# Report on the Clusters at Fermilab

Don Holmgren
USQCD All-Hands Meeting
Fermilab
May 1-2, 2015

# Hardware – Current and Next Clusters

| Name | CPU | Nodes | Cores | Network | DWF | HISQ | Online |
|---|---|---|---|---|---|---|---|
| Ds (2010) (2011) | Quad 2.0 GHz Opteron 6128 (8 Core) | 421 | 13472 | Infiniband Quad Data Rate | 51.2 GFlops per Node | 50.5 GFlops per Node | Dec 2010 Aug 2011 |
| Dsg (2012) | NVIDIA M2050 GPUs + Intel 2.53 GHz E5630 (quad core) | 76 | 152 GPUs 608 Intel | Infiniband Quad Data Rate | 29.0 GFlops per Node (cpu) | 17.2 GFlops per Node (cpu) | Mar 2012 |
| Bc (2013) | Quad 2.8 GHz Opteron 6320 (8 Core) | 224 | 7168 | Infiniband Quad Data Rate | 57.4 Gflops per Node | 56.2 Gflops per Node | July 2013 |
| Pi0 (2014) (2015) | Dual 2.6 GHz Xeon E2650v2 (8 core) | 314 | 5024 | Infiniband QDR | 69.0 Gflops per Node | 53.4 Gflops per Node | Oct 2014 Apr 2015 |
| Pi0g (2014) | NVIDIA K40 | 32 | 128 GPUs 512 cores | Infiniband QDR | 69.0 Gflops per Node (cpu) | 53.4 Gflops per Node (cpu) | Oct 2014 |

# Storage

- Global disk storage:

  - 964 TiB  Lustre filesystem at /lqcdproj

  - ~ 6 TiB  "project" space at /project  (backed up nightly)

  - ~ 6 GiB per user at /home on each cluster  (backed up nightly)

- Robotic tape storage is available via *dccp* commands against the dCache filesystem at /pnfs/lqcd

  - Some users will benefit from direct access to tape by using *encp* commands on lqcdsrm.fnal.gov

- Worker nodes have local storage at /scratch

# Storage

- Two Globus Online (GO) endpoints:

  - usqcd#fnal – for transfers directly into our out of FNAL's robotic tape system.  Use DOE or OSG certificates, or Fermilab KCA certificates.  You must become a member of either the FNAL LQCD VO or the ILDG VO. There continue to be compatibility issues between GO and "door" nodes; globus-url-copy or gridftp may be a better choice for some endpoints.

  - lqcd#fnal – for transfers into our out of our Lustre file system (/lqcdproj).  You must use a FNAL KCA certificate.  See http://www.usqcd.org/fnal/globusonline.html

- Two machines with 10 gigE connections:

  - lqcdgo.fnal.gov – used for Globus Online transfers to/from Lustre (/lqcdproj), not available for interactive use

  - lqcdsrm.fnal.gov – best machine to use for moving data to/from tape.

# Storage – Lustre Statistics

- 964 TiB capacity, 793 TiB currently used, 138 disk pools
  (2013: 847 TiB capacity, 773 TiB used in 130 pools)

- 108M files  (85M last year)

- File sizes: 1.9 TiB maximum (an eigensystem file)
  7.67 MiB average  (9.52 MiB last year)

- Directories: 797K  (479K last year)
  801K files in largest directory

# Storage – Planned Changes

Our current Lustre software (1.8.8) is essentially End-of-Life (maintenance releases only), so we have started a second Lustre instance (2.5.3)

- Best feature of the new Lustre is data integrity checking and correction because of the backing ZFS filesystem

- We have also just completed the migration of our /project workflow area to ZFS

- By late summer we will migrate all existing data to the new Lustre

- Migrations to 2.5.3 will be done project-by-project

- We will attempt to make this as transparent as possible, but it might require a short break in running a given project's jobs

- We have hardware in hand to expand /lqcdproj by about 380 TB, but we will also retire old RAID arrays this year containing 100 to 150 TB

# Storage – Date Integrity

- Some friendly reminders:

  – Data integrity is your responsibility

  – With the exception of home areas and /project, backups are not performed

  – Make copies on different storage hardware of any of your data that are critical

  – Data can be copied to tape using dccp  or encp commands.  Please contact us for details. We have never lost LQCD data on Fermilab tape (~ 4 PiB and growing, up from 2.28 PiB last year).

  – At 138 disk pools and growing, the odds of a partial failure will eventually catch up with us

# Statistics

- May 2014 through April 2015 including JPsi, Ds, Dsg, Bc, Pi0, Pi0G
  - 360K  jobs
  - 271.6M JPsi-core-hours
  - 2.13 GPU-MHrs

- USQCD users submitting jobs:
  - FY10: 56
  - FY11: 64
  - FY12: 59
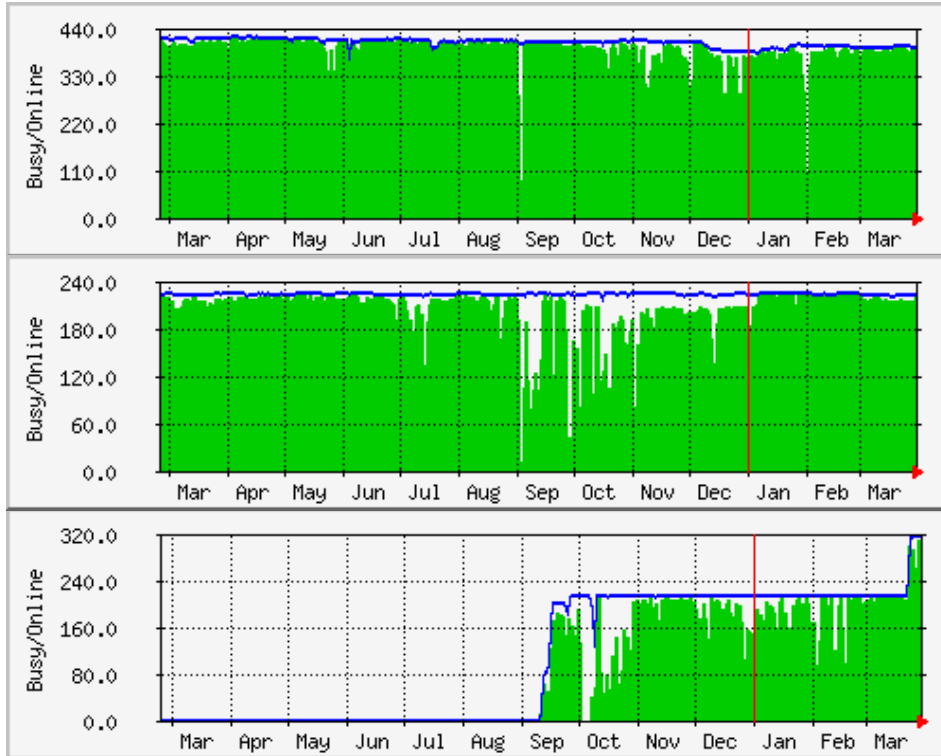  - FY13: 60
  - FY14: 53
  - FY15: 48 thru April

# Progress Against Allocations

- Total Fermilab allocation: 272.7M JPsi core-hrs
  2815 GPU-KHrs

- Delivered to date: 231.9M (85.0%, at 82.5% of the year)
  2083 GPU-KHrs (74.0%)

  – Does not include disk and tape utilization (roughly 14M + 1.5M)

  – Class A (19 total): 5 finished, 2 at or above pace (213M, 1729K)

  – Class B (4 total): 3 finished, 0 at or above pace (6.9M, 34.6K)

  – Class C: 8 for conventional, none for GPUS (2.8M, 0K)

  – Opportunistic: 5 conventional (9.5M), 4 GPU (320K)

- As was the case last year, a high number of Class A projects started
  late and/or are running at a slow pace

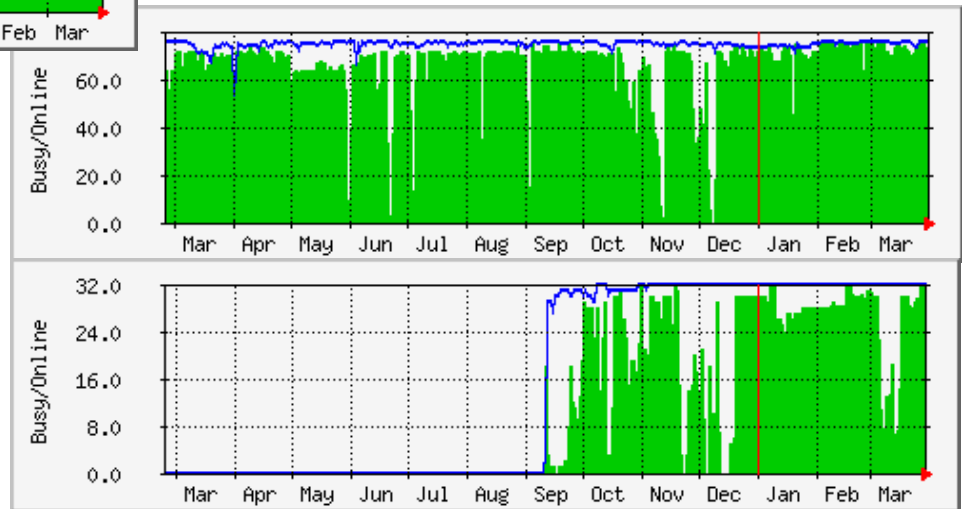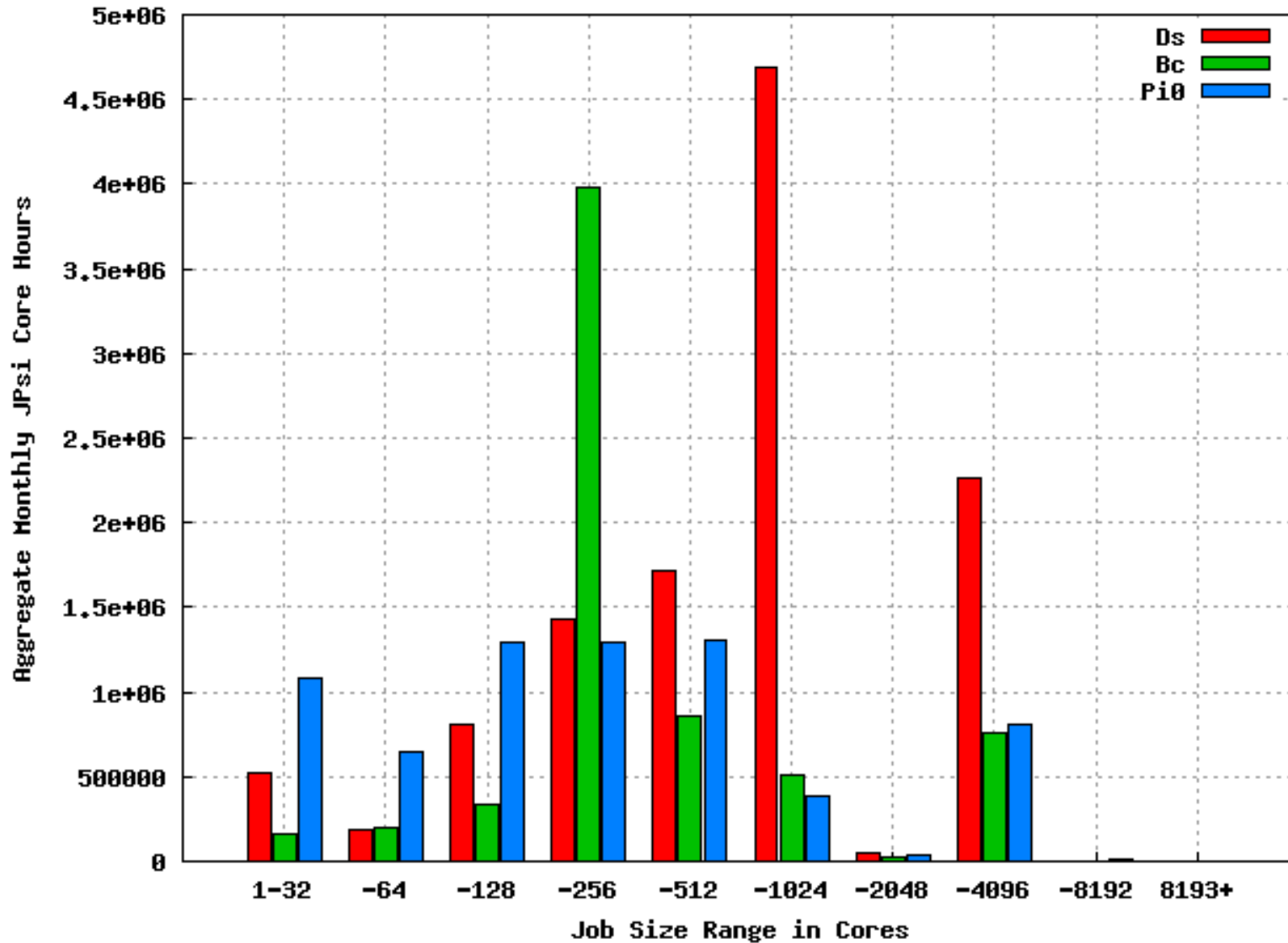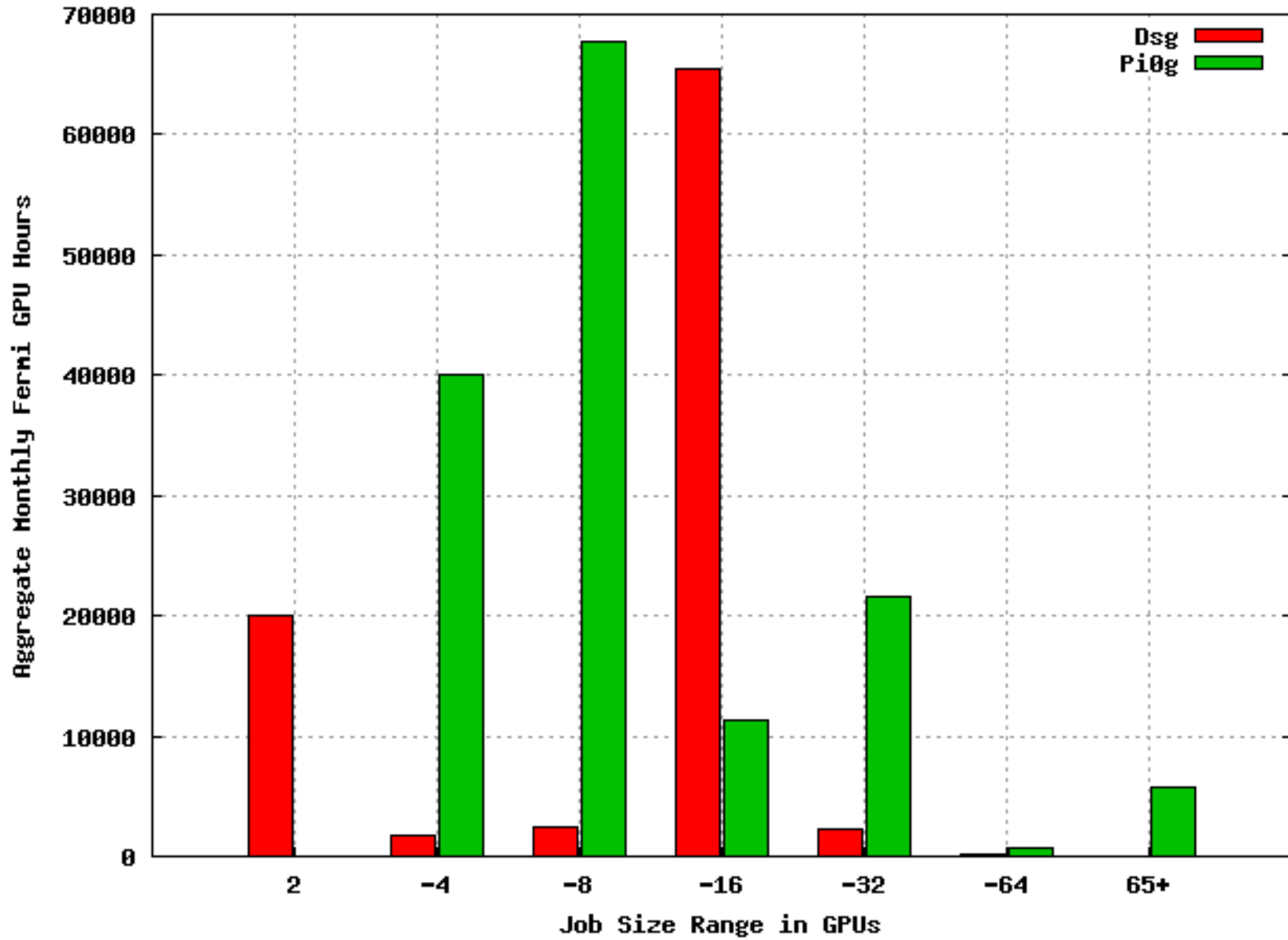# Fermilab LQCD Cluster FY14-FY15 Utilization

# Utilization



Ds

Bc

Pi0

Dsg

Pi0g

Fermilab May 2014 – April 2015 Job Statistics

Fermilab May 2014 - April 2015 Job Statistics

# Planned Operating System Upgrades

- Ds, Dsg, Bc all run Scientific Linux 5.  Pi0 and Pi0g use Scientific Linux 6.

- At the beginning of the program year (July 1), we will move Ds, Dsg, and Bc to SL6 to match Pi0 and Pi0g

  – User binaries will have to be rebuilt

- Once Lustre (/lqcdproj) has been migrated to the new version (2.5.3, late summer) we will upgrade the Infiniband fabric software

  – This allows use of Lustre 2.5.3 clients (1.8.8 clients during migration), which gives us more flexibility to move to newer Lustre versions

  – It will also enable all available Mellanox/NVIDIA GPUDirect optimizations on Pi0g

  – This IB upgrade will force us to rebuild MPI libraries.  Binaries that use shared MPI libraries may not require rebuilding

# User Support

Fermilab points of contact:

– Don Holmgren, djholm@fnal.gov

– Amitoj Singh, amitoj@fnal.gov

– Sharan Kalwani, sharan@fnal.gov

– Alexei Strelchenko, astrel@fnal.gov (GPUs)

– Alex Kulyavtsev, aik@fnal.gov    (Tape and Lustre)

– Yujun Wu, yujun@fnal.gov        (Globus Online)

– Jim Simone, simone@fnal.gov

– Ken Schumacher, kschu@fnal.gov

– Rick van Conant, vanconant@fnal.gov

– Paul Mackenzie, mackenzie@fnal.gov

– **Please use lqcd-admin@fnal.gov for requests and problems**

# Questions?