# Data Management Services for SRCNet v0.1
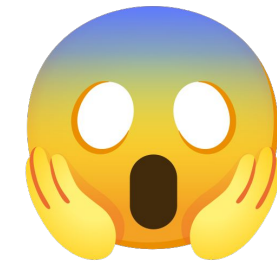
## Rob Barnsley

SKAO

# Problem reduction

SRCNet v0.1 DM

+ A&A + everything else 😱

# **A better angle in:** High level DM related flows <u>for v0.1</u>

- Consider three distinct stakeholders:
  - A **user** has a level of access in accordance with data policies
  - A **site administrator** has the required privileges to manage data at their node
  - A **service operator** has full access to all Rucio functionality; they will be responsible for the day to day running of Rucio and all that entails, and as such require the lowest level (root) access.

- Within the context of an SRCNet **user**, we need to think about:
  - Data discovery (what data is there?)
  - Data location (where is the data?)
  - Data access (how do I get a local copy of the data?)
  - Data staging (how can I stage data for remote work on a particular compute node?)

# A better angle in: High level DM related flows <u>for v0.1</u>
(continued)

- Within the context of an **operator**/**site admin of an SRCNet node**, we need to (additionally) think about:
  - Data logistics e.g. I need N copies of the data at sites X, Y and Z
  - Data ingestion
  - Data curation e.g. update metadata for a data product

- Can we use internal Rucio systems alone for all these flows? Yes, but…
  - At time of implementation, OIDC token support was in its infancy
  - Current permissions system doesn't quite satisfy need for fine grained control over e.g. listing data and metadata
  - For consistency, want permissions handled at a layer higher than Rucio & through a centralised permissions system that is used for all services
  - Don't want to force users (or site administrators) to have to learn the Rucio ecosystem and associated tooling

- To empower users and site administrators we have decided to abstract services via a set of APIs
  - This includes Rucio

# v0.1 APIs

- **What?**

  Form a significant part of the public facing component of the SRCNet that an SRCNet stakeholder will utilise to perform actions, either directly through their REST interfaces or via command line clients

- **Why?**

  Abstract interfaces to SRCNet services allow signatures to be predetermined and technology to be switched out at a later date if required

# v0.1 APIs (continued)

| Data management | Site capabilities |
|---|---|
| • Data discovery<br>• Data location<br>• Data access<br>• Data staging<br>• Data logistics<br>• Data curation | • Listing basic attributes of sites in the SRCNet<br>• Listing available storages & supported storage protocols e.g. https/xroot<br>• Listing available compute & associated services e.g. Rucio, SI, Dask, Jupyterhub |
| **Permissions** | **Authentication** |
| • Authorising access to an API's route<br>• Authorising a token exchange for a particular service<br>• Authorising access to a service | • Requesting tokens<br>• Exchanging tokens for access to different services |

# v0.1 APIs (continued)

| Data management | Site capabilities |
|---|---|
| • Data discovery<br>• Data location<br>• Data access<br>• Data staging<br>• Data logistics<br>• Data curation | • Listing basic attributes of sites in the SRCNet<br>• Listing available storages & supported storage protocols e.g. https/xroot<br>• Listing available compute & associated services e.g. Rucio, SI, Dask, Jupyterhub |
| **Permissions** | **Authentication** |
| • Authorising access to an API's route<br>• Authorising a token exchange for a particular service<br>• Authorising access to a service | • Requesting tokens<br>• Exchanging tokens for access to different services |

# Data Management API (DM API)

 + FastAPI + RapiDoc



Operator portal

User portal

# Permissions API (& interaction with DM API)



Permissions API

RBAC

GET /data/list/{namespace}

request

DM API

authz check

Rucio

DDM request

**Groups**
- data
- data/namespaces
- data/namespaces/testing
- data/namespaces/testing/owner
- services/data-management-api
- services/data-management-api/roles
- services/data-management-api/roles/admin
- services/rucio
- services/rucio/roles
- services/rucio/roles/admin

POST /authorise/exchange/{service}  Authorise Service Exchange
POST /authorise/route/{service}  Authorise Service Route
POST /authorise/plugin/{service}  Authorise Service By Plugin
GET /policies  List Policies
GET /policies/types  List Policy Types
GET /policies/{type}/  List Policies By Type
GET /policies/{type}/{service}  List Policies By Service Name
GET /policies/{type}/{service}/{version}  Get Policy By Service
GET /ping  Ping
GET /health  Health

```
1   {
2     "name": "data-management-api",
3     "type": "route",
4     "iam_subgroup_name": "data-management-api",
5     "expected_token_issuer": "https://ska-iam.stfc.ac.uk",
6     "expected_token_audience": "data-management-api",
7     "expected_service_token_scope": "data-management-api-service",
8     "version_number": 1,
9     "description": "Permissions policy defining how to authorise routes for the data management API.",
10    "roles": {
11      "any": [],
12      "namespace-viewer": [
13        "data/namespaces/{namespace}/viewer"
14      ],
15      "namespace-editor": [
16        "data/namespaces/{namespace}/editor"
17      ],
18      "namespace-owner": [
19        "data/namespaces/{namespace}/owner"
20      ],
21      "admin": [
22        "{root_group}/roles/admin"
23      ],
24      "developer": [
25        "{root_group}/roles/developer"
26      ]
27    },
28    "routes": {
29      "/data/download/{storage_area_uuid}/{namespace}/{name}": {
30        "GET": "admin or namespace-viewer or namespace-editor or namespace-owner"
31      },
32      "/data/list": {
33        "GET": "admin or developer"
34      },
35      "/data/list/{namespace}": {
36        "GET": "admin or namespace-viewer or namespace-editor or namespace-owner or developer"
37      },
38      "/data/locate/{namespace}/{name}": {
39        "GET": "admin or namespace-viewer or namespace-editor or namespace-owner or developer"
40      },
41      "/data/move": {
42        "POST": "admin or namespace-viewer or namespace-editor or namespace-owner"
43      }
```

# A (very) brief overview of (some of the) services that build off of this

# Data Discovery
## IVOA DaCHS

A software suite implementing various International Virtual Observatory Alliance (IVOA) protocols, e.g. Simple Cone Search (SCS)



### SKAO Rucio SCS

SCS query service running against an ObsCore table with a view on the Rucio database.

| Position/Name | 202.295 42.3359 |
| --- | --- |
| | Coordinates (as h m s, d m s or decimal degrees), or SIMBAD-resolvable object |
| Search radius | 1 |
| | Search radius in arcminutes |
| Table | Sort by _r ASC |
| | Limit to 100 items. |
| Output format | HTML  More output fields |

Go

**Result**

Matched: 1

Send via SAMP   Quick Plot

| Dist. [arcsec] | Obs_publisher_did | Obs_title | Obs_creator_did | Target_name | Target_class | T_exptime [s] | T_min | T_max | S_region | Em_min [m] | Em_max [m] |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| 1.08 | ivo://test.skao/~? sp3531_soda:2023-09-22-14-07-00_LOTSS-DR2_P39Hetdex19_mosaic-blanked.fits | N/A | N/A | M 51 | N/A | N/A | N/A | N/A | [204.711 47.405 | 17.84 | 24.98 |

Metadata from external postgres instance managed by Rucio plugin system

# Data Location
## IVOA Datalink

A    service    acting    as    the    glue    between    IVOA    services    and    Rucio    (via    the    DM    API)

**e.g. GET https://datalink.ivoa.srcdev.skao.int/rucio/links?id=testing:PTF10tce.fits**



Geographically nearest (or any other metric) replica

Services available at node

# Data Access (also discovery)
## astroquery extension

An extension to a widely used Python package offering astronomers a unified interface to query diverse astronomical databases using IVOA standards and protocols

```
>>> from astroquery.srcnet import SRCNet
>>> srcnet=SRCNet(verbose=True)
>>> srcnet.login()

-------------------------------------------------

Scan the QR code, or using a browser on another
device, visit https://ska-iam.stfc.ac.uk/device
and enter code XXXYYY
```

[QR code image]

```
-------------------------------------------------

Polling for token... (3/60)

Successfully polled for token. You are now logged in.

DEBUG: Access token: <redacted>
DEBUG: Refresh token: <redacted> [astroquery.srcnet.core]
DEBUG: Persisting access token to: /tmp/access_token [astroquery.srcnet.core]
DEBUG: Persisting refresh token to: /tmp/refresh_token [astroquery.srcnet.core]
```

### query_region
Query for results around a region.

```
>>> from astroquery.srcnet import SRCNet
>>> srcnet=SRCNet()
>>> srcnet.query_region(coordinates='82.1deg 12.58deg', radius=0.01)
>>>
>>> <Table length=1>
>>> dataproduct_type dataproduct_subtype calib_level obs_collection      obs_id     ... em_ucd
>>>                                                                                 ...
>>>      object              object          int16        object         object    ... object
>>> --------------- ------------------- ----------- -------------- --------------- ... ------
>>>           image                           2                   RACS RACS-DR1_0528+12A ...
```

### get_data
Get data from the datalake given a namespace and name.

```
>>> from astroquery.srcnet import SRCNet
>>> srcnet=SRCNet(verbose=True)
>>> srcnet.get_data(namespace='testing', name='PTF10tce.fits')

>>> INFO: Exchanged authn-api service token for data-management-api service [astroquery.srcnet.c
>>> DEBUG: Access token: <redacted>
>>> DEBUG: Refresh token: <redacted>
>>> DEBUG: Persisting access token to: /tmp/access_token [astroquery.srcnet.core]
>>> DEBUG: Persisting refresh token to: /tmp/refresh_token [astroquery.srcnet.core]
>>> DEBUG: Access token is valid, will not attempt token refresh. [astroquery.srcnet.core]
>>> 8248KB downloaded
```

# Data Curation and Logistics
## srcnet-oper

A command line tool with a focus on **high level admin/operator** flows, e.g.:

```
eng@dev:~$ srcnet-oper metadata set --namespace testing --name PTF10tce.fits --metadata '{"some_key": "some_value"}'
{'successful': True}
eng@dev:~$ srcnet-oper metadata get --namespace testing --name PTF10tce.fits --store science
+-----------------+------------------+-------------------------------------------------------------------------+
| Store           | Key              | Value                                                                   |
+-----------------+------------------+-------------------------------------------------------------------------+
| POSTGRES_JSON   | s_ra             | 349.7905833                                                             |
| POSTGRES_JSON   | s_dec            | 9.1960000                                                               |
| POSTGRES_JSON   | obs_id           | testing:PTF10tce.fits                                                   |
| POSTGRES_JSON   | some_key         | some_value                                                              |
| POSTGRES_JSON   | access_url       | https://ivoa.datalink.srcdev.skao.int/rucio/links?id=testing:PTF10tce.fits |
| POSTGRES_JSON   | calib_level      | 1                                                                       |
| POSTGRES_JSON   | access_format    | application/x-votable+xml                                               |
| POSTGRES_JSON   | obs_collection   | collection_testing_test                                                 |
| POSTGRES_JSON   | obs_publisher_did| testing                                                                 |
+-----------------+------------------+-------------------------------------------------------------------------+
```

```
eng@dev:~$ srcnet-oper token request
------------------------------------------------
Scan the QR code, or using a browser on another
device, visit https://ska-iam.stfc.ac.uk/device
and enter code RSVBXE
```

```
------------------------------------------------
Polling for token... (2/60)

Successfully polled for token. You are now logged in.
```

TBD:

**Data Logistics** ⌄

**POS** Request movement of existing data to another storage area

**GET** Get the status of a data movement request

**GET** Inspect a data movement request
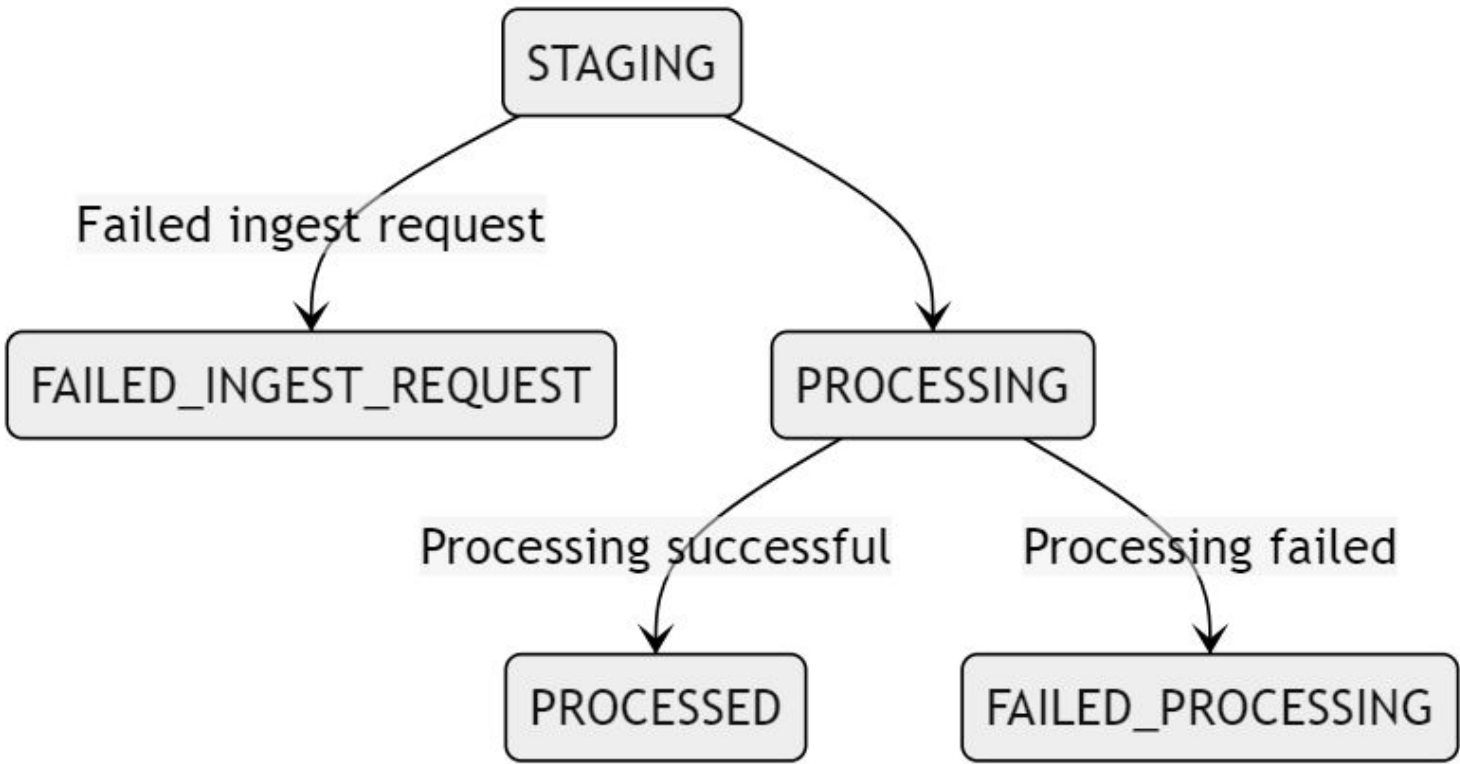
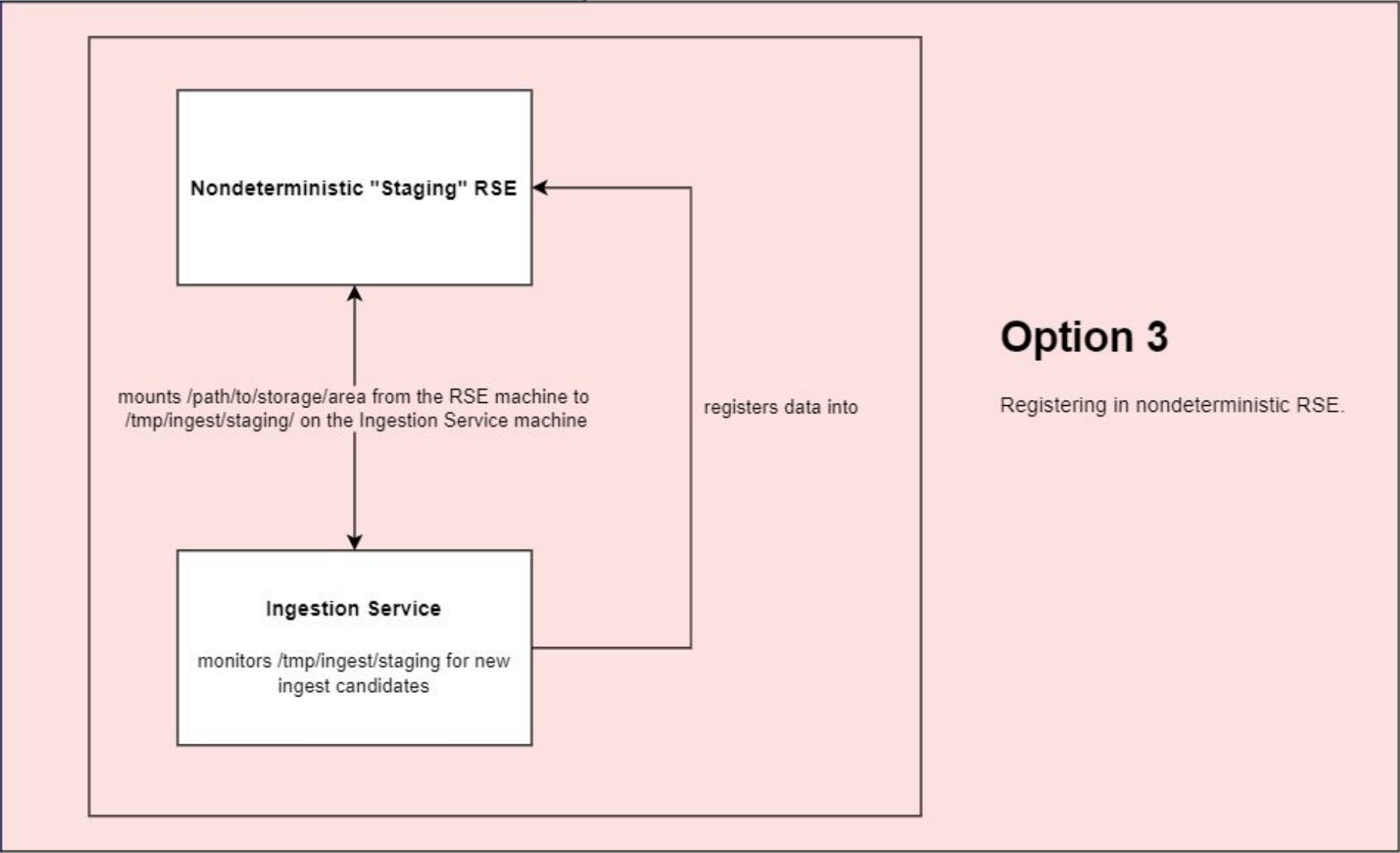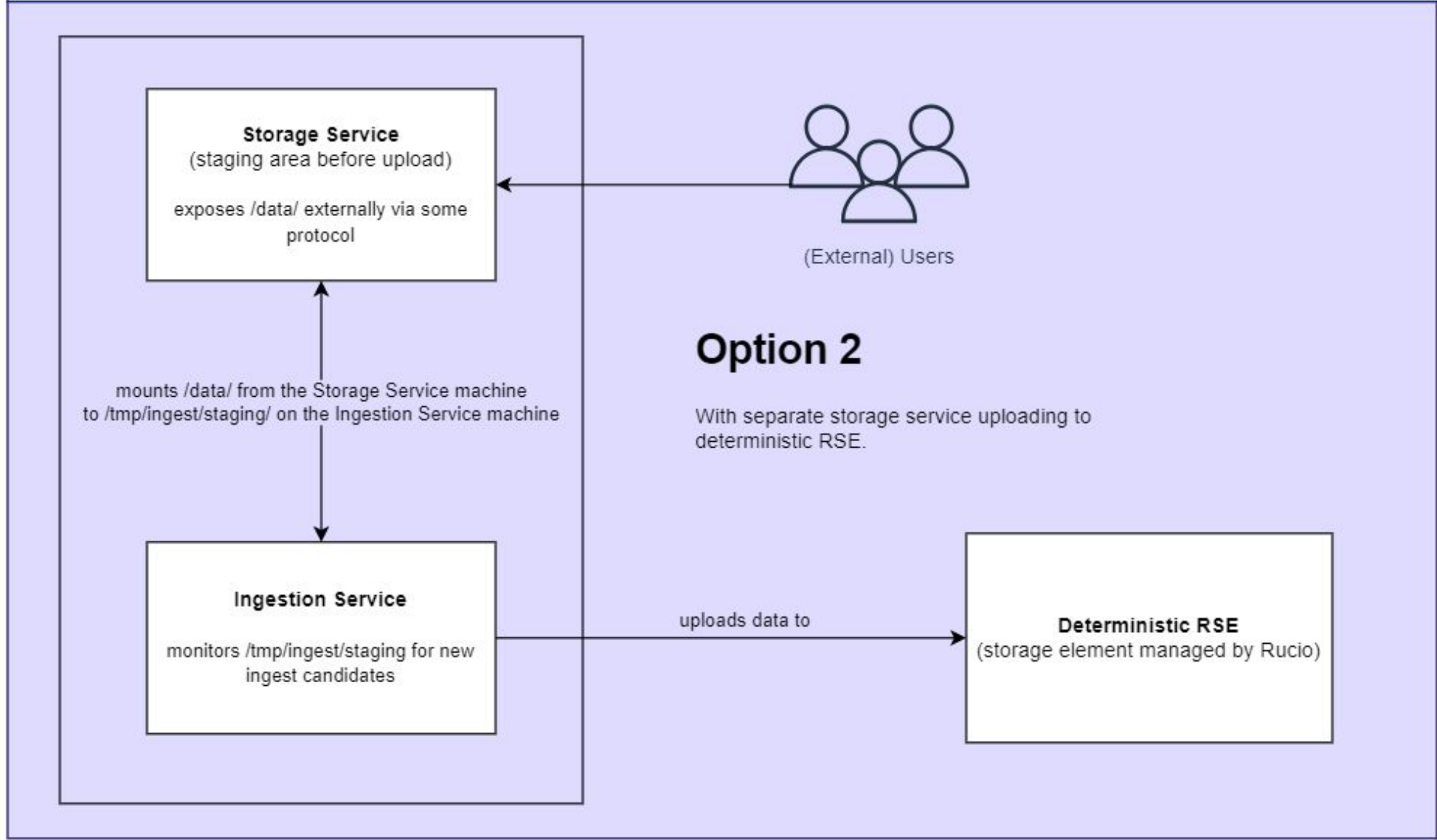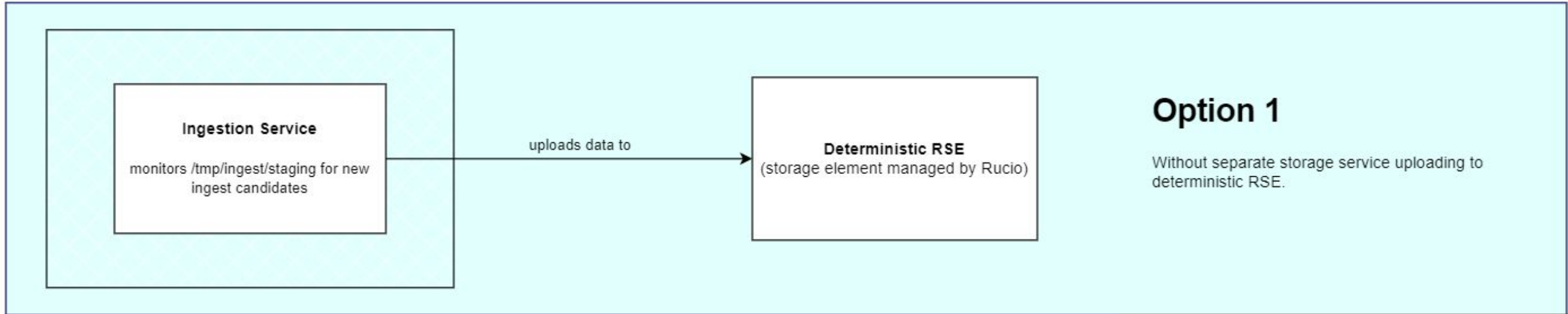**DEL** Remove data from a storage area

# Data Ingestion

Data Ingestor Service 🌐

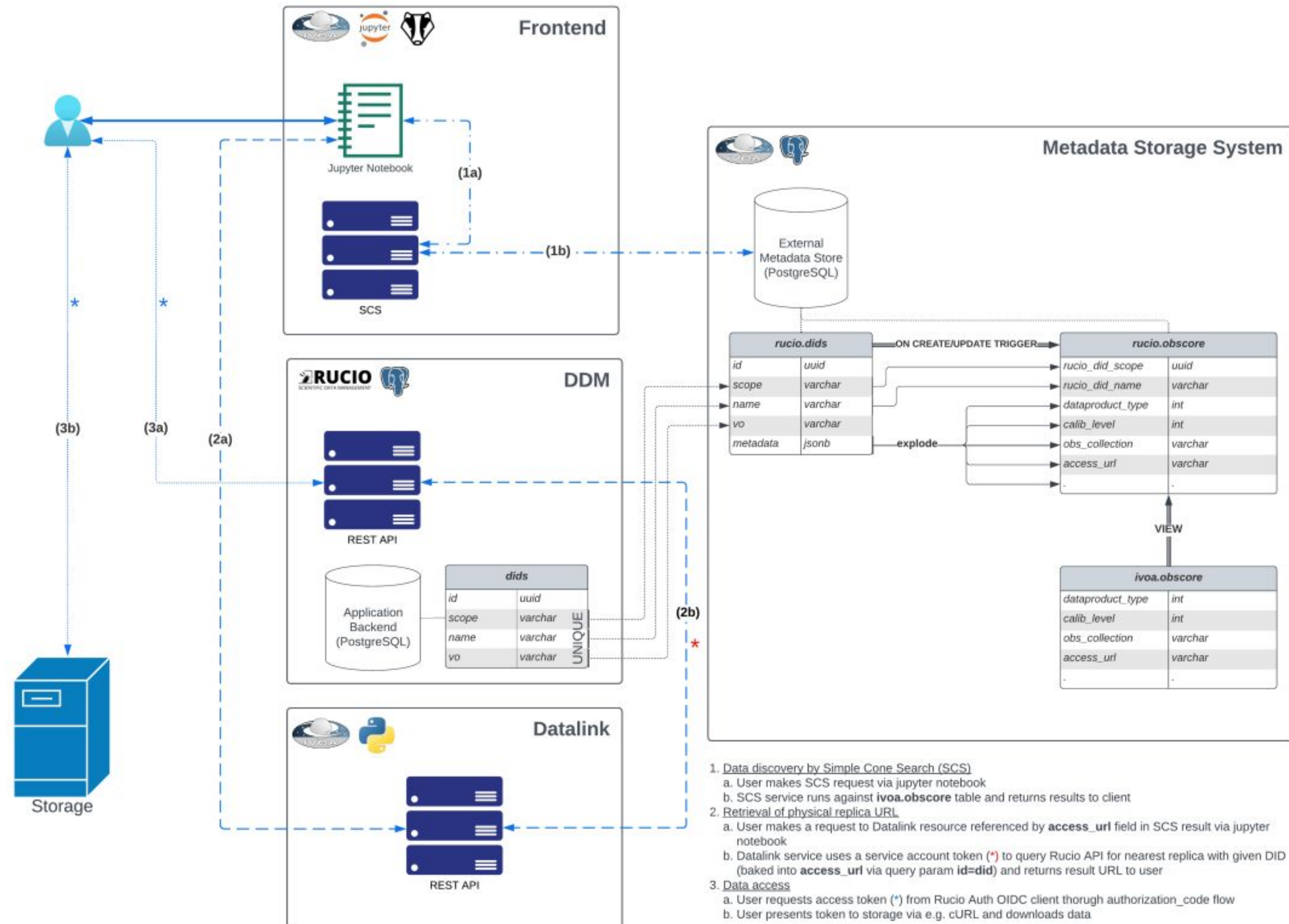A service to ingest data products into the datalake

# Conclusions

- Rucio will form the backbone of the SRCNet v0.1 DDM component
- Functionality will be hidden behind a "Data Management" API hooked into a RBAC Permissions system
- Slowly moving towards meeting the required v0.1 DM functionality
- Rapidly losing more hair (possibly related)
- Thanks!

# v0.1 DM Architecture