



# Computational structural biology

**Matteo Dal Peraro**

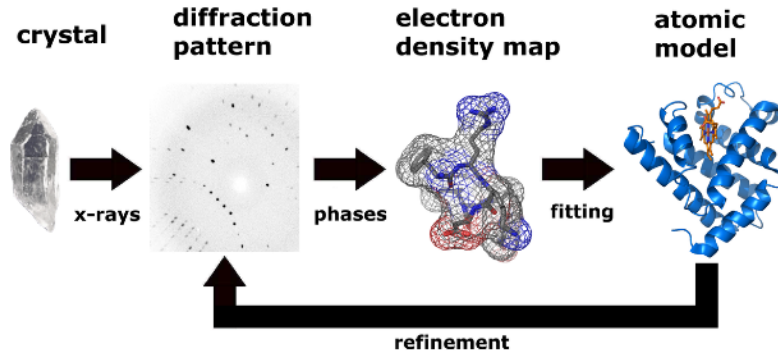
Institute of Bioengineering, School of Life Sciences

*made by GPT4 - DALL-E*

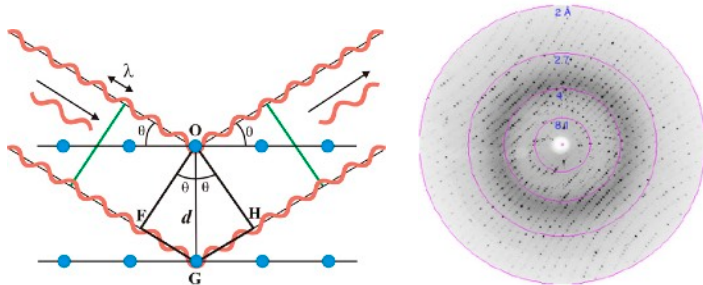
[matteo.dalperaro@epfl.ch](mailto:matteo.dalperaro@epfl.ch)

# Lecture 11 – Quick Summary

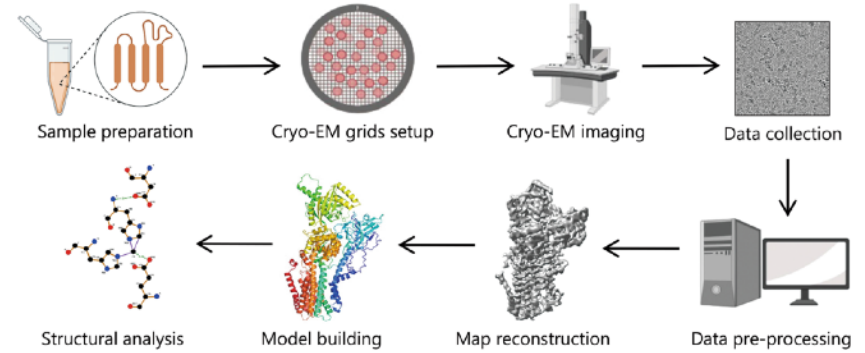
## X-ray crystallography



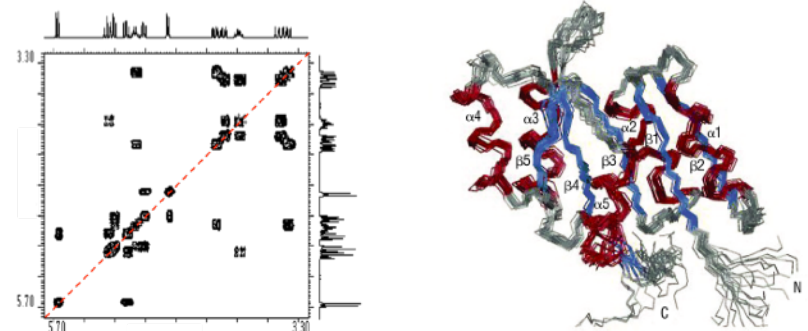
Bragg's law  $n\lambda = 2d \sin\theta$



## Cryo-EM

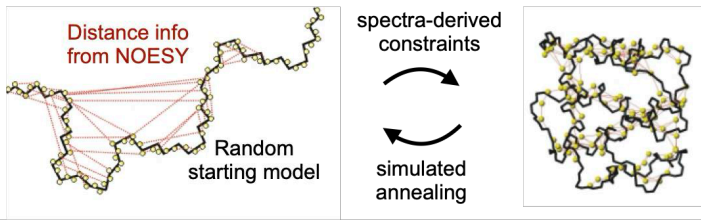


## NMR



# Methods for determining biomolecular structures

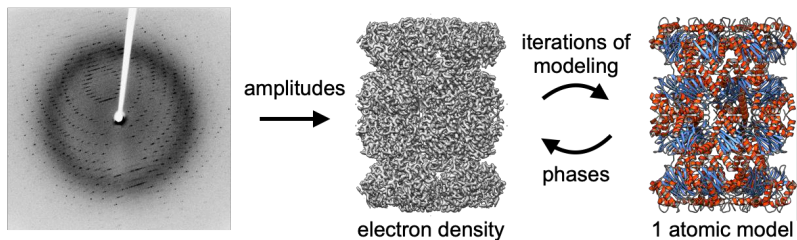
## NMR



## no dynamics !

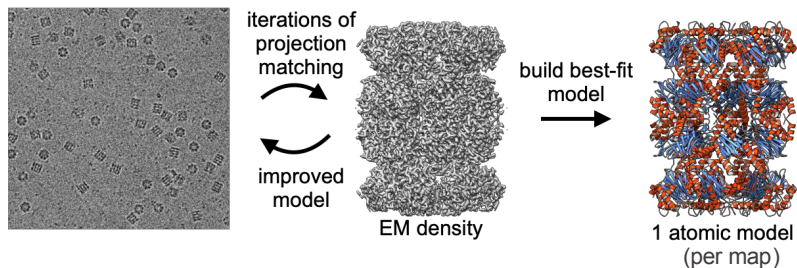
- Versatile tool for studying protein structure and dynamics
- Computationally light
- Full structural analysis limited to smaller proteins (<50kDa)
- Requires isotopic labeling
- Results in model ensemble

## X-ray



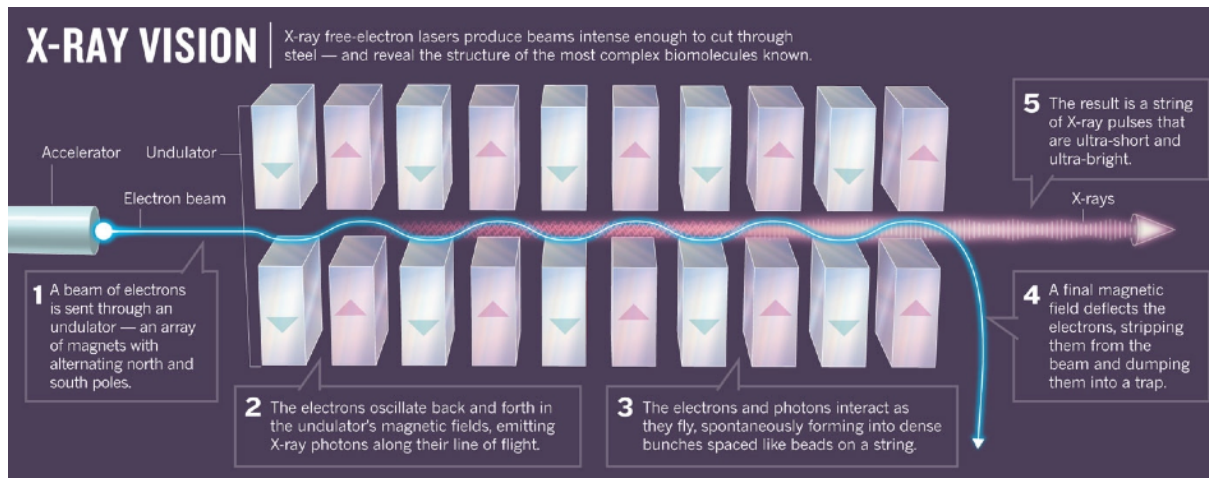
- Gold-standard method for solving protein structures
- **Not limited in size or achievable resolution**
- Computationally light
- Requires highly homogenous, crystallizable sample
- Requires screening of crystallization conditions
- Phase problem
- Results in a single model

## CryoEM



- Versatile tool for studying protein assembly, structure, dynamics
- Limited to proteins >40kDa
- No requirement for protein labeling
- **Does not require homogenous samples**
- Grid preparation procedure requires screening
- Real space imaging – no phase problem
- Can be used to study protein dynamics
- Can be expanded to larger assemblies (e.g., viruses and cells)
- Results in 1 or more models per dataset
- Computationally heavy (TBs of data + requirement for GPU processing)

- X-ray free electron laser (XFEL)



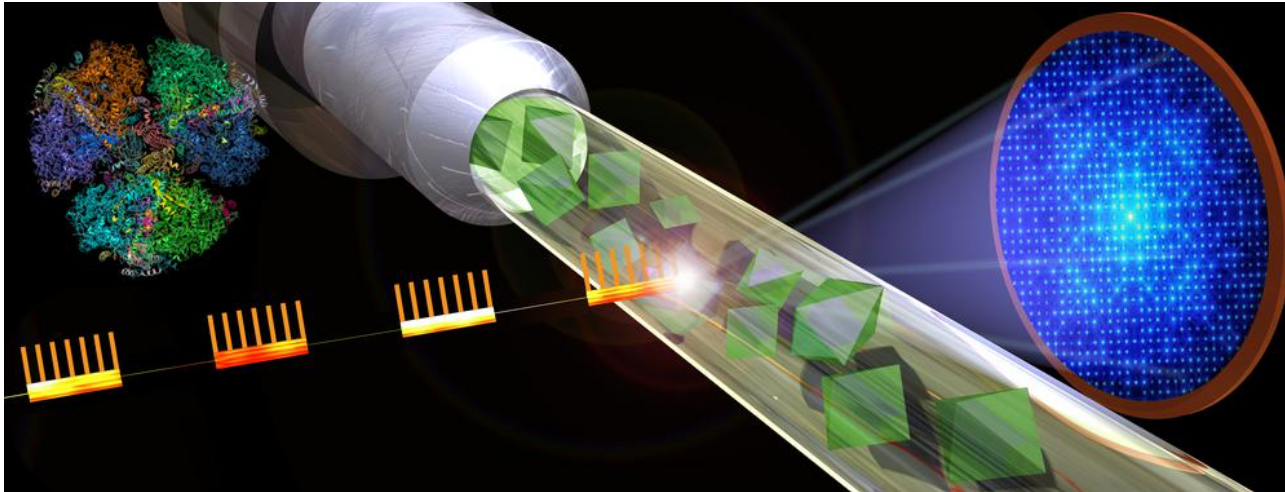
use of nanocrystals,  $\sim 10,000$  pulses per second, big data problem, promising for membrane proteins, look at molecular dynamics in real time

**100 years of crystallography:** in 1914, Max von Laue won the Nobel Prize in Physics for discovering how crystals can diffract X-rays

>> <http://www.nature.com/news/specials/crystallography-1.14540>



- X-ray free electron laser (XFEL)

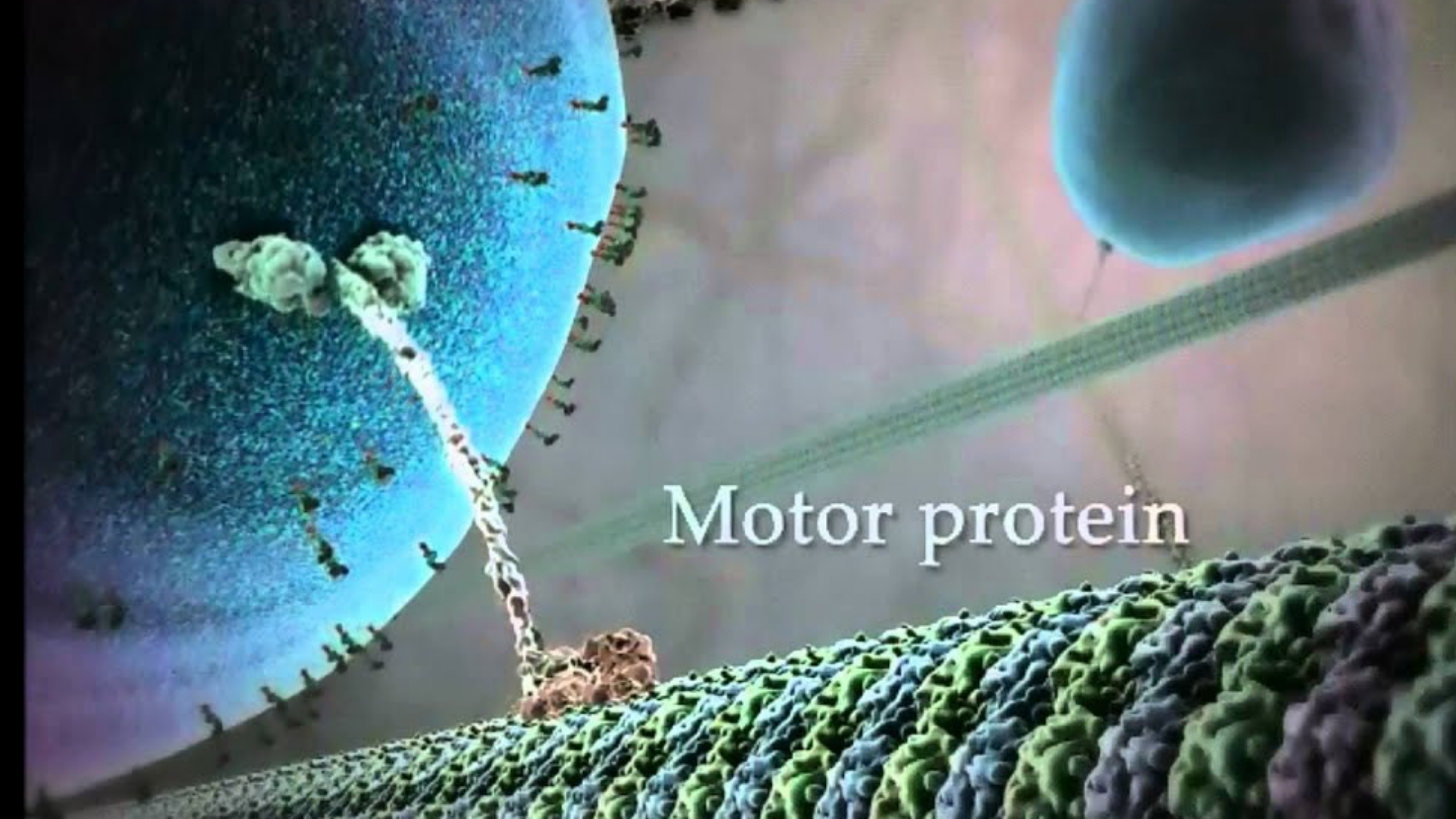


use of nanocrystals,  $\sim 10,000$  pulses per second, big data problem, promising for membrane proteins, look at molecular dynamics in real time

**100 years of crystallography:** in 1914, Max von Laue won the Nobel Prize in Physics for discovering how crystals can diffract X-rays

>> <http://www.nature.com/news/specials/crystallography-1.14540>

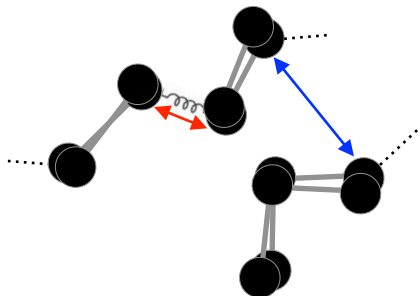




Motor protein

The force field:  $U_{MM} = U_{bonded} + U_{non-bonded}$

- representation of **intra-** and **inter-**molecular interactions
- approximations using simple concepts of **classical mechanics** (i.e. balls and springs, harmonic oscillators)
- large number of parameters fitted to represent experimental data or QM calculated quantities
- “trial and error” or least-squares fitting methods to converge to a consistent set of parameters
- assumption that parameters can be transferable to different contexts (specialized vs. generalized FF)



$$\begin{aligned}
 &= \sum_{\text{All Bonds}} \frac{1}{2} K_b (b - b_0)^2 + \sum_{\text{All Angles}} \frac{1}{2} K_\theta (\theta - \theta_0)^2 \\
 &+ \sum_{\text{All Torsion Angles}} K_\phi [1 - \cos(n\phi + \delta)] \\
 &+ \sum \epsilon \left[ \left( \frac{r_0}{r} \right)^{12} - 2 \left( \frac{r_0}{r} \right)^6 \right] \\
 &+ \sum \frac{332 q_i q_j}{r}
 \end{aligned}$$

$$\frac{d^2 x_i}{dt^2} = \frac{F(x_i)}{m_i} = -\frac{1}{m_i} \frac{dU(x_i)}{dx_i}$$

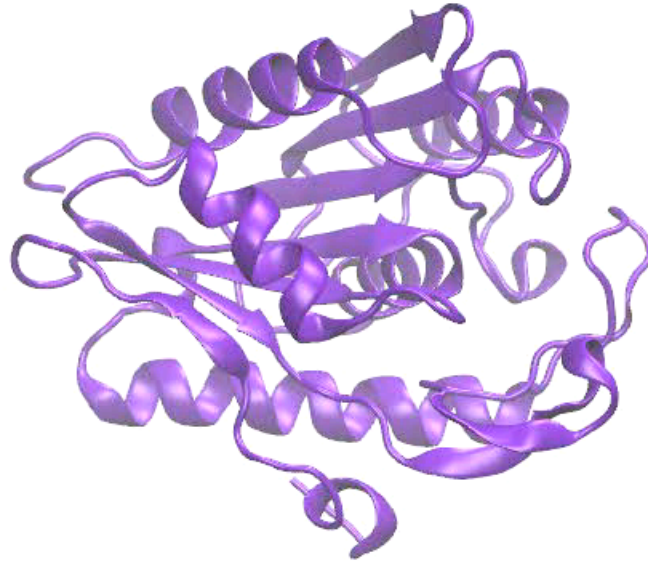
Newton's EoMs

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

molecular movies

## X-ray crystallography

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

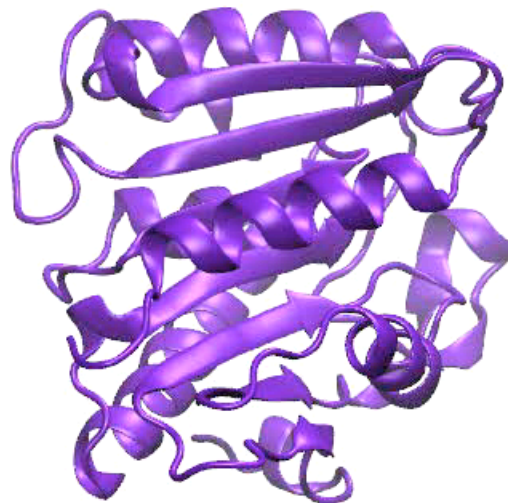


(human acyl-protein thioesterase)

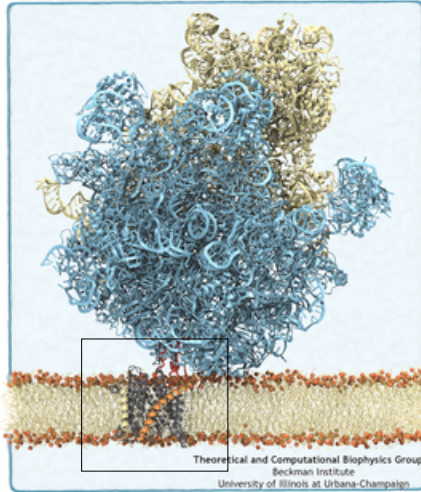
# molecular modeling and simulations

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

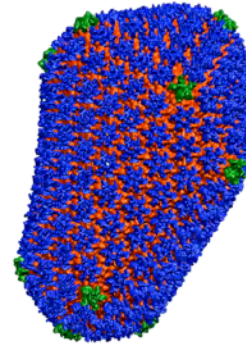
solvation  
pH  
post-translational modifications  
interactions network  
temperature effects ( $k_B T$ )  
.....



- up to  $10^2$  millions of atoms (e.g. viruses, ribosome)



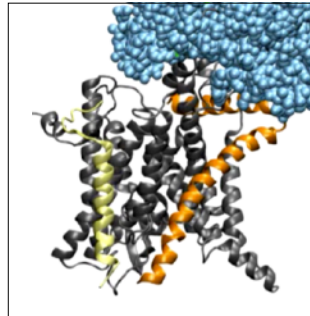
**HIV-1  
capsid**



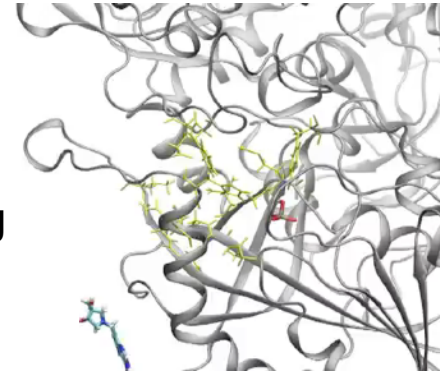
Zhao et al. *Nature*, 497:643-646, 2013  
<http://www.youtube.com/watch?v=pupVZI347H0>

James Gumbart, et al.  
*Structure*, 17:1453-1464, 2009.

**protein  
translocation**

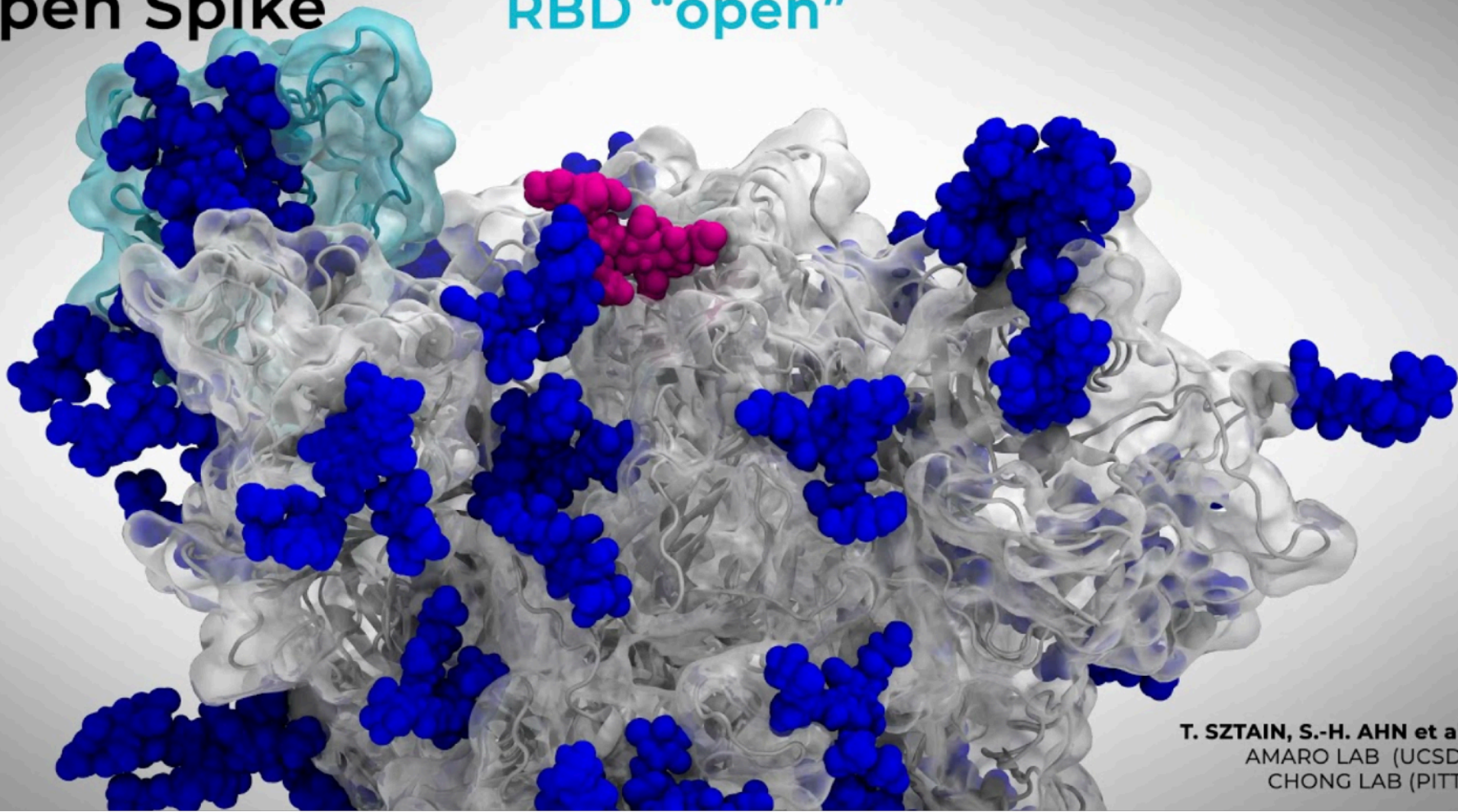


**drug binding  
on a kinase**

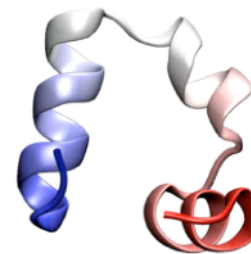
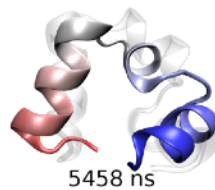
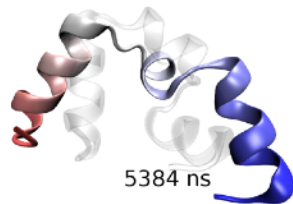
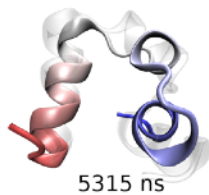


# Open Spike

RBD "open"

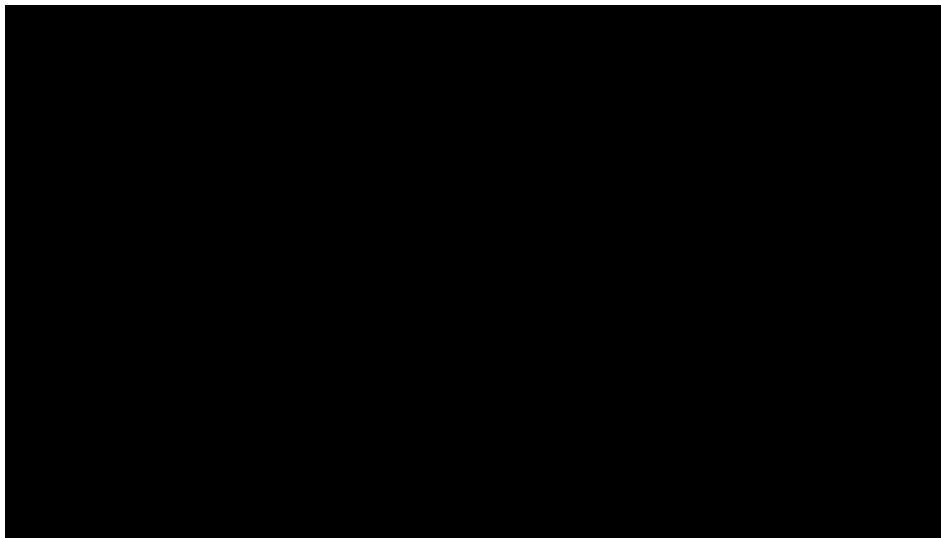


- up to the **millisecond** timescale



**villin folding**

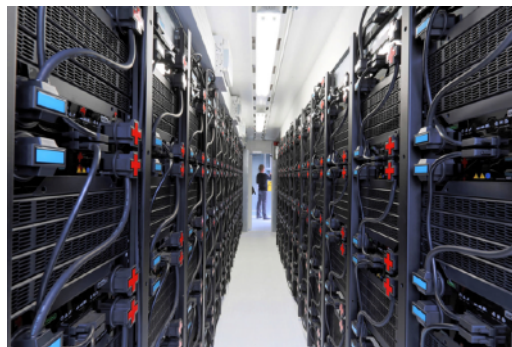
Freddolino, et al.. *Biophysical Journal*, 94:L75-L77, 2008.



Voelz et al. *J. Am. Chem. Soc.*, 2010, 132, 1526



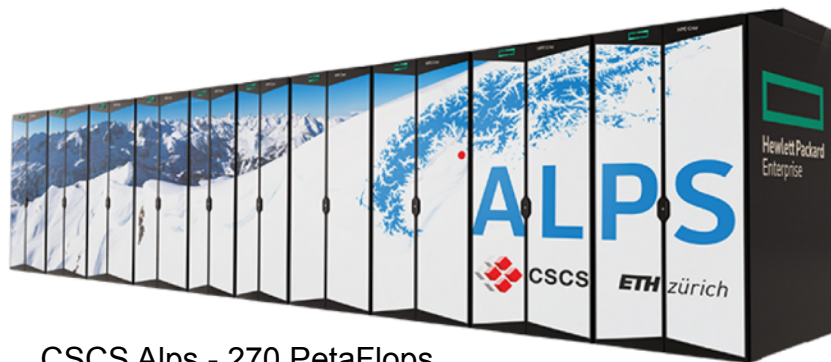
Frontier - Oak Ridge 1.3 exaFlops  
<http://www.top500.org/lists>



HPC@EPFL Kuma 12 petaFlops



Anton D.E. Shaw Research



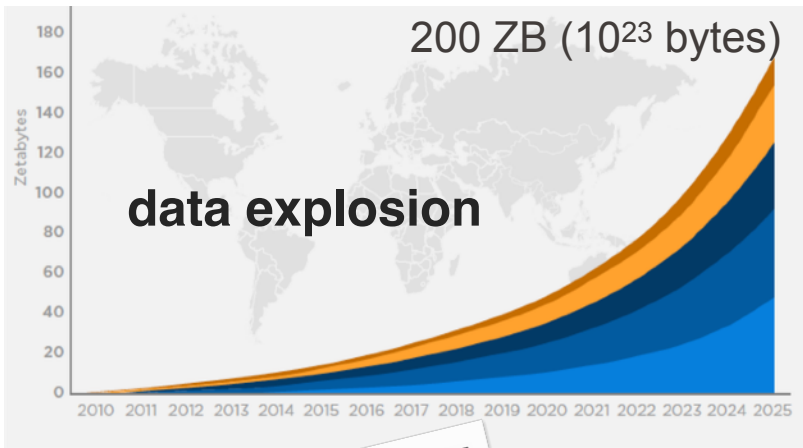
CSCS Alps - 270 PetaFlops

Molecular simulations are computational expensive and intensive because you need to integrate billions of times the equations of motion

The integration time has to be short to ensure

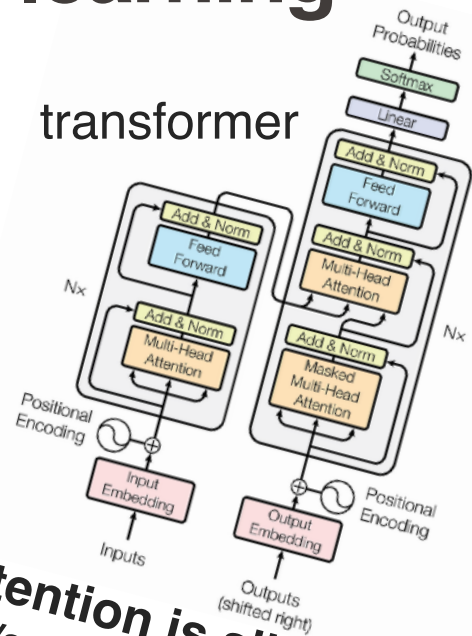
- a frame per  $\sim 1$  fs
- for 1 microsec  $\sim 10^9$  frames

# Data revolution — machine learning



powerful hardware (GPU)

transformer



“Attention is all you need”  
Vaswani et al. arXiv 2017

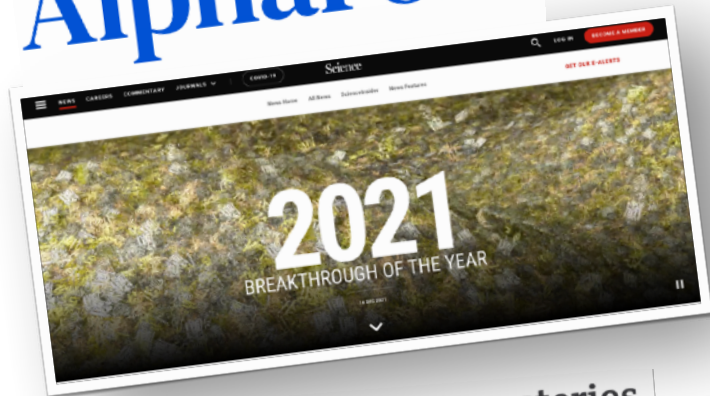


ChatGPT

# Machine learning in biology



## AlphaFold



One of biology's biggest mysteries 'largely solved' by AI

By Helen Briggs  
BBC science correspondent

**'It will change everything':  
DeepMind's AI makes gigantic leap  
in solving protein structures**

Google's deep-learning program for determining the 3D shapes of proteins stands to transform biology, say scientists.

FOCUS | 11 JANUARY 2022

### **Method of the Year 2021: Protein structure prediction**

Protein structure prediction is our Method of the Year 2021, for the remarkable levels of accuracy achieved by deep learning-based methods in predicting the 3D structures of proteins and protein complexes, essentially solving this long-standing challenge.



NEWS | BIOLOGY

### **'The game has changed!' AI triumphs at solving protein structures**

In milestone, software predictions finally match structures calculated from experimental data

**'THE ENTIRE PROTEIN UNIVERSE':  
AI PREDICTS SHAPE OF NEARLY  
EVERY KNOWN PROTEIN**

DeepMind's AlphaFold tool has determined around 200 million protein structures, which are now available to scientists in a database.

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

# CASP: Critical Assessment of Techniques for Protein Structure Prediction (now CASP16)



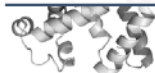
Protein Structure Prediction Center



## Menu

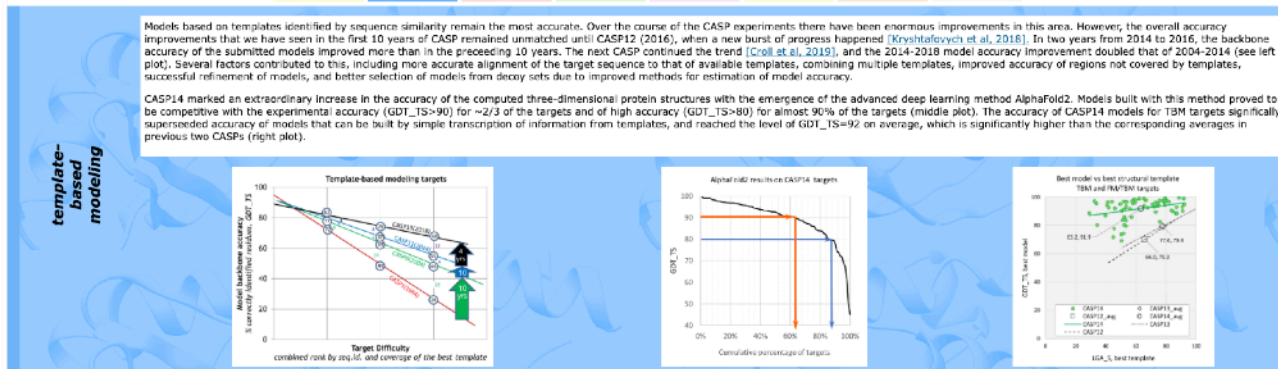
- Home
- PC Login
- PC Registration
- ▼ CASP Experiments
- CASP15 (2022)
- CASP14 (2020)
- CASP Commone (COVID-19, 2020)
- CASP13 (2018)
- CASP12 (2016)
- CASP11 (2014)
- CASP10 (2012)
- CASP9 (2010)
- CASP8 (2008)
- CASP7 (2006)
- CASP6 (2004)
- CASP5 (2002)
- CASP4 (2000)
- CASP3 (1998)
- CASP2 (1996)
- CASP1 (1994)
- Initiatives
- Data Archive
- CASP Measures
- Feedback
- Assessors
- People
- Community Resources
- Job Fair

- Initiatives
- Data Archive
- CASP Measures
- Feedback
- Assessors
- People
- Community Resources
- Job Fair



## Success Stories From Recent CASPs

- assembly
- template-based modeling
- ab initio modeling
- contact prediction
- help structural biologists
- refinement
- alpha-assisted modeling



## Welcome to the Protein Structure Prediction Center!

Our goal is to help advance the methods of identifying protein structure from sequence. The Center has been organized to provide the means of objective testing of these methods via the process of blind prediction. The Critical Assessment of protein Structure Prediction (CASP) experiments aim at establishing the current state of the art in protein structure prediction, identifying what progress has been made, and highlighting where future effort may be most productively focused.

There have been fourteen previous CASP experiments. The fifteenth experiment is planned to start in Spring 2022. Description of these experiments and the full data (targets, predictions, interactive tables with numerical evaluation results, dynamic graphs and prediction visualization tools) can be accessed following the links:

CASP1 (1994) | CASP2 (1996) | CASP3 (1998) | CASP4 (2000) | CASP5 (2002) | CASP6 (2004) | CASP7 (2006) | CASP8 (2008) | CASP9 (2010) | CASP10 (2012) | CASP11 (2014) | CASP12 (2016) | CASP13 (2018) | CASP14 (2020) | CASP15 (2022)

Raw data for the experiments held so far are archived and stored in our [data archive](#).

Details of the experiments have been published in a scientific journal *Proteins: Structure, Function and Bioinformatics*. **CASP proceedings** include papers describing the structure and conduct of the experiments, the numerical evaluation measures, reports from the assessment teams highlighting state of the art in different prediction categories, methods from some of the most successful prediction teams, and progress in various aspects of the modeling.

Prediction methods are assessed on the basis of the analysis of a large number of blind predictions of protein structure. Summary of numerical evaluation of the tertiary structure prediction methods tested in the latest CASP experiment can be found [on this web page](#). The main numerical measures used in evaluations, data handling procedures, and guidelines for navigating the data presented on this website are described in [1].

Some of the best performing methods are implemented as [fully automated servers](#) and therefore can be used by public for protein structure modeling.

To proceed to the latest CASP experiment click on the logo below:

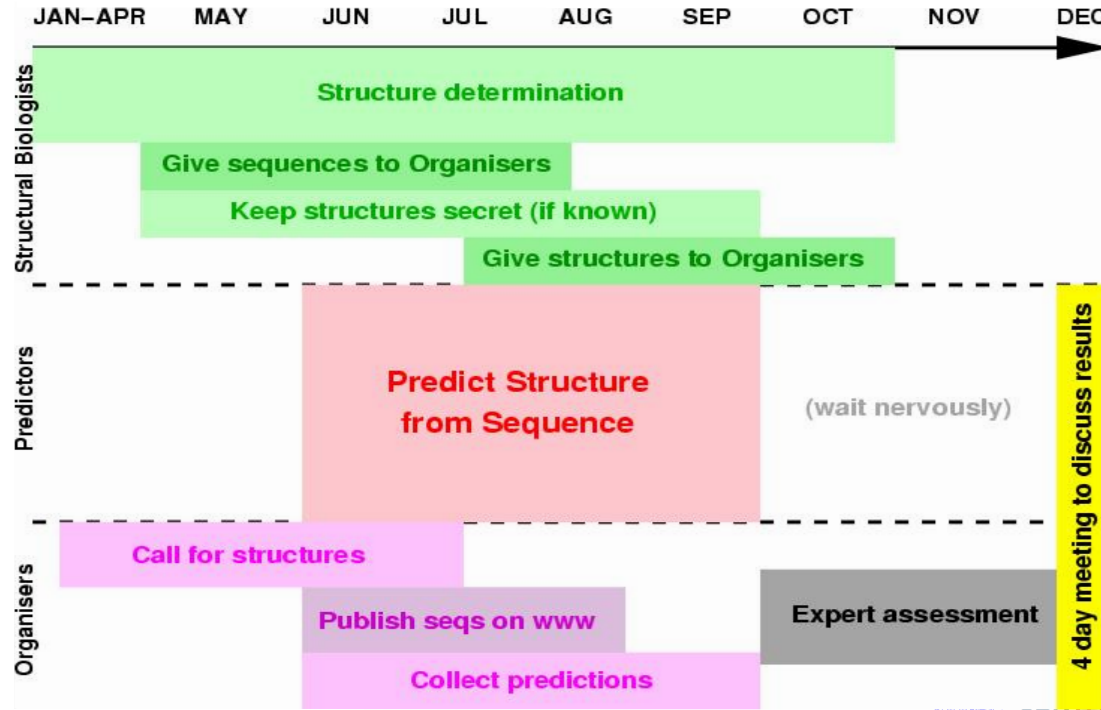
## Message Board

**CASP Special Interest Group on Ensembles of Alternative Conformations**  
 We have organized a CASP Special Interest Group on Ensembles of Alternative Conformations to provide discussions of methods for modeling alternative conformations of proteins and nucleic acids, and ...

**Reposting: Postdoctoral position at the Protein Structure Prediction Center, UC Davis**  
 We are still looking for a postdoc to help with CASP-related issues - see the posting below. --- A one-year position focused on large scale analysis of model accuracy with respect to applications ...

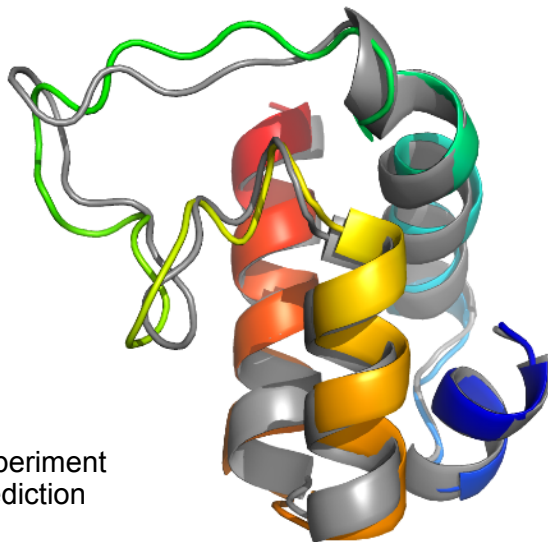
**CASP Special Interest groups**  
 One outcome of the CASP15 conference was a decision to create CASP special interest groups (SIGs) with the goal of increasing communication and discussion of new CASP-related developments. This message ...

**EPFL** **CASP: Critical Assessment of Techniques for Protein Structure Prediction (now CASP16)**



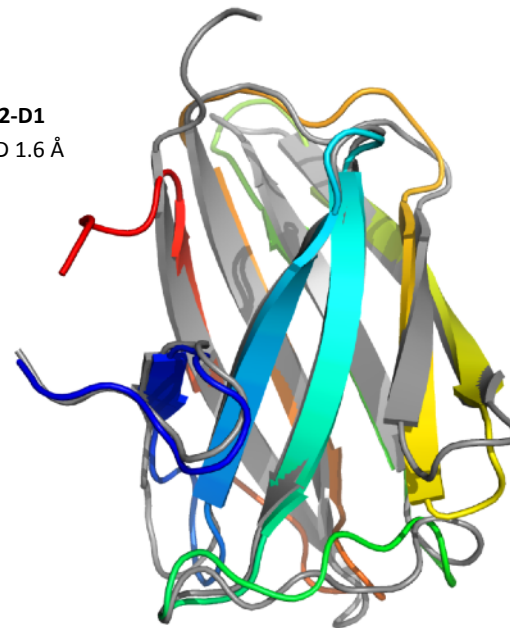
EPFL **CASP: Critical Assessment of Techniques for Protein Structure Prediction (now CASP16)**

T0990-D1  
RMSD 1.6 Å



gray = experiment  
rainbow = prediction

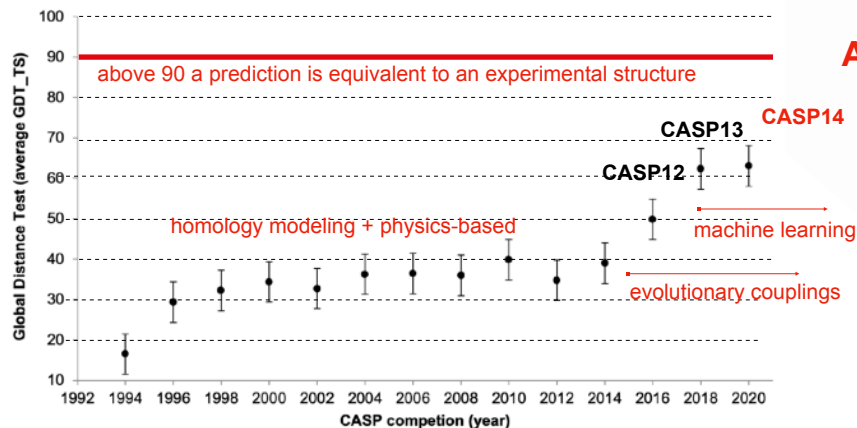
T0992-D1  
RMSD 1.6 Å



**Root Mean Square Displacement :: RMSD** defines a measure for similarity:

$$\text{RMSD} = \sqrt{\frac{1}{N} \sum_{i=1}^N d_i^2}$$

## Critical Assessment of protein Structure Prediction — CASP

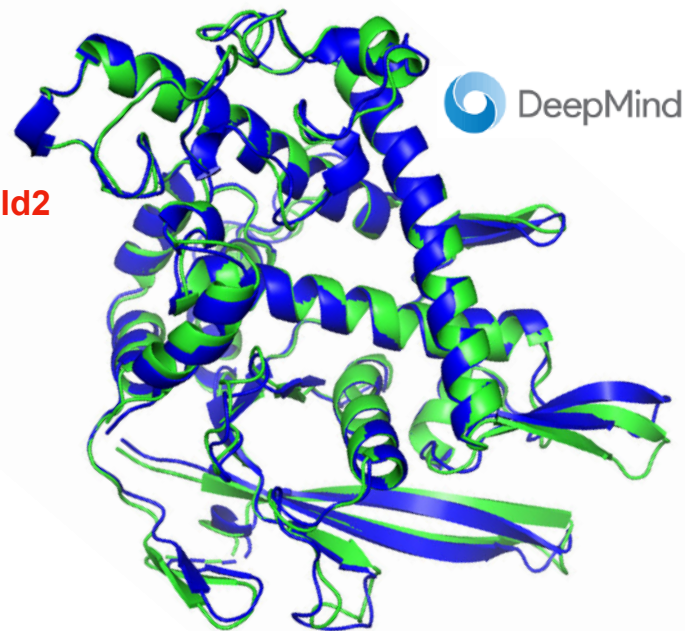


Abriata and Dal Peraro, *Proteins* 2019

Abriata, Tamo' and Dal Peraro, *Proteins* 2018

**Luciano Abriata**

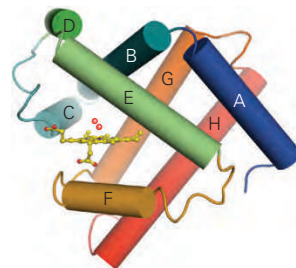
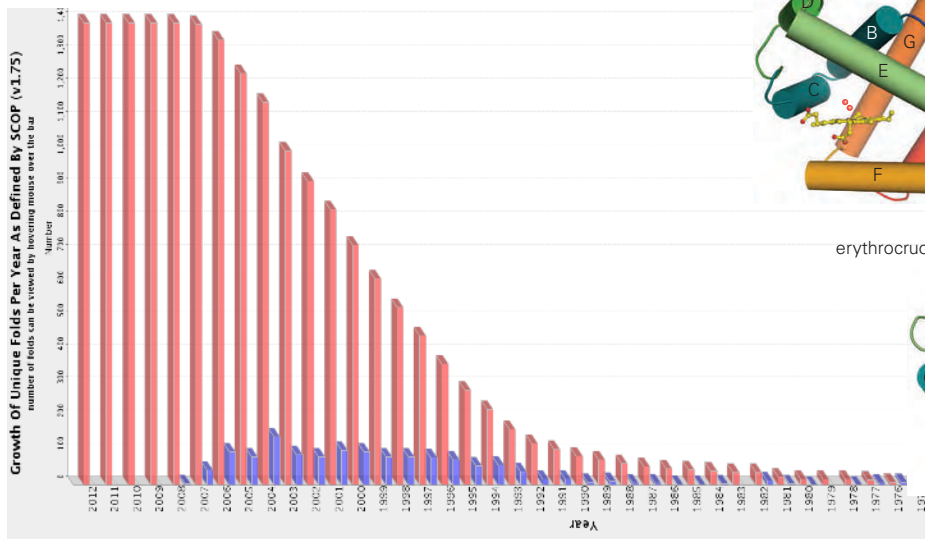
- experimental-like accuracy
- >200 M predicted models available in UniProt



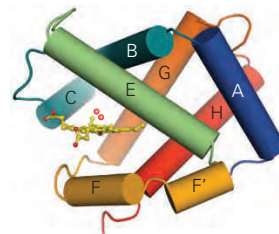
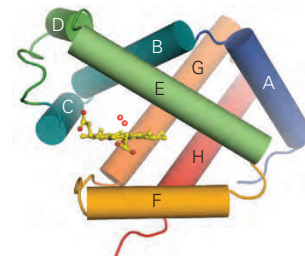
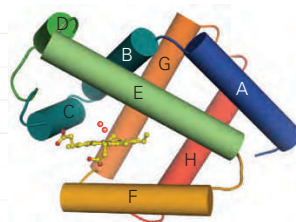
— experimental structure  
— computational prediction

# Structure prediction by homology

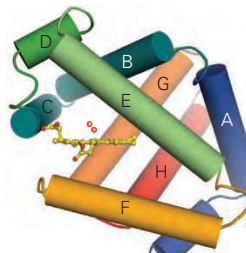
structure is more conserved than sequence



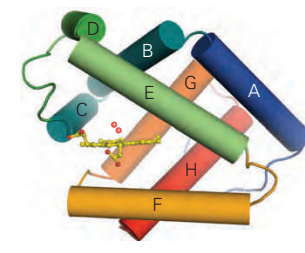
myoglobin

hemoglobin  $\alpha$  chainhemoglobin  $\beta$  chain

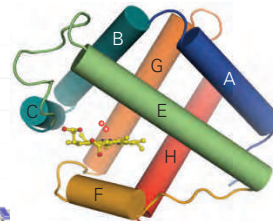
erythrocrucrin



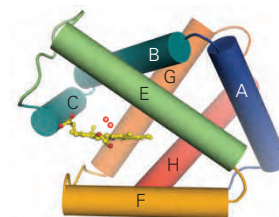
clam hemoglobin



worm hemoglobin



leghemoglobin



Glycera hemoglobin

**globin fold**



# Clustal Omega

[Input form](#) | [Web services](#) | [Help & Documentation](#) | [Bioinformatics Tools FAQ](#)

 [Feedback](#)

[Tools](#) > [Multiple Sequence Alignment](#) > [Clustal Omega](#)

## Multiple Sequence Alignment

Clustal Omega is a new multiple sequence alignment program that uses seeded guide trees and HMM profile-profile techniques to generate alignments between **three or more** sequences. For the alignment of two sequences please instead use our [pairwise sequence alignment tools](#).

**Important note:** This tool can align up to 4000 sequences or a maximum file size of 4 MB.

### STEP 1 - Enter your input sequences

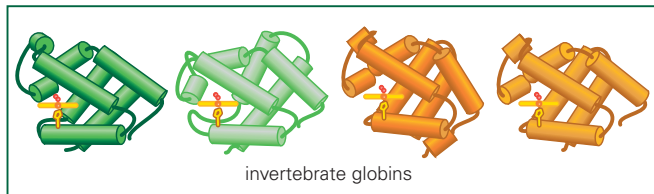
Enter or paste a set of

PROTEIN 

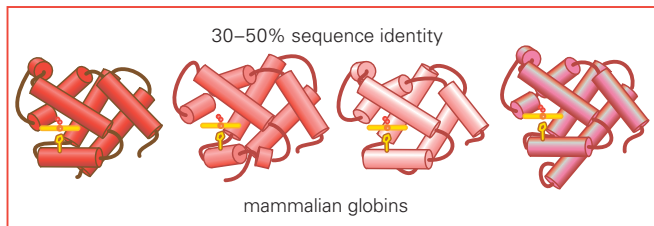
sequences in any supported [format](#):

# Structure prediction by homology

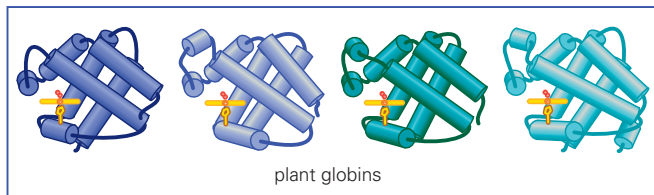
structure is more conserved  
than sequence



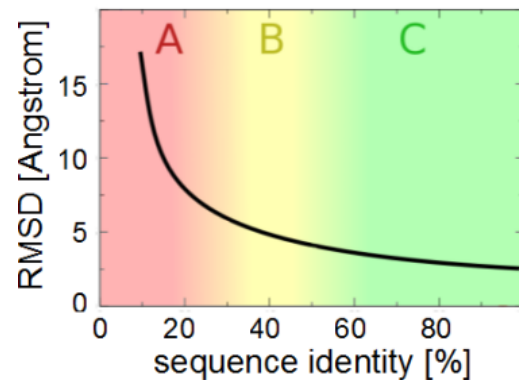
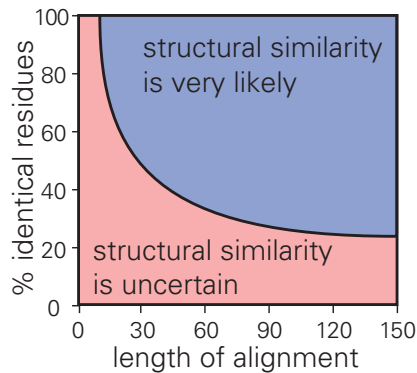
↕ 10–20% sequence identity



↕ 10–20% sequence identity



**~30% sequence identity** is required to generate an useful model by homology



# Automated homology modeling servers

<http://swissmodel.expasy.org/>

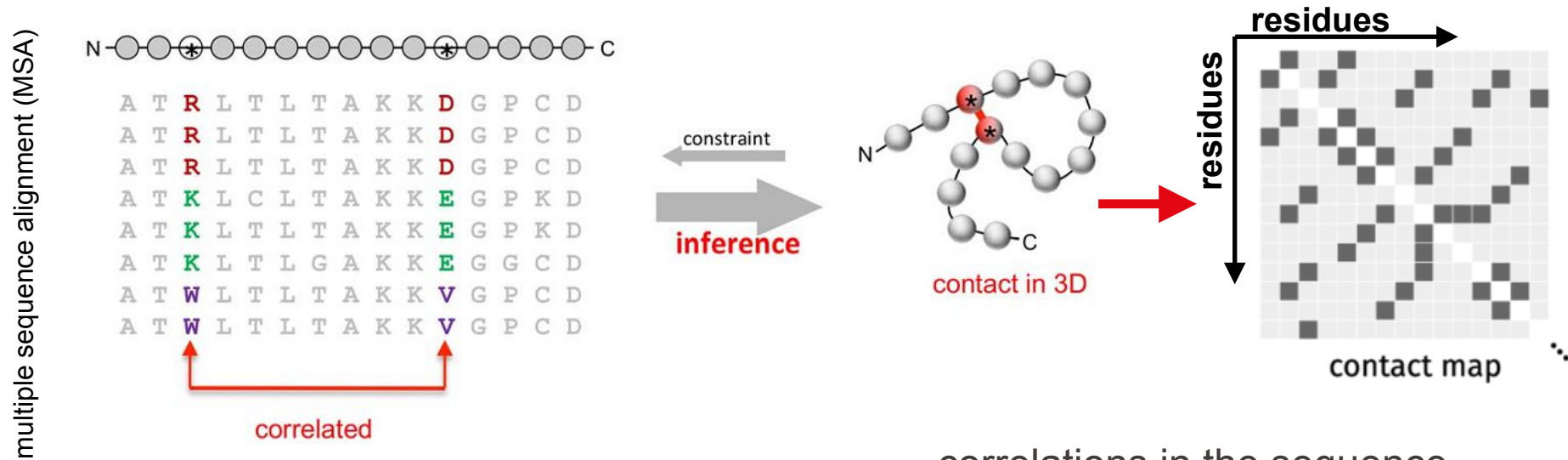
PSI | The Protein Model Portal

YLDVGFDTTRVAVIQIVLVSSA  
SDFSNDVFPFADRSGLRIR  
SVVWKRGGAVPIGIGADADTIS

[www.proteinmodelportal.org/](http://www.proteinmodelportal.org/)

<http://modbase.compbio.ucsf.edu/>

# EPFL Evolutionary couplings for protein prediction



Marks, D. S.; Colwell, L. J.; Sheridan, R.; Hopf, T. A.; Pagnani, A.; Zecchina, R.; Sander, C. *PLoS One* **2011**, *6*, e28766.

- correlations in the sequence space give structural information
- if you have enough predicted contacts you can fold a protein (similar to NMR)

## Direct-coupling analysis (DCA)

Calculate covariance matrix for each pair of sequence positions for all pairs of amino acids (A,B)

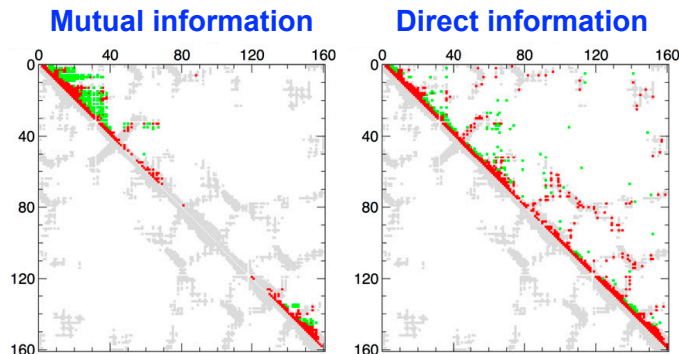
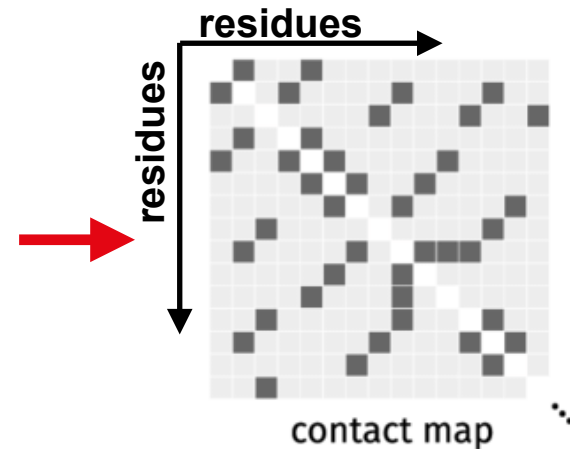
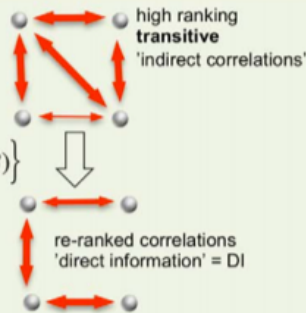
$$C_{ij}(A,B) = f_{ij}(A,B) - f_i(A)f_j(B)$$

$$C_{ij}^{-1}(A,B) = -e_{ij}(A,B)_{i \neq j}$$

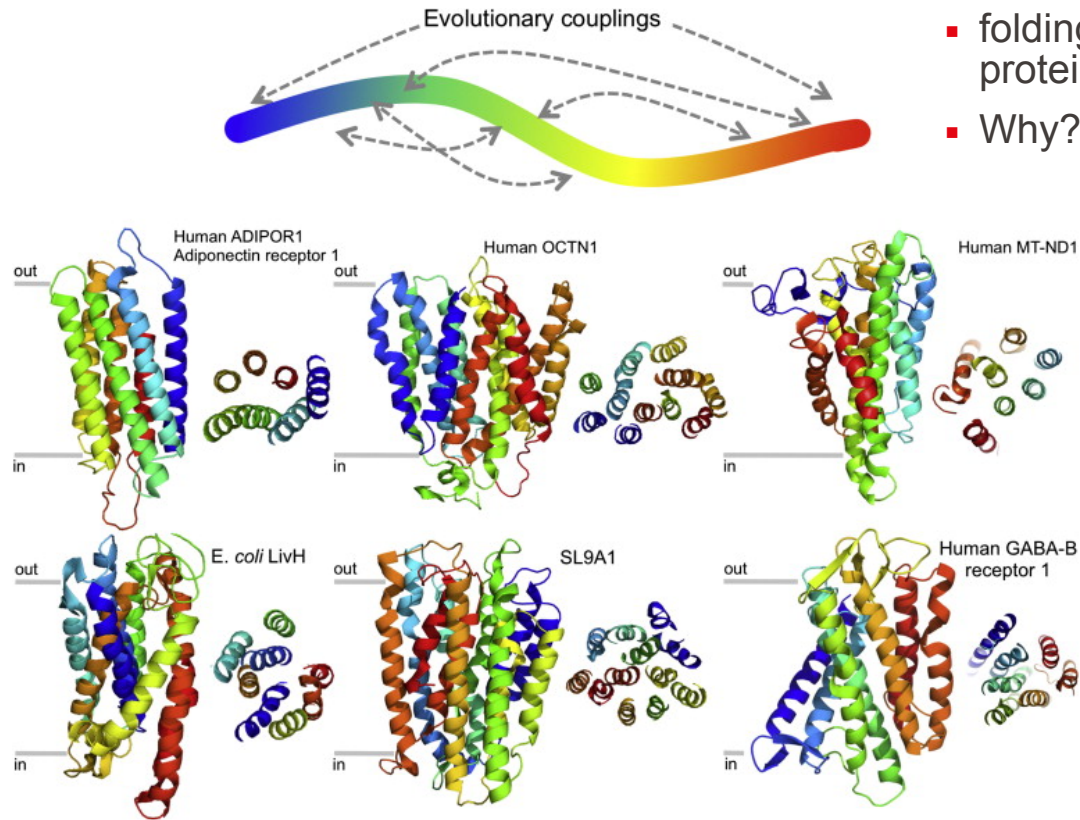
Identify maximally informative pair couplings using **statistical model** of entire protein to infer residue-residue co-evolution

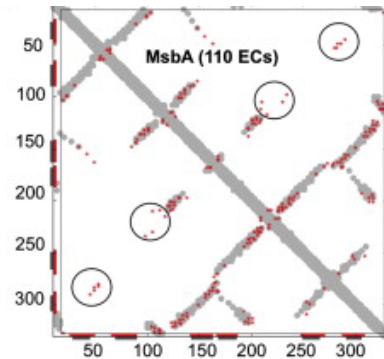
$$P_{ij}^{Dir}(A,B) = \frac{1}{Z} \exp\{e_{ij}(A,B) + \tilde{h}_i(A) + \tilde{h}_j(B)\}$$

$$DI_{ij} = \sum_{A,B=1}^q P_{ij}^{Dir}(A,B) \ln \frac{P_{ij}^{Dir}(A,B)}{f_i(A)f_j(B)}$$



- a statistical physics method were used to crack the problem (ie Potts model)

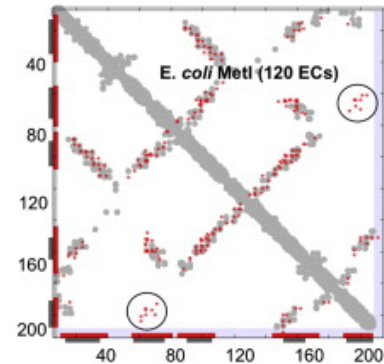
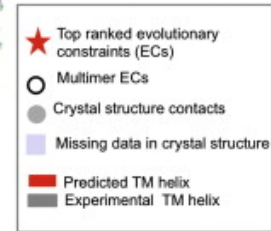
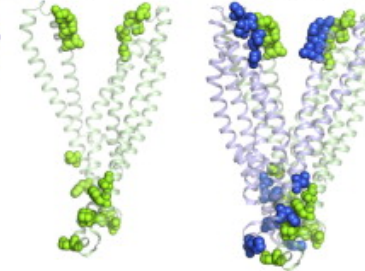




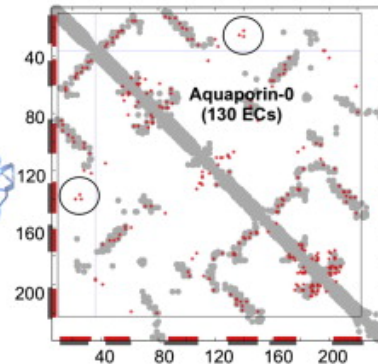
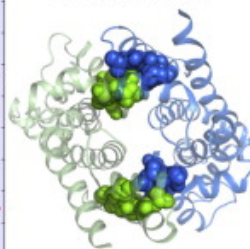
MsbA monomers



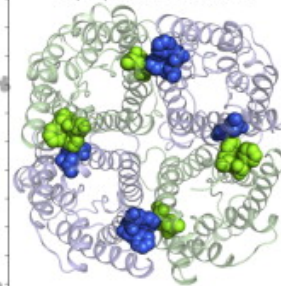
MsbA dimer



E. coli MetI dimer



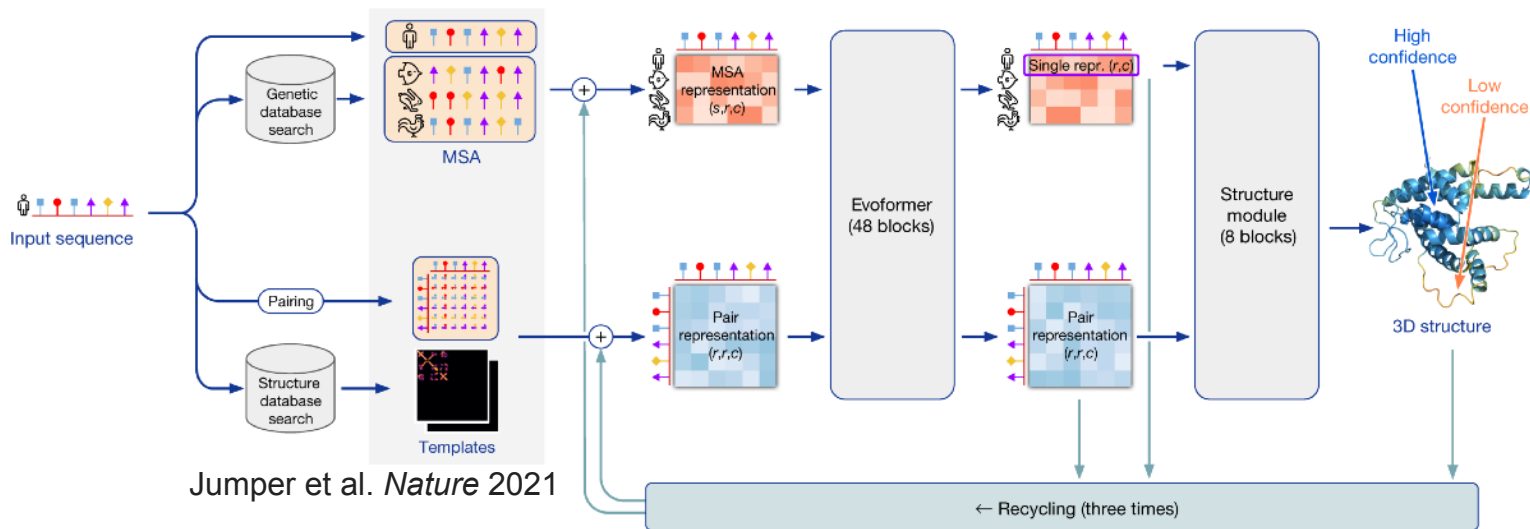
Aquaporin-0 tetramer



- ECs are found not only within the same protein but also among protein interfaces
- this can happen for homo or hetero multimers

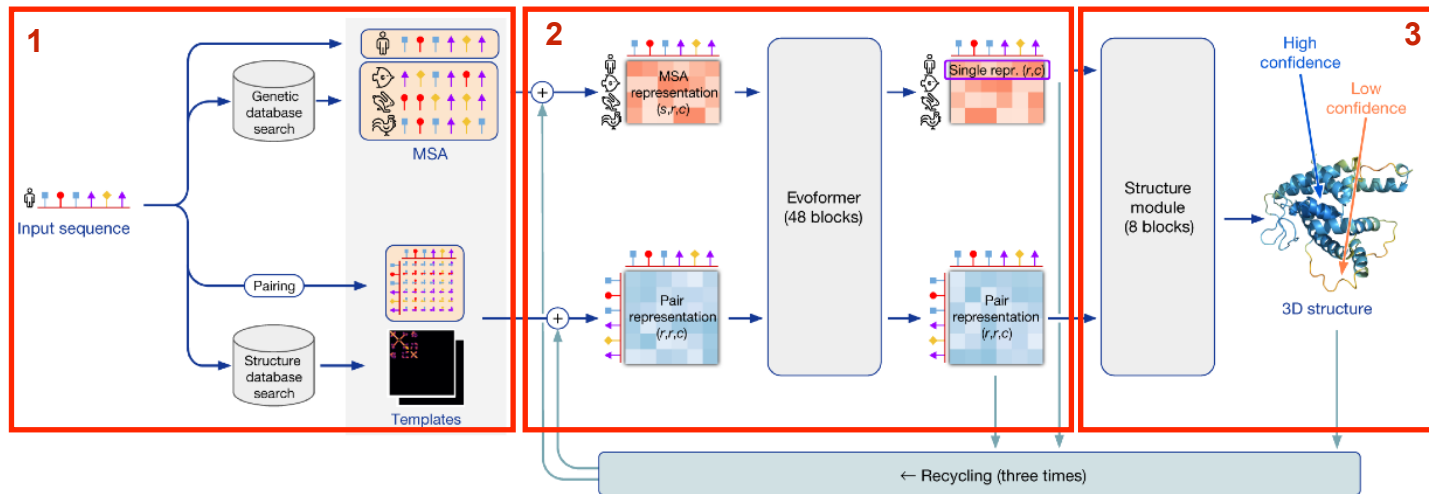
# AlphaFold2 architecture

- “AlphaFold greatly improves the accuracy of structure prediction by incorporating novel neural network architectures and training procedures based on the evolutionary, physical and geometric constraints of protein structures”.



- trained on sequence similarity and structural templates from databases (UniProt/metagenomics and PDB)
- end-to-end model produces prediction in one shot using *transformers*

# AlphaFold2 architecture



## 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)

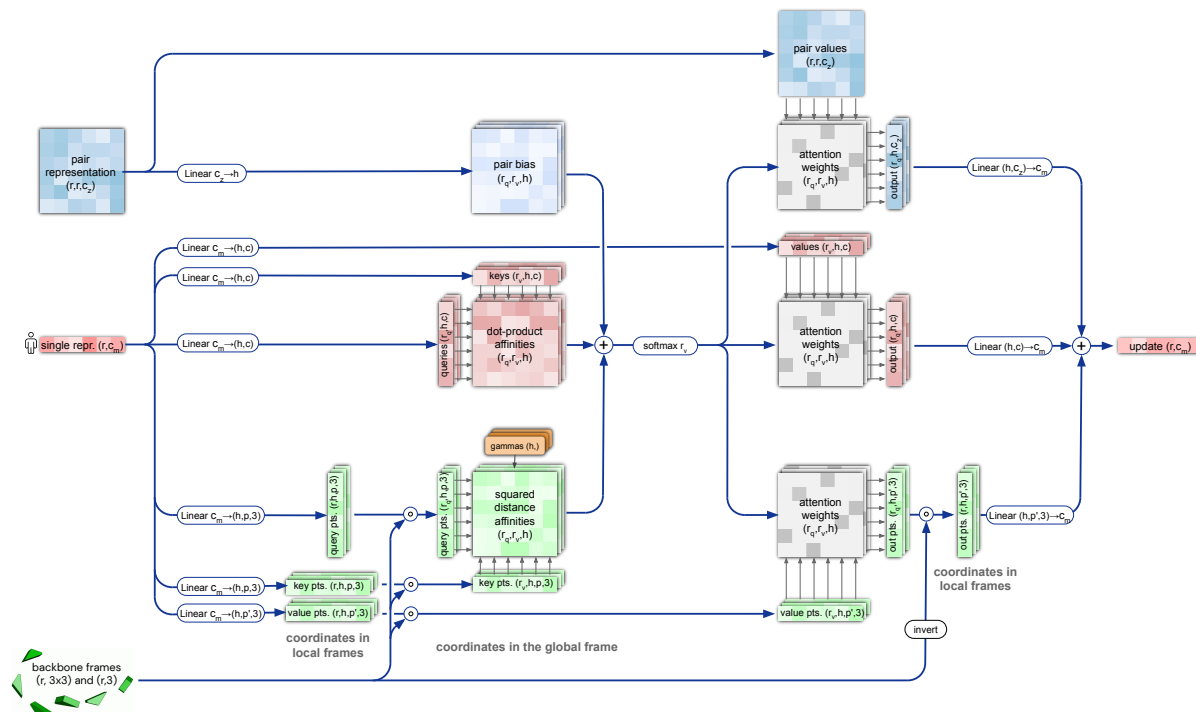
## Second module:

Evoformer transformer which refine the MSA and pair interactions

## Third module:

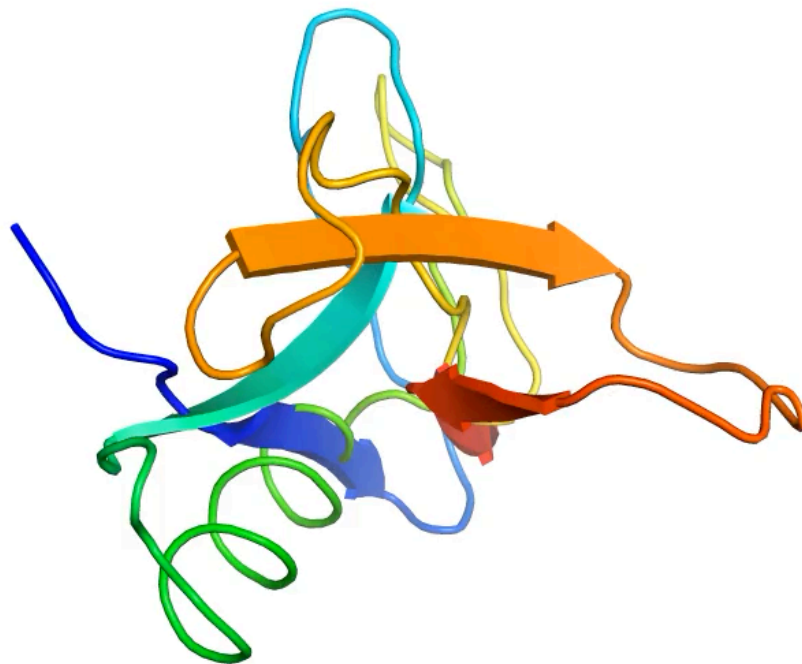
The structure module build the 3D structure based on the MSA and pair interactions information

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



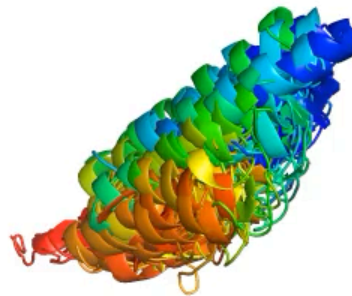
**Supplementary Figure 8** | Invariant Point Attention Module. **(top, blue arrays)** modulation by the pair representation. **(middle, red arrays)** standard attention on abstract features. **(bottom, green arrays)** Invariant point attention. Dimensions: r: residues, c: channels, h: heads, p: points.

# AlphaFold2 at work



Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction

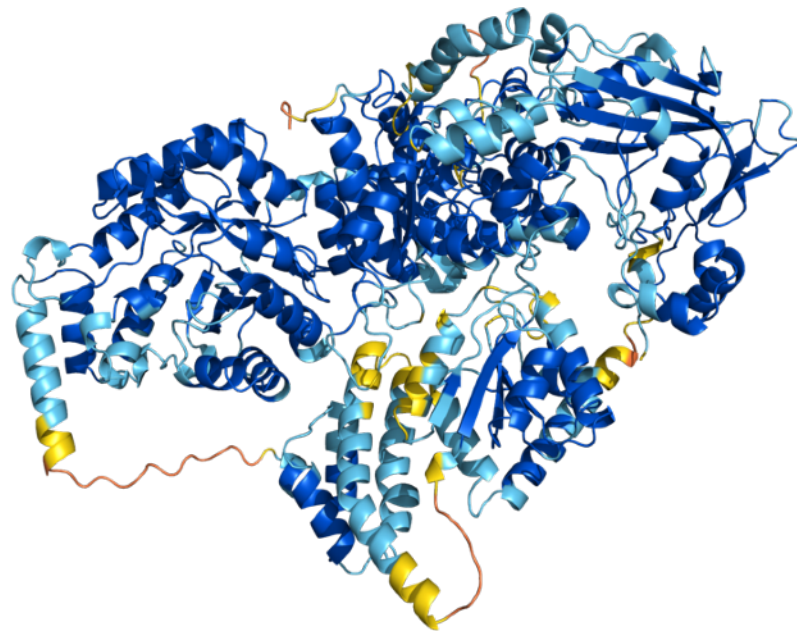
# AlphaFold2 at work



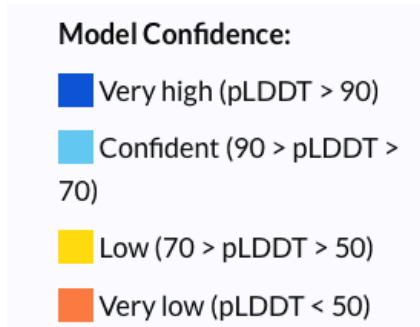
Recycling iteration 0, block 01  
Secondary structure assigned from the final prediction

# AlphaFold2 database

- > 200 million protein structure predictions
- Almost all catalogued proteins (UniProt)
- Over 1 million organisms
- Freely and openly available
- Collaboration DeepMind and EMBL-EBI
- 35.2% predictions with mean pLDDT > 90
- 79.1% predictions with mean pLDDT > 70
  
- <https://www.alphafold.ebi.ac.uk/>
- <https://uniprot.org/>



# AlphaFold2 is self-assessing



- **pLDDT**: predicted local distance difference test score
- prediction of the local distances between pairs of residues in the predicted structure compared to a reference or ground truth structure.
- low score (<50) indicates that the region is disordered or AF2 does not have enough information

Search for protein, gene, UniProt accession or organism

BETA

Search

Examples: Free fatty acid receptor 2 At1g58602 Q5VSL9 E. coli

See search help

## PAE: predicted alignment error

### 3D viewer

Sequence of AF-P35247-F1 Chain 1: Pulmonary su

MLLFLLSALVLLTQPLGYLEAEMKTYSHRTNPSACTLVMCSSVESGLPGRDGRDREGPRGEK  
 GDPGLPGAAGQAGMPGOAGPVGPKGDNGSVGEPGPKGDYGPSPGPPGVPVPGPAGREGPLGKQ  
 GNIGPQGPVKGEAGPKGEVGAPEGMQGSAGARLAGPKGERGVPGERGVPGNTGAAGSAGAM

AF-P35247-F1

Type Model

Nothing Focused

Quick Styles

Default Stylized Illustrative

Components AF-P35247-F1

Preset + Add

Polymer Cartoon

Measurements

+ Add

Export Animation

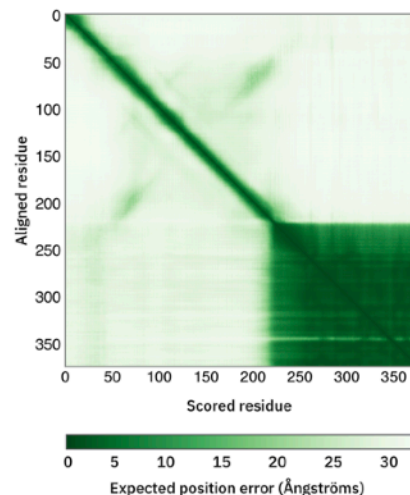
Export Geometry

### Model Confidence

- Very high (pLDDT > 90)
- High (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

AlphaFold produces a per-residue model confidence score (pLDDT) between 0 and 100. Some regions below 50 pLDDT may be unstructured in isolation.

### Predicted aligned error (PAE)

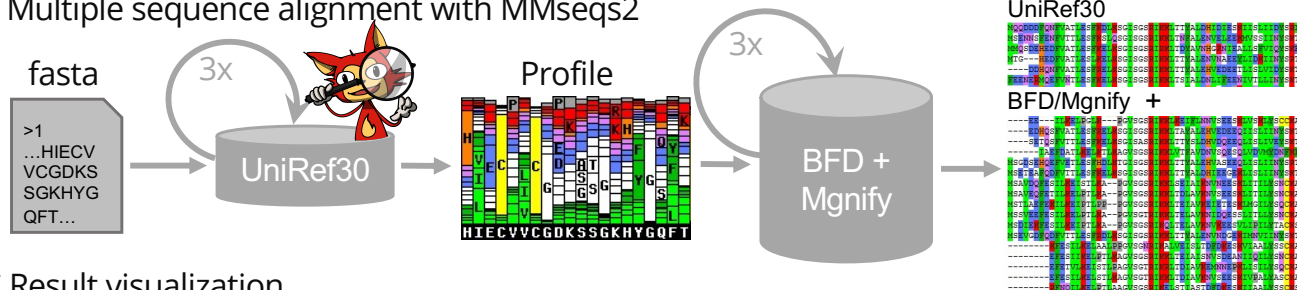


Click and drag a box on the PAE viewer to select regions of the structure and highlight them on the 3D viewer.

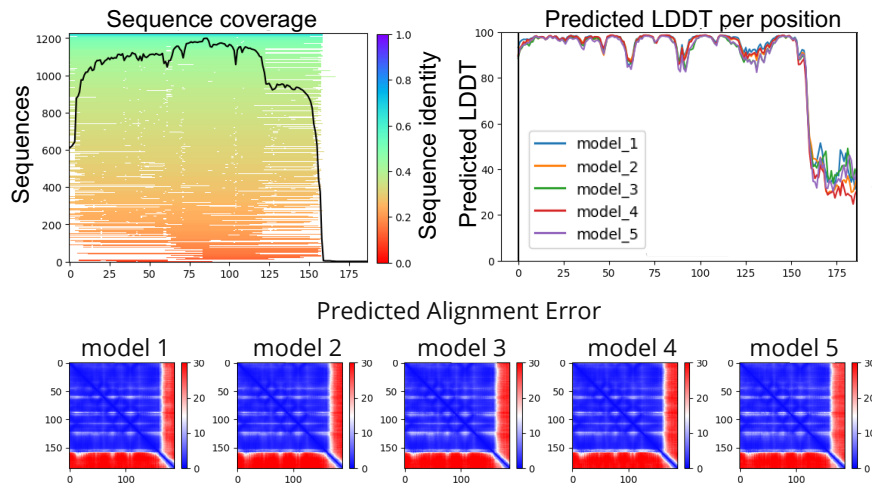
PAE data is useful for assessing inter-domain accuracy – [go to Help section below](#) for more information.

- 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

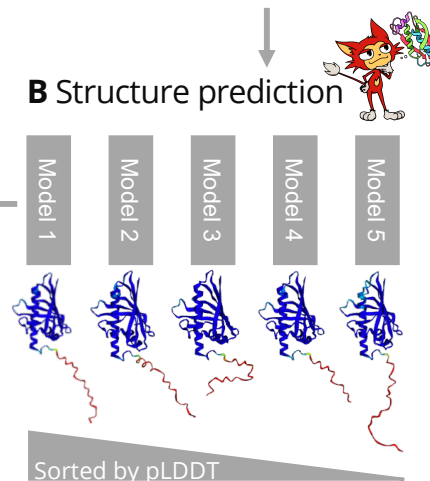
## A Multiple sequence alignment with MMseqs2



## C Result visualization



## B Structure prediction

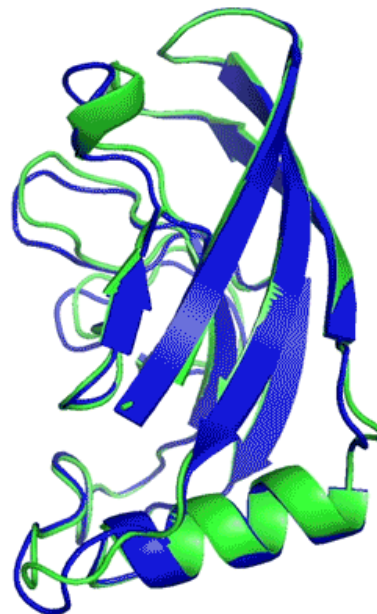


# AlphaFold2 is a new SB methods

- Experimental result
- Computational prediction



T1037 / 6vr4  
90.7 GDT  
(RNA polymerase domain)



T1049 / 6y4f  
93.3 GDT  
(adhesin tip)

~200K (PDB)

>200M (UniProt)

# EPFL AlphaFold2 is not the ultimate oracle

## ▪ Limitations

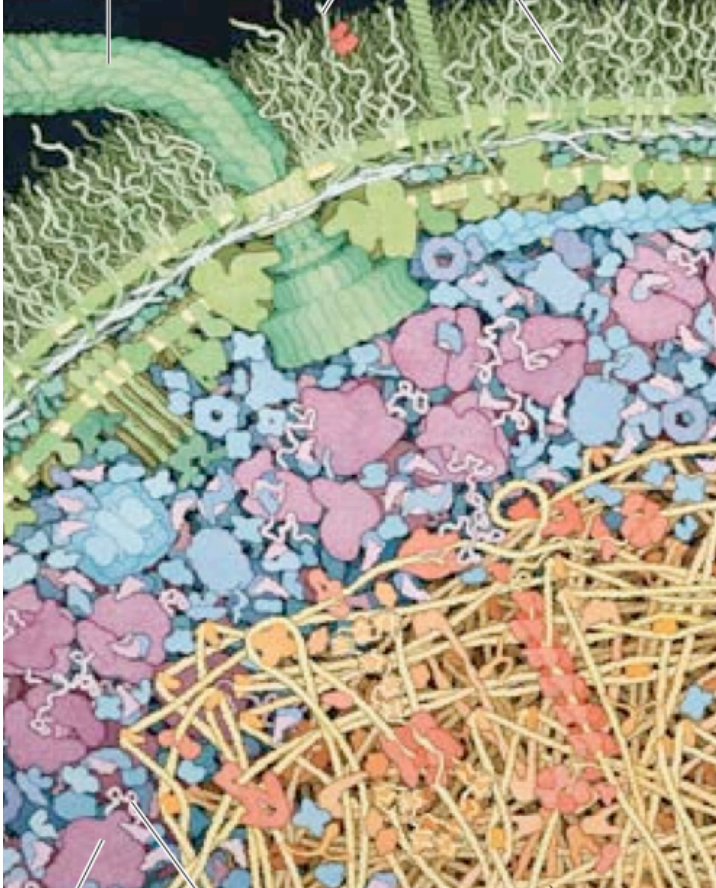
- need for a deep MSA (>30 sequences) to create accurate models
- not all the models are highly accurate as an experimental structure
- it does not account for dynamics and multiple states
- does not account for the post-translational modifications

## ▪ Other potential benefits

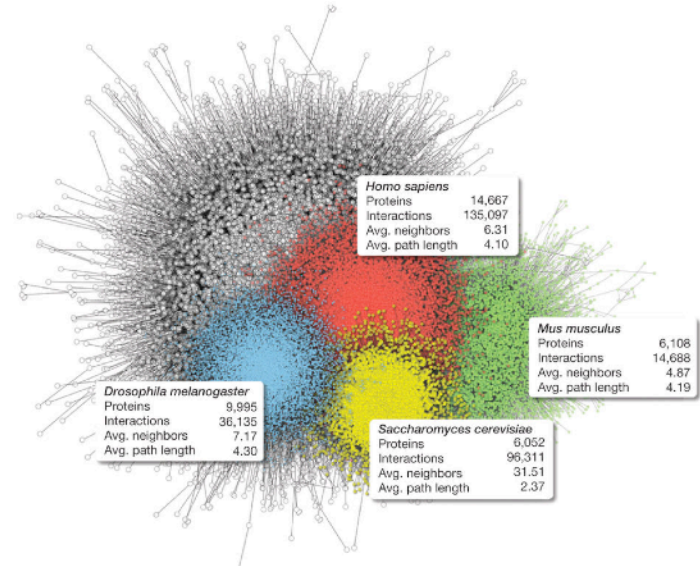
- AF triggered many other developments
- can assist experimental structure determination
- eg, in molecular replacement in X-ray crystallography
- eg, in cryoEM fitting and model building
- it is a means to look at protein-protein networks

▪

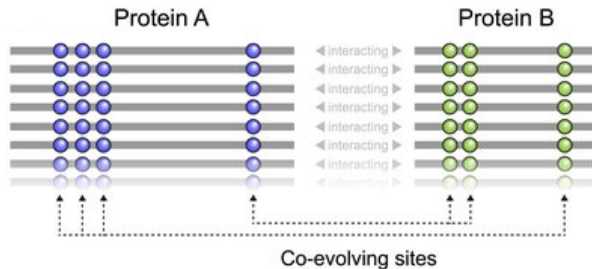
# Proteins form assemblies and networks



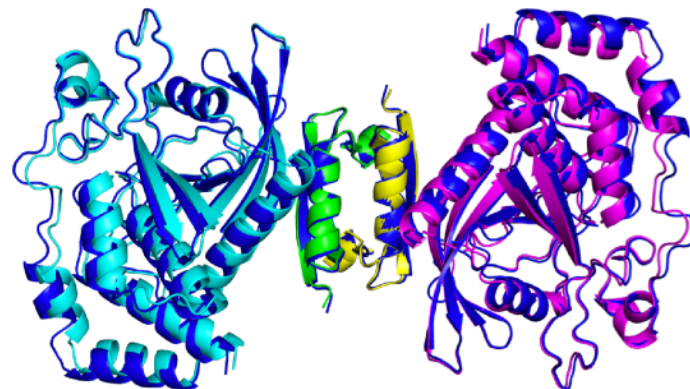
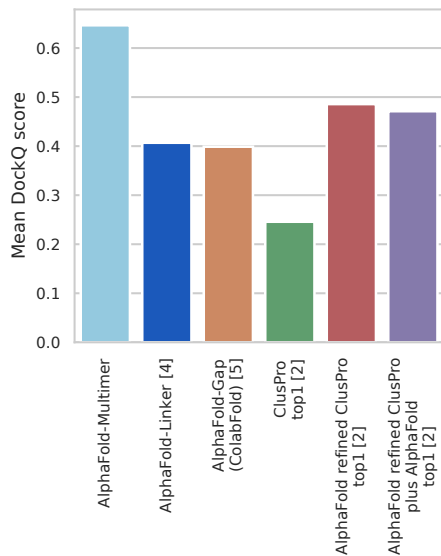
- Proteins can then arrange into supramolecular assemblies
- Interacting with nucleic acid, metabolites, membranes, etc.
- They create large network of interactions



coevolution coupling holds for protein-protein interactions



- + Support for multiple chains
- + Multi-chain features
- + Various architectural modifications
- + Paired MSAs
- + Training on complexes



Protein complex prediction with AlphaFold-Multimer  
 Richard Evans, ..., John Jumper, Demis Hassabis  
 bioRxiv 2021.10.04.463034; doi: <https://doi.org/10.1101/2021.10.04.463034>

Article

<https://doi.org/10.1038/s41594-022-00910-8>

# Towards a structurally resolved human protein interaction network

Received: 11 February 2022

Accepted: 14 December 2022

Published online: 23 January 2023

Check for updates

David F. Burke<sup>1,9</sup>, Patrick Bryant<sup>2,3,9</sup>, Inigo Barrio-Hernandez<sup>1,9</sup>, Danish Memon<sup>1,9</sup>, Gabriele Pozzati<sup>2,3,9</sup>, Aditi Shenoy<sup>2,3</sup>, Wensi Zhu<sup>2,3</sup>, Alistair S. Dunham<sup>1</sup>, Pascal Albanese<sup>4,5</sup>, Andrew Keller<sup>6</sup>, Richard A. Scheltema<sup>4,5</sup>, James E. Bruce<sup>6</sup>, Alexander Leitner<sup>7</sup>, Petras Kundrotas<sup>2,3,8</sup>, Pedro Beltrao<sup>1,7</sup> & Arne Elofsson<sup>2,3</sup>✉

Article

<https://doi.org/10.1038/s41594-024-01791-x>

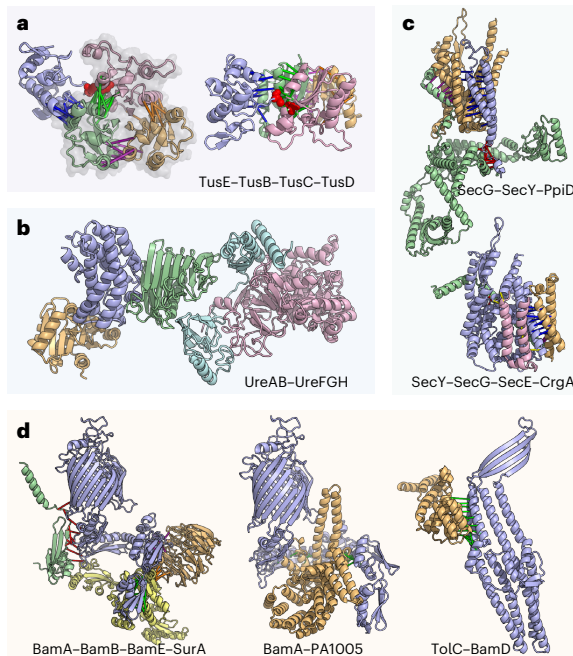
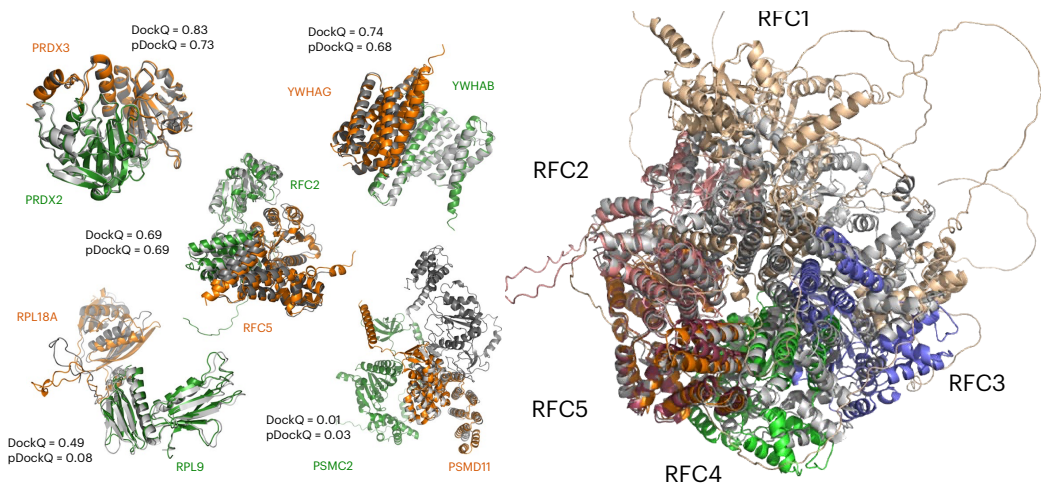
# Protein interactions in human pathogens revealed through deep learning

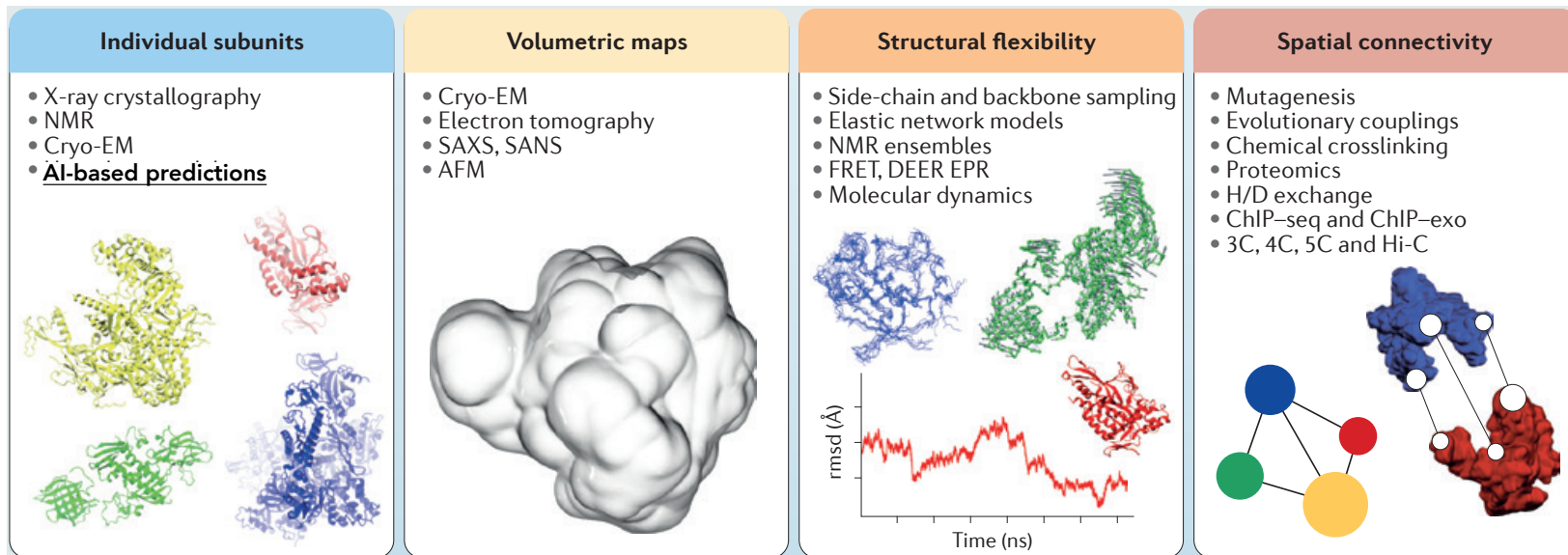
Received: 28 April 2023

Accepted: 23 July 2024

Published online: 18 September 2024

Ian R. Humphreys<sup>1,2,13</sup>, Jing Zhang<sup>3,4,5,13</sup>, Minkyung Baek<sup>6,13</sup>✉, Yaxi Wang<sup>7,13</sup>, Aditya Krishnakumar<sup>1,2</sup>, Jimin Pei<sup>3,4,5</sup>, Ivan Anishchenko<sup>1,2</sup>, Catherine A. Tower<sup>7</sup>, Blake A. Jackson<sup>7</sup>, Thulasi Warriar<sup>8,9,10</sup>, Deborah T. Hung<sup>8,9,10</sup>, S. Brook Peterson<sup>7</sup>, Joseph D. Mougous<sup>7,11,12</sup>, Qian Cong<sup>3,4,5</sup>✉ & David Baker<sup>1,2,11</sup>✉

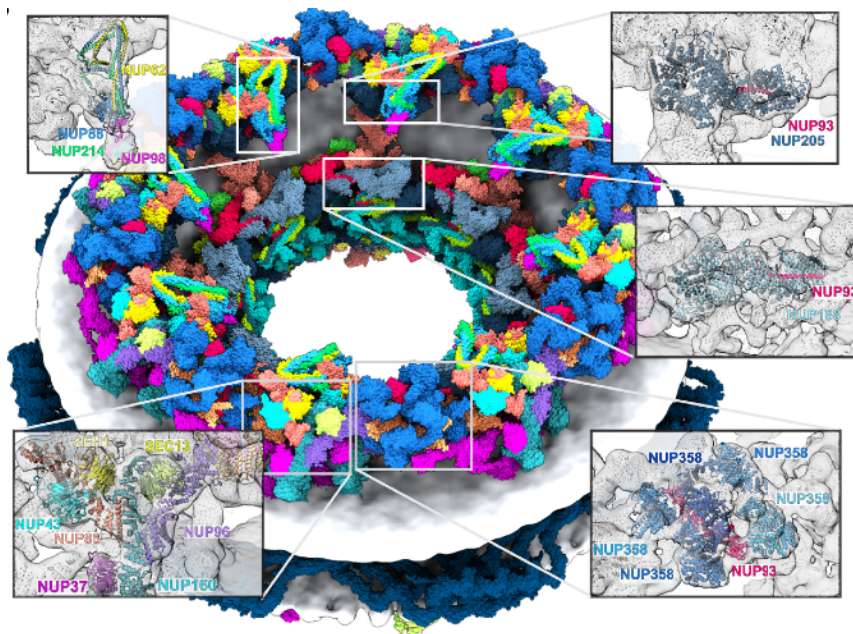
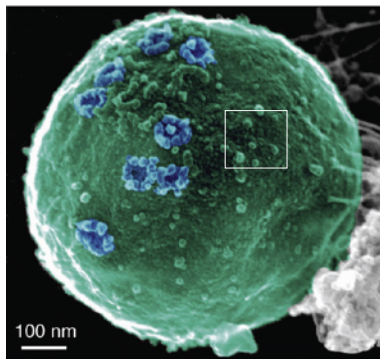




**data integration and  
model building**

# AF2 affects integrative structural biology

e.g., the Nuclear Pore Complex



- yeast: ~52 MDa, ~550 proteins
- human: ~120 MDa, ~100 proteins

# AlphaFold Server

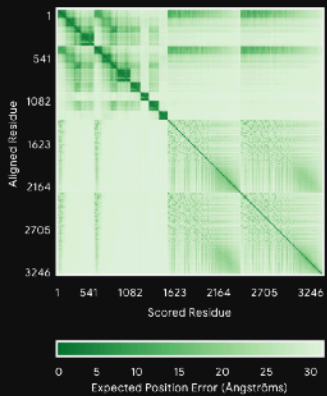
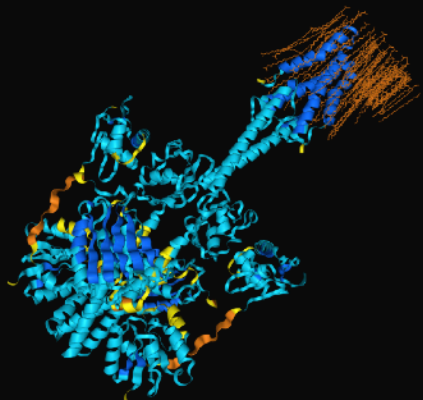
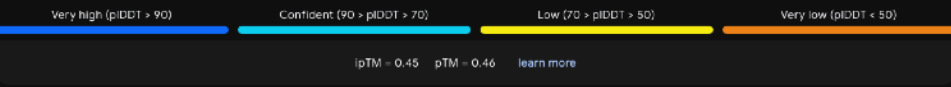
Powered by AlphaFold 3

Abramson, J et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. Nature (2024)

AlphaFold Server BETA Server About FAQs

## PhoR-PhoB

← Back Download Clone and reuse Feedback on structure



Type	Copies	Sequence
Protein	2	<pre> 10      20      30      40 MLERLSWKRL VLELLLCCLP AFILGAFFGY LPWFLLASV 70      80      90 RSMTPPPGRG SWEPLLYGLH QQLRNKKRR RELGNLIKR 130     140     150 FWCNGLAQOI LGLRWPEDNG QNILNLLRYP EFTQYLKTR 190     200     210     220 PYTHKQLLMV ARDVTQMHQL EGARRNFFAN VSELRTPL 250     260     270 KALHTMREQT QRMEGLVKQL LTLSKIEAAP THLLNEKVD 310     320     330 FTFEIDNGLK VSGNEDQLRS AISNLVYNAP NHTPEGTHI 370     380     390 IAPEHIPRLT ERFYRVDKAR SRQTGGSGLG LAIVKHAVN 430     431 IPERLIAKNS D </pre>
Ligand	50	PLM – Palmitic acid
Ligand	50	OLA – Oleic acid
Protein	2	<pre> 10      20      30      40 MARRILVVEE EAPIREMVCF VLEQNGFQPV EAEDYDSAV 70      80      90 GIOFIKHLKR ESMTRDIPVV MLTARGEED RVRGLETGA 130     140     150 RRISPNAVEE VIEMQGLSLD PTSHRVMAGE EPLEMGPTT 190     200     210     220 NHVWGTVVYV EDRTVDVHIR RLRKALEPGG HDRMVQTVR </pre>

PTMs: 213H: ND1-Phosphochistidine

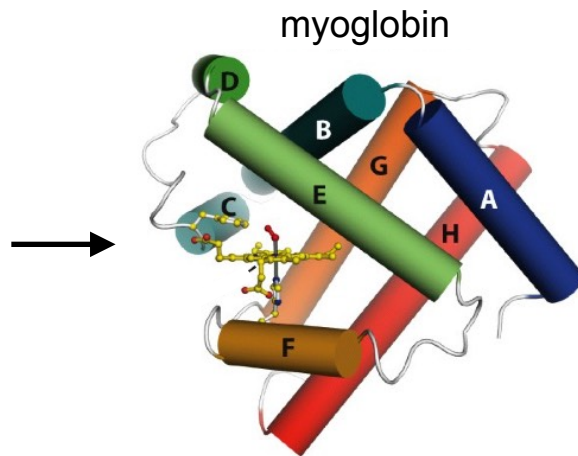
Seed: 1026411006

# The folding paradigm

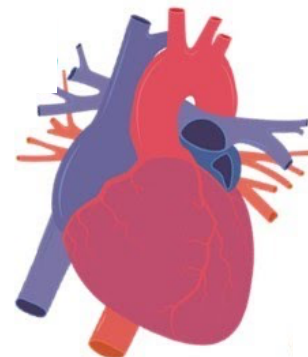
```

MGLSDGEWQLVLNVWG
KVEADIPGHGQEVLR
LFKGHPELLEKFDKFK
HLKSEDEMKALEDLKK
HGATVLTALGGILKKK
GHEAEIKPLAQSHAT
KHKIPVKYLEFISECI
IQVLQSKHPGDFGADA
QGAMNKALELFRKDMA
SNYKELGFQGG
  
```

sequence



structure

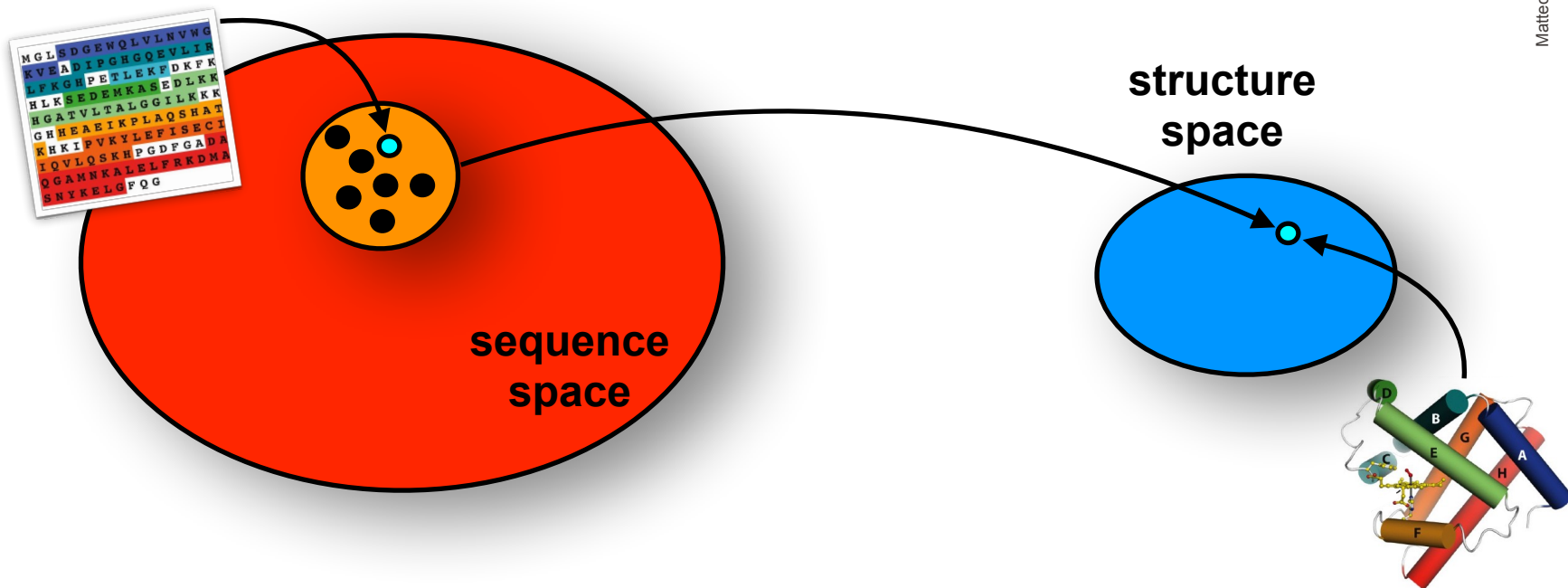


function

evolution (billion year)

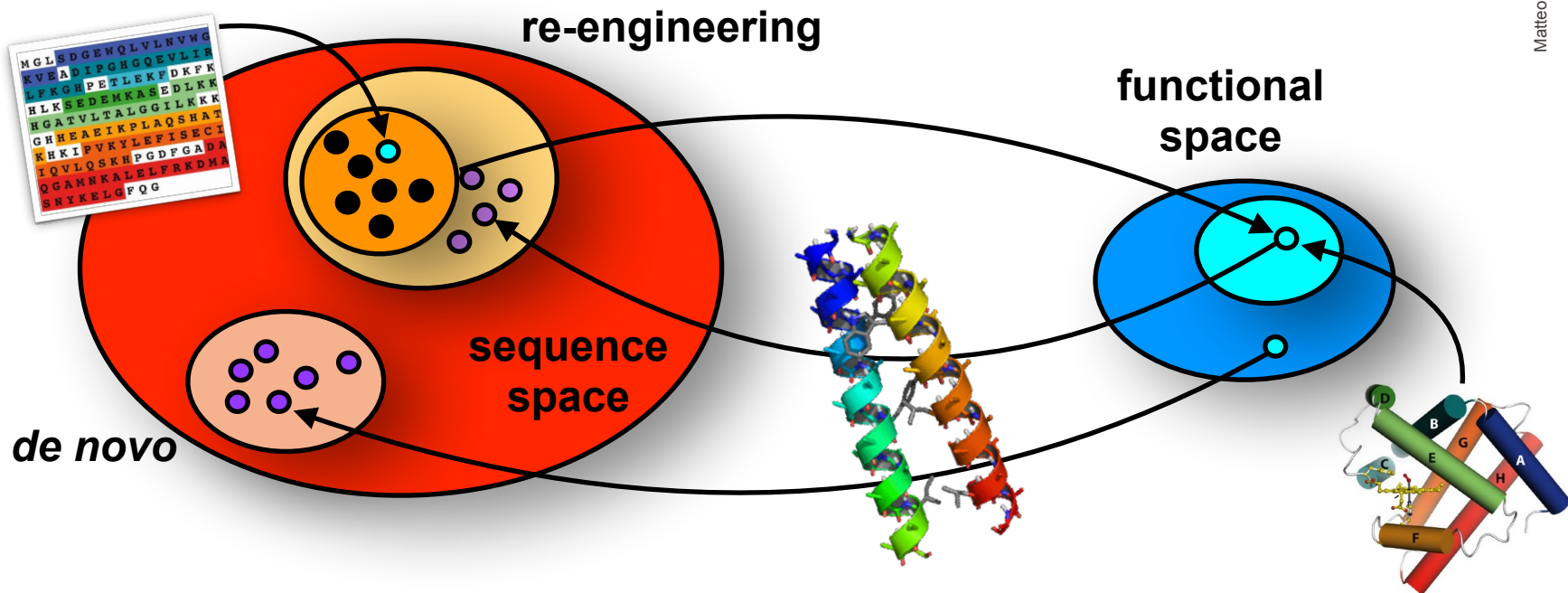
- Prediction of final structure and binding helps discovering new biology
- Not all the questions are answered though by AF2 !!

# 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

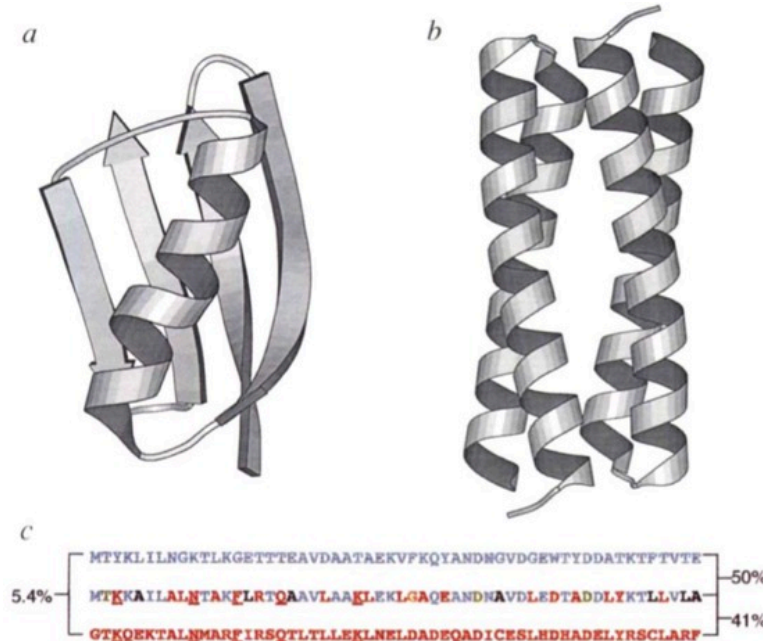
# EPFL The inverse folding problem — design



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

# EPFL The origins: the Paracelsus challenge ('94)

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

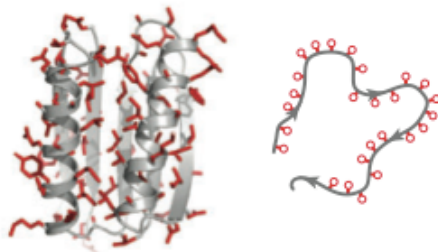


**Fig. 1** Ribbon representation<sup>29</sup> of the folds of **a**, the B1 domain of IgG-binding protein G<sup>5</sup> and **b**, Rop<sup>6</sup>. **c**, An alignment of the sequences of the B1 domain (blue), Rop (red) and Janus. Residues in Janus are coded as follows: blue, residues from B1; red, residues from Rop; underlined red, RNA-binding residues in Rop<sup>13</sup>; 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 Rop<sup>30</sup>. 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 Rop<sup>31</sup>.

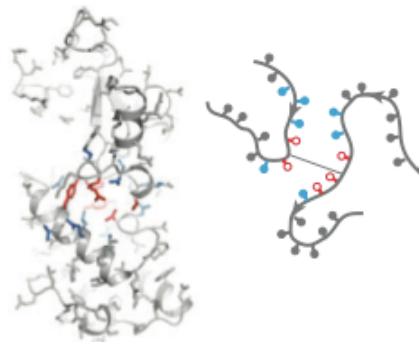
# Multiple tasks for protein design

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

Protein design

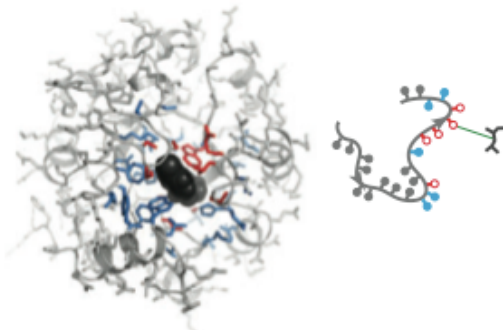


Protein-protein interface design



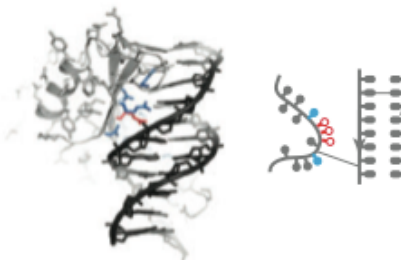
- create high affinity binders
- therapeutic biologics
- artificial sensors/probes

Enzyme design



- tailor enzymatic function
- improve thermostability
- catalyse new reactions

Protein-DNA interface design

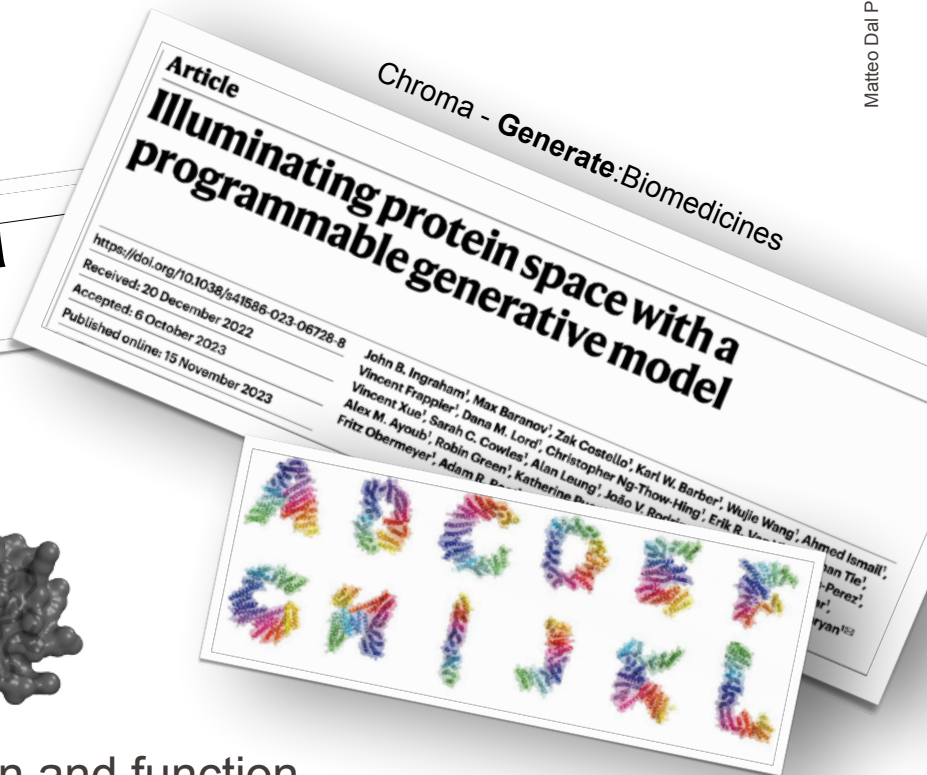
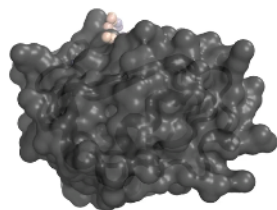
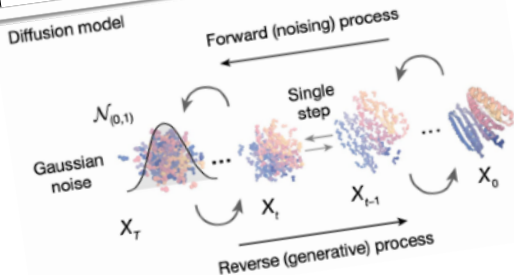
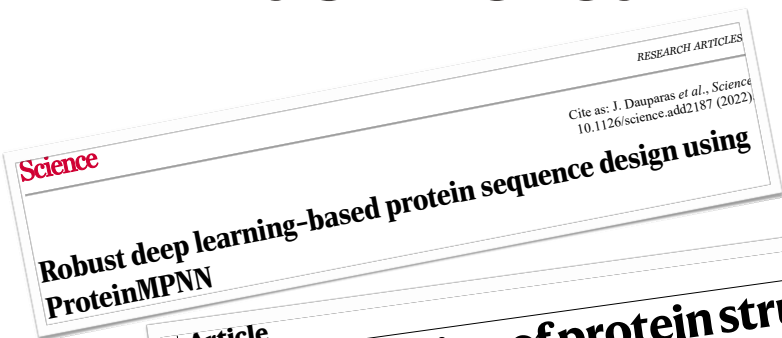


- explore DNA interactions
- new therapeutic solutions

▪

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

# EPFL Machine learning for protein design



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

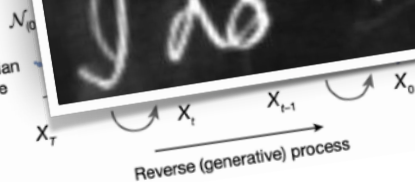
# ... leading to molecular engineering



Article  
**Den  
 func**

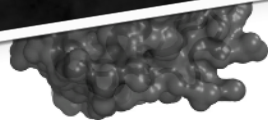
Diffusion mo

Gaussian  
 noise



What I cannot create,  
 I do not understand.

R. Feynman



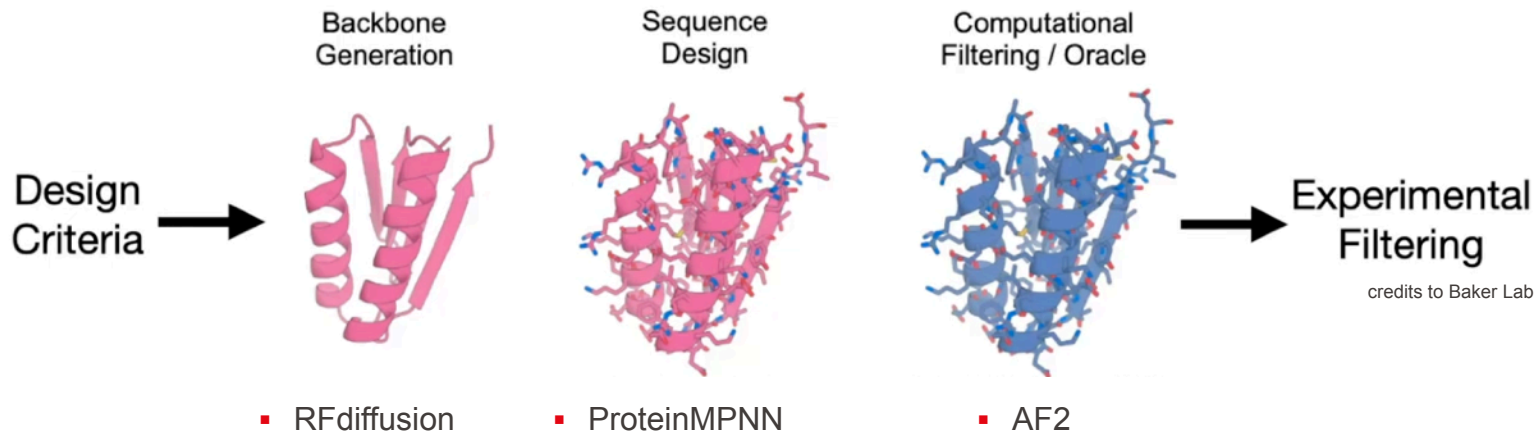
Generate:Biomedicines

with a  
 model

Stello', Karl W. Barber', Wujie Wang', Ahmed Ismail',  
 Christopher Ng-Thow-Hing', Erik R. ...  
 ...an Leung', João V. ...  
 ... Katherine ...  
 ...han Tie', ...Perez',  
 ...r',  
 ...yan<sup>1,2</sup>

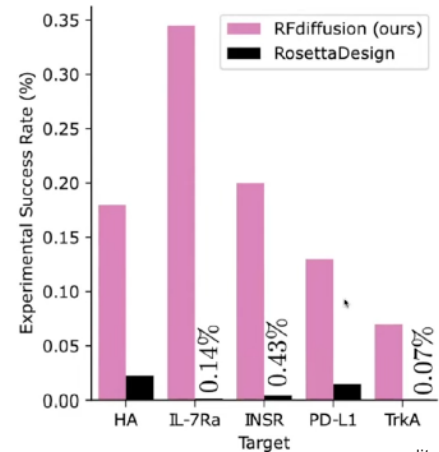
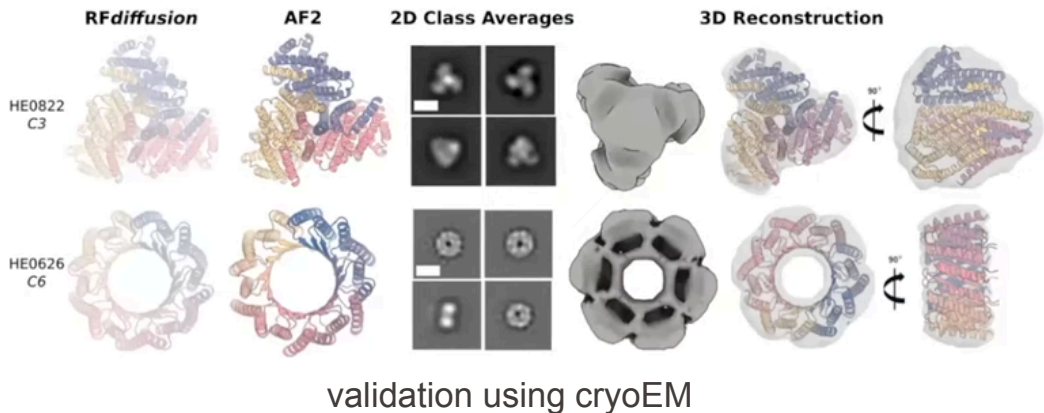
- Application to study protein evolution and function
- Protein engineering for therapeutics, synthetic biology and biotechnology

# 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



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



# De novo design of protein structure and function with RFdiffusion

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

Received: 14 December 2022

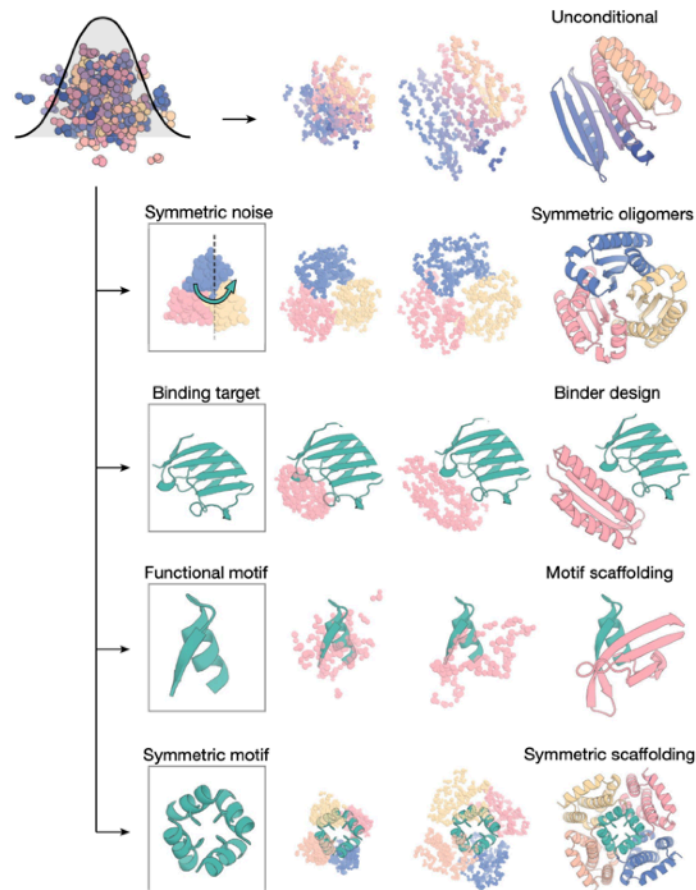
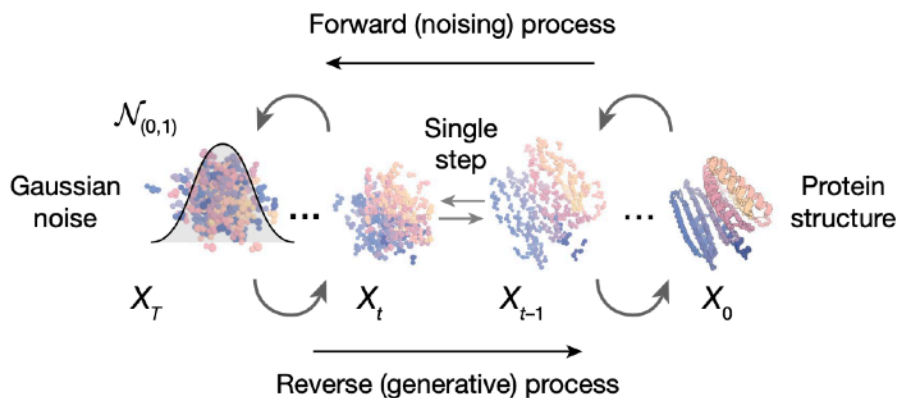
Accepted: 7 July 2023

Published online: 11 July 2023

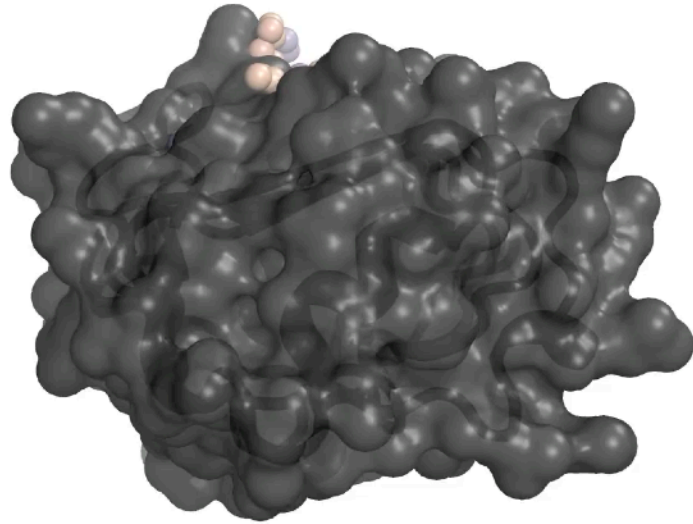
Open access

Joseph L. Watson<sup>1,2,15</sup>, David Juergens<sup>1,2,3,15</sup>, Nathaniel R. Bennett<sup>1,2,3,15</sup>, Brian L. Trippe<sup>2,4,5,15</sup>, Jason Yim<sup>2,6,15</sup>, Helen E. Eisenach<sup>1,2,15</sup>, Woody Ahern<sup>1,2,7,15</sup>, Andrew J. Borst<sup>1,2</sup>, Robert J. Ragotte<sup>1,2</sup>, Lukas F. Milles<sup>1,2</sup>, Basile I. M. Wicky<sup>1,2</sup>, Nikita Hanikel<sup>1,2</sup>, Samuel J. Pellock<sup>1,2</sup>, Alexis Courbet<sup>1,2,8</sup>, William Sheffler<sup>1,2</sup>, Jue Wang<sup>1,2</sup>, Preetham Venkatesh<sup>1,2,9</sup>, Isaac Sappington<sup>1,2,9</sup>, Susana Vázquez Torres<sup>1,2,9</sup>, Anna Lauko<sup>1,2,9</sup>, Valentin De Bortoli<sup>8</sup>, Emile Mathieu<sup>10</sup>, Sergey Ovchinnikov<sup>1,12</sup>, Regina Barzilay<sup>6</sup>, Tommi S. Jaakkola<sup>6</sup>, Frank DiMaio<sup>1,3</sup>, Minkyung Baek<sup>1,3</sup> & David Baker<sup>1,2,14,15</sup>

## 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).

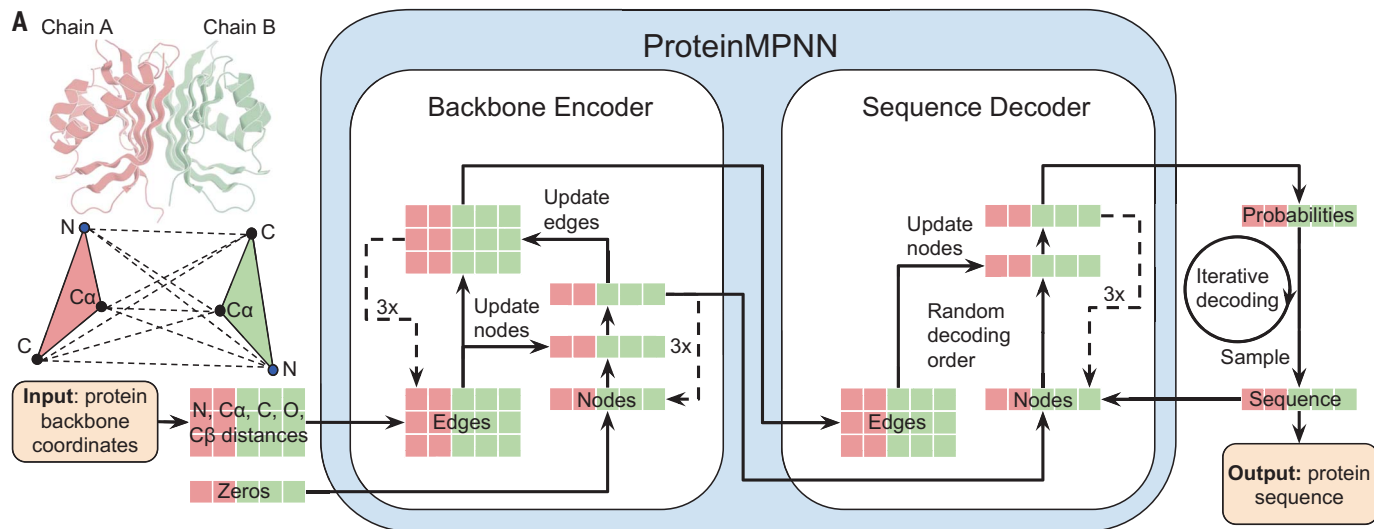
## Robust deep learning–based protein sequence design using ProteinMPNN

J. Dauparas<sup>1,2</sup>, I. Anishchenko<sup>1,2</sup>, N. Bennett<sup>1,2,3</sup>, H. Bai<sup>1,2,4</sup>, R. J. Ragotte<sup>1,2</sup>, L. F. Milles<sup>1,2</sup>, B. I. M. Wicky<sup>1,2</sup>, A. Courbet<sup>1,2,4</sup>, R. J. de Haas<sup>4</sup>, N. Bethel<sup>1,2,4</sup>, P. J. Y. Leung<sup>1,2,5</sup>, T. F. Huddy<sup>1,2</sup>, S. Pellock<sup>1,2</sup>, D. Tischer<sup>1,2</sup>, F. Chan<sup>1,2</sup>, B. Koepnick<sup>1,2</sup>, H. Nguyen<sup>1,2</sup>, A. Kang<sup>1,2</sup>, B. Sankaran<sup>6</sup>, A. K. Bera<sup>1,2</sup>, N. P. King<sup>1,2</sup>, D. Baker<sup>1,2,4\*</sup>

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

\*Corresponding author. Email: dabaker@uw.edu

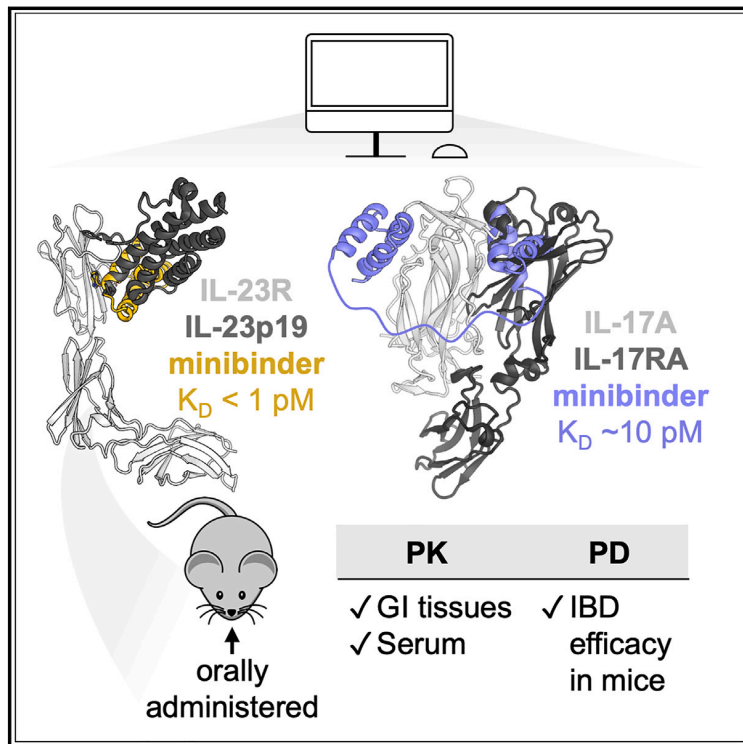
While deep learning has revolutionized protein structure prediction, almost all experimentally characterized de novo protein designs have been generated using physically based approaches such as Rosetta. Here we describe a deep learning–based protein sequence design method, ProteinMPNN, with outstanding 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

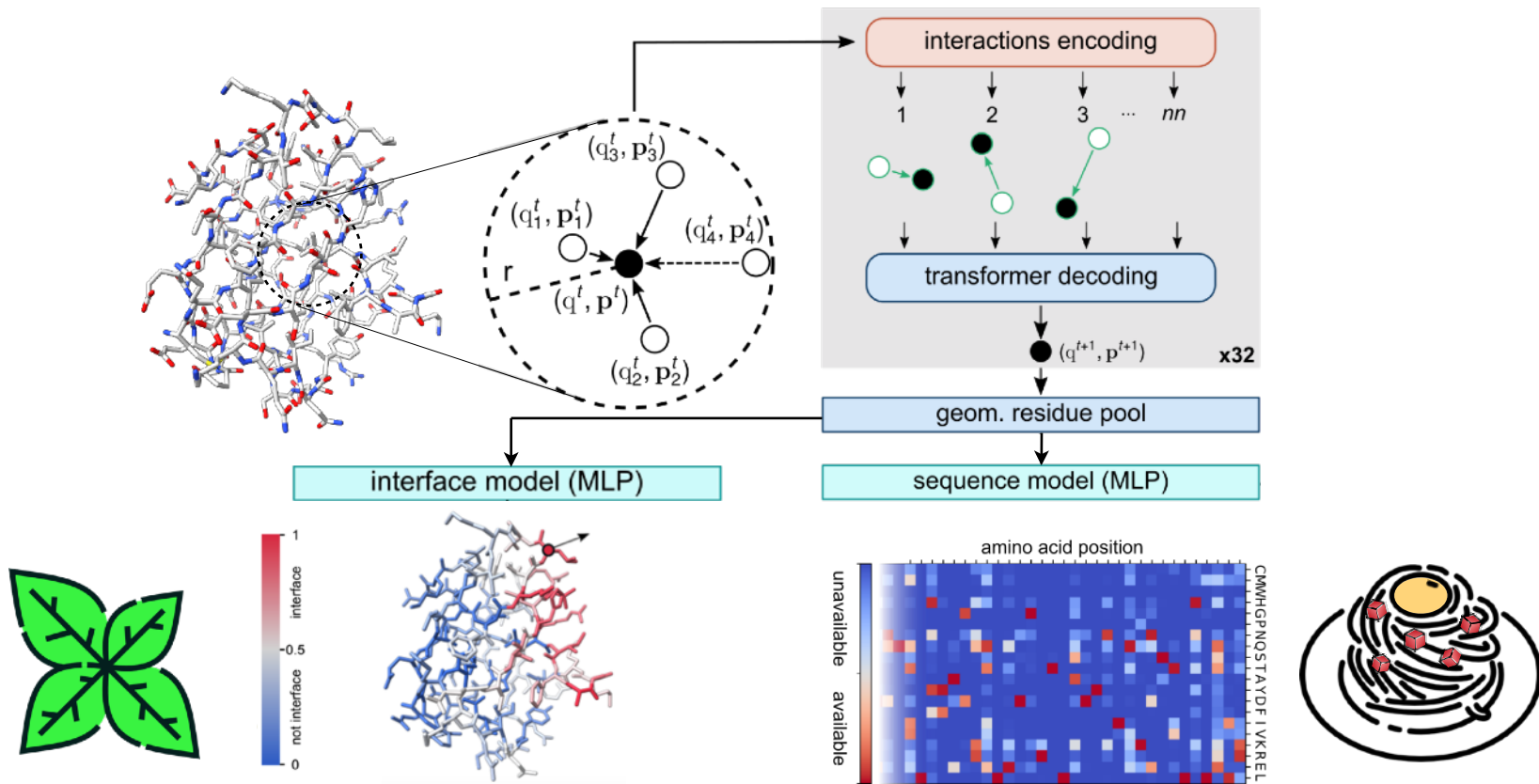
berger389@gmail.com (S.B.),  
dabaker@uw.edu (D.B.)

## 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 Protein Structure Transformer @LBM



**PeSTo: binding interfaces**

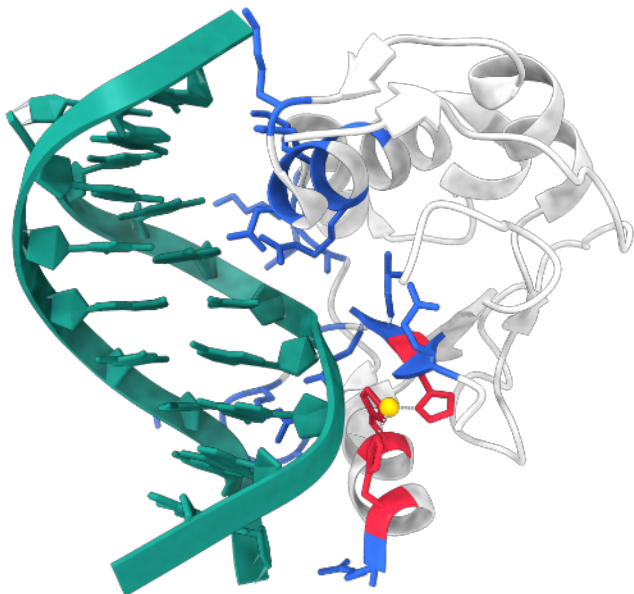
Krapp et al. *Nat Comms* 2023

**CARBonAra: molecular design**

Krapp et al. *Nat Comms* 2024

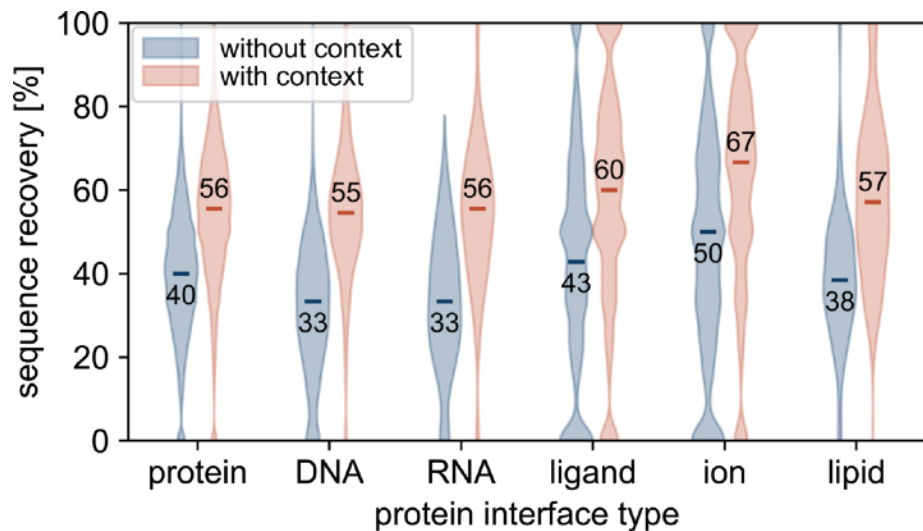
# EPFL Unique ability — context awareness

- example with context



colicin E7

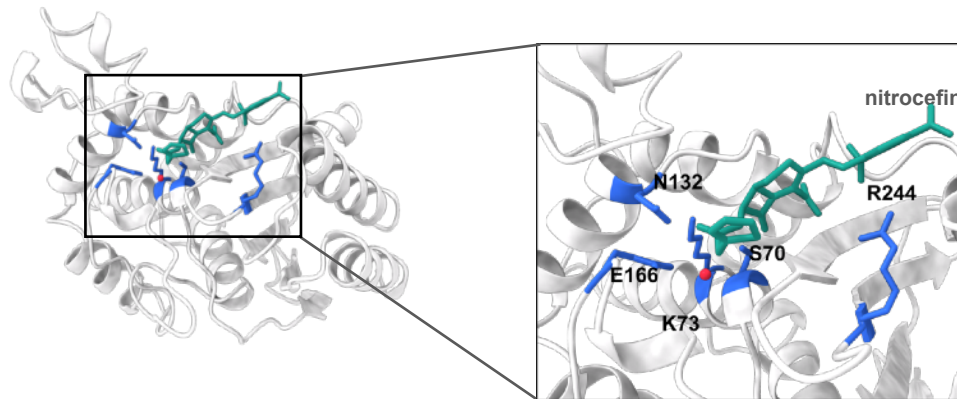
- large-scale benchmark



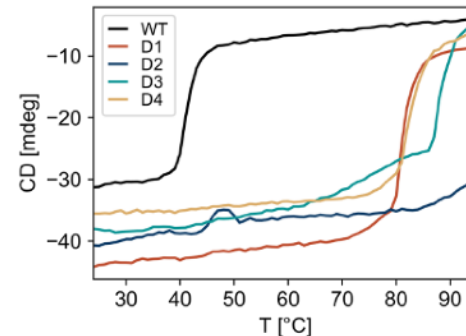
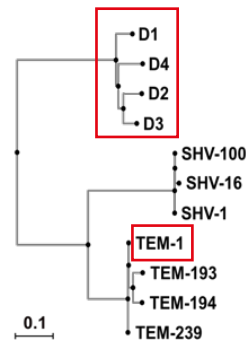
1000 structures sampled with maximum 30% sequence identity and separate C.A.T.H. classification from training set

# Can we re-engineer an enzyme?

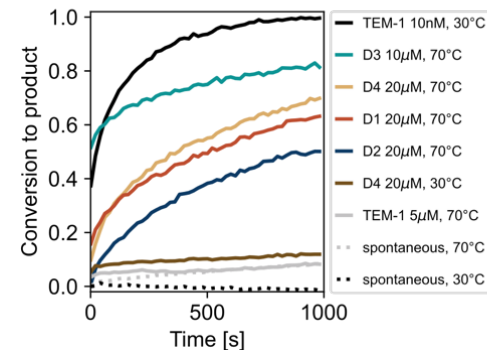
## ■ TEM-1 serine $\beta$ -lactamase



only 50% sequence identity



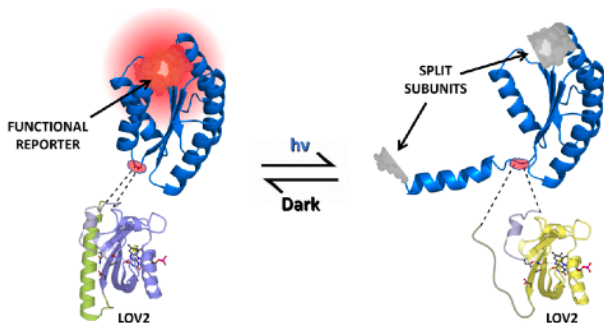
- sequences generation with substrate as constraint
- selected 10 top-ranked predictions based on pIDDT
- 4/10 designs are soluble and monomeric
- they are folded and more thermostable than wild-type TEM-1
- catalytically active at high T - not as the wild-type yet
- represent a separate subclass of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



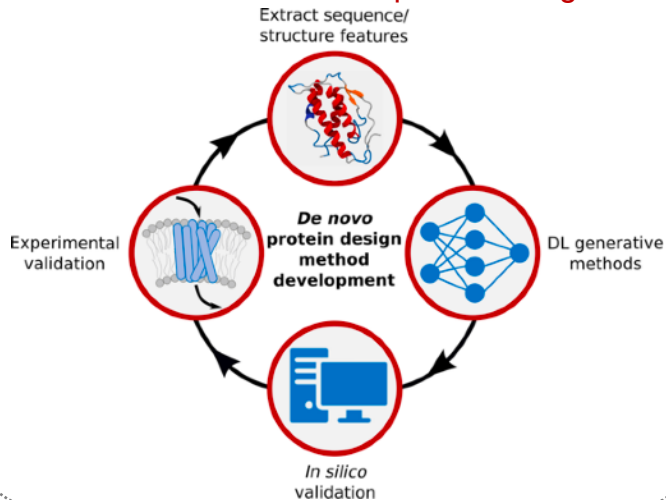
# Laboratory of Protein and Cell Engineering

AI FOR PROTEIN DESIGN

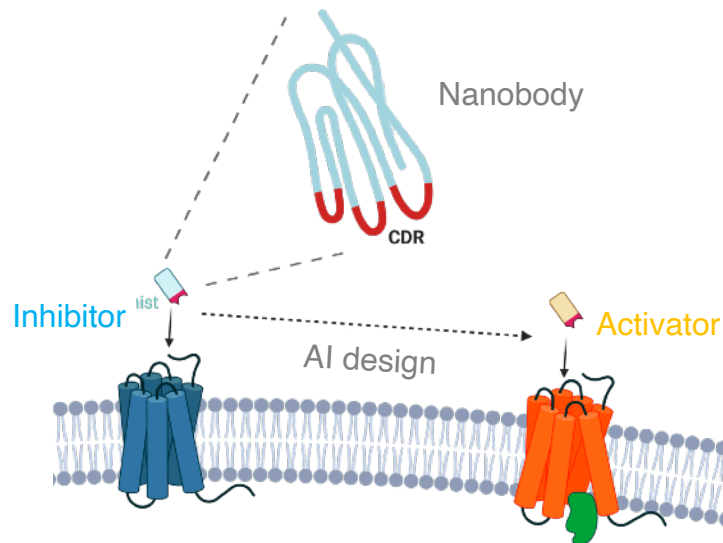
## De novo optogenetic switches



## De novo membrane protein design



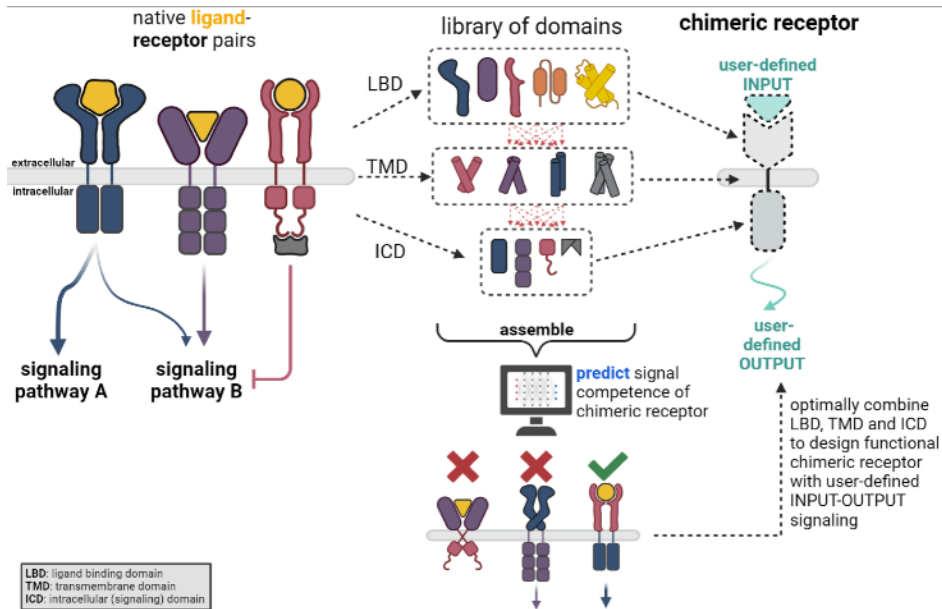
## De novo therapeutic protein ligands



Barth Lab

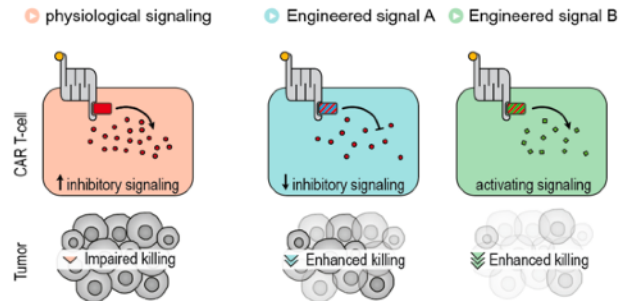
# Laboratory of Protein and Cell Engineering

## Allosteric biosensors for engineered cell therapies

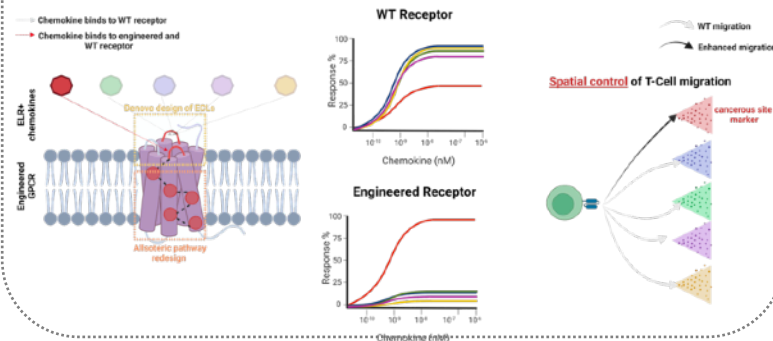


IMMUNOENGINEERING

## G-protein signal rewiring to boost CAR T-cells



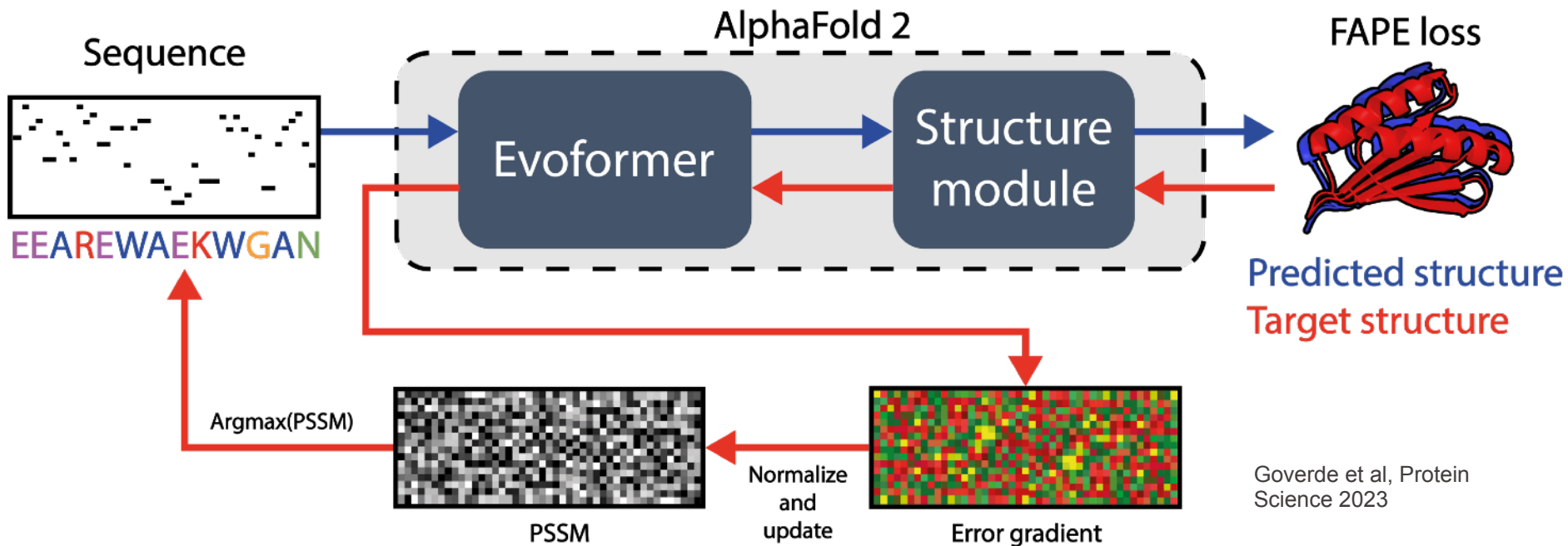
## Redirecting cell migration with designed receptors



Barth Lab

# EPFL Inverting AlphaFold for protein design

Bruno Correia, Laboratory of Protein Design and Immunoengineering

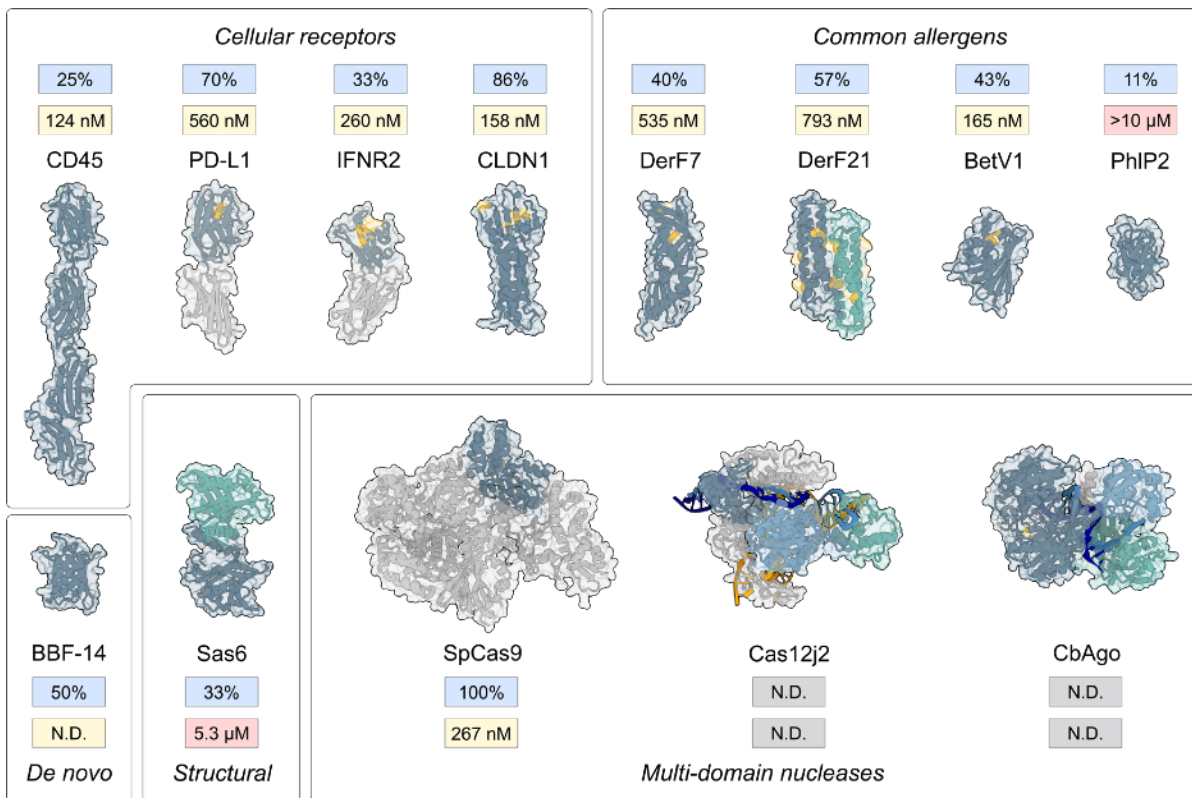


-Final sequences designed with proteinMPNN on AF2 generated backbones

■

# High experimental success rates in binder design

Bruno Correia, Laboratory of Protein Design and Immunoengineering



Experimental success rate

Highest affinity binder (no experimental optimisation)

BindCraft

# We weren't alone !!!!

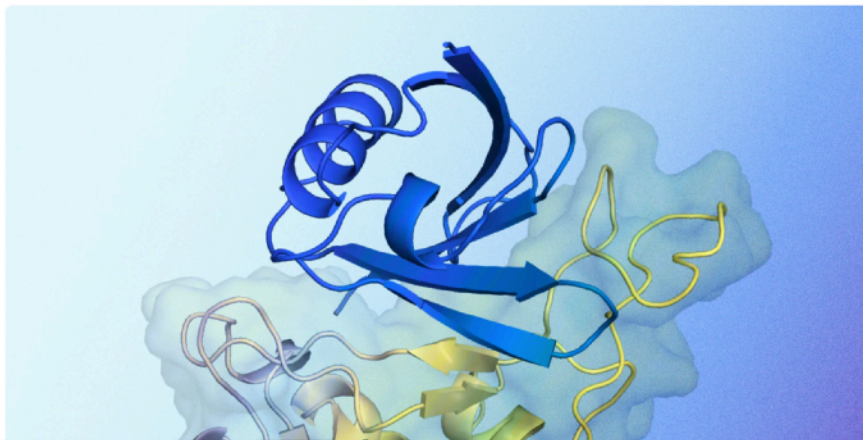
Google DeepMind About ▾ Research Technologies ▾ Impact Discover ▾

## AlphaProteo generates novel proteins for biology and health research

5 SEPTEMBER 2024

Protein Design and Wet Lab teams

[← Share](#)



- ... but also
- **BoltzGen**
  - **Chai-1**
  - **ProtGPT2**
  - **ESMDesign**
  - **and many others**

# The future is bright and exciting ...

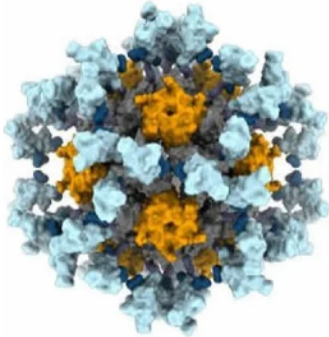
... biomolecular design will address many societal needs

## ■ Medicine

vaccines & antivirals

smart medicines

drug delivery



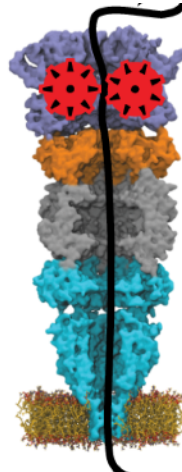
■ SARS-CoV-2 RBD  
nanoparticle immunogen (Cell 2020)

## ■ Biotechnology

protein-silicon devices

bio-based computers

nanoscale manufacturing



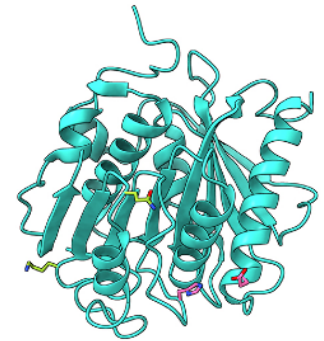
SM proteomics with  
biological nanopores  
(Nat Chem 2021)

## ■ Sustainability

artificial photosynthesis

CO<sub>2</sub> sequestration

plastic degradation



FAST-PETase  
(Nature 2022)

# 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](#)

