

A Large Ion Collider Experiment

---



**ALICE**

# ALICE HLT Cluster operation during ALICE Run 2

Johannes Lehrbach  
for the ALICE Collaboration



**FIAS** Frankfurt Institute  
for Advanced Studies



Bundesministerium  
für Bildung  
und Forschung



# Outline

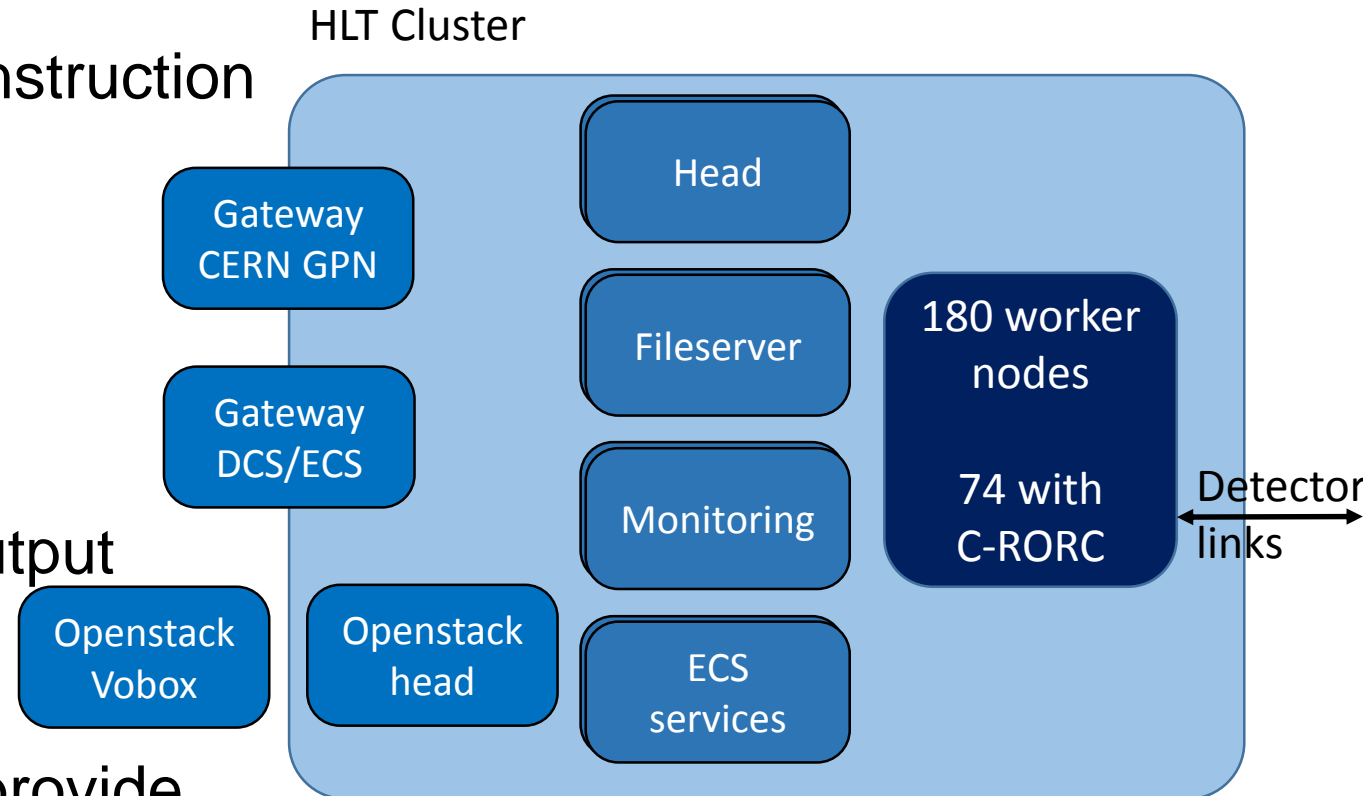
1. Introduction
  - High Level Trigger Cluster
2. Administration
  - Provisioning & configuration
  - Performance & stability
  - Monitoring & logging
3. WLCG operation
  - Grid setup



# High Level Trigger Cluster

Data compression and online reconstruction facility

- 180 compute nodes
- Gigabit Ethernet for provisioning, configuration, monitoring
- Infiniband for data transport
- Detector Data Links for Input / Output
- FPGAs and GPUs for hardware acceleration
- Additional infrastructure servers provide services like DNS, DHCP, NFS

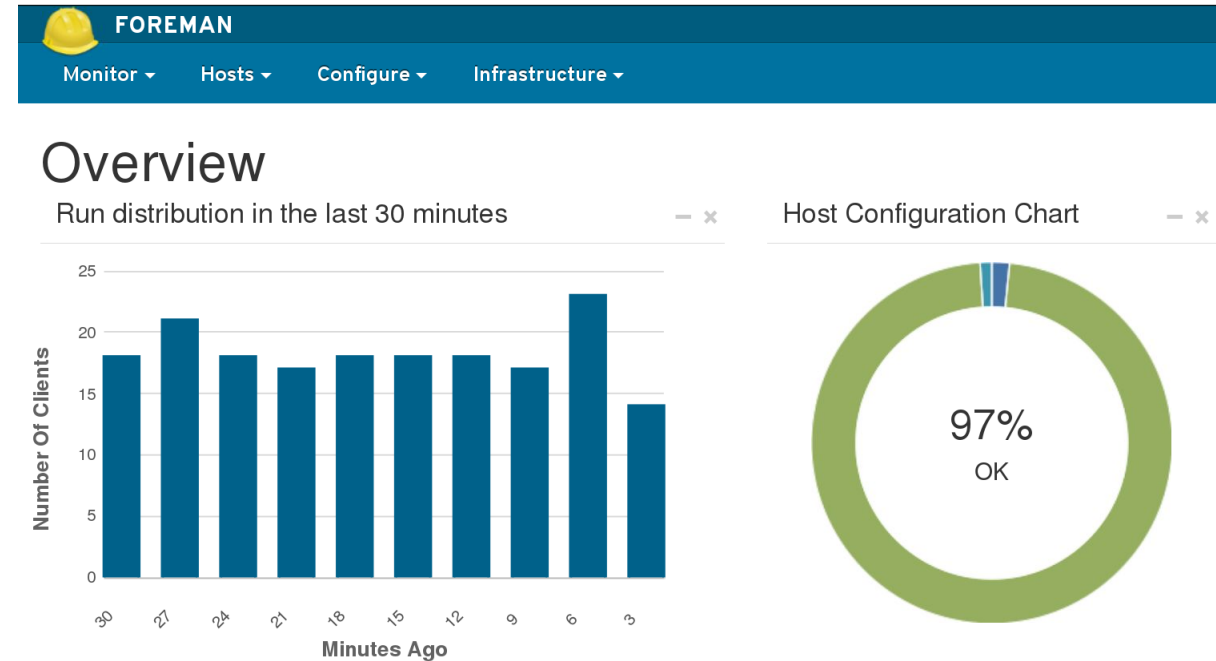


(GPU: D. Rohr, track 1; FPGA: H.Engel, track 1; Performance: M. Krzewicki, track 1)



# Cluster administration

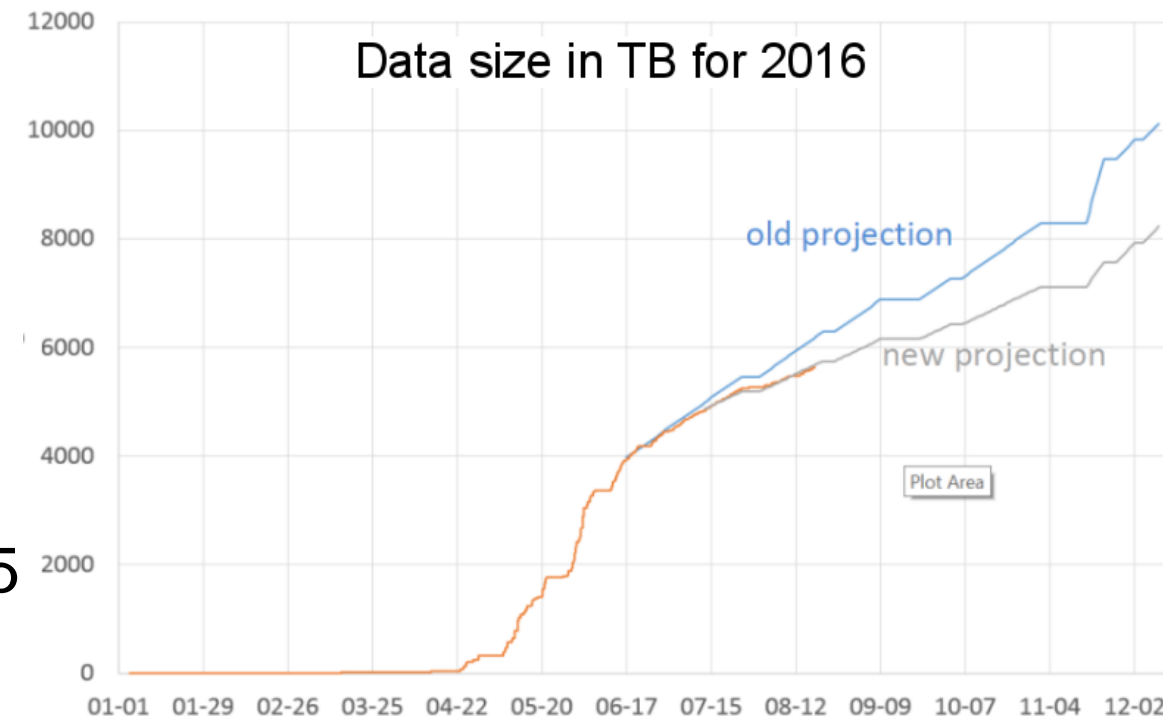
- Main goal automatize everything
- Provisioning via PXE over Ethernet with Foreman
- Multiple host groups / environments for different server roles
- Configuration of all servers done via Puppet
- Environments match Git branches
- No manual configuration
- Cluster rebuild easily possible





# Performance and stability

- After validation of the new TPC readout we were participating in every physics run during LHC collisions
- 665 runs with a total of 1070 hours
- 13 of 665 runs ended by HLT
  - 5 hardware failures
  - 5 software errors (fixed by now)
  - 3 operator requests
- Current TPC compression factor of  $\sim 5.5$ 
  - HLT improvements save  $\sim 20\%$  tape

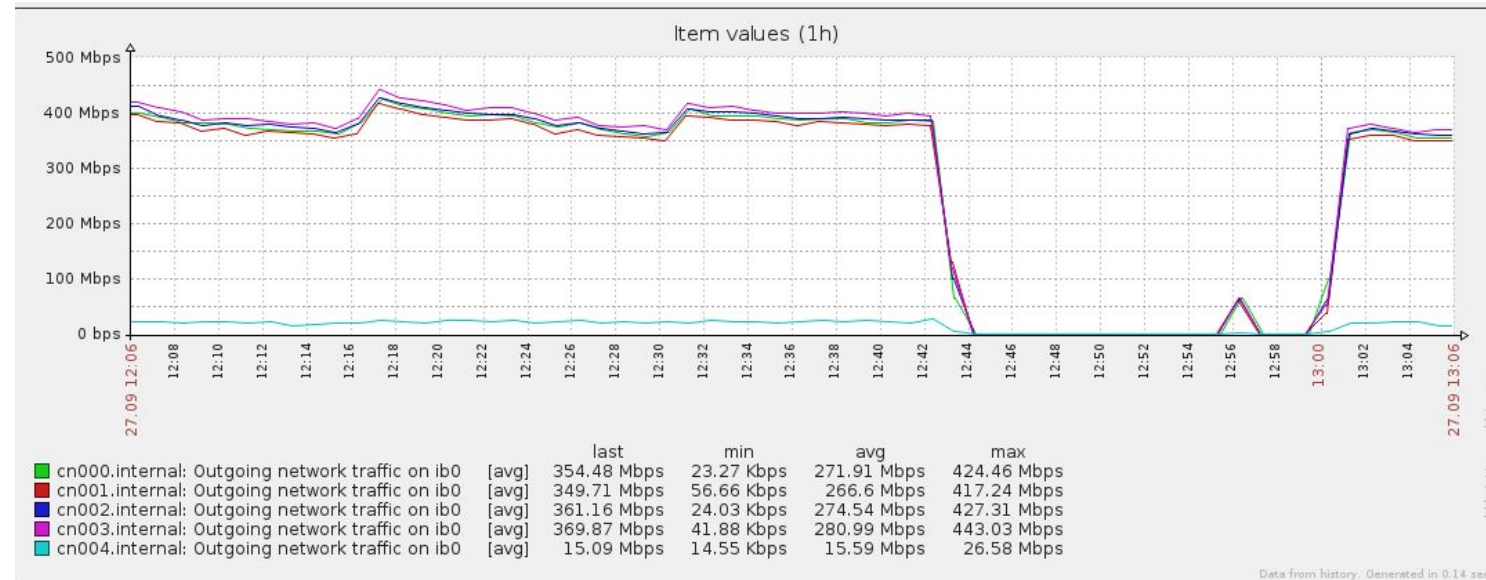


(see talk of M. Krzewicki, track 1)



# Monitoring

- Zabbix server / agent is used for monitoring
- 100 metrics per server e.g. load, temperature network traffic
- IPMI values included for some redundancy
- Daily cluster status mail with open issues
- Automated actions like node shutdown in case of temperature alert



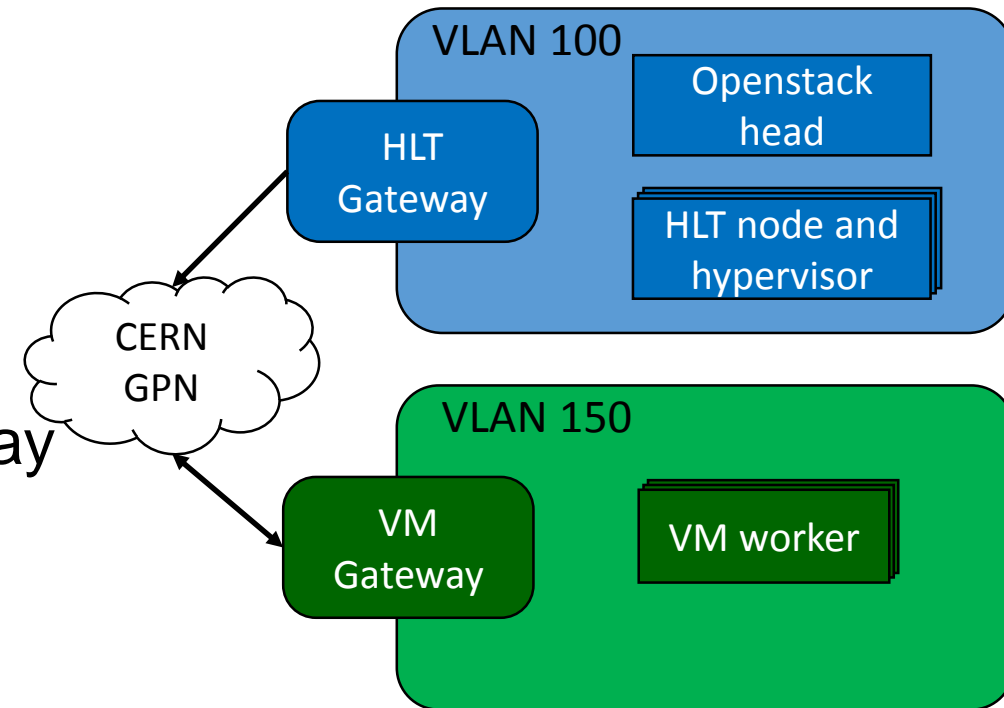


# Logging and additional notifications

- HLT framework logging via fork of ALICE DAQ Infologger
- Log messages in internal mysql database
- Analysis of log messages to push important information to the central Infologger → visible for shifter
- Email notifications in case of serious errors in our framework
- Usable by Detector code running inside HLT

# Worldwide LHC Computing Grid setup

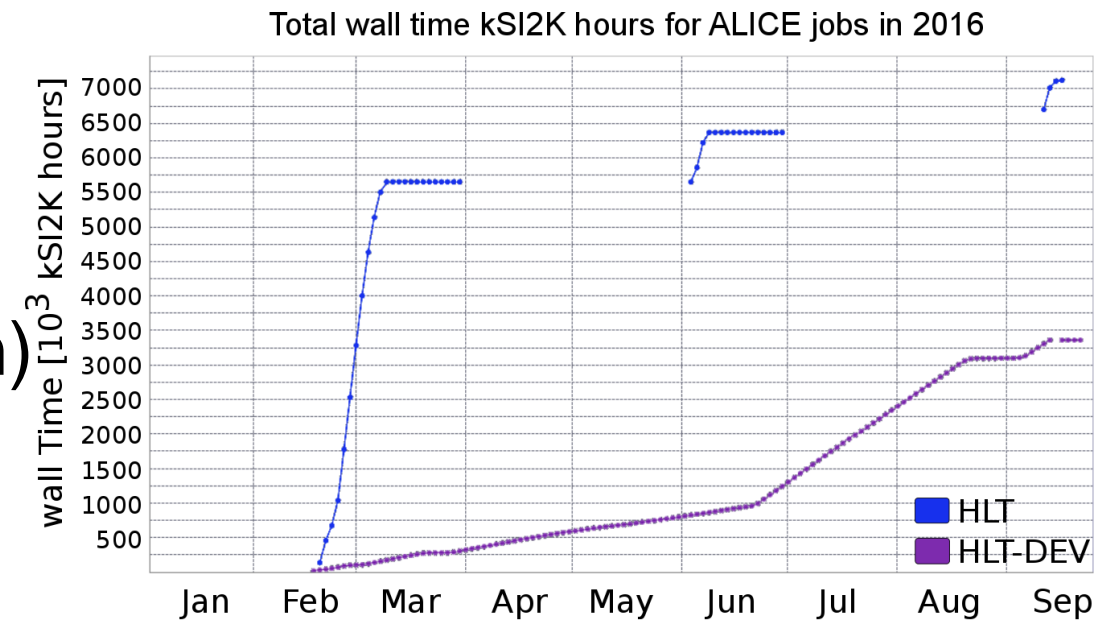
- In collaboration with our offline group we configured our cluster as a WLCG site
- Clear separation of the openstack VMs with separated VLANs on switch level
- IP table rules on hosts only allow communication between VMs and VM gateway
- No automatic start of openstack services at boot
- Only MC jobs are running due to limited I/O





# Worldwide LHC Computing Grid setup

- Running grid jobs on production cluster during technical stops
- Development cluster running grid jobs while no HLT development is done
- Management scripts allow flexible addition and removal of nodes
- Fast startup of VMs possible (~5 min)
- ~2.5% of ALICE MC jobs done by opportunistic usage of our clusters





# Summary

- HLT is capable to cope with the higher TPC readout rate after the upgrade
- Improving compression rate during Run 2
- We are running quite stable
- Development ongoing, already testing features for Run 3
- Opportunistic WLCG computation works well
- Next year longer WLCG operation planned during YETS



Thank you



# Backup

# HLT Performance / Compression

Main HLT tasks: Online reconstruction and data compression

- Improved compression factor from 4.3 to 5.5 this year
  - Reject QA data for clusters since histograms are available
  - Switched to differential Huffman compression
- Preparations for run 3 with further improvements
- Online reconstruction using FPGAs for Cluster-Finding and GPUs for Tracking
- Online monitoring of Detector output for QA
- Testing Online calibration which is needed for run 3



# HLT run 2 production Cluster Hardware

- HLT Production cluster Hardware for each compute node:
- Two Intel Xeon E5-2697 v2 CPUs with 12 cores @2,7GHz (+hyper threading)
- 128 GB DDR3 RAM
- Infiniband FDR 56GBits
- Gigabit Ethernet for monitoring and configuration
- AMD Firepro S9000 graphics card
- Two Intel SSDs in RAID 1 for operating system
- One additional SSD for openstack
- 74 of these nodes have an additional C-RORC for DDL input / output to DAQ
  
- 8 infrastructure nodes providing different services with slightly different CPU (E5-2690)



# Cluster administration

- Similar tool-set as CERN IT
- Complete configuration stored in Git
- Cluster rebuild easily possible in ~3 hours for all 180 compute nodes including final configuration
- Upgrade to CERN CentOs 7 this year
- Clear separation of production and development cluster