



# Recent Developments in Panda

---

Torre Wenaus

Thanks to Tadashi, Paul, Alden, Jose, Kaushik, Sergey

US ATLAS Tier2/Tier3 Workshop

UTA

November 10, 2009



# Checked off the April Priority List...

September 1, 2009

- Done - finish Oracle migration
- Done - adapt AKTR+monitor to same Oracle interface code as server/bamboo
- Done - use bind variables in monitor
- Done - pilot release-candidate testing framework
- Done - pilot code refactoring to introduce glexec with file management either at prod-proxy or user-proxy level
- Done - http based scheduler interaction with DB, no direct access
- Done - analysis expansion to all clouds
- Done - system for establishing and maintaining site/cloud membership with associated rights
- Done - dataset browser overhaul. New browser deployed
- Done - Amazon cloud as Panda server platform

# The August (& Current) Priority List



- Done - New cache(+release) brokerage scheme
  - New swreleases table with releases+caches available at sites
- Done - Activate checksum in PandaMover
- Partly done - New job stats system with site, user analysis diagnostics
- Pending - Migration of Panda server to 3 dedicated machines
- Pending - PandaMover monitoring
- Pending - Event/alarm logging in monitor (eg. logging site state changes)
- Pending - New schedconfig management system
- Pending - autopyfactory: pyfactory + autopilot's DB-based config/monitoring



# 'In the Wings'

- Will move on these when there's a need...
  - glexec usage. Capability is now there in the production pilot to use glexec for analysis jobs.
    - Will deploy when there's a place supporting it and wanting it
  - Validation/deployment of dq2 based site mover
    - Highly desirable to get this deployed and used. Requires attention from very busy people (pilot, DQ2ers)
  - Panda support for athena running in multi-core mode
    - Not difficult to implement in Panda, mainly athena-level and batch-level issues. Can pursue once athena and site(s) support it
  - pcache based scheme for WN-level data brokerage
    - Scheme devised in August; see Charles' talk
  - Realignment of monitoring effort (AGIS, Dashboard, GWT)
    - Promising discussions and plans in August; fairly quiet since



# In the Future: Panda in the Cloud

- To date:
  - Panda server/monitor setup created and maintained in Amazon cloud (Sergey)
  - Proof-of-principle VM-based Panda pilot processing real (evgen) workloads (CERNVM team, ~1 year ago)
- New:
  - CERNVM team beginning a renewed effort in making VMs a real supported environment for Panda pilots
  - Doing the work themselves, asking for support/consultation from us
    - Which we gladly give (Jose Caballero will be first point of contact/support -- current status: contact established)
- Objective: Have VMs/clouds as an available capability
  - Such that it is a matter of policy, not capability, how/whether we use the cloud
- cf. recently announced large DOE investment in VM clusters at ANL/LBNL

# Steady State Activities



- Always, steady state activities driven by operations and users consume the most time
  - ongoing brokerage tuning
  - ongoing analysis functionality extensions
  - ongoing pilot changes to accommodate site, athena, middleware, ... issues
  - ongoing monitoring improvements



# Panda Server

- Migration of Panda servers to new dedicated machines progressing (with hiccups) but not complete
  - Isolate them from the Panda monitors
  - Current showstopper is `mod_gridsite` which doesn't work on SL5
- CERN single sign-on capability added
  - Can now leverage CERN authentication in Panda
- Task brokerage speeded up by increasing number of threads
  - We were running brokerage with only one thread up to now
  - Lots of scaling headroom by increasing the thread count
- ...and the usual brokerage tuning

# PandaMover



- Use of tape-staging priority to expedite throughput
- File status reported out of PandaMover into Panda file table to more clearly indicate missing files to ADC debugging
- Checksum verification enabled (required VDT upgrade to middleware supporting it)
- Priority introduced to handle reprocessing jobs faster than simulation





# Panda Based Analysis

- SLC5 support in Panda client & analysis jobs
- athena FileStager support
  - Athena will copy and process its input files using remote I/O from local SE
  - In Ganga domain, better performance than Ganga's standard posix-level remote I/O
- Output file transfer for analysis
  - Traditionally, Panda analysis involves no data movement
  - To support (particularly US) Tier 3s, there's now an exception
  - Panda can transfer analysis output files to Tier 3
    - using DQ2 as usual
  - Supports analysis workflow of Tier2 => Tier3 based analysis steps



# Panda Based Analysis (2)

- Rebrokerage for analysis jobs implemented and under test
  - Queues at popular sites are very long
  - While you wait, your data may become available at another quiet site, or the site where you're waiting may go down
  - User-triggered (or, in future, automated) rebrokerage asks the Panda server to redo its brokerage decision for your jobs
    - Assigning to the optimal site at the current moment



# Job Statistics Monitor

- New monitor gathering statistics on usage, timing, job mix for performance analysis
  - Without overtaxing the job DB
- Binned, averaged/summed data, 4hr bins for recent history (~week), 24hr bins kept indefinitely
- Provide interface, data for other analyses
  - text mode coming
- Consolidate and clean up related monitor functions
  - Pilot rates, user stats

[Configuration](#)

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

**[Panda monitor](#)**  
Times are in UTC

## Quick guide to the Panda monitor

[Panda info and help](#)

For Panda documentation and information on support and problem reporting see the [Panda info and user support](#) page.

**Monitor instances**

**[CERN](#): Primary production monitor at CERN**

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



## Panda statistics on site usage, performance and analysis activity

First set up the filters you want, then show the data.

*Current filters (click to remove):* [cloud:US](#) [sitetype:analysis](#)

[Clear all filters](#)

**Available filters:** (you can add them cumulatively, and multiple values of the same filter param are ORed)

Filter on cloud: [CA](#) [CERN](#) [DE](#) [ES](#) [FR](#) [IT](#) [ND](#) [NL](#) [TW](#) [UK](#) [US](#)

Filter on site type: [analysis](#) [production](#) [software](#)

Filter on site: [ANALY\\_AGLT2](#) [ANALY\\_BNL\\_ATLAS\\_1](#) [ANALY\\_BNL\\_LOCAL](#) [ANALY\\_GLOW-ATLAS](#) [ANALY\\_HU\\_ATLAS\\_Tier2](#)  
[ANALY\\_IllinoisHEP](#) [ANALY\\_LONG\\_BNL\\_ATLAS](#) [ANALY\\_LONG\\_BNL\\_LOCAL](#) [ANALY\\_MWT2](#) [ANALY\\_NET2](#) [ANALY\\_OU\\_OCHEP\\_SWT2](#)  
[ANALY\\_SLAC](#) [ANALY\\_SWT2\\_CPB](#) [ANALY\\_UTA](#)

Filter on job type: [ddm](#) [managed](#) [panda](#) [user](#) [software](#)

Filter on processing type: [digit](#) [pile](#) [pathena](#) [pandamover](#) [stresstest](#) [reco](#) [hammercloud](#) [merge](#) [reprocessing](#) [evgen](#) [simul](#) [ganga](#) [prun](#)  
[validation](#)

Detail level of tabulated data: [summary](#) [details](#) [all](#)

Time bin widths (hrs) to display if available (default is all): [4](#) [24](#)

---

Show data:

[Today](#)

[Yesterday](#)

[48 hours](#)

[72 hours](#)

[1 week](#)



## Panda usage stats

[Click for help](#)

[Go to Panda statistics dashboard](#)

Current filters (click to remove): [cloud:US](#)

[Clear all filters](#)

Summary of this selection: (Ordering is by job count, with (finished/failed) counts shown. Click on item to add to filter.)

Time interval: 2009-11-10 00:00:00 to 2009-11-10 14:10:00

cloud: [US \(40363/10805\)](#)

site: [ANALY\\_BNL\\_ATLAS\\_1 \(9324/4485\)](#) [ANALY\\_MWT2 \(9492/1264\)](#) [ANALY\\_LONG\\_BNL\\_ATLAS \(4455/4668\)](#) [BNL\\_ATLAS\\_1 \(4288/162\)](#) [AGLT2 \(3454/5\)](#) [MWT2\\_UC \(1435/44\)](#) [SLACXRD \(1187/152\)](#) [HU\\_ATLAS\\_Tier2 \(1325/3\)](#) [BU\\_ATLAS\\_Tier2o \(1222/8\)](#) [SWT2\\_CPB \(1149/1\)](#) [MWT2\\_IJ \(814/0\)](#) [BNL\\_ATLAS\\_ODM \(739/4\)](#) [ANALY\\_SWT2\\_CPB \(617/0\)](#) [UTD-HEP \(290/2\)](#) [OU\\_OCNEP\\_SWT2 \(208/2\)](#) [IU\\_OSG \(135/0\)](#) [IllinoisHEP \(107/0\)](#) [ANALY\\_SLAC \(82/0\)](#) [ANALY\\_AGLT2 \(18/3\)](#) [ANALY\\_IllinoisHEP \(16/0\)](#) [ANALY\\_NET2 \(6/2\)](#)

jobtype: [user \(12304/3853\)](#) [managed \(15614/379\)](#) [ddm \(739/4\)](#) [panda \(200/9\)](#)

proctype: [simul \(14513/375\)](#) [pathena \(10651/3592\)](#) [prun \(880/191\)](#) [pandamover \(739/4\)](#) [ganga \(125/548\)](#) [reprocessing \(36/0\)](#) [pile \(0/2\)](#)

username: [karsevan \(14513/377\)](#) [Dawe \(414/0\)](#) [Annovi \(283/0\)](#) [Kuehn \(89/0\)](#) [boyd \(36/0\)](#) [Strandberg \(25/0\)](#) [Boelaert \(21/0\)](#) [Elmsheuser \(15/4\)](#)

release: [15.3.1 \(14513/375\)](#) [None \(739/4\)](#) [15.5.1 \(47/395\)](#) [15.5.0 \(419/4\)](#) [15.0.0 \(283/0\)](#) [14.5.1 \(89/0\)](#) [15.1.0 \(46/0\)](#) [15.3.1 15.3.1 \(0/2\)](#)

time bin widths (hrs): [4](#)

### Job data:

Detail level: [summary](#) [details](#) [all](#)

Cloud: Site	Time	Hrs	JobType	ProcType	Release	User	Nuser	Nworker	Njob	Nget	Nupdate	Nfinish	Nfail	Nfail app	Nfail sys	Nfail dot	Nfail time	Eff	Utiliz	Twait min	TfullJob min	Tget sec	Tin sec	Trun min	Tout sec
US: ANALY_AGLT2	20091110-10	4	panda(2) user(1)	ganga(2) pathena(1)	15.5.0(2) 15.1.0(1)	Elmsheuser(2) Boelaert(1)	2	2	3	0	0	2	1	0	0	0	0	0	0	12	7	36	20	23	17
US: ANALY_AGLT2	20091110-02	4	panda(2) user(2)	pathena(2) ganga(2)	15.5.1(2) 15.1.0(2)	Elmsheuser(2) Boelaert(2)	2	4	4	0	0	4	0	0	0	0	0	0	0	26	11	78	19	9	28
US: ANALY_BNL_ATLAS_1	20091110-10	4	user(1020) panda(92)	pathena(745) prun(265) ganga(102)	15.5.1(215) 15.4.0(205) 14.5.0(201) 14.5.1(134) 15.5.2(107) 15.3.1(97)...	Mete(205) Liu(141) Lai(133) Hasegawa(114) BEHERA(106) Bitenc(101)...	13	43	1112	0	0	920	192	0	0	0	0	0	0	11	7	23	198	19	14
US: ANALY_BNL_ATLAS_1	20091110-06	4	user(1188) panda(22)	pathena(1146) prun(63) ganga(1)	15.5.0(516) 15.5.2(316) 14.2.25(228) 14.5.0(90) 15.5.1(60)	Cheung(453) BEHERA(314) Canelli(228) Liu(90) Dawe(63) Hasegawa(59)...	8	47	1210	0	0	709	501	0	0	0	0	0	0	78	51	28	254	22	24
US: ANALY_BNL_ATLAS_1	20091110-02	4	user(2257) panda(24)	pathena(1676) ganga(396) prun(209)	15.5.1(883) 15.1.0(399) 15.5.0(333) 14.5.0(211) 14.2.25(201) 14.5.1(113)...	Tabrizi(399) Goy(395) Rezvani(260) Liu(211) Dawe(209) Melachrinou(201)...	16	49	2281	0	0	1479	802	0	0	0	0	0	0	67	14	18	173	13	11
US: ANALY_LONG_BNL_ATLAS	20091110-10	4	user(1280) panda(10)	pathena(1101) ganga(156) prun(33)	14.5.1(684) 15.3.1(292) 15.5.1(141) 15.0.0(115) 15.6.0(16)	Kuehn(677) Sbrizzi(287) robinson(115) Schouten(78) Kleine-Limberg(35)	18	104	1290	0	0	503	787	0	0	0	0	0	0	268	14	5	143	218	9

# Panda usage stats



[Click for help](#)  
[Go to Panda statistics dashboard](#)

Current filters (click to remove): [cloud:US](#) [detail level:summary](#) [sitetype:analysis](#) [timebin:24](#)  
[Clear all filters](#)

Summary of this selection: (Ordering is by job count, with (finished/failed) counts shown. Click on item to add to filter.)

Time interval: 2009-11-03 00:00:00 to 2009-11-10 00:00:00

cloud: US (78675/65372)

site: [ANALY\\_BNL\\_ATLAS\\_1 \(37633/23783\)](#) [ANALY\\_LONG\\_BNL\\_ATLAS \(10080/35603\)](#) [ANALY\\_MWT2 \(18898/3726\)](#) [ANALY\\_IllinoisHEP \(6/0\)](#)

jobtype:  
 proctype:  
 username:  
 release:  
 time bin widths (hrs): [24](#)

Job data:  
 Detail level: [summary](#) [details](#)

Cloud: Site	Time	Twait min	TfullJob min	Tget sec	Tin sec	Trun min	Tout sec
US: ANALY_BNL_ATLAS_1	20091105-12	1393	66	14	264	19	11
US: ANALY_SLAC	20091108-12	1347	22	5	32	19	17
US: ANALY_LONG_BNL_ATLAS	20091106-12	1306	68	3	203	264	18
US: ANALY_LONG_BNL_ATLAS	20091105-12	1271	1271	3	230	251	49
US: ANALY_BNL_ATLAS_1	20091103-12	503	22	6	19	19	15
US: ANALY_NET2	20091103-12	480	6	3	187	25	13
US: ANALY_BNL_ATLAS_1	20091104-12	409	68	13	153	21	13

Twait min	TfullJob min	Tget sec	Tin sec	Trun min	Tout sec
64	30	14	1389	5	5
44	9	35	646	20	53
45	5	31	551	273	44
13	7	53	426	57	71
1393	66	14	264	19	11
103	11	7	248	7	28



Current filters (click to remove): [cloud:US](#) [detail\\_level:details](#) [jobtype:user](#) [site:ANALY\\_MWT2](#)  
[Clear all filters](#)

Summary of this selection: (Ordering is by job count, with (finished/failed) counts shown. Click on item to add to filter.)

Time interval: 2009-11-10 00:00:00 to 2009-11-10 14:15:00

cloud: US (6328/752)

site: ANALY\_MWT2 (6328/752)

jobtype: user (6328/752)

proctype: pathena (4673/752)

username:

release:

time bin widths (hrs): 4

Job data:

Detail level: [summary](#) [details](#) [all](#)

Cloud: Site	Time	Hrs	JobType	ProcType	Release	User	Nuser	Nworker	Njob	Nfinish	Nfail	Nfail app	Nfail sys	Nfail dat	Nfail time	Eff	Utiliz	Twait min	TfullJob min	Tget sec	Tin sec	Trun min	Tout sec
US: ANALY_MWT2	20091110-10	4	user	pathena(1377) prun(278)	14.2.25(1240) 15.0.0(278) 15.5.1(137)	Feickert(657) Tuggle(389) Annovi(278) Shochet(193) Tompkins(137) Dunford(1)...	6	204	1655	1655	0	0	0	0	0	0	0	710	29	16	90	46	39
US: ANALY_MWT2	20091110-06	4	user	pathena	14.2.25(1897) 15.5.1(492)	Shochet(796) Tuggle(526) Tompkins(492) Feickert(356) Melachrinos(199) Boveia(14)...	7	221	2389	2131	258	0	0	0	0	0	0	649	45	8	103	51	36
US: ANALY_MWT2	20091110-10	4	user	pathena	14.2.25(1240) 15.5.1(137)	Feickert(657) Tuggle(389) Shochet(193) Tompkins(137) Dunford(1)	5	188	1377	1377	0	0	0	0	0	0	0	710	29	9	106	50	44
US: ANALY_MWT2	20091110-02	4	user	pathena	14.2.25(1333) 15.5.1(326)	Feickert(348) Tuggle(330) Tompkins(326) Melachrinos(201) Kapliy(165) Boveia(161)...	7	223	1659	1165	494	0	0	0	0	0	0	301	43	11	163	46	33

# Job Statistics Monitor ToDo



- Still in progress:
  - Error categorization (user/panda/application/data)
  - Dump mode for binned and job-level data
  - Integration over bins
  - ... adding and adjusting based on usage and feedback
    - feedback appreciated!



# schedconfig & pilotController



- schedconfig (and other related tables) contain the queue and site configurations for
  - ATLAS production, analysis, sw installation sites
  - non-ATLAS (so far) EGEE sites (everything in BDII)
  - OSG sites: ITB, sites not (yet) used by ATLAS
- Merges information from Tiers of Atlas, BDII, OSG info (soon MyOSG) with its own information
- pilotController populates and maintains the schedconfig tables
  - ensures (so long as all changes go into pilotController) that errors can be backed off and DB restored

# pilotController Issues



- pilotController has grown ‘organically’ (bad) from the times when the Panda queue set was much smaller and was defined through several configuration files
- Difficult to track down changes to the DB through complex and confusing code
- Backup and recovery only through svn maintenance of source (a good mechanism iff all DB changes are made through pilotController)
- Adding queues and making changes hits the live DB; can be dangerous



# newController

- Under development by Alden
- First presentation at the August SW week, feedback/iteration
- Config changes done entirely outside the codebase, stored in a separate SVN for easy versioning and reversion
- DB backups both in the config files and stored as SQL code for quick restores
- Quick and easy queue creation in a text editor

# newController Timeframe



- Working code is in the process of being tested and run through the final stages of debugging
- Code will evolve slightly until the end of the year (refinements)
- As soon as it is deemed reliable, it will go live and replace pilotController

# Pilot Status

- Current pilot version: 40b
- 67 pilot updates so far in 2009 (108 in 2008)
  - A major release (*Na*) is often followed by several minor releases (*Nb-..*) with additional urgent new features or corrections to the previous release that can not wait until the next major release
- Each new pilot version is tested thoroughly first on a range of resources (both on production and analysis sites) using test jobs for different releases, and secondly on a massive scale (all sites) using the Pilot Testing Framework
  - NB. Still impossible to catch all types of errors, system extremely complex at this point. Any problem (pilot bug, trf or site incompatibility, ..) is usually corrected immediately

# Recent\* new pilot features (part 1)

## ■ SL5/SLC5 modifications

- The pilot has been successfully tested with release 15.3.1 jobs using cmtconfig i686-slc5-gcc43-opt at CERN and BNL on special queues
- For CERN and BNL tests the compiler setup was hardwired, awaiting new installation procedures for the releases
- Currently testing final setup at INFN-ROMA1 (additional pilot modifications were needed – no hardwiring of anything anymore, but AtlasSettings and AtlasLogin not installed at the ‘right’ place, new set of cmt tags)
  - Pilot modifications ready November 6, first successful test job:  
<http://panda.cern.ch:25980/server/pandamon/query?job=1027586197>

\* ) Since October

# Recent new pilot features (part 2)

- Direct access support for local and Castor site movers (LocalSiteMover, CastorSiteMover, castorSvcClassSiteMover, rfcplFCSiteMover)
  - Input files are accessed directly from the SE by athena instead of being handled by the pilot (root files only)
  - Direct access can be used with dCache, xrootd, RFIO (assuming correct schedconfig.copyssetupin)
- File stager support for testing
  - Input files are staged in and processed by athena one at a time
  - Successful tests at ANALY\_SLAC and some European sites

# Recent new pilot features (part 3)

- New LFC path convention added to the pilot
- Now using proper DQ2 site names from ToA instead of PanDA site names in DQ2 tracing report
- Various script and command verifications, and post-processing by the pilot
  - E.g. copysetup which can have complex format, can be updated by the pilot when file stager is used
  - Copysetup scripts are only used if they actually exists! Not always the case, e.g. Frontier-Squid scripts recently added to all EGEE sites
- ...and lots of code optimizations and re-factoring



# Pilot todo list... (part 1)

- Minor/major issues/requests are added more or less on a daily basis
- Final merging of current pilot version with glExec enabled pilot tested by Jose (Paul and Jose)
- Complete re-write of job recovery algorithm
  - Recent progress: JobRecovery class defined and half-way implemented. Much simplified as compared to old version. Constantly being delayed however due to the endless stream of new feature requests
  - Estimated time for code complete: 2 working days, testing not included

# Pilot todo list... (part 2)

- Ability to process several jobs sequentially by the pilot
  - Minor re-write of the pilots' main loop needed
- pCache integration / pilot based brokering
  - Pilot will scan the WN pCache folder and send the file list to the dispatcher when asking for a job (Jose)
- Maximum input file size on a per site basis
  - Pilot will use a new schedconfig field to allow certain tasks to use larger than normal maximum input file sizes on some sites (> 15 GB)



# Summary

- Panda stably functional
- But still much development driven by ops/needs
  - Tracking the environment changing around us
    - Especially in the pilot, our insulation layer from the real world of the grid
  - Improving/extending analysis capability
  - Scalability, performance testing/analysis/improvements
  - Security extensions
  - Functional extensions, new usage modes (esp. analysis)
    - eg. exploiting local disk in concert with pcache
  - New tools/technologies to leverage & support
    - VMs, cloud, many-core, ...