# Data Lifecycle

Anirban Mandal,
Assistant Director of Network Research and Infrastructure,
RENCI, UNC – Chapel Hill
http://nrig.renci.org

# Data Lifecycle (DLC) for Scientific Facilities



| Capture | Initial Processing | Central Processing | Archiving / Storage | Dissemination |
|---|---|---|---|---|
| some type of sensor or instrument (e.g. GRAPEs, telescope, DOMs) | often at the sensor site, or nearby | main data center — secondary data center(s) | main data center — secondary data center(s) | scientists / public |

Different forms of transmission (e.g., plane, satellite, cables), redundant networks, ...

**Data Movement**

**Cross-cutting CI elements (e.g., Disaster Recovery, Identity Management)**

# Data Lifecycle (DLC): Goals

- Understand and document the cyberinfrastructure (CI) best practices and solutions for data life cycle (DLC) for Large Facilities (LFs).

- Can a generalized DLC abstraction help us understand the diverse CI landscape for LFs ?

  - Can it be ***ONE way to learn/catalog the CI functionalities*** at each stage of data operation for LFs ?

  - What *services are offered* by each DLC stage ?

  - What *CI architectural elements* support each DLC stage ?

- Study the end-to-end life cycle for data as it traverses different CI entities inside a LF and then catalog the underlying services/tools/platforms for the life cycle stages.

# Data Lifecycle (DLC): Participant Questions

- Best-practices for getting researchers to curate their data, in particular taking it off of working storage attached to HPC
- I want to learn the security check points needed for protected data
- How to get researchers to pay for the proper data management over the life cycle?
- I saw Anirban's presentation on NEON Data Lifecycles and subsequently started working on data lifecycle documentation for RCRV.  Would be great to check in with Anirban and review the RCRV Lifecycle docs looking for problem areas or gaps.

# Data Lifecycle (DLC): Other Relevant Questions

- Does the DLC figure capture all the steps that data goes through in your organization? Are any critical stages missing?

- What technical functions or tasks are performed at each DLC stage?

- Can you describe how the data flows between stages in your DLC?

- Which DLC stages or cross-cutting elements (e.g., disaster recovery, identity management) do you experience challenges with and might appreciate external expertise/help with?

- Do you use academic resource providers (OSG, XSEDE, National Laboratories etc.) ?

- Do you use Cloud technologies to support DLC stages (computing/storage/backup) ?

- Do you use any data movement services (Globus, CloudConnect etc.) ?

# Questions?

Contact: anirban@renci.org

**See you at 1PM EST for CI/CS Workshop's Welcome and Introductions**

CI/CS WORKSHOP · ResearchSOC | CI CoE PILOT

Backup Slides

# IceCube Facility DLC: Data Capture

# IceCube Facility DLC: Initial Processing/Filtering

**Initial processing, filtering: South Pole**

- Data is received by DOMHubs in the IceCube Lab (surface of South Pole)

- ~500 core filtering cluster; ~100 machines for detector readout

- Hits are output as events. Internal PnF system selects events based on their usefulness for a particular analysis. It also creates event metadata and reduces data volume before it is transmitted away from the South Pole.

- **Alert production** is an important process that happens in this stage of the DLC.

# IceCube Facility DLC: Central Processing

**Central processing: UW-Madison** processes what is sent from the South Pole to a "science ready level" up to level 3.

- UW-Madison: 7600 core, 400 GPU cluster, ~10 PB storage. PFFILT → L2 and L3.

- Additional downstream processing happens using a mix of resources: DESY, OSG, IceCube Grid (campus clusters, contributed resources, etc.), XSEDE allocations, DOE resources (e.g. NERSC).

- Increased demand for GPU resources.

- PyGlidein + HTCondor based distributed computing middleware.

- Exploring cloud resources for CPU, GPU, ML.

# IceCube Facility DLC: Data Movement

1. Hits at DOMs → DOMHubs → Data Acquisition System (DAQ) → Events (PFRAW)
2. Sent to Processing and Filtering System (PnF) - PFRAW made ready for analyses
3. Sent to South Pole Station JADE for archival storage to disk (PFRAW and PFFILT/Level 1)
4. JADE transmits via satellite to UW-Madison (PFFILT)
5. PFFILT sent to DESY and PFRAW sent to NERSC for additional tape backups

In addition, Alerts are sent out using GCN (Gamma Ray Coordination Network - operated by NASA) or Astronomical telegrams along with initial estimate of PFRAW data sample via satellite link to UW

- Limited bandwidth of ~125 GB/day from South Pole to UW; 3TB/day raw data is filtered down to ~80GB/day and transmitted via satellite from South Pole Station to UW
- Once a year, raw data from the South Pole is sent via plane, boat in disks to UW
- UW connected to SciDMZ through Starlight-ESNet for connection to DOE facilities
- Leverages GridFTP for data transfers from UW-Madison to DESY/NERSC/OSG

# IceCube Facility DLC: Data Storage/Archiving

**JADE (archival system)** exists in ~3 locations

- South Pole JADE - writes 2 copies to disk (3 TB/day)

- JADE North (UW) - warehouses the data to disk (~200 TB/yr)

- JADE Long Term Archive (LTA) in DESY – keeps replicas of Level 1 and 2 data

NERSC archives PFRAW (the raw data) in tape archives.

# IceCube Facility DLC: Access/Publishing/Distribution

**Dissemination of Alerts**

- Alerts happen at the South Pole during Level 1 processing.

- Alerting systems detect events and then an immediate alert is sent out using GCN (Gamma Ray Coordination Network - operated by NASA) or Astronomical telegrams along with initial estimate/small portion of PFRAW data sample via satellite link to UW.

- When a full PFFILT (Level 1) data set is available at UW later, a refinement of the first alert is sent.

researchSOC | CICoE

# IceCube Facility DLC: Access/Publishing/Distribution
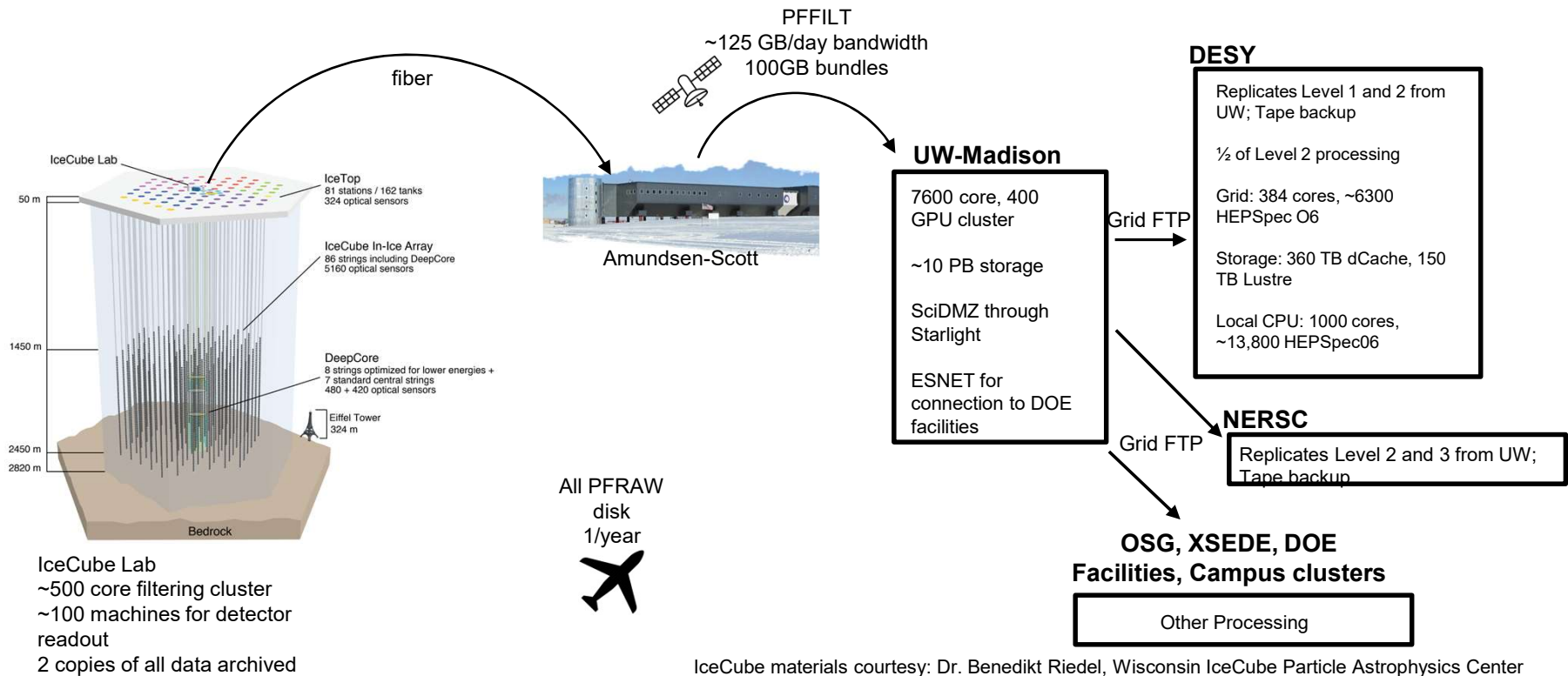
3 forms of **data access** for other types of data:

- Be a member of the IceCube collaboration
- Be an "associate member" – one applies for use of the data for a particular purpose but is not required to fulfill collaboration obligations
- Public web portal

Planned enhancements for data organization, management, access, and data catalog

- Xrootd-based solution, Ceph/www

Data is released to members and associates. When the data has been analyzed and those analyses published, it becomes available for release to others.

ResearchSOC | CICoE

# IceCube Facility DLC: Logical Architecture



IceCube Lab
~500 core filtering cluster
~100 machines for detector readout
2 copies of all data archived

PFFILT
~125 GB/day bandwidth
100GB bundles

All PFRAW
disk
1/year

Amundsen-Scott

**UW-Madison**

7600 core, 400 GPU cluster

~10 PB storage

SciDMZ through Starlight

ESNET for connection to DOE facilities

Grid FTP

**DESY**

Replicates Level 1 and 2 from UW; Tape backup

½ of Level 2 processing

Grid: 384 cores, ~6300 HEPSpec O6

Storage: 360 TB dCache, 150 TB Lustre

Local CPU: 1000 cores, ~13,800 HEPSpec06

**NERSC**

Replicates Level 2 and 3 from UW; Tape backup

**OSG, XSEDE, DOE Facilities, Campus clusters**

Other Processing

IceCube materials courtesy: Dr. Benedikt Riedel, Wisconsin IceCube Particle Astrophysics Center

researchSOC | CICoE

# Heading

- Bullet 1

- Bullet 2

- Bullet 3

# Heading

- Bullet 1

- Bullet 2

- Bullet 3