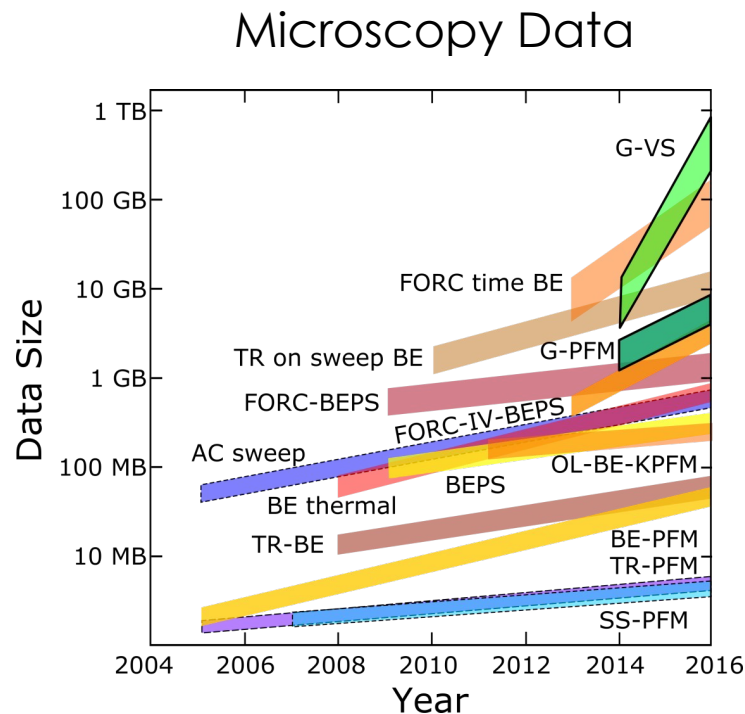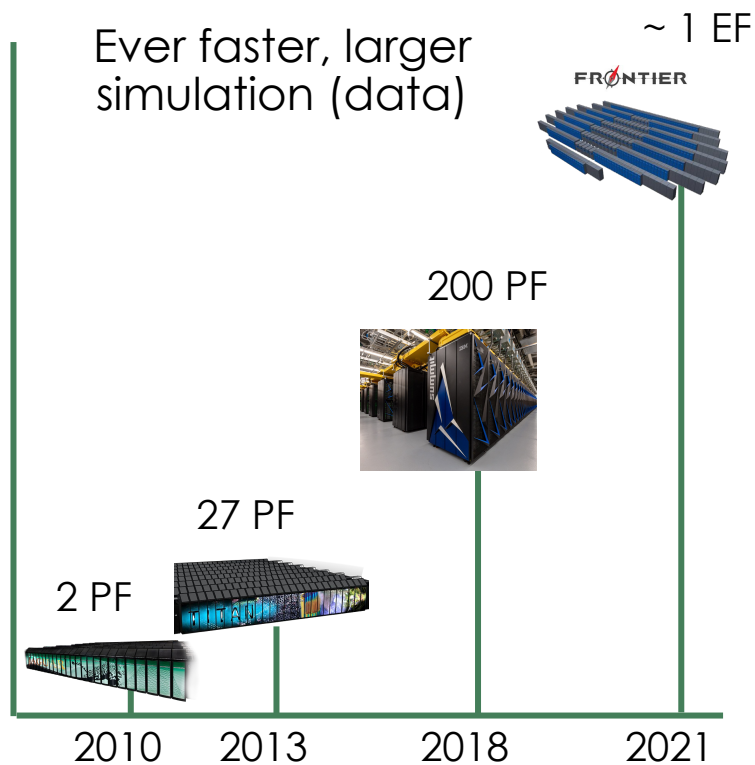# A Systemic Approach to Facilitating Reproducibility via Federated, End-to-End Data Management

Dale Stansberry, Suhas Somnath, Gregory Shutt, and Mallikarjun Shankar

Advanced Data and Workflows Group
Oak Ridge Leadership Computing Facility
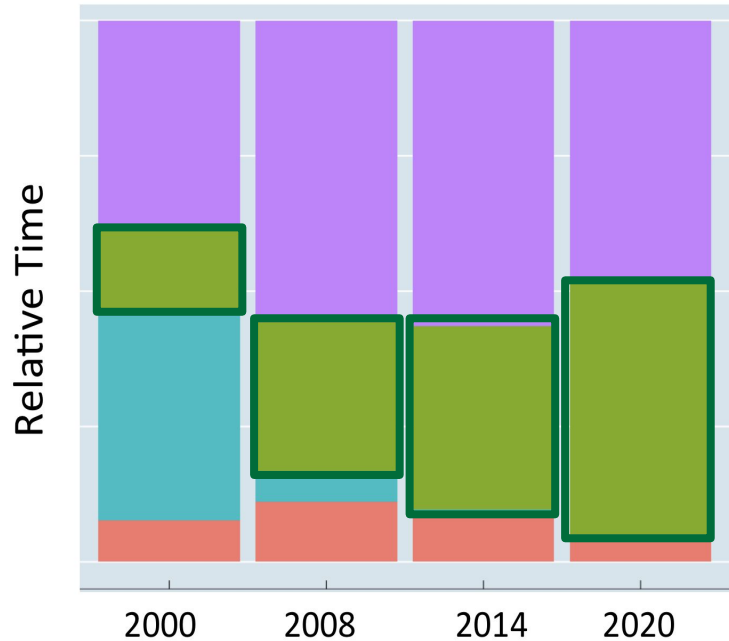
# Explosion in Data Volume, not (necessarily) Quality



Ever faster, larger simulation (data)

~ 1 EF

2 PF — 2010
27 PF — 2013
200 PF — 2018
~ 1 EF — 2021

Microscopy Data

Kalinin et al., *ACS Nano*, 9068-9086, 2015

# Explosion in Time Spent on Data Management

**Experimental Time Breakdown**

■ Experiment Design ■ Management ■ Measurements ■ Post Processing

Relative Time

2000    2008    2014    2020

* at light sources

Figure from Kevin Mader
https://rawgit.com/4Quant/SRI2015/master/SRIPres.html#/
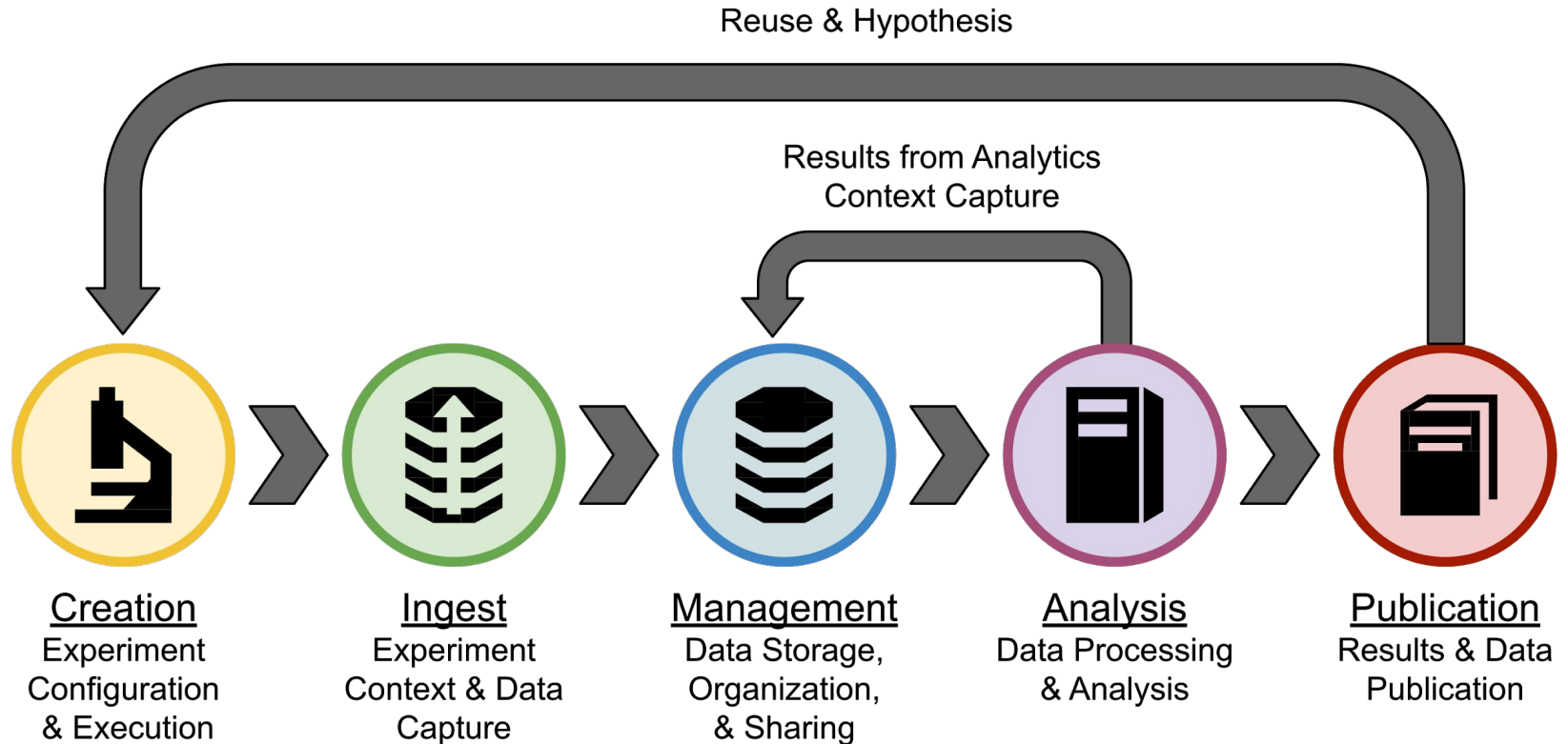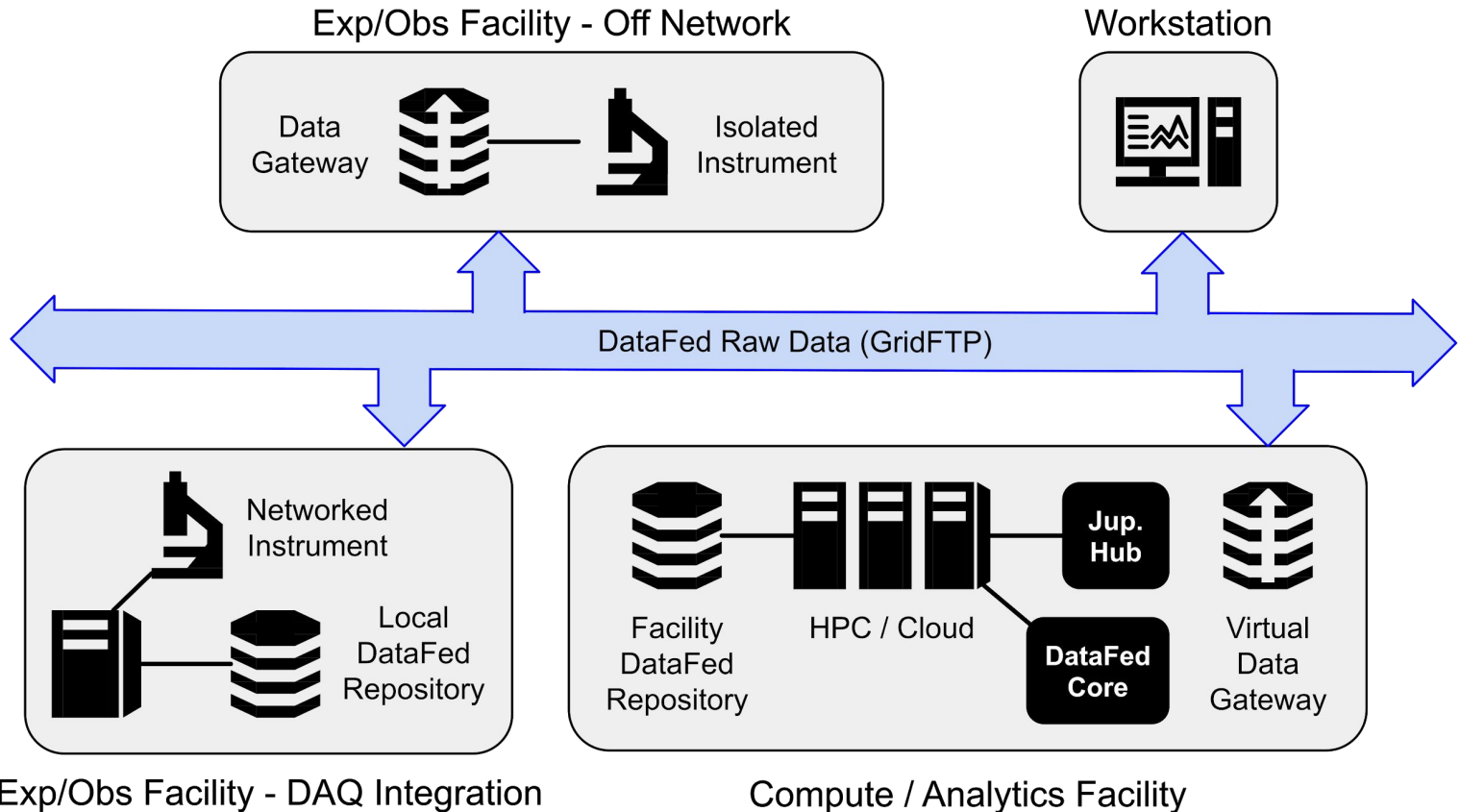
# Lack of Data Infrastructure & ad-hoc Practices

- Metadata inadequately captured when generating data
  - Physical / electronic lab notebooks - rarely reconciled / findable
- Filesystem for data management
  - Metadata embedded into file paths
  - Sharing, searching, organization
- Thousands of data and metadata formats
- Emails, dropboxes, and portable storage for sharing data
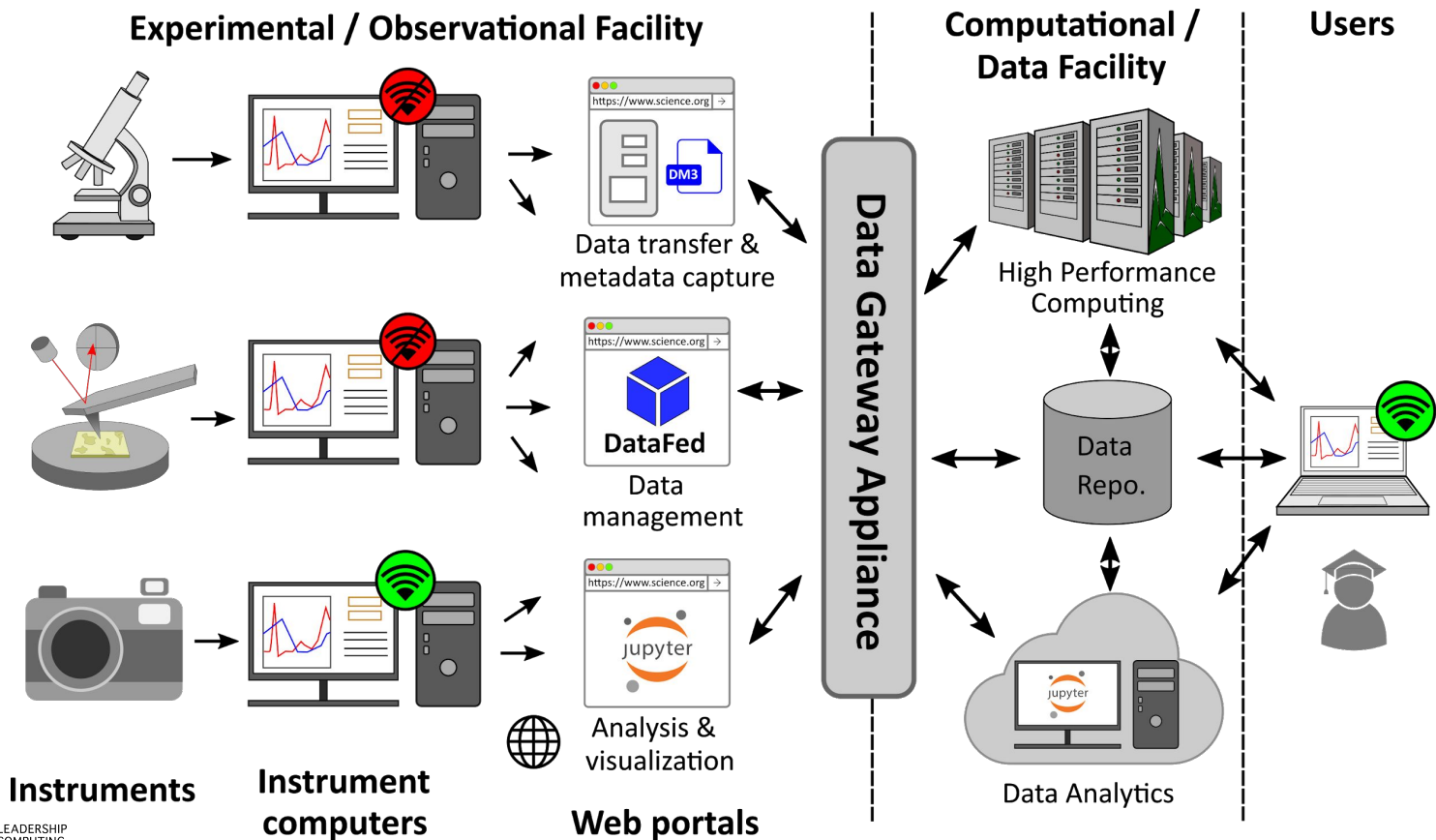- Dark data - majority of data never published / shared

OAK RIDGE | LEADERSHIP
National Laboratory | COMPUTING
FACILITY

# Lifecycle of Scientific Datasets



Reuse & Hypothesis

Results from Analytics
Context Capture

**Creation**
Experiment
Configuration
& Execution

**Ingest**
Experiment
Context & Data
Capture

**Management**
Data Storage,
Organization,
& Sharing

**Analysis**
Data Processing
& Analysis

**Publication**
Results & Data
Publication

# Systemic Approach to Data Management

Exp/Obs Facility - Off Network

Workstation

Data Gateway — Isolated Instrument

DataFed Raw Data (GridFTP)

Networked Instrument

Local DataFed Repository

Exp/Obs Facility - DAQ Integration

Facility DataFed Repository

HPC / Cloud

Jup. Hub

DataFed Core

Virtual Data Gateway

Compute / Analytics Facility

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Data Gateway - Data Ingest Tool

# Data Gateway - Future Development

- Containers for data pre-processing
    - Extracting metadata, translating from proprietary formats, etc.
    - Isolates science code from Gateway
    - Repository of vetted codes / containers
- Virtual Gateway - Centralized deployment for users outside facilities
- Exposed API for automated ingest of data - long experiments
- Remote control / monitoring of instruments
- Trigger compute jobs for data processing on HPC

**OAK RIDGE** | LEADERSHIP COMPUTING FACILITY
National Laboratory

# DataFed - Scientific Data Management System

- Data handling application
- Unique identifiers for each record, collection, user, etc.
- Abstracts file system complexities (directories, paths, files, ….)
- Indexes metadata
  - General - author, creation date, size, provenance, etc.
  - Domain-specific
- Powerful searches
- Fine-grained access control
- Globus for:
  - Data movement (Grid-FTP)
  - Authentication (federated identity management)

**OAK RIDGE**
National Laboratory | LEADERSHIP COMPUTING FACILITY

# DataFed Core Services and Example Repository

## DataFed Core Services



DataFed Core Services

API

Web App

DataFed Core Service

Metadata DB

Globus

DataFed Repository Service

Data Store

DataFed Data Repository

- DataFed Core Services:
  - Control servers
  - Web servers
  - Database
  - DataFed handles authorization, concurrency control, access control, etc.

- Example DataFed Data Repository:
  - Any file-system supported by Globus
  - Raw binary data stays in repository
  - Metadata only in DataFed Database

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

# DataFed Interface - Web Portal and CLI



Modern web application



Command-line & Python package

Interactive & non-interactive scripting
(e.g. - HPC environments)

# General metadata



- Owner
- Repository
- Size
- Date / time
- Keywords
- Tags
- Title
- Description
- more...
- Searchable

# Domain-specific Metadata



- Arbitrary tree
- String, numeric, array values
- Searchable
- Community defined schemas upcoming
- No need to embed metadata in file names

# Search

# Provenance Capture

- Currently supports:
  - "Derived from"
  - "Component of"
  - "New version of"
- More coming soon
- User-defined relationships
  - "Calibration associated with"
  - ...

# DataFed Applied to Simulations / Modelling

```
1 datafed data get input_parms_record ./parameters.txt
2 simulation run --input ./parameters.txt
3 context.json = parse(parameters)  AND/OR
3 context.json = extract_metadata(results_files)
4 datafed data create \
            --metadata context.json \
            --raw-data-file results_files.tar \
            "New_sim_results"
```

- Line 1 - Unambiguous identification of input files
- Line 4 - One line to create DataFed record, provide metadata, upload data
- Works for parameter sweeps - multiple similar simulation runs
- Towards reproducible simulation workflow
- Ease collection of training data for machine learning / deep learning

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

# DataFed Applied to Artificial Intelligence

```
1 datafed data get training_coll ./train
2 datafed data get input_parms_record ./parameters.txt
3 python script.py --train ./train --parms ./parms.json
4 context.json = parse(output.log)
5 datafed data create \
            --metadata context.json \
            --raw-data-file results_files.tar \
            "New_sim_results"
```

- Line 1 - Collaborative collection of datasets
- Line 1 - Easily staging data located at multiple repositories at file-system
- Line 5 - One line to create DataFed record, provide metadata, upload data
- Assemble training datasets using tags and keywords of data records

OAK RIDGE | LEADERSHIP
National Laboratory | COMPUTING
FACILITY

# Data Infrastructure Applied to Experimental Sciences

- Data Gateway servers facilitate (at instrument):
  - Data Ingest
    - Consistent capture of metadata / context wrt experiments
    - Drag-and-drop upload to data repository
    - Automated upload of data for long-running experiments
  - Data Management web portal piped in
  - Data Analytics web portal (JupyterHub) piped in
  - Automated data processing jobs on HPC
  - Instrument remote control
- DataFed - multi-user, -modal, -instrument data correlations

OAK RIDGE
National Laboratory | LEADERSHIP COMPUTING FACILITY

# Summary



Exp/Obs Facility - Off Network

Data Gateway — Isolated Instrument

Workstation

DataFed Raw Data (GridFTP)

Networked Instrument

Local DataFed Repository

Exp/Obs Facility - DAQ Integration

Facility DataFed Repository — HPC / Cloud — Jup. Hub — DataFed Core — Virtual Data Gateway

Compute / Analytics Facility

DataFed:
- Data backplane
- Federation of repositories
- Rich metadata, provenance
- Data sharing, search, movement, staging, etc.

Data Gateway:
- Data upload
- Metadata capture
- Portal to data services

OAK RIDGE National Laboratory | LEADERSHIP COMPUTING FACILITY

# Acknowledgements

**OAK RIDGE** | LEADERSHIP
National Laboratory | COMPUTING
FACILITY

# Questions?

- Try out DataFed:
  - https://datafed.ornl.gov
  - Requires Globus account
  - Contact: stansberrydv@ornl.gov
- somnaths@ornl.gov