

Rucio & Cloud Storage

Mario Lassnig

with content from Fernando, Rob, Johannes, James, Mihai, Cedric, Alba, Tobi

5th Rucio Community Workshop

10 - 11 November 2022, Lancaster UK

<https://indico.cern.ch/event/1185600/>



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	-> <i>grid-style</i> 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 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**

You cannot TPC from cloud storage to cloud storage

Credentials need to be inserted in FTS configuration

Secrets

`:8446/config/cloud_storage`

Insert entry in specific format

S3.atlas-amazon-cloud.cern.ch Service API

User	VO roles	Access token	Access secret	Request token	Request secret
.	/atlas/Role-production /Capability=NULL				

Save Delete

GFAL Configuration `:8449/fts3/ftsmon/#/config/gfal2`

Cannot be edited directly, has to be set by FTS admin

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

HTTP Configuration `:8449/fts3/ftsmon/#/config/http_plugin.so`

Cannot be edited directly, has to be set by FTS admin

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

Integration with commercial clouds 1/2



Google Cloud Storage

- First integration R&D project incredibly painful
- Shoehorn X.509 certificates into commercial cloud
- Friendly administrators at sites
- CERN-provided certificate injected into Google loadbalancer
- Custom proxy rules to accommodate our typical Tier-1 storage setup

Frontend

Protocol ↑	IP:Port	Certificate	SSL Policy	Network Tier ?
HTTPS		atlas-google-europe-west1-cern-provided	GCP default	Premium

Routing rules

Hosts ↑	Paths	Backend
All unmatched (default)	All unmatched (default)	atlas-europe-west1-litter
*	/*	atlas-europe-west1-litter
*	/atlascratchdisk/*	atlas-europe-west1-scratchdisk
*	/atlasdatadisk/*	atlas-europe-west1-datadisk

<input type="checkbox"/>	Name ↑	Created	Location type	Location	Default storage class ?	Last modified ?	Public access ?
<input type="checkbox"/>	atlas-europe-west1-datadisk	Sep 21, 2022, 8:26:54 AM	Region	europe-west1	Standard	Nov 1, 2022, 3:35:01 PM	Not public
<input type="checkbox"/>	atlas-europe-west1-litter	Sep 21, 2022, 10:47:42 AM	Region	europe-west1	Archive	Sep 21, 2022, 10:47:42 AM	Not public
<input type="checkbox"/>	atlas-europe-west1-scratchdisk	Sep 21, 2022, 8:27:21 AM	Region	europe-west1	Standard	Sep 21, 2022, 8:27:21 AM	Not public

SEAL Storage Technology

- Distributed cloud storage, offered 10PB of storage to ATLAS for a long-term R&D project
- Integration went relatively smooth with standard URL signature mechanism
- Same trick used for integration: SEAL administrators injected CERN-provided certificate in their loadbalancer

Rucio ROOT Direct-IO mechanism



For interactive analysis and other processing cases, do not download file but do remote reads

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 required CloudFront CDN anyway, 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* maybe 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 in Rucio 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 transfertool, 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 (Tobi's PhD)

Temporary cloud bursting to improve workflows needing tape recalls

Demonstrates 15% improvement in job times

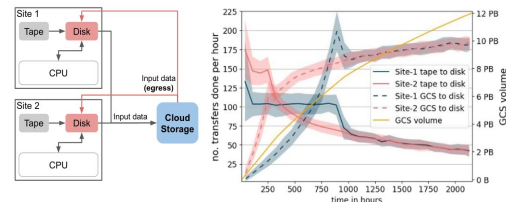


Fig. 8 The solid blue and red line show the number of transfers from tape to disk per hour for each site. The dashed lines show the number of transfers from GCS to disk per hour for each site. The orange line shows the GCS volume used.