

Virtual Research Environment updates

EOSC-Future DM monthly meeting 14.06.2022

Elena Gazzarrini



Updates

Amsterdam ESCAPE DIOS in-person meeting

- <u>https://indico.in2p3.fr/event/26928/</u>
- Analysis facilities discussion
- Storage facilities connections

Solved DLaaS issues related to new version of Rucio-jupyterlab package

• Improved Rucio knowledge

Current state of DLaaS: certificates expired!

• SSL certificates passed as secrets to Rucio authentication server expire after 2 years → dependent on CERN admin, working on it



Main infrastructure - ESCAPE Rucio instance

The ESCAPE infrastructure on which the DLaaS sits on is composed of:

- Main Server
 - handles REST requests to the resources (Apache HTTP Server)
- . Authorization Server
 - handles REST authentication/authorization requests (Apache HTTP Server)
- . WebUI Server
 - Rucio GUI (Apache HTTP Server)
- . Daemons
 - . Python modules that interact with the DB







Cluster set-up

- Cluster creation on Openstack (magnum)
 - github public repository, bootstrapped flux on it
- Secrets management with Mozilla SOPS
 - Rucio expects secrets (certificates, DB passwords, OIDC client ID, etc.) before starting the service
 - .p12 certificates, split into host and key files
 - gridCA certificates
 - TLS certificates
 - client_id and client_secret of the Rucio Admin account created with the Identity Provider (ESCAPE IAM for us) → needed for JSON web tokens (JWTs) and OAuth2.0 authentication and authorization with Rucio
 - Database credentials (Oracle, PostgreSQL, MySQL/MariaDB are currently supported)
- Network
 - 2 nodes of the cluster are set as K8s Ingress controllers
 - They accept traffic from outside, and load balance it to pods (containers) running inside the platform
 - Set NGINX as Ingress controller, as it is most popular and open source way to have a reverse proxy (to protect the server) + supports X509
 - External traffic
 - LanDB-alias is set for eosc-auth.cern.ch, eosc-main.cern.ch, eosc-webui.cern.ch (by default, the CERN outer perimeter firewall blocks incoming access to systems on the CERN site → need to request to open for the lanDB-alias property configuration)
- RSE (storage) management and configuration
 - RSEs can be configured in CRIC for easier RSE management
 - script to sync the service with new RSE creation via JSON request
 - One RSE on EULAKE-1 at CERN
 - One with CNAF-STORM
- Monitoring
 - Logging producer requested for eosc-future cluster
 - Logs are injected into Grafana for monitoring dashboards
 - FTS service is already configured to push data into CERN Monit
- Helm installing server, daemons, webui
 - The rucio helm charts can be found in the <u>rucio repo</u>
 - Each Helm Release will start a deployment of each of:
 - Server
 - Webui
 - Daemons figuring out minimal ones needed to have the service running
- Jupyter notebook + rucio jupyterlab extension



DLaas Feature Highlights

The goal of the service is to abstract the complexities of the Data Lake from the scientists. This way, scientists can focus their time on doing science instead of data procurement.

- Multiple notebook <u>environment</u> selection
- Rucio data browser (with scope browser and wildcard search)
- "Add to shopping cart" for data catalogue
 - DID is attached as a metadata in the Notebook file
- Injects a variable containing local file path, ready to be used
- Direct file upload to Rucio
- Scratch space for large files (EOS FUSE mount)



Deployment

- Deployed in Kubernetes @ CERN Openstack, using <u>Zero-to-JupyterHub Helm chart</u>.
 - <u>https://escape-notebook.cern.ch</u>
- CI/CD
 - o <u>Gitlab CI</u> Container build
 - Flux2 Kubernetes manifest
- OAuth authentication using ESCAPE IAM.
- Uses Rucio JupyterLab Extension in Replica mode
 - Connected to ESCAPE Data Lake (escaperucio.cern.ch; rucio_host)
 - Automatically preconfigured to use OIDC authentication (RUCIO_DEFAULT_AUTH_TYPE)
 - Has a FUSE mount to EULAKE-1 RSE (EOS; RUCIO_RSE_MOUNT_PATH)
 - Making files available means creating a replication rule to move files to EULAKE-1 (RUCIO_DESTINATION_RSE)



(*) <u>CONFIGS</u>



Alternative Rucio instances - inspiration

- SKAO: <u>https://gitlab.com/ska-telescope/src/ska-rucio-prototype</u>
- ATLAS: <u>https://gitlab.cern.ch/atlas-adc-ddm/rucio-k8s-setup/-/blob/master/README.md#L48-118</u>
- Our DL:

https://indico.cern.ch/event/1069544/contributions/4497649/attachments/2311244/3933259/Data%20Lake%20as%20a%20Service%2 Ofor%20Open%20Science.pdf

- ASTRON replicating DLaaS:

https://git.astron.nl/groups/astron-sdc/escape-wp5/-/wikis/Meeting-Notes/Spring-2022-Busy-Week/Replicating-Data-Lake-as-a-Service

Reading list

- The ESCAPE Data Lake: The machinery behind testing, monitoring and supporting a unified federated storage infrastructure of the exabyte-scale https://www.epj-conferences.org/articles/epjconf/abs/2021/05/epjconf_chep2021_02060/epjconf_chep2021_02060.html
- ESCAPE Data Lake: Next-generation management of cross-discipline Exabyte-scale scientific data
 <a href="https://www.epj-conferences.org/articles/epjconf/abs/2021/05/epjconf_chep2021_02056/epjconf_chep2020056/epjconf_chep2020056/epjconf_chep2020056/epjconf_chep2020056/epjconf_chep2020056/epjconf