



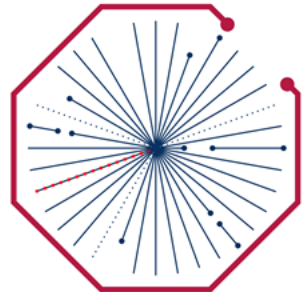
ALICE

US sites - operations and resources planning

T1/T2 Workshop, Seoul, Korea

18 April 2024

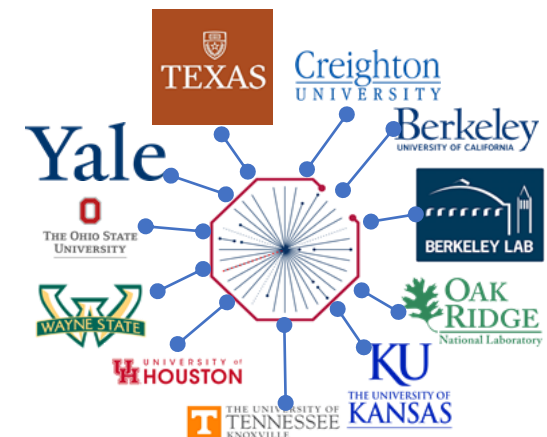
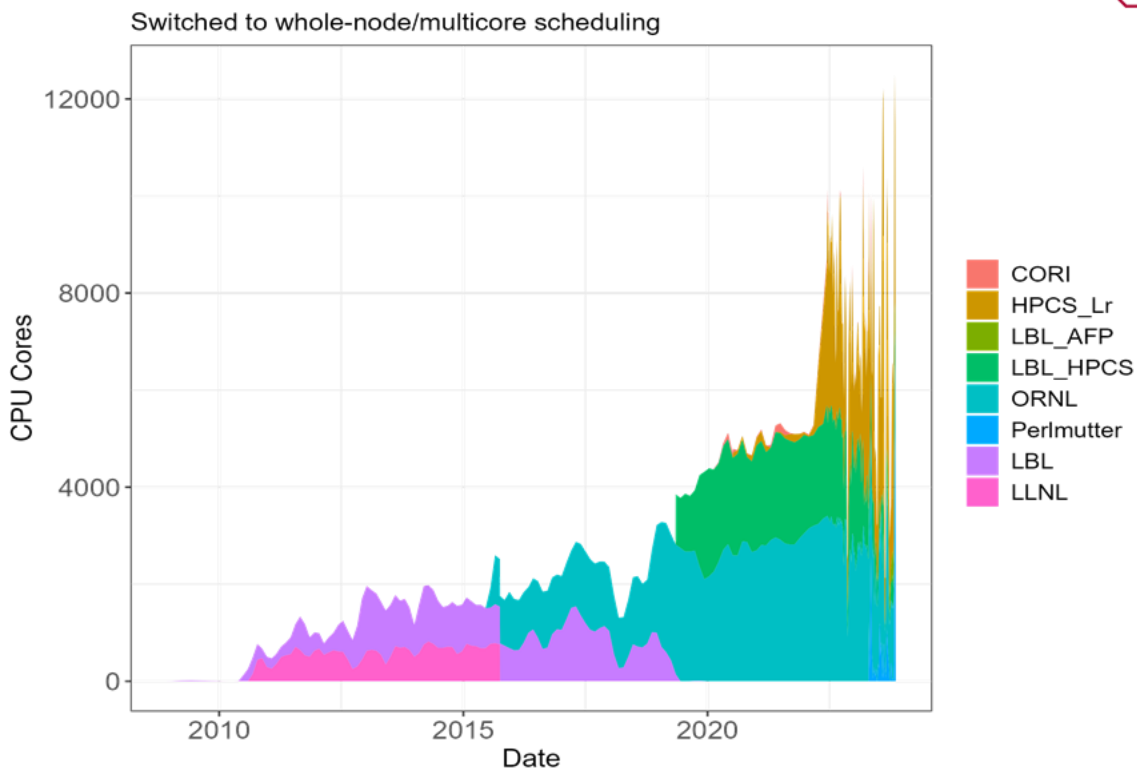
Irakli Chakaberia

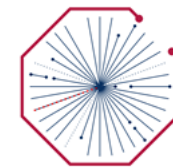


US Department of Energy Supported ALICE-USA Group



- ALICE-USA Computing project provides and maintains compute and storage resources for ALICE
- Fulfills DOE funded MoU-based ALICE USA obligations for compute and storage resources to ALICE
- Operates ALICE grid facilities at 2 DOE labs
- Project was proposed in 2009
- Operational changes 2017-2021
 - ORNL T2 relocated within ORNL, personnel retained
 - New LBNL/ITD cluster - HPCS, PDSF retired in 2019
 - Report all resources to the WLCG under the US_LBNL_ALICE federation
 - Funded R&D use of HPC resources
 - Opportunistic use of Lawrencium and Perlmutter supercomputers at LBNL





US-ALICE Review Meetings

- External project review in July 2021
- We hold annual meetings at LBNL and ORNL

- Last annual meeting at ORNL in October of 2023

29-30 July 2021
Online
America/Los_Angeles timezone

External DOE Review of LICE-USA Computing Project

- Overview
- Timetable**
- Scientific Programme
- Contribution List
- My Conference
- My Contributions
- Participant List
- Organization
- Contact**
- ✉ iraklic@lbl.gov

Timetable

< Thu 29/07 Fri 30/07 All days >

Print PDF Full screen Detailed view Filter

07:00	Executive Session	
	Online	07:30 - 08:00
08:00	Introductions	Jefferson Porter
	Online	08:00 - 08:05
	ALICE-USA	Mateusz Ploskon

ALICE-USA Computing Meeting @ ORNL

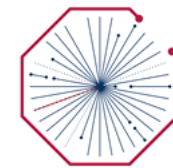
3-5 Oct 2023
Oak Ridge National Laboratory
America/Los_Angeles timezone

Enter your search term

- Overview
- Timetable**
- Contribution List
- My Conference
- My Contributions
- Registration
- Participant List
- Videoconference
- Contact**
- ✉ iraklic@lbl.gov
- ✉ moultonsa@ornl.gov
- ☎ 865 386 3733

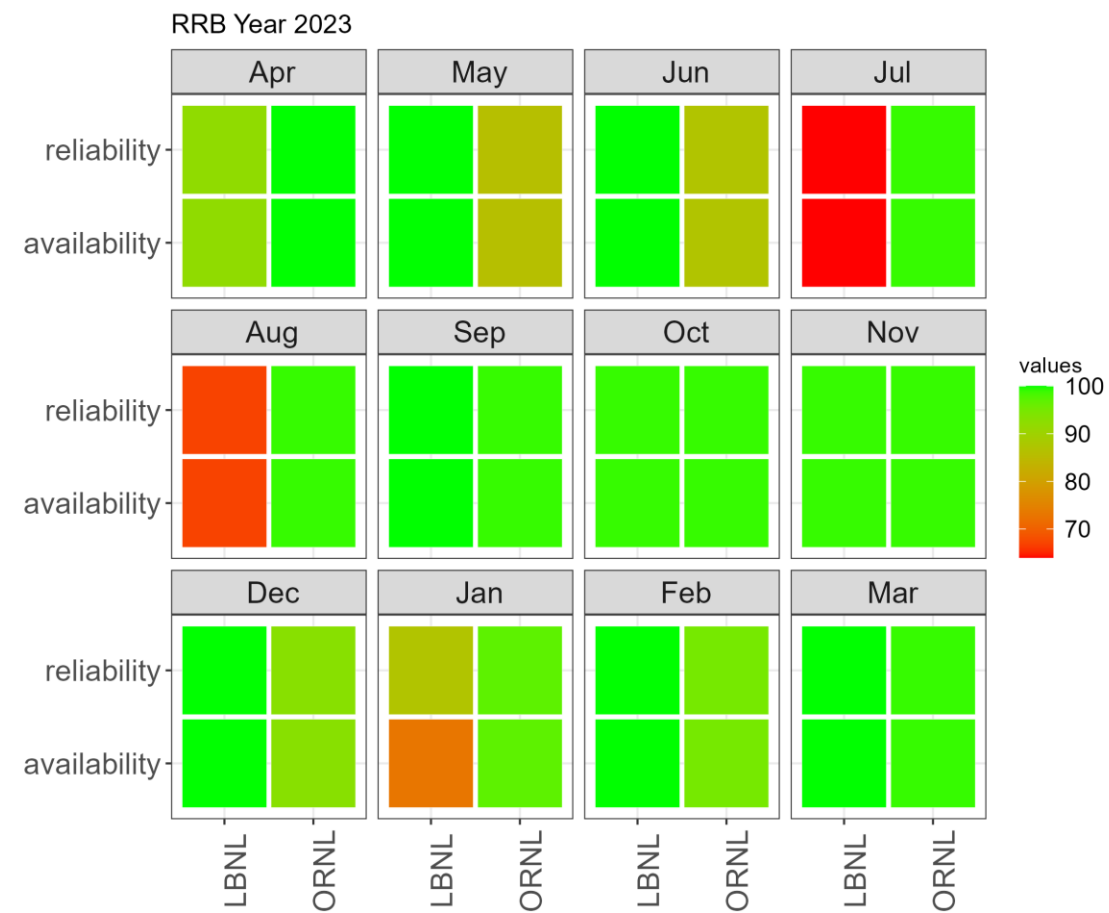
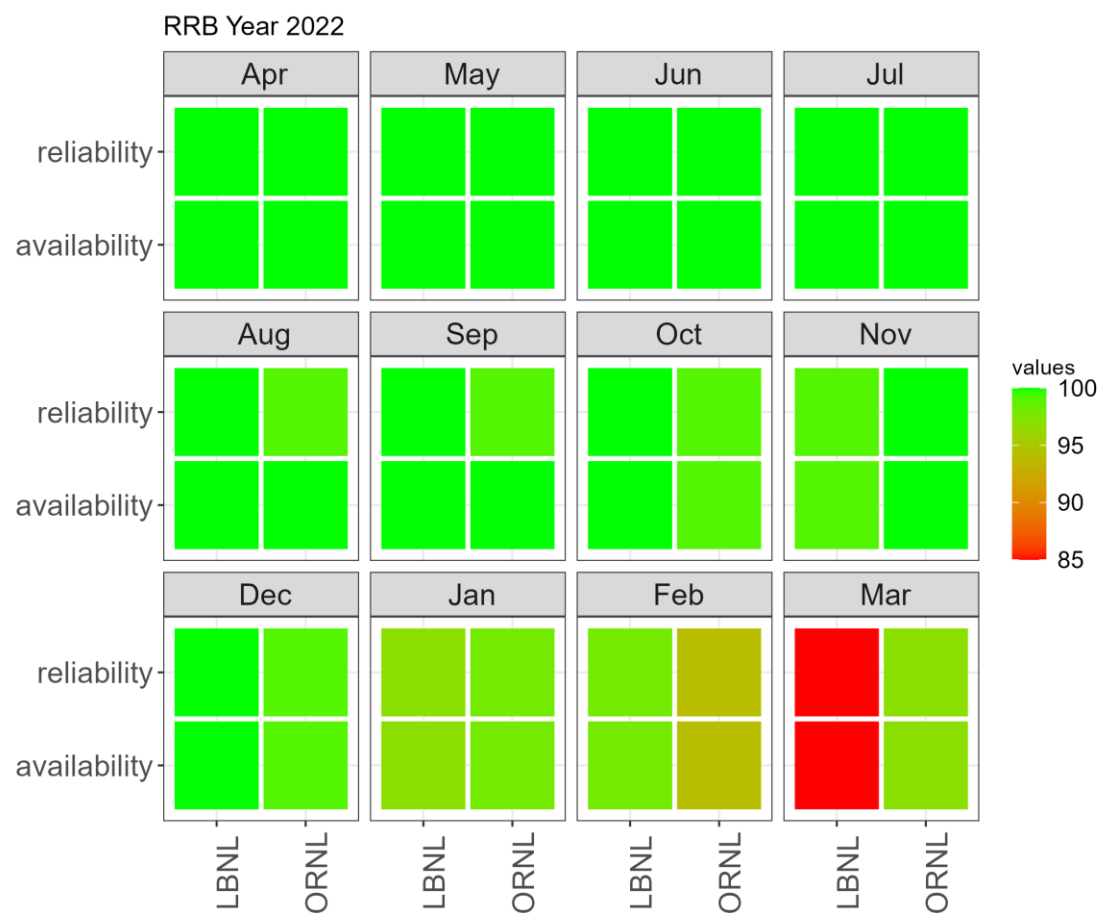
Annual meeting hosted by the Oak Ridge National Laboratory to discuss progress and plan for the future of the ALICE-USA Computing Project.

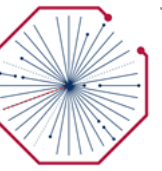




WLCG Summary | RRB 2022 & RRB 2023

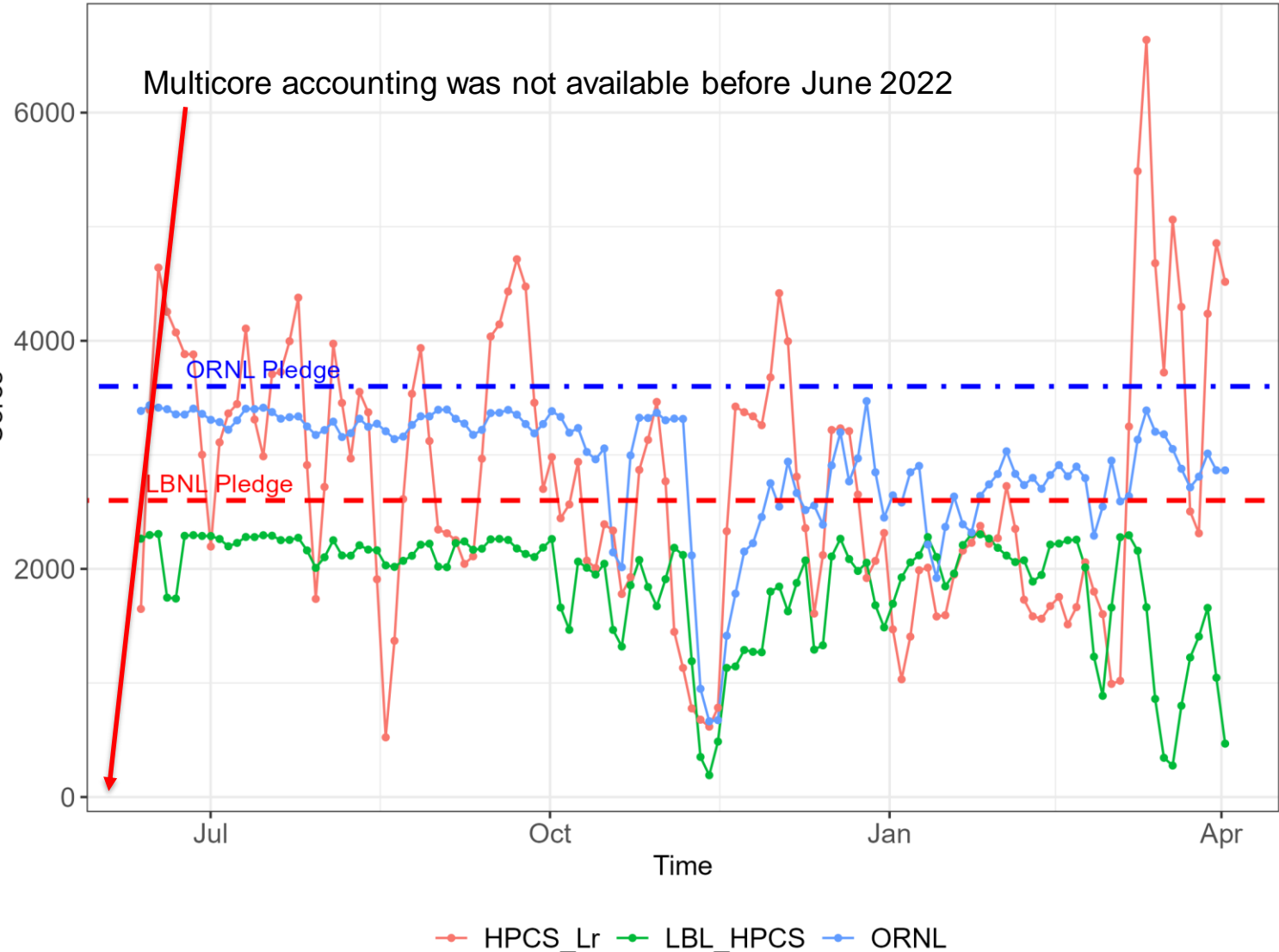
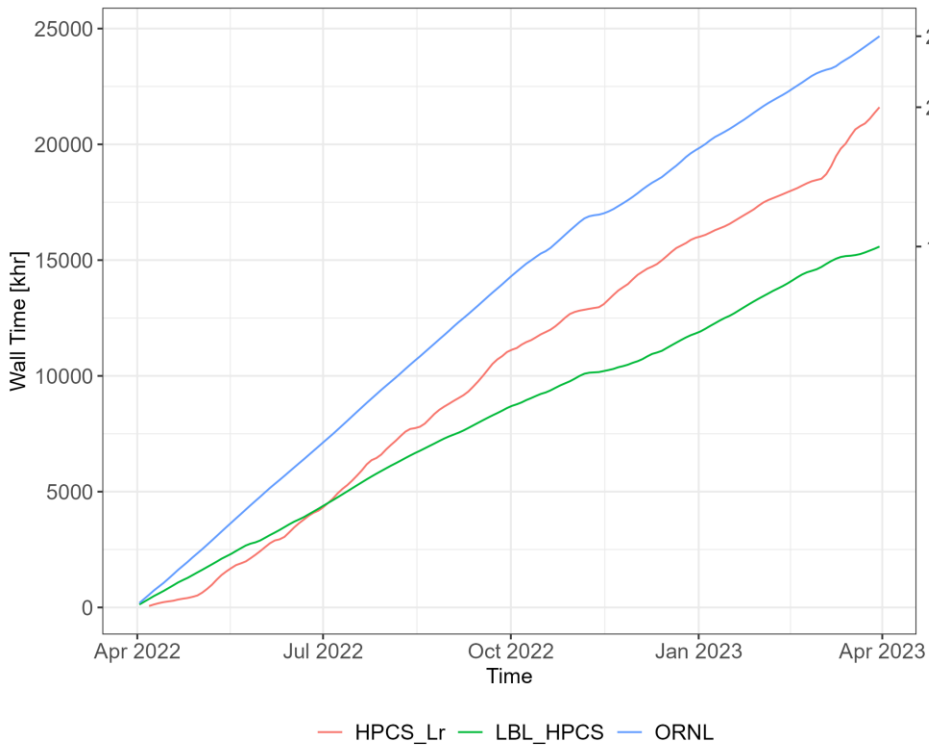
- Target Availability for each site is 97.0%.
- Availability Algorithm: $@ALICE_CE * @ALICE_VOBOX * \text{all AliEn-SE}$
- The deviation from 100% shows below is due to the SE problems (all addressed)

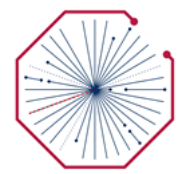




CPU Performance | RRB 2022

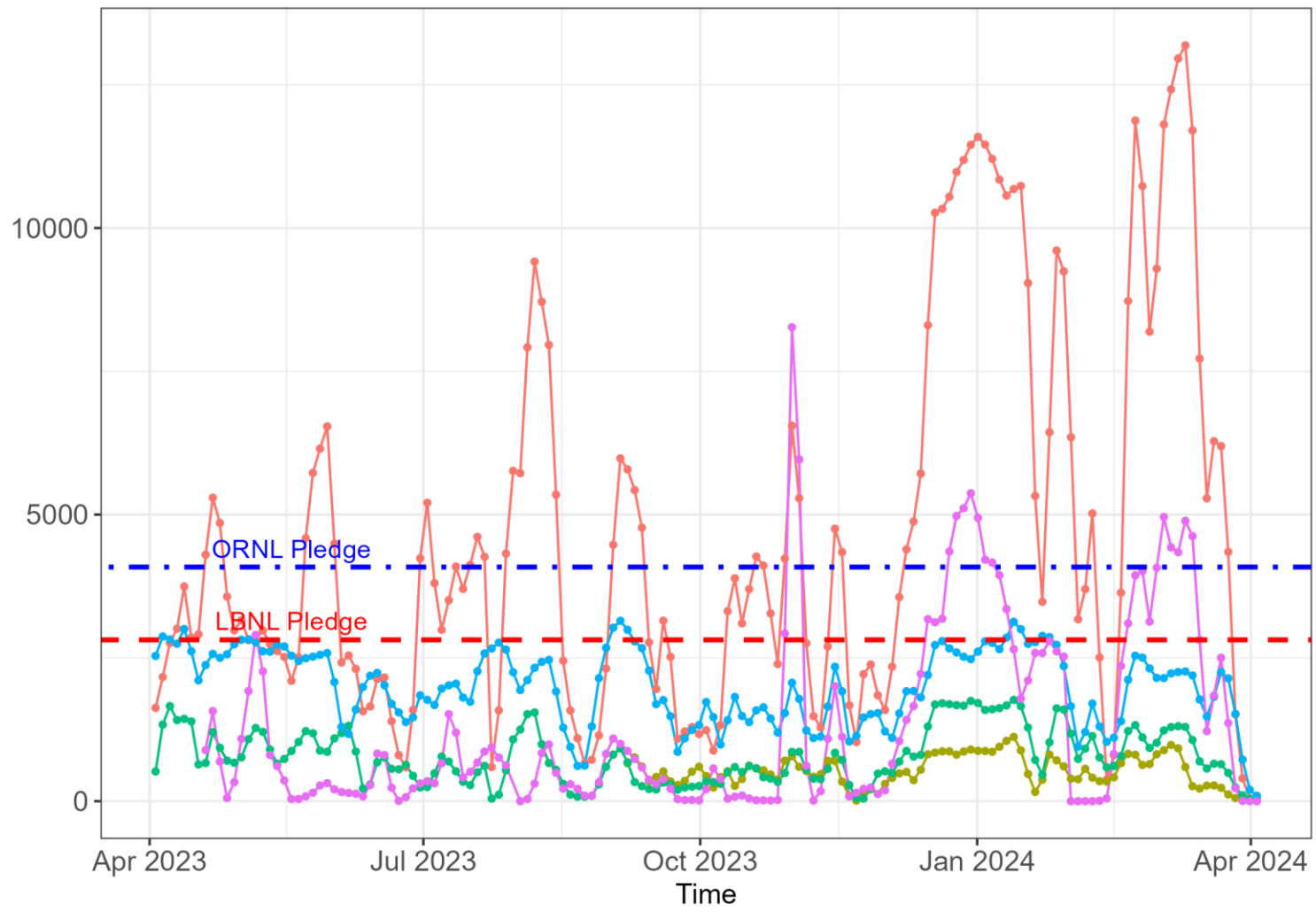
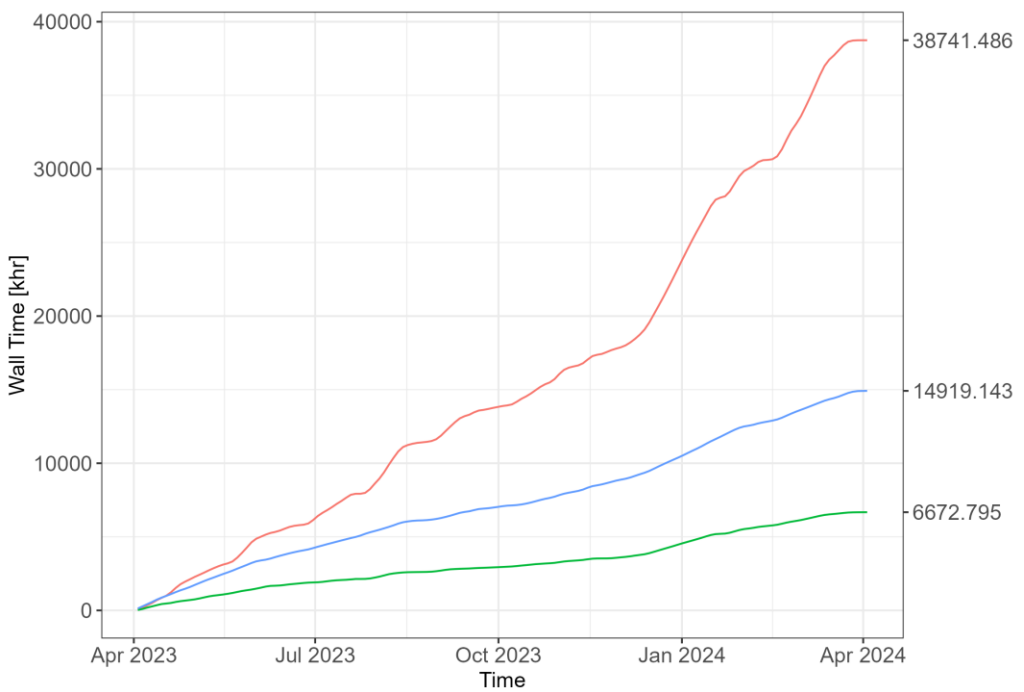
- Lawrencium a full member of the family
- Transition period from Cori to Perlmutter





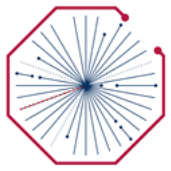
CPU Performance | RRB 2023

- Added new partitions to Lawrencium scavenging queue for alice
- Perlmutter is now in
- New AF prototype also visible



— HPCS_Lr — LBL/AFP — LBL_HPCS — ORNL — Perlmutter

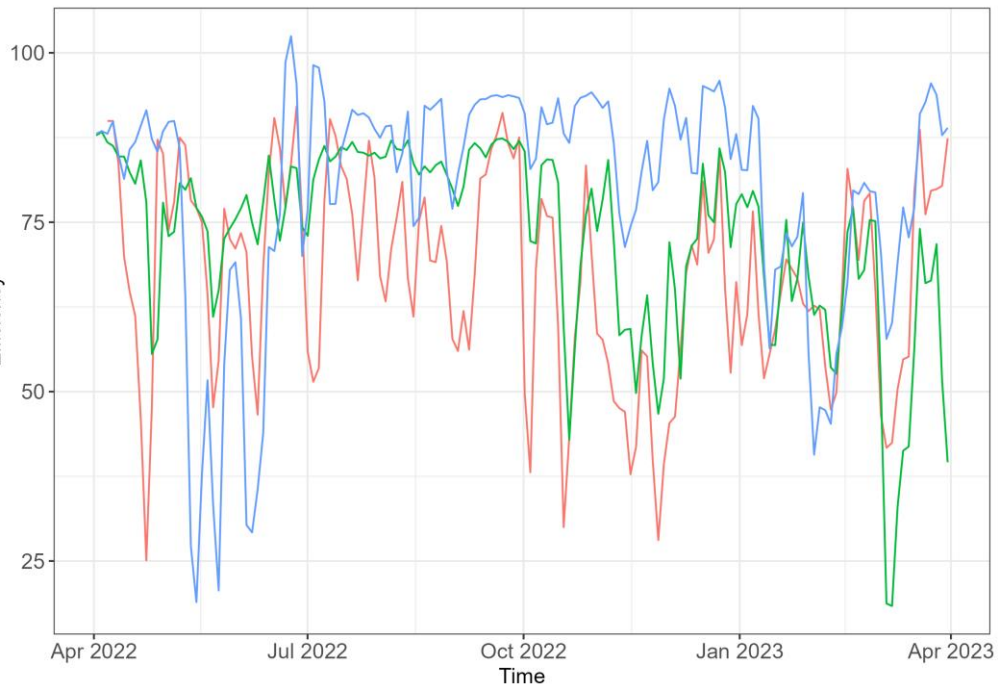
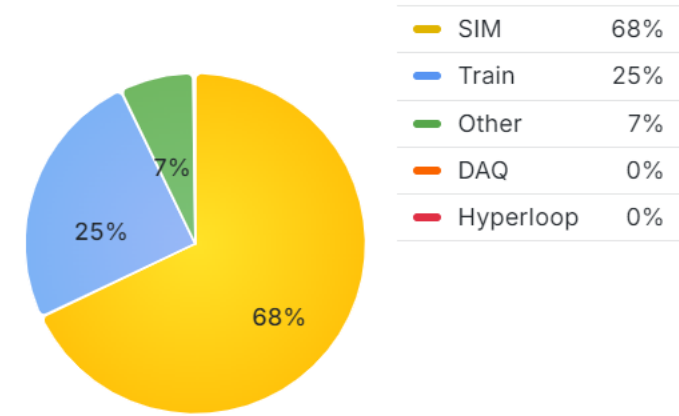
— HPCS_Lr — LBL_HPCS — ORNL



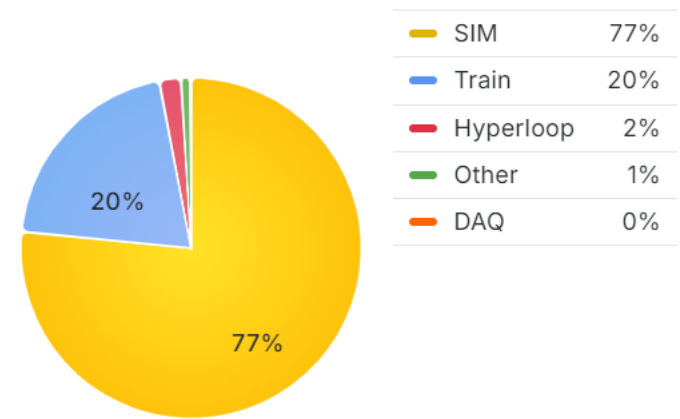
Efficiency | RRB 2022

- JobMix largely defines the average CPU efficiency delivered
- However, with the whole node scheduling this figure of merit need to be taken with a grain of salt

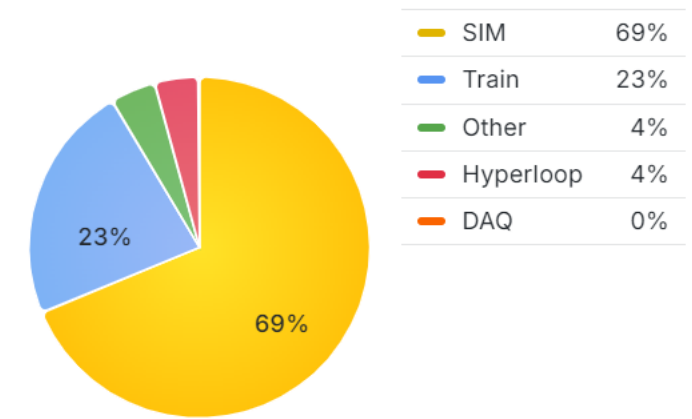
Average Job Mix [LBL_HPCS] ⓘ

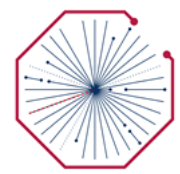


Average Job Mix [HPCS_Lr] ⓘ



Average Job Mix [ORNL] ⓘ

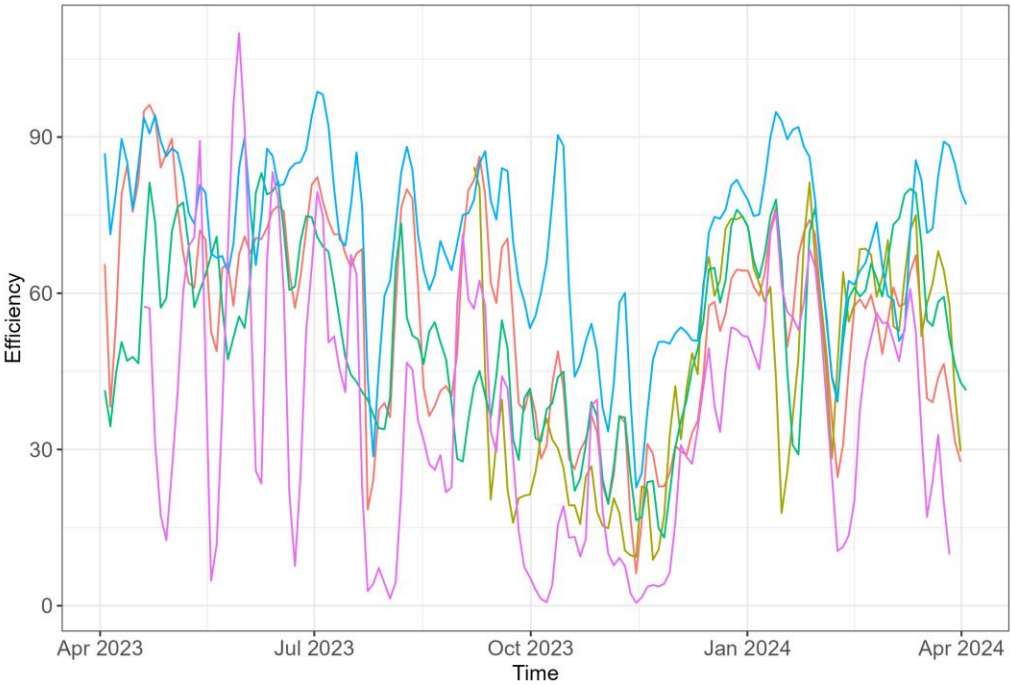
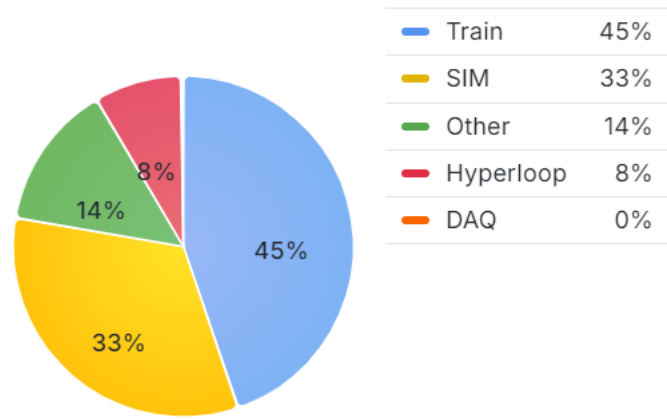




Efficiency | RRB 2023

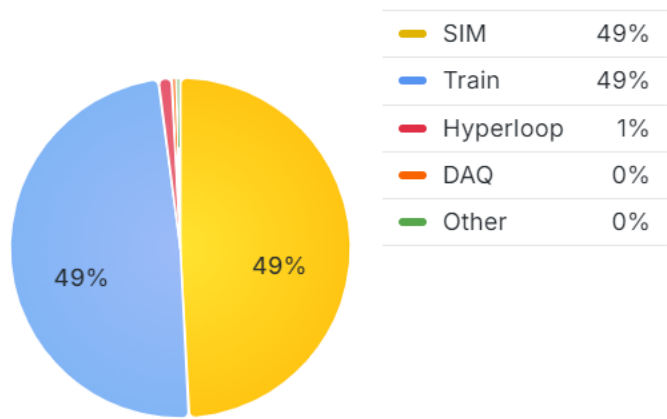
- JobMix largely defines the average CPU efficiency delivered
- Substantial drop in fraction of MC jobs compared to 2022 (or earlier)

Average Job Mix [LBL_HPCS] ⓘ

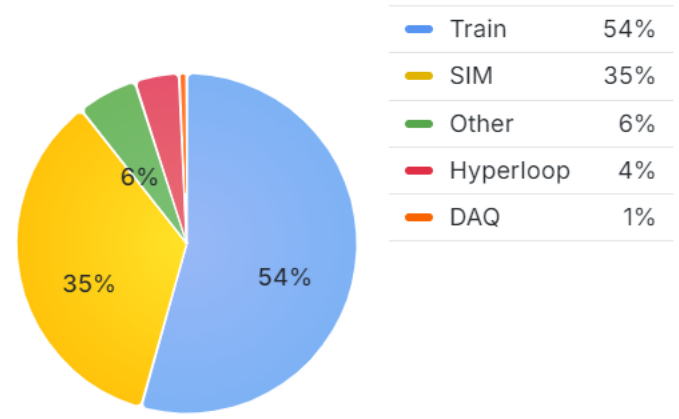


— HPCS_Lr — LBL_AFP — LBL_HPCS — ORNL — Perlmutter

Average Job Mix [HPCS_Lr] ⓘ



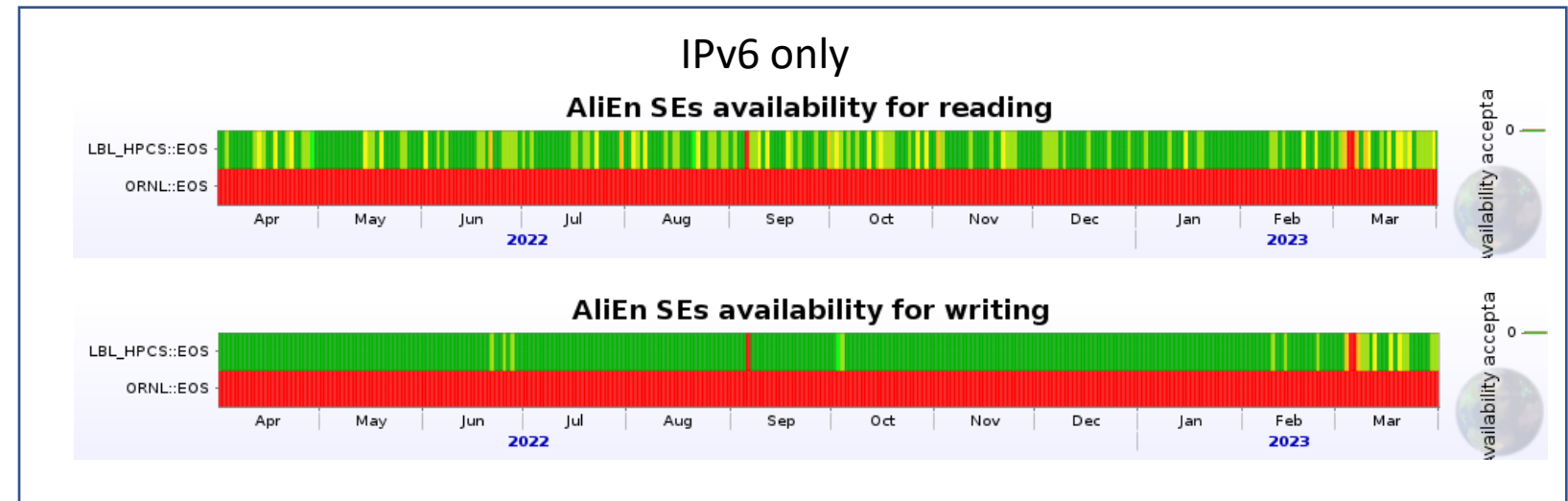
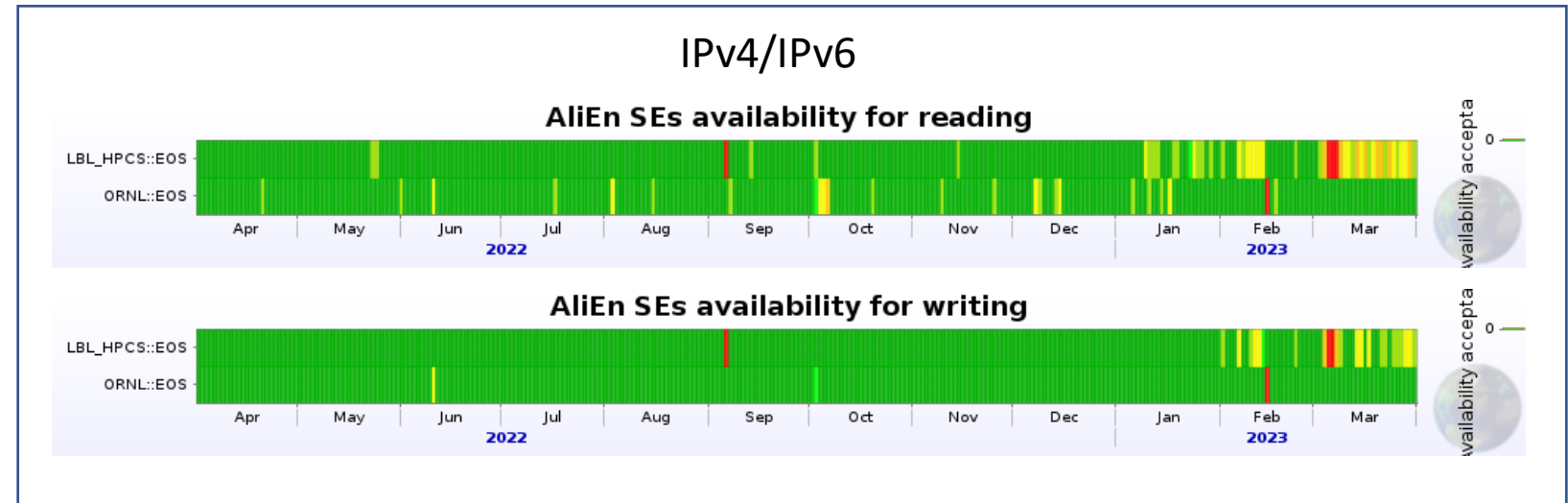
Average Job Mix [ORNL] ⓘ



Storage Performance | RRB 2022

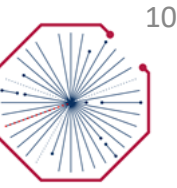


- SE availability has been very high – over 98% on average
- IPv6 is up at LBNL since RRB2021
- IPv6 @ ORNL is largely configured and soon to be extended onto our SE

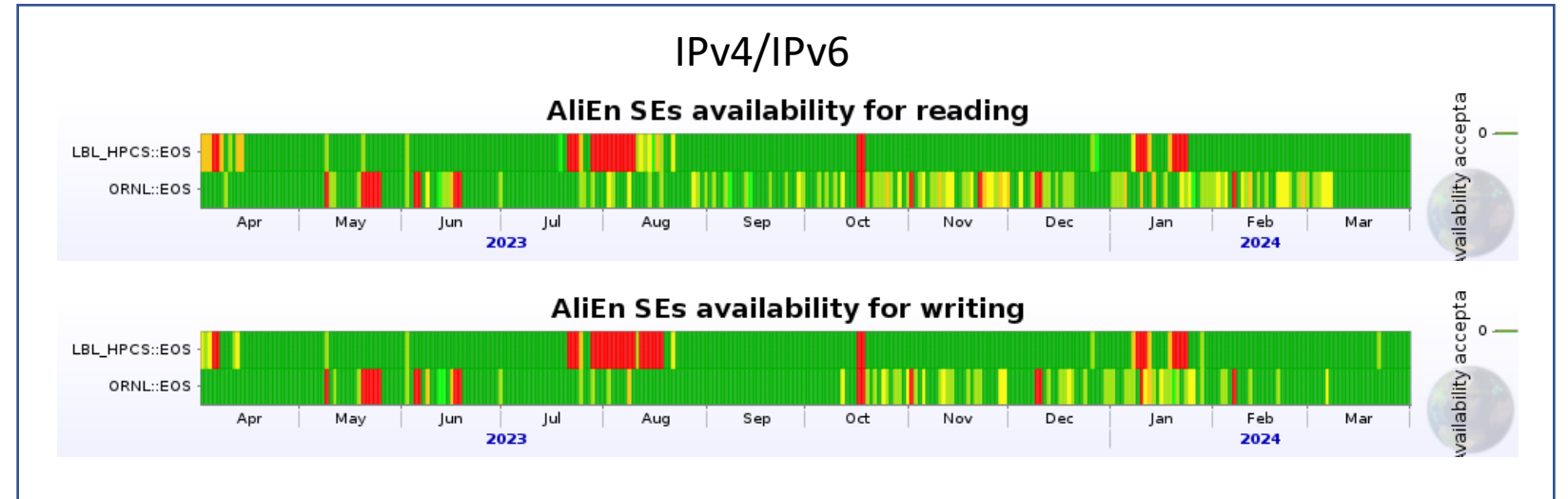


	Read	Write
LBNL	98.16	98.79
ORNL	99.43	99.77

Storage Performance | RRB 2023



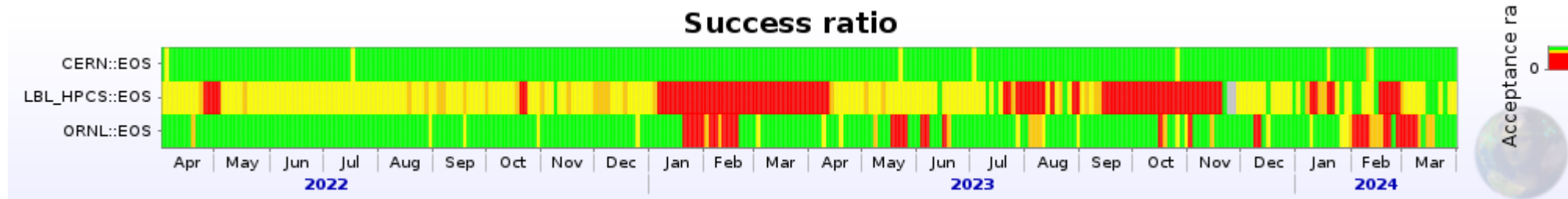
- SE availability has been very high – over 90% on average
- Downtime and some issues during EOS 4 to EOS 5 switch
- The problems are fixed, and we are back to near 100% availability in the past months



	Read	Write
LBNL	91.84	90.38
ORNL	95.06	95.50

Some EOS Problems

- At LBNL we notice that number of corrupted files were pretty large
- Investigations were hard and after involving EOS experts ((Huge thank you to Alvin) a
- About 1.4 million files from about 62 million files reported rep_missing_n
- No fsck repair possible due to single-disk fsid setup
- clearing rep_missing_n was unsuccessful - we had not migrated FMD to FSTs (REQUIRED FOR 5.2!!!!)
- Once FMD move was done, the prior work evacuating the rep_missing_n files directly via qdb commands was successful
- We also had some EOS problems at ORNL after migrating to EOS 5



ALICE-USA Obligations

- ALICE-USA share of resources in respect to WLCG / C-RCG recommendations for the next few years and planned hardware acquisition strategy table
- Starting 2025 we accounted for the share redistribution of the lost Russian resources

Year	FY2021	FY2022	FY2023	FY2024	FY2025	FY2026	FY2027
ALICE Requirements							
CPU (kHS06)	852	1013	1164	1280	1475	1549	1626
Disk (PB)	91	104	121	138	160	168	176
ALICE-USA Participation							
(ALICE - CERN) M&O-A	561	571	571	571	526	526	526
ALICE-USA M&O-A	46	46	46	45	48	48	48
ALICE-USA/ALICE (%)	8.20%	8.06%	8.06%	7.88%	9.13%	9.13%	9.13%
ALICE-USA Obligations							
CPU (kHS06)	69.9	81.6	93.8	100.9	134.6	141.3	148.4
Disk (PB)	7.5	8.4	9.7	10.9	14.6	15.3	16.1

ALICE-USA Acquisition Strategy

- We have delayed 2023 CPU procurement, but not cancelled
- The 2023 plan is not fully implemented but is well underway (CPU part)
- Due to the ample amount of opportunistic resources this is not a problem

Resource	Installed	FY2023	FY2024	FY2025	FY2026	FY2027
LBNL HW & Costs						
CPU change (+/- kHS06)		-5.0+11.0	-5.0+5.0	0.0+12.0	0.0+5.0	-5.0+5.0
CPU Installed (kHS06)	47	47	47	59	64	64
Disk change (+/- PB)		0.0+0.0	-1.6+3.0	0.0+1.5	0.0+1.5	0.0+1.5
Disk installed (PB)	4.75	4.75	6.15	7.65	9.15	10.65
ORNL HW & Costs						
CPU change (+/- kHS06)		-11.0+14.0	0.0+5.0	-7.0+12.0	0.0+5.0	-5.0+5.0
CPU Installed (kHS06)	49	45.5	50.5	55.5	60.5	60.5
Disk change (+/- PB)		0.0+1.5	0.0+3.0	0.0+1.5	0.0+3.0	0.0+1.5
Disk installed (PB)	4.50	4.50	7.50	9.00	12.00	13.50

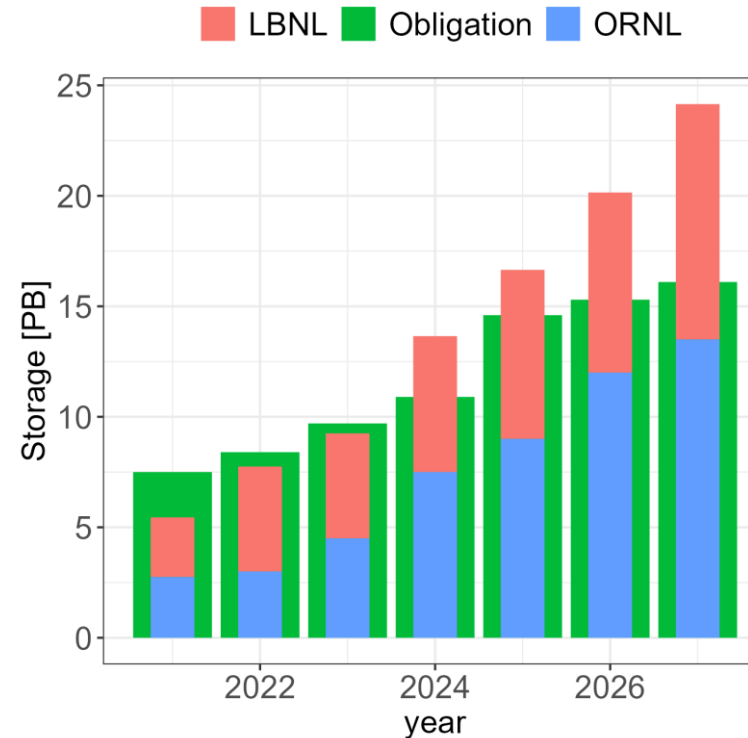
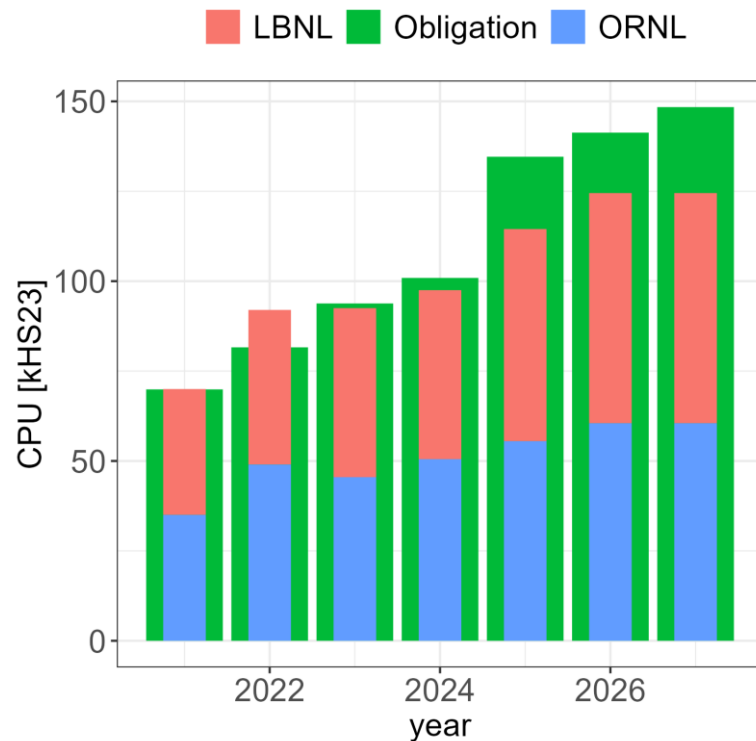
ALICE-USA Deployment Strategy

- Due to the recent heavy priority switch to the storage our plan envision to rebalance CPU resources towards the storage while still delivering pledged CPU with some reliance on the opportunistic resources
- The accounted 20 kHS23 reliance of the opportunistic resources is extremely conservatory

Resource	FY2023	FY2024	FY2025	FY2026	FY2027
ALICE-USA Obligations					
CPU (kHS23)	93.8	100.9	134.6	141.3	148.4
Disk (PB)	9.7	10.9	14.6	15.3	16.1
ALICE-USA Plan					
Opportunistic CPU Resource (kHS23)	20.0	20.0	20.0	20.0	30.0
CPU (kHS23)	92.5	97.5	114.5	124.5	124.5
% CPU obligation	120.0%	116.5%	100%	102.2%	104.1%
Disk (PB)	9.3	13.7	16.7	21.2	24.2
% Disk obligation	94.9%	125.5%	114.0%	138.0%	150.0%

CPU Acquisition

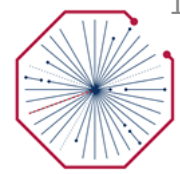
- The CPU plan below does not have 20 kHS contribution from HPC resources
- All e19 nodes coming up at ORNL very soon
- E18 switch at LBNL
- Storage plan presented for 2026 and 2026 is just based on the flat storage addition and throwing the CPU resources at the storage
- Storage plan for 2026 and beyond will be revised as ALICE requirements for those years will be assessed



Resource Delivery

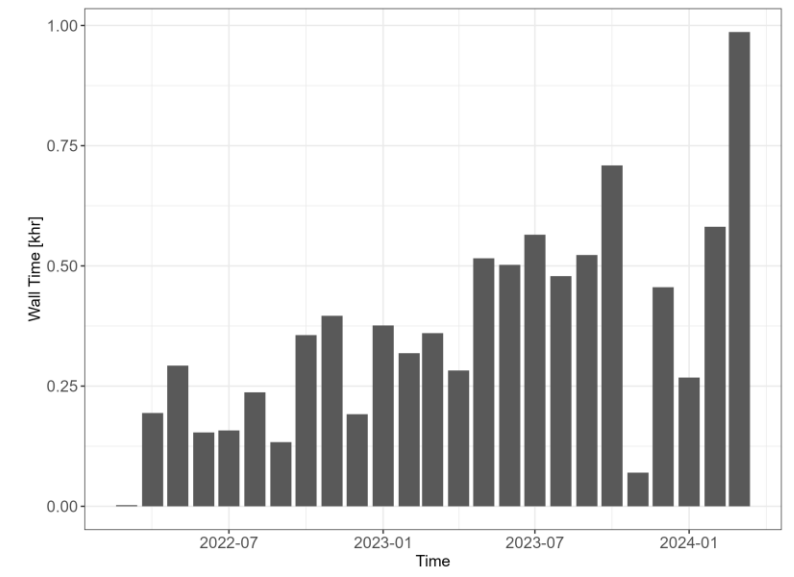
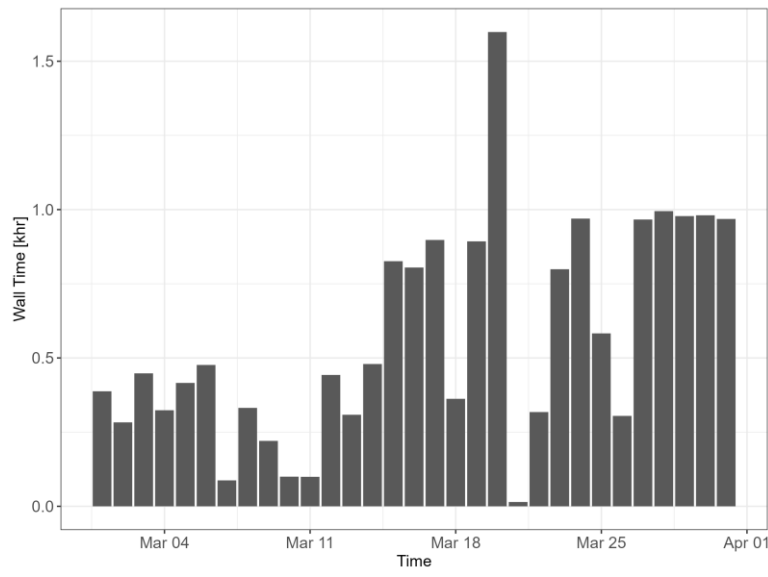
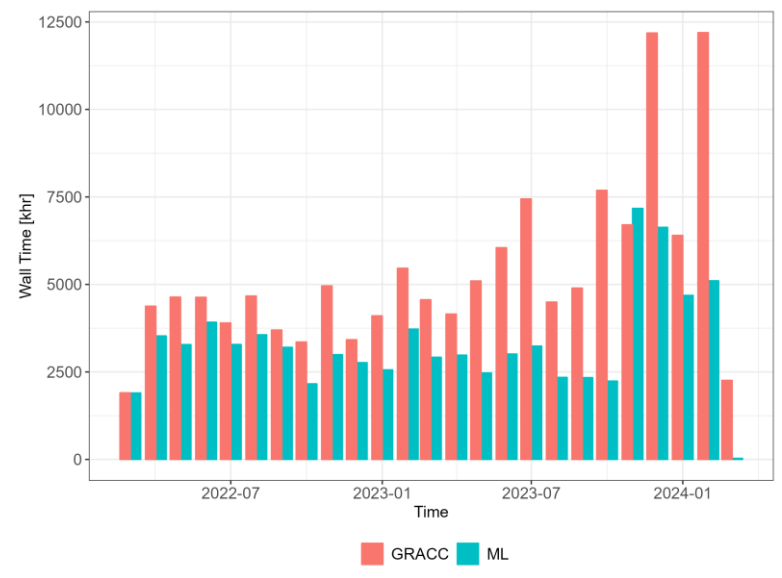
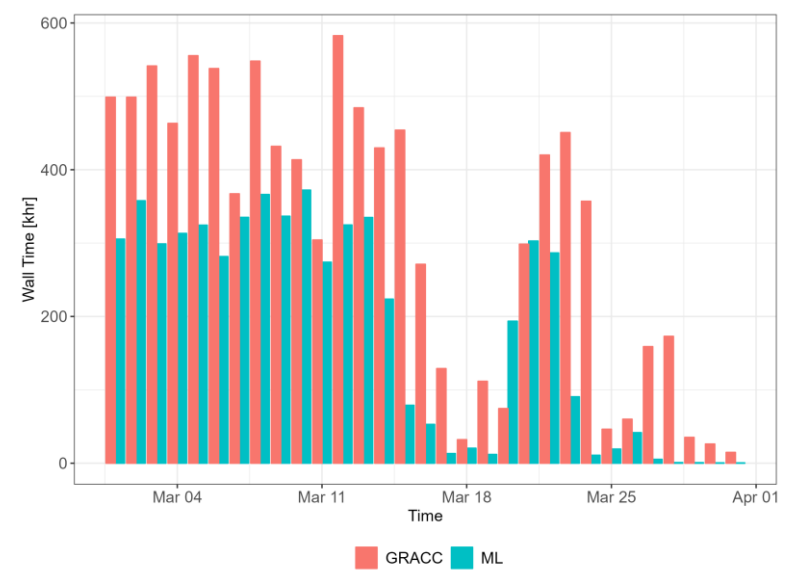
- With the Lawrence Livermore resources being part of the accounting ALICE USA easily delivers to the obligation
- Below is the report for the last quarter

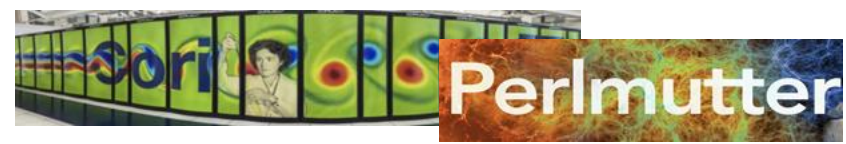
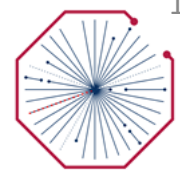
T2 Site and HPCS_Lr	CPU/Core [HS06/Core]	ALICE-USA Obligation [MHS06 x hr]	CPU Delivered [MHS06 x hr] (ML reporter)	CPU Delivered [MHS06 x hr] (WLCG reporter)	Delivered per Obligation [%] (ML reported)	Delivered per Obligation [%] (WLCG reported)
LBNL	13.00	83.27	36.17	361.40	43.44	434.00
ORNL	11.50	88.56	51.80	83.15	58.49	93.89
Total		171.83	87.97	444.55	51.19	258.71



Peculiarity in reporting

- Since the move to the whole-node scheduling there is a substantial difference in SLURM and ML reporting walltime
- The difference is understandable but perhaps its size is something to be looked at in detail

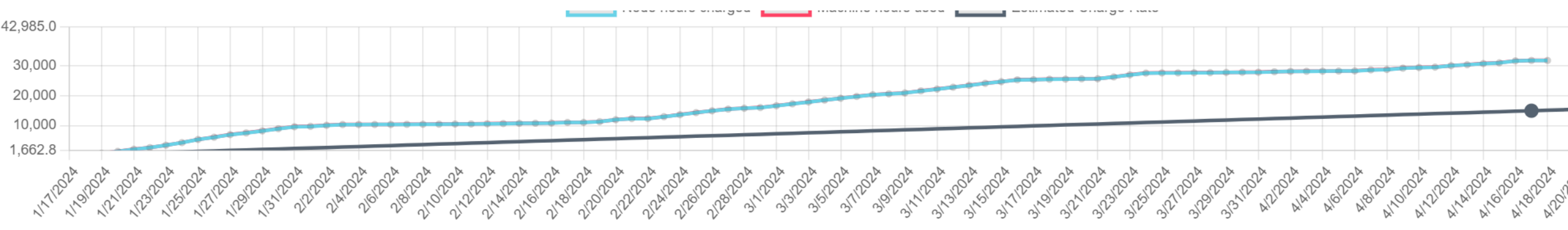


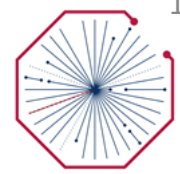


NERSC Allocation

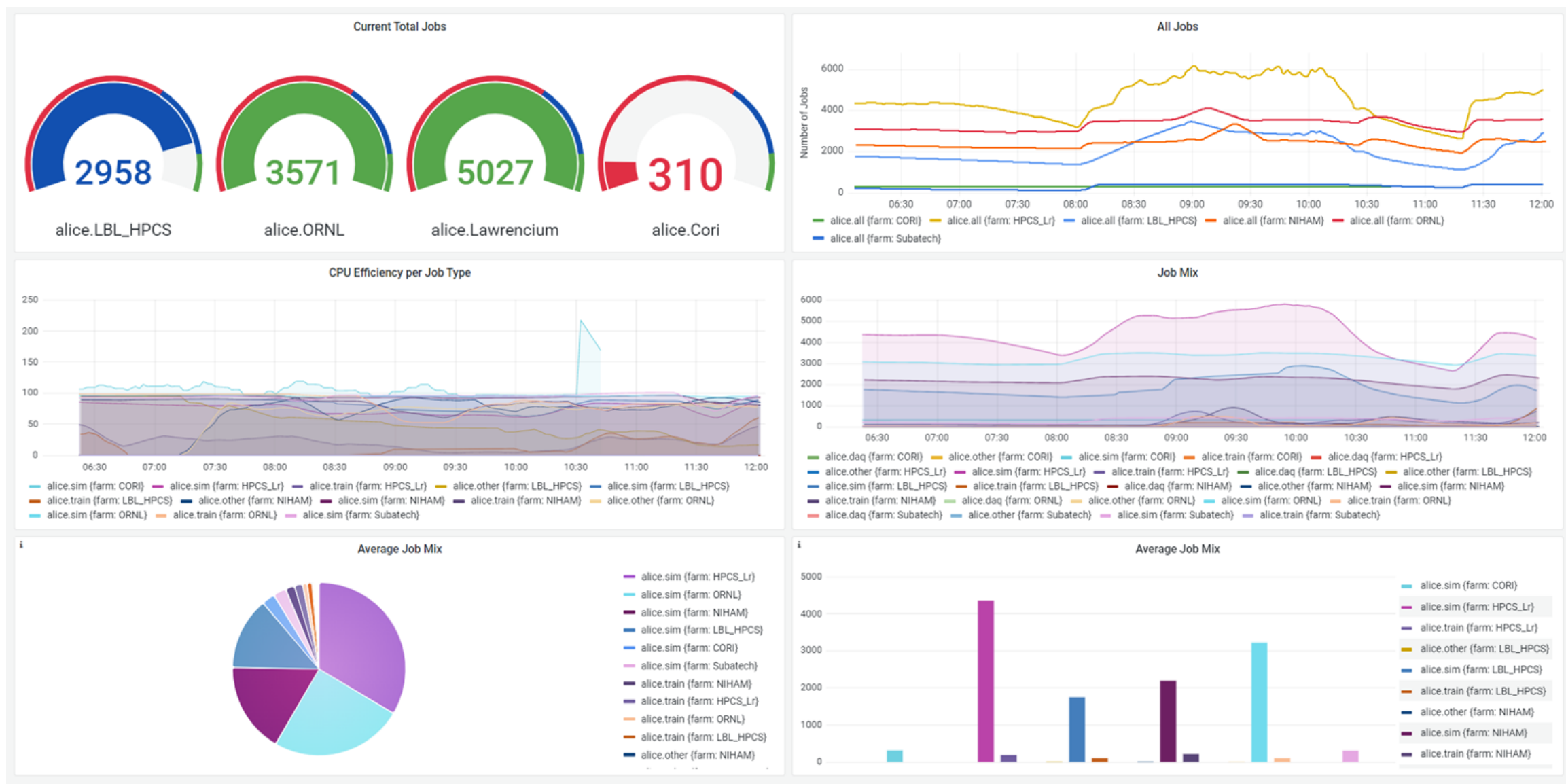
- As part of the project, we request allocation on NERSC HPCs
- Very good review of our NERSC HPC work by Sergiu yesterday
- We will continue maintaining this resource and have ambitious plans towards it as presented also yesterday

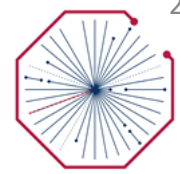
10:00	The streamlined T1 at LBNL - project proposal	Irakli Chakaberia
	31st Floor "Mozart Hall", Hotel President	10:00 - 10:30
	Perlmutter HPC integration and operation	Sergiu Weisz
	31st Floor "Mozart Hall", Hotel President	10:30 - 11:00



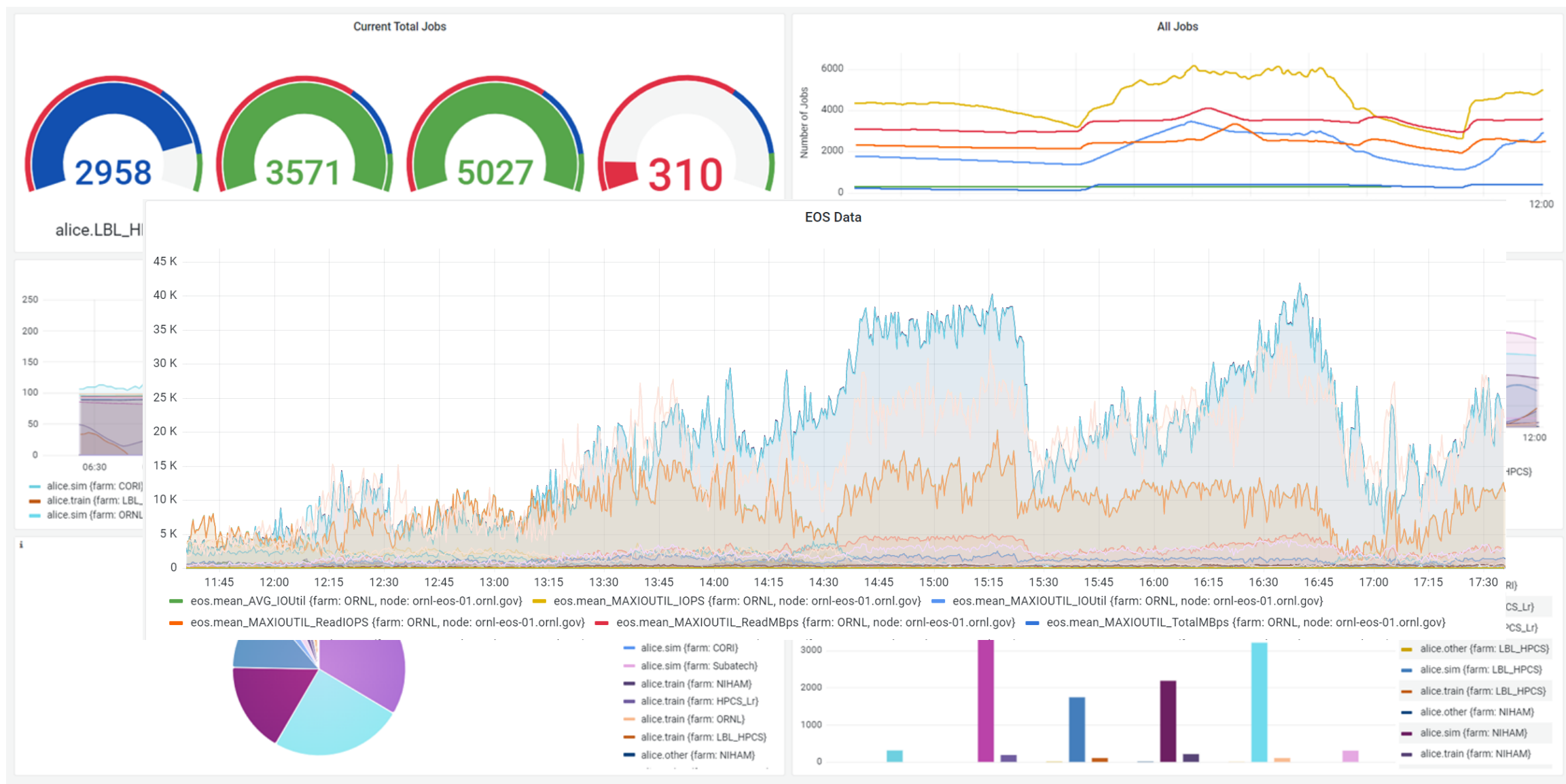


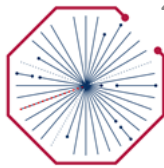
ALICE-USA Local cluster monitoring





ALICE-USA Local cluster monitoring

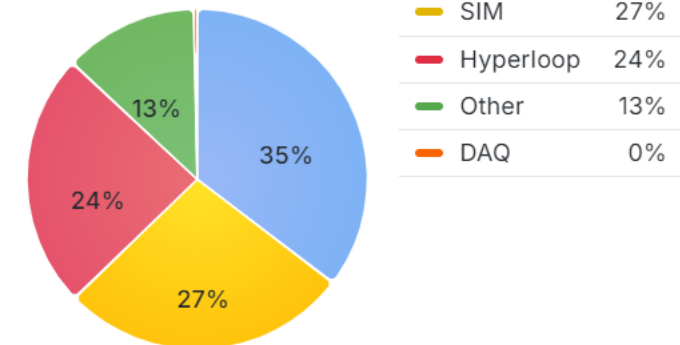




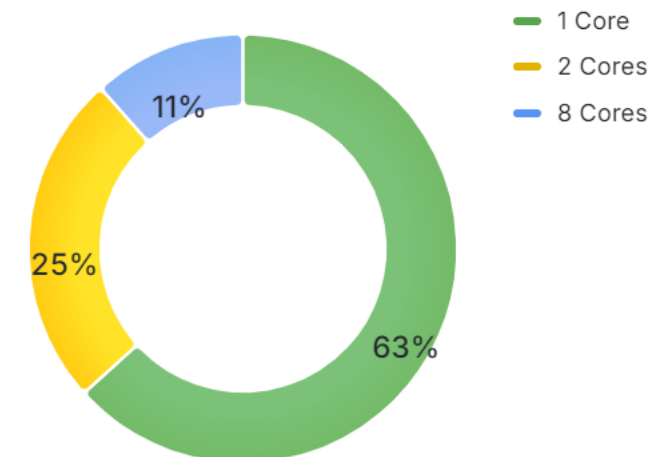
ALICE-USA Analysis Facility

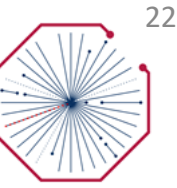
- During last T1/T2 we reported a plan for the US AF
 - In September of 2023 we deployed a prototype AF with out-of-warranty hardware
 - It has been fully operational since
 - We did learn some lessons running it for past half a year
 - Especially the drastic differences to the data access of the very “hot” datasets
 - At the moment we provide 640 CPU cores and 1.1 PB of storage
-
- The proposal to host such an AF by the ALICE-USA was presented to the DOE
 - ~2000 cores for analysis jobs
 - ~2.5 PB of usable disk storage
 - The initial proposal was well received, and we were advised to resubmit after a proof-of-concept prototype
 - The proposal will be resubmitted in the coming months

Average Job Mix [LBL_AFP] ⓘ



Core Mix [LBL_AFP] ⋮





Summary

- ALICE-USA is meeting the obligation and plans to grow accordingly
- We currently operate two T2 sites, one at LBL and one at ORNL
- Sites continue to deliver pledged resources with very high availability and reliability
- Both sites operated in whole-node scheduling mode, and we often provide platform for R&D
- Substantial resources added to the grid by utilizing Lawrenceium and Perlmutter HPCs
- We successfully deployed an AF prototype and are seeking approval for its growth