

# Site report for UiO

NDGF All-hands meeting 22.10.2019

Maiken, Vincent, Darren

# NDGF All-hands meeting April Ljubljana

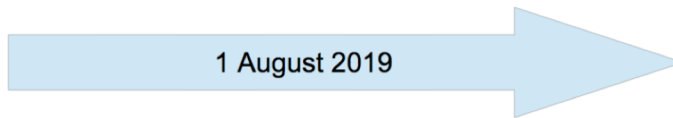
## UiO team



Darren Starr



2019-04



NDGF All-Hands 2019-1, Ljubljana



6

# Tier 1 @ Oslo

- Abel HPC is being decommissioned
  - ce01.grid.uio.no put out of service as of 15.09.2019
  - Low priority queue (ce03.grid.uio.no) is still going strong - will be running until Abel completely shuts off
- Soon...
  - Plan to take out a rack, upgrade to Centos7 & run grid-jobs there ...
  - Size: not sure yet

# New Tier1: UIO\_CLOUD - Grid on OpenStack @ UiO

- Openstack AMD Epyc machines - HEPSPEC06 12.5 w/hyperthreading
- ~3.7GB RAM per vCPU
- Instances with 8 vCPU's
- Size: 30 instances from ca June
- Limited by disk - therefore running ARC in pilot mode
- /scratch disk on each compute node - 200GB, 40 GB for cvmfms
- pilots download data directly to /scratch/<arcjobid> directory

October: got some more disk, and eventually used it as ARC cache

- Extended cluster with 10 instances
- Wed 16.10: Switched to Nordugrid mode with ARC datadelivery service and ARC cache
- 20T x 3 ARC caches
- 3 Datadelivery host machines

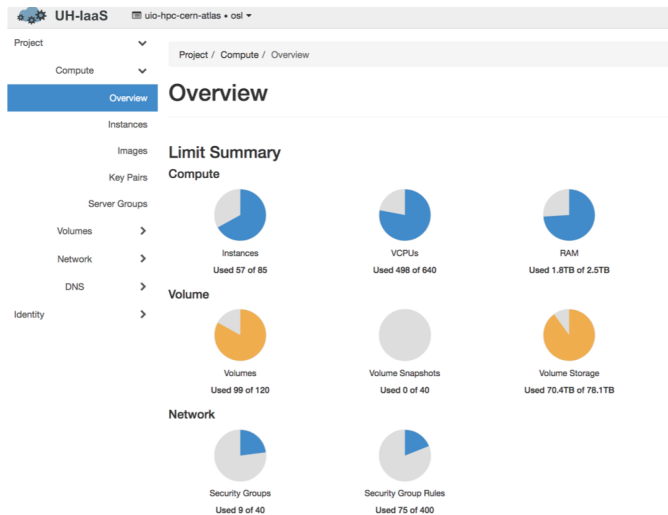
More disk is available now

- Extend cluster with 30 more instances (30\*8vCPU=240 vCPUs) ⇒ Total vCPUS in cluster then 560
- Possibly create more ARC datadelivery service hosts

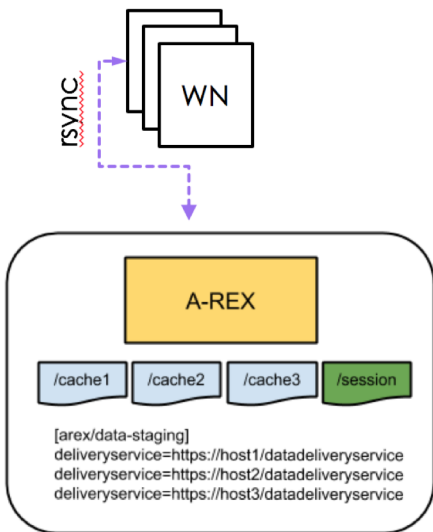
More servers have arrived and will be installed in OpenStack

- AMD Rome - 12 servers of 2(sockets)x2(threads)x48 cores - each with 512GB RAM ⇒ ~2.6 GB per thread after reserving some memory for the host
- 12x4x48=2304 vCPUs

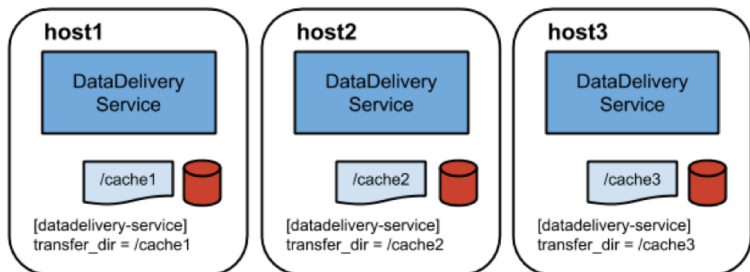
⇒ Covers (more than) pledge (22.4kSpec/12.5=1792 cores)



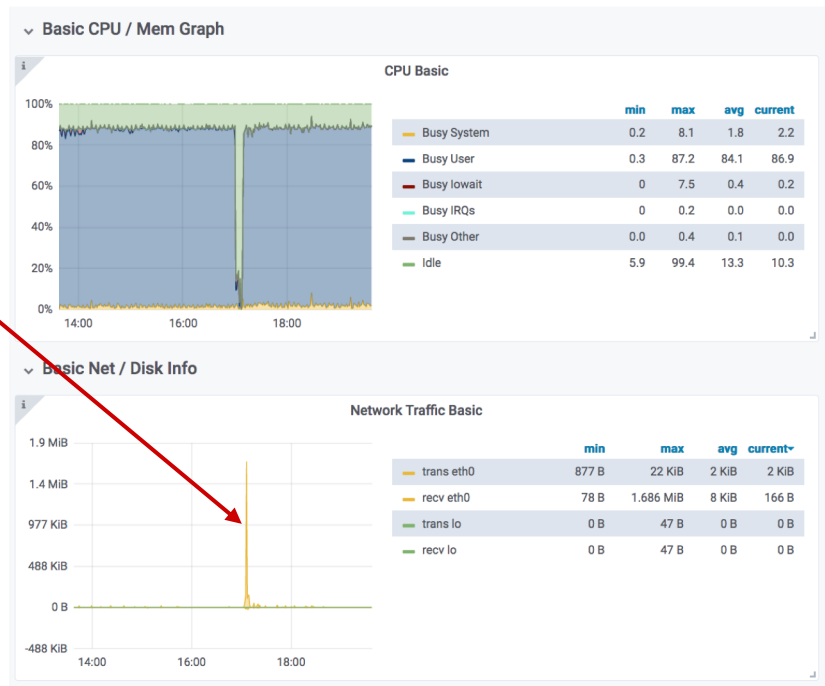
# Current setup of UIO\_CLOUD



rsync  
of 11 input  
files of  
6.3GB each  
(~70 GB tot)



- Wanted to avoid having to NFS mount cache dirs to worker nodes
- To get input data from ARC frontend sessiondir to workernode - rsync in submit\_common.sh

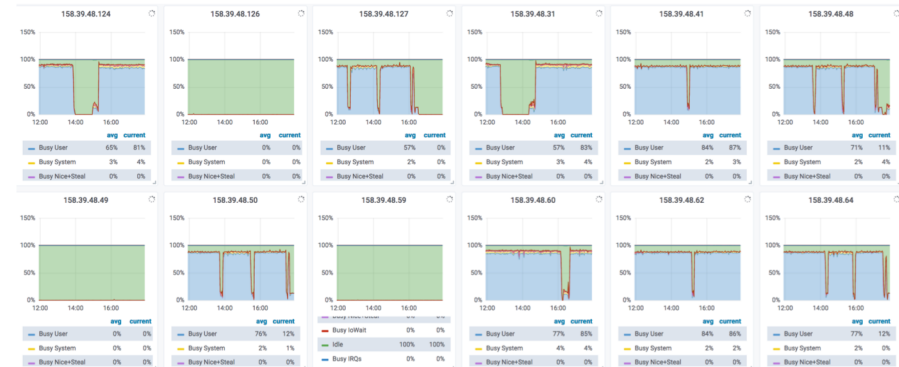


# ARC Datastaging

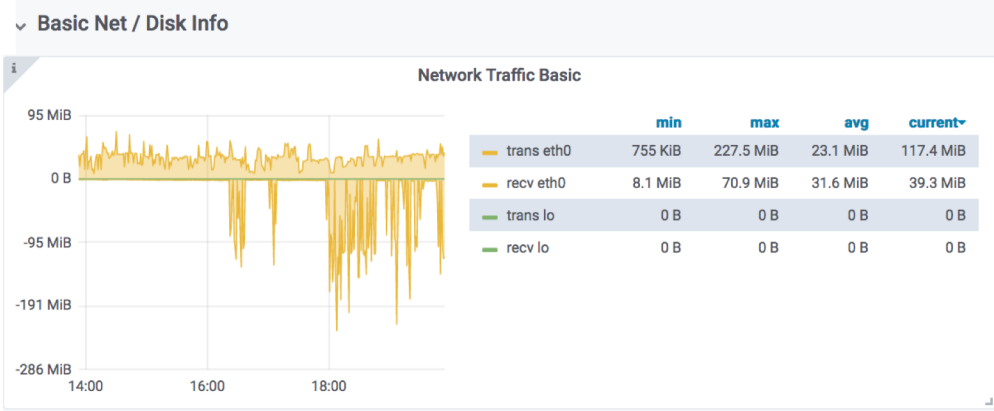
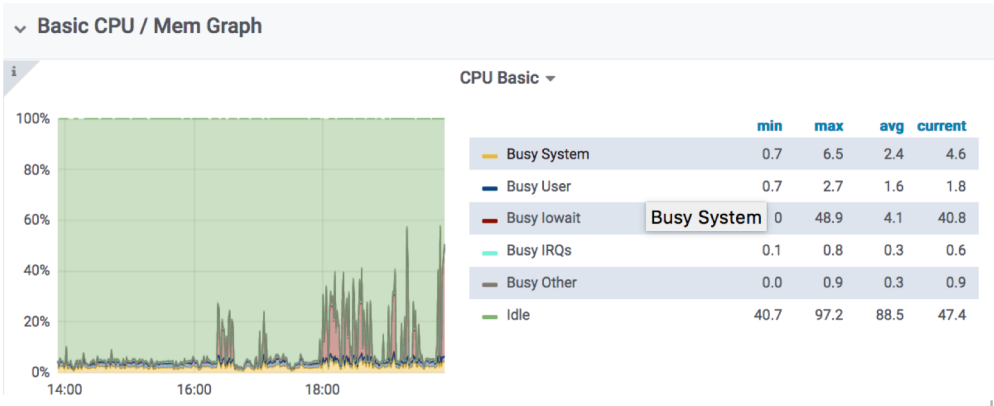
- Working on finding the optimal settings in order not to starve compute nodes
- io-wait is an issue
  - how will scaling the cluster up go?

```

PARTITION AVAIL  TIMELIMIT  NODES  STATE NODELIST
main*      up          infinite   1    mix compute006
main*      up          infinite  22    alloc compute[001-013,015-017,028],computeb[003,007-010]
main*      up          infinite  17    idle compute[014,018-027,029-030],computeb[001-002,004-005]
    
```



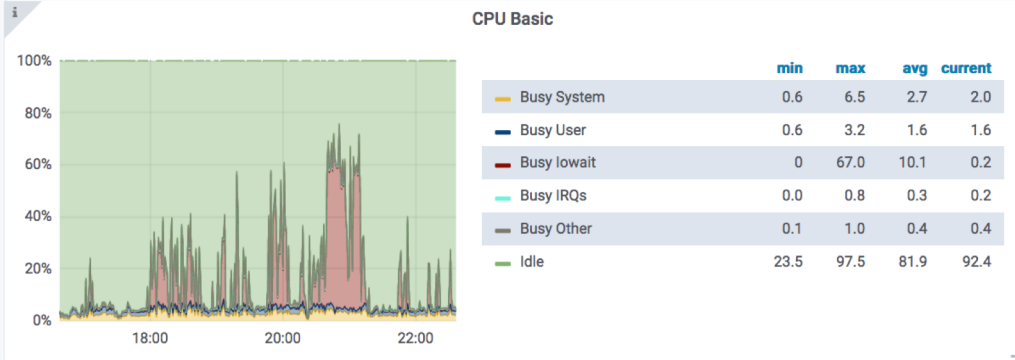
Host: arc-cache-1.grid.uiocloud.no ▾



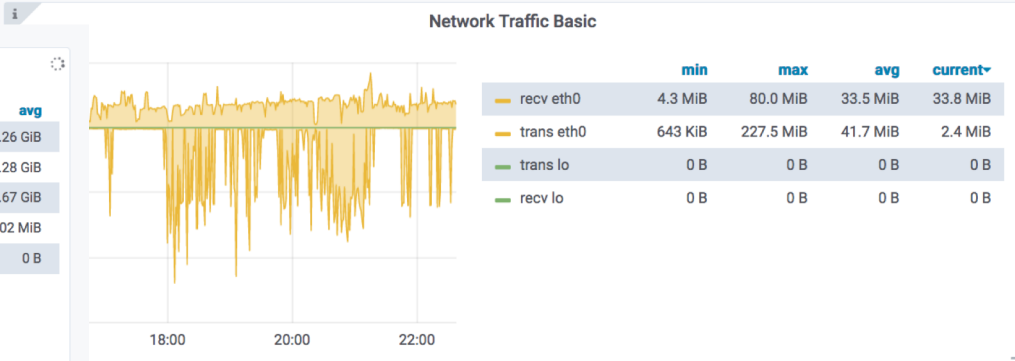
# Datadelivery host

- When there are a lot of reads in addition to writes, the IOWait rockets
- RAM Cache+Buffer filled up

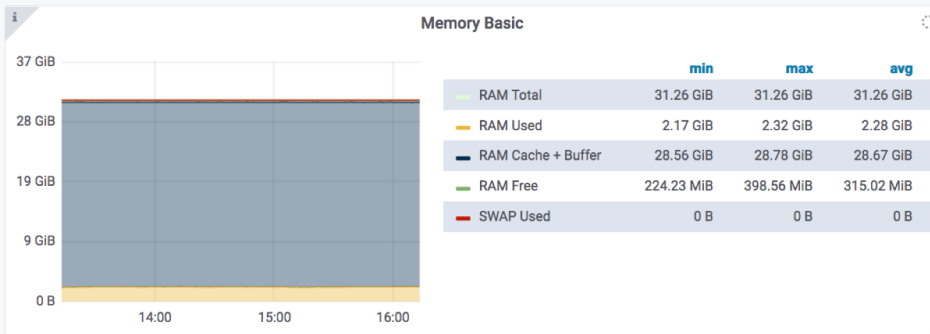
## Basic CPU / Mem Graph



## Basic Net / Disk Info

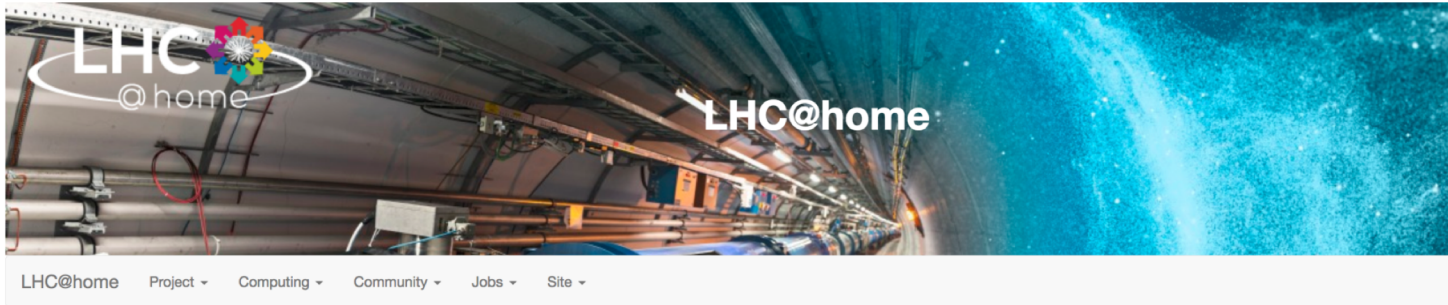


## Memory Basic



# UiO@Oracle

- Have been running for free in Oracle since before summer
  - UiO has an Oracle license which comes with CPU-hours which no one was using
- Running boinc on these (Intel/AMD - 100 instances of 8 vCPUS)



Rank	Name	Recent average credit	Total credit	Country	Participant since
1	AGLT2	4,019,061	1,653,943,055	United States	23 Jun 2014, 2:32:15 UTC
2	Agile Boincers	1,922,654	3,512,216,430	Switzerland	20 Sep 2012, 13:19:40 UTC
3	BlueHat	1,530,422	50,921,898	International	30 Aug 2019, 15:24:52 UTC
4	NDGF-T1	1,098,536	174,073,806	Norway	26 Feb 2019, 12:43:24 UTC
5	TRIUMF-LCG2	952,548	337,575,165	Canada	15 Mar 2018, 21:05:31 UTC
6	wHewitt	767,070	81,085,838	Canada	19 May 2014, 22:33:39 UTC
7	YellowHat	721,876	37,109,442	International	22 Jun 2019, 10:01:10 UTC
8	Toby Broom	620,175	333,788,357	Switzerland	27 Sep 2008, 16:22:03 UTC



# Storage @ Oslo

- Current dCache grid storage pools will be decommissioned in March 2020
- The new dCache grid storage pools are on **Ceph**
  - **2PB** of available space
  - Erasure Code
- New pools are being commissioned
  - Already in preproduction for local UiO physics group

# Tape @ Oslo

Slides from last all-hands with tape and disk status:

[https://indico.neic.no/event/73/contributions/282/attachments/51/86/UiO\\_site\\_report.pdf](https://indico.neic.no/event/73/contributions/282/attachments/51/86/UiO_site_report.pdf)

Extra....

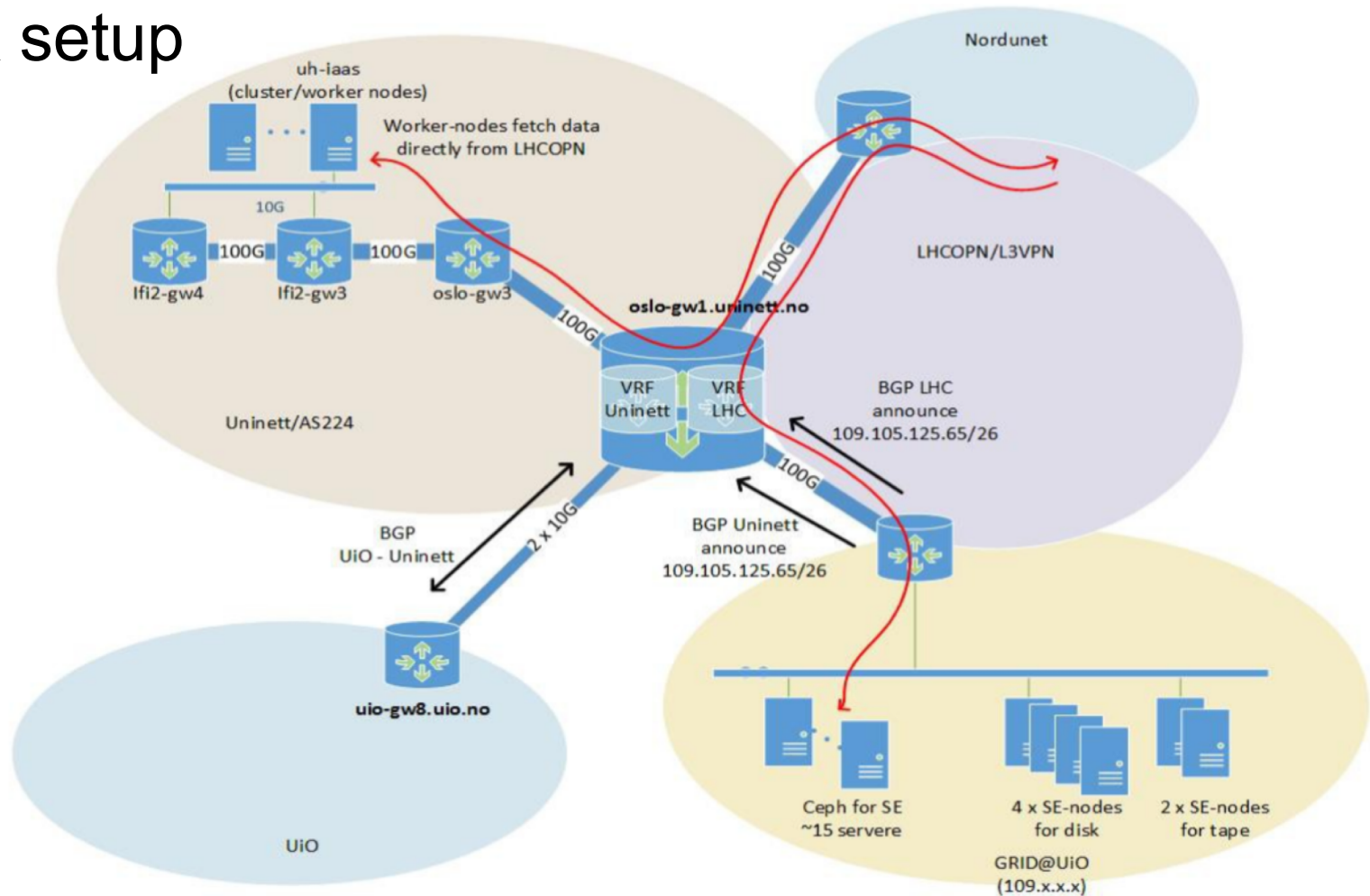
# Issues faced w/ OpenStack

- Instances got killed due to lack of memory
  - Solved from OpenStack team by reserving more memory for the host after cpu/thread pair and pinning on Numa node
- Jobs were using many times the max walltime set in the job description
  - Found out that if job asked for 8 cores, it started 10 athena.py processes where actually 9 were needing 100% cpu - caused extremely slow system
  - Solved by announcing only 7 cores per node to ATLAS - and actually running with 8 threads (lua plugin for SLURM)
- The ARC frontend suddenly got corrupted after switching to ARC cache
  - Got overloaded since I had wrongly configured the datastaging machines - frontend was doing all the data staging
  - Luckily managed to take a snapshot of the instance and fire up a new one which is working w/o problems.
- Things are more or less going well now :)

# Some details OpenStack configuration

- Instances get dedicated CPU and thread pairs which means that they then are pinned to one Numa-node
- Hugepages on
- 16 GB memory reserved to the host - 496 available to the instances
- Disk: block storage - volume service

# Network setup



# NorduNet load map

Map ndgf for 2019-10-21 showing peak traffic  
NDGF Tier-1 private network

