



ALICE experiences with ARM

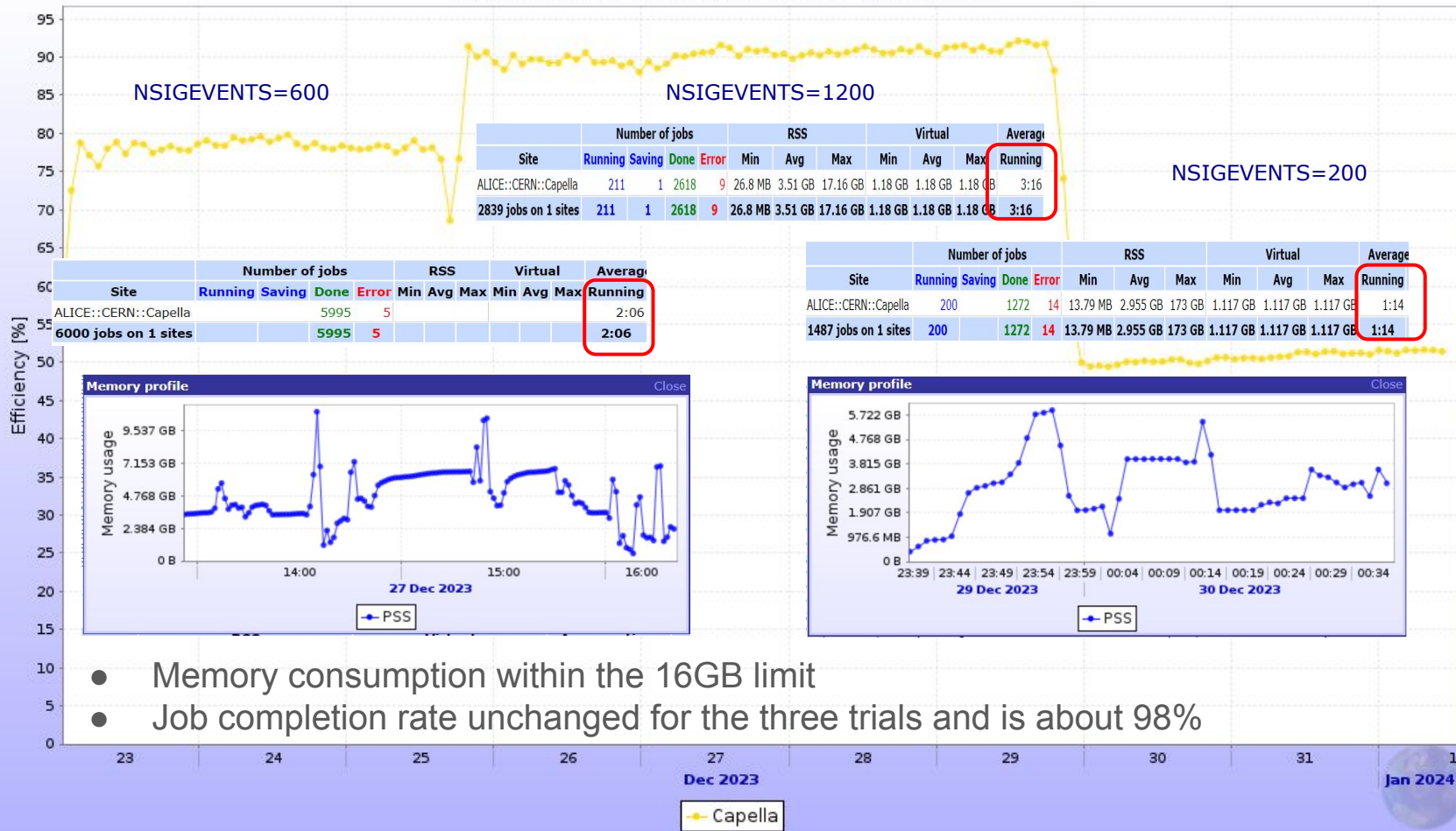
L. Betev

Acknowledgement and disclaimers

- Many thanks to the UKI-SCOTGRID-GLASGOW site (Dwayne Spiteri and colleagues) for the super-stable setup provided for the ALICE first tests of ARM processing on a large scale
- What was done
 - Extensive test of MC simulation jobs - comparison of run times, efficiency, memory use
- What wasn't done
 - Data processing and analysis workflows
 - Multiple software releases (this will become important later)



Jobs efficiency (cpu time / wall time)



- Memory consumption within the 16GB limit
- Job completion rate unchanged for the three trials and is about 98%

Jobs efficiency (cpu time / wall time)

NSIGEVENTS=200

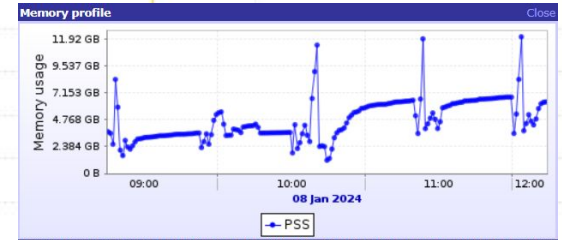
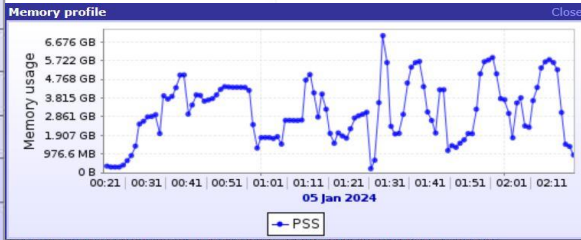
NSIGEVENTS=10

NSIGEVENTS=1200

Site	Number of jobs				RSS			Virtual			Average
	Running	Saving	Done	Error	Min	Avg	Max	Min	Avg	Max	Running
ALICE::CERN::Capella	198		2175	84	13.13 MB	2.939 GB	188.7 GB	8 KB	231.1 MB	1.126 GB	2:05
2457 jobs on 1 sites	198		2175	84	13.13 MB	2.939 GB	188.7 GB	8 KB	231.1 MB	1.126 GB	2:05

Site	Number of jobs				RSS			Virtual			Average
	Running	Saving	Done	Error	Min	Avg	Max	Min	Avg	Max	Running
ALICE::CERN::Capella	220		293		173.5 MB	3.921 GB	17.05 GB	4 KB	17.89 MB	54.15 MB	5:02
513 jobs on 1 sites	220		293		173.5 MB	3.921 GB	17.05 GB	4 KB	17.89 MB	54.15 MB	5:02

Efficiency [%]



- Memory is not a problem on ARM (?)
- MC can be run in almost any configuration
- TF generation comprises about 2/3 of the time @20% efficiency

Site	Number of jobs				RSS			Virtual			Average
	Running	Saving	Done	Error	Min	Avg	Max	Min	Avg	Max	Running
ALICE::CERN::Capella	188		993	88	16.39 MB	2.477 GB	7.677 GB				1:26
1271 jobs on 1 sites	188		993	88	16.39 MB	2.477 GB	7.677 GB				1:26

Capella



Comparison of p-p, high event number MC job behaviour

x86: 65% job efficiency, 78% job success rate

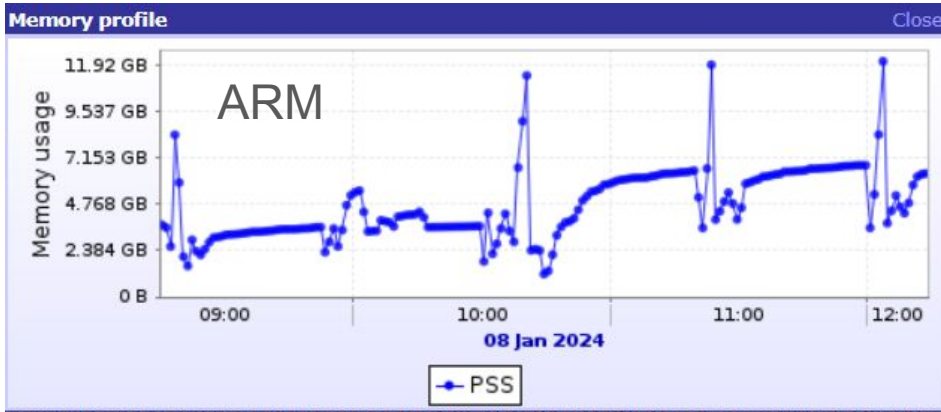
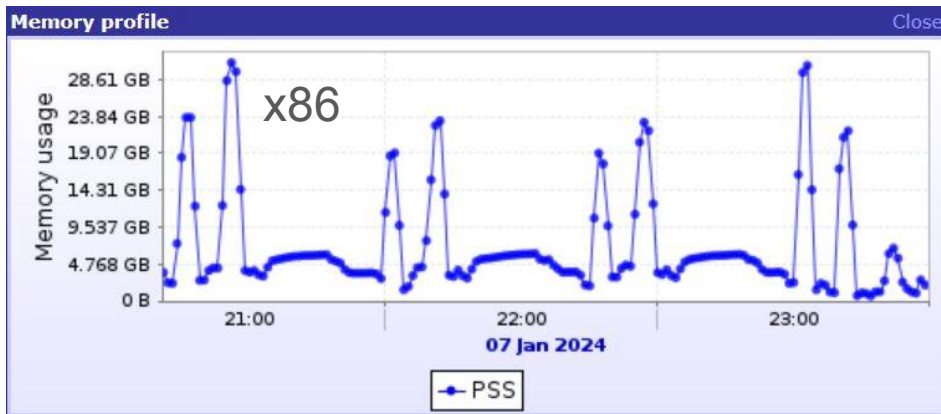


ARM: 50% job efficiency, 99% job success rate



- Average running time on x86 (14HS23/core) is ~20% faster than ARM Ampere Altra Q80-30 (80 cores);
- Different ARM memory management avoids the high memory spikes, which kill a large fraction (~20%) of all jobs on average x86 site, while on ARM 99% of jobs are successful
- CPU efficiency on x86 site for the same type of job is 15% higher than on ARM, both efficiencies are low for MC-type jobs (this is a code issue)

Comparison of Pb-Pb, high event number MC job behaviour



- x86 MC is hampered by memory spikes and is thus limited in number of TFs - less probability of job being killed
- 8% job success rate (all other killed for memory overuse)
- The successful jobs ran only at sites with very generous memory allocation policies (40GB+/8 core slot)
- MC on ARM can be of arbitrary length and duration, up to the 24TTL limit, with what appears to be a slight memory leak
- Well below the 16GB memory limit of the job slot

Operation on ARM

- Software build and test - ALICE standard CI (aarch64)
 - Multiple OS versions (EL7, AlmaLinux9, Ubuntu24)
 - Distribution through CVMFS
- Containers for payload execution
 - Same as above
- Detection of platform and CPU type
 - Transparent setup of SW build and container - users do not differentiate between CPU types - same treatment as AMD, Intel CPU
 - However users can specifically select ARM
- Automatic detection and setup is the only viable option for adding new CPU types to the Grid sites

Issues and further tests

- After the December 2023 successful test
 - Subsequent software release produced a (quite unhelpful) segmentation fault, only on ARM and only with a specific (the default for Alma9) kernel version
 - We now have a reproducer and know with a relatively high level of certainty where the issue is
- Currently testing the behaviour of data processing and analysis jobs
 - ARM resources available at CNAF, CERN and ALICE local
 - Should gain experience and understanding of (subtle) differences for wider adoption

Conclusions and recommendations

- Current status of ALICE tests on ARM
 - Very good performance of the CPU, compatible with the average HS23 on the Grid
 - Possible improvements of memory management observed
 - Integration of ARM clusters/queues into the ALICE build and Grid machinery is complete and was not problematic
 - Subtle differences of code behaviour on ARM point to the need for more and detailed tests before wide and uncontested adoption
- Management of ARM machines/queues on the sites
 - We recommend that ARM machines are either segregated behind a separate CE or
 - In a mixed queue there is an option, which enables/disables ARM hosts