# Extending WLCG Tier-2 Resources using HPC and Cloud Solutions

**Jiří Chudoba**[*]
*Institute of Physics of the CAS, Prague*
*E-mail:* Jiri.Chudoba@cern.ch

**Michal Svatoš**
*Institute of Physics of the CAS, Prague*
*E-mail:* Michal.Svatos@cern.ch

Available computing resources limit data simulation and processing of LHC experiments. WLCG Tier centers connected via Grid provide majority of computing and storage capacities, which allow relatively fast and precise analyses of data. Requirements on the number of simulated events must be often reduced to meet installed capacities. Projection of requirements for future LHC runs shows a significant shortage of standard Grid resources if a flat budget is assumed. There are several activities exploring other sources of computing power for LHC projects. The most significant are big HPC centers (supercomputers) and Cloud resources provided both by commercial and academic institutions. The Tier-2 center hosted by the Institute of Physics (FZU) in Prague provides resources for ALICE and ATLAS collaborations on behalf of all involved Czech institutions. Financial resources provided by funding agencies and resources provided by IoP do not allow to buy enough servers to meet demands of experiments. We extend storage resources by two distant sites with additional finance sources. Xrootd servers in the Institute of Nuclear Physics in Rez near Prague store files for the ALICE experiment. CESNET data storage group operates dCache instance with a tape backend for ATLAS (and Pierre Auger Observatory) collaboration. Relatively big computing capacities could be used in the national supercomputing center IT4I in Ostrava. Within the ATLAS collaboration, we explore two different solutions to overcome technical problems arising from different computing environment on the supercomputer. The main difference is that individual worker nodes do not have an external network connection and cannot directly download input and upload output data. One solution is already used for HPC centers in the USA, but until now requires significant adjustments of procedures used for standard ATLAS production. Another solution is based on ARC CE hosted by the Tier-2 center at IoP and resubmission of jobs remotely via ssh.

---

[*]Speaker.

## 1. Introduction

LHC experiments use a world-wide network of computing centers connected by several flavors of grid middleware to cover their computing requirements. Centers pledge available resources to given experiments. Because computing resources are often shared among several different user groups (Virtual Organizations), many VOs can use resources above pledges if their computing systems continuously submit enough tasks. This is also the case of LHC experiments, where computing requirements are higher than pledges. A projection of requirements for future LHC runs shows a significant shortage of standard Grid resources if a flat budget is assumed [[1]]. With expected 20% annual increase of capacity we see almost a factor of about 6 between requirements and pledges for HL-LHC run 4 (Figure 1b). Therefore exploration of additional resources like HPC centers and public or private clouds is crucial for success of the HL-LHC experimental program.



(a) Expected increase of collected luminosity of the ATLAS experiment

(b) Expected CPU resources needed until the end of Run4 of the LHC

Figure 1: Projection of future computing needs of the ATLAS experiment.

## 2. Tier-2 resources

The only Tier-2 center in the Czech Republic is located in the Computing Center of the Institute of Physics of the Czech Academy of Sciences in Prague. Computing resources are shared among several Virtual Organizations; the highest share is for the ATLAS experiment (almost 60%) followed by astroparticle projects Pierre Auger Observatory and Cherenkov Telescope Array (together 20%), neutrino experiment NOvA (10%) and another LHC experiment ALICE (10%). Storage resources are separated; the ALICE VO uses dedicated xrootd servers and the ATLAS VO uses reservations via spacetokens in the DPM system shared with astroparticle projects.

### 2.1 Farm

The farm contains about 7000 cores. The praguelcg2 uses ARC-CE [2] as a compute element and HTCondor [3] as batch system to access them. Grid jobs use ARC-CE+HTCondor while jobs of local users are submitted directly via HTCondor.
For the ATLAS experiment, several different queues are defined fulfilling various requirement. Queue praguelcg2_fzu_SCORE is a single core queue, i.e. it accepts jobs which require one core and maximum of 2 GB of RAM. Queue praguelcg2_fzu_MCORE accepts jobs which require eight

cores and maximum of 16 GB of RAM (thus fullfilling the WLCG requirement of 2 GB of RAM per core). These two queues are available for "production" jobs, i.e. grid jobs submitted by AT-LAS production managers to produce data to be analysed by physicists. Queue ANALY_FZU is accepting jobs which require one core and maximum of 2 GB of RAM. It is used for "analysis" jobs, i.e. grid jobs run by physicists on pre-prepared input datasets. Anselm_MCORE and praguelcg2_IT4I_MCORE are HPC queues which will be described in detail 4.

Computing power of the farm is illustrated on the following figures. Figure 4 shows that site is able to run more than 2k jobs using more than 4k cores. CPU consumption can reach more than 300M seconds per day (Figure 5a). CPU efficiency fluctuate between 0.7 and 1 (Figure 5b). Size-wise, most of input is processed by praguelcg2_fzu_MCORE and ANALY_FZU (Figure 6a). Most of output is produced by praguelcg2_fzu_MCORE and praguelcg2_fzu_SCORE (Figure 6b). This can be explained by the fact that analysis jobs running on the ANALY_FZU process a lot of data but filter out most of it and produce only small files used to make final plots for analysis. On the other hand, event generation, which often needs no input files, runs often on the praguelcg2_fzu_SCORE.

To manage the storage, site uses Disk Pool Manager (DPM) [4]. It currently controls about 2.5 PB of disk space but the size will increase to about 4 PB later this year. Most of the space is assigned to ATLAS datadisk (see figure 2) which contains inputs and outputs of ATLAS production jobs.

## 3. CESNET e-Infrastructure

Distributed capacities require a reliable and performant network connection. CESNET is the Czech NREN (National Research and Education Network) provider and operates network connections for the Tier-2 center. The Tier-2 center at FZU has a dedicated 2x10 Gbps network link to LHCONE [6] and 10 Gbps link to standard Geant network. Local Czech network Czechlight connects several high energy physics institutions in Prague and Nuclear Physics Institute in Rez close to Prague. This connection enabled to add xrootd disk servers running in Rez to the xrootd cluster at FZU. Users see just one xrootd instance and they do not have to care about physical disk server location. This extension adds additional storage capacity which could not be hosted directly at FZU. The Czechlight network is also used to include compute servers located at the Faculty of Mathematics and Physics at Charles University. Jobs to these servers are directly submitted by HTCondor instance at FZU and process mostly jobs for NOvA experiment with lower input and output requirements. This solution eliminate need for operations of another small cluster with its own batch system and full grid services.

Sufficient network connection was important for usage of storage capacities of the CESNET DataCare department. This group operates facilities now in 4 different locations in the Czech Republic with a total capacity over 21 PB. They installed dCache [7] headnode on a server in Pilsen and disk servers (dCache poolnodes) in Pilsen and later in Jihlava. This dCache instance is published to the EGI grid infrastructure via Tier-2 site at FZU because the CESNET DataCare unit does not operate own grid site. Total available disk capacity for the ATLAS experiment 20 TB was used by Czech ATLAS users when the disk space at Tier-2 for local users was small. Now users take advantage of tape backends of the DataCare department facilities exported as ATLASLO-CALGROUPTAPE spacetoken. They use this spacetoken for a backup of private (it means not

produced by the central ATLAS team or on behalf of any ATLAS group) datasets with only one copy on disks on ATLASLOCALGROUPDISK spacetoken at the Tier-2. A planned movement of DataCare facilities from Pilsen to Ostrava will be fully transparent for users.
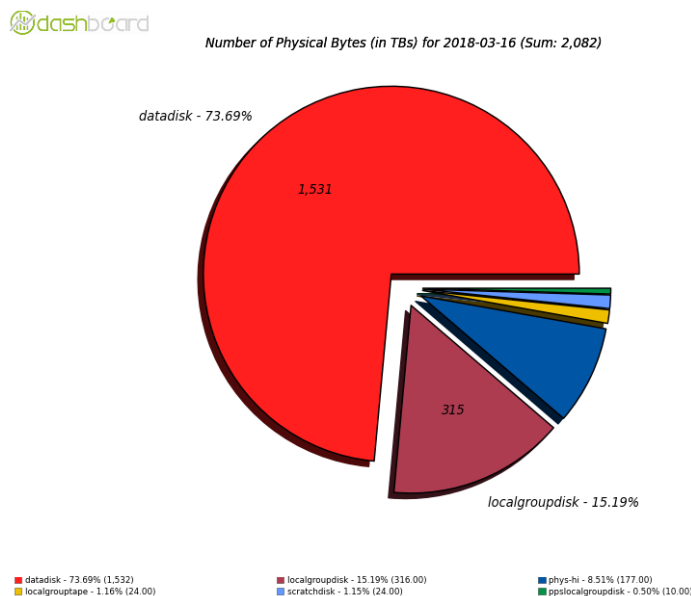


Figure 2: Storage used at FZU Tier-2 by ATLAS spacetokens. LOCALGROUPTAPE and PPSLO-CALGROUPDISK are provided by the CESNET DataCare group.

## 4. HPC e-Infrastructure

Czech national supercomputer center IT4Innovations (IT4I) in Ostrava provides access to two HPC systems - Anselm and Salomon. They run CentOS with PBSpro as a batch system and Lustre for shared filesystem. Access is provided via login nodes. Worker nodes have very limited connectivity to outside world (only http is allowed). They are used to run ATLAS production jobs. Anselm was build in 2013 providing 93 TFLOPs. It consists of 209 worked nodes, providing 16 cores and 64 GB of RAM per node. It has Infiniband QDR and Gigabit Ethernet. It uses job submission system originally developed for Titan [8] via Anselm_MCORE queue. This system will not be described here. Salomon was build in 2017 providing 2 PFLOPs in peak. It consists of 1008 compute nodes providing 24 cores and 128 GB of RAM per node. Nodes are connected by Infiniband (56 Gbps). We use ARC-CE installed at FZU to submit jobs to Salomon cluster. Details of the submission system follow.

### 4.1 Job submission

Jobs are submitted to Salomon via ARC-CE installed at praguelcg2. The ARC-CE accepts job from ARC Control Tower (aCT), authorize the user and translates the job requirements into a script runnable in PBSpro. It also downloads input files which are put into the session directory together

with run script and later uploads output files. The submission scripts of ARC-CE were modified so it can submit jobs via ssh to login node. Session and runtime directories of the ARC-CE are mounted via sshfs to dedicated scratch space on Salomon's Lustre filesystem. When ARC-CE executes job submission, script with workload is submitted to the PBSpro via ssh on the login node. PBSpro takes input from the session directory and puts output there when the job finishes.

## 4.2 Software installation

ATLAS jobs use many software packages. At a grid site, the whole software stack is accesible via CVMFS [9]. The HPC provides access only to login nodes, therefore there are no ATLAS specific services (like CVMFS) running. To work around this problem, CVMFS is mounted at the ARC-CE and the software is rsynced to shared Lustre directory on Salomon. As all ATLAS software releases represent huge amount of data, only sub-release 21.0 is rsynced once a day to Salomon. This represents about half TB of data in about 10M files.

## 4.3 Consumed resources

The PBSpro configuration of Salomon allows maximum of 100 jobs submitted to the used queue free, i.e. 100 jobs in any state (running, queueing, etc.). This can be seen on figure 3 where number of jobs running in the ARC-CE reaches maximum of 100 and then forms a plateau.
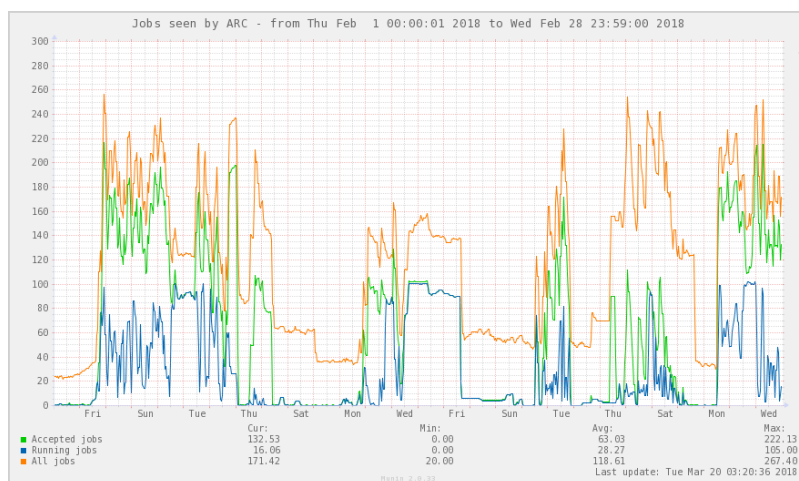


Figure 3: Job statistics provided by the ARC-CE. Running jobs are jobs submitted into the batch system. Accepted jobs are running jobs and jobs which are being prepared by the CE. All jobs include also deleted jobs (but with logs still available).

While number of grid jobs running at the praguelcg2_IT4I_MCORE is low (Figure 4a) the amount of cores and CPU time used (Figure 4b and 5a) is significant in comparison with praguelcg2. This is because Salomon's worker nodes provide 24 cores while praguelcg2 provides only 8 cores. The CPU efficiency is also comparable (Figure 5b).
ATLAS uses many workloads but only simulation is used on Salomon. Simulation workload has best ratio of CPU utilization to I/O. While figure 4b and 5a show significant contribution of computing resources, the amount of input files processed (Figure 6a) and output files produced (Figure 6b) is very small.

(a) Amount of running jobs



(b) Amount of slots of running jobs

Figure 4: Jobs and jobslots during February 2018



(a) CPU consumption of jobs



(b) CPU efficiency as a function of number of cores job uses during February. One core is used by praguelcg2_fzu_SCORE and ANALY_FZU, eight cores by praguelcg2_fzu_MCORE and 24 cores by praguelcg2_IT4I_MCORE (HPC queue).

Figure 5: CPU consumption and efficiency of jobs during February 2018

5

(a) Amount of input processed by jobs during February



(b) Amount of output produced by jobs during February

Figure 6: Size of processed input and produced output of jobs during February

## 5. Summary and conclusions

Computing resources of the Czech Tier-2 site were extended by several types of external resources. The storage include xrootd servers of the Nuclear Physics Institute in Rez and dCache servers (including tape backend) of the CESNET DataCare department. Some ATLAS jobs are transparently submitted via ARC CE at FZU to the national HPC center IT4I in Ostrava. The possibility to use cloud resources provided by CESNET will be investigated later when a planned migration from OpenNebula to OpenStack is finished.

## Acknowledgments

## References

[1] https://twiki.cern.ch/twiki/bin/view/AtlasPublic/ComputingandSoftwarePublicResults

[2] M.Ellert et al., *Advanced Resource Connector middleware for lightweight computational Grids*, Future Generation Computer Systems 23 (2007) 219-240.

[3] D. Thain, T. Tannenbaum, and M. Livny. *Distributed Computing in Practice: The Condor Experience*. Concurrency and Computation: Practice and Experience, Vol. 17, No. 2-4, pages 323-356

[4] http://lcgdm.web.cern.ch/dpm

[5] http://boinc.berkeley.edu/

[6] https://twiki.cern.ch/twiki/bin/view/LHCONE/WebHome

[7] `https://www.dcache.org/`

[8] K. De et al., *Integration of PanDA workload management system with Titan supercomputer at OLCF*, 2015 J. Phys. Conf. Ser. **664** 092020

[9] J. Blomer et al., *Distributing LHC application software and conditions databases using the CernVM file system*, 2011 J. Phys.: Conf. Ser. **331** 042003

[10] A. Filipcic at al., *The ATLAS ARC backend to HPC*, 2015 J. Phys.: Conf. Ser. **664** 062057

[11] M. Hostettler et al., *ATLAS computing on CSCS HPC*, 2015 J. Phys.: Conf. Ser. **664** 092011