

Topological modelling and analysis for biomolecular data: 1 from graph to simplicial complex

Kelin Xia

*School of Physical and Mathematical Sciences,
Nanyang Technological University*

应用拓扑短期课程与报告会, 12月13-16, 2020

Fund: NTU-JSPS(2019), MOE-Tier 1(2018,2019), MOE-Tier 2(2018,2021), Alibaba-NTU(2020), Merlion(2020)

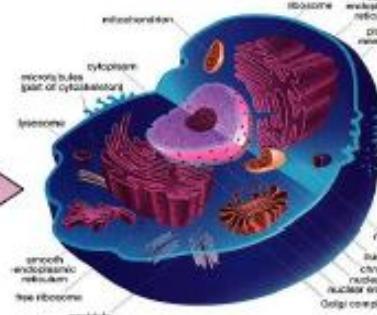
Species



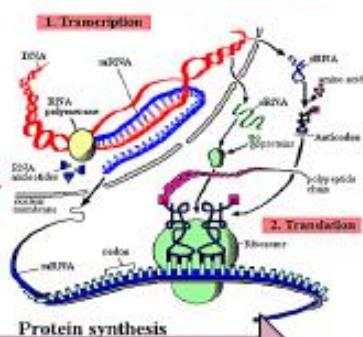
Organs



Cells



Molecules



About 20 orders in spatial scales

About 20 orders in time scales

Evolutionary biology

Reaction diffusion
Stochastic models
Kinetic models
Delayed ODEs
Discrete models
Homology models
Machine learning

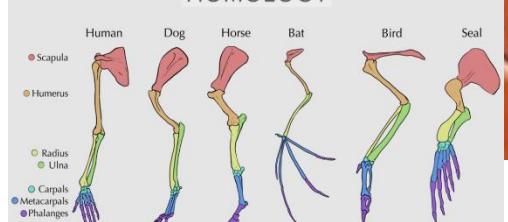
Developmental biology Physiology Biomechanics

Continuum models
Mechanical models
Navier-Stokes
(Non-) linear elasticity
Maxwell's equation
Thermal models
Rheological models
Hodgkin-Huxley model
Lattice models
Neural networks
Geometric models
Topological models

Cellular biology Systems biology Cellular mechanics

Chemical kinetics (ODEs)
Gene regulatory network
Protein network
Neural networks
Hodgkin-Huxley model
FitzHugh-Nagumo model
Mechanical models
Reaction diffusion
Phase field models
Stochastic models
Statistical models
Monte Carlo
Combinatory
Topological models
Machine learning

HOMOLOGY



Molecular biology Biochemistry Biophysics

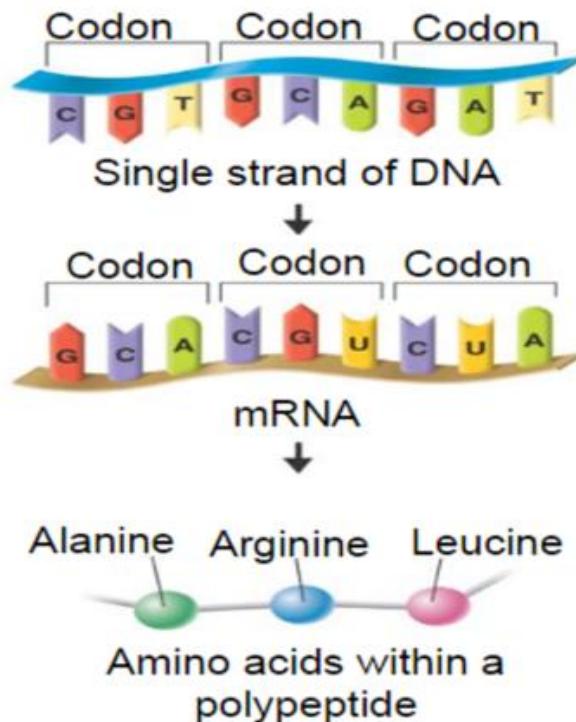
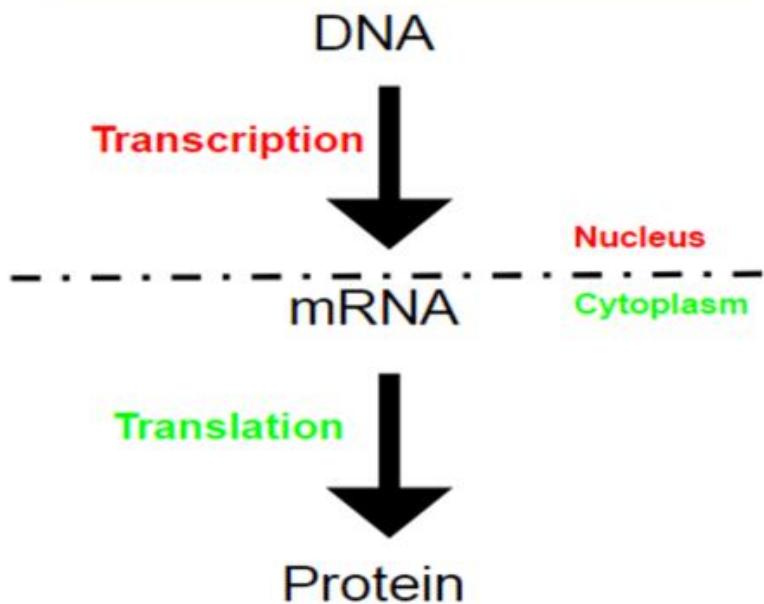
Molecular dynamics
Thermal dynamics
Brownian dynamics
Lagevin dynamics
Quantum models
QM/MM
Electrostatics
Implicit models
Boltzmann equation
Vlasov-Boltzmann
Fokker-Planck
Monte Carlo
Master equation
Homology models
Knot theory

Molecular biology

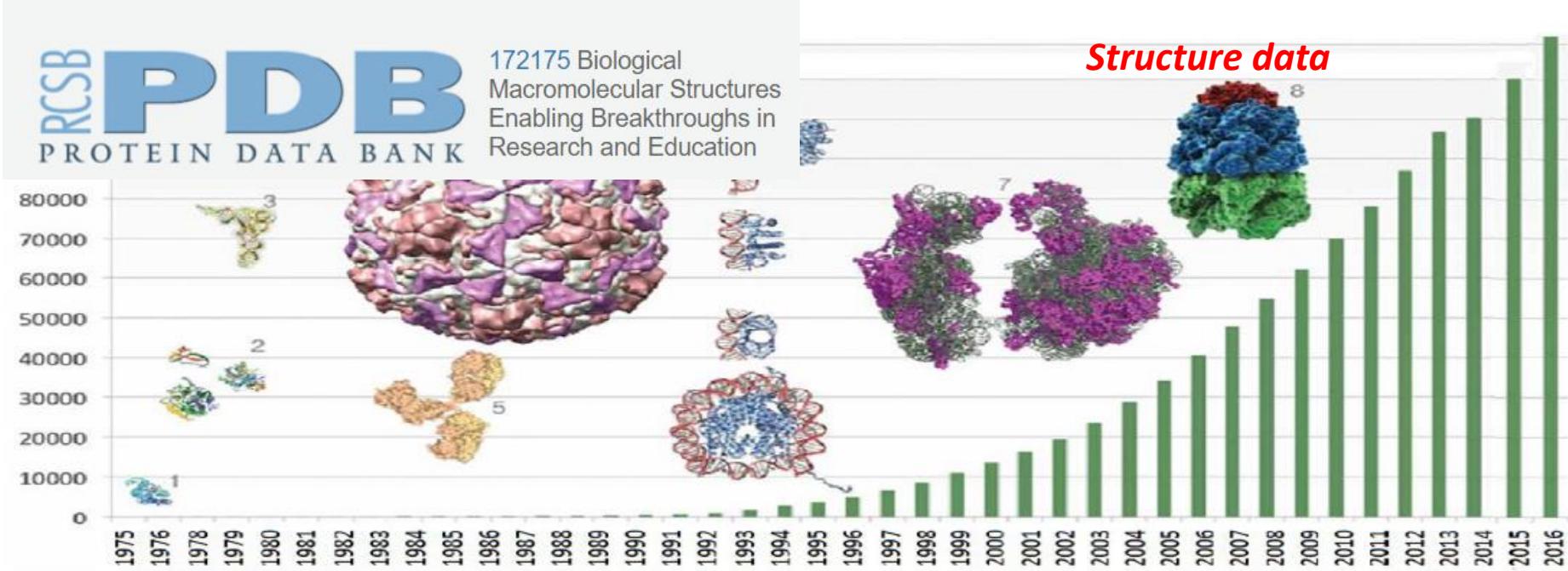
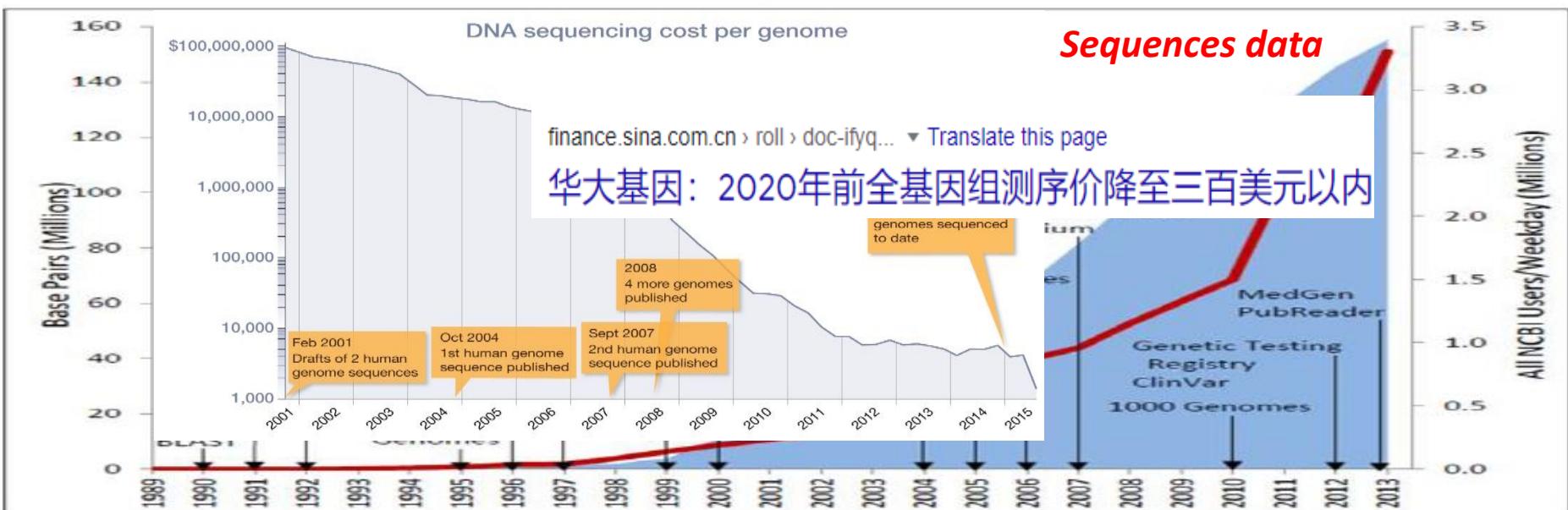


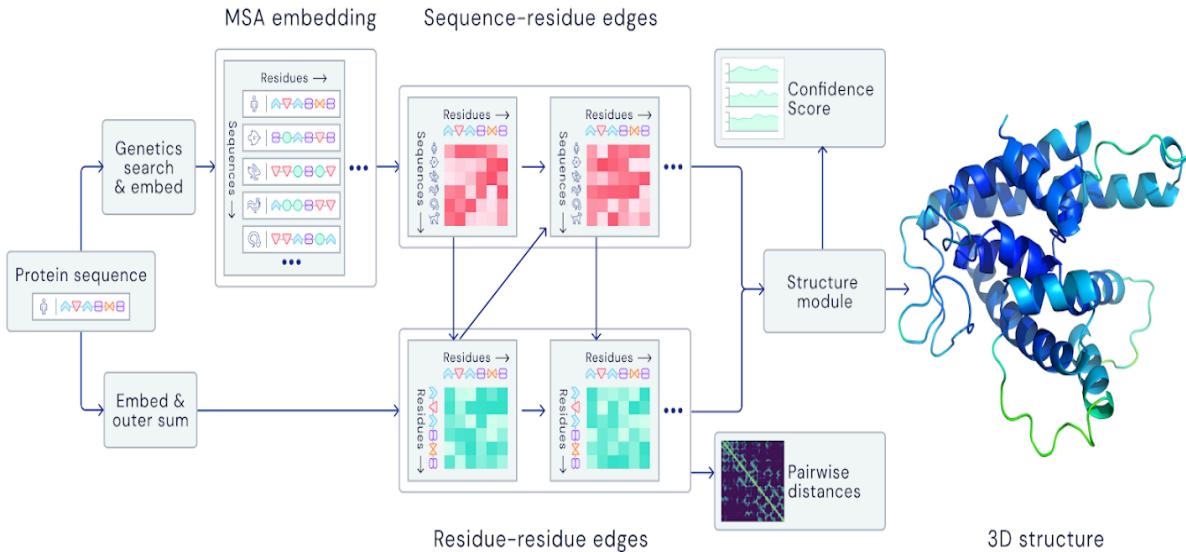
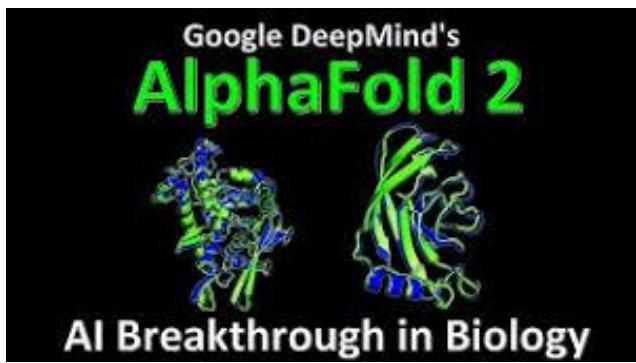
Francis Crick (1916 - 2004) and James Watson (1928 –) Nobel Prize in Physiology or Medicine in 1962 "for discoveries concerning the molecular structure of nucleic acids and its significance for information transfer in living material", i.e., the [Central Dogma](#) of molecular biology:

Central Dogma in a nutshell

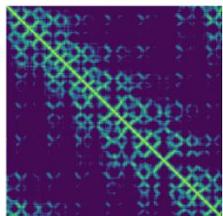


Biological data

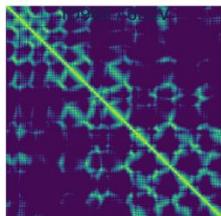




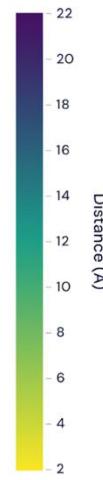
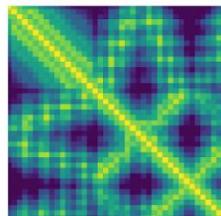
T0954 / 6CVZ



T0965 / 6D2V



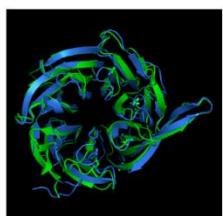
T0955 / 5W9F



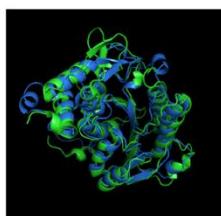
Ground truth

Average predicted distance

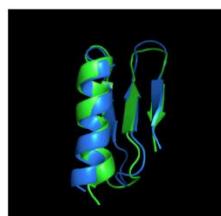
T0954 / 6CVZ



T0965 / 6D2V



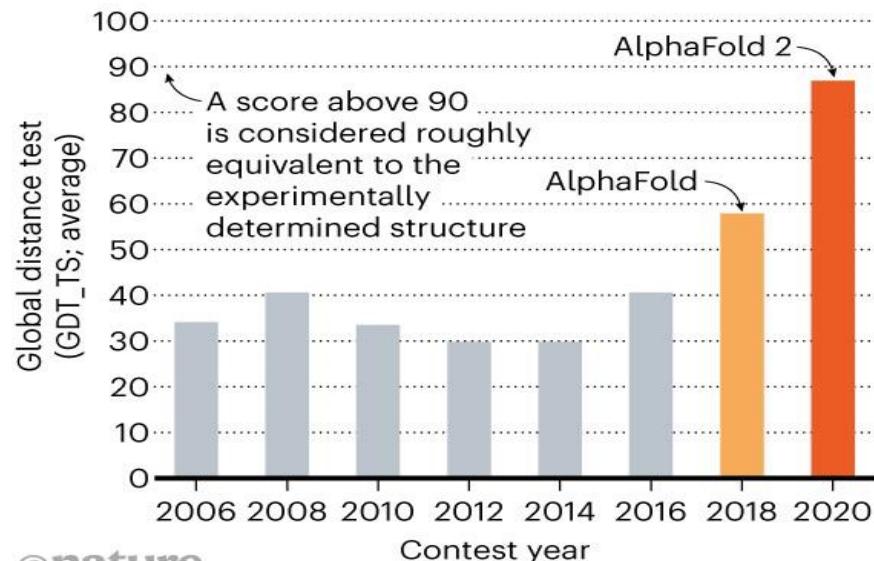
T0955 / 5W9F

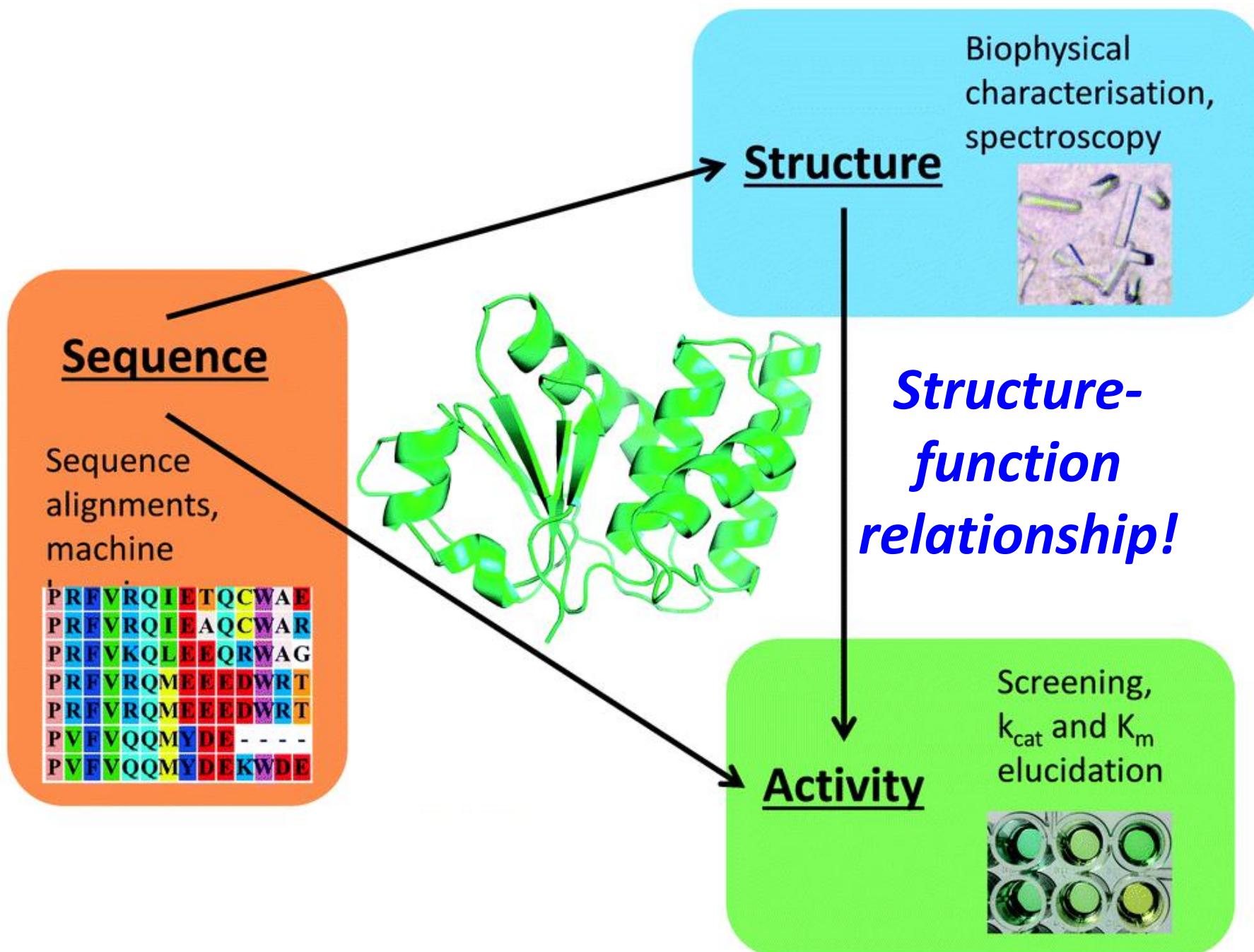


Structures:
Ground truth (green)
Predicted (blue)

STRUCTURE SOLVER

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.

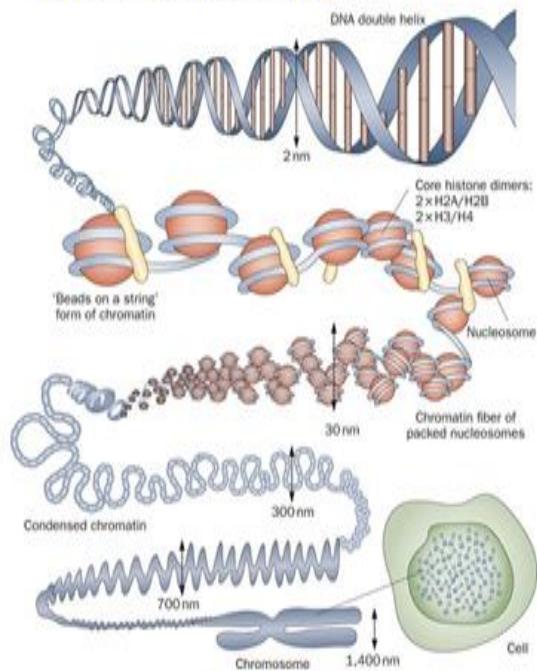




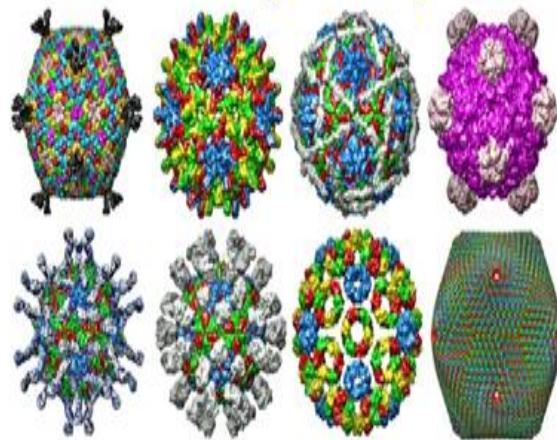
***Structure-
function
relationship!***

1. Protein and DNA structures

Chromosome structure

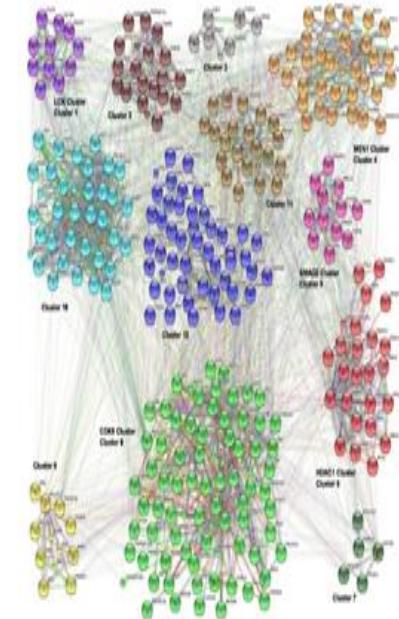


Virus self-assembly



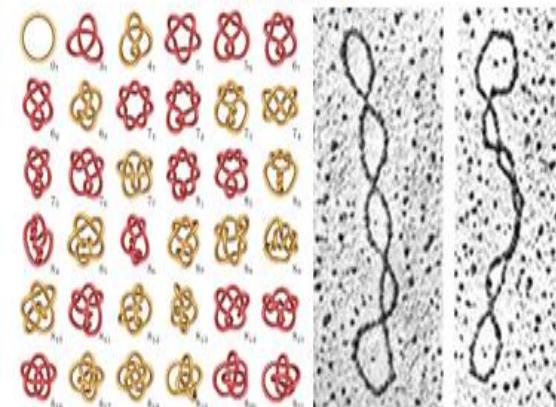
Molecular dynamics

Protein-protein Interactions

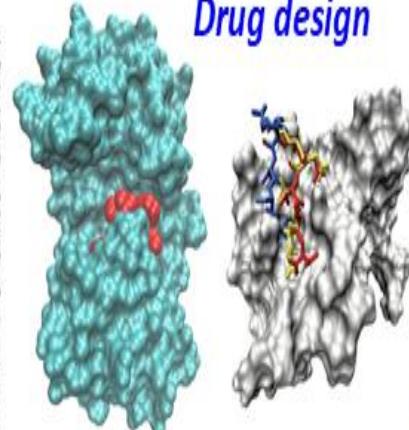


Biomolecular Topology

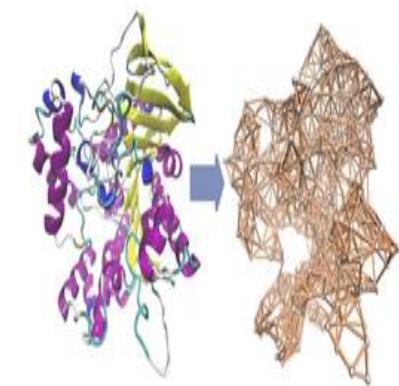
DNA knots



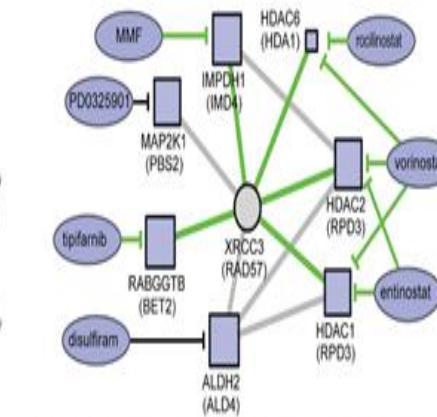
Drug design



Protein structure network

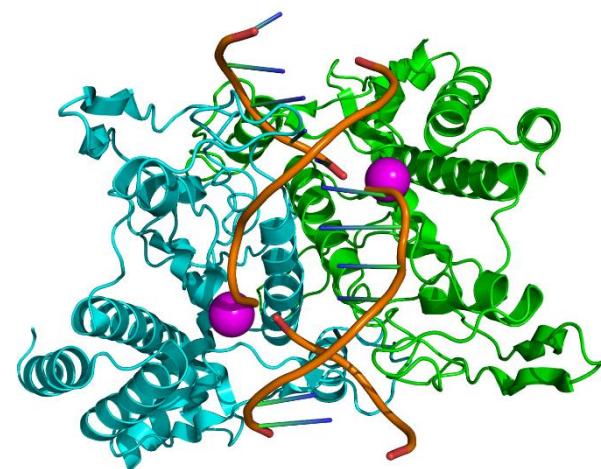
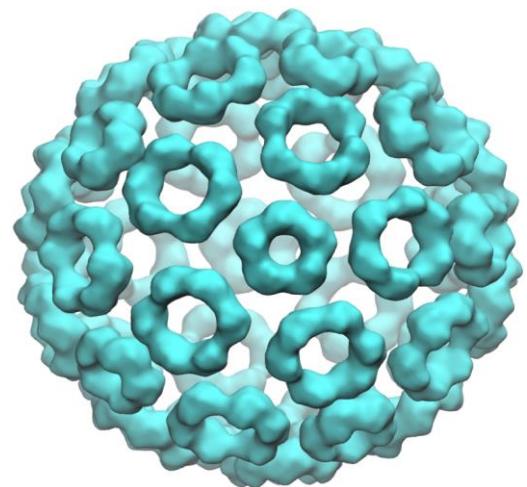
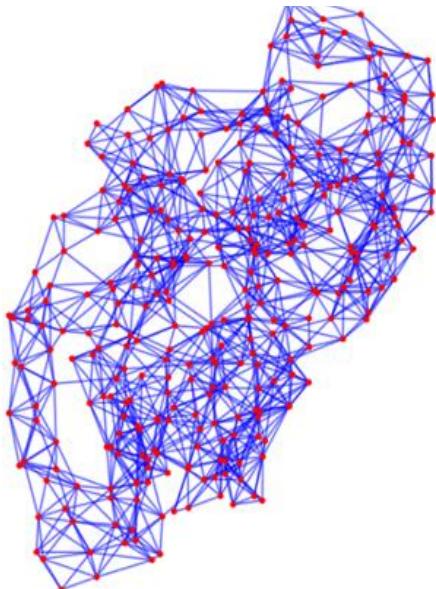
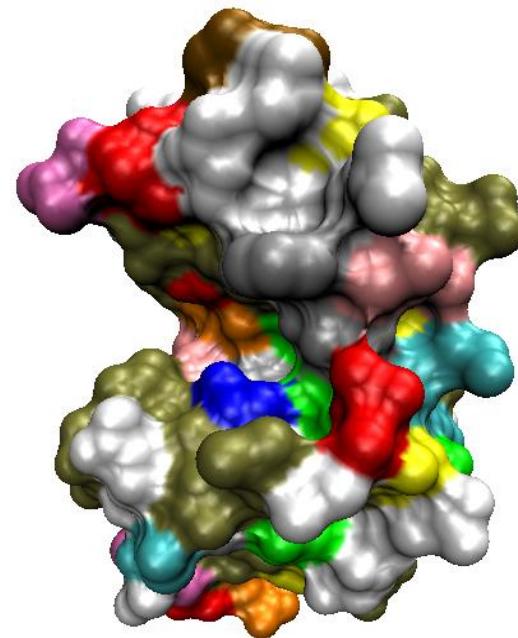
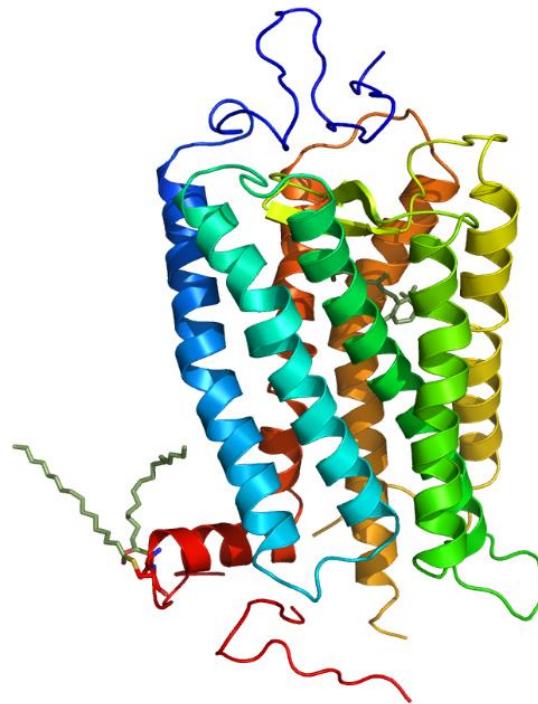
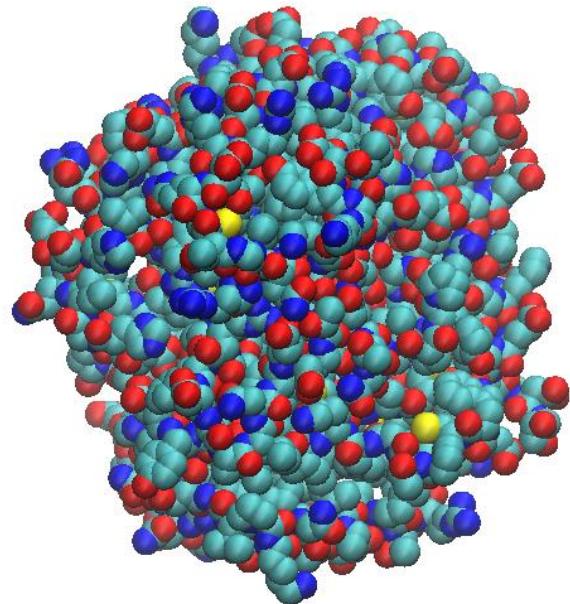


Biomolecular regulatory network

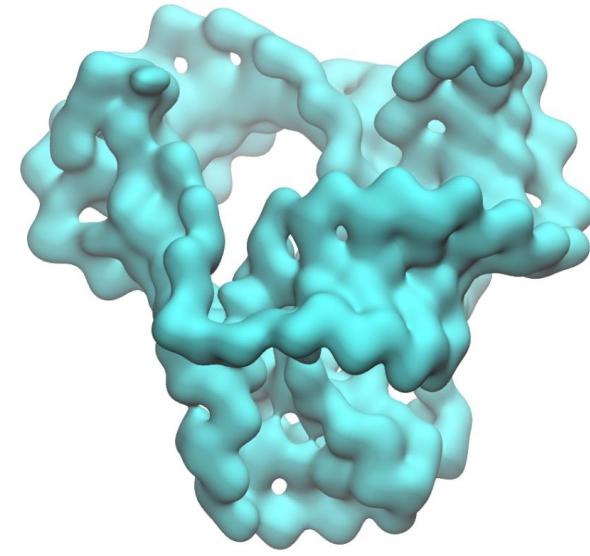
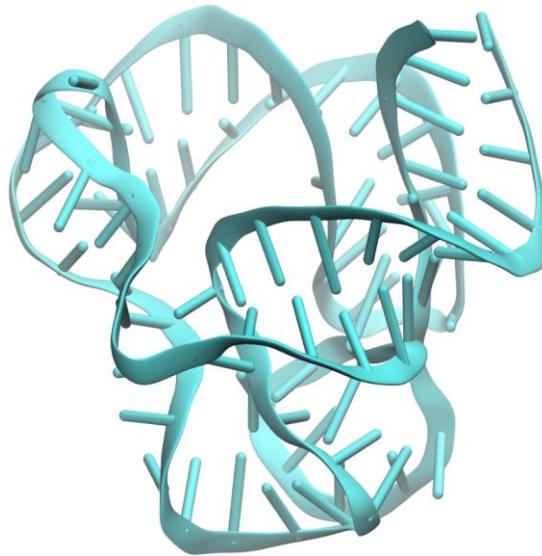
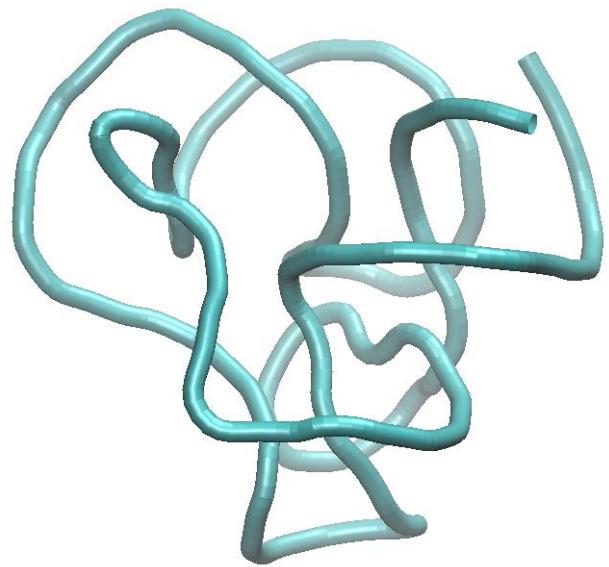
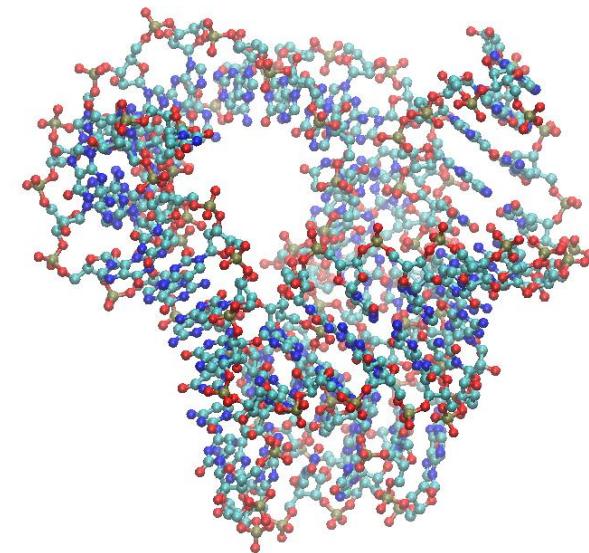
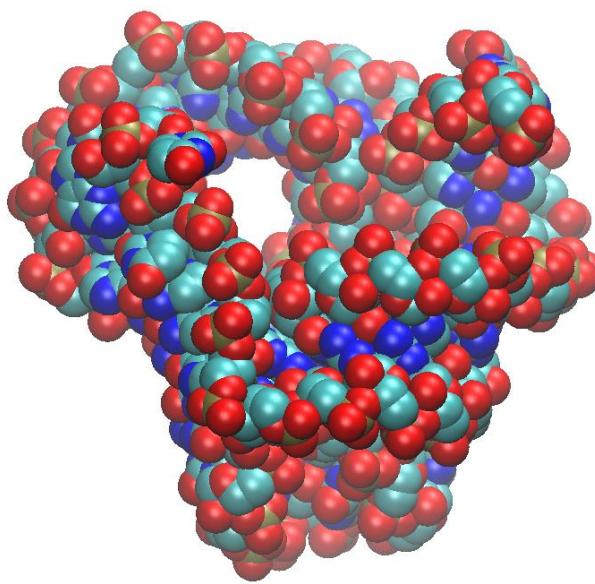
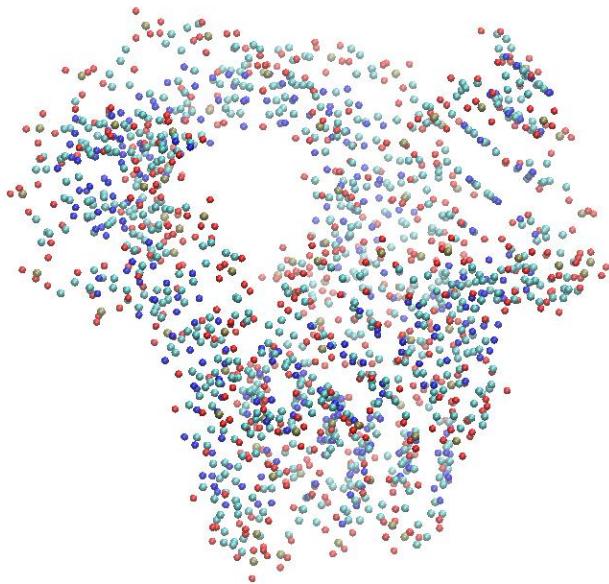


Biomolecular structure-function relationship!

All proteins?



Same biomolecule?



Protein Structure

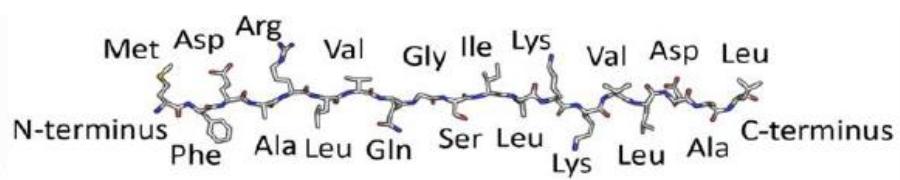
Primary structure
Secondary structure
Tertiary structure
Quaternary structure

Primary

Secondary

Tertiary

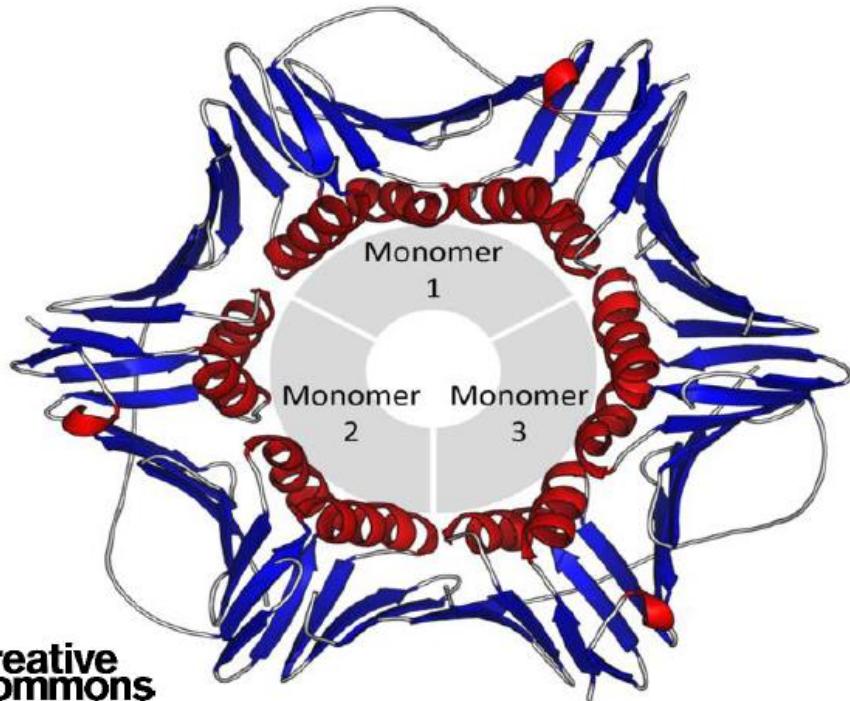
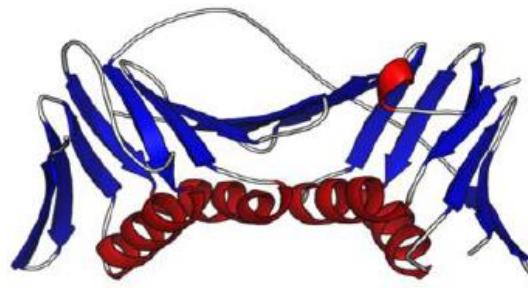
Quaternary



β -Sheet (3 strands)



α -helix

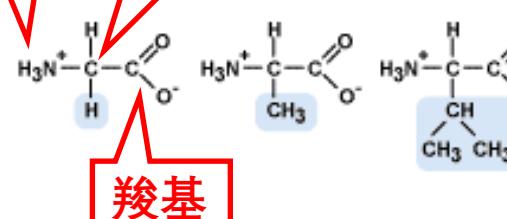


Amino Acids

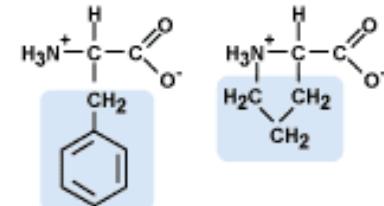
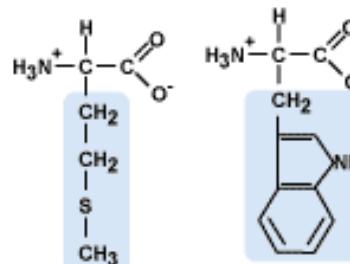
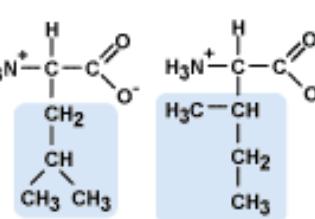
氨基

C_alpha

NONPOLAR

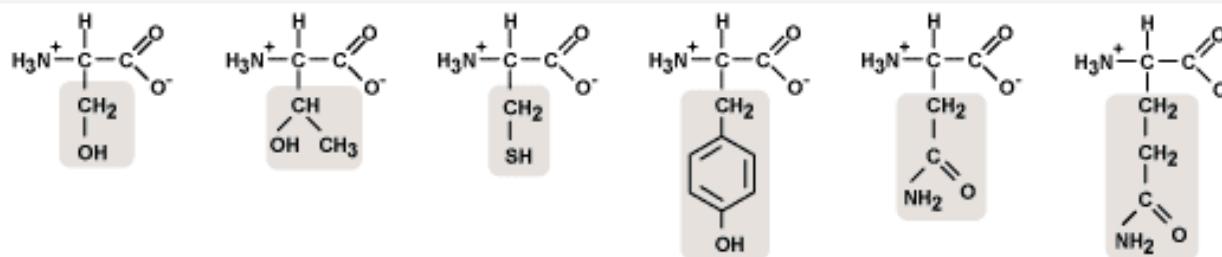


Glycine (Gly) Alanine (Ala) Valine (Val) Leucine (Leu)



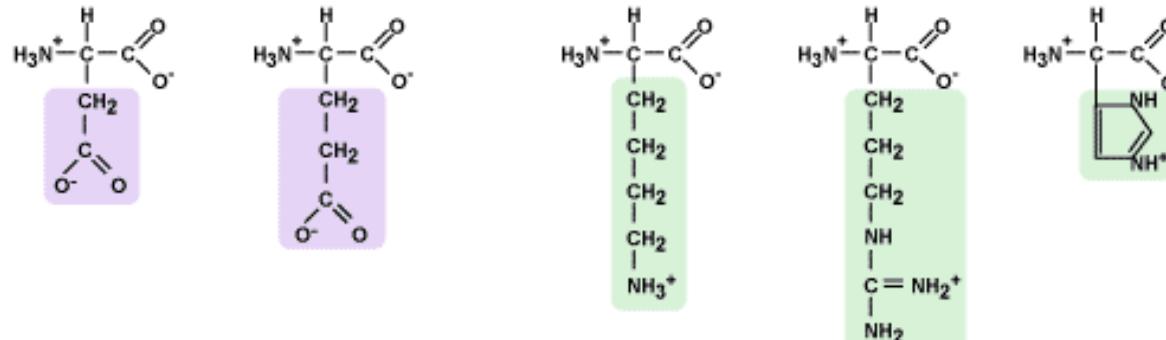
Tryptophan (Trp) Phenylalanine (Phe) Proline (Pro)

POLAR

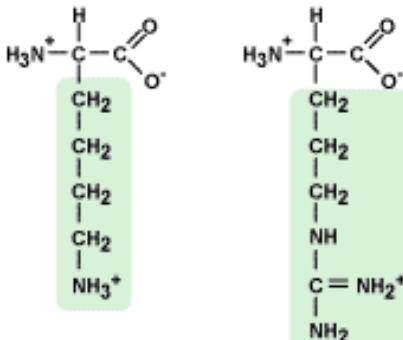


Serine (Ser) Threonine (Thr) Cysteine (Cys) Tyrosine (Tyr) Asparagine (Asn) Glutamine (Gln)

Electrically Charged



Aspartic Acid (Asp) Glutamic Acid (Glu)



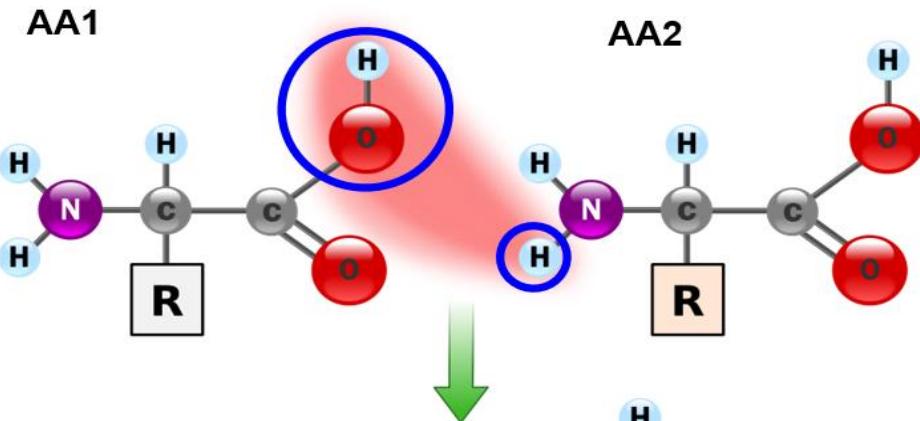
Lysine (Lys) Arginine (Arg) Histidine (His)

22 genetically encoded

Acidic

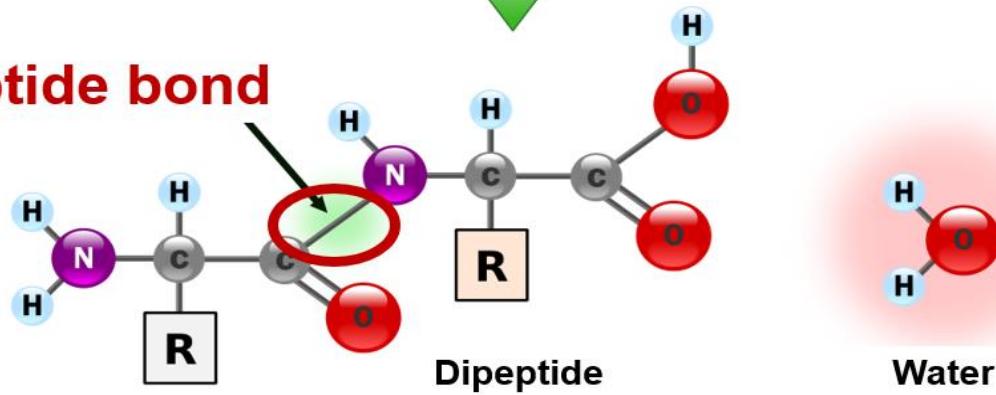
Basic

Peptide chain

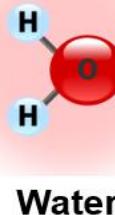


Peptide bond:
C-N ~100kcal/mol

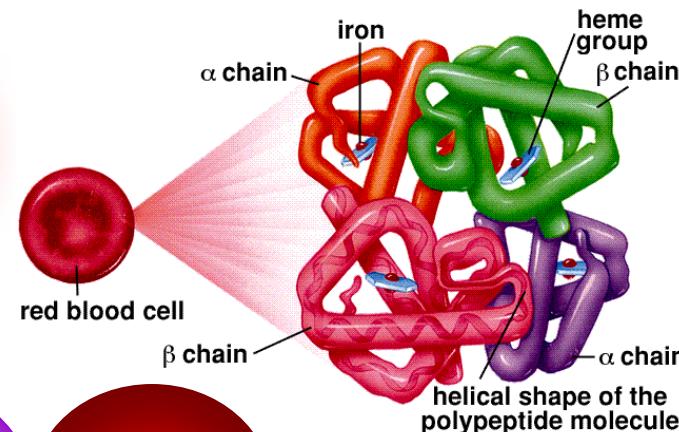
Peptide bond



Dipeptide



Water



hemoglobin

Valine

Histidine

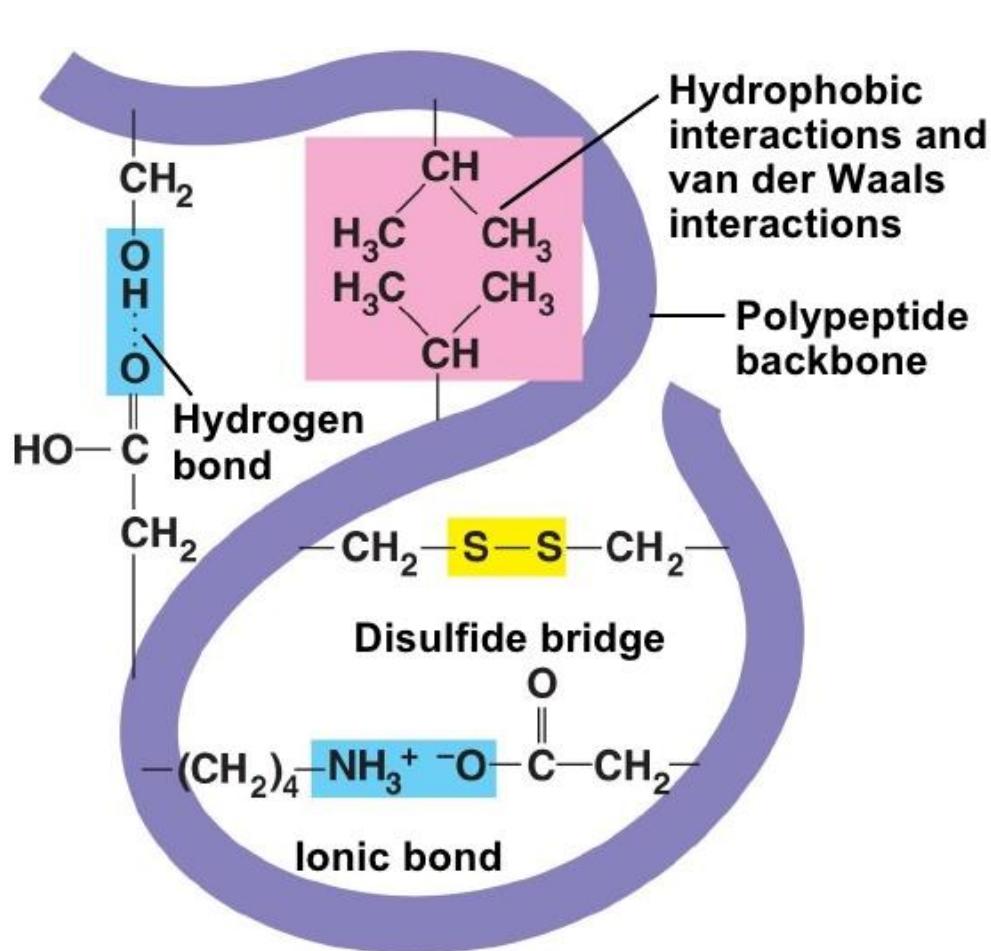
Leucine

Threonine

Proline

Glutamic Acid

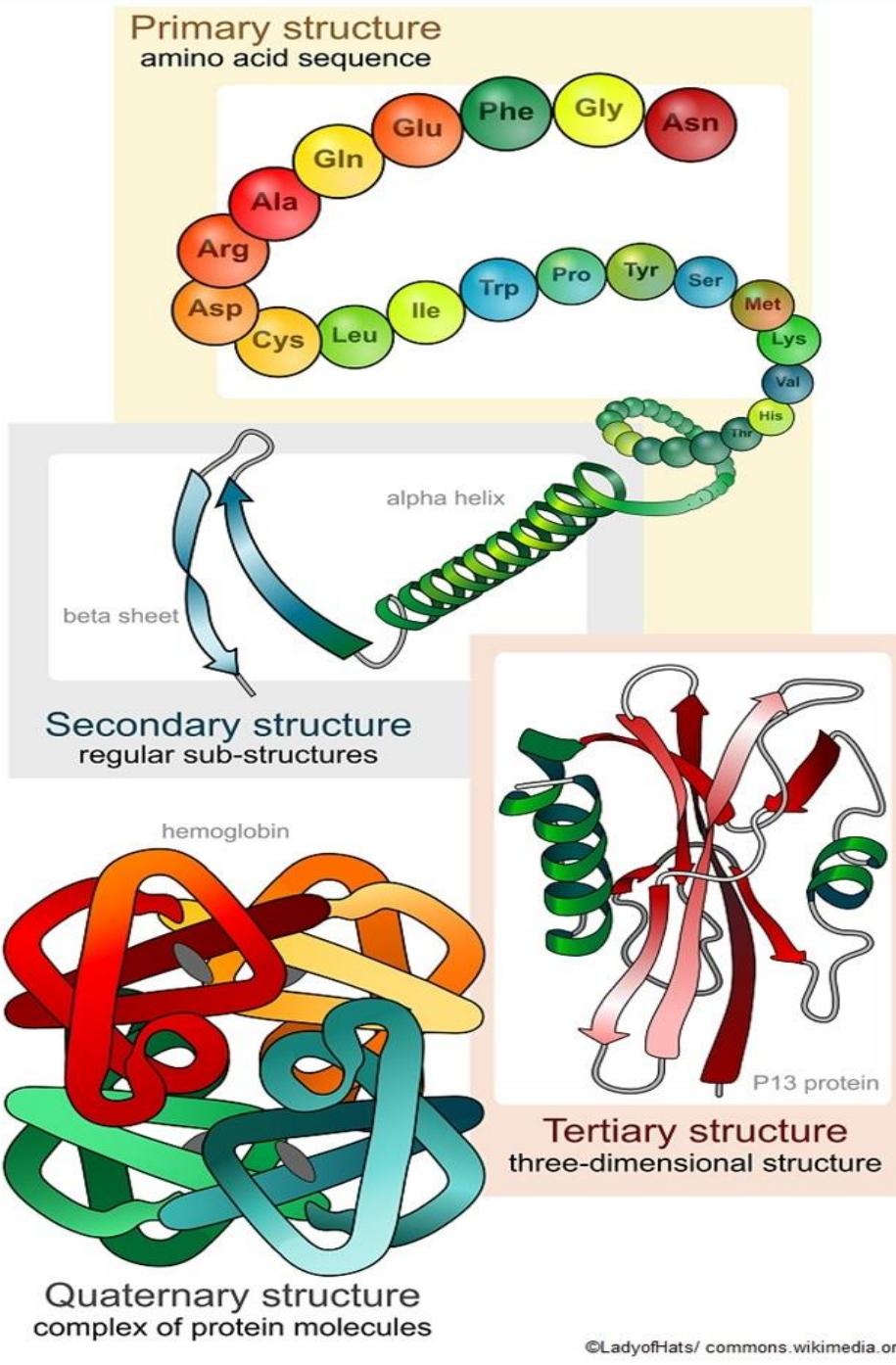
Protein bond types



- ❖ **Disulfide bond**
~60kcal/mol
- ❖ **Salt (Ionic) bond**
~20kcal/mol
- ❖ **Hydrogen bond**
~10kcal/mol
- ❖ **Hydrophobic interaction and van der Waals**
~1kcal/mol

Covalent bond: C-C ~100kcal/mol

Protein structure



Nuclei Acid Structure

Nitrogenous base (Adenine, Guanine, Cytosine, Thymine (in DNA), Uracil (in RNA))

5-carbon sugar called deoxyribose (found in DNA) and ribose (found in RNA).

One or more phosphate groups.

Primary structure
Secondary structure
Tertiary structure
Quaternary structure

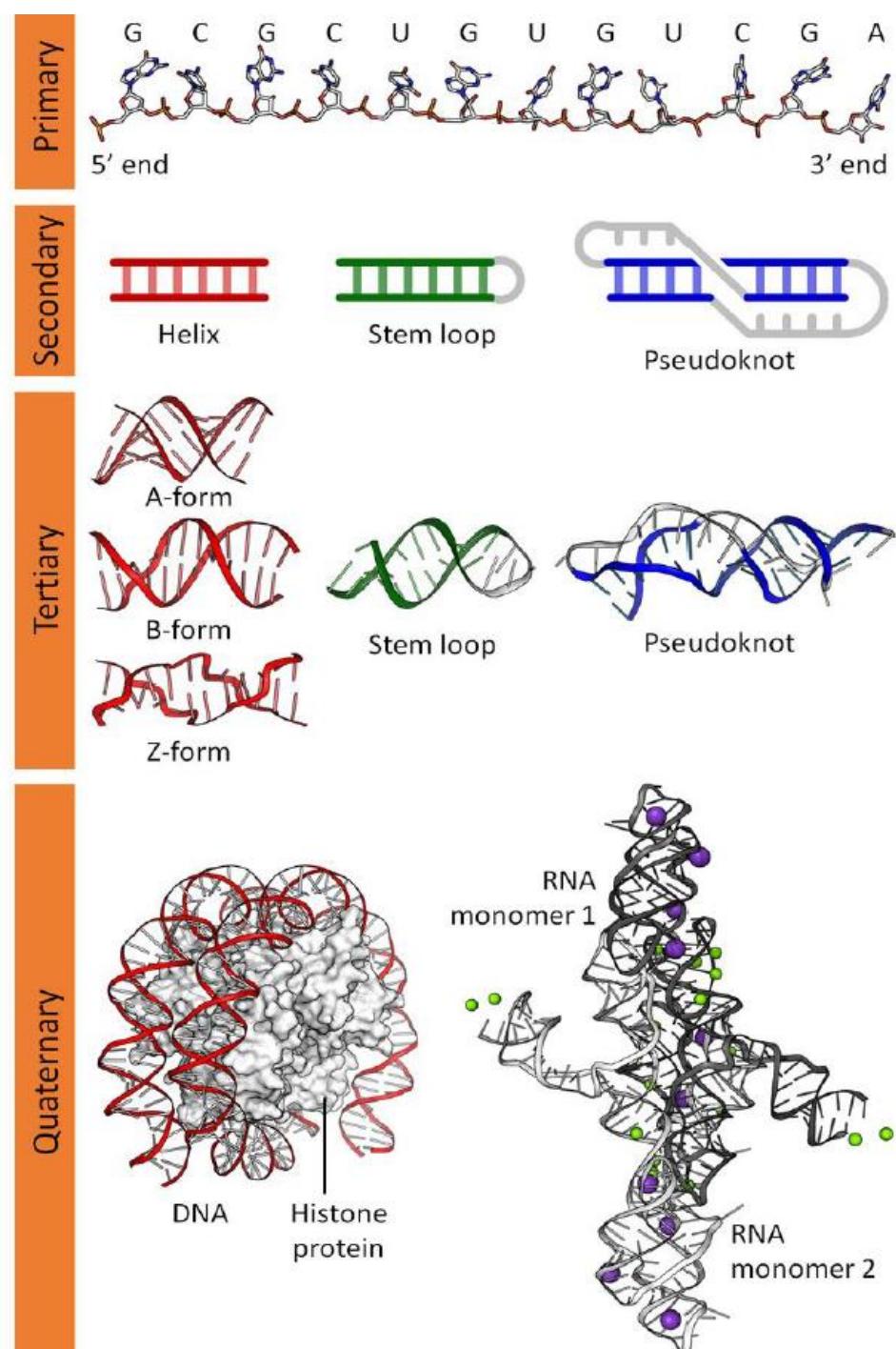
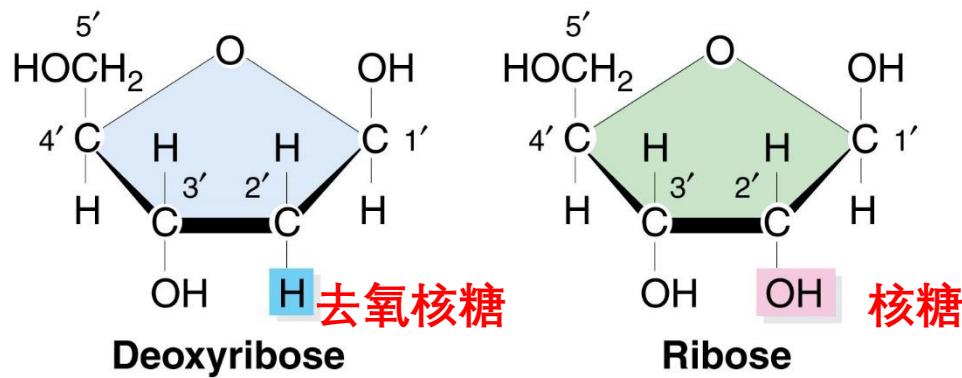
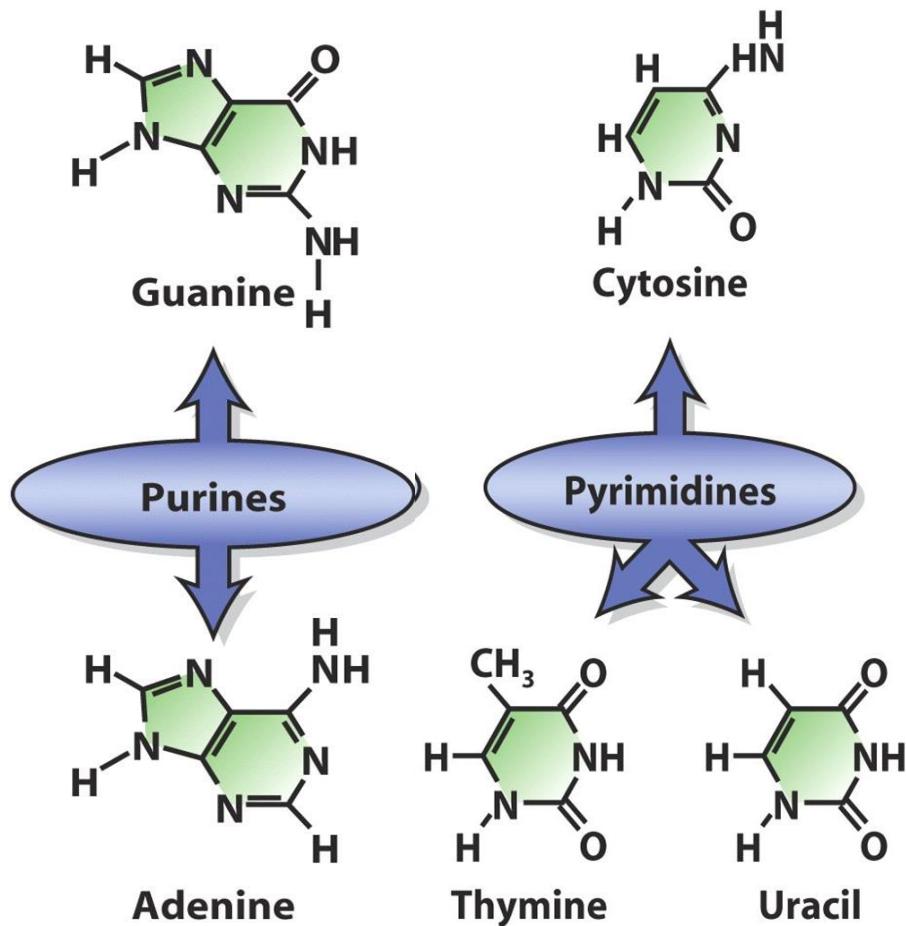
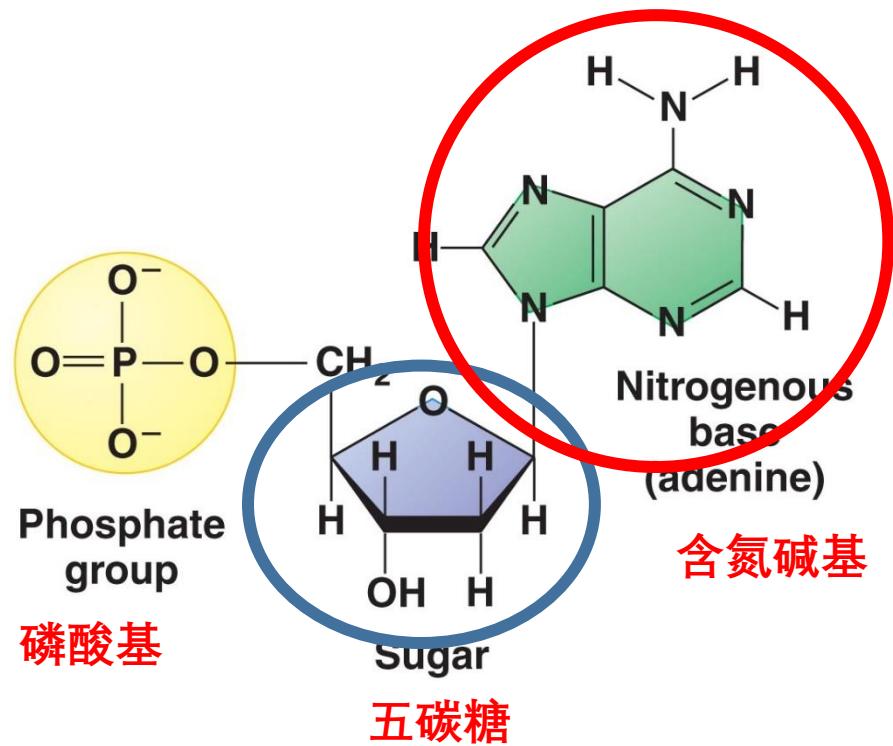
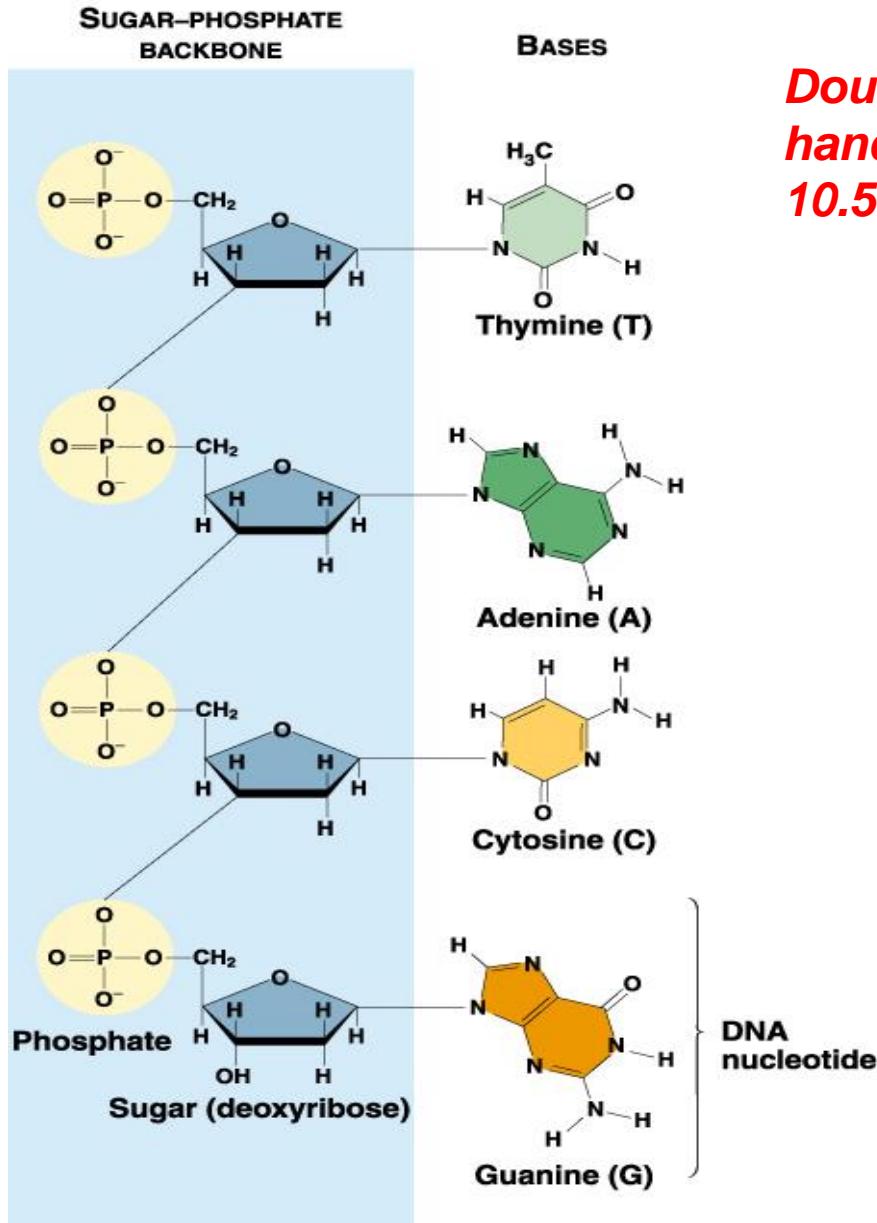


Image credit: Thomas Shafee

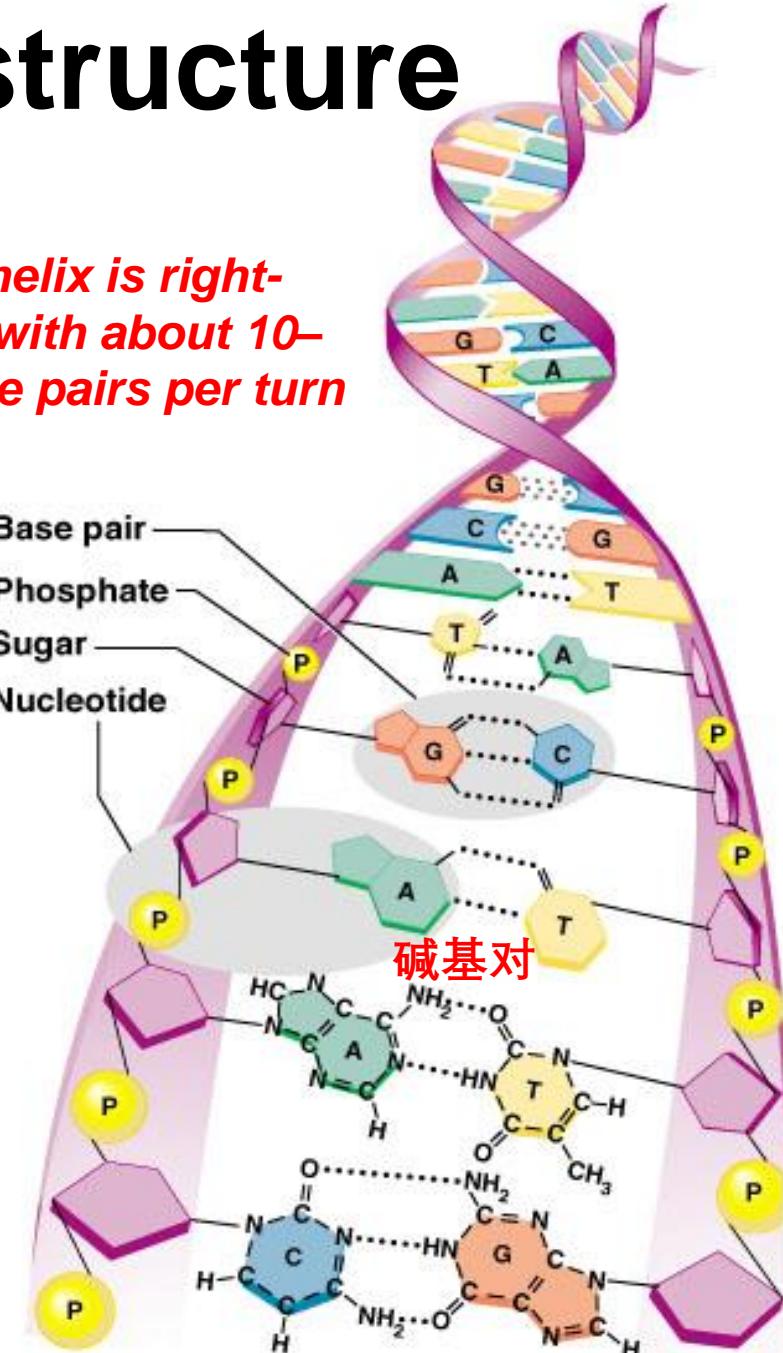
Nucleic acid structure



Nucleic acid structure



Double helix is right-handed with about 10–10.5 base pairs per turn



Experimental tools and datasets

- **Experimental tools**
 - X-Ray Crystallography
 - NMR Spectroscopy
 - *Cryogenic electron microscopy (Cryo-EM)*
- **Repositories**
 - Protein Data Bank
 - Cryo-Electron Microscopy Databank
- **Classification Data Bank**
 - CATH (Class, Architecture, Topology, Homologous superfamily)
 - SCOP (Structural Classification Of Proteins)
 - FSSP (Fold classification based on Structure-Structure alignment of Proteins)

Structure of a PDB file

2. Graph representation for biomolecules

Graph theory for molecular bioscience

- Structural stability and flexibility analysis
- Surface modeling
- Visualization
- Biomolecular domain analysis and hinge detection
- Entropy estimation
- Modeling of a wide range of biomolecular interactions
- Prediction of a wide variety of chemical and biological properties, including binding affinity, solubility, participation coefficient, mutation impact, reaction rates, toxicity, ordered-disordered transition, ...

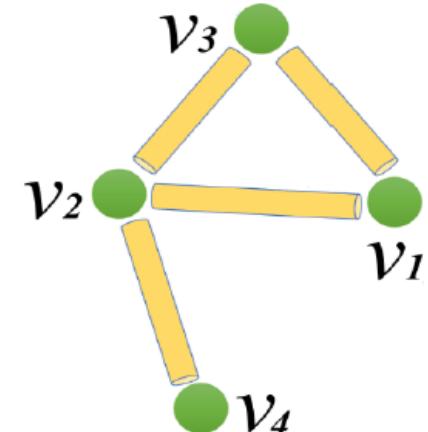
- ❖ ***Molecular descriptors***
- ❖ ***Molecular fingerprints***
- ❖ ***Graph neural network***

Algebraic graph Theory

Graph representations: Degree matrix (D),
Laplacian matrix (L) and Adjacency matrix (A).

Example: A simple (undirected) graph:

$$G = (V, E), \quad V = \{v_1, v_2, v_3, v_4\}, \\ E \{\{v_1, v_2\}, \{v_1, v_3\}, \{v_2, v_3\}, \{v_2, v_4\}\}$$



$$D = \begin{bmatrix} 2 & 0 & 0 & 0 \\ 0 & 3 & 0 & 0 \\ 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \quad A = \begin{bmatrix} 0 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 0 & 1 & 0 & 0 \end{bmatrix}$$

$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}$$

$$L = D - A$$

$$L = B^* B^T$$

B:boundary operator!

$$\sum_i d_i = 2 + 3 + 2 + 1 = 2|E| = 8$$

Degree of node i

Total # of edges

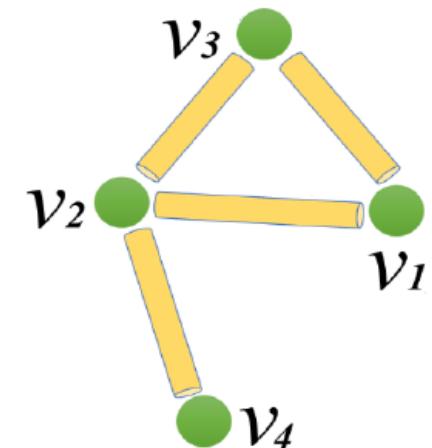
Algebraic Graph Theory

Laplacian matrix (L_G) is symmetric and has real valued entries. So it is self-adjoint and thus has real, non-negative eigenvalues:

- $0 \leq \lambda_1^L \leq \lambda_2^L \dots 0 \leq \lambda_{\text{Max}}^L$
- $\lambda_1^L = 0$ for L_G
- $\lambda_2^L > 0$ if G is connected
- Multiplicity of 0 as an eigenvalue of L_G is equal to the number of connected components of G (the topology).

Let L_G be a symmetric matrix with eigenvalues $\lambda_1^L \leq \lambda_2^L \dots 0 \leq \lambda_{\text{Max}}^L$. Then

- $\lambda_1^L = \min_{x \neq 0} \frac{x^T L_G x}{x^T x}$
- $\lambda_2^L = \min_{x \neq 0, x \perp x_1^L} \frac{x^T L_G x}{x^T x}$
- $\lambda_{\text{Max}}^L = \max_{x \neq 0} \frac{x^T L_G x}{x^T x}$



$$L = \begin{bmatrix} 2 & -1 & -1 & 0 \\ -1 & 3 & -1 & -1 \\ -1 & -1 & 2 & 0 \\ 0 & -1 & 0 & 1 \end{bmatrix}$$

Graph Theory

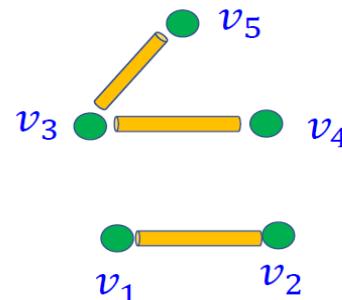
Graph partition: Partition a graph $G = (V, E)$ into smaller components with certain properties.

Fiedler eigenvalue and eigenvector: the second smallest eigenvalue (λ_2^L) provides a lower bound on ratio-cut partition: $c \geq \frac{\lambda_2^L}{|E|}$. The associated eigenvector, Fiedler vector bisects the graph into two sections based on the sign of the eigenvector (i.e., spectral bisection based on algebraic connectivity).

Example I:

$$\text{Vec} = \begin{bmatrix} -\frac{1}{\sqrt{2}} & 0 & 0 & -\frac{1}{\sqrt{2}} & 0 \\ -\frac{1}{\sqrt{2}} & 0 & 0 & \frac{1}{\sqrt{2}} & 0 \\ 0 & -\frac{1}{\sqrt{3}} & 0 & 0 & -\frac{2}{\sqrt{6}} \\ 0 & -\frac{1}{\sqrt{3}} & -\frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{6}} \\ 0 & -\frac{1}{\sqrt{3}} & \frac{1}{\sqrt{2}} & 0 & \frac{1}{\sqrt{6}} \end{bmatrix} \quad L = \begin{bmatrix} 1 & -1 & 0 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 2 & -1 & -1 \\ 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 0 & 1 \end{bmatrix}$$

Eig = [0, 0, 1, 2, 3]

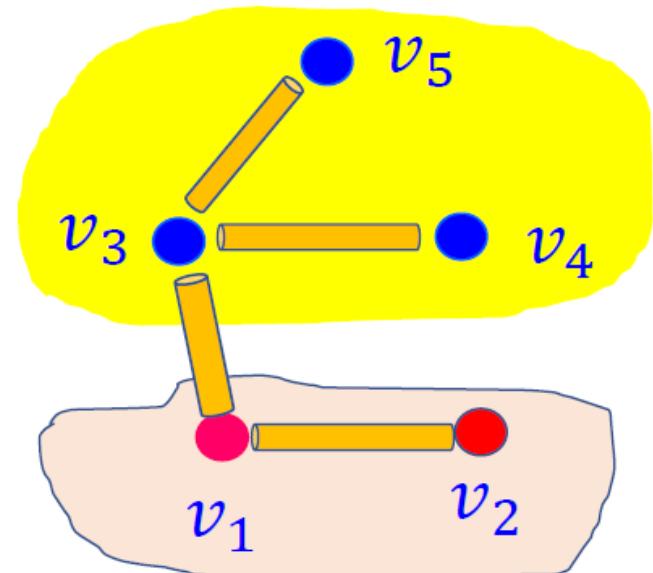


Two disconnected components (harmonic part due to the topology)

Spectral bisection

Example II:

$$L = \begin{bmatrix} 2 & -1 & -1 & 0 & 0 \\ -1 & 1 & 0 & 0 & 0 \\ -1 & 0 & 3 & -1 & -1 \\ 0 & 0 & -1 & 1 & 0 \\ 0 & 0 & -1 & 0 & 1 \end{bmatrix}$$



$$\text{Vec} = \begin{bmatrix} 0.45 & 0.34 \\ 0.45 & 0.70 \\ 0.45 & -0.20 \\ 0.45 & -0.42 \\ 0.45 & -0.42 \end{bmatrix} \quad \begin{bmatrix} 0 & -0.70 & 0.44 \\ 0 & 0.54 & -0.14 \\ 0 & -0.32 & -0.81 \\ -0.70 & 0.24 & 0.26 \\ 0.70 & 0.24 & 0.26 \end{bmatrix}$$

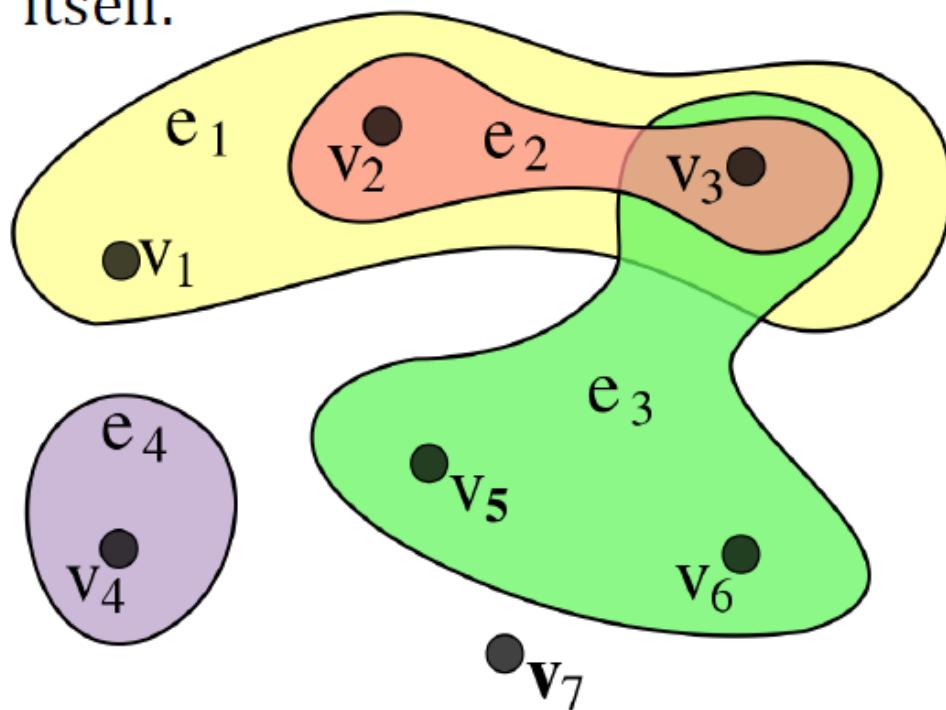
$$\text{Eig} = [0, 0.52, 1.00, 2.31, 4.17]$$

Graph Theory

Hypergraph: A hypergraph (H) is a generalization of a graph in which an edge can join any number of vertices.

$H = (X, E)$, where X is a set of nodes or vertices and E is a set of hyperedges, which are a subset of the power set such that .
 $E \subseteq \wp(S) \setminus \{\emptyset\}$.

A power set $\wp(S)$ of set S is the set of all subsets of S , including the empty set $\{\emptyset\}$ and S itself.



Gaussian Network Model (GNM)—Laplacian model

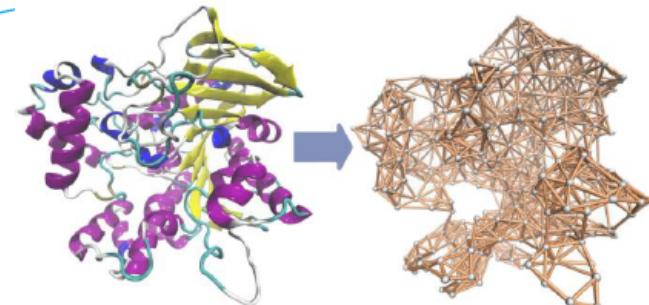
(Bahar, Atilgan and Erman, FD, 1997)

$$B_j^{\text{GNM}} = \alpha (L^{-1})_{jj},$$

$$L_{ij} = \begin{cases} -1, & i \neq j, r_{ij} \leq r_c \\ 0, & i \neq j, r_{ij} > r_c \\ -\sum_{j,j \neq i} L_{ij}, & i = j \end{cases}$$

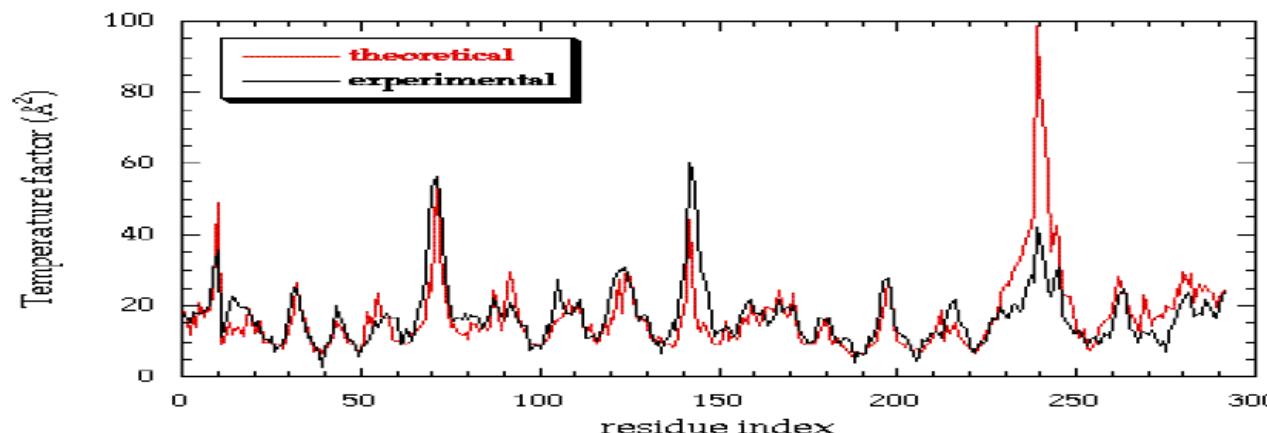
$$(L^{-1})_{jj} = \sum_{k=2}^N \frac{1}{\lambda_k} [u_k u_k^T]_{jj}$$

L is a Laplacian matrix, $O(N^3)$ method



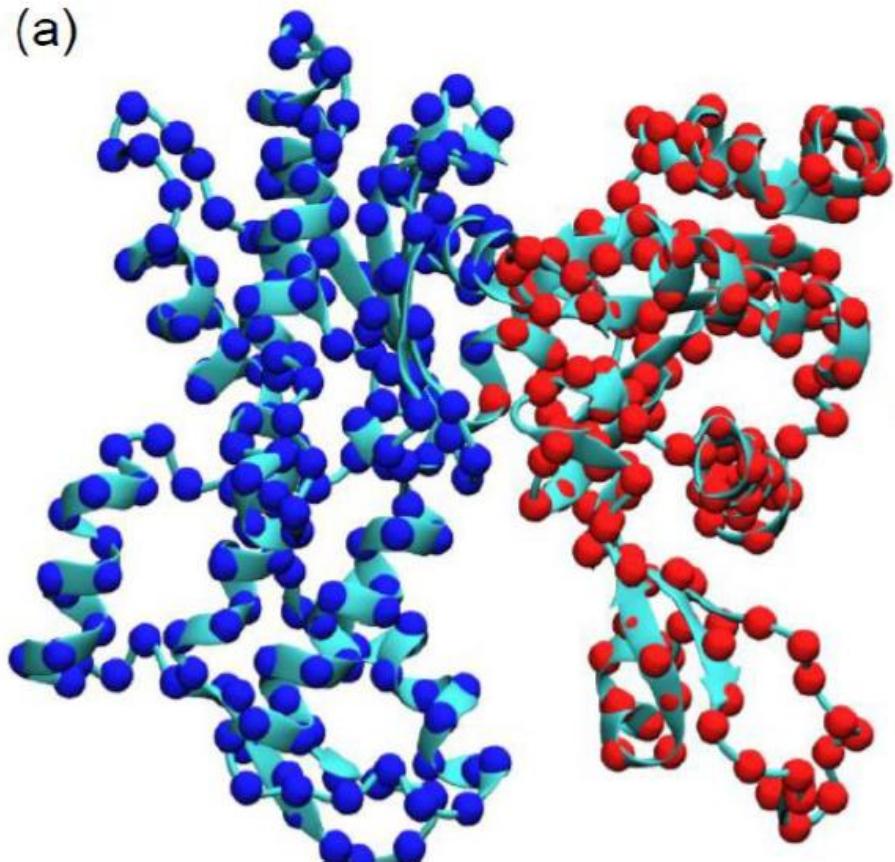
Moore–Penrose pseudoinverse

where α is fitting parameter, r_c a cutoff distance (7 Å is often used for C_α networks), λ_k the kth eigenvalue and u_k the kth eigenvector.

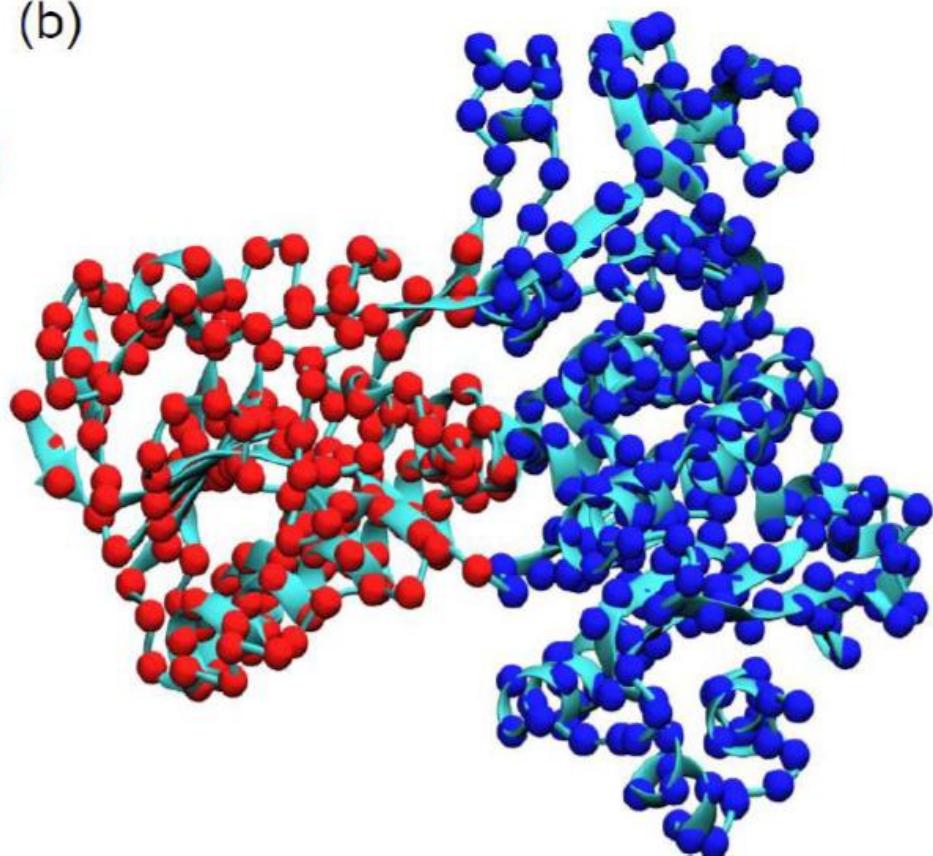


Multiscale GNM based domain analysis

(a)



(b)



Protein domain decomposition with Type-1 mGNM. The first non-zero eigenvector (Fiedler vector) is used to decompose the protein into two domains. (a) Protein 1ATN (chain A); (b) protein 3GRS.

Anisotropic network model

Potential function:

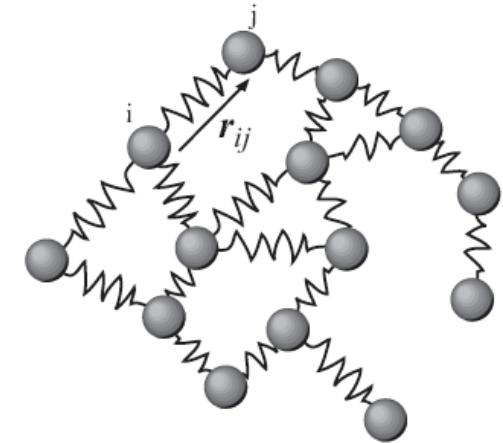
$$r_{ij}^d = |\mathbf{r}_{ij}^d| \quad r_{ij} = |\mathbf{r}_{ij}|$$

$$\Delta \mathbf{R} = \{\Delta x_1, \Delta y_1, \Delta z_1, \dots, \Delta x_N, \Delta y_N, \Delta z_N\}$$

$$V^{\text{ANM}} = \gamma \sum_{i,j}^N (r_{ij}^d - r_{ij})^2 f(r_{ij}) = \frac{\gamma}{2} \Delta \mathbf{R}^T H \Delta \mathbf{R}.$$

$$H_{ij} = -\frac{1}{r_{ij}^2} \begin{bmatrix} (x_j - x_i)(x_j - x_i) & (x_j - x_i)(y_j - y_i) & (x_j - x_i)(z_j - z_i) \\ (y_j - y_i)(x_j - x_i) & (y_j - y_i)(y_j - y_i) & (y_j - y_i)(z_j - z_i) \\ (z_j - z_i)(x_j - x_i) & (z_j - z_i)(y_j - y_i) & (z_j - z_i)(z_j - z_i) \end{bmatrix} \quad i, j = 1, 2, \dots, N, i \neq j \text{ and } r_{ij} \leq r_c.$$

$$H_{ii} = - \sum_{i \neq j} H_{ij}, \quad \forall i = 1, 2, \dots, N.$$



eigenvector

eigenvalue

Moore-Penrose pseudoinverse:

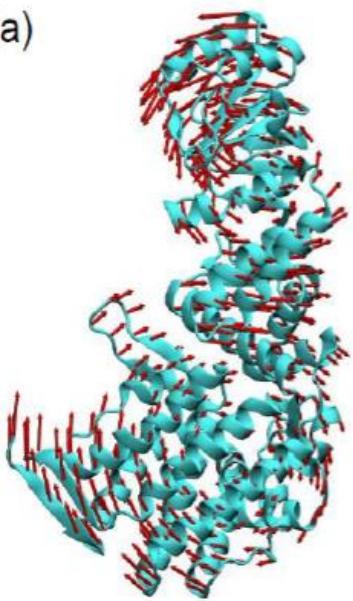
$$(H^{-1})_{ii} = \sum_{k=1}^{3N} \lambda_k^{-1} [\mathbf{v}_k \mathbf{v}_k^T]_{ii}$$

Predicted b-factor:

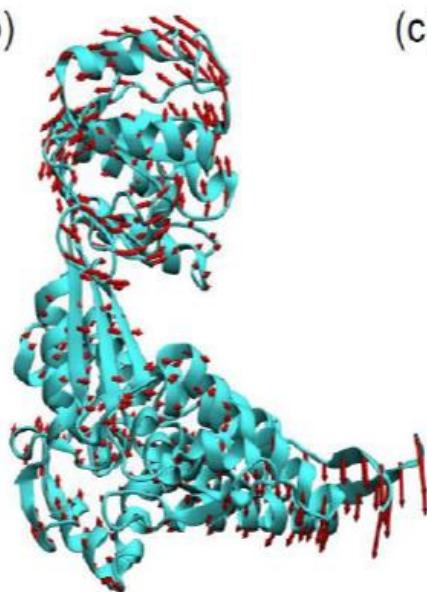
$$B_i^{\text{ANM}} = \frac{8\pi^2}{3} \sum_{j=3i-2}^{3i} \langle \Delta \mathbf{R}_j \cdot \Delta \mathbf{R}_j \rangle, \quad \forall i = 1, 2, \dots, N.$$

$$\langle \Delta \mathbf{R}_i \cdot \Delta \mathbf{R}_j \rangle = \frac{3k_B T}{\gamma} (H^{-1})_{ij}, \quad \forall i, j = 1, 2, \dots, 3N.$$

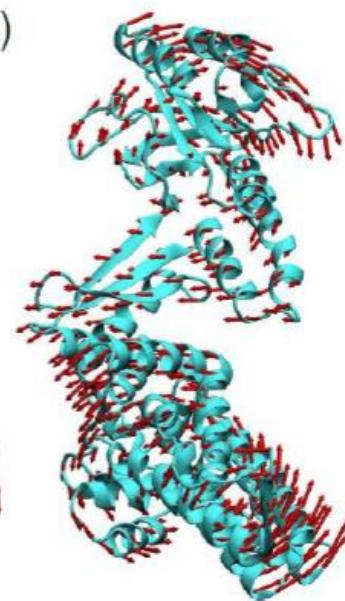
(a)



(b)

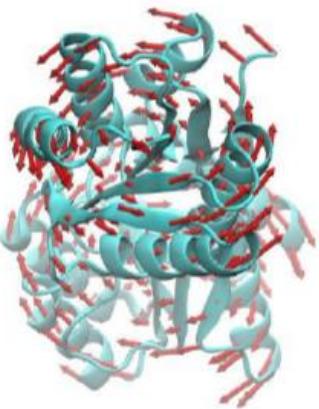


(c)

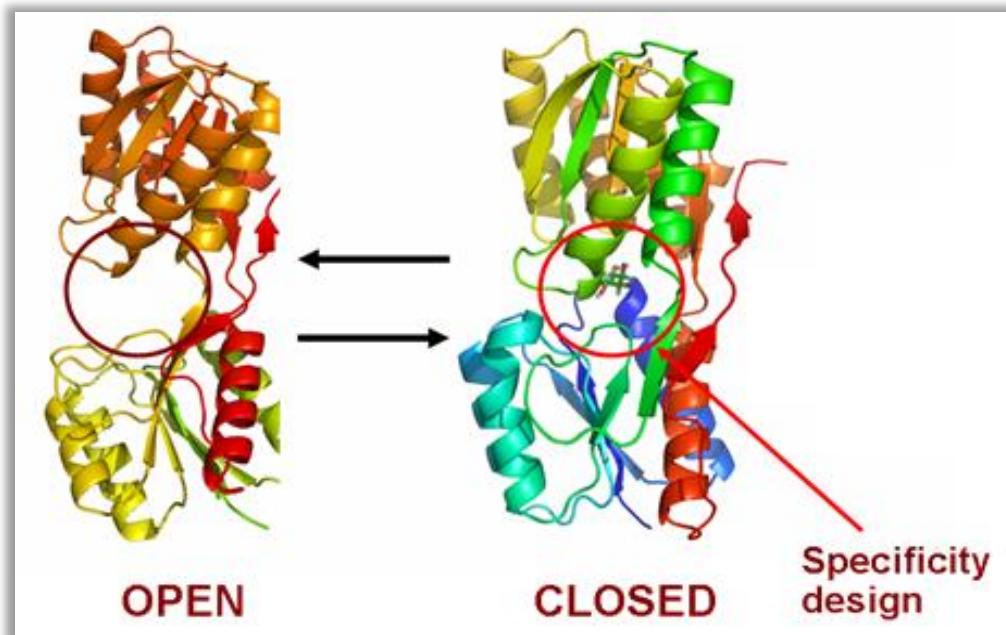
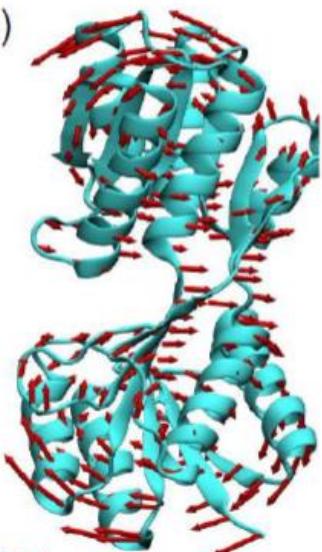


The motions of 1GRU (chain A). The 7th, 8th, and 9th nANM modes are demonstrated in (a)-(c), respectively

(a)

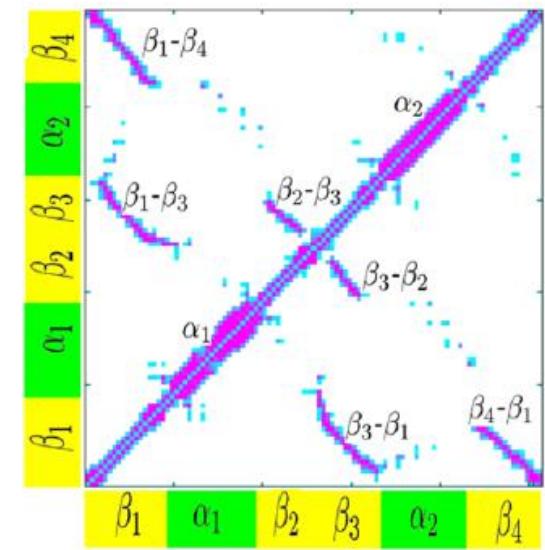
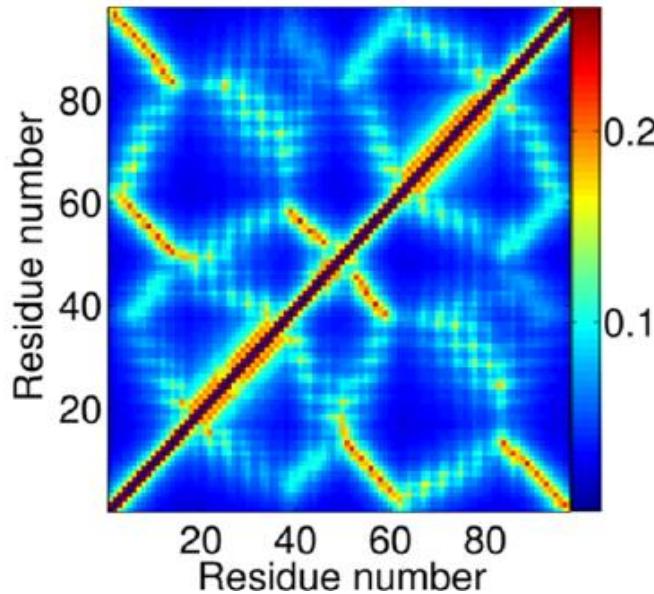
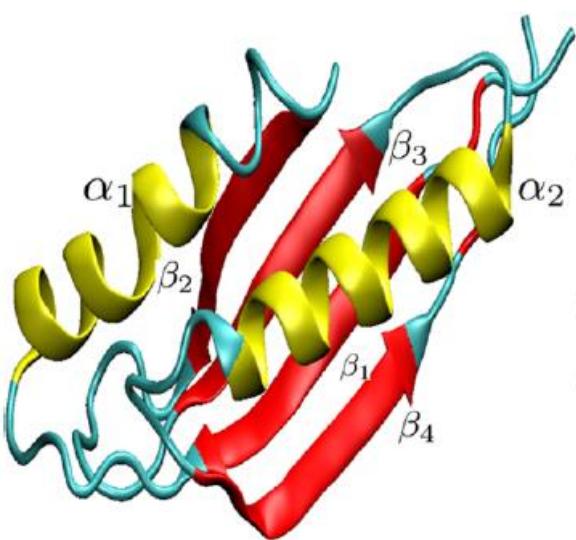


(b)



(Xia, Opron and Wei, JCP, 2016)

Flexibility rigidity index



Connectivity matrix:

$$A_{ij} = \begin{cases} \phi(\| r_i - r_j \|; \eta), & i \neq j; \\ -\sum_{i \neq j} A_{ij}, & i = j. \end{cases}$$

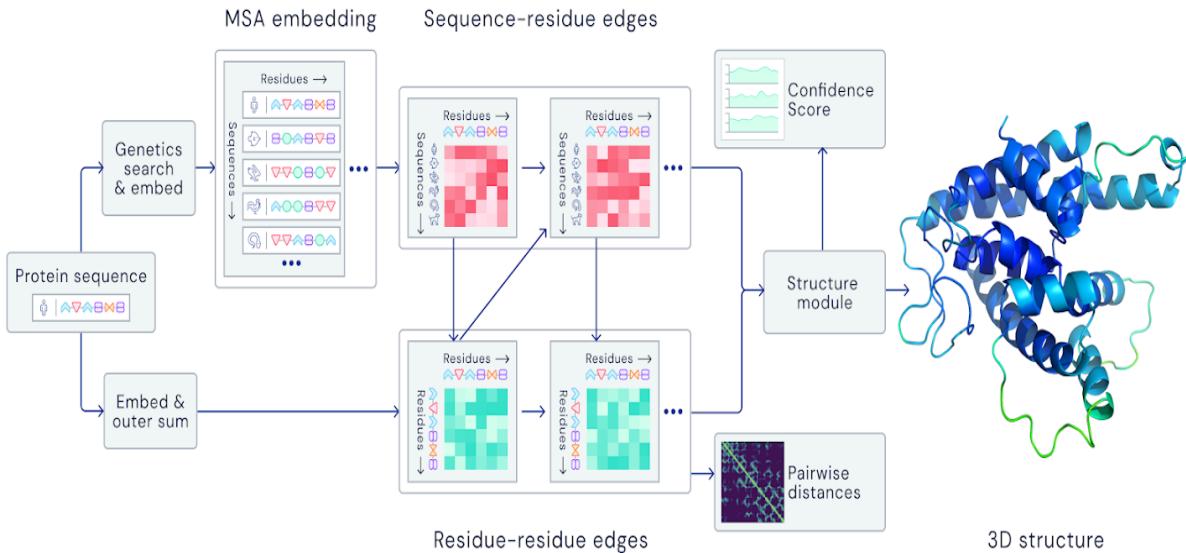
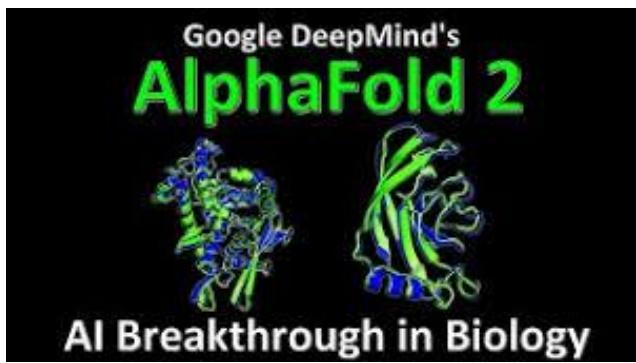
Kernel function:

$$\phi(\| r_i - r_j \|; \eta) = 1, \text{ as } \| r_i - r_j \| \rightarrow 0$$

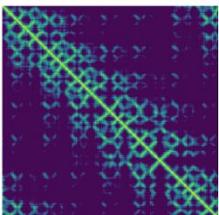
$$\phi(\| r_i - r_j \|; \eta) = 0, \text{ as } \| r_i - r_j \| \rightarrow \infty$$

Generalized exponential: $\phi(\| r_i - r_j \|; \eta) = e^{-(\| r_i - r_j \| / \eta)^\kappa}$

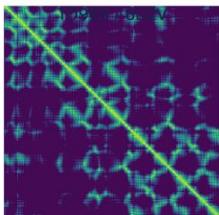
Generalized Lorentz: $\phi(\| r_i - r_j \|; \eta) = \frac{1}{1 + (\| r_i - r_j \| / \eta)^\nu}$



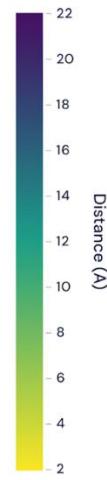
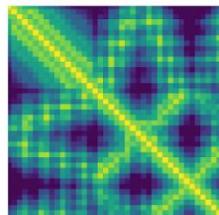
T0954 / 6CVZ



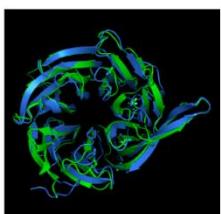
T0965 / 6D2V



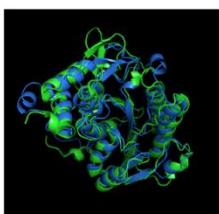
T0955 / 5W9F



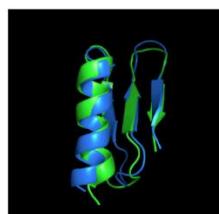
T0954 / 6CVZ



T0965 / 6D2V



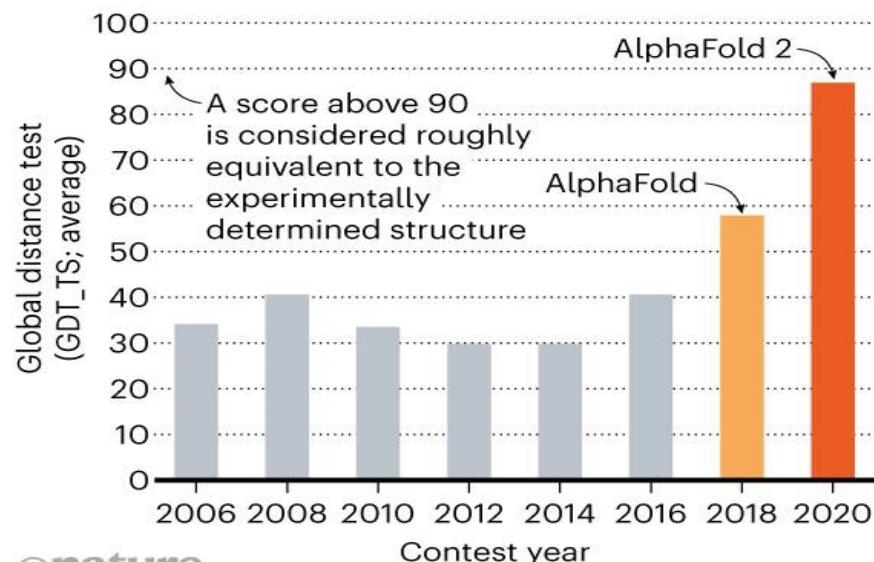
T0955 / 5W9F



Structures:
Ground truth (green)
Predicted (blue)

STRUCTURE SOLVER

DeepMind's AlphaFold 2 algorithm significantly outperformed other teams at the CASP14 protein-folding contest — and its previous version's performance at the last CASP.



Flexibility rigidity index (FRI)

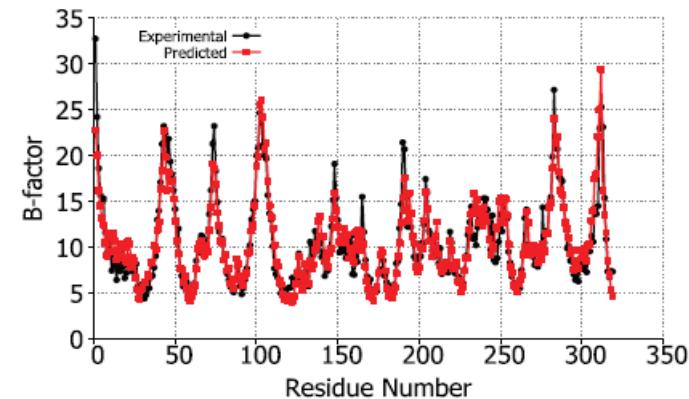
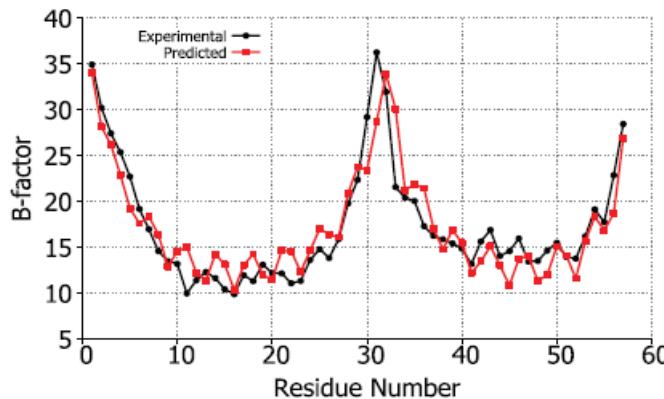
Rigidity index:

$$\mu_i = \sum_{j \neq i} w_j \phi(\| r_i - r_j \|; \eta)$$

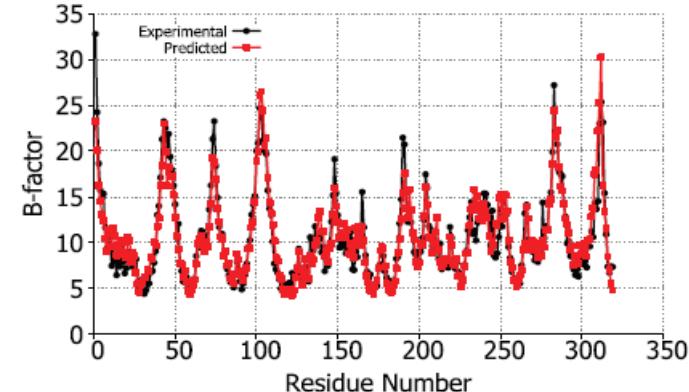
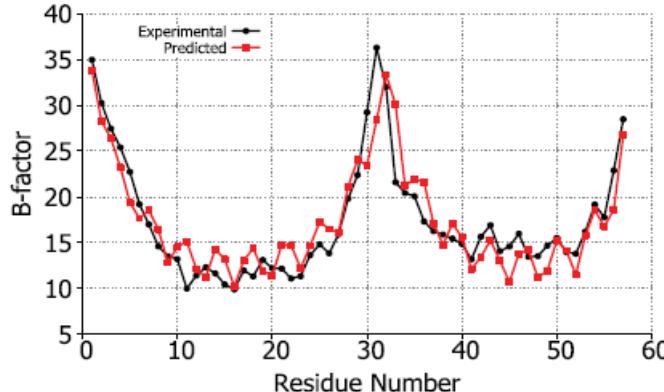
Flexibility index:

$$f_i = \frac{1}{\mu_i} = \frac{1}{\sum_{j \neq i} w_j \phi(\| r_i - r_j \|; \eta)}$$

**Upper parts:
Lorentz kernel**



**Lower parts:
exponential kernel**



Protein ID: 1DF4

Protein ID: 2Y7L

Parameter testing

Correlation
function

Parameter
range

Average correlation
coefficient

$$e^{-(r/\eta)^\kappa}$$

$$1.0 \leq \eta \leq 10.0$$

$$0.5 \leq \kappa \leq 10.0$$

$$0.676$$

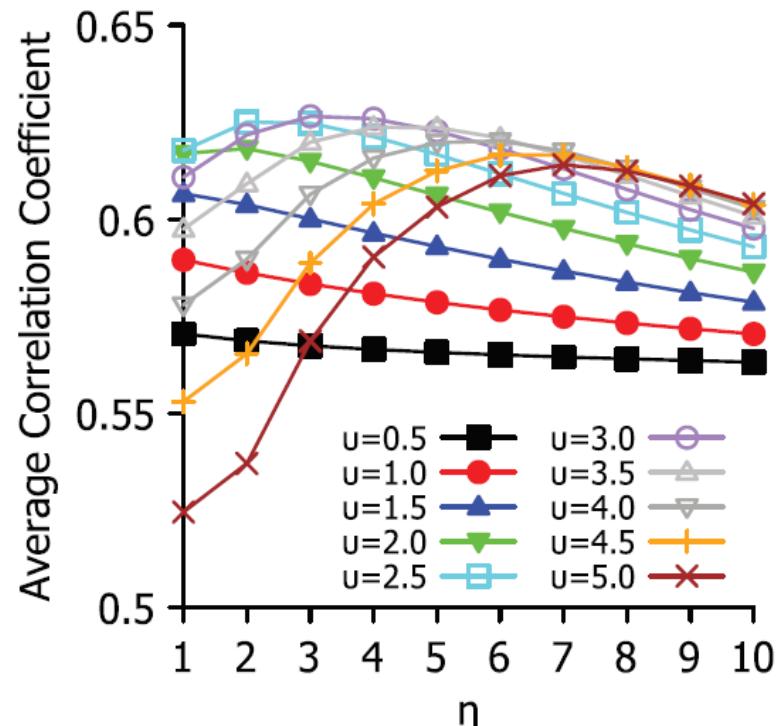
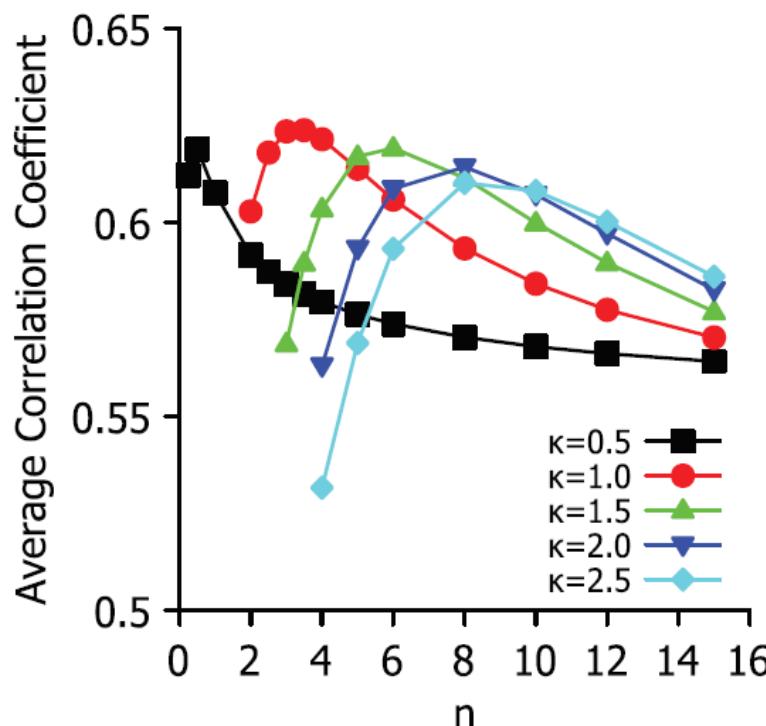
$$\frac{1}{1+(r/\eta)^\nu}$$

$$1.0 \leq \eta \leq 10.0$$

$$0.5 \leq \nu \leq 10.0$$

$$0.673$$

$$C_c = \frac{\sum_{i=1}^N (B_i^e - \bar{B}^e)(B_i^t - \bar{B}^t)}{\left[\sum_{i=1}^N (B_i^e - \bar{B}^e)^2 \sum_{i=1}^N (B_i^t - \bar{B}^t)^2 \right]^{1/2}}$$



Parameter testing for exponential (left chart) and Lorentz (right chart) using the dataset with 365 proteins.

Performance of our FRI

Accuracy: (10% improvement)

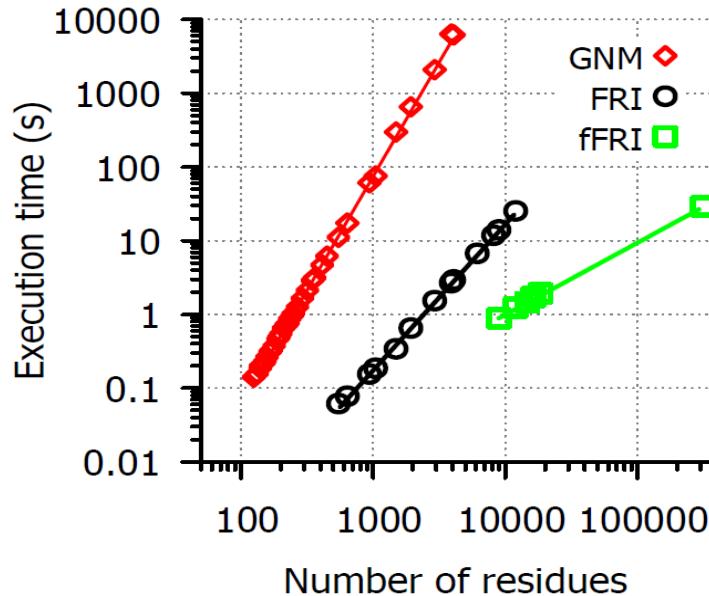
.. atomic mean-square
displacements is essentially
determined by spatial variations in
local packing density.."

Bertil Halle, PNAS, Vol. 99, No.3, 1274-1279, 2002

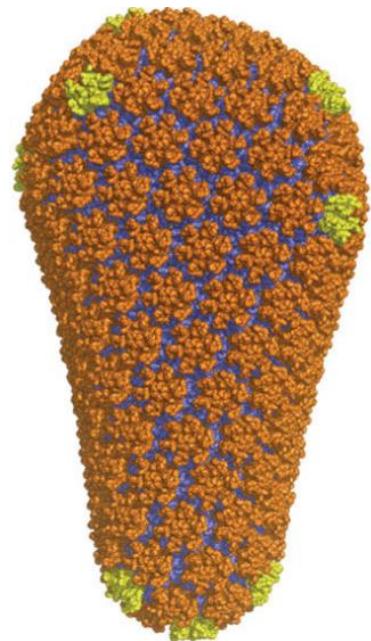
PDB set	pfFRI	GNM	NMA
Small	0.594	0.541	0.480
Medium	0.605	0.550	0.482
Large	0.591	0.529	0.494
Superset	0.626	0.565	NA

Exponential parameters	Avg. CC	Lorentz parameters	Avg. CC
$\kappa=0.5, \eta=0.5$	0.615 (8.8%)	$v=2.5, \eta=2.0$	0.622 (10.1%)
$\kappa=1.0, \eta=3.0$	0.623 (10.3%)	$v=3.0, \eta=3.0$	0.626 (10.8%)
$\kappa=1.5, \eta=6.0$	0.619 (9.6%)	$v=3.5, \eta=4.0$	0.623 (10.3%)

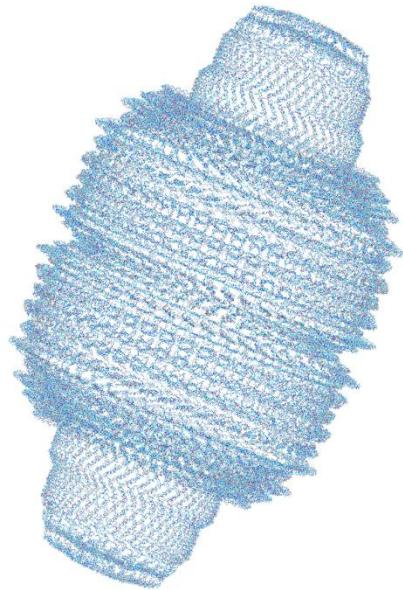
Time: fFRI $O(N)$



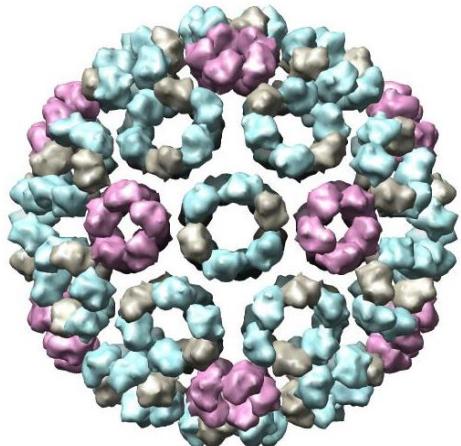
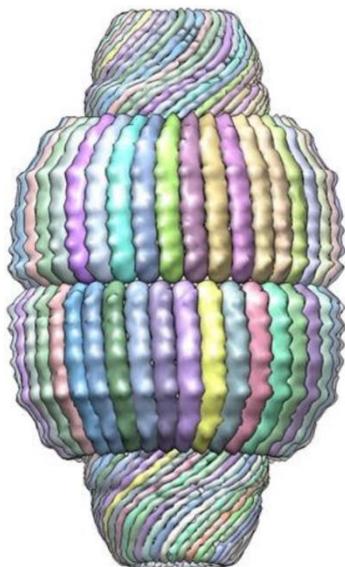
HIV Virus capsid
(313 236
residues) in less
than 30 seconds



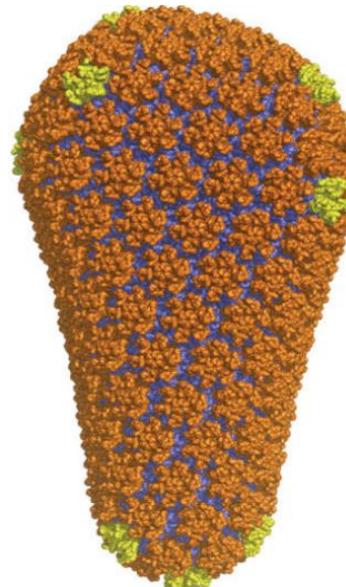
Extremely large biomolecules



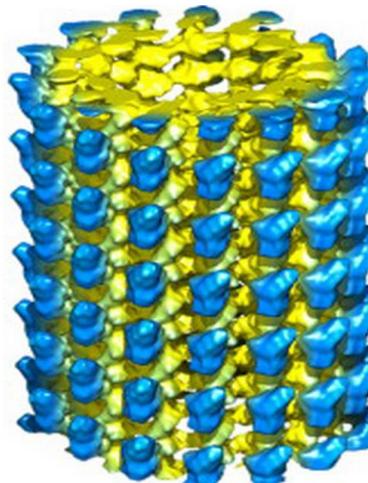
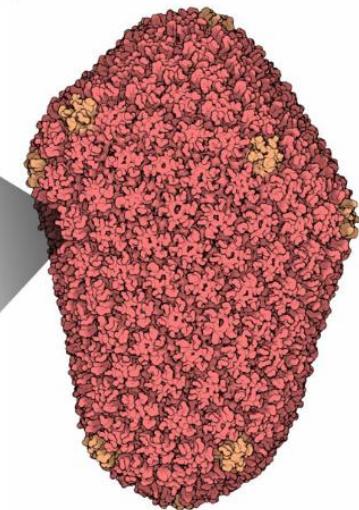
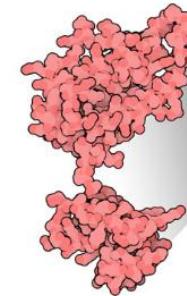
Vault particle



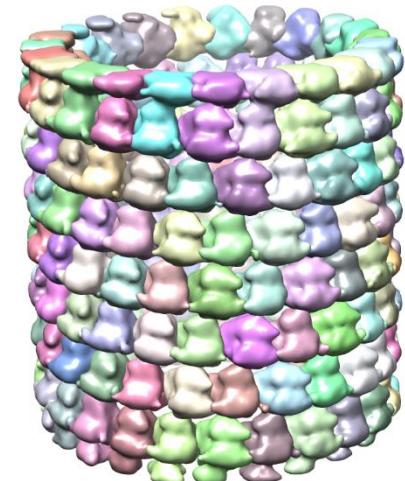
Poliovirus capsid



HIV virus capsid

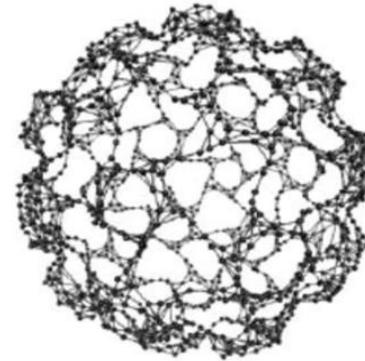
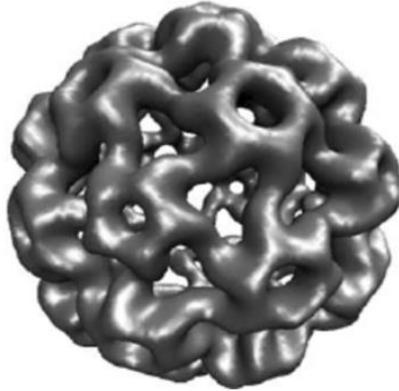
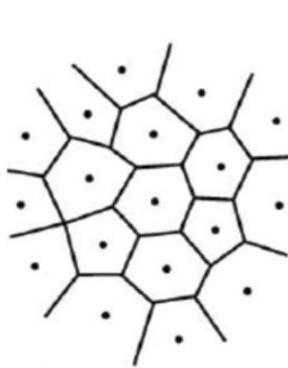


microtubule



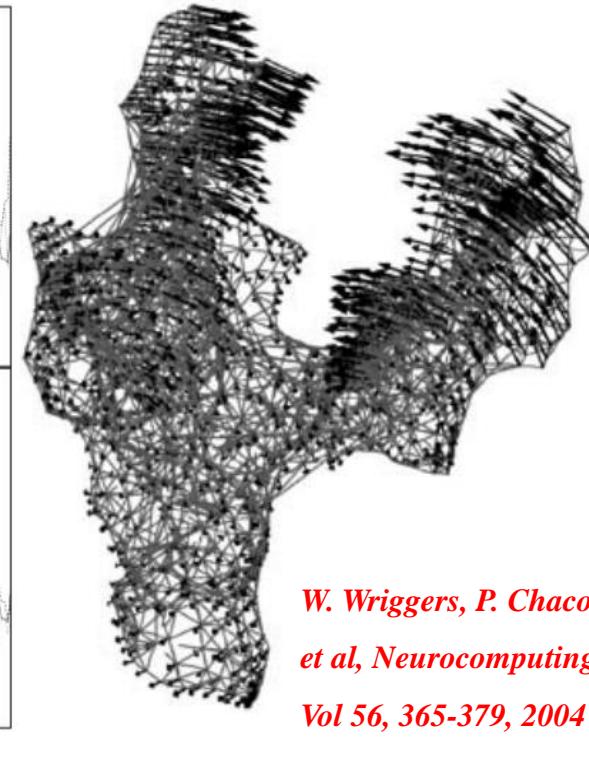
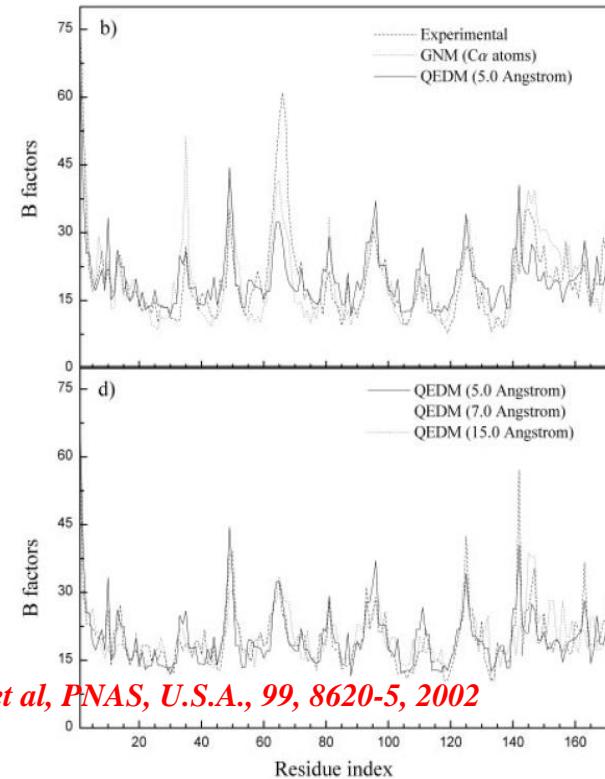
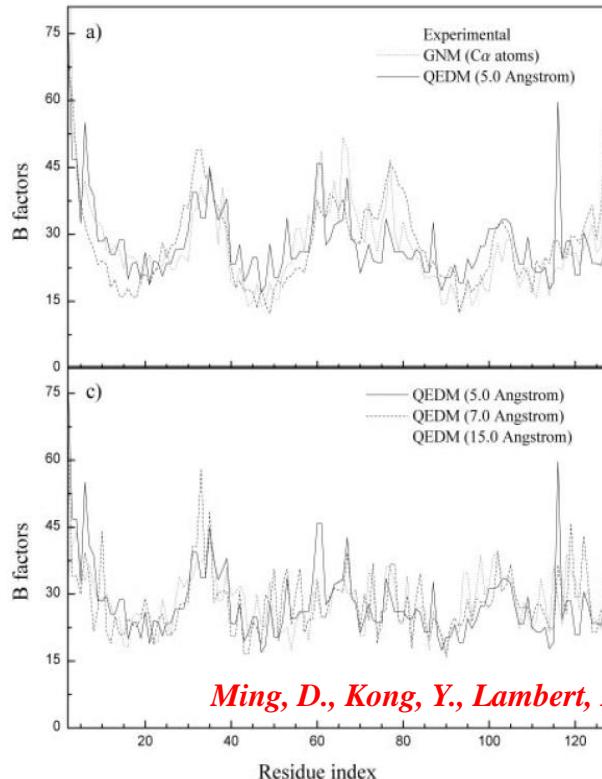
Quantized elastic deformational model (QEDM)

Voronoi Tessellation (Vector quantization)



W. Wriggers, P. Chacon, et
al, Neurocomputing, Vol
56, 365-379, 2004

Deformational motions are determined by GNM and ANM



Ming, D., Kong, Y., Lambert, M.A. et al, PNAS, U.S.A., 99, 8620-5, 2002

Multiscale virtual particle model

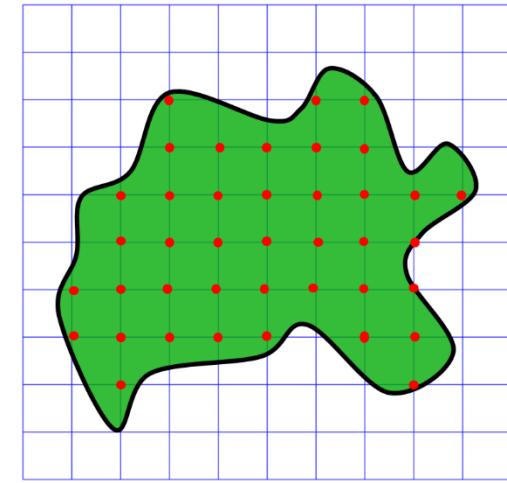
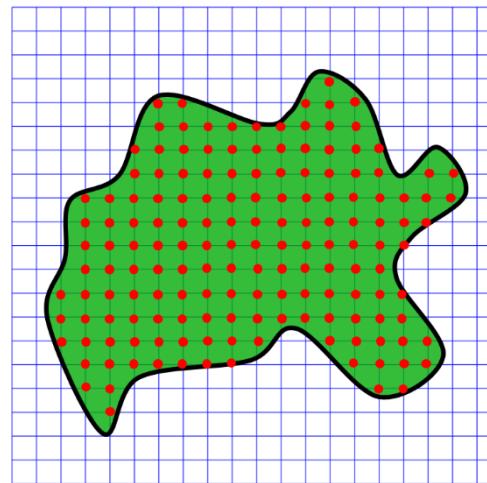
Virtual particle generation:

Various types of meshes:

Cartesian grid;

Tetrahedral mesh;

Hexahedron; Voronoi tessellation, etc.



Connection between particles:

$$\gamma(\mathbf{r}_I, \mathbf{r}_J, \Omega_I, \Omega_J, \mu^s(\mathbf{r}), \eta^{\text{MVP}}) = \gamma_1(\Omega_I, \Omega_J, \mu^s(\mathbf{r})) \cdot \gamma_2(\mathbf{r}_I, \mathbf{r}_J, \eta^{\text{MVP}})$$

Density

contribution:

$$\gamma_1(\Omega_I, \Omega_J, \mu^s(\mathbf{r})) = \left(1 + a \int_{\Omega_I} \mu^s(\mathbf{r}) d\mathbf{r}\right) \left(1 + a \int_{\Omega_J} \mu^s(\mathbf{r}) d\mathbf{r}\right)$$

Distance

contribution:

$$\gamma_2(\mathbf{r}_I, \mathbf{r}_J, \eta^{\text{MVP}}) = e^{-\left(\|\mathbf{r}_I - \mathbf{r}_J\|/\eta^{\text{MVP}}\right)^\kappa}, \quad \kappa > 0.$$

Multiscale virtual particle based Gaussian network model (MVP-GNM)

Potential function:

$$V^{\text{MVP-GNM}} = \frac{1}{2} \Delta \mathbf{r}^T L^{\text{MVP-GNM}} \Delta \mathbf{r}$$

$$L_{ij}^{\text{MVP-GNM}} = \begin{cases} -\gamma(\mathbf{r}_I, \mathbf{r}_J, \Omega_I, \Omega_J, \mu^s(r)) & i \neq j \\ -\sum_{i \neq j}^N L_{ij} & i = j \end{cases}.$$

Moore-Penrose pseudoinverse:

$$\sum_{k=2}^N \lambda_k^{-1} [\mathbf{v}_k \mathbf{v}_k^T]_{ii}$$

eigenvalue

Predicted b-factor:

$$\frac{8\pi^2}{3} < \Delta \mathbf{r}_i \cdot \Delta \mathbf{r}_i >, \quad \forall i = 1, 2, \dots, N$$

$$< \Delta \mathbf{r}_i \cdot \Delta \mathbf{r}_j > = \frac{3k_B T}{\gamma} (L^{-1})_{ij}, \quad \forall i = 1, 2, \dots, N$$

Multiscale representation of biomolecules

Kernel function:

$$\phi(\| r - r_j \|; \eta) = 1, \text{ as } \| r - r_j \| \rightarrow 0$$

We use kernel

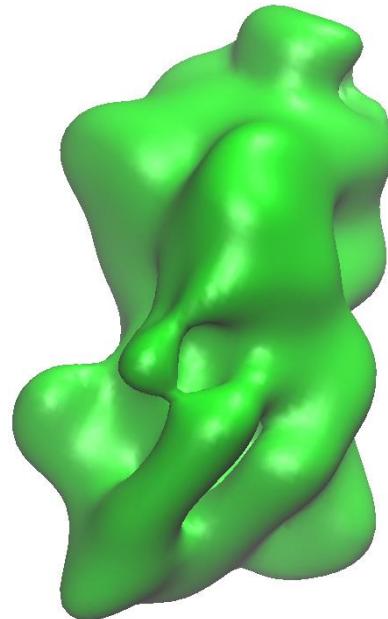
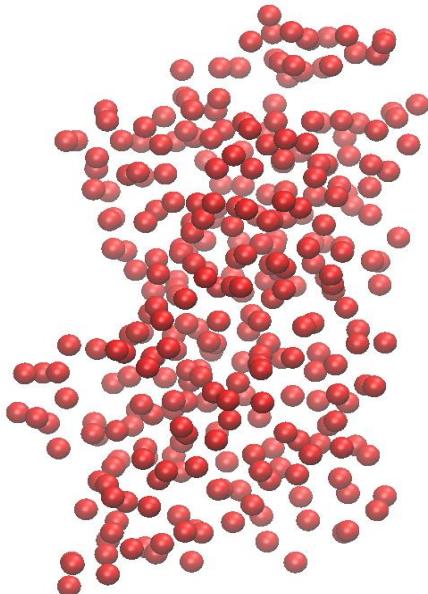
$$\phi(\| r - r_j \|; \eta) = e^{-(\| r - r_j \| / \eta)^2}$$

$$\phi(\| r - r_j \|; \eta) = 0, \text{ as } \| r - r_j \| \rightarrow \infty$$

Rigidity function:

$$\mu(r) = \sum_j^N w_j \phi(\| r - r_j \|; \eta)$$

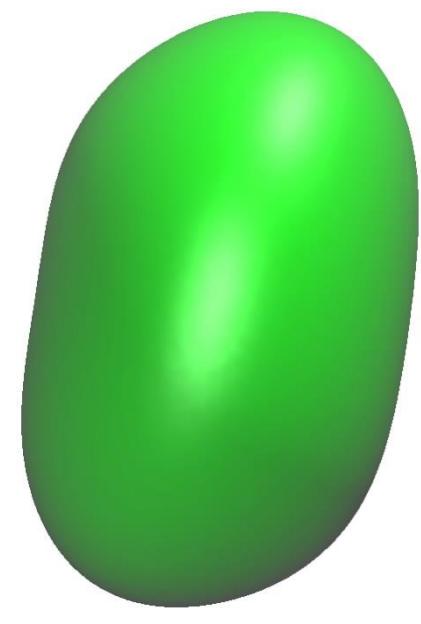
Protein ID: 2ABH



Resolution: 5Å



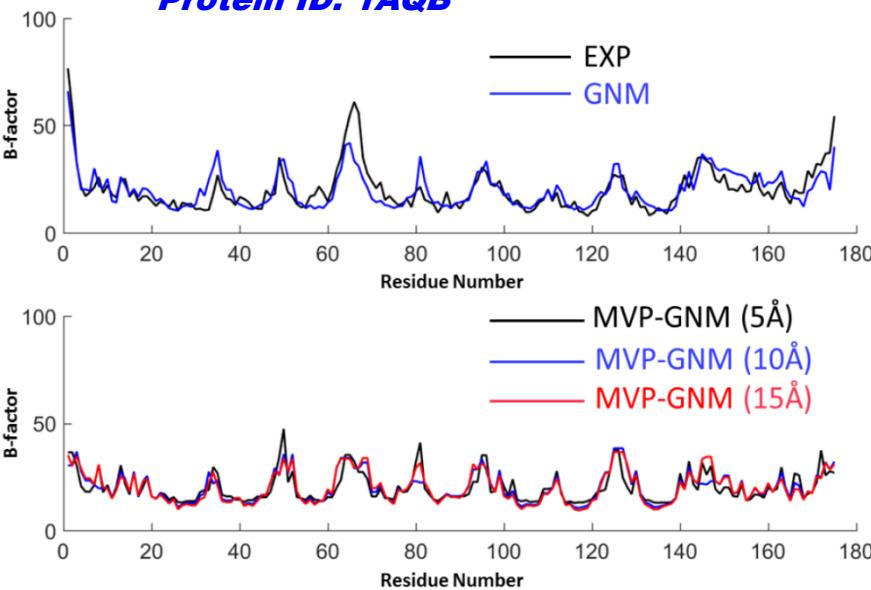
Resolution: 10Å



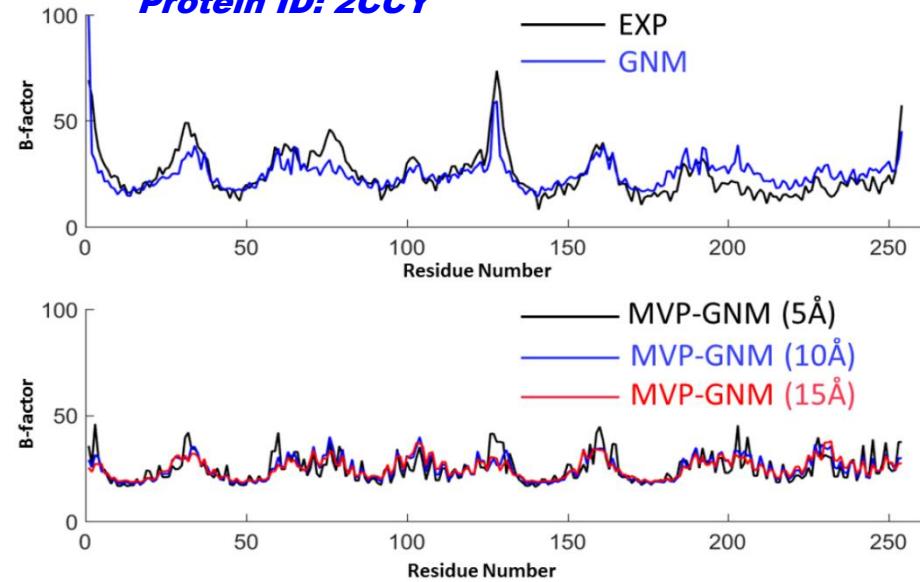
Resolution: 15Å

Validation of MVP-GNM

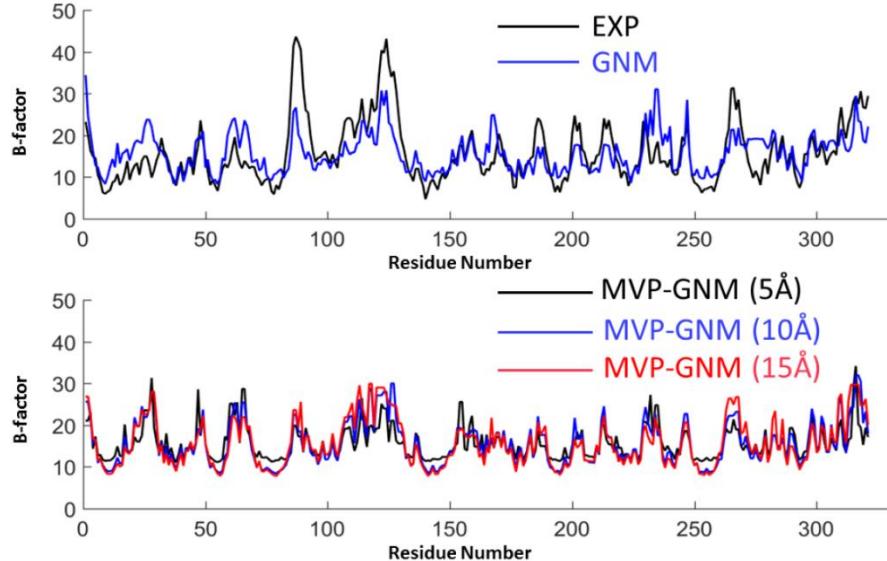
Protein ID: 1AQB



Protein ID: 2CCY



Protein ID: 2ABH



Cartesian grid with grid spacing 4 Å!!

	GNM	MVP-GNM (5 Å)	MVP-GNM (10 Å)	MVP-GNM (15 Å)
1AQB	0.822	0.666	0.657	0.699
2CCY	0.739	0.623	0.507	0.439
2ABH	0.647	0.550	0.731	0.775

Multiscale virtual particle based anisotropic network model (MVP-ANM)

Potential function:

$$V^{\text{MVP-ANM}} = \frac{1}{2} \Delta \mathbf{R}^T H^{\text{MVP-ANM}} \Delta \mathbf{R}$$

$$H_{IJ}^{\text{MVP-ANM}} = -\frac{\gamma_{IJ}}{r_{ij}^2} \begin{bmatrix} (x_J - x_I)(x_J - x_I) & (x_J - x_I)(y_J - y_I) & (x_J - x_I)(z_J - z_I) \\ (y_J - y_I)(x_J - x_I) & (y_J - y_I)(y_J - y_I) & (y_J - y_I)(z_J - z_I) \\ (z_J - z_I)(x_J - x_I) & (z_J - z_I)(y_J - y_I) & (z_J - z_I)(z_J - z_I) \end{bmatrix} \quad I \neq J.$$

$$H_{II}^{\text{MVP-ANM}} = - \sum_{I \neq J} H_{IJ}^{\text{MVP-ANM}}, \quad \forall i = 1, 2, \dots, N.$$

Moore-Penrose pseudoinverse:

$$\sum_{k=1}^{3N} \lambda_k^{-1} [\mathbf{v}_k \mathbf{v}_k^T]_{ii}$$

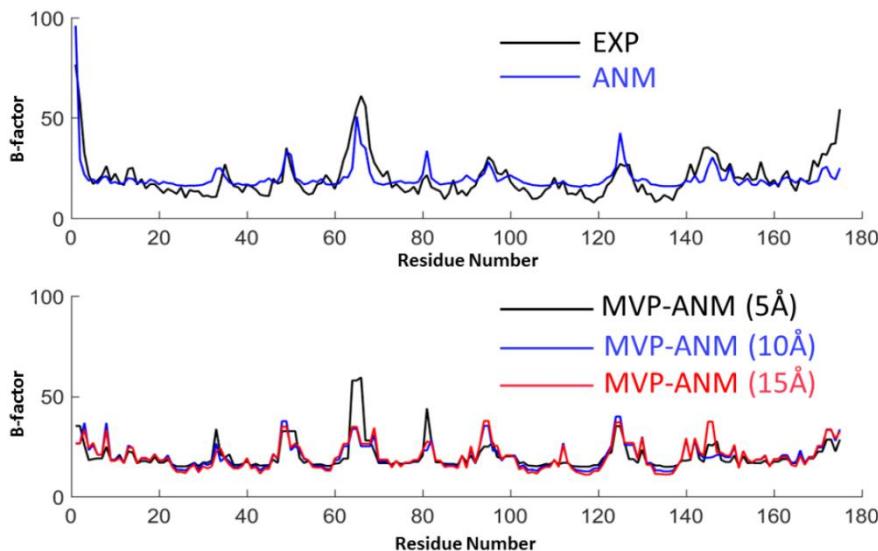
Predicted b-factor:

$$\frac{8\pi^2}{3} \sum_{j=3i-2}^{3i} < \Delta \mathbf{R}_j \cdot \Delta \mathbf{R}_j >, \quad \forall i = 1, 2, \dots, N.$$

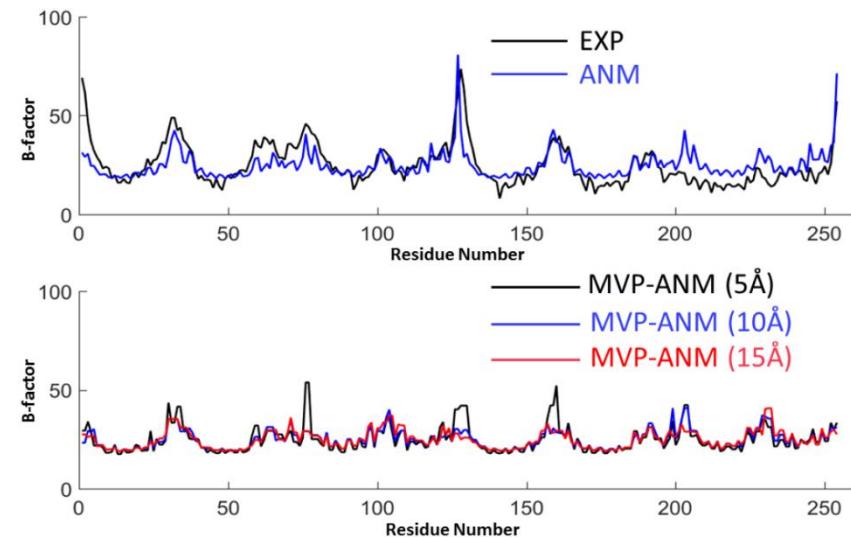
$$< \Delta \mathbf{R}_i \cdot \Delta \mathbf{R}_j > = \frac{3k_B T}{\gamma} (H^{-1})_{ij}, \quad \forall i, j = 1, 2, \dots, 3N.$$

Validation of MVP-ANM

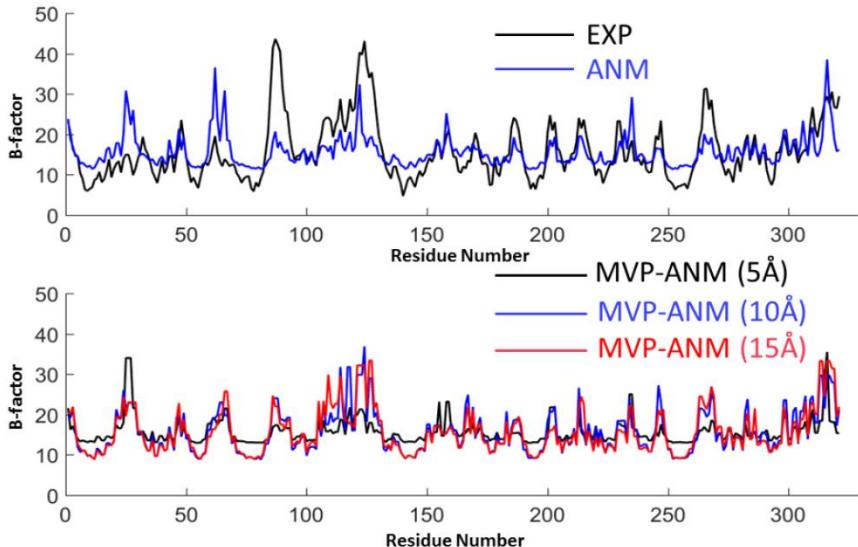
Protein ID: 1AQB



Protein ID: 2CCY



Protein ID: 2ABH



We use kernel

$$\phi(\| r - r_j \|; \eta) = e^{-(\| r - r_j \|/\eta)^2}$$

Cartesian grid with grid spacing 5 Å

	GNM	MVP-ANM (5 Å)	MVP-ANM (10 Å)	MVP-ANM (15 Å)
1AQB	0.725	0.696	0.593	0.646
2CCY	0.664	0.627	0.450	0.435
2ABH	0.548	0.442	0.743	0.760

ANM

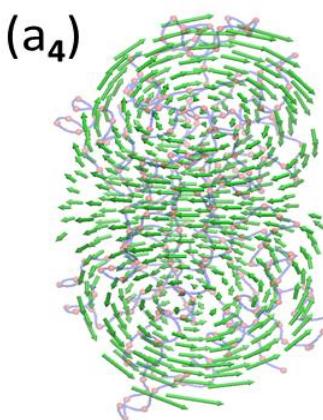
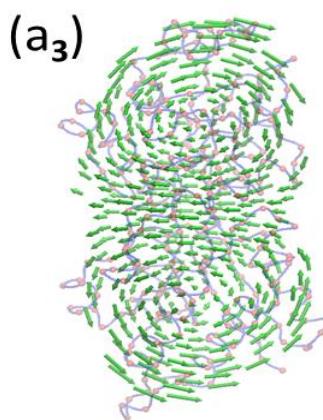
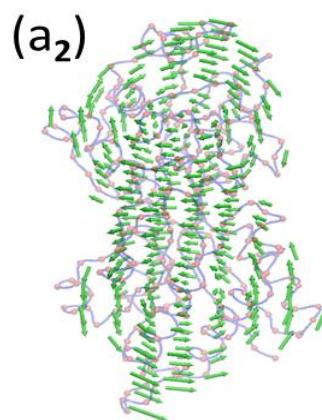
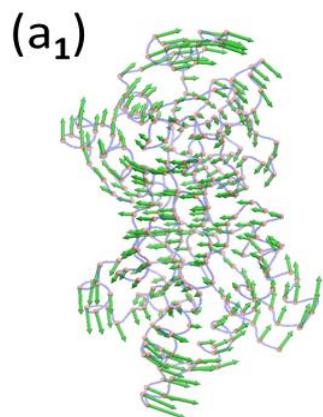
MVP-ANM

(5Å)

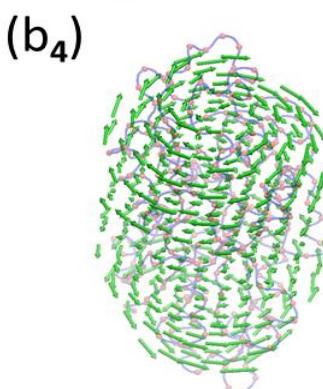
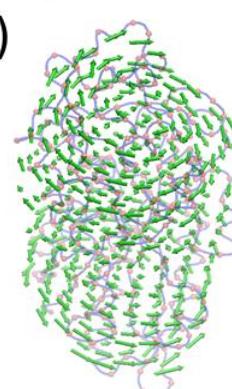
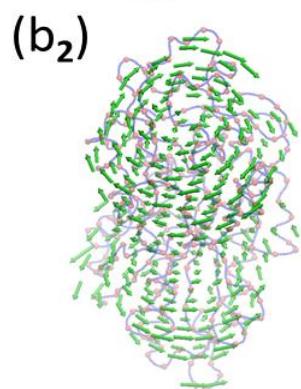
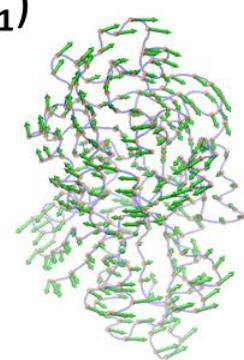
(10Å)

(15Å)

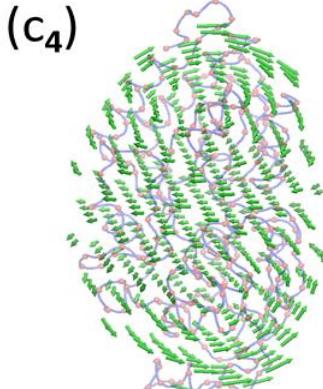
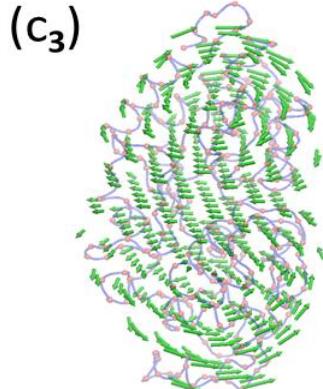
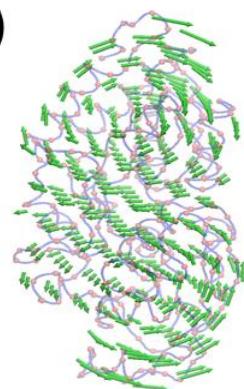
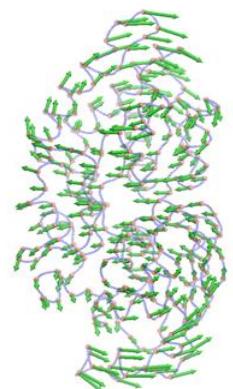
Mode 7



Mode 8

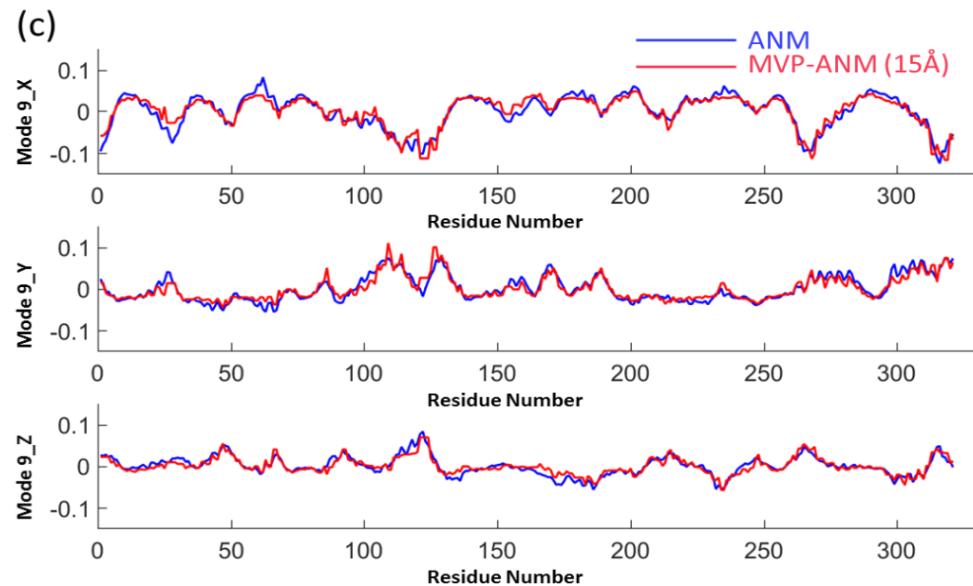
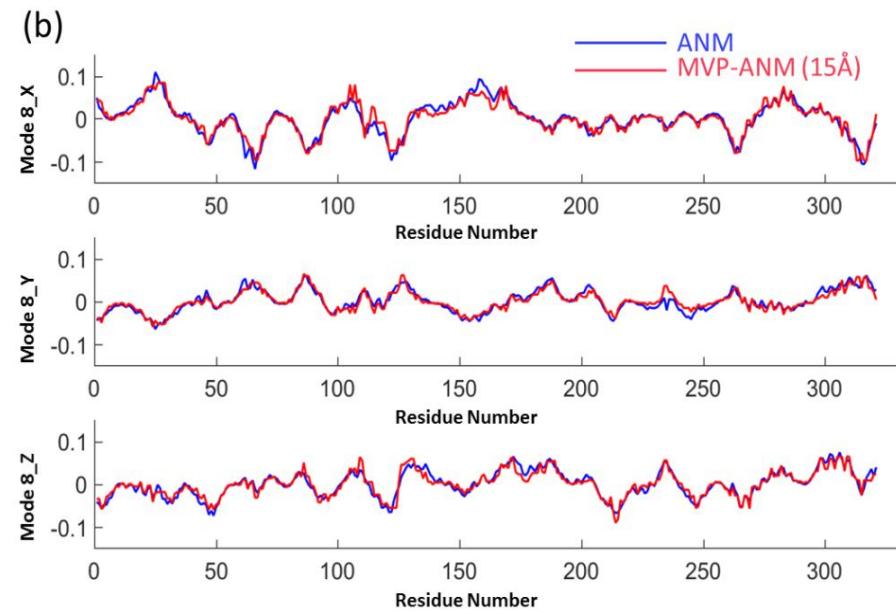
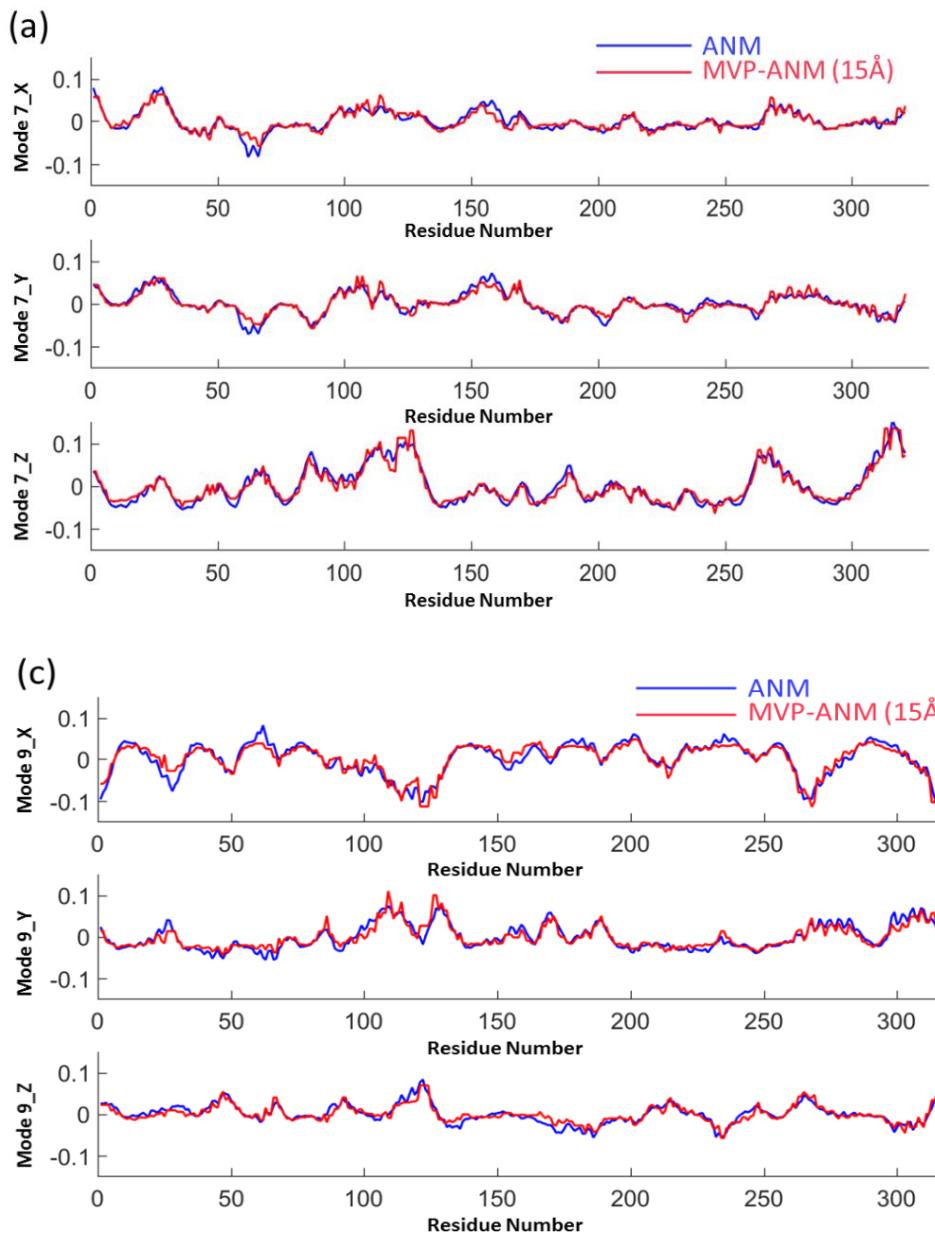


Mode 9



Protein ID: 2ABH

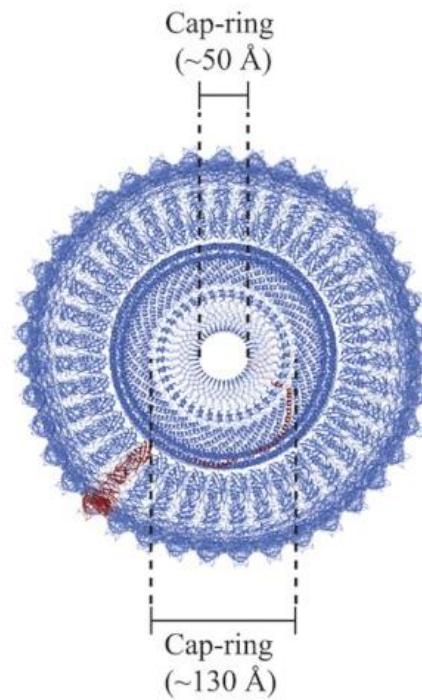
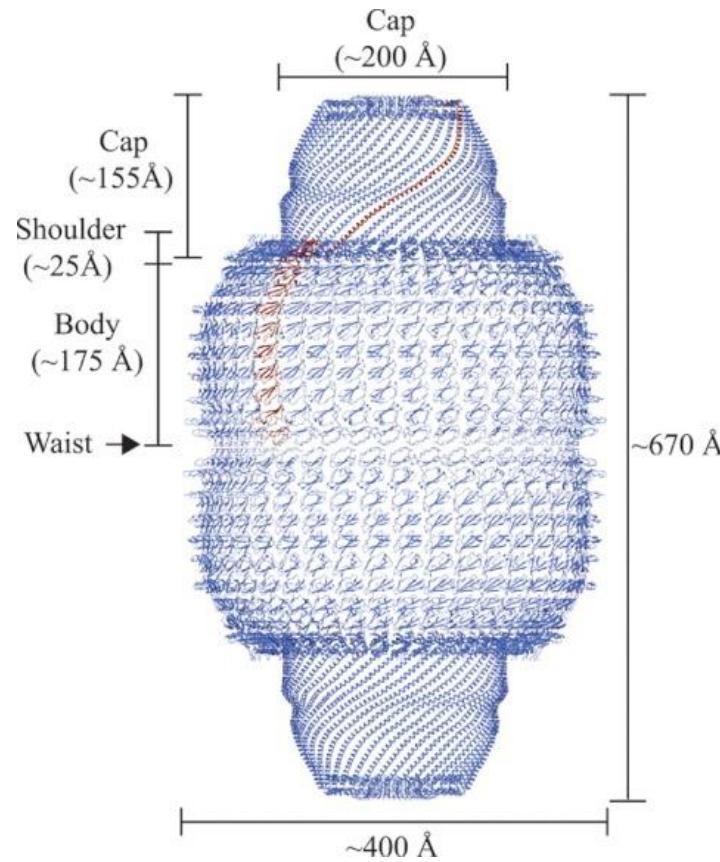
Normal modes for protein 2ABH



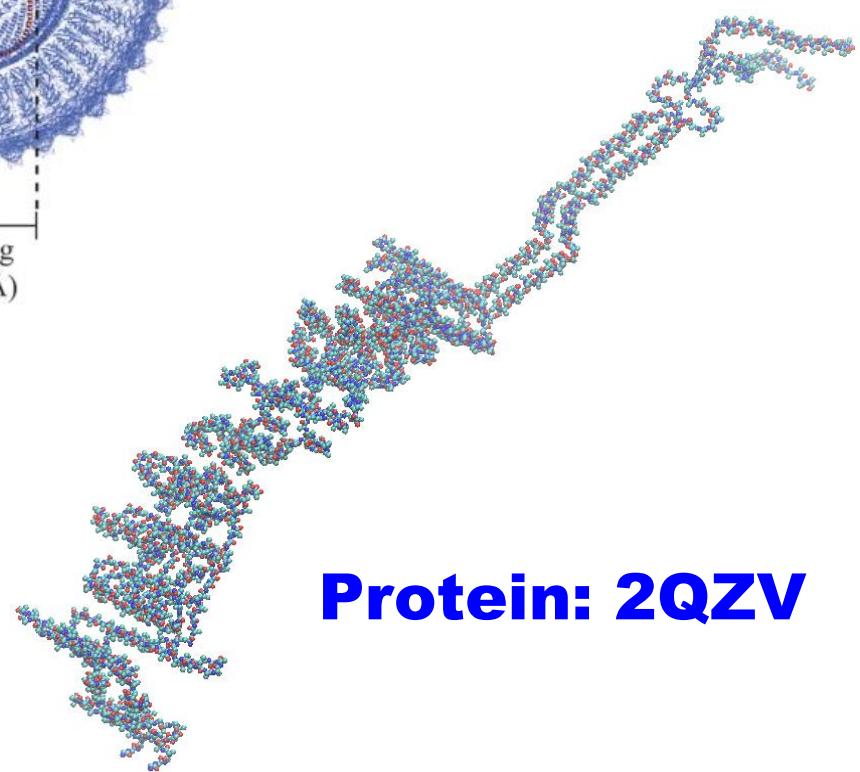
Resolution parameter is 15 Å

	MVP-ANM (X)	MVP-ANM (Y)	MVP-ANM (Z)
Mode 7	0.914	0.929	0.960
Mode 8	0.948	0.930	0.943
Mode 9	0.936	0.921	0.906

The dynamics of Vault Shell

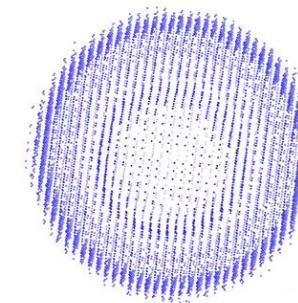
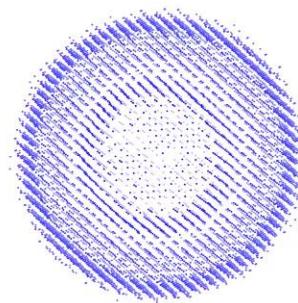
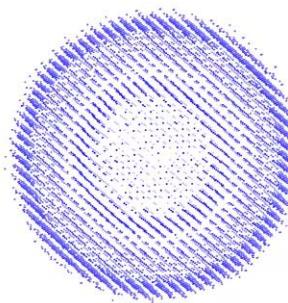
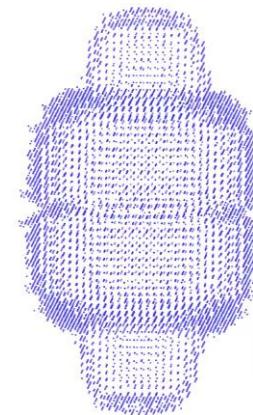
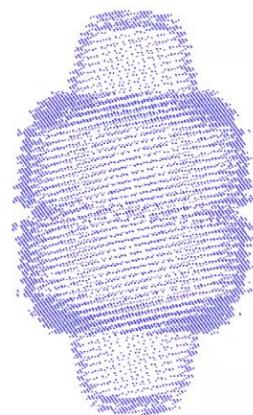
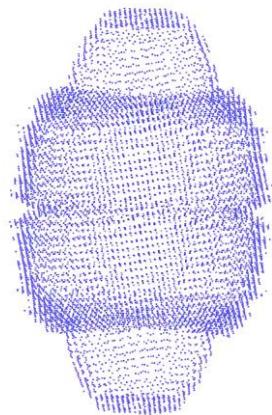


**96 copies of vault
protein 2QZV**



Protein: 2QZV

Multiscale virtual particle based elastic network model of Vault



Group members



Grant support

NTU-JSPS (2019-2022)

Alibaba-NTU (2020-2021)

Merlion (2020-2022)

MOE-Tier 1 (2018-2021, 2019-2022)

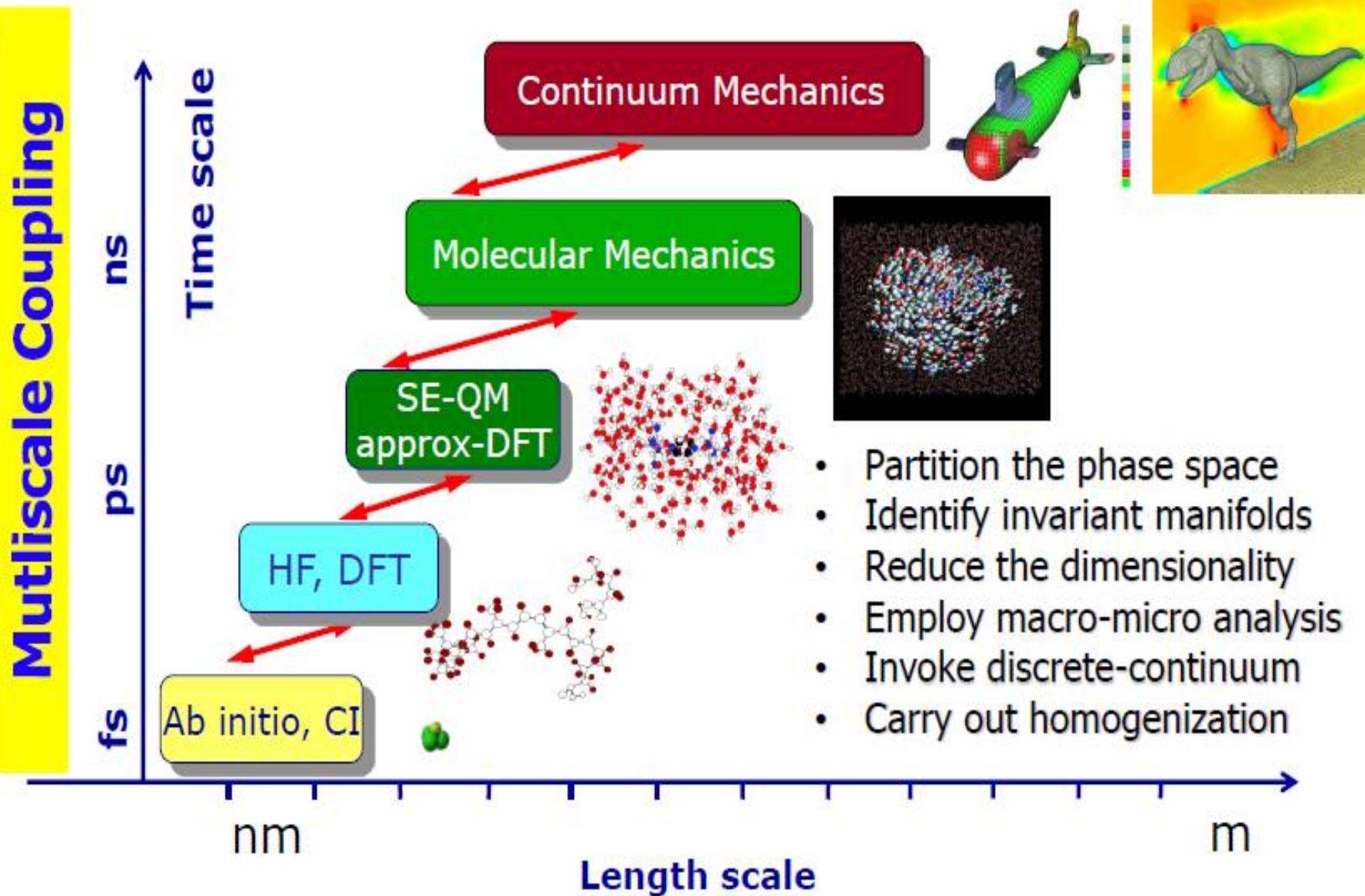
MOE-Tier 2 (2018-2021, 2021-2024)



Biophysics

- 1) Biophysics is an interdisciplinary science that applies approaches and methods traditionally used in physics to study biological phenomena.
- 2) Biophysics covers all scales of biological organization, from molecular to organismic and populations
- 3) Molecular biophysics applies physical approach to model biomolecular systems and understand their interactions and structure-function relationship.
- 4) Unlike data-driven bioinformatics and knowledge-based systems biology, biophysics is mechanistic.

Hierarchy of Methods



Molecular Mechanics (MM)

Using classical particle assumption.

Newton's second law: $m_i \frac{d^2}{dt^2} \mathbf{r}_i = \mathbf{F} = -\nabla U(\mathbf{r}_1(t), \mathbf{r}_2(t), \dots, \mathbf{r}_n(t))$

Approximation: $U = \sum_{\text{bonds}} K_d (d - d_0)^2 + \sum_{\text{angle}} K_\theta (\theta - \theta_0)^2$

+ $\sum_{\text{dihedrals}} K_\chi (1 + \cos(n\chi - \delta))$

+ $\sum_{\text{nonbond}} \left\{ \epsilon_{ij} \left[\left(\frac{R_{ij}^{\min}}{\mathbf{r}_{ij}} \right)^{12} - \left(\frac{R_{ij}^{\min}}{\mathbf{r}_{ij}} \right)^6 \right] + \frac{q_i q_j}{\epsilon \mathbf{r}_{ij}} \right\}$

Langevin equation: $m \ddot{\mathbf{r}} = -\nabla V(\mathbf{r}) - \gamma \dot{\mathbf{r}} + \sqrt{2\gamma k_B T} R(t)$,

where, $\langle \mathbf{R}(t) \rangle = 0$, $\langle \mathbf{R}(t)\mathbf{R}(t') \rangle = \delta(t - t')$

Explicit solvent/Implicit solvent/Coarse Grained

MM Software: AMBER/CHARMM/NAMD/TINKER/GROMOS

Research issues: high-order force fields, coarse grained methods, implicit MD, etc.

Foundations of Biophysics

Continuum Mechanics (CM)

Hydrodynamics (HD)

Electrodynamics (ED)

Thermodynamics (TD)

Molecular Mechanics (MM)

Kinetic Theory (KT)

Statistical Mechanics (SM)

Quantum Mechanics (QM)

Computability

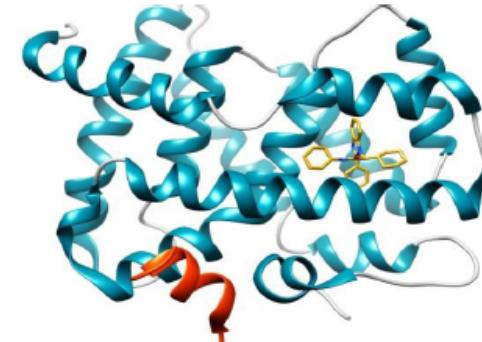
Reliability

Bioinformatics

- 1) **Sequence analysis** (DNA sequencing, sequence assembly, genome annotation, comparative genomics, pan genomics, computational evolution, genetics, cancer mutation, etc.).
- 2) **Gene and protein expression** (Gene expression analysis, protein expression analysis, gene regulation, genotype–phenotype map, etc.).
- 3) **Systems biology** (Pan networks, integrated systems analysis).
- 4) **Cellular organization** (Microscopy and image analysis, protein localization, membrane mechanics, chromatin analysis, etc.).
- 5) **Structural bioinformatics** (Biomolecular structure and interaction, structure-function relationship, protein folding, protein design, etc..)
- 6) **Database and software.**
- 7) Unlike biophysics, bioinformatics is data-driven.

Protein Structure Prediction

SVYDAAAQLTADVKKDLRDSW
KVIGSDKKGNGVALMTTLFAD
NQETIGYFKRLGNVSQGMAND
KLRGHSITLMYALQNFIDQLD
NPDSL DVCS



- 1) Understand protein structure-function relationship
- 2) Design protein with desired function
- 3) Drug development
- 4) Methods (knowledge-based):

Template-based modeling (homology modeling) is used when there is one or more similar known structures in PDB.

Ab initio structure prediction (e.g., Rosetta) is used when one cannot find any similar structure.

Deep learning (e.g., AlphaFold, CNN, RNN)

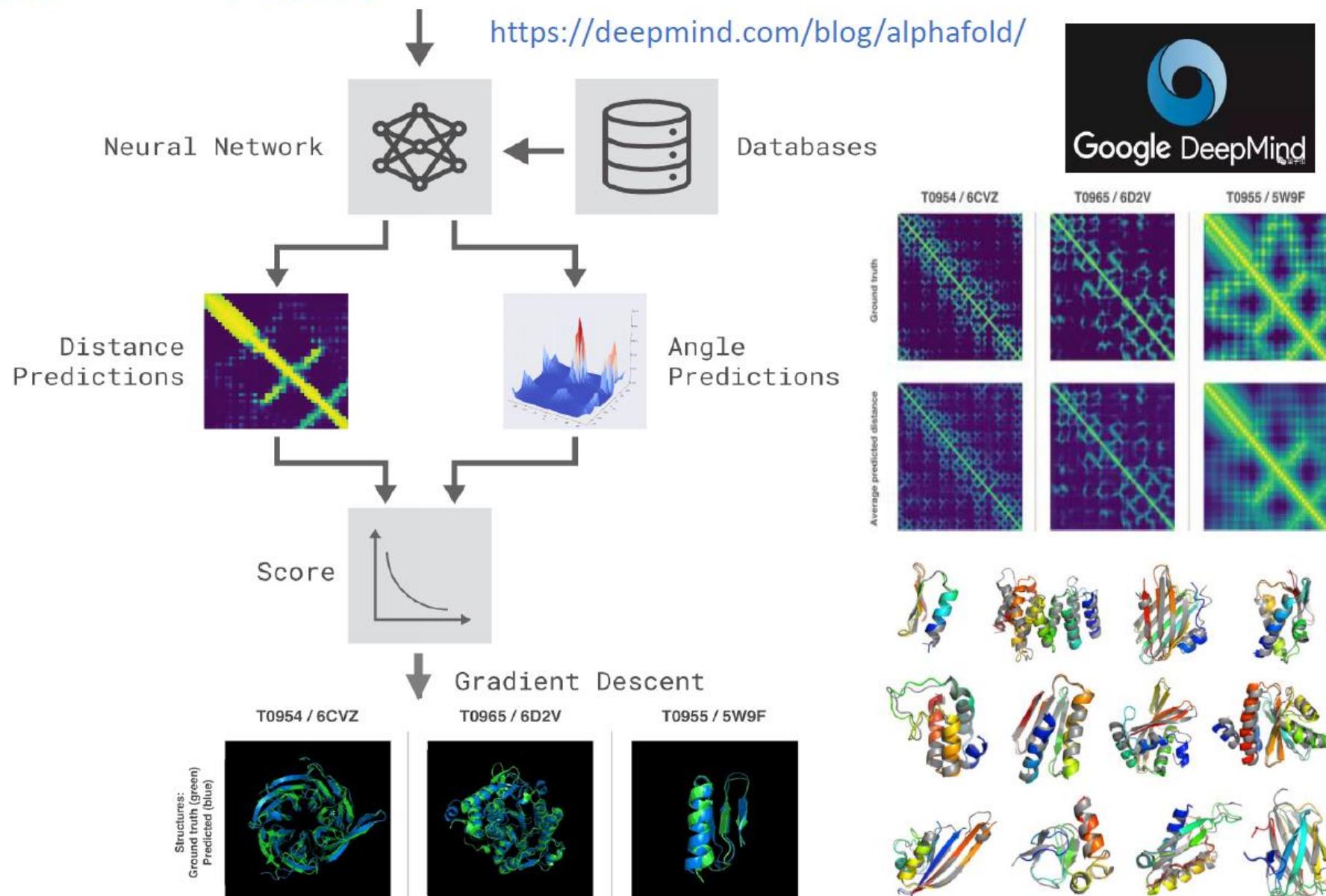
- 5) Evaluation:

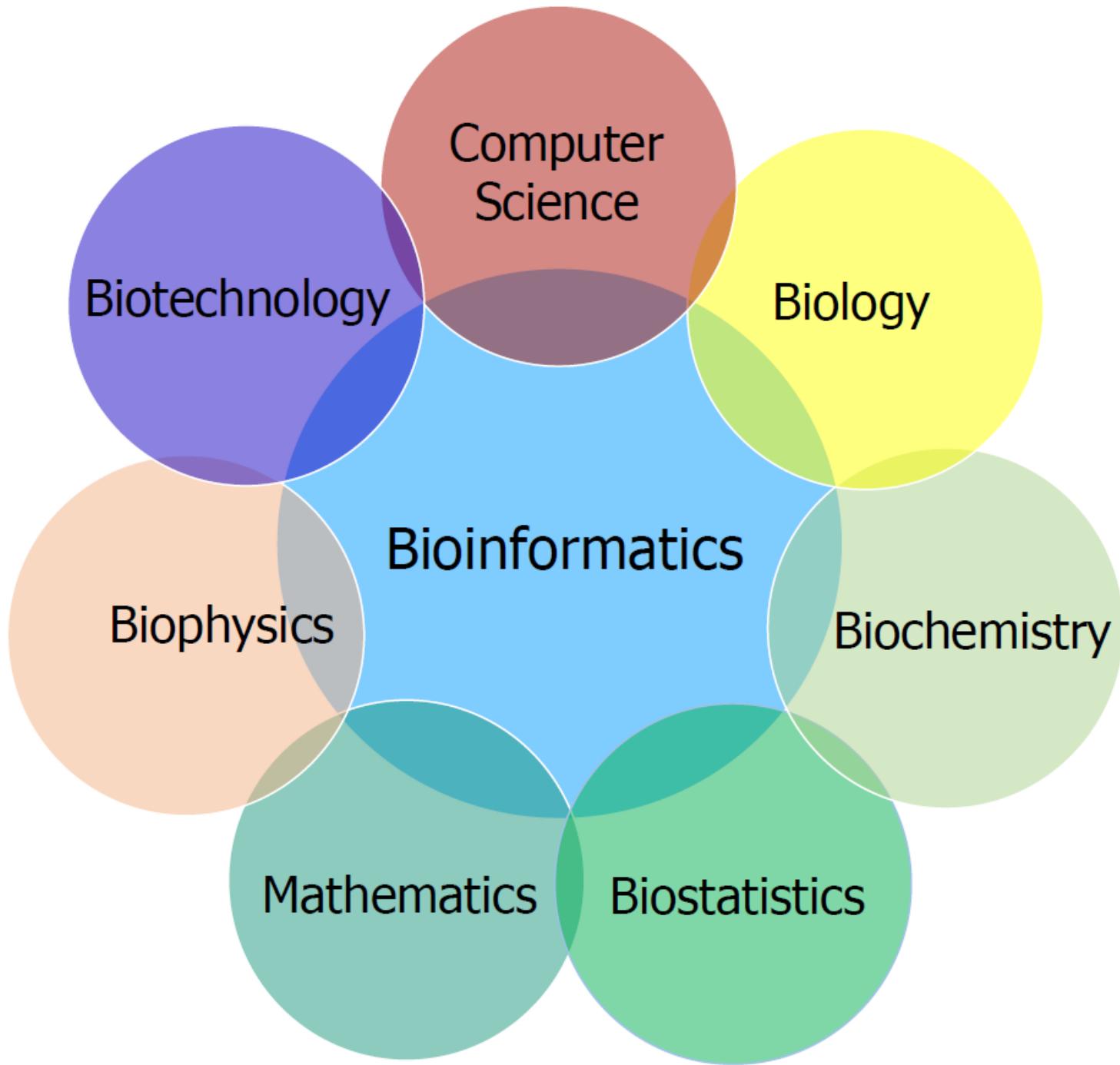
Critical Assessment of protein Structure Prediction (CASP)
Knowledge-based methods win; QM/MM do not work well.



Protein Sequence

SQETRKKCTEMKKKFKNCEVRCDESNHCVENVRCSDTKYTL

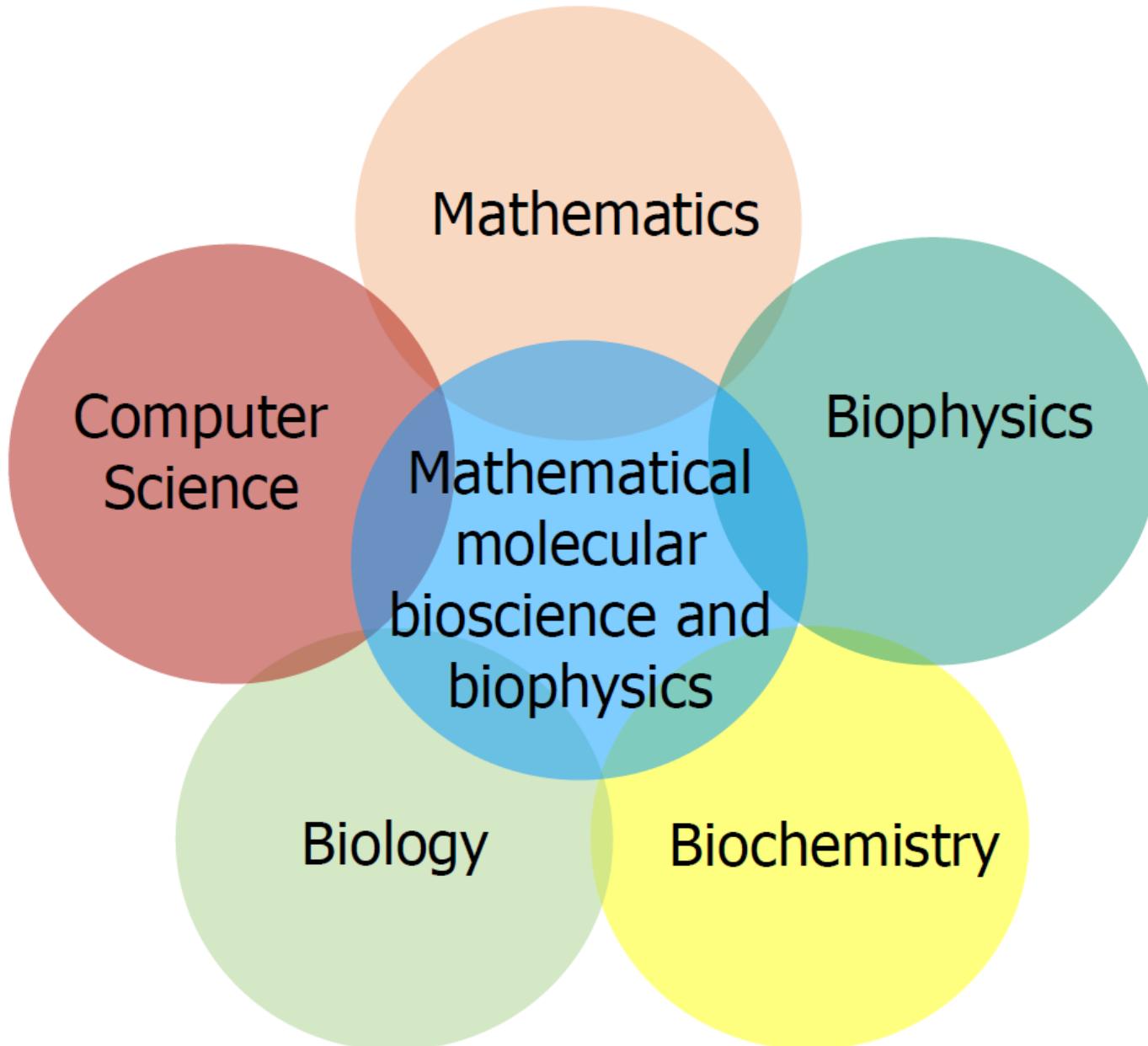




New Trends in Biological Science

- 1) New bioscience is based on molecules and/or omics.
- 2) Integrative biology (from molecules, organism to environment).
- 3) Integration of biophysics, systems biology, and bioinformatics.
- 4) Integration of mathematics, data science, theoretical biology and experimental biology.
- 5) Mathematical molecular bioscience and biophysics.
- 6) Quantitative systems pharmacology (from systems biology, biomechanics, systems physiology to systems pharmacology).
- 7) Personalized medicine (precision medicine).

Mathematical Molecular Bioscience and Biophysics



Mathematical Molecular Bioscience and Biophysics

- 1) It concerns the mathematical foundation of biological science.
- 2) It is based on molecular bioscience and omics in contrast to macroscopic biosciences.
- 3) It overlaps with molecular biophysics, systems biology and bioinformatics but is distinguished from any mathematical biology that is macroscopic and phenomenological.
- 4) It exploits existing mathematics for describing biological observations and dynamics.
- 5) It makes use of computational algorithms and methods from mathematics, machine learning and statistics.
- 6) It has applications to a wide range of biological problems, including protein design, drug discovery, precision medicine, to mention only a few.
- 7) It generates new mathematics from biological challenges.

Mathematics Commonly Used in Molecular Bioscience

Geometry	D.W. Sumners, Isabel .K. Darcy , Mariel Vazquez, Dorothy Buck, Tamar Schlick, Erica Flapan, Christian Reidys, Yusu Wang, Peter Rogen, Jack Quine,
Topology	
Algebra	Christine Heitsch, David Murrugarra, Reidun Twarock
Group theory	Natasha Jonoska, R Brijder, HJ Hoogeboom, Julie Mitchell
Combinatorics	
Analysis	M Karplus, M Levitt, A Warshel, B Honig, E Alxov, A Onufriev,...
Calculus/Variation	B.S. Eisenberg, Chun Liu, Weishi Liu, Yun Kyong Hyon, TC Lin, JL Liu, TL Horng, YN Young, HX Huang, Lei Zhang, Tom Chou
ODE and PDE	J.A. McCammon, Michael Holst, Jingfang Huang, Benzhuo Lu, Nathan Baker, Bo Li, LT Cheng, MX Chen, Shenggao Zhou,
Numerical analysis	Keith Promislow, Shibin Dai, Nir Gavish, Robert Krasny, DX Xie, LR Scott, Wei Cai, ZL Xu, Amit Singer, D. Kozakov, R Rizzo, D. Green, R Ryham, LJ Cowen, ...