi$ in the polypeptide chain



**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 $\phi$ and $\psi$.*
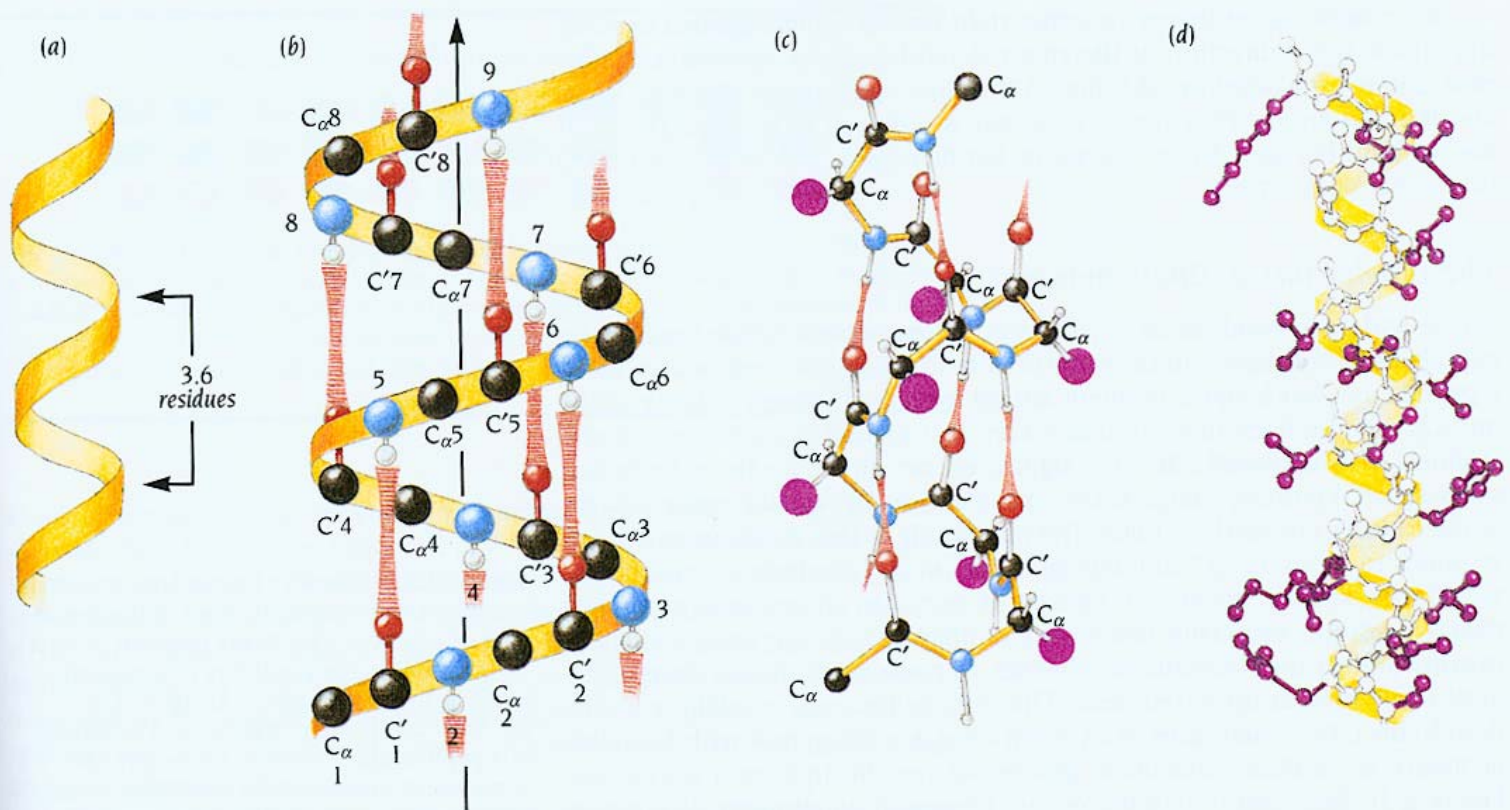
**Figure 2.2** The α helix is one of the major elements of secondary structure in proteins. Main-chain N and O atoms are hydrogen-bonded to each other within α helices. (a) Idealized diagram of the path of the main chain in an α helix. Alpha helices are frequently illustrated in this way. There are 3.6 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.

# More on Secondary Structures

- **$\alpha$-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.

- **$\beta$-strand**
  - Stability provided by H-bonds with one or more $\beta$-strands, forming $\beta$-sheets. Needs a $\beta$-turn.

# Secondary Structure Prediction Software



**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).

# Active Sites

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



**Figure 4.13** (a) The active site in open twisted α/β domains is in a crevice outside the carboxy ends of the β strands. This crevice is formed by two adjacent loop regions that connect the two strands with α helices on opposite sides of the β sheet. This is illustrated by the curled fingers of two hands (b), where the top halves of the fingers represent loop regions and the bottom halves represent the β strands. The rod represents a bound molecule in the binding crevice.

# 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.

# Motifs in Protein Sequences

**Motifs** are combinations of secondary structures in proteins with a specific **structure** and a specific **function**. They are also called **super-secondary structures**.

Examples: Helix-Turn-Helix, Zinc-finger, Homeobox domain, Hairpin-beta motif, Calcium-binding motif, Beta-alpha-beta motif, Coiled-coil motifs.

Several motifs may combine to form **domains**.
- Serine proteinase domain, Kringle domain, calcium-binding domain, homeobox domain.

# Motif Detection Problem

**Input:** Set, S, of known (aligned) examples of a motif M,
A new protein sequence, P.

**Output:** Does P have a copy of the motif M?

Example: Zinc Finger Motif
...**Y**KC**G**L**C**ERS**F**VEKSA**L**SR**H**ORV**H**KN...
     3     6                    19        23

**Input:** Database, D, of known protein sequences,
A new protein sequence, P.

**Output:** What interesting patterns from D
are present in P?

# Helix-Turn-Helix Motifs

- Structure
  - 3-helix complex
  - Length: 22 amino acids
  - Turn angle

- Function
  - Gene regulation by binding to DNA

# DNA Binding at HTH Motif



**Figure 7.10** The helix-turn-helix motif in lambda Cro bound to DNA (orange) with the two recognition helices (red) of the Cro dimer sitting in the major groove of DNA. The binding model, suggested by Brian Matthews, is shown schematically in (a) with connected circles for the Cα positions as they were model built into regular B-DNA. A schematic diagram of the Cro dimer is shown in (b) with different colors for the two subunits. A schematic space-filling model of the dimer of Cro bound to a bent B-DNA molecule is shown in (c). The sugar-phosphate backbone of DNA is red, and the bases are yellow. Protein atoms are colored red, blue, green, and white. [(a) Adapted from D. Ohlendorf et al., *J. Mol. Evol.* 19: 113, 1983. (c) Courtesy of Brian Matthews.]

# HTH Motifs: Examples

| Loc | Protein Name | Helix 2 | | | | | | | | | Turn | | | | Helix 3 | | | | | | | | |
|-----|------|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| | | -1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
| 14 | Cro | F | G | Q | E | K | T | A | K | D | L | G | V | Y | Q | S | A | I | N | K | A | I | H |
| 16 | 434 Cro | M | T | Q | T | E | L | A | T | K | A | G | V | K | Q | Q | S | I | Q | L | I | E | A |
| 11 | P22 Cro | G | T | Q | R | A | V | A | K | A | L | G | I | S | D | A | A | V | S | Q | W | K | E |
| 31 | Rep | L | S | Q | E | S | V | A | D | K | M | G | M | G | Q | S | G | V | G | A | L | F | N |
| 16 | 434 Rep | L | N | Q | A | E | L | A | Q | K | V | G | T | T | Q | Q | S | I | E | Q | L | E | N |
| 19 | P22 Rep | I | R | Q | A | A | L | G | K | M | V | G | V | S | N | V | A | I | S | Q | W | E | R |
| 24 | CII | L | G | T | E | K | T | A | E | A | V | G | V | D | K | S | Q | I | S | R | W | K | R |
| 4 | LacR | V | T | L | Y | D | V | A | E | Y | A | G | V | S | Y | Q | T | V | S | R | V | V | N |
| 167 | CAP | I | T | R | Q | E | I | G | Q | I | V | G | C | S | R | E | T | V | G | R | I | L | K |
| 66 | TrpR | M | S | Q | R | E | L | K | N | E | L | G | A | G | I | A | T | I | T | R | G | S | N |
| 22 | BlaA Pv | L | N | F | T | K | A | A | L | E | L | Y | V | T | Q | G | A | V | S | Q | Q | V | R |
| 23 | TrpI Ps | N | S | V | S | Q | A | A | E | Q | L | H | V | T | H | G | A | V | S | R | Q | L | K |

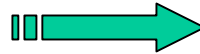# Basis for New Algorithm

- Combinations of residues in specific locations (may not be contiguous) contribute towards stabilizing a structure.

- Some reinforcing combinations are relatively rare.

# New Motif Detection Algorithm

Pattern Generation:

Aligned Motif Examples → Pattern Generator

Pattern Dictionary

Motif Detection:

New Protein Sequence → Motif Detector → Detection Results

# Patterns

| Loc | Protein Name | -1 | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | Helix 2 | | | | | Turn | | | | | | Helix 3 | | | | | |
| 14 | **Cro** | F | G | Q | E | K | T | A | K | D | L | G | V | Y | Q | S | A | I | N | K | A | I | H |
| 16 | **434 Cro** | M | T | Q | T | E | L | A | T | K | A | G | V | K | Q | Q | S | I | Q | L | I | E | A |
| 11 | **P22 Cro** | G | T | Q | R | A | V | A | K | A | L | G | I | S | D | A | A | V | S | Q | W | K | E |
| 31 | **Rep** | L | S | Q | E | S | V | A | D | K | M | G | M | G | Q | S | G | V | G | A | L | F | N |
| 16 | **434 Rep** | L | N | Q | A | E | L | A | Q | K | V | G | T | T | Q | Q | S | I | E | Q | L | E | N |
| 19 | **P22 Rep** | I | R | Q | A | A | L | G | K | M | V | G | V | S | N | V | A | I | S | Q | W | E | R |
| 24 | **CII** | L | G | T | E | K | T | A | E | A | V | G | V | D | K | S | Q | I | S | R | W | K | R |
| 4 | **LacR** | V | T | L | Y | D | V | A | E | Y | A | G | V | S | Y | Q | T | V | S | R | V | V | N |
| 167 | **CAP** | I | T | R | Q | E | I | G | Q | I | V | G | C | S | R | E | T | V | G | R | I | L | K |
| 66 | **TrpR** | M | S | Q | R | E | L | K | N | E | L | G | A | G | I | A | T | I | T | R | G | S | N |
| 22 | **BlaA Pv** | L | N | F | T | K | A | A | L | E | L | Y | V | T | Q | G | A | V | S | Q | Q | V | R |
| 23 | **TrpI Ps** | N | S | V | S | Q | A | A | E | Q | L | H | V | T | H | G | A | V | S | R | Q | L | K |

- Q1 G9 N20
- A5 G9 V10 I15
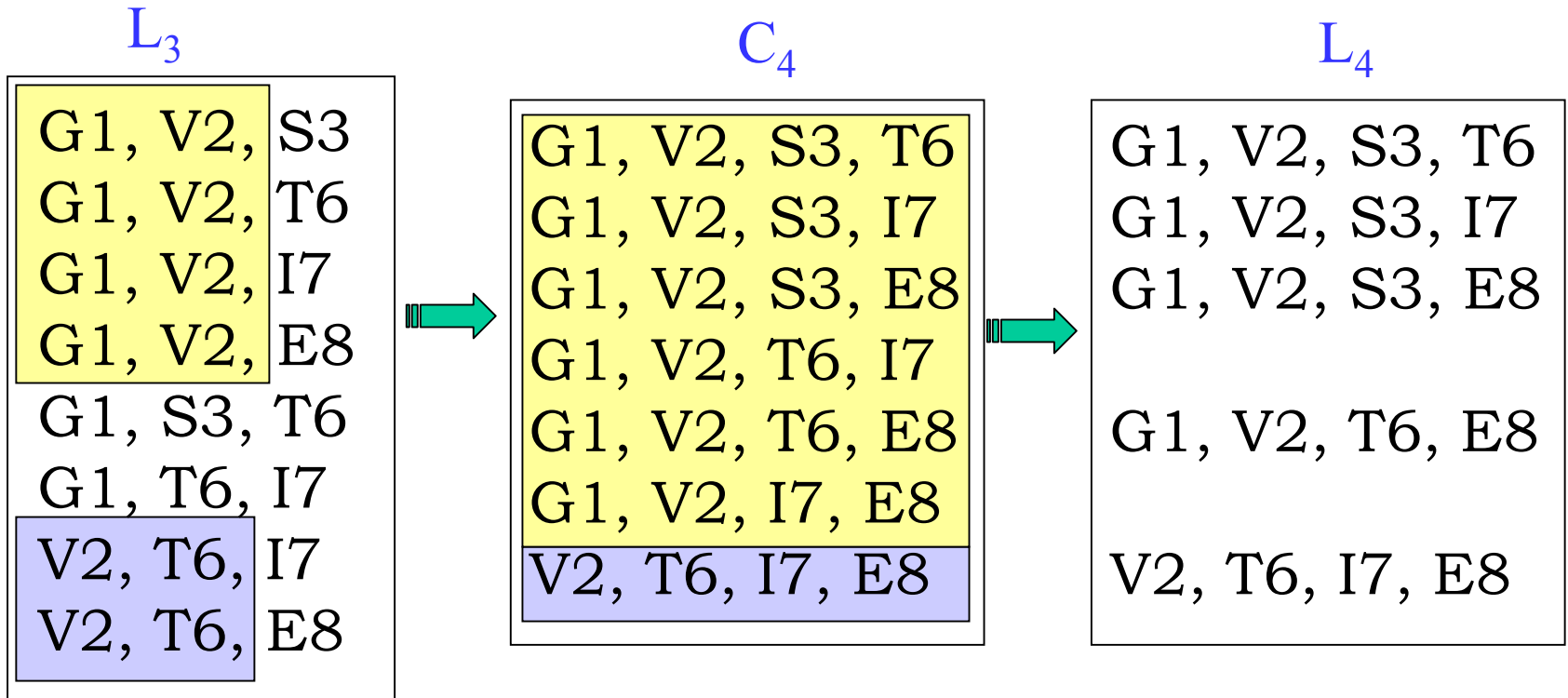
# Pattern Mining Algorithm

**Algorithm Pattern-Mining**
**Input**: Motif length $m$, support threshold $T$,
list of aligned motifs $M$.
**Output**: Dictionary $L$ of frequent patterns.

1. $L_1$ := All frequent patterns of length 1
2. **for** $i = 2$ **to** $m$ **do**
3. $C_i$ **:= Candidates**$(L_{i-1})$
4. $L_i$ := Frequent candidates from $C_i$
5. **if** $(|L_i| <= 1)$ **then**
6. **return** $L$ as the union of all $L_j$, $j <= i$.

# **Candidates** Function

$L_3$

| G1, V2, | S3 |
| G1, V2, | T6 |
| G1, V2, | I7 |
| G1, V2, | E8 |
| G1, S3, T6 |
| G1, T6, I7 |
| V2, T6, | I7 |
| V2, T6, | E8 |

$C_4$

G1, V2, S3, T6
G1, V2, S3, I7
G1, V2, S3, E8
G1, V2, T6, I7
G1, V2, T6, E8
G1, V2, I7, E8
V2, T6, I7, E8

$L_4$

G1, V2, S3, T6
G1, V2, S3, I7
G1, V2, S3, E8

G1, V2, T6, E8

V2, T6, I7, E8

# Motif Detection Algorithm

**Algorithm Motif-Detection**

**Input** :  Motif length $m$, threshold score $T$, pattern dictionary $L$, and input protein sequence $P[1..n]$.
**Output** : Information about motif(s) detected.

**1.** **for** each location i **do**
2.       S := **MatchScore**(P[i..i+m-1], L).
3.       **if**  (S > T) **then**
4.           Report it as a possible motif

# Experimental Results: **GYM 2.0**

| Motif | Protein Family | Number Tested | GYM = DE Agree | Number Annotated | GYM = Annot. |
|-------|----------------|---------------|----------------|------------------|--------------|
| HTH Motif (22) | Master | 88 | 88 (100 %) | 13 | 13 |
| | Sigma | 314 | 284 + 23 (98 %) | 96 | 82 |
| | Negates | 93 | 86 (92 %) | 0 | 0 |
| | LysR | 130 | 127 (98 %) | 95 | 93 |
| | AraC | 68 | 57 (84 %) | 41 | 34 |
| | Rreg | 116 | 99 (85 %) | 57 | 46 |
| | Total | 675 | 653 + 23 (94 %) | 289 | 255 (88 %) |

# Experiments

- Basic Implementation (Y. Gao)
- Improved implementation & comprehensive testing
                                        (K. Mathee, GN).
- Implementation for homeobox domain detection (X. Wang).
- Statistical methods to determine thresholds (C. Bu).
- Use of substitution matrix (C. Bu).
- Study of patterns causing errors (N. Xu).
- Negative training set (N. Xu).
- NN implementation & testing (J. Liu & X. He).
- HMM implementation & testing (J. Liu & X. He).