**Proteins: Evolution, and Analysis**
Lecture 7
9/15/2009

---

**(1)**



Amino acids with nonpolar side chains (G, A, V, L, I, M, P, F, W)

---

**(2)**



Amino acids with uncharged polar side chains (S, T, N, Q, Y, C)

---

**(3)**



Amino acids with charged polar side chains (K, R, H, D, E)

---

## The Fischer Convention

Absolute configuration about an asymmetric carbon

related to glyceraldehyde

(+) = D-Glyceraldehyde

(-) = L-Glyceraldehyde



L-Glyceraldehyde    L-α-Amino acid

**Geometric formulas**

Fischer projection

Mirror plane

L-Glyceraldehyde    D-Glyceraldehyde

---

## Cahn - Ingold - Prelog system

Can give absolute configuration nomenclature to multiple chiral centers.

Priority

Atoms of higher atomic number bonded to a chiral center are ranked above those of lower atomic number with lowest priority away from you R highest to lowest = clockwise, S highest to lowest = counterclockwise

$SH > OH > NH_2 > COOH > CHO > CH_2OH > C_6H_5 > CH_3 > H$

## Slide 1 (top left)

CHO
HO—C—H
CH₂OH

L-Glyceraldehyde

$CHO_{(X)}$
$H_{(Z)}$ —OH$_{(W)}$
$CH_2OH_{(Y)}$

(S)-Glyceraldehyde

$COO^-$
$H_3\overset{+}{N}$—C—H
$CH_3$

L-Alanine

$^-OOC_{(X)}$
$H_{(Z)}$ —$NH_3^{+}{}_{(W)}$
$H_3C_{(Y)}$

(S)-Alanine

## Slide 2 (top right)

- A *projection formula* representing the spatial arrangement of bonds on two adjacent atoms in a molecular entity.



- The structure appears as viewed along the bond between these two atoms, and the bonds from them to other groups are drawn as projections in the plane of the paper.
- The bonds from the atom nearer to the observer are drawn so as to meet at the centre of a circle representing that atom.
- Those from the further atom are drawn as if projecting from behind the circle.

## Slide 3 (middle left)

The major advantage of the CIP or RS system is that the chiralities of compounds with multiple asymmetric centers can be unambiguously described

H
$H_3C$ —— OH
$^-OOC$ —— $NH_3^+$
H

(2S, 3R)-Threonine

H
$H_3C$ —— $CH_2CH_3$
$^-OOC$ —— $NH_3^+$
H

(2S, 3S)-Isoleucine

## Slide 4 (middle right)

Structural Hierarchy in proteins



(a) – Lys – Ala – His – Gly – Lys – Lys – Val – Leu – Gly - Ala –
Primary structure (amino acid sequence in a polypeptide chain)

Secondary structure (helix)

Tertiary structure: one complete protein chain (β chain of hemoglobin)

Quaternary structure: the four separate chains of hemoglobin assembled into an oligomeric protein

Figure 6-1 Fundamentals of Biochemistry, 2/e

## Slide 5 (bottom left)

**Overview of Protein Sequencing**

**(1) Purify Protein**

Protein (two different polypeptide chains linked by disulfide bonds)

**(2) Determine number of PP**
End group analysis
(Dansyl chloride rxn)

Reduce disulfide bonds
Separate the chains

**(3) Fragment PP into smaller peptides**
Enzymes (Trypsin, Chymotrypsin, etc.)
Chemical (CNBr)

Use chemical or enzymatic methods to break each polypeptide into smaller peptides

Use different methods to generate a different set of peptide fragments

Determine the sequence of each peptide fragment

Determine the sequence of each peptide fragment

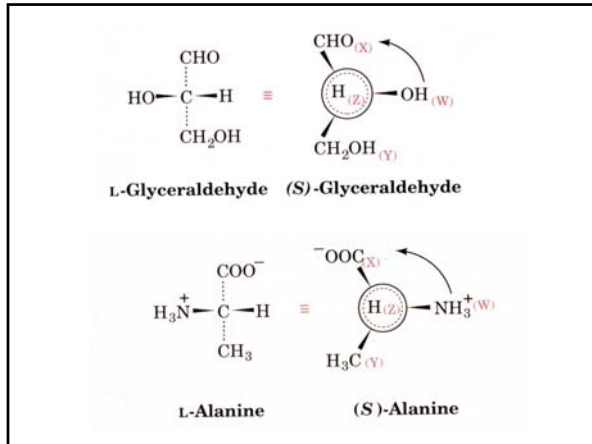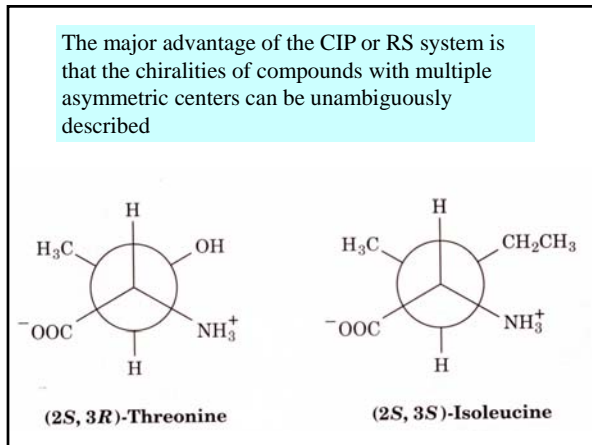**(4) Determine the sequence**
Edman Degradation with PITC

**(5) Assemble a sequence**
Use the two sets of overlapping peptide sequences to reconstruct the sequence of each polypeptide

Repeat fragmentation without breaking disulfide bonds to identify the Cys-containing sequences involved in the disulfide linkages

**(6) Elucidate S-S bonds**
Amino acid composition

## Slide 6 (bottom right)

**Long peptides have to be broken to shorter ones to be sequenced**



Table 5-3 Specificities of Various Endopeptidases

| Enzyme | Source | Specificity | Comments |
|---|---|---|---|
| Trypsin | Bovine pancreas | $R_{n-1}$ = positively charged residues: Arg, Lys; $R_n \neq$ Pro | Highly specific |
| Chymotrypsin | Bovine pancreas | $R_{n-1}$ = bulky hydrophobic residues: Phe, Trp, Tyr; $R_n \neq$ Pro | Cleaves more slowly for $R_{n-1}$ = Asn, His, Met, Leu |
| Elastase | Bovine pancreas | $R_{n-1}$ = small neutral residues: Ala, Gly, Ser, Val; $R_n \neq$ Pro | |
| Thermolysin | Bacillus thermoproteolyticus | $R_n$ = Ile, Met, Phe, Trp, Tyr, Val; $R_{n-1} \neq$ Pro | Occasionally cleaves at $R_n$ = Ala, Asp, His, Thr; heat stable |
| Pepsin | Bovine gastric mucosa | $R_n$ = Leu, Phe, Trp, Tyr; $R_{n-1} \neq$ Pro | Also others; quite nonspecific; pH optimum = 2 |
| Endopeptidase V8 | Staphylococcus aureus | $R_{n-1}$ = Glu | |

Table 5-3 Fundamentals of Biochemistry, 2/e
© 2006 John Wiley & Sons

Table 5-3 Specificities of Various Endopeptidases

| Enzyme | Source | Specificity | Comments |
|---|---|---|---|
| Trypsin | Bovine pancreas | $R_{n-1}$ = positively charged residues: Arg, Lys; $R_n \neq$ Pro | Highly specific |
| Chymotrypsin | Bovine pancreas | $R_{n-1}$ = bulky hydrophobic residues: Phe, Trp, Tyr; $R_n \neq$ Pro | Cleaves more slowly for $R_{n-1}$ = Asn, His, Met, Leu |
| Elastase | Bovine pancreas | $R_{n-1}$ = small neutral residues: Ala, Gly, Ser, Val; $R_n \neq$ Pro | |
| Thermolysin | Bacillus thermoproteolyticus | $R_n$ = Ile, Met, Phe, Trp, Tyr, Val; $R_{n-1} \neq$ Pro | Occasionally cleaves at $R_n$ = Ala, Asp, His, Thr; heat stable |
| Pepsin | Bovine gastric mucosa | $R_n$ = Leu, Phe, Trp, Tyr; $R_{n-1} \neq$ Pro | Also others; quite nonspecific; pH optimum = 2 |
| Endopeptidase V8 | Staphylococcus aureus | $R_{n-1}$ = Glu | |

**Q9**. You must cleave the following peptide into smaller fragments. Which of the proteases listed in the table would be likely to yield the most fragments? The fewest?

NMTQGRCKPVNTFVHEPLVDVQNVCFKE

---

## Cyanogen bromide cleavage of a polypeptide



Figure 5-16 Fundamentals of Biochemistry, 2/e
© 2006 John Wiley & Sons

---

### Reconstructing the protein's sequence

#### Specific chemical cleavage reagents

Cleave the large protein using i.e. trypsin, separate fragments and sequence all of them. (We do not know the order of the fragments!!)

Cleave with a different reagent i.e. Cyanogen Bromide, separate the fragments and sequence all of them. Align the fragments with overlapping sequence to get the overall sequence.



Phe—Trp—Met—Gly—Ala—Lys—Leu—Pro—Met—Asp—Gly—Arg—Cys—Ala—Gln

Figure 5-20 Fundamentals of Biochemistry, 2/e
© 2006 John Wiley & Sons

---



Polypeptide fragment containing disulfide bond

Reduce disulfide and block with iodoacetate

Separate and sequence the polypeptides

Figure 5-21 Fundamentals of Biochemistry, 2/e
© 2006 John Wiley & Sons

Determining the positions of disulfide bond

---

### How to assemble a protein sequence

1. Write a blank line for each amino acid in the sequence starting with the N-terminus.

2. Follow logically each clue and fill in the blanks.

3. Identify overlapping fragments and place in sequence blanks accordingly.

4. Make sure logically all your amino acids fit into the logical design of the experiment.

5. Double check your work.

---

|  | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

$H_3N^+$—_ _ _ _ _ _ _ _ _ _ _ _ _ _ _—COO⁻

A - **T**

F - M - A - **T**

A - K - F - M

Q - M - A - K

D - I - K - Q - M

G - M - D - I - K

Y - R - G - M

Y - R

| Cyanogen Bromide (CNBr) Cleaves after Met i.e M - X | Trypsin cleaves after K or R (positively charged amino acids) |
|---|---|
| D - I - K - Q - M | Q - M - A - K |
| A - **T** | G - M - D - I - K |
| A - K - F - M | F - M - A - **T** |
| Y - R - G - M | Y - R |

3

**Q11**. Separate cleavage reactions of a polypeptide by CNBr and chymotrypsin yield fragments with the following amino acid sequences. What is the the sequence of the intact polypeptide?

| CNBr treatment | Chymotrypsin |
| --- | --- |
| 1. Arg-Ala-Tyr-Gly-Asn | 1. Met-Arg-Ala-Tyr |
| 2. Leu-Phe-Met | 2. Asp-Met-Leu-Phe |
| 3. Asp-Met | 3. Gly-Asn |

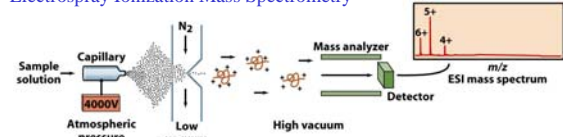**Q13.** Treatment of a polypeptide with 2-mercaptoethanol yields two PP:
1. Ala-Val-Cys-Arg-Thr-Gly-Cys-Lys-Asn-Phe-Leu
2. Tyr-Lys-Cys-Phe-Arg-His-Thr-Lys-Cys-Ser

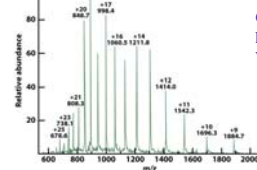Treatment of the intact PP with trypsin yields fragments with the following aa compositions:

| | |
| --- | --- |
| 3. (Ala, Arg, Cys2, Ser, Val) | 4. (Arg, Cys2, Gly, Lys, Thr, Phe) |
| 5. (Asn, Leu, Phe) | 6. (His, Lys, Thr) |
| 7. (Lys, Tyr) | |

---

## Sequencing by Mass Spectrometry

Electrospray Ionization Mass Spectrometry



ESI-MS spectrum of horse heart apomyoglobin

**Q**. Two successive peaks in the mass spectrum have measure m/z ratios of 1414.0 and 1542.3. What is the original apomyoglobin molecule?

$p1 = (M+z)/z$
$p2 = (M+z-1)/z-1$

M= 16,975D (16,951 D in table 5-1)

---



Tandem Mass Spectrometry in amino acid sequencing

---



---

## Protein Evolution

### Species variation in homologous proteins

The primary structures of a given protein from related species closely resemble one another. If one assumes, according to evolutionary theory, that related species have evolved from a common ancestor, it follows that each of their proteins must have likewise evolved from the corresponding ancestor.

A protein that is well adapted to its function, that is, one that is not subject to significant physiological improvement, nevertheless continues to evolve.

**Neutral drift**: changes not effecting function

---

### Homologous proteins

(evolutionarily related proteins)

Compare protein sequences:

Conserved residues, i.e invariant residues reflect chemical necessities.

Conserved substitutions, substitutions with similar chemical properties (Asp for Glu), (Lys for Arg), (Ile for Val)

Variable regions, no requirement for chemical reactions etc.

**Amino acid difference matrix for 26 species of cytochrome c**

```
Man,chimp       0
Rhesus monkey   1   0                        Average differences
Horse          12  11   0
Donkey         11  10   1   0                        10.0
cow,sheep      10   9   3   2   0
dog            11  10   6   5   3   0
gray whale     10   9   5   4  *2   3   0            5.1
rabbit          9   8   6   5   4   5   2   0
kangaroo       10  11   7   8   6   7   6   6   0
Chicken        13  12  11  10   9  10   9   8  12   0
penguin        13  12  12  11  10  10   9   8  10   2   0     9.9
Duck           11  10  10   9   8   7   6   5  11   3   3   0
Rattlesnake    14  15  22  21  20  21  19  18  21  19  20  17   0   12.6
turtle         15  14  11  10   9   9   8   9  11   8   8   7  22   0
Bullfrog       18  17  14  13  11  12  11  11  13  11  12  11  24  10   0
Tuna fish      21  21  19  18  17  18  17  18  17  18  17  26  18  15   0   18.5
worm fly       27  26  22  22  22  21  22  21  24  23  24  22  29  24  22  24   0
silk moth      31  30  29  28  27  25  27  26  28  28  27  27  31  28  29  32  14   0   25.9
Wheat          43  43  46  45  45  44  44  44  47  46  46  46  46  48  49  45  45   0
Bread mold     48  47  46  46  46  46  46  49  47  48  48  48  48  41  47  54   0   47.0
 Yeast         45  45  46  45  45  45  45  46  46  45  46  47  49  47  47  45  47  41   0
Candida k.     51  51  51  50  50  49  50  50  51  51  50  51  51  53  51  48  47  47  50  42  27   0
```

Man,chimp / monkey / Horse / Donkey / cow,sheep / dog / gray whale / rabbit / kangaroo / Chicken, / penguin / Duck / Rattlesnake / turtle / Bullfrog / Tuna fish / worm fly / silkworm / Wheat / Bread mold / Yeast / Candida

---



---



---

## Phylogenetic tree

- Indicates the ancestral relationships among the organisms that produced the protein.
- Each branch point indicates a common ancestor.
- Relative evolutionary distances between neighboring branch points are expressed as the number of amino acid differences per 100 residues of the protein.

  PAM units

  or

  **P**ercentage of **A**ccepted **M**utations



---



PAM values differ for different proteins.

★ Although DNA mutates at a assumed constant rate. Some proteins cannot accept mutations because the mutations kill the function of the protein and thus are not viable. ★

---

Mutation rates appear constant in time

Although insects have shorter generation times that mammals and many more numbers of replication, number of mutations appear to be independent of the number of generations but dependent upon time

Cytochrome c amino acid differences between mammals, insects and plants note the similar distances
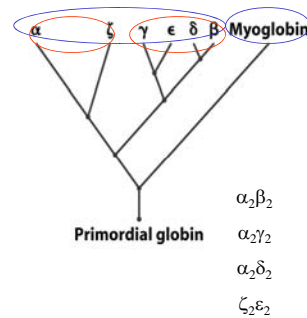
## Evolution through gene duplication

Many proteins within an organism have sequence similarities with other proteins.

• These are called gene or protein families.

• The relatedness among members of a family can vary greatly.

• These families arise by gene duplication.

• Once duplicated, individual genes can mutate into separate genes.

• Duplicated genes may vary in their chemical properties due to mutations.

• These duplicate genes evolve with different properties.

• Example the globin family.

## Genealogy of the globin family



**Hemoglobin:**
• is an oxygen transport protein
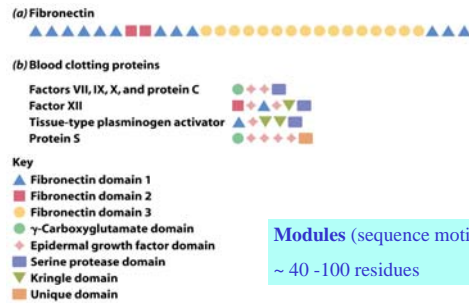• it must bind and release oxygen as the cells require oxygen

**Myoglobin:**
• is an oxygen storage protein
• it binds oxygen tightly and releases it when oxygen concentrations are very low

$\alpha_2\beta_2$
$\alpha_2\gamma_2$
$\alpha_2\delta_2$
$\zeta_2\epsilon_2$

## The globin family history

1. Primordial globin gene acted as an Oxygen-storage protein.

2. Duplication occurred 1.1 billion years ago.

   lower oxygen-binding affinity, monomeric protein.

3. Developed a tetrameric structure two $\alpha$ and two $\beta$ chains increased oxygen transport capabilities. ($\alpha_2\beta_2$).

4. Mammals have fetal hemoglobin with a variant $\beta$ chain i.e. $\gamma$ ($\alpha_2\gamma_2$).

5. Human embryos contain another hemoglobin ($\zeta_2\epsilon_2$).

6. Primates also have a $\delta$ chain with no known unique function. ($\alpha_2\delta_2$).

## Modular Construction of some proteins



(a) Fibronectin

(b) Blood clotting proteins

Factors VII, IX, X, and protein C
Factor XII
Tissue-type plasminogen activator
Protein S

Key
▲ Fibronectin domain 1
■ Fibronectin domain 2
● Fibronectin domain 3
● γ-Carboxyglutamate domain
◆ Epidermal growth factor domain
■ Serine protease domain
▼ Kringle domain
■ Unique domain

**Modules** (sequence motifs):

~ 40 -100 residues

# Lecture 8

(9/17/2009)

Chapter 6 - Proteins: 3-D structure
6-1. Secondary Structure