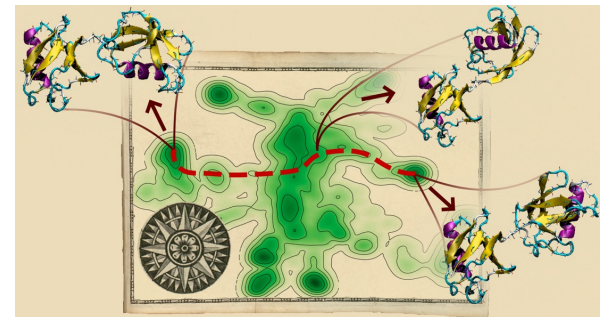


# Machine learning in biomolecular simulations: from characterizing conformational free energy landscapes to scale bridging

Christine Peter

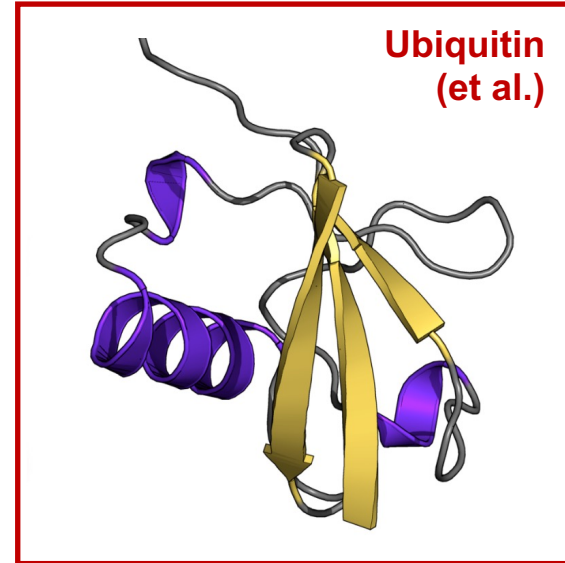
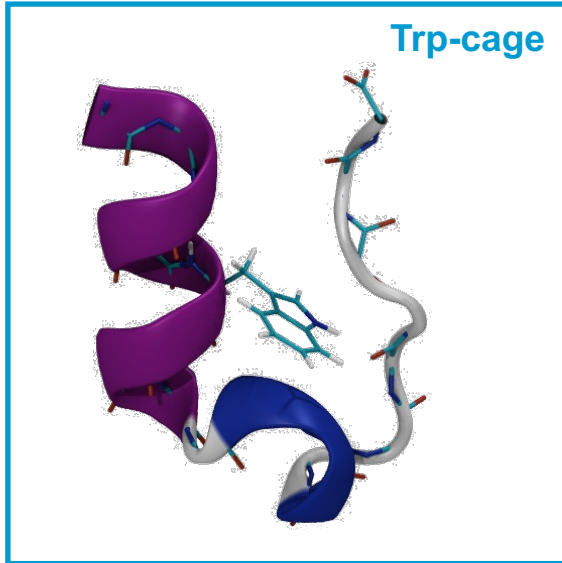
ECT\* Workshop, Trento 2024



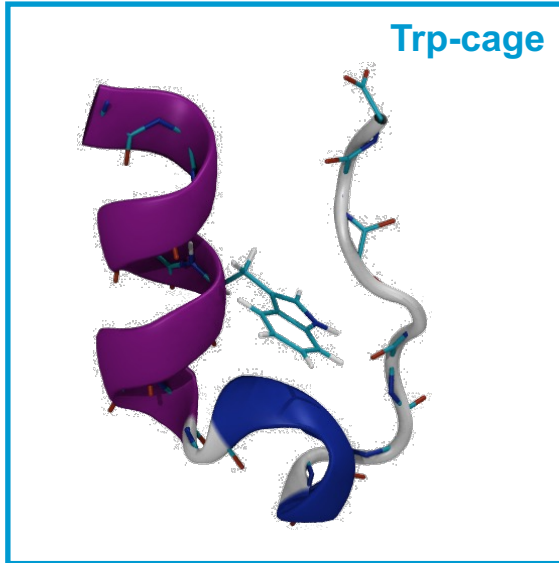
# Outline

- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations

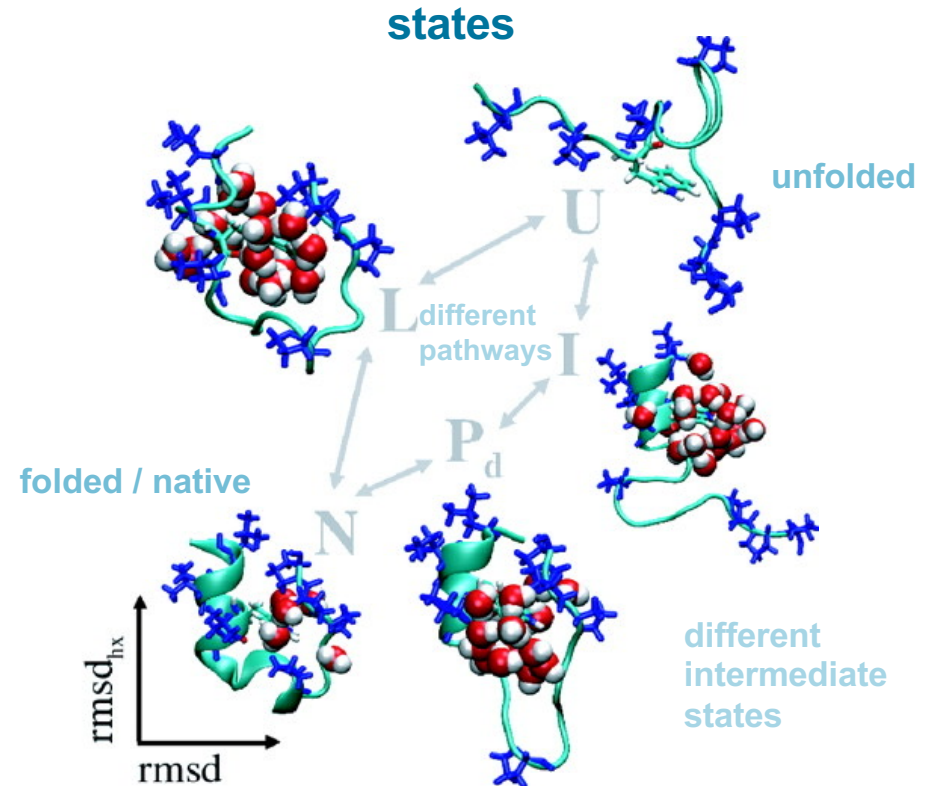
## Two recurring systems



# Two recurring systems

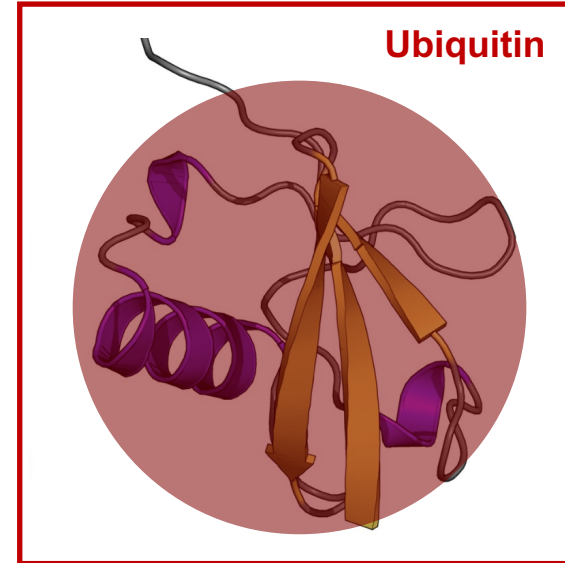
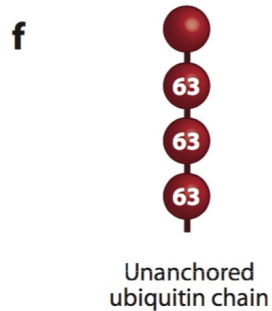
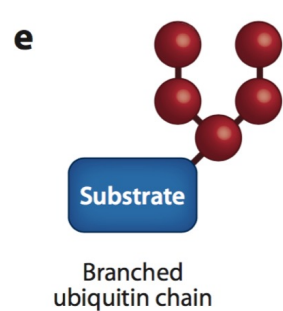
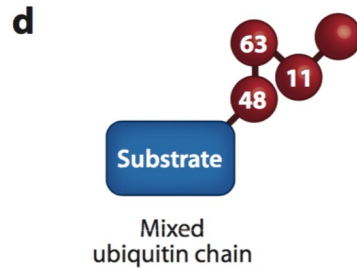
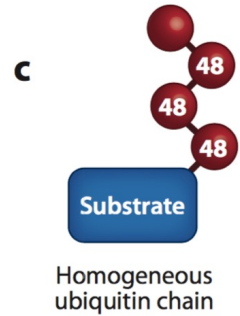
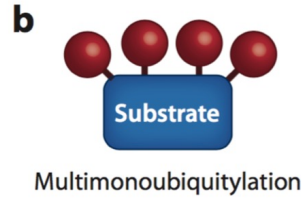
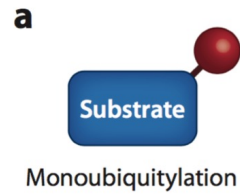


- folding
- states
- features & CVs
- low-dimensional representations



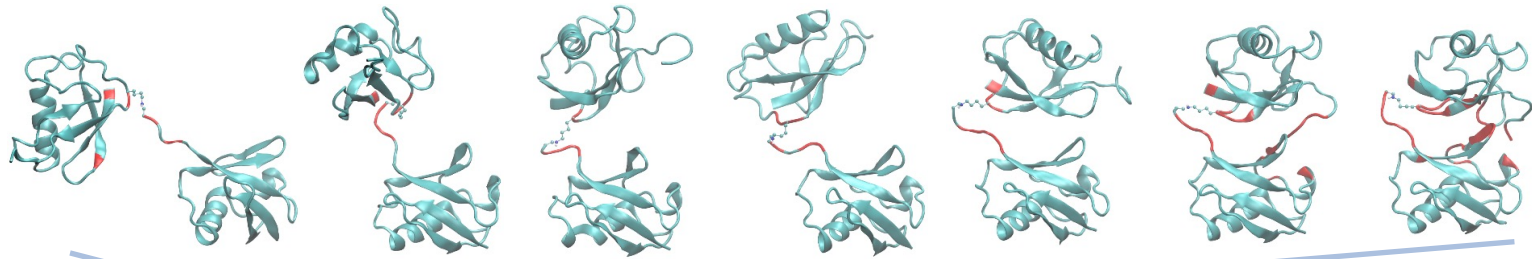
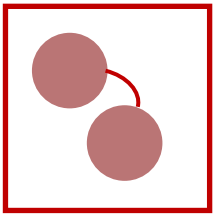
Juraszek, J. and P. G. Bolhuis.  
*PNAS* 103.43 (2006): 15859-15864.

# Two recurring systems: Trp-cage and Ubiquitin et al.

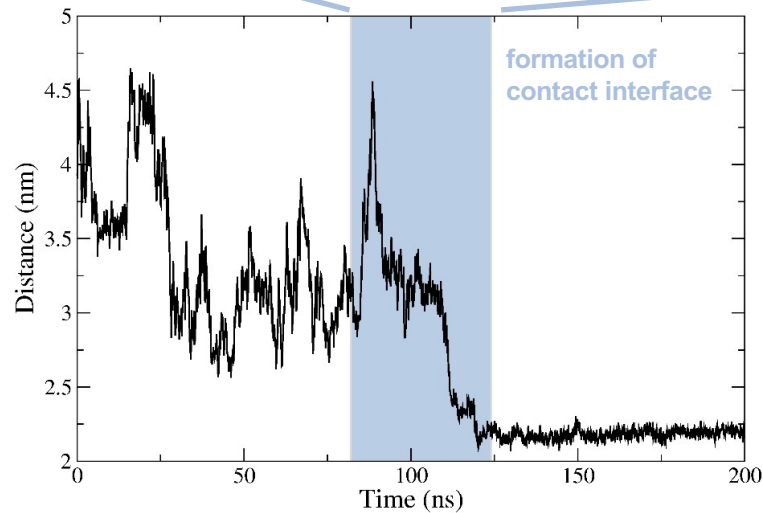


- protein-protein interactions
- protein interfaces
- sampling / scale-bridging

# An atomistic simulation of K48-linked di-Ubiquitin

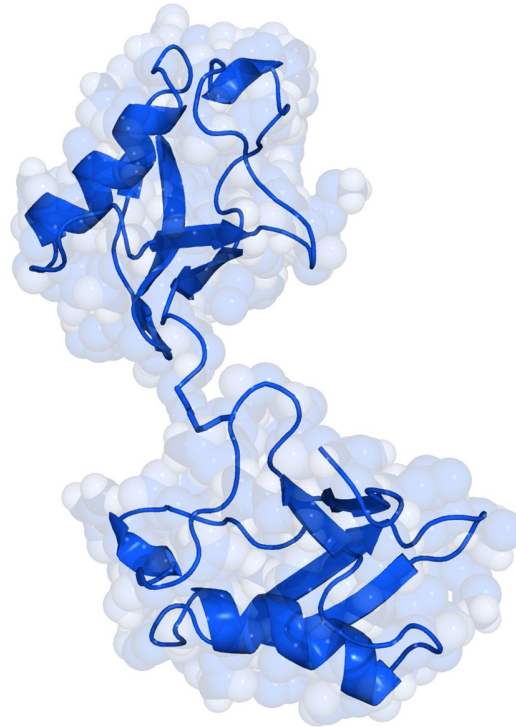
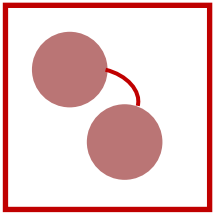


protein interfaces



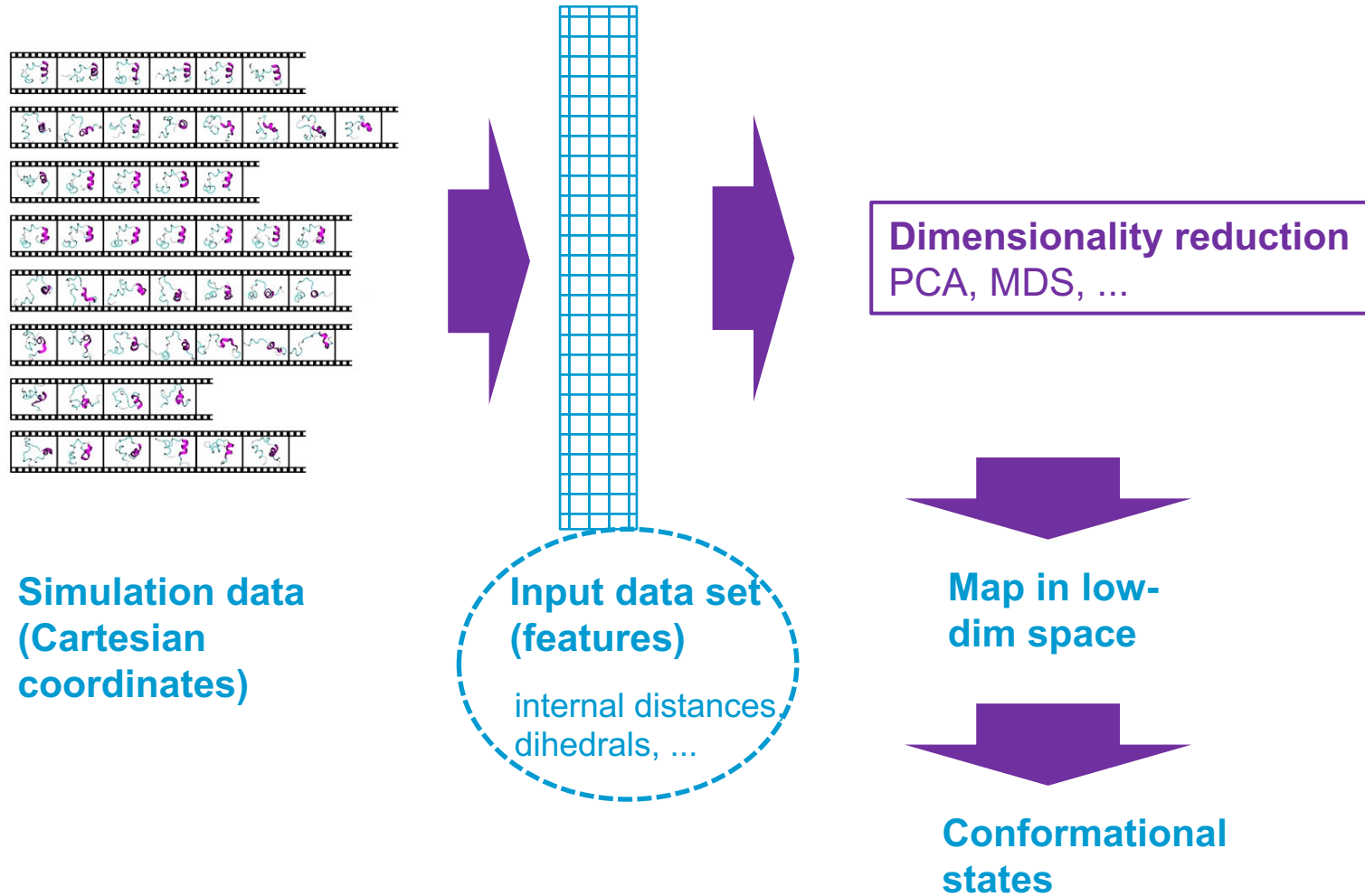
the accessible timescales are too short  
→ need for enhanced sampling

# A coarse grained simulation of K48-linked di-Ubiquitin



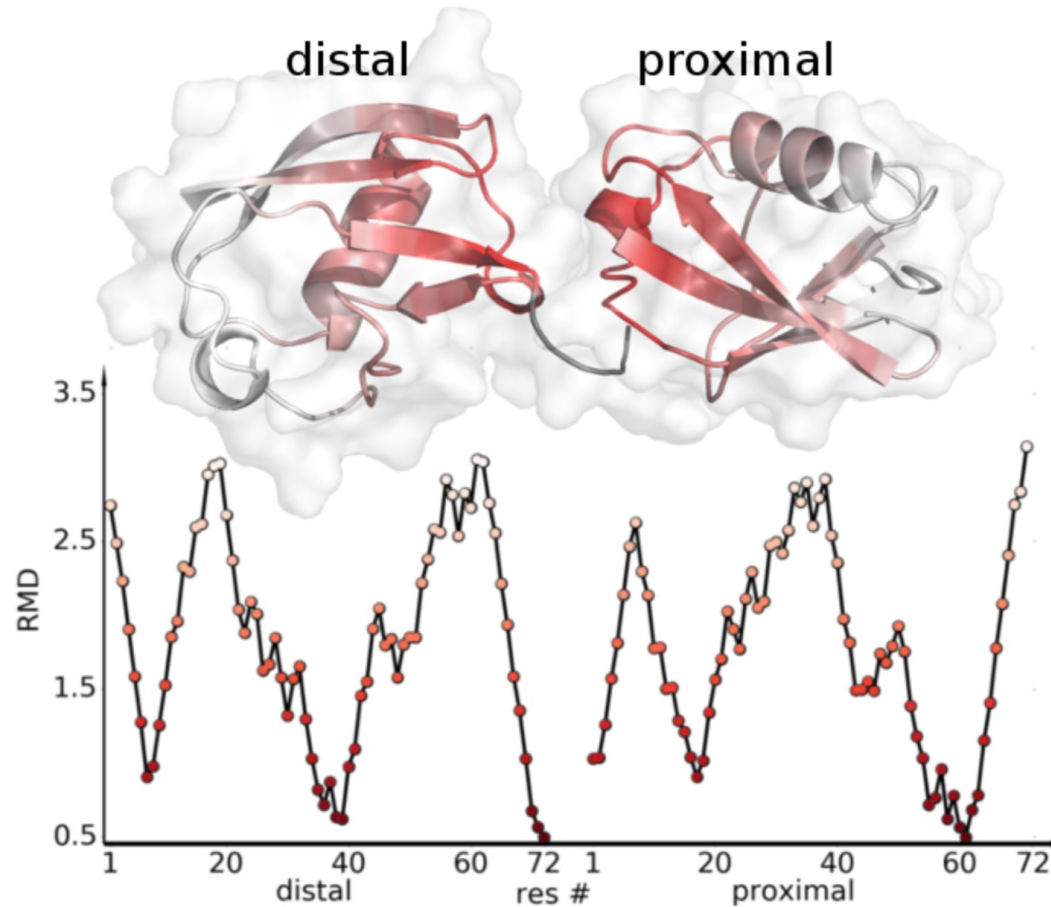
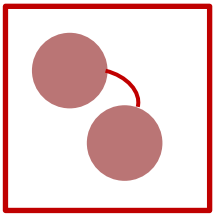
→ Need for methods to characterize conformational phase space

# Low-dimensional representations of conformational phase space



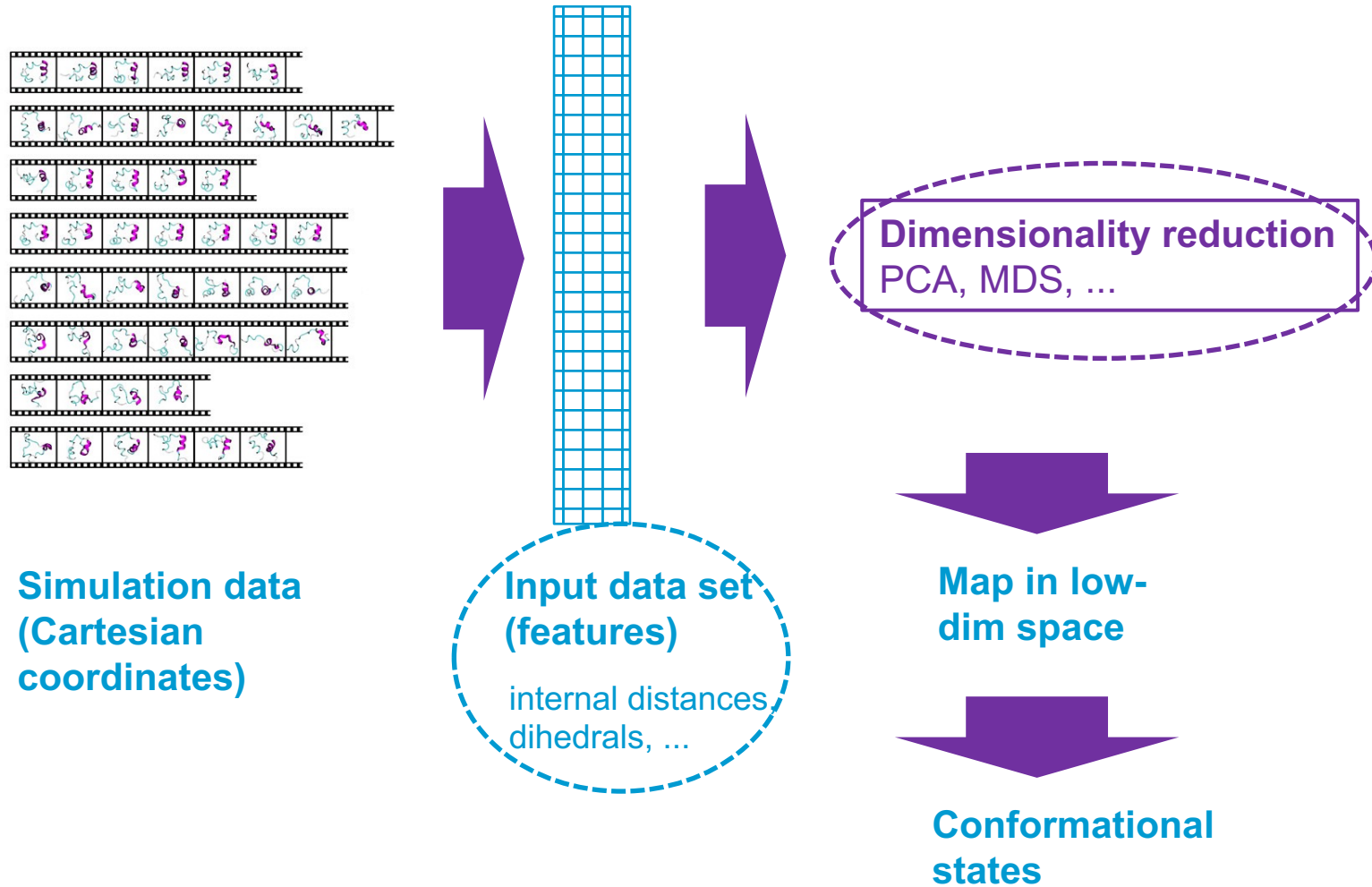


# Collective variables that characterize dimer structures



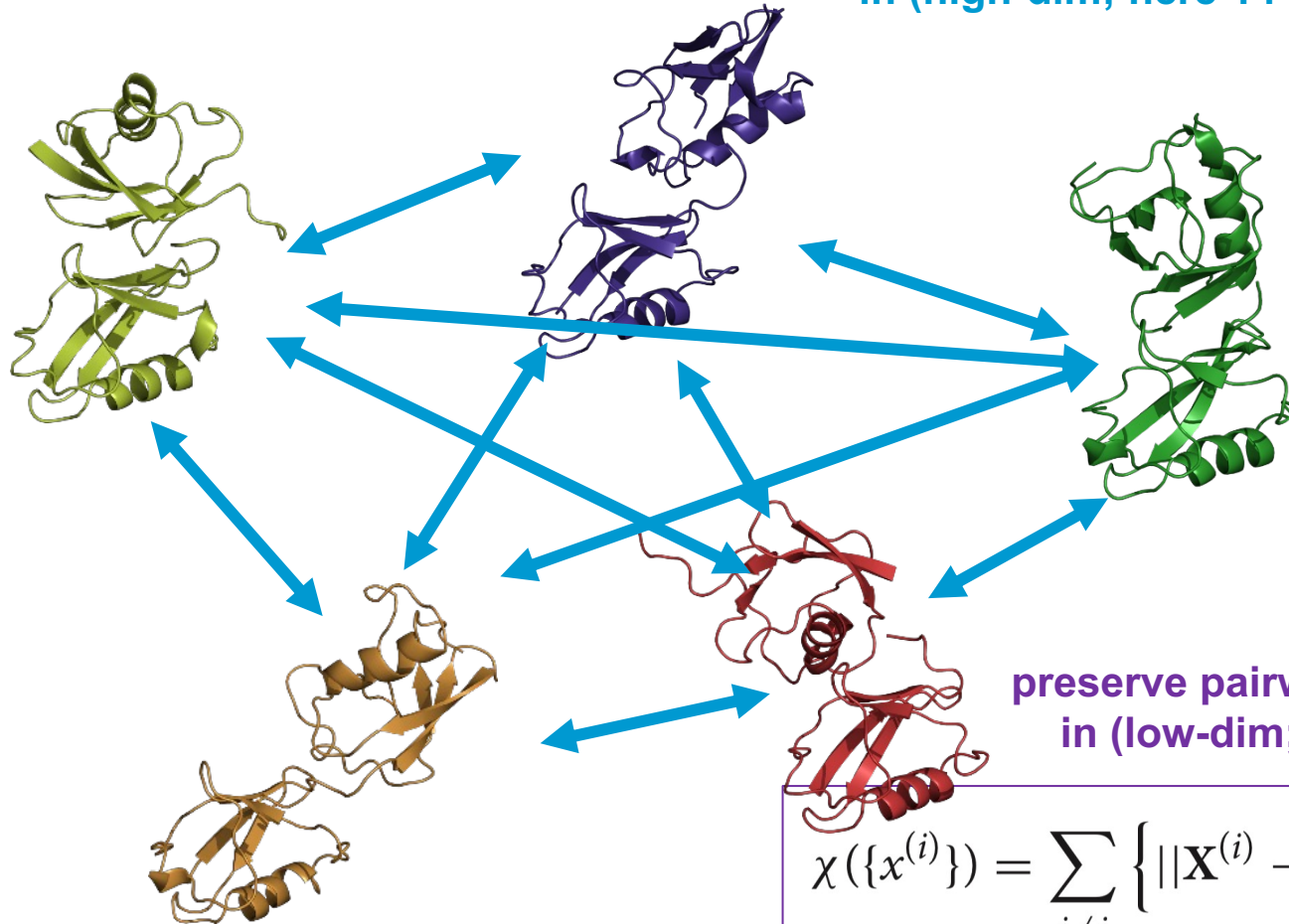
- 144 residue-wise minimum distances (RMD)
- fingerprints of each structure

# Low-dimensional representations of conformational phase space



# The idea behind metric multidimensional scaling

pairwise differences/distances between structures  
in (high-dim; here 144-dim) feature space



$$\chi(\{x^{(i)}\}) = \sum_{i \neq j} \left\{ \|\mathbf{X}^{(i)} - \mathbf{X}^{(j)}\| - \|x^{(i)} - x^{(j)}\| \right\}^2$$

# The idea behind SketchMap (Ceriotti et al J. Chem. Theory Comput. 2013)

pairwise differences/distances between structures

distances in high- and low-dimensional spaces are transformed by a **sigmoid function** (44-dim) feature space

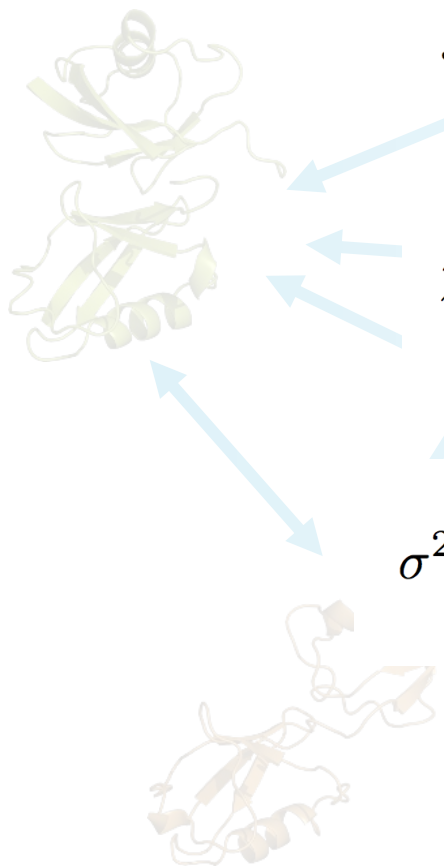
$$s(r, a, b) = 1 - (1 + (2^{a/b} - 1)(r/\sigma)^a)^{-b/a}$$

produce mapping by minimizing stress function for **landmark** structures:

$$\chi^2 = \sum_{i \neq j} [s(R_{ij}, A, B) - s(r_{ij}, a, b)]^2$$

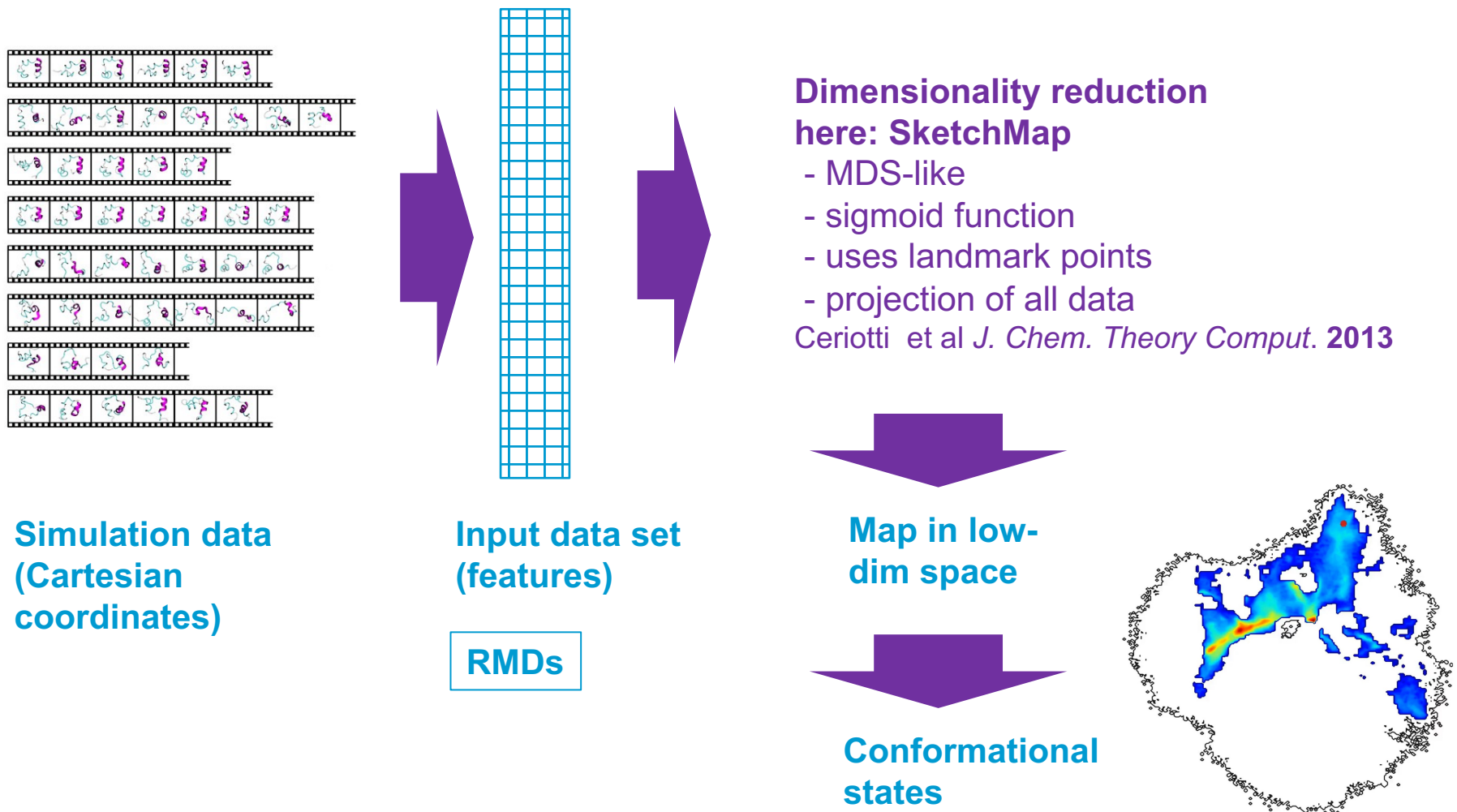
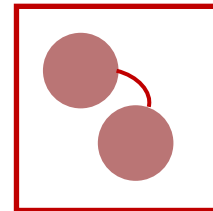
calculate projection  $\mathbf{x}$  of any high-dimensional point  $\mathbf{X}$  by minimizing:

$$\sigma^2(\mathbf{x}) = \sum_{i=1}^N [s(R_i(\mathbf{X}), A, B) - s(r_i(\mathbf{x}), a, b)]^2$$

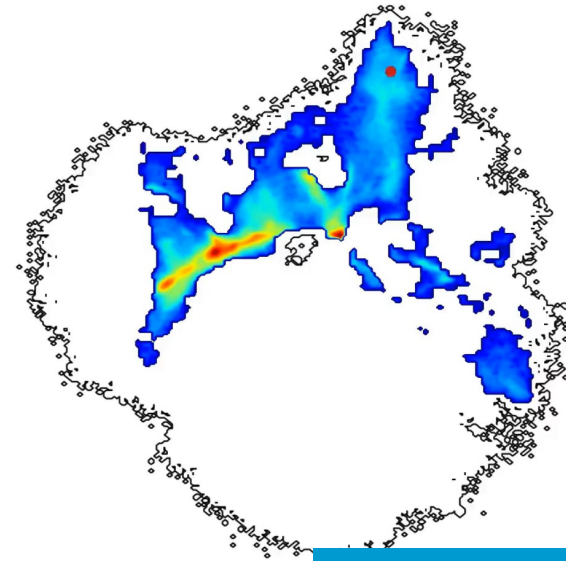
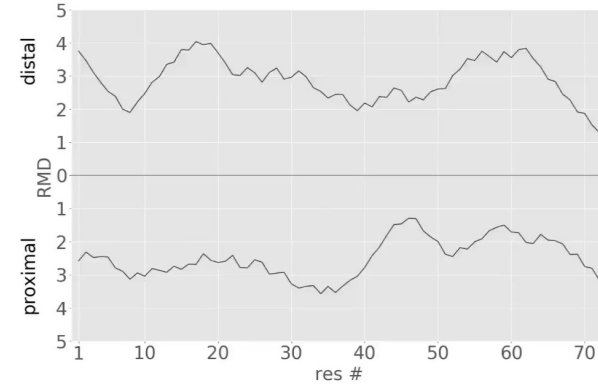
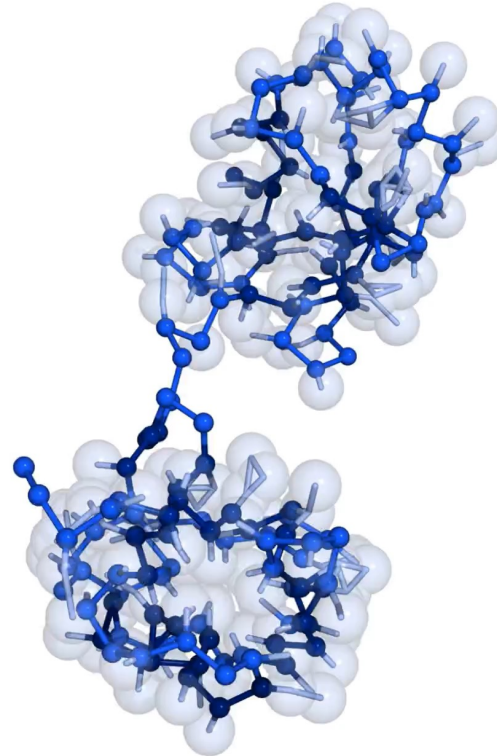
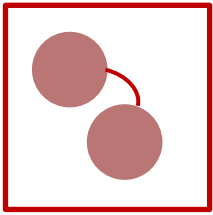


$$\chi(\{x^{(i)}\}) = \sum_{i \neq j} \left\{ \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\| - \|x^{(i)} - x^{(j)}\| \right\}^2$$

# Low-dimensional representations of conformational phase space



## Residue-wise minimum distances as input features: Illustration for one simulation

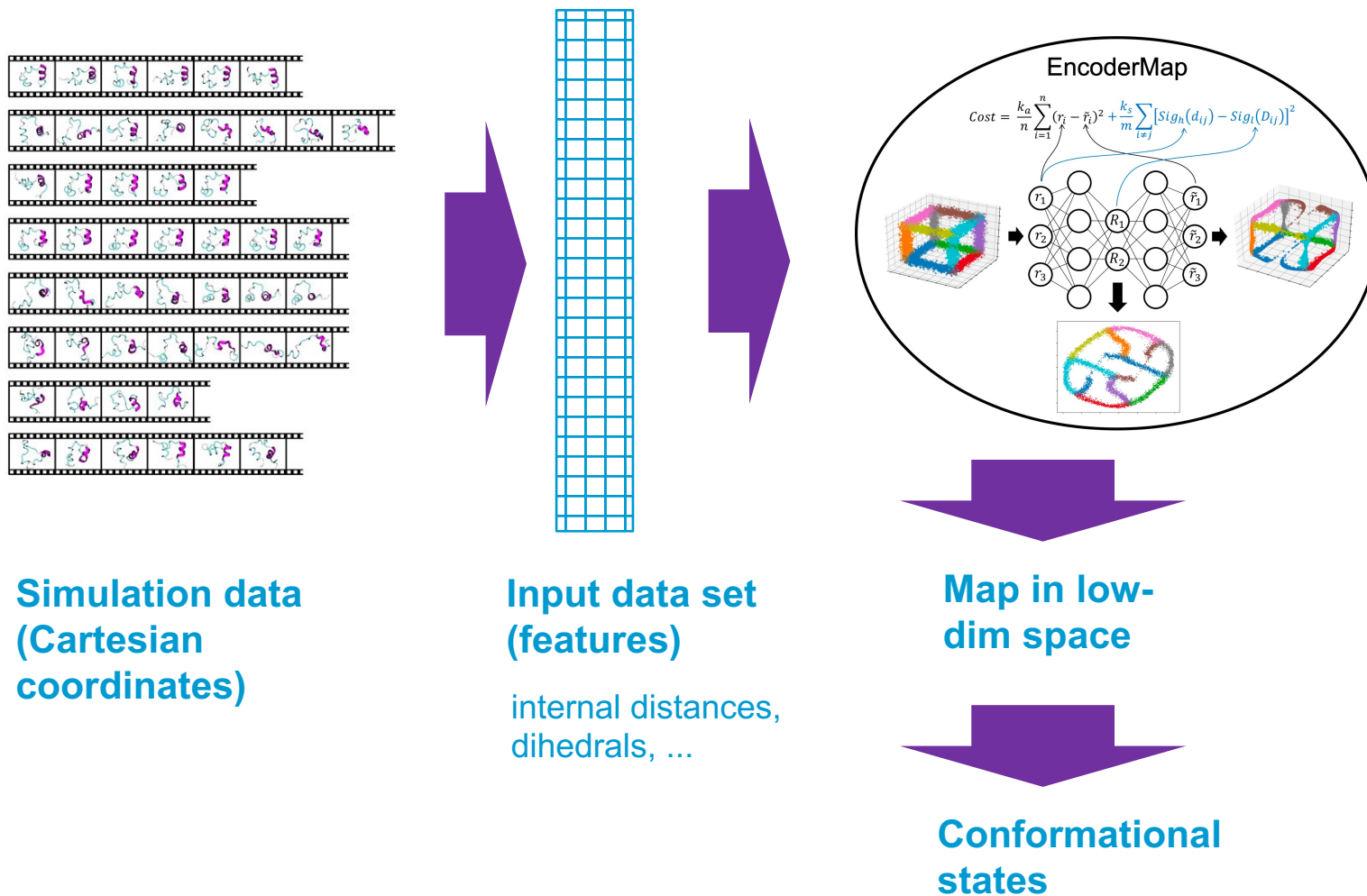


→ need for methods to process huge datasets

# Outline

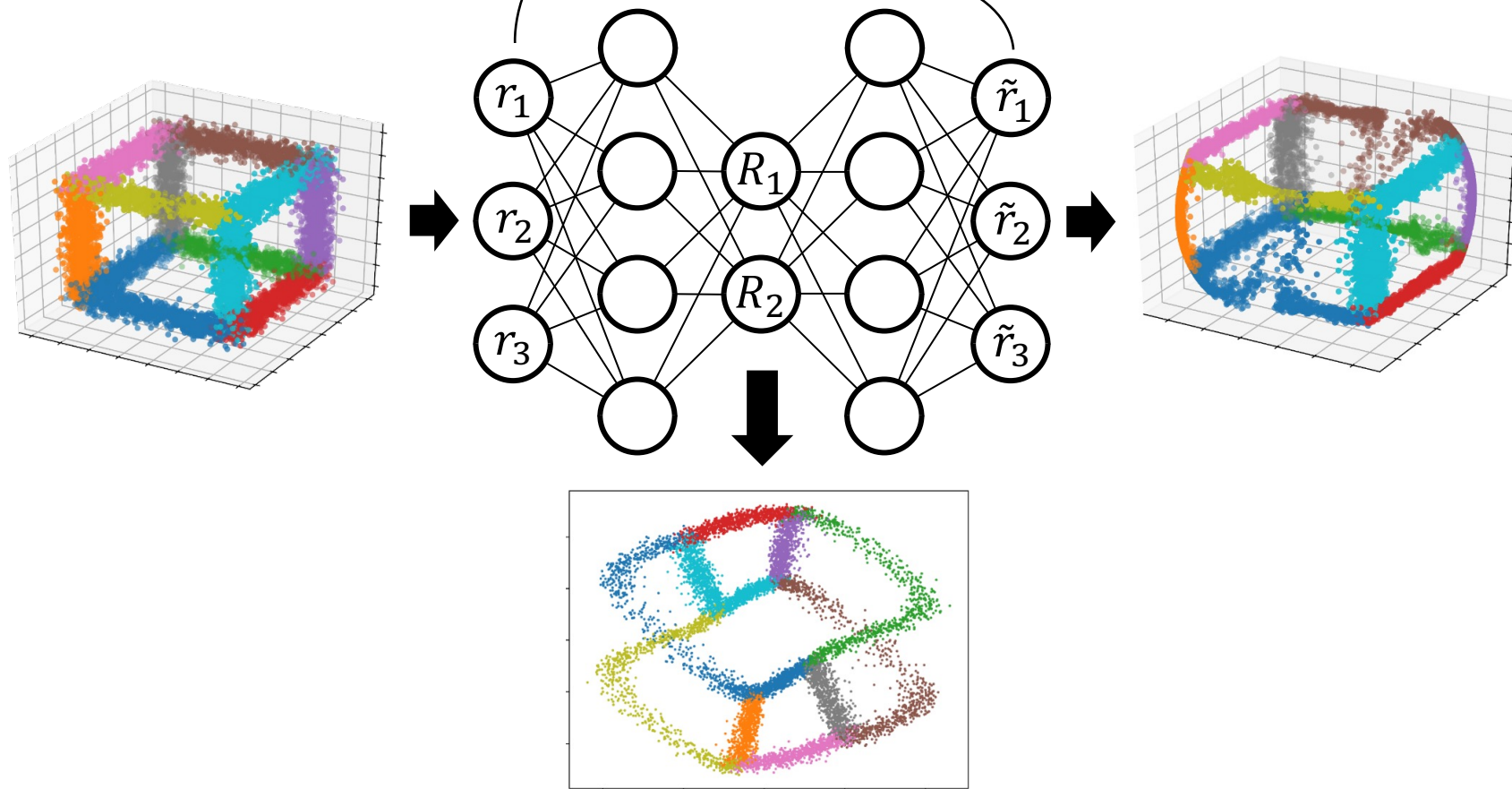
- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations

# Low-dimensional representations of conformational phase space

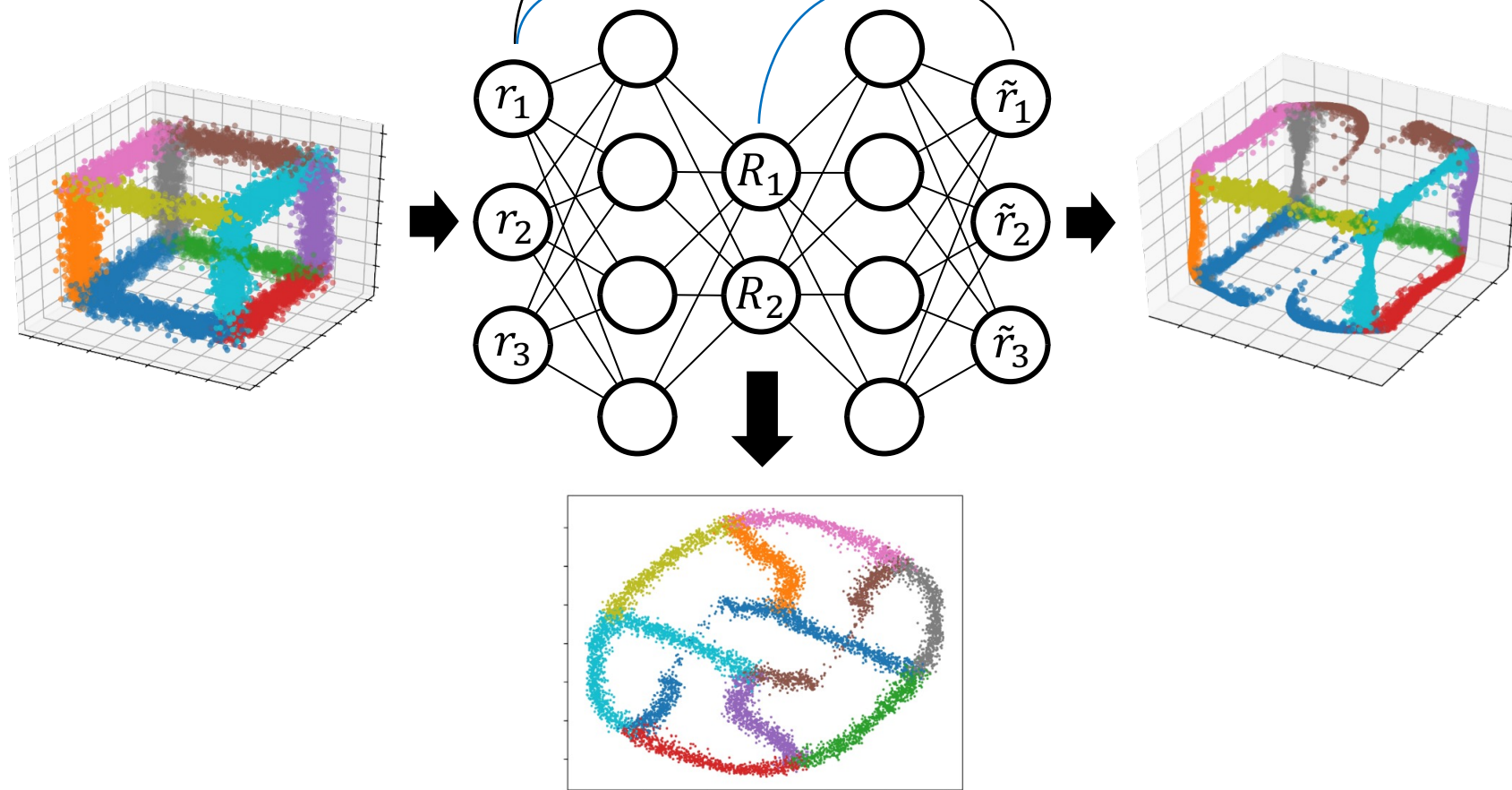




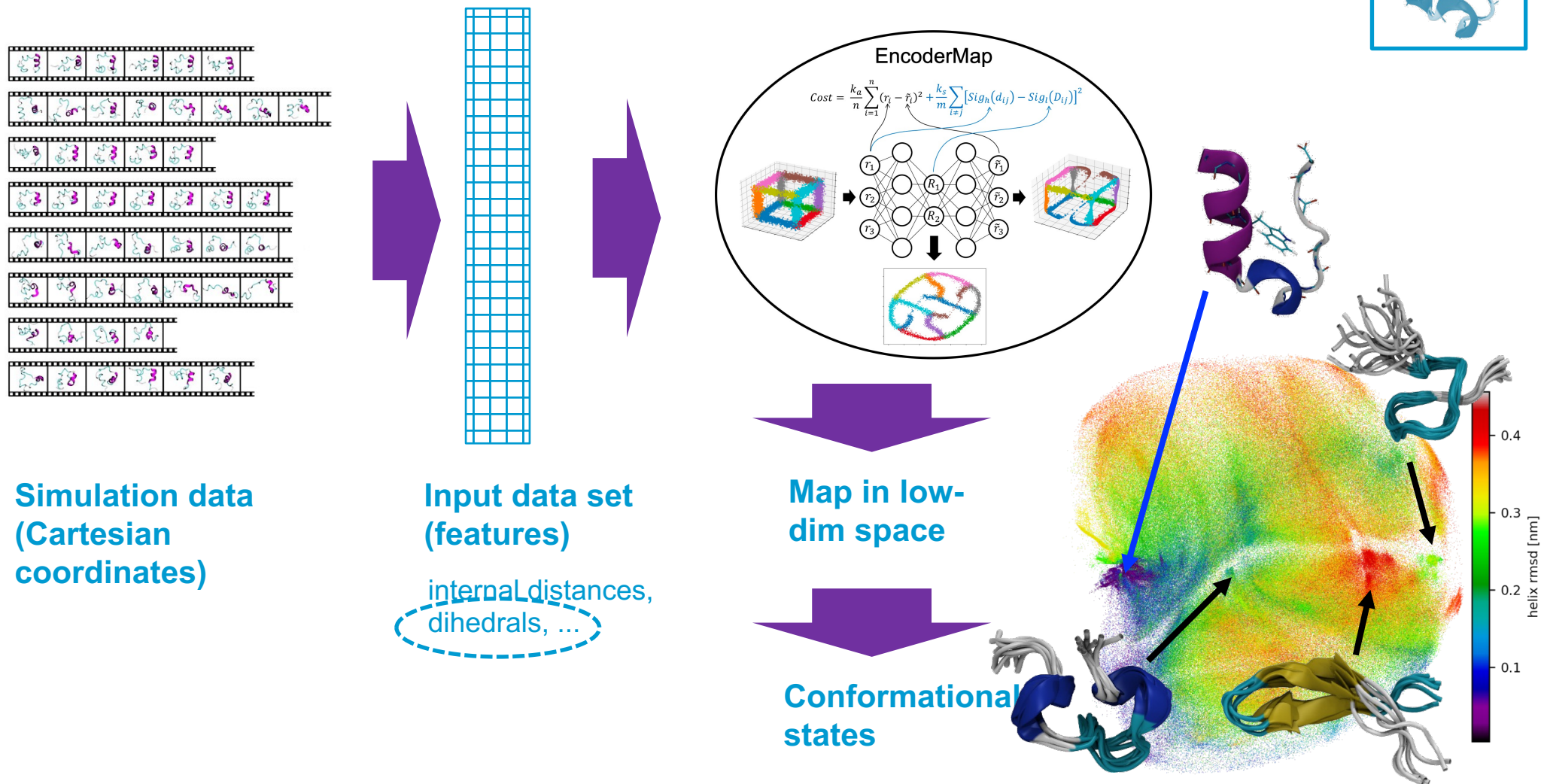
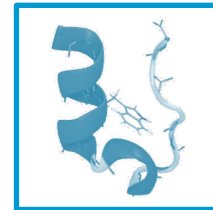
$$Cost = \frac{k_a}{n} \sum_{i=1}^n (r_i - \tilde{r}_i)^2$$



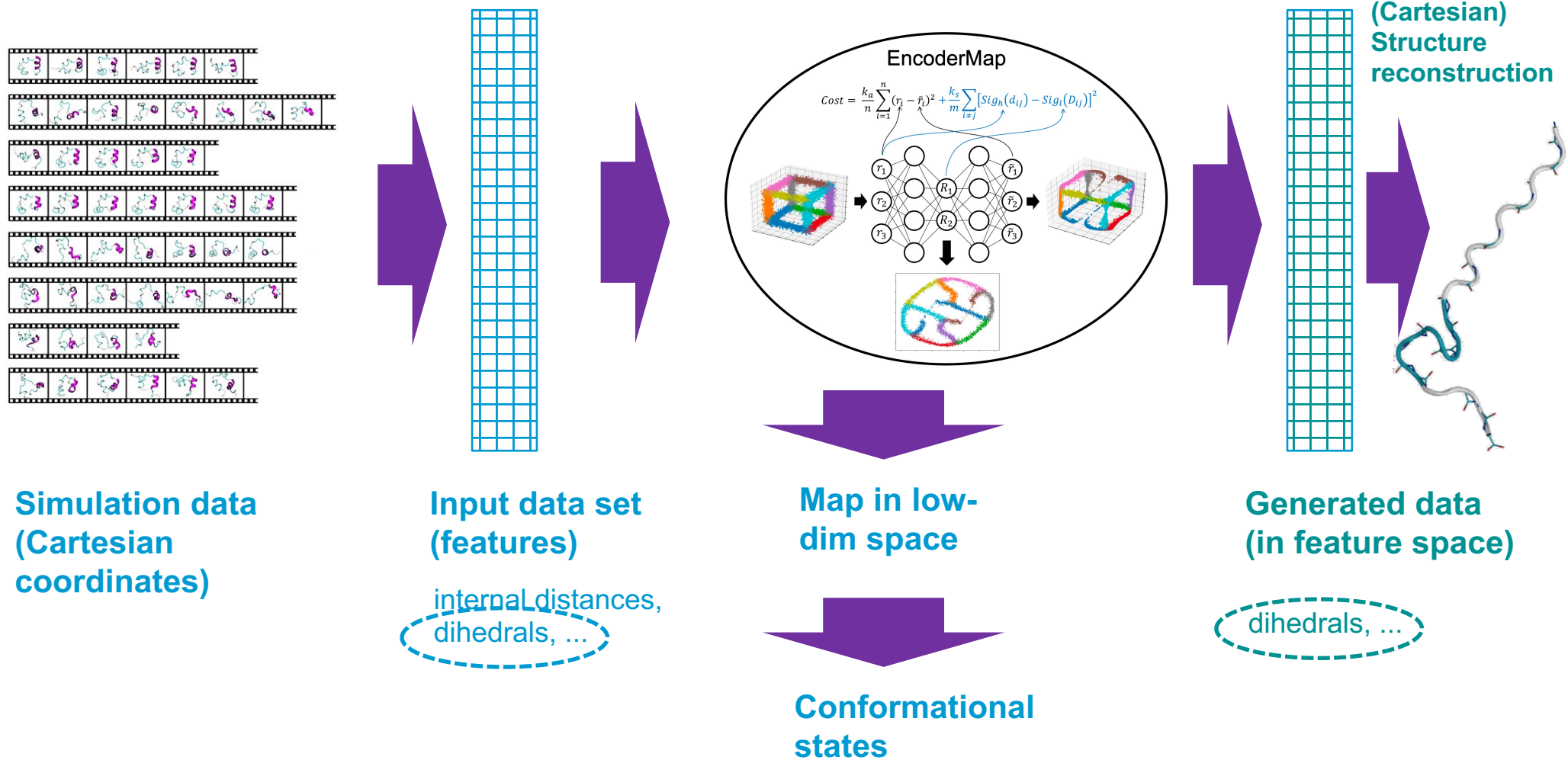
$$Cost = \frac{k_a}{n} \sum_{i=1}^n (r_i - \tilde{r}_i)^2 + \frac{k_s}{m} \sum_{i \neq j} [Sig_h(d_{ij}) - Sig_l(D_{ij})]^2$$

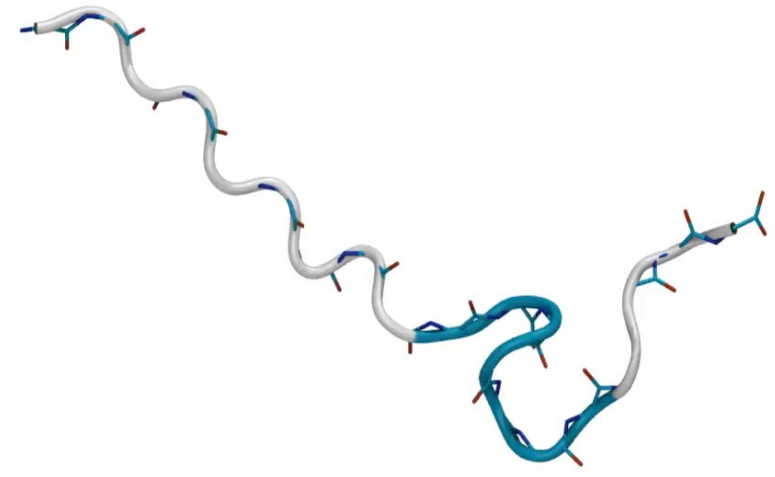
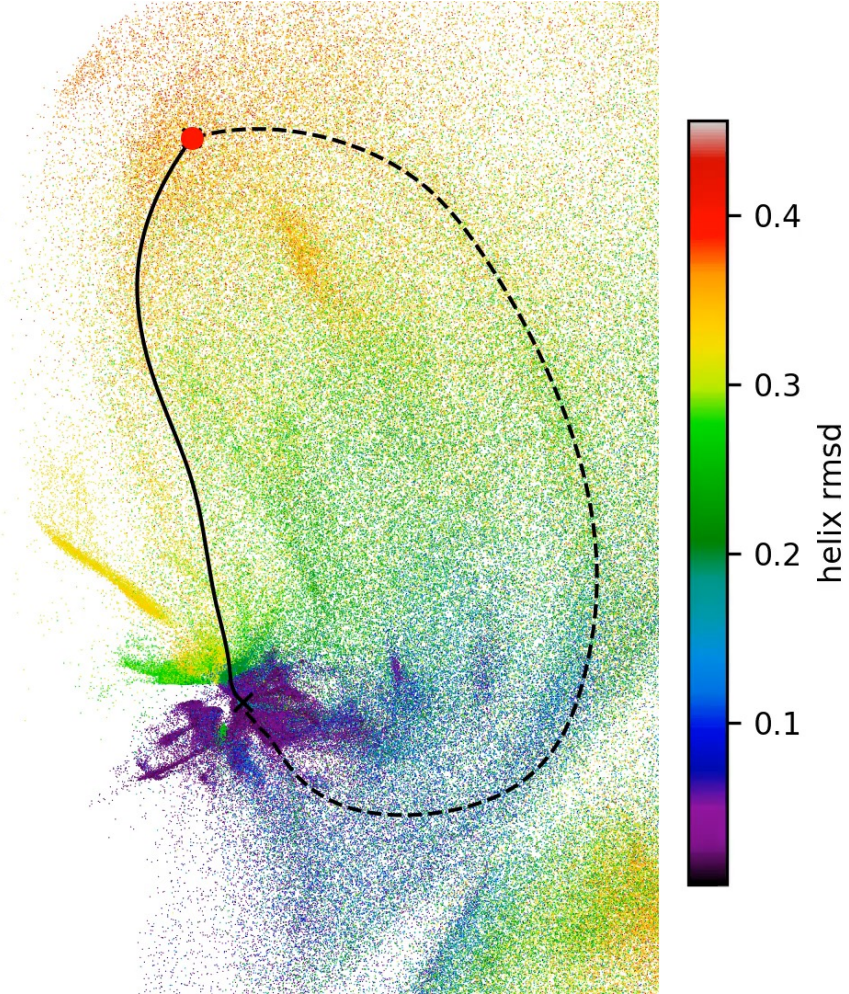


# EncoderMap example 1: Trp-cage

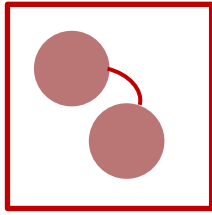
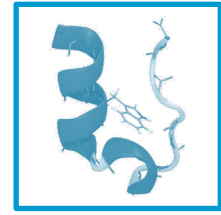
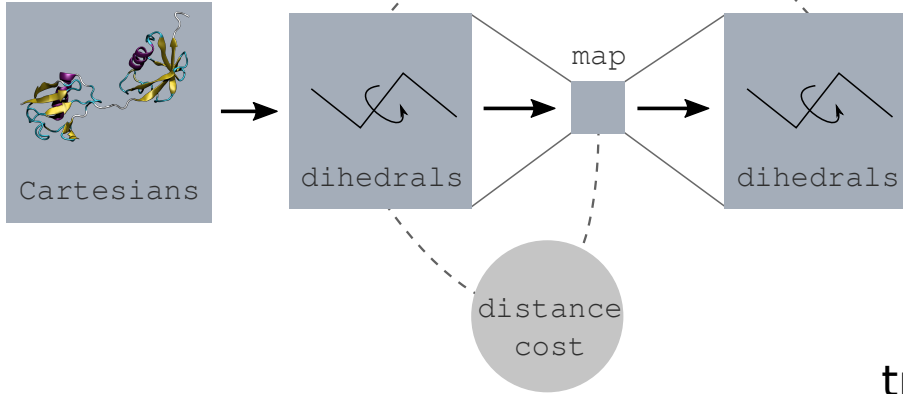


# EncoderMap example 1: Trp-cage

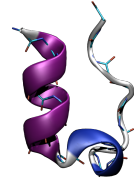




But ...

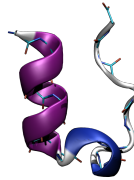


trp-cage

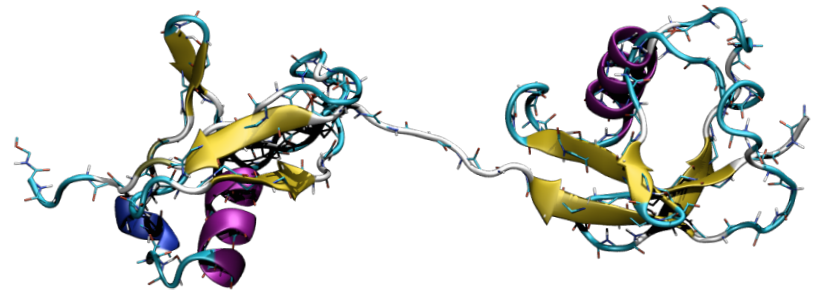
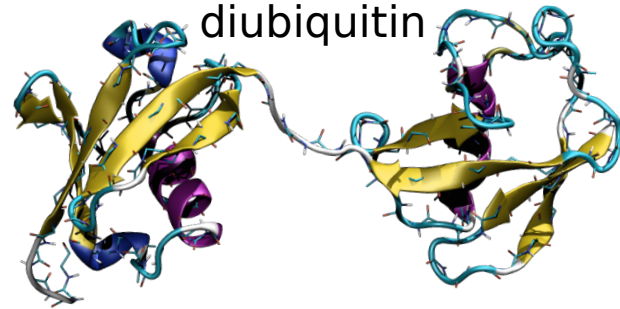


original

generated

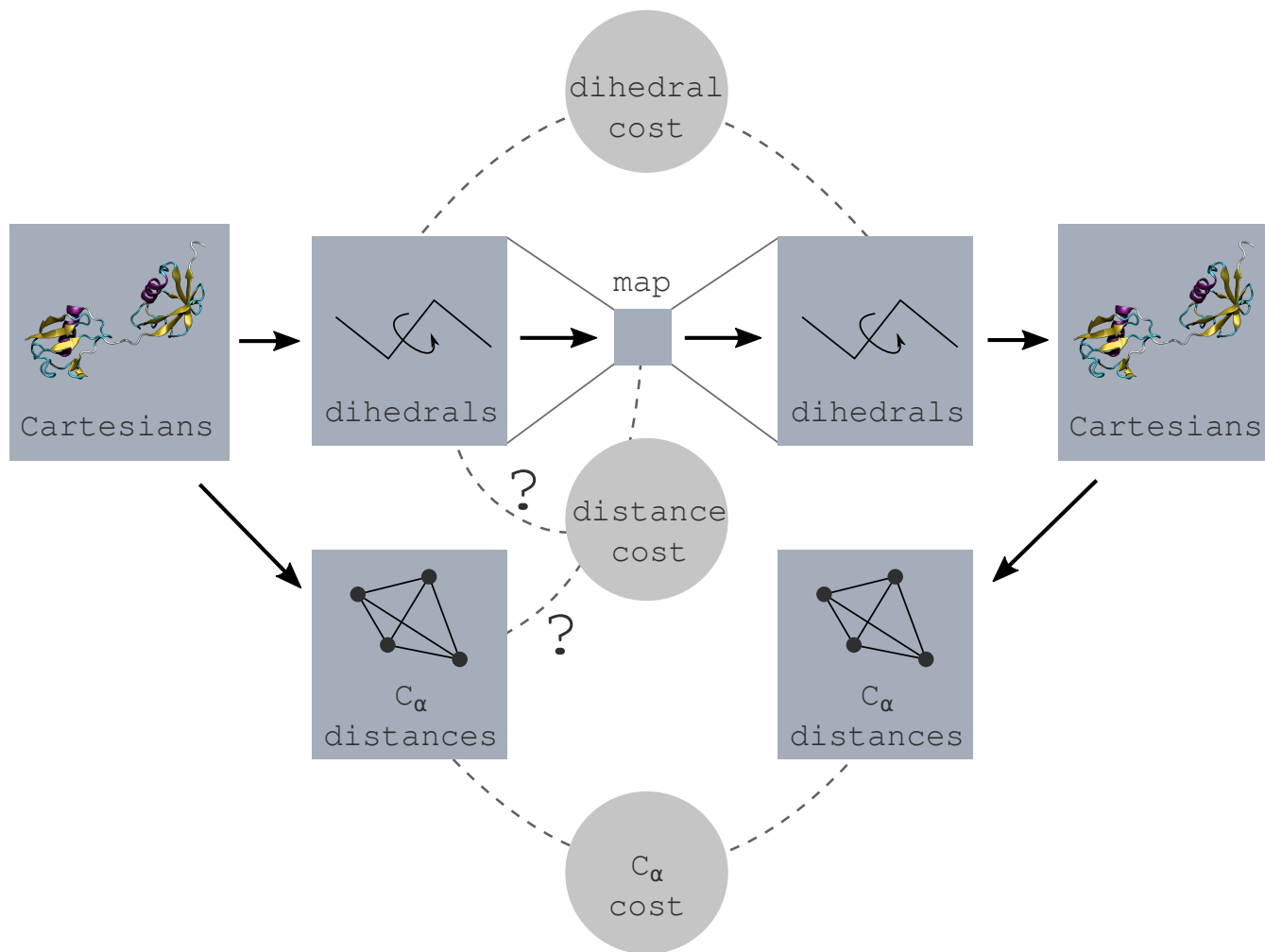


diubiquitin

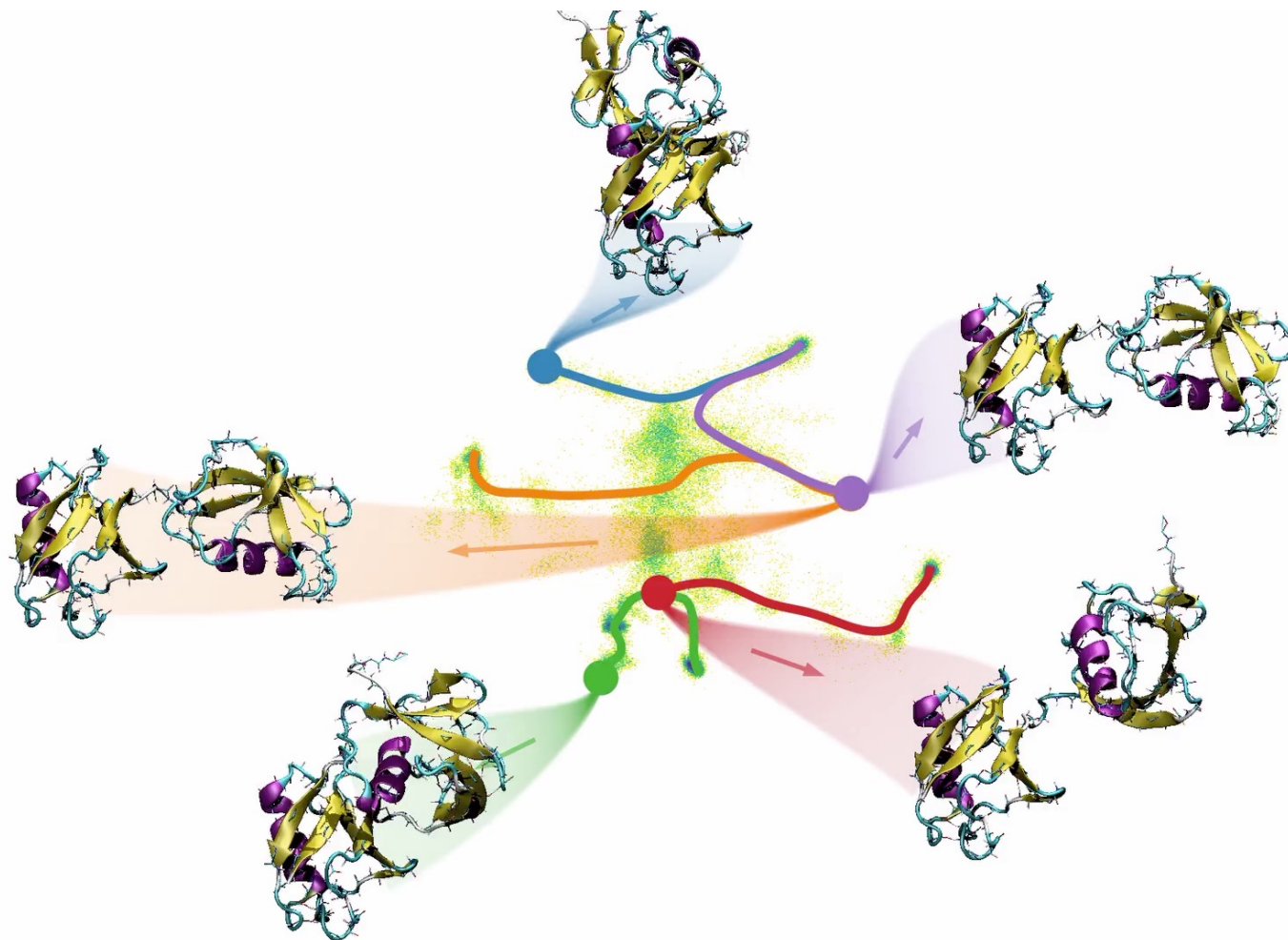
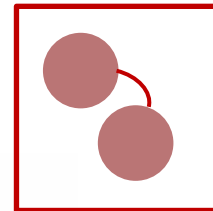


## EncoderMap (II)

$$C = k_{\text{dih.}} C_{\text{dihedral}} + k_{C_{\alpha}} C_{C_{\alpha}} + k_{\text{dist.}} C_{\text{distance}}$$

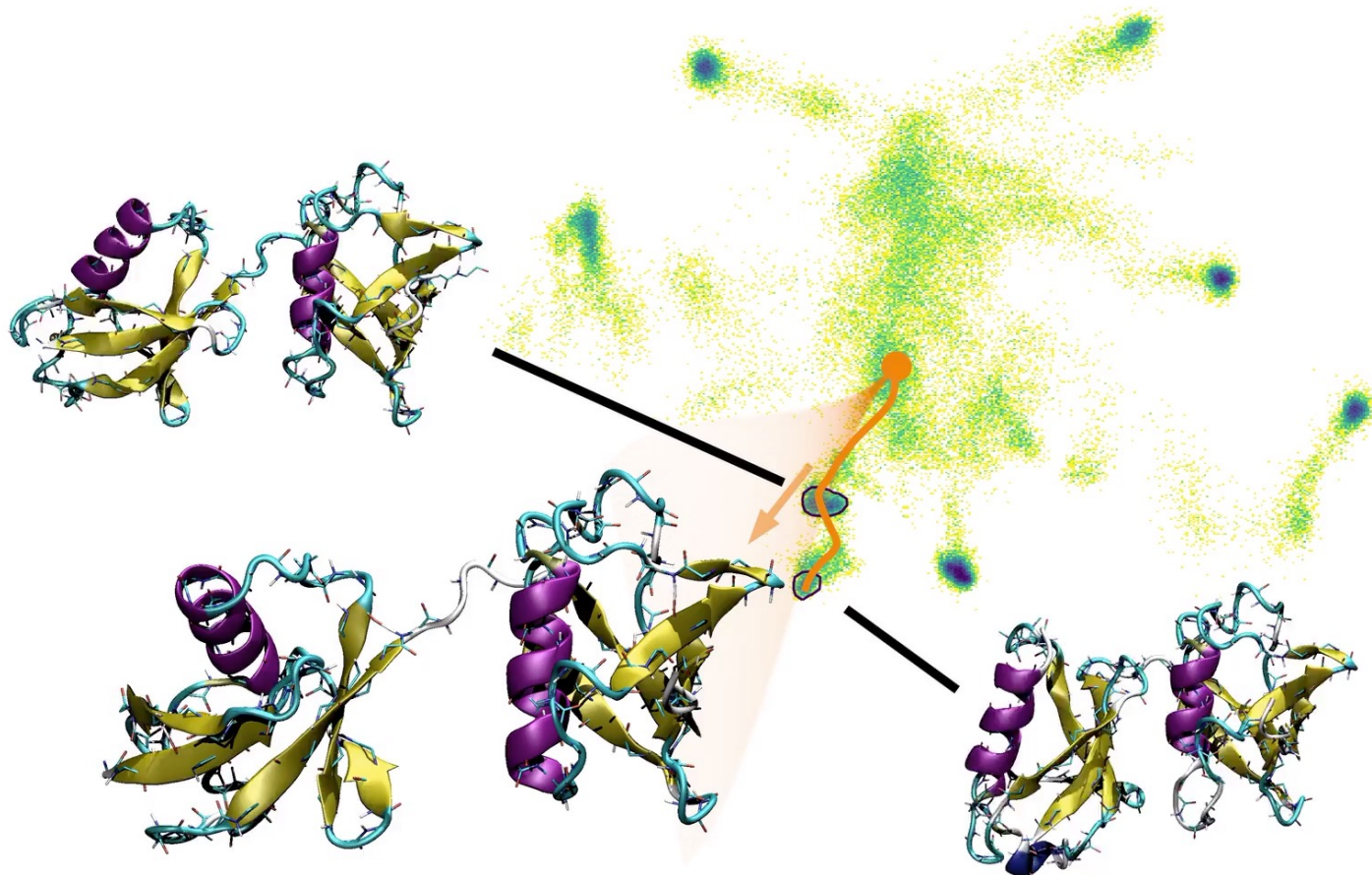
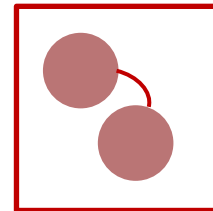


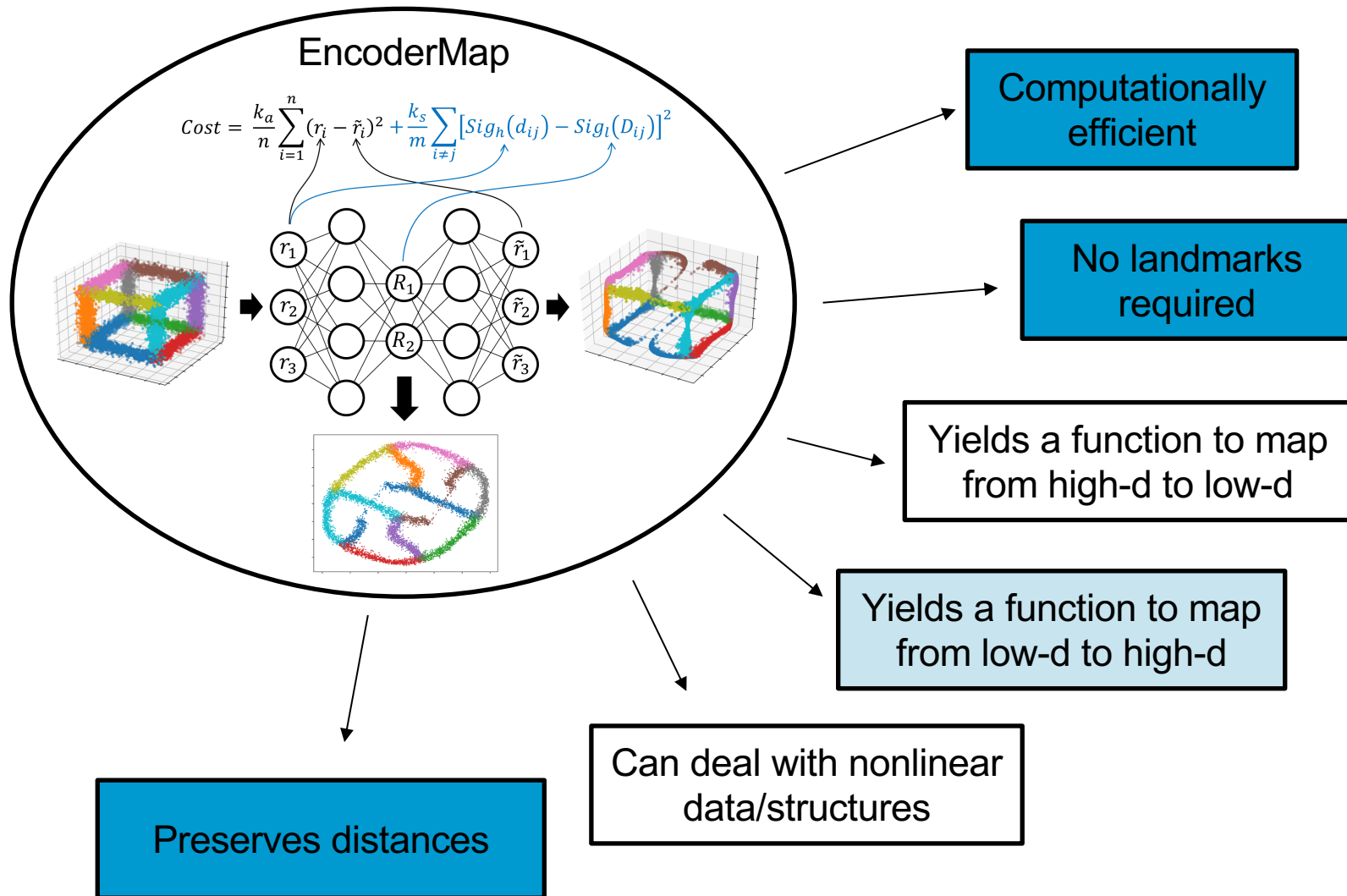
## Example 2: M1-linked di-Ubiquitin – generating paths





## Example 2: M1-linked di-Ubiquitin – visualizing important motions





→ EncoderMap: combines advantages of a NN autoencoder and SketchMap/MDS

# Outline

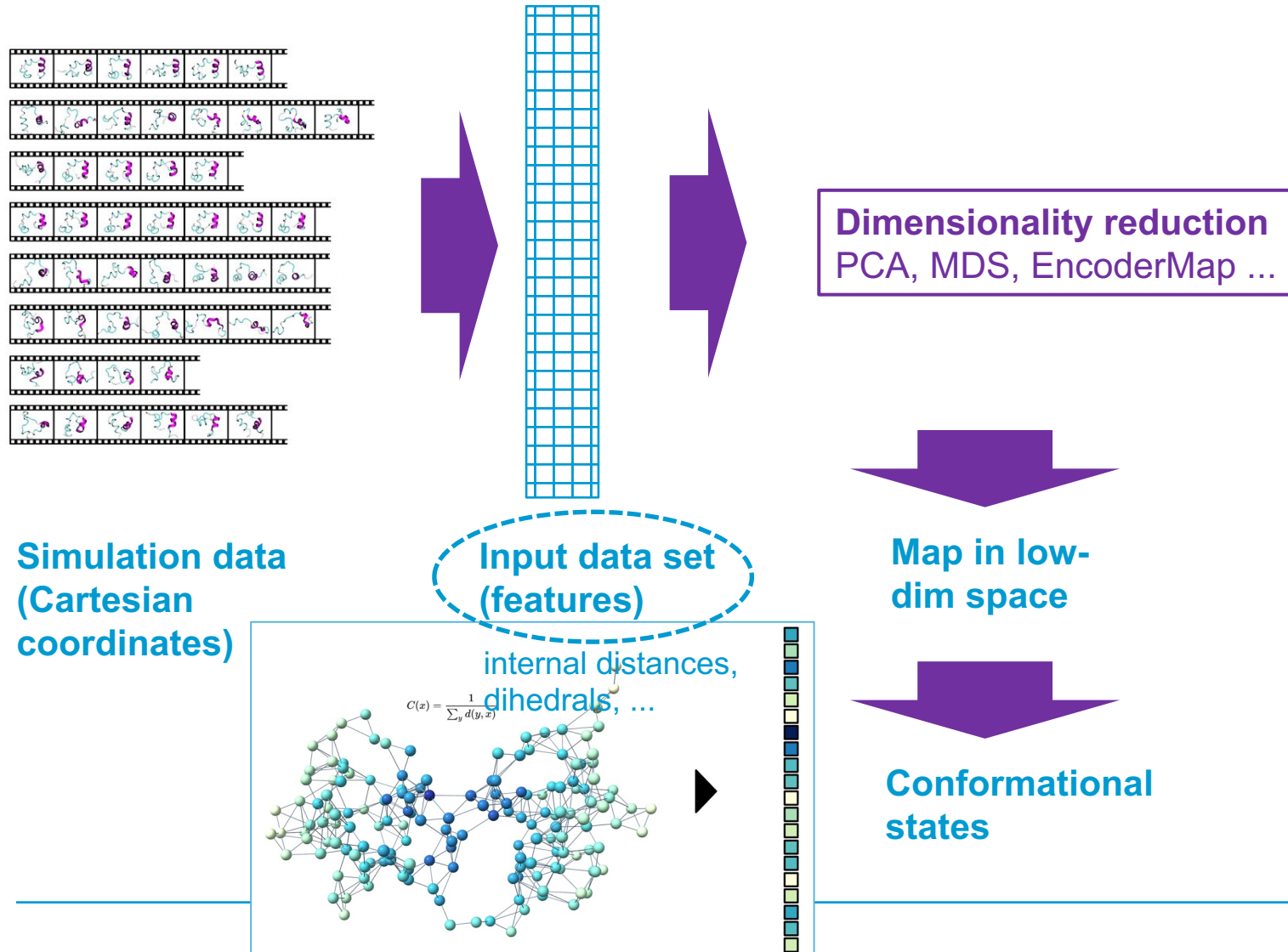
- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations

→ the input features matter

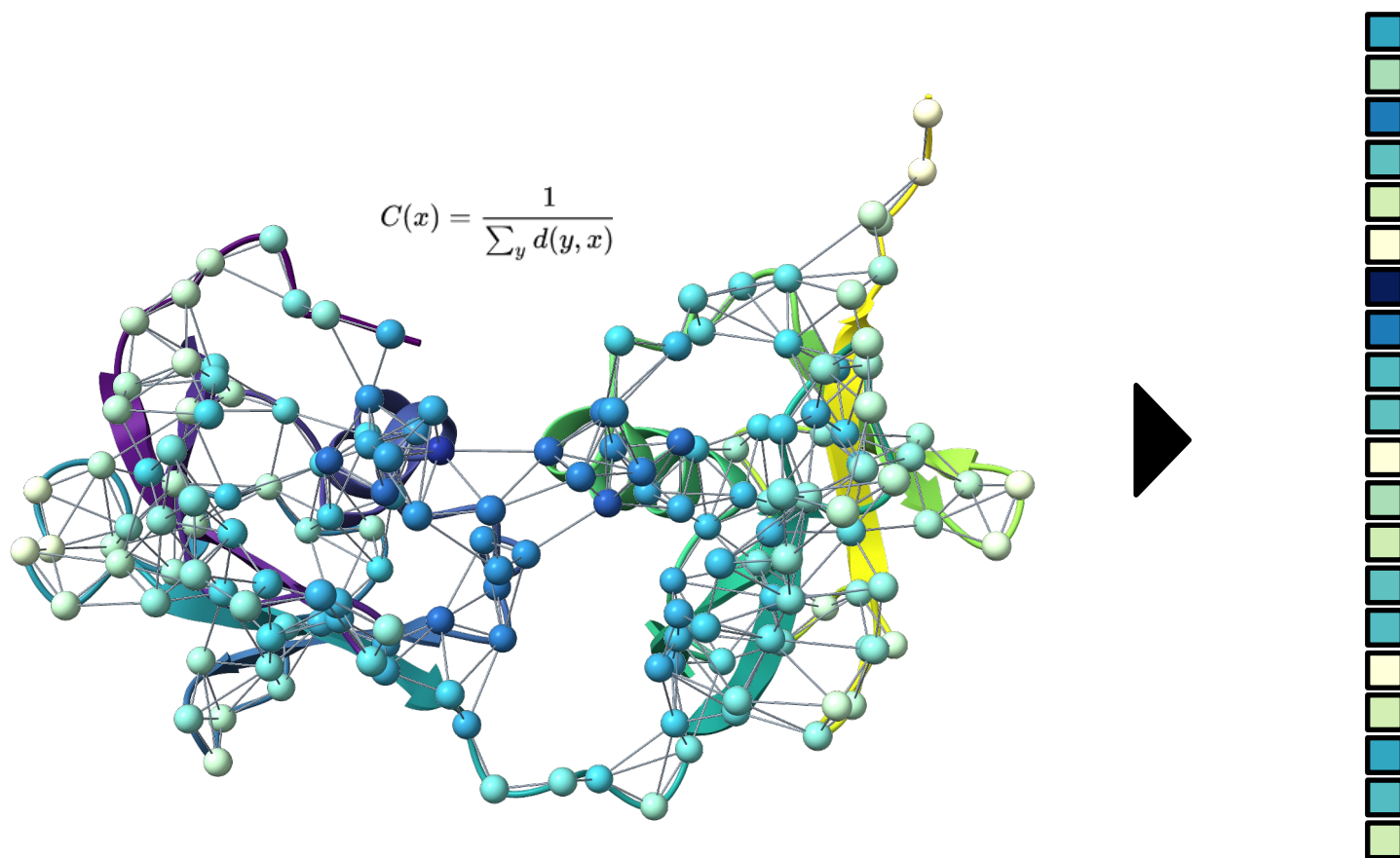
# Outline

- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations

# Extracting meaningful feature sets from protein graph representations



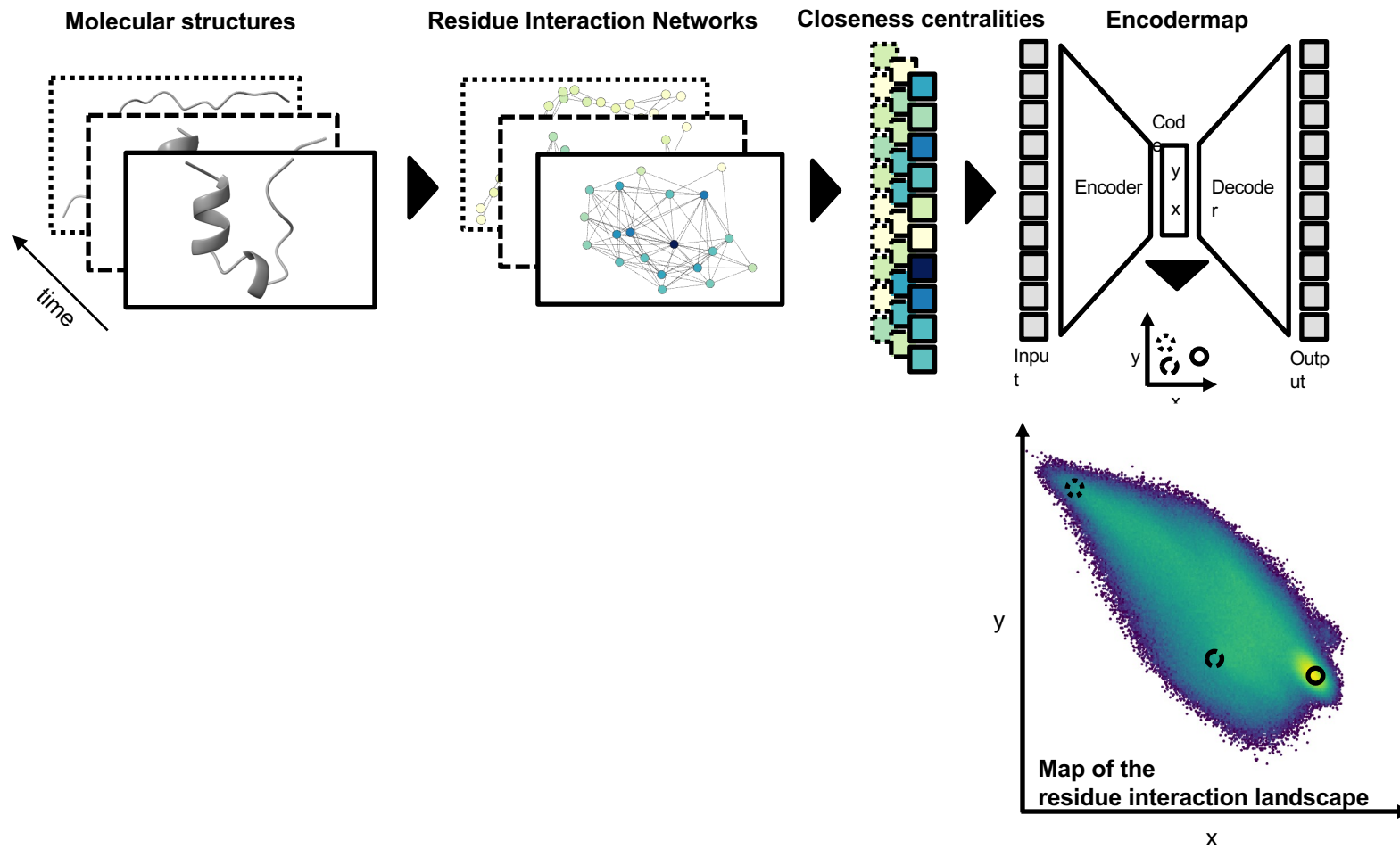
# Extracting meaningful feature sets from protein graph representations



→ Closeness centrality: „How close is a node to all other nodes in the network?“

→ Feature set that captures role of each residue in the protein structure

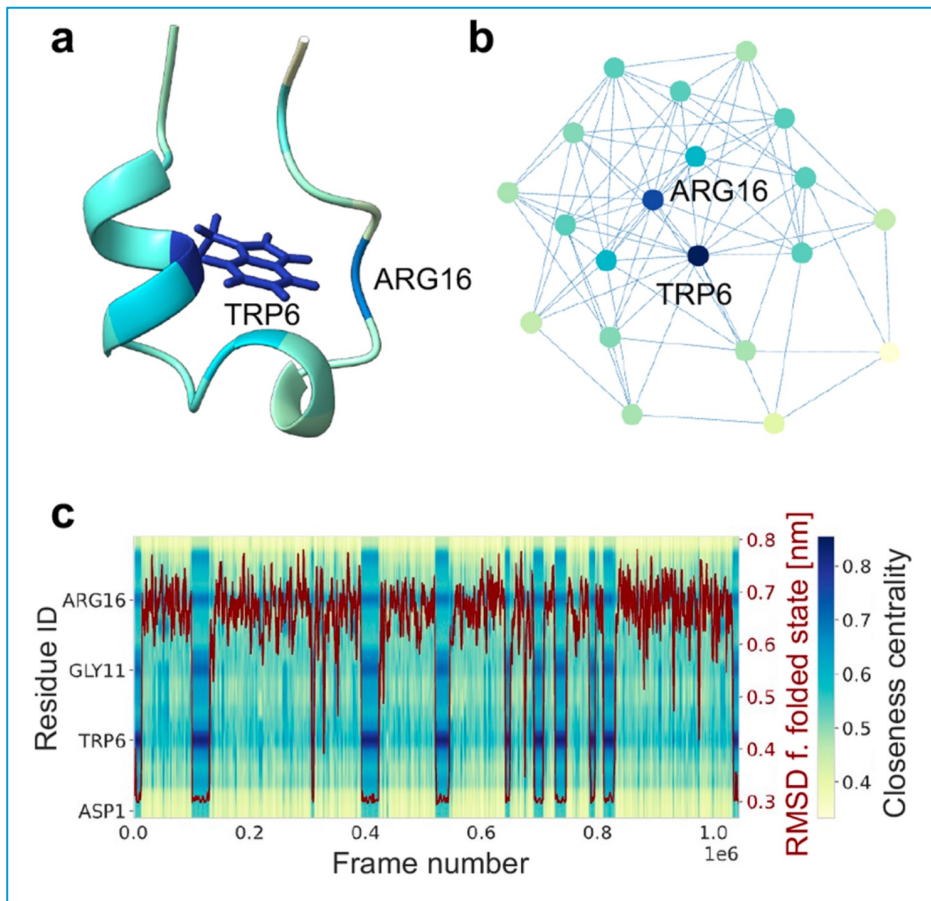
# Workflow



# Example 1: Trp-cage

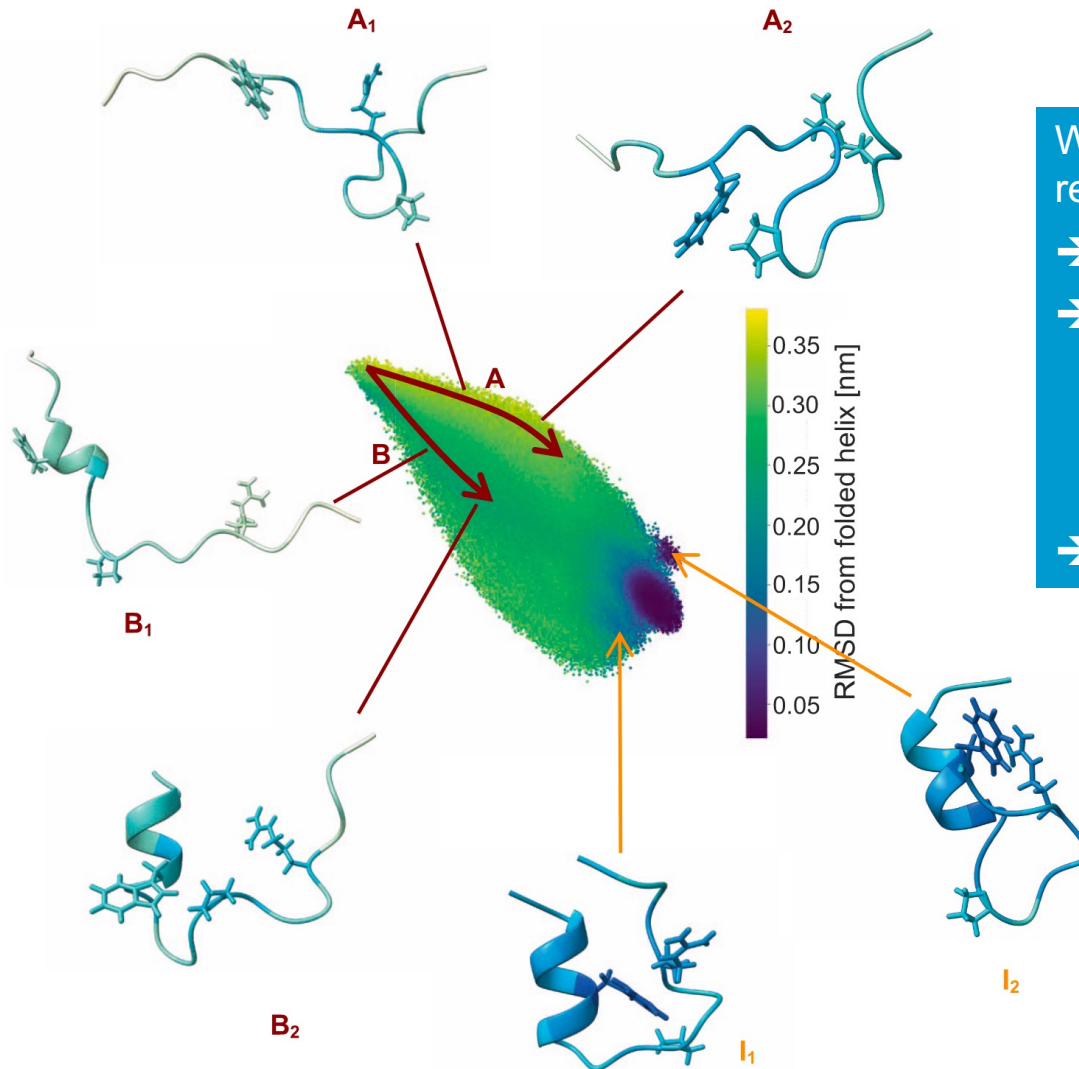


Trp-Cage



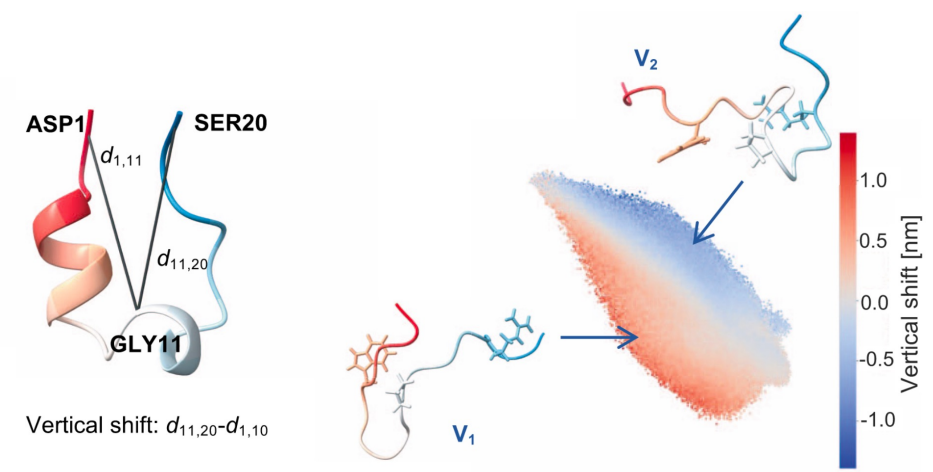


# Residue interaction landscape of Trp-cage

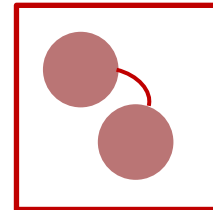


Without protein-specific prior input, the map resolves known folding behaviors

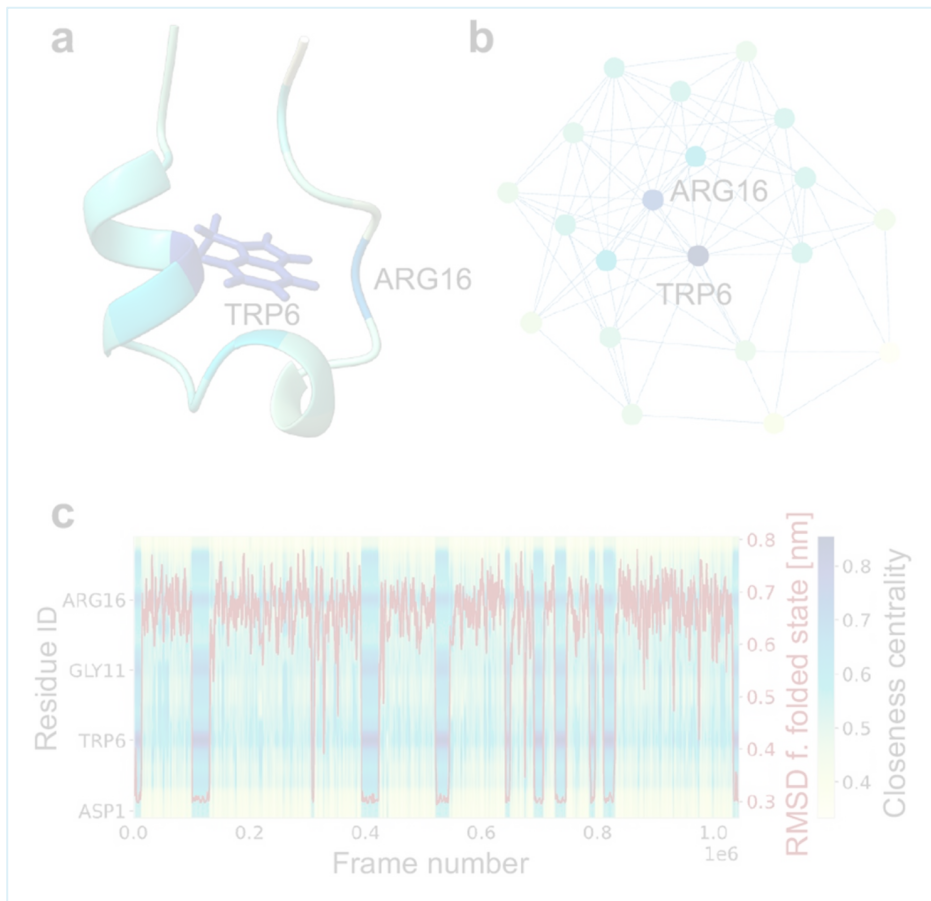
- Separation of folded and unfolded states
- Two folding pathways:
  - A: Nucleation-condensation; hydrophobic collapse first
  - B: Diffusion-collision; helix forms first
- Near-native intermediates



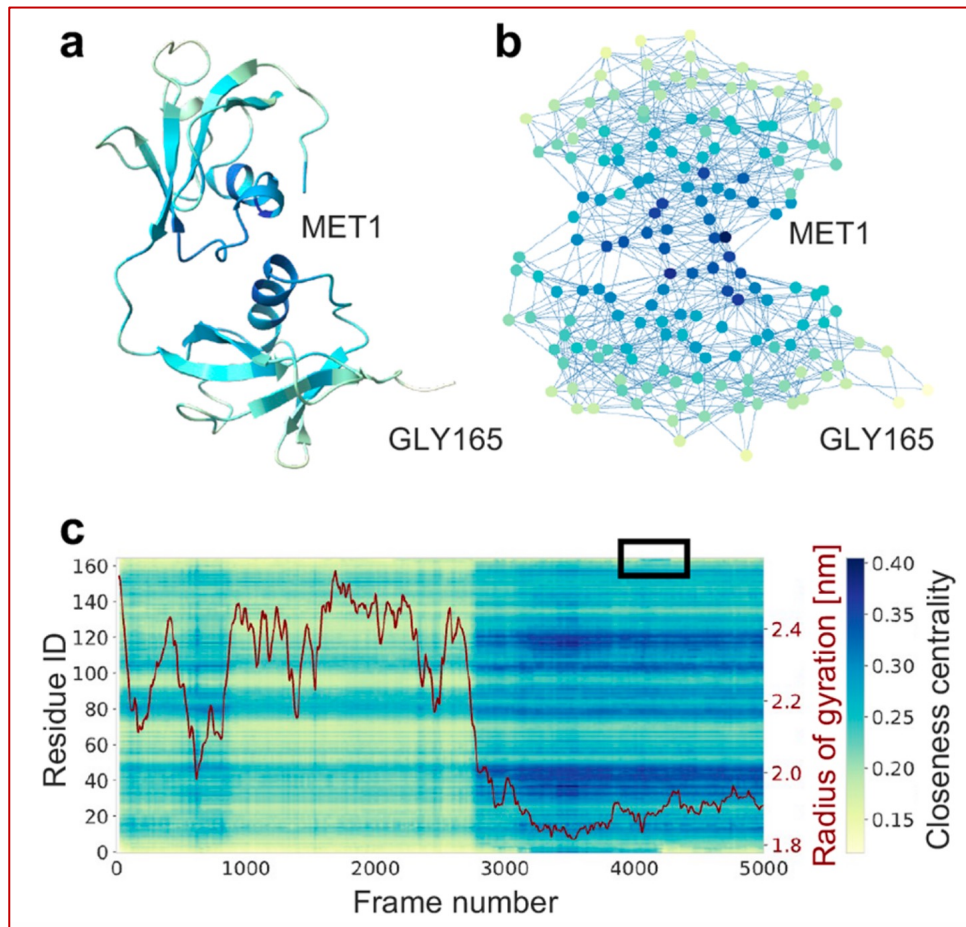
## Example 2: FAT10



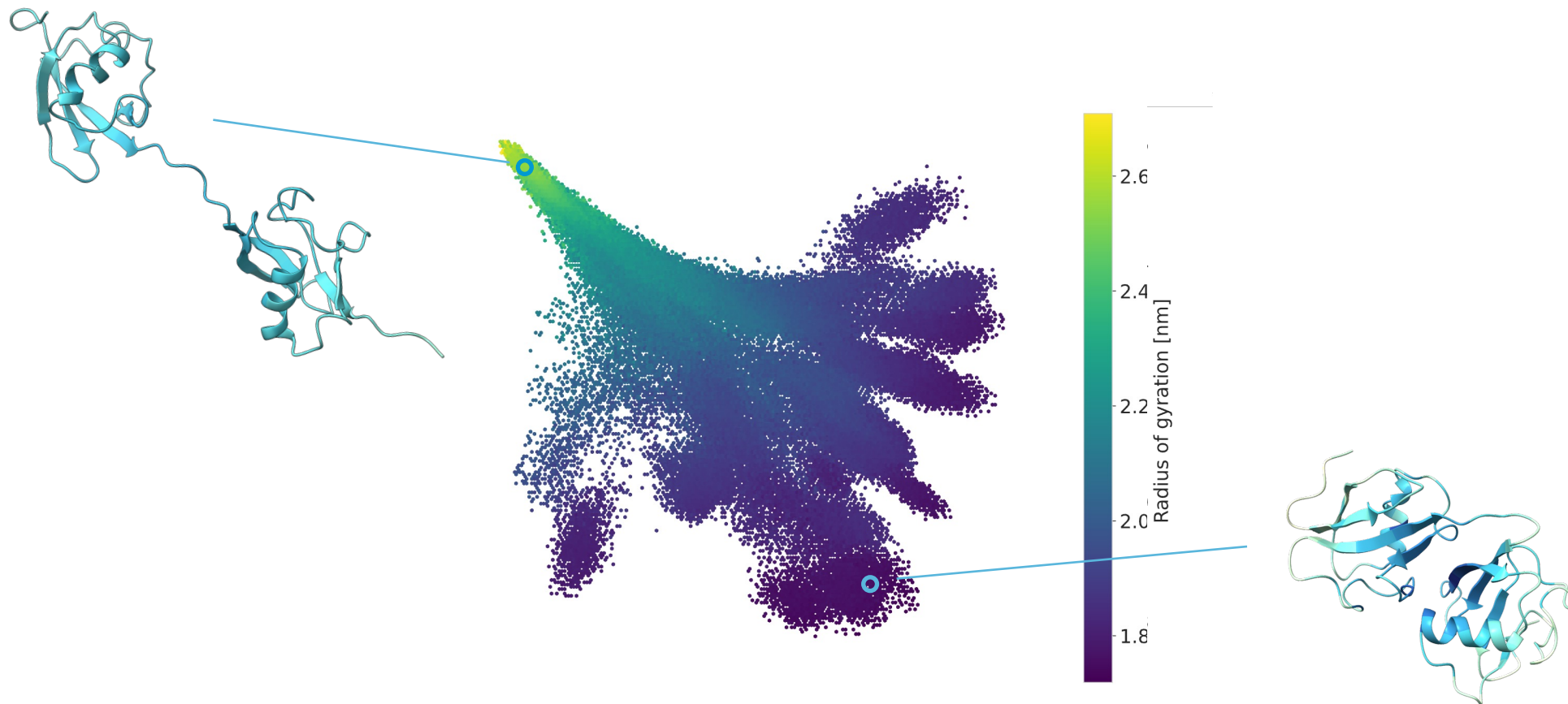
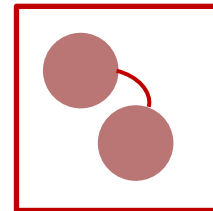
### Trp-Cage



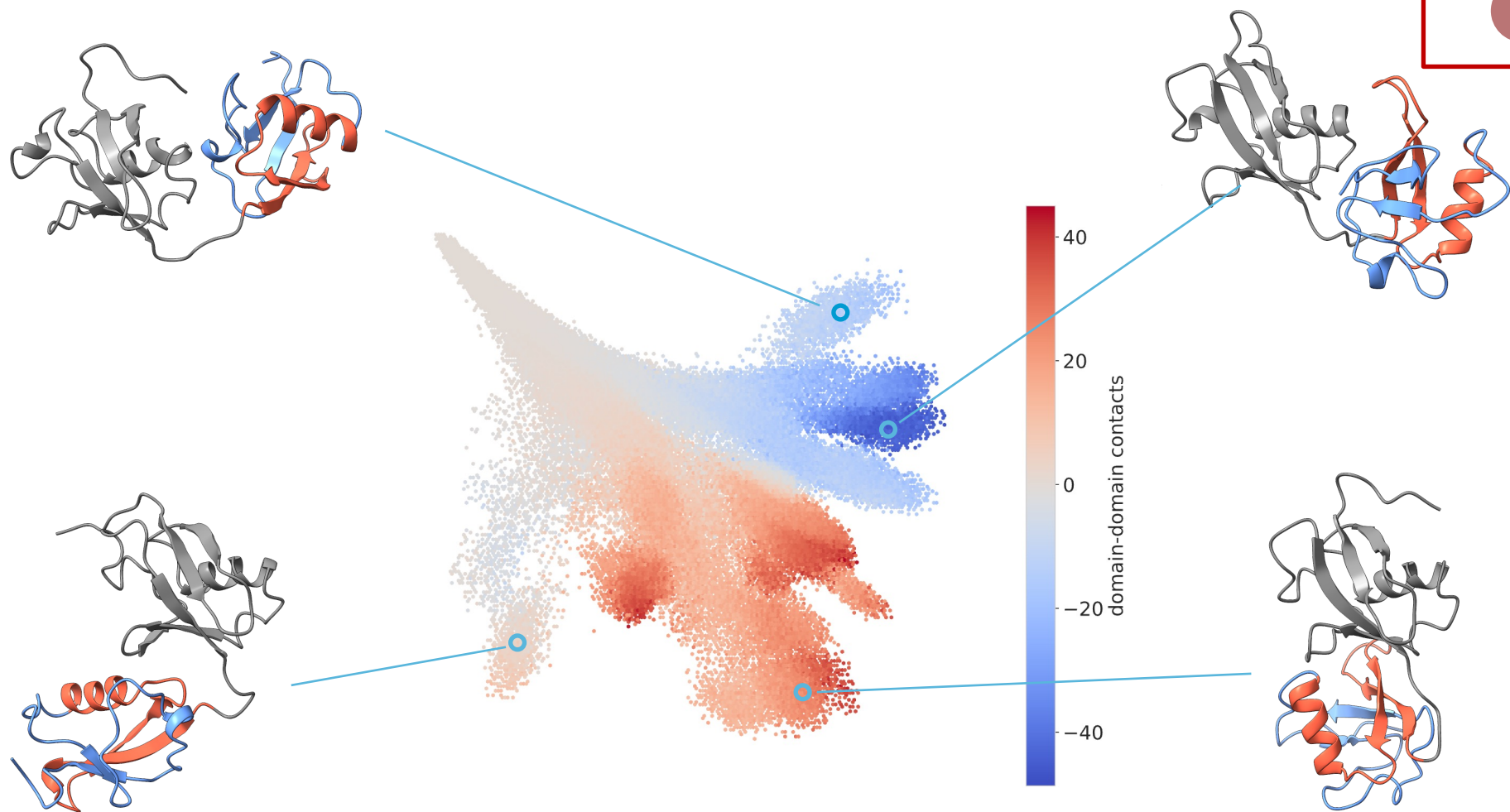
### FAT10



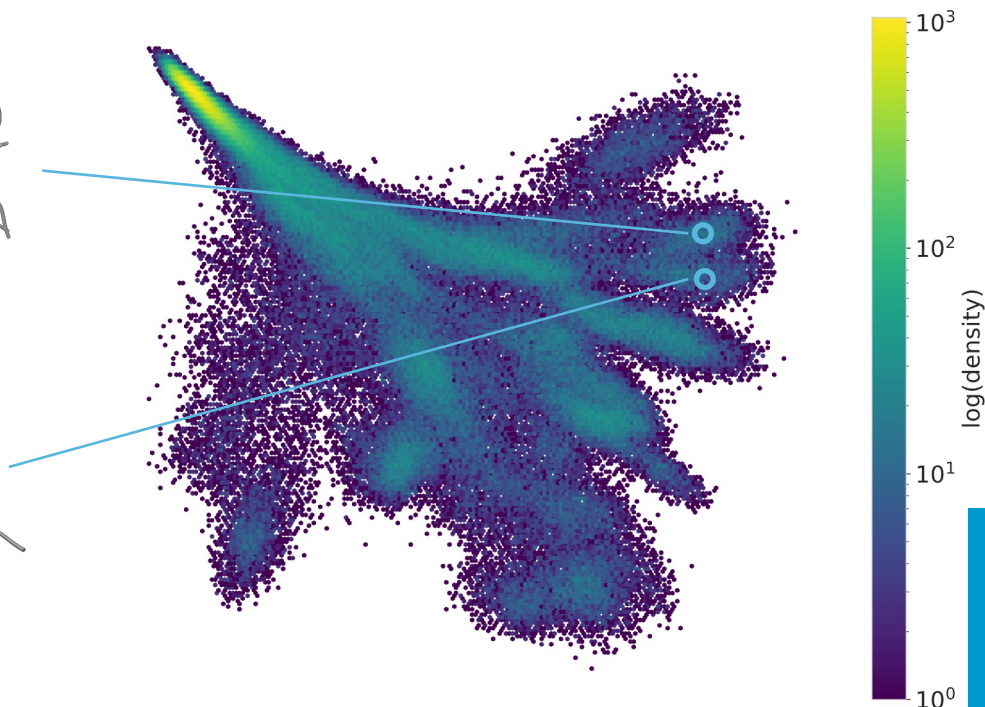
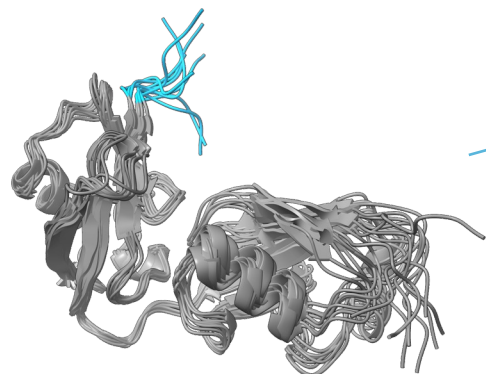
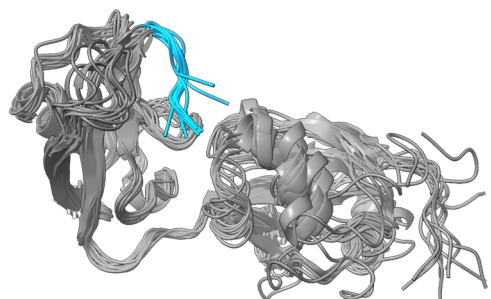
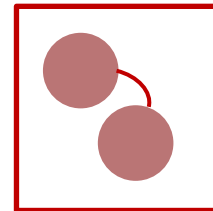
# Residue interaction landscape of FAT10



# Residue interaction landscape of FAT10



# Residue interaction landscape of FAT10

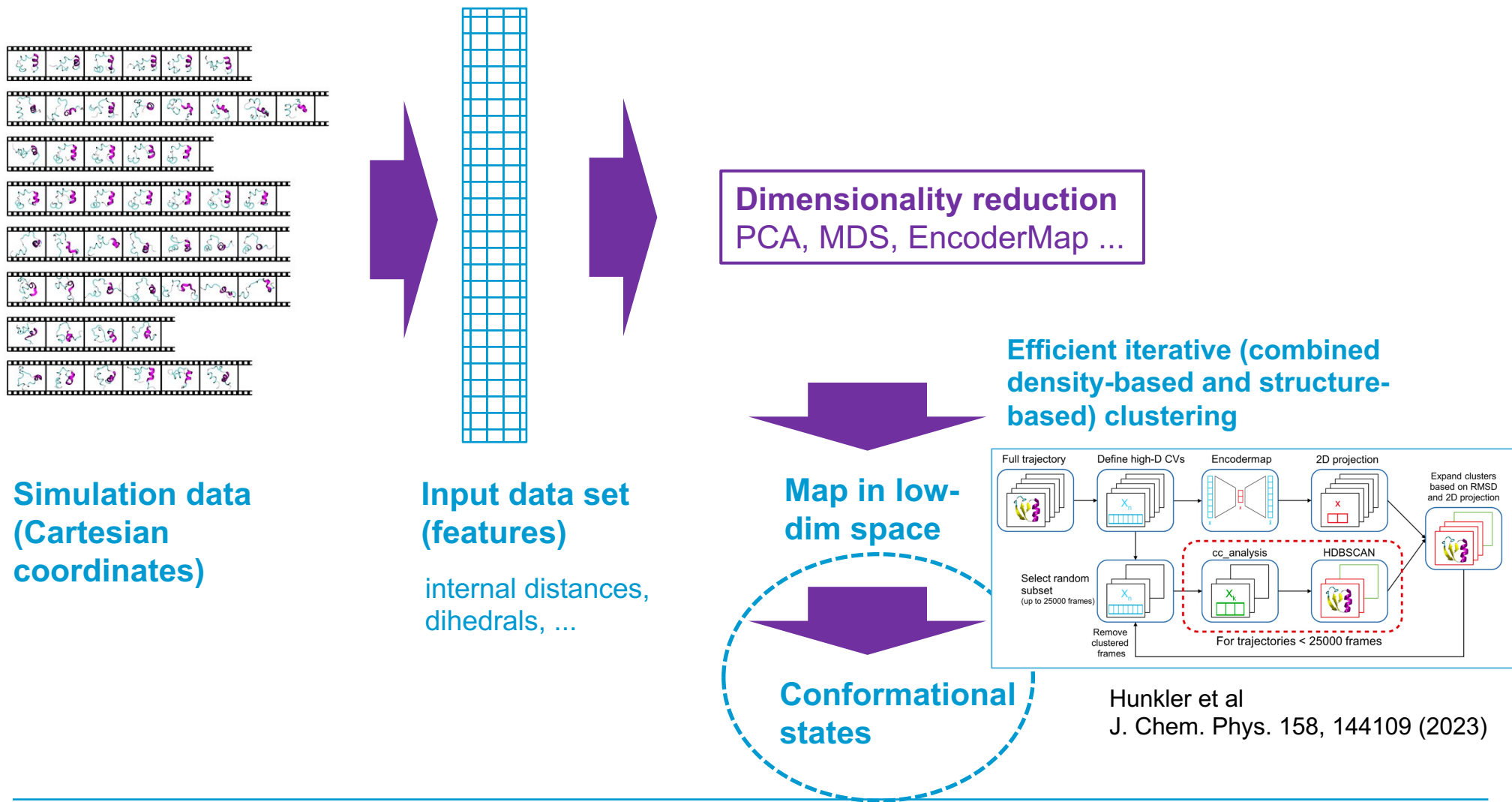


- global interpretability
- resolves meaningful details
- useful for visualization and further down-stream processing

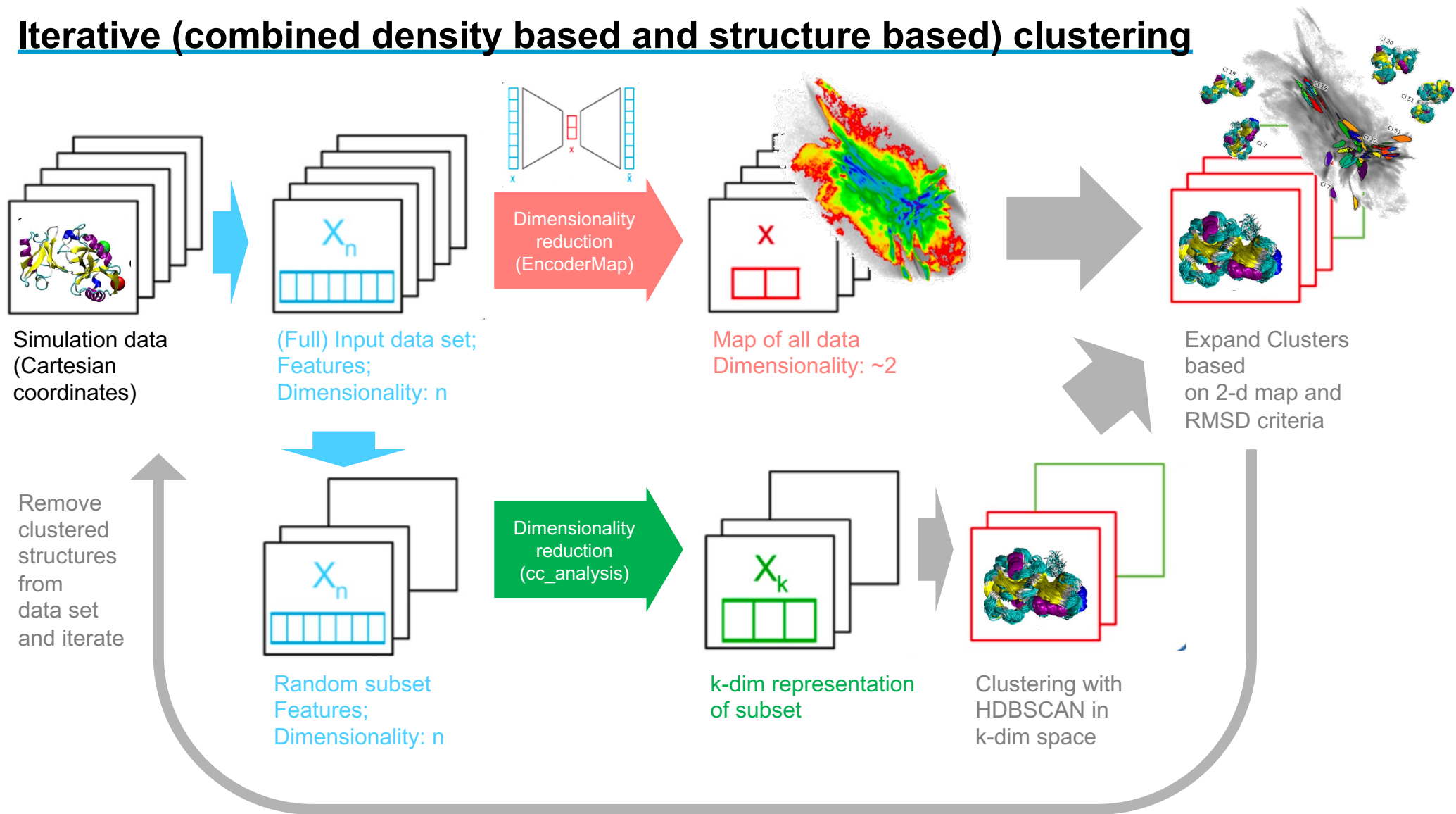
# Outline

- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations

# Clustering

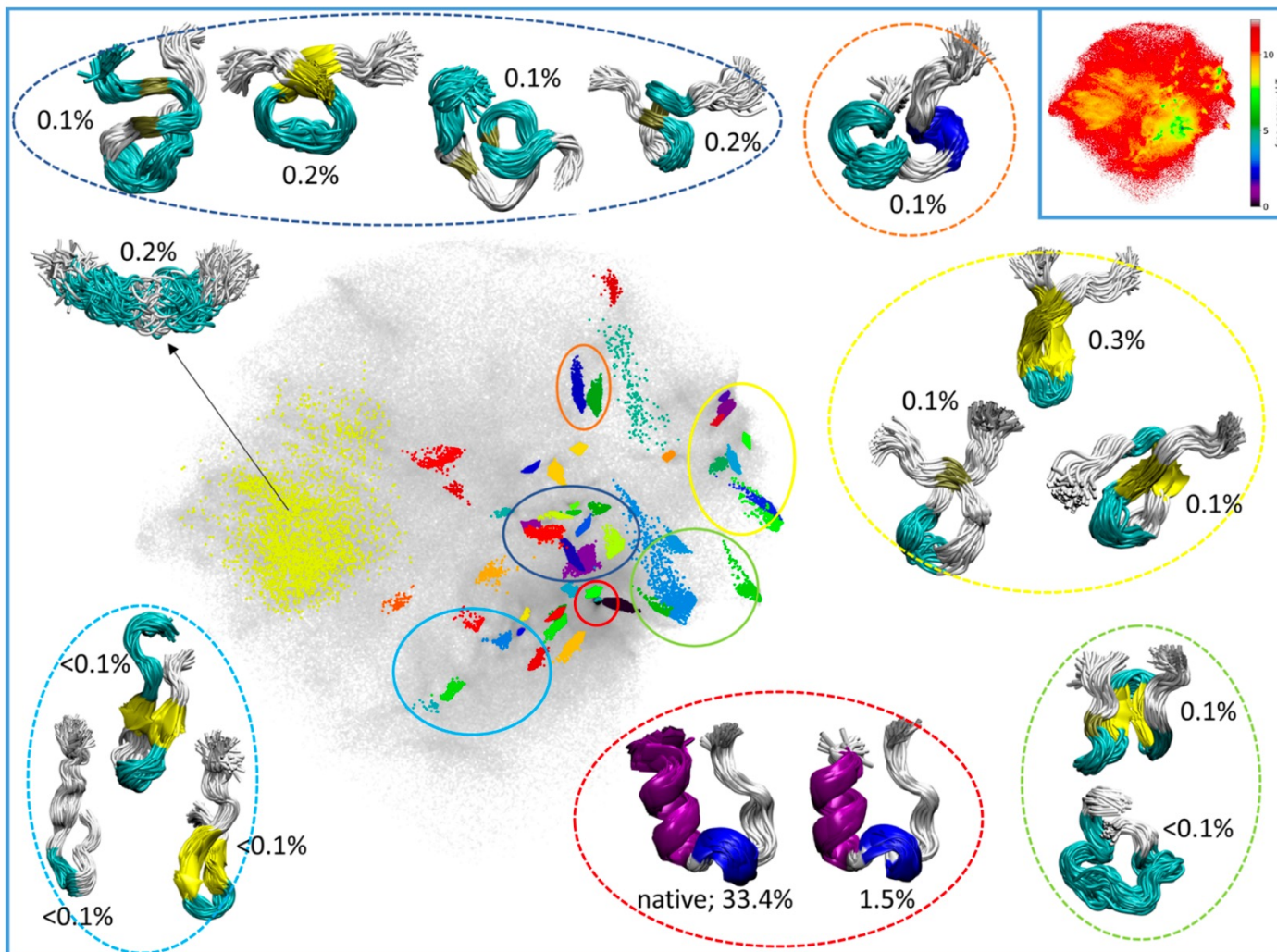


# Iterative (combined density based and structure based) clustering



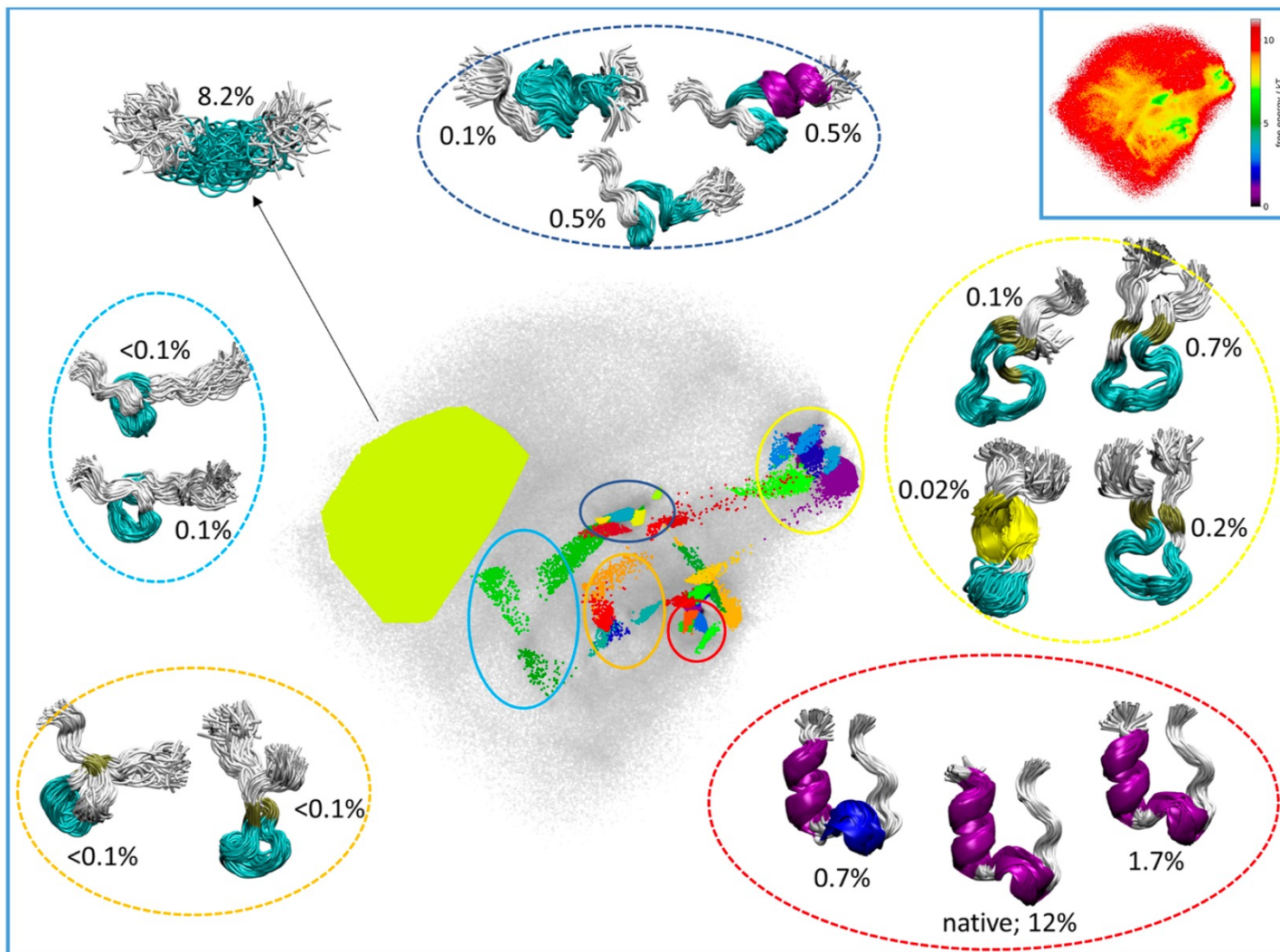


# Example 1: Trp-cage



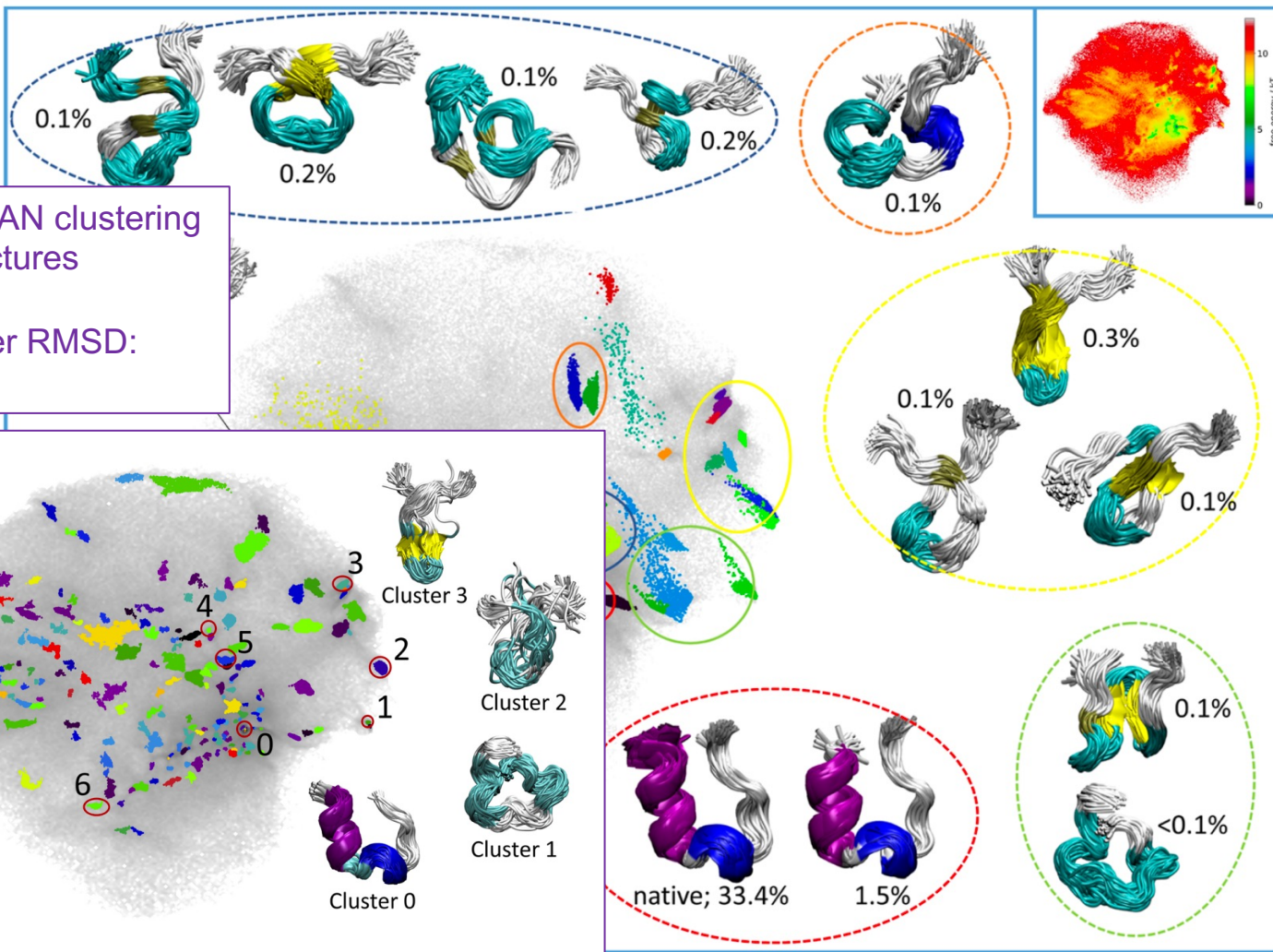
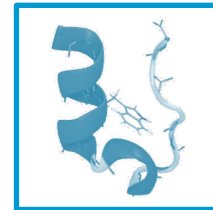
TC5b (40 temperature RE trajectories)

# Example 1: Trp-cage



TC10b  
(DE Shaw trajectories)

# Example 1: Trp-cage



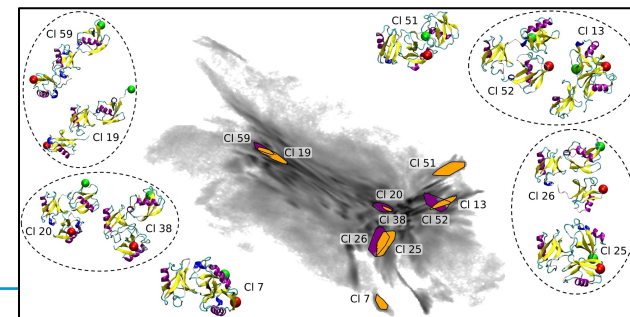
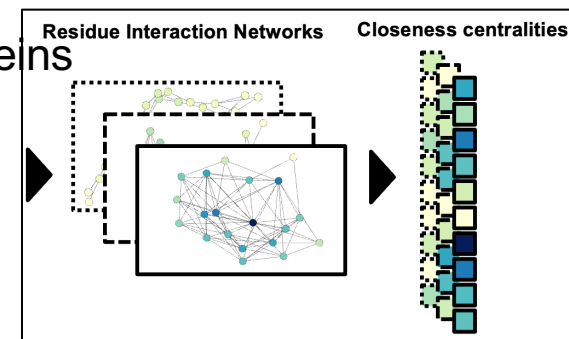
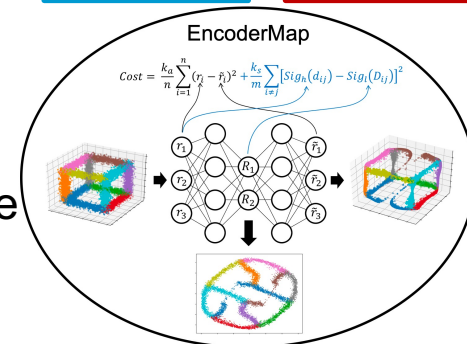
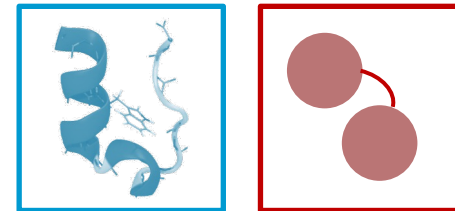
direct HDBSCAN clustering  
< 15% of structures assigned;  
average cluster RMSD:  
2.25 Å

~60% of structures assigned;  
average cluster RMSD:  
1.34 Å

TC5b (40 temperature RE trajectories)

# Outline

- Setting the stage
- EncoderMap
  - Learning meaningful representations of conformational phase space
  - Generating protein conformations and visualizing molecular motion
- Extracting meaningful feature sets from graph representations of proteins
  - Generating residue interaction landscapes
- Utilizing low-dimensional embeddings for clustering
  - Identifying conformational states
- Backmapping based sampling
  - Linking scales through low-dimensional representations



# Acknowledgements

Andrej Berg, Teresa Buhl, Christoph Globisch, Leon Franke, Simon Hunkler,

Oleksandra Kukharenko, Tobias Lemke, Kevin Sawade

Collaborators: Tobias Schneider, Michael Kovermann, Martin Scheffner, Kay Diederichs

