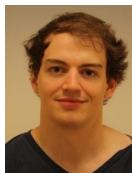# FAIR Sequencing Data Repository based on iRODS



Felipe O. Gutierrez

AMC - Academic Medical Center - Amsterdam, Netherlands

A.C.Camargo Cancer Center - São Paulo, Brazil

F. Oliveira Gutierrez

P.F.G. De Geest

Aldo Jongejan

Sjoerd Repping
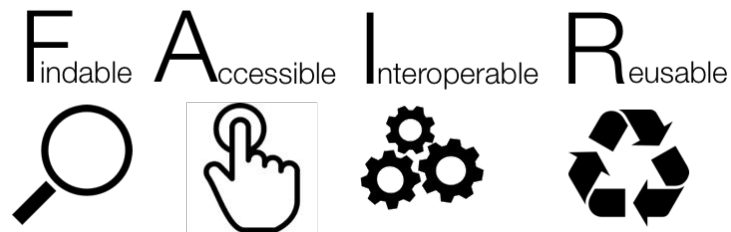
Diogo Ferreira Patrão

A.H.C. van Kampen

J.T. van den Berg

Silvia D. Olabarriaga

# Problem

- Inadequate RDM (Research Data Management) solution for NGS data (Next Generation Sequencing):
    - Individual storage and backup
    - Dispersed datasets
    - Disconnected from metadata
    - Not FAIR

F<sub>indable</sub> A<sub>ccessible</sub> I<sub>nteroperable</sub> R<sub>eusable</sub>
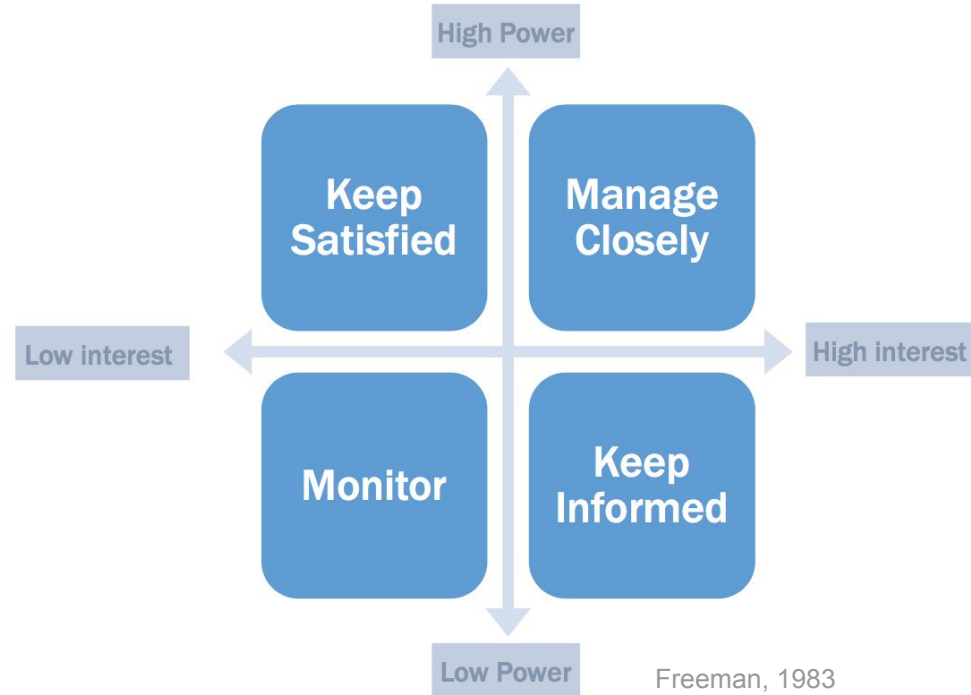
# Considerations

Fit within organization

- ICT culture
- Research culture
- Sustainability vision

Adhere to international community best practices

Reuse and extend existing solutions



High Power

Keep Satisfied

Manage Closely

Low interest

High interest
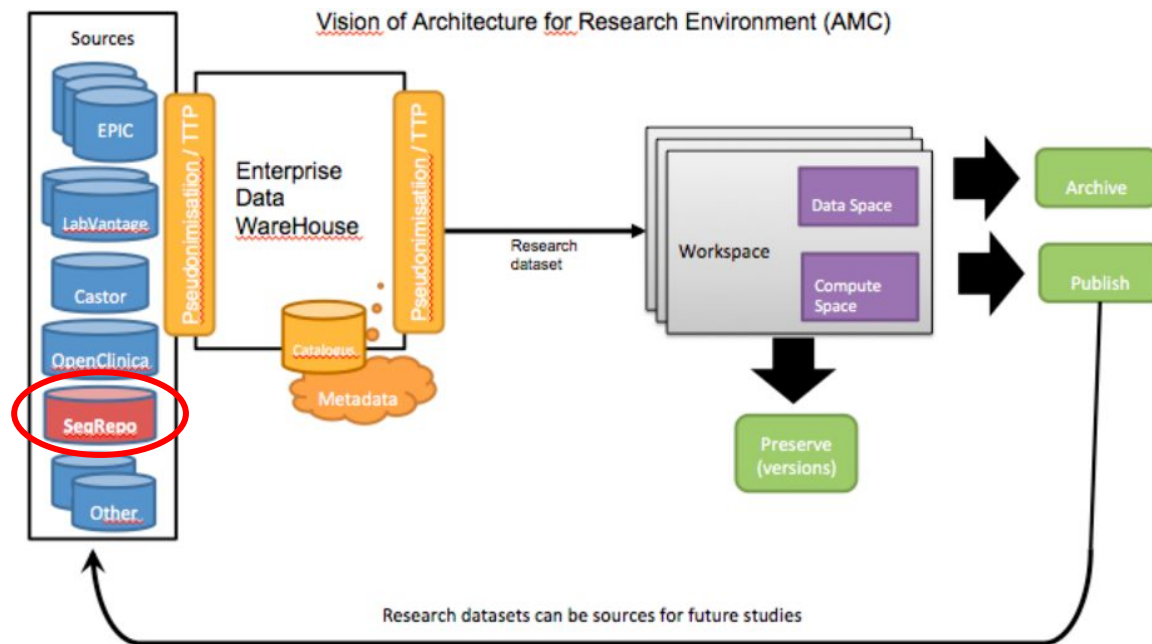
Monitor

Keep Informed

Low Power

Freeman, 1983

# Fit into AMC Vision for RDM

Based on NFU Data4Lifesciences WP2

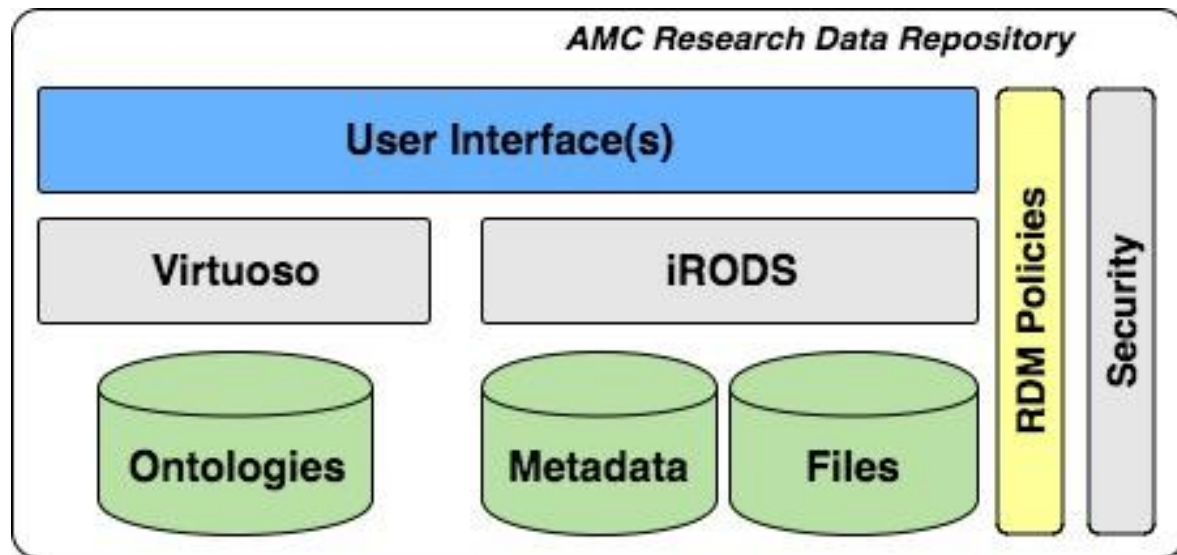An NGS repository that is:

- Part of an ecosystem
- Controlled by AMC
- Distributed
- Scalable
- FAIR compliant
- Easy to use



Vision of Architecture for Research Environment (AMC)
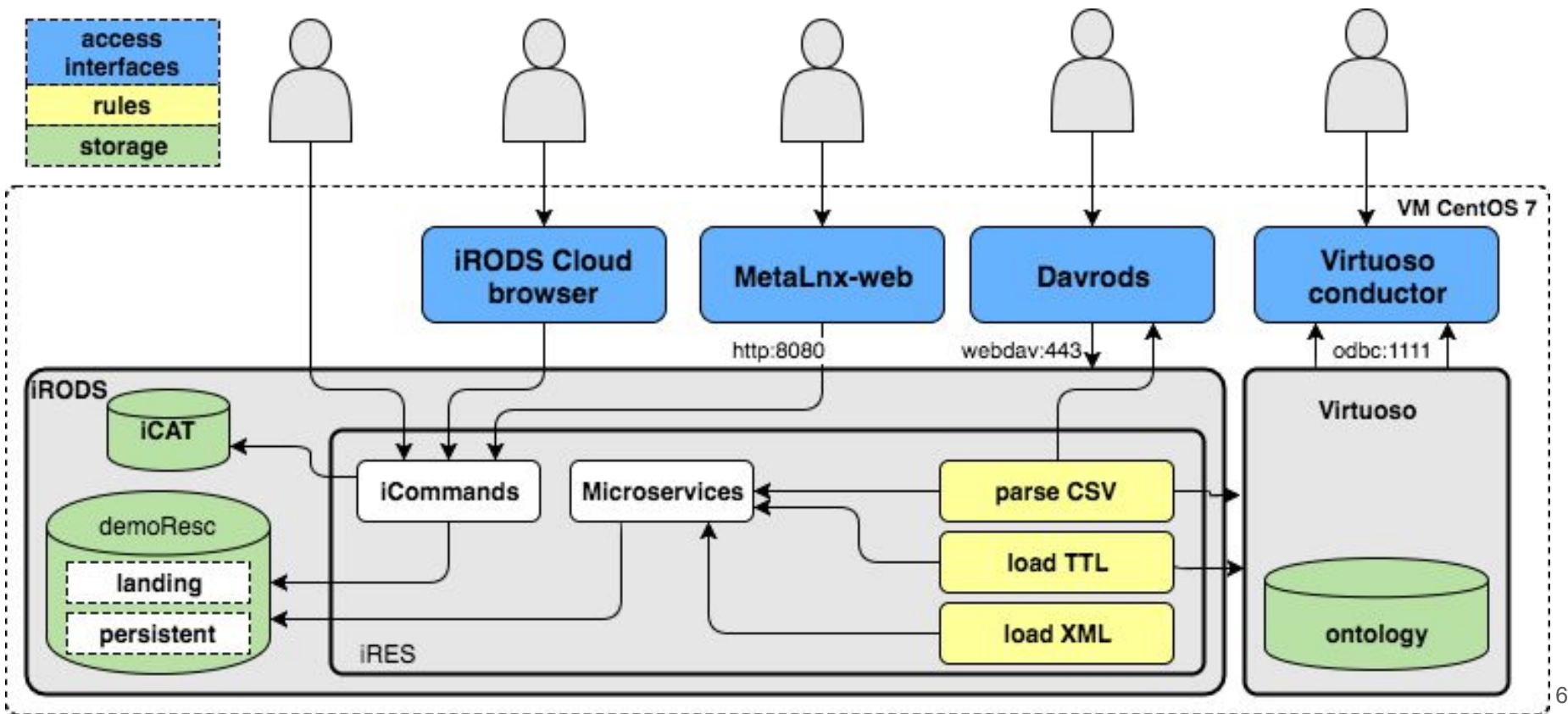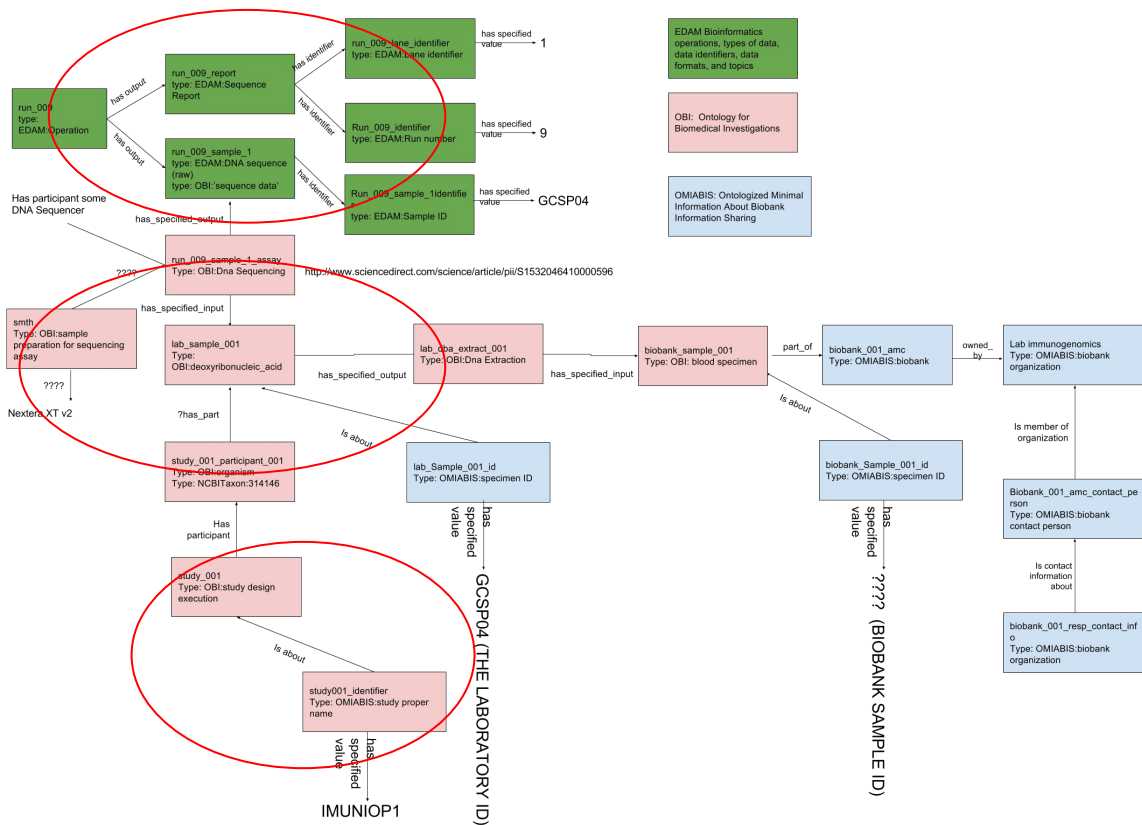
# System Design

- iRODS 4.1.10
  - Middleware
  - Data virtualization
- Virtuoso 7.2
  - Triplestore
  - Supports ontologies
- User interfaces:
  - Metalnx web
  - Davrods 4.1
  - iCommands

**AMC Research Data Repository**

User Interface(s)

Virtuoso | iRODS

Ontologies | Metadata | Files
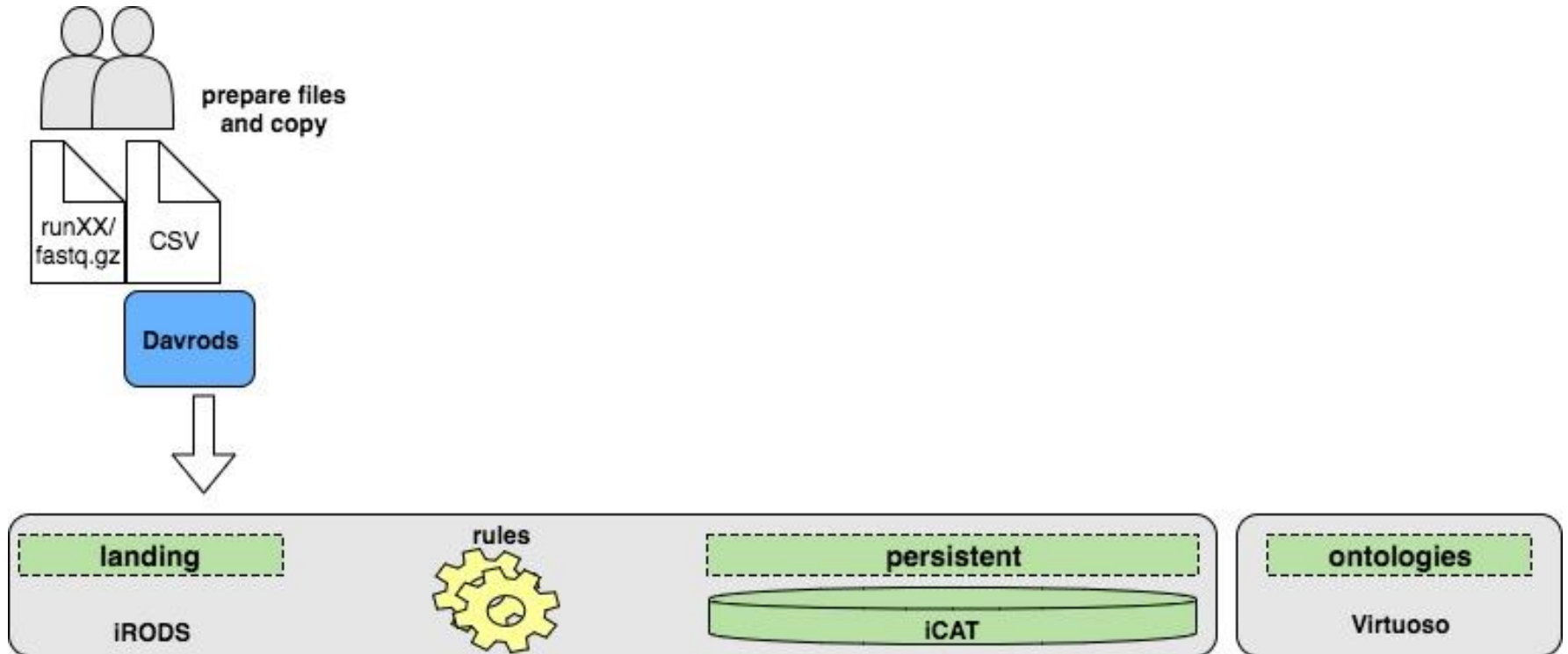
RDM Policies

Security

# System Architecture
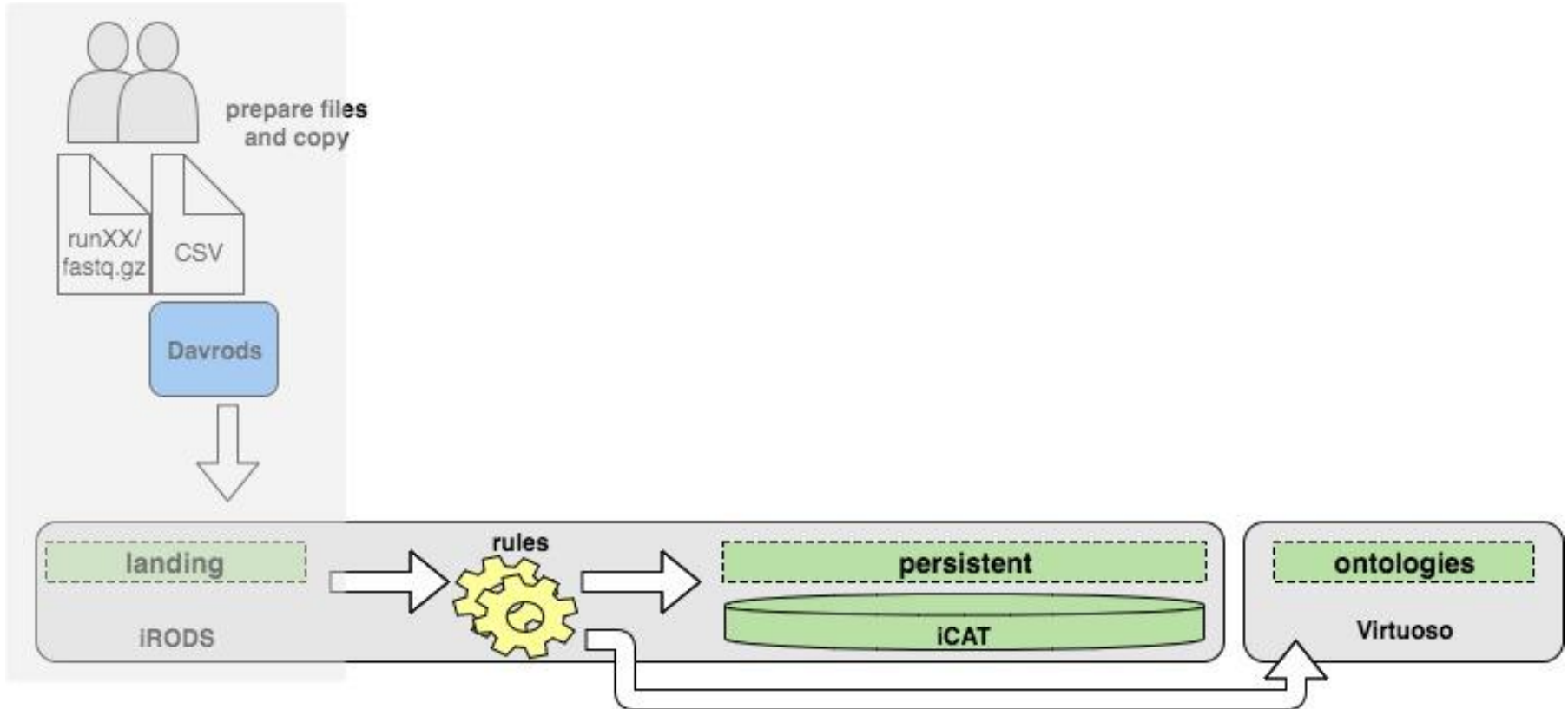
# Stewardship: Ontologies



- EDAM
  Ontology for bioinformatics operations, types of data, data identifiers, data formats, and topics
- OMIABIS
  Ontologized Minimum Information About Biobank data Sharing (MIABIS)
- OBI
  Ontology for Biomedical Investigations
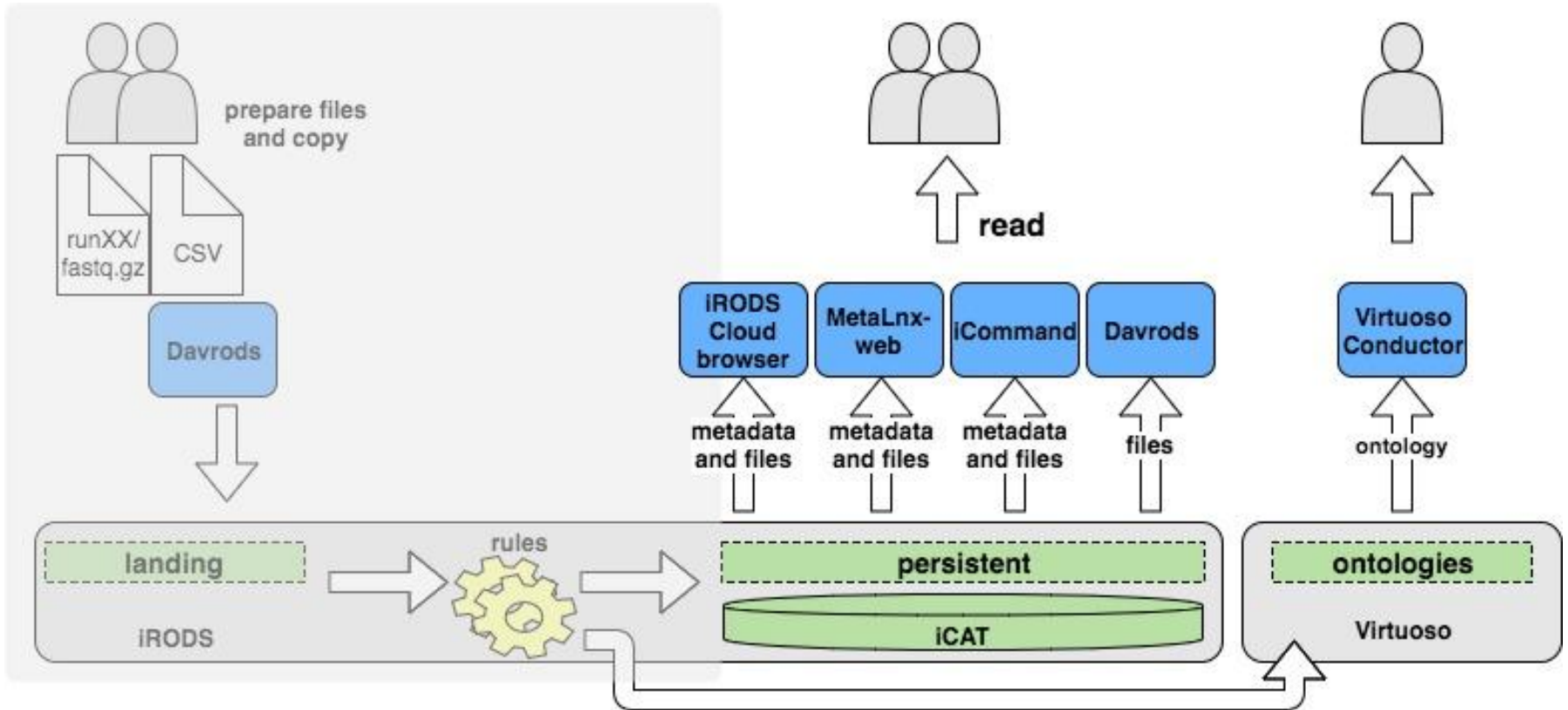- EFO
  Experimental Factor Ontology

# Workflow: Data Ingestion

# Workflow: (meta)data Registration

# Workflow: (meta)data Retrieval

# Access and Security

| functions | iRODS roles | | | | |
|---|---|---|---|---|---|
| | Davrods | iRODS Cloud Browser | Metalnx web | iCommand | Virtuoso |
| **Data management** | user | user | user | user | user |
| **Metadata management** | user * | user | user | user | user |
| **User/Group management** | | | PI | PI | |
| **Access control** | | | PI | PI | |
| **(meta)data curation** | | Data steward | Data steward | Data steward | |
| **Policies and rules** | | | Admin | Admin | |
| **Security** | | | Admin | Admin | |

# Report file



**Report file of the FAIR-RDM workflow**

| Sample sheet | Final status | Date |
|---|---|---|
| samplesheet_2016_12_20_Joana_s1.csv | Workflow process completed successfully | 2017-04-11 |

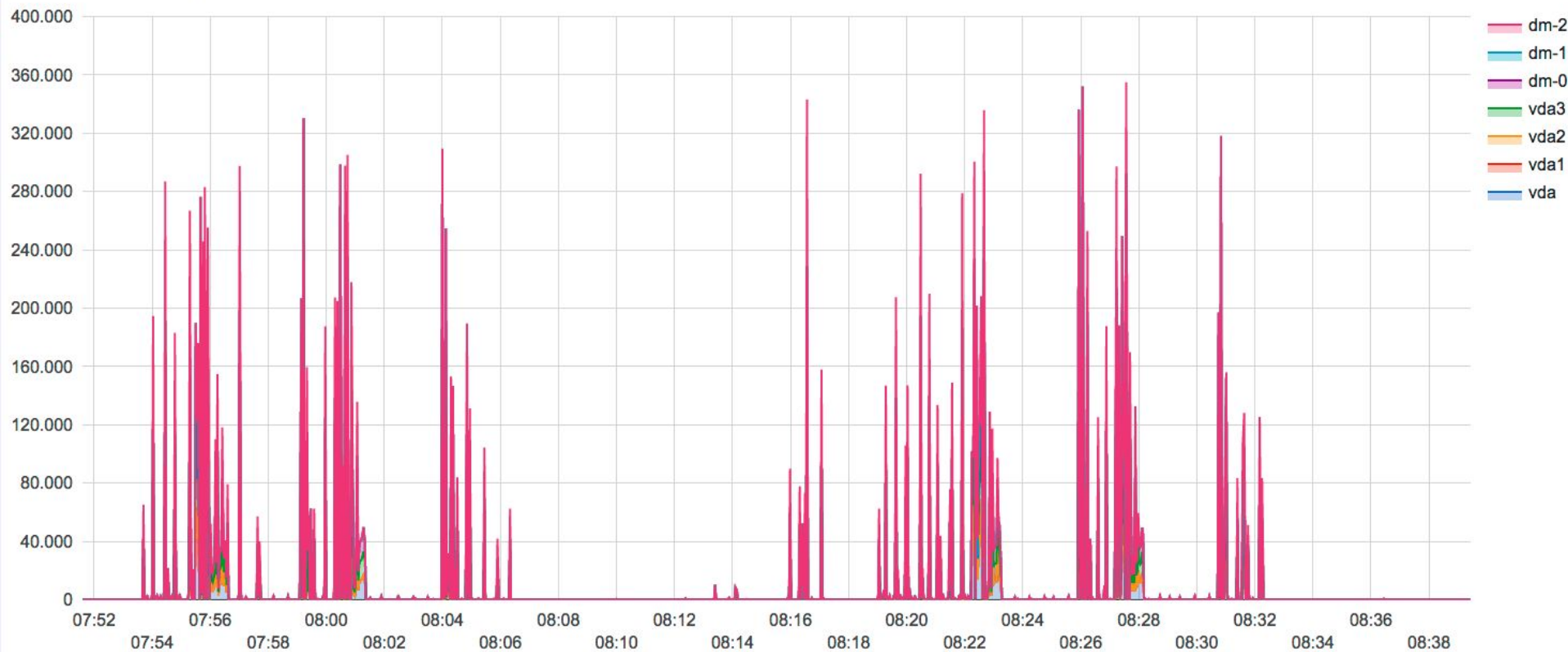| Step | Status | Code | Timestamp | Description |
|---|---|---|---|---|
| 010 | Process on going | 010 | 2017-04-11:11:02:44 | Workflow started successfuly. |
| 010 | OK | 011 | 2017-04-11:11:02:45 | Changed the permission of the landing directory to read-write successfuly. |
| 012 | OK | 014 | 2017-04-11:11:03:05 | The number of fastq files matches with the [samplesheet_2016_12_20_Joana_s1.csv] file. |
| 011 | OK | 012 | 2017-04-11:11:03:05 | The [samplesheet_2016_12_20_Joana_s1.csv] file name matches with regex expression (sa |
| 013 | OK | 016 | 2017-04-11:11:03:05 | The experiment from the [samplesheet_2016_12_20_Joana_s1.csv] file is new. |
| 020 | OK | 020 | 2017-04-11:11:03:05 | Connected to Virtuoso and get unique ID. |
| 021 | OK | 022 | 2017-04-11:11:03:05 | Created new experiment directory. |
| 021 | OK | 027 | 2017-04-11:11:03:20 | Fastq files restructured inside the experiment. |
| 021 | OK | 024 | 2017-04-11:11:04:04 | Created XML file successfuly. |
| 021 | OK | 023 | 2017-04-11:11:06:39 | Created TTL file successfuly. |
| 021 | OK | 028 | 2017-04-11:11:06:39 | Experiment [experiment_51] copied to the persistent directory. |
| 040 | OK | 040 | 2017-04-11:11:07:30 | A XML [samplesheet_2016_12_20_Joana_s1.xml] was upload. |
| 041 | OK | 041 | 2017-04-11:11:07:32 | XSD for XML [samplesheet_2016_12_20_Joana_s1.xml] exists. |
| 042 | OK | 043 | 2017-04-11:11:07:33 | XML [samplesheet_2016_12_20_Joana_s1.xml] validate by XSD schema. |
| 043 | OK | 045 | 2017-04-11:11:08:12 | XML [samplesheet_2016_12_20_Joana_s1.xml] load proccess OK. |
| 046 | OK | 047 | 2017-04-11:11:08:14 | Changing the permission of the persistent directory to read-only. This is the last operation o |

13

# nmon read KB/s


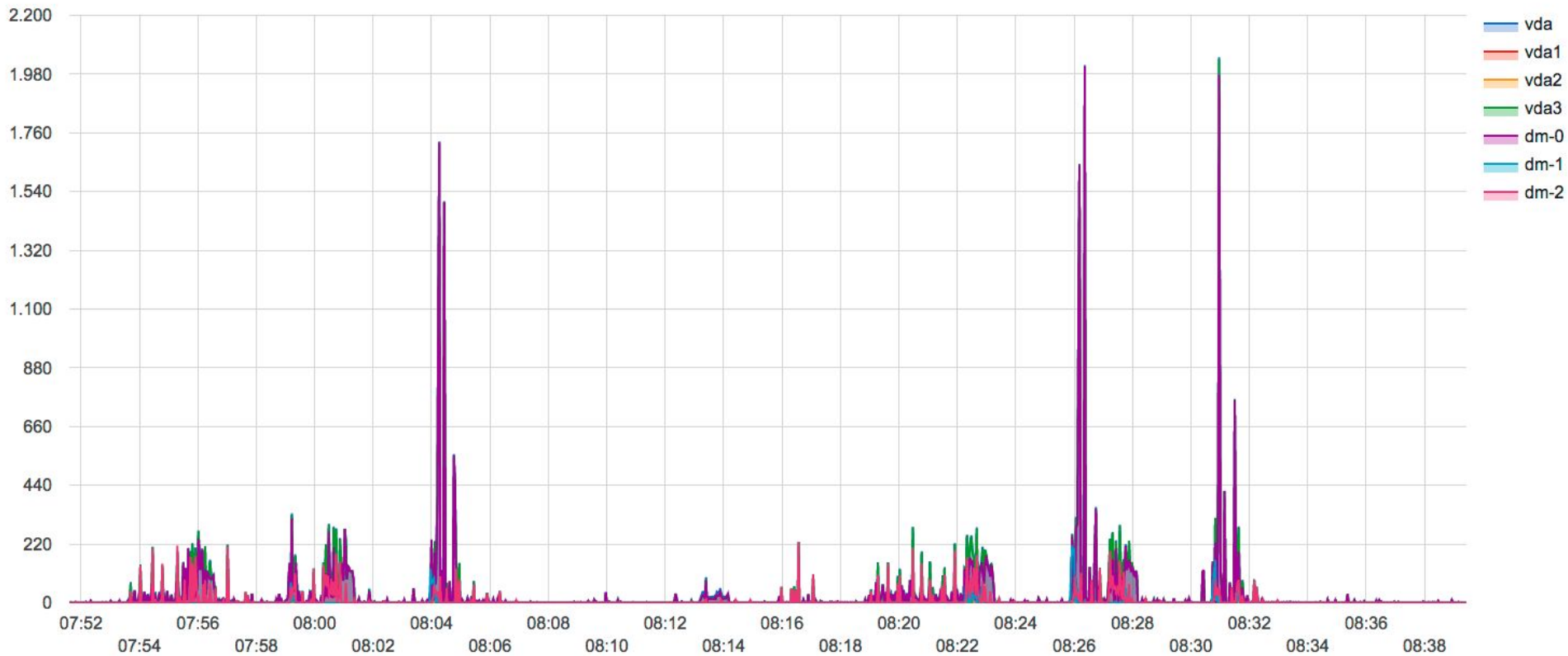
Disk Read KB per second (Stacked)

# nmon write KB/s



Disk Write KB per second (Stacked)

# nmon IOPs



Disk Transfers per second

# Qualitative & Quantitative questions

- (meta)data preparation? Clear, doable, easy, ...
- (meta)data upload? Type, size, quantity, integrity, ...
- Rule processing? Report file clear and easy, system delay feedback, ...
- (meta)data retrieval? Findable, Accessible, Organized, Interoperable, Reusable, ..
- Concurrent users, variation on the number and size of files.

# Acknowledgements