

# CNAF: stato Tecnopolo, Pledge 2024, Esito Gare, CINECA

D. Cesini – INFN-CNAF

# Resources@T1 2024-2025



ALL VO No Cloud	2024	Delta 2025-2024 (Preventivi)	2025 (preliminary)
Pledge CPU (HS06)	<b>792000 (plan)</b> 703000(with OF) 844000 (w/o OF)	CSN1-LHC: 127210	<b>1022780 (no OF)</b>
		CSN1-NoLHC: 8000	
		CSN2: 12970	
		CSN3-Alice: 22500	
		CSN3-NoLHC: 8100	
		<b>TOTALE: 178780</b>	
Pledge disk (TBN)	<b>82949</b>	CSN1-LHC: 13350	<b>103627</b>
		CSN1-NoLHC: 450	
		CSN2: 4066	
		CSN3-Alice: 2700	
		CSN3-NoLHC: 112	
		<b>TOTALE: 20678</b>	
Pledge tape (TB)	<b>193581</b>	CSN1-LHC: 45080	<b>254026</b>
		CSN1-NoLHC: 1060	
		CSN2: 5750	
		CSN3-Alice: 6990	
		CSN3-NoLHC: 65	
		<b>TOTALE: 60445</b>	

+ cloud e HPC da CSN2,3,5

# Richieste 2025 – Cloud e HPC

[DarkSide] CPU-Cloud: 1000 HS06	10.00
[Limadou] CPU-Cloud: 200 HS06	2.00
[QUAX] CPU-Cloud: 300 HS06	3.00
[DarkSide] DISK-Cloud: 100 TB	10.00
[Limadou] DISK-Cloud: 10 TB	1.00
[Auger - DataCenter] GPU: 8640 GPU hours	6.00
[CYGNO] GPU: 2400 GPU hours	1.50
[Limadou] GPU: 2037 GPU hours	1.50
[XRO] GPU: 4320 GPU hours	3.00
[ET] HPC: 100k core*hours	0.50
[Euclid] HPC: 7M core*hours	26.00
[LSPE] HPC: 500k core*hours	2.00
[QUBIC] HPC: 200k core*hours	0.50
SPES-MED - CPU 3000HS06 (Cloud?)	30.00

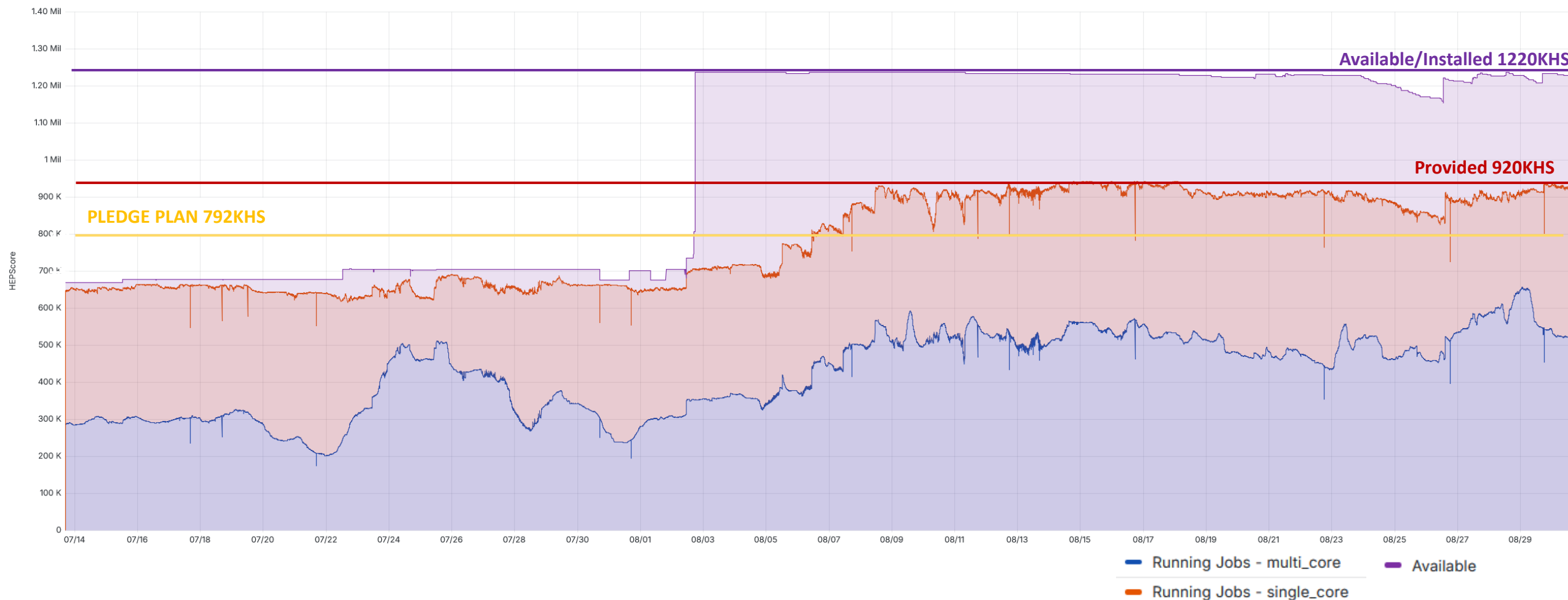
ARTEMIS - GPU: 2000 GPU-hours GPU A100 (> 32 GB di VRAM) per training modelli AI	1.50
AIM_MIA - GPU: 9000 GPU-hours suddivise su 2x GPU A100 (> 40GB di VRAM) per training modelli AI	6.50
Geant4INFN - GPU: 10000 GPU-hours tramite piattaforma AI_INFN	7.00
GEANT4INFN - CPU: 120k core-hour per simulazioni geant4_DNA richiesta una macchina con almeno 128 GB di RAM e accesso tramite coda condor su tier1	1.50
GEANT4INFN - Storage 2 TB accessibile posix da infn cloud	0.50
BRAINSTAIN - CPU: 400k core-hours	4.50
BRAINSTAIN - GPU: 10k GPU-hours	7.00
BRAINSTAIN - Storage: 8 TB	1.00
FRIDA - CPU: 1M core-hours (2GB/core) per attivita' Roma1: 200k core-hours (6GB/core) per attivita' Roma3-TIFPA	13.50
FRIDA - GPU: 1440 GPU-hours (GPU NVIDIA A100 con VRAM ~ 10 - 20 GB) per simulaz. FAST MC	1.00
SPRITZ - CPU: 125000 core-hours per simulazioni multithread (40-80 cores) con shared-memory (256-512 GB).	1.50
SPRITZ - Storage: 150 TB con protocolli di accesso standard basati su ssh	15.00
PLASMA4BEAM2 - Storage: 100 TB	10.00
PLASMA4BEAM2 - CPU: 500000 core-hours	6.00
SEGNAR - CPU: 100000 core-hours	1.00
FUSION - CPU: 240000 core-hour	3.00
FUSION - GPU: 20000 GPU-hours	14.00
FUSION - Storage: 1 TB	0.50
VITA - CPU: 1400000 core-hours	16.00
VITA - GPU: 3000 GPU-hour	2.00
VITA - Storage: 700 TB (una parte per dati sensibili)	70.00
MIRO - CPU: 1.2M core-hours [500k core-hours per simulaz. MC (Catania); 500k core-hours per simulaz. dinamica molecolare (Pisa); 200k core-hours per simulaz. MC e analitiche multiscala (TIFPA)]	13.50
MIRO - GPU: 2200 GPU-hours GPU NVIDIA A100 (> 32 GB di VRAM) per training modelli AI	1.50
SEGNAR - Storage: 10 TB	1.00

# T1 CPU

# CPU - Farm

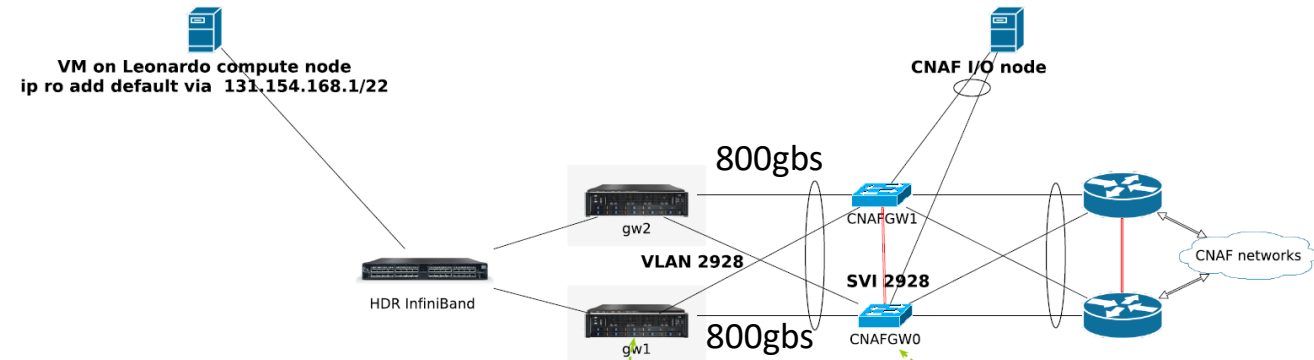
- Pledge 2024: 703kHS06 (w/o OVERLAP- 843k) → **CNAF TOTAL PLAN 792kHS06 - Potenza installata Totale: 1220KHS06**
- Da metà luglio abbiamo iniziato ad inserire i nodi Leonardo – 200 inseriti di cui 100 in draining (2880HS/nodo)
  - Per aggiungere altri nodi aspettiamo di potenziare il link da 200Gbit/s a 1600Gbit/s – in arrivo le ottiche , ordine effettuato a Luglio
- **Tutta la Farm ad HTCondor 23 – dal 20/09 partiamo con aggiornamento a Alma9+IPv6 sui WN**

HEPScore status for prod Cluster



# Set-up Leonardo GP

- Strategia che abbiamo implementando:
  - WN creati tramite job SLURM “infiniti” che instanziano machine virtuali WholeNode gestite da noi
    - Immagini VM create da noi e disponibili via shared fs
  - PCI passthrough per scheda infiniband
  - IP pubblico su interfaccia infiniband
  - Accesso inbound/outbound via NVIDIA Skyway collegati direttamente a nostri apparati

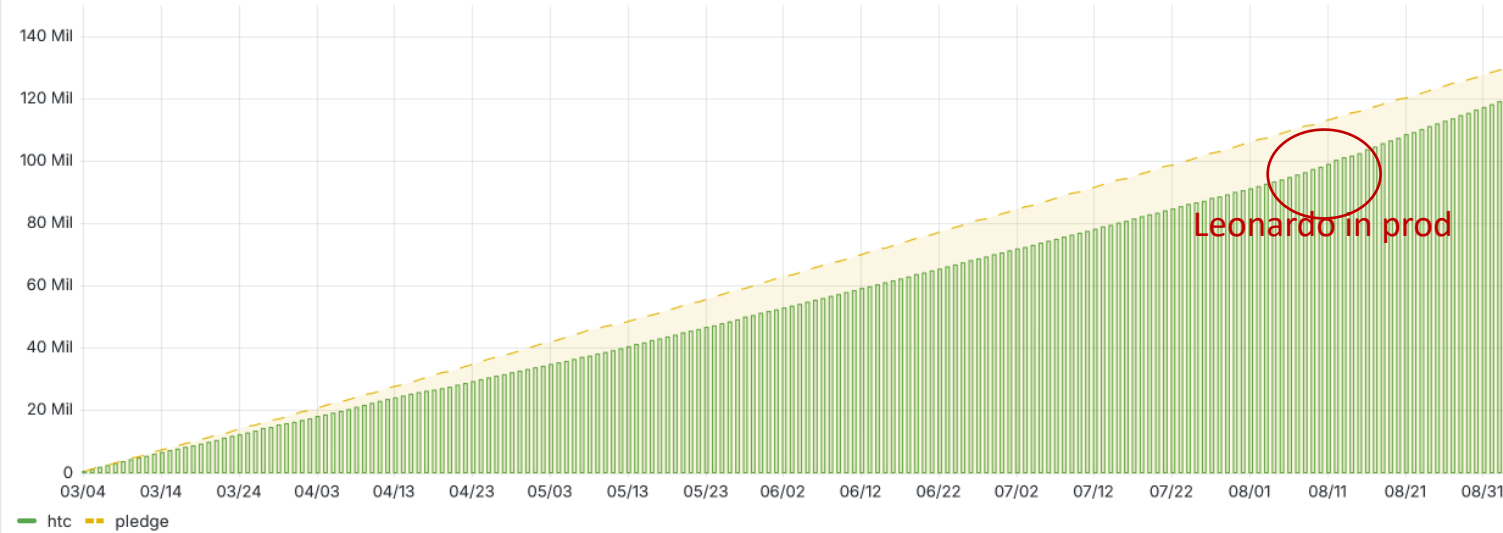


- **Come sta andando la produzione su Leonardo?**
  - Da Luglio non ci sono più stati down globali della GP, ma....
    - In media 6-7 machine al giorno si bloccano
      - Nodi random
      - Per cause sconosciute, sotto investigazione
    - Circa 1000 job/day che perdiamo senza motivo

# HS06 Integrati – ultimi 6 mesi

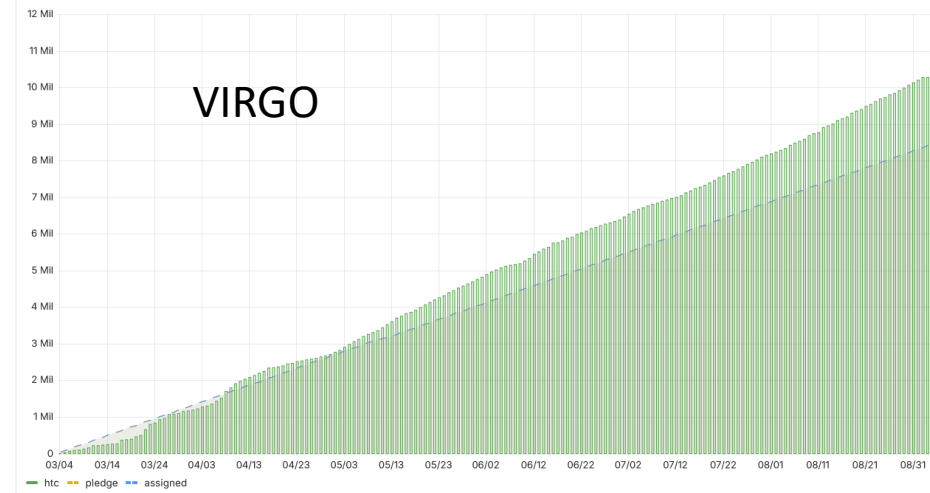
Total HS06 cumulative [HS06\*day]

TOTAL

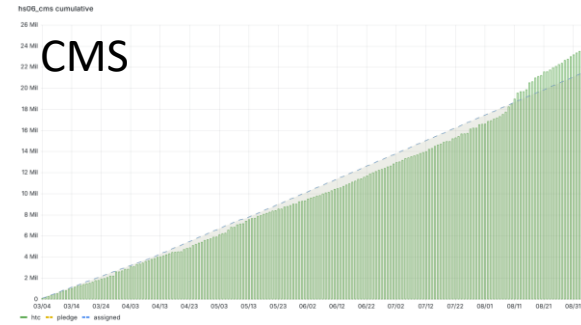


hs06\_virgo cumulative

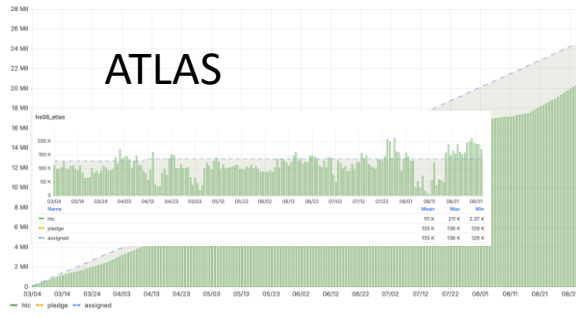
VIRGO



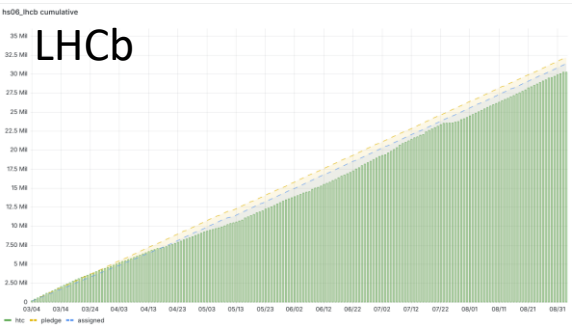
CMS



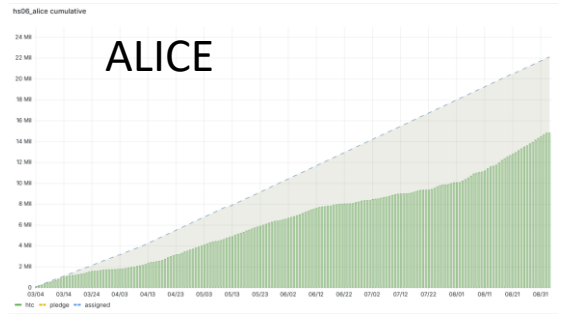
ATLAS



LHCb

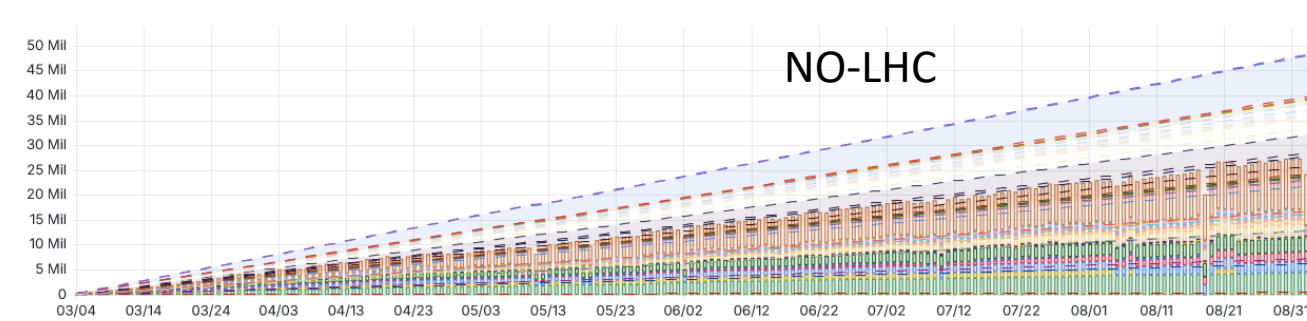


ALICE



HS06\_per\_group cumulative

NO-LHC

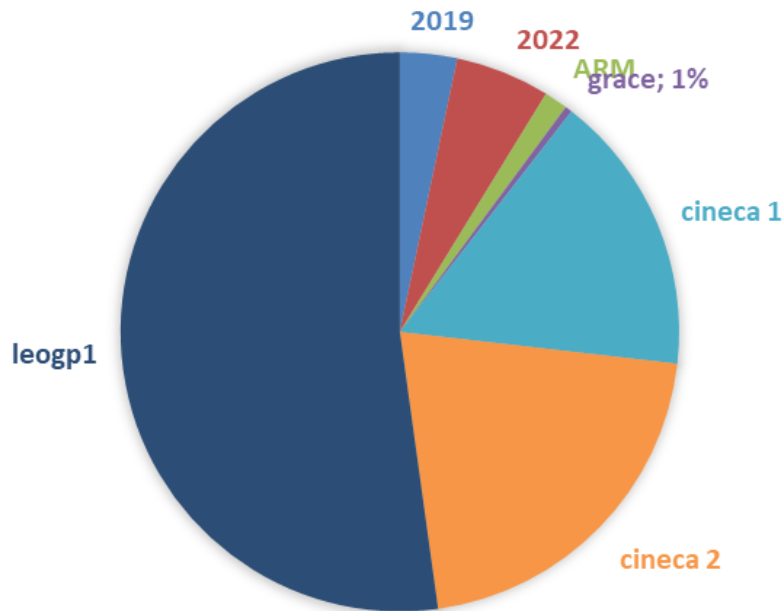


04/09/2024

CNAF: stato Tecnopolo, Pledge, Esito gare, CINECA

# Composizione della farm vs Pledge 2025

FARM POWER PER TENDER



- Considerando 1.22MHS
  - Il 33% è installato a CINECA-Casalecchio
    - 411kHS
  - Il 46% è su CINECA-Leonardo
  - Il 9% va dismesso (gare 2015-2016)
    - 110 kHS
  - 12% al CNAF(TP)
    - 1% è su ARM
- Aggiunti i nodi ARM
  - Ampere (4 nodi)
    - 3754HS06/nodo
    - 3.74 hs06/W (vs 2.64 HS06/W gara2022)
  - Grace (1 nodo)
    - 4459 HS06/nodo
    - 4.67 HS06/W (vs 2.64 HS06/W gara2022)

} 79% su CINECA

Gara 2022, 2019, 2017 portate al tecnopolo  
2015 e 2016 speriamo di poterle dismettere

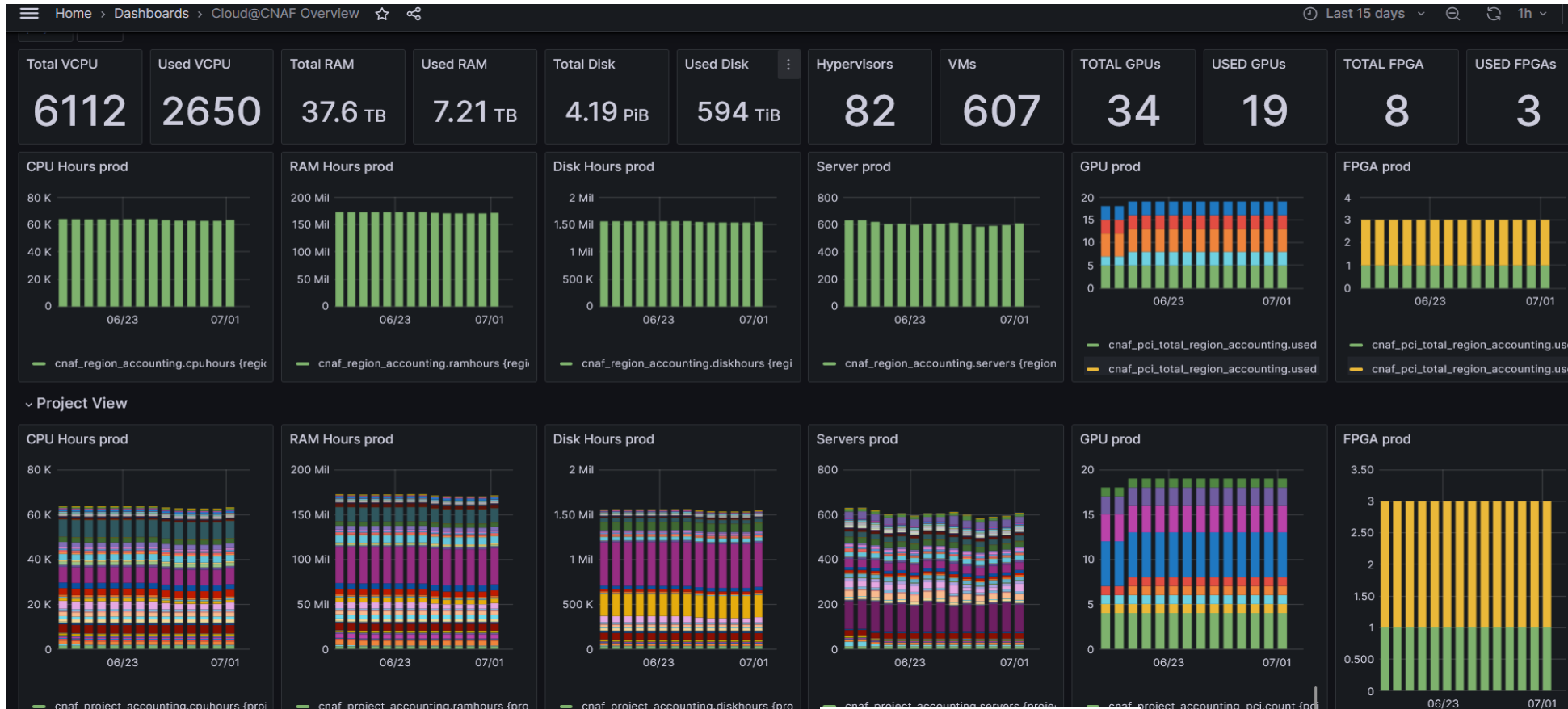
**1220kHS – 411kHS(casalecchio) – 110kHS(dismisisoni) → 700kHS**  
**Pledge stimato 2025 → 970kHS**

Bubble CPU-Only non cosiderate

- Equivalenti a:
  - 270kHS da acquisire
  - 270KHS da non spegnere a Casalecchio (o spostare al tecnopolo)
  - altri 100 nodi Leonardo



# Stato Cloud@CNAF



- Parte dei pledge «HTC» assegnati su Cloud@CNAF per accesso interattivo o piccoli cluster dedicati
  - AGATA, NTOF, etc..
- VIRGO Low Latency Cluster on K8s

- Circa 100 tenant configurati
  - Cloud@CNAF
  - INFN CLOUD
- Pledge assegnato a tutti gli esperimenti dei referaggi 2022 e 2023 con label "CLOUD"

	CPU (HS06)	Disk (TB-N)
QUAX	100	130
AMS-02	200	
HERD	1.000	100
SWG0	40	
Fermi	1.100	
AUGER	80	
Cygn0	160	10
<b>Totale</b>	<b>2.680</b>	<b>240</b>

	Crescita netta	
	CPU (HS06)	Disk (TB-N)
Cygn0	2.800	125
Darkside	200	100
NEWS	50	10
QUAX	400	130
<b>Totale</b>	<b>3.450</b>	<b>365</b>
<b>Totale effettivo</b>	<b>2.760</b>	<b>365</b>



Il totale effettivo con OF per ora non implementato

# T1 DISK and TAPE

# Disk storage in produzione

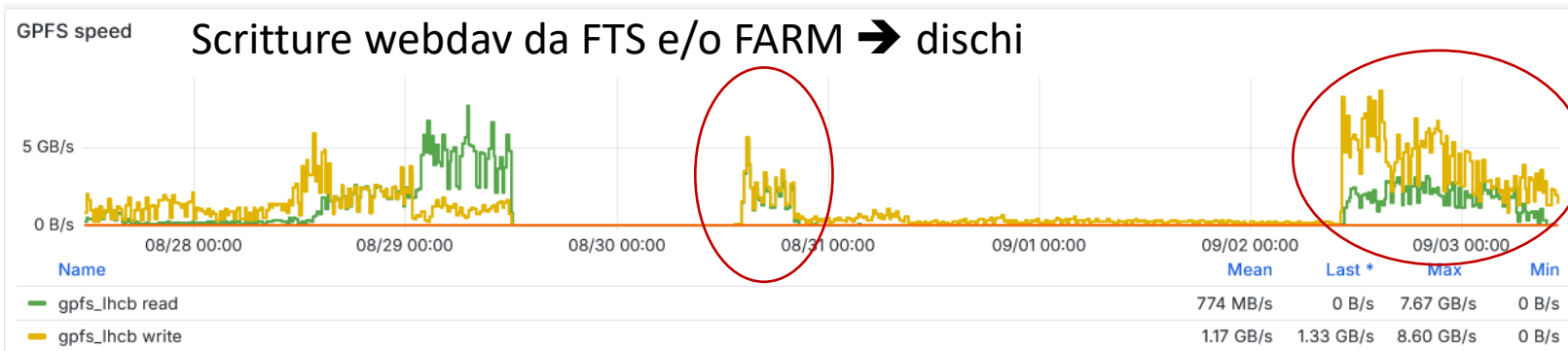
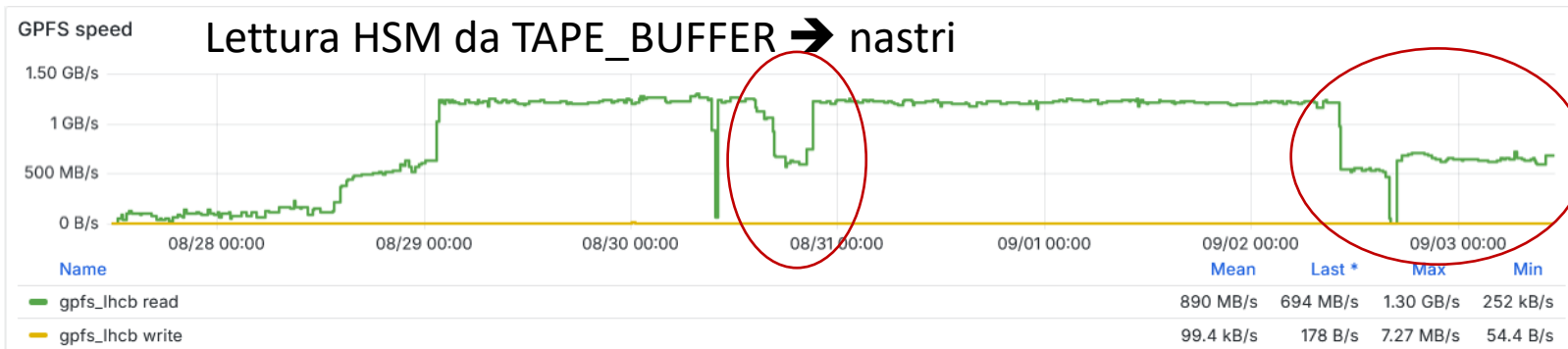
Installed: 113PB-33 da dismettere=80PB, Pledge 2024: 82.1 PB, Used: 61.7 PB

	Storage system	Model	Net capacity, TB	Experiment	End of support		
2015	ddn-10, ddn-11	DDN SFA12k	10120	ALICE, AMS	12/2022	Da Rimpiazzare con AQ 23-24	
	os6k8	Huawei OS6800v3	3400	GR2, Virgo	07/2024		
2016	md-1,md-2,md-3,md-4	Dell MD3860f	2308	DS, Virgo, Archive	12/2024		
	md-5, md-6 e md-7	Dell MD3820f	50	metadati, home, SW	11/2023 e 12/2024		
2017	os18k1, os18k2	Huawei OS18000v5	7800	LHCb	7/2024		
2018	os18k3, os18k5, os18k5	Huawei OS18000v5	11700	CMS	6/2024		
	ddn-12, ddn-13	DDN SFA 7990	5840	GR2,GR3	2025		Da spostare al Tecnopolo
	ddn-14, ddn-15	DDN SFA 2000NV	24	metadati	2025		
	os5k8-1,os5k8-2	Huawei OS5800v5	8999	ATLAS	2027		
	Cluster CEPH	12xSupermicro SS6029	3400	ALICE, cloud, etc.	2027		
	od1k6-1,2,3,4,5,6	Huawei OD1600	60000	ALICE,ATLAS,LHCb, CMS	2031		

Svuotato e riportato su GPFS

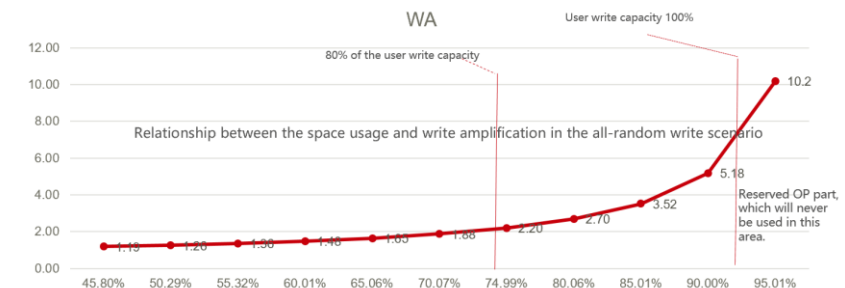
**Mancano 14PB da Pledge Storage 2022 – bloccati nella relativa Gara → In arrivo entro Settembre AQ 2023-2024 - fondi per secondo AS in arrivo a Settembre (ICSC) 16PB + eventuale 6/5 16PB**

# Problemi performance nuovo Storage Huawei



- Analisi in corso con vendor (HUAWEI)
- Soluzioni alternative in fase di valutazione se HUAWEI non risolve:
  - Spostare il buffer su sistemi NVMe
    - Ma troppi piccoli per le necessità di LHCB (2 settimane di dati)
  - Riutilizzare i vecchi sistemi HUAWEI estendendo la manutenzione
    - Almeno per parte delle necessità di LHCB (buffer\_tape?)
  - Utilizzare i sistemi DDN/Lenovo che stanno per arrivare (metà-fine Settembre)

- Calo drastico delle performance del sistema che ospita il TAPE\_BUFFER
  - Quando il buffer è quasi pieno
  - Quando ci sono scritture concorrenti sul disco da FTS e/o Farm via StoRM-WebDAV
- Storm-WebDAV va in crisi per alto numero di thread concorrenti
- Innesca effetto valanga

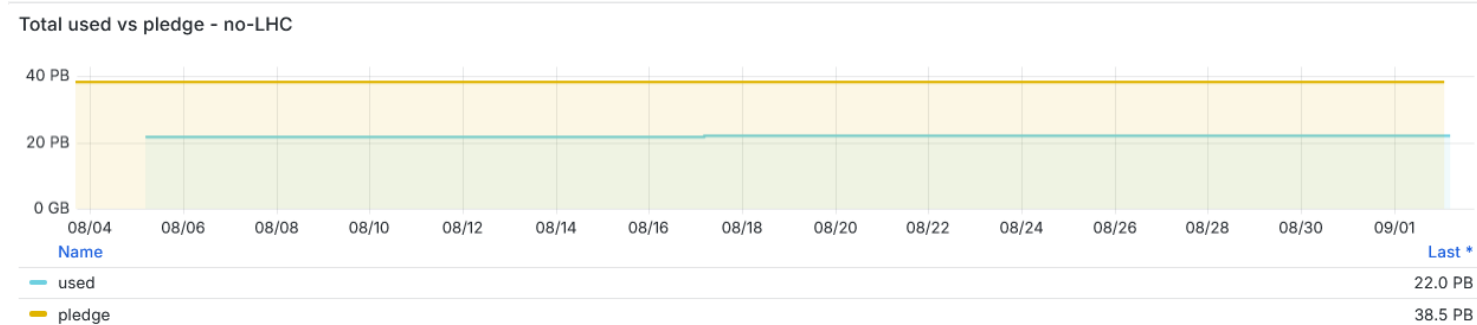
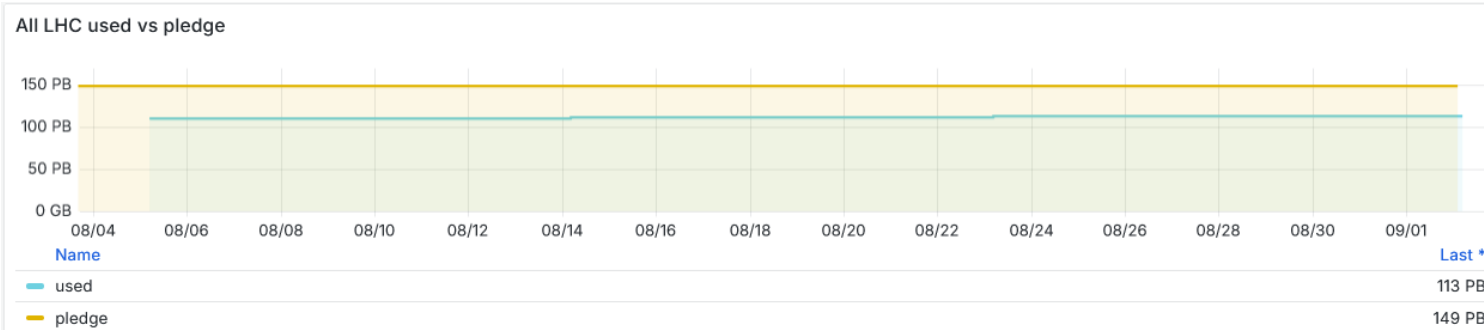
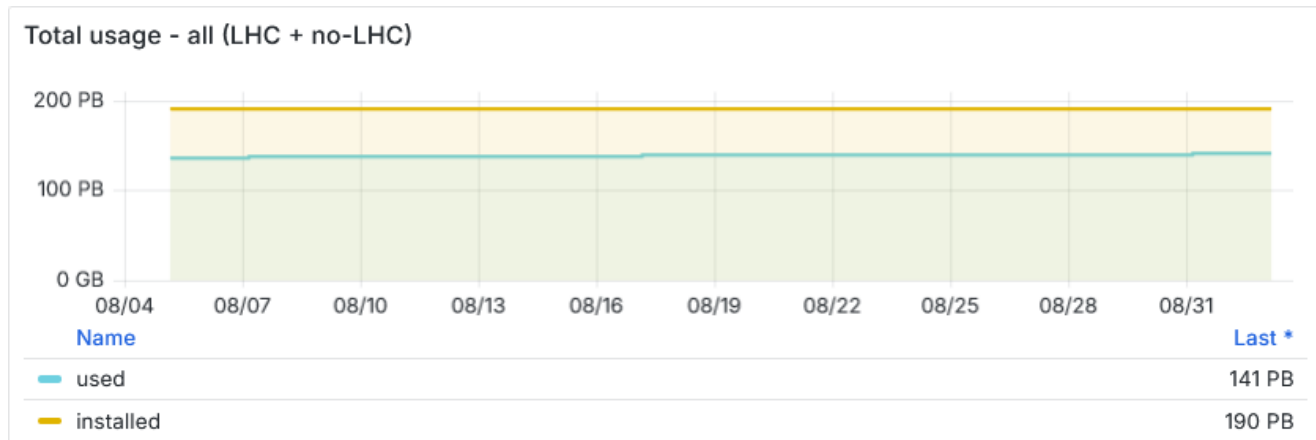


# Acquisti storage recenti e futuri

- Gara storage 2022 (14PB netti)
  - Nuova proposta con apparati DDN SFA7990X
  - In attesa per la consegna **entro fine Settembre**
- AQ storage 2023-2024 (Terabit+ICSC)
  - Huawei OceanStore Micro 1500/1600
    - 8 sistemi di 10PB + 40 server
  - Installazione e collaudo del primo AS 64PB effettuato
    - 60PB in prod
  - In produzione da fine Luglio
  - **Secondo AS da 16 PB a Settembre/Ottobre??**
  - **6/5?? occorre decidere se acquistarlo su ICSC**
- Tape Library (ICSC)
  - Nuova libreria Installata
  - Da completare collaudo e messa in produzione
- Gare nastri (ICSC)
  - Acquistati 14PB (JE e JF)
  - **URGENTE: Nuova gara di acquisto tape JF da 96PB**
    - Spedita in AC a inizio Luglio
      - In attesa del bando
    - Fondi ICSC arrivati – RDA inserita
    - Pledge+Overpledge+ICSC+Repack (2024) vecchia Oracle da dismettere

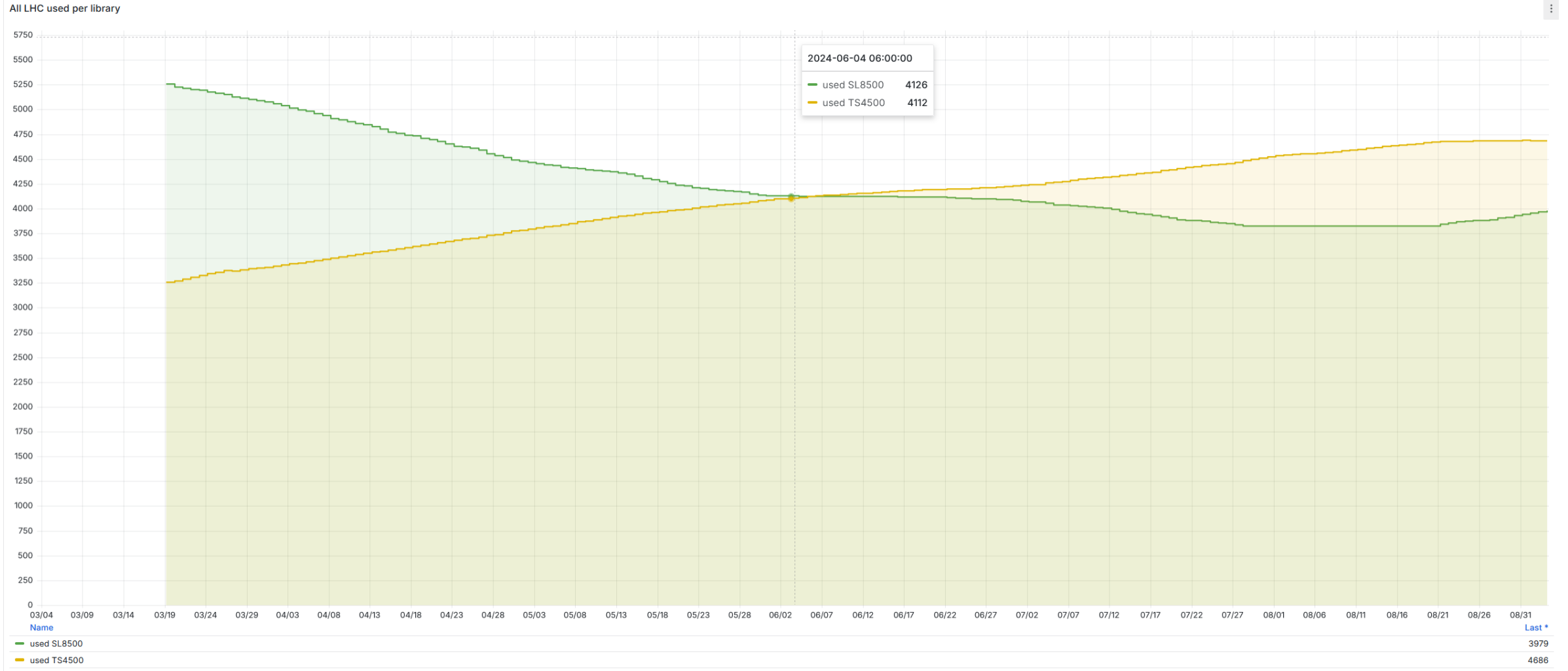


# Stato Tape

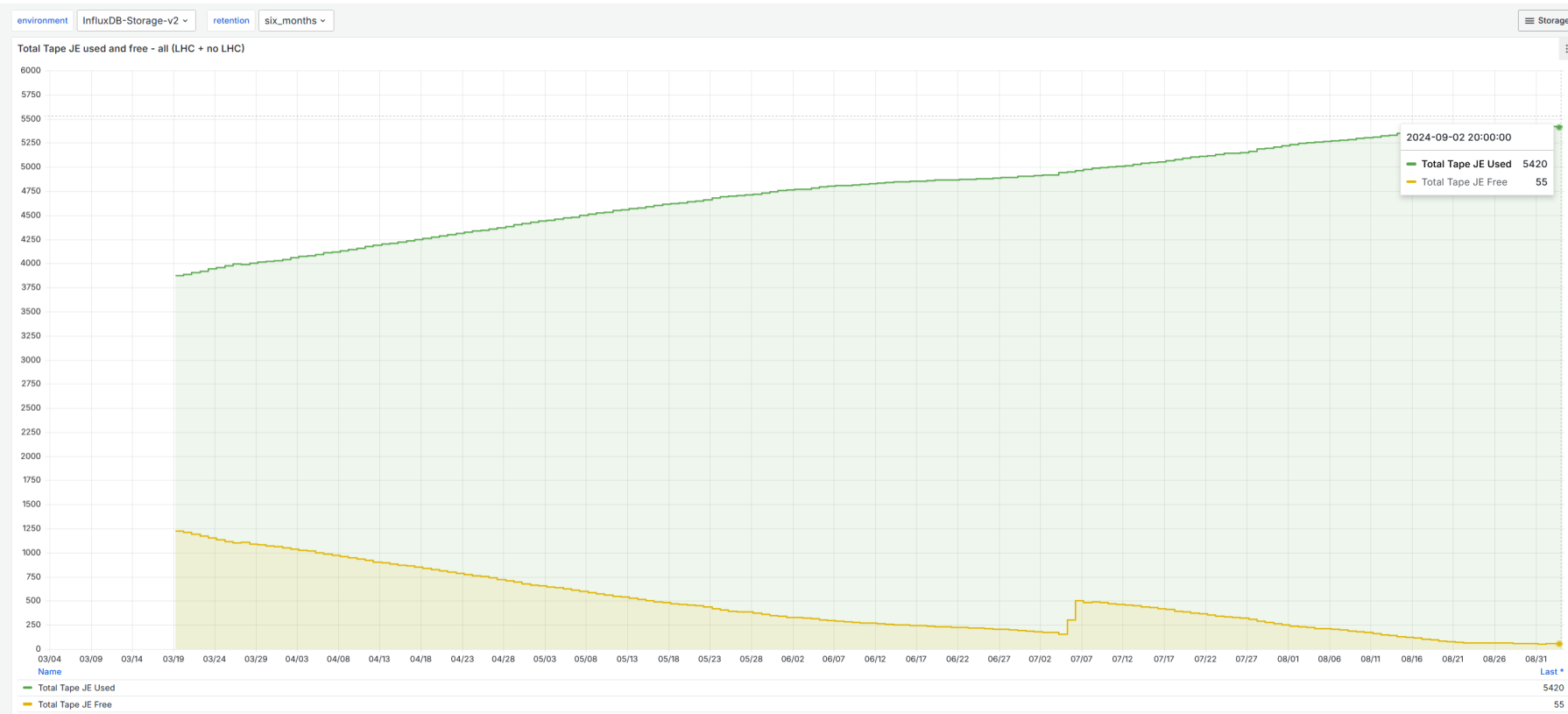


- Pledge 2024: 190PB
- Usato: 141PB
- Nuova Libreria già installata direttamente al Tecnopolo, non ancora in produzione → 8PB liberi
- Repack vecchia Oracle bloccato per risparmiare nastri
- Oracle di nuovo in produzione

# Uso nastri per libreria



# Cassette libere



- Pledge 2024: 190PB
- Usato: 141PB
- 8PB Liberi nella nuova Libreria
- 55 cassette da 20TB nella vecchia IBM in produzione
- + 8PB su nuova libreria
- **DIFFICILEMENTE ARRIVEREMO ALLA CONCLUSIONE DELLA PROCEDURA DA 96PB senza usare la Oracle**



# Stato trasferimento al Tecnopolo

- Da spostare sistemi storage “svuotati” (Ottobre 2024)
  - TD in corso
- Da svuotare un sistema storage (5 PB, 1 rack – DDN12-13)
  - Serve sistema vuoto al Tecnopolo (Novembre 2024)
- Da rilocare 1 tape library - “vecchia IBM”
  - Spostamento entro Q4 2024 (da concordare down con esperimenti)
- Da migrare 3 rack zona certificata ISO (EPIC) + 3 rack SSNN
  - Dipendenza da rete dedicata Tecnopolo-LNL (GARR)
  - Q4 2024?
- Da spostare risorse in housing (gennaio 2025)
  - 8 rack di INGV 3 rack DiFA di UNIBO
  - 3 rack Sezione INFN-BO
  - 3 rack servizi GARR



# Gara "HPC Bubbles"

- **Accordo Quadro Nazionale**
  - Listino prezzi per nodi + accessori
  - 2 anni di validità
  - Lotto1
    - CPU, GPU, FPGA
  - Lotto2
    - Storage
  - Sedi Coinvolte: CNAF, BARI, MI-BI, PI, TO, LNGS, NA, RM1, PD/LNL
- **Stato gara**
  - **Ordini inviati a parte 6/5**
    - **MI e LNL anticipati su capienza ordinaria**
  - **HW arrivato**
    - **CNAF**
  - **HW installato**
    - **PD, Torino**
  - **Mancano su 6/5**
    - **CT L1+L2, LNFESA L2, ROMA1 L1, NA L2**

## Quantità nodi con fondi Terabit-ICSC-DARE

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

Core: 30 kcore fisici  
Circa 34 HS/core

GPU: 244 NVIDIA H100  
40 FPGA  
InfiniBAnd 400Gbs

45 PB RAW



## \* Quantità nodi con fondi DARE – Terabit per Spoke8

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



# 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

# Stato zona ISO27001 - EPIC

- Siamo in attesa della risorse HBD acquistate da ACC
  - Ritardi molti significativi – devono rifare tutta la procedura
- Il nuovo cluster virtualizzazione in fase di installazione (Pichat)
  - fondi per svecchiamento 2024 – contratto S.Orsola
- Nuovi ToR acquistati su consip (in consegna - TP)
  - Fondi HBD
- nuovi concentratori (RDA in corso- TP)
  - Fondi HBD
- Trasloco rimandato a Q4 2024
  - Dipendenza da configurazione rete
  - 10 rack già allocati e predisposti

# Commenti/Richieste

- Per il 2025 sarebbe molto comodo avere le tabelle non solo dei delta ma anche del complessivo sul T1 per tutti gli esperimenti
- Anche per il 2024 I fondi referati non sono stati assegnati ad inizio anno
  - Fondi per OdF storage
  - Fondi per Nastri
- Avere il 90% delle risorse della farm HTC del T1 fuori dal T1 è sicuramente una strategia rischiosa
  - oltre che poco “gratificante”