

# Structural bioinformatics

## KFC/STBI

What is structural bioinformatics?

Karel Berka  
Miroslav Krepl

# Requirements

- **Project:**
  - Structure analysis, docking, comparison of proteins, prediction of properties from structure, ...
  - 1(max. 2) page-long report with
    - Hypothesis
    - Brief Methodology
    - Conclusions

ev. ChannelsDB – doplnění 5 struktur do databáze
- **Exam:**
  - Project-like Questions – problem + discussion about its possible resolution from you side

# Content

- Structural bioinformatics, Biomolecules, Structural hierarchy
- Structure determination (X-Ray,NMR,EM), Structure file formats
- Structural databases (PDB, CATH, SCOP, Drugbank)
- Visualization of structure, structural alignment
- Structure prediction, CASP, AlphaFold ML revolution
- Function prediction, CASA
- Binding prediction – protein-ligand and protein-protein docking
- Challenges of structural bioinformatics - membrane proteins, nucleic acids, protein-protein interactions prediction
- Examples: SARS-CoV-2, Switchable proteins

# Bioinformatics

(*Molecular*) **bio** – informatics: bioinformatics is conceptualising **biology in terms of molecules** (in the sense of physical chemistry) and applying "**informatics techniques**" (derived from disciplines such as applied maths, computer science and statistics) to understand and **organise** the **data and information** associated with these molecules, on a **large scale**.

In short, bioinformatics is a management information system for molecular biology and has many **practical applications**.

# Structural bioinformatics

## Use of structure

- Databases, classification
  - proteins, NA, drugs
- Patterns
  - Active sites, allosteric sites, ...
- Prediction
  - structure, function, active site, channels...
- Docking
  - Fitting of small molecules into the active site  
-> in silico drug design
- Simulations
  - What if...

# Problems of structural bioinformatics

- Structural data are hard to work with:
  - Nonlinear
  - Imprecise from experiment (resolution of structure)
  - 3D representation (3D search)
  - Visualization is not trivial
  - More conserved than sequence data (genomics)
  - Structural genomics prepare structures without annotation
  - Most structures are water soluble globular proteins (most drug targets are membrane proteins)

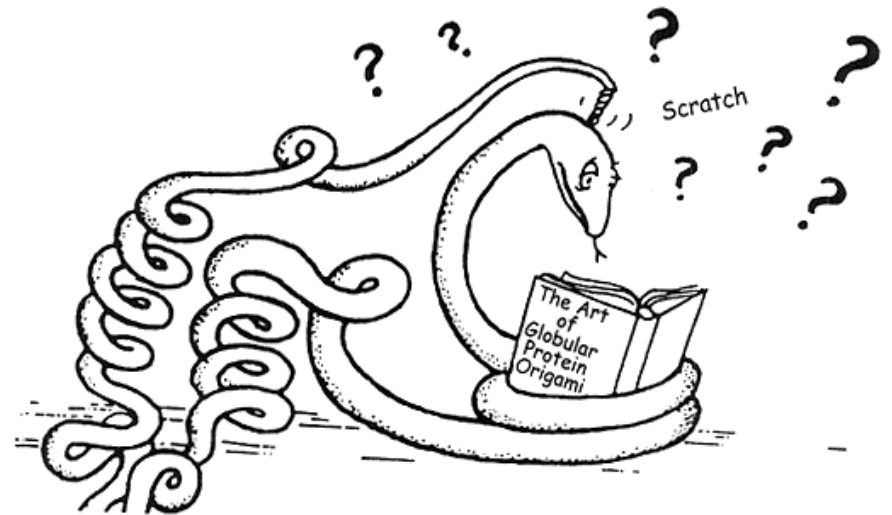
# Challenges

- Target selection
  - Large structures are resource intensive, maybe just one domain might be enough
- Structure methods
  - XRay – crystallisation is not easy
  - NMR – size problem – indistinguishable peaks
  - EM – only recently with atomistic detail
- Validation and Annotation
- Databases
- Correlation of structural data with experimental data

# Example 1 : Prediction of protein structure

- Tertiary structure
  - Fold recognition
  - Homolog modelling
    - Structural alignment
  - ab initio modelling
  - ML methods
- Function prediction
  - active sites, channels, pores, allosteric sites, conformations...

"Now collapse down hydrophobic core, and fold over helix 'A' to dotted line, bringing charged residues of 'A' into close proximity to ionic groups on outer surface of helix 'B' ..."

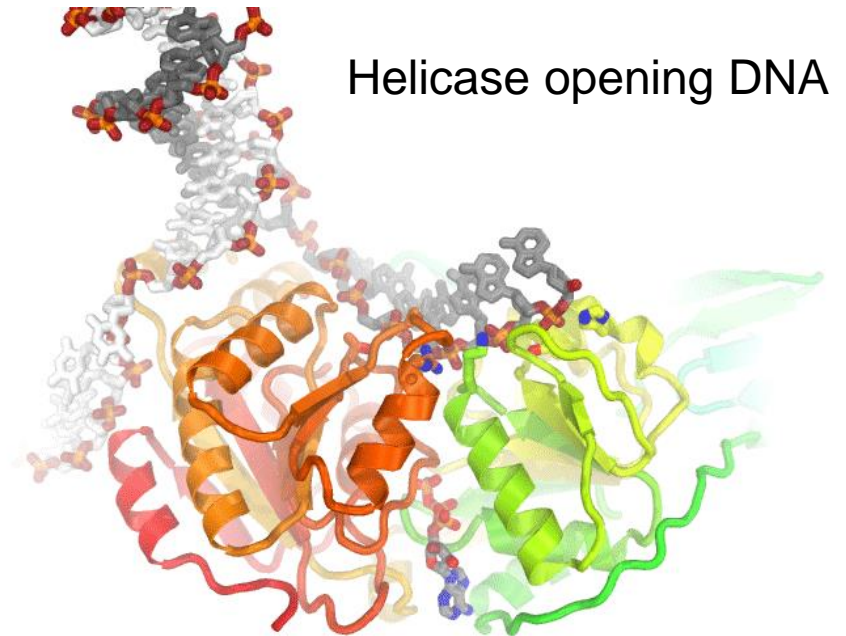


Reproduced in U. Tollemar, "Protein Engineering i USA", Sveriges Tekniska Attach er, 1988

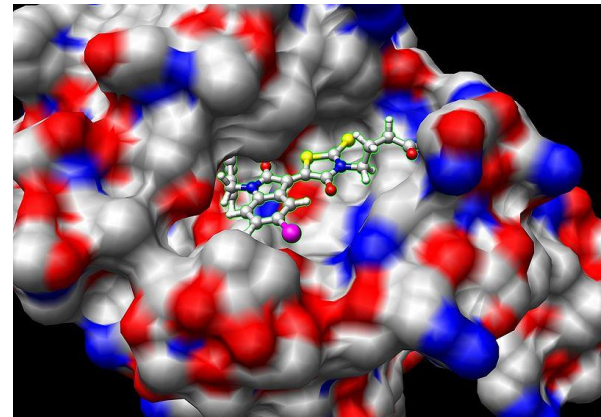


# Example 2: Molecular graphics

- We make nice figs!
- Simulations
  - Structure => Energy
  - Time => Dynamics
- Docking – binding
  - ligands
  - Protein-protein



GOLD docking of compound to acetyltransferase



# Structure Description

## Coordinate systems

- XYZ (cartesian)
- Inner coordinates (bond lengths, bond angles, torsion angles)
- object representation (secondary structure)

## Structure comparison:

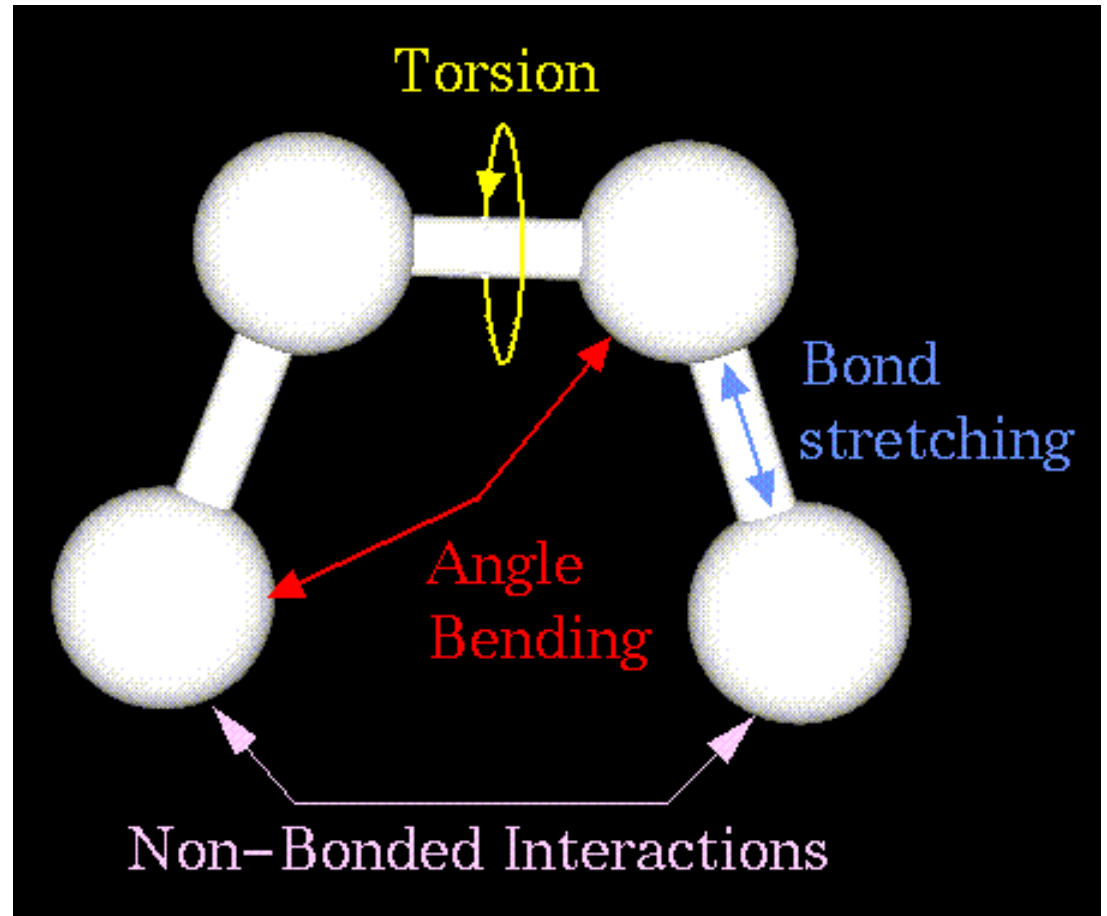
**RMSD** – root mean square distance

# Typical geometrical operations

Bond lengths

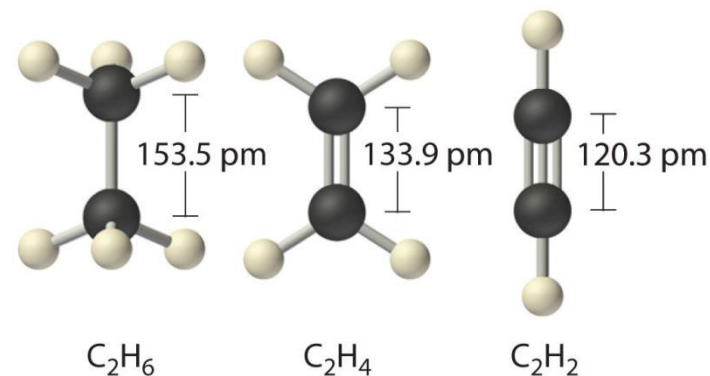
Bond angles

Torsions  
(dihedral angles)



# Bond Lengths

- function of position of 2 atoms
- Bond length is almost constant



- Type of bond
  - simple C-C
  - double C=C
  - triple C≡C

- Minimal - 1.09 Å (C–H)
- Typical - 1.54 Å (C–C)
- Longer – heteroatoms (sulphur, halogens, metal ions)

**Table 9.2** Average Bond Energies (kJ/mol) and Bond Lengths (pm)

Bond	Energy	Length	Bond	Energy	Length	Bond	Energy	Length	Bond	Energy	Length
<b>Single Bonds</b>											
H–H	432	74	N–H	391	101	Si–H	323	148	S–H	347	134
H–F	565	92	N–N	160	146	Si–Si	226	234	S–S	266	204
H–Cl	427	127	N–P	209	177	Si–O	368	161	S–F	327	158
H–Br	363	141	N–O	201	144	Si–S	226	210	S–Cl	271	201
H–I	295	161	N–F	272	139	Si–F	565	156	S–Br	218	225
			N–Cl	200	191	Si–Cl	381	204	S–I	~170	234
C–H	413	109	N–Br	243	214	Si–Br	310	216			
C–C	347	154	N–I	159	222	Si–I	234	240	F–F	159	143
C–Si	301	186							F–Cl	193	166
C–N	305	147	O–H	467	96	P–H	320	142	F–Br	212	178
C–O	358	143	O–P	351	160	P–Si	213	227	F–I	263	187
C–P	264	187	O–O	204	148	P–P	200	221	Cl–Cl	243	199
C–S	259	181	O–S	265	151	P–F	490	156	Cl–Br	215	214
C–F	453	133	O–F	190	142	P–Cl	331	204	Cl–I	208	243
C–Cl	339	177	O–Cl	203	164	P–Br	272	222	Br–Br	193	228
C–Br	276	194	O–Br	234	172	P–I	184	246	Br–I	175	248
C–I	216	213	O–I	234	194				I–I	151	266
<b>Multiple Bonds</b>											
C=C	614	134	N=N	418	122	C≡C	839	121	N≡N	945	110
C=N	615	127	N=O	607	120	C≡N	891	115	N≡O	631	106
C=O	745	123	O <sub>2</sub>	498	121	C≡O	1070	113			

(799 in CO<sub>2</sub>)

# Calculation of atom distance

In Cartesian coordinates:

For two points with coordinates  $(x_1, y_1, z_1)$  and  $(x_2, y_2, z_2)$

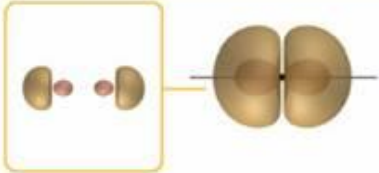
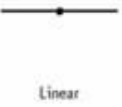
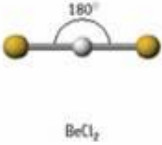
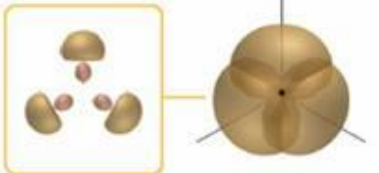

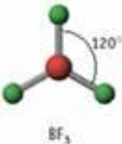
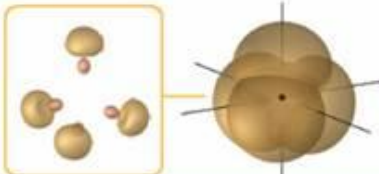


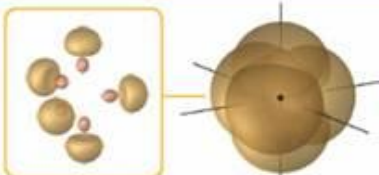

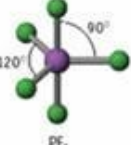
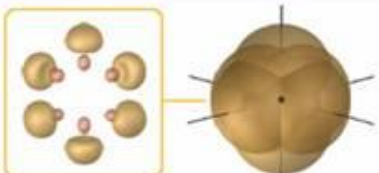

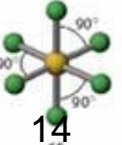
$$d_{2-1} = \sqrt{[(x_2 - x_1)^2 + (y_2 - y_1)^2 + (z_2 - z_1)^2]}$$

Some distances within protein backbone are **constant** even if not in direct bond:

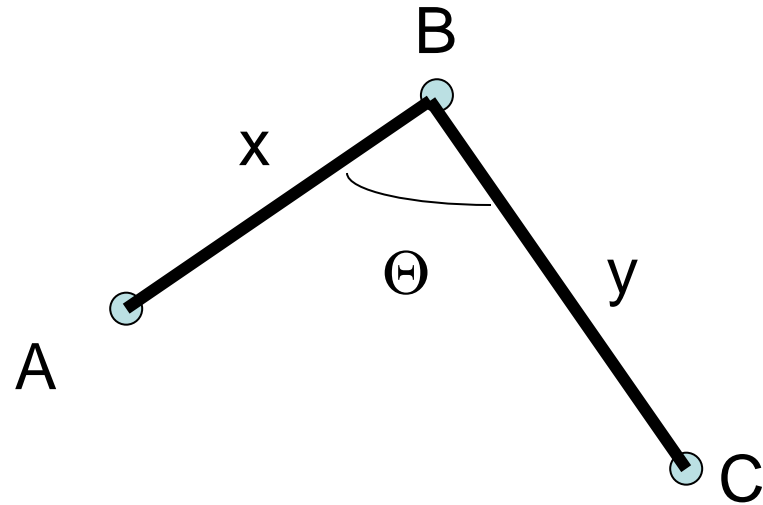
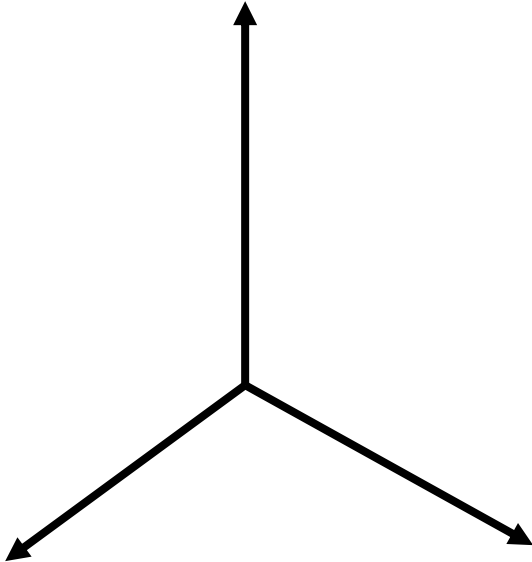
$C\alpha - C\alpha$  distance between consecutive amino acids is 3.8 Å

# Bond Angles

- function of position of 3 atoms
- Almost constant for given combination of type of atoms
- Depend on atom type and number of electrons in bonding
- Interval from 90 to 180

Arrangement of Hybrid Orbitals	Geometric figure	Example
Two electron pairs $sp$		 Linear  $180^\circ$ BeCl <sub>2</sub>
Three electron pairs $sp^2$		 Trigonal-planar  $120^\circ$ BF <sub>3</sub>
Four electron pairs $sp^3$		 Tetrahedral  $109.5^\circ$ CH <sub>4</sub>
Five electron pairs $sp^3d$		 Trigonal-bipyramidal  $90^\circ$ $120^\circ$ PF <sub>5</sub>
Six electron pairs $sp^3d^2$		 Octahedral  $90^\circ$ $90^\circ$ SF <sub>6</sub>

# Calculation of bonding angle



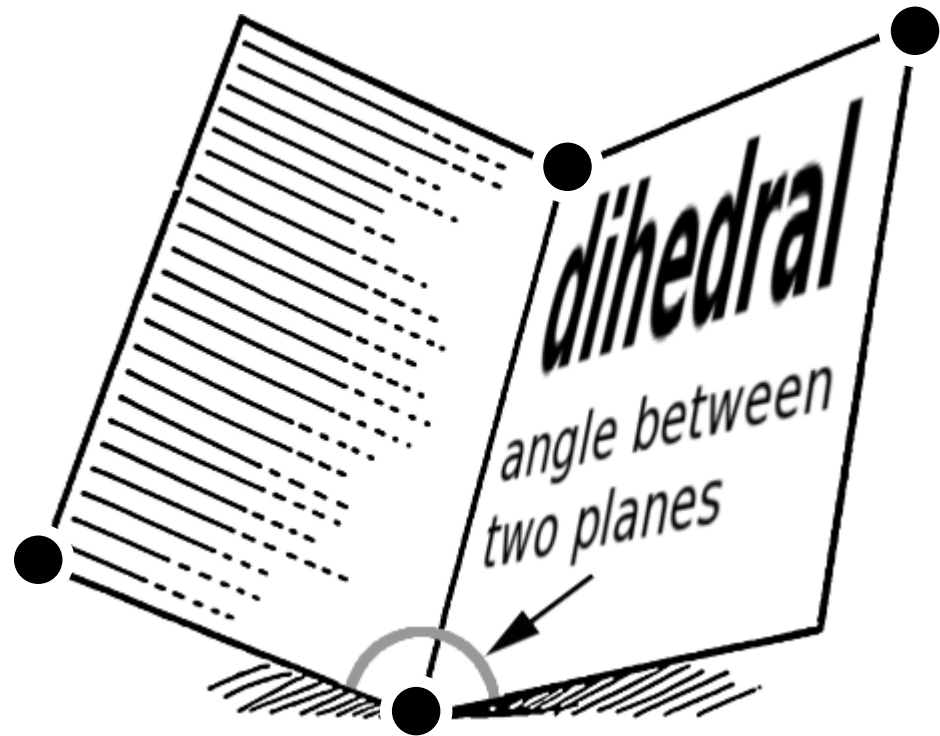
$$X \cdot Y = |X| \cdot |Y| \cdot \cos(\Theta)$$

$$\Theta = \arccos(X \cdot Y / |X| \cdot |Y|)$$

Arccosin of angle between two vectors BA and BC

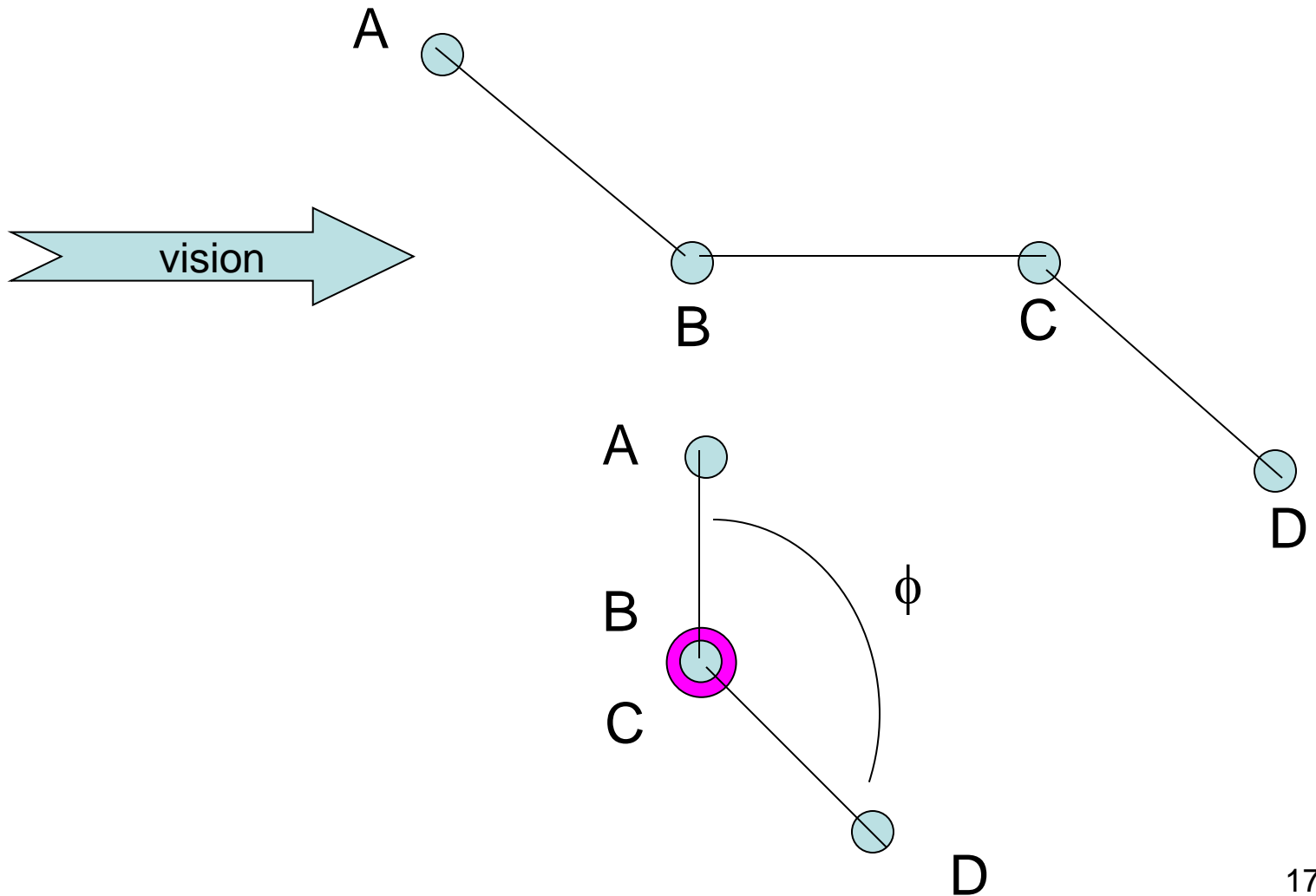
# Dihedral Angle

- function of position of 4 atoms
- Quite variable (0 to  $360^\circ$ )
- its change change conformations





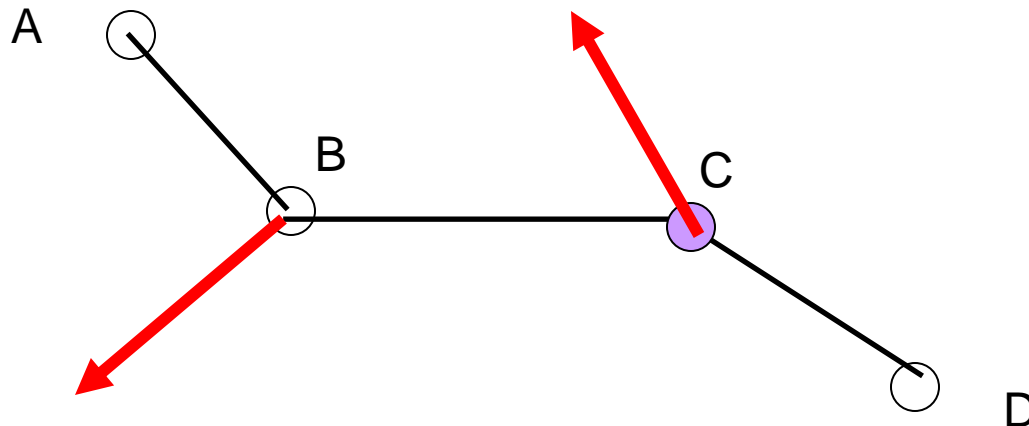
# Dihedral Angle



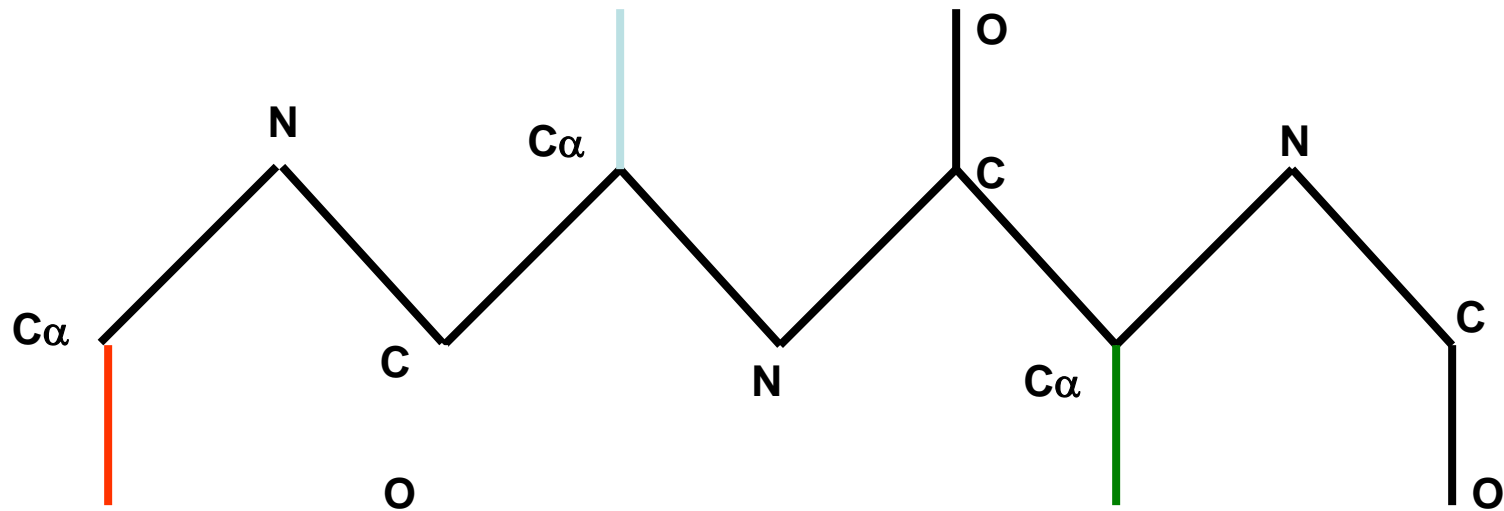
# Calculation of dihedral angle

Dihedral angle = Angle between vectors orthogonal to planes defined by vectors:

- 1) Plane 1 - Vectors BA and CB
- 2) Plane 2 - Vectors CB and DC



# Important dihedral angles in proteins



omega  $\omega$

$C\alpha - C\alpha$



psi  $\psi$

N - N



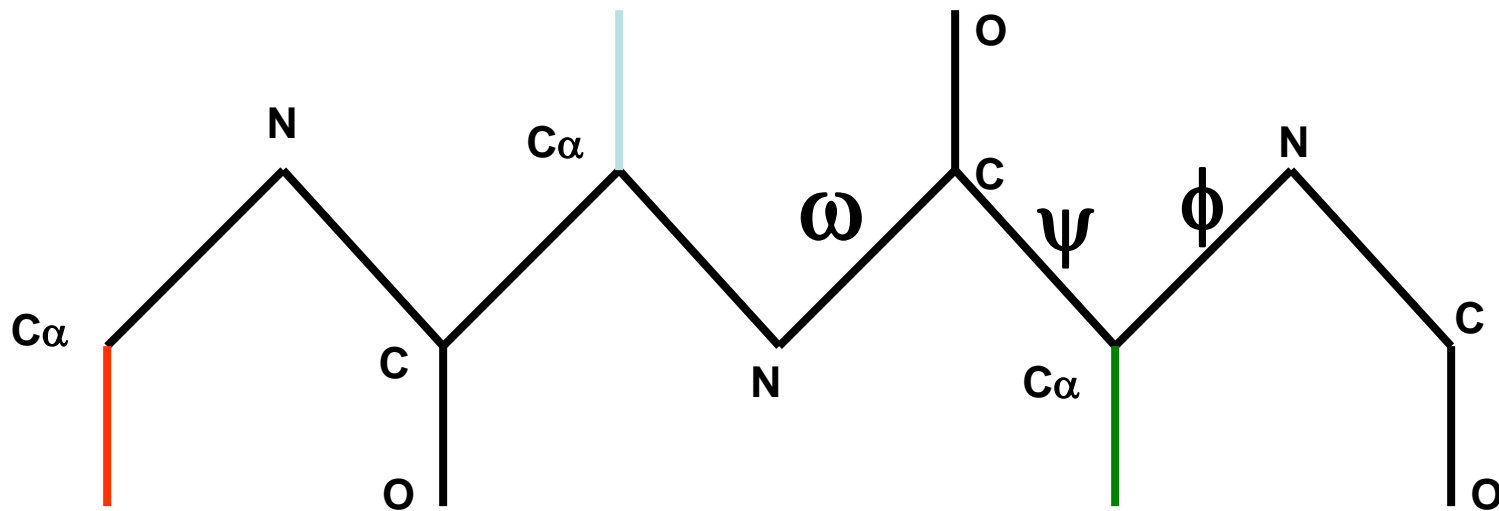
phi  $\phi$

C - C



# Important dihedral angles in proteins

- Omega  $\omega$  is constant = 180 (C-N do not rotate)
- Phi  $\Phi$ , Psi  $\Psi$  intervals (C $\alpha$ -N, C-C $\alpha$  can rotate) restricted to certain areas due to following amino acids



# Ramachandran plot

- Typical values of dihedral angles define individual secondary structure elements:
  - $\alpha$ -helix                       $\phi = -57, \psi = -47$
  - 3-10 helix                         $\phi = -49, \psi = -26$
  - Parallel  $\beta$ -sheet                 $\phi = -119, \psi = 113$
  - Antiparallel  $\beta$ -sheet            $\phi = -139, \psi = 135$

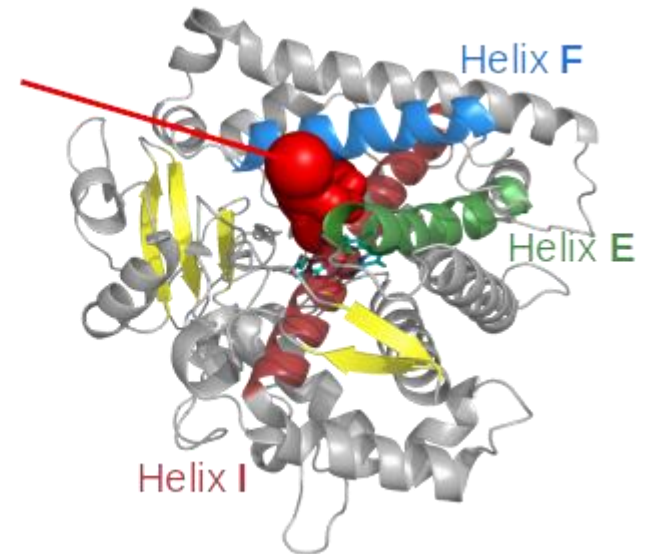
# Secondary structure

Helices and  $\beta$ -strands = Secondary Structure Elements (SSEs)

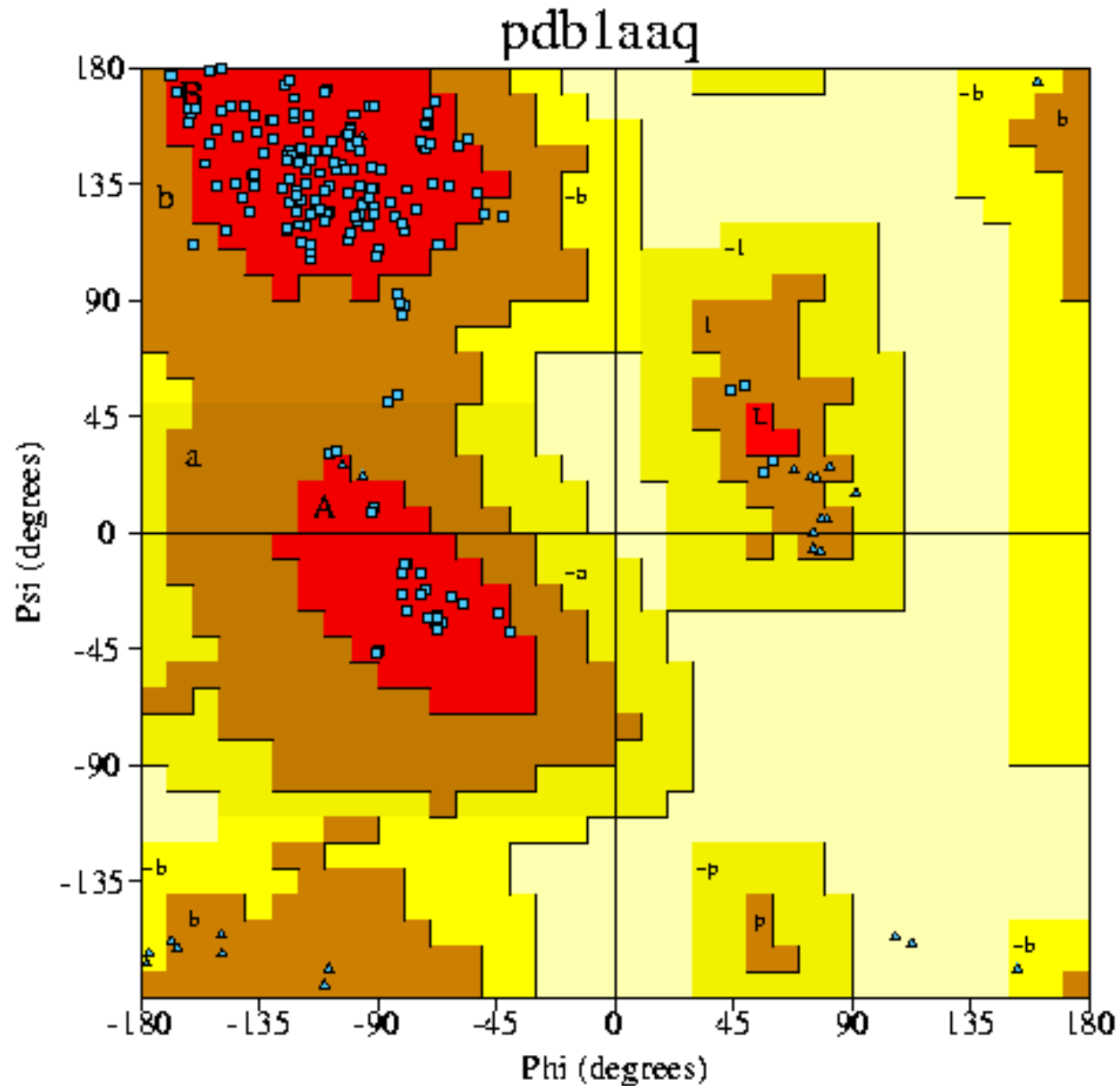


- Quite conserved arrangement within a protein family
- Can serve as landmarks, which
  - Help us orient in the structure
  - Help us locate the key regions (active sites, channels...)

Solvent  
channel



# Ramachandran plot



# Other Coordinate Systems

**Cartesian coordinates** are orthogonal (x,y,z)

-> used most often

If bond lengths and bond angles are constant ->  
reduction of coordinates -> only dihedral angles =>

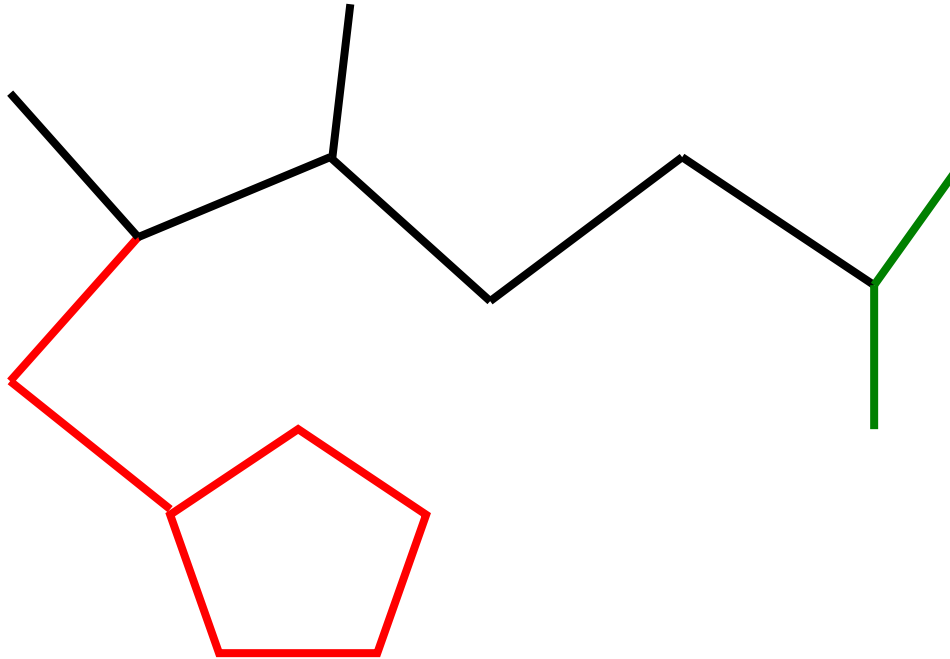
**Inner coordinates**

If some part of structure can be defined by “rigid”  
structural element -> solid objects =>

**Object-based coordinates**



# Advantages of Inner Coordinates



3 peptide units = 12 atoms = 36 coordinates OR 6 dihedral angles  
3 sidechains = 12 atoms = 36 souřadnic OR 5 dihedral angles

*72 cartesians versus 11 inners*

# Disadvantages of Inner Coordinates

**Some calculations are more difficult**

*Atom-atom distance*

*Closest atoms toward a point in space*

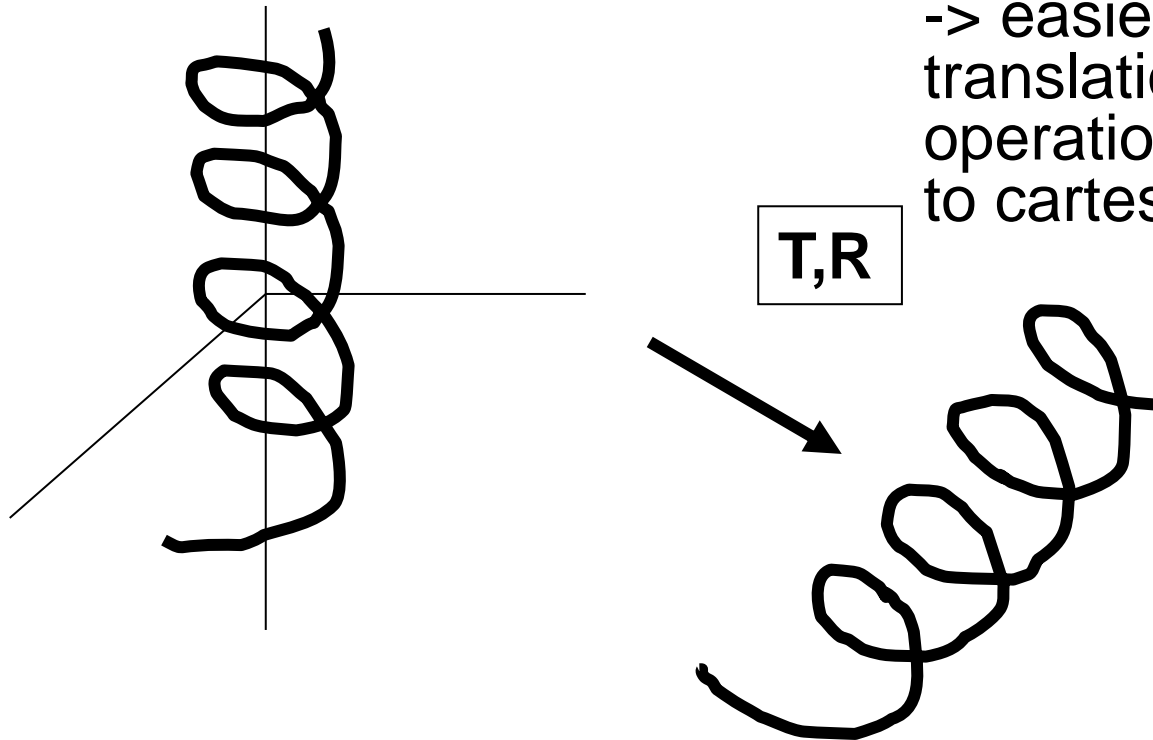
**Hard comparison of independent objects (two molecules)**

**Nonlinear relationships between coordinates => problem for optimizations and simulations**

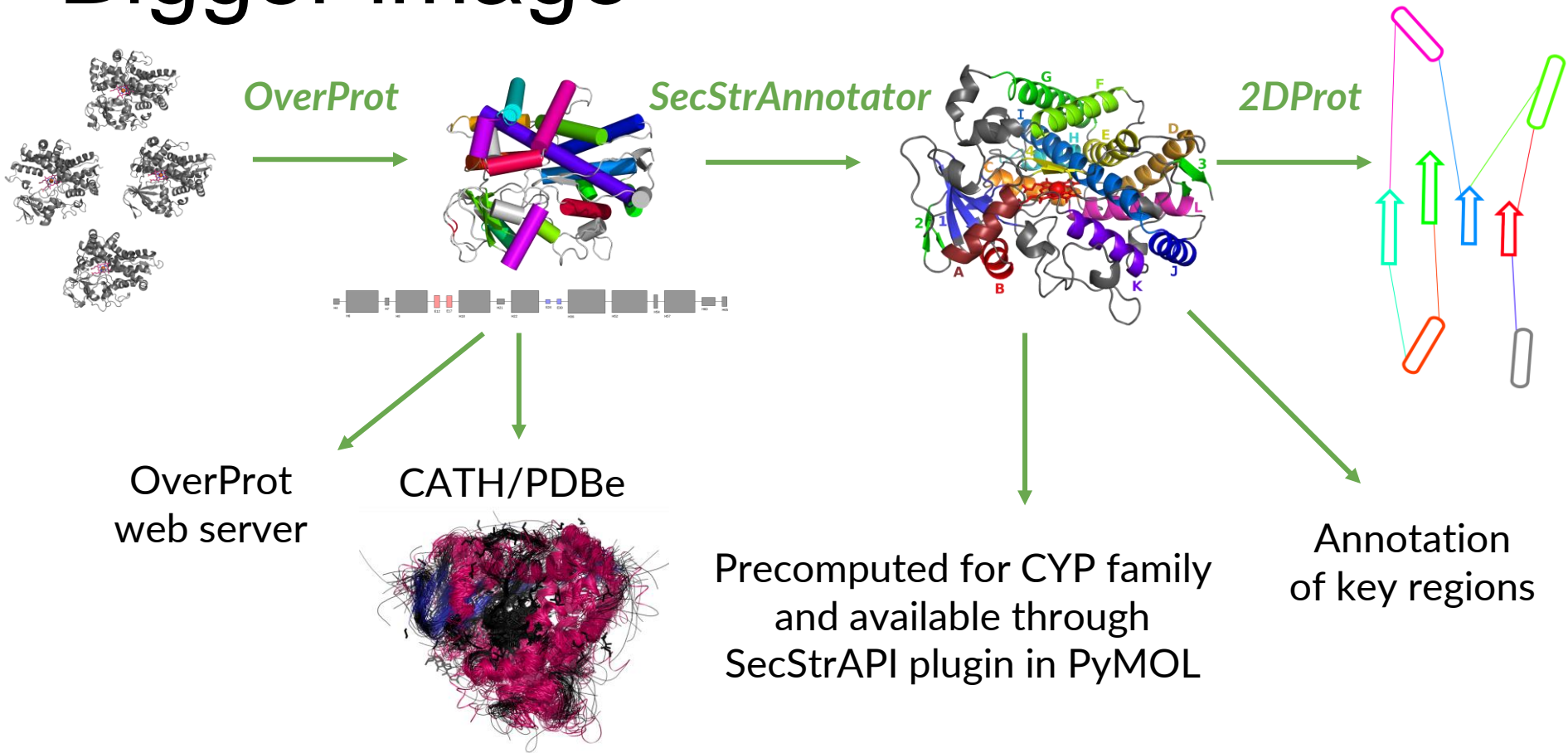
# Object-based coordinates

Use of larger objects – secondary structure, subset of atoms...

Example -> Helix can be represented as a vector with just 6 coordinates  
-> easier operations such as translation and rotation (final operation is later transferred to cartesian coordinates)



# Bigger image



Midlik A, Hutařová Vařeková I, Hutař J, Charehneou A, Berka K, Svobodová R: **OverProt**: secondary structure consensus for protein families, *Bioinformatics*, 38(14), July 2022, 3648–3650

Midlik A, Navrátilová V, Moturu TR, Koča J, Svobodová R, Berka K: Uncovering of cytochrome P450 anatomy by **SecStrAnnotator**. *Sci Rep* 11, 2021, 12345

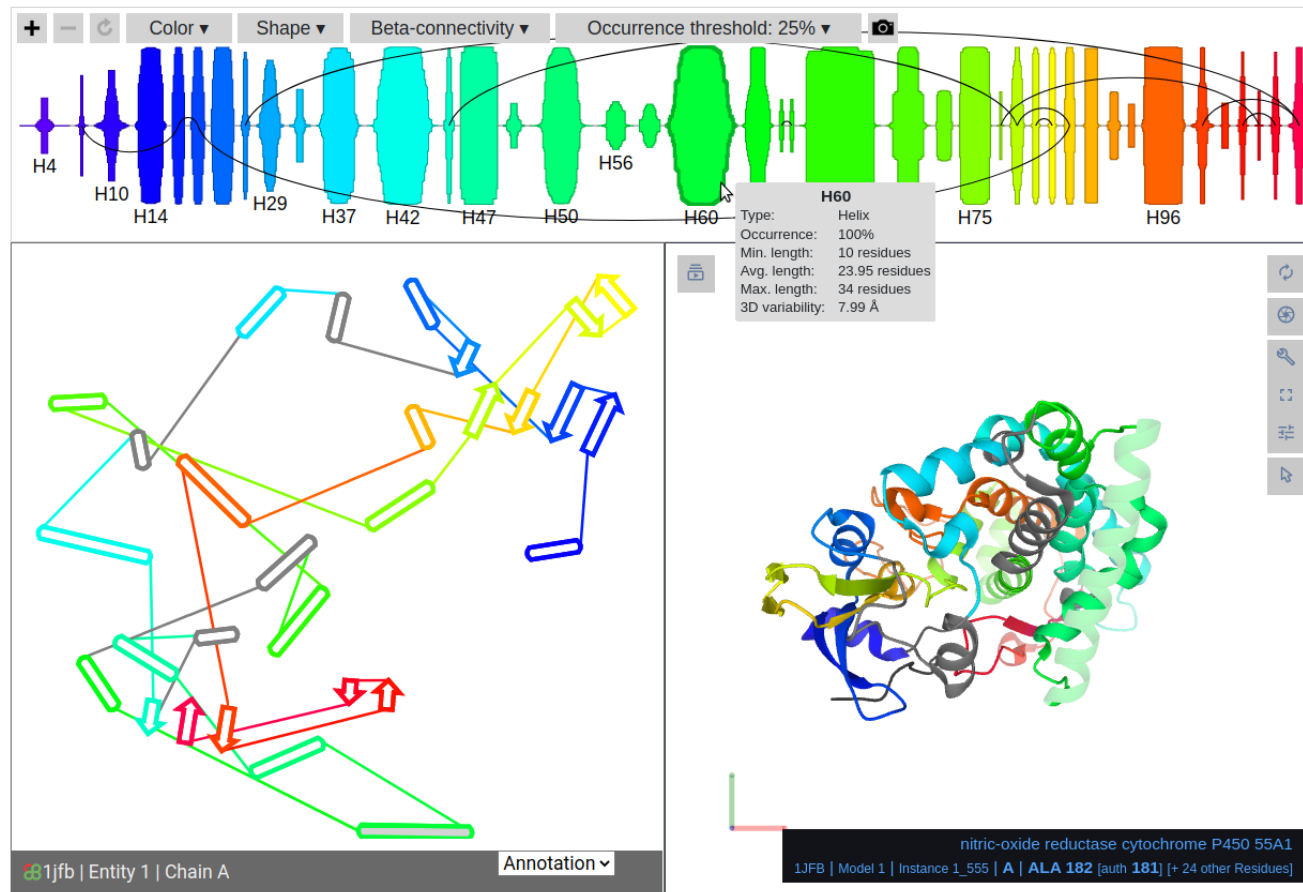
Hutařová Vařeková I, Hutař J, Midlik A, Horský V, Hladká E, Svobodová R, Berka K, **2DProts**: database of family-wide protein secondary structure diagrams, *Bioinformatics*, 37(23), 2021, 4599–4601,

# OverProt Server – Interactive view

- 1D of the family linked to 2D and 3D of a domain

Family: 1.10.630.10 *Cytochrome P450*

Domain: 1jfbA00



# Structure Comparison

For comparison of two structures A and B we need:

- 1. Which atom from A corresponds to which atom from B**  
=> alignment
- 2. Atom localization**  
=> PDB files
- 3. Comparison criteria**  
RMSD, energy

# RMSD = Root Mean Square Deviation

- Atoms from A and B are taken as equivalent
- Superposition and calculation of differences in distance

$$\text{RMSD} = \sqrt{\frac{\sum d_i^2}{N}}$$

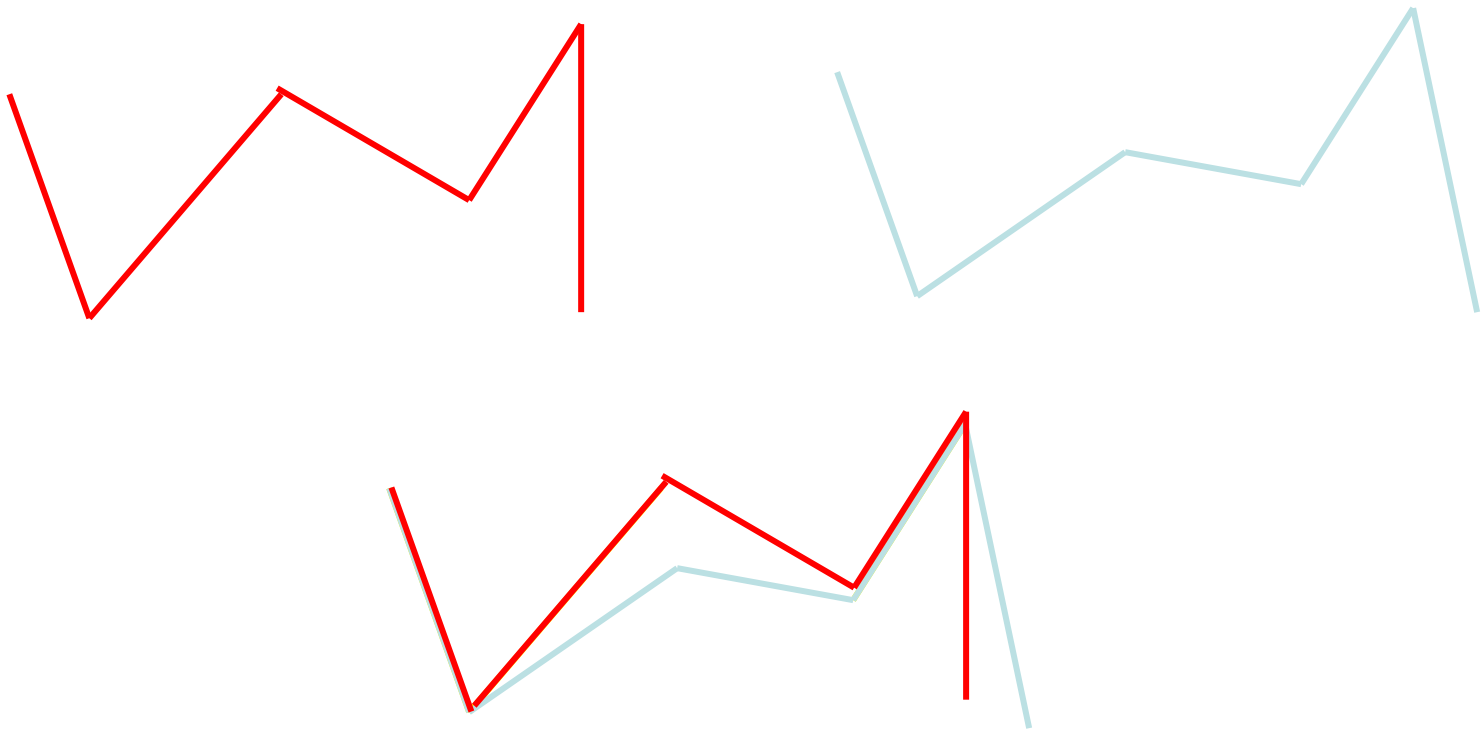
- If are structures identical -> RMSD = 0
- With more differences between structures -> RMSD increases

N – number of atoms

$d_i$  – distance of two atoms with index  $i$  from A and B

# Structure Comparison

To find minimal RMSD





# Calculation of RMSD

- translate and rotate one structure with respect to the other to minimize the RMSD
- Centroid-based solutions  
(Huang, Blostein, Margerum)
- Quaternion-based solutions
  - (rotation-translation) that minimizes the RMSD between two sets of vectors  
(Faugeras a Hebert, Petitjean)
- Matrix Singularity-based methods  
(Arun, Huang, Blostein)

# Arun algorithm

- Matrices of  $p_i' = R.p_i + T + N_i$ 
  - $p_i$  – 3x1 column matrix of positions
  - $R$  – rotation matrix
  - $T$  – translation vector 3x1 column matrix
  - $N$  – noise vector
- 1) Translation over **centroids**
- 2) Singular value decomposition of matrix to obtain **rotation**
  
- Arun algorithm is optimal, universal and not iterative

# Kabsch algorithm

- 1) Translation over **centroids**
  - 2) computation of a **covariance matrix**,
  - 3) the computation of the **optimal rotation matrix**.
- 
- Kabsch algorithm is widely used as *fit* function in PyMol, or within VMD
  - Algorithm do not recognise similar pairs of residues – these have to be defined iteratively (typically  $C\alpha$ )

Kabsch, W (1976): A solution for the best rotation to relate two sets of vectors. *Acta Cryst.* **A32** (5): 922.  
With a correction in Kabsch, W (1978). ["A discussion of the solution for the best rotation to relate two sets of vectors"](#). *Acta Cryst.* **A34** (5): 827–828.

# Advantages and Disadvantages of RMSD

Good behavior, identical structures RMSD = 0

Simple calculation in Cartesian coordinates

Natural units (Ångstroms)

Experience (similar structures have RMSD ~ (1 – 3 Å))

**Weight of all atoms is the same**

however hydrogens have much smaller effect in practice  
→ RMSD only for backbone or C $\alpha$

**Prone to extremities**

**RMSD of larger protein is larger even if the structure is almost identical**

RMSD of 3 Å for 100 residue protein is really bad,  
for 1000 residue protein it is sensible.

# Other measures

- **global distance test (GDT)**

- largest set of amino acid residues' C $\alpha$  atoms in the model structure falling within a defined distance cutoff of their position in the experimental structure.

⇒ Used in structure prediction assessment (CASP)

- **template modeling score (TM-score)**

- difference between two structures by a score between (0,1]

$$\text{TM-score} = \max \left[ \frac{1}{L_{\text{target}}} \sum_i^{L_{\text{aligned}}} \frac{1}{1 + \left( \frac{d_i}{d_0(L_{\text{target}})} \right)^2} \right]$$

- TM-score = 1 - perfect match between two structures

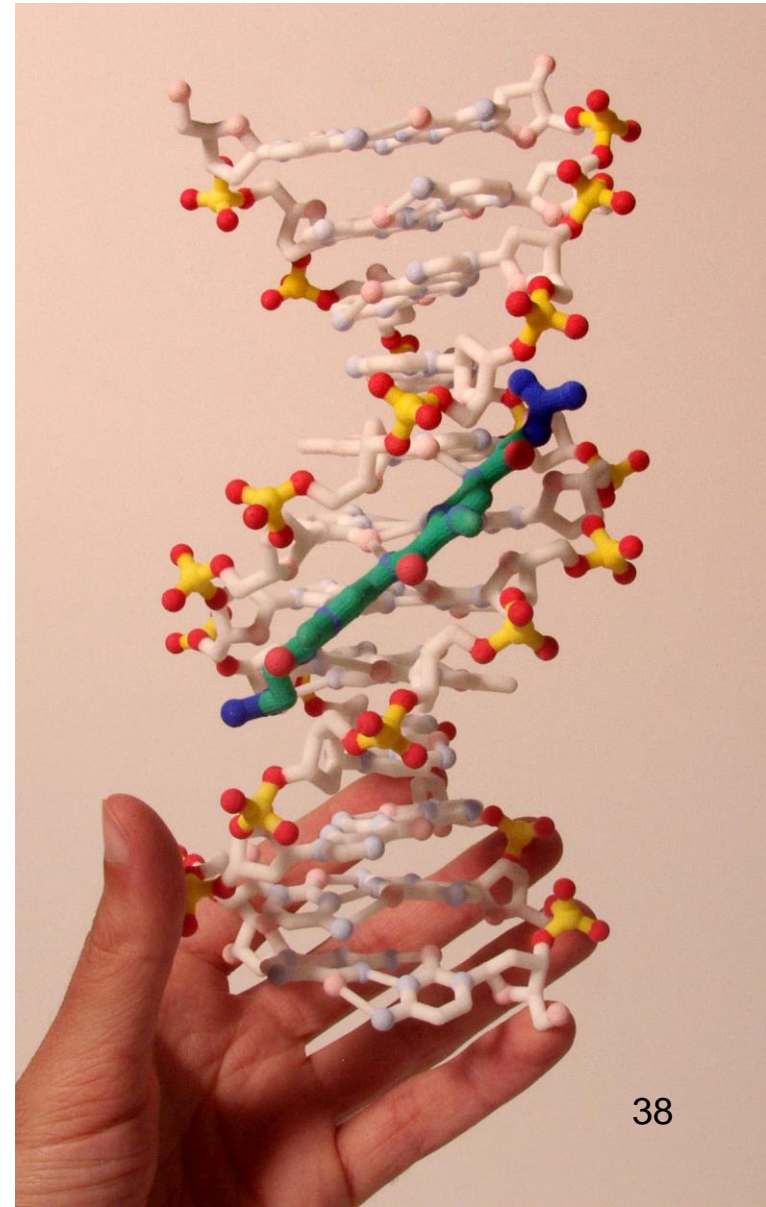
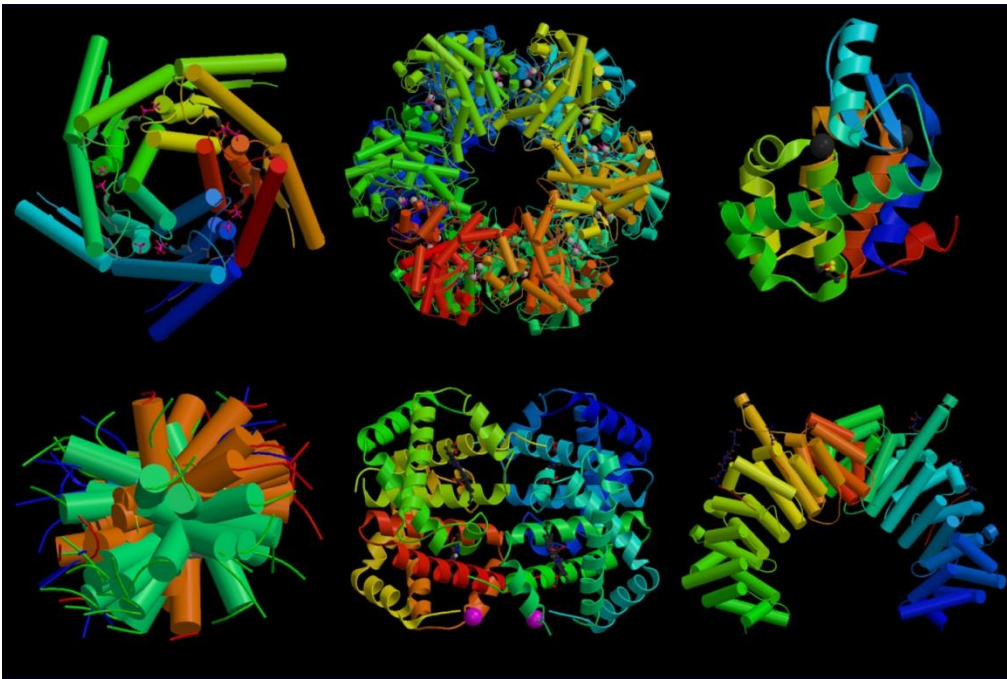
- TM-score > 0.5 assume roughly the same fold

- TM-score < 0.20 - randomly chosen unrelated proteins

⇒ Used in structure prediction assessment (CASP)

# Biomolecules

- proteins
- NA – DNA, RNA
- lipids
- polysaccharides
- Small molecules (hormones, drugs)



# Structural Hierarchy

## **MOLECULAR STRUCTURE**

Primary (sequence)



Secondary (local folding)



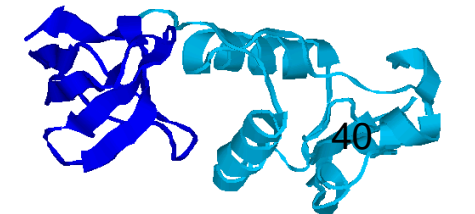
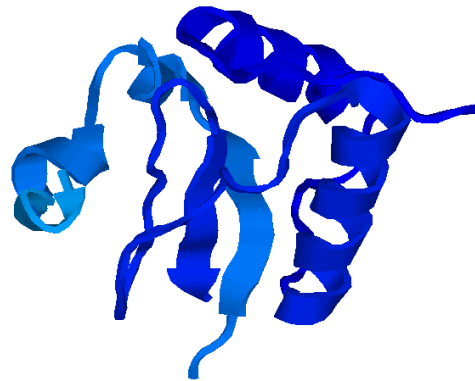
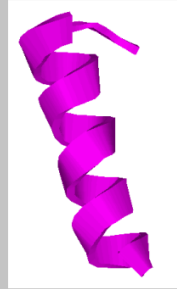
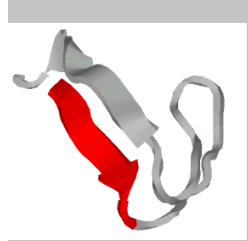
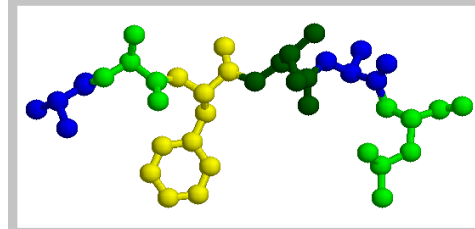
Tertiary (long-range folding)



Quaternary (multimeric organization)

# Proteins

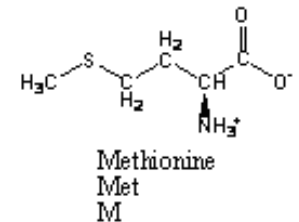
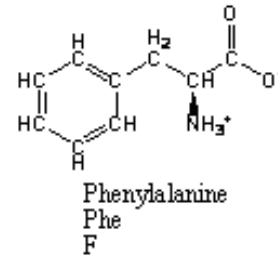
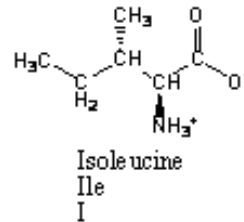
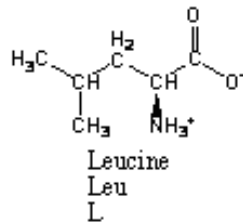
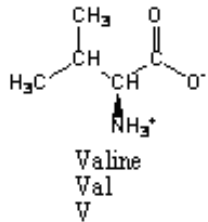
- Amino acids
- Backbone and Sidechains
- Primary structure
  - sequence of amino acids
- Secondary structure
  - Local structural patterns
- Tertiary structure
  - Domain Fold
- Quarternary structure
  - Multichain organization



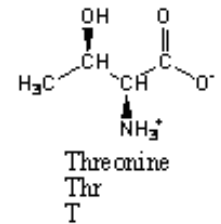
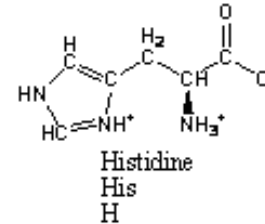
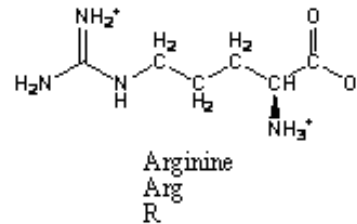
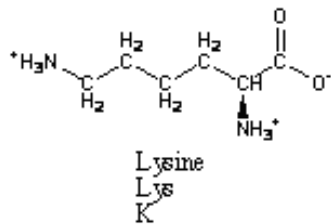
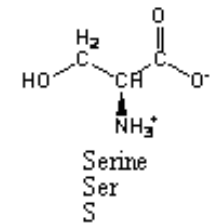
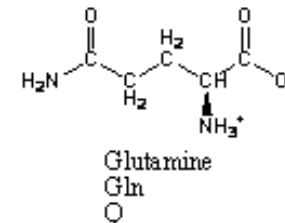
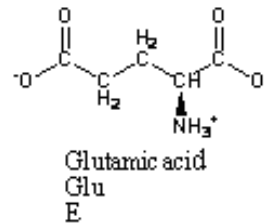
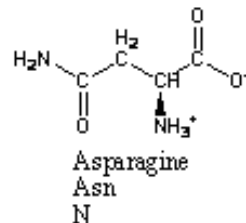
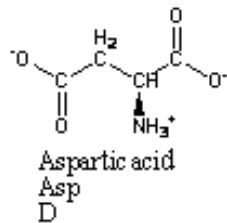


# Amino acids

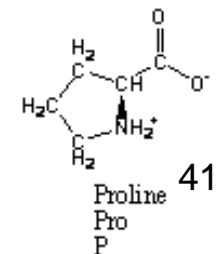
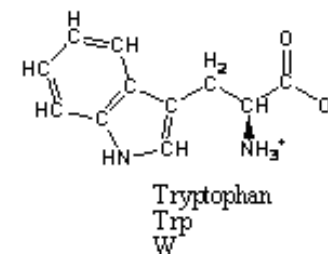
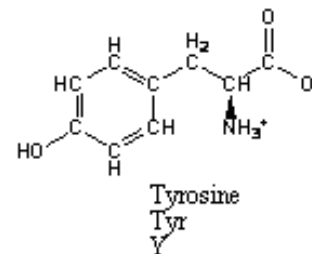
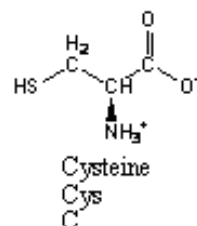
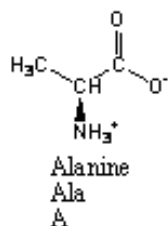
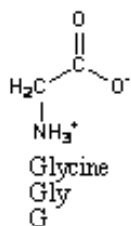
## Amino acids with hydrophobic side chains



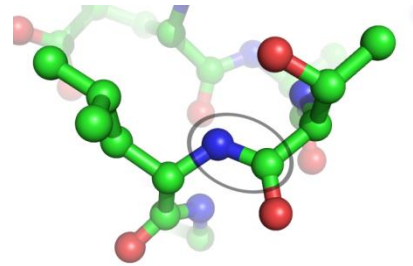
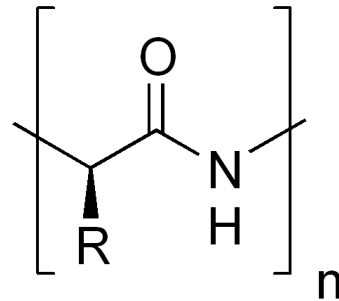
## Amino acids with hydrophilic side chains



## Amino acids with intermediate side chains



# Primary Structure of Protein



## AMINO ACID

## SIDE CHAIN

Aspartic acid	Asp	D	negative
Glutamic acid	Glu	E	negative
Arginine	Arg	R	positive
Lysine	Lys	K	positive
Histidine	His	H	positive
Asparagine	Asn	N	uncharged polar
Glutamine	Gln	Q	uncharged polar
Serine	Ser	S	uncharged polar
Threonine	Thr	T	uncharged polar
Tyrosine	Tyr	Y	uncharged polar

## POLAR AMINO ACIDS

## AMINO ACID

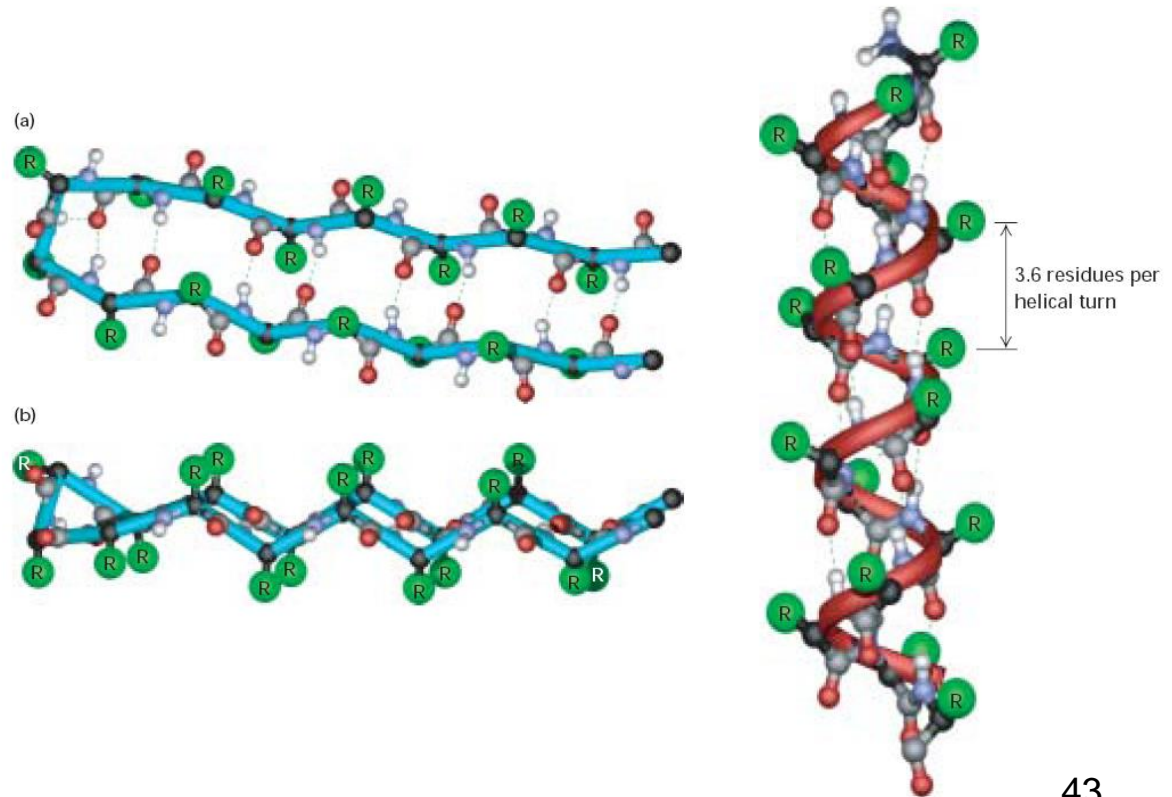
## SIDE CHAIN

Alanine	Ala	A	nonpolar
Glycine	Gly	G	nonpolar
Valine	Val	V	nonpolar
Leucine	Leu	L	nonpolar
Isoleucine	Ile	I	nonpolar
Proline	Pro	P	nonpolar
Phenylalanine	Phe	F	nonpolar
Methionine	Met	M	nonpolar
Tryptophan	Trp	W	nonpolar
Cysteine	Cys	C	nonpolar

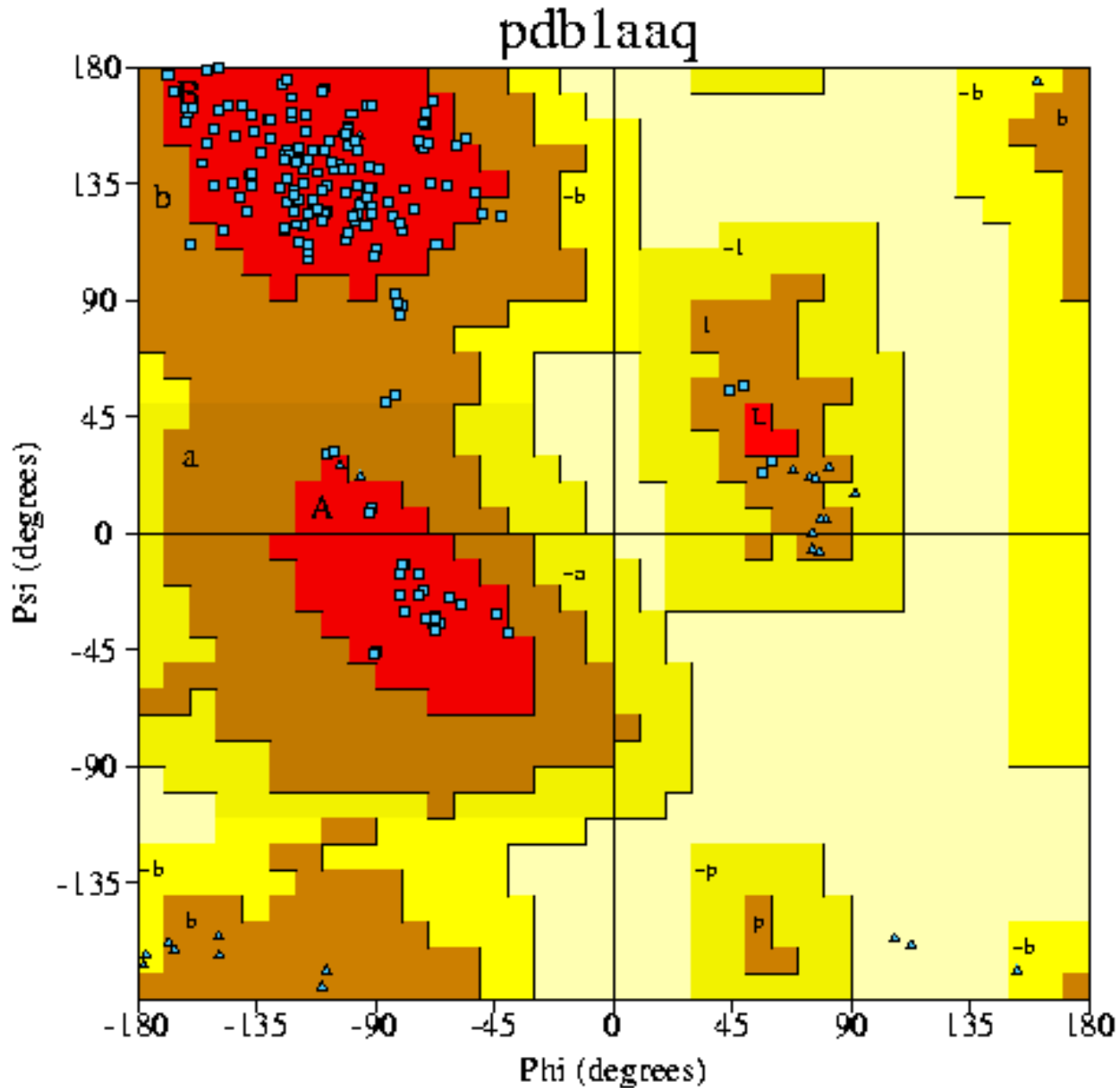
## NONPOLAR AMINO ACIDS

# Secondary structure of Proteins

- Local folding
- Secondary structure depends on amino acid sequence
  - $\alpha$ -helix
  - 3-10 helix
  - $\beta$ -sheet
  - $\beta$ -turn, loop



# Ramachandran plot



# PROCHECK summary for 1aaq

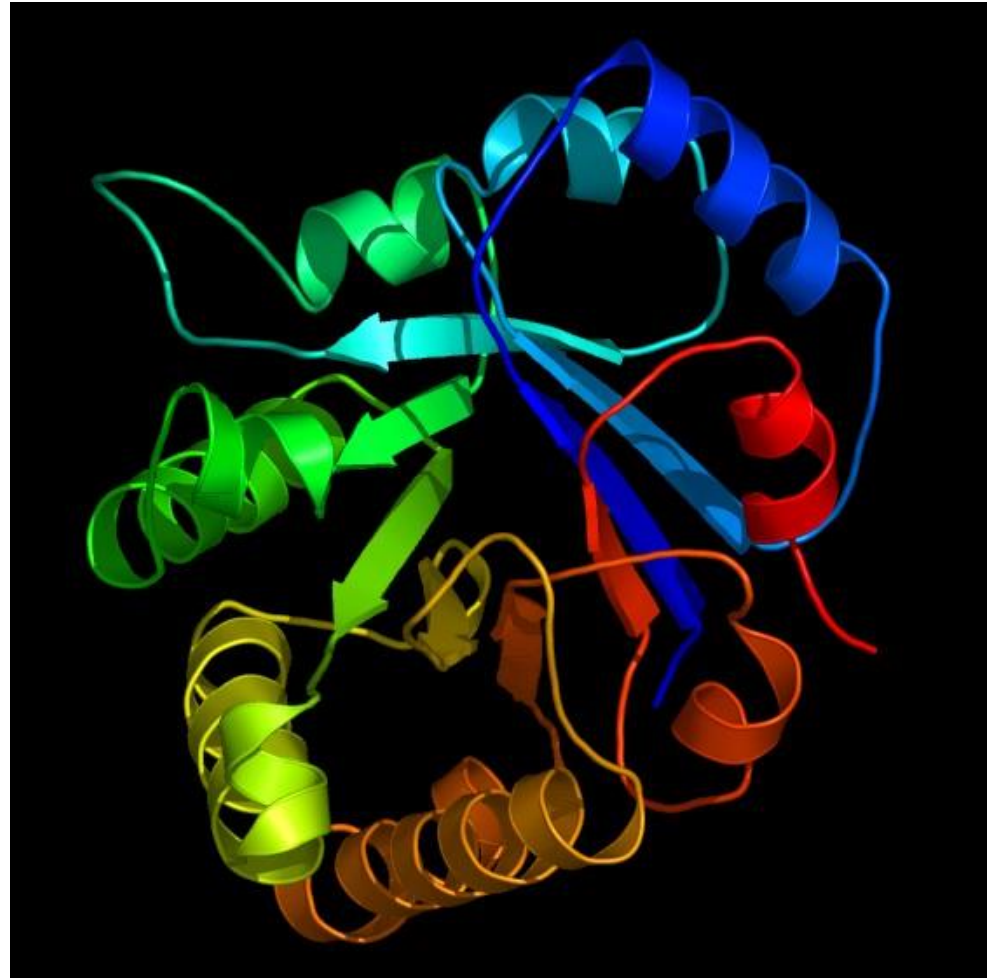
## PROCHECK statistics

### Ramachandran Plot statistics

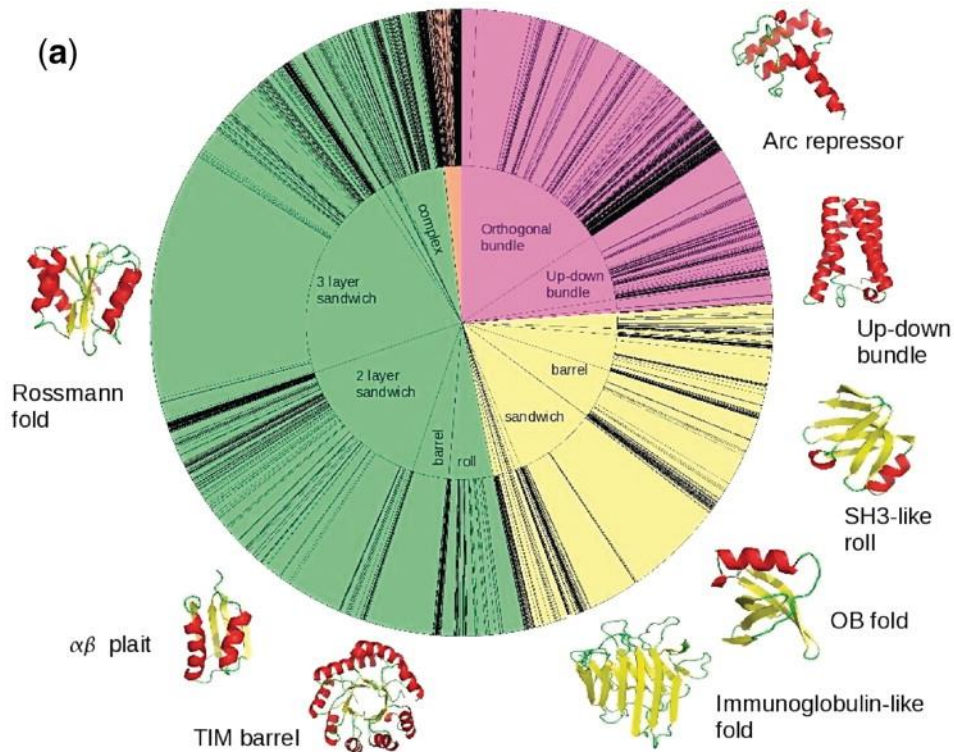
	No. of residues	%-tage
	-----	-----
Most favoured regions [A,B,L]	146	92.4%
Additional allowed regions [a,b,l,p]	12	7.6%
Generously allowed regions [~a,~b,~l,~p]	0	0.0%
Disallowed regions [XX]	0	0.0%
	-----	-----
Non-glycine and non-proline residues	158	100.0%
End-residues (excl. Gly and Pro)	2	
Glycine residues	26	
Proline residues	12	
	----	
Total number of residues	198	

# Tertiary Structure

- fold
  - globular
  - membrane
  - Fibrillar
  - IUP
- Necessary for **FUNCTION**
- domains



# 'CATHerine wheels'.



The distribution of all non-homologous structures (2386) within CATH v3.3

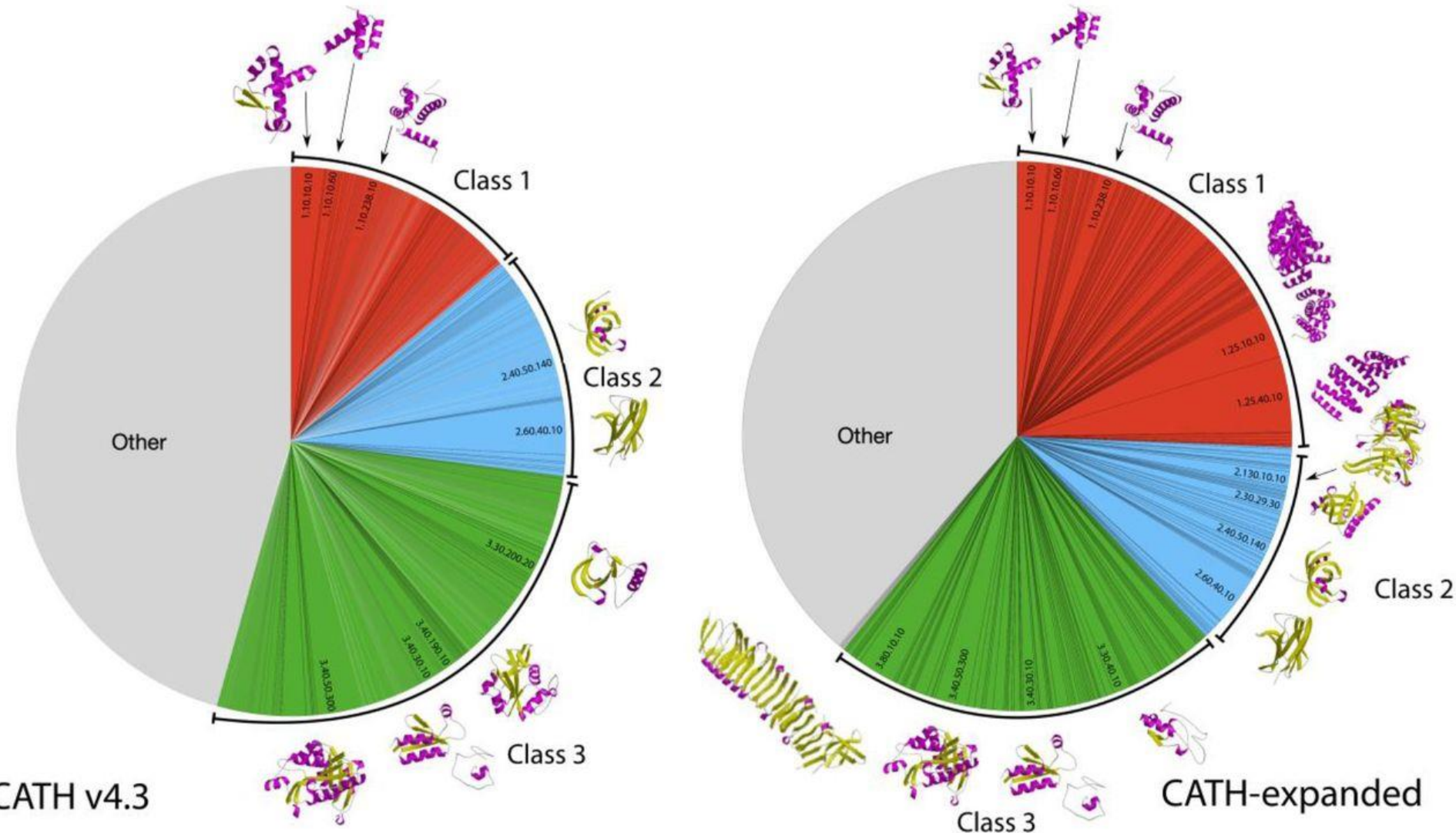
## Classes:

pink (mainly  $\alpha$ ),  
yellow (mainly  $\beta$ ),  
green ( $\alpha\beta$ )  
brown (little secondary structure).

**Proportion** of structures within any given architecture (inner circle)

**Fold group** (outer circle).

# CATH update 2022



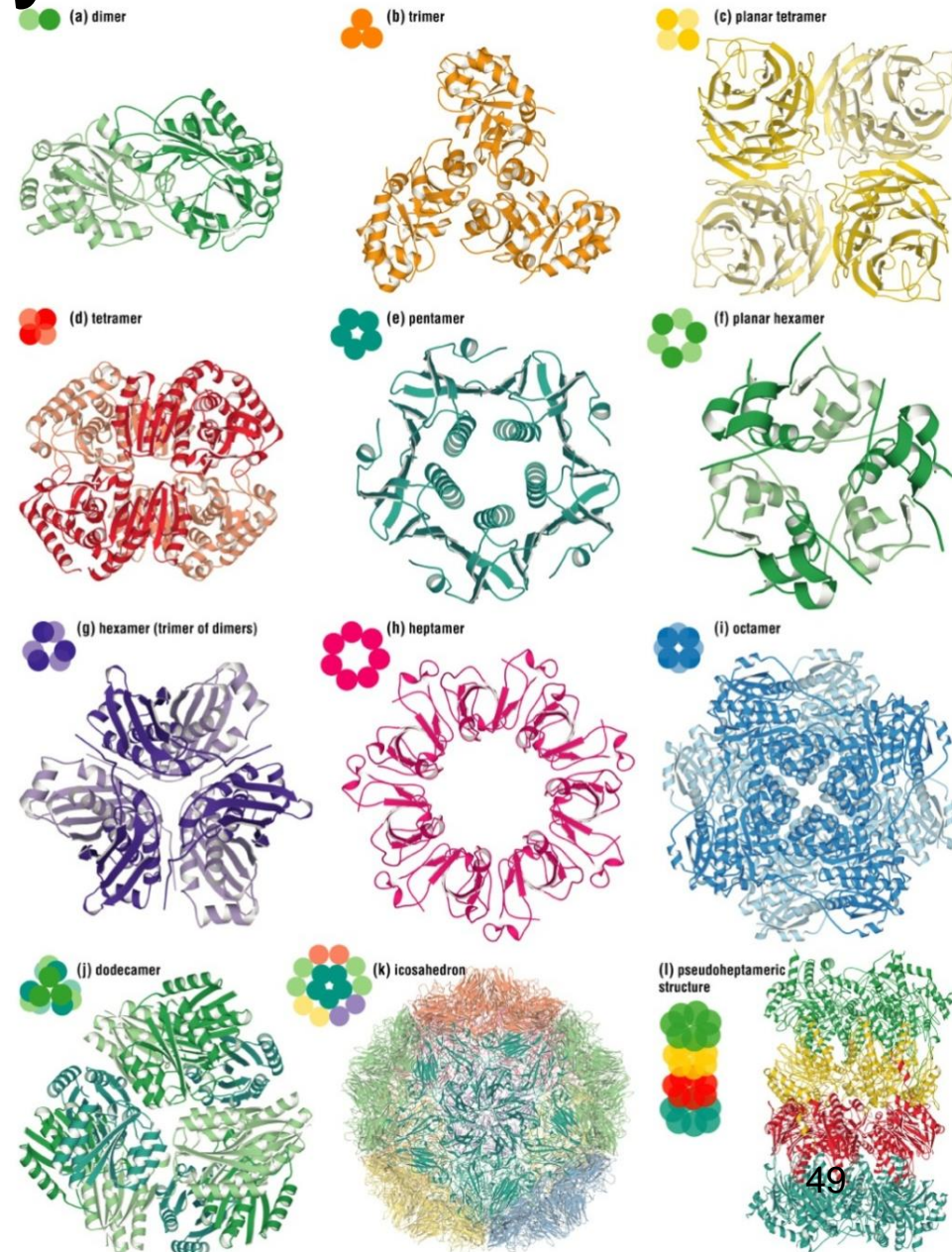
Structural diversity in CATH Superfamilies (left) and expanded by AF2 models (right).

Bordin N et al.: AlphaFold2 reveals commonalities and novelties in protein structure space for 21 model organisms. *bioRxiv* 2022, doi:10.1101/2022.06.02.494367v1.full

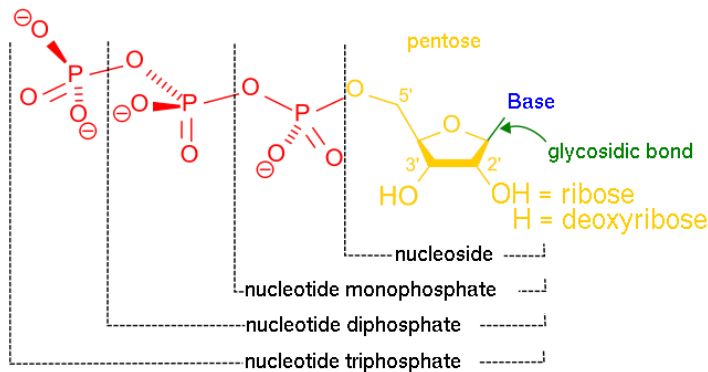


# Quarternary Structure

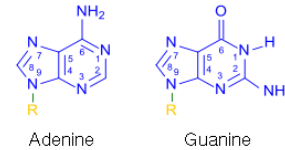
- asociace více řetězců:
  - Kooperativita  
(asociace zesílí vazebné vlastnosti)  
hemoglobin
  - Kolokalizace funkce  
(každá podjednotka dělá něco jiného)  
tryptophansyntáza
  - Kombinace podjednotek  
(přizpůsobování)  
imunoglobuliny
  - Skládání větších struktur  
(podjednotky uspořádávají procesem self-assembly)  
aktin,  
virové kapsidy



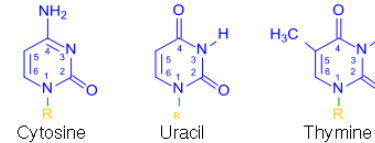
# Nucleic Acids (NA)



## Purines



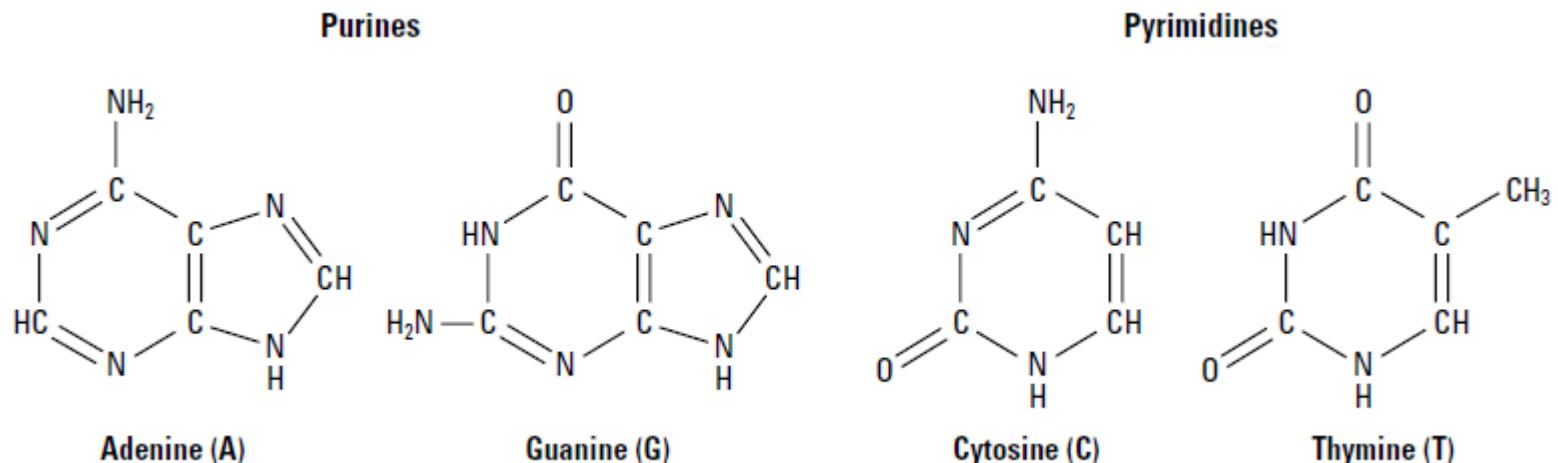
## Pyrimidines



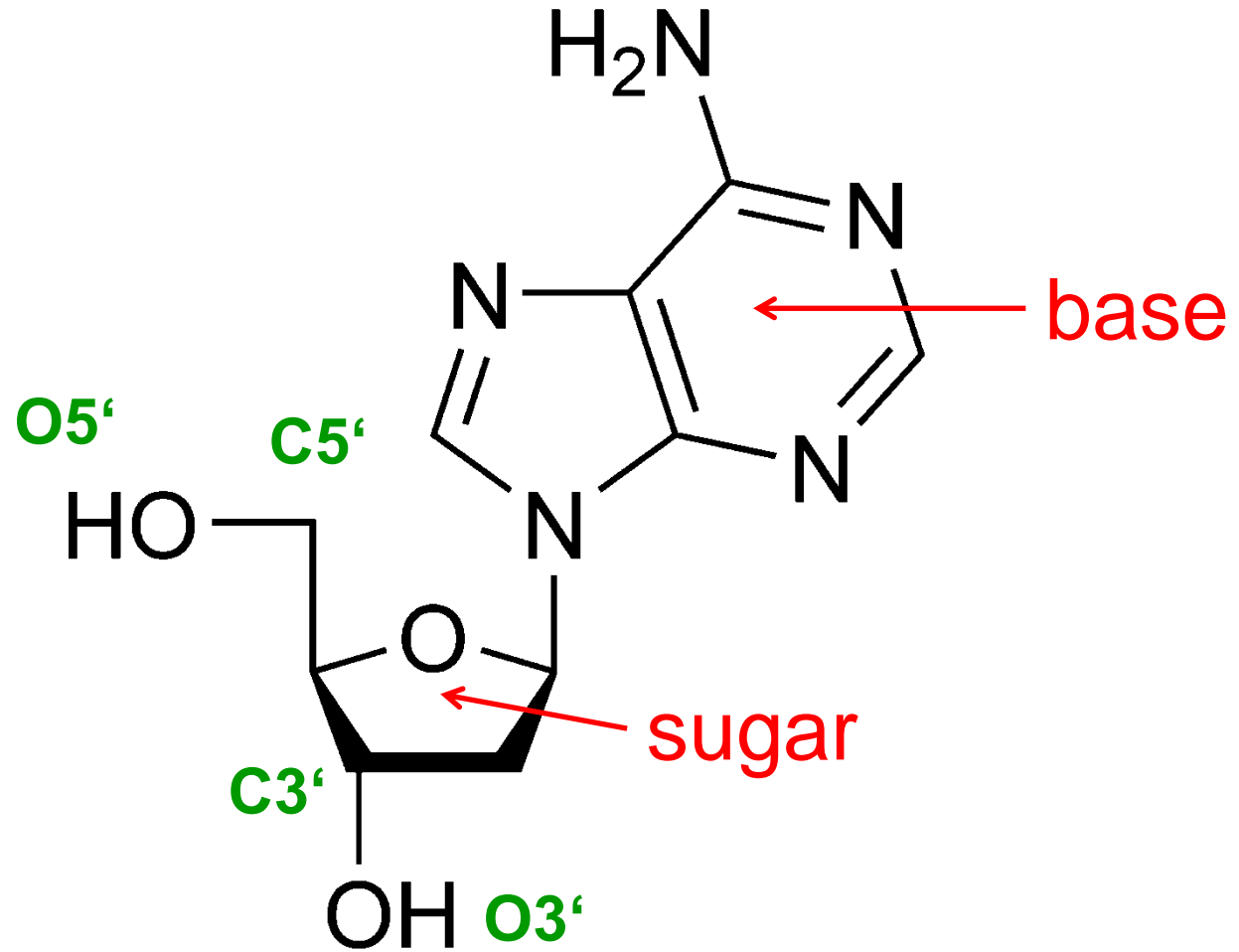
- Primary structure
  - sequence of NA basis in chains
- Secondary structure
  - set of interactions between nucleic basis
- Tertiary structure
  - 3D localization of atoms
- Quarternary structure
  - Higher organization levels
    - DNA in chromatin
    - Interaction of RNA units in ribosome or spliceosome.

# DNA – DeoxyriboNucleic Acid

- bases, deoxyribose sugar, phosphate – nucleotide
- Bases are flat → stacking
- pYrimidines – C, T
- puRines – A, G

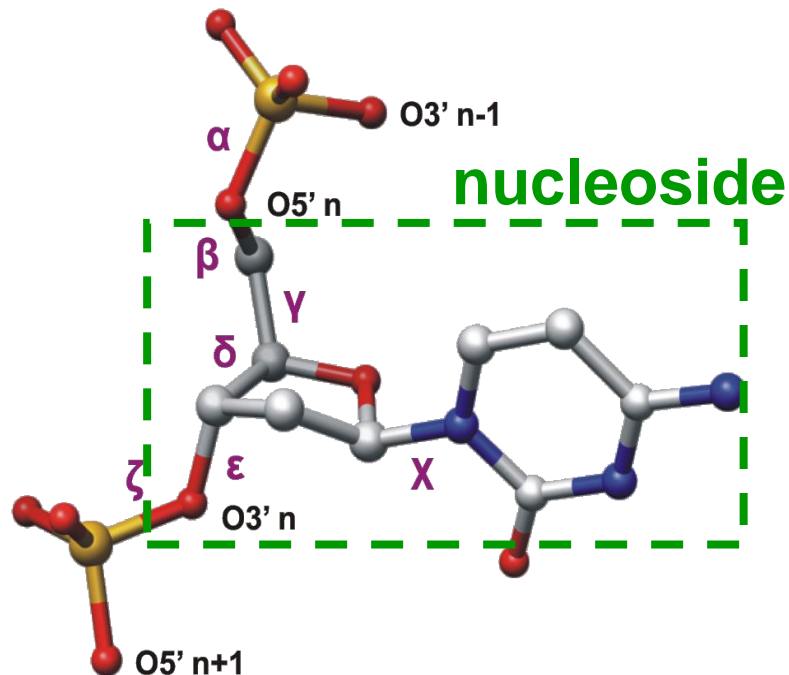


# Nucleoside



# Nucleotide

- nucleosides are interconnected by phosphodiester bond
- nucleotide monophosphate

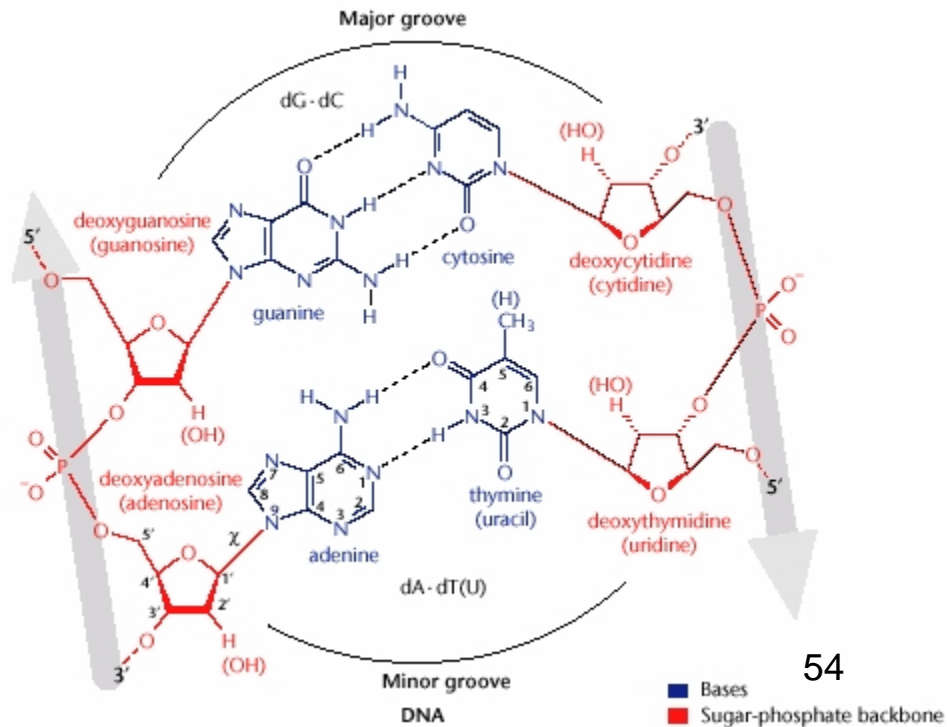
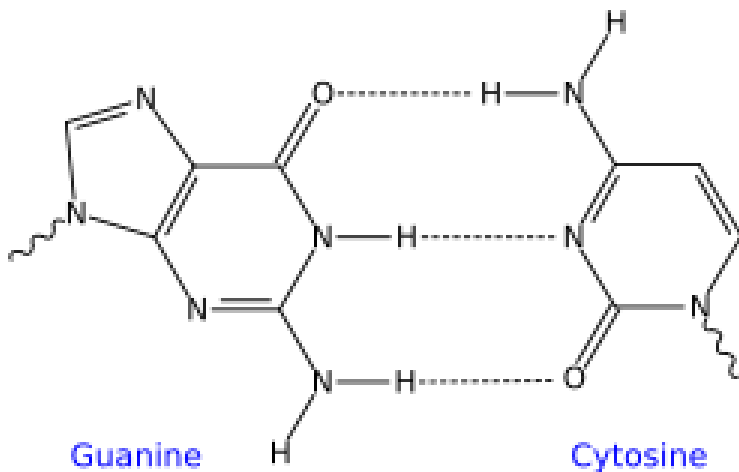
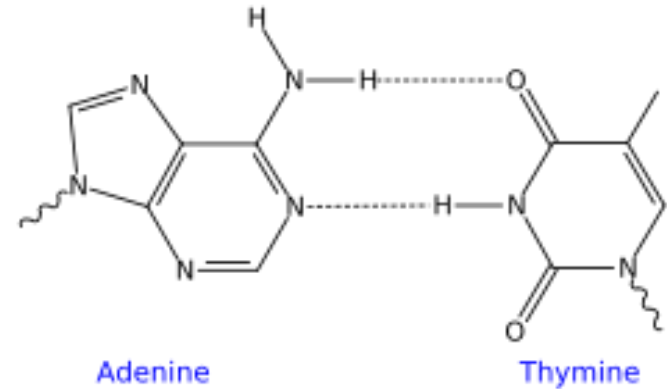


# Watson-Crick pairing

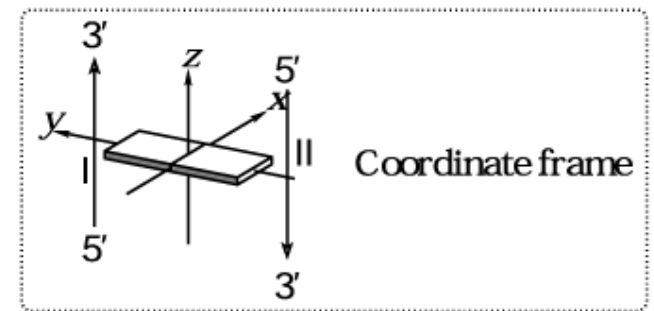
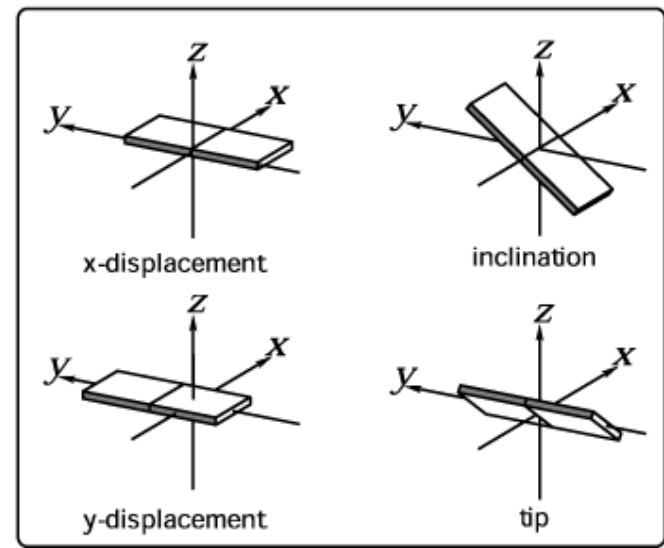
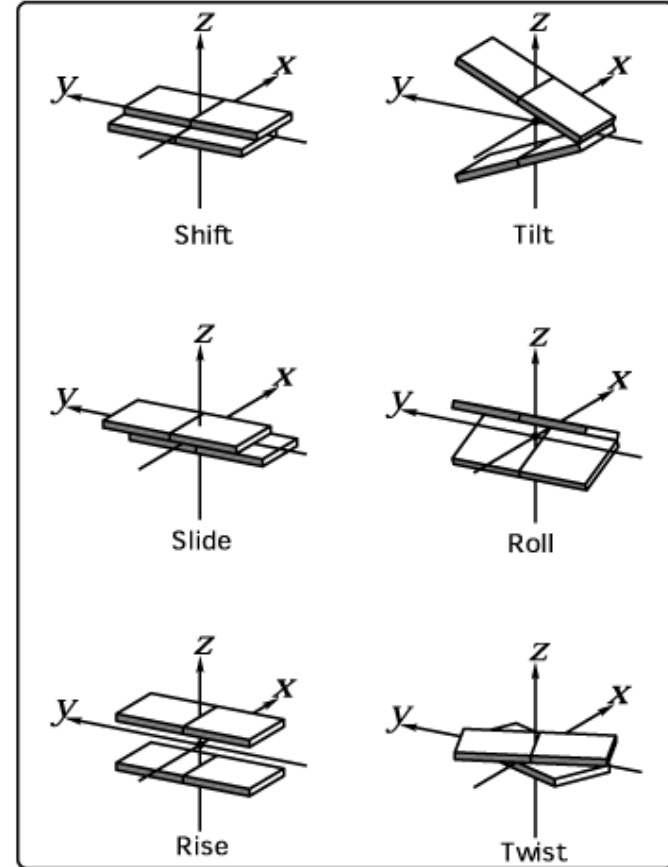
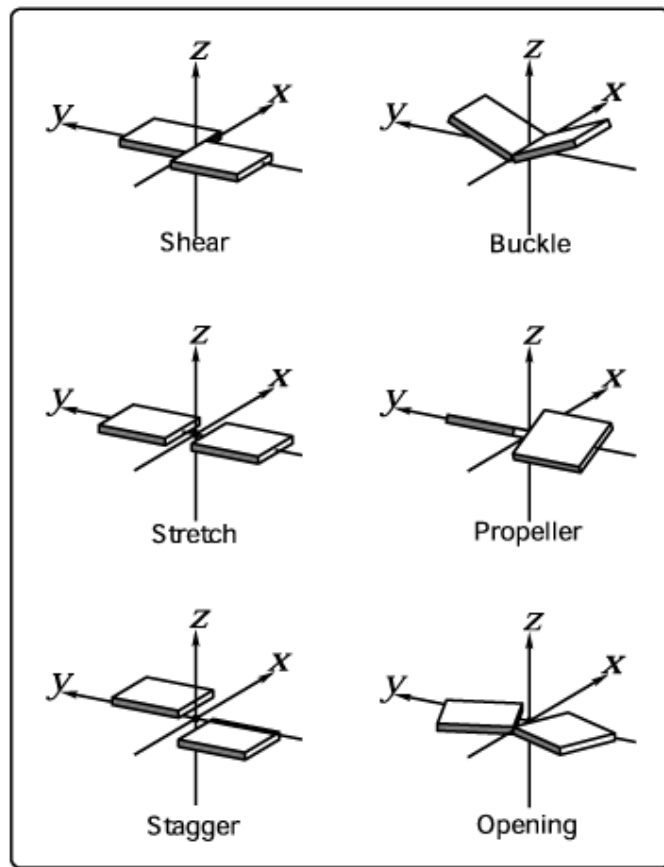
Bases complement each other.

## Chargaff's rules

- amount of G = C
- amount of A = T

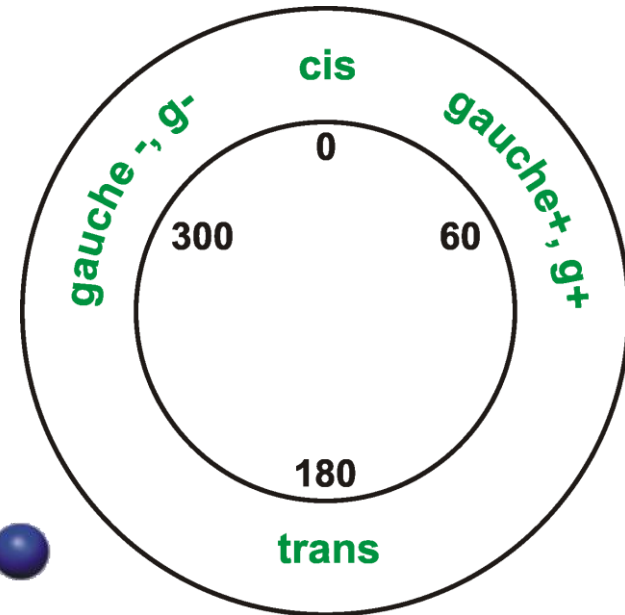
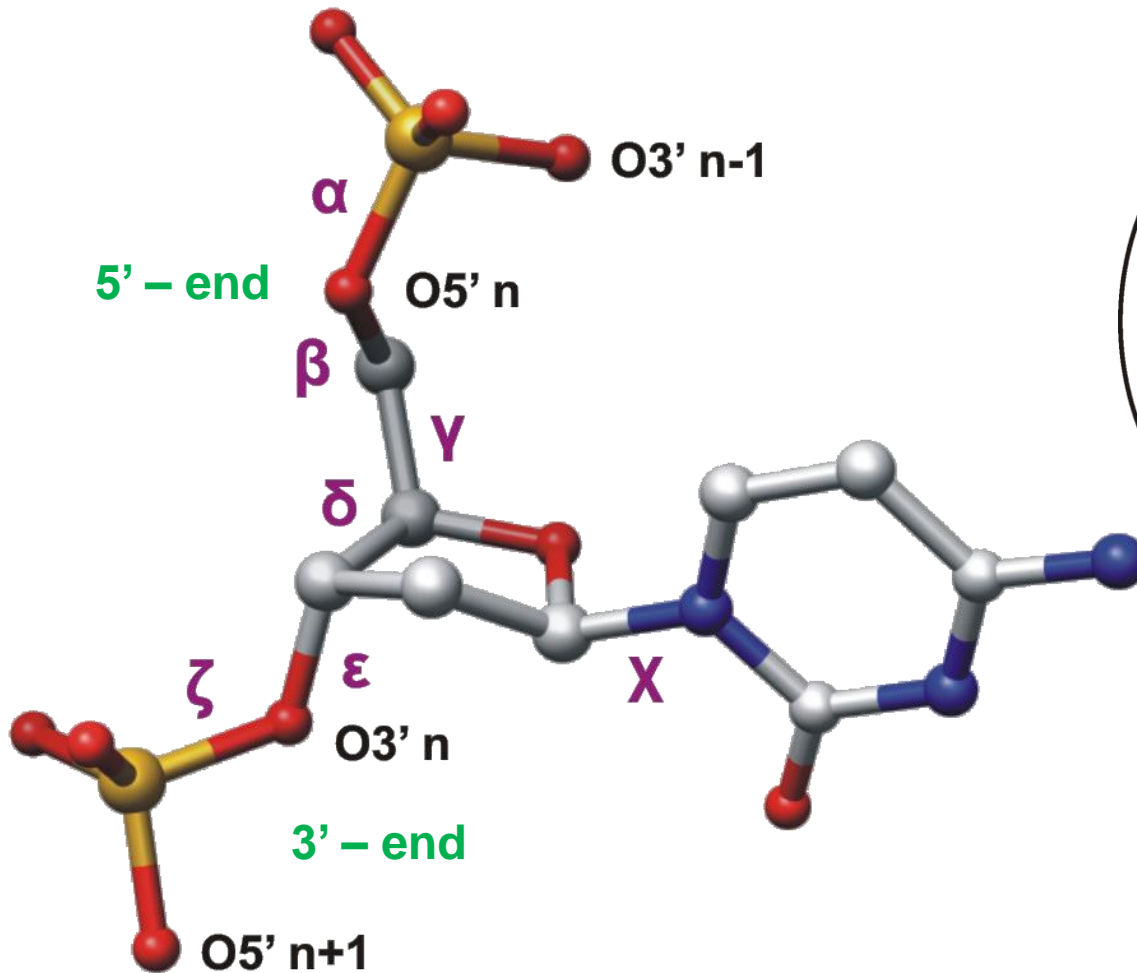


# Párování



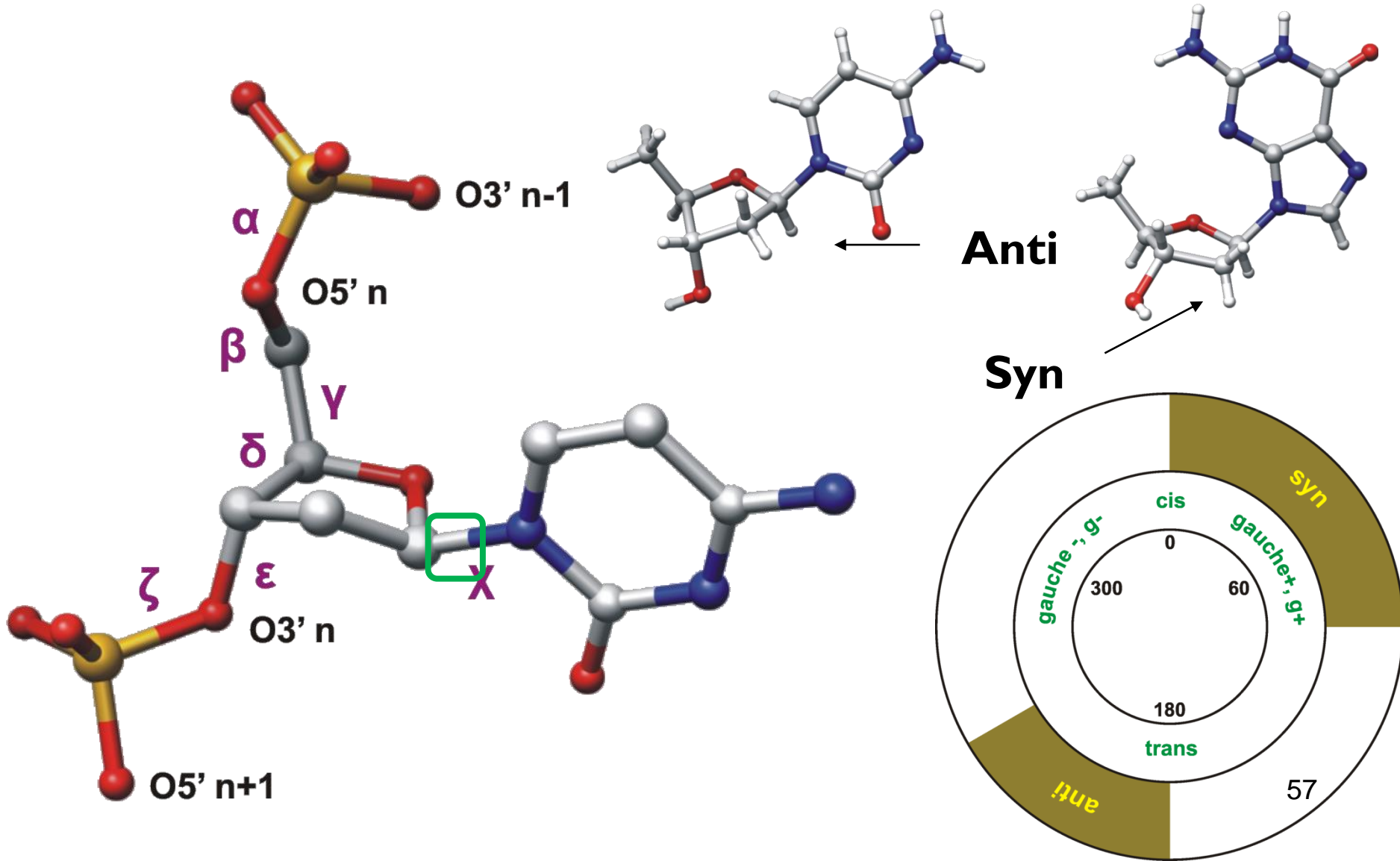
Images created with 3DNA illustrating positive values of designated parameters

# DNA backbone

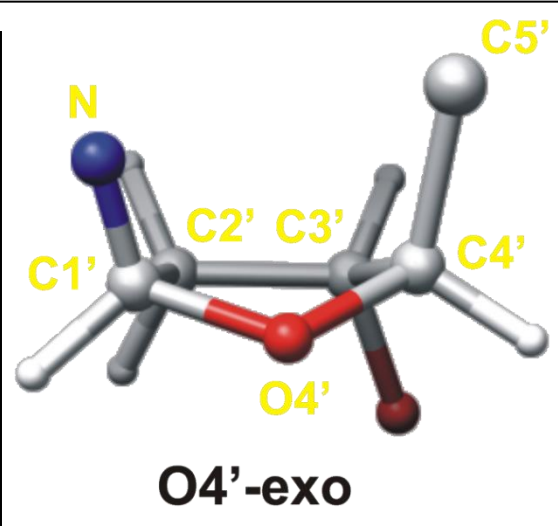
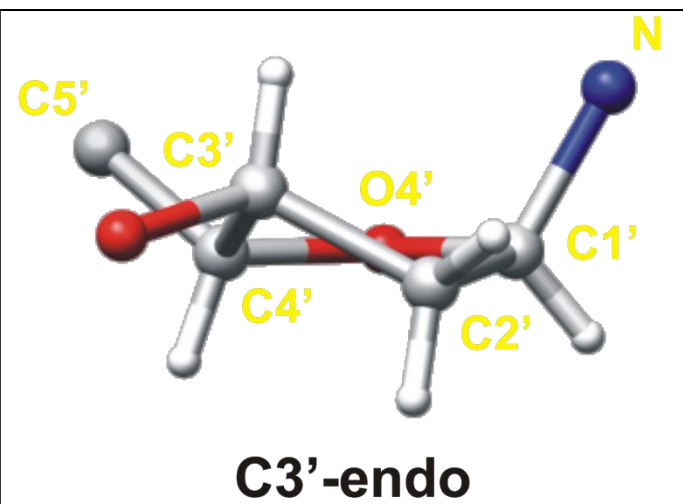
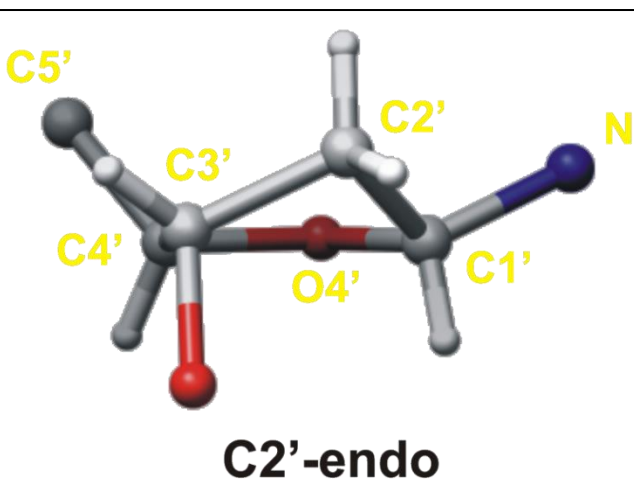




# Base at sugar dihedrals

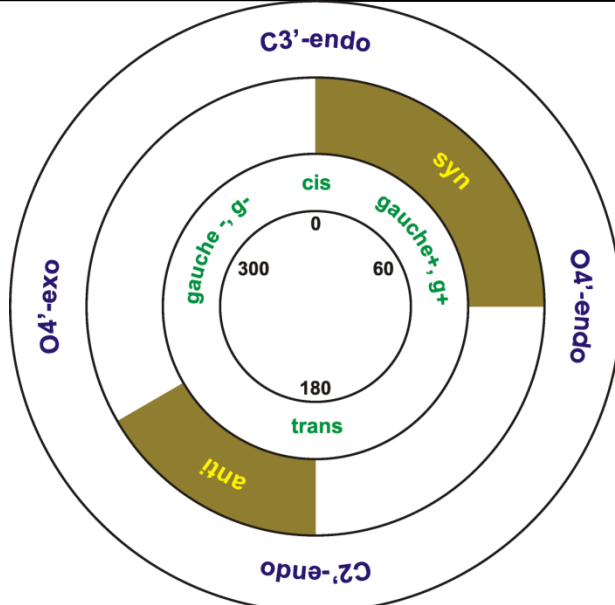


# Sugar conformation

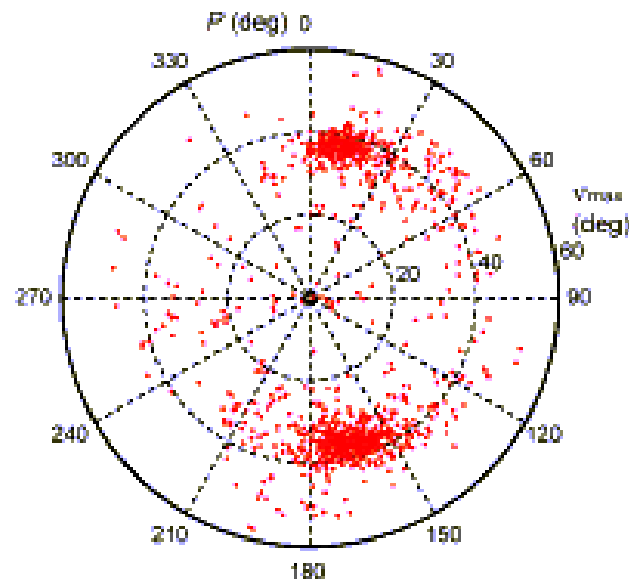
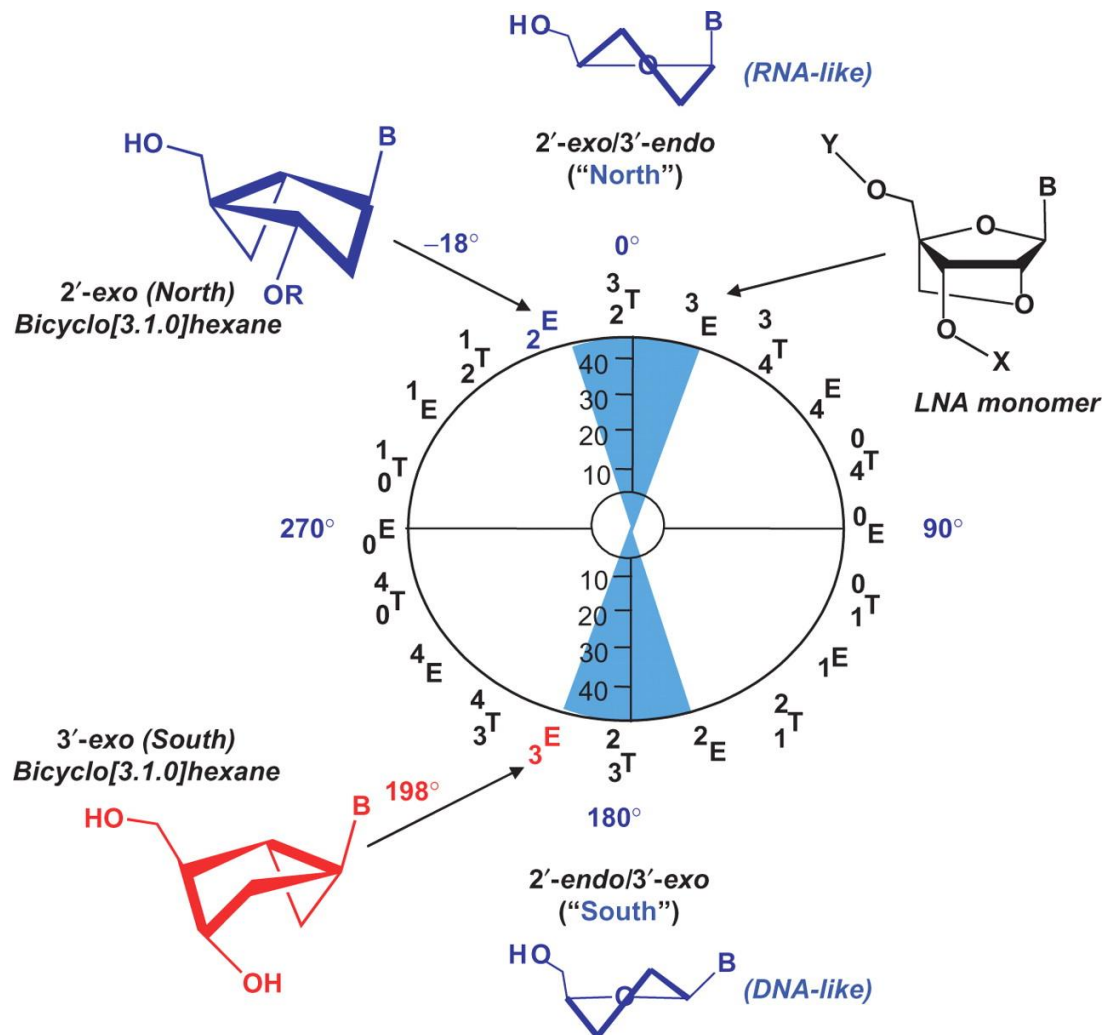


orientation with respect to C5'

- same side – endo
- opposite side – exo

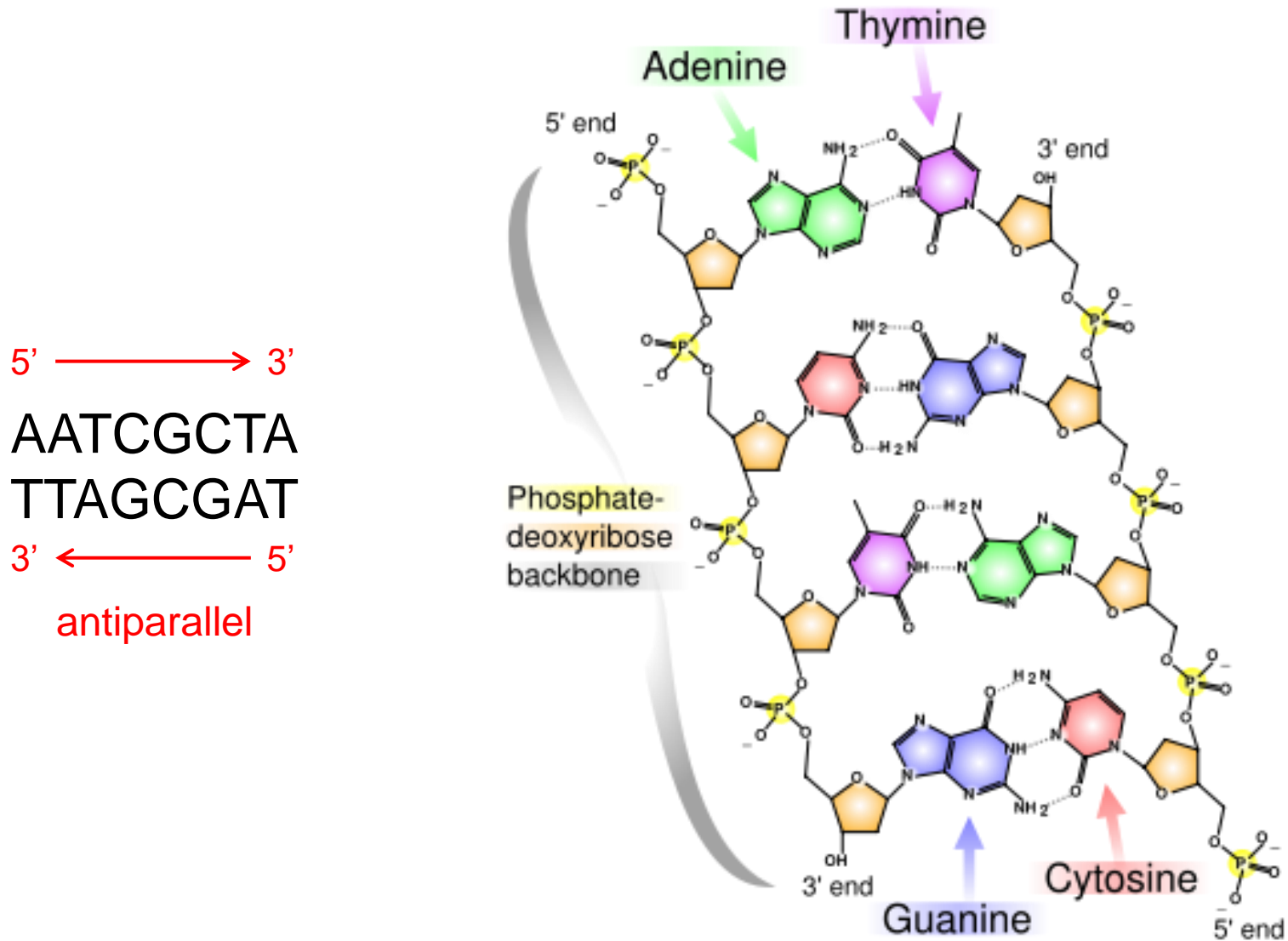


# Pseudorotational cycle for furanose ring puckers.

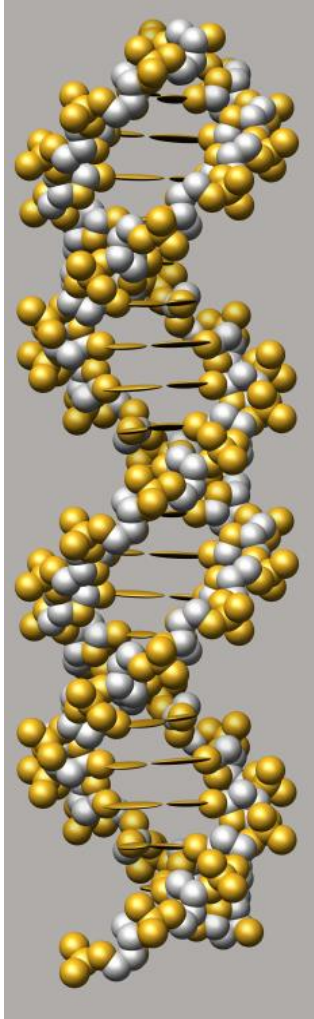


Pucker conformation of sugars in CSD database from PROSIT server

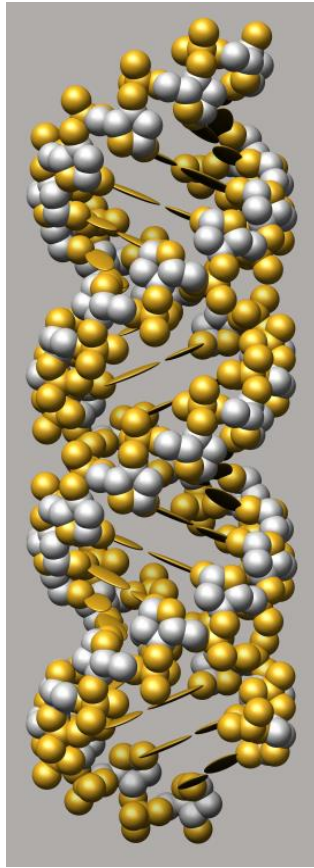
# DNA Double helix



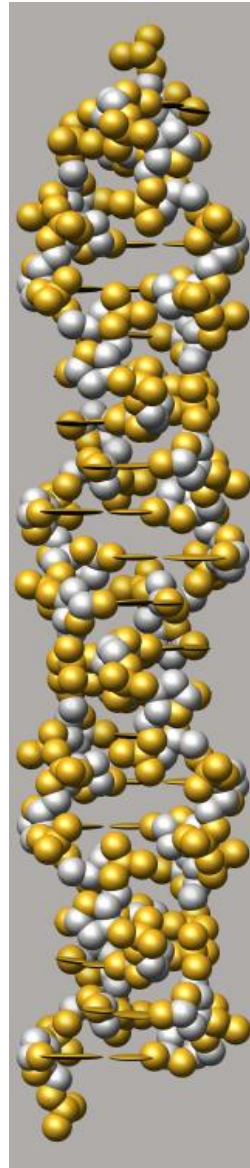
# Types of DNA



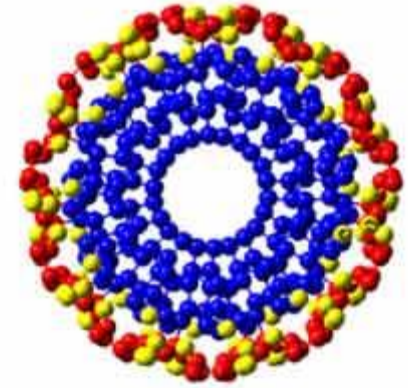
**B-DNA**



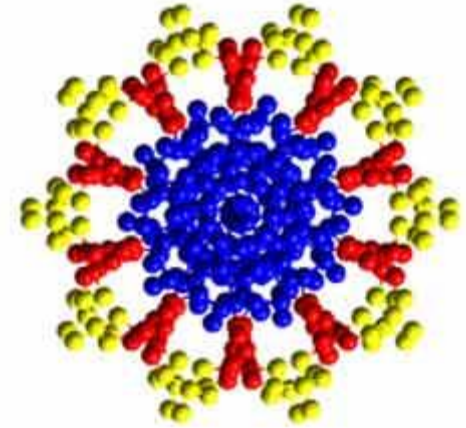
**A-DNA**



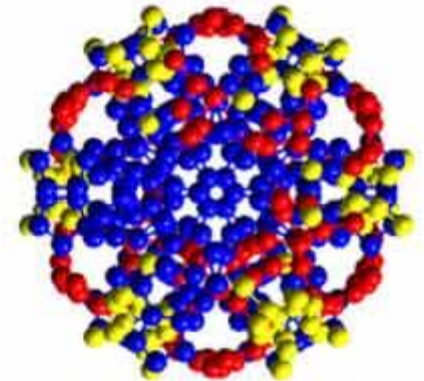
**Z-DNA**



**A**



**B**



**Z**

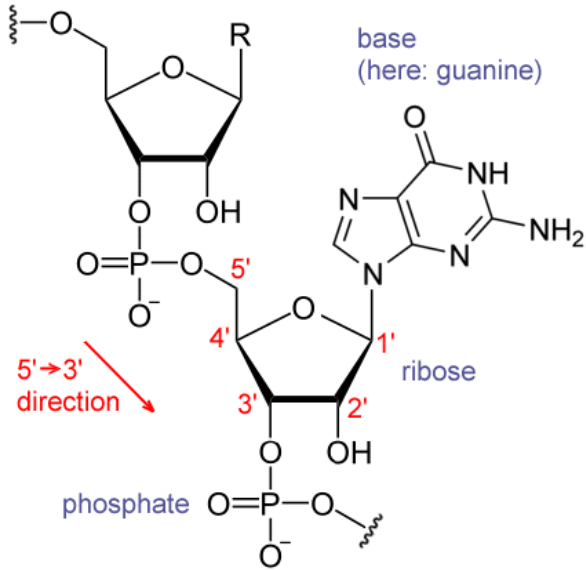
# Biological role of different DNAs

- B-DNA
  - canonical DNA
  - predominant
- A-DNA
  - Conditions of lower humidity, common in crystallographic experiments. However, they're artificial.
  - In vivo – local conformations induced e.g. by interaction with proteins.
- Z-DNA
  - No definite biological significance found up to now.
  - It is commonly believed to provide torsional strain relief (supercoiling) while DNA transcription occurs.
  - The potential to form a Z-DNA structure also correlates with regions of active transcription.

# Different sets of DNA

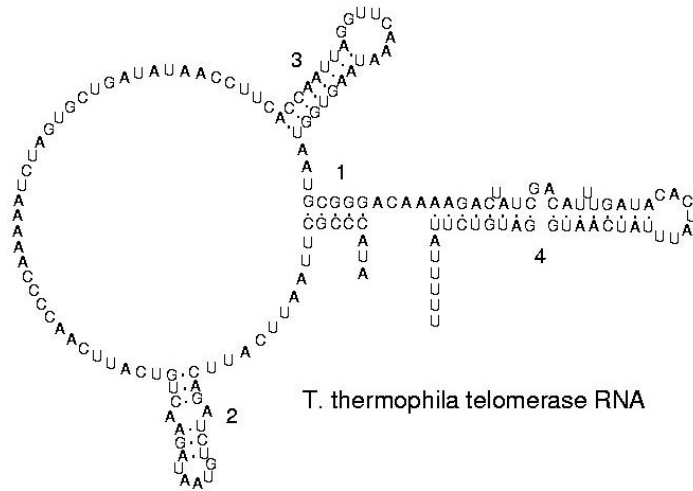
- nuclear DNA
  - cell's nucleus
  - majority of functions cell carries out
  - sequencing the genome – scientists mean nuclear DNA
- mitochondrial DNA
  - *mtDNA*
  - circular, in human very short (17 kbp) with 37 genes (controlling cellular metabolism)
  - all *mtDNA* comes from mom
- chloroplast DNA
  - *cpDNA*
  - circular and fairly large (120 – 160 kbp), with only 120 genes
  - inheritance is either maternal, or paternal

# RNA - ribonucleic acid



primární struktura

sekundární struktura



terciární struktura



hammerhead  
ribozyme 2GOZ



## RNAs involved in protein synthesis

Type	Abbr.	Function	Distribution	Ref.
Messenger RNA	mRNA	Codes for protein	All organisms	
Ribosomal RNA	rRNA	Translation	All organisms	
Signal recognition particle RNA	7SL RNA or SRP RNA	Membrane integration	All organisms	[1]
Transfer RNA	tRNA	Translation	All organisms	
Transfer-messenger RNA	tmRNA	Rescuing stalled ribosomes	Bacteria	[2]

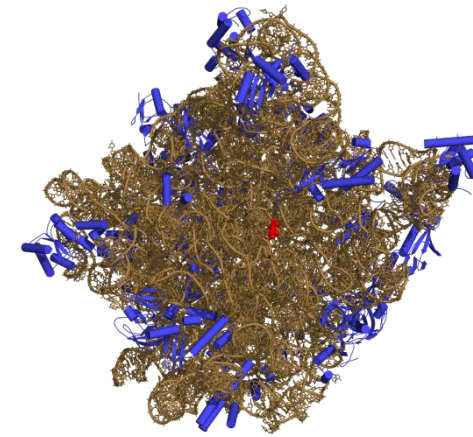
# RNA

## RNAs involved in post-transcriptional modification or DNA replication

Type	Abbr.	Function	Distribution	Ref.
Small nuclear RNA	snRNA	Splicing and other functions	Eukaryotes and archaea	[3]
Small nucleolar RNA	snoRNA	Nucleotide modification of RNAs	Eukaryotes and archaea	[4]
SmY RNA	SmY	mRNA trans-splicing	Nematodes	[5]
Small Cajal body-specific RNA	scaRNA	Type of snoRNA; Nucleotide modification of RNAs		
Guide RNA	gRNA	mRNA nucleotide modification	Kinetoplastid mitochondria	[6]
Ribonuclease P	RNase P	tRNA maturation	All organisms	[7]
Ribonuclease MRP	RNase MRP	rRNA maturation, DNA replication	Eukaryotes	[8]
Y RNA		RNA processing, DNA replication	Animals	[9]
Telomerase RNA		Telomere synthesis	Most eukaryotes	[10]



pre-mRNA hairpin



50S-ribozome

## Regulatory RNAs

Type	Abbr.	Function	Distribution	Ref.
Antisense RNA	aRNA	Transcriptional attenuation / mRNA degradation / mRNA stabilisation / Translation block	All organisms	[11][12]
Cis-natural antisense transcript		Gene regulation		
CRISPR RNA	crRNA	Resistance to parasites, probably by targeting their DNA	Bacteria and archaea	[13]
Long noncoding RNA	Long ncRNA	Various	Eukaryotes	
MicroRNA	miRNA	Gene regulation	Most eukaryotes	[14]
Piwi-interacting RNA	piRNA	Transposon defense, maybe other functions	Most animals	[15][16]
Small interfering RNA	siRNA	Gene regulation	Most eukaryotes	[17]
Trans-acting siRNA	tasRNA	Gene regulation	Land plants	[18]
Repeat associated siRNA	rasRNA	Type of piRNA; transposon defense	Drosophila	[19]
7SK RNA	7SK	negatively regulating CDK9/cyclin T complex		



hammerhead ribozyme

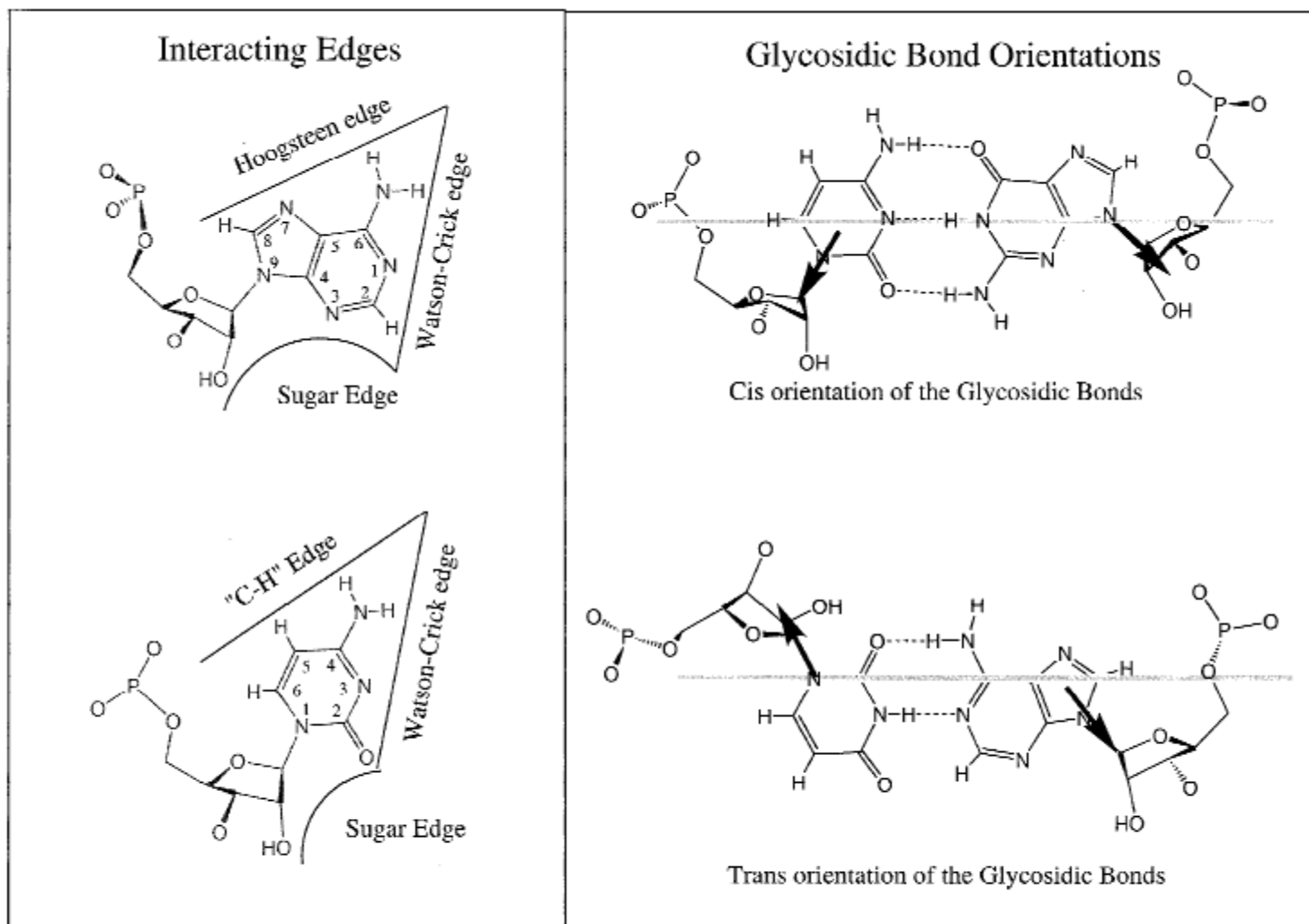
2GOZ

## Parasitic RNAs

Type	Function	Distribution	Ref.
Retrotransposon	Self-propagating	Eukaryotes and some bacteria	[20]
Viral genome	Information carrier	Double-stranded RNA viruses, positive-sense RNA viruses, negative-sense RNA viruses, many satellite viruses and reverse transcribing viruses	
Viroid	Self-propagating	Infected plants	[21]
Satellite RNA	Self-propagating	Infected cells	

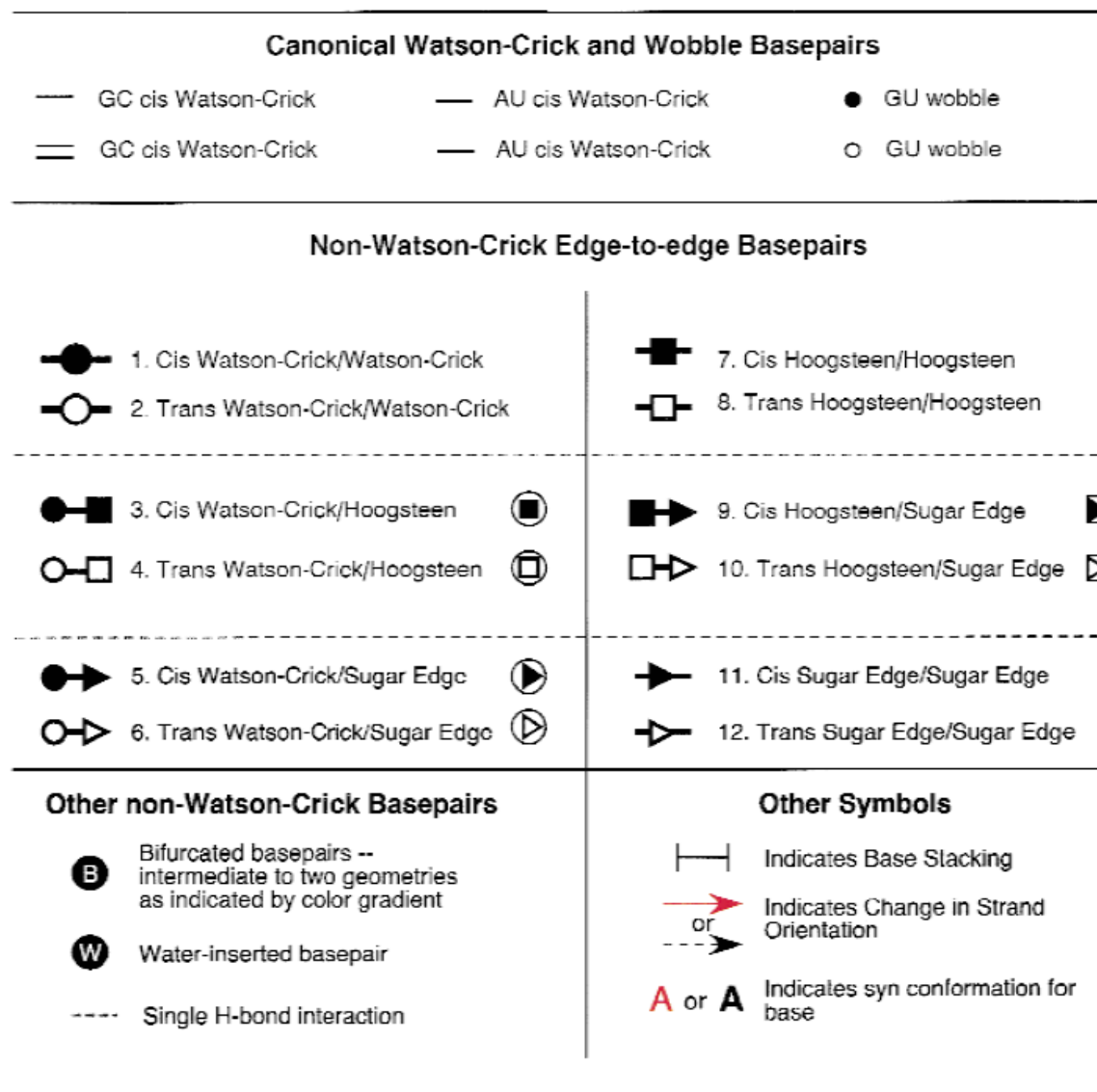
## Other RNAs

Type	Abbr.	Function	Distribution	Ref.
Vault RNA	vRNA	Expulsion of xenobiotics, maybe		[22]



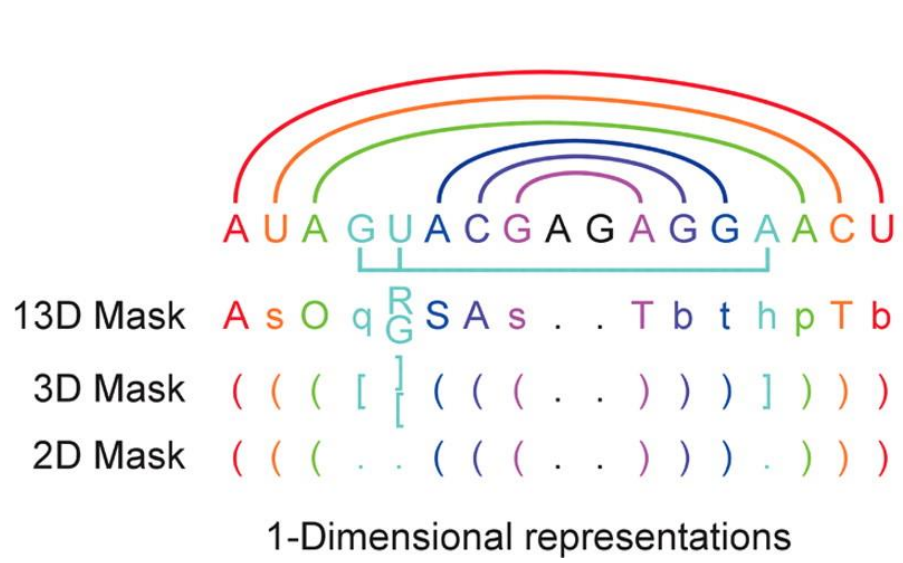
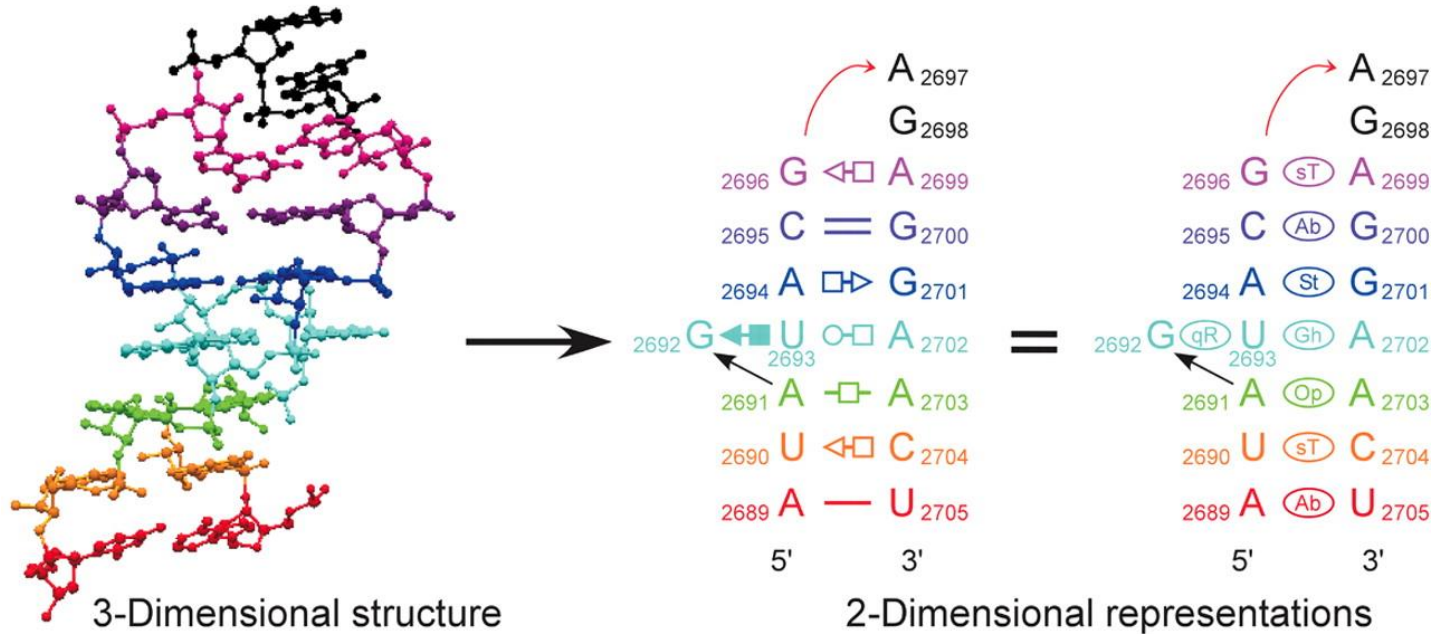
**FIGURE 1.** Left panel: Purine (A or G, indicated by "R") and pyrimidine (C or U, indicated by "Y") bases provide three edges for interaction, as shown for adenosine and cytosine. The Watson-Crick edge comprises A(N6)/G(O6), R(N1), A(C2)/G(N2), U(O4)/C(N4), Y(N3), and Y(O2). The Hoogsteen edge comprises A(N6)/G(O6), R(N7), U(O4)/C(N4), and Y(C5). The Sugar-edge comprises A(C2)/G(N2), R(N3), Y(O2), and the ribose hydroxyl group, O2'. Right panel: The *cis* and *trans* orientations are defined relative to a line drawn parallel to and between the *base-to-base* hydrogen bonds in the case of two hydrogen bonds or, in the case of three hydrogen bonds, along the middle hydrogen bond.

# Annotation of 2D RNA Structures

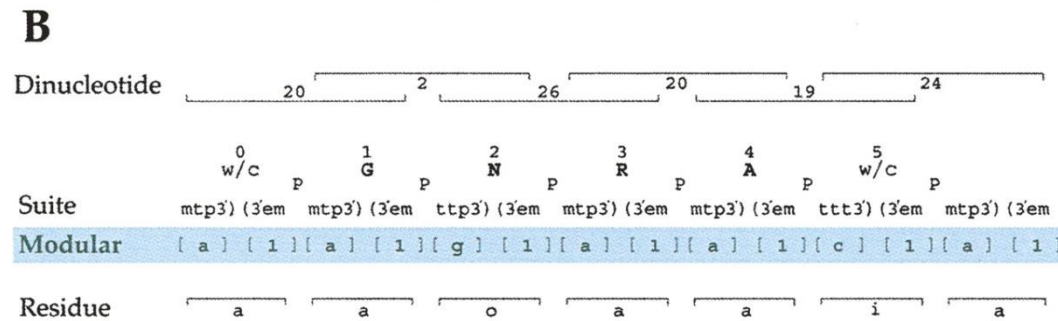
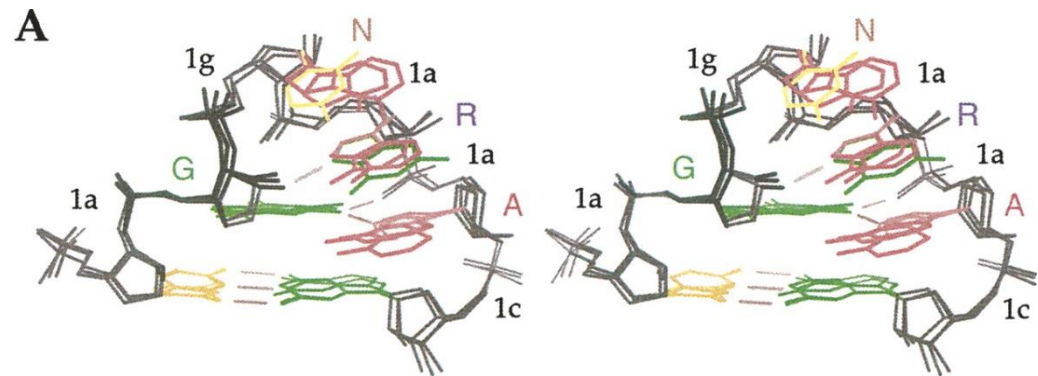
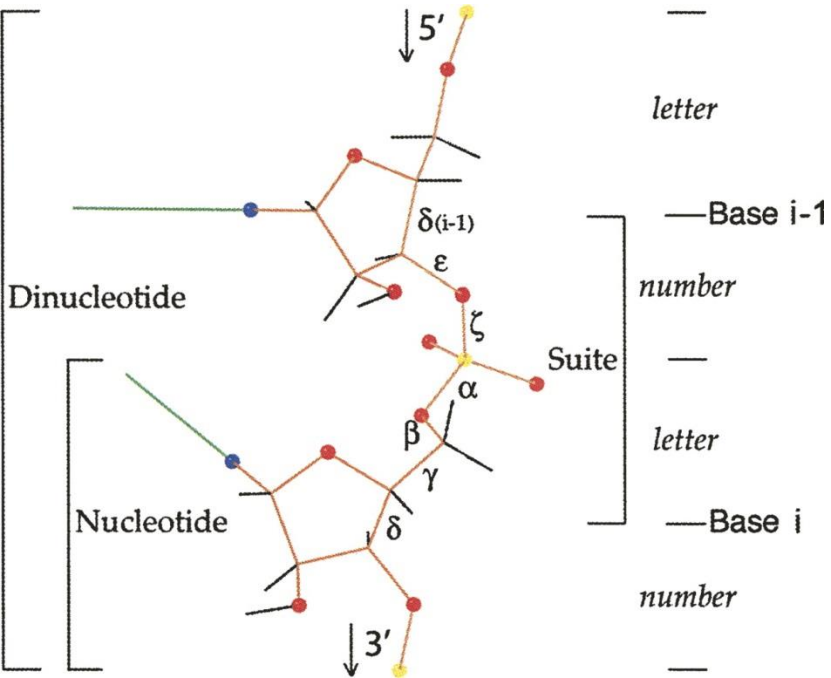


**FIGURE 6.** Suggested symbols for indicating tertiary interactions and other three-dimensional structural features in two-dimensional representations of RNA structures.

# RNA Representations

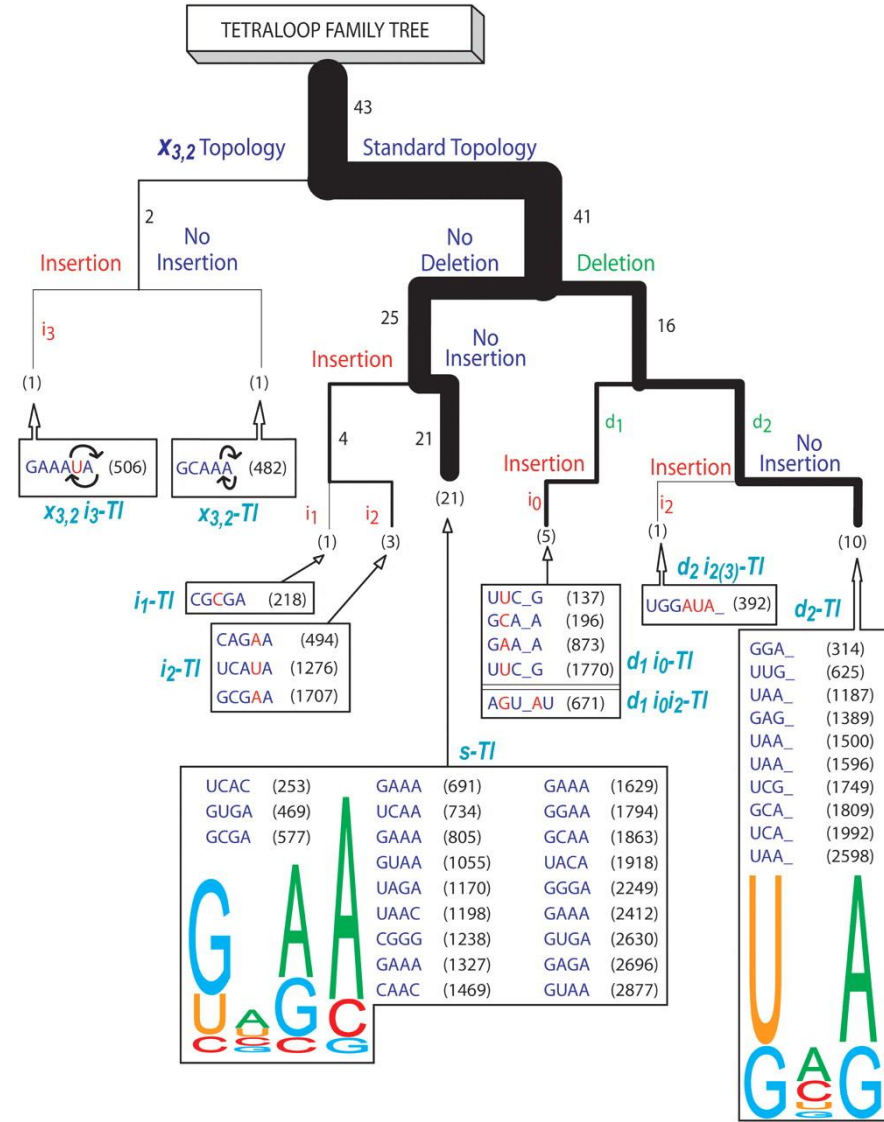
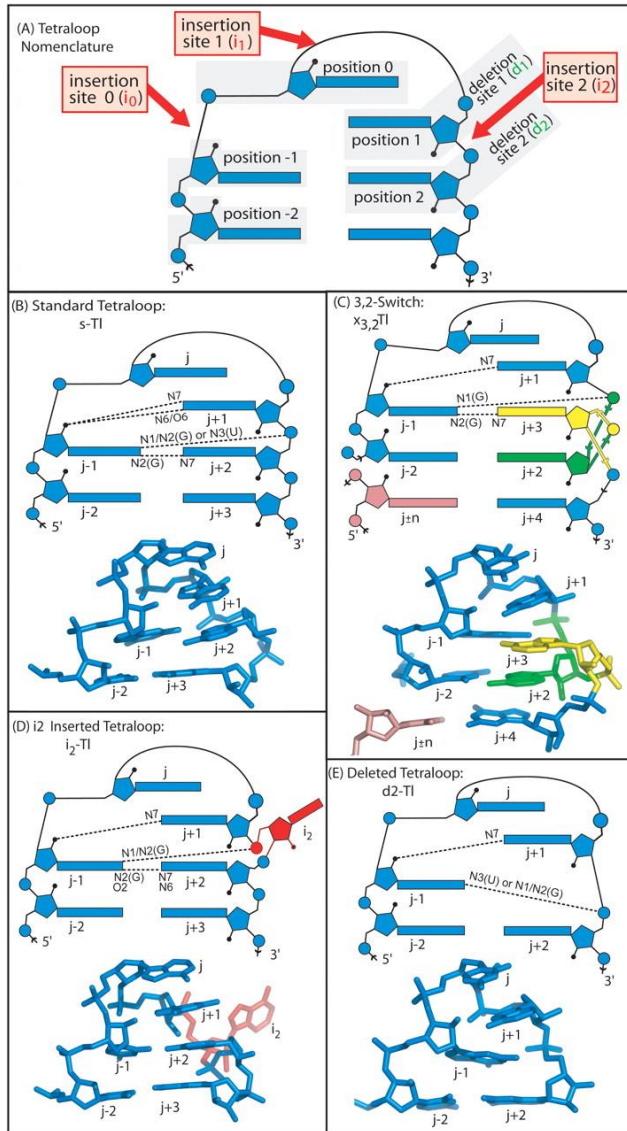


# RNA Backbone



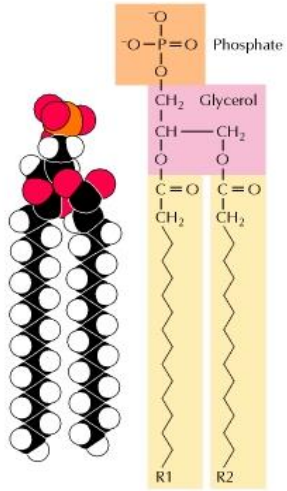
sequence/conformation string: N1aG1gN1aR1aA1cN1a

# RNA Tetraloop Family Tree.

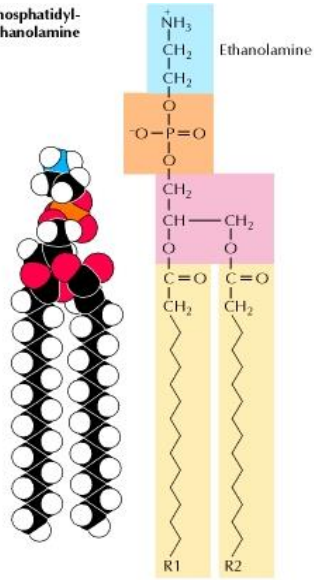


# Lipids

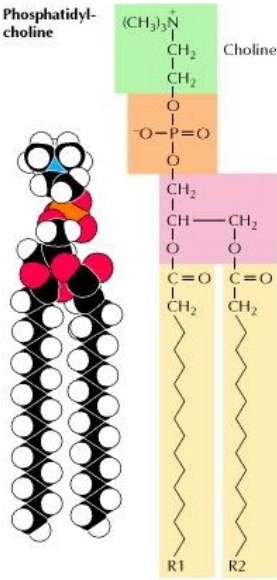
Phosphatidic acid



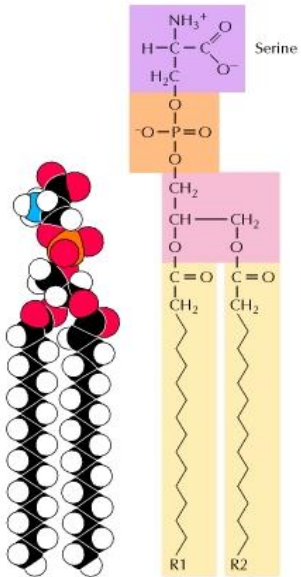
Phosphatidylethanolamine



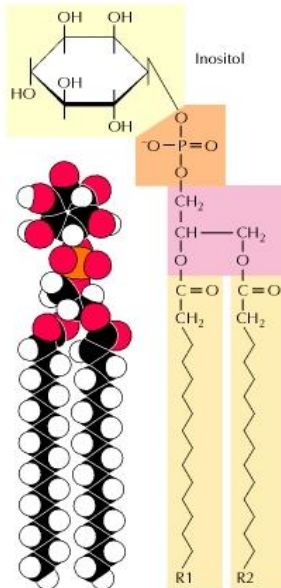
Phosphatidylcholine



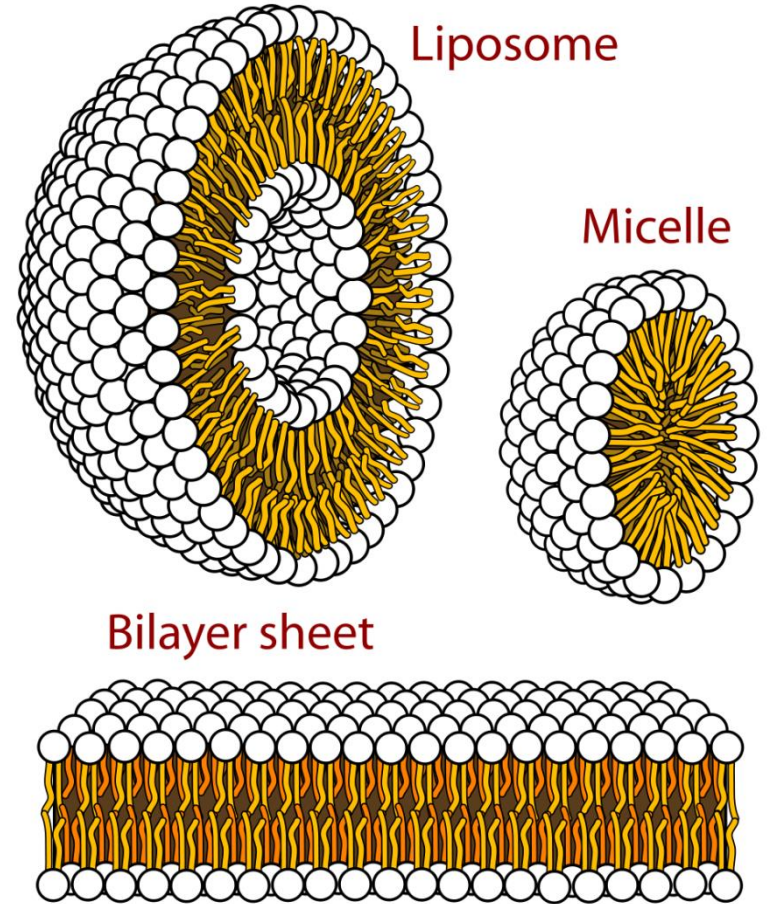
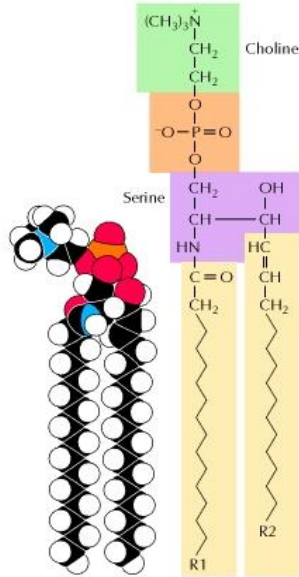
Phosphatidylserine



Phosphatidylinositol



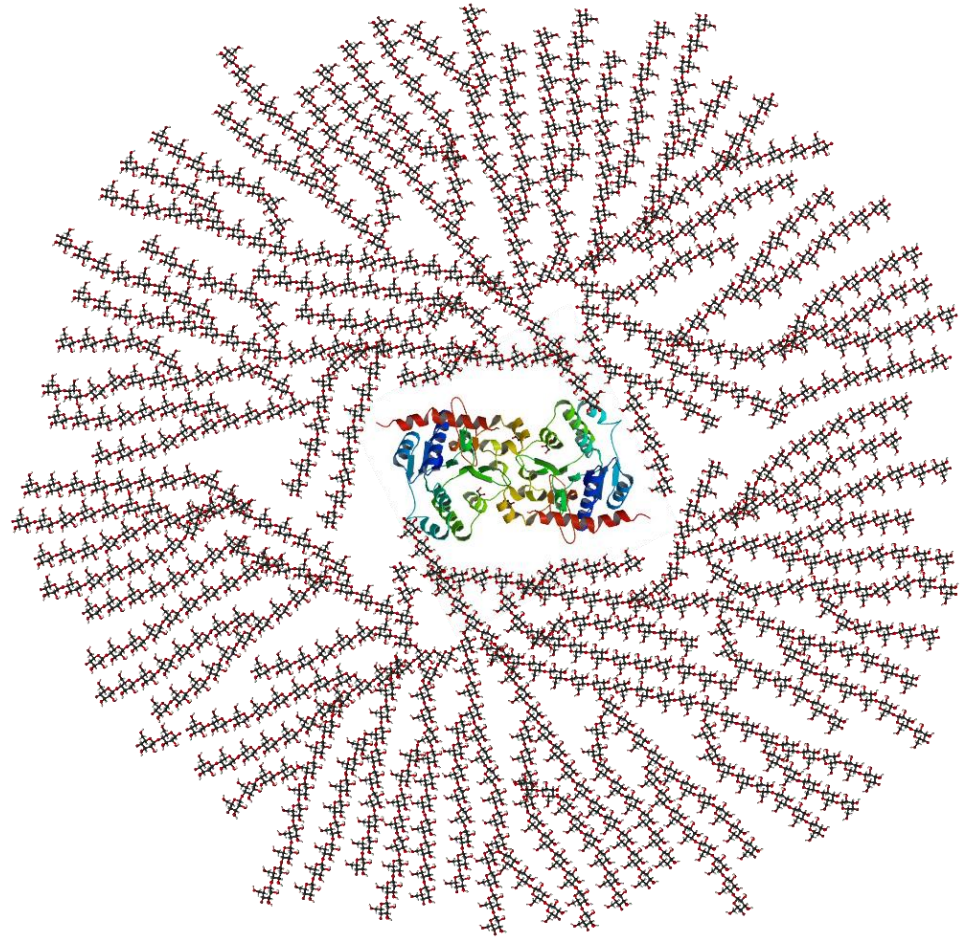
Sphingomyelin



main phospholipids

# Polysaccharides

- role:
  - Energy storage
  - Molecular recognition
- Harder to read in sequences than NA or proteins
- Quite often on extracellular proteins



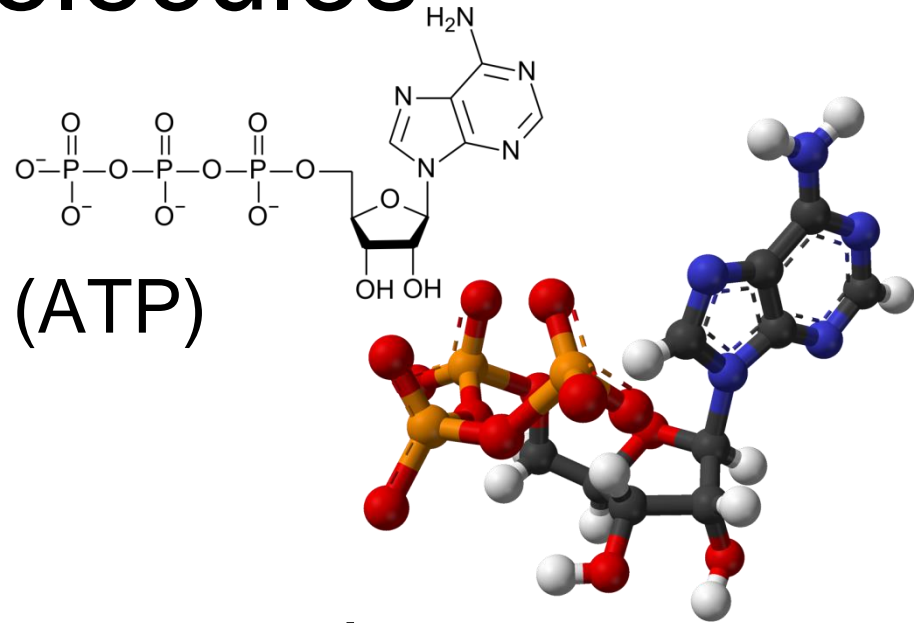
glycogen



# Small molecules

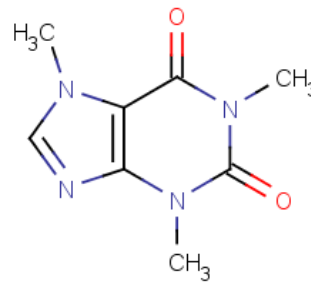
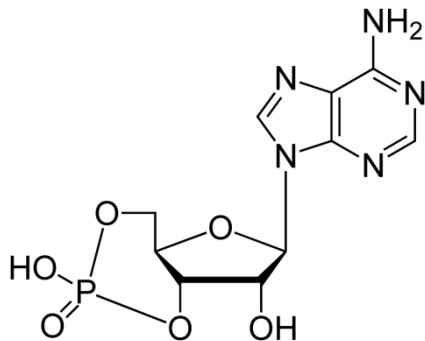
- NTP

- Cell energy transporter (ATP)
- Basic stones for NA

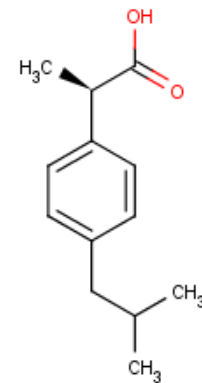


- Messengers, Agonists, antagonists

- (cAMP, xenobiotics)



caffeine



ibuprofen