# iRODS

**Data Management Technology Driven and Sustained by the eResearch Community**

David Fellinger
Storage Scientist
Data Management Technologist
iRODS Consortium
11 February 2022

eResearch NZ 2022
9-11 FEBRUARY, CHRISTCHURCH & ONLINE

# What is iRODS?

- The iRODS (Integrated Rule-Oriented Data System) technology is an open source data management platform.

- The iRODS Consortium was formed in 2013 funded by a group of technology companies and universities based on government funded work dating back to 1995.

- The iRODS community is comprised of over 30 members spanning a variety of disciplines worldwide.

- The Integrated Rule-Oriented Data System (iRODS) has been designed by the iRODS Consortium with 4 key functionalities;

**DATA VIRTUALIZATION**  **DATA DISCOVERY**  **WORKFLOW AUTOMATION**  **SECURE COLLABORATION**

## iRODS is:

- Open Source
- Distributed
- Data Centric
- Metadata Driven

# Why is the Focus on Metadata?

- User defined metadata allows researchers to describe their work.
  - Metadata is a descriptive communication tool that can categorize a research work in a way that is relevant to other researchers in the same field.

- The cataloging of metadata as a file description allows;
  - Discovery
  - Data grouping based on content to enable analysis
  - Data movement to analytic platforms
  - An HSM (Hierarchical Storage Management) based on data content and other factors that go far beyond just file extension and date.

- iRODS can extract metadata from the file either in flight or in a static storage location.
  - Metadata extraction can be taken from a file header or, it can be extracted from the file contents
  - Rules can be established based upon site policies to catalog this metadata and control all aspects of collection management.
  - Metadata is dynamic and can change based on citations or access patterns.
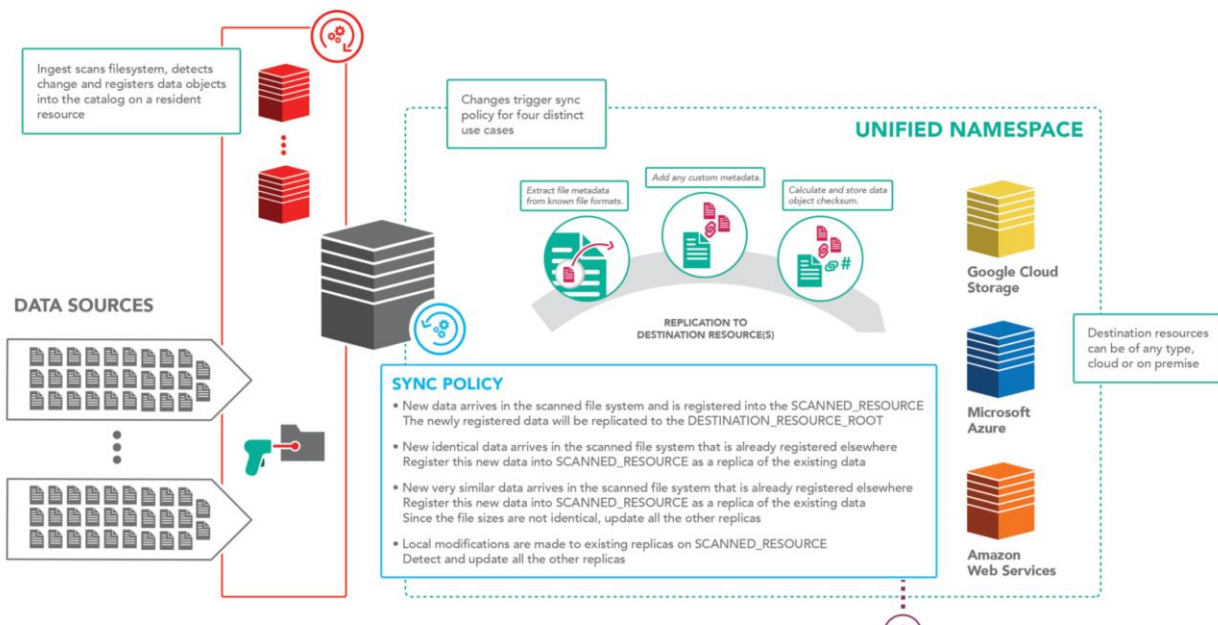
# What is New in iRODS Technology?

**iRODS**

- Logical locking has been implemented to assure file consistency.
  - While file systems implement locks, iRODS virtualizes multiple file systems to enable geographically diverse storage
  - Files can be replicated over diverse file systems so locks and checksums assure accurate synchronization
- Pass-through streaming has been added to the S3 resource plug-in.
  - iRODS can move data to the "cloud" faster than the AWS CLI
  - This is tested against AWS, GCS, Ceph, MinIO, and Fujifilm
- Glacier capability has been added to the S3 resource plug-in.
  - This has been tested against AWS and Fujifilm
- A partnership has been announced with Globus.
  - A connector has been designed that enables iRODS access at Globus endpoints
- GUIs (Graphical User Interfaces) have been improved with additional functionality responding to community requirements.
  - Metalnx has been updated and a gallery image viewing capability has been added to make discovery intuitive
  - A management GUI is being evolved to enable web-enabled configuration

iRODS
CONSORTIUM

# iRODS is a Product of the Community

- Monthly Planning Committee meetings allow each member to report on progress, goals, and problems.
    - Future development work is discussed and priorities are set by consensus
- Working groups are established so that the end product reflects "real world" solutions to a dynamic environment. These groups currently focus on:
    - Metadata templates
    - Authentication and authorization
    - S3 interfaces and protocols
    - Imaging
- An active Google Group is utilized to answer questions and concerns in real time
- An annual User Group Meeting brings the community together and many members and users describe their particular deployments.
    - The 2021 UGM featured 41 presentations and many lightning talks
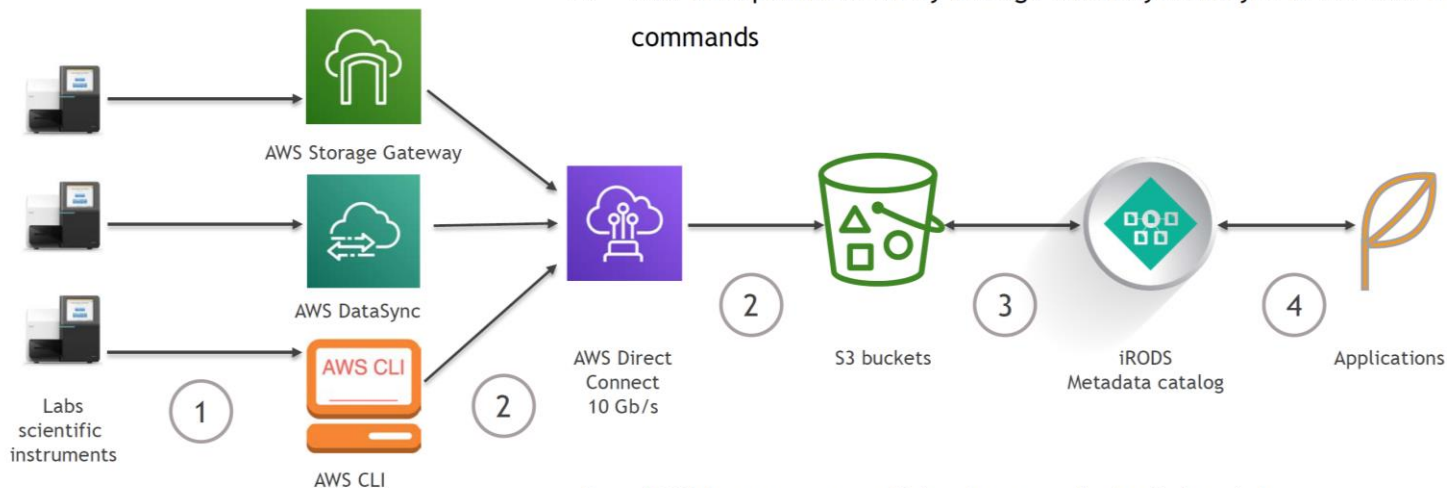    - All of these presentations can be found at: https://irods.org/ugm2021

# Worldwide Data Synchronization Fully Enabled

Ingest scans filesystem, detects change and registers data objects into the catalog on a resident resource

Changes trigger sync policy for four distinct use cases

UNIFIED NAMESPACE

Extract file metadata from known file formats.

Add any custom metadata.

Calculate and store data object checksum.

REPLICATION TO DESTINATION RESOURCE(S)

DATA SOURCES

Google Cloud Storage

Microsoft Azure

Amazon Web Services

Destination resources can be of any type, cloud or on premise

**SYNC POLICY**

• New data arrives in the scanned file system and is registered into the SCANNED_RESOURCE
The newly registered data will be replicated to the DESTINATION_RESOURCE_ROOT

• New identical data arrives in the scanned file system that is already registered elsewhere
Register this new data into SCANNED_RESOURCE as a replica of the existing data

• New very similar data arrives in the scanned file system that is already registered elsewhere
Register this new data into SCANNED_RESOURCE as a replica of the existing data
Since the file sizes are not identical, update all the other replicas

• Local modifications are made to existing replicas on SCANNED_RESOURCE
Detect and update all the other replicas

Public cloud and private storage facilities can be combined to enable collaboration

# Typical data flow diagram

1. Instruments writes raw data into local scratch space
2. Raw data pushed to S3 by Storage Gateway/DataSync or via AWS CLI S3 commands

AWS Storage Gateway

AWS DataSync

AWS CLI

Labs scientific instruments

① 

AWS Direct Connect 10 Gb/s

② S3 buckets

③ iRODS Metadata catalog

④ Applications

3. iRODS system scans S3 buckets regularly via Lambda
4. Applications request data via iRODS metadata catalog

Presentation available from: https://irods.org/uploads/2021/Konduri-Khavich-BMS-Leveraging_iRODS_for_Scientific_Applications_in_AWS_Cloud-slides.pdf accessed 13 September 2021

**iRODS** CONSORTIUM

Business Process & Solution: Lab Data Hub Architecture

# Deployment: KU Leuven

## System architecture (pilot environment)

**Clients**
- YODA
- metalnx
- DAVRods / Cyberduck
- iCommands
- Python Client

**Go front-end**

**iRODS**
- HA Proxy
- iRODS servers (VM)
- iCAT
- HOST
- MySQL Database Cluster

1 zone
All containerized (docker/nomad):
- 6 iRODS server containers
- Metalnx containers
- 2 WebDav containers
- ...

**Storage servers**

DC Heverlee
iRODS storage servers
Posix Resc1
iRODS primary
1A 1B 3A 3B
Ceph Resc3 Testing

3km

DC Leuven
2A 2B 4A 4B
iRODS storage servers
Posix Resc2
iRDS Replica Resc1

VLAAMS SUPERCOMPUTER CENTRUM

Presentation available from: https://irods.org/uploads/2021/Barcena-KULeuven-A_Year_of_iRODS_Lessons_Learned-slides.pdf accessed 13 September 2021

## TAPE LIBRARY AND ARCHIVING DATA FLOW

- The SURF Data Archive service offers storage space on a Tape Library
- The Service is accessible directly via command line clients, but sometimes it is configured as iRODS resource and accessed only through iRODS.

Presentation available from: https://irods.org/uploads/2021/Cacciari-SURF-Archiving_Off-line_and_Beyond_Using_the_BagIt_Format_and_the_bdbag_Library-slides.pdf accessed 13 September 2021
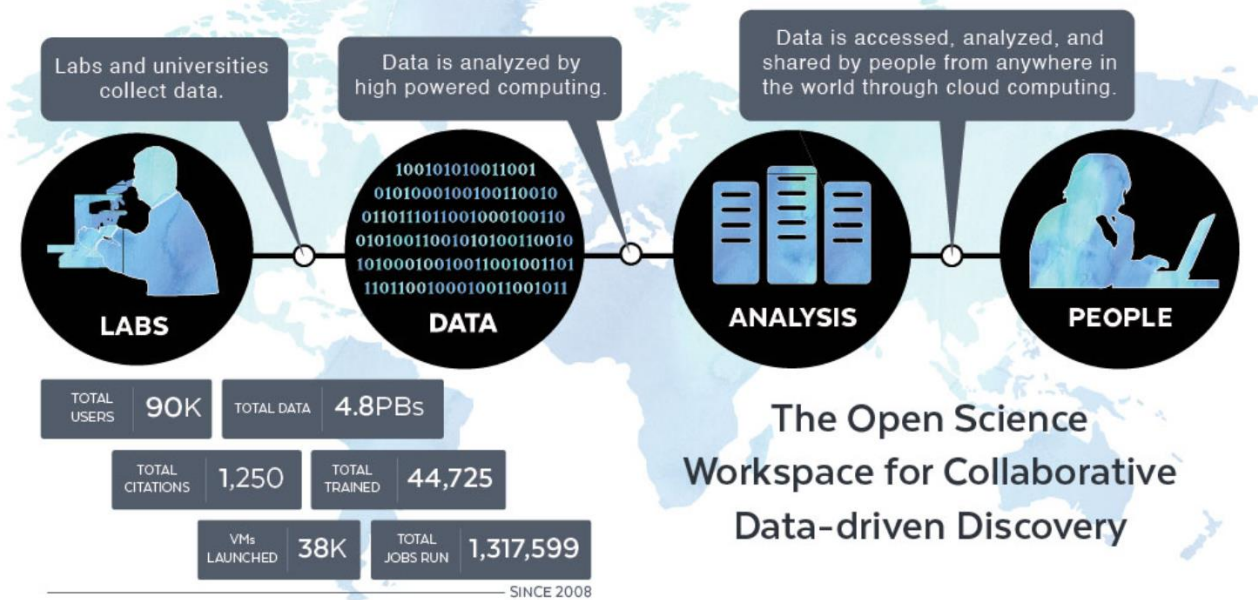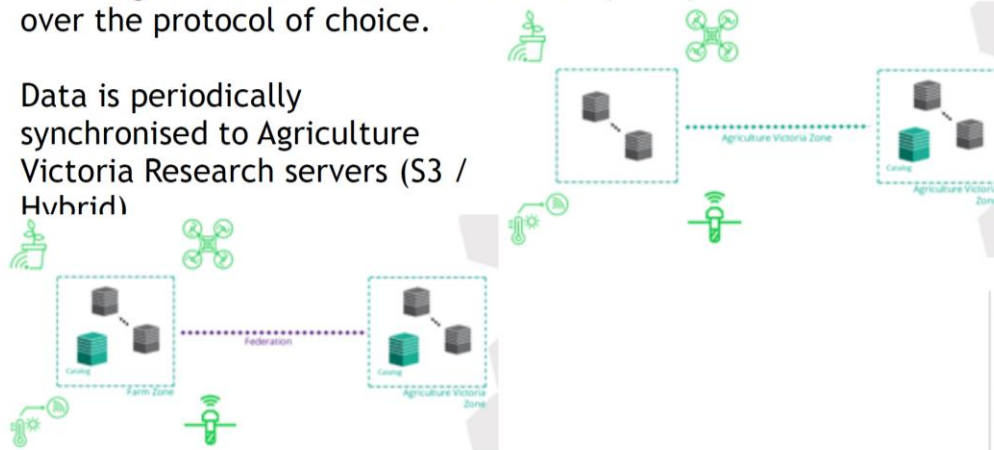
Diagram available from: https://cyverse.org/about accessed 13 September 2021

**iRODS**

## Emerging SmartFarm Data Infrastructure

Each SmartFarm may host their own application (iRODS) to manage metadata description and catalogue for each UAV trial.

Data is gathered from the UAV over the protocol of choice.

Data is periodically synchronised to Agriculture Victoria Research servers (S3 / Hybrid)

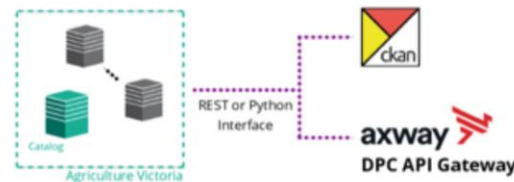SmartFarm hosts Agriculture Victoria Research servers (S3 / Hybrid)

Data is periodically **replicated** to Agriculture Victoria Research Servers (BASC)

Once data is at rest in the Agriculture Victoria Research namespace i.e. Horsham_UAV_AVR_Plot1

Data may be replicated to HPC storage for analytics.

Data may be published to CKAN or made accessible via the API gateway

Data may be shared over an IRODS interface : WebDAV, MetaInx, NFS, Command Line.

# Partner: Globus



iRODS + Globus deployment

Presentation available from: https://irods.org/uploads/2021/Vasiliadis-Globus-Automating_Data_Management_Flows_with_iRODS_and_Globus-slides.pdf accessed 13 September 2021

**iRODS**

- The iRODS Consortium serves the user community by building a product that consistently satisfies evolving data management requirements.

- Use cases span research sites and disciplines worldwide.

- Users can build iRODS rules that enable the auditable requirements of FAIR principles and site policy adherence.

- iRODS can enable complete workflow control, data lifecycle management, and present discoverable data sets with assured traceability and reproducibility.

# The iRODS Consortium (iRODS.org)

**iRODS**

The iRODS Consortium

- Leads software development and support of iRODS
- Hosts iRODS Events
- Tiered membership model

# Questions?

Additional use cases can be found in the
proceedings of the 2021 iRODS User Group Meeting;
https://irods.org/ugm2021

Thank you!
David Fellinger

davef@renci.org