# Data organization and data management

Artem Petrosyan (MLIT JINR)
October 17, 2023

# Our storage resources

- Disc: EOS

- Tapes: CTA

- External (bright, but near future)

# EOS

- Users dir, there are already some data

  - /eos/nica/spd/users

- Prod dir, there will be some data, hopefully, pretty soon

  - /eos/nica/production

- It is wrong to think that EOS storage is infinite

```
lxui02:~ > eos quota /eos/nica/spd

By group:
┌─> Quota Node: /eos/nica/spd/
```

| group | used bytes | logi bytes | used files | aval bytes | aval logib | aval files | filled[%] | vol-status | ino-status |
|-------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|
| project | 287.28 TB | 279.33 TB | 661.26 K | 1.00 PB | 500.00 TB | 0 | 28.73 % | ok | ignored |

- Bright future: a dedicated endpoint and EOS instance for our data, say /eos/spd

- Far bright future: to have a separated quota for users and for production

# CTA

- The CERN Tape Archive (CTA) is the tape backend to EOS

- There is ongoing work and tests here in MLIT to enable CTA

- Once its done we'll be able to write our the most valuable data there

- Write access will be granted only for a service user

- We expect to have quota at least 10 times larger than the one that we have at EOS from the beginning

- In the future, for production, we expect to use EOS as a disk pool for the data during its processing, not as a long term storage

# External storage

- There are already propositions to store some our data on the external storages, for example, at Minsk

- In order to start doing this we must build a data catalog to know where and which our data is stored to avoid creating a grey data

- The natural way to manage data on the external source is to do it through the data management service

# Data management service

- We've deployed an instance of Rucio data management system and now we're in the middle of the configuration process

- Rucio provides not only a data catalog, but also a metadata catalog, can manage replicas and data integrity on the different storages, allows to define a lifetime of the data basing on its type, etc.

- Manpower: 1 person holds a laborant position at MLIT till the end of June, 2024

# Summary

- We have everything to store our data at the service level: disk, tape storages

- We need to make efforts in the following directions: data and metadata catalog preparation, data types definition, data lifetime definition, data management service configuration

- The most important task now is to define a "business processes", identify data types and define lifetime for each type; basing on data type we'll organize its storage on the suitable storage type

# Thank you!