# SYSTEM ADMINISTRATION

## ALICE, ATLAS, CMS & LHCB
## JOINT WORKSHOP ON DAQ@LHC

- ❑ **Introduction**
- ❑ **Configuration**
- ❑ **Monitoring**
- ❑ **Virtualization**
- ❑ **Security and access**
- ❑ **Support**
- ❑ **Next steps and conclusions**

**Diana Scannicchio**
on behalf of
**ALICE, ATLAS, CMS, LHCb**
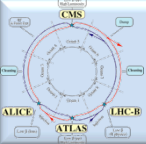**System Administration**

1

# Introduction: run efficiency

**The usage of big farms of computers is needed to take data (run)**

❑ ALICE:
~450 PCs

❑ ATLAS:
~3000 PCs, ~150 switches

❑ CMS:
~2900 PCs, ~150 switches

❑ LHCb:
~2000 PCs, ~200 switches

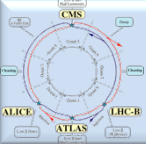**Achieve a good efficiency within the limits of available hardware, manpower, cost, ...**

❑ High availability, from the system administration (not DAQ) point of view:

★ minimize the number of single points of failure
➢ critical systems are unavoidable

★ have a fast recovery to minimize the downtime
➢ usage of configuration management tools and monitoring systems

❑ Complementing DAQ capability of adapting to the loss of nodes
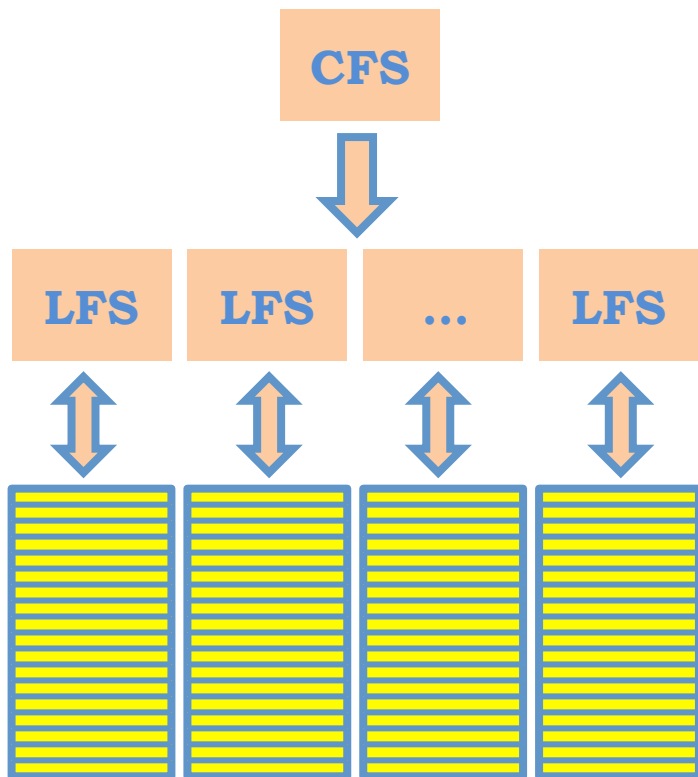
**The common goal is Run Efficiency**

# *Run*

❑ The farms are composed by nodes fulfilling various functions
  ★ Trigger and Data Acquisition
  ★ Detector Control Systems
  ★ Services
    ➢ monitoring, authorization, access, LDAP, NTP, MySQL, Apache, …
  ★ Control Rooms

❑ Run should survive GPN disconnection
  ★ any vital IT service is duplicated (DNS, NTP, DHCP, LDAP, DC)
  ★ event data can be locally stored for 1-2 days
    ➢ ATLAS and CMS

# Farm Architecture - ATLAS
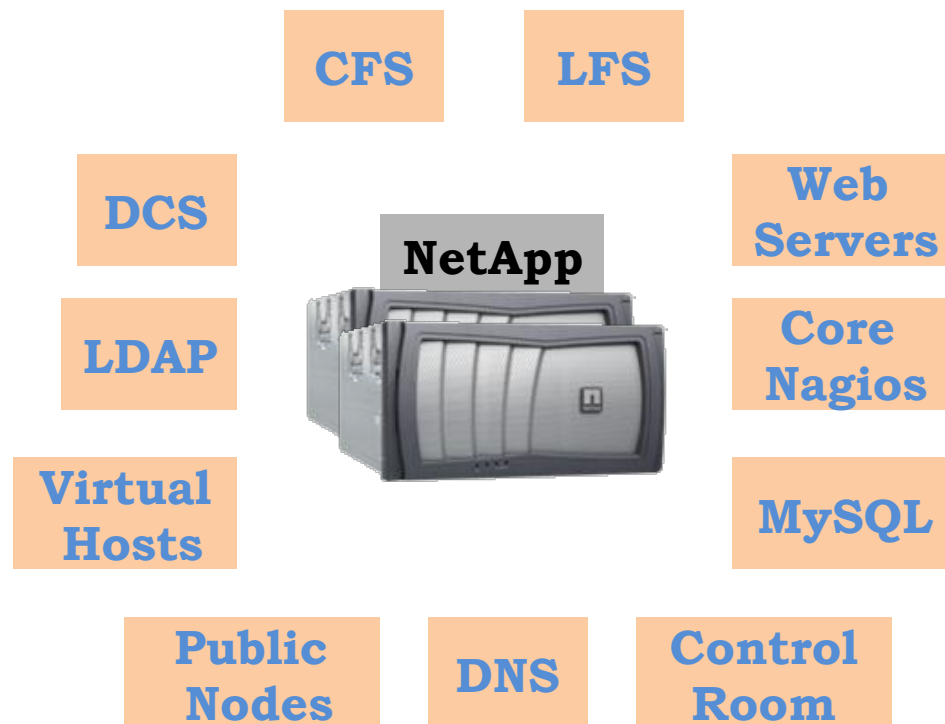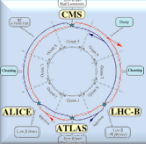
□ **Hierarchical structure**
- ★ Central File Server (CFS)
- ★ Local File Server (LFS)
- ★ netbooted nodes

□ **Flat structure**
- ★ local installed
- ★ NetApp: centralized storage
  - ➢ home directories and different project areas
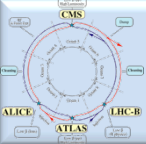  - ➢ 84 disks (6 spares), ~10 TB

# *Farm Architecture*

## CMS

➢ Flat structure

★ all nodes are local installed

★ NetApp: centralized storage

✓ home directories and different project areas

✓ ~17 TB

## ALICE
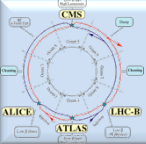
➢ Flat structure

★ all nodes are local installed

## LHCb

➢ Hierarchical structure

★ all nodes are netbooted

# *Efficiency*

❑ Single points of failure are impossible to avoid

★ ATLAS: DCS, ROS, NetApp (but it is redundant)

★ CMS: during LS1 DCS will move to blades for a large portion, with failover to a blade on surface

❑ Core services: DNS/DHCP/kerberos, LDAP, LFS are redundant

❑ Fast recovery

★ needed especially to recover a "single point of failure" system

★ monitoring is a fundamental tool

➢ to get promptly informed about failure or degradation

★ configuration management

➢ to quickly (re-)install a machine as it was, e.g. on new hardware (20~40 min.)

★ moving DNS alias (~15 min., due to propagation, caches)

★ diskless nodes have no re-install downtime (~5 min.) (ATLAS, LHCb)

➢ flexible system designed in-house to configure diskless nodes

➢ redundant boot servers to serve boot images, NFS shares, ...

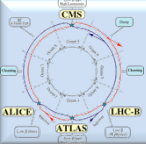❑ Efficiency loss due to hardware failures has been negligible compared to operator errors or detector failures

# Configuration Management

# *Configuration*

❑ Central configuration management is needed to speed up and keep under control the installation (OS and other software) on
  ★ local installed nodes
  ★ netbooted nodes

❑ Various configuration management tools are available, the ones used are:
  ★ Quattor
    ➢ CERN IT standard Configuration Management Tool
      ✓ being dismissed in favour of Puppet
    ➢ tight control on installed packages
    ➢ lack of flexibility for complex configuration and service dependencies
  ★ Puppet
    ➢ high flexibility
    ➢ active development community

# *Quattor and Puppet*

❑ Quattor
- ★ CMS
- ★ LHCb
- ★ ATLAS
  - ➢ still nodes configured by mixing with Puppet
  - ➢ finalizing the dismissing of Quattor in the next months
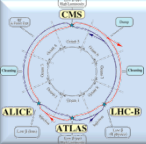
❑ Puppet
- ★ ALICE
  - ➢ the first configuration is done through kickstart, then puppet
- ★ ATLAS
  - ➢ in use for ~3 years, ~15000 LOC
  - ➢ complicated servers have been the first to be managed by Puppet
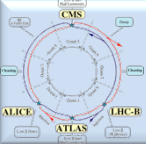  - ➢ on the HLT farm is complementing Quattor

# *Packages and updates*

❑ Software distribution and package management
  ★ SLC and other public RPMs from CERN repositories
    ➢ ALICE, ATLAS and CMS also have repositories mirrored in P2, P1 and P5
  ★ Trigger and DAQ software packaged as RPMs
    ➢ ALICE and CMS: installed locally on each node
    ➢ ATLAS: installed from CFS, synchronized to LFS, NFS-mounted on clients
    ➢ LHCb: in-house package distribution systems (Pacman, same as for GRID)
❑ Update policy
  ★ ATLAS
    ➢ snapshot of yum repositories, versioned test/production/… groups
    ➢ Quattor clients receive version list based on repository group
    ➢ Puppet clients pull directly from assigned repository group
  ★ CMS
    ➢ Quattor/SPMA controlled, updates are pushed as needed
  ★ ALICE
    ➢ updates are propagated at well-defined moments
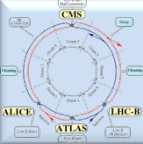  ★ LHCb
    ➢ updates are deployed at well-defined moments
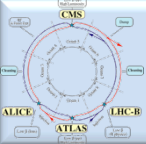
**See next Thursday for detailed news**

# Monitoring

# *Monitoring and alerting*

❑ **Large infrastructure must be monitored automatically**
  ★ proactively warned of any failure or degradation in the system
  ★ avoid or minimize downtime

❑ **What does monitoring mean?**
  ★ data collection
  ★ visualization of collected data (performance, health)
  ★ alert (sms, mail) on collected data

❑ **Various monitoring packages are available, the ones in use are:**
  ★ Icinga
  ★ Ganglia
  ★ Lemon
  ★ Nagios
  ★ Zabbix

# *Current monitoring tools*

- ❑ Lemon is used by Alice for metrics retrieval and display, and alerting
  - ★ monitoring Linux generic hosts and remote devices using SNMP
  - ★ retrieving DAQ-specific metrics (rates, software configuration, etc)
  - ★ reporting/alerting

- ❑ Nagios (v2) was used by CMS and is used by ATLAS
  - ★ problem with scaling in growing cluster
  - ★ configuration is distributed over more servers in order to scale

- ❑ Ganglia is used by ATLAS to provide detailed performance information on interesting servers (e.g. LFS, virtual hosts, …)
  - ★ no alert capabilities

- ❑ Icinga is already being used by CMS and LHCb
  - ★ configuration is compatible with the Nagios one, so it is "easy" to migrate
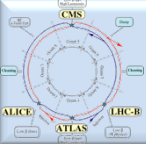  - ★ data collection is performed using Gearman/mod_gearman (queue system) to distribute the work load

# *Future monitoring tools*

❑ ALICE will replace Lemon with Zabbix

❑ ATLAS will complete the migration to Icinga complementing the information with GANGLIA

★ Gearman/mod_gearman to reduce workload on the monitoring server and improve scaling capabilities

❑ LHCb will also use GANGLIA

**See next Thursday for detailed news**

# Virtualization

# *Virtualization in the present*

**ALICE**

➢ none

**ATLAS**

➢ gateways
➢ domain controllers
➢ few windows services
➢ development web servers
➢ core Nagios servers
➢ Puppet and Quattor servers
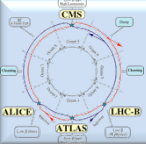➢ one detector machine
➢ public nodes

**CMS**

➢ domain controllers
➢ Icinga workers and replacement server
➢ few detector machines

**LHCb**

➢ web services
➢ infrastructure services
  ★ DNS, Domain Controller, DHCP, firewalls
  ★ always a tandem for critical systems: one VM, one real
➢ few control PCs

# *Virtualization in the future*

❑ Virtualization is a very fertile playground
- ★ Everyone thinking how to exploit

❑ Offline software (analysis and simulation) will run on virtual machines on the ATLAS and CMS HLT farms
- ★ OpenStack is used for management

**ALICE**
- ➢ Control Room PCs
- ➢ Event Builders

**LHCb**
- ➢ general login services
  - ★ gateways and windows remote desktop
- ➢ all control PCs
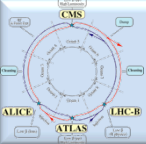  - ★ PVSS, linux, windows, specific HW issues (CANBUS)

**ATLAS**
- ➢ DCS windows systems

**CMS**
- ➢ servers
  - ★ DNS, DHCP, kerberos, LDAP slaves
- ➢ DAQ services

**See next Thursday for detailed news**

# Security and Access Management

# *Authentication*

## ALICE

- internal usernames/passwords used for detector people
  - ★ no sync with NICE users/ passwords
- RFID/Smartcard authentication after LS1
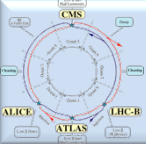  - ★ still no access to/from outside world

## ATLAS

- local LDAP for account information
  - ★ usernames and local password if needed (e.g. generic accounts)
- NICE authentication using the CERN Domain Controllers mirrors inside P1

## CMS

- local kerberos server
  - ★ same usernames and userID as in IT
- LDAP is used to store user info and user to group mappings

## LHCb

- Local LDAP
- Local Domain Controllers
- UIDs, usernames and user info are in sync with the CERN LDAP
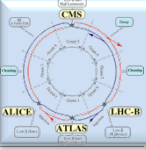
# Security and Access Restriction

❑ Web pages and Logbooks are
  ★ accessible from outside CERN and secured through CERN SSO
  ★ firewalls and reverse proxies also used

❑ The networks are separated from GPN and TN (for ATLAS, CMS, LHCb)
  ★ exceptions are implemented via CERN LanDB Control Sets

**ALICE**
➢ no external/GPN access to any DAQ services

**LHCb**
➢ no external/GPN access to any DAQ services
  ★ access is possible only with an LHCb account through the linux gateways or windows terminal servers
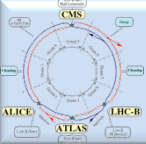
# Security and Access Restriction

## ATLAS

➢ **access to the ATLAS network is controlled**
  - ★ RBAC (Role Based Access Control) mechanism in place to restrict user access to nodes and resources (i.e. Access Manager)
  - ★ during Run Time the access is only authorized by ShiftLeader, and it is time limited
  - ★ sudo rules define limited administration privileges for users

➢ **two steps for a user to login on a P1 node**
  - ★ first step on the gateway where roles are checked before completing the connection
  - ★ second step to the internal host, managed by login script

## CMS

➢ **access to the CMS network via boundary nodes (user head nodes) is not blocked at any time, any valid account can login**
  - ★ nodes are not restricted either (anyone can log into any machine)
  - ★ sudo rules are restrictive to the types/uses of nodes
  - ★ access is through password authentication only for the peripheral nodes (SSH keys not allowed)

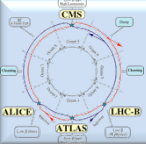➢ **The boundary nodes are fully fledged nodes similar to general nodes on the network**

# Support

# *Workload and requests management*

❑ Ticket systems are used to track issues and requests
  ★ ALICE and CMS use Savannah and will move to Jira
  ★ ATLAS uses Redmine for 3 years (before Jira availability)
  ★ LHCb uses OTRS and has installed Redmine

❑ Urgent matters are managed via on-call with different philosophies
  ★ ALICE: DAQ on-call and the other DAQ experts as needed
  ★ ATLAS: direct call to TDAQ SysAdmins
  ★ CMS and LHCb: DAQ on-call is the first line, then SysAdmins
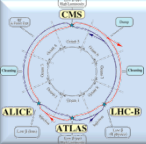
# Next Steps and Conclusions

# *Next steps*
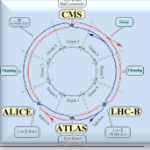
A lot of work is planned by all experiments during LS1

❑ Updating the Operating Systems to
  ★ SLC6 on both local installed and netbooted nodes
  ★ Windows Server 2008 or later

❑ Complete the migration to new configuration management tool

❑ Upgrading and improving the monitoring systems

❑ Looking more and more at virtualization
  ★ HLT Farms will be used as virtual machines to run offline software

# *Conclusions*

- ❑ Systems are working: we happily ran and took data
  - ★ complex systems
  - ★ 24x7 support
- ❑ Interesting and proactive "Cross Experiment" meetings to
  - ★ share information
  - ★ compare solutions and performances
- ❑ Converging on using the same or similar tools for "objective" tasks
  - ★ e.g. for monitoring and configuration management
- ❑ Appropriate tools are now available to deal with big farms
  - ★ big farms are now available outside in the world
  - ★ CERN is no more a peculiarity
- ❑ Differences observed for "subjective" tasks
  - ★ restrict access or not
  - ★ uniformity (netbooted) vs. flexibility (local installed)
- ❑ Improvement is always possible... unfortunately it depends on costs, time and manpower

# *Thanks to...*

❑ ALICE
  ★ Adriana Telesca
  ★ Ulrich Fuchs

❑ CMS
  ★ Marc Dobson

❑ LHCb
  ★ Enrico Bonaccorsi
  ★ Christophe Haen
  ★ Niko Neufeld