

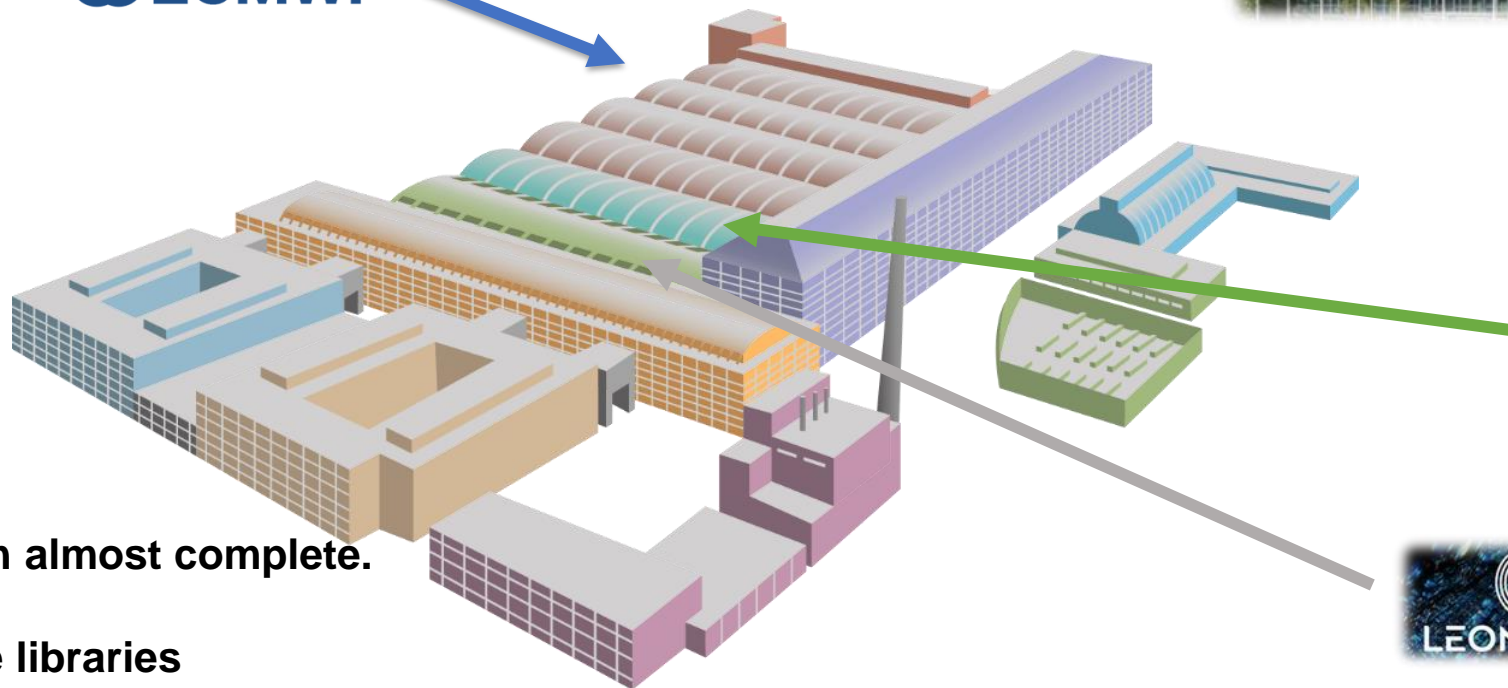


CNAF Tier1 Status Update

ATLAS Italia Computing Annual Meeting - 19/11/2024

D.Cesini – INFN-CNAF

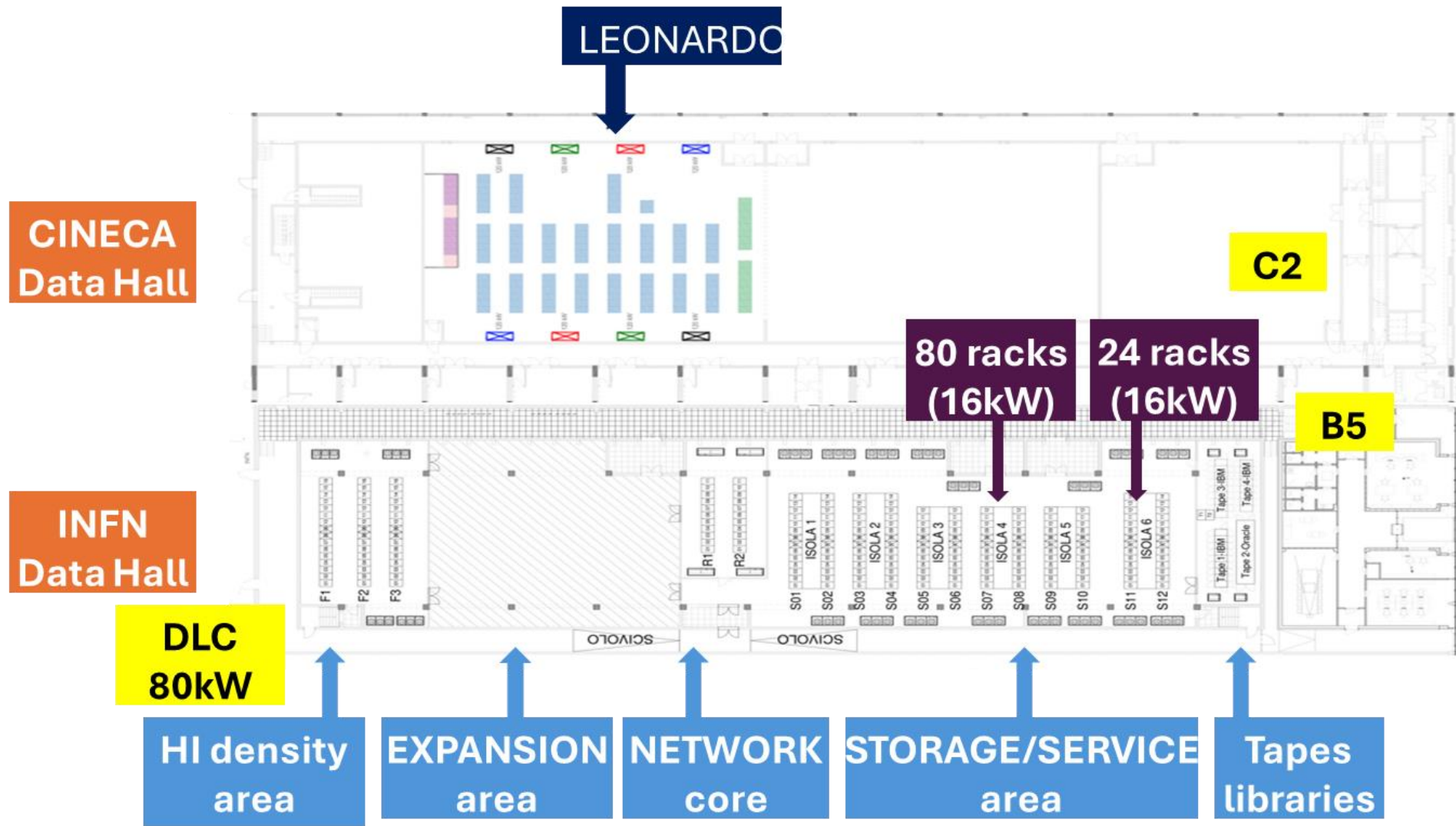
INFN-T1 new Data Center



Migration almost complete.
Missing:

- 2 tape libraries
- 1 storage system
- ISO27001 certified zone
- SSNN
- Decommissioning of the old DC

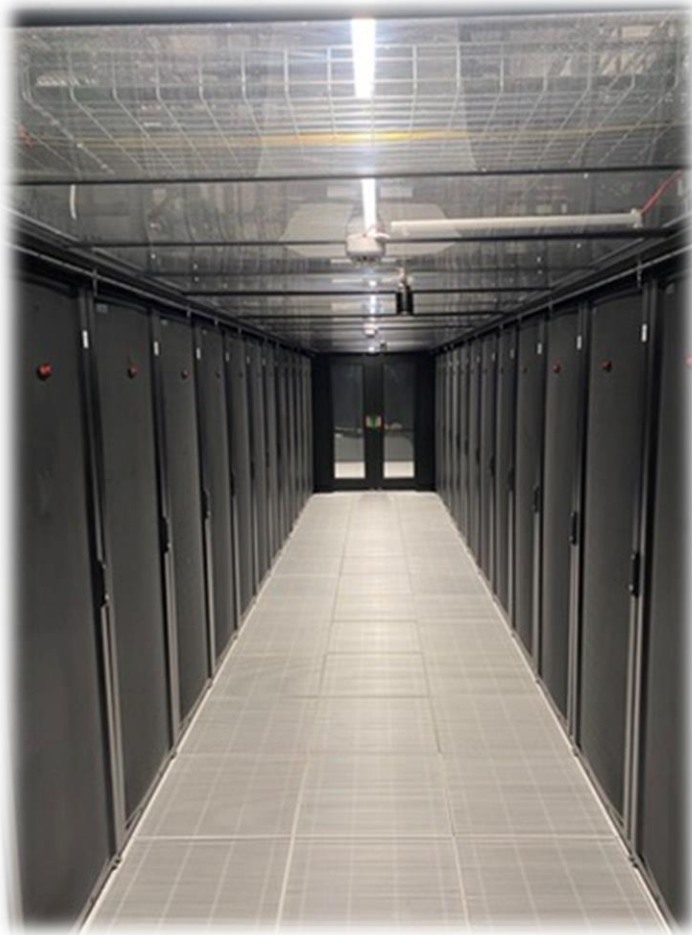
Layout of the new Data Center



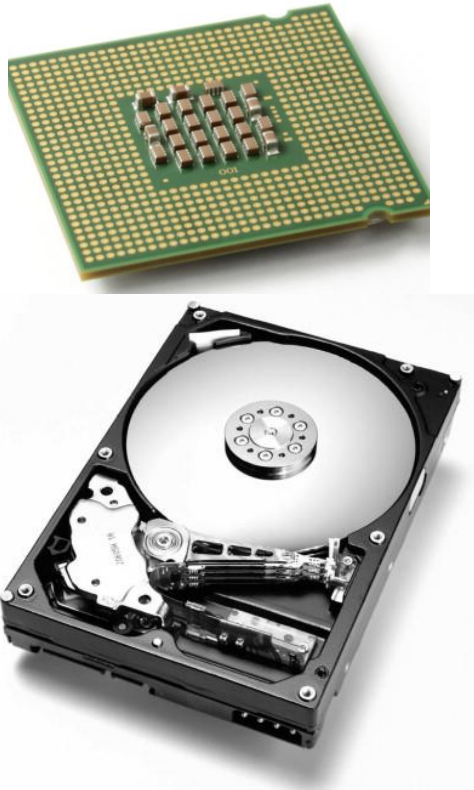
Resources:
80.000 CPU cores
80PB Disk
200PB Tape

20% year-over-year growth rate

- Air cooled Cold Corridor aisles
- DLC in **High Density**
- 3+1 redundancy in all infrastructure facilities



Available Tier1 resources



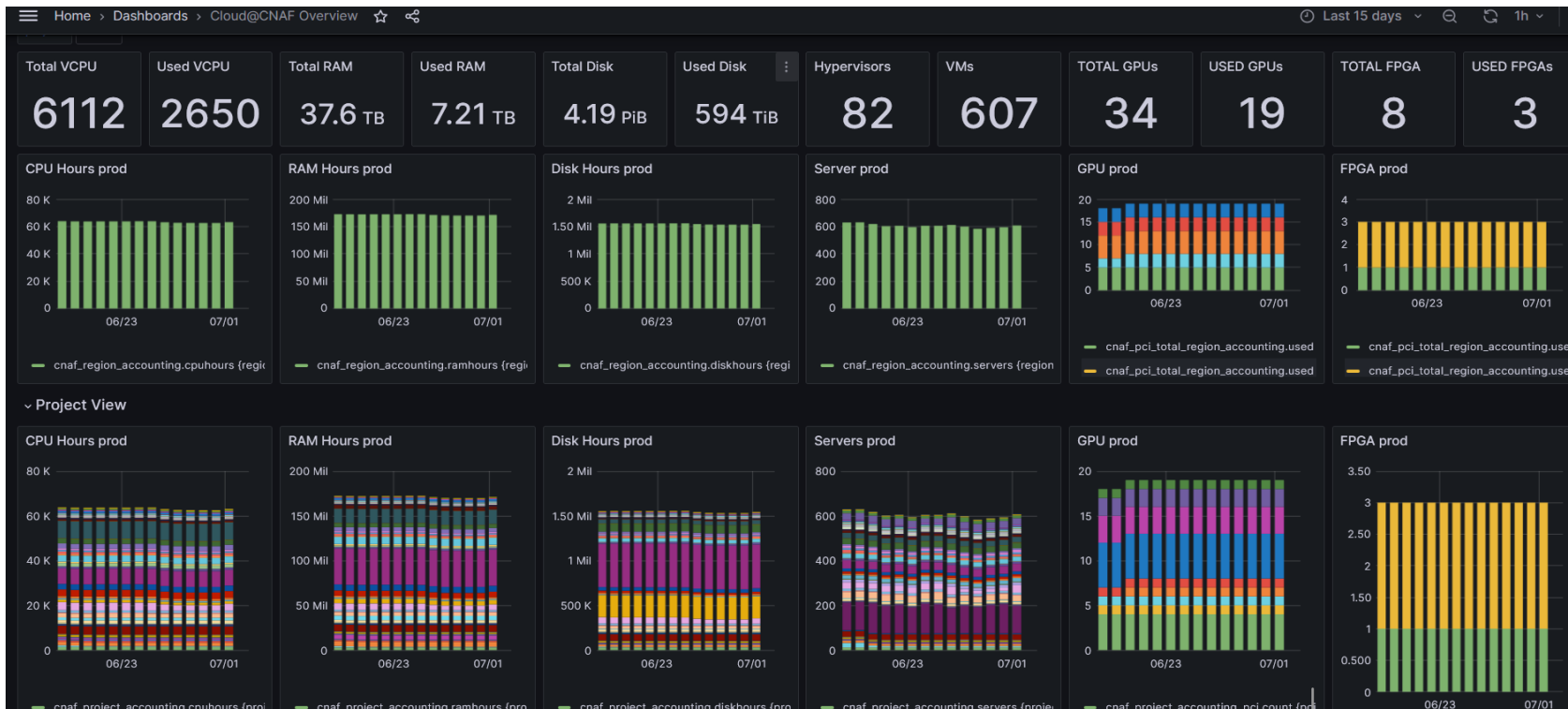
TOTAL	2024	2025
CPU Pledge (HS06)	792.000	825.000
Disk Pledge (TBN)	82.949	102.000
Tape Pledge (TB)	193.581	233.000



ATLAS@CNAF	2024	2025
CPU Pledge (HS06)	136.440	147.125
Disk Pledge (TBN)	14.670	16.740
Tape Pledge (TB)	40.680	50.490

No GPU/FPGA pledged, yet some already available via cloud.
No ARM CPU pledged, yet already available via HTCondor

- About 100 tenants configured
 - Cloud@CNAF
 - INFNCLOUD
- Big expansion thanks to Next gen EU funding (not in plot)
 - 26 nodes with 192 cores and 1.5TB RAM
 - 30 nodes with 4 H100 GPU
 - 2 FPGA “bubbles”
 - 52 storage nodes (Ceph), with obj storage, rados, filesystem
 - Installation and configuration in progress
 - Completed by the end of the year
 - Some resources will be provided via an HPC cluster

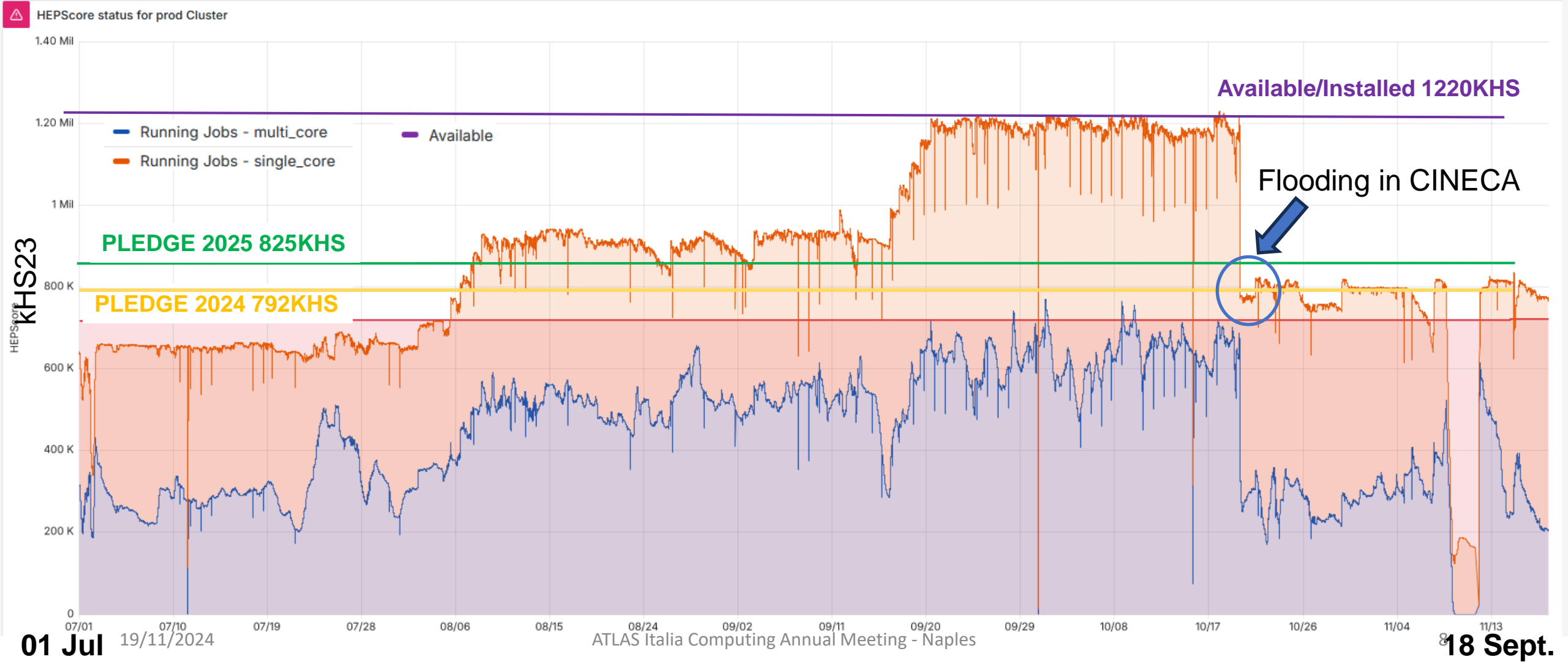


T1 CPU

CPU - Farm

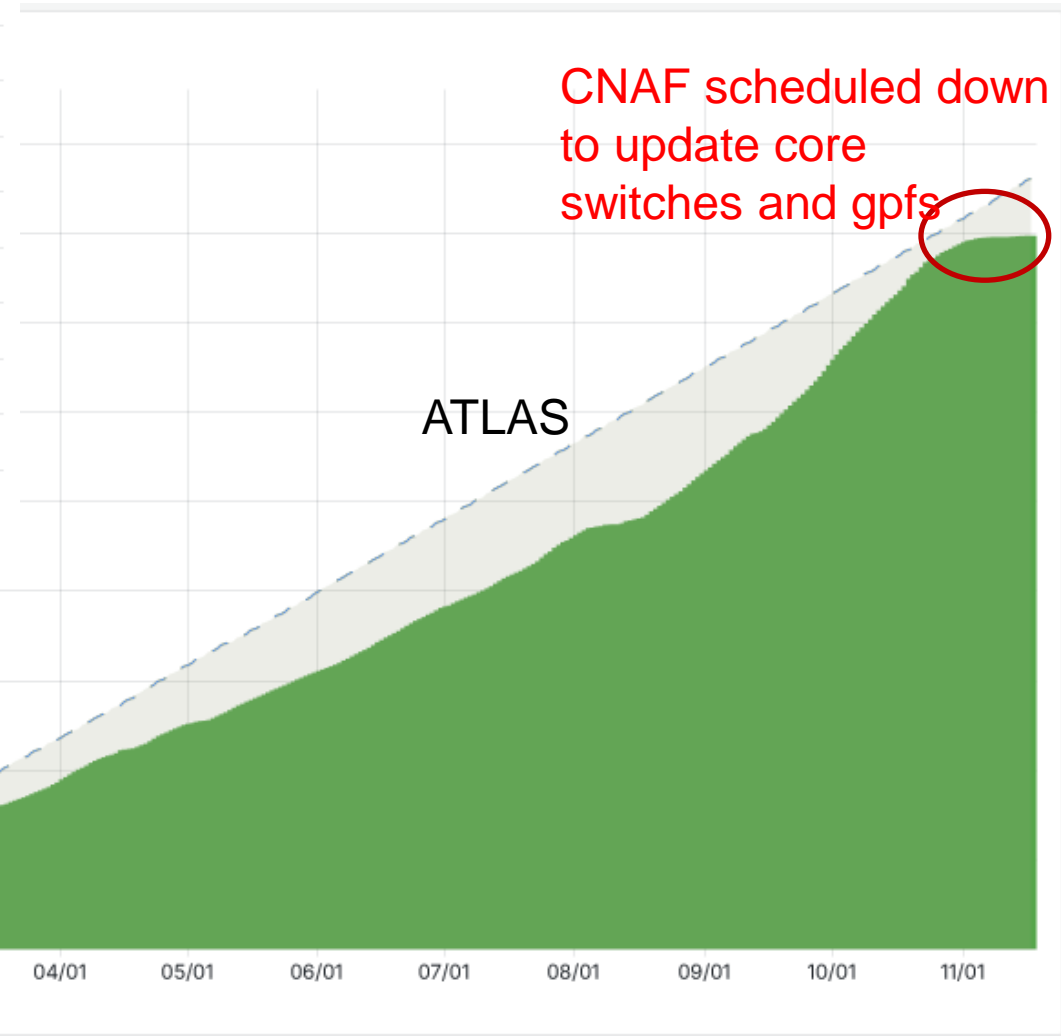
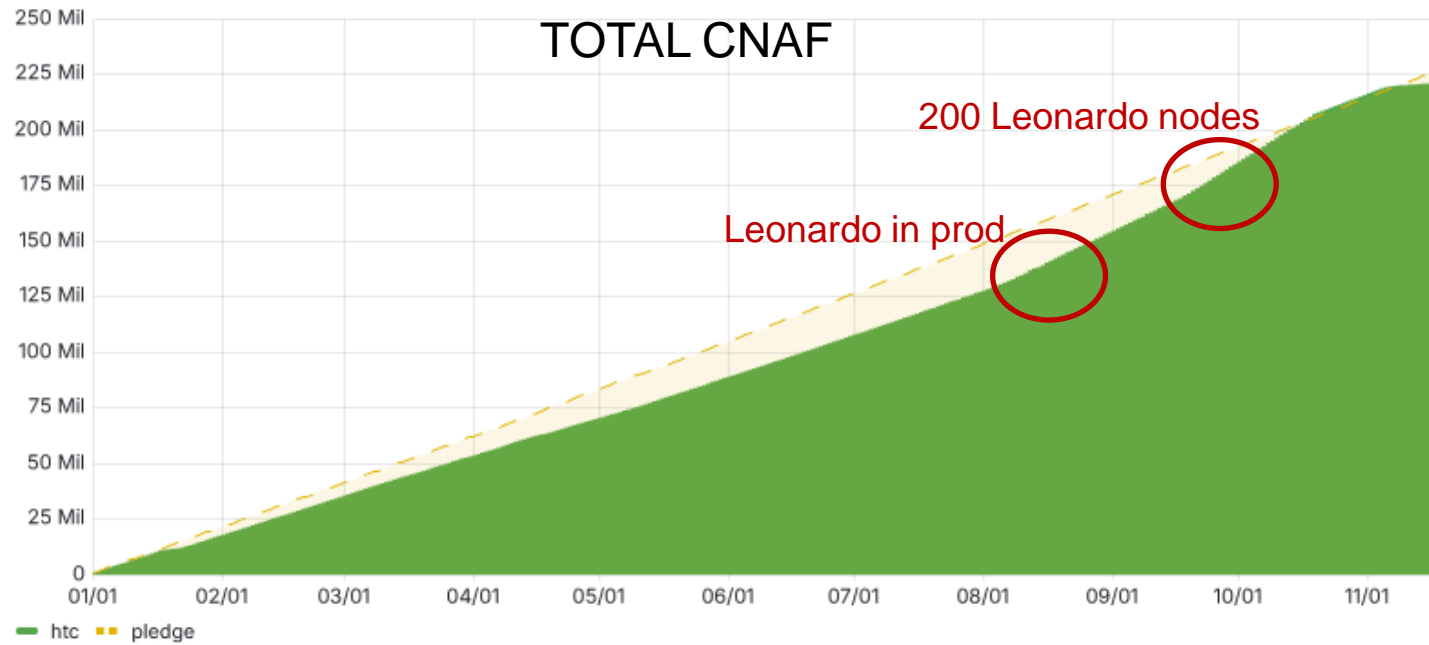
- **Pledge 2024: 792kHS06 – Pledge 2025: 825kHS => Total Installed: 1220KHS06**
- Starting from mid July Leonardo nodes gradually in production; 200 nodes, 2880HS/node

Cluster prod • **The whole farm is now HTCondor 23 and Alma9**

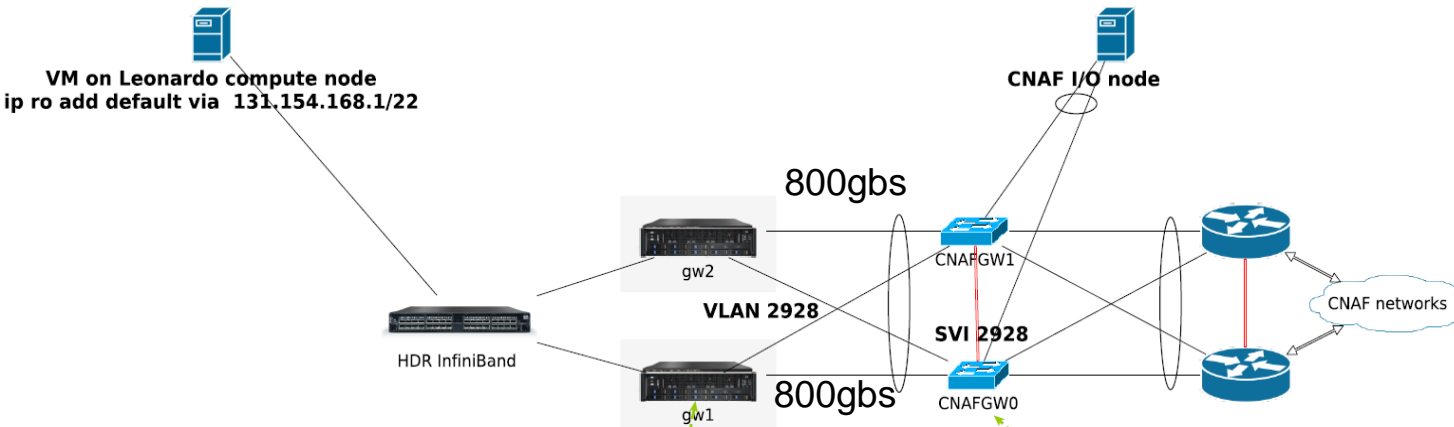


Integral HS06 – 2024

Total HS06 cumulative [HS06*day]



Leonardo GP set-up

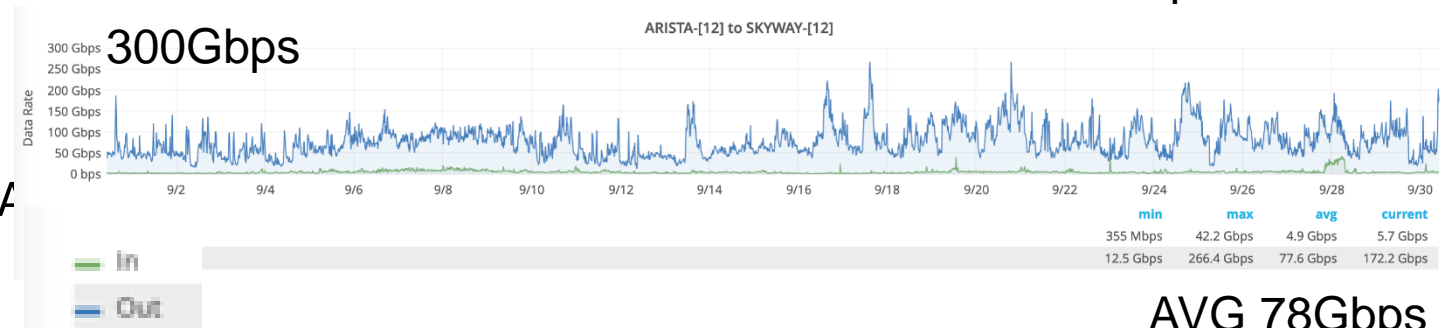


How is is going on Leonardo?

- Since July no downtimes on GP, but....
 - Next one on Nov 26th
 - Periodically we have crashes on some nodes
 - CINECA support responsive
 - IPv6 not yet supported on the SKYWAY
 - Interaction with NVIDIA ongoing

- HTCondor WN created via “infinite” CINECA SLURM whole-node jobs on the Leonardo General Partition (CPU-only)
- PCI passthrough for the infiniband NIC
- Public IP on IB
- Inbound/Outbound connectivity via NVIDIA Skyway directly attached to our core switches → 1.6Tb/s

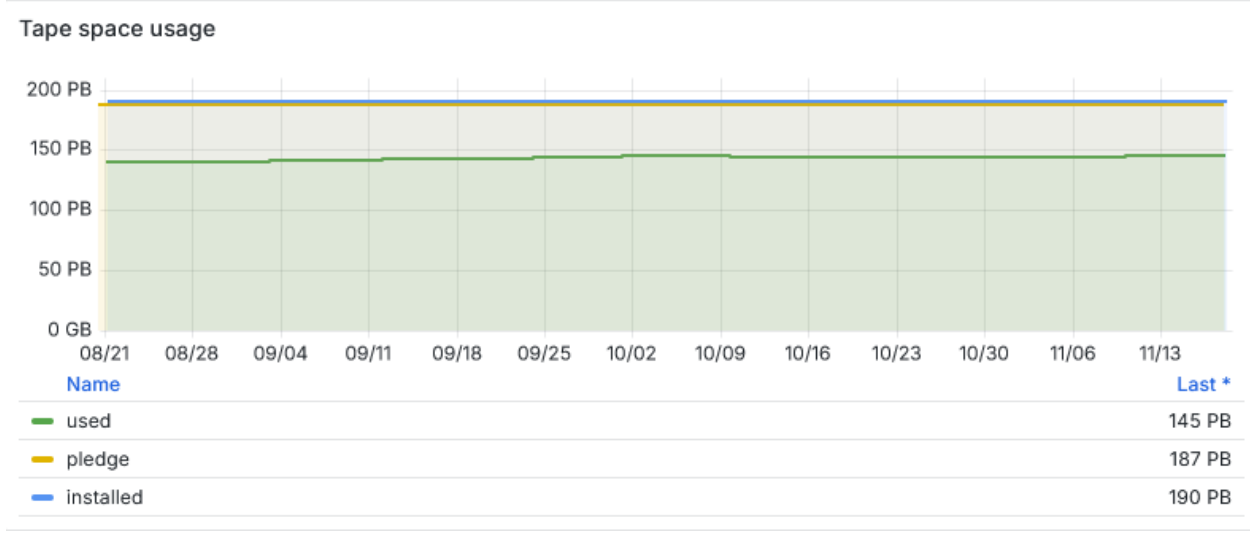
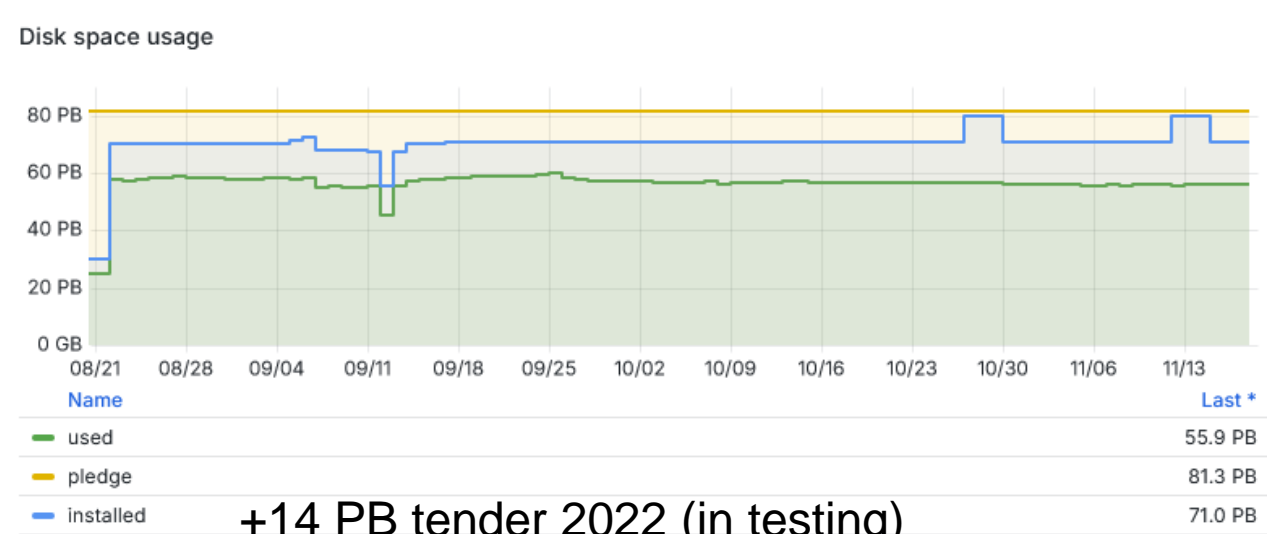
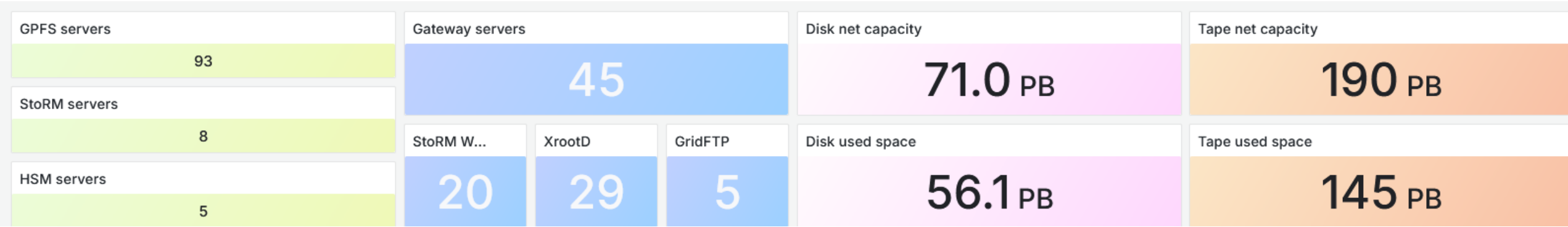
SKYWAY BW CNAF<>Leonardo in september



AVG 78Gbps

T1 DISK and TAPE

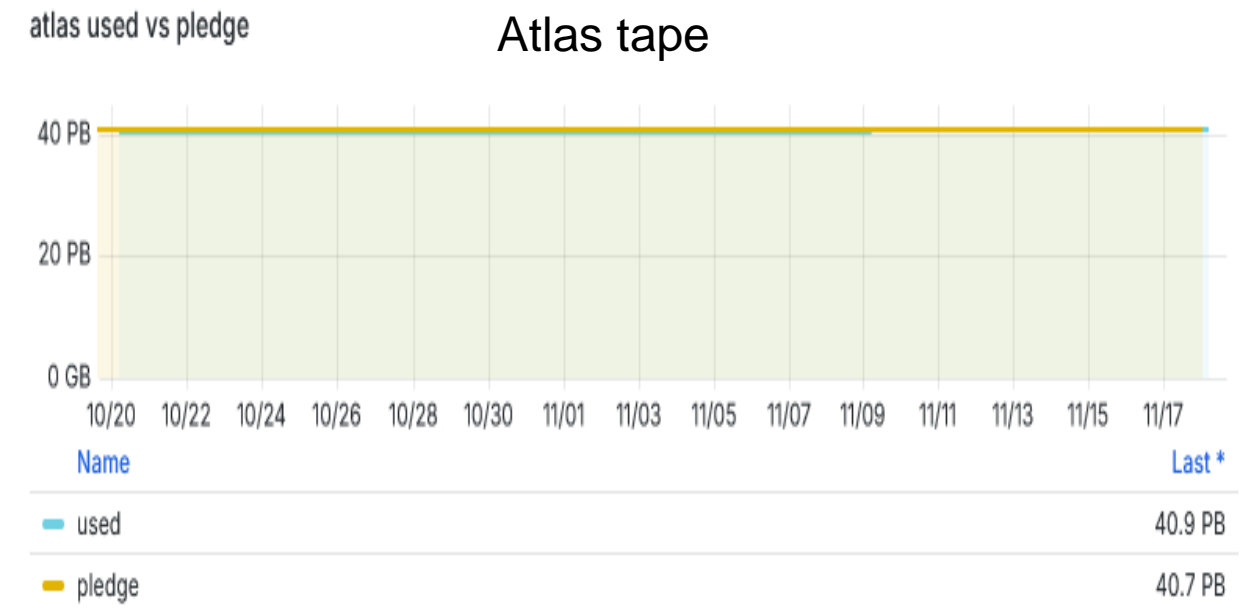
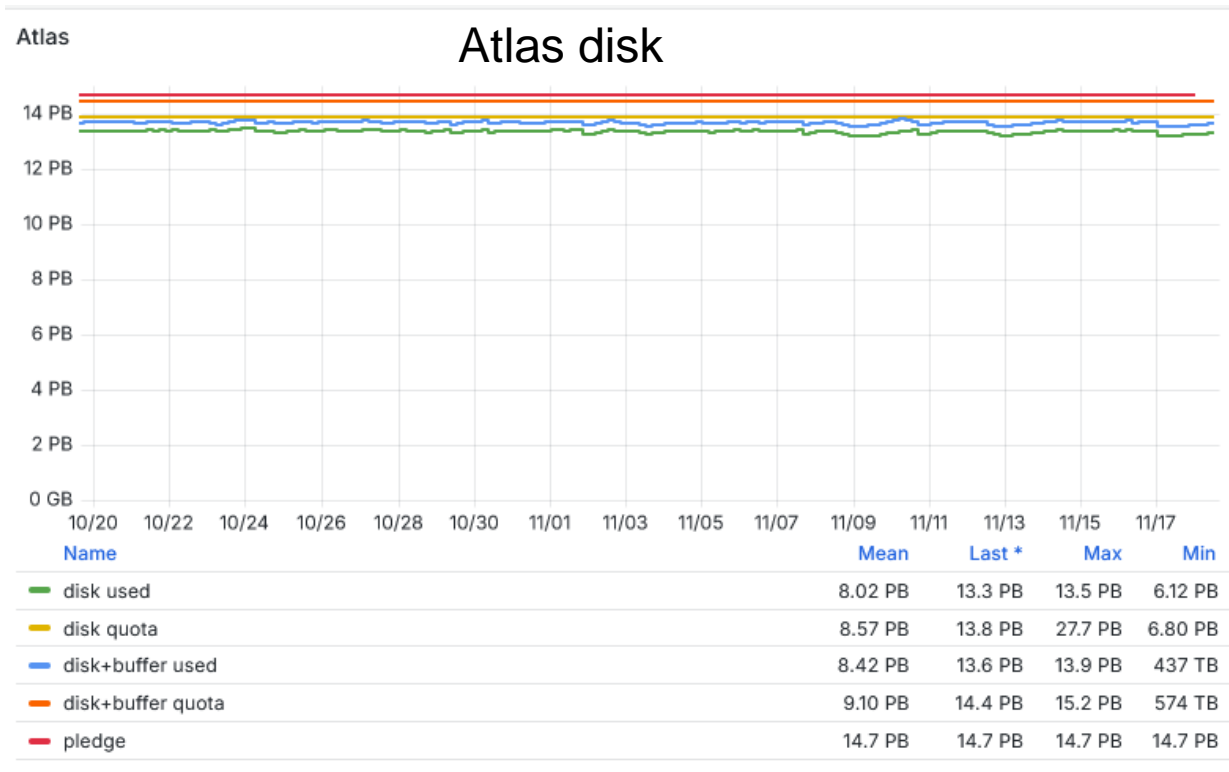
Tier1 Pledged Storage



+14 PB tender 2022 (in testing)
 +16 PB new acquisition on AQ23-24
 - 8 PB new system inefficiency

+96 PB new tender
 - 80 PB old Oracle Library decommissioning

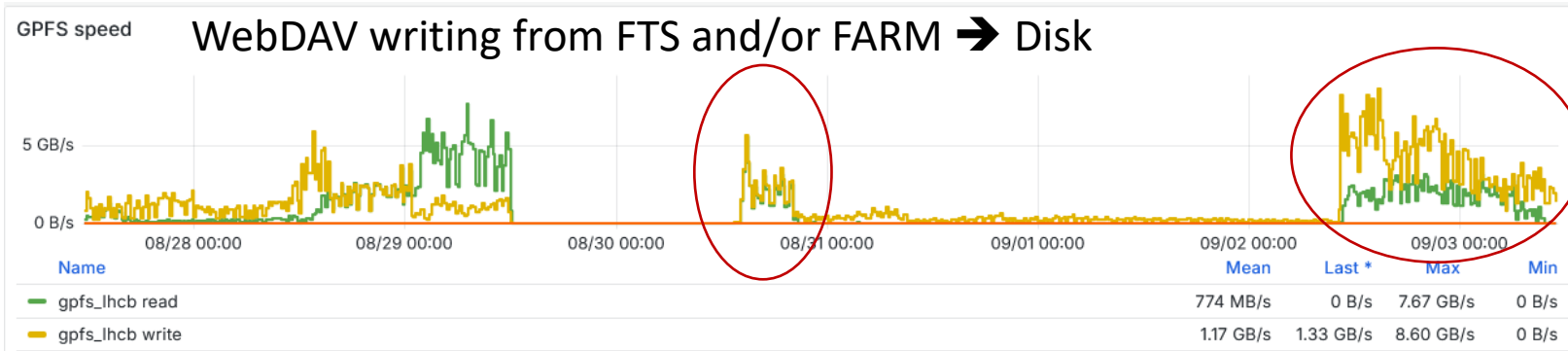
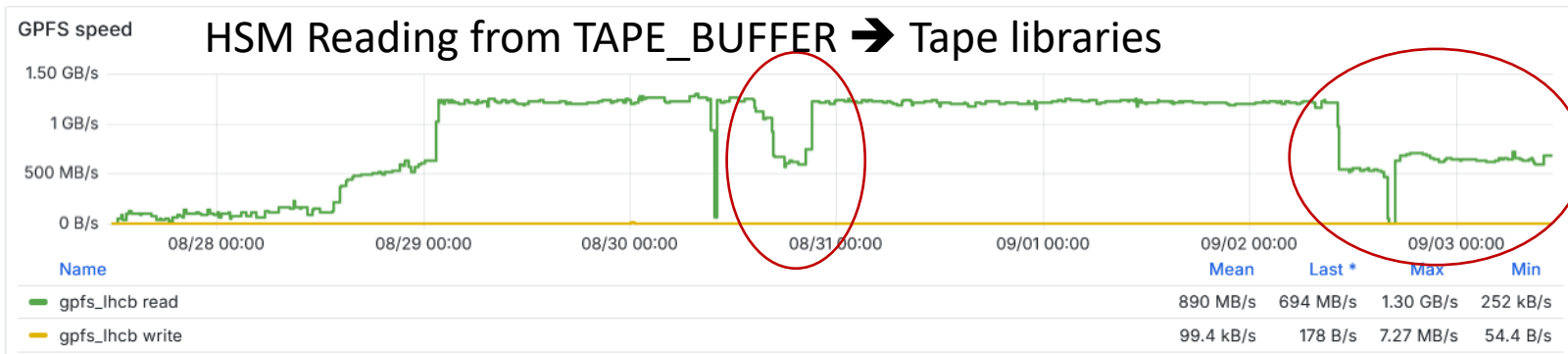
ATLAS DISK AND TAPE USAGE



ATLAS delta tape (2025-2024) → 10PB

ATLAS Delta disk (2025-2024) → 2PB

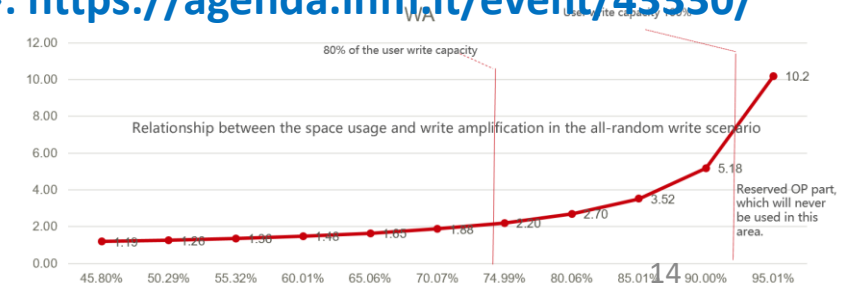
New Storage performance issue



- Mainly affected LHCb....
- ...but ATLAS is on the same system
- Huge performance drop
 - In particular for storage areas TAPE_BUFFER and DISK_BUFFER
 - Ehen the buffer is almost full
 - During cuncurrent access from FTS and from the Farm WNs (via StoRM-WebDAV)
- Storm-WebDAV gets stuck saturating the max I/O threads number
- Avalanche effect on all the StoRM-WebDAV nodes

- Long debugging phase with the vendor
 - See the post-mortem here: https://twiki.cern.ch/twiki/pub/LCG/WLCGServiceIncidents/Post-mortem_of_the_incident_on_the_CNAF_StoRM_occurred_on_August2024-1.pdf
- LHCb TAPE_BUFFER and DISK_BUFFER moved su older, but more performant, systems

See Vladimir Talk @CdG 20 Sept. 2024: «State of Storage»: <https://agenda.infn.it/event/43330/>



2024 WLCG Data Challenge

Slides from: Plenary talk di [Katy](#), brand new; [Talk](#) Track 1 di A. Forti (ATLAS); [Talk](#) Track 1 di Wissing (CMS)

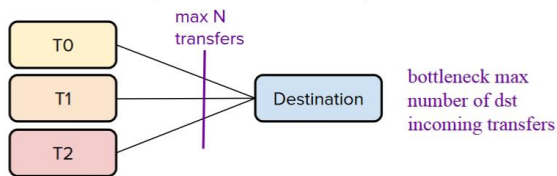
General results

- 3 large T1s had either hardware, network or MW problems
 - These problems became apparent with extra rates
 - T0 rates affected by this
- Day 8 was affected by FTS operations
- Second week was affected by the really large number of transfers

Day	Scenario	BNL-ATLAS	FZK-LCG2	IN2P3-CC	INFN-T1	NDGF-T1	pic
1	T0 → T1	25.68	N/A	29.76	N/A	35.6	N/A
2	T0 → T1	35.1	N/A	13	N/A	41	N/A
3	T0 → T1 → T1 → T2	61.6	67.1	47.4	42.2	43.8	39.3
4	T0 → T1 → T1 → T2	65.3	79.7	61.8	58.5	64.6	47.2
5	T0 → T1 → T1 → T2	83	116	81.3	78.4	75.6	56.6
6	T0 → T1 → T1 → T2	73.7	98.9	85	77.9	71.1	51
7	T0 → T1 → T1 → T2	65.7	84	79.6	102	63.6	44.6
8	T0 → T1 → T1 → T2 → T0	82.8	77.3	59.5	56.5	38.9	33.7
9	T0 → T1 → T1 → T2 → T0	87.5	80.7	51.6	63.6	40.1	34.5
10	T0 → T1 → T1 → T2 → T0	90	85.9	43.7	97.5	39.6	36.8
11	T0 → T1 → T1 → T2 → T0	110	96.8	58.8	82.1	42.1	44.6
12	T0 → T1 → T1 → T2 → T0	89.8	84.2	52.4	51.8	34	38.7

Results explained

- FTS orchestrates **transfers per link over many links**
 - Doesn't orchestrate throughput
 - To increase throughput we had to increase the number of allowed parallel transfers by an over an order of magnitude
- Has a concept of fair share per activity
 - Doesn't have a concept of links priorities within an activity, i.e. all links are equally treated T0-T1 same level as T2-T2
 - Could prioritise faster transfers or more important channels
- Testing also new authz system with tokens put further load on the system



- Agreed in 2021 FTS problems to solve first for next challenge

Comparison: Rates achieved vs. targeted – Tier 1s

Day	Scenario	JINR		FNAL		IN2P3		RAL		PIC		KIT		CNAF	
		DEST	SRC	DEST	SRC	DEST	SRC	DEST	SRC	DEST	SRC	DEST	SRC	DEST	SRC
1	T0 Export	1.42	N/A	1.13	N/A	1.09	N/A	0.76	N/A	1.18	N/A	1.16	N/A	1.17	N/A
2	T0 Export	1.46	N/A	1.12	N/A	1.10	N/A	0.50	N/A	1.17	N/A	0.94	N/A	1.17	N/A
3	T0Export, T1Export	1.31	0.62	1.08	0.88	1.33	1.03	0.72	0					1.28	0.93
4	T1 Export	N/A	0.37	N/A	0.91	N/A	1.12	N/A	0					N/A	1.00
5	T1-Export, Prod-out	1.18	1.72	1.15	0.87	1.25	0.89	0.98	1.01	1.21	1.09	1.23	0.77	1.17	0.77
6	T1-Export, Prod-out	1.14	2.42	1.18	0.88	1.47	0.88	0.72	0.81	1.17	1.03	1.19	0.76	1.18	0.95
7	T1-Export, Prod-out	1.19	2.19	1.15	0.87	1.22	0.87	0.81	1.04	1.20	0.98	1.21	0.73	1.16	1.02
8	AAA	1.30	N/A	N/A	1.10	1.39	N/A	1.31	N/A	1.31	N/A	1.70	N/A	1.32	N/A
9	All	0.38	0.34	0.87	0.84	0.57	0.57	0.5					0.56	0.65	0.25
10	All	0.70	0.34	0.98	0.74	0.58	0.65	0.56		0.70	0.6	0.03	0.98	0.63	0.28
11	All	0.63	0.33	0.91	0.73	0.43	0.76	0.77	1.05	1.09	0.84	0.91	1.09	0.69	0.24
12	All	0.40	0.54	0.92	0.86	0.89	1.00	0.85	1.15	1.21	0.87	1.13	0.89	0.78	0.29

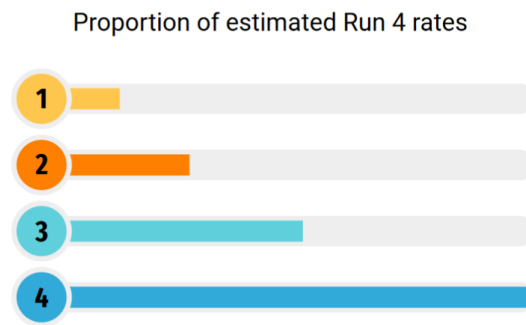
Ratio = observed/targeted
 Green - ratio > 0.9
 Yellow - 0.9 > ratio > 0.7
 Orange - 0.7 > ratio > 0.5
 Red - ratio < 0.5

Typically storage overloaded

Tuning of FTS

Typically storage overloaded

Data challenge series



- ✓ 10% 2021
- ✓ 25% 2024
- 50% ? / 2027
- 100% ? / 2029

Intermediate mini data challenges focused on sites, technologies

Run 4 is scheduled to start in 2030

Metropolitan Tape Area Network



- 2 libraries at CNAF
- 1 new library at the Tecnopole
- About **7 km** of fiber to connect the 2 datacenters
 - **yellow + red** paths
- 2 fiber pairs dedicated to extend the fiberchannel TAN
 - Brocade optics for 10km distance

BROCADE
A Broadcom Company

Product Brief

Brocade® 32Gb/s LWL
(10 km) SFP+

Optimized, Certified Optical Transceivers for
Extending Service Provider and Data Center
Networks

Overview

Today's enterprise data centers are undergoing an infrastructure transformation, requiring higher speeds, greater scalability, and higher levels of performance and reliability to better meet the demands of business. As speed and performance needs increase, optical transceivers—once considered a generic component of Fibre Channel switching technologies—have become an integral part of overall system design.

The Brocade® 32Gb/s Long Wavelength (LWL) 10 km SFP+, part of the

Highlights

- Provides high system reliability through rigorous qualification and certification processes.
- Leverages unique design parameters to provide the highest performance with industry-leading Brocade

CNAF Tape libraries and drives

- **1 x Oracle SL8500**

- **1 tape library with 16 tape drives T10000D**
(8.5TB/cartridge)
- 80PB installed, 64PB USED
- **Repack on the other libraries ongoing**
 - After completion of repack this library will be dismissed
- **Due to time constraints for decommissioning we move this to the new location (December 2024)**



- **2 x IBM TS4500**

- **1 tape library with 19 tape drives TS1160** (20TB/cartridge)
 - 102 PB Installed, 50PB USED
 - cannot be further extended due to physical constraints in the current room
 - **This library moved to the new data center (November 2024)**
- **1 tape library with 18 tape drives TS1170** (50TB/cartridge) acquired and installed at new data center Q2 2024

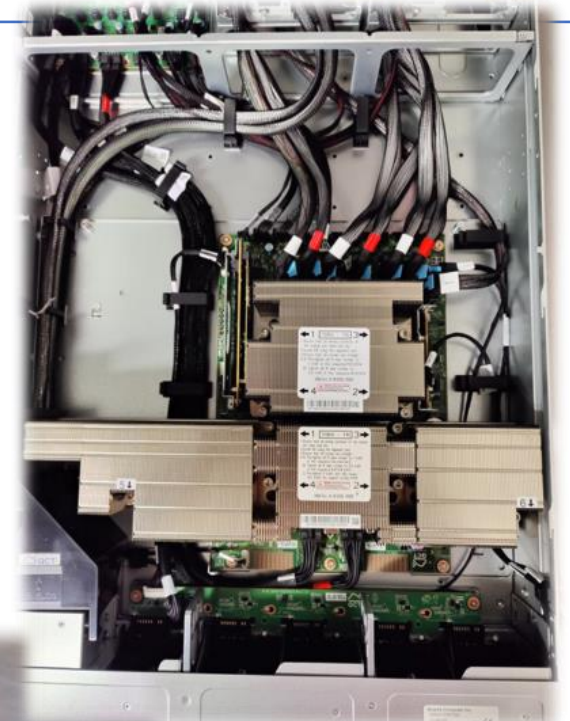


Heterogeneous computing

- Started a **technology tracking** program (C3SN WG), mainly focused on computing
 - Investigate new processors technologies
 - Power consumption comparisons
 - CPU architectures
 - Understand middleware and general **software readiness**

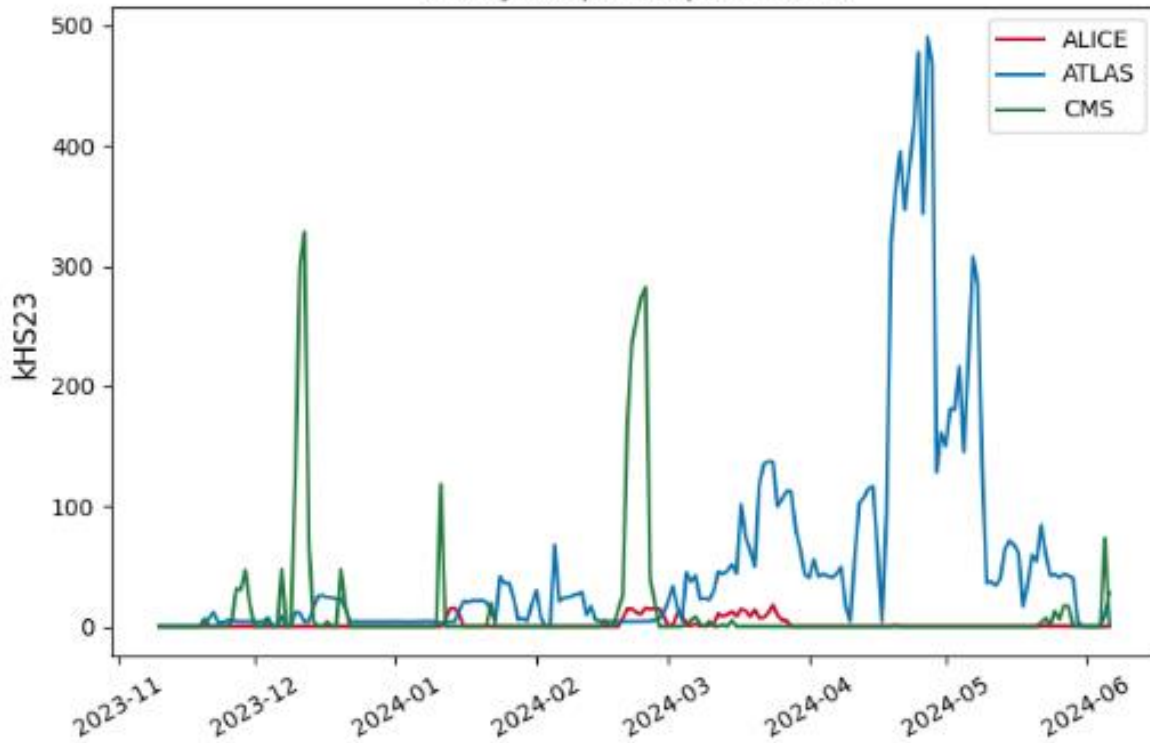


See our presentation at last CHEP!

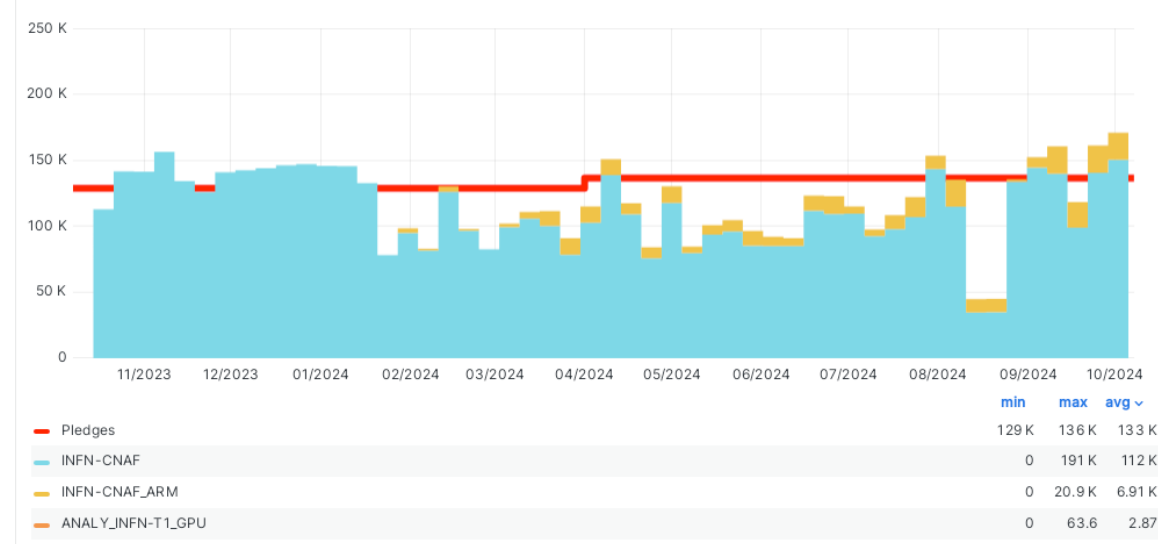


ARM CPU usage at INFN-T1

ARM jobs per experiments



Slots of Running jobs (HS23)



TICKETS RELVALS DASHBOARD

Logged in as Daniele Spiga

Workflows (Jobs in ReqMgr2)

- pdmvserv_RVCMSSW_14_0_0_pre3RunJetMET2023D_CNAFARM_ReVal_2023D_240215_092819_7275 open in: Stats2 status: normal-archived

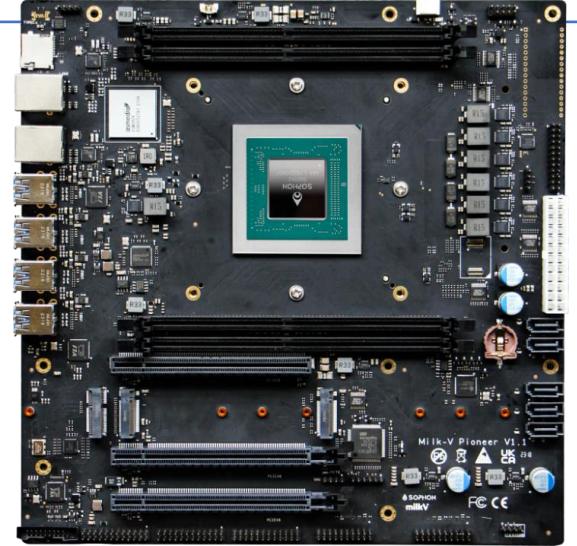
 - datatier: FEVTDEBUGHLT, completed: 109.72%, events: 1,266,761, type: **VALID**
 - datatier: AOD, completed: 106.33%, events: 1,227,653, type: **VALID**
 - datatier: MINIAOD, completed: 106.33%, events: 1,227,653, type: **VALID**
 - datatier: NANOAO, completed: 106.33%, events: 1,227,653, type: **VALID**
 - datatier: DQMIO, completed: 0.00%, events: 0, type: **VALID**
- pdmvserv_RVCMSSW_14_0_0_pre3RunDisplacedJet2023D_CNAFARM_ReVal_2023D_240215_092843_9419 open in: Stats2 status: normal-archived

 - datatier: FEVTDEBUGHLT, completed: 123.17%, events: 788,261, type: **VALID**
 - datatier: AOD, completed: 120.25%, events: 769,596, type: **VALID**
 - datatier: MINIAOD, completed: 119.45%, events: 764,476, type: **VALID**
 - datatier: NANOAO, completed: 120.25%, events: 769,596, type: **VALID**
 - datatier: DQMIO, completed: 0.00%, events: 0, type: **VALID**

CMSSW_14_0_0_pre3_Data_2023_CNAFARM-RunDisplacedJet2023D-00001

RISC-V CPU usage at INFN-T1

- 2x Milk-V Pioneer
 - 64 Core RISC-V CPU up to 2GHz
 - 128GB DDR4 3200
 - 1TB NVMe disk
 - 2x10Gbps network cards
 - Pre-installed OS: Fedora 38



CMS Point-of-view

- db12 CMS benchmark
 - milk-v: 378.3, **5.8 per core**
 - e5-2640v3: 248, **3.8 per core**
- Milk-v today performs better than a 2016 CPU...
- If we take out the hyperthreading, milk-v is 2-3 times slower than a modern xeon
- **64 real cores**, not so common on modern CPUs
- Trade-off between core number and power (taking into account power consumption too)



HPC Bubbles



Nodo CPU

192 core fisici
1.5TB RAM DDR5
IB NDR 400G
20TBL (SSD) + dischi di sistema



Nodo GPU

Come CPU + 4x NVIDIA H100 SXM5 con minimo 80GB e memoria HBM2e



Nodo FPGA

32core
RAM 768GB DDR5
IB NDR 440G
4 x XILINX U55C o 4 x TerasicP0701



Nodo Storage (CEPH Bricks)

64 core fisici
1TB RAM DDR5
384 TBL HDD + 25.6 TBL NVMe



Accessori

Switch IB, Switch ETH
Cavi IB, Cavi ETH
Transceiver vari
Assistenza 3+2



“HPC Bubbles”

(*) NODES in EPIC – ISO certified zone

	Nodo CPU	Nodo GPU	Nodo FPGA Xilinx	Nodo FPGA Terasic	Nodo storage
BA_DARE	12	0	0	0	6
BA_TerabitS8	0	0	0	0	0
CNAF_DARE	10	9	0	0	16
CNAF_TerabitS8	0-8?	0-8?	0	0	0-6?

TOTAL NODES funded by Terabit-ICSC-DARE

	Nodo CPU	Nodo GPU	Nodo FPGA Xilinx	Nodo FPGA Terasic	Nodo storage
BA	24*	6	0	0	32 *
CNAF	26*	30 *	2	2	52 *
MIB	0	0	2	2	0
NA	18	1	2	0	8
PD	6	6	0	0	0
PI	20	0	0	0	0
RM1	12	0	0	0	0
TO	14	6	0	0	0
LNGS	0	6	0	0	12
CT	12	0	0	0	8
LNF	12	0	0	0	0
LNFESA	8	6	0	0	6
LNL	4	0	0	0	0
MI	4	0	0	0	0
TOTALE	160	61	6	4	118

Core: 30 kcore fisici
Circa 34 HS/core

GPU: 244 NVIDIA H100
40 FPGA
InfiniBAnd 400Gbs

45 PB RAW