

SPD data organization and naming convention

Artem Petrosyan on behalf of the SPD offline computing team
SPD Physics and MC meeting N35
November 22, 2023

Our storage resources

- JINR
 - Discs: EOS
 - Tapes: CTA
- It was announced recently that JINR will be able to cover 25-30% of our declared storage resources needs
- External (bright, but near future) both discs and tapes — 75%

EOS

- Users dir, there are already some data
 - /eos/nica/spd/users
- 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

- We work to organize a dedicated endpoint and EOS instance for our data, /eos/spd
- Some bright future: to have a separated quota for users and for production

CTA

- The CERN Tape Archive (CTA) is the tapes 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
- We expect to have quota at least 10 times larger than we have at EOS from the beginning (already declared in the latest version of the SPD TDR)
- In the future, we expect to use EOS **only** as a disk pool for the data during its processing, not as a long term storage

External storage(s)

- There are already propositions to store some our data at the external storages, for example, at PNPI, SPbU, INP BSU (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 dark 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 (<https://rucio.cern.ch/>) 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 different storages, allows to define a lifetime of the data basing on its type, activity, etc.
- Rucio can manage quotas at the level of individual users, groups: we invite groups to identify themselves and get much larger quota ;-)
- We plan to use Rucio for managing not only data at external storages, but inside JINR as well to build a single namespace for our data

Datasets naming convention

- In order to ease metadata catalog navigation, data filtration, identification, etc., we propose the following naming convention for datasets (a set of files which represents results of organized calculations)
- 2050.DATA.250LT.minbias.27189.RAW.636763fd78df7d.0
- Open questions/issues:
 - Do we need data taking period?
 - Special runs: alignment, calibration? Should be marked in [Desc] — (convention in convention)?
 - Are dataset names case-sensitive?
 - Allowed symbols for text filed names

Grouping tier	Field	Description	Example
0	[YEAR]	Main Scope - the year of data production	2050
1	[MC DATA]	Real data or simulated data	DATA
2	[energy][polarization]		250LT
3	[desc]	Short name of physics aim	minbias
4	[RunNumber]	Run number for DATA, ID for MC	27189
5	[data type]	EVGEN, SIMUL, RECO....	RAW
6	[<u>DatasetUID</u>]	unique ID of the dataset	636763fd78df7d
7	[Version]	for reprocessing	0

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
- We need to
 - Fix dataset naming convention (it can be changed later, but we need to start from something)
 - Identify group activities and define their quotas in Rucio
 - Define a “business processes” for each group activity to understand data flows
 - Identify data types and define lifetime for each type globally to define storage rules
- Basing on data type we’ll organize its storage on the most suitable storage resource

Thank you!