

# **ALICE USA Computing Project Status Report**

R. Jeff Porter

ALICE T1/T2 Workshop, Strasbourg Fr

May 3rd, 2017

# Outline of Project Status Report

---

- **Overview of ALICE-USA Computing project**
- **2016 Operations**
- **Technical Items**
- **Project Changes for 2017 & Beyond**
- **HW Deployment Scenario Relative to US Obligations**

# Section I

---



- **Overview of ALICE-USA Computing Project**

# ALICE-USA Computing Project



- **Original 2009 Project Proposal**

- Goal to fulfill MoU-base ALICE USA obligations for compute resources to ALICE
- Funded by US Department of Energy (DOE)
- Operate facilities at 2 DOE labs
  - NERSC @ LBNL
  - Livermore Computing @ LLNL
- LBNL as the host institution

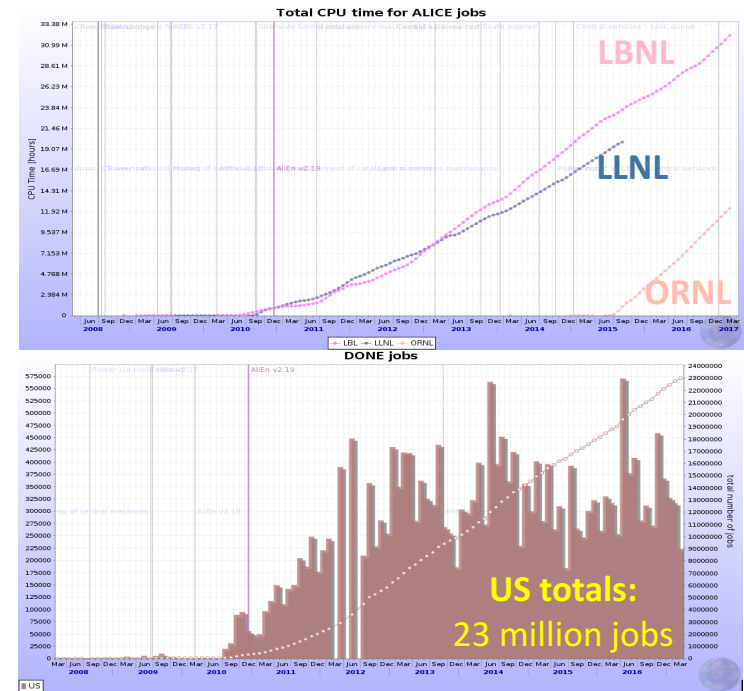
- **In operational since 2010**

- **New Project Proposal in 2014**

- Approved in Sept. 2014
- Replaced LLNL/LC with ORNL/CADES
- ORNL CADES T2 fully operational in 2015

- **Project working documents:**

- Project SLA: Institutions & roles
- Project Execution & Acquisition Plan: → **PEAP updated to DOE annually for funding approval**



# Project Organization & Computing Steering Committee

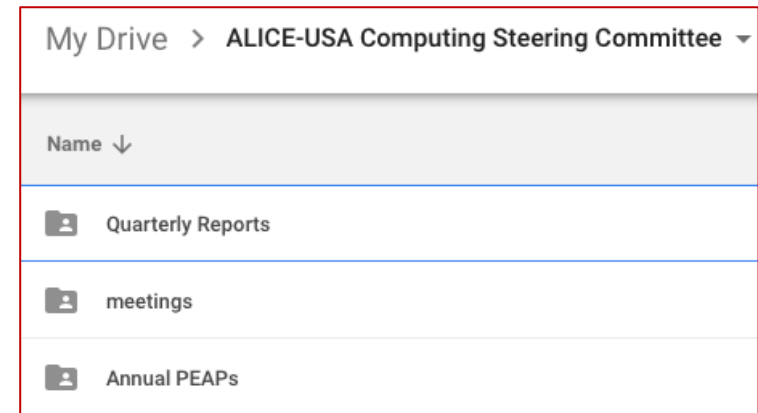
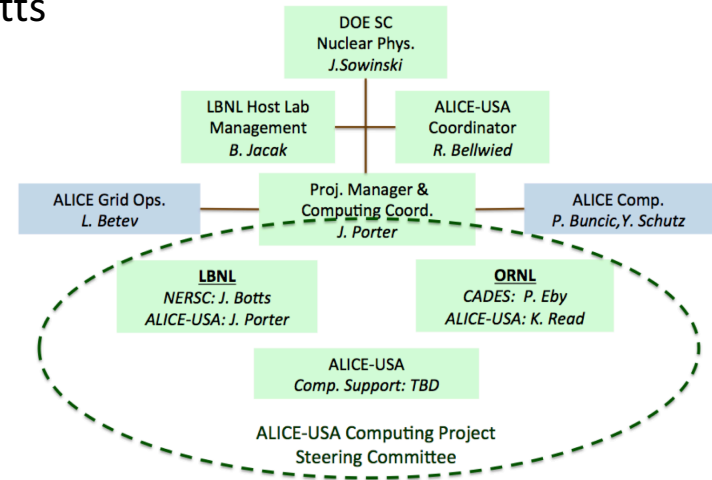


- **Project Steering Committee:**

- Currently: J.Porter, K.Read, P.Eby, M.Galloway, J.Botts
- Local documentation on LBNL Google docs:
  - Project Document repository
  - Monthly Meetings & minutes
- Project email list

- **Connection to ALICE Grid Operations**

- Alice-grid-task-force email list
- Annual US meeting with CERN team since 2012
- Annual ALICE T1/T2 workshops
  - 2012 @ KIT Germany: I. Sakrejda & J. Cunningham
  - 2013 Lyon, Fr: J. Cunningham & J. Porter
  - 2014 Tsukuba, Jp: J. Cunningham & J. Porter
  - 2015 Torino, Italy: J. Porter & P. Eby
  - 2016 Bergen, Norway: J. Porter, P. Eby & M. Galloway
  - 2017 Strasbourg, Fr: J. Porter & M. Galloway
- Annual AliEn Developers Workshops
  - 2010 – 2012, J.Porter
  - 2013, J. Porter & B. Nilsen



# ALICE-USA T2 Sites

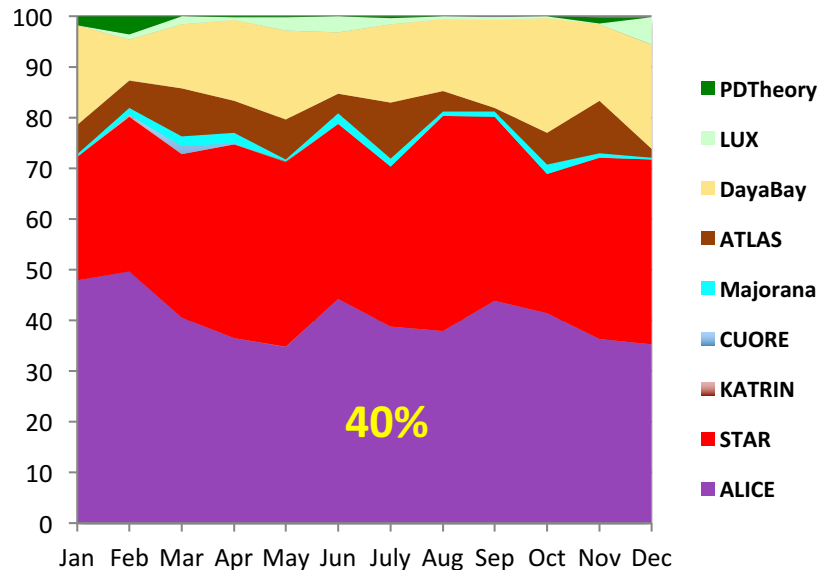


- **LBNL T2 at NERSC/PDSF**

- Shared HENP facility
- Operated by NERSC:
  - National User facility

<http://www.nersc.gov/users/computational-systems/pdsf/>

## Usage by Group



- **ORNL T2 at CADES**

- Single use (ALICE) Cluster
- Operated by CADES:
  - An institutional compute facility

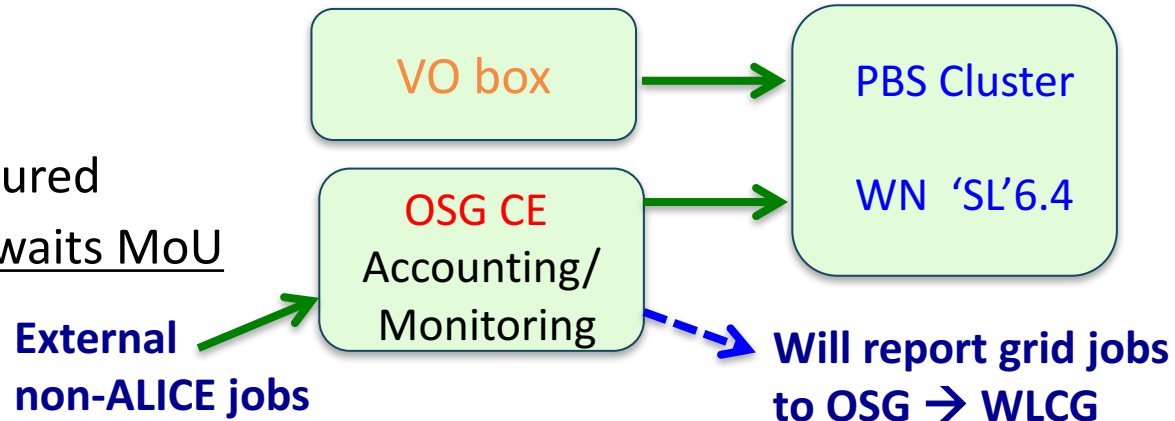
<http://cades.ornl.gov>



# US Site Configurations with OSG

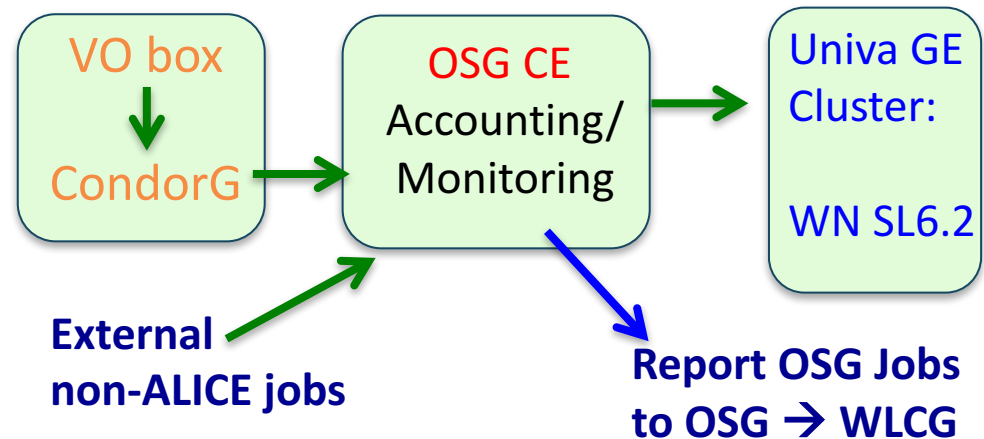
## ORNL CADES

- VObox submits to PBS
- OSG CE is being configured
- WLCG reporting still awaits MoU



## LBNL NERSC PDSF

- VObox submits to CondorG
- CondorG submits to OSG-CE
  - New HTCondor-CE
- OSG-CE submits to UGE
  - Only UGE site in OSG !?!
- OSG Accounting
  - Monitors batch logs



# Section II

---



- **2016 Operations**



# Site Job Profiles



## Ave. Running Jobs:

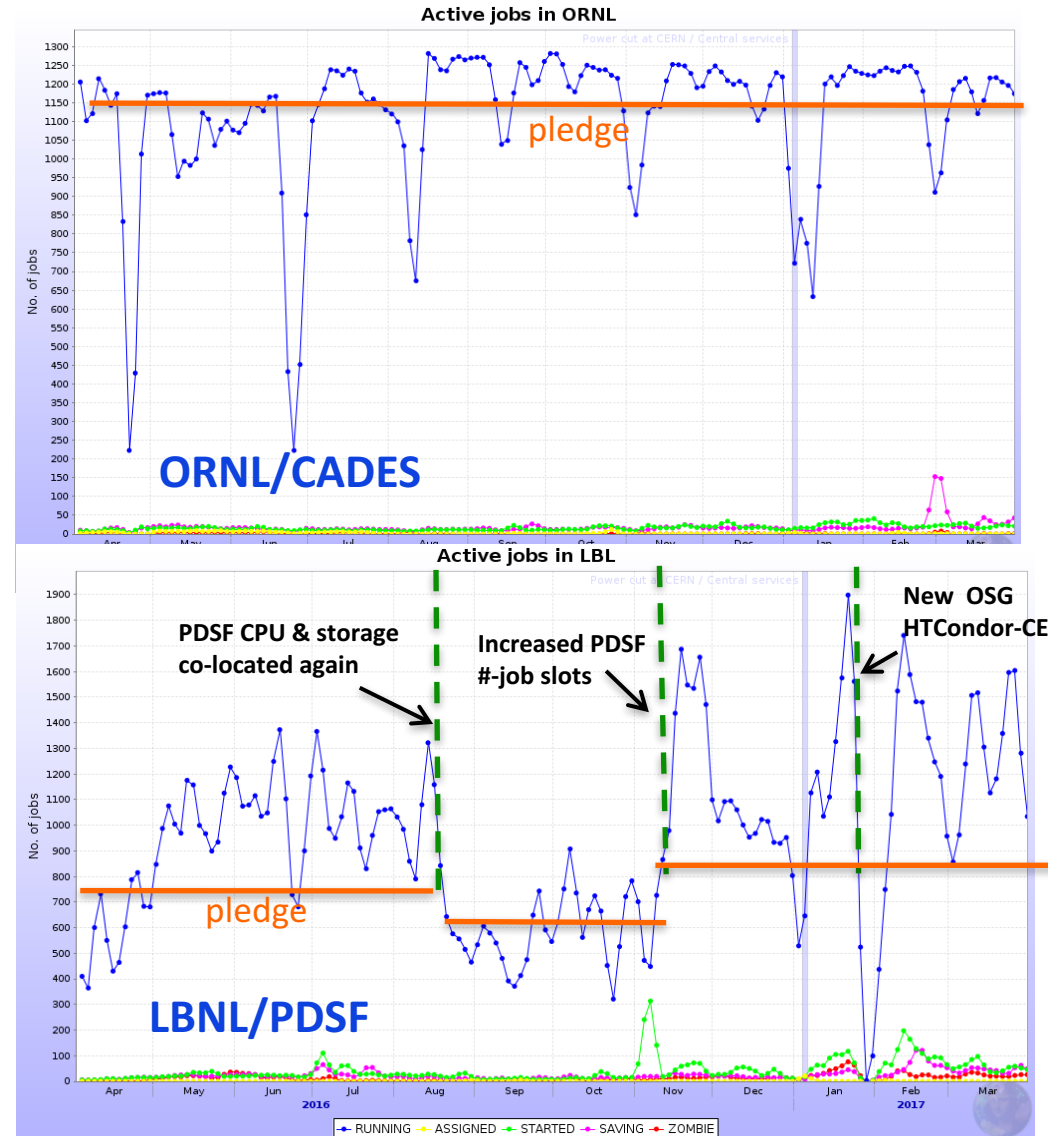
ORNL ~ 1150

LBL ~ 900 – 1500

## Zombie Rate remains low:

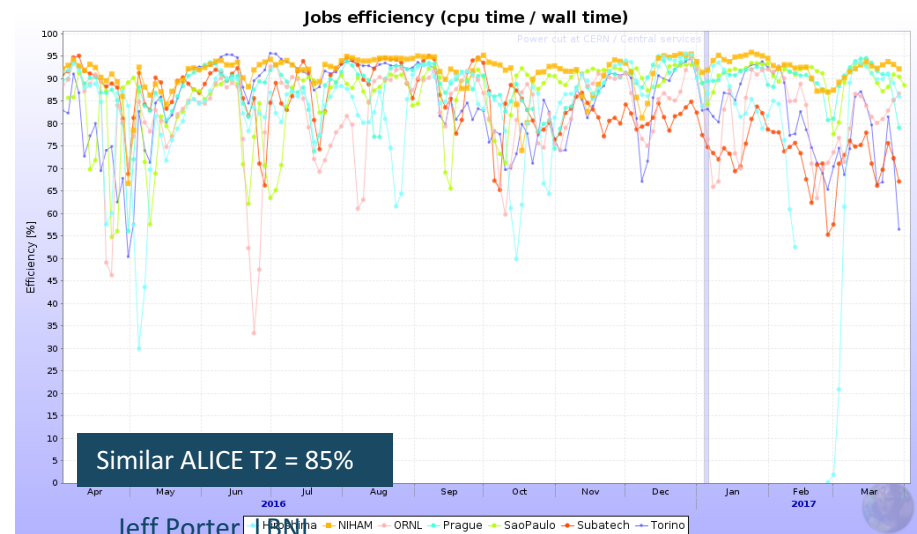
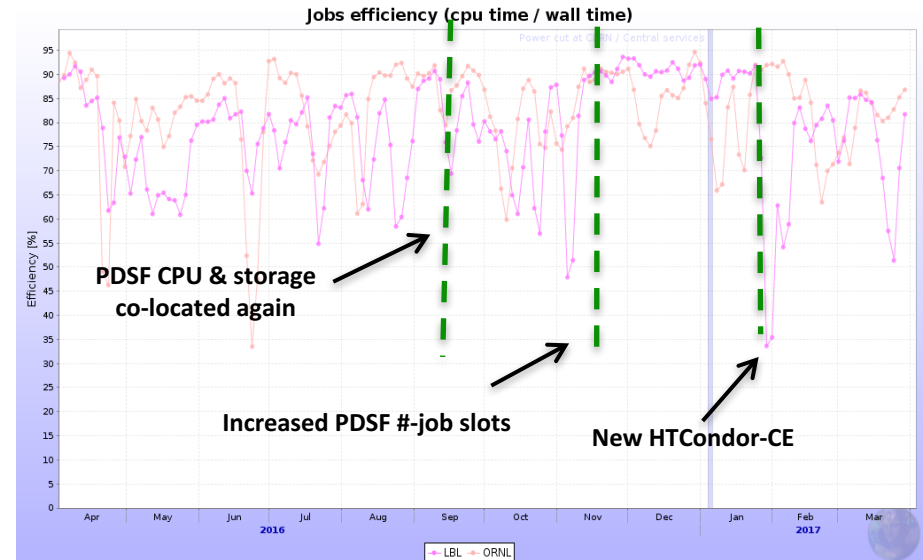
LBL ~1.0%

ORNL ~0.1%



# Site Efficiencies: cpu-time/wall-time

- **Ave Site Efficiencies**
  - ORNL 82.5%
  - LBL 79%
- **Results track other ALICE T2s**
  - Similar T2s ~ 85%
  - All ALICE T2s ~ 79%
- **Specific issues:**
  - PDSF operated CPU & storage at different sites until end of Aug
  - PDSF new HTCondor-CE problems



# CPU Delivered to ALICE Grid Relative to Pledged Obligations



## CPU Delivered Per US Site



| Site                      | <per-core> capacity (HS06/core) | CPU delivered (Mhrs)  | CPU delivered (MHS06-hrs) | US Obligation (pledge*hrs*0.70) | % delivered |
|---------------------------|---------------------------------|-----------------------|---------------------------|---------------------------------|-------------|
| <b>LBNL</b><br>12 kHS06   | 16.6, 19.8, 14.5                | 2.27+1.00+3.05 = 5.85 | 37.7+19.8+44.2= 101.7     | 73.8                            | 138%        |
| <b>ORNL</b><br>16.5 kHS06 | 14.0                            | 2.74+2.15+3.07 = 7.96 | 111.4                     | 101.5                           | 110%        |
| <b>Totals</b>             |                                 |                       | 213.1                     | 175.3                           | 122%        |

# Storage Capacity & Utilization

- **Storage Deployment History**

- LBNL NERSC

- LBL::SE XRootD based SE retired in September 2016
    - LBL::EOS installed 910TB, commissioned in Aug 2016
    - Project plans to add 600TB in FY17

- ORNL CADES

- ORNL::EOS installed 1090 TB in June 2015
    - 500 TB installed in 2016, added to ORNL::EOS in 2017
    - Project plans to add 300 TB more in FY17

**ALICE-USA Storage Elements Capacities & Usage: 04/2016**

| ALICE SE      | #-servers | Space (PB) | Used Space (PB) | % Used |
|---------------|-----------|------------|-----------------|--------|
| LBL::EOS      | 3         | 0.91       | 0.62            | 68     |
| ORNL::EOS     | 4         | 1.59       | 1.03            | 65     |
| In Production | 7         | 2.50       | 1.65            | 66     |

# EOS SE Availability



- **Writing**

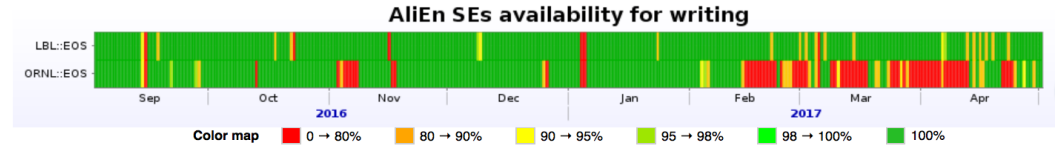
- LBL ~98%
- ORNL ~88%
  - >90% for year

- **Reading**

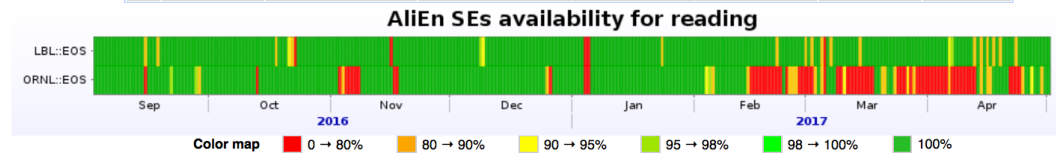
- LBL ~98%
- ORNL ~ 88%

- **ORNL EOS Problems began as SE became full**

- Extremely large thread count leading to high loads
- “solution” : auto restart FST triggered on high thread counts
- “solution” : add more storage



| Statistics  |                   |                   |                                     |        |               |         |
|-------------|-------------------|-------------------|-------------------------------------|--------|---------------|---------|
| Link name   | Data              |                   | Individual results of writing tests |        |               | Overall |
|             | Starts            | Ends              | Successful                          | Failed | Success ratio |         |
| ⚠ LBL::EOS  | 02 Sep 2016 00:18 | 01 May 2017 22:18 | 2856                                | 45     | 98.45%        | 98.44%  |
| ⚠ ORNL::EOS | 02 Sep 2016 00:22 | 01 May 2017 22:22 | 2544                                | 358    | 87.66%        | 87.69%  |



| Statistics  |                   |                   |                                     |        |               |         |
|-------------|-------------------|-------------------|-------------------------------------|--------|---------------|---------|
| Link name   | Data              |                   | Individual results of reading tests |        |               | Overall |
|             | Starts            | Ends              | Successful                          | Failed | Success ratio |         |
| ⚠ LBL::EOS  | 02 Sep 2016 00:18 | 01 May 2017 22:18 | 2853                                | 45     | 98.45%        | 98.44%  |
| ⚠ ORNL::EOS | 02 Sep 2016 00:22 | 01 May 2017 22:22 | 2542                                | 357    | 87.69%        | 87.73%  |

# Section III

---

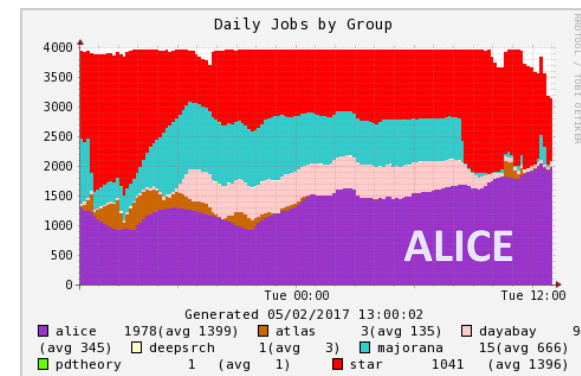
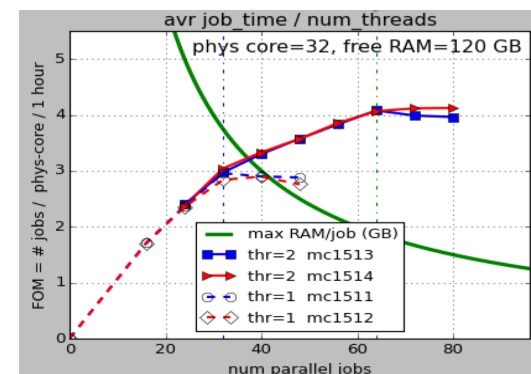
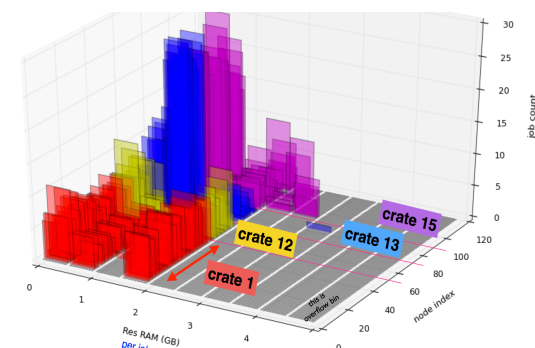


- **Technical Items**

# Maximizing the CPU capacity relative to available memory on PDSF



- Study of memory use on PDSF showed that we always had available memory
  - Due to the job mix from different user groups
- Nearly doubled # job slots (HT)
  - Per-core capacity dropped ~32%
  - Per-node CPU increased by ~25%
  - CPU efficiency was stable
- Also done at ORNL where we install 4GB/core but run #-Jobs at ~3GB/job-slot
  - We include a lot of swap space ... just in case
- Larger core count is not fully accessible by ALICE as there is a throughput limitation in the OSG HTCondor-CE
  - May considering removing that leg



# Working with EOS at ORNL

- **Life cycle management: How best to add new and retire old HW?**
  - Short answer – work with EOS team (Andreas)
  - Longer answer – use “eos scheduling groups” to carve up storage and retain generation coherence
- **Quotas & EOS headroom**
  - Needed to implement thread-monitor once FSTs became nearly full
  - & what do you know - it just crashes at 100% full. What to headroom to recommend for configuration?
- **Some good news:**
  - ORNL turned on ZFS compression on newest storage
  - See consistent 6%-7% compression rates without other observable impact
- **Lots of good work by ORNL folks:**
  - Slide deck from Pete Eby attached to the materials of current timeslot



# Section IV

---



- **Project Changes for 2017 & Beyond**

# Major Project Change in 2017:

---

- **PDSF Operational changes**

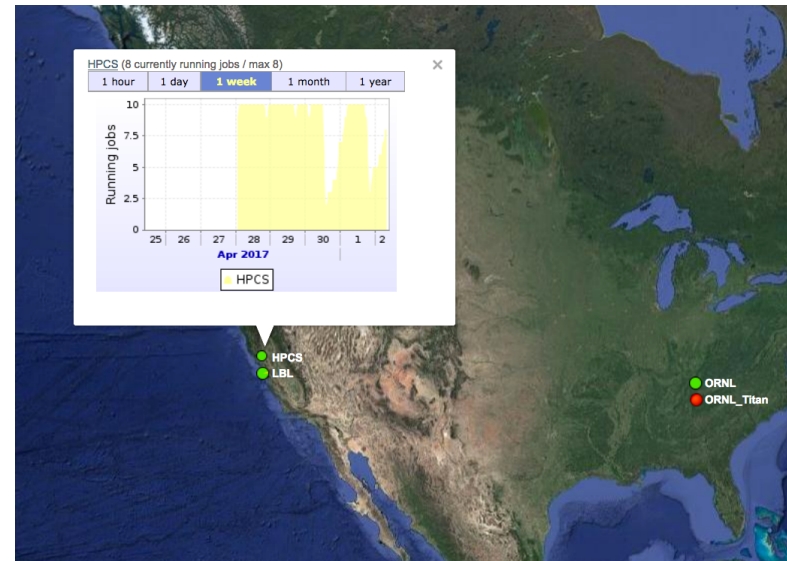
- Scheduled move from UGE to SLURM
  - Due to site-wide support
- Change virtualization: CHOS (chroot) to Shifter (docker-based container)
  - CHOS is no longer maintained
- Combined changes will make PDSF environ like NERSC HPC systems

- **PDSF Lifetime:**

- Cluster has existed for HENP Experiments since ~1997.
- Jan 2017 news: no room for a PDSF type cluster @ NERSC by ~2020.
- Mar 2017 news: no new hardware will be added to the cluster
  - HW capacity at PDSF will begin dropping next year, phase out by end of 2019
  - NERSC would like to support our work on their HPC systems
    - We do not view this as a current option

# Proposed New ALICE-USA T2 Facility at LBNL

- **SCS – Institutional Scientific Computing at LBNL**
  - <http://scs.lbl.gov>
- **Accelerated Timeline**
  - Began discussions in January
  - Met with the ALICE CERN team in March
  - Proof on Concept T2 establish in April
    - VObox & cluster on SL7
  - HW pricing is nearly finished
    - Expect initial HW purchase soon
- **Expectations**
  - SLURM-based, non-OSG stand-alone cluster
  - EOS SE independent of current LBL::EOS
- **New project plan is being formed:**
  - Transition from PDSF to HPCS over the next 2+ years
    - First as a stand alone cluster but evaluate evolution to a shared cluster
  - Maintain ALICE footprint at NERSC to use allocated HPC resources
    - A ~3.5 kHS06 allocation exists this year & expect similar allocations in coming years.



- **IPv6**

- All sites are ready @ border but EOS doesn't support IPv6

- **LHCONE**

- ORNL CADES

- local network team & ESnet LHCONE POC have completed talks
    - Still to do is for project to engage with LHCONE collaboration

- LBNL PDSF

- postponed, originally delayed by building move

- LBNL HPCS

- Being planned from the start

# Section V

---



- **HW Deployment Plans Relative to US Obligations**

# ALICE-USA Obligation Evaluation



- **ALICE Computing Requirements**

- Established Annually, reported to the ALICE Computing Board & approved by WLCG

**Table 1.** ALICE Computing requirements and corresponding ALICE-USA obligations.

| Year                           | FY2016 | FY2017<br>Apr 2016 | FY2017 | FY2018<br>Apr 2016 | FY2018 |
|--------------------------------|--------|--------------------|--------|--------------------|--------|
| <b>ALICE Requirements</b>      |        |                    |        |                    |        |
| CPU (kHS06)                    | 394    | 496                | 622    | 604                | 744    |
| Disk (PB)                      | 38.1   | 53.3               | 53.8   | 70.7               | 74.9   |
| <b>ALICE-USA Participation</b> |        |                    |        |                    |        |
| ALICE Total-CERN Ph.D.         | 573    | 585                | 585    | 585                | 585    |
| ALICE-USA Ph.D.                | 40     | 44                 | 44     | 44                 | 44     |
| ALICE-USA/ALICE (%)            | 7.0    | 7.5                | 7.5    | 7.5                | 7.5    |
| <b>ALICE-USA Obligations</b>   |        |                    |        |                    |        |
| CPU (kHS06)                    | 28.4   | 37.2               | 46.7   | 45.3               | 55.8   |
| Disk (PB)                      | 3.2    | 4.0                | 4.0    | 5.3                | 5.6    |

FY17 PEAP Update  
Submitted to DOE in  
Dec. 2016

- **ALICE-USA Obligations:**

- Fraction of ALICE Requirements Defined by proportion of ALICE-USA to ALICE

26% jump in CPU  
from original 2017 plan

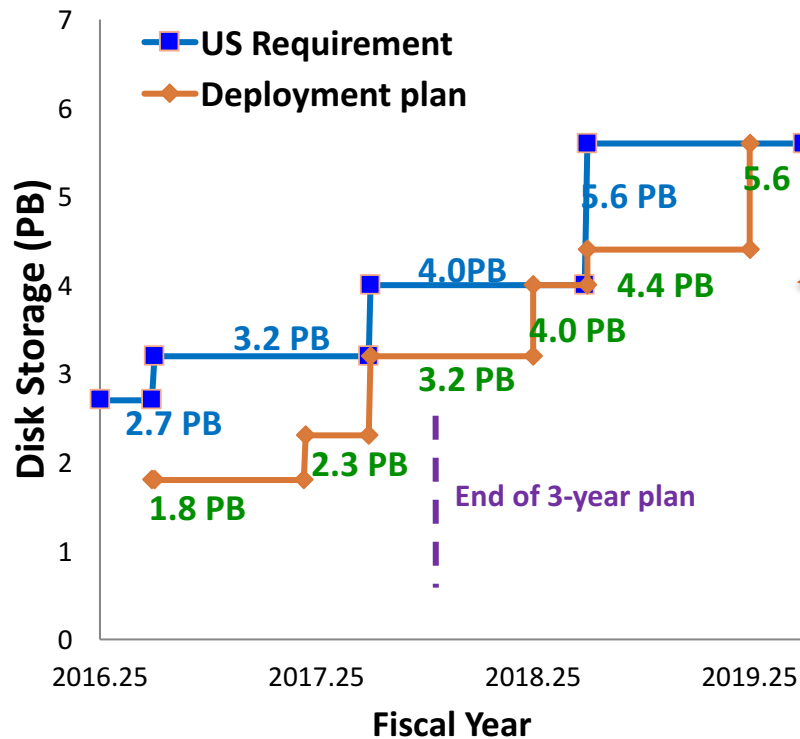
# 2017 PEAP Plan – Hardware

| Resource       | Currently Installed | 2017 Original | 2017 Dec. Plan |
|----------------|---------------------|---------------|----------------|
| <b>LBNL HW</b> |                     |               |                |
| CPU +/- kHS06  |                     | +5.0          | +10.0          |
| CPU installed  | 12.0                | 17.0          | 22.0           |
| Disk +/-       |                     | +0.6          | +0.6           |
| Disk installed | 0.9                 | 1.5           | 1.5            |
| <b>ORNL HW</b> |                     |               |                |
| CPU +/- kHS06  |                     | +3.5          | +7.5           |
| CPU installed  | 17.0                | 20.5          | 24.5           |
| Disk +/- (PB)  |                     | +0.3          | +0.3           |
| Disk installed | 1.6                 | 1.9           | 1.9            |

# 2017 PEAP Plan – Hardware

- **Target:**
  - 100% CPU on time
  - Disk lags with utilization

| Resource                     | Installed/FY16 | FY2017<br>Apr. 2016 | FY2017 | FY2018 |
|------------------------------|----------------|---------------------|--------|--------|
| <b>ALICE-USA Obligations</b> |                |                     |        |        |
| CPU (kHS06)                  | 28.4           | 37.3                | 46.7   | 55.8   |
| Disk (PB)                    | 3.2            | 4.0                 | 4.0    | 5.3    |
| <b>ALICE-USA Plan</b>        |                |                     |        |        |
| CPU (kHS06)                  | 29.0           | 37.5                | 46.5   | 54.0   |
| % CPU obligation             | 102%           | 100%                | 100%   | 97%    |
| Disk (PB)                    | 2.3            | 3.2                 | 3.2    | 4.4    |
| % Disk obligation            | 72%            | 80%                 | 80%    | 80%    |
| Disk deficit (PB)            | 0.9            | 0.8                 | 0.8    | 1.1    |



**Project funding is per US Fiscal Year, Oct-Sept, which is offset of RRB Year by 6 months**



# 2017 Hardware update

---

- **Currently behind schedule due to uncertainty on where to deploy HW, but expect to make HW procurements at both sites in May**
  - Full CPU target at LBNL
  - 80% Storage target at LBNL, remainder in fall 2017
  - Full CPU target at ORNL
  - Remainder of storage target later this summer

# 2017 Hardware update

- Bonus candy for Latchezar from ORNL after April's maintenance

