

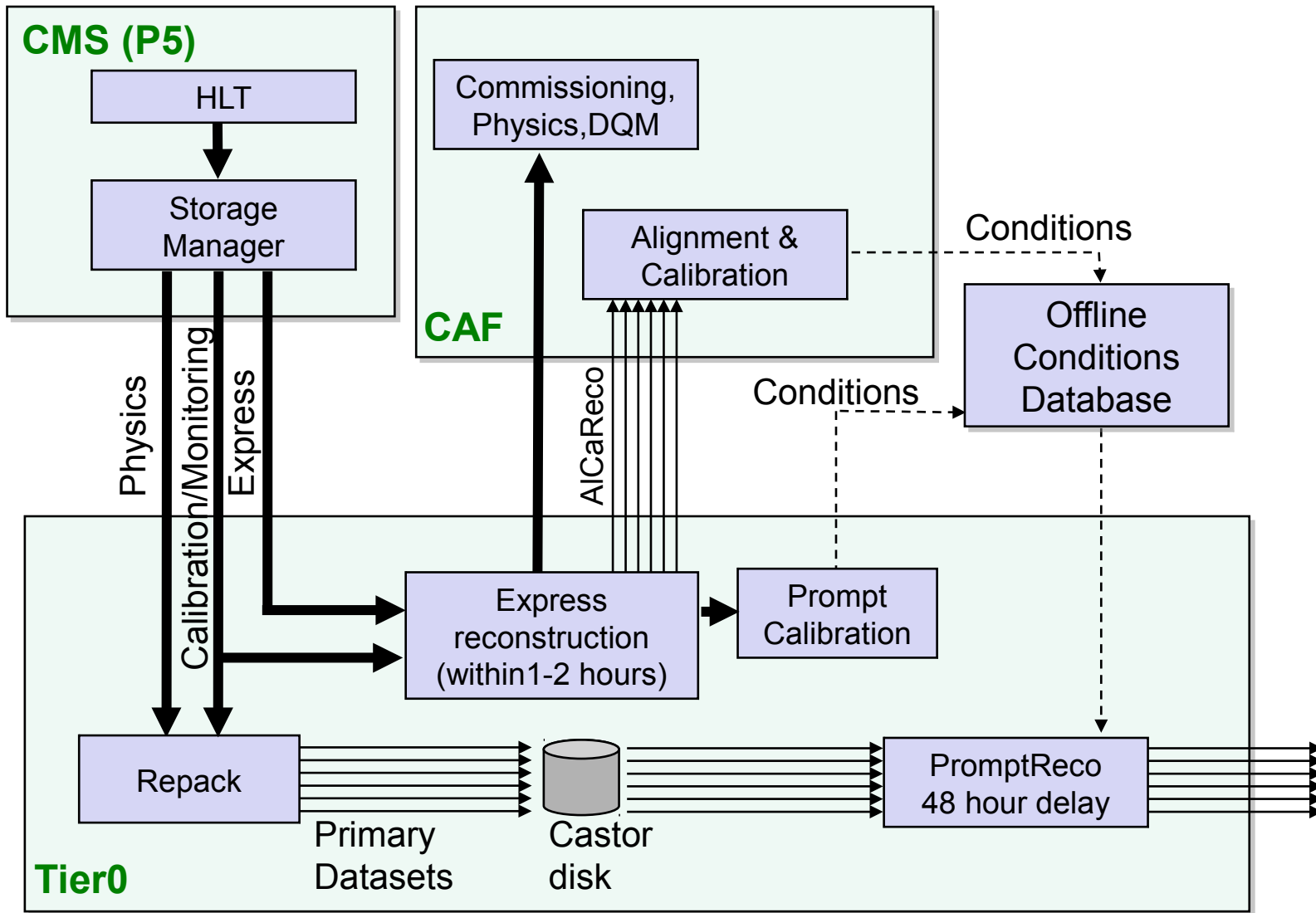
The CMS Tier0 goes Cloud and Grid for LHC Run 2

Dirk Hufnagel (FNAL)
for CMS Computing

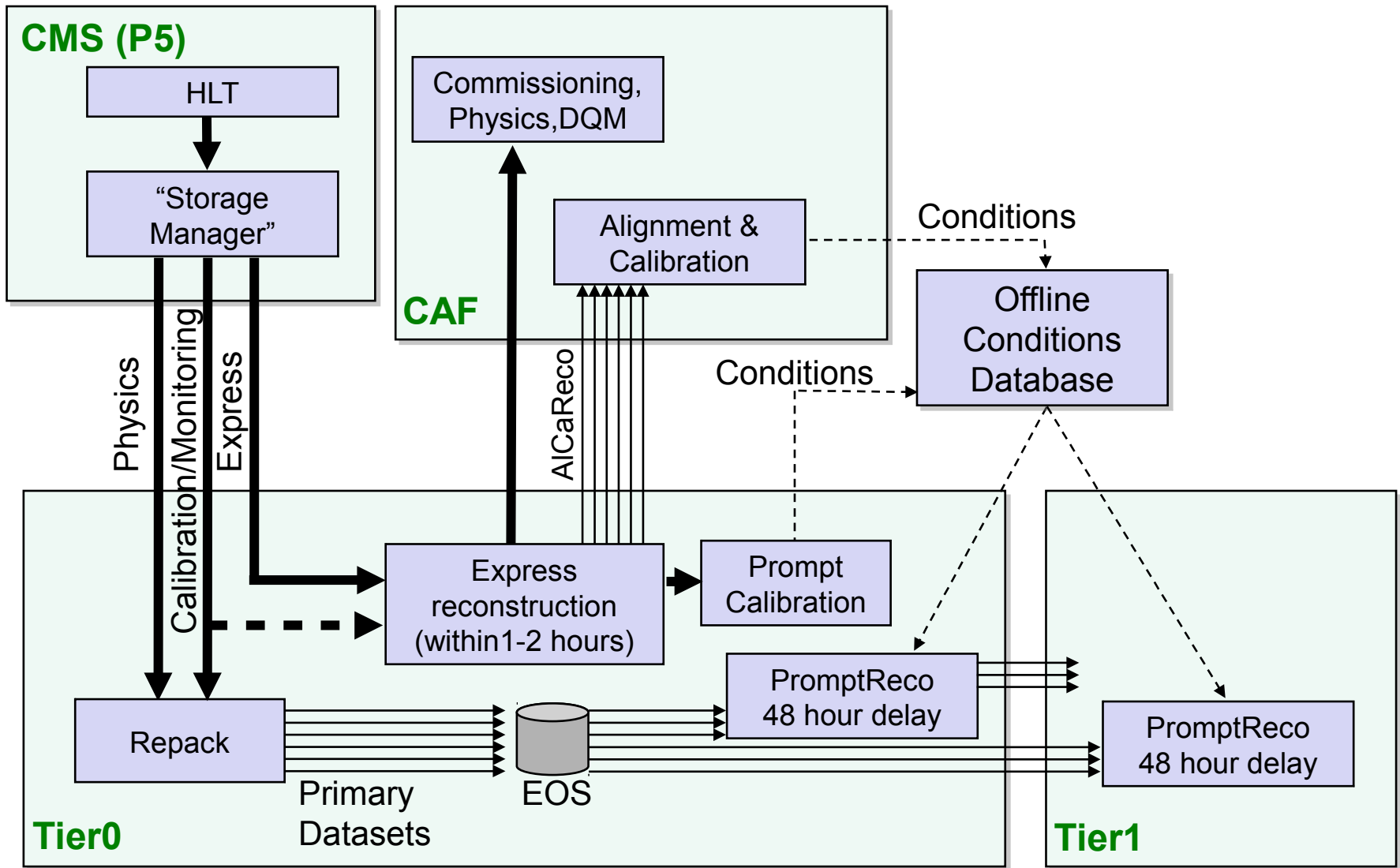
Overview

- Changes for the Tier0 between Run 1 and Run 2
- CERN Agile Infrastructure (in GlideInWMS)
- Operational experience / scale tests

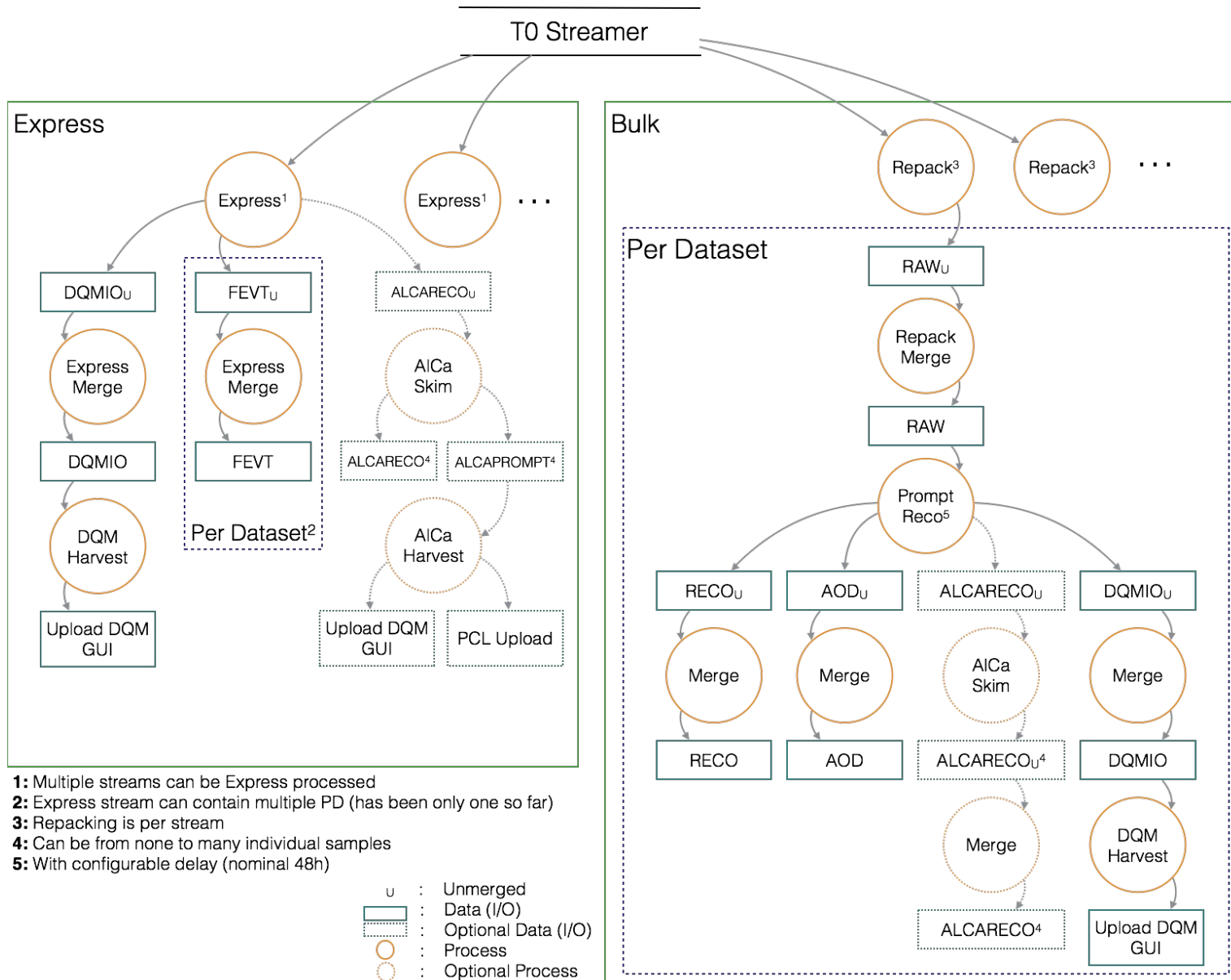
Tier0 Big Picture Run1



Tier0 Big Picture Run 2



Tier0 Processing steps



Tier0 Run 1 => Run 2 changes

- Use same job submission system as rest of CMS

LSF (direct bsub) => HTCondor/GlideInWMS

- Use sharable resources

dedicated LSF => CERN AI cloud (Openstack)

- Integrate Tier1 resources into Tier0 processing

CERN => CERN + fraction of T1 for PromptReco

- Disk/Tape separation

write to castor disk/tape

=> write to EOS disk, separate transfer to tape

CERN Agile Infrastructure (in GlideInWMS)

- It is a CERN internal (Openstack based) cloud.
- It is fully integrated into GlideInWMS system, demand for more pilots triggers creation of VM (with pilot inside) in the CERN AI Tier0 Project.
- We use custom VM (no integration with CERN setup).
- No dynamic resource scheduling on CERN AI:
 - => configure pilot/VM with 1 month lifetime (might be revised downwards)
 - => expect mostly steady state operations (#VM == project quota), only replacing expiring VM

GlidenInWMS setup from a Tier0 viewpoint

- CERN resources used by the Tier0 are exclusively custom CMS VM in CERN AI, we do not use grid/LSF.
- These resources are in a separate HTCondor Tier0 pool.
- In general all (other) CMS resources are in a Global pool, the Tier0 will submit to Tier1 resources via flocking to the Global pool.
- Separate Tier0 pool is a “safety” measure in case there are (scaling) problems with the Global pool. Long term we might move the CERN Tier0 resources into the Global pool too.

Multicore

- CMS has switched to a multi-threaded framework and all PromptReco is expected to run in multicore mode.
- As such, the Tier0 needs access to multicore resources at CERN and the Tier1.
- CERN resources are multicore VM, each VM corresponds to a single multicore partitionable pilot.
 - => no problem, we have full control here
- The Tier1 were asked to provide at least 50% of their pledged resources via multicore pilots as well.
 - => not complete for all Tier1, still work in progress

Operational experience

- We have run Tier0 production instances submitting jobs through glideInWMS to the CERN AI since July 2014 for CMS Cosmics data taking, testing was ongoing even earlier than this.
- The current production deployment for the Tier0 is on its final 2015 hardware (in terms of head node, database instance, etc - although not with the final resource allocation in CERN AI) and has been supporting data taking continuously since beginning of February 2015.

Scale test CERN / T1

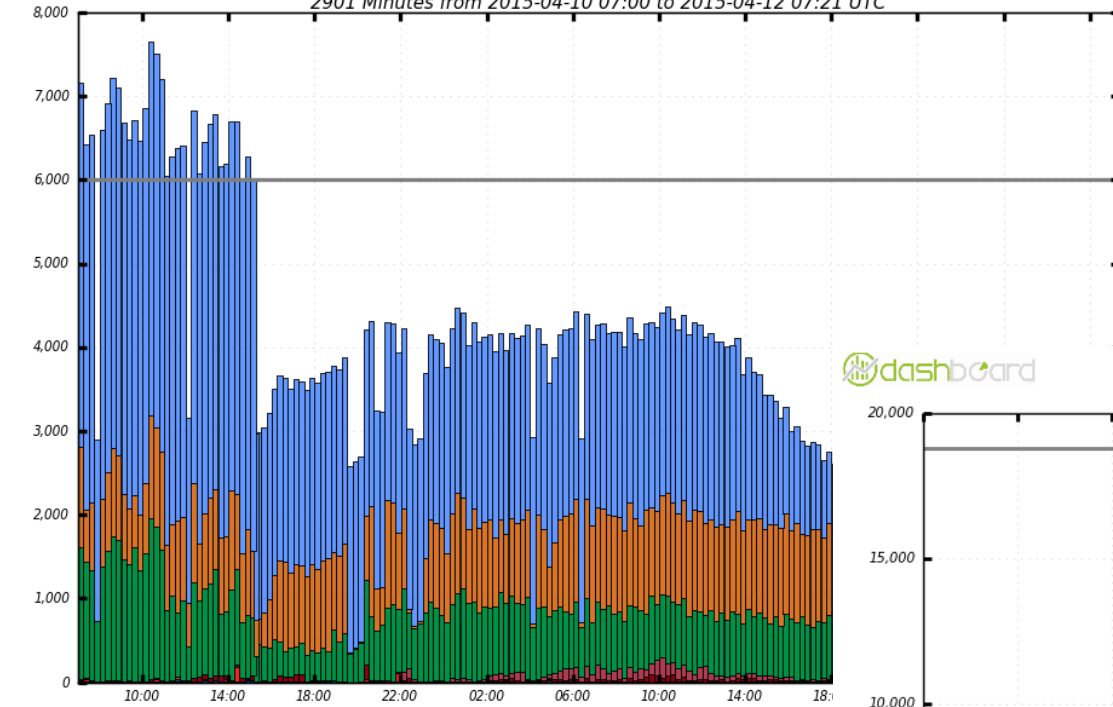
- In November/December 2014 we ran the system flat out replaying lots of 2012/2013 data that was saved for this purpose.
- We have utilized our current allocation of 9000 cores on CERN AI (just in the last week).
- In March 2015 we ran scale tests that replayed the same data but submitted all PromptReco jobs to all CMS T1.
- We reached a peak of almost 50% of all CMS T1 resources, although the per site fractions had a wide variety.

Scale test CERN / T1



Running jobs

2901 Minutes from 2015-04-10 07:00 to 2015-04-12 07:21 UTC



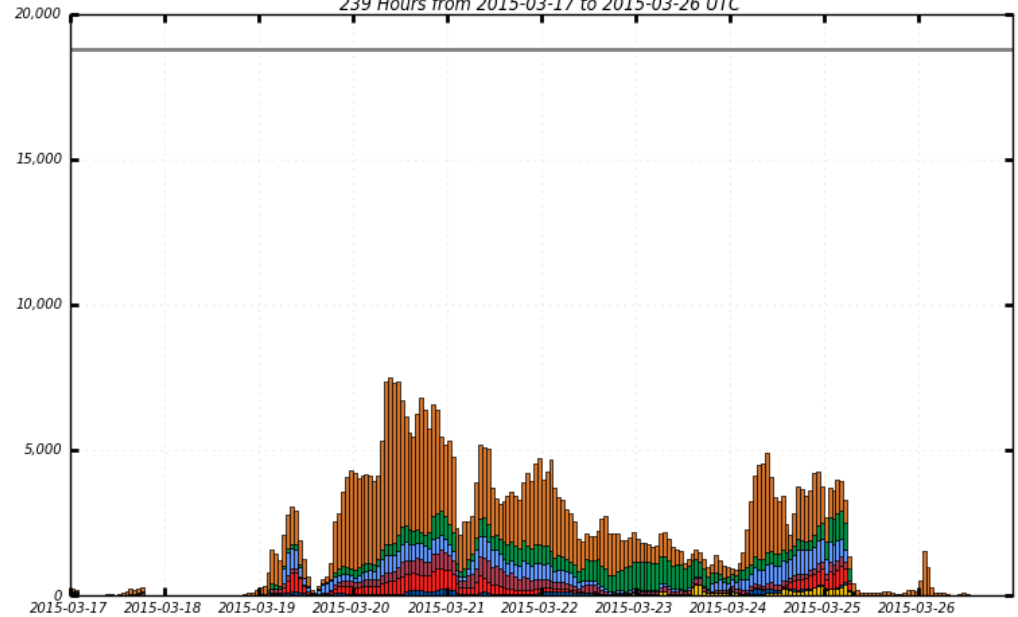
■ Reco ■ Repack ■ Merge ■ Express
■ Cleanup ■ LogCollect ■ Harvest

Maximum: 7,654 , Minimum: 0.00 , Average: 3,952 , Current: 2,842



Slots of Running jobs

239 Hours from 2015-03-17 to 2015-03-26 UTC



■ T1_US_FNAL ■ T1_RU_JINR ■ T1_UK_RAL ■ T1_ES_PIC ■ T1_DE_KIT
■ T1_IT_CNAF ■ T1_FR_CCIN2P3

Maximum: 7,504 , Minimum: 0.00 , Average: 2,067 , Current: 64.00

Upcoming scale test CERN + T1

- Keep pushing scale tests on the 9000 cores in the CERN AI Tier0 project (final quota expected to be in the 12k to 15k core range).
- We are working on improving the T1 situation, addressing some problems that prevented us from utilizing more resources in the last test.
- Upcoming scale tests will combine PromptReco at CERN **AND** T1, running it from the same Tier0 instance.

Summary

- Much has changed for the CMS Tier0 between Run 1 and Run 2. The new setup is more integrated into the general CMS production infrastructure, which will allow a more flexible use of resources.
- We are supporting ongoing cosmics and beam splash data taking with a fully functional Tier0.
- We are scaling up tests to make sure we are ready for Run 2 collision data.
- We make direct use of the new CERN AI cloud based resources, via direct connection of the glideInWMS system with the Openstack cloud controller.

Backup

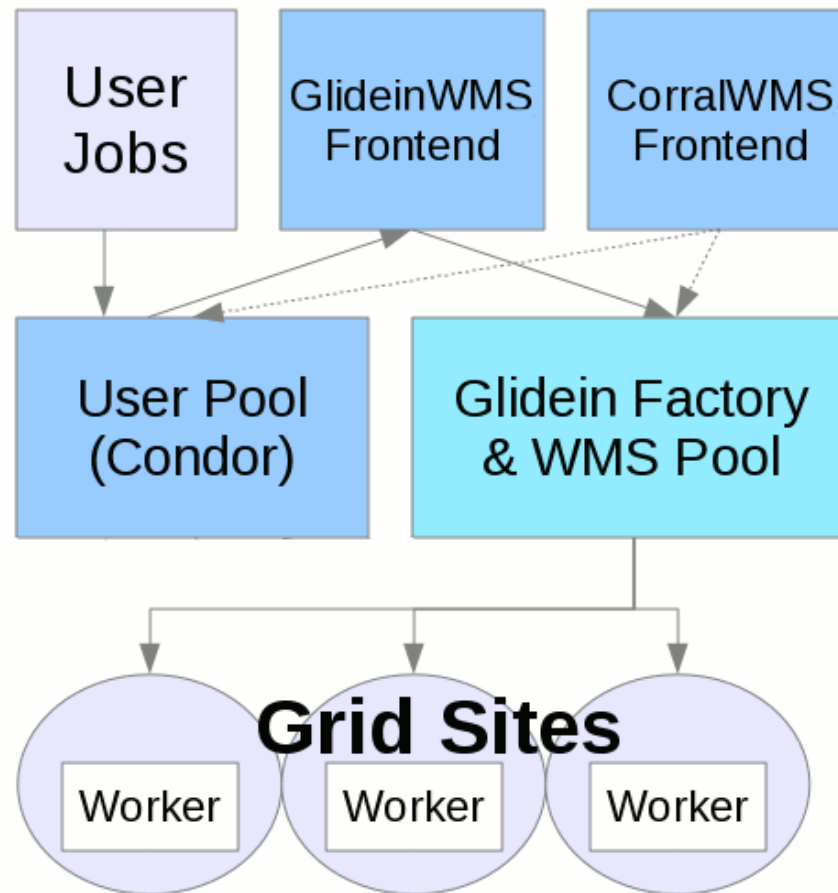
Introduction of terms

- GlideInWMS : pilot based workflow management system used by CMS, makes grid (and other) resources accessible via HTCondor (see next slide)
- Repack : CMS workflow converting detector output (which is a transient data format) to RAW (our custodial raw data format)
- Express : fast processing of a fraction of the data for quick feedback monitoring, prompt calibration etc
- PromptReco : prompt reconstruction of raw data 48 hours after data taking

GlideInWMS

<http://www.uscms.org/SoftwareComputing/Grid/WMS/glideinWMS/doc.prd/index.html>

Follow the link and
check out animation



Why 1 month lifetime for CERN AI VM ?

- Tests in 2014 showed that Openstack would only support VM instantiation through the EC2 interface every few minutes.
- Long VM lifetime ensured enough running VM.
- Current Openstack setup supports VM instantiation through the EC2 interface every few seconds. We have adjusted our timeouts to 10 seconds, but haven't changed the VM lifetime yet.
- No dynamic AI resource scheduling, therefore not many advantages to cycling VM fast. Will see how this goes operationally and if 1 month causes problems (could depending on how often we need to update the VM).