
The First Disk-based Custodial Storage for the ALICE experiment

AHN SANG-UN @ EOS WORKSHOP, 1-4 MARCH 2021

Outline

- Introduction
- System Architecture
- QRAIN Layout
- Current Status
- ALICE Integration
- Monitoring
- Plan
- Summary

Introduction

- Replacing tape library (+3PB) with disk-only storage for archiving @ KISTI for ALICE experiment
 - Simpler architecture, less operational efforts, cost-effective comparable to tape
- Found domestic suppliers of high-density(> 60 disks/box) JBOD models
- Relying on EOS erasure coding implementation (RAIN layout) for data protection
- About 1M CHF budget (2019) included
 - 18 High-density JBOD boxes (84 disks/box \simeq 18PB raw capacity)
 - 9 Servers for EOS front-end nodes (12Gbps SAS HBAs, 40Gbps uplinks + switches)
- Providing production service to ALICE before the start of RUN3 (by June 2021) ← **POSTPONED**

Recap CHEP2019

<https://doi.org/10.1051/epjconf/202024504001>

- EOS Erasure Coding implementation => RAIN layout
 - 2 FSTs (data server) on a single FE node to maximize usable space (~70%) out of raw capacity
- Upper cap on total throughput limit by PCI-e 3.0 (~6GB/s)
- Power consumption ~ 1.75W/TB (JBOD+Server+S/W)
 - Enterprise Disk storage: 5 ~ 9W/TB
 - Tape library ~ 0.5W/TB

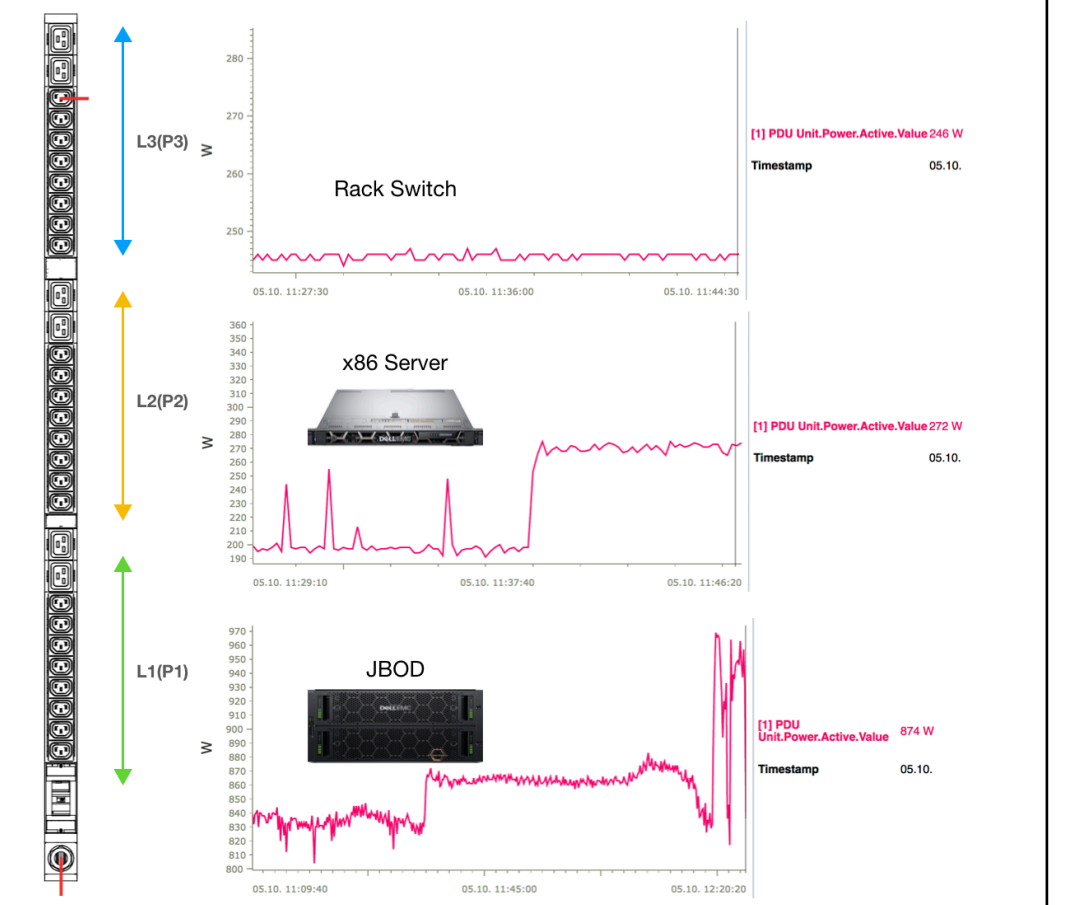
I/O Test: Read/Write

- XFS read/write performance (simultaneous read and/or write from 70 disks)
 - **VDBench** shows full read/write transfer performance @ transfer size >= 2048k (6GB/s)
 - **IOZone** shows full read/write transfer performance @ transfer size ~ 2048k (6GB/s)



Power Consumption

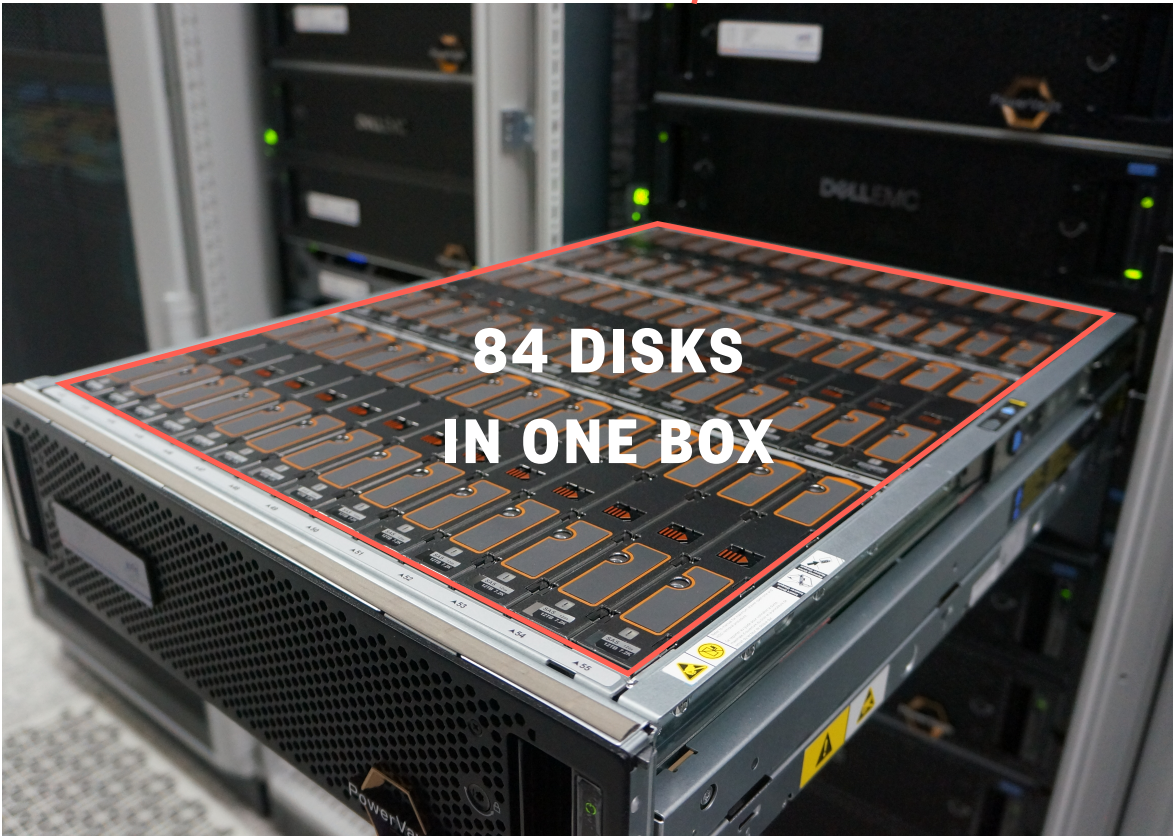
- JBOD Test Equipment (70 Disks)
 - JBOD (DELL ME484): idle = 830W; load = 860W (Max 960) (**1.12W/TB**)
 - Server: idle = 200W; load = 270W
 - Switch: idle = 246W; load = 246W
 - **1.75W/TB** including JBOD, Server and Switch
- Disk Storages (Full Load)
 - Dell EMC SC7020, 2.5PB - 12,120W (**4.8W/TB**)
 - EMC Isilon, 16 Nodes, 2.95 PB - 13,730W (**4.6W/TB**)
 - EMC VNX, 12 Nodes, 2.36 PB - 5,100W (**2.2W/TB**)
 - HITACHI VSP, 2 PB - 18,300W (**9.15W/TB**)
 - EMC Isilon, 15 Nodes, 1.43 PB - 12,880W (**9W/TB**)
 - EMC CX4-960, 1.5PB - 14,900W (**9.9W/TB**)
- Tape Library (Full Load)
 - IBM TS3500 5-Frame (3.2PB) - 1,600W (**0.5W/TB**)



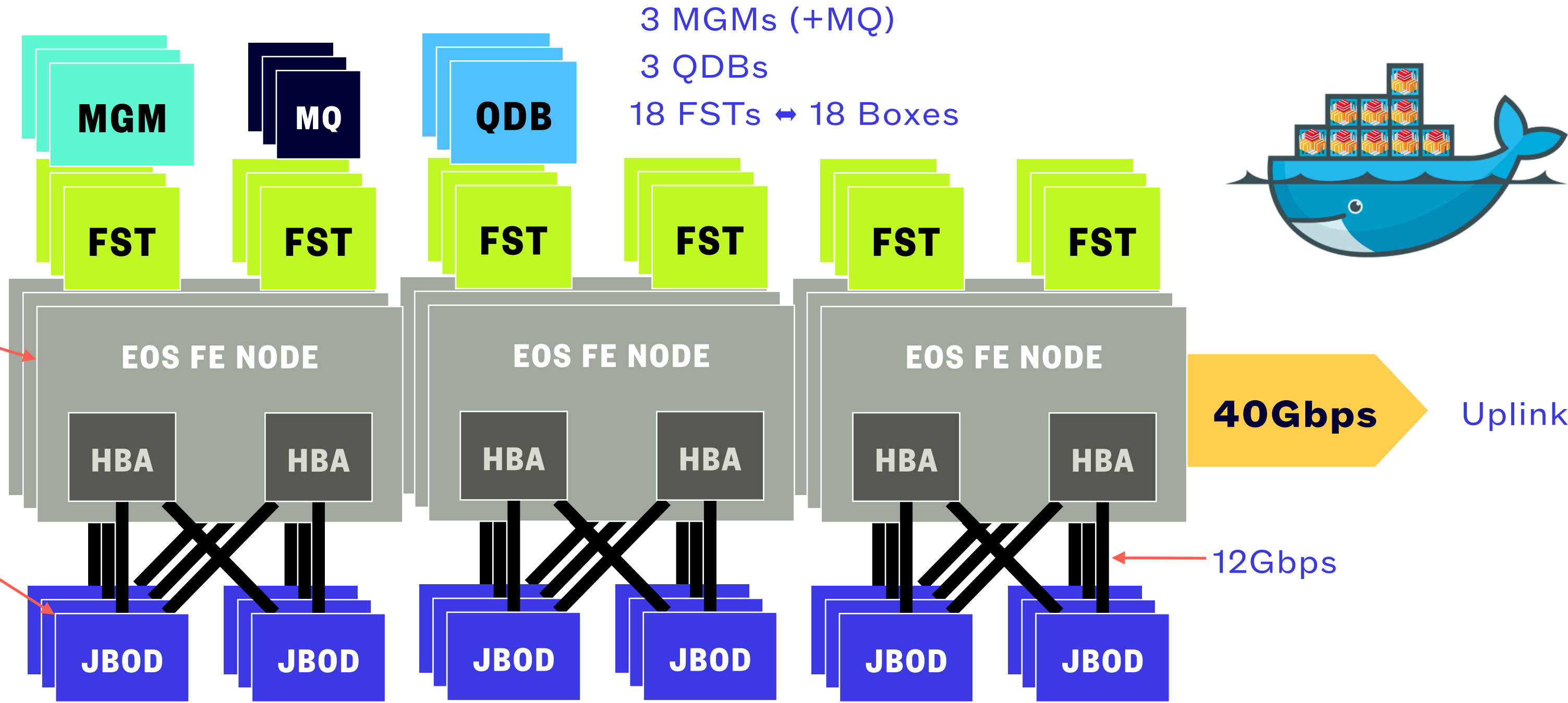
System Architecture



9 servers
18 boxes



84 DISKS
IN ONE BOX

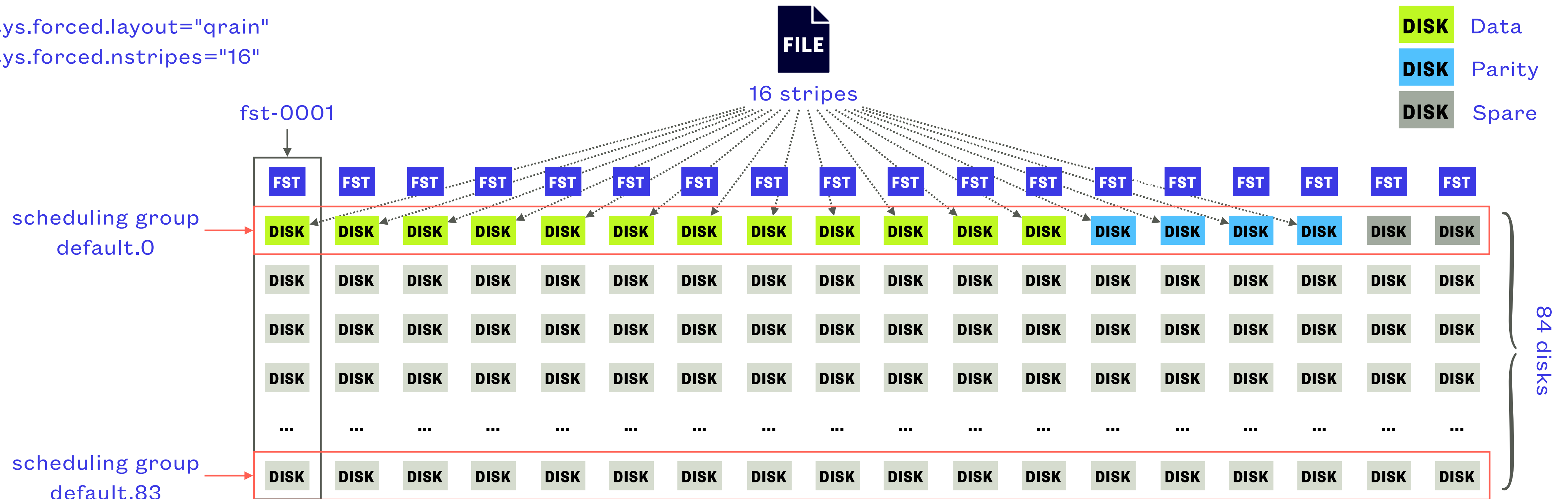


- Total raw capacity = 18,144TB (= 12TB * 84 disks * 18 boxes)
- EOS version = 4.8.25 (released on 2020.11.10)
- EOS components are running on containers (a fork of EOS-Docker project)
 - Ansible playbook available at <https://github.com/jeongheon81/gsd-eos-docker>

QRAIN Layout



sys.forced.layout="grain"
sys.forced.nstripes="16"



- Thanks to spare FSTs,
 - Data are still accessible if 6 FSTs are offline
 - Data can be written if 2 FSTs are offline
 - One node (= 2 FSTs) can be turned off for maintenance at any time
- Data loss rate in a year is $\approx 8.6 \times 10^{-5}\%$, where 5 disks are failed simultaneously, considering 1.17% of AFR in practice
cf. vendor published AFR is 0.35% (AFR = Annualized Failure Rate)

QRAIN Configuration

- 'eos attr' command
 - One can have different layouts on different directories in an EOS instance
 - Available layouts = plain (default, 1 single copy); replica (2 copies); raid6, raiddp (2 parities); archive (3 parities); qrain (4 parities), ...

```
eos attr -r set default=raid6 /eos/gsdctestarea/raid6
eos attr -r set default=archive /eos/gsdctestarea/archive
eos attr -r set default=replica /eos/gsdctestarea/replica
eos attr -r set default=qrain /eos/gsdctestarea/rain12
```

- # of stripes can be changed, e.g. 3 copies, 16 stripes...

```
eos attr -r set default=qrain /eos/gsdctestarea/rain16
eos attr -r set sys.forced.nstripes=16 /eos/gsdctestarea/rain16
```

```
sh-4.2# eos attr ls /eos/gsdctestarea/rain16
sys.eos.btime="1605069261.927407367"
sys.forced.blockchecksum="crc32c"
sys.forced.blocksize="1M"
sys.forced.checksum="adler"
sys.forced.layout="qrain"
sys.forced.nstripes="16"
sys.forced.space="default"
sys.recycle="/eos/gsdctestarea/recycle/"
```

Fileinfo



EOS fileinfo command

```
sh-4.2# eos fileinfo /eos/gsdctestarea/rain16/testfile.10G
```

```
File: '/eos/gsdctestarea/rain16/testfile.10G' Flags: 0640
Size: 10485760000
Modify: Thu Oct 22 00:01:35 2020 Timestamp: 1603324895.724750000
Change: Thu Oct 22 00:00:51 2020 Timestamp: 1603324851.619542497
Birth: Thu Oct 22 00:00:51 2020 Timestamp: 1603324851.619542497
CUid: 0 CGid: 0 Fxid: 0000159b Fid: 5531 Pid: 40 Pxid: 00000028
XStype: adler XS: a1 1c 00 01 ETAGs: "1484716507136:a11c0001"
```

Layout type
of stripes
of replica

```
Layout: grain Stripes: 16 Blocksize: 1M LayoutId: 40640f52 Redundancy: d5::t0
#Rep: 16
```

File chunk location
Scheduling group
Filesystem status

no.	fs-id	host	schedgroup	path	boot	configstatus	drain	active	geotag
0	995	jbod-mgmt-06.sdfarm.kr	default.70	/jbod/box_12_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02
1	1499	jbod-mgmt-09.sdfarm.kr	default.70	/jbod/box_18_disk_070	booted	rw	nodrain	online	kisti::gsdc::g03
2	659	jbod-mgmt-04.sdfarm.kr	default.70	/jbod/box_08_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02
3	407	jbod-mgmt-03.sdfarm.kr	default.70	/jbod/box_05_disk_070	booted	rw	nodrain	online	kisti::gsdc::g01
4	827	jbod-mgmt-05.sdfarm.kr	default.70	/jbod/box_10_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02
5	491	jbod-mgmt-03.sdfarm.kr	default.70	/jbod/box_06_disk_070	booted	rw	nodrain	online	kisti::gsdc::g01
6	1079	jbod-mgmt-07.sdfarm.kr	default.70	/jbod/box_13_disk_070	booted	rw	nodrain	online	kisti::gsdc::g03
7	71	jbod-mgmt-01.sdfarm.kr	default.70	/jbod/box_01_disk_070	booted	rw	nodrain	online	kisti::gsdc::g01
8	743	jbod-mgmt-05.sdfarm.kr	default.70	/jbod/box_09_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02
9	1247	jbod-mgmt-08.sdfarm.kr	default.70	/jbod/box_15_disk_070	booted	rw	nodrain	online	kisti::gsdc::g03
10	155	jbod-mgmt-01.sdfarm.kr	default.70	/jbod/box_02_disk_070	booted	rw	nodrain	online	kisti::gsdc::g01
11	1415	jbod-mgmt-09.sdfarm.kr	default.70	/jbod/box_17_disk_070	booted	rw	nodrain	online	kisti::gsdc::g03
12	911	jbod-mgmt-06.sdfarm.kr	default.70	/jbod/box_11_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02
13	1331	jbod-mgmt-08.sdfarm.kr	default.70	/jbod/box_16_disk_070	booted	rw	nodrain	online	kisti::gsdc::g03
14	239	jbod-mgmt-02.sdfarm.kr	default.70	/jbod/box_03_disk_070	booted	rw	nodrain	online	kisti::gsdc::g01
15	575	jbod-mgmt-04.sdfarm.kr	default.70	/jbod/box_07_disk_070	booted	rw	nodrain	online	kisti::gsdc::g02

```
*****
```

Current Status

- EOS version installed: 4.8.31
 - Automated deployment via Ansible playbook
- Public DNS name pointing to 3 MGMs
- IPv4/IPv6 dual stack configured
- ALICE Integration
 - Enabling Token-based AuthN/AuthZ
 - Enabling ApMon daemons on all EOS FSTs for ALICE MonALISA monitoring
 - Allowing Third-Party Copy by disabling sss enforcement on FSTs

CDS Ready for RUN3

ALICE Integration Done !!

ALICE::KISTI_GSDC::CDS integrated as Custodial Storage Elements (former "Tape Storage Elements") in the ALICE experiment

<http://alimonitor.cern.ch/stats?page=SE/table>

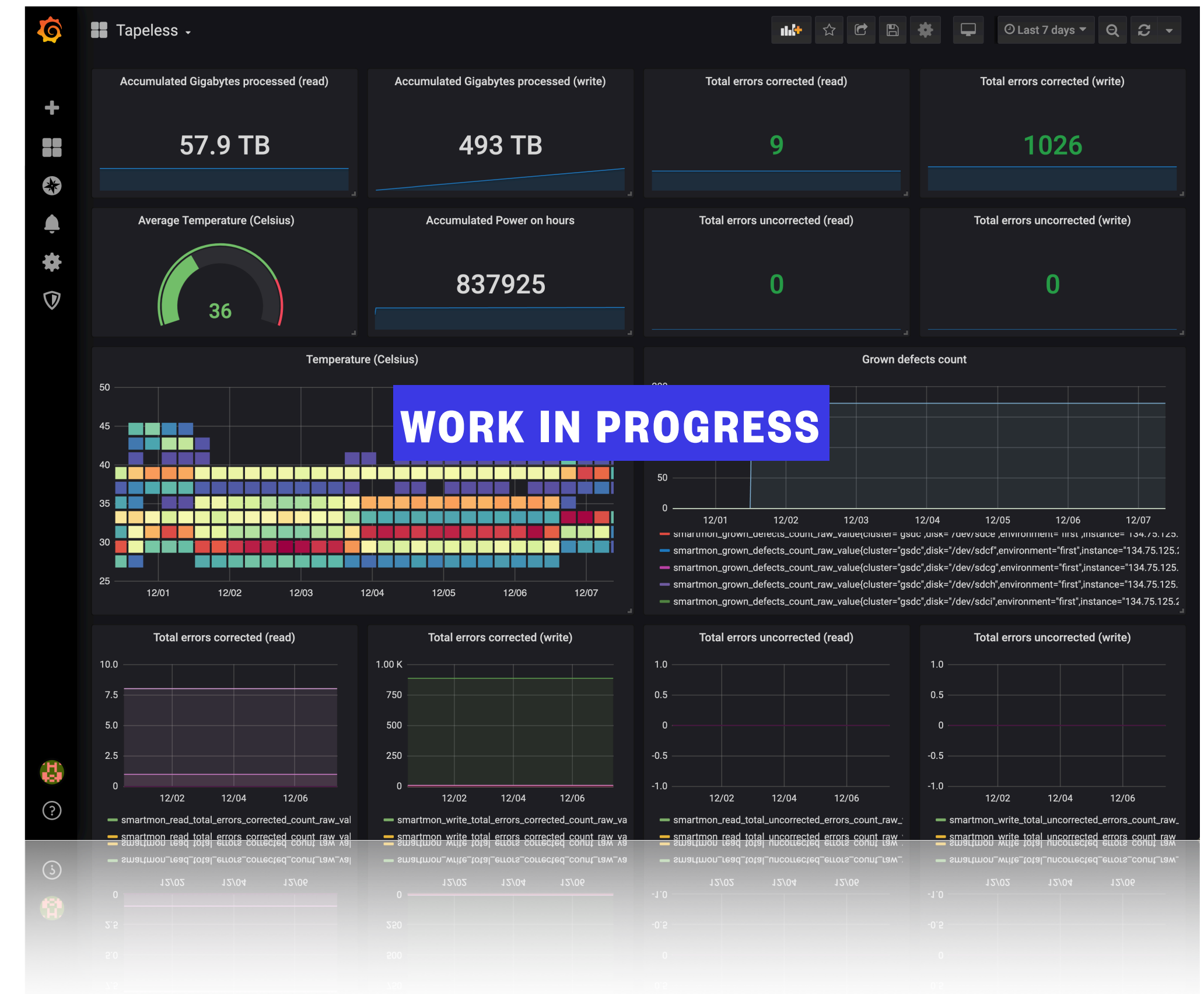
Custodial storage elements																							
AliEn SE		Catalogue statistics							Storage-provided information					Functional tests				Last day add tests		Demotion	IPv6		
SE Name	AliEn name	Tier	Size	Used	Free	Usage	No. of files	Type	Size	Used	Free	Usage	Version	EOS Version	add	get	rm	3rd	Last OK add	Successful	Failed	factor	add
1. CCIN2P3 - TAPE	ALICE::CCIN2P3::TAPE	1	327.4 TB	3.174 PB	-	992.6%	2,242,109	FILE	213.8 TB	208.1 TB	5.783 TB	97.3%	Xrootd v4.12.6			Test...		17.02.2021 10:15	24	0	0		
2. CERN - CTA	ALICE::CERN::CTA	0	4.252 PB	39.17 PB	-	921.3%	38,603,801	CTA	4.188 PB	4.152 PB	36.35 TB	99.15%	Xrootd v4.12.5					17.02.2021 10:28	24	0	1.488%	Test...	
3. CNAF - TAPE	ALICE::CNAF::TAPE	1	1.283 PB	10.29 PB	-	802.3%	6,646,782	FILE	521.6 TB	441.6 TB	79.96 TB	84.67%	Xrootd v4.8.4					17.02.2021 10:13	23	0	0		
4. FZK - TAPE	ALICE::FZK::TAPE	1	601.5 TB	8.528 PB	-	1451%	5,572,494	FILE	601.5 TB	184.6 TB	416.9 TB	30.7%	Xrootd v4.12.2					17.02.2021 10:12	24	0	0	Test...	
5. KISTI_GSDC - CDS	ALICE::KISTI_GSDC::CDS	1	12 PB	0	12 PB	-	0	FILE	15.79 PB	1.233 TB	15.78 PB	0.008%	Xrootd v4.12.5					17.02.2021 10:09	25	0	0		
6. KISTI_GSDC - TAPE	ALICE::KISTI_GSDC::TAPE	1	12.38 PB	3.814 PB	8.564 PB	30.82%	2,878,899	FILE	384.6 TB	339 TB	45.66 TB	88.13%	Xrootd v4.12.4					17.02.2021 10:20	22	0	0		
7. NDGF - DCACHE_TAPE	ALICE::NDGF::DCACHE_TAPE	1	93.13 TB	1.584 PB	-	1741%	1,218,044	SRM	-	-	-	-	dCache 6.2.11				Test...	17.02.2021 10:21	26	0	0		
8. RAL - TAPE	ALICE::RAL::TAPE	1	420 TB	803.2 TB	-	191.2%	543,932	CASTOR	-	-	-	-	Xrootd v4.10.0					17.02.2021 10:16	25	0	0.444%	Test...	
9. RRC_KI_T1 - DCACHE_TAPE	ALICE::RRC_KI_T1::DCACHE_TAPE	1	100 TB	2.605 PB	-	2667%	1,804,048	FILE	-	-	-	-	dCache 5.2.35				Test...	17.02.2021 10:21	25	0	0	Test...	
10. SARA - DCACHE_TAPE	ALICE::SARA::DCACHE_TAPE	1	492 TB	672.7 TB	-	136.7%	393,104	SRM	-	-	-	-	dCache 6.0.29				Test...	17.02.2021 10:22	23	0	0		
Total			31.9 PB	70.61 PB	20.56 PB		59,903,213		21.65 PB	5.299 PB	16.35 PB				10	10	9	7					6

Passed periodic functional tests as well as IPv6 tests

Monitoring



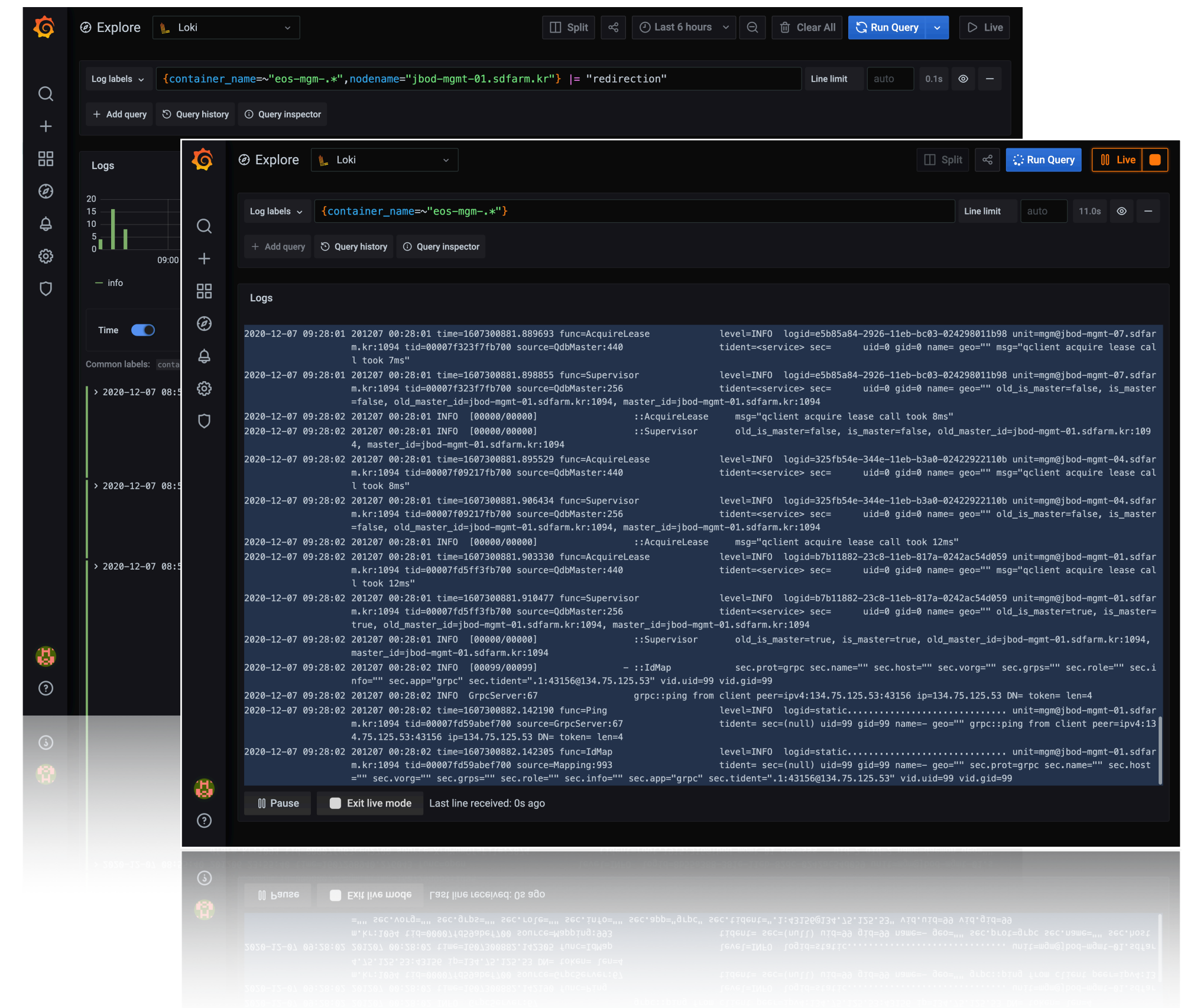
- Prometheus node_exporter + Grafana dashboard
 - Hardware level monitoring using *smartmon* tools
 - Corrected/uncorrected error counts
 - Temperature, defects count, etc.
 - Learning metrics to identifying and predicting disk failures
- Custom script parsing JBOD enclosures status from *sg_ses* command output



Debugging



- EOS services log dump using *loki*, *promtail*
 - Easy access to container level logs (*MGMs*, *FSTs*, ...) by name
 - Live stream (like *'tail -f /var/log/eos/...xrd.log'*) from multiple nodes
 - Complementary to EOS CLI tools
 - Remote support by EOS developers



Issues Fixed

- When read a file with "eos cp" (or eoscp) from RAIN layout with more than 12 stripes
 - Redirection information over 2kB truncated
 - <https://gitlab.cern.ch/dss/eos/-/commit/4cb0f733650e041a3153da60610a8d5a0e4672f4>
 - <https://gitlab.cern.ch/dss/eos/-/commit/03bcc4b55556bc0d7b3160670a41a98bfa50a941>
- Failed to read a file with "eos cp" through MGM secondaries
 - Quota information propagation corrupting namespace
 - <https://gitlab.cern.ch/dss/eos/-/commit/32b93012459f73409b45fcea3c575cc6c47f421>

SPECIAL THANKS TO ELVIN ALIN SINDRILARU !!!

Plan

- Keep monitoring the storage system (and updating EOS if must)
- Developing hardware-level monitoring and power consumption measurement system
- Working on maintenance and operation scheme
 - Disk replacement, JBOD and server maintenance, rolling update of EOS components

Thank you
