

4CeeD: Real-Time Operating Infrastructure for Capturing, Curating, Correlating, Coordinating and Distributing Materials- related Data

Klara Nahrstedt (klara@illinois.edu)

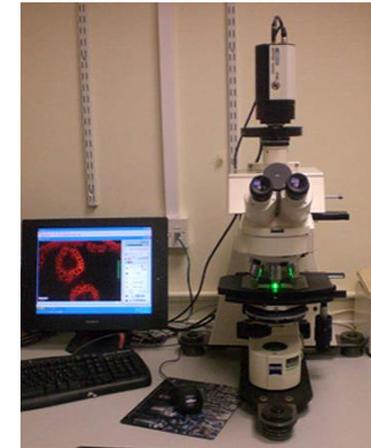
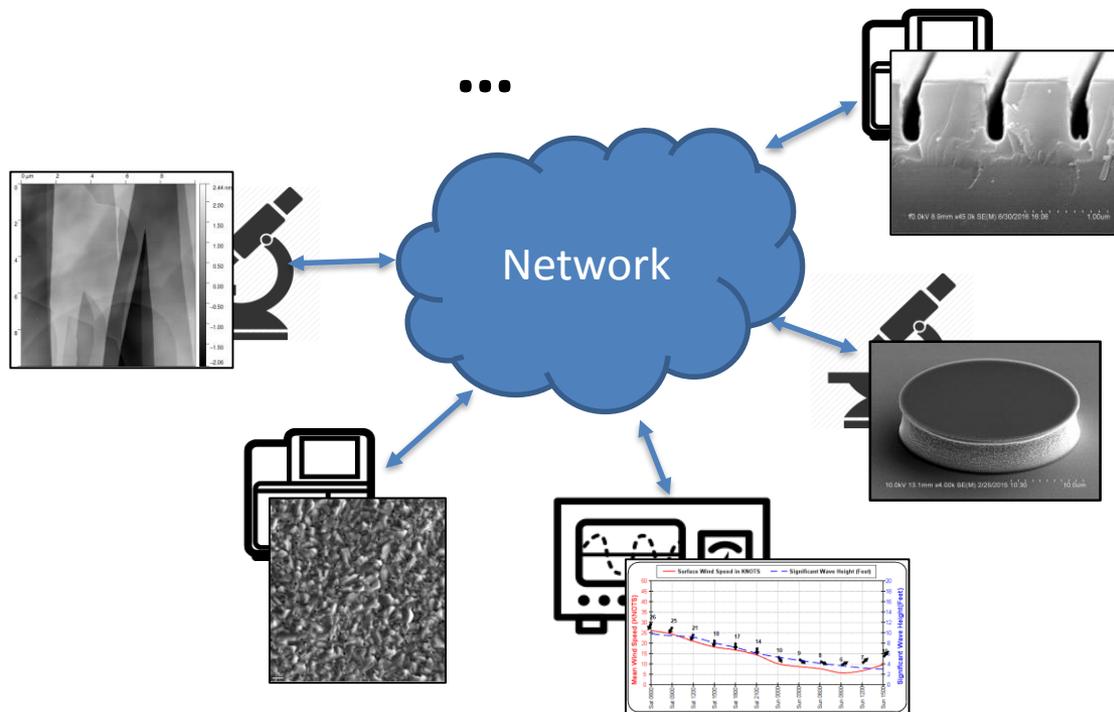
Coordinated Science Laboratory

University of Illinois at Urbana-Champaign



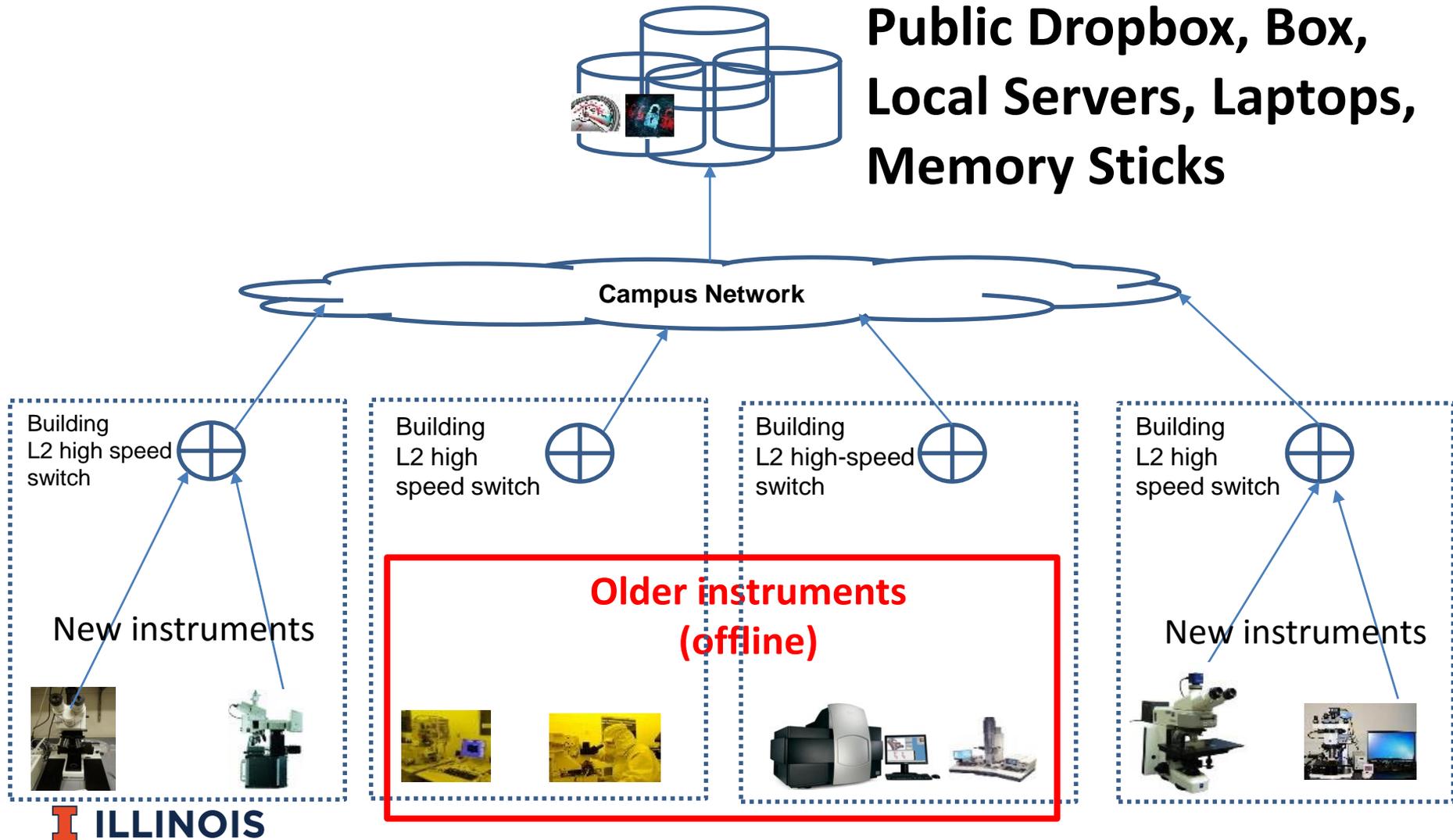
Background: Increasingly data-driven and interdisciplinary scientific research

- *Key enabling factor*: Network connected scientific instruments capable of real-time data capture

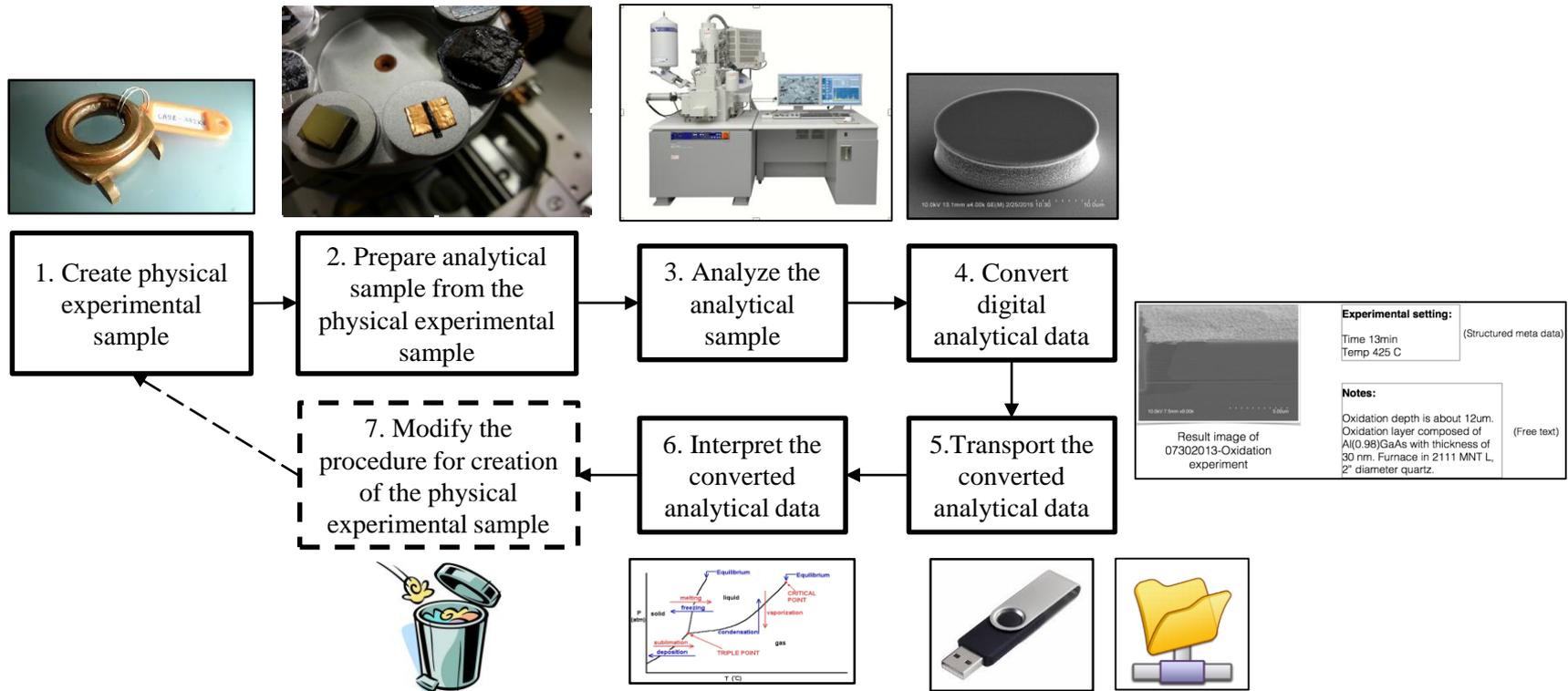


Digital microscope

Current situation in campus cyberinfrastructure



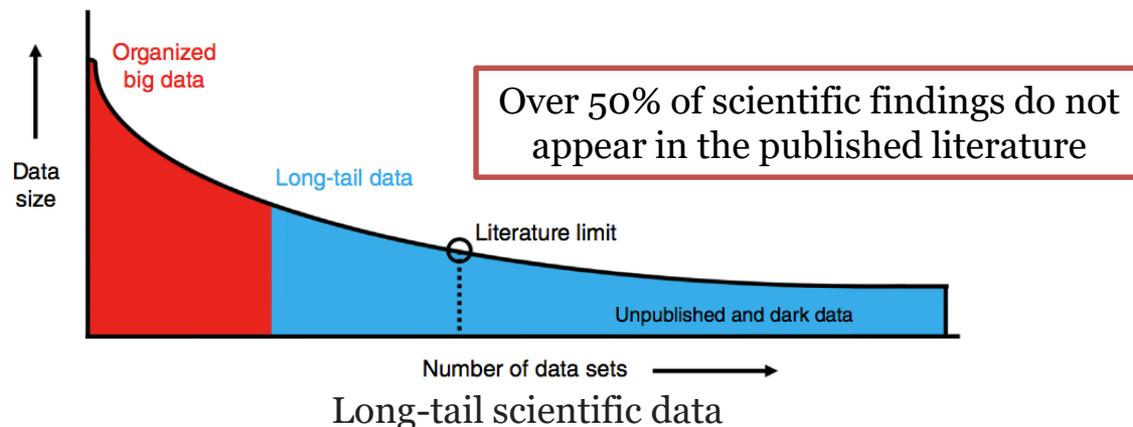
Example: Typical experimental process in material science research



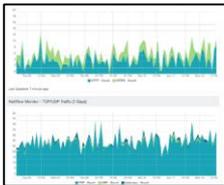
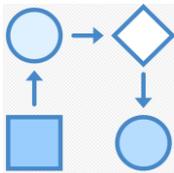
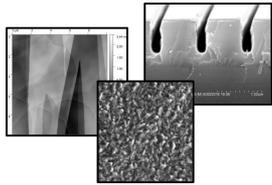
It typically takes 20 years to go from the discovery of new materials to fabrication of new and next-generation devices*

Motivation: Needs for advanced cyber-infrastructure for long-tail scientific data

- Related efforts mainly focus on **homogenous, well-organized data** in an offline or batch manner
- Much less effort has been on **long-tail scientific data**:
 - Small/medium sized data sets collected during day-to-day research
 - “Dark data”, e.g., unpublished data of failed experiments



Challenges in General



➤ **Heterogeneous** scientific data management and processing

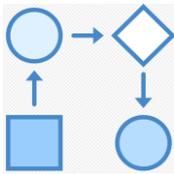
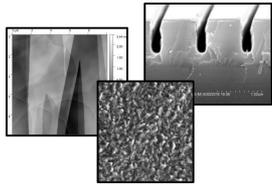
➤ Support ad hoc and complex data analysis **workflows**

➤ Shorten **time** from digital capture to interpretation & insights

➤ **Real-time data capture and acquisition**

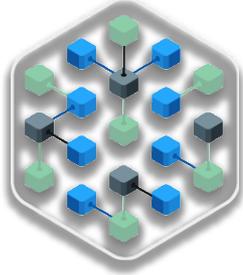
➤ **Analytics** support to gain insights from data

Challenges on Campuses



- Very **diverse scientific instruments** in Materials Research Lab (MRL), Micro-Nano-Technology Lab (MNTL), other labs
- Support very **different user groups** that collect and analyze data
- **Relations between students, faculty staff** and academic cycles are **different** than in industry, impacting how insights are gained
- Rules on campuses regarding **secure access to data and metadata** in scientific labs vary
- **Analytics tools support** in scientific campus labs vary

Our approach



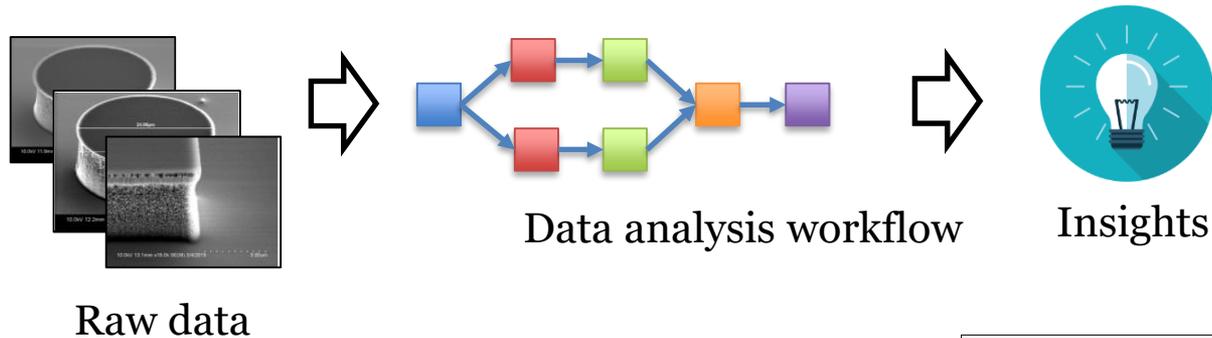
- ✓ Micro-service private cloud execution environment for instrument data curation and coordination (4CeeD)



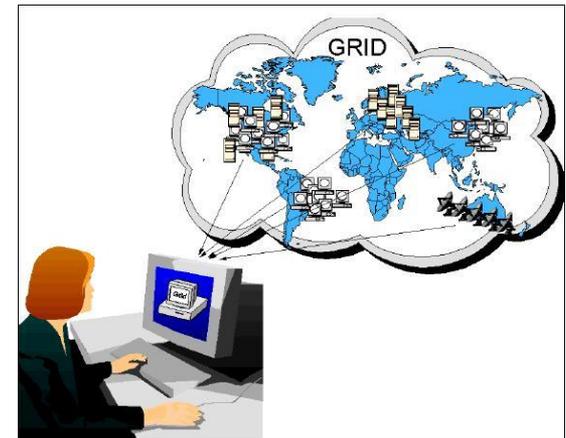
- ✓ Data acquisition from aging instruments (BRACELET)

Long-tail scientific data processing challenges

- **Challenges:** Support execution of heterogeneous types of data processing & analysis workflows

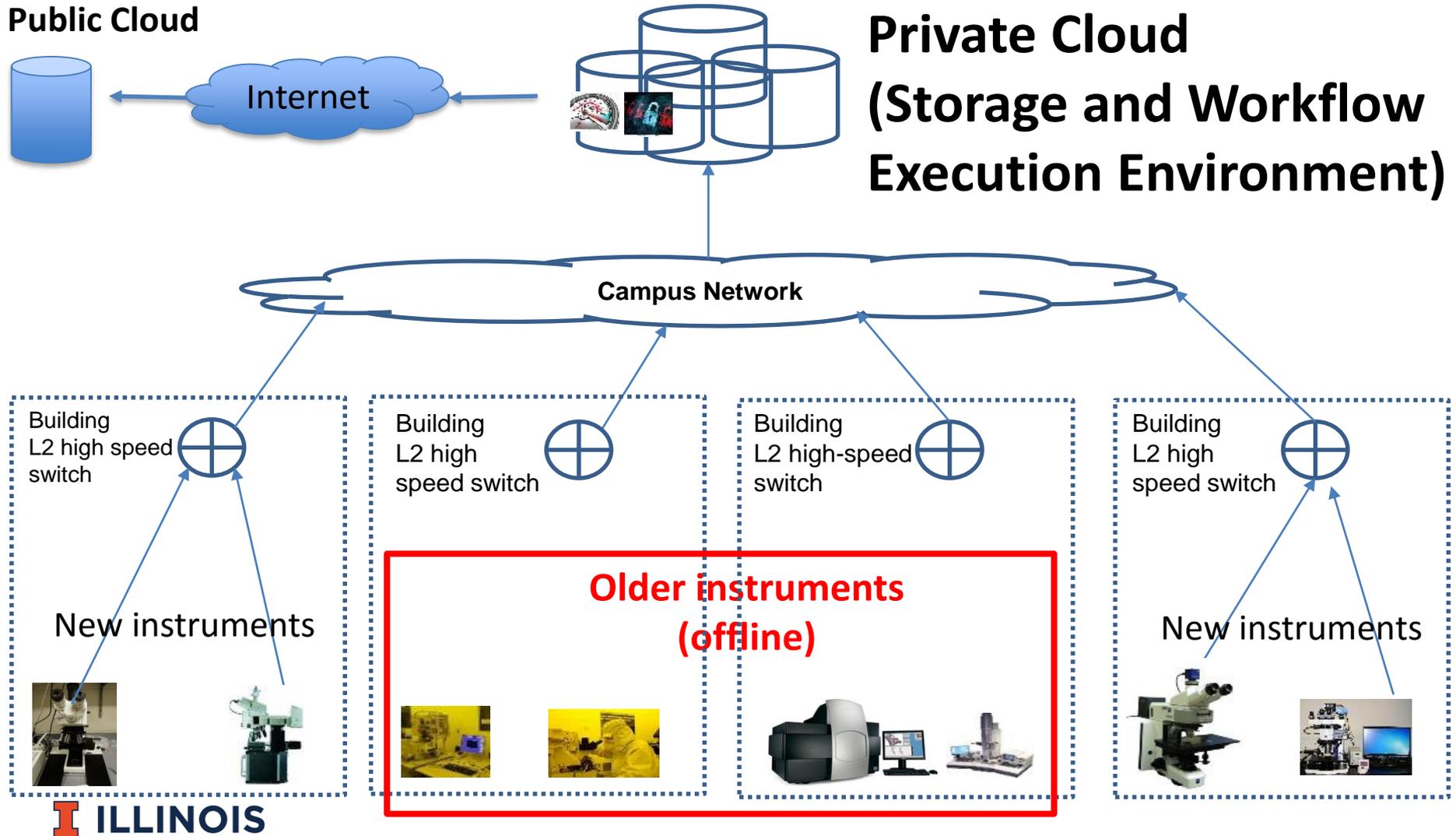


- Previous work often employs a monolithic approach in workflow implementation and execution
 - E.g.: Pegasus, Taverna, Kepler, etc.
 - Run on large-scale & homogeneous datasets



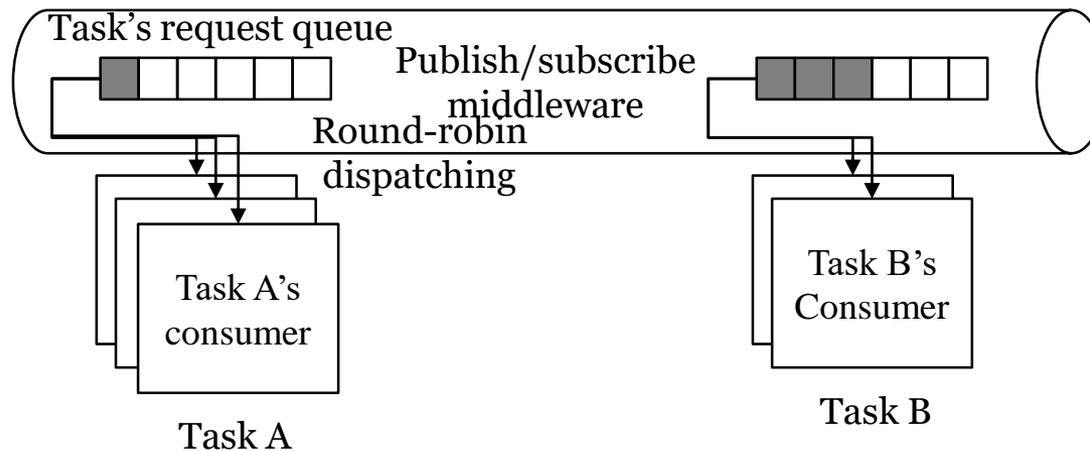
Executing workflows on grid infrastructure

Our Goal for Campus Cyberinfrastructure regarding New Scientific Instruments (1)

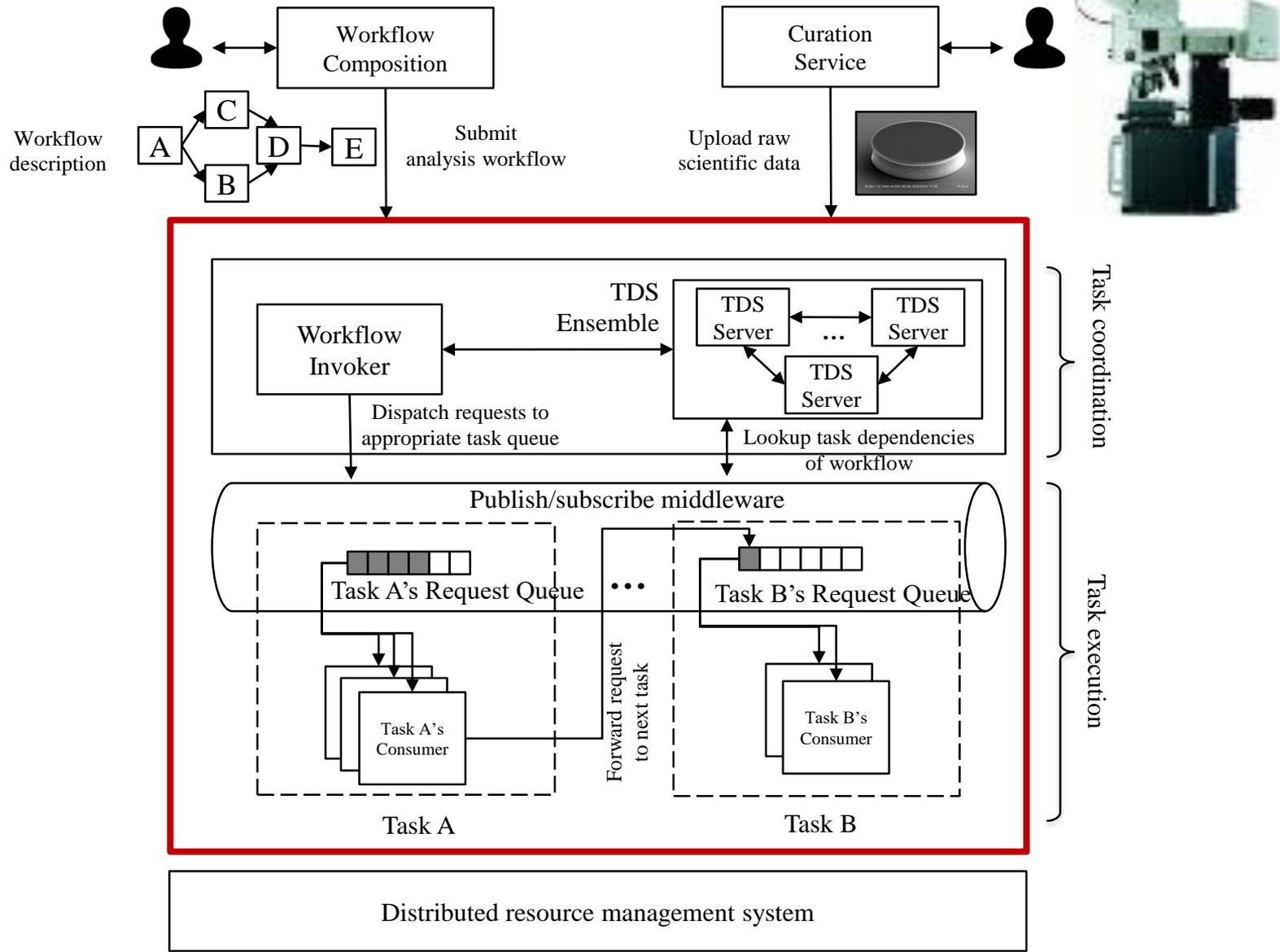


Our Approach: Micro-service execution environment in Private Cloud

- **Micro-services over monoliths:** Each task is modeled as a micro-service
 - Use publish-subscribe middleware to connect between micro-services

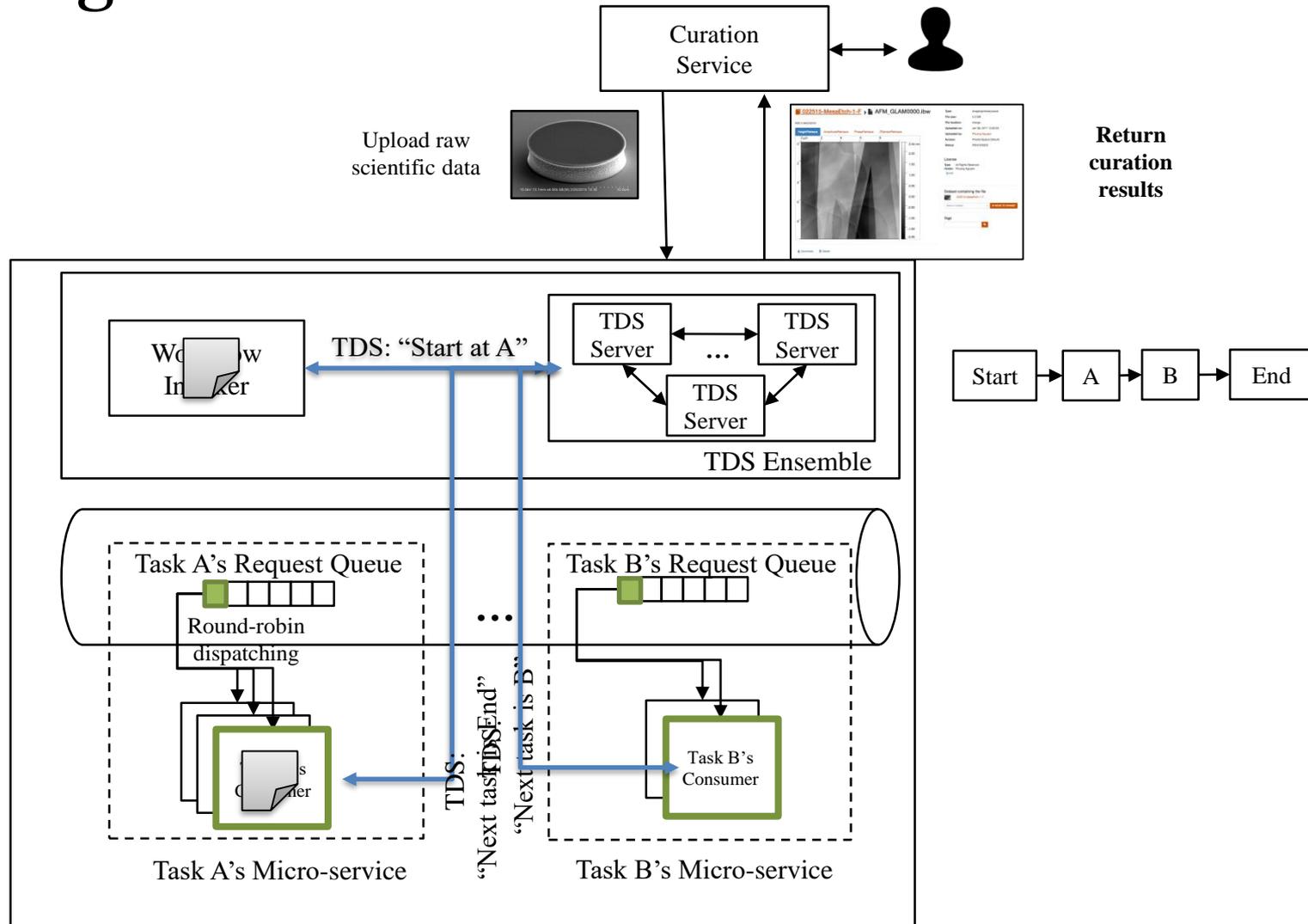


- Separate task dependencies from task implementation & deployment
 - Enable flexible workflow composition
 - Task-level resource provisioning

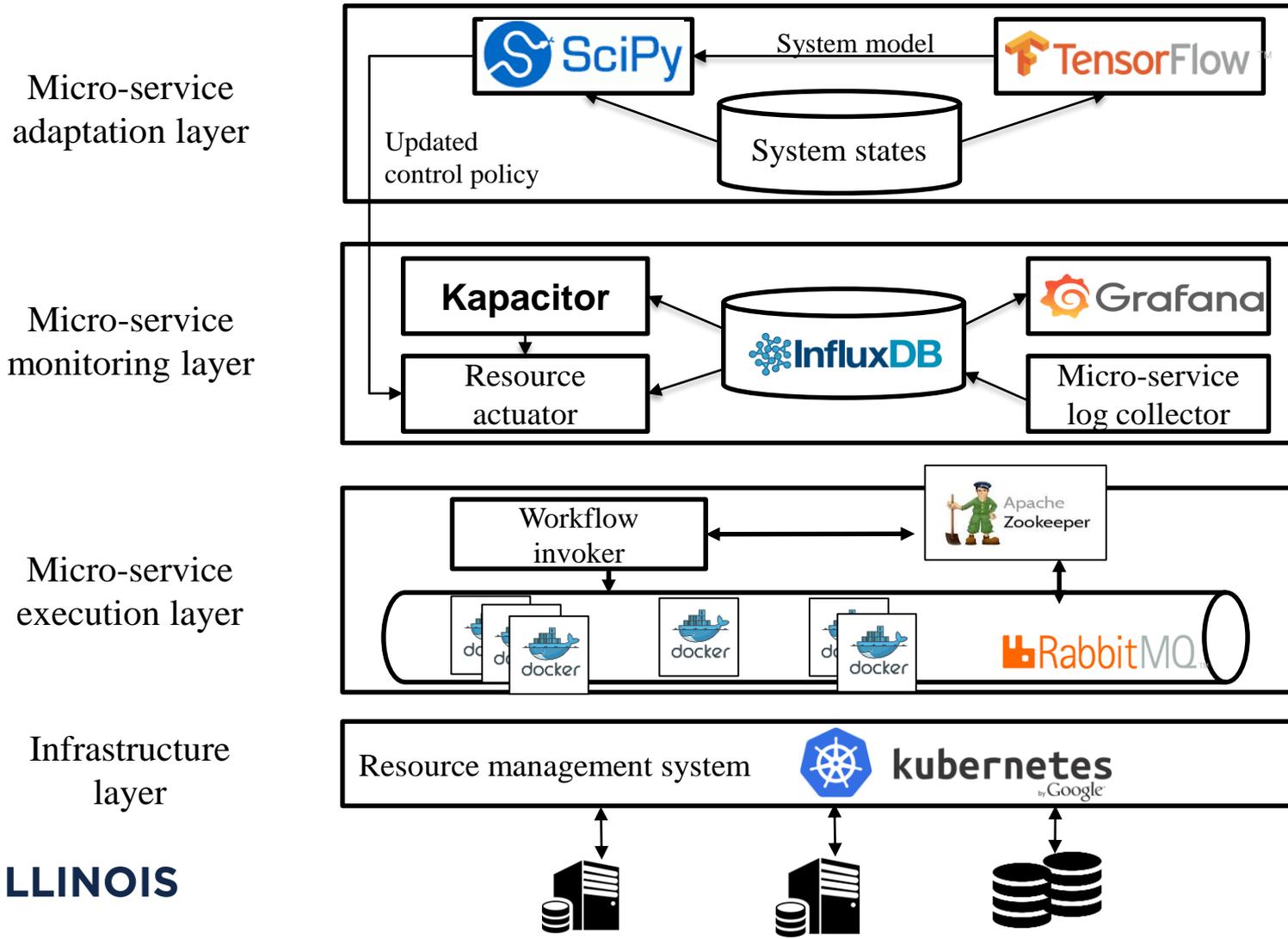


Phuong Nguyen et al., "Resource Management for Elastic Publish Subscribe Systems: A Performance Modeling-based Approach" (IEEE CLOUD 2016)

Example: Executing scientific data processing workflow

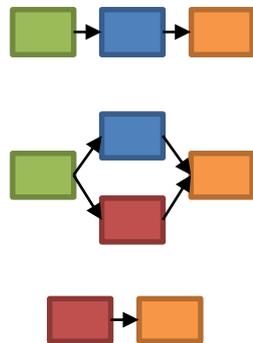


Adaptive micro-service system implementation

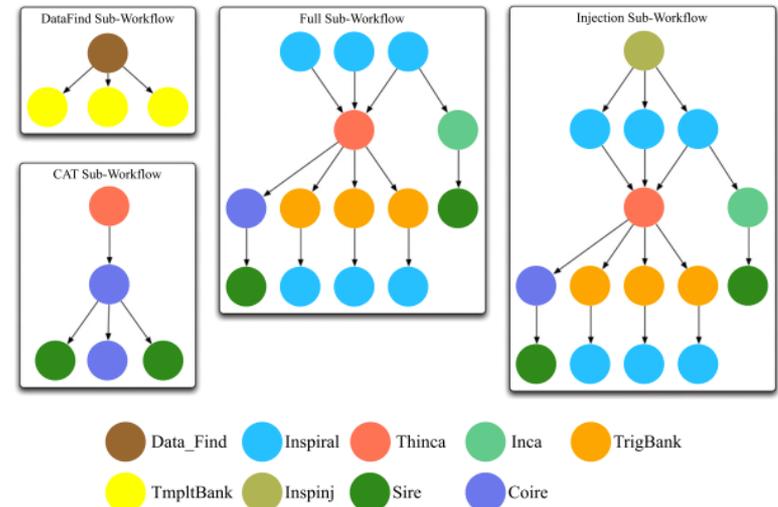


Evaluation #1: Micro-service resource adaptation

- Data processing workflows:
 - MDP: material data processing workflows (to process output of digital microscopy, such as DM3, AFM, etc.)
 - LIGO: analyze data to study stars and black holes

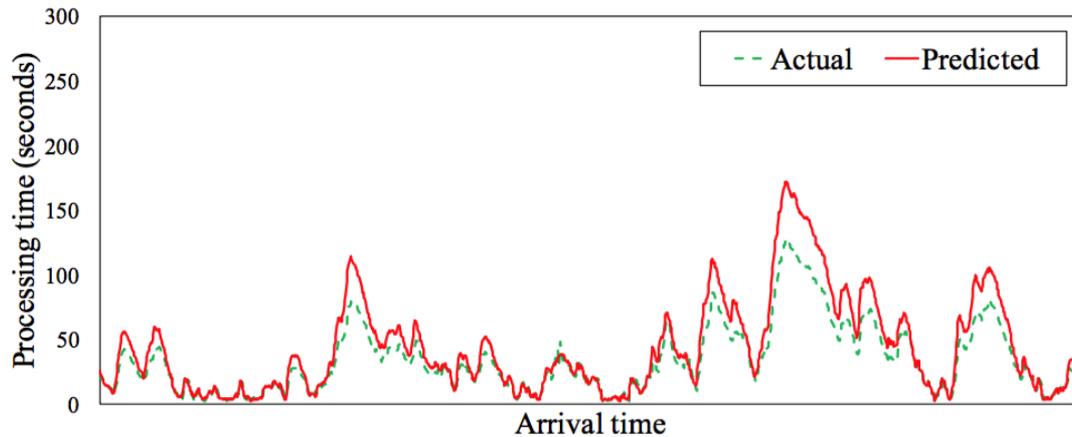


MDP workflows

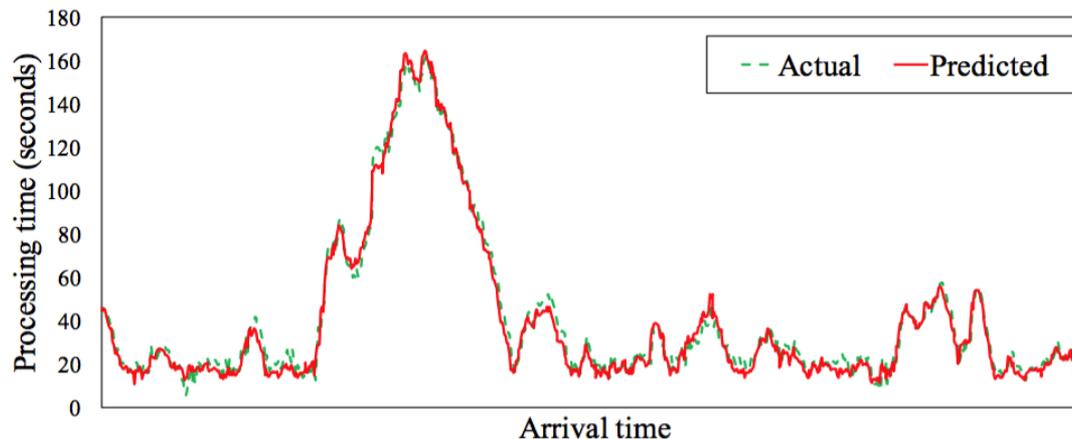


LIGO workflows

Effectiveness of neural network-based system identification



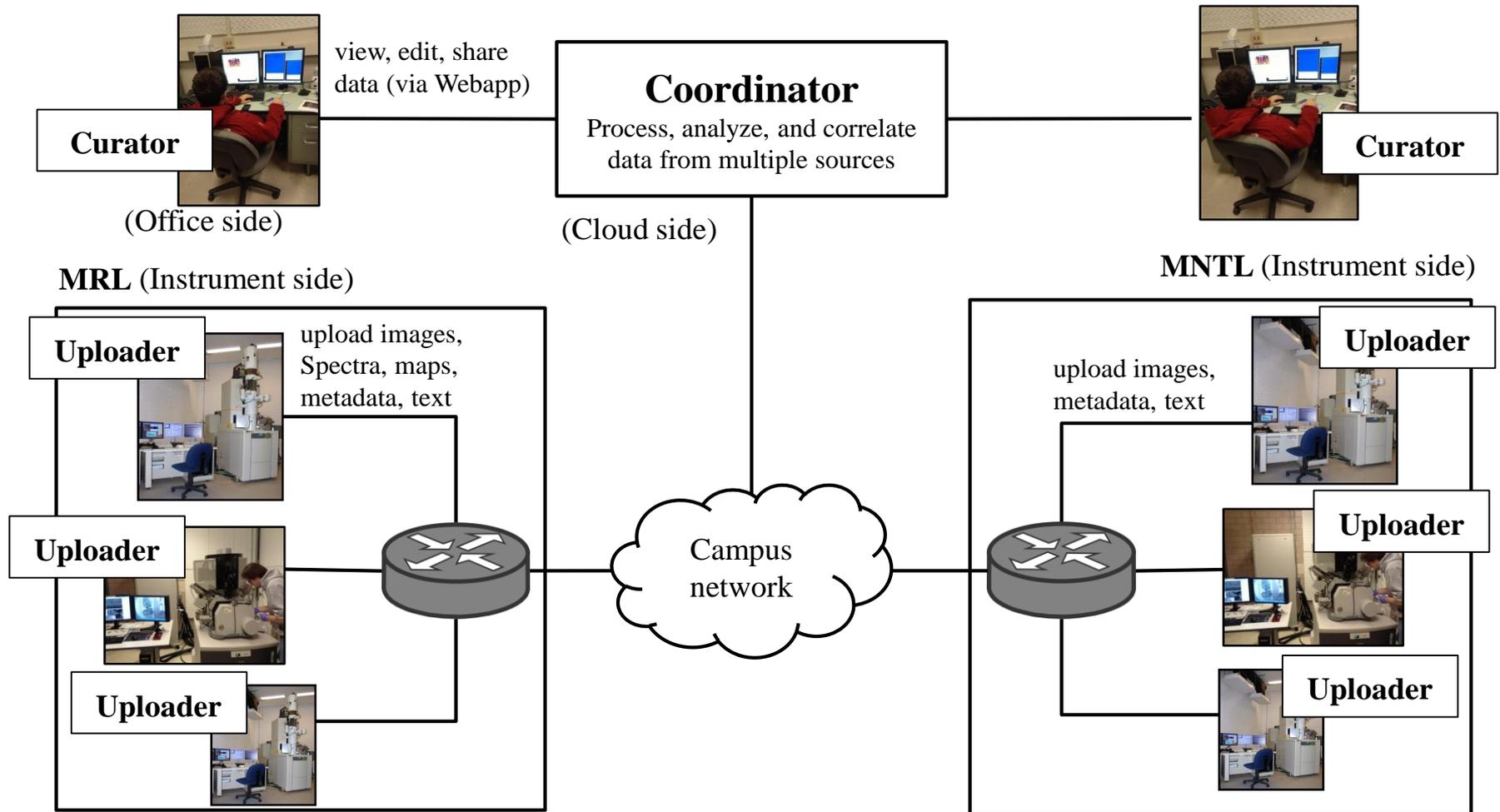
(a) MDP



(b) LIGO

CASE STUDY: 4CeeD - Real-time Acquisition and Analysis Framework for Materials-related Cyber-Physical Environments

4CeeD: Extending micro-service infrastructure to material-related environment



Phuong Nguyen et al., "4CeeD: Real-time Acquisition and Analysis Framework for Materials-related Cyber-Physical Environments". In CCGrid 2017, **Best Paper Award**

4CeeD Uploader Service (Simple and Speed-Up Usage at Microscopes)

Simple steps, with support for advanced usage

01 Choose a collection... [what's this?](#)

Existing collections

New Root Collection

1 Choose a name for the new collection:

2 Choose a description for the new collection:

3

02 Choose a dataset... [what's this?](#)

Existing Datasets

New Dataset

Basic Load Template Create Template Load Previous

1 My Templates: **1** Global Templates: **1** Template Tag Search:

2 Choose a name for your dataset:

3 Dataset Description:

4

5

Name:	Unit Type:	Data Type:	Value:	Required:
<input type="text" value="Brij mass"/>	<input type="text" value="mg"/>	<input type="text" value="Number"/>	<input type="text"/>	<input checked="" type="checkbox"/> Yes <input type="button" value="Remove"/>
Name:	Unit Type:	Data Type:	Value:	Required:
<input type="text" value="What's internalized"/>	<input type="text"/>	<input type="text" value="String"/>	<input type="text"/>	<input type="checkbox"/> No <input type="button" value="Remove"/>
Name:	Unit Type:	Data Type:	Value:	Required:
<input type="text" value="Mass of internalized mole"/>	<input type="text" value="mg"/>	<input type="text" value="Number"/>	<input type="text"/>	<input type="checkbox"/> No <input type="button" value="Remove"/>

03 Click browse or drag and drop files...

1 Drag & Drop Files

2
File Comments:

3

1. Choose or select a collection.
2. Load template and enter user defined metadata to create a dataset.
3. Upload files to cloud coordinator.

4CeeD: Two Types of Uploaders

- Smart Dropbox-like Uploaders
- Two types of Uploaders exist in 4CeeD to make it easier to import data.
 - **Standard Upload** allows for templates to be used
 - **Zip Upload** allows for amounts of data to be uploaded with file structure kept intact.

The screenshot shows the 4CeeD web interface for creating a dataset. The browser address bar shows the URL <https://www.4ceed.illinois.edu/t2c2/uploader>. The page has a navigation bar with the 4CeeD logo, user information, and a search bar. The main content area is titled "My Templates:" and includes a dropdown menu for "RBS Data". Below this, there are sections for "Global Templates:" and "Template Tag Search:". A section titled "Choose a name for your dataset:" contains a text input field with the value "2018-05-11". The "Dataset Description:" section has a text input field with the value "RBS Data". Below the description are two buttons: "ADD NEW FIELD" and "CLEAR TEMPLATE". The main part of the interface is a table with columns for Name, Value, Units, Data Type, and Required. Each row represents a field in the dataset, with a "REMOVE" button next to it.

Name:	Value:	Units:	Data Type:	Required:	
Beam Ion	He+		String	Yes	REMOVE
Beam Energy	2.024	MeV	Number	Yes	REMOVE
Beam Current	100	nA	Number	Yes	REMOVE
alpha (incident angle)	22.5	degrees	Number	No	REMOVE
beta (exit angle)	52.5	degrees	Number	No	REMOVE
theta (scattering angle)	150	degrees	Number	No	REMOVE

At the bottom of the form is a large orange button labeled "CREATE DATASET".

4CeeD Curator Service (Speed-Up Curation)

File View

4CeeD You - Shared - Create - Help - Search

Demo Dataset Name > 2016_04_14_Gd-filled micelle_0008.dm3

Type: image/digitalmicrograph
File size: 17.4 MB
File location: mongo
Uploaded on: Oct 07, 2016 19:04:53
Uploaded by: Steve K
Access: Private (Space Default)
Status: PROCESSED

License
Type: All Rights Reserved
Holder: Steve K
[Edit](#)

Dataset containing the file
Demo Dataset Name
Select a Dataset [Move to Dataset](#)

Tags

Metadata
- Extracted by <http://clowder.ncsa.illinois.edu/extractors/deprecatedapi> on Oct 7, 2016

Microscope Info Indicated Magnification: 10000.0
Microscope Info Magnification Interpolated: False
Acquisition Parameters High Level Shutter Pre Exposure Compensation (s): 0.0
Acquisition Frame Intensity Range Dark Current (counts/s): 0.0
Acquisition Frame Sequence Exposure Time (ns): 500003080.0

[Preview, annotate, download,
extracted metadata]

Dashboard View

4CeeD You - Shared - Create - Help - Search

2 3

[Profile](#) [Create Space](#) [Create Dataset](#) [Create Collection](#) [Template Management](#)

Activity Tree View My Spaces **My Datasets** My Collections Followers

Create datasets to upload and publish data. Further organize your data using folders and assign metadata at both the file and dataset level.

4 See More

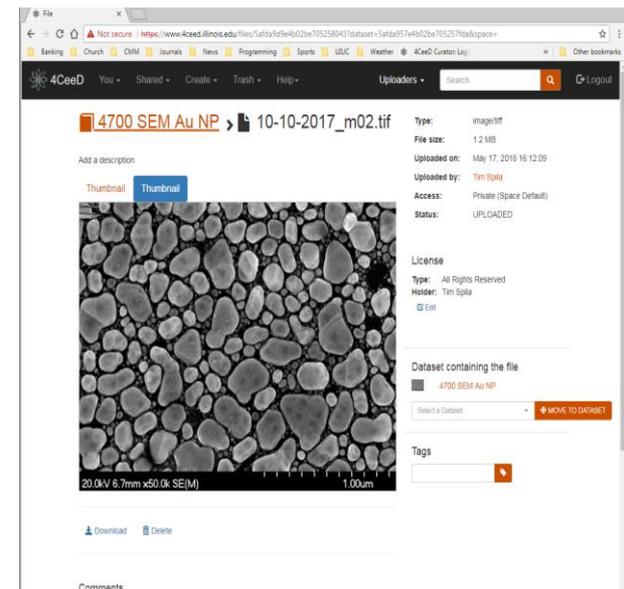
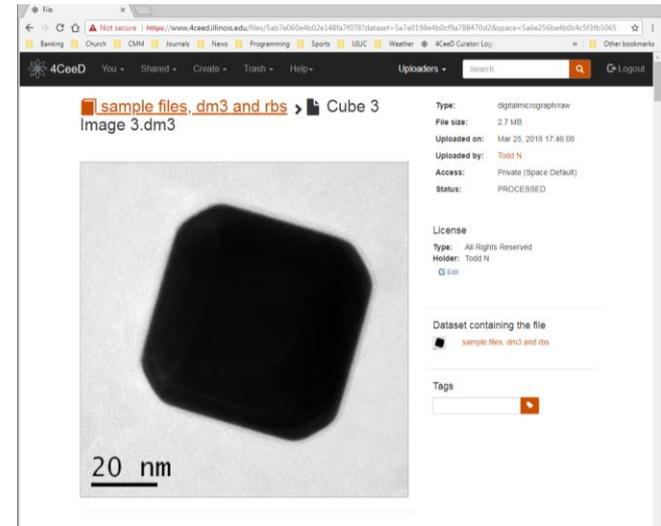
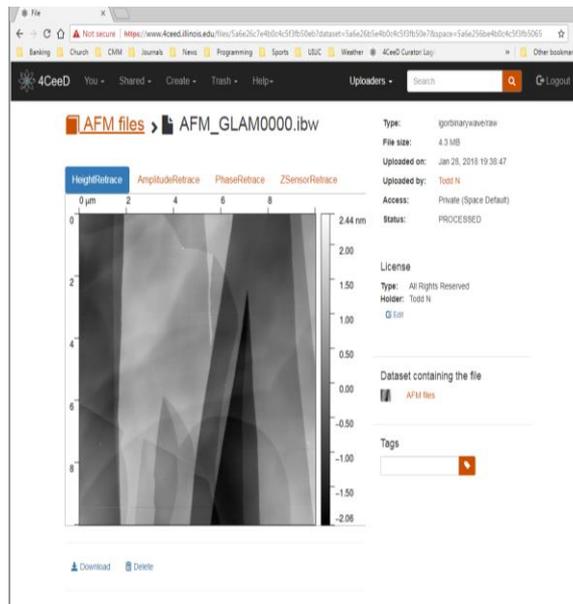
Demo Dataset Name
Demo Dataset Description
0 1 0 0 0 1

demo dataset
0 1 0 0 0 1

[Dashboard management]

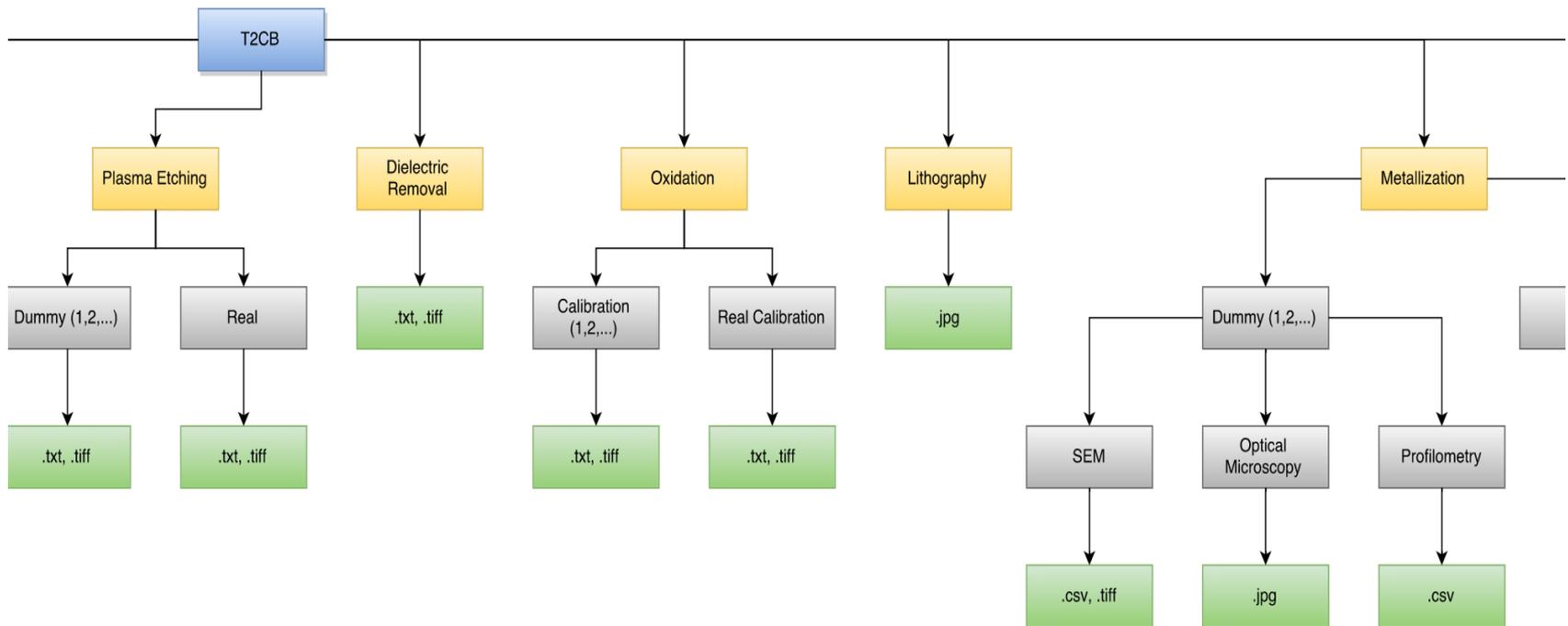
4CeeD: Extractors as Micro-Services at Cloud Side

- TEM Extractor
 - Works with DM3 Files
- SEM Extractor
- AFM Extractor

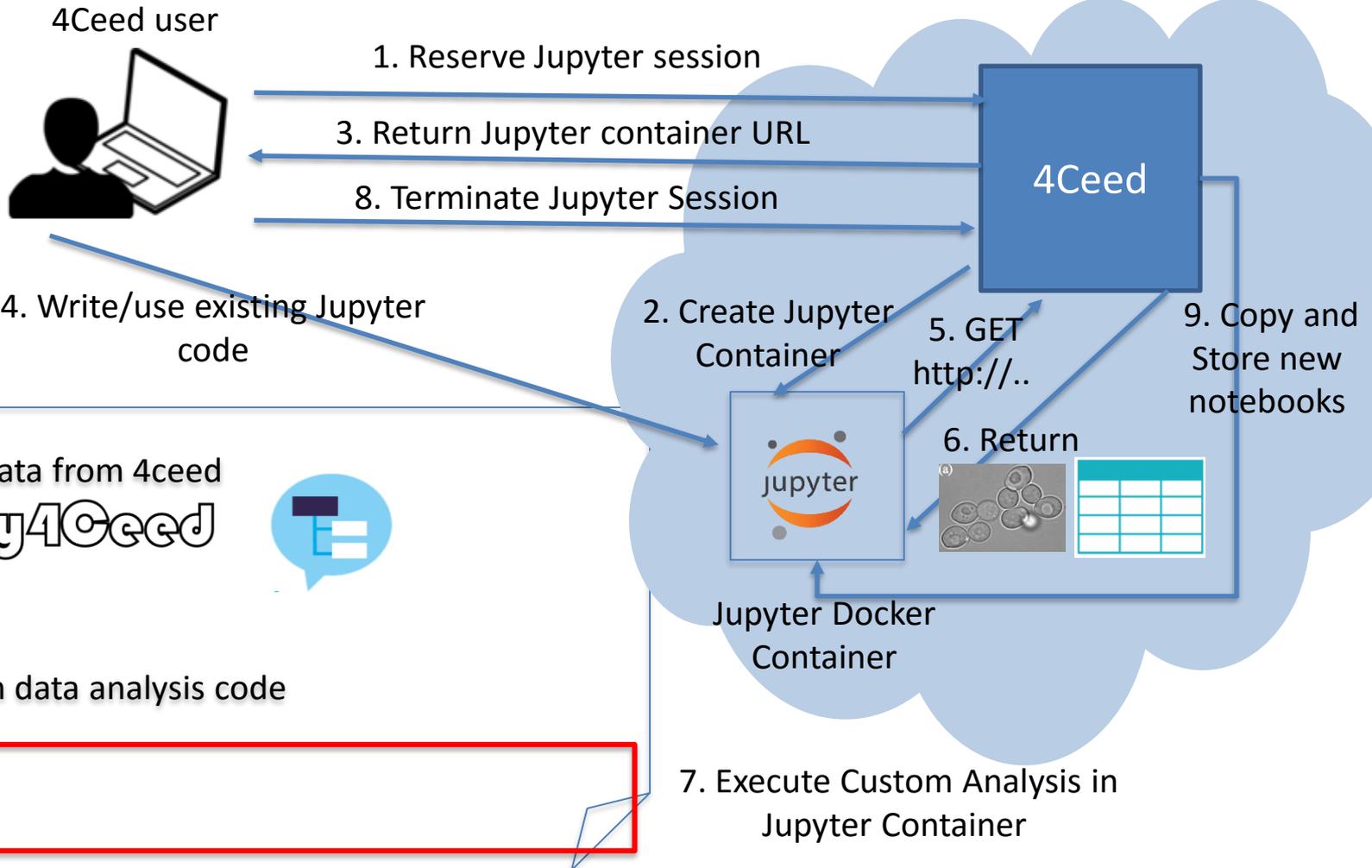


4CeeD Smart **Data** Management

4CeeD Data Model organizes projects into collections, datasets, and files. These can then be shared in spaces. 4CeeD utilizes and modifies NCSA Clowder data management system.



4Ceed++: Jupyter Notebook Integration



4CeeD Production System

Goals:

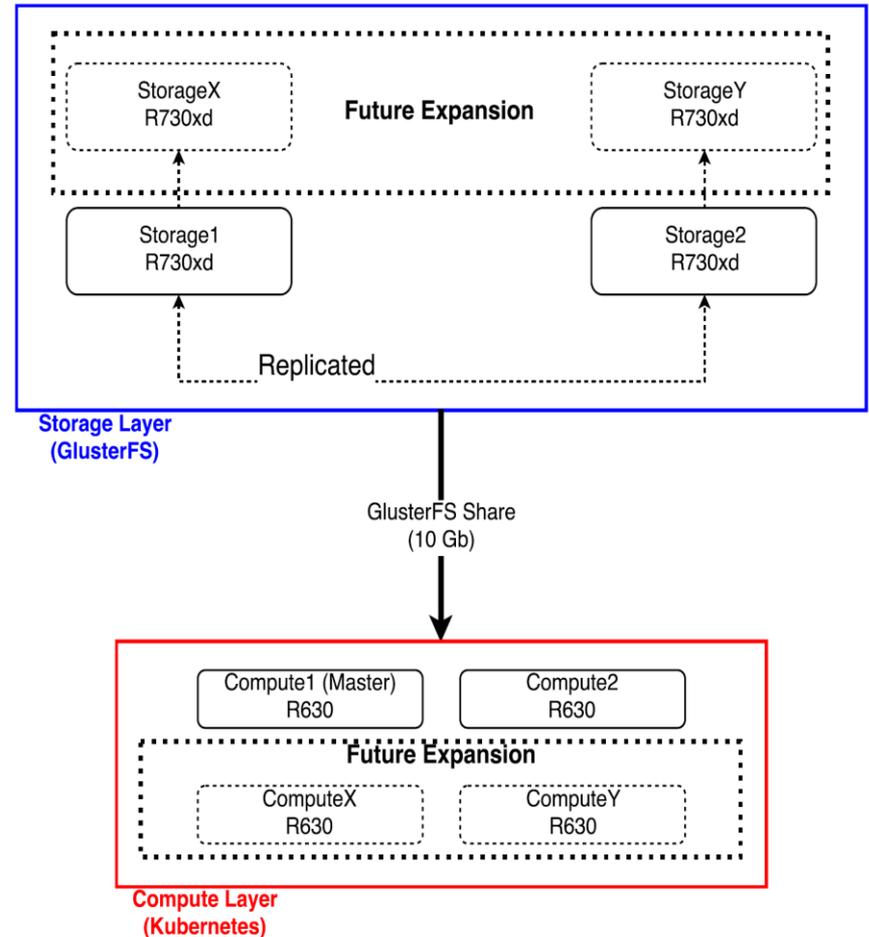
- Redundancy
- Availability
- Scalability

Storage Layer:

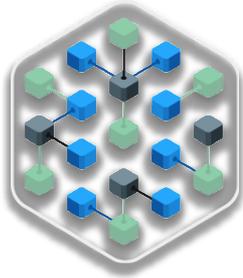
- **40 TB (20 TB per investor)**
- Replicated for redundancy

Compute Layer:

- Docker container orchestration (Kubernetes)
- **Single master**
(High Available masters in future)



Our approach

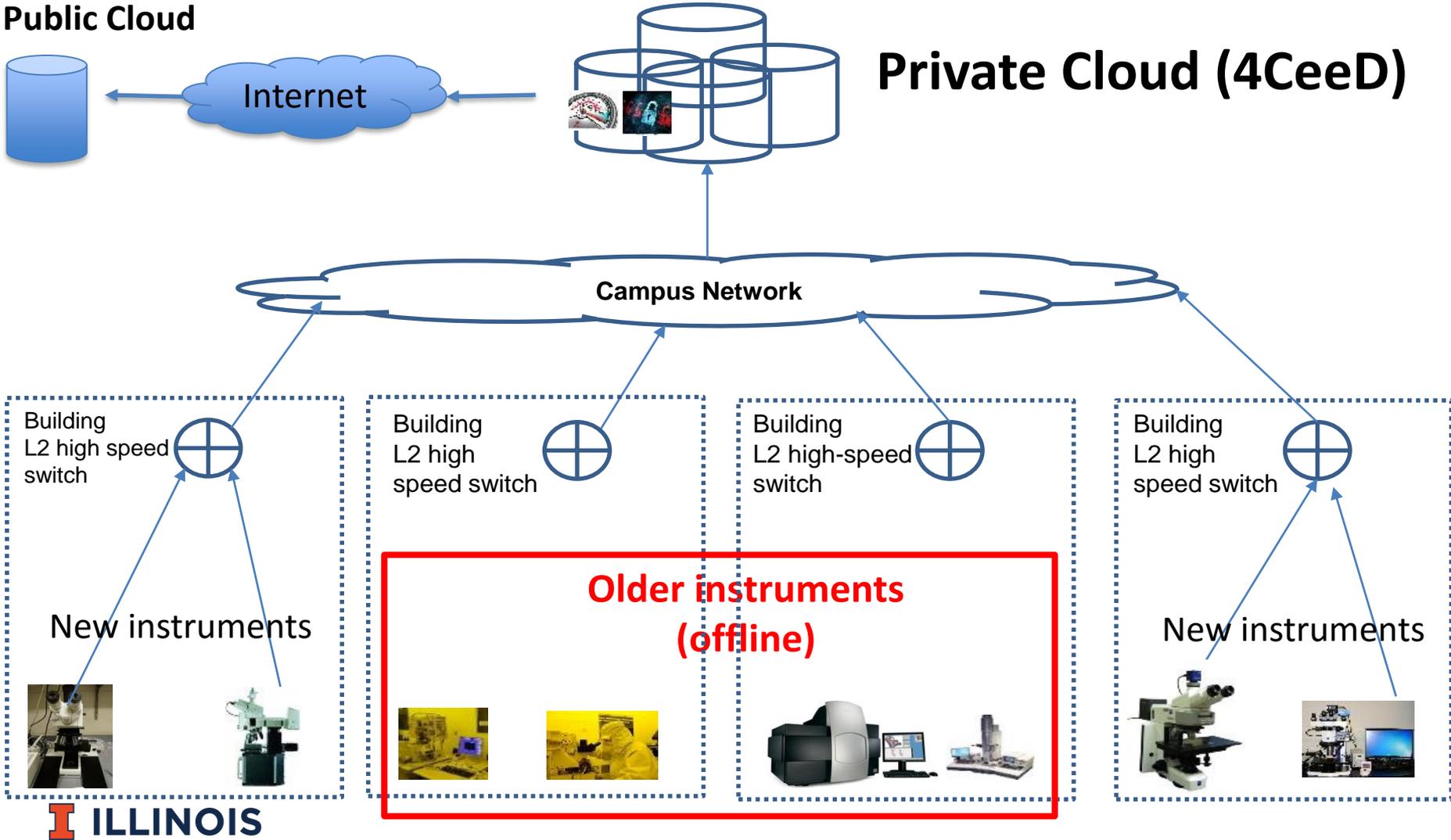


- ✓ Micro-service private cloud execution environment for instrument data curation and coordination (4CeeD)

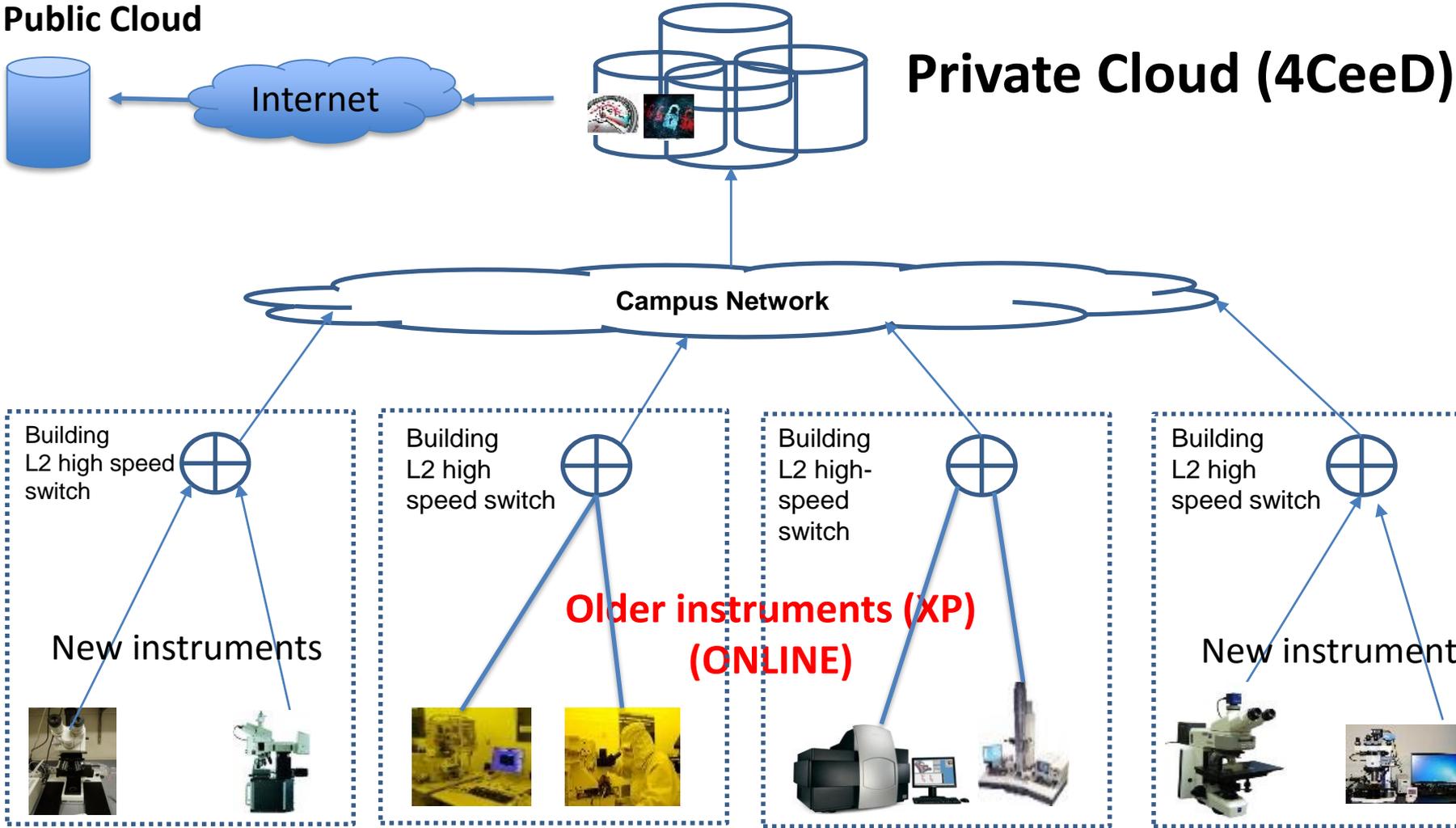


- ✓ Data acquisition from aging instruments (BRACELET)

Current Status for Campus Cyberinfrastructure regarding Aging Scientific Instruments



Our Goal for Campus Cyberinfrastructure regarding Aging Scientific Instruments (2)



Challenges of connecting offline older instruments



- **Performance mismatch:** Older instruments' Windows XP runs network protocols at lower bandwidth speeds (10Mbps or 100Mbps)

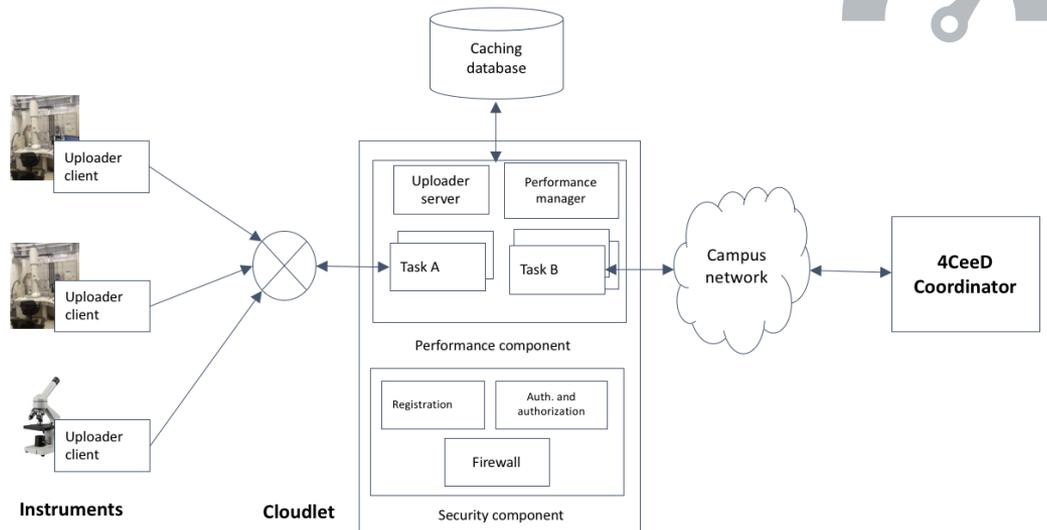


- **Obsolete security:** Older devices and their OS systems cannot be patched, hence being vulnerable & taken offline

BRACELET: Putting edge device between older instruments and private cloud

Performance:

- Have two network interfaces configured at different speeds
- Traffic shaping & offloading between edges & cloud

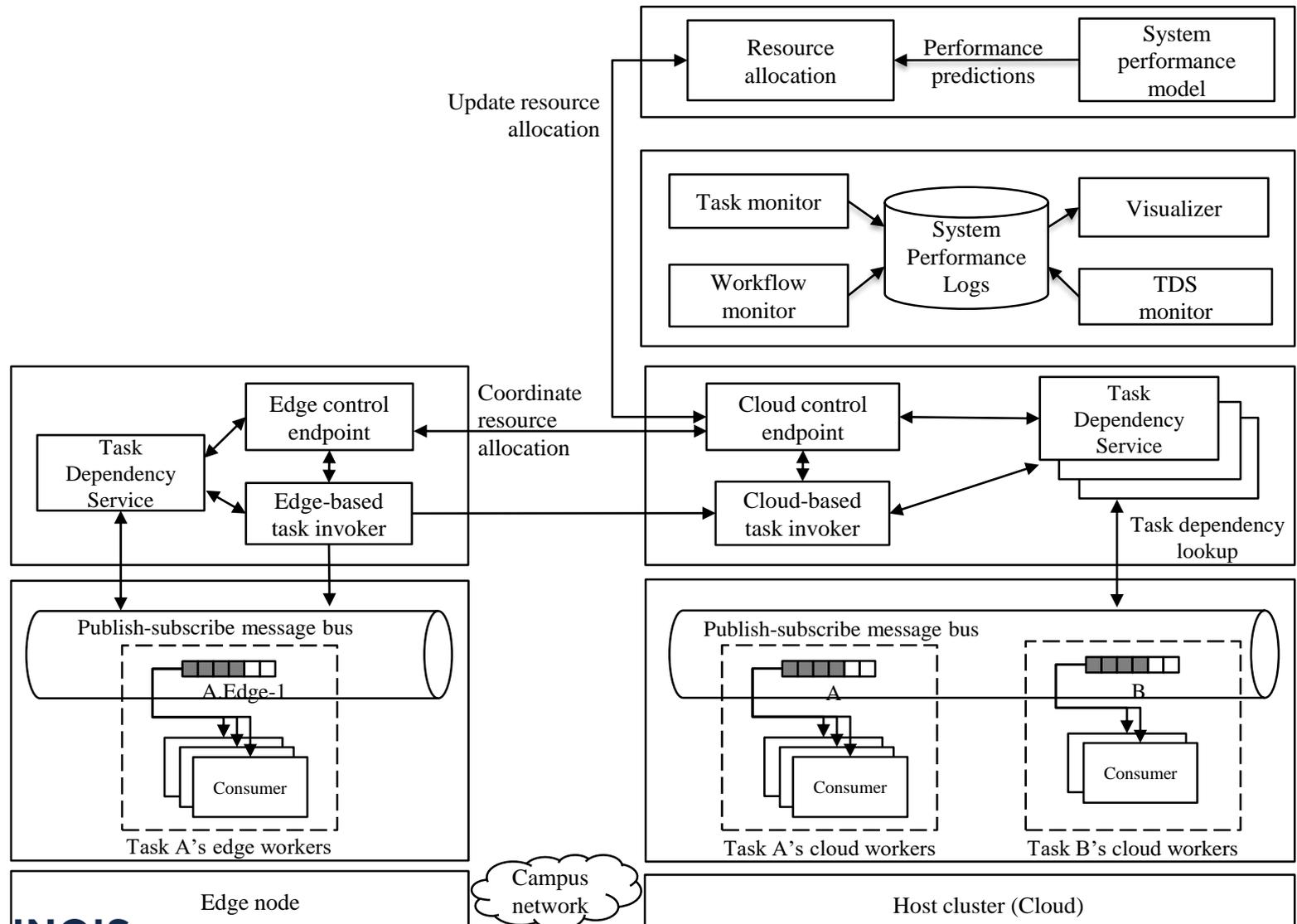


BRACELET in 3-tier architecture

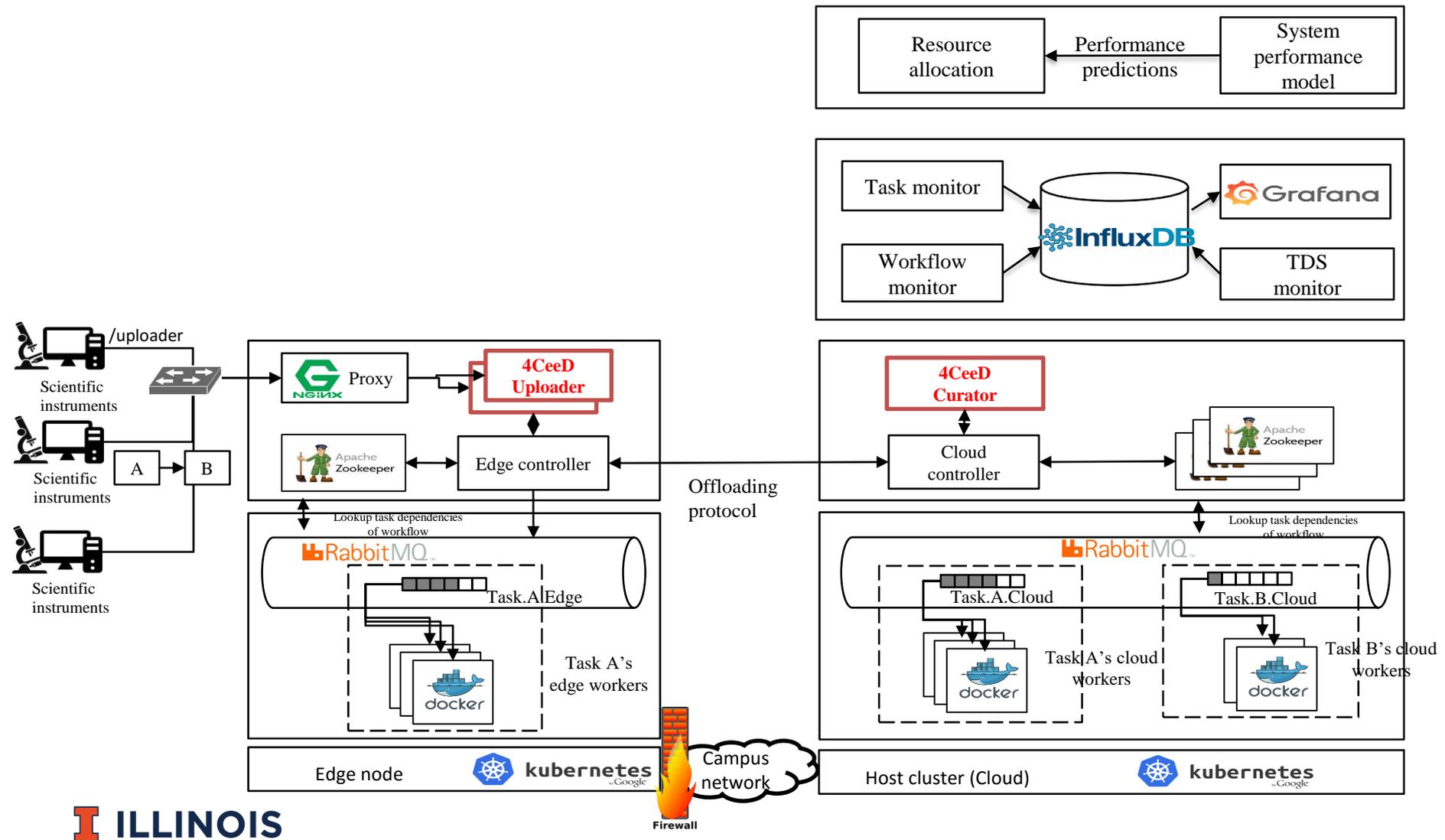
Security:

- User & instrument registration
- Data encryption during upload
- Firewall to protect against external threats

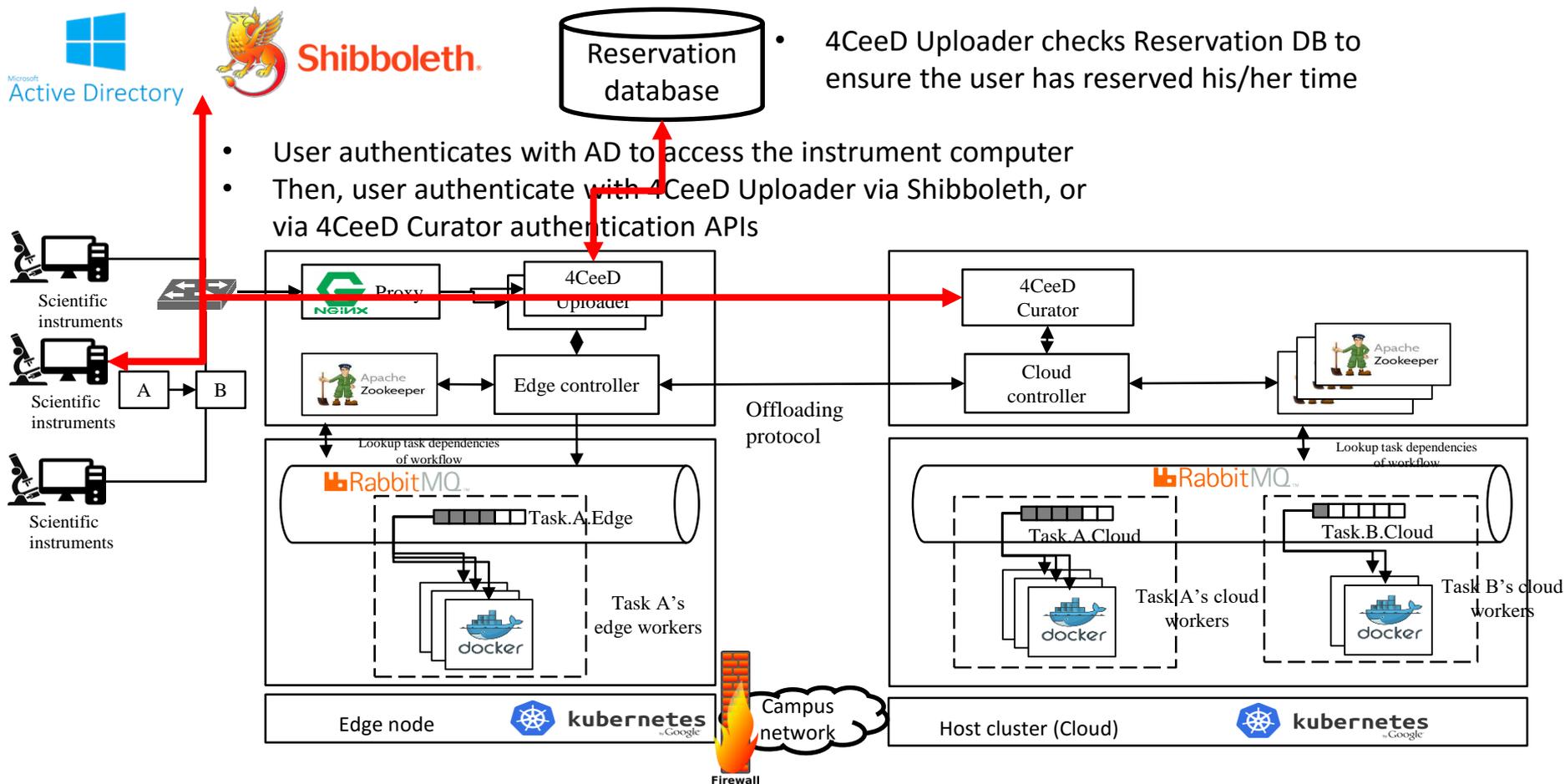
BRACELET: Extending cloud-based architecture to the edges for seamless integration



BRACELET's implementation & integration with 4CeeD

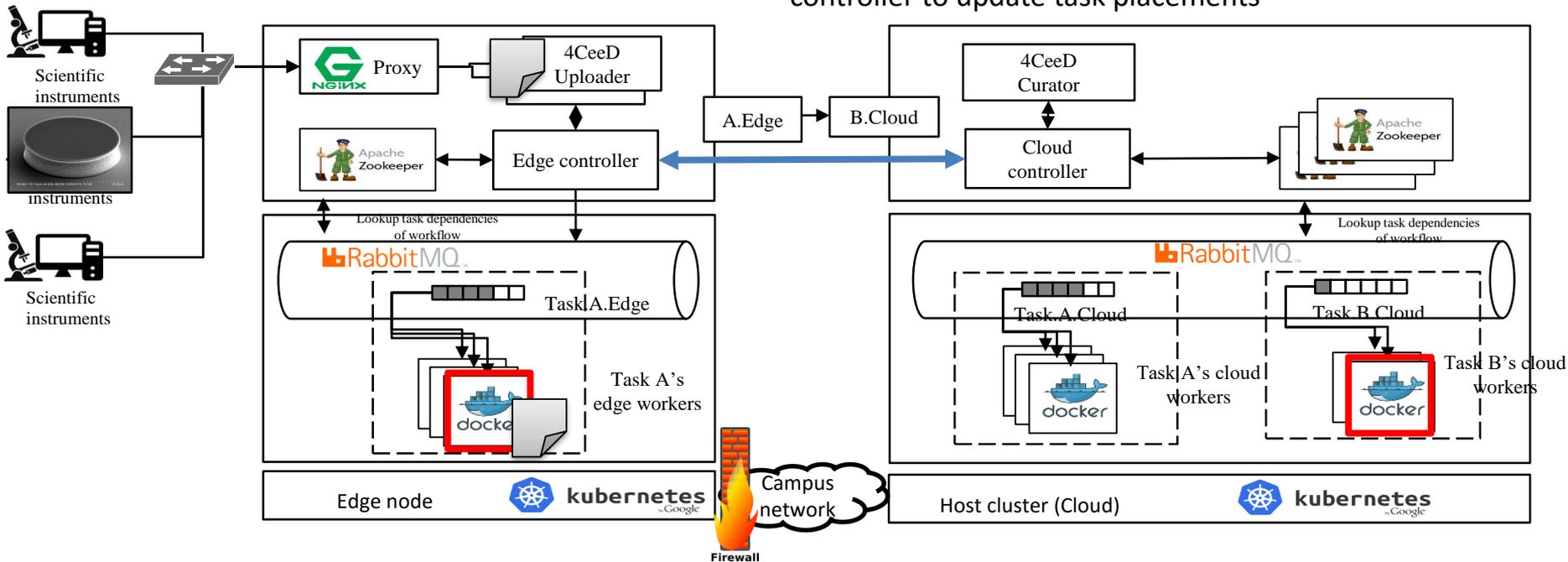


User authentication from instruments via BRACELET

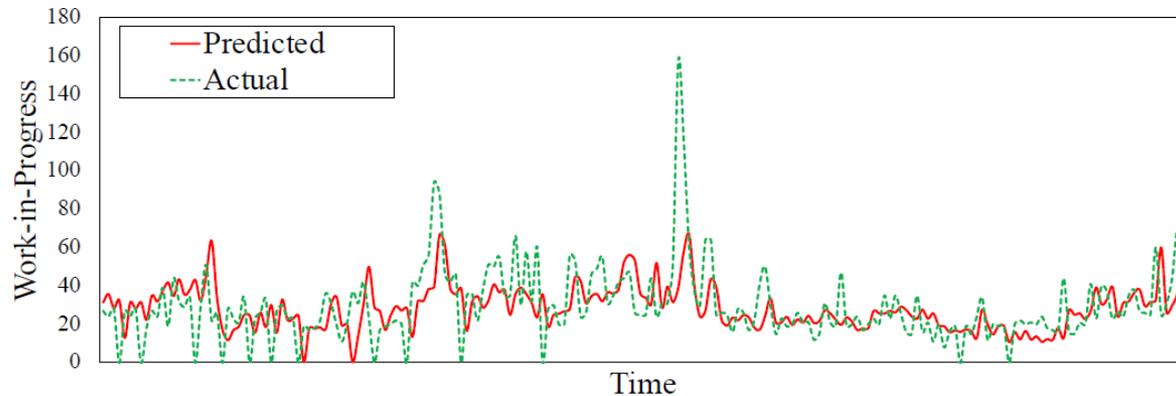


Computation offloading between edge & cloud

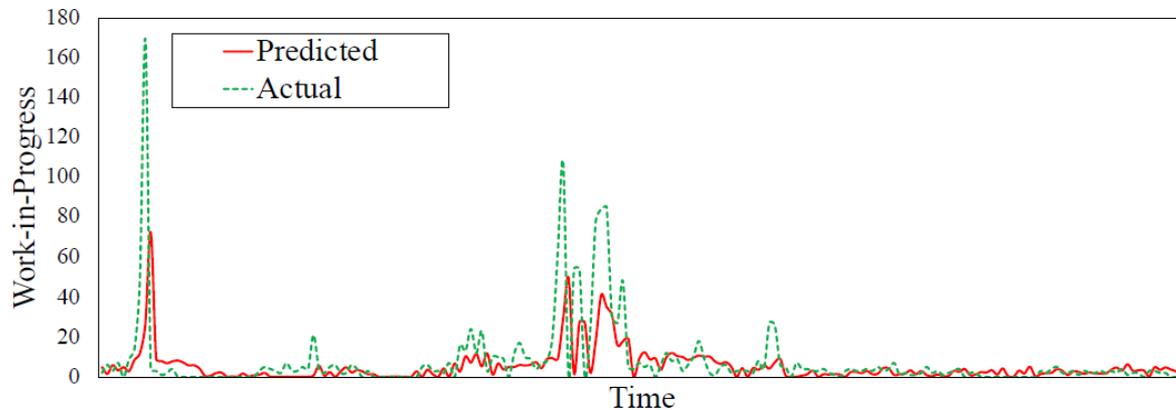
- After learning about the request, the data consumer forwards processing requests to the next task (following current placement)
- 4CeeD Uploader communicates with local Edge controller to learn about where to send request to Edge controller periodically communicates with cloud controller to update task placements



Evaluation of 4Ceed-BRACELET micro-service performance prediction



(a) Cloud-based micro-service D



(b) Edge-based micro-service C on cloudlet E1

Software Availability

- All publications are available in IEEE Digital Libraries
- All software systems, 4CeeD, BRACELET are open source
- All Projects are described at <https://t2c2.csl.illinois.edu/>
- 4CeeD System is available for download <https://github.com/4ceed>
 - Contact: Steve Konstanty (stevek@illinois.edu) and Todd Nicholson (tcnichol@illinois.edu)
- Bracelet system is in testing phase. Software will be released by December 2019.
 - Contact: Steve Konstanty (stevek@illinois.edu)

Lessons Learned

- We have explored **novel cloud system approaches** that lend themselves well for real-time and trustworthy materials-to-device data and metadata storage, management, and computing of workflows over these data.
- Lightweight micro-service cloud architecture for materials genomic challenge is the way to go, including the **three tier approach** (instrument - private cloud – public cloud).
- Bringing **aging instruments online** proved to be very challenging for very old OS systems
- Hardest Part is convincing experimentalists/scientists to use new data management techniques and new cyber-infrastructure
 - Continuous training and inclusion of modern data management systems into experimental instrumentation classes will be needed.
- To achieve **sustainability**, it is crucial to work with **college and campus IT teams!**

Acknowledgement

- Research and Development Team:
 - Phuong Nguyen (CS), Tarek Elgamal (CS), Zhe Yang (CS), Tuo Yu (CS), Xiaoyuan Wang (CS), Steve Konstanty (Senior Research Programmer), Todd Nicholson Research Programmer), Patrick Su (ECE), Robert Kaufman (ECE), Tommy O'Brien (ECE).
- NSF Funding: ACI DIBBS and CC Programs
 - *NSF ACI DIBBS Award 1443013 – 4CeeD (T2C2 Project)*
 - *NSF OAC CC Award 1659293 - BRACELET*
- Co-PIs and Collaborators:
 - Roy Campbell (CS/CSL/Engineering IT), Indranil Gupta (CS), Paul Braun (MRL), Brian Cunningham (MNTL/ECE), Greg Pluta (MNTL), Tim Spila (MRL), Michael Chan (Engineering IT), Tracy Smith (IT Tech Services), Kenton McHenry (NCSA), John Dallesasse (ECE/MNTL), Mark McCollum (MNTL), Gianni Pezzarossi (Engineering IT), Stuart Turner (Engineering IT), Laura Herriott (Engineering IT)

Publications

- Phuong Nguyen, Klara Nahrstedt, “Resource Management for Elastic Publish Subscribe Systems: A Performance Modeling-based Approach”, **IEEE International Conference on Cloud Computing (CLOUD 2016)**, San Francisco, CA, June 2016
- Phuong Nguyen, Steven Konstanty, Todd Nicholson, Thomas O'Brien, Aaron Schwartz-Duval, Timothy Spila, Klara Nahrstedt, Roy Campbell, Indranil Gupta, Michael Chan, Kenton McHenry and Normand Paquin, “4Ceed: Real-Time Data Acquisition and Analysis Framework for Material-related Cyber-Physical Environments”, **IEEE/ACM 17th IEEE/ACM International Symposium on Cluster, Cloud and Grid Computing**, Madrid, Spain, May 14-17, 2017 (**Best Paper Award**)
- Phuong Nguyen, Klara Nahrstedt, “MONAD: Self-adaptive Micro-service Infrastructure for Heterogeneous Scientific Workflows”, **14th IEEE International Conference on Autonomous Computing (ICAC 2017)**, July 17-21, 2017, Columbus, Ohio
- Phuong Nguyen, Tarek Elgamal, Steve Konstanty, Todd Nicholson, Stuart Turner, Patrick Su, Michael Chan, Klara Nahrstedt, Tim Spila, Kenton McHenry, John Dallesasse, Roy Campbell, “BRACELET: Edge-Cloud Micro-service Infrastructure for Aging Scientific Instruments”, **IEEE International Conference on Computing, Networking, and Communications (ICNC) 2019**, Hawaii, February 2019.
- Zhe Yang, Phuong Nguyen, Haiming Jin, Klara Nahrstedt, “MIRAS: Model-based Reinforcement Learning for Microservice Resource Allocation over Scientific Workflows”, **IEEE International Conference on Distributed Computing Systems (ICDCS 2019)**, July 2019, Dallas, TX.
- Zhe Yang, Patrick Su, Robert Kaufman, Steve Konstanty, John Dallesasse, Klara Nahrstedt, “SENSELET: Sensory Network Infrastructure for Scientific Lab Environments”, **ACM Practice and Experience in Advanced Research Computing Conference Series (PEARC 2019)**, Chicago, IL (Poster), August 2019