

Protein evolution

Chapter 3

Proteins

TE Creighton

Evolution

- In nature the macroscopic diversity we see stems from a molecular unity.
 - All organisms use the same 20 aminoacids
 - All organisms use the same nucleotides for DNA, RNA
 - The aminoacids are all L amino acids
 - Ribose units of DNA and RNA are all D isomers.
 - The basic properties of the proteins in E.coli and Humans are the same.
 - The function of DNA and RNA pertains, but the differences between prokaryotes and eukaryotes becomes significant

Evolution

- The biochemical similarities extend to even higher levels of organisation
 - Biochemical pathways
 - Metabolism
- All this can be explained so far, with the presence of a common ancestor for all living organisms.
- This cell had already acquired all basic biochemical features.

Evolution

- The diversity at higher level of organisation has arisen
- From Darwinian divergence
- But the most basic functions of an organism have been conserved.
- Reason for this: any alteration in these functions
- LETHAL !!!!!!!

Evolution

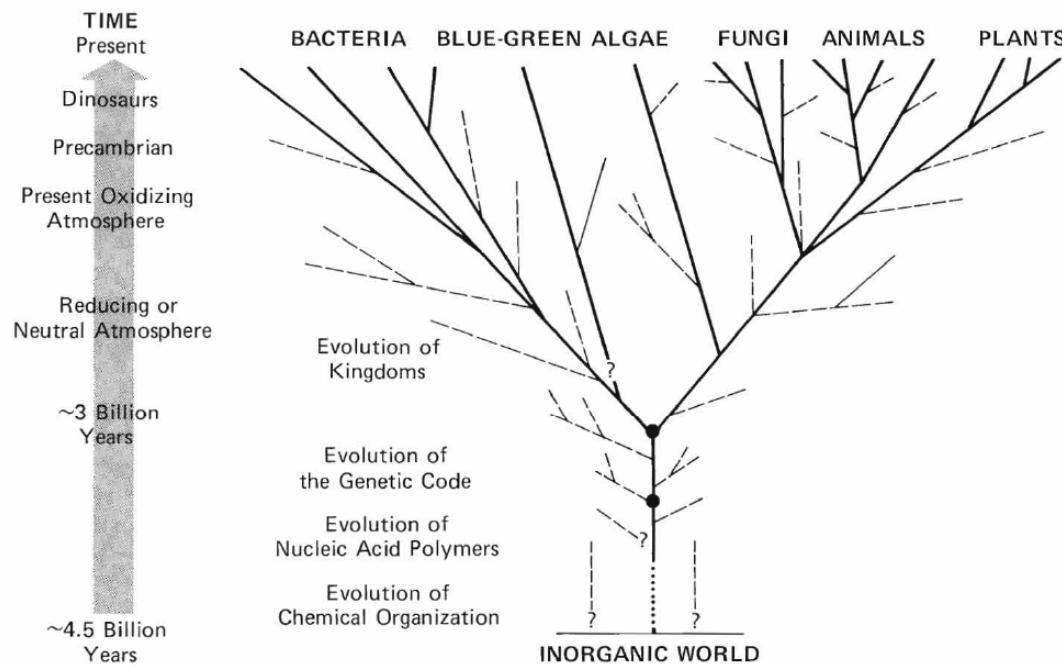


FIGURE 3.1

General phylogenetic tree of how life on earth might have evolved. The dashed lines represent hypothetical lineages that have not survived to the present day. (Adapted from M. O. Dayhoff, *Atlas of Protein Sequence and Structure*, National Biomedical Research Foundation, Washington, D.C., 1972.)

Evolution

- Comparisons at the macroscale:
- comparisons of the organisms,
- the fossil record
- But even more compelling is the comparison on the protein level.
- These studies provide deeper insights into how evolution works

Evolution

- E.g. A DNA chain of only 1000 nucleotides
- Can exist with any one of 4^{1000} different nucleotide sequences.
- A small protein of 100 aminoacids can have 20^{100} different sequences.
- Compared to the mere numbers of possible proteins the existing proteins are determined by their ancestors
- An existing protein serves its function. The variation is within the limit of function.'

Primordial origins of life

- Life started out of a primordial soup
- Where all components have been made out of simple molecules
- Forming the primordial soup
- All the processes have been demonstrated in the laboratory
- The chicken-egg problem of molecular biology:
- Who comes first DNA, RNA or protein

The accepted picture

- Since the finding of Ribozymes:
- First RNA and Ribozymes
- Thereafter proteins
- And DNA since DNA is chemically more stable as RNA
- Remainders of this RNA world are:
- mRNA
- T-RNA
- Ribosomes (hybrids RNA-Protein)
- Coenzymes: ATP, NAD⁺, FAD composed of RNA moieties

Evolutionary divergence of proteins

- Similar sequences mean descent from a common ancestor.
- The number of possible RNA,DNA or protein sequences is so great that it is implausible that similar long sequences could have arisen by any mechanism other than divergence from the same ancestor.
- Proteins and nucleic acids that have evolved from a common ancestor are homologous.
- After the first replication, mutations have started in the two sequences

Homologous genes and proteins

- Proteins and genes are either Homologous or NOT.
- Either or not from a common ancestor.

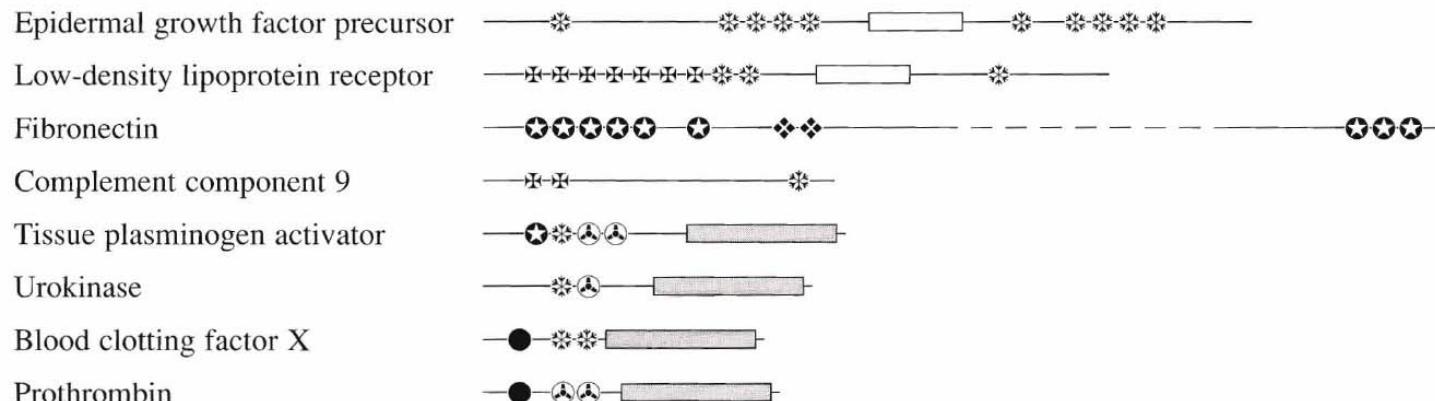


FIGURE 3.16

Mosaic structures of some vertebrate proteins. Each symbol represents members of a set of homologous segments of polypeptide chain, which in these examples are also independent structural units, or domains, of the protein. Those with distinct structural or functional properties are the serine protease domains (shaded rectangles), Ca^{2+} -binding domains containing γ -carboxy-Glu residues (●), so-called kringle domains (○), and epidermal growth factor (EGF) domains (※).

Homologous Genes and Proteins

- The other explanation Than divergence is
- Convergence
- Convergence means two different sequences have acquired the same function ality and share a certain degree of homology.
- It exists at the macroscopic level but there are no instances where it could have been shown conclusively for nucleic acids and proteins.
- But it can not be dismissed

Sequence alignments

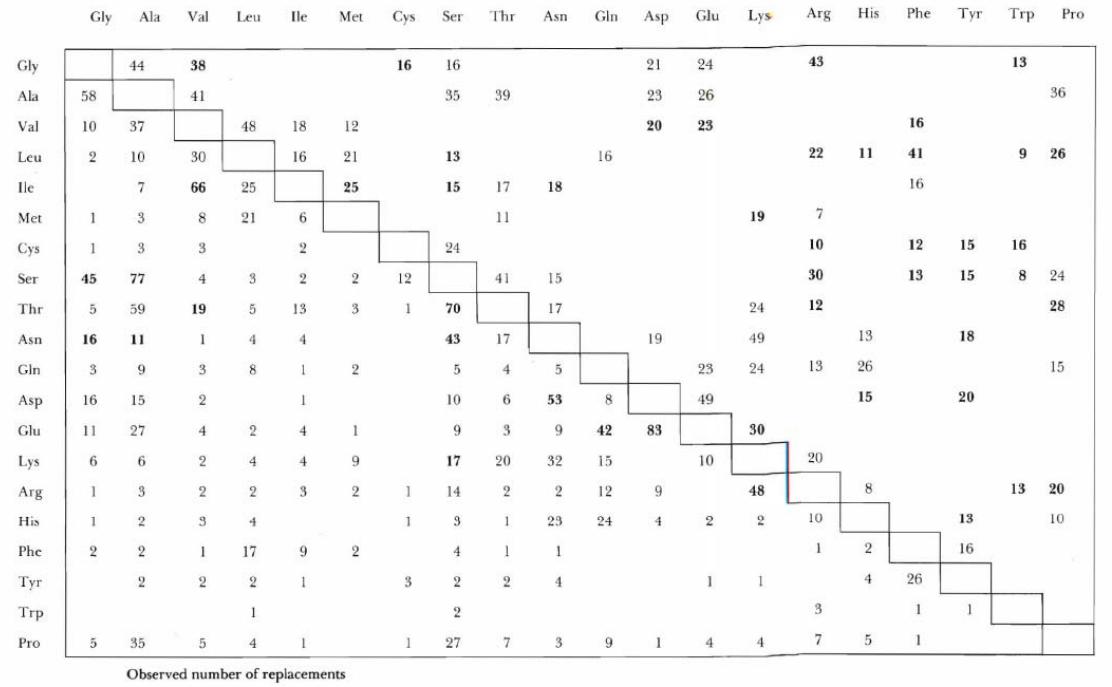
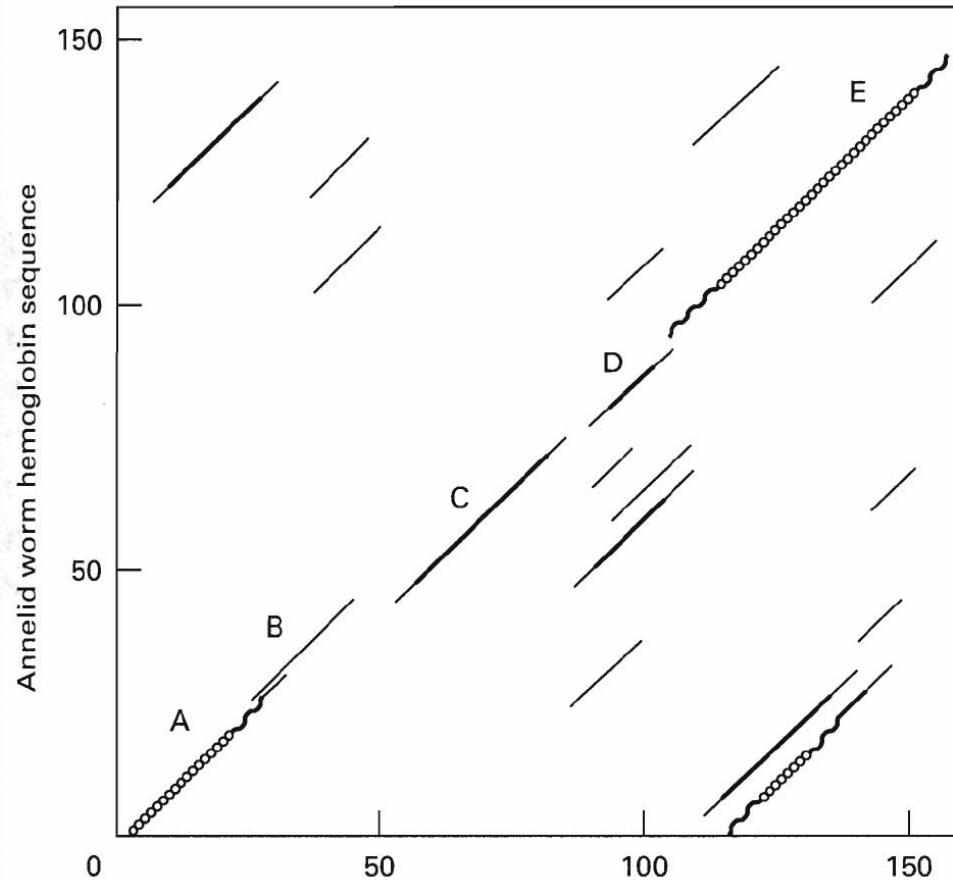


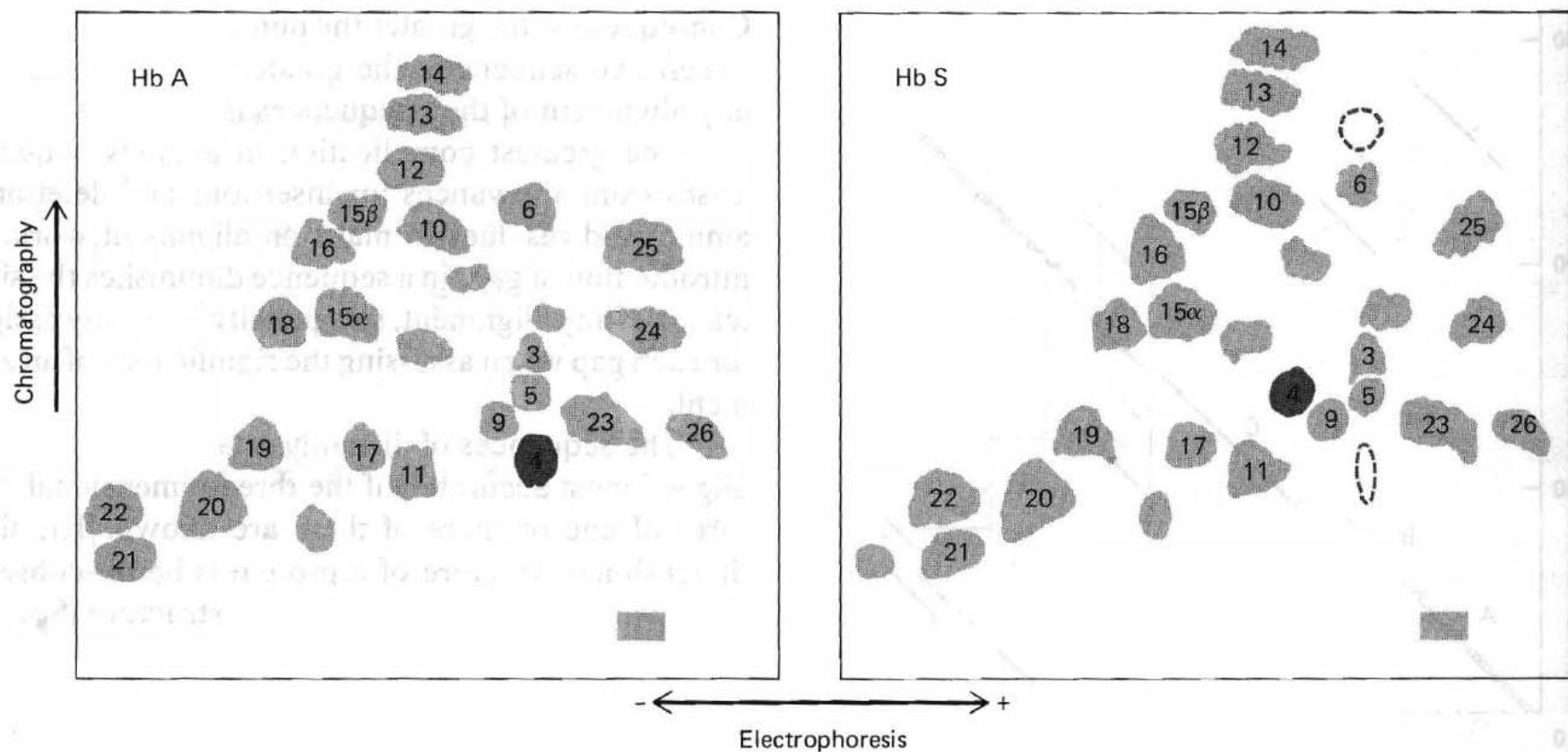
FIGURE 3.2

Relative frequencies of amino acid replacements in a total of 1572 examples between closely related proteins that are observed (bottom, left) and those expected for random single-nucleotide mutations (top, right). The greatest discrepancies between the observed and random replacements are shown in boldface type. Replacements involving chemically similar amino acids are generally observed to be much more frequent than expected for random mutations. (Observed replacements compiled by M. O. Dayhoff, *Atlas of Protein Sequence and Structure*, vol. 5, suppl. 3, National Biomedical Research Foundation, Washington, D.C., 1978.)

Diagonal dot blot



Mutation and Protein structure



Mutation and protein structure

Cranston	-Lys -Ser -Ile -Thr -Lys -Leu -Ala -Phe -Leu -Leu -Ser -Asn -Phe -Tyr -terminate AAG AGU AUC ACU AAG CUC GCU UUC UUG CUG UCC AAU UUC UAU UAA
Normal	
Tak	-Lys -Tyr -His -terminate AAG UAU CAC UAA GCU CGC UUU CUU GCU GUC CAA UUU CUA UUA A
	
	-Lys -Tyr -His -Thr -Lys -Leu -Ala -Phe -Leu -Leu -Ser -Asn -Phe -Tyr -terminate AAG UAU CAC ACU AAG CUC GCU UUC UUG CUG UCC AAU UUC UAU UAA

FIGURE 3.5

Insertion mutations that alter the reading frame of the gene for the β chain of hemoglobin. The normal sequences of mRNA and of the protein product are shown, as well as the sequences for two mutants, *Cranston* and *Tak*. The nucleotides inserted in the mutants are in boldface type; note that both insertions apparently occurred by duplication of the two preceding nucleotides. (Sequences from N. Proudfoot and G. Brownlee, *Brit. Med. Bull.* 32:251–256, 1976.)

Reconstructing evolution

	1	10	20	30	40	48	
Human and chimpanzee	G D V E K G K K I F I M K C S Q C H T V E K G G K H K T G P N L H G L F G R K T G Q A P G Y S Y						
Pig, bovine, and sheep	G D V E K G K K I F V Q K C A Q C H T V E K G G K H K T G P N L H G L F G R K T G Q A P G F S Y						
Gray kangaroo	G D V E K G K K I F V Q K C A Q C H T V E K G G K H K T G P N L N G I F G R K T G Q A P G F T Y						
Chicken and turkey	G D I E K G K K I F V Q K C S Q C H T V E K G G K H K T G P N L H G L F G R K T G Q A E G F S Y						
Snapping turtle	G D V E K G K K I F V Q K C A Q C H T V E K G G K H K T G P N L N G L I G R K T G Q A E G F S Y						
Puget Sound dogfish	G D V E K G K K I F V Q K C A Q C H T V E N G G K H K T G P N L S G L F G R K T G Q A E G F S Y						
Pacific lamprey	G D V E K G K K V F V Q K C S Q C H T V E K A G K H K T G P N L S G L F G R K T G Q A P G F S Y						
Garden snail	G Z A Z K G K K I F T Q K C L Q C H T V E A G G K H K T G P N L S G L F G R K Q G Q A P G F A Y						
Screw-worm fly	G V P A G D V E K G K K I F V Q R C A Q C H T V E A G G K H K V G P N L H G L F G R K T G Q A A G F A Y						
Tobacco hornworm moth	G V P A G N A D N G K K I F V Q R C A Q C H T V E A G G K H K V G P N L H G F F G R K T G Q A P G F S Y						
<i>Candida krusei</i>	P A P F E Q G S A K K G A T L F K T R C A Q C H T I E A G G P H K V G P N L H G I F S R H S G Q A E G Y S Y						
Rust fungus	G F E D G D A K K G A R I F K T R C A Q C H T L G A G E P N K V G P N L H G L F G R R S G T V E G F S Y						
Rape and cauliflower	A S F D E A P P G N S K A G E K I F K T K C A Q C H T V D K G A G H K Q G P N L N G L F G R Q S G T T A G Y S Y						
	49	60	70	80	90	100	104
Human and chimpanzee	T A A N K N K G I I W G E D T L M E Y L E N P K K Y I P G T K M I F V G I K K K E E R A D L I A Y L K K A T N E						
Pig, bovine, and sheep	T D A N K N K G I T W G E E T L M E Y L E N P K K Y I P G T K M I F A G I K K K G E R E D L I A Y L K K A T N E						
Gray kangaroo	T D A N K N K G I I W G E D T L M E Y L E N P K K Y I P G T K M I F A G I K K K G E R A D L I A Y L K K A T N E						
Chicken and turkey	T D A N K N K G I T W G E D T L M E Y L E N P K K Y I P G T K M I F A G I K K K G E R A D L I A Y L K K A T N E						
Snapping turtle	T E A N K N K G I T W G E E T L M E Y L E N P K K Y I P G T K M I F A G I K K K G E R A D L I A Y L K K A T N E						
Puget Sound dogfish	T D A N K S K G I T W Q Q E T L R I Y L E N P K K Y I P G T K M I F A G L K K K S E R Q D L I A Y L K K T A A S						
Pacific lamprey	T D A N K S K G I V W N Q E T L F V Y L E N P K K Y I P G T K M I F A G I K K E G E R K D L I A Y L K K S T S E						
Garden snail	T D A N K G K G I T W K N Q T L F E Y L E N P K K Y I P G T K M V F A G L K B Z T E R V D L I A Y L Z Z A T K K						
Screw-worm fly	T N A N K A K G I T W Q D D T L F E Y L E N P K K Y I P G T K M I F A G L K K P N E R G D L I A Y L K S A T K						
Tobacco hornworm moth	S N A N K A K G I T W Q D D T L F E Y L E N P K K Y I P G T K M V F A G L K K A N E R A D L I A Y L K Q A T K						
<i>Candida krusei</i>	T D A N K R A G V E W A E P T M S D Y L E N P K K Y I P G T K M A F G G L K K A K D R N D L V T Y M L E A S K						
Rust fungus	T D A N K K A G Q V W E E T F L E Y L E N P K K Y I P G T K M A F G G L K K E K D R N D L V T Y L R E E T K						
Rape and cauliflower	S A A N K N K A V E W E E K T L Y D Y L L N P K K Y I P G T K M V F P G L K K P Q D R A D L I A Y L K E A T A						

FIGURE 3.6

Amino acid sequences of some cytochromes c from various eukaryotes. The one-letter code (see Table 1.1) is used; Z is either E or Q (Glu or Gln), B is either D or N (Asp or Asn). Residues that are identical to those in the human protein are boldface; chemically similar residues are in italics.

Reconstructing evolution

	Chimpanzee	Sheep	Rattlesnake	Carp	Snail	Moth	Yeast	Cauliflower	Parsnip
Human	0	10	14	18	29	31	44	44	43
Chimpanzee		10	14	18	29	31	44	44	43
Sheep			20	11	24	27	44	46	46
Rattlesnake				26	28	33	47	45	43
Carp					26	26	44	47	46
Garden snail						28	48	51	50
Tobacco hornworm moth							44	44	41
Baker's yeast (iso-1)								47	47
Cauliflower									13

FIGURE 3.7

Amino acid difference matrix for cytochromes *c* of a few representative species. These proteins all consist of at least 104 amino acid residues, so at least half the residues of each pair are identical.

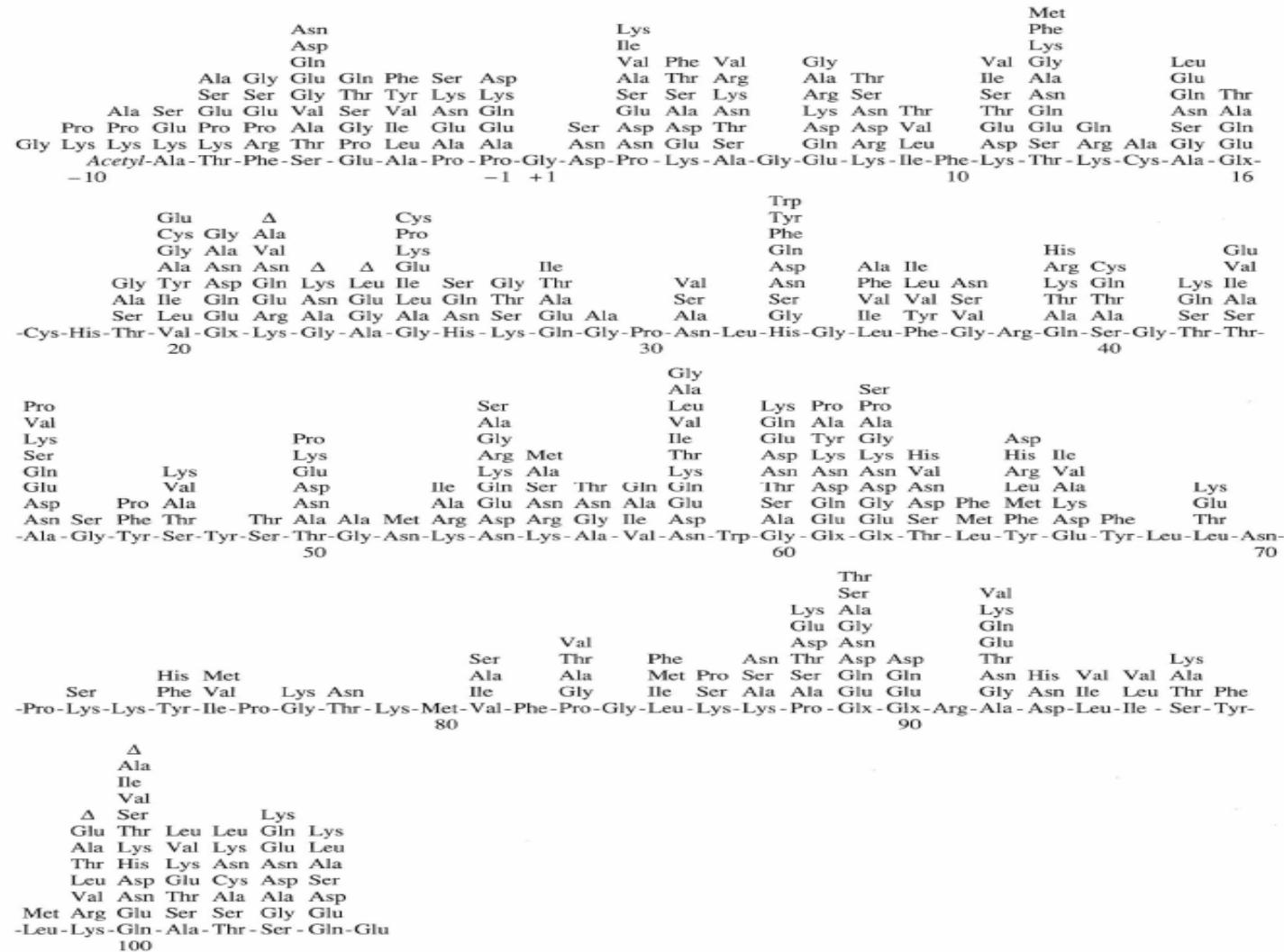


FIGURE 3.8

Composite amino acid sequence of cytochromes *c* from 92 eukaryotic species. The continuous sequence of 113 residues is the longest cytochrome *c*, that of *Ginkgo bilboa*. The numbering of the residues is that of the mammalian proteins, which start at position 1 and end at 104. The proteins included here start at positions -11 to +1 and end at positions 101 to 105. Deletions of residues are indicated by Δ. (The sequences were obtained from M. O. Dayhoff, *Atlas of Protein Sequences and Structure*, National Biomedical Research Foundation, Washington, D.C., 1972; and from D. M. Hampsey et al., *J. Biol. Chem.* 261:3259–3271, 1986.)

Reconstructing evolution

Table 3.1 Relative Variabilities of Amino Acid Residues during Divergence of Homologous Proteins^a

Residue	Variability	Residue	Variability	Residue	Variability
Asn	100	Met	70	Gly	37
Ser	90	Gln	69	Tyr	31
Asp	79	Val	55	Phe	31
Glu	76	His	49	Leu	30
Ala	75	Arg	49	Cys	15
Thr	72	Lys	42	Trp	13
Ile	72	Pro	42		

^a The number of times that a given amino acid residue has changed during the evolution of various proteins has been divided by the number of times the residue occurred. These values have been normalized by setting the highest to 100.

Adapted from M. O. Dayhoff, *Atlas of Protein Sequence and Structure*, vol. 5, suppl. 3, National Biomedical Research Foundation, Washington, D.C., 1978.

Phylogenetic trees

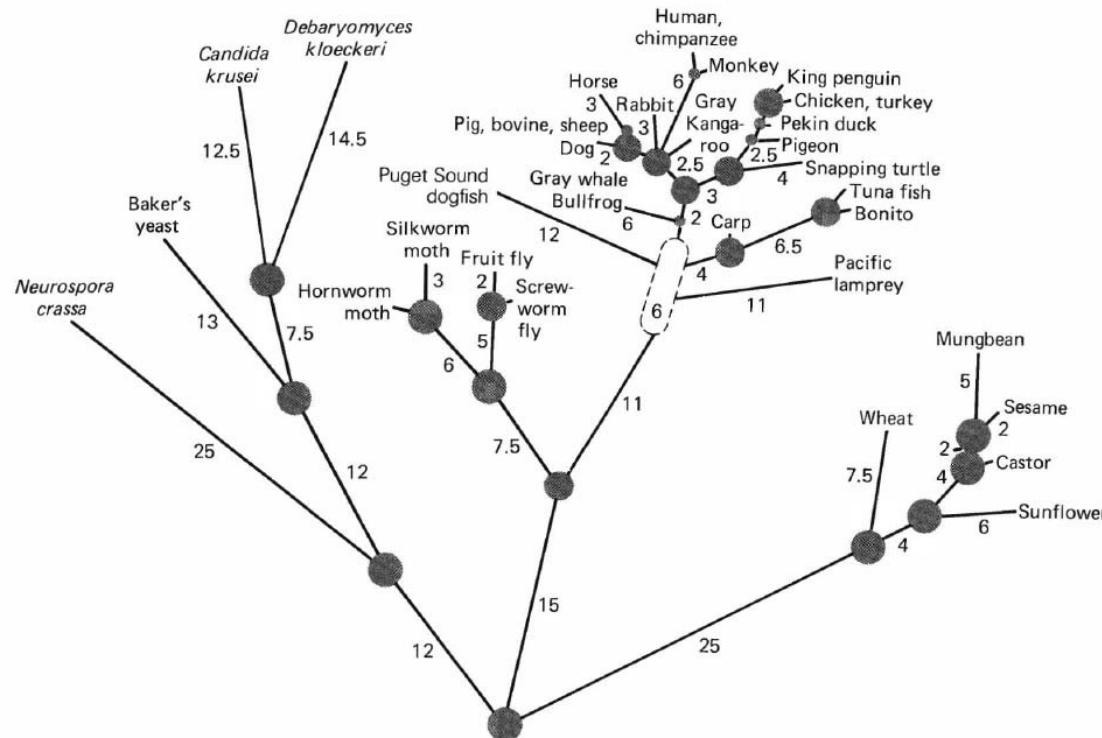


FIGURE 3.9

Phylogenetic tree constructed from the sequences of cytochromes *c* by minimizing the total number of amino acid replacements. The sequences of the ancestors at each branch point were inferred from the present-day sequences. The length of each branch is proportional to the indicated number of amino acid changes that are believed to have occurred. The branch points within the dashed oval were not adequately defined by the sequences. (Adapted from M. O. Dayhoff, *Atlas of Protein Sequence and Structure*, National Biomedical Research Foundation, Washington, D.C., 1972.)

Phylogenetic trees

Ancestral 1 5 10 15 20 25
Pro-Ala-Gly-Asp- ? -Lys-Lys-Gly-Ala-Lys-Ile-Phe-Lys-Thr- ? -Cys-Ala-Gln-Cys-His-Thr-Val-Glu- ? -Gly-Gly- ? -
Human Gly-Asp-Val-Glu-Lys-Gly-Lys-Ile-Phe-Ile -Met-Lys-Cys-Ser-Gln-Cys-His-Thr-Val-Glu-Lys-Gly-Gly-Lys-

 30 35 40 45 50
His-Lys-Val-Gly-Pro-Asn-Leu-His-Gly-Leu-Phe-Gly-Arg-Lys- ? -Gly-Gln-Ala- ? -Gly-Tyr-Ser-Tyr-Thr-Asp-
His-Lys-Thr-Gly-Pro-Asn-Leu-His-Gly-Leu-Phe-Gly-Arg-Lys-Thr-Gly-Gln-Ala-Pro-Gly-Tyr-Ser-Tyr-Thr-Ala-

 55 60 65 70 75
Ala-Asn-Lys-Asn-Lys-Gly- ? - ? -Trp- ? -Glu-Asn-Thr-Leu-Phe-Glu-Tyr-Leu-Glu-Asn-Pro-Lys-Lys-Tyr-Ile-
Ala-Asn-Lys-Asn-Lys-Gly-Ile-Ile-Trp-Gly-Glu-Asp-Thr-Leu-Met-Glu-Tyr-Leu-Glu-Asn-Pro-Lys-Lys-Tyr-Pro-

 80 85 90 95 100
Pro-Gly-Thr-Lys-Met- ? -Phe- ? -Gly-Leu-Lys-Lys- ? - ? -Asp-Arg-Ala-Asp-Leu-Ile-Ala-Tyr-Leu-Lys- ? -
Pro-Gly-Thr-Lys-Met-Ile-Phe-Val-Gly-Ile -Lys-Lys-Lys-Glu-Glu-Arg-Ala-Asp-Leu-Ile-Ala-Tyr-Leu-Lys-Lys-

Ala-Thr-Ala-
-Ala-Thr-Asn-Glu

Phylogenetic trees

Table 3.2 Relative Extents of Divergence for Several Proteins from Different Species

Protein	Residues Differing from Human Protein (%)								
	Rhesus monkey	Cow	Pig	Rabbit	Chicken	Frog	Fish	Fruit fly	Yeast
Cytochrome c	1	10	10	9	13	17	20	27	41
Hemoglobin α chain	3	12	13	18	25		50		
Hemoglobin β chain	5	17	16	10	26	46			
Fibrinopeptides A and B	30	70	67	70					

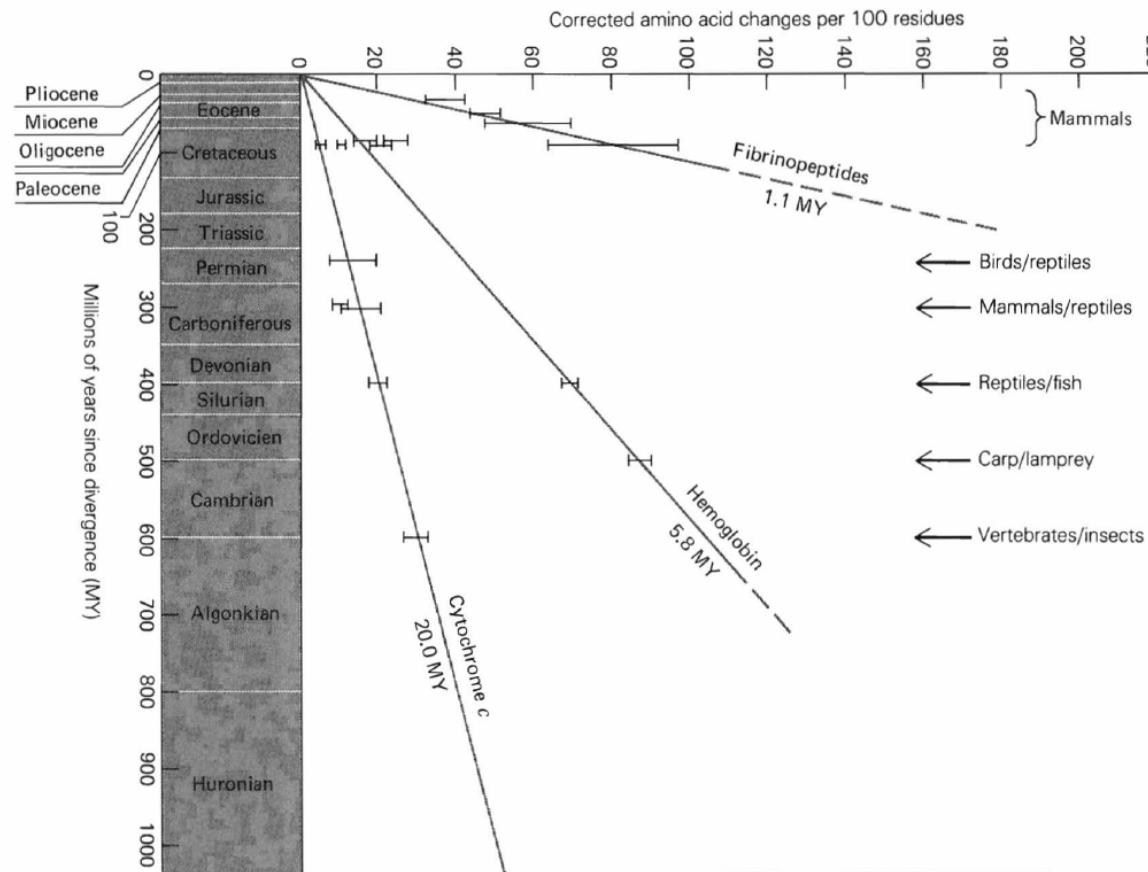
Rates of Divergence

Table 3.3 Rates of Evolution of Some Proteins

Protein	Accepted point mutations [(number)/(100 residues) · (10 ⁸ years)]	Protein	Accepted point mutations [(number)/(100 residues) · (10 ⁸ years)]
Histones		Hormones <i>cont.</i>	
H4	0.25	Proparathyrin	14
H3	0.30	Prolactin	20
H2A	1.7	Growth hormone	25
H2B	1.7	Lutropin β chain	33
H1	12	Insulin C peptide	53
Fibrous proteins		Oxygen-binding proteins	
Collagen (α -1)	2.8	Myoglobin	17
Crystallin (α A)	4.5	Hemoglobin α chain	27
Intracellular enzymes		Hemoglobin β chain	30
Glutamate dehydrogenase	1.8	Secreted enzymes	
Triosephosphate dehydrogenase	5.0	Trypsinogen	17
Triosephosphate isomerase	5.3	Lysozyme	40
Lactate dehydrogenase H	5.3	Ribonuclease A	43
Lactate dehydrogenase M	7.7	Immunoglobulins	
Carbonic anhydrase B	25	κ chains (V region)	100
Carbonic anhydrase C	48	κ chains (C region)	111
Electron carriers		λ chains (V region)	125
Cytochrome <i>c</i>	6.7	γ chains (V region)	143
Cytochrome b_5	9.1	Snake venom toxins	
Plastocyanin	14	Long neurotoxins	111
Ferredoxin	17	Cytotoxins	111
Hormones		Short neurotoxins	125
Glucagon	2.3	Other proteins	
Corticotropin	4.2	Parvalbumin	20
Insulin	7.1	Albumin	33
Thyrotropin β chain	11	α -Lactalbumin	43
Lipotropin β chain	13	Fibrinopeptide A	59
Lutropin α chain	14	Casein κ chain	71
		Fibrinopeptide B	91

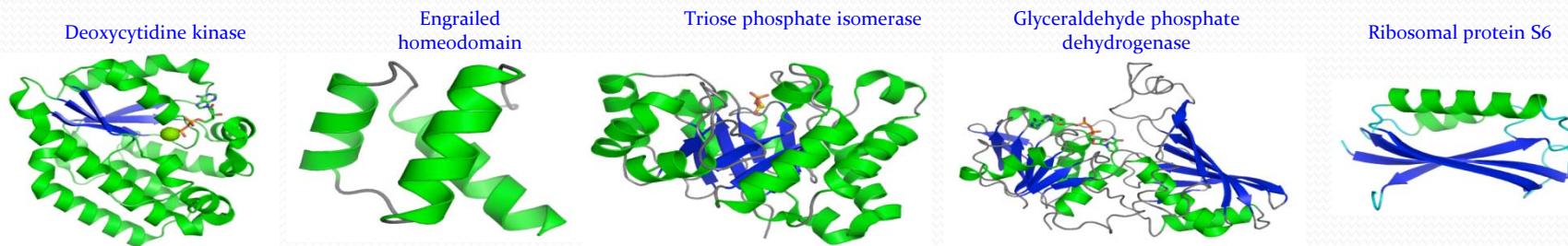
From A. C. Wilson et al., *Ann. Rev. Biochem.* 46:573–639 (1977).

Rates of divergence

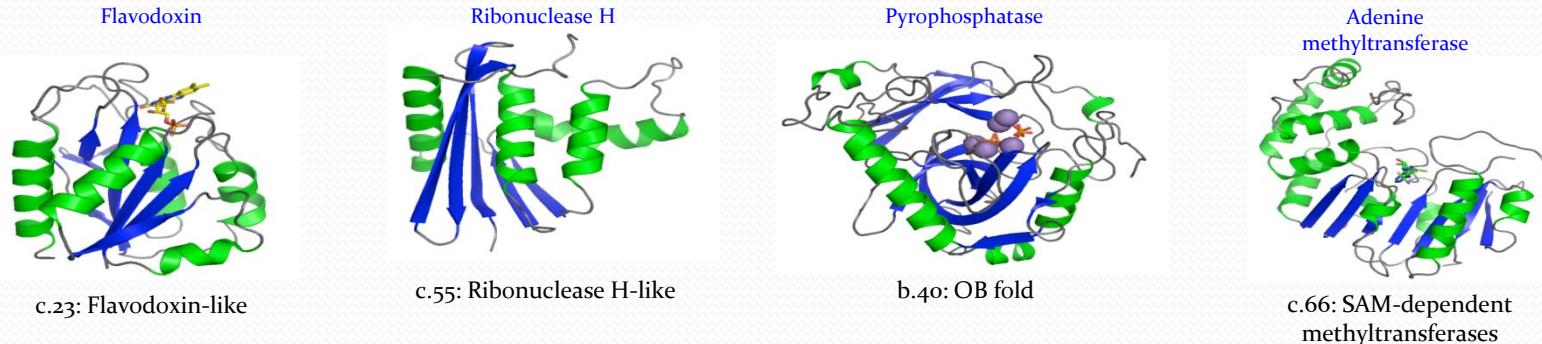


The nine most ancient protein folds

SCOP label	Fold
c.37	P-loop containing nucleoside triphosphate hydrolases 3 layers with $\alpha/\beta/\alpha$ arrangement, parallel or mixed β -sheets of variable sizes
a.4	DNA/RNA-binding 3-helical bundle Core: 3-helices; closed or partly opened bundle, right-handed twist; up-and-down
c.1	TIM β/α -barrel Closed barrel with parallel β -sheet and strand order 12345678; $n=8$, $S=8$
c.2	NAD(P)-binding Rossmann-fold domains Core: 3 layers in $\alpha/\beta/\alpha$ arrangement; parallel β -sheet of 6 strands, order 321456
d.58	Ferredoxin-like Core: 3 helices; closed or partly opened bundle, right-handed twist; up-and-down
c.23	Flavodoxin-like 3 layers with $\alpha/\beta/\alpha$ arrangement; parallel β -sheet of 5 strands, order 21345
c.55	Ribonuclease H-like motif 3 layers with $\alpha/\beta/\alpha$ arrangement; mixed β -sheet of 5 strands, order 32145 with strand 2 antiparallel to the rest
b.40	OB-fold Closed or partly opened barrel, $n=5$, $S=10$ or $S=8$ with greek-key motif
c.66	S-adenosyl-L-methionine-dependent methyltransferases Core: 3 layers with $\alpha/\beta/\alpha$ arrangement; mixed β -sheet of 7 strands, order 3214576 with strand 7 antiparallel to the rest



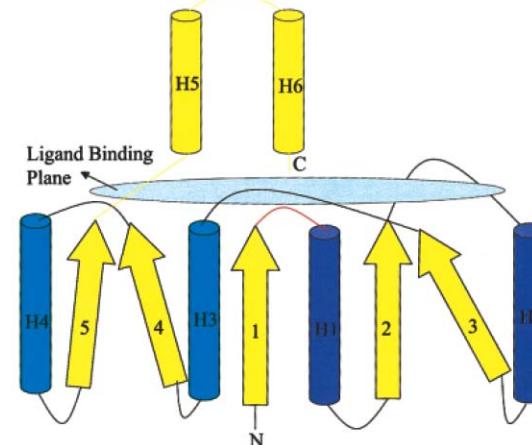
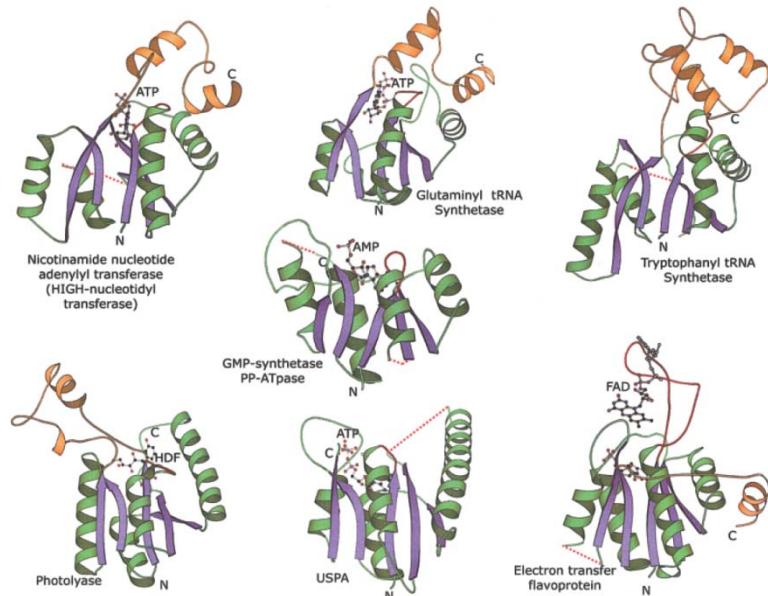
c.37: P-loop containing NTP hydrolase
 a.4: DNA/RNA binding 3 helix bundle
 c.1: 8-stranded α/β barrel
 c.2: NAD(P) binding Rossmann fold
 d.58: Ferredoxin-like



Not all widely diversified and broadly functional folds are ancient

- Adenine nucleotide binding protein

Scop C.28: also called 'HUP' domains



Gene rearrangements

Table 3.4 Number of Amino Acid Identities between Various Members of the Human Globin Family

	Hemoglobin				
	α	β	γ	δ	ϵ
Myoglobin	38	36	36	37	34
Hemoglobin α		64	59	63	55
Hemoglobin β			107	136	110
Hemoglobin γ				105	116
Hemoglobin δ					106

Gene rearrangements

Myoglobin α β γ^{G} δ ϵ	1 10 20 30 40 50 60 70 80 90 100 110 120 130 140 150	Gly - - - Leu-Ser -Asp-Gly-Glu-Trp-Gln-Leu-Val-Leu-Asn- Val -Trp-Gly-Lys-Val-Glu-Ala-Asp- Ile -Pro-Gly-His-Gly-Gln-Glu-Val- Val - - - Leu-Ser -Pro-Ala-Asp-Lys-Thr-Asn-Val-Lys-Ala-Ala -Trp-Gly-Lys-Val-Gly-Ala-His-Ala-Gly-Gln-Tyr-Gly-Ala-Glu-Ala- Val-His-Leu-Thr-Pro-Glu-Glu-Lys-Ser-Ala-Ala -Val-Thr-Ala- Leu -Trp-Gly-Lys-Val- Asn - - - - - Val-Asp-Glu-Val-Gly-Gly-Glu-Ala- Gly-His-Phe -Thr-Glu-Glu-Asp-Lys-Ala-Thr- Ile -Thr-Ser- Leu -Trp-Gly-Lys-Val- Asn - - - - - Val-Glu-Asp-Ala-Gly-Gly-Glu-Thr- Val-His-Leu-Thr-Pro-Glu-Glu-Lys-Thr-Ala-Val-Asn-Ala -Leu-Trp-Gly-Lys-Val- Asn - - - - - Val-Asp-Ala-Val-Gly-Gly-Glu-Ala- Val-His-Phe -Thr-Ala-Glu-Glu-Lys-Ala-Ala -Val-Thr-Ser- Leu -Trp-Ser-Lys-Met- Asn - - - - - Val-Glu-Glu-Ala-Gly-Gly-Glu-Ala- -Leu-Ile -Arg-Leu-Phe-Lys-Gly-His -Pro-Glu -Thr-Leu-Glu-Lys-Phe-Asp-Lys-Phe-Lys-His -Leu-Lys-Ser-Glu-Asp-Glu-Met-Lys -Ala-Ser -Glu- -Leu-Glu-Arg-Met-Phe -Leu-Ser-Phe-Pro-Thr -Thr-Lys -Thr -Tyr -Phe-Pro -His-Phe - - Asp-Leu-Ser- Ile - - - - - - - Gly-Ser -Ala - -Leu-Gly-Arg-Leu-Leu-Val-Val-Tyr-Pro-Trp-Thr-Gln-Arg-Phe-Phe-Glu -Ser-Phe-Gly-Asp-Leu-Ser-Thr-Pro-Asp-Ala-Val -Met-Gly-Asn-Pro- -Leu-Gly-Arg-Leu-Leu-Val-Val-Tyr-Pro-Trp-Thr-Gln-Arg-Phe-Phe-Glu -Ser-Phe-Gly-Asn-Leu-Ser-Ser-Ala-Ser -Ala-Ile -Met-Gly-Asn-Pro- -Leu-Gly-Arg-Leu-Leu-Val-Val-Tyr-Pro-Trp-Thr-Gln-Arg-Phe-Phe-Glu -Ser-Phe-Gly-Asp-Leu-Ser-Ser-Pro-Asp-Ala-Val -Met-Gly-Asn-Pro- -Leu-Gly-Arg-Leu-Leu-Val-Val-Tyr-Pro-Trp-Thr-Gln-Arg-Phe-Phe-Glu -Ser-Phe-Gly-Asp-Leu-Ser-Ser-Pro-Ser-Ala-Ile -Leu-Gly-Asn-Pro- -Asp-Leu-Lys-Lys-His-Gly-Ala-Thr-Val-Leu-Thr-Ala-Leu-Gly-Gly-Ile -Leu-Lys-Lys-Lys-Gly-His -His -Glu-Ala -Glu -Ile -Lys-Pro-Leu-Ala- -Gln-Val-Lys-Gly-His-Gly-Lys-Lys-Val-Ala -Asp-Ala-Leu-Thr-Asn-Ala- Val -Ala-His -Val -Asp-Asp-Met-Pro -Asn-Ala -Leu-Ser -Ala-Leu-Ser- -Lys-Val-Lys-Ala-His-Gly-Lys-Lys-Val-Leu-Gly-Ala-Phe -Ser-Asp-Gly-Leu-Ala-His-Leu-Asp-Ala- Leu -Lys-Gly-Thr-Phe-Ala-Thr-Leu-Ser- -Lys-Val-Lys-Ala-His-Gly-Lys-Lys-Val-Leu-Thr-Ser-Leu-Gly-Asp-Ala-Ile -Lys -His -Leu-Asp-Asp-Leu-Lys-Gly-Thr-Phe-Ala-Gln-Leu-Ser- -Lys-Val-Lys-Ala-His-Gly-Lys-Lys-Val-Leu-Gly-Ala-Phe -Ser-Asp-Gly-Leu-Ala-His -Leu-Asp-Asn-Leu-Lys-Gly-Thr-Phe-Ser-Gln-Leu-Ser- -Lys-Val-Lys-Ala-His-Gly-Lys-Lys-Val-Leu-Thr-Ser-Phe-Gly-Asp-Ala-Ile -Lys -Asn-Met-Asp-Asn-Leu-Lys-Pro -Ala -Phe-Ala-Lys-Leu-Ser- -Gin-Ser -His-Ala -Thr -Lys-His -Lys -Ile -Pro -Val-Lys-Tyr -Leu -Glu -Phe -Ile -Ser -Glu -Cys -Ile -Ile -Gln -Val -Leu -Gin -Ser -Lys -His -Pro -Gly- -Asp -Leu -His -Ala -His -Lys -Leu -Arg -Val -Asp -Pro -Val -Asn -Phe -Lys -Leu -Ser -His -Cys -Leu -Leu -Val -Thr -Leu -Ala -Ala - His -Leu -Pro -Ala- -Glu -Leu -His -Cys -Asp -Lys -Leu -His -Val -Asp -Pro -Glu -Asn -Phe -Arg -Leu -Leu -Gly -Asn -Val -Leu -Val -Cys -Val -Leu -Ala -His - His -Phe -Gly -Lys- -Glu -Leu -His -Cys -Asp -Lys -Leu -His -Val -Asp -Pro -Glu -Asn -Phe -Arg -Leu -Leu -Gly -Asn -Val -Leu -Val -Thr -Val -Leu -Ala -Ile - His -Phe -Gly -Lys- -Glu -Leu -His -Cys -Asp -Lys -Leu -His -Val -Asp -Pro -Glu -Asn -Phe -Arg -Leu -Leu -Gly -Asn -Val -Leu -Val -Cys -Val -Leu -Ala -Arg -Asn -Phe -Gly -Lys- -Glu -Leu -His -Cys -Asp -Lys -Leu -His -Val -Asp -Pro -Glu -Asn -Phe -Lys -Leu -Leu -Gly -Asn -Val -Met -Val -Ile -Ile -Leu -Ala -Thr - His -Phe -Gly -Lys
--	---	--

FIGURE 3.12

Amino acid sequences of members of the human globin family, myoglobin plus the hemoglobin polypeptide chains. Positions where three or more of the amino acid residues are identical are in boldface; chemically similar residues are in italics.

Gene rearrangements

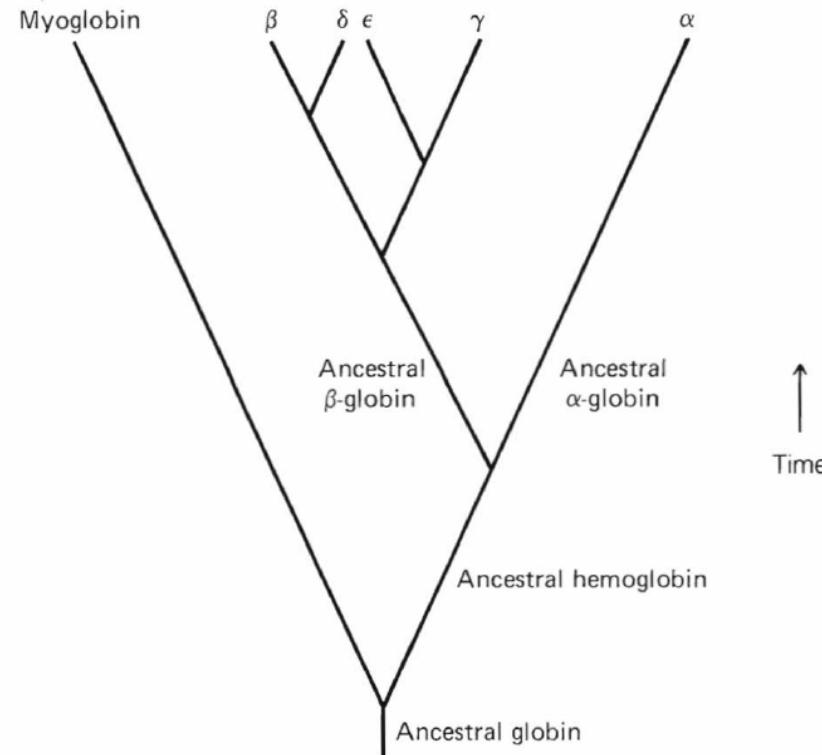
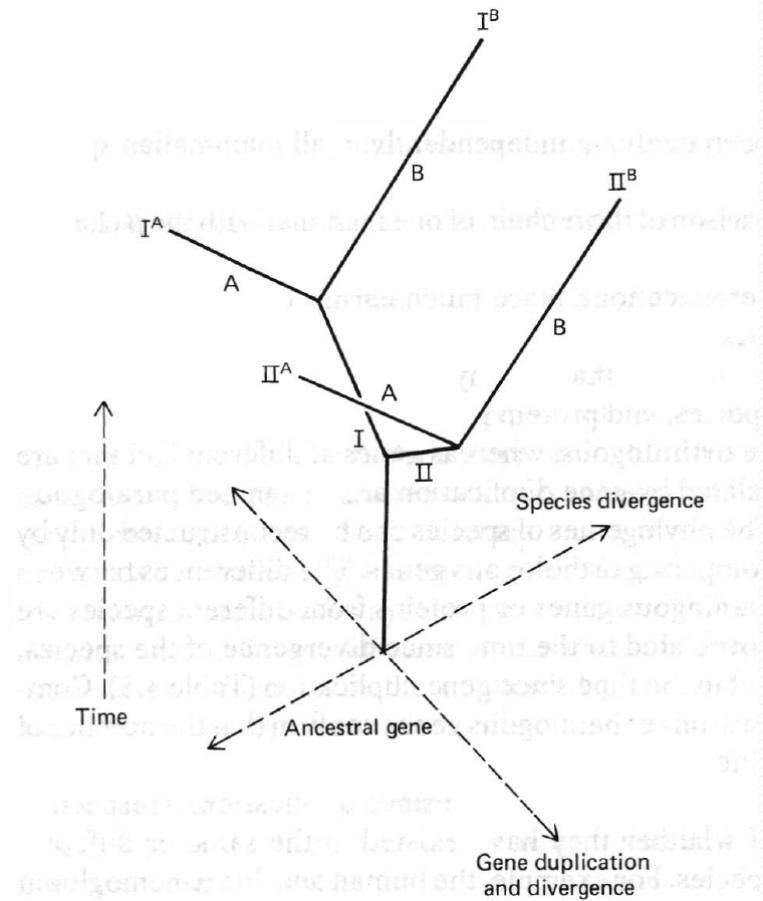


FIGURE 3.13

Evolutionary scheme for the divergence of the globin chains, inferred from the relative similarities of their gene and protein sequences. Time increases from the bottom to the top of the figure. The separation of the $^A\gamma$ and $^G\gamma$ genes is not depicted here. (Adapted from A. Efstratiadis et al., *Cell* 21:653–668, 1980.)

Gene rearrangements



Paralogous comparison

Table 3.5 Paralogous Comparisons of Globin Chains^a

Comparison protein	Number of Differences in Amino Acid Sequence	
	Human α -globin	Dog α -globin
Human β -globin	84	84
Human γ -globin	89	83
Human myoglobin	115	119
Dog β -globin	88	87
Dog myoglobin	115	118

^a The globin genes duplicated and diverged long before the human and dog species diverged, so comparisons of the different globins are relevant to the time since the genes duplicated. Consequently, very similar results are obtained irrespective of which species is used.

Elongation by intragene duplication

1		10	
Ala-Tyr	-Lys-Ile-	-Ala-Asp-Ser-Cys-Val-Ser	-Cys-Gly-Ala
-Ile	-Phe-Val-Ile-Asp-Ala-Asp-Thr-Cys-Ile	-Asp-Cys-Gly-Asn	
	30		40
		20	
-Cys-Ala-Ser		-Glu-Cys-Pro-Val-Asn-Ala-Ile	-Ser-Gln-Gly-Asp-Ser
-Cys-Ala-Asp-Val-Cys-Pro-Val-Gly		-Ala-Pro-Val-Gln-Glu	
		50	55

FIGURE 3.15

Homology between the two halves of the primary structure of *Clostridium pasteurianum* ferredoxin. Identical residues are in boldface, similar residues in italics.

Gene fusion and division

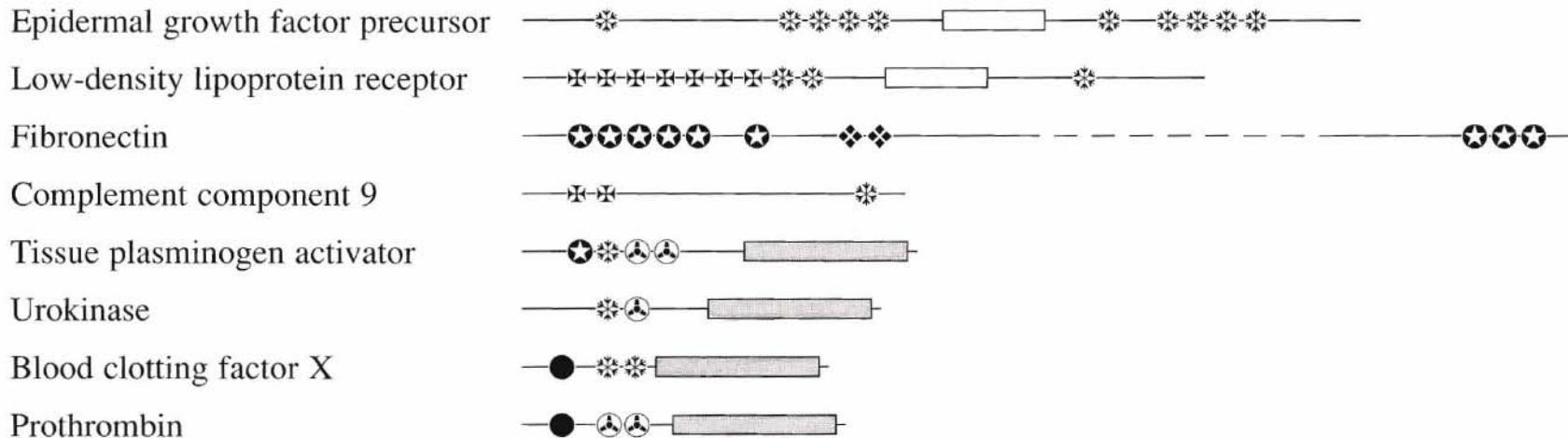


FIGURE 3.16

Mosaic structures of some vertebrate proteins. Each symbol represents members of a set of homologous segments of polypeptide chain, which in these examples are also independent structural units, or domains, of the protein. Those with distinct structural or functional properties are the serine protease domains (shaded rectangles), Ca^{2+} -binding domains containing γ -carboxy-Glu residues (●), so-called kringle domains (○), and epidermal growth factor (EGF) domains (※).

Hybrid protein by mismatch recombination

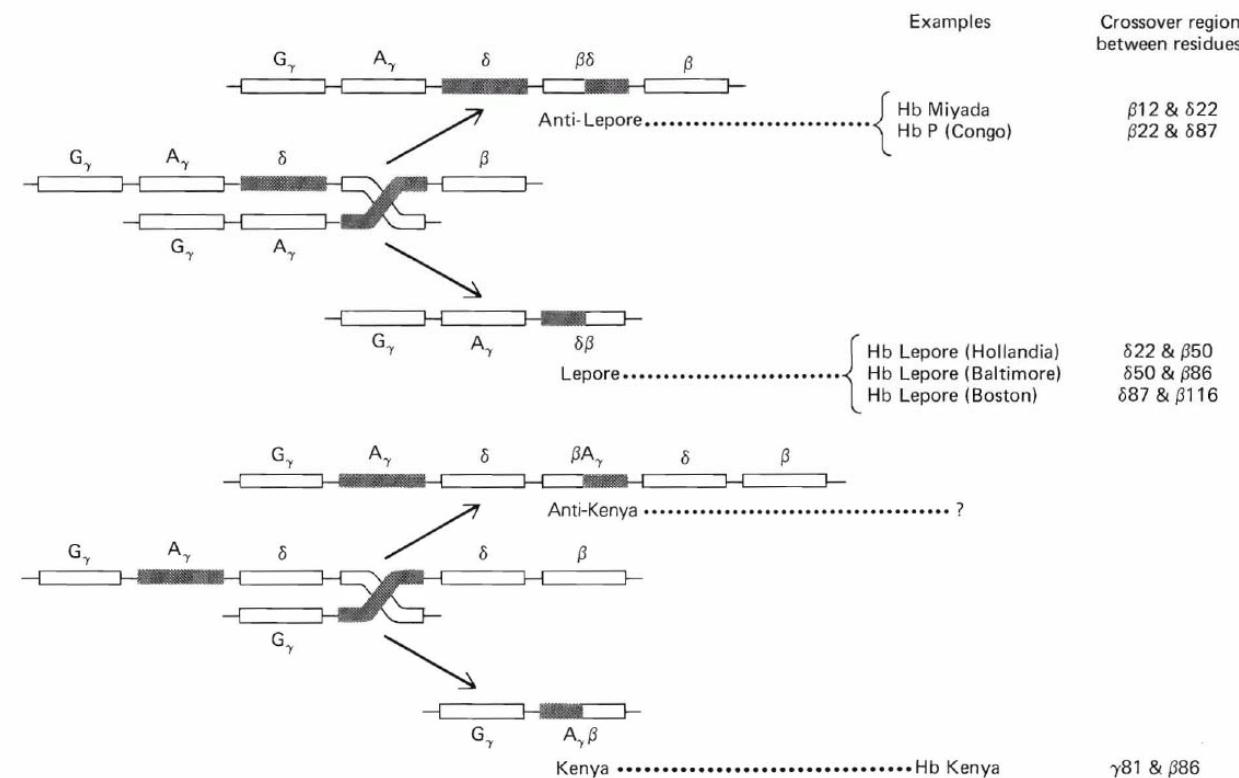


FIGURE 3.17

Genetic mechanism for production of hybrid proteins by mismatched recombination between duplicated genes. The chromosomal arrangements of the β -like globin genes are depicted, with recombination between the δ and β genes at the top and between the A_{γ} and β genes at the bottom. These events are inferred from the amino acid sequences of the hybrid proteins produced. (From D. J. Weatherall and J. B. Clegg, *Cell* 16:467–479, 1979.)