# 4CeeD Backend Services

**4CeeD Backend Services**

Robert Kaufman (rbkaufm2@Illinois.edu),  Leah Espenhahn (leahe2@illinois.edu),
Beitong Tian (beitong2@illinois.edu),  **Prof. Klara Nahrstedt (klara@illinois.edu)**
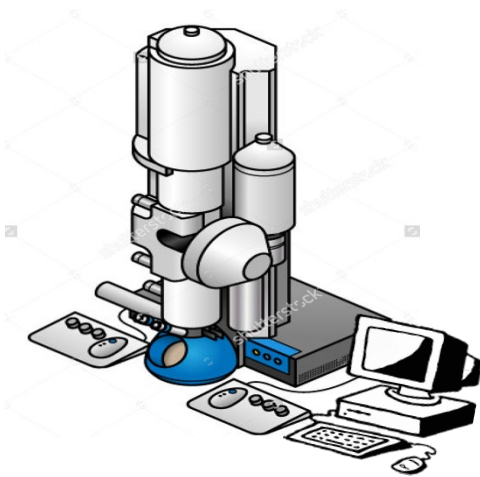
# A timely and trusted curator and coordinator of scientific data
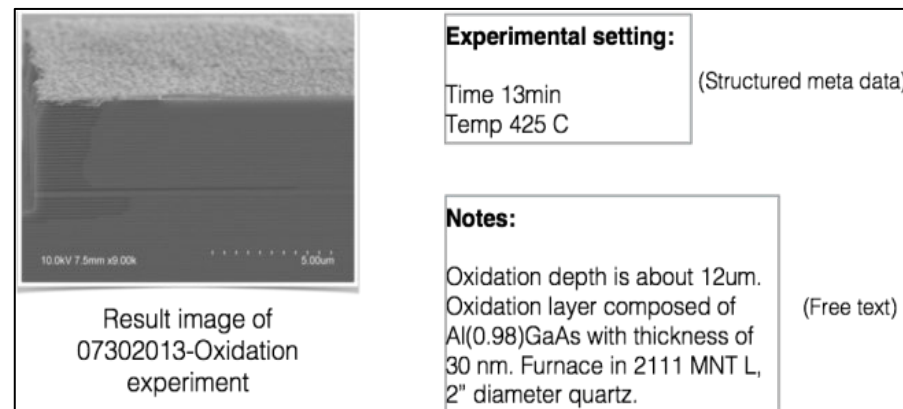
# Outline

- 4CeeD Distributed Architecture, Backend Cloud Concepts and Services
  - What is 4Ceed and its goals
  - What is behind the 4CeeD Dashboard
  - 4CeeD Cloud Design and Deployment
  - How to deal with Aging Scientific Instrument

# What is 4CeeD and its goals?

- Address Scientific Digital Data Acquisition, Curation and Sharing prior to Scientific Publication of Results via Private Cloud Storage Facility



**Instrument**
**(in MRL/HMNTL/BI)**



Sample output data from SEM microscopy
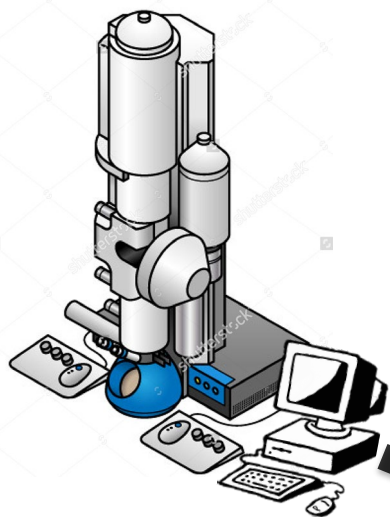
4CeeD

ILLINOIS

3

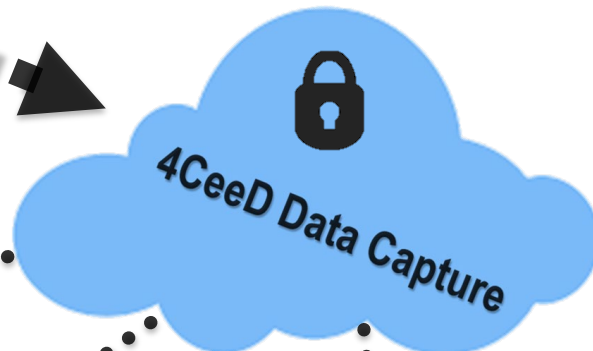# How this looks from 4CeeD [Datasets]



- 4CeeD is designed to present only pertinent information for quick understanding of the experiment

# Scenario with 4CeeD Integration

- Fabricate experimental sample
- Prepare analytical sample
- Bring sample to instrument for analysis
- **Extract data (NO FILE CONVERSION)**
- **Transport data to office computer (DIRECT)**
- **Analyze data (REAL TIME)**
- Repeat per iteration

**Instrument**
**(MRL/MNTL)**

**4CeeD Data Capture**

**Laptop**

**Campus PC**

**Collaborators**

**Office**
**(MRL/MNTL Office)**

**Benefits Data Interface**

- Merge sample file with data
- Aggregated reporting
- Easy data/annotation/search

**4CeeD**

ILLINOIS

# Outline

- 4CeeD Distributed Architecture, Backend Cloud Concepts and Services
  - What is 4Ceed and its goals
  - What is behind the 4CeeD Dashboard
  - 4CeeD Cloud Design and Deployment
  - How to deal with Aging Scientific Instrument

4CeeD

ILLINOIS

# Increasingly data-driven and interdisciplinary scientific research in Physical Sciences and Live Sciences

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



Digital microscope

# 4CeeD Design Considerations - Distributed View



Private Cloud (4CeeD)

Here is your instrument data stored

Campus Network

Building
L2 high speed
switch

MRL instruments
And offices

Building
L2 high speed
switch

MNTL instruments
And offices

4CeeD

ILLINOIS

# 4Ceed Design Considerations – Component View



User

Curator

view, edit, share data
(via Webapp)

**Coordinator**
Process, coordinate, correlate, store data from multiple sources

User

Curator

**MRL**

upload DM3, images, metadata, text

Uploader/Curator

Uploader/Curator

Uploader/Curator

**Cloudlet**

bulk data transfer
(via API)

**MNTL**

upload DM3, images, metadata, text

Uploader/Curator

Uploader/Curator

Uploader/Curator

**Cloudlet**

4CeeD

ILLINOIS

# 4CeeD Design Considerations - Multimodal data format View



Result image of 07302013-Oxidation experiment

**Experimental setting:**

Time 13min
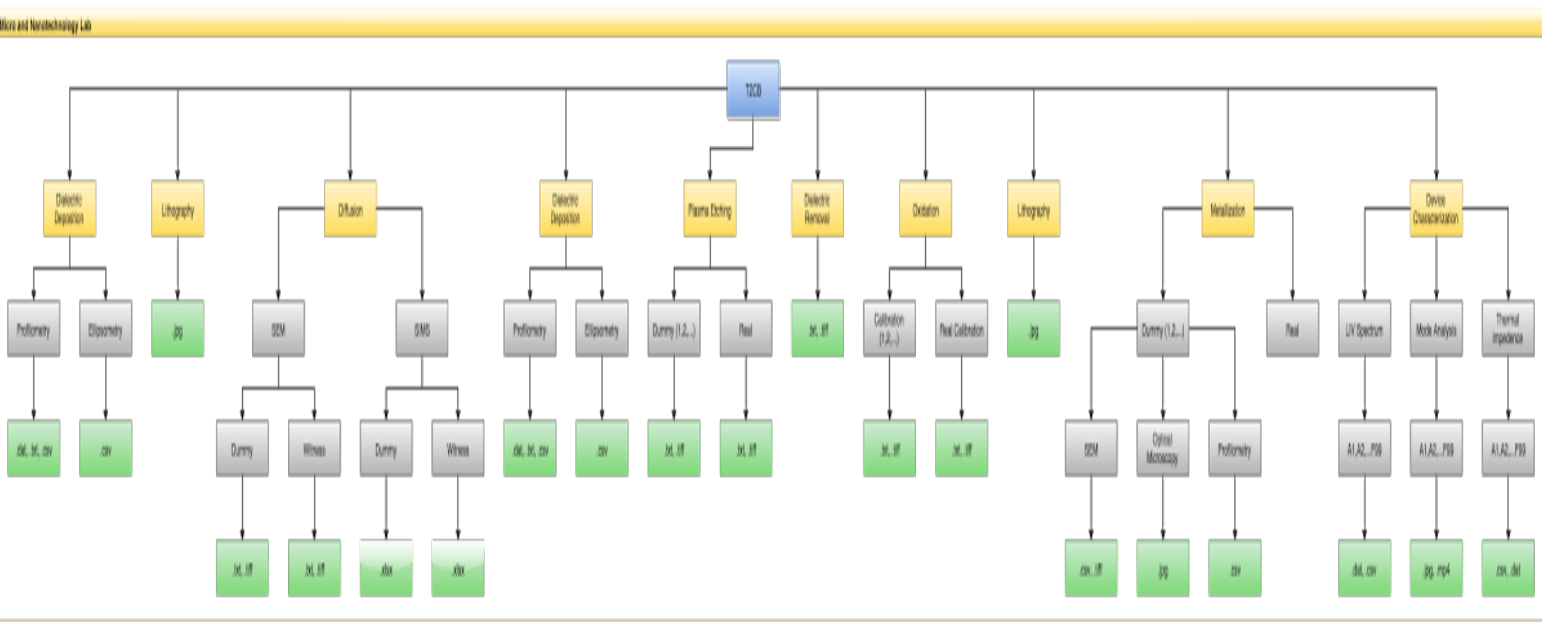Temp 425 C

(Structured meta data)

**Notes:**

Oxidation depth is about 12um. Oxidation layer composed of Al(0.98)GaAs with thickness of 30 nm. Furnace in 2111 MNT L, 2" diameter quartz.

(Free text)

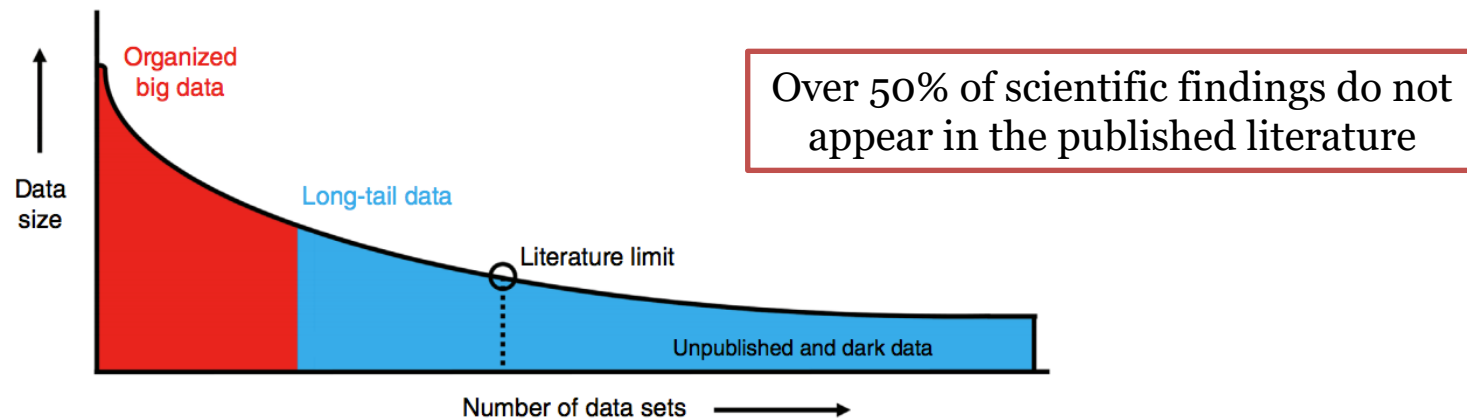A lot of useful information is hidden in unstructured text

Example of multimodal experimental



Heterogeneity of experimental data (Spaces, Collections Datasets)

4CeeD

# 4CeeD Design Considerations - 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
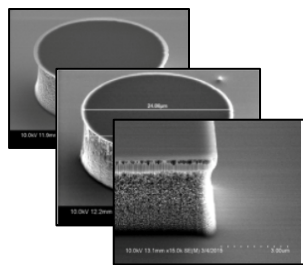  - "Dark data", e.g., unpublished data of failed experiments

Over 50% of scientific findings do not appear in the published literature
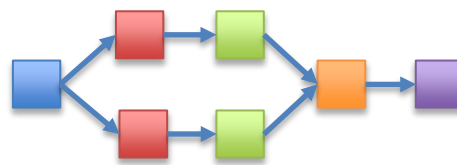
Long-tail scientific data

*Source*: Ferguson et al. *Nature neuroscience* 17.11 (2014)

# 4CeeD Design Considerations - Long-tail scientific data processing challenges

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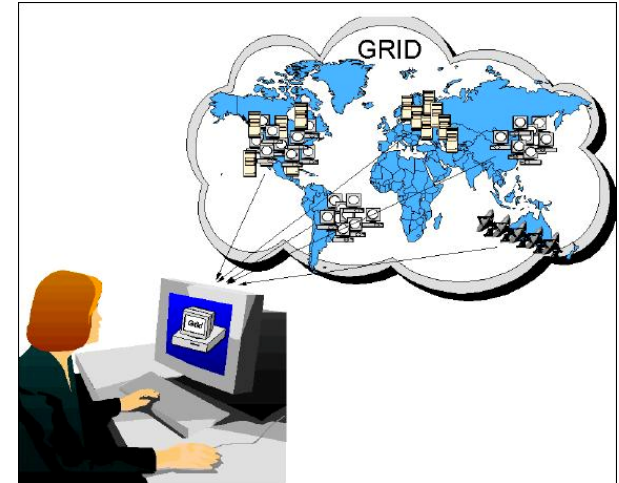


Data analysis workflow       Insights

Raw data

- 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

# 4CeeD Design Considerations – Task Workflows

- Application is a Computational Workflow

- Workflow is Set of Tasks (e.g., A, B, C, D) executing over materials data



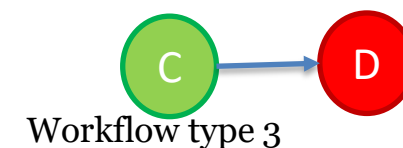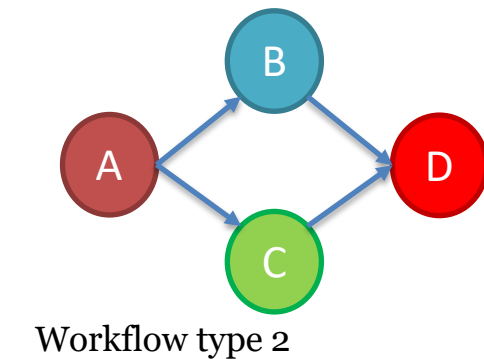Workflow type 1

1. Example of a Task C: "Plotting a graph"

```
In [5]:   metadata = py4ceed.get_metadeta()
          metadata.plot(x='Pressure', y='Etch_Rate')
          plt.show()
```

2. Example of a Task D: "Filter Data"

**In [6]: metadata[metadata ['Pressure'] >=7]**

- Other examples of tasks: Extraction of features from an image, compression of image, …



Workflow type 2



Workflow type 3

# Summary of 4CeeD Design Challenges



➢ 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

# Outline

- 4CeeD Distributed Architecture, Backend Cloud Concepts and Services
  - What is 4Ceed and its goals
  - What is behind the 4CeeD Dashboard
  - 4CeeD Cloud Design and Deployment
  - How to deal with Aging Scientific Instrument

4CeeD

# 4CeeD Cloud Design

- ✓ Cloud Concept

- ✓ Micro-service execution environment

- ✓ Data Management

# Cloud Computing Concept



Figure Source: Wikipedia

4CeeD

ILLINOIS

17

# Cloud Computing Concept



Cloud computing

Cloud Clients
Web browser, mobile app, thin client, terminal emulator, ...

SaaS
CRM, Email, virtual desktop, communication, games, ...

PaaS
Execution runtime, database, web server, development tools, ...

IaaS
Virtual machines, servers, storage, load balancers, network, ...

**4CeeD Cloud**

# Private and Public Clouds



4CeeD is Private Cloud

Hybrid

Private/Internal

Public/External

The Cloud

On Premises / Internal

Off Premises / Third Party

Cloud Computing Types

CC-BY-SA 3.0 by Sam Johnston

Figure Source: Wikipedia

# Example of Cloud Components



Figure Source: Wikipedia

4CeeD

ILLINOIS

# 4CeeD Cloud Components



4CeeD uses **Queue** RabitMQ)

4CeeD is **Web-based**

Cloud Service (eg Queue)

Cloud Platform (eg Web Frontend)

Cloud Infrastructure (eg Billing VMs)

Cloud Storage (eg Database)

4Ceed uses **Containers And Micro-service** Infrastructure (Docker and Kubernates/Docker Compose)

4CeeD uses **Database** (MongoDB) and Filesystem (HFS)

Figure Source: Wikipedia

# Hardware Virtualization

- Two types of hardware virtualization
  - Emulation-based virtualization
  - Container-based virtualization

| App 1 | App 2 | | App 1.1 App 1.2 | App 2.1 App 2.2 |
|---|---|---|---|---|
| Guest OS 1 | Guest OS 2 | | | |
| Hypervisor | | | Container Engine | |
| Host OS | | | OS | |
| Hardware | | | Hardware | |

# Container

- ## Container – Software Unit that bundles its own software, libraries and configuration files

  – Containers are isolated from one another and can communicate with each other through well-defined channels.

  – All containers are run by a single operating system kernel and therefore use fewer resources than virtual machines.

  – Virtual Container, called Docker, is professional software package developed by *Docker Inc*. as part of PaaS.



Source: Wikipedia

# Micro-Service

- Microservice
  - a software development technique (a variant of the service-oriented architecture (SOA) structural style)
  - an application is arranged via microservices as a collection of loosely coupled services.
- In a microservices architecture, services are _fine-grained_ and the protocols are _lightweight_.

# 4CeeD Cloud Architecture Components – Putting it Together



Web Browser (Client App)

Web Browser (Client App)

Client Side

Server Side

Web Server Application

**4CeeD-Clowder Data Management and Storage**

**MongoDB**

**File System**

**4CeeD Cloud Execution Environment** (based on Microservice Architecture)

4CeeD

ILLINOIS    25

# 4CeeD Cloud Design

✓ Cloud Concept

✓ Micro-service execution environment

✓ Data Management

# In Cloud - Micro-service execution environment

- *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

# 4CeeD Executing scientific data processing workflow



**4CeeD Cloud Execution Environment**

Curation Service

Upload raw scientific data

Return curation results

Workflow Invoker

TDS: "Start at A"

TDS Server ... TDS Server

TDS Server

TDS Ensemble

Task Graph

Start — A — B — End

Task A's Request Queue

Round-robin dispatching

Task A's Consumer

Task A's Micro-service

TDS: "Next task is End"

TDS: "Next task is B"

Task B's Request Queue

Task B's Consumer

Task B's Micro-service

**TDS** – Task Dependency Service
**Task A: Example** - Color Extraction from TEM Image
**Task B: Example** - Detect Edge in TEM Image

4CeeD

ILLINOIS

28

# 4CeeD Cloud Design

✓ Cloud Concept

✓ Micro-service execution environment

✓ Data Management

# 4CeeD Data Management and Storage

- 4CeeD uses NoSQL database to store spaces, collection and dataset metadata and some data

- MongoDB is open-source NoSQL database

  – Non-relational database (NoSQL), i.e., data storage and retrieval are not organized in tabular relations

  – Developed due to the limits of relational databases and their scalability to very large datasets (scale was limited because of the requirement for consistency in relational databases)

  – 4 models of NoSQL

    • key-value stores,

    • graph stores,

    • column stores,

    • document stores

4CeeD

# 4CeeD-Clowder Data Management and Storage (2)

- **Document Store Model**
  - Store data in semi-structured form, called documents
  - Documents encoded in standardized format such as
    - XML format
    - Javascript Object Notation (JSON)
- **Example of Document store database**

**Collection of Documents**

```
{
"firstName":"John",
"lastName":"Smith",
"age": 45
}
```

```
{
"firstName":"Helen",
"lastName":"Second",
"age": 54
}
```

```
{
"firstName":"Robert",
"lastName":"Third",
"age": 23
}
```

Source: P. Bajcsy et al. "Web Microanalysis
Of Big Image Data", Spring, 2018

4CeeD

ILLINOIS

31

# 4CeeD Data Management and Storage (3)

- 4CeeD uses MongoDB

- In MongoDB
  - Documents are stored in a JSON-like format

- Example of JSON-like Format

- 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.

```
{ "first name": "John",
"last name": "Smith",
"age": 25,
"address": {
    "street address": "21
2nd Street",
    "city": "New York",
    "state": "NY",
    "postal code": "10021"
},
"phone numbers":[
    {
        "type": "home",
        "number": "212 555-
1234"
    },
    {
    "type": "fax",
    "number": "646 555-4567"
    }
],
    "sex":
    {
        "type": "male"
    }
}
```

Source: wikipedia

# 4CeeD Smart Data Management

**Collection:** T2CB; **Datasets**: PlasmaEtching, …., Metalization
**Folders:** Calibration, SEM, Optical Microscopy…, **Files**: txt files, tiff files, …



4CeeD

ILLINOIS

# 4CeeD Deployment – Cloud Production System

## 4CeeD Cloud

**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)



Storage Layer (GlusterFS):
- StorageX R730xd
- StorageY R730xd
- Future Expansion
- Storage1 R730xd
- Storage2 R730xd
- Replicated

GlusterFS Share (10 Gb)

Compute Layer (Kubernetes):
- Compute1 (Master) R630
- Compute2 R630
- Future Expansion
- ComputeX R630
- ComputeY R630

# 4CeeD Micro-service implementation system (in Compute Layer)



Micro-service execution layer

Infrastructure layer

# Outline

- 4CeeD Distributed Architecture, Backend Cloud Concepts and Services
  - What is 4Ceed and its goals
  - What is behind the 4CeeD Dashboard
  - 4CeeD Cloud Design and Deployment
  - How to deal with Aging Scientific Instrument

4CeeD

# Current situation in campus cyberinfrastructure



**Private Cloud (4CeeD)**

**Campus Network**

Building L2 high speed switch

Building L2 high speed switch

Building L2 high-speed switch

Building L2 high speed switch

New instruments

**Older instruments (offline)**

New instruments

4CeeD

ILLINOIS

# Challenges of connecting offline older instruments

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

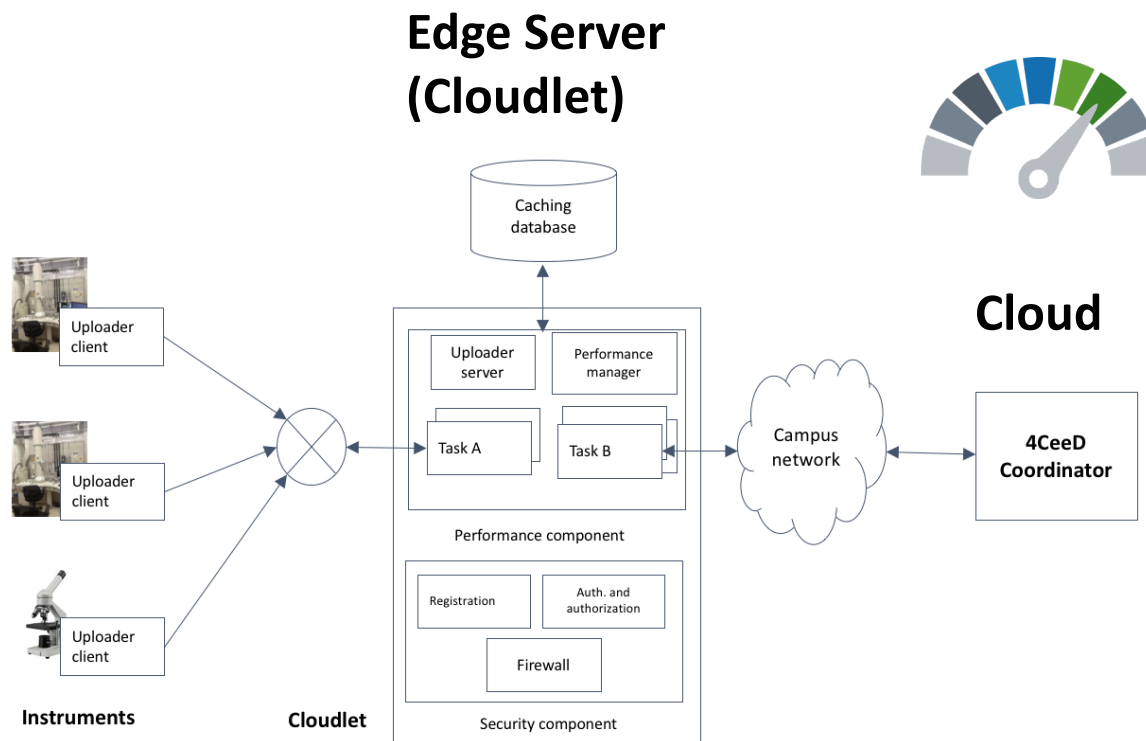- **Obsolete security**: Older devices and their OS systems cannot be patched, hence being vulnerable & taken offline

# BRACELET: Putting edge device between older instruments and private cloud



**Edge Server (Cloudlet)**

**Cloud**

BRACELET in 3-tier architecture

**Performance**:

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

**Security:**

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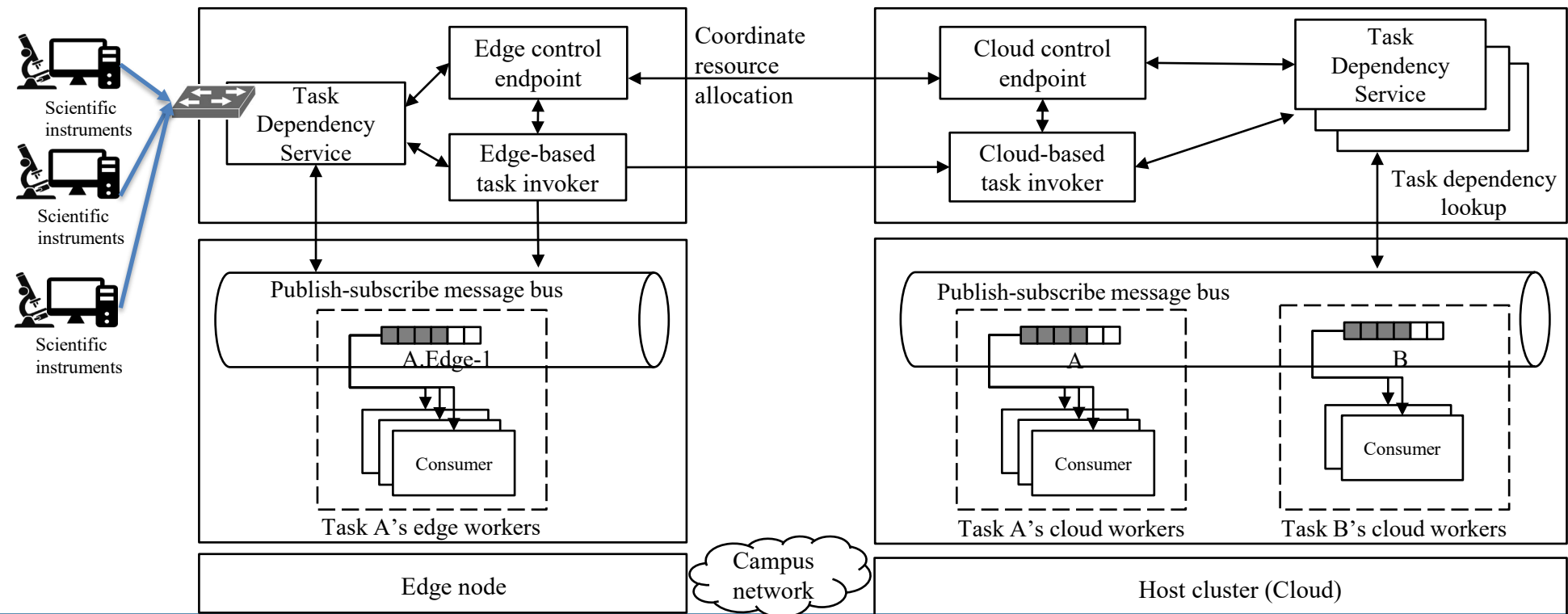– Firewall to protect against external threats

4CeeD

ILLINOIS

39

# BRACELET Design

## Edge Server

- Security service
  - Check equipment address
  - Authenticate user and his reservation
- Compute/Transport service
  - Forward and upload data

## Cloud

- Compute/Data service
  - Compute tasks/workload
  - Store/Retrieve metadata, data
- Security service
  - Authenticate user, access control

# User authentication from instruments via BRACELET



**Bracelet Edge Server**

**4CeeD Cloud**

- 4CeeD Uploader checks Reservation DB to ensure the user has reserved his/her time

- User authenticates with AD to access the instrument computer
- Then, user authenticate with 4CeeD Uploader via Shibboleth, or via 4CeeD Curator authentication APIs

**4CeeD**

# Transport service between edge & cloud
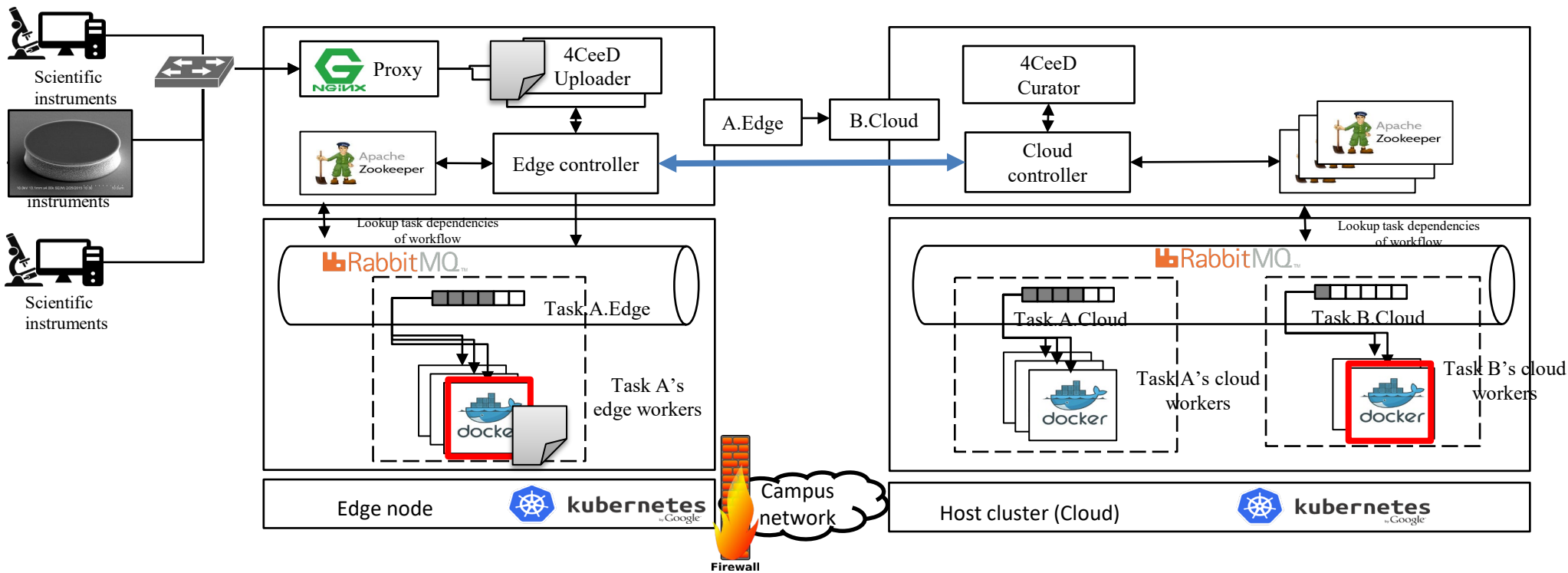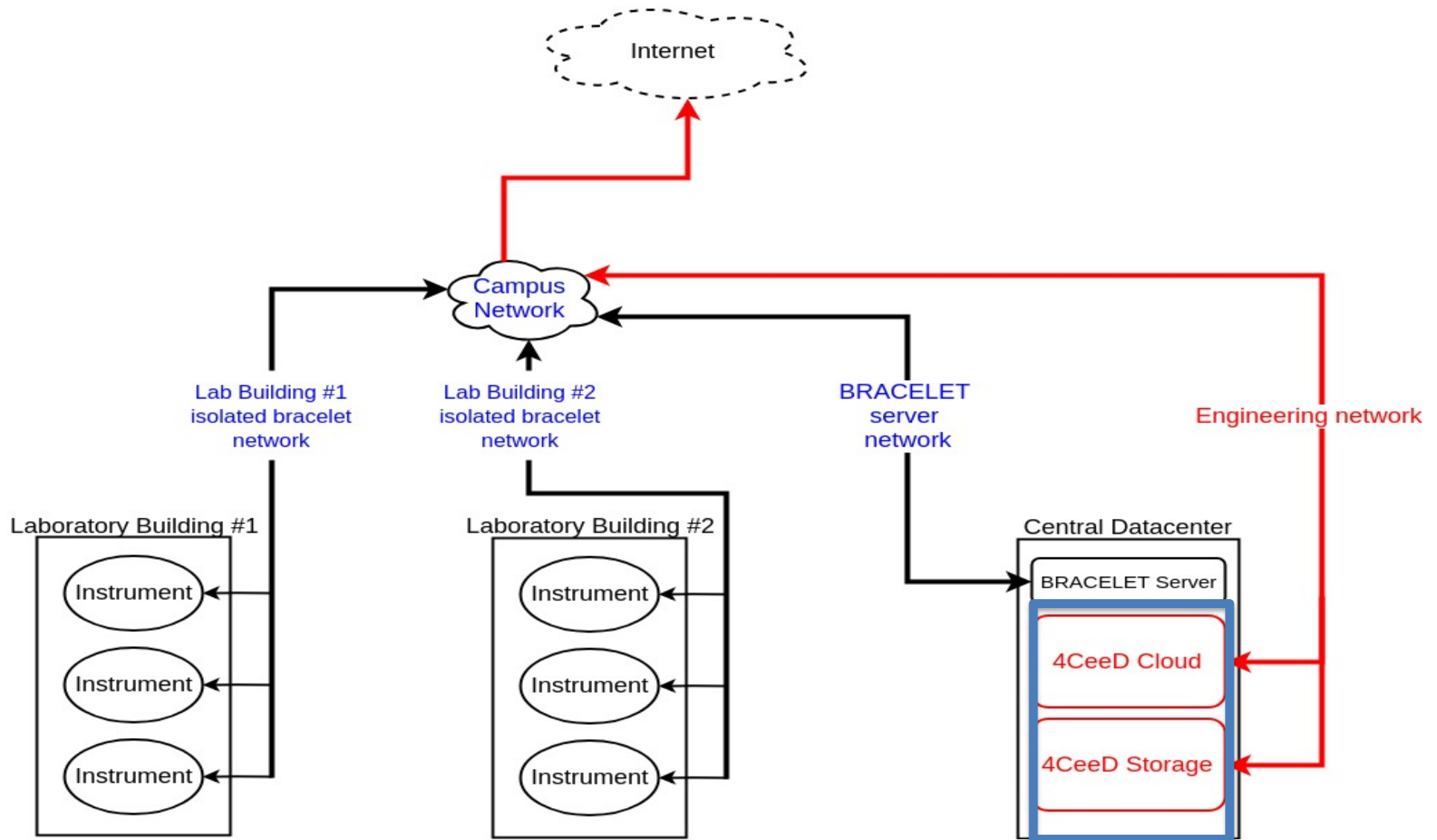
- After processing request, the task consumer forwards request to the next task (following current placement)

- After learning about the placement, data processing request is sent to the first task

- 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

# BRACELET Deployment



BRACELET Network Architecture

# 4CeeD Summary

- Lightweight microservice cloud architecture for materials genomic challenge

- Real-time cloud service for
  – Curation Service
  – Data Analysis (Jupyter Notebook)

- Smart data management system for materials data

- Novel usage of edge computing for aging IoT devices to enable security

- Sources (code and project description):

- https://4ceed.github.io/

- http://t2c2.csl.illinois.edu/

# Publications

- 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

- 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; DOI: 10.1109/ICDCS.2019.00021

- 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 Microservice Infrastructure for Aging Scientific Instruments", **IEEE International Conference on Computing, Networking, and Communications (ICNC)** 2019, Hawaii, February 2019.