



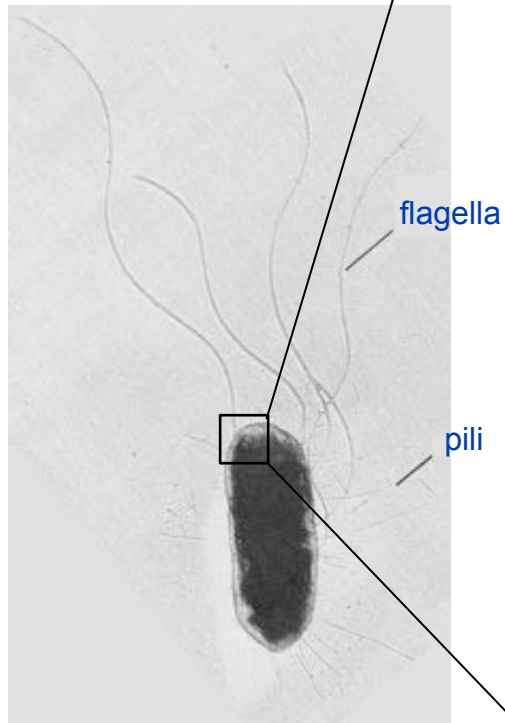
Proteins: an introduction

Life is the result of the chemical activity of proteins

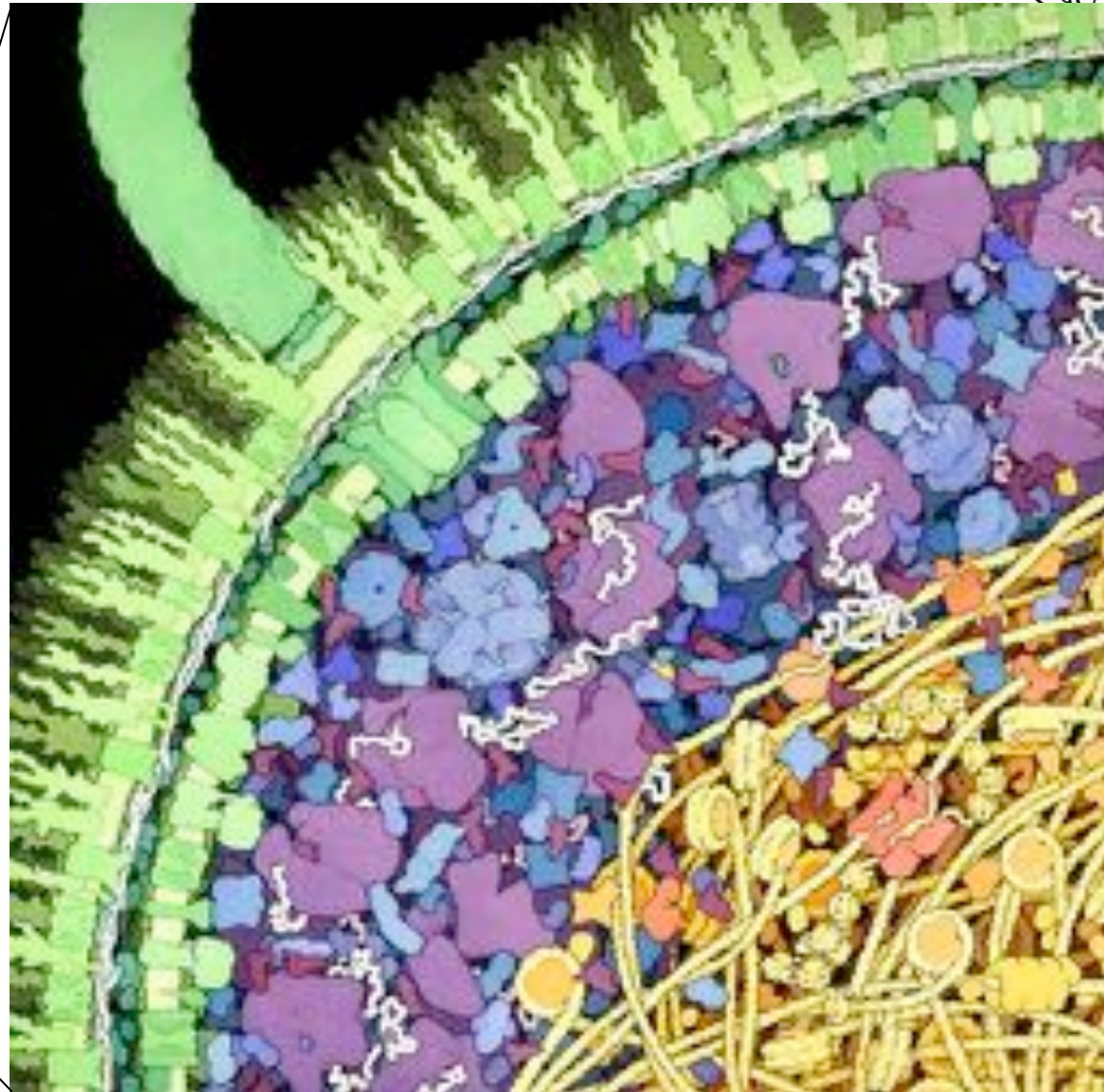
(with nucleic acids essentially confined to encoding and producing them)



A cross-section through *Escherichia coli* shows the diversity and density of macromolecules forming a cell



Howard Berg



David Goodsell

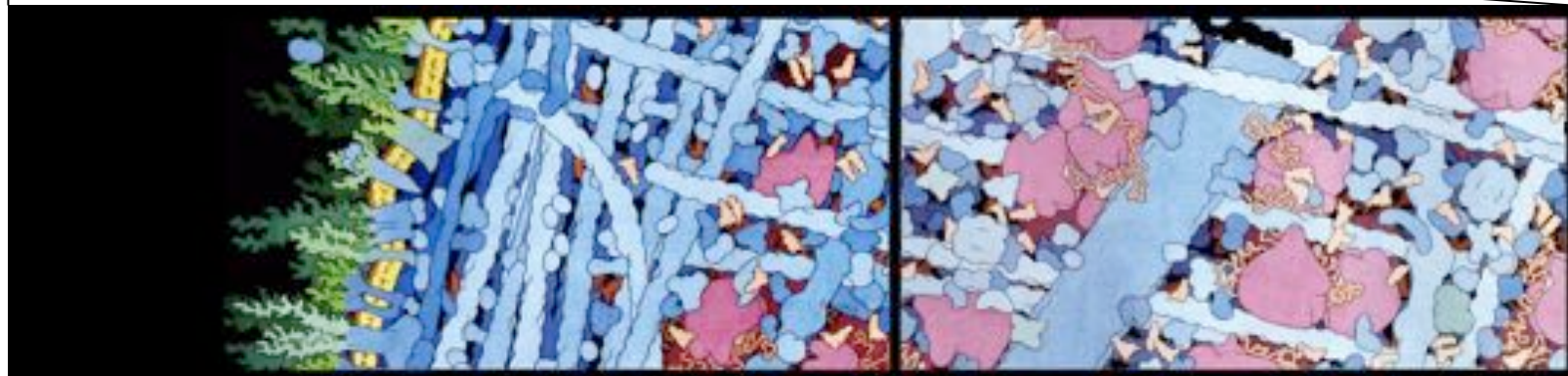
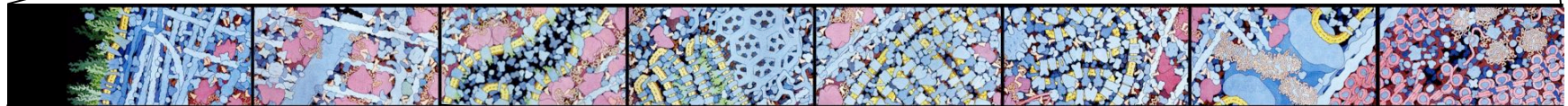
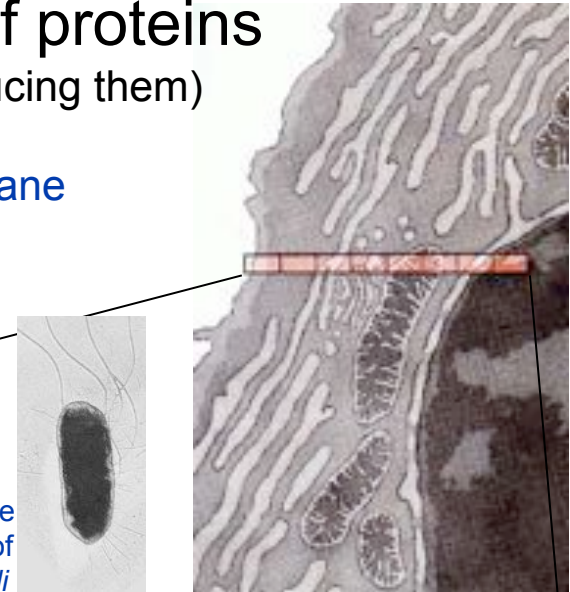
Life is the result of the chemical activity of proteins

(with nucleic acids essentially confined to encoding and producing them)

A cross-section through a human cell, from the plasma membrane to the nucleus

The enlargement shows cell-surface glycoproteins, the sub-membraneous cytoskeleton (primarily actin filaments), a microtubulus, ribosomes (in pink; they are everywhere except the nucleus)

relative size of *E. coli*

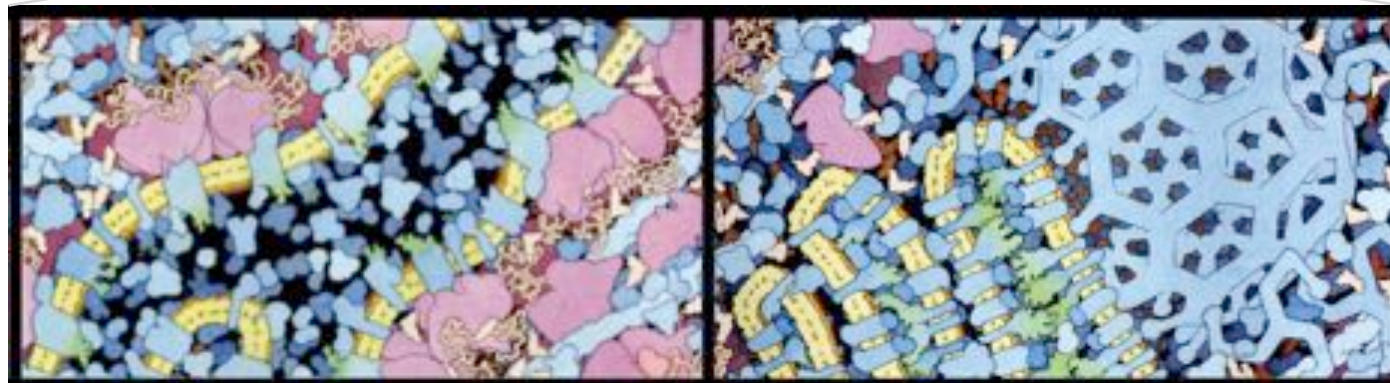
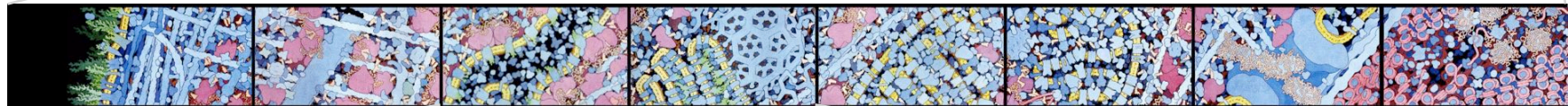


Life is the result of the chemical activity of proteins

(with nucleic acids essentially confined to encoding and producing them)

A cross-section through a human cell, from the plasma membrane to the nucleus

The enlargement shows the endoplasmic reticulum (,rough ER'), Golgi stacks, and a clathrin cage assembling for vesicle budding

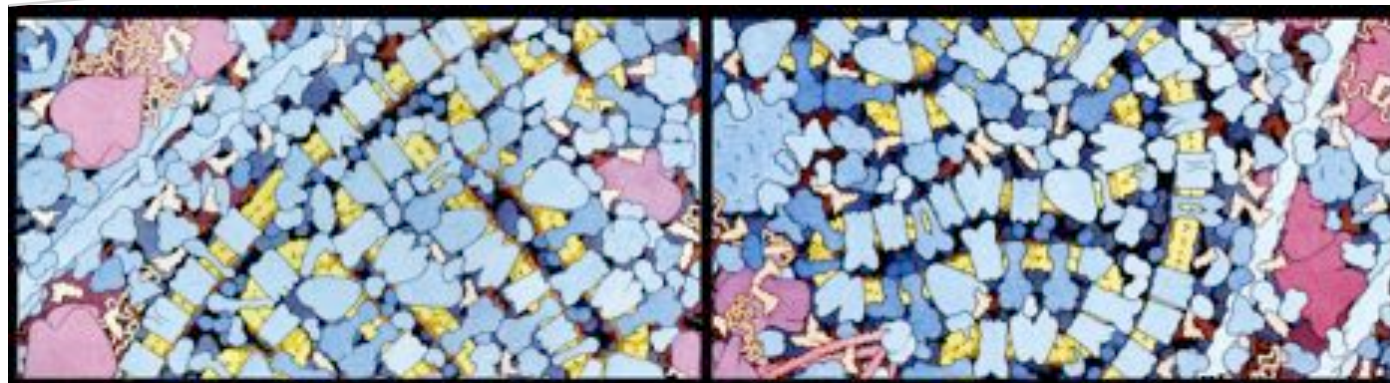
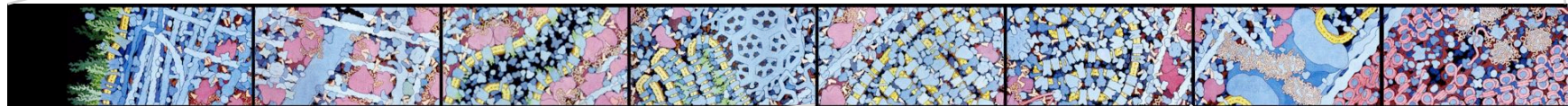


Life is the result of the chemical activity of proteins

(with nucleic acids essentially confined to encoding and producing them)

A cross-section through a human cell, from the plasma membrane to the nucleus

The enlargement shows a mitochondrion



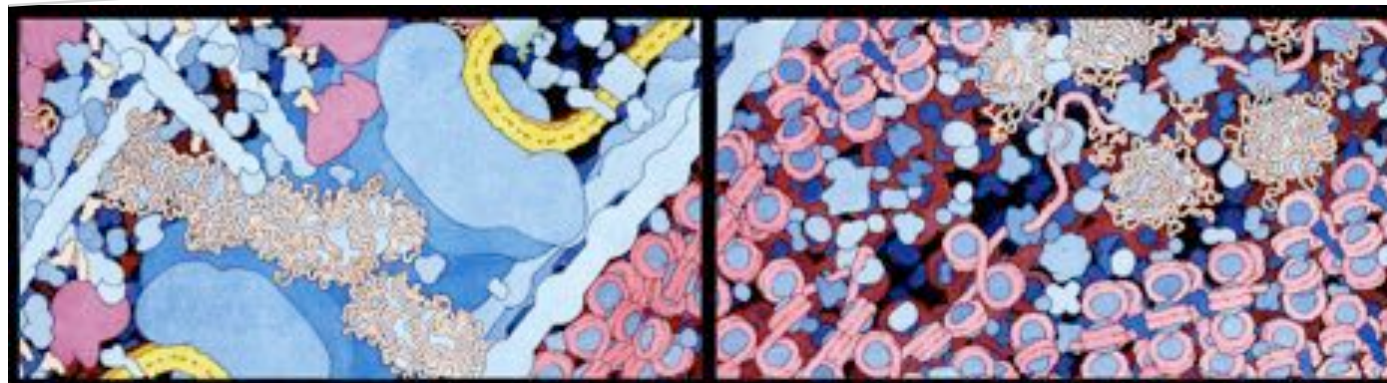
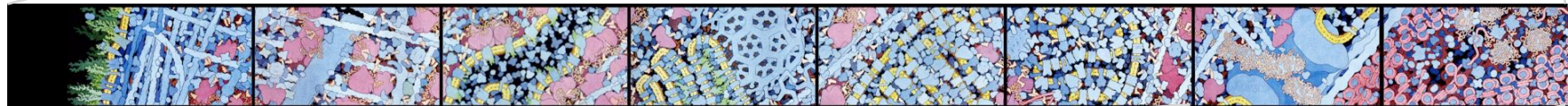
David Goodsell

Life is the result of the chemical activity of proteins

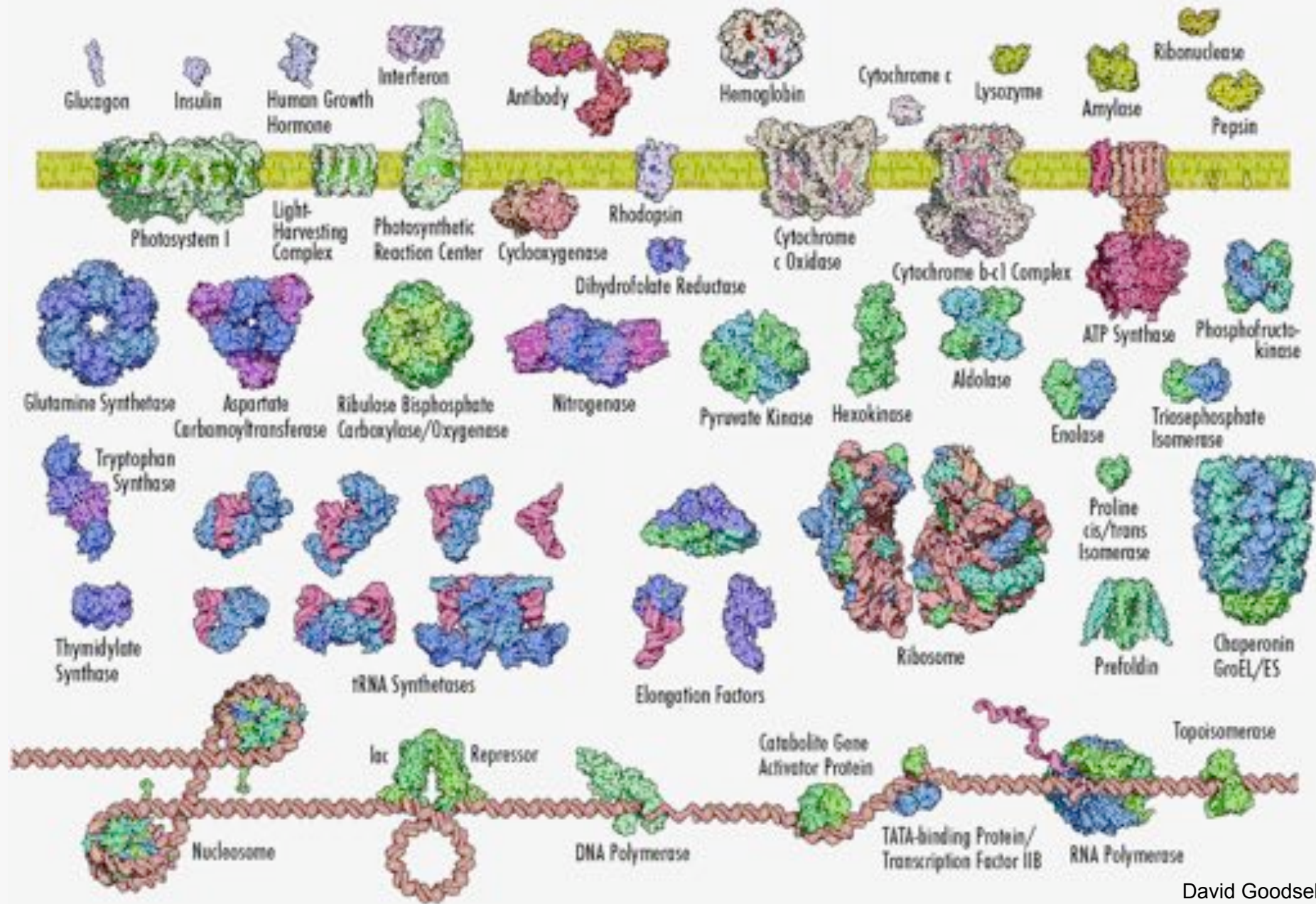
(with nucleic acids essentially confined to encoding and producing them)

A cross-section through a human cell, from the plasma membrane to the nucleus

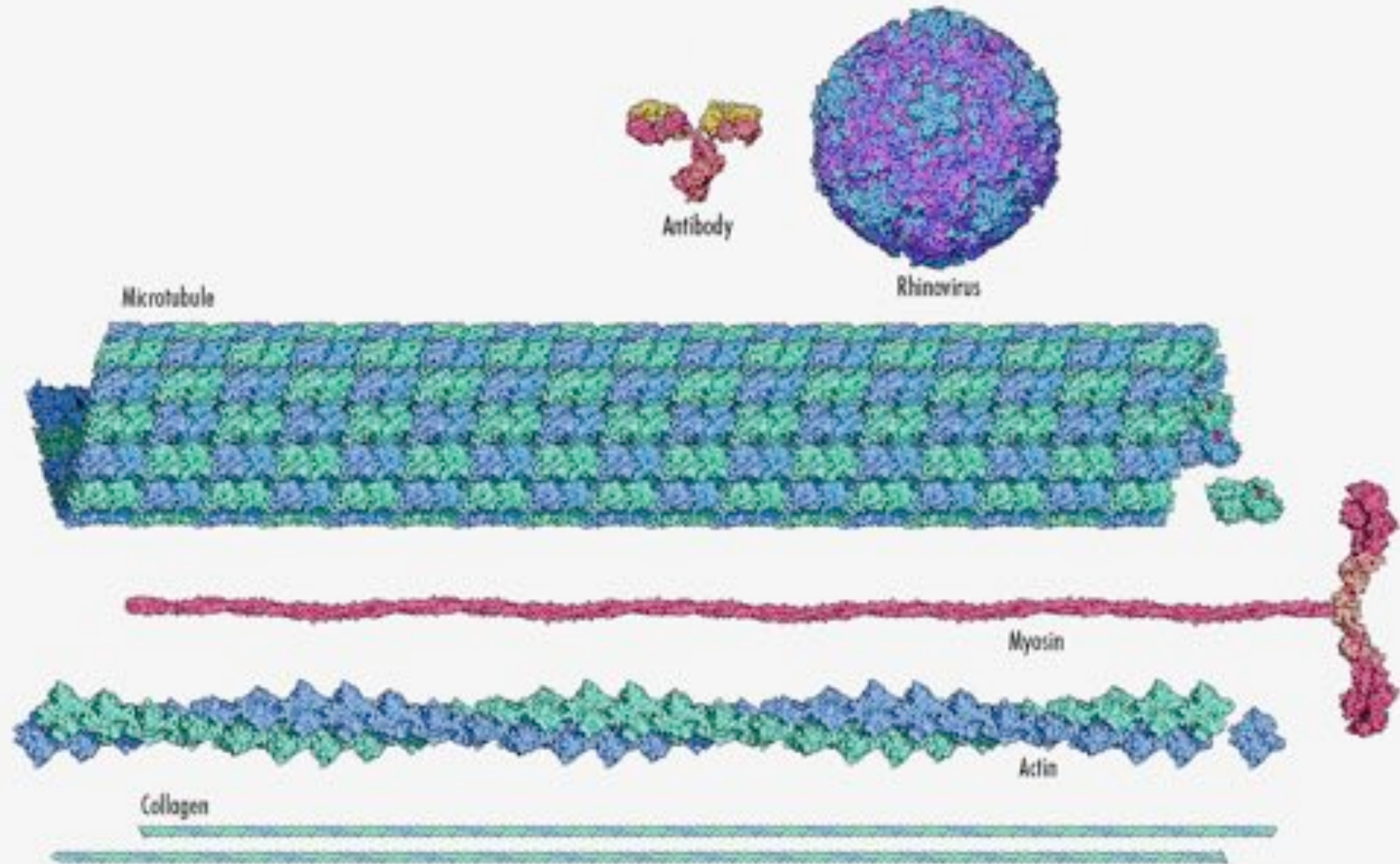
The enlargement shows a nuclear pore, the perinuclear and intranuclear cytoskeleton, chromatin, spliceosomes and other RNP particles



Proteins are extremely diverse



Proteins are extremely diverse



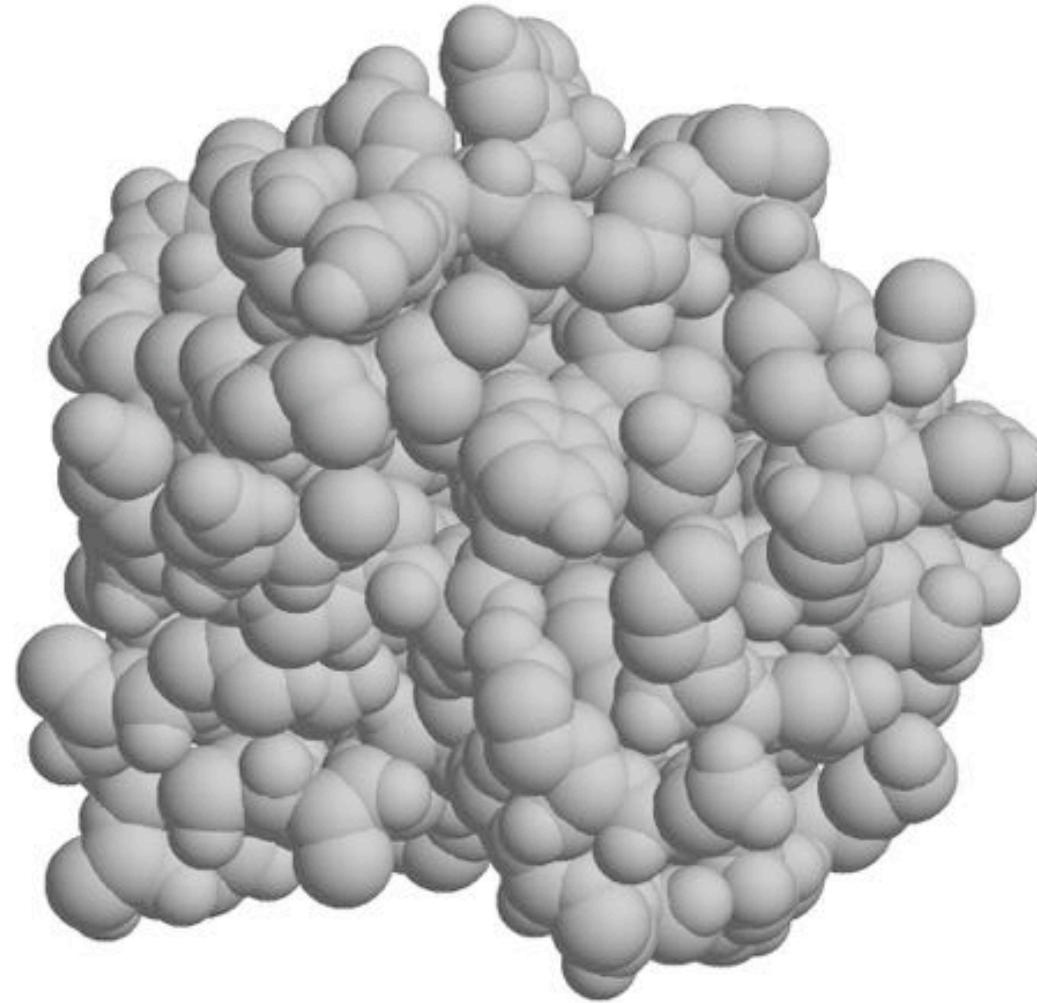
What is a protein?



- *protein* - 1844, from Fr. protéine, coined 1830s by Dutch chemist Mulder from Gk. proteios "the first quality," from protos "first." Originally a theoretical substance thought to be essential to life, the modern use is from Ger. Protein, borrowed in Eng. 1907. (Etymology Online)
- Unbranched polymer chain of α -L-amino acids connected by peptide bonds (condensation between the carboxyl- and amino-groups of consecutive monomers)
- "The term protein is usually reserved for those chains with a specific sequence, length, and **folded conformation**." (Creighton, Proteins, W.H.Freeman and Co. 1984)

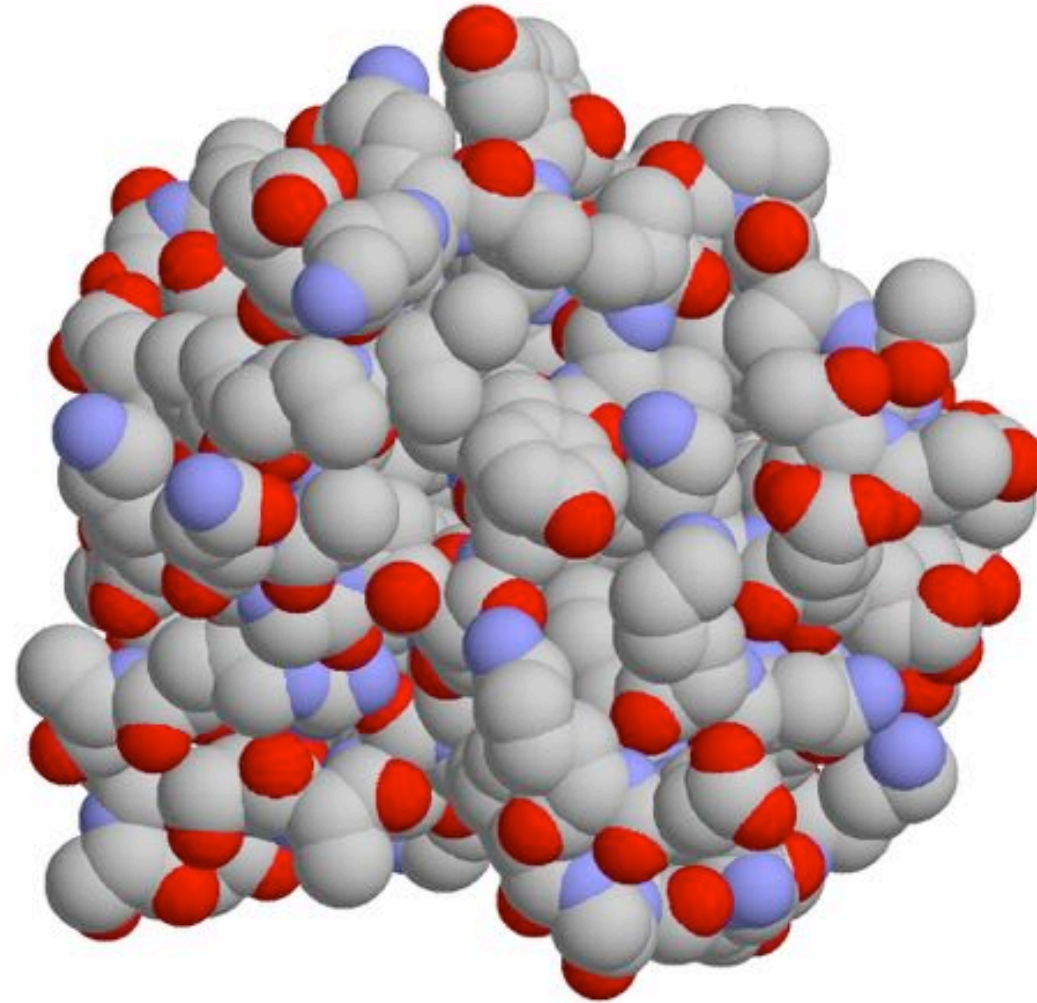
Schematic representations of protein structure

- spacefilling representation



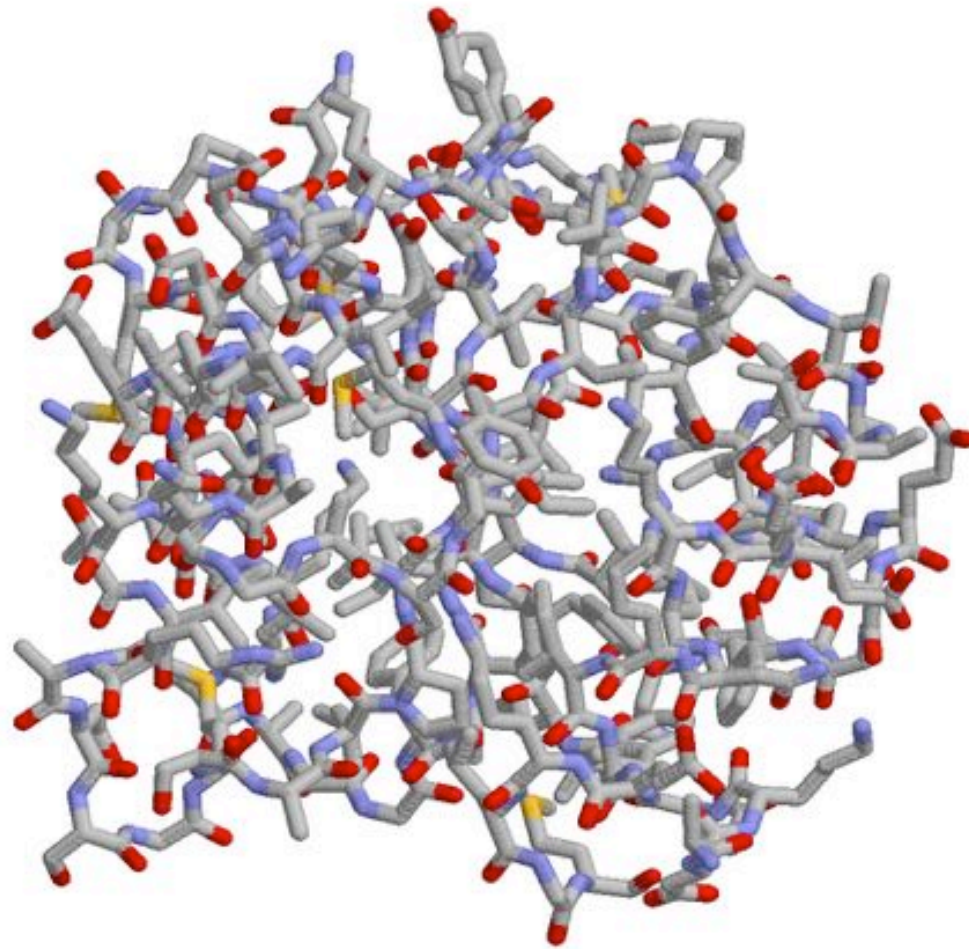
Schematic representations of protein structure

- spacefilling representation with CPK (Corey-Pauling-Koltun) colors
(C - grey, N - blue, O - red, S - yellow, P - purple, halogens - green)



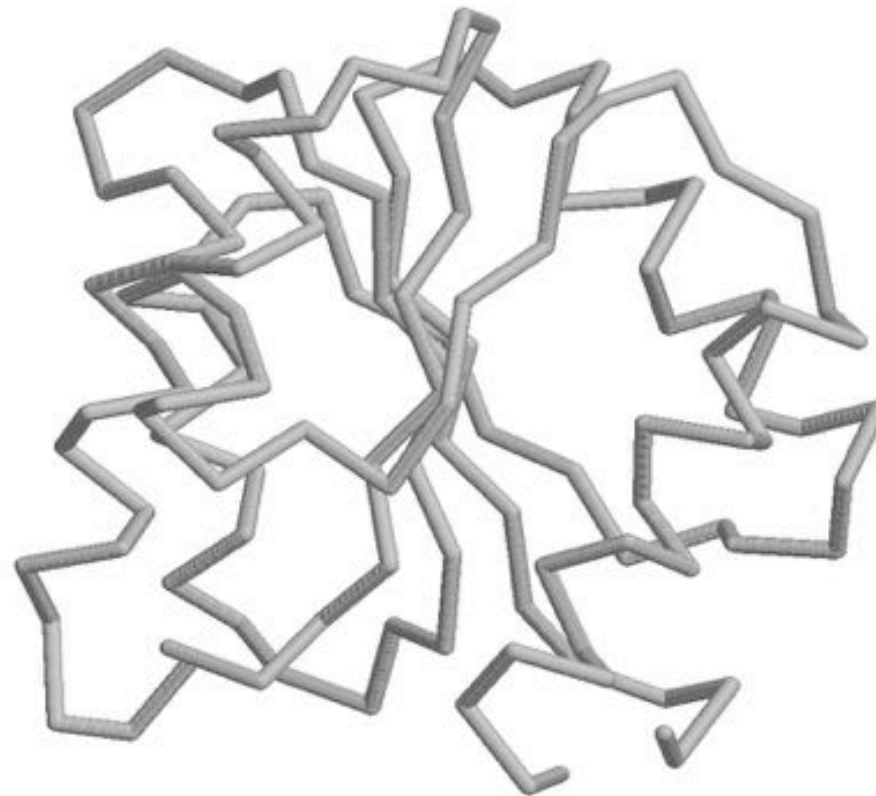
Schematic representations of protein structure

- stick representation with CPK colors



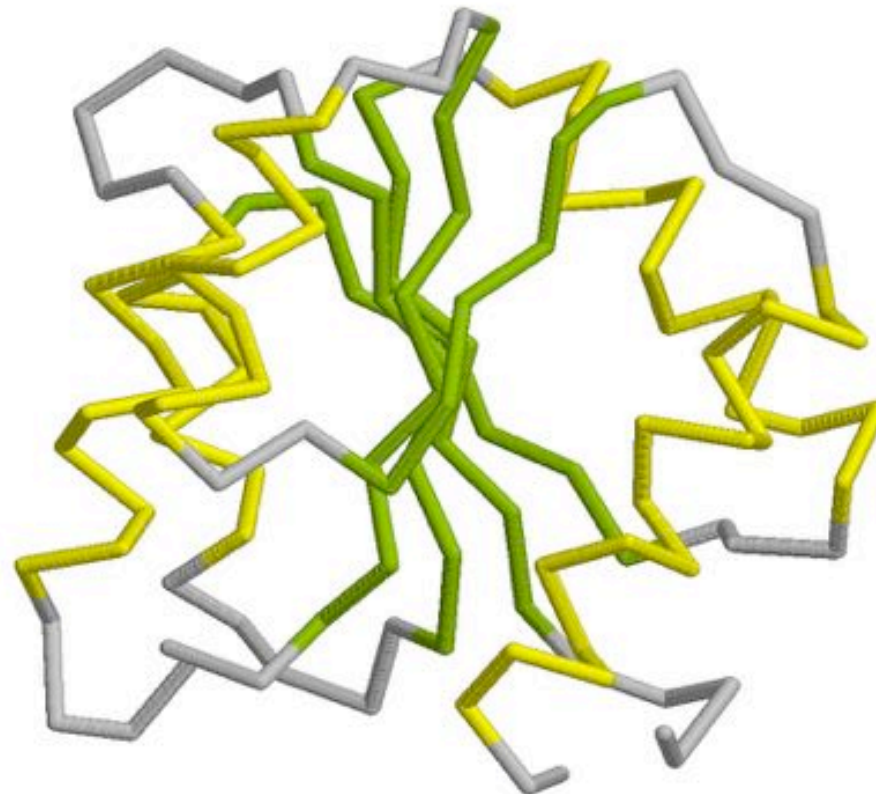
Schematic representations of protein structure

- C α trace



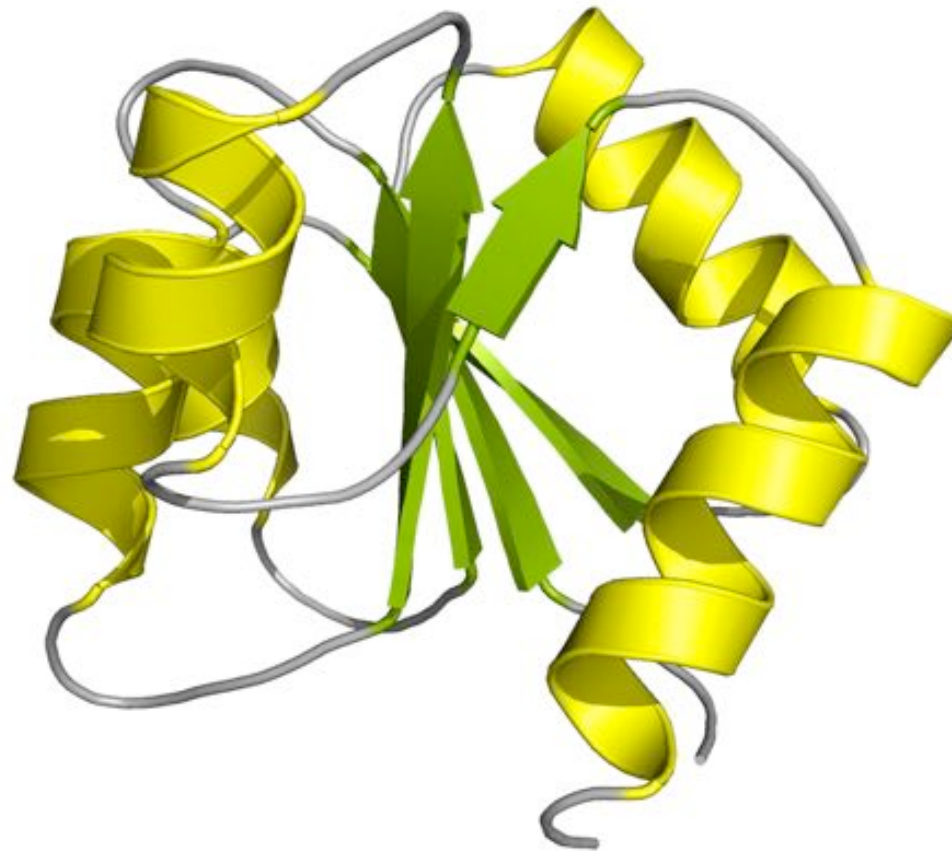
Schematic representations of protein structure

- C α trace colored by secondary structure (strand - green, helix - yellow)



Schematic representations of protein structure

- cartoon representation colored by secondary structure

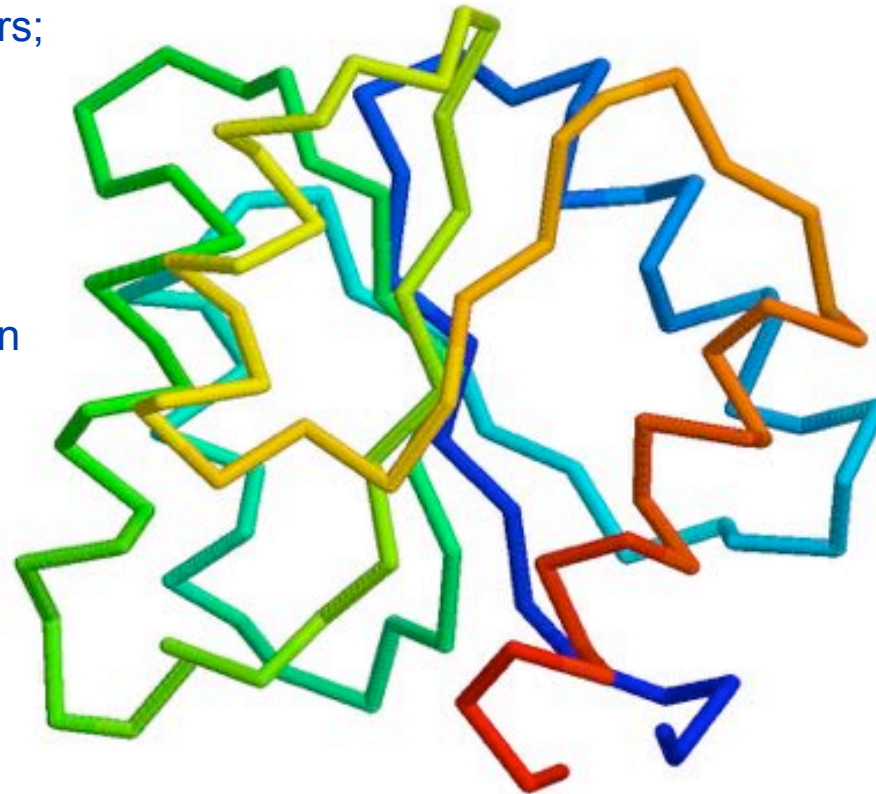


The path of a polypeptide chain

- C α trace in chainbow coloring
(from blue at the N-terminus to red at the C-terminus)



Despite the convoluted way in which their backbone folds, proteins are unbranched polymers; there are however instances of branching through covalent sidechain-sidechain (e.g. disulfide bonds) and sidechain-mainchain (e.g. ubiquitination) bonds introduced post-translationally



Levels of protein structure

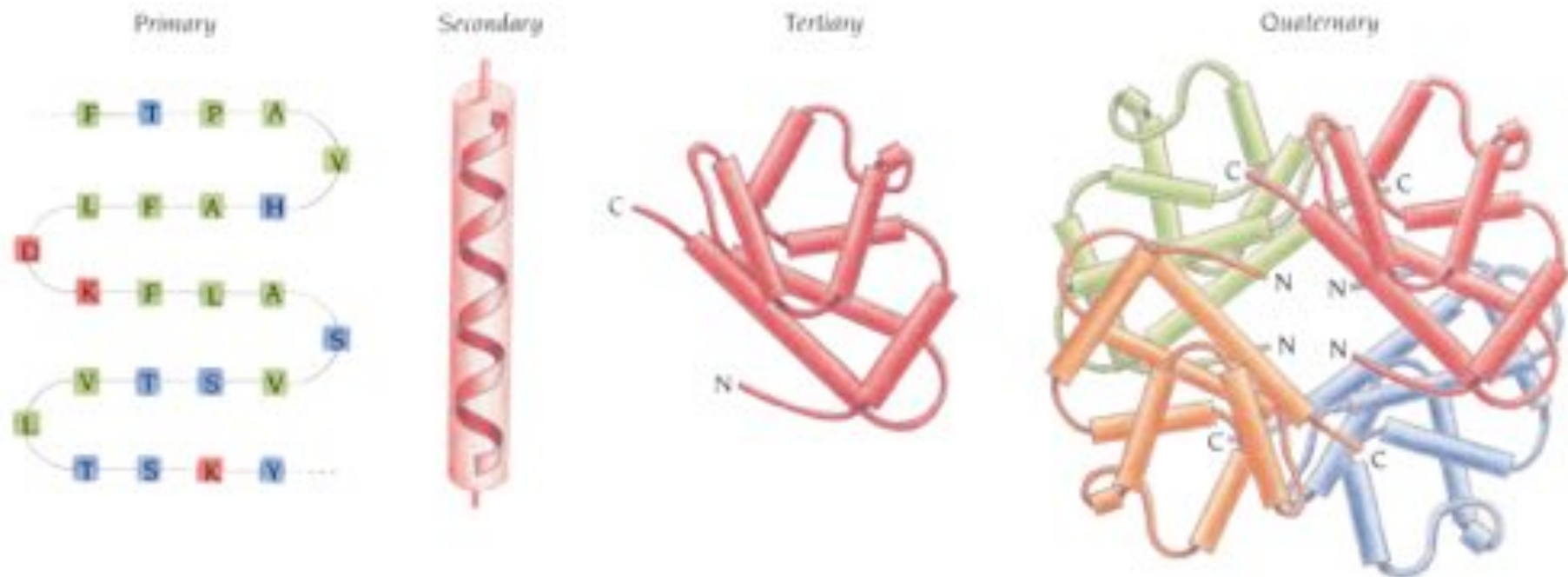


Primary - the succession of amino acids in the polypeptide chain (N to C)

Secondary - local, hydrogen-bonded configurations of the polypeptide chain

Tertiary - topological arrangement of secondary structures in 3D

Quaternary - non-covalent association of tertiary structures into a complex

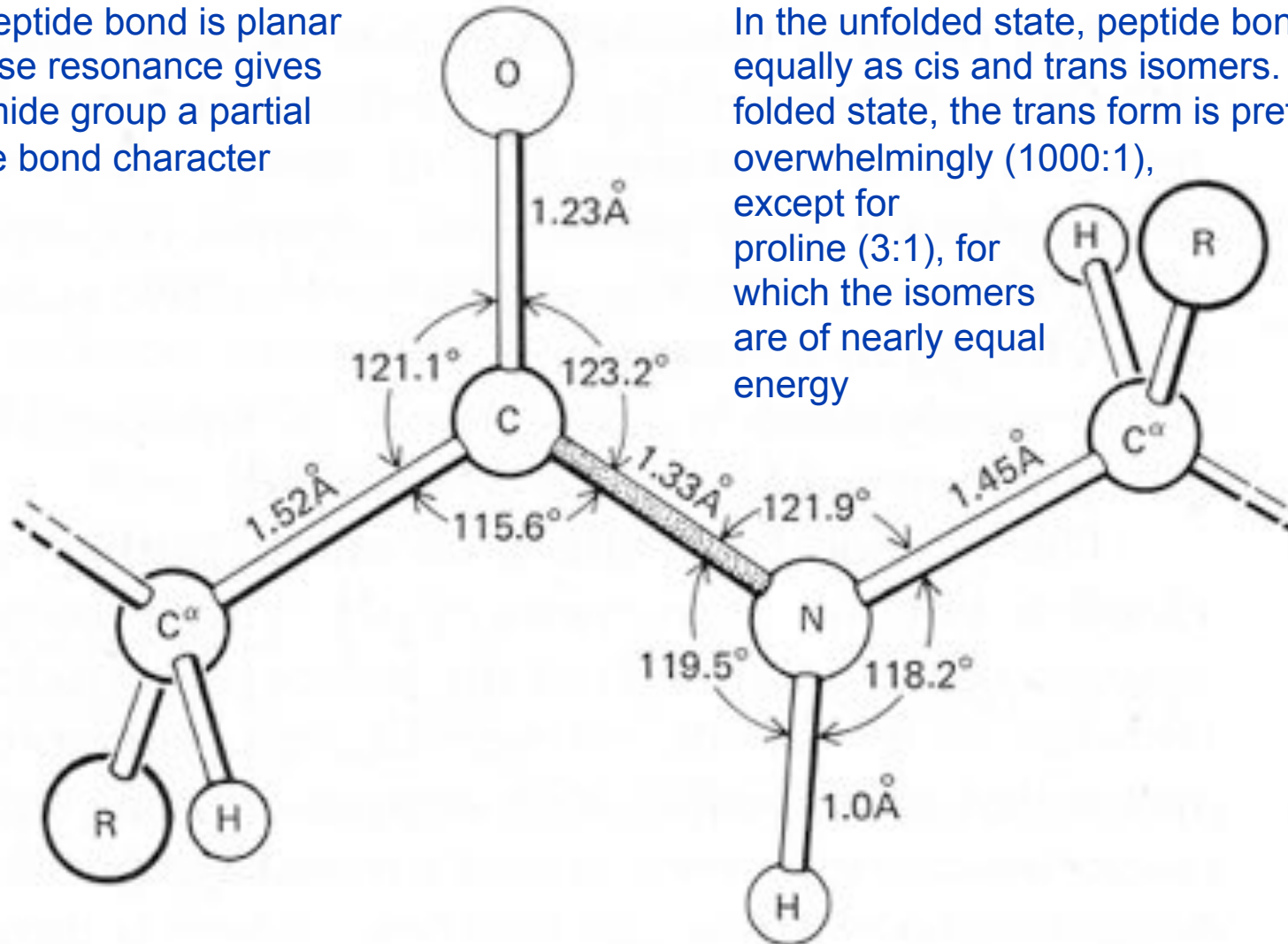


Primary structure - the peptide bond



The peptide bond is planar because resonance gives the amide group a partial double bond character

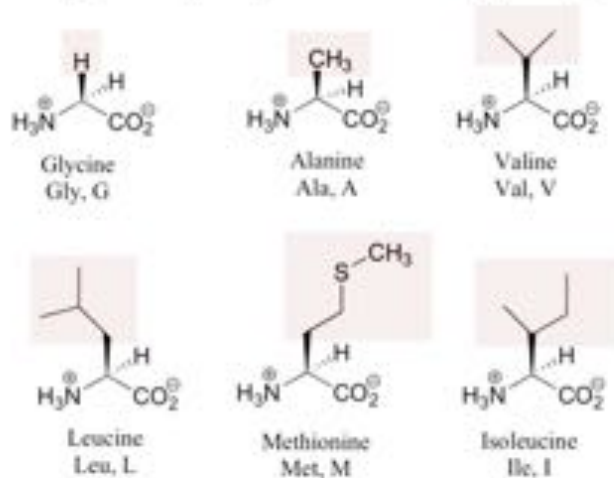
In the unfolded state, peptide bonds occur equally as cis and trans isomers. In the folded state, the trans form is preferred overwhelmingly (1000:1), except for proline (3:1), for which the isomers are of nearly equal energy



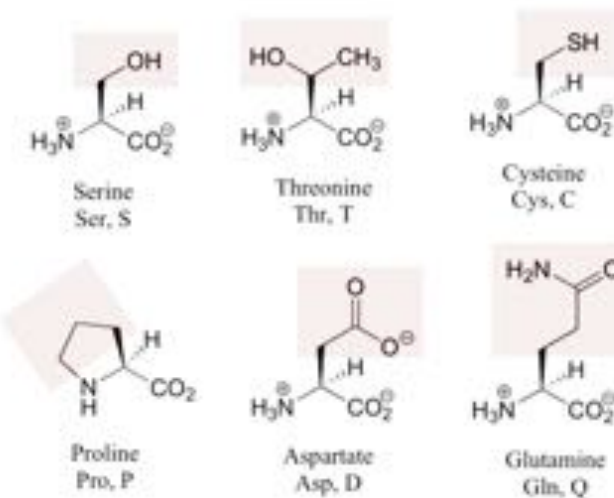
Primary structure - amino acid sidechains (= residues)



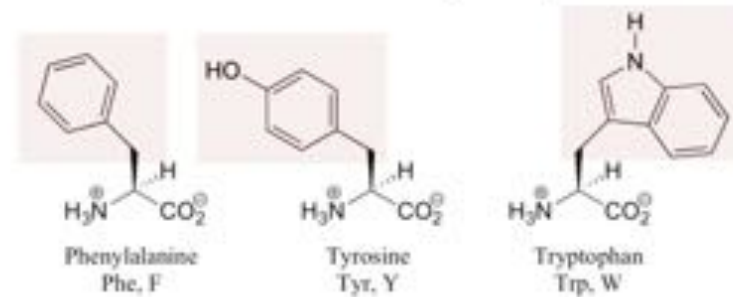
Nonpolar, aliphatic side groups



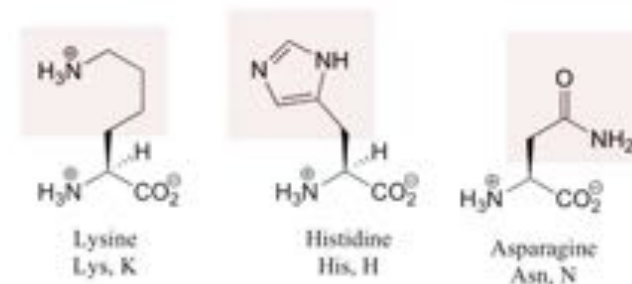
Polar, uncharged side groups



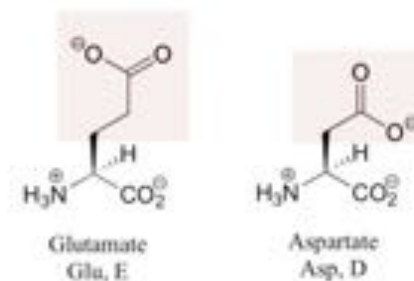
Aromatic side groups



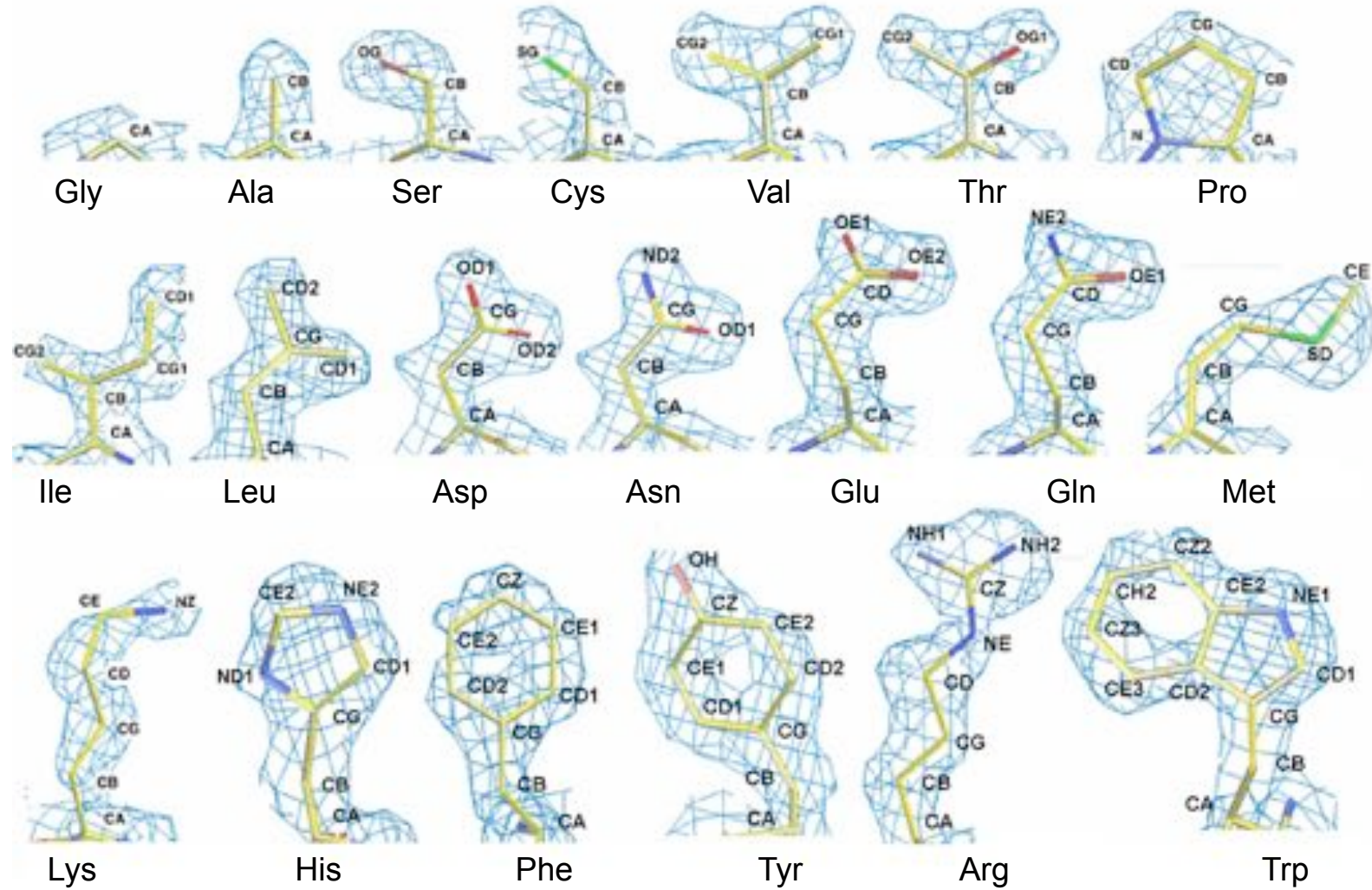
Positively charged side groups



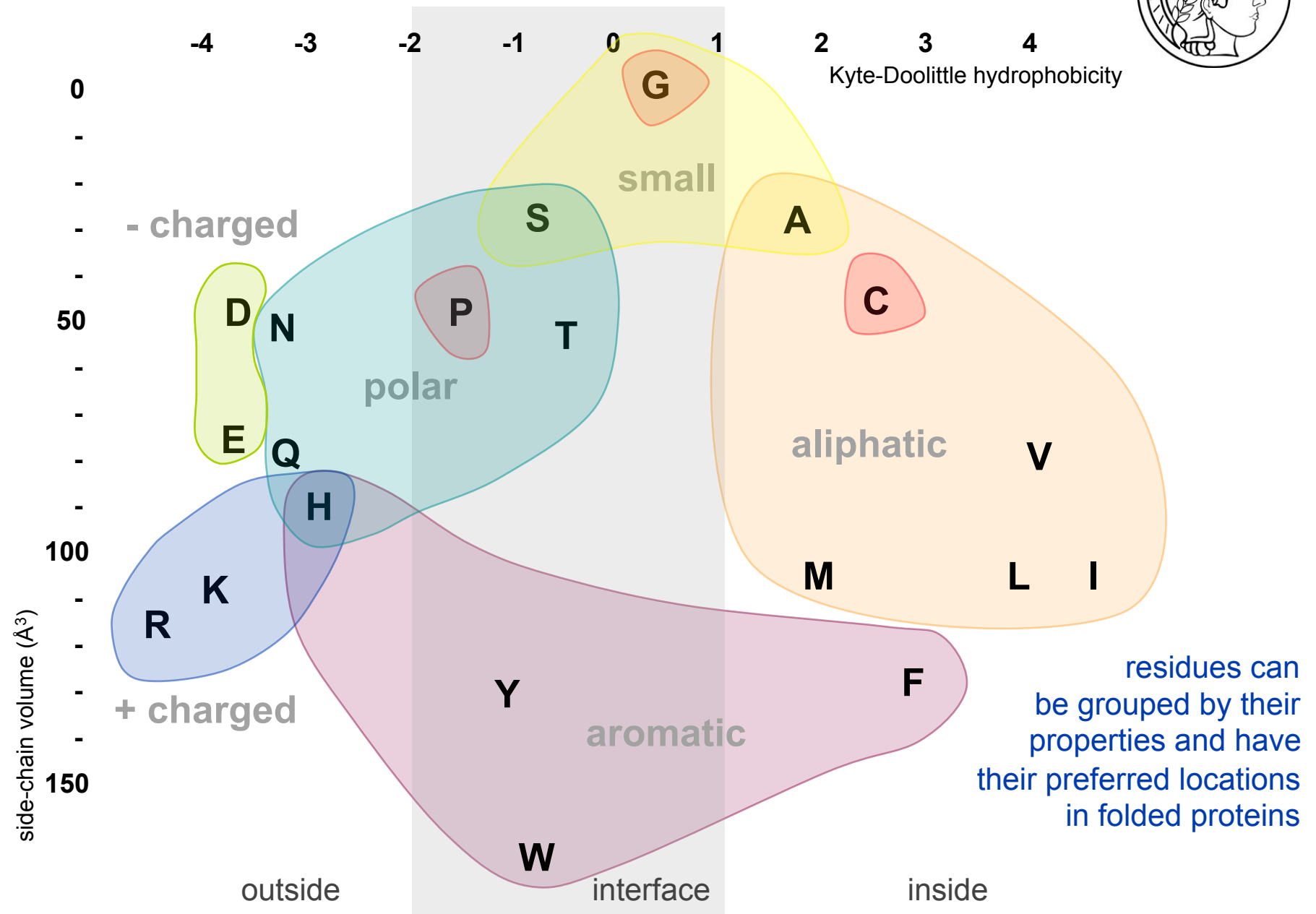
Negatively charged side groups



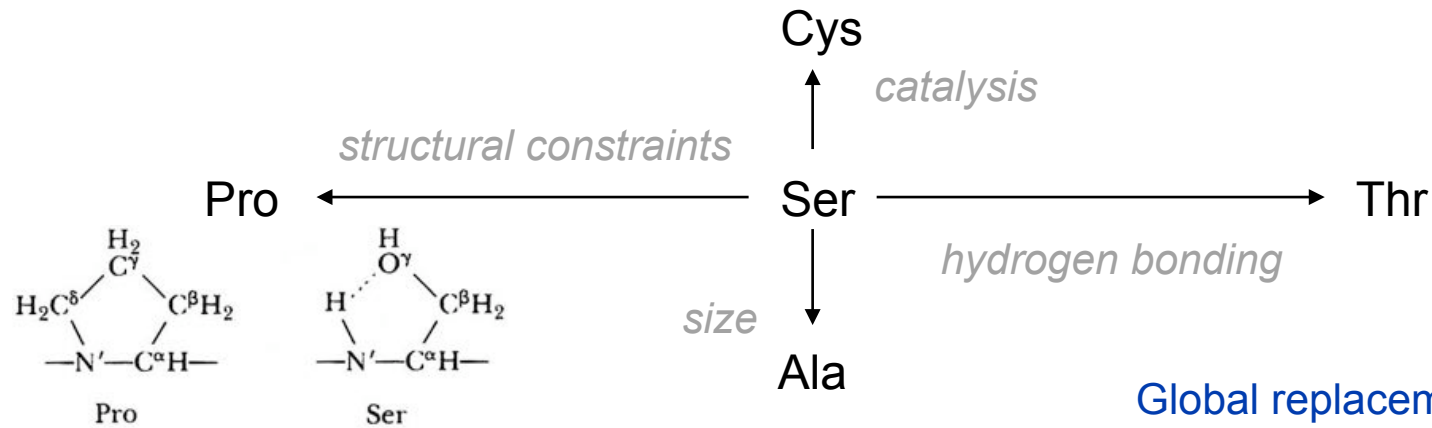
Primary structure - amino acid sidechains



Primary structure - amino acid sidechains



Primary structure - amino acid substitutions



	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	
C	9																				C
S	-1	4																			S
T	-1	1	5																		T
P	-3	-1	-1	7																	P
A	0	1	0	-1	4																A
G	-3	0	-2	-2	0	6															G
N	-3	1	0	-2	-2	0	6														N
D	-3	0	-1	-1	-2	-1	1	6													D
E	-4	0	-1	-1	-1	-2	0	2	5												E
Q	-3	0	-1	-1	-1	-2	0	0	2	5											Q
H	-3	-1	-2	-2	-2	-2	1	-1	0	0	8										H
R	-3	-1	-1	-2	-1	-2	0	-2	0	1	0	5									R
K	-3	0	-1	-1	-1	-2	0	-1	1	1	-1	2	5								K
M	-1	-1	-1	-2	-1	-3	-2	-3	-2	0	-2	-1	-1	5							M
I	-1	-2	-1	-3	-1	-4	-3	-3	-3	-3	-3	-3	-3	1	4						I
L	-1	-2	-1	-3	-1	-4	-3	-4	-3	-2	-3	-2	-2	2	2	4					L
V	-1	-2	0	-2	0	-3	-3	-3	-2	-2	-3	-3	-2	1	3	1	4				V
F	-2	-2	-2	-4	-2	-3	-3	-3	-3	-3	-1	-3	-3	0	0	0	-1	6			F
Y	-2	-2	-2	-3	-2	-3	-2	-3	-2	-1	2	-2	-2	-1	-1	-1	-1	3	7		Y
W	-2	-3	-2	-4	-3	-2	-4	-4	-3	-2	-2	-3	-3	-1	-3	-2	-3	1	2	11	W
	C	S	T	P	A	G	N	D	E	Q	H	R	K	M	I	L	V	F	Y	W	

Global replacement matrices (e.g. BLOSSUM62) assign the same value to a substitution irrespective of its location and have global penalties for indels

Position-specific scoring matrices (PSSMs) derive specific substitution values from positional frequencies in multiple alignments, but maintain global penalties for indels

Profile Hidden Markov Models (HMMs) derive position-specific substitution and indel values from multiple alignments

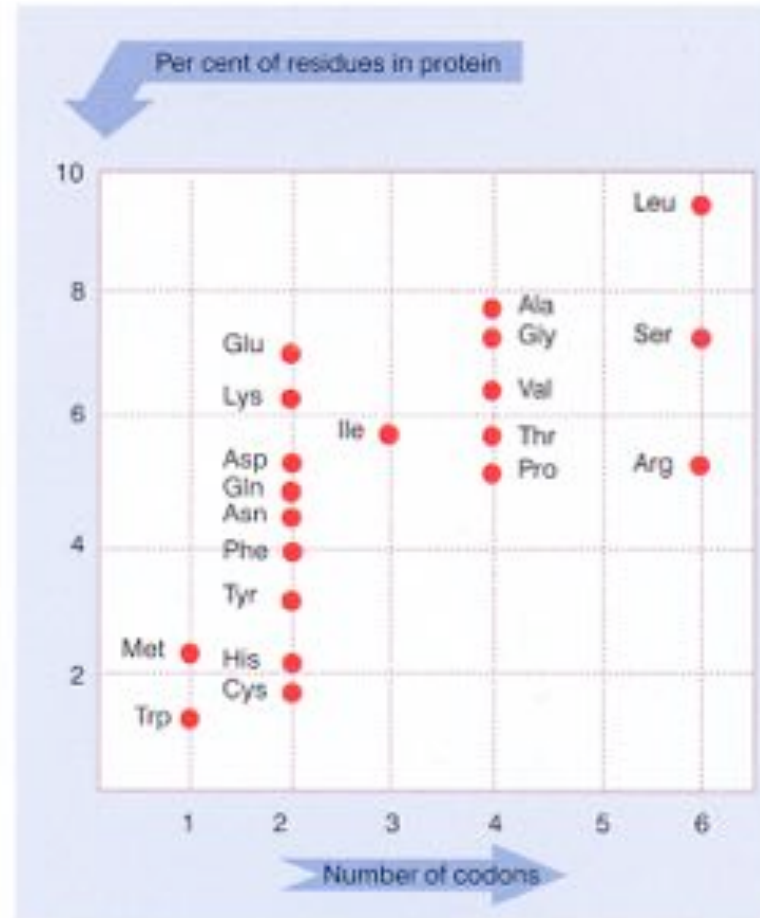
Primary structure - the linearity of the peptide chain originates from the linear nature of the gene



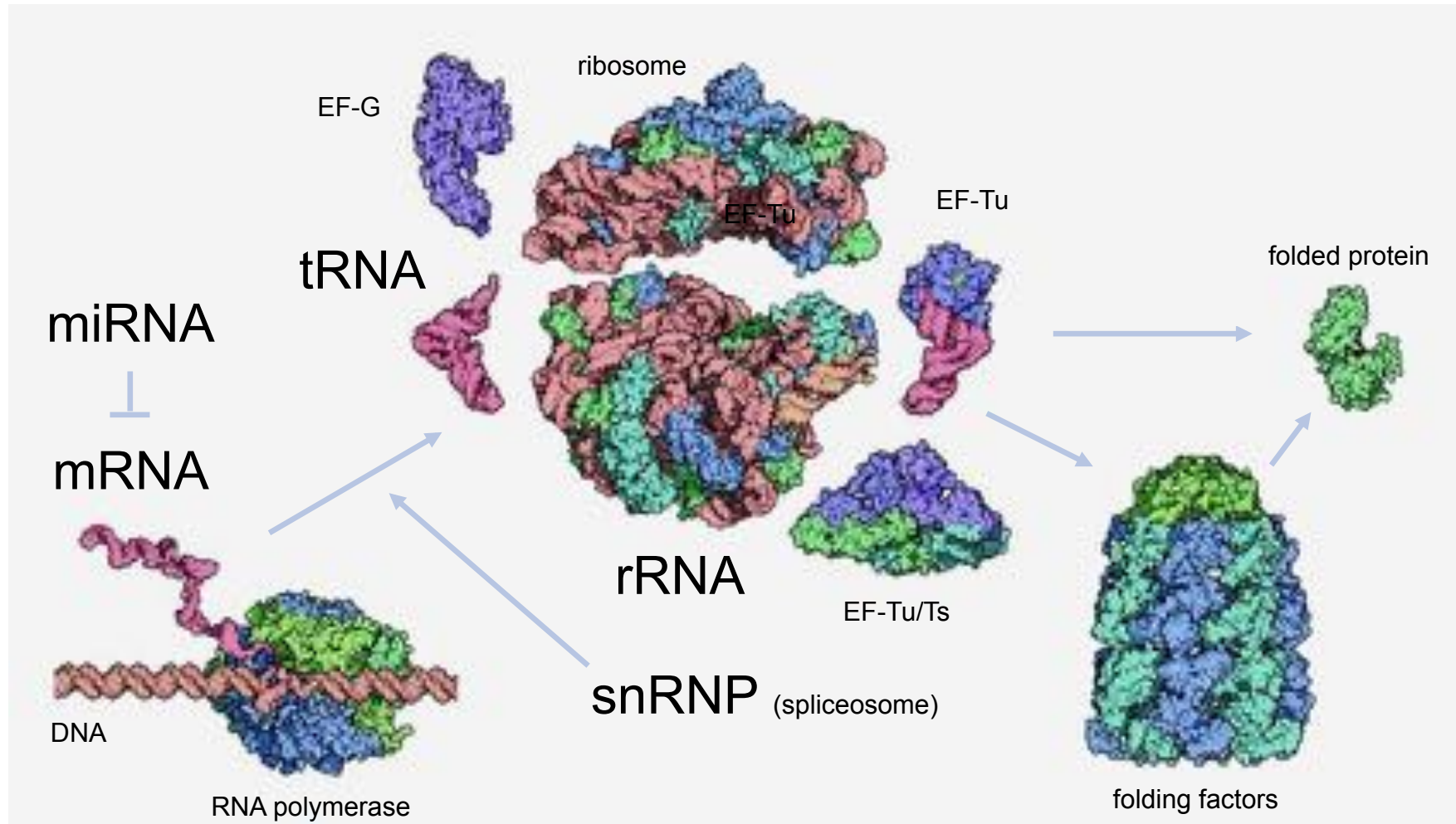
The DNA code

	U	C	A	G
U	UUU } Phe UUC } UUA } Leu UUG }	UCU } UCC } Ser UCA } UCG }	UAU } Tyr UAC } UAA } STOP UAG }	UGU } Cys UGC } UGA } STOP UGG } Trp
C	CUU } CUC } Leu CUA } CUG }	CCU } CCC } Pro CCA } CCG }	CAU } His CAC } CAA } Gln CAG }	CGU } CGC } Arg CGA } CGG }
A	AUU } Ile AUC } AUA } AUG } Met	ACU } ACC } Thr ACA } ACG }	AAU } Asn AAC } AAA } Lys AAG }	AGU } Ser AGC } AGA } Arg AGG }
G	GUU } GUC } Val GUA } GUG }	GCU } GCC } Ala GCA } GCG }	GAU } Asp GAC } GAA } Glu GAG }	GGU } GGC } Gly GGA } GGG }

Frequency of residues as a function of their number of codons

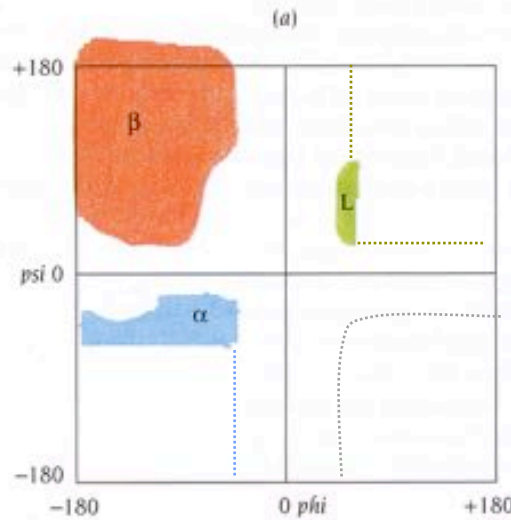
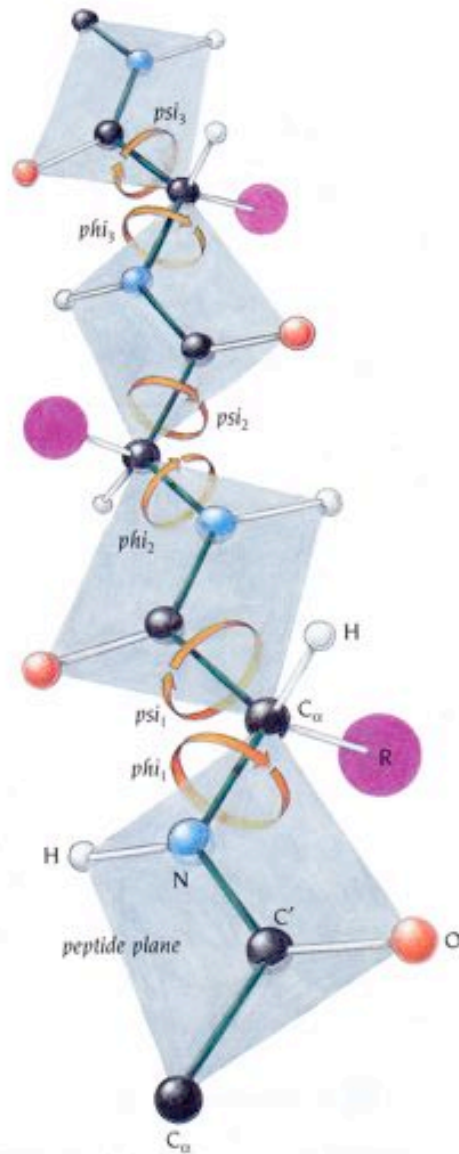


Primary structure - the central dogma



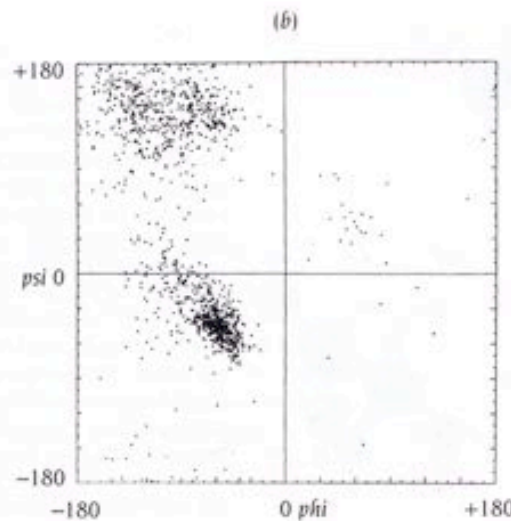
RNA molecules are the key players in converting genetic information into protein
→ remnants of an **RNA world** ?

Secondary structure - the Ramachandran plot



Because of its partial double bond nature, the peptide bond is always planar ($\omega = 180^\circ$ (trans) or 0° (cis)).

Rotation is only possible around the N-C α (ϕ) and C α -C' (ψ) bonds.



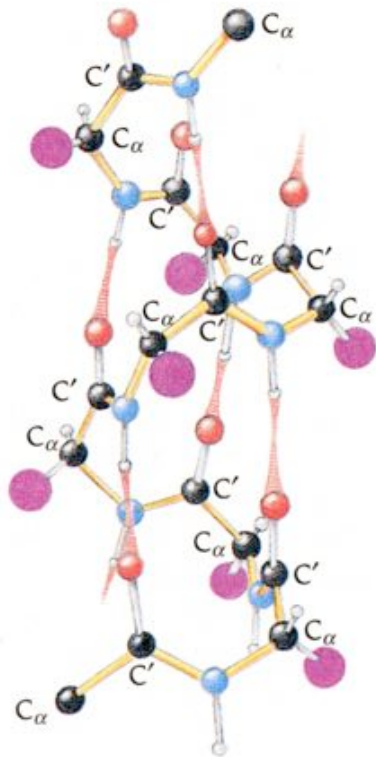
The Ramachandran plot shows the ϕ and ψ angles that can be assumed by the peptide chain.

Because of the bulkiness of the C β carbon, angles of $\phi > 80$ and $\psi < -80$ are severely disfavored, except for glycine

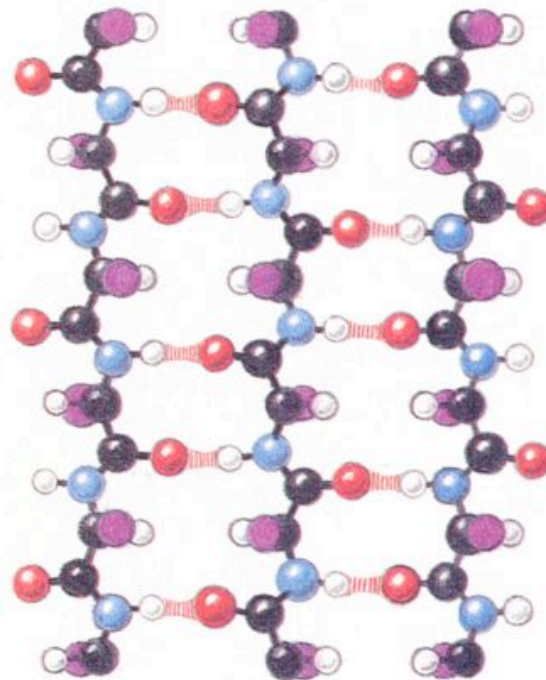
Secondary structure - helices and sheets



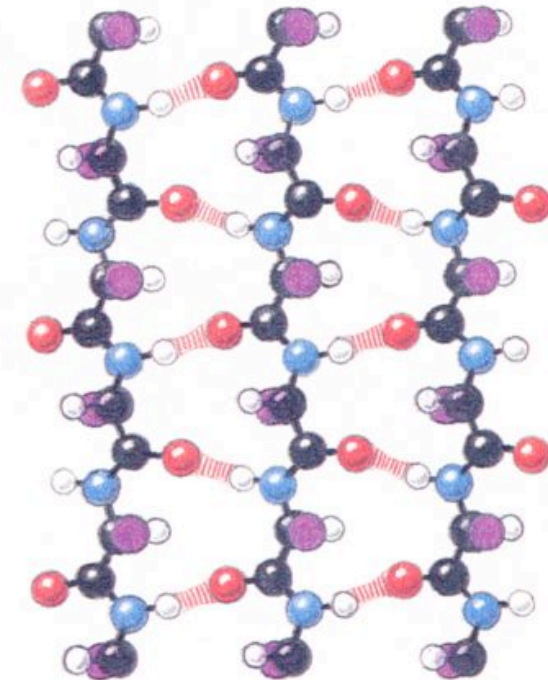
α -helix



antiparallel β -sheet



parallel β -sheet



Secondary structure - prediction: Quick2D



HOPE | Login | FDSAlert | Personal Databases | Contact | Tool Versions | Imprint | Disclaimer | Help

Bioinformatics Toolkit
Max-Planck Institute for Developmental Biology

Search | Alignment | Sequence Analysis | **Secondary Structure** | Secondary Structure | Classification | Utilities

Quickfinder

Quick2D - Results Job-ID: 4934381 Date: 12:39 on Sep 15 2012 [Help](#)

Save Export

gi|145534|gb|AAA23577.1| CheY [Escherichia coli]

```

      50
      MADRILKPLVSDPFTNRKTVWKLLEKLGPNVVFRAEDGVSALEKKGAGCYCPVTEINRMPHMDCLLLFTTRAGAMSA
SS PSIPKED  [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS JNET  [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS Prod (Ovall) [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS Prod (Root) [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
CC Colla
TH TM0000P
TH N000A1-SVM
TH P00000
DO DISOPRED0
DO ITPRED
SO Prod (Root)
SO JNET
  0 1000 0 0 20 20 0 0 0 0 20 0 0 1000 0 0 0 0
  0 100000 10 10 10 0 0 0 10 0 10 10 1000000 0 1000 0 0
  
```

```

      100
      LPVLRVYVZAKKSTIAAGAGAGYVVEPPTAATLSEKLERIFKLGK
SS PSIPKED  [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS JNET  [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS Prod (Ovall) [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
SS Prod (Root) [alpha-helix] [beta-strand] [beta-strand] [beta-strand] [beta-strand] [beta-strand]
CC Colla
TH TM0000P
TH N000A1-SVM
TH P00000
DO DISOPRED0
DO ITPRED
SO Prod (Root)
SO JNET
  0 100 10000 100 0 0 0 0 0
  100000 10 10 1000000 0 10 10 10 0
  
```

SS = **alpha-helix** **beta-strand** Secondary structure
 CC = **Coiled coil**
 TH = **Transmembrane** ('='-outside, '='-inside)
 DO = **Disorder**
 SO = **Solvent accessibility** (A buried residue has at most 25% of its surface exposed to the solvent.)

Secondary structure - prediction: Ali2D



HOME | Login | Feedback | Personal Databases | Contact | Tool Versions | Imprint | Disclaimer | Help

Bioinformatics Toolkit

Max-Planck Institute for Developmental Biology

Search Alignment Sequence Analysis **Secondary Structure** Secondary Structure Classification Utility

Ali2D [Help](#) Quick2D

[New job](#) [Resubmit](#) [Log](#) [Input-params](#) [Delete](#)

Ali2D - Results

Job-ID: 1463618 Date: 12:45 on Sep 15 2012 [Help](#)

Results Alignment Viewer Applet **Colored Results**

Colored Alignment with Confidence Colored Alignment Black/White Alignment

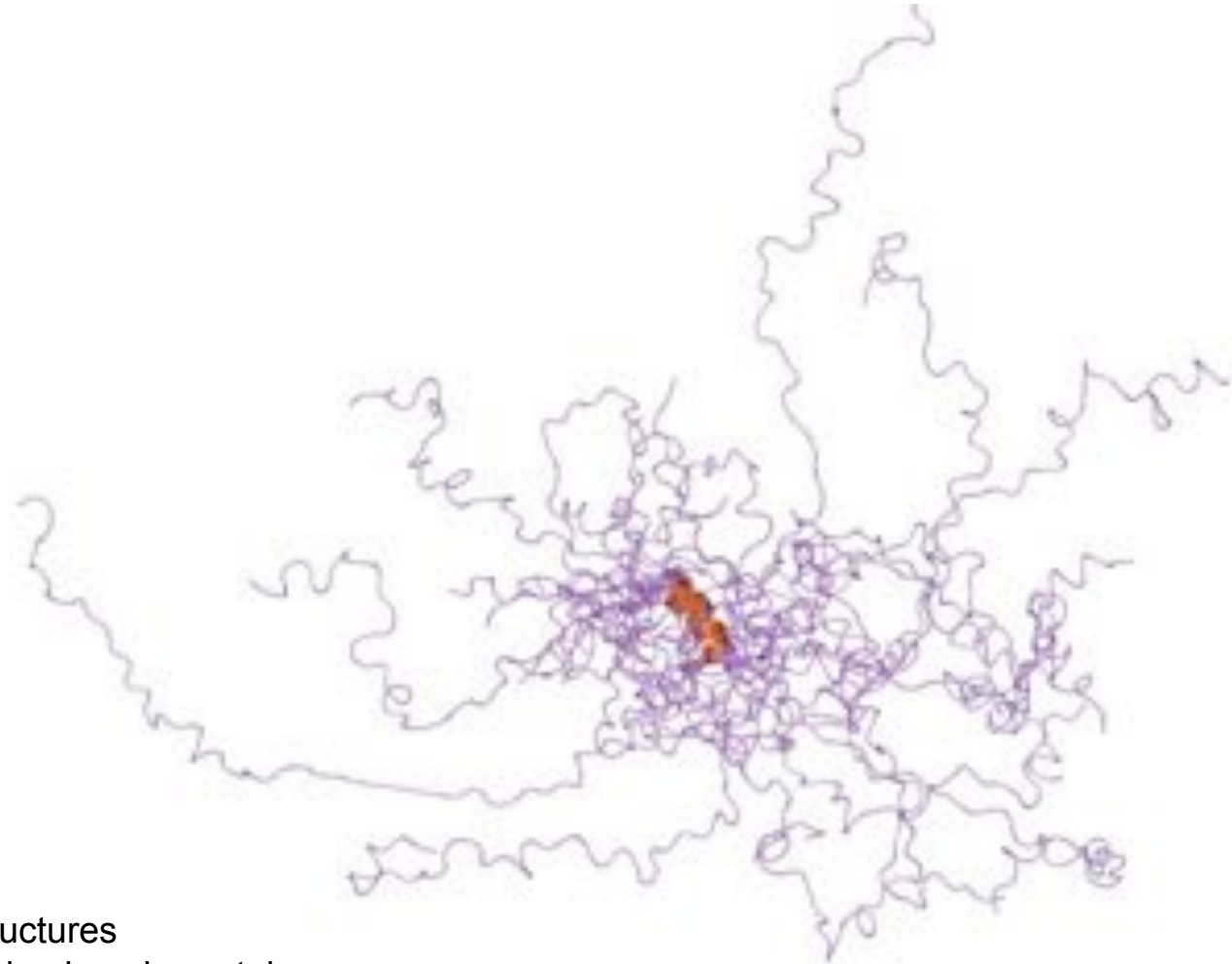
```
0000 ----TPOIGSRVATLLEISAGS-----GTTTSSSRRTKSLTQHPFLSD----WTFPLETIP--WVFAFSE-SGRRTPLKASLRKGLLDMVSLQD
0001 -LIEVHTLRQVAGSLKRVVTFN-----TLFRRKPRKLPQRRKTLIE-----VLSLTKRQD--LQTERKQD--RDRFPRLAPLVTTLQISPLRQD
0002 VTRKIDRIGKLTIFLLRFTDQD-----LPLKTLKQVQSLQVSKAKRQV--WLESLCCVQ--GIGVKAQD--PQIRSLQMLEIKKIMQDQTYNKAQ
0003 WACPLGIDSRVYSLLRKPKK-----TRPFRFTFAYDQKCRGCRKQVQ--TKLTKLPLSALS-WVQGTACQD--SRAGTTLKNS
0004 --LDIRKPTISGLRPLAKQNS-----FVLLKQLPLVFDQKTPSTQVQKGLIDISGLVSKNIVTLDQA--QVQRKTLSSKPKKQTRSEL
0005 --LQCVYFAGMVELFQNLQI-----TLLPGRVYRSTYRATQCRLESDVFLAAGQSDAIE--GLVDRKLT--QDREKTVLQF
0006 QVMSVYFAGQCSLSEKATQD-----SSQKRLQGLSLRQKCFQDQRTLVQVQATAVELLP--TVVLRKQD--QATKSLQKNS
0007 --RFTKLCYVQYTKSRQPSASVYKIVRSEVELTRQQVYVPIKSLRQK--AKYKQSL-----KALLANQD--RLEVYQD--GIDKSL
0008 --TRKRTIILVQGLDQKPK-----YVREKSLKCOLSTTRKLPF-----LISVVADE--LPLIKQPTQKQEVSDMS
```

© 2008, Dept. of Protein Evolution at the Max-Planck Institute for Developmental Biology, Tübingen Release-2.17.0

Secondary structure - natively unstructured proteins



Sometimes proteins lack secondary structure.



Ensemble of NMR structures
of the Thylakoid soluble phosphoprotein
TSP9, which shows a largely flexible protein chain (PDB 2fft)

Secondary structure - natively unstructured proteins



Sometimes proteins lack secondary structure. Few proteins, however, are entirely unstructured. Usually, proteins contain unstructured segments between folded regions.

Consensus secondary structure prediction for GW182, a protein with large unstructured segments



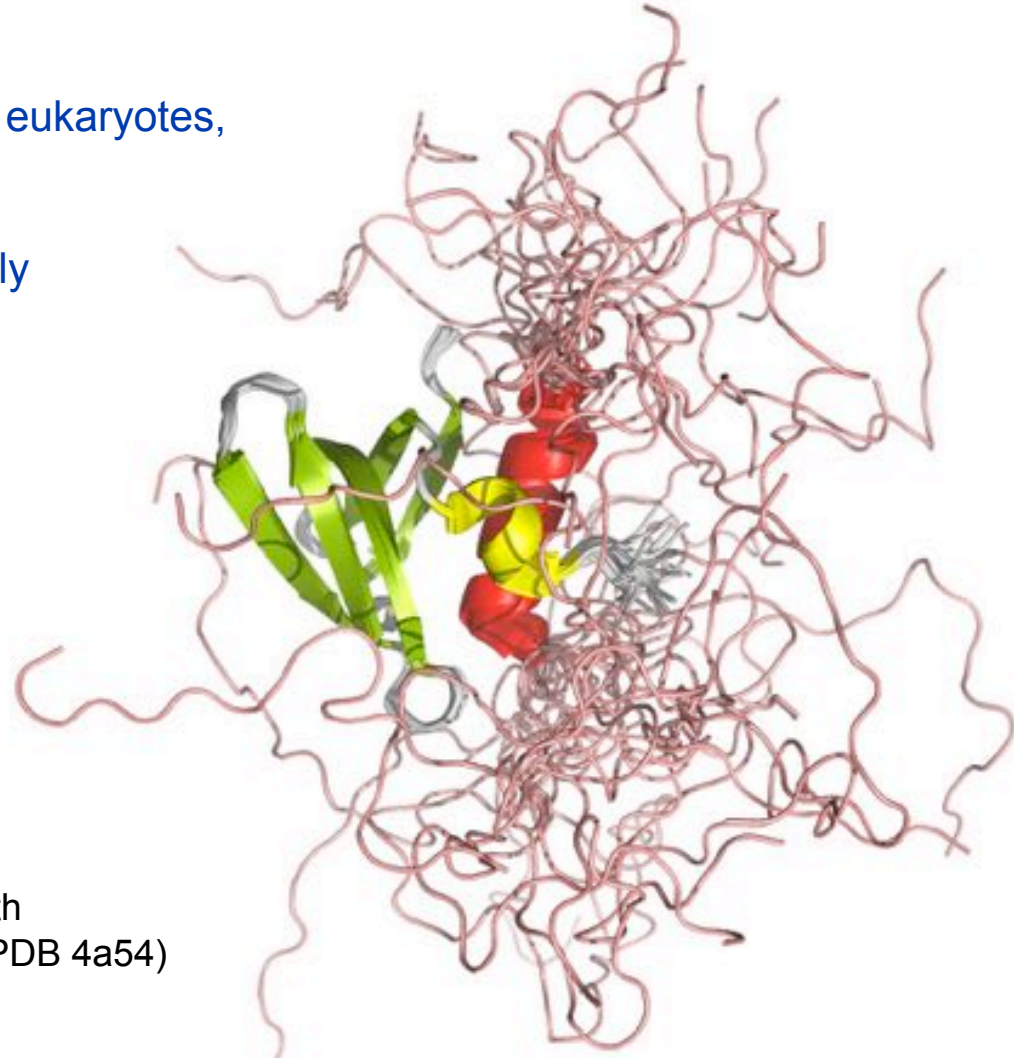
Secondary structure - natively unstructured proteins



Sometimes proteins lack secondary structure. Few proteins, however, are entirely unstructured. Usually, proteins contain unstructured segments between folded regions.

Unstructured segments are common in eukaryotes, but rare in prokaryotes.

Unstructured segments almost invariably assume a defined structure in complex with a structured partner.

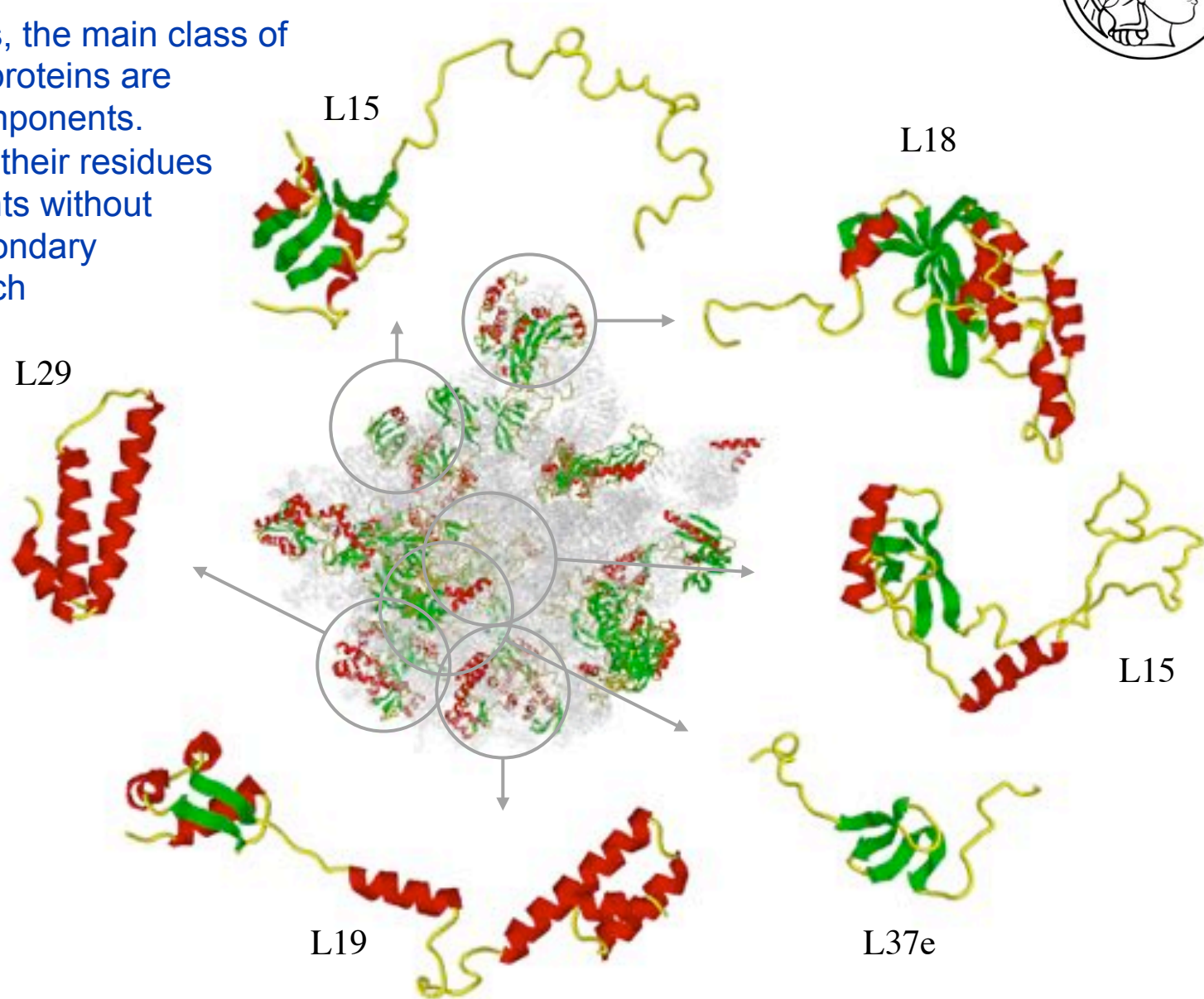


yeast Edc3 LSm domain in complex with a leucine-rich motif (HLM) from Dcp2 (PDB 4a54)

Secondary structure - natively unstructured proteins



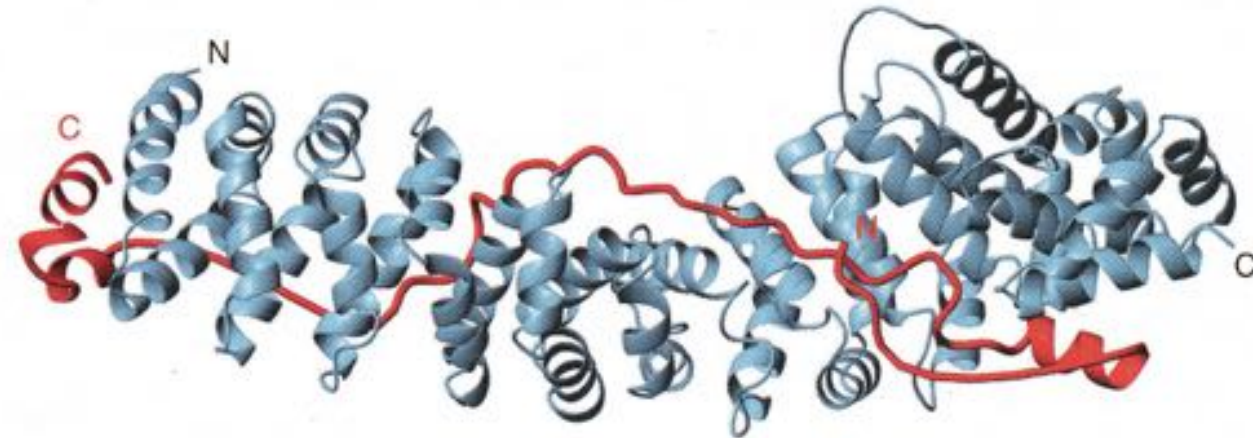
In prokaryotes, the main class of unstructured proteins are ribosomal components. About 60% of their residues are in segments without a defined secondary structure, which in complex with the RNA.



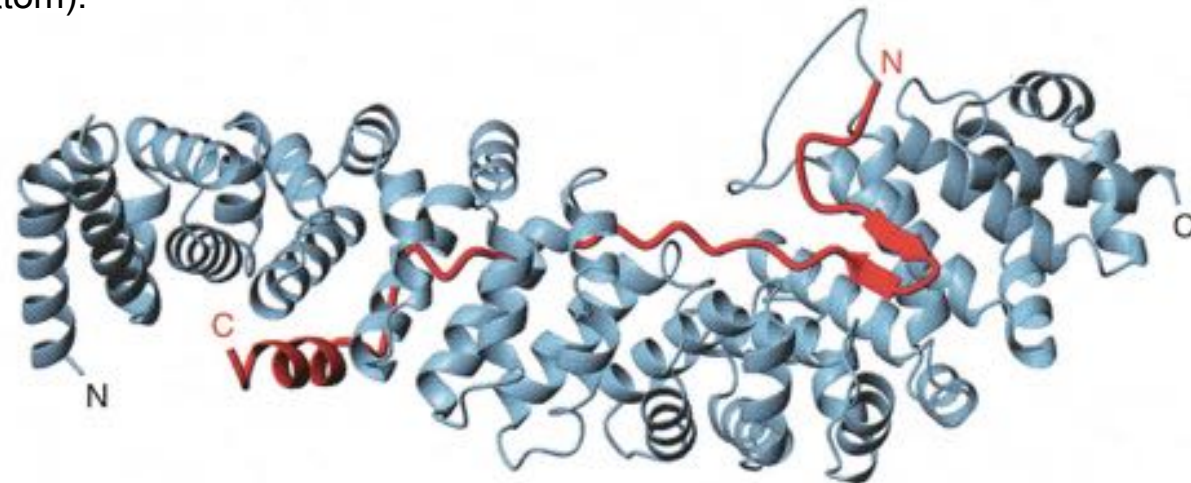
Secondary structure - natively unstructured proteins



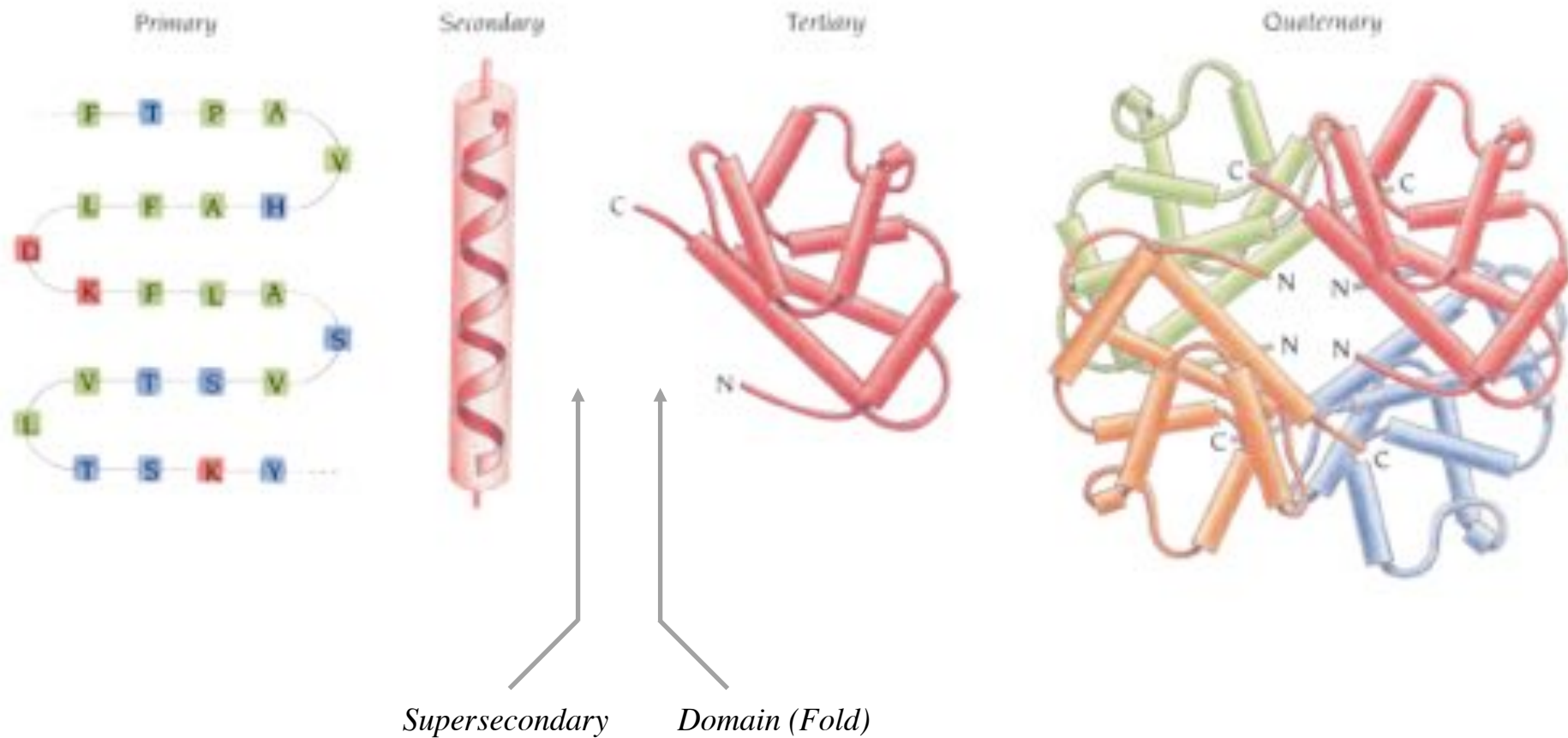
In eukaryotes, unstructured segments are the main location for linear motifs, short sequences that mediate the specific association to folded ('scaffold') proteins.



The armadillo repeat domain of β -catenin in complex with E-cadherin (top) and Tcf3 (bottom).



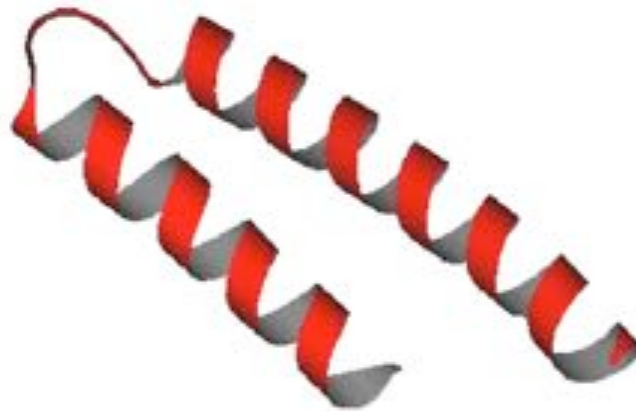
From secondary to tertiary structure



The main supersecondary structure elements



$\beta\beta$ -hairpin



$\alpha\alpha$ -hairpin

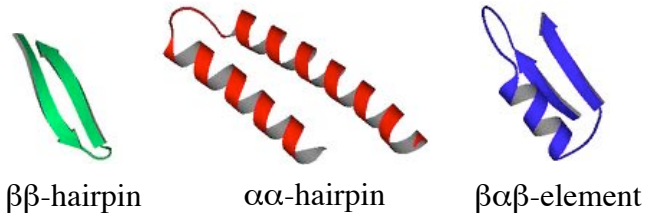


$\beta\alpha\beta$ -element

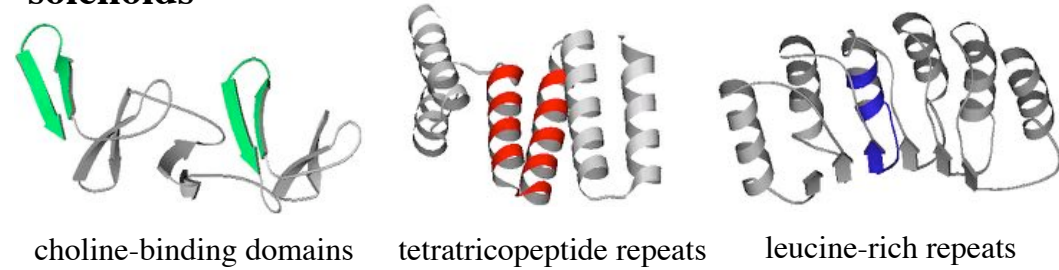
Supersecondary structures are at the core of domains - the hierarchy of fold complexity



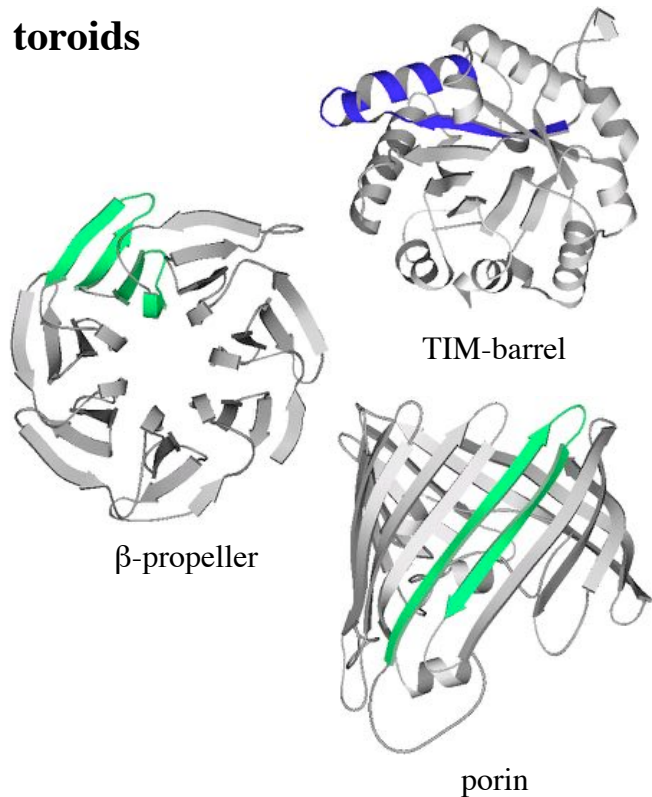
supersecondary structures



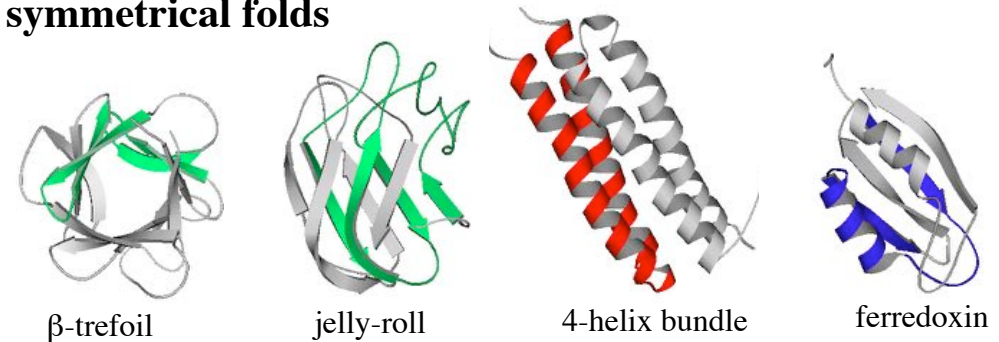
solenoids



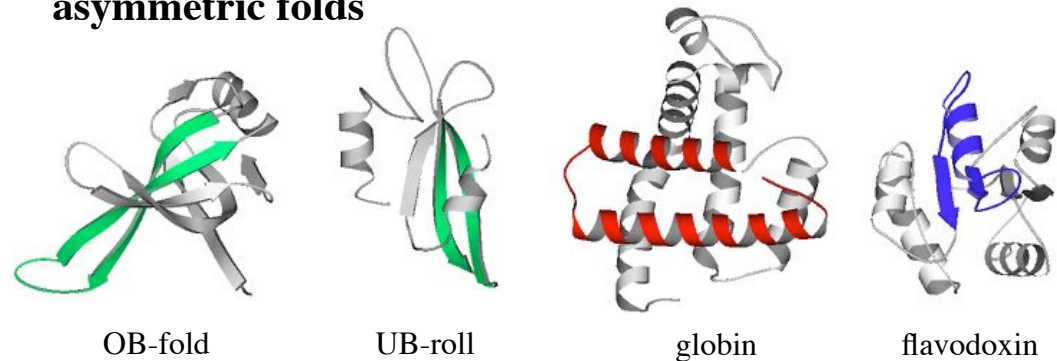
toroids



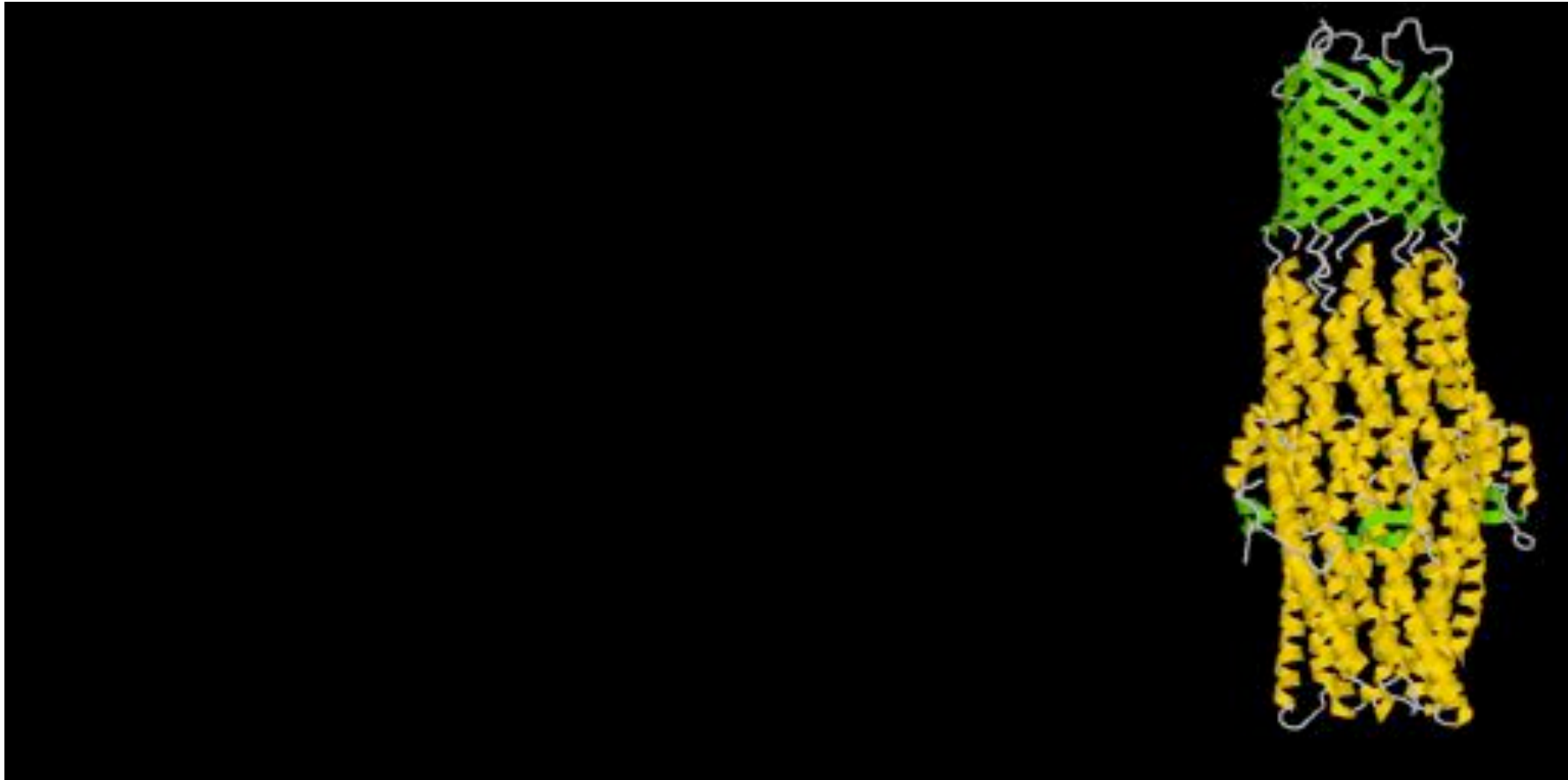
symmetrical folds



asymmetric folds

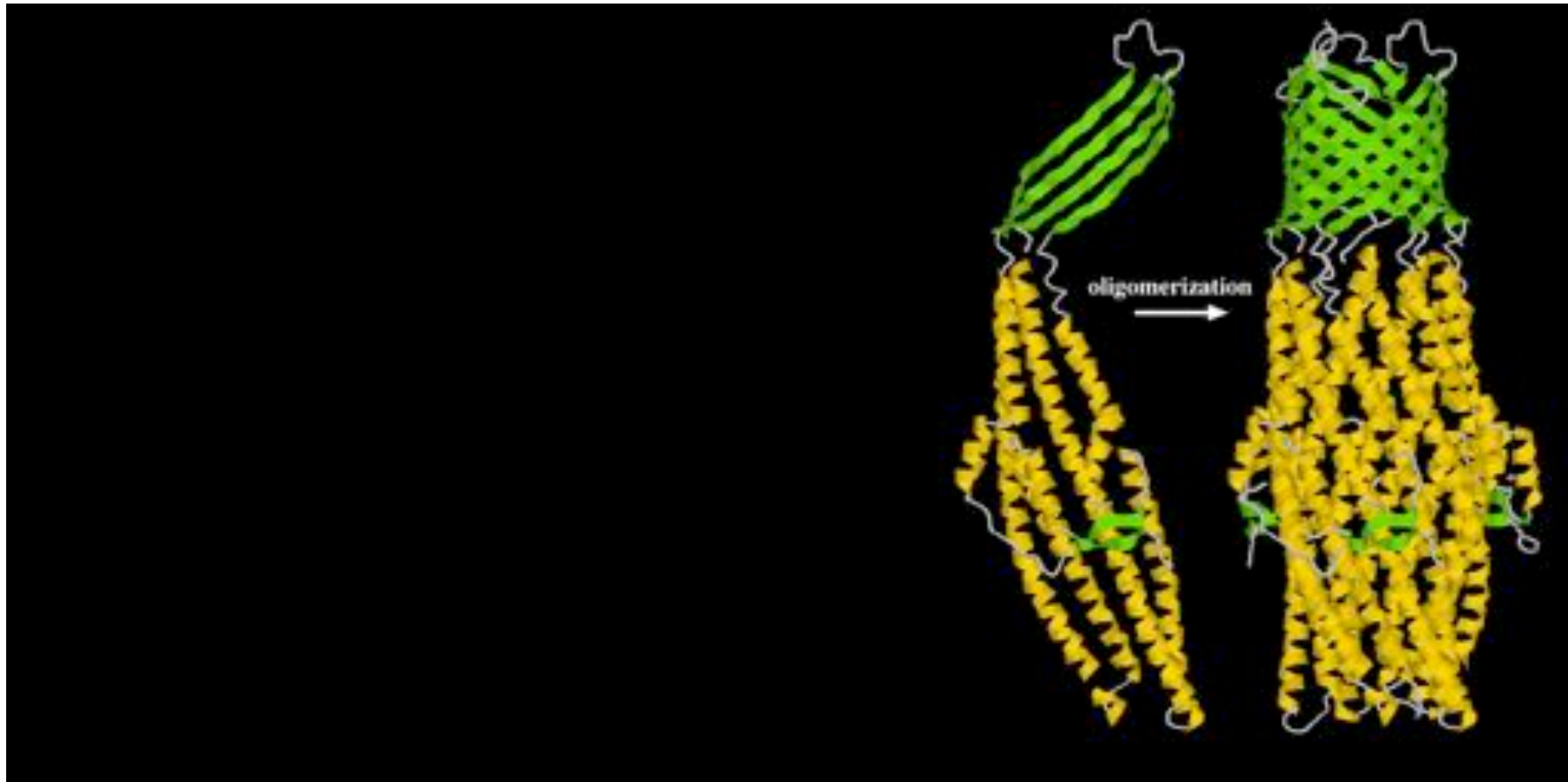


Deconstructing a protein into supersecondary structures



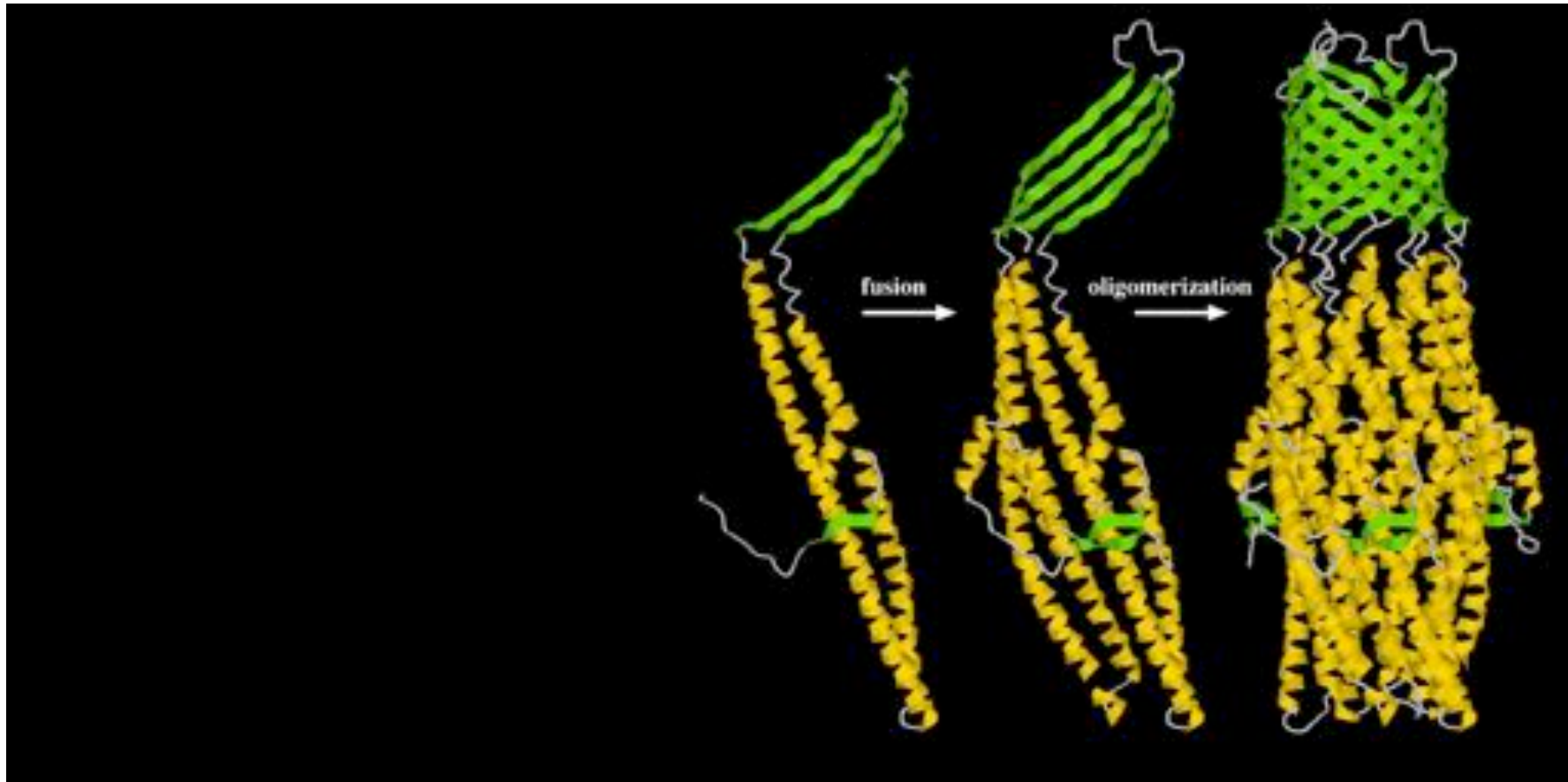
TolC

Deconstructing a protein into supersecondary structures



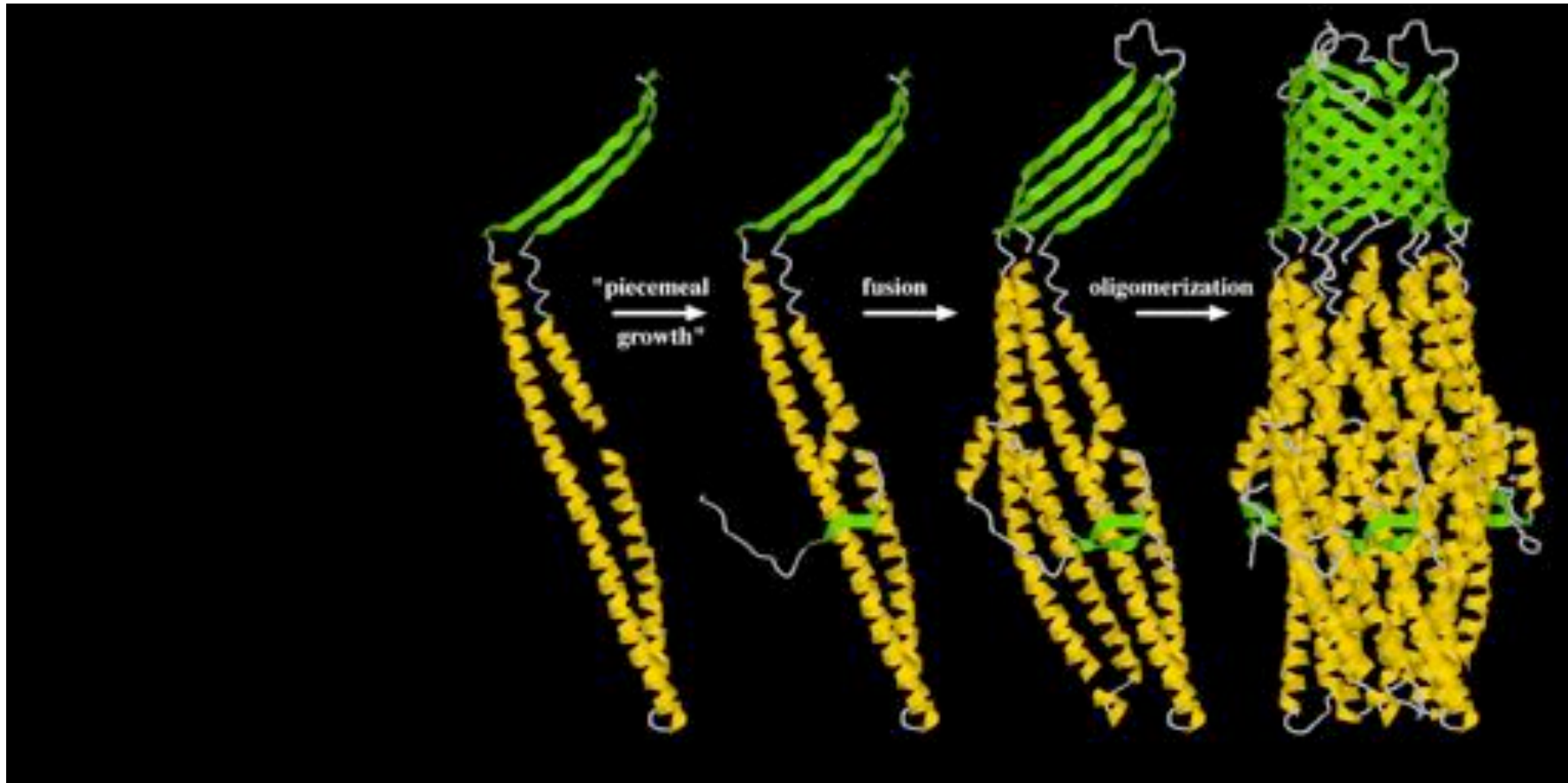
ToIC

Deconstructing a protein into supersecondary structures



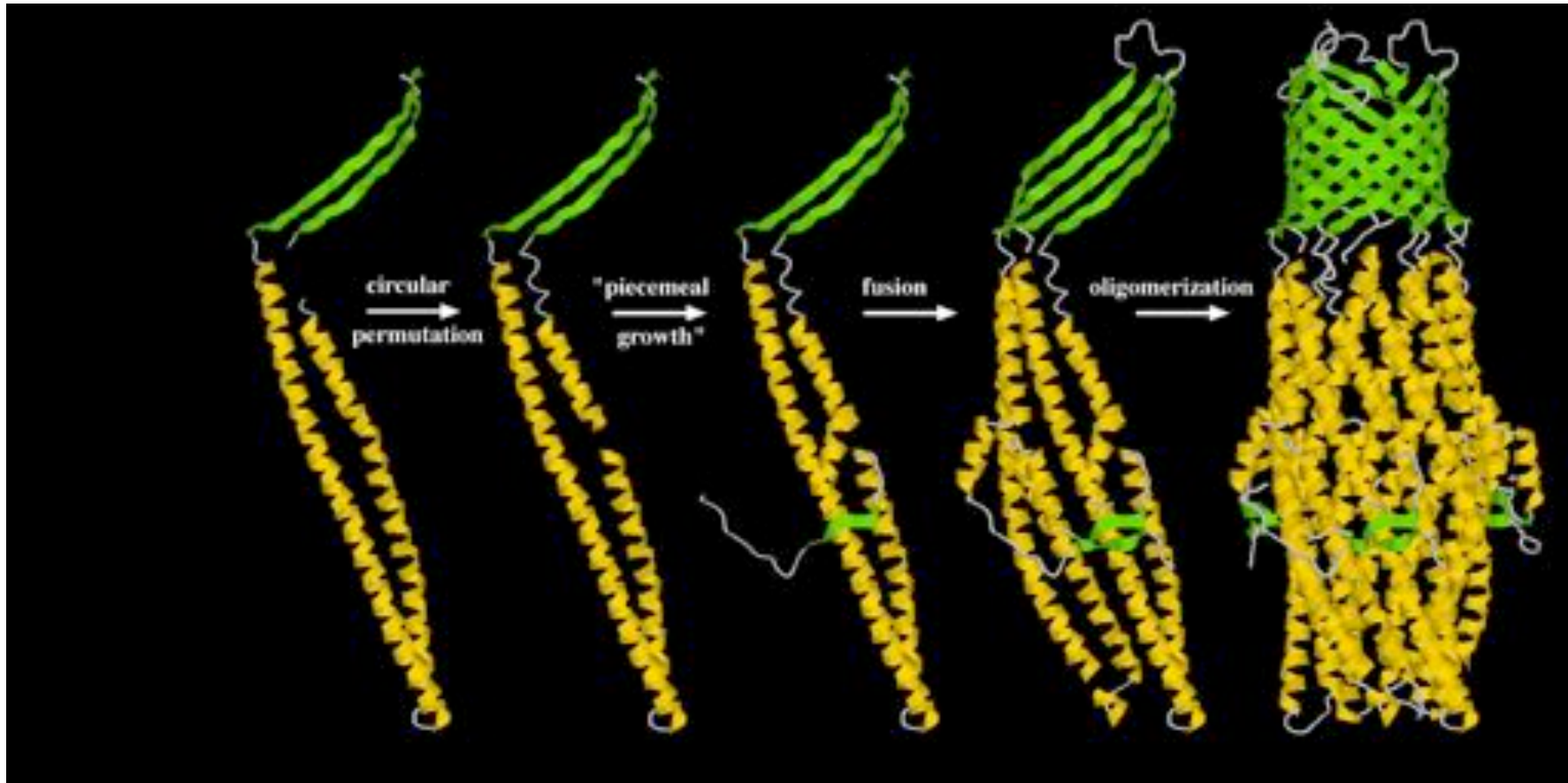
ToIC

Deconstructing a protein into supersecondary structures



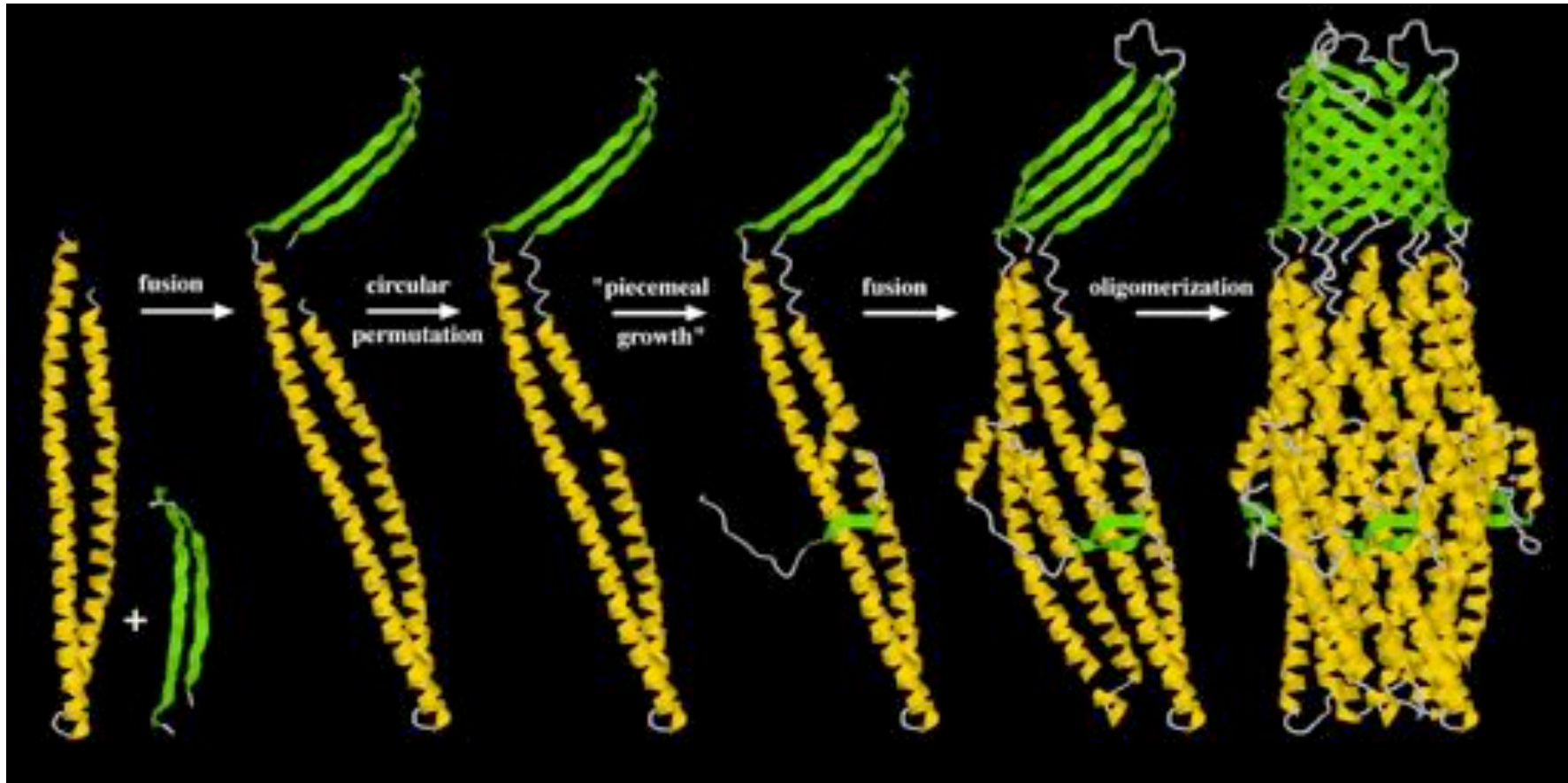
ToIC

Deconstructing a protein into supersecondary structures



ToIC

Deconstructing a protein into supersecondary structures



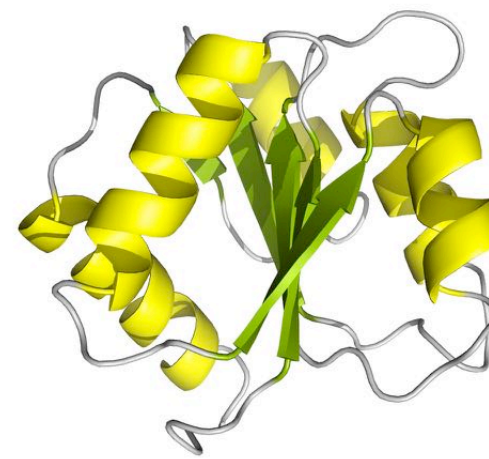
ToIC

Tertiary structure - the fold classes

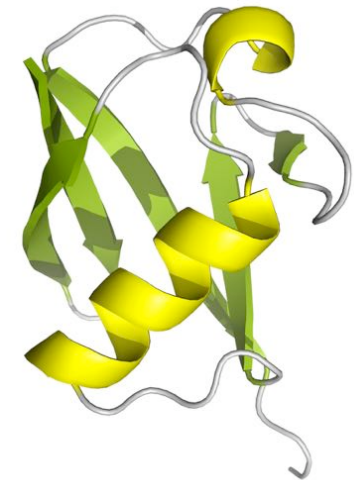


α (formed of α -helices)

β (formed of β -strands)



α/β (formed of alternating α -helices and β -strands; strands mainly parallel)

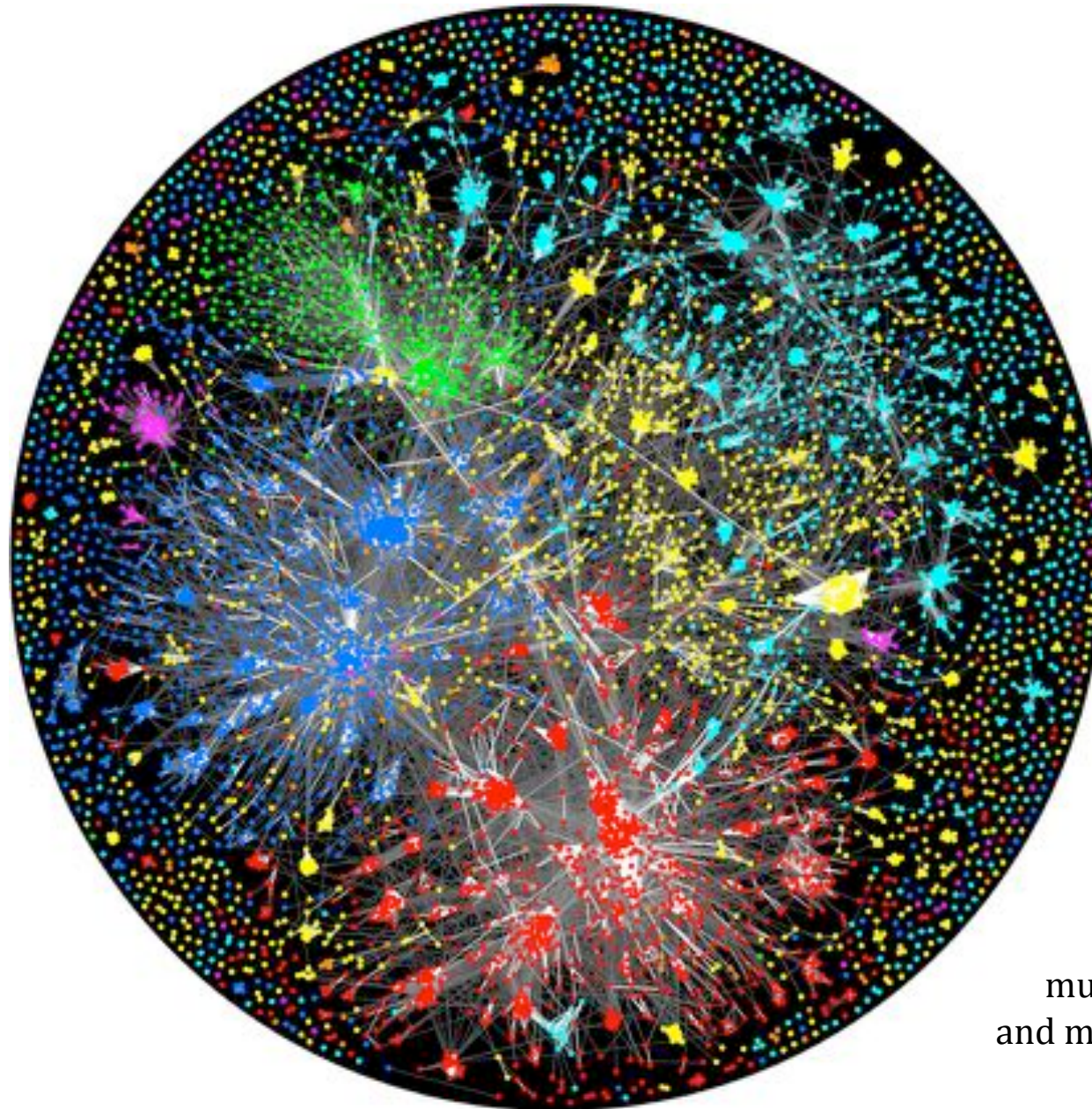


$\alpha+\beta$ (formed of irregular arrangements of α -helices and β -strands; strands mainly antiparallel)

Tertiary structure - the fold classes

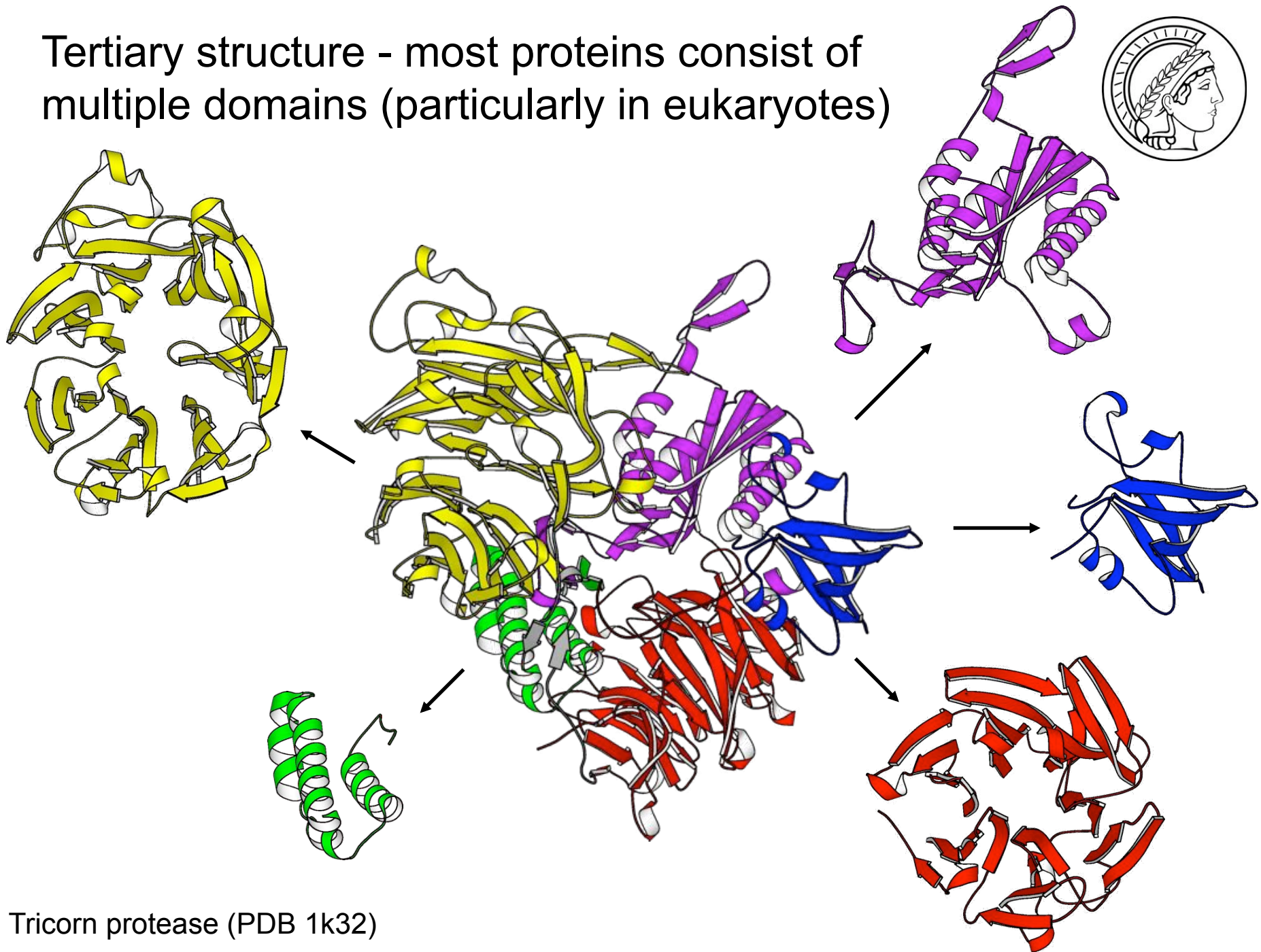


Structure classes have preferred residue composition



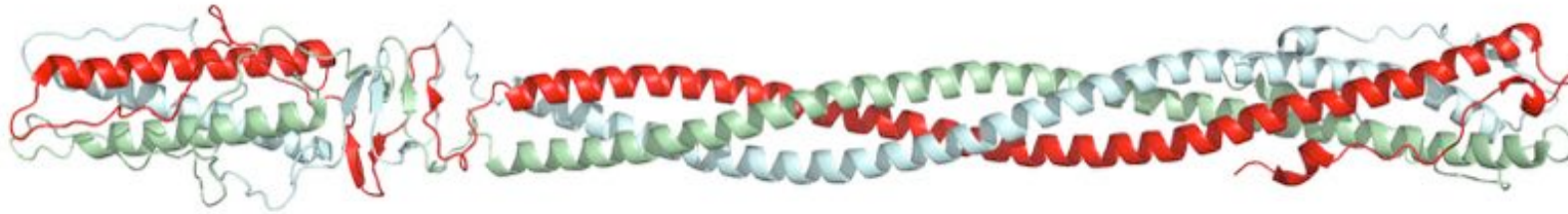
Proteins of known structure, clustered by sequence similarity and colored by structural class: all- α (blue), all- β (cyan), α/β (red), $\alpha+\beta$ (yellow), small proteins (green), multi-domain proteins (orange), and membrane proteins (magenta).

Tertiary structure - most proteins consist of multiple domains (particularly in eukaryotes)

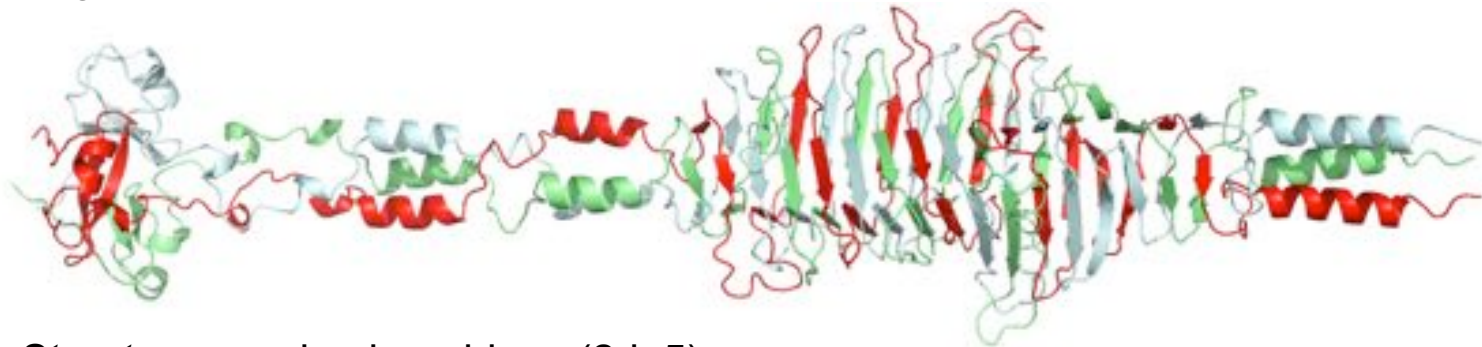


Tricorn protease (PDB 1k32)

From secondary to quaternary structure - protein fibers



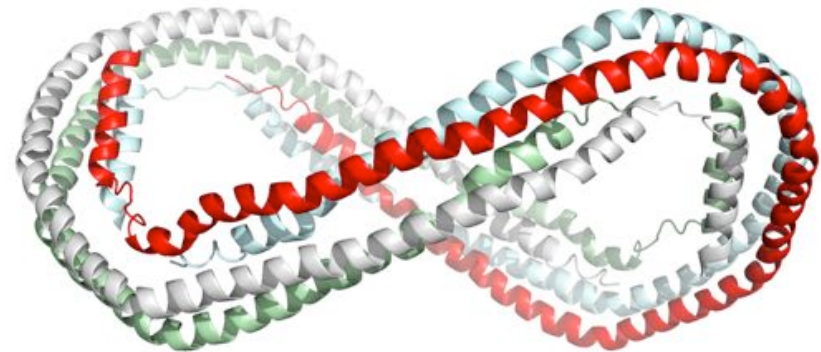
Salmonella phage P22 tail spike protein (2poh)



Streptococcus hyaluronidase (2dp5)

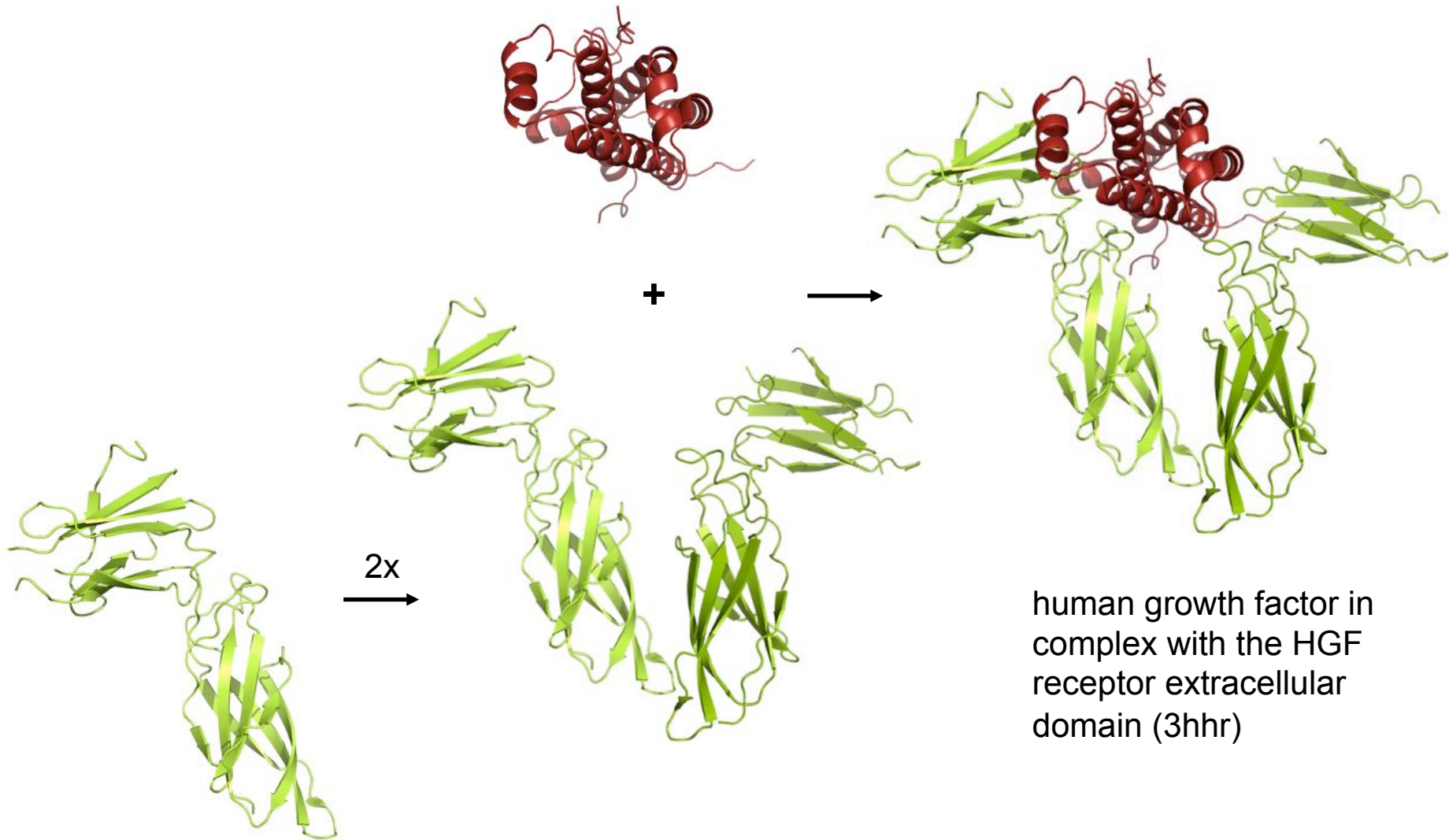


BadA adhesin head segment (3d9x)



apolipoprotein A-I (1av1)

Quaternary structure - homo- and hetero-oligomers



Quaternary structure - complex assemblies

