

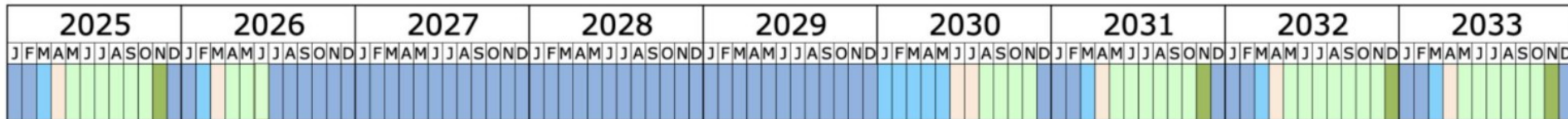
ATLAS status report

Frédéric Derue, LPNHE Paris

LCG France meeting, CC-IN2P3 Lyon

5th december 2024

- **Run 3 and LS3 program in 2025-2026**
- **Organisation and communication**
- **Ongoing ADC topics**
- **Heterogeneous resource evolution**
- **WLCG CPU/storage in ATLAS France in 2024**



- **Draft LHC schedule for 2025 [v1.0, 25 Nov 2024]**
For approval by CERN Research Board on Dec 4th and LPC on 5th
- **Main features (partly as presented in S&C week in October)**
 - extended running in 2025, another month compared to original schedule
 - see also *slides* from Laurent yesterday

Remaining Run3 schedule (2025/2026), **proposal**

Proposal, up for discussion.

- Experiments feedback highly desired.

Main features:

- Single ion run end of 2025
- No pPb in Run3
- Cooldown before LS3 using MD days
- ATLAS/CMS interested in ≥ 1 week of low-pileup data taking?
- Additional MD block for cooldown before VdM
 - Highly desired to remove VdM
 - pp reference in 2024 sufficient?
 - Removing it would allow reaching the 5.3 /nb to PbPb target for Run3.



- Pb(p)Pb runs: ATLAS prefers single long PbPb run, and more PbPb luminosity over pPb data, but if pPb, should be at 5.36 TeV CM energy (2016 was 8.16 TeV and 184 nb⁻¹), discussions ongoing (likely: split scenario with 18 ion physics days in each 2025 & 2026)
- short YETS, something like 7-9 weeks (depending on if you include Christmas)
- end of Run 3 on June 30th

- **In the next round (March 2025) re-adjust for updated HL-LHC schedule**
 - LS3 from July 2026 to May 2030 (likely means very little “real” data in 2030)
 - no more LS5 - running from ~2036 to 2041 (6 straight years!)

C-RSG Report Sept. 2024 - Preliminary Request for 2026

ATLAS Report to the C-RSG, September 2024
3rd September 2024, 1.0

ATLAS Computing Status and Plans Report to the C-RSG, September 2024

Contents

1. Introduction	3
2. Resources in 2024	4
3. Computing activities in 2026 Scenario A: No collisions	5
3.1 Data reconstruction in 2026	6
3.2 Monte Carlo sample production	7
3.3 Preparation of samples for physics analysis	7
3.4 Upgrade studies	8
3.5 Summary of data processing activities for Scenario A	8
4. Computing activities in 2026 Scenario B: Full year of collisions	9
4.1 LHC and data-taking assumptions	9
4.2 Data reconstruction in 2026	10
4.3 Monte Carlo sample production	10
4.4 Preparation of samples for physics analysis	10
4.5 Run 3 heavy ion data processing	11
4.6 Upgrade studies	11
4.7 Summary of data processing activities for Scenario B	12
5. Computing model parameters	13
5.1 Event processing times	14
5.2 Event sizes	14
5.3 Disk, tape, and replication factors	14
6. Resource requests for 2026 and outlook for 2027 and beyond	17
6.1 Resources requirements for 2026	18
6.2 Tier-0 requirements	20
6.3 Unpledged resources	20
6.4 Resources request for 2026	20
6.5 Outlook beyond 2026	20
7. C-RSG Recommendations	20

- **Report from Sep. is available here**
 - for a full presentation on this round see this ATLAS Weekly talk

- **Unusual and rather difficult report to write**
 - at the time of writing: unclear LHC schedule for 2026
 - unclear which resources are available in LS3 (HLT)
 - needed to consider multiple scenarios

- **It does contain meaningful and important inputs**
 - reduced RAW size 1.6 MB/event (already expected in 2025)
 - updated several other data sizes based on production observations HITS (-10%), RDO (+10%), Data AOD (+10%), MC AOD (-10%), PHYS (-10%)
 - simulation time further reduced 20% of the last report
 - gentle updates to Generator time (-15%), MC digi+reco (- ~10%) to align with what is observed
 - revised and updated resource requirement model
 - first experiment to accept ARM as pledge (see later)

2026 is an (almost) regular data taking year

- For data taking, pretty typical increases for sites
 - 2025 final increase requests were 10% / 14% / 24% for CPU / disk / tape
 - mostly driven by larger total data set size and matching increases in MC

		2025 Agreed @ April 2024 RRB	2026 Scenario B Prel. Request @ September 2024 RRB	Balance 2026 Scenario B wrt 2025
CPU	T0 (kHS23)	1100	1265	15.0%
CPU	T1 (kHS23)	1635	1787	9.3%
CPU	T2 (kHS23)	1998	2184	9.3%
CPU	SUM (kHS23)	4733	5235	10.6%
Disk	T0 (PB)	56	65	16.1%
Disk	T1 (PB)	186	201	7.9%
Disk	T2 (PB)	227	245	8.0%
Disk	SUM (PB)	469	511	8.9%
Tape	T0 (PB)	258	320	24.0%
Tape	T1 (PB)	561	707	25.9%
Tape	SUM (PB)	819	1026	25.3%

Final numbers will be different again in Spring 2025, with a shorter 2026 than considered here

Continuing 45/55% split between Tier 1/Tier 2

With data taking, tape increase driven by data

Table 9: Summary of the preliminary requests for computing resources for 2026 Scenario B.

- **Scenario A1: HLT farm successfully relocated and running for a fraction of 2026***
- **Scenario A2: No successful relocation of the HLT farm in 2026**
 - increase depends critically on HLT farm availability: 2.3M HS23, turned off “on Day 1” of LS3
 - with 30% HLT farm availability, 5% CPU increase, well below ‘flat budget’ scenarios (15%)
 - without the HLT farm, ~20% CPU increase

		2025 Agreed @ April 2024 RRB	2026 Scenario A1 Prel. Request @ September 2024 RRB	Balance 2026 Scenario A1 wrt 2025	2026 Scenario A2 Prel. Request @ September 2024 RRB	Balance 2026 Scenario A2 wrt 2025
CPU	T0 (kHS23)	1100	1210	10.0%	1265	15.0%
CPU	T1 (kHS23)	1635	1719	5.1%	1968	20.4%
CPU	T2 (kHS23)	1998	2101	5.2%	2405	20.4%
CPU	SUM (kHS23)	4733	5030	6.3%	5639	19.1%
Disk	T0 (PB)	56	62	10.7%	62	10.7%
Disk	T1 (PB)	186	194	4.5%	194	4.5%
Disk	T2 (PB)	227	238	4.7%	238	4.7%
Disk	SUM (PB)	469	494	5.3%	494	5.3%
Tape	T0 (PB)	258	268	3.8%	268	3.8%
Tape	T1 (PB)	561	644	14.7%	644	14.7%
Tape	SUM (PB)	819	911	11.3%	911	11.3%

*using 30% as our benchmark availability fraction, similar to what we get during operation

Continuing 45/55% split between Tier 1/Tier 2

No data taking; tape increase driven by MC

Table 8: Summary of the preliminary requests for computing resources for 2026 Scenario A. Two sets of numbers are shown, Scenarios A1 and A2, which differ according to the availability of the decommissioned Run 3 HLT farm for offline use in 2026.

PROBLEM POSTPONED BY ONE YEAR

HLT farm availability will nevertheless be a concern in 2027

Still need to find a solution, but with an extra year’s time (might be easier to get it done) and hardware that is a year older (might be harder to find a taker)

- **Scen. A2 (or similar) would clearly be a substantial request for 2027**
 - modest disk and tape requirements offset the big demand for CPU
 - request should be below flat-budget, but sites not anticipating significant spending at start of LS3
- **ATLAS has been working on a number of options to relocate / redeploy the hardware: It is a lot: 30% of our worldwide compute, >50 racks, 700 kW; out of warranty, ~6-7 years old; (some grid CPUs >10 year old)**
 - single site solution if possible (many sites means the hardware may never appear in the monitoring, and we wouldn't have much follow up). Two stand out:
 - LHCb containers at CERN: IT has sent cost estimates for move and maintenance (not too bad) Need to loop back with LHCb on container availability / schedule. Starts to look promising
 - US solution: Might be able to take all the hardware, with US-ATLAS retaining some (significant) fraction, which we would use opportunistically. Need to discuss wrt the new schedule
 - update in the final 2026 request in April and have a plan in place for preliminary 2027 request in September
- **If cannot be solved, we will have to look for mitigation strategies**
 - try to get small parts of the HLT hardware relocated, look for more unpledged CPU
 - some French sites declared their interest (LPNHE, IJCLab, IRFU (?)) for a couple of racks
 - transport (how much ?) will/would be paid by sites
 - reduce/delay/spread out MC campaigns
 - could impact the physics programme of the experiment
 - more optimisation of our workloads to reduce CPU demands

- **ADC Coordination Board, ADC Weekly and ADC Operations daily meeting including an additional meeting in US time zone**
 - area coordinators meet weekly in a closed meeting
 - sites and production managers meet weekly
 - operation team meets daily: two meetings at 9am (CET) and 3:30pm (for US-TZ)
- **Communication with grid clouds and sites**
 - several cloud and sites have requested improved communication between sites and ADC regarding issues
 - communication from ADC to the sites: Decide on a single centralized notification way to inform sites of outages, avoiding unnecessary local efforts on known external issues
 - communication from sites to ADC: Decide on a single centralized means for sites to contact ADC to report site-specific issues, local problems, or ask questions
 - **ADC use in general either GGUS tickets either mail to atlas-support-cloud-fr. I/Fred then fwd either to all FR-sites or to specific sites**
ADC thought reaching all FR sites – but only few persons (doing the support) are in this list
 - **established a ~monthly ATLAS Fabrics meeting for ADC, Sites & Service Providers, which consistently saw strong site participation, averaging 30-40 participants**
 - **[indico, last meeting](#)**
 - ADC Site Jamboree during next S&C week [[indico](#)]
 - ADC Technical Interchange Meeting (TIM) in Stony Brook end of January [[indico](#)]
- **Documentation**
 - Gradual transition from ATLAS [Twiki](#) to <https://atlas-computing.docs.cern.ch>

● Lack of central effort for central operations

- similar reports from other areas
- ADC are deeply concerned about the lack of central effort. DDM ops is critical
 - it also feels like sites have retracted into themselves.
 - We've asked them to make explicit some of the opportunities and costs of this
 - existing people (mostly on soft money) are being squeezed more and more.
 - US-ATLAS continues to provide significant effort, although sometimes whilst we can hang on to someone longer, ADC retains a smaller fraction than before
- ATLAS Computing infrastructure cannot be operated efficiently
 - from fractions of people distributed worldwide
 - from short term appointments
 - from people fighting to remain at CERN and to be stabilised
- failures on infrastructure affect expensive resources worldwide
- some resources need to be filled continuously to reach the pledge, other resources filled intensively in a time frame to reach allocation, other resources have to stop after a budget cost has reached a limit
- Importance to have focused people on ADC tasks with significant presence at CERN
- discussion with Funding Agencies (ATLAS ICB), inside ATLAS France

Note: Due to lack of personpower, DDM Ops will only handle already agreed high priority tasks (such as ensuring data taking or running the lifetime model).

All other requests are summarized [here](#): 16 tasks among which 3 for French sites :

- 2 for GRIF – lost files & unable to store files,
- 1 for LAPP – decommissioning of FR-ALPES_*_DATALAKES

These will start to be processed in January 2025.

● Site evolution

- we are observing an increasing number of fluctuations in resources deriving in long unavailability
 - the power unit of a Tier 2 site broke mid July and the new unit are expected back mid October
- geo-political decisions
 - we stop using some resources
- end of funding
 - sites are decommissioned, sometimes without warning to ATLAS
- in the case of CPU we are less affected but for storage this needs a lot of effort
 - unique data can be trapped
 - secondary copies are also concerned because the copy is counted by Rucio but the data is not available
 - we may lose data because we have an incident in a site and then the secondary copy is unavailable

● Recovering of unavailable DAOD trapped at sites

- operational issues trigger some development to mitigate them
- we are now able to recreate input files or even full datasets automatically if the file is lost or unavailable by rerunning the jobs
- the files in the dataset are replaced and Panda refreshes the list input files while tasks are running

● Memory optimization

- ATLAS require to sites to provide resources with minimum 2 GB of RAM per running core
- some workflows require higher memory requirements per core
- now most of the sites have part of the resources with 3 GB/core or 4 GB/core
- up to now we had the limit of memory per core in CRIC declared by the site
- now we ask the sites to provide the average of memory usage by all jobs submitted to them and the maximum memory per core they could accept if necessary
 - allowed maxrss/core increased from 2 to 6GB on 4th Nov
but the target meanRSS was left at 2GB/core
 - see ticket [ggus 169241](#) for an example of side effect (and solution ?) at GRIF-LPNHE
- this allows sites to accept a lower number of high memory jobs in exchange of jobs with lower than average
- this allows us to run high memory jobs worldwide without affecting too much the sites

● Tape resources

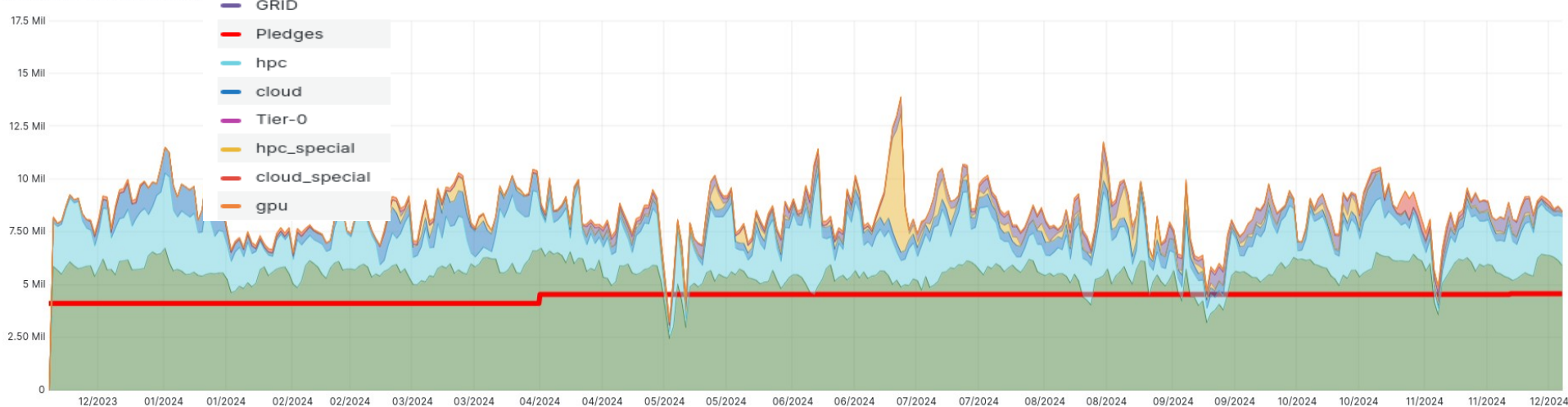
- an increasing number of issues related to tape storage
- originally tape was designed as a distributed long term archive for RAW data
- tapes are not a replacement for disks
- tape systems are subject to different constraints compared to disks, such as tape remounting limits, sensitivity to humidity
- tape storage is oriented to backup and disaster recovery
- tape storage is usually a shared resource between ATLAS, other experiments and IT services
- this means that if the price of tape storage for archival is lower than disk, the price of a tape storage for random data serving skyrockets
 - this is one of the cases where the cost for TB is sensitive to the usage
- there is a development to mitigate these issues which is the Data Carousel
 - available for production tasks
 - requests to tape files are grouped to optimize the tape usage
- there is another development in the pipeline to use metadata to group collections of datasets/files in the same set of tapes

● Computing usage by resource type

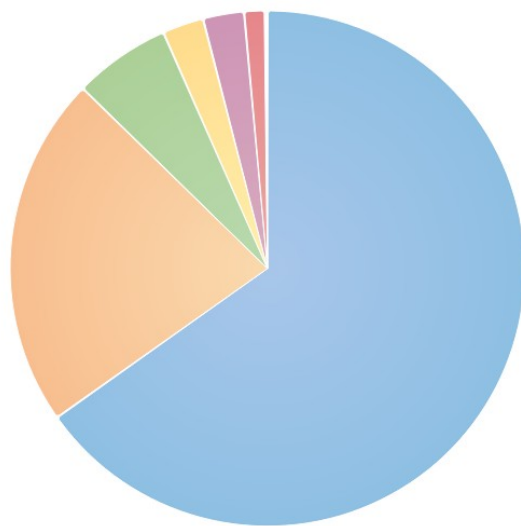
[link](#)

Slots of running jobs (HS23)
since a year

Slots of Running jobs (HS23) by AI



[link](#)



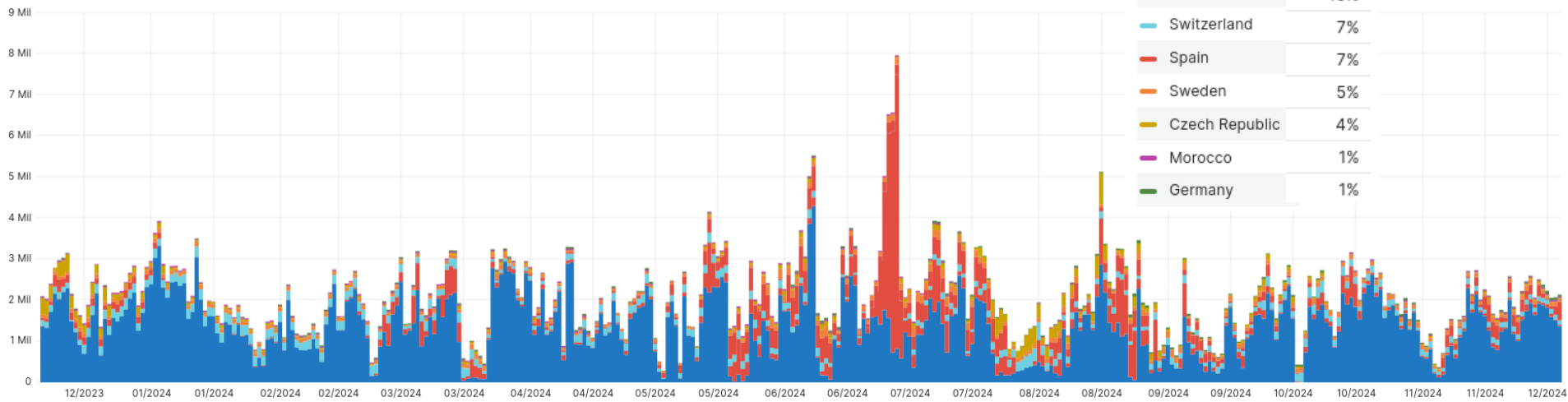
	Value	Percent
GRID	2.03 Bil	66%
hpc	685 Mil	22%
cloud	187 Mil	6%
hpc_special	78.3 Mil	3%
Tier-0	77.6 Mil	3%
cloud_special	37.6 Mil	1%
UNKNOWN	79.3 K	0%
gpu	17.6 K	0%

● HPC

[link](#)

Slots of running jobs (HS23)
since a year

Slots of Running jobs (HS23) ⓘ



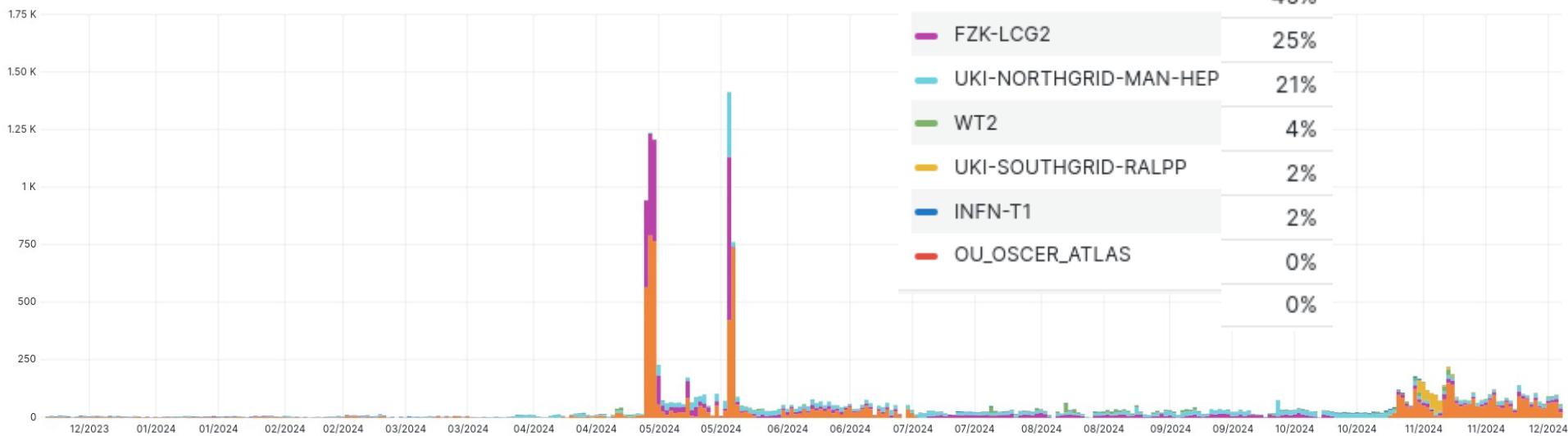
- seen as free resources for funding agencies
- a workhorse for simulation
- we have a significant usage of HPC resources (~25% of computing usage)
- we cannot run all the workflows in HPC
- we cannot influence security and architecture requirements
- allocations are a fixed number of CPU hours within a time frame, for ex. 3 months
- advantages
 - usually room for high memory per core
 - we can run all CPU allocation in a limited time. Good for peaks of simulation needs

● GPU

[link](#)

Slots of running jobs (HS23)
since a year

Slots of Running jobs (HS23) ⓘ



- able to run jobs using GPUs in the ATLAS production system
- up to now the usage of GPU resources is low
- GPU resource need currently limited to 10s of GPUs:
 - however, if we progress like CMS (10% of reco task on GPU) may not have access to enough
 - FastCaloSimGPU large scale PhysVal (code ready and queues exist), expected until Q2 2024
- several technical issues to solve (licensing software)
- software is not ready yet for major workflows
- new HPCs offer more resources for GPUs than for CPUs
- training, user analysis - not used in production ⇒ do not appear in plot above !
 - each ATLAS France group is using GPUs at CC-IN2P3, locally (lab, university), JeanZay, ...

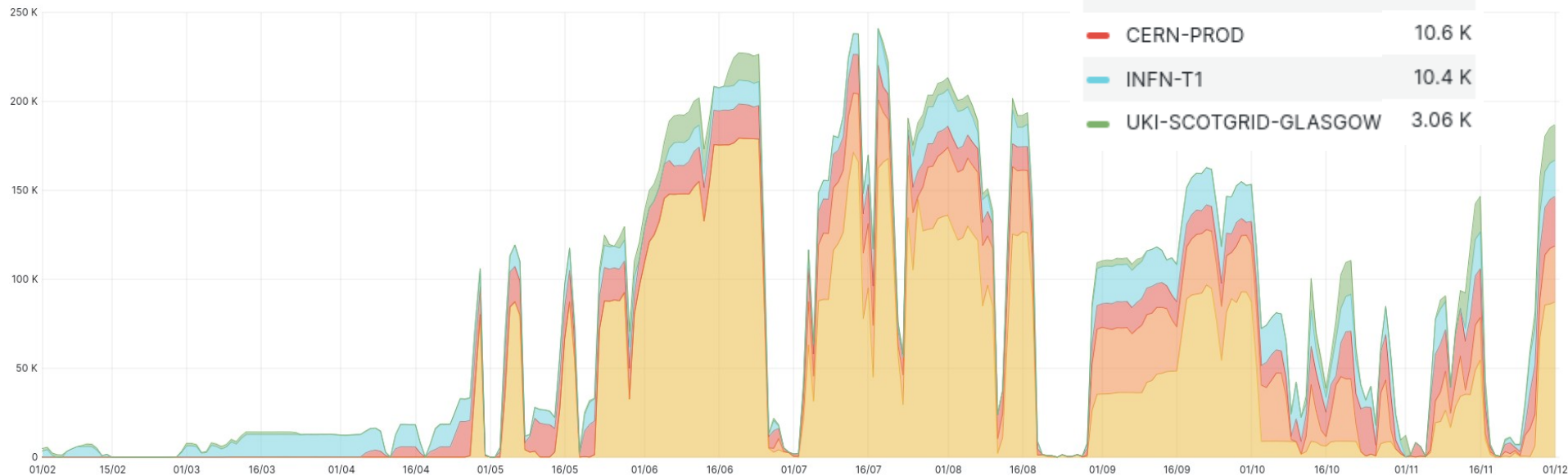
● ARM

[link](#)Slots of running jobs (HS23)
since a year

avg ▾

SWT2_GOOGLE	43.2 K
FZK-LCG2	11.1 K
CERN-PROD	10.6 K
INFN-T1	10.4 K
UKI-SCOTGRID-GLASGOW	3.06 K

Slots of Running jobs (HS23) by ADC activity ⓘ



- first experiment to accept ARM as pledge
 - up to 50% at any single site may be provided using ARM processors
- ARM looking good (PhysVal for both sim and reco) -
 - if successful, can/will embrace ARM for production - What about user analysis?
 - expect a reasonable amount of sites to deploy ARM due to expected energy saving
- represent 1% of computing usage in 2024

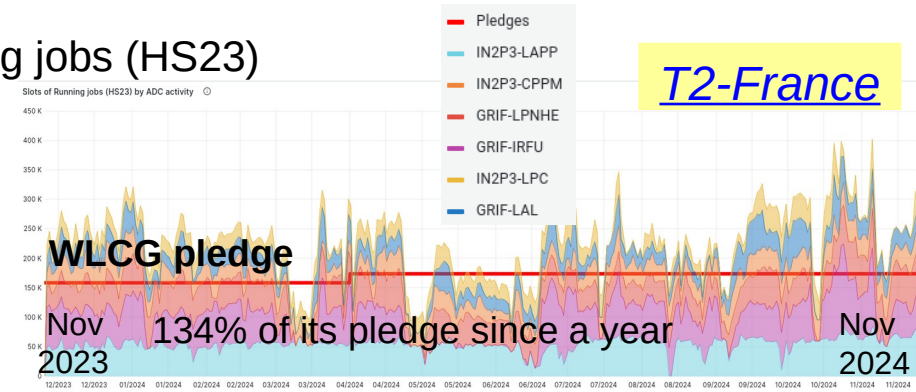
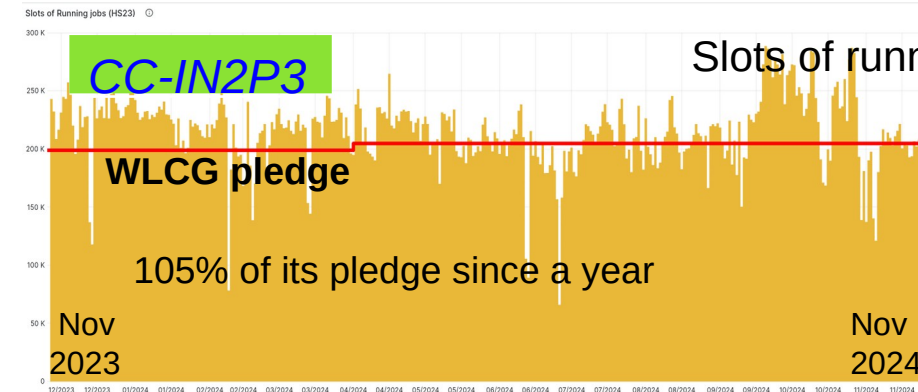
● Pledges (see *cric*)

- CPU
 - CC : 13.5% of T1s, +4% in 2025
 - T2s : 7.5% of T2s, +7% in 2025
- Storage
 - CC : disk=13.5%, tape=15.6% of T1s
in 2025 : +14% for disk, +24% for tapes
 - T2s = 9.5% of T2s, +9.5% in 2025

Site	CPU Pledge 2024 (HS23)	Disk Pledge 2024 (TB)	Tape Pledge 2024 (TB)
IN2P3-CC	204660	22005	65540
GRIF	69527	7221	-
IN2P3-CPPM	24000	2200	-
IN2P3-LAPP	55000	7000	-
IN2P3-LPC	25000	2914	-

● CPU realized

- CPU : CC : 15% of T1s, T2s : 15% of T2s



- French pledges (in %) remain at same level as last years
- + « local » resources in Lyon and labs.
- no HPC cpu resources in France
- recent increase of price of TB/cpu ⇒ need of R&D

- **S&C in ATLAS France**

- smooth operation of CC-IN2P3 as a Tier-1 and our Tier-2s
but less person-power in sites, FR-support and ADC/DDM central Ops
⇒ life is harder and we can feel it !

**A BIG thanks to CC-IN2P3, all our Tier-2s colleagues,
& LCG-FR management for the operation, maintenance
and development of our computing infrastructure**

**R&D for storage/analysis evolution requires person power
– syst admins and physicists –
in order to use solutions which fit our needs !**