

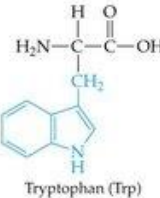
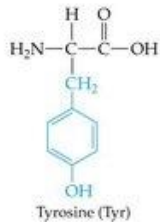
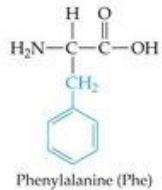
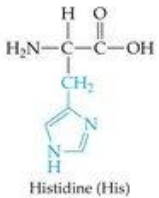
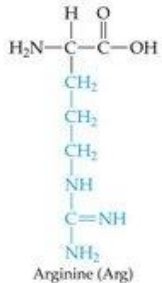
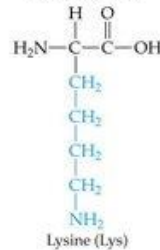
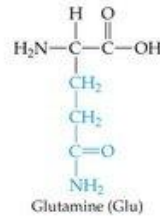
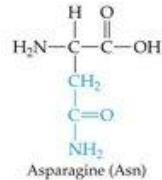
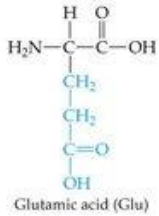
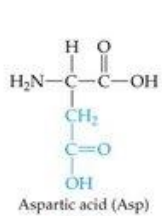
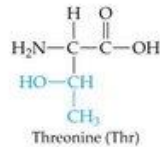
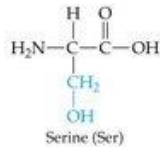
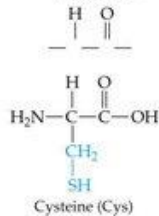
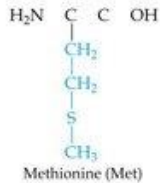
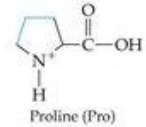
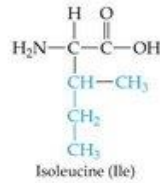
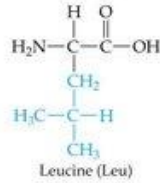
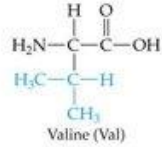
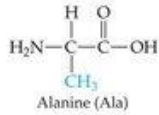


Free University of Bolzano - Bozen  
Faculty of Engineering

SFSCON

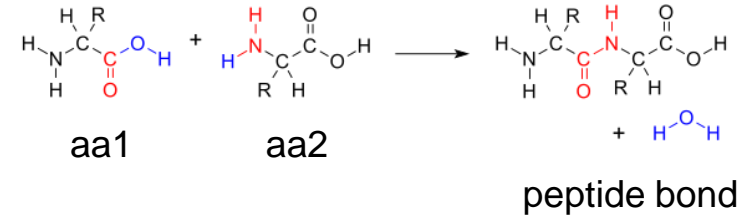
# Machine learning-driven simulation of protein folding atomistic trajectories

Alan Ianeselli  
SFSCON 2023

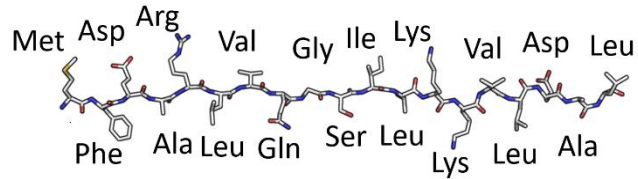


~20 amino acids

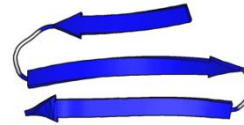
They link by peptide bonds to form a polymer



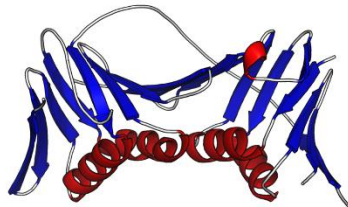
A protein is a polymer of tens, hundreds or thousands of amino acids



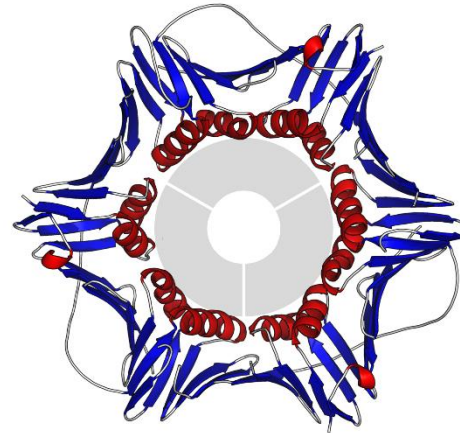
Primary structure  
(aa sequence)



Secondary structure  
(local structures)

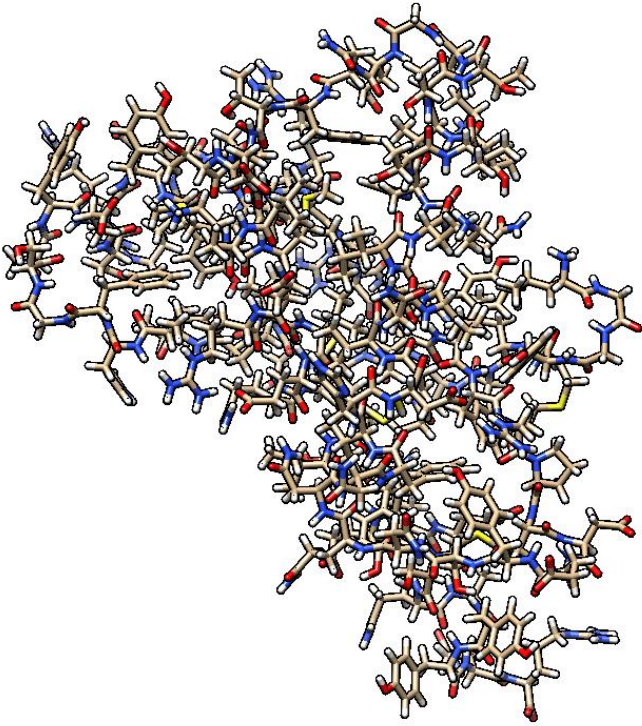


Tertiary structure  
(3D conformation)

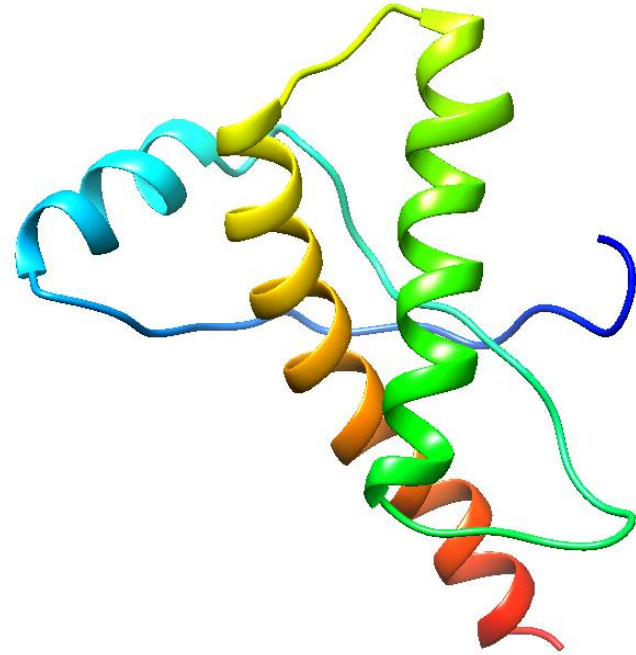


Quaternary structure  
(protein complexes)

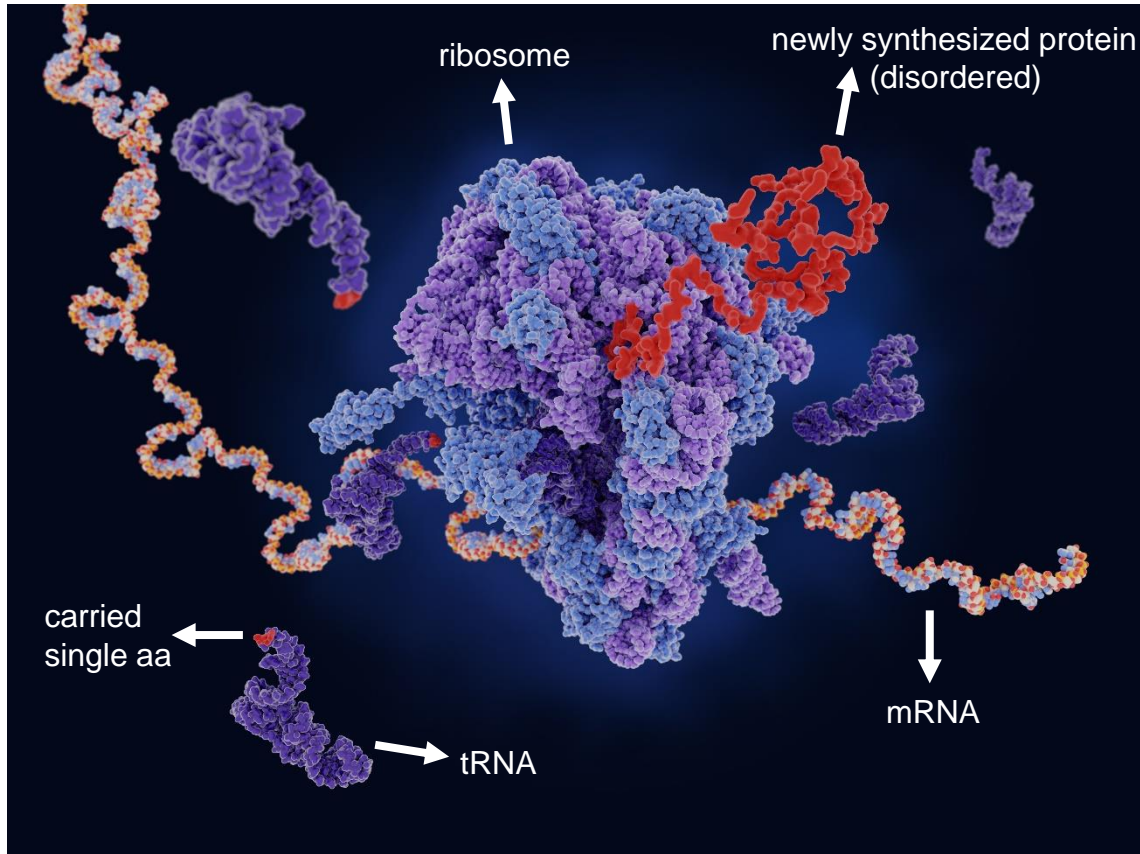
→ Unique conformation given a specific aminoacidic sequence  
= the **protein folding problem**



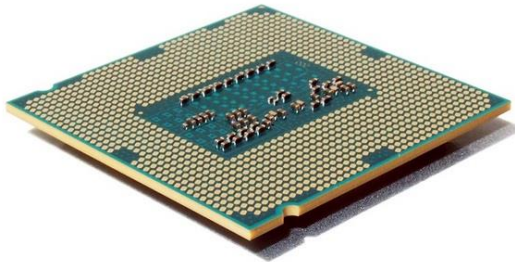
All-atom  
representation



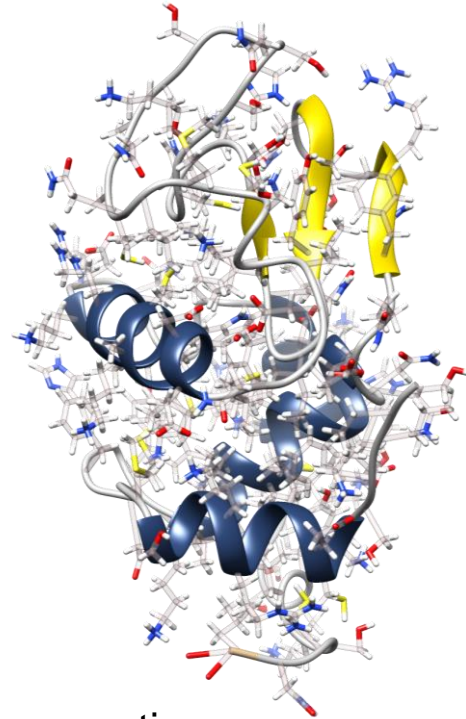
Ribbon  
representation



→ How does the newly synthesized disordered protein achieve its final conformation?



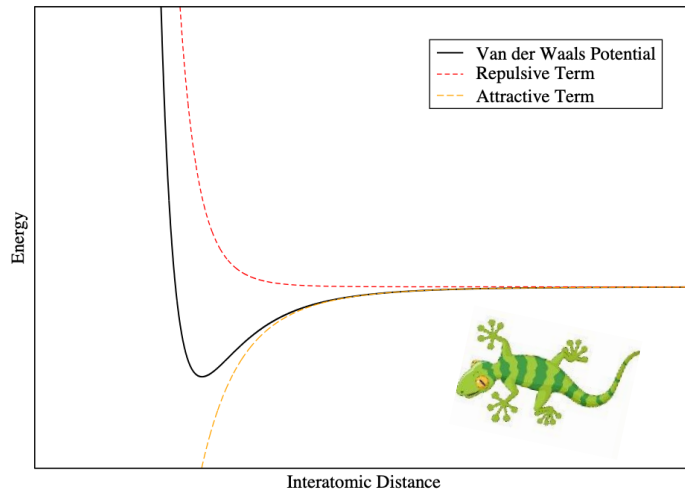
$$\vec{F} = m\vec{a}$$



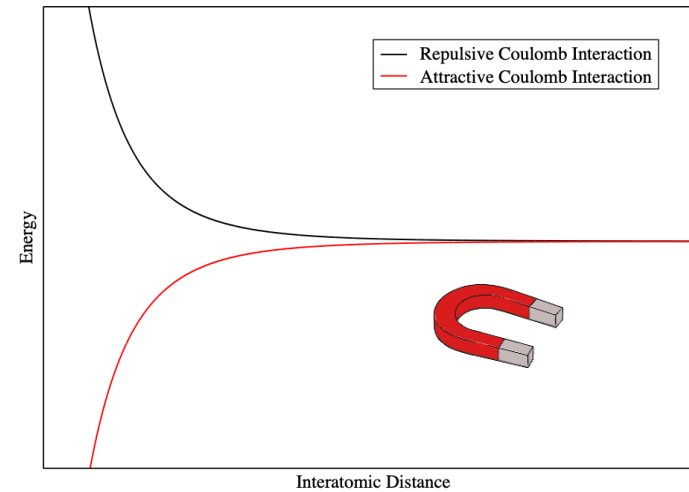
- ▶ Numerically solve Newton's equations of motion over time for each atom of the system

- Forces are computed from *force fields*

## Non-bonded interactions



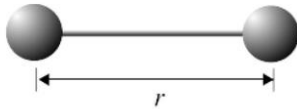
Van der Waals



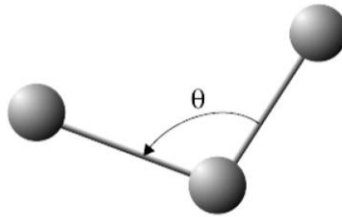
Coulomb

- Forces are computed from *force fields*

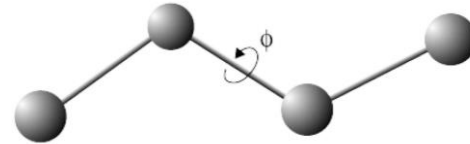
### Bonded interactions



Bond stretching

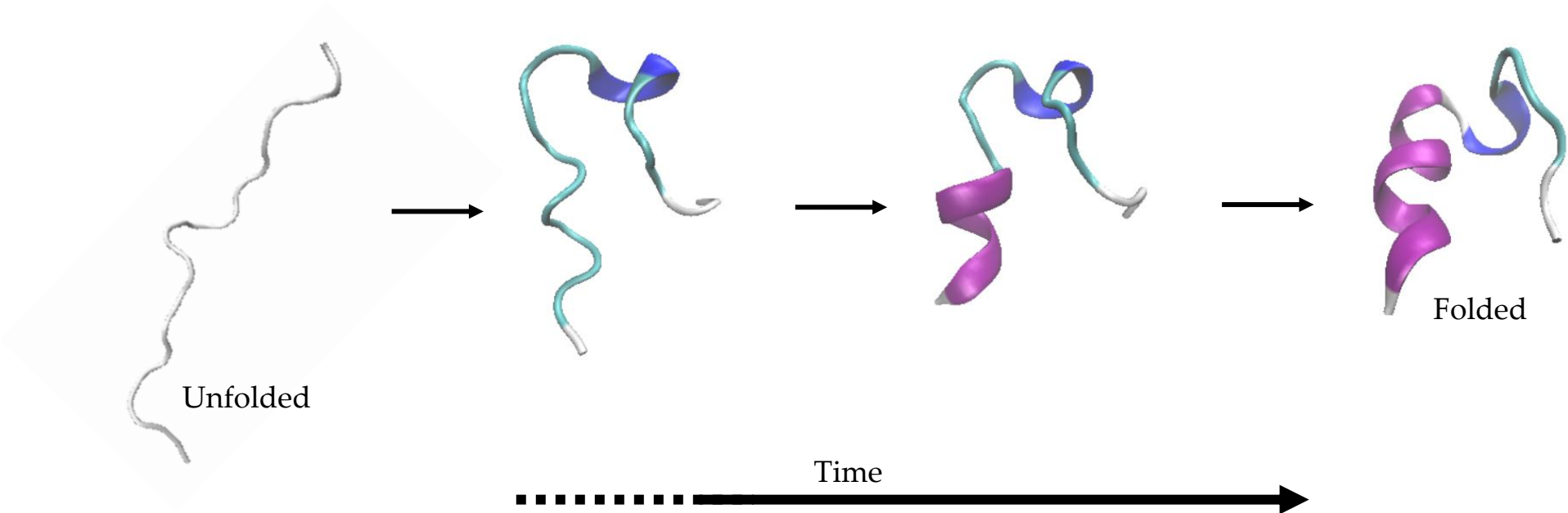


Angle bending



Dihedral torsions





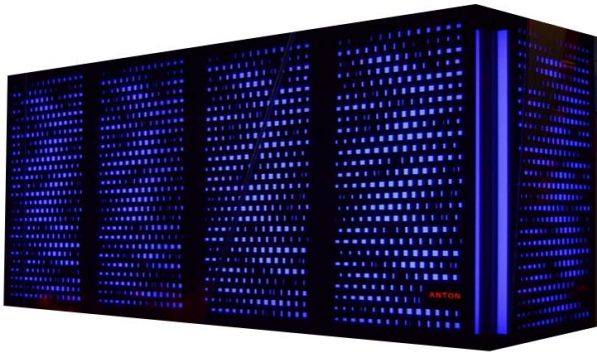
Unfolded

Folded

Time

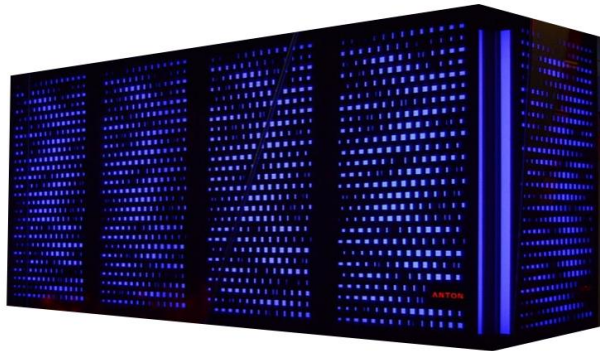
Folding time ~ microseconds (for small proteins)  
= Weeks of simulation on supercomputers

*Anton* Supercomputer  
~50  $\mu$ s/day for ~100'000 atoms



DE Shaw et al., 2009

*Anton* Supercomputer  
~50  $\mu$ s/day for ~100'000 atoms



DE Shaw et al., 2009

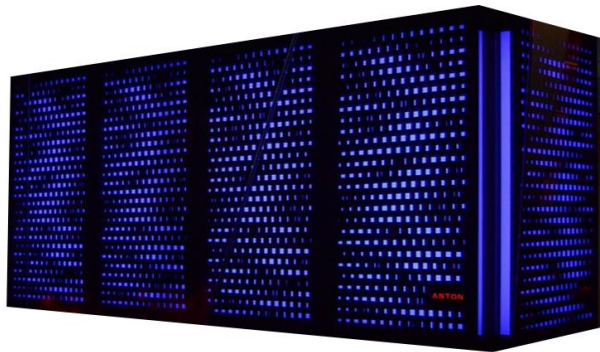
For example, Lysozyme in water (~100'000 atoms)  
requires SECONDS to fold



~100 years of simulation!!



*Anton* Supercomputer  
~50  $\mu$ s/day for ~100'000 atoms



DE Shaw et al., 2009

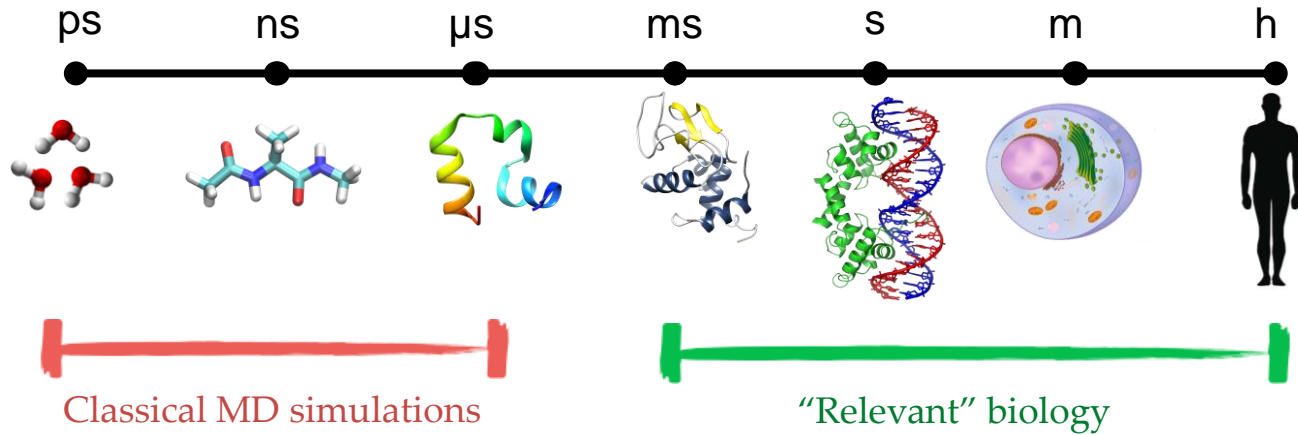
For example, Lysozyme in water (~100'000 atoms)  
requires SECONDS to fold

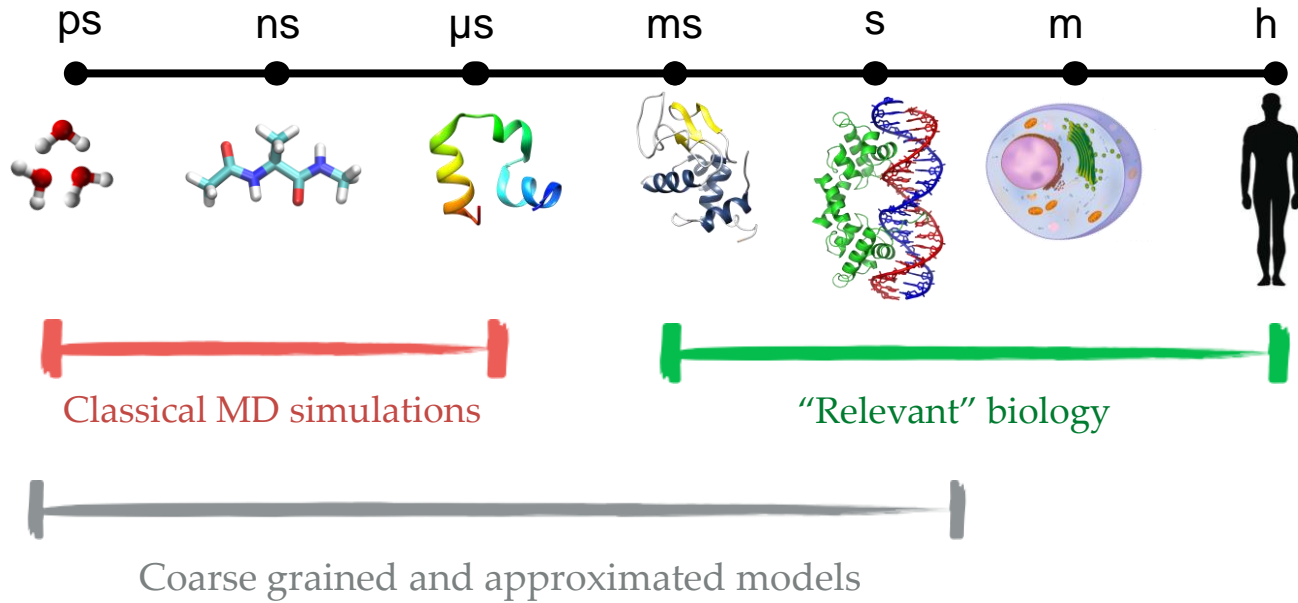


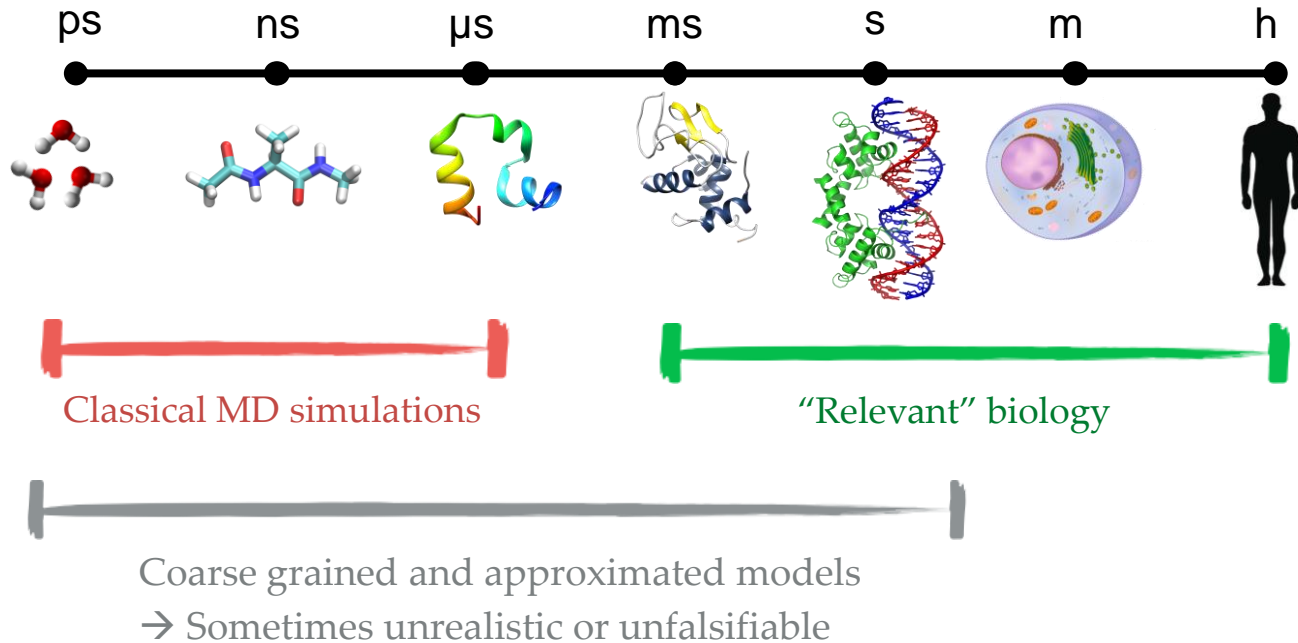
~100 years of simulation!!



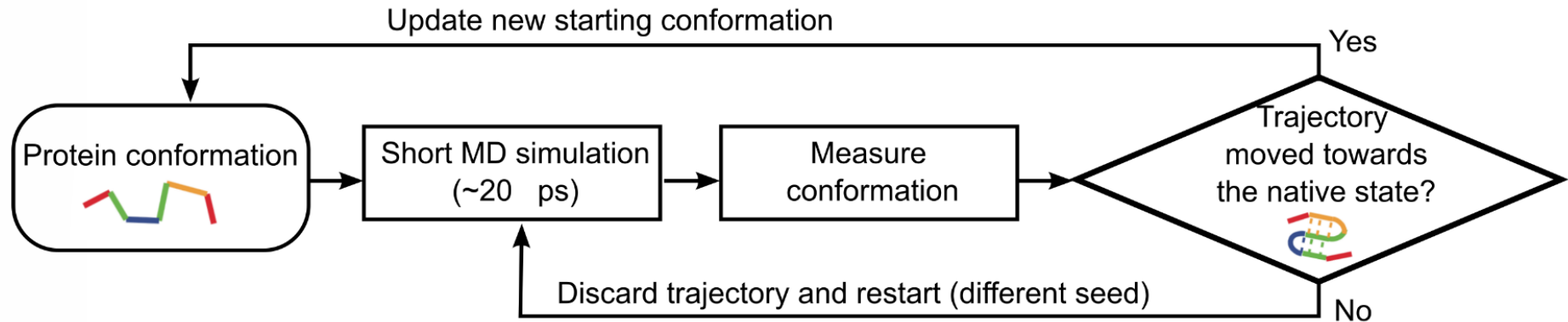
▸ Conventional MD approaches are unfeasible





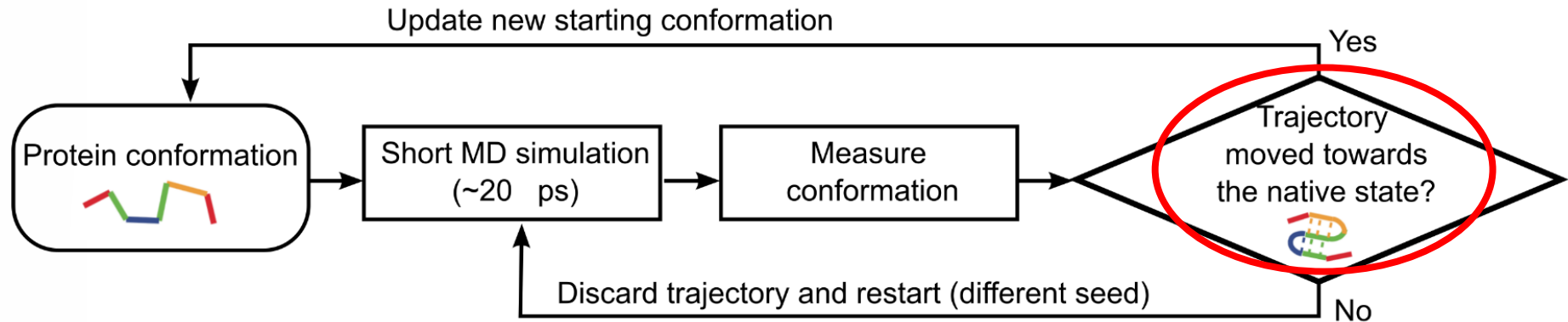


A smart **algorithm** to study protein folding trajectories

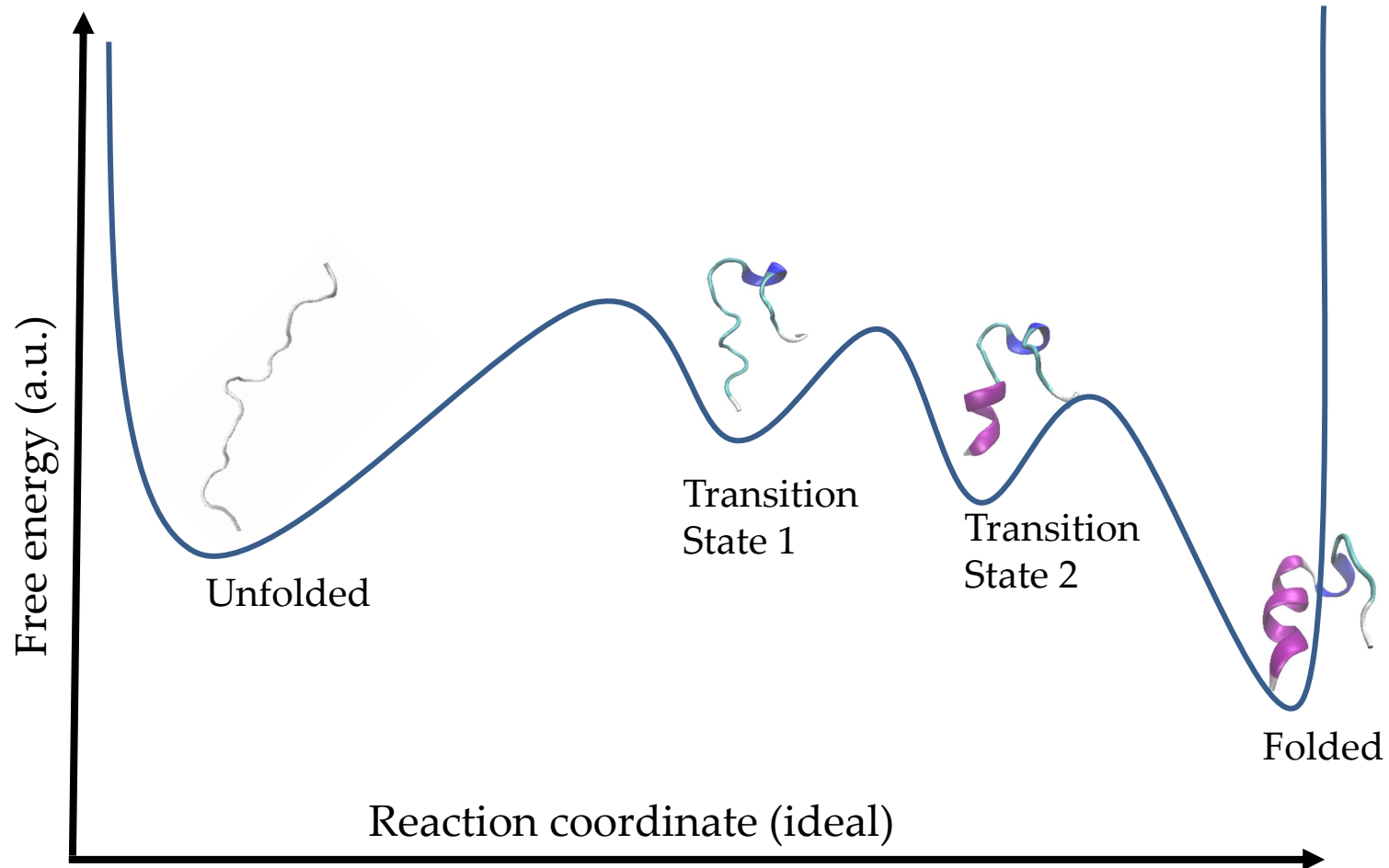


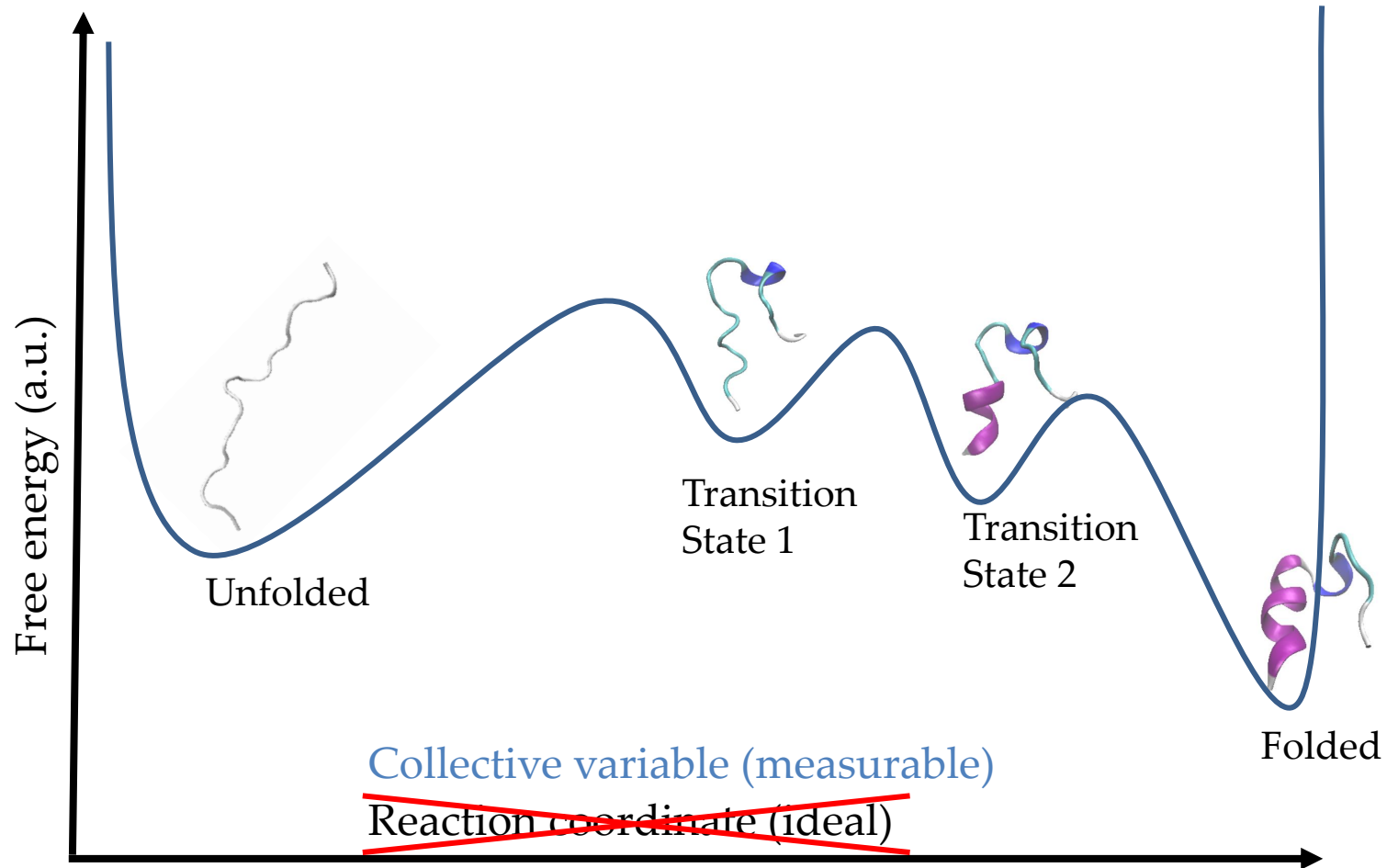


A smart **algorithm** to study protein folding trajectories

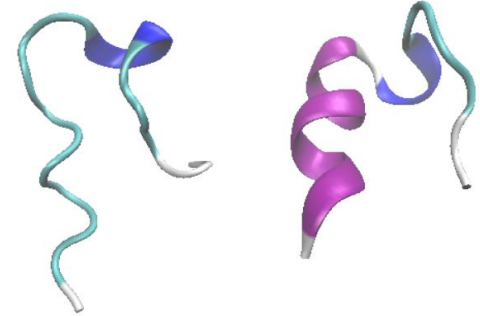


How do you measure if it went “forward”?





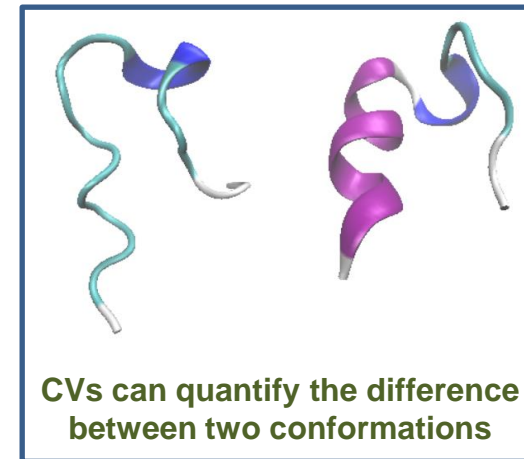
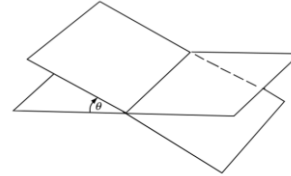
„Best“ collective variables (CVs)  
for protein folding



CVs can quantify the difference  
between two conformations

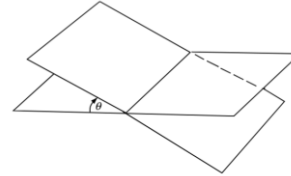
## 1. Dihedral angles deviation from the native state

Syzonenko *et al*,  
J. Chem. Inf. Model. 2020



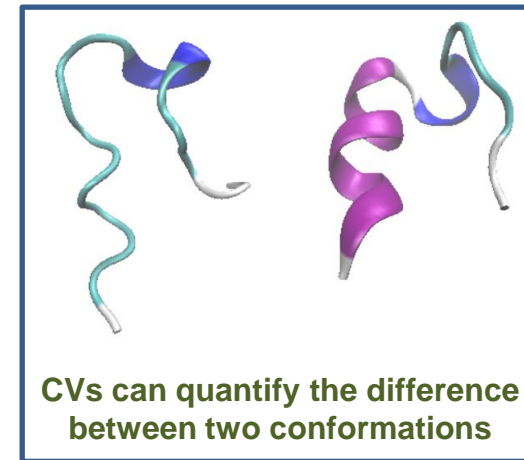
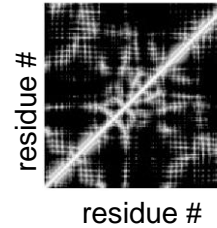
# 1. Dihedral angles deviation from the native state

Syzonenko *et al*,  
J. Chem. Inf. Model. 2020



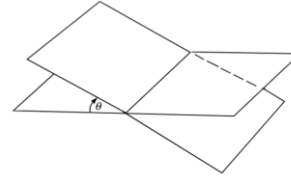
# 2. Inter-aa contact deviation from the native state

Beccara *et al*,  
Phys. Rev. Lett. 2015



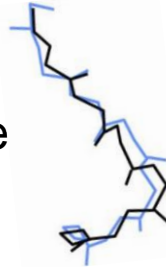
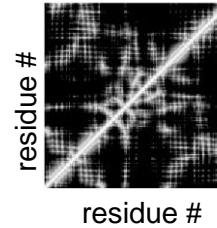
## 1. Dihedral angles deviation from the native state

Syzonenko *et al*,  
J. Chem. Inf. Model. 2020

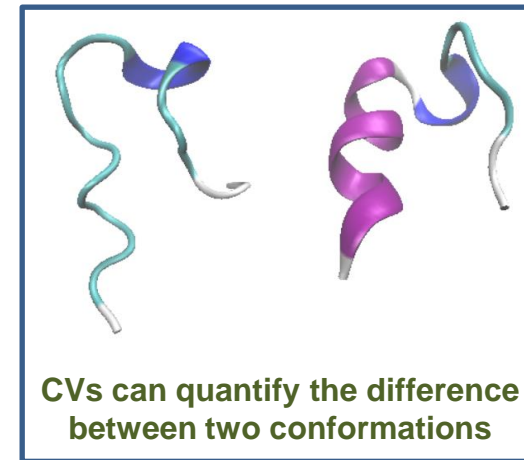


## 2. Inter-aa contact deviation from the native state

Beccara *et al*,  
Phys. Rev. Lett. 2015

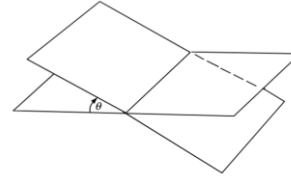


## 3. Geometrical difference from the native state



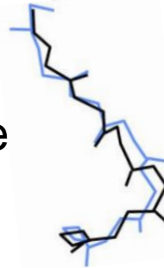
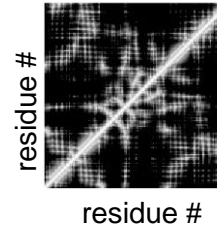
## 1. Dihedral angles deviation from the native state

Syzonenko *et al*,  
J. Chem. Inf. Model. 2020



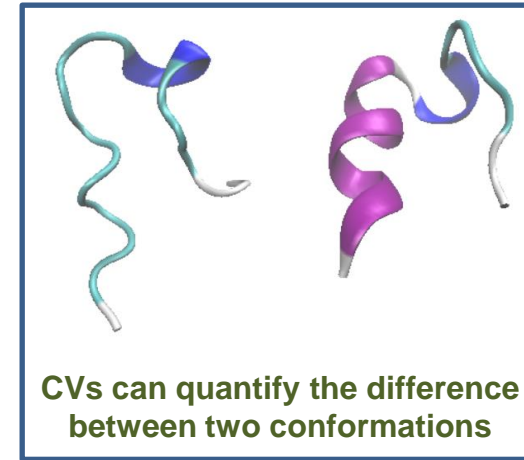
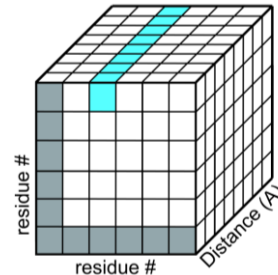
## 2. Inter-aa contact deviation from the native state

Beccara *et al*,  
Phys. Rev. Lett. 2015



## 3. Geometrical difference from the native state

## 4. Deviation from Google's Deepmind AlphaFold tensor





#### 4. Deviation from **Google's Deepmind AlphaFold** tensor

Google's Deepmind AlphaFold:

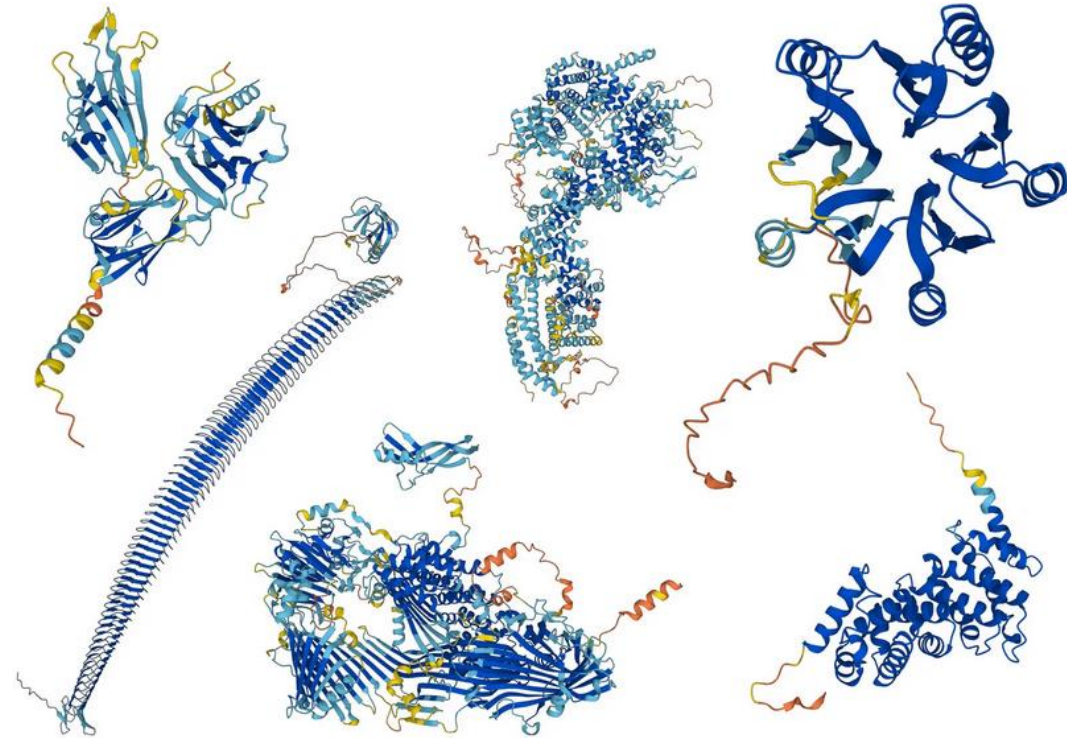
Latest AI milestone in the protein folding field

Training set: **170k** protein structures

Able to predict more than **200 million** of structures

Unprecedented accuracy

→ able to predict the final conformation of any aminoacidic sequence



#### 4. Deviation from **Google's Deepmind AlphaFold** tensor

Google's Deepmind AlphaFold:

-**Input** = aa sequence (text string)

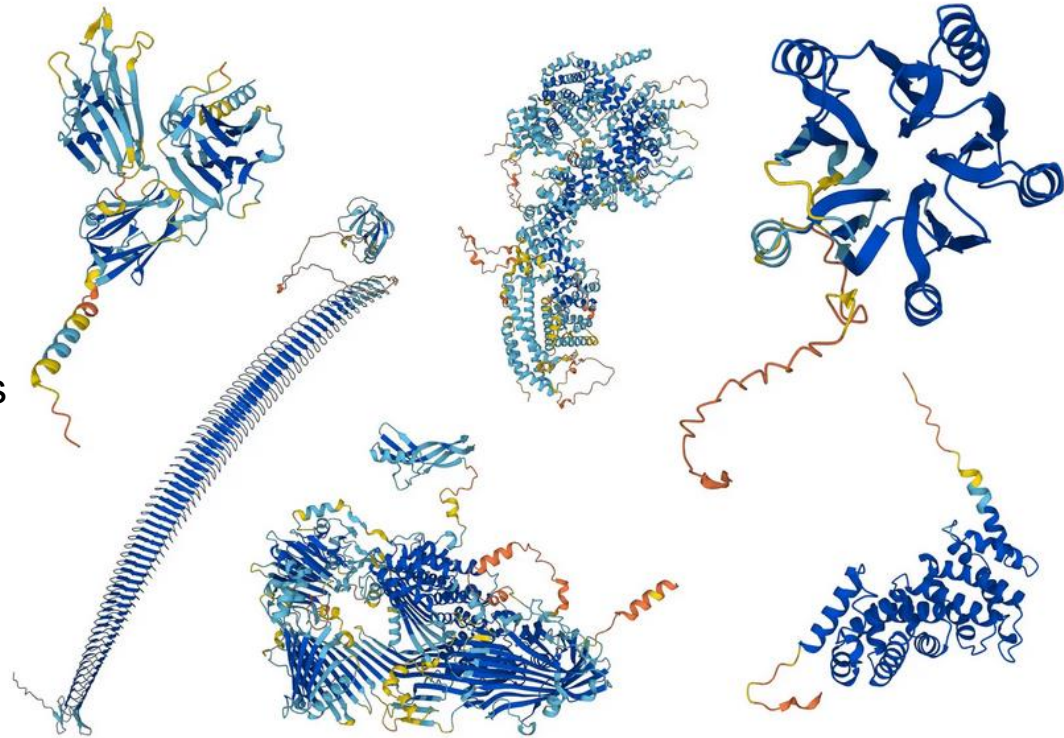
...

-Sequence alignment (database comparison)

-Prediction of distance and angle between aa pairs

...

-**Output** = 3D protein structure

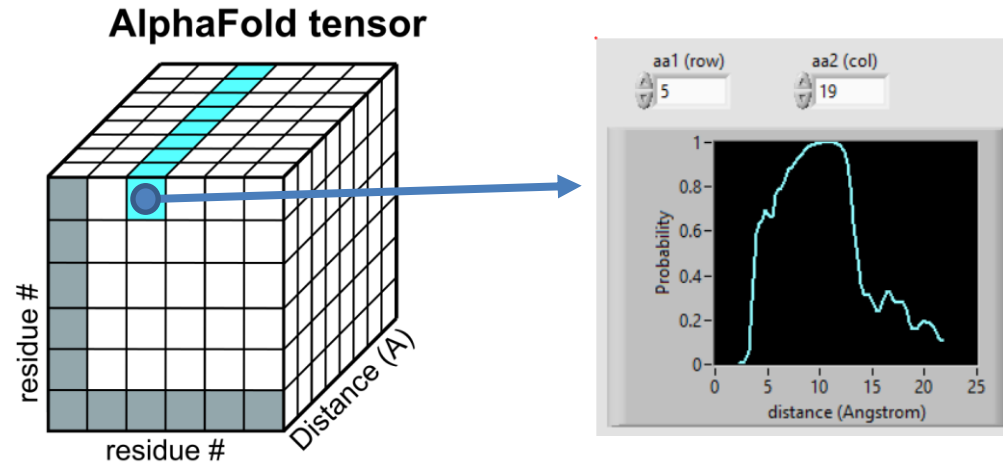


#### 4. Deviation from Google's Deepmind AlphaFold tensor

One of the outputs of AlphaFold is the so-called **distogram**:  
**Tensor** of distance bins x aa x aa

→ **probability over distance** between pairs of aa

example below: aa at position 5 (Asparagine) vs aa at position 19 (Aspartic acid)

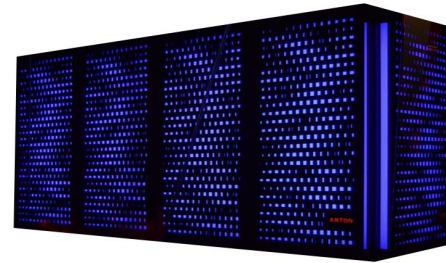


→ machine-learned 170k protein conformations

→ corresponds to a **quasi-chemical potential**

Training set:

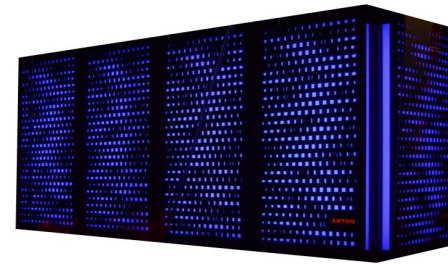
**Very long** folding trajectories obtained by the most powerful supercomputer (Anton)  
→ **200 $\mu$ s** of trajectories (Villin and Fip35 proteins)



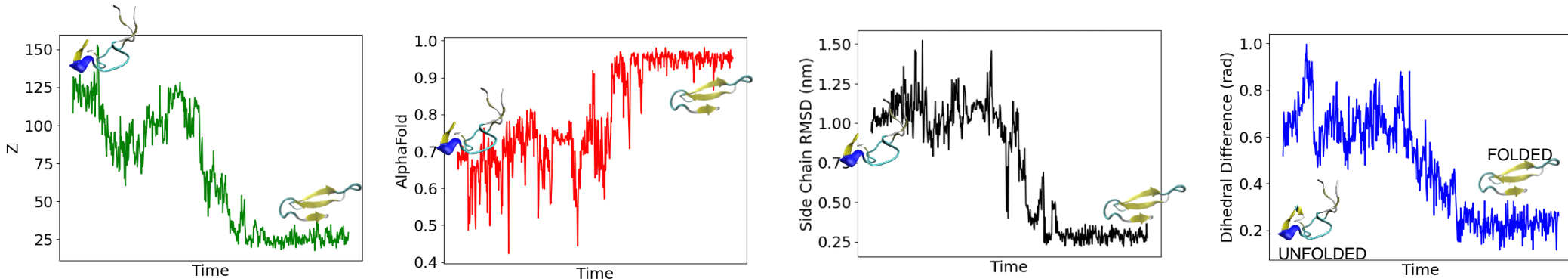
Anton

Training set:

**Very long** folding trajectories obtained by the most powerful supercomputer (Anton)  
 → **200 $\mu$ s** of trajectories (Villin and Fip35 proteins)



Anton



Training set:

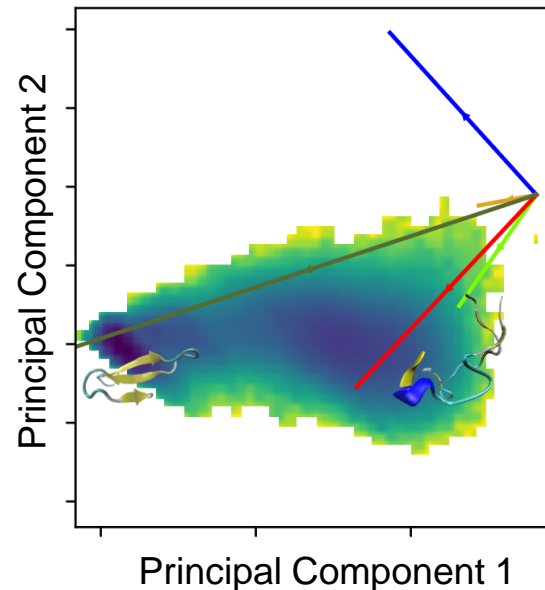
**Very long** folding trajectories obtained by the most powerful supercomputer (Anton)  
→ **200 $\mu$ s** of trajectories (Villin and Fip35 proteins)

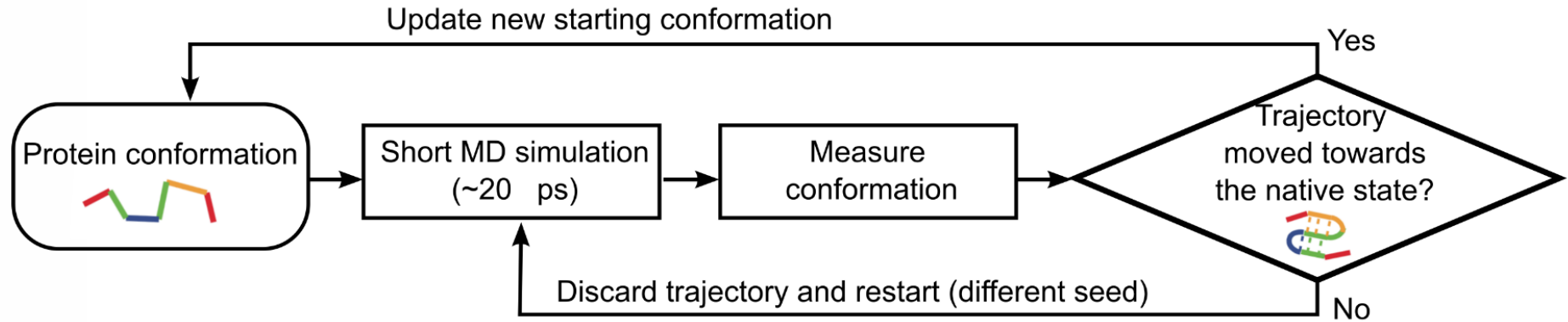
Training set:  
396k rows  
4 features

Machine learning (PCA)



**Optimal CV identified**



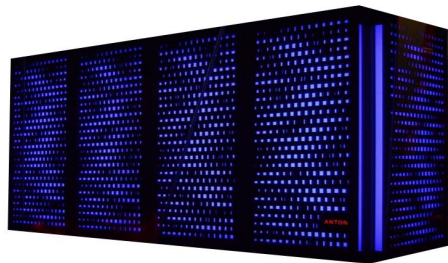


Towards the native state = **along the optimal CV**

**NOTE:** trajectories are unbiased

Obtained the **folding trajectories** of 4 small proteins

My simulation time = 1 day per trajectory on a weak laptop  
(Anton supercomputer would need weeks)



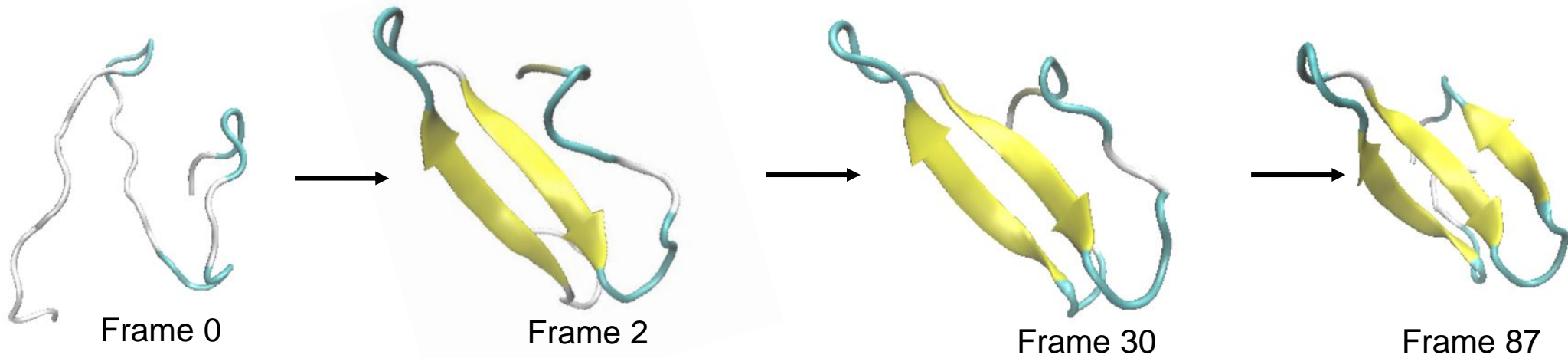
Brute force MD  
~weeks

VS

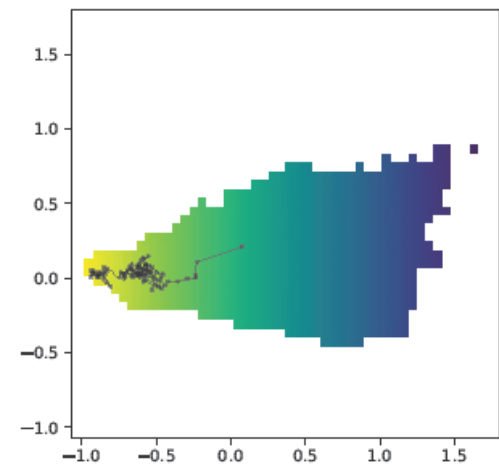


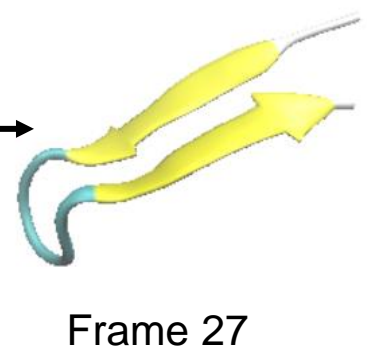
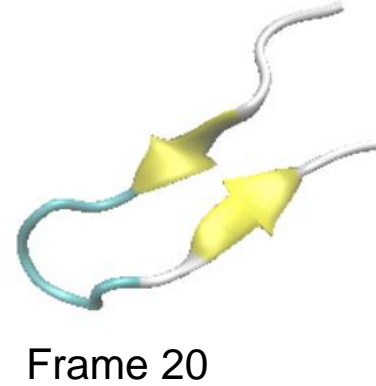
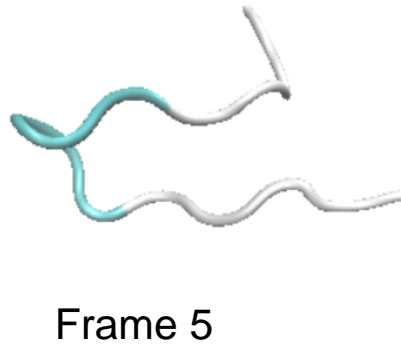
Smart algorithm  
~hours



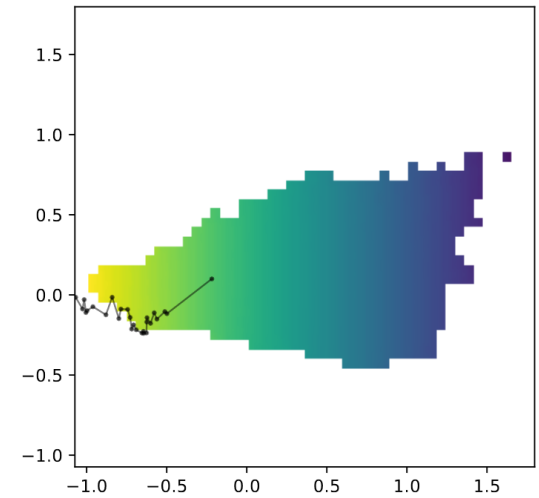


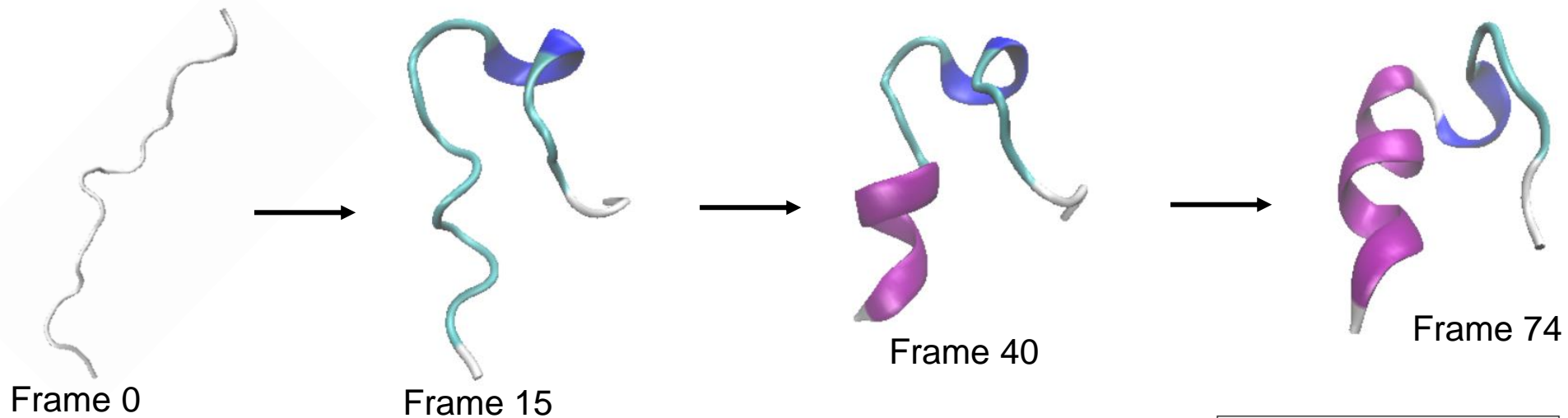
Fip35 (35 aa, 562 atoms): 16 trajectories, 1 folded



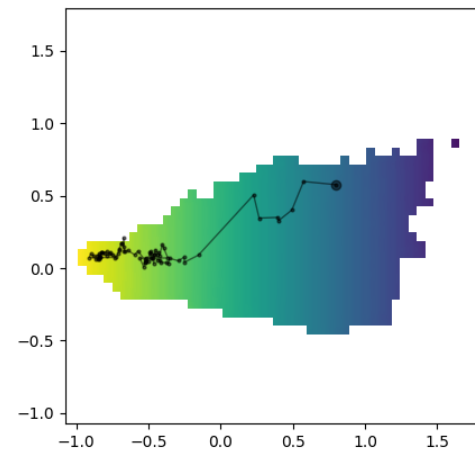


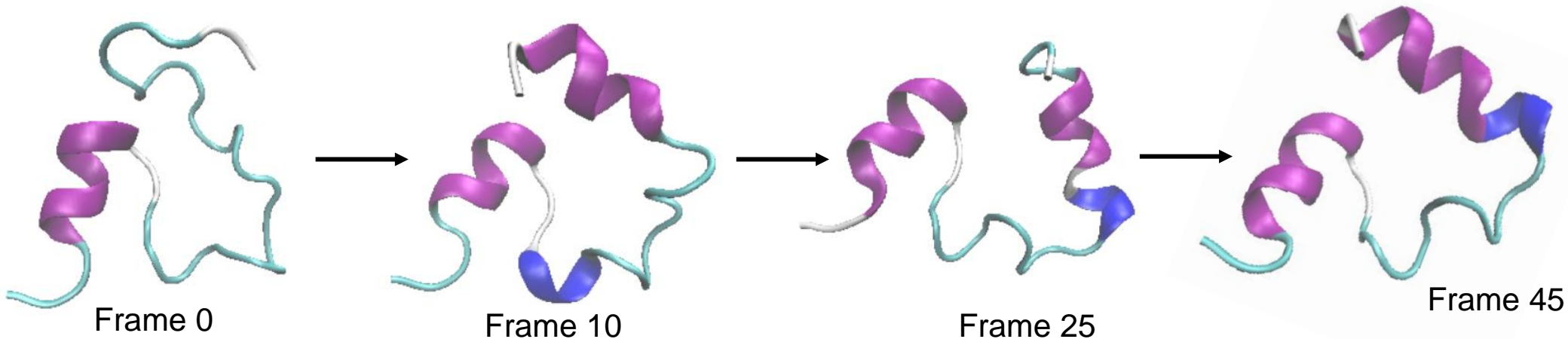
Beta Hairpin (16 aa, 247 atoms): 3 trajectories, 1 folded



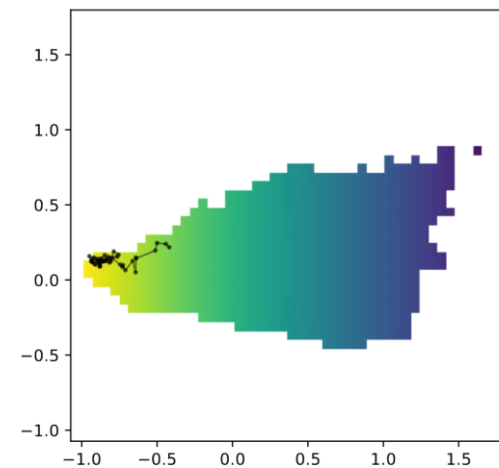


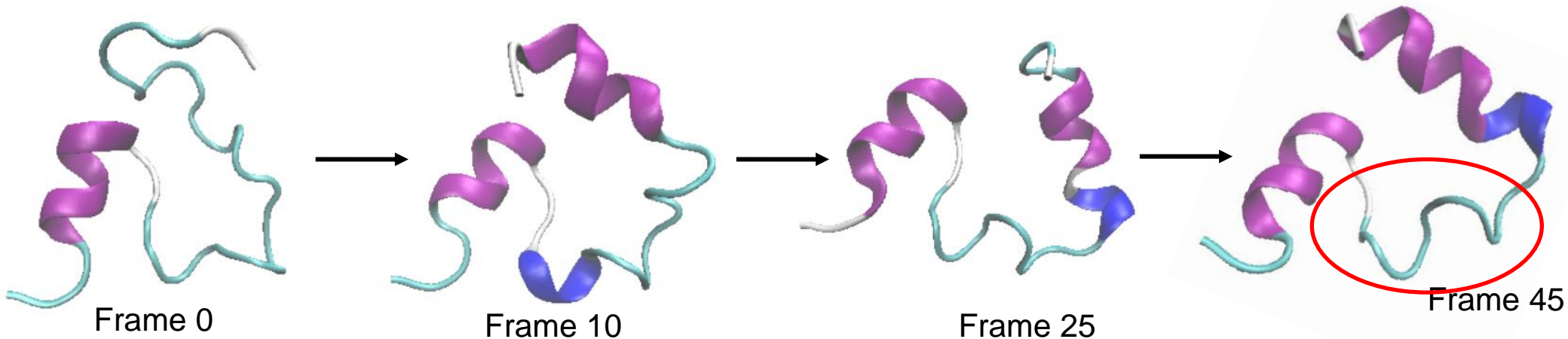
TrpCage (20 aa, 304 atoms): 12 traj, 4 folded





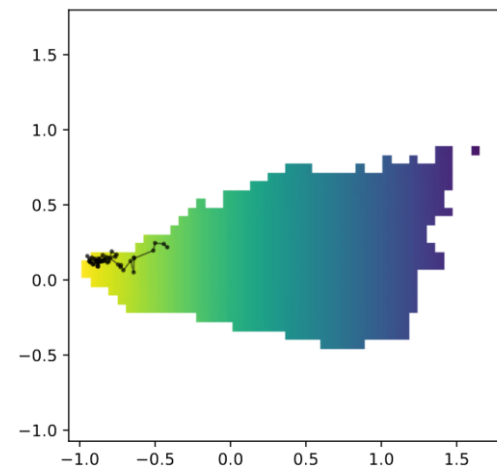
Villin (35 aa, 583 atoms): 16 traj, 4 “almost” folded



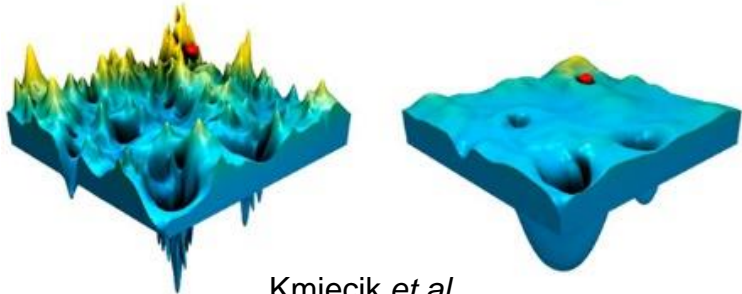
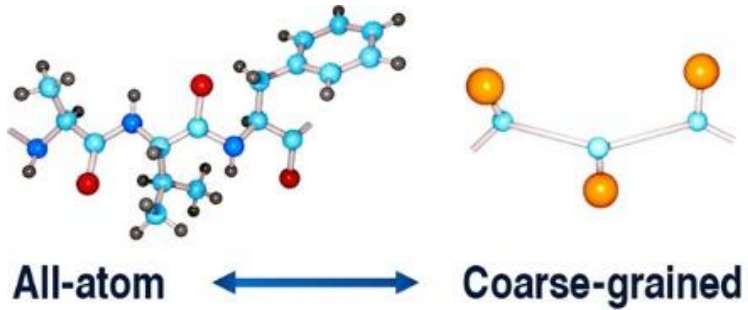


Villin (35 aa, 583 atoms): 16 traj, 4 “almost” folded

→ High energy barrier for the last helix?

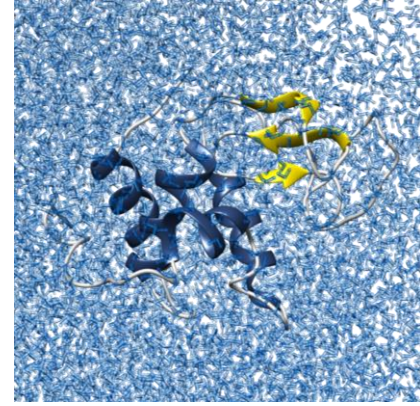
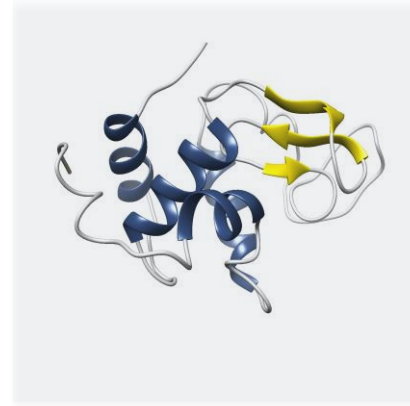


## Coarse graining

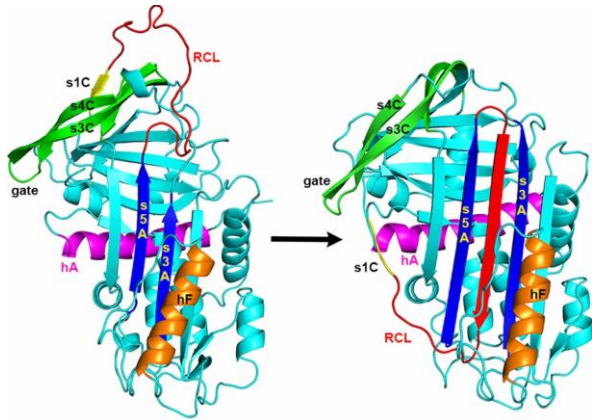


Kmiecik *et al*,  
*Chem. Rev.* 2016,

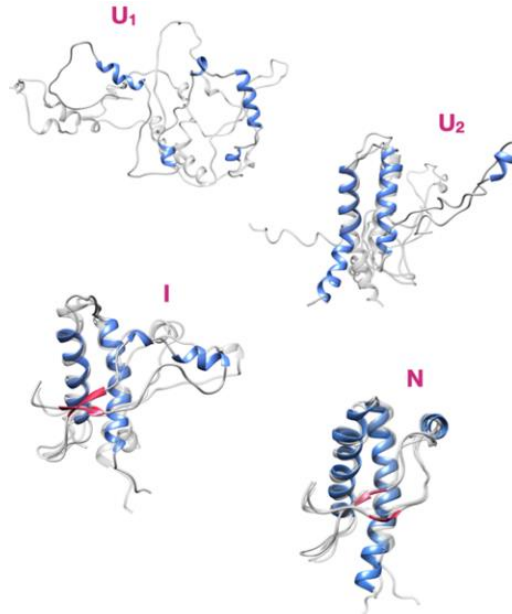
## Explicit solvent



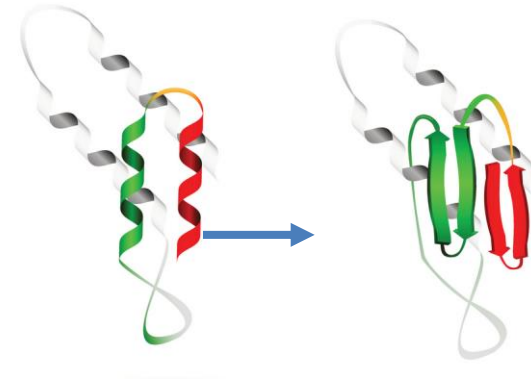
## Conformational transitions and point mutations

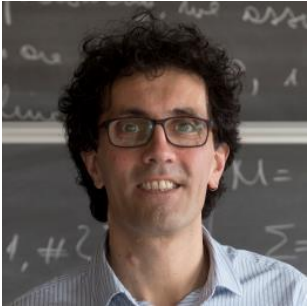


## Identify intermediate conformations



## Misfolding





**Prof. Diego Calvanese**  
Smart Data Factory  
University of Bolzano



**Prof. Emiliano Biasini**  
Collaborator  
University of Trento



**Prof. Pietro Faccioli**  
Master's supervisor  
University of Milan

