# Site Report HPC2N

NDGF All Hands Bergen 2025-05-13

## Network

- UMU finally at a 100G based core network
- HPC2N finally got its 100G core switch
  - With 100G uplink
- LHCOPN 100G
  - In production since late September 2023
  - Finally moved to a 2x100G LACP trunk connected to one router to have less panic if there's fibre/optics issues
  - SUNET wants us to dabble with BGP instead
    - BGP stub router (ie only default/single route) isn't too complex
    - Looking at doing 2x2x100G via BGP/ECMP as a pilot project

# Tape/backup

- Same as last time
- IBM TS4500 library (2550 slot capacity)
- 6x TS1155 tape drives (JD tapes, 15T, 360 MB/s)
- 6x TS1170 tape drives (JF tapes, 50T, 400 MB/s)
- Dell R750
  - 2x100G Ethernet
  - 4x32G FC
  - Approx 30T NVMe for DB and incoming stgpool
  - A few TB of SAS SSD for log mirrors etc
  - Approx 250T spinning disk for on-disk backup storage
  - 256 G RAM, 2xIntel Gold 6334 (total 16 cores @ 3.6 GHz)

#### New tape pools

- DELL PowerEdge R6615
- AMD EPYC 9124 3.0GHz, 16-core
- 64 GB RAM
- BOSS for OS (400G mirror)
- 5x3.5 TB NVMe in RAID5 via PERC H965i HW RAID controller
- 25 GbE network
- Four in total (ATLAS+ALICE READ+WRITE)
- Installed with Ubuntu 24.04 Noble

## $OS - Ubuntu : focal \rightarrow noble$

The data staging machines was upgraded (Jan-Feb)

 Smooth sailing except for ... ZFS mailing about missing devices @ reboot
 But no actual problems seen later

O Fix: zpool export arcpool; zpool import -a -d /dev/disk/by-path/

dCache disk pools upgraded to Noble

 Weird issues after a while with 100x slower read performance from disk, recovered after yet another while?

- Only one minor service host (grid-home & user mapping) still focal • Needs some fixing with python code so it works with 3.12
- Most of WLCG compute (CE, nodes) are Jammy

 Might do an upgrade of this to noble together with some major downtimes (like next week or upgrade of slurm version)

## WLCG compute – Aigert (aka g-ce01)

- Same as last time
  - 31 PowerEdge R6625, 256 cores/node, 3GB/core, 3.2TB/node
    - 25.4 HS23/core (~6510/node)
    - Only using 152 cores at the moment\*
  - OS Ubuntu 20.04 Jammy
- Only change is upgrade to ARC 7



### Upgrade ARC 7

- Mostly straight forward but ...
- Complicated by local puppet code for ARC CE not particularly ... nice?
- ENV/SINGULARITY a local version of upstream RTE
  - Can select image based on role
  - Was not updated/corrected for ARC7 change of path to job.local
  - Got strange failures with jobs *not* running in containers
- prometheus-arc-exporter
  - Puppet module created
  - Uses own modified arc-exporter code
    - Most changes could possibly be pushed upstream
    - Some changes should possible not be upstreamed

#### 256 cores is a lot of cores

- Up  $_{\tt nfiles}$  for cvmfs to match number of cores: 4096 × 256
- Increased cvmfs client cache size from 50 GB to 90GB (74GB quota)
  - Gianfranco Sciacca recommends 80GB quota, I modified it to 74GB so the increase would be 40GB
  - If the cache hitrate > 99%, 74GB should be sufficient. Otherwise, some investigation on the 99% number might be needed, and maybe increase it to 80GB.
    - ~ ½ nodes got just an increase to 61G, because they were running jobs
  - Added usage-graphs: prometheus-cvmfs-client
    - Uses node-exporter text plugin
    - 1st with cron-job, but that failed because of jobs not finishing in time which led to MASSIVE amounts of cron-job processes
    - Use systemd-unit instead
- Changed cleaning up of arc cache to keep up with increase of jobs
  - From cron-job to systemd timer
  - Want to avoid filling 10% of filesystem per cleaning interval
  - Within 1 minute, the maximum downloaded bytes are (based on a 100Gb/s interface)
    - 60 × 100 Gb = 6000 Gb = 750 GB = 0.75 TB
  - FS size is 12 TB, so 0.75TB is 6.25% < 10%
    - 2 minutes would be too long
  - max = floor(100 6.25)% = 93%, min = ceil(max 6.25/2)% = 90%

#### Problems - WLCG compute

- Missing lsc-files
  - Jobs started failing to map DNs for users
  - Forgot to update lsc-files
    - Manually entered yaml-entries to puppet
  - Added nagios-probe check missing lsc-files
    - Basically grep 'The lsc file .\* does not exist'
- Datastaging still stalling
  - Nordugrid bug 4191
  - Improved with ARC 7 but still lots of transfers busy doing nothing
  - Reduced core count since load seems to exacerbate the problem

<pre>dteam: host: voms2.hellasgrid.gr port: 15004 base: /C=GR/0=HellasGrid/OU=hellasgrid.gr ca: /C=GR/0=HellasGrid/OU=Certification Authorities/CN=HellasGrid CA 2016 ca: /C=GR/0=HellasGrid/OU=Certification Authorities/CN=HellasGrid CA 2016 host: voms-dteam-auth.cern.ch port: 443 file: voms-dteam-auth.cern.ch.vomses file: voms-dteam-auth.cern.ch extra:</pre>	
ot exist'	
e <mark>rs busy doing n</mark> othing	
acerbate the problem	