# Use of resources with ARM architectures in LHC experiments

L. Rinaldi, L. Anderlini, T. Boccali, C. Bozzi, L. Carminati, F. Noferini, D. Spiga, M. Veltri

22/05/2024 CCR Workshop 2024

# Outline

- Why are HEP experiments interested in ARMs ?
- Available resources @ CNAF
- Experiments usage:
  - ALICE
  - ATLAS
  - CMS
- Outlook and future plans

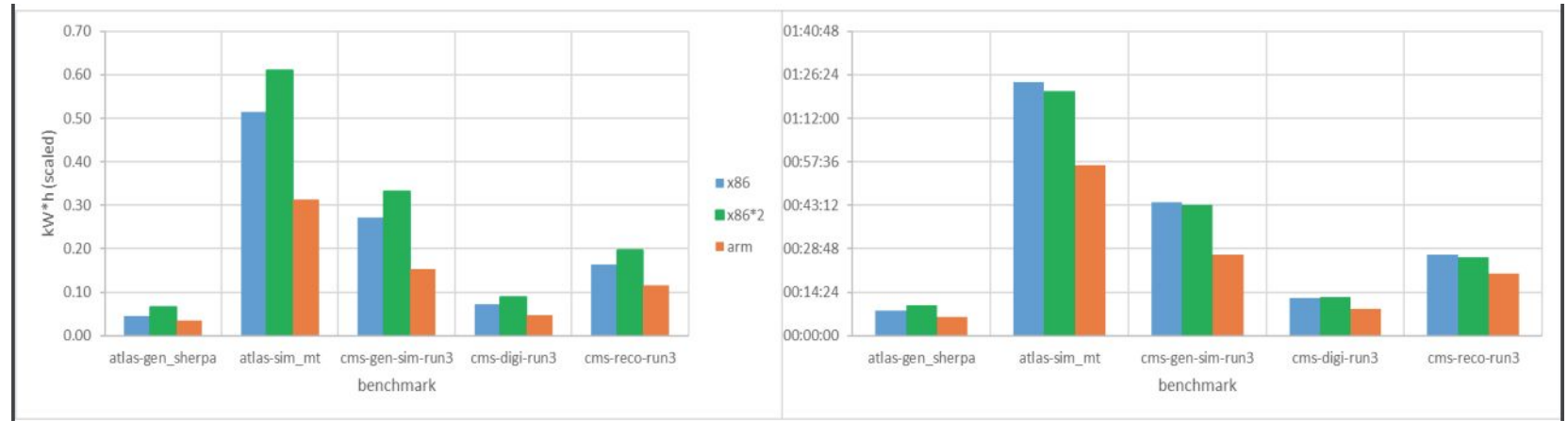# Why are HEP experiments interested in ARMs

- The ARM architecture is becoming a competitive and energy efficient alternative

- Some surveys indicate its increased presence in HPCs and commercial clouds, and some WLCG sites have expressed their interest

- Chip makers are also developing their next generation solutions on ARM architectures, sometimes combining ARM and GPU processors in the same chip

- It is important that the HEP exp software embraces the change and is able to successfully exploit this architecture.

# ARM Flagship (use-case) in CN-HPC Spoke 2

- At the moment, all the 4 major LHC experiments are providing software builds able to run with this architecture,

- A complete validation in full GRID submission chain is still missing

- The goal of this Flagship use case provide the hardware and the software infrastructure to enable a validation of the software used by experiments to process data and Monte Carlo (MC), at least for two major LHC experiments, using resources from the ICSC datalake and in particular hosted at the INFN CNAF Tier-1 facility.

- Lower Total Cost of Ownership (TCO) and move towards more sustainable computing models in terms of energy saving

# ARM Flagship (use-case) in CN-HPC Spoke 2

Performances on ARM architectures allow to reduce power consumption at the level of 30-60%, with respect to AMD cpus, depending on the HEPscore benchmark used.
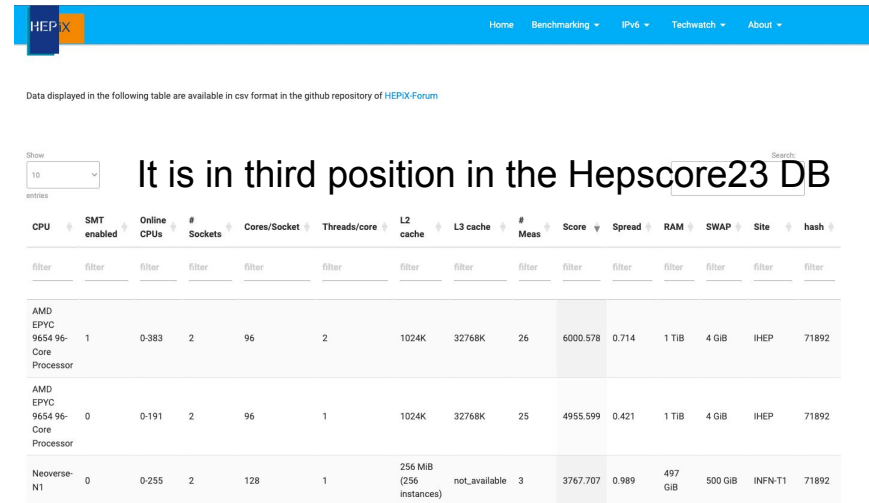
https://indico.cern.ch/event/1106990/contributions/4991256/attachments/2534801/4362468/PoW_ACAT2022.pdf

# Available resources @ CNAF

- Resource procurement started in the second half of 2023 when CNAF acquired two ARM nodes allowing the experiments to loging on those machines and running preliminary tests. By the end of 2023 CNAF acquired two additional nodes and setup a GRID-HTCondor queue.
- 4 ARM nodes: Neoverse-N1, 256 core, 1TB RAM, 8TB NVMe

Current setting (still work in progress)

- Cvmfs available
- Network: access to external network
- Gpfs client -> not yet available for ARM
- Condor/GRID -> in production

**Many thanks to the CNAF crew for support!**

It is in third position in the Hepscore23 DB

Data displayed in the following table are available in csv format in the github repository of HEPiX-Forum

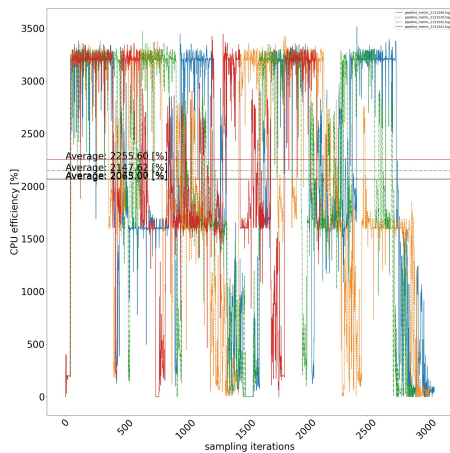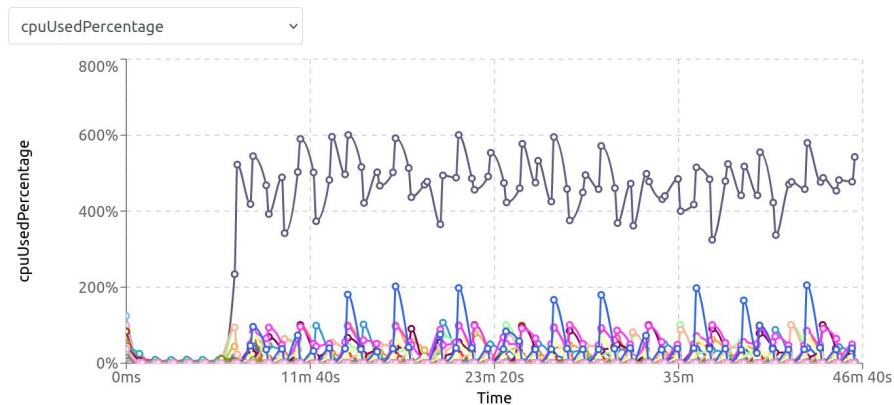| CPU | SMT enabled | Online CPUs | # Sockets | Cores/Socket | Threads/core | L2 cache | L3 cache | # Meas | Score | Spread | RAM | SWAP | Site | hash |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| filter | filter | filter | filter | filter | filter | filter | filter | filter | filter | filter | filter | filter | filter | filter |
| AMD EPYC 9654 96-Core Processor | 1 | 0-383 | 2 | 96 | 2 | 1024K | 32768K | 26 | 6000.578 | 0.714 | 1 TiB | 4 GiB | IHEP | 71892 |
| AMD EPYC 9654 96-Core Processor | 0 | 0-191 | 2 | 96 | 1 | 1024K | 32768K | 25 | 4955.599 | 0.421 | 1 TiB | 4 GiB | IHEP | 71892 |
| Neoverse-N1 | 0 | 0-255 | 2 | 128 | 1 | 256 MiB (256 instances) | not_available | 3 | 3767.707 | 0.989 | 497 GiB | 500 GiB | INFN-T1 | 71892 |

# Experiment usage: ALICE

Results from preliminary tests with resource provided by E4 presented at [CCR23](CCR23)

- ## Simulation



4 concurrent 16 workers (+ overbooking)
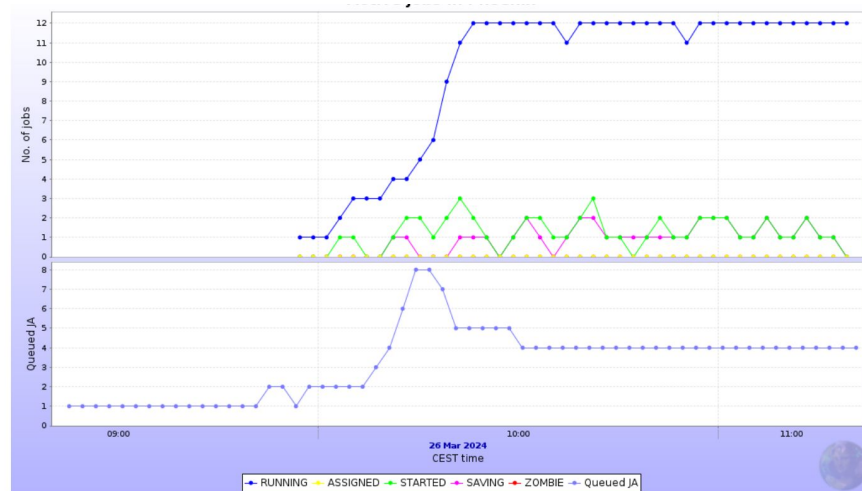MC job

- ## Data reconstruction



- GRID setting tuned for 8 cores per node
- CPU efficiency consistent with what observed in the GRID node
- Physics validation on the output not yet done

# Experiment usage: ALICE

ARM@CNAF

- At the end of 2023 ALICE set the queue at CNAF (an instance was already running at Glasgow)
  - ARM build on el7-aarch64 → still very unstable (software issues)
- So far we managed to run only few simulation jobs but not reconstruction jobs
- From May 2024 the builds was moved to slc9 (new tests to be started soon).
- From now on CNAF is the unique site providing ARM resources to ALICE
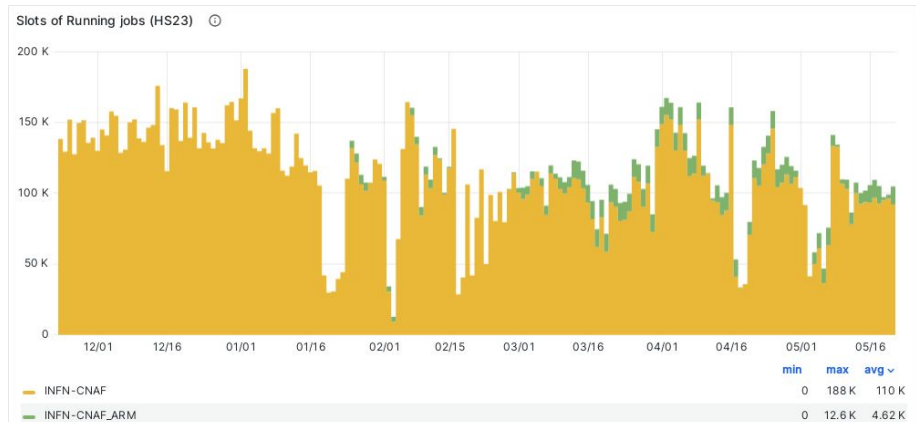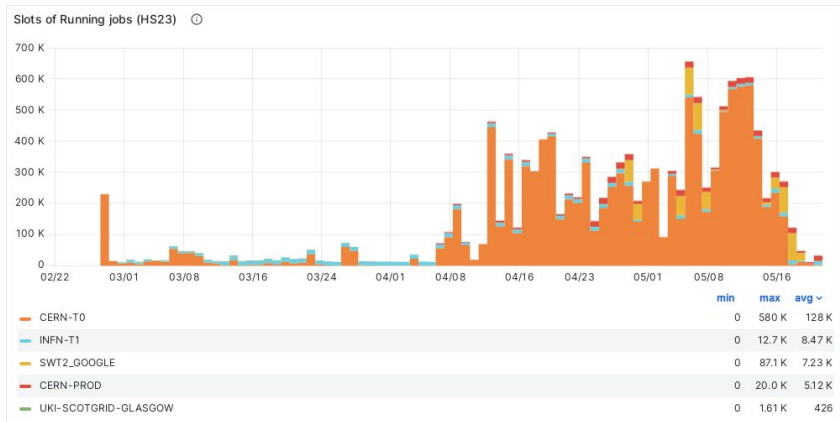
# Experiment usage: ATLAS

ATLAS SW already validated for

- T0-processing
- MC-SIMULATION (FULL-FAST) and MC-RECO→~60% of total atlas computing power)
- Some user analysis workflows

ARM@CNAF

Running since nov 2023, on el9 container, with PanDA production system
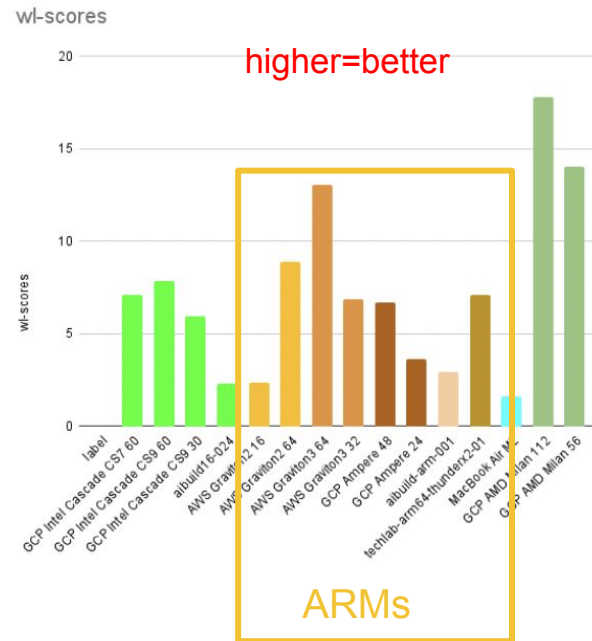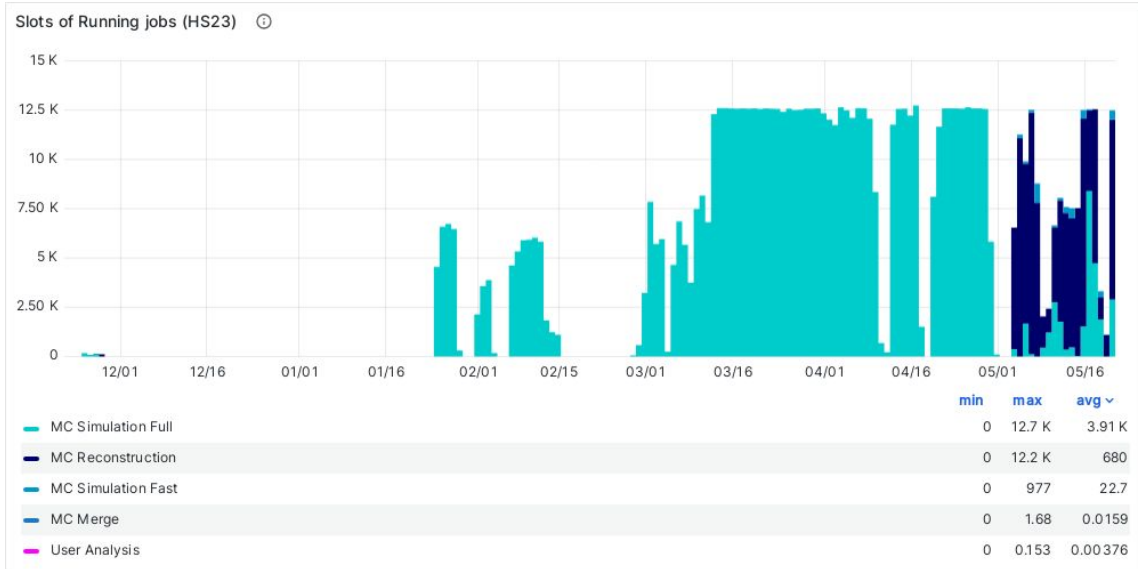
# Experiment usage: ATLAS



Slots of Running jobs (HS23)

|  | min | max | avg |
|---|---|---|---|
| CERN-T0 | 0 | 580 K | 128 K |
| INFN-T1 | 0 | 12.7 K | 8.47 K |
| SWT2_GOOGLE | 0 | 87.1 K | 7.23 K |
| CERN-PROD | 0 | 20.0 K | 5.12 K |
| UKI-SCOTGRID-GLASGOW | 0 | 1.61 K | 426 |



Slots of Running jobs (HS23)

|  | min | max | avg |
|---|---|---|---|
| INFN-CNAF | 0 | 188 K | 110 K |
| INFN-CNAF_ARM | 0 | 12.6 K | 4.62 K |

Where?
- Mainly at T-0, followed by CNAF

How much?
- At peak: 2.6% world-wide pledge

@CNAF
- At peak: 12.5 kHS23
  - 9% INFN-T1 pledge

10

# Experiment usage: ATLAS



Slots of Running jobs (HS23)

|  | min | max | avg |
|---|---|---|---|
| MC Simulation Full | 0 | 12.7 K | 3.91 K |
| MC Reconstruction | 0 | 12.2 K | 680 |
| MC Simulation Fast | 0 | 977 | 22.7 |
| MC Merge | 0 | 1.68 | 0.0159 |
| User Analysis | 0 | 0.153 | 0.00376 |

wl-scores

higher=better

ARMs

What (@CNAF)?
- MC FULL Simulation

How many events per HS23/hour?
- 3 events/hour/HS23 (7.5 for other machine@cnaf) VERY PRELIMINARY!
- Results are comparable with other ATLAS measurements

# CMS: Status and motivations

Software: CMS continuously invest a lot of effort to evolve and improve the CMS-SW (CMSSW). Since a few years, among many others, it support ARM architectures.

Computing: In the recent years, thanks to activities carried on within INFN with Marconi 100 at CINECA (a Power9 machine) the computing stack of the experiment fully enabled multiple architectures

Currently several opportunities for exploiting ARM are arising and in the near future this might increase.
- We envision ARM to become more widely deployed especially due to its superior power efficiency.
    - A few examples: **UK** with Isambard 3; **ES** with HPC @ BSC will host a second partition based on ARM CPU partition (Nvidia Grace chips), **US** with ARM for HPC at Ookami…

To use ARM based system in production, the physics validation is mandatory. Having a "stable" allocation is a key to the success in this respect.
- The physics validation means comparing the distributions of physics observables regarding any significant deviations.
- In turn this means producing several twisted samples (both on x86 and ARM)... repeating until all is completely understood!

**This is why the access to the two ARM nodes at CNAF represents a fundamental enabler**

# CMS: Status of the ARM validation

**The ARM nodes at T1_IT_CNAF have been integrated as a sub site of the regular Tier1 and thus accessed via GRID**
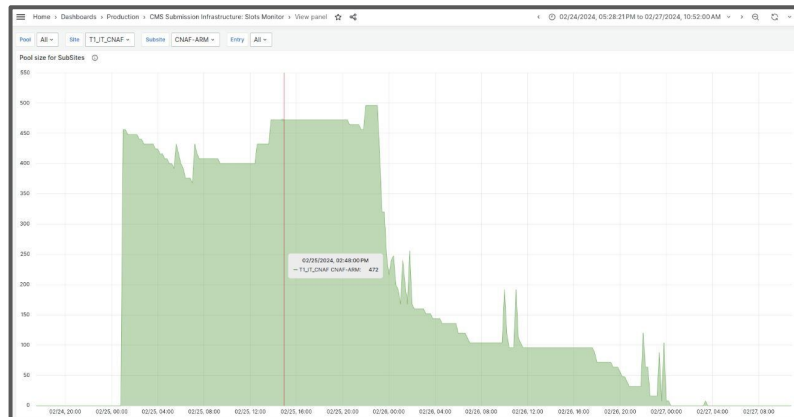
- Minimize the effort on CMS and simplify the site admins life
- At the moment data are accessed via Xrootd protocol. The plan is to provide direct access also via GPFS.

**Technical validation fully done.**

- Being the first allocation at a Tier site, the CNAF nodes where essential to finally validate that CMS Computing is multi-architectures enabled.

**Physics validation in progress:**

- most of the subsystem reports green light (especially when looking at MC).
- While some discrepancies spotted on DATA require further analysis for a better understanding
- This step has been carried on both at CNAF and on Glasgow temporary allocation

# CMS: Next steps

Interpretation of the validation results is not straightforward. Although most of the reports green for MC. Few failures for data need better understanding

- Moreover one of the issue running in Glasgow was that non-negligible fraction of failed jobs on ARM resources for data workflow. Statistics is missing

The main priority now is to finalize the validation thanks to the ARM allocation @CNAF, currently the only stable one at CMS



**C-RSG Question #7: ARM Performance**

We would really need to run real production workflows to answer this question. Fortunately, there are now ARM resources (~1K cores) newly available at CNAF. The core software team is being offered access. Thank you!

**[CMS-7] Page 13: Are the CPU efficiency of ARM based machines different from IA64 machines? And are there any notable learnings wrt. performance.**

CMS physics validation workflows, which were run on the ARM resources at Glasgow, are not optimized for CPU efficiency. We do not have any measurements running standard production workflows on ARM resources. We are, however, currently running a small validation campaign on ARM resources at the CNAF Tier-1 site. The results of the physics validation are not yet complete.

*We had an interesting discussion with the chief LHCC Referee about accessing ARM resources in **commercial clouds** as a way to perform validations on new architectures without waiting for our sites to purchase new machines. Are we (DRP?) interested in revisiting this capability? See two ATLAS CHEP presentations:*

- https://indico.jlab.org/event/459/contributions/11636/
- https://indico.jlab.org/event/459/contributions/11553/

# Outlook and future plans

- An important piece is missing: measurement of **power consumption**
  - Define a strategy for the measurements and a good metric (e.i. evnt/sec/watt ?)

- The validation is on the way..
  - LHCB is performing test on Glasgow resources, CNAF will be used in next phase for DIRAC provisioning

- The use case would enable production workflows on ARM at the INFN CNAF center, a Tier-1 in the WLCG hierarchy and a node of the ICSC datalake

- On top of this, a successful validation would pave the way to more communities (and in particular smaller HEP and Astro experiments)

- In perspective, the path could enable the Italian Computing Infrastructure (INFN and ICSC, in this particular case) to provision ARM machines in the near future, with a sizable reduction of computing TCO.

Thank you!