

Optimization of the job management in a multi-queue environment

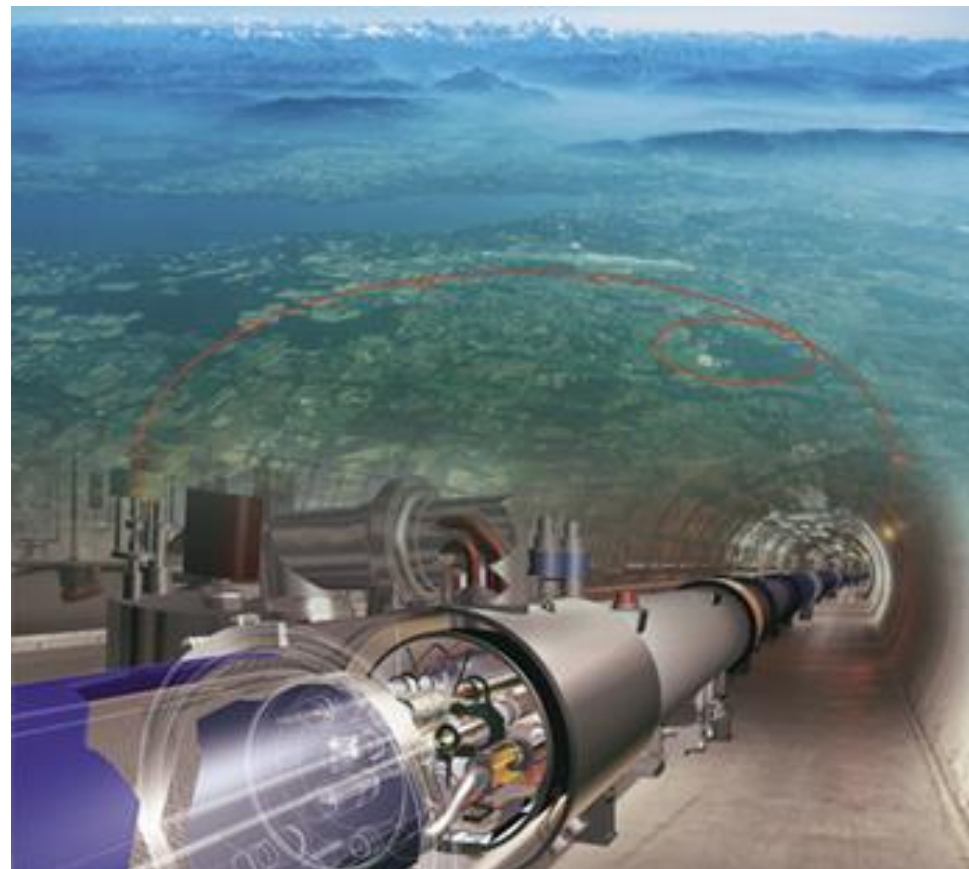
Dr. Mihai Ciubancan IFIN-HH, Dr. Mihnea Dulea IFIN-HH

Optimization of the job management in a multi-queue environment

Worldwide LHC Computing Grid (WLCG) – provides the HTC infrastructure for the 4 main experiments (ALICE, ATLAS, CMS, LHCb) at the the Large Hadron Collider (LHC) in CERN

2017 report:

- 40PB raw data stored in WLCG
- 70PB raw+ simulations + analysis stored in WLCG
- Global transfer rate up to 40GB/s



Optimization of the job management in a multi-queue environment

INTRO:

- RO-07-NIPNE part of RO-LCG Federation – member in WLCG collaboration
- Computing resources dedicated to 3 LHC VOs: ALICE, ATLAS, LHCb
- Storage resources dedicated to 3 LHC VOs: ALICE, ATLAS, LHCb(EOS+DPM)
- 3 different resource managers: PBS/Torque+Maui, SLURM,HTCondor
- 5 subclusters , 8 queues ,3 multicore queues of 8cores
- The single Romanian site running HTCondor and Docker
- The single Romanian site providing EOS storage
- DPM storage used for Romanian ATLAS diskless sites
- Part of LHCONE network (20Gbps connectivity)

Optimization of the job management in a multi-queue environment

RO-07-NIPNE: HARDWARE

Computing infrastructure

- APC InRow Chilled Water Cooling
- 160KVA UPS
- Around 4380 CPU(~230 nodes)
- Blade + "pizza boxes"
- 8,12, 16, 20, 32 cores/server

DC1->



Optimization of the job management in a multi-queue environment

RO-07-NIPNE: HARDWARE

Storage infrastructure

- 4x80KVA UPS Emerson
- 10 DPM servers+1 EOS server
- 2,1PB total capacity
- 1,8PB used capacity

Network infrastructure:

- 100Gbps link between DC1 and DC2

DC2->



Optimization of the job management in a multi-queue environment

RO-07-NIPNE: SOFTWARE

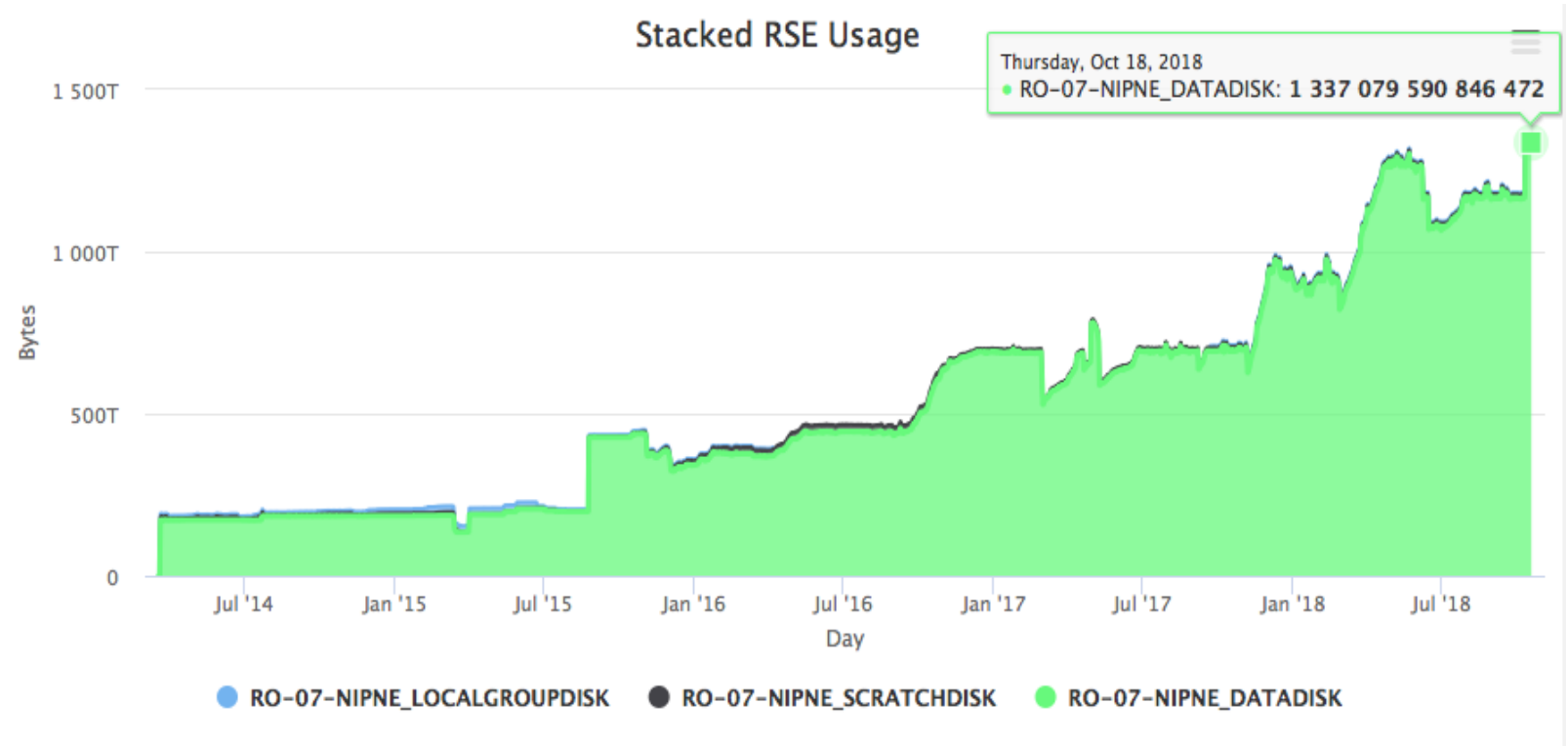
- Scientific Linux 6 /Centos7, UMD3 middleware ,
- 3 CREAM + 2 ARC-CE as job management service
- 8 queues (PBS/Torque + MAUI, SLURM, HTCondor)
- Disk Pool Manager(DPM) for 10 disk storage –shared between ATLAS and LHCb
- EOS - 1FST(File Server storage) –dedicated to ALICE
- Top BDII,
- Site BDII
- VOMS (for local VOs, ex ELI-NP)
- VOBOX(ALICE)
- CVMFS for all WLCG VOs

Optimization of the job management in a multi-queue environment

2018 upgrades:

Processing: ARC-CE
HTCondor + Docker –
dedicated to ATLAS, 2
queues (multicore
processing, analysis) -
~800 cores

Storage: DPM-Disk –
configured with puppet –
400TB – dedicate to ATLAS
Network: 100Gbps link
between DC1 and DC2
(DELL S4128)



Optimization of the job management in a multi-queue environment

2018 problems:

- storage hardware failure –lost 20TB –affected ATLAS and LHCb
- data reading errors from the storage –because of the increased number of ATLAS analysis jobs
- IOwait reached 90% of CPU time on heavy load - led to large amount of failed analysis jobs

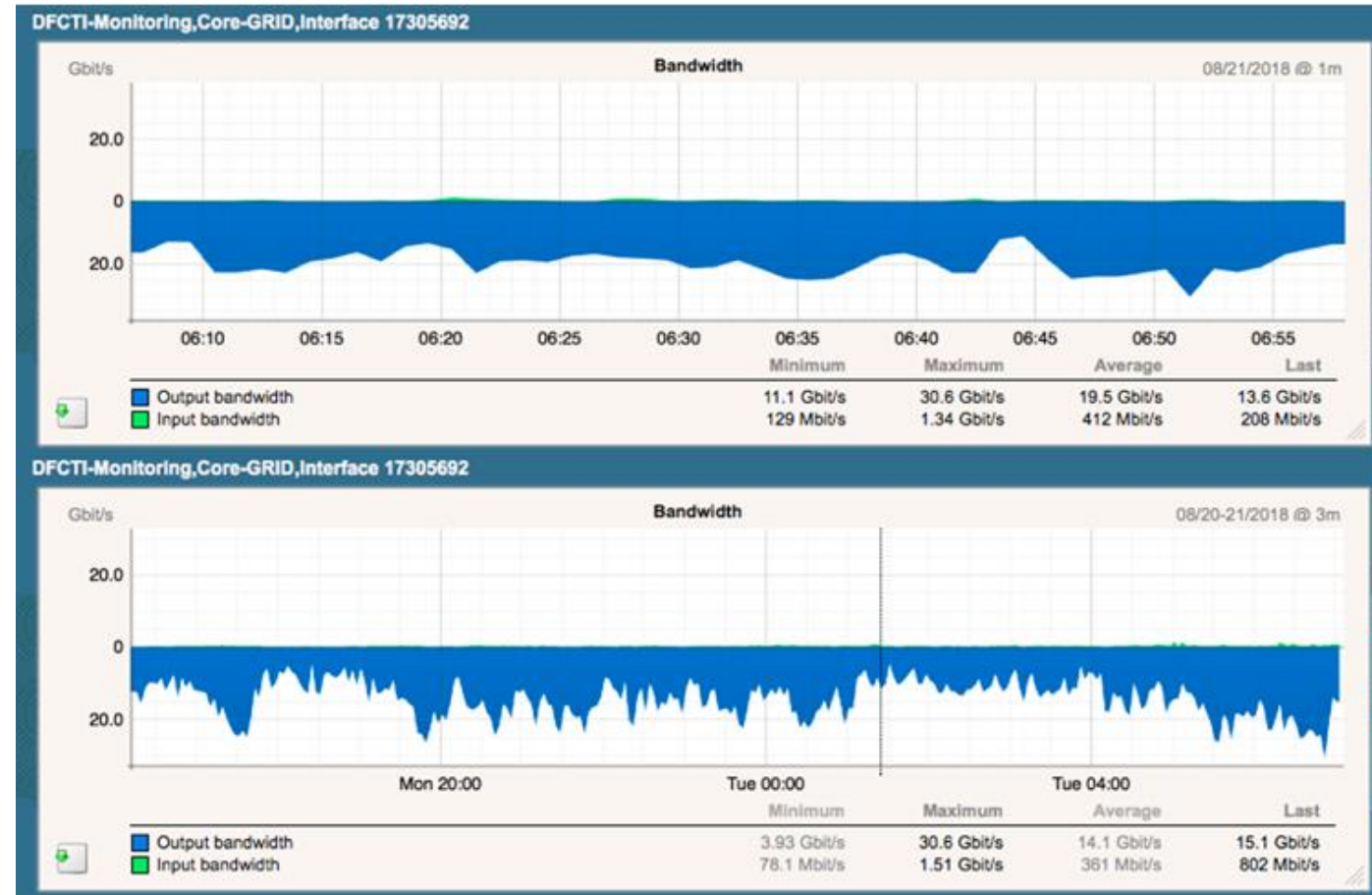
Solutions:

- changing the I/O schedulers of the O.S.
- using *deadline scheduler* instead of *CFQ(Complete Fair Queing) scheduler*(the O.S. default scheduler), providing a higher throughput rate
- increased the maximal I/O queue size for the storage devices (`nr_request` - specifies the maximum number of read and write requests that can be simultaneously queued) – from 128 to 1024
- iowait time has decreased from 90% to 40%-45% - done/failed jobs ratio have significantly improved

Optimization of the job management in a multi-queue environment

Peak traffic between DC1 and DC2: 30.6Gbps

Average traffic for ~1h: 19.5Gbps

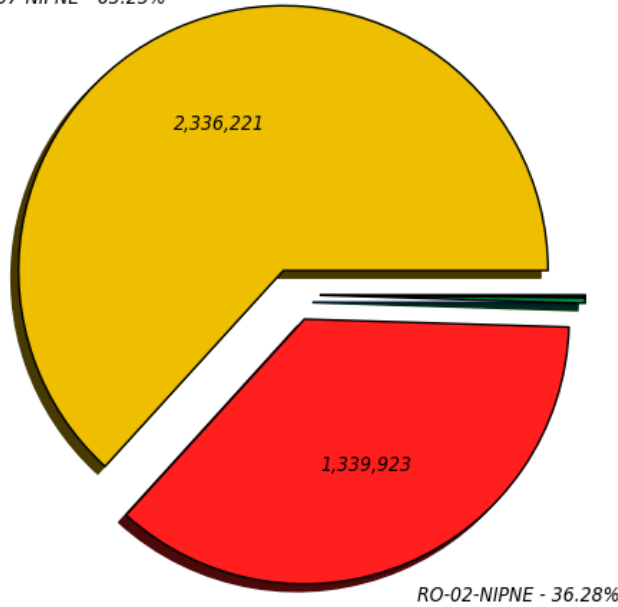


Optimization of the job management in a multi-queue environment



NBytes Processed in GBs (Pie Graph) (Sum: 3,693,729)

RO-07-NIPNE - 63.25%



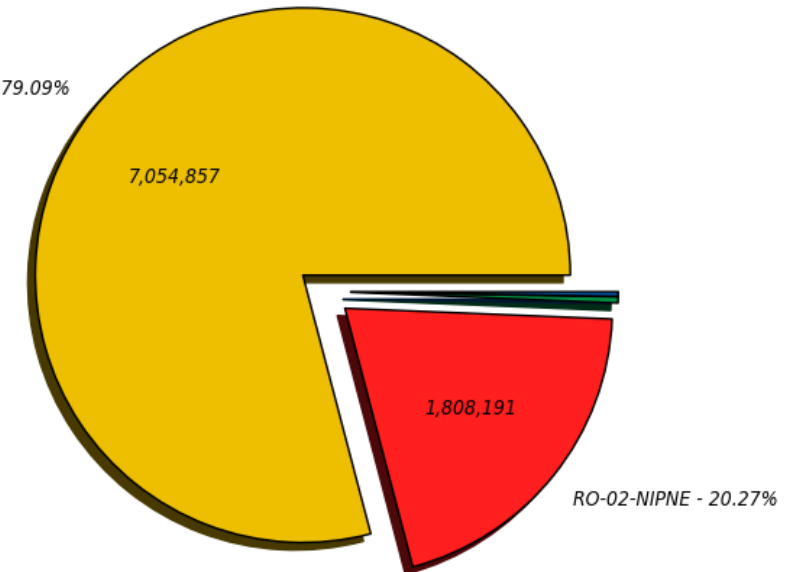
2017
2,33PB

■ RO-07-NIPNE - 63.25% (2,336,222)
■ RO-02-NIPNE - 36.28% (1,339,924)
■ RO-16-UAIC - 0.33% (12,328)
■ RO-14-ITIM - 0.14% (5,255)



NBytes Processed in GBs (Pie Graph) (Sum: 8,920,084)

RO-07-NIPNE - 79.09%



2018
7PB

■ RO-07-NIPNE - 79.09% (7,054,858)
■ RO-02-NIPNE - 20.27% (1,808,191)
■ RO-16-UAIC - 0.37% (32,728)
■ RO-14-ITIM - 0.27% (24,307)

Optimization of the job management in a multi-queue environment

CONCLUSIONS & PLANS:

ALICE & ATLAS & LHCb:

- Move from CREAMCE with Torque+Maui clusters to ARC-CE with HTCondor(in a virtualized environment)
- Increase the computing resources for the VOs
- Increase the storage capacity for the VOs
- Deploy the IPv6 stack for the storage services

**THANK YOU
FOR YOUR ATTENTION!**