



# Toward the solution of the protein structure prediction problem

Received for publication, April 20, 2021, and in revised form, June 7, 2021 Published, Papers in Press, June 11, 2021,  
<https://doi.org/10.1016/j.jbc.2021.100870>

**Robin Pearce<sup>1</sup> and Yang Zhang<sup>1,2,\*</sup>**

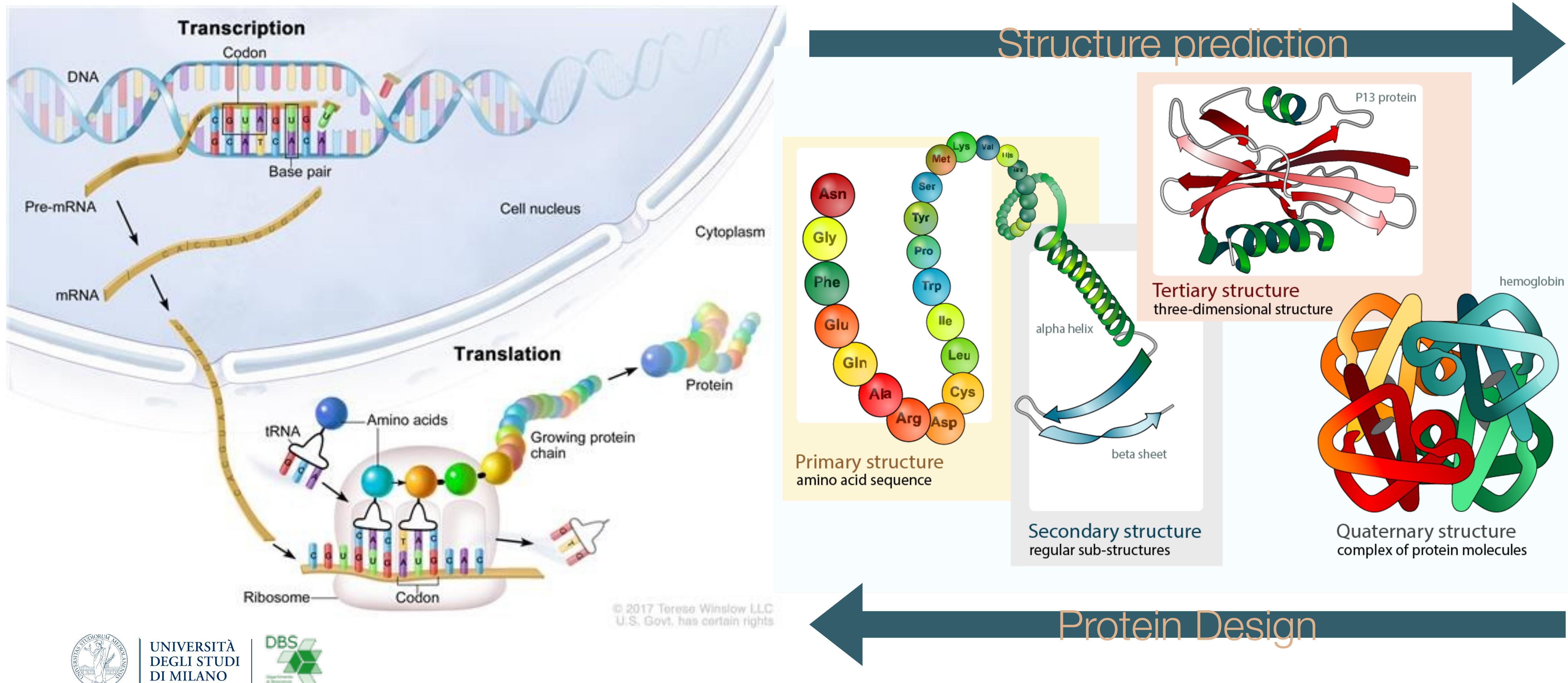
*From the <sup>1</sup>Department of Computational Medicine and Bioinformatics, <sup>2</sup>Department of Biological Chemistry, University of Michigan, Ann Arbor, Michigan, USA*

Edited by Wolfgang Peti

Structures predictions and  
molecular docking

Structural Bioinformatics

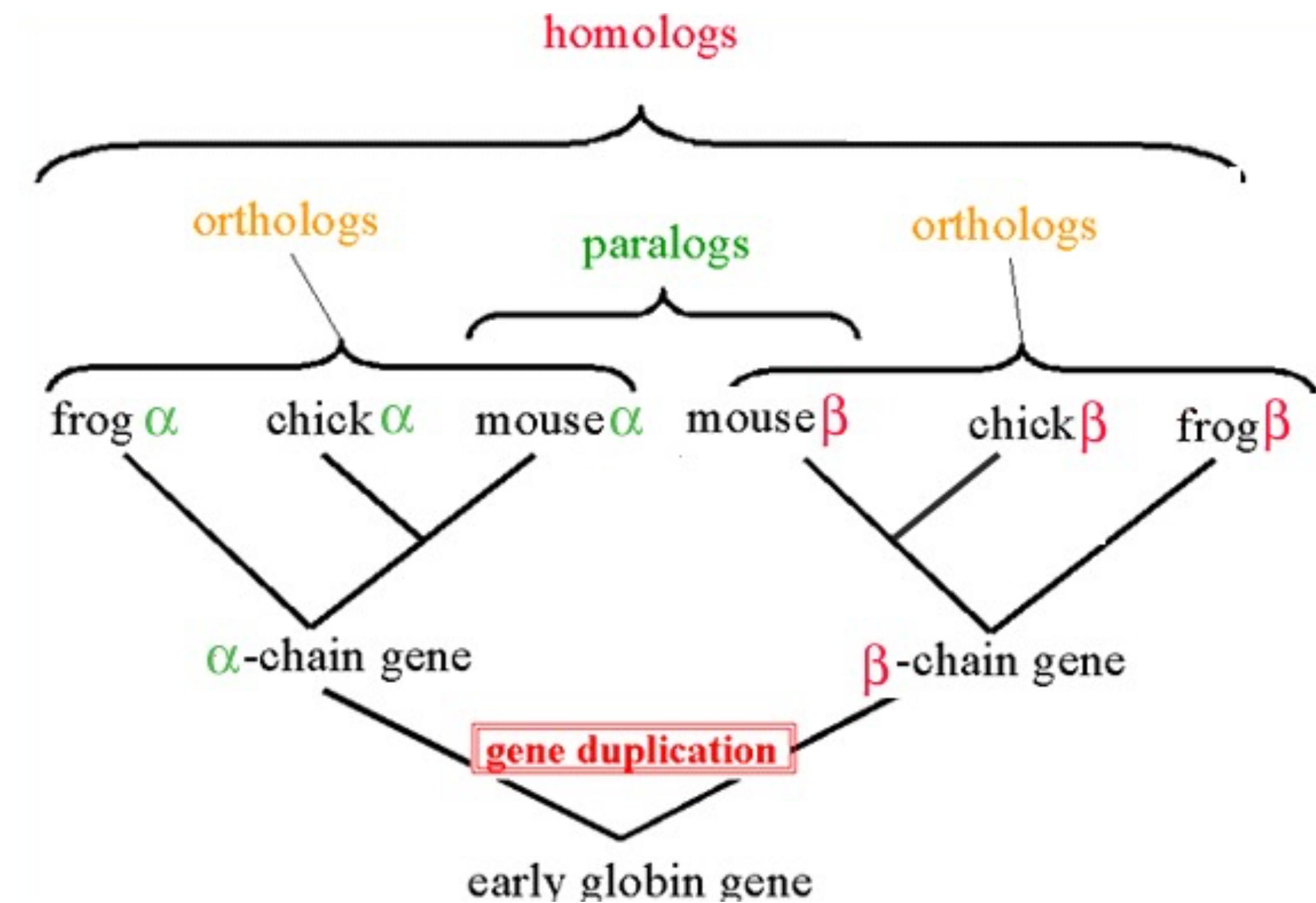
Proteins are encoded in DNA, produced by ribosomes as linear polymers of amino acids, and have a hierarchical structure.



# Protein structure and function are intimately linked

Upon natural evolution, through sequence modifications, structural motifs performing specific functions have been reused and adapted to work in different contexts.

As a result if one takes proteins performing specific functions, even in very unrelated organisms, they still show high structural similarity and some sequence similarity.

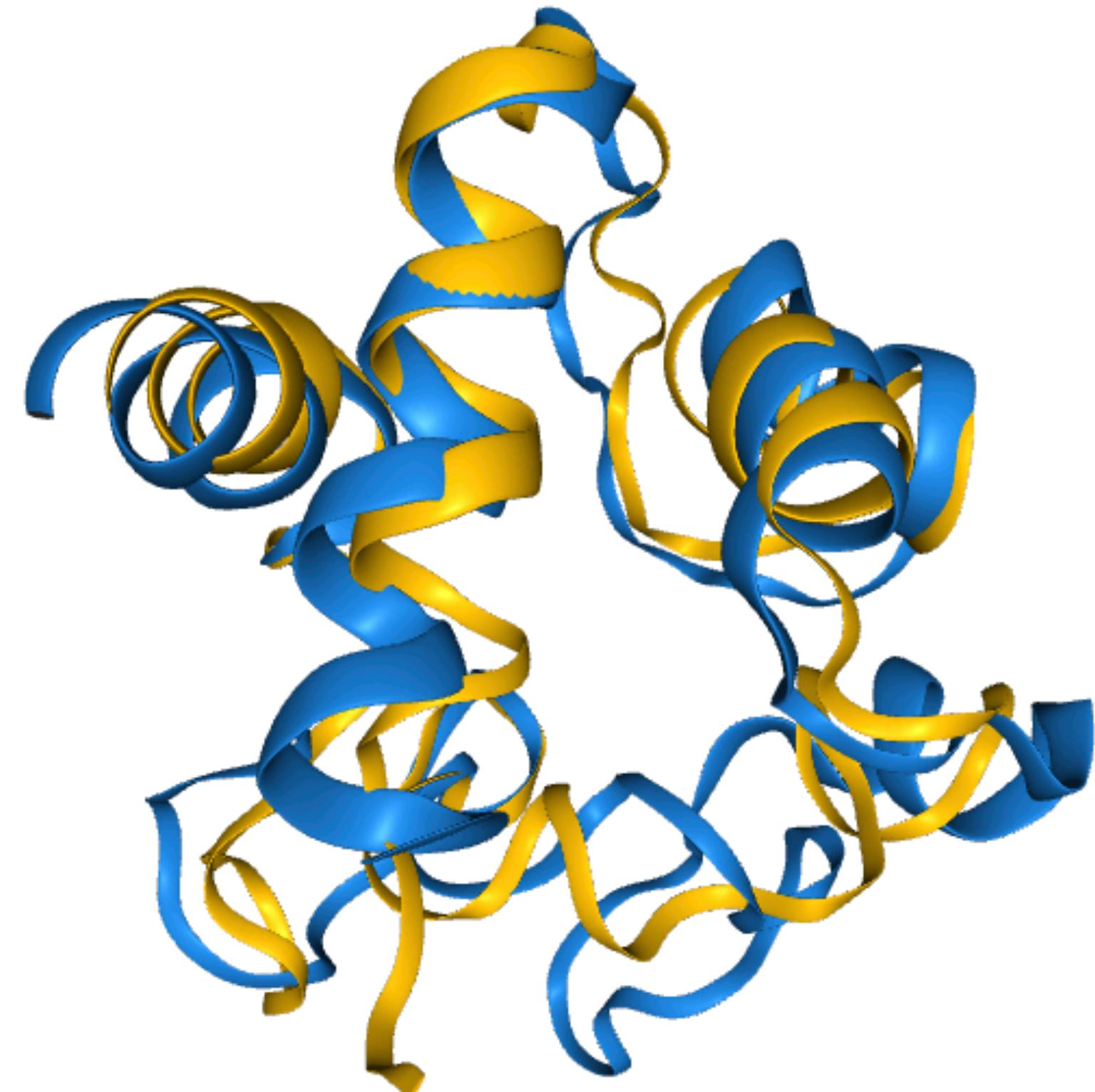


# Cytochrome C: *human* vs *thermus thermophilus* (23% seq identity)

Q	4	EKGKKIFIMKCSQCHTVE--KGGKHKTGPNLHGLFGRKTGQAPGYSYTAANKNKGII-WGEDTLMEYLENPKKYIPGTM
		E+G+++F C+ CH V GP L R + GI+ + L ++ +P PG KM
T	219	ERGQQVFFQQNCAACHGVARSMPAV-IGPELGLWGNRTSL-----GA--GIVENTPFNI KAWTRDPAGMKPGVKM
Q	81	IFVGI--KKKEERADLIALYLKKATNE
		G +E+ L+ YL+ E
T	285	--PGFPQLSEEDLDALVRYLEGLKVE

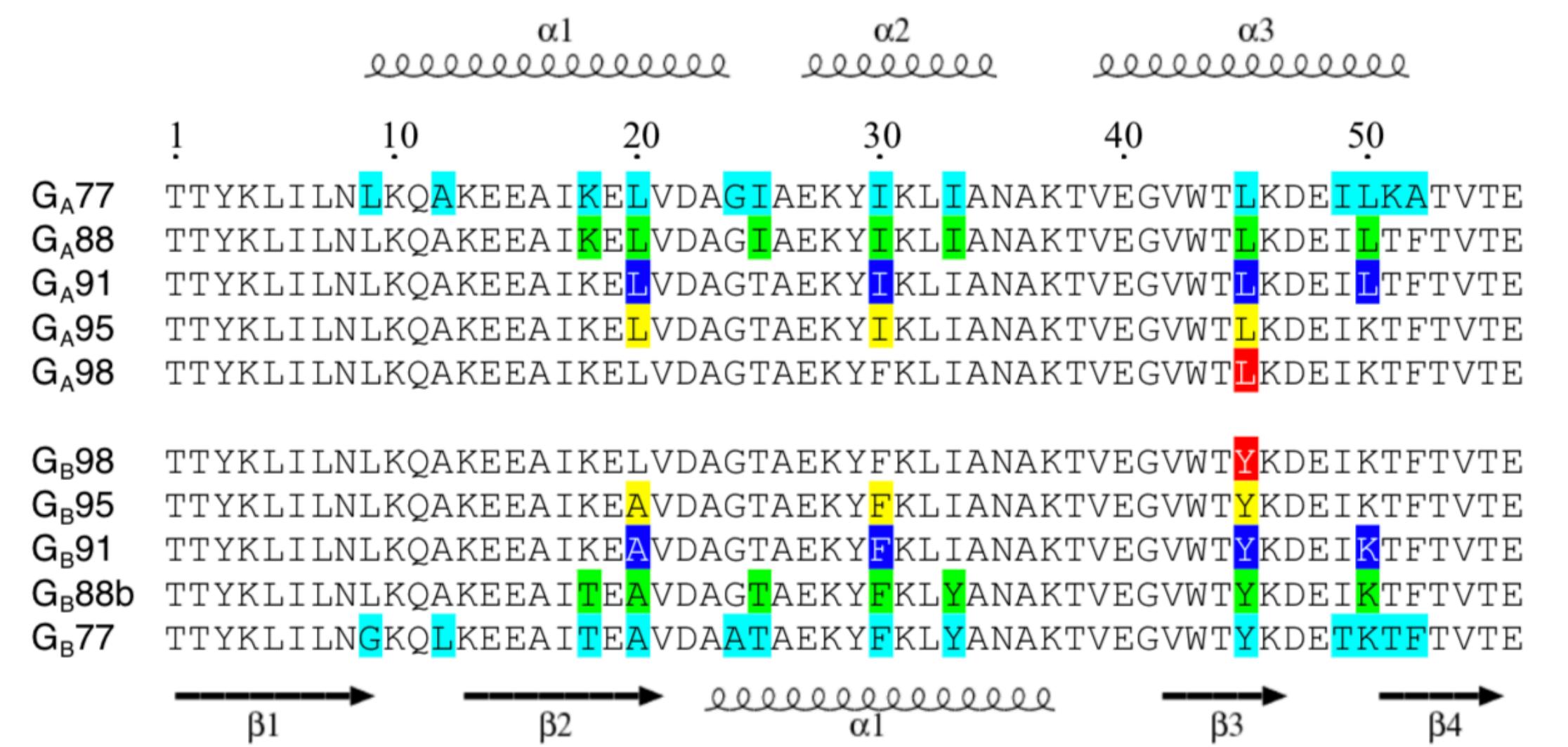
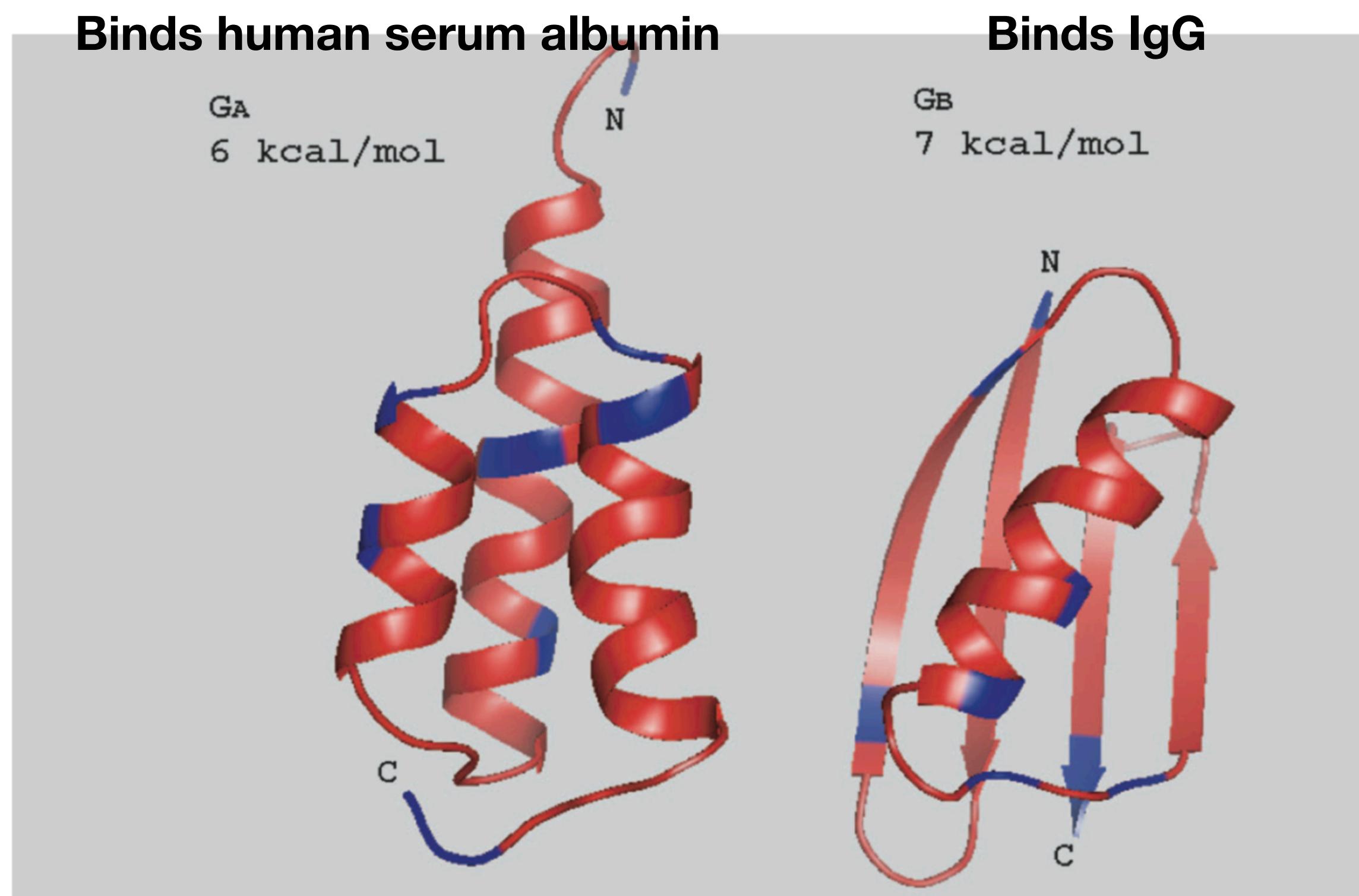
**Ideally a structure superposition should allow us to observe how sequences have evolved, in practice since we do not know enough structure, we align sequences and when aligned sequences are similar we hypothesize structure and function similarity.**

More Conserved -> Function <-> Structure <->  
Sequence <- Less Conserved



...yet small sequence variations can lead to dramatic differences.

The correlation of protein structures and sequences in nature is not a necessary condition imposed by chemistry, indeed it is possible to design proteins with different structures and very similar sequences:



**1AA is enough to determine the fold! And the function?**

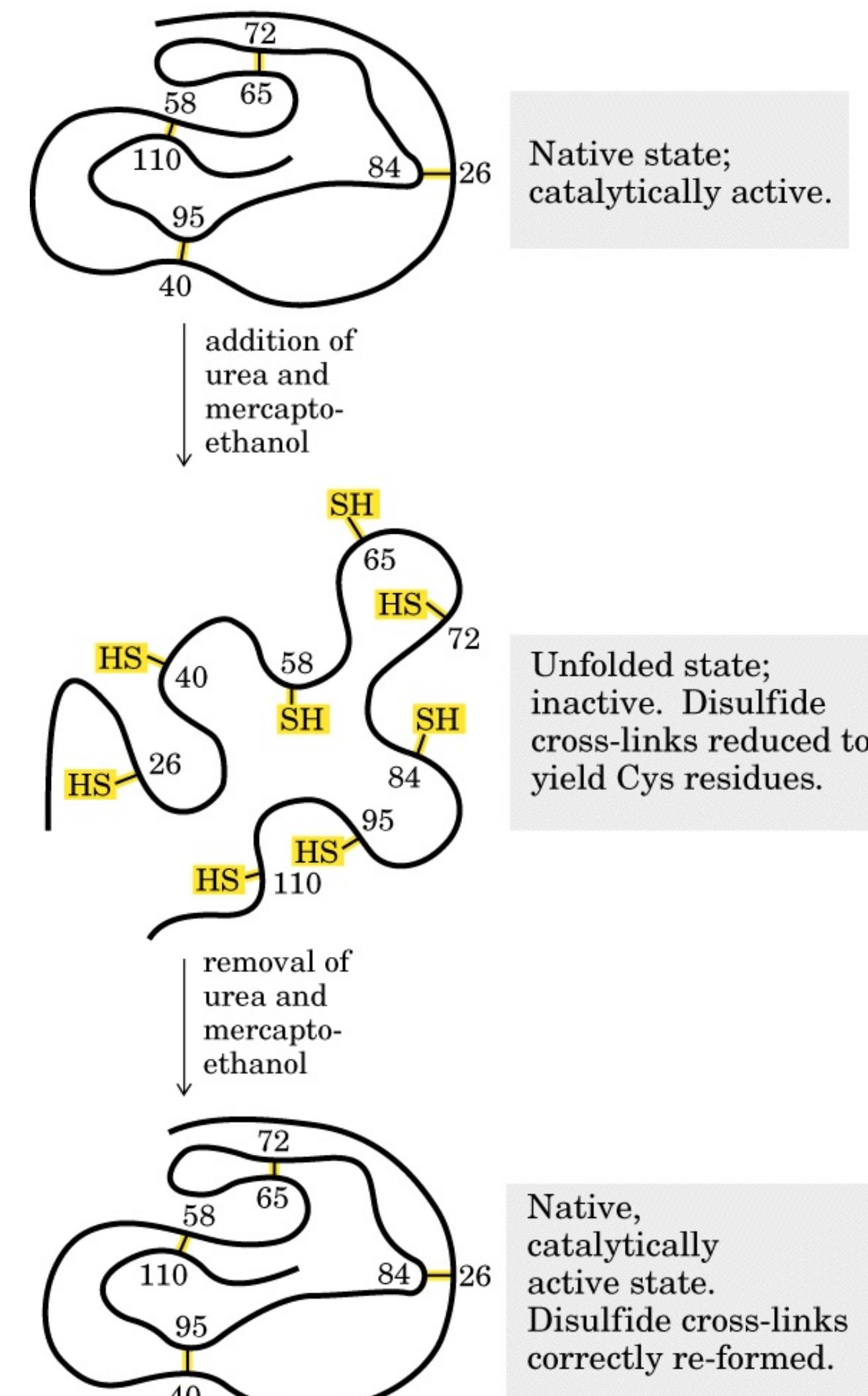
**G<sub>A</sub>98 exhibits diminished affinity for HSA but has acquired affinity for IgG. G<sub>B</sub>98 binds tightly to IgG but not HSA.**



# Anfinsen: The native structure of a protein is dictated by its sequence

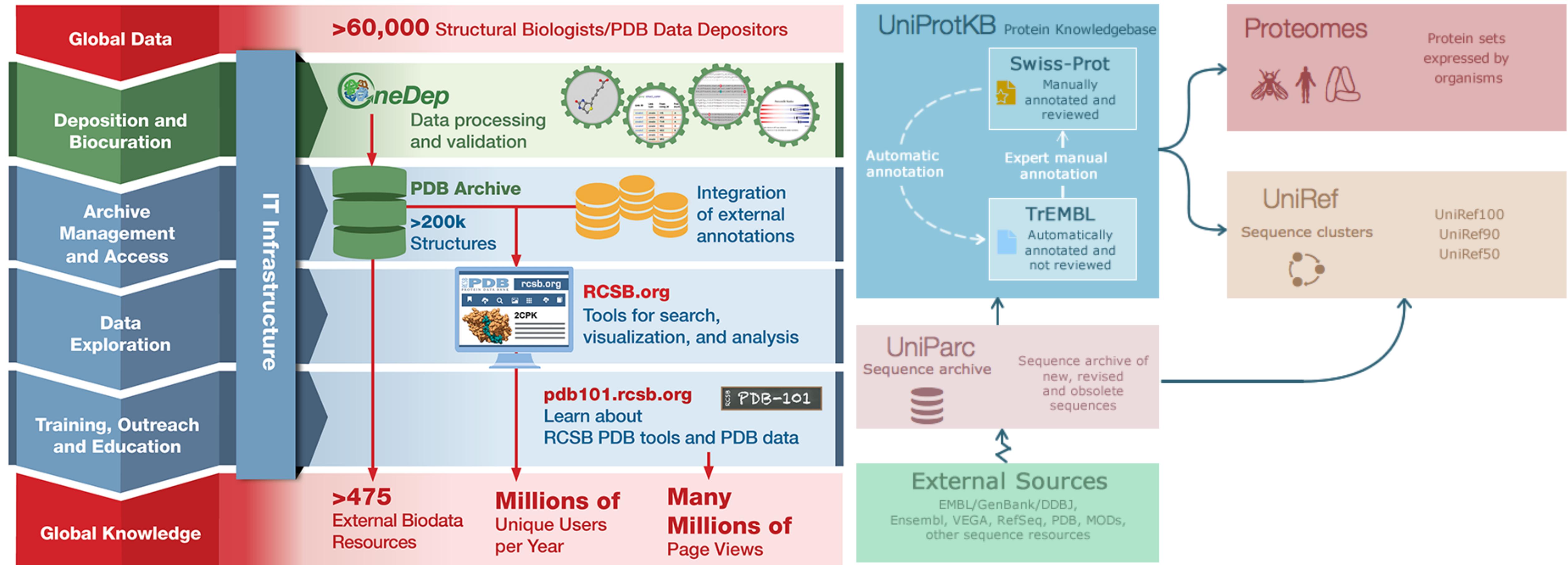
- Small protein can fold independently (Anfinsen's experiment with ribonuclease A)
- Large proteins fold with the assistance of molecular chaperons

*Principles That Govern Protein Folding*  
**Anfinsen C.B.**  
*Science (1974) 181, 223-230.*

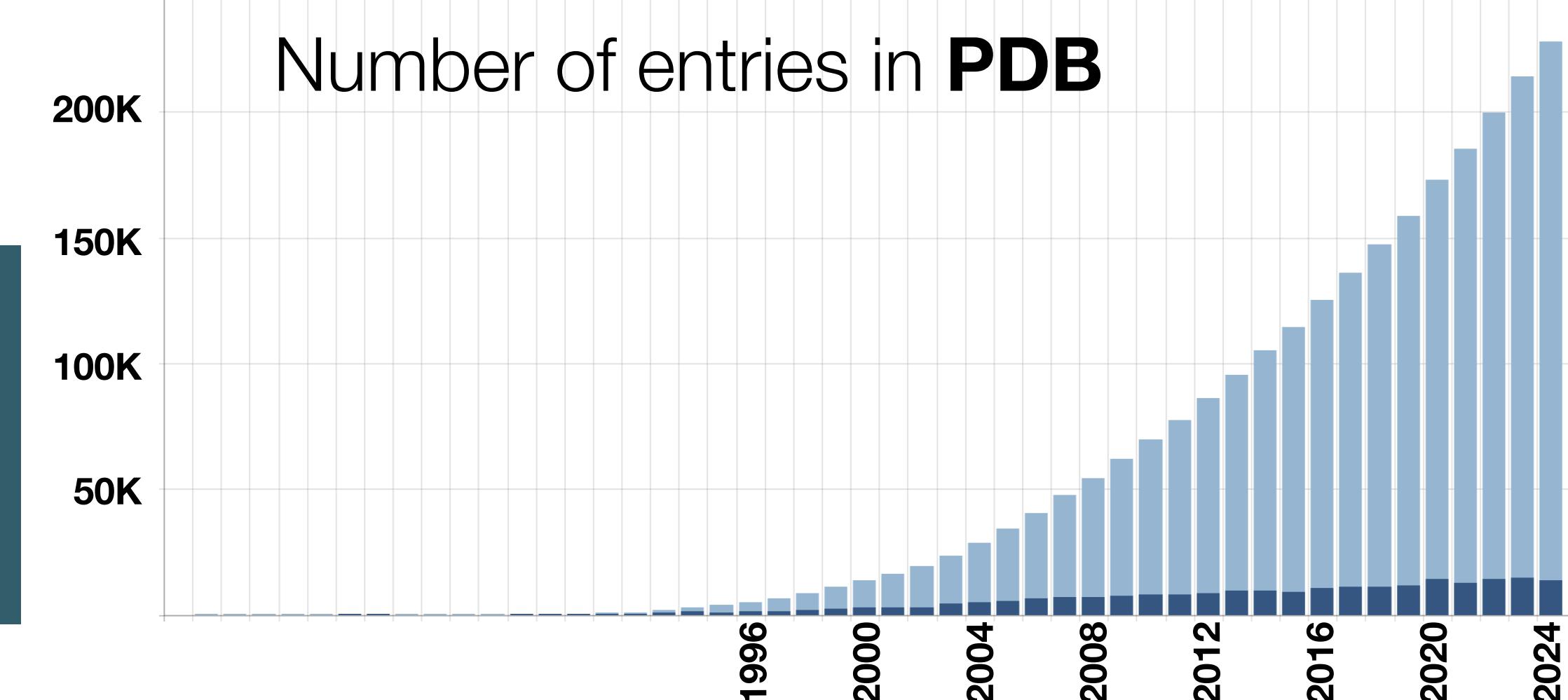
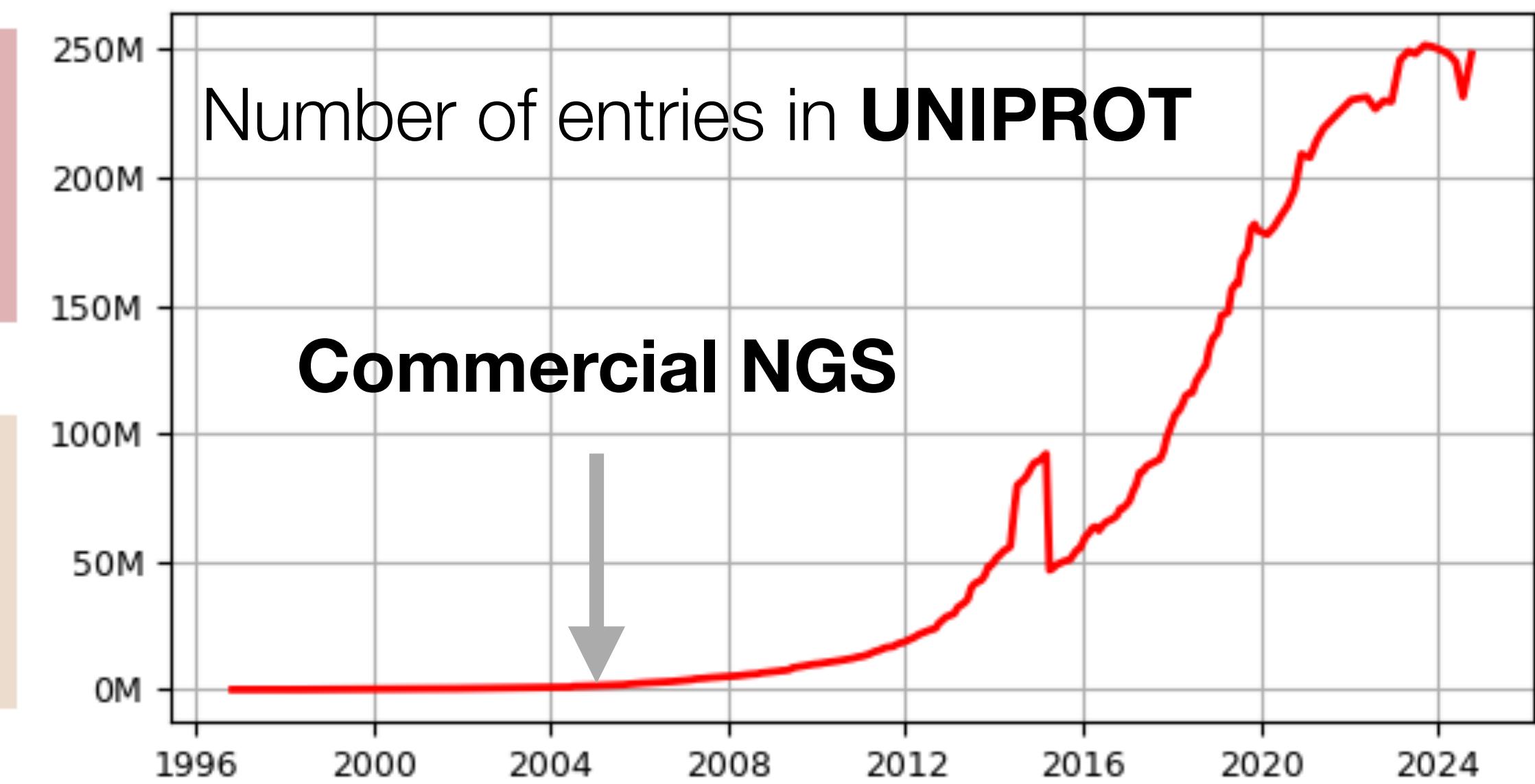
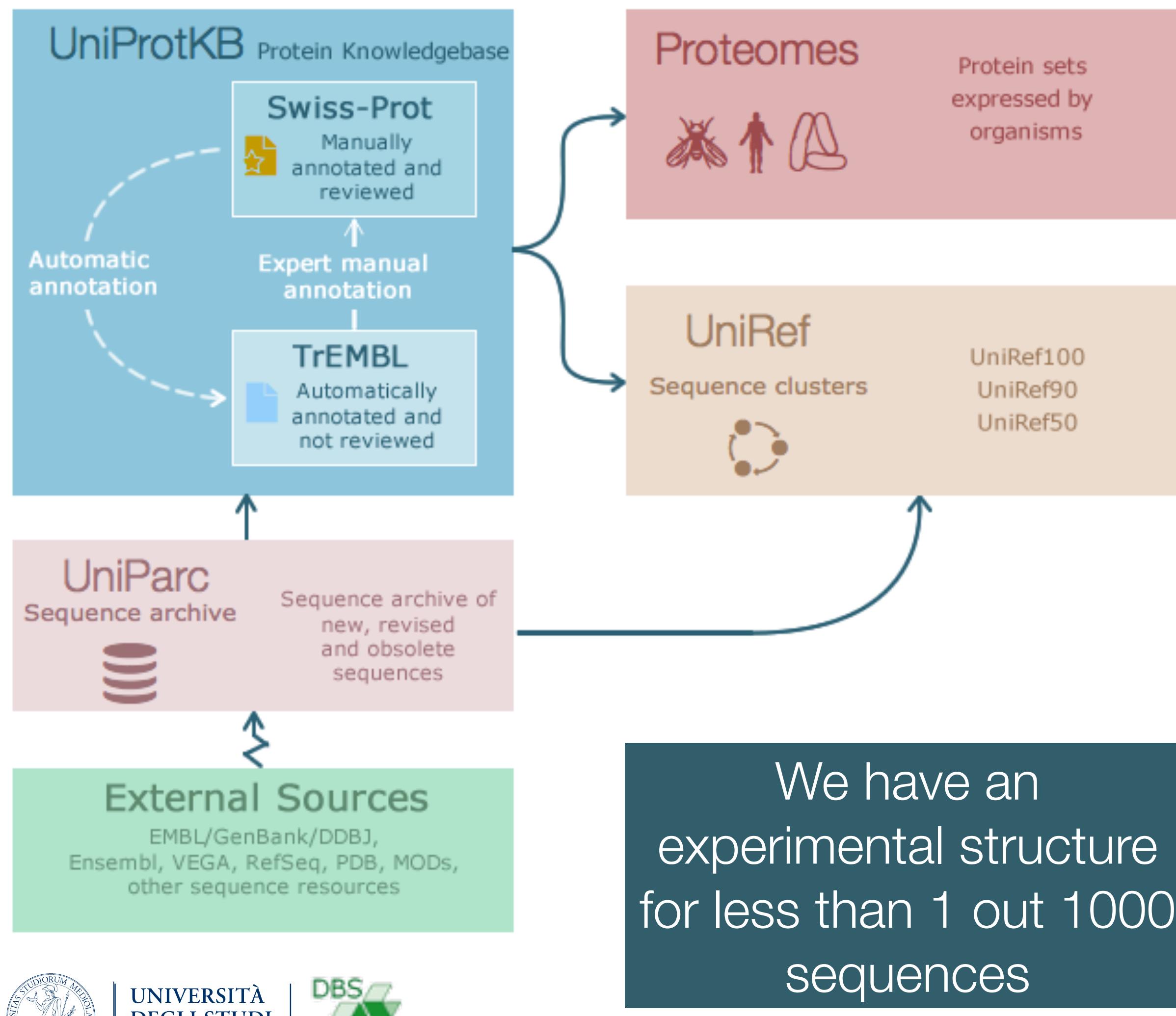


If proteins can refold spontaneously *in vitro*, this suggests that all the information needed to fold is stored in their sequence.

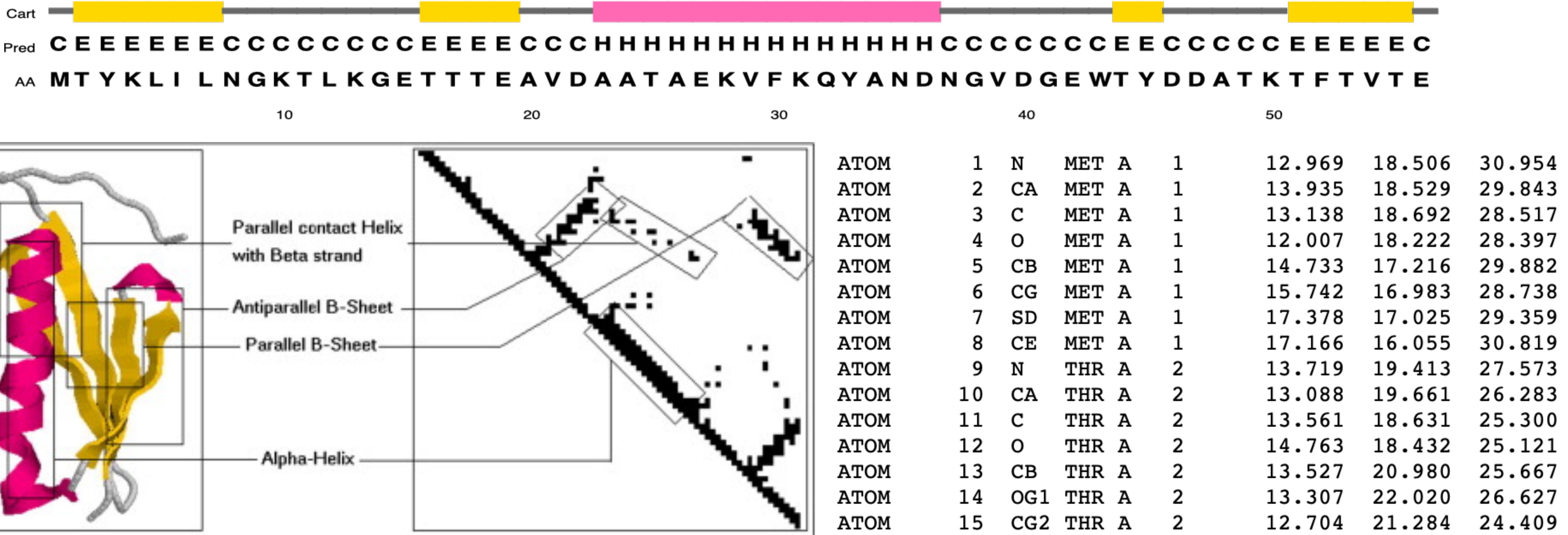
# The data: public, curated databases



# UNIPROT, Next Generation Sequencing and the need for protein structure prediction



# Data representation: sequence, secondary structure, contact maps, 3D structures



Contact or distance maps are matrices reporting either a 0/1 or a distance between two amino acids (usually considering the C<sub>b</sub> carbon)

# Secondary structure prediction is relatively easy but secondary structure can depend on tertiary and quaternary structure

## Useful as:

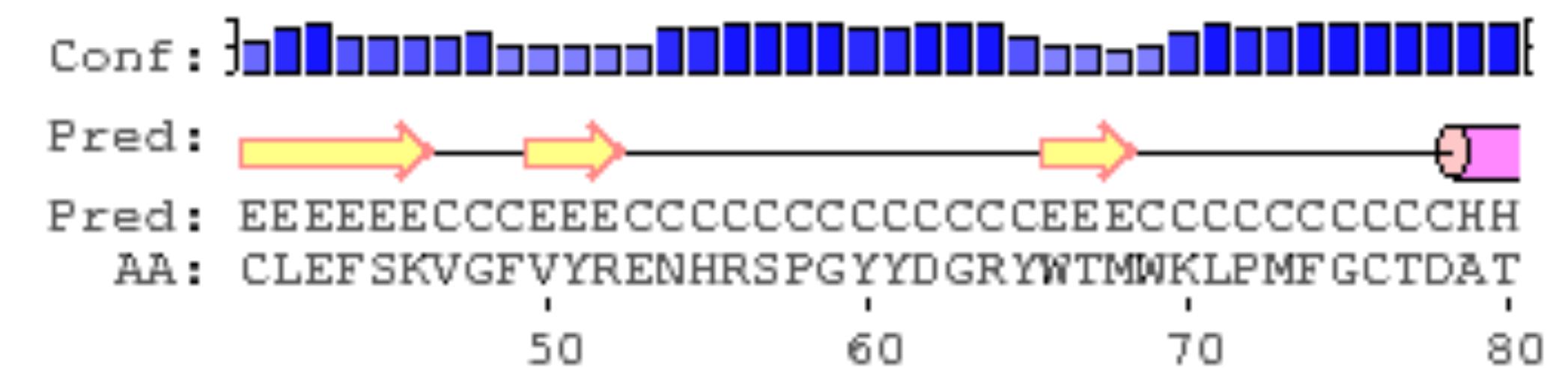
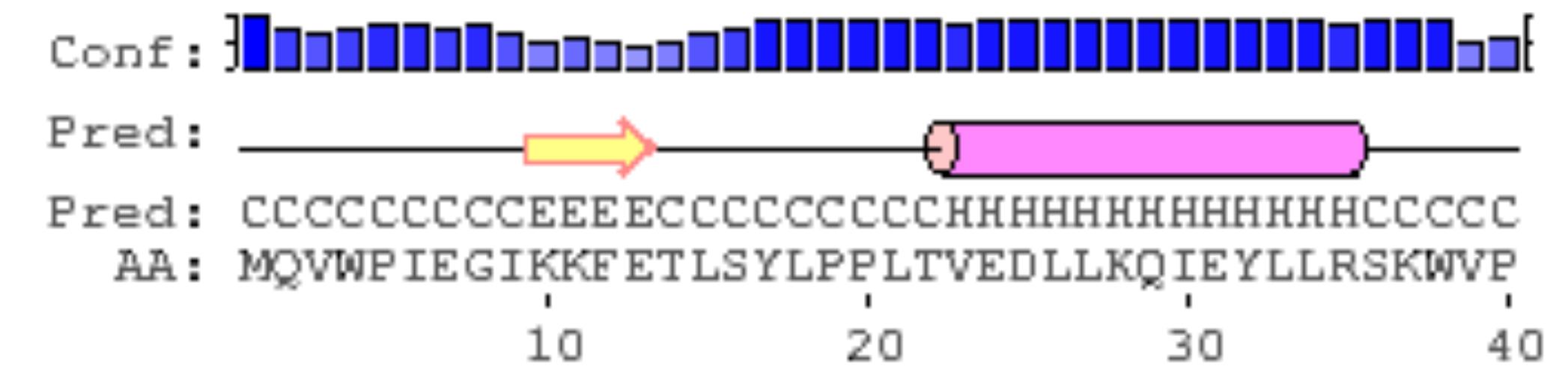
- a step towards tertiary structure prediction
- Sequence alignment
- Structure determination at intermediate resolution
- Prediction of other properties (aggregation, trans-membrane, disorder)

First generation of predictors (~1970) were based on single residues frequencies, second generation (~1990) took into account patterns for segments. Both biased by the small size of the PDB and the lack of systematic assignment.

Accuracy ~60%

Sequence alignment is performed (psi-blast) and neural networks are trained to assign secondary structures (from one letter sequence to one letter secondary structure)

Accuracy ~81%



Conf: [blue bars]  
Pred:  
Pred: CC  
AA: GC

<http://bioinf.cs.ucl.ac.uk/psipred/>  
<http://www.compbio.dundee.ac.uk/jpred>

Phys-chem based methods like MD simulations can predict protein structures, but are limited by size and folding time scale.

**Classical molecular dynamics simulations of a protein in explicit environment with an energy function that tries to approximate the interactions (Force Field):**

$$V(r) = \sum_{bonds} k_b(b - b_0)^2 + \sum_{angles} k_\theta(\theta - \theta_0)^2 + \sum_{torsions} k_\phi[\cos(n\phi + \delta) + 1] \\ + \sum_{nonbond\ pairs} \left[ \frac{q_i q_j}{r_{ij}} + \frac{A_{ij}}{r_{ij}^{12}} - \frac{C_{ij}}{r_{ij}^6} \right]$$

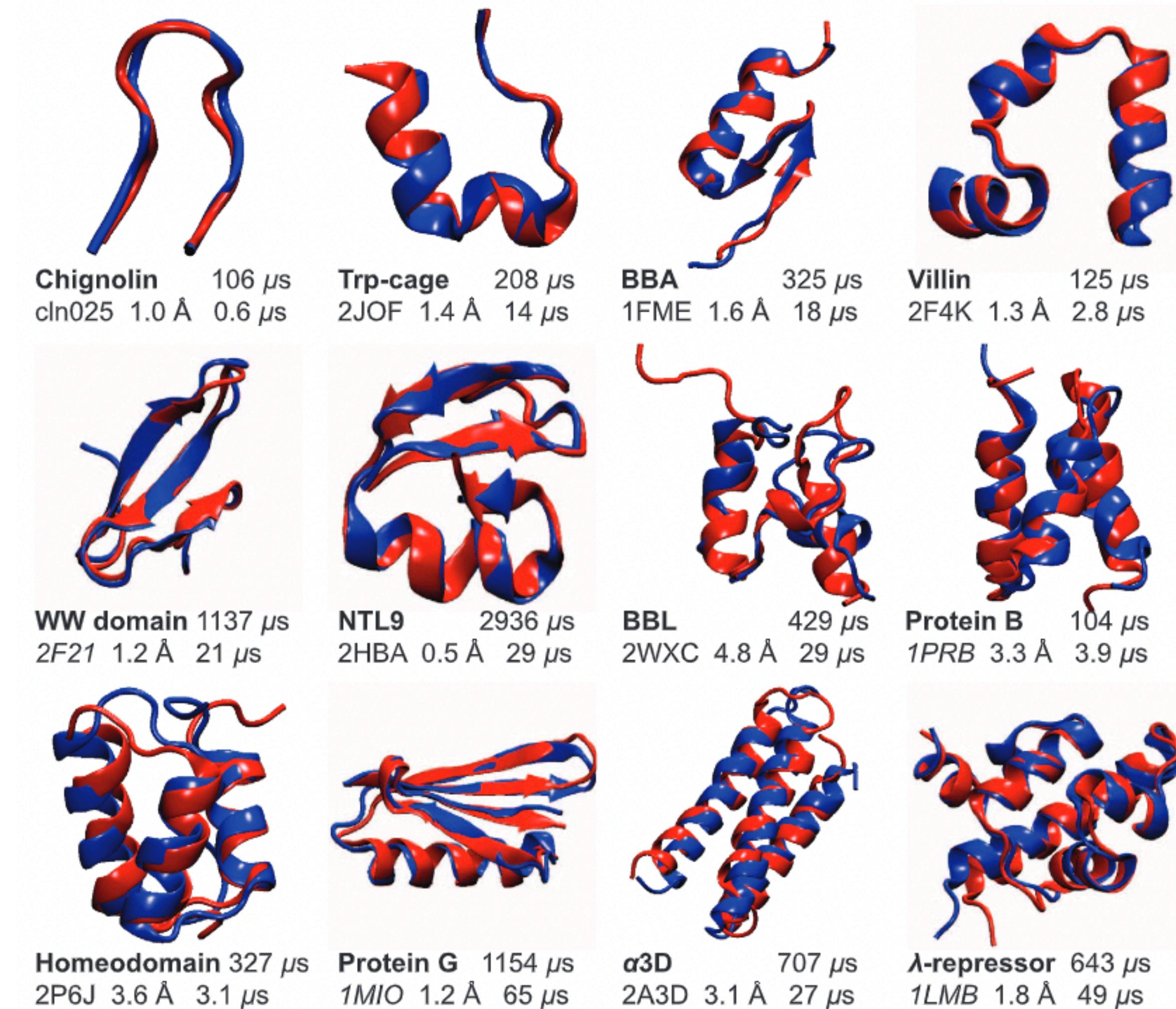
first approx for vibrations  
(anharmonic potential can be used for more accurate vibrations)

geometrical consideration  
pi-bonds, etc

point charge Coulomb

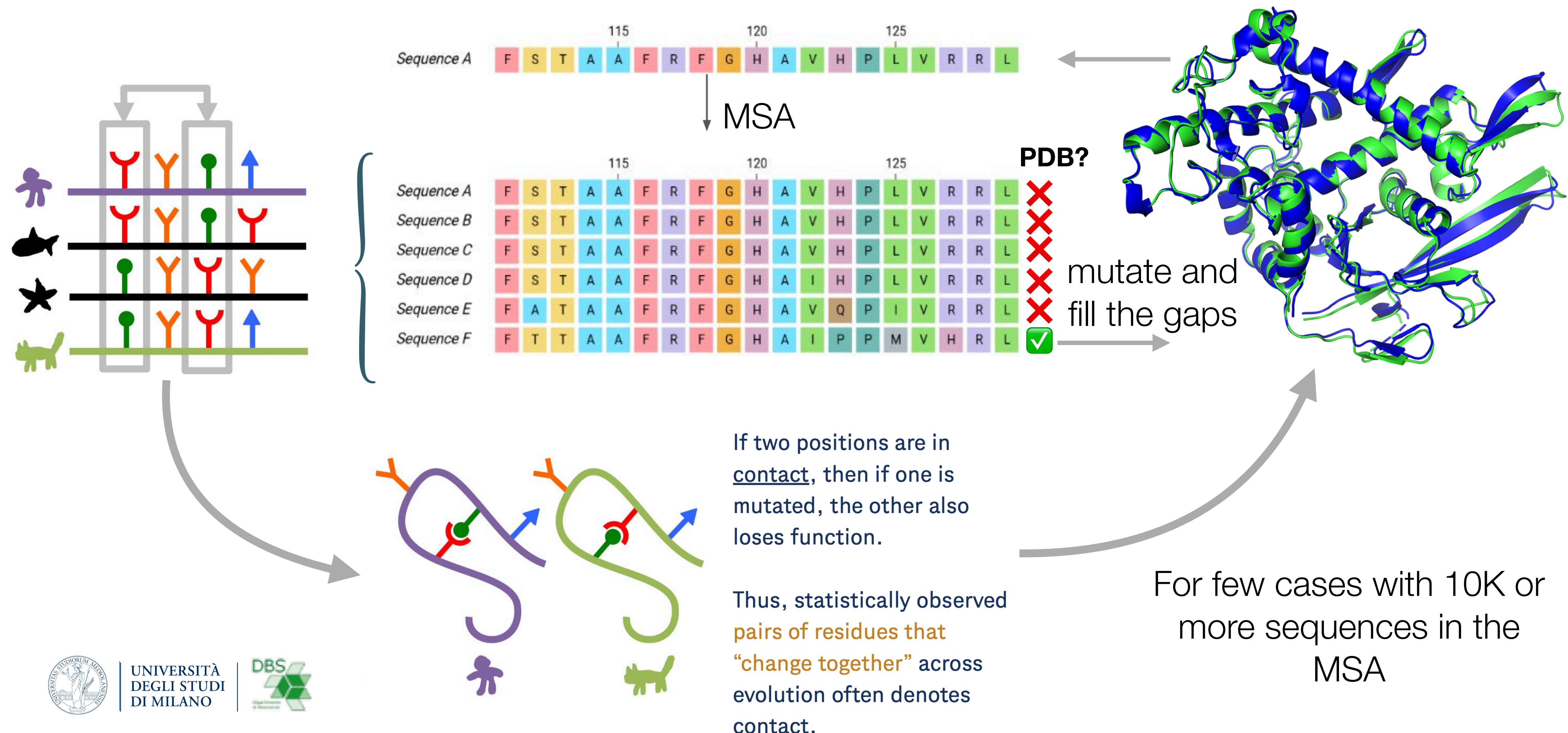
excluded volume

Dispersions (interactions of neutral molecules)

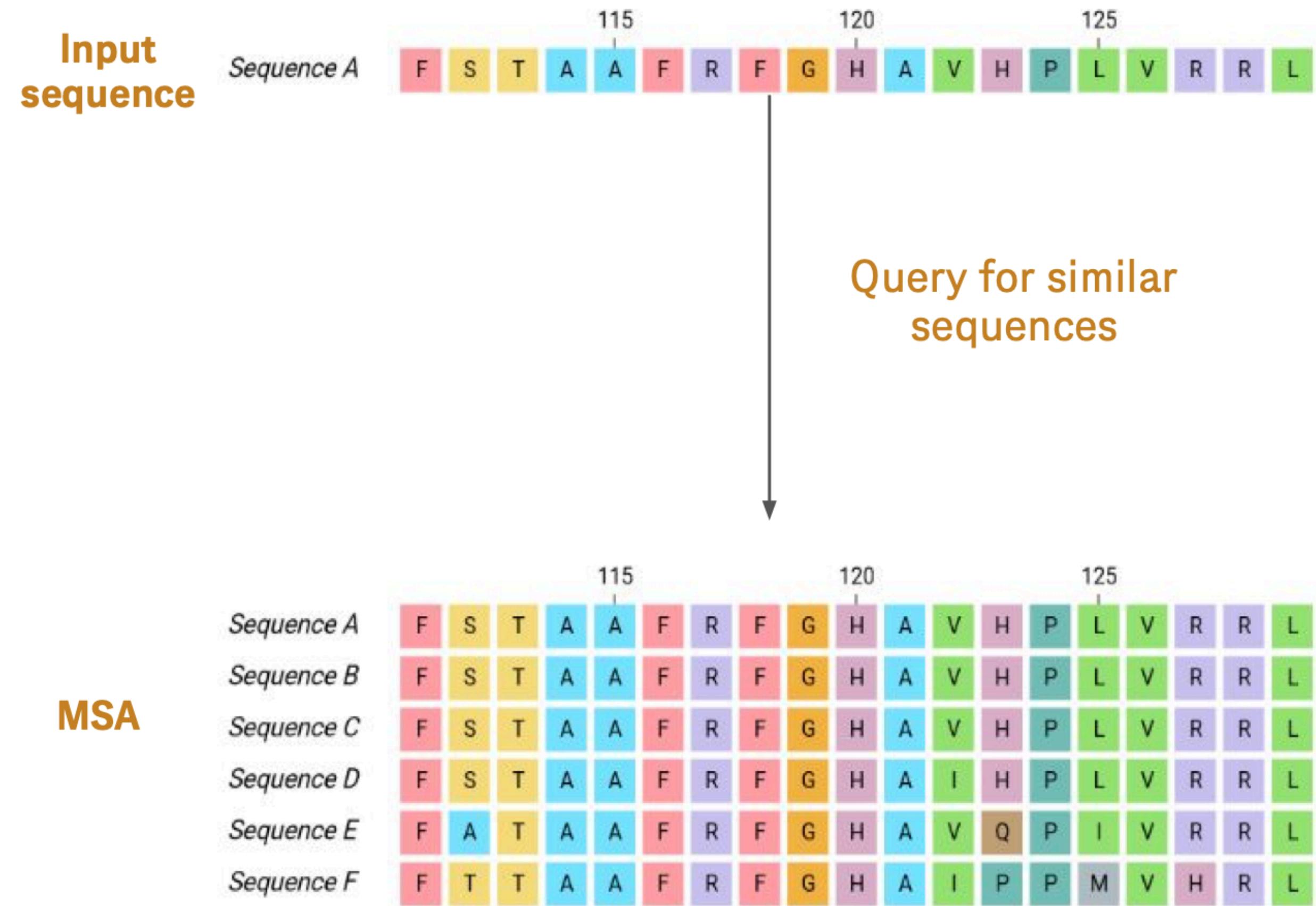


How fast-folding proteins fold. *Sci New York N Y* 334, 517–20 (2011).

# Structure prediction: homology modelling and coevolution analysis



# Search for patterns in protein sequences:



From an MSA, assuming that similar sequences are evolutionary related, we can learn a lot:

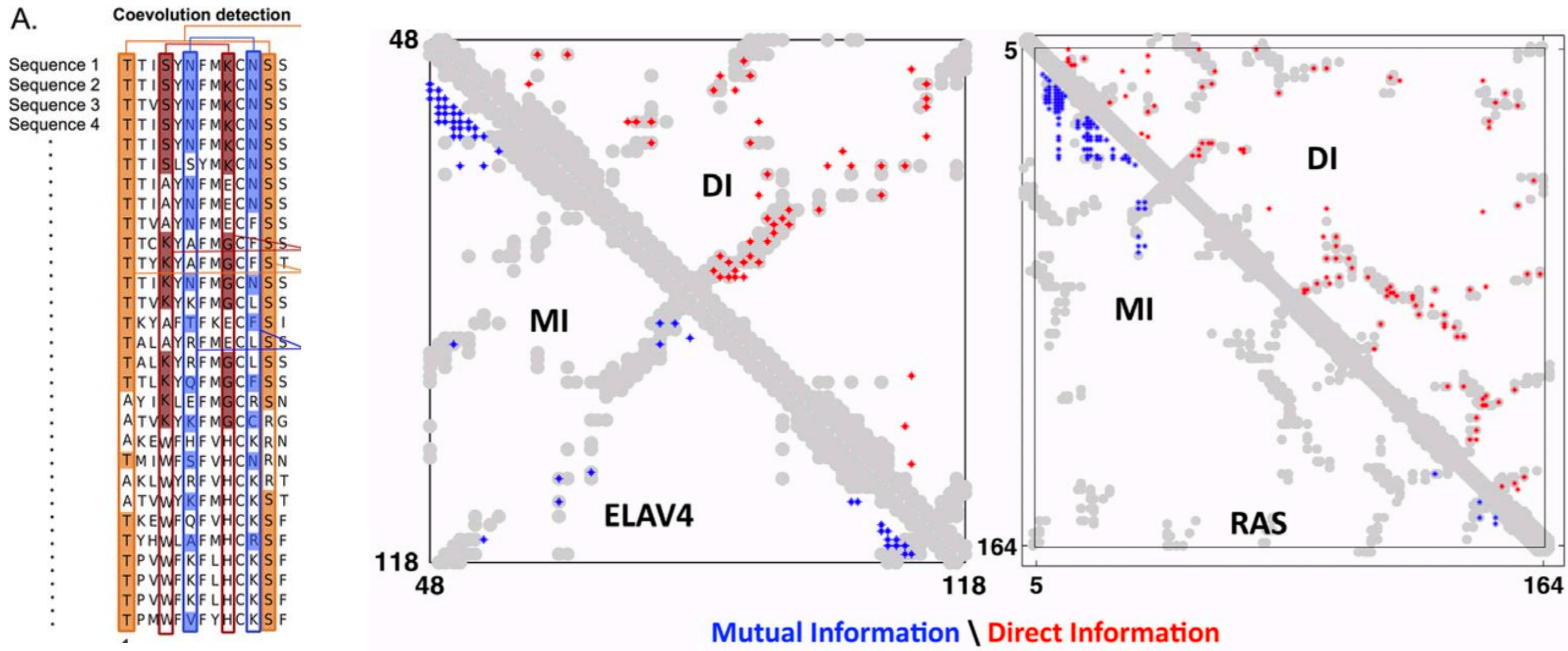
1. Sequence conservation: The frequency of each aminoacid in each position:  $P(F1)$ ,  $P(S2)$ , ...

2. The frequency of pair of AA:  $P(F1, S2)$ ,  $P(F1, T3)$ ,  $P(F1, A4)$ , ...

In principle we could get probability of triplets, etc, but already for pairs we need huge MSA.



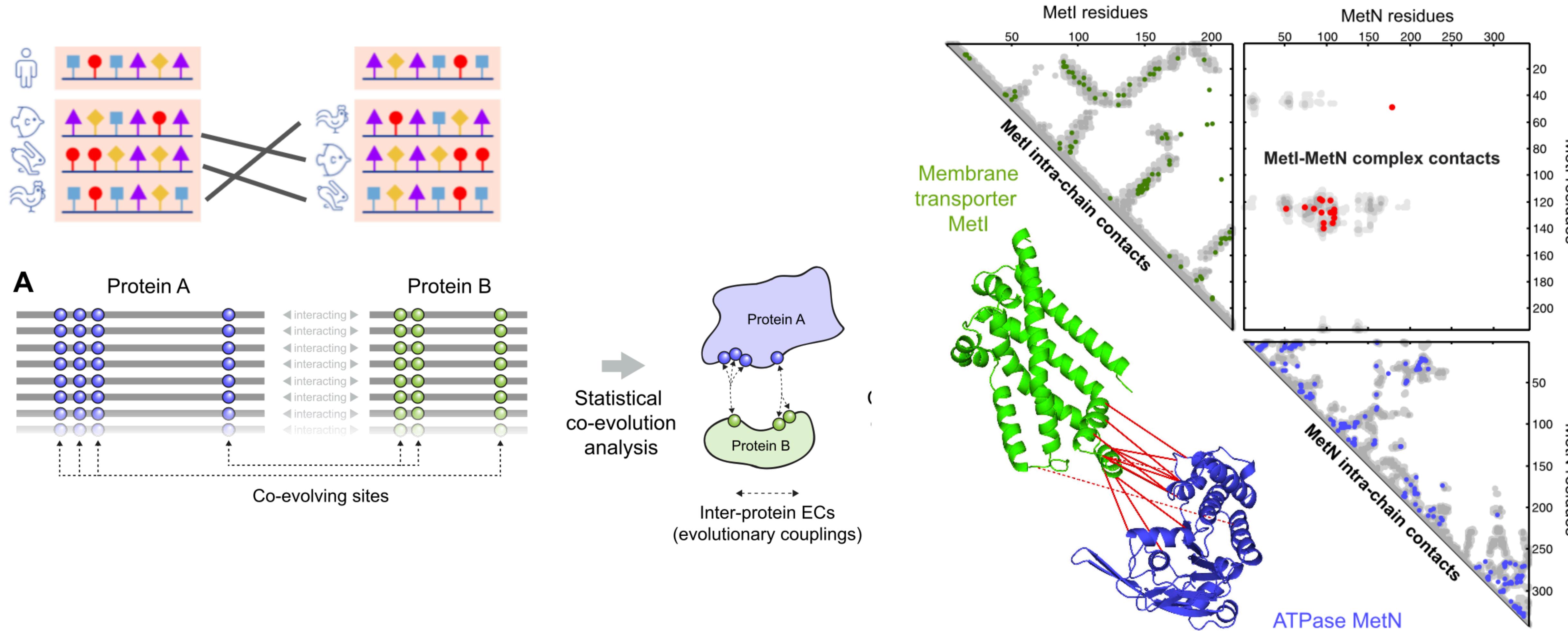
# A properly performed Coevolution analysis from LARGE MSA can provide reliable structural data



UNIVERSITÀ  
DEGLI STUDI  
DI MILANO



# Coevolution also works for protein complexes



# Summarising:

---

Chemistry does not enforce sequence-structure similarity in proteins.

"Our" proteins are the result of natural evolution. Sequences have generally evolved under the constraint of preserving/repurposing/improving function. Function can be associated with structure. Thus, in our experience, there is a sequence-structure similarity relationship.

MSA is a way to analyze sequences and identify homology. From MSA it is possible to find patterns related to the structure/function of a system.

Keep in mind that 30% of proteins/protein regions are unstructured, so the above assumption does not hold because function is not associated with structure. Evolution worked differently there.



# From COEVOLUTION to MACHINE LEARNING

---

Instead than solving the coevolution analysis case by case, that is very expensive and requires large MSA, would it be possible to use Machine Learning on many MSA and structures to train a single, reusable, network that is able to extract the structural information present in a MSA also when limited number of sequences are available?



# Not all ML are the same

The architecture is the creative, trial-and-error, part. Not all architecture are equally good at learning from a limited amount of data. Architecture design imply making strong hypothesis on how we want to interpret the data.

The weights are obtained for a given architecture and are the result of a more or less costly computer optimisation.



To train an architecture, that means to find the weights, you need a SCORING function that is a measure in some way of the distance of the output from the Training set.

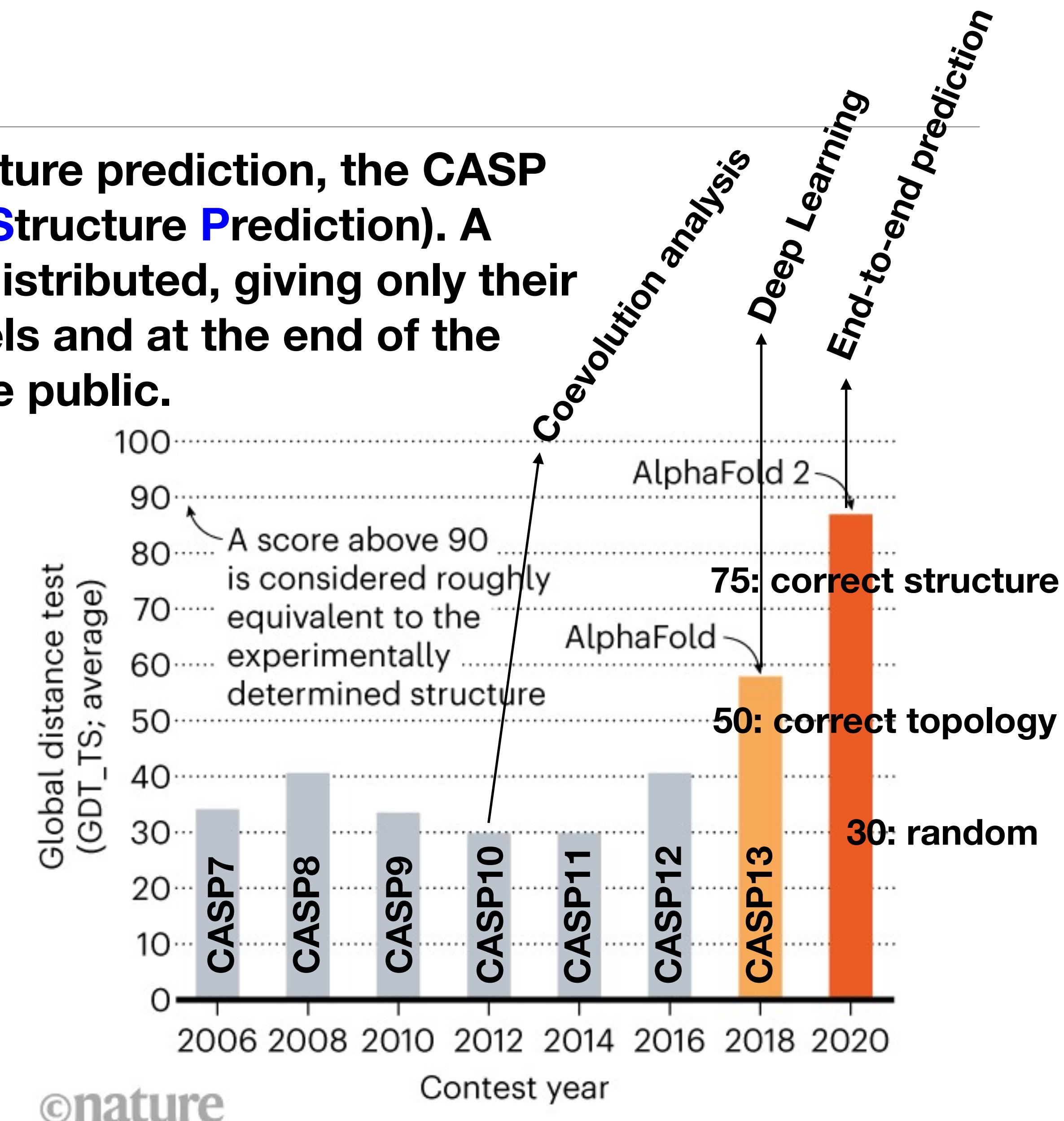
It could be something as simple as  $\text{sum}[(\text{output}-\text{training})^2]$

# Ab-initio protein structure predictions

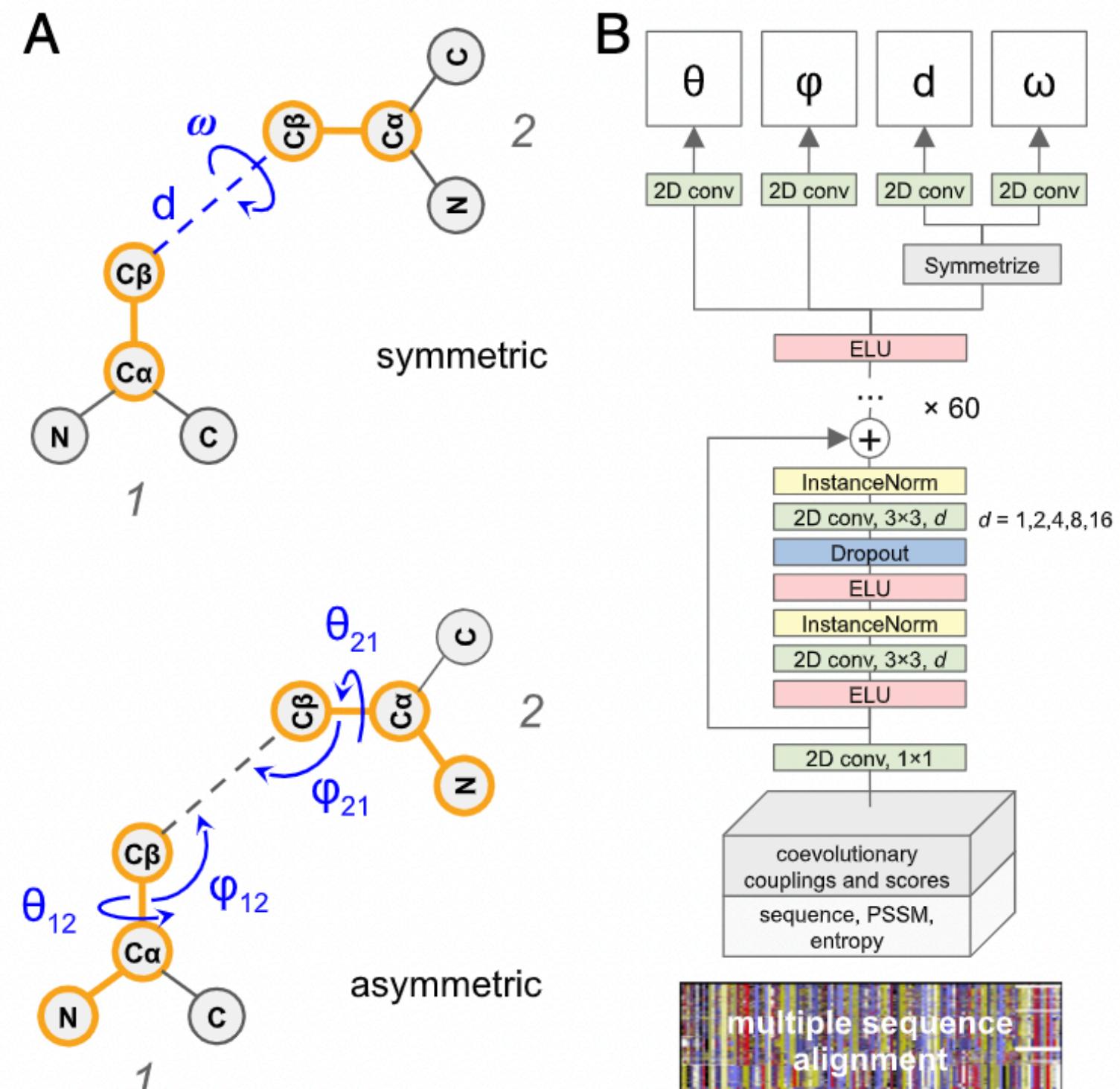
From 1994 there is a competition for protein structure prediction, the CASP (Critical Assessment of techniques for protein Structure Prediction). A number of structures recently determined are not distributed, giving only their sequences. Participants can deposit their models and at the end of the competition the results are made public.

From the early 2000 the way to predict the structure of a protein has been to follow a “fragment replacement” strategy:

1. Generate many configuration using local information like secondary structure prediction as well as contact restraint using physico-chemical principles.
2. Refine the final structure at atomic resolution and use some fancy “scoring function”



# CASP13 - AlphaFold, RaptorX and TrRosetta: deep-learning distances and orientations from coevolution features



~15,000 proteins structures PDB  
Histograms for the distances among all couples of amino acids (distance distributions):

ALA-ALA: distances from all couples in all selected PDBs

ALA-CYS:  
ALA-ASP:....

All sequences are also associated to their relative MSA derived covariance matrix

A deep neural network is trained to go from the covariance matrix to the distance distribution

**NOTE:** we move from contact, binary, maps to distance distribution maps, this allow doing direct minimisation of the structure using the propagation property of the network

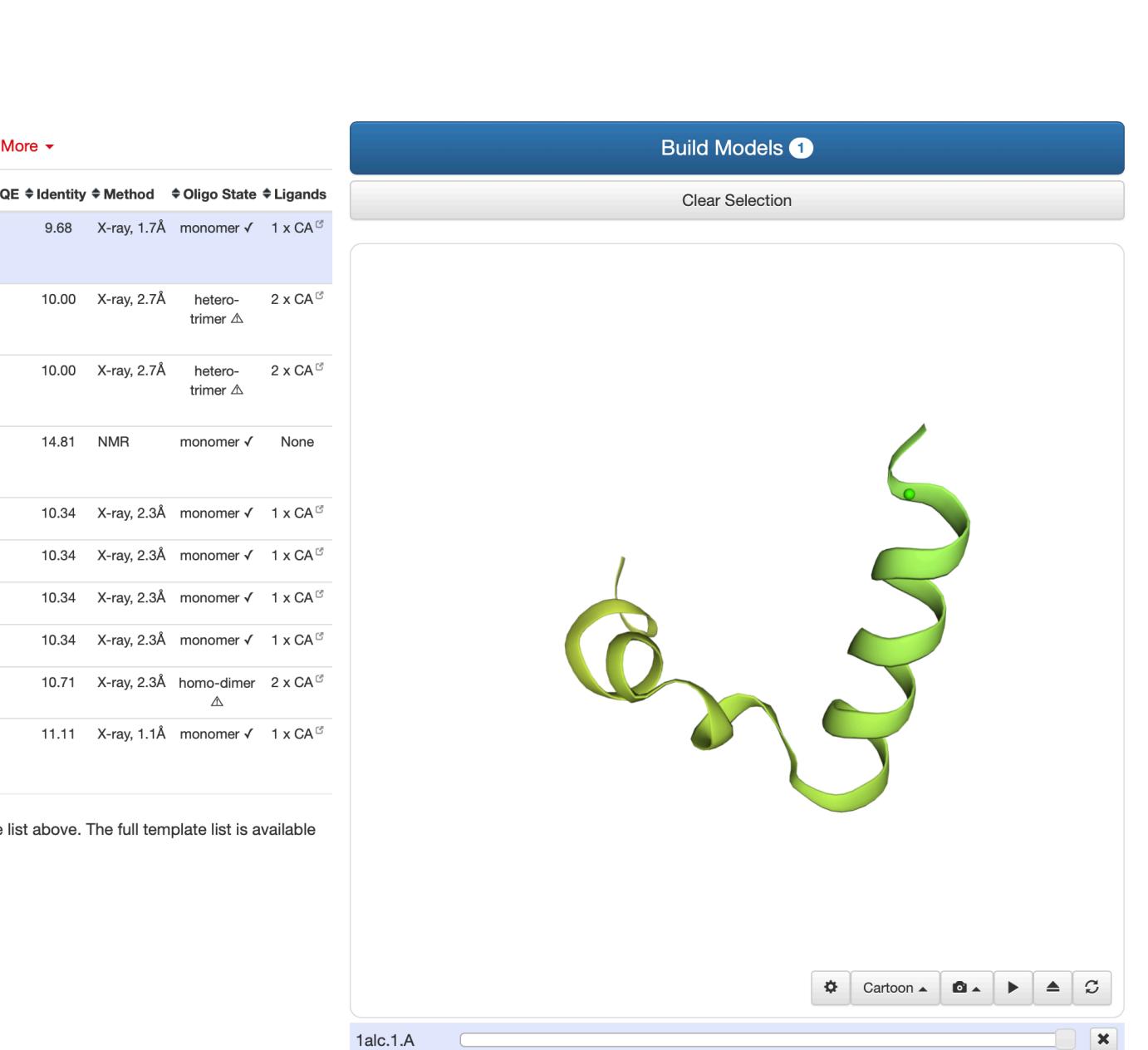


# MIZ1: an example from my own experience

**MIZ1 is a protein that is studied here in the department for its functional role in plants were it plays a key role in root development. Its structure is unknown and we wanted to determine it, or at least to help in design a construct for further experiments, because the full length protein gets into the inclusion bodies in E. Coli.**

MVPYQELTLQRSFSYNSRKINPVTSPARSSHVRSPSSALIPSIPEHEFLVPCRRCSYV  
PLSSSSSASHNIGKFHLKFSLLRSFINIINIPACKMLSLPSPPPSSSSVSNQLISLVTG  
GSSSLGRRVTGTYGHKRKGHTFSVQYNQRSDPVLLLLAMSTATLVKEMSSGLVRIALE  
CEKRHRSGTKLFQEPKWTMYCNGRKCGYAVSRGGACTTDWRLNTSRVTVGAGVIPTP  
KTIDDVSGVGSGTELGELLYMRGKFERVVGSRDSEAFYMMNPDKNGGPELSIFLLRI

**First test - homology modelling - but there are no homologues with known structure:**

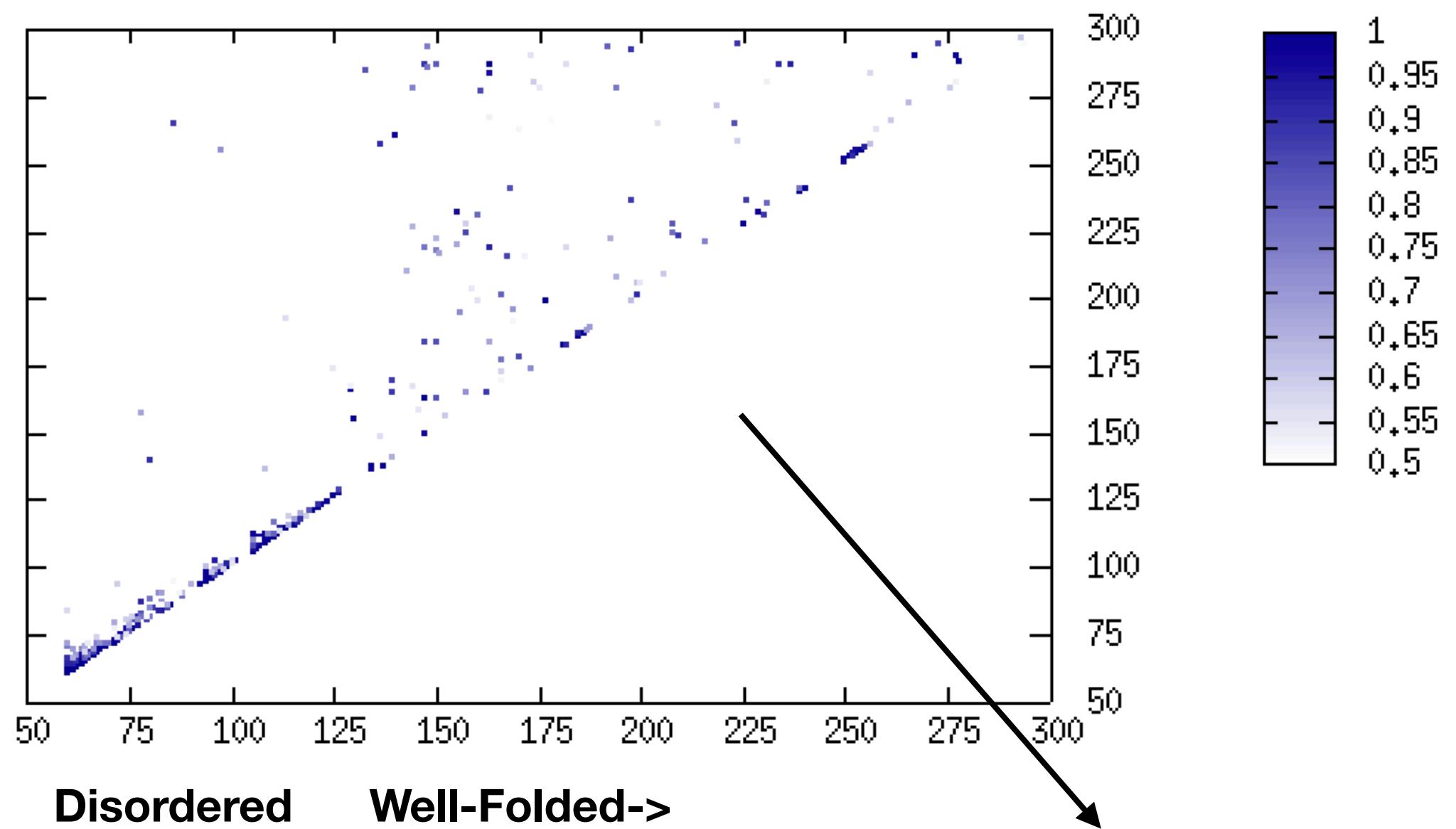


**NOPE!**



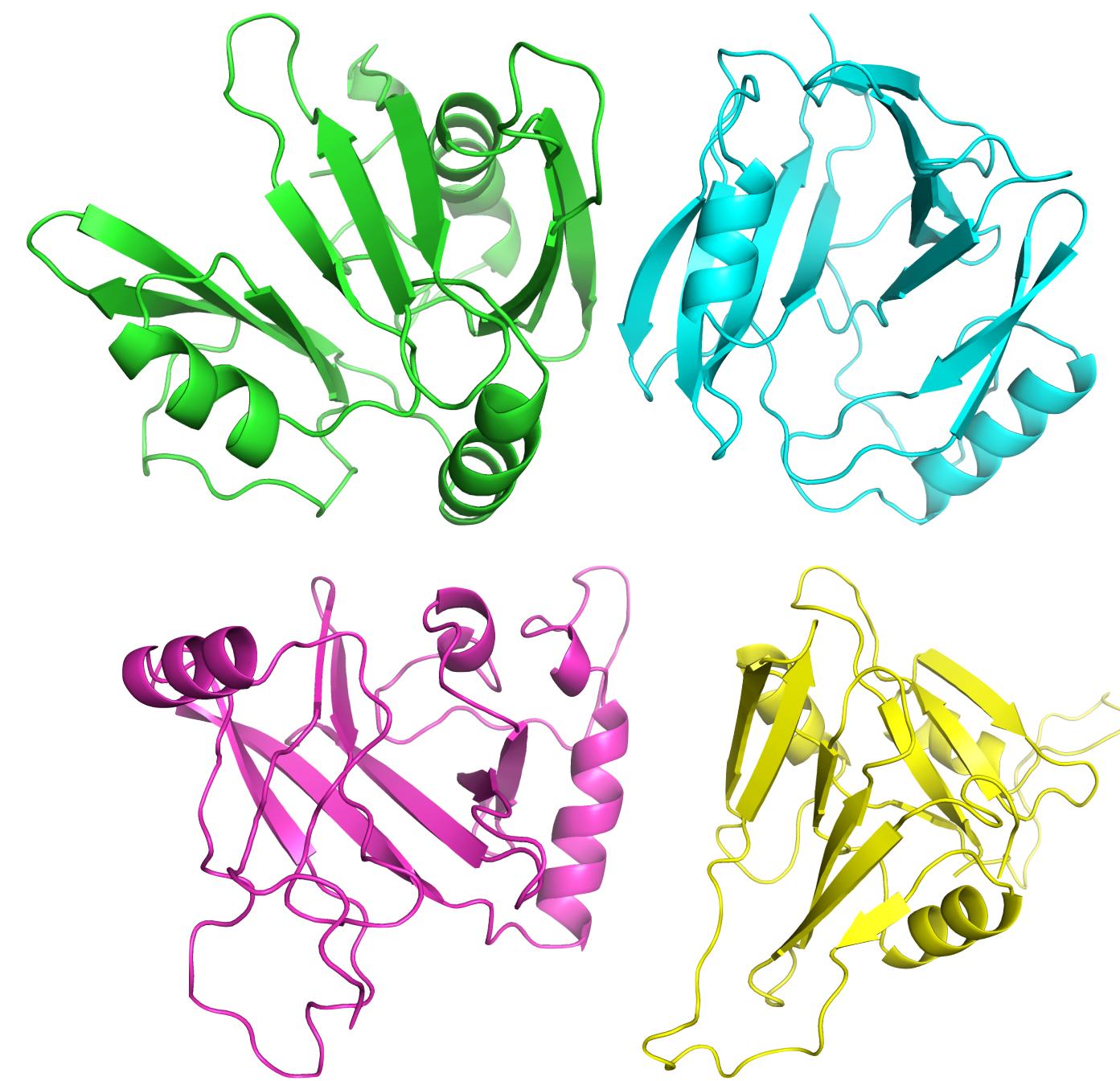
# MIZ1: an example from my own experience

**MSA and coevolution analysis, we found ~700 homologues sequences, that is a small number for pre-AF methods.**



**Yet, we used the map to design a construct to produce the protein that actually worked well**

**Robetta + ev-couplings - ‘old’ ab-initio structure prediction:**



**Four structures with nothing in common... nope.**

# MIZ1: an example from my own experience

RaptorX

(first time of ev-couplings+deep learning):

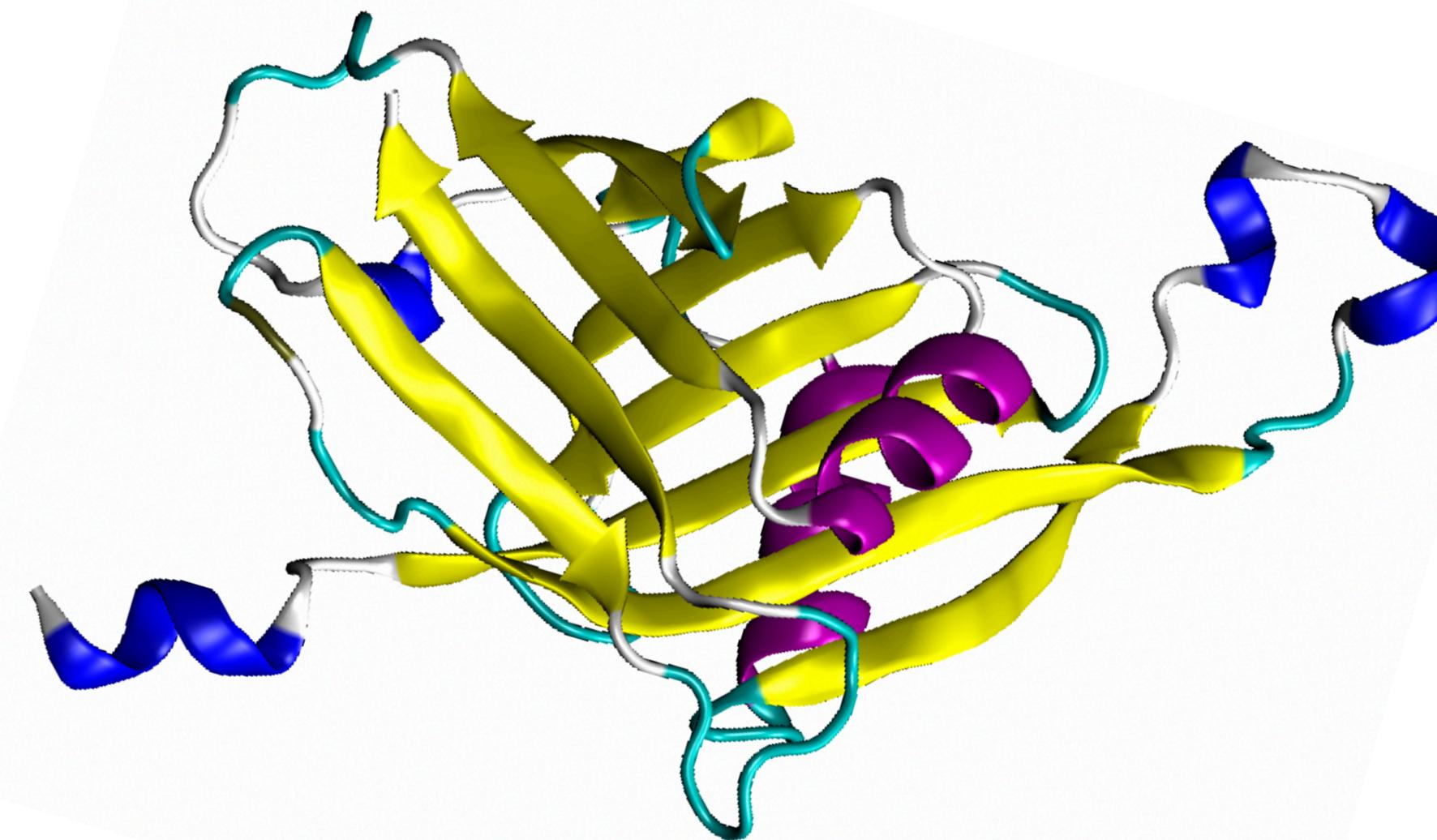


Five structures with something in common,  
but when looked in the details they are  
unphysical

**NOTE 1:** this is a case with ~700 hundred homologue sequences, this is still a relatively large number

trRosetta

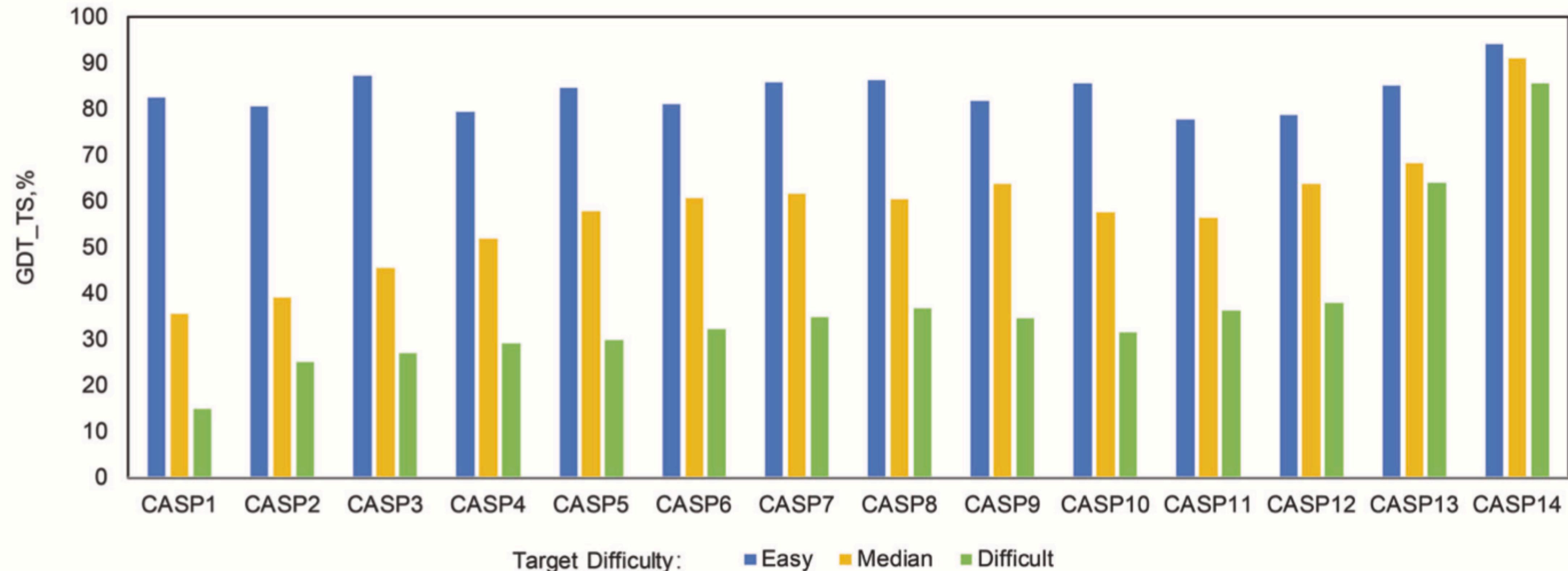
(post AF ev-couplings+deep learning):



Well converged, well-defined unique structural model with perfectly reasonable features.

**NOTE 2:** I am talking about trRosetta and not AF, because the first AlphaFold as never been made publicly available

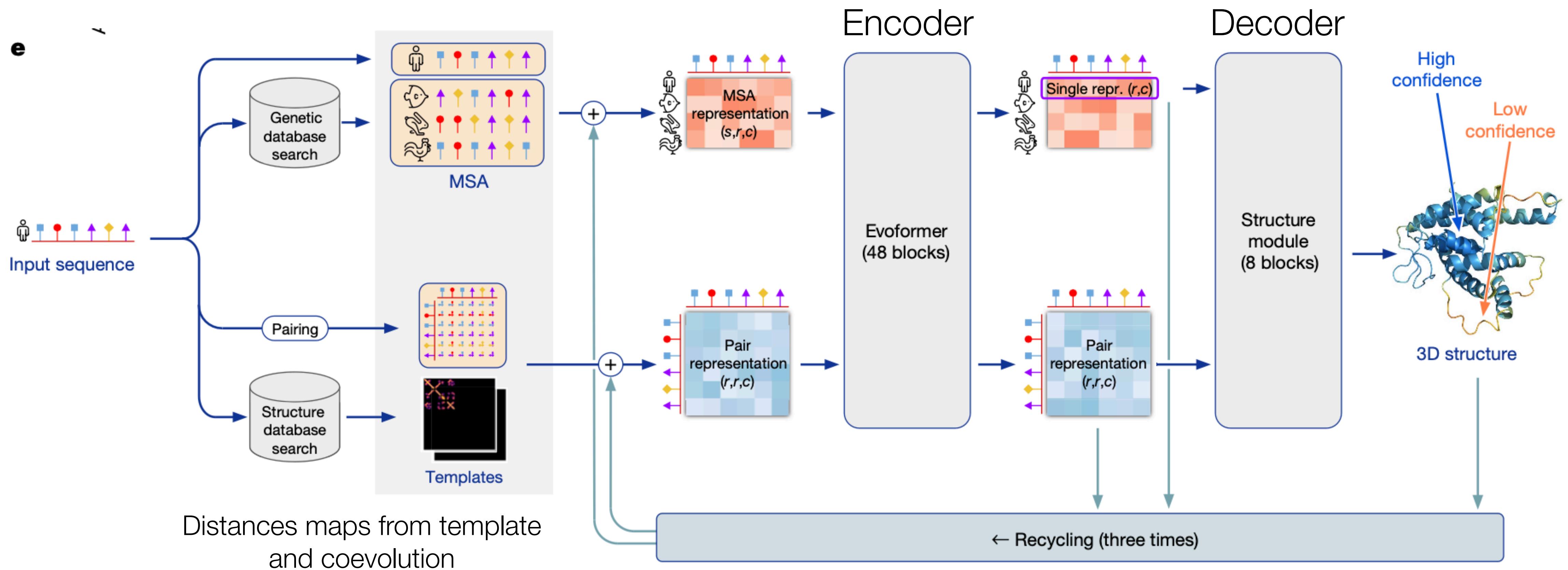
# CASP14 - AlphaFold2: not just a technical improvement but a completely ad hoc newly designed network architecture



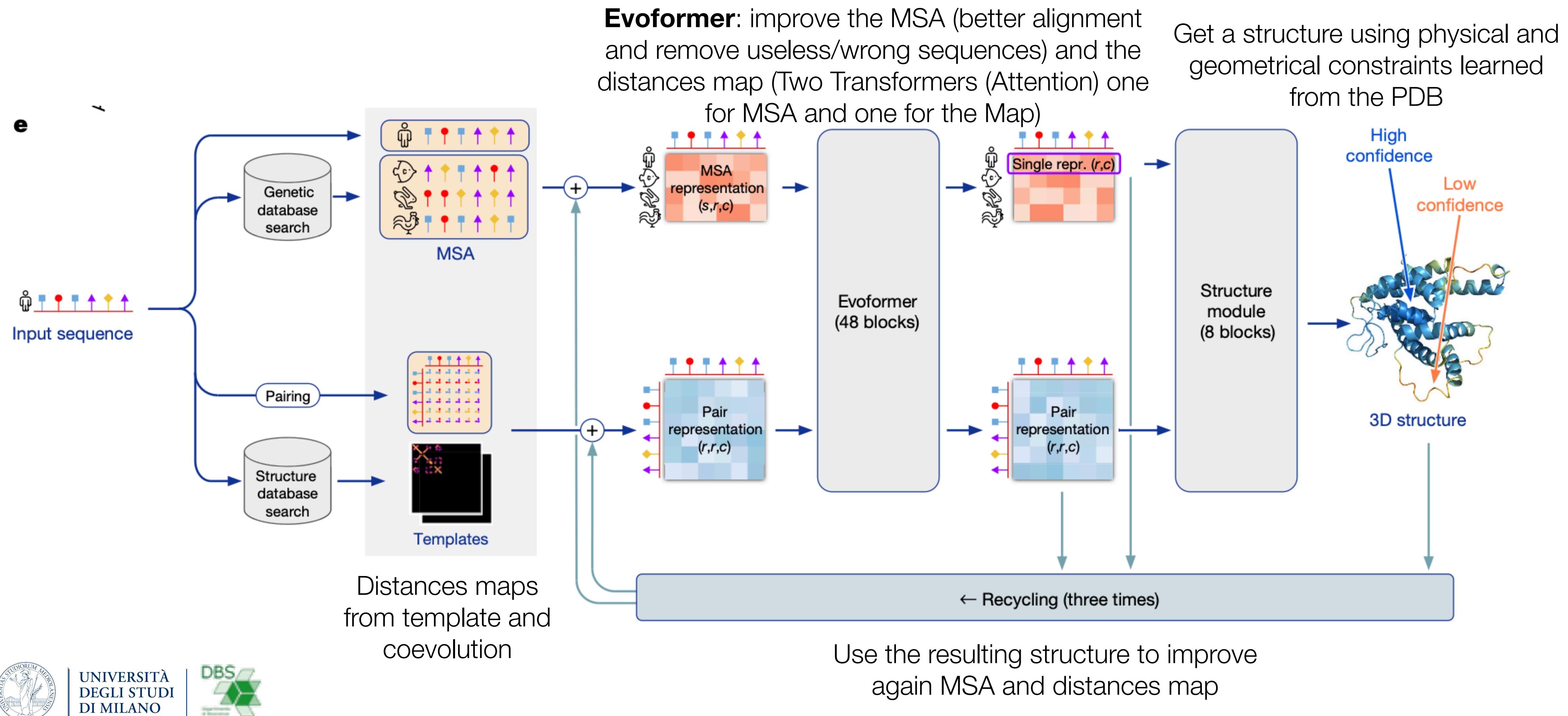
# AlphaFold2: some ideas on the architecture

In pre AF2 approaches the network allowed to increase the amount of information extracted from the MSA.

**AF2 first addition is to improve the quality of the MSA itself. The second addition is to have an end-to-end network** instead than using the network to generate the distance distribution map and then applying it to a conventional fragment replacement method.

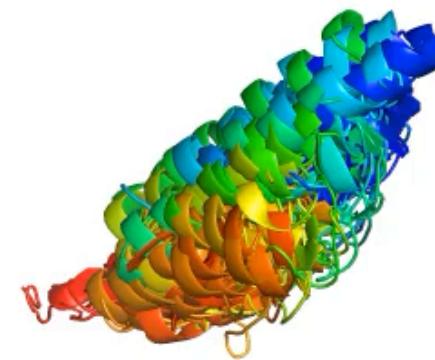


# AlphaFold2: some ideas on the architecture

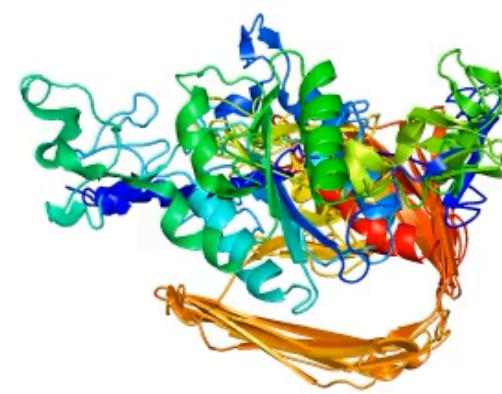


# Examples:

---



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



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



UNIVERSITÀ  
DEGLI STUDI  
DI MILANO

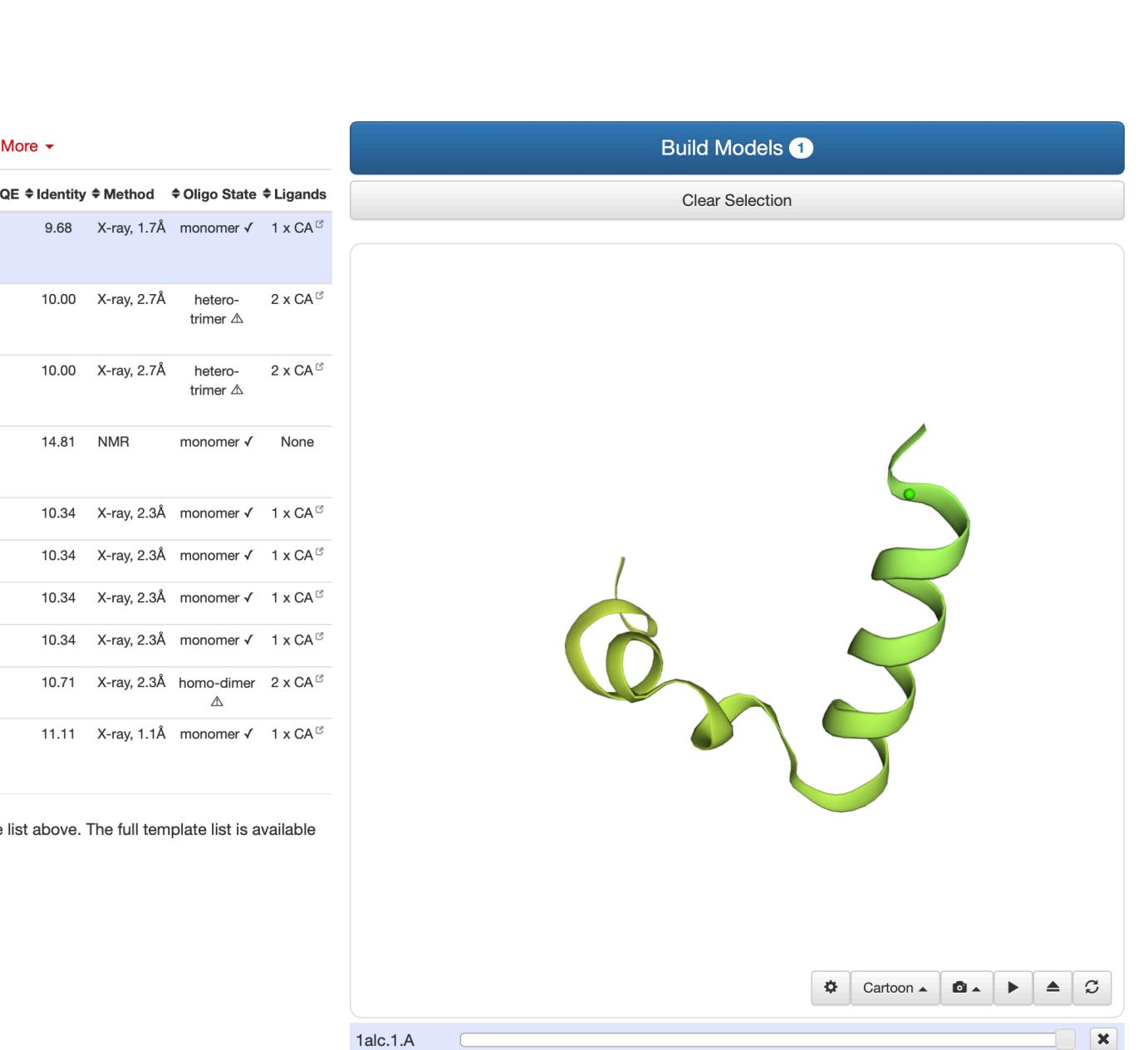


# MIZ1: an example from my own experience

**MIZ1 is a protein that is studied here in the department for its functional role in plants were it plays a key role in root development. Its structure is unknown and we wanted to determine it, or at least to help in design a construct for further experiments, because the full length protein gets into the inclusion bodies in E. Coli.**

MVPYQELTLQRSFSYNSRKINPVTSPARSSHVRSPSSALIPSIPEHEFLVPCRRCSYV  
PLSSSSSASHNIGKFHLKFSLLRSFINIINIPACKMLSLPSPPPSSSSVSNQLISLVTG  
GSSSLGRRVTGTYGHKRKGHTFSVQYNQRSDPVLLLLAMSTATLVKEMSSGLVRIALE  
CEKRHRSGTKLFQEPKWTMYCNGRKCGYAVSRGGACTTDWRLNTSRVTVGAGVIPTP  
KTIDDVSGVGSGTELGELLYMRGKFERVVGSRDSEAFYMMNPDKNGGPELSIFLLRI

**First test - homology modelling - but there are no homologues with known structure:**

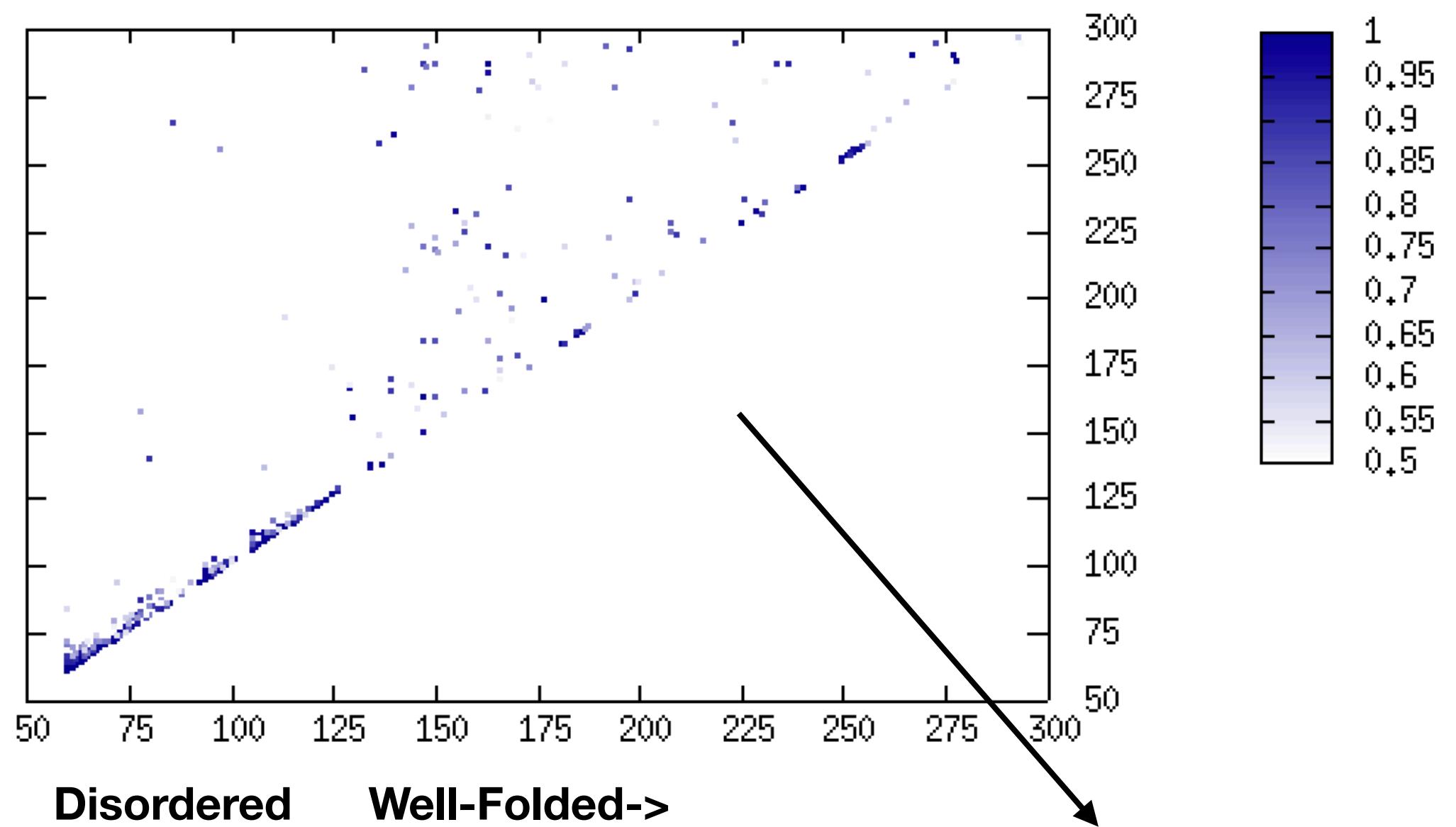


**NOPE!**



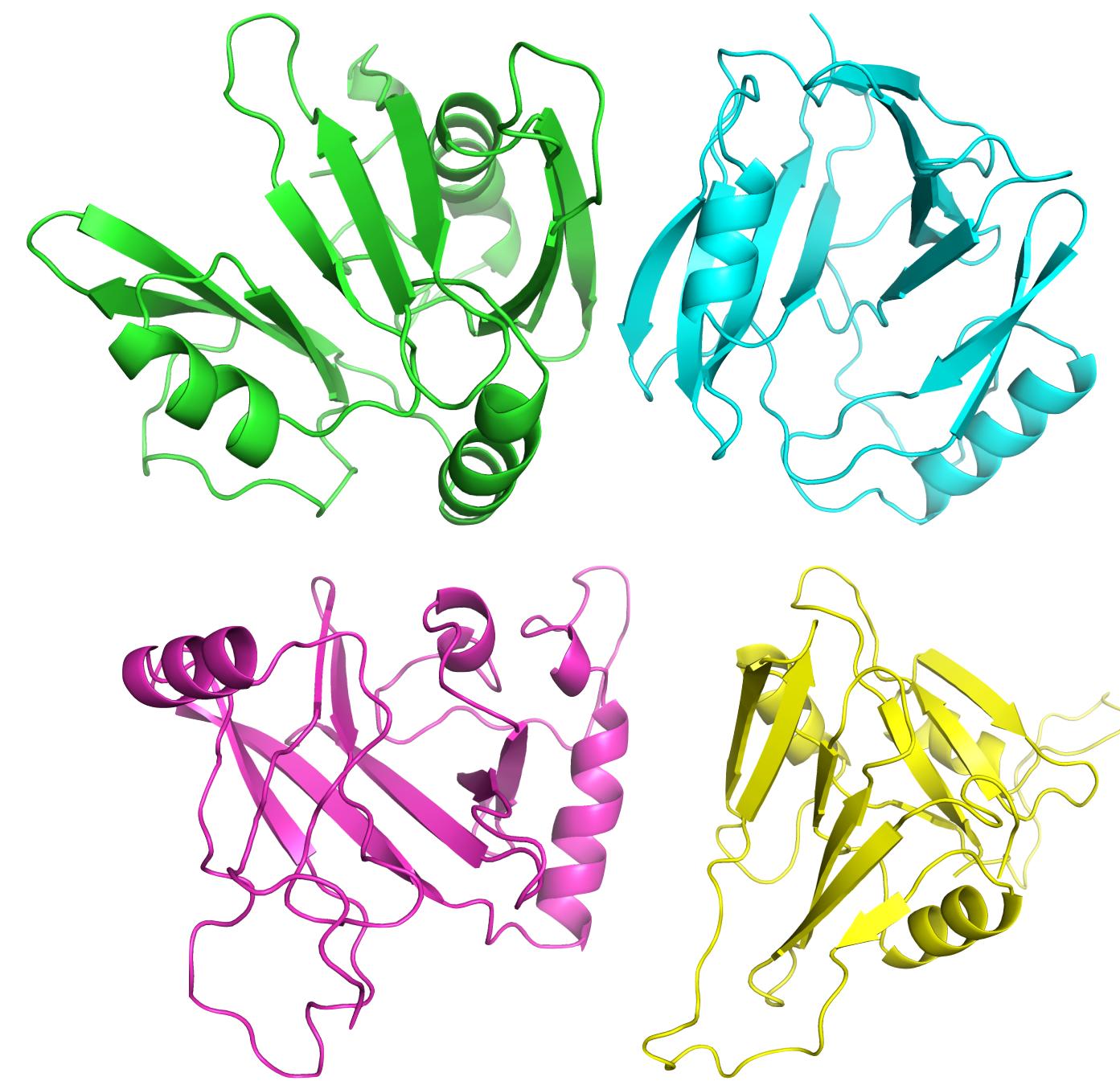
# MIZ1: an example from my own experience

**MSA and coevolution analysis, we found ~700 homologues sequences, that is a small number for pre-AF methods.**



**Yet, we used the map to design a construct to produce the protein that actually worked well**

**Robetta + ev-couplings - 'old' ab-initio structure prediction:**



**Four structures with nothing in common... nope.**

# MIZ1: an example from my own experience

RaptorX

(first time of ev-couplings+deep learning):

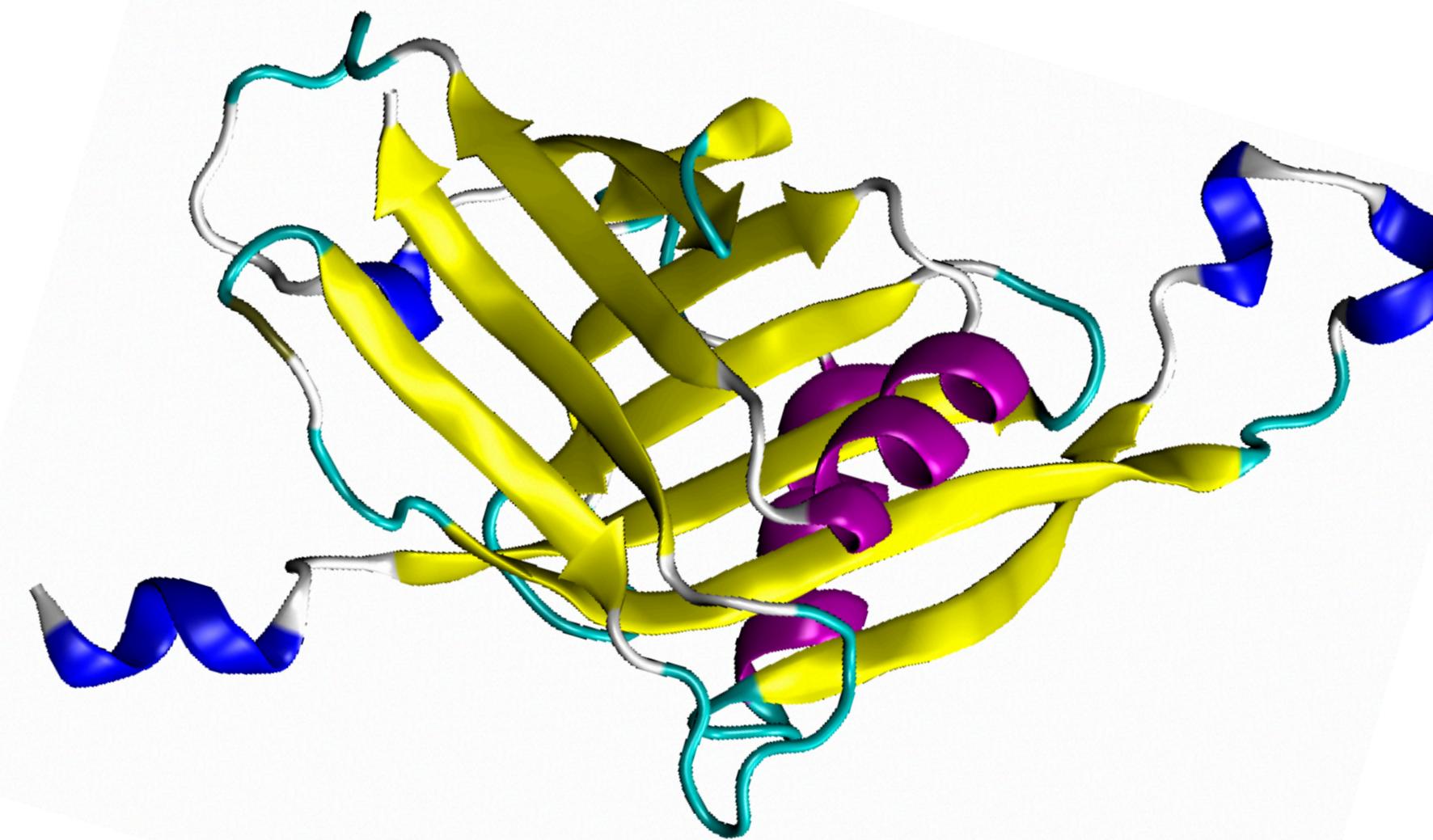


**Five structures with something in common,  
but when looked in the details they are  
unphysical**

**NOTE 1:** this is a case with ~700 hundred homologue sequences, this is still a relatively large number

trRosetta

(post AF ev-couplings+deep learning):



**Well converged, well-defined unique structural model with perfectly reasonable features.**

**NOTE 2:** I am talking about trRosetta and not AF, because the first AlphaFold as never been made publicly available

# AlphaFold Protein Structure Database

Developed by Google DeepMind and EMBL-EBI

Search for protein, gene, UniProt accession or organism or sequence search

BETA

**Search**

Examples: MENFQKVEKIGEGTYGV...

Free fatty acid receptor 2

At1g58602

Q5VSL9

E. coli

See search help 

Go to online course 

See our updates – September 2024

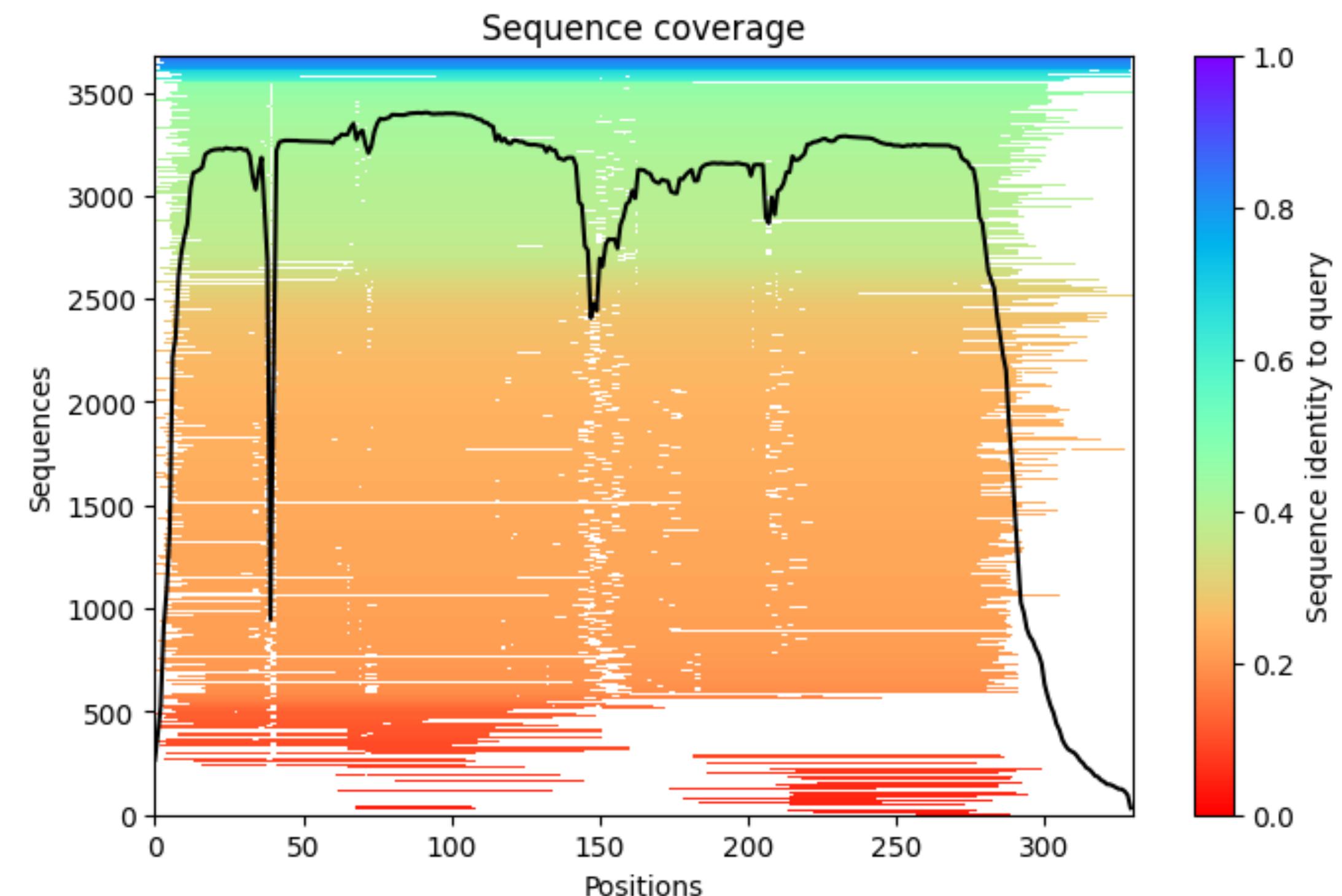
AlphaFold DB provides open access to over 200 million protein structure predictions to accelerate scientific research.

# Interpreting an AF2 prediction: from COLABFOLD

It is possible to obtain AF2 predictions either precalculated from the AFDB or running from COLABFOLD (and references therein)

The first data obtained are about the MSA: in this case I have many many sequences mostly with high identity and coverage: very good.

As an example I took from uniprot “A0A2K5XT84” the free fatty acid receptor 2 of “*mandrillus leucophaeus*” and I ran it on COLABFOLD



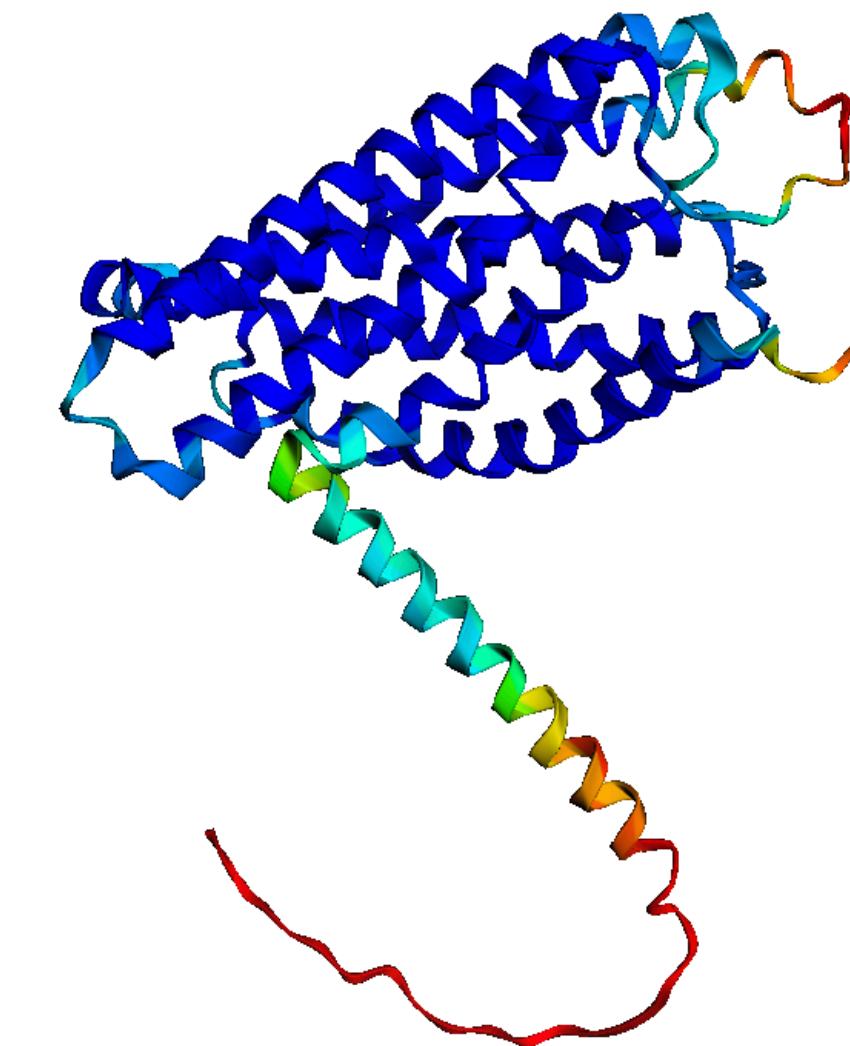
UNIVERSITÀ  
DEGLI STUDI  
DI MILANO



# Interpreting an AF2 prediction: from COLABFOLD

Once the run is done I look at the results: all the predicted structures look the same. This is a very good indication about prediction robustness, if AF2 is not convinced about the prediction why I should be?

The I can look at the **PLDDT** (that is used to colour the strucutre). This the **KEY INDICATOR** of prediction quality for **TERTIARY STRUCTURE**.



■ Very low (<50) ■ Low (60) ■ OK (70) ■ Confident (80) ■ Very high (>90)



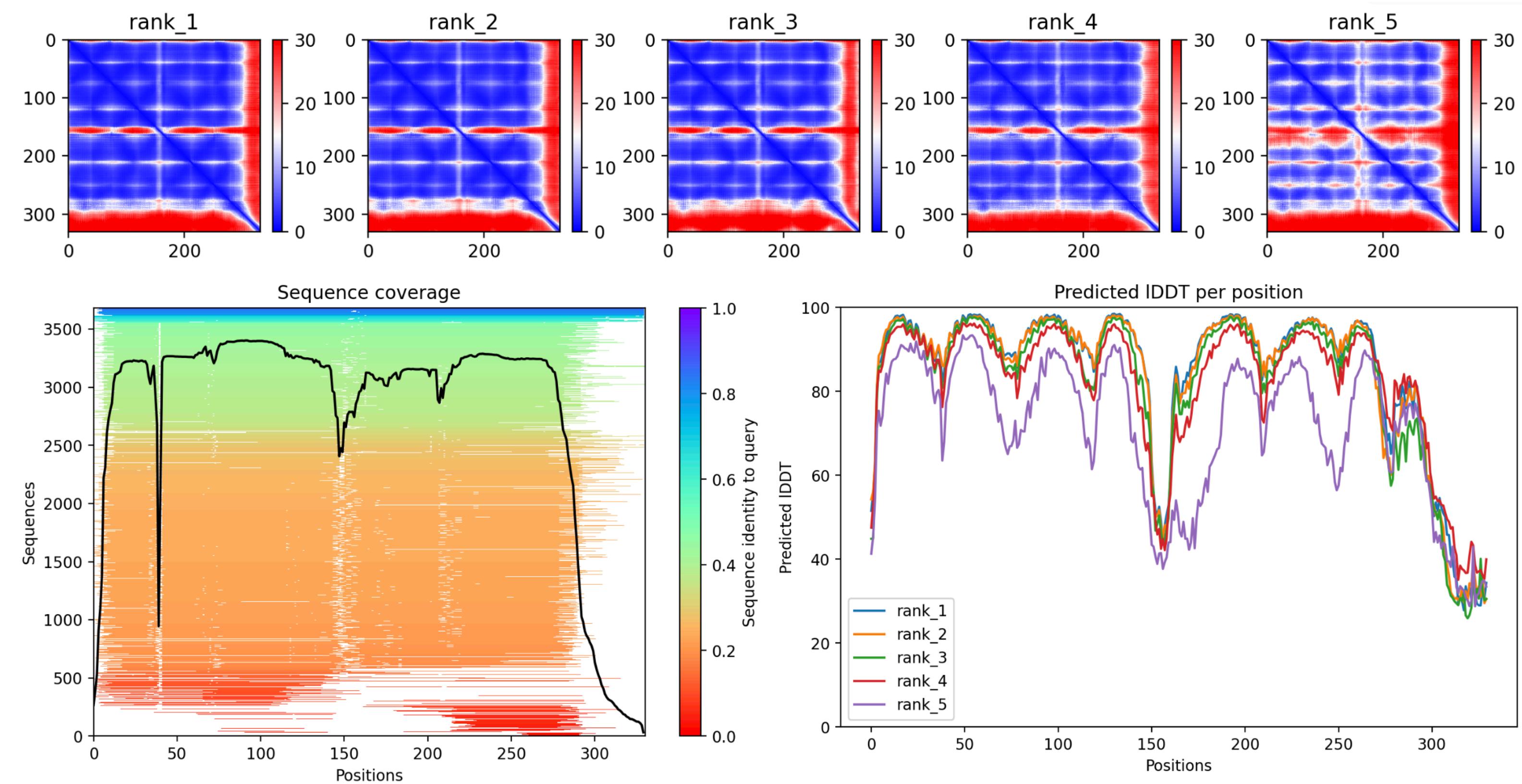
UNIVERSITÀ  
DEGLI STUDI  
DI MILANO



# Interpreting an AF2 prediction: from COLABFOLD

On TOP we have the PAE indicator that is relevant for super tertiary/quaternary structure and does report about the quality of relative organisation of tertiary units.

pLDDt is correlated with MSA, but running a prediction is a random process, so one structure is less good than the other even in this optimal case

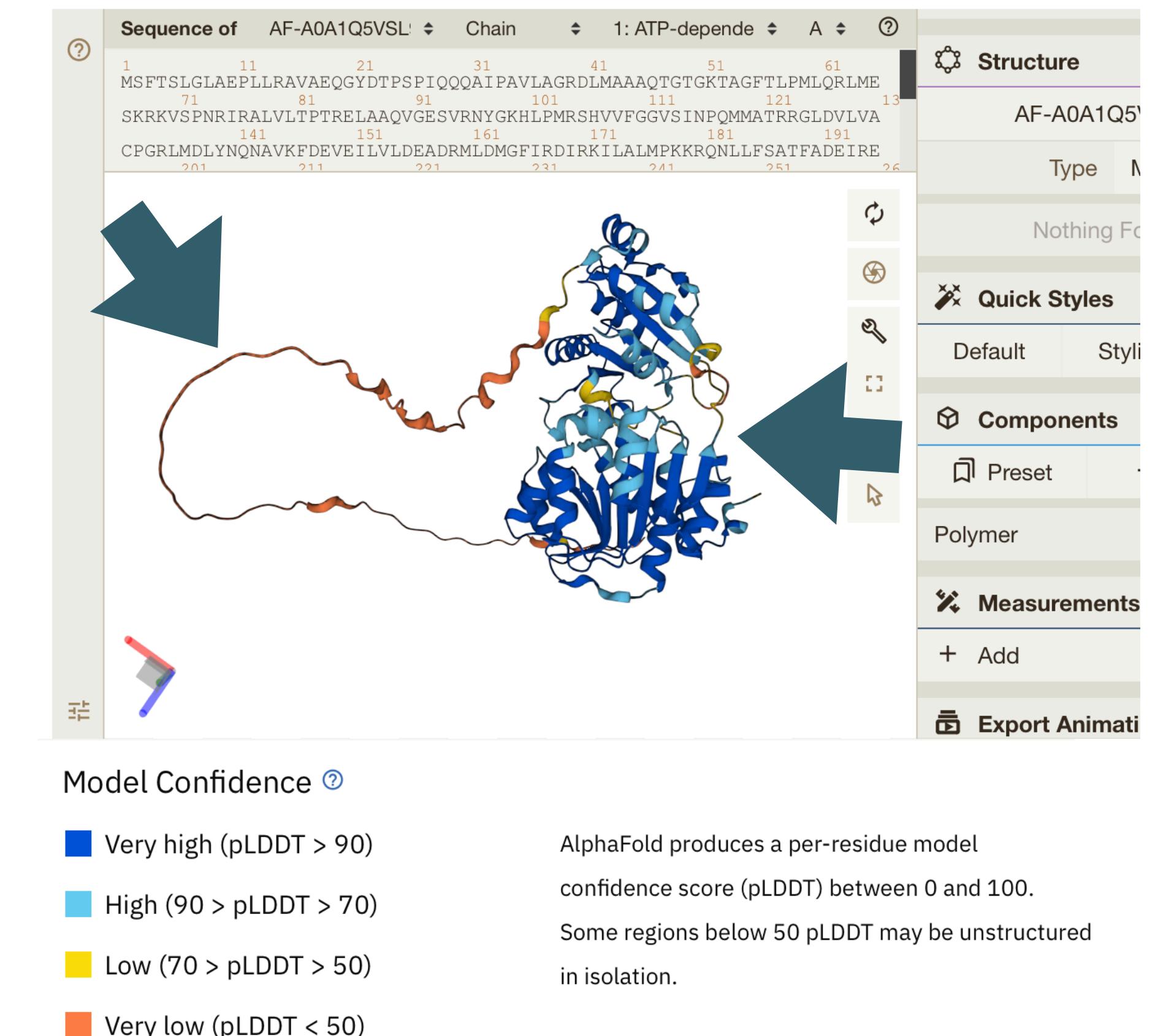


# Interpreting an AF2 prediction: from AFDB

AFDB does not provide informations about the MSA and the consistency of the prediction: so we need to FOCUS on **PLDDT** and PAE

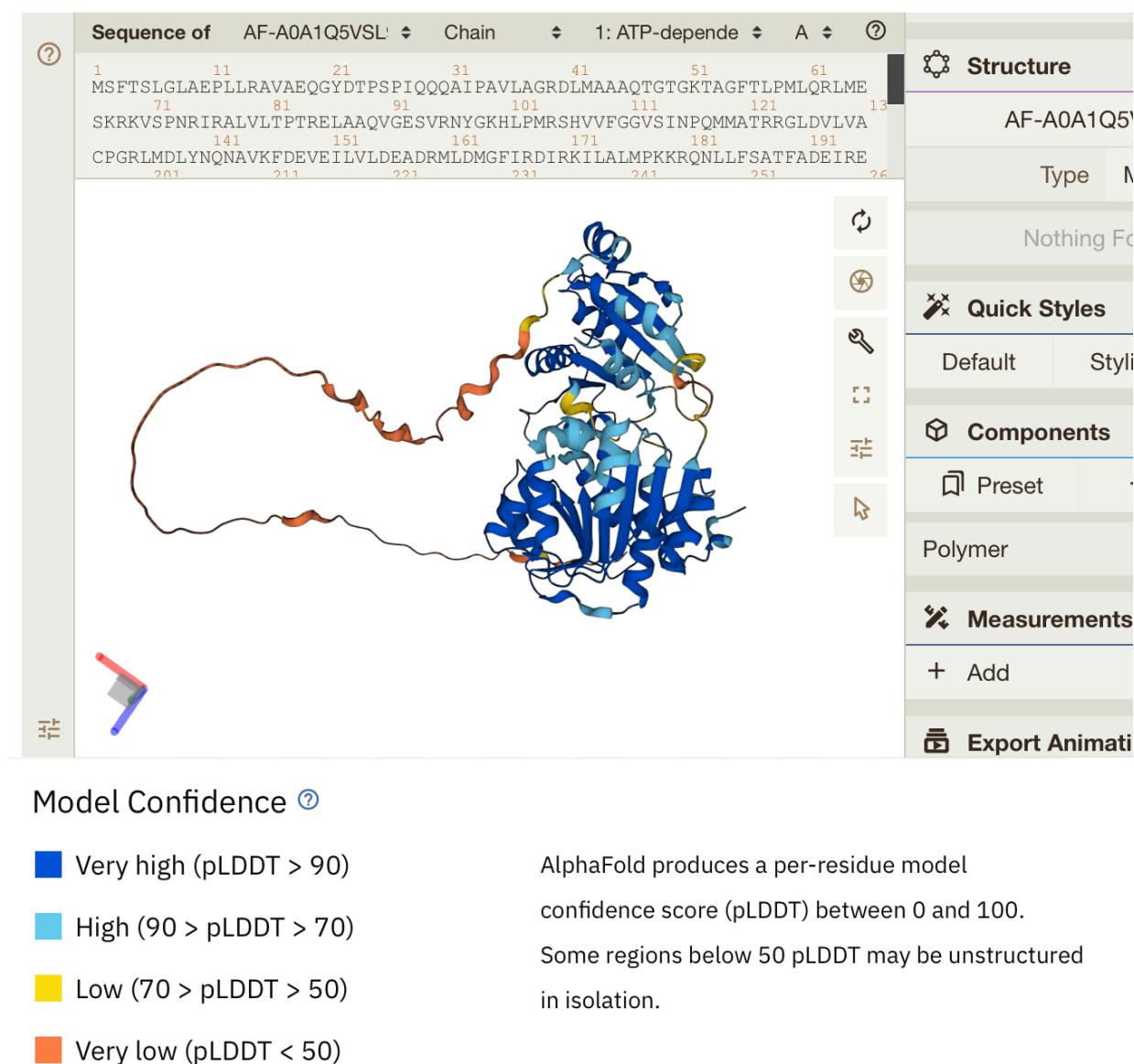
Lets' look at “Q5VSL9” from AFDB that is a prediction for the “ATP-dependent RNA helicase RhIE” from “*Aeromonas allosaccharophila*”:

**PLDDT is HIGH** everywhere but in the DISORDERED REGIONS



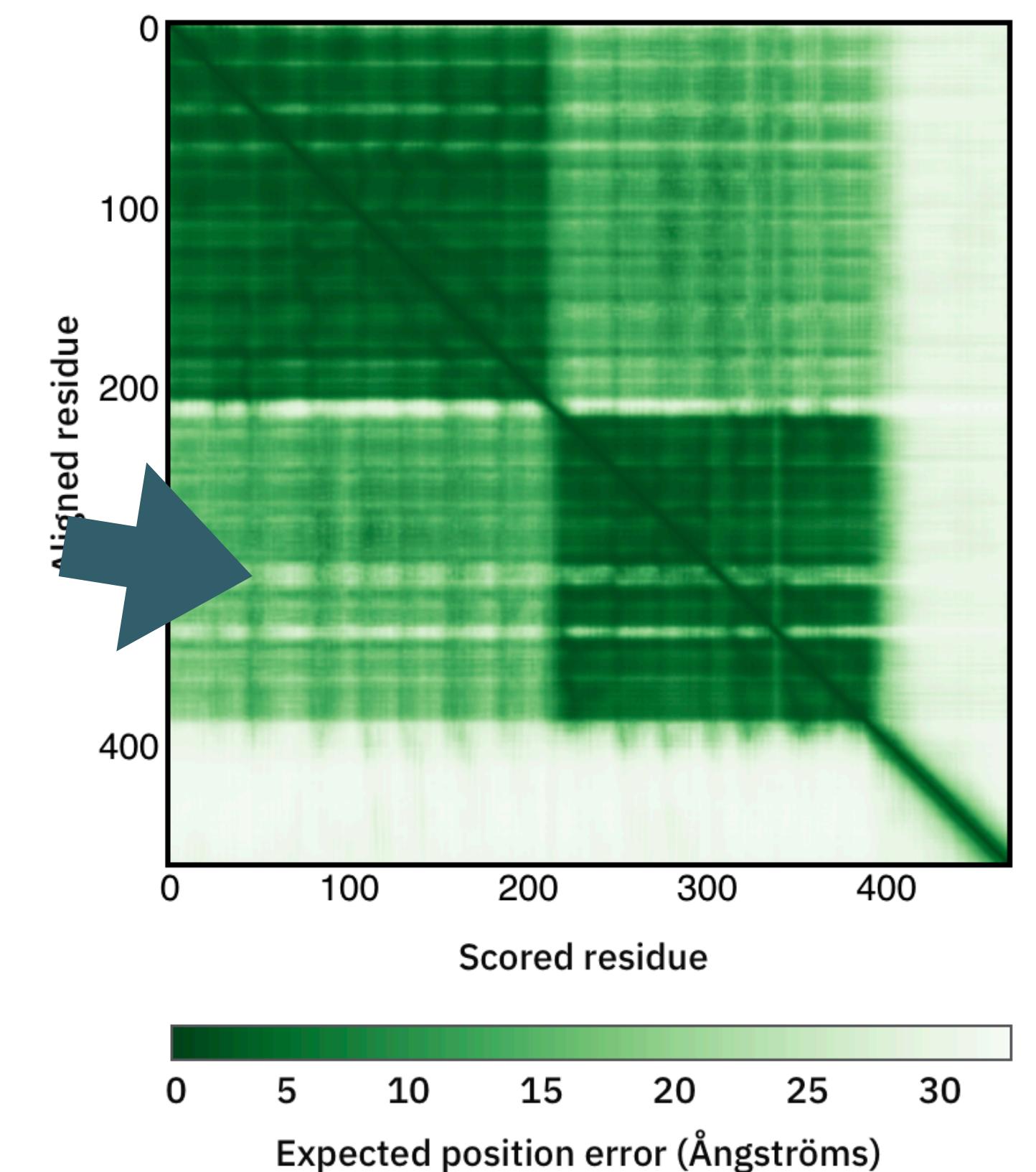
# Interpreting an AF2 prediction: from AFDB

The low **PLDDT** in the region linking the two domains is reflected in the lower confidence reported for their relative orientation



**HINT:** when a protein has long disordered regions is better to cut them out from the prediction. Long stretch of IDP regions hamper the prediction.

Predicted aligned error (PAE)



UNIVERSITÀ  
DEGLI STUDI  
DI MILANO



# Summary

---

AF2 and similar methods can provide PREDICTIONS of structure and CONFIDENCE metrics. The latter are very important to evaluate the quality (local and global) of the prediction. Even for structure in AFDB, you can run it again to be sure and get more INFO.

DISORDERED REGIONS will not be predicted realistically and may interfere with the rest of the prediction.

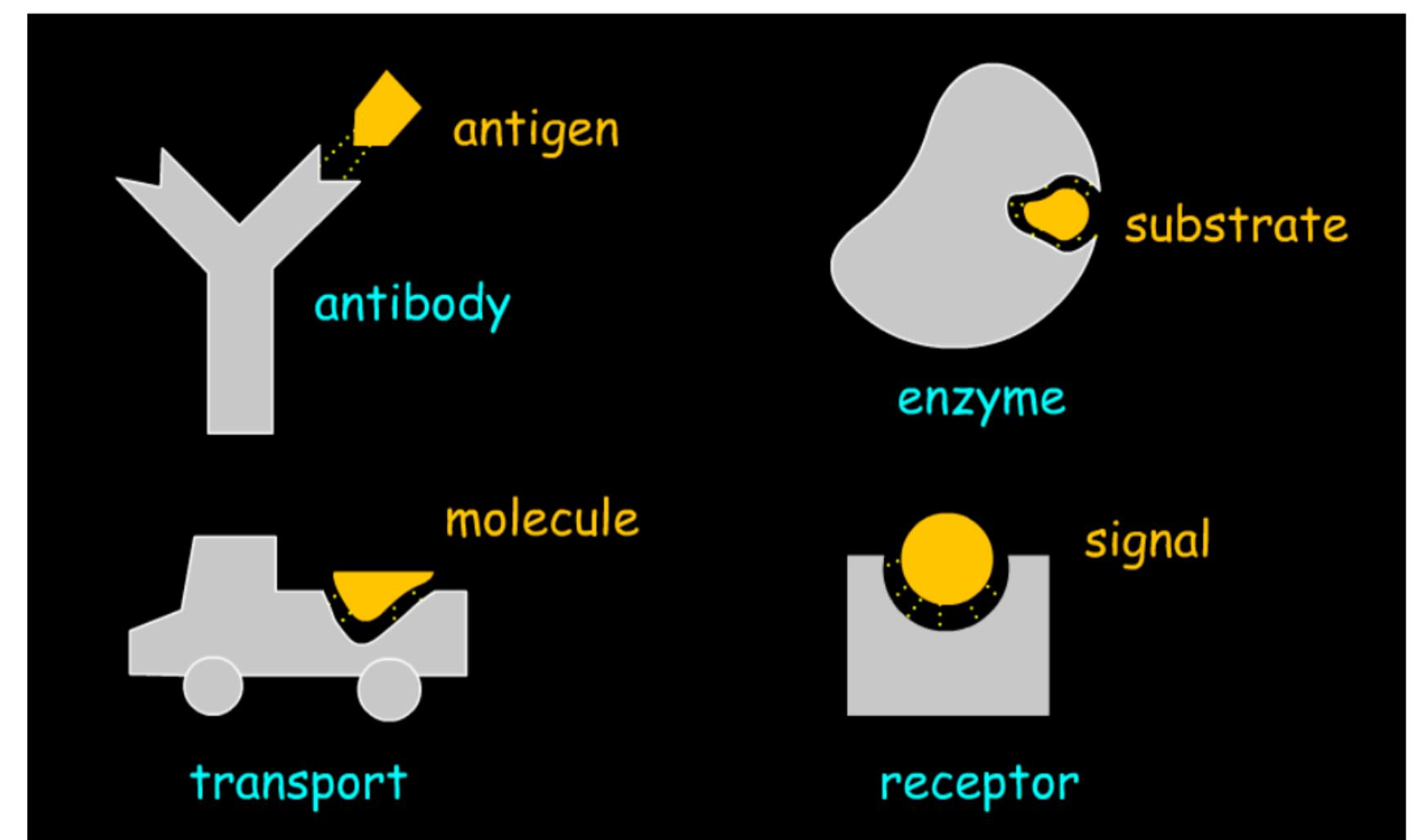
Limited MSA may be reflected in a poor PLDDT score. Remember that AF2 and related methods are designed to make predictions based on MSA.



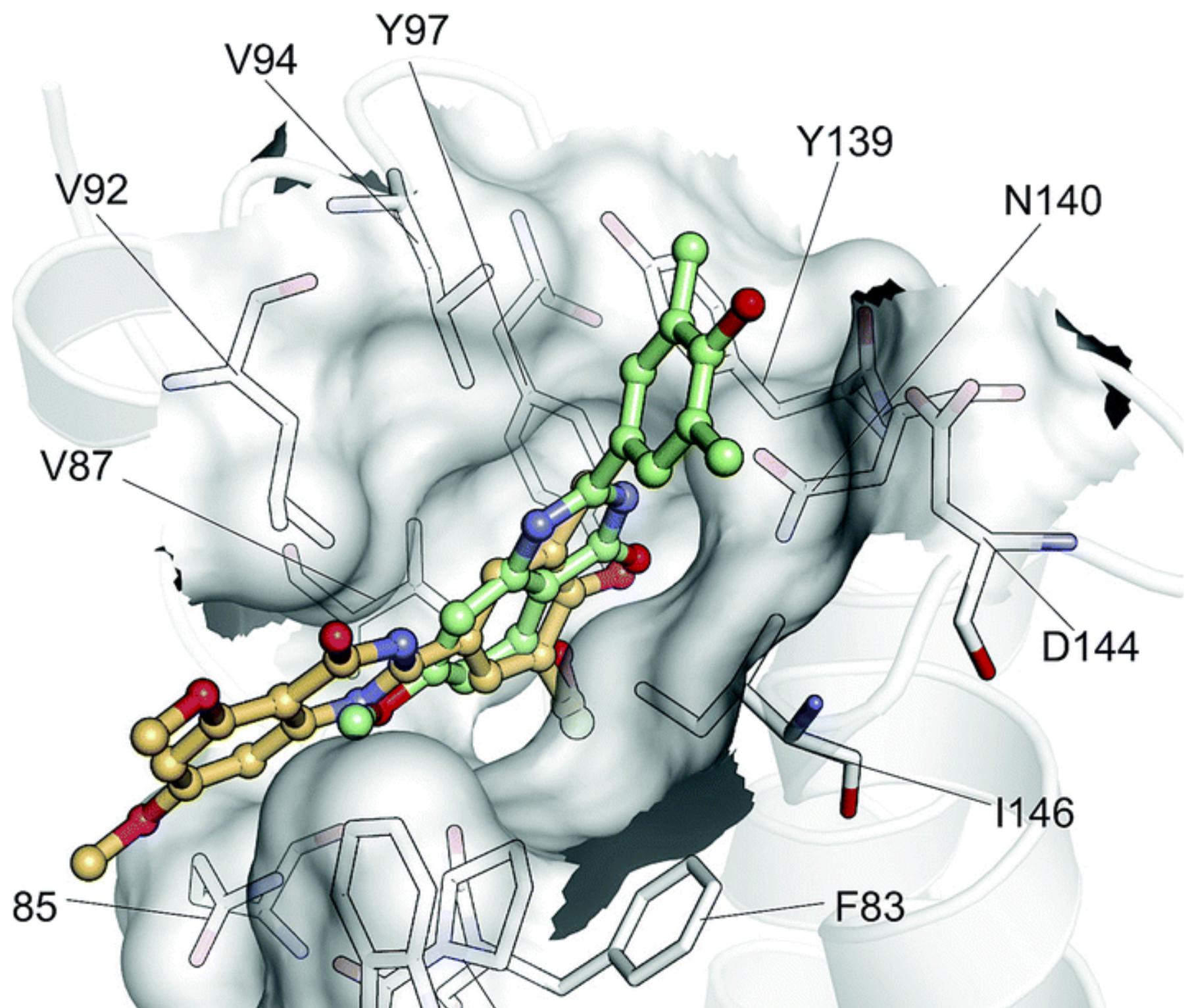
# Molecular Recognition and Docking

**Molecular recognition** is the ability of biomolecules to recognize other biomolecules and selectively interact with them in order to promote fundamental biological events such as transcription, translation, signal transduction, transport, regulation, enzymatic catalysis, viral and bacterial infection and immune response.

**Molecular docking** is the process that involves placing molecules in appropriate configurations to interact with a receptor. Molecular docking is a natural process which occurs within seconds in a cell.



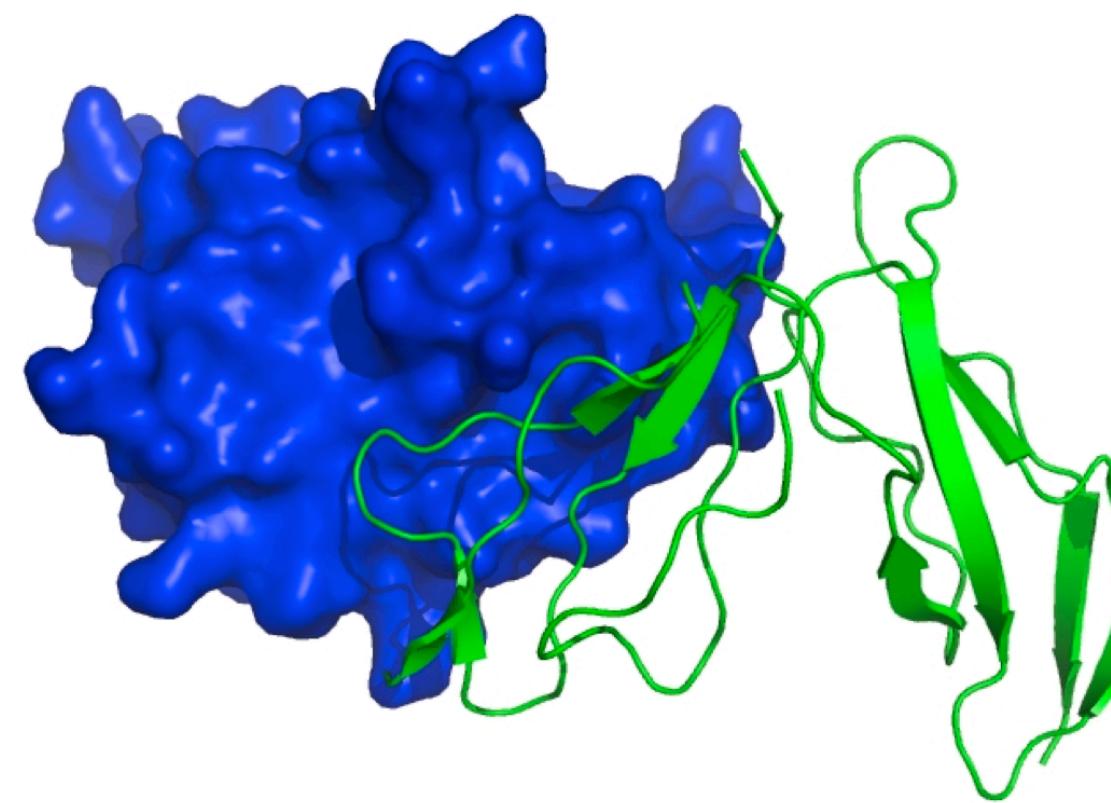
# Molecular Docking: binding pose and affinity



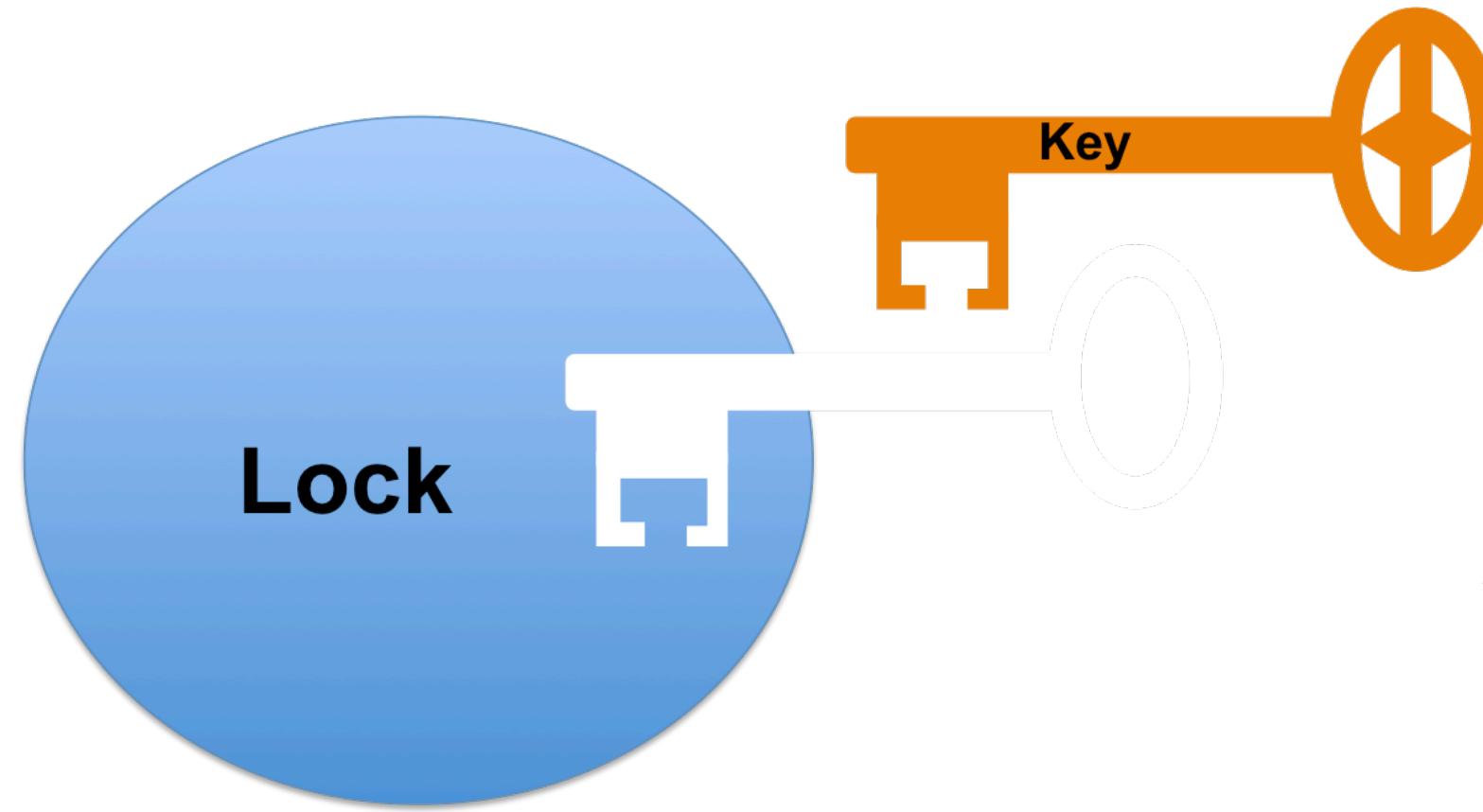
**Prediction of the bound structure and binding strength for**

- **Small molecules (virtual screening, drug discovery, ..)**
- **Macromolecular complexes (protein-protein, protein-nucleic, ...)**

**Complexation can be associated with conformational changes**



# Lock and Key



Emil Fischer (1894)

Specificity in enzyme-substrate recognition

Generally speaking the idea is that specificity in molecular recognition is the result of rigid surface that as a consequence can bind together only when exactly complementary (**Lock and Key model**)

## Docking algorithm: (0 level)

- analyse surface of the two molecules
- find possible regions of binding (surface compatibility -> shape and electrostatic interactions)
- minimise the energy.

Actually in many cases recognition mechanisms are more complex, with in principle flexibility on both the receptor and the ligand (**Flexible Docking**)



# Docking mechanisms

---

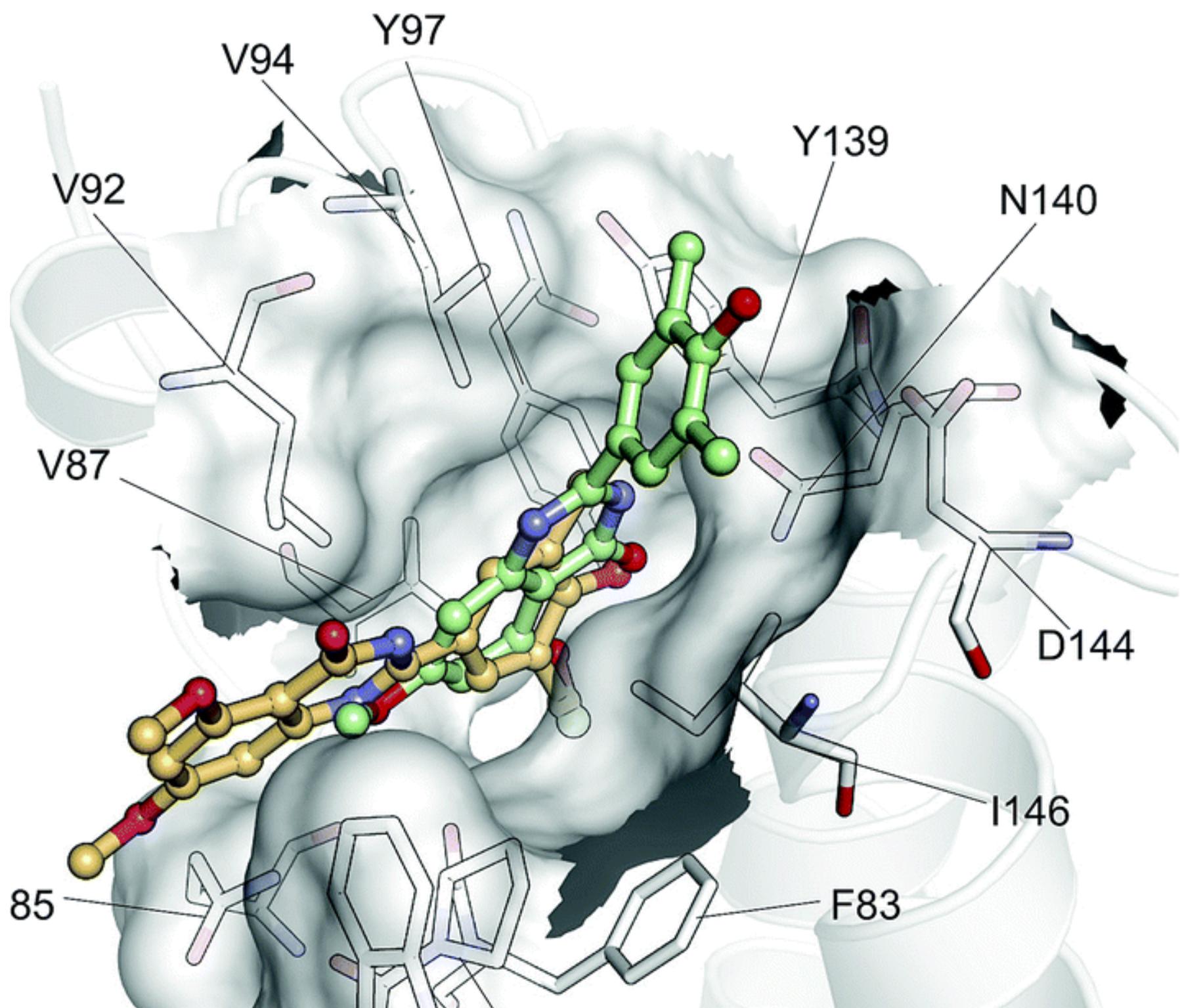
**Lock-and-Key:** this is simple, the assumption is that both the receptor and the ligand are rigid and their structure is ideal for binding. The next step is to accept that the ligand (this is more specific for the binding of small molecules) can be flexible, while the receptor is rigid.

**Induced Fit:** In 1958 Daniel Koshland introduced the "induced-fit theory". The basic idea is that in the recognition process, both ligand and target mutually adapt to each other through small conformational changes, until an optimal fit is achieved. This is traduced in what is usually called the flexible-docking where both active site region and the ligand are allowed a certain degree of flexibility.

**Conformational Selection:** this recognises that both the ligand and the receptor can be in many different configurations, and that binding can only occur when relatively optimal configurations are populated. This is traduced in the so called ensemble docking, where multiple configurations of the receptor are used at the same time each can also be considered locally flexible to account for additional induced fit.

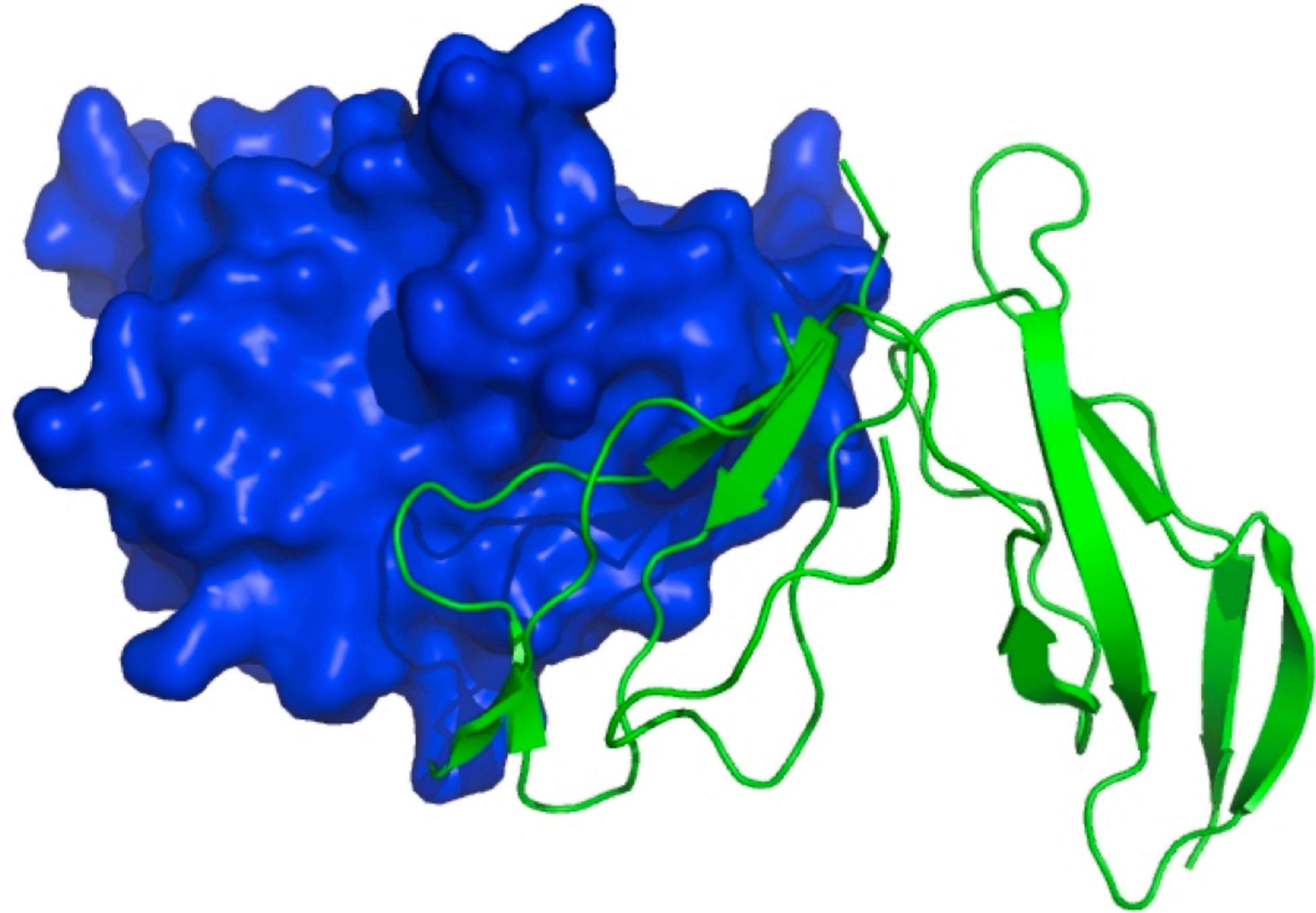


# Molecular Docking: softwares



**DiffDock**

...

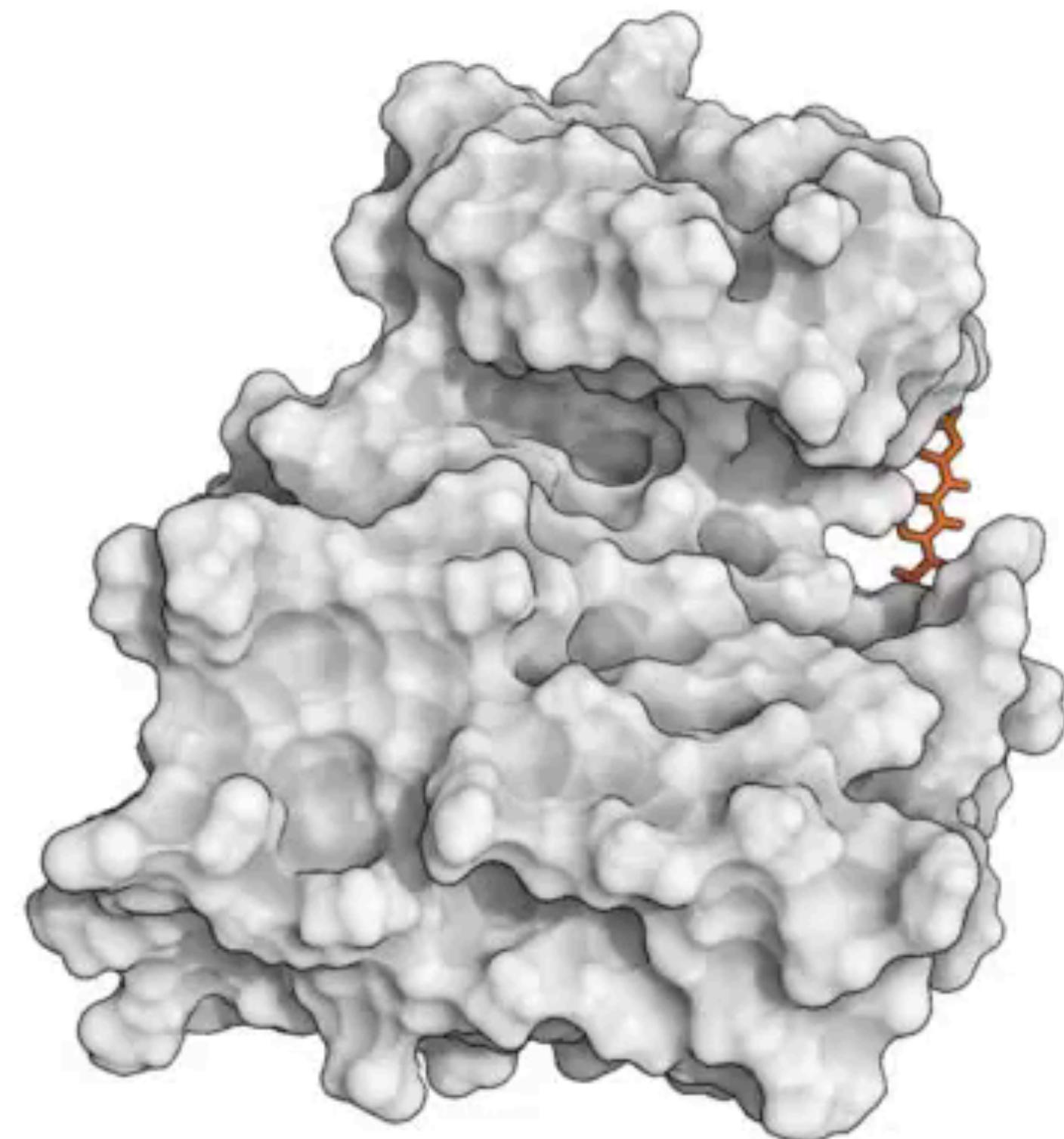


...

...

# MD simulations

---



In principle MD simulations can be used to search for ligand binding, but there are problems:

- 1) Parametrizing force-fields for all the possible chemicals is hard (this is more relevant for ligand binding than for protein complexes)
- 2) The binding time scale can be slow in particular when associated with conformational changes (the problem of sampling). This is even more relevant for protein complexes where the systems become large and so the simulations become slow.



# Rigid or semi-flexible docking: a simplified simulation approach

---

**There are two problems:**

- 1. How to generate quickly as many ‘reasonable’ poses as possible (including or not flexibility)**
- 2. How to distinguish (‘score’) a good pose from a bad one (will the true complex score better than all other possible complexes?)**

A number of poses can be generated at random by rotation and translation of a molecule around the other (using Monte Carlo or genetic algorithms) and then they can be scored for example using a physicochemical inspired scoring function (a simplified force-field). Protein Flexibility can be introduced using multiple structures, while ligand flexibility can be introduced allowing the molecule to rotate around “rotatable” bonds.



# VINA scoring function

This is what is usually called a “scoring function” that is a mathematical object that given some numbers return a single number, e.g. all the x, y, z coordinates of a given configuration put into the scoring function return a single number that is the score for that configuration.

$$\Delta G_{binding} = \Delta G_{gauss} + \Delta G_{repulsion} + \Delta G_{hbond} + \Delta G_{hydrophobic} + \Delta G_{tors}$$

$\Delta G_{gauss}$

Attractive term for dispersion, two gaussian functions

$\Delta G_{repulsion}$

Square of the distance if closer than a threshold value

$\Delta G_{hbond}$

Ramp function - also used for interactions with metal ions

$\Delta G_{hydrophobic}$

Ramp function

$\Delta G_{tors}$

Proportional to the number of rotatable bonds

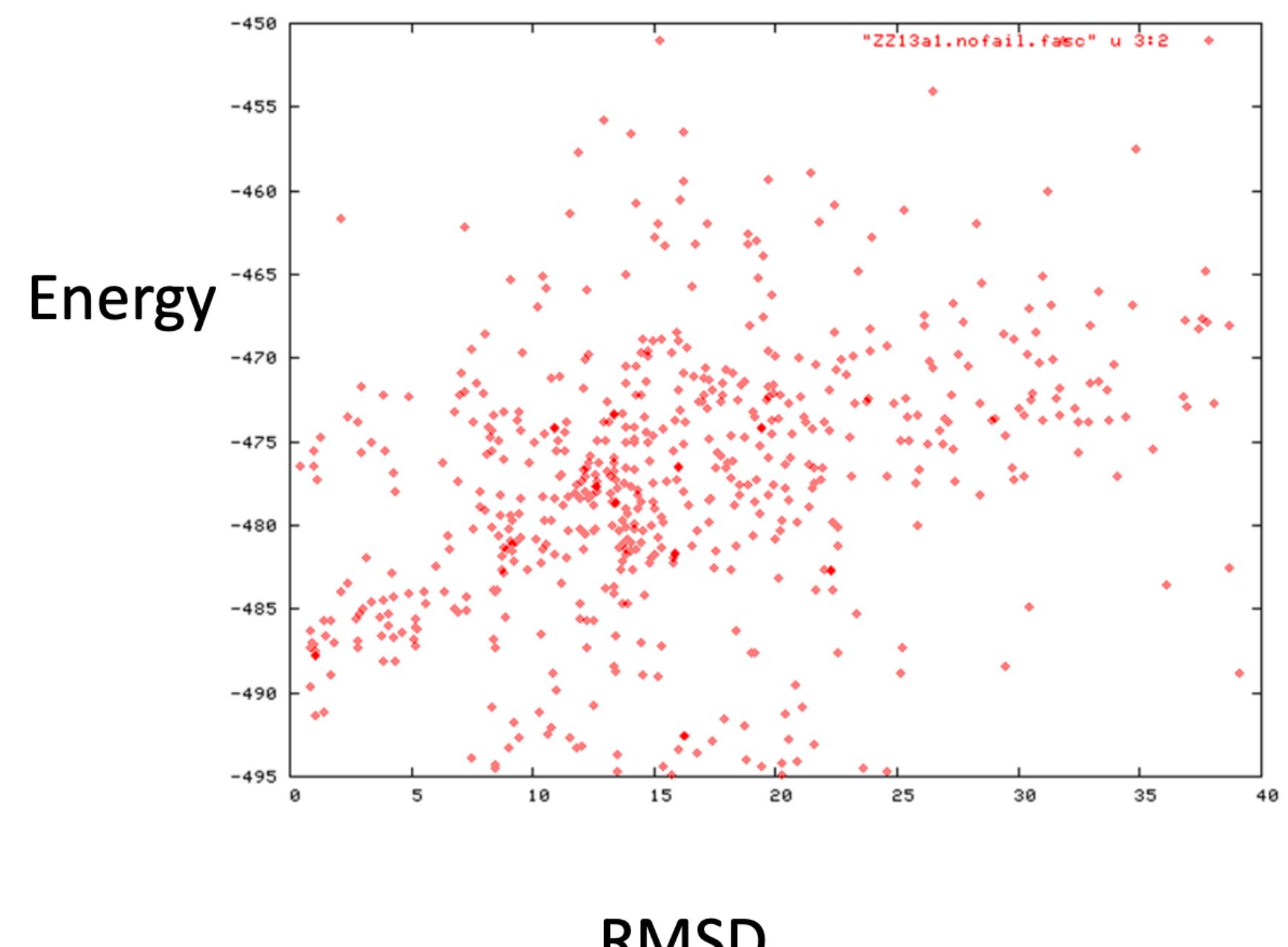
These are all related to intermolecular interactions. The hydrophobic term can be seen as the one accounting for solvation.



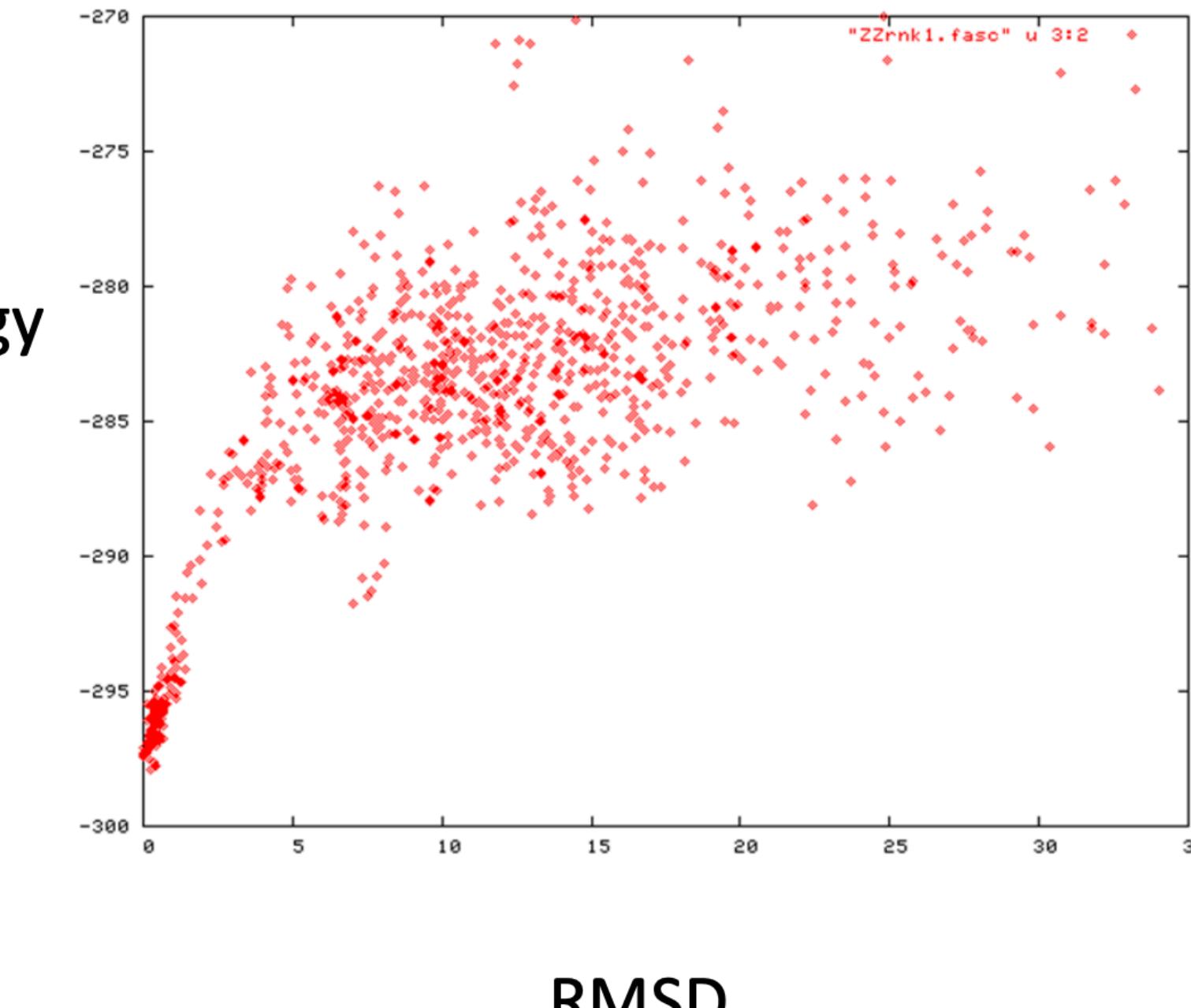
This is then only source  
of internal energy



# Scoring function and convergence



Bad



Good



# Monte Carlo sampling

**In principle one could generate configurations at random, evaluate their scoring function, and keep the lowest energy one. The limit of this approach is that most of the time will be spent evaluating scoring functions for bad scoring configurations.**

**Calculating  $P(x)$  would mean calculating it for all possible configurations, but the probability is simply proportional to:**

$$f(x) = \exp\left[\frac{-U(x)}{k_B T}\right]$$

$$f(x) \propto P(x)$$

**Starting from a configuration we can calculate its energy or score, then we can generate a second configuration and get its score and we need to decide if to keep it or not. How? Their relative probability is:**

$$\frac{f(x_{new})}{f(x)} = \exp\left(\frac{-U(x_{new}) + U(x)}{k_B T}\right)$$

**Metropolis Monte-Carlo (but often referred simply as Monte Carlo) is an algorithm to evaluate a probability distribution.**

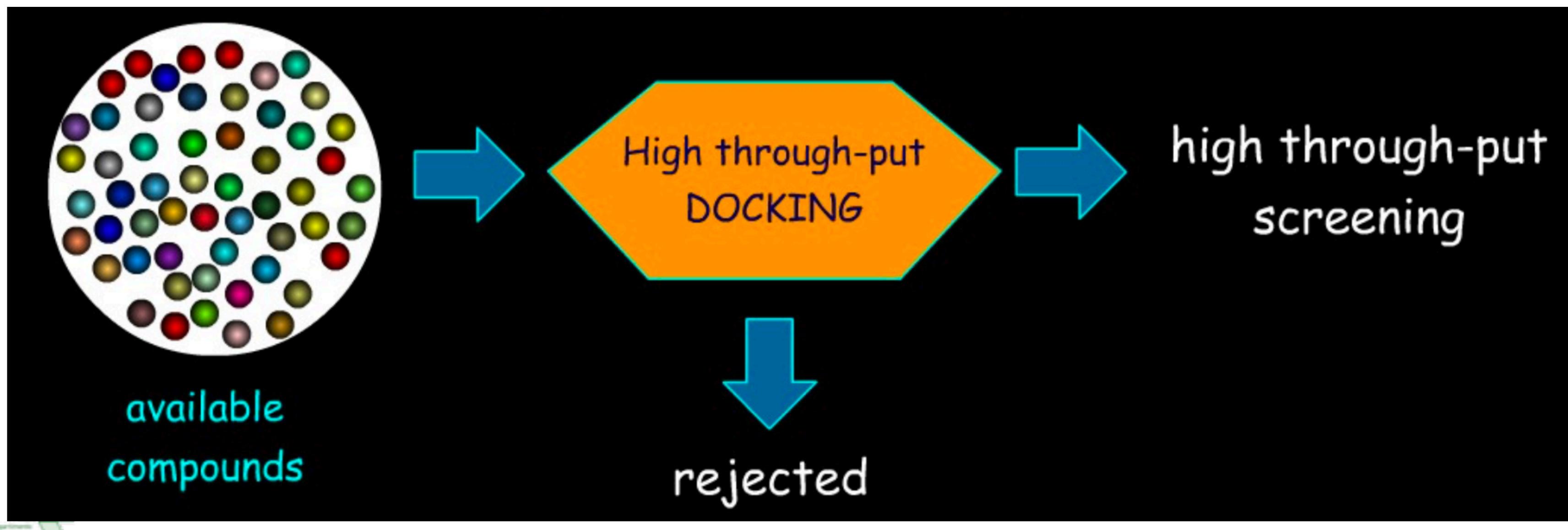
$$P(x) = \frac{\exp\left[\frac{-U(x)}{k_B T}\right]}{\sum \exp\left[\frac{-U(x)}{k_B T}\right]}$$

**Do we keep the new one or not? We generate a random number (0..1), if it is smaller than the above ratio we keep the new configuration, otherwise the old.**



# Virtual Screening

- When the goal of docking is to dock all the compounds of a library (the molecules being available or not yet synthesized), the process is **called virtual screening** or **high throughput docking**
- Virtual screening identifies active compounds in a large database and ranks them by their affinity to the receptor
- The method is not used to recognize active molecules but to **eliminate those that are likely to be inactive**



# Complementary approaches

---

- Pocket identification algorithms: what are the features of a good binding site? (P2rank, ...)
- Ligand design methods (which chemistry is optimal to bind in a given pocket?) here generative AI can help a lot (DiffSBDD, ...)
- Structure prediction methods are trying to implement also the prediction of chemicals: AlphaFold3, RosettaFold-AllAtom, ...
- Prediction of ADME/PK (Absorption, Distribution, Metabolism, Excretion and Pharmacokinetics) refers to the complete set of properties of a drug molecule describing its entry into the body, residence time within the body, distribution to organs, metabolic transformations involved in its clearance, and routes of elimination

