

# Distributed Monitoring in the OSG

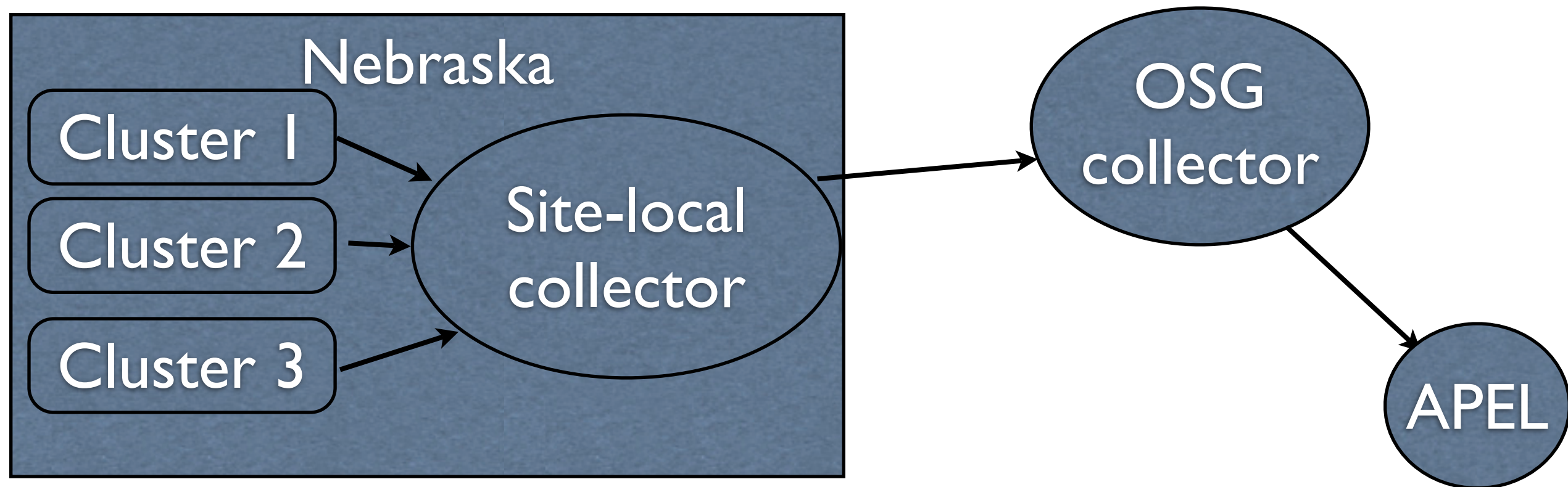
Brian Bockelman, GDB, 7 April 2009

Present the projects of many teams, including  
contributions from the OSG GOC, Metrics and  
Measurements, CMS, and ATLAS

# Data on the OSG

- Principle is that all data is available and can be displayed locally.
  - Accounting: Gratia.
    - Now with transfers!
  - Service Monitoring: RSV
  - Information services: CEMon & BDII
    - CEMon data can actually check many attributes
    - Historical BDII information is kept by the Metrics team; used widely for reports, experts-only web interface
  - Site registration: OIM
  - Troubleshooting: :(
  - Job-Level monitoring: :(

Gratia is not just an accounting system, but also includes a transport mechanism.



Looks a lot like the message brokers! Implementation (for better or worse) predates the maturity of open source solutions.

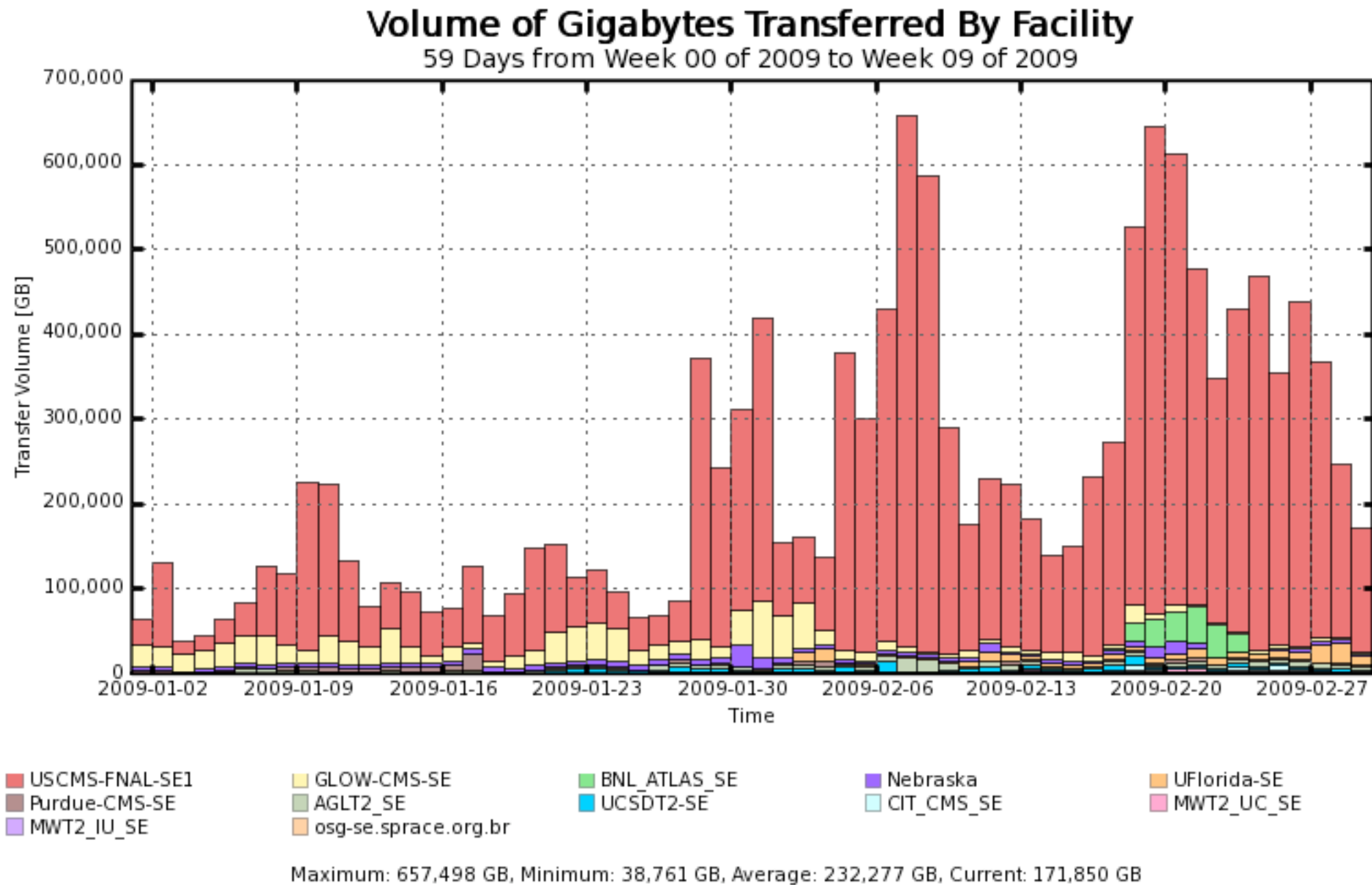
# Monitoring using Accounting Data

- Gratia is general-purpose (@ Nebraska, we use it to monitor our non-grid clusters), and you can choose what data you forward.
- The Gratia transport mechanism forwards XML messages through a network of “collectors” -- somewhat unique twist is that the clients persist messages until they are successfully sent to the next broker.

# Transfer Accounting

- Pieces required for transfer accounting have been coming together in the last three months.
  - Displays for folks to see
  - Probes for all OSG-supported systems
  - Critical mass of sites wanting to show up in the graphs.

# Transfer Accounting



Note: don't use this to compare sites, for demonstration purposes only (specifically, BNL's probe is dealing with issues and only has reported a handful of data ... but wow, FNAL moves a lot of data; Feb 19 averaged ~55Gbps)

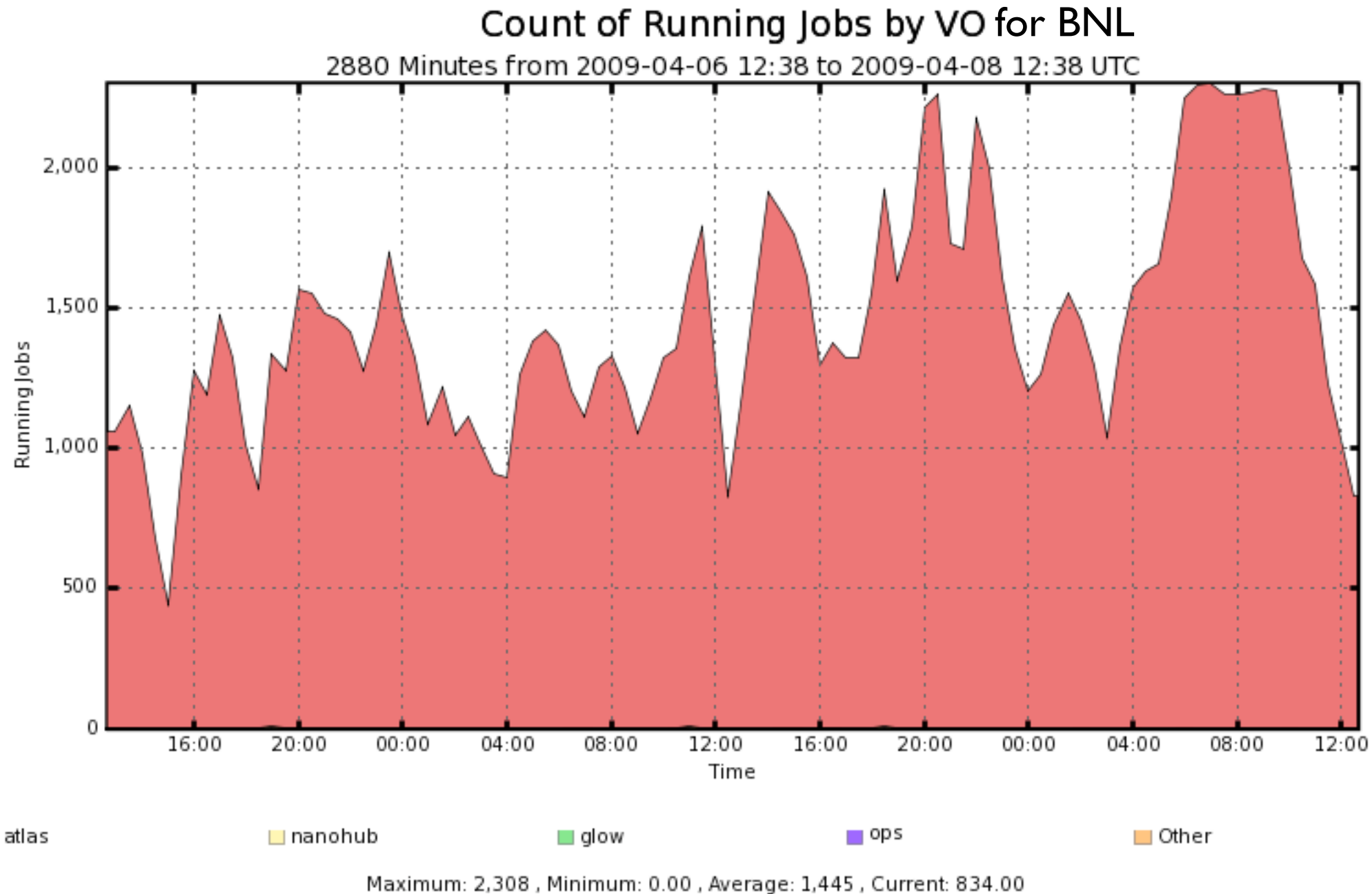


# Service Monitoring

- Based on probes which run at the site through a Condor-based cron-like mechanism.
- Because RSV uses the Gratia transport infrastructure, we can forward many site level probes to a site collector, then forward the site collector to OSG and the OSG data to WLCG.
- Next RSV version will give the option of running RSV probes in Nagios.
- Similar to the proposed EGEE architecture!



# Information Services



Like gstat; based on BDII information and can also apply (OSG) filters and interpretations. It's now being integrated into the MyOSG pages

# Panda & GlideInWMS

- Both GlideInWMS (used by CDF & CMS) and Panda (ATLAS) are developing sophisticated job-level monitoring tools to complement their submission tools.
- This is always complemented by the Dashboard project @ CERN.
- Really, quite lovely.

# Panda Snapshot

Panda Based Distributed Analysis Dashboard - Mozilla Firefox

File Edit View History Bookmarks Tools Help

http://panda.cern.ch:25880/server/pandamon/query?dash=analysis

Most Visited Getting Started Latest Headlines Customize Links Windows Marketplace Windows Media Windows

[CERN monitor](#) [Production](#) [Clouds](#) [DDM](#) [PandaMover](#) [AutoPilot](#) [Sites](#) [Analysis](#) [Physics data](#) [Usage](#) [Plots](#) [ProdDash](#) [DDMDash](#)

16 min old [Update](#) Not logged in. [List users](#)

## Panda monitor

Times are in UTC

[Panda info and help](#)

### Panda Based Distributed Analysis Dashboard

Information and tools for distributed analysis with Panda

**Jobs - [search](#)**  
Recent [running](#), [activated](#), [waiting](#), [assigned](#), [defined](#), [finished](#), [failed](#) jobs  
Select [analysis](#), [prod](#), [install](#), [test](#) jobs

**Quick search**  
Job   
Dataset   
Task request   
Task status   
File

**Summaries**  
Blocks:  days  
Errors:  days  
Nodes:  days  
[Daily usage](#)

**Tasks - [search](#)**  
[Generic Task Req](#)  
[EvGen Task Req](#)  
[CTBSim Task Req](#)  
[Task list](#)  
[New Tag](#)  
[Bug Report](#)

**Datasets - [search](#)**  
[Dataset browser](#)  
[Aborted MC datasets](#)  
[Panda subscriptions](#)

**Datasets Distribution**  
[DDM Req](#)  
[Req list](#)  
[AODs](#)  
[EVNTs](#)  
[Conditions\\_DS](#)

**Documentation on user analysis with Panda:**  
[Distributed analysis on Panda - overview page](#)  
[Client tools for Panda analysis jobs](#)  
[pathena: how to submit athena analysis jobs](#)  
[prun: how to submit ROOT and general jobs](#)  
[pbook: bookkeeping for Panda analysis jobs](#)  
[psequencer: how to perform sequential jobs/operations](#)

**Frequently asked questions:**  
[Full FAQ](#)  
[How is job priority calculated?](#)

**Status of pathena analysis queues:** See the wiki page [PathenaAnalysisQueues](#)

**Analysis jobs:** [Listing of analysis jobs](#). To look up a particular Panda job by ID use the quick search at left or click a PandaID in the job listing.

**Analysis users:** [User list](#) (also linked at top right, or above if you've logged in) shows analysis usage, ordered by most recent. From there you can go to your page (you're on the list if you've run a Panda job); if you log in you'll get easier access to your page from a new menu at the top of the page.

**Groups:** [Groups](#) are supported to organize users by role, physics working groups etc. and support collaborative work, accounting rights etc. (Not much used yet.)

**Data access:** See the [physics data](#) page linked above for information on data location, requesting replication of data, and staging data from tape to disk.

#### Analysis Summary By Cloud

World Wide - analy\_running - day

#### Analysis Summary By Site

US - running - day



# GlideInWMS Snapshot

UCSD Farm Job Monitor

https://glidein-mon.t2.ucsd.edu/jobmon/ucsd/

Identity: /DC=org/DC=doegrids/OU=People/CN=Brian Bockelman 504307#[cms]

JobId

Job Info

Local JobID	131362.0
Status	Running
Local User	uscms1483 (cms)
Queue	cms
Submitted at	Mon Apr 6, 2009 19:57:55 hrs
Started at	Mon Apr 6, 2009 19:58:40 hrs
Finished at	n/a
Exit Status	n/a
CPU Time	23:48 hrs
Wall Time	24:10 hrs
Execution Host	slot8@cabinet-4-4-4

CPU Efficiency/Memory

The top graph, titled 'CPU Efficiency for 131362.0@slot8@cabinet-4-4-4', plots CPU Efficiency (0.0 to 1.4) against Time (20:00 to 13:10). The efficiency starts at 0.0, rises to approximately 0.8 by 20:30, and then fluctuates between 0.8 and 1.0 for the remainder of the run.

The bottom graph, titled 'Memory/Disk usage for 131362.0@slot8@cabinet-4-4-4', plots Usage (MB) (0 to 1500) against Time (20:00 to 13:10). Memory usage (orange line) rises sharply to about 1100 MB by 20:30 and remains stable. Disk usage (blue line) stays near 0 MB until about 12:00, then rises to approximately 200 MB.

Admin Job ps WN top Work Dir Job Dir Job Output Job Error

Grid JobID	osg-gw-2.t2.ucsd.edu_131362.0_1239073075
CE ID	osg-gw-2.t2.ucsd.edu/jobmanager-condor-cms
RB/WMS	T2_US_UCSD
Subject	/C=UK/O=eScience/OU=Bristol/L=IS/CN=james jackson
Proxy Validity	60:12 hrs
Role	group_cms.uscms1483

Show Grid ID 23 Entries

State Selection Diagnosis Query Builder

Running

# Things we wish we had

- We wish we had better troubleshooting tools.
  - Yes, this is part of distributed monitoring.
  - Currently envisioned architecture would be a central site log collector for all OSG-supported applications.
    - Working with Globus, BestMan, dCache, and others to unify logging formats for machine analysis of logging data.
  - If we can centrally log data at sites, then sites can quickly give experts access to information they need.
    - Reduces the “send expert email”, “read response and send log snippet”, “read second response and send log snippet” cycle.
- Work is nascent, but we’re hopeful!
  - We probably won’t be centralize logging at the grid level.

# Thoughts

- On the OSG, we try to keep as few central services as possible - this is why our accounting was distributed from the design stages.
- RSV continued this line of thinking for service monitoring - and we're headed toward better site Nagios integration.
- Job-level monitoring is highly desired, but not delivered by the middleware - VOs have done a good job in delivering specialized solutions. What can we generalize?
- It appears that troubleshooting and logging could go the same route as service monitoring in using a distributed infrastructure.
  - But this is just in the "big idea" stage -- nothing concrete as of today.