# Technology Upgrades for 2015

Brian Bockelman
USCMS Tier-2 Spring F2F

# Introduction and Disclaimer

- This talk is personal opinion of:

  - What we *should* do

  - What we *might* do

  - What we *likely won't* do

- During 2015.

# 2015 is not 2014

- The focus for March 2014 - March 2015 was **upgrade upgrade upgrade** before Run2 started.

  - Major projects involving production components will take longer than 2014.  We simply can't do as many as we did before.

- For reference, 2014 upgrade table:

  - https://twiki.cern.ch/twiki/bin/view/CMSPublic/USCMSTier2Upgrades

# Benchmarking

- As budgets are getting tighter in future years, it becomes important to better understand "hardware on the floor".

    - In particular, we've done a poor job of HS06 benchmarking.

    - There are many issues we all have with HS06, but it is the currency of the realm.

- Each site should make sure they run HS06 benchmarks on each of their worker nodes.

    - Note HS06 is a *per machine* benchmark, not a *per batch slot*. You need to take the host out of the batch system to test it.

- We'll work with you to do sanity checks across sites and make sure the math is right.

# Improved Accounting

- We'd like to capture as much Tier-2 usage in the "normal" accounting tools such as Gratia and Dashboard.  Particularly:

  - Make sure non-grid CMS usage is reported to Gratia.

  - If you have a local site data processing tool (like UW's *farmout*), then integrate it with Dashboard.

# Improved Local Access

- The CMS Submit Infrastructure team is working to help prioritize jobs belonging local users.

  - We've done this at various times in the past, but haven't emphasized this much. Standing queues were relatively short in Run1 and non-existent in LS1! We expect this to be more important in resource-constrained Run2!

- Three classes of work:

  - Pledge: Pilots with VOMS Role=pilot. Target fairshare is precisely the WLCG pledge. Set quota to (# of batch slots) / (site HS06) * (WLCG pledge).

  - US analysis: Pilots from the VOMS /cms/uscms group. Should get remaining batch slots.

  - Local users: Pilot (DN still TBD) that only runs jobs from local users. Batch system quota needs discussion!

- Local user case is under development; I'm looking for volunteer sites!

# Local Access: Site-customized Glideins

- **Basic idea:**

  - Site provides a list of "local users" for their site.

  - We send a special pilot to the site that *only* runs the jobs from "local users" at that site.

    - Sites can prioritize this pilot's DN appropriately.

    - Local user jobs will run on either the special or regular pilot at the site.

- Advantages: Simple for site and central CMS.  Just provide the list of users!

- Disadvantages:

  - Site cannot prioritize between multiple local users.

  - Sites must provide a CE (although no BDII needed!).

# Local Access: More Details

- Sites would provide:

  - The list of CMS user names in a flat-formatted list by placing an appropriate file in SITECONF. These users are added to group $SITENAME.

  - The list of groups considered "local" to the site (also placed in SITECONF).

    - Allows T2_US_Nebraska to state "all T1_US_FNAL users are considered local".

- At schedd, we would map job->group; pilot would just maintain list of allowed groups.

- The "local user" pilot is the same DN at all sites. Uses command-line arguments to determine local site.

  - Prevents "T1_US_FNAL" local users from running at "T2_US_Nebraska".

  - Unlike national VOMs, only one new grid proxy is needed for the whole setup - significantly less work than having 50 pilots for 50 sites.

  - In gWMS terminology, this setup would be run with a single group, not a group-per-site. There is *no* work to add new sites past the first site.

# The Long Slog of IPv6

- As CERN continues to run low on IPv4 addresses for worker nodes, sites should continue deploying services with IPv6. Priority order:

  - Xrootd servers.

    - IPv6-only worker nodes will otherwise not be able to access your site via AAA!

  - SRM/GridFTP servers.

  - Worker nodes.

  - CE.

# Shutdown older services

- All USCMS T2 sites have at least one HTCondor-CE present.

    - It's now time to make sure you've decommissioned GRAM.

- BDII is only use for SAMv3.

    - However, we really SAMv3 to test the same endpoint as the glideinWMS factory.

    - In the next few weeks, Dashboard team is supposed to roll out a new update to their topology script to additionally take.

    - By May, it should be possible to turn off GIP / osg-info-services.

# Xrootd Testing & 4.2

- Xrootd 4.2 is "just around the corner".

  - Filesystem throttle code has been merged into release.

  - Extensive python bindings - you can now do simple testing of your server without forking clients.

  - Many improvements to the caching proxy.

- As AAA becomes more important during Run2 - and possibly further in 2016 if we start offering a more unified disk service - we should redo the 10Gbps test from 2014.

# HTCondor-CE

- HTCondor-CE updates were fast and furious throughout 2014 as new features came online.

  - Thanks to all who were on the leading edge of the transition!

- For 2015, critical updates should be less frequent.

  - I'm only one major bug for CMS - fix for HTCondor #4915 - that we'll bother everyone to fix in the next 1-2 months.

- As typical, new releases may be very interesting to subsets of sites.  There will likely be no need to organize upgrades.

# HTCondor 8.4

- The 8.3 series has focused on a few items:

  - Scaleability improvements.  Significant decrease in shadow memory use.  Significant decrease in schedd<->startd communication.

  - Python bindings: almost all day-to-day admin activities can be done from a python script.

  - The usual range of small fixes and improvements that accumulate over a year.

- 8.4.0 will land in June or July; everyone should upgrade during 2015.

- See http://research.cs.wisc.edu/htcondor/manual/v8.3/10_3Development_Release.html

# PerfSonar

- There will be a new perfsonar release later this year, resulting in a coordinated upgrade.

  - From James - how should we make better use of this tool?

  - Any ideas besides SAM tests?

- Almost every site has issues to address.

# HDFS Upgrade? MAYBE

- The latest upstream HDFS has many features we don't leverage (HA, ACLs, extended attributes, snapshots, zero-downtime / rolling upgrades).

  - All good things, but perhaps not enough to upgrade.

- The under-development version of HDFS features erasure coding; duplication overhead is 40% instead of 100%.

  - This is sufficiently interesting to trigger an upgrade.

  - Not clear when this work will be released. My estimate is Fall 2015.

- In the meantime, consider installing the HDFS healer from the UCSD team (duplication overhead goes to 0% for some subset of the files).

# Remove SRM?  MAYBE

- About 2 years ago, Nebraska attempted to remove all use of SRM from the site.

  - Technically, everything seemed to work - but we never put things into production.  The setup has bit-rotted since then.

  - Since then, bestman2 hasn't gotten any younger (or really received more maintenance).

- It may be the time to try again.  Looking for an interested site; will require an admin who is not afraid to write some python patches for WMCore and SAMv3 tests.

# Docker Support?  MAYBE

- Docker support takes the existing HTCondor containerization features to the next level.

  - Adds network isolation, which has not been merge to HTCondor master branch.

  - More "familiar" in tech literature than HTCondor's container work.

  - **Importantly**, does a marvelous job in helping to create and manage runtime environments.

  - Looks like the "docker universe" will land in ~8.3.6.

- Assuming the underlying HTCondor support, integration with the CE is straightforward - simply need to change the routes.

  - Sites may prefer to wait for RHEL7 worker nodes for better OS integration.

  - Looking for volunteers.

# Other Items for Not This Year

- Things we shouldn't expect for this year:

  - Widespread RHEL7 support in CMS. If you want to use RHEL7, you'll need to run CMS jobs inside a chroot.

    - CMS will validate RHEL6-jobs-on-RHEL7 sometime this year.

  - Widespread multicore usage of Tier-2s. Multicore focus is on Tier-1 for now as only RECO use case is supported.

    - If desired, we can do multicore pilots with single core jobs (as at Nebraska and Purdue).

    - We may do RECO tests at Tier-2s sometime this year.

# Proposed 2015 Checklist

- Here's the summary:

  - Benchmarking

  - Review accounting

  - Implement local access

  - Shutdown GRAM & BDII

  - IPv6 (esp. Xrootd)

  - HTCondor 8.4 upgrade

  - Track down missing / broken perfSonar data.

  - AAA scale test (if possible?)