



Structural Biology - BIO315

**Master SV - Spring Semester
Lecture 9**

Matteo Dal Peraro

matteo.dalperaro@epfl.ch

AAB 048, phone: 31681

Outline of lecture 9:

- **multiscale simulations**
- **integrative modeling**
- **structure-based drug design**
- **protein design**

molecular modeling and simulations

$$\{x_i(t), y_i(t), z_i(t)\}_{i=1,\dots,N}$$

solvation

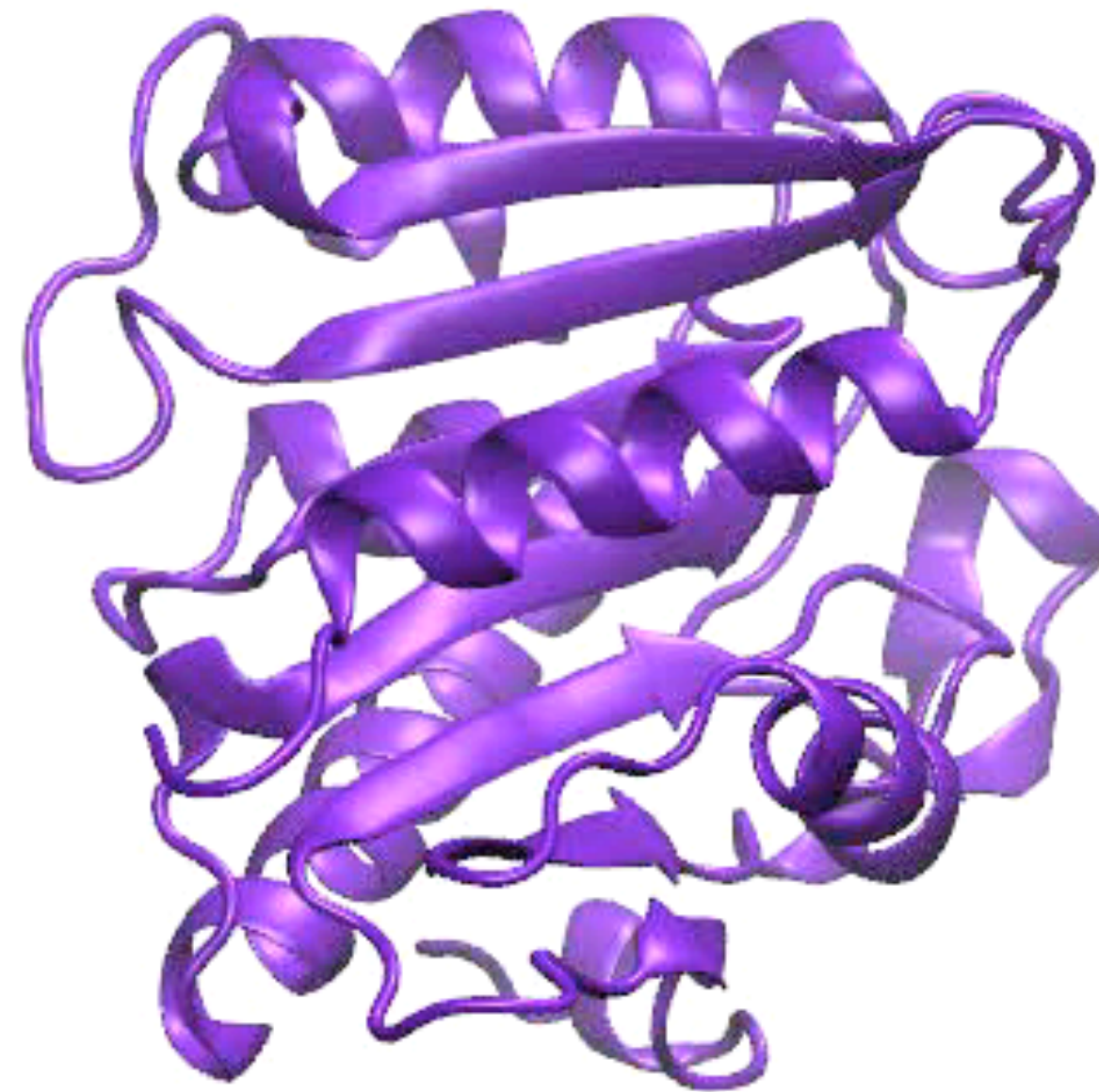
pH

post-translational modifications

interactions network

temperature effects ($k_B T$)

.....

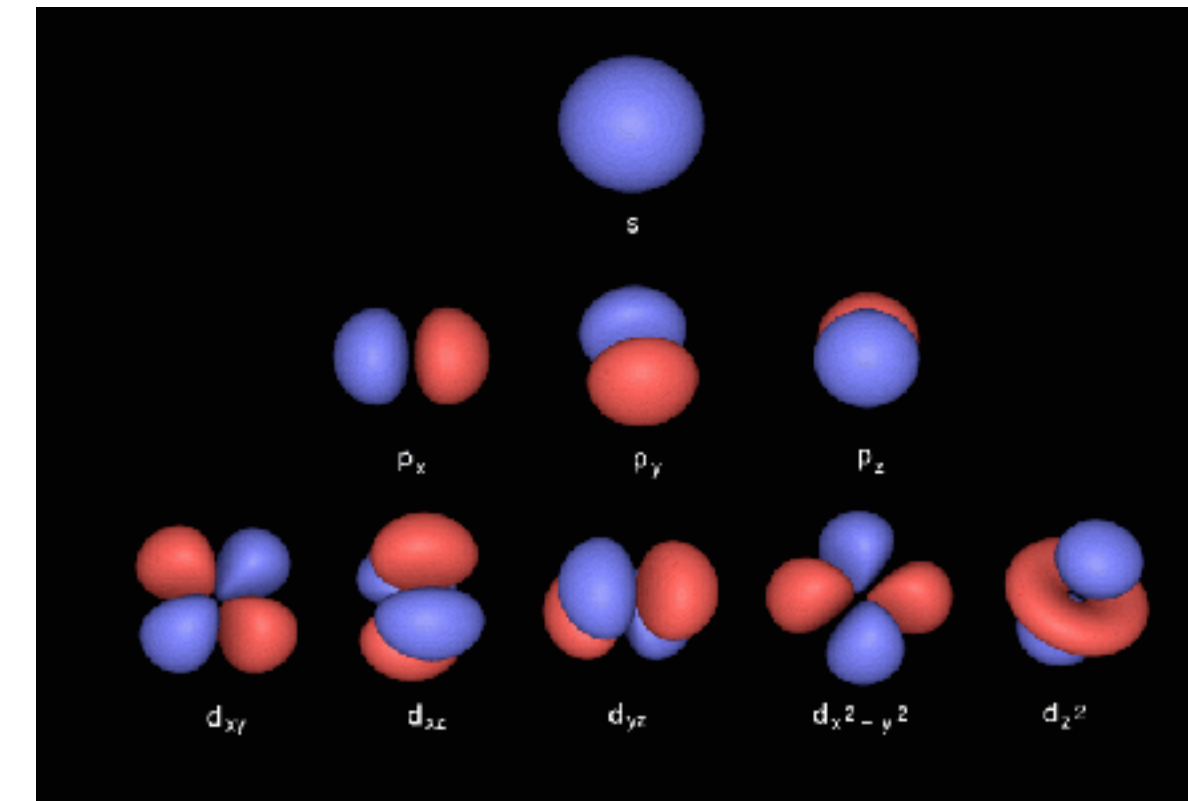
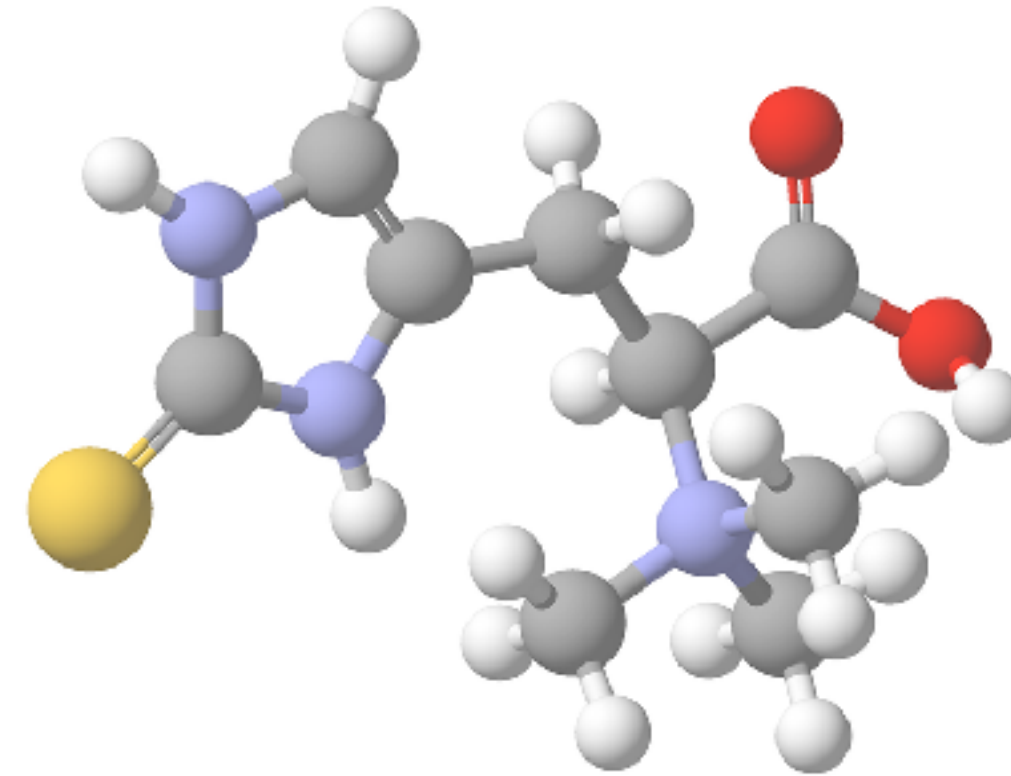


Current common MD engines

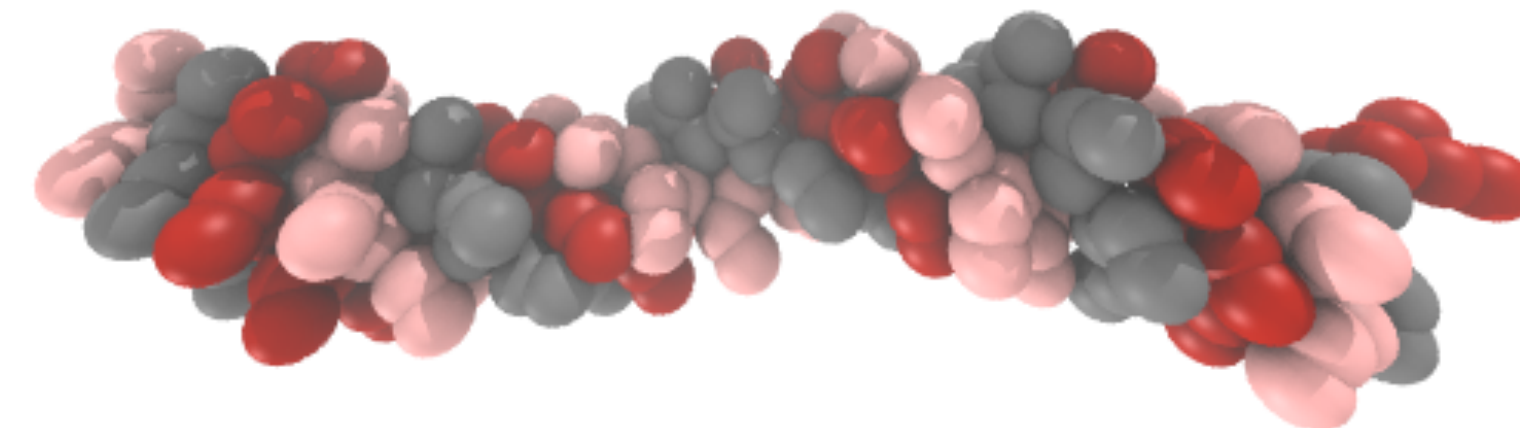
- **CHARMM**: Karplus Harvard, <http://www.charmm.org/>
- **AMBER**: Kollman UCSF, <http://ambermd.org/>
- **GROMOS**: van Gunsteren, ETHZ, www.igc.ethz.ch/GROMOS/index
- **DESMOND**: Shaw, <http://www.deshawresearch.com/>
- **GROMACS**: <http://www.gromacs.org>
- **LAMMPS**: <http://lammps.sandia.gov>
- **ACEMD**: <http://multiscalelab.org/acemd>
- **NAMD**: <http://www.ks.uiuc.edu/Research/namd/>

Multiscale resolution in modeling

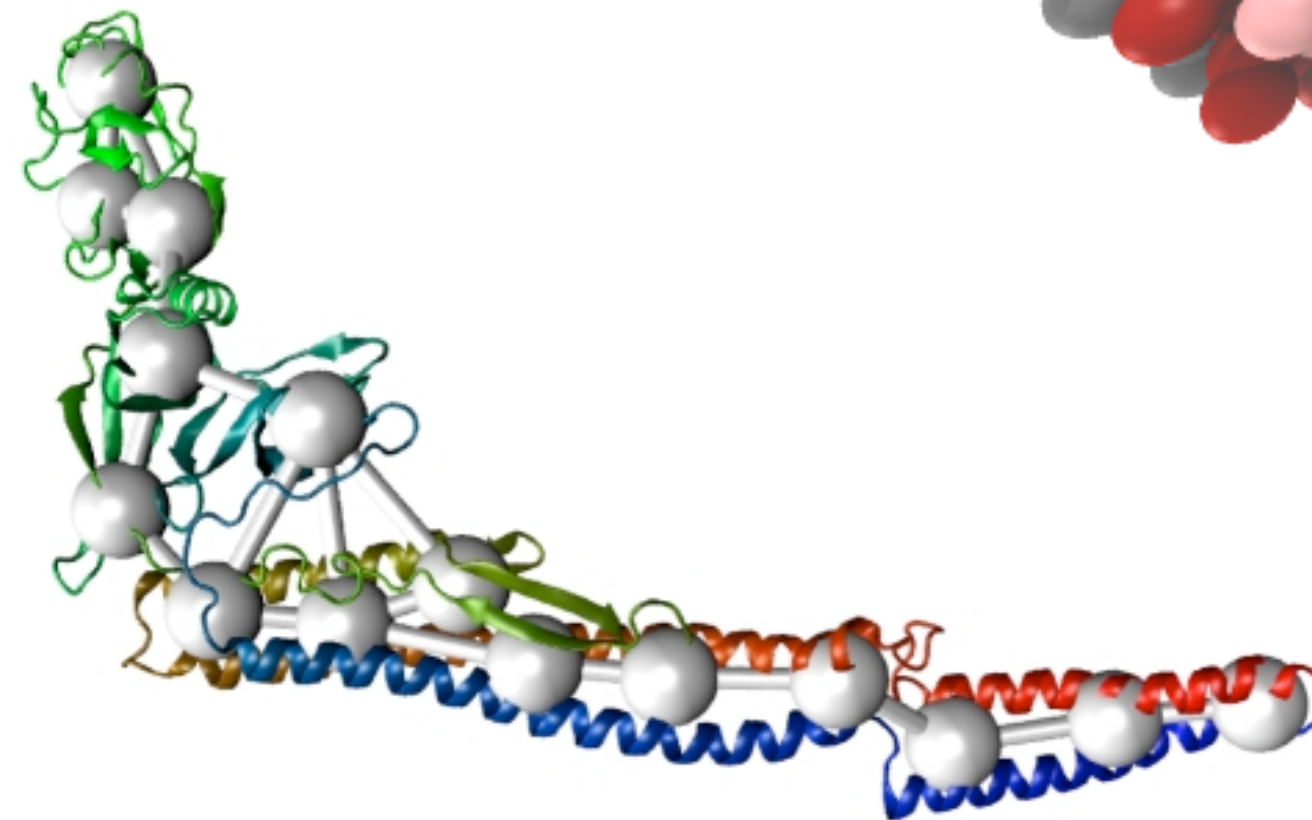
- electrons



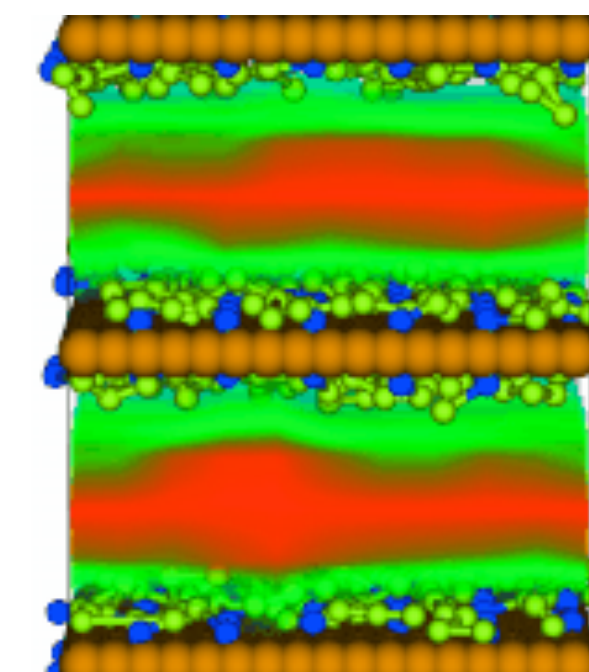
- atoms



- amino-acids

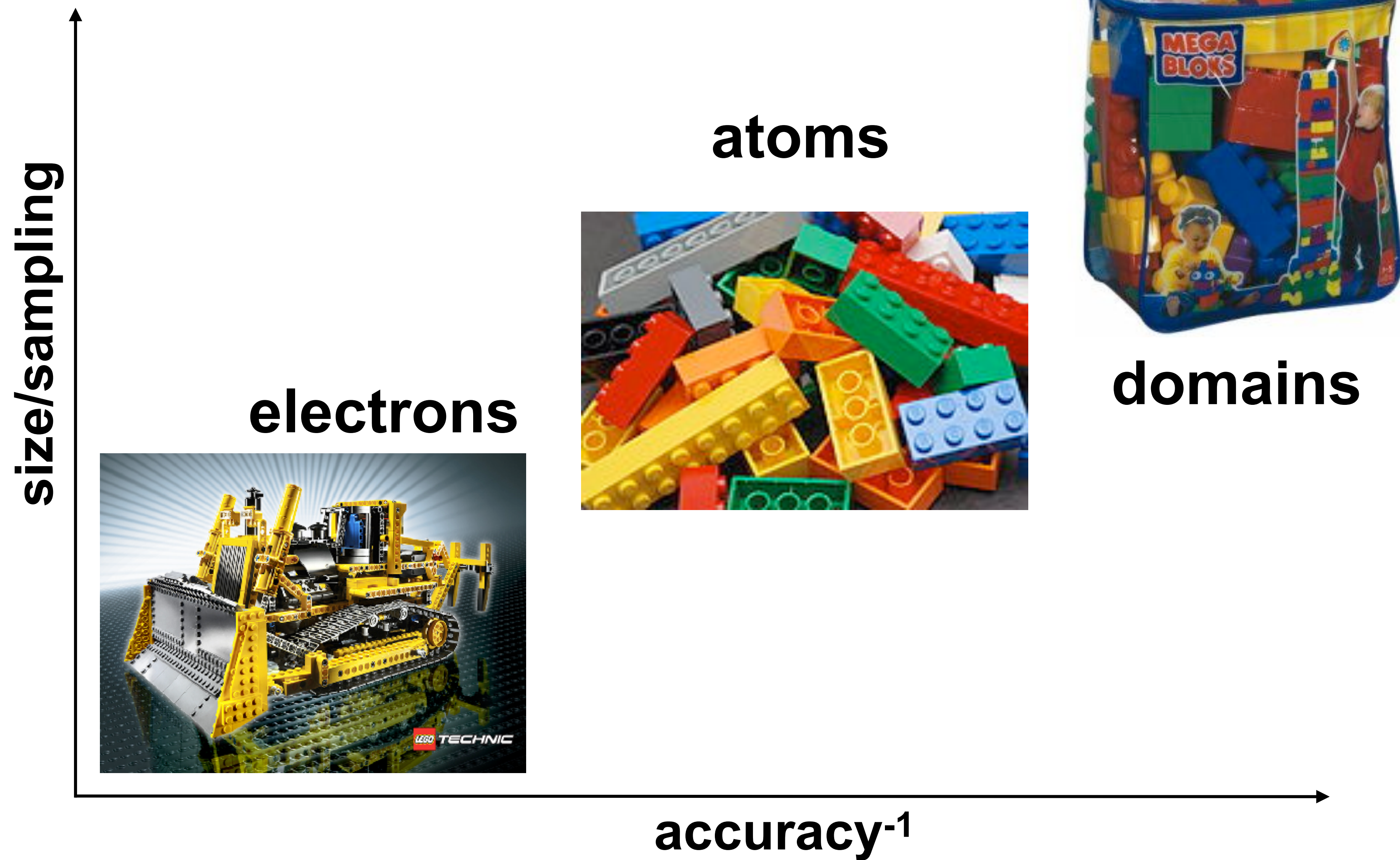


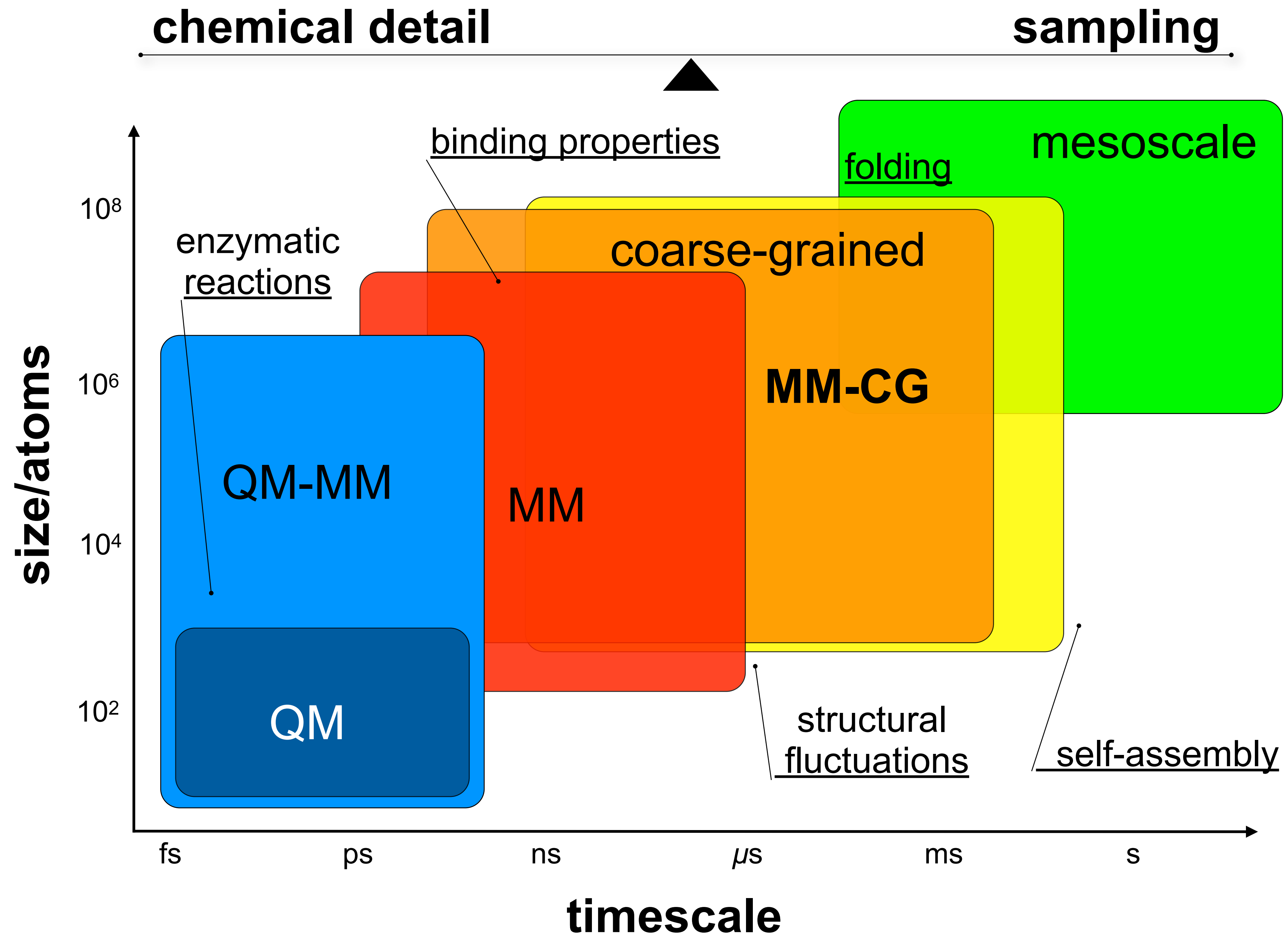
- domains



- mesoscopic to continuum

Building blocks





Speeding up timescales of Chemical Reactions

- **Enzymes** enhance the rate of chemical reactions by several orders of magnitude (e.g. arginine decarboxylase, alkaline phosphatase, staphylococcal nuclease **up to 10^{14} fold**)
- the transition rate depends on the activation barrier

$$\Gamma_{\text{reactants} \rightarrow \text{products}} \propto e^{-G_{\text{barrier}}/k_B T}$$

- and enzymes affect this, not the R and P states

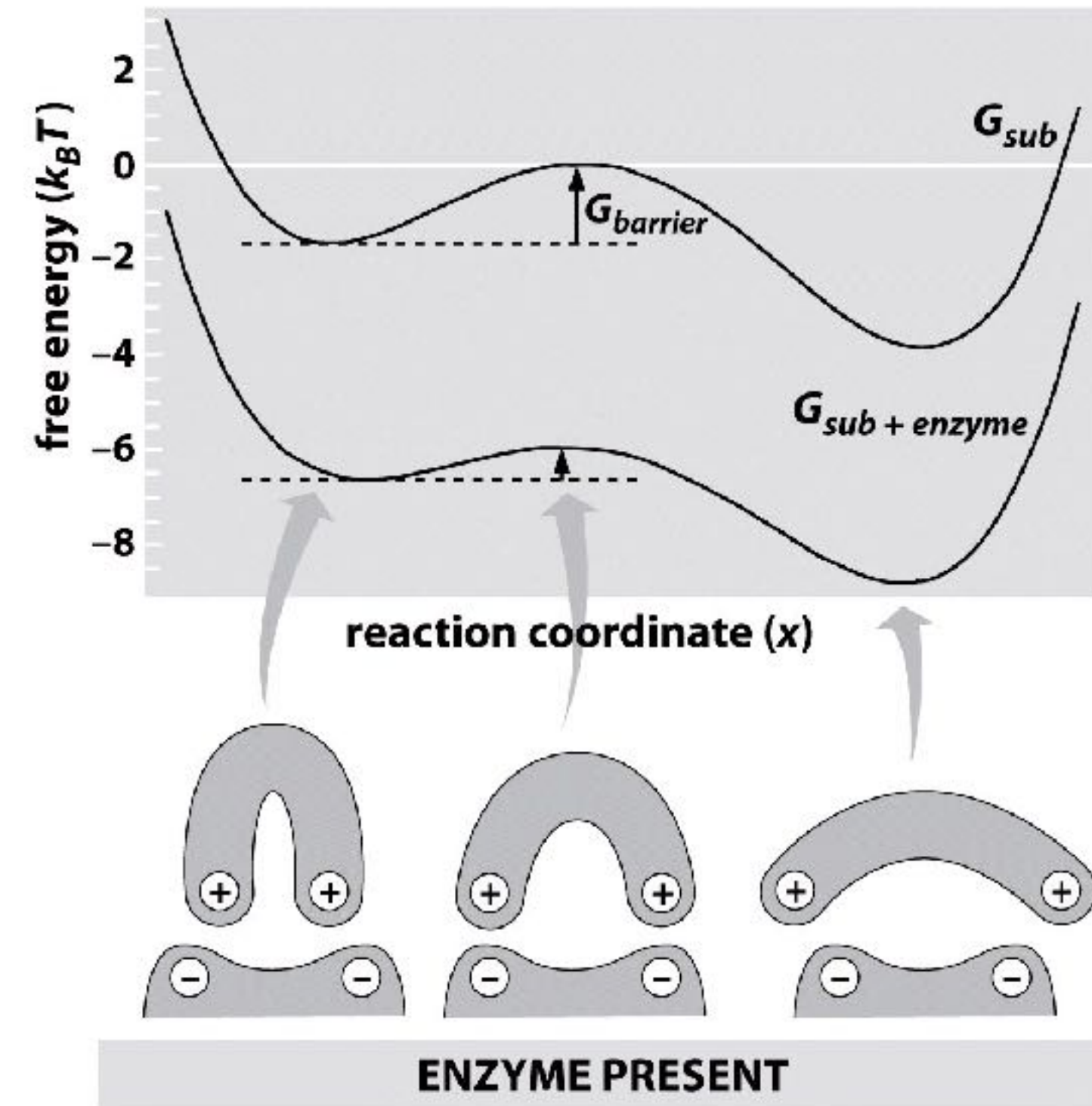
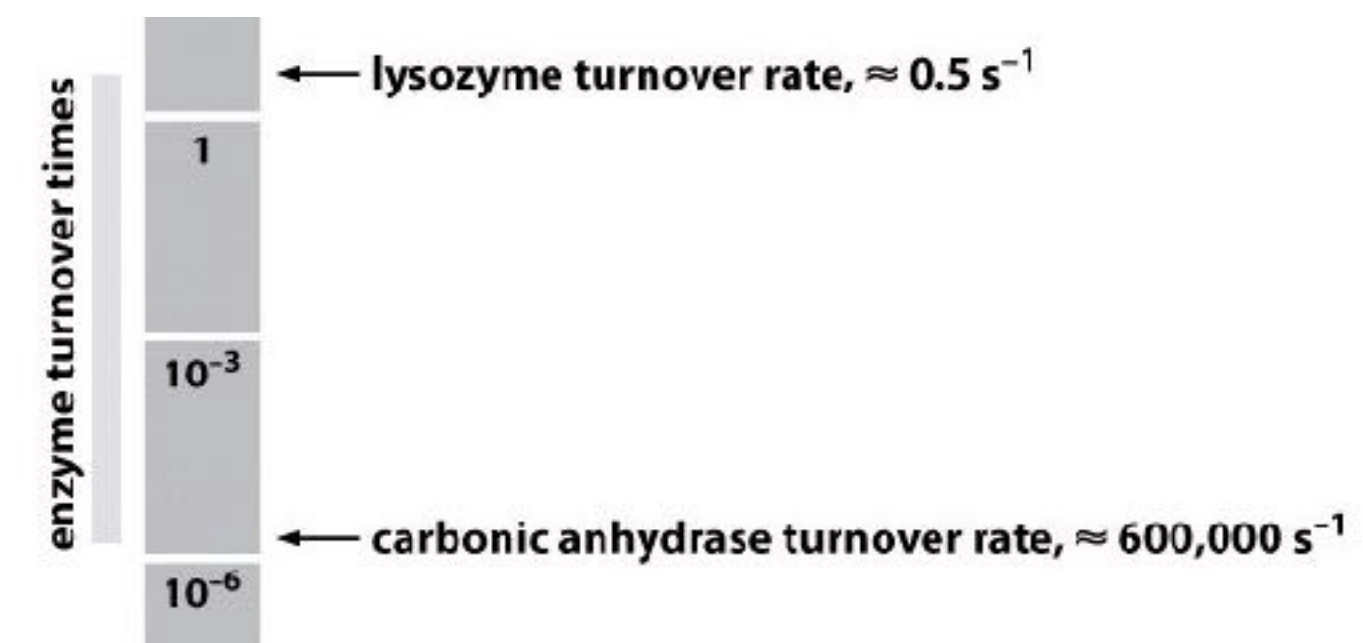


Figure 3.24b Physical Biology of the Cell (© Garland Science 2009)

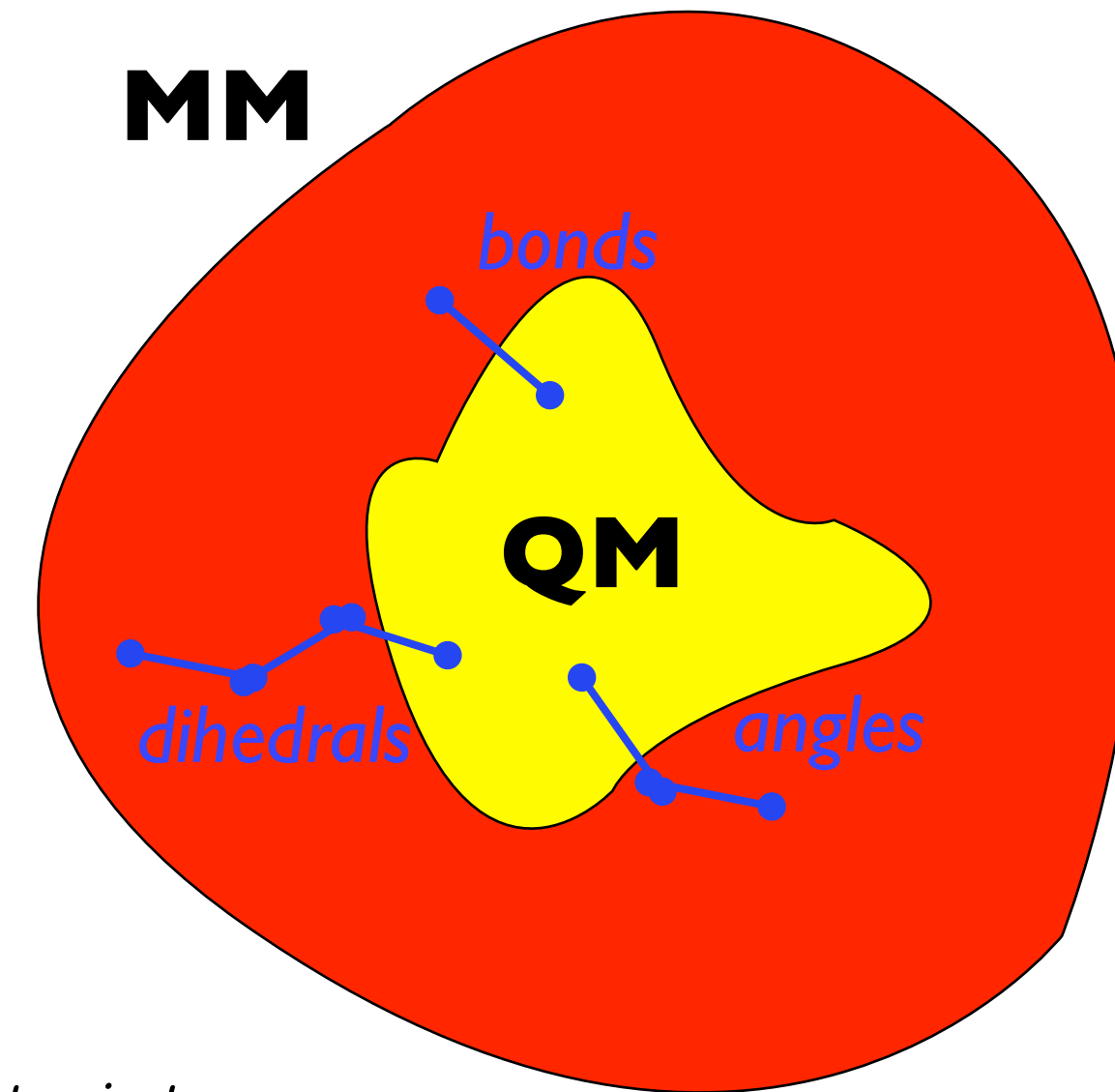


Hybrid QM/MM molecular dynamics

$$H = H_{QM} + H_{MM} + \underbrace{H_{QM/MM}}_{\text{coupling term}}$$

QM: First principles Density functional theory MD

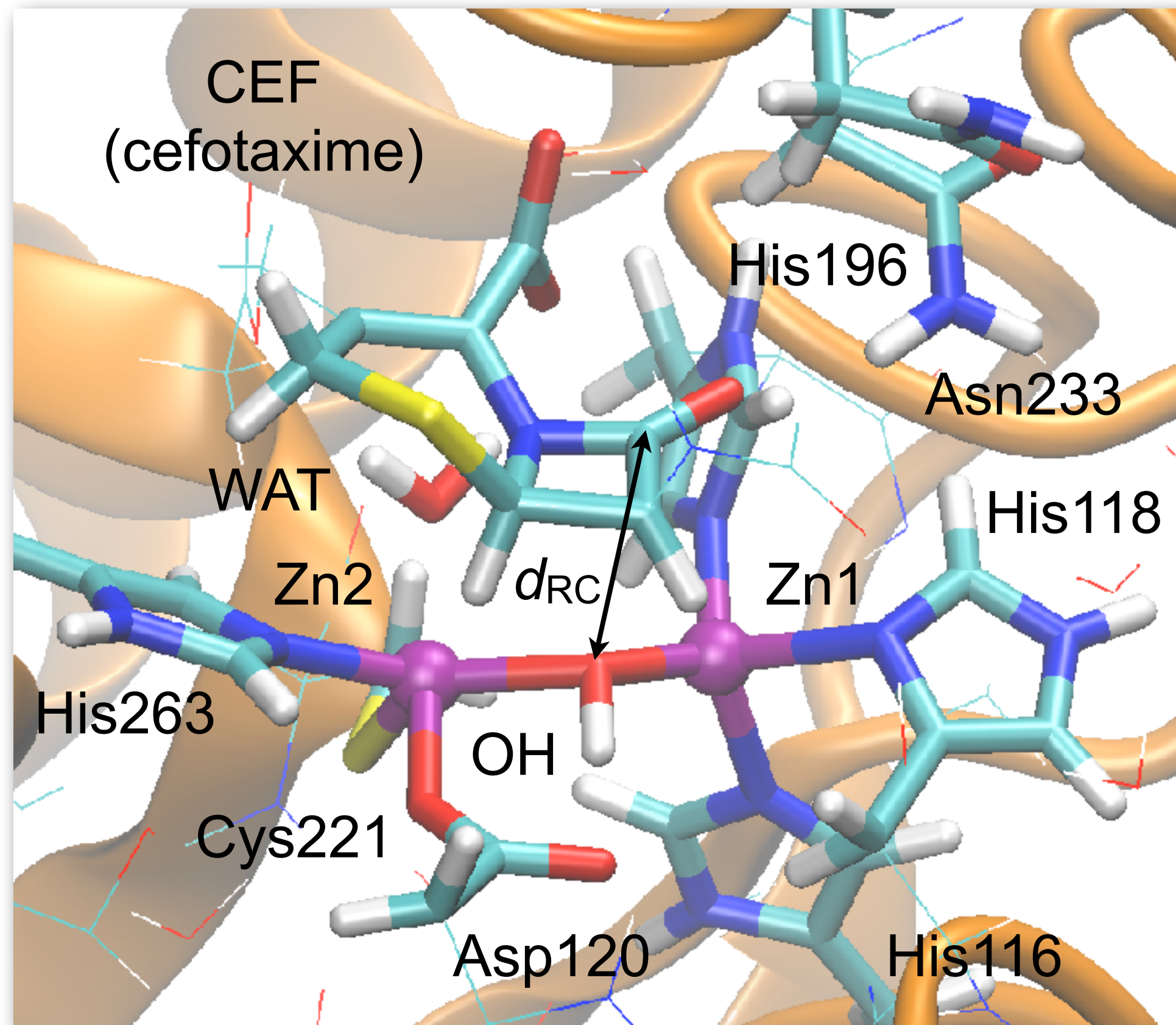
$$\mathcal{L}_{CP} = \underbrace{\sum_I \frac{1}{2} M_I \dot{\mathbf{R}}_I^2 + \sum_i \frac{1}{2} \mu_i \langle \dot{\psi}_i | \dot{\psi}_i \rangle}_{\text{kinetic energy}} - \underbrace{\langle \Psi_0 | \mathcal{H}_e | \Psi_0 \rangle}_{\text{potential energy}} + \underbrace{\text{constraints}}_{\text{orthonormality}}$$



MM: Classical molecular dynamics (e.g. AMBER, Gromos force fields)

QM/MM: - boundary atom (*ad hoc* monovalent pseudopotential or H capping)
- hierarchical scheme to compute Coulomb interactions

CcrA M β L from *Bacteroides fragilis*



Thermodynamic integration along the reaction coordinate d_{RC}
DFT-BLYP, Martins-Troullier PPs, 70 Ry cutoff,
Nose' thermostat at 300 K,
2 reactions pathways for a total of ~150 ps trajectory

Reactant state

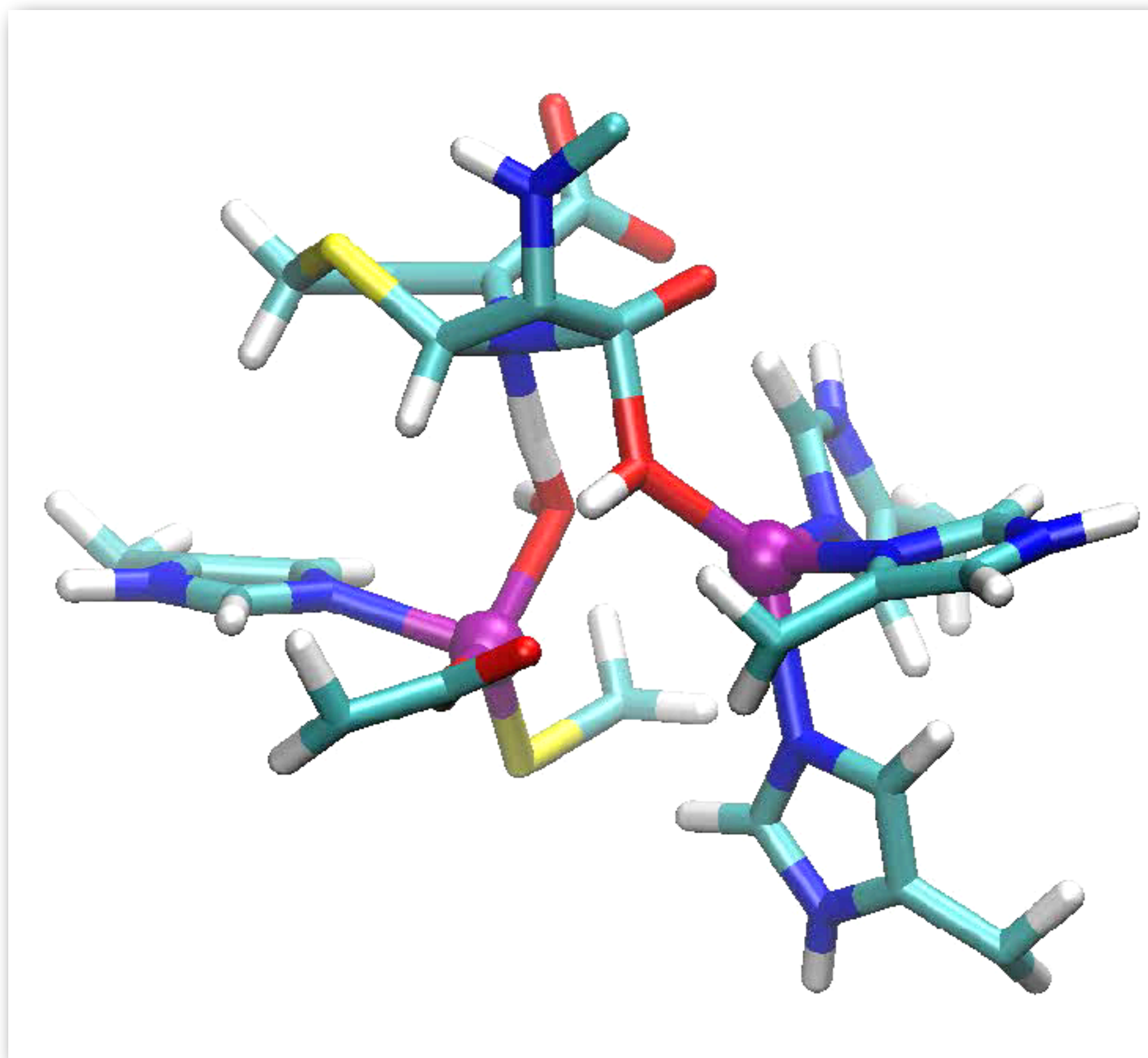
CcrA complexed with cefotaxime

- stable Michaelis complex
OH- β -lactam distance=3.3(2) \AA
during 5 ns MD and 20ps QM/MM

- Zn2-bound WAT is the only water between the zinc center and CEF in 5 \AA

➡ Classical force-field based MD is used as a tool to sample conformational space within the nanosecond timescale

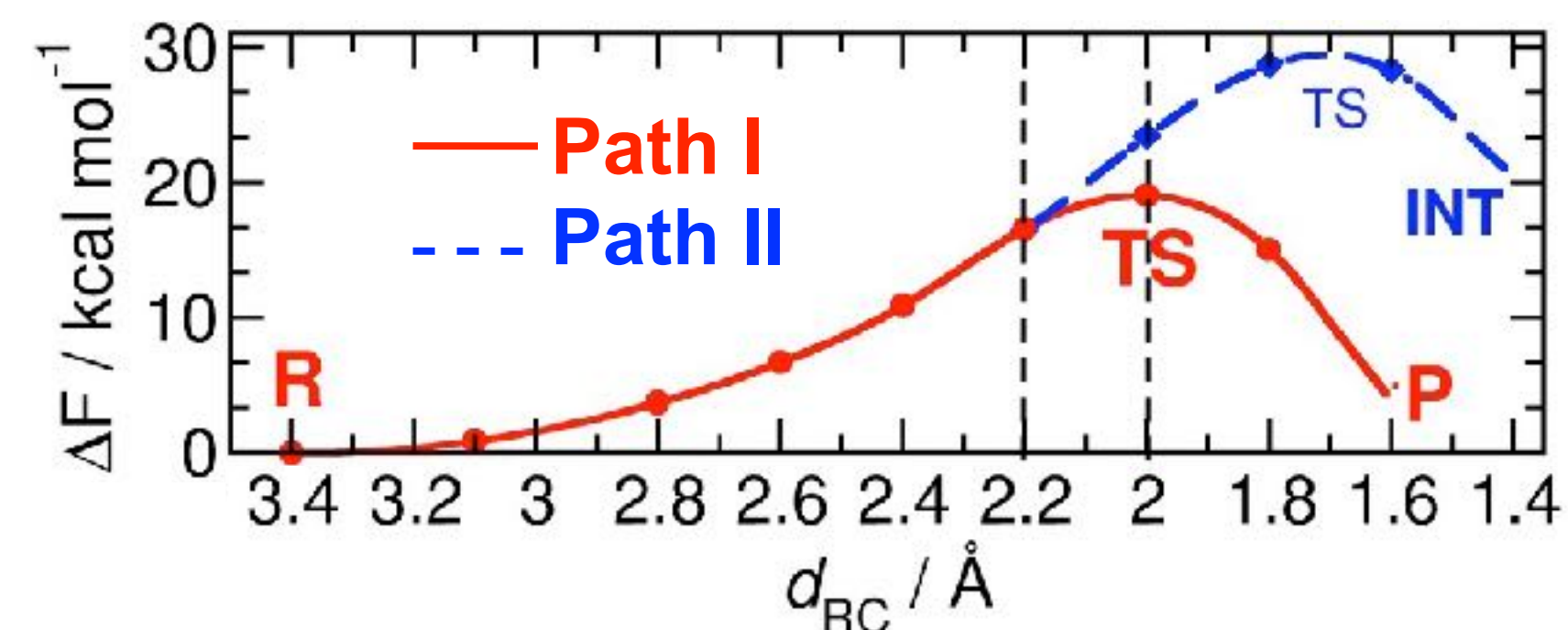
... from transition state to products

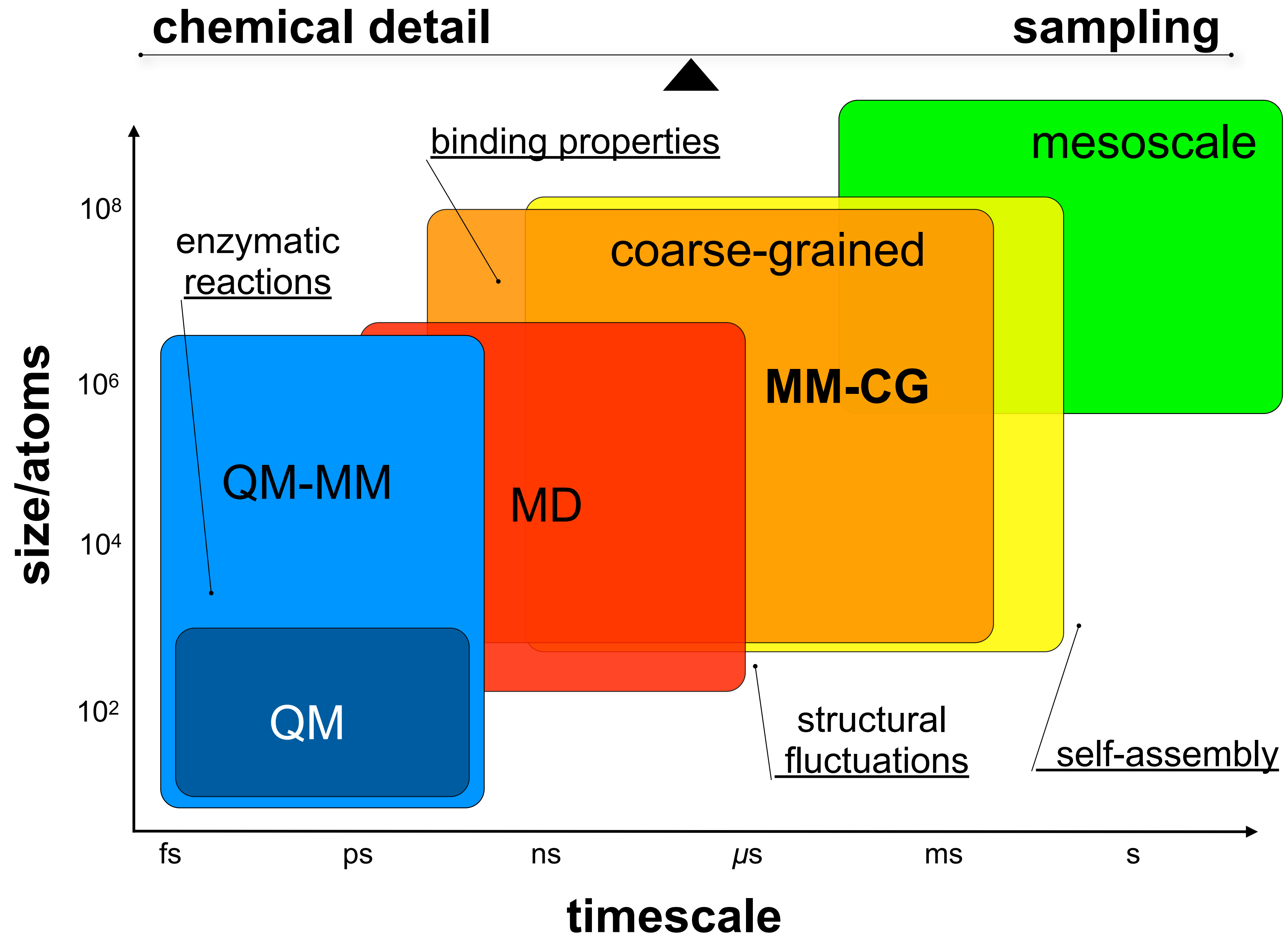


water-mediated single-step

- OH⁻ loses Zn² coordination
- Zn1, Zn2 flexibility
- WAT protonates β-lactam N
- N-C β-lactam bond breaks
- WAT replaces OH⁻ as an hydroxide

- **$\Delta F = 18(2)$ kcal/mol** is in good agreement with experiments
- if Asn233 *does* H-bond β-lactam: formation of a high unfavorable intermediate (Path II)





Coarse-graining degrees of freedom

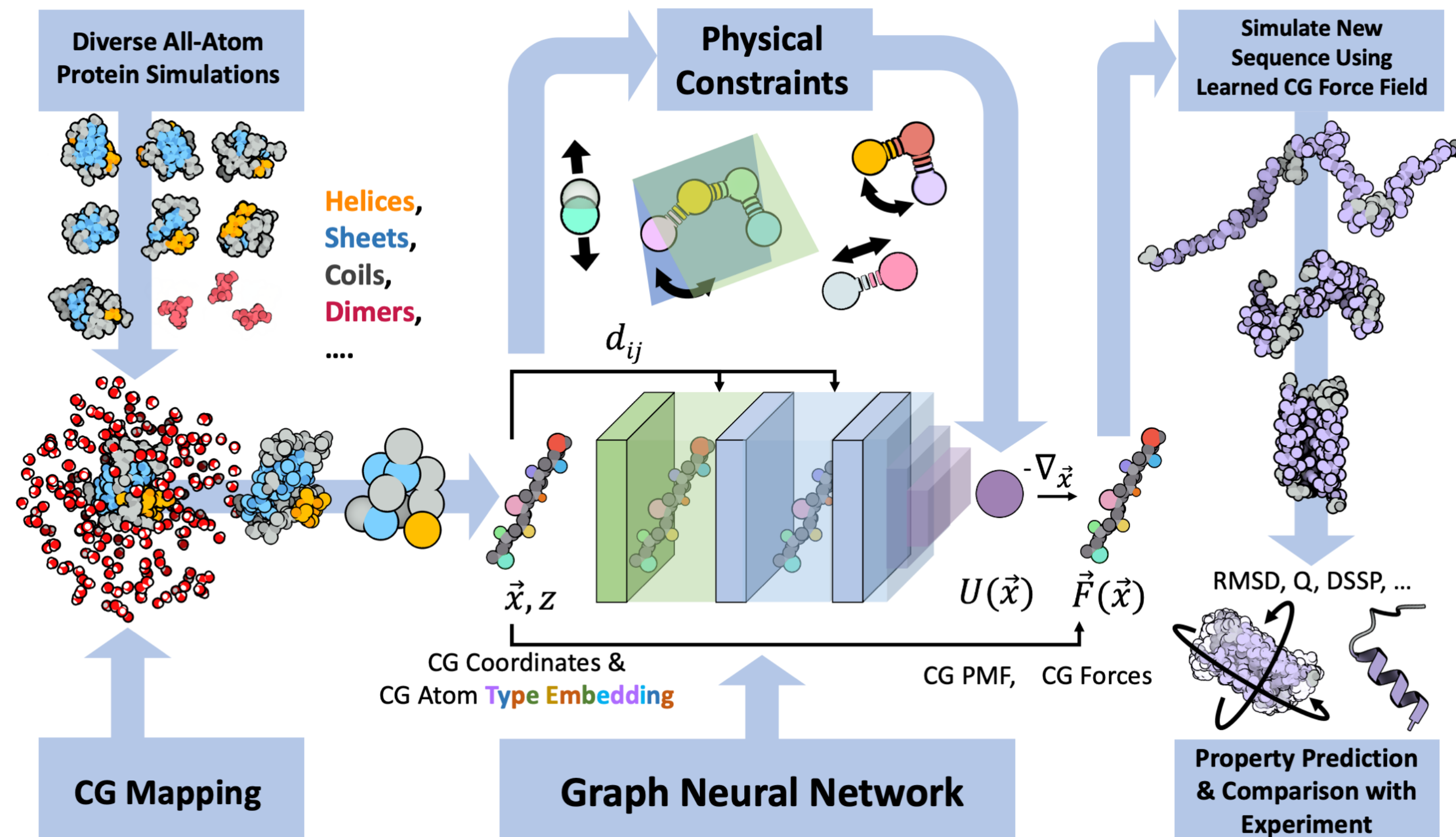
- **CG** is the process of consistently reduce the complexity of your problem integrating out degrees of freedom which can be in principle neglected for your system.

$$V_{QM} \rightarrow V_{MM} \rightarrow V_{CG-MM} \rightarrow V_{mesoscopic}$$

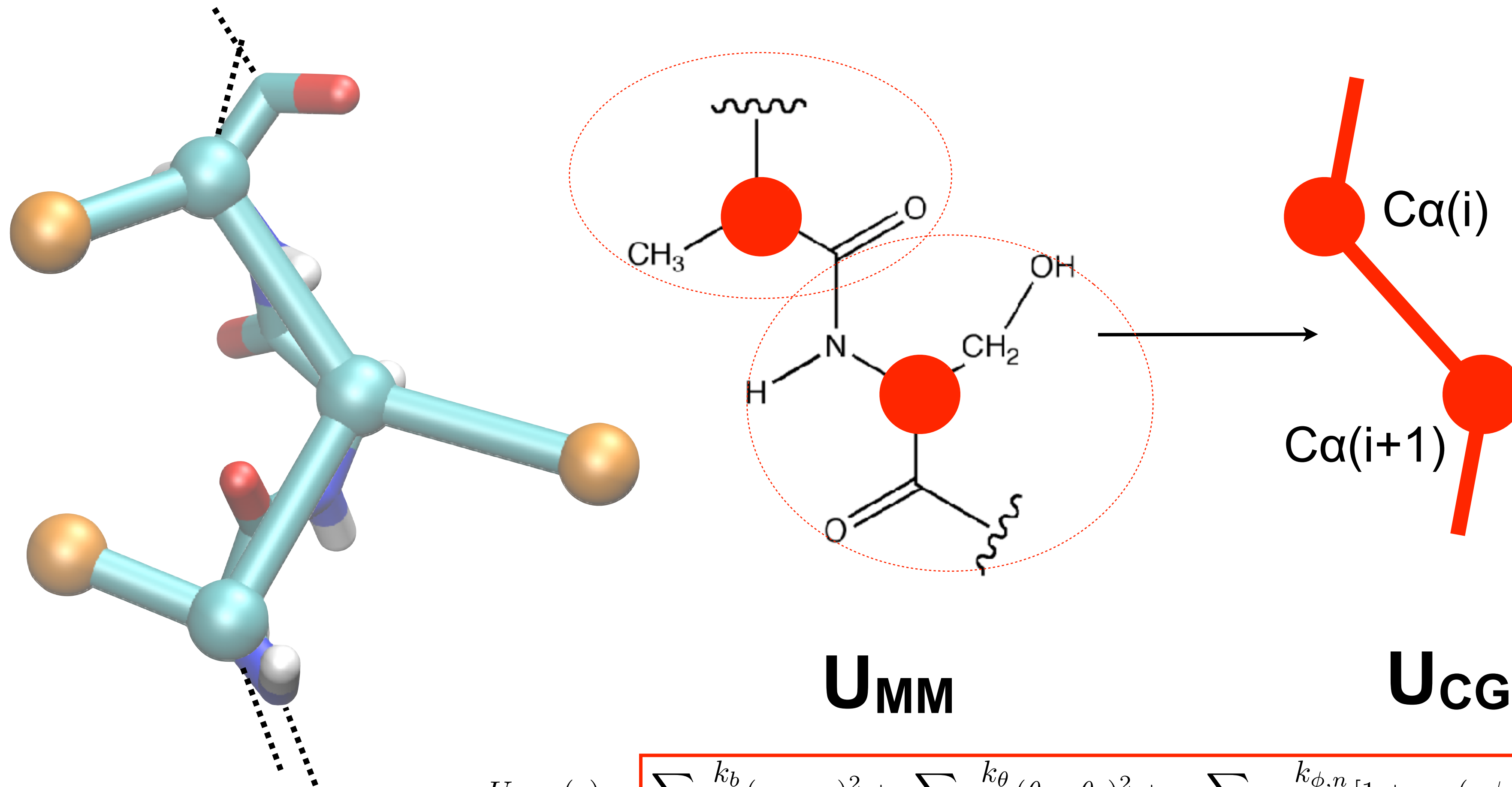
- the CG process implies a **simplification** of your potential that is not always rigorous and includes **approximations**
- what you obtain is an **effective** potentials which is parametrized to reproduce given properties

New directions

universal and computationally efficient machine-learned CG model for proteins



Coarse-graining degrees of freedom



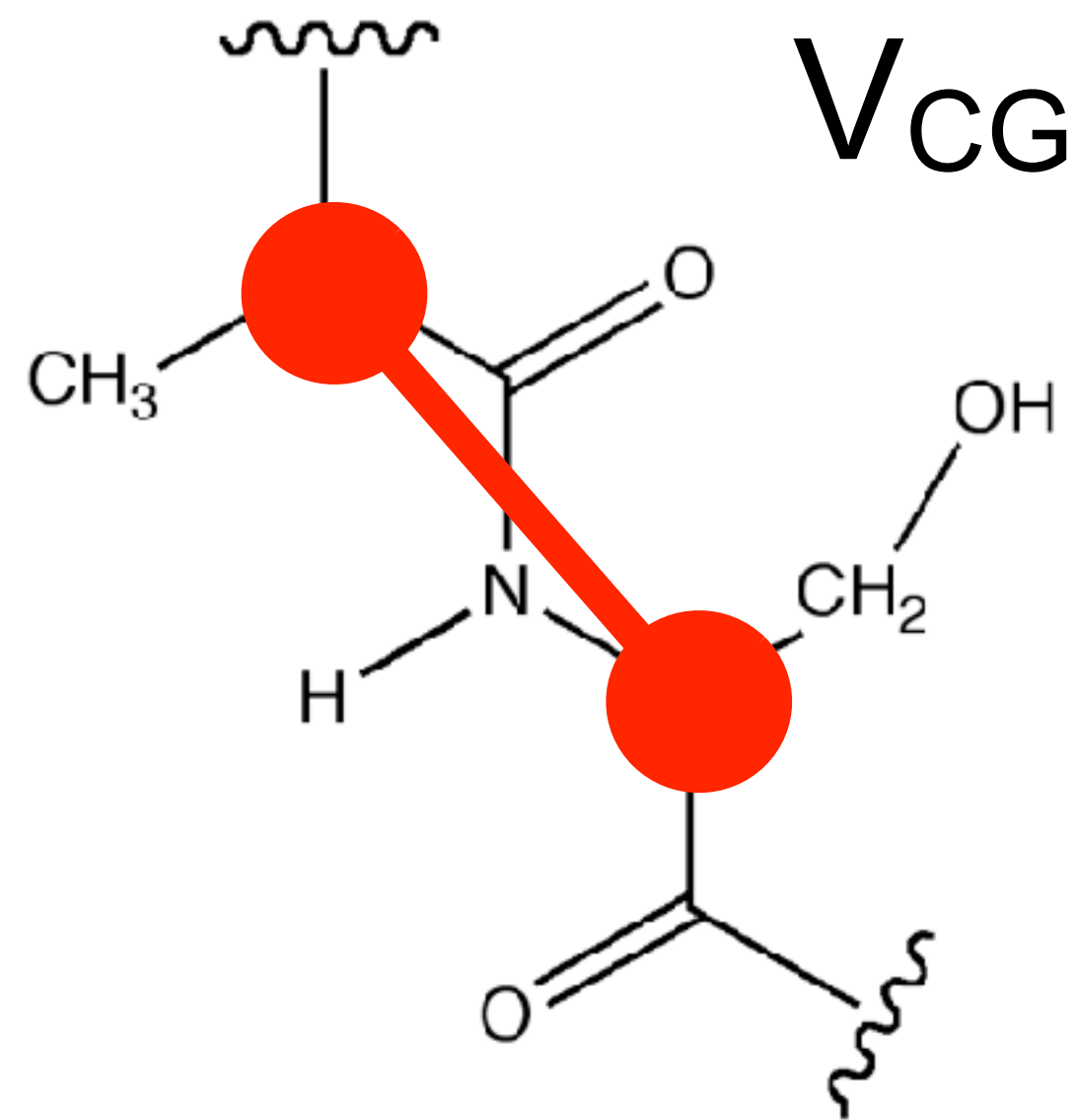
U_{MM}

U_{CG}

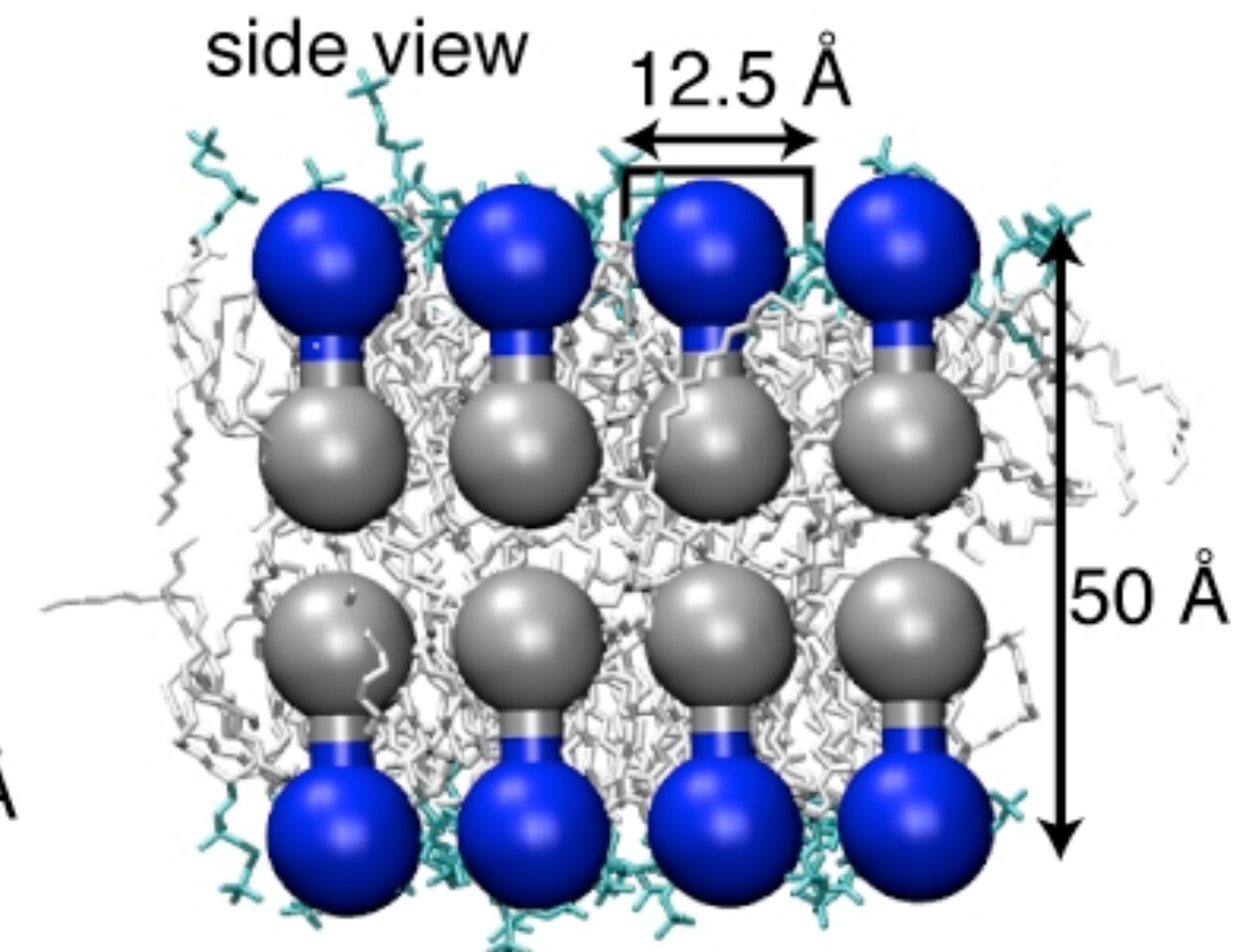
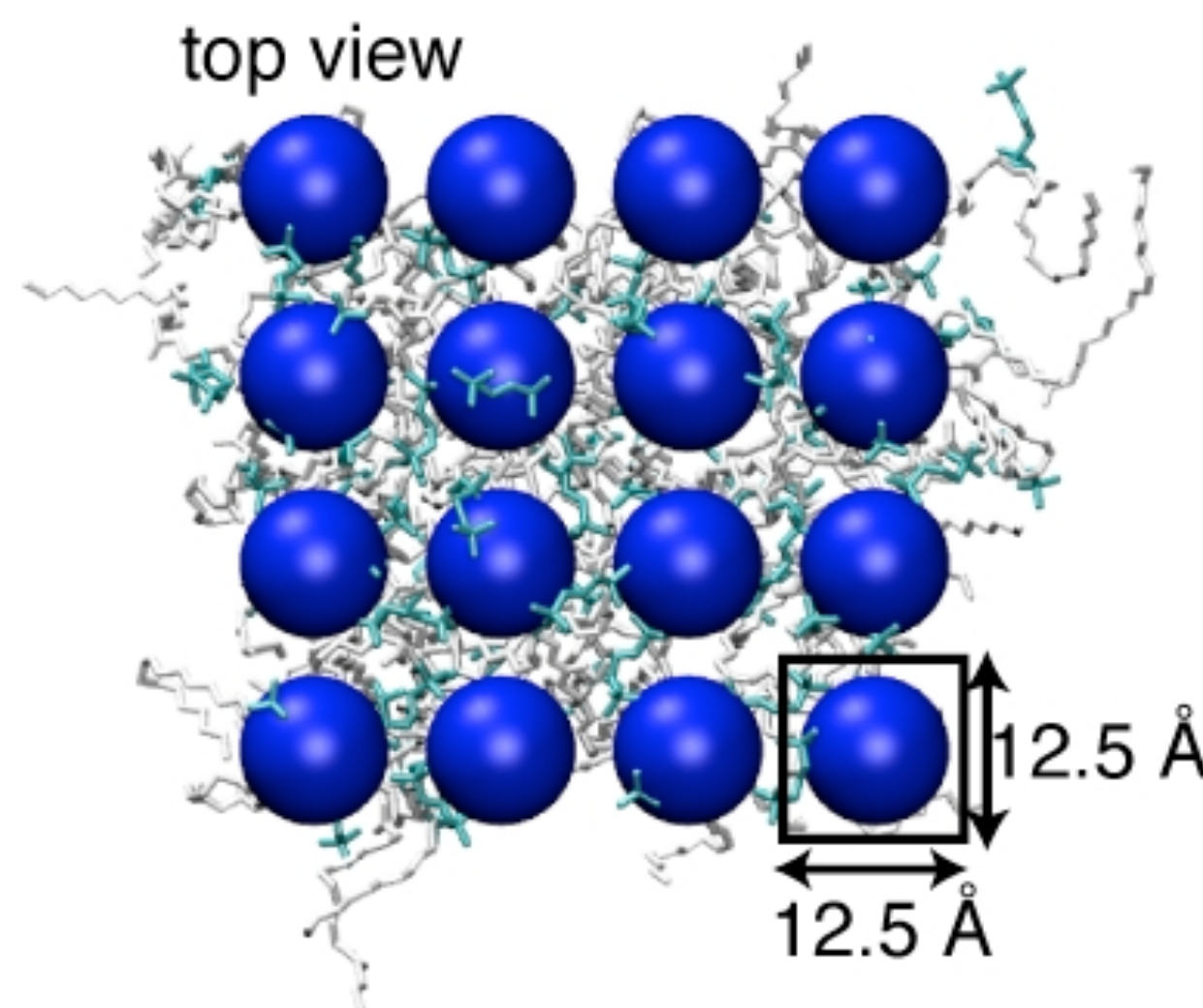
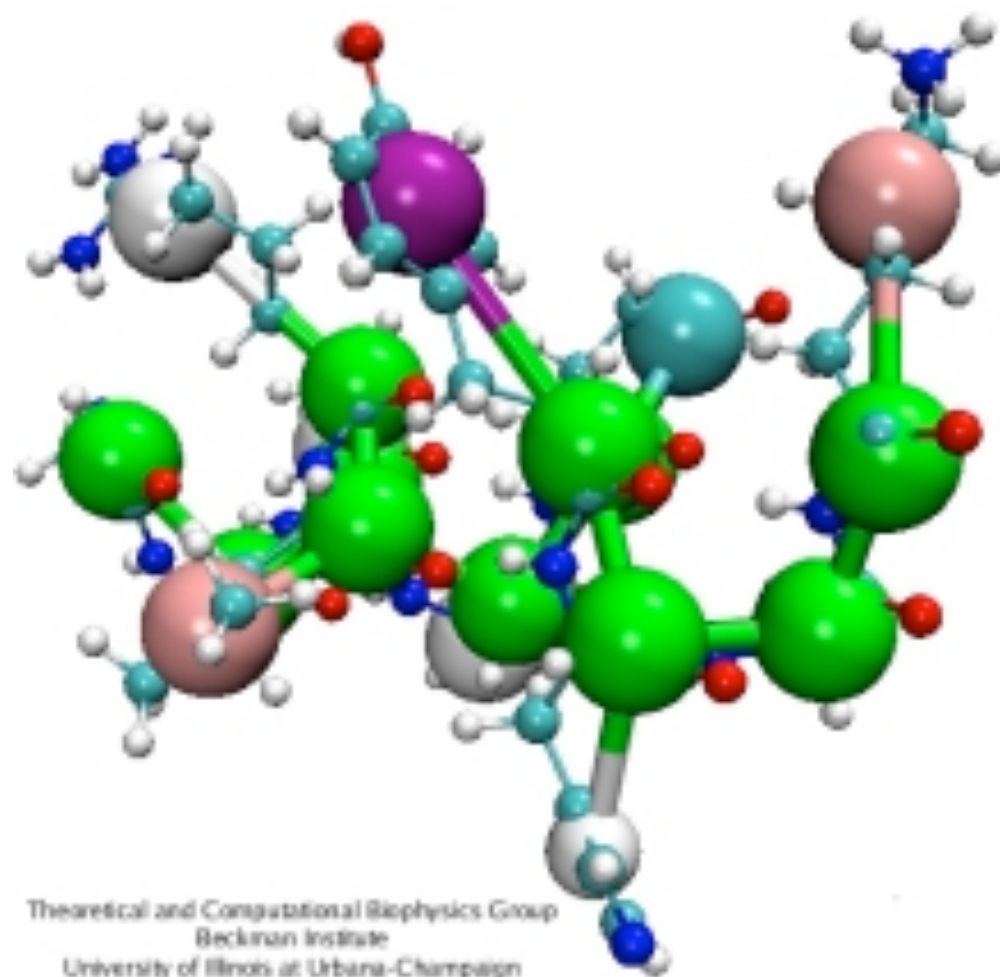
$$U_{MM}(r) = \sum_{bonds} \frac{k_b}{2} (r - r_0)^2 + \sum_{angles} \frac{k_\theta}{2} (\theta - \theta_0)^2 + \sum_{torsions, n} \frac{k_{\phi, n}}{2} [1 + \cos(n\phi - \delta)] +$$

$$+ \sum_{i > j}^N \left(\frac{A}{r_{ij}^{12}} - \frac{C}{r_{ij}^6} \right) + \sum_{i > j}^N \frac{1}{4\pi\epsilon_0} \frac{q_i q_j}{r_{ij}}$$

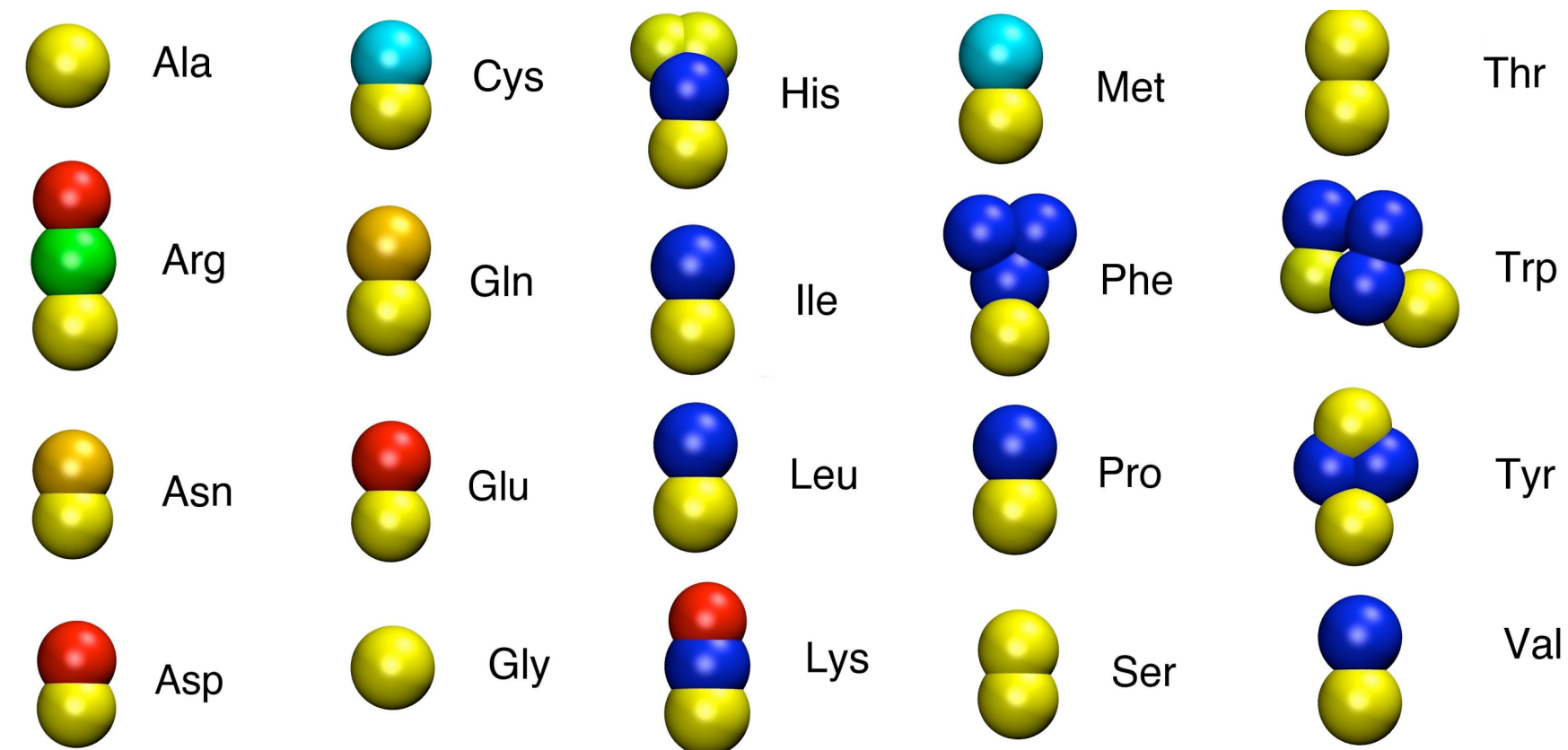
Coarse-grained force fields



- CG FF models are not topologically biased on the native structure
- softer interactions allow for **longer** timestep in MD simulations
- sampling on the **millisecond** timescale
- accuracy can be a problem (e.g. **no explicit electrostatic** contribution)
- biases on the secondary structures

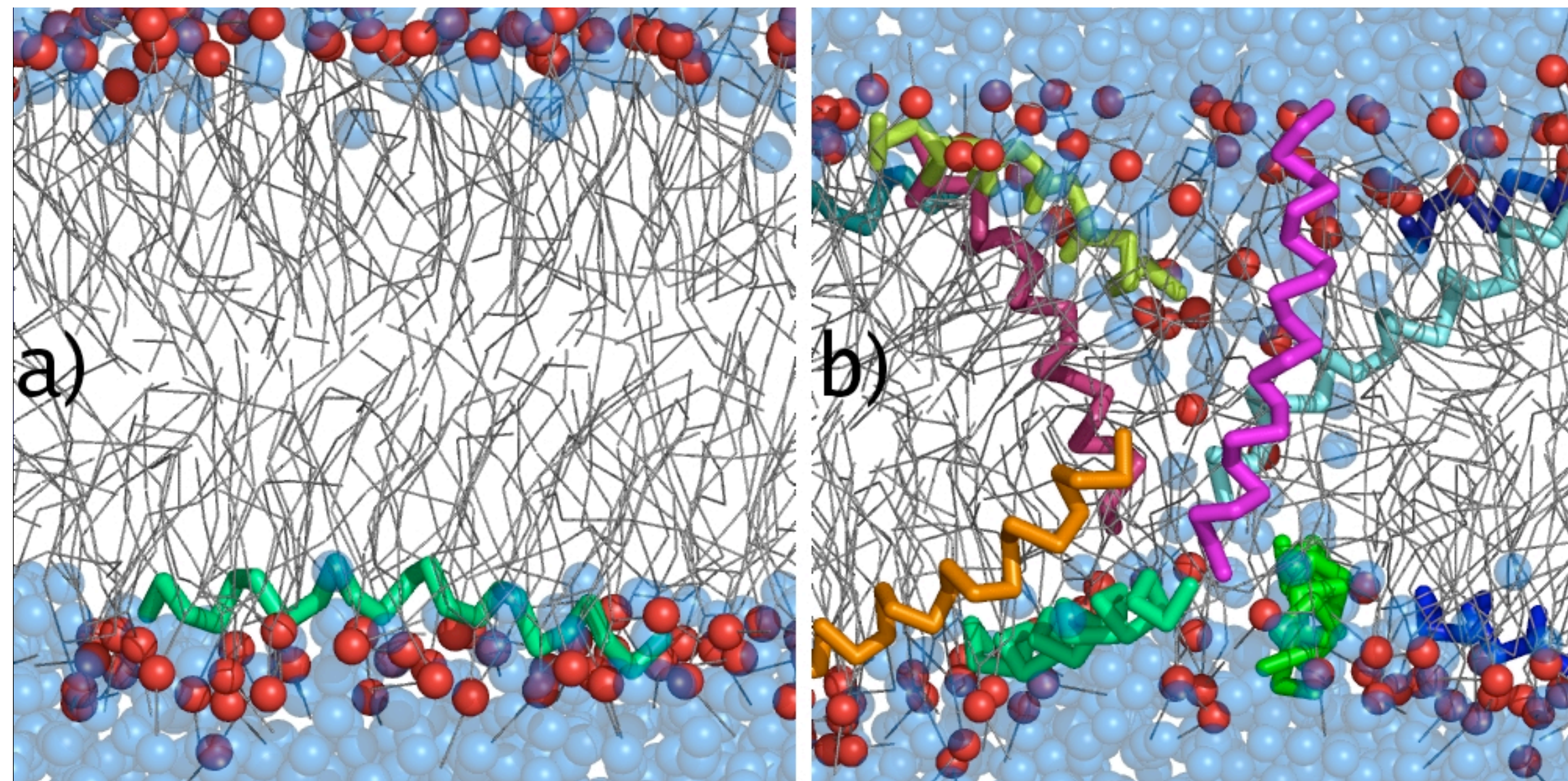


Coarse-grained MARTINI FF

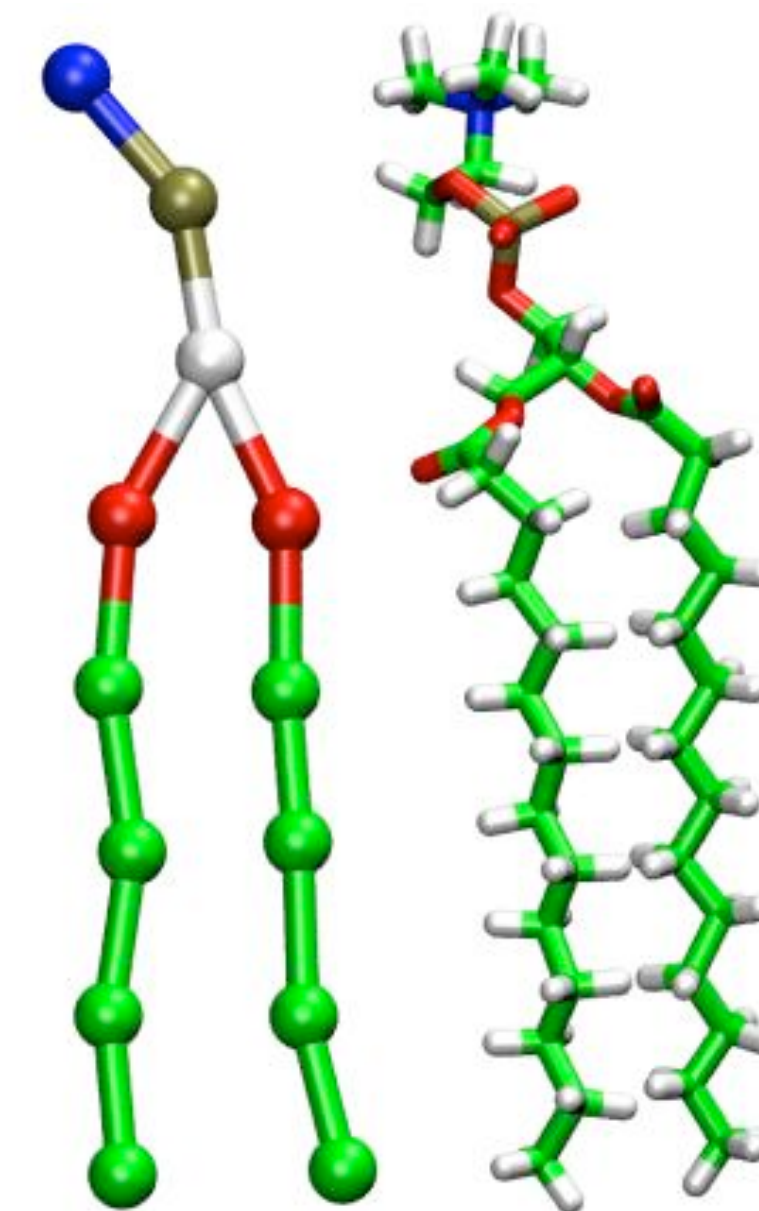


- MARTINI CG FF has functional form similar to MM FF
- 4-to-1 mapping from MM to CG
- very convenient for membranes and peptide-membrane interactions

Monticelli et al, JCTC 2008
Klein and coworkers



Magainin H2 in a DPPC bilayer, at low concentration (a) and high concentration

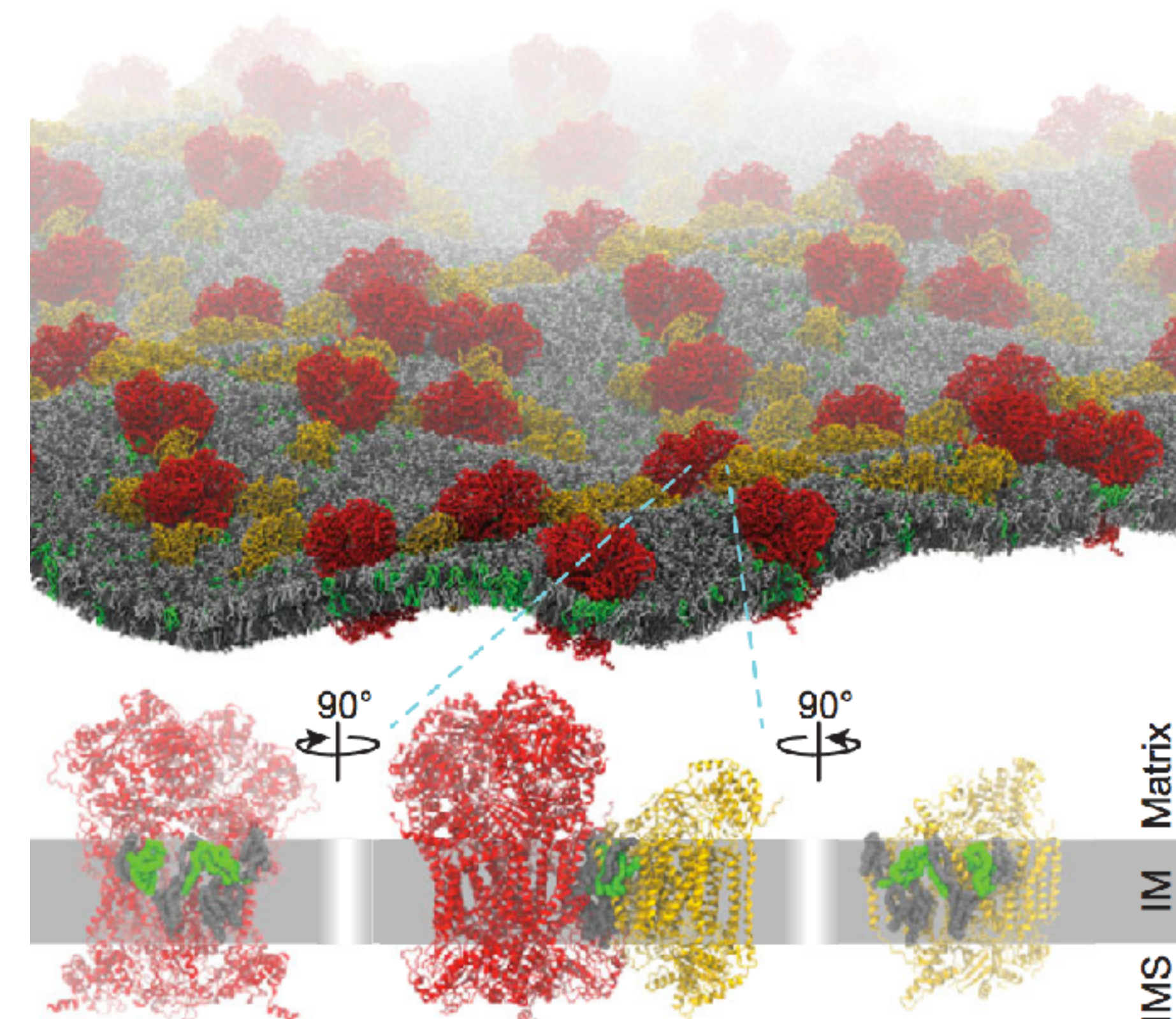
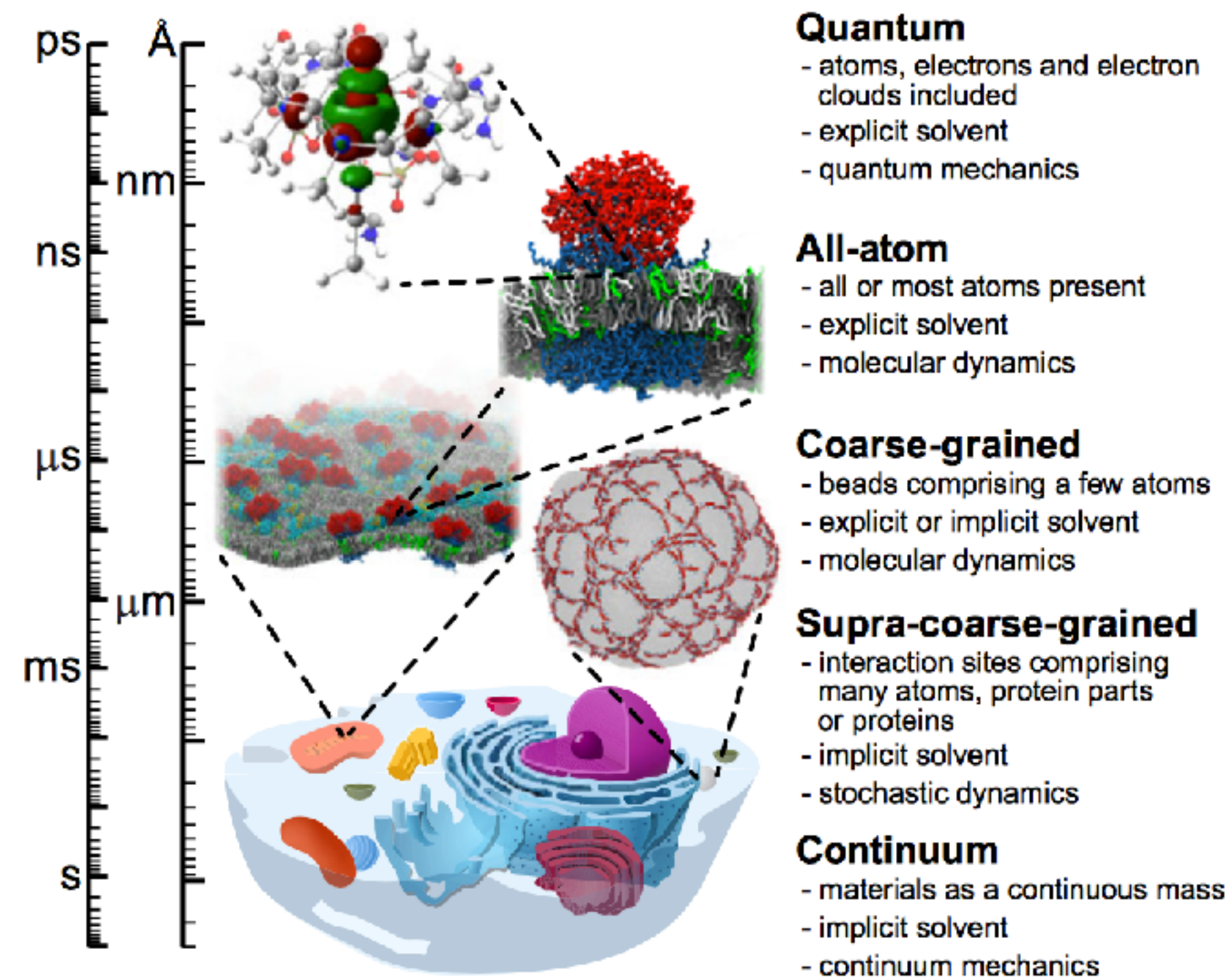


COMMENTARY

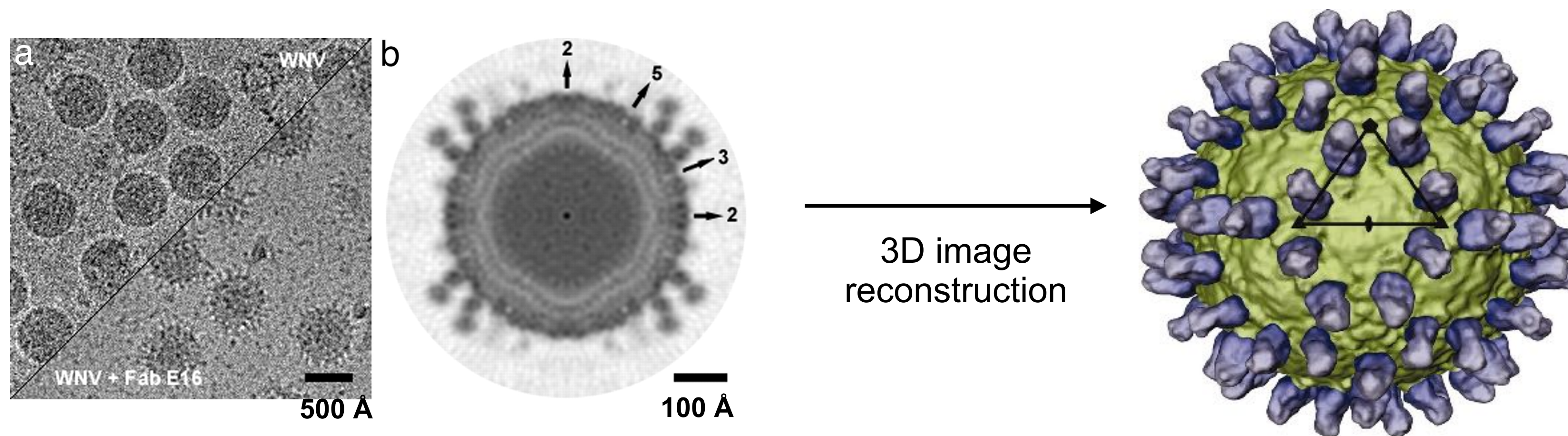
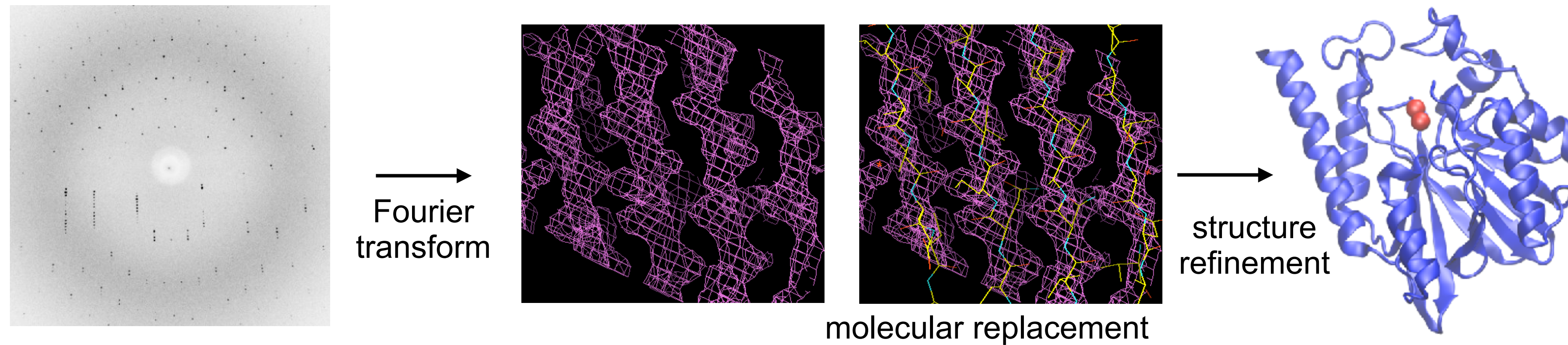
ARTICLE SERIES: IMAGING

Computational ‘microscopy’ of cellular membranes

Helgi I. Ingólfsson, Clément Arnarez, Xavier Periole and Siewert J. Marrink*

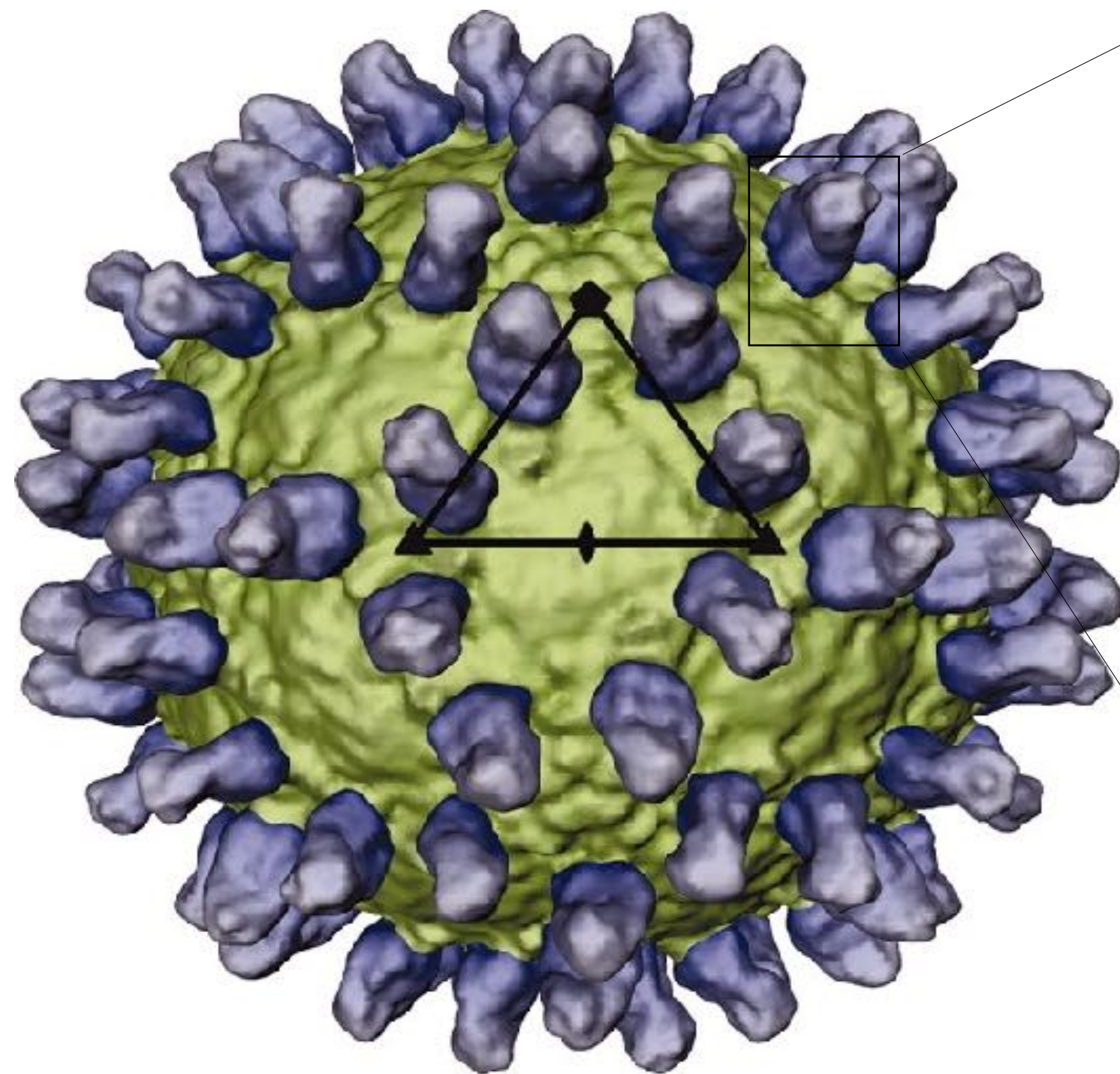


Structural biology methods are strongly based on theory and computation!

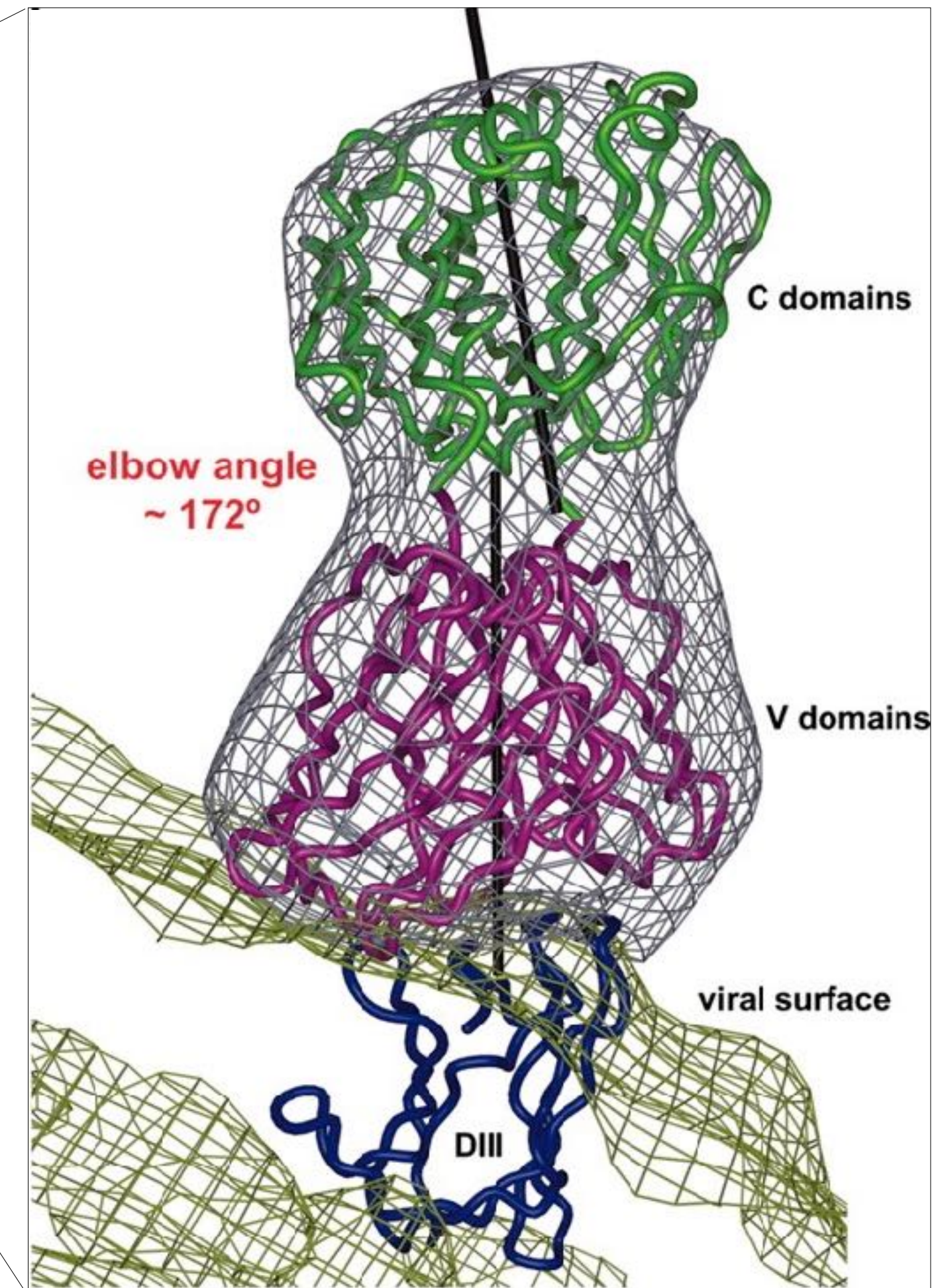


Combine cryo-EM and X-ray structures

surface rendering of the 3D image
reconstruction of WNV (green) in complex
with Fab E16 (blue) at 15-Å resolution



integrative modeling

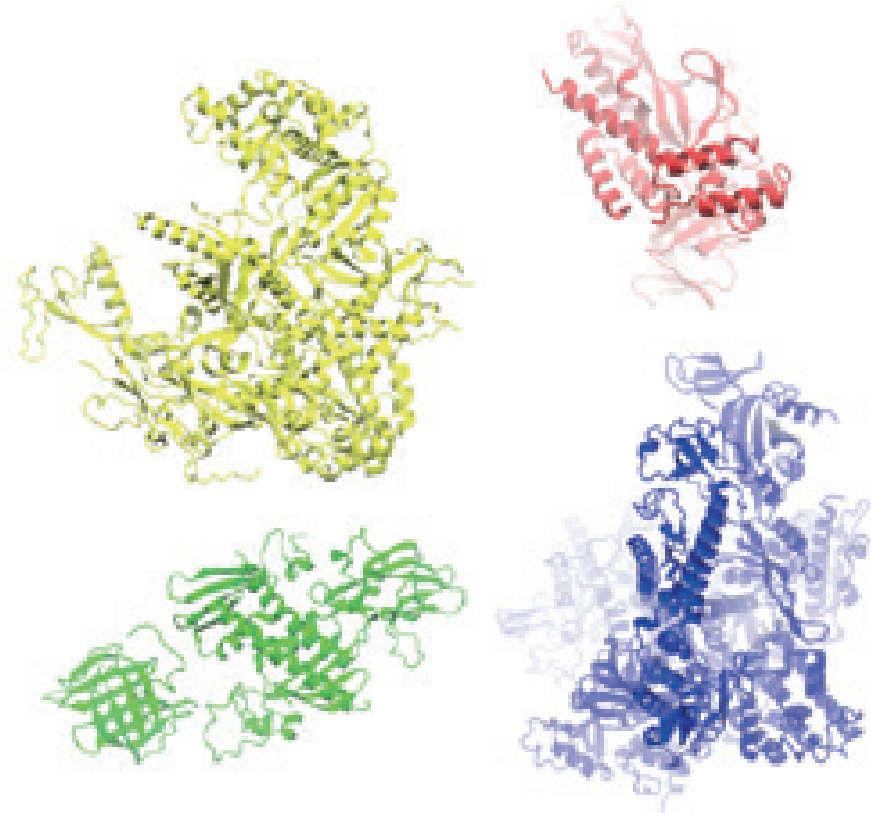


Kaufmann B et al. PNAS 2006;103:12400-12404

Integrative Modeling

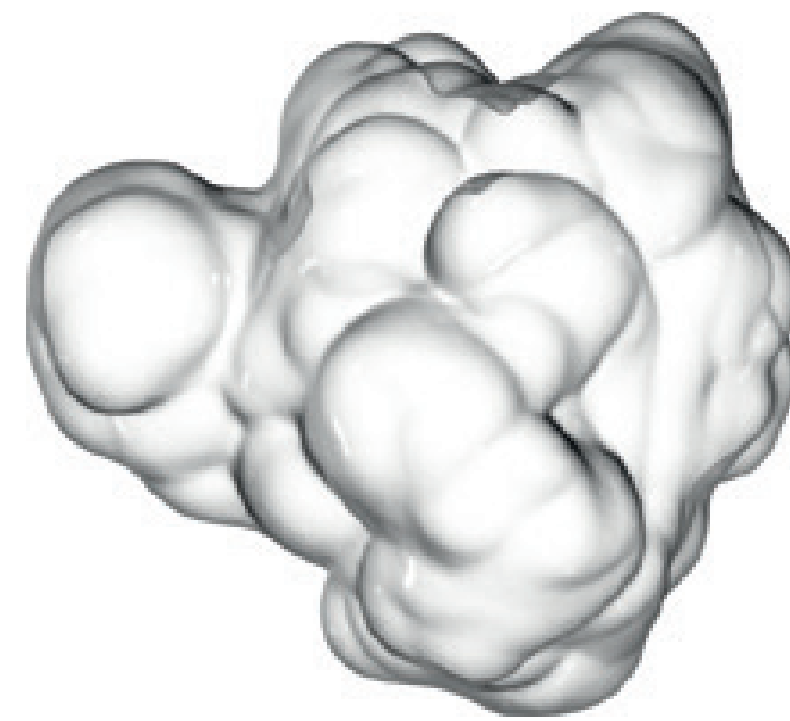
Individual subunits

- X-ray crystallography
- NMR
- Cryo-EM
- Homology models



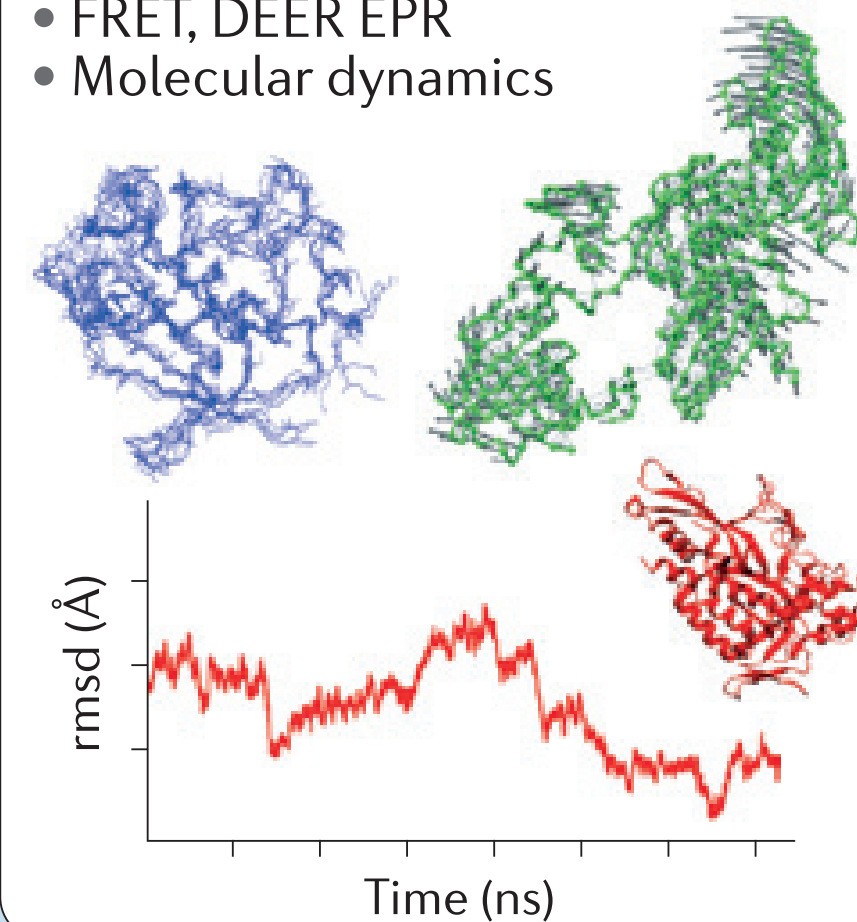
Volumetric maps

- Cryo-EM
- Electron tomography
- SAXS, SANS
- AFM



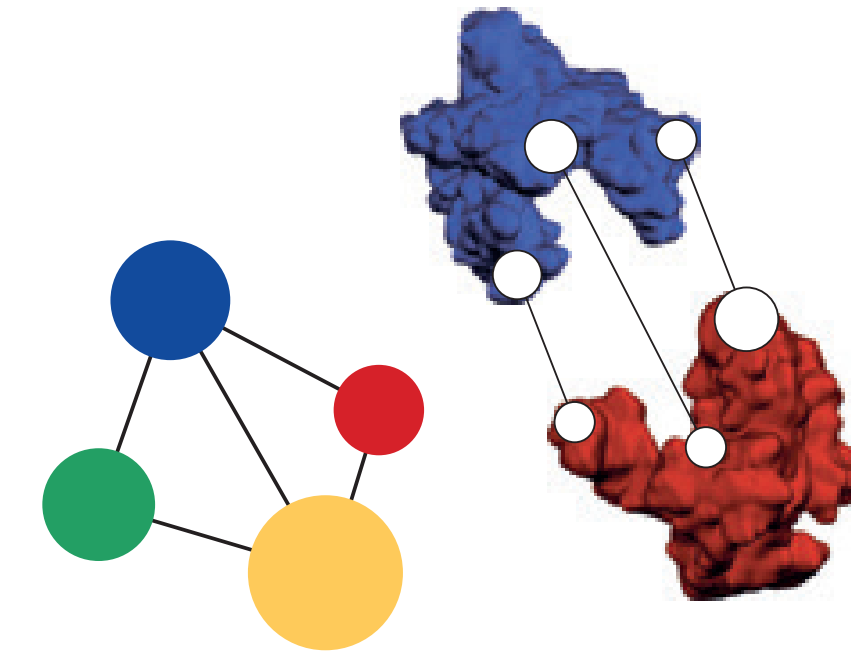
Structural flexibility

- Side-chain and backbone sampling
- Elastic network models
- NMR ensembles
- FRET, DEER EPR
- Molecular dynamics

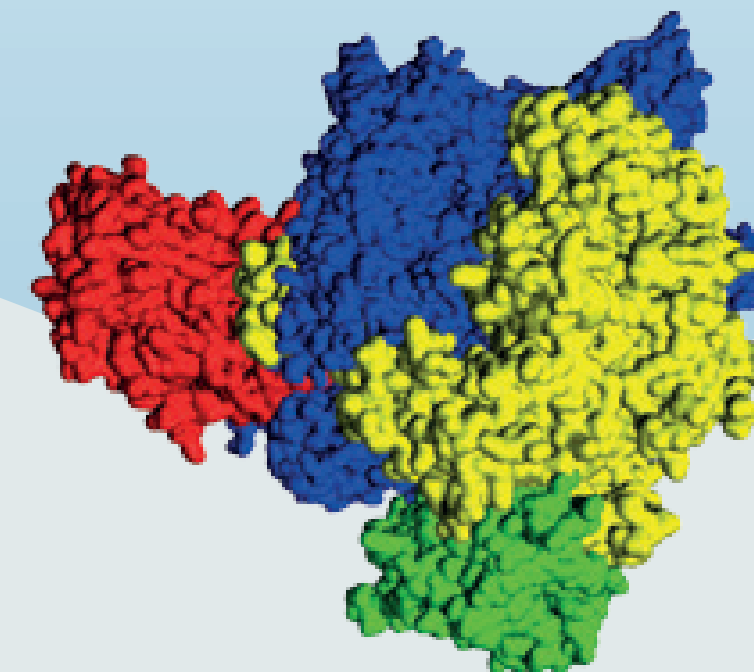


Spatial connectivity

- Mutagenesis
- Evolutionary couplings
- Chemical crosslinking
- Proteomics
- H/D exchange
- ChIP-seq and ChIP-exo
- 3C, 4C, 5C and Hi-C

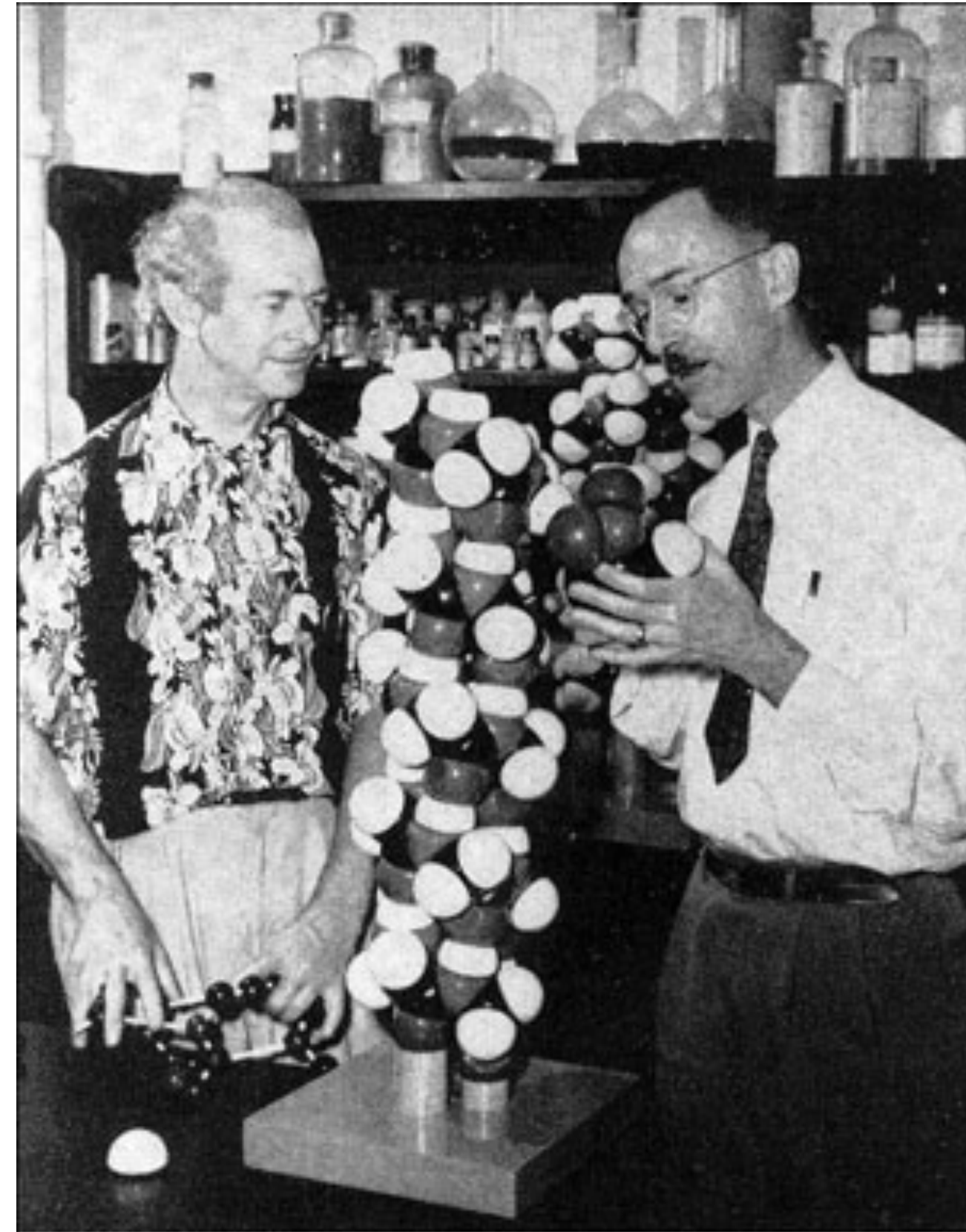
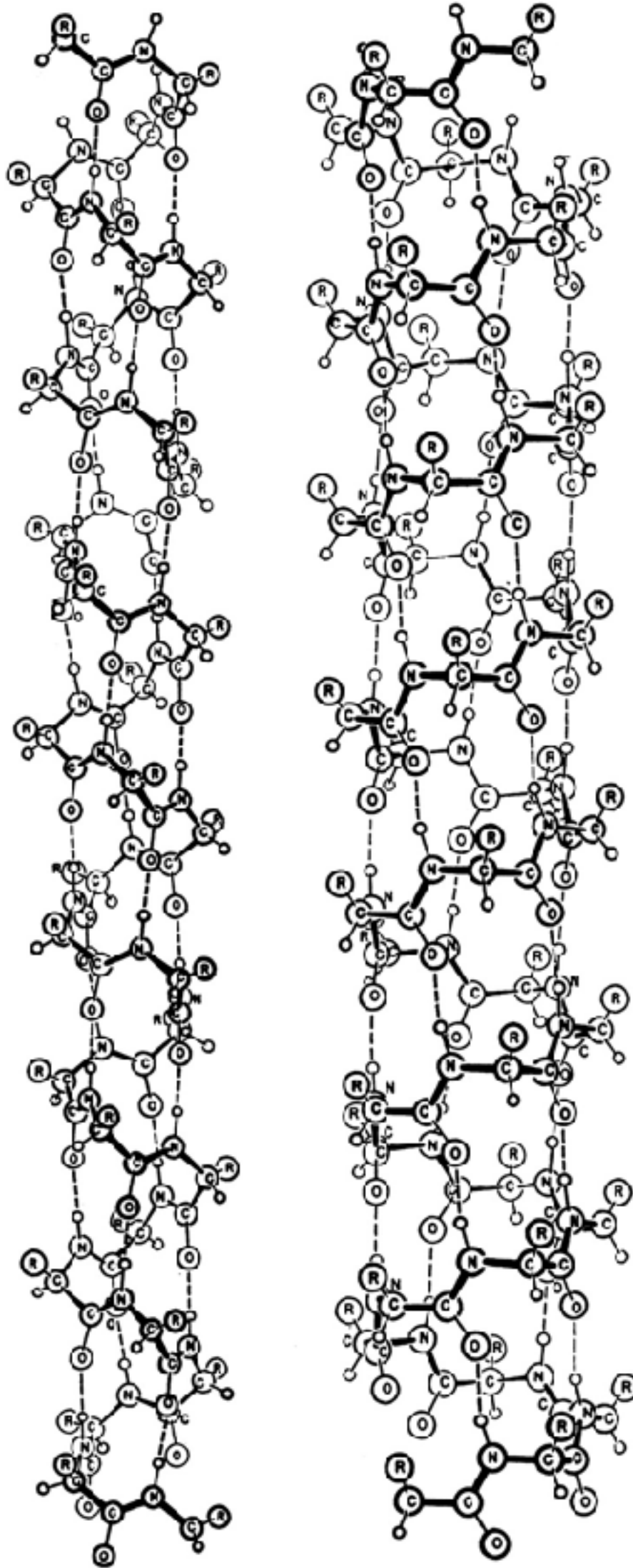


3C, chromatin conformation capture; 4C, circularized 3C; 5C, carbon-copy 3C; AFM, atomic force microscopy; ChIP-exo, ChIP-seq with an exonuclease sample preparation step; ChIP-seq, chromatin immunoprecipitation followed by sequencing; DEER EPR, double electron-electron resonance electron paramagnetic resonance; FRET, fluorescence resonance energy transfer; H/D exchange, hydrogen-deuterium exchange; NMR, nuclear magnetic resonance; Hi-C, genome-wide 3C; rmsd, root-mean-square deviation; SANS, small-angle neutron scattering; SAXS, small-angle X-ray scattering.



Near-atomic-resolution structure of supramolecular assemblies

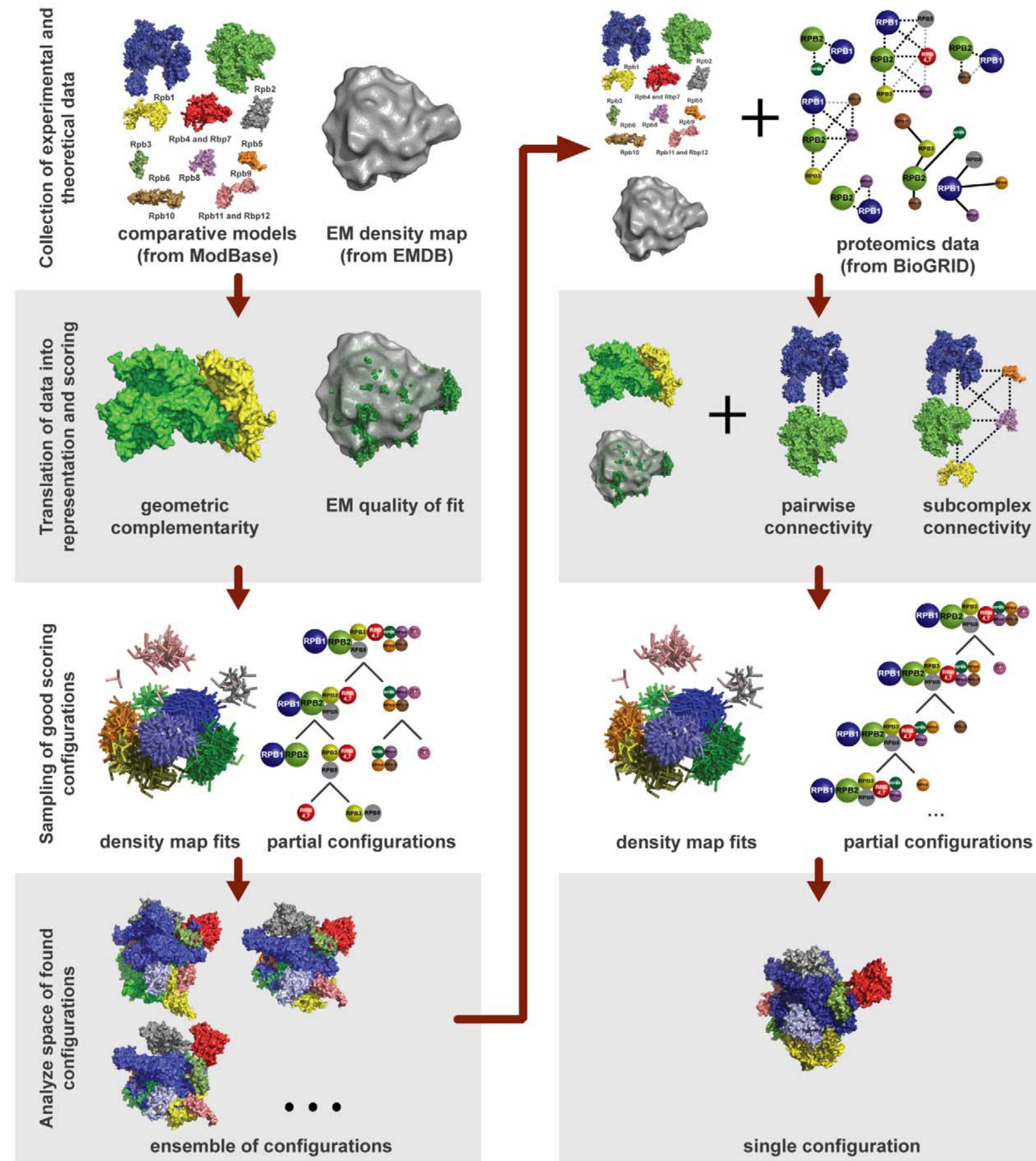
(the dawn of) Integrative modeling



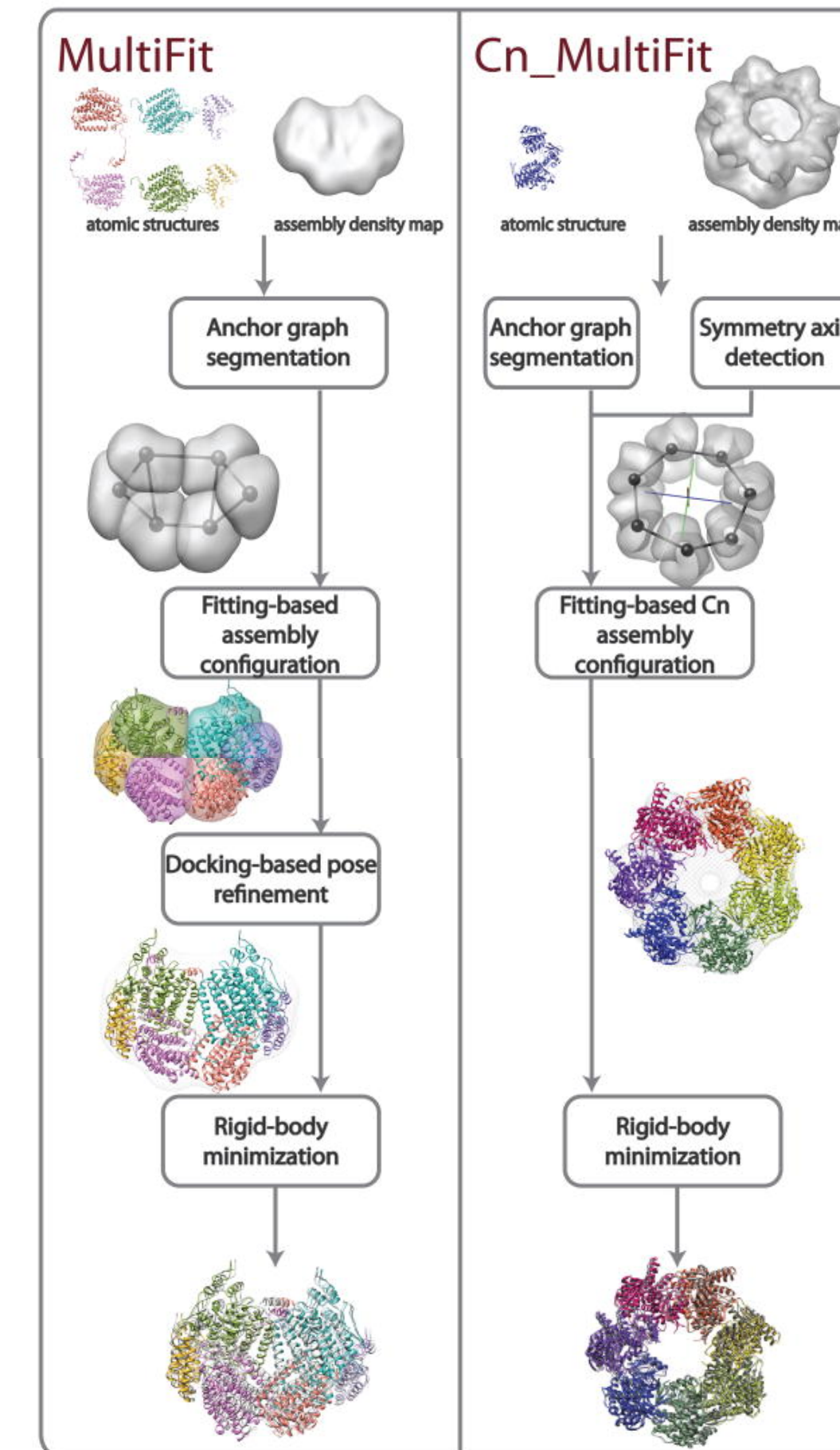
Pauling and Corey protein structure papers (1951):

X-ray myoglobin 1959 (Kendrew and Perutz)

Integrative modeling



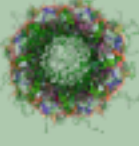
IMP: integrative modeling platform
<http://salilab.org/imp/>



<https://modbase.compbio.ucsf.edu/multifit/>


Russel et al. (2012) Putting the Pieces Together: Integrative Modeling Platform Software for Structure Determination of Macromolecular Assemblies. PLoS Biol 10(1): e1001244. doi:10.1371/journal.pbio.1001244

https://pdb-dev.wwpdb.org

**PDB-Dev**
Prototype Archiving System for Integrative Structures

Released Entries: 75

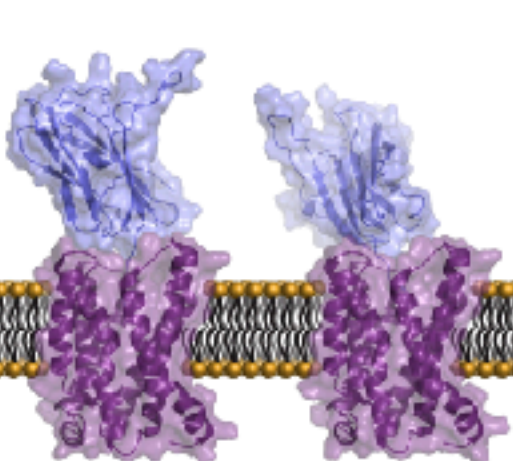
[Home](#) [About](#) [Deposit](#) [Contact](#) [FAQ](#)

[Browse Structures](#)  [Search Tips](#)

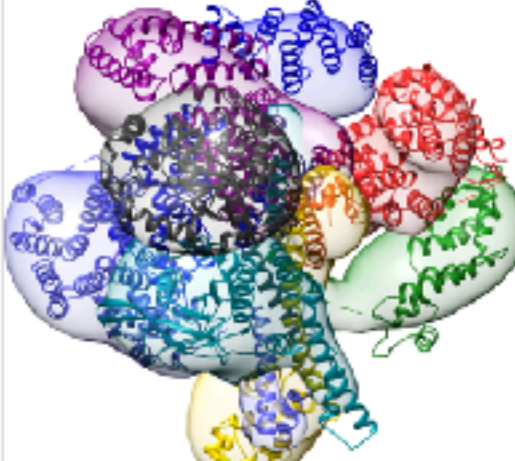
Welcome to PDB-Dev

PDB-Dev is a prototype archiving system for structural models obtained using integrative or hybrid modeling and is funded by the NSF ABI Development Program. Structural characterization of many complex macromolecular assemblies is increasingly carried out using integrative modeling, where a combination of complementary experimental and computational techniques is used to determine the structure. The structural models obtained through integrative modeling are collected, archived and disseminated to the public through PDB-Dev. Once the mechanisms for processing integrative models are fully established through PDB-Dev, the key components will be integrated with the [wwPDB OneDep](#) system and the PDB-Dev holdings will be moved into the PDB.

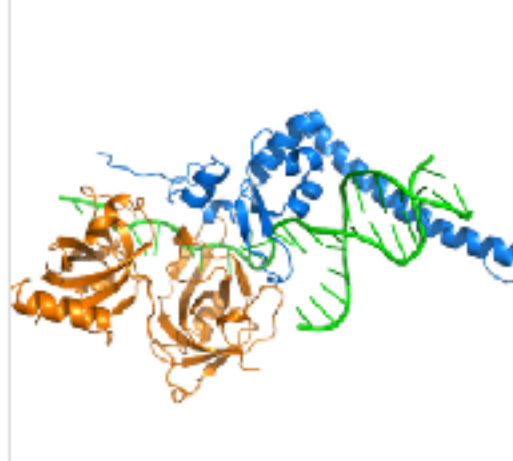
Released PDB-Dev Structures



PDBDEV_00000036
Myeloid-derived growth factor with KDEL receptor
Release Date: 2019-12-18
Publication doi: [10.1038/s41467-019-13677-5](#)



PDBDEV_00000037
Human COP9 Signalosome
Release Date: 2020-01-30
Publication doi: [10.1073/pnas.1915542117](#)



PDBDEV_00000039
Nucleotide excision repair complex
Release Date: 2020-02-07
Publication doi: [10.1093/nar/gkz1231](#)

News

[All News](#)

PDB-Dev Publication

Recent updates to PDB-Dev, including the development of a new data harvesting system, have been published in *Acta Crystallographica Section D*: Vallat B et al., New system for archiving integrative structures, *Acta Cryst. 2021*; D77: 1486-1496. doi:10.1107/S2059798321010871

New Data Harvesting System

We have developed a new [Data Harvesting System](#) that provides a web interface for depositors to assemble all the information required for archiving integrative structures in PDB-Dev and to create a compliant mmCIF file. This includes the submission of integrative structures, associated spatial restraints and starting models used, modeling protocols and metadata information. [Read more...](#)

BioExcel Webinar

Recently, the PDB-Dev team participated in the [BioExcel webinar series](#). The webinar presentation was titled "PDB-Dev: A prototype system for archiving integrative structures". The recorded version of the webinar is available on [youtube](#).

Visualization of Structures using Molstar


3D visualization of structures using [Molstar](#) is now available on PDB-Dev. Structures in PDB-Dev can be directly visualized from the respective entry pages. Molstar can visualize atomic and multi-scale structures. [Read more...](#)

Welcome to the Updated PDB-Dev website

The PDB-Dev web interface has been revamped to provide dynamic, responsive and mobile-friendly web pages. PDB-Dev website now includes a new service that facilitates search and retrieval of integrative structures archived in PDB-Dev. [Read more...](#)

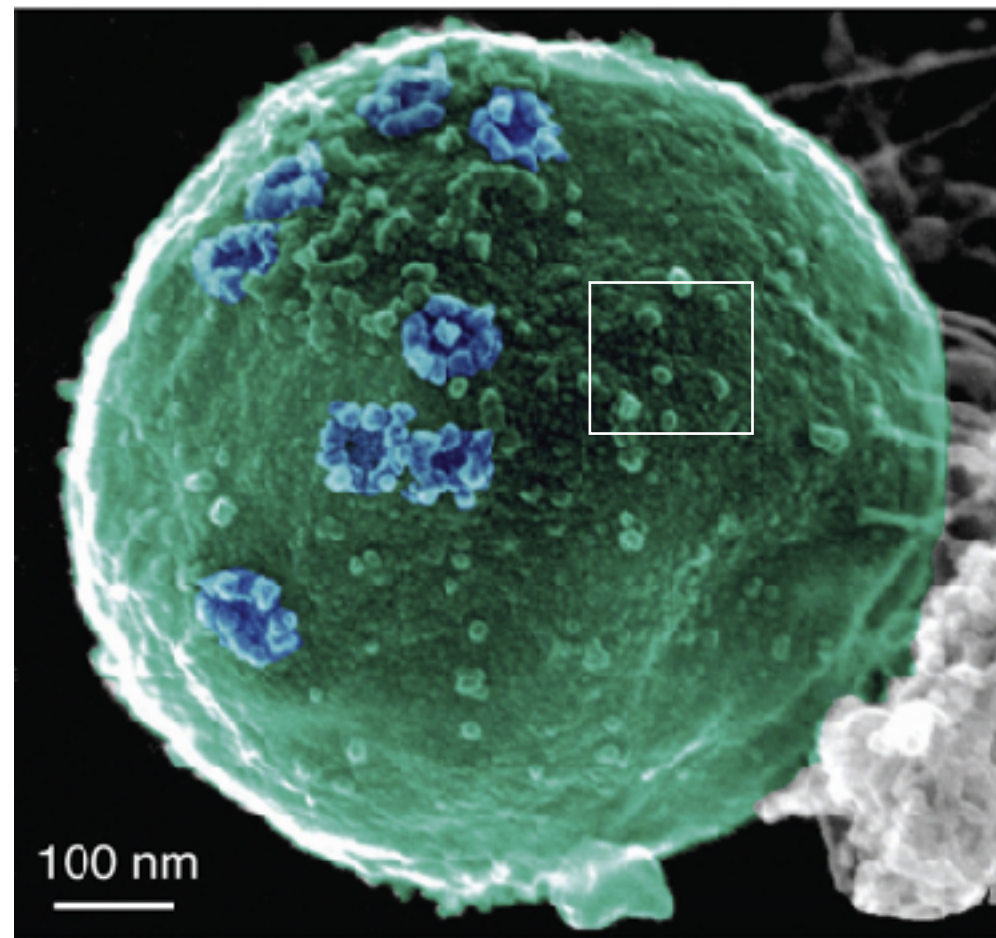
Whitepaper from the Integrative Modeling Community

In March 2019, a satellite workshop titled "Working towards federating structural models and data" was held at the Biophysical Society Annual meeting in Baltimore, Maryland. A whitepaper summarizing the outcomes of the workshop was recently published in the journal *Structure*. [Read more...](#)

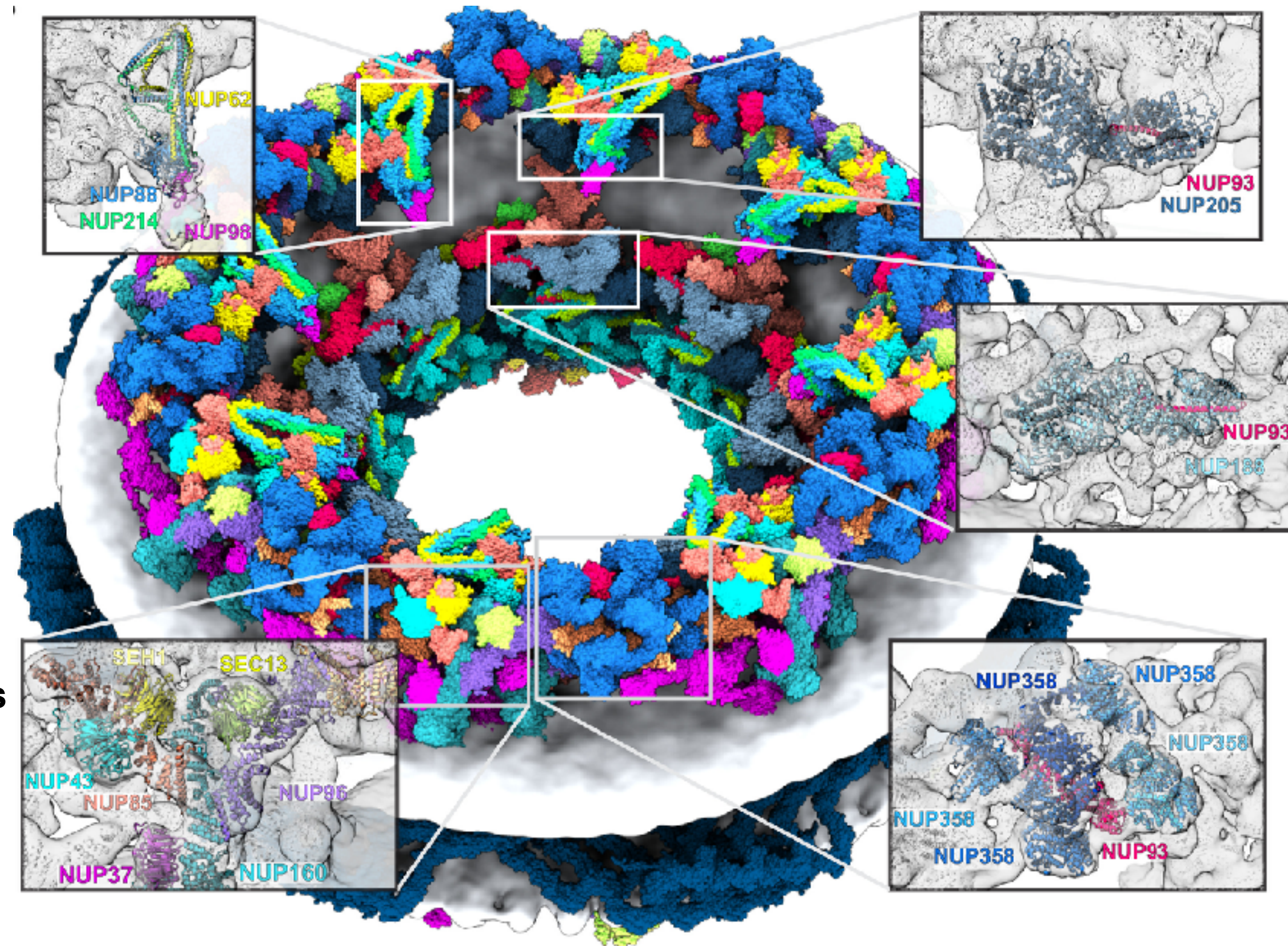


Supported by
National Science Foundation

e.g., the Nuclear Pore Complex ... combining AF and cryoEM/ET



yeast NPC :
~52 MDa complex
~550 protein subunits of ~30 different types



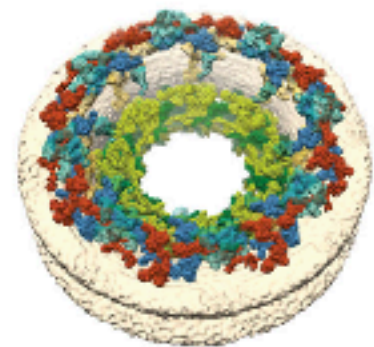
Assembline

Assembline is an assembly line of macromolecular assemblies!

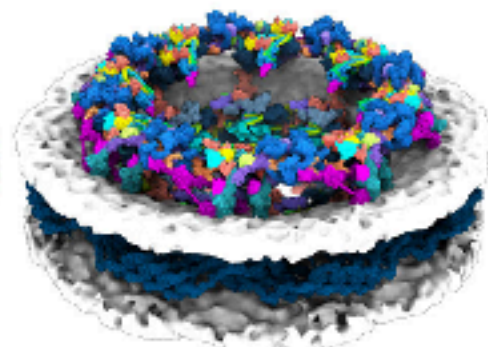
Assembline is a multi-step protocol for integrative structural modeling of macromolecular complexes based on electron microscopy, cross-linking mass spectrometry and other data. The protocol is based on our [Xlink Analyzer](#) and external software: [Integrative Modeling Platform \(IMP\)](#), [Python Modeling Interface \(PMI\)](#), [UCSF Chimera](#), and [imp-sampcon](#). Comparing to other methods, Assembline enables efficient sampling of conformational space through a multi-step procedure, provides new modeling restraints, and includes a unique configuration system for setting up the modelling project.

Complexes modeled using Assembline

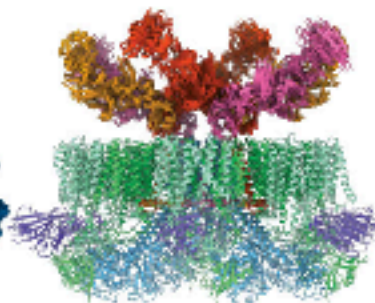
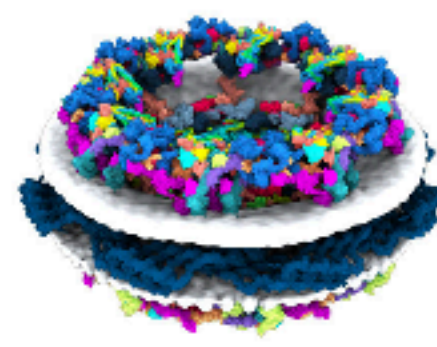
Human pore complex
(Science, 2016)



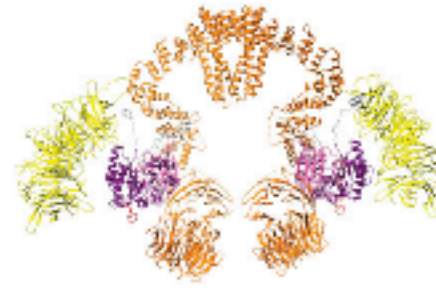
Human pore complex (bioRxiv, 2021)



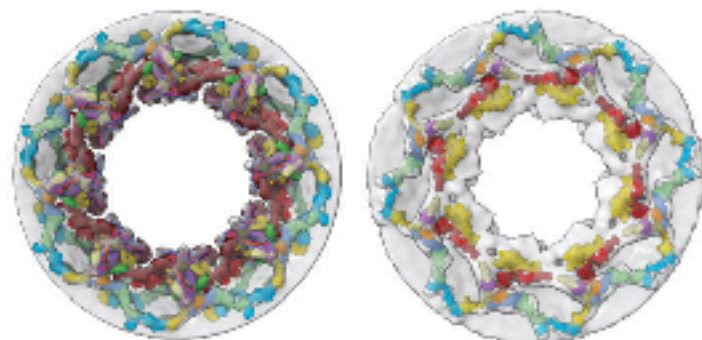
Type VII secretion system
(Science Advances, 2021)



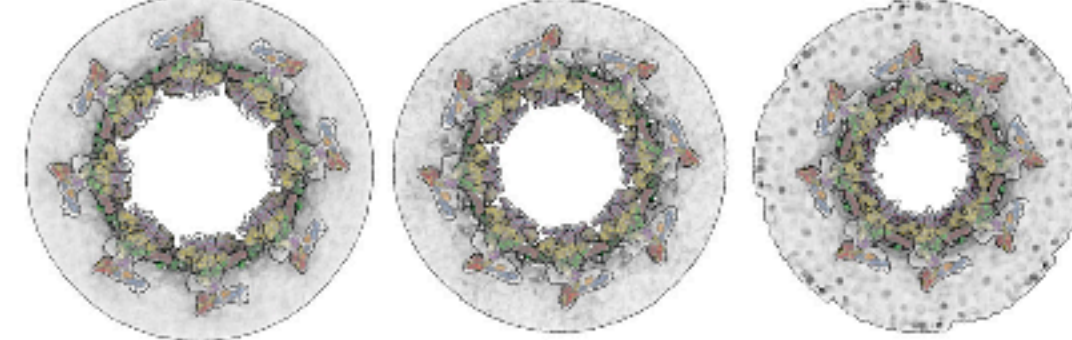
Elongator complex
(EMBO Reports, 2017)



Budding yeast nuclear pore complex (Nature, 2020)



Fission yeast nuclear pore complex (Science, 2021)



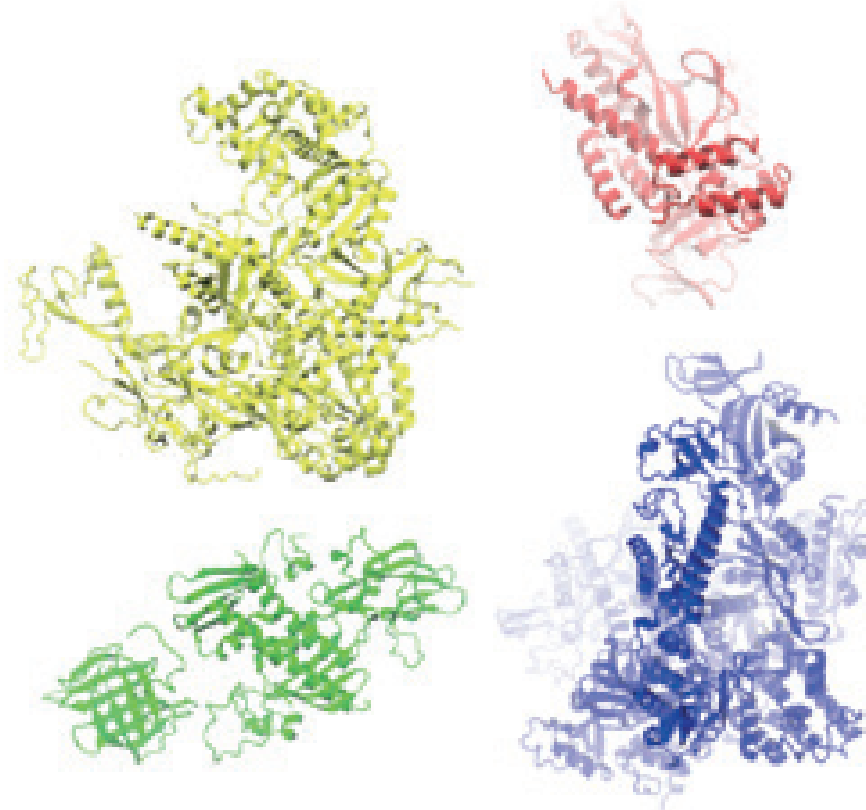
Code

- [Assembline installation package for Anaconda](#)
- [Assembline code](#)
- [Installation instructions and manual](#)

Integrative Modeling

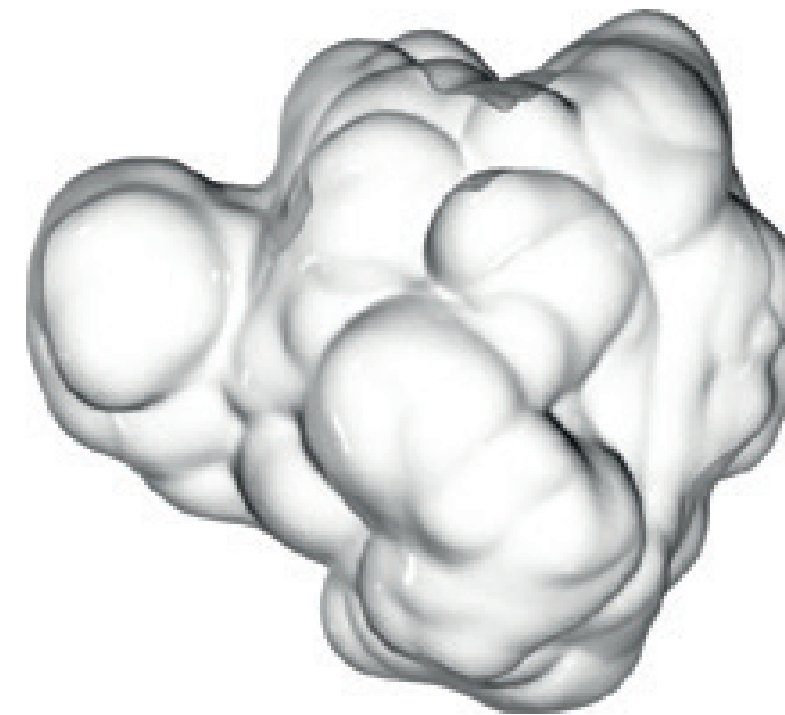
Individual subunits

- X-ray crystallography
- NMR
- Cryo-EM
- Homology models



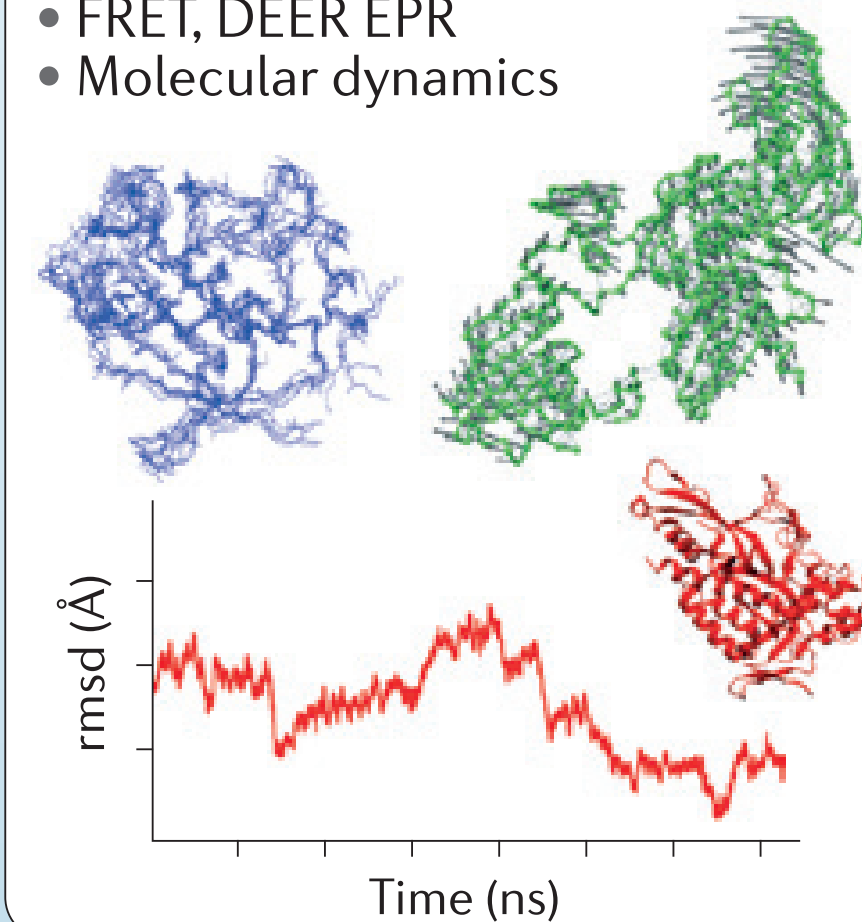
Volumetric maps

- Cryo-EM
- Electron tomography
- SAXS, SANS
- AFM



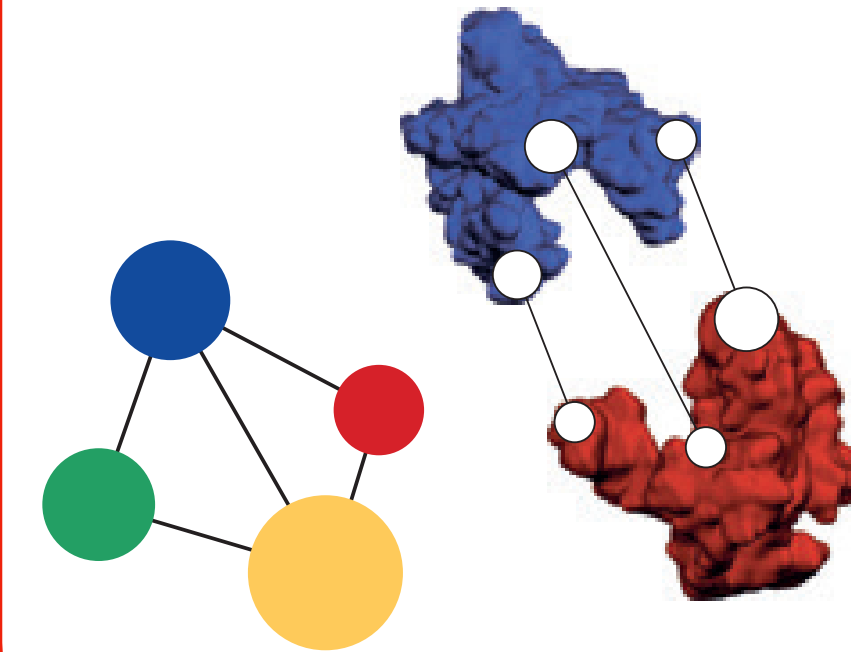
Structural flexibility

- Side-chain and backbone sampling
- Elastic network models
- NMR ensembles
- FRET, DEER EPR
- Molecular dynamics

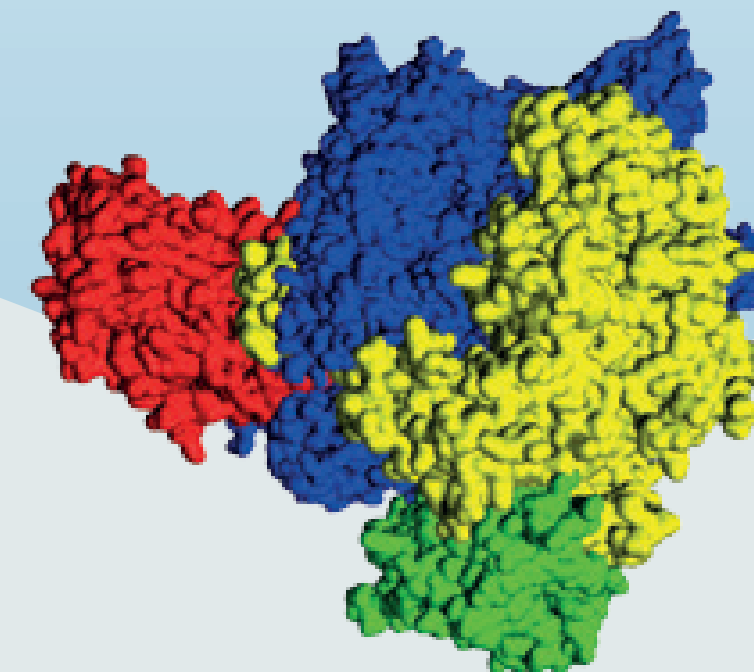


Spatial connectivity

- Mutagenesis
- Evolutionary couplings
- Chemical crosslinking
- Proteomics
- H/D exchange
- ChIP-seq and ChIP-exo
- 3C, 4C, 5C and Hi-C

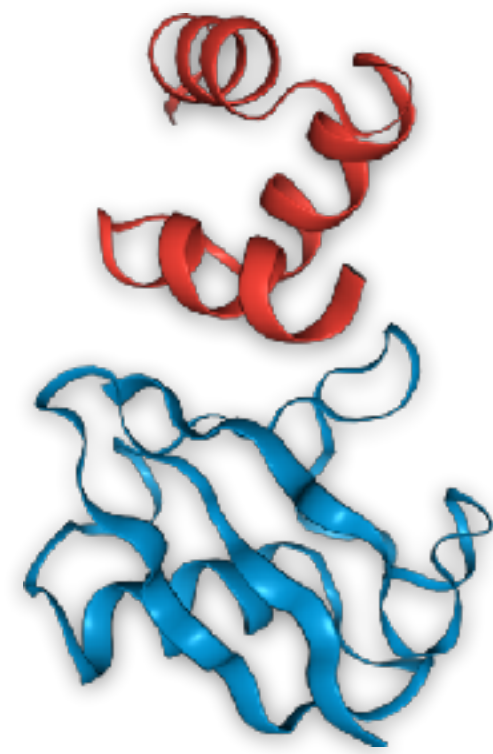


3C, chromatin conformation capture; 4C, circularized 3C; 5C, carbon-copy 3C; AFM, atomic force microscopy; ChIP-exo, ChIP-seq with an exonuclease sample preparation step; ChIP-seq, chromatin immunoprecipitation followed by sequencing; DEER EPR, double electron-electron resonance electron paramagnetic resonance; FRET, fluorescence resonance energy transfer; H/D exchange, hydrogen-deuterium exchange; NMR, nuclear magnetic resonance; Hi-C, genome-wide 3C; rmsd, root-mean-square deviation; SANS, small-angle neutron scattering; SAXS, small-angle X-ray scattering.

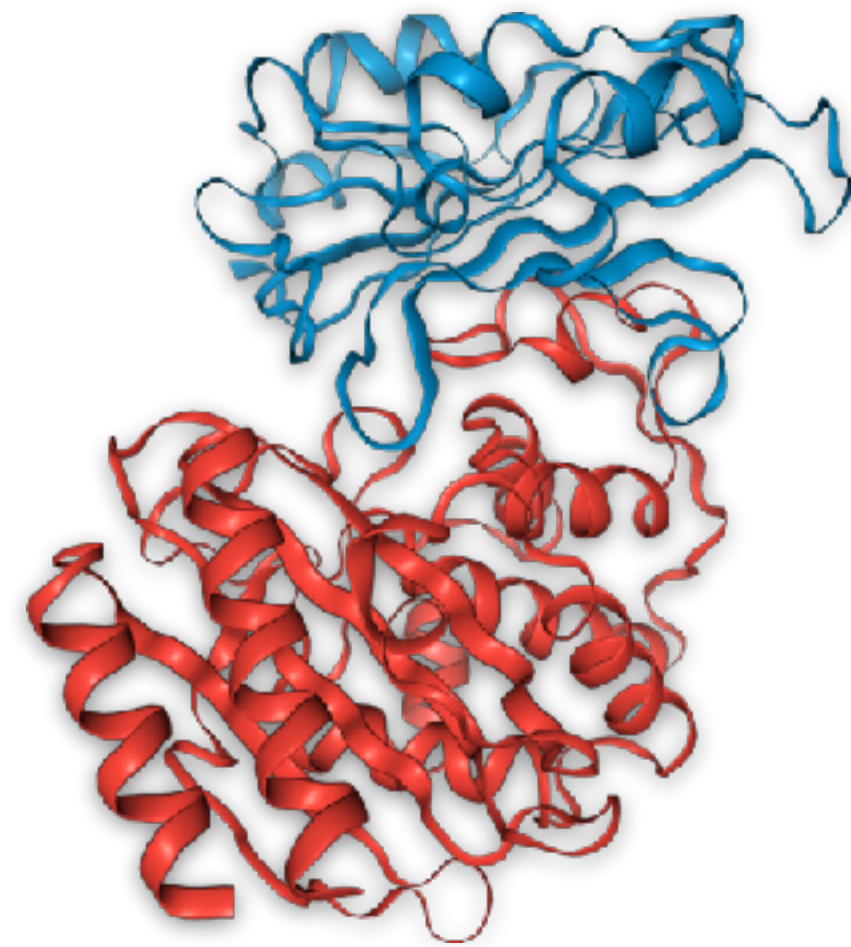


Near-atomic-resolution structure of supramolecular assemblies

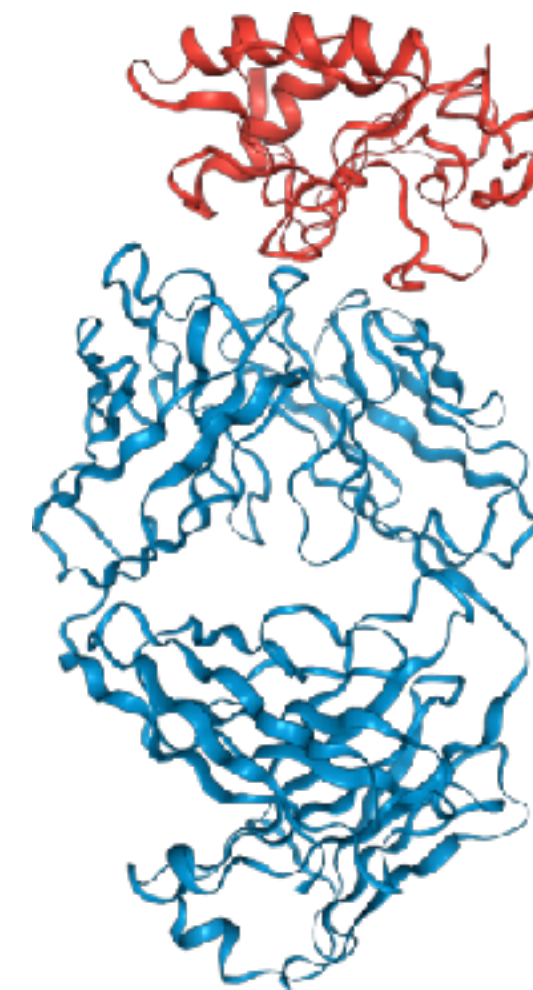
Protein-protein interaction interfaces



Enzyme-Substrate
(2OOB)
Ubiquitin bound to
ubiquitin ligase



Enzyme-Inhibitor (1JTG)
Beta-lactamase bound to
beta-lactamase inhibitor



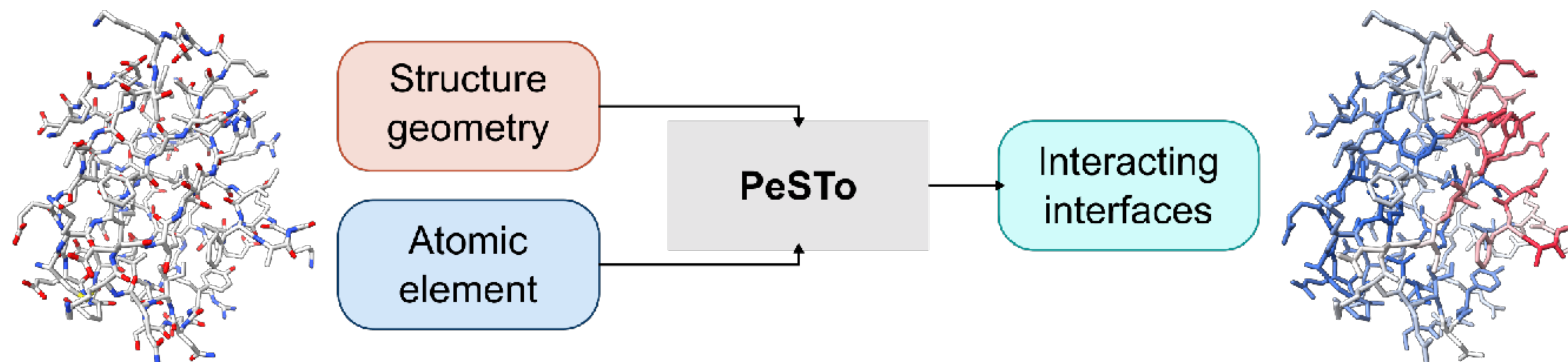
Antibody-Antigen
(3MXW)
Sonic hedgehog
bound to the 5E1
fab fragment



PeSTo: Protein Structure Transformer



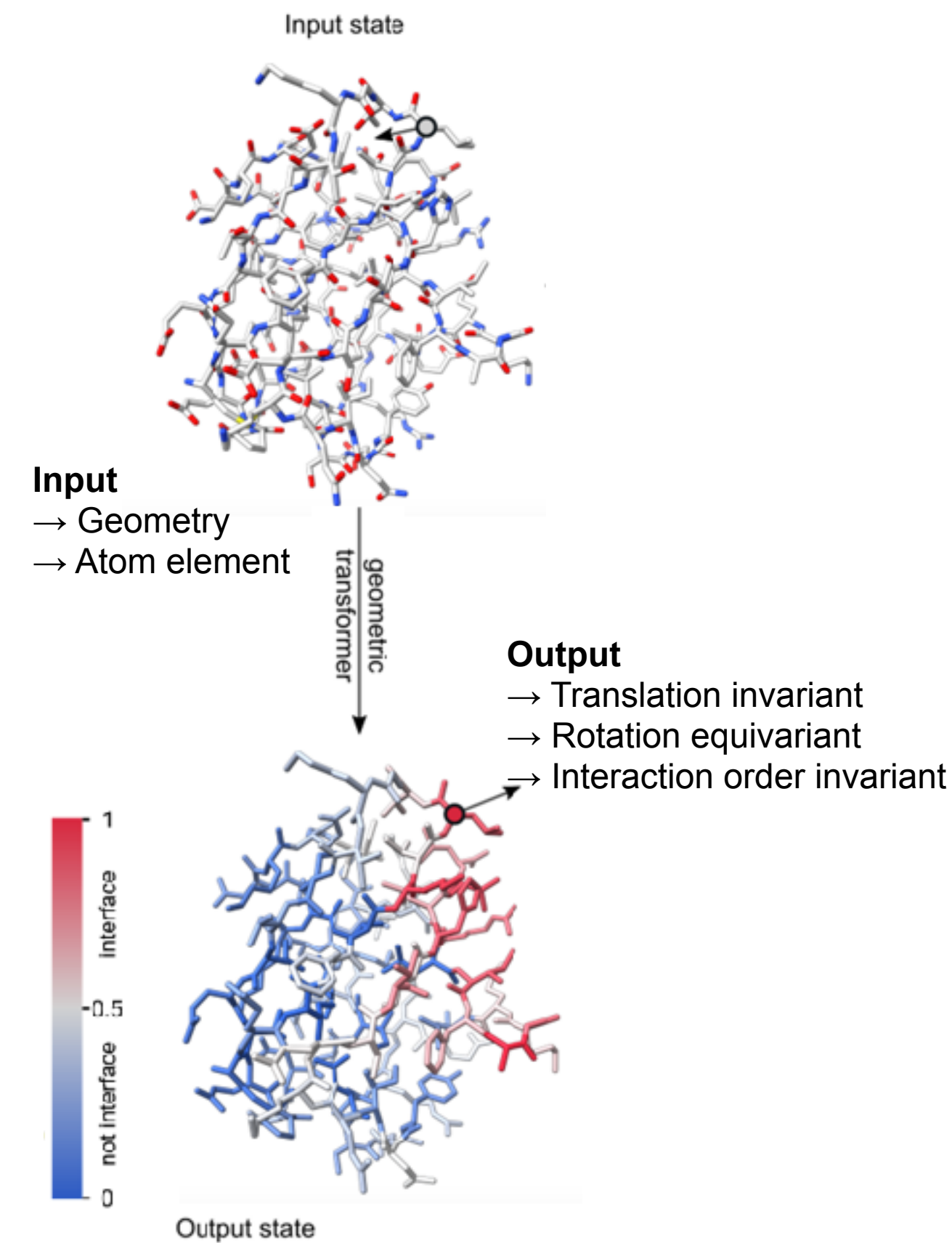
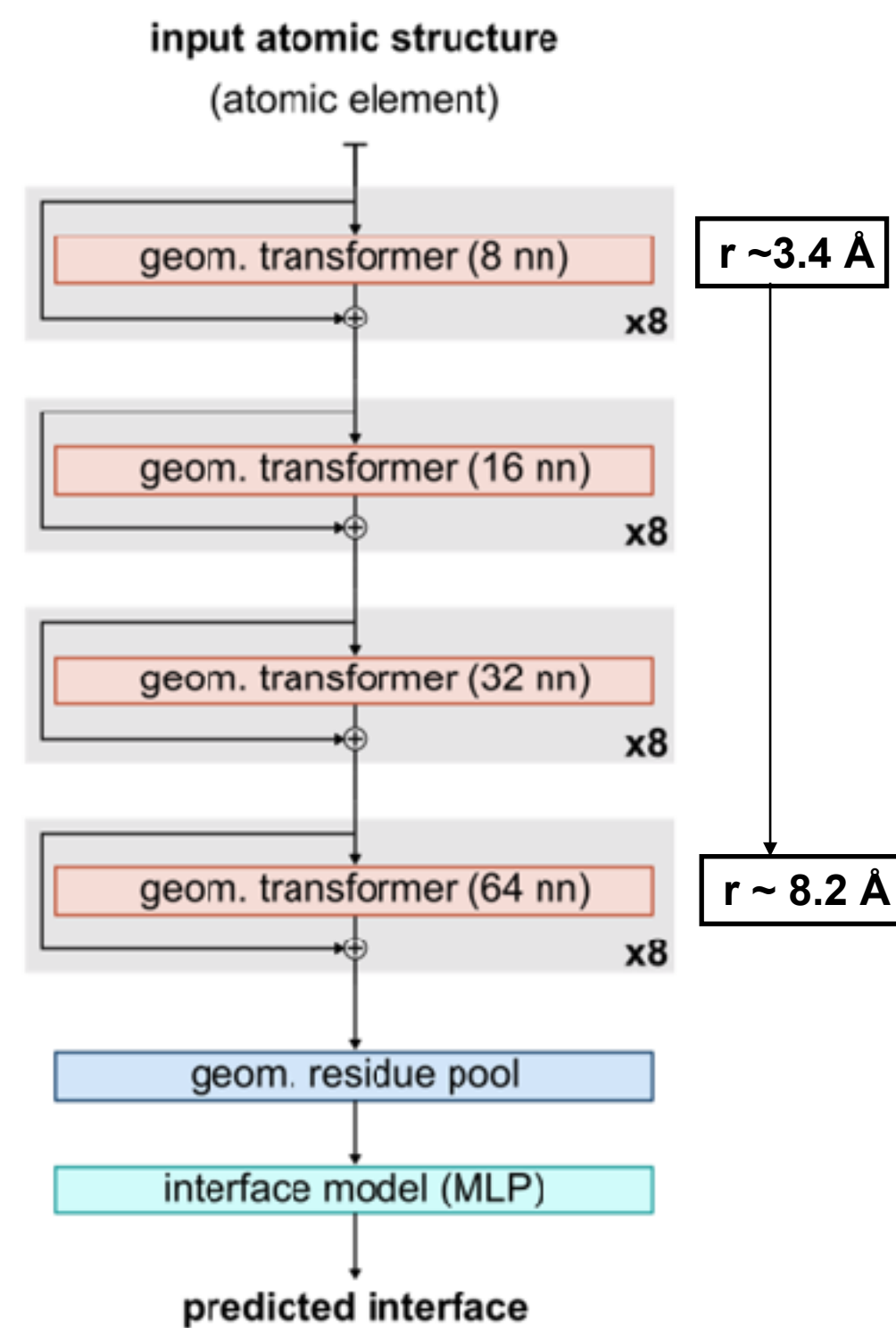
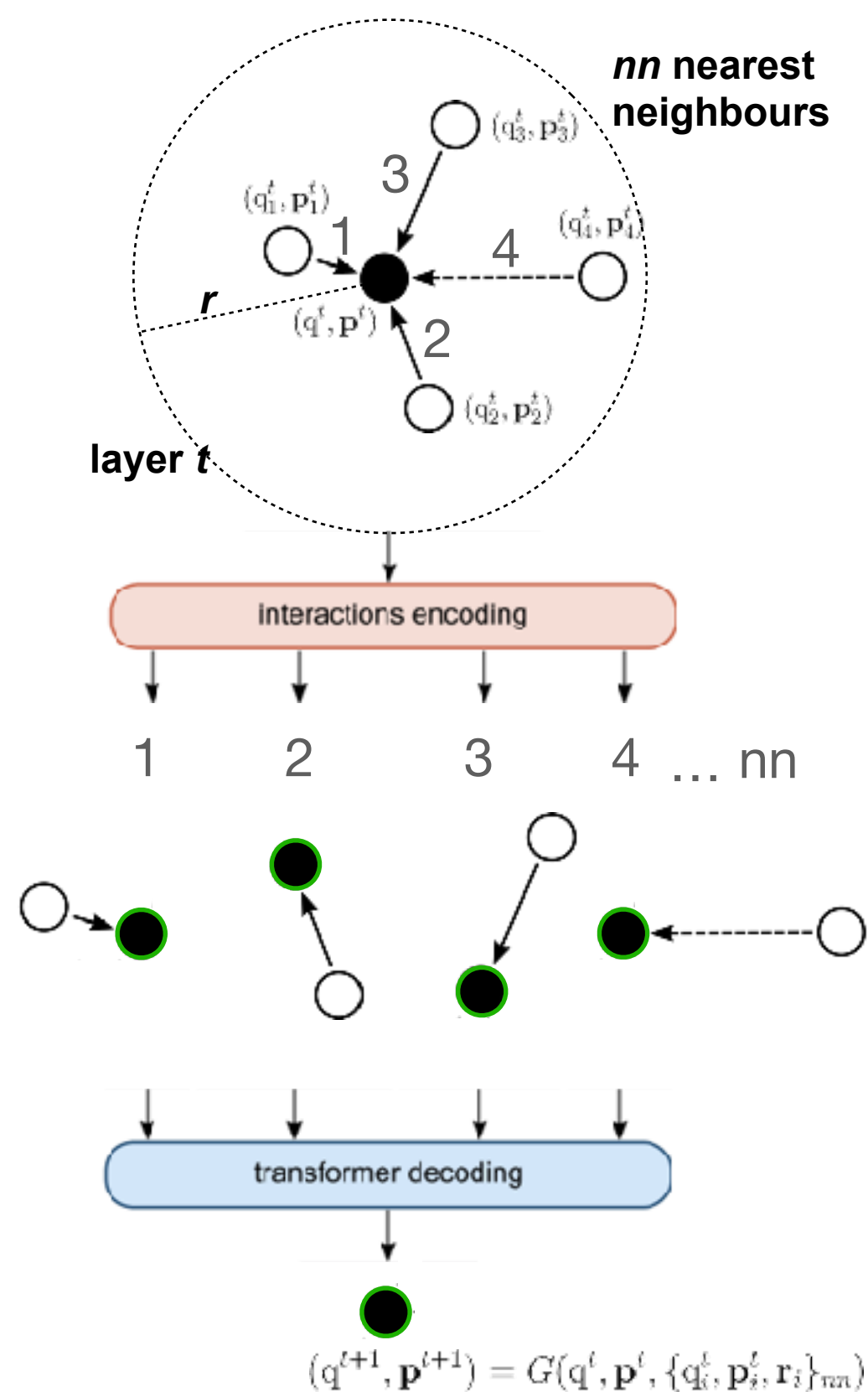
Lucien Krapp



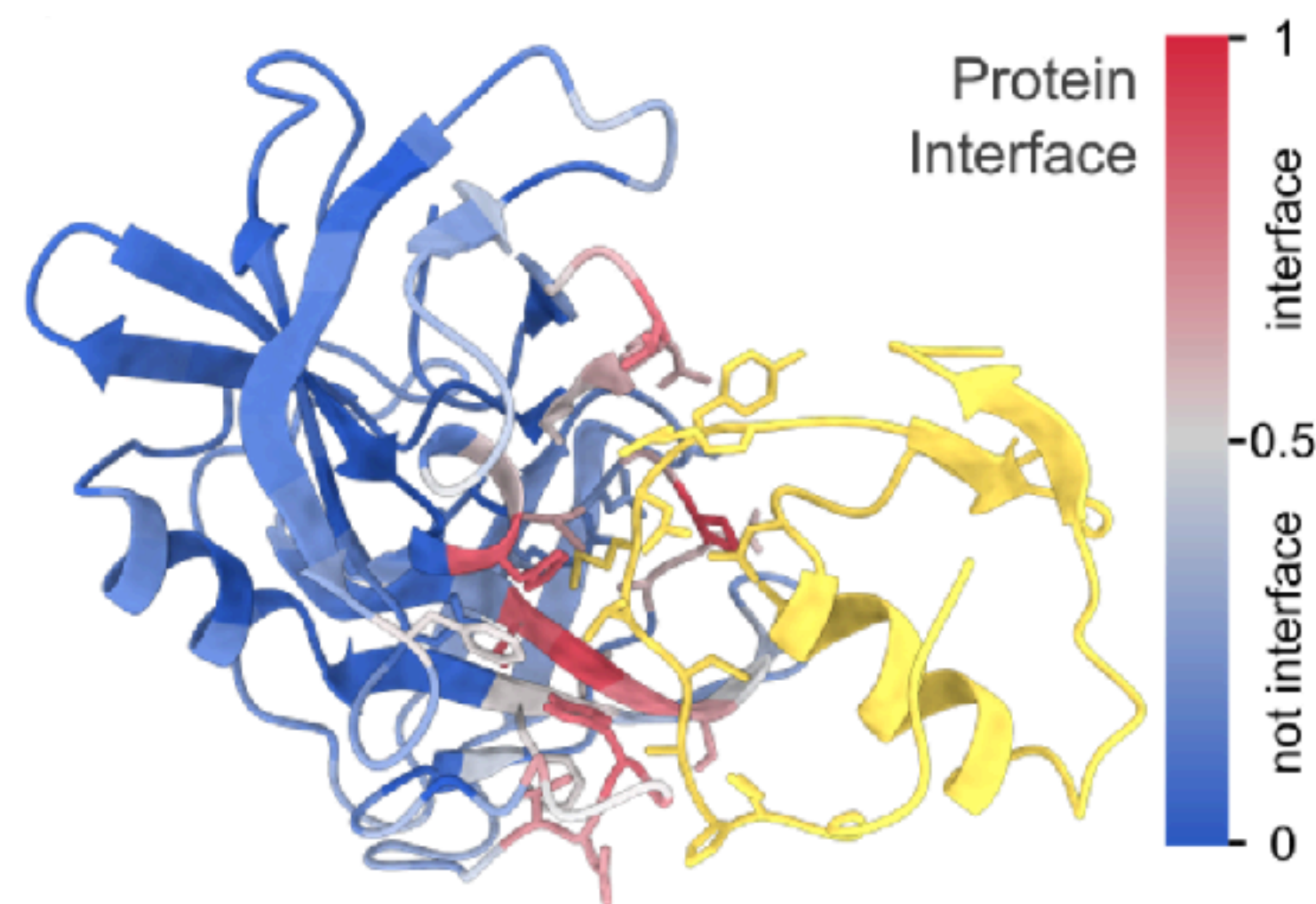
PeSTo: parameter-free geometric deep learning for accurate prediction of protein binding interfaces
Krapp L. F. et al. Nature Communications, 2023



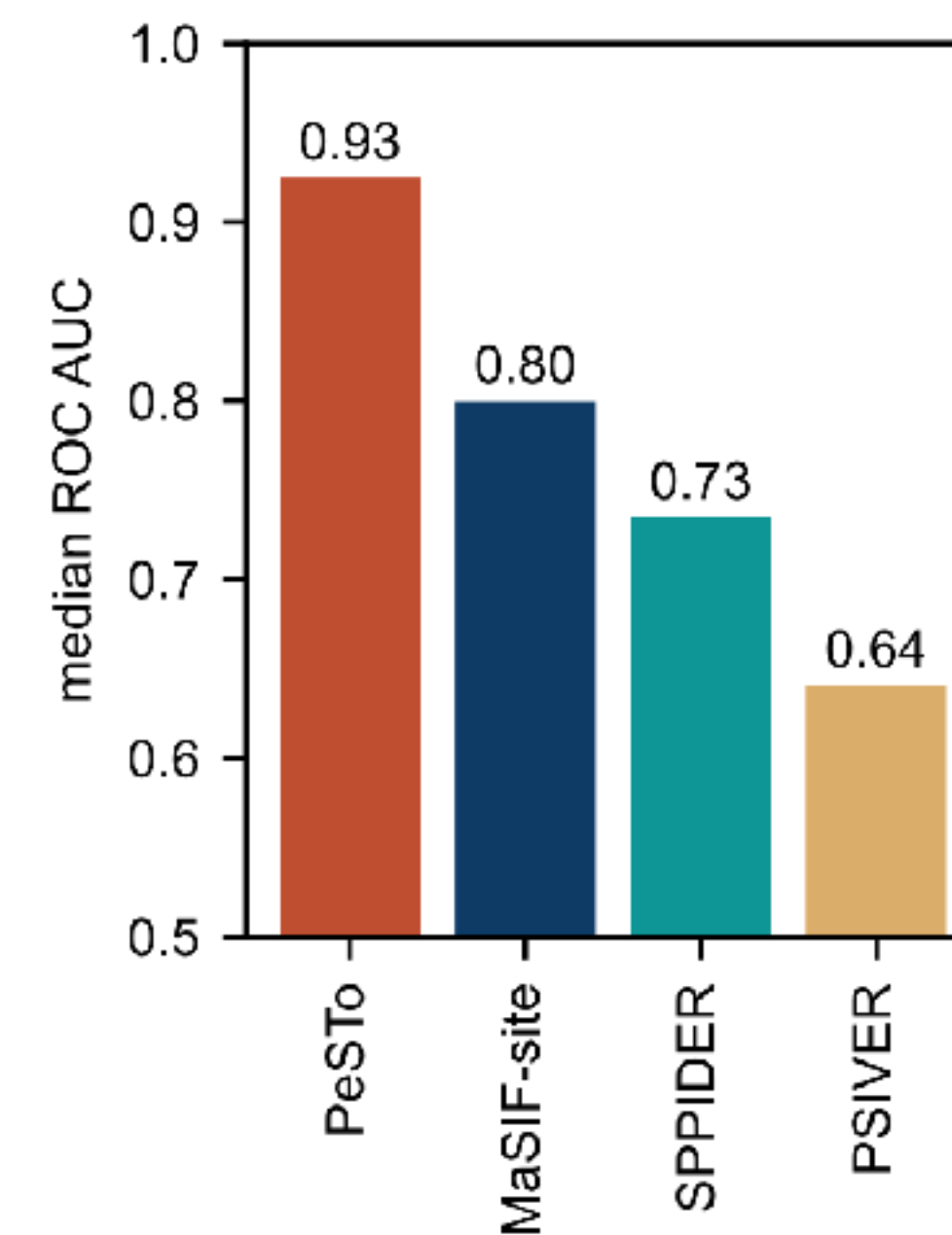
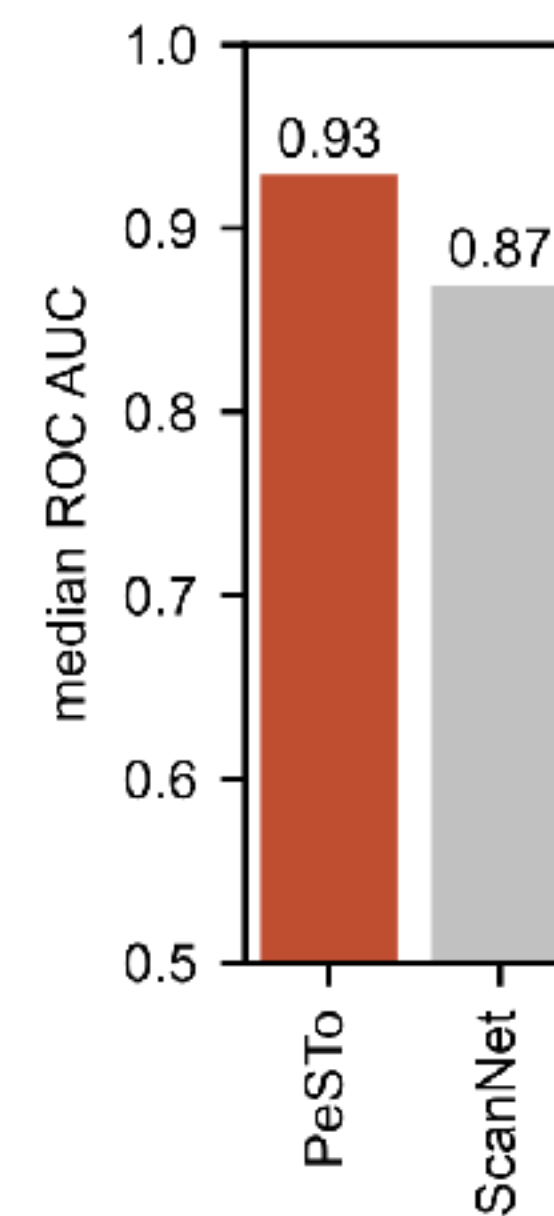
PeSTo: Protein Structure Transformer



PeSTo for protein-protein interfaces prediction ... comparison with state of the art

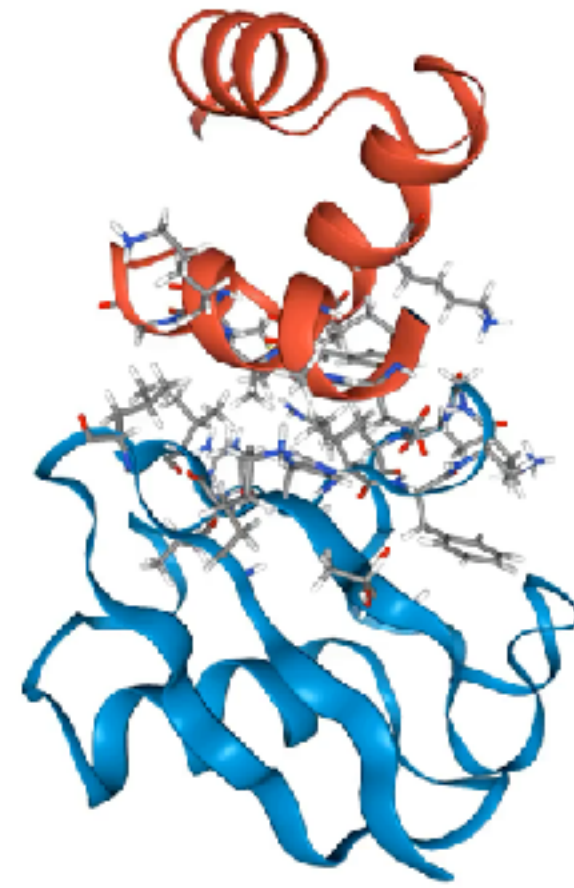


Streptogrisin B with ovomucoid - unbound conformation (0.93 Å RMSD) with a ROC AUC of 96%

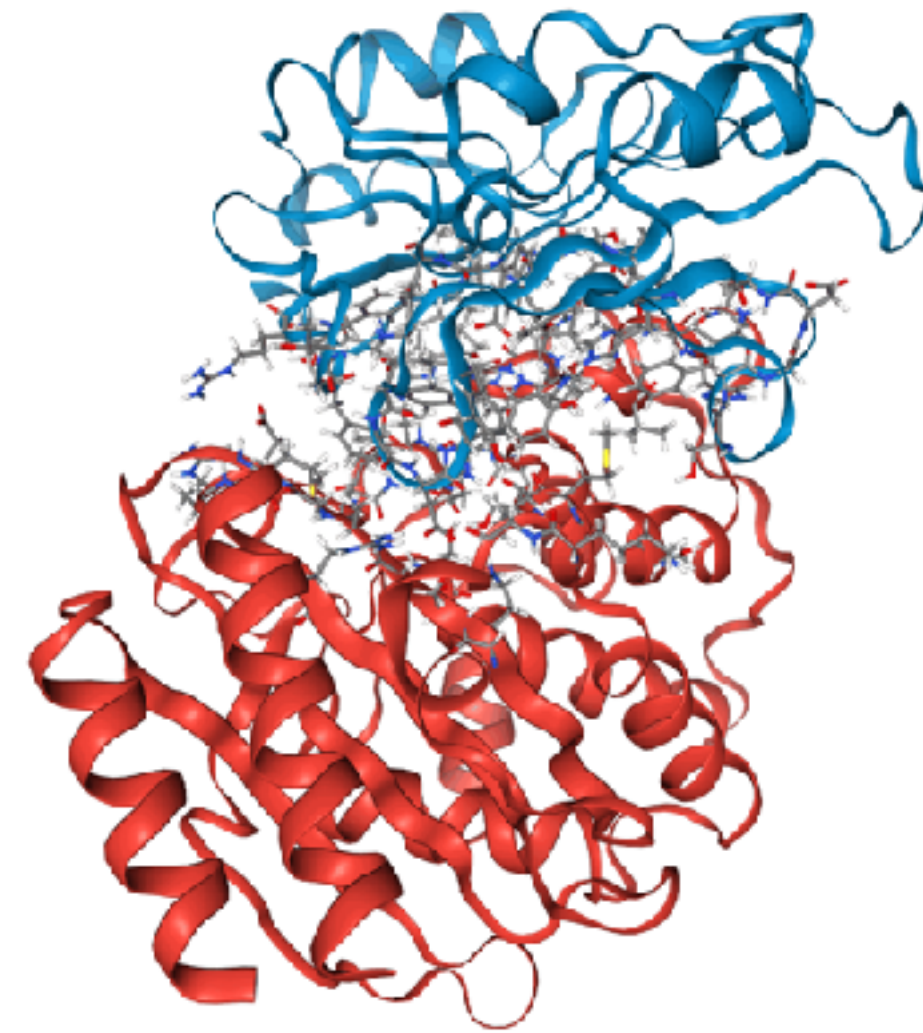


Gainza P. et al., *Deciphering interaction fingerprints from protein molecular surfaces using geometric deep learning*, Nature Methods, 2020
Tubiana J. et al., *ScanNet: an interpretable geometric deep learning model for structure-based protein binding site prediction*, Nature Methods, 2022

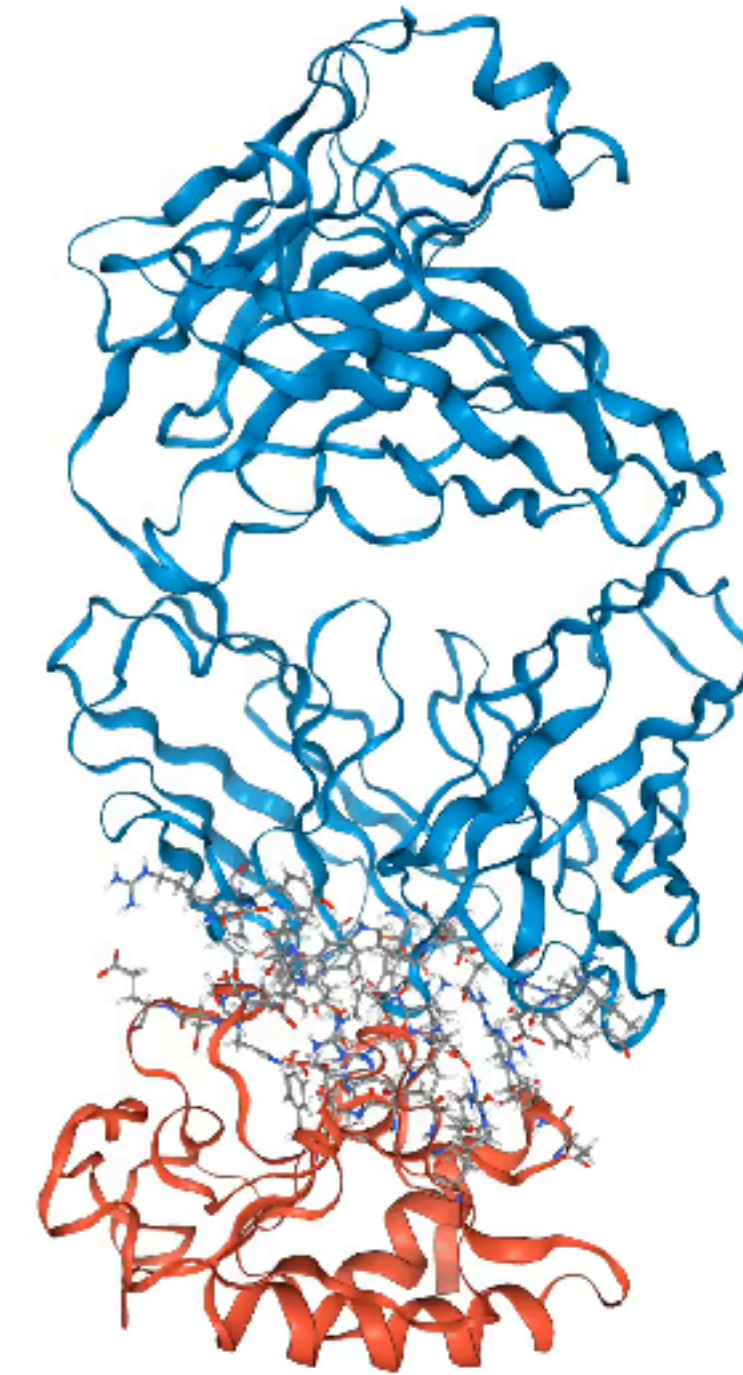
Interfaces are dynamic



Enzyme-Substrate
(2OOB)
Ubiquitin bound to
ubiquitin ligase



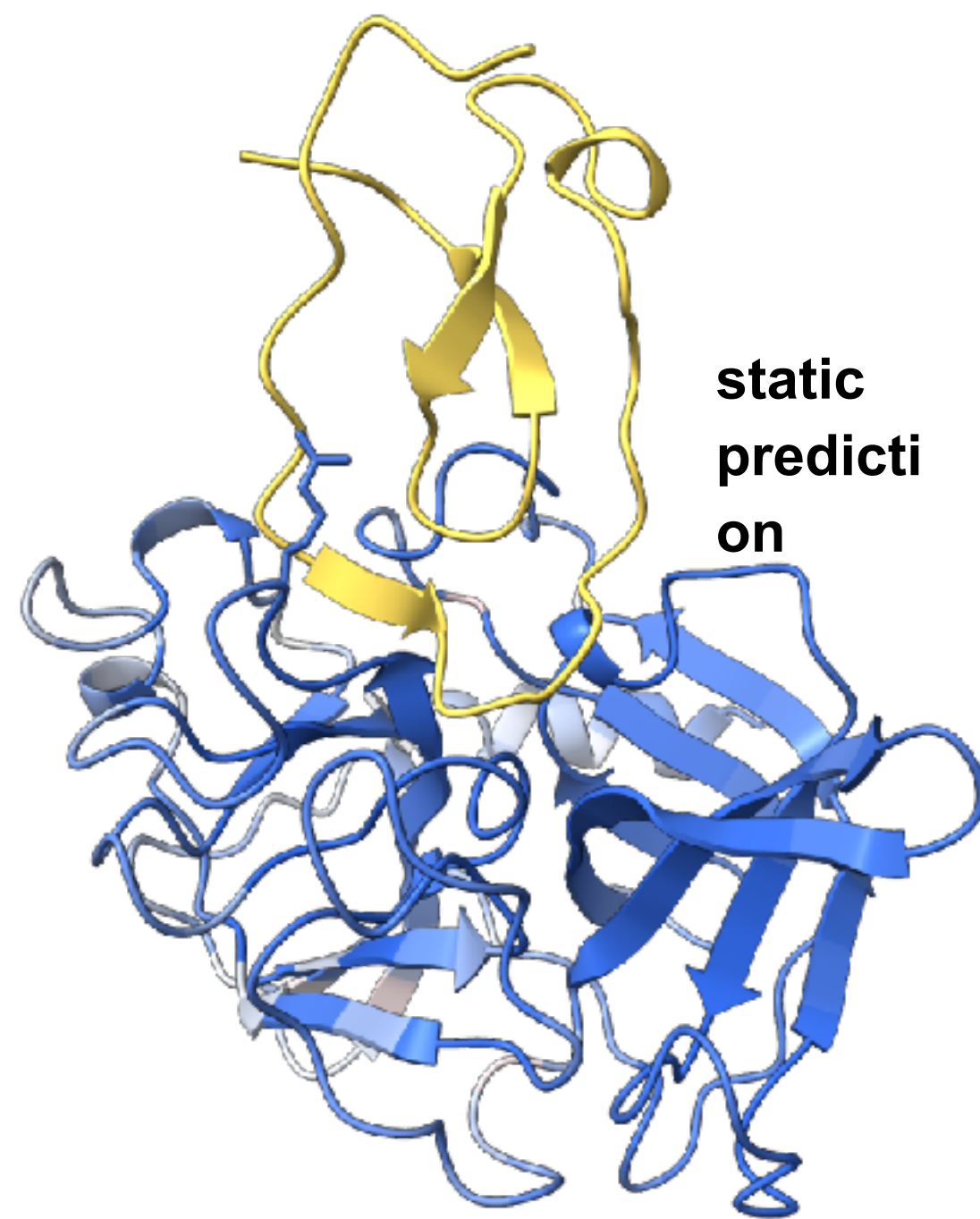
Enzyme-Inhibitor (1JTG)
Beta-lactamase bound to
beta-lactamase inhibitor



Antibody-Antigen
(3MXW)
Sonic hedgehog
bound to the 5E1
fab fragment

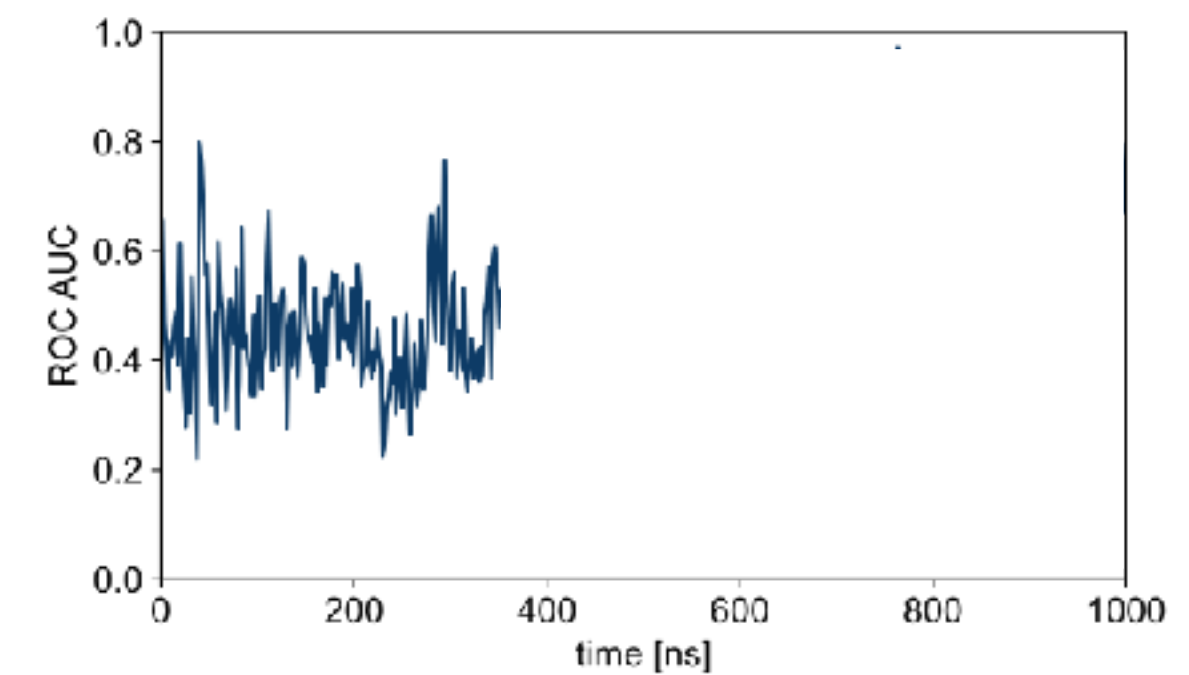
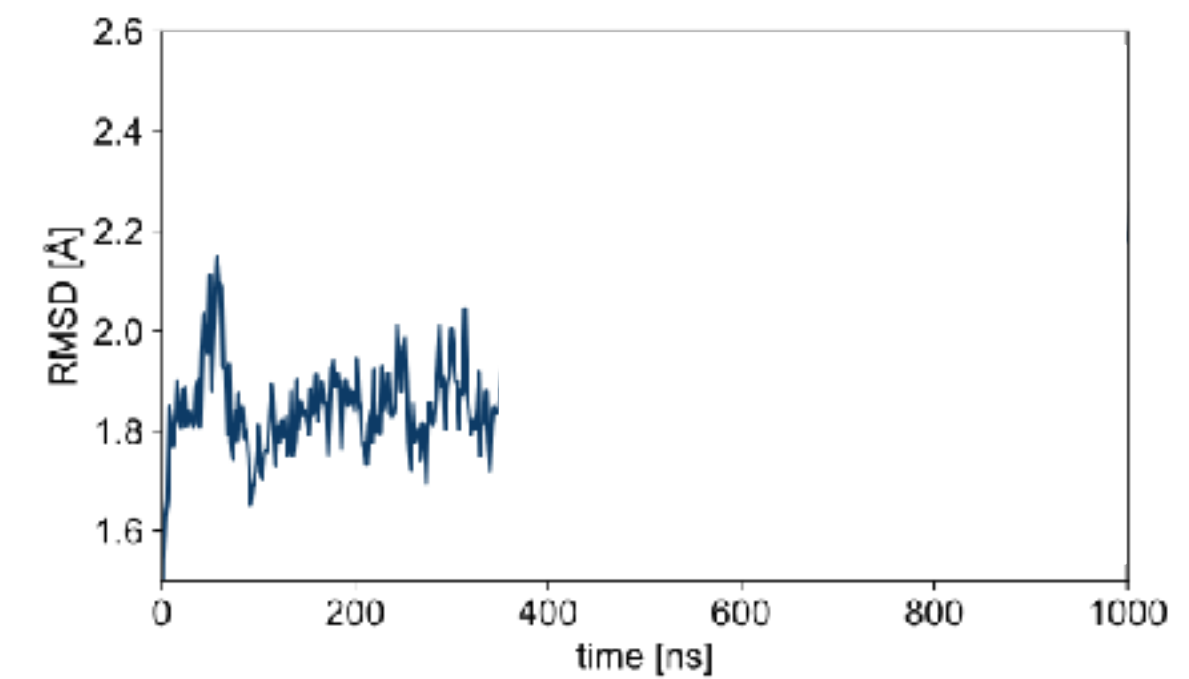
Interface recovery through molecular simulation

Elafin complexed with porcine pancreatic elastase (1FLE)



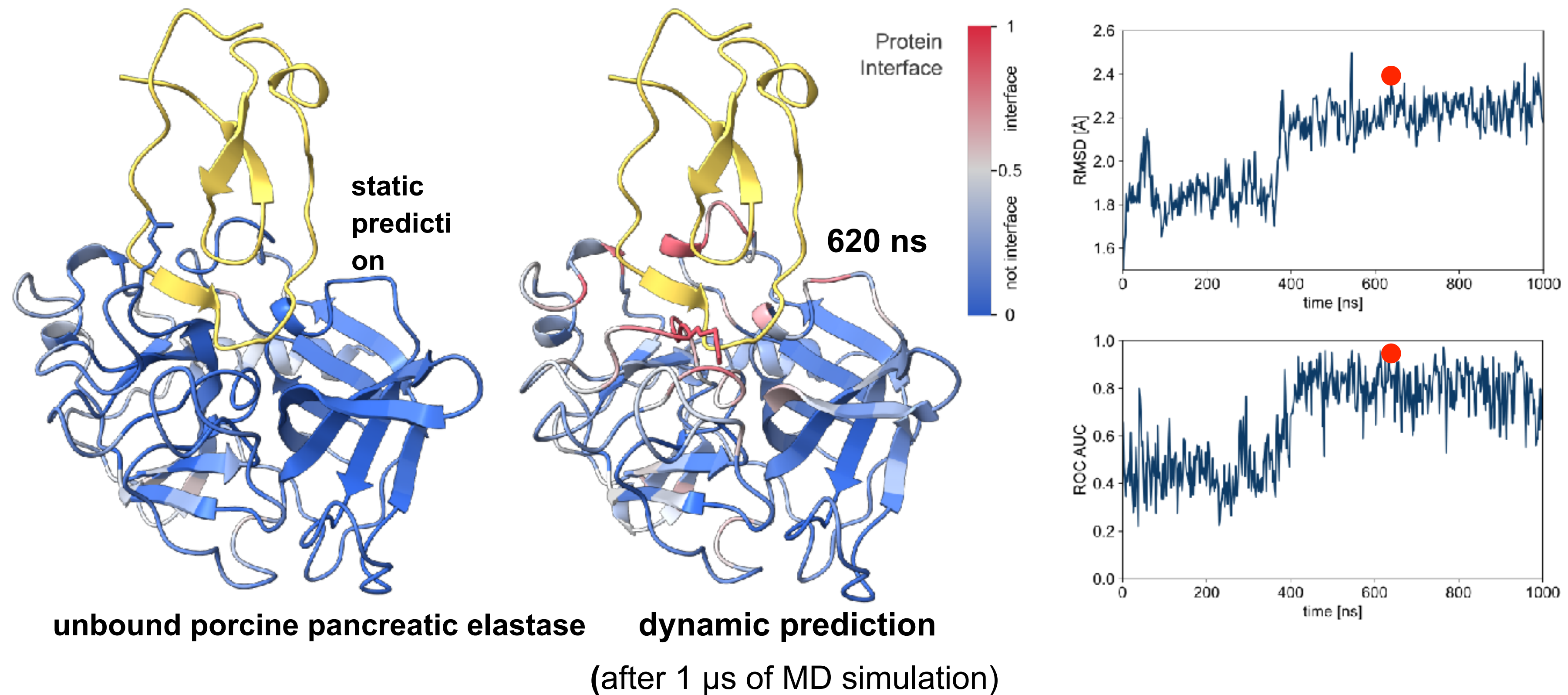
static
predicti
on

unbound porcine pancreatic elastase



Interface recovery through molecular simulation

Elafin complexed with porcine pancreatic elastase (1FLE)

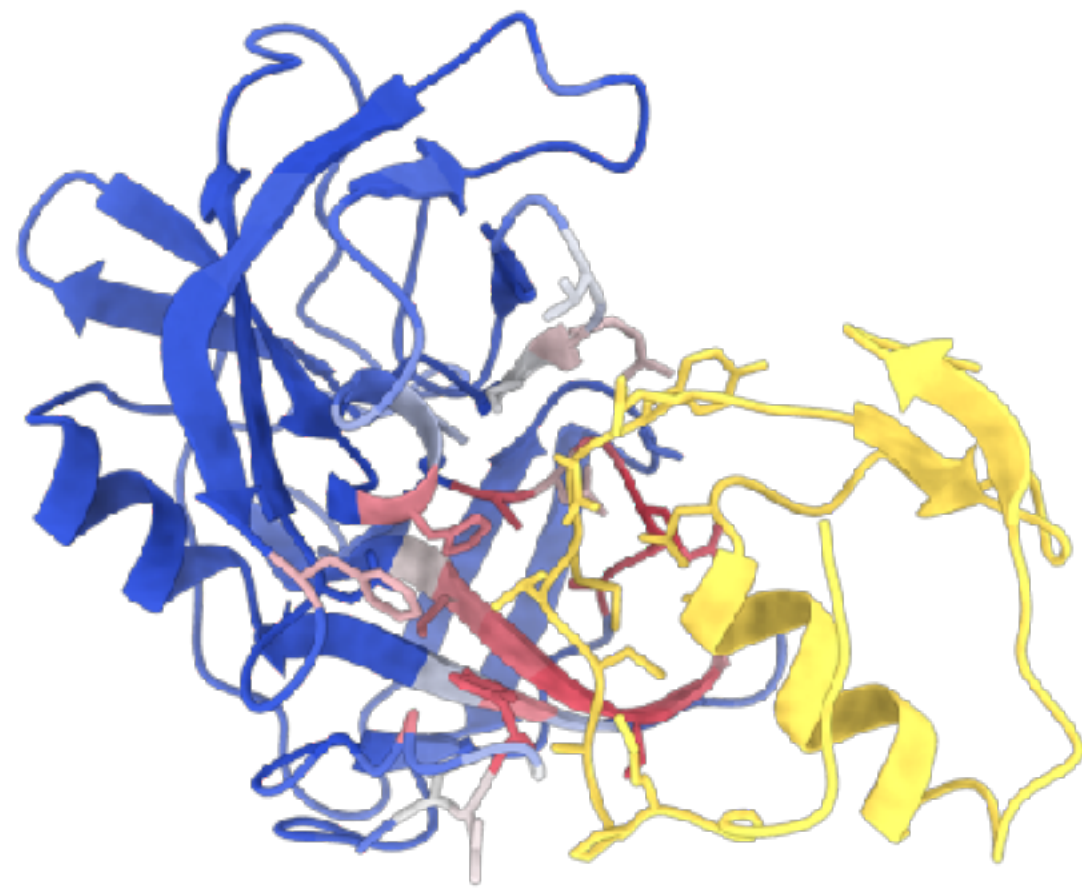


revealing cryptic binding sites
and potential allosteric mechanisms

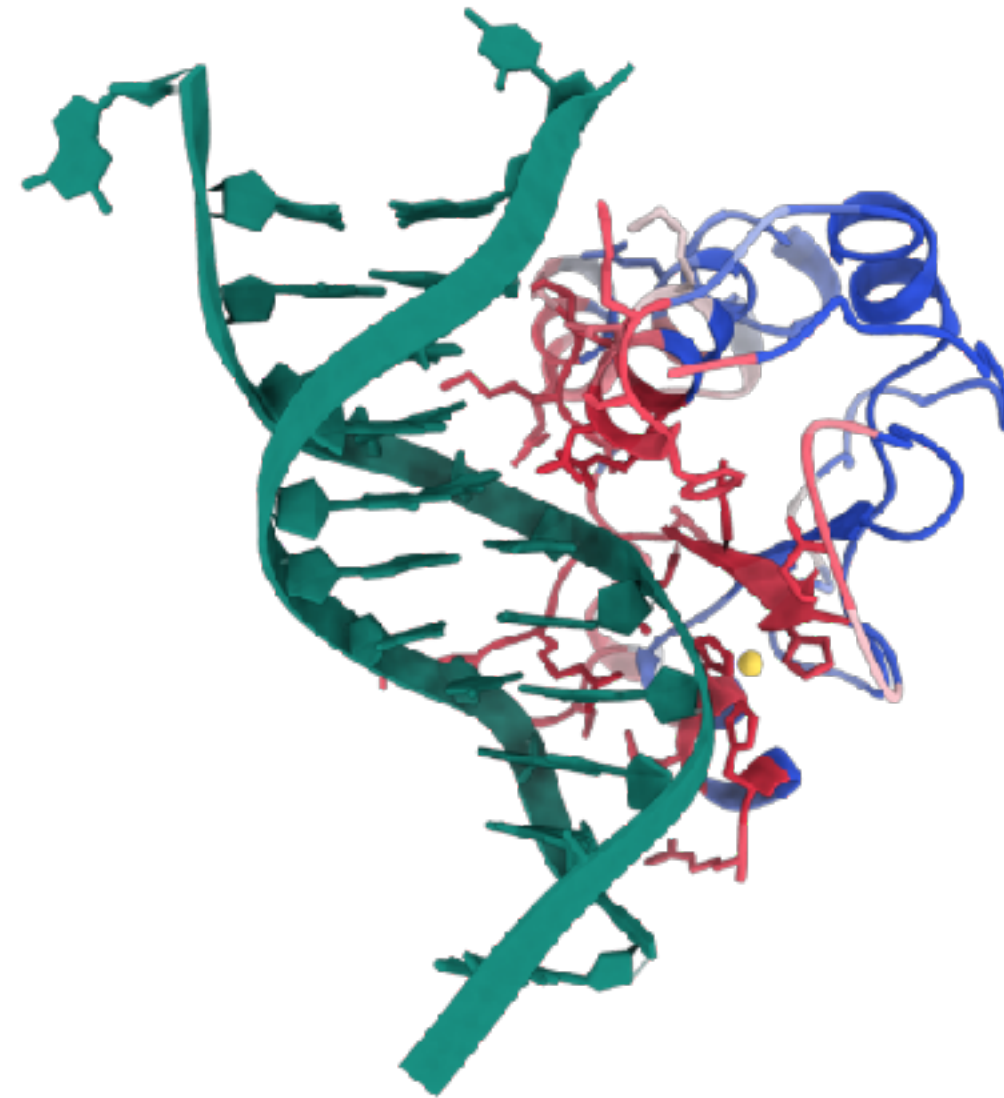


Other binding interfaces of proteins

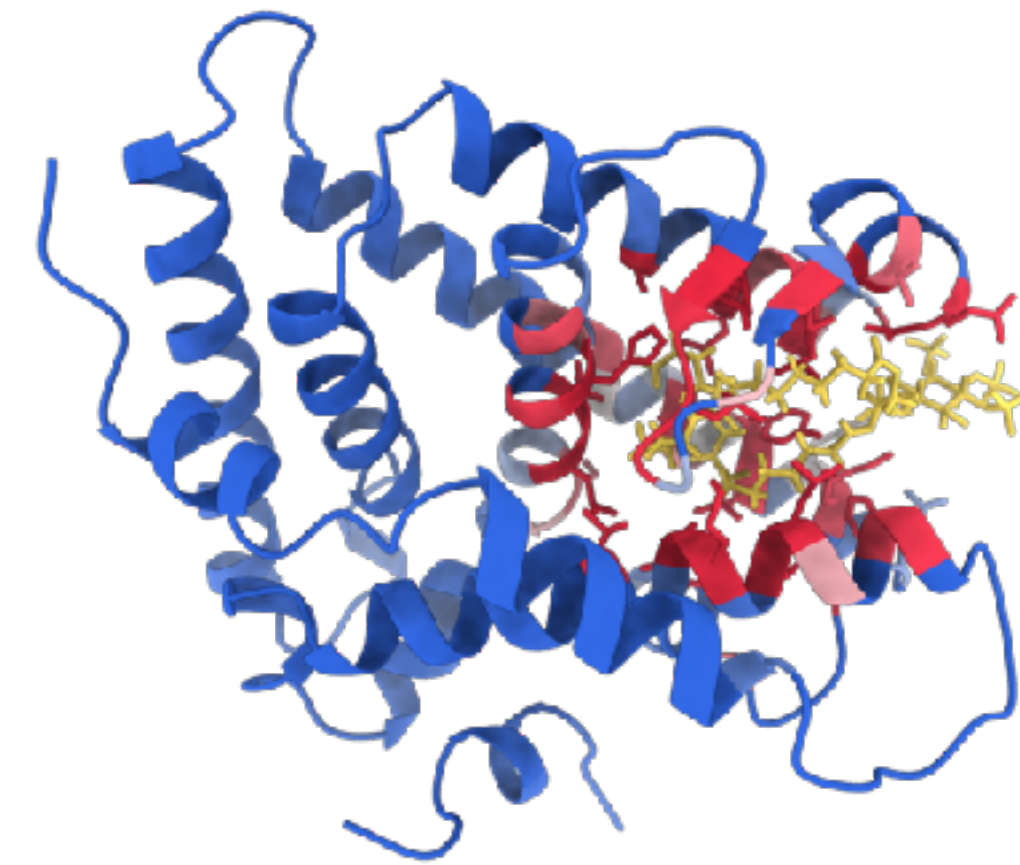
Protein-protein



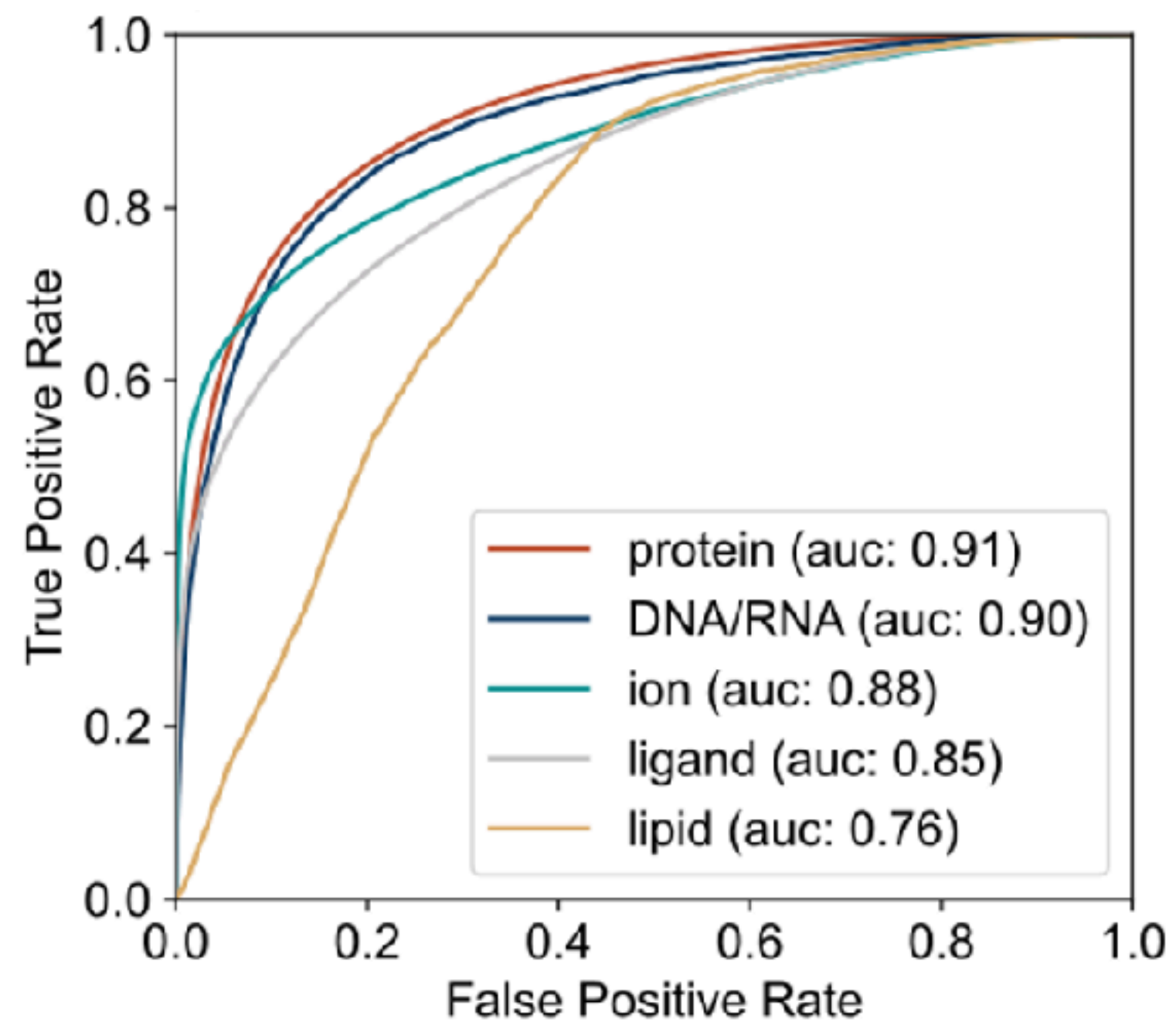
Protein-DNA-ions



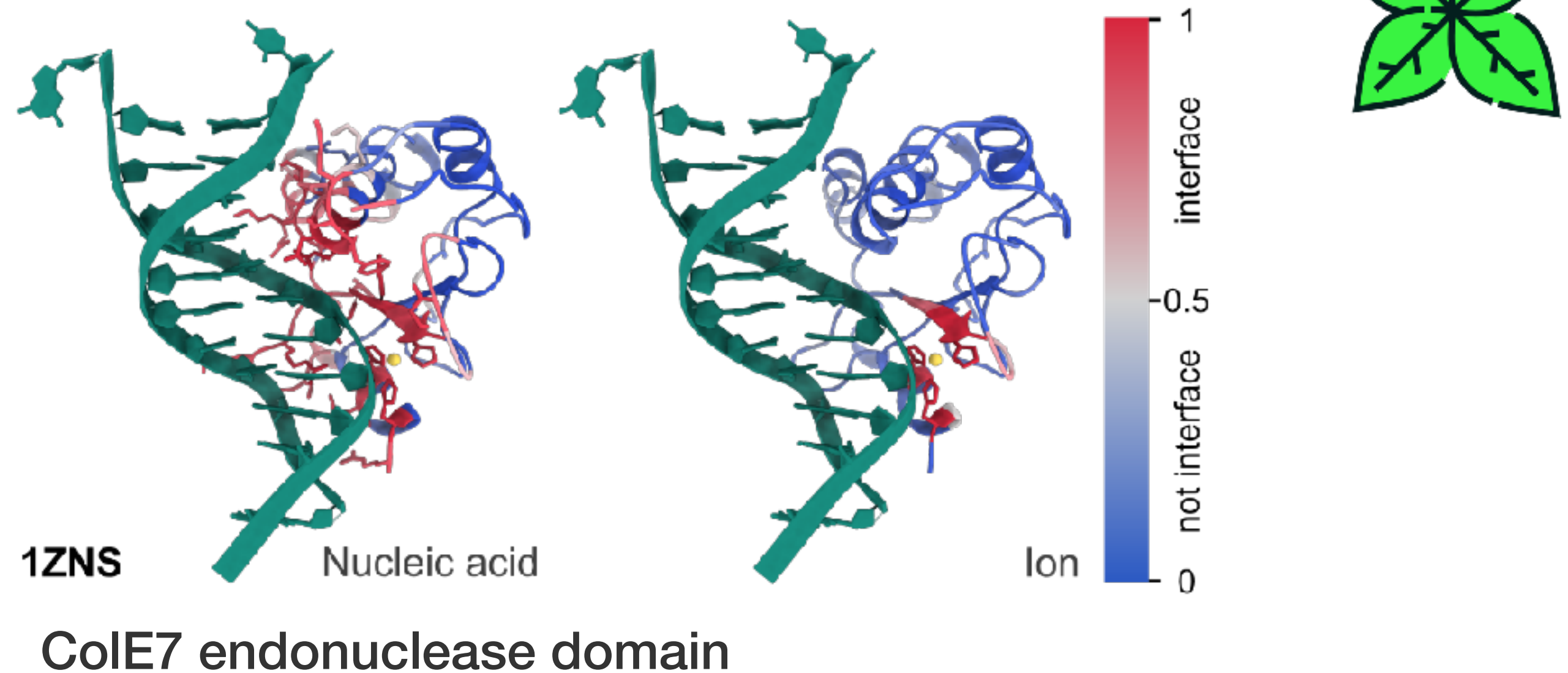
Protein-lipid



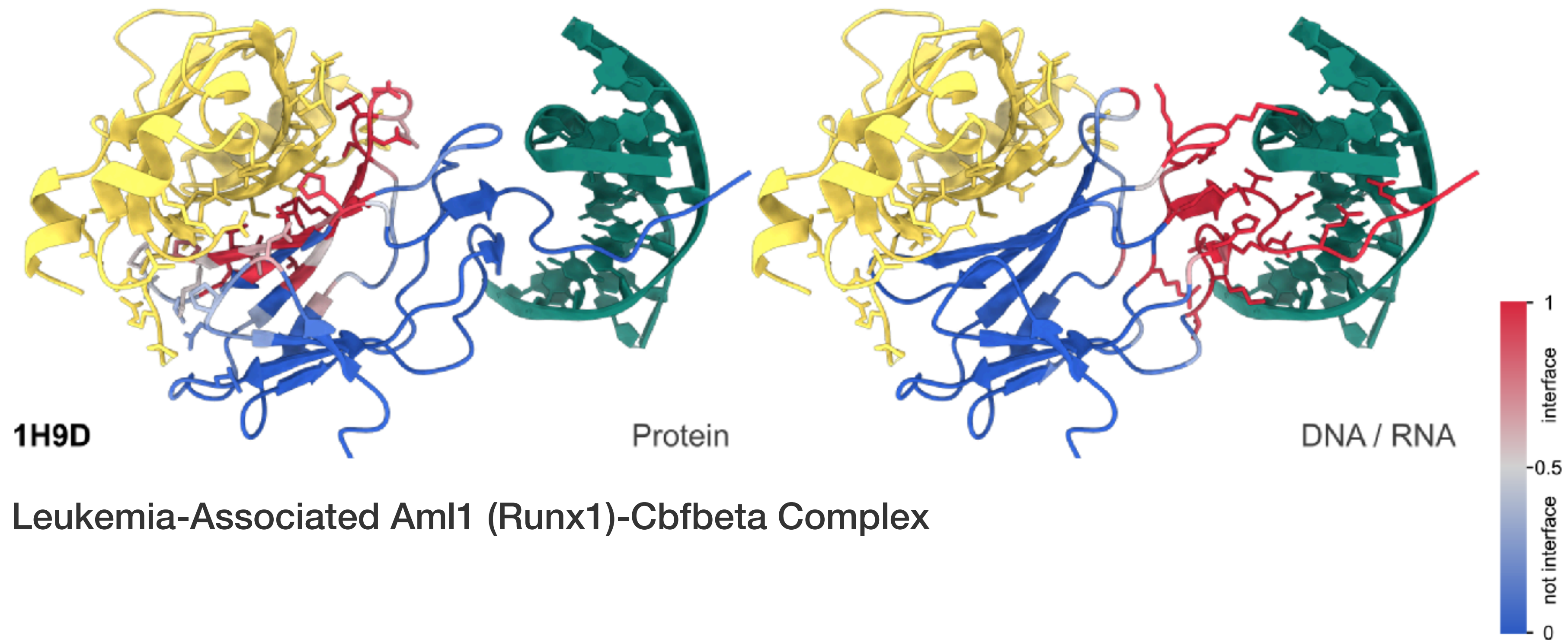
Extending prediction to other interacting interfaces



nucleic acid & ion interface



Protein-nucleic acid interface predictions



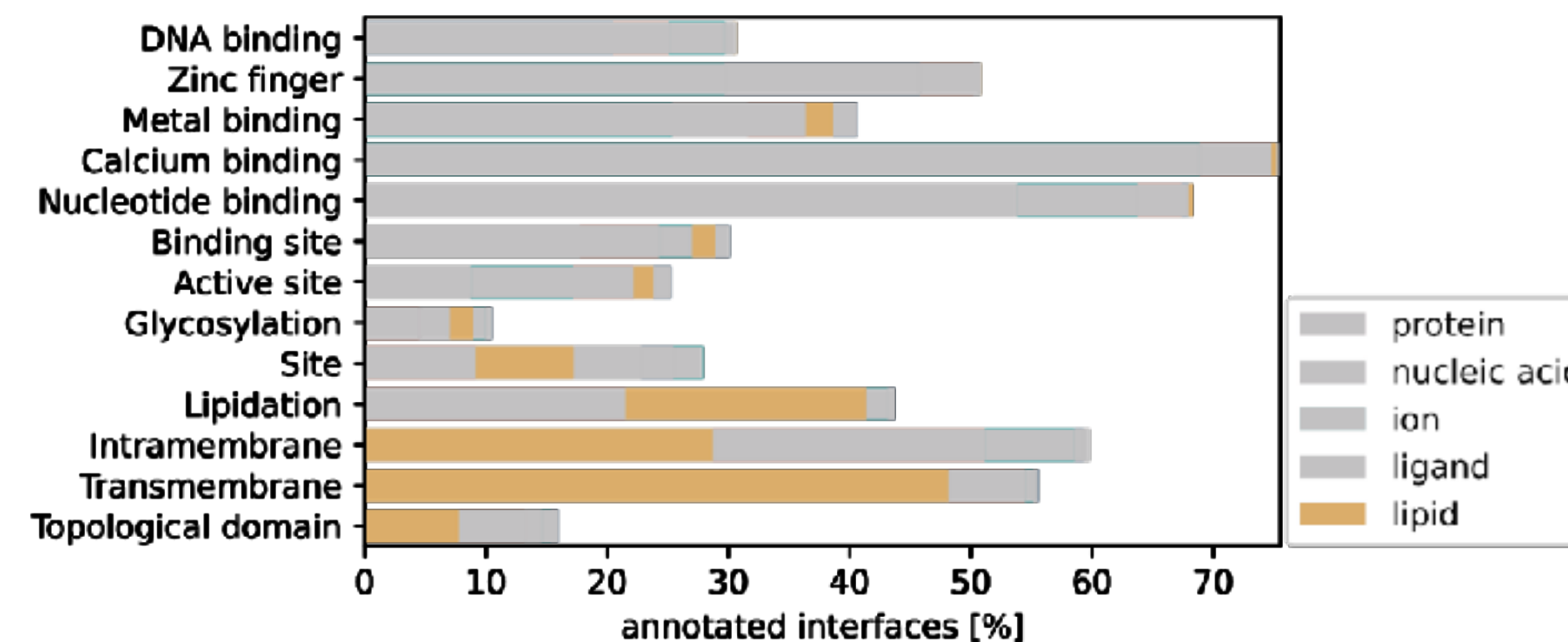
Leukemia-Associated Aml1 (Runx1)-Cbfbeta Complex

Proteome-wide interface prediction - aka *interfaceome*

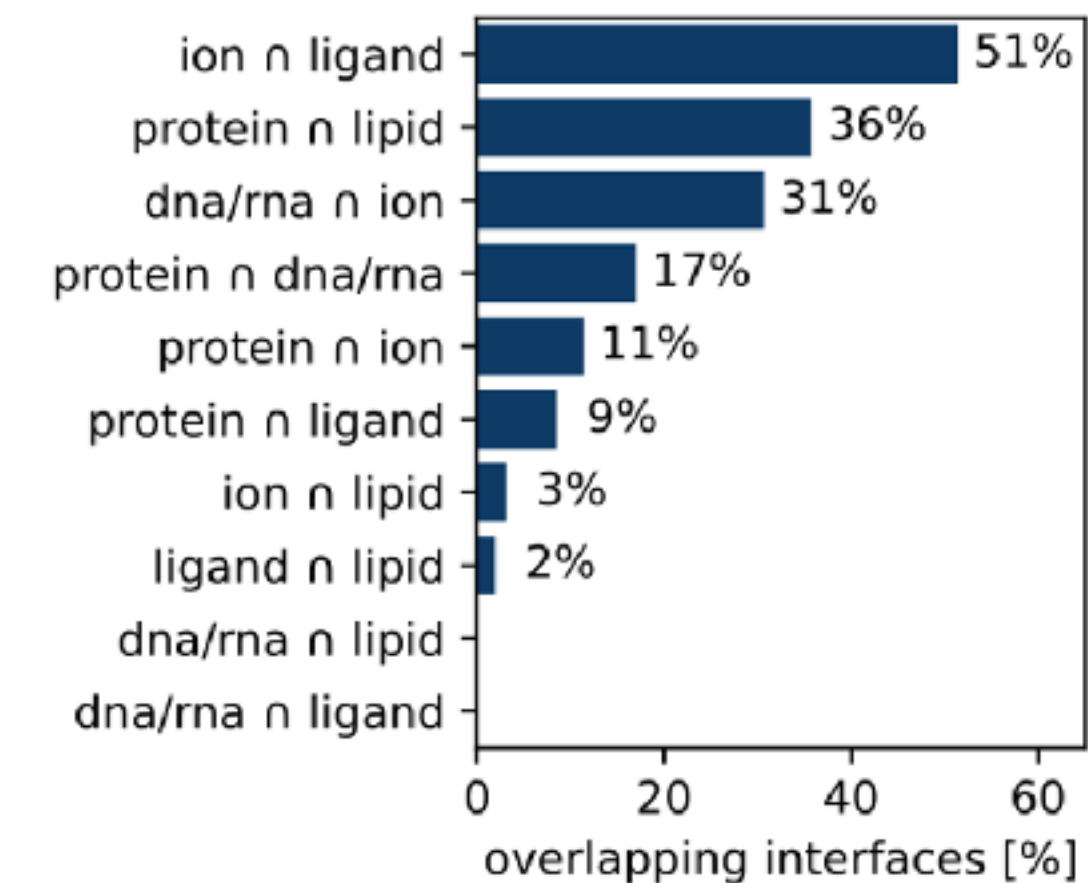
Human proteome prediction using the AF-EBI database

(only high-quality models, 7464 of 20504)

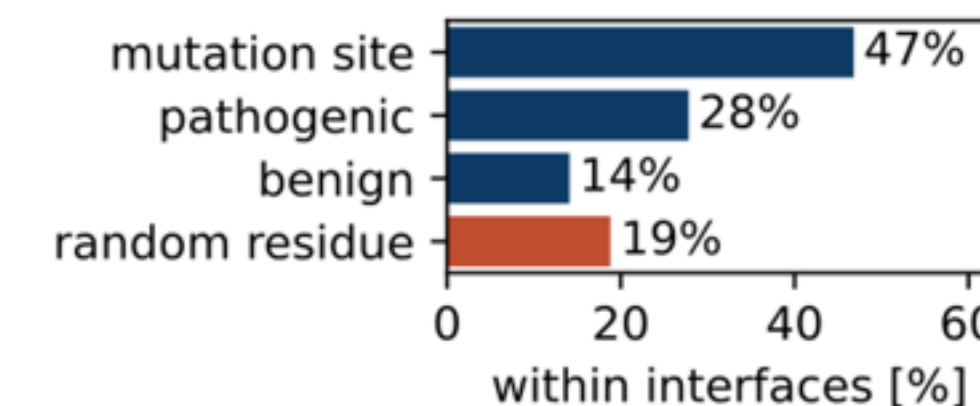
correlation with protein function and features



interface cross-talk



correlation with protein mutations



@ <http://pesto.epfl.ch>

Fabio Cortés



PeSTo

PeSTo (Protein Structure Transformer) is a parameter-free geometric deep learning method to predict protein interaction interfaces from a protein structure. It is available for free without registration as an online tool. A manuscript of the method is in preparation and will be available soon.

Learn more about this project in this [preprint](#) at [Biorxiv](#).

How to use

Copy-paste your atomic coordinates in PDB format, or upload a PDB file from your drive, or fetch a protein structure/model from:

- The protein data bank by typing a PDB ID. Example: 2CUA
- The AlphaFold-EBI database by typing a Uniprot ID. Example: P27695
- Upload your own PDB formatted structure

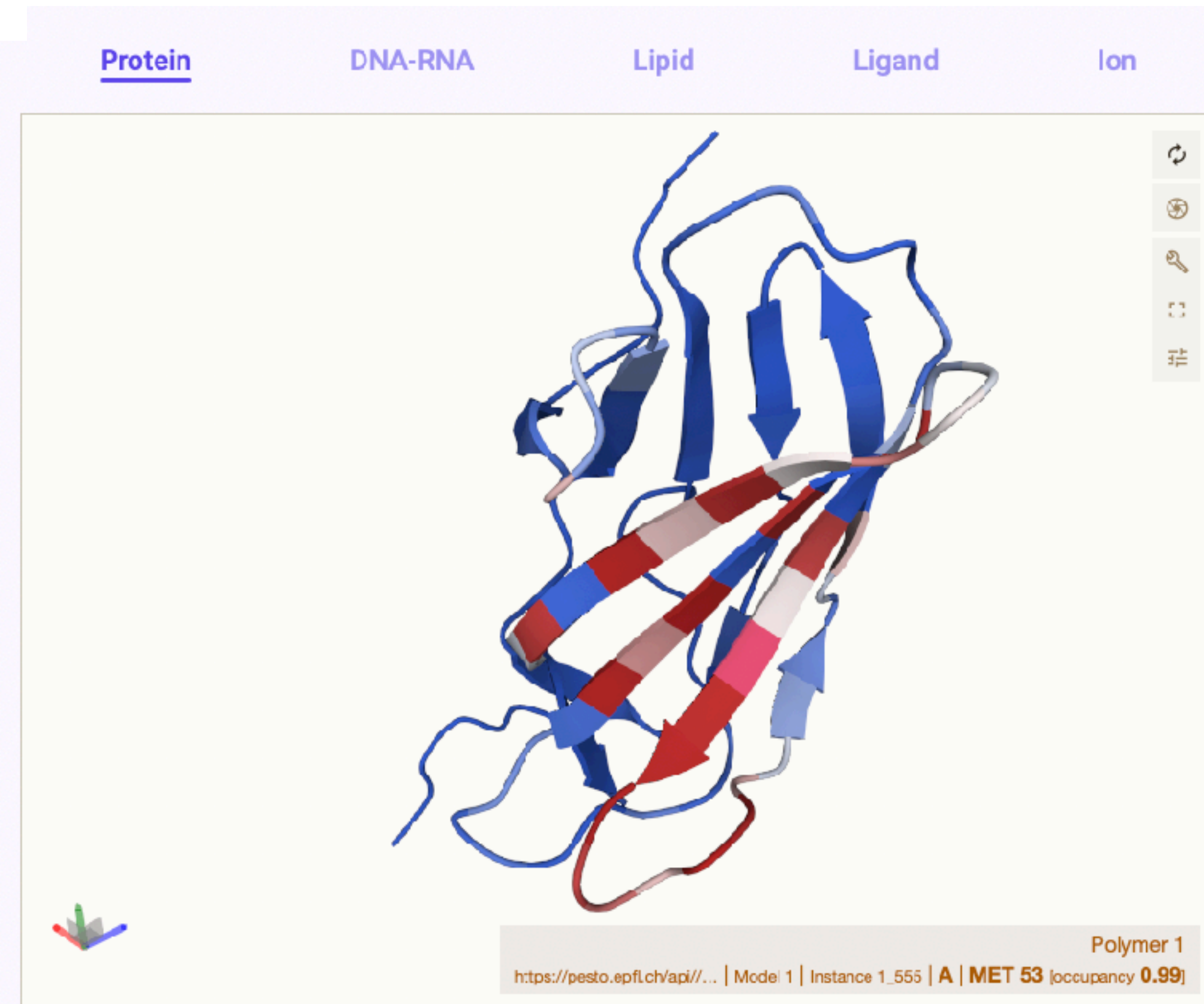
Then click "Detect chains", select one or more, and submit your job to run the prediction. Your results should be available in less than a minute. If an error occurs, the PDB file might be not correctly formatted or the input structure is too big

2CUA

Fetch PDB/AF-EBI

Upload PDB

Copy-Paste molecule here



Chain	Res Name	Res ID	Prediction
A:0	ASP	13	0.64
A:0	THR	51	0.91
A:0	VAL	52	0.54
A:0	MET	53	0.99
A:0	ALA	54	0.96
A:0	GLY	55	0.97
A:0	ASN	56	0.98
A:0	ASP	57	0.91

Paths to drug discovery and the role of computational chemistry

1. The process of Drug Discovery (& Development)

2. Structure-based Drug Design

courtesy of Marco De Vivo, PhD

Laboratory of Molecular Modeling and Drug Discovery
Istituto Italiano di Tecnologia- Genoa, Italy

Drug Discovery & Development

FDA approvals in 2020 - *Some facts* -

53 New Drugs Approved

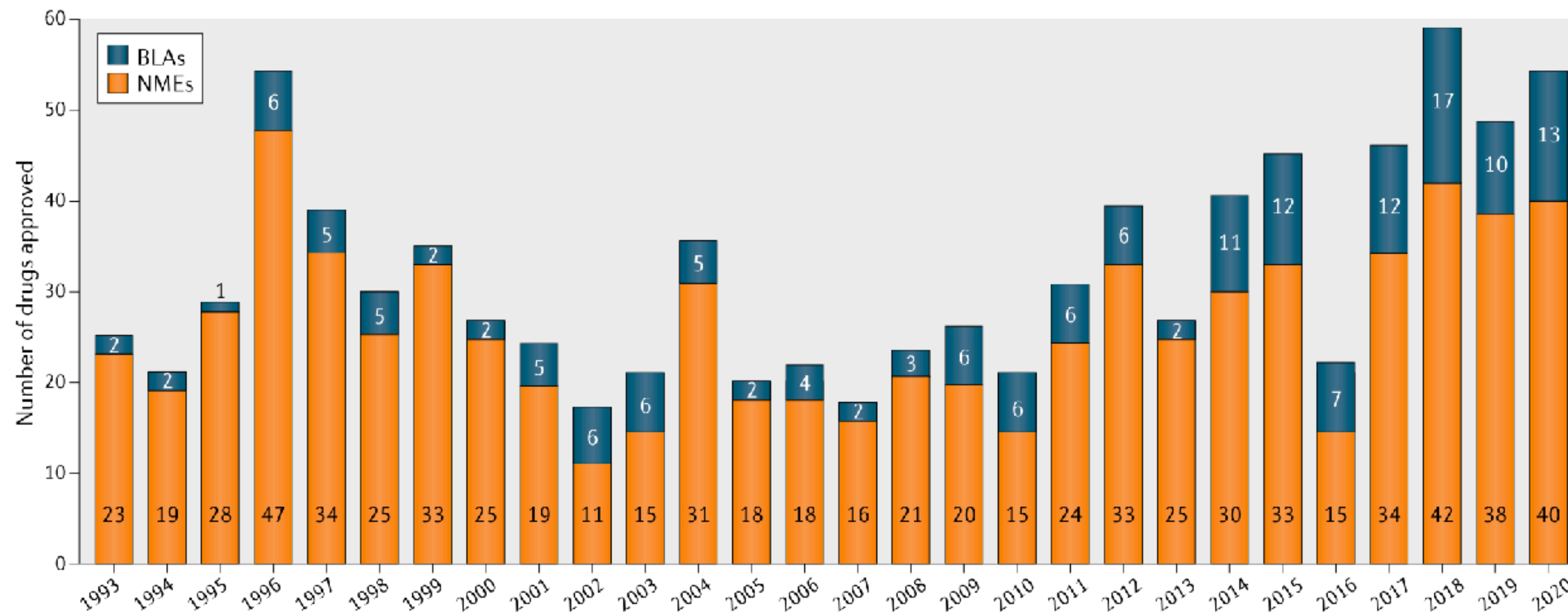


Fig. 1 | **Novel FDA approvals since 1993.** Annual numbers of new molecular entities (NMEs) and biologics license applications (BLAs) approved by the FDA's Center for Drug Evaluation and Research (CDER). See TABLE 1 for

new approvals in 2020. Approvals by the Center for Biologics Evaluation and Research (CBER), for products such as vaccines and gene therapies, are not included in this drug count (see TABLE 2). Source: FDA.

Source:

Asher Mullard,
Nature Reviews Drug Discovery, February (2021)

Drug Discovery & Development

- Some facts -

Approval by therapeutic area 2020

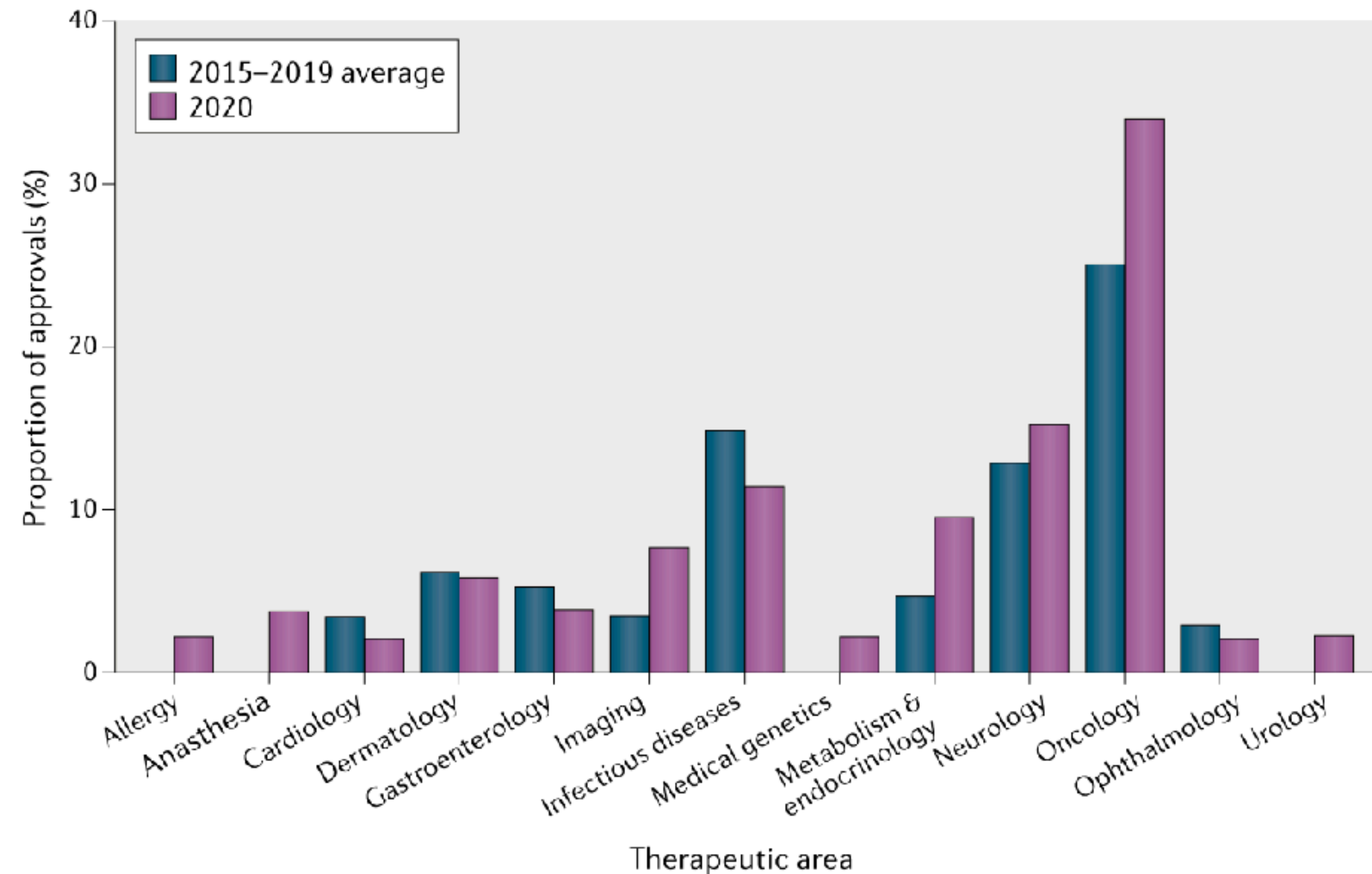


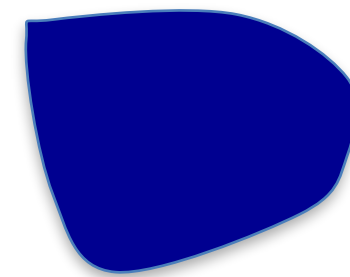
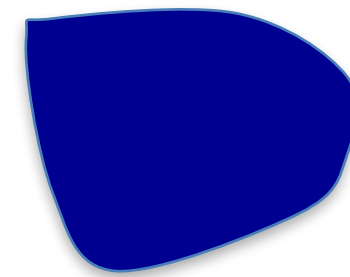
Fig. 2 | **CDER approvals by selected therapeutic areas.** Source: *Nature Reviews Drug Discovery*, FDA.

Source:

Asher Mullard,
Nature Reviews Drug Discovery, February (2021)

Enzymes and catalytic activity

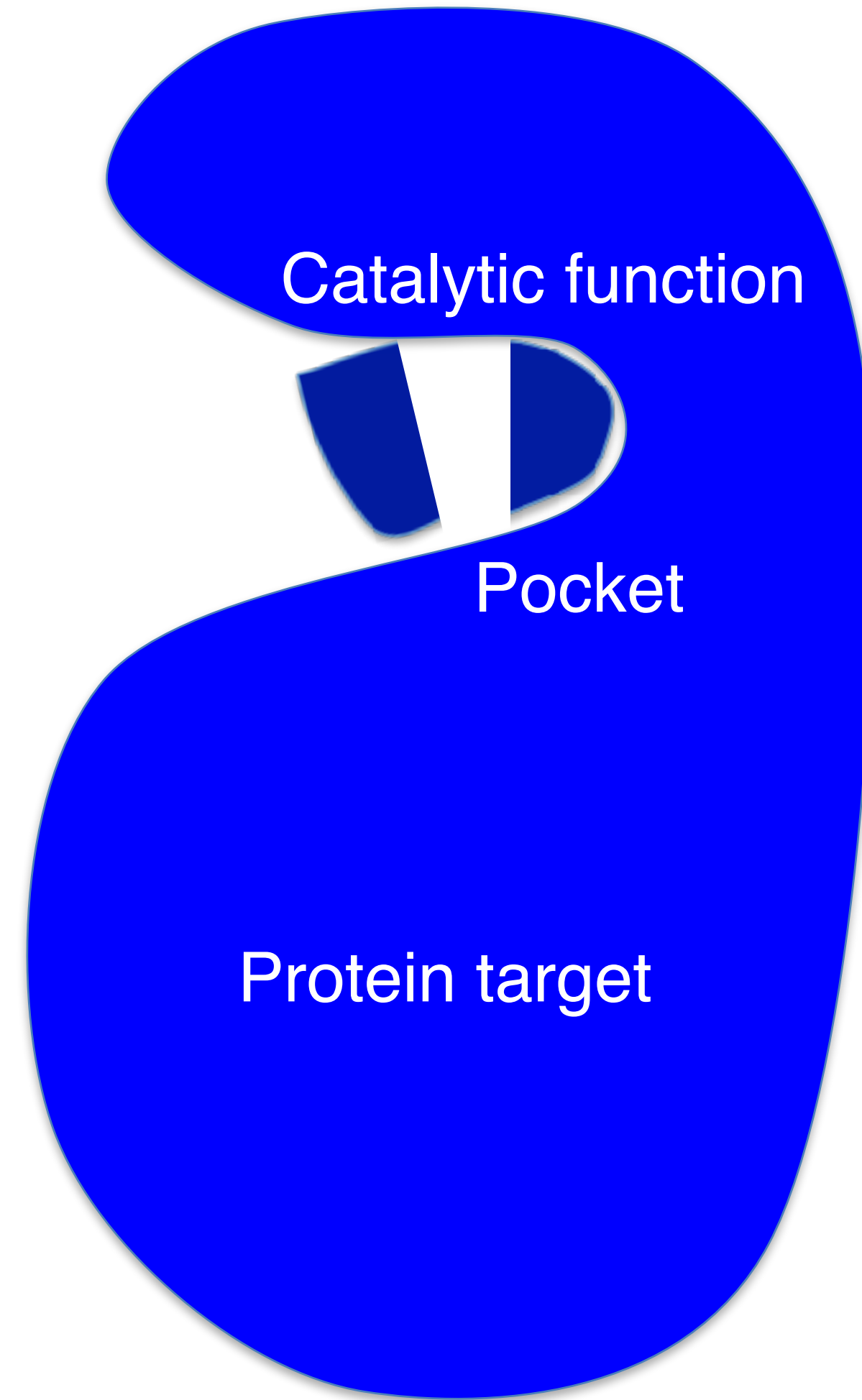
Protein substrate



Catalytic function

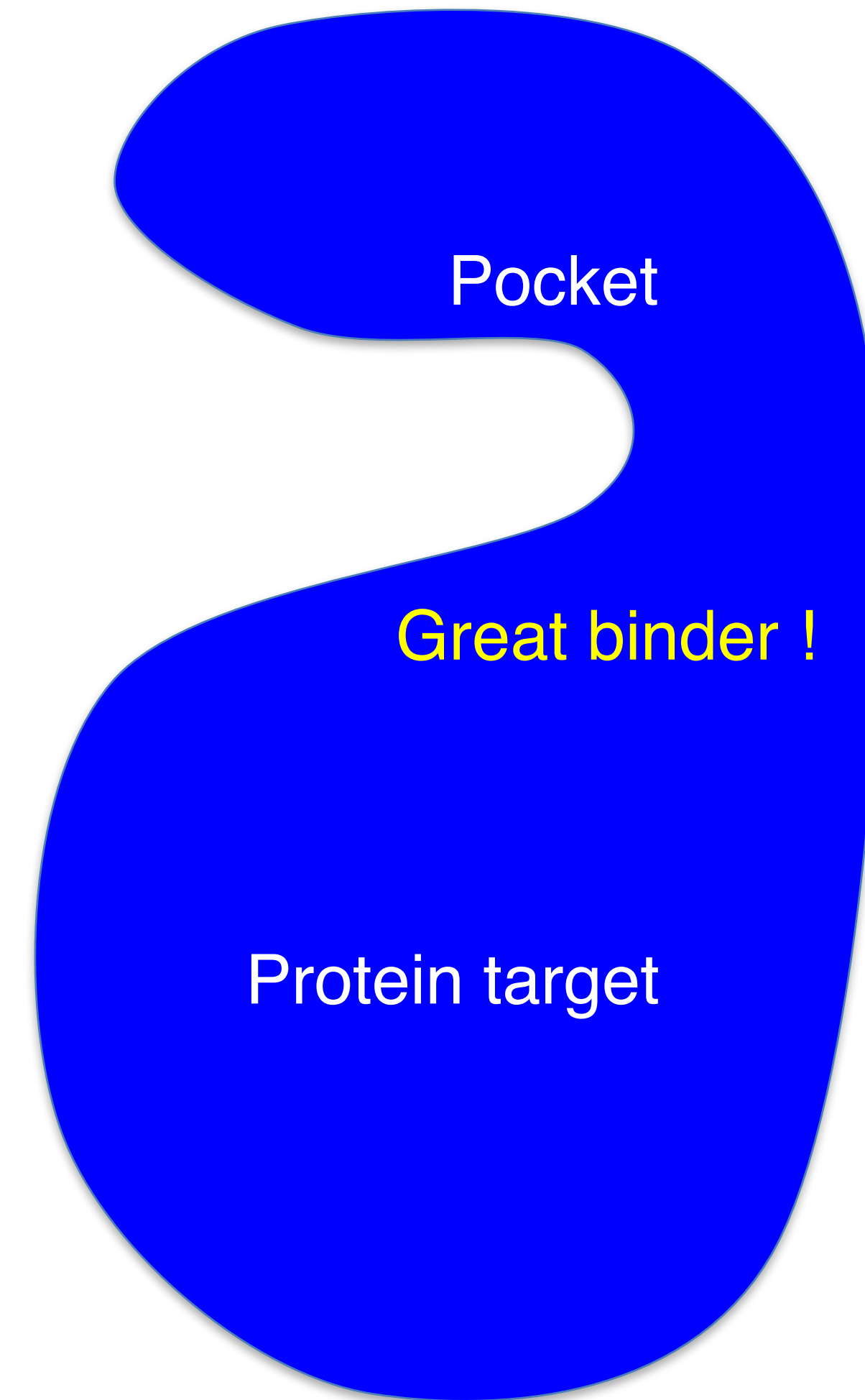
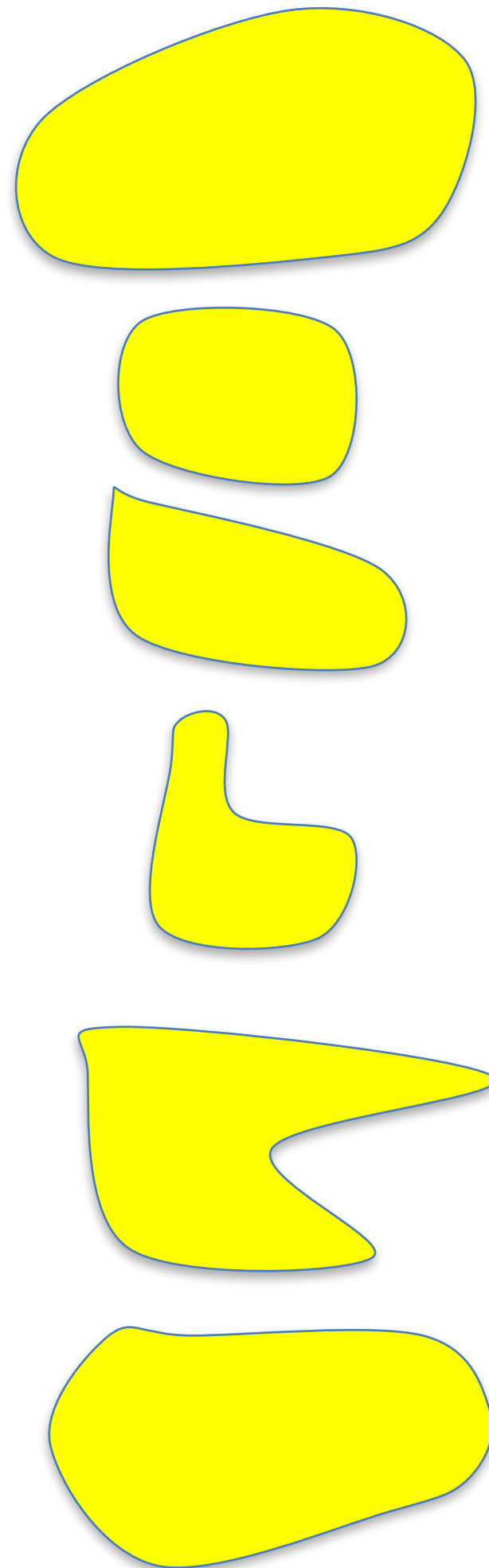
Pocket

Protein target

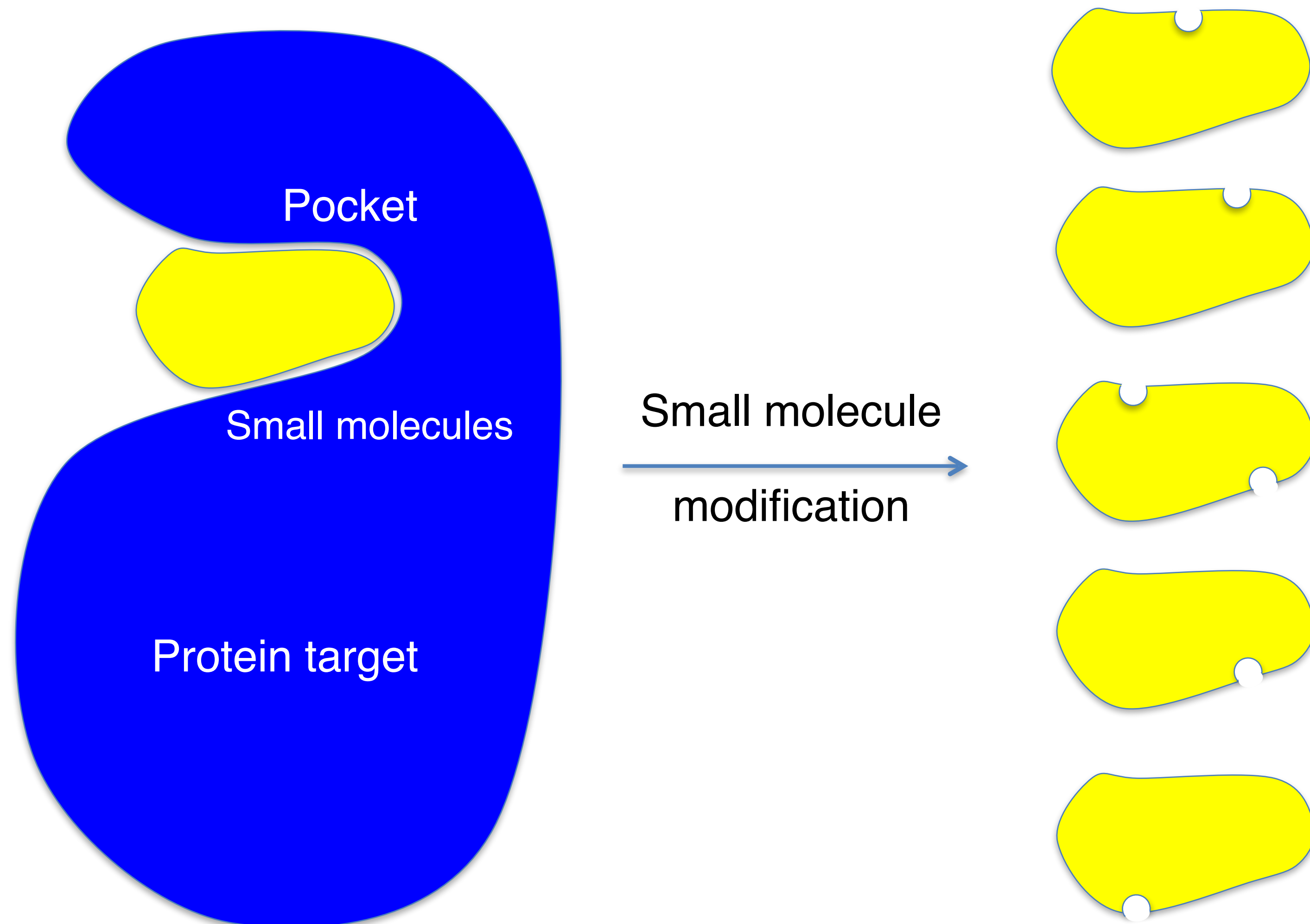


Protein binding and inhibition

Small molecules

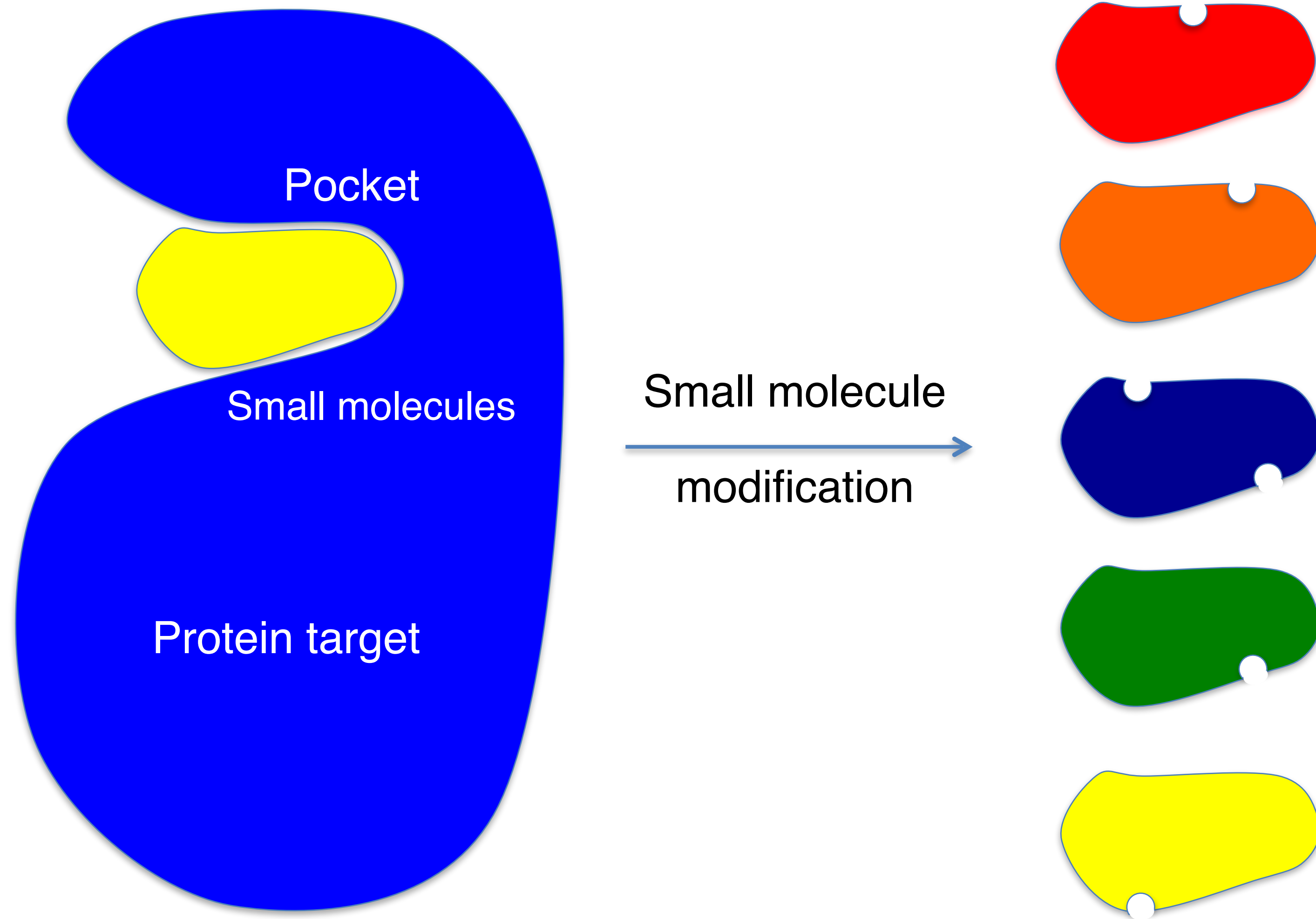


Drug design



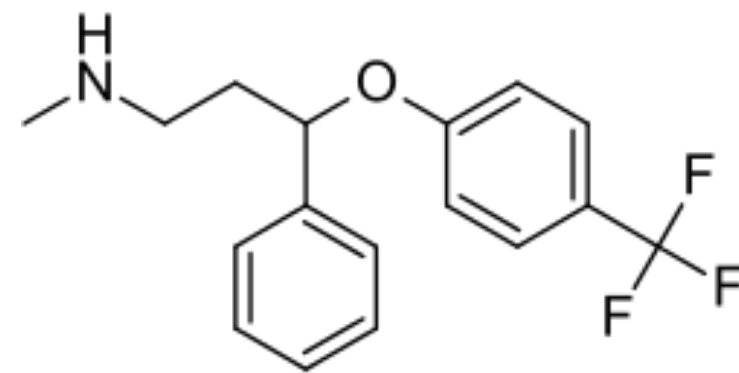
From a good inhibitor to a potential drug

Drug design

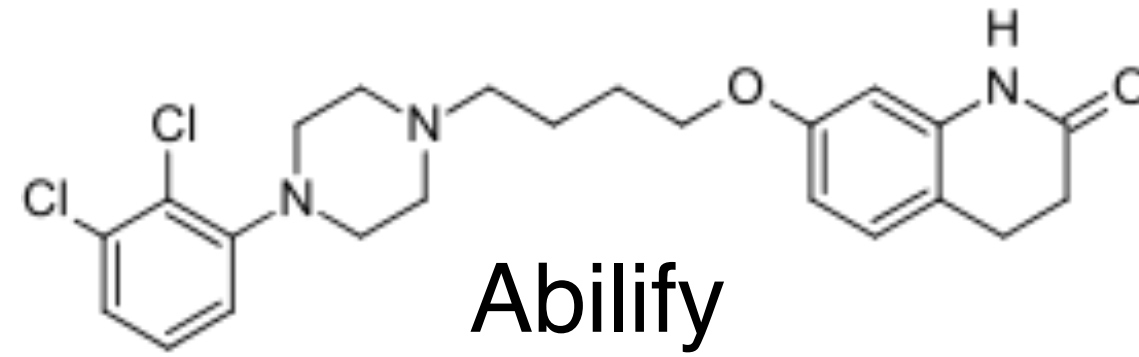


Similar binding but DIFFERENT properties!

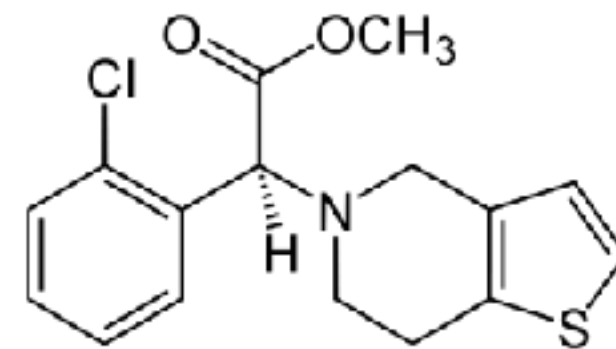
Small molecules drugs



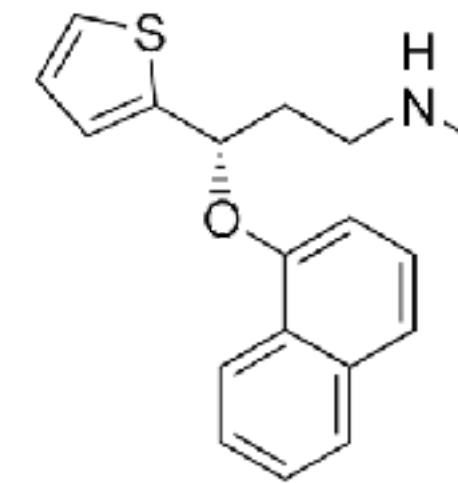
Prozac
(antidepressant)



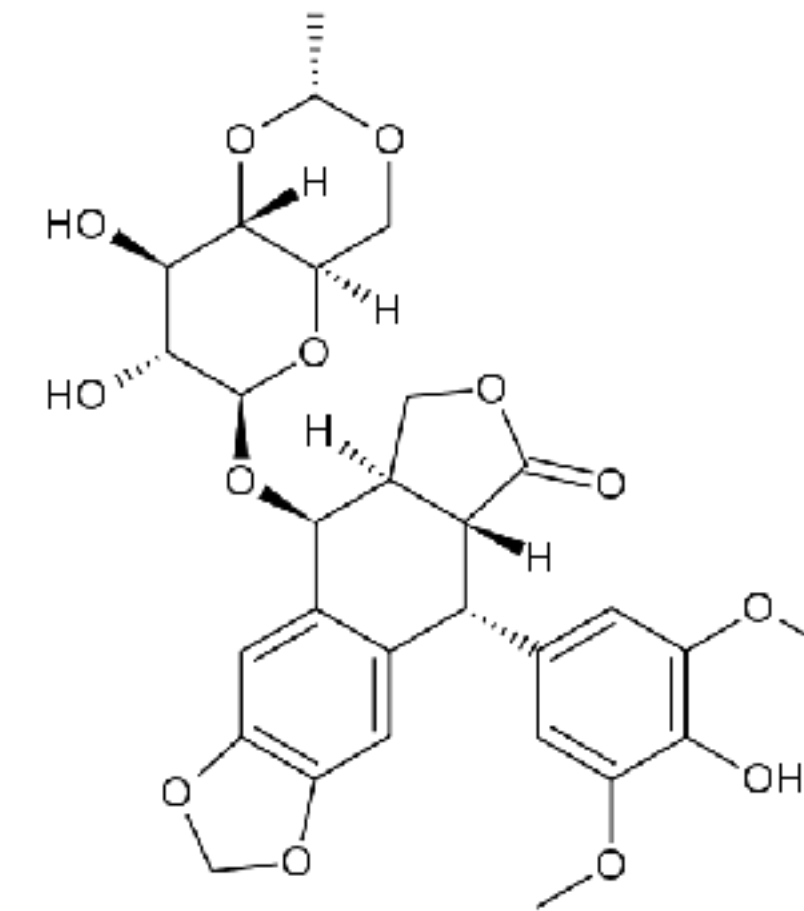
Abilify
(antipsychotic)



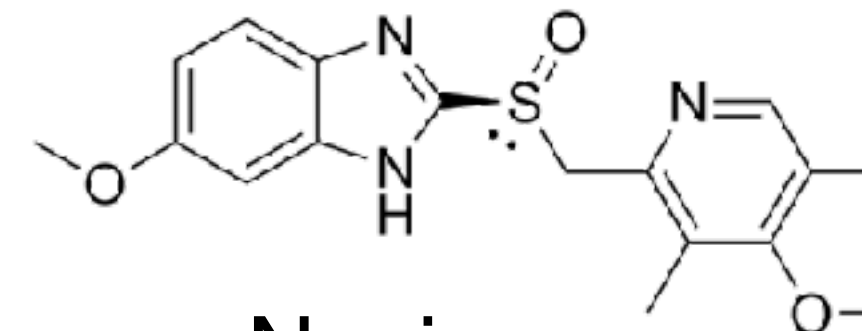
Plavix
(antiplatelet agent)



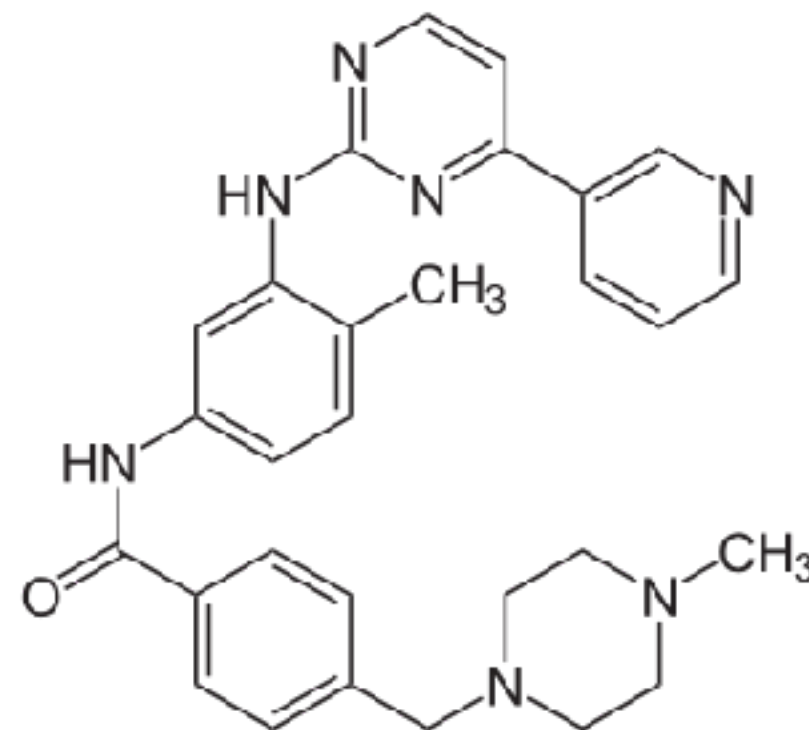
Cymbalta
(pain and anxiety)



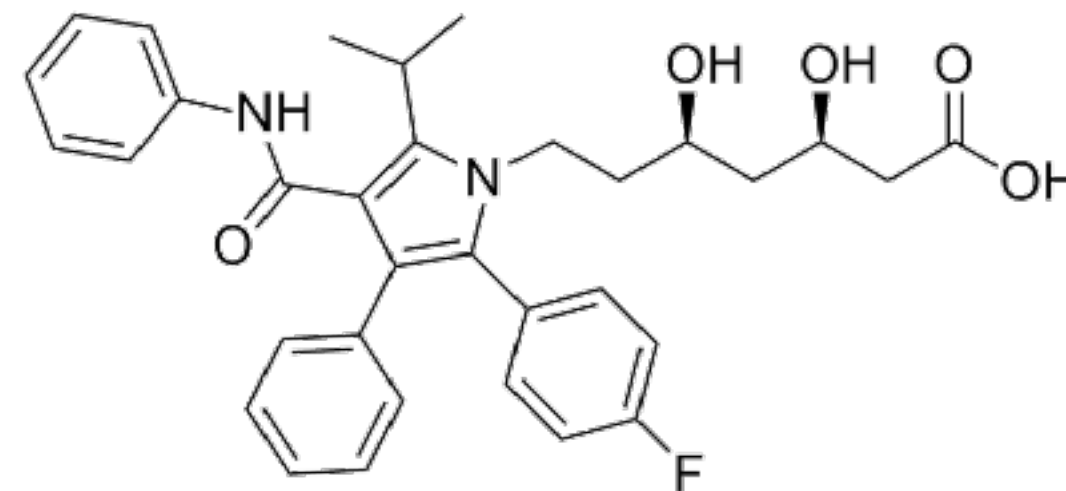
Etoposide
(cancer)



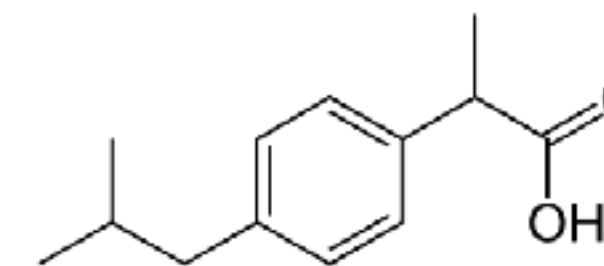
Nexium
(gastric acid)



Gleevec
(cancer)



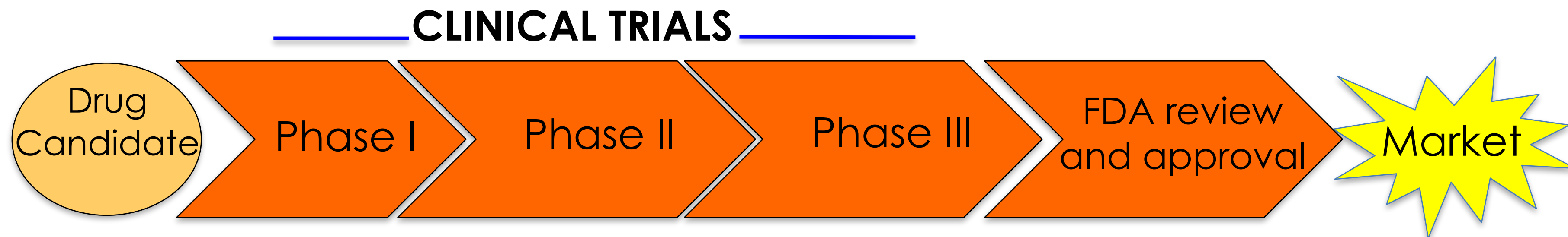
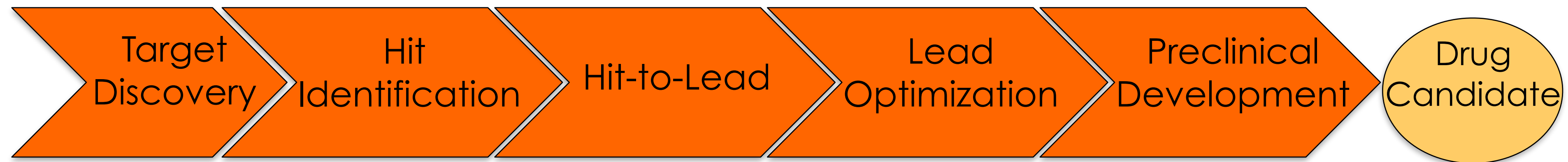
Lipitor
(Cholesterol)



Ibuprofen
(inflammation)

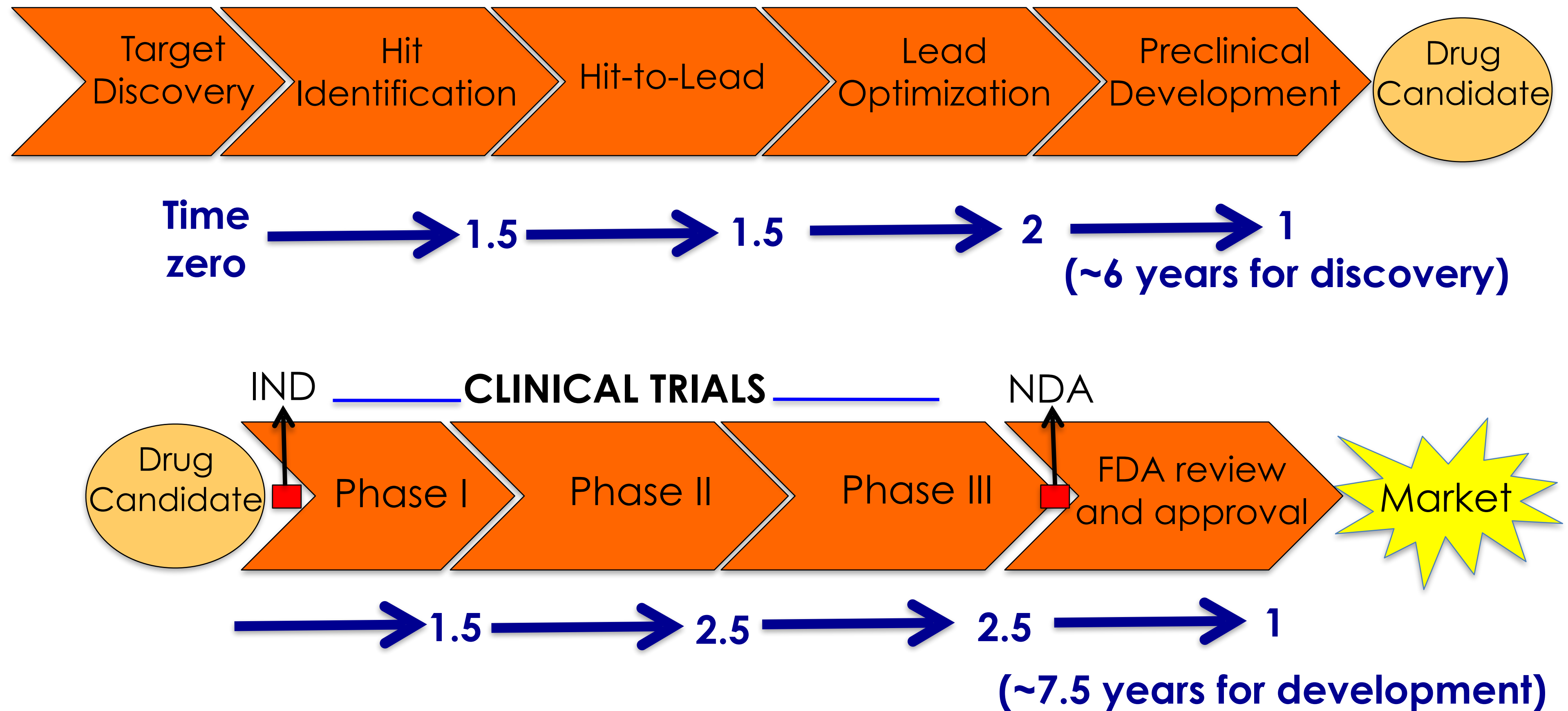
Drug Discovery & Development

Sequence of steps



Drug Discovery & Development

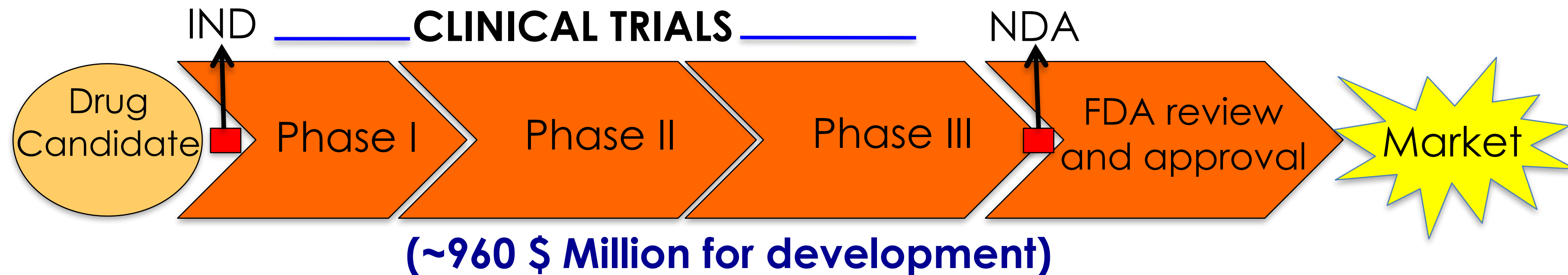
Average time requested for it



Total time (on average) = ~13.5 years

Drug Discovery & Development

Average cost requested for it



Total cost on average = ~1.78 \$ Billion for one NME

Source: How to improve R&D productivity: the pharmaceutical industry's grand challenge
Steven M. Paul, *Nat. Review Drug Discovery* March Vol. 9 **2010**

Computational Drug Design

Two approaches:

◇ Structure-Based Drug Design - SBDD

Drug design based on the interaction of the **ligand** with the 3D dimensional structure of the **receptor**

◇ Ligand-Based Drug Design - LBDD

Unknown structure of the receptor.

Drug design based on the key features of active **compounds**.

Hypothesis:

Ligands similar to an active ligand are more likely to be active than random ligands. (pharmacophore models)

Computational Drug Design

Approaches and methods:

◇ Structure-Based Drug Design - SBDD

Docking (Glide, Dock, Autodock, ICM... etc)

Kitchen D., *Nat. Review Drug Discovery* Vol. 3 Nov. **2004**

De novo design (BOMB, SMOG, BREED.. etc)

Gisbert Schneider and Uli Fechner, , *Nat. Review Drug Discovery* Vol. 4 Aug. **2005**

◇ Ligand-Based Drug Design - LBDD

Quantitative Structure-Activity Relationship (QSAR)

Markus A. Lill, *Drug Discovery Today*, Vol. 12 Dec. **2007**

Ligand similarity approaches (2D or 3D)

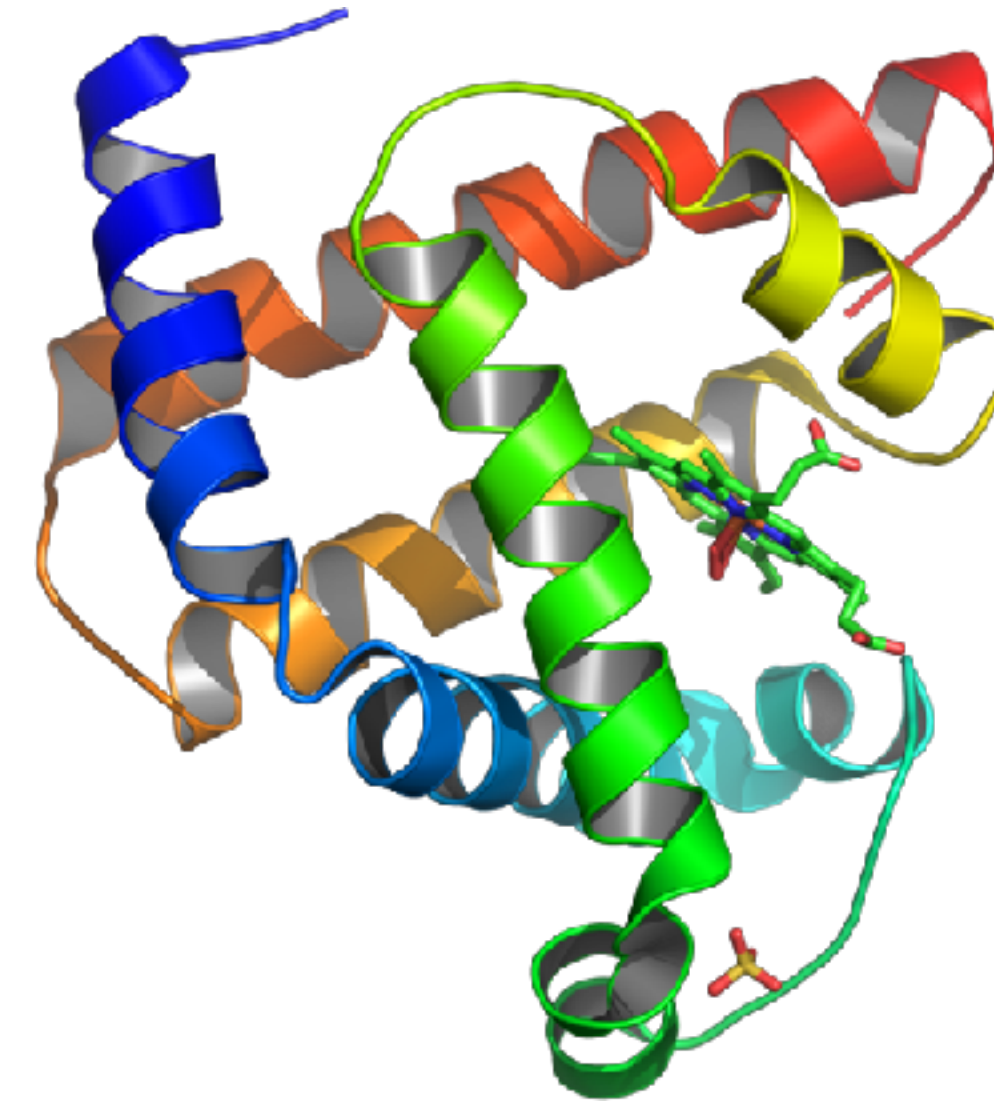
Johann Gasteiger, *J. Med. Chem.* 49, 22, **2006**.

Structure-Based Drug Design

MUST: 3-dimensional structure of the target.

Sources of structures:

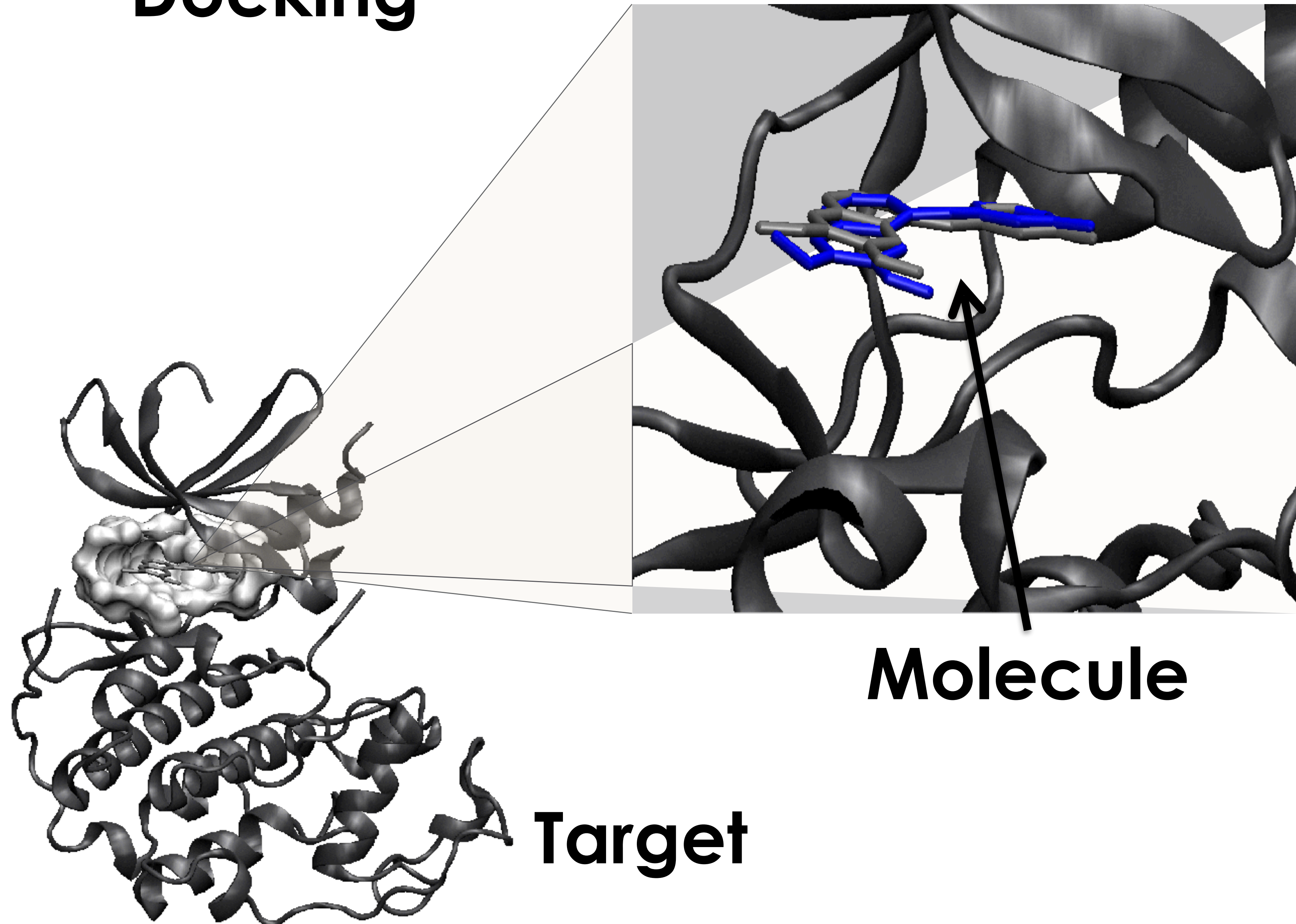
1. Crystallography
2. NMR structures
3. Cryo-EM
4. Homology or AF structures



```
AAB24882    TYHMCQFHCRCYVNNHSGEKLIECNERSKAFSCPSHLQCHKRRQIGKTHEHNQCGKAFPT 60
AAB24881    -----YECNQCGKAFAQHSSLKCHYRTHIGKPYECNQCGKAFSK 40
              *****: .***:  * *:*** * :*****.:* *****..

AAB24882    PSHLQYHERTHTGKPYECHQCGQAFKKCSLLQRHKRTHTGKPYE-CNQCGKAFAQ- 116
AAB24881    HSHLQCHKRTHTGKPYECNQCGKAFSQHGLLQRHKRTHTGKPYMNVINMVKPLHNS 98
              ***** *:*****:*****:*. : .*****:      *:.: :
```

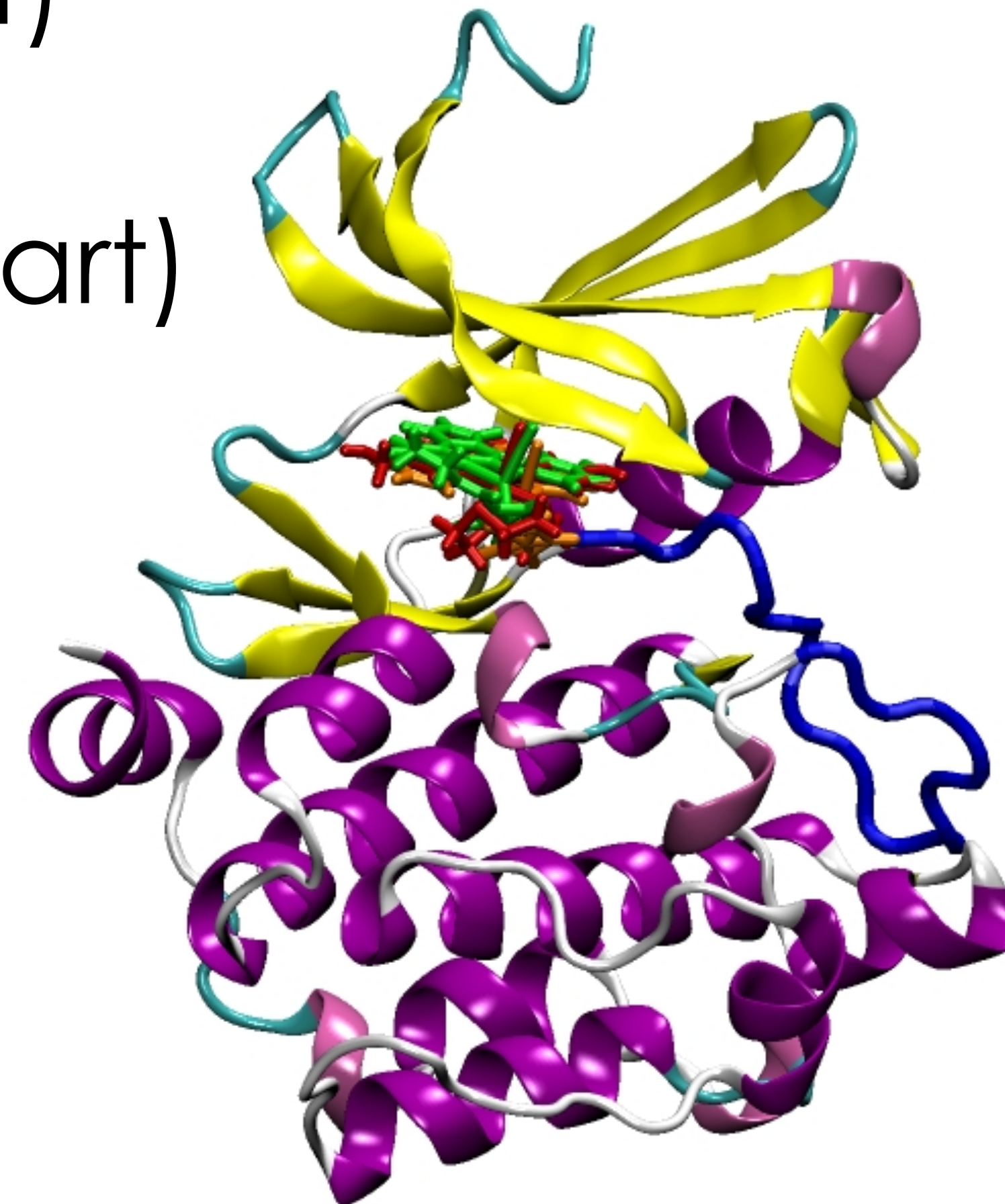
Docking



Ligand Docking

Two problem to solve:

1. Posing (the easy part)
1. Scoring (the tough part)



Ligand Docking - Scoring

The quantitative modeling of receptor – ligand interactions can be achieved by determining the equilibrium binding constant K_{eq} . The binding constant K_{eq} is directly related to the Gibbs free energy:

$$\Delta G_{bind} = -RT \ln K_{eq}$$

$$\Delta G_{bind} = \Delta H - T \Delta S = G_{complex} - (G_{receptor} + G_{ligand}).$$

Why it is so difficult to score compounds:

Experimental range of binding affinities: from 10^{-2} M (mM) to 10^{-12} M (pM)

At $T=298$ K the enthalpic contribution to the $\Delta G_{binding}$ is between -2.4 kcal/mol and -16.7 kcal/mol

In other words, a change in binding (free) energy of ~ 1.5 kcal/mol alters the binding affinity of one order of magnitude ($T=298$) !!!!

Ligand Docking - Scoring

Scoring Functions FOR DOCKING CALCULATIONS:

- Force-Field Based (Physics-based)
- Empirical
- Knowledge-based
- Descriptor-based

Classification of Current Scoring Functions

Liu and Wang, JCIIM **2015**, 55 (3), pp 475–482

Consensus Scoring

Consensus scoring are more often applied in **Virtual Screening**

Ligand Docking - Scoring

$$\Delta G_{bind} = E_{MM} - T\Delta S_{solute} + \Delta G_{solvent}$$

E_{MM} from FF:

$$E_{MM} = \sum_{bonds} K_r (r - r_{eq})^2 + \sum_{angles} K_{\vartheta} (\vartheta - \vartheta_{eq})^2 + \sum_{dihedrals} \frac{V_n}{2} [1 + \cos(n\phi - \gamma)] \\ + \sum_{i < j} \left[\frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]. \quad 3.$$

ΔS_{solute} : The solute entropy consists of four terms, namely translational, rotational, vibrational, and conformational entropy.

$\Delta g_{solvent}$: The solvent free energy consists of the two terms: 1) a nonpolar and a polar term.

Ligand Docking - Scoring

DOCK (v4.0)	$E_{vdW} + E_{electrostatic} = \sum_{prot} \sum_{lig} \left[\left(\frac{A_{ij}}{d_{ij}^a} + \frac{B_{ij}}{d_{ij}^b} \right) + 332.0 \frac{q_i q_j}{\epsilon(d_{ij}) d_{ij}} \right]$
------------------------	--

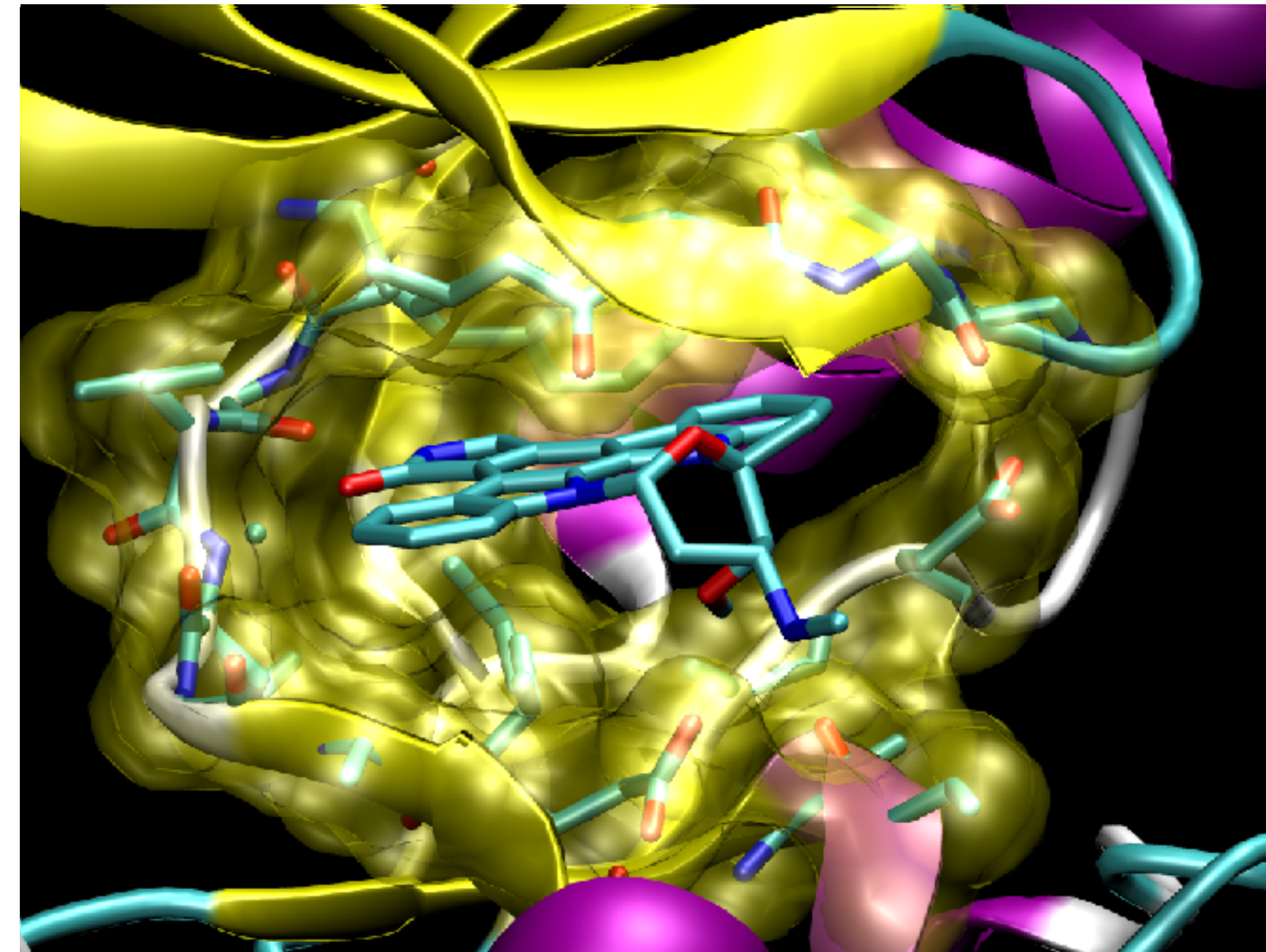
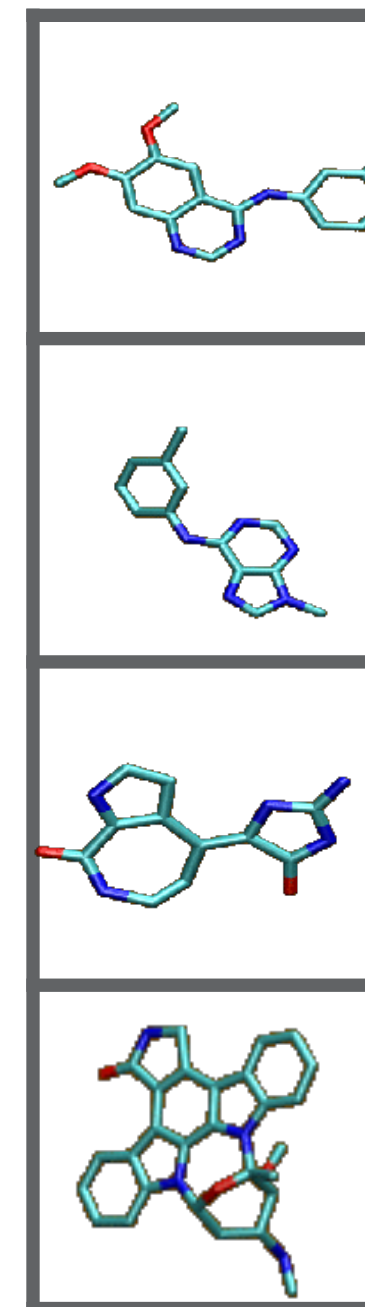
Total energy is given by the sum of energy terms.

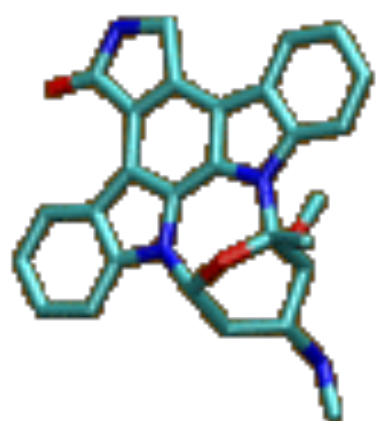
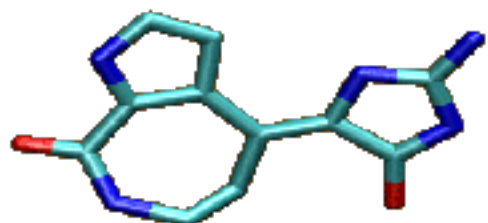
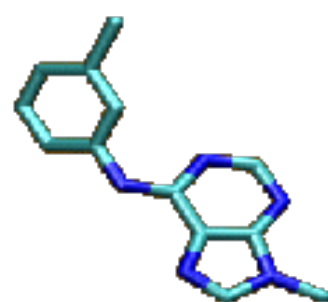
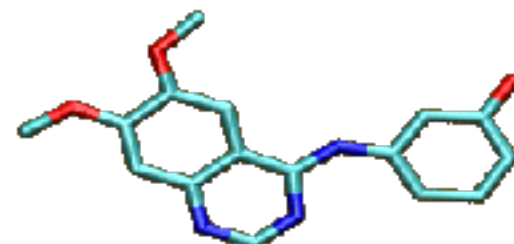
For two atoms i and j , A_{ij} and B_{ij} are van der Waals parameters for given atom types, d_{ij} is the interatomic distance, q_i and q_j are atomic partial charges, and $\epsilon(d_{ij})$ is a distance-dependent dielectric function.

Virtual screening

An exercise carried out by computational means aimed at predicting which molecules from an ensemble will likely display some activity against a target.

VLS is usually implemented as an iterative docking simulation at the target binding site.



1st -37.9 kcal/mol	2nd -29.7 kcal/mol	3rd -27.3 kcal/mol	4th -9.8 kcal/mol
			

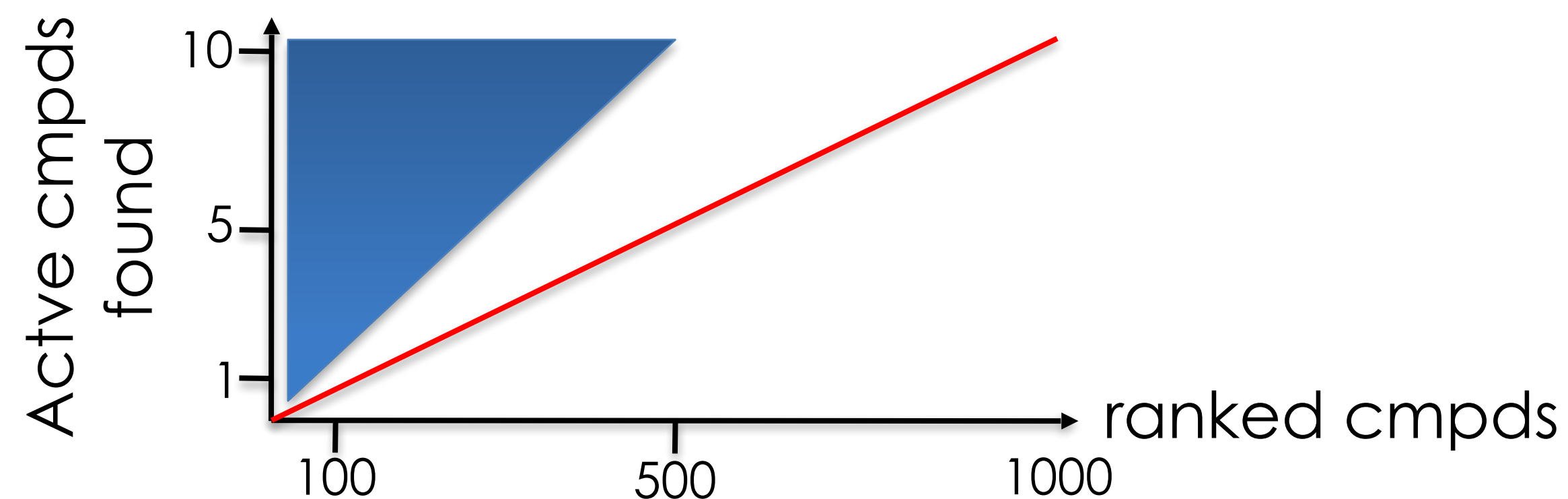
How can we evaluate the virtual screening performance?

Enrichment factor :

Example:

Suppose to have 10 active cmpds in a database of 1000 cmpds (1% of active compounds)

Random pick: 1 out of 100 should be active. (1% chance)



Ligand Docking - Scoring

Speed

Docking
Virtual Screening
Etc ..



Accuracy

Molecular Dynamics
Quantum Mechanics
Etc..

New directions

Published as a conference paper at ICLR 2023

DIFFDOCK: DIFFUSION STEPS, TWISTS, AND TURNS FOR MOLECULAR DOCKING

Gabriele Corso*, **Hannes Stärk***, **Bowen Jing***, **Regina Barzilay** & **Tommi Jaakkola**
CSAIL, Massachusetts Institute of Technology

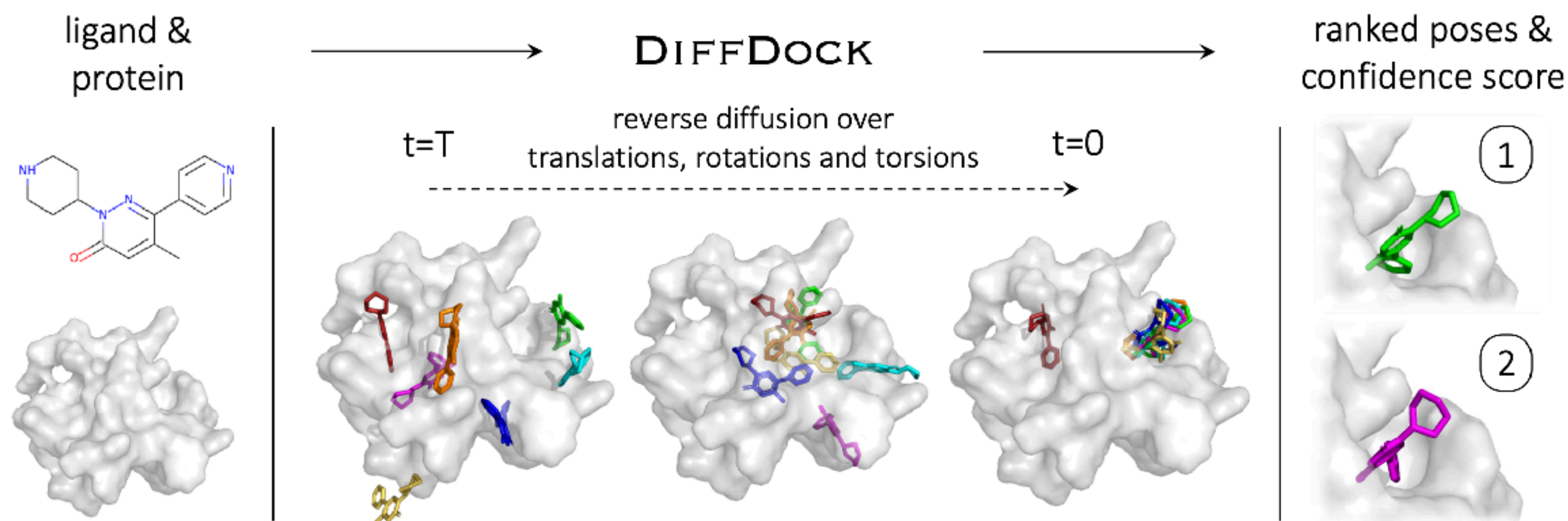
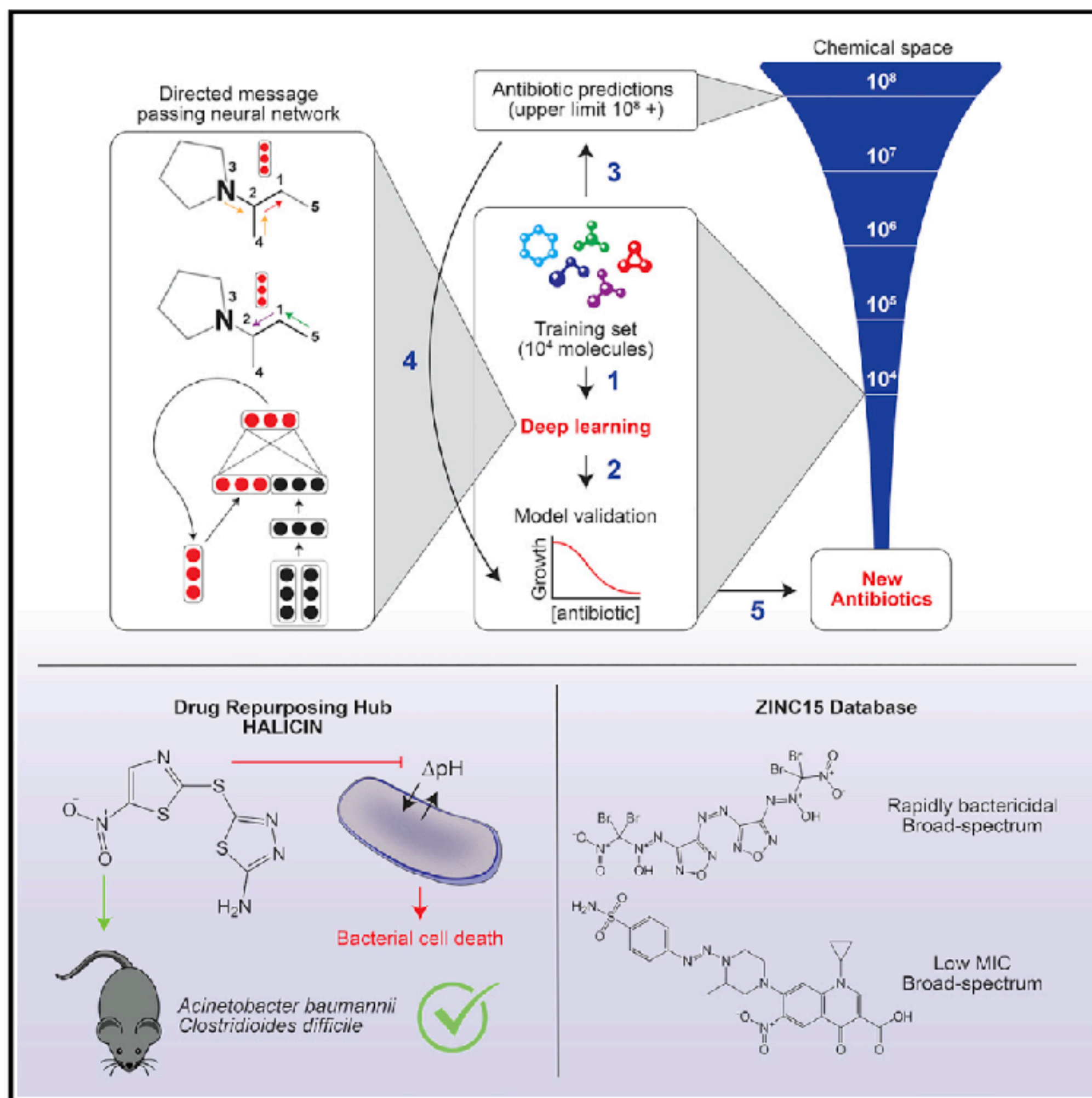


Figure 1: Overview of DIFFDOCK. *Left:* The model takes as input the separate ligand and protein structures. *Center:* Randomly sampled initial poses are denoised via a reverse diffusion over translational, rotational, and torsional degrees of freedom. *Right:* The sampled poses are ranked by the confidence model to produce a final prediction and confidence score.

A Deep Learning Approach to Antibiotic Discovery

Graphical Abstract



Authors

Jonathan M. Stokes, Kevin Yang,
Kyle Swanson, ..., Tommi S. Jaakkola,
Regina Barzilay, James J. Collins

Correspondence

regina@csail.mit.edu (R.B.),
jimjc@mit.edu (J.J.C.)

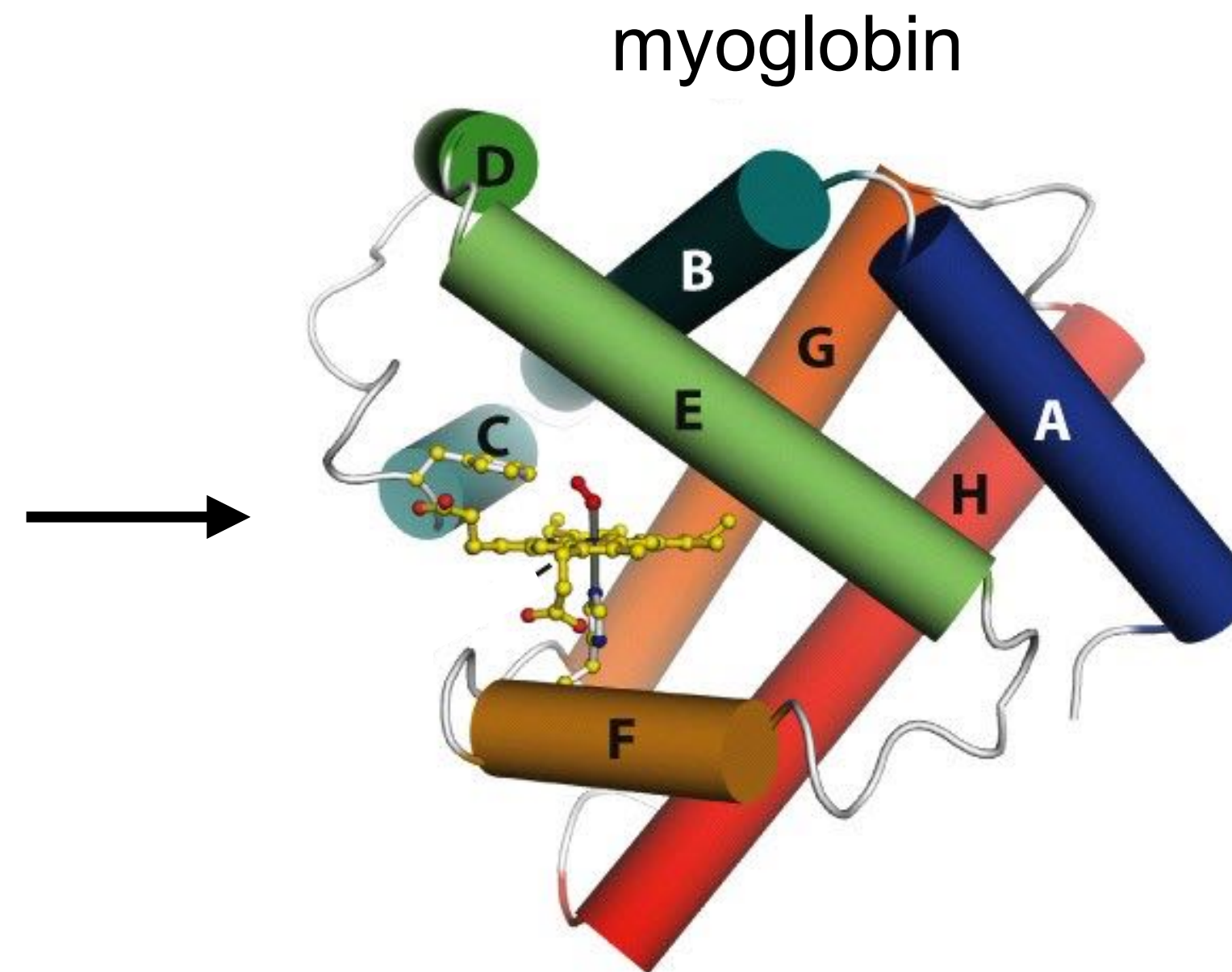
In Brief

A trained deep neural network predicts antibiotic activity in molecules that are structurally different from known antibiotics, among which Halicin exhibits efficacy against broad-spectrum bacterial infections in mice.

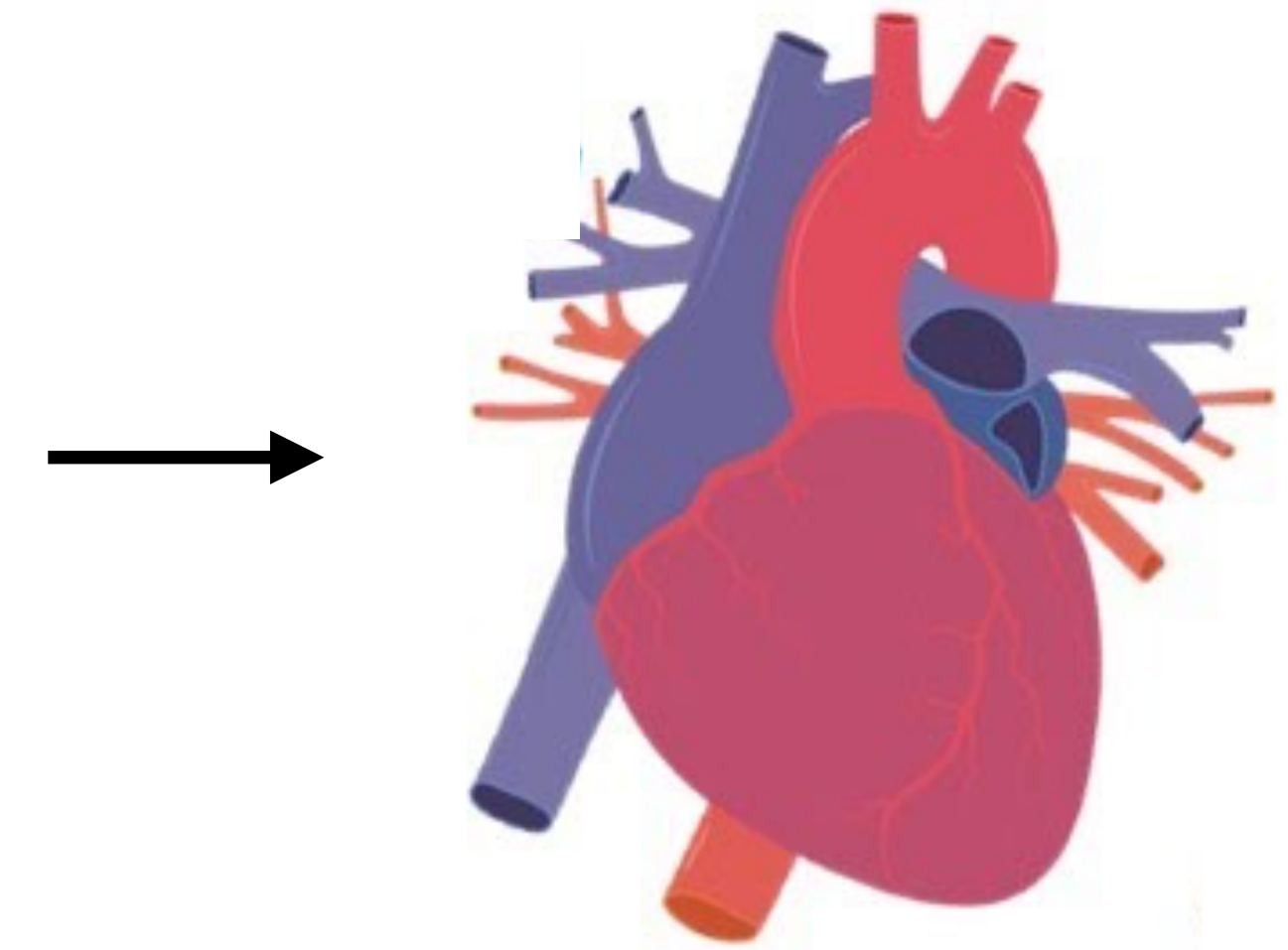
The folding paradigm

M G L S D G E W Q L V L N V W G
K V E A D I P G H G Q E V L I R
L F K G H P E T L E K F D K F K
H L K S E D E M K A S E D L K K
H G A T V L T A L G G I L K K K
G H H E A E I K P L A Q S H A T
K H K I P V K Y L E F I S E C I
I Q V L Q S K H P G D F G A D A
Q G A M N K A L E L F R K D M A
S N Y K E L G F Q G

sequence



structure

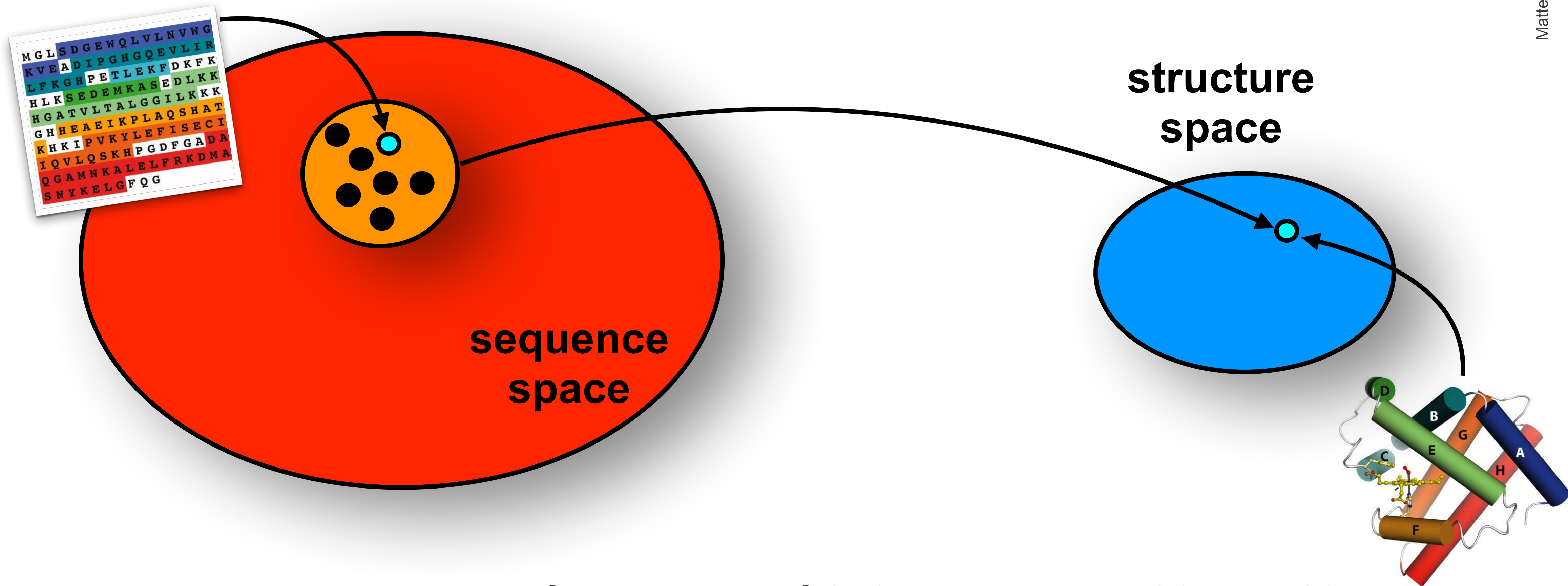


function

evolution (billion year)

- Prediction of final structure and binding helps discovering new biology
- Not all the questions are answered though by AF2 !!

The sequence space is enormous

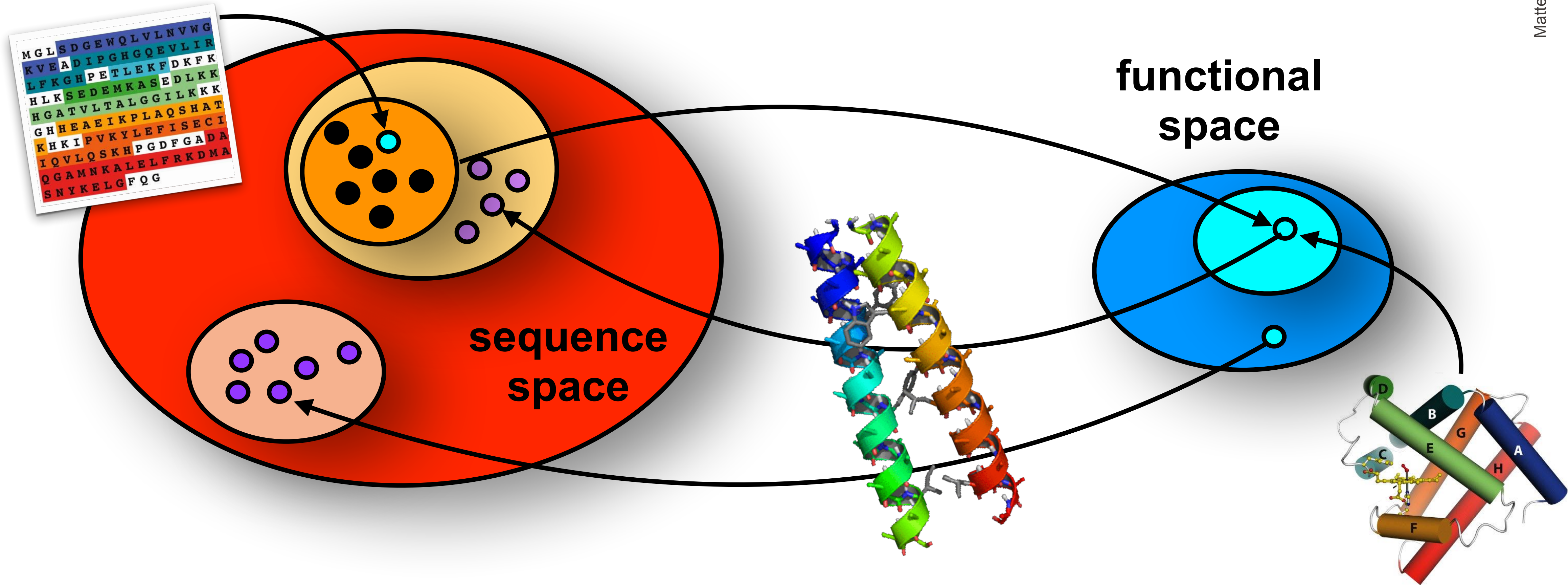


- potential sequence space for proteins of 150 amino acids $20^{150} \sim 10^{195}$
- atoms in the observed universe $\sim 10^{80}$
- the sequences explored by evolution are much less ($\sim 10^{10-20}$), structures lesser

EPFL The inverse folding problem — design

67

Matteo Dal Peraro

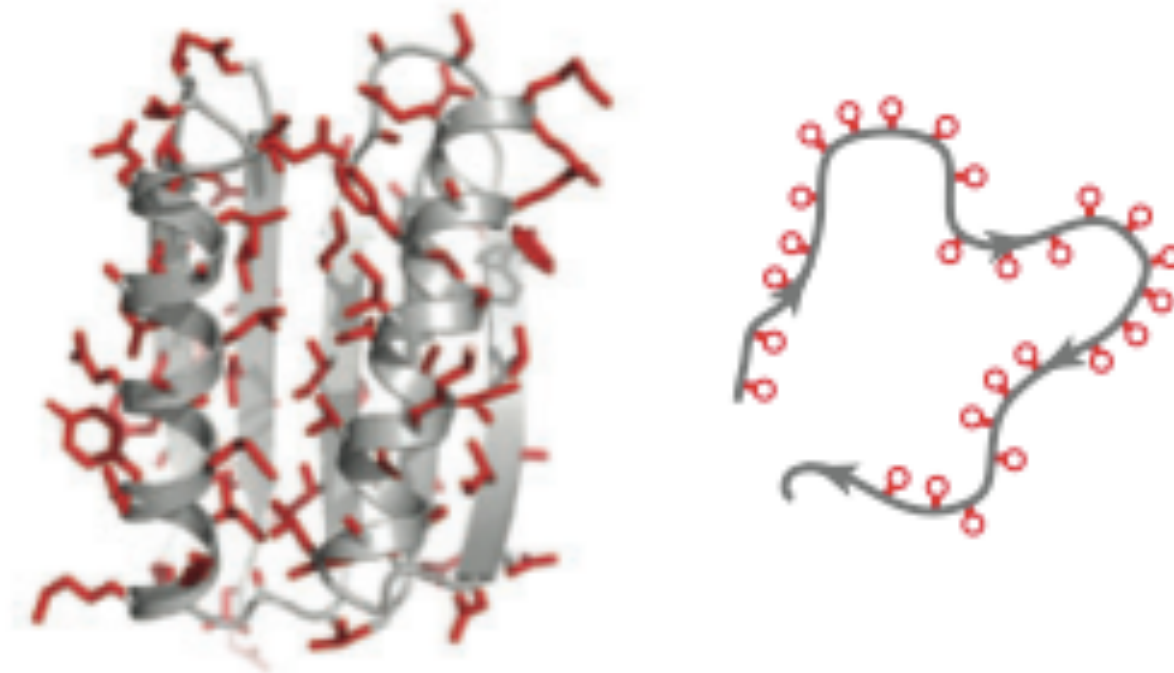


- Application to study protein evolution and function
- Protein engineering for therapeutics, synthetic biology and (bio)technology

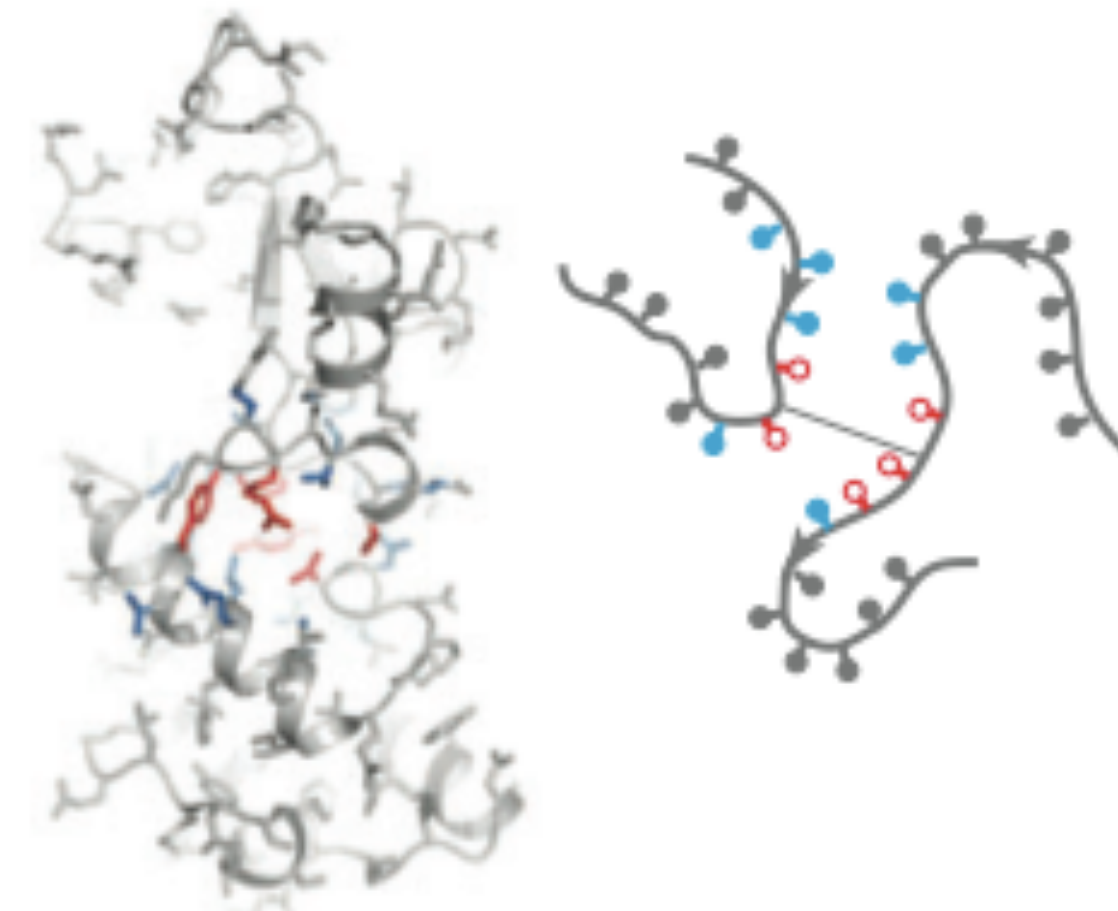
Multiple tasks for protein design

- create de novo proteins
- explore new folds
- embed new functions

Protein design

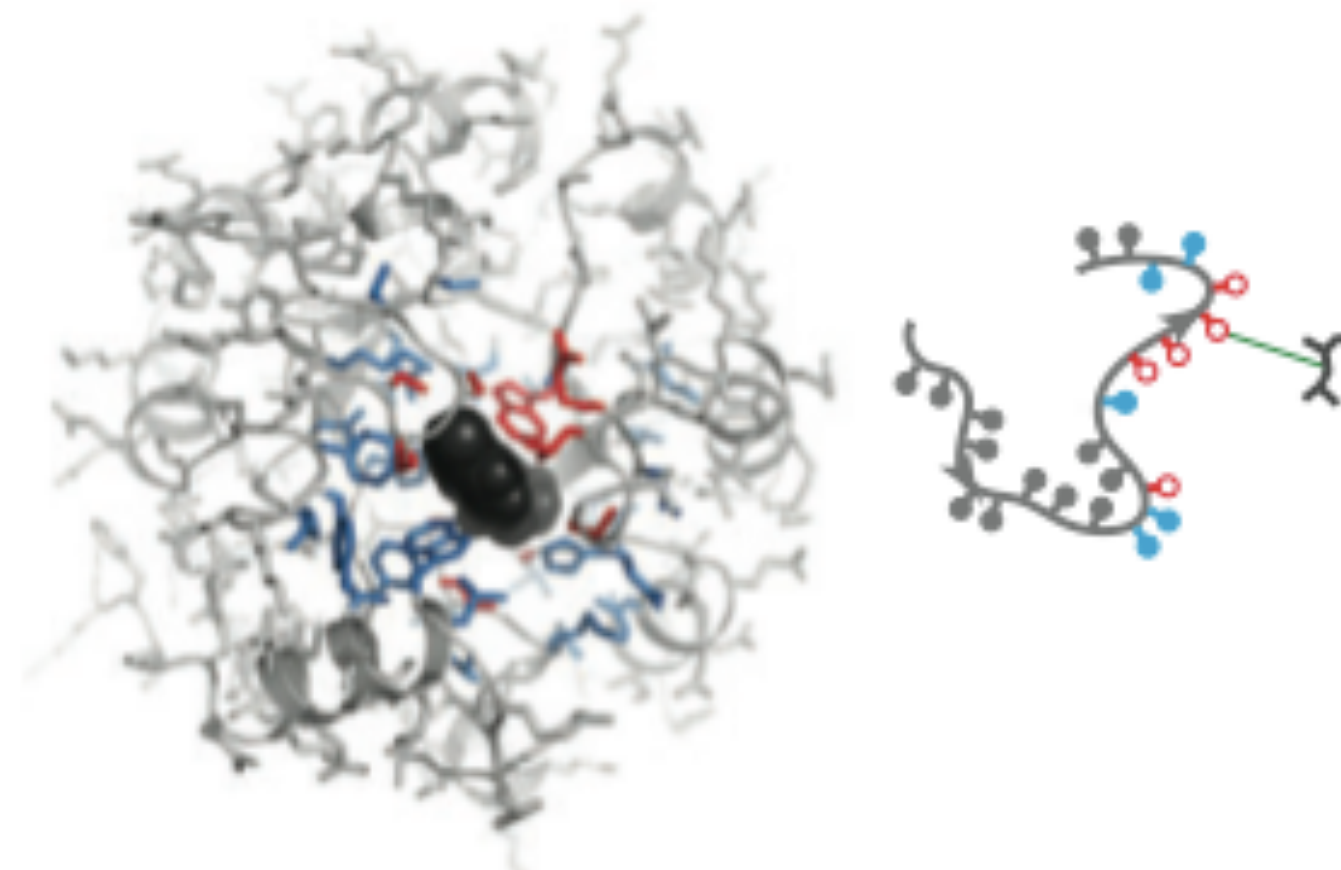


Protein-protein interface design



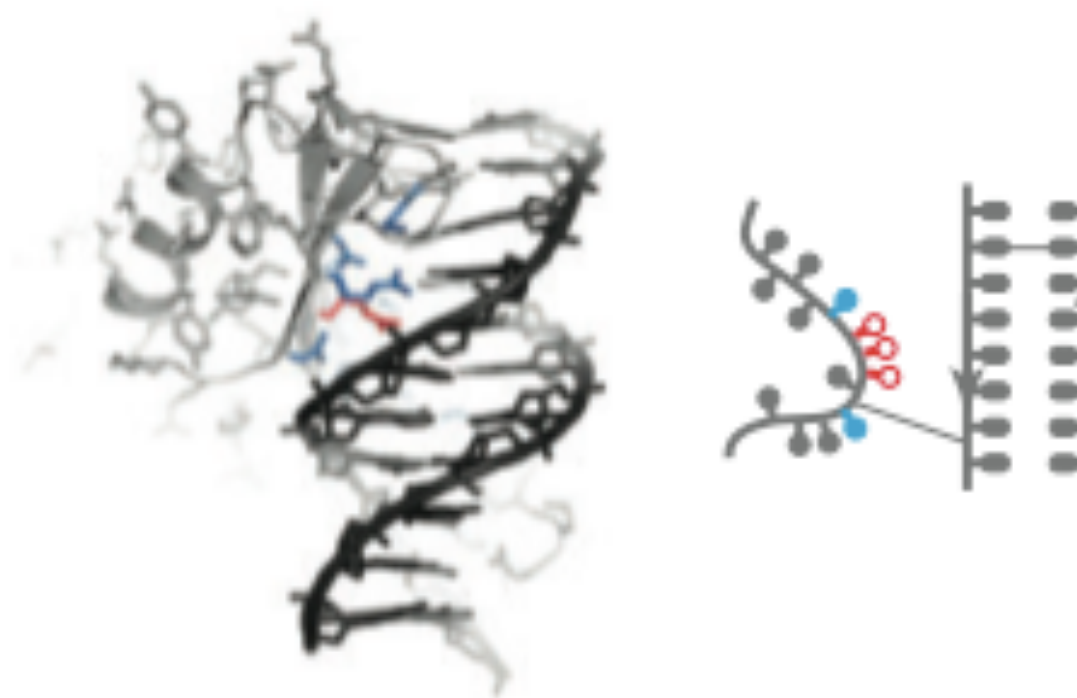
- create high affinity binders
- therapeutic biologics
- artificial sensors/probes

Enzyme design



- tailor enzymatic function
- improve thermostability
- catalyse new reactions

Protein-DNA interface design



- explore DNA interactions
- new therapeutic solutions

• Filled colored circles - flexible side chains
○ empty colored circles – flexible amino acid: design

EPFL The origins: the Paracelsus challenge ('94)

- Rose and Creamer: convert a protein to another fold changing no more than 50% of its sequence

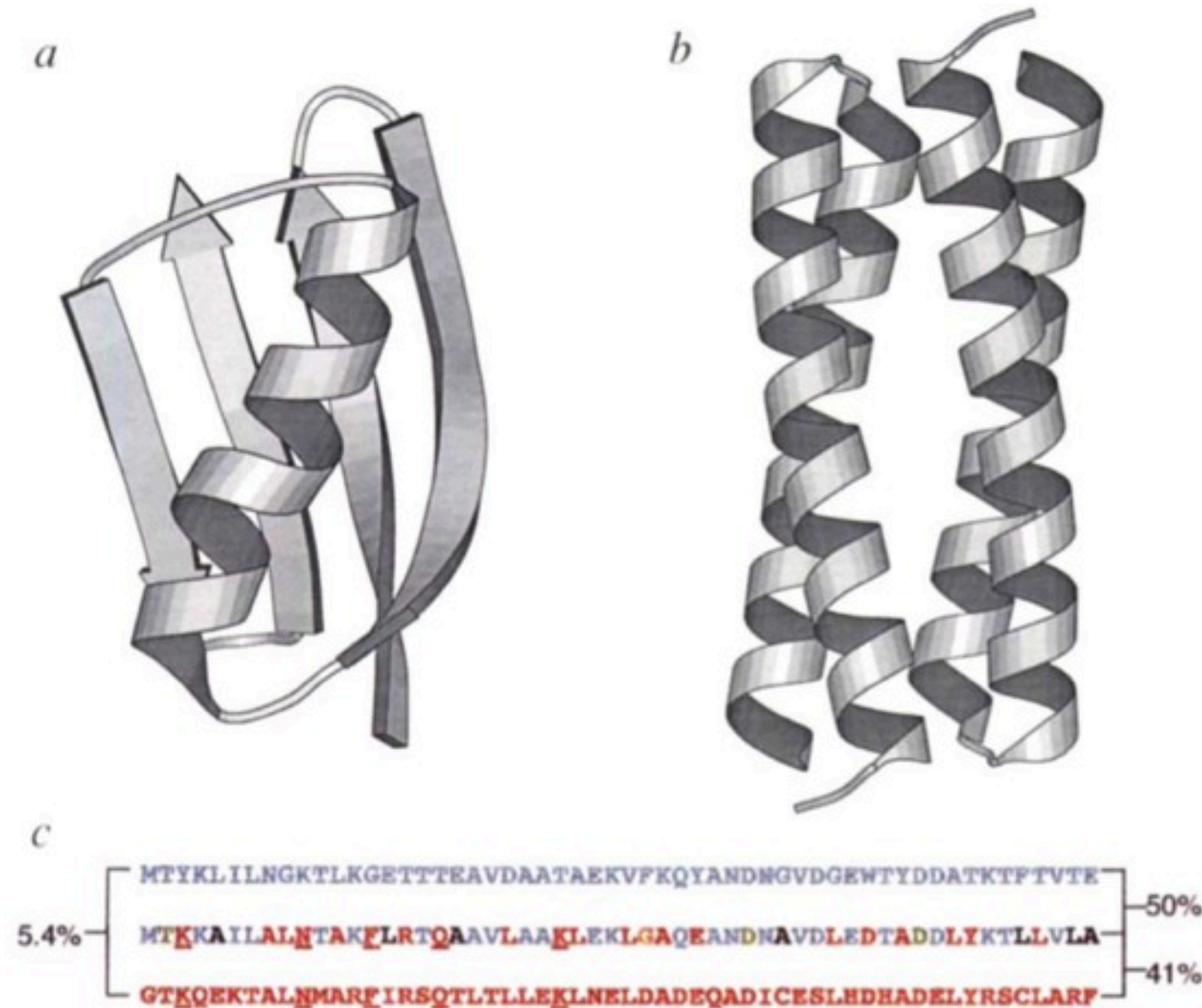


Fig. 1 Ribbon representation²⁹ of the folds of **a**, the B1 domain of IgG-binding protein G⁵ and **b**, Rop⁶. **c**, An alignment of the sequences of the B1 domain (blue), Rop (red) and Janus. Residues in Janus are coded as follows: blue, residues from B1; red, residues from Rop; underlined red, RNA-binding residues in Rop¹³; green, residues that are conserved in both Rop and B1; black, 'a' and 'd' position residues that are different from those in wild-type Rop; orange, the first residue of the turn between Helix 1 and Helix 2. The D30G mutation was introduced in the turn of Janus because a previous study demonstrated that this point mutation increases the stability of Rop³⁰. The percent identity between the different sequences are indicated. The seven amino acid, unstructured C-terminal tail of Rop (Gly-Asp-Asp-Gly-Glu-Asn-Leu) extends beyond the sequence depicted for both Rop and Janus and is also not shown in (b). It was retained in Janus because it increases the solubility of wild type Rop³¹.

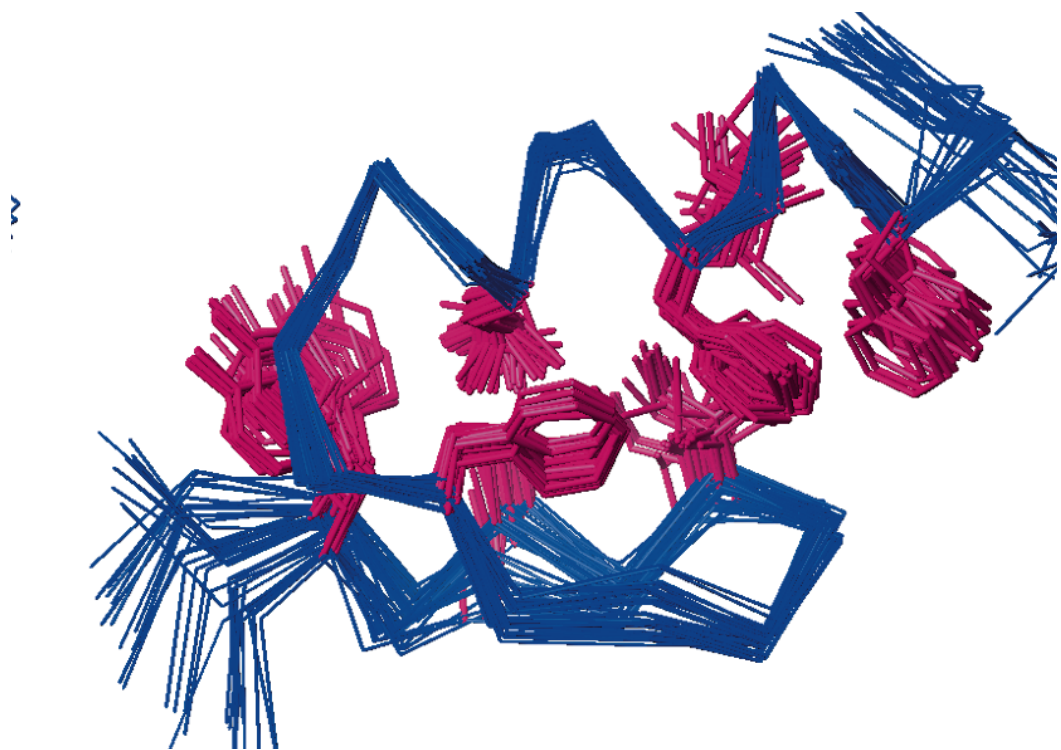
Design of a zinc-less zinc finger

De Novo Protein Design: Fully Automated Sequence Selection

Bassil I. Dahiyat† and Stephen L. Mayo*

The first fully automated design and experimental validation of a novel sequence for an entire protein is described. A computational design algorithm based on physical chemical potential functions and stereochemical constraints was used to screen a combinatorial library of 1.9×10^{27} possible amino acid sequences for compatibility with the design target, a $\beta\beta\alpha$ protein motif based on the polypeptide backbone structure of a zinc finger domain. A BLAST search shows that the designed sequence, full sequence design 1 (FSD-1), has very low identity to any known protein sequence. The solution structure of FSD-1 was solved by nuclear magnetic resonance spectroscopy and indicates that FSD-1 forms a compact well-ordered structure, which is in excellent agreement with the design target structure. This result demonstrates that computational methods can perform the immense combinatorial search required for protein design, and it suggests that an unbiased and quantitative algorithm can be used in various structural contexts.

Dahiyat, BI, and SL Mayo.
Science 278, 5335 (3 October 1997): 82-7



FSD-1 NMR determination

comparison of FSD-1 designed and NMR

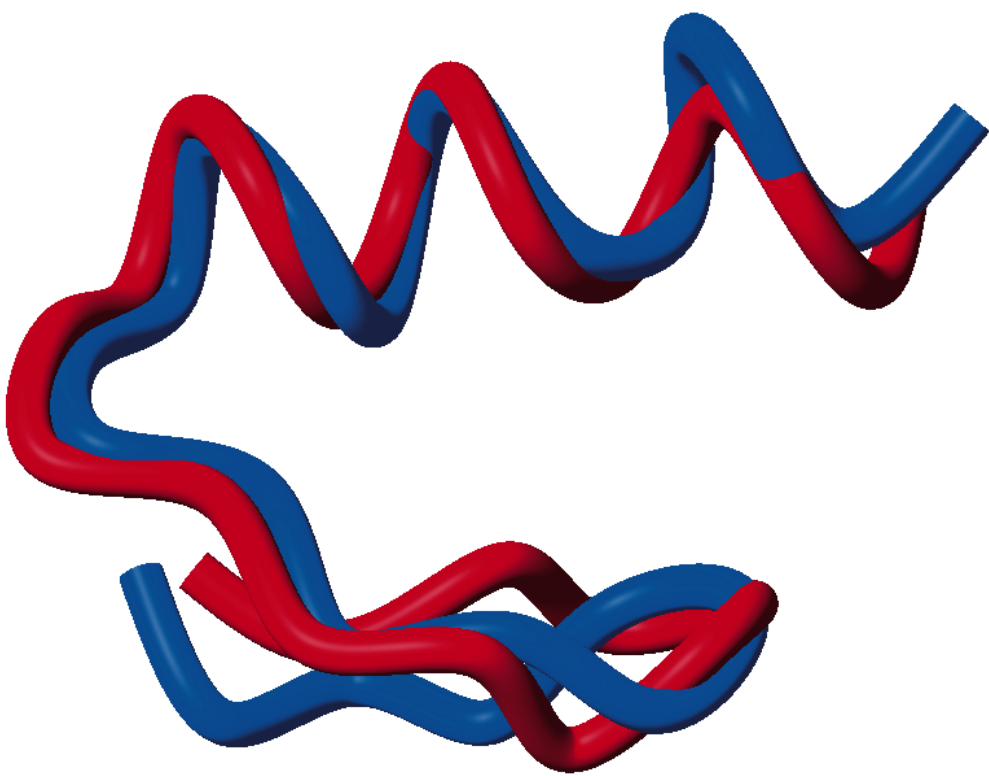
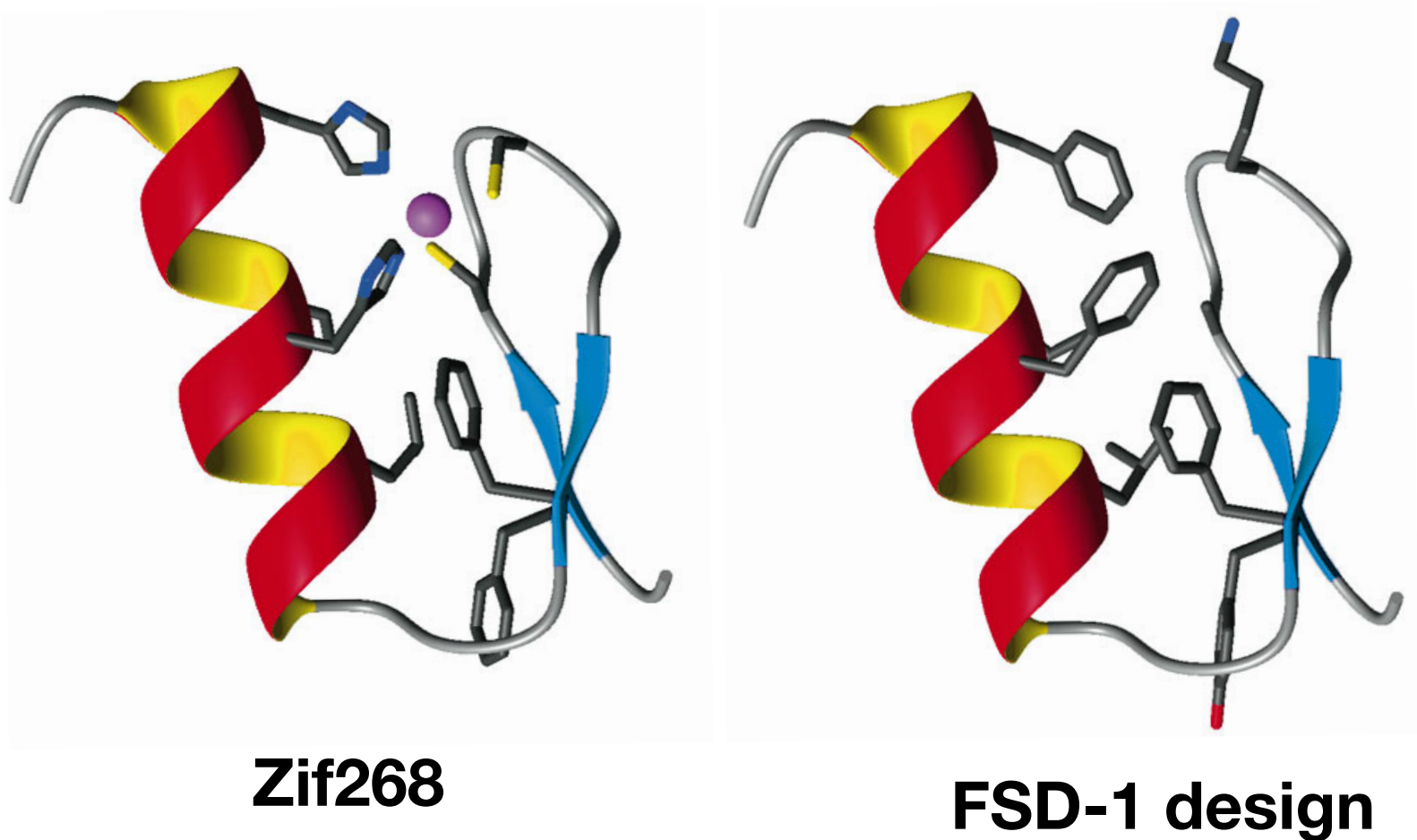
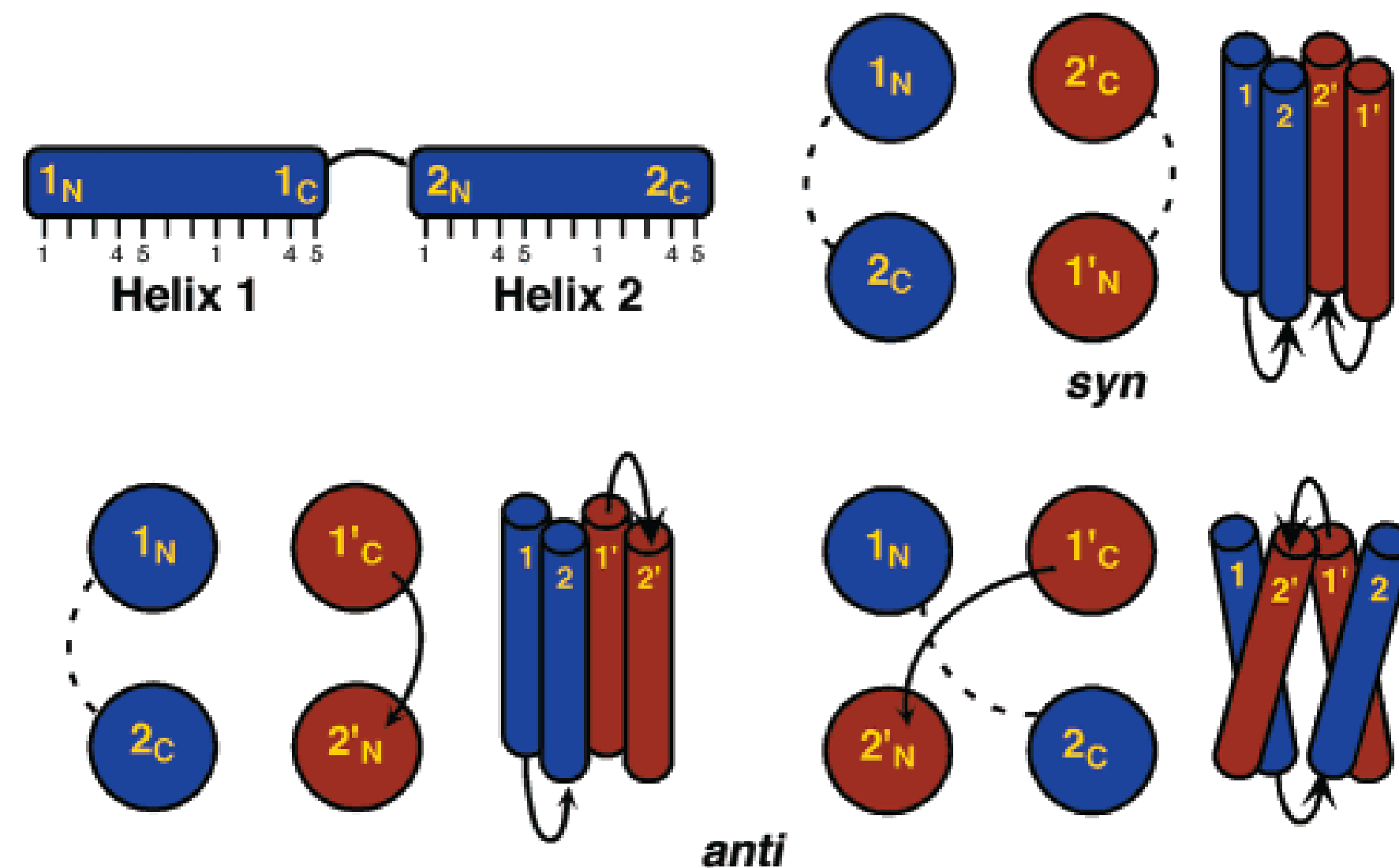


Table 2. Comparison of the FSD-1 experimentally determined structure and the design target structure. The FSD-1 structure is the restrained energy minimized average from the NMR structure determination. The design target structure is the second DNA binding module of the zinc finger Zif268 (9)

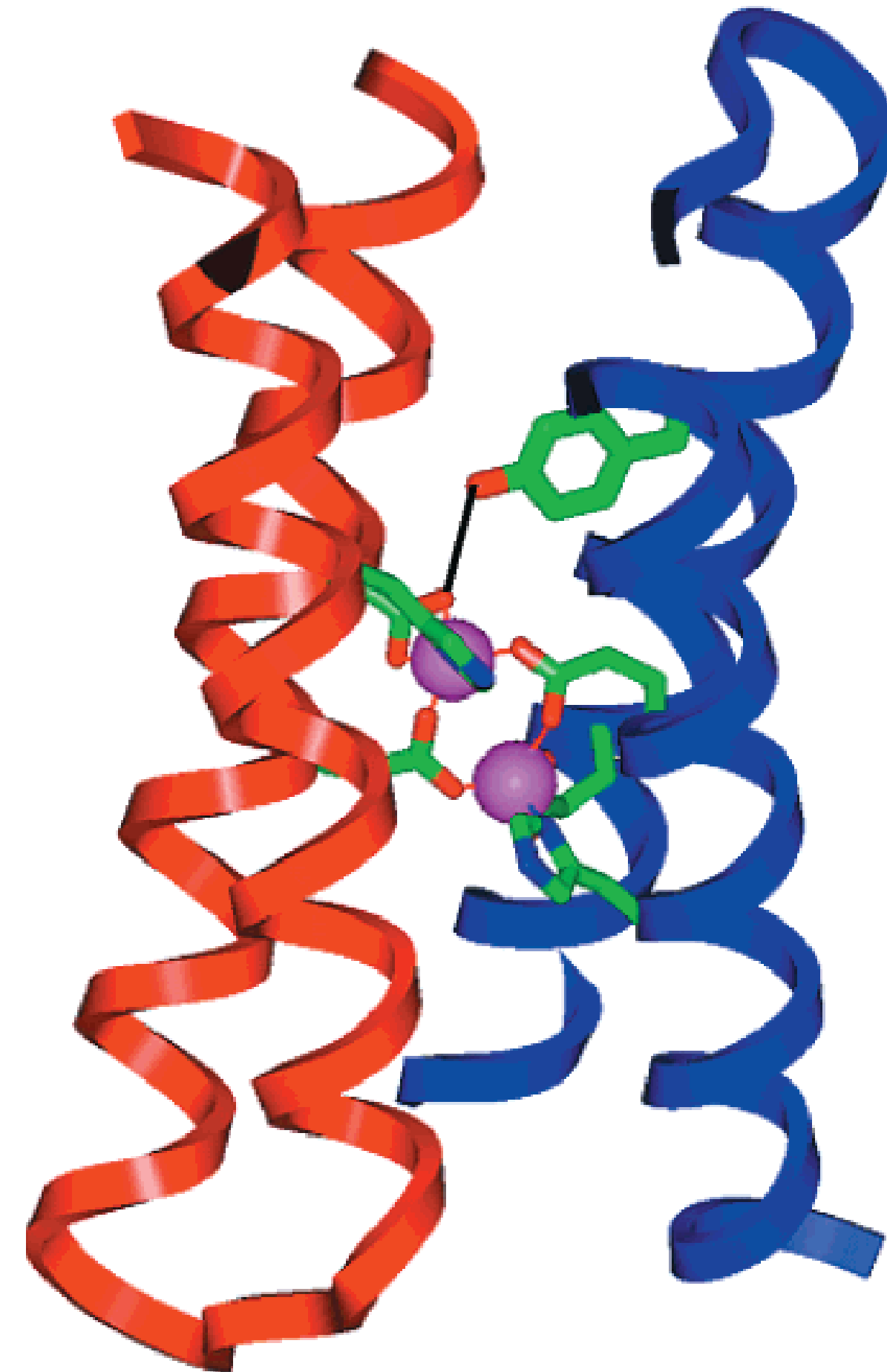
Atomic rms deviations (Å)		
Backbone, residues 3 to 26	1.98	
Backbone, residues 8 to 26	0.98	
Super-secondary structure parameters*		
	FSD-1	Design target
h (Å)	9.9	8.9
θ (degrees)	14.2	16.5
Ω (degrees)	13.1	13.5

De novo design of helical bundles


- use knowledge on natural occurring helical bundles to assemble and functionalize new folds and functions

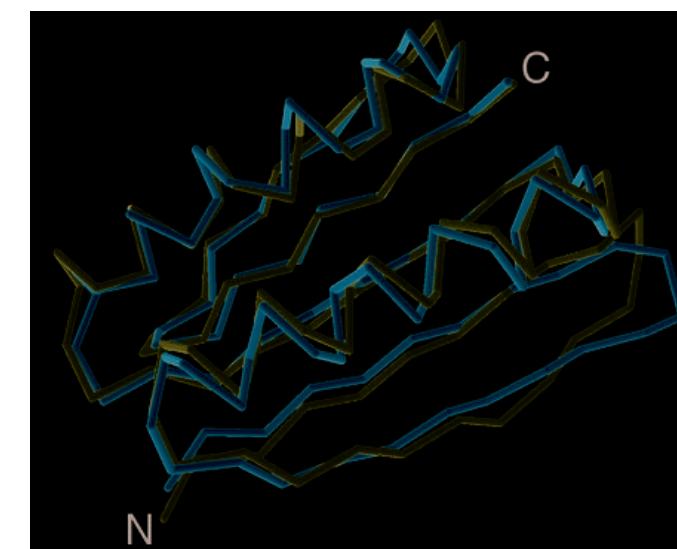
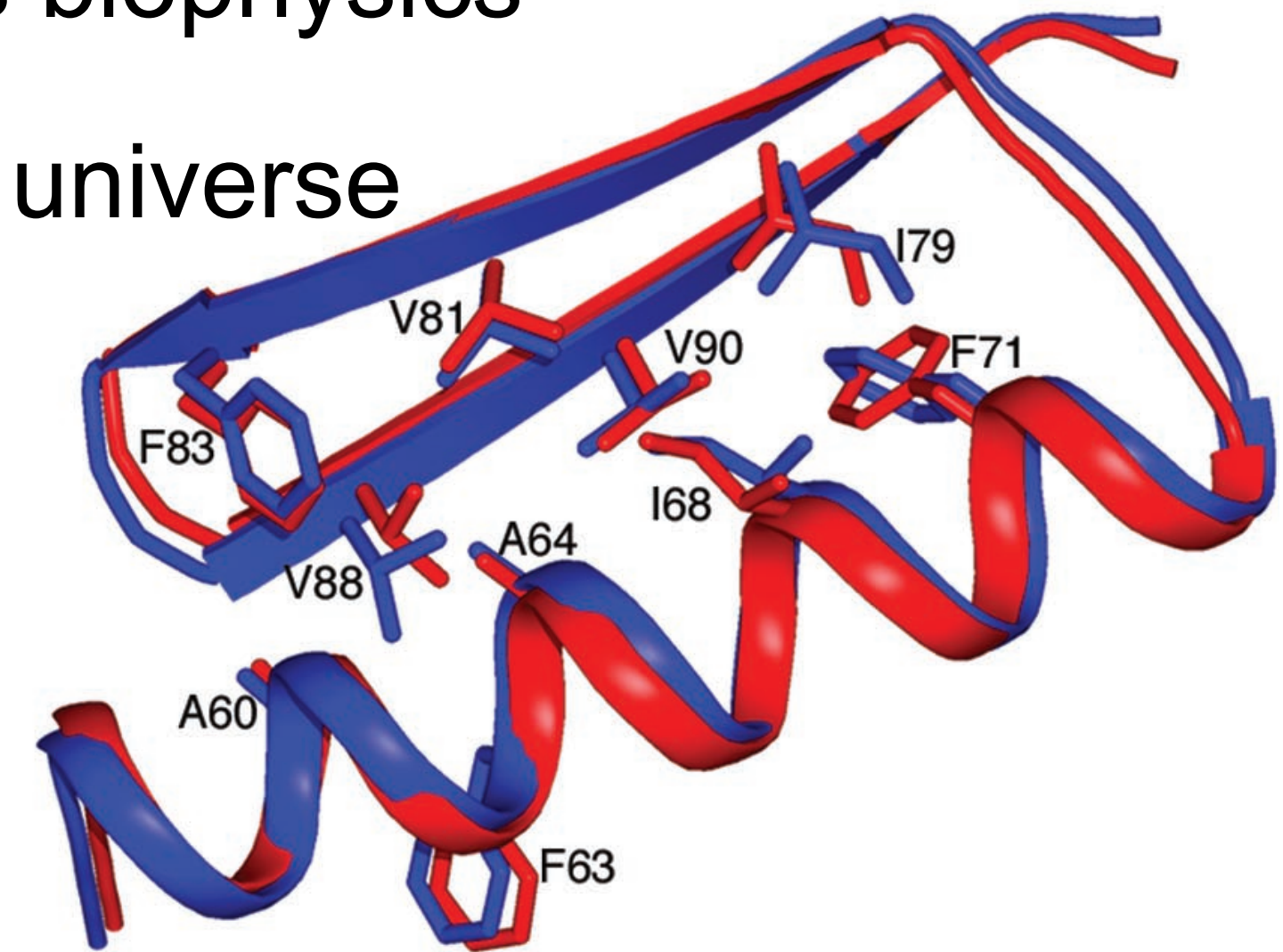
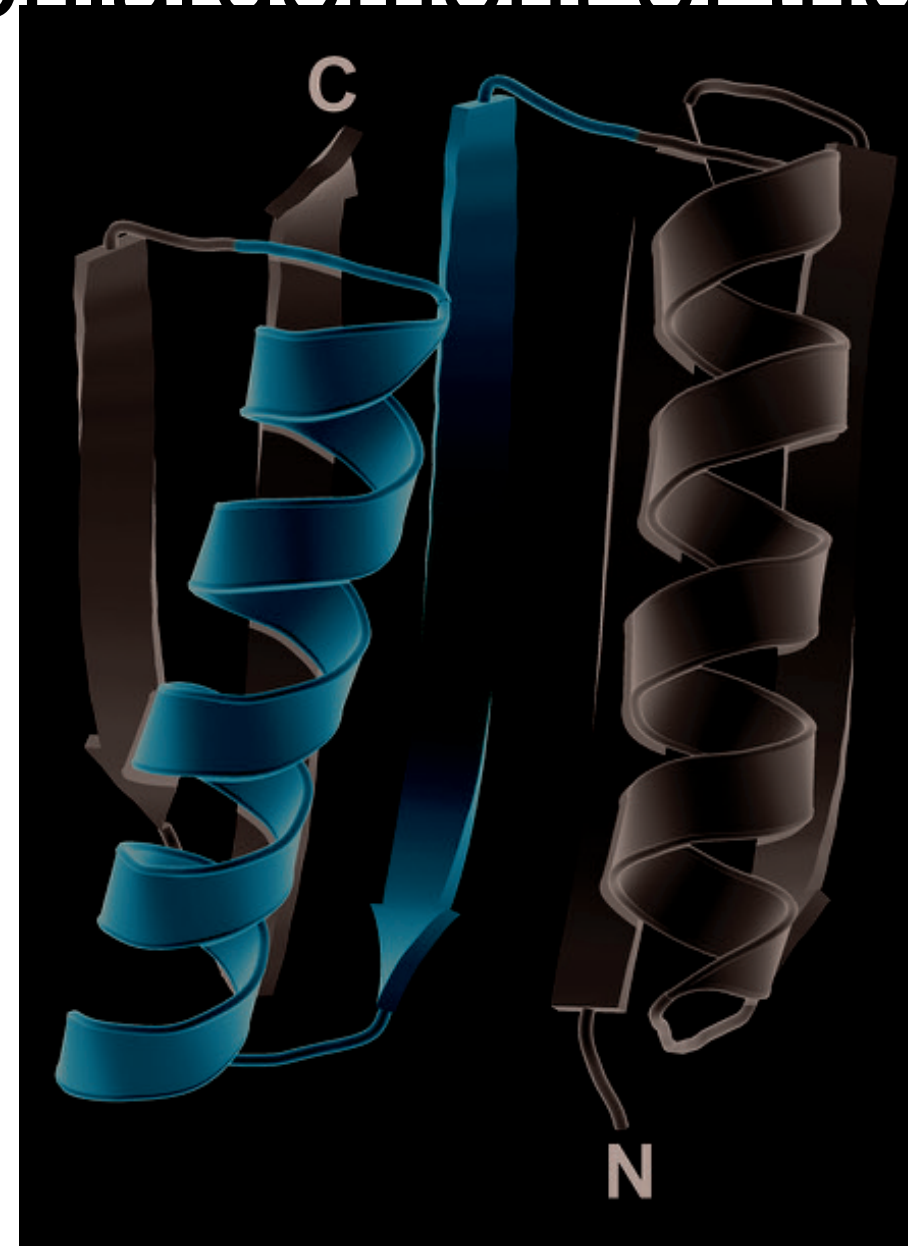
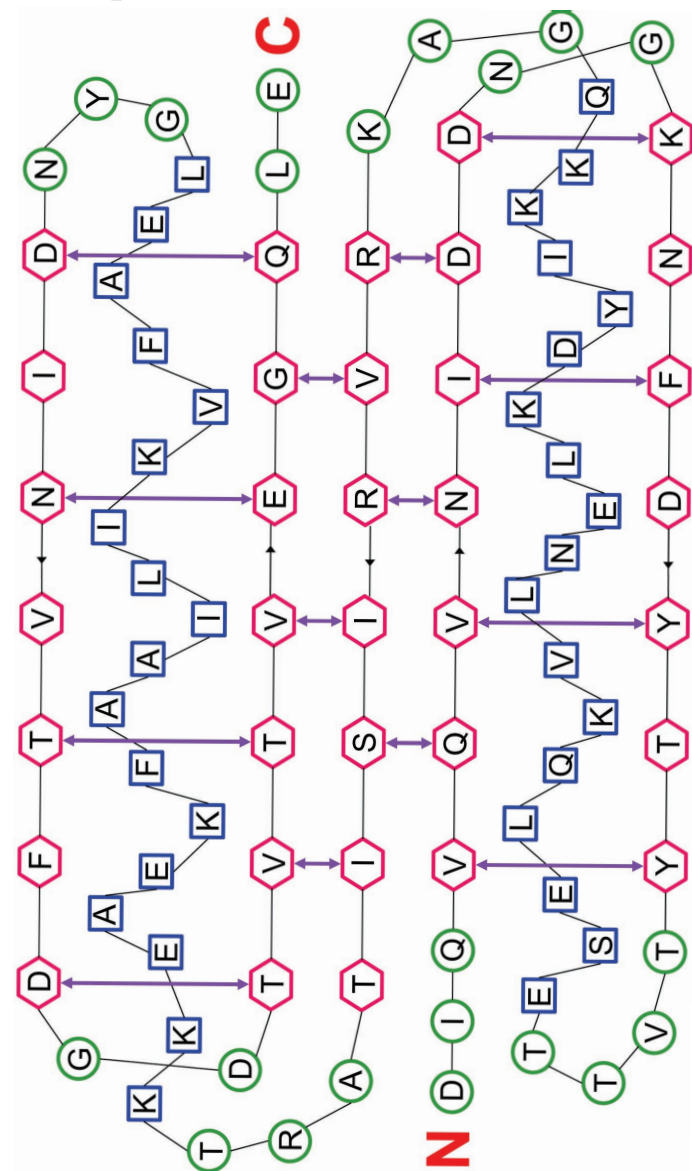


The 'Due Ferri' (two-iron; DF) family



Design of a new protein fold Top7

- define a new topology > design sequence to fit backbone > relax the backbone to fit the sequence > iterate (using ROSETTA)
 - top7 is soluble, stable, monomeric, unfolds cooperatively
sequence is unique (no good homolog by BLAST search)
 - energy function is good for representing its biophysics
 - implications for enlargement of the protein universe
- 



RMSE: 1.2 Å
(much better than prediction)

Kuhlman, B, G Dantas, GC Ireton, G Varani, BL Stoddard, and D Baker.
"Design of A Novel Globular Protein Fold with Atomic-level Accuracy."
Science 302, no. 5649 (21 November 2003): 1364-8.

The Nobel Prize in Chemistry 2024 was divided, one half awarded to David Baker "for computational protein design", the other half jointly to Demis Hassabis and John M. Jumper "for protein structure prediction"



Ill. Niklas Elmehed © Nobel Prize Outreach

David Baker

Prize share: 1/2



Ill. Niklas Elmehed © Nobel Prize Outreach

Demis Hassabis

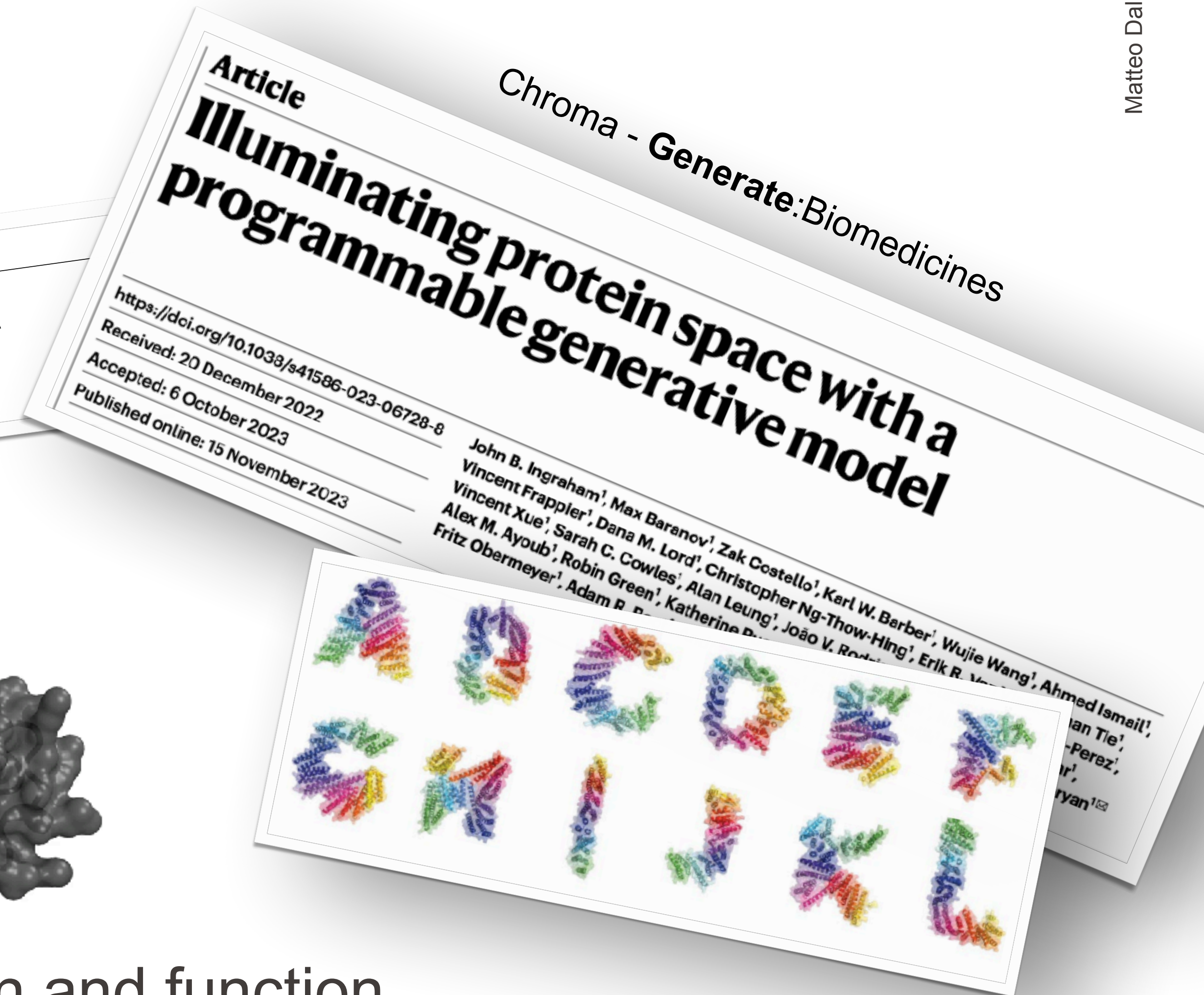
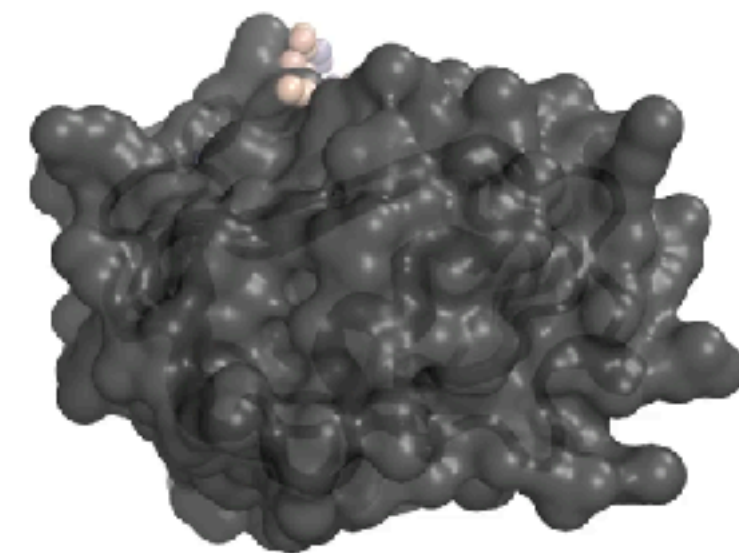
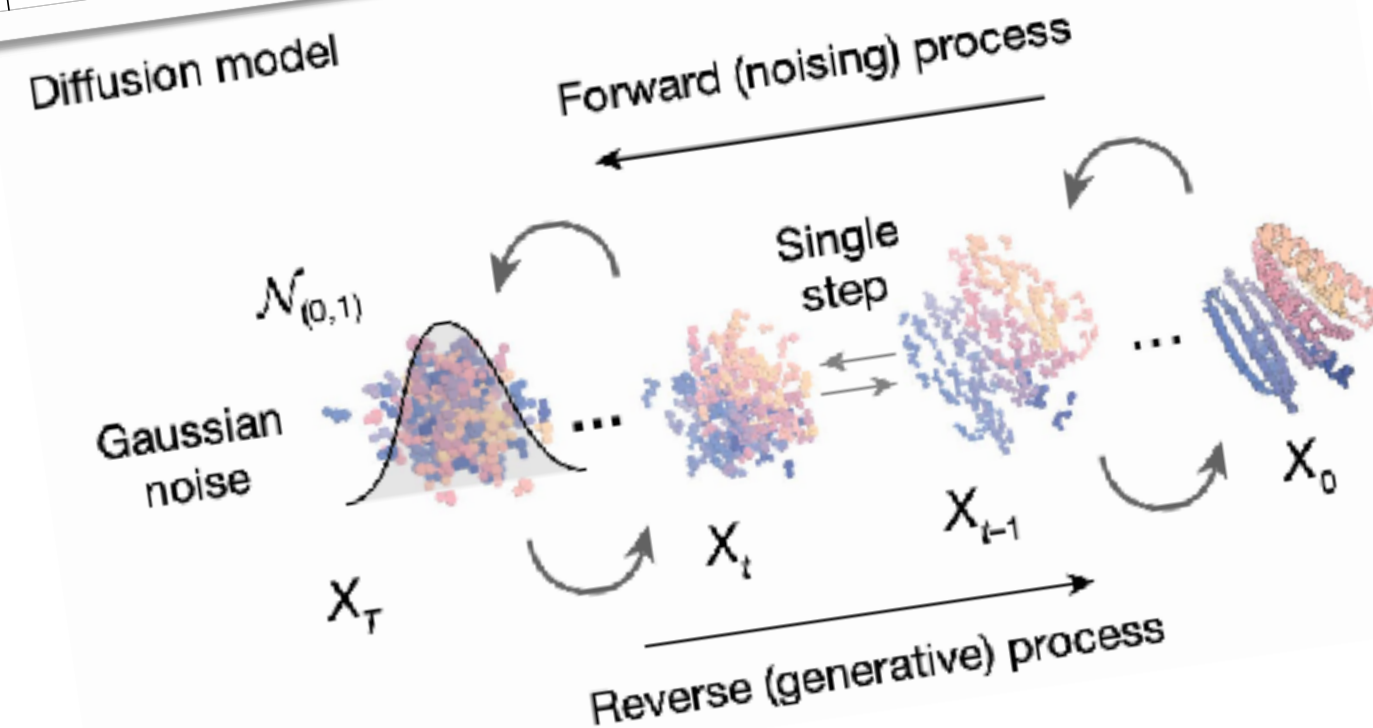
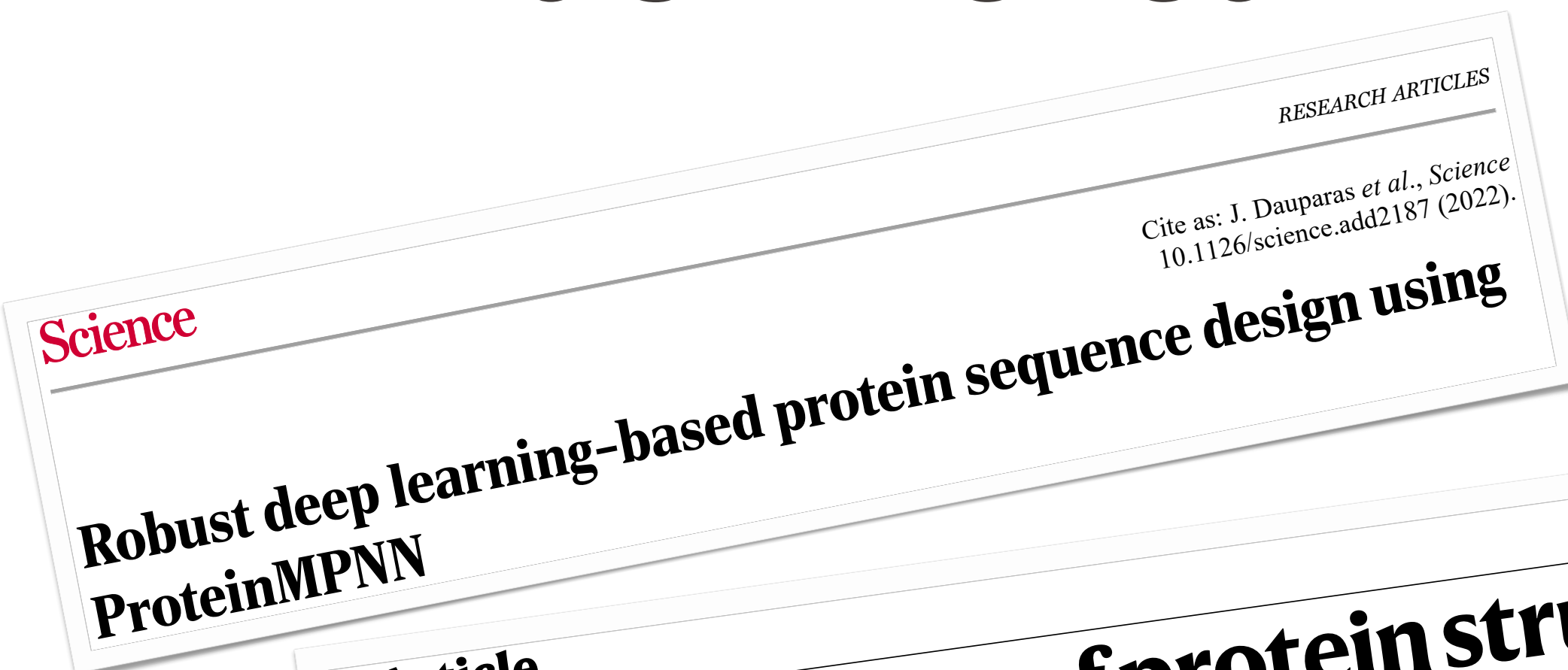
Prize share: 1/4



Ill. Niklas Elmehed © Nobel Prize Outreach

John M. Jumper

Prize share: 1/4



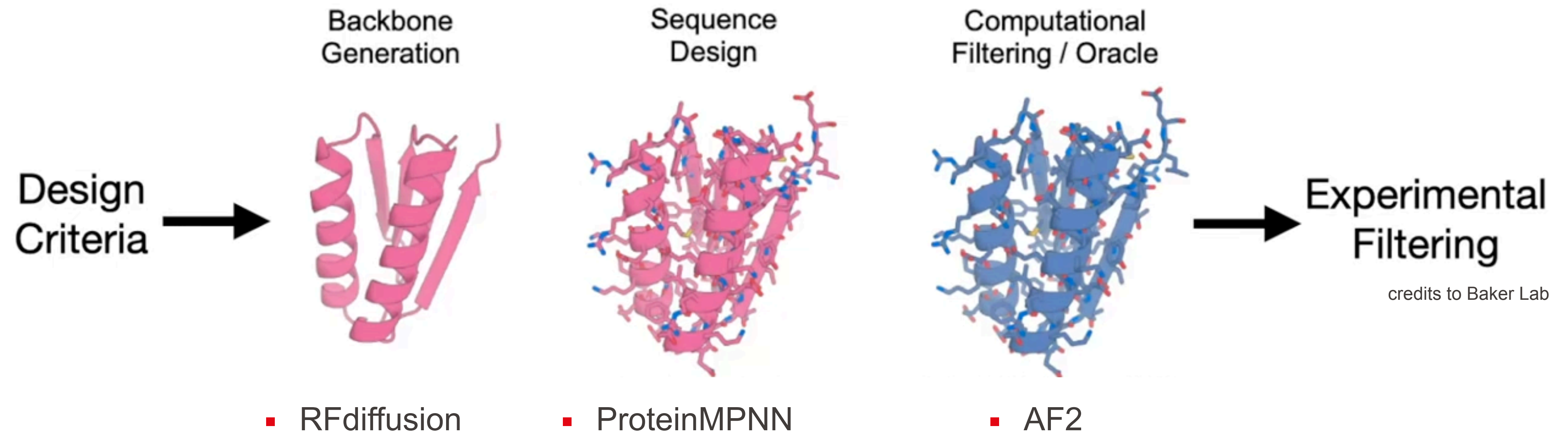
- Application to study protein evolution and function
- Protein engineering for therapeutics, synthetic biology and biotechnology

... leading to molecular engineering



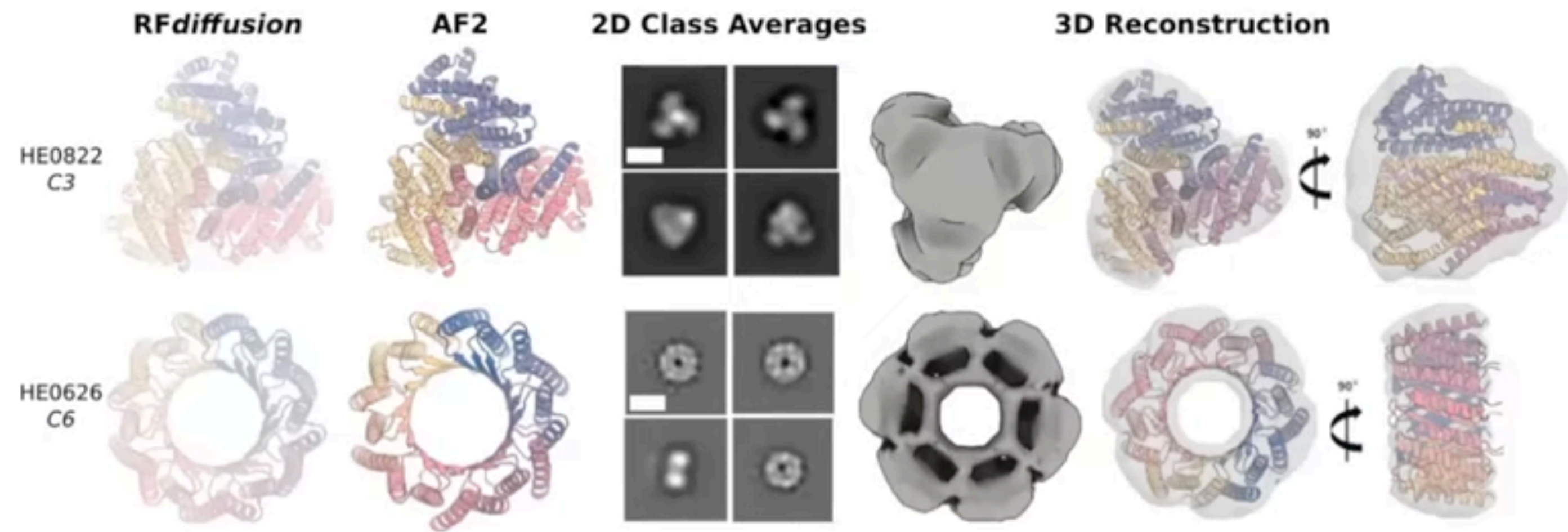
- Application to study protein evolution and function
- Protein engineering for therapeutics, synthetic biology and biotechnology

Pipeline of today's protein design

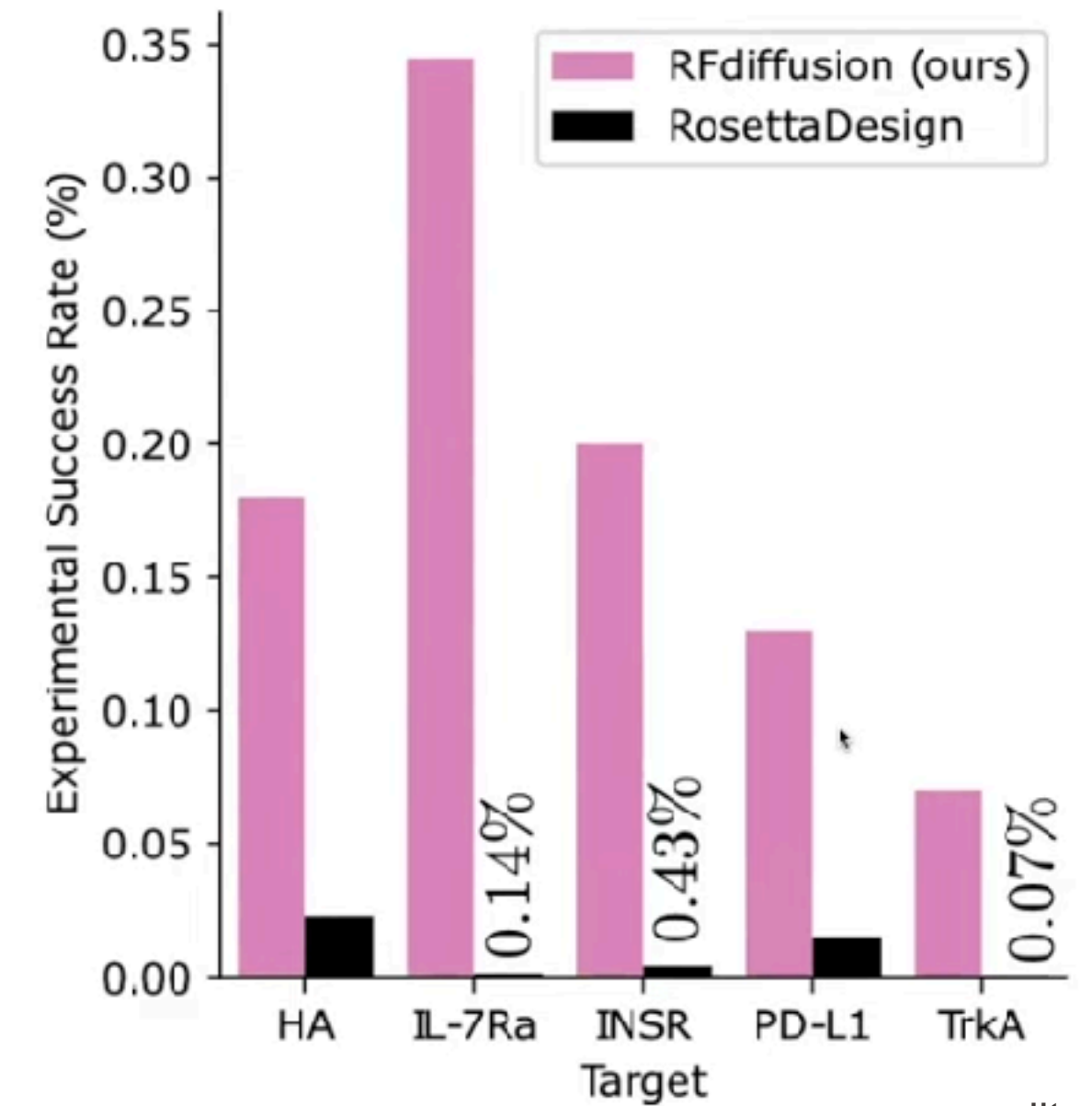


- AF2 has been key to filter potentially good protein designs
- Experimental testing is the ultimate validation of designs
- AI methods enhanced the experimental rate of success
- Protein engineering is now feasible for therapeutics, synthetic biology and biotechnology

Pipeline of today's protein design



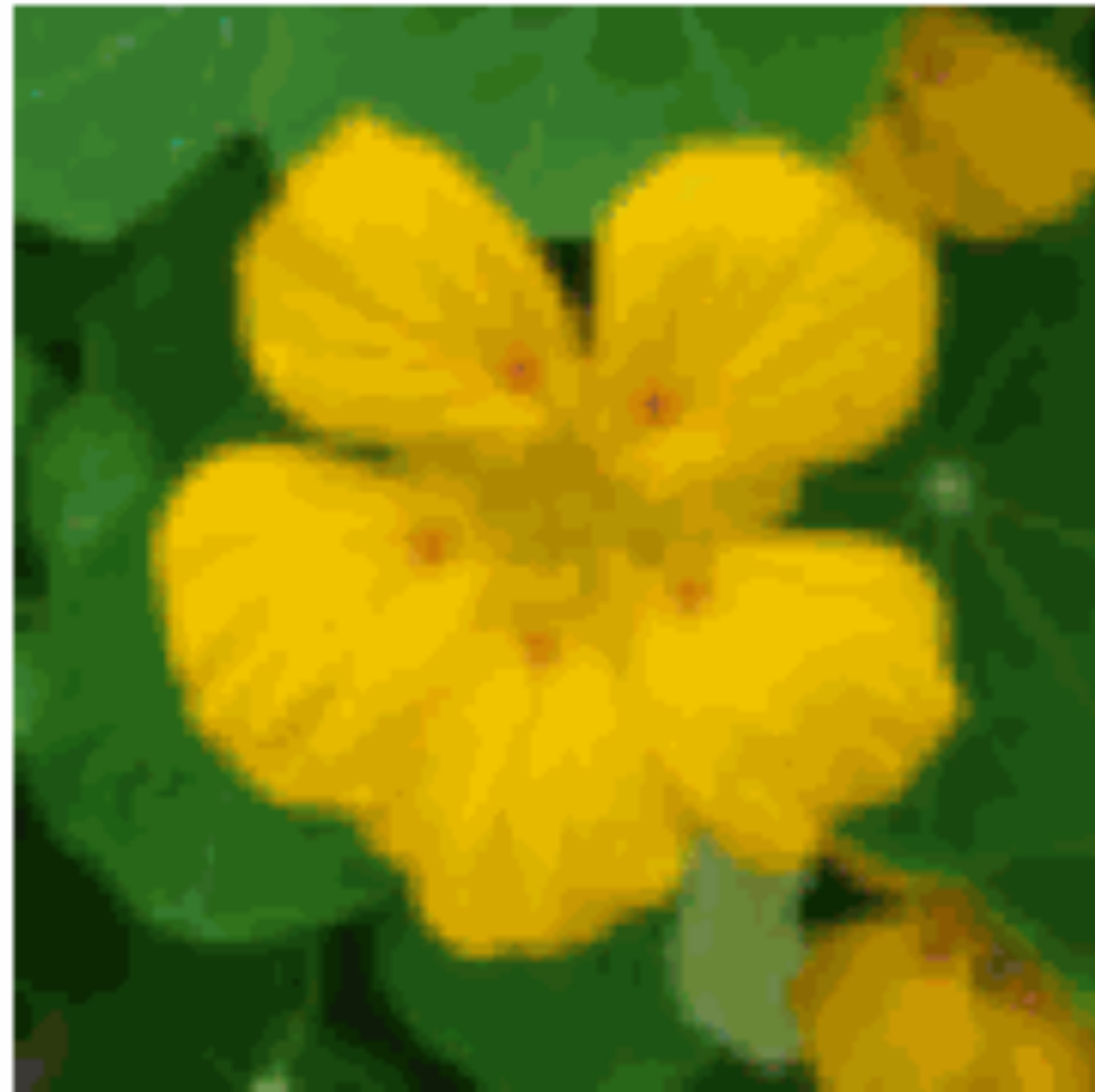
validation using cryoEM



credits to Baker Lab

- AF2 has been key to filter potentially good protein designs
- Experimental testing is the ultimate validation of designs
- AI methods enhanced the experimental rate of success
- Protein engineering is now feasible for therapeutics, synthetic biology and biotechnology

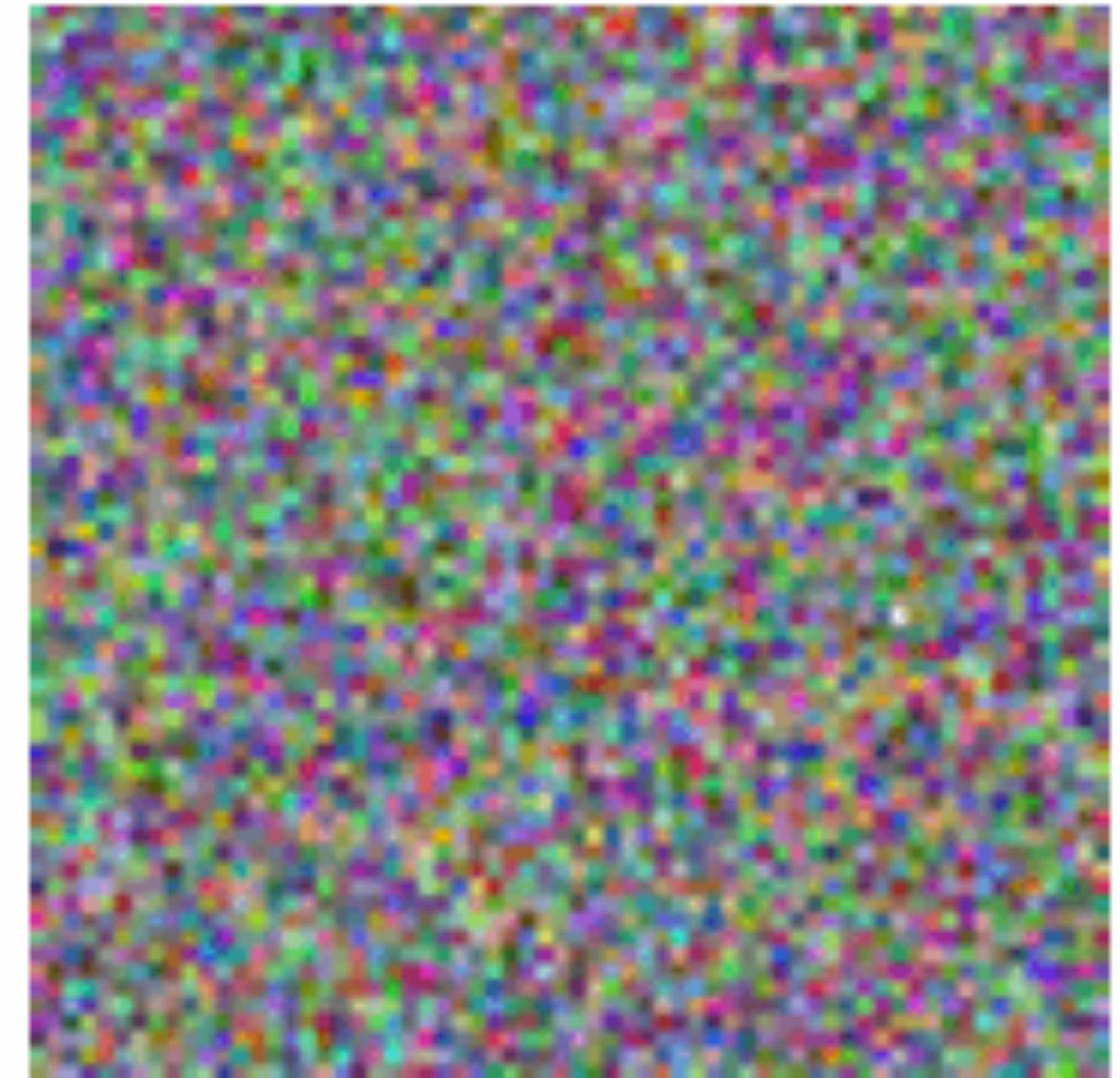
image



Forward diffusion
noisy image



noise



András Béres

- Denoising Diffusion Models - as those used in DALL-E
- Trained to denoise noisy images, they can generate images by iteratively denoising pure noise

De novo design of protein structure and function with RFdiffusion

<https://doi.org/10.1038/s41586-023-06415-8>

Received: 14 December 2022

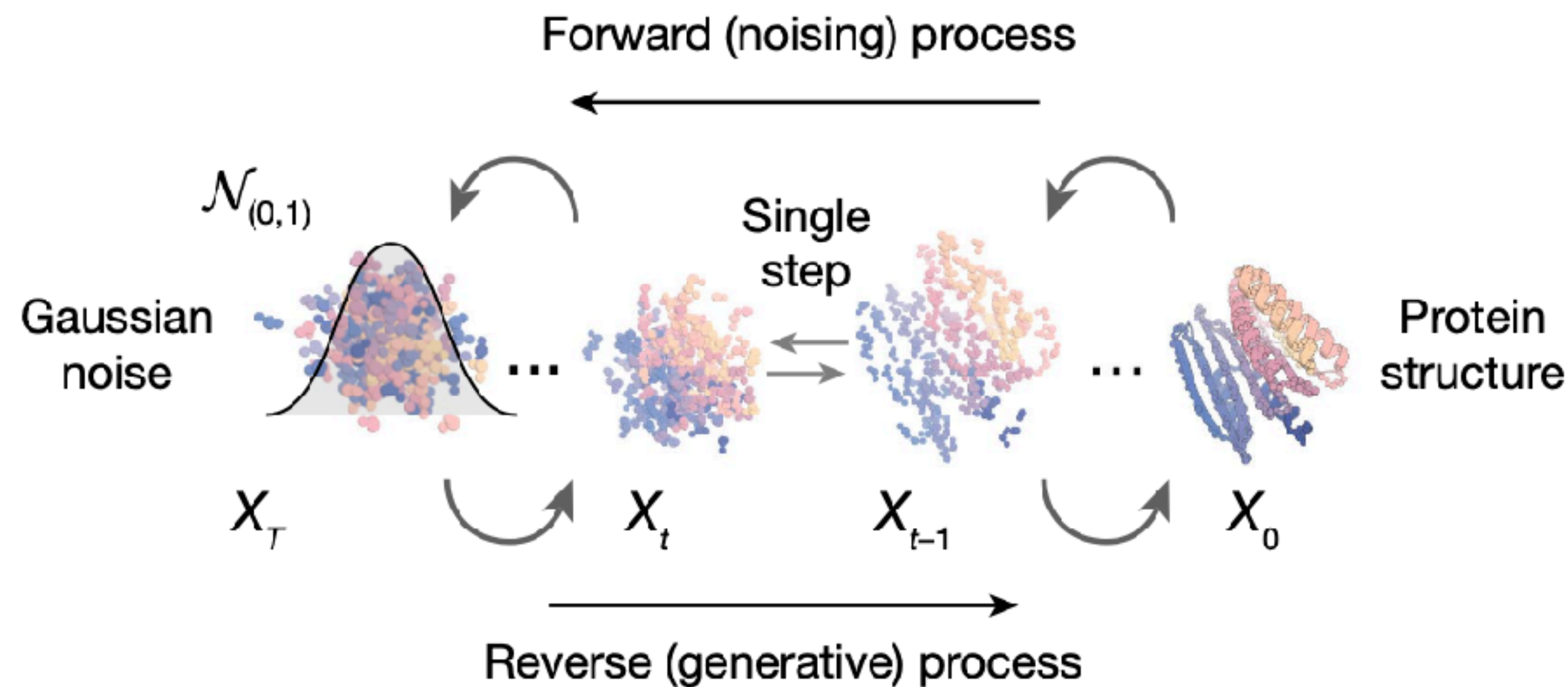
Accepted: 7 July 2023

Published online: 11 July 2023

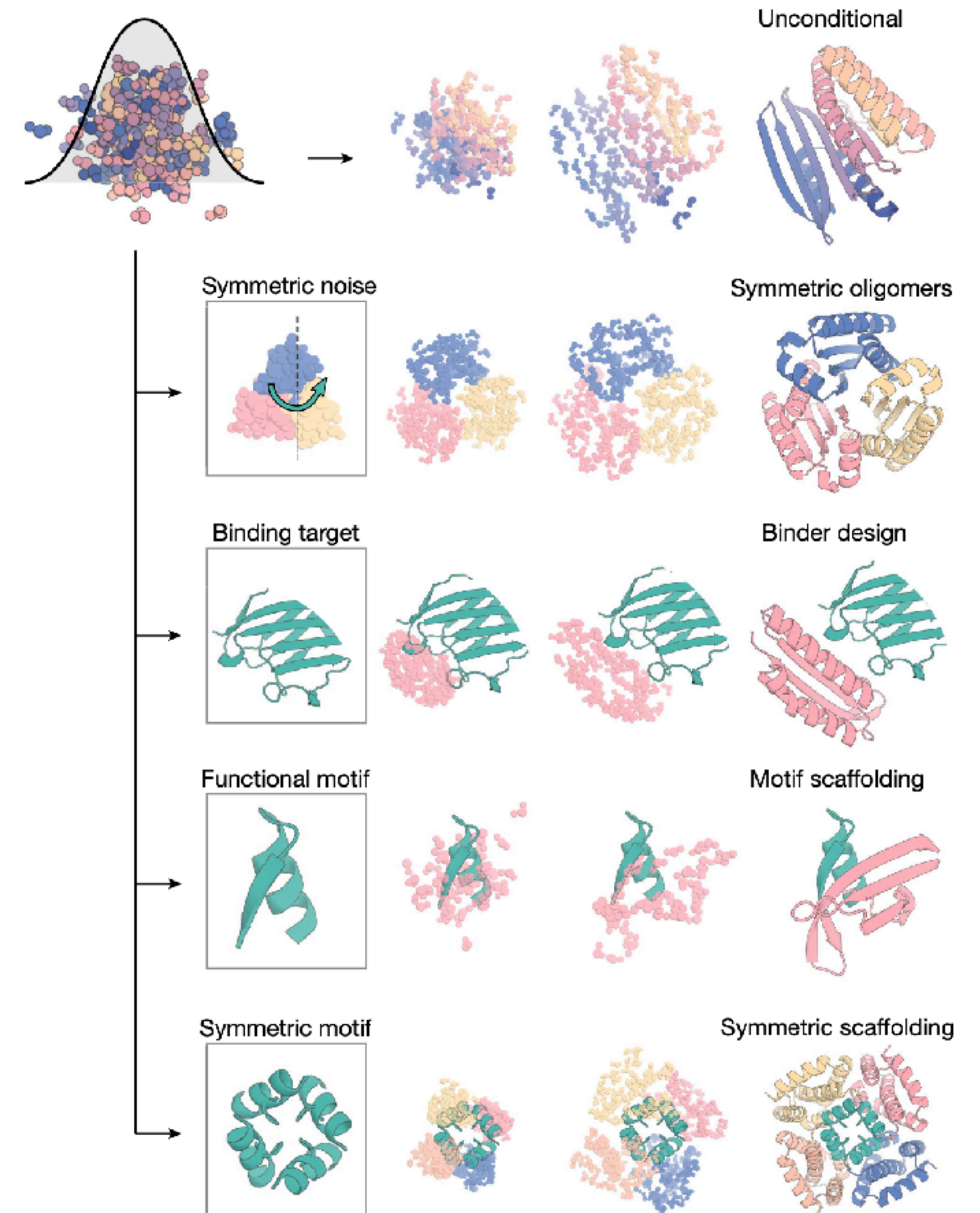
Open access

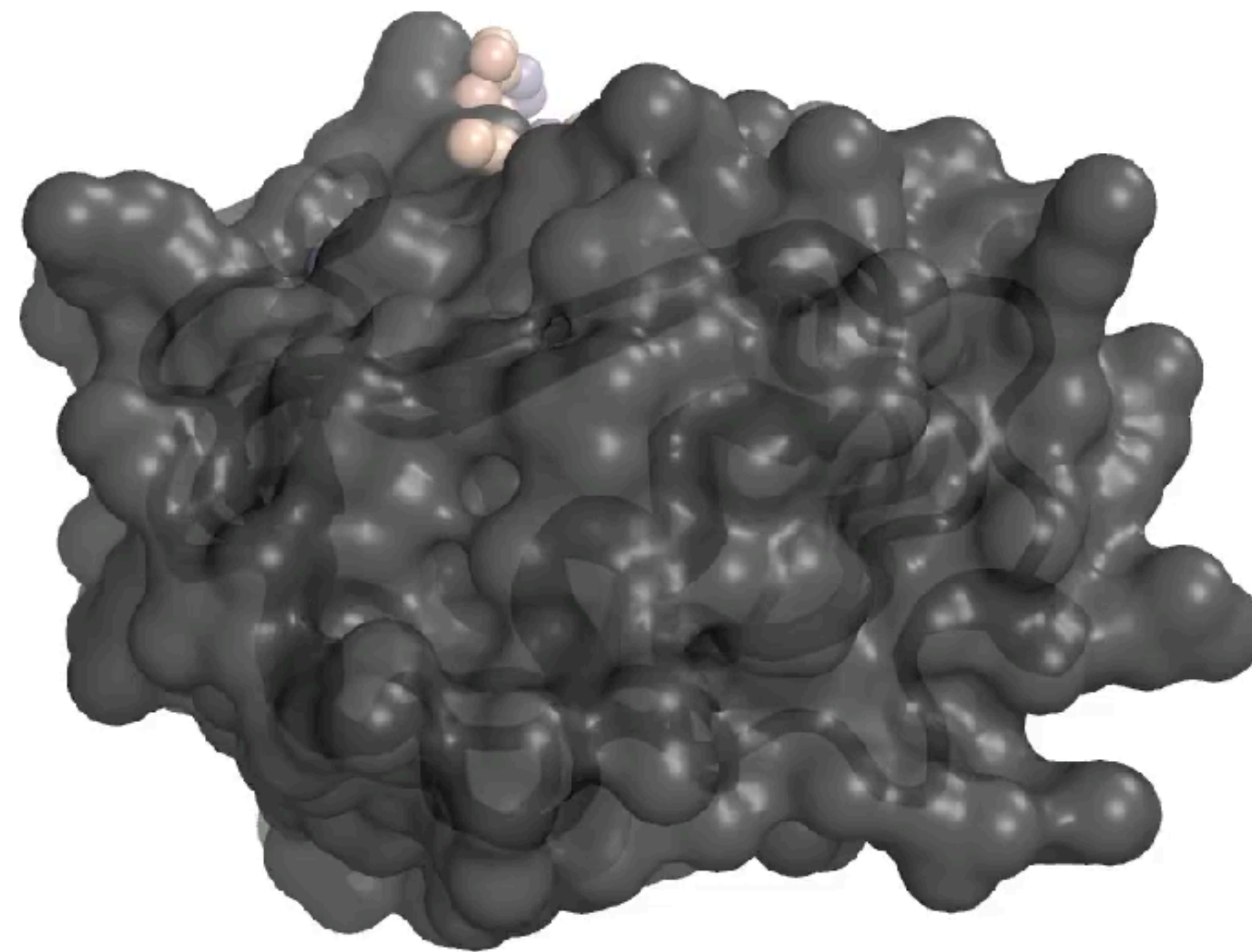
Joseph L. Watson^{1,2,15}, David Juergens^{1,2,3,15}, Nathaniel R. Bennett^{1,2,3,15}, Brian L. Trippe^{2,4,5,15}, Jason Yim^{2,6,15}, Helen E. Eisenach^{1,2,15}, Woody Ahern^{1,2,7,15}, Andrew J. Borst^{1,2}, Robert J. Ragotte^{1,2}, Lukas F. Milles^{1,2}, Basile I. M. Wicky^{1,2}, Nikita Hanikel^{1,2}, Samuel J. Pellock^{1,2}, Alexis Courbet^{1,2,8}, William Sheffler^{1,2}, Jue Wang^{1,2}, Preetham Venkatesh^{1,2,9}, Isaac Sappington^{1,2,9}, Susana Vázquez Torres^{1,2,9}, Anna Lauko^{1,2,9}, Valentin De Bortoli⁸, Emile Mathieu¹⁰, Sergey Ovchinnikov^{11,12}, Regina Barzilay⁶, Tommi S. Jaakkola⁶, Frank DiMaio^{1,2}, Minkyung Baek¹³ & David Baker^{1,2,14}✉

Diffusion model



- the reverse process is learned using a neural network
- its loss function encourages the reverse process to accurately estimate how the data transitions from one noisy step to the previous step.





Cite as: J. Dauparas *et al.*, *Science*
10.1126/science.add2187 (2022).

Robust deep learning-based protein sequence design using ProteinMPNN

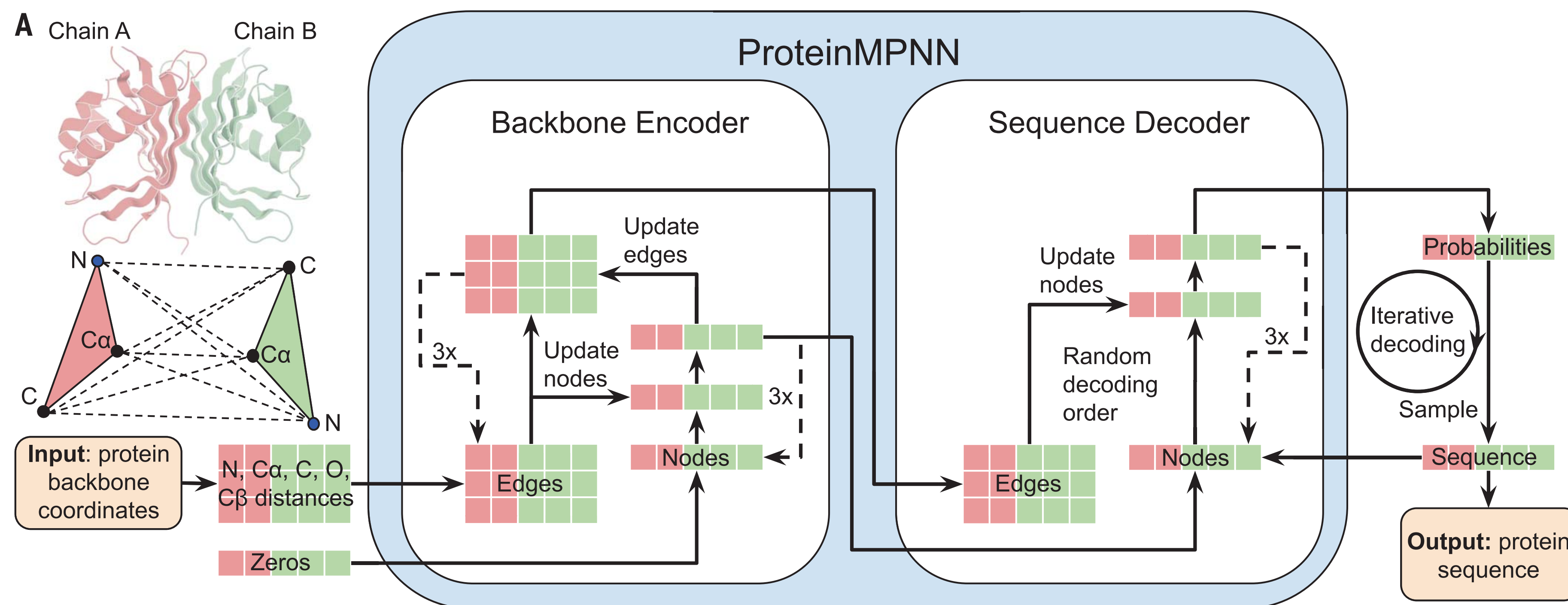
J. Dauparas^{1,2}, I. Anishchenko^{1,2}, N. Bennett^{1,2,3}, H. Bai^{1,2,4}, R. J. Ragotte^{1,2}, L. F. Milles^{1,2}, B. I. M. Wicky^{1,2}, A. Courbet^{1,2,4}, R. J. de Haas⁵, N. Bethel^{1,2,4}, P. J. Y. Leung^{1,2,3}, T. F. Huddy^{1,2}, S. Pellock^{1,2}, D. Tischer^{1,2}, F. Chan^{1,2}, B. Koepnick^{1,2}, H. Nguyen^{1,2}, A. Kang^{1,2}, B. Sankaran⁶, A. K. Bera^{1,2}, N. P. King^{1,2}, D. Baker^{1,2,4*}

¹Department of Biochemistry, University of Washington, Seattle, WA, USA. ²Institute for Protein Design, University of Washington, Seattle, WA, USA. ³Molecular Engineering Graduate Program, University of Washington, Seattle, WA, USA. ⁴Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA. ⁵Department of Physical Chemistry and Soft Matter, Wageningen University and Research, Wageningen, Netherlands. ⁶Berkeley Center for Structural Biology, Molecular Biophysics and Integrated Bioimaging, Lawrence Berkeley Laboratory, Berkeley, CA, USA.

*Corresponding author. Email: dabaker@uw.edu

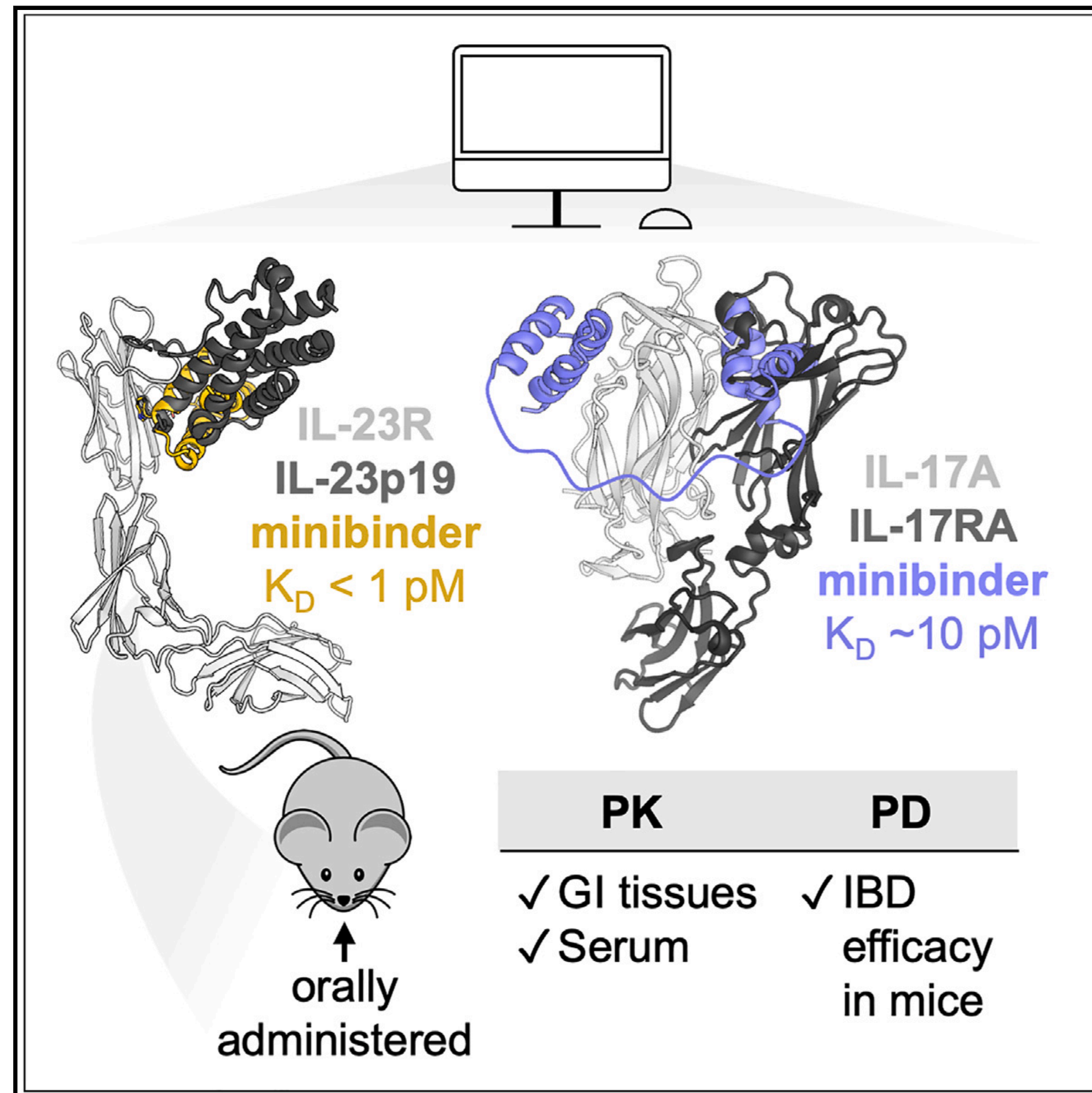
While deep learning has revolutionized protein structure prediction, almost all experimentally characterized de novo protein designs have been generated using physically based approaches such as Rosetta. Here we describe a deep learning-based protein sequence design method, ProteinMPNN, with outstanding performance in both in silico and experimental tests. On native protein backbones, ProteinMPNN has a sequence recovery of 52.4%, compared to 32.9% for Rosetta. The amino acid sequence at different positions can be coupled between single or multiple chains, enabling application to a wide range of current protein design challenges. We demonstrate the broad utility and high accuracy of ProteinMPNN using X-ray crystallography, cryoEM and functional studies by rescuing previously failed designs, made using Rosetta or AlphaFold, of protein monomers, cyclic homo-oligomers, tetrahedral nanoparticles, and target binding proteins.

- Backbone distances are encoded and processed using a message-passing neural network (Encoder) to obtain graph node and edge features.
- The encoded features, together with a partial sequence, are used to generate amino acids iteratively in a random decoding order.



Preclinical proof of principle for orally delivered Th17 antagonist miniproteins

Graphical abstract



Authors

Stephanie Berger, Franziska Seeger, Ta-Yi Yu, ..., Matthias Siebeck, Roswitha Gropp, David Baker

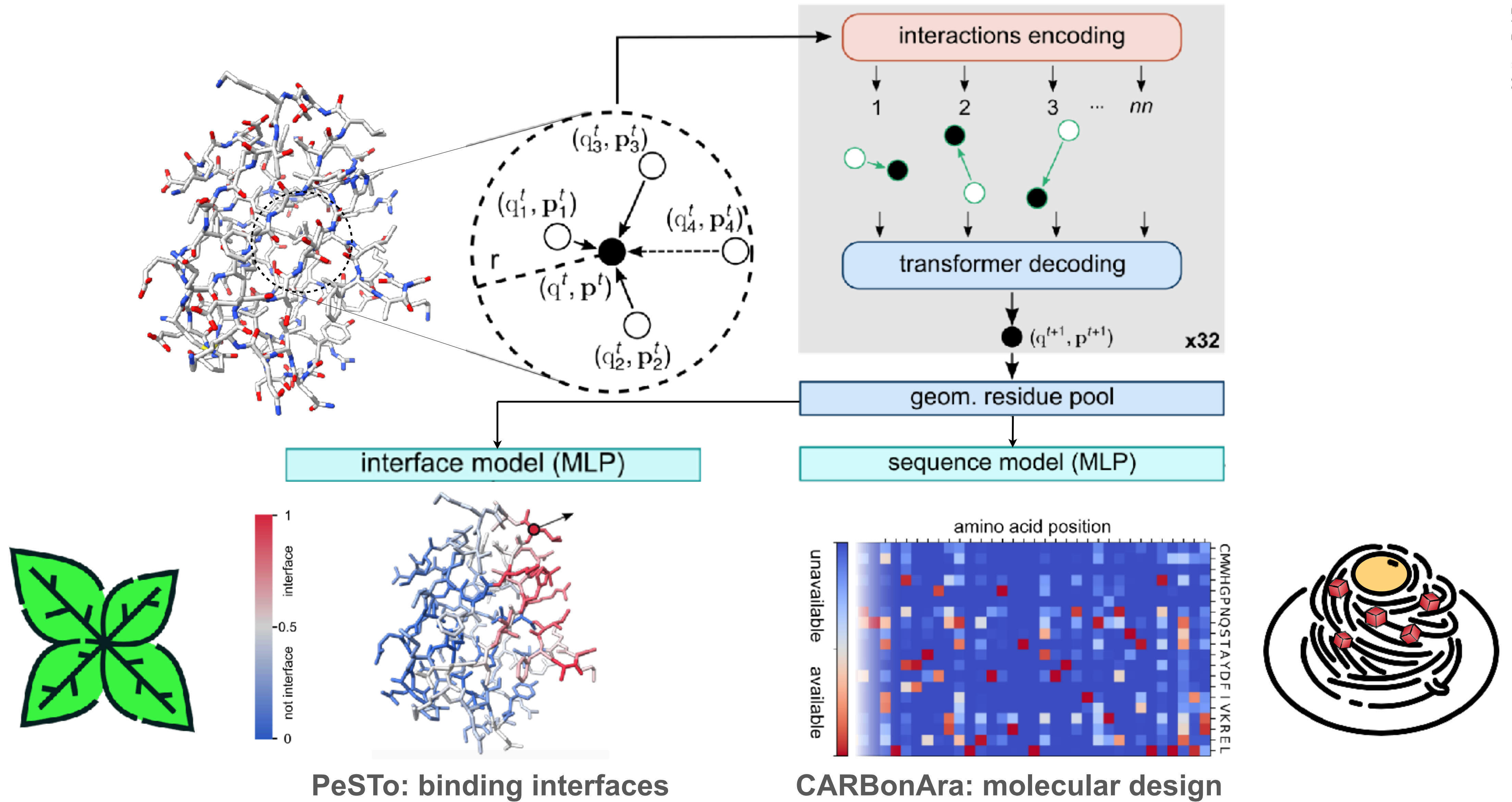
Correspondence

berger389@gmail.com (S.B.),
dabaker@uw.edu (D.B.)

Highlights

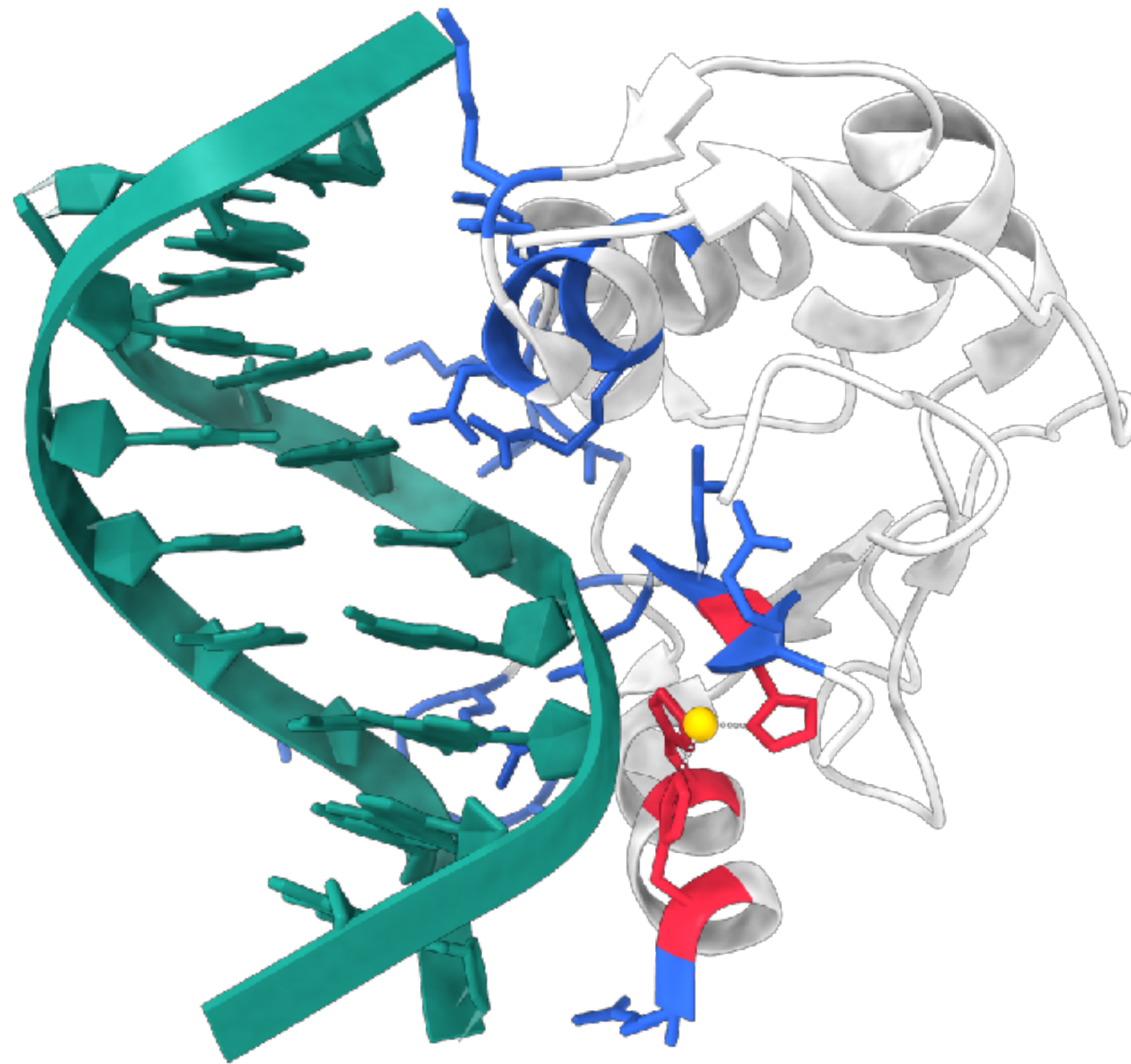
- Computational design yielded low- and sub-pM minibinders of IL-17A and IL-23R
- IL-23R minibinders are extremely resistant to heat, acid, and proteolysis
- Oral IL-23R minibinder is as effective as a clinical mAb in mouse colitis

Berger et al., 2024, Cell 187, 4305–4317
August 8, 2024 © 2024 The Author(s). Published by Elsevier Inc.
<https://doi.org/10.1016/j.cell.2024.05.052>



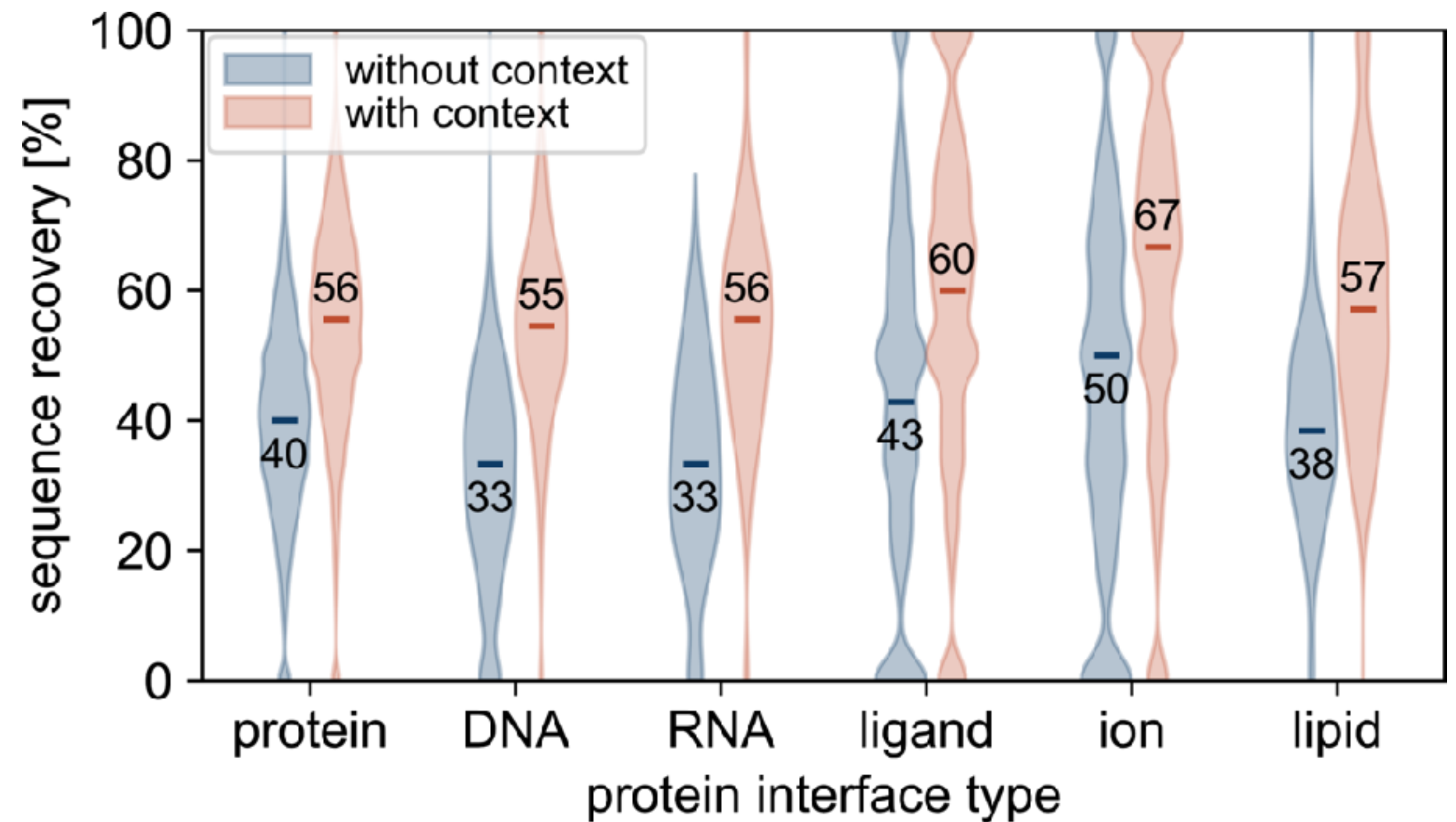
EPFL Unique ability — context awareness

- example with context



colicin E7

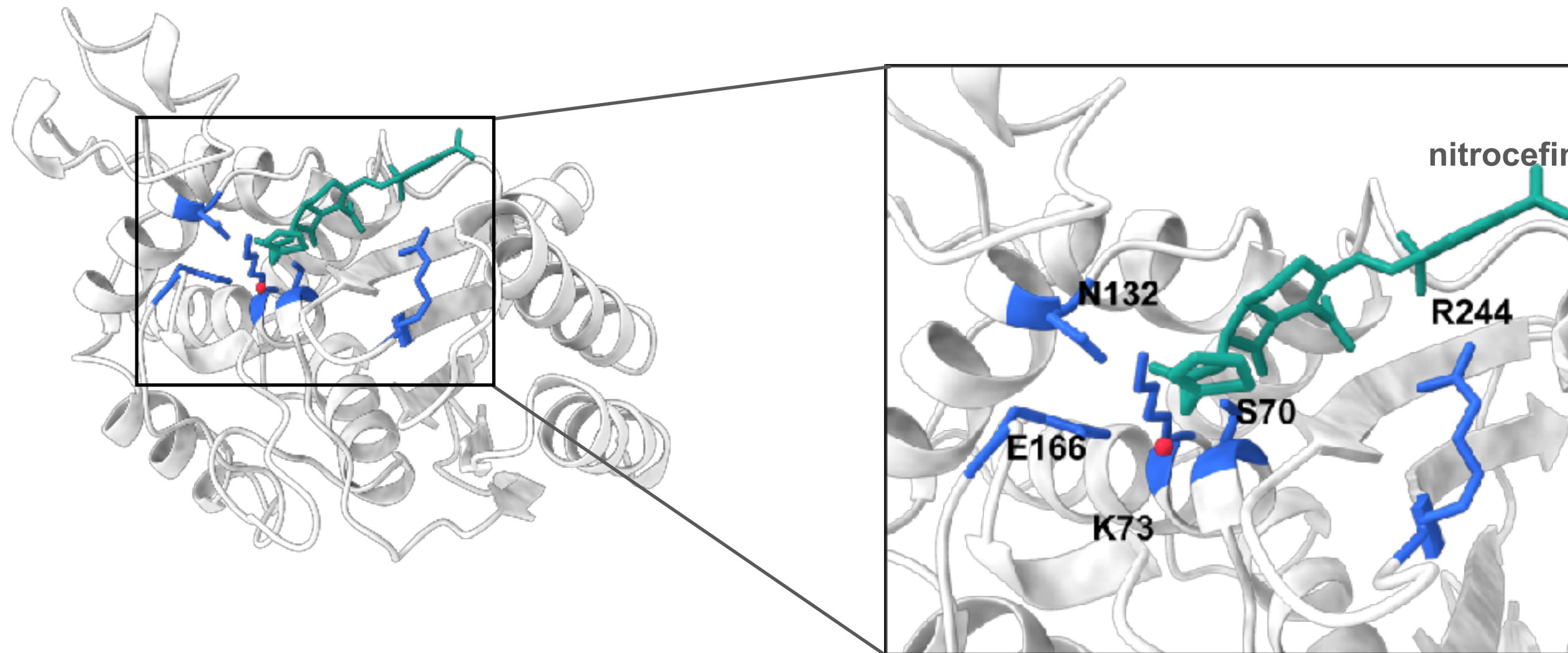
- large-scale benchmark



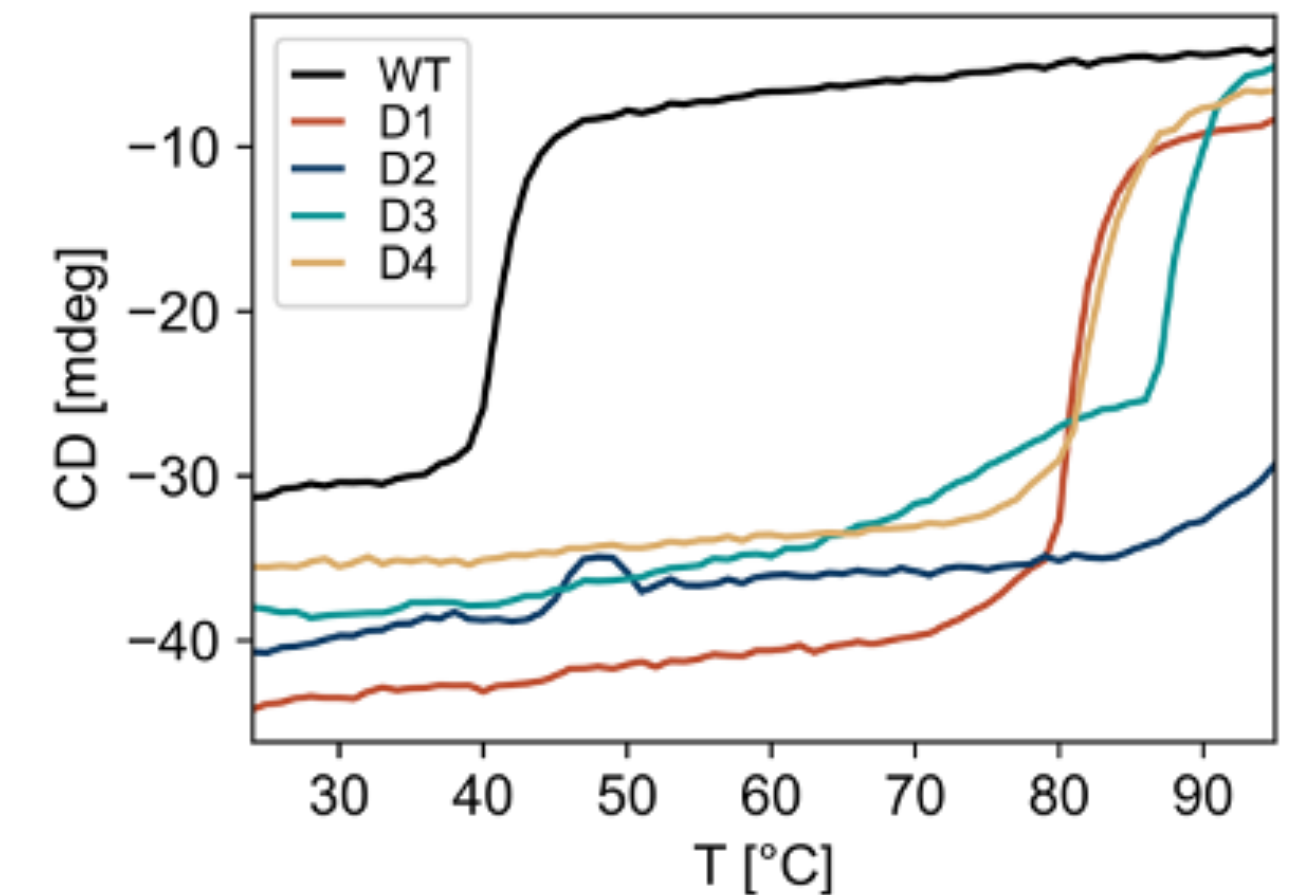
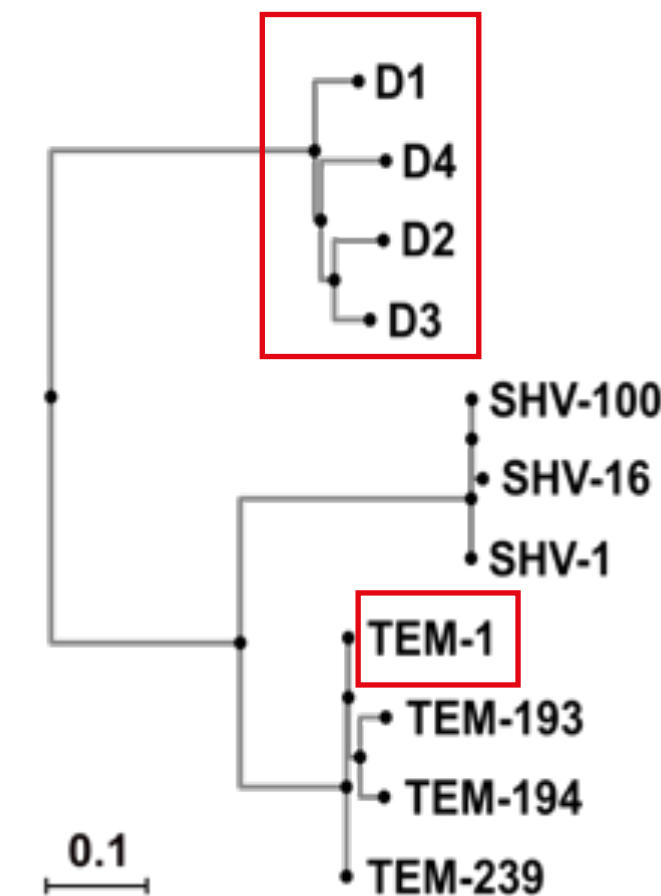
1000 structures sampled with maximum 30% sequence identity and separate C.A.T.H. classification from training set

Can we re-engineer an enzyme?

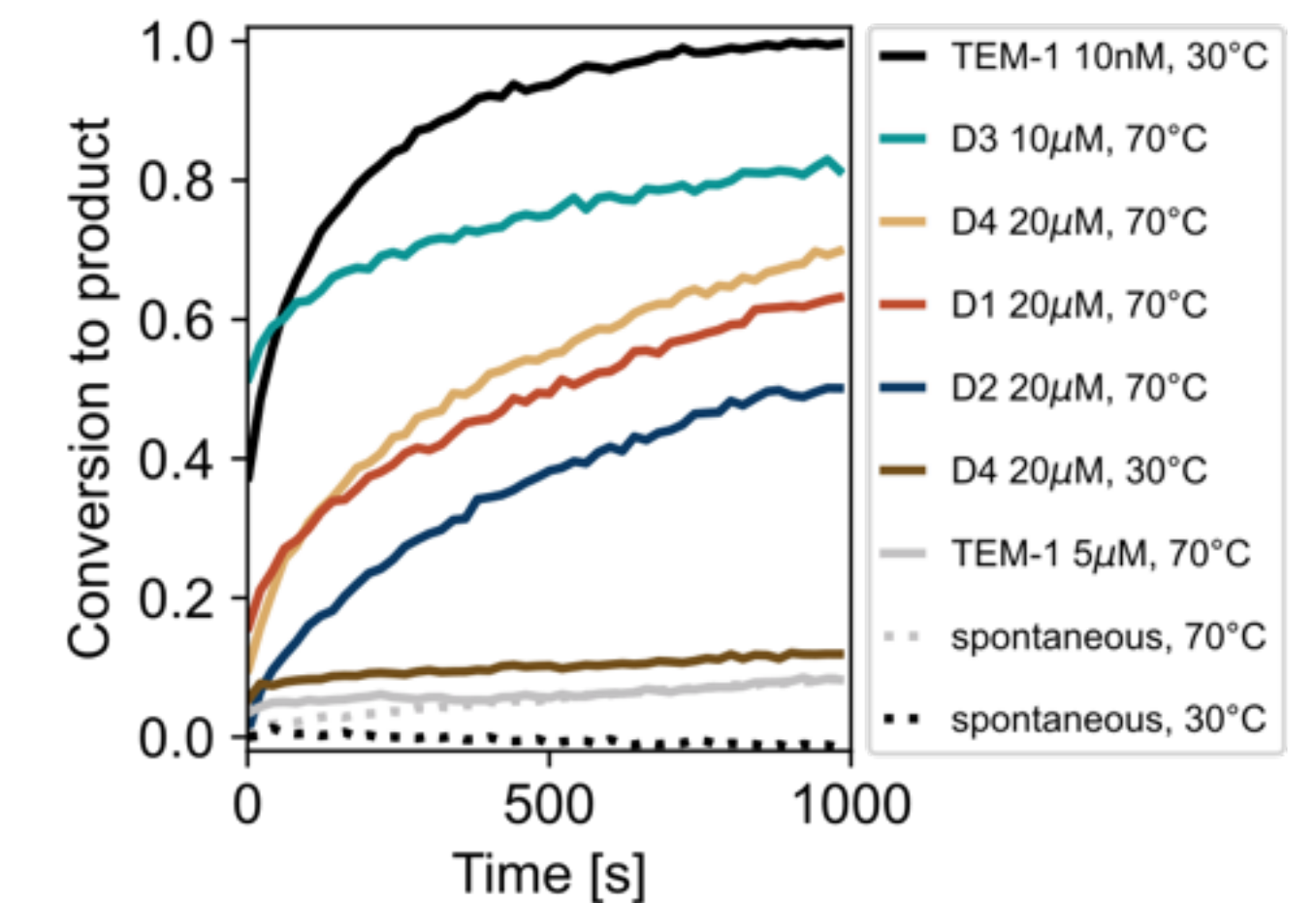
■ TEM-1 serine β -lactamase



only 50% sequence identity



- sequences generation with substrate as constraint
- selected 10 top-ranked predictions based on pLDDT
- 4/10 designs are soluble and monomeric
- they are folded and more thermostable than wild-type TEM-1
- catalytically active at high T - not as the wild-type yet
- represent a separate subclass of β -lactamases



The future is bright and exciting ...

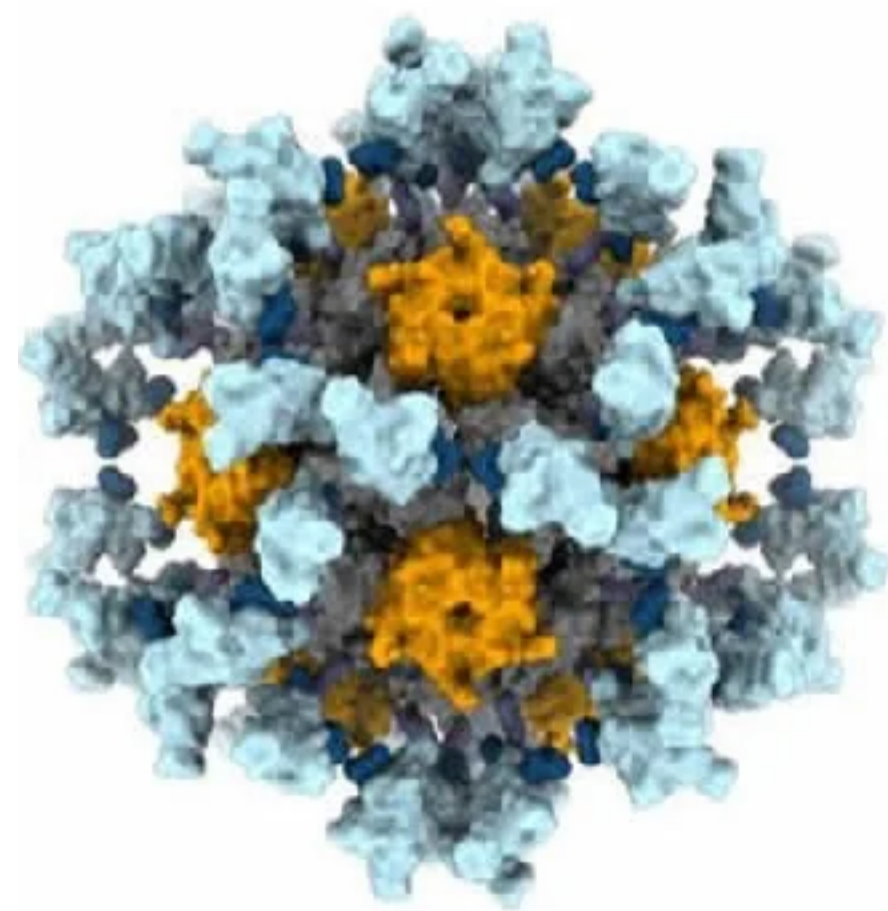
... biomolecular design will address many societal needs

■ Medicine

vaccines & antivirals

smart medicines

drug delivery



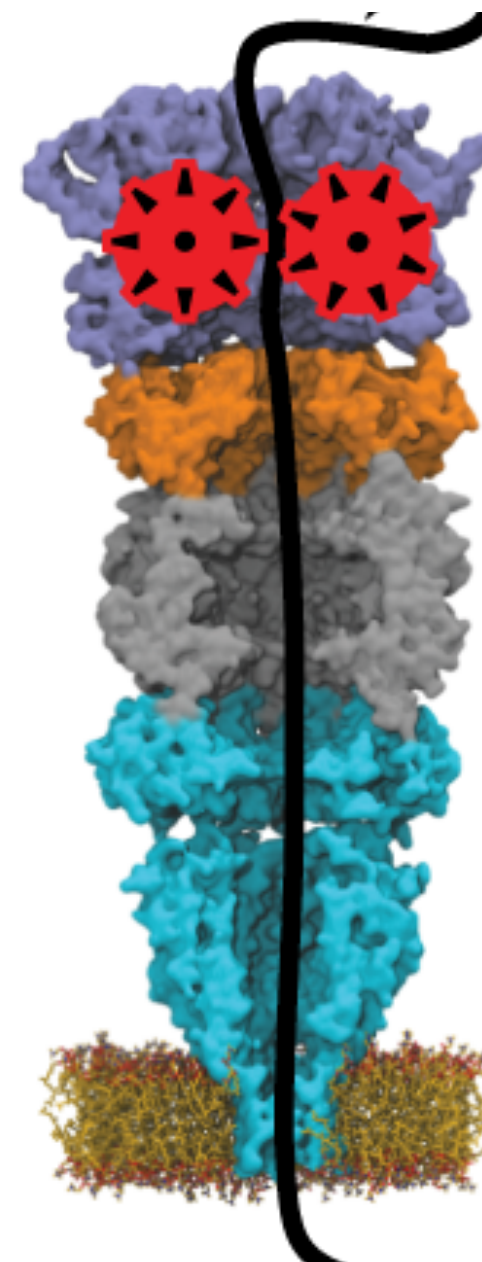
■ SARS-CoV-2 RBD
nanoparticle immunogen (Cell 2020)

■ Biotechnology

protein-silicon devices

bio-based computers

nanoscale manufacturing



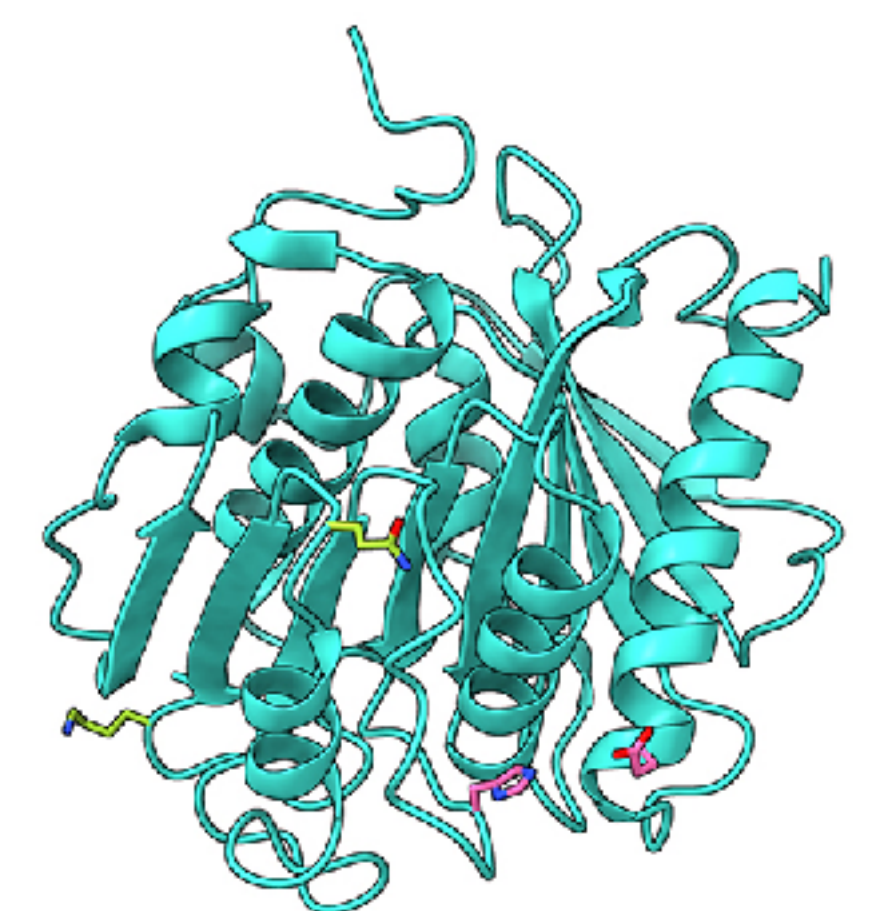
SM proteomics with
biological nanopores
(Nat Chem 2021)

■ Sustainability

artificial photosynthesis

CO₂ sequestration

plastic degradation



FAST-PETase
(Nature 2022)

Google DeepMind

About ▾

Research

Technologies ▾

Impact

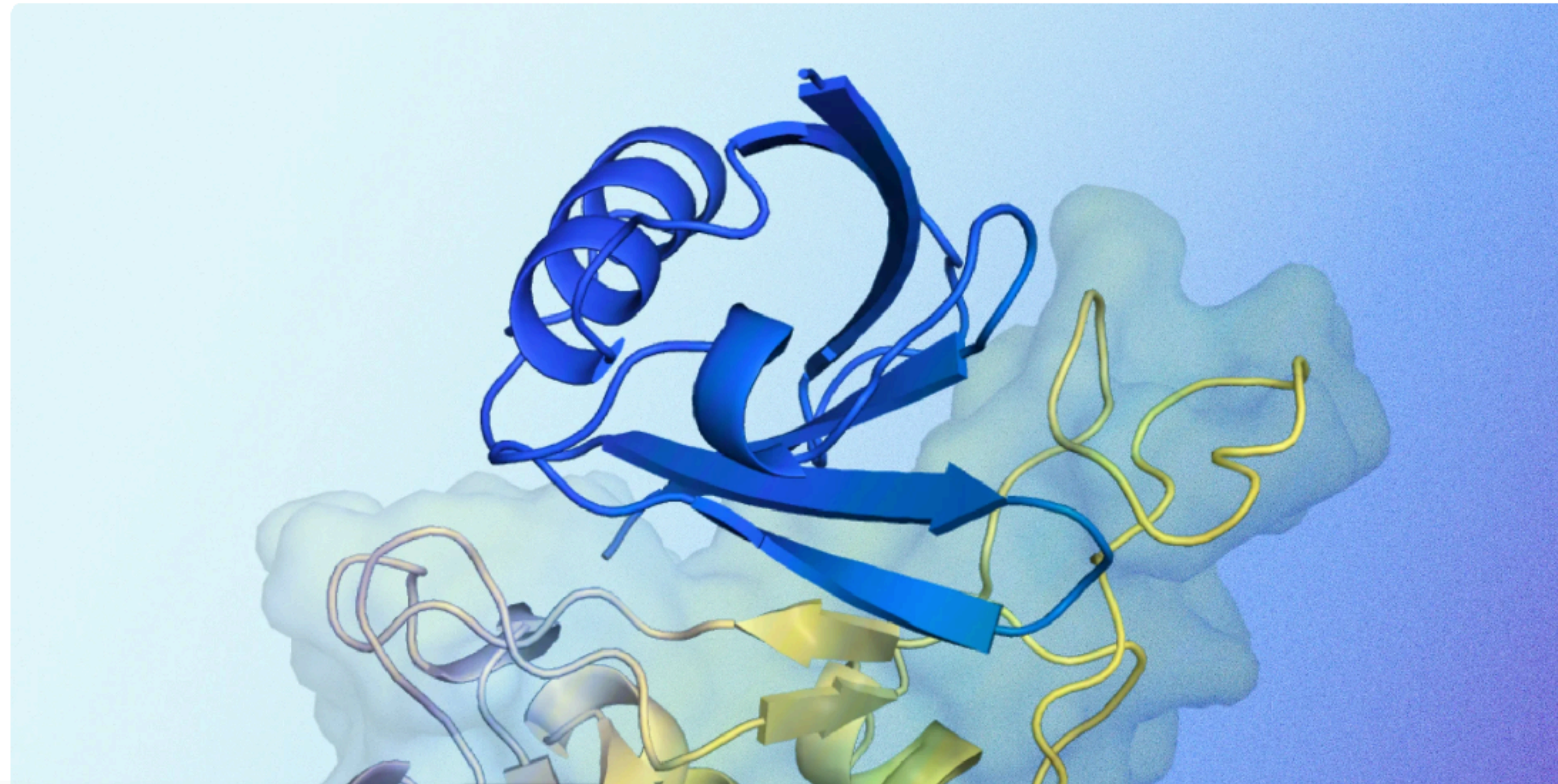
Discover ▾

AlphaProteo generates novel proteins for biology and health research

5 SEPTEMBER 2024

Protein Design and Wet Lab teams

[Share](#)



Designing Life with AI

We're thrilled to introduce "Designing Life with AI" at EPFL, where AI and protein design intersect, involving faculty, professors, and 40 students collaborating on topics like binder design and phosphosite engineering to kinase remodeling. After a year of innovative research, our projects are now being tested in the wet-lab, and we're working on creating a pipeline and resources for new students, aiming to expand our project and make EPFL a hub for protein design.

EPFL

MAKE!
USEFUL. CREATIVE. SUSTAINABLE.

<https://www.designinglifewithai.com/>

[contact the MAKE team for ongoing projects offered by labs](#)

