
The First Disk-based Custodial Storage for the ALICE experiment

SANG UN AHN, JEONGHEON KIM, HEEJUNE HAN, SEUNG HEE LEE, HEEJUN YOON @ VCHep2021, 18 MAY 2021

Outline

- Introduction
- System Architecture
- QRAIN Layout
- EOS Deployment
- ALICE Integration
- Monitoring
- Conclusion

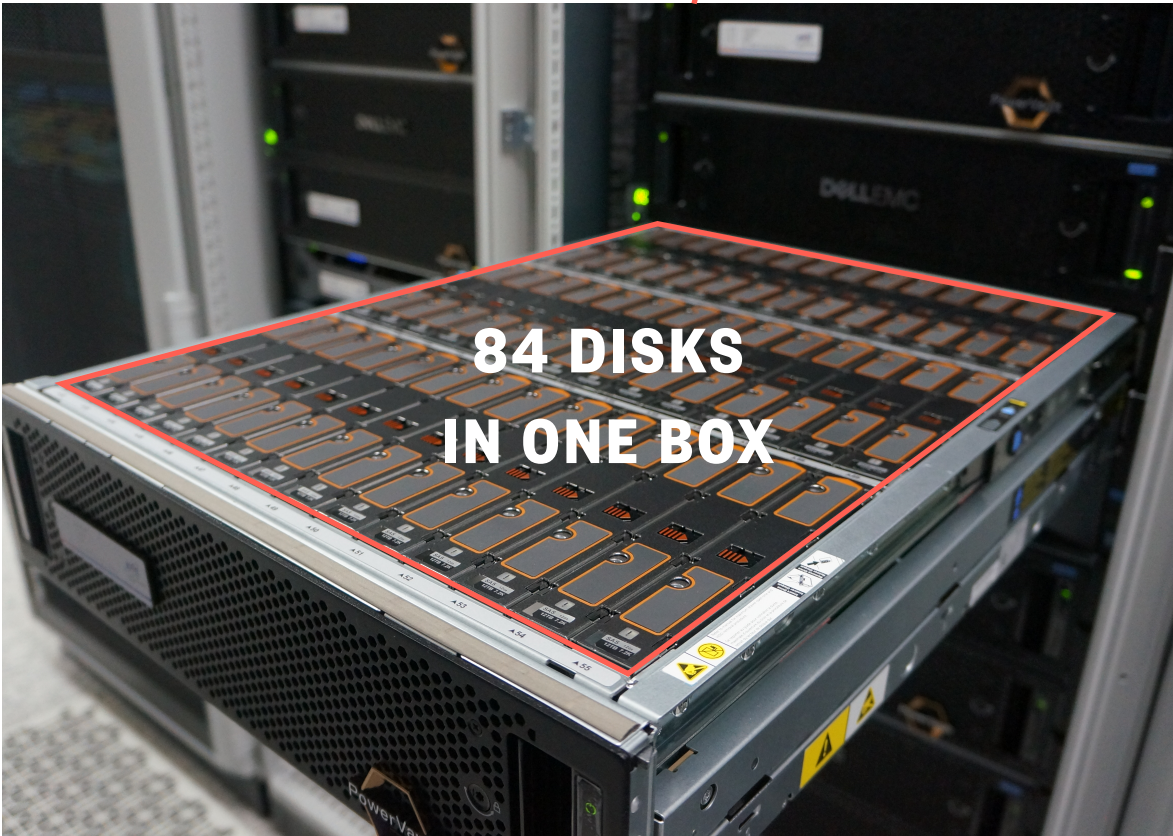
Introduction

- 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**

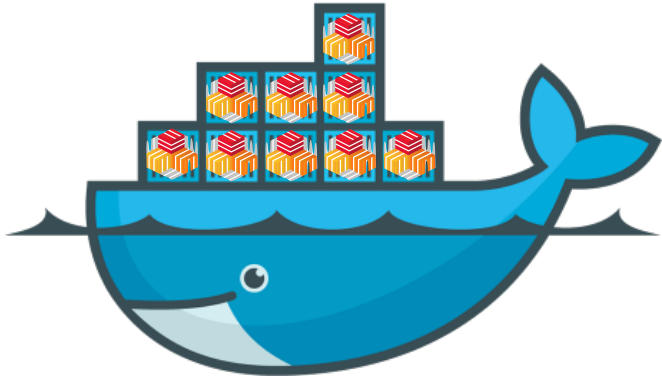
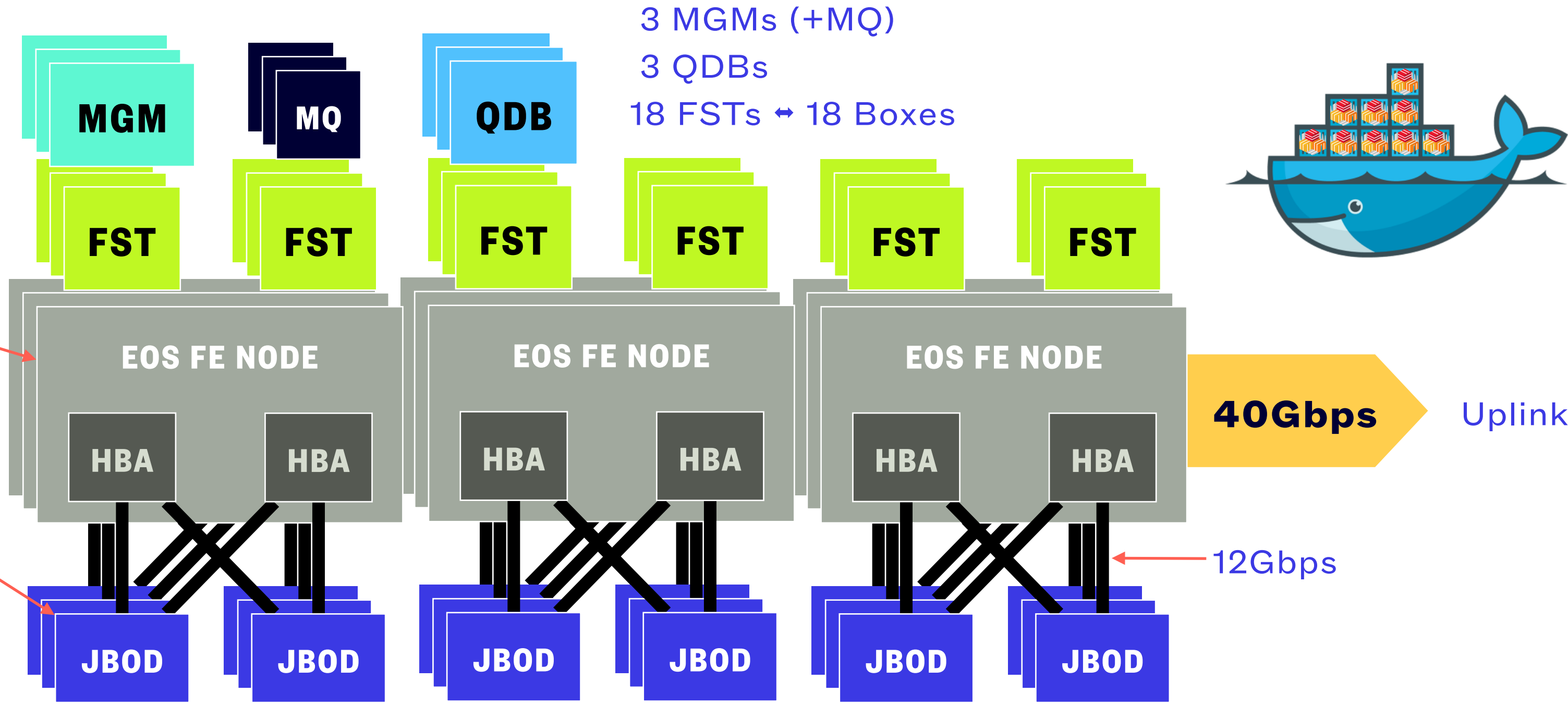
System Architecture



9 servers
18 boxes



84 DISKS
IN ONE BOX

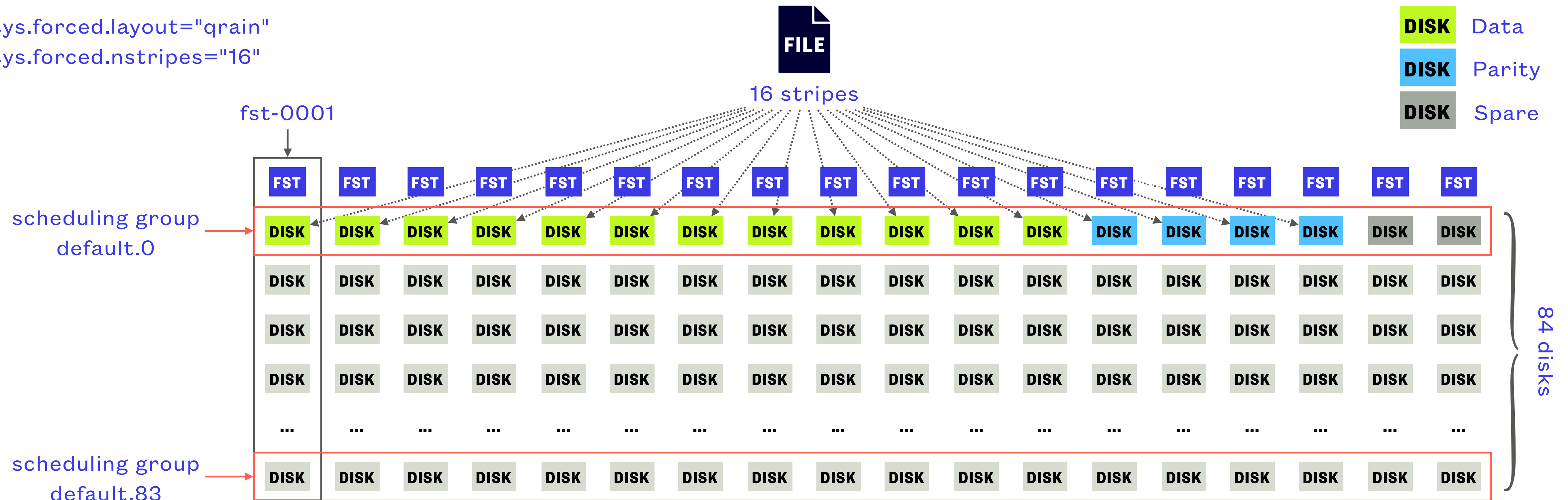


- Total raw capacity = 18,144TB (= 12TB * 84 disks * 18 boxes)
- EOS version = 4.8.31
- EOS components are running on containers (a fork of EOS-Docker project)
 - Ansible playbook available at <https://github.com/jeongheon81/gsd-eos-docker>

QRAIN Layout



sys.forced.layout="qrain"
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)

EOS Deployment

- EOS version installed: 4.8.31
 - Automated deployment via Ansible playbook
 - Issues fixed:
 - Redirection information over 2kB was truncated when read a file with "eos cp" (or eoscp) from RAIN layout with more than 12 stripes
 - Quota information propagation corrupted namespace that fails to read a file with "eos cp" through MGM secondaries
- Public DNS name pointing to 3 MGMs
- IPv4/IPv6 dual stack configured

ALICE Integration

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																						
SE Name	AliEn SE		Catalogue statistics						Storage-provided information					Functional tests				Last day add tests		Demotion	IPv6	
	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
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

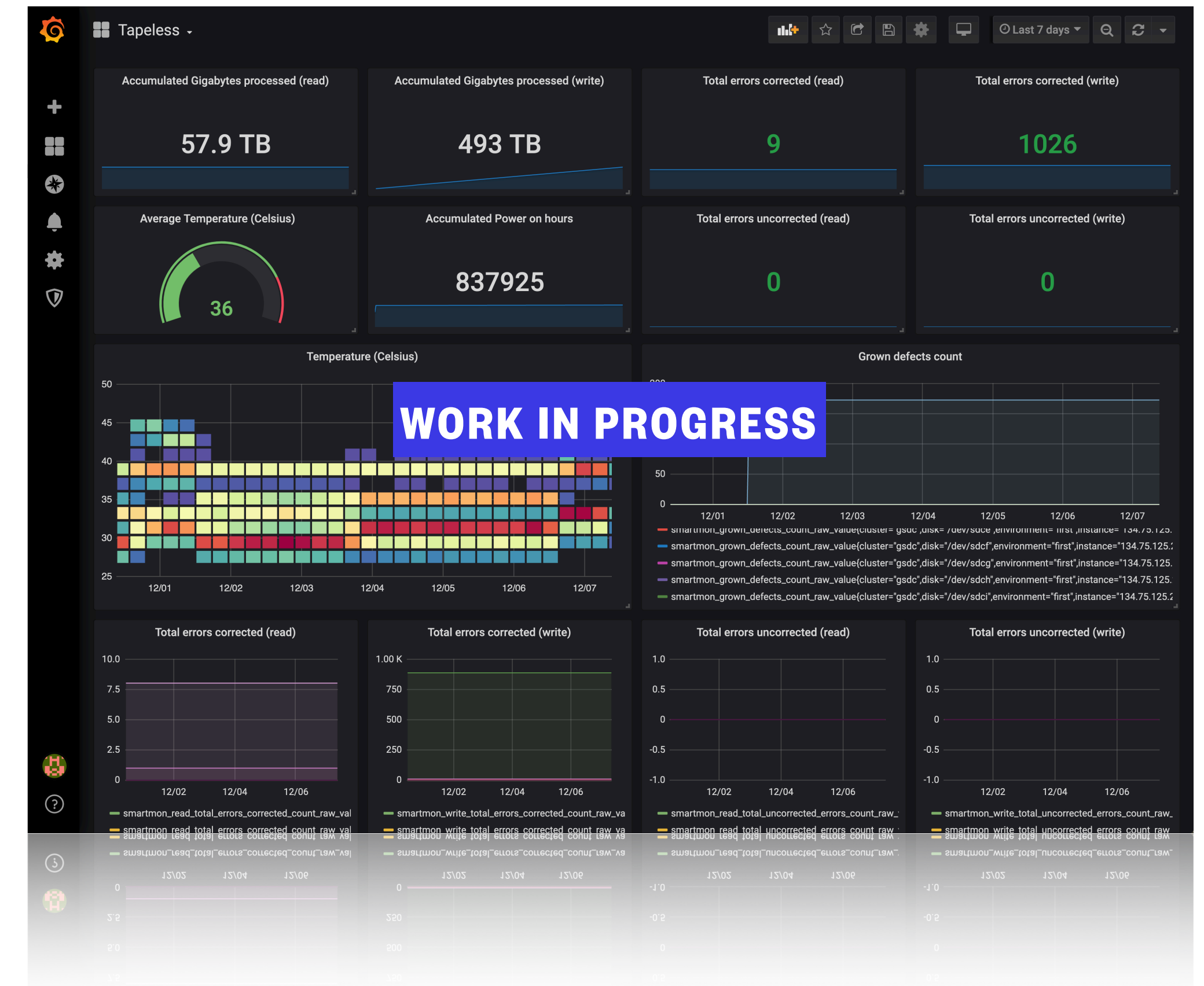
Passed periodic functional tests as well as IPv6 tests

- 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

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
- EOS services log dump using *loki*, *promtail*



Conclusion

- A disk-based custodial storage system as an alternative to tape for preserving data produced from the ALICE experiment at CERN was proposed in order to avoid the potential tape market risks
- The custodial storage system is based on the cheap high-density JBOD enclosures and EOS QRAIN layout to reduce cost and to achieve high enough data protection comparable to tape
- The system is currently integrated and is being commissioned as custodial storage element for the ALICE experiment
- We plan to make the storage in production by July 2021

Thank you
