

# **HEPiX Spring 2019 Workshop**

Monday 25 March 2019 - Friday 29 March 2019

SDSC Auditorium

## **Book of Abstracts**



# Contents

The Experience and Challenge in Grid Computing at KEK 1 . . . . .	1
CERN Site Report 2 . . . . .	1
INFN-T1 Site Report 3 . . . . .	1
Diamond Light Source Site Report 4 . . . . .	1
Cost and system performance modelling in WLCG and HSF: an update 5 . . . . .	2
Computing/Storage/Networking for next generation photon science experiments @DESY 6 . . . . .	2
Computer Security Update 7 . . . . .	2
KEK Site Report 8 . . . . .	3
Prague Site Report 9 . . . . .	3
Text Classification via Supervised Machine Learning for an Issue Tracking System 10 . .	4
OpenAFS Release Team report 11 . . . . .	4
GSI Site Report 12 . . . . .	4
Endpoint user device and app provisioning models 13 . . . . .	5
Evolution of interactive data analysis for HEP at CERN –SWAN, Kubernetes, Apache Spark and RDataFrame 14 . . . . .	5
Swiss HPC Tier-2 Computing @ CSCS 15 . . . . .	6
6 years of CERN Cloud - From 0 to 300k cores 17 . . . . .	6
BNL Site Report 18 . . . . .	6
Creating an opportunistic OSG site inside the PRP Kubernetes cluster 19 . . . . .	7
How Fair is my Fair-Sharing? Exposing Some Hidden Behavior Through Workload Analy- sis 20 . . . . .	7
Keeping Pace with Science: How a Modern Filesystem Can Accelerate Discovery 21 . . .	8
PDSF Site Report 22 . . . . .	8
Public cloud for high throughput computing 23 . . . . .	9

Benchmarking Worrrking Group - Status Report 24 . . . . .	9
Omni, N9 and the Superfacility 25 . . . . .	10
Nikhef Site Report 26 . . . . .	10
Storage services at CERN 27 . . . . .	10
CERN DNS and DHCP service improvement plans 28 . . . . .	11
Welcome to SDSC and UCSD 29 . . . . .	11
Using the dynafed data federation as site storage element 30 . . . . .	11
University of Wisconsin-Madison CMS T2 site report 31 . . . . .	12
HEPiX Technology Watch working group 33 . . . . .	12
Update of Canadian T1 / T2 37 . . . . .	13
CloudScheduler version2 38 . . . . .	13
Jupyter at SDCC 39 . . . . .	13
RAL Site Report 40 . . . . .	13
The difference in network equipment in the 25/100/400G era and how to test/break them 41 . . . . .	14
OSiRIS: Open Storage Research Infrastructure 42 . . . . .	14
IPv6 & WLCG - an update from the HEPiX IPv6 Working Group 43 . . . . .	14
Storage management in a large scale at BNL 44 . . . . .	15
Developing for a Services Layer At The Edge (SLATE) 45 . . . . .	15
Config Management and Deployment Setup at KIT 46 . . . . .	16
Virtualization for Online Storage Clusters 47 . . . . .	16
Network Functions Virtualisation Working Group Update 48 . . . . .	16
WLCG/OSG Network Activities, Status and Plans 49 . . . . .	17
AGLT2 Site Report 50 . . . . .	18
RACF/SDCC Datacenter Transformation within the Scope of BNL CFR Project and Beyond 51 . . . . .	18
How to make your Cluster look like a Supercomputer (for Fun and Profit) 52 . . . . .	18
Tokyo Tier-2 Site Report 53 . . . . .	19
BNL activities on federated access and Single Sign-On 54 . . . . .	19
BEIJING Site Report 55 . . . . .	19

DPM DOME 57 . . . . .	20
Addressing the Challenges of Executing Massive Computational Clusters in the Cloud 58	20
IntegratingHadoop Distributed File System to Logistical Storage 59 . . . . .	21
Deep learning in a container, experience and best practices 60 . . . . .	21
Developments in disk and tape storage at the RAL Tier 1 61 . . . . .	22
DESY site report 62 . . . . .	22
Token Renewal Service (TRS) at SLAC 63 . . . . .	22
University of Nebraska CMS Tier2 Site Report 64 . . . . .	23
UW CENPA site report 65 . . . . .	23
JLab Site Report 66 . . . . .	24
Fermilab Site Report 67 . . . . .	24
Changes to OSG and how they affect US WLCG sites 68 . . . . .	24
The glideinWMS system: recent developments 69 . . . . .	24
Technolog Watch WG Reports 70 . . . . .	25



**Grid, Cloud and Virtualization / 1****The Experience and Challenge in Grid Computing at KEK**

**Authors:** Go Iwai<sup>1</sup>; Hiroyuki Matsunaga<sup>2</sup>; Tomoaki Nakamura<sup>2</sup>; Takashi Sasaki<sup>2</sup>; Soh Suzuki<sup>1</sup>; Wataru Takase<sup>2</sup>

<sup>1</sup> *KEK*

<sup>2</sup> *High Energy Accelerator Research Organization (JP)*

**Corresponding Authors:** tomoaki.nakamura@kek.jp, go.iwai@kek.jp, hiroyuki.matsunaga@kek.jp, takashi.sasaki@kek.jp, wataru.takase@kek.jp, soh.suzuki@kek.jp

The KEK Central Computer System (KEKCC) is a service, which provides large-scale computer resources, Grid and Cloud computing, as well as common IT services. The KEKCC is entirely replaced every four or five years according to Japanese government procurement policy for the computer system. Current KEKCC has been in operation since September 2016 and decommissioning will start in early 2020.

In this talk, we would like to share our experiences and challenges for the security, operation, and some applications dedicated to each experiment. In particular, we report several improvements on the Grid computing system for the Belle experiment based on the nearly three years operational performance of the KEKCC. We also introduce a prospect for the next KEKCC which is planned to be launched in September 2020.

**Site Reports / 2****CERN Site Report**

**Author:** Andrei Dumitru<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** andrei.dumitru@cern.ch

News from CERN since the HEPiX Autumn/Fall 2018 workshop in Barcelona.

**Site Reports / 3****INFN-T1 Site Report**

**Author:** Andrea Chierici<sup>1</sup>

<sup>1</sup> *Universita e INFN, Bologna (IT)*

**Corresponding Author:** andrea.chierici@cern.ch

An update on what's going on at the Italian Tier1 center

**Site Reports / 4****Diamond Light Source Site Report**

**Author:** James Thorne<sup>1</sup>

<sup>1</sup> *Diamond Light Source*

**Corresponding Author:** j.i.thorne@gmail.com

Diamond Light Source is an X-ray synchrotron light source co-located with STFC RAL in the UK. This is the first site report from Diamond at HEPiX since 2015. The talk will discuss recent changes, current status and future plans as well as the odd disaster story thrown in for good measure.

Diamond has a new data centre, new storage and new compute as well as new staff and a few forays into various cloud providers.

## IT Facilities & Business Continuity / 5

### Cost and system performance modelling in WLCG and HSF: an update

**Authors:** Jose Flix Molina<sup>1</sup>; Markus Schulz<sup>2</sup>; Andrea Sciabà<sup>2</sup>

<sup>1</sup> *Centro de Investigaciones Energéticas Medioambientales y Tecnológicas*

<sup>2</sup> *CERN*

**Corresponding Authors:** jose.flix.molina@cern.ch, markus.schulz@cern.ch, andrea.sciaba@cern.ch

The HSF/WLCG cost and performance modeling working group was established in November 2017 and has since then achieved considerable progress in our understanding of the performance factors of the LHC applications, the estimation of the computing and storage resources and the cost of the infrastructure and its evolution for the WLCG sites. This contribution provides an update on the recent developments of the working group activities, with a special focus on the implications for computing sites.

## Computing & Batch Systems / 6

### Computing/Storage/Networking for next generation photon science experiments @DESY

**Authors:** Martin Gasthuber<sup>1</sup>; Sergey Yakubov<sup>1</sup>

<sup>1</sup> *DESY*

**Corresponding Authors:** martin.gasthuber@cern.ch, sergey.yakubov@desy.de

We will briefly show the current onsite accelerator infrastructure and their resulting computing and storage usage and future requirements. The second section will discuss the plans and work done regarding the hardware infrastructure, the system level middleware (i.e. container, storage connection, networks) and the higher level middleware (under development) covering low latency data access and selection (including metadata queries) directly connection the DAQ at detector level to the data processing code (mostly developed by experimentators running on multiple nodes in parallel). The last section will shortly discuss initial result (technical and non-technical) and possible collaborations with other sites with similar challenges.

## Networking & Security / 7



## Computer Security Update

**Authors:** Stefan Lueders<sup>1</sup>; Computer Security<sup>None</sup>

<sup>1</sup> CERN

**Corresponding Authors:** computer.security@cern.ch, stefan.lueders@cern.ch

This presentation provides an update on the global security landscape since the last HEPiX meeting. It describes the main vectors of risks to and compromises in the academic community including lessons learnt, presents interesting recent attacks while providing recommendations on how to best protect ourselves. It also covers security risks management in general, as well as the security aspects of the current hot topics in computing and around computer security.

This talk is based on contributions and input from the CERN Computer Security Team.

### Site Reports / 8

## KEK Site Report

**Author:** Tomoaki Nakamura<sup>1</sup>

**Co-authors:** Go Iwai <sup>2</sup>; Koichi Murakami ; Tadashi Murakami <sup>2</sup>; Takashi Sasaki <sup>1</sup>; Soh Suzuki ; Wataru Takase <sup>1</sup>

<sup>1</sup> High Energy Accelerator Research Organization (JP)

<sup>2</sup> KEK

**Corresponding Authors:** go.iwai@kek.jp, wataru.takase@kek.jp, koichi.murakami@kek.jp, soh.suzuki@kek.jp, takashi.sasaki@kek.jp, tadashi.murakami@kek.jp, tomoaki.nakamura@kek.jp

We would like to report an update of the computing research center at KEK including the Grid system from the last HEPiX Fall 2018 for the data taking period of SuperKEKB and J-PARC experiments in 2019. The network connectivity of KEK site has been improved by the replacement of network equipment and security devices in September 2018. The situation of the international network for Japan will also be introduced. In addition to the status report, we will present on the preparation for procurement of the next system.

### Site Reports / 9

## Prague Site Report

**Authors:** Jiri Chudoba<sup>1</sup>; Martin Adam<sup>1</sup>

**Co-authors:** Alexandr Mikula <sup>1</sup>; Jana Uhlirova ; Petr Vokac <sup>2</sup>; Dagmar Adamova <sup>1</sup>

<sup>1</sup> Acad. of Sciences of the Czech Rep. (CZ)

<sup>2</sup> Czech Technical University (CZ)

**Corresponding Authors:** alexandr.mikula@cern.ch, dagmar.adamova@cern.ch, jiri.chudoba@cern.ch, martin.adam@cern.ch, petr.vokac@cern.ch

We will give an overview of the site including our recent network redesign. We will dedicate a part of the talk to disk servers: report on the newest additions as well as upgraded old hardware. We will also share experience with our distributed HT-Condor batch system.

**End-User IT Services & Operating Systems / 10****Text Classification via Supervised Machine Learning for an Issue Tracking System****Author:** Martin Kandes<sup>1</sup><sup>1</sup> *Univ. of California San Diego (US)***Corresponding Author:** mkandes@sdsc.edu

Comet is SDSC's newest supercomputer. The result of a \$27M National Science Foundation (NSF) award, Comet delivers over 2.7 petaFLOPS of computing power to scientists, engineers, and researchers all around the world. In fact, within its first 18 months of operation, Comet served over 10,000 unique users across a range of scientific disciplines, becoming one of the most widely used supercomputers in NSF's Extreme Science and Engineering Discover Environment (XSEDE) program ever.

The High-Performance Computing (HPC) User Services Group at SDSC helps manage user support for Comet. This includes, but is not limited to, managing user accounts, answering general user inquiries, debugging technical problems reported by users, and making best practice recommendations on how users can achieve high-performance when running their scientific workloads on Comet. These interactions between Comet's user community and the User Service Group are largely managed through email exchanges tracked by XSEDE's internal issue tracking system. However, while Comet is expected to maintain a 24x7x365 uptime, user support is generally only provided during normal business hours. With such a large user community spread across nearly every timezone, the result is a number of user support tickets submitted during non-business hours waiting between 12 hours to several days for responses from the User Services Group.

The aim of this research project is to use supervised machine learning techniques to perform text classification on Comet's user support tickets. If an efficient classification scheme can be developed, the User Services Group may eventually be able to provide automated email responses to some of the more common user issues reported during non-business hours.

**Storage & Filesystems / 11****OpenAFS Release Team report****Authors:** Michael Meffie<sup>1</sup>; Stephan Wiesand<sup>2</sup>; Benjamin Kaduk<sup>3</sup>; Mark Vitale<sup>1</sup><sup>1</sup> *Sine Nomine*<sup>2</sup> *DESY*<sup>3</sup> *Formerly MIT***Corresponding Authors:** stephan.wiesand@desy.de, mvitale@sinenomine.net, kaduk@mit.edu, mmeffie@sinenomine.net

A report from the OpenAFS Release Team on recent OpenAFS releases and development branch updates. Topics include acknowledgement of contributors, descriptions of issues fixed, updates for new versions of Linux and Solaris, changes currently under review, and an update on the new RXGK security class for improved security.

**Site Reports / 12****GSI Site Report**

**Author:** Thomas Roth<sup>1</sup>

<sup>1</sup> GSI

**Corresponding Author:** t.roth@gsi.de

Ongoing developments at GSI/FAIR: diggers, Lustres, procurements, relocations, operating systems

## End-User IT Services & Operating Systems / 13

### Endpoint user device and app provisioning models

**Author:** Michal Kwiatek<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Author:** michal.kwiatek@cern.ch

Over the last years, there has been a number of trends related to how devices are provisioned and managed within organizations, such as BYOD - “Bring Your Own Device” or COPE: “Company Owned, Personally Enabled”. In response, a new category of products called “Enterprise Mobility Management Suites”, which includes MDM - “Mobile Device Management” and MAM - “Mobile Application Management” emerged on the market. Vendors like VMWare, MobileIron, Microsoft and Citrix now all provide more or less comprehensive systems in this category. But how do these commercial systems correspond to the needs of the Scientific Community?

This talk will summarize current status of device management practices and the strategy for provisioning of devices and applications at CERN. It will also attempt to initiate wider discussion within the community.

## Computing & Batch Systems / 14

### Evolution of interactive data analysis for HEP at CERN –SWAN, Kubernetes, Apache Spark and RDataFrame

**Author:** Piotr Mrowczynski<sup>1</sup>

**Co-authors:** Prasanth Kothuri<sup>1</sup>; Enric Tejedor<sup>1</sup>

<sup>1</sup> CERN

**Corresponding Authors:** piotr.mrowczynski@cern.ch, prasanth.kothuri@cern.ch

This talk is focused on recent experiences and developments in providing data analytics platform SWAN based on Apache Spark for High Energy Physics at CERN.

The Hadoop Service expands its user base for analysts who want to perform analysis with big data technologies - namely Apache Spark –with main users from accelerator operations and infrastructure monitoring. Hadoop Service integration with SWAN Service offers scalable interactive data analysis and visualizations using Jupyter notebooks, with computations being offloaded to compute clusters - on-premise YARN clusters and more recently to cloud-native Kubernetes clusters. The ROOT framework is most widely used tool for high-energy physics analysis. Its integration with SWAN allows physicists to perform web-based interactive analysis using standard tools and libraries, in the cloud.

The first part of presentation will focus on integration of Spark on Kubernetes into SWAN service, which allows to offload computations to elastic, virtualized and container-based infrastructure in the private or public clouds, compared to complex to manage and operate on-premise Hadoop clusters.

The second part will focus on evolutions in exploiting analytics infrastructure - namely new developments in ROOT framework - Distributed RDataFrame - which would allow interactive, parallel and distributed analysis on large physics datasets by transparently exploiting dynamically pluggable resources in SWAN, e.g. Hadoop or Kubernetes clusters.

## Computing & Batch Systems / 15

### Swiss HPC Tier-2 Computing @ CSCS

**Author:** Dino Conciatore<sup>1</sup>

<sup>1</sup> *CSCS (Swiss National Supercomputing Centre)*

**Corresponding Author:** dino.conciatore@cscs.ch

For the past 10 years, CSCS has been running compute capability in the WLCG Tier-2 for ATLAS, CMS and LHCb on standard commodity hardware (a cluster named Phoenix). Three years ago, CSCS began providing this service on the flagship High Performance Computing (HPC) system, Piz Daint (a Cray XC40/50 system). Piz Daint is a world-class HPC system with over 1800 dual-processor multicore nodes and more than 5700 hybrid compute nodes with GPU accelerators. Piz Daint currently holds the 5th position on the Top500 List and is the most powerful HPC system in Europe.

In preparation for future challenges that the HL-LHC will impose on the computing sites, CSCS is in the process of decommissioning the Phoenix cluster and fully consolidating the Tier-2 compute load onto Piz Daint.

In this presentation, the critical milestones to achieve on the road to a successful migration to Piz Daint will be explained.

## Grid, Cloud and Virtualization / 17

### 6 years of CERN Cloud - From 0 to 300k cores

**Author:** Belmiro Moreira<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** belmiro.moreira@cern.ch

CERN, the European Laboratory for Particle Physics, is running OpenStack for its private Cloud Infrastructure among other leading open source tools that helps thousands of scientists around the world to uncover the mysteries of the Universe.

In 2012, CERN started the deployment of its private Cloud Infrastructure using OpenStack. Since then we moved from few hundred cores to a multi-cell deployment spread between two data centres. After 6 years deploying and managing OpenStack at scale, we now look back and discuss the challenges of building a massive scale infrastructure from 0 to +300K cores.

With this talk we will dive into the history, architecture, tools and technical decisions behind the CERN Cloud Infrastructure.

**Site Reports / 18****BNL Site Report**

**Authors:** Ofer Rind<sup>1</sup>; Tony Wong<sup>1</sup>

<sup>1</sup> *Brookhaven National Laboratory*

**Corresponding Authors:** tony@bnl.gov, rind@bnl.gov

News and updates from BNL activities since the Barcelona meeting

**Grid, Cloud and Virtualization / 19****Creating an opportunistic OSG site inside the PRP Kubernetes cluster**

**Author:** Igor Sfiligoi<sup>1</sup>

**Co-authors:** Edgar Fajardo Hernandez <sup>2</sup>; Dima Mishin <sup>1</sup>

<sup>1</sup> *UCSD*

<sup>2</sup> *Univ. of California San Diego (US)*

**Corresponding Authors:** edgar.mauricio.fajardo.hernandez@cern.ch, dmishin@ucsd.edu, isfiligoi@sdsc.edu

The Pacific Research Platform (PRP) is operating a Kubernetes cluster that manages over 2.5k CPU cores and 250 GPUs. Most of the resources are being used by local users interactively starting directly Kubernetes Pods.

To fully utilize the available resources, we have deployed an opportunistic HTCondor pool as a Kubernetes deployment, with worker nodes environment being fully OSG compliant. This includes both the OSG client software and CVMFS. A OSG HTCondor-CE is available for OSG users to access the resources as any other OSG site. The first user of the new site is the IceCube collaboration, which is using the available GPUs.

In this presentation we will describe the steps (and challenges) involved in creating the opportunistic OSG site in the Kubernetes cluster and the experience of running GPU jobs of the IceCube collaboration.

**Computing & Batch Systems / 20****How Fair is my Fair-Sharing? Exposing Some Hidden Behavior Through Workload Analysis**

**Authors:** Frédéric AZEVEDO<sup>1</sup>; Dalibor Klusáček<sup>2</sup>; Frederic Suter<sup>3</sup>

<sup>1</sup> *CC-IN2P3*

<sup>2</sup> *CESNET*

<sup>3</sup> *CNRS / CC-IN2P3*

**Corresponding Authors:** klusacek@cesnet.cz, frederic.azevedo@cc.in2p3.fr, frederic.suter@cc.in2p3.fr

Monitoring and analyzing how a workload is processed by a job and resource management system is at the core of the operation of data centers. It allows operators to verify that the operational

objectives are satisfied, detect any unexpected and unwanted behavior, and react accordingly to such events. However, the scale and complexity of large workloads composed of millions of jobs executed each month on several thousands of cores, often limit the depth of such analysis. This may lead to overlook some phenomena that, while they are not harmful at the global scale of the system, can be detrimental to a specific class of users.

In this talk, we illustrate such a situation by analyzing the large High Throughput Computing (HTC) workload trace coming from the Computing Center of the National Institute of Nuclear Physics and Particle Physics (CC-IN2P3) which is one of the largest academic computing centers in France. The batch scheduler implements the classical Fair-Share algorithm which ensures that all user groups are fairly provided with an amount of computing resources commensurate to their expressed needs for the year. However, the deeper we analyze this workload's scheduling, especially the waiting times of jobs, the clearer we see a certain degree of unfairness between user groups. We identify some of the root causes of this unfairness and propose a drastic reconfiguration of the quotas and scheduling queues managed by the job and resource management system. This modification aims at being more suited to the characteristics of the workload and at improving the balance across user groups in terms of waiting. We evaluate the impact of this modification through detailed simulations. The obtained results show that it still guarantees the satisfaction of the main operational objectives while significantly improving the quality of service experienced by the formerly unfavored users.

## Storage & Filesystems / 21

### Keeping Pace with Science: How a Modern Filesystem Can Accelerate Discovery

**Authors:** Andy Watson<sup>1</sup>; David Hiatt<sup>1</sup>

<sup>1</sup> *WekaIO*

**Corresponding Authors:** dave@weka.io, watson@weka.io

In November 2018, running on a mere half-rack of ordinary SuperMicro servers, WekaIO's Matrix Filesystem outperformed 40 racks of specialty hardware on Oak Ridge National Labs' Summit system, yielding the #1 ranked result for the IO-500 10-Node Challenge. How can that even be possible?

This level of performance becomes important for modern use cases whether they involve GPU-accelerated servers for artificial intelligence and deep learning or traditional CPU-based servers at massive scale. Teams of researchers and data scientists should be free to focus on their work and not lose precious time waiting for results caused by IO bottlenecks. An example use case within HEP where this technology may be most useful is the production of pre-mixing libraries in experiments like CMS. CMS uses at present a 600TB "library" to simulate overlapping proton proton collisions during its simulation campaigns. The production of this library is an IO limited workflow on any filesystem in use within the experiment today.

In this tech-talk, the architecture of the Matrix filesystem will be put under the microscope, explored and discussed. This talk will include real-world examples of data intensive workloads along with a variety of benchmark results that show the filesystem's versatility and ability to scale.

## Site Reports / 22

### PDSF Site Report

**Author:** Georg Rath<sup>1</sup>

**Co-author:** Jan Balewski<sup>2</sup>

<sup>1</sup> *Lawrence Berkeley National Laboratory*

<sup>2</sup> *Lawrence Berkeley National Lab. (US)*

**Corresponding Authors:** balewski@lbl.gov, gbrath@lbl.gov

PDSF, the Parallel Distributed Systems Facility has been in continuous operation since 1996, serving high energy physics research. The cluster is a tier-1 site for Star, a tier-2 site for Alice and a tier-3 site for Atlas.

We'll give a status report of the PDSF cluster and the migration into Cori, the primary computing resource at NERSC. We'll go into how we tried to ease the process by providing a stepping stone environment as intermediary between a commodity cluster and a supercomputer. Updates on NERSC systems will be given as well.

## Grid, Cloud and Virtualization / 23

### Public cloud for high throughput computing

**Author:** Gregory Parker<sup>1</sup>

<sup>1</sup> *Entonos*

**Corresponding Author:** gregory.parker@entonos.com

The vast breadth and configuration possibilities of the public cloud offer intriguing opportunities for loosely coupled computing tasks. One such class of tasks is simply statistical in nature requiring many independent trials over the targeted phase space in order to converge on robust, fault tolerant and optimized designs. Our single threaded target application (50-200 MB) solves a stochastic non-linear integro-differential equation relevant for read/write simulations of heat assisted magnetic recording (HAMR) for high areal density hard disk drives (HDD). Here, the phase space is multi-dimensional in physical parameters and potential recording schemes. Furthermore, for any one such point in phase space, 100s of simulations must be repeated due to the stochastic nature of the physical simulation.

In this talk, we show that a simple abstraction layer between the target application and cloud vendor provided batch systems can be easily constructed thus avoiding changes to the underlying simulation and workflow. With some planning, this abstraction layer is portable between three available cloud providers: Amazon Web Services, Microsoft Azure and Google Cloud. This abstraction layer is required to be light weight and not introduce significant overhead and was implemented as simple Bash scripts. To reduce cost, it was critical to test the application under multiple configurations (e.g. instance types and compiling options), avoid local block storage and minimize network traffic. Fleets of 100,000 concurrent simulations are easily achieved with over 99.99% of the cost just for compute (versus storage or network). By implementing a third party grid engine, 1,000,000 concurrent simulations were achieved with no modifications to the abstraction layer.

Best practices and design principles for HTC in public cloud will be discussed with emphasis on robustness, cost and horizontal scale and unique challenges encountered in this migration.

## Computing & Batch Systems / 24

### Benchmarking Working Group - Status Report

**Authors:** Michele Michelotto<sup>1</sup>; Manfred Alef<sup>2</sup>; Domenico Giordano<sup>3</sup>

<sup>1</sup> *Università e INFN, Padova (IT)*

<sup>2</sup> *Karlsruhe Institute of Technology (KIT)*

<sup>3</sup> CERN**Corresponding Authors:** domenico.giordano@cern.ch, manfred.alef@kit.edu, michele.michelotto@cern.ch

The Benchmarking Working Group has been very active in the last months. The group observed that SPEC CPU 2017 is not very different from SPEC CPU 2006. On the worker node available the two benchmark are highly correlated. Analysis with Trident shows that the hardware counters usage is rather different from the HEP applications. So the group started to investigate the usage of real applications running inside docker. The result are very promising. The current efforts are in the directions of having a suite very simple that can be distributed and runs everywhere without any knowledge of the applications, so that it can be given to a WLCG data center, a supercomputer center or a vendor for procurement procedure.

## IT Facilities & Business Continuity / 25

### Omni, N9 and the Superfacility

**Author:** Cary Whitney<sup>1</sup><sup>1</sup> LBNL**Corresponding Author:** clwhitney@lbl.gov

I will be presenting how we are using our data collection framework (Omni) to help facilitate the installation of N9 (our new system) and how this all ties together with the Superfacility concept which mentioned in the fall.

## Site Reports / 26

### Nikhef Site Report

**Author:** Dennis Van Dok<sup>None</sup>**Corresponding Author:** dennisvd@nikhef.nl

Site report with updates at Nikhef since last year.

## Storage & Filesystems / 27

### Storage services at CERN

**Author:** Enrico Bocchi<sup>1</sup><sup>1</sup> CERN**Corresponding Author:** enrico.bocchi@cern.ch

The Storage group of the CERN IT department is responsible for the development and the operation of petabyte-scale services needed to accommodate the diverse requirements for storing physics data generated by LHC and non-LHC experiments as well as supporting users of the laboratory in their day-by-day activities.



This contribution presents the current operational status of the main storage services at CERN, summarizes our experience in operating largely distributed systems and highlights the ongoing efforts for the evolution of the storage infrastructure.

It reports about EOS, the high-performance distributed filesystem developed at CERN designed to store all the physics data and to operate at the high rates demanded by experiments data taking. EOS is also used as the storage backend for CERNBox, the cloud storage synchronization and sharing service for users' personal files. CERNBox provides uniform access to storage from all modern devices and represents the data hub for integration with various applications ranging from office suites (Microsoft Office 365, OnlyOffice, Draw.io) to specialized tools for data analysis (SWAN).

Besides storage for physics data and personal files, the Storage group runs multiple large Ceph clusters to provide the storage backbones for the OpenStack infrastructure and the HPC facility, and to offer an S3 service and a CephFS/Manila shares for other internal IT services. Also, the Storage group operates the release managers, replica servers and caches of CVMFS (a fundamental WLCG service used for software distribution) in collaboration with the SoFTware Development for Experiments (CERN EP-SFT) department.

## Networking & Security / 28

### CERN DNS and DHCP service improvement plans

**Author:** Quentin Barrand<sup>1</sup>

**Co-authors:** David Gutierrez Rueda <sup>1</sup>; Veronique Lefebure <sup>1</sup>

<sup>1</sup> CERN

**Corresponding Authors:** [quentin.barrand@cern.ch](mailto:quentin.barrand@cern.ch), [veronique.lefebure@cern.ch](mailto:veronique.lefebure@cern.ch), [david.gutierrez@cern.ch](mailto:david.gutierrez@cern.ch)

The configuration of the CERN IT central DNS servers, based on ISC BIND, is generated automatically from scratch every 10 minutes using a software developed at CERN several years ago. This in-house set of Perl scripts has evolved and is reaching its limits in terms of maintainability and architecture. CERN is in the process of reimplementing the software with a modern language and is taking the opportunity to redefine the DNS service architecture by introducing a redundant solution for the master DNS. Meanwhile, Anycast is being evaluated in order to increase the DNS service robustness and scalability. Finally, CERN is considering the possibility of moving from a static to a dynamic zone for the cern.ch domain to allow immediate commissioning while controlling the update process.

Concerning the DHCP services, ISC DHCP has been the software of choice to support dynamic host configuration for almost 20 years. However system provisioning has massively scaled in the last years and DHCP software shortcomings have lead ISC to develop Kea. CERN intends to modernize the service replacing ISC DHCP with Kea, which will allow the implementation of a highly available and geographically dispersed DHCP service, as well as a fast provisioning so that changes in the network database are immediately propagated to the DHCP servers.

## Welcome to SDSC and UCSD / 29

### Welcome to SDSC and UCSD

## Storage & Filesystems / 30

### Using the dynafed data federation as site storage element

**Authors:** Marcus Ebert<sup>1</sup>; Frank Berghaus<sup>2</sup>; Kevin Casteels<sup>1</sup>; Colson Driemel<sup>1</sup>; Fernando Fernandez Galindo<sup>3</sup>; Colin Roy Leavett-Brown<sup>2</sup>; Michael Paterson<sup>4</sup>; Rolf Seuster<sup>2</sup>; Randy Sobie<sup>2</sup>; Reda Tafirout<sup>3</sup>

<sup>1</sup> *University of Victoria*

<sup>2</sup> *University of Victoria (CA)*

<sup>3</sup> *TRIUMF (CA)*

<sup>4</sup> *U*

**Corresponding Authors:** tafirout@triumf.ca, frank.berghaus@cern.ch, mhp@uvic.ca, rsobie@uvic.ca, casteels@uvic.ca, colin.roy.leavett-brown@cern.ch, mebert@uvic.ca, rolf.seuster@cern.ch, fernandezgalindo@triumf.ca, colsond@uvic.ca

We describe our experience and use of the Dynafed data federator with cloud and traditional Grid computing resources as an substitute for a traditional Grid SE.

This is an update of the report given at the Fall HEPiX meeting of 2017 where we introduced our use case for such federation and described our initial experience with it.

We used Dynafed in production for Belle-II since late 2017 and also in testing mode for Atlas. We will report on changes we made since then to our setup and also report on changes in Dynafed itself that makes it more suitable as a site SE. We will also report on a new monitoring system we developed for such data federation and report also on a way to use such data federation by anyone who uses distributed compute and storage but with the need to read/write from a local file system.

## Site Reports / 31

### University of Wisconsin-Madison CMS T2 site report

**Authors:** Ajit Kumar Mohapatra<sup>1</sup>; Sridhara Dasu<sup>1</sup>; Daniel Charles Bradley<sup>1</sup>; Carl Vuosalo<sup>1</sup>; Chad W Seys<sup>2</sup>

<sup>1</sup> *University of Wisconsin Madison (US)*

<sup>2</sup> *University of Wisconsin-Madison*

**Corresponding Authors:** cwseys@physics.wisc.edu, carl.vuosalo@cern.ch, daniel.bradley@cern.ch, dasu@hep.wisc.edu, ajit.kumar.mohapatra@cern.ch

As a major WLCG/OSG T2 site, the University of Wisconsin-Madison CMS T2 has consistently been delivering highly reliable and productive services towards large scale CMS MC production/processing, data storage, and physics analysis for last 13 years. The site utilizes high throughput computing (HTCondor), highly available storage system (Hadoop), scalable distributed software systems (CVMFS), and provides efficient data access using xrootd/AAA. The site fully supports IPv6 networking, and is a member of the LHCONE community with 100Gb WAN connectivity. An update on the activities and developments at the T2 facility over the past 1.5 years (since the KEK meeting) will be presented.

## IT Facilities & Business Continuity / 33

### HEPiX Technology Watch working group

**Author:** Helge Meinhard<sup>1</sup>

<sup>1</sup> *CERN*

**Corresponding Author:** helge.meinhard@cern.ch

A short report on what has happened, how we have organised ourselves, how we intend to present results etc. Note that the findings themselves will be discussed in other contributions - this is about how the group works.

**Site Reports / 37****Update of Canadian T1 / T2**

**Author:** Rolf Seuster<sup>1</sup>

<sup>1</sup> *University of Victoria (CA)*

**Corresponding Author:** rolf.seuster@cern.ch

I will present recent developments of the Canadian T1 and T2.

**Grid, Cloud and Virtualization / 38****CloudScheduler version2**

**Authors:** Rolf Seuster<sup>1</sup>; Colin Roy Leavett-Brown<sup>1</sup>; Michael Paterson<sup>2</sup>; Kevin Casteels<sup>3</sup>; Colson Driemel<sup>3</sup>; Marcus Ebert<sup>3</sup>; Danika MacDonell<sup>1</sup>; Randy Sobie<sup>1</sup>

<sup>1</sup> *University of Victoria (CA)*

<sup>2</sup> *U*

<sup>3</sup> *University of Victoria*

**Corresponding Authors:** colin.roy.leavett-brown@cern.ch, rsobie@uvic.ca, mebert@uvic.ca, danikam1@uvic.ca, casteels@uvic.ca, mhp@uvic.ca, rolf.seuster@cern.ch, colsond@uvic.ca

I present the recent developments for our cloudschedduler, which we use to run HEP workloads on various clouds in North America and Europe. We are working on a complete re-write utilizing modern software technologies and practices.

**Computing & Batch Systems / 39****Jupyter at SDCC**

**Author:** William Strecker-Kellogg<sup>1</sup>

<sup>1</sup> *Brookhaven National Lab*

**Corresponding Author:** willsk@bnl.gov

...Placeholder...

**Site Reports / 40****RAL Site Report**

**Author:** Martin Bly<sup>1</sup>

<sup>1</sup> *STFC-RAL*

**Corresponding Author:** martin.bly@stfc.ac.uk

An update on activities at STFC-RAL.

## Networking & Security / 41

### The difference in network equipment in the 25/100/400G era and how to test/break them

**Author:** Tristan Suerink<sup>1</sup>

<sup>1</sup> *Nikhef National institute for subatomic physics (NL)*

**Corresponding Author:** t.suerink@cern.ch

The network market has changed a lot compared with a decade ago. Every hardware vendor sells their own switches and routers. Most of the switches and routers are based on the same merchant silicon that is available on the market.

Therefore the amount of real choices is limited because what is inside is the same for most of them. This talk will tell about the differences that are still there and what are the risks for choosing certain solutions.

What will vendors really allow you to do in their “Open Networking” strategy?

Why is knowing how many packets per second more important than how much bandwidth a network device can process?

How do you test this and what type of effects can you expect when reaching the limits of the equipment?

Why aren't commercial network testers the best way of testing network equipment?

Is building your own network test machine expensive?

## Storage & Filesystems / 42

### OSiRIS: Open Storage Research Infrastructure

**Author:** Benjeman Jay Meekhof<sup>1</sup>

**Co-author:** Shawn Mc Kee<sup>1</sup>

<sup>1</sup> *University of Michigan (US)*

**Corresponding Authors:** bmeekhof@umich.edu, shawn.mckee@cern.ch

OSiRIS is a pilot project funded by the NSF to evaluate a software-defined storage infrastructure for our primary Michigan research universities and beyond. In the HEP world OSiRIS is involved with ATLAS as a provider of Event Service storage via the S3 protocol as well as experimenting with dCache backend storage for AGLT2. We are also in the very early stages of working with IceCube and the nationwide Open Storage Network. Our talk will cover current status on these projects and the latest details of how we use Ceph, HAProxy, NFSv4, LDAP, CManage, Puppet and other tools to provision, manage, and monitor storage services to federated users.

## Networking & Security / 43

## IPv6 & WLCG - an update from the HEPiX IPv6 Working Group

**Authors:** David Kelsey<sup>1</sup>; Andrea Sciabà<sup>2</sup>

<sup>1</sup> *Science and Technology Facilities Council STFC (GB)*

<sup>2</sup> *CERN*

**Corresponding Authors:** andrea.sciaba@cern.ch, david.kelsey@stfc.ac.uk

The transition of WLCG storage services to dual-stack IPv4/IPv6 is progressing well, aimed at enabling the use of IPv6-only CPU resources as agreed by the WLCG Management Board and presented by us at previous HEPiX meetings.

The working group, driven by the requirements of the LHC VOs to be able to use IPv6-only opportunistic resources, continues to encourage wider deployment of dual-stack services and has been monitoring the transition. During recent months we have also started to investigate in more detail the reasons for various edge cases where the fraction of data transferred over IPv6 is lower than expected.

This talk will present the current status of the transition to IPv6 together with some of the common reasons for sites that have not yet been able to move to dual-stack operations. Some issues related to unexpected monitoring results for IPv6 versus IPv4 will also be discussed.

### Storage & Filesystems / 44

## Storage management in a large scale at BNL

**Author:** Robert Hancock<sup>1</sup>

<sup>1</sup> *Brookhaven National Laboratory*

**Corresponding Author:** nps2tls@gmail.com

Brookhaven National Laboratory stores and processes large amounts of data from the following: PHENIX, STAR, ATLAS, Belle II, Simons, as well as smaller local projects. This data is stored long term in tape libraries but one working data is stored in disk arrays. Hardware raid devices from companies such as Hitachi Ventara are very convenient and require minimal administrative intervention. However, they are very expensive relative the alternatives. BNL is making a move toward JBOD (Just a Bunch of Disk) arrays with Linux based software raid. The performance is comparable and sometimes better than the hardware cousins but the cost is less than half. However, the construction and administration is more complex. This requires more hours of skilled manpower from staff to install and maintain. I am developing software at BNL to automate these processes to the level of hardware raid in order to reduce this burden while allowing cost savings.

### Grid, Cloud and Virtualization / 45

## Developing for a Services Layer At The Edge (SLATE)

**Authors:** Shawn Mc Kee<sup>1</sup>; Robert William Gardner Jr<sup>2</sup>; Joe Breen<sup>3</sup>; Ben Kulbertis<sup>3</sup>

<sup>1</sup> *University of Michigan (US)*

<sup>2</sup> *University of Chicago (US)*

<sup>3</sup> *University of Utah*

**Corresponding Authors:** shawn.mckee@cern.ch, ben.kulbertis@utah.edu, joe.breen@utah.edu, robert.w.gardner@cern.ch

Modern software development workflow patterns often involve the use of a developer's local machine as the first platform for testing code. SLATE mimics this paradigm with an implementation of a light-weight version, called MiniSLATE, that runs completely contained on the developer local machine (laptop, virtual machine, or another physical server). MiniSLATE resolves many development environment issues by providing an isolated and local configuration for the developer. Application developers are able to download MiniSLATE which provides a fully orchestrated set of containers on top of a production SLATE platform, complete with central information service, API server, and a local Kubernetes cluster. This approach mitigates the overhead of a hypervisor but still provides the requisite isolated environment. They are able to create the environment, iterate, destroy it, and repeat at will. A local MiniSLATE environment also allows the developer to explore the packaging of the edge service within a constrained security context in order to validate its full functionality within limited permissions. As a result, developers are able to test the functionality of their application with the complete complement of SLATE components local to their development environment without the overhead of building a cluster or virtual machine, registering a cluster, interacting with the production SLATE platform, etc.

## Basic IT Services / 46

### Config Management and Deployment Setup at KIT

**Author:** Andreas Petzold<sup>1</sup>

<sup>1</sup> *KIT - Karlsruhe Institute of Technology (DE)*

**Corresponding Author:** andreas.petzold@cern.ch

For several years, the GridKa Tier-1 center, the Large Scale Data Facility and other infrastructures at KIT have been using Puppet and Foreman for configuration management and machine deployment. We will present our experiences, the workflows that are used and our current efforts to establish a completely integrated system for all our infrastructures based on Katello.

## Storage & Filesystems / 47

### Virtualization for Online Storage Clusters

**Author:** Jan Erik Sundermann<sup>1</sup>

<sup>1</sup> *Karlsruhe Institute of Technology (KIT)*

**Corresponding Author:** jan.sundermann@kit.edu

The computing center GridKa is serving the ALICE, ATLAS, CMS and LHCb experiments as Tier-1 center with compute and storage resources. It is operated by the Steinbuch Centre for Computing at Karlsruhe Institute of Technology in Germany. In its current stage of expansion GridKa offers the HEP experiments a capacity of 35 Petabytes of online storage. The storage system is based on Spectrum Scale as software-defined-storage layer. Its storage servers are inter-connected via two redundant infiniband fabrics and have ethernet uplinks to the GridKa backbone network. In this presentation we discuss the use of virtualization technologies in the context of the described storage system, including hardware virtualization of the infiniband and ethernet interfaces.

**Networking & Security / 48****Network Functions Virtualisation Working Group Update****Authors:** Marian Babik<sup>1</sup>; Shawn Mc Kee<sup>2</sup><sup>1</sup> *CERN*<sup>2</sup> *University of Michigan (US)***Corresponding Authors:** shawn.mckee@cern.ch, marian.babik@cern.ch

High Energy Physics (HEP) experiments have greatly benefited from a strong relationship with Research and Education (R&E) network providers and thanks to the projects such as LHCOPN/LHCONE and REN contributions, have enjoyed significant capacities and high performance networks for some time. RENs have been able to continually expand their capacities to over-provision the networks relative to the experiments needs and were thus able to cope with the recent rapid growth of the traffic between sites, both in terms of achievable peak transfer rates as well as in total amount of data transferred. For some HEP experiments this has led to designs that favour remote data access where network is considered an appliance with almost infinite capacity. There are reasons to believe that the network situation will change due to both technological and non-technological reasons starting already in the next few years. Various non-technological factors that are in play are for example anticipated growth of the non-HEP network usage with other large data volume sciences coming online; introduction of the cloud and commercial networking and their respective impact on usage policies and securities as well as technological limitations of the optical interfaces and switching equipment.

As the scale and complexity of the current HEP network grows rapidly, new technologies and platforms are being introduced that greatly extend the capabilities of today's networks. With many of these technologies becoming available, it's important to understand how we can design, test and develop systems that could enter existing production workflows while at the same time changing something as fundamental as the network that all sites and experiments rely upon. In this talk we'll give an update on the working group's recent activities, updates from sites and R&E network providers as well as plans for the near-term future.

**Networking & Security / 49****WLCG/OSG Network Activities, Status and Plans****Authors:** Shawn Mc Kee<sup>1</sup>; Marian Babik<sup>2</sup>; Brian Paul Bockelman<sup>3</sup>; Robert William Gardner Jr<sup>4</sup><sup>1</sup> *University of Michigan (US)*<sup>2</sup> *CERN*<sup>3</sup> *University of Nebraska Lincoln (US)*<sup>4</sup> *University of Chicago (US)***Corresponding Authors:** robert.w.gardner@cern.ch, brian.bockelman@cern.ch, marian.babik@cern.ch, shawn.mckee@cern.ch

WLCG relies on the network as a critical part of its infrastructure and therefore needs to guarantee effective network usage and prompt detection and resolution of any network issues, including connection failures, congestion and traffic routing. The OSG Networking Area is a partner of the WLCG effort and is focused on being the primary source of networking information for its partners and constituents. We will report on the changes and updates that have occurred since the last HEPiX meeting.

The primary areas to cover include the status of and plans for the WLCG/OSG perfSONAR infrastructure, the WLCG Throughput Working Group and the activities in the IRIS-HEP and SAND projects.

**Site Reports / 50****AGLT2 Site Report**

**Authors:** Wenjing Wu<sup>1</sup>; Shawn Mc Kee<sup>2</sup>; Philippe Laurens<sup>3</sup>

<sup>1</sup> *University of Michigan*

<sup>2</sup> *University of Michigan (US)*

<sup>3</sup> *Michigan State University (US)*

**Corresponding Authors:** philippe.laurens@cern.ch, wuwj@umich.edu, shawn.mckee@cern.ch

We will present an update on AGLT2, focusing on the changes since the Fall 2018 report. The primary topics to cover include the update on VMware, update of dCache, status of new purchased hardware, encountered problems and solutions on improving the CPU utilization of our HT-Condor system.

**IT Facilities & Business Continuity / 51****RACF/SDCC Datacenter Transformation within the Scope of BNL CFR Project and Beyond**

**Author:** Alexandr Zaytsev<sup>1</sup>

<sup>1</sup> *Brookhaven National Laboratory (US)*

**Corresponding Author:** alezayt@bnl.gov

The BNL Computing Facility Revitalization (CFR) project aimed at repurposing the former National Synchrotron Light Source (NSLS-I) building (B725) located on BNL site as the new datacenter for BNL Computational Science Initiative (CSI) and RACF/SDCC Facility in particular. The CFR project is currently wrapping up the design phase and expected to enter the construction phase in the first half of 2019. The new B725 data center is to become available in early 2021 for ATLAS compute, disk storage and tape storage equipment, and later during the year of 2021 - for all other collaborations supported by the RACF/SDCC Facility, including but not limited to: STAR and PHENIX experiments at RHIC collider at BNL, Belle II Experiment at KEK (Japan), and BNL CSI HPC clusters. Migration of the majority of IT payload from B515 based datacenter to the B725 datacenter is expected to begin even earlier, as the central networking systems and first BNL ATLAS Tier-1 Site tape robot are to be deployed in B725 starting from early FY21, and expected to continue throughout the period of FY21-23, leaving the B515 datacenter physically reduced down to a subset of areas it is currently occupying, and drastically reducing its power profile. In this talk I am going to highlight the main design features of the new RACF/SDCC datacenter, summarize the preparation activities already underway in our existing datacenter since FY18 needed to ensure a smooth transition B515 and B725 datacenters inter-operation period in FY21, discuss the planned sequence of equipment migration between these two datacenters in FY21 and gradual equipment replacement in FY21-24, and also show the expected state of occupancy and infrastructure utilization for both datacenters in FY25.

**End-User IT Services & Operating Systems / 52****How to make your Cluster look like a Supercomputer (for Fun and Profit)**

**Author:** Georg Rath<sup>1</sup>

**Co-author:** Ershaad Basheer<sup>2</sup>



<sup>1</sup> *Lawrence Berkeley National Laboratory*

<sup>2</sup> *LBNL/NERSC*

**Corresponding Authors:** ebasheer@lbl.gov, gbrath@lbl.gov

During the last two years, the computational systems group at NERSC, in partnership with Cray, has been developing SMWFlow, a tool that makes managing system state as simple as switching branches in git. This solution is the cornerstone of collaborative systems management at NERSC and enables code-review, automated testing and reproducibility.

Besides supercomputers, NERSC hosts Mendel, a commodity meta-system containing multiple clusters, among them PDSF, used by the HEP community, and Denovo, used by the Joint Genome Institute, which uses a custom management stack built on top of xCAT and cfengine.

To merge efforts, provide a consistent user experience, and to leverage the work done on SMWFlow, we will talk about how we adapted the Cray imaging and provisioning system to work on a system on an architecture like Mendel and therefore reap the benefits of a modern systems management approach.

## Site Reports / 53

### Tokyo Tier-2 Site Report

**Author:** Tomoe Kishimoto<sup>1</sup>

<sup>1</sup> *University of Tokyo (JP)*

**Corresponding Author:** tomoe.kishimoto@cern.ch

The Tokyo Tier-2 center, which is located in the International Center for Elementary Particle Physics (ICEPP) at the University of Tokyo, is providing computing resources for the ATLAS experiment in the WLCG. Almost all hardware devices of the center are supplied by a lease, and are upgraded in every three years. This hardware upgrade was performed in December 2018. In this presentation, experiences of the system upgrade will be reported. The configuration of the new system will also be shown.

## Networking & Security / 54

### BNL activities on federated access and Single Sign-On

**Author:** Tejas Rao<sup>1</sup>

<sup>1</sup> *Brookhaven National Laboratory*

**Corresponding Author:** raot@bnl.gov

Various High energy and nuclear physics experiments already benefit from using the different components of Federated architecture to access storage and infrastructure services. BNL moved to Identity management (Redhat IPA) in late 2018 which will serve as the foundation to move to Federated authentication and authorization. IPA provides central authentication via Kerberos or LDAP, simplifies administration,

has a rich CLI and a web based user interface. This presentation describes how federated authn/authz will be enabled in the near future at the level of individual applications like Globus online, Invenio, BNLbox, Indico, Web services and Jupyter.

## Site Reports / 55

## BEIJING Site Report

**Author:** Qiulan Huang<sup>1</sup>

<sup>1</sup> *Institute of High Energy Physics, Chinese Academy of Science*

**Corresponding Author:** huangql@ihep.ac.cn

News and updates from IHEP since the last HEPiX Workshop. In this talk we would like to present the status of IHEP site including computing farm, HPC, IHEPcloud, Grid, data storage ,network and so on.

### Storage & Filesystems / 57

## DPM DOME

**Author:** Petr Vokac<sup>1</sup>

**Co-authors:** Martin Adam <sup>2</sup>; Alexandr Mikula <sup>2</sup>; Jiri Chudoba <sup>2</sup>

<sup>1</sup> *Czech Technical University (CZ)*

<sup>2</sup> *Acad. of Sciences of the Czech Rep. (CZ)*

**Corresponding Authors:** petr.vokac@cern.ch, jiri.chudoba@cern.ch, martin.adam@cern.ch, alexandr.mikula@cern.ch

DPM (Disk Pool Manager) is mutli-protocol distributed storage system that can be easily used within grid environment and it is still popular for medium size sites. Currently DPM can be configured to run in legacy or DOME mode, but official support for the legacy flavour ends this summer and sites using DPM storage should think about their upgrade strategy or coordinate with WLCG DPM Upgrade task force.

We are going to present our almost a year long experience with DPM running in DOME mode on our production storage hosting several petabytes of data for different VOs. Our experience can help others to avoid common problems and also choose right protocols to get best performance from DPM storage. DOME provides support for SRM-less site configuration, but SRM can be still used if necessary and we'll show advantages and/or disadvantages that comes from such configuration. New DPM features are developed only for DOME mode. We would like to summarize which improvements are available in DOME and these features will never be available for legacy flavour. Running DOME brings greatly improved support for non-GridFTP protocols like full support for transfer checksums, storage resource reporting and most importantly third-party-copy (TPC). We are going to describe DPM TPC support including various credentials delegation mechanisms for XRootD and WebDAV protocols and interoperability with other storage implementations.

### Grid, Cloud and Virtualization / 58

## Addressing the Challenges of Executing Massive Computational Clusters in the Cloud

**Authors:** Alexander Herzog<sup>1</sup>; Amy Apon<sup>1</sup>; Boyd Wilson<sup>2</sup>; Brandon Posey<sup>3</sup>; Christopher Gropp<sup>1</sup>

<sup>1</sup> *Clemson University*

<sup>2</sup> *Omnibond*

<sup>3</sup> *BMW*

**Corresponding Authors:** boydw@omnibond.com, brandon@omnibond.com, aherzog@clemson.edu, cgropp@clemson.edu, aapon@clemson.edu

This talk will discuss how we worked with Dr. Amy Apon, Brandon Posey, AWS and the Clemson DICE lab team dynamically provisioned a large scale computational cluster of more than one million cores utilizing Amazon Web Services (AWS). We discuss the trade-offs, challenges, and solutions associated with creating such a large scale cluster with commercial cloud resources. We utilize our large scale cluster to study a parameter sweep workflow composed of message-passing parallel topic modeling jobs on multiple datasets.

At peak, we achieve a simultaneous core count of 1,119,196 vCPUs across nearly 50,000 instances, and are able to execute almost half a million jobs within two hours utilizing AWS Spot Instances in a single AWS region.

Additionally we will discuss a follow on project that the DICE Lab is currently working on in the Google Cloud Platform (GCP) that will enable a Computer Vision analytics system to concurrently processes hundreds of thousands of hours of highway traffic video providing statistics on congestions, vehicle trajectories and neural net pre-annotation. We will discuss how this project will differ from the previous one and how additional boundaries are being pushed.

Relevant Papers:

<https://ieeexplore.ieee.org/abstract/document/8411029>

[https://tigerprints.clemson.edu/computing\\_pubs/38/](https://tigerprints.clemson.edu/computing_pubs/38/)

## Storage & Filesystems / 59

### Integrating Hadoop Distributed File System to Logistical Storage

**Author:** Shunxing Bao<sup>1</sup>

**Co-authors:** Alan Tackett<sup>1</sup>; Andrew Malone Melo<sup>2</sup>

<sup>1</sup> *Vanderbilt University*

<sup>2</sup> *Vanderbilt University (US)*

**Corresponding Authors:** [andrew.malone.melo@cern.ch](mailto:andrew.malone.melo@cern.ch), [alan.tackett@accre.vanderbilt.edu](mailto:alan.tackett@accre.vanderbilt.edu), [onealbao@gmail.com](mailto:onealbao@gmail.com)

Logistical Storage (LStore) provides a flexible logistical networking storage framework for distributed and scalable access to data in both an HPC and WAN environment. LStore uses commodity hard drives to provide unlimited storage with user controllable fault tolerance and reliability. In this talk, we will briefly discuss LStore's features and discuss the newly developed native LStore plugin for the Apache Hadoop ecosystem. The Hadoop Distributed File System (HDFS) will directly access LStore using this plugin allowing users to create Hadoop clusters on the fly in an HPC environment. The primary benefit of the plugin is that it avoids the need for data redundancy across a traditional Hadoop and HPC cluster. Moreover, the on the fly Hadoop clusters created in the HPC environment can be scaled as needed and tune the hardware requirements to the analysis - large memory needs, GPU, etc.

We will show several empirical results using the plugin in both a traditional HPC environment and utilizing a high-latency WAN connection. The proposed plugin is compared with two current LStore interfaces: LStore command line interface and LStore FUSE mounted client interface.

## Computing & Batch Systems / 60

### Deep learning in a container, experience and best practices

**Authors:** Bertrand Rigaud<sup>1</sup>; Frederic Suter<sup>None</sup>

<sup>1</sup> CC-IN2P3 / CNRS

**Corresponding Authors:** frederic.suter@cc.in2p3.fr, bertrand.rigaud@cc.in2p3.fr

Deep Learning techniques are gaining interest in the High Energy Physics, following a new and efficient approach to solve different problems. These techniques leverage the specific features of GPU accelerators and rely on a set of software packages allowing users to compute on GPUs and program Deep Learning algorithms. However, the rapid pace at which both the hardware and the low and high level libraries are evolving poses several operational issues to computing centers such as the IN2P3 Computing Center (CC-IN2P3 – <http://cc.in2p3.fr>).

In this talk we present how we addressed these operational challenges thanks to the use of container technologies. We show that the flexibility offered by containers comes with no overhead while allowing users to benefit of the better performance of compiled from sources versions of popular deep learning frameworks. Finally, we detail the best practices proposed to the users of the CC-IN2P3 to prepare and submit their deep learning oriented jobs on the available GPU resources.

## Storage & Filesystems / 61

### Developments in disk and tape storage at the RAL Tier 1

**Author:** Rob Appleyard<sup>1</sup>

**Co-authors:** Tom Byrne <sup>1</sup>; Aidan McComb <sup>1</sup>; George Patargias <sup>1</sup>

<sup>1</sup> STFC

**Corresponding Authors:** rob.appleyard@stfc.ac.uk, george.patargias@stfc.ac.uk, aidan.mccomb@stfc.ac.uk, tom.byrne@stfc.ac.uk

RAL's Ceph-based Echo storage system is now the primary disk storage system running at the Tier 1, replacing a legacy CASTOR system that will be retained for tape. This talk will give an update on Echo's recent development, in particular the adaptations needed to support the ALICE experiment and the challenges of scaling an erasure-coded Ceph cluster past the 30PB mark. These include the smoothing of data distribution, managing disk errors, and dealing with a very full cluster.

In addition, I will discuss the completed project to remodel RAL's CASTOR service from a combined disk and tape endpoint to a low-maintenance system only providing access to tape.

## Site Reports / 62

### DESY site report

**Author:** Peter van der Reest<sup>1</sup>

<sup>1</sup> DESY

**Corresponding Author:** peter.van.der.reest@desy.de

Brief overview of activities at DESY since the last site report

## Basic IT Services / 63

### Token Renewal Service (TRS) at SLAC

**Author:** Andrew May<sup>1</sup>

<sup>1</sup> *SLAC National ACCELERATOR LABORATORY*

**Corresponding Author:** amay@stanford.edu

The token renewal service (TRS) has been used at SLAC National Accelerator Laboratory since the late 1990s. In 2018 it was found to be lacking in some critical areas (encryption types used and other basic mechanism would no longer be available for post Red Hat 6 systems.)

1-to-1 replacement areas:

Running Batch Jobs:

Our local solution to batch jobs (LSF): The need for TRS was already resolved with the movement to a new mechanism used for re-authorizing in IBM's LSF version ???.

Our remote solution for batch (like) jobs: OSGrid and CVMFS services have no requirement for a TRS type solution. In the area of using remote (full or partial-stack) computing resources from the industry titans, (e.g. Azure, AWS and Google), SLAC-OCIO does not actively use those resources as of March 2019.

Users that run TRScron jobs: a distributed cron service which leverages the token renewal service, the need remained.

Questions:

What adaptations should be pondered in the face of future computing workloads?

How does a new vision of TRS compare to past and present mechanisms used to provide renewable secure access to a distributed service. What are we missing? Comments, questions, and concerns?

<https://confluence.slac.stanford.edu/display/TRS/TRS+Home>

## Site Reports / 64

### University of Nebraska CMS Tier2 Site Report

**Author:** Garhan Attebury<sup>1</sup>

<sup>1</sup> *University of Nebraska Lincoln (US)*

**Corresponding Author:** garhan.attebury@cern.ch

Updates on the activities at T2\_US\_Nebraska over the past year. Topics will cover the site configuration and tools we use, troubles we face in daily operation, and contemplation of what the future might hold for sites like ours.

## Site Reports / 65

### UW CENPA site report

**Author:** Duncan Prindle<sup>1</sup>

<sup>1</sup> *University of Washington*

**Corresponding Author:** prindle@npl.washington.edu

At CENPA at the University of Washington we have a heterogeneous rocks 7 cluster of about 135 nodes containing 1250 cores. We will present the current status and issues.

## Site Reports / 66

### JLab Site Report

**Author:** Sandra Philpott<sup>None</sup>

**Corresponding Author:** philpott@jlab.org

Updates from JLab since the Autumn 2018 meeting at PIC in Barcelona.

## Miscellaneous / 67

### Fermilab Site Report

**Author:** Bo Jayatilaka<sup>1</sup>

<sup>1</sup> *Fermi National Accelerator Lab. (US)*

**Corresponding Author:** bo.jayatilaka@cern.ch

Status, ongoing activities, and future directions at Fermilab.

## Grid, Cloud and Virtualization / 68

### Changes to OSG and how they affect US WLCG sites

**Author:** Jeffrey Michael Dost<sup>1</sup>

<sup>1</sup> *Univ. of California San Diego (US)*

**Corresponding Author:** jdost@ucsd.edu

In the spring of 2018, central operations services were migrated out of the Grid Operations Center of Indiana into other participating Open Science Grid institutions. This talk summarizes how the migration has affected the services provided by the OSG, and gives a summary of how central OSG services interface with US WLCG sites.

## Grid, Cloud and Virtualization / 69

### The glideinWMS system: recent developments

**Author:** Marco Mascheroni<sup>1</sup>

<sup>1</sup> *Univ. of California San Diego (US)*

**Corresponding Author:** marco.mascheroni@cern.ch

GlideinWMS is a workload management and provisioning system that lets you share computing resources distributed over independent sites. A dynamically sized pool of resources is created by GlideinWMS pilot Factories, based on the requests made by GlideinWMS Frontends. More than 400 computing elements are currently serving more than 10 virtual organizations through glideinWMS. This contribution will give an overview of the glideinWMS setup, and will present the recent developments in the project, including the addition of the singularity support, and the improvements to minimize resource wastages. Future enhancements for automatizing the generation of factory configurations will also be outlined.

## IT Facilities & Business Continuity / 70

### Technolog Watch WG Reports

**Authors:** Andrea Chierici<sup>1</sup>; Martin Gasthuber<sup>2</sup>; Michele Michelotto<sup>3</sup>; Rolf Seuster<sup>4</sup>

<sup>1</sup> *INFN-CNAF*

<sup>2</sup> *DESY*

<sup>3</sup> *Università e INFN, Padova (IT)*

<sup>4</sup> *University of Victoria (CA)*

**Corresponding Authors:** michele.michelotto@cern.ch, chierici@cnaif.infn.it, martin.gasthuber@cern.ch, rolf.seuster@cern.ch

We will report on the findings of the technology watch working group concerning CPUs, storage, networks and related fields