Glasgow's ARM cluster: **Experience and Findings**

Prof David Britton Dr Gordon Stewart Dr Samuel Skipsey Dr Emanuele Simili Dr Bruno Borbely **Dr Dwayne Spiteri**









Journées LCG-France, CC-IN2P3 - 29/11/2023





GridPP and Glasgow

 GridPP is the UK's branch of the WLCG. Its 18 Tier-2 sites and 1 Tier-1 (RAL) are split into 4 groups (ScotGrid, NorthGrid, SouthGrid and LondonGrid).

Glasgow is the largest Tier-2 site in ScotGrid.









 After hearing about our work, AMPERE got in touch with us and kindly agreed to donate some ARM servers to us.

CPU	Ampere Altra Q80-30
Sockets	2
Max Frequency	3.0 GHz
Memory	512 GB

 With several 2U servers we now have 1760 active cores.





How do we get VO's to see and use our ARM resource?

 This is a small visualisation of our x86 compute, which we this as a model for our ARM resources.



 \approx 1000 x86 cores **x86**

> Homogeneous x86 resources, managed by a single queue.



4

How do we get VO's to see and use our ARM resource?



Glasgow ARM Cluster: Experience and Findings



Technical setup differences between ARM and x86

- the rest of the site remains at Centos7 / HTC8 (this will change in the near future).
- EL7 to EL9. By and large the ARM nodes are treated the same as x86 ones.
- For Provisioning, we currently use a bespoke PXE- / kickstart-based system (think Cobbler but targeted a little more closely at our needs) and that works across all our systems.
- Configuration management is via Ansible, and works across all hosts.
- possible through Ansible (SSH and Python underneath).

• Our ARM nodes are on Rocky9 / HTCondor 10, attached to their own CE and Condor manager, while

• While there are configuration changes between x86 to ARM, some changes are due the move from

• Automatic updates are run across the site, with a couple of exceptions, and this is also true of the ARM nodes - patching normally involves checking applied updates rather than pushing them out. All



Setup: ATLAS-side

PanDA setup stores data locally (now)

Points to our temporary CE and our condor_arm queue in the ARC

aarch architecture only

• Now we just have to run some proper jobs on them

PanDA Queue: UKI-SCOTGRID-GLASGOW_ARM

Object details

PanD Desc Default Core Core Coree

State

Object State con Last modification

Associated DD

ODMEndpoint

UKI-SCOTGRID-0

IKI-SCOTGRID-0

Showing 1 to 2 of 2

Associated Qu



Showing 1 to 1 of

Associated Pa



Dwayne Spiteri, University of Glasgow

ID	930										
Name BanDA Site	0	UKI-SCOTGRID-GLASG	OW_A								
Description	ด	UNI-SCOTGRID-GLASG	000-0	CFN							
Default object	õ	UKI-SCOTGRID-GLASG	ow_v								
Corecount	8	overwritten									
Corepower	0	8.74 inherited.rcsite									
Coreenergy	0	10.0 Inherited.rcsite									
Object state	0	ACTIVE									
dification date	202	ect was cloned from UKI 3-08-11 13:26:51.509877	-SCO1 7	GRID-0	iLASGOW_C	EPF	I VIA WEDUI				
	noi	ate									
	poi	115									
										Searah	
										Search.	
Endpoint			11	Туре	Jt	0	Experiment Site	J1	Activities		
TGRID-GLASGO	W-C	EPH_DATADISK		DATAD	DISK	UK	KI-SCOTGRID-GLASGOW		pl/0, read_lan/0	, write_lan/0	
TGRID-GLASGO	W-C	EPH_SCRATCHDISK		SCRAT	TCHDISK	UK	KI-SCOTGRID-GLASGOW		write_lan_analysi	s/0, read_lan/1	, write_
to 2 of 2 entries											
ted Overver											
ted Queues											
										Search:	
	-	10			• • • •				0	A	<u> </u>
fa Ta	QL	ieue 🚛 flavour 🚛	versi	on 💵	O qstatus	ΨI.	site	1I	Coutime	WClock I	() ETF
a.scotgrid.ac.uk	со	ndor_arm>ARC-CE	None	1	production		UKI-SCOTGRID-GLASG	WC	0	0	×
to 1 of 1 entries											
ted PandaQu	ieue	e architectures									
A bit				Vanda			la chu chi c no		Madal		
Archit	ectu	69		venuo	Л		manucuona		Model		
aarche	4) ex	cl									(
Glasgow ARI	A CI	uster: Experience ar	nd Fir	ndings	; ;						



Site Monitoring

- Monitoring based around Prometheus. Managed to find version of it compiles for aarch64.
- Data relayed from ARM nodes into the same Grafana dashboard as the rest of the site with minimal changes.
- Some metrics are IPMI based here some changes were required, but technical issues are similar to the ones whenever new hardware is installed.



Dwayne Spiteri, University of Glasgow

8

ATLAS Testing - ARM Work

- The ARM nodes ran jobs starting with 50M events and with simulation, reconstruction and derivation tasks.
- Meant to run for about ~21 days, ended up taking 58 days.



34427583 1	task															E	stimated CC based on the regional e AT	02 emission weighted electricity : TLAS Grid	n in ave sou coi
Task ID	Campaign	Request	Туре	Processing type	Working Group	User	Nucleus	Status	N input files finished	N input events used	N output events	HS23s Expected Total done failed	Time stamps: created last modified	Cores	Priority: original current	Attempt	Parent	Tracker	d fi
34427583	MC23c	50572	prod	merge	AP_MCGN	janders	RAL-LCG2	finished	24,989 (100%) 24,989	49,978,000 (100%) 49,978,000	49,978,000	187.0M 179.3M 166.6M 12.6M	2023-08-21 07:10:11 2023-10-07 15:49:31	1	<mark>885</mark> 900	0	34412357	JIRA	1 1 1



ngleLep.py	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_
0	events: 50000000
4811 p5667 p5669	partially_submitted
orted	Produced events: 49978000
ngleLep.py	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_
o contraction of the second seco	events: -1 (5000000)
4811 p5667 p5669	partially_submitted
orted	Produced events: 49970000
ngleLep.py	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_
0	events: -1 (5000000)
4811 p5667 p5669	partially_submitted
ished	Produced events: 49970000
ngleLep.py	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_
D	events: -1 (5000000)
4811 p5822	submitted
finished	Produced events: 49970000









Dwayne Spiteri, University of Glasgow





Dwayne Spiteri, University of Glasgow



Dwayne Spiteri, University of Glasgow











Dwayne Spiteri, University of Glasgow



ATLAS Testing - x86 Work

- Using the same s-, r- and p-tags as the ARM work, jobs were also sent our x86 nodes to get some metrics for direct comparison.
- Only 200k events were sent in this way, site still live and taking our usual mix of work from other VO's
- Took ~26 days to complete the full chain.



Task ID	Campaign	Request	Туре	Processing type	Working Group	User	Nucleus	Status	N input files finished	N input events used	N output events	HS06*sec Expected Total done failed	Time stamps: created last modified	Cores	Priority: original current	Attempt	Tracker	gC do fai
34714849	MC23c	50572	prod	simul	AP_MCGN	janders	SWT2_CPB	done	100 (100%) 100	200,000 (100%) 200,000	200,000	223,400,000 175,231,218 174,063,834 1,167,384	2023-09-07 07:45:09 2023-09-08 11:05:42	8	220 900	0	JIRA	15, 15, 10,

34714849 task

ingleLep.py ep	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_ events: 200000 (50000000)
14811 p5822 p5669	partially_submitted
one done	Produced events: 200000
ingleLep.py ep	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_ events: 200000 (50000000)
14811 p5667 p5669	partially_submitted
	Produced events: None
ingleLep.py ep	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_ events: 200000 (50000000)
14811 p5667 p5669	partially_submitted
	Produced events: None
ingleLep.py ep	mc23_13p6TeV.601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_ events: 200000 (50000000)
14811 p5667 p5669	partially_submitted
	Produced events: 200000



Site-side metrics

- The ARM nodes have been largely exclusively running ATLAS arm work.
- Didn't have IPMI reporting power to Grafana when this was started.
- Had to use exported metrics from running machines every minute (timestamp, temperature, instantaneous) power, frequency) to estimate the power we've dissipated as a site running this work.
- Many different ways of evaluating power could be used. Maximum power is reported regularly, but we feel it's not the most representative figure.
 - To compare to figures we have for x86, I'll estimate the peak plateau power

UNIXTIME(s)	CPU1TEMP(oC)	CPU2TEMP(oC)	POWER(W)	FREQUENC
1695827622	84.000	83.000	544.000	3.00
1695827682	83.000	85.000	544.000	3.00
1695827742	83.000	84.000	544.000	3.00
1695827803	82.000	84.000	528.000	3.00
1695827863	83.000	82.000	520.000	3.00
1695827923	81.000	82.000	504.000	3.00
1695827983	80.000	81.000	520.000	3.00
1695828043	81.000	80.000	512.000	3.00
1695828103	80.000	79.000	520.000	3.00
1695828163	81.000	81.000	528.000	3.00
1695828223	82.000	81.000	528.000	3.00
1695828283	79.000	81.000	512.000	3.00
1695828343	81.000	82.000	520.000	3.00
1695828403	83.000	83.000	552.000	3.00
1695828463	80.000	82.000	528.000	3.00
1695828523	80.000	83.000	536.000	3.00
1695828583	81.000	82.000	520.000	3.00
1695828643	81.000	81.000	512.000	3.00
1695828703	81.000	83.000	528.000	3.00
1695828763	82.000	81.000	536.000	3.00
1695828823	78.000	82.000	512.000	3.00
1695828883	80.000	81.000	504.000	3.00
1695828943	80.000	80.000	488.000	3.00
1695829003	80.000	79.000	480.000	3.00
1695829063	79.000	79.000	480.000	3.00



Reported power per node

- IPMI Power reporting on some nodes was not working correctly, small hacks required.
- Not able to use this to measure power draw for ATLAS work, but can estimate.
- For nodes where the reporting works, power draw looks consistent with what we expect.
- Fixed for future runs with automatic reporting







Power with Time for ATLASARM012 df 500 Ñ 400 de Rep 100 -12-10-2023 16-10-2023 08-20-2023 2809-2023 04-2023 04-20-2023 Time









Frequency vs Temperature

- During RAL's fixed-frequency ARM tests, the frequency was altered automatically to cope with the servers operating at high temperatures.
- This behaviour is not seen at our site. Temperature and **Frequency** seem uncorrelated.
- Samples were only taken once a minute, phenomena could potentially only occur at higher resolutions.



24

Comparison Methodology

 Use PanDA data taken on the ATLAS-side to estimate usage metrics.

- Estimate how many cores were used for how long to calculate the average CPUhours taken for each type of ATLAS job (area of orange to the right)
- Then use a combination of reported and measured data to calculate HEPScore/Watt

States of job	os in this	task [drop mo	ode]				
р	ending		defined	d	wait	ing		
Run								
lah liat Cart	hu Dondo				ata akan			
JOD list. Sort	by Panda	iD, tim	e since i	last st	ate chang	ge, ascend		
PanDA ID Attempt#	Owne Group	r)	Reque Task I	est D	Transfo	ormation		
	jander AP_M	's CGN	50572 34714	849	Sim_tf.p	ру		
5957143510 Attempt 2	Job na	ıme: <mark>m</mark> a	c23_valid	1.6012	29.PhPy8	EG_A14_tt		
	Datase Out: m	ets: In: i nc23_va	mc23_13 alid.6012	p6TeV 29.Phf	/:mc23_13 Py8EG_A1	Bp6TeV.601		
;	1							
Job name: n	n 23_vali	d.6012	229.PhP	y8EG	_A14_ttb	ar_hdam		
PanDA ID		Own	er	WG		Req Tasl		
5956895785		jand	ers	AP_	MCGN	505 347		
Datasets:		In: mc23_13p6TeV:mc23_13p6TeV.6 Out: mc23_valid.601229.PhPy8EG_						
Files summa	iry:	input: 1, size: 579.80 (MB); log: 1; o						
Logs		Go	o to	-				

ed th	nrottled	activated	t l	sent	starting	running	holding	transferring	merging	fir	nished	failed	cancelled
										10	10 4		
										CONSTRUCTION OF			
creation time, descending creation time, ascending mod time, descending mod time, priority, attemptor, ascending duration, descending duration													
Created	Tin sta d:h	ne to art a:m:s	Duration d:h:m:s	Mod		Cloud Site				Priority	N input eve files)	nts (N input	Max PSS/c GB
2023-09-08 02:36:02	8 0:0	:04:33	0:4:46:37	2023-0 10:27:0	19-08 V 09 r	WORLD UKI-SCOTGF rules defined	RID-GLASGOW_CE	EPH online no active black	listing	900	2000 (1)		0.34
	ed th ion time, desce c Created d 2023-09-0 02:36:02	ed Introttled throttled	ed throttled activated	ed throttled activated Interval activated Interval a	ed throttled activated sent in throttled Interview in throttled Interview Interview Duration d:h:m:s Duration d:h:m:s Mod 0:0:0:4:33 0:4:46:37	ed throttled activated sent starting in time, descending creation time, ascending mod time, descending mod time, descendi	ed throttled activated sent starting running is ent is ent starting running is entime, descending mod time, descending mod time, priority, attemptor, is entime, descending mod time, priority, attemptor, is entime, descending mod time, descending mod time, priority, attemptor, is entime, descending mod time, priority, attemptor, is entime, descending mod time, descending mod time, priority, attemptor, is entime, priority, attemptor, is entit, priority, attemptor, is entime, priori	ed throttled activated sent starting running holding In the second in time, descending creation time, ascending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In the tom time, descending creation time, ascending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending creation time, ascending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending creation time, ascending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s In time tom time, descending mod time, descending mod time, descending mod time, priority, attemptor, ascending duration dith:m:s	ed throttled activated sent starting running holding transferring in throttled in throttled activated sent starting running holding transferring in throttled in time, descending creation time, ascending mod time, descending mod time, descending mod time, priority, attemptint, ascending duration, descending duration in Created Time to start din:m:s Duration dth:m:s Mod Cloud Site Cloud Site d $2023 - 09 - 08$ $0:3::0^2$ $0:0:04:33$ $0:4:46:37$ $2023 - 09 - 08$ $10:27:09$ WORLD UKI-SCOTGRID-GLASGOW_CEPH online no active black	ed throttled activated sent starting running holding transferring merging in throttled is throttled activated sent starting running holding transferring merging is throttled Image: Sent three sentimes according mod time, according mod time, according mod time, priority, attemptor, ascending duration, descending duration desc	Image: Problem of the second se	Image: Problement of the starting of the second start of the starting of the start of	Image: Normal line line line line line line line lin

oar_hdamp258p75_SingleLep.simul.e8514_e8528_s4159.5956895702 #2

229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00 || Rucio link amp258p75_SingleLep.simul.log.e8514_e8528_s4159_tid34714849_00

o258p75_SingleLep.simul.e8514_e8528_s4159.5956895785

uest ID	Status	Туре	Transformation	Created Last modified	Time to start Duration [d:h:m:s]	Site	Harvester instance Worker ID	Cores	Priority
2 4849	finished	simul	Sim_tf.py	2023-09-07 20:40:37 2023-09-08 11:05:27	0:1:35:29 0:12:18:23	UKI-SCOTGRID-GLASGOW_CEPH	CERN_central_B 468057624	8	900

601229.PhPy8EG_A14_ttbar_hdamp258p75_SingleLep.merge.EVNT.e8514_e8528_tid33116249_00 _A14_ttbar_hdamp258p75_SingleLep.simul.log.e8514_e8528_s4159_tid34714849_00

utput: 1; pseudo_input: 1





Estimates of HEPScore/Watt



- The x86 bars here are weighted averages from the different x86 machines we have at our cluster.
- HEPScore/Watt better for ARM for reconstruction and simulation work
- ARM load order of magnitude larger, expect merge steps (I/O bound processes) of jobs to more CPU intensive relative to the jobs for x86.







CMS Testing

- CMS currently running validation tasks on output datasets.



• Power Reporting in Grafana working at this point - Expect average usage ~6.2kW for efficient usage

Power Consumption breakdown







Secondary Testing

- We've done some provisional benchmarking with the ARM Farm and compared it to other arm and x86 machines.
- Found that the dual socket machine wasn't twice as performant as the single socket (orange to red).
- Cache coherency overheads on the ARM dual-CPU were at times very high compared to other Intel/AMD CPUs.
- The ARM CPU's have two different protocols for cache coherency: in-chip and between-chip. Transferring between these makes it less efficient. (for further reading see this article)
- Bergamo outperforms all other systems in terms of HEPScore but...





Tes Std

AM

Test

Far

Am

David Britton's Talk: Carbon and Sustainability in ALTAS S&C Week - 2nd October

HEPScore Measurements

Machine		CPU	Th	reads	HEPScore			
t AMD	1x EPYC	. 7643 HT		96	1359			
AMD WN	2x EPYC	. 7513 HT		128	1887	,		
D Bergamo	2x EPYC	. 9754 HT	1	512	7497	← ^{X4}		
t ARM	1x ALTR	A 80		80	1513			
m ARM	2x ALTR	A 80		160	2619 ??	← Not	: x2 ??	
pere Max	1x ALTR	A Max		128	2065			
	_				HEP-Score			
		8,000 -						
		7,000 -						
		6,000 - O						
		5,000 -						
		± 2,000 -						
		1,000	T					
		0 -	lest 1x7643	WN: 2*75	Bergamo	lest Altra80	ARM-Farm	Altra
itton, University of Gla	sgow		AMD - 96HT	AMD - 2*6	4HT AMD - 2*256HT	ARM - 80c	ARM - 2*80c	ARM





Future Glasgow Farm Plans What will we fill our racks with? (hyperlink)

... as it consumes a lot of power, it's less efficient at providing work per unit of power.

Bergamo vs Altra Max! 4 AltraMax CPUs and 2 Bergamo CPUs have comparable numbers of cores (512) and should draw roughly the same amount of power. ARM would still be more energy efficient.

 In this configuration Bergamo would take up less rack space. So you can get more threads per rack with the Bergamo (in theory).



Da



ow max (W)	Avg. Pow (W)	Machine	CPU	Threads	HS/Watt	
451	367	Test AMD	1x EPYC 7643 HT	96	3.7	
512	449	Std AMD WN	2x EPYC 7513 HT	128	4.2	
1347	1198	AMD Bergamo	2x EPYC 9754 HT	512	6.3	50% improvement from 7513
348	270	Test ARM	1x ALTRA 80	80	5.6	
381	294	Farm ARM	2x ALTRA 80	160	5.2 ??	
664	502	Ampere Max	1x ALTRA Max	128	7.0	25% improvement from 80c
					HEPScore/Watt (∝ to Events/Joule)

David Britton's Talk: Carbon and Sustainability in ALTAS S&C Week - 2nd October



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



Future Glasgow Farm Plans What will we fill our racks with?

Data-centres are usually limited not by space, but by power^{1,2}. To fill a rack with Bergamo would tax our power budget before we filled the space. ARM servers are better in this regard.

High core density could mean that cooling may become an issue. Our racks are designed to be low density (<10kW per rack), and a Bergamo machine may leave more gaps.

 What you choose depends on your data-centre design and what metrics you value. For energy efficiency ARM seems to be better.





Dav

Pow max (W)	Avg. Pow (W)	Machine	CPU	Threads	HS/Watt	
451	367	Test AMD	1x EPYC 7643 HT	96	3.7	
512	449	Std AMD WN	2x EPYC 7513 HT	128	4.2	
1347	1198	AMD Bergamo	2x EPYC 9754 HT	512	6.3	50% improvement from 7513
348	270	Test ARM	1x ALTRA 80	80	5.6	
381	294	Farm ARM	2x ALTRA 80	160	5.2 ??	
664	502	Ampere Max	1x ALTRA Max	128	7.0	25% improvement from 80c
		0	HEPScore/Watt (∝ to Events/Joule)			
			8			

31

(hyperlink)



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



Takeaways and Conclusions

- errors in provisioning, but even with these we seem to be running well.
- since we expect most sites to become heterogeneous in the coming years.
- still present in a "real" environment, but smaller in magnitude.
- them (want > 1000 cores AltraMax M128-30).
- We've leaned a lot running both ATLAS and CMS work. Hopefully we'll be running ALICE work soon!



• Our experience running these machines is still not as mature as x86. Still hitting minor silly bugs on occasion and

Moving forward experiments may to be smart about how they advertise these resources and submit jobs to them

Preliminary comparisons show that the HEPScore/Watt performance increases from ARM compared to x86 are

• As a result of our tests; we decided to purchase more ARM machines and we have an ongoing procurement for







Additional Tips/Comments Technical setup differences between ARM and x86

and added in restrictions based on architecture (in the same way our Ansible

 We've encountered several bugs in the ARM version of packages (e.g. Grub2) which not clear whether it's actually an ARM issue or whether it's really an EL9 issue.

 As we'd never really considered the situation where we'd have architecture other than x86 in the site, PXE and Ansible required a one-off process where we went through configuration has already branched based on OS). This wasn't a lot of effort, though.

aren't in the x86 versions, but nothing we couldn't work around. Again, it's sometimes



Best way of advertising resources in the future?

At some point we will want to get rid of the test CE

Option 1

Modify exiting Queue

Option 2

 Add a queue pointing to the same CE's

UKI-SCOTGRID-GLASGOW_CEPH

ceO1.gla.scotgrid.ac.uk (x86) ceO2.gla.scotgrid.ac.uk (x86)

ceO3.gla.scotgrid.ac.uk (x86)

ce04.gla.scotgrid.ac.uk (x86)

ceO5.gla.scotgrid.ac.uk (aarch64)

- Dangerous, will impact the workflow of many VO's. Will have all the ARM traffic on one CE - not scalable in the future?
- All the CE's are on x86, they just front different architectures

UKI-SCOTGRID-GLASGOW_CEPH

ceO1.gla.scotgrid.ac.uk (x86) ceO2.gla.scotgrid.ac.uk (x86) ceO3.gla.scotgrid.ac.uk (x86)

ceO4.gla.scotgrid.ac.uk (x86)

UKI-SCOTGRID-GLASGOW_ARM

ceO1.gla.scotgrid.ac.uk (aarch64) ceO2.gla.scotgrid.ac.uk (aarch64) ceO3.gla.scotgrid.ac.uk (aarch64) ceO4.gla.scotgrid.ac.uk (aarch64)

• If Job requirements can be successfully injected this seems the safest option

Option 3

 Read architecture from the jobs themselves.

Option 4

 Pilots report to VO what architecture it's running on.



UKI-SCOTGRID-GLASGOW_CEPH

ceO1.gla.scotgrid.ac.uk (x86, aarch64) ceO2.gla.scotgrid.ac.uk (x86, aarch64) ceO3.gla.scotgrid.ac.uk (x86, aarch64) ceO4.gla.scotgrid.ac.uk (x86, aarch64)

- Condor_submit has architecture flags. Could try to pulling the architecture flag from condor_submit into the ARCsub, maybe modify in job description language (JDL)?
- Would potentially need to set up/inject default architecture so that "standard" x86 jobs don't get sent to ARM cores.
- That VO sends jobs of that type. Would require every VO to add this functionality to their pilots, not really default-able.
- Potentially wasteful if a site gets a pilot running on an ARM server and has no ARM work, long term solution?



35

Best way of advertising resources in the future?

At some point we will want to get rid of the test CE

Option 1

Modify exiting Queue

• Dangerous, will impact the workflow of many VO's. Will have all the ARM traffic on one CE - not scalable in the future?

This option preferred by ATLAS, certainly the simplest to implement

Option 2

 Add a queue pointing to the same CE's

UKI-SCOTGRID-GLASGOW_CEPH

ceO1.gla.scotgrid.ac.uk (x86) ceO2.gla.scotgrid.ac.uk (x86) ceO3.gla.scotgrid.ac.uk (x86) ceO4.gla.scotgrid.ac.uk (x86)

UKI-SCOTGRID-GLASGOW_ARM

ceO1.gla.scotgrid.ac.uk (aarch64) ceO2.gla.scotgrid.ac.uk (aarch64) ceO3.gla.scotgrid.ac.uk (aarch64) ceO4.gla.scotgrid.ac.uk (aarch64)

• If Job requirements can be successfully injected this seems the safest option



Option 3

Read architecture from the jobs themselves.

Option 4

Pilots report to VO what architecture it's running on.

- Condor_submit has architecture flags. Could try to pulling the architecture flag from condor_submit into the ARCsub, maybe modify in job description language (JDL)?
- Would potentially need to set up/inject default architecture so that "standard" x86 jobs don't get sent to ARM cores.
- That VO sends jobs of that type. Would require every VO to add this functionality to their pilots, not really default-able.
- Potentially wasteful if a site gets a pilot running on an ARM server and has no ARM work, long term solution?



