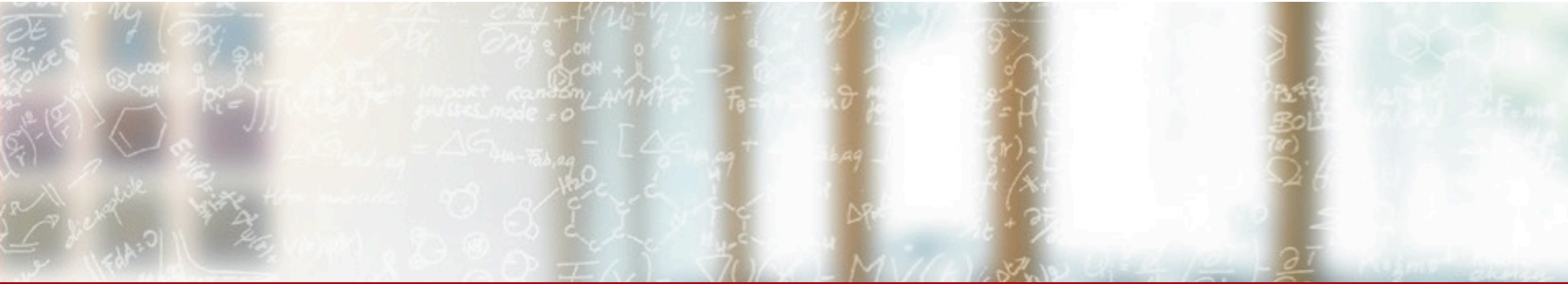




CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETHzürich



LHConCRAY

Acceptance Tests 2017 – Run5 System Report (Aug 03 2017 – Aug 31 2017)

Miguel Gila, CSCS

September 01, 2017

Table of Contents

1. Current configuration
2. System statistics
3. Proposition for Run6



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

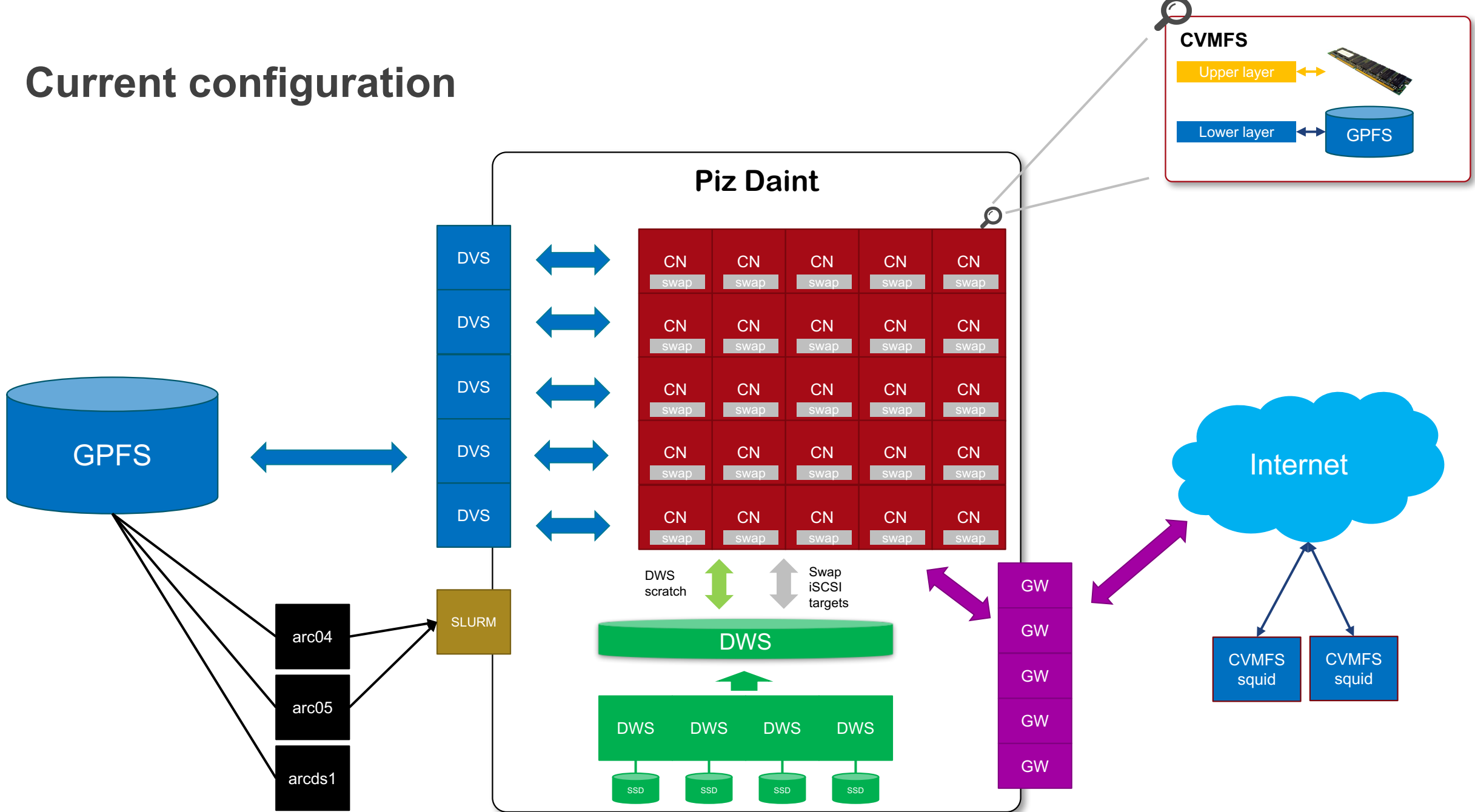
ETH zürich

Current configuration

Current configuration

- 25 compute nodes (72core, 128GB RAM, diskless)
- 1 production ARC + 1 data staging ARC + 1 test ARC (internal)
- Dedicated GPFS filesystem shared with Phoenix
- 5 DVS nodes exposing GPFS to CNs
- CVMFS running natively on CNs using tiered cache and workspaces
 - Upper layer: 6 GB in-RAM shared to all experiments
 - Lower layer: preloaded cache on GPFS, mount on CNs RO with caching enabled
- Memory limits NOT really enforced
 - Hard limit of 6000MB/core to catch rogue jobs
- Swap on DataWarp enabled (64GB per node)
- ARC caching not enabled (ATLAS)
 - Each job has a copy of the files, even if they're the same on multiple jobs.
- [Arc delegation database converted to sqlite for better performance \(22.08.2017\)](#)
- [New blackhole detector that will drain a node if more than 5 jobs have failed in the last 5min, but only if more than 10 jobs finished in the period \(28.08.2017\)](#)
- [Improved node auto-drain mechanism to be more permissive and drain less often \(28.08.2017\)](#)

Current configuration





CSCS

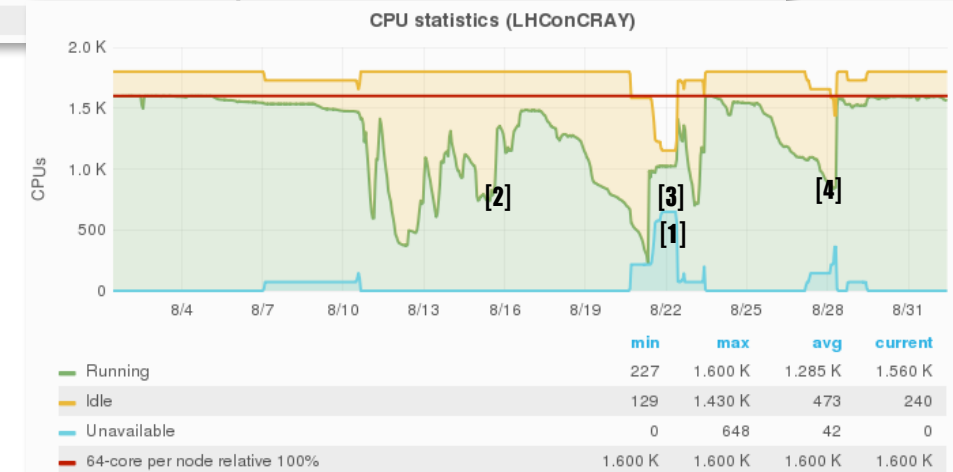
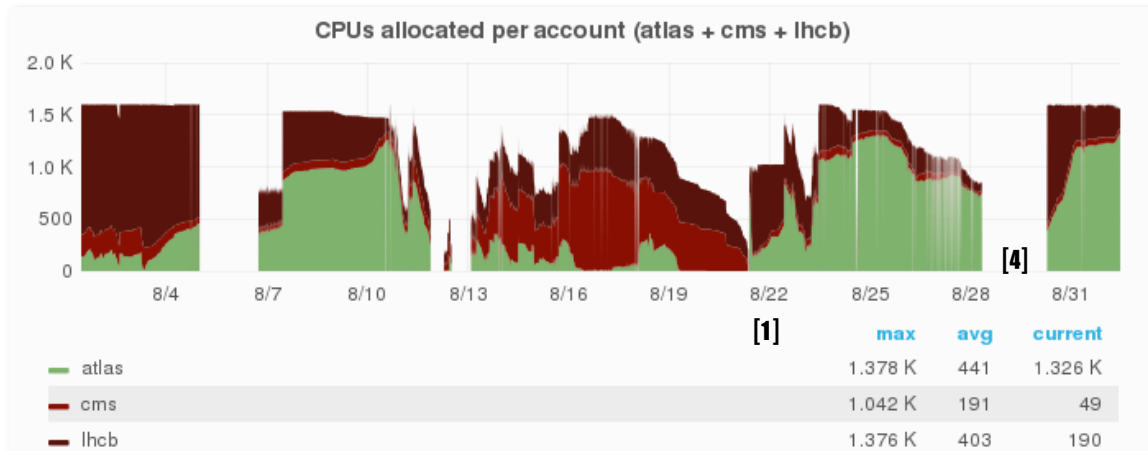
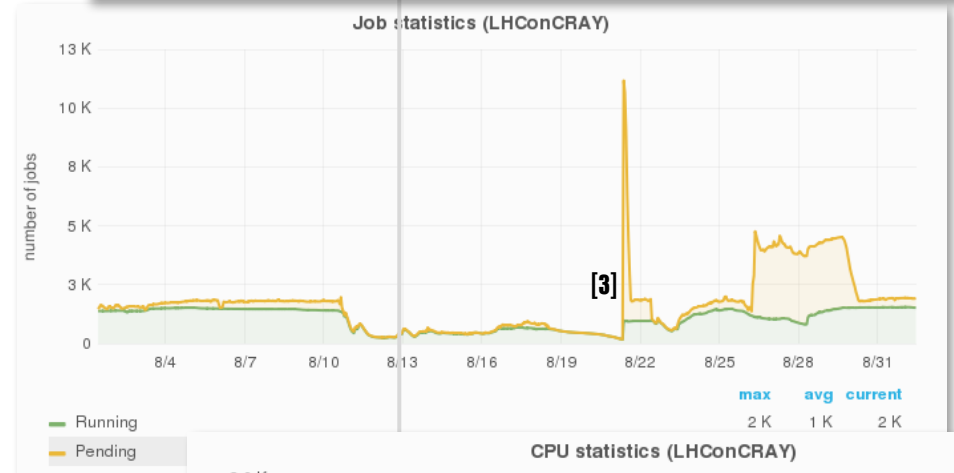
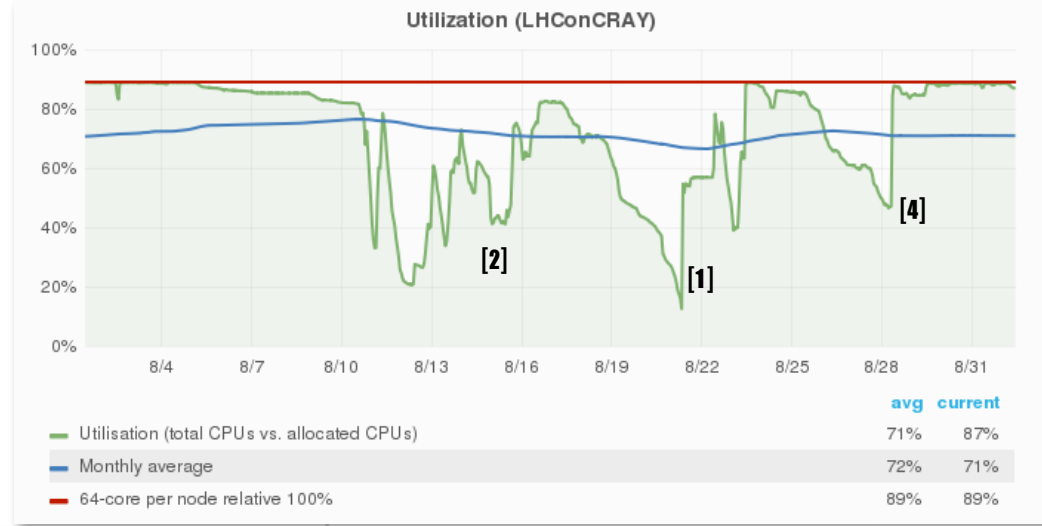
Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich

System statistics

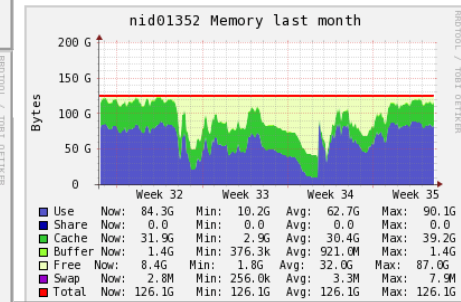
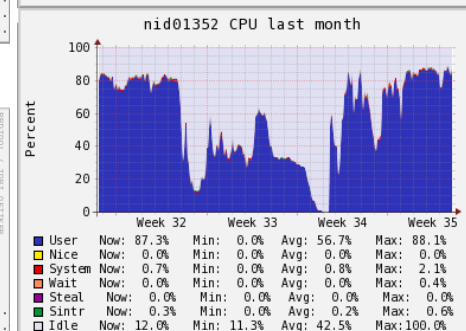
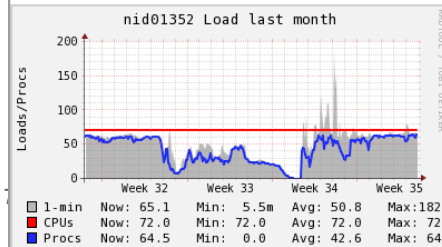
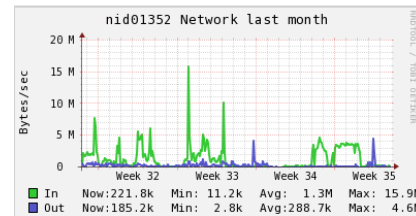
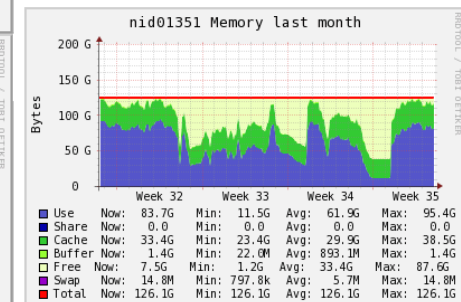
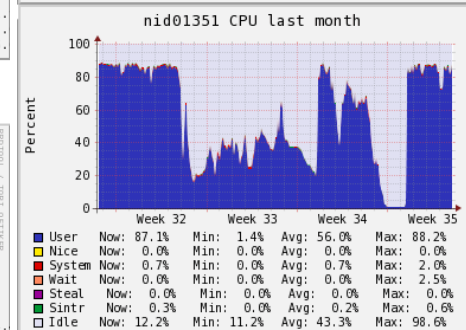
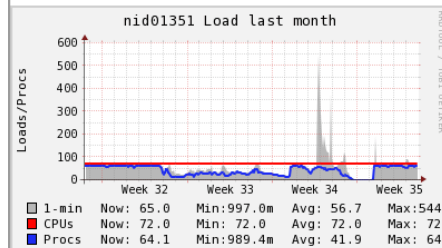
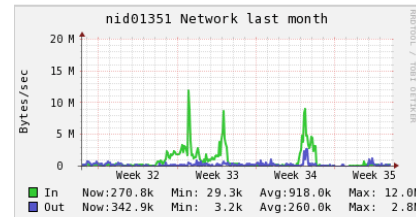
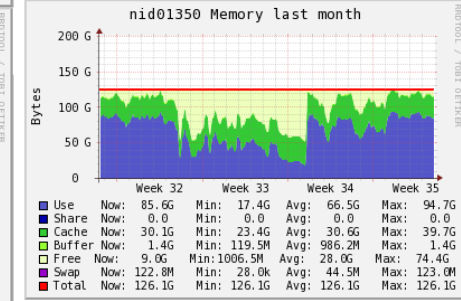
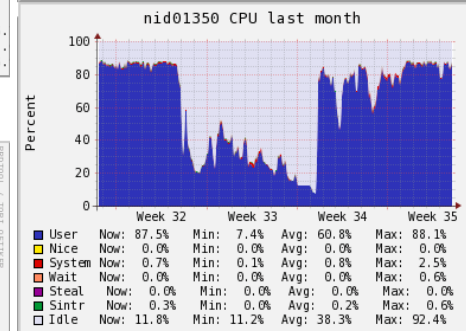
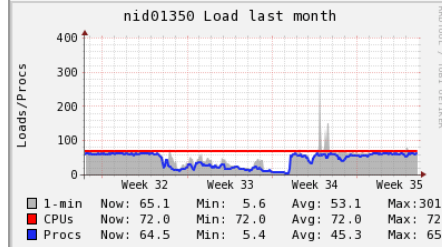
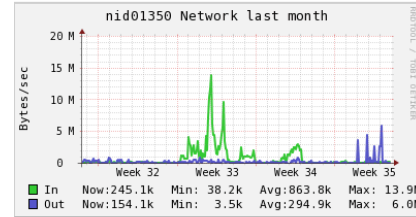
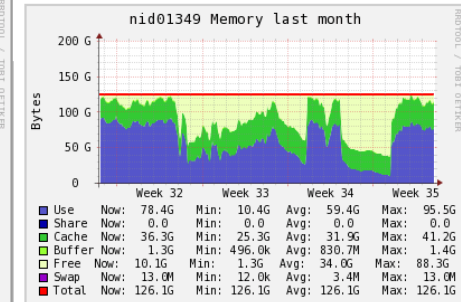
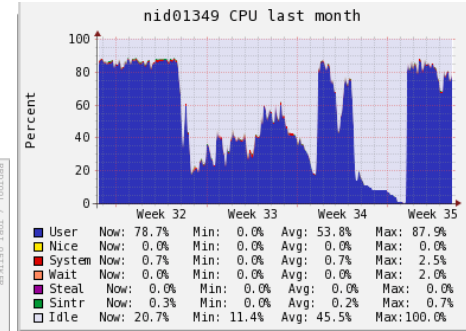
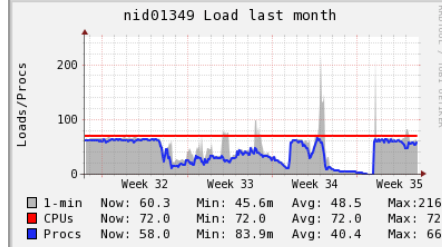
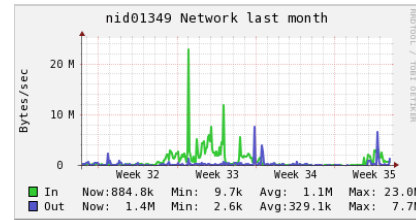
System utilization and issues

- Core allocation up to 100% relative with with 64core/node (out of 72) for long periods of time
- Encountered certain issues with ARC delegations [1] and nodes becoming silently blackholes [2]
- LHCb submitted ~10K jobs because of a problem with ARC BDII [3]
- Non LHC users hammered Slurm consistently and this affected scheduling for a while [4]
- ATLAS has picked up on LHCb and CMS seems to be consistently running a low number MC of jobs



Node statistics

- Load
 - Number of procs in line with load
 - Some load peaks due to IO
- CPU utilization
 - Almost flat on ~85%
 - IO wait negligible
- Memory utilization
 - About 30GB in cache
 - About 1GB free on average
- Network
 - No significant activity



Report

- Relatively stable operation, all VOs capable of running jobs
- Overall utilization reaching relative maximum
- Swap not really used so far
- There seems to be room to allocate more cores/node
- CVMFS in RAM seems to work quite well, not a single issue in the period
- DVS and node load high at times due to IO
- ATLAS and CMS have picked up CPU hours to LHCb
- About 9% of the total CPUhours available, CPUs were unavailable due to auto-drain or maintenance

Piz Daint	ATLAS	408'706	45%
Piz Daint	CMS	152'226	17%
Piz Daint	LHCb	355'457	39%
Piz Daint	TOTAL	916'389	

916'389 is **85%** of the total available time (1'075'200) !



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich

Proposal for Run6

Proposal for Run6

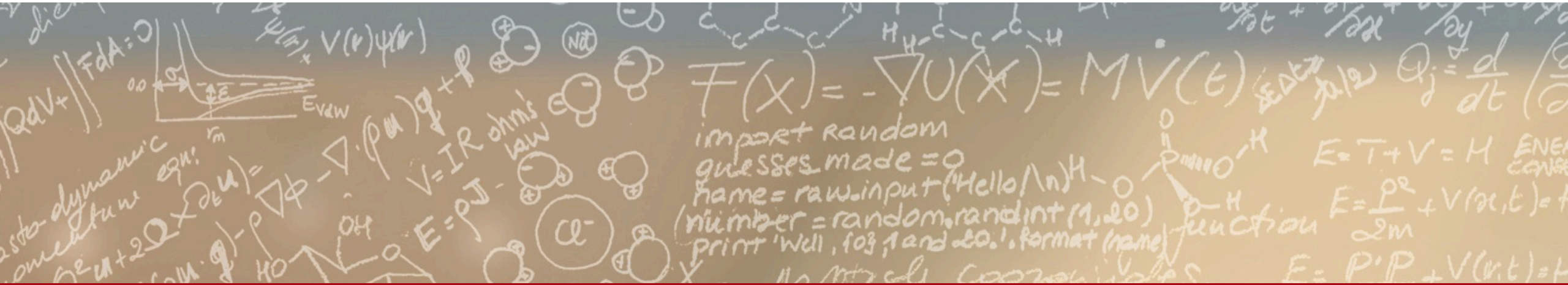
- Run6 between Sept 04 and Sept 27
- (from last meeting) Memory utilization suggests that there might be room to squeeze in a few cores and allocate ~66 or 68 cores instead of 64 cores
 - This could be positive (more CPU used!)
 - Or negative if nodes start swapping
- IO is affecting performance on GPFS over DVS; DataWarp nodes are patched. We propose to move jobs' runtime directory to SSDs on DataWarp
 - We could limit this per VO or per user
 - Only 15TiB available; current session directory in GPFS is about 2TiB
 - There seems to be enough space, but need to be careful to avoid exhausting it and killing jobs
- Extended maintenance around Sep. 27 2017
 - 3-day maintenance (from September 27 at 5:15 AM until September 29 at 4 PM)
 - Upgrade to CLE 6.0.UP04 and SLURM (minor)



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Thank you for your attention.