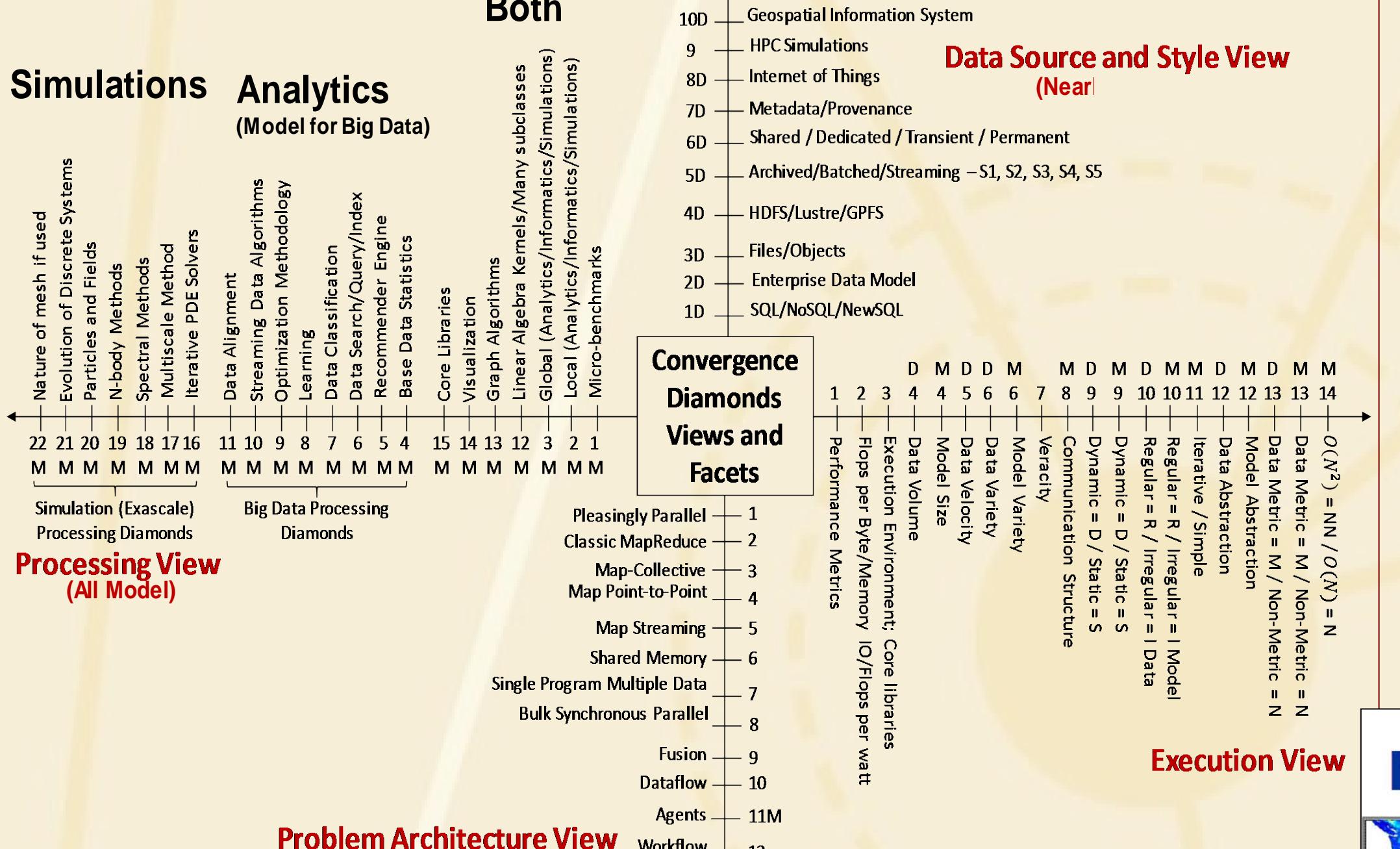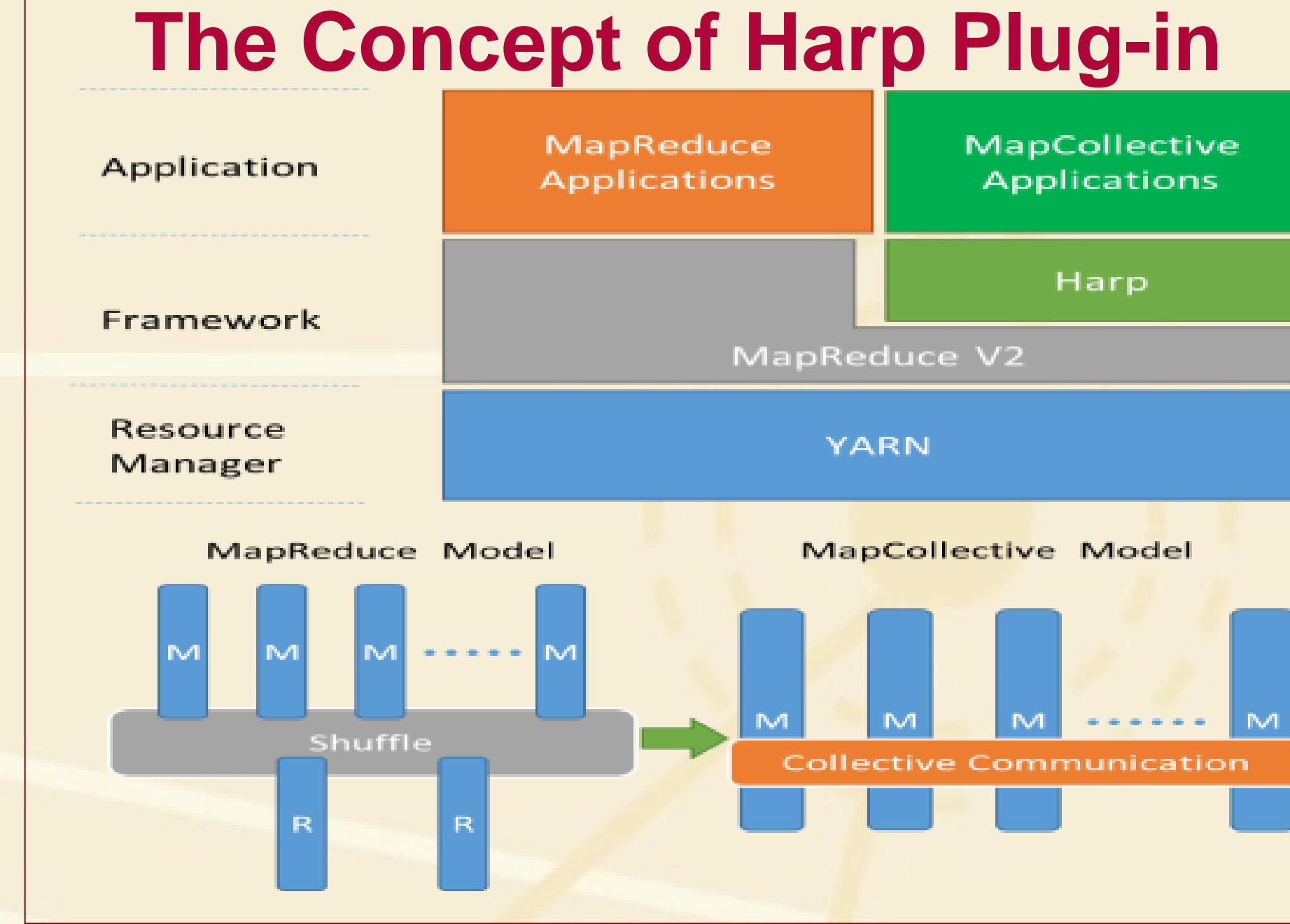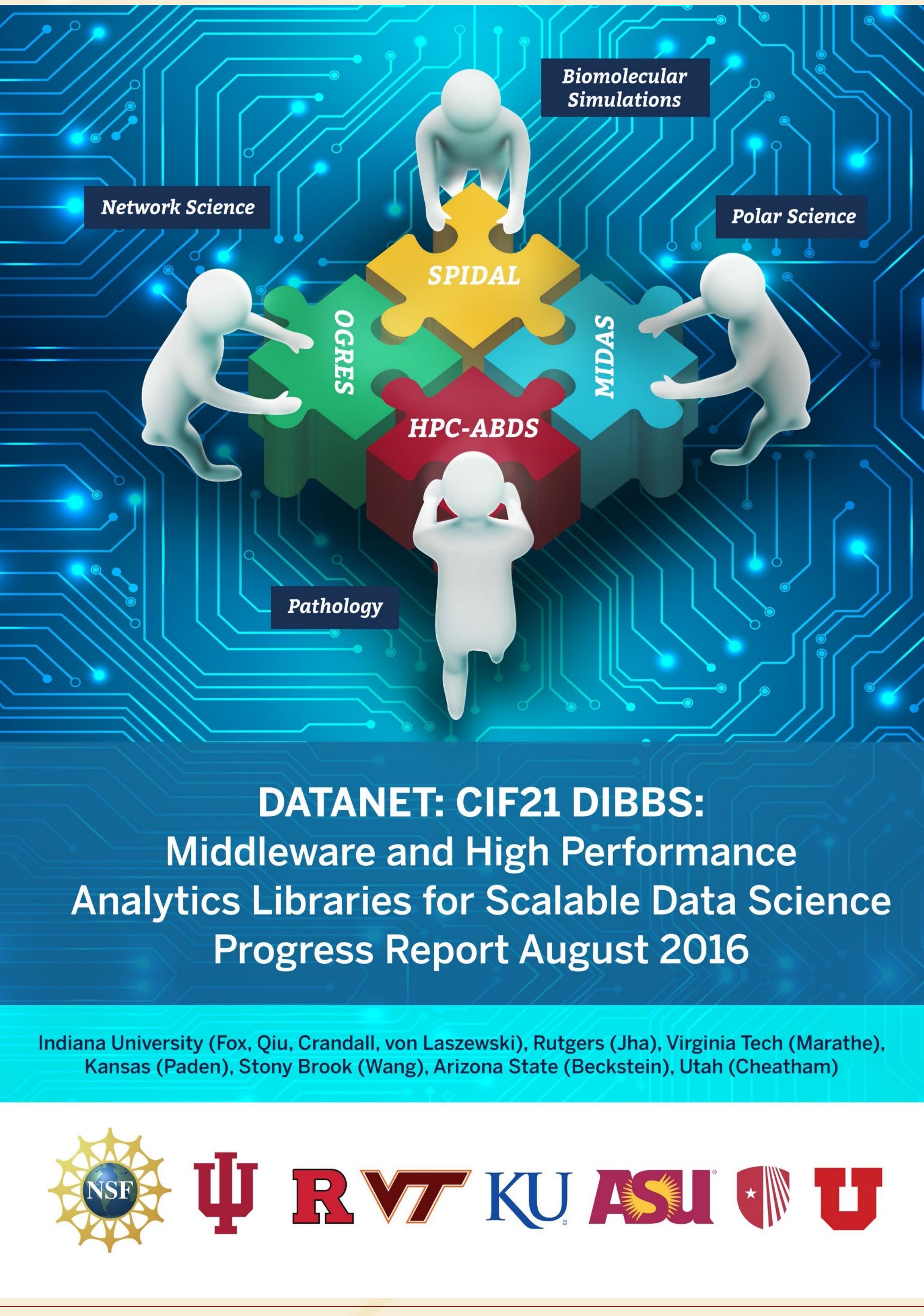## Components of NSF 144305

- Design and build **SPIDAL** – a Scalable Parallel Interoperable Data Analytics Library
  - Domain specific libraries – mainly from project
  - Core Machine Learning Libraries
- High Performance for Java, Libraries and MIDAS
- NIST **Big Data Application Analysis**: Features of data intensive Applications deriving 50 Ogres and 64 Convergence Diamonds
- HPC-ABDS: Cloud-HPC Interoperable software with performance of HPC and rich functionality of commodity Apache Stack
- Implementations: HPC and Clouds with DevOps
- Applications: Biomolecular Simulations, Network Science, Epidemiology, Computer Vision, Spatial Geographical Information Systems, Remote Sensing for Polar Science and Pathology Informatics.
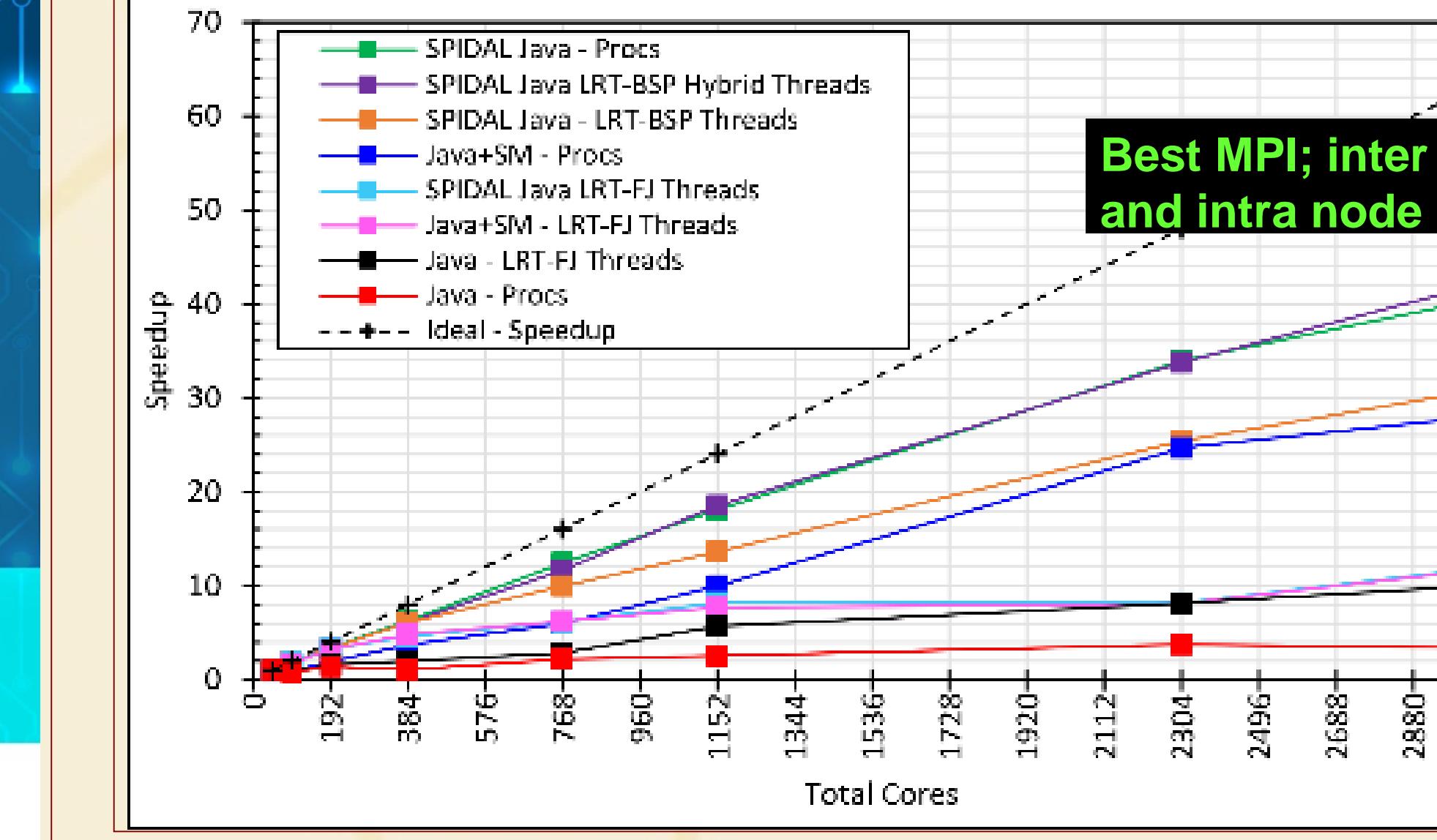
## Building Blocks of Proposal

### 64 Features in 4 views for HPC-Big Data Converged Classification



## HPC-ABDS Apache Big Data Stack



21 layers
Over 350
Software
Packages

January
29
2016



### DATANET: CIF21 DIBBS:
Middleware and High Performance
Analytics Libraries for Scalable Data Science
Progress Report August 2016

Indiana University (Fox, Qiu, Crandall, von Laszewski), Rutgers (Jha), Virginia Tech (Marathe), Kansas (Paden), Stony Brook (Wang), Arizona State (Beckstein), Utah (Cheatham)

## MIDAS and SPIDAL Java High Performance Middleware and Language

http://dsc.soic.indiana.edu/publications/SPIDAL-DIBBSreport_July2016.pdf

http://spidal.org

## Applications in action with MIDAS/SPIDAL

### Multi-Scale Imaging and Spatial Data Analytics



Earth   Humans   Patients   Vessels and Cells

Data Integration   Data Extraction

2D and 3D image analysis

2D geospatial data   2D imaging spatial data   3D imaging spatial data

#### Hadoop-GIS Query Engine
Global Spatial Indexes
RESQUE
Tile Spatial Indexes
Spatial Index Builder   Spatial Query Processor   Boundary Handling
Hadoop/MapReduce/Spark

2D Spatial queries and analytics

#### 3D Spatial Data Partitioning
Fixed Grid   Octree   Binary Space   Hilbert Curve
3D Spatial Query Engine
Global Indexes   On-Demand Spatial Query Engine
Cuboid Indexes   Object Indexes   Structural Indexes
3D Spatial Query Processor
Hadoop/MapReduce

3D Spatial queries and analytics

## The Concept of Harp Plug-in



## DA-MDS speedup for 200K with different optimization techniques
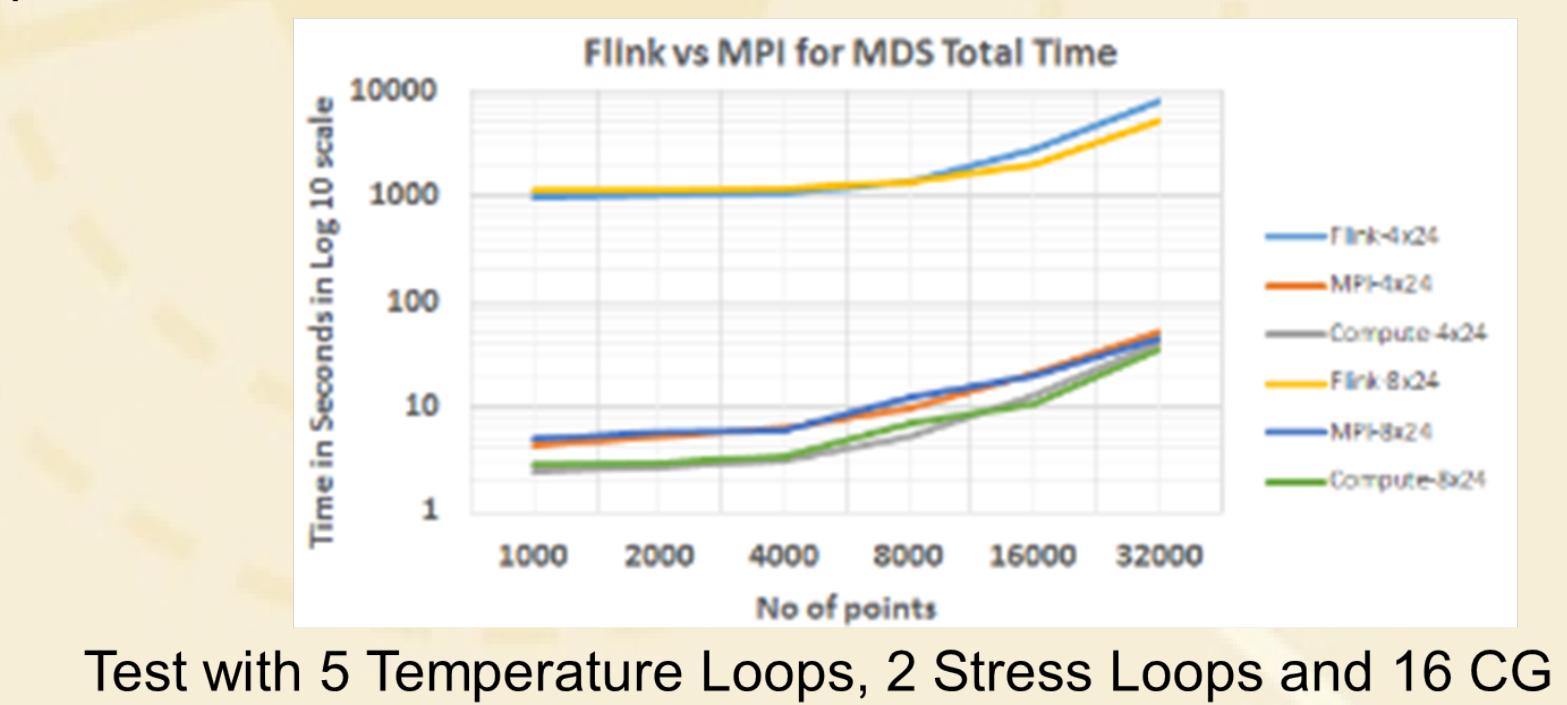


Best MPI; inter and intra node

## Multidimensional Scaling with Flink

- Projects NxN distance matrix to NxD matrix where N is the number of points and D is the target dimension.
- Row wise partitioning of NxN matrix for parallelization
- Our algorithm uses deterministic annealing and has three nested loops that called Temperature, Stress and CG (Conjugate Gradient) in that order.

### Flink vs MPI DA-MDS Performance

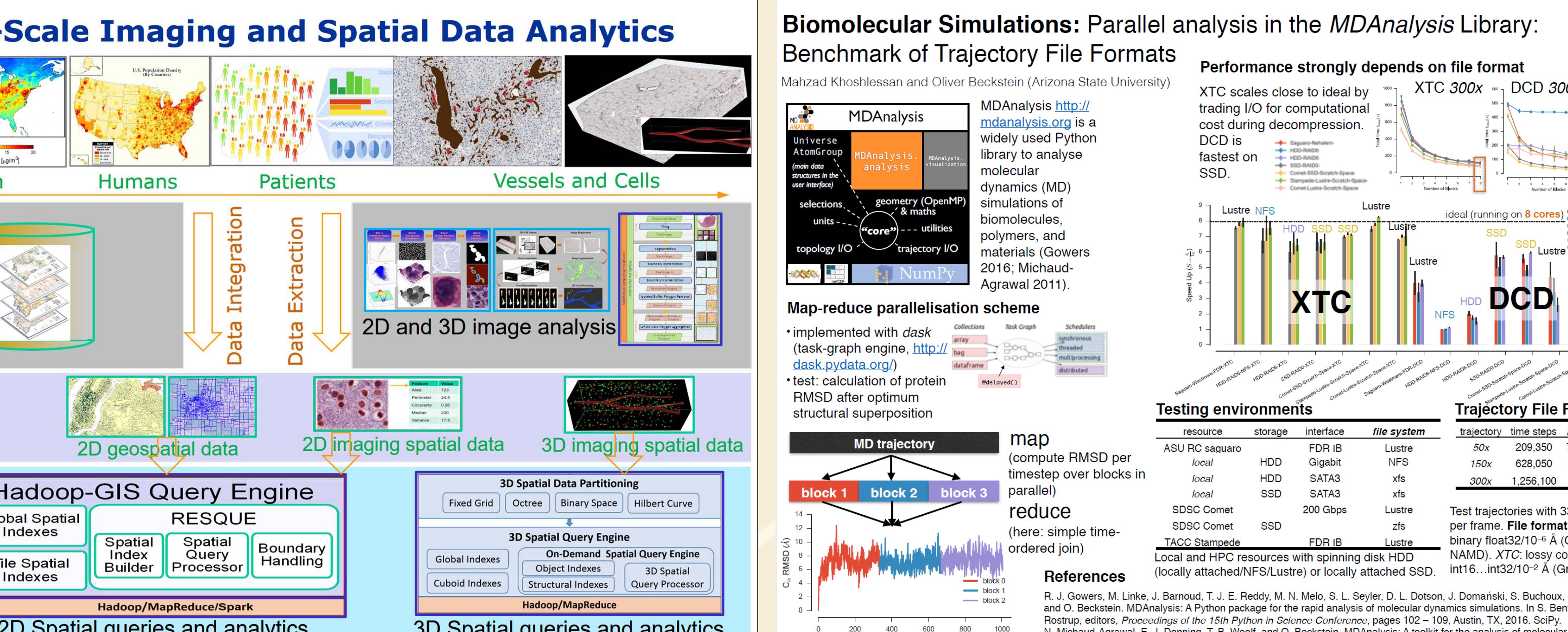Total time of MPI Java and Flink MDS implementations for 96 and 192 parallel tasks



Test with 5 Temperature Loops, 2 Stress Loops and 16 CG Loops.

## Biomolecular Simulations: Parallel analysis in the *MDAnalysis* Library: Benchmark of Trajectory File Formats

Mahzad Khoshlessan and Oliver Beckstein (Arizona State University)



MDAnalysis http://mdanalysis.org is a widely used Python library to analyse molecular dynamics (MD) simulations of biomolecules, polymers, and materials (Gowers 2016; Michaud-Agrawal 2011).

### Performance strongly depends on file format

XTC scales close to ideal by trading I/O for computational cost during decompression. DCD is fastest on SSD.

### Map-reduce parallelisation scheme

- implemented with *dask* (task-graph engine, http://dask.pydata.org)
- test: calculation of protein RMSD after optimum structural superposition

MD trajectory
map (compute RMSD per timestep over blocks in parallel)
reduce (here: simple time-ordered join)

block 1   block 2   block 3

#### Testing environments

| resource | storage | interface | file system |
|---|---|---|---|
| ASU RC saguaro | FDR IB | gigabit | NFS |
| | HDD | SATA3 | xfs |
| local | HDD | SATA3 | xfs |
| | SSD | SATA3 | xfs |
| SDSC Comet | 200 Gbps | | zfs |
| SDSC Comet | SSD | | xfs |

Local and HPC resources with spinning disk HDD (locally attached/NFS/Lustre) or locally attached SSD.

#### Trajectory File Formats

| trajectory | time steps | DCD | XTC |
|---|---|---|---|
| f40x | 209,350 | 7.9G | 2.6G |
| f40x | 628,050 | 24G | 7.5G |
| 300x | 1,256,100 | 47G | 15G |

Test trajectories with 3341 atoms per frame. File formats: DCD: binary float32/10^-6 Å (CHARMM, NAMD). XTC: lossy compressed int16...int32/10^-3 Å (GROMACS)

#### References
R. J. Gowers, M. Linke, J. Barnoud, T. J. E. Reddy, M. N. Melo, S. L. Seyler, D. L. Dotson, J. Domanski, S. Buchoux, I. M. Kenney, and O. Beckstein. MDAnalysis: A Python package for the rapid analysis of molecular dynamics simulations. In S. Benthall and S. Rostrup, editors, *Proceedings of the 15th Python in Science Conference*, pages 102 – 109, Austin, TX, 2016. SciPy.
N. Michaud-Agrawal, E. J. Denning, T. B. Woolf, and O. Beckstein. MDAnalysis: A toolkit for the analysis of molecular dynamics simulations. *J Comp Chem*, 32:2319–2327, 2011.

## Applications in action

### MIDAS and Biomolecular Simulations

- **Previous Work (Middleware):** We have created a hybrid execution framework by coupling the Pilot-abstraction with Apache Hadoop and Spark. This enhances the ability of applications to utilizing HPC resources in conjunction with Hadoop frameworks concurrently or sequentially as needed.
- **Current Work: (Applications)** [i] We run Molecular Dynamics trajectory analysis using the hybrid execution framework, on a set of heterogeneous resources using task level parallel implementation and MapReduce programming through Spark. [ii] We are investigating the characteristics of these type of applications and the performance trade-off. Initial results show that these types of applications can greatly benefit from such an approach.
- **Future Work (Streaming):** Large volumes of streaming data that are real time or close to realtime need to be processed or saved for later analysis and thus need high performance processing capabilities. We are integrating Kafka to handle messages, Spark for processing the data to RADICAL-Pilot to handle the deployment and the job submissions on HPC systems (Pilot-Streaming) and will employ these for HPC simulations and analysis.
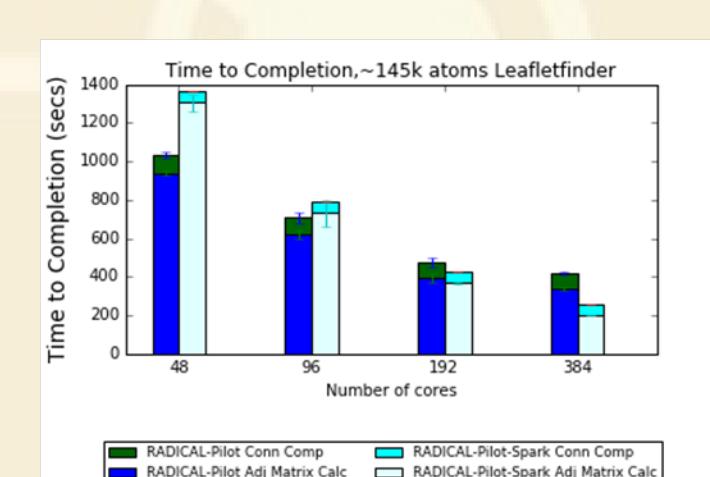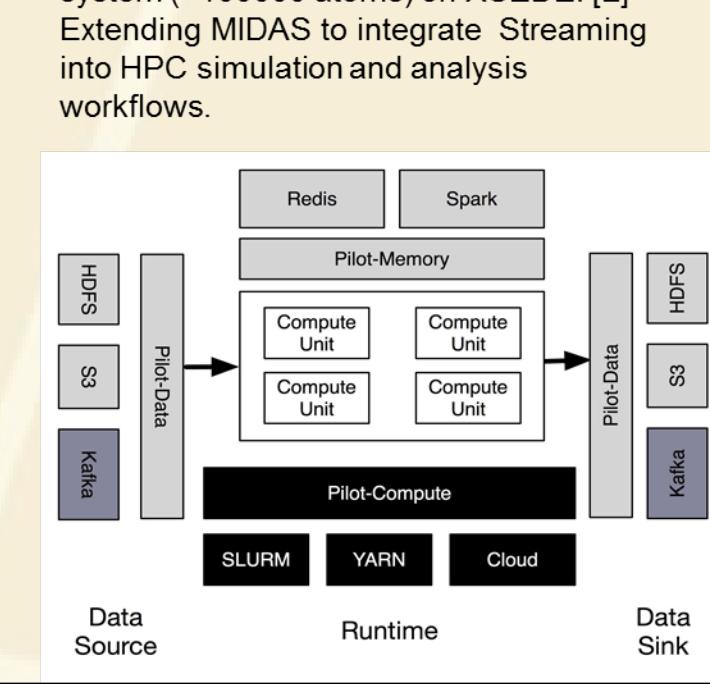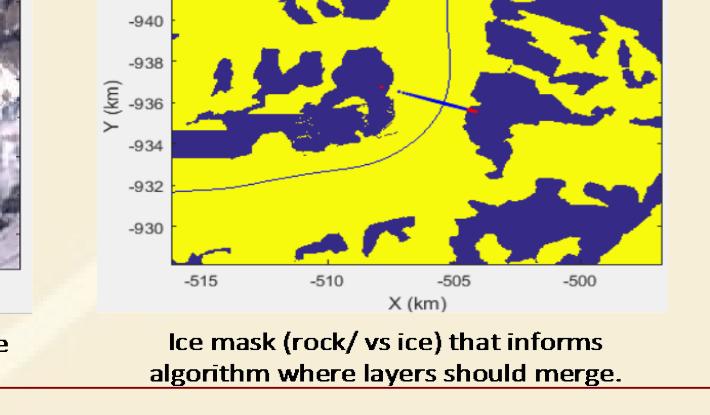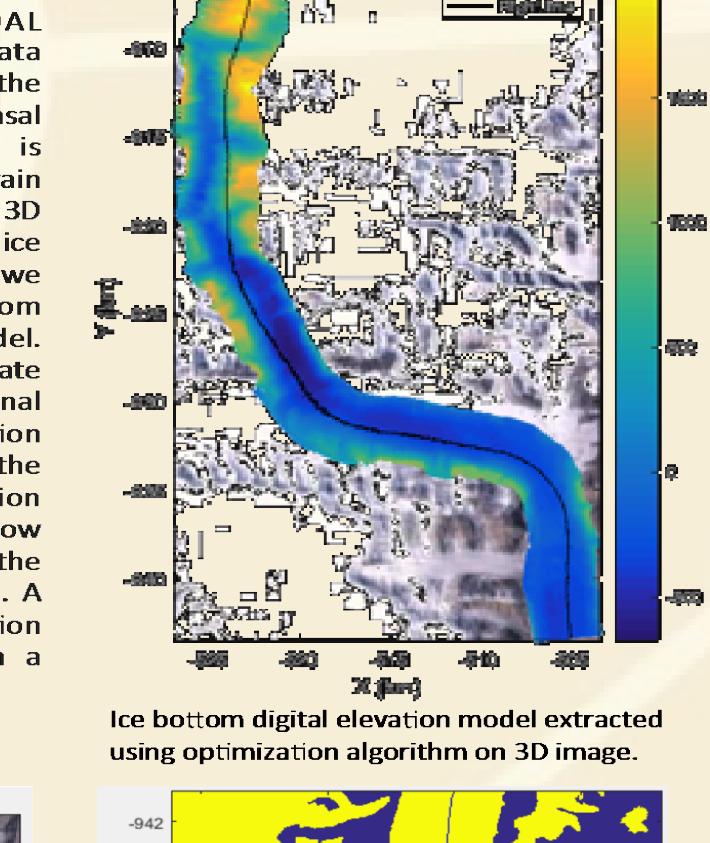


Figure Caption: [U] Data from MD Trajectory Analysis using MIDAS on SPARK for long trajectory for large system (>100000 atoms) on XSEDE. [L] Extending MIDAS to integrate Streaming into HPC simulation and analysis workflows.

### Polar Remote Sensing Algorithms

Many problems can be posed as mathematical optimization tasks, where the goal is to find the values for a set of variables that minimize a particular energy function of those variables. In applications, these tasks often arise in the context of fitting a model to data, for instance, nearly all machine learning algorithms are simply optimizing a set of model parameters by minimizing an objective function that measures the error of the model on a set of labeled training examples. Meanwhile data mining algorithms like K-means and LDA are similarly fitting model parameters to data to minimize some residual error. Many of the imaging applications in Section 5, for example, are simply finding a "simple" model to explain complicated image data, e.g. a model that compactly describes a noisy radar echogram in terms of ice layer structure, or a segmentation that compactly describes a high-resolution CT scan. this section describes our approach to many classes of optimization.
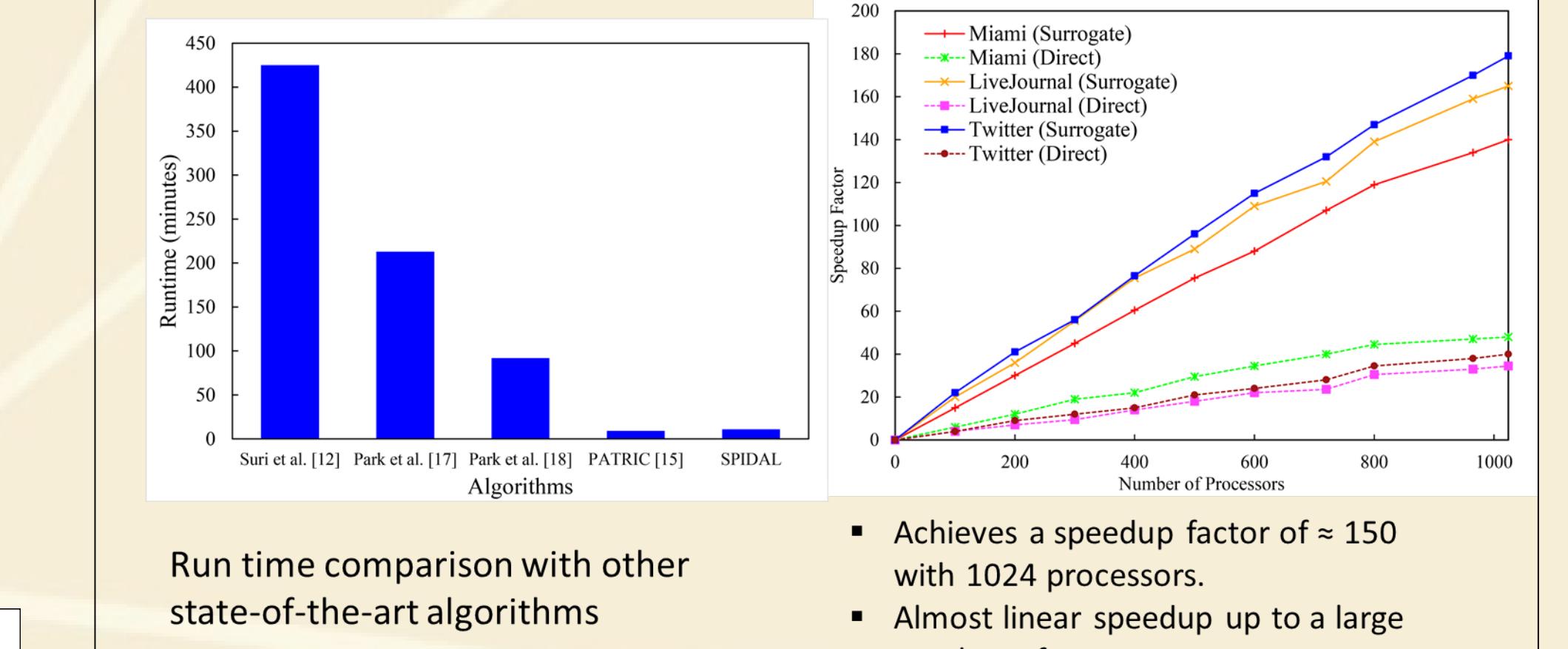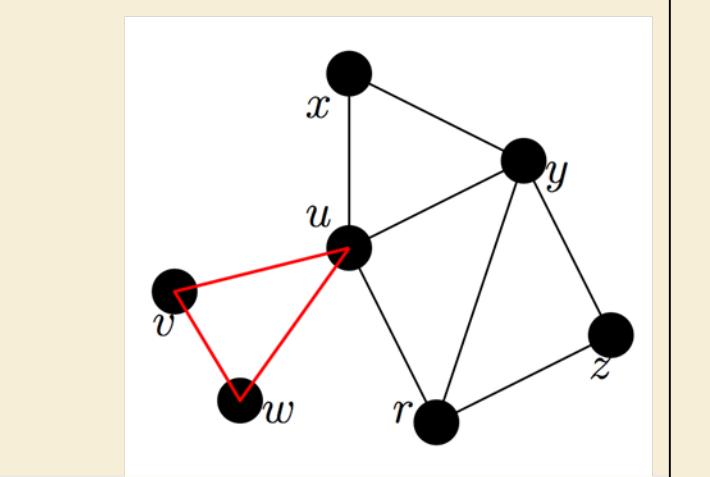
We have successfully applied SPIDAL optimization algorithms to 3D radar data collected by NASA Operacon IceBridge over the Canadian Arctic Archipelago ice caps. The basal topography underlying these ice caps is unknown for many of the glaciers that drain them. The high volume of data produced by 3D imaging renders manual tracking of the ice bottom impractical. To solve this problem, we used SPIDAL algorithm to extract ice bottom surfaces using a probabilistic graphical model. We first estimate layer boundaries to generate a seed surface, and then incorporate additional evidence, such as ice masks, surface elevation models, and feedback from users, to refine the surface in a discrete energy minimization formulation. The sequence of figures below shows the user interface for examining the optimization algorithm's inputs and outputs. A sample result of an ice bottom digital elevation model is shown to the right, overlaid on a Landsat-8 image.



GUI to browse ice surface and bottom layer results, enter ground truth, and rerun optimization.

Map view showing projection of 3D image slice on a geographic map.

Ice bottom digital elevation model extracted using optimization algorithm on 3D data.

Ice mask (rock/ ice ice) that informs algorithm where layers should merge.

### Counting Triangles in Massive Networks

Many applications in data mining, network analysis, social science, and database systems



Run time comparison with other state-of-the-art algorithms

- Achieves a speedup factor of ~150 with 1024 processors.
- Almost linear speedup up to a large number of processors

### WebPlotViz – Browser Visualization of High Dimensional Data

WebPlotViz is a 2D/3D data point browser that can visualize very large volumes of 2D or 3D data, as points in a virtual space and enable users to explore the virtual space interactively. WebPlotViz also includes support for Time Series Data plots



446K sequences and ~100 clusters visualized in WebPlotViz