

# KISTI-GSDC Report

Sang-Un Ahn

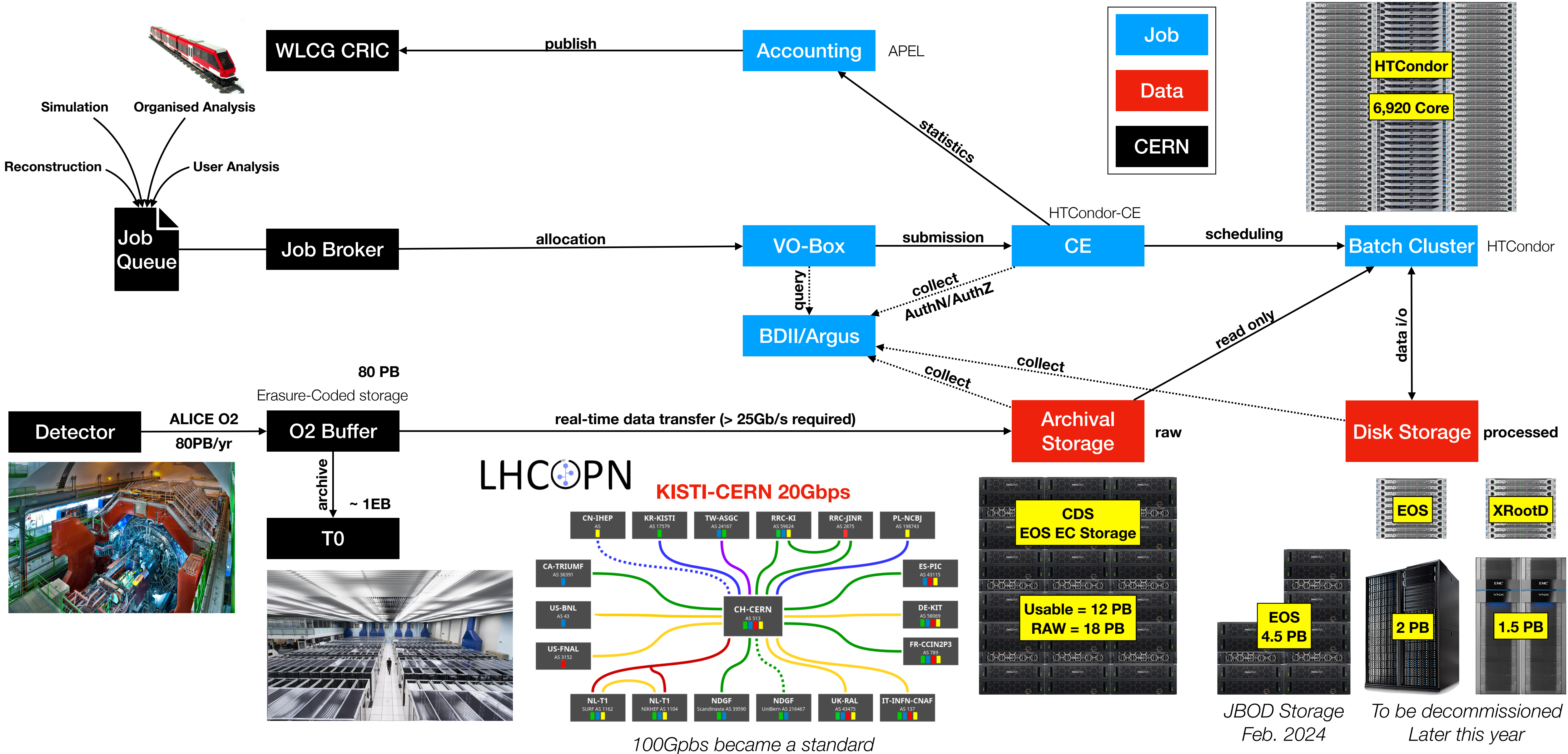
On behalf of KISTI-GSDC



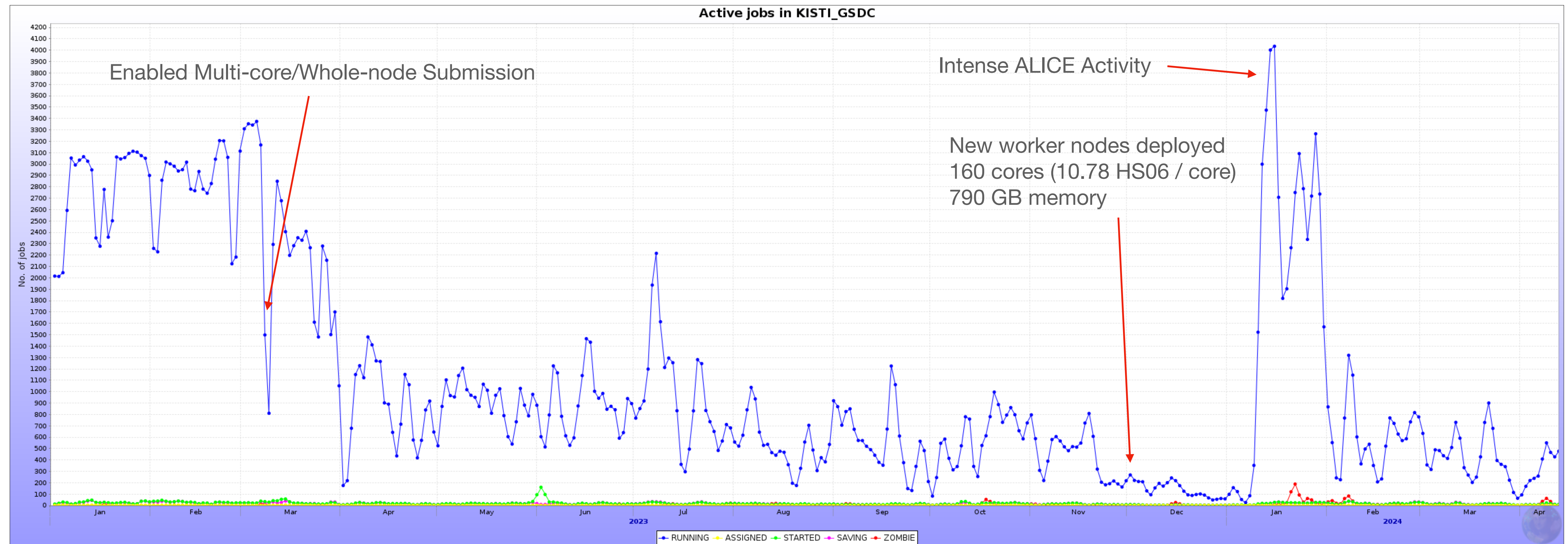
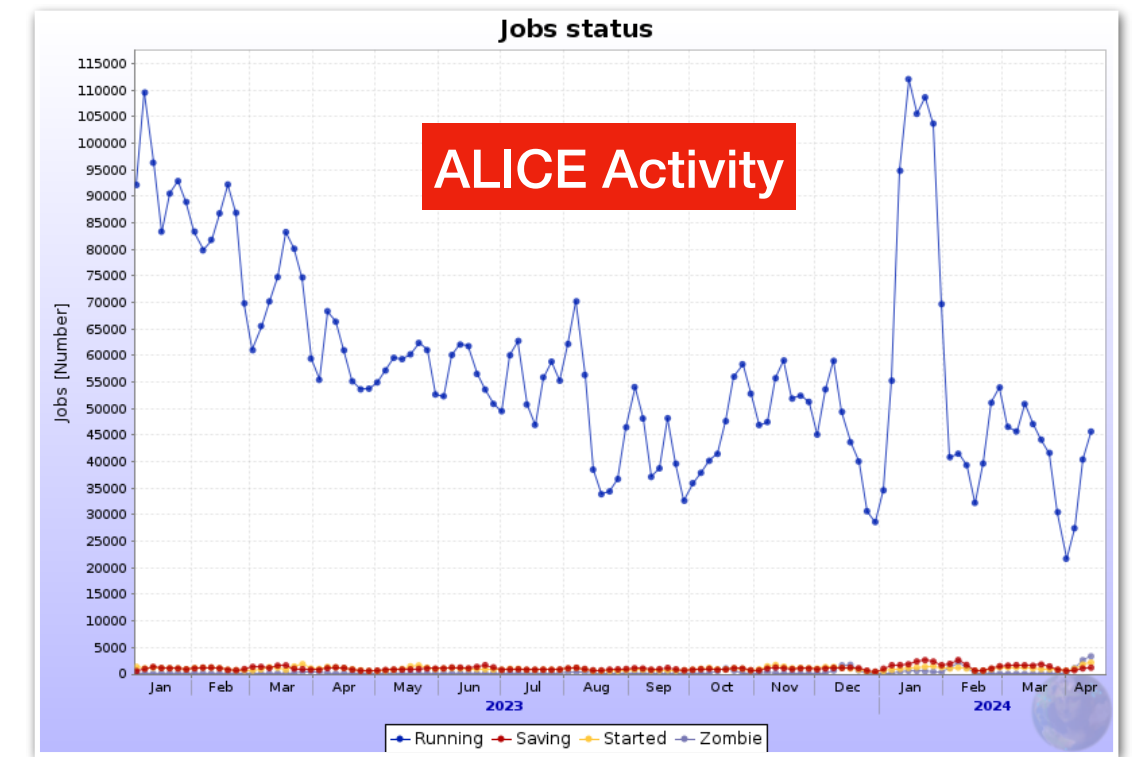
16-18 April 2024 @ ALICE Tier-1/Tier-2 Workshop in Seoul

# Operations

# KISTI ALICE T1 Structure Overview



# Active Jobs

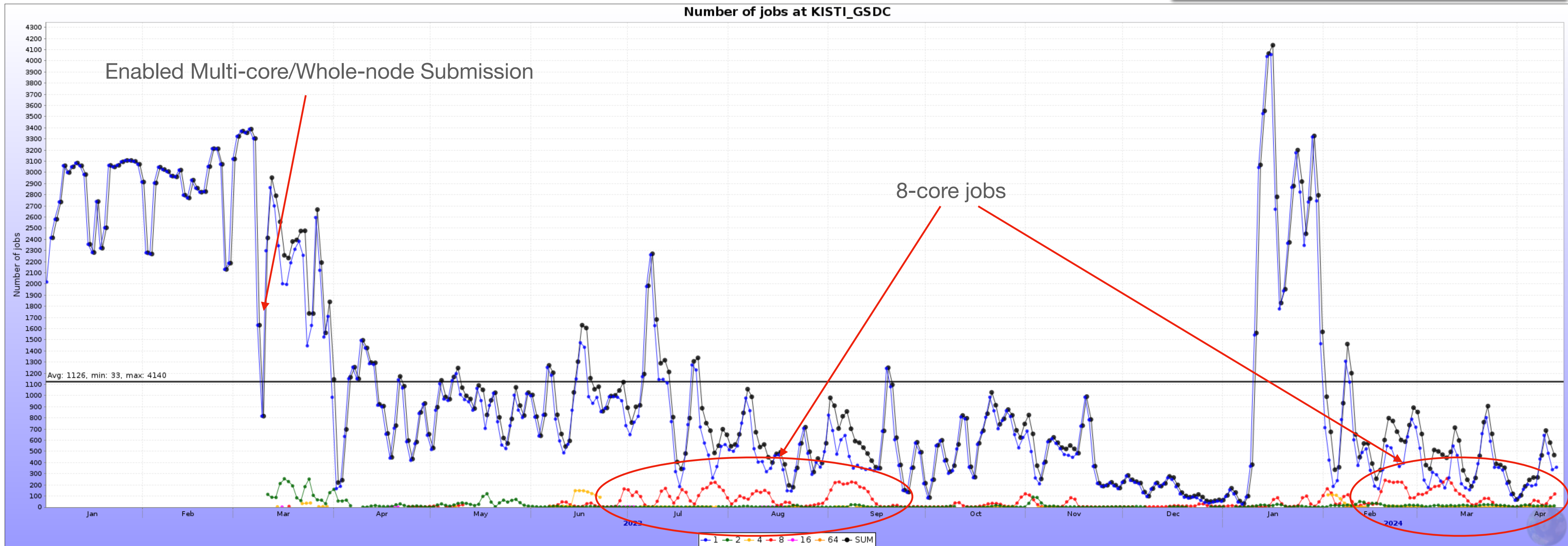


# Multicore Jobs

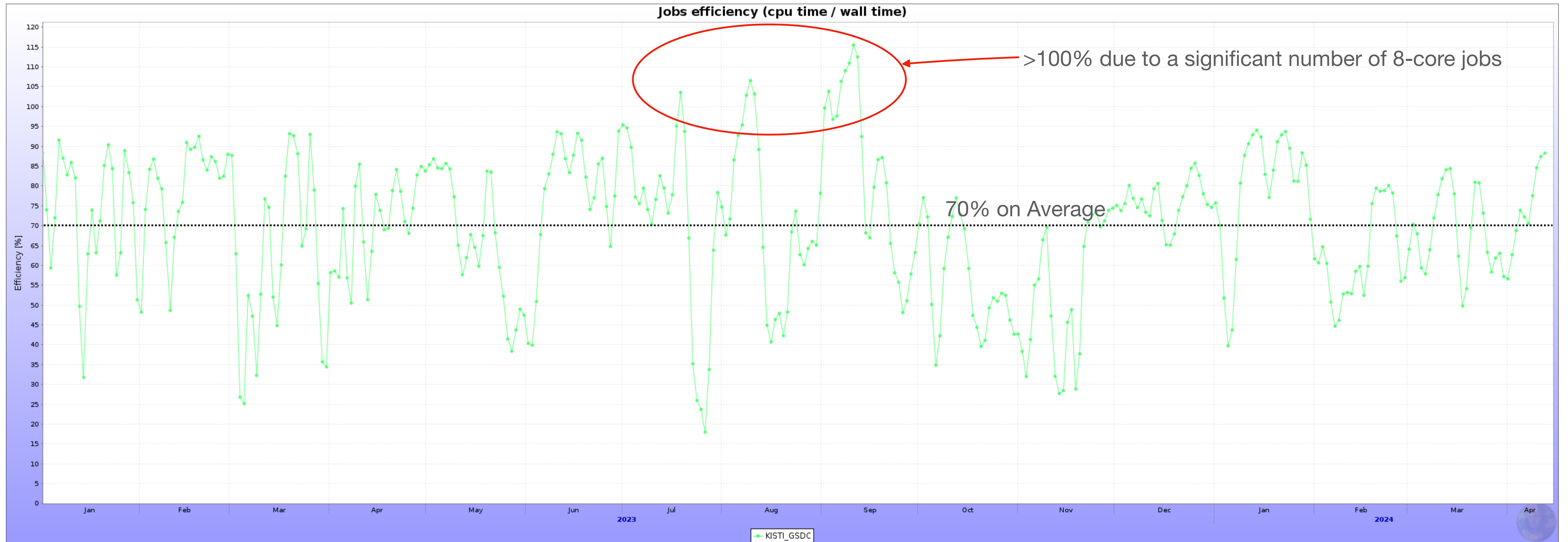
Number of jobs at KISTI\_GSDC

	Series	Last value	Min	Avg	Max
1.	1	358	0	1069	5801
2.	2	12	0	18.39	874
3.	4	1	0	15.16	300
4.	8	119	0	54.35	549
5.	16	1	0	0.619	11
6.	64	2	0	4.2	20

Number of jobs at KISTI\_GSDC



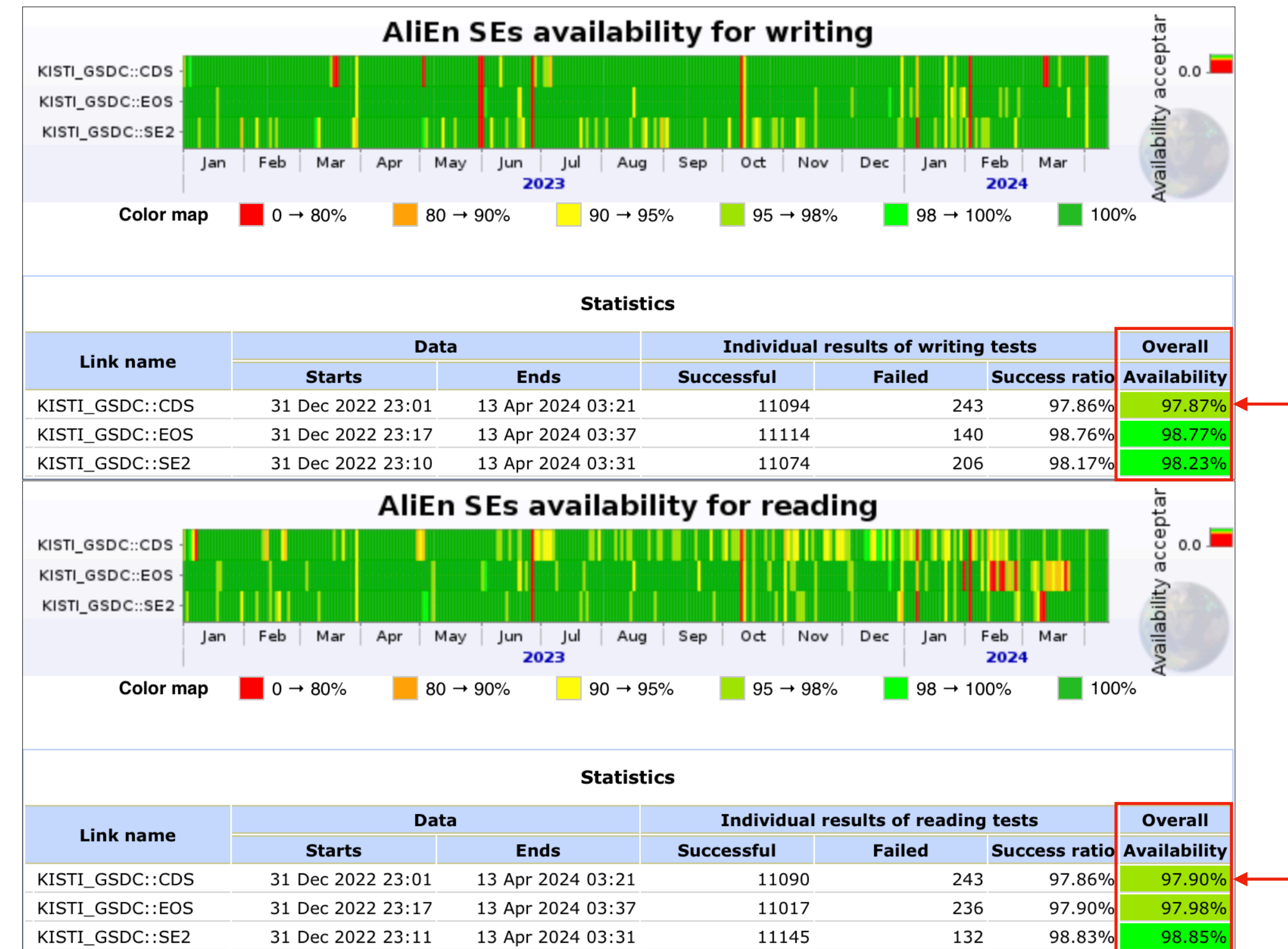
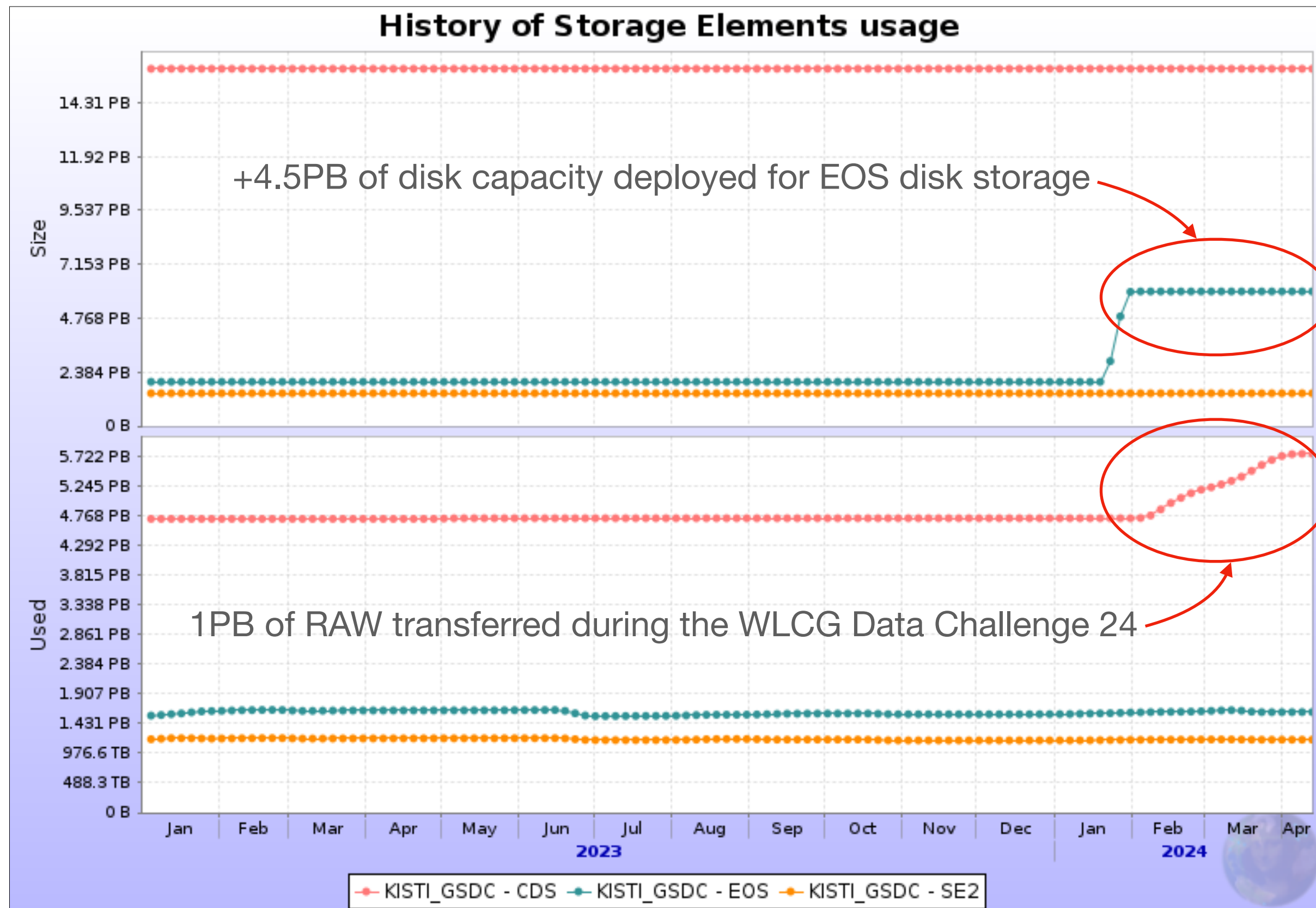
# Jobs Efficiency



# SE Status

	Used (PiB)	Total (PiB)	Usage (%)
CDS	5.766	15.78	36.54
EOS	1.62	5.948	27.24
SE2	1.172	1.446	81.05

> 97% of overall SE availability for writing/reading



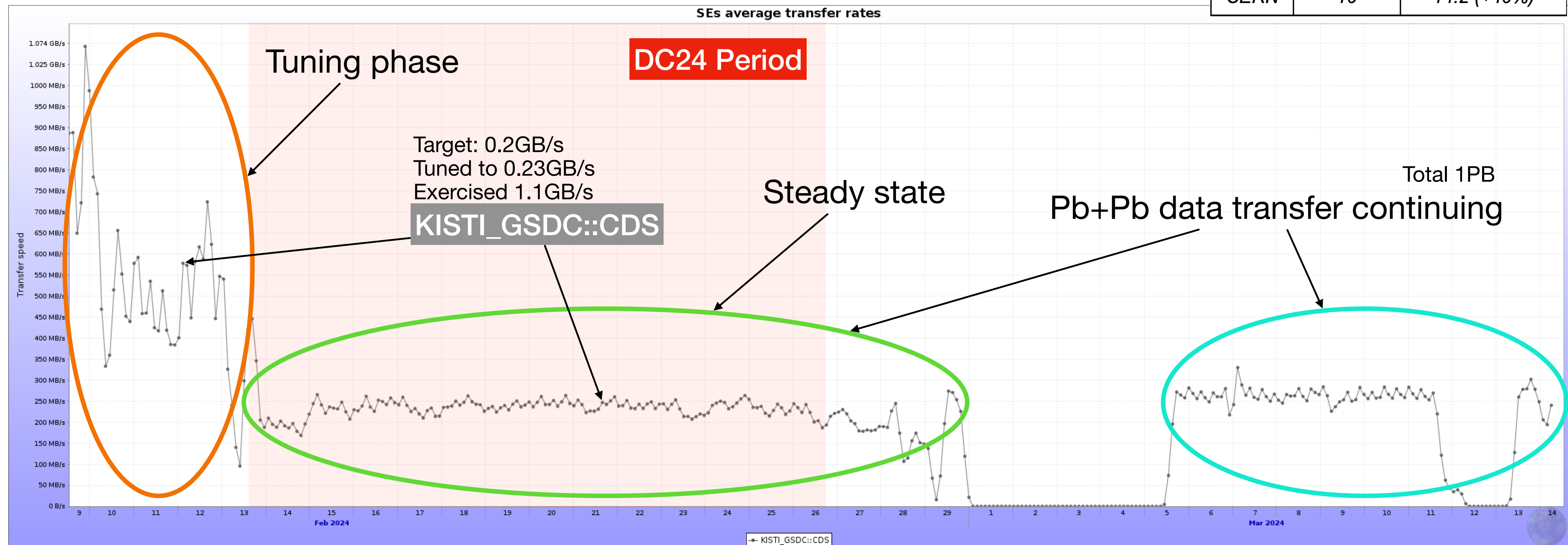
CDS : Custodial Disk Storage  
 EOS : EOS Disk Storage  
 SE2 : XRootD Disk Storage

# WLCG Data Challenge 24

## CDS Participation as a Tape

- Transfer of real Pb+Pb data collected in 2023, 34PB in total
- 1PB of data being transferred after the challenge, ETA end of March

Centre	Target rate GB/s	Average achieved GB/s
CNAF	0.8	0.98 (+20%)
IN2P3	0.4	0.6 (+40%)
KISTI	0.2	0.25 (+22%)
GridKA	0.6	1.12 (+90%)
NDGF	0.3	0.35 (+15%)
NL-T1	0.1	0.25 (+150%)
RAL	0.1	0.58 (+500%)
CERN	10	14.2 (+40%)



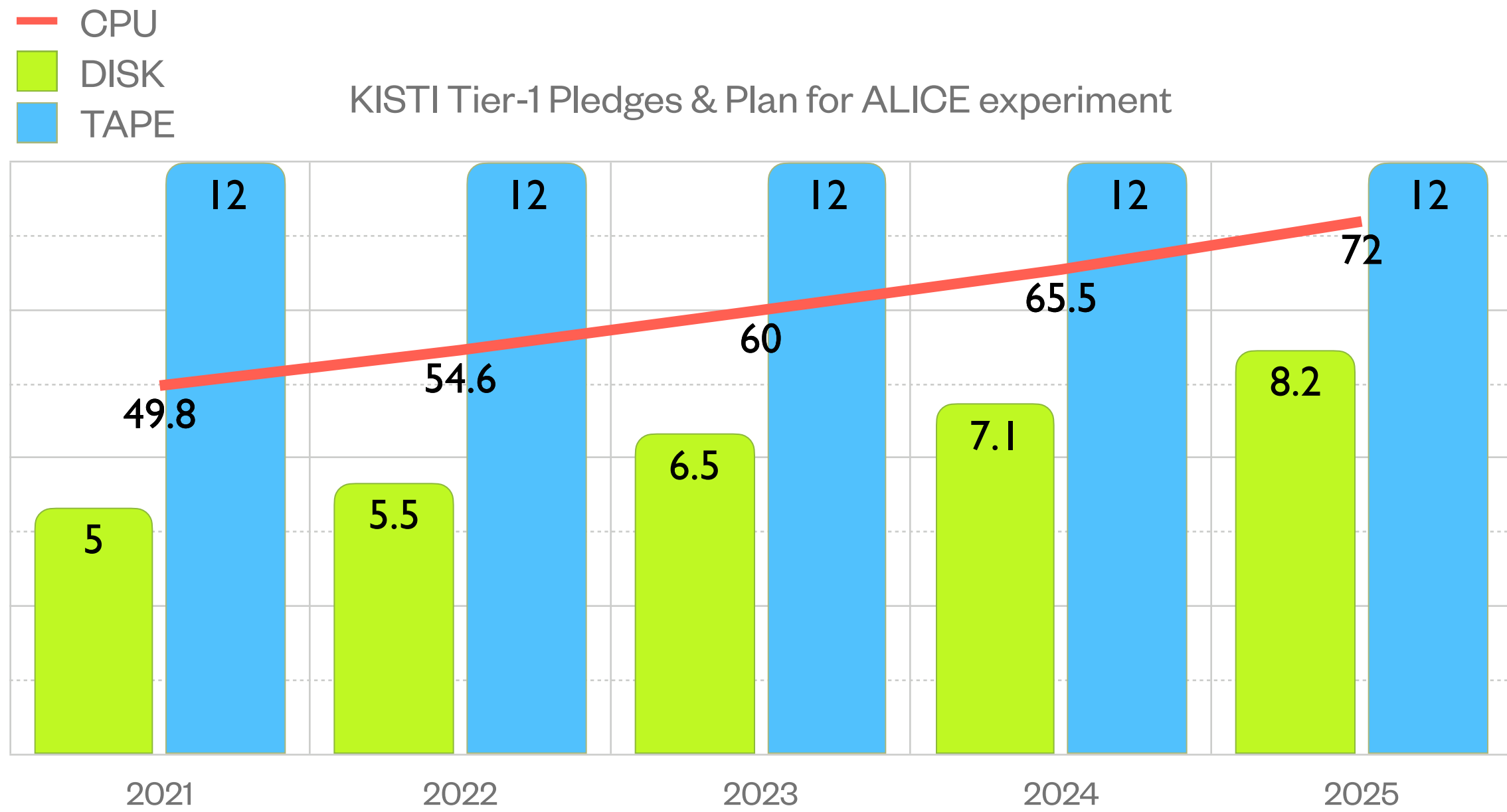


# Site Availability

	2023												2024	
	Jan	Feb	Mar	Apr	May	Jun	Jul	Aug	Sep	Oct	Nov	Dec	Jan	Feb
<b>Availability</b>	98	97	99	100	98	94	98	98	99	98	95	97	97	97
<b>Overall</b>	97.5												97	

- Misconfiguration on LHCOPN backup link has affected significantly the site availability, persisting from the second half of last year until recently
  - Close collaboration with the KREONet service team to address the issue
- An unknown timeout issue during the EOS deployment on CentOS 9 Stream degraded availability in February and March of this year
  - Migrating to AlmaLinux 9.3 resolved the issue, attributed to the newer kernel on CS9 and complementary packages accompanying infrastructure provisioning for security, accounting, auditing and monitoring
- Availability re-computation requests vis GGUS tickets were approved and applied to WLCG monthly reports

# Pledges



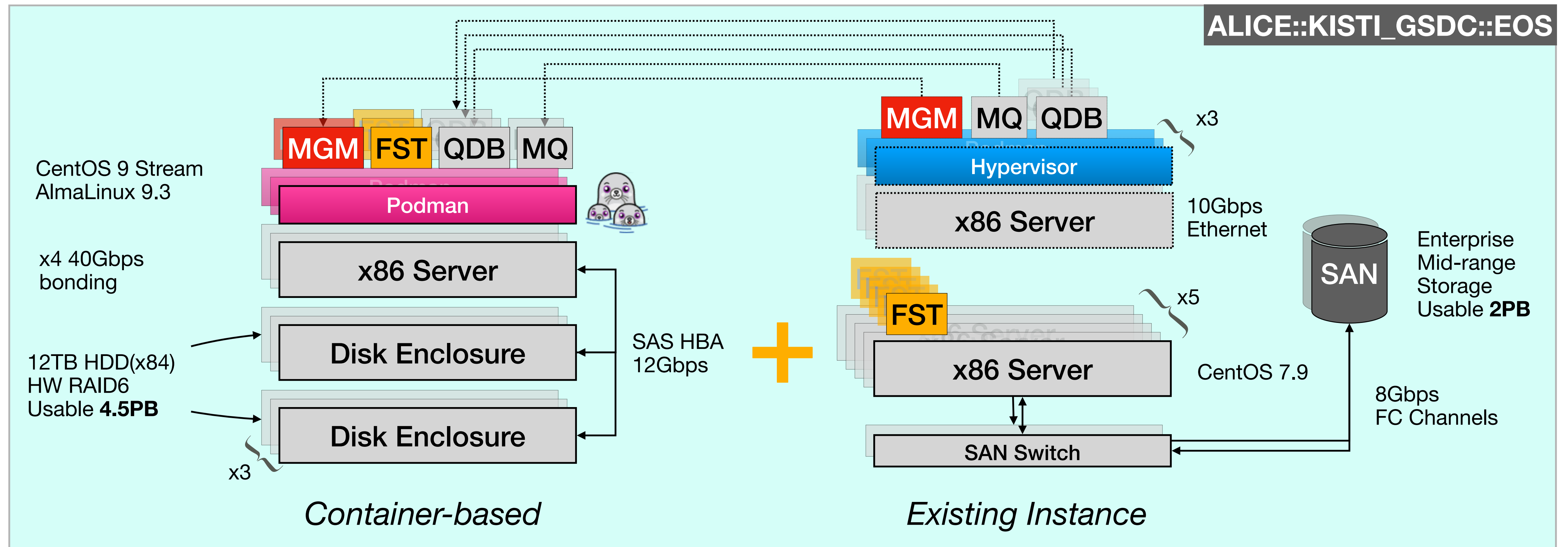
Installed (Plan)	2021	2022	2023	2024	(2025)
<b>CPU (cores)</b>	3,880	3,880	3,880	6,920	8,040
<b>DISK (TB)</b>	4,500	4,500	6,500	7,100	8,200
<b>TAPE (TB)</b>	12,000	12,000	12,000	12,000	12,000

- Current deployment
  - CPU : 72 kHS23 (6.9k Cores)
  - Disk : 8 PB (7.4 PiB)
  - CDS (TAPE) : 12 PB
- Plan
  - CPU : 84 kHS23 (+12 kHS23 later this year)
  - Disk : 7.5 PB
    - SE2 **-1.5 PB** (decommissioning this year)  
8 PB → 6.5 PB
    - EOS **+3.0 PB** to replace old HW **-2.0 PB**  
6.5 PB → 7.5 PB
  - CDS : 15 PB (tentatively)
    - Expanding in 2025 with 2024 procurement

EOS

# EOS Deployments (1/2)

## Disk Storage

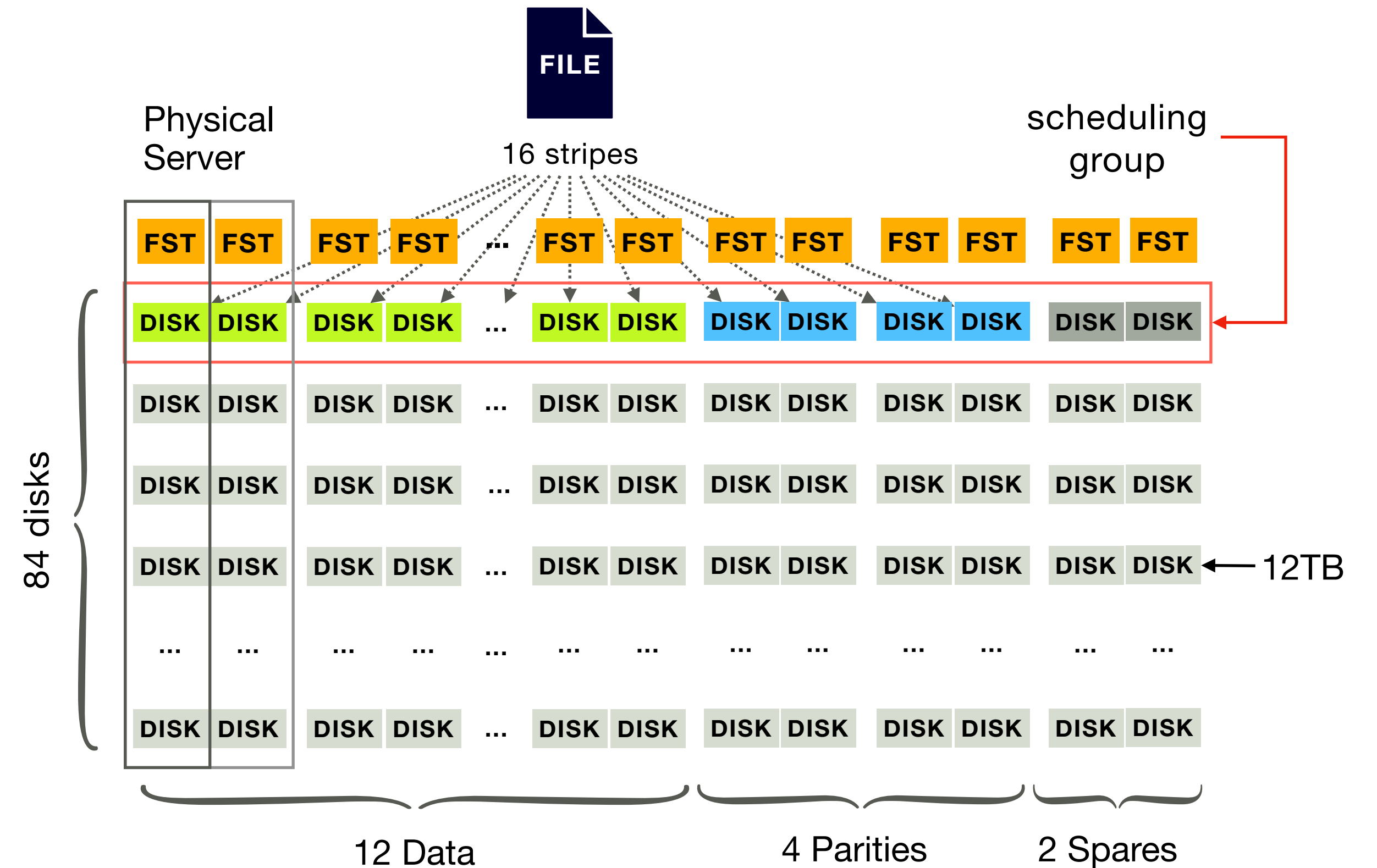
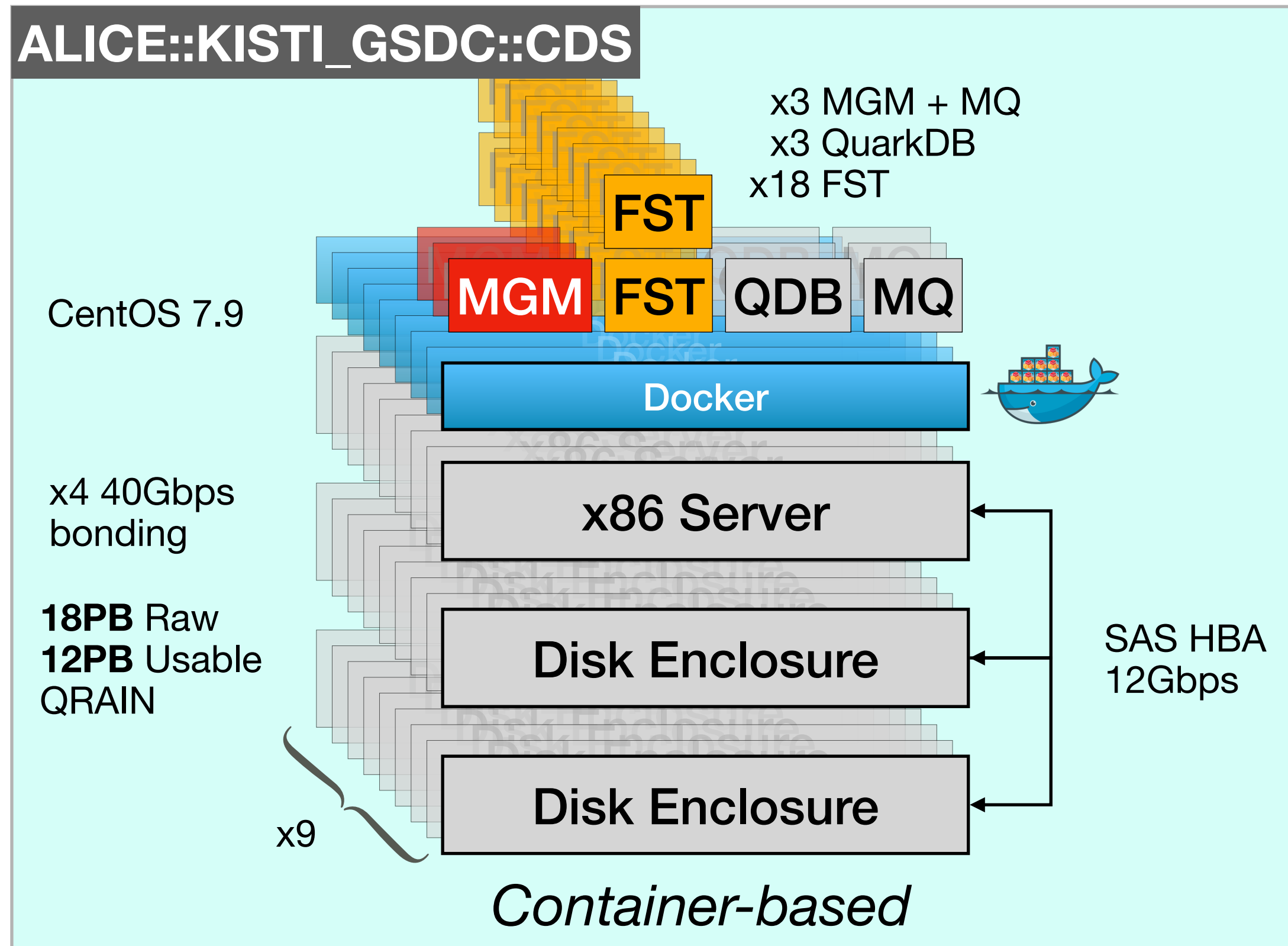


- Transparent transition of MGM and QuarkDB clusters from VMs to Containers
- EOS upgrade from 5.1.22 to 5.2.16 for existing setup, FMD migration from LevelDB completed beforehand
- Expanded to 6.5PB

# EOS Deployments (2/2)

## Custodial Storage

```
[root@jbod-mgmt-07 MGM_MASTER=true /]# eos attr ls /eos/gsdcd/grid
sys.eos.btime="1612374338.811408574"
sys.forced.blockchecksum="crc32c"
sys.forced.blocksize="1M"
sys.forced.checksum="adler"
sys.forced.layout="grain"
sys.forced.nstripes="16"
sys.forced.space="default"
```



- Disk-based Raw Archive storage for ALICE in production since 2021 deployed using Docker Container
- Comparable level of data protection provided by QRAIN Layout (12 stripes + 4 parities + 2 spares)
- Successful upgrade to v5.1.22 from v4.8.82 (May 2023)

# EOS @ KISTI for ALICE

## ALICE::KISTI\_GSDC::EOS

### Disk storage elements

KISTI_GSDC - EOS																							
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. KISTI_GSDC - EOS	ALICE::KISTI_GSDC::EOS	1	5.948 PB	1.639 PB	4.309 PB	27.55%	50,149,564	FILE	5.948 PB	1.74 PB	4.208 PB	29.25%	Xrootd 5.6.7	5.2.16					14.03.2024 14:43	25	0	0	
<b>Total</b>			<b>5.948 PB</b>	<b>1.639 PB</b>	<b>4.309 PB</b>		<b>50,149,564</b>		<b>5.948 PB</b>	<b>1.74 PB</b>	<b>4.208 PB</b>				<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>					<b>1</b>

## ALICE::KISTI\_GSDC::CDS

### Custodial storage elements

CDS																							
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. KISTI_GSDC - CDS	ALICE::KISTI_GSDC::CDS	1	15.79 PB	5.378 PB	10.41 PB	34.06%	10,959,791	FILE	15.76 PB	7.909 PB	7.856 PB	50.17%							14.03.2024 14:27	24	0	4.706%	
<b>Total</b>			<b>15.79 PB</b>	<b>5.378 PB</b>	<b>10.41 PB</b>		<b>10,959,791</b>		<b>15.76 PB</b>	<b>7.909 PB</b>	<b>7.856 PB</b>				<b>1</b>	<b>1</b>	<b>1</b>	<b>1</b>					<b>1</b>

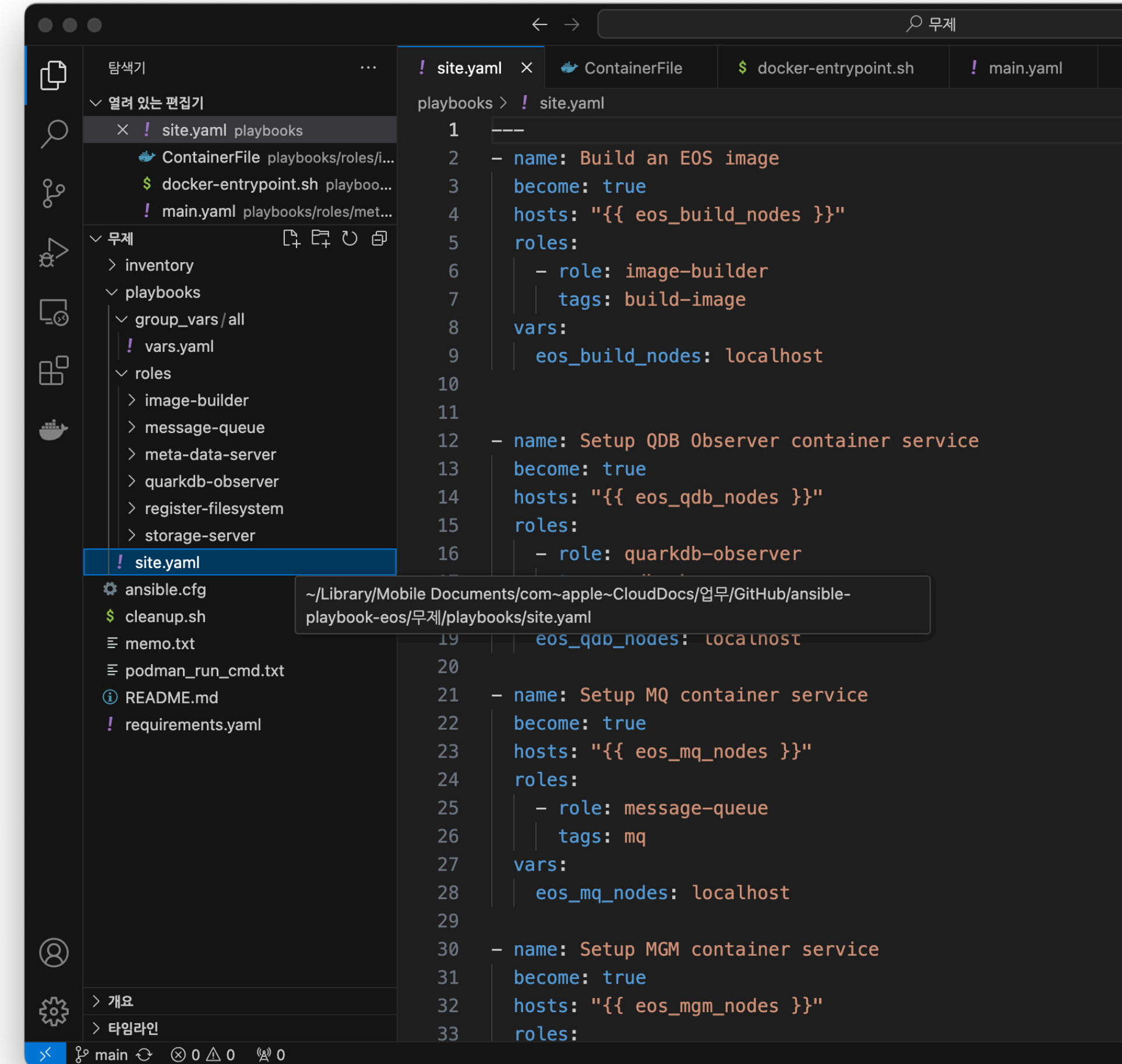
- IPv4/IPv6 Dual Stack
- ALICE-Specific Token Authentication/Authorization, HTTP(S), Third-Party Copy enabled
- MLSensor (successor of EOS Apmon) deployed for monitoring (not yet for CDS)



# EOS v5 Container on EL9: Practices (2/3)

## Automation via Ansible Playbook

- Playbook structure:
  - *site.yaml* - tags defined to perform specific operation (role) in an automated way
  - *group\_vars*
    - *vars.yaml* - key-value variables for group parameters such as *eos\_instance\_name*, *eos\_geotag*, ports, master/slave MGM FQDNs, QDB cluster and FST data directories
  - roles
    - *image-builder* | *message-queue* (MQ) | *meta-data-server* (MGM) | *quarkdb-observer* (QDB) | *register-filesystem* | *storage-server* (FST)
    - *handlers* defined to invoke *firewalld* policy implementation and *systemd* integration
    - Creating essential configuration files by templating *xrd.cf*.{*qdb*|*mq*|*mgm*|*fst*}, *eos\_env*, *scitokens*, ALICE-specific (*TkAuthz.Authorization* & *mlsensor*), etc.



The screenshot shows a code editor with a file explorer on the left and a code editor on the right. The file explorer shows a directory structure for an Ansible project, including `site.yaml`, `ContainerFile`, `docker-entrypoint.sh`, and `main.yaml`. The code editor displays the content of `site.yaml`, which defines several Ansible plays for building and setting up EOS containers.

```
1 ---
2 - name: Build an EOS image
3   become: true
4   hosts: "{{ eos_build_nodes }}"
5   roles:
6     - role: image-builder
7       tags: build-image
8   vars:
9     eos_build_nodes: localhost
10
11
12 - name: Setup QDB Observer container service
13   become: true
14   hosts: "{{ eos_qdb_nodes }}"
15   roles:
16     - role: quarkdb-observer
17
18   eos_qdb_nodes: localhost
19
20
21 - name: Setup MQ container service
22   become: true
23   hosts: "{{ eos_mq_nodes }}"
24   roles:
25     - role: message-queue
26       tags: mq
27   vars:
28     eos_mq_nodes: localhost
29
30 - name: Setup MGM container service
31   become: true
32   hosts: "{{ eos_mgm_nodes }}"
33   roles:
```



# EOS v5 Container on EL9: Practices (3/3)

## Systemd Integration

- Systemd service file for each of EOS components manipulating podman commands in such a way that it invokes *podman {run|rm|stop|...}*

- E.g. */etc/systemd/system/{qdb|mq|fst|mgm}-container.service*

<...>

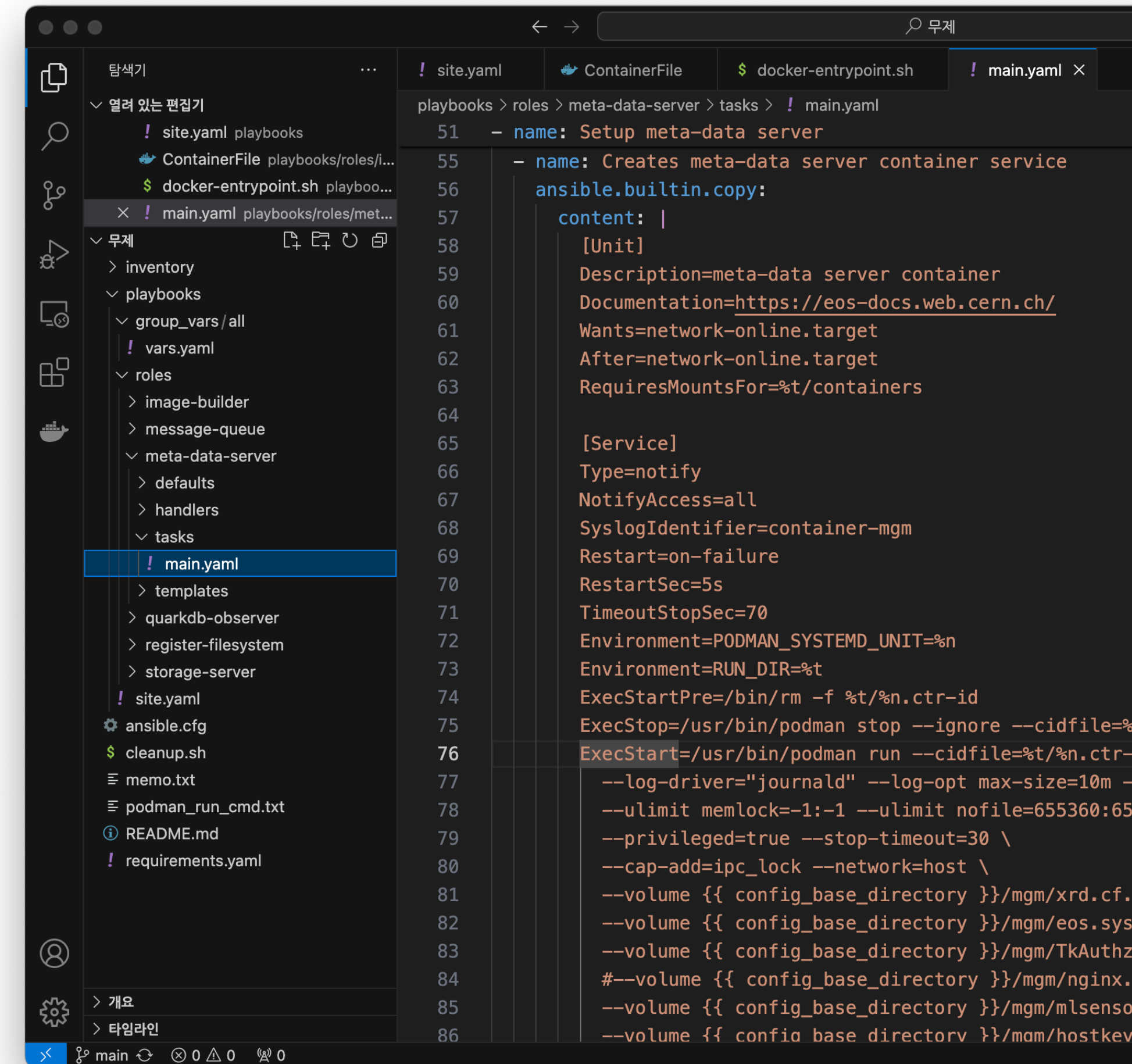
*ExecStart=/usr/bin/podman run < parameters >*

*ExecStop=/usr/bin/podman stop*

*ExecStopPost=/usr/bin/podman rm*

<...>

- *systemdctl {start|stop|restart} {qdb|mq|fst|mgm}-container.service*
  - *syslog (journalctl)* traces container logs (= podman logs)
- Service update as well as roll-back can be quick and easy
  - Update images (pulling from registries or uploading from local one)
  - *systemctl restart \*-container.service*



The screenshot shows a code editor with a file explorer on the left and a code editor on the right. The file explorer shows a directory structure for Ansible playbooks, including `site.yaml`, `ContainerFile`, `docker-entrypoint.sh`, and `main.yaml`. The code editor displays the content of `main.yaml`, which defines a task to create a meta-data server container service. The task is named `Setup meta-data server` and uses the `ansible.builtin.copy` module to create a systemd service file. The service file content is as follows:

```
playbooks > roles > meta-data-server > tasks > ! main.yaml
51 - name: Setup meta-data server
55 - name: Creates meta-data server container service
56   ansible.builtin.copy:
57     content: |
58       [Unit]
59       Description=meta-data server container
60       Documentation=https://eos-docs.web.cern.ch/
61       Wants=network-online.target
62       After=network-online.target
63       RequiresMountsFor=%t/containers
64
65       [Service]
66       Type=notify
67       NotifyAccess=all
68       SyslogIdentifier=container-mgm
69       Restart=on-failure
70       RestartSec=5s
71       TimeoutStopSec=70
72       Environment=PODMAN_SYSTEMD_UNIT=%n
73       Environment=RUN_DIR=%t
74       ExecStartPre=/bin/rm -f %t/%n.ctr-id
75       ExecStop=/usr/bin/podman stop --ignore --cidfile=%
76       ExecStart=/usr/bin/podman run --cidfile=%t/%n.ctr-
77         --log-driver="journald" --log-opt max-size=10m --
78         --ulimit memlock=-1:-1 --ulimit nofile=655360:65
79         --privileged=true --stop-timeout=30 \
80         --cap-add=ipc_lock --network=host \
81         --volume {{ config_base_directory }}/mgm/xrd.cf.
82         --volume {{ config_base_directory }}/mgm/eos.sys
83         --volume {{ config_base_directory }}/mgm/TKAuthz
84         #--volume {{ config_base_directory }}/mgm/nginx.
85         --volume {{ config_base_directory }}/mgm/mlsenso
86         --volume {{ confia base directorv }}/mam/hostkev
```

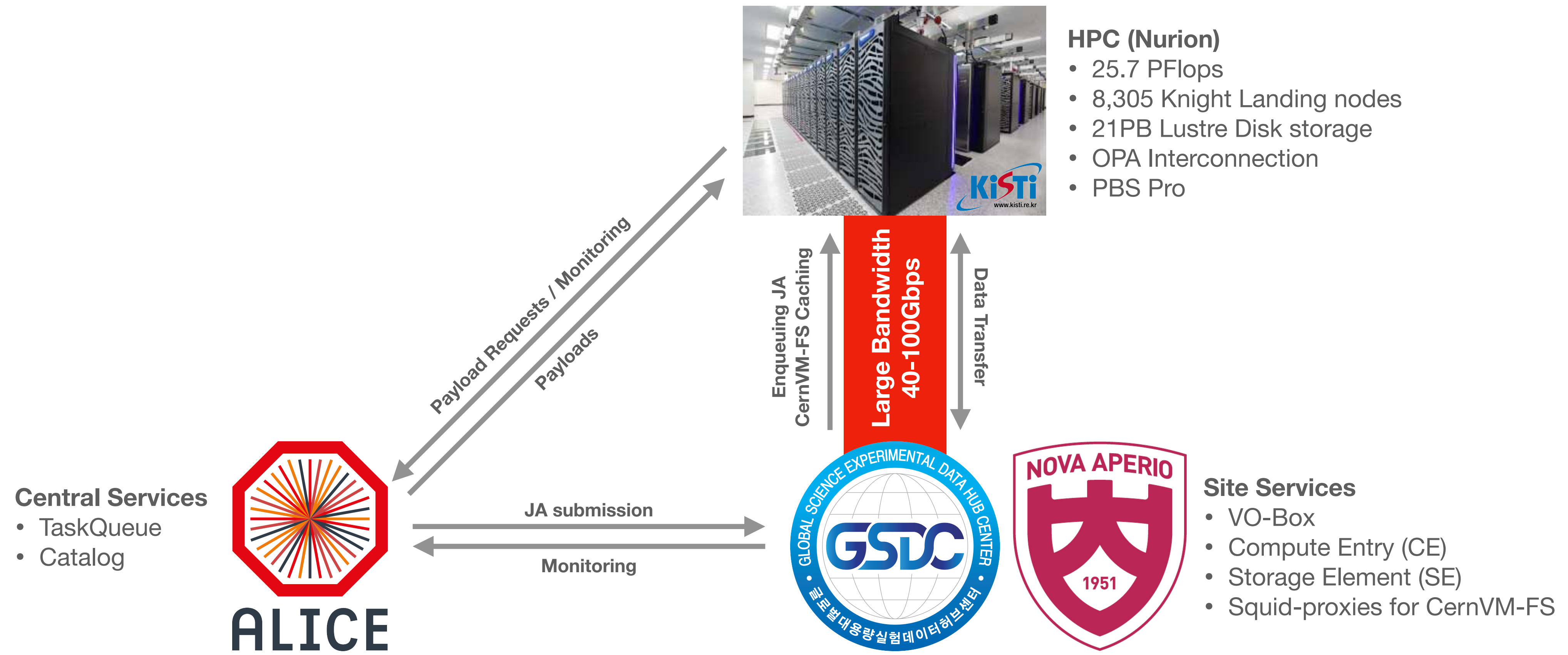
# Plan

- Further work on EOS deployment playbook to run on AWX system
- Expanding EOS Disk for ALICE further up to 7.5 PB to meet pledges
  - FST nodes running on bare metals (2PB) to be decommissioned
    - Group draining could help to vacate there FSTs
- Updating EOS CDS to v5 as well as upgrading to EL9 flavour
  - Heavy revisions required on CDS Docker deployment

HPC

# HPC for ALICE @ KISTI

Collaboration with KISTI Nurion Team & CBNU (a member of KoALICE)



# Considerations

- A "Knight Landing" node has 68 cores and 96GB memory (1.4GB RAM per core)
- Assuming 100MB/s bandwidth per single job slot, (may) start with 1,020 cores (15 nodes) for testing

## ALICE HPC Requirements

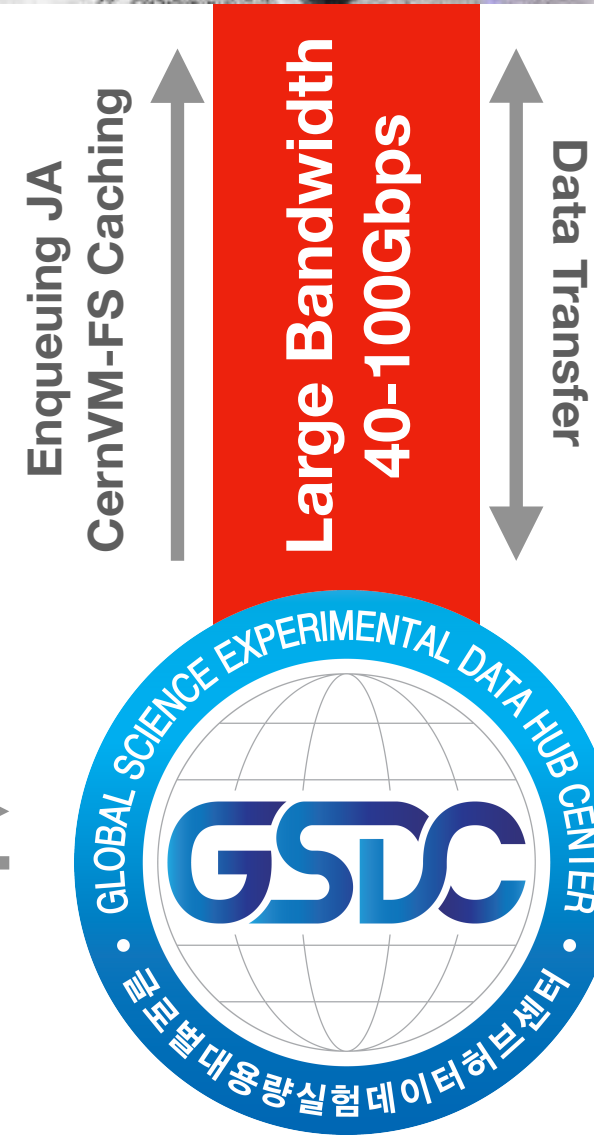
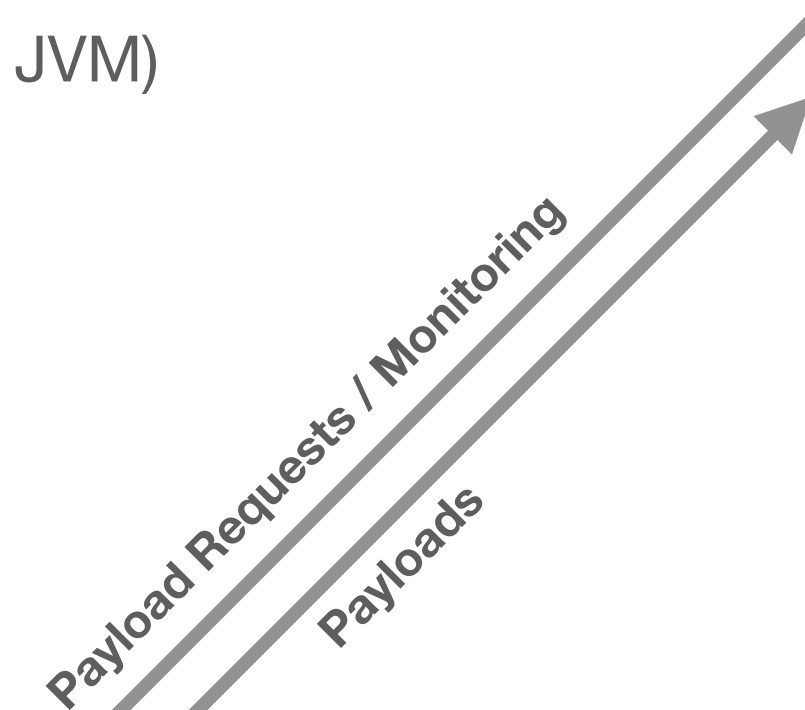
- Preemption policy allowing backfilling
- Java compatibility (ALICE Job Agent runs on JVM)
- CernVM-FS
- 2GB RAM & 10GB space per single core job (resources proportional to # of cores)

## Central Services

- TaskQueue
- Catalog



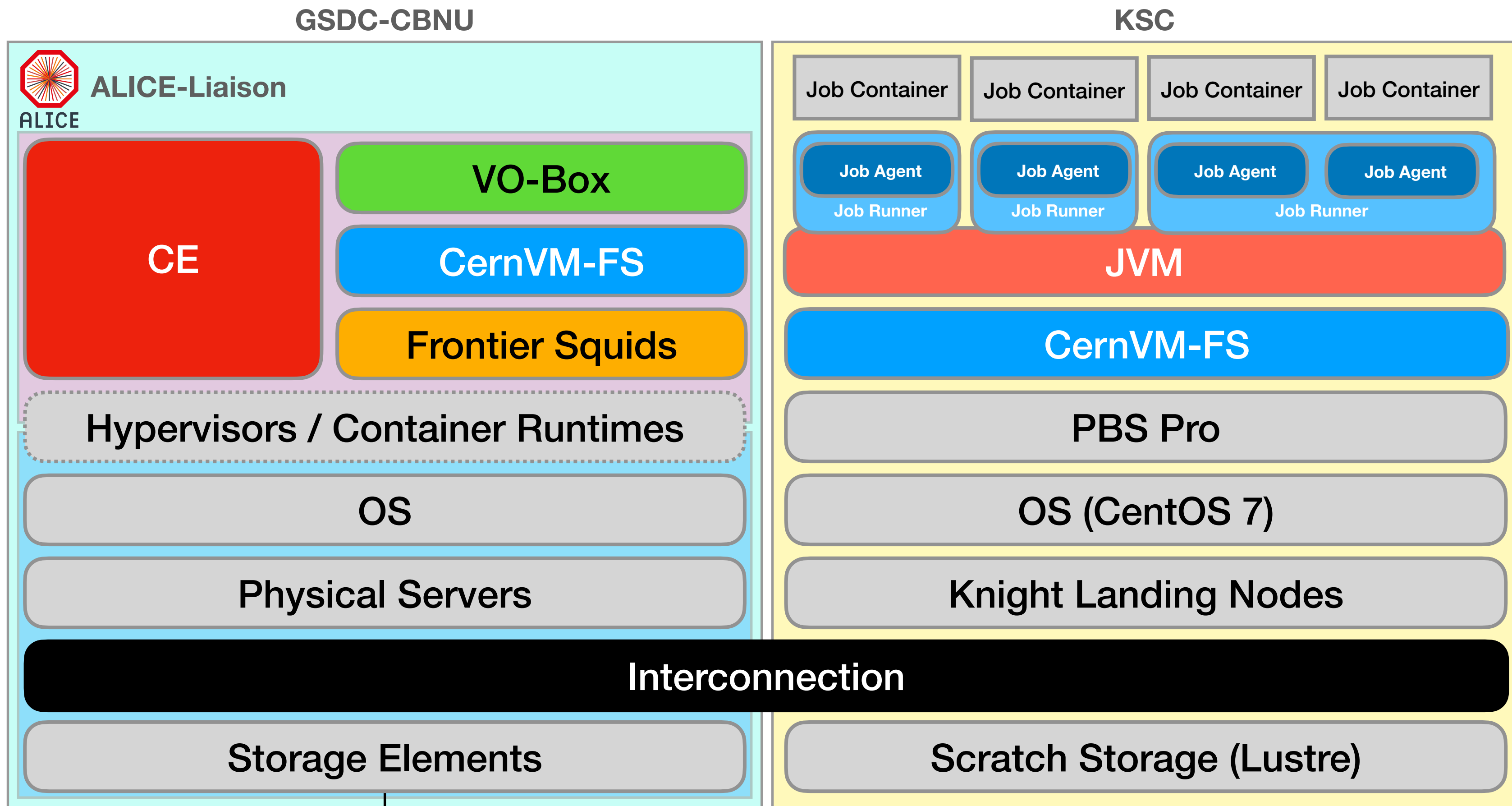
- Bypassing external authentication procedure that allows (direct) submission from CE (or VO-Box) @ GSDC
- Pool accounts (or dedicated one to be mapped)
- Scratch storage (spooling input files, storing output)
- autofs for CernVM-FS mount
- Java Runtime to run JVM binaries on CernVM-FS
- Allowing out-going connectivity
- ...



- Standalone VO-Box (independent of T1)  
Unique Site identifier (e.g. KISTI\_Nurion)
- Standalone CE (independent of T1)  
Compatible to HPC batch scheduler  
PBS Pro + ARC-CE (or direct submission?)
- Accessible SEs @ GSDC
- Standalone (load-balanced) Squid-proxy
- ...

# Prototyping Site Services

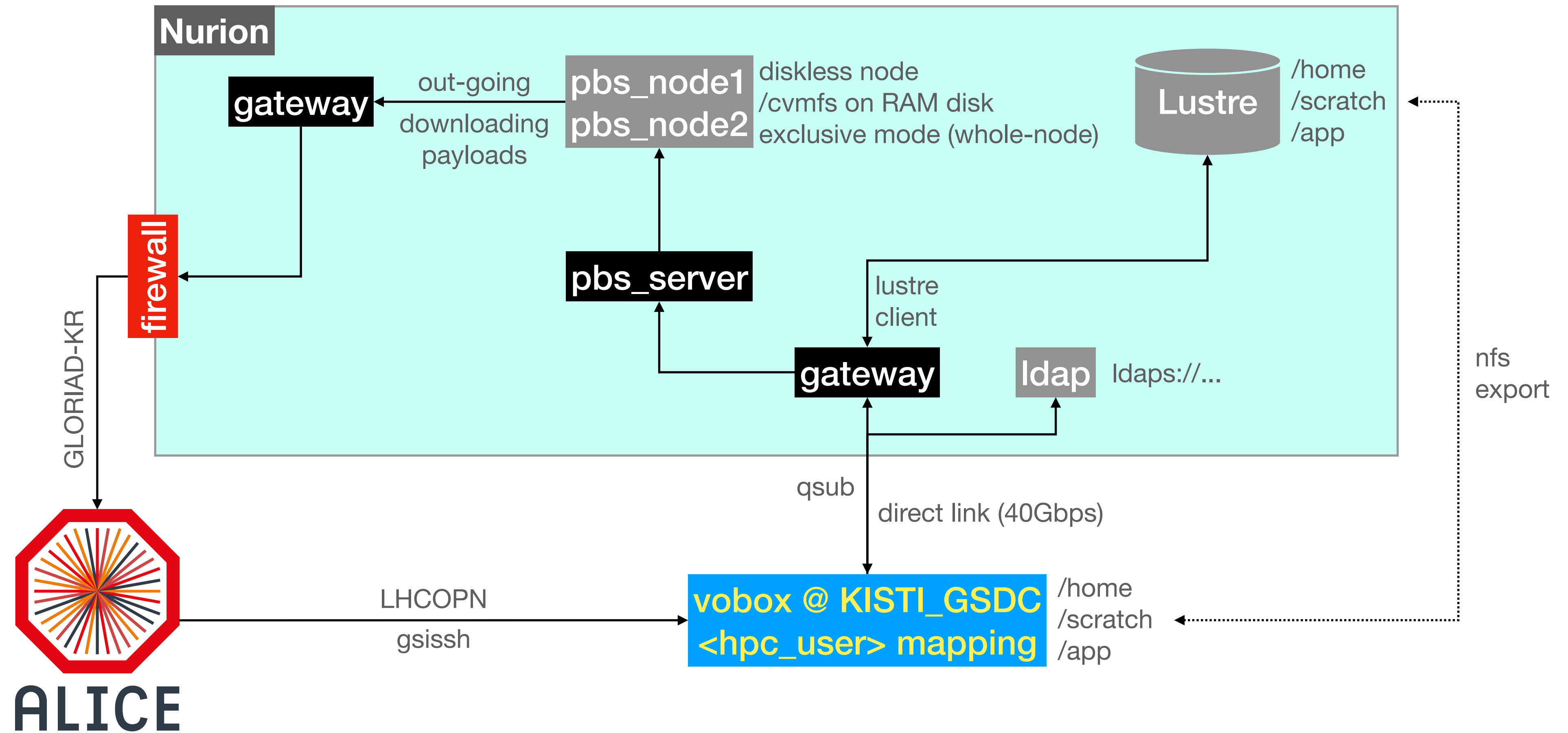
## KISTI\_GSDC\_Nurion



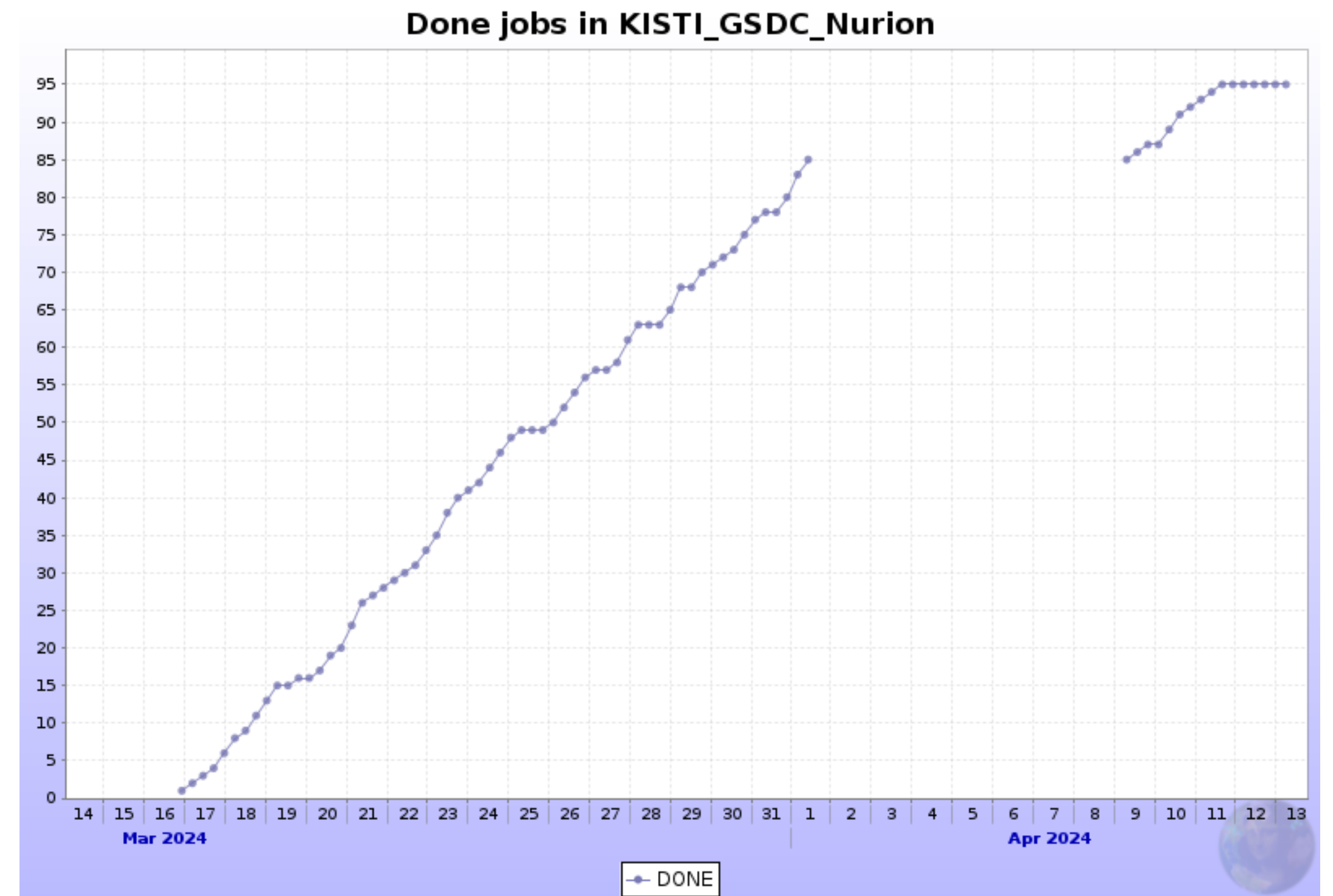
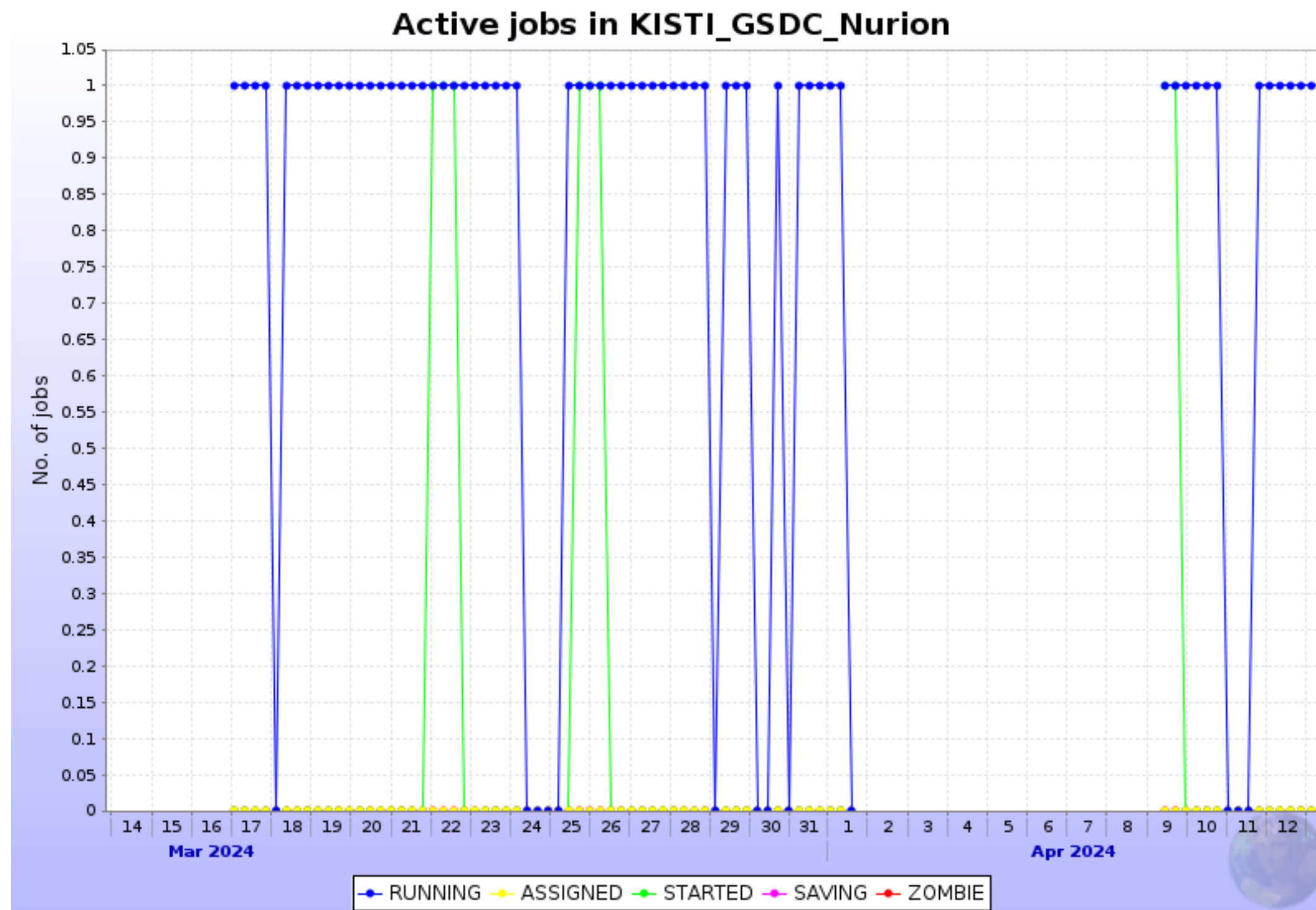
ALICE::KISTI\_GSDC::SE2  
 ALICE::KISTI\_GSDC::EOS  
 ALICE::KISTI\_GSDC::CDS



# Current Test-bed Status



# Successful Jobs





# Summary

# Summary

- KISTI Tier-1 for ALICE experiment has been operating without critical issues
  - Configuration change made for HTCondor to accept multi-core as well as whole-node submission jobs
  - New and powerful machines were deployed to meet CPU pledges
  - Successful participation to WLCG DC24 in early this year
  - OPN backup network and storage OS issues affected site availability (re-computation approved)
- Major EOS disk deployment with container technology on EL9 flavour has been successful
  - With efforts on deployment automation and systemd integration
- KISTI HPC Project has conducted in collaboration with a KoALICE member and KISTI Nurion Team
  - Successful jobs have been observed since March this year