# Extending Rucio with modern cloud storage support

Mario.Lassnig@cern.ch

with content from *Martin Barisits, Bob Barnsley, Fernando Barreiro Megino, Johannes Elmsheuser, Mihai Patrascoiu, James Perry, Cedric Serfon, Alba Vendrell Moya, Tobias Wegner*

RUCIO
SCIENTIFIC DATA MANAGEMENT

# Cloud storage?

In recent years there has been **significant work** done **integrating Rucio with cloud storage**

**Two major angles** to consider when discussing clouds

**Technical**          Access tools, transfer protocols, monitoring, authn/z, accounting, billing, storage, …

**Organisational**     Deployed on-site or off-site
                       Centralised or distributed
                       Open or closed source software
                       Public (institute, laboratory, …) or commercial
                       In-kind contribution or paid service

It can get **complicated quickly,** e.g. …

Self-hosted MinIO S3 server on a CERN data centre VM using a centrally managed CephFS volume
WebDAV portal to self-hosted Nextcloud on a commercial hoster which points to free-tier AWS S3 storage
Experiment collaborates with commercial cloud provider and gets free storage with S3v4 protocol support

From a Rucio point of view, cloud storage is **storage that requires URL-based signatures**

Putting CephFS on top of RADOS                          -> requires some sort of storage system on top (*grid-style* storage)
Putting Ceph Object Gateway S3 API on top of RADOS      -> cloud storage

# Rucio credential mechanism

For namespace (*listing replicas*) and storage operations (*rucio upload/download*)

      Generate URL signatures **at the time of execution** of the command

      URL signatures are **generated server-side** by the Rucio server

          No deployment of secrets necessary to clients

      The account must have **schema permission** (`perm_get_signed_url`) and **account attribute** (`sign_url`)

      The Rucio Storage Element (RSE) must have several configurations applied

| | | | |
|---|---|---|---|
| scheme | https | | |
| impl | rucio.rse.protocols.gfal.NoRename | | |
| attributes | *sign_url:* s3 \| gcs \| swift | *verify_checksum:* False | *s3_url_style:* path |
| | *skip_upload_stat:* True | *strict_copy:* True | |

Credential secrets configuration

      For S3 and SWIFT compatible interfaces (e.g. MinIO, Amazon, Ceph S3 Gateway), requires an entry in `rse-account.cfg`

      For Google Cloud Storage requires the JSON credential file from Google Cloud Console

```
"d87c29b7e3294df5eacc154effd99bae": {
    "access_key": "...",
    "secret_key": "...",
    "signature_version": "s3v4",
    "region": "us-west-2"
},
```

```
{
"type": "service_account",
"project_id": "rucio-test",
"private_key_id": "be5e4aa2a0fc07e672d4051d8582c45fc630bc77",
"private_key": "-----BEGIN PRIVATE KEY-----...",
"client_email": "rucio-test@rucio-test.iam.gserviceaccount.com",
"client_id": "123456",
"auth_uri": "https://accounts.google.com/o/oauth2/auth",
"token_uri": "https://accounts.google.com/o/oauth2/token",
"auth_provider_x509_cert_url": "https://www.googleapis.com/oauth2/v1/certs",
"client_x509_cert_url": "https://www.googleapis.com/robot/v1/metadata/x509/rucio-test%40rucio-test.iam.gserviceaccount.com"
}
```

# FTS credential mechanism

When adding rules for third-party-copy, the URL signatures are generated by FTS when needed

We don't know how long transfer jobs will be in the **queue of FTS**

URL signatures are **time-limited**

**No universal TPC method** for **cloud storage to cloud storage**

Credentials need to be inserted in FTS configuration

**Secrets**              *:8446/config/cloud_storage*
                          Insert entry in specific format

**GFAL Configuration**   *:8449/fts3/ftsmon/#/config/gfal2*
                          Cannot be edited directly, has to be set by FTS admin

**HTTP Configuration**   *:8449/fts3/ftsmon/#/config/http_plugin.so*
                          Cannot be edited directly, has to be set by FTS admin

```
[S3:ATLAS-SEAL-CLOUD.CERN.CH]
ALTERNATE=true
REGION=dummy
```

```
# GCloud related options
[GCLOUD]
JSON_AUTH_FILE=/etc/fts3/gcloud_atlas.json
```

# Commercial clouds :: Google

## Google Cloud Storage

Long-term ATLAS R&D project to evaluate a grid site in the cloud
- Shoehorning X.509 certificates into commercial clouds
- Friendly administrators at sites were required
- CERN-provided certificate injected into new Google loadbalancer

Custom proxy rules to accommodate our typical Tier-1 storage
- Didn't properly work out, had to return to legacy Google loadbalancer
- Running stable since then with jobs on Google Compute Engine

Space occupancy model moved to greedy deletion

Activation of greedy deletion

|  | min | max | avg | current |
|---|---|---|---|---|
| Persistent | 0 B | 527 TB | 416 TB | 279 TB |
| Cached | 0 B | 5.54 PB | 1.66 PB | 27.6 TB |
| Temporary | 0 B | 10.3 TB | 1.88 TB | 411 GB |

**Ingress to GCS (~1Gbps / 100k files per day)**

|  | avg | total |
|---|---|---|
| User Subscriptions | 35.2 TB | 3.17 PB |
| Production Download | 33.2 TB | 2.99 PB |
| Production Input | 29.2 TB | 2.63 PB |
| Analysis Download | 23.0 TB | 2.07 PB |
| Analysis Input | 6.04 TB | 543 TB |
| Production Upload | 3.59 TB | 323 TB |

**GCS internal usage (~500Mbps / 50k files per day)**

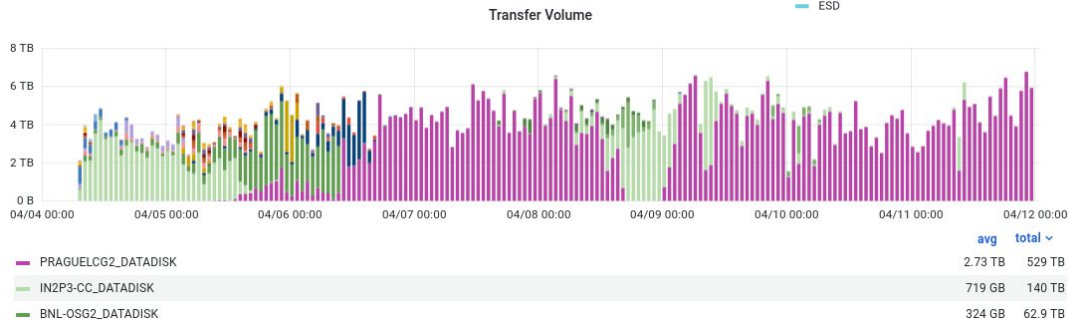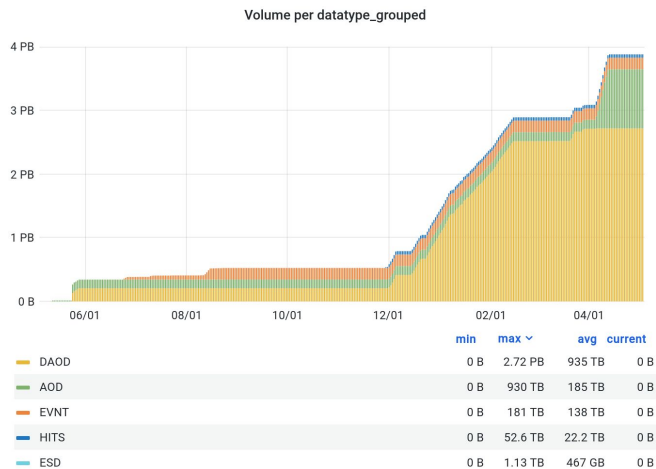|  | avg | total |
|---|---|---|
| Production Download | 33.2 TB | 2.99 PB |
| Analysis Download | 23.0 TB | 2.07 PB |
| Production Upload | 3.59 TB | 323 TB |
| Production Output | 3.56 TB | 320 TB |
| Analysis Upload | 553 GB | 49.8 TB |
| Analysis Download Direct IO | 345 GB | 31.1 TB |

# Commercial clouds :: SEAL

**SEAL Storage Technology**

Distributed cloud storage backed by Interplanetary File system (IPFS) and Filecoin (FIS)

Offered 10PB of storage to ATLAS for a long-term R&D project

Sealing process of data for long-term archival and safe-keeping

**Rucio Integration**

Very smooth integration with standard URL signature mechanism

Same trick used: SEAL administrators injected CERN-provided certificate in their loadbalancer

Gradual selection and transfer of datasets

**SEAL is funding a full-time development position to improve cloud support in Rucio**

# Commercial clouds :: SEAL

**SEAL Storage Technology**
> Distributed cloud storage backed by Interplanetary File system (IPFS) and Filecoin (FIS)
> Offered 10PB of storage to ATLAS for a long-term R&D project
> Sealing process of data for long-term archival and safe-keeping

**Rucio Integration**
> Very smooth integration with standard URL signature mechanism
> Same trick used: SEAL administrators injected CERN-provided certificate in their loadbalancer
> Gradual selection and transfer of datasets

**SEAL is funding a full-time development position to improve cloud support in Rucio**

### Volume per datatype_grouped



| | min | max | avg | current |
|---|---|---|---|---|
| DAOD | 0 B | 2.72 PB | 935 TB | 0 B |
| AOD | 0 B | 930 TB | 185 TB | 0 B |
| EVNT | 0 B | 181 TB | 138 TB | 0 B |
| HITS | 0 B | 52.6 TB | 22.2 TB | 0 B |
| ESD | 0 B | 1.13 TB | 467 GB | 0 B |

### Transfer Volume



| | avg | total |
|---|---|---|
| PRAGUELCG2_DATADISK | 2.73 TB | 529 TB |
| IN2P3-CC_DATADISK | 719 GB | 140 TB |
| BNL-OSG2_DATADISK | 324 GB | 62.9 TB |

# Commercial clouds :: AWS

## Now… Amazon was a different story

This is where it gets complicated

It worked out of the box for a year (thanks to DigiCert) until they changed to their own custom CA

In ATLAS there's a US Tier-3 (FRESNO) with a considerably sized investment

Setting this up was… challenging: 6+ months of trial & error lead to this short [document](document)

# Rucio ROOT Direct-IO mechanism

For interactive analysis and other stream processing cases: remote reads are used

> The path returned from list-replicas usually can be fed straight into `TFile::Open()`
>
> `TFile::Open("https://mycloud:443/file.root?url_signature=1234");`
>
> S3 protocol **does not provide multi-range** byte requests
> Amazon requires CloudFront CDN, which does **multi-range translation**
> Others, e.g., Google Cloud Storage or MinIO, do not have this translation layer

Workaround is simply to disable multi-range requests through Davix

> Have to append URL options to emulate:  `#multirange=false&nconnections=30`
>
> This is highly client dependent, *one size fits all* not really applicable
> We will have to investigate if we should simply make Rucio reply with these options
> Would require a potential hint to list-replicas (`--use-for-direct-io=30`) or similar solution

# Future work

Configuration / Setup
> Complicated, but grew organically from the ongoing Cloud R&D projects
> Needs a complete overhaul: esp. naming of attributes

Already identified features that we will need for production-level integration
> Access control right now is all-or-nothing, needs to be more fine grained
> Smarter peering mechanism
>> Static multihop distance config vs. dynamic cloud regions
>> The concept of cloud regions is missing completely
> Security considerations
>> Right now completely dependent on X.509 with DNS-injection trick
>> Clouds typically support OpenID/OAuth2 flows, should be helpful for token migration work
> Throughput and cost control not yet implemented, if you have the access rights you get the "full cloud power"
> Bucket-copy transfer tools, no need to go through FTS for this
> Cloud boosting option: Dynamically spend currency for extra throughput/storage
> Data lifetime considerations / different cloud QoS costs

Theoretical R&D studies:     Simulation and evaluation of cloud storage caching
                             (Tobias Wegner's PhD)
> Temporary cloud bursting to improve workflows needing tape recalls
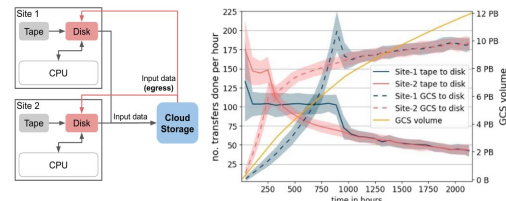> Demonstrates 15% improvement in job times



Fig. 8 The solid blue and red line show the number of transfers from tape to disk per hour for each site. The dashed lines show the number of transfers from GCS to disk per hour for each site. The orange line shows the GCS volume used.