EPFL

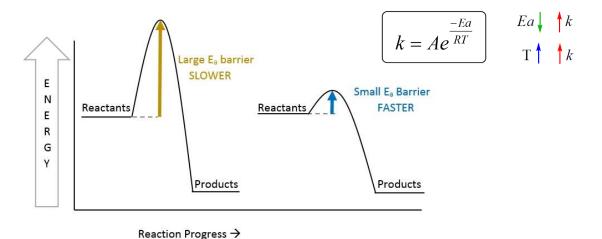
BIO-212 - Lecture 14 Biomolecule Engineering



 École polytechnique fédérale de Lausanne

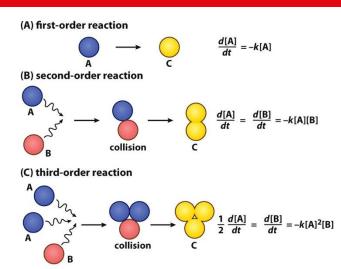
Lecture 13 - Summary

Reaction rates and Activation Energies

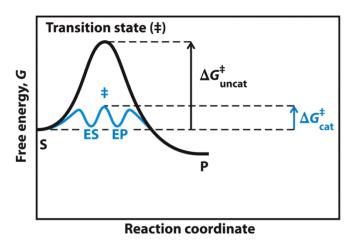


Reaction order

- Order of the reaction depends on the number of molecules whose quantities impact the reaction rate
- Rate constant units depend on the order



Catalysis and enzymes

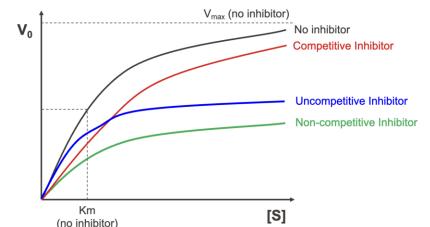


$$E + S \rightleftharpoons ES \rightleftharpoons EP \rightleftharpoons E + P$$

$$V_0 = \frac{V_{\text{max}}[S]}{[S] + K_m}$$

- Michaelis-Menten equation

Enzyme inhibition mechanisms



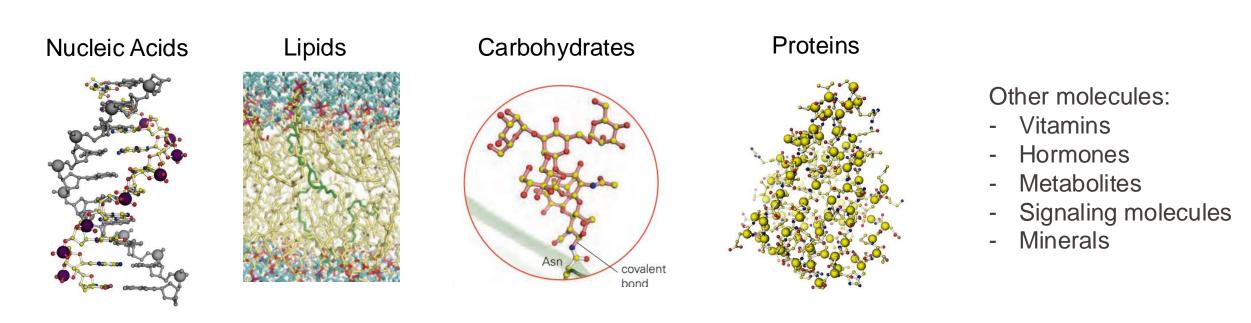
- Inhibition depends on the K₁ and the amount of inhibitor
- Often expressed as α

$$\alpha = 1 + \frac{[1]}{k_l}$$



Biomolecule diversity

- The adult human body is composed of ~37*10¹² cells, and at least as many as part of the human microbiota, with the total number of molecules estimated at ~10¹⁰-10¹² per cell (including water)
- Each main type of biomolecule features chemical and functional diversity. For example, there are >10³ different lipid species in eucaryotic membranes, and ~10⁴ different proteins per human cell.



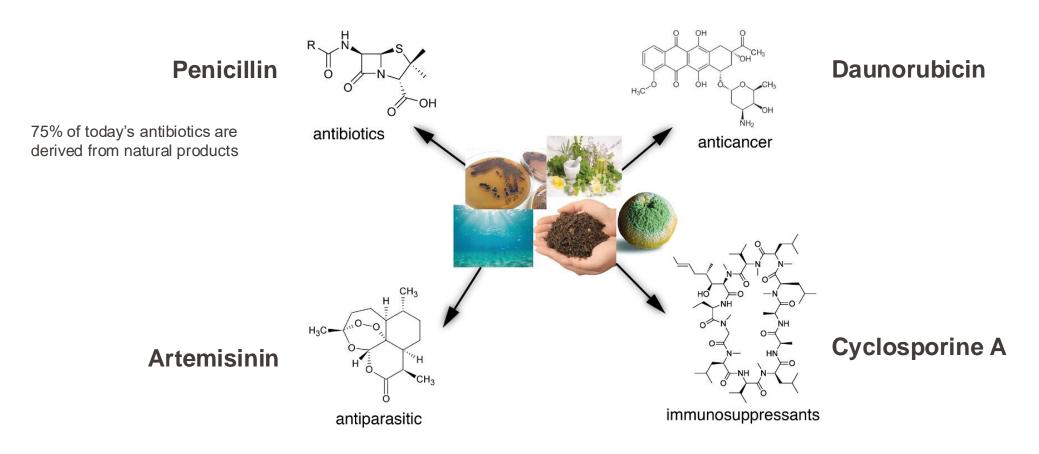
• Compositional and functional diversity further increases across different cells, tissues, organisms, altogether comprising a massive library of biomolecules which is constantly updated through evolution

3



Natural products are important therapeutics

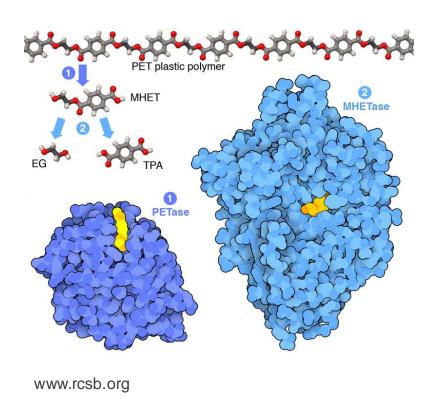
- In addition to enriching the repertoire of biomolecules of similar functions, great biological diversity can result in novel molecules with unique functions
- Natural products are an important group of biomolecules (usually metabolites) that are widely exploited for therapeutic purposes

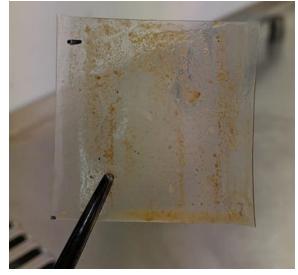




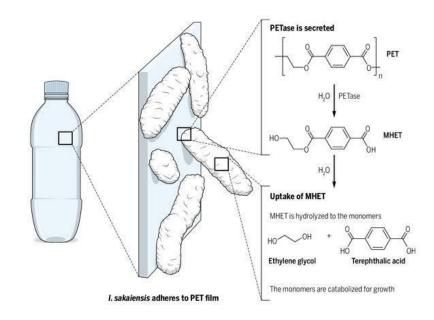
Evolution as a constant driver of molecule diversity

- Plastic-degrading microbes (e.g., *Ideonella sakaiensis*) have been discovered at a plastic-recycling facility in Japan
- PETase and MHETase enzymes drive the process of degradation
- Likely evolved from lipase, esterase and/or endopeptidase enzymes





Bacteria on a plastic sheet



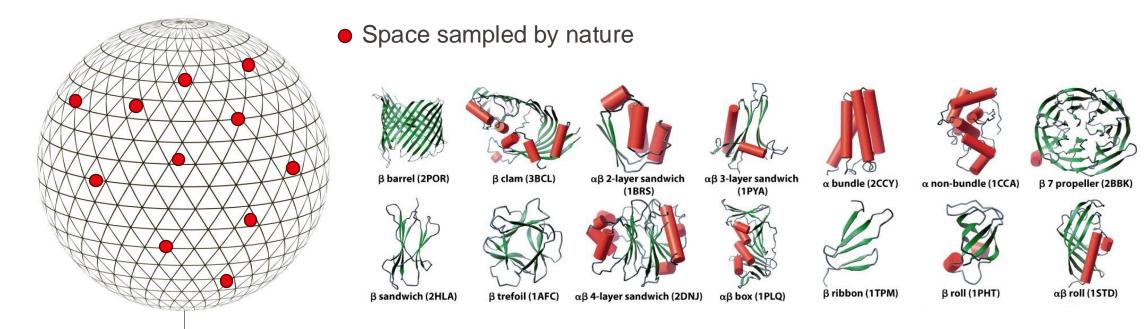


Possible Chemical Space >> Natural Diversity

• One good example is **protein sequence space**:

Protein Sequence Space

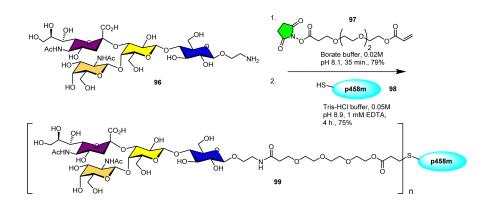
- For 10 amino acids the number of possible linear peptides are -10^{13}
- For 50 amino acids the number of possible linear peptides are -10^{65}
- For 100 amino acids the number of possible linear peptides are -10^{130}
- For 300 amino acids the number of possible linear peptides are -10^{325}
- The number of atoms in observable universe is estimated at 1080



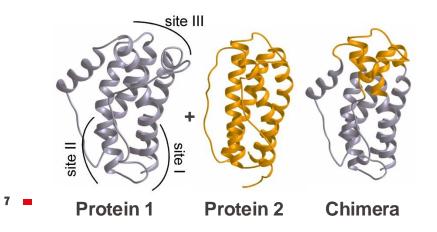


Engineering novel biomolecule functions

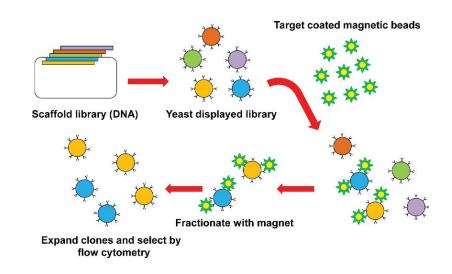
Chemical biology methods



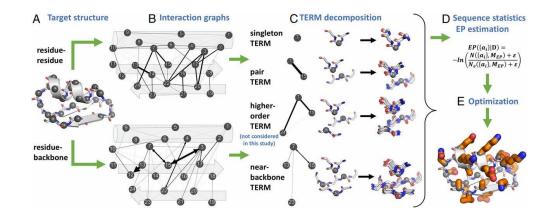
Simple biochemical manipulations



Library screening methods



Computational design methods





Chemical biology tools for biomolecule engineering

- Chemical biology is the study of the chemical groups and chemical reactions involved in biological processes, incorporating the disciplines of bioorganic chemistry, cell biology and pharmacology
- Through application of chemical toolsets, biomolecules with different functional groups can be engineered, allowing to modulate their properties

Sulfonated coumarin cage, DIPEA DCM, 24h, rt., 27-68 % DAG Transbilayer movement (ms-min) B before uncaging 6s after uncaging 200 s after uncaging C1-EGFP 10 jum

Vaccine design **MBS (66)** AcHN Tumor-associated sugar (poorly immunogenic) **KLH**

KLH conjugation makes it highly

immunogenic upon vaccination



Bioorthogonal chemistry

- Proteins and several other classes of biomolecules are sensitive to buffer components, pH, and/or temperature. Therefore, chemical reactions need to be performed under "mild" conditions (i.e., aqueous buffer, neutral pH, room temp).
- Bioorthogonal chemistry represents a class of high-yield chemical reactions that proceed rapidly and selectively in biological environments without side reactions towards endogenous functional groups. It partially overlaps with "click chemistry" high-yielding reactions wide in scope and simple to perform.

2022 Nobel Prize in Chemistry





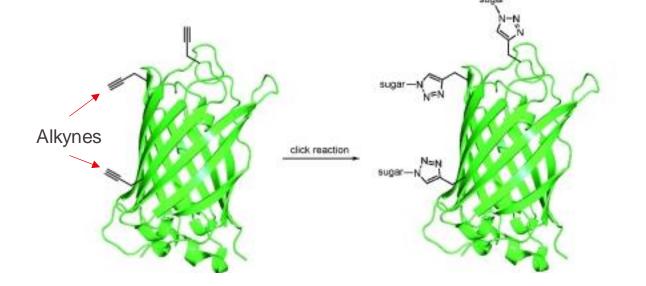


© Nobel Prize Outreach. Photo: Stefan Bladh S Carolyn R. Bertozzi

© Nobel Prize Outreach. Photo: Stefan Bladh Morten Meldal

© Nobel Prize Outreach. Photo Stefan Bladh K. Barry Sharpless

Example: Azide - Alkyne reaction

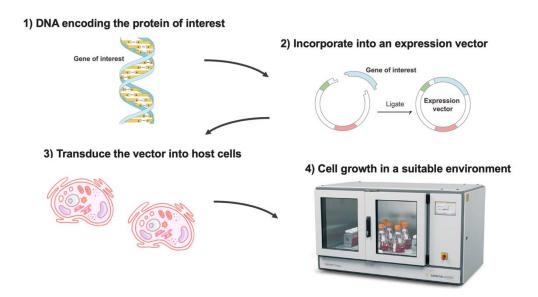




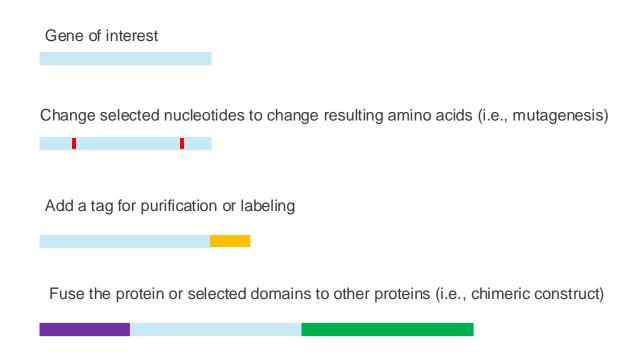
Simple biochemical manipulations

- A category of methods that include manual, structure- or sequence-assisted modifications of biomolecules using biochemical and molecular biology tools. Usually low-throughput due to manual engineering component.
- Examples include mutagenesis, generation of chimeric proteins, tagging, enzymatic modifications in vitro...

Protein production workflow:



Common molecular biology modifications:



• For proteins, the modifications are typically done on the genetic level (e.g., by PCR) followed by expression in a suitable cell system

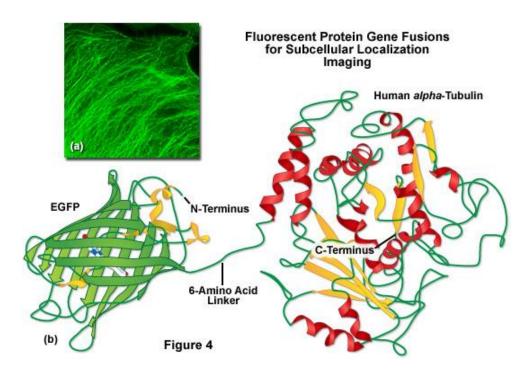


Simple biochemical manipulations

- Mutagenesis is routinely applied to validate the importance of certain molecular interactions (e.g., in the binding site between two proteins) by disrupting/enhancing them
- Fusion between biomolecules can be introduced to combine their functions or for the purpose of labeling.

Protein A Protein B

Genetic fusion to GFP



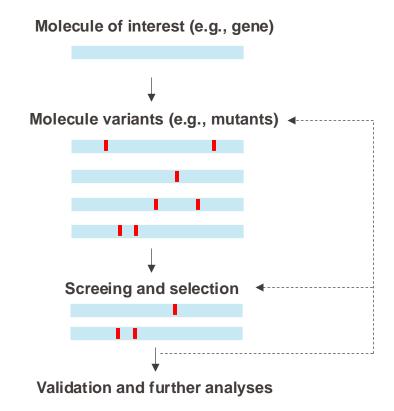
• These molecule modifications are not always "simple" and usually require screening of experimental conditions

EPFL

Library screening methods

- These methods are based on the generation of large libraries of molecules with partially or completely random chemical properties (e.g., amino-acid sequence). The library is then evaluated in a high-throughput manner to identify biomolecules capable of performing certain molecular functions (e.g., binding, enzymatic activity).
- There are conceptually 3 phases:
 - Library generation
 - Library screening (one or multiple rounds)
 - Validation of hits

- Typical library sizes (# of unique molecule variants):
 - Small molecules: ~104 106
 - Peptide arrays: ~10⁵ 10⁶
 - Genetic libraries: ~106 108
- Much bigger libraries can also be generated but they present problems in the form of QC and workload

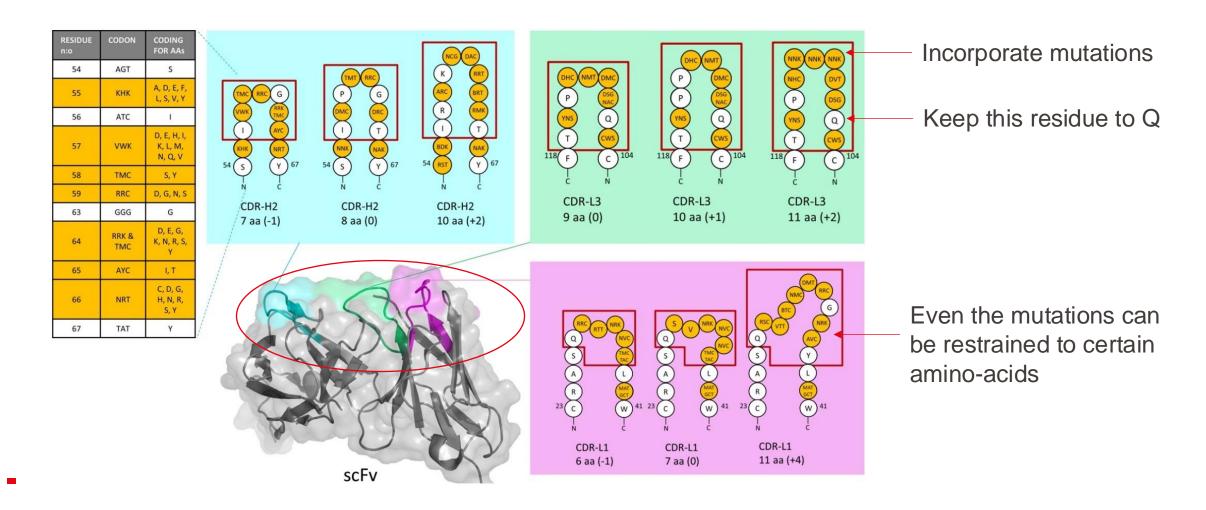


• Desired molecular properties are generated by cumulative selection of modifications that lead to desired function (as measured by some type of biophysical method)



Library screening methods

- Libraries can be tailored/customized to include modifications at only selected areas or to be completely random
- For example, if designing antibody libraries the focus is on complementarity determining regions (CDRs) which regulate antibody-antigen interactions, while the remainder of the scaffold is kept constant

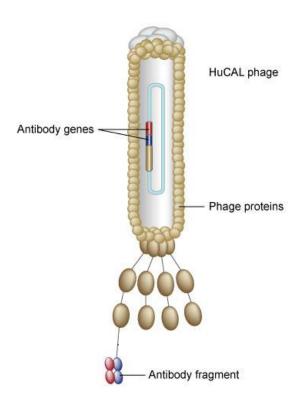


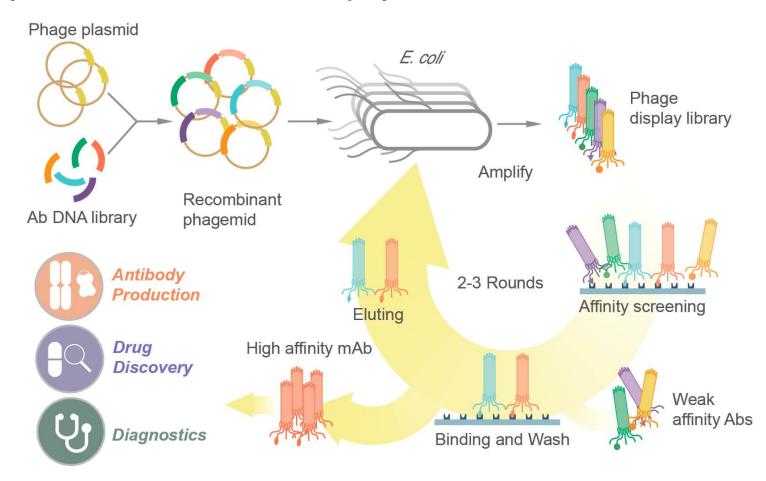
EPFL

Protein display systems

- There are several scalable systems for display and **screening of protein binders** (e.g., antibodies)
- The most common ones are **phage**, **yeast and mammalian cell display**

Phage display:





• Yeast and mammalian display are used in cases of more challenging proteins that cannot be readily produced by bacterial cells (e.g., requiring post-translational modifications)



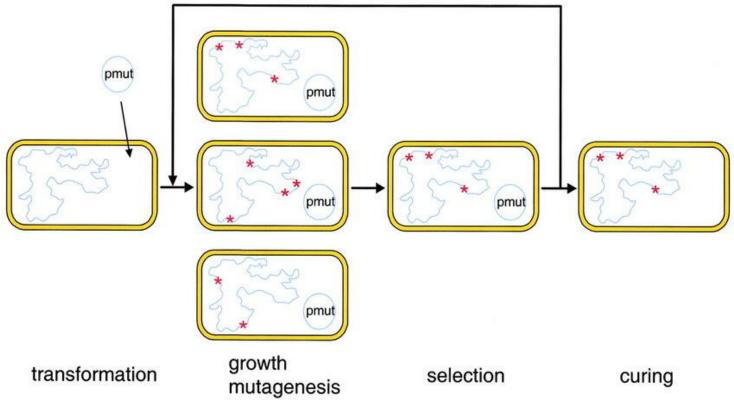
Directed evolution systems

- Directed evolution is a laboratory process by which biological entities with desired traits are created through iterative rounds of genetic diversification and library screening or selection.
- Different display methods fall under the Directed Evolution umbrella, but this general approach also offers to randomly generate new molecule variants through *in vivo* gene diversification, genetic recombination, CRISPR and retroelement based approaches.

Mutators are elements that increase the error rate during gene replication which results in random point mutations

Normal error rates are ~10⁻⁹-10⁻¹⁰ errors per replicated base, while with mutator the rate can be as high as 10⁻⁵

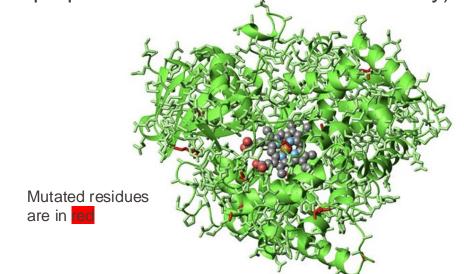


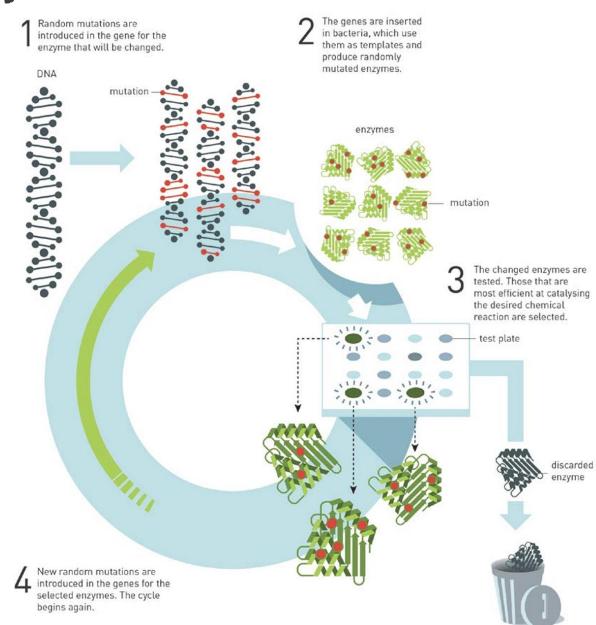




Directed evolution of enzymes

- By accumulation of mutations in the target genes, the resulting proteins change their molecular properties and in some cases exhibit the desired enzymatic activity
- The screening/selection method allow for retention of the desired genetic constructs and removal of the ones that are unfit.
- F. Arnold's P450 variant that catalyzes cyclopropanation reaction (new functionality)







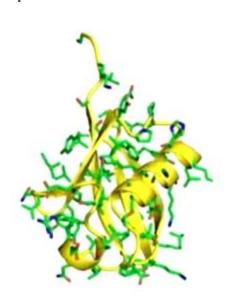
Computational biomolecule design

- Application of computational tools to engineer biomolecules with certain molecular properties (e.g., domain organization, binding to other biomolecules, enzyme specificity etc.)
- Protein design has underwent a quantum leap over the last 20 years.

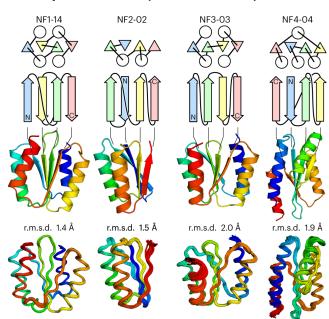


David Baker

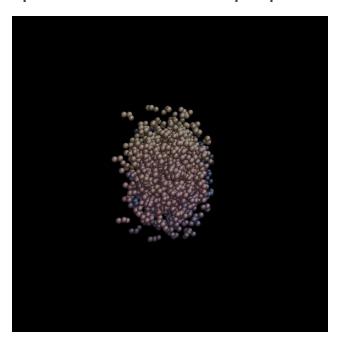
Modifications of existing protein domains



De novo engineering of proteins (new folds)



Al-enhanced design of proteins for desired purposes



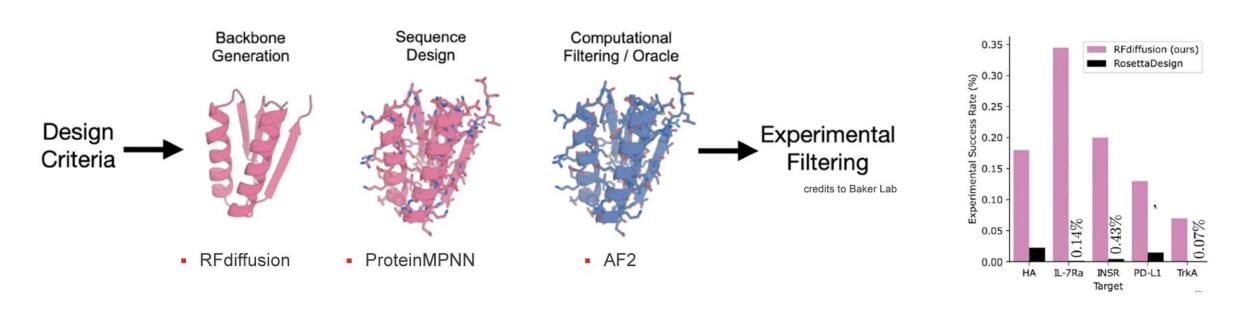
Next-generation methods

Previous generation methods such as Rosetta Design

EPFL

Al methods enhance the experimental rate of success

- Rosetta Design and similar methods incorporated the calculations of energy terms (dG) for different interactions (hydrophobic, hydrogen bonding, electrostatic) that were used to estimate stability and binding.
- However, these energy calculations were challenged by to the complexity of biomolecule interaction networks which lead to relatively low success rate (i.e., there was still a need to screen a library of designed molecules)
- Al methods have been trained on the wealth of available structural data to predict "viable" protein designs, and while they don't provide explicit interaction energy values their success rate is significantly higher.



• Experimental validation is necessary to confirm the functionality of the designed protein, and measure the



Biomolecule design applications

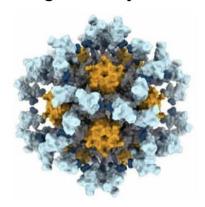
• Computational biomolecule design has many uses in fundamental and applied science

Medicine

vaccines & antivirals

<u>smart medicines</u>

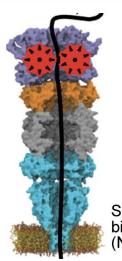
drug delivery



SARS-CoV-2 RBD nanoparticle immunogen (Cell 2020)

Biotechnology

protein-silicon devices
bio-based computers
nanoscale manufacturing

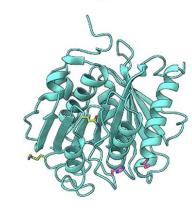


SM proteomics with biological nanopores (Nat Chem 2021)

Sustainability

artificial photosynthesis
CO₂ sequestration

<u>plastic degradation</u>



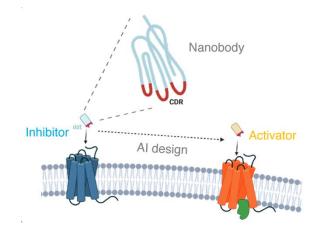
FAST-PETase (Nature 2022)



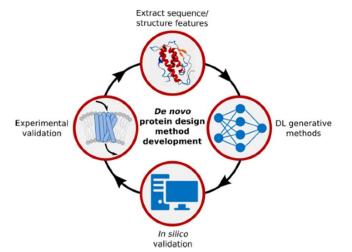
Biomolecule design applications

• Engineering cellular receptors and responses to molecular queues

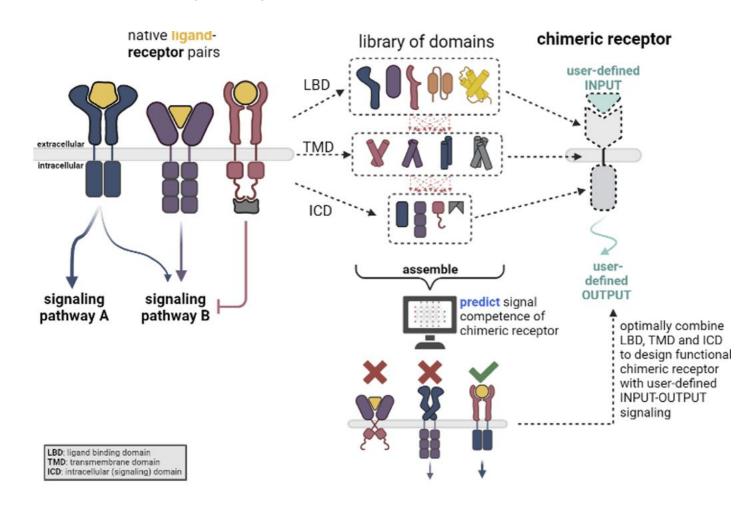
Engineering membrane protein activators



Engineering membrane protein receptors



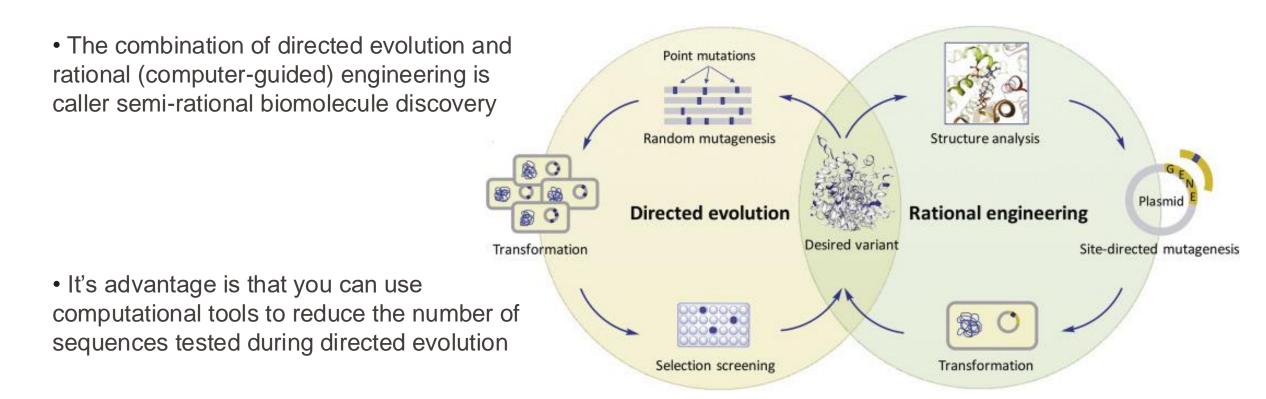
Engineering cellular responses to molecular stimuli





The complementarity of different approaches

- Chemical, biochemical, directed evolution and computational methods for biomolecule engineering are not mutually exclusive and often you have to apply them in combination to reach the desired outcome
- The main consideration is to find the most economical (time- and reagent-wise) approach to create a desired molecule. Each method requires extensive training and access to equipment.





Summary

- Biomolecule engineering allows to move past the natural biochemical diversity and develop molecules with altered or novel functions
- Chemical biology methods use bio-organic chemistry toolkit to generate new chemical species. Bioorthogonal and click-chemistry approaches allow to perform chemical reactions under mild conditions
- Manual manipulations using basic biochemistry and molecular biology tools are commonly used to produce novel biomolecule variants, but are limited in throughput.
- Directed evolution approaches allow to generate and screen very large libraries of biomolecule variants to identify the ones that feature desired properties
- Computational biomolecule design is a powerful tool to engineer proteins and inhibitors. With the advent of Al-enhanced methods the success rate of these methods has drastically improved.



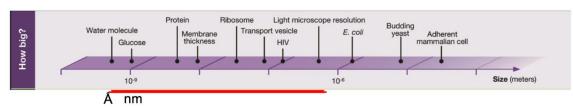
Course Recap



The scales and composition of biomolecules

• Biomolecules on the scales of life:

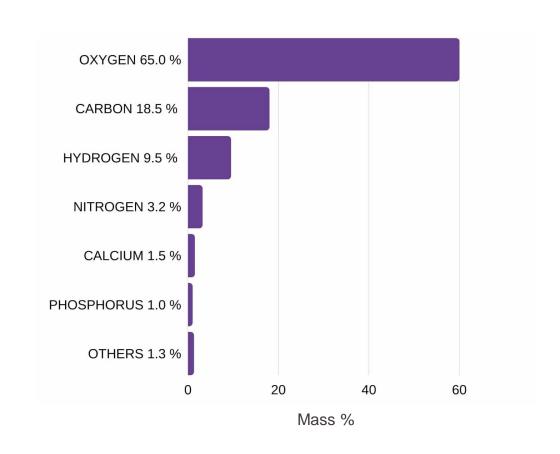
Size



Protein diffusion Step of RNA across E. coli polymerase Molecular motor 1 µm transport across HeLa cell mRNA half life in E. coli E. coli yeast HeLa cell mRNA half life in E. coli yeast HeLa cell mRNA half life in E. coli yeast HeLa cell yeast HeLa cell mRNA half life in E. coli yeast HeLa cell years HeLa cell years

Quantities One molecule in an Signaling proteins in E. coli volume One molecule in an Signaling proteins in E. coli volume Concentration (Molar)

Atomic composition of human bodies by mass:



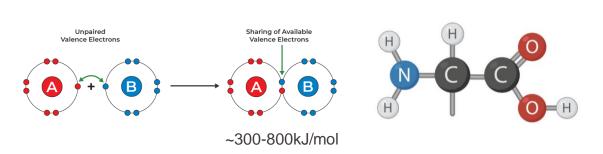


Molecular composition of biological systems

• Atomic and molecular interactions in biomolecules

Covalent bonds

• Molecular composition of human bodies by mass:

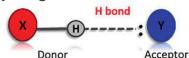


Non-covalent interactions

van der Waals interactions

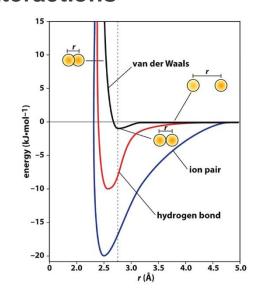


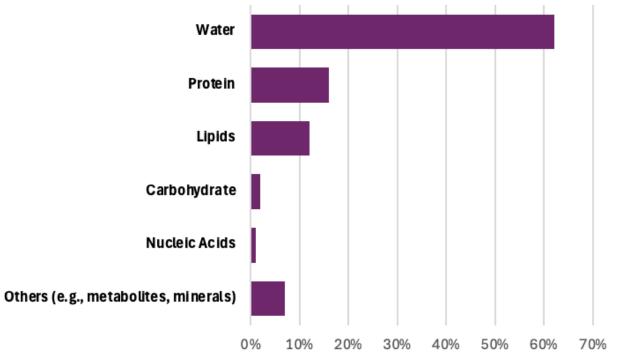
Hydrogen bonds



Ionic interactions



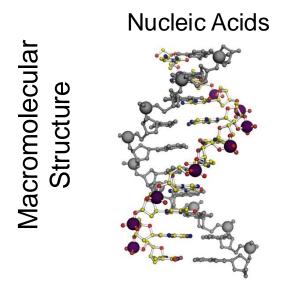




Cell and tissue variation (e.g., bone, adipose cells)



The basic molecules of life

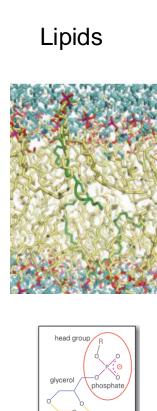


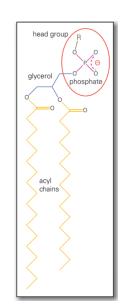
nucleotide = nucleoside + phosphate

NH2
6
N1
9
NH2
6
N1
phosphate

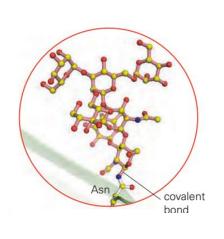
CA'
H
C3'
C2'
H

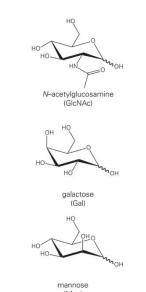
nucleoside = sugar + base

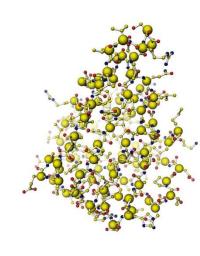




Carbohydrates



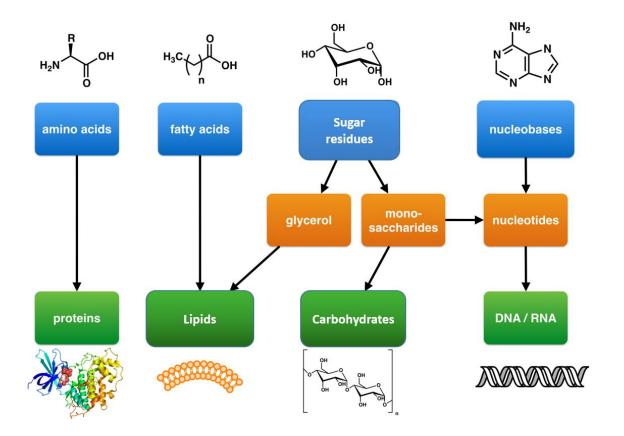




Building Block



How to approach learning about biomolecules



What are the building blocks?

What are the most important covalent bonds?

What are the most important non-covalent interactions?

How do they assemble in 3D space?

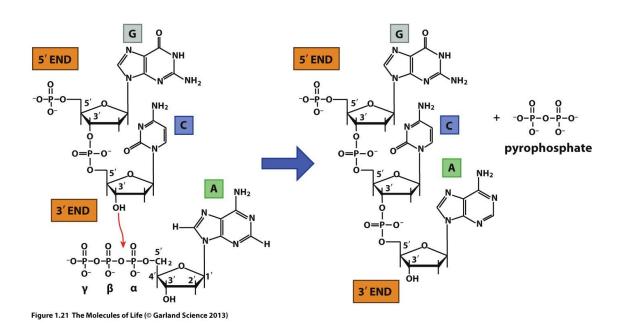
What function do they serve in cells?

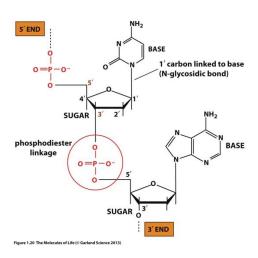
What are some important examples?



Nucleic acids are polymers of nucleotides

- The synthesis of new molecules of DNA and RNA involves the stepwise addition of nucleotide to one end of the chain.
- The triphosphate group is high in energy and its hydrolysis drives the reaction





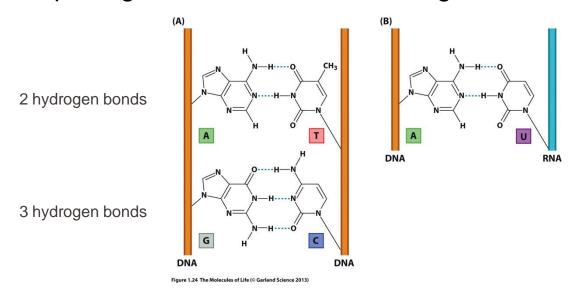
• DNA and RNA synthesis are template directed – DNA polymerases use a template strand to select each nucleotide to be added to the growing chain

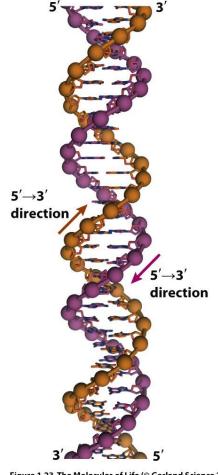
- 3'->5' phosphodiester linkage imposes directionality
- By convention DNA sequences are written from 5' to the 3' end



3D assembly of DNA

- DNA forms a double helix with antiparallel strands
- Two strands together wind up to form a right-handed double-helix
- Bases are on the inside of the helix and the phosphate backbone group are on the outside. Allowing for interactions with ions and water and minimizing repulsion between phosphates
- Base pairing holds the DNA strands together and is strictly complementary





igure 1.23 The Molecules of Life (© Garland Science 2

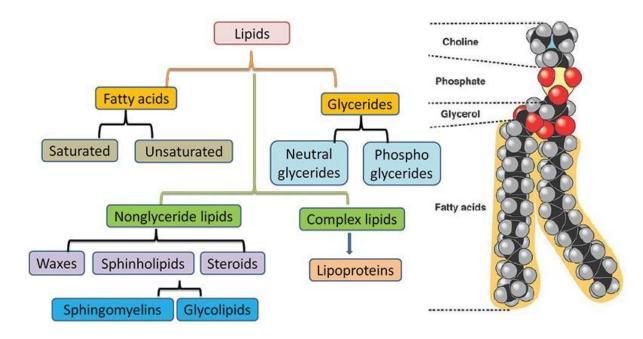
Phosphate groups in spheres



Lipids are diverse molecules with amphipathic properties

Main lipid types and roles

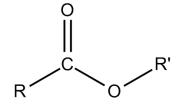
- Main lipid types based on chemical properties



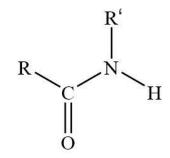
- There are >1000 building blocks which can assemble in different ways (very diverse)
- Their main roles include energy storage, assembly of biological membranes, cell and hormone signaling

Important bonds and interactions

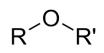
Ester bonds (triglycerides)

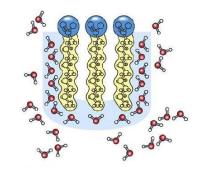


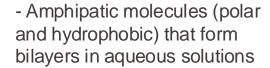
Amide bonds (sphingolipids)

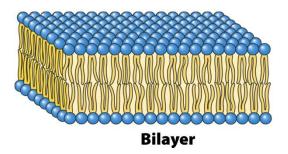


Ether bonds (head domains)









- Van der Waals interactions in the tail

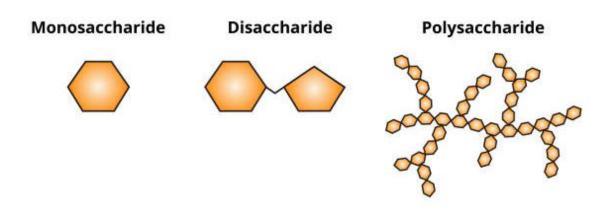
 Hydrogen bonding and charged interactions using the head domain

EPFL Carbohydrates

- Carbohydrates, glycans, sugars or saccharides are a diverse group of biomolecules that represent 2-10% of dry matter in humans and as high as 60-90% in plants (depending on the source).
- They serve many roles in cells:
 - energy metabolism and storage (e.g., glucose, glycogen, and amylose)
 - markers of cellular identity (e.g., glycolipids and glycoproteins),
 - structural components (e.g., cellulose in plants),
 - constituents of nucleotides (e.g., ribose in RNA, deoxyribose in DNA)



- General Formula: C_nH_{2n}O_n (due to the 1:1 ratio of C to H₂O, the name **carbohydrates** was proposed)
- Like nucleic acids they also form polymers from smaller building blocks (monosaccharides)





Categories of larger carbohydrate chains

- Disaccharide: 2 monosaccharides linked by O-glycosidic bond (e.g. sucrose lactose)
- Oligosaccharides: 3-12 monosaccharides linked by O- or N-glycosidic bonds (e.g. glycoprotein or glycolipid moieties)

ÓН

• Polysaccharides: very large number of linked monosaccharides (e.g. starch, cellulose)

Disaccharide (Maltose) CH2OH OH OH OH OH OH

Oligosaccharide (Oligofructose)

l _{300–600}

CH₂OH CH₂OH CH₂OH OH OH

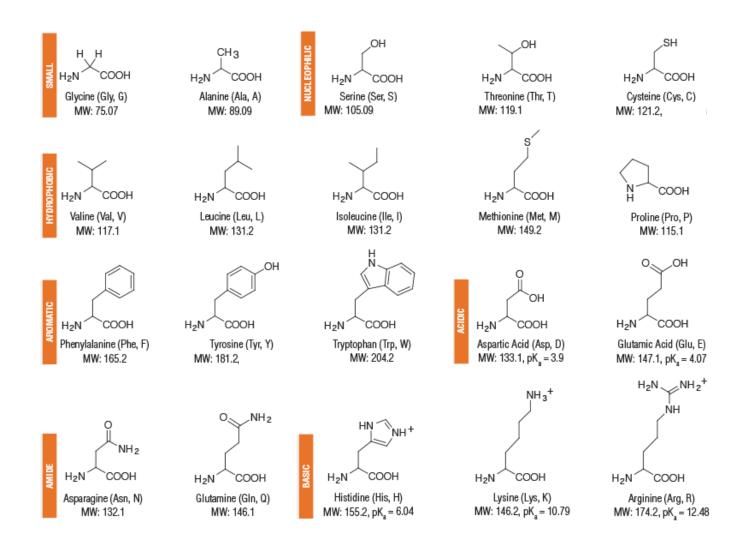
Branched polysaccharide (Amylopectin)

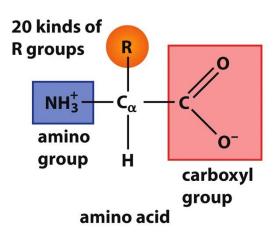
Linear polysaccharide (Amylose)

Components of starch

EPFL

Amino acids are the building blocks of proteins



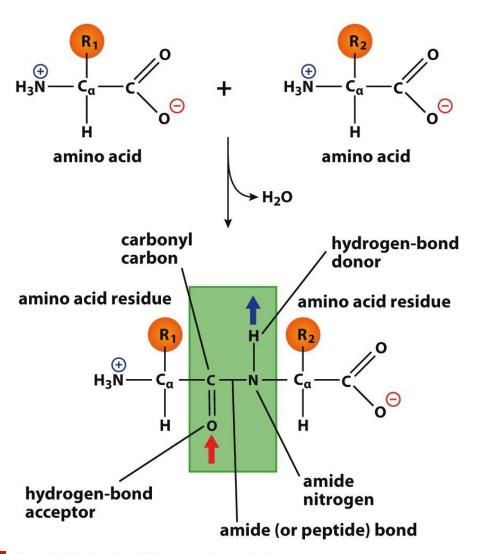


- There are 20 primary amino acids that assemble all proteins
- Amino acids have the same scaffold but differ only in the side chain group (R)
- Note 3- and 1- letter abbreviations for each amino acid

EPFL

Proteins are polymers of amino acids

• Proteins are assembled from amino acids through translation of genetic material



• The synthesis of proteins involves a **condensation** reaction in which the amino group of one amino acid combines with the carboxyl of another.

 This reaction forms a peptide bond and the elimination of a water molecule

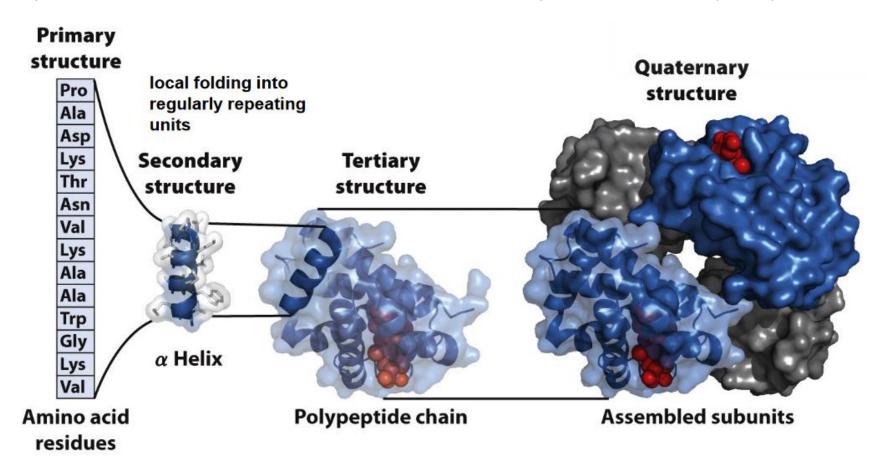
• This reaction is catalyzed by a large assembly of proteins and RNA called the **ribosome**

• The end product is a **polypeptide chain** where amino acids are added sequentially (no branching)



Hierarchical organization of protein structure

• From polypeptide chain sequence to full protein assembly into functional (multi)molecular complexes

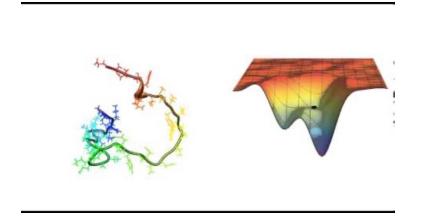


• Protein folding is guided by hydrogen bonding, hydrophobic, VDW and electrostatic interactions

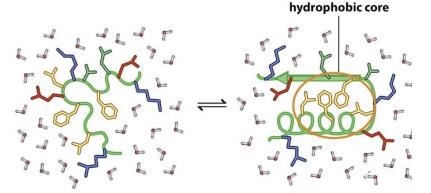


Thermodynamics of protein folding

- Thermodynamics of protein folding
- Proteins folding is about reaching the state of minimum free energy

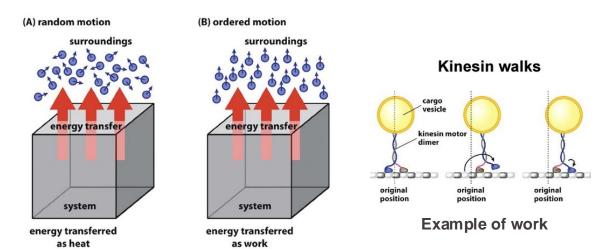


- Hydrophobic core formation as an essential driver



Work and heat in the context of biological systems

- Energy released by chemical reactions is converted into heat and work



- Free energy, internal energy, enthalpy and enthropy

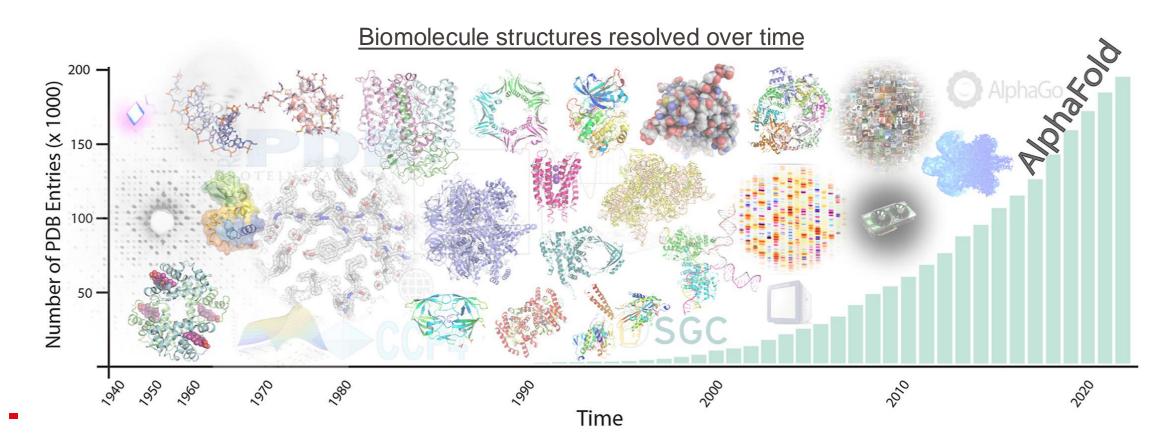
$$G = H - TS$$

- Heat capacity of macromolecules changes depending on the conformational (energy) state



Structural biology - A subfield of Biochemistry

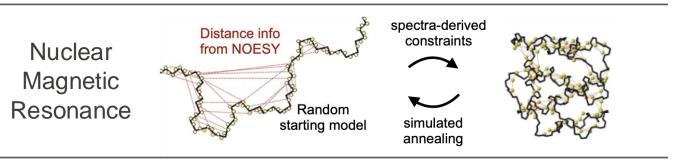
- Structural biology is the study of molecular structure and dynamics of biological macromolecules, particularly proteins and nucleic acids, and how alterations in their structure affect their function.
- Total number of available experimentally determined structures to date is ~225k
- Computational tools (e.g., AlphaFold) can predict macromolecular structures *in silico* allowing to expand the available structure space



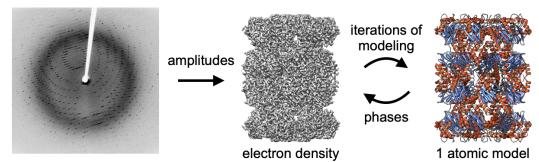


The main structural biology methods

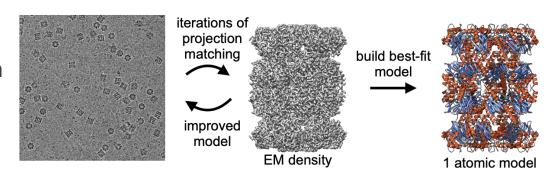
• 3 main methods for experimental structure determination:



X-ray Crystallography



Cryo-electron Microscopy



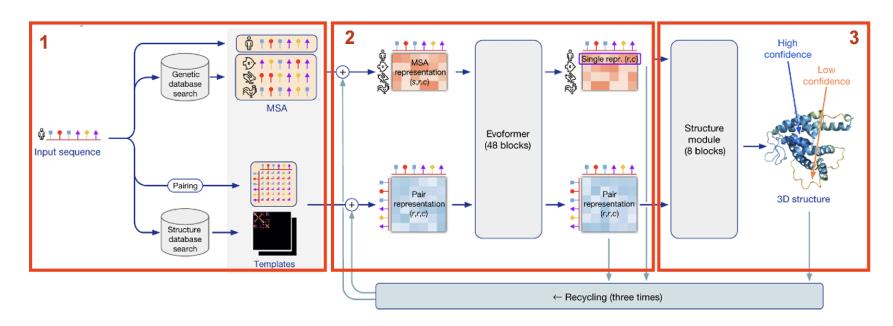
Molecular Type ↓↑	X-ray↓ 	EM↓↑	NMR↓↑
Protein (only)	166,790	15,369	12,516
Protein/Oligosaccharide	9,624	2,600	34
Protein/NA	8,710	4,654	286
Nucleic acid (only)	2,867	137	1,507
Other	170	10	33
Oligosaccharide (only)	11	0	6
Total	188,172	22,770	14,382

- Most structures came from X-ray crystallography since it is historically the oldest method
- The outputs of computational (AI) prediction of biomolecular structures are usually not referred to as "structures" but rather as "models"
- This is because they have not been experimentally verified



Protein structure prediction using AlphaFold2

Basic architecture of the program:



First module:

gather available information like sequence similarity (MSA) and structural templates from databases (UniProt/ metagenomics and PDB) to create a pair representation (which aa are likely in contact with each other)

Second module:

Evoformer transformer which refine the MSA and pair interactions

Third module:

The structure module build the 3D structure based on the MSA and pair interactions information

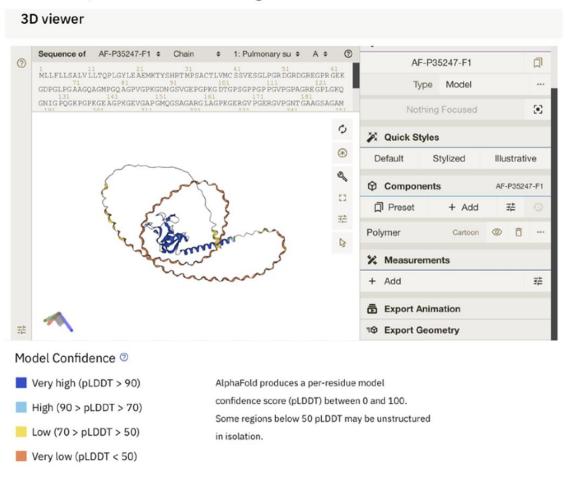
- end-to-end model produces prediction in one shot
- recycling (3X) to refine further prediction
- huge engineering effort



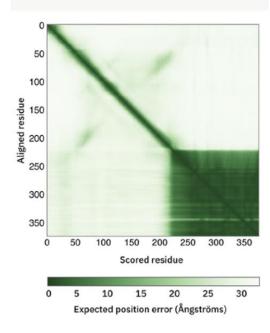
Protein structure prediction using AlphaFold2

• Different scoring metrics:

PAE: predicted alignment error



Predicted aligned error (PAE)



Click and drag a box on the PAE viewer to select regions of the structure and highlight them on the 3D viewer.

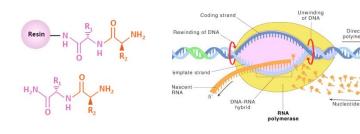
PAE data is useful for assessing inter-domain accuracy – go to Help section below for more information.

- pLDDT is a good score only at short/local distances
- it cannot give you good estimation of the quality of a prediction with different domains
- their reciprocal orientation cannot be estimated by pLDDT
- PAE is the solution for this scenario



Biomolecule production and purification

• Biomolecule production methods



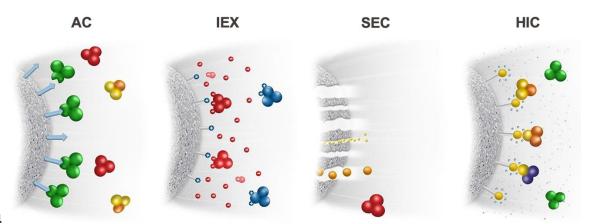


Chemical

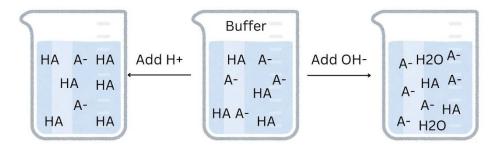
Enzymatic

Cell-based

• Biomolecule purification - Liquid Chromatography

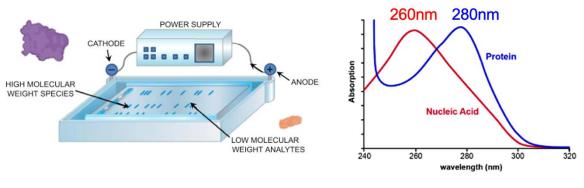


Buffers and buffer components



- Maintaining pH and ionic strength of the solution
- Other components can be added for LC or stability

Evaluating purity and quantity

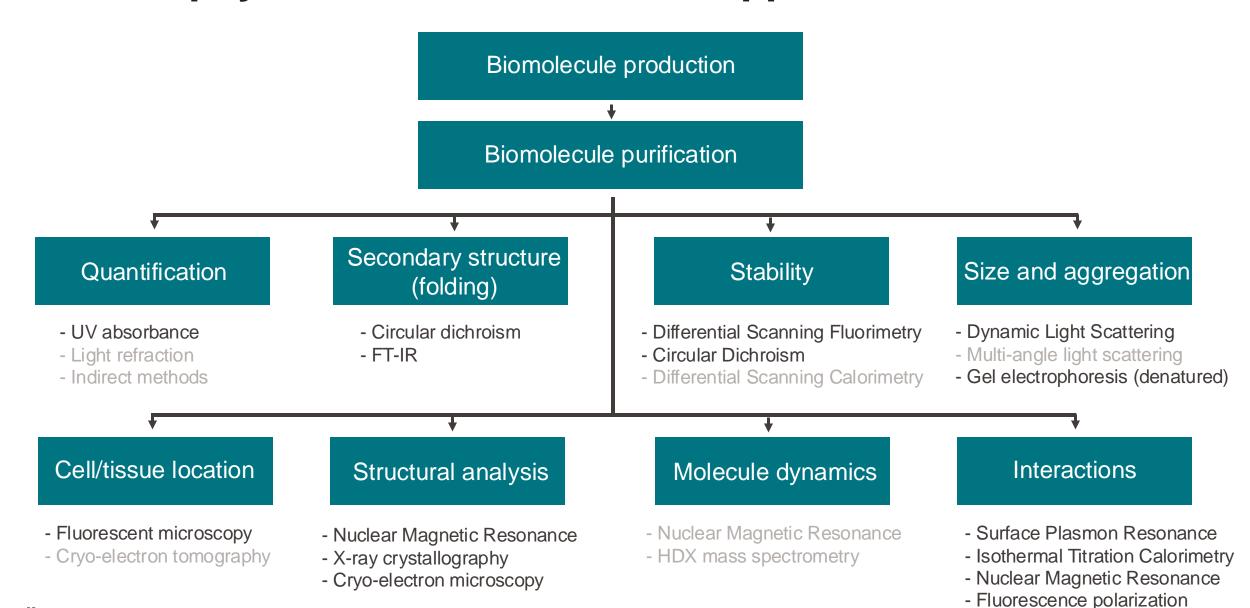


Gel electrophoresis

UV absorbance



Biophysical methods and their applications

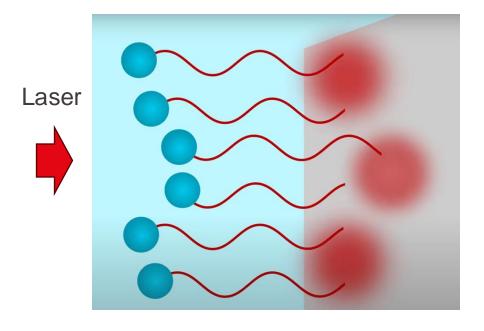




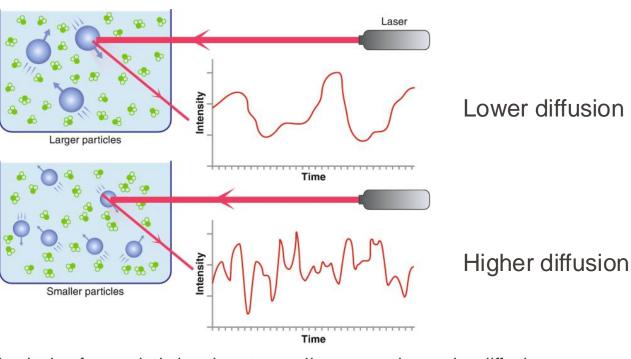
Biophysical methods - What to focus on?

- Primary questions:
 - How do different biophysical methods work?
 - What are they used for?
 - Do they have any limitations in their applicability?

Example: Dynamic Light Scattering



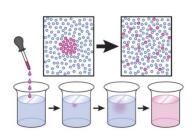
Fluctuations in signal intensity correlated to molecular size and movement (diffusion)



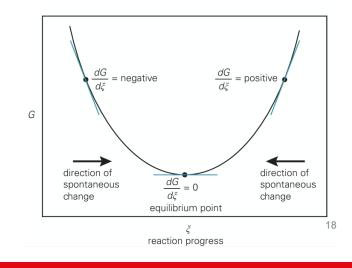
Analysis of recorded signal patterns allows to estimate the diffusion coefficient of the underlying molecules at a given temperature

Energetics of molecular interactions

Gibbs free energy and equilibrium



$$dG = dH - TdS$$



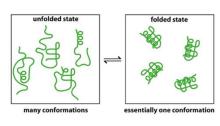
Equilibrium in different contexts

Chemical reaction

$$ATP + H_2O \rightarrow ADP + P_i$$

$$K = \frac{[ADP][P_i]}{[ATP]}$$

Protein Folding



$$K_{\text{folding}} = \frac{[F]}{[U]}$$

Acid-Base Eq.

$$HA \rightleftharpoons H^+ + A^-$$

$$K_{a} = \frac{[H^{+}][A^{-}]}{[HA]}$$

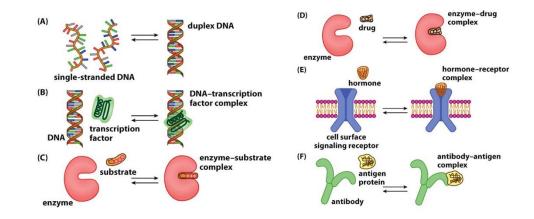
Equilibrium constants

$$v_{A}A + v_{B}B \rightleftharpoons v_{C}C + v_{D}D$$

$$K_{\text{eq}} = \frac{[C]_{\text{eq}}^{v_{\text{C}}} [D]_{\text{eq}}^{v_{\text{D}}}}{[A]_{\text{eq}}^{v_{\text{A}}} [B]_{\text{eq}}^{v_{\text{B}}}} \quad \Delta G^{\text{o}} = -RT \ln K_{\text{eq}}$$

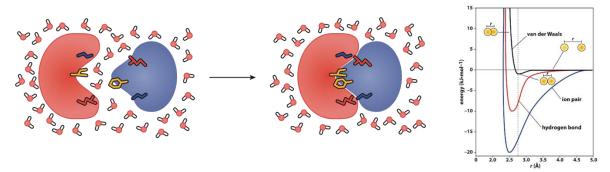
Reaction quotient at chemical equilibrium

Molecular interactions



Quantification of biomolecular interactions

Biomolecular interactions and binding



- Small energy contributions from hydrogen bonds, vdW, ionic and hydrophobic interactions
- Water can have a positive and negative impact on binding
- Dissociation constants and affinity

$\Delta G^{\circ} = RT \ln K_D = \Delta H - T\Delta S$

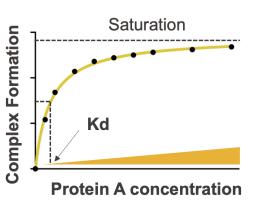
$$\frac{\mathbf{k}_{d}}{\mathbf{k}_{a}} = \frac{k_{off}}{k_{on}} = \frac{[P][L]}{[P \bullet L]} = K_{D} = \frac{1}{K_{A}}$$

- Thermodynamic and kinetic analysis of binding

Experimental methods for Kd determination

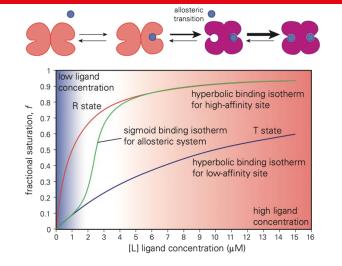
- Titrate one binding partner and measure complex formation

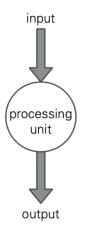




- Calorimetric (ITC) or spectroscopic (FP, SPR, NMR) measurements

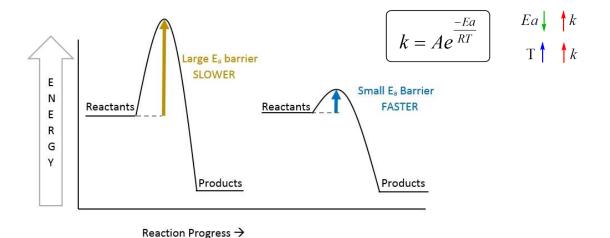
Cooperativity and Allostery





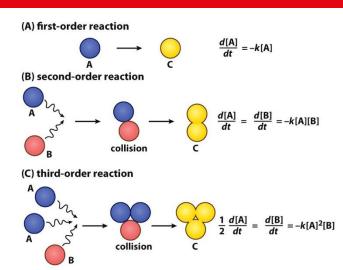
Enzymes and kinetics

Reaction rates and Activation Energies

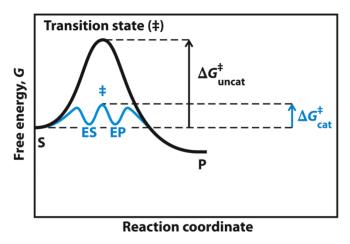


Reaction order

- Order of the reaction depends on the number of molecules whose quantities impact the reaction rate
- Rate constant units depend on the order



Catalysis and enzymes



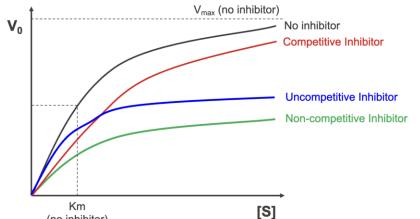
$$E + S \rightleftharpoons ES \rightleftharpoons EP \rightleftharpoons E + P$$

$$V_0 = \frac{V_{\text{max}}[S]}{[S] + K_m}$$

- Michaelis-Menten equation

• Enzyme inhibition mechanisms

(no inhibitor)



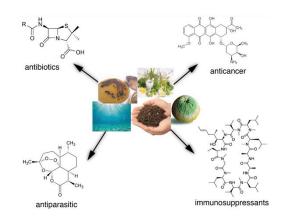
- Inhibition depends on the K_I and the amount of inhibitor
- Often expressed as α

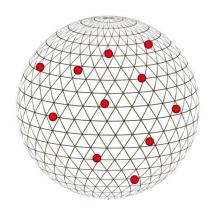
$$\alpha = 1 + \frac{[I]}{k_I}$$



Biomolecule engineering

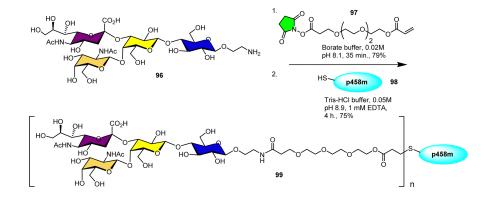
Natural diversity of biomolecules



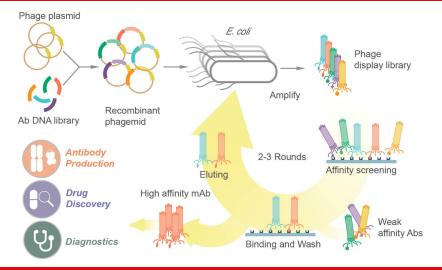


Space sampled by nature

Chemical biology to engineer novel function



Library based methods (Directed Evolution)



Computational tools for engineering proteins

