# In-situ Data Analysis of Protein Folding Trajectories

Travis Johnston[1], Boyu Zhang[1], Adam Liwo[2], Silvia Crivelli[3], Michela Taufer[1]

[1]U. Delaware Global Computing Lab, [2]U. Gdansk, Poland, [3]Lawrence Berkeley National Lab
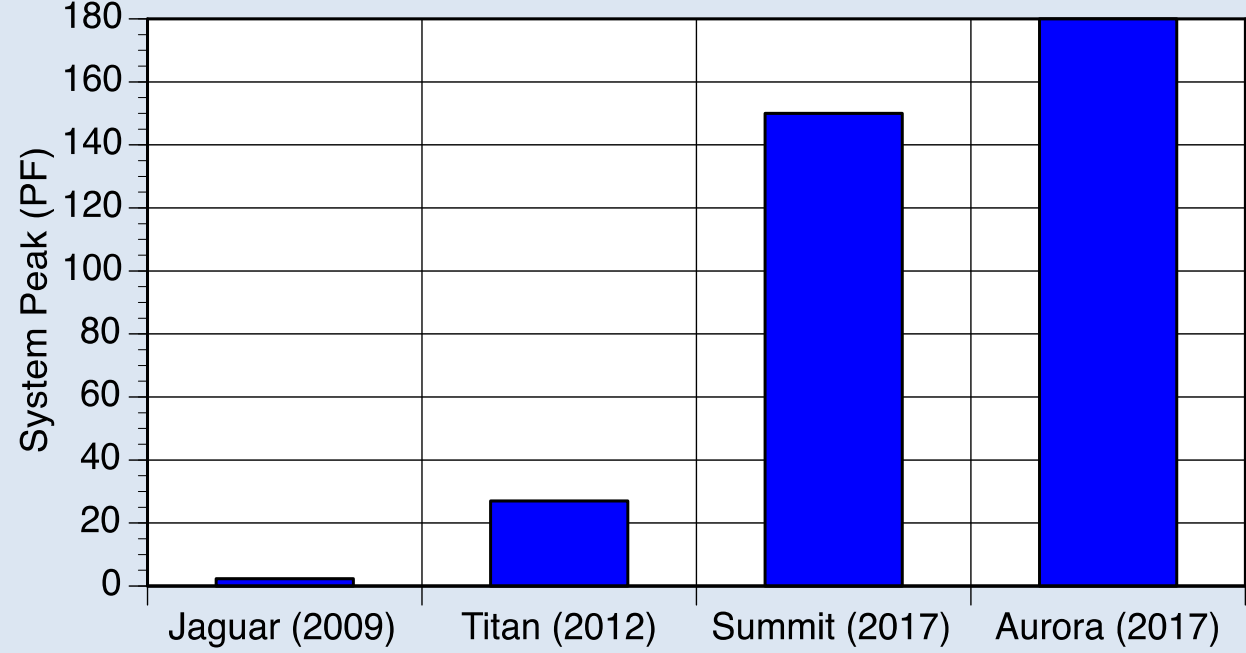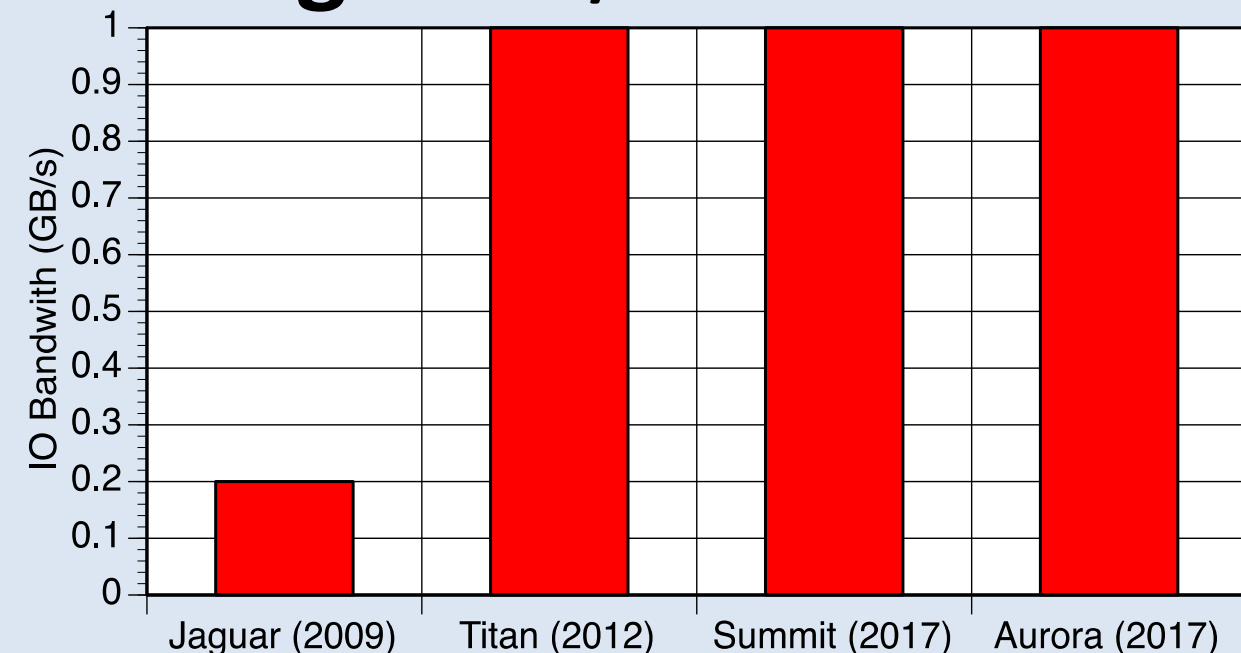
UNIVERSITY OF DELAWARE.

## Motivation

The transition from petascale to exascale computing will be accompanied by substantial changes in computer architectures and technologies. The research community relying on computational simulations is being forced to revisit the algorithms for data generation and analysis due to various concerns, such as higher degrees of concurrency, deeper memory hierarchies, substantial I/O and communication constraints. Simulations today typically save all data for post simulation analysis. Simulations at the exascale will require us to analyze data as it is generated and save only the results that enhance our scientific understanding. The analysis of this data will need to primarily be accomplished *in-situ*, i.e. executed sufficiently fast locally, using very limited amounts of memory and disk space, and avoiding large data movement.
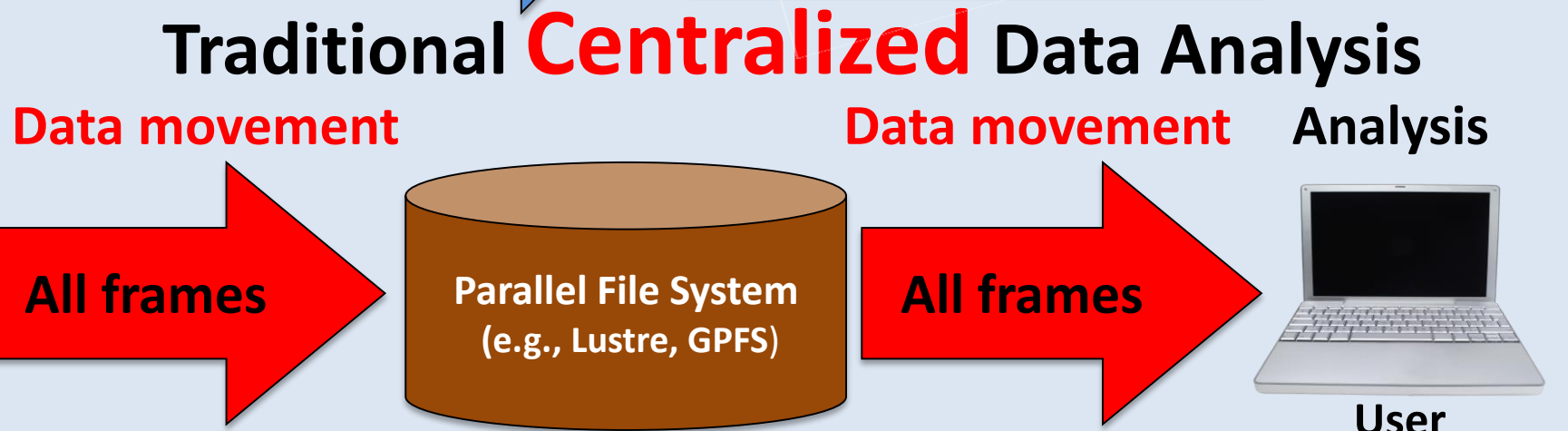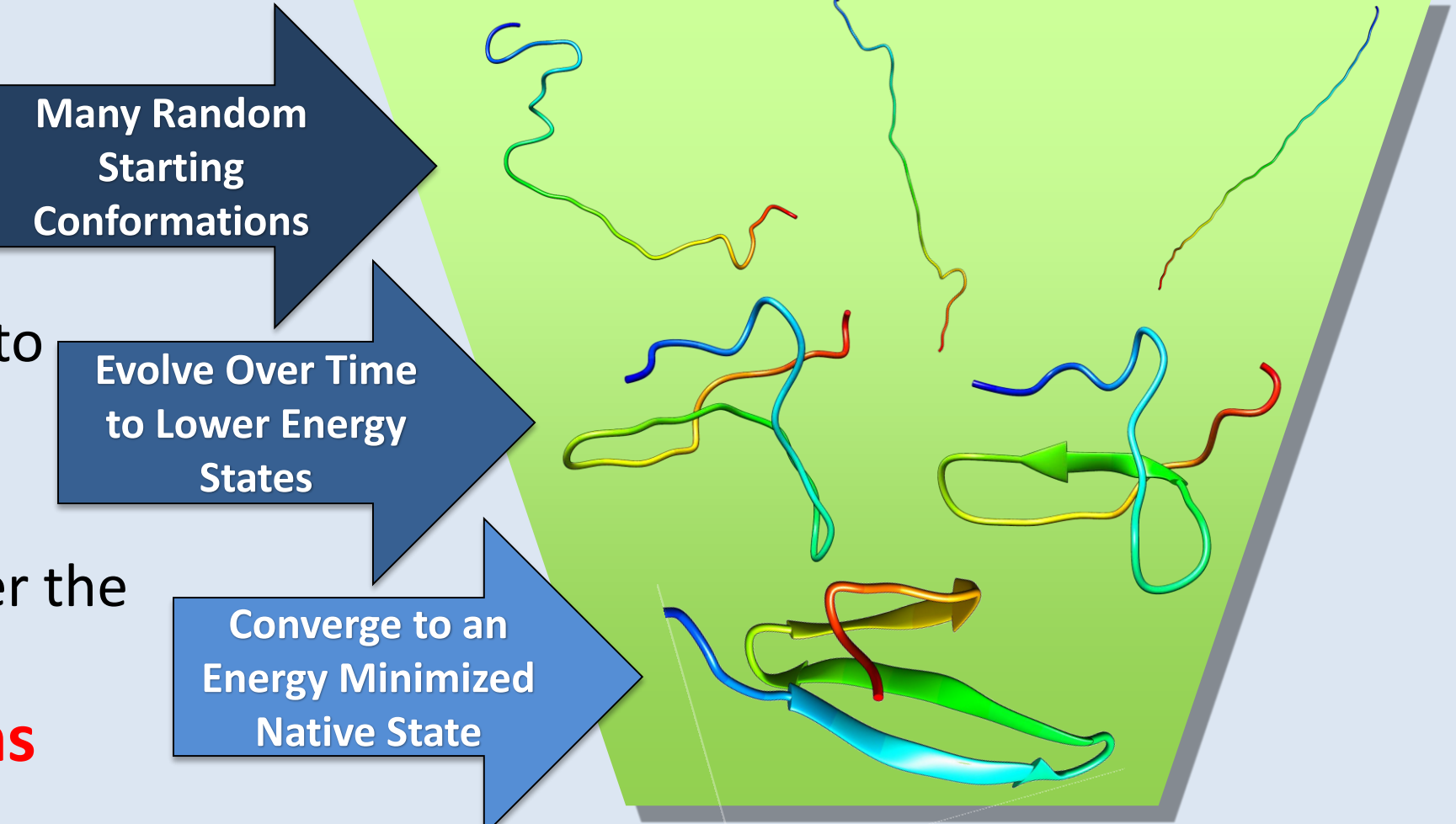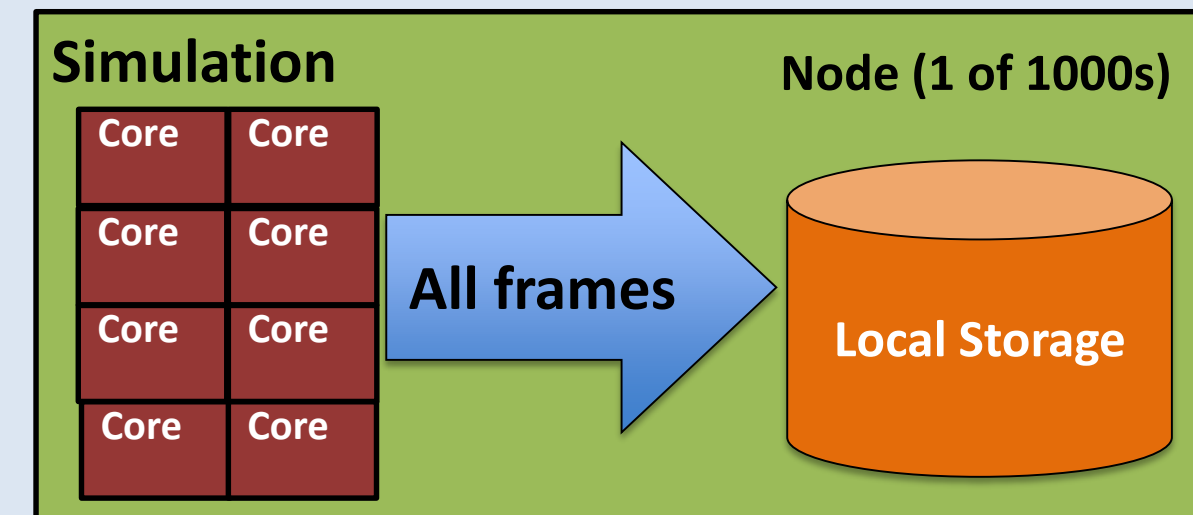
### Increasing Peak Performance



### Stagnant I/O Bandwidth



## Protein Folding

- Start from many unfolded conformations of a protein with correct chemical bonds but random torsion angles
- Run simulations on supercomputers to generate conformations (frames) of multiple folding trajectories
- Store all frame and analysis them after the simulation completes
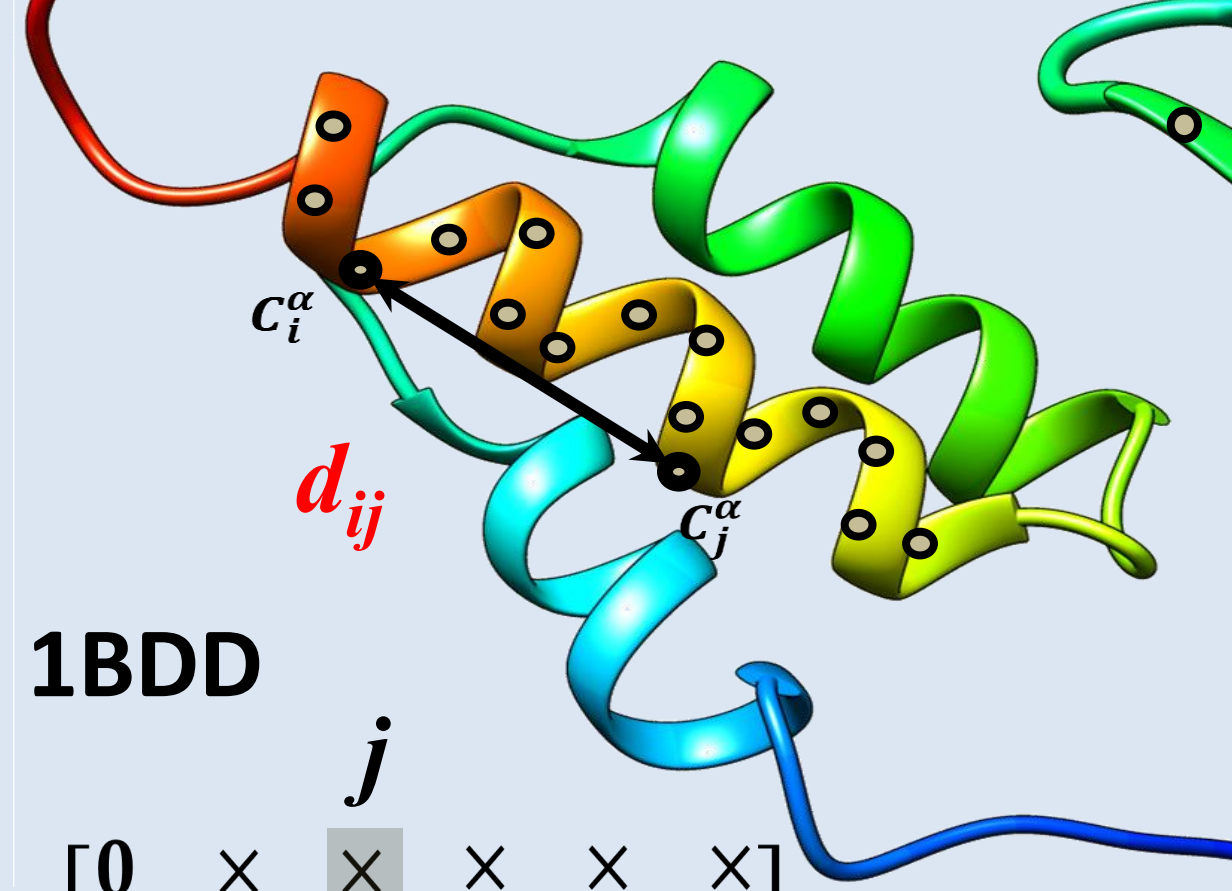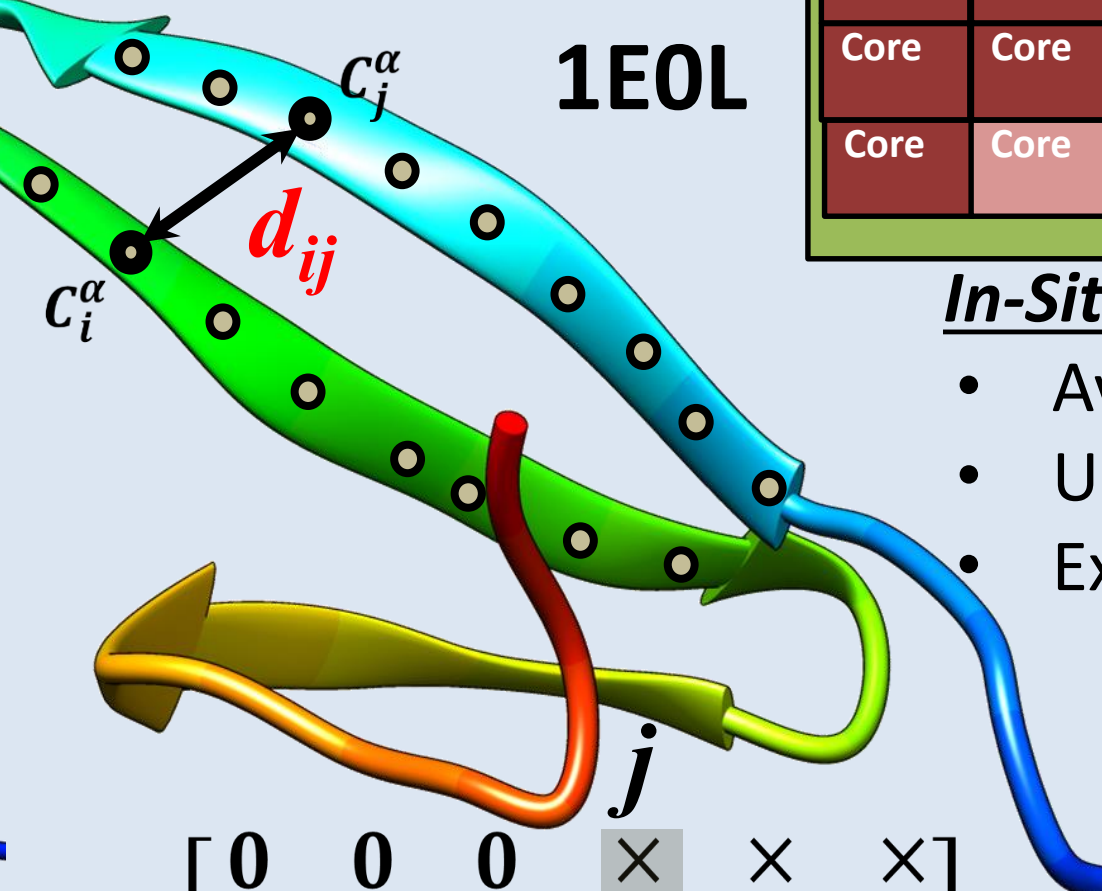
**Does NOT scale to exascale systems**



Traditional **Centralized** Data Analysis

## Method for In-situ Data Analysis

### Algorithm

1. Discretize the protein by extracting the positions of α-Carbons and β-Carbons
2. Build a **Euclidean Distance Matrix**, D, (single structure) or a **Bipartite** distance matrix (two structures)
3. Associate the **largest eigenvalue** of D with the protein conformation as **metadata**
4. Use the metadata to identify stable and transition states during the simulation

Individual Substructures — **1BDD**

Pairs of Substructures — **1E0L**

$$D = \begin{pmatrix} 0 & \times & \times & \times & \times & \times \\ \times & 0 & d & \times & \times & \times \\ \times & d & 0 & \times & \times & \times \\ \times & \times & \times & 0 & \times & \times \\ \times & \times & \times & \times & 0 & \times \\ \times & \times & \times & \times & \times & 0 \end{pmatrix} \Rightarrow \lambda_1$$

$$D = \begin{pmatrix} 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & d & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ \times & d & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \end{pmatrix} \Rightarrow \lambda_1$$

Novel *In-Situ* Data Analysis

*In-Situ* Data Analysis
- Avoid data movement among nodes
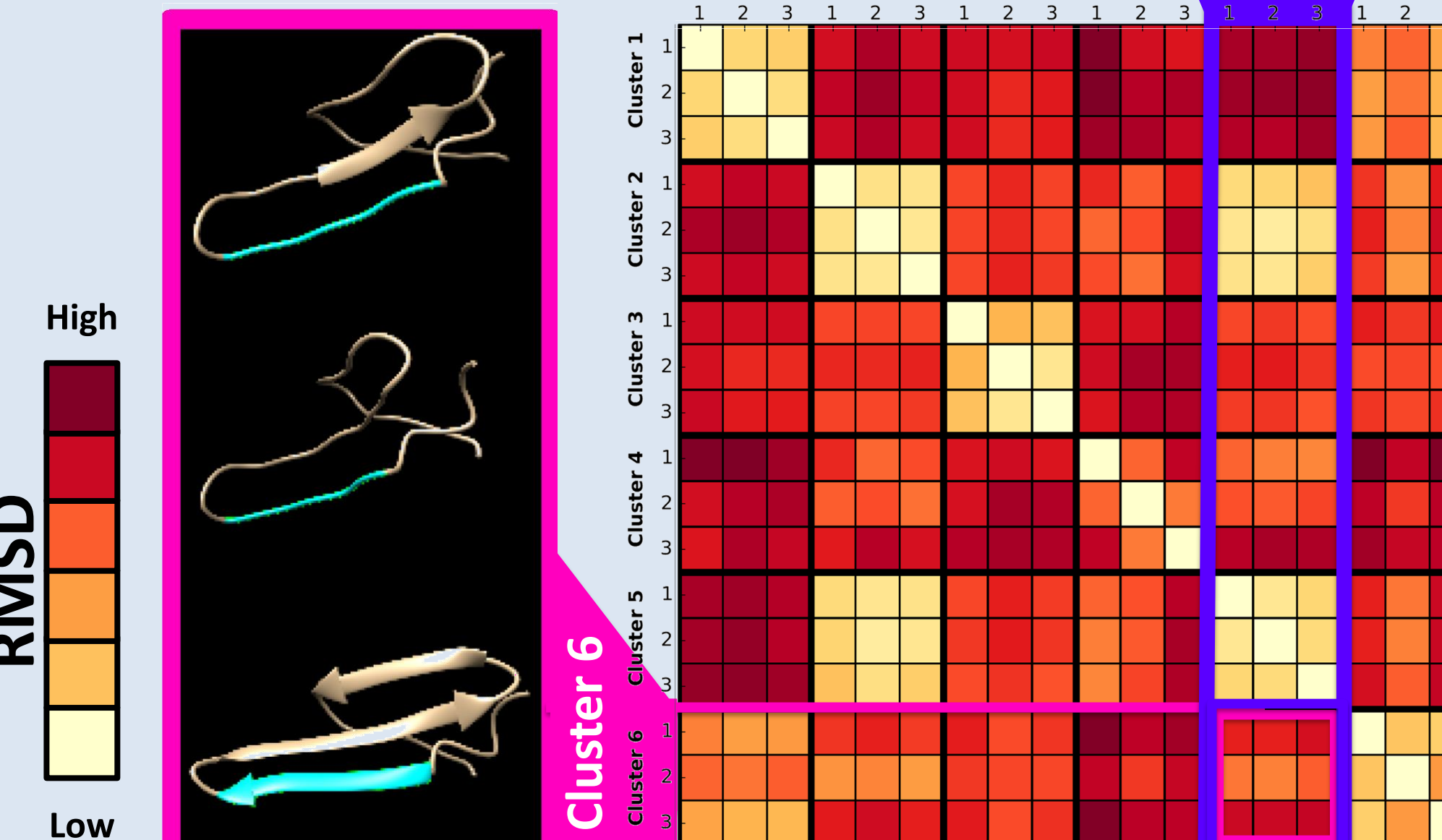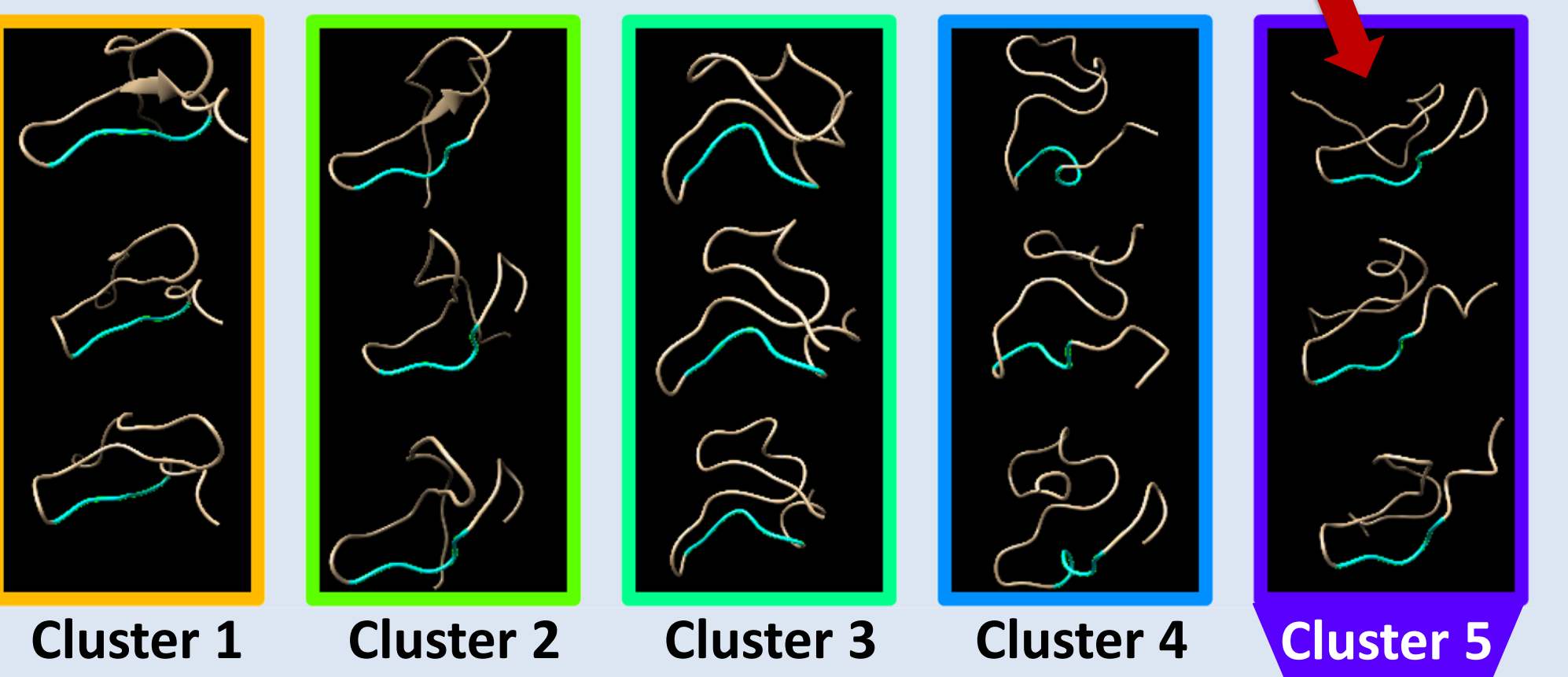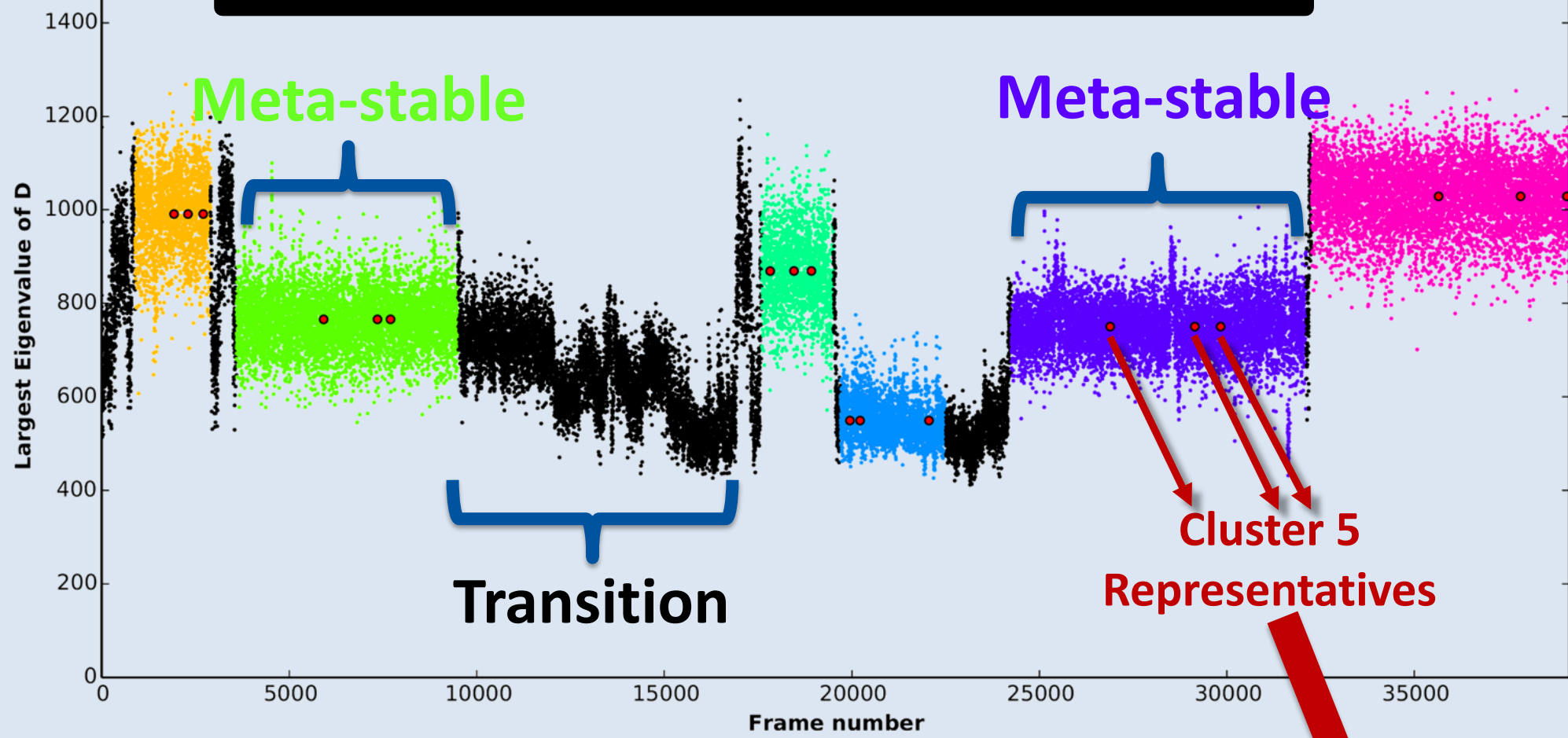- Use a small amount of memory
- Execute sufficiently fast

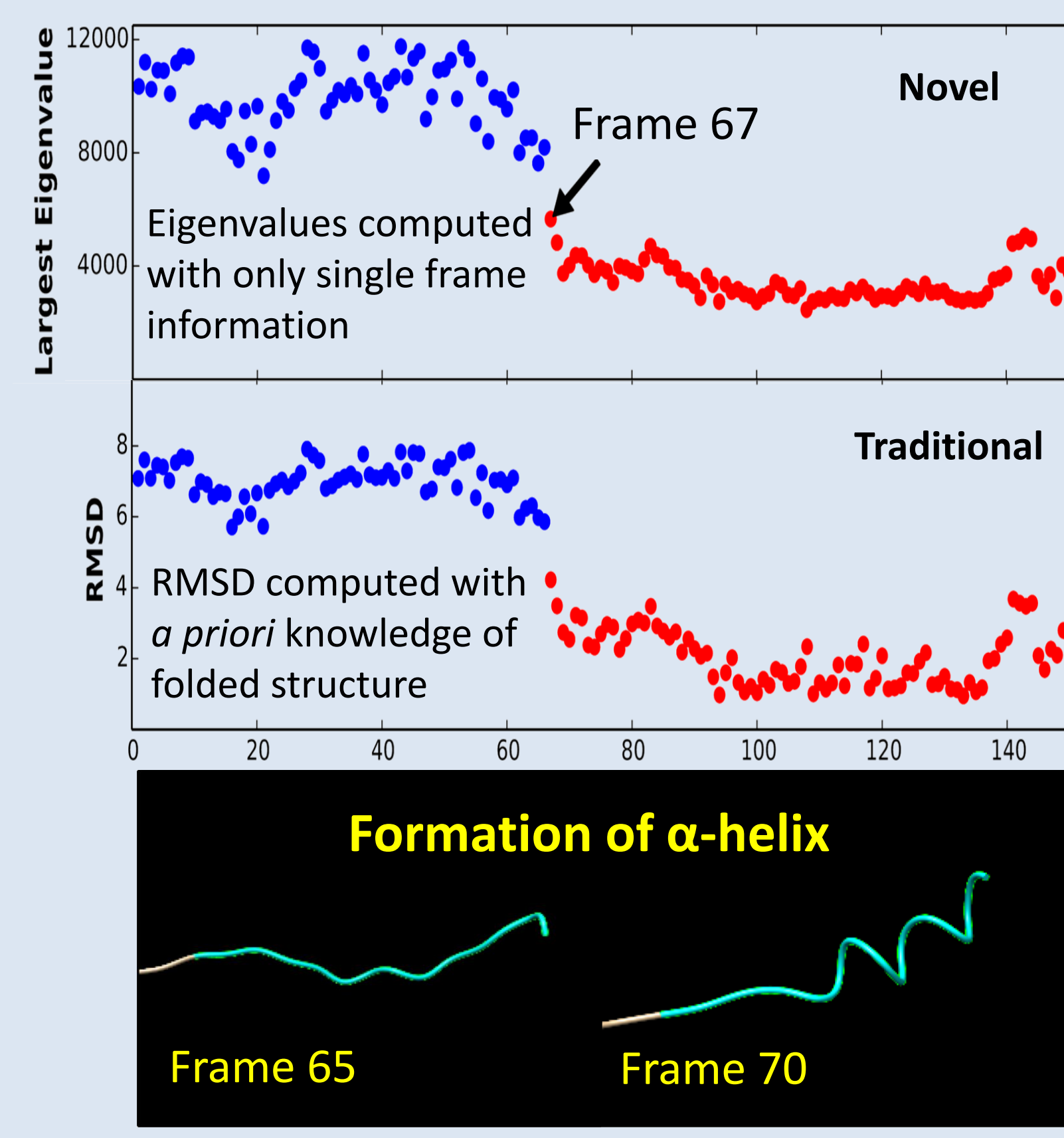**Does scale to exascale systems**

## Our Contribution

To satisfy I/O bandwidth constraints at exascale, we propose a method for *in-situ* data analysis:
- Accurately captures features of protein trajectory
- Performs light computation and does not interfere with ongoing folding simulation
- Significantly reduces the amount of data written to disk (from about 40k frames to a few 10s of frames)
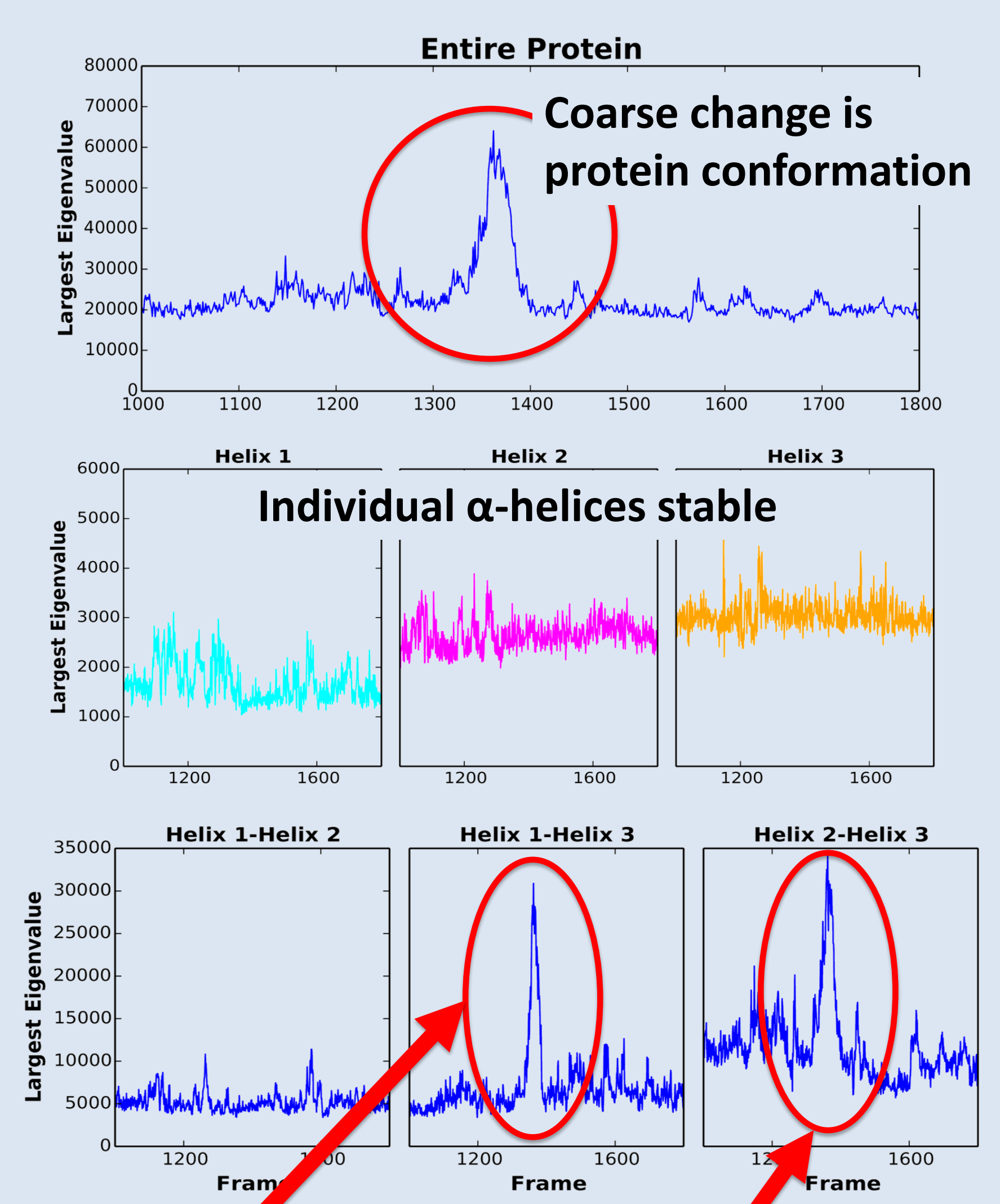
## 1E0L: Single β-Strand



Meta-stable — Transition — Cluster 5 Representatives

Cluster 1, Cluster 2, Cluster 3, Cluster 4, Cluster 5, Cluster 6

RMSD: High / Low

## 1BDD: Single α-Helix



Frame 67 — Eigenvalues computed with only single frame information — **Novel**

RMSD computed with *a priori* knowledge of folded structure — **Traditional**

**Formation of α-helix**

Frame 65 — Frame 70

## 1BDD: Pair of α-Helices



Entire Protein — **Coarse change is protein conformation**

Individual α-helices stable — Helix 1, Helix 2, Helix 3

Helix 1-Helix 2, Helix 1-Helix 3, Helix 2-Helix 3

**Movement between 1st and 3rd α-helix**

Frame 1300 — Frame 1370 — Frame 1410

## Conclusions

We propose a novel method for *in-situ* data analysis of protein folding trajectories. We validate our metadata mapping method by applying it to two 40k frame trajectories: one trajectory of 1BDD (containing 3 α-helices), and one trajectory of 1E0L (containing 3 β-strands). Our metadata mapping enabled us to observe metastable states, transition states, the formation of an individual substructure and the repositioning of substructure relative to others.