# EIC Rucio Schema

**Anil Panta**
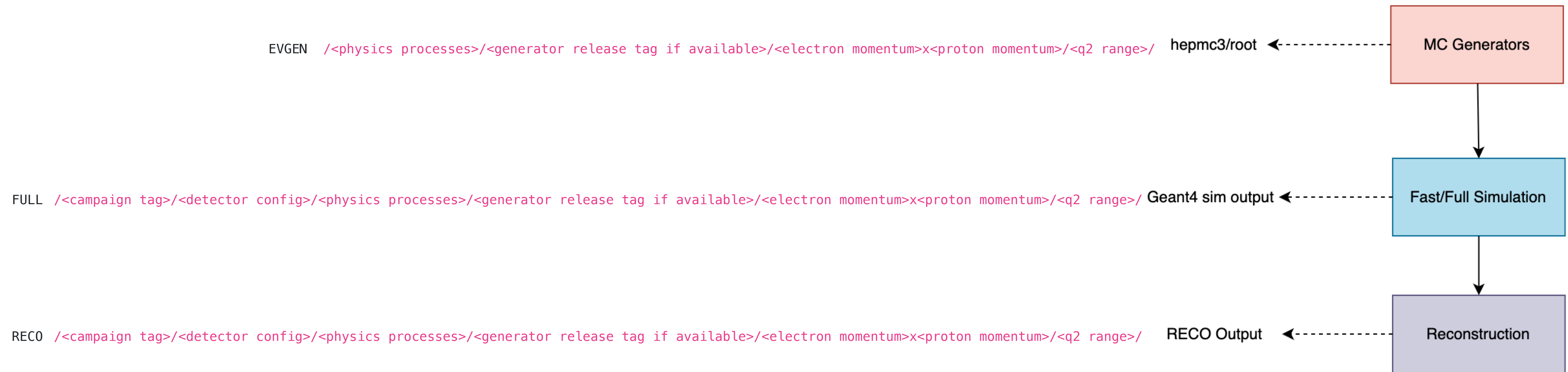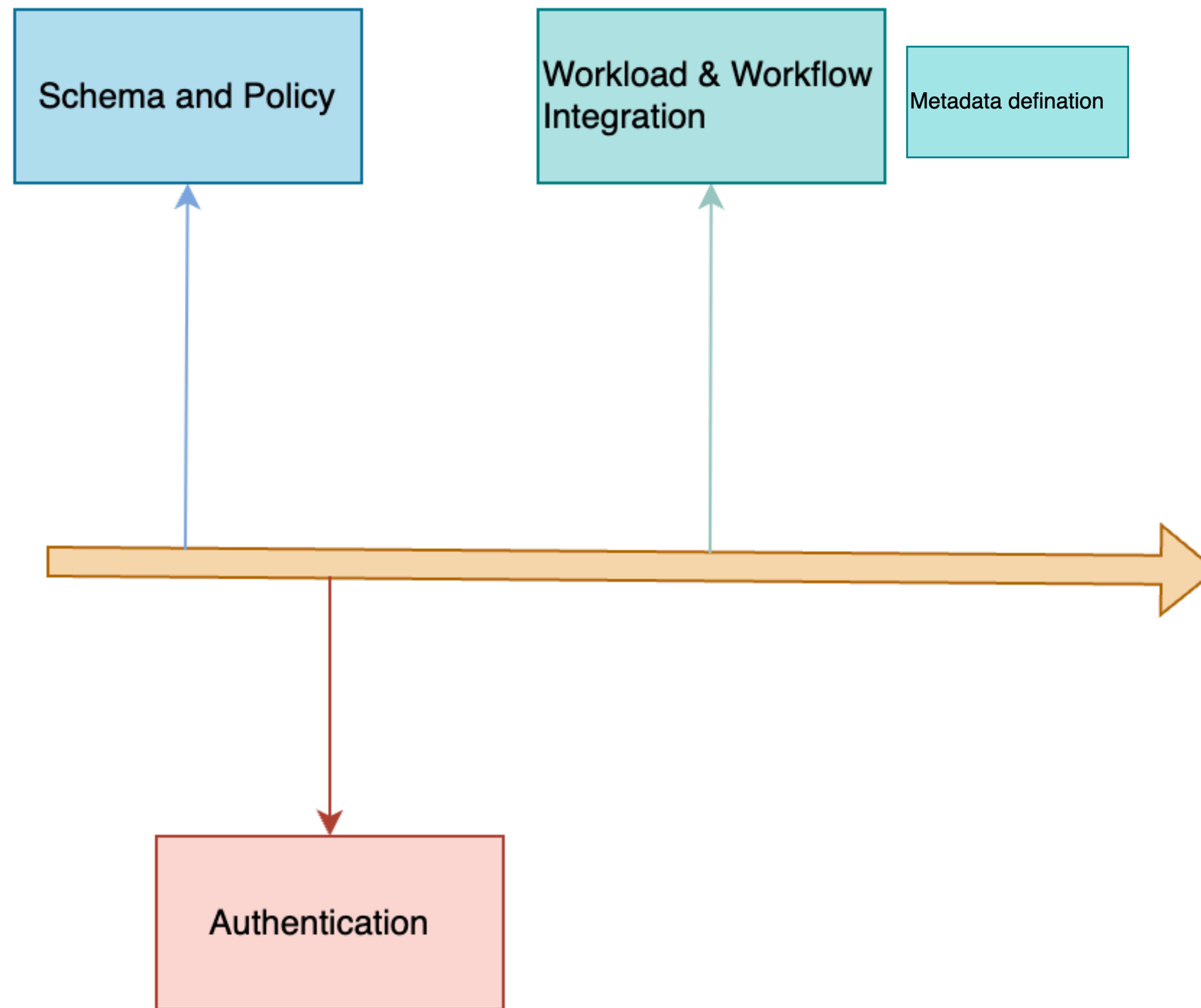
Feb 19, 2024

Jefferson Lab

# Current MC structure

- Three file type to register to storage:
  - Output of MC generator
  - Geant4 Simulation output.
  - Reconstruction output.
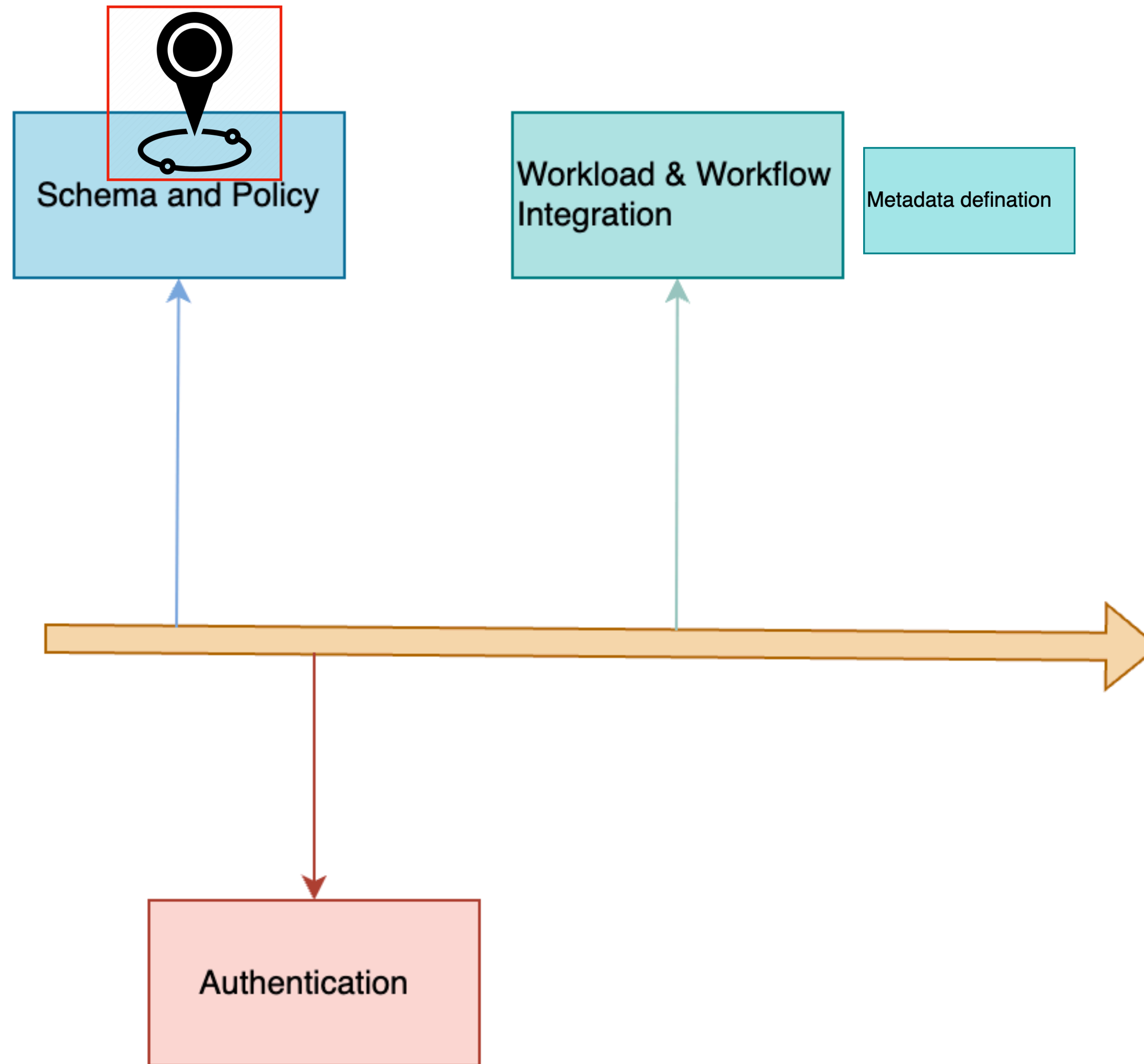
# Things to do:

# Things to do:
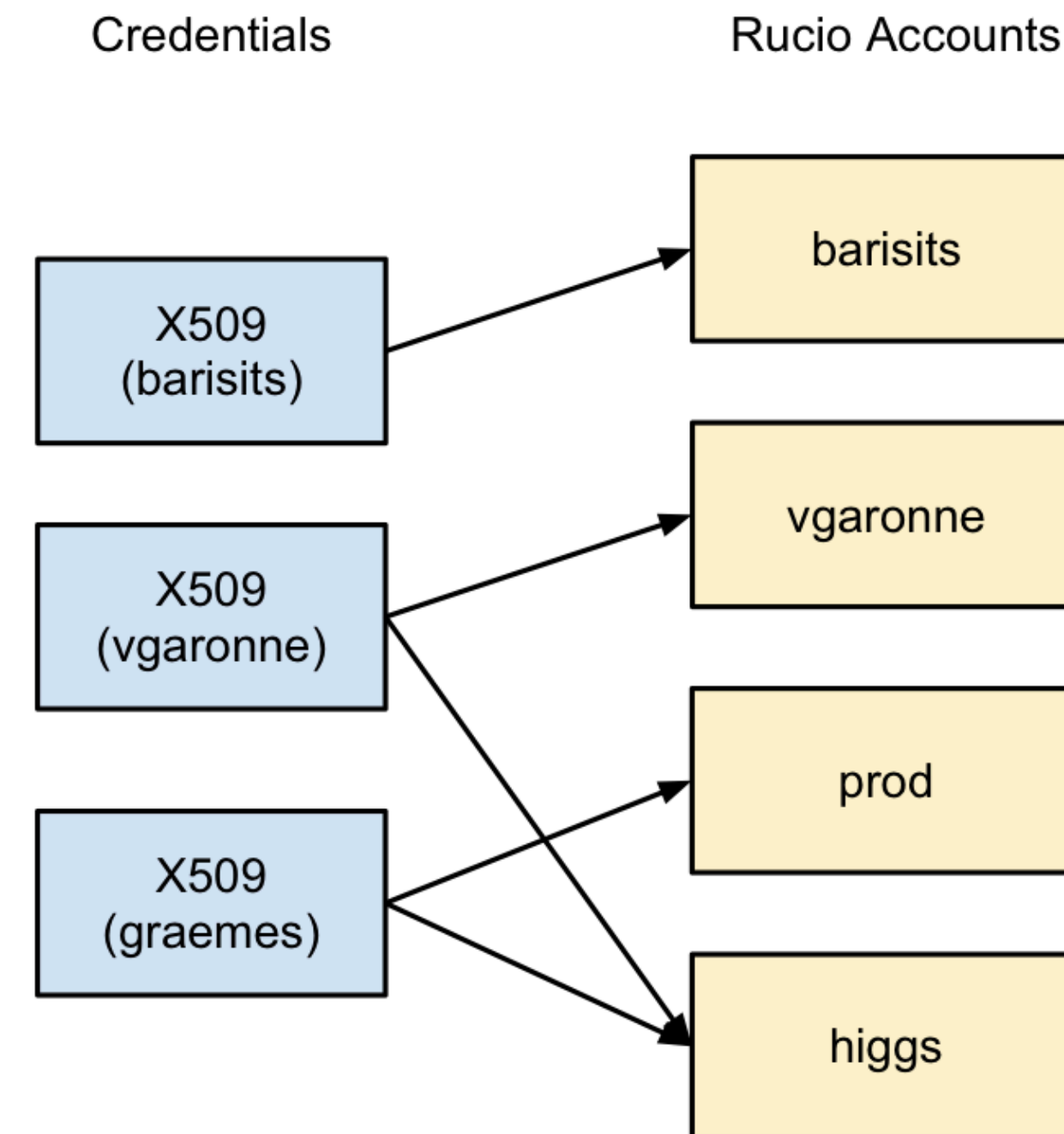


Schema and Policy

Workload & Workflow Integration

Metadata defination

Authentication

Jefferson Lab

# Account

- Two types:
  - User account : for each user
  - Group account : for central activity

- Rucio user = "account" + "identity" (N:M mapping)
- "account" = nickname
- "identity" = specification of authentication type + user identifier

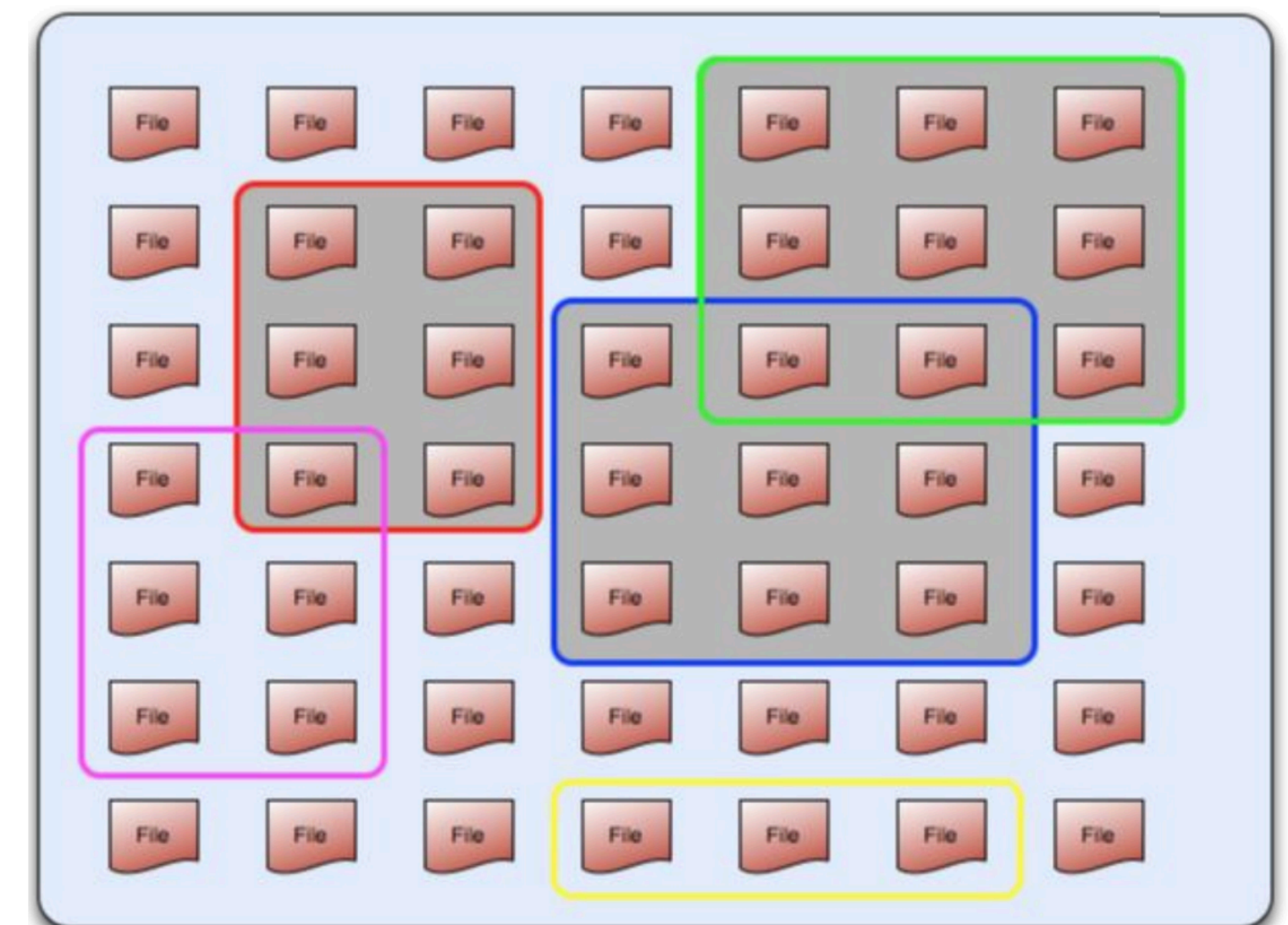More on Accounts during Authorization Discussion

Credentials          Rucio Accounts

X509
(barisits)            barisits

X509                  vgaronne
(vgaronne)

                      prod
X509
(graemes)
                      higgs

Jefferson Lab

# Rucio Naming Schema

- Any path is referred to as DID (Data Identifier)
            **scope:name**

- Scope is for Namespace separation.

- **DID can be flat** or hierarchal (partial or full)

- Types of did:

  - File : corresponding to actual file.

  - Dataset : Grouped set of files.

  - Container : Group set of datasets/container

# Atlas/rucio-default naming schema

- Flat naming.

- Scope is fine-grained.

**User data**

Scope: User name, e.g.,  user.jdoe:this.is.my.test.file001

user.jdoe:this.is.my.dataset1

**Group data**

Scope: Group name, e.g.,       group.thebest:this.is.our.test.file001

group.thebest:this.is.my.dataset1

**Detector data**

Scope: data + <year>, e.g.,     data18:16TeV.00199926.calibration.daq.RAW

**Reprocessed data**

Scope: Campaign +<year> + real data scope, e.g., Repro18Data18:00169783.AOD.r2059

dataset    `data23_13p6TeV:data23_13p6TeV.00450445.physics_Main.f1342_m2112`

files    `data23_13p6TeV:data23_13p6TeV.00450445.physics_Main.f1342_m2112._lb0110.01`    `data23_13p6TeV:data23_13p6TeV.00450445.physics_Main.f1342_m2112._lb0110.01`

Jefferson Lab

# Atlas/rucio-default naming schema

- Flat naming.

- Scope is fine-grained.

**User data**

Scope: User name, e.g., user.jdoe:this.is.my.test.file001

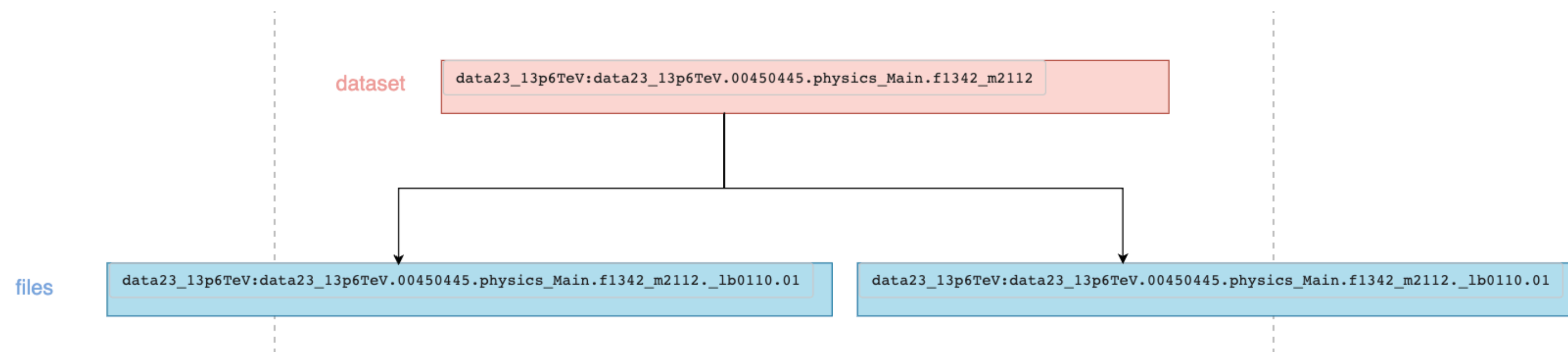user.jdoe:this.is.my.dataset1
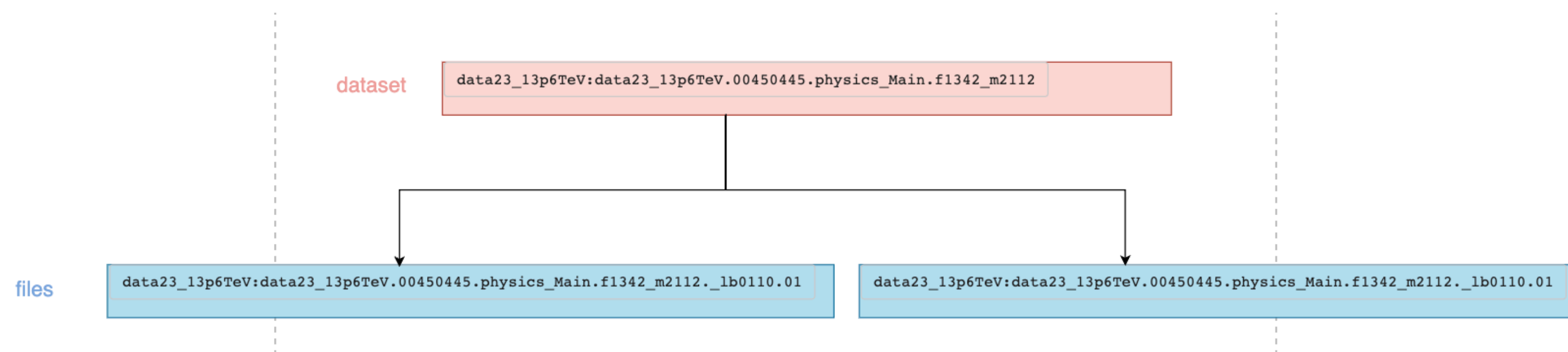
**Group data**

Scope: Group name, e.g., group.thebest:this.is.our.test.file001

group.thebest:this.is.my.dataset1

**Detector data**

Scope: data + <year>, e.g., data18:16TeV.00199926.calibration.daq.RAW

**Reprocessed data**

Scope: Campaign +<year> + real data scope, e.g., Repro18Data18:00169783.AOD.r2059
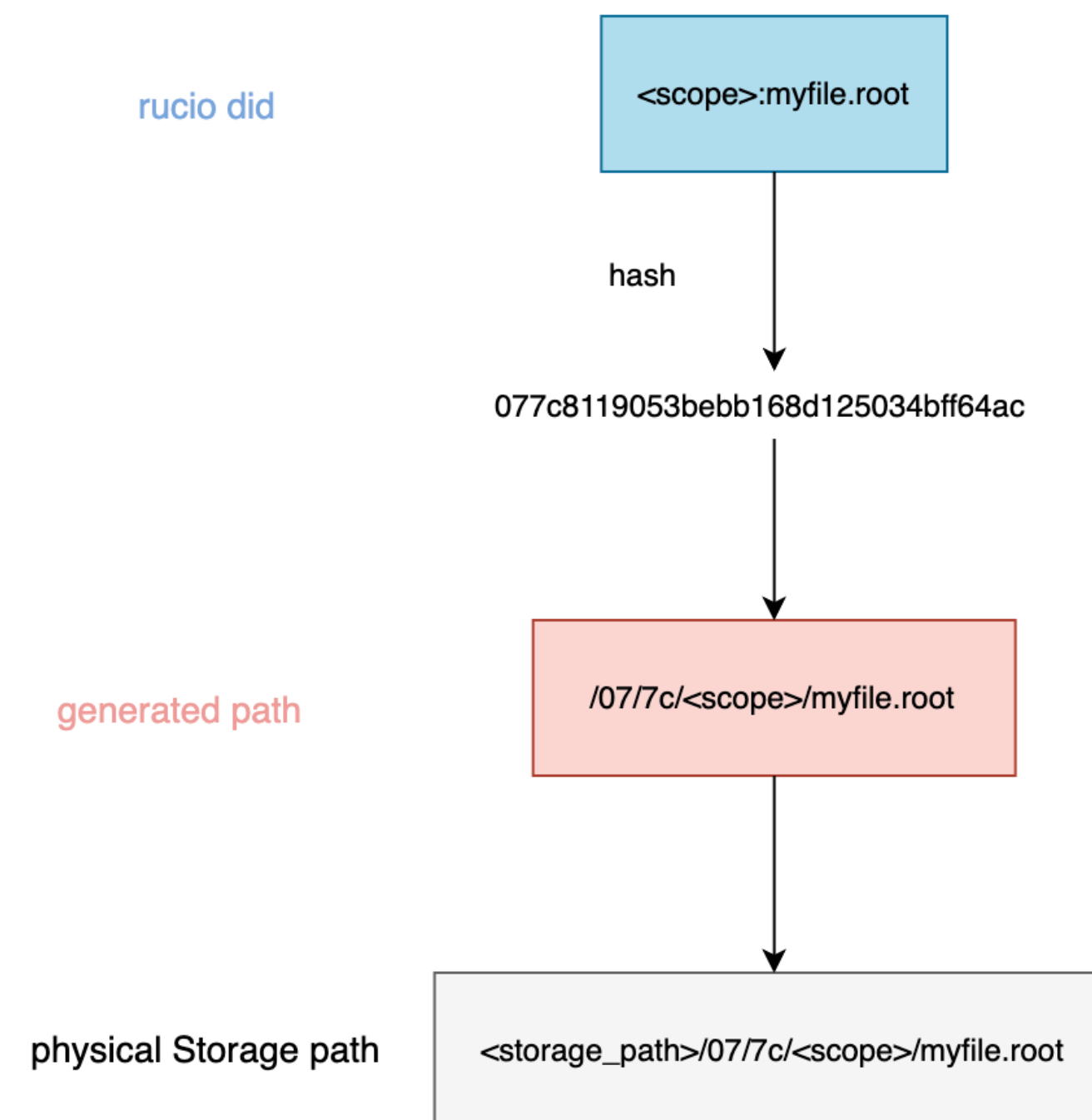
# Atlas Dataset naming schema

- Flat naming.

| Dataset type | Nomenclature | Example |
|---|---|---|
| Monte Carlo Datasets | mcNN_subProject.datasetNumber.physicsShort.prodStep.dataType.Version | mc08.105010.J1_pythia_jetjet.recon.ESD.e344_s456_r456 |
| Real Data (Primary) | DataNN_subProject.runNumber.streamName.prodStep.dataType.Version | data08_cos.00079123.physics_HLT_Cosmics_NIM4.daq.RAW<br>data08_cos.00079123.physics_HLT_Cosmics_NIM4.daq.AOD.f35 |
| Physics Container datasets | Project.runRange.StreamName.PhysCont.dataType.version | |
| Calibration dataset | dataNN_calib.xxxxxxxx.calibration_DetectorPart-meta-information-field.daq.RAW | data08_calib.00654321.calibration_LArElec-Pedestal-Medium-EM.daq.RAW |
| User dataset | user.userName.[otherFields] | |
| Group dataset | group.groupName.[otherFields] | |
| Conditions dataset | Project.internalCondNumber.shortComment.COND | |
| Database Release datasets | ddo.NNNNNN.[otherFields].vDBReleaseVersion | ddo.000001.Atlas.Ideal.DBRelease.v09010207 |
| SW Release datasets | sitNN.nnnnnn.AtlasSWRelease.PAC.vMMmmp[cc] | |

Jefferson Lab

# Atlas/rucio-default naming schema: Physical file

- Each file did is mapped to storage path.
  - Based on some policy.
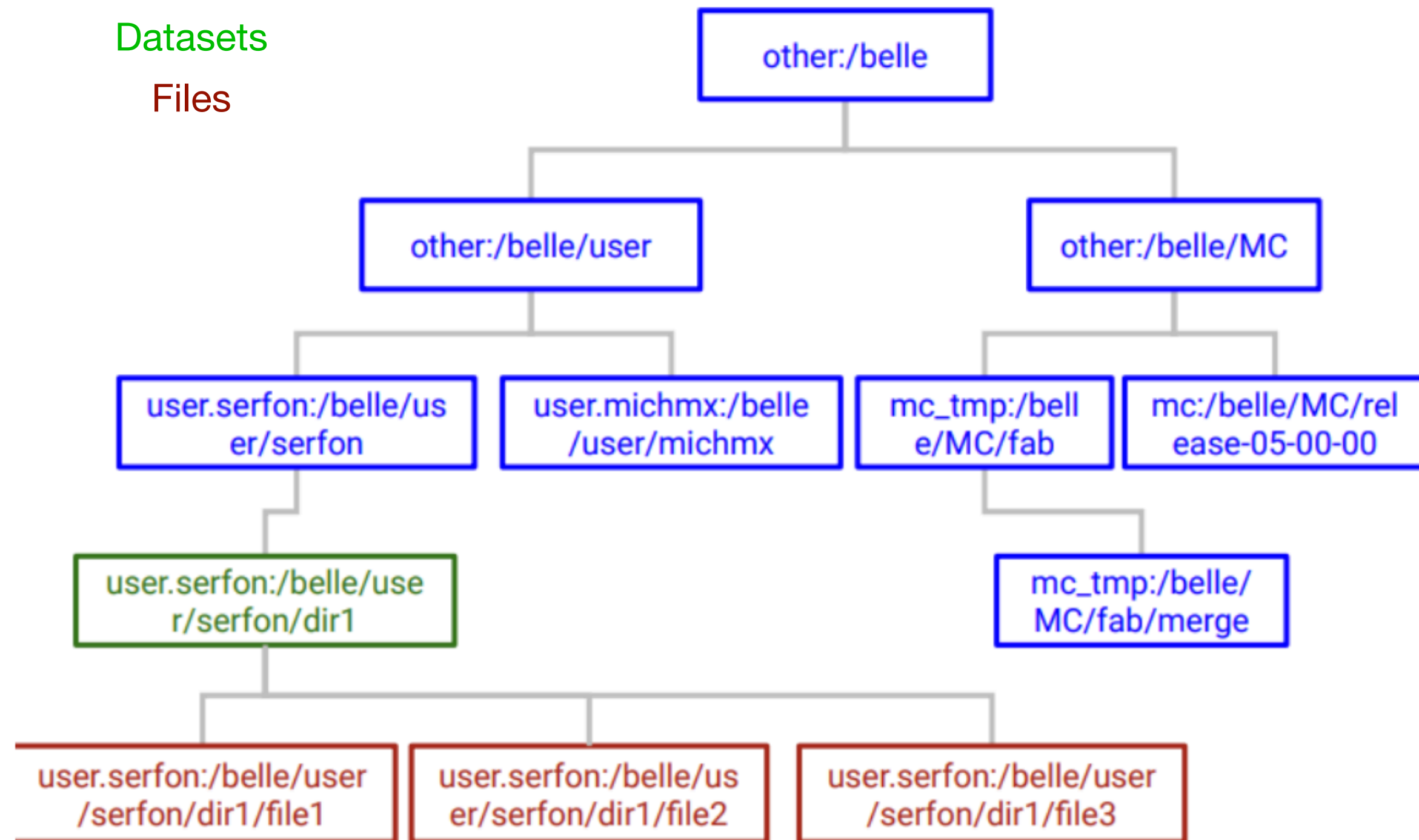- "Hash" is default policy.

<storage_path> = root://dtn-eic.jlab.org:1094/work/eic2

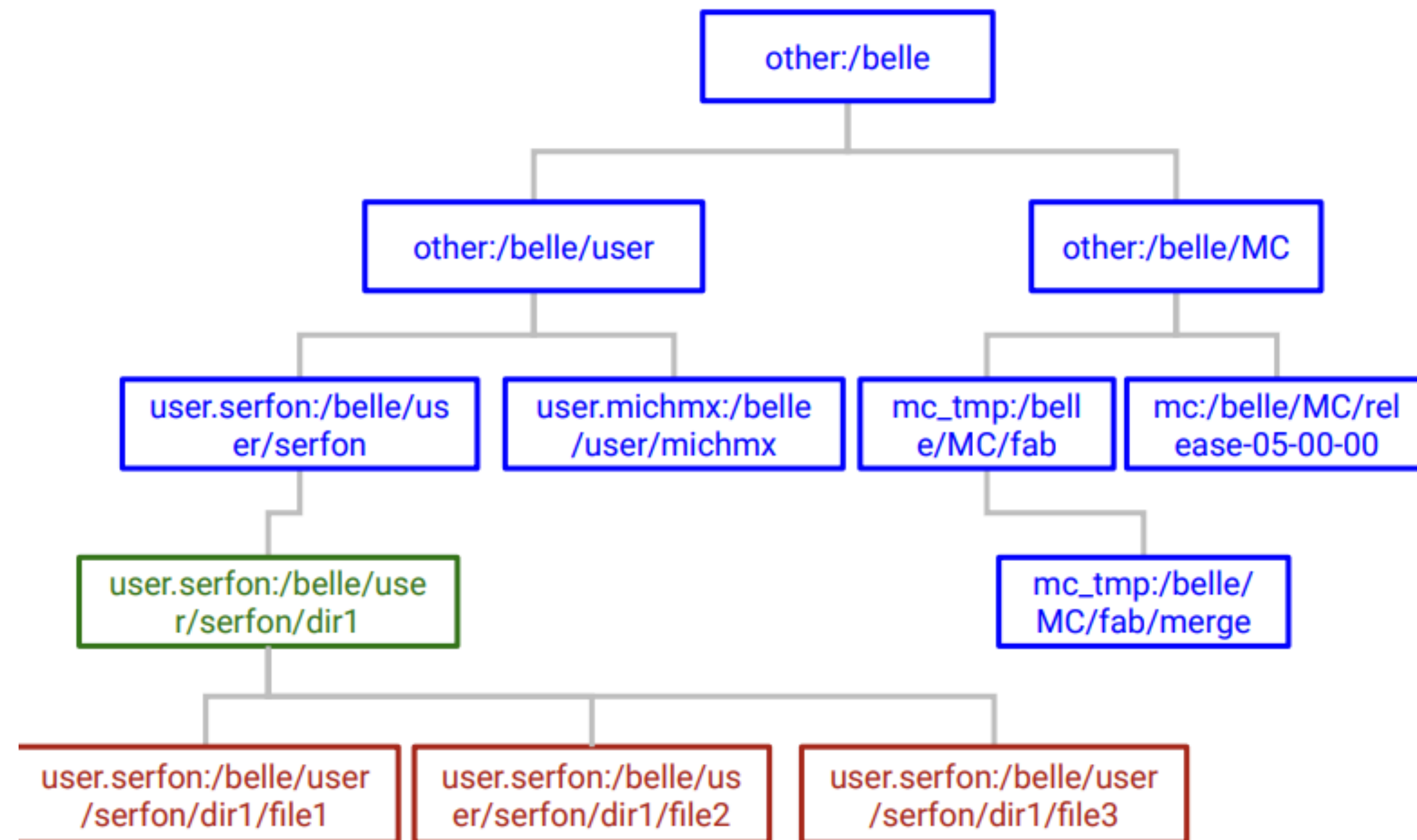# Purely Hierarchal Naming schema: Belle II

Containers

Datasets

Files



- POSIX like structure
- Pure hierarchy.
  - One DID sits at top.
  - There is chain of child dids.
- Datasets is unit of data management

Jefferson Lab

# Purely Hierarchal Naming schema: Belle II
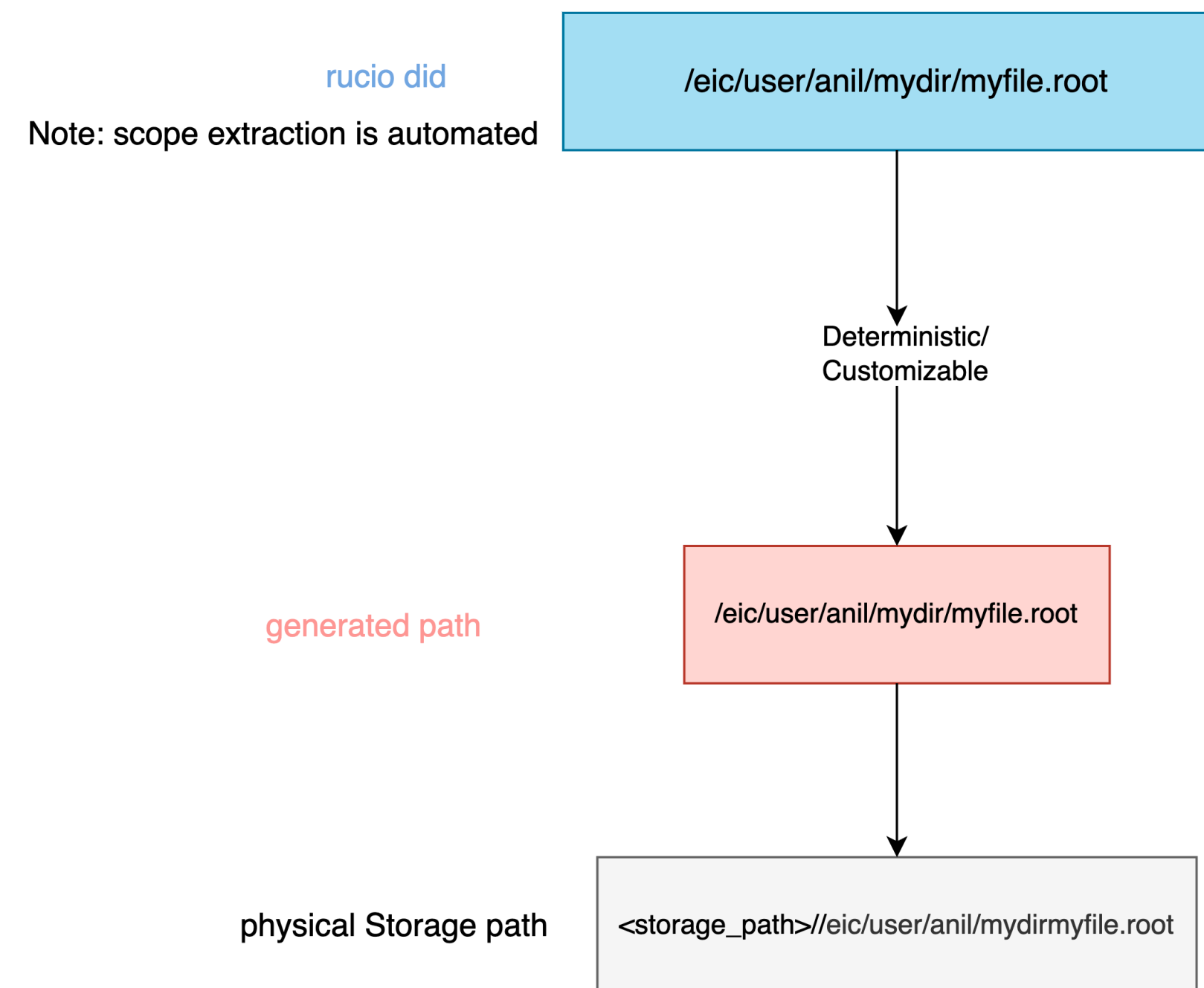


- Example : Creation of files /belle/user/serfon/dir1/file{1-3} on SE A
- The addfile method is a bulk and atomic method that :
  - Creates all the non-existing parent directories in the hierarchy
  - Creates the directory containing the files if it doesn't exists and create a Rucio rule (RSE expression ANY and grouping NONE)
  - Create the files and their replicas on A

Jefferson Lab

# Purely Hierarchal Naming schema: Physical file

- Each file did is mapped to storage path.
  - Based on some policy.
- We define what policy can be.
- Example:

rucio did
Note: scope extraction is automated

/eic/user/anil/mydir/myfile.root

Deterministic/
Customizable

generated path

/eic/user/anil/mydir/myfile.root

physical Storage path

<storage_path>//eic/user/anil/mydirmyfile.root

<storage_path> = root://dtn-eic.jlab.org:1094/work/eic2

Jefferson Lab

# Data types: speculative for EIC

- **Raw** : data from detector

- **Data** : reconstructed data
  - Full reconstructed
  - Skimmed per WG or per physics or….

- **MC**
  - Generator output
  - Geant4 Simulation output.
  - Reconstructed
  - Skimmed

- **Group** dataset.
  - Files within each working group for their own work.

- **User** output.
  - Subdivision by username and whatever sub-division workflow/workload management requires.

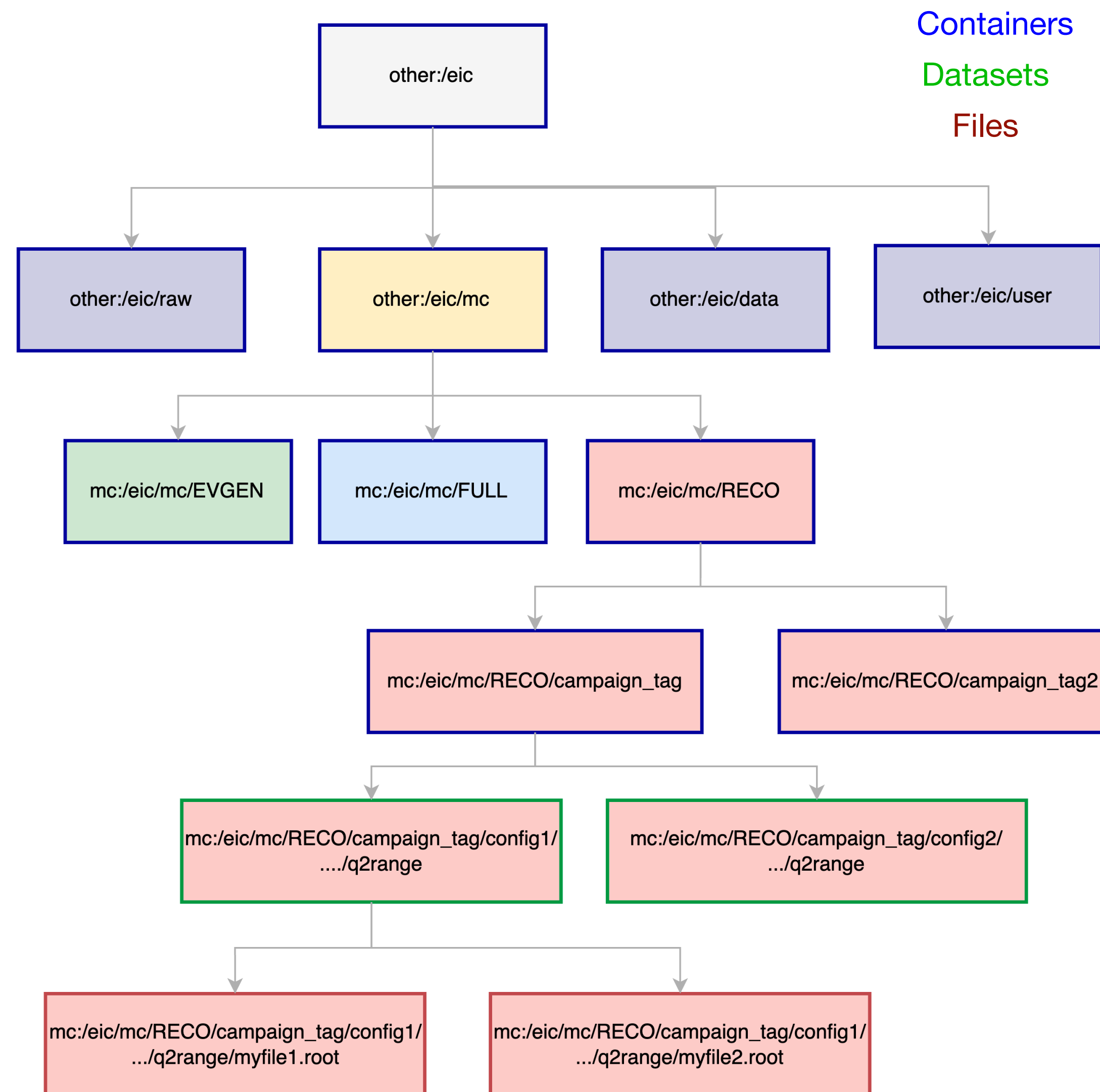Jefferson Lab

# Flat namespace/semi-hierarchal for EIC

- **Scope:name**

- Fine grained scope.

- Admin needs to be aware of scopes and create new scopes as needed.

- Workflow/central system.

  - Think about how you will mange both scope and name during registration from central system.

- User need to know the scope as well as name.

  - Think about how this info will managed/propagated in future.

- What are steps for registration.

  - One call to rucio for add-replicas file.

  - Next call to create datasets (if not created).

  - Another call to attach file to dataset

Note: I haven't worked with ATLAS distributed computing

Jefferson Lab

# Hierarchal schema for EIC:

- **Scope is deterministic.**
- **Only did's name to deal with.**
- **User just need to know the name.**
- **Storage mapping is one to one.**
- **Purely POSIX namespace**
- **Already proven to work in rucio in very large scale by Belle-II.**

Note: Name length is max 250 character.

<span style="color:blue">Containers</span>
<span style="color:green">Datasets</span>
<span style="color:red">Files</span>

```
other:/eic
├── other:/eic/raw
├── other:/eic/mc
│   ├── mc:/eic/mc/EVGEN
│   ├── mc:/eic/mc/FULL
│   └── mc:/eic/mc/RECO
│       ├── mc:/eic/mc/RECO/campaign_tag
│       │   ├── mc:/eic/mc/RECO/campaign_tag/config1/..../q2range
│       │   │   ├── mc:/eic/mc/RECO/campaign_tag/config1/.../q2range/myfile1.root
│       │   │   └── mc:/eic/mc/RECO/campaign_tag/config1/.../q2range/myfile2.root
│       │   └── mc:/eic/mc/RECO/campaign_tag/config2/.../q2range
│       └── mc:/eic/mc/RECO/campaign_tag2
├── other:/eic/data
└── other:/eic/user
```
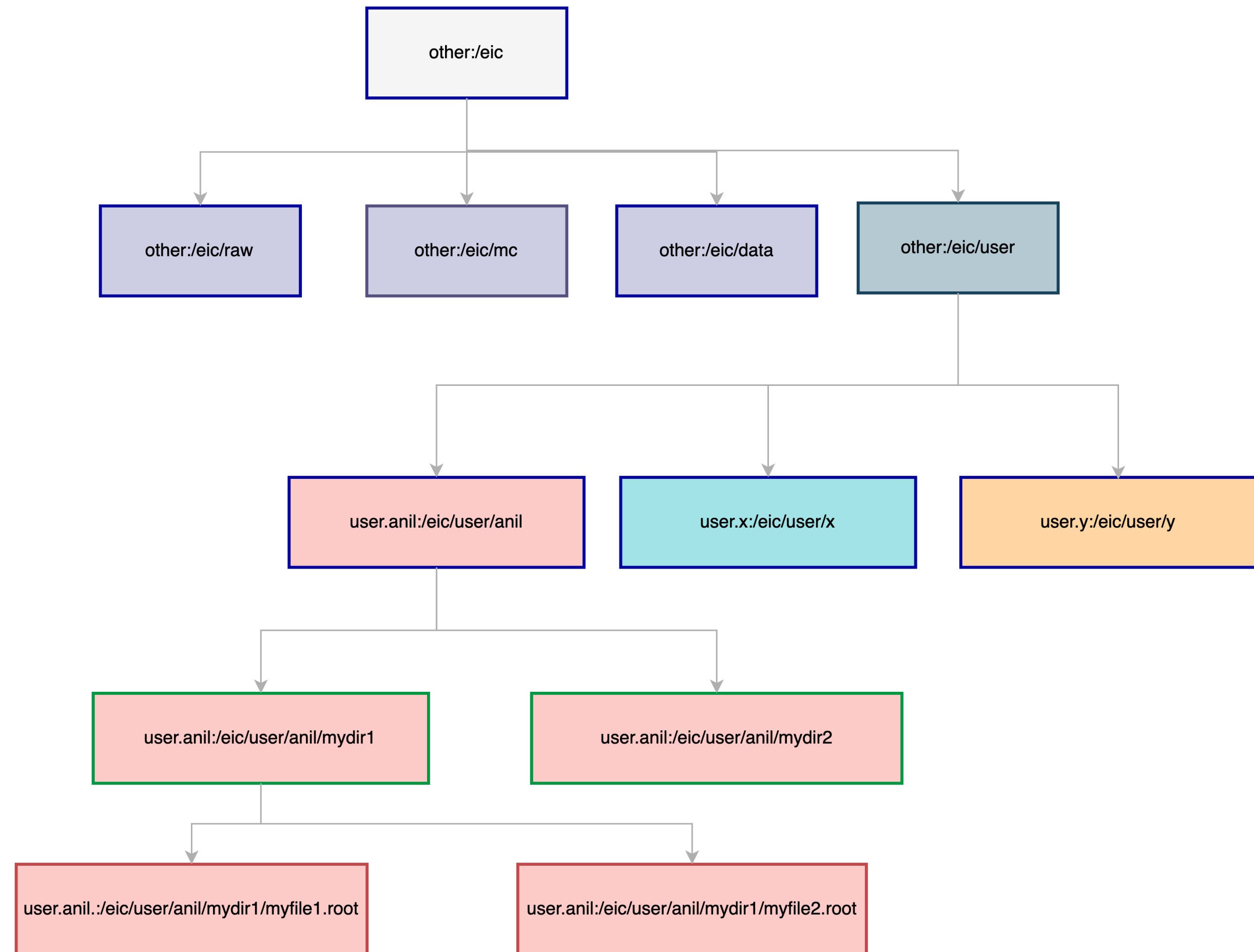
Jefferson Lab

# Hierarchal schema for EIC:

Containers
Datasets
Files

Currently only Relevant one are:
• MC
• User ?

# EIC naming: Implementation in Rucio

- Rucio provides plugins.

- Experiments can define their own policy.

- Policy Package.

  - schema: did, scope naming regex.

  - extract_scope: scope2name

  - lfn2pfn: did to storage path mapping for deterministic.

  - construct_surl : did to storage path map for Non-Deterministic.

  - permission: account and auth

- **EIC rucio policy package created and tested.**
  **JeffersonLab/eic_rucio_policy_package**

  - **Will update once the decision is taken.**

  - **Then send this to the rucio admin to add to rucio path.**
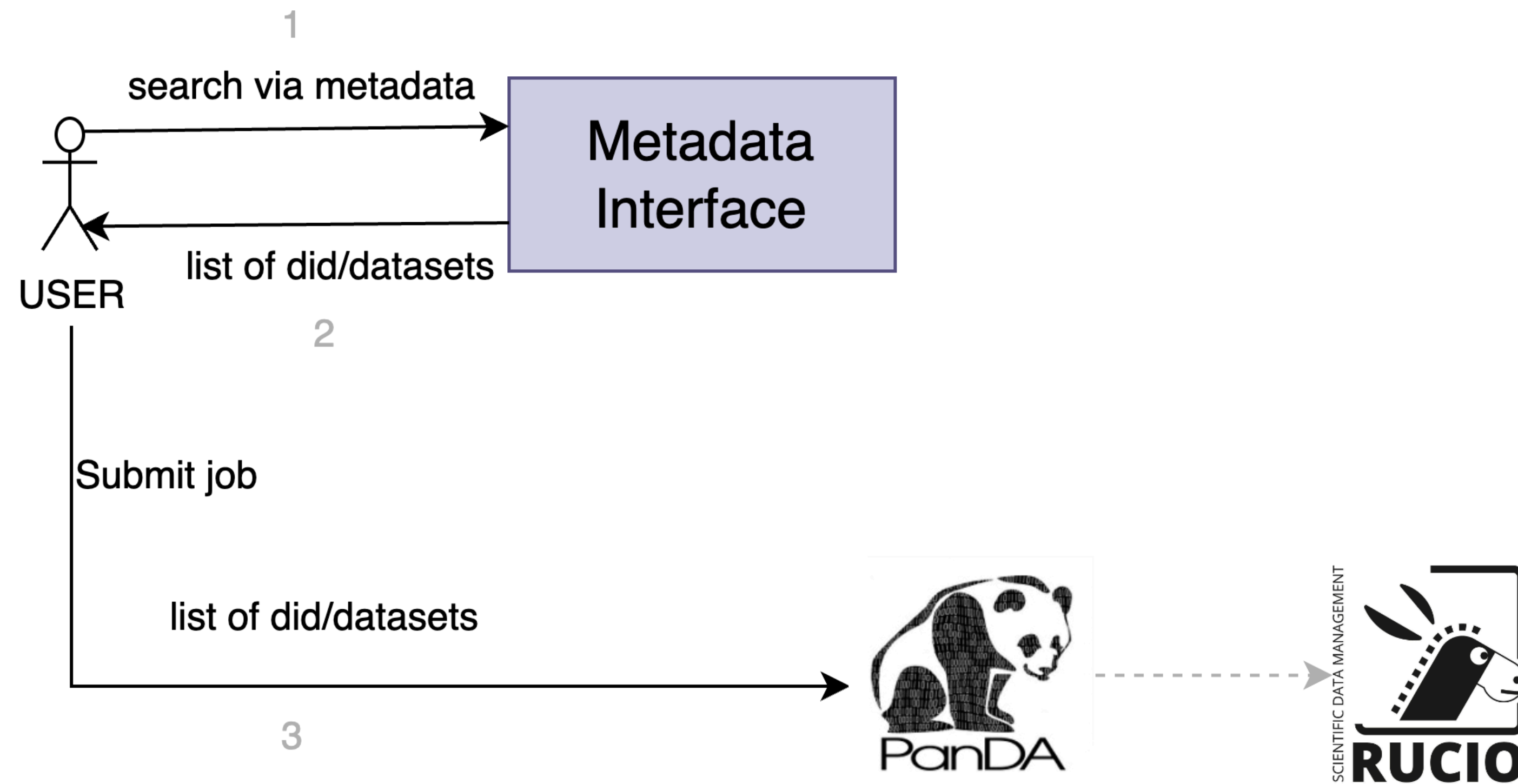
Jefferson Lab

# Backup

Jefferson Lab

# Metadata in Rucio

- Each DID in rucio have metadata associated with it.
- Metadata canoe classified.
  - Column Metadata: for processing and others ATLAS define
  - Other metadata: User define metadata.
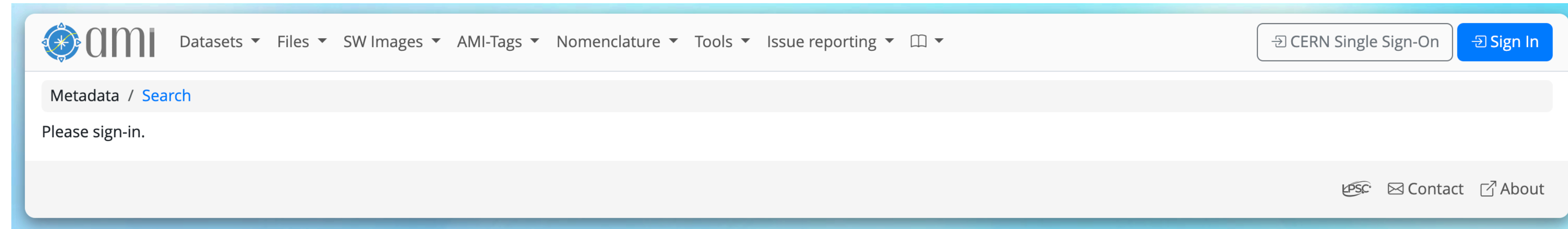    - JSON
    - External like mongoDB .

Jefferson Lab
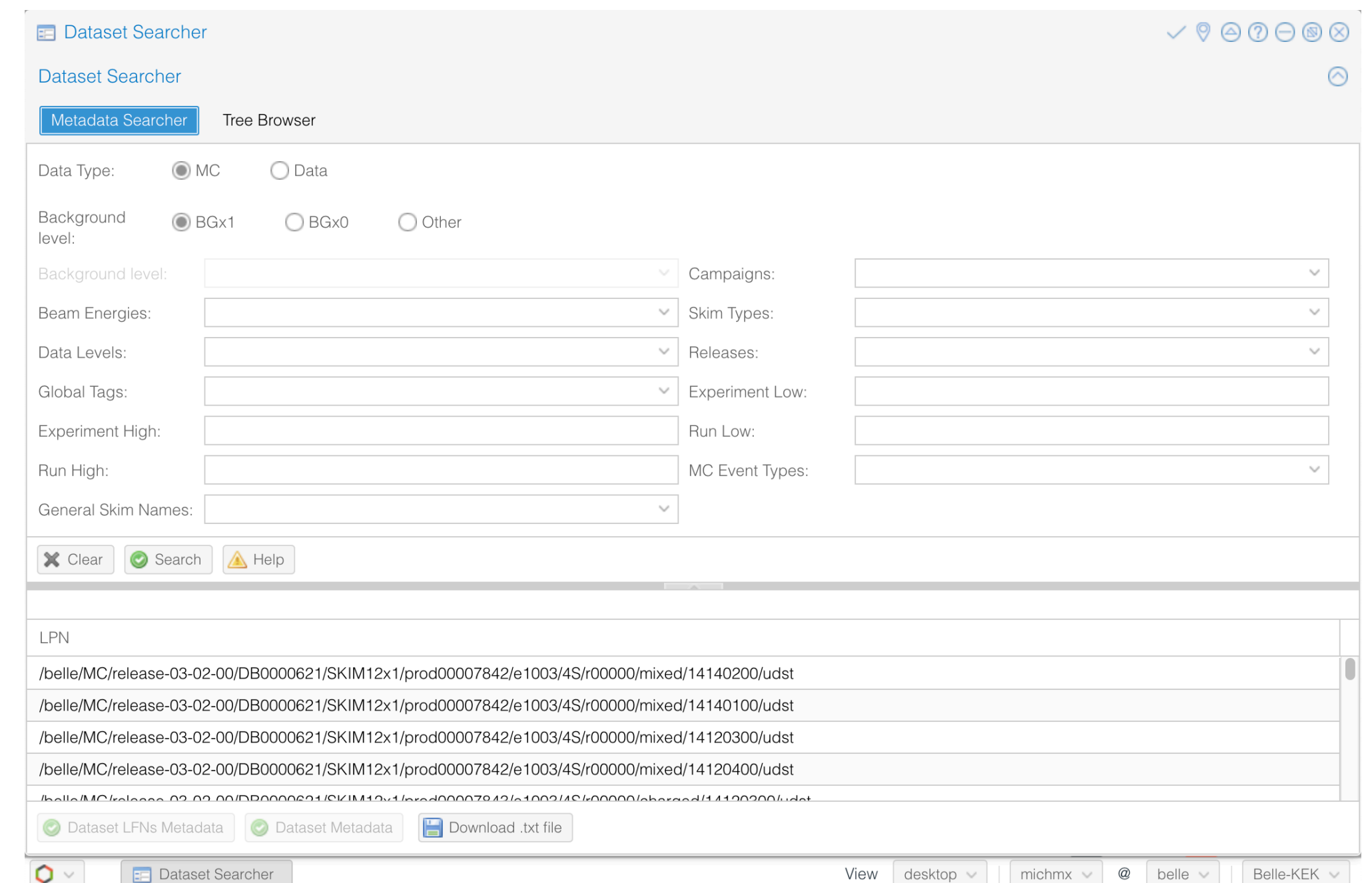
# User workflow

**Metadata Interface should be build for EIC**

# Metadata Interface : experiment example

- Atlas has **ATLAS Metadata Interface (AMI)**
  - https://cds.cern.ch/record/2649430/files/Fulltext.pdf



- Belle II has **dataset-searcher, DSS**.
  - It is working on scale for Belle II users.

# Dev for Hierarchal Namespace: Rucio side

- One PR made 2 week ago and available in 33.3.0.
  - Atomic did and metadata registration for hierarchal namespace.
- Rucio Upload for hierarchal namespace.
  - PR opened will be available soon.
  - https://github.com/rucio/rucio/pull/6492

Jefferson Lab