# Scientific Applications and Heterogeneous Architectures –
# Data Analytics and the Intersection of HPC and Edge Computing

*Michela Taufer*

THE UNIVERSITY OF
**TENNESSEE**
KNOXVILLE

**BIG ORANGE. BIG IDEAS.**®
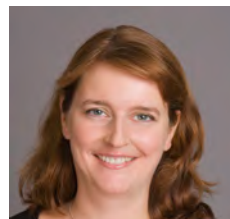
# Acknowledgements



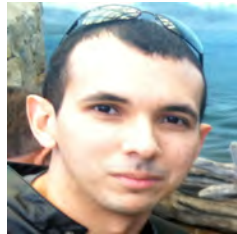T. Estrada

H. Weinstein
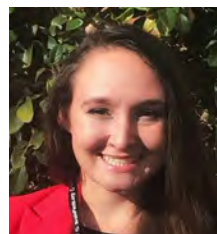
M. Cuendet

E. Deelman
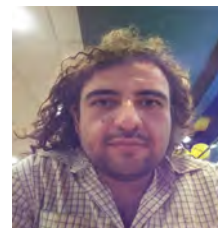
R. Vargas

R. da Silva

T. Johnston

T. Do

B. Mulligan

D. Rorabaugh

S. Thomas

H. Carrillo
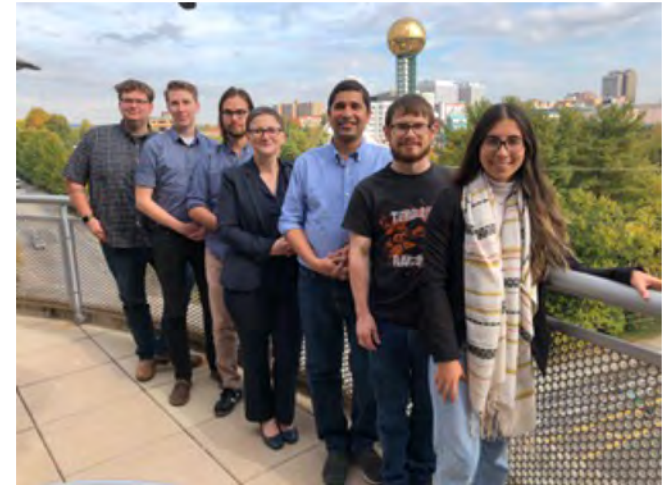
A. Razavi

R. LLamas

M. Guevara

# Trends in Next-Generation Systems: IO Gap and Ensembles

### Widening IO Gap



Source: Lucy Nowell (DOE)

### Rising Importance of Ensembles



Source: https://wci.llnl.gov/simulation/computer-codes/uncertainty-quantification

3

BIG ORANGE
BIG IDEAS

# Trends in Workflows: Compute + Analytics + Data



Pegasus LIGO PyCBC Workflow

Laser Interferometer Gravitational-Wave Observatory (LIGO)

Source: Ewa Deelman, ISI - USC

# Extending HPC to Connect to the "Edge"

| Simulations | Data Analytics | Real System with Sensors |
|:---:|:---:|:---:|

Image Source: https://iiot-world.com/digital-disruption/the-right-representation-of-digital-twins-for-data-analytics/

BIG**ORANGE**
BIG**IDEAS**

# Two Use Cases

- Extending HPC to integrate data analytics
  - Next generation MD workflows
  - Molecular structures
  - ***Data transformation – i.e.,*** *capturing information*
  - ***Dataflow modeling –*** *i.e., lost information*
- Extending HPC to connect to the "Edge"
  - Next generation precision farming
  - Soil moisture data
  - ***Data prediction*** *– i.e., from coarse- to fine-grained information*
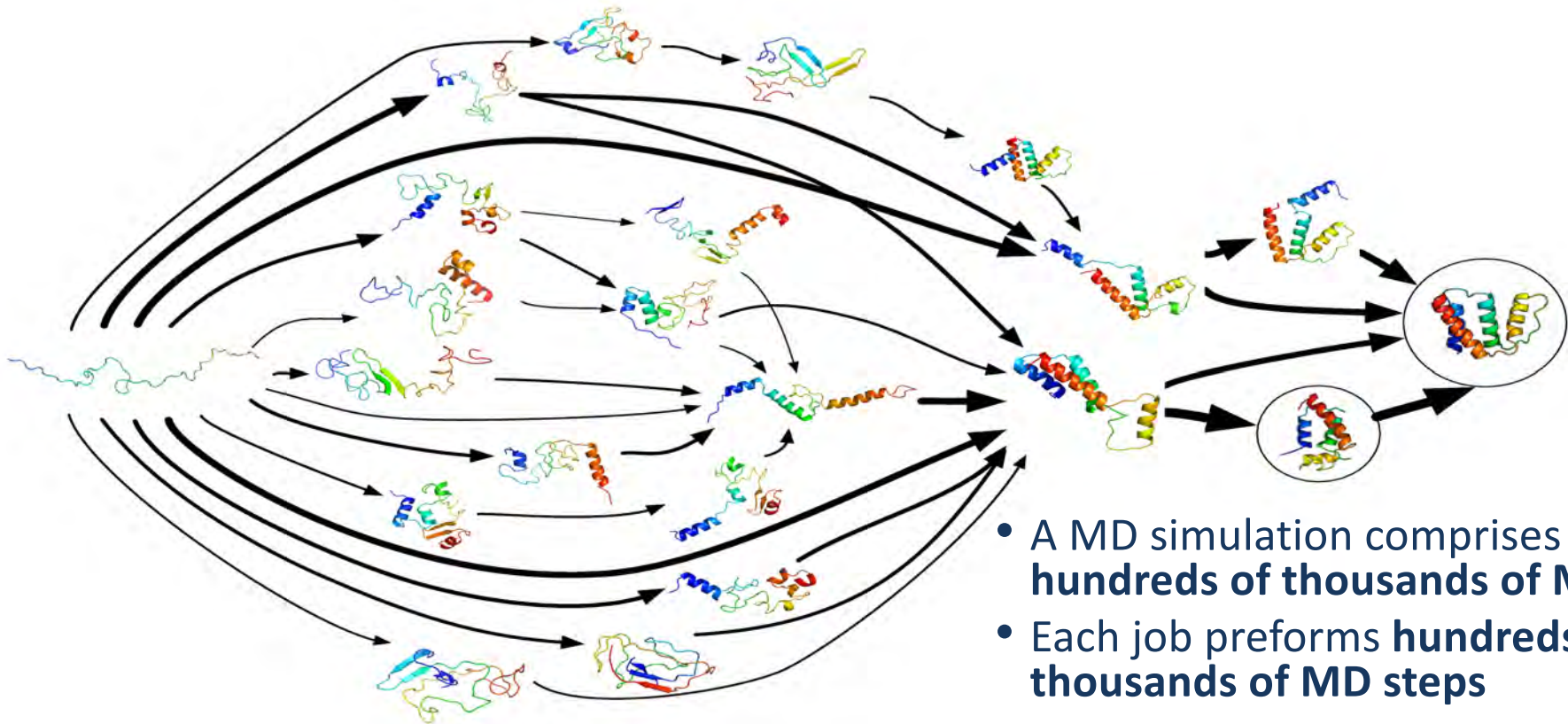
BIG**ORANGE**BIG**IDEAS**

# Two Use Cases

- Extending HPC to integrate data analytics
  - Next generation MD workflows
  - Molecular structures
  - ***Data transformation – i.e.,*** *capturing information*
  - ***Dataflow modeling –*** *i.e., lost information*
- Extending HPC to connect to the "Edge"
  - Next generation precision farming
  - Soil moisture data
  - ***Data prediction*** *– i.e., from coarse- to fine-grained information*
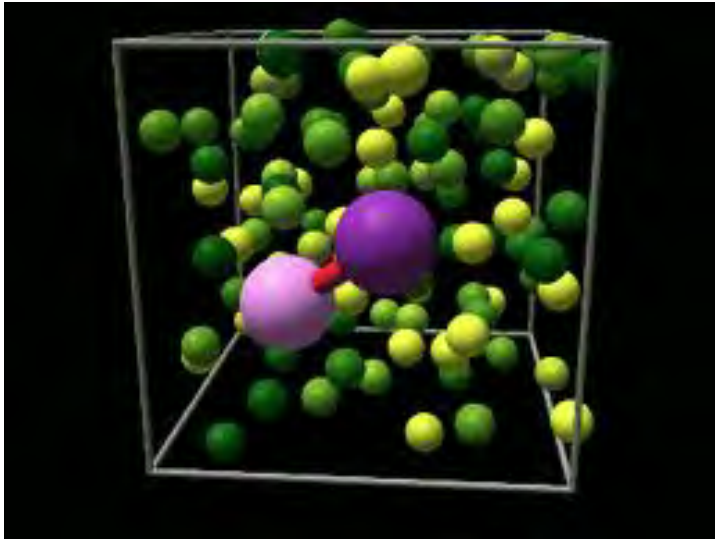
BIG**ORANGE**
BIG**IDEAS**

# Classical Molecular Dynamics Simulations



- A MD simulation comprises of **hundreds of thousands of MD job**
- Each job preforms **hundreds of thousands of MD steps**

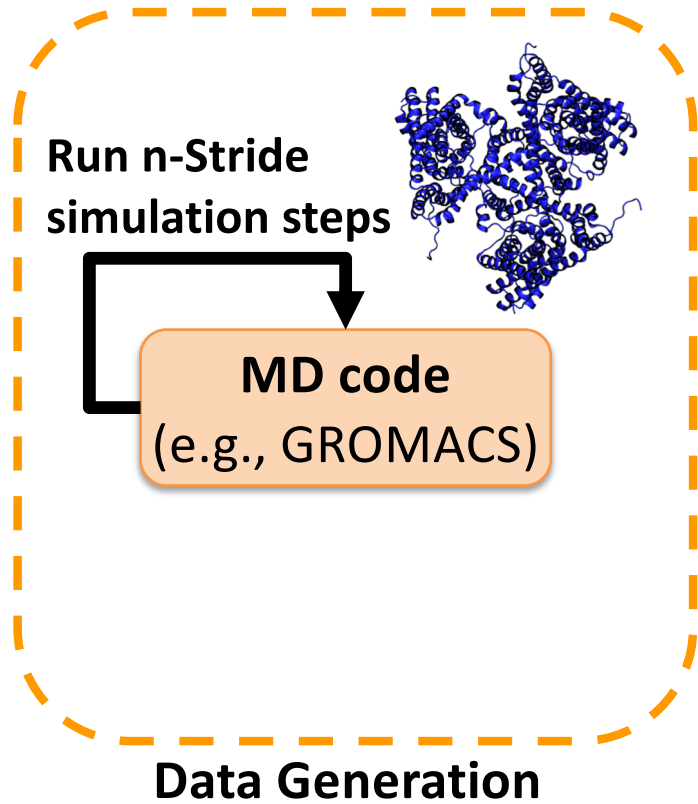# Classical Molecular Dynamics Simulations



→ Forces on single atoms
  → Acceleration
    → Velocity
      → Position

- MD step computes **forces** on single atoms (e.g., bond, dihedrals, nonbond)
- Forces are added to compute **acceleration**
- Acceleration is used to update **velocities**
- Velocities are used to update the **atom positions**
- Every *n* steps (Stride)
  → ***Store 3D snapshot or frame***

# Building a Closed-loop Workflow

**Run n-Stride simulation steps**

**MD code**
(e.g., GROMACS)

**Data Generation**

# Building a Closed-loop Workflow



**Run n-Stride simulation steps**

MD code (e.g., GROMACS)

Plumed

**Dataflow**

Ingestor

**Data Generation**

Parallel File System

Burst Buffer

**In-memory Staging Area** *DataSpaces*

**Data Storage**

BIG**ORANGE** BIG**IDEAS**

# Building a Closed-loop Workflow



**Run n-Stride simulation steps**

**MD code** (e.g., GROMACS)

Plumed

**Dataflow**

Ingestor

**Data Generation**

Parallel File System

Burst Buffer

**In-memory Staging Area** *DataSpaces*

**Data Storage**

ML-inferred algorithms Collective variables

**Dataflow**

Retriever

**Analytics representations + algorithms**

**Data Analytics**

BIG**ORANGE** BIG**IDEAS**

# Building a Closed-loop Workflow



**Data Feedback**

Run n-Stride simulation steps

**MD code** (e.g., GROMACS)

Plumed

**Dataflow**

Ingestor

**Data Generation**

Parallel File System

Burst Buffer

**In-memory Staging Area** *DataSpaces*

**Data Storage**

ML-inferred algorithms Collective variables

**Dataflow**

Retriever

**Analytics representations + algorithms**

**Data Analytics**
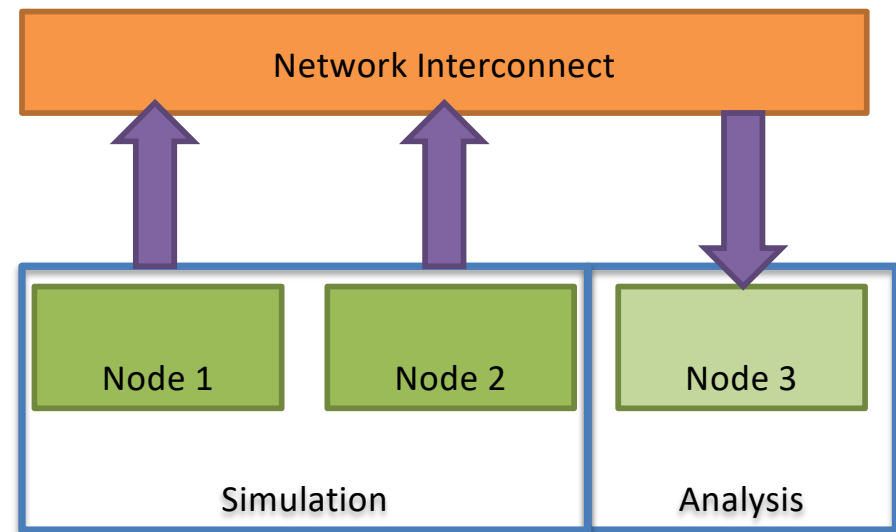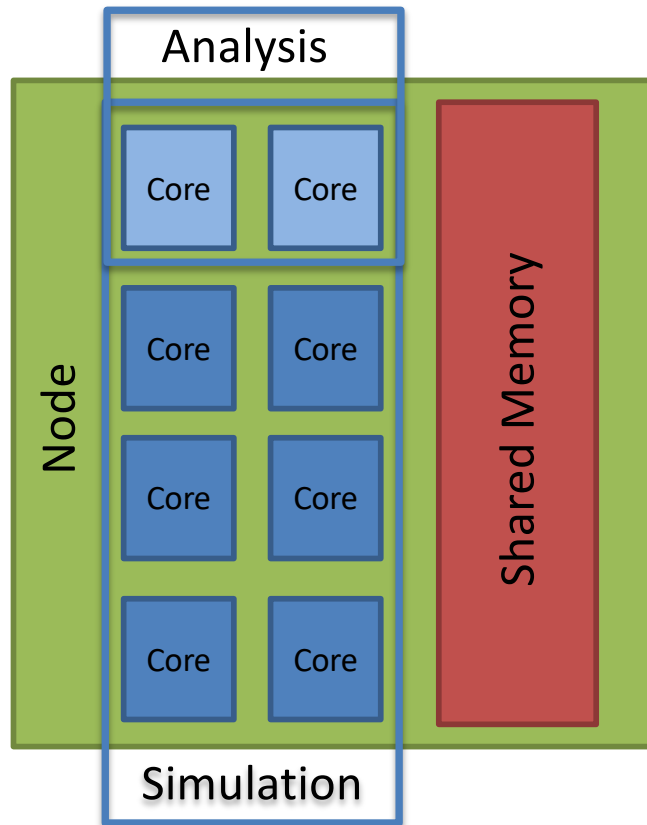
# Extending HPC to Integrate Data Analytics

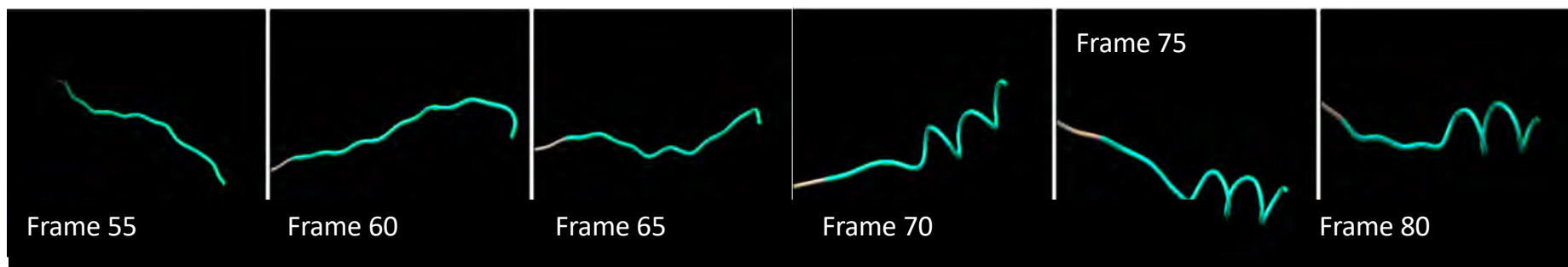# *Augmenting HPC with In Situ* and *In Transit* Analytics



Example of tools:
- DataSpaces (Rutgers U.)
- DataStager (GeorgiaTech)

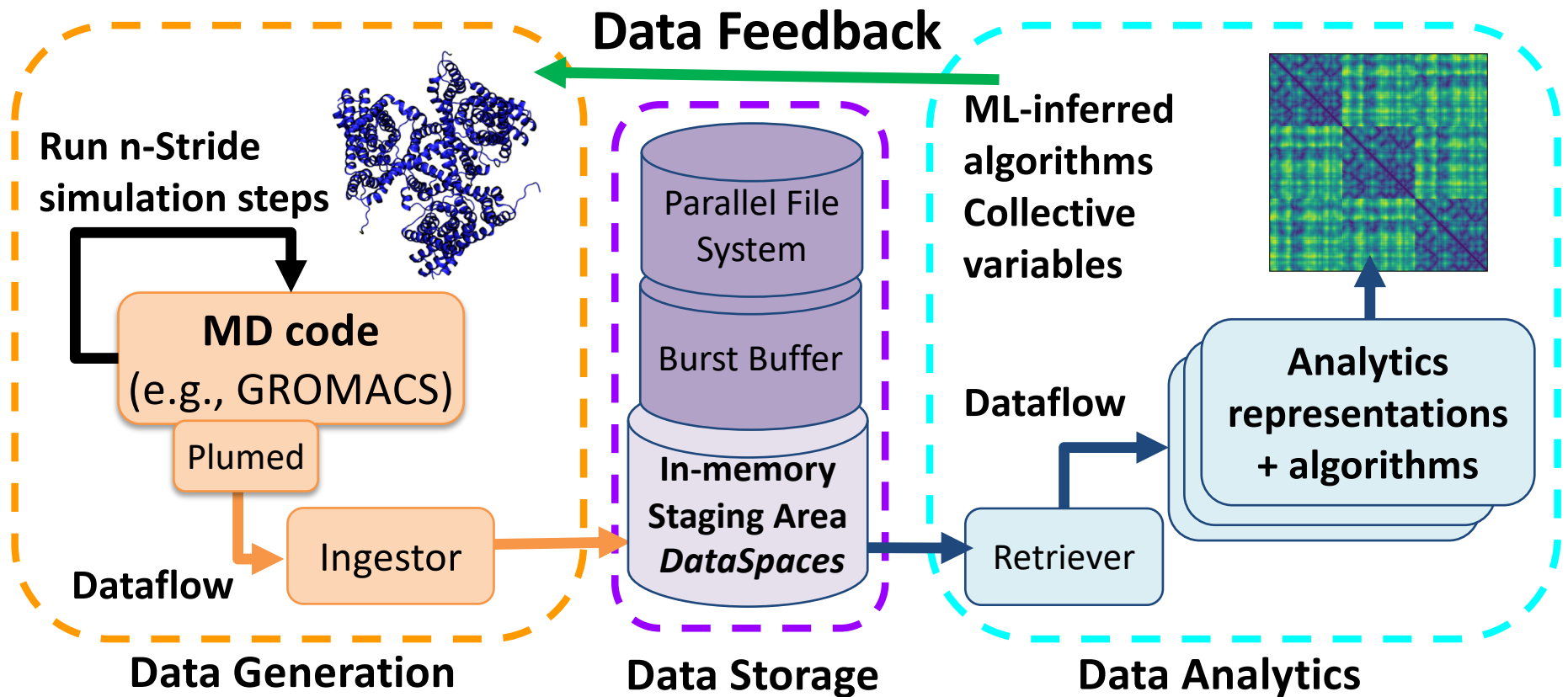# *In Situ* and *In Transit* Analytics for MD Simulations

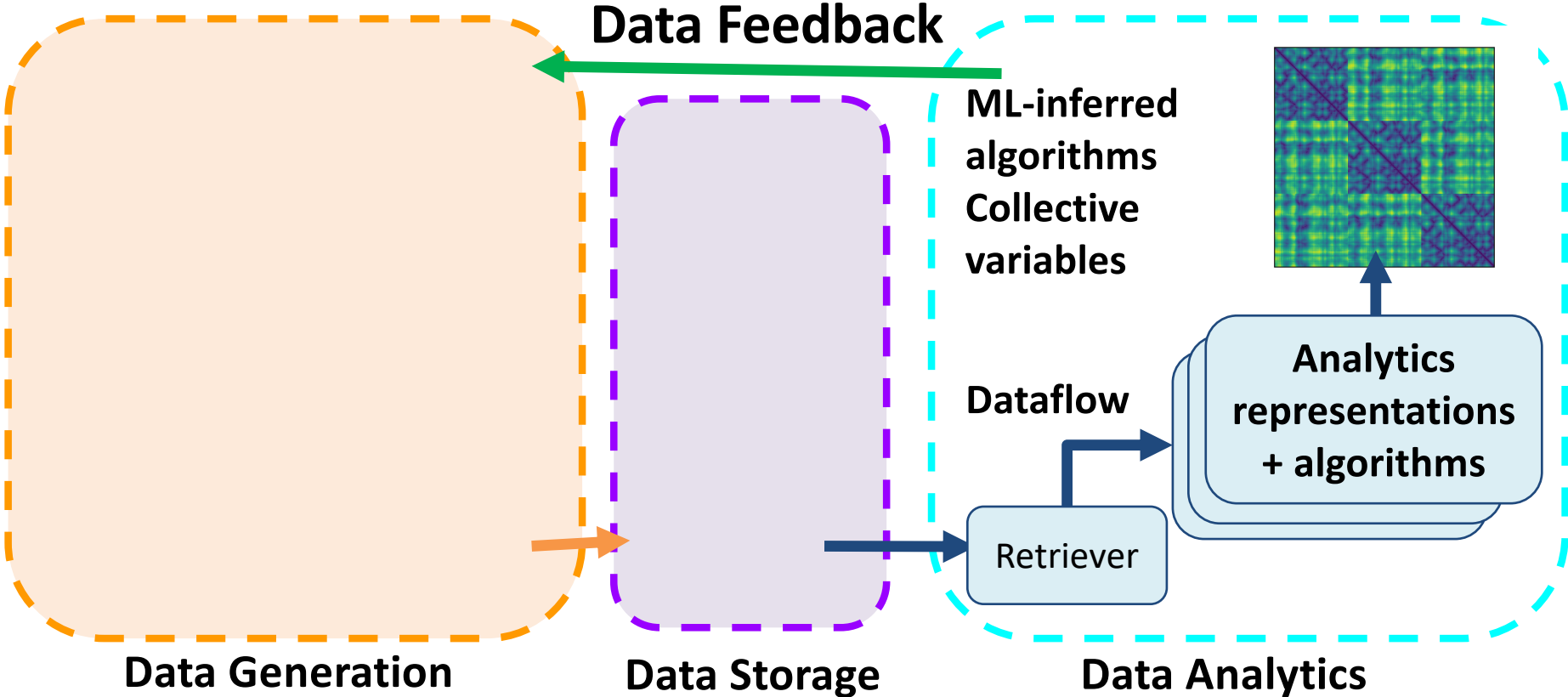Frames (or snapshots) of an MD trajectory with a stride of 5 steps:



- 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

# Building a Closed-loop Workflow



**Data Feedback**

Run n-Stride simulation steps

MD code (e.g., GROMACS)

Plumed

Dataflow

Ingestor

**Data Generation**

Parallel File System

Burst Buffer

In-memory Staging Area *DataSpaces*

**Data Storage**

ML-inferred algorithms Collective variables

Dataflow

Retriever

Analytics representations + algorithms

**Data Analytics**

# Building a Closed-loop Workflow



**Data Feedback**

ML-inferred algorithms
Collective variables

Dataflow

Retriever

Analytics representations + algorithms

Data Generation

Data Storage

Data Analytics

BIG**ORANGE** BIG**IDEAS**
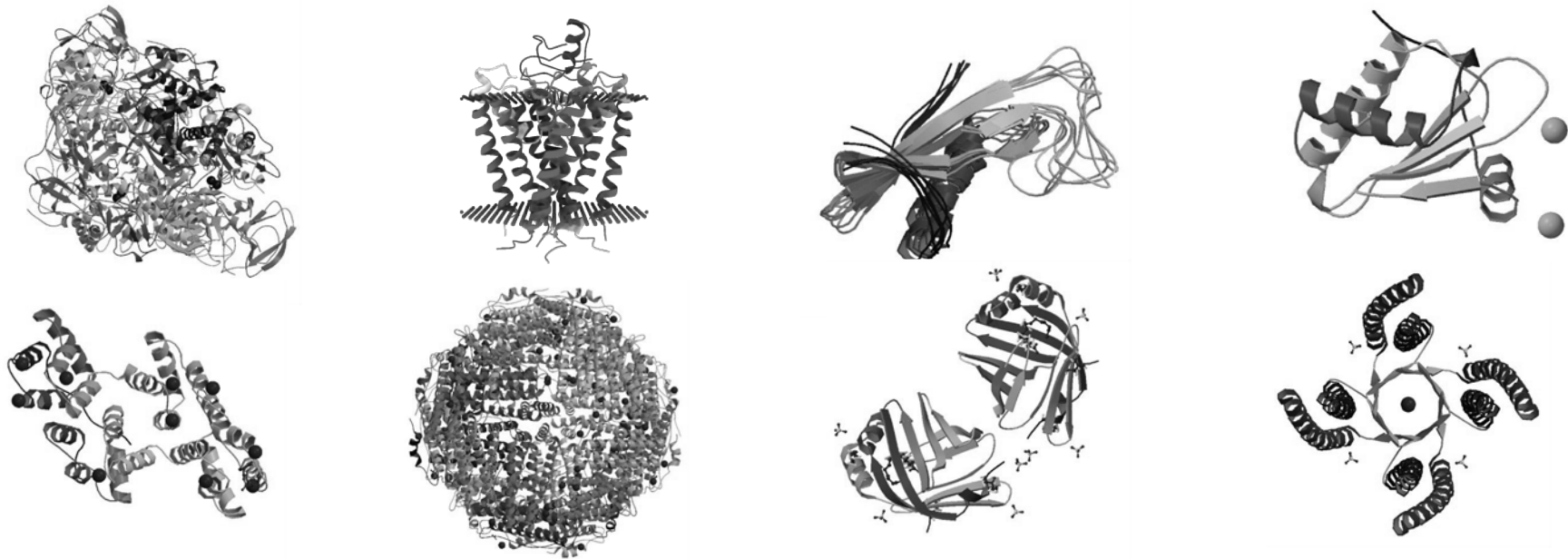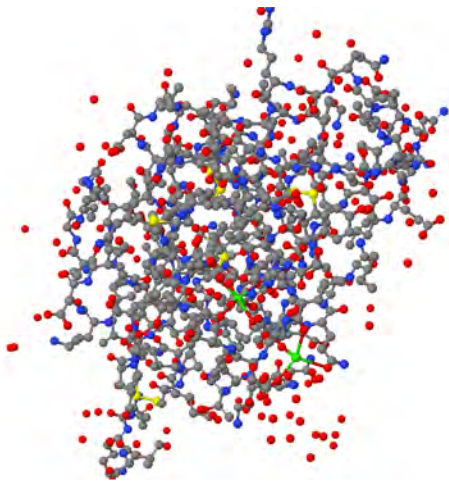
# Proteins with Similar Functions

Key principle: proteins with similar structures have similar functions

- Measure millions of protein variants expressed from yeast or bacteria
- Structure proteins to produce desired properties (protein engineering)
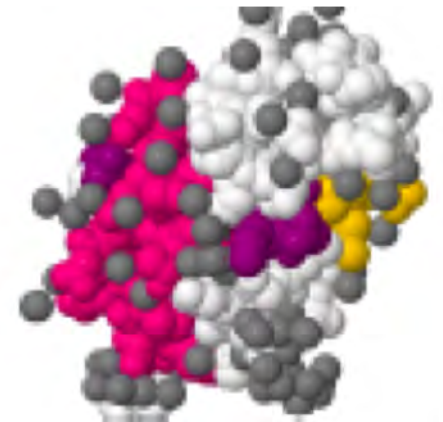
BIG**ORANGE**
BIG**IDEAS**

# Protein Representations



3D Cartesian representation

Multi-fold representation

Surface representation

BIG**ORANGE** BIG**IDEAS**

# From Multi-fold Representation to Image Encoding



Backbone dihedral angles

Original 3D protein

Atoms' Cartesian coordinates*

1.- Ramachandran Plot

2.- Distance Matrix

Channel:1

Channel:2

Channel:3

3.- Channel Encoding

4.- Final Encoding

Every channel encodes information associated with particular secondary structures and their spatial relationship

T. Estrada, J. Benson, H. Carrillo-Cabada, A. Razavi, M. Cuendet, H. Weinstein, E. Deelman, and M. Taufer. *Graphic Encoding of Proteins for Efficient High-Throughput Analysis*. ICPP 2018.

# From Multi-fold Representation to Image Encoding

T. Estrada, J. Benson, H. Carrillo-Cabada, A. Razavi, M. Cuendet, H. Weinstein, E. Deelman, and M. Taufer. *Graphic Encoding of Proteins for Efficient High-Throughput Analysis*. ICPP 2018.

# From Multi-fold Representation to Image Encoding

T. Estrada, J. Benson, H. Carrillo-Cabada, A. Razavi, M. Cuendet, H. Weinstein, E. Deelman, and M. Taufer. *Graphic Encoding of Proteins for Efficient High-Throughput Analysis*. ICPP 2018.

# High-Throughput Protein Analysis

- Eight biological processes from biological process taxonomy in RCSB-PDB
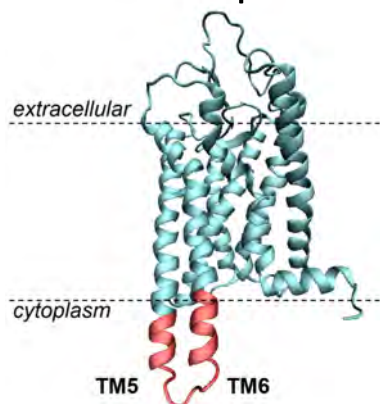- 62,991 proteins from the PDB

Proteins as 3D tens

convolutional neural network

Google's Inception-v3, Gem-Net

Normalized Confusion Matrix - **Accuracy 80.66%**



Ground true

Predictions

24    T. Estrada, J. Benson, H. Carrillo-Cabada, A. Razavi, M. Cuendet, H. Weinstein, E. Deelman, and M. Taufer.
*Graphic Encoding of Proteins for Efficient High-Throughput Analysis*. *ICPP 2018*.
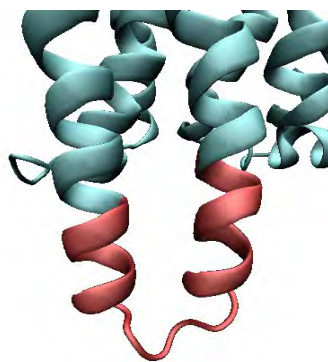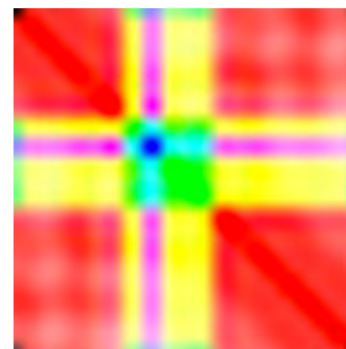
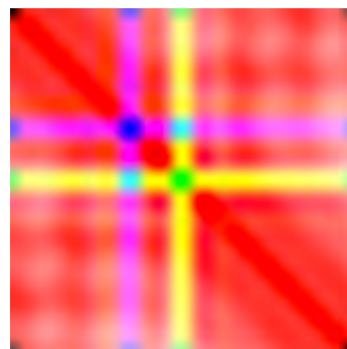# Capturing Changes in Folding with Transfer Learning

Protein: Opsin                Frame 50            Frame 1500           Frame 1950

T. Estrada, et al. **A Graphic Encoding Method for Quantitative Classification of Protein Structure and Representation of Conformational Changes**. 2019
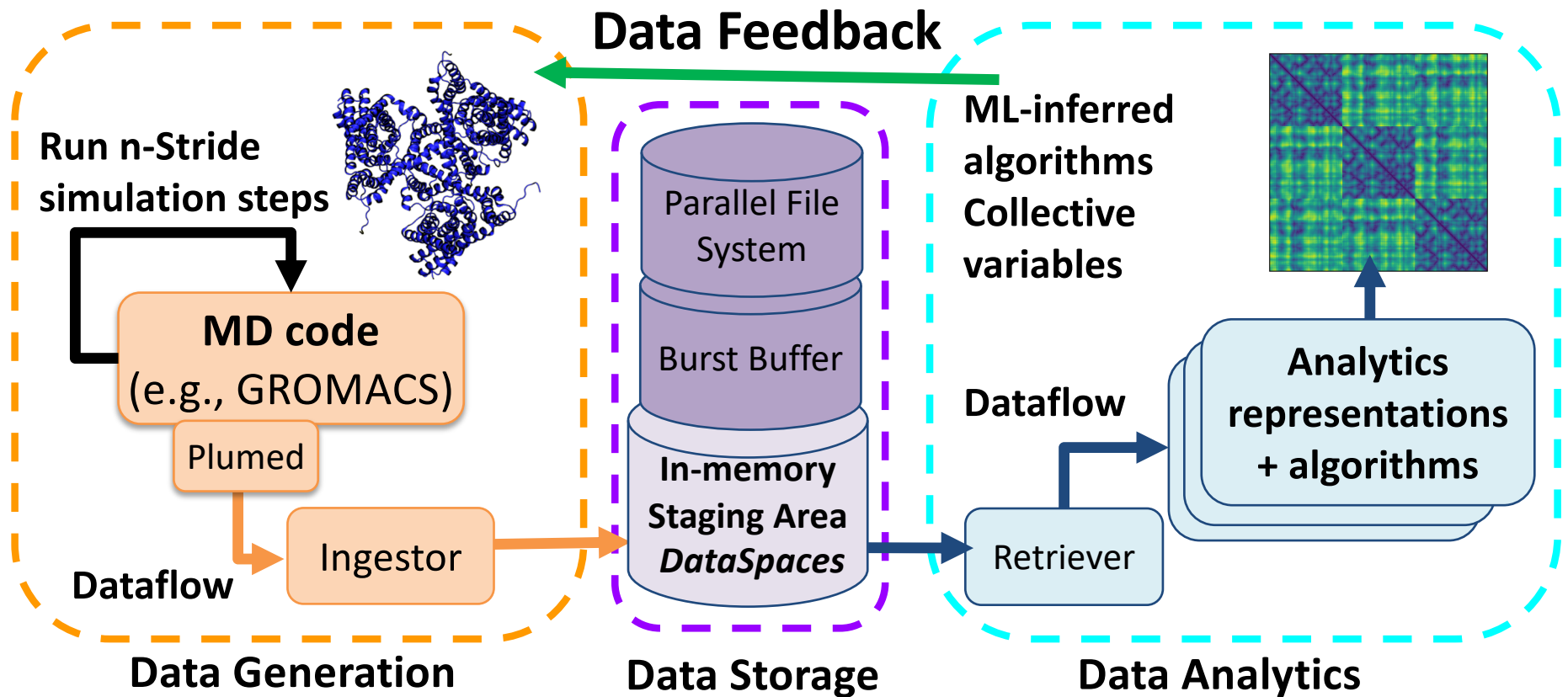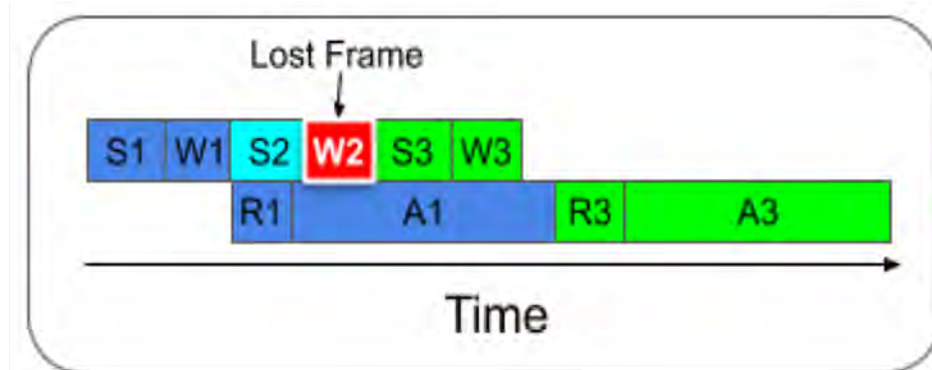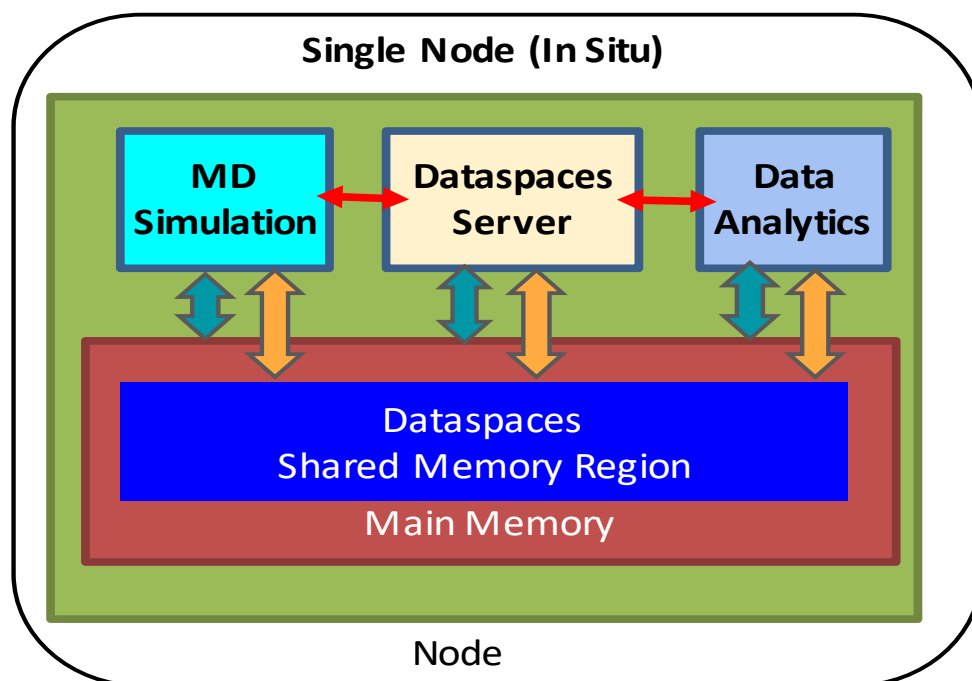
# Two Use Cases

- Extending HPC to integrate data analytics
  - Next generation MD workflows
  - Molecular structures
  - *Data transformation – i.e., capturing information*
  - **Dataflow modeling – *i.e., lost information***
- Extending HPC to connect to the "Edge"
  - Next generation precision farming
  - Soil moisture data
  - *Data prediction – i.e., from coarse- to fine-grained information*
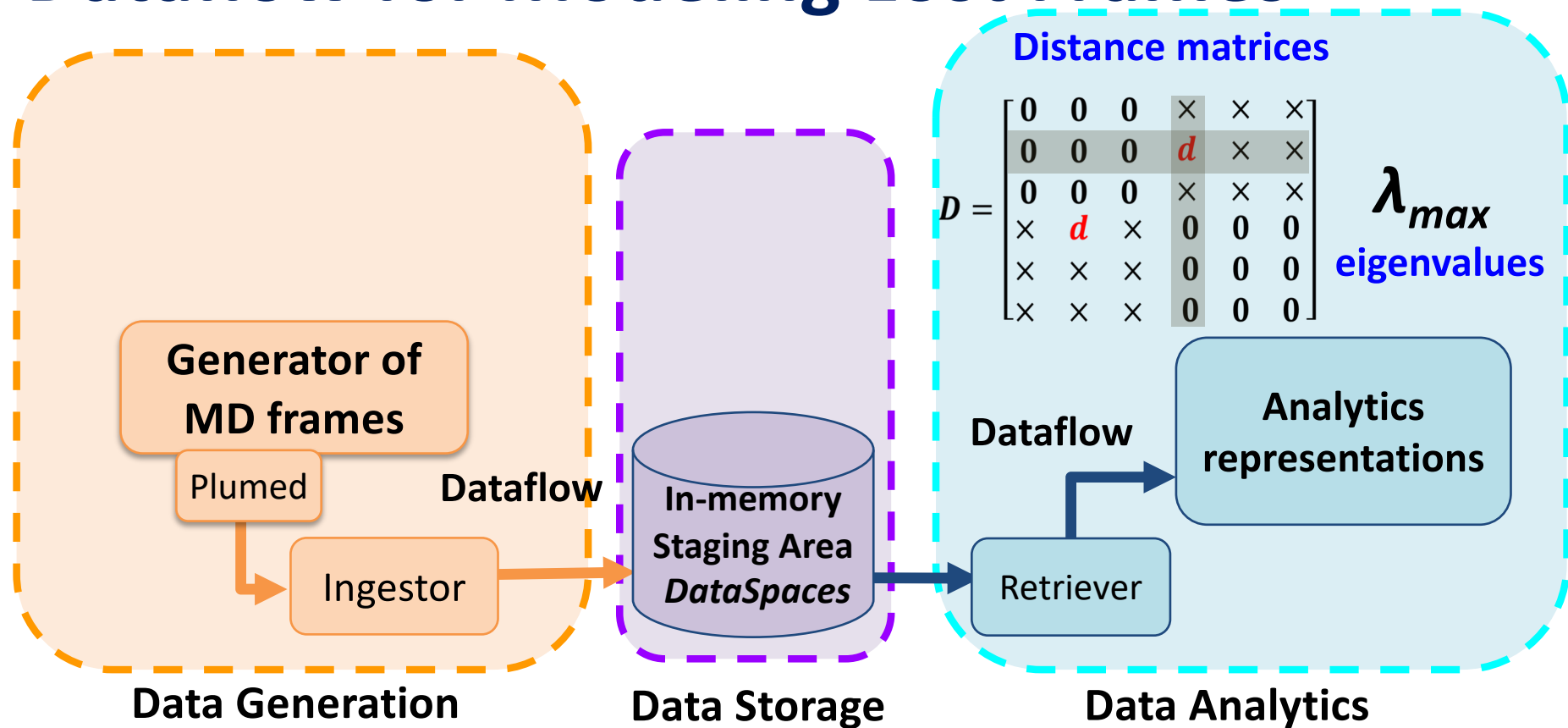
# Building a Closed-loop Workflow



**Data Feedback**

**Run n-Stride simulation steps**

**MD code** (e.g., GROMACS)

Plumed

**Dataflow**

Ingestor

Parallel File System

Burst Buffer

**In-memory Staging Area** *DataSpaces*

ML-inferred algorithms Collective variables

**Dataflow**

Retriever

**Analytics representations + algorithms**

**Data Generation**          **Data Storage**          **Data Analytics**

BIG**ORANGE** BIG**IDEAS**

# Modeling Lost Frames



**Single Node (In Situ)**

MD Simulation ↔ Dataspaces Server ↔ Data Analytics

Dataspaces Shared Memory Region

Main Memory

Node



Lost Frame

S1 | W1 | S2 | **W2** | S3 | W3

R1 | A1 | R3 | A3

Time
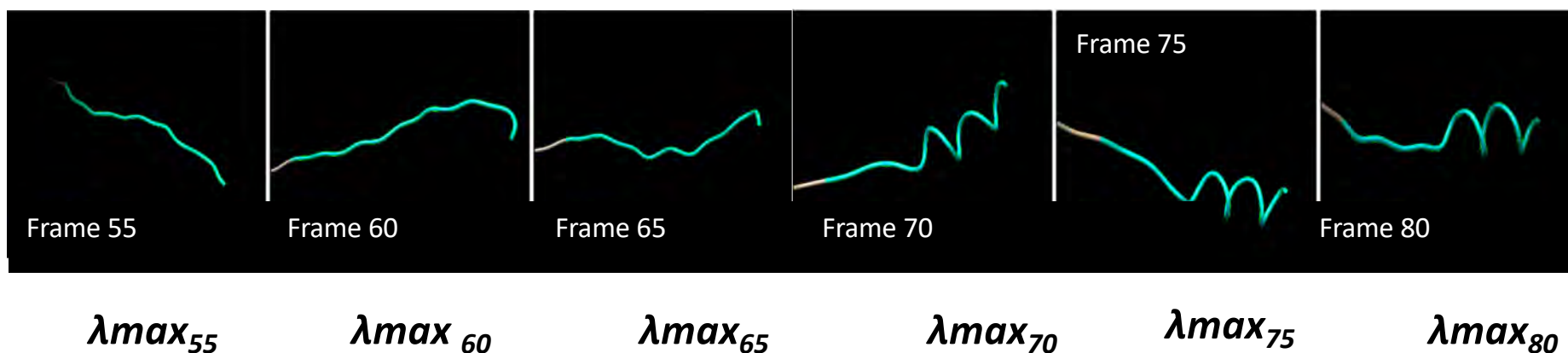
S1, S2, S3: generate MD frame
W1, **W2**, W2: write to shared memory
R1, R2, R3: read from shared memory
A1, A3: analyze frame

BIG ORANGE BIG IDEAS

# Dataflow for Modeling Lost Frames



**Distance matrices**

$$D = \begin{bmatrix} 0 & 0 & 0 & \times & \times & \times \\ 0 & 0 & 0 & d & \times & \times \\ 0 & 0 & 0 & \times & \times & \times \\ \times & d & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \\ \times & \times & \times & 0 & 0 & 0 \end{bmatrix}$$

$\lambda_{max}$ **eigenvalues**

**Generator of MD frames**

Plumed

Ingestor

**Dataflow**

**In-memory Staging Area** *DataSpaces*

**Dataflow**

Retriever

**Analytics representations**

**Data Generation**          **Data Storage**          **Data Analytics**

T. Johnston et al. In-Situ Data Analytics and Indexing of Protein Trajectories. *Journal of Computational Chemistry (JCC),* 38(16):1419-1430, 2017.

BIG**ORANGE** BIG**IDEAS**

# Eigenvalues: Proxy for Structural Changes

Frames (or snapshots) of an MD trajectory with a stride of 5 steps:



Frame 75

Frame 55   Frame 60   Frame 65   Frame 70   Frame 80

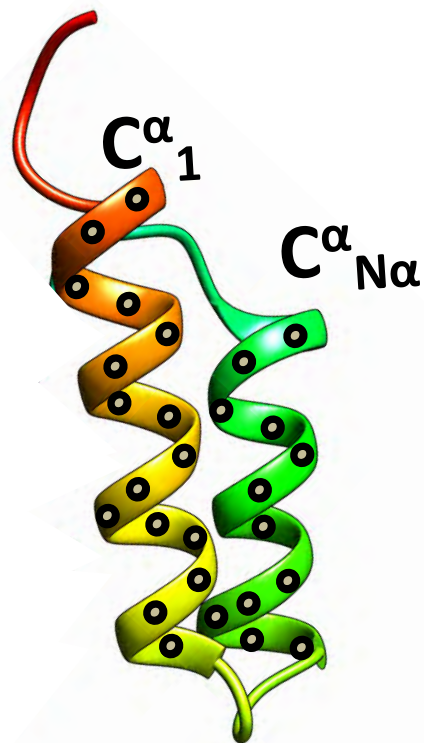$\lambda max_{55}$   $\lambda max_{60}$   $\lambda max_{65}$   $\lambda max_{70}$   $\lambda max_{75}$   $\lambda max_{80}$

*The distance between two max eigenvalues can serve as a proxy for distance between the two associated conformations*

BIG**ORANGE** BIG**IDEAS**

Single frame at time *t*
Nα Cα atoms

$C^{\alpha}_{1}$

$C^{\alpha}_{N\alpha}$

S. Thomas et al. Characterization of In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-generation Supercomputers. *eScience,* 2019.

BIGORANGE
BIGIDEAS

Single frame at time $t$
$N\alpha$ $C^{\alpha}$ atoms

$C^{\alpha}_{1}$

$C^{\alpha}_{N\alpha}$

Distance of two segments with segment length:

$N\alpha/2$ x $C^{\alpha}$ atoms

$(C^{\alpha}_{1} - C^{\alpha}_{N\alpha/2 - 1}) (C^{\alpha}_{N\alpha/2} - C^{\alpha}_{N\alpha/2})$

Single $\lambda_{max}$

|  | $C^{\alpha}_{1}$ | $C^{\alpha}_{N\alpha}$ |
|---|---|---|
| $C^{\alpha}_{1}$ $C^{\alpha}_{N\alpha/2 - 1}$ | 0 | $d_{ij}$ |
| $C^{\alpha}_{N\alpha/2}$ $C^{\alpha}_{N\alpha}$ | dji | 0 |

S. Thomas et al. Characterization of In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-generation Supercomputers. *eScience,* 2019.

BIG ORANGE
BIG IDEAS

Single frame at time $t$
Nα $C^\alpha$ atoms



Distance of two segments with segment length:
$$N\alpha/2 \times C^\alpha \text{ atoms}$$
$$(C^\alpha_1 - C^\alpha_{N\alpha/2 - 1})(C^\alpha_{N\alpha/2} - C^\alpha_{N\alpha/2})$$
Single $\lambda_{max}$

Distances of Nα/2 segments with segment length: $2 \times C^\alpha$ atoms
$$(C^\alpha_1 - C^\alpha_2)(C^\alpha_3 - C^\alpha_4)$$
$$(C^\alpha_5 - C^\alpha_6)(C^\alpha_7 - C^\alpha_8)$$
....
$$(C^\alpha_{N\alpha/2-3} - C^\alpha_{N\alpha/2-2})(C^\alpha_{N\alpha/2-1} - C^\alpha_{N\alpha/2})$$
$\lambda_{max, 1} \lambda_{max, 2} \dots \lambda_{max, N\alpha/2}$

S. Thomas et al. Characterization of In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-generation Supercomputers. *eScience*, 2019.

BIG ORANGE
BIG IDEAS

Segment size = proxy of **number of matrices** and **matrix sizes**

$(\frac{N_\alpha}{m})\frac{(\frac{N_\alpha}{m}-1)}{2}$

$(2m)^2$

# of Matrices

Matrix Size (# of Elements)

Segment Length (m)

Many small matrices

Few large matrices

$C^\alpha_1$     $C^\alpha_{N\alpha}$

$C^\alpha_1$

$C^\alpha_{N\alpha/2}-1$

$C^\alpha_{N\alpha/2}$

$C^\alpha_{N\alpha}$

| 0 | $d_{ij}$ |
| dji | 0 |

$C^\alpha_1$   $C^\alpha_4$

$C^\alpha_1$
$C^\alpha_2$
$C^\alpha_3$
$C^\alpha_4$

| 0 | dij |
| dji | 0 |

S. Thomas et al. Characterization of In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-generation Supercomputers. *eScience,* 2019.

BIG**ORANGE** BIG**IDEAS**

# 2-step Model: Fraction of Analyzed Frames

Observables: small segment lengths
(i.e., 2 , 4 , 6 , 8 , 10 , 12 , 14, 16, and 18 )

Trp cage 12,619 atoms    T cell receptor  81,092 atoms

Gltph 270,088  atoms

BIG ORANGE
BIG IDEAS

# 2-step Model: Fraction of Analyzed Frames

Observables: segment lengths
(i.e., 2 , 4 , 6 , 8 , 10 , 12 , 14, 16, and 18 )

Model: polynomial model of degree 2
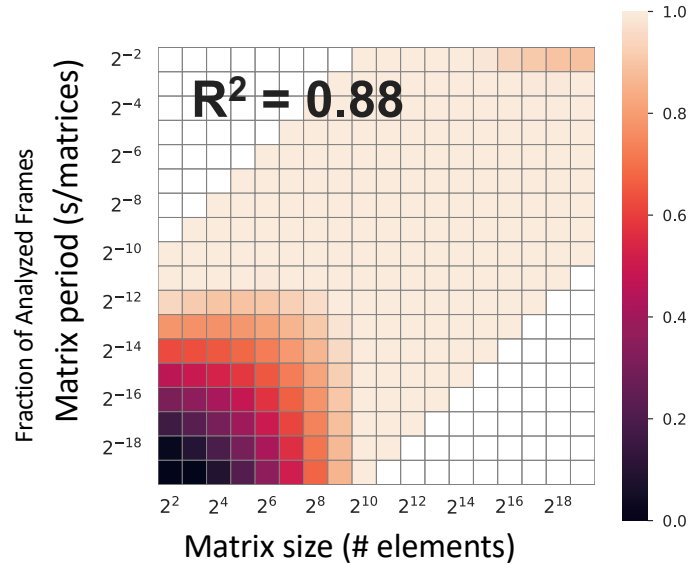obtained by least-square fitting observables



$R^2 = 0.88$

S. Thomas et al.  Characterization of In Situ and In Transit Analytics of Molecular Dynamics Simulations for Next-generation Supercomputers. *eScience*, 2019. *(Submitted)*

BIG**ORANGE**
BIG**IDEAS**

# 2-step Model: Fraction of Analyzed Frames

Observables

Degree 2 polynomial model

$R^2 = 0.88$

Absolute error between data and fitting model

# 2-step Model: Frames Distribution

- Given a trajectory, we model the proportions $p$ and $q$ of analyzed frames ($f$) with periods $k$ and $k+1$

  - Example: Gltph (27,000 atoms and TPS 318), trajectory of 1,000 frames.
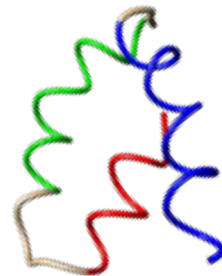
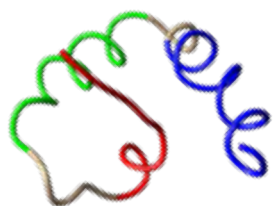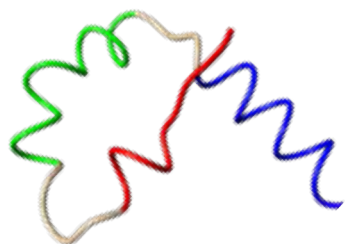# Case Study: 1BDD Protein Conformations

Frame 1330        Frame 1360        Frame 1390
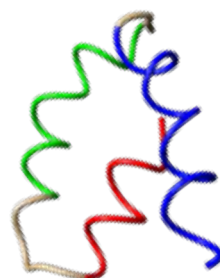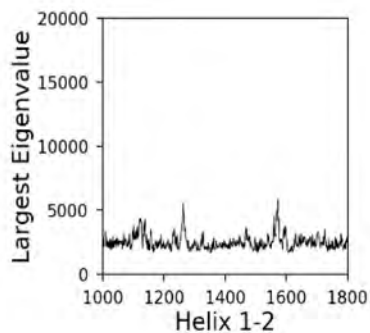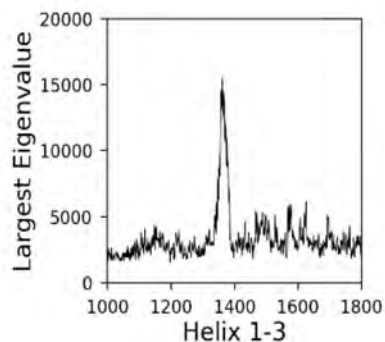
# Case Study: 1BDD Protein Conformations

Frame 1330          Frame 1360          Frame 1390



Helix 1-2: $\lambda_{max}$     Helix 1-3: $\lambda_{max}$     Helix 2-3: $\lambda_{max}$

BIG ORANGE
BIG IDEAS

# Case Study: 1BDD Protein Conformations
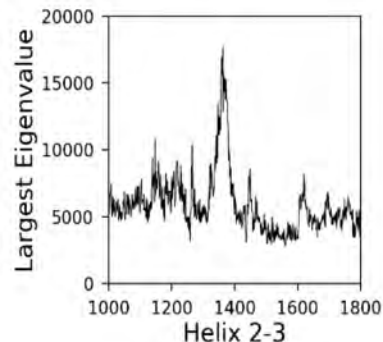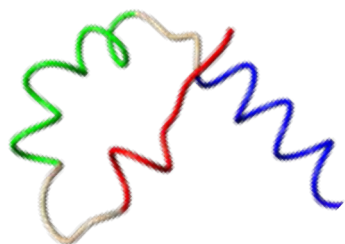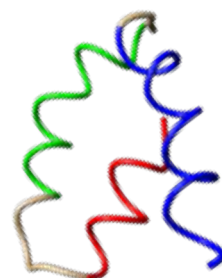
Frame 1330

Frame 1360

Frame 1390



Helix 1-2: $\lambda_{max}$

Helix 1-3: $\lambda_{max}$

Helix 2-3: $\lambda_{max}$

# Two Use Cases

- Extending HPC to integrate data analytics
  - Next generation MD workflows
  - Molecular structures
  - ***Data transformation – i.e.,*** *capturing information*
  - ***Dataflow modeling –*** *i.e., lost information*

- Extending HPC to connect to the "Edge"
  - Next generation precision farming
  - Soil moisture data
  - ***Data prediction*** *– i.e., from coarse- to fine-grained information*

# Soil Moisture Data for Precision Farming

- Satellites collect raster data across the surface of the Earth



Image Source: http://www.esa-soilmoisture-cci.org/

# Soil Moisture: Incomplete Data



Soil moisture (m³/m³)

High : 0.86247

December 2000 Averages

Low : 0

Image source: Ricardo Llamas, University of Delaware
Data source: ESA-CCI soil moisture database (http://www.esa-soilmoisture-cci.org/)

BIG ORANGE BIG IDEAS

# Soil Moisture: Incomplete Data



**Soil moisture** (m³/m³)

High : 0.86247

December 2000
Averages

Low : 0

Causes of missing data:
- dense vegetation
- snow/ice cover
- frozen surface
- extremely dry surface

Image source: Ricardo Llamas, University of Delaware
Data source: ESA-CCI soil moisture database (http://www.esa-soilmoisture-cci.org/)

BIG**ORANGE**
BIG**IDEAS**

# Soil Moisture: Coarse-grained Data



Original Resolution
27 km × 27 km

Desired Resolution
1 km × 1 km

BIG ORANGE
BIG IDEAS

# Building a Closed-loop Workflow



Weather Data
NOAA

Landscape
Surface DSM

Soil Moisture
ESA-CCI

Course-grained,
incomplete data

Satellite and
sensors data

**Data Generation**

# Building a Closed-loop Workflow

**Analytics representations + algorithms**

Data predictions

High : 0.35
Low : 0.00
Ecoregion 8.5.1
Mid Atlantic Coastal Plains

Weather Data NOAA

Landscape Surface DSM

Soil Moisture ESA-CCI

Satellite and sensors data

Course-grained, incomplete data

**Data Analytics**

**Data Generation**

# Building a Closed-loop Workflow



**Fire Dynamics Simulator**

Soil moisture integrated in FDS for:
- Controlled combustion;
- Wildfire propagation

**Compute**

Shifter and Docker

Singularity

Kubernetes

Fine-grained, complete data

**Analytics representations + algorithms**

Data predictions

**Data Analytics**

Weather Data NOAA

Landscape Surface DSM

Soil Moisture ESA-CCI
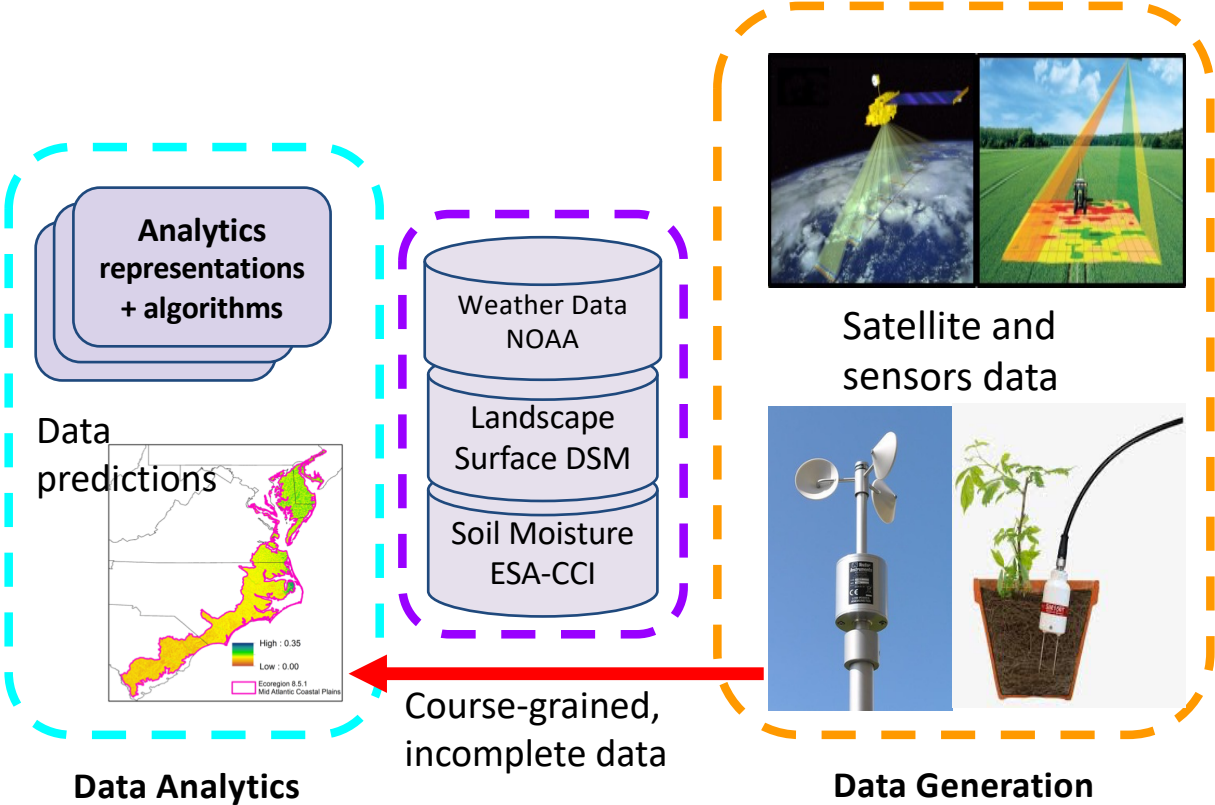
Course-grained, incomplete data

Satellite and sensors data

**Data Generation**

# Building a Closed-loop Workflow

**Data Feedback**

Fire Dynamics Simulator

Soil moisture integrated in FDS for:
- Controlled combustion
- Wildfire propagation

**Compute**

Shifter and Docker

Singularity

Kubernetes

**Analytics representations + algorithms**

Data predictions

Fine-grained, complete data

Weather Data NOAA

Landscape Surface DSM

Soil Moisture ESA-CCI

Course-grained, incomplete data

High : 0.35
Low : 0.00
Ecoregion 8.5.1
Mid Atlantic Coastal Plains

**Data Analytics**

Satellite and sensors data

**Data Generation**

# Augmenting HPC with the Cloud



Source: https://www.nextplatform.com/2018/02/26/adaptive-approach-bursting-hpc-cloud/

# Building a Closed-loop Workflow

**Data Feedback**

Fire Dynamics Simulator



Soil moisture integrated in FDS for:
- Controlled combustion
- Wildfire propagation

**Shifter and Docker**

**Singularity**

**Kubernetes**

**Analytics representations + algorithms**

Data predictions



**Weather Data NOAA**

**Landscape Surface DSM**

**Soil Moisture ESA-CCI**



Satellite and sensors data



Fine-grained, complete data

Course-grained, incomplete data

**Compute**

**Data Analytics**

**Data Generation**

# Hybrid Algorithms for Analytics



**Data Feedback**

**Analytics representations + algorithms**

Fine-grained, complete data

Course-grained, incomplete data

**Compute**

**Data Analytics**

**Data Generation**

BIG**ORANGE** BIG**IDEAS**
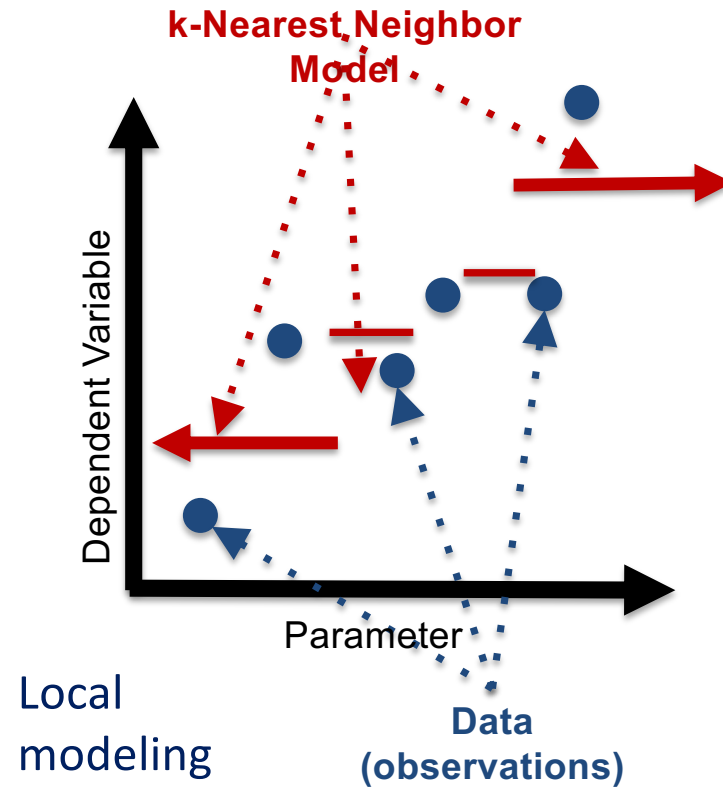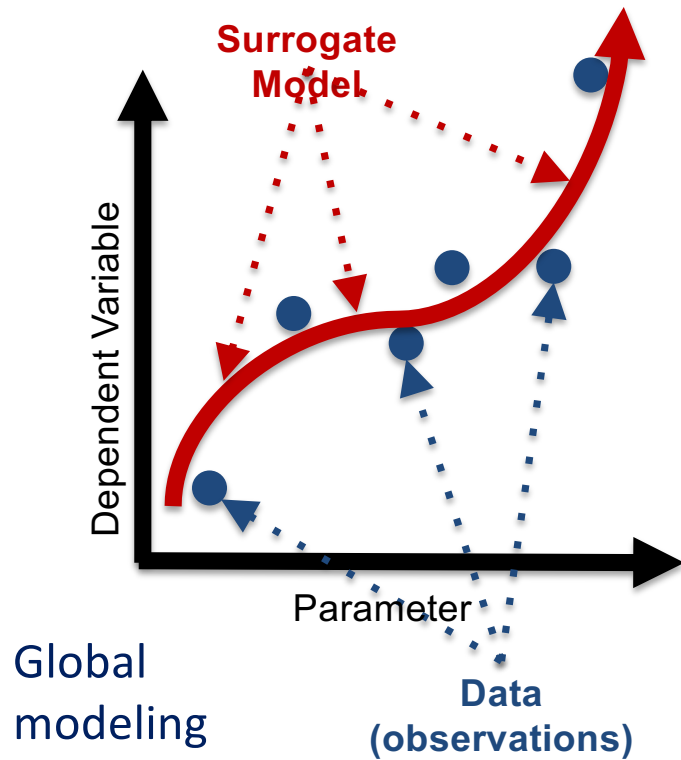
# Global versus Local Data Modeling

# Global versus Local Data Modeling

# SBM vs. k-NN

**Observations**
**Piecewise linear data**
**Surrogate-based model**
**2-Nearest Neighbor Model**



Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33
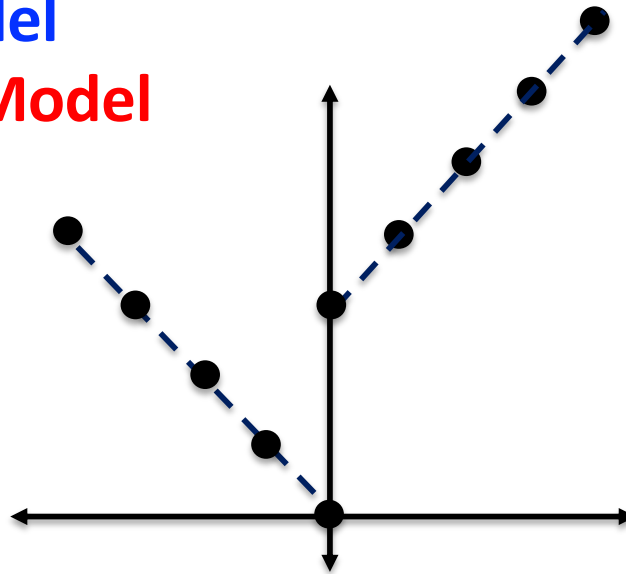
# SBM vs. k-NN

**Observations**
**Piecewise linear data**
**Surrogate-based model**
**2-Nearest Neighbor Model**

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

BIG**ORANGE**
BIG**IDEAS**

# SBM vs. k-NN

**Observations**
**Piecewise linear data**
**Surrogate-based model**
**2-Nearest Neighbor Model**

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

BIG**ORANGE**
BIG**IDEAS**

# SBM vs. k-NN

**Observations**
**Piecewise linear data**
**Surrogate-based model**
**2-Nearest Neighbor Model**

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

# SBM vs. k-NN

**Observations**
**Piecewise linear data**
**Surrogate-based model**
**2-Nearest Neighbor Model**

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

# Hybrid Piecewise Polynomial Modeling (HYPPO)

k Nearest Neighbors
- Use **local** data
- Compute the average
  (*many simple local models*)

Surrogate-Based Modeling
- Use **all** sampled data
- Construct **one** **polynomial**
  (*single complex global model*)

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

BIG**ORANGE**
BIG**IDEAS**

# Hybrid Piecewise Polynomial Modeling (HYPPO)
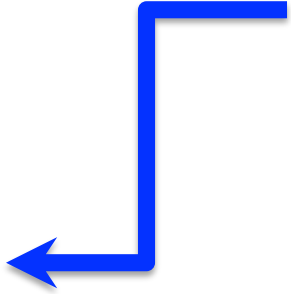
k Nearest Neighbors
- Use **local** data
- Compute the average
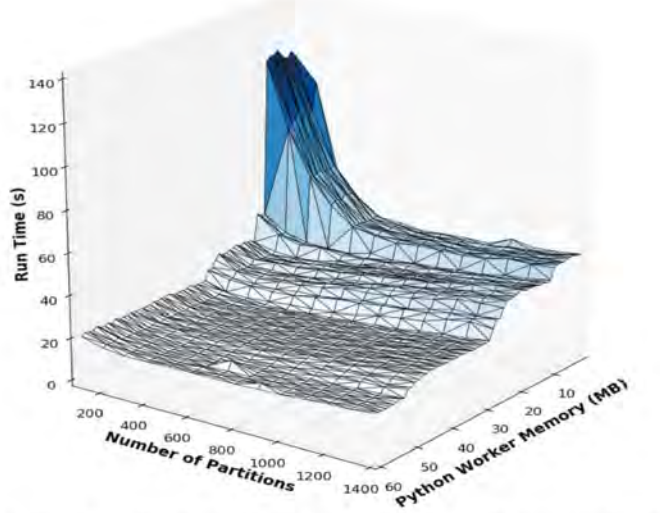  (*many simple local models*)

Surrogate-Based Modeling
- Use **all** sampled data
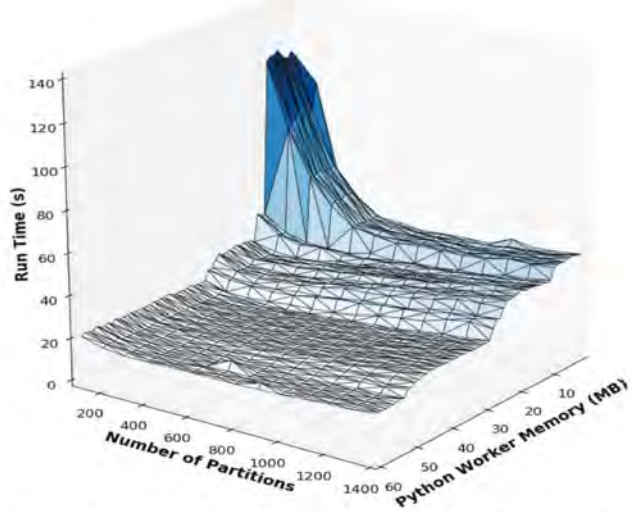- Construct **one** **polynomial**
  (*single complex global model*)

## Hybrid modeling → HYPPO
- Use **local** data
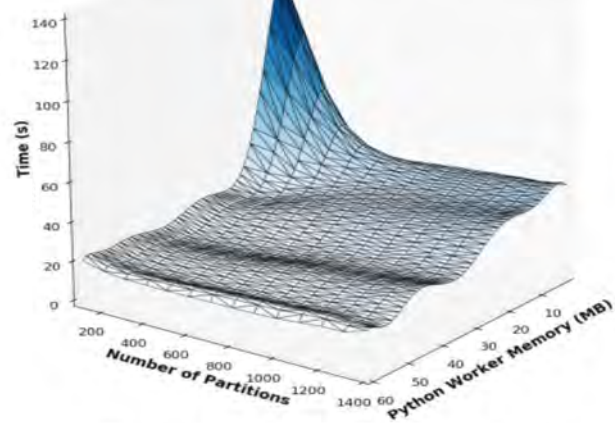- Construct many **polynomials**
  (*many complex local models*)

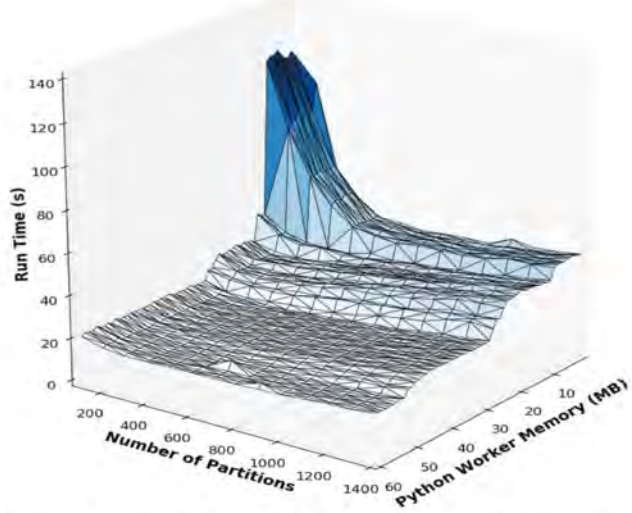Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

BIG**ORANGE**
BIG**IDEAS**

## Observations

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

Observations

Surrogate-based Model

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33
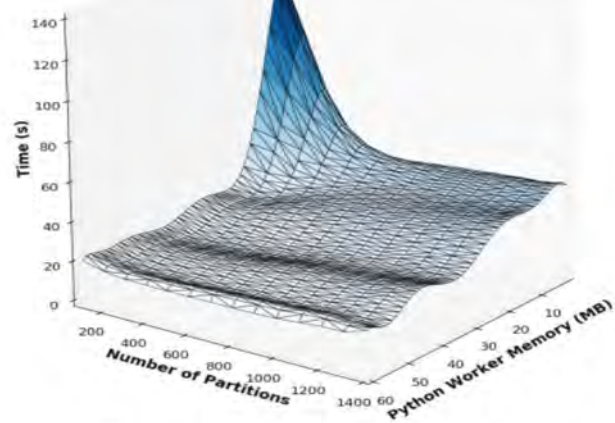
BIG**ORANGE**
BIG**IDEAS**

Observations

Surrogate-based Model

k Nearest Neighbors Model

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33
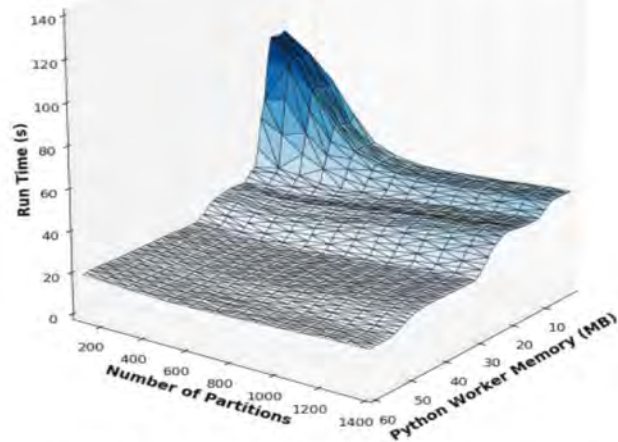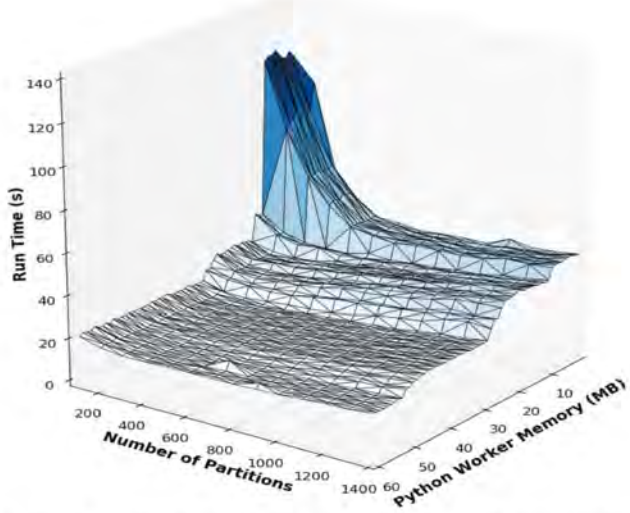
BIG ORANGE
BIG IDEAS

Observations
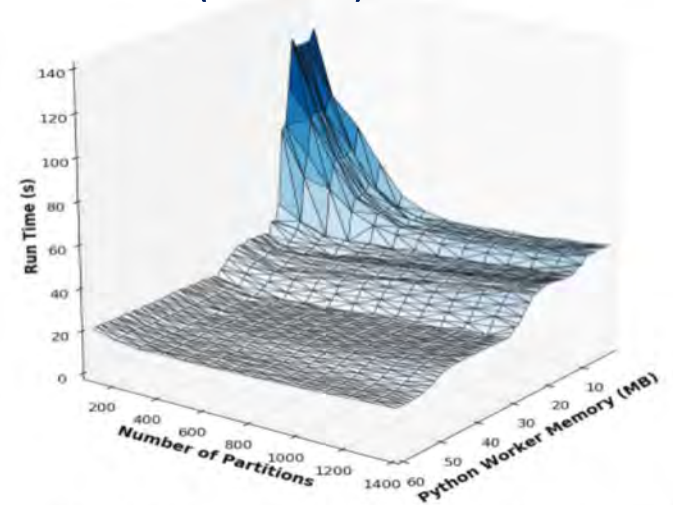
Surrogate-based Model

k Nearest Neighbors Model

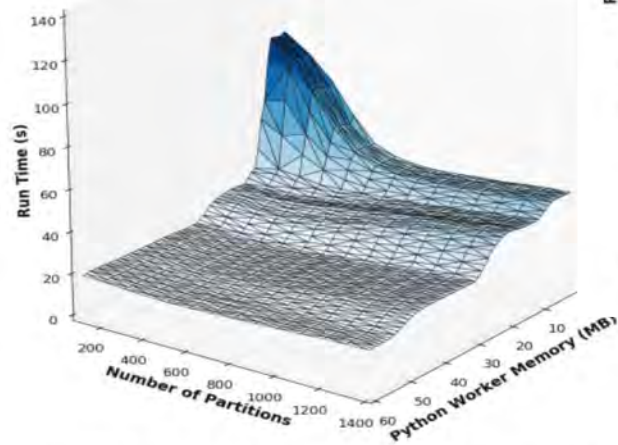Hybrid Piecewise Polynomial Model (HYPPO)

Johnston et al., "HYPPO: A Hybrid, Piecewise Polynomial Modeling Technique for Non-Smooth Surfaces." SBAC-PAD 2016: 26-33

# Case Study: Fine-grained Modeling of Mid-Atlantic Region



Soil moisture predictions

Model degrees

67    D. Rorabaugh, M. Guevara, R. Llamas, J. Kitson, R. Vargas, M. Taufer. SOMOSPIE: A modular SOil MOisture SPatial Inference Engine based on data driven decision. **arXiv:1904.07754, 2019**

# Challenges and Opportunities (I)

- Two trends in HPC are impacting scientific applications:
    - Convergence of simulations and data analytics
    - Emergence of edge computing
- Applications in precision medicine and precision farming are leveraging these trends
- There is the need to further integrate these trends in HPC
    - New challenges and opportunities for the HPC community

# Challenges and Opportunities (II)

- *Efficiency:* Optimize performance and power usage associated to data generation, movement, and analytics

- *Non-invasive:* Capture knowledge from data without rewriting simulations' legacy codes or simulations' scripts

- *Generality:* Build workflows that support different types of analytics across different applications and different data

- *Portability:* Execute combined compute and analytics across different systems, including the edge, and with heterogenous resources

- *Scalability:* Design methods for knowledge discovery at scale (e.g., scalable ML algorithms) for "compute + analytics + data" workflows