



home.cern

HEP*iX* Fall 2016 Summary

<https://indico.cern.ch/event/531810>

Lawrence Berkeley National Laboratory

Jérôme Belleman · Marco Guerri · Adam Sosnowski

Outline

Jérôme

Fall 2016 Meeting & General HEPiX News

Site Reports

Computing & Batch Services

Basic IT Services

Marco

Storage & Filesystems

Grid, Cloud & Virtualisation

IT Facilities & Business Continuity

Adam

Security & Networking

End-User IT Services & Operating Systems

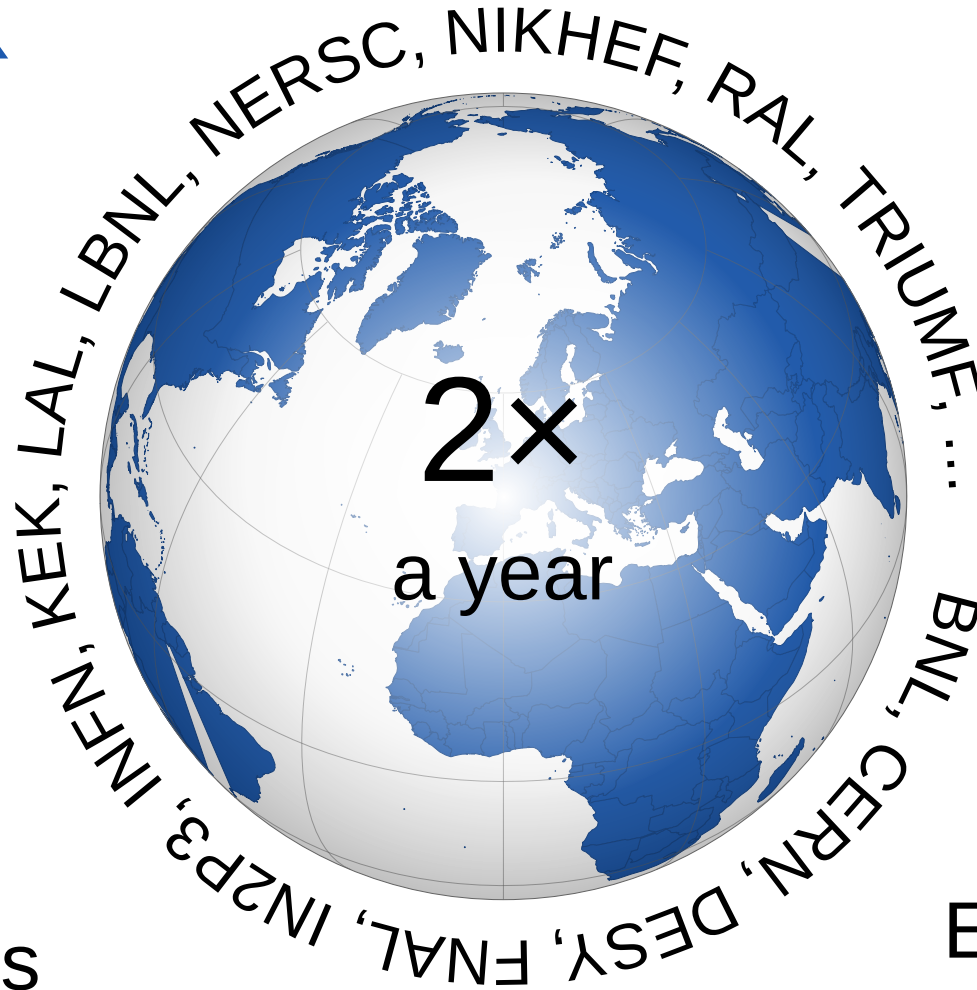
Miscellaneous

HEPiX

Future plans

Working groups

Challenges



Recent work

Status reports

Experiences

www.hepix.org

Lawrence Berkeley National Lab



Lawrence Berkeley National Lab



Lawrence Berkeley National Lab

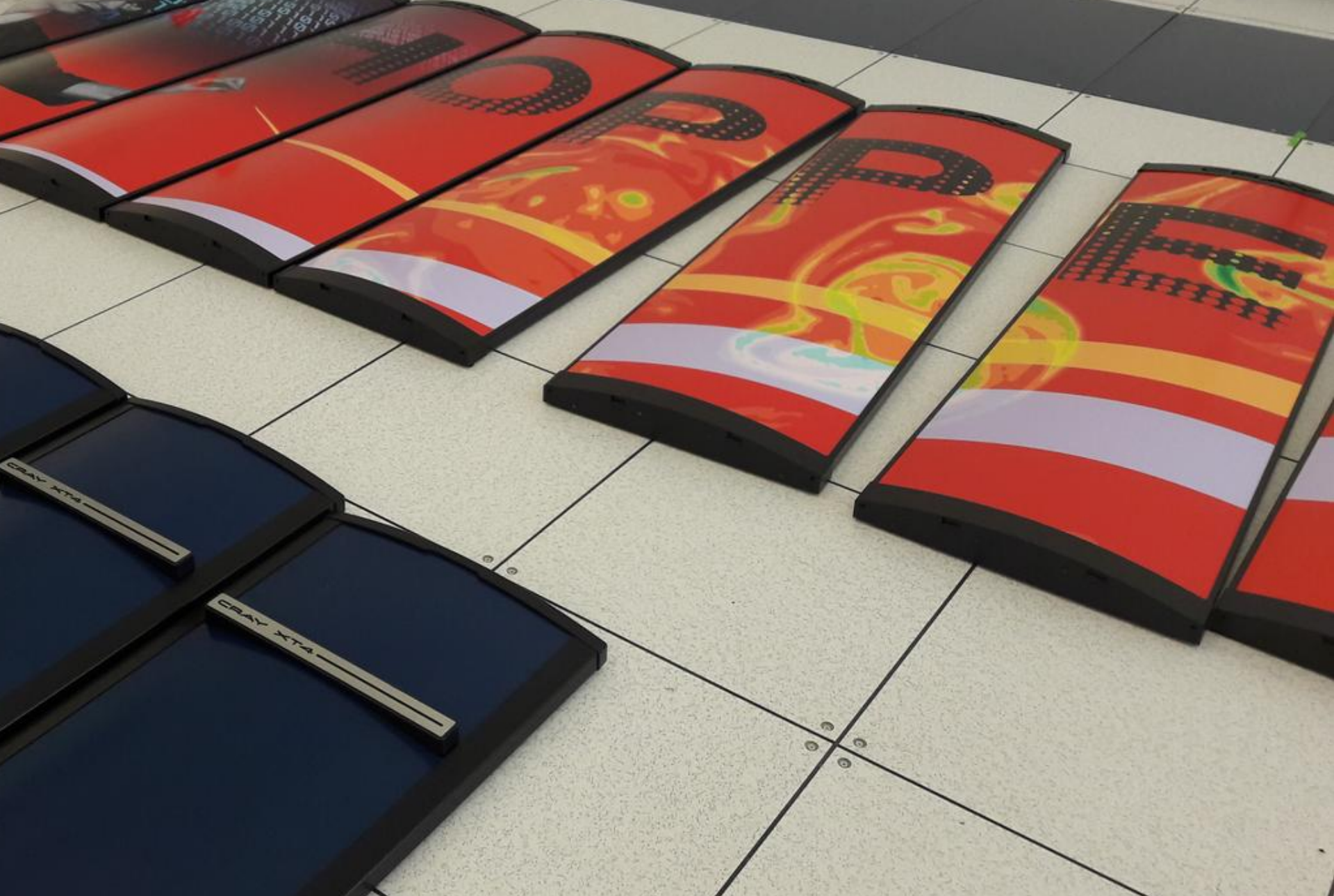


Lawrence Berkeley National Lab









Site Reports (I)

The rise of SLURM

Flexible

Scalable

Integrates with HTCondor

HTCondor

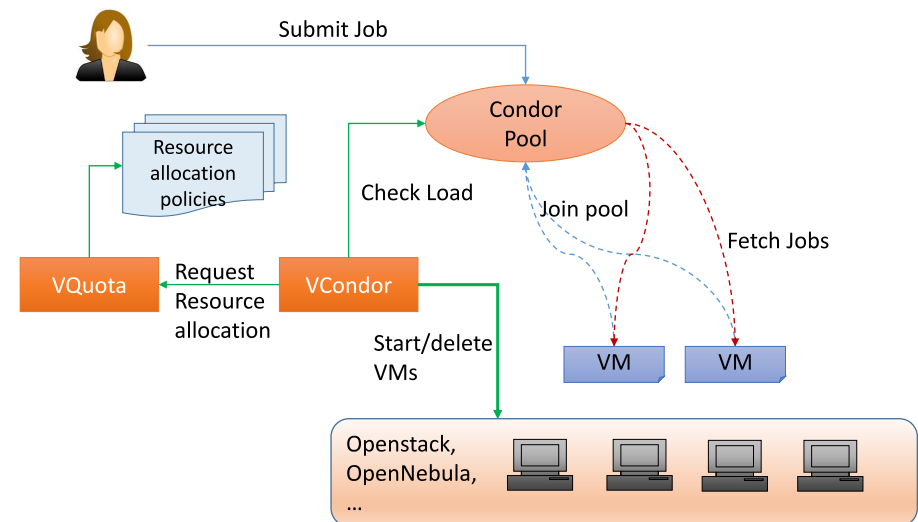
... and HTCondor-CE

More multicore

VCondor

LSF

SLAC testing 10.1



Github: <https://github.com/hep-gnu/VCondor>
Yaodong Cheng, <https://indico.cern.ch/event/531810/contributions/2316213>

Site Reports (II)

Monitoring

Elasticsearch

Logstash

Kibana

Time series DBs

OpenTSDB

InfluxDB/Grafana

Prometheus

Config management

Puppet 3 → 4

Salt

Ansible

Site Reports (III)

Service Management

ServiceNow → Service taxonomy → Evergreen

Lustre problems

Small I/O

Lost files

IPv6

More full-site IPv6

Often dual-stack

Docker

In batch jobs

CentOS 7, SL6 containers

Site Reports (IV)



RHEL, RHEV, Nagios, Splunk, home-made
config management, Lustre, JIRA

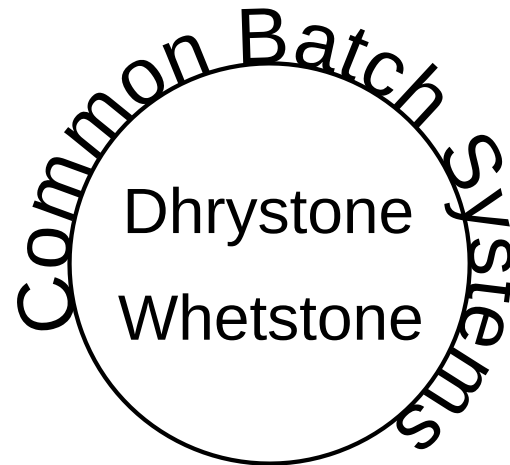
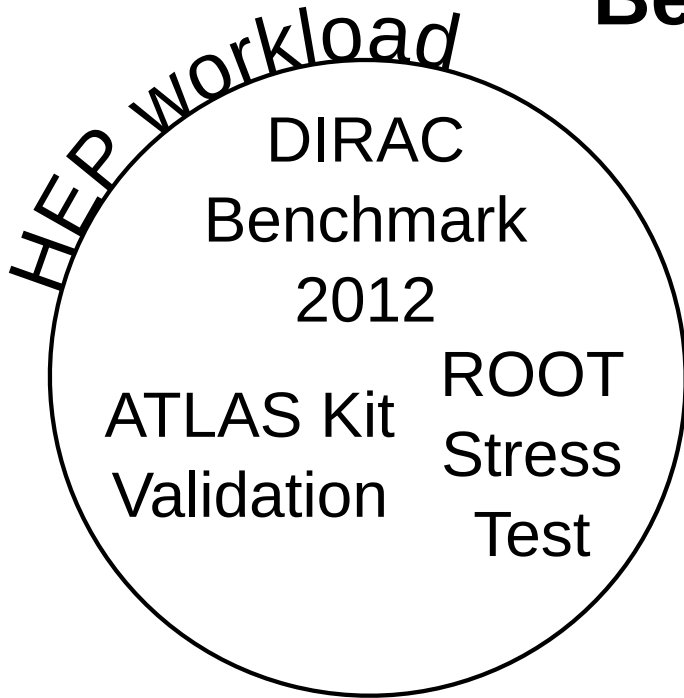
Lana Abadie, <https://indico.cern.ch/event/531810/contributions/2298936>

ITER

Computing & Batch Services (I)

Benchmarking Working Group

Long-running benchmark Fast benchmark



Next chip generations boosting

only HEP apps **only benchmark**

Computing & Batch Services (II)

Big data in physics and genomics

- **Physics:**
 - 10-100 PB
- **Growth rate:**
 - ~10's PB per year
 - Essentially linear
- **Known years in advance**
 - fixed by accelerator & experiment design
 - Can't build the experiment without some idea what the data looks like
- **Genomics:**
 - 10-100 PB
- **Growth rate:**
 - Doubling every 12-18 months
- **Unpredictable future**
 - Cheaper, faster sequencing
 - New sequencing methods
 - Lower bound: scarily fast



Tony Wildish, <https://indico.cern.ch/event/531810/contributions/2298944>

Computing & Batch Services (III)

Condor

condor_annex

More scalable → 500k cores

Faster schedd

Docker Universe

Grid Universe

Transform job ad on submit

Kerberos

More support for Grid Universe

HTCondor View

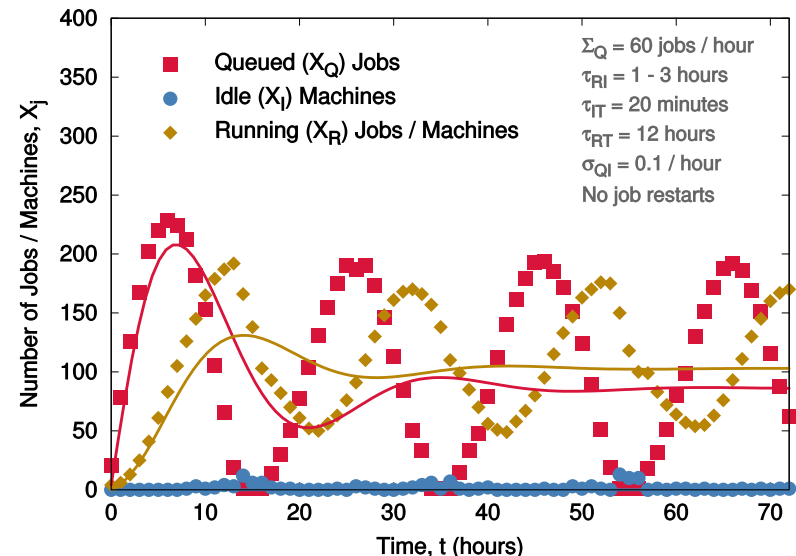
IPv6

Singularity

DAGMan

New command output

Submission Response Test: Exp. Results v. PM1 (ODEs)



Martin Kandes, <https://indico.cern.ch/event/531810/contributions/2311410>

Computing & Batch Services (III)

Condor

condor_annex

More scalable → 500k cores

Faster schedd

Docker Universe

Grid Universe

Transform job ad on submit

Kerberos

More support for Grid Universe

HTCondor View

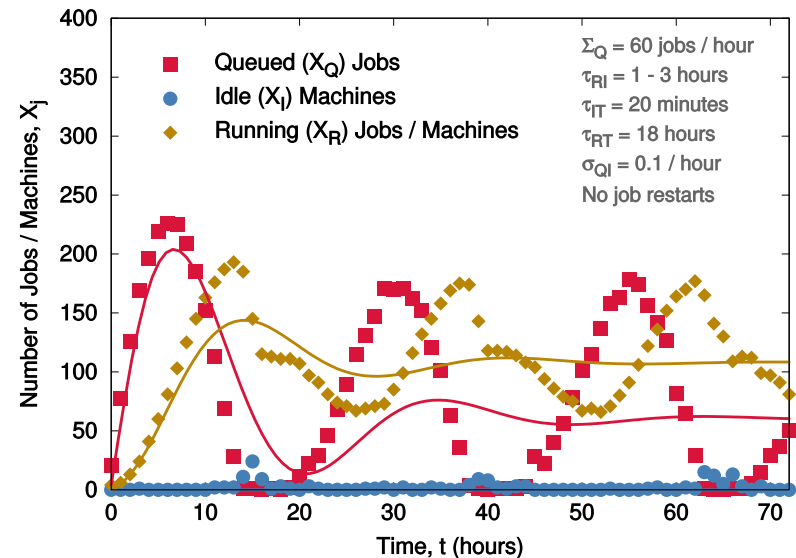
IPv6

Singularity

DAGMan

New command output

Submission Response Test: Exp. Results v. PM1 (ODEs)



Martin Kandes, <https://indico.cern.ch/event/531810/contributions/2311410>

Computing & Batch Services (III)

Condor

condor_annex

More scalable → 500k cores

Faster schedd

Docker Universe

Grid Universe

Transform job ad on submit

Kerberos

More support for Grid Universe

HTCondor View

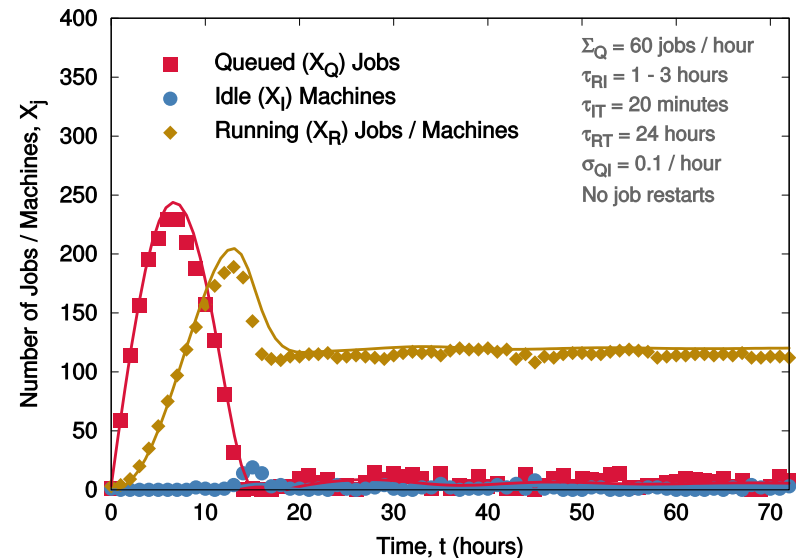
IPv6

Singularity

DAGMan

New command output

Submission Response Test: Exp. Results v. PM2 (DDEs)



Martin Kandes, <https://indico.cern.ch/event/531810/contributions/2311410>

Basic IT Services (I)

Monitoring

collectd,
Diamond

Graphite,
Grafana

Fifemon

Batch

Monitoring
Working Group

KEK: Kibana/Elasticsearch
access control with
Kerberos, patched Kibana,
Search Guard

Basic IT Services (II)

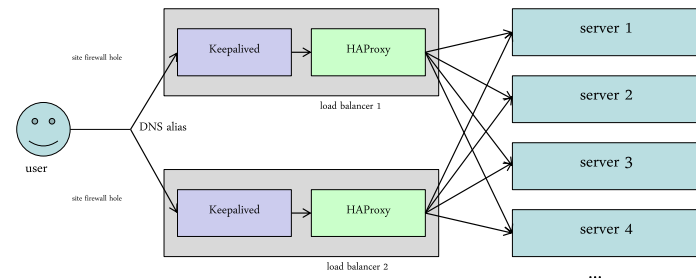
Architecture at RAL

And also...

RAL: load balancing with HAProxy, Keepalived

Fermilab: system management with Red Hat Satellite 6

- 2 floating IP addresses associated with the DNS entry for FTS3
 - traffic normally flows through both HAProxy instances



Ian Collier, <https://indico.cern.ch/event/531810/contributions/2314174>

Satellite Architecture

Open source upstream is Katello. It aggregates various open source products into a single collected workflow.

Includes:

- Puppet 3 (configuration management)
- Pulp (repository management)
- Foreman (External Node Classifier)
- Open SCAP (Auditing and compliance)
- Candlepin (subscription management)
- IPMI Web console
- System Administration Job scheduler

Rennie Scott & Patrick Riehecky, <https://indico.cern.ch/event/531810/contributions/2326347>

Outline – Part 2

Storage and Filesystem

- 12 presentations
- 3 presentations from CERN
 - EOS, DPM and FTS developments and plans (Andrea)
 - Update from Database Services (Katarzyna)
 - AFS phase-out at CERN (Jan)

Grid, Cloud and Virtualization

- 8 presentations
- 1 presentation from CERN
 - Update on HNSciCloud project (Helge)

IT Facilities and Business Continuity

- 4 presentations
- 2 presentations from CERN
 - Deploying Open Compute hardware at CERN (Marco)
 - CERN Computing Facilities Evolution (Wayne)

Storage and Filesystem (I)

Deep dive into Spectrum Scale – IBM

- Clustered filesystem, formerly GPFS, which exposes various interfaces (Spark, HDFS API, Object/Block/FS like access)
- Version 4.2.1 introduces **performance auto-tuning**
- Single and multithreaded performance improved (e.g. refactored communication code to be **lock free and improved NUMA awareness**: single thread RPC latency went down 50%, x4 single client throughput)
- Flash device be used as part of **pagepool kernel cache**

Experience of Development and Deployment of a Large-Scale Ceph-Based Data at RAL

- 13.6 PB RAW
- Access via GridFTP and xrootd using Ceph plugins. Identified and fixed several bugs in xrootd plugin
 - Incorrect **checksum on large files due to race condition in xrootd**
- Update to Jewel was problematic due to SL 7 requirement, infrastructure was is still not completely ready for SL7
- **Bug in kmem_alloc in XFS fixed in SL7**, filesystem on data disk needs to be re-created (technical support provided by Dan).

Storage and Filesystem (II)

CEPHFS: a new generation storage platform for Australian high energy physics

- CephFS to address the requirement for distributed POSIX-compliant filesystem
- Pre-production single threaded fio test showed that stripe unit can be selected to optimized performance if workload is known (random, sequential, block size)
- Production cluster now in Jewel: 305TB of raw storage, running under SL6, **compiled with gcc 4.8 due to C++11 deps**
- Clients use ceph-fuse implementation (more flexible then kernel implementation)
 - Had to disable fuse cache and set **fuse_disable_pagecache** to enable direct I/O (bypassing kernel pagecache) for consistent I/O on the same file from multiple clients
 - Tweak of **glibc memory pools** (MALLOC_ARENA_MAX)
- Had to tweak MDS configuration to cope with large directories

Ceph Based Storage Systems at the RACF (BNL)

- Two clusters, 1.2PB RAW on v0.87.2 and 1.8PB RAW on v9.2.1 (bound to SL6)
- One 800GB NVMe drive on each FE
- **FC attached storage arrays + Shared NVMe array over Infiniband** (24 x 400GB)

Storage and Filesystem (III)

Highly Available dCache (NDGF - NeIC)

- Upcoming version is 3.0 instead of 2.17
- **CEPH support**, storing files directly as Rados Block Device image, **support for HA configuration via HAProxy** with multiple SRM frontends/backends.

Effective Data Retrieval from Massive Amounts of Tape-Resident Data (BNL)

- ERADAT, file retrieval scheduler developed in-house to optimize tape mounts, forwards, rewinds
- **Requests aggregation and scheduling** (FIFO, On-demand)
- Very reliable, test randomly restoring 704x 10GB files out of 21 tapes: 270min, 34 mounts without ERADAT, 70 minutes 21 mounts with ERADAT

EOS, DPM and FTS developments and plans (Andrea, CERN)

- **Abstracted EOS namespace** interface to allow for multiple implementations
 - Persistent mechanism using xrootd and RocksDB (access to xrootd via Redis protocol, plugin upstreamed)
- New DPM operational mode that dropped SRM support, added caching
- New developments in FTS: improved optimizer, **added support for importing and exporting to object stores**, including cloud providers

Storage and Filesystem (IV)

ZFS on Linux (University of Edinburgh)

- 1PB of disk space under ZFS in GridPP
- Performance comparison RAIDZ2 vs HW RAID + ext4 or XFS
 - ZFS shows better performance for parallel read/write while hardware caching implemented by RAID controller yields better single file read performance
 - lz4 compression does not affect performance

OSiRIS: One Year Update (University of Michigan)

- Ceph based infrastructure for Michigan research universities
- Stressed Ceph installation when recovering data in-between two sites: **simulated increasing latency via Linux traffic control and netem**
 - Infrastructure halts at 320ms RT, recovery happens only at 80ms
 - More tests to be done with installation at SC16

Update from Database services (Katarzyna, CERN)

- Oracle plans: supporting UNICODE, migration to new hardware, maintain support for Oracle 11.1 and 11.2
- DB on demand: InfluxDB pilot service will move to production in Q1 2017, upgrade to MYSQL CE 5.7 in Q1 2017, PostgreSQL 9.6 in Q4 2016
- Hadoop ecosystem: Apache Kafka pilot service in place, Hadoop data backup in Q1 2017

Storage and Filesystem (V)

The future of AFS filesystem (Auristor Inc.)

- OpenAFS development activity has been decreasing, "struggling to stay afloat"
- AuristorFS as "2038 safe" AFS implementation, Zero Flag day conversion from OpenAFS, supports for recent features (IPv6, multifactor authentication).

AFS phaseout at CERN (Jan, CERN)

- Project in slow decline but still functional, phaseout considered necessary but in a controller fashion
- No single replacement, need to analyze case by case. Possible alternatives are: CERNBOX, EOS-FUSE, EOS, CVMFS, CASTOR.
- "Easy" use cases addressed in 2016, e.g. projects space, harder use cases to be tackled by the end of 2019

Grid, Cloud and Virtualization (I)

Chameleon: A Computer Science Testbed as Application of Cloud Computing (Argonne National Laboratory)

- Environment that allows to provision bare metal and virtualized resources for experimental purposes providing high degree of configurability (e.g. custom kernel version)
- Based on Openstack (e.g. Ironic for bare metal provisioning) and Grid'5000.
- Infrastructure now counts ~650 nodes (14500 cores), 5PB over 2 sites (100GB link)
- Various flavors of bare metal resources: Infiniband nodes, high memory nodes, NVMe, GPUs, plans for FPGAs and ARM

On-demand provisioning of HEP compute resources on cloud sites and shared HPC centers (KIT)

- HTCondor to tie resources at different German sites together, ROCED to include resource scheduling in the cloud
- Very successful project (tested with 1&1 cloud resources), but open questions: how to treat I/O intensive jobs, how to account for cloud resources utilization on experiment basis

Update on HNSciCloud project (Helge, CERN)

- Implementation of a hybrid cloud platform, potentially partially substituting in-house resources, but technical (e.g. data caching, network connectivity) and non-technical (procurement) hurdles
- HELIX NEBULA project in its initial phase, contracts awarded for design phase (prototype and pilot are the next steps)

Grid, Cloud and Virtualization (II)

Extending the farm to external sites: the INFN Tier-1 experience (INFN)

- Dynamic extension of the farm to commercial providers and remote INFN resources to cope with increasing demand (planned local resources are enough for Run3).
- Tested dynamic extension to Italian cloud provider (only on idle clock cycles): very good job efficiency on Monte Carlo jobs, low in average
- Tested static extension to remote INFN resources
 - Some hurdles along the way, e.g. GPFS caching for POSIX access to filesystem

The advances in IHEP Cloud facility

- HTCondor based batch system which interfaces with VCondor, in-house software that does dynamic Virtual worker nodes provisioning on the Openstack cloud

Running HEP Workloads on the NERSC HPC Systems

- Containers as a way to provide static execution environment to HEP workload on HPC systems (e.g. Cori, Cray XC, 1630 Haswell nodes + 9300 KNL compute nodes)
- Shifter preferred solution over Docker (“simplified” containerization, still using docker images)
- Cori “Burst Buffer” blades (high performance SSDs) provide buffering of data coming from I/O node, now non completely transparent (jobs must be annotated)
- SLURM with SDNs allows to configure network optimally based on job placement

Grid, Cloud and Virtualization (III)

Container Orchestration at RAL

- Mesos to orchestrate containers on 6500 cores
- Distributed private registry with Swift and Ceph as back end
- Integrated container orchestration in batch infrastructure based on HTCondor
- Related work: use Kubernetes to unify the way resources are provisioned on multiple commercial clouds, avoiding cloud provider specific APIs

CSNS Computing Environment Based on OpenStack

- Openstack, GlusterFS for volumes, live migration supported without any interruption
- Difficulties in scaling RabbitMQ, moved to broker-less architecture with 0MQ

IT Facilities and Business Continuity (I)

CERN Computing Facilities Evolution (Wayne, CERN)

- Second Network Hub currently being built in Prévessin, due to be ready in summer 2017
- 2nd Data Centre Project covered in-depth during ITTF on 23rd September
- Reorganization of the main room: installation of wider and deeper racks, 11 racks per row instead of 15

The role of dedicated computing centers in the age of cloud computing (BNL)

- Little space and power left at BNL Data Center (2.3MW) to support new programs
- Analysis of two solutions: **cloud resources vs repurposing different building (over 3 years)**
- Estimate based on AWS prices (Jul 2016) :
 - In-house is more cost effective for both computing and storage (-\$0.6M, -\$3.3M)
 - Considered spot market prices, which are not guaranteed resources

GreenITCube - Status & Monitoring

- Still only 2/6 levels equipped, no cooling problem apart from one incident with a cooling tower that was detected by the monitoring system
- 13k cores and 20 PB over 50 racks by end of 2016
- Current monitoring, reporting and alerting infrastructure being redesigned

Outline



- **Miscellaneous**
 - **3 talks in total**

Miscellaneous



- Plans to Support Data-Intensive Computing on the NERSC 8 System (NERSC)
 - NERSC - The National Energy Research Scientific Computing Center (NERSC) is the primary scientific computing facility for the Office of Science in the U.S. Department of Energy.
 - NERSC currently deploys separate Compute Intensive and Data Intensive Systems
 - Dramatically growing data sets require Petascale+ computing for analysis
 - Need to couple large-scale simulations and data analysis
 - Significant investments on Cori to support data intensive science
 - High bandwidth external connectivity to experimental facilities from compute nodes (Software Defined Networking)
 - NVRAM Flash Burst Buffer as I/O accelerator
 - Virtualization capabilities (Docker)
 - More login nodes for managing advanced workflows
 - Support for real time and high-throughput queues
 - Big Data Software
 - Progress is being made on a wide range of codes
 - Focus is on concurrency and locality
 - Increasing engagement from scientific facilities
 - Burst buffer
 - Real time queues
 - Software defined networking
 - Focus on improved scalability for deep learning
 - Including data benchmarks and workflows in planning for NERSC-9 (2020)

Outline



- **End User Services & Operating Systems**
 - **2 talks in total**

End User Services & OS(1)



- Scientific Linux Status Update (Fermilab)
 - SL 5 - less than 6 months remain - 31 March 2017 - current status is Production Phase 3
 - SL 6 - current version 6.8 from July 15
 - Upstream changes:
 - Updates to pacemaker clustering
 - OpenSCAP updates
 - Hardware support updates
 - Includes "Relax-and-Recover" backup tools
 - SL 7 - RHEL 7.3 in private BETA
 - Custodia
 - Distributes secrets safely to hosts
 - Targeted at cloud images, should work for anything
 - XFS Updates
 - XFS Statistics in /sys/fs
 - XFS tools rebased to much newer version
 - v3.2.2 -> v4.5.0
 - Enables Metadata CRC by default
 - Significant performance updates for auditd
 - New Setting: incremental_async should boost performance
 - Allow audit traps by process name
 - BETA Kernel has GPIO support enabled
 - No specific hardware drivers packaged

End User Services & OS(2)

- An e-mail quarantine with open source software (DESY)
 - DESY is hosting 70+ e-mail domains
 - Mixed environment of open source software and commercial products
 - Currently ~6.500 fully-fledged mailboxes
 - Daily ~300.000 delivered e-mails
 - Problems with MS-Office documents and macros
 - Mostly open source software solution (commercial virus scan engine still needed)
 - Using amavis, qpsmtpd and MariaDB for e-mail filtering
 - e-mail transport - Postfix
 - e-mail decomposition - amavis
 - classification - amavis
 - virus scanning - ClamAV

Outline



- **Security and Networking**
 - **11 talks in total, 6 from CERN**
 - **Adam K:** SDN-enabled Intrusion Detection System
 - **Adam S:** Pre-Studies for Wi-Fi service enhancement at CERN
 - **Vincent:** Wi-Fi service enhancement at CERN
 - **Tony:** Cloud Services – Network realities
 - **Hannah:** Security Update
 - **Hannah:** Can we trust eduGAIN?

Security & Networking (1)



- Platform Providing Network Awareness to ATLAS and Beyond (University of Chicago)
 - Aggregate and index network related data of interest not only to ATLAS but also WLCG, OSG communities
 - Provide a generalized network analytics platform
 - Network anomaly detection, alarm and alert system
 - Serve derived network analytics (eg. to ATLAS production, DDM & analysis clients)
 - Tools:
 - Each source has dedicated python collector
 - Elasticsearch
 - Fast visualization in Kibana
 - Data analysis on a co-located Jupyter cluster
 - perfSONAR
 - FTS
 - FAX
 - SVM for anomaly detection

Security & Networking (2)



- Upgrade of network connection between KEK and SINET (KEK)
 - Most of the connectivity for HEP researchers in Japan is provided by SINET
 - Migration from 10Gx2 for SINET4 to 100G+10G for SINET5
 - 10G is assigned ordinary Internet access from the campus network
 - 100G for inter-lab VPNs including LHCONE
 - New KECC is now connected to LHCONE
 - Improved FTS3 throughput especially for EU sites

Security & Networking (3)



- SDN-enabled Intrusion Detection System (CERN)
 - Setup: A traffic mirrored at the CERN firewall is distributed across a pool of 16 IDS servers
 - Both the bypass and the firewalled traffic are monitored
 - Decoupling the network control plane (decision logic) from the forwarding plane
 - Logic centralized in the SDN controller
 - Switching hardware programmed by external software
 - Improved flexibility and programmability
 - Features:
 - Symmetrical load-balancing
 - Traffic shunting - filtering out TCP data packets belonging to trusted flows
 - Selective mirroring
 - Building blocks:
 - OpenFlow
 - OpenDaylight
 - Brocade Flow Optimizer (BFO)
 - Project in a test phase

Security & Networking (4)



- SDN Implementation in IHEP (CC-IHEP)
 - Simple, flexible, robust, high performance and central—controlled network environment
 - High performance controller cluster is under researching
 - Improve the data exchange performance, based on the current network infrastructure and applications
 - Overlay: use IPv4 & IPv6 network link/Separately or Aggregately
 - Automatically and Dynamically network path choosing based on the application requirements and network performance status
 - SDN solution @datacenter - integration with openstack (own cloud platform based on OpenStack)
 - Future plans:
 - More network security devices will be connected to SDN switch
 - Connect all the network security devices including behavior auditing and flow analysis to SDN switch
 - Research on the high performance multi-controller solution and deploy it
 - Research on the Load balance between ODL controller and switches
 - Develop and deploy monitoring pages for the system

Security & Networking (5)



- Plans to support IPv6-only CPU on WLCG - an update from the HEPiX IPv6 Working Group (STFC-RAL)
 - Most storage solutions and protocols now work in dual-stack mode
 - dCache, DPM, StoRM
 - XrootD4, GridFTP, http
 - Several sites have been running dual-stack for some time
 - Production data transfers over IPv6 are happening
 - WLCG IPv6 deployment strategy:
 - Provide a viable migration path for sites needing to switch to IPv6
 - Allow sites to make long term planning decisions regarding their network setup
 - Allow VOs to make use of IPv6-only CPU resources should they become available in future
 - Plans to fully support dual-stack among Tier-1's by April 2018
 - Plans to have a large number of sites migrated their storage to IPv6 by end of Run2

Security & Networking (6)



- Security Update (CERN)
 - Zombie-Trojans attack Switzerland
 - Academia as a target
 - Dridex
 - Arrest made over kernel.org
 - Problems and challenges with macro's detection
 - Key is good malware detection, trigger malicious behaviour in an advanced test environment to analyse
 - Make use of the intelligence collected by external partners
 - What is done at CERN
 - Physical appliance forming a secure gateway to our mail servers
 - Appliance able to analyse attachments and links
 - Quarantines suspicious emails
 - How to protect an infrastructure and users
 - Sirtfi
 - MISP
 - WLCG SOC WG

Security & Networking (7)



- Pre-Studies for Wi-Fi service enhancement at CERN (CERN)
 - WiSE project at CERN
 - Controller-based Wi-Fi solution as modern way of design of Wi-Fi network
 - Rules for a planning of new CERN Wi-Fi network
 - Dedicated APs in meeting/conference rooms
 - APs installed in offices with a density of ~1 AP per 3 offices
 - Roaming provided inside buildings and across selected buildings complexes
 - New Wi-Fi infrastructure at CERN - planning, validation and deployment
 - RF planning and simulation
 - Preparation of plans: AutoCAD & WallMAN
 - RF Simulation: ProMAN
 - Site survey
 - Deployment
 - Cabling
 - APs installation
 - Post deployment site survey

Security & Networking (8)



- Wi-Fi service enhancement at CERN (CERN)
 - WiSE project at CERN
 - New, fast Wi-Fi infrastructure for CERN campus
 - Project organization
 - Long process...
 - Pilot in IT in November
 - General deployment starts beginning of 2017
 - Technical evaluations
 - Market Survey
 - Tests
 - Two vendors
 - Two months of tests
 - New Wi-Fi network infrastructure design
 - Controller-based Wi-Fi infrastructure, integrated with our current tools
 - Seamless roaming of user devices
 - Guest Wi-Fi

Security & Networking (9)



- Cloud Services – Network realities (CERN)
 - The problem is that many research sites have better network connections to NRENs than to the commercial internet. But some NRENs won't carry traffic from commercial cloud service providers...
 - There is also security aspect...
 - Some solutions are proposed by e.g. GÉANT and ESnet. They are good but not perfect...
 - High capacity network connectivity to cloud service providers is not a solved problem
 - Cloud service providers need a fat pipe to an exchange point where (a friendly) [N]REN is present
 - Start to think about how you can secure, control and isolate cloud resources

Security & Networking (10)



- Can we trust eduGAIN? (CERN)
 - Federated Identity Management (FIM) is the concept of groups of Service Providers (SPs) and Identity Providers (IdPs) agreeing to interoperate under a set of policies
 - eduGAIN is a form of interfederation inside FIM
 - eduGAIN
 - 38 federations
 - >3000 entities
 - Effective Security Incident Response
 - Operational Support
 - Shared Policies
 - Trust
 - Trust is not inherently present in eduGAIN...
 - Sirtfi
 - Security
 - Incident
 - Response
 - Trust Framework for
 - Federated
 - Identity

Security & Networking (11)



- Effective and non-intrusive security within NERSC's Open Science HPC environment (NERSC Security Group)
 - NERSC network
 - > 100Gbit Network
 - Thousands of users
 - SSH access and shell accounts for everyone
 - Passwords/Keys for unprivileged authentication
 - Highly diverse code base
 - Users can run what they want
 - Small number of incidents, but potentially high impact
 - Security tools and technics
 - Bro
 - IDS
 - Event Correlation
 - Monitoring
 - iSSHD

Next Meetings

WIGNER Research
Centre for Physics
Hungary
24 to 28 April 2017

www.hepixon.org

KEK
Japan
16 to 20 October 2017

Questions?





home.cern