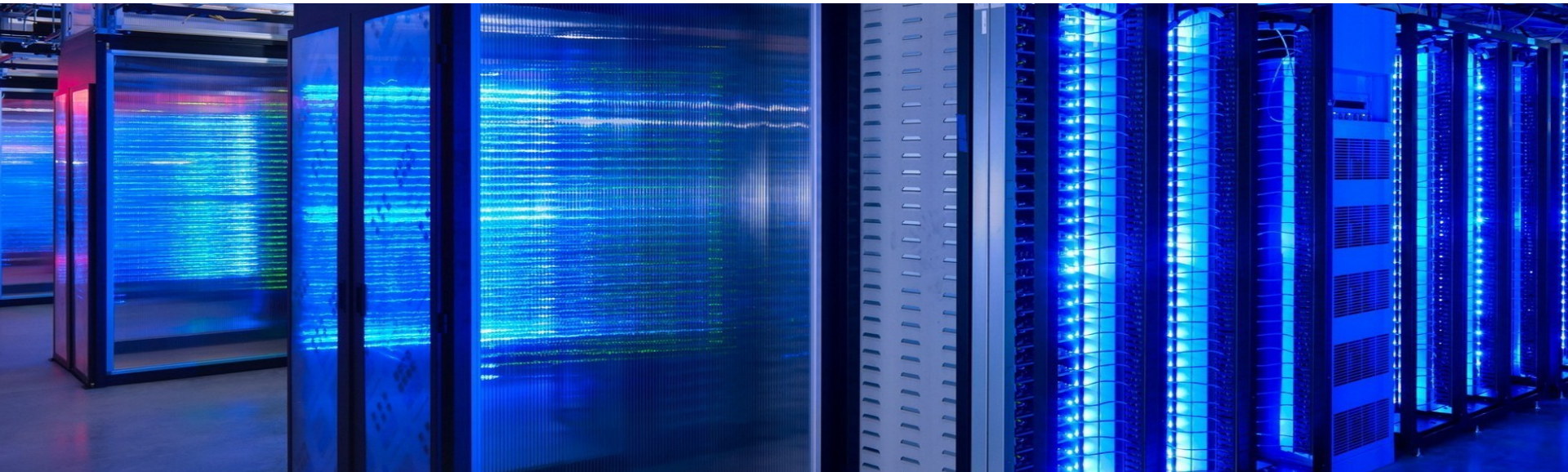# HTCondor as the Job Gateway: HTCondor CE

**Brian Lin, Marian Zvada**

OSG-AHM 2015
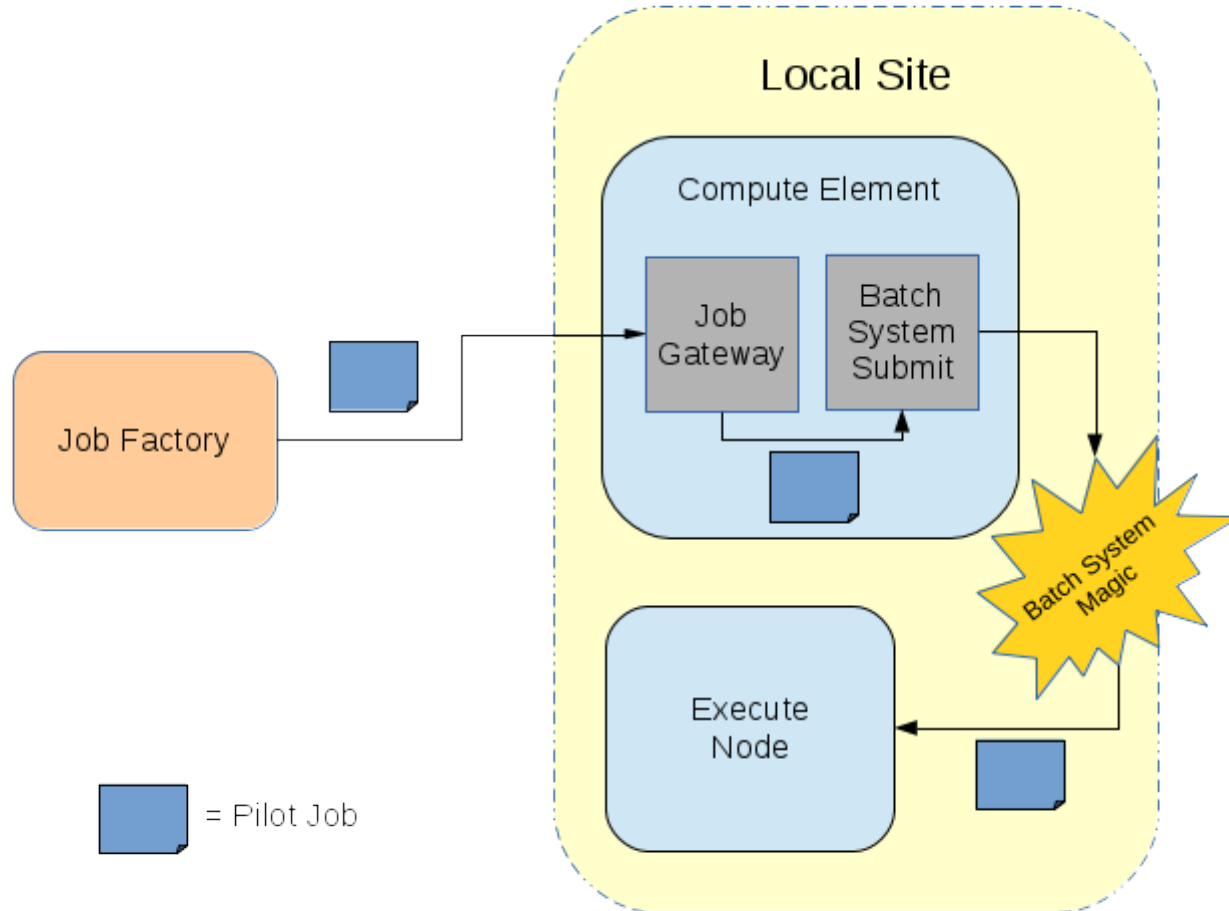Northwestern University, March 26

# What is a Compute Element?

- A Compute Element (CE) is OSG's entry point to a site's local resources

- Job gateway software is the main driver:

  - Job routing

  - Remote job submission

  - Job authorization

*Also known as a gatekeeper in the Globus world

# Submitting Pilot Jobs in the OSG

# Why use HTCondor CE as your Job Gateway?

- Scalability – supports job workloads of large sites

- Job routing as separate configuration

- Debugging tools – Built-in HTCondor and CE specific tools

# HTCondor CE (I)

- Currently, Globus GRAM provides the abstraction, sandbox movement, and remote submission layers for the OSG CE.

# HTCondor CE (II)

- Currently, Globus GRAM provides the abstraction, sandbox movement, and remote submission layers for the OSG CE.

- Why HTCondor-CE?

    - With the HTCondor team, the OSG has been working to provide an alternate job gateway implementation, the HTCondor CE.

    - The HTCondor CE is a special configuration of the HTCondor software which provides the three core pieces of functionality described previously.

# HTCondor CE overview (I)

- Special configuration of HTCondor

- Sits on CE* of each cluster (submit host)

- Allows:
  - Remote job submission and management
  - Strong authentication (GSI/VOMS)
  - Logging and monitoring
  - Scalability
  - Work with existing batch systems

# HTCondor CE overview (II)

- HTCondor CE provides (based on three fundamentals of the CE concept):

- Remote access

  - Based on the internal CEDAR protocol.
  - CEDAR provides a RPC and messaging mechanism over UDP or TCP, and can provide various levels of integrity or encryption based upon the session parameters.

- Authentication and authorization

  - Based on Globus libraries for GSI and authorization callout.

- Resource allocation

  - Grid jobs are taken and transformed to local jobs using the JobRouter component.

    - Any software HTCondor can interact with is a potential backend. This includes EC2, OpenStack, or even another HTCondor CE!
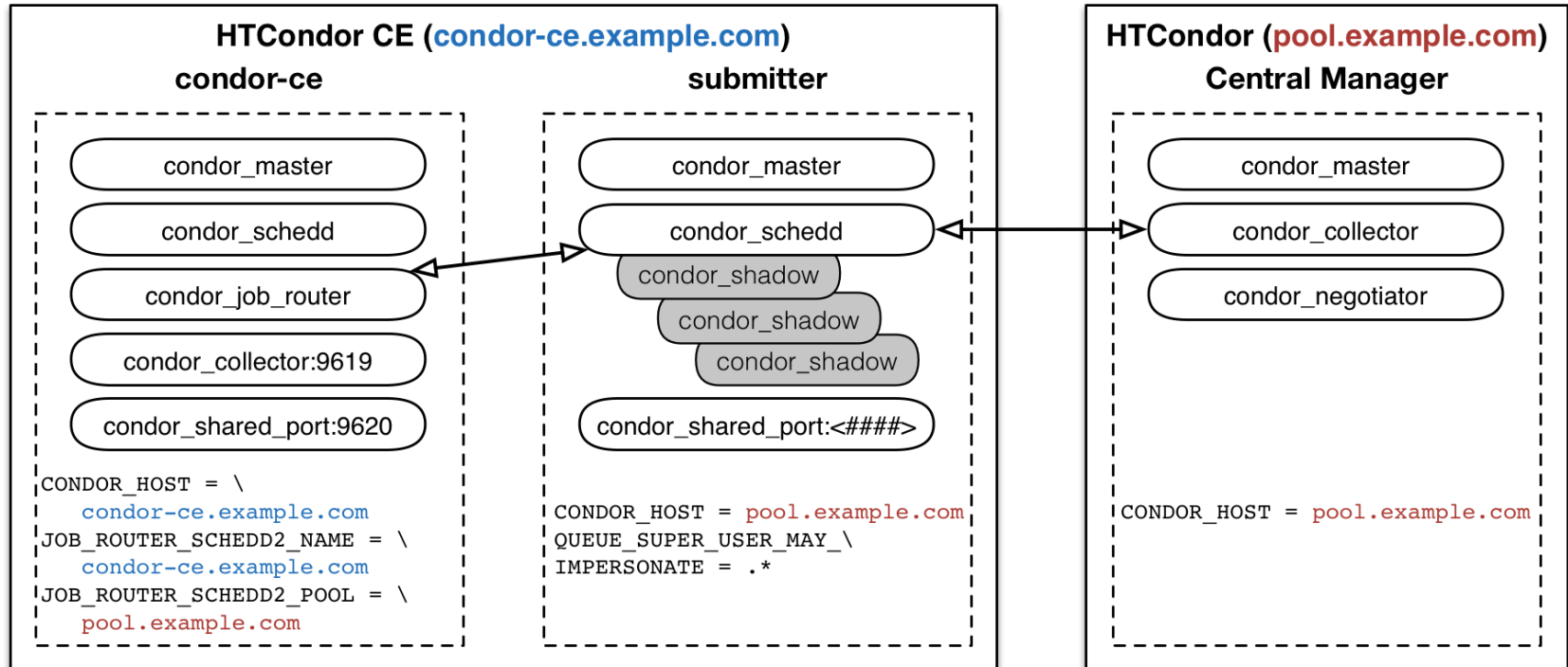
# HTCondor-CE Building Blocks

- ## HTCondor-C
  - Submit jobs from one HTCondor scheduler to another (submit machine to CE)

- ## Job Router
  - Transform jobs (localize jobs at CE)

- ## BLAHP
  - Submit jobs to non-HTCondor batch systems (PBS, SGE, SLURM, etc.)
  - blahp is the executable which then calls, for example, qstat / qsub / qdel.
  - blahp has another layer of customization if, for example, you need to tweak qsub arguments. Most useful things can be done via the JobRouter transform.

# HTCondor-CE Building Blocks

- ## HTCondor-C
  - Submit jobs from one HTCondor scheduler to another (submit machine to CE)

- ## Job Router
  - Transform jobs (localize jobs at CE)

- ## BLAHP
  - Submit jobs to non-HTCondor batch systems (PBS, SGE, SLURM, etc.)
  - blahp is the executable which then calls, for example, qstat / qsub / qdel.
  - blahp has another layer of customization if, for example, you need to tweak qsub arguments. Most useful things can be done via the JobRouter transform.

- ## HOW IT WORKS?
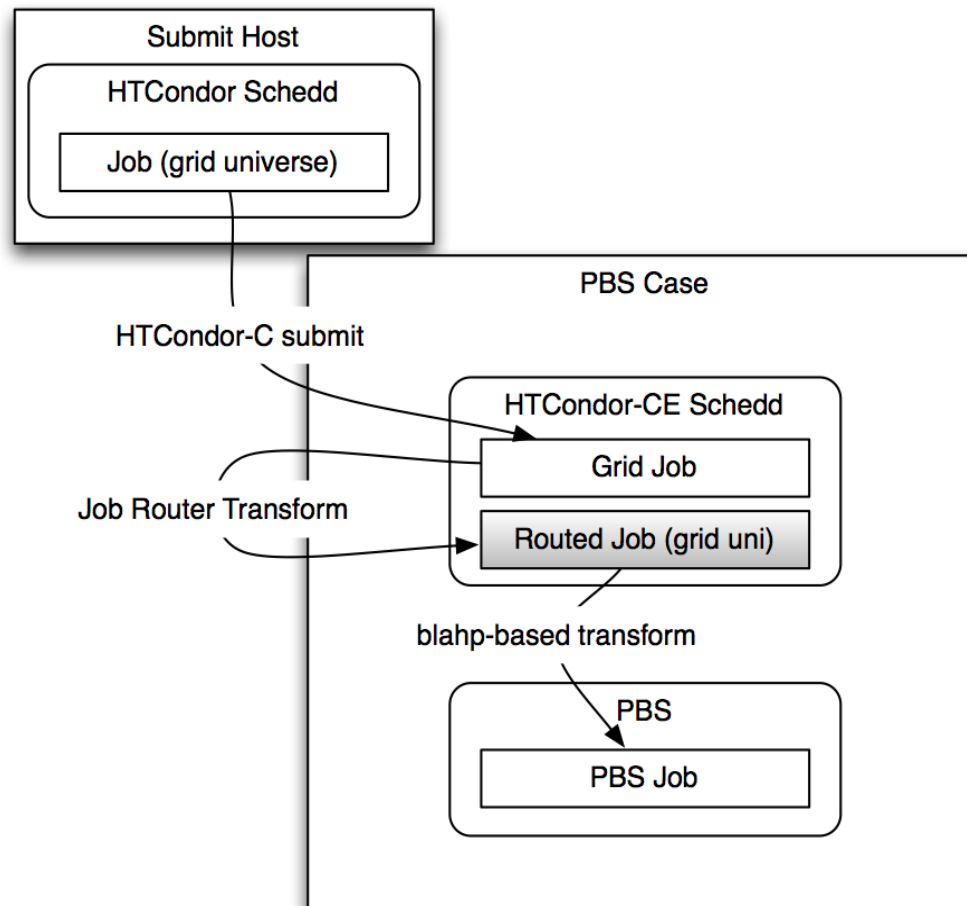  - Submit workflow for the HTCondor CE running on the site with the:

B.Lin, M.Zvada – HTCondor as the Job Gateway: HTCondor CE UWISC, UNL, OSG

# HTCondor CE: How it works?

- HTCondor batch system

# HTCondor CE: How it works?

- PBS batch system

# HTCondor CE: always room for improvement

- **Harden and Scale**
  - BLAHP
    - Improved file cleanup
    - Better error messages on failure
    - Handle errors more gracefully
  - HTCondor-C

- **Security Audit**
  - Record actions by the user that affect the job queue
    - Submission, removal, modification
  - Record how the user was authenticated
  - Record job credential files
  - Time-based rotation

B.Lin, M.Zvada – HTCondor as the Job Gateway: HTCondor CE

# HTCondor CE: Job Router (I)

- A key technology is the Job Router, which creates a copy of the job and transforms it according to a set of rules.

# HTCondor CE: Job Router (I)

- <span style="color:red">A key technology is the Job Router, which creates a copy of the job and transforms it according to a set of rules. In other words:</span>

    - we use the <span style="color:red">condor_jobrouter</span> daemon for transforming the job for the local site.
    - this daemon creates a copy of the job and applies a set of admin-prescribed transformations aka routes.
        - These can either be done via a ClassAd policy (declarative way) or a script callout.
        - The site customizations will no longer be overwritten by RPM upgrades!
    - The JobRouter can create the job copy directly in a site schedd, doing the site batch system submission for HTCondor sites.

# HTCondor CE: JobRouter ClassAd Policy

```
JOB_ROUTER_ENTRIES = \
  [ \
    GridResource = "batch pbs"; \
    TargetUniverse = 9; \
    name = "Local_PBS_cms"; \
    default_queue = "cms"; \
    Requirements = target.x509UserProxyVOName =?= "cms"; \
  ] \
  [ \
    GridResource = "batch pbs"; \
    TargetUniverse = 9; \
    name = "Local_PBS_other"; \
    default_queue = "other"; \
    Requirements = target.x509UserProxyVOName =!= "cms"; \
  ]
```

- More details/recipes for the routes:
  https://twiki.grid.iu.edu/bin/view/Documentation/Release3/JobRouterRecipes

# HTCondor CE: Job Router (II)

- Previously (GRAM), job transformations were specified in an imperative language (perl). The JobRouter includes a "hook" which allows the sysadmin to specify a script in any language.
  - e.g. JobRouter script JOB_ROUTER_DEFAULTS (python)


- NEW PHILOSHOPHY
  - The pilot describes the resources it needs and the site implementation details are hidden by the JobRouter.
  - Sites have the option of exposing internal configurations, but we'd like to encourage VOs to get to "site-independent pilot submission" - only the endpoint name is different!
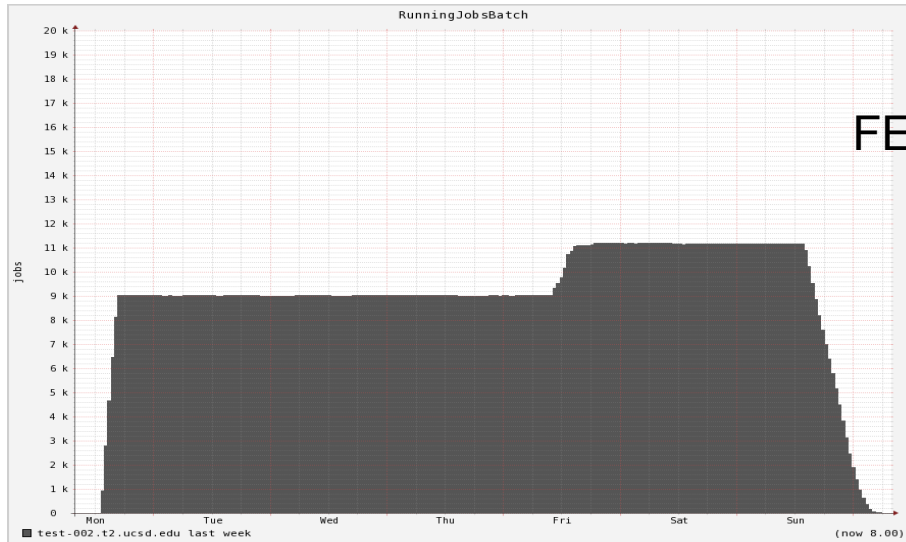
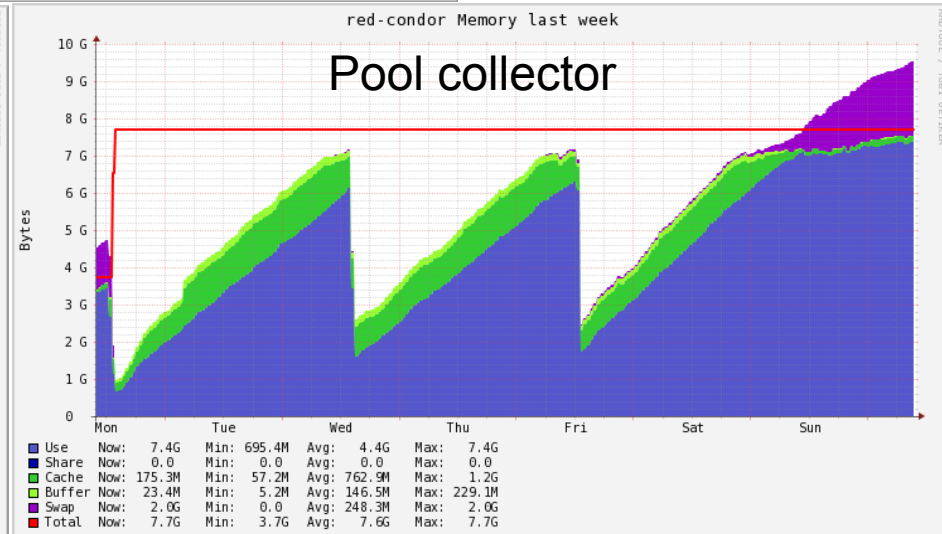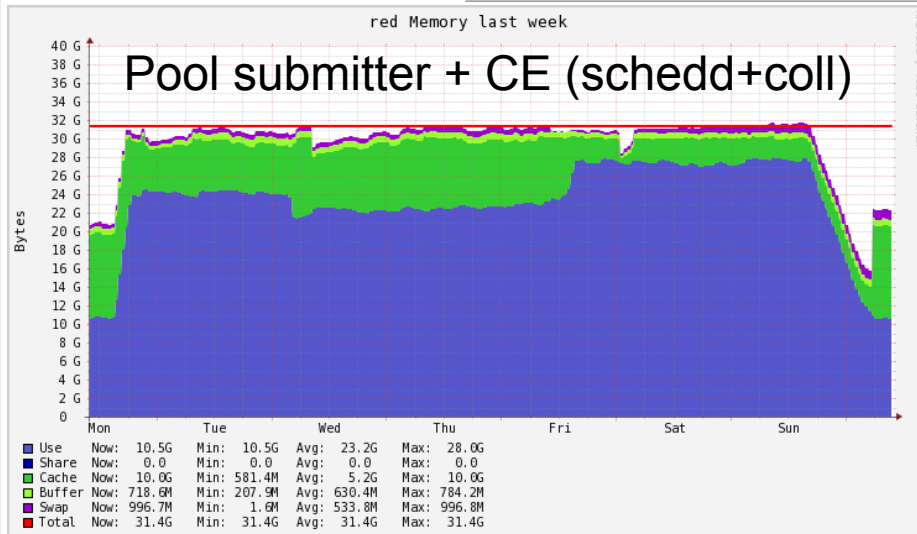# HTCondor CE: Hardware vs. Software limits?

- Well, yes :)

# HTCondor CE: Hardware vs. Software limits?

- <span style="color:red">Well, yes :)</span>

  - Depends on underlaying batch system you use – non-HTCondor sites might expect less system resource usage

  - Assuming your HTCondor cluster is well tuned you shouldn't meet any troubles, check it out:

    https://htcondor-wiki.cs.wisc.edu/index.cgi/wiki?p=LinuxTuning

  - let's have a look at some interesting plots anyway...

# HTCondor+HTCondor CE: HW vs. SW limits?



FE (user submit host)

Pool submitter + CE (schedd+coll)

Pool collector

https://htcondor-wiki.cs.wisc.edu/index.cgi/wiki?p=LinuxTuning
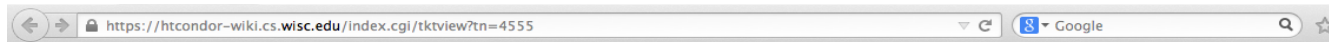
# HTCondor CE: not a flawless product

- OSG Technology and Software team tracks issues:

  - OSG JIRA open tickets: htcondor-ce component; mostly for configuration-related ticket and bugs not directly (sometimes) related to the HTCondor; testing and release promotion of new features

  - HTCondor project wiki: https://htcondor-wiki.cs.wisc.edu/index.cgi/tktview?tn=4555



**Ticket #4555: Parent tickets for HTCondor-CE items**

This is a parent ticket to help coordinate OSG issues with the HTCondor-CE.

Whenever possible we'd like to backport simple patches to v8.2 stable.

EMBED: Status Table of Derived Tickets

TN: 4555 [Go]

| Developer | # | Status | Type | Pri | Created | Changed | Target | Due Date | Days left | Title |
|---|---|---|---|---|---|---|---|---|---|---|
| johnkn | 4433 | resolved | enhance | 2 | 2014 Jul | Nov 17 | v080203 | | | Add syntax to allow users to append to configuration arrays |
| jfrey | 4720 | review | defect | 2 | Nov 13 | Nov 17 | | | | Files fail to spool using python bindings to an 8.0.x host |
| | 4719 | abandoned | defect | 2 | Nov 13 | Nov 17 | | | | JobRouter fails to write userlog due to bad permissions |
| jfrey | 3379 | new | defect | 2 | 2012 Dec | Nov 17 | | | | Race condition in file stageout |
| danb | 3072 | new | defect | 2 | 2012 Jun | Nov 17 | | | | JobRouter should forcex jobs that can't be removed |
| zmiller | 4557 | review | defect | 2 | Aug 25 | Nov 03 | v080204 | | | HTCondor-C should gracefully handle authentication issues |
| tannenba | 3056 | active | enhance | 2 | 2012 Jun | Oct 29 | v080302 | | | Schedd should automatically trigger a reschedule on job submission |
| jfrey | 4570 | resolved | enhance | 2 | Sep 04 | Sep 22 | v080203 | | | Allow override of hostname |
| johnkn | 4576 | abandoned | enhance | 2 | Sep 06 | Sep 17 | | | | JobRouter requires GridResource to be set |
| johnkn | 4569 | new | enhance | 3 | Sep 04 | 08:46 | v080309 | | | Debugging tools for JobRouter. |
| | 4683 | new | defect | 3 | Oct 28 | Oct 28 | | | | condor_status -schedd returns incorrect number of jobs in a CE |
| johnkn | 4590 | resolved | enhance | 3 | Sep 16 | Oct 10 | v080302 | | | Improve diagnostic output of JobRouter |
| tlmiller | 4599 | resolved | defect | 3 | Sep 18 | Oct 02 | v080204 | | | condor_router_q doesn't work in the HTCondor-CE / condor context |
| jfrey | 4573 | new | enhance | 3 | Sep 05 | Sep 29 | | | | Migrate HTCondor-CE classad functions to HTCondor proper |
| jfrey | 4574 | new | enhance | 3 | Sep 05 | Sep 17 | | | | Provide ClassAd functions for merging environments |
| tim | 4556 | docpending | defect | 4 | Aug 25 | Nov 05 | v080302 | | | UDP invalidations can be sent to hosts without a UDP command socket |
| jfrey | 4592 | resolved | defect | 4 | Sep 16 | Sep 30 | v080203 | | | Blahp lsf_status bug fixes |
| johnkn | 4609 | new | enhance | 4 | Sep 24 | Sep 24 | | | | Add 'analyze' functionality to condor_job_router_tool |

# HTCondor CE: Information services (CE Collector)

- Courtesy slide of Matyas

## Purpose of Information Services

- Clusters have machines that vary in power, policy, etc.

- Need to know about these differences to send jobs

- Send Glideins from a central Factory to shield users from this complexity

- Collect machine information in a central location; Factory can query it to determine where to send Glideins

- More technicalities presented by Matyas on Wednesday here

# HTCondor CE: Information services (CE Collector)

- Courtesy slide of BrianB

  - We have been doing information services wrong throughout the life of the grid.

    - Projects like the BDII have been an attempt to generalize the description of the state and queues of an LSF system.

      - You can generalize for other batch systems, but at the core, it is still optimized for a particular use case; works poorly for our needs.

  - The CE collector publishes a description of the HTCondor-CE, the resources accessible, and how to access the resources.

    - No monitoring information!  No hardcoded concept of the queue!

    - Currently, our one focus is to get core use case - providing the information necessary for provisioning systems - right on HTCondor-CE.  Additional use cases may follow.

Brian's slides here from yesterday...
...and on Tuesday "Upcoming improvements to the HTCondor-CE"

# HTCondor CE: WLCG world and SAMv3 job support

- sam_uri in OIM – htcondor://<your_ce_host_name>

B.Lin, M.Zvada – HTCondor as the Job Gateway: HTCondor CE
UWISC, UNL, OSG

# HTCondor CE: Troubleshooting Tools

- Diagnose communication problems

- Detailed diagnosis of failures
    - Can you connect to the server?
    - Can you authenticate with the server?
    - Are you authorized by the server?
        - …
        - …

- Troubleshooting data and list of tools
    https://twiki.opensciencegrid.org/bin/view/Documentation/Release3/TroubleshootingHTCondorCE

B.Lin, M.Zvada – HTCondor as the Job Gateway: HTCondor CE