NDGF All Hands 2021-2

HPC2N Site Report



Replacing old WLCG/NDGF-T1 dCache disk pools

- Likely more Dell R740xd2
 - 8 servers each with:
 - 25 GigE networking
 - 26 of 18TB SAS 7kRPM HDDs
 - PERC H730P RAID-controller (had hoped to get H740, but not possible)
 - BOSS controller for OS storage (ie mirrored M.2 sticks)
 - 2 of Intel Xeon Silver 4215 2.5G, 8core
 - Needed to populate 2 sockets due to PCIe slot layout (motherboard is based on R540)
 - 96G RAM
- Also got offers from HPE and Lenovo, higher price per TB (but likely better RAID controllers).

New ARC-cache

- 4 HPE DL325 Gen10+ with AMD EPYC 7302P (16 core, 3.0 3.3 GHz)
- 128GB memory (8 × 16GB)
- $8 \times 1.6T$ NVME Mixed-Use (MU) SSD:s in each
- Separate HW-RAID1 for OS
 - $\circ ~~2\times 480~\text{GB}~\text{M.2}~\text{SSDs}$
- 100Gb EN Mellanox card running ethernet
 - Directly connected at 40Gb to the same switch as grid nodes
- 25Gb Marvell cards for LHCOPN connection

But HPE said they might have some supply issues with the EPYC cores which might lead to long lead time. So update order...

Actual new ARC-cache

- 4 HPE DL325 Gen10+ with AMD EPYC 7402P (24 core, 2.8 3.35 GHz)
- 128GB memory (8 × 16GB)
- $8 \times 1.6T$ NVME Mixed-Use (MU) SSD:s in each
- Separate RAID1 for OS
- 200Gb IB/HDR EN Mellanox card running ethernet
 - Still running at 40Gb
- 25Gb Marvell cards for LHCOPN connection

Delivery ...

- Was expected to be just after midsummer.
- Did not happen, but at least before Nikke went on vacation?
- The claim was that they were delivered on July 13th, but not to us.
- Hunted for the packages on Campus.
 - Was only able to find the local delivery firm that supposedly had delivered it.
- Finally came the day before Roger went on vacation.
- Was able to unpack and rack them.
- But the network was not working. And then it was time for vacation.

Namn Postnr/Ort Tidigast/Datum Senast/Datum Land Umeå Terminalen 90132 Umeå 2021-07-13 00:59 2021-07-13 08:00 Sweden

INFORMATION

Fraktsedelsnummer:	201532892	
Littera:	UT	
Datum:	2021-07-12	
Kundreferens	EDI	
Tjänst:	Frakter	
Kollin:	4	
Status:	Lossad	
Leveransanvisningar:	Consignmentnr: 800846098748	
	Startavsändare: MODRICE CZ	

Namn UMEI UNIVERSITET, MIT MARK FOR: Postnr/Ort 90736 UMEA Tel 070-5782767 Tidigast/Datum 2021-07-14 16:00 Land Sweden

T&T ORDERLOGG

Tid	Status	Kvittens	POD Nam n
2021-07-1 3 09:59	Levererad	2	godsmott agning
2021-07-1 2 08:26	Avgångsregistrerad U meå Terminal		
2021-07-1 2 08:26	Lastad på distribution sbil		
2021-07-0 9 09:25	Order skapad		

Installing & Testing

- Others got network working
 - Network card was at the wrong speed 10Gb instead of auto/25Gb
- Install and testing got delayed because Roger got sick.
- Nikke benchmarked out that ZFS would be good enough to use
 - In fact it was fastest but probably because ZFS avoid VFS which has some issues fixed in later kernels
 - ZFS is nice since it has built in checksumming
- While doing some simple benchmarking one machine hung or went into goslow mode
 - dd if=/dev/zero of=/dev/nvmeNn1 bs=1024K count=1200000



Fixing?

- Initial thought that it was the NVME drive problem
- Drive replacement did not fix.
- Next logical thing would be drive backplane. HPE thought system board/CPU
- System board/CPU replaced. Got DIMM errors
- HPE tries to fix DIMM errors with 2 more system boards
- Next thing they will try is DIMM reseat, DIMM, CPU reseat, CPU, system board
 - Supposed to have happened yesterday

HPE support

- HPE likes to follow their support script.
 - Send us AHS (and later SOS report)
 - Update to newest FW
 - The OS is the problem!
 - This with 4 identical machines with only one misbehaving
- HPE must not understand where Umeå is
 - Unable to fulfill 8x5 service
 - Parts come 1-2 days later than they say it will arrive at the latest
- They like to mention component shortage

Current state of new ARC-cache

- 4 HPE DL325 Gen10+ with AMD EPYC 7402P (24 core, 2.8 3.35 GHz)
 - \circ Only three working! >
- 128GB memory (8 × 16GB)
- $8 \times 1.6T$ NVME Mixed-Use (MU) SSD:s in each
 - Reformatted to 4k (native) block size
 - ZFS RAID0
- Separate RAID1 for OS
- 40 Gb ethernet from each ARC cache host towards the grid nodes
 - Actually IB HDR/EN 200Gb Mellanox card running ethernet
 - Directly connected to same switch as grid nodes
- 25Gb direct connection to LHCOPN

More LHCOPN bandwidth

- After lots of harassing the network people we finally got 4x10G to LHCOPN
- Less of a bottleneck now, but not a lot of headroom during peaks
- 100G (or actually 4x25G) "any year now" :/



hourly average + peak interval in hour

Future plans

- 100G LHCOPN on HPC2N segment (never ending story?)
- Purchase and commission replacement dCache disk pools
- Decommission old dCache disk pools (EOS 2022-03)
- Tape library drive upgrade
 - When next-gen drives are available, end of 2022?
 - No more purchases of current-gen tapes if we can avoid it
 - NSC takes up the HPC2N pledge slack for Sweden
 - 2022 Sweden pledge is 9.2PB ATLAS tape
 - HPC2N has 4.35 PB estimated capacity

The End