

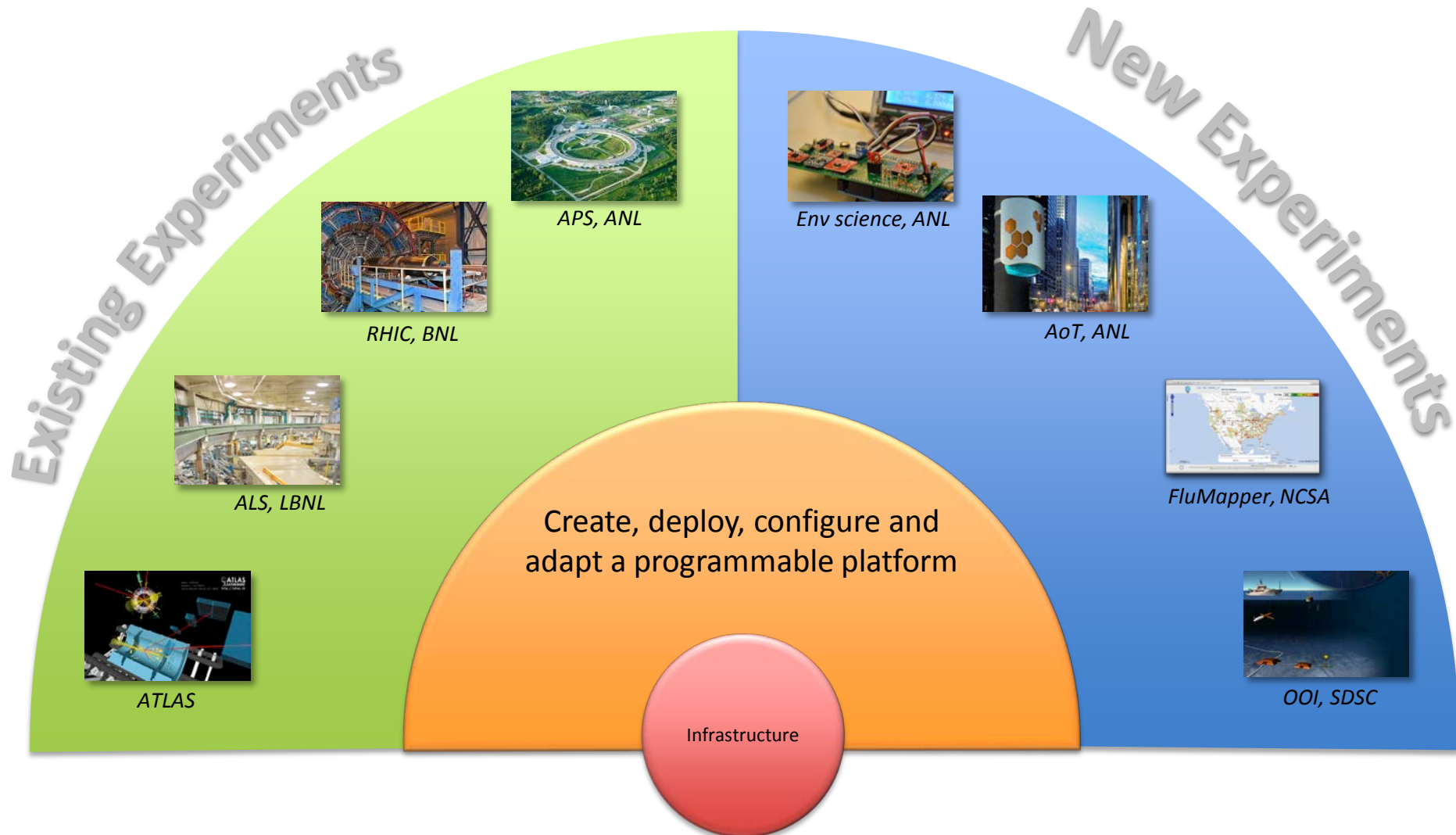


# IRMO: Programmable Platforms for Scientific Applications

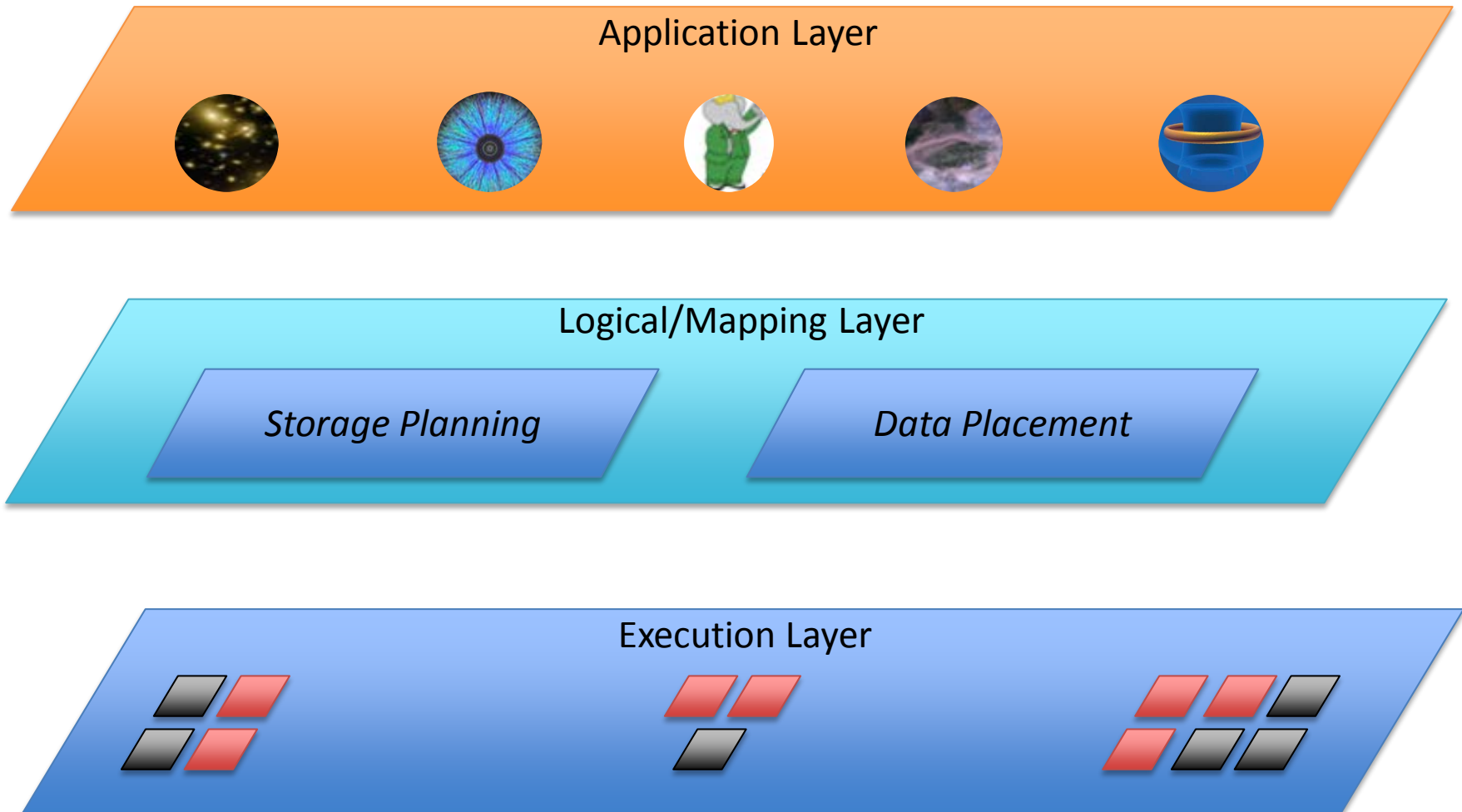
*Kate Keahey, Lavanya Ramakrishnan*

*<http://press3.mcs.anl.gov/irmo/>*

# Emergent Needs for Infrastructure Strategy

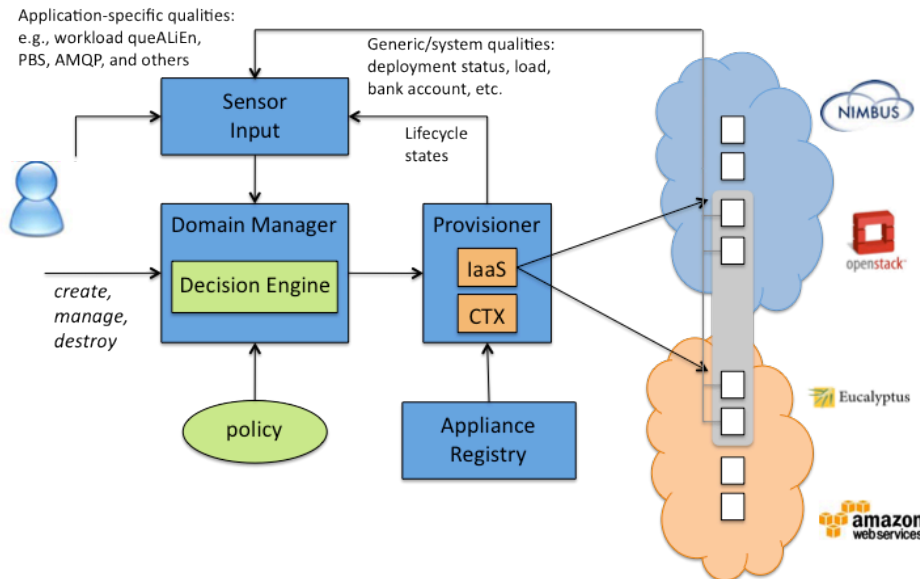


# Approach

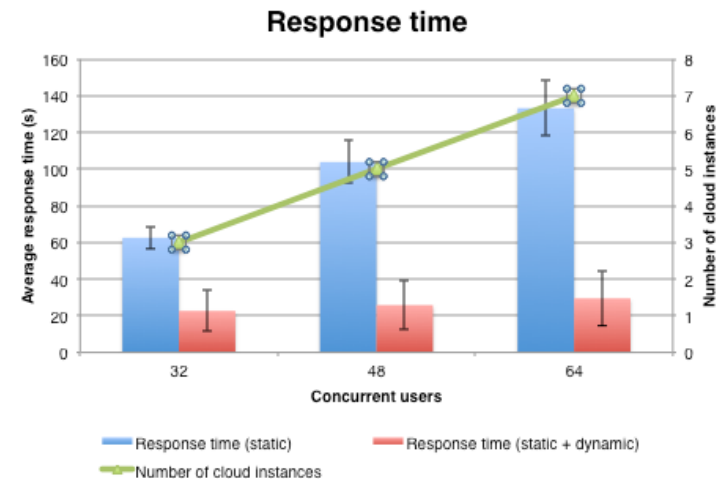


# Computation to Fit

Compute environment that automatically scales to fit the evolving need of application or community over a federation of resources



*Scaling based on system and application factors*



*Response time of a CyberGIS application  
as a result of scaling (ScienceCloud 2014)*

*Paper: "Infrastructure Outsourcing in Multi-Cloud Environment", Cloud Services and Federation 2012*

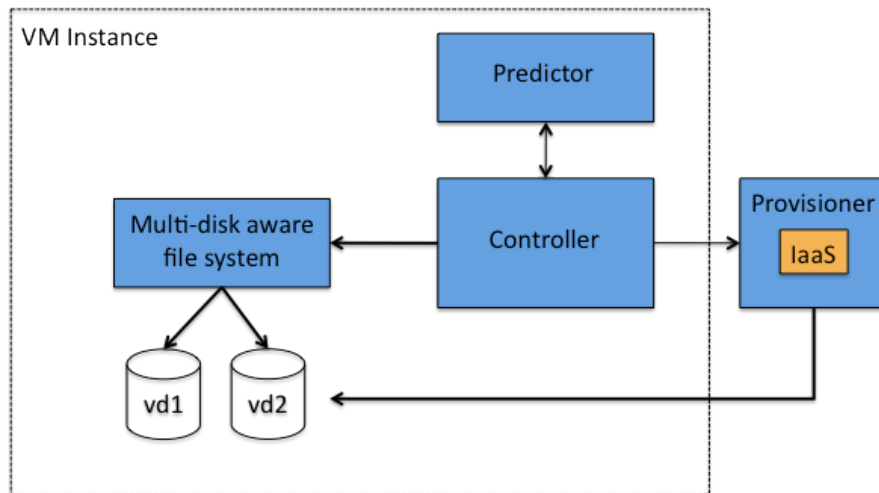
*Paper: "Rebalancing in a Multi-Cloud Environment", ScienceCloud 2013*

*Paper: "A Cloud Computing Approach to On-Demand and Scalable CyberGIS Analytics", ScienceCloud 2014*

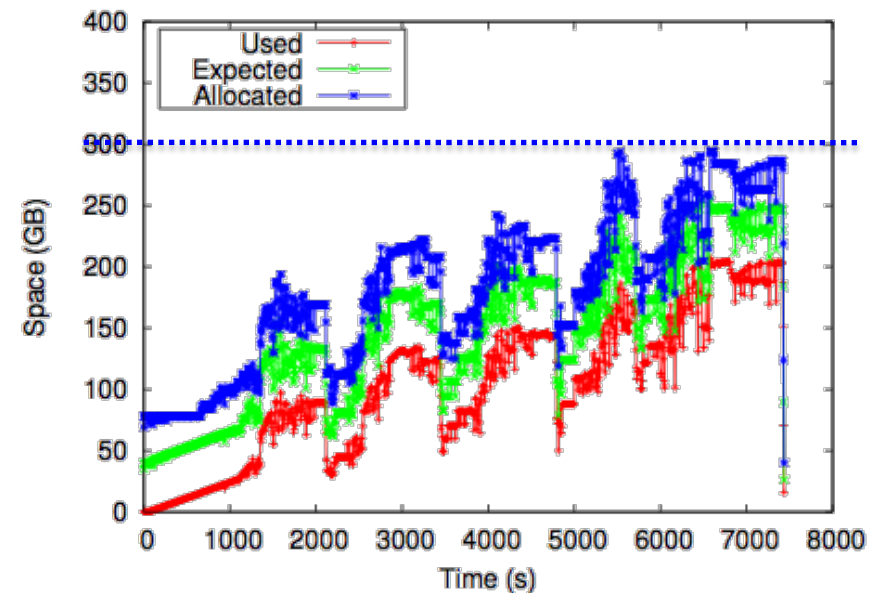
*... and others at [www.nimbusproject.org](http://www.nimbusproject.org)*

# Storage to Fit

Storage that automatically scales to fit the evolving application needs both in terms of size and type (performance).



*Adaptive storage scaling system*

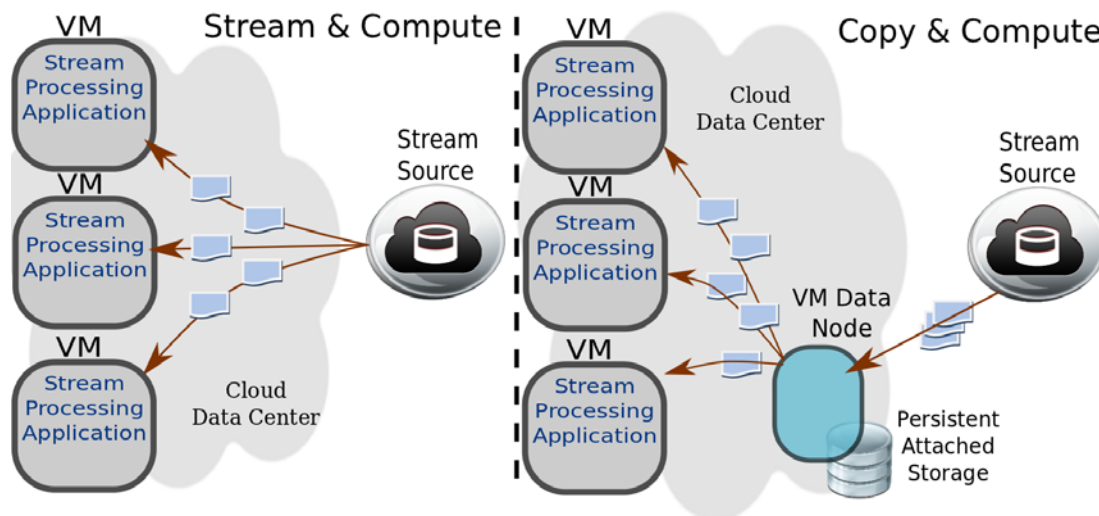


*Predictive storage scaling for K-means (IPDPS 2014)*

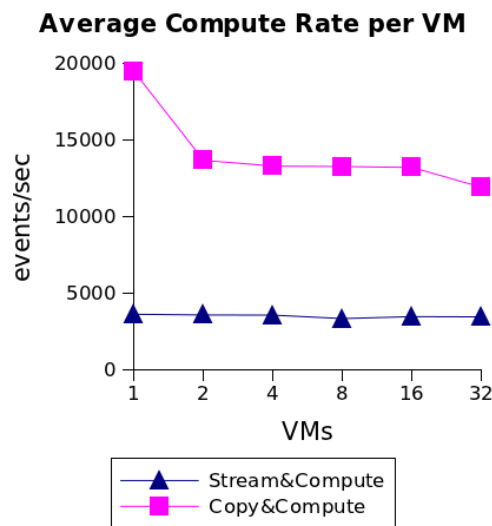
*Paper: Bursting the Cloud Data Bubble: Towards Transparent Storage Elasticity in IaaS Clouds, IPDPS 2014*  
*Paper: Transparent Throughput Elasticity for Cloud Storage using Guest-side Block-level Caching, UCC 2014*

# Network to Fit

Network that tells me how far it can scale to feed the resources I acquire and how to best use them.



*Two streaming strategies for transfer to the cloud*



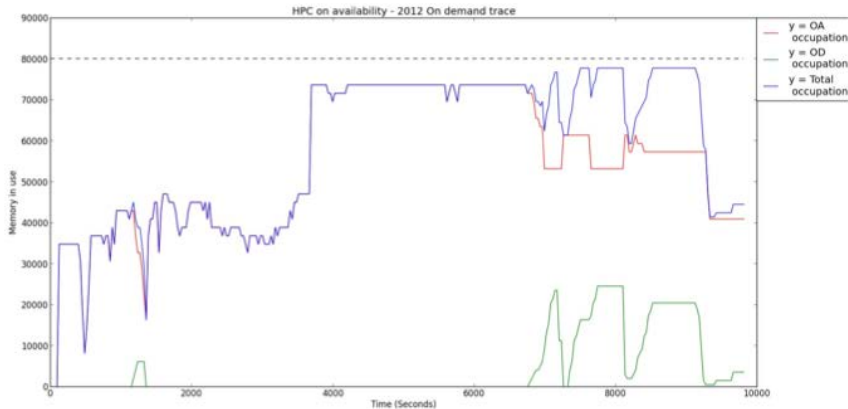
*Comparison of compute rates resulting from different streaming scenarios*

Collaboration with "Next Generation Workload and Analysis System for Big Data"

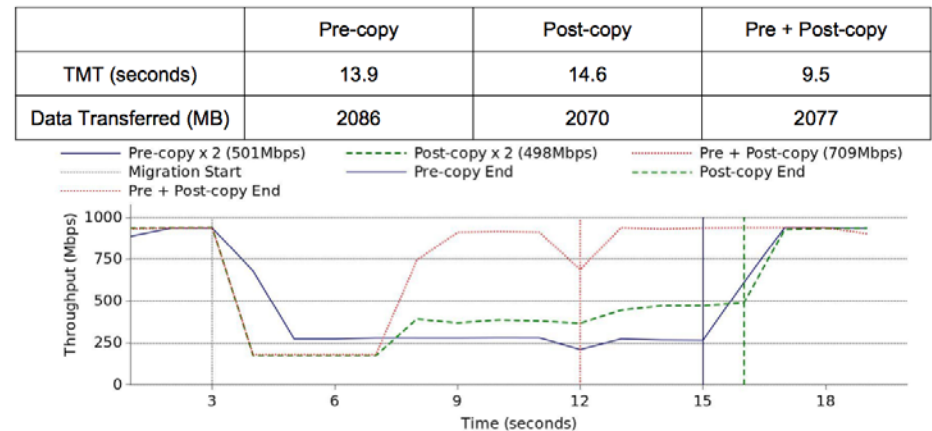
*Paper: Evaluating Streaming Strategies for Event Processing across Infrastructure Clouds, CCGrid 2014*

# A Fitting Data Center

- How does a scientific data center need to be organized to support such programmable resources?
- Defining and interleaving different types of leases to optimize both provider's and user's objectives
- Exploring efficient techniques to implement lease semantics

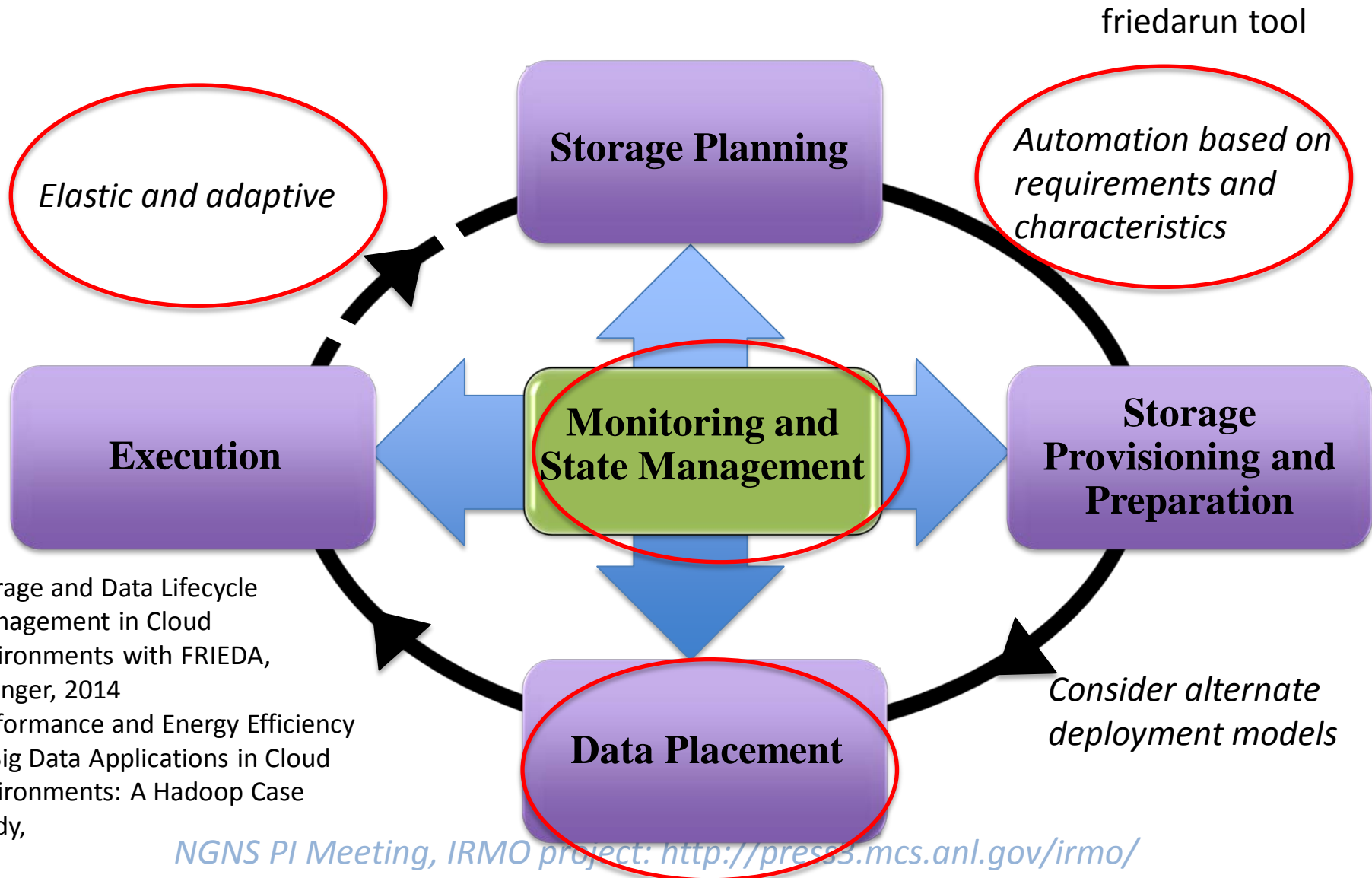


*Defining and interleaving various types of leases*



*Exploring efficient techniques to implement lease semantics*

# FRIEDA – Focus This Year



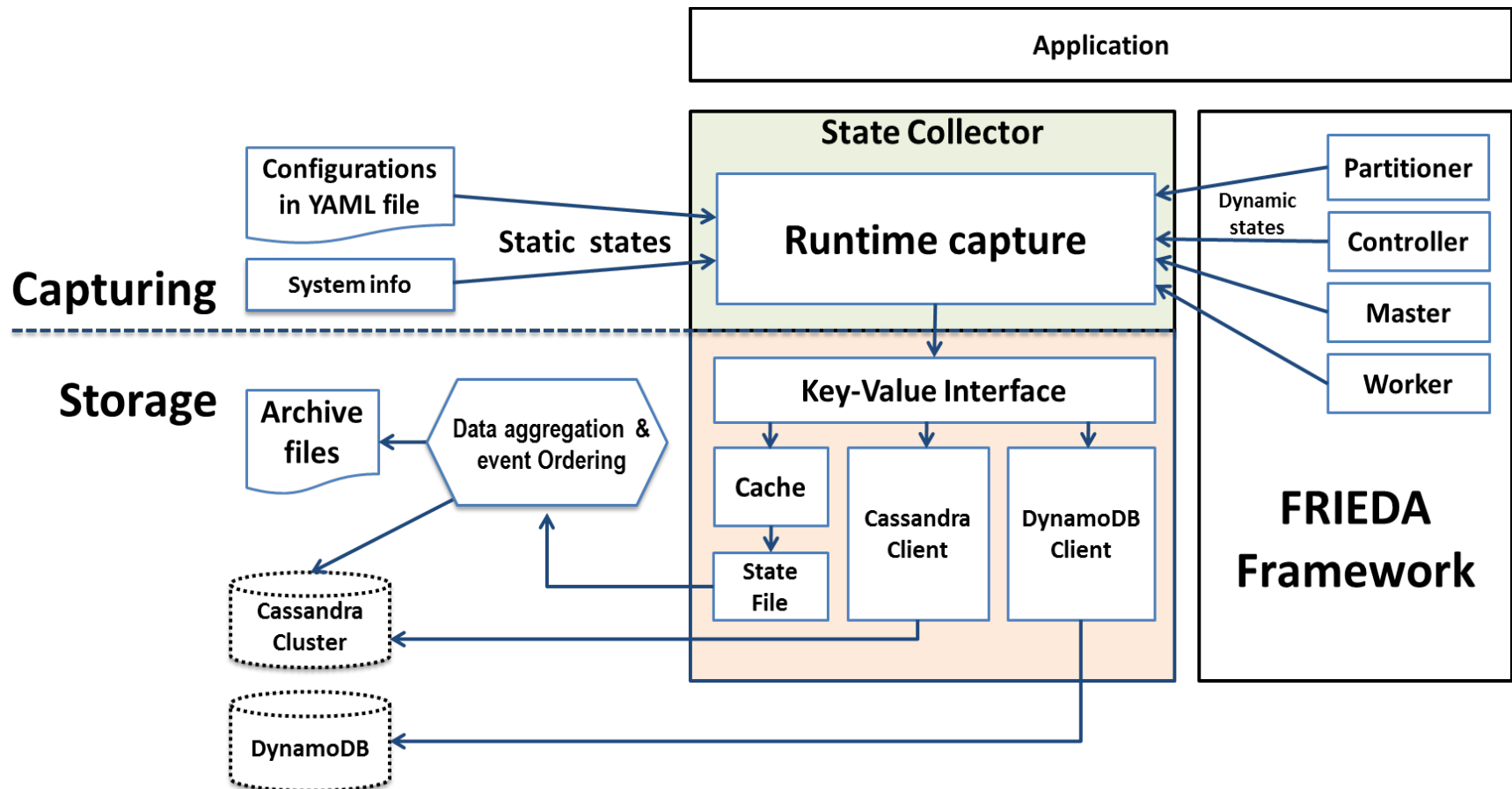


# Collaboration with ATLAS

- Collaboration with LBNL team (Beate Heinemann, Mike Hance and Sourabh Dube )
  - Science: Why is the Higgs mass so low?
  - Using FRIEDA to manage their computation and data on Amazon Web Services Grant

# FRIEDA-State: Provenance tracking

How do we capture state that allows reconstruction of events from transient resources?  
Challenges – event sequencing and data collection

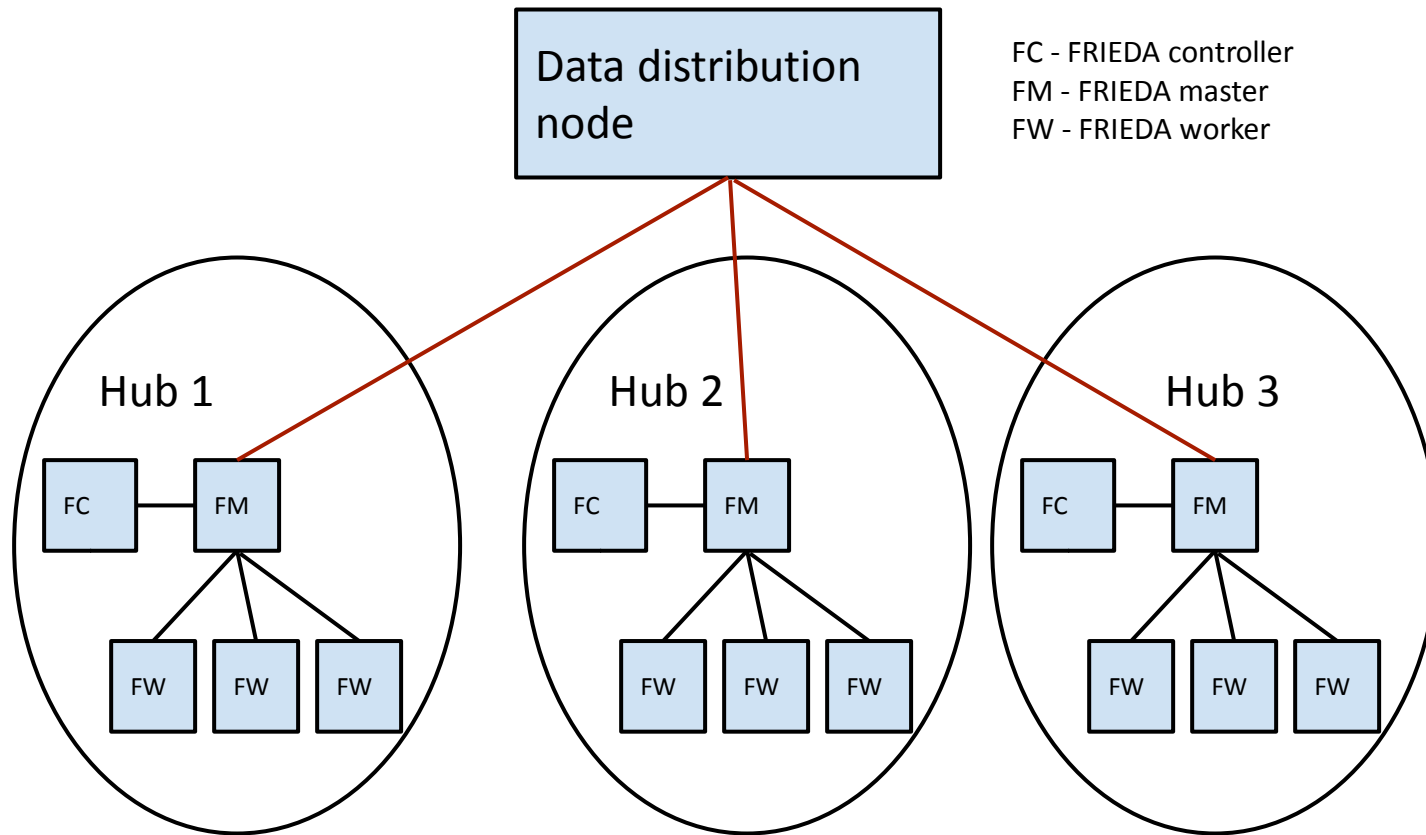


Scalable State Management for Scientific Applications in the Cloud, IEEE Big Data, 2014

NGNS PI Meeting, IRMO project: <http://press3.mcs.anl.gov/irmo/>

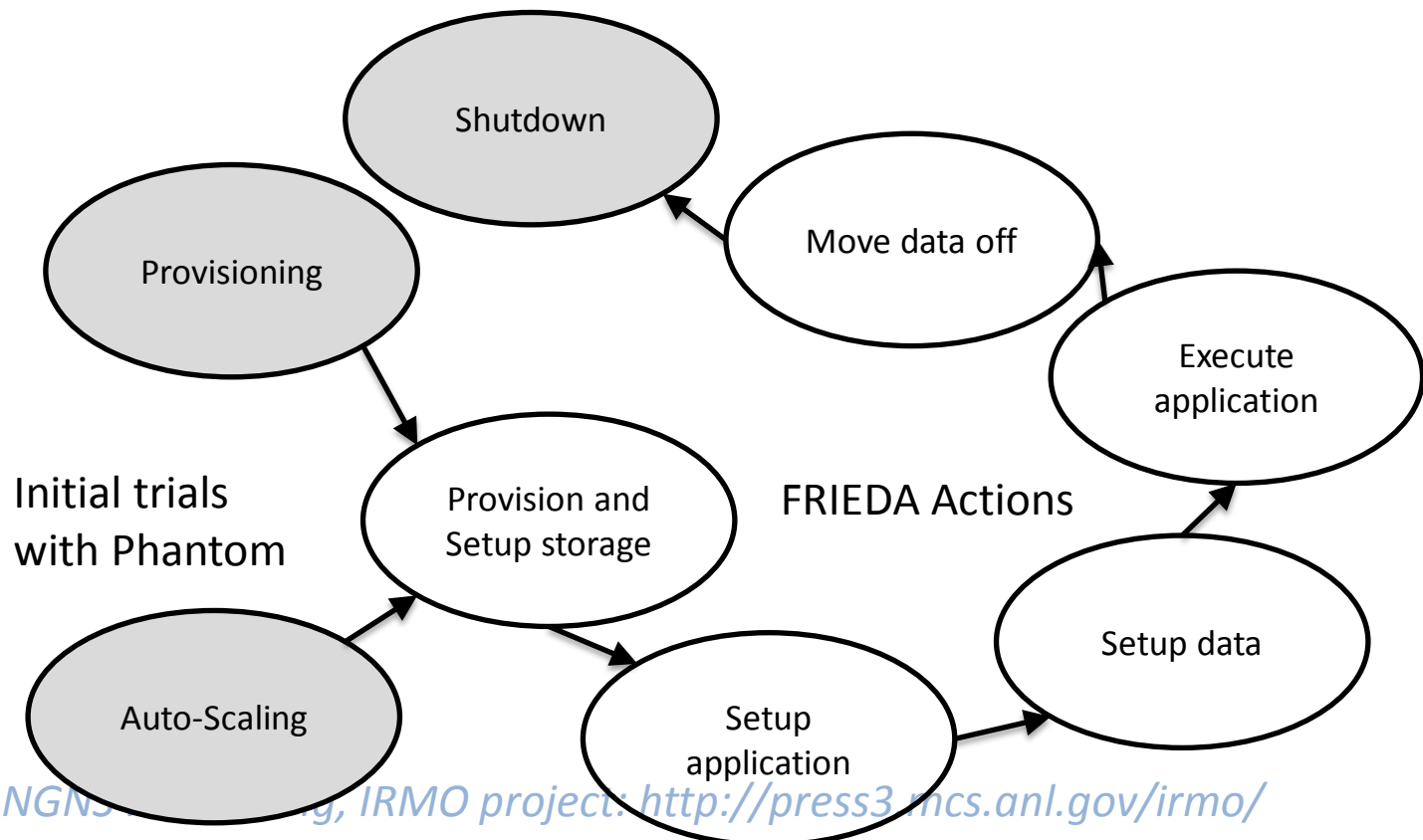
# Multi-Site Data Placement

How do we manage data across multiple sites (e.g. across experiment site and computational facility)?



# Scaling and Storage and Data Management

- **External:** Impact of data arrival rates on storage provisioning mechanisms
- **Application Execution:** correlation of I/O load with scaling and data management strategies
- **Auto-scaling:** impact on storage and data-management



# Future Work

- Programmable platforms
  - To what extent can we define them? What technology is missing? How will the underlying infrastructure have to change? What information/models need to be exposed to the user? How do we do it efficiently?
- Dynamic shaping/scaling
  - What are the best methods to scale/adapt programmable platforms automatically?
- Sharing and incentives
  - How do I need to organize/manage resources to support efficient sharing? What cost models are best/fair? Sharing vs redundancy vs energy management?