# PNNL DataHub: Building a Central Data Capability

May 27, 2022

**Michael Hofmockel**
**Ian Smith**
**Shannon Sheridan**
**Miriam Blake**
Research computing, PNNL

# Problem

**Estimate "data of record" creation at ~10PB/2,000 projects/year**
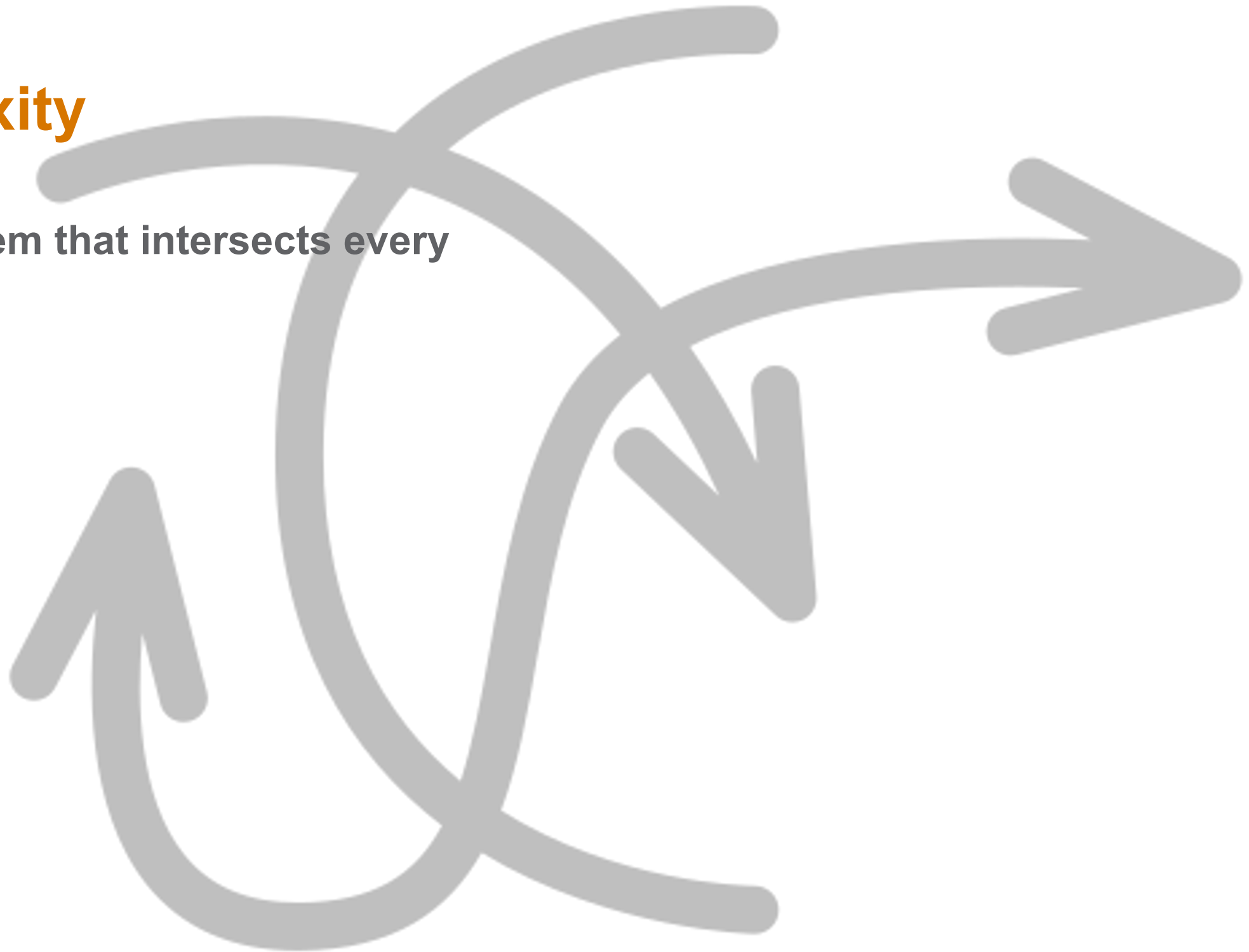
- Institution cannot enumerate its research data

- Projects roll their own solutions/policies

- Preservation is limited to project lifetime and PI knowledge

- Commercial data management solutions are focused on back-office data that doesn't translate to research data.

# Complexity

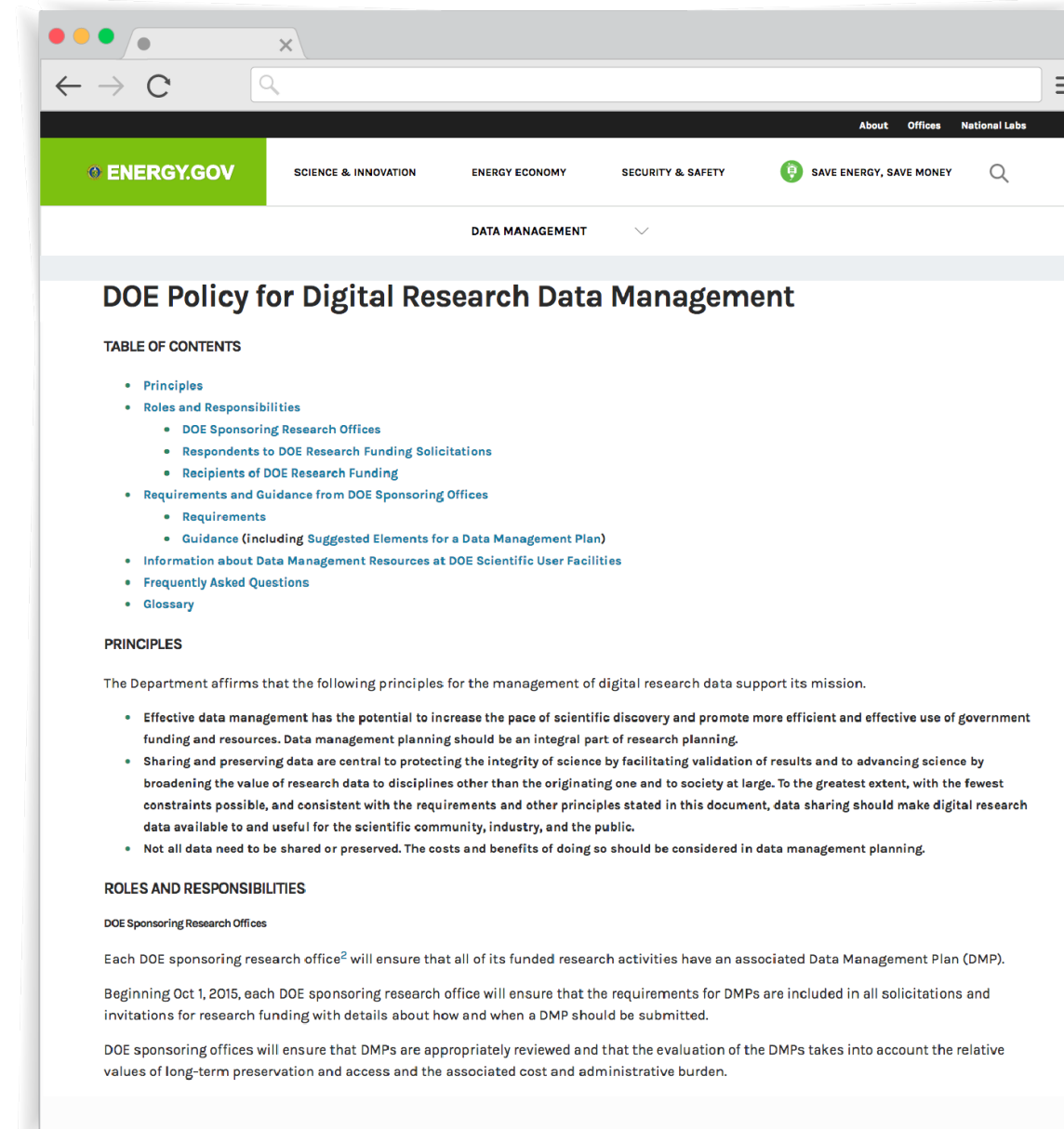**Multi-dimensional problem that intersects every part of the institution**

- Commitment
- Understanding
- Standards
- Implementation
- Over generalization
- Education
- Training
- Sustainability

# Our Sponsor's Policy Should Be Our Policy

*"Sharing and preserving data are central to protecting the integrity of science by facilitating validation of results and to advancing science by broadening the value of research data to disciplines other than the originating one and to society at large."*

**- Department of Energy**

*DOE Policy for Digital Research Data Management*

# Research Data Management Vision

*Increase the impact and scientific value of our research by elevating research data to be a principal PNNL product.*

This will require transformation of our institutional culture, policies, and infrastructure to fully support the Research Data Life Cycle.
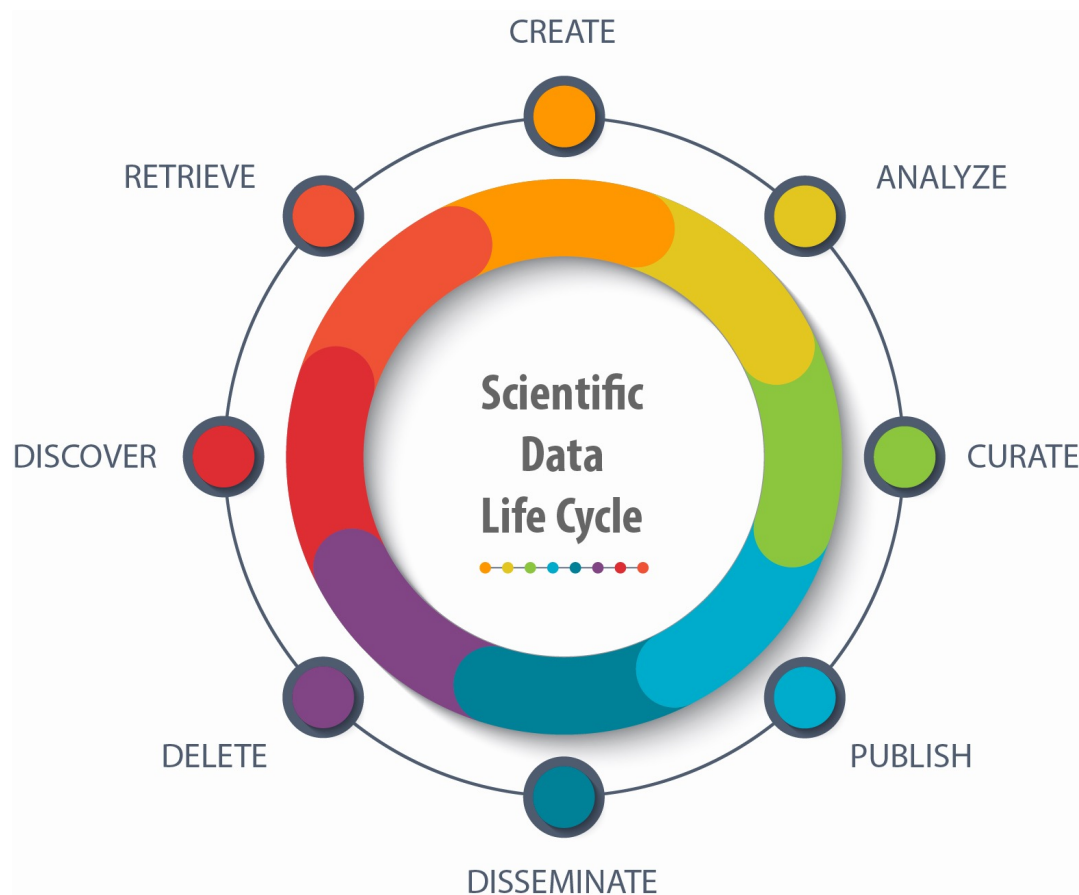
# DataHub Advance Team - DATa

**Data stewardship** is the practice that ensures data is accessible, usable, safe, and trusted. It includes overseeing every aspect of the data lifecycle, promoting data quality and integrity.

The **DataHub Advance Team** is the source for data stewardship support within Research Computing at PNNL. Utilizing DataHub, the Advance Team helps researchers address their institutional projects' data stewardship. Bringing together a wide variety of subject matter experts, the team can provide curated assistance.

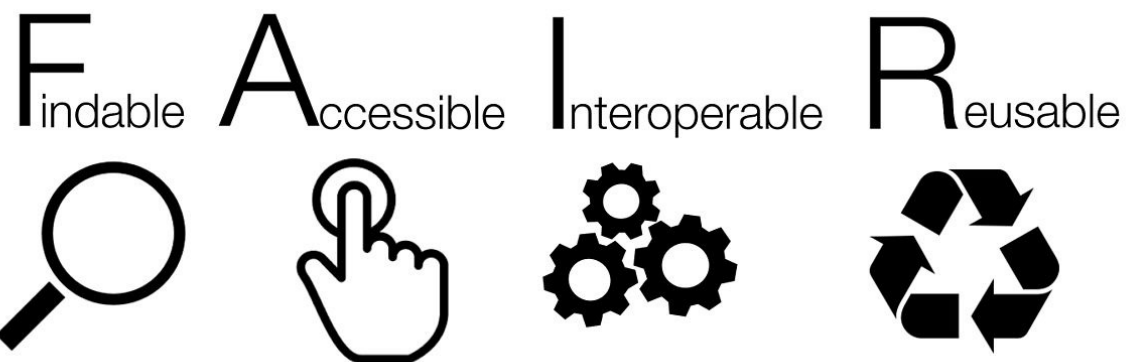Research groups around the lab have utilized the Advance Team to:
- Empower project team data stewards to create data policies for entire projects
- Develop data catalogs to track all data outputs from multiple collaborators internally
- Create project-wide naming conventions and directory structures
- Deploy a public web presence that semantically links data, publications, research staff, code, and scientific instruments

# PNNL DataHub Mission



Scientific Data Life Cycle

CREATE · ANALYZE · CURATE · PUBLISH · DISSEMINATE · DELETE · DISCOVER · RETRIEVE

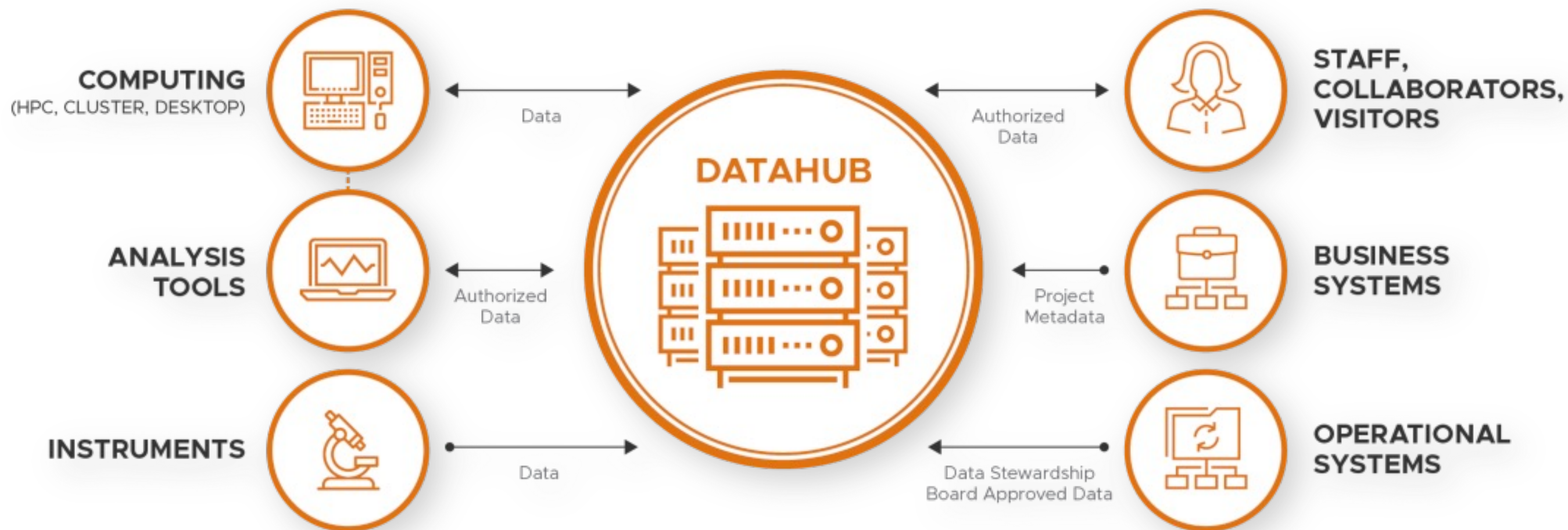**F**indable **A**ccessible **I**nteroperable **R**eusable

- Singular Registry for ALL data
- Easy and secure access
- Interoperable and standards based
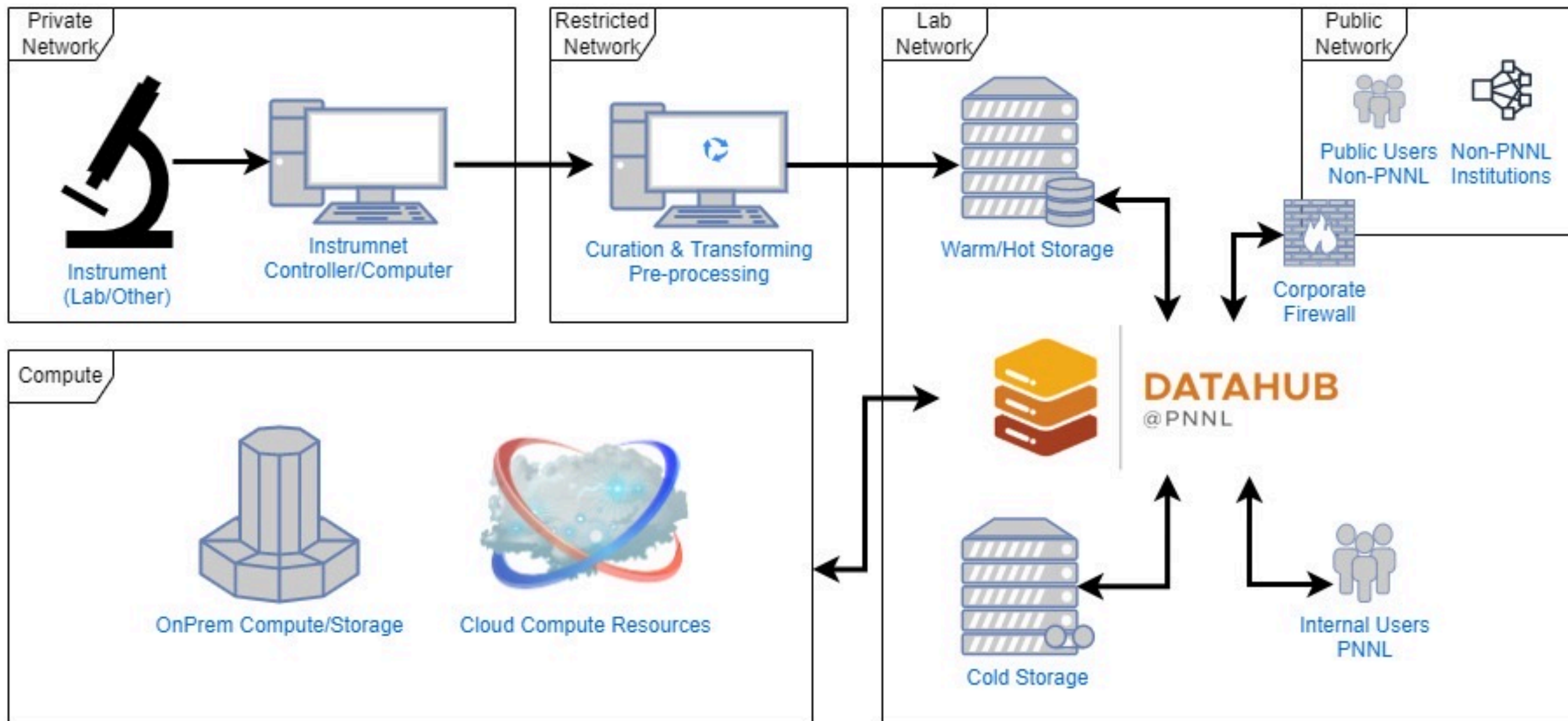- Research is reusable AND **reproducible**

**Research Data Management and data literacy is embedded in institutional process and culture**
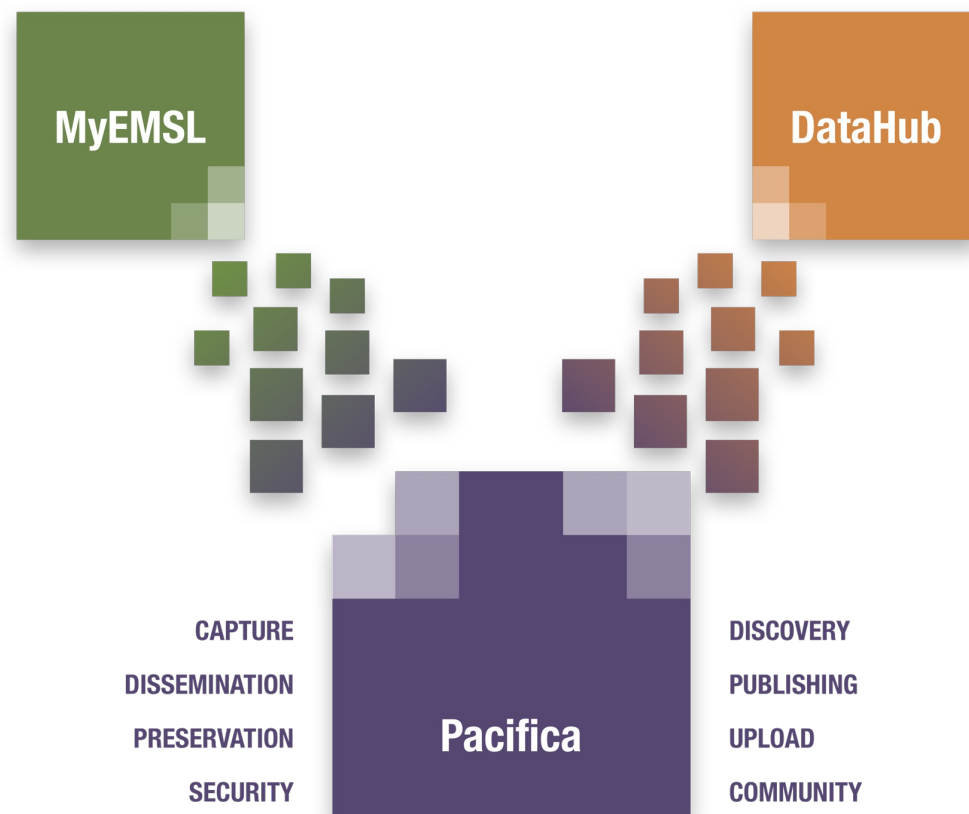
# DataHub System at a glance

# Example Data Pipeline

# Lessons Learned



- Change is difficult
- Everyone wants access NOW!
- Starting simple is ok
- Document decisions and communicate
- Agree that a problem exists
- Commit to change
- Start a community around data

# What are we working on now…

- Automation
- Standards
- Workforce
- Education
- Community
- Influencing Behaviors

# Thank you

This research was funded by the National Nuclear Security Administration, Defense Nuclear Nonproliferation Research and Development (NNSA DNN R&D). The authors acknowledge important interdisciplinary collaboration with scientists and engineers from LANL, LLNL, MSTS, PNNL, and SNL.