



UNIVERSITY *of* DELAWARE

# ***Challenges in Big Data Computing on HPC Platforms***

Michela Taufer

Computer and Information Sciences

University of Delaware

Newark, Delaware, USA



## Acknowledgements



Travis J.



Boyu Z.



Trilce E.



Adam L.



Silvia C.



Ewa D.



Michel C.



Dong A.



Don L.



Rafael D.



Stephen H.

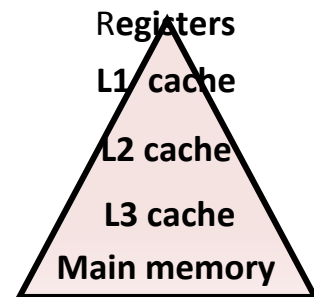
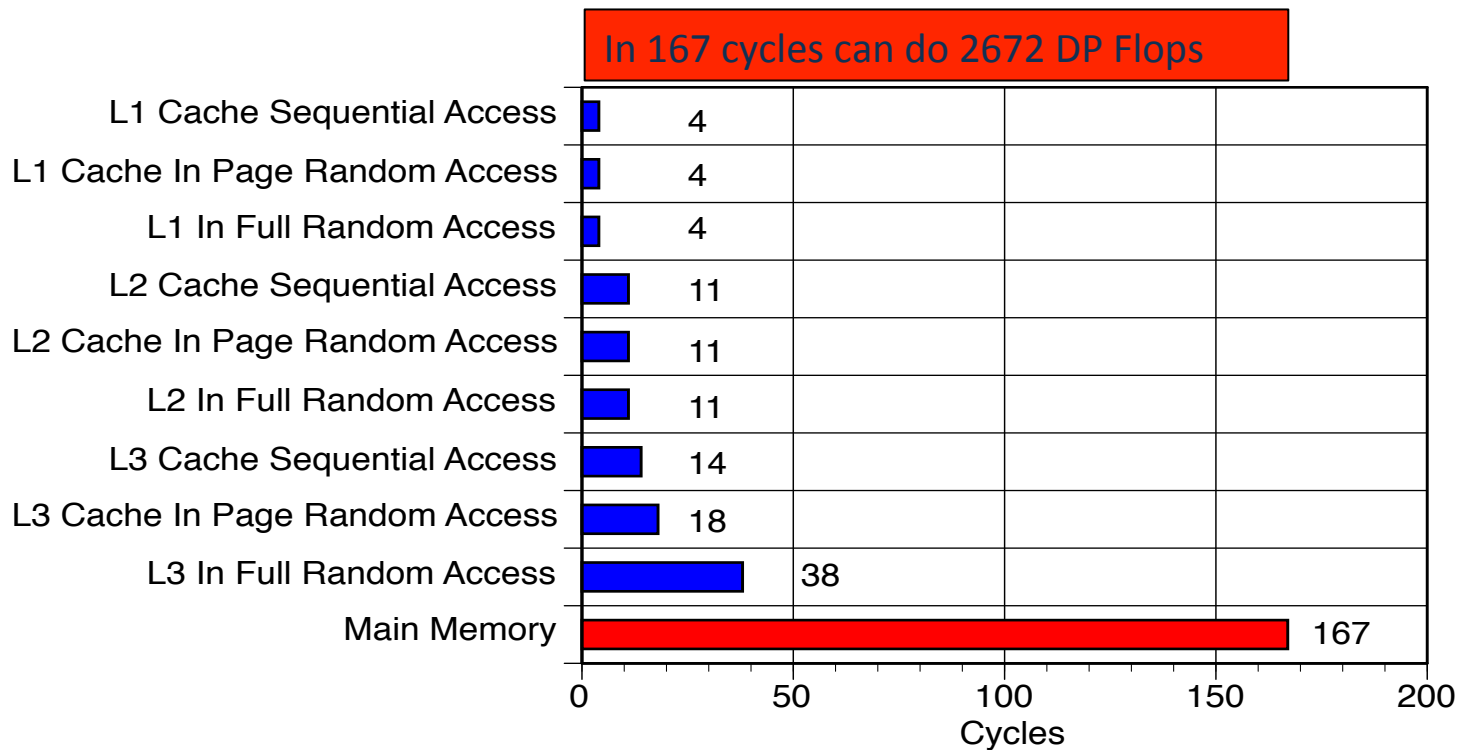
Sponsors:





## The Cost of Data Movement

- Today floating point operations are inexpensive



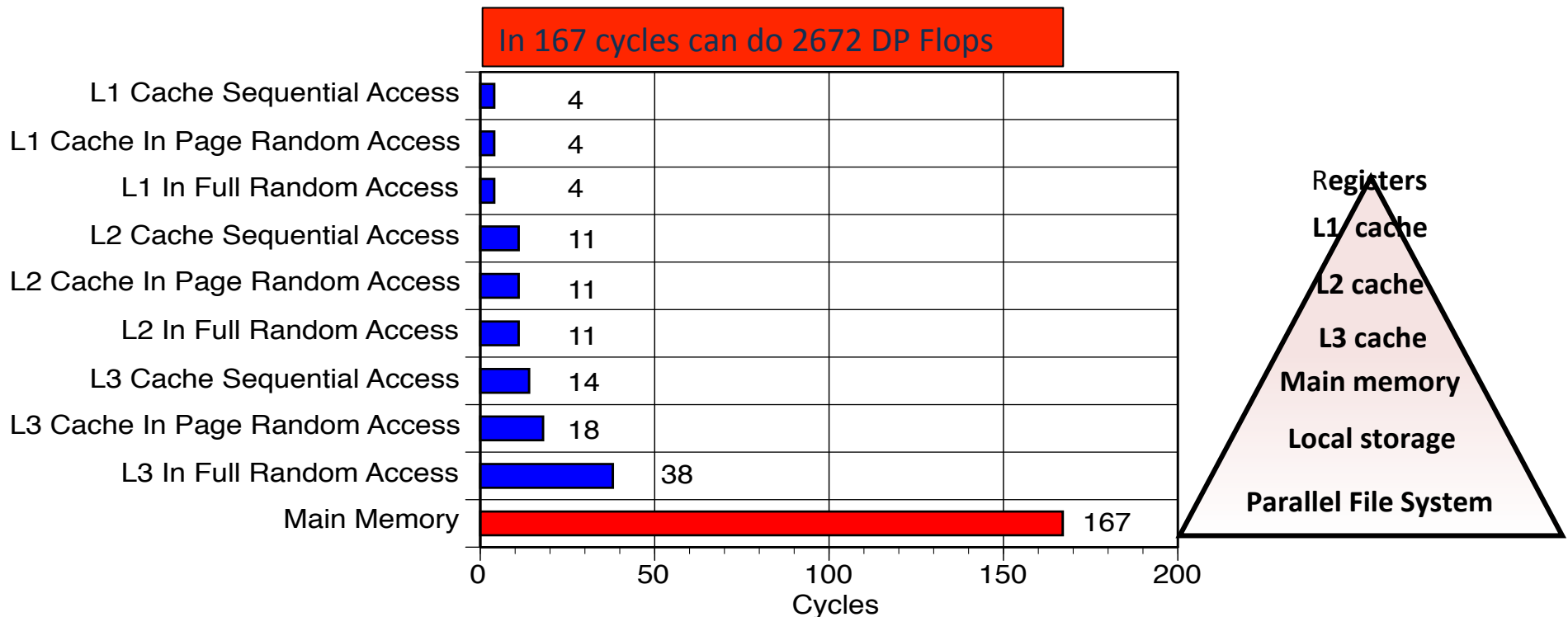
- Data movement is very expensive

*Courtesy of Jack Dongarra, UTK and ORNL*



## The Cost of Data Movement

- Today floating point operations are inexpensive



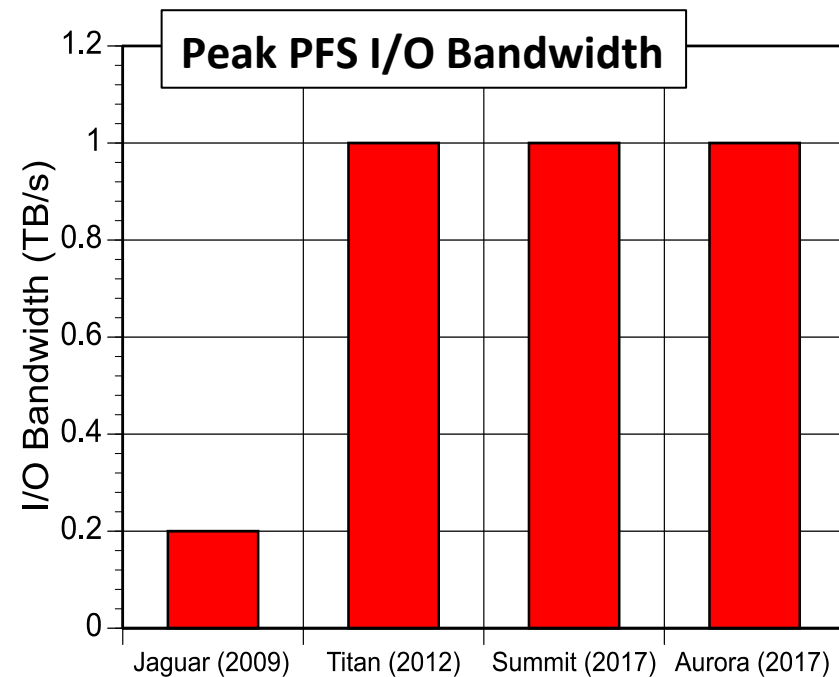
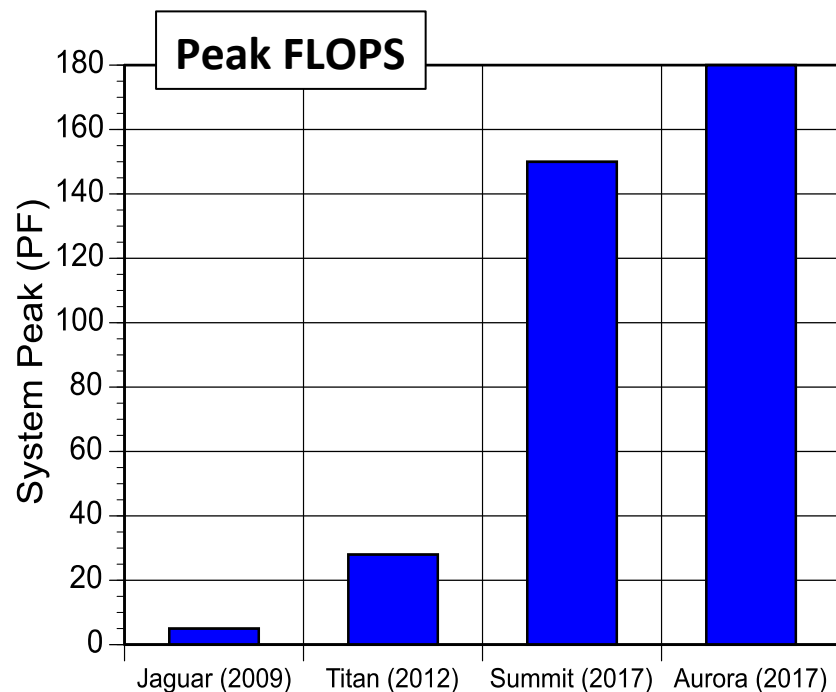
- Data movement is very expensive

*Courtesy of Jack Dongarra, UTK and ORNL*



## The Cost of Data Movement

- Floating point operations will further increase



- Speed to move data down the memory hierarchy is stagnant



## Perspective

The scientist:

*“Storage technologies are advancing [...] and it is really not clear at all [to me] that especially distributed storage platforms would not be able to handle [...] petabyte data sets”*

*Anonymous Feedback*

The computer architect:

*“[...] there will be burst buffers on the DOE machines which will give applications much faster I/O [...]”*

*Anonymous Feedback*



## Burst Buffers

*Many have heard about it,  
few have seen real machines with it,  
even fewer have ran applications on those machines ...*



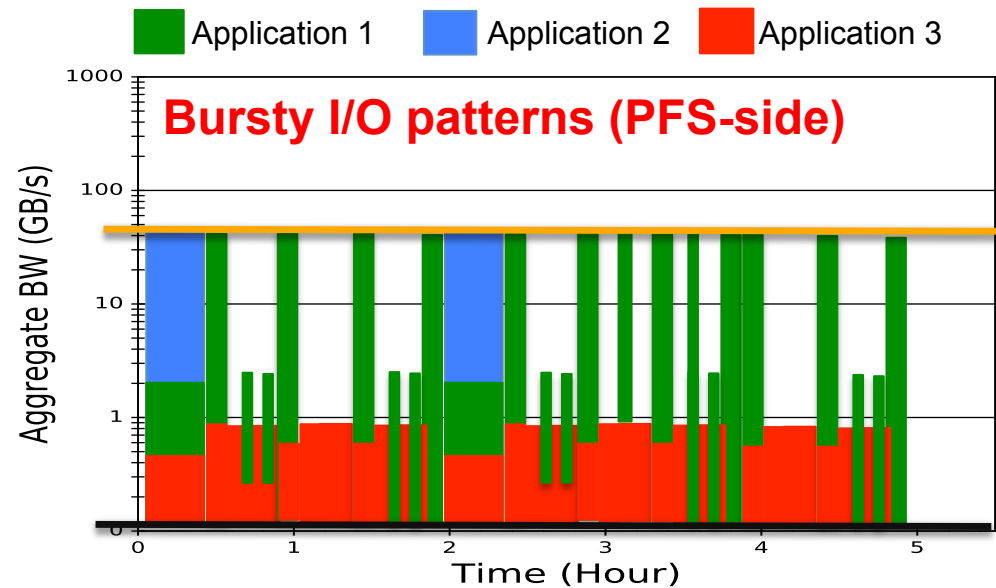
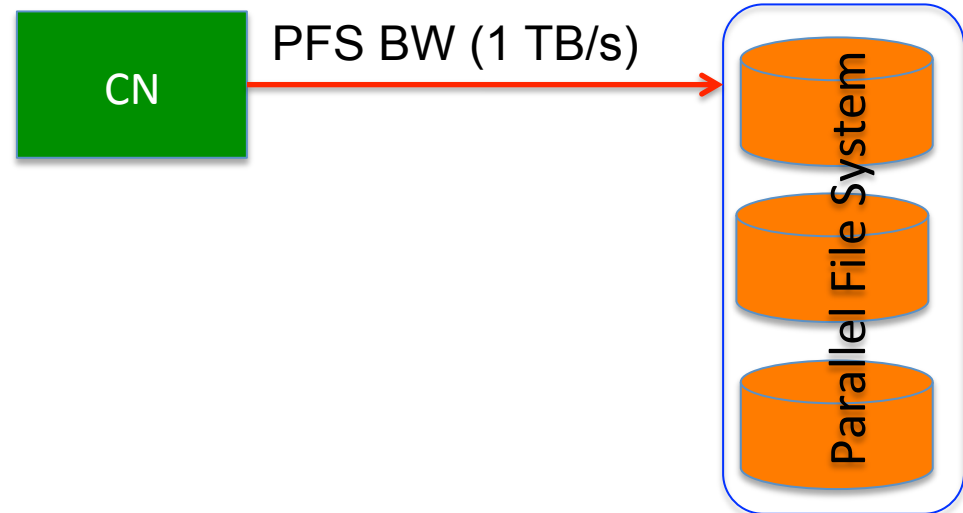
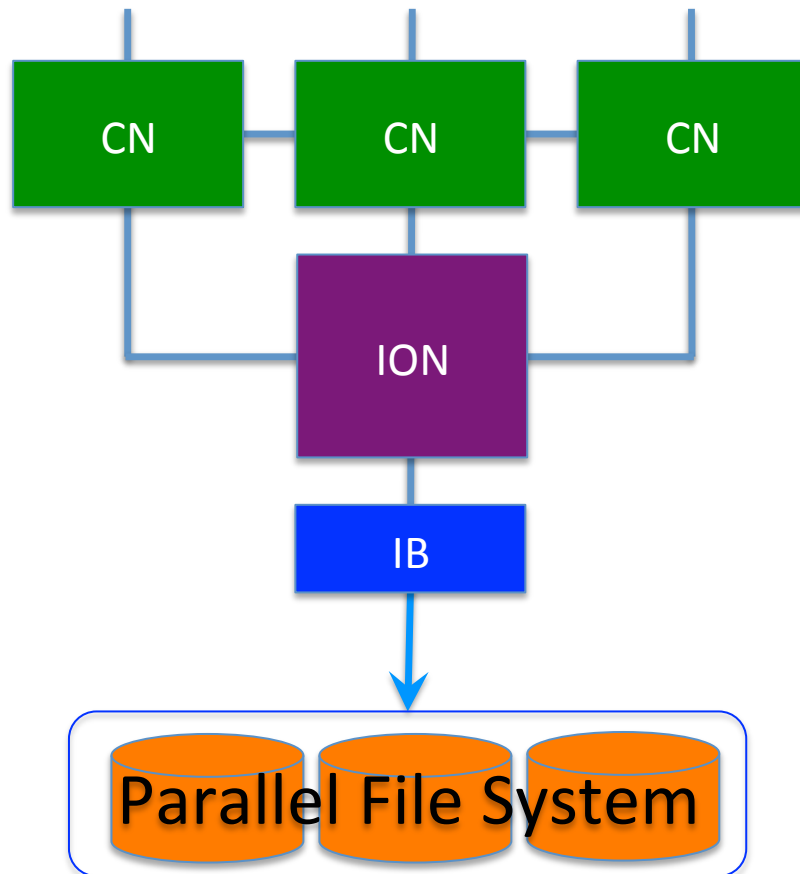
## Challenges

- Burst Buffers are not the magic I/O silver bullet
  - I/O contention still a problem if we exceed the burst buffer capability
  - Burst buffers improve offloading bandwidth but do **NOT** help uploading data from storage for analysis and visualization





# Traditional System

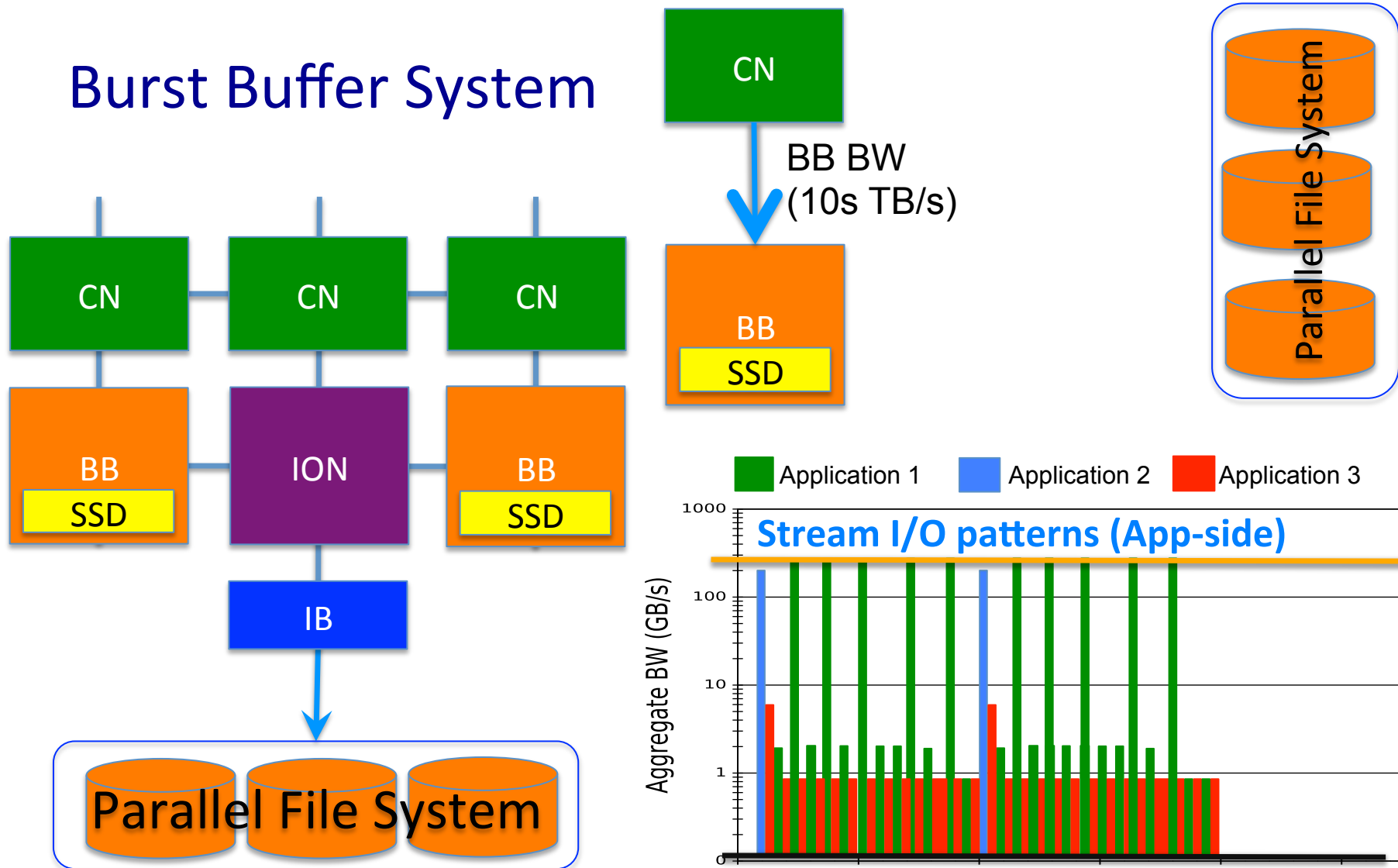


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Based on: Liu, N, Cope, J, Carns, P, Carothers, C, Ross, R, Grider, G, Crume, A, Maltzahn, C. "On the Role of Burst Buffers in Leadership-class Storage Systems" MSST/SNAPI 2012



# Burst Buffer System

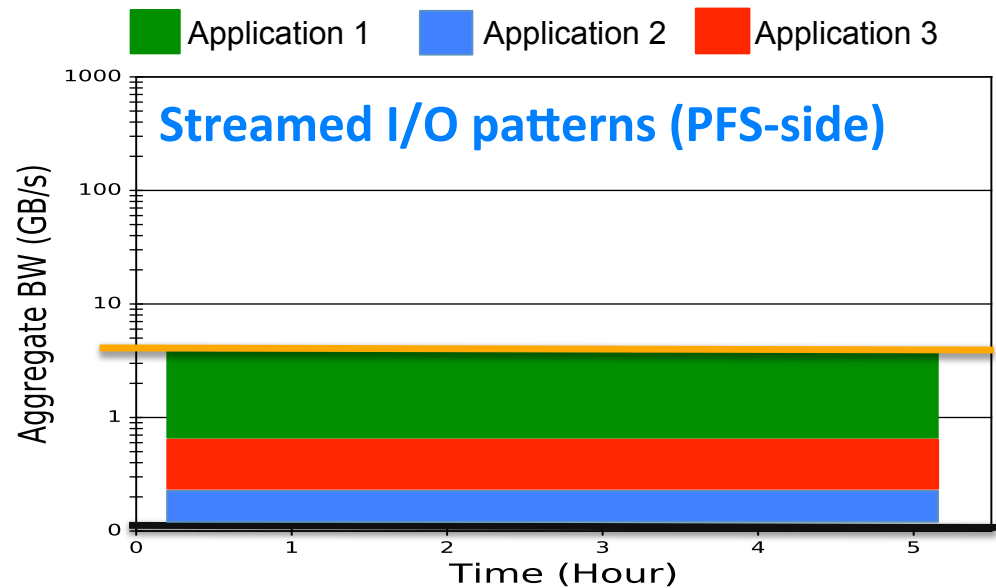
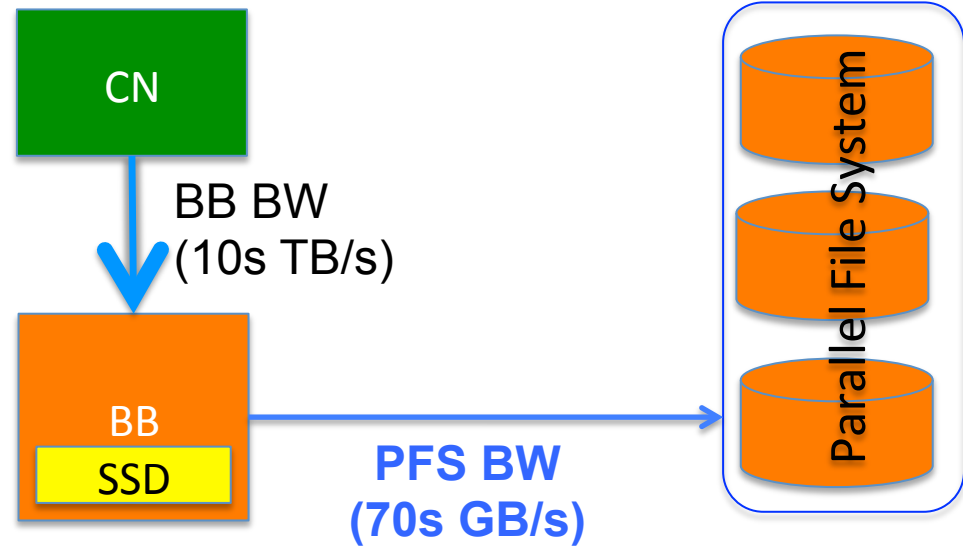
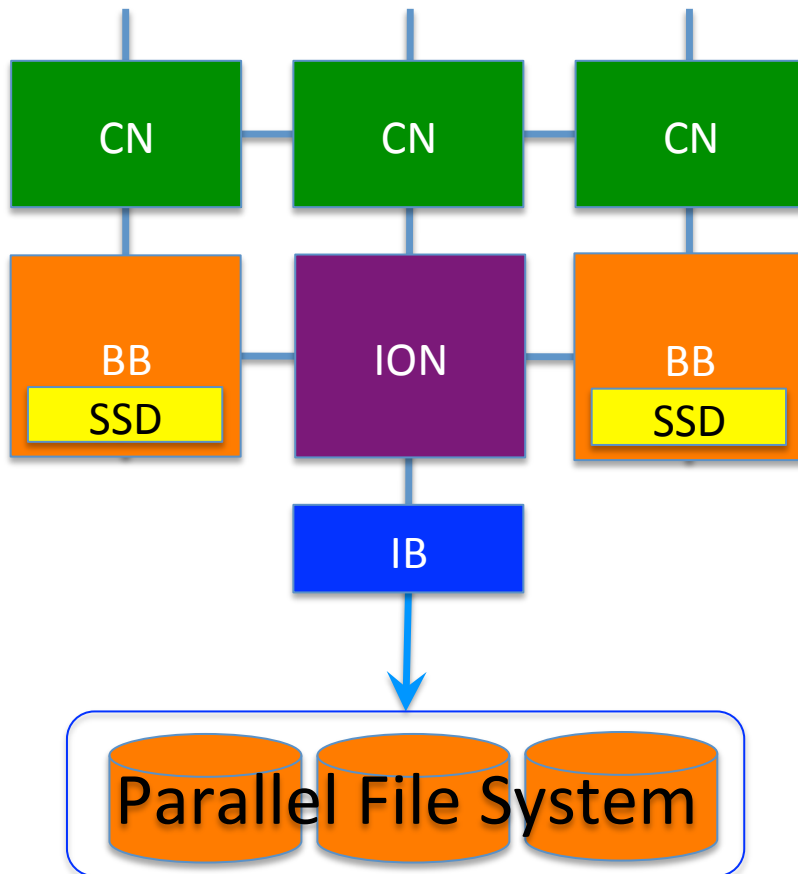


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Based on: Liu, N, Cope, J, Carns, P, Carothers, C, Ross, R, Grider, G, Crume, A, Maltzahn, C. "On the Role of Burst Buffers in Leadership-class Storage Systems" MSST/SNAPI 2012



# Burst Buffer System

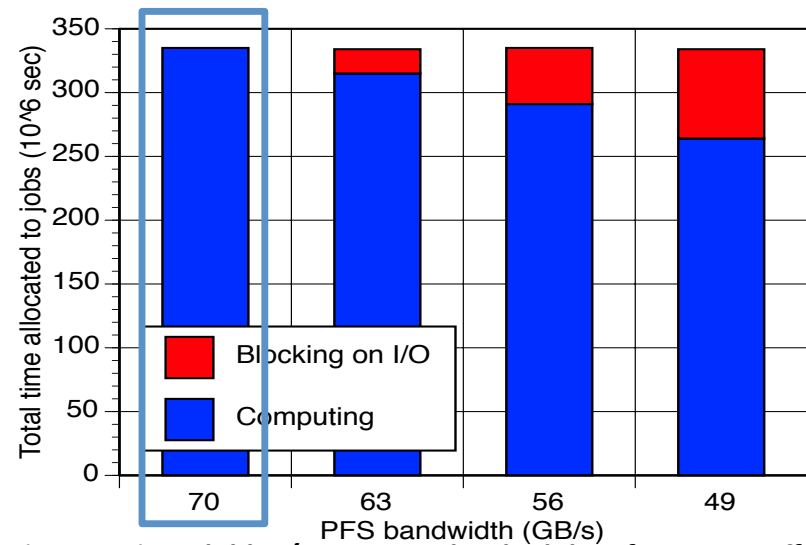
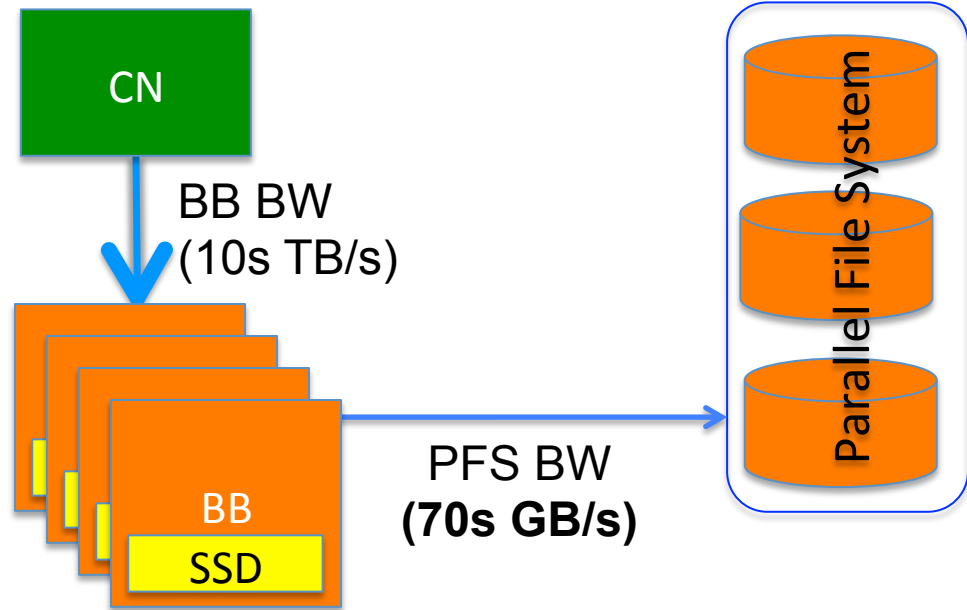
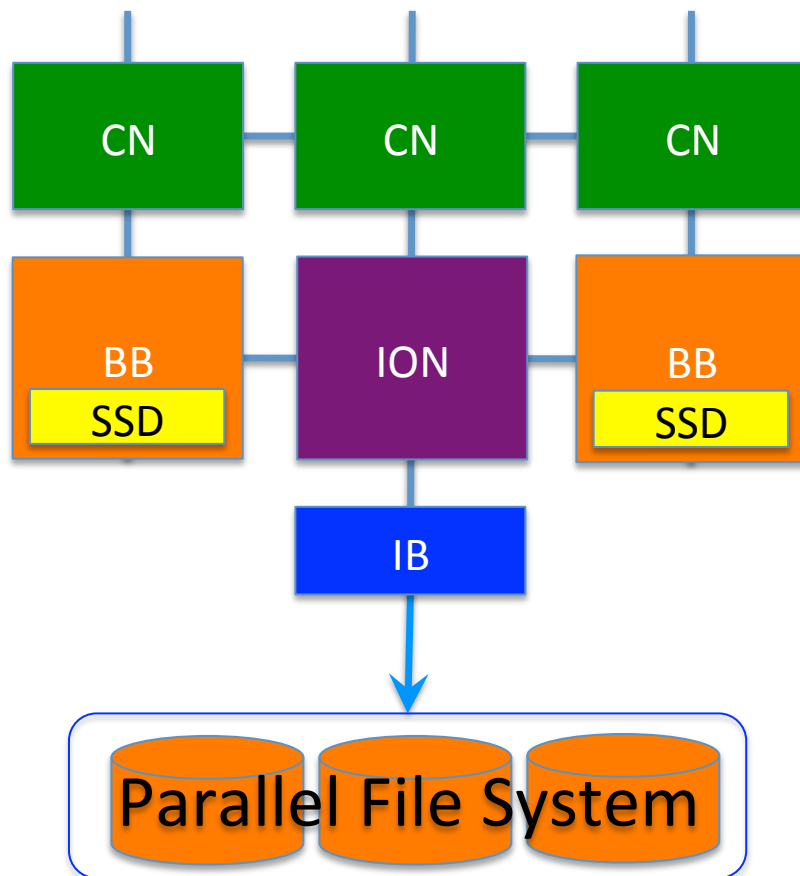


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Based on: Liu, N, Cope, J, Carns, P, Carothers, C, Ross, R, Grider, G, Crume, A, Maltzahn, C. "On the Role of Burst Buffers in Leadership-class Storage Systems" MSST/SNAPI 2012



# PFS Bottleneck

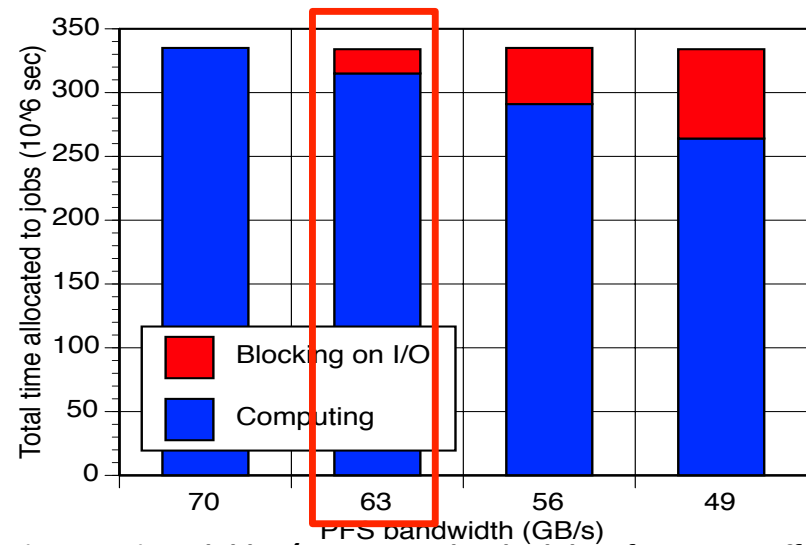
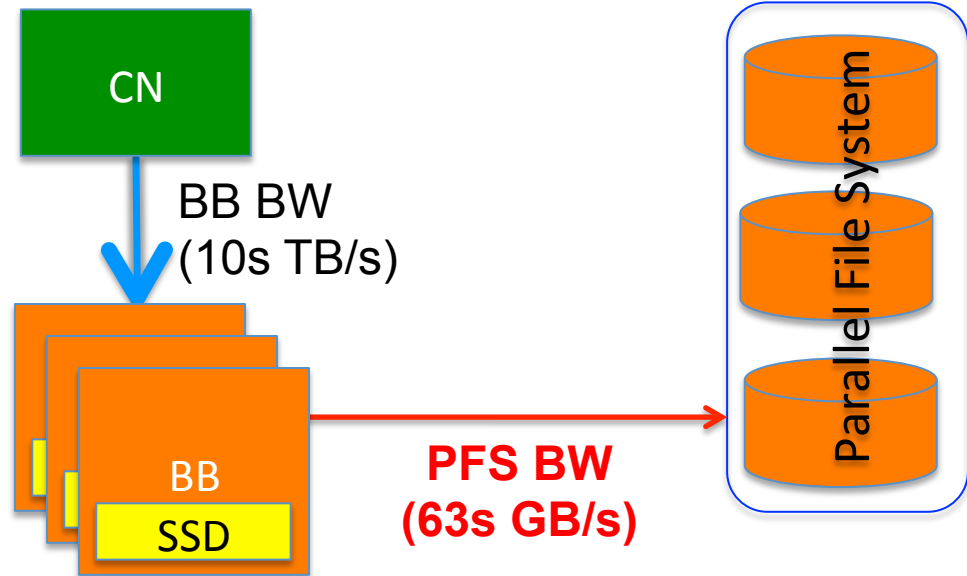
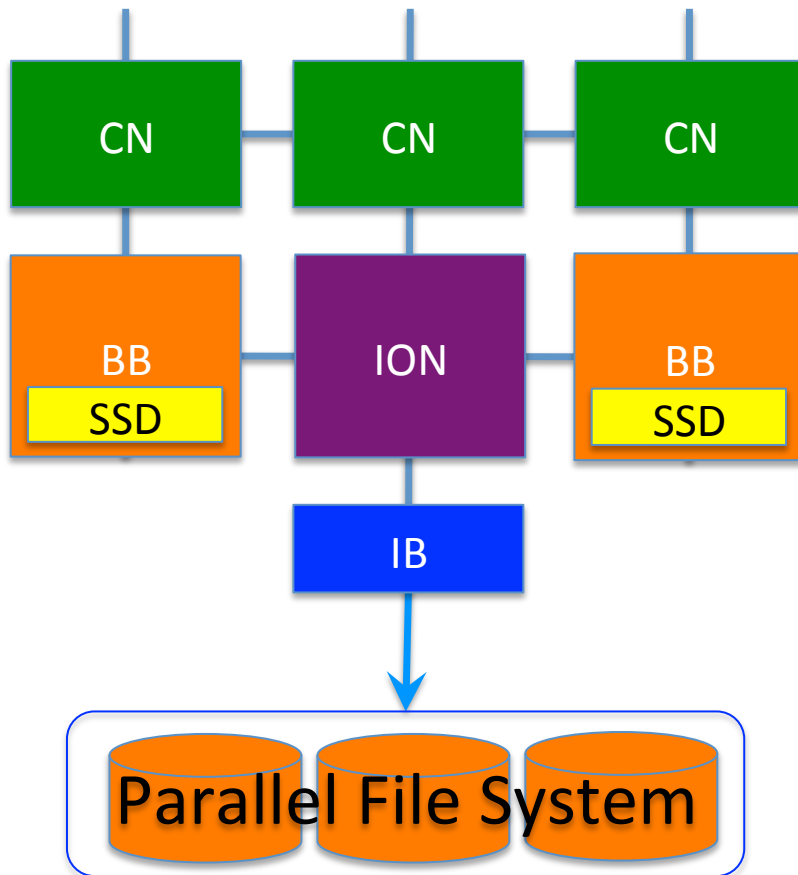


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Herbein et al. *Scalable I/O-aware Job Scheduling for Burst Buffer Enabled HPC Clusters*, HPDC 2016.



# PFS Bottleneck

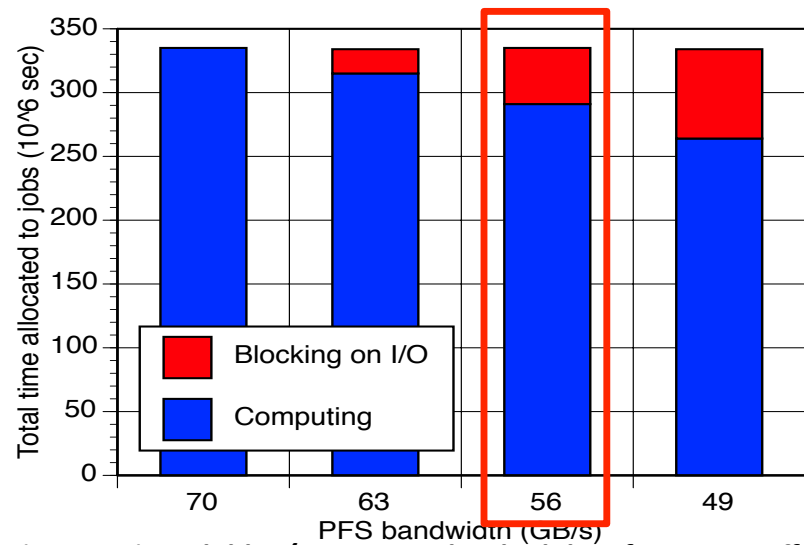
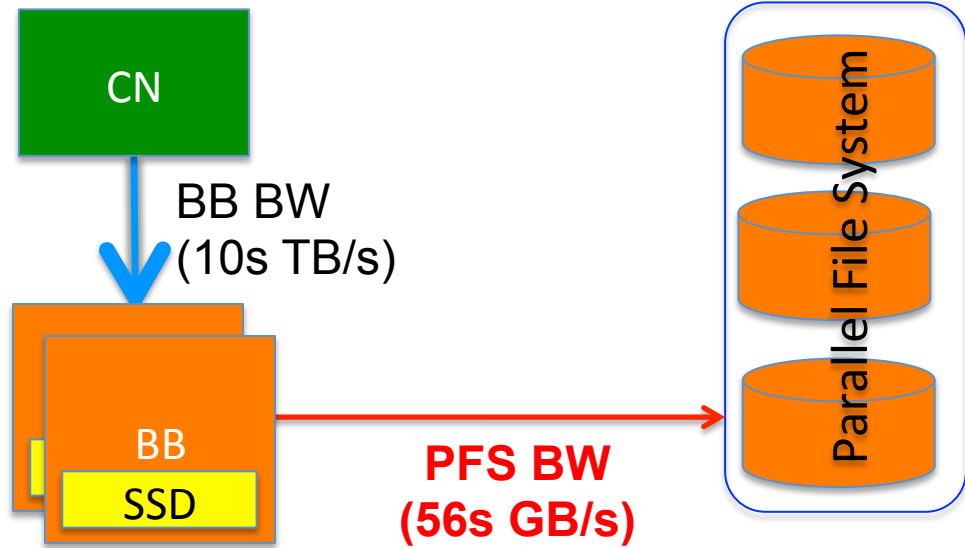
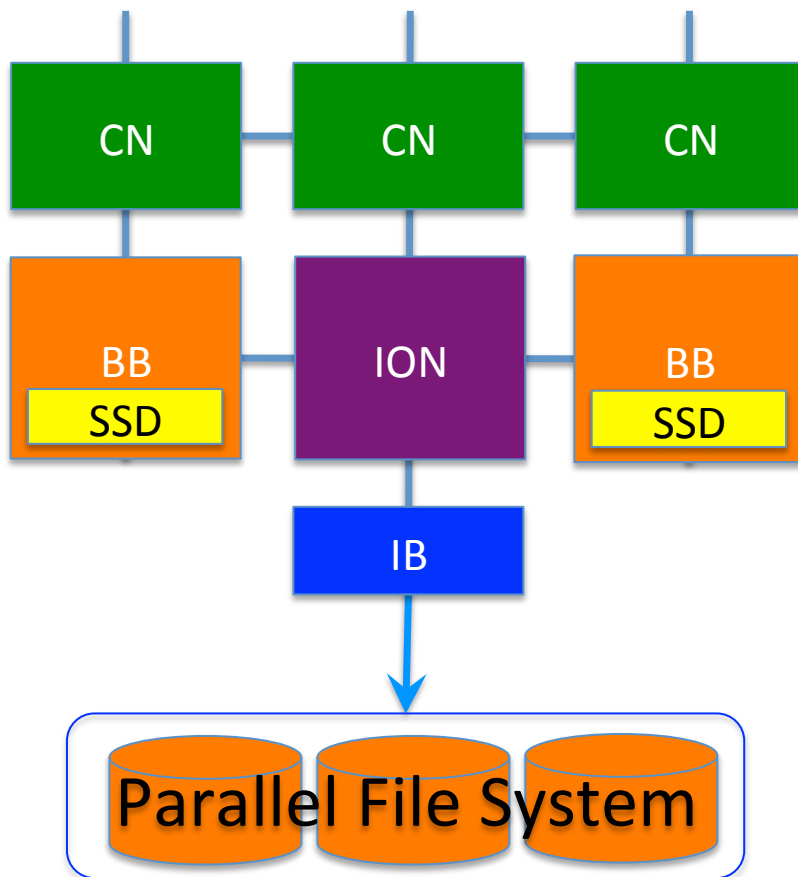


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Herbein et al. *Scalable I/O-aware Job Scheduling for Burst Buffer Enabled HPC Clusters*, HPDC 2016.



# PFS Bottleneck

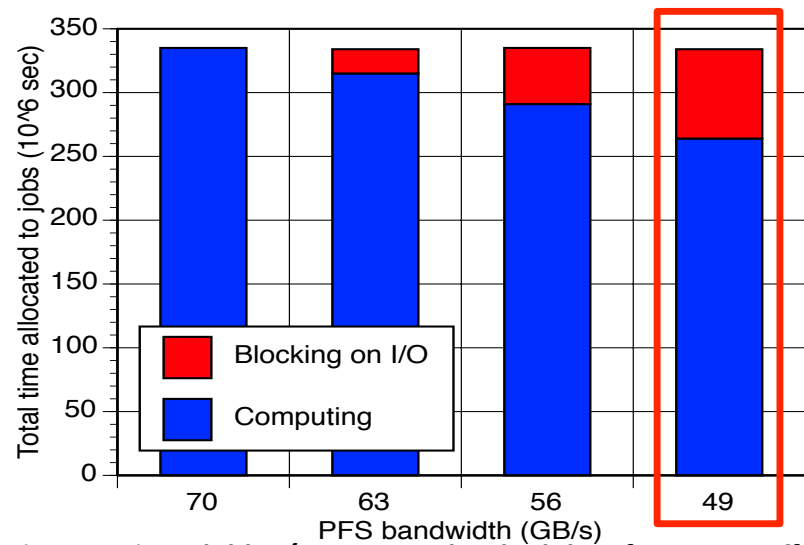
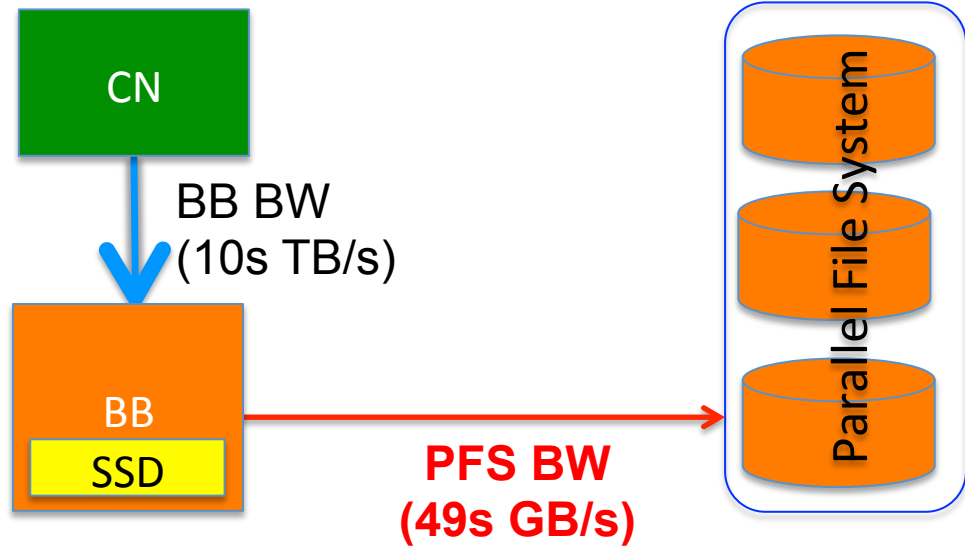
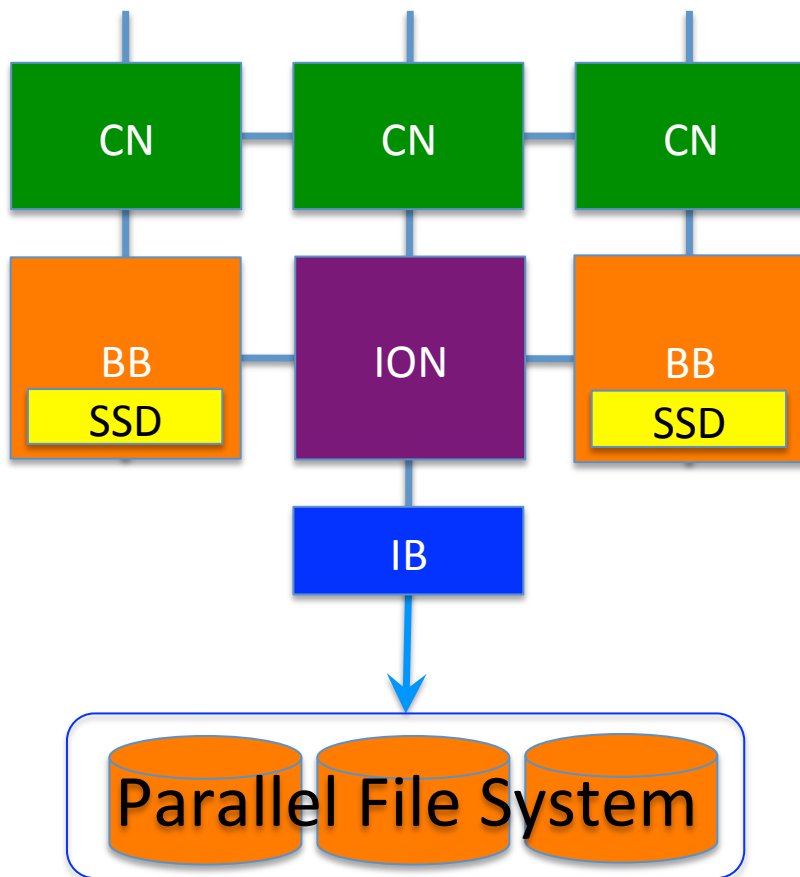


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Herbein et al. *Scalable I/O-aware Job Scheduling for Burst Buffer Enabled HPC Clusters*, HPDC 2016.



# PFS Bottleneck

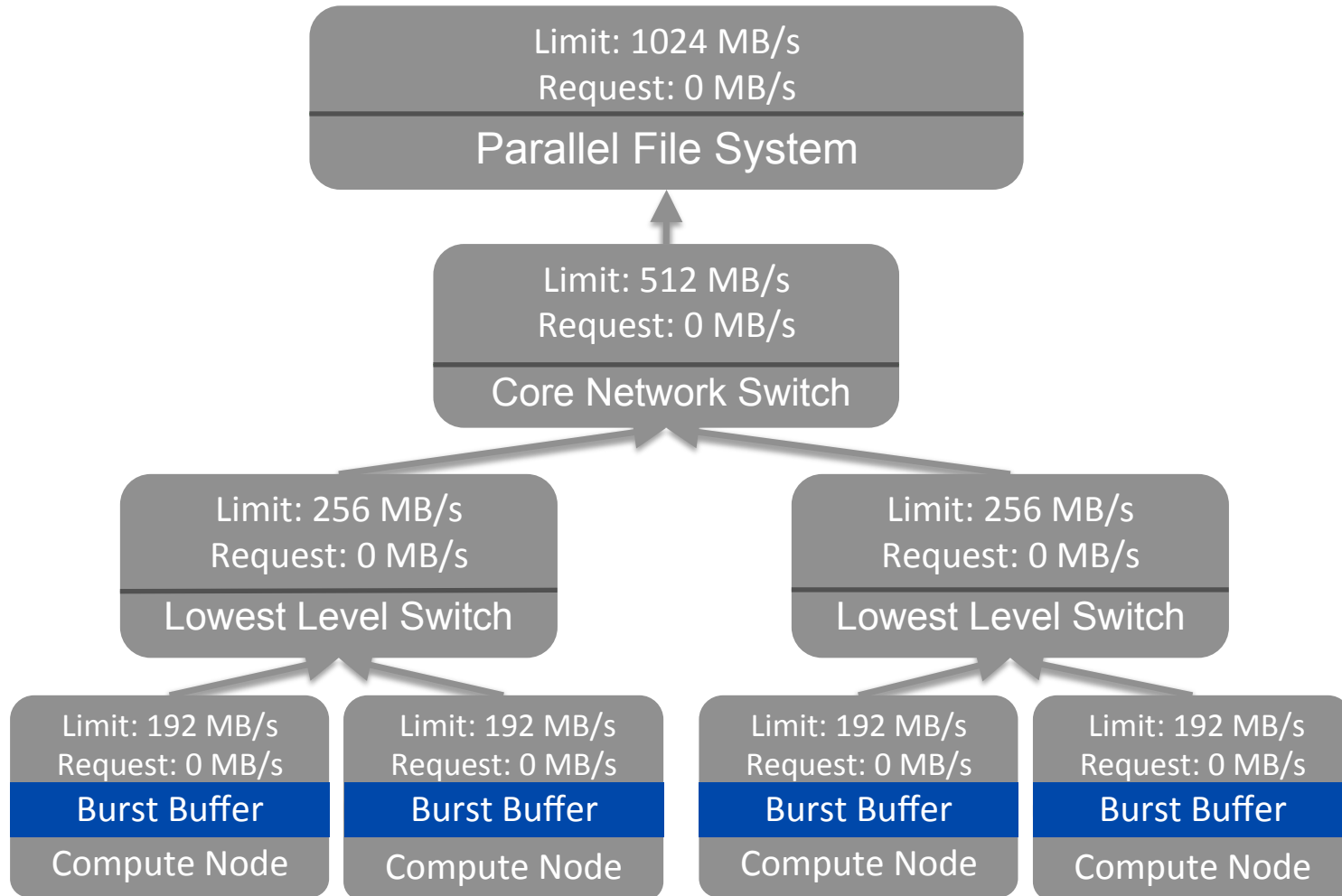


Based on: <http://www.nersc.gov/users/computational-systems/cori/burst-buffer/burst-buffer/>

Herbein et al. *Scalable I/O-aware Job Scheduling for Burst Buffer Enabled HPC Clusters*, HPDC 2016.



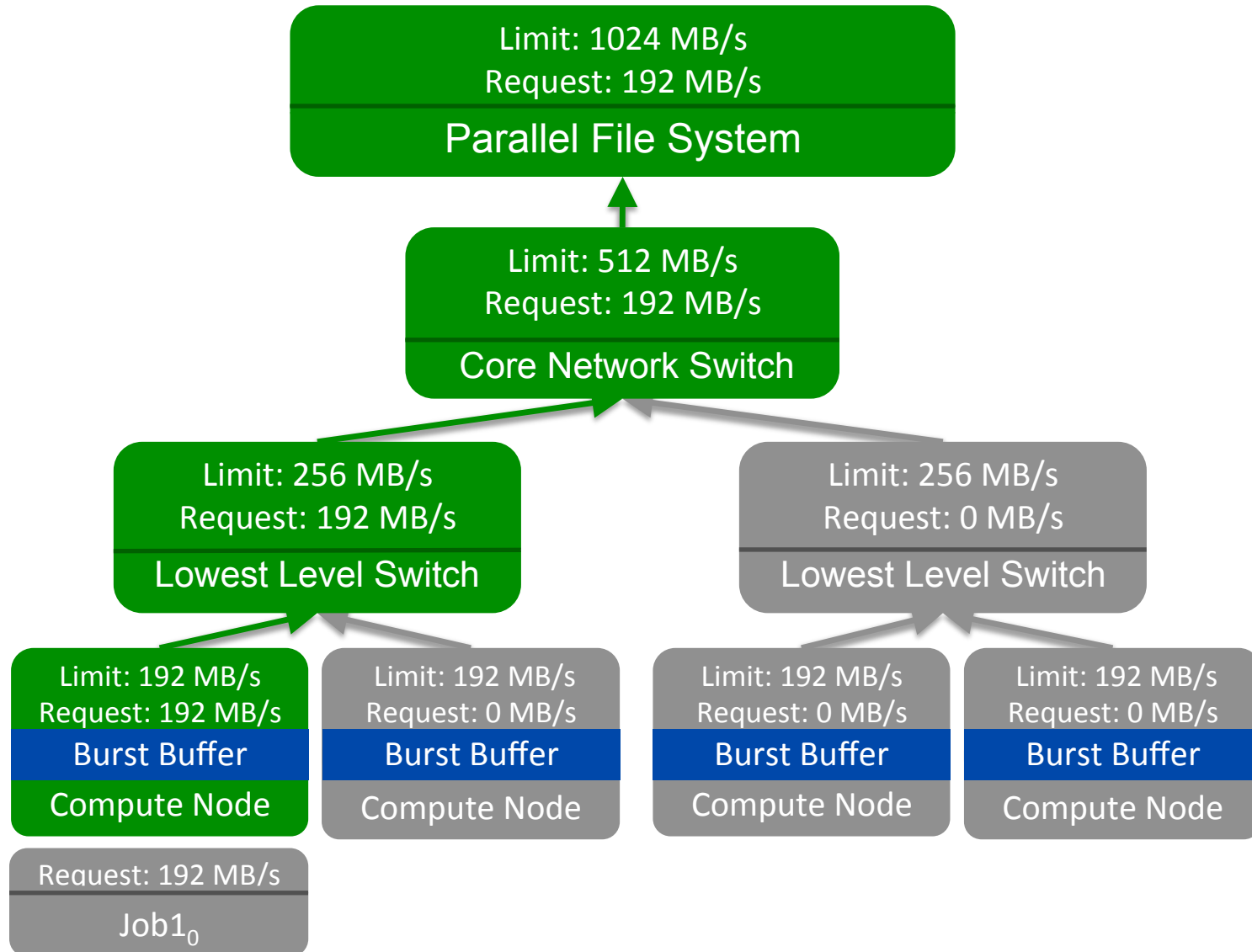
# Integrate I/O-awareness in Flux Scheduler





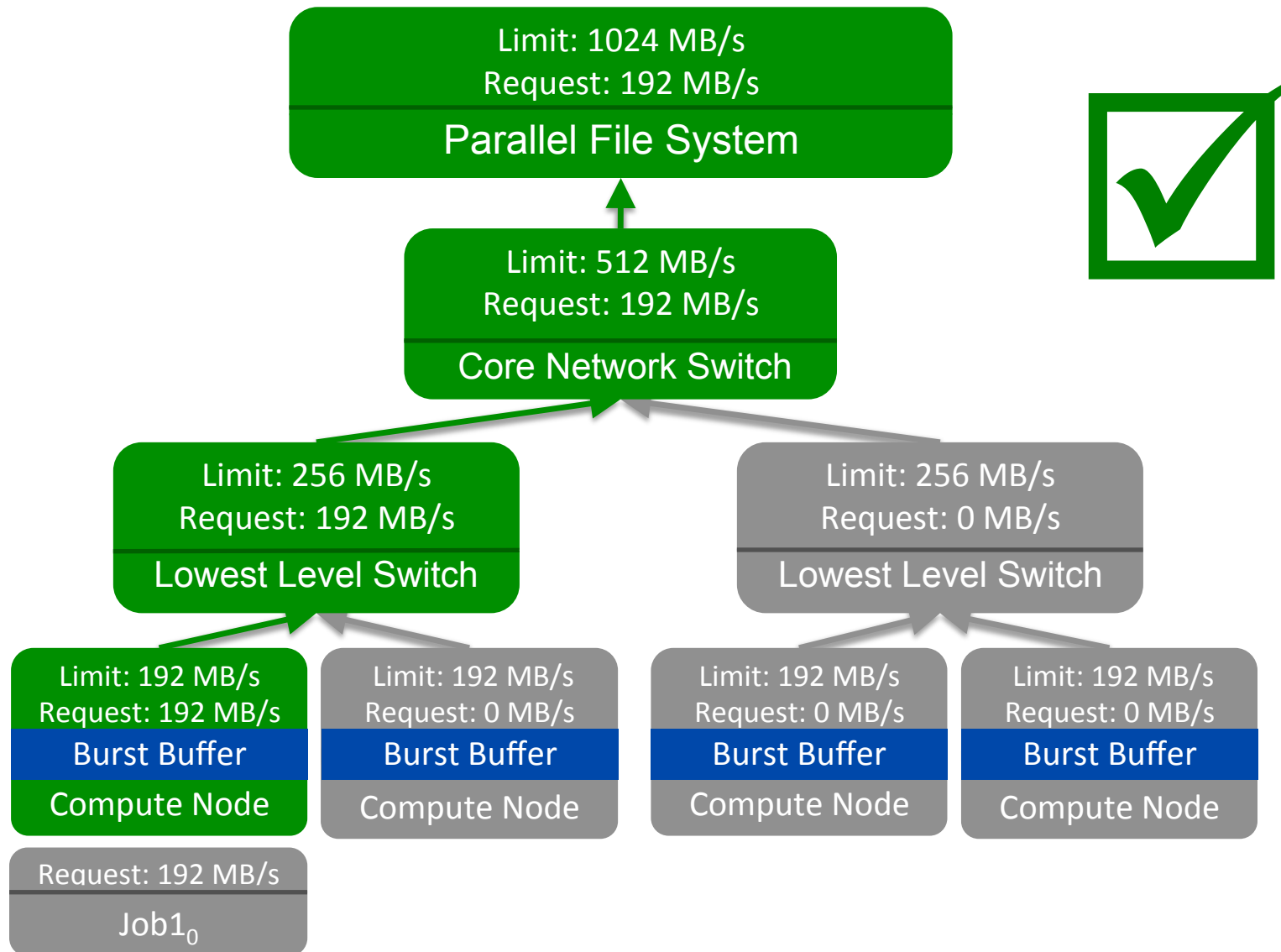


# Integrate I/O-awareness in Flux Scheduler



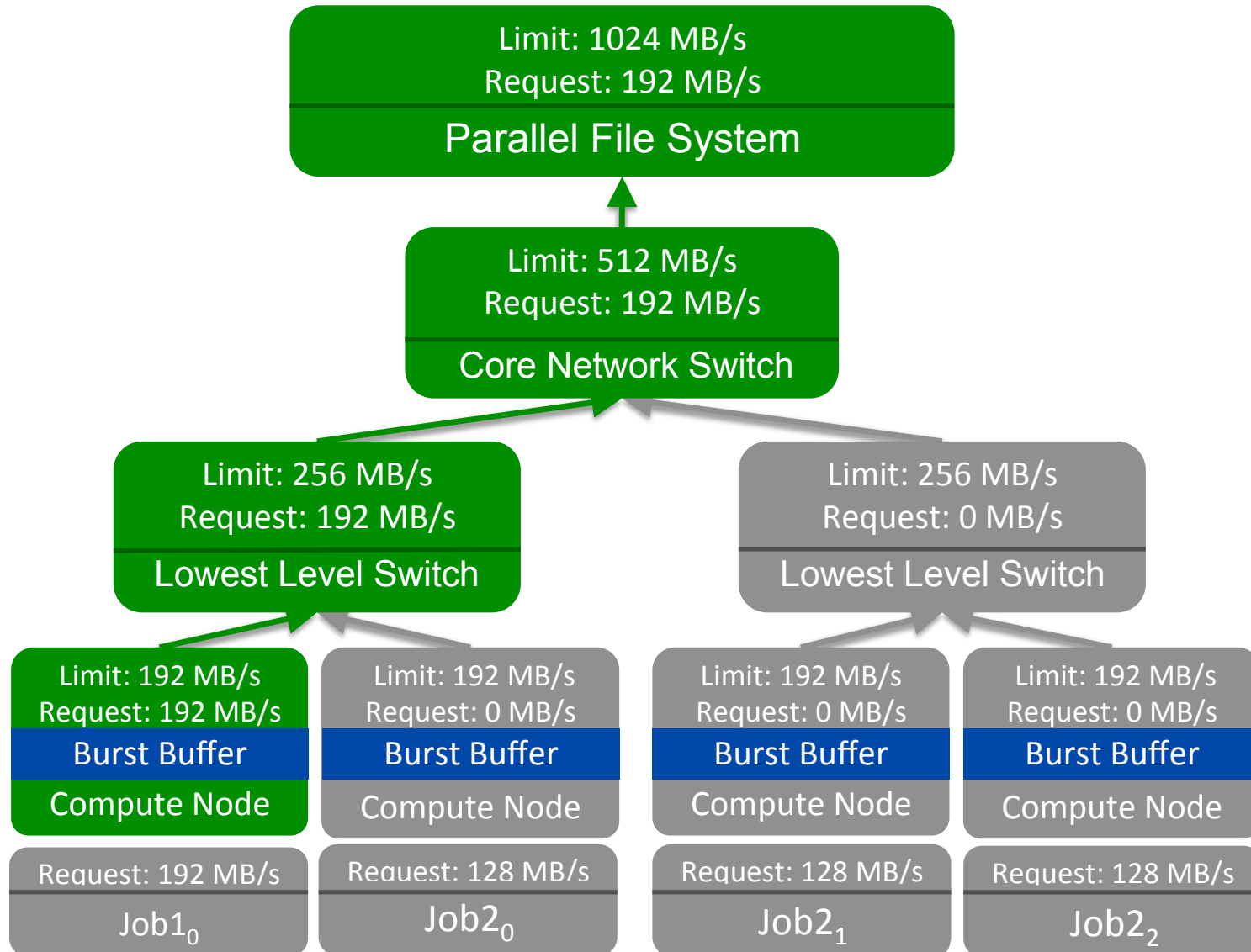


# Integrate I/O-awareness in Flux Scheduler



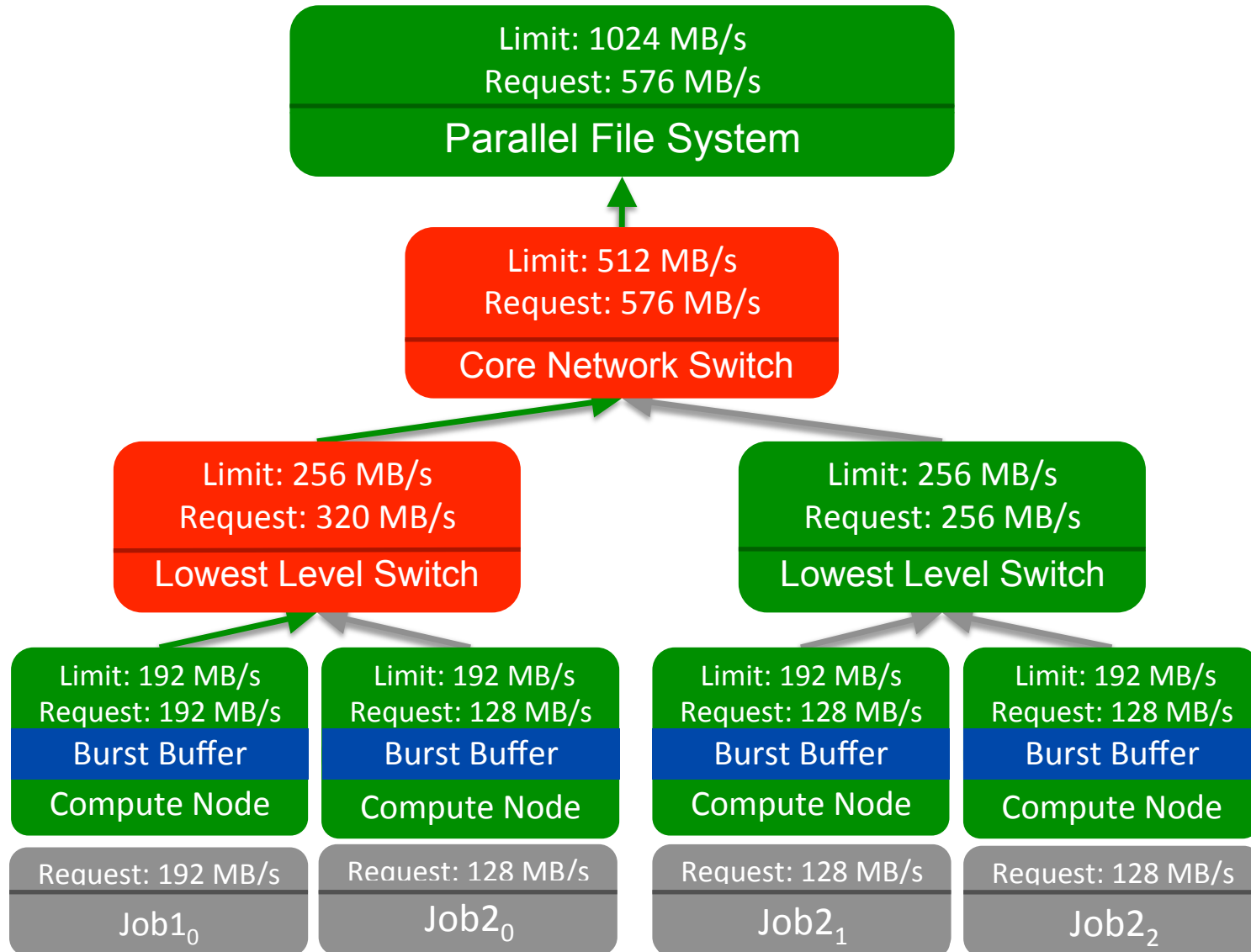


# Integrate I/O-awareness in Flux Scheduler



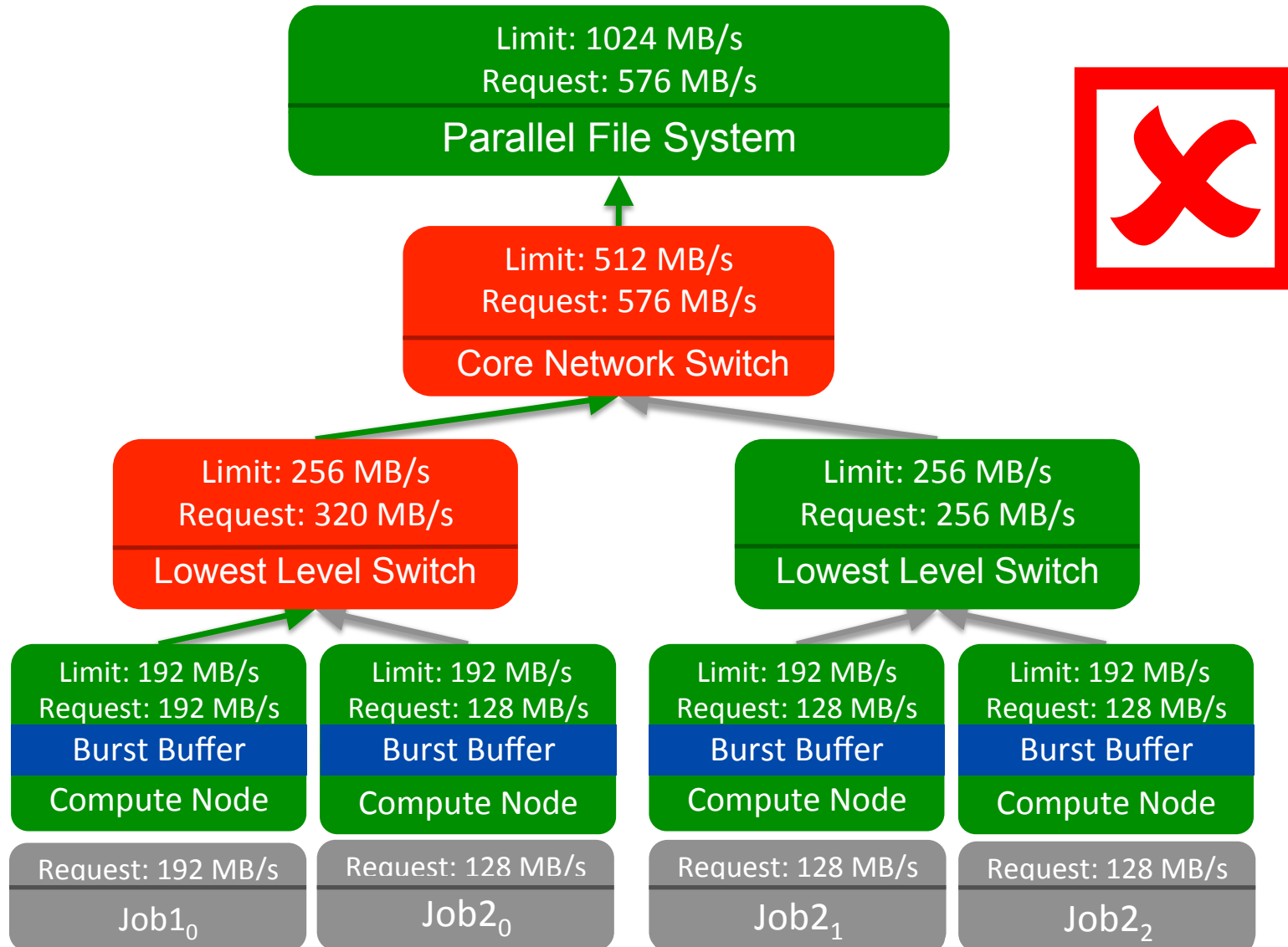


# Integrate I/O-awareness in Flux Scheduler





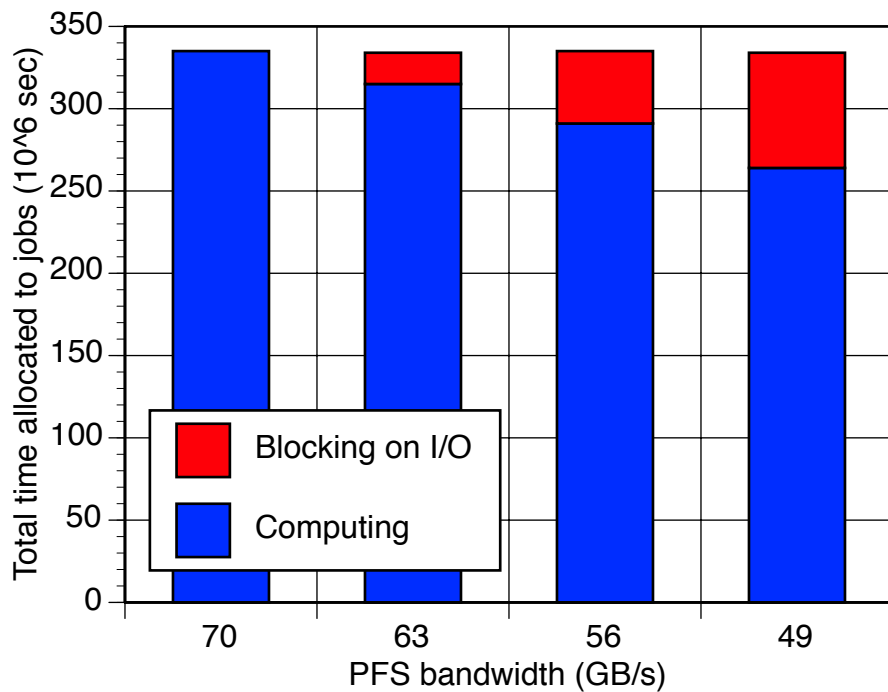
# Integrate I/O-awareness in Flux Scheduler



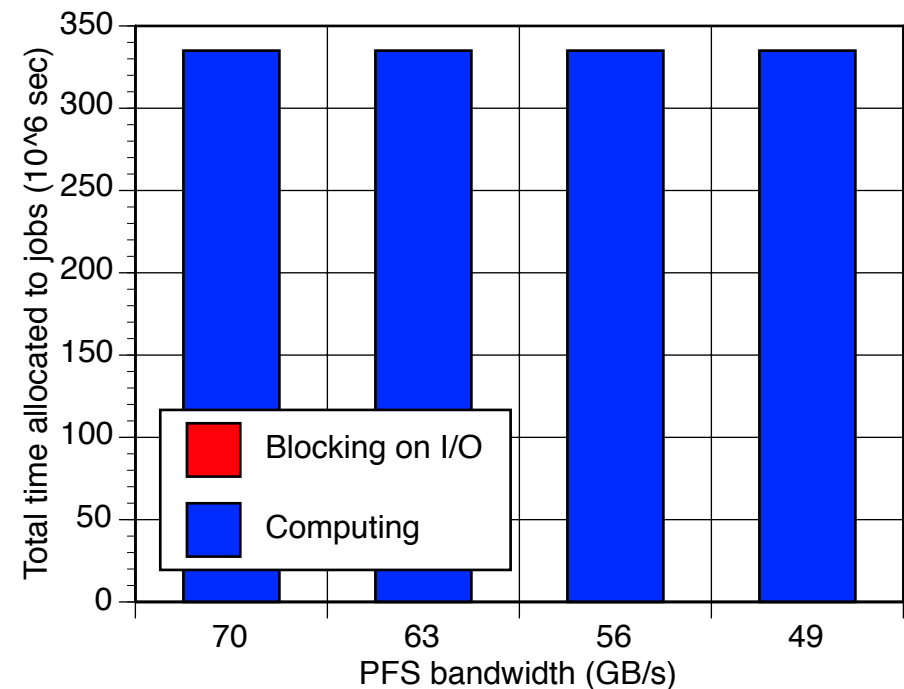


## I/O-Ignorant vs. I/O-Aware Scheduling in Flux

**I/O-ignorant** scheduler



**I/O-aware** scheduler



***I/O-Aware scheduling results in 100% of application time to be spent in computation***



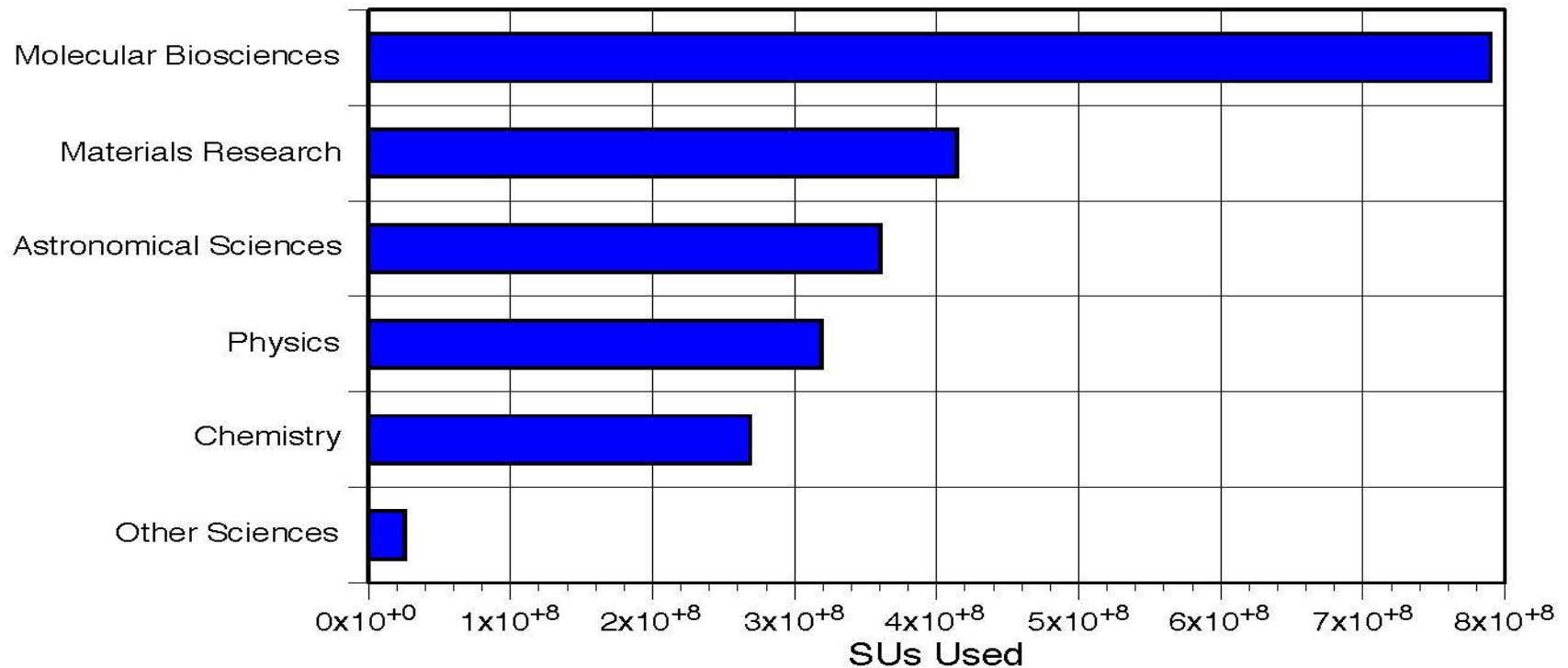
## Challenges

- Burst Buffers are not the magic I/O silver bullet
  - I/O contention still a problem if we exceed the burst buffer capability
  - Burst buffers improve offloading bandwidth but do **NOT** help uploading data from storage for runtime analysis



## MD Simulations are Alive and Kicking!

XSEDE SUs used by type of targeted science over the past 6 months (March 1, 2016 - August 31, 2016)

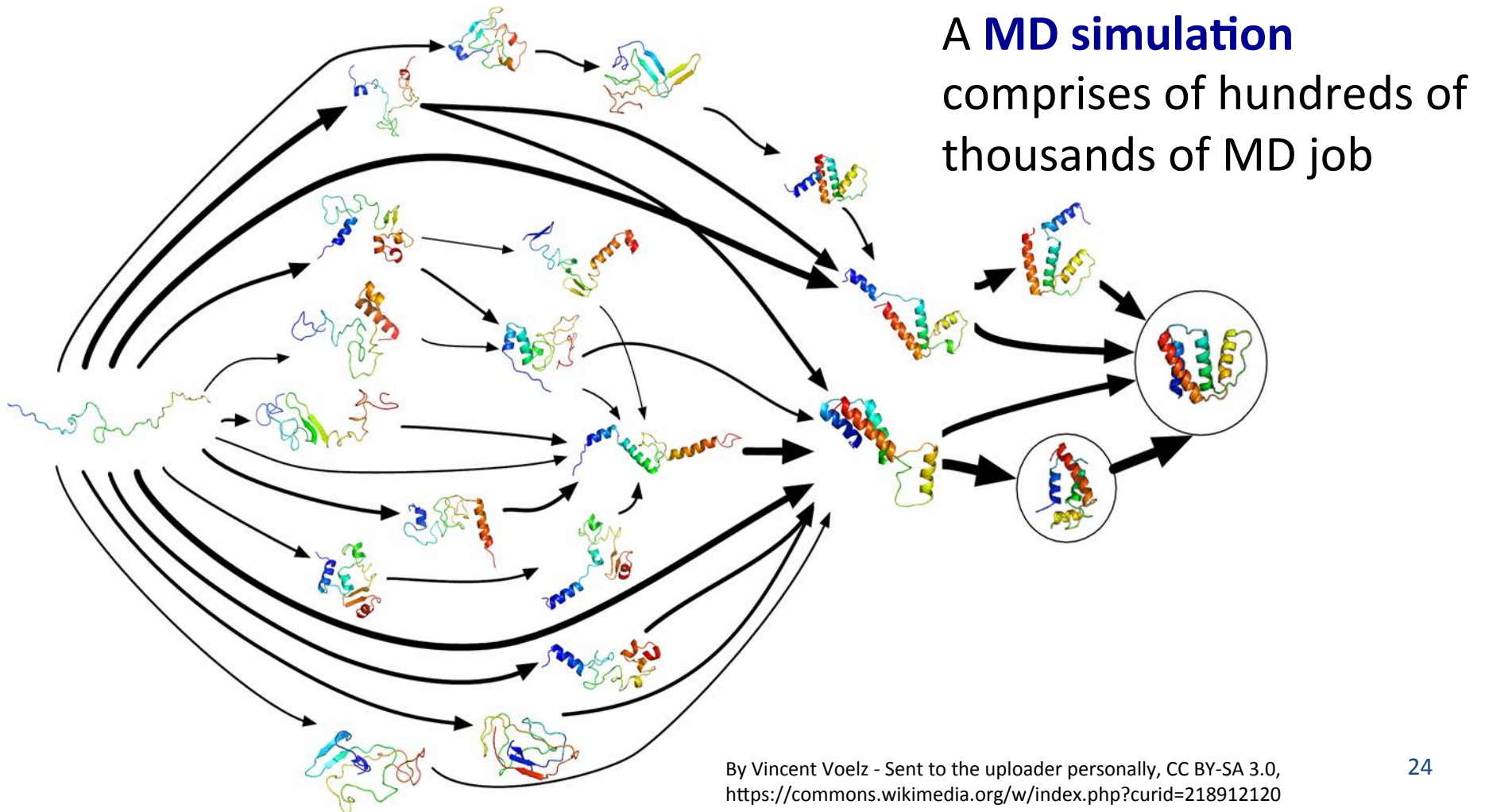


Four of the top 10 XSEDE users run molecular simulations (i.e., Schulten at UIUC, Feig at Michigan State U, Voth at U Chicago, and Case at Rutgers U)





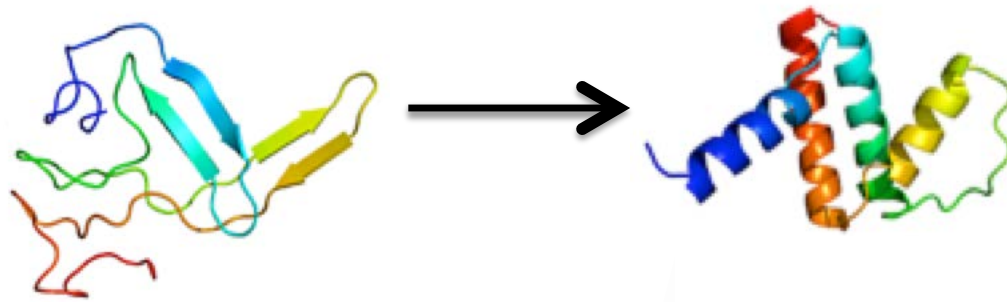
## MD Simulations as an Ensemble of HPC Jobs



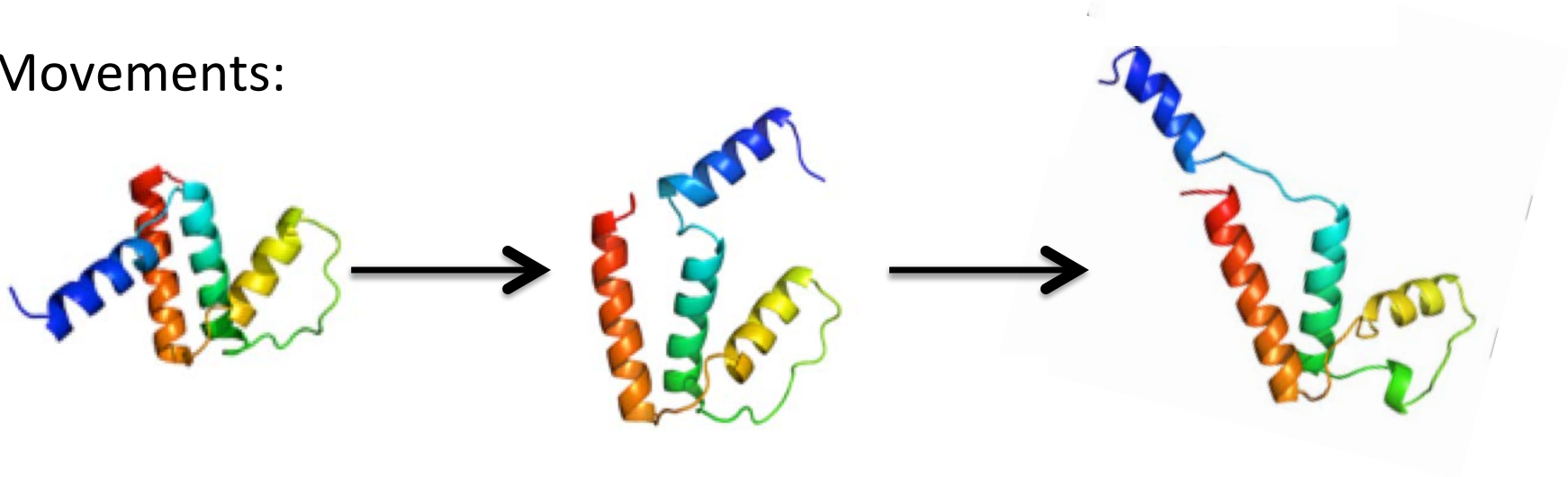


## Capturing Rare Events

Transformations:

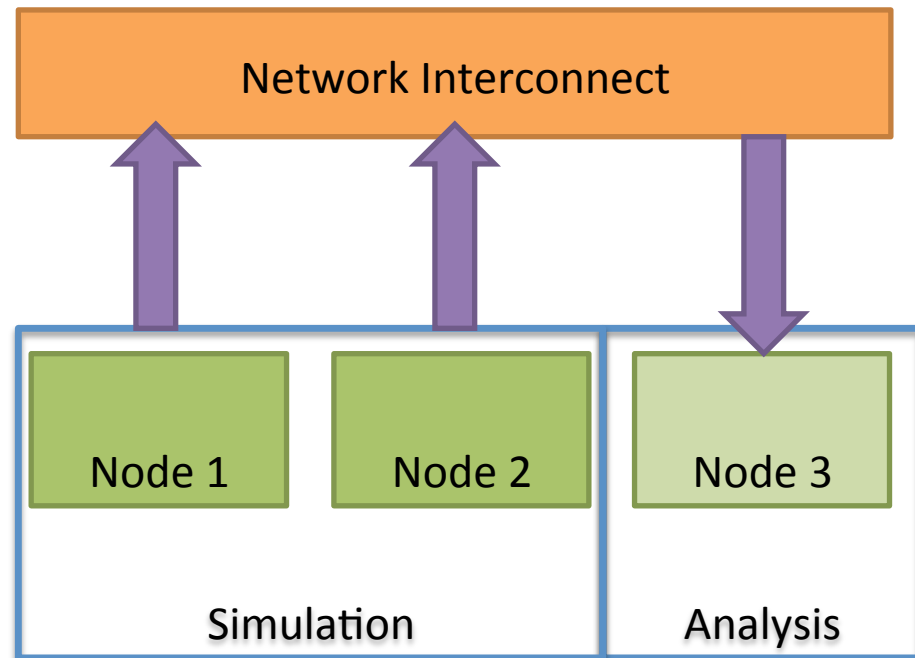
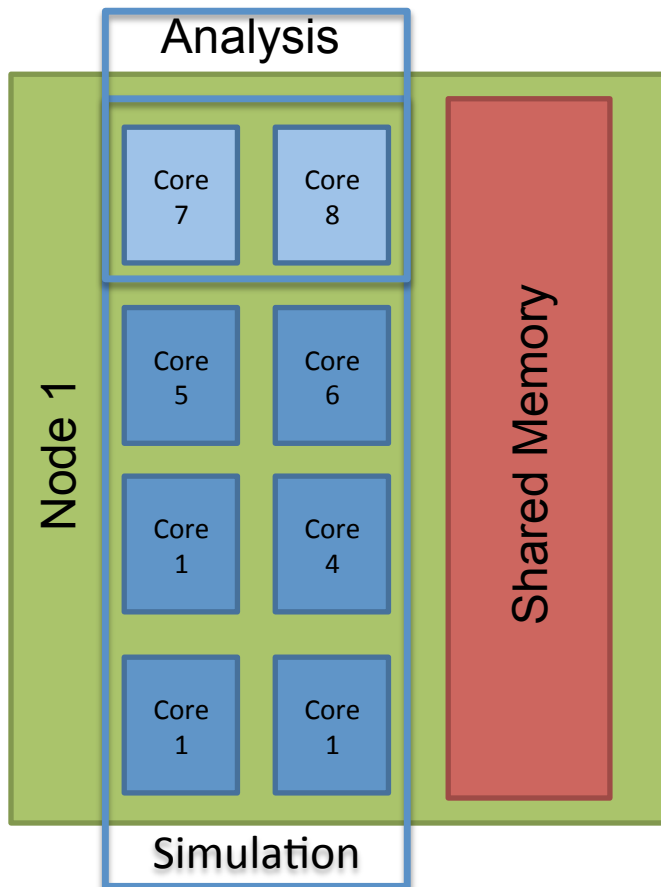


Movements:





## *In-situ and In-transit Analysis*



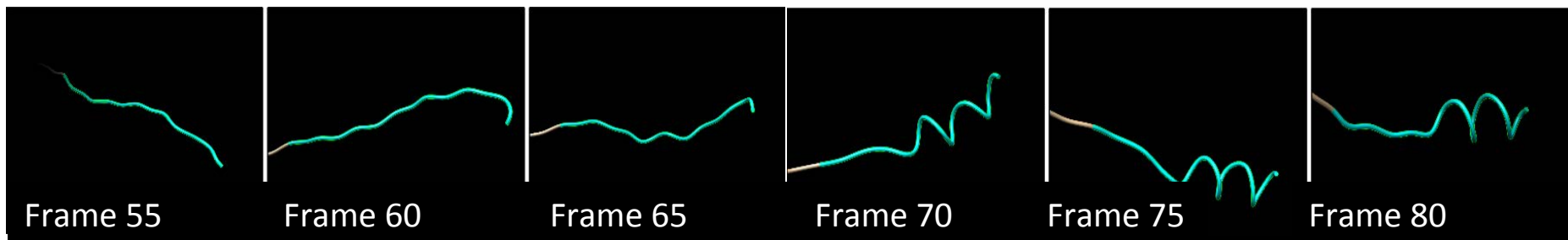
Example of tools:

- DataSpaces (Rutgers U.)
- DataStager (GeorgiaTech)



## Requirements to Capture Rare Events *In Situ*

Frames (or snapshots) of an MD trajectory:

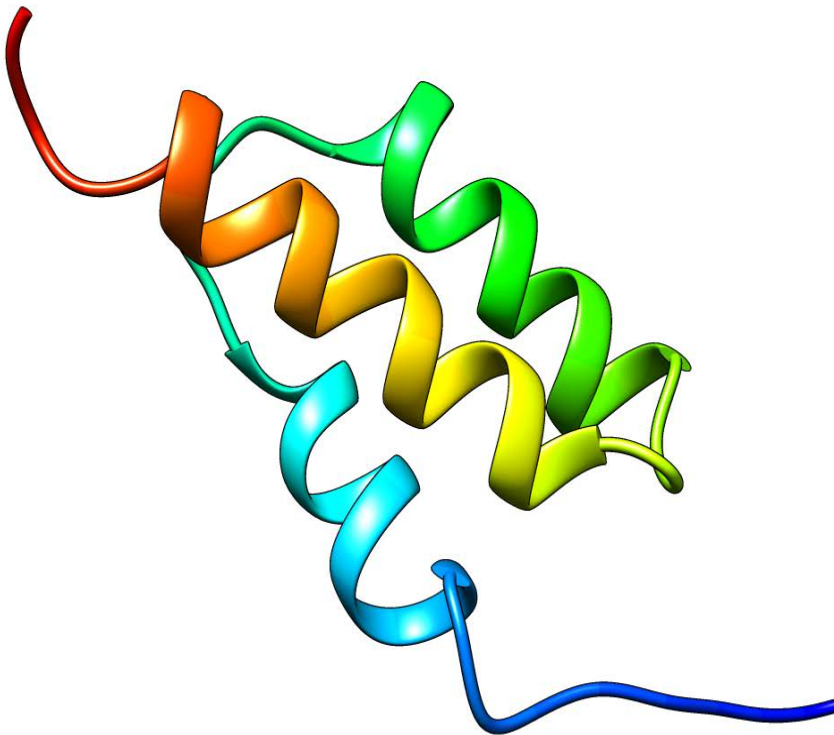


- We want to capture what is going on in each frame **without**:
  - Disrupting the simulation (e.g., stealing CPU and memory on the node)
  - Moving all the frames to a central file system and analyzing them once the simulation is over
  - Comparing each frame with past frames of the same job
  - Comparing each frame with frames of other jobs



## Capturing the Transformations in a Structure

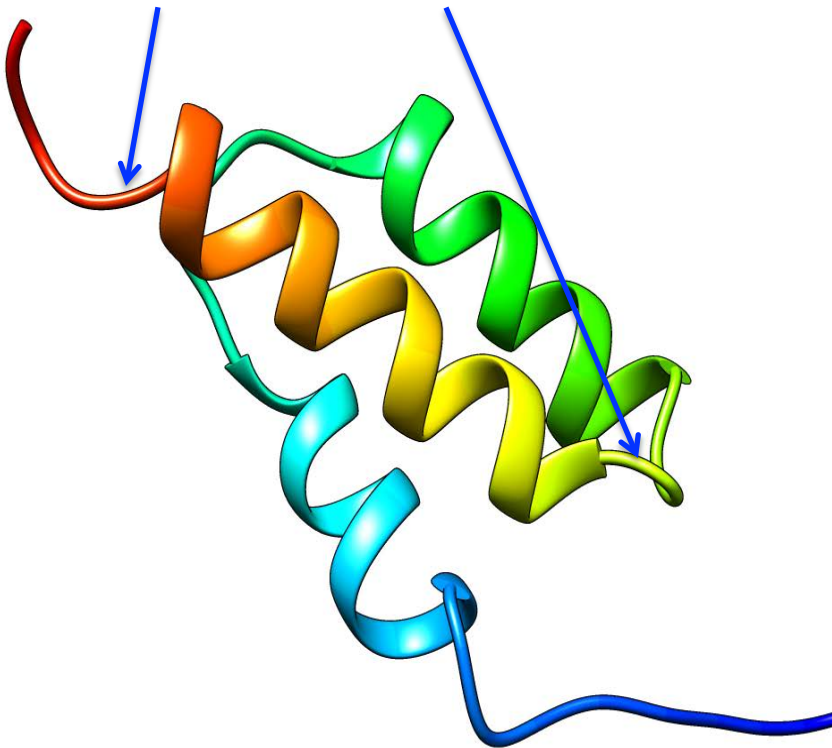
Given a **frame** of an  
MD job **at time  $t$**





## Capturing the Transformations in a Structure

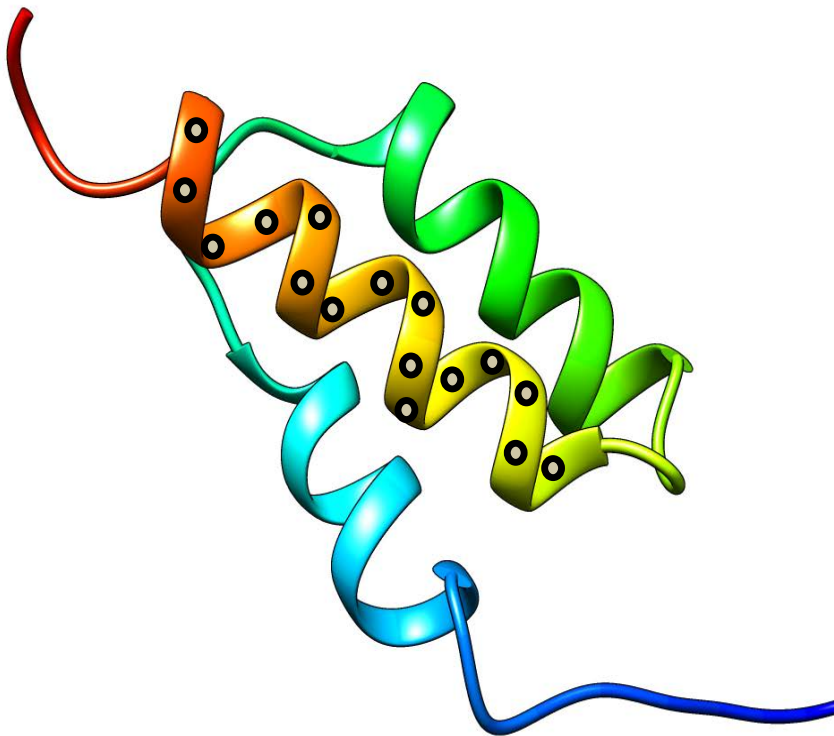
Define the substructure:  
**start** and **stop** amino acids





## Capturing the Transformations in a Structure

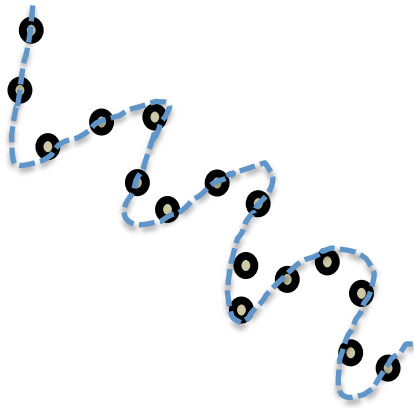
Drop all but not the backbone atoms of the structure ( $C^\alpha$  atoms)





## Capturing the Transformations in a Structure

Drop all but not the backbone atoms of the structure ( $C^\alpha$  atoms)

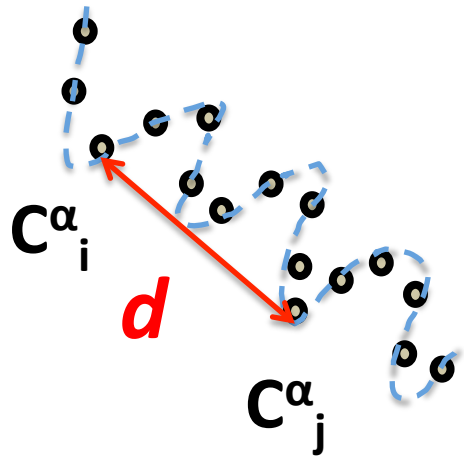






## Capturing the Transformations in a Structure

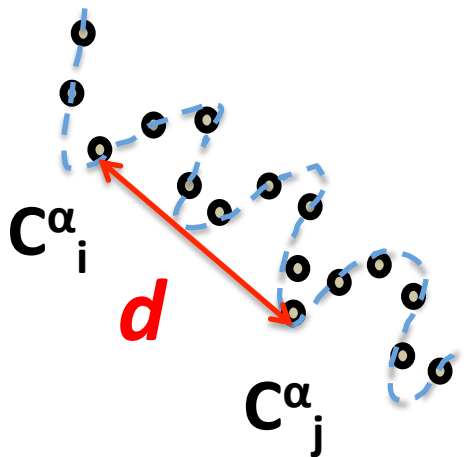
Measure the distance  
between  $C^{\alpha}_j$  and  $C^{\alpha}_i$





## Capturing the Transformations in a Structure

Measure the distance between  $C^{\alpha}_j$  and  $C^{\alpha}_i$



Build the **substructure Euclidean Distance Matrix (D)**

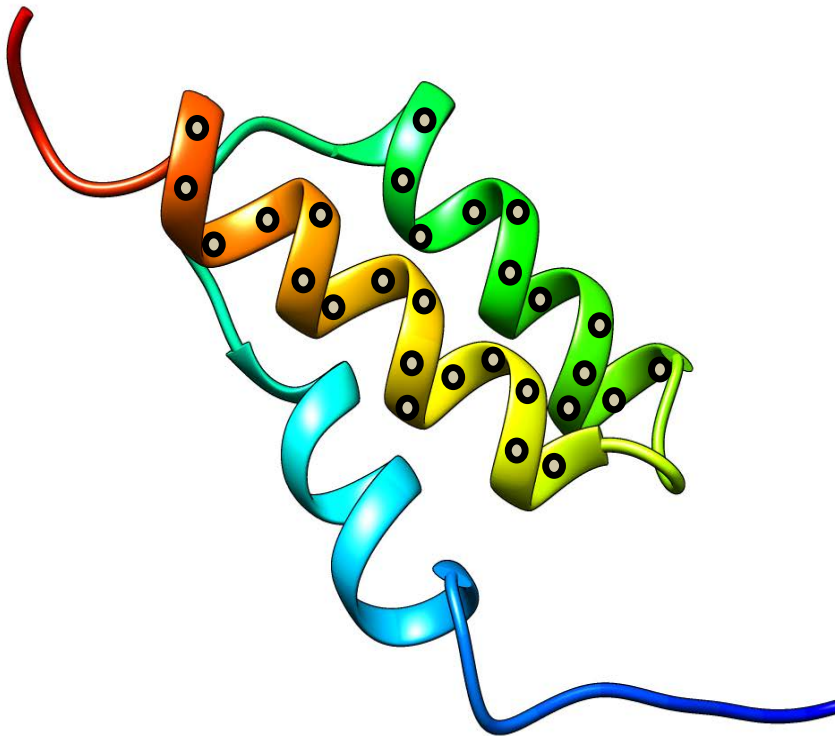
$$D = \begin{matrix} & & C^{\alpha}_i & & & & \\ C^{\alpha}_j & \begin{bmatrix} 0 & \times & \times & \times & \times & \times \\ \times & 0 & d & \times & \times & \times \\ \times & d & 0 & \times & \times & \times \\ \times & \times & \times & 0 & \times & \times \\ \times & \times & \times & \times & 0 & \times \\ \times & \times & \times & \times & \times & 0 \end{bmatrix} \end{matrix}$$

Compute largest eigenvalue  $\rightarrow \lambda_{max}$



## Capturing Movements between Structures

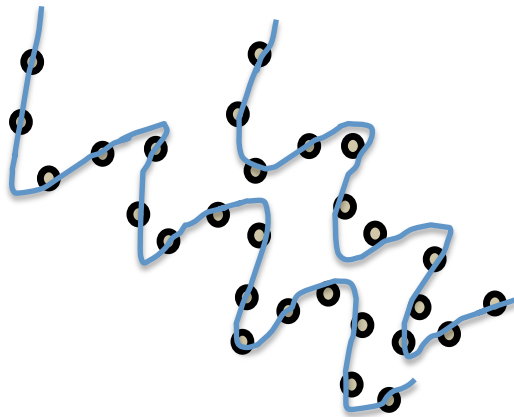
Drop all but not the backbone atoms of the two structure





## Capturing Movements between Structures

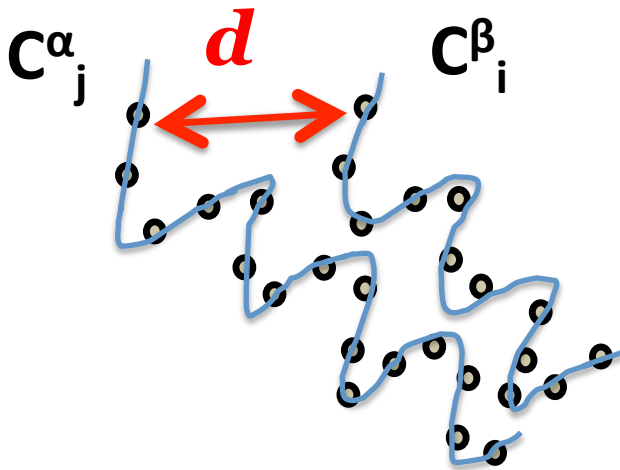
Drop all but not the backbone atoms of the two structure





## Capturing Movements between Structures

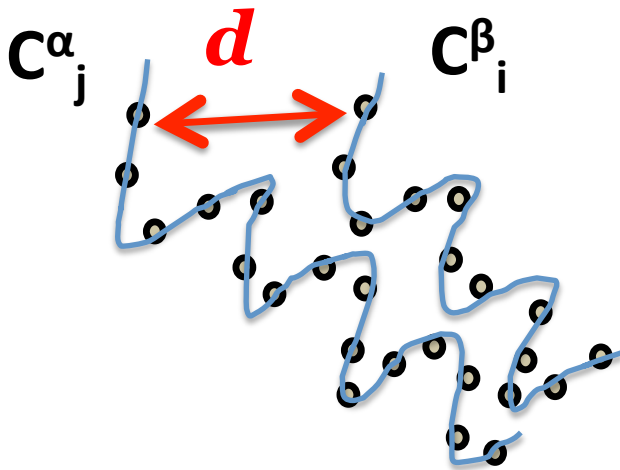
Measure the distance  
between  $C_j^\alpha$  and  $C_i^\beta$





## Capturing Movements between Structures

Measure the distance between  $C^{\alpha}_j$  and  $C^{\beta}_i$



Build a **bipartite distance matrix** by comparing two substructures

$$D = \begin{matrix} & & & i & & & \\ \begin{matrix} j \\ 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & d & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ \times & d & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \end{matrix} & \end{matrix}$$

Compute largest eigenvalue  $\rightarrow \lambda_{max}$



## Proxy for Rare Events

Frames of an MD job:

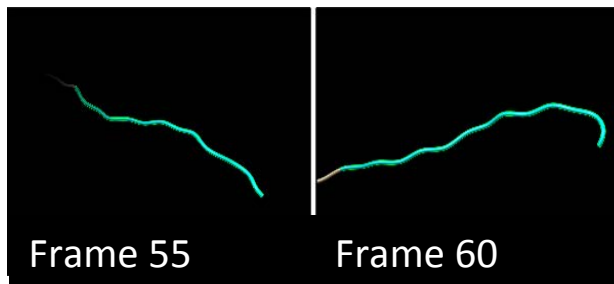


$\lambda_{55}$



## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

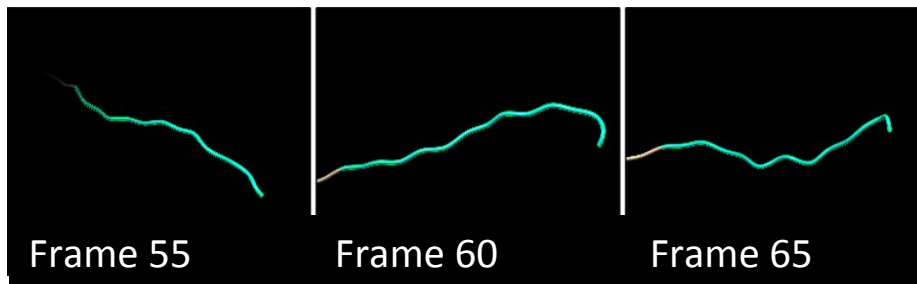
$\lambda_{60}$





## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

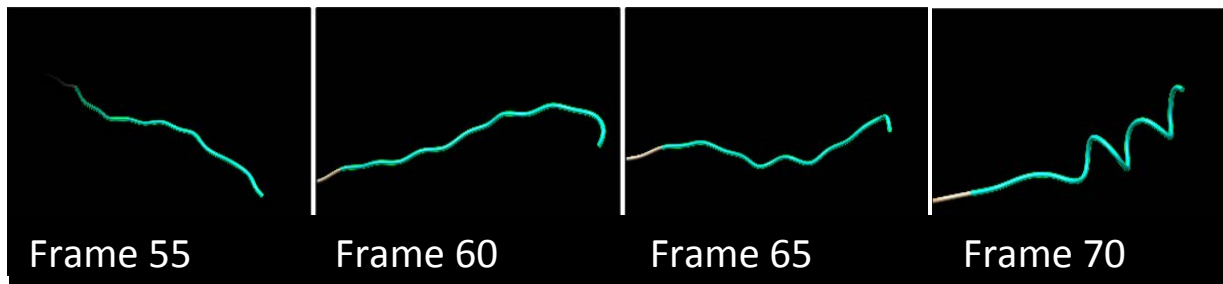
$\lambda_{60}$

$\lambda_{65}$



## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

$\lambda_{60}$

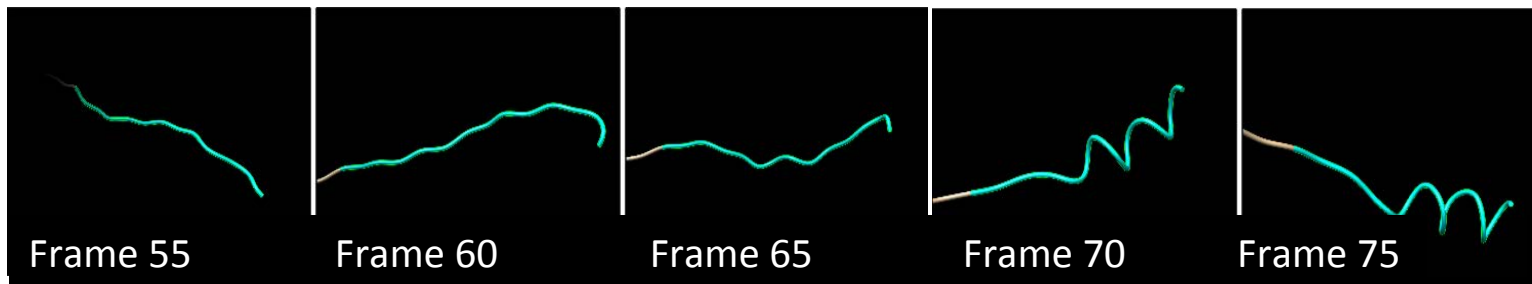
$\lambda_{65}$

$\lambda_{70}$



## Proxy for Rare Events

Frames of an MD job:



Frame 55

Frame 60

Frame 65

Frame 70

Frame 75

$\lambda_{55}$

$\lambda_{60}$

$\lambda_{65}$

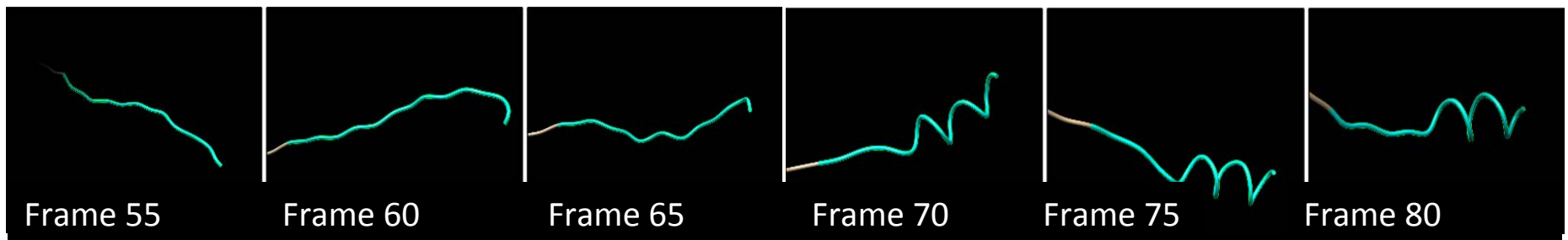
$\lambda_{70}$

$\lambda_{75}$



## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

$\lambda_{60}$

$\lambda_{65}$

$\lambda_{70}$

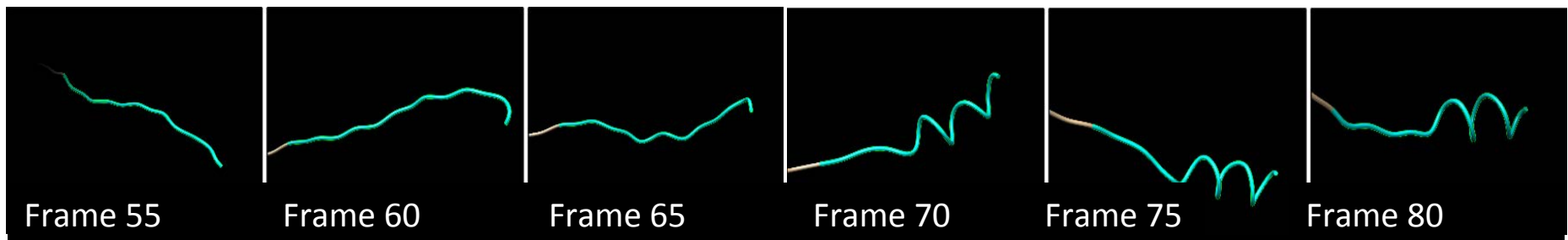
$\lambda_{75}$

$\lambda_{85}$



## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

$\lambda_{60}$

$\lambda_{65}$

$\lambda_{70}$

$\lambda_{75}$

$\lambda_{85}$

***Can the distance between two max eigenvalues serve as a proxy for distance between the two associated conformations?***



## Reasons to Love Symmetric Matrices

***Can the distance between two max eigenvalues serves as a proxy for distance between the two associated conformations?***

- Euclidean distance matrix D is symmetric
- Eigenvalues of symmetric, real matrices are stable
  - Small perturbations of D result in only small changes in the eigenvalues
  - Euclidean distance matrix is insensitive to rigid transformation
- Use only largest eigenvalue in distance matrix

$$\lambda_{max} = \lambda_1 < \lambda_2 < \lambda_3 < \lambda_4 < \lambda_5 = \lambda_{min}$$

$$\lambda_1 + \lambda_2 + \lambda_3 + \lambda_4 + \lambda_5 = 0$$

$$\lambda_1 \gg \lambda_2 \sim \lambda_3 \sim \lambda_4 \sim 0$$

$$\lambda_{max} = \lambda_1 \sim -\lambda_5 = -\lambda_{min}$$

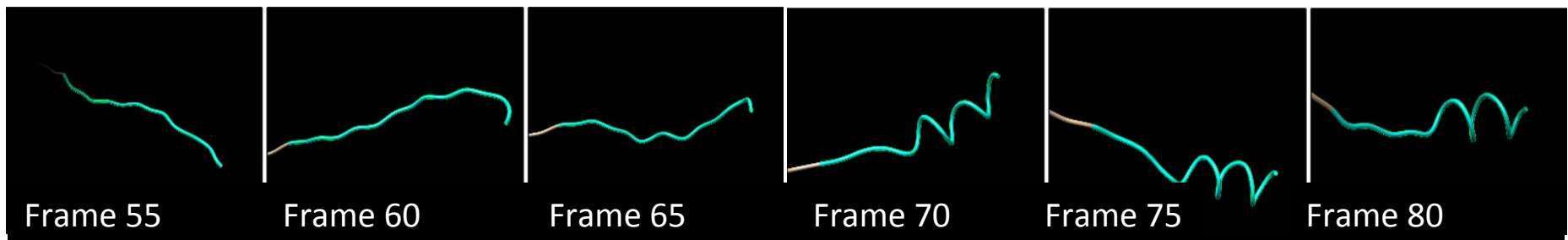
	α-carbon										
α-carbon	0	x	x	x	x	x	x	x	x	x	x
	x	0	x	x	x	x	x	x	x	x	x
	x	x	0	x	x	x	x	x	x	x	x
	x	x	x	0	x	x	x	x	x	x	x
	x	x	x	x	0	x	x	x	x	x	x
	x	x	x	x	x	0	x	x	x	x	x
	x	x	x	x	x	x	0	x	x	x	x
	x	x	x	x	x	x	x	0	x	x	x
	x	x	x	x	x	x	x	x	0	x	x
	x	x	x	x	x	x	x	x	x	0	x
	x	x	x	x	x	x	x	x	x	x	0

***“In-Situ Data Analysis and Indexing of Protein Trajectories,”*** Travis Johnston, Buyu Zhang, Adam Liwo, Silvia Crivelli, and Michela Taufer. JCC 2017.



## Proxy for Rare Events

Frames of an MD job:



$\lambda_{55}$

$\lambda_{60}$

$\lambda_{65}$

$\lambda_{70}$

$\lambda_{75}$

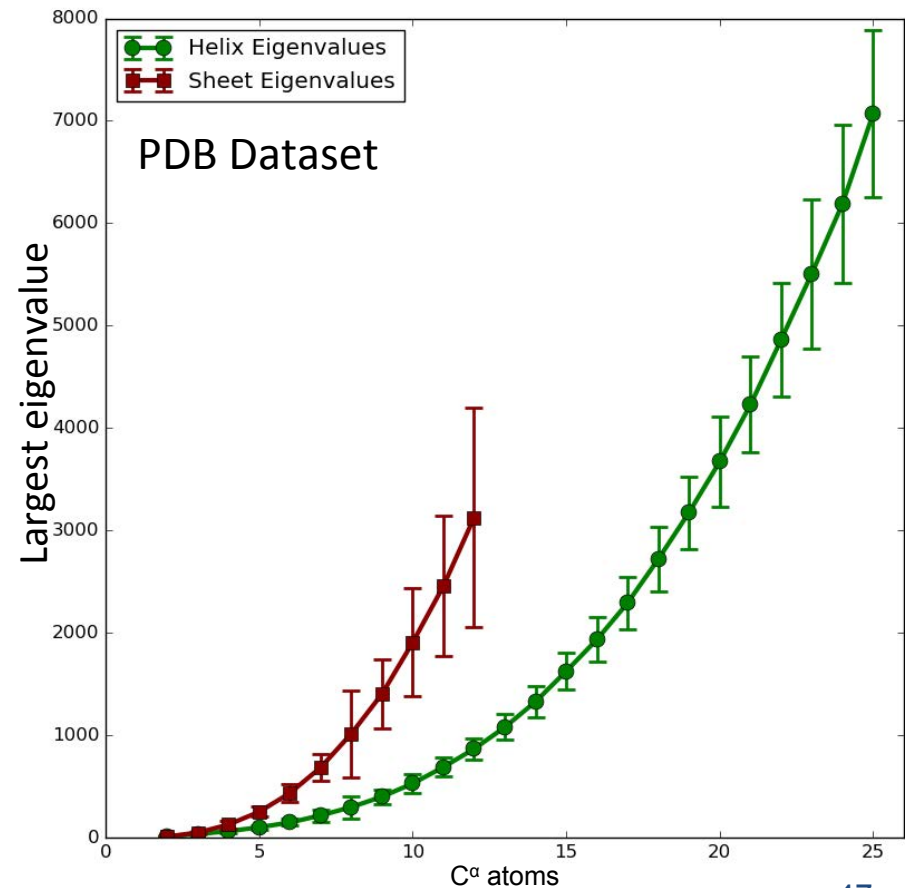
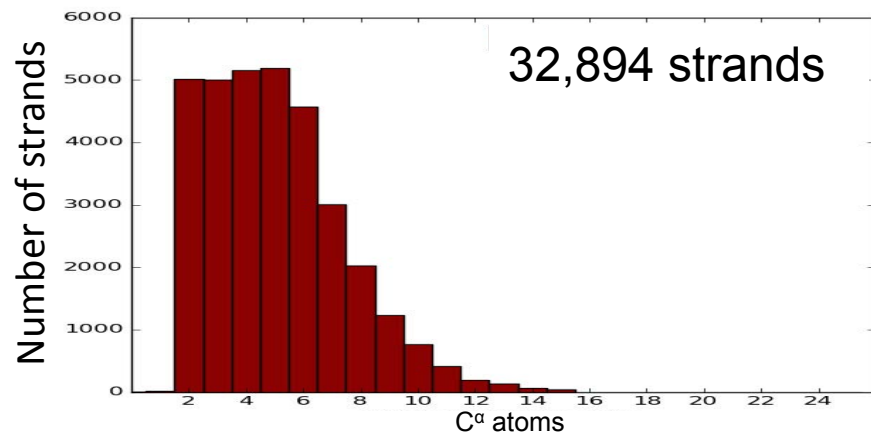
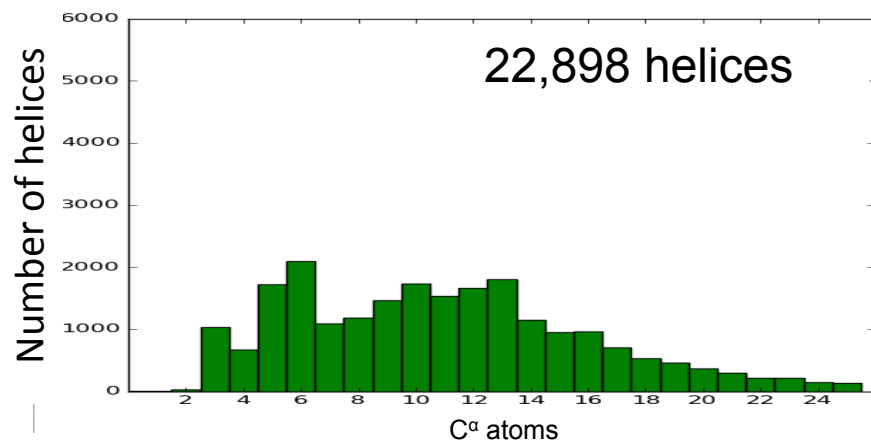
$\lambda_{85}$

***Yes, the distance between two max eigenvalues serves as a proxy for distance between the two associated conformations!***



# Mapping Largest Eigenvalues to Structures

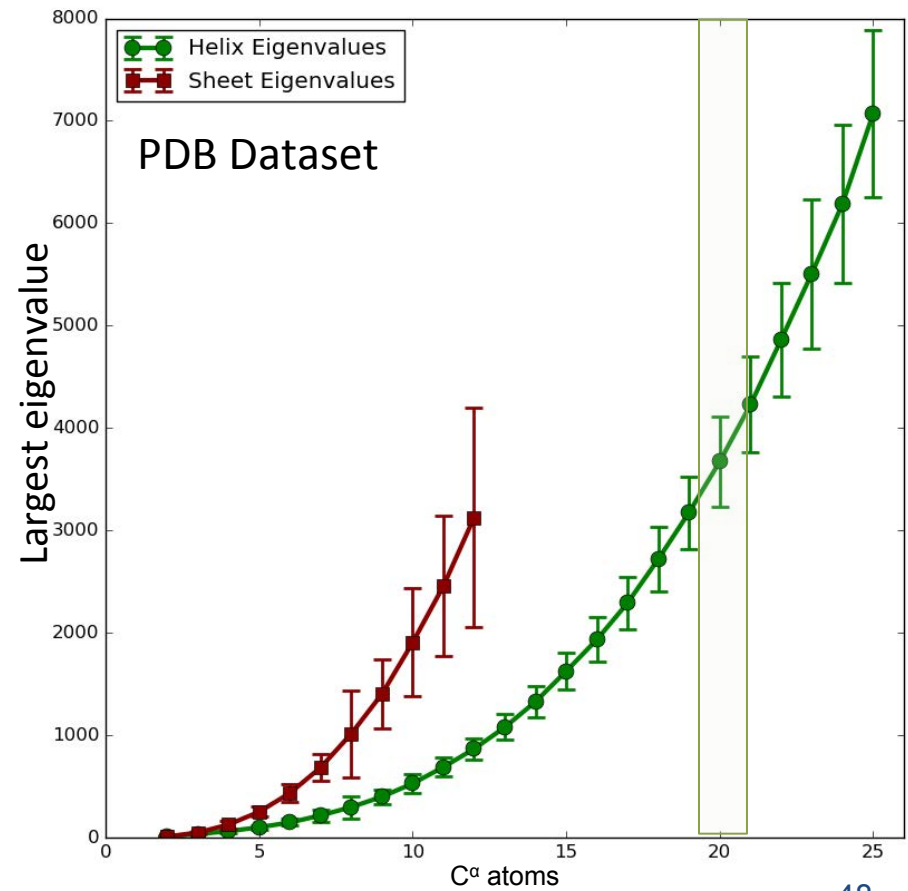
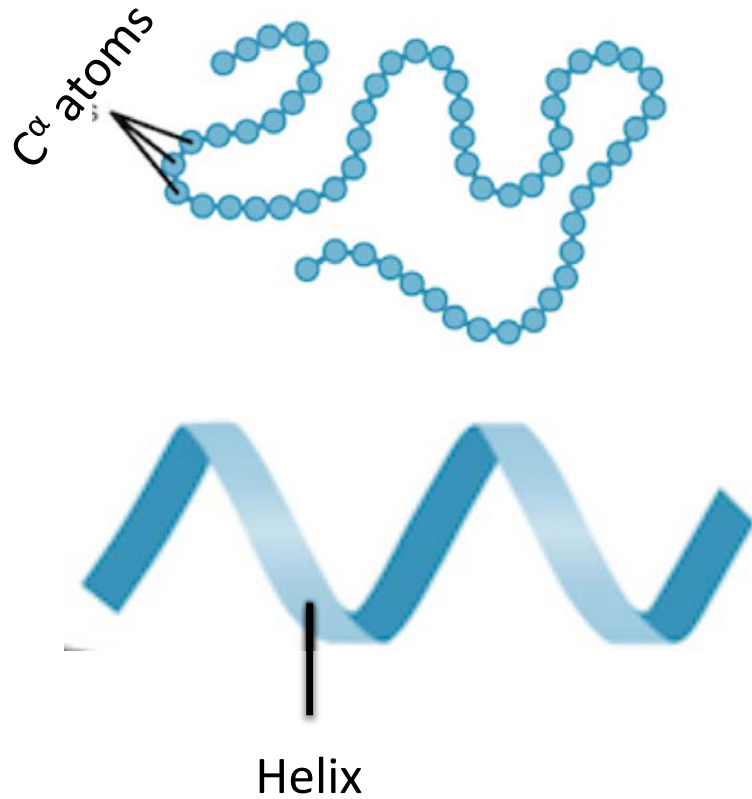
**PDB dataset:** 3,197 different proteins including 22,898 helices and 32,894 strands





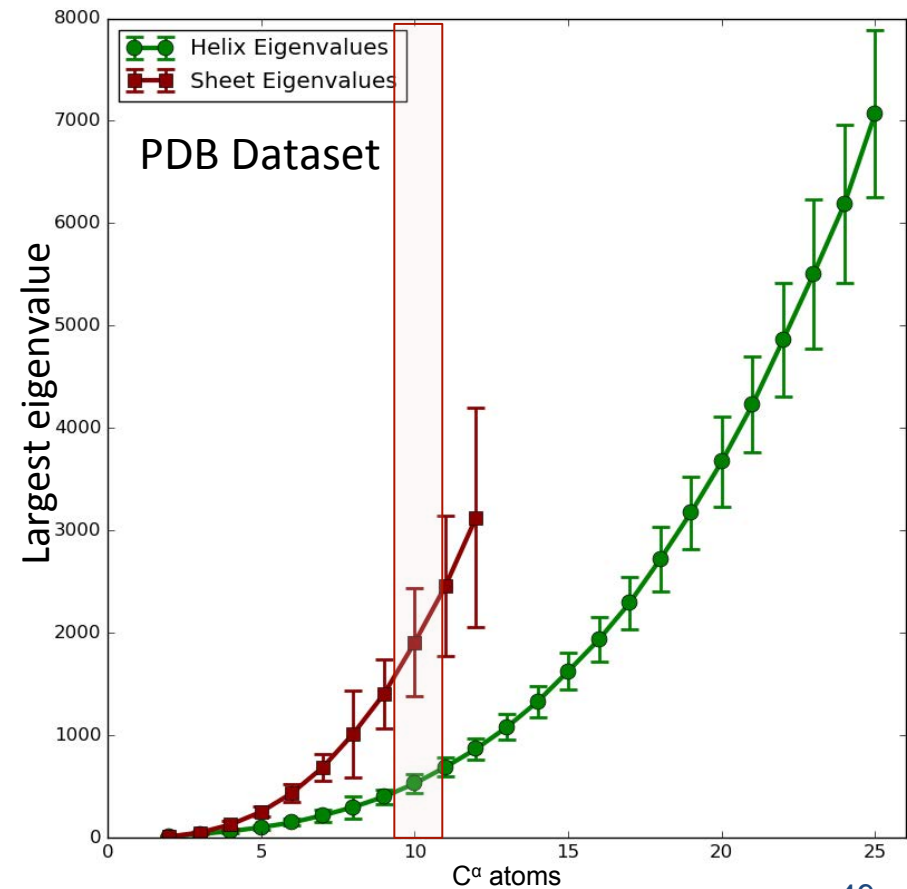
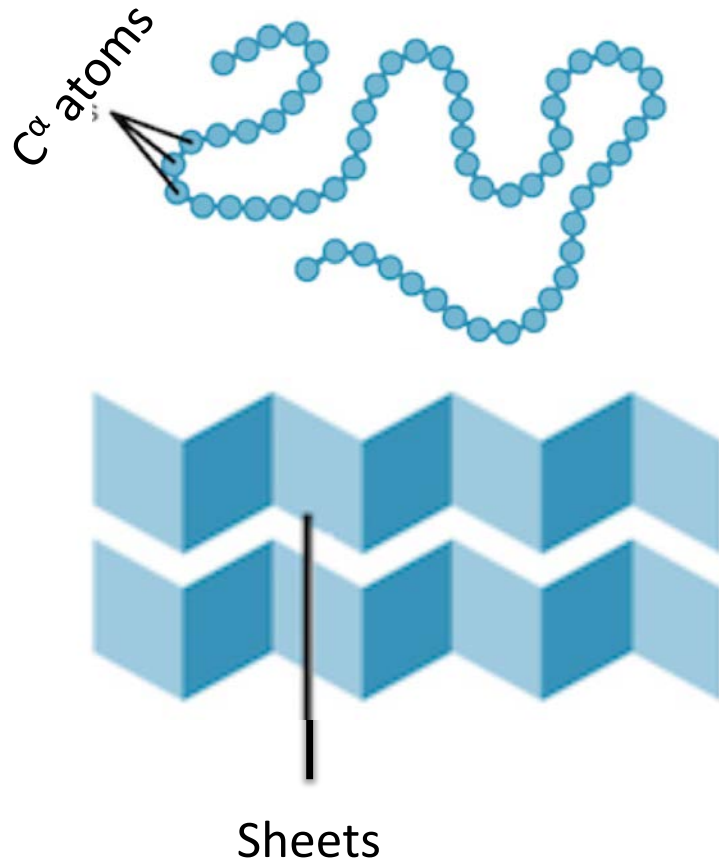


# Mapping Largest Eigenvalues to Structures





# Mapping Largest Eigenvalues to Structures

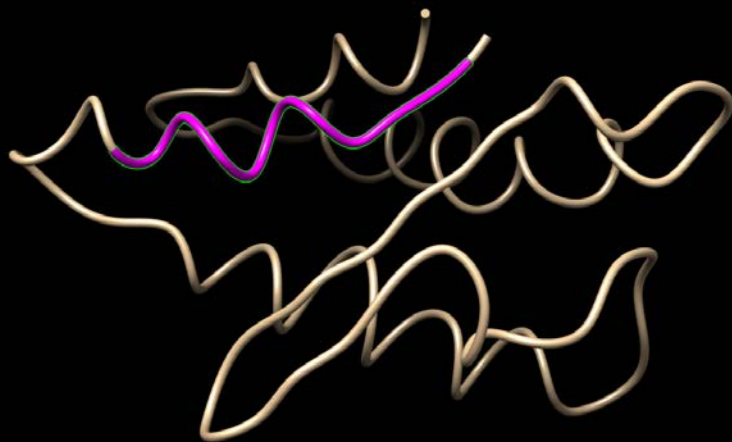




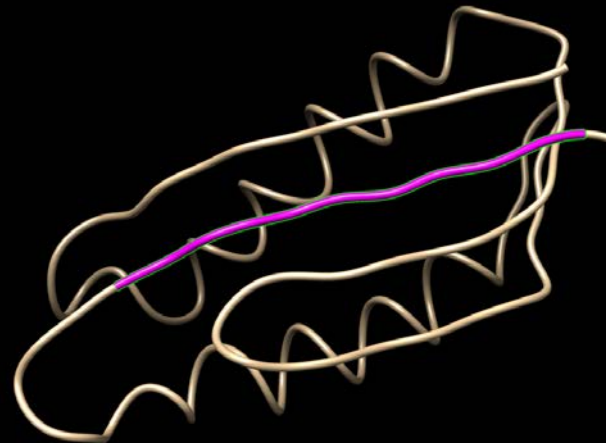
## Case Study I: 2MQ8 Protein

- Canonical simulation of 2MQ8 protein including both  $\alpha$  helices and  $\beta$  strands
  - After  $\sim 9$ M steps  $\alpha$  helices pack tighter and change into  $\beta$  strands

Frame 7686



Frame 8925

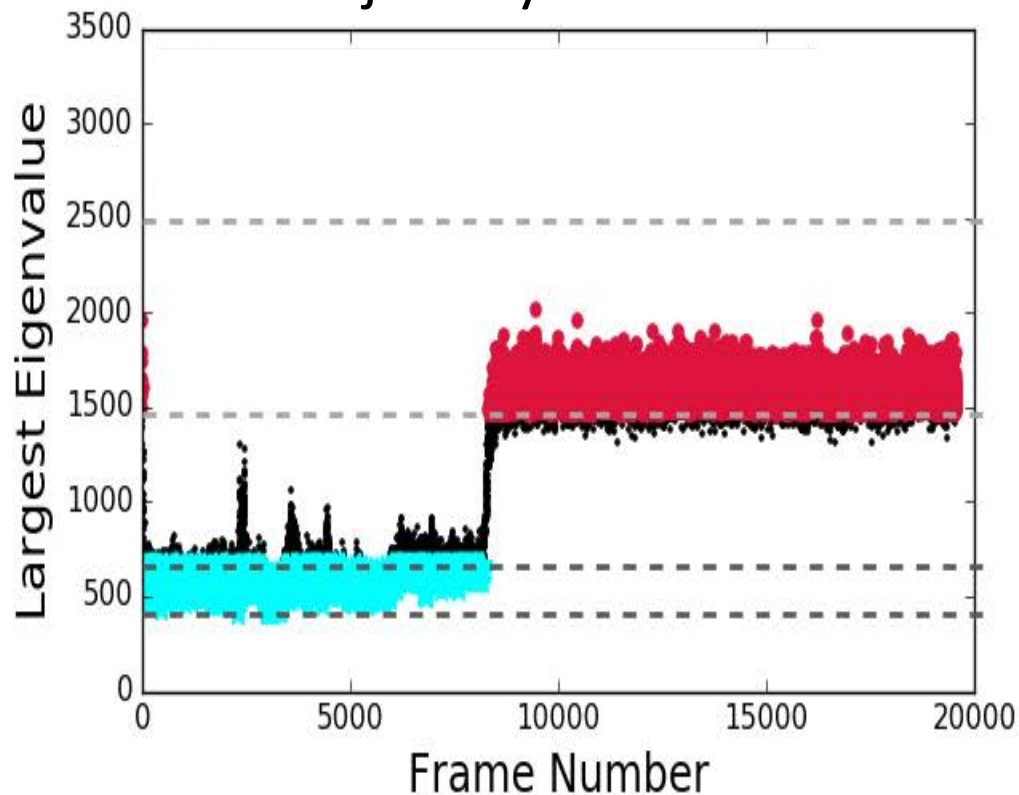


***Can the eigenvalue analysis capture the conformational change?***



## Case Study I: 2MQ8 Protein

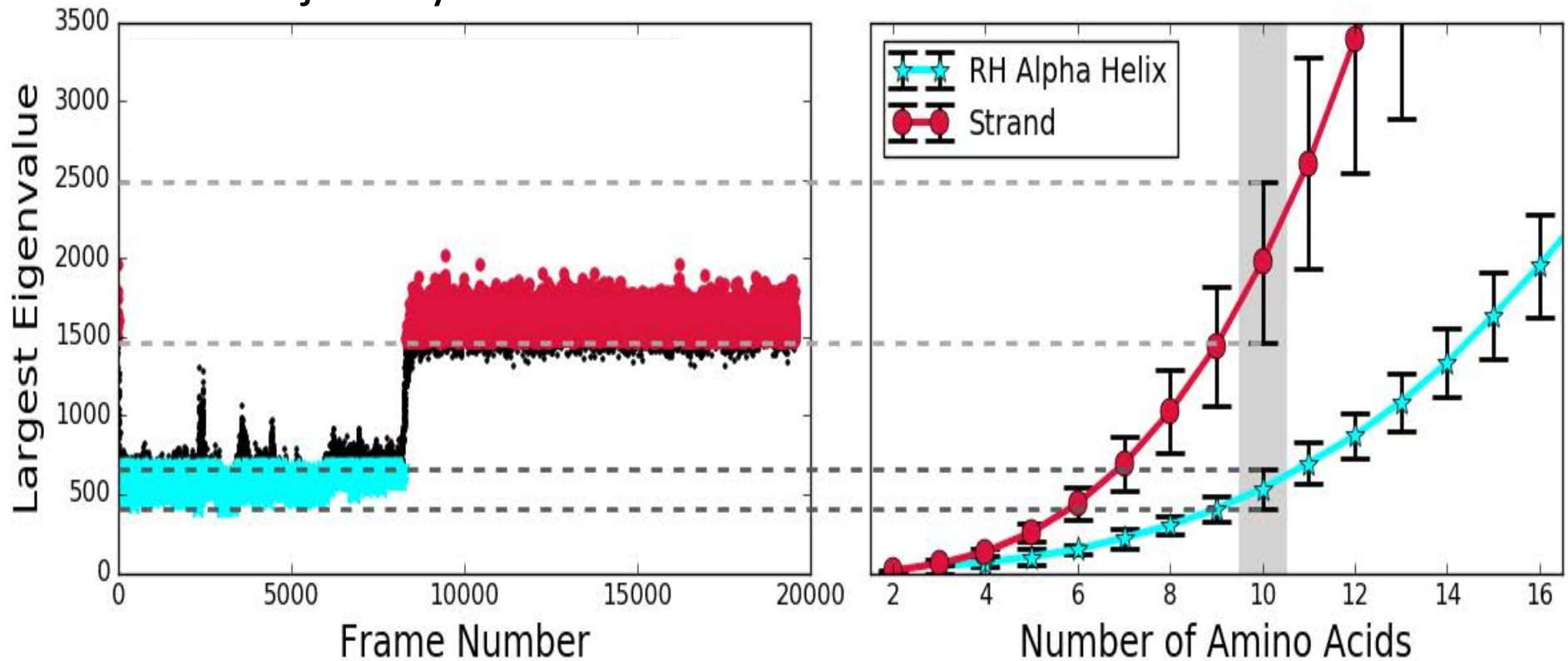
Compute largest eigenvalue of 3<sup>rd</sup> strand (10 amino acids) for each trajectory frame





## Case Study I: 2MQ8 Protein

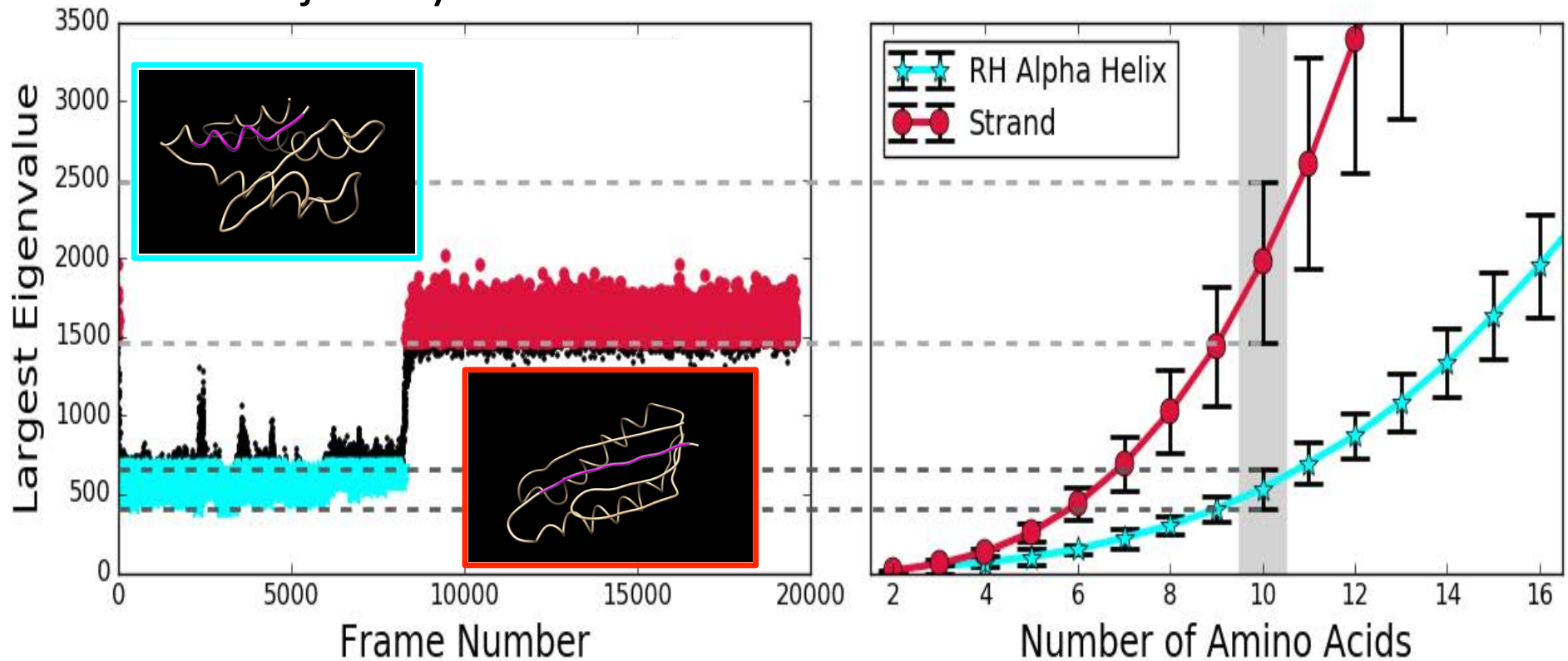
Compute largest eigenvalue of 3<sup>rd</sup> strand (**10 amino acids**) for each trajectory frame





## Case Study I: 2MQ8 Protein

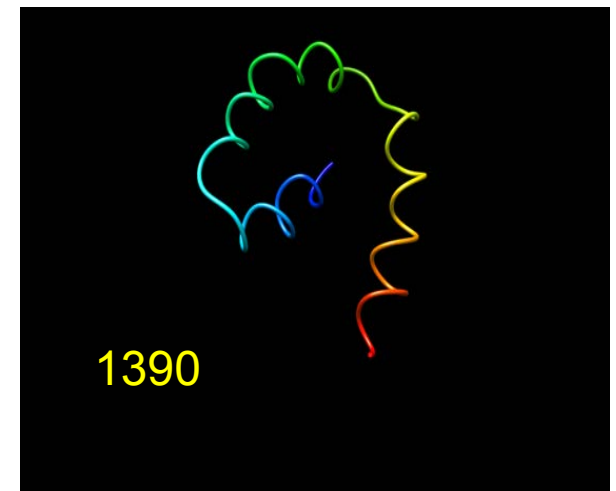
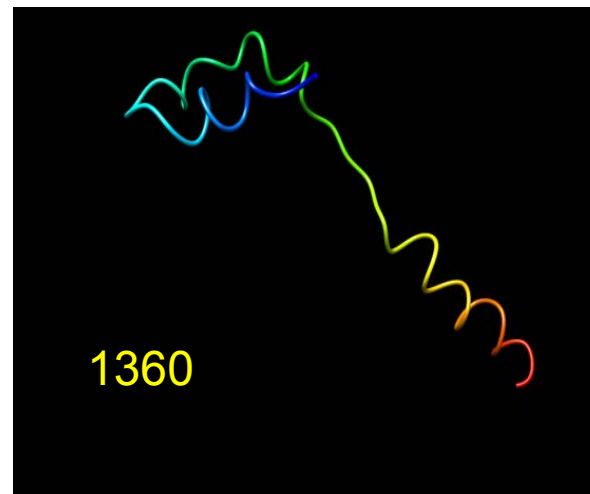
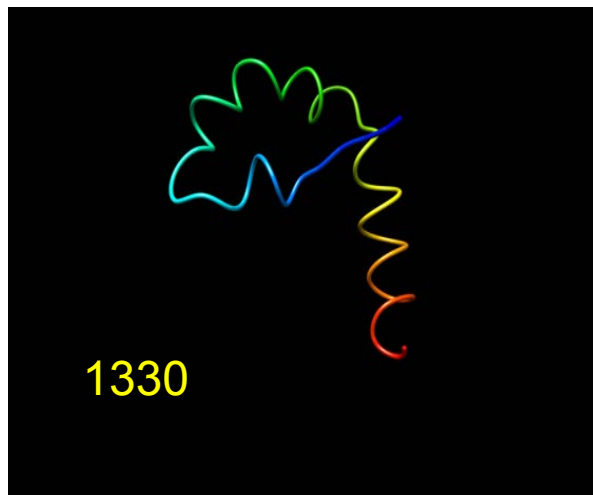
Compute largest eigenvalue of 3<sup>rd</sup> strand (10 amino acids) for each trajectory frame





## Case Study II: Capturing Movement of $\alpha$ -helices

Capture movement of structures with respect to each other

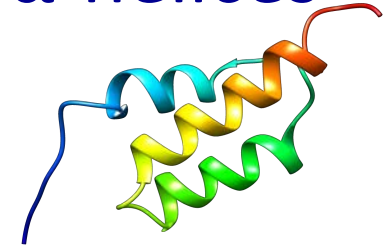


***Can the eigenvalue analysis capture the movement of helices ?***

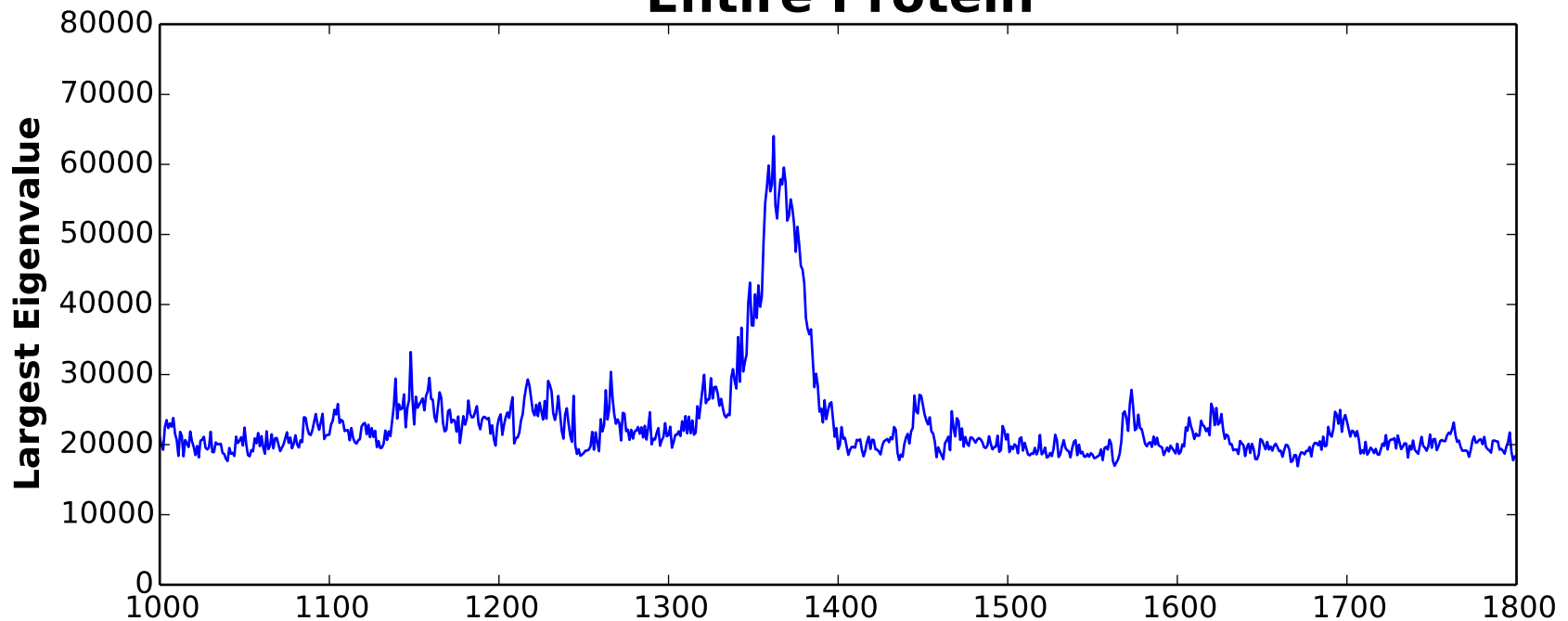


## Case Study II: Capturing Movement of $\alpha$ -helices

Monitor largest eigenvalue of entire protein



**Entire Protein**

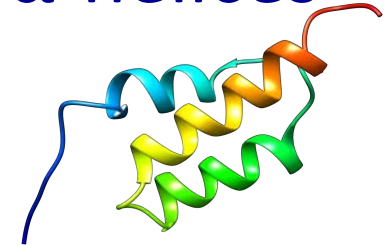




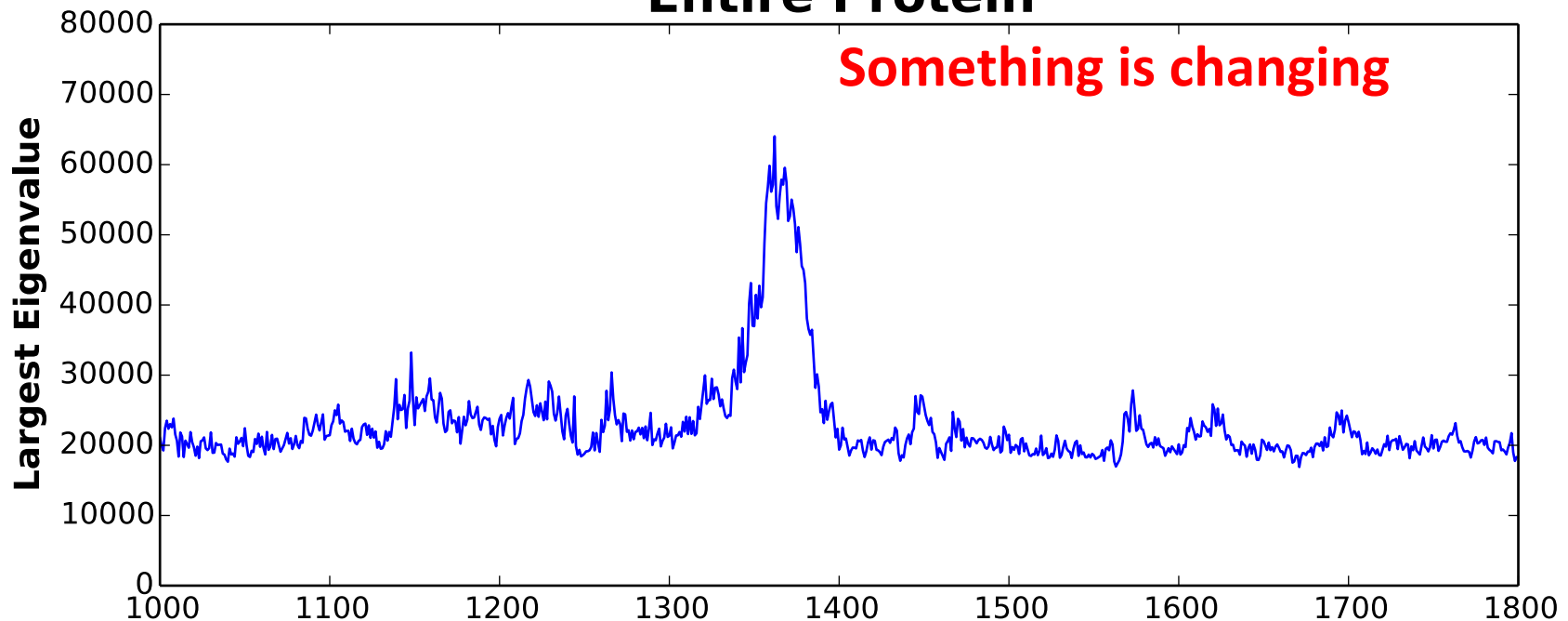


## Case Study II: Capturing Movement of $\alpha$ -helices

Monitor largest eigenvalue of entire protein



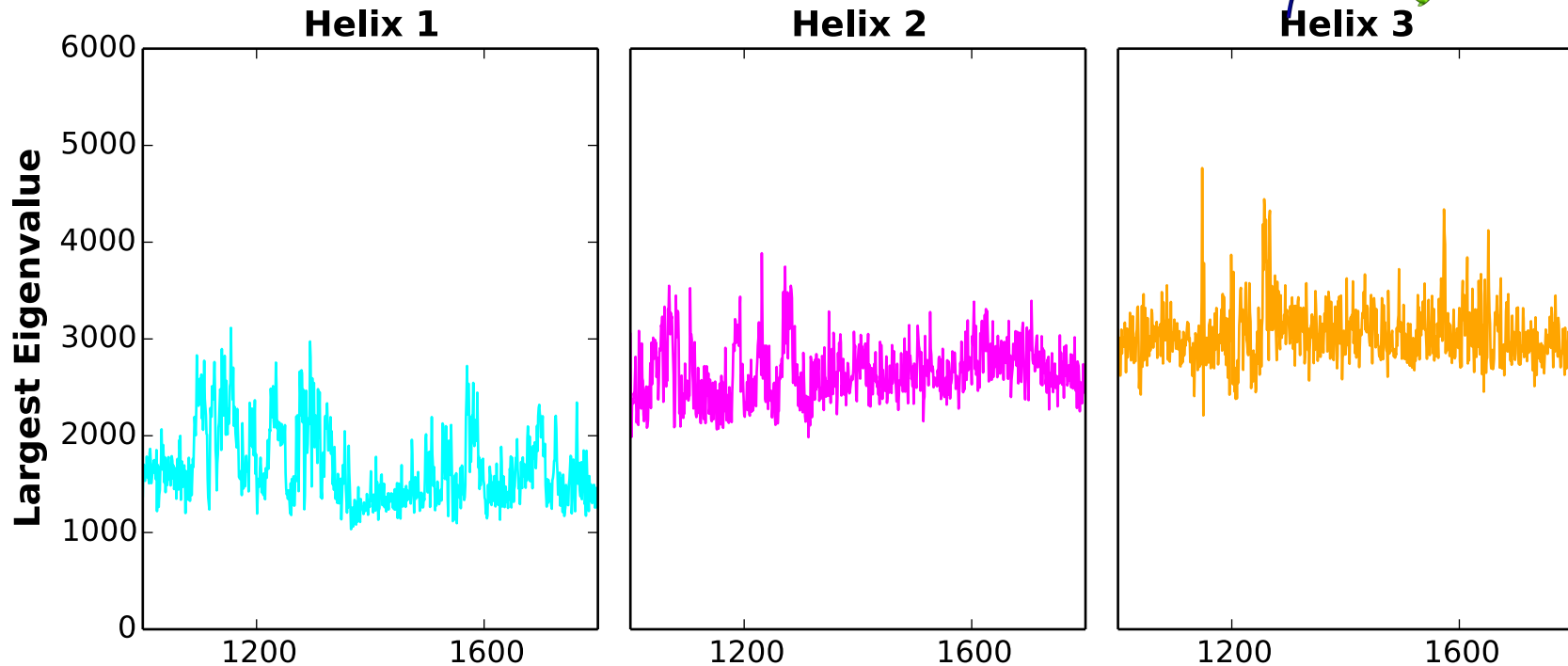
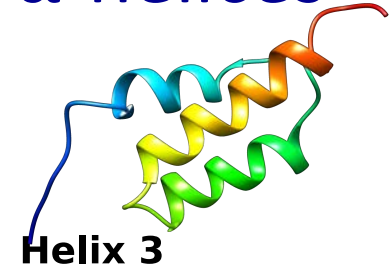
**Entire Protein**





## Case Study II: Capturing Movement of $\alpha$ -helices

Monitor largest eigenvalue of single helices

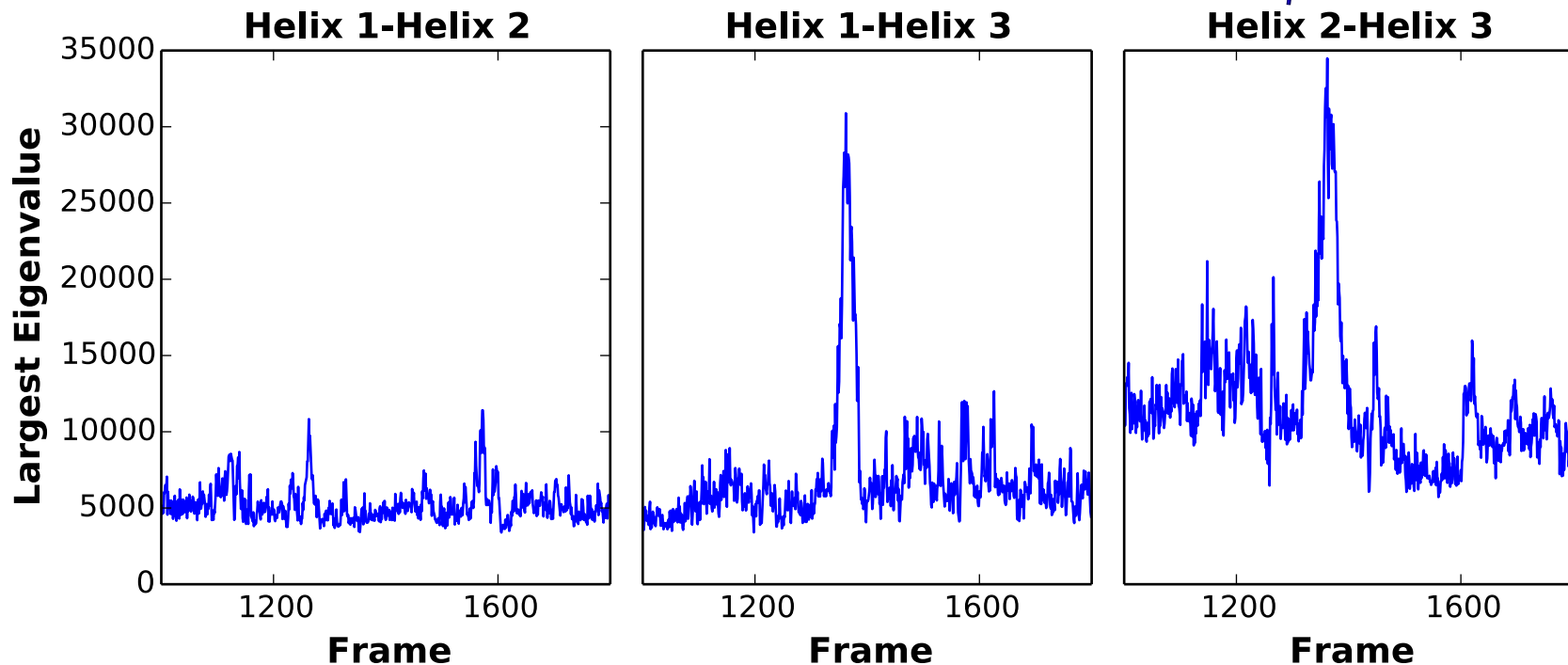
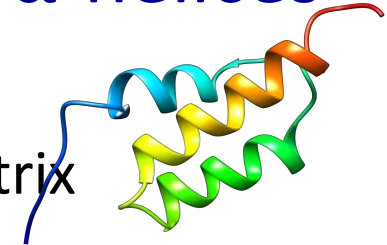


Individual  $\alpha$ -helices (Helix 1, Helix 2, and Helix 3) appear stable



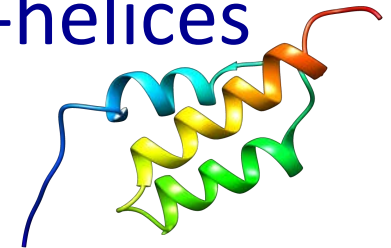
## Case Study II: Capturing Movement of $\alpha$ -helices

Monitor largest eigenvalue of bipartite distance matrix

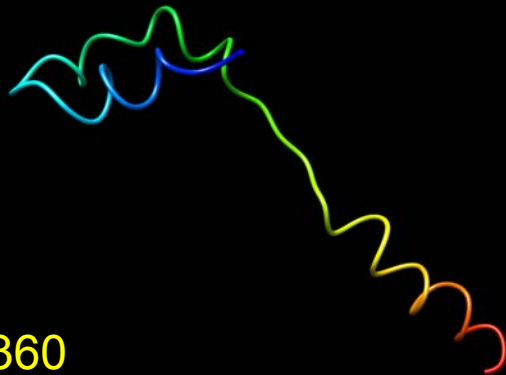


First and second  $\alpha$ -helices appear stable; third helix moves

## Case Study II: Capturing Movement of $\alpha$ -helices



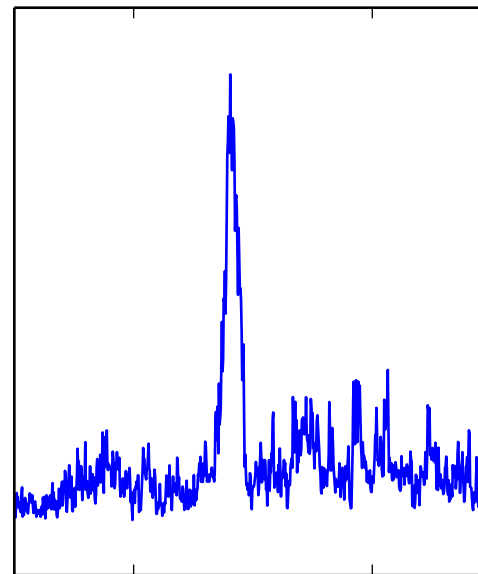
1330



1360

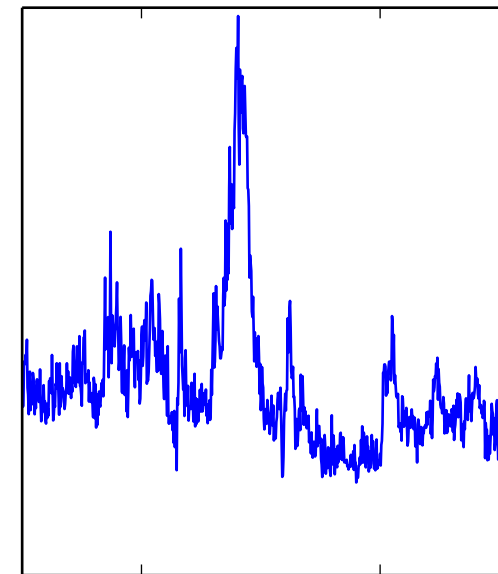


1390

**Helix 1-Helix 3**

1200 1600

Frame

**Helix 2-Helix 3**

1200 1600

Frame

Large relative change between  
two pairs of  $\alpha$ -helices



*“Storage technologies are advancing [...] and it is really not clear at all [to me] that especially distributed storage platforms would not be able to handle [...] petabyte data sets”*

*Anonymous Feedback*

**Yes, new technologies will be able to handle data at the extreme scale but *only* if we integrate new software paradigms.  
I/O-aware schedulers are *a must!*  
In-situ and in-transit analysis are *here to stay!***