



ATLAS usage of ARM

Frédéric Derue, LPNHE Paris

LCG France meeting, CC-IN2P3 Lyon 28st november 2023

Largely based on presentation of J. Elmsheuser (BNL) at CHEP2023 (<u>link</u>) + recent presentations at ATLAS Software & Computing week



Access to ARM resources

• CERN (usually single instances)

- Ixplus8/9-arm (in Oracle Cloud)
- ATLAS build machine provisioned by CERN IT and Openlab on prem and in Oracle Cloud (Neoverse and Cavium ThunderX2)
- AWS (Graviton2 and 3)
 - fully integrated in ATLAS PanDA/Rucio based production system
- Google cloud (Ampere Altra/Neoverse-N1)

Current and future large scale resources

- Fugaku HPC at Riken Center for Computational Science, Japan is #2 in TOP500 supercomputer list of November 2022 (link)
- \circ a few other HPCs plan to have ARM partitions
- grid sites are starting to be interested in ARM when the experiment software is validated (e.g Glasgow, CNAF)

J. Elmsheuser (BNL) at CHEP2023 (link)



ATLAS Panda and Rucio setup on AWS



J. Elmsheuser (BNL) at CHEP2023 (link)



- self-container, cloud-native, vendor-agnostic, auto-scaled, auto-healing
 - Rucio
 - object store + signed URL + http protocol
 - integrated in Rucio and WLCG MW

∘ PanDA

- Harvester submitting to K8S directly
- CVMFS K8S plugin
- Frontier squid either in K8S cluster or as separate load balanced VM group
- auto-scaling : no jobs, almost no cost

- Full aarch64 grid setup available with OS container, middleware, Kubernetes etc.
- More details in talks from Fernando Barreiro (<u>link</u>) and Mario Lassnig (<u>link</u>) at CHEP 2023 conference

ATLAS setup in Glasgow







ATLAS Panda setup in Glasgow





ATLAS usage of ARM, LCG France meeting, 28/11/2023



ATLAS testing in Glasgow



ATLAS Testing

- The ARM nodes are running a job starting with 50M events and with simulation, reconstruction and derivation tasks.
- Running since Late August, and still running!

Evgen	Evgen Merge	Simul	Merge	Digi	Reco	Rec Merge	Atlfast	Atlf Merge	TAG	Deriv	Deriv Merge		
												Replace empty	
+ 601	229/mc.P	hPy8EG	_A14_ttba	ar_hdam	p258p75	_SingleL	ep.py					mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00	0
(Fullsim)(0)PhPy8E	G_A14_t	tbar_hda	mp258p	75_Single	eLep						events: 5000000	Cloned
		s4159	s4114		r14799	r14811		1		p5667	p5669	partially_submitted	edit (saved
T:		running	running		aborteo	d aborte	d					Produced events: 39370000	
+ 601	229/mc.P	hPy8EG	A14_ttb	ar_hdam	p258p75	_SingleL	.ep.py					mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00	C
(Fullsim)(0)PhPy8E	G_A14_t	tbar_hda	mp258p	75_Single	eLep						events: -1 (5000000)	Cloned
	-	1	s4114	-	r14799	r14811		1	1	p5667	p5669	partially_submitted	edit (saved
T:					running	aborte	d					Produced events: 36570000	
+ 601	229/mc.P	hPy8EG	A14_ttb	ar_hdam	p258p75	SingleL	ep.py					mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00	C
(Fullsim)(0)PhPy8E	G_A14_t	tbar_hda	mp258p	75_Single	eLep						events: -1 (5000000)	
	-	1	-	-	r14799	r14811	-	1	1	p5667	p5669	partially_submitted	edit (saved
T:						running	3					Produced events: 35740000	
+ 601	229/mc.P	hPy8EG	_A14_ttba	ar_hdam	p258p75	_SingleL	.ep.py					mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00	0
(Fullsim)(0)PhPy8E	G_A14_t	tbar_hda	mp258p	75_Single	eLep						events: -1 (5000000)	
	-	1	1	-	-	r14811	-	-	1	p5822	-	submitted	edit (saved
T				-					-	running	a	Produced events: 34520000	

34412357 t	ask														E	stimated C based on the regional	O ₂ emissio e weighted electricity TLAS Grid	n in grams, average of sources for computing
Task ID	Campaign	Request	Туре	Processing type	Working Group	User	Nucleus	Status	N input files finished failed	N input events used	N output events	HS06*sec Expected Total done failed	Time stamps: created last modified	Cores	Priority: original current	Attempt	Tracker	gCO ₂ total done failed
34412357	MC23c	50572	prod	simul	AP_MCGN	janders	UKI-SCOTGRID- GLASGOW	running	25,000 (94%) 23,564 (0%) 11	50,000,000 (94.3%) 47,128,000	47,128,000	54,150,000,000 37,482,088,818 36,757,221,900 724,866,918	2023-08-18 12:45:09 2023-10-03 11:13:09	8	220 341	0	JIRA	3,128,346 3,065,700 62,646

Nightly builds on ARM/aarch64



ARM	master_AnalysisBase_aarch64-centos7-gcc11-opt	2023-04-18T0220	18-APR 10:34	0 (1)	0 (0)	N/A	N/A	18-APR 11:11
ARM	master_Athena_aarch64-centos7-gcc11-opt	2023-04-18T2300	19-APR 06:16	0 (1)	4 (4)	N/A	N/A	19-APR 09:21
ARM	master_AthSimulation_aarch64-centos7-gcc11- opt	2023-04-18T2101	18-APR 21:37	0 (0)	1 (1)	N/A	N/A	18-APR 22:11
ARM	master_DetCommon_aarch64-centos7-gcc11-opt	2023-04-18T2000	18-APR 20:02	0 (0)	0 (0)	N/A	N/A	18-APR 20:21

- 4 nightly builds for Athena, AthSimulation, AnalysisBase and DetCommon (<u>link</u>) projects fully integrated in standard ATLAS build system and available on CVMFS
- Selected stable Athena releases like 23.0.3 and 23.0.14 installed on CVMFS used in physics validation
- since December 2022 dedicated build machine provided by CERN IT hosted in Oracle Cloud, 40 vCPUs (Ampere Altra A1), 250 GB RAM, special manual CentOS7 installation
- used before techlab-arm64-thunderx2-01 provided by CERN OpenLab and shared with CERN SFT
- LCG_102b_ATLAS_* / LCG_103 layers and tdaq/tdaq-common builds available and used in ARM/aarch64 builds
- can easily build stand-alone docker/podman container for e.g. AthSimulation

ARM/aarch64 ATLAS builds and caveats



J. Elmsheuser (BNL) at CHEP2023 (link)

• Build flags

- using Armv8 defaults (gcc 11.2 allows up to armv8.6-a, gcc docu link)
- no special arch option and no special compilation flags set apart from different linker flag in max-page-size
- able to build Athena/AthenaExternals with clang16 as well

• Floating Point Exception (FPE)

- Athena FPE auditor code not working on ARM/aarch64 since it uses x86 specifics
- discussions about different FPE behaviours of x86_64 vs. ARM/aarch64 compilers: StackOverflow 1, StackOverflow 2, ARM compiler

Potential numerical differences

 \circ e.g. due to different floating point libraries used - see StackOverflow link \circ some fluctuations in physics objects at the level of (10⁻⁴ – 10⁻⁶) \circ NB. small numerical differences or Intel vs. AMD if IntelMathFunction used

Reconstruction and physics validation



• using Athena 23.0.14

ATLAS

J. Elmsheuser (BNL) at CHEP2023 (link)

- March 2023 successfully passed reconstruction physics validation
- standard procedure of 13 MC physics processes (100k events each), no pile-up
- Digitization+Reconstruction step on AWS Graviton2 PanDA queue (130 jobs, 215 GB output, 1.5 days, 450 USD in total)
- compared with same events produced on x86_64 grid sites all subsequent steps (merging, histogramming) on x86_64 - task speed comparison in the backup









Perfect agreement

Example plots from calorimeter cluster validation

HEPScore integration and task speed



wl-scores





Ratio ARM vs. x86 tasks for walltime/event per physics task

J. Elmsheuser (BNL) at CHEP2023 (link)

- HepScore numbers for ATLAS
 reconstruction
- Node fully packed with n×4 threads
 - → different number of available cores explains different scales !
- Orange/Brown/Blue: ARM flavours
- Green: Intel/AMD flavours
- Some hyperthreading or IO related differences
- Similar trends for Event generation and simulation

Ratio ARM vs. x86 Grid job averages:

- Walltime/event: 0.92
- CPU time/event: 0.81
- HS06sec/event: 0.79
 (with ARM HS06=10)

HEPScore measurements



D. Britton, Glasgow

Machine	CPU	Threads	HEPScore	
Test AMD	1x EPYC 7643 HT	96	1359	
Std AMD WN	2x EPYC 7513 HT	128	1887 _	
AMD Bergamo	2x EPYC 9754 HT	512	7497 🔦	
Test ARM	1x ALTRA 80	80	1513	→
Farm ARM	2x ALTRA 80	160	2619 ?? 🔺	<u>Not x2 ??</u>
Ampere Max	1x ALTRA Max	128	2065	



HEPScore/Watt



Machine	СРО	Threads	HS/	D. Britton, Glasg			
			Watt				
Test AMD	1x EPYC 7643 HT	96	3.7				
Std AMD WN	2x EPYC 7513 HT	128	4.2				
AMD Bergamo	2x EPYC 9754 HT mo		6.3	50% improvement from 7513 WN			
Test ARM	1x ALTRA 80	80	5.6				
Farm ARM	2x ALTRA 80	160	5.2 ??	25% improvement from 80c			
Ampere Max	1x ALTRA Max	128	7.0	25% improvement nom ooc			

Note added: We use Average Watts to calculate HS/Watt. Other people also use Max Watts, so care needs to be taken comparing numbers.



Energy efficiency vs clock







- Running one step below max clock frequency seems to be the sweet-spot for all CPUs tested.
- On ARM128, frequency reduction lowers HEPScore by 6.5% but saves another 9% of the energy.
- No data on Bergamo (bug?)



ATLAS usage of ARM, LCG France meeting, 28/11/2023



• Further plans:

[©]∕ATLAS

migrate to AlmaLinux9 (all pieces already available)

- investigate further production workflows (e.g. MC generators)
- \circ power consumption tricky to assess on virtualized hardware

• Summary

- full ARM/aarch64 ATLAS nightly builds on CVMFS through regular nightly ATLAS build system
- Geant4 and reconstruction physics/technical validation successfully passed using AWS through PanDA/Rucio production system
- Athena,23.0.3 sherpa, simulation and reco workflows integrated into new HepScore benchmark

• Energy efficiency

 128-core ALTRA-Max is most energy efficient CPU tested by Glasgow
 AMD has made big improvements to energy efficiency with Bergamo.
 running CPU one step below max-frequency seems to be sweet-spot in terms of energy efficiency

Project is part of ATLAS HL-LHC R&D projects and will be part of ATLAS computing technical design report for HL-LHC