



# CIF21 DIBBs: Middleware and High Performance Analytics Libraries for Scalable Data Science

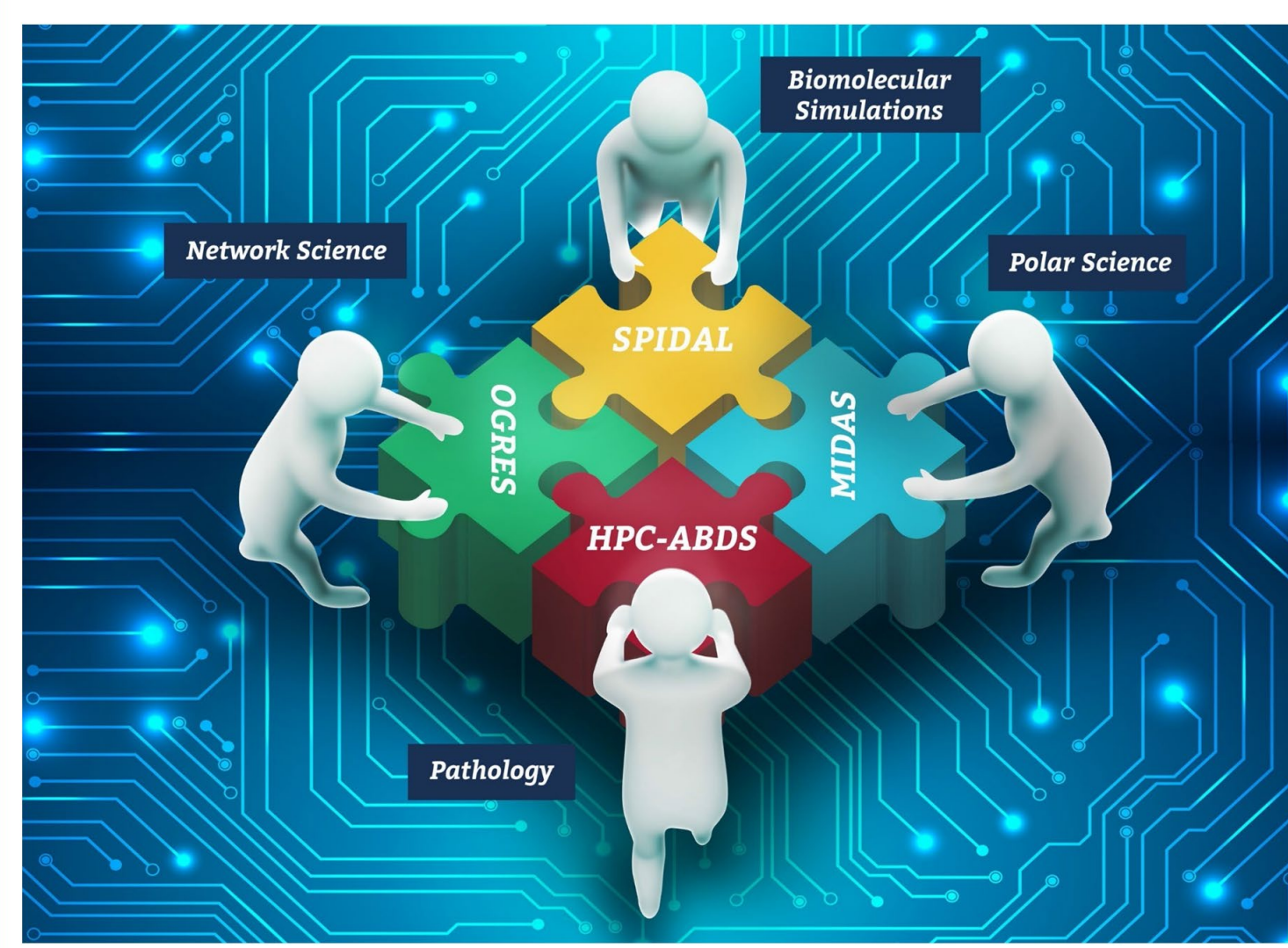
PI: G. Fox, Co-PIs: Madhav Marathe, Shantenu Jha, Judy Qiu, Fusheng Wang

Institutions: Arizona State, Indiana (lead), Kansas, Rutgers, Stony Brook, Virginia, and Utah.

Award #: 1443054

## MLforHPC and HPCforML

- Project mainly addresses **HPCforML** enabling high performance big data
- MLaroundHPC** broadly applicable but current use is **nonuniform** across domains
  - we need to improve Cyberinfrastructure to make MLforHPC more effective for **more users**
- Use of modest DL network to **map material/potential drug structure to properties** (generalized QSAR) with simulation and observation: Advanced Progress
- Learn **surrogates** for **large scale simulations**: initial small scale results
- Use of MLforHPC in **agent-based systems** (learn agents): Very promising but few results
- Macroscopic** Structure as in learn complex multi-particle **potentials** scaling to  $N^7$ : many great successes
- Learn **Collective coordinates and guide ensemble** computations: dramatic progress with speedups up to  $10^8$
- Microscale**; learn dynamics of small scale such as clouds, turbulence: Interesting results but much more to do
- Use of **Recurrent NN's** to represent dynamics (**learn numerical differential operators**): Promising but only studied in small problems



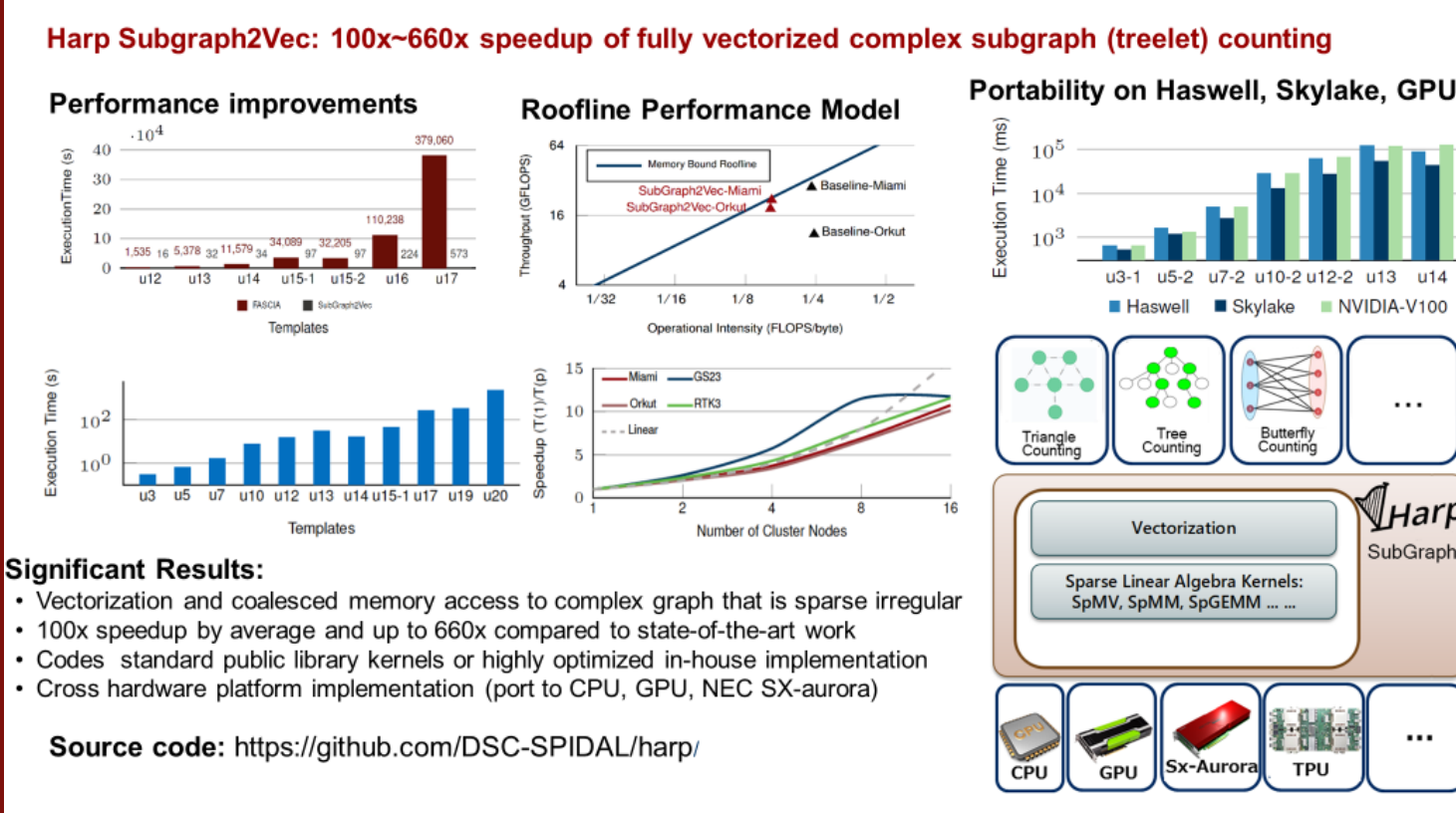
## Community Driven High-Performance Big Data for bio-physical applications based on HPC, distributed systems, network science, GIS and machine/deep learning

### Completed Activities

- MIDAS Pilot Jobs HPC Task Management
- High Performance for Java/C++ for Machine Learning
- NIST Big Data Application Analysis: Features of data intensive Applications deriving 64 Convergence Diamonds
- HPC-ABDS: Cloud-HPC Interoperable software with performance of HPC and rich functionality of commodity Apache Stack
- Implementation of HPC and Clouds with DevOps
- Image Processing and Machine Learning Toolkit <http://hpcanalytics.org>

## HPC Big Data Cloud Convergence

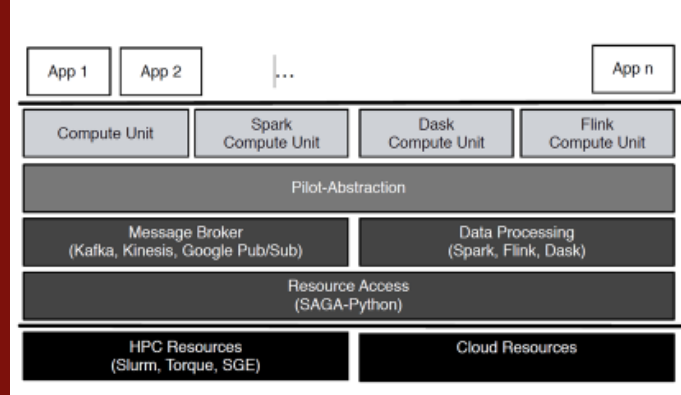
### Harp-DAAL: A HPC-Cloud convergence Data Analytics framework



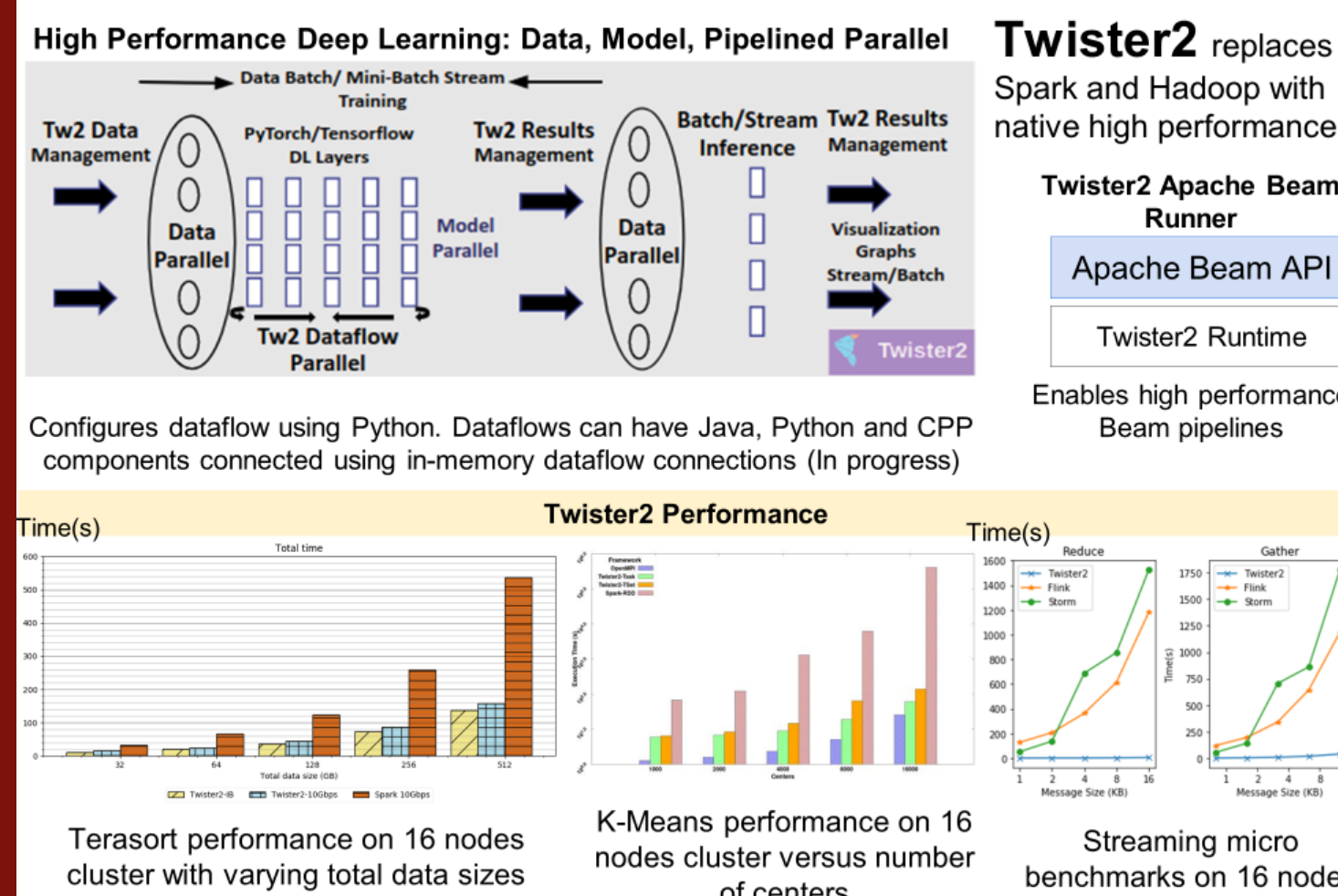
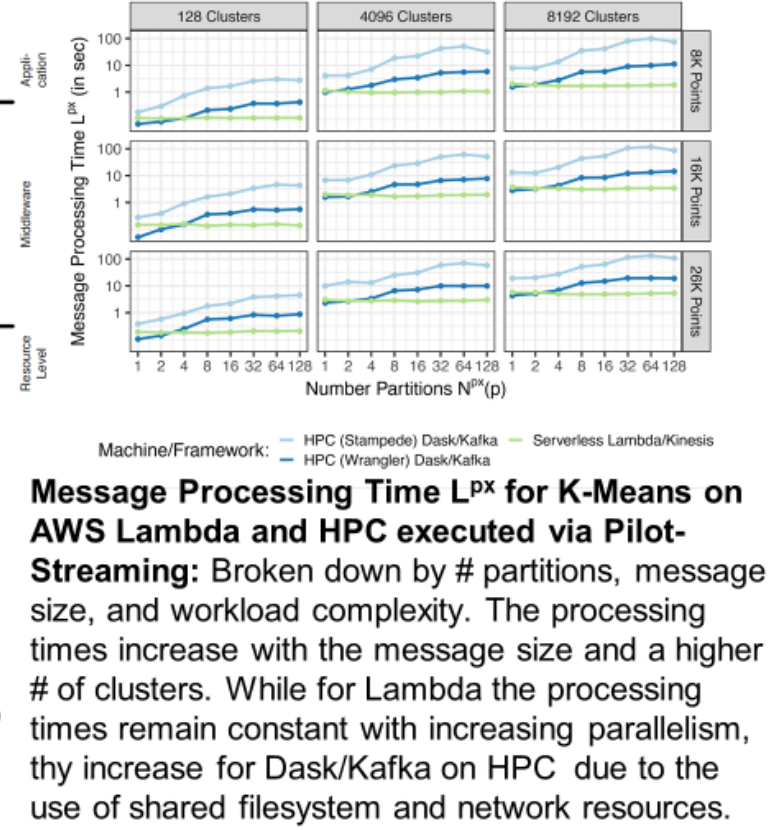
## Streaming Data Systems

### Pilot-Streaming (deployed)

Pilot-Streaming provides a unified abstraction for resource management for HPC, cloud, and serverless, and allocates resource containers independent of the application workload removing the need to write resource-specific code.



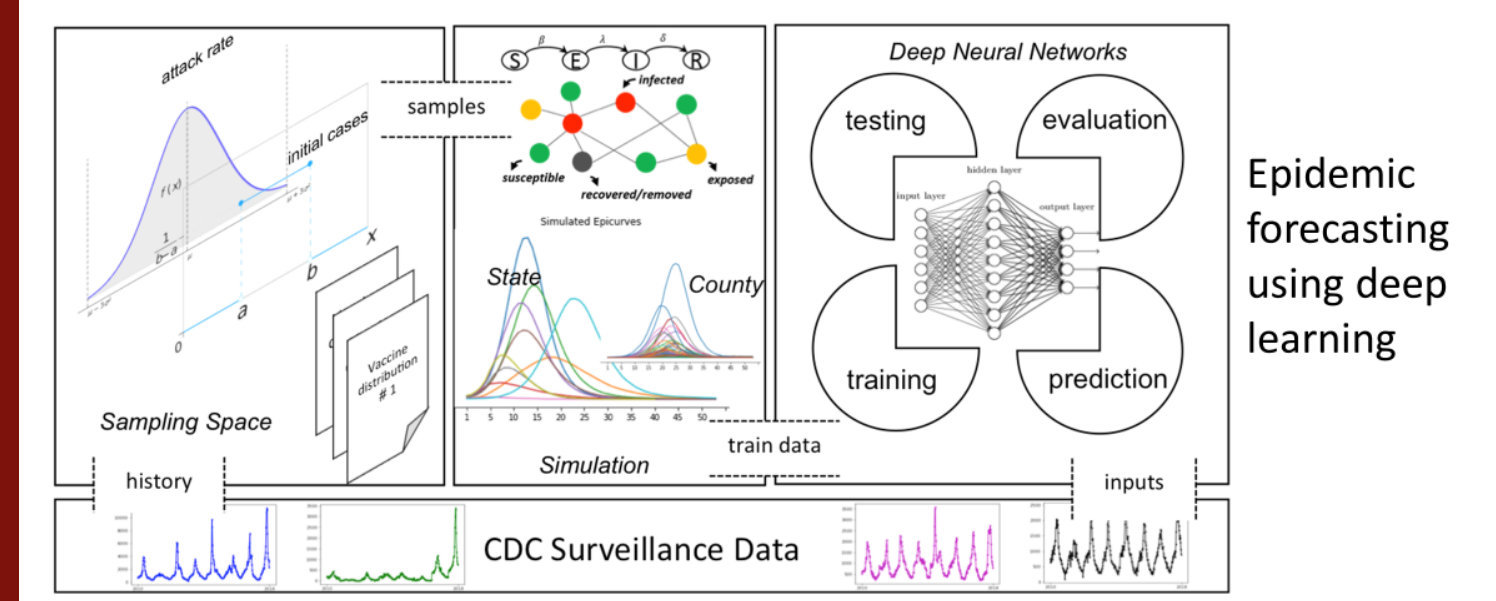
**Pilot-Streaming on Serverless High-Level Architecture and Interactions:** Pilot-Streaming enables the unified management of batch and streaming compute tasks on serverless infrastructure. It orchestrates the setup of Kinesis and Lambda and manages the streaming workload.



## Network Science Community

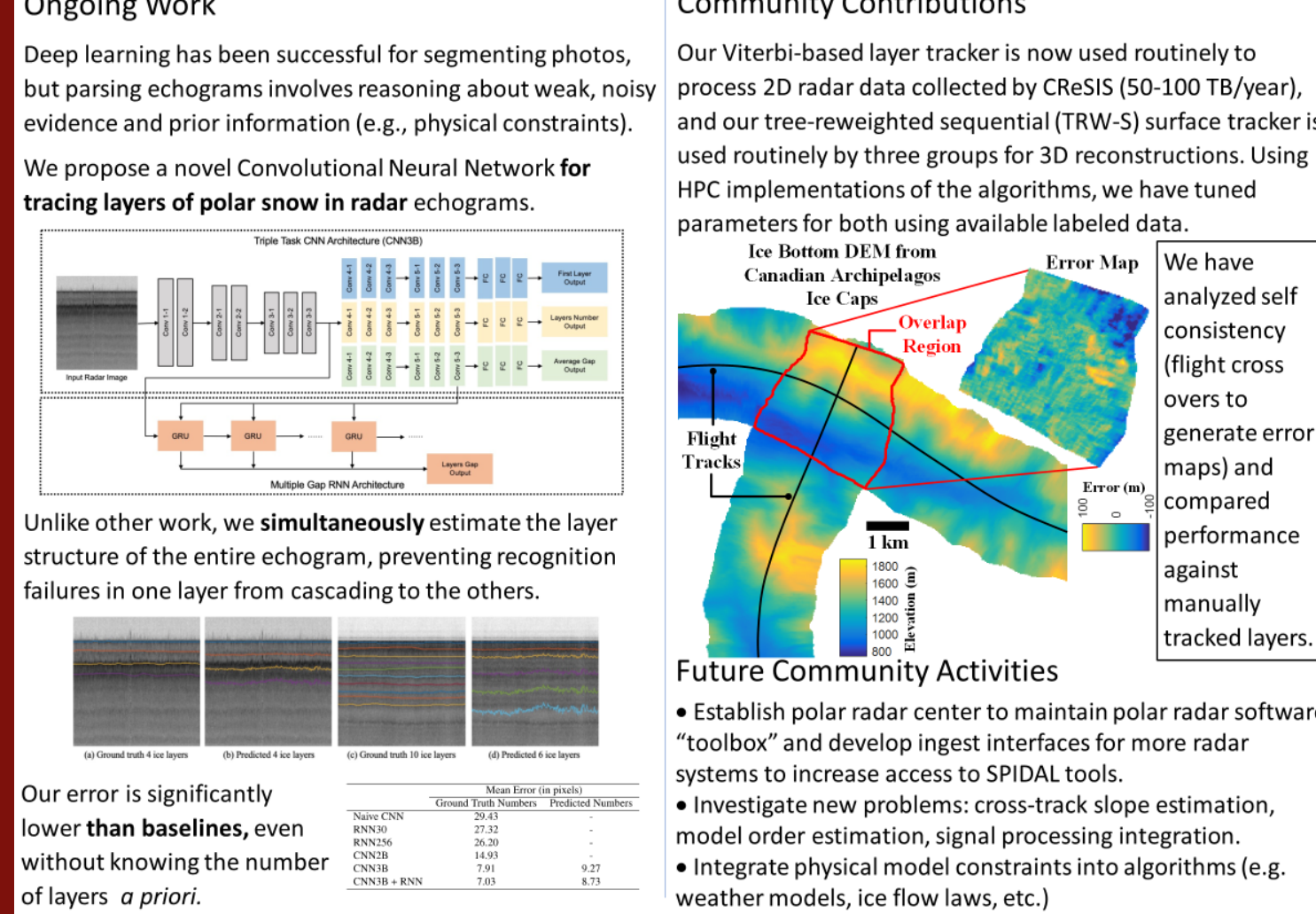
### Network Analytics

- Substantial Community Contributions of Project:** New advances in sequential and parallel algorithms for a broad class of network problems
  - Subgraph detection and counting, anomaly detection, dense subgraph and community detection
  - Parallel network generation: generating very large instances from many random graph models and by parallel edge switching, with different kinds of constraints
  - Applications: public health, social network analysis, scaling studies of network algorithms, sensitivity analysis, synthetic population generation
- Ongoing Project work:** Algorithms for dense subgraphs, communities in dynamic and temporal networks, and parallel network generation
- Future Community Activities:** Integration of algorithms within CINES cyberinfrastructure



## Polar Science Community

### Structure of polar ice sheets from radar informatics



## Health Science Community

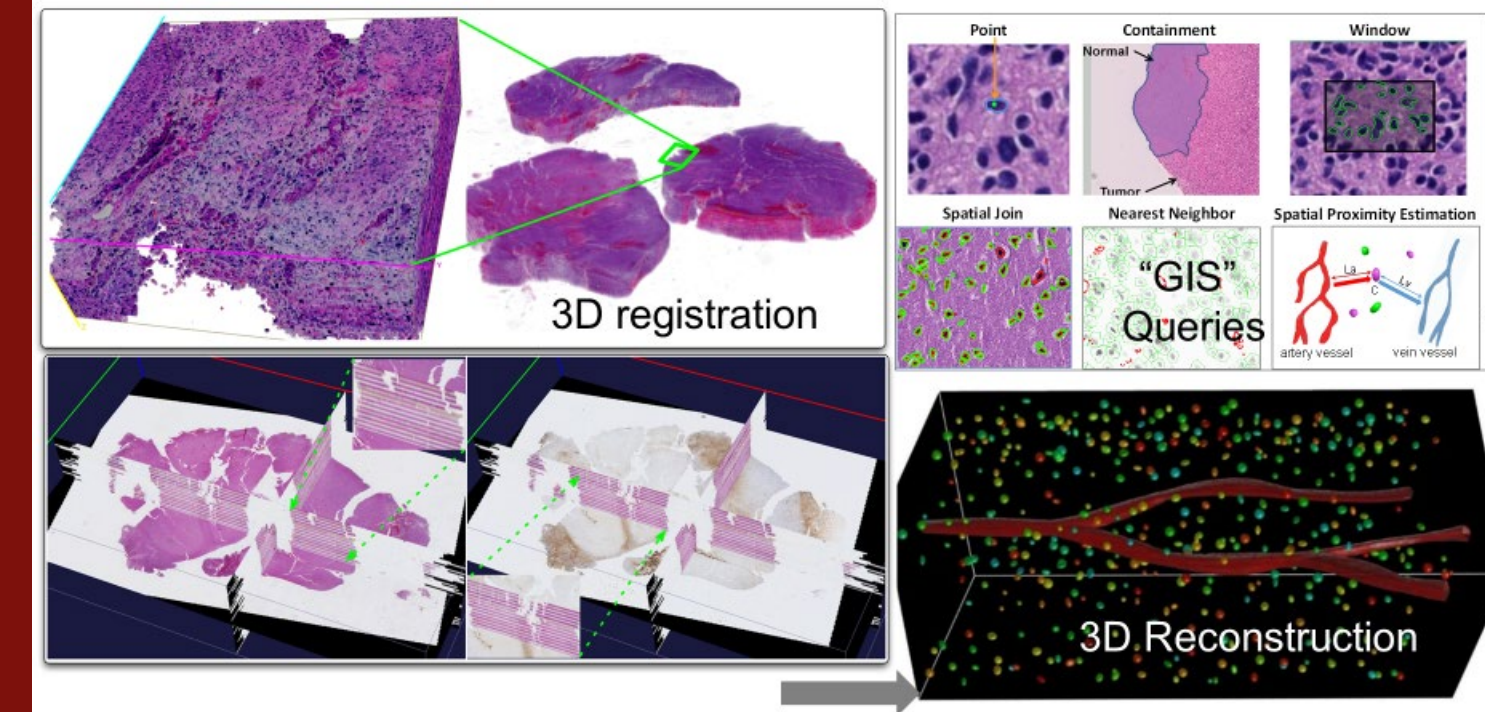
### Spatial Big Data Query Systems

- Managing and querying big spatial data is challenged by explosion of data, multi-dimensions and the complexity of geometric computation
- Developed scalable and efficient spatial big data querying systems running on big data platforms, for both 2D and 3D
  - Hadoop-GIS, SparkGIS and iSPEED (3D)
- Integrative CPU-GPU based high performance 3D spatial queries
- The platforms can be used to support geospatial applications, big data driven public health studies, and digital pathology
- Future:** Spatial analysis for understanding resource availability for opioid epidemic prevention and intervention



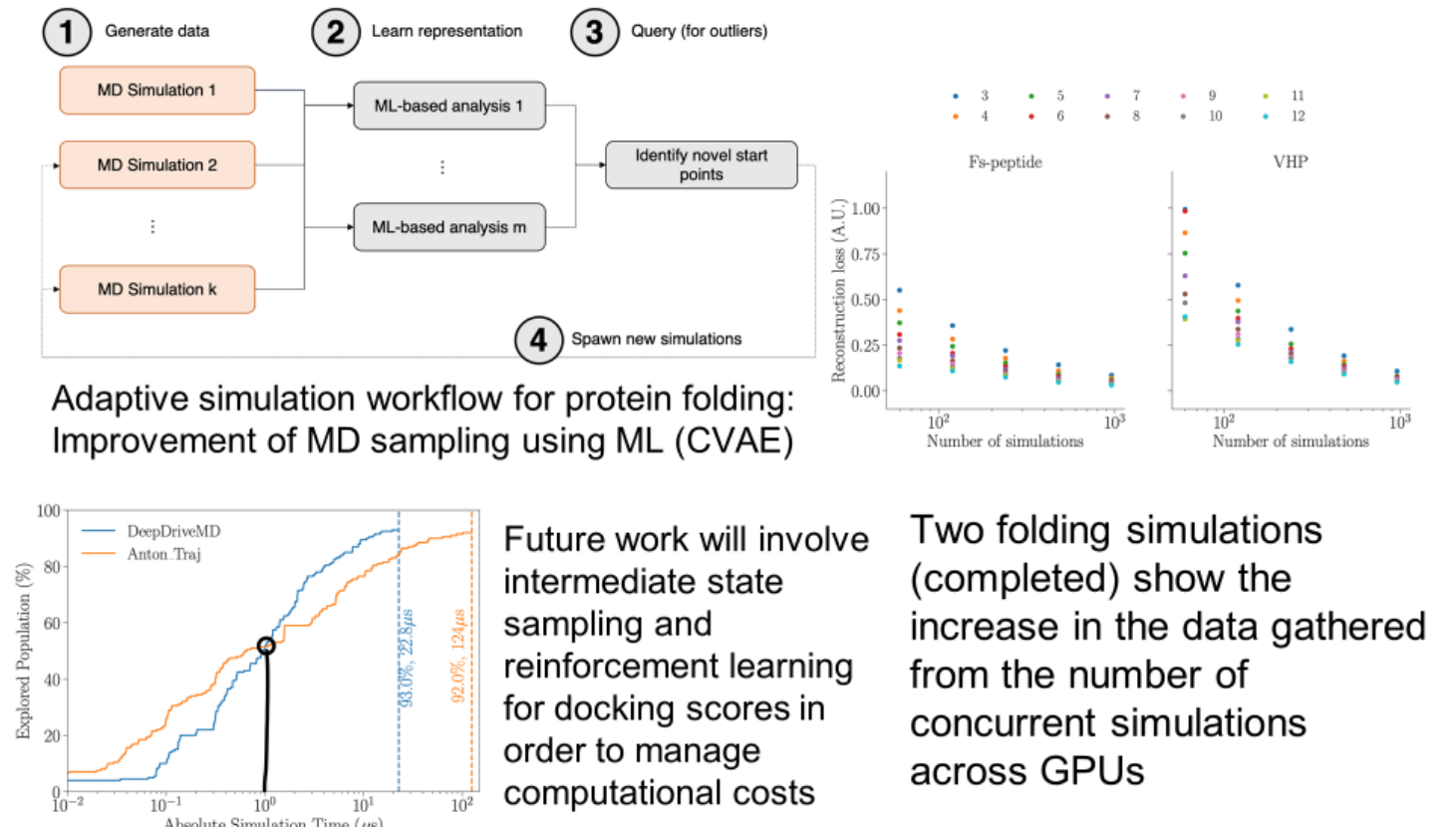
### Digital pathology image analysis

- Created a suite of deep-learning based image analysis tools for level-set based image segmentation, 3D registration, reconstruction, and spatial analysis
- Future: Data analysis for understanding breast cancer treatment and prognosis using 3D digital pathology



## Biomolecular Science Simulation Community

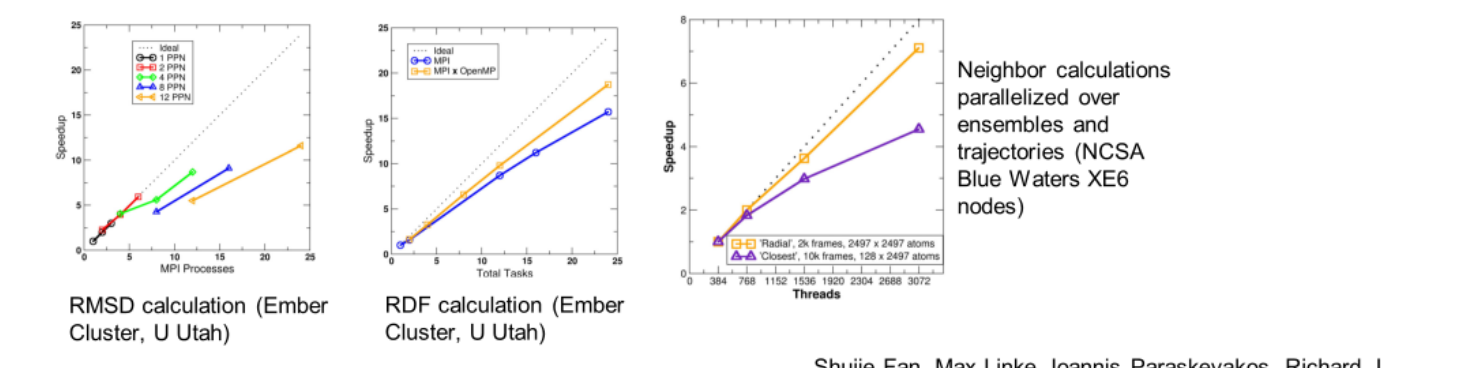
### DeepDriveMD: AI driven biomolecular simulations on HPC



### SUBSTANTIAL COMMUNITY CONTRIBUTIONS OF PROJECT

#### Released software

**cpptraj:** MD trajectory analysis tool (MPI and OpenMP parallelization) <https://github.com/Amber-MD/cpptraj>



**PMDA:** parallel molecular dynamics analysis library (task-parallelization) <https://www.mdanalysis.org/pmda/>

**alchemy:** general library for processing and analyzing free energy data from MD and MC simulations <https://github.com/alchemy/alchemy>

### Parallel Molecular Dynamics Analysis (PMDA)

