

Yves Kemp, C. Beyer, T. Hartmann, S. Dietrich, C. Voss NUC, 16.02.2023

NAF special 'incidences' since last NUC

No 'incidences' as such but new challenges

Energy crisis is a new reality

•expected compute cluster usage low between christmas and first week of new year - time for some powersaving by shutting down older nodes from 23-12 to 09-01

•NAF

•73 of 332 WNs (older systems, bought 2012 - 2014)

•Energy savings ca. 13 kw of 80 kw regular usage (16%)

•Number of slots lowered by 830 of 7.500 regular slots (11%)

•GRID

•287 of 391 WNs (older systems bought 2012-2016)

•Energy savings ca. 100 kw of 150 kw (66%)

Compute loss 137 kHS06 of 317 kHS06

•Pledges (ATLAS & CMS & BELLE = 164kHS06) fullfilled at all times (180kHS06)

Maxwell





Reduced Power Consumption

Over x-mas

- PoC: shutting off old nodes over the Christmas break
- Grid: shedding ~40% of HS06, saved ~60% of energy
- ~35 MWh saved over two week





Things we learned on the way & future plans

it's complicated & idle workers are evil

•Measuring power consumption can be tricky

•The base power consumption is much higher than we anticipated

•We developed a reliable monitoring of actual and possible max & min power consumption

•Older machines use more power in every sense

•Idle cores and even more idle workernodes are a big luxury

•Horizontal filling of the pool is to be reconsidered

•was debatable before as more the jobs are spread the more IPs are stressing the storage

•We need a more dynamic workernode management in order to shut-off

idle workernodes and 'wake' them when needed

•Preparation is done and a single test-workernode behaves accordingly to the plans

.Some integration in monitoring etc. needs to be tested



Future plans concerning power management

Previously: running HW till it dies

- As long as power came from the wall sockets, running old nodes made sense(?)
- Performance as [Benchmark per Watt] has grow ~4x over the last decade



Future plans concerning power management

Old HW tail

- cut off the tail of old nodes?
- Significant Compute Power...
 - ...that's quite inefficient compared to newer nodes



Raising awareness with users & working on code side

•We write to users once per week about their batch usage, and translate this to CO2 equivalent

•No negative feedback so far, some positive

•We think about a "storage CO2 footprint" ... More difficult

•Management also things about training for scientists ... Will see efforts in futre

NAF special 'incidences' since last NUC

Max_job_per_user_limit on sched14

•bird-htc-sched14 currently scheduler for the smaller VOs

•Total overloaded by astrop & belle users peak of 180k waiting jobs (end of january)

•5 to 10k jobs are considered 'healthy' by the software developers

.We have tweaked the scheds and run huge hardware

•180k waiting jobs in principal can be OK but one additional problem on top e.g. 10k jobs with a typo in the path lead to never

ending database updates and housekeeping

•This time someone was issueing 80 parallel condor_q requests in miliseconds ir

•Deleted all idle jobs while another submit of 50k jobs arrived

•Hence limited the max jobs per user to 5k similar to the CMS sched

•No complains so far and scheduler runs smooth



•See <u>https://confluence.desy.de/display/IS/large+number+of+jobs</u> for strategies **DESY.** | NUC slides | 16-02-2023 | Yves Kemp

New Hardware

Installed in january

- 12 x AMD EPYC
- 74F3 CPUs with 48 Cores & 256GB memory (4.3 GB/core)
- Comparison to old HW e.g. bird060 ~1 HS06/Watt new HW ~3.8 HS06/Watt !
- 253.692 HS06 added in total



The future of LINUX in NAF

CentOS7 vs. ALMA8 vs. ALMA9

•AlmaLinux 8 WGS hosts deployed

•Naf-alps-el8/naf-cms-el8/naf-ilc-el8/naf-belle-el8/naf-atlas-el8

•CentOS Linux 7 end-of-life roughly in 17 months

•AlmaLinux 8 support until 2029



•But full support for RHEL8 and thus for AlmaLinux 8 ends around the same time as EOL of CentOS Linux 7

•Means no new functionality and hardware support from that point on

•We propose to deploy AlmaLinux 9 during 2023 in the NAF, GRID & Maxwell (inline with CERN&WLCG)

•[.*]8 support through Apptainer AlmaLinux 9 nativ

 'RequestOS' feature for easy switching between OS versions without 'naming' an image path

DESY. | NUC slides | 16-02-2023 | Yves Kemp

DUST Status

Access Troubles & Faster Ethernet

- Access troubles on 2023-01-13 and 2023-01-26
 - Nodes waiting for each other, long waiters/deadlock situation
 - Solved by rebooting the affected NFS and storage server
 - Analyzed with vendors known issue w.r.t. to RPC message handling
 - Efix request rejected, only option: full upgrade to a new release
 - Online upgrade performed on 2023-01-31 without issues
- NFS servers: Ethernet upgrade 4x10GbE -> 2x100GbE
 - Node by node upgrade without downtime
 - Solves suboptimal traffic distribution across links
 - See plot: 3x10GbE saturated, 1x10GbE underutilized
- Reminder: DUST supports NFS4 ACLs
 - New ACL setting available, Linux *umask* now honored!



NAF Storage

GPFS/DCACHE

- Largely migrated old hardware for ATLAS/CMS/Belle; smaller VOs on DESY dCache still ongoing
- Observe increased hot-spotting with ATLAS and CMS causing problems with throughput and stability on NAF compute nodes (remedies under investigation)
- Plan to add additional NFS-doors for the NAF (currently 1 Door for all WNs) to decrease effect of single hotspots affecting all WNs
- Issue with one CMS storage node, offline for more than one week due to delays with Dell support (hopefully finally fixed by Friday 17th Feb 2023)
- Preparations for Disaster-Recovery largely in place; final preparations in the making; first packing/flushing to tape to take place soon