

# Assessment of Analysis Failures, DAST Perspectives

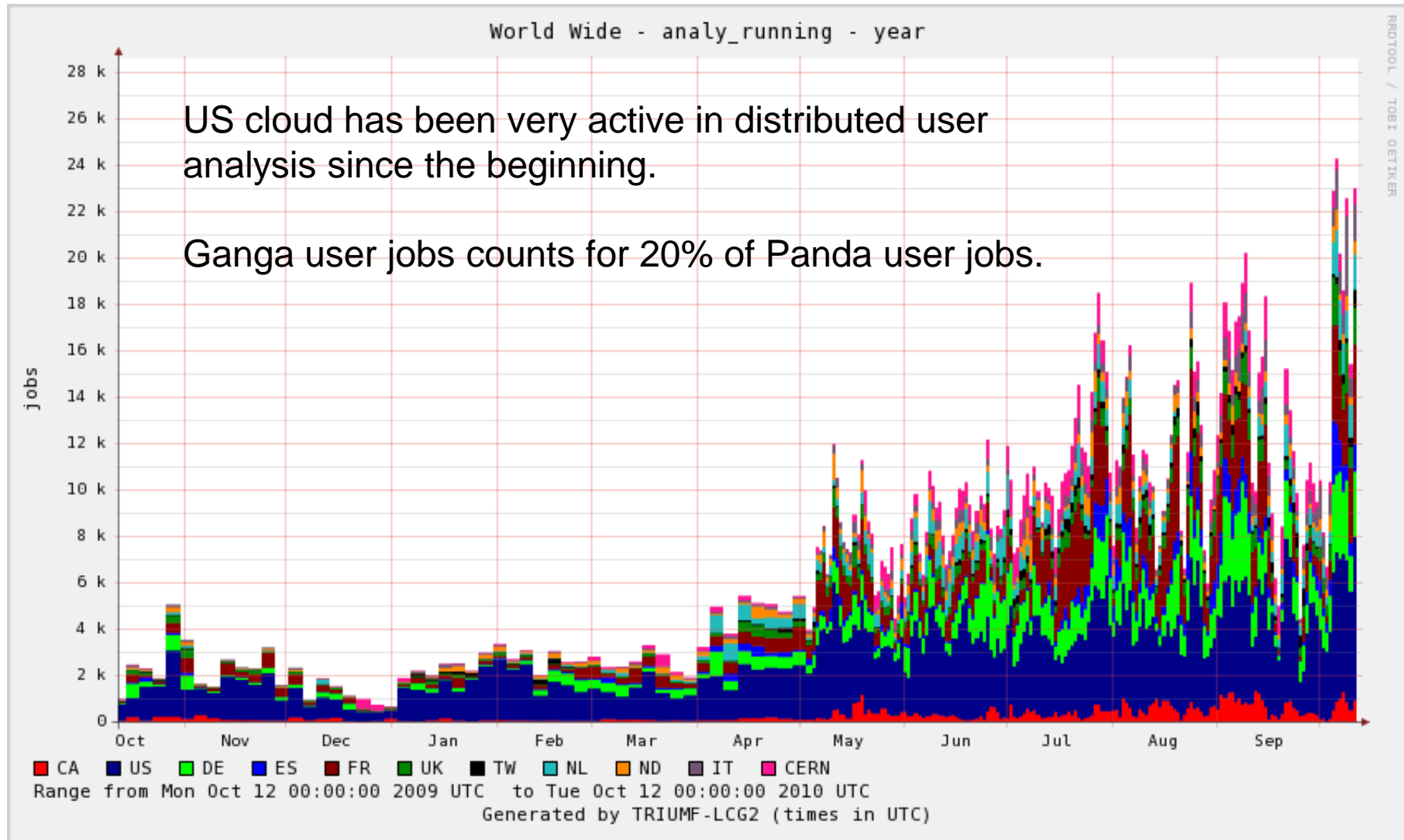
**Nurcan Ozturk**

**University of Texas at Arlington**

US ATLAS Distributed Facility Workshop at SLAC

12-13 October 2010

# Distributed Analysis Panda User Jobs in All Clouds

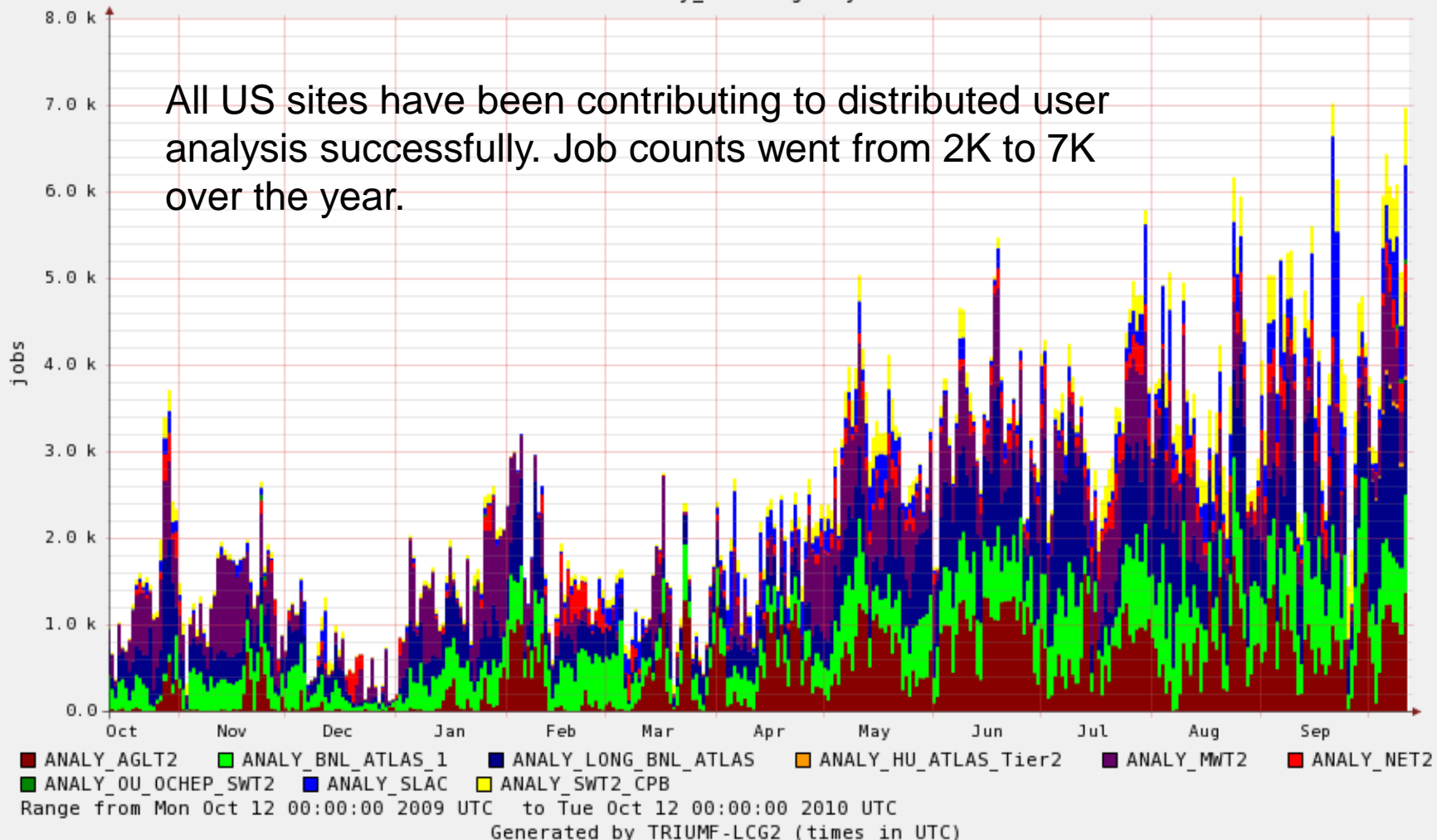


# Distributed Analysis Panda User Jobs in US Cloud



US - analy\_running - year

All US sites have been contributing to distributed user analysis successfully. Job counts went from 2K to 7K over the year.



# Analysis job summary in all clouds, last 12 hours (@12.50pm CERN time)



Cloud	Pilots	Latest	defined	assigned	waiting	activated	sent	running	holding	transferring	finished	failed	cancelled	%fail
<a href="#">ALL</a>			<a href="#">30989</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">63402</a>	<a href="#">4</a>	<a href="#">19407</a>	<a href="#">2107</a>	<a href="#">13</a>	<a href="#">157435</a>	<a href="#">26867</a>	<a href="#">28939</a>	15%
<a href="#">CA</a>	532	10-12 10:20	<a href="#">23</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">6645</a>	<a href="#">1</a>	<a href="#">793</a>	<a href="#">61</a>	<a href="#">0</a>	<a href="#">2685</a>	<a href="#">402</a>	<a href="#">11</a>	13%
<a href="#">CERN</a> (brokeroff)	1198	10-12 10:20	<a href="#">68</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1367</a>	<a href="#">0</a>	<a href="#">998</a>	<a href="#">143</a>	<a href="#">0</a>	<a href="#">4027</a>	<a href="#">288</a>	<a href="#">9</a>	7%
<a href="#">DE</a>	3547	10-12 10:20	<a href="#">7400</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">9446</a>	<a href="#">0</a>	<a href="#">2976</a>	<a href="#">543</a>	<a href="#">0</a>	<a href="#">28783</a>	<a href="#">2924</a>	<a href="#">2545</a>	9%
<a href="#">ES</a>	831	10-12 10:20	<a href="#">1186</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">3905</a>	<a href="#">0</a>	<a href="#">1192</a>	<a href="#">273</a>	<a href="#">0</a>	<a href="#">9625</a>	<a href="#">1975</a>	<a href="#">73</a>	17%
<a href="#">FR</a>	4883	10-12 10:20	<a href="#">4306</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">9805</a>	<a href="#">1</a>	<a href="#">5415</a>	<a href="#">184</a>	<a href="#">0</a>	<a href="#">16496</a>	<a href="#">3453</a>	<a href="#">2990</a>	17%
<a href="#">IT</a>	1123	10-12 10:20	<a href="#">3643</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">985</a>	<a href="#">0</a>	<a href="#">1235</a>	<a href="#">47</a>	<a href="#">0</a>	<a href="#">6593</a>	<a href="#">1689</a>	<a href="#">5703</a>	20%
<a href="#">ND</a>	194	10-12 10:20	<a href="#">6452</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1</a>	<a href="#">151</a>	<a href="#">144</a>	<a href="#">13</a>	<a href="#">3946</a>	<a href="#">1633</a>	<a href="#">49</a>	29%
<a href="#">NL</a>	1304	10-12 10:20	<a href="#">623</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">8067</a>	<a href="#">0</a>	<a href="#">1064</a>	<a href="#">83</a>	<a href="#">0</a>	<a href="#">6525</a>	<a href="#">5345</a>	<a href="#">1945</a>	45%
<a href="#">TW</a>	344	10-12 10:20	<a href="#">3</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">5849</a>	<a href="#">0</a>	<a href="#">99</a>	<a href="#">3</a>	<a href="#">0</a>	<a href="#">8599</a>	<a href="#">489</a>	<a href="#">2486</a>	5%
<a href="#">UK</a>	1556	10-12 10:20	<a href="#">1740</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">2930</a>	<a href="#">0</a>	<a href="#">1370</a>	<a href="#">231</a>	<a href="#">0</a>	<a href="#">8177</a>	<a href="#">1116</a>	<a href="#">1536</a>	12%
<a href="#">US</a>	2602	10-12 10:20	<a href="#">5545</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">14403</a>	<a href="#">1</a>	<a href="#">4114</a>	<a href="#">395</a>	<a href="#">0</a>	<a href="#">61979</a>	<a href="#">7553</a>	<a href="#">11592</a>	11%

## Analysis job error report, last 12 hours

Job wall time: 286395 hrs Error losses: trans: 12282 (4.3%) panda: 24156 (8.4%) ddm: 8623 (3.0%) athena: 24464 (8.5%) user: 202 (0.1%) other: 8338 (2.9%)

US cloud has been running the most jobs among all clouds. A stable cloud with low job failure rate. **Note:** BNL queues are being drained for the dCache intervention scheduled for today.

# Analysis job summary in US cloud, last 12 hours (@12.50pm CERN time)



US sites	Pilots	Latest	defined	assigned	waiting	activated	sent	running	holding	transferring	finished	failed	cancelled	%fail
<a href="#">ANALY_AGLT2</a> ✓	520	10-12 10:20	<a href="#">867</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">2995</a>	<a href="#">0</a>	<a href="#">1749</a>	<a href="#">94</a>	<a href="#">0</a>	<a href="#">7917</a>	<a href="#">11</a>	<a href="#">0</a>	0%
<a href="#">ANALY_ANLASC</a> (offline) ✓			<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_BNL_ATLAS_1</a> ✓	2	10-12 10:20	<a href="#">1459</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">30</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1</a>	<a href="#">0</a>	<a href="#">7152</a>	<a href="#">4950</a>	<a href="#">2151</a>	41%
<a href="#">ANALY_BNL_LOCAL</a> ✓	2	10-12 09:00	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_DUKE</a> (brokeroff) ✓	3	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_GLOW-ATLAS</a> (offline) ✓			<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_HU_ATLAS_Tier2</a> ✓	40	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">851</a>	<a href="#">0</a>	<a href="#">50</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">52</a>	<a href="#">203</a>	<a href="#">0</a>	80%
<a href="#">ANALY_IllinoisHEP</a> ✓	26	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">117</a>	<a href="#">0</a>	<a href="#">43</a>	<a href="#">4</a>	<a href="#">0</a>	<a href="#">279</a>	<a href="#">0</a>	<a href="#">500</a>	0%
<a href="#">ANALY_LONG_BNL_ATLAS</a> ✓	638	10-12 10:20	<a href="#">1326</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1530</a>	<a href="#">0</a>	<a href="#">752</a>	<a href="#">9</a>	<a href="#">0</a>	<a href="#">9781</a>	<a href="#">450</a>	<a href="#">316</a>	4%
<a href="#">ANALY_LONG_BNL_LOCAL</a> ✓	2	10-12 09:00	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_MWT2</a> ✓	430	10-12 10:20	<a href="#">506</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">3637</a>	<a href="#">0</a>	<a href="#">622</a>	<a href="#">80</a>	<a href="#">0</a>	<a href="#">19782</a>	<a href="#">607</a>	<a href="#">367</a>	3%
<a href="#">ANALY_MWT2_X</a> (brokeroff) ✓	17	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_NET2</a> ✓	190	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">2035</a>	<a href="#">0</a>	<a href="#">219</a>	<a href="#">43</a>	<a href="#">0</a>	<a href="#">2733</a>	<a href="#">113</a>	<a href="#">2</a>	4%
<a href="#">ANALY_OU_OCCEP_SWT2</a> ✓	47	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1055</a>	<a href="#">0</a>	<a href="#">44</a>	<a href="#">5</a>	<a href="#">0</a>	<a href="#">274</a>	<a href="#">0</a>	<a href="#">0</a>	0%
<a href="#">ANALY_SLAC</a> (offline) ✓	289	10-12 10:20	<a href="#">1086</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">1428</a>	<a href="#">0</a>	<a href="#">1</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">7159</a>	<a href="#">463</a>	<a href="#">4</a>	6%
<a href="#">ANALY_SLAC_LMEM</a> (offline) ✓	4	10-12 08:30	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_SWT2_CPB</a> ✓	363	10-12 10:20	<a href="#">300</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">725</a>	<a href="#">1</a>	<a href="#">634</a>	<a href="#">159</a>	<a href="#">0</a>	<a href="#">6850</a>	<a href="#">756</a>	<a href="#">8252</a>	10%
<a href="#">ANALY_Tufts_ATLAS_Tier3</a> (brokeroff) ✓	29	10-12 10:20	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_UTA_T3</a> (brokeroff) ✓			<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	
<a href="#">ANALY_WISC-ATLAS</a> (offline) ✓			<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	<a href="#">0</a>	

BNL queues are being drained for the dCache intervention scheduled for today.

# Quick Look at BNL and HU Site Failures



ANALY_BNL_ATLAS_1		defined:1459 assigned:0 waiting:0 activated:30 sent:0 running:0 holding:1 transferring:0 finished:7152 failed:4950 (40.9%)									
pilotErrorCode (4049)	2	1.9	10-12 02:02	<a href="#">1099</a>	Get error: Staging input file failed						
pilotErrorCode (4049)	1	1.7	10-12 00:05	<a href="#">1104</a>	User work directory too large						
pilotErrorCode (4049)	3969	1642.0	10-12 06:00	<a href="#">1165</a>	Put error: Local output file missing <-- User code problem						
pilotErrorCode (4049)	77	115.8	10-12 00:56	<a href="#">1169</a>	Put error: LFC registration failed						
transExitCode (901)	1	0.2	10-12 02:06	<a href="#">20</a>	Unknown error code						
transExitCode (901)	900	73.3	10-12 05:50	<a href="#">40</a>	Athena crash - consult log file <-- User code problem						

ANALY_HU_ATLAS_Tier2		defined:0 assigned:0 waiting:0 activated:851 sent:0 running:50 holding:0 transferring:0 finished:52 failed:198 (79.2%)									
pilotErrorCode (31)	27	33.8	10-11 23:32	<a href="#">1099</a>	Get error: Staging input file failed						
pilotErrorCode (31)	4	9.2	10-11 23:30	<a href="#">1151</a>	Get error: Input file staging timed out						
transExitCode (167)	167	100.9	10-12 01:04	<a href="#">40</a>	Athena crash - consult log file <-- User code problem						

# Closer Look to Analysis Job Error Report (1)



## Analysis job error report, last 12 hours

Job wall time: 279524 hrs Error losses: trans: 11597 (4.1%) panda: 24176 (8.6%) ddm: 9039 (3.2%) athena: 23437 (8.4%) user: 239 (0.1%) other: 8188 (2.9%)

- **trans:**
  - Unspecified error, consult log file
  - TRF\_SEGVIO - Segmentation violation
- **panda:**
  - Put error: Error in copying the file from job workdir to localSE
  - Lost heartbeat
  - No space left on local disk
  - Get error: Staging input file failed
  - Put error: Failed to add file size and checksum to LFC
  - Payload stdout file too big
  - Exception caught by pilot
  - Put error: Failed to import LFC python module
  - Bad replica entry returned by lfc\_getreplicas(): SFN not set in LFC for this guid
  - wget command failed to download trf
  - Missing installation
  - etc

“panda’ errors (pilotErrorCode) are relevant for the attention of the US Facility, however difficult to observe a pattern.

# Closer Look to Analysis Job Error Report (2)



- **ddm:**
  - Could not add output files to dataset
- **athena:**
  - Athena ran out of memory, Athena core dump
  - ATH\_FAILURE - Athena non-zero exit
  - Athena core dump or timeout, or conddb DB connect exception
  - Athena crash - consult log file
- **user:**
  - User work directory too large
- **other:**
  - Unknown error code

Note: above is my classification of errors, have not double checked with Torre as to how they are classified in the analysis job error report on Panda monitor.



# User problems reported to DAST for US sites

# DAST – Distributed Analysis Support Team



- DAST started in September 2008 for a combined support of pathena and Ganga users.
- First point of contact for distributed analysis questions.
- All kinds of problems are discussed in the DA help forum (hn-atlas-dist-analysis-help@cern.ch) not just pathena and Ganga related ones:
  - athena
  - physics analysis tools
  - conditions database access
  - site/service problems
  - dq2-\* tools
  - data access at sites
  - data replication
  - etc
- DAST helps directly by solving the problem or escalating to relevant experts

# Team Members



## EU time zone

## NA time zone

-----

Daniel van der Ster

Nurcan Ozturk (now in EU time zone)

Mark Slater

Alden Stradling

Hurng-Chun Lee

Sergey Panitkin

Bjorn Samset

Bill Edson

Christian Kummer

Wensheng Deng

Maria Shiyakova

Shuwei Ye

Jaroslava Schovancova

Nils Krumnack

Manoj Jha

Woo Chun Park

Elena Oliver Garcia

Karl Harrison

Frederic Brochu

Daniel Geerts

Carl Gwilliam

Mohamed Goughri

Borge Gjelsten

Katarina Pajchel

red: trainee

Eric Lancon

A small team in NA time zone, most people are the ones who are fully/partly supported by the USATLAS Physics Support and Computing Program, thus close ties to USATLAS Facilities and Support.

# User problems reported to DAST for US sites (1)



By looking at the recent DAST shift reports, here are some of the issues reported by users on US sites:

- A wrong update concerning the analysis cache 15.6.13.1.1. Affected several sites including US ones.
- Error at BNL: **pilot: Get error: lsm-get failed.** The pool hosts the input files was not available due to machine reboot. The pool is back shortly after.
- MWT2 site admin reporting user analysis failures (HammerCloud and user code problems). Very useful for heads-up and understanding the site performances better.
- Production cache 16.0.1.1 missing at MWT2. Installed by Xin, available at several analysis sites but not all, no indication of usage in physics analyses at: <http://atlas-computing.web.cern.ch/atlas-computing/projects/releases/status/>

# User problems reported to DAST for US sites (2)



- dq2-get error, 'lcg\_cp : Communication error on send' at MWT2 – Transient error.
- Why do Australian grid certs not work on BNL? dq2-get difficulties by some Australian users. Shifter has not escalated this issue to BNL yet
- Tier3 issue – User can't use WISC\_LOCALGROUP disk as a DaTRI source. DDM people reported that it is in the works.

# User File Access Pattern (1)



- Pedro examined the user access pattern at BNL. His observations as follows:
- Took a 1.5h sample, 29 files per second were accessed just for reading. 56% of the files are user library files. This suggested that users are splitting jobs into many subjobs (each subjob needs to access the library file).
- **Outcome:** Lots of reads and writes to the userdisk namespace database, the writes clog up and timeouts on the whole pnfs namespace. Heavy usage of the same file, possibility to break the storage server, the files will not be accessible since there is no redundancy (this is not a hotdisk area). In this case he estimated a failure of more than 2/sec and depending on the problem, the recovery can take 2min-15min (240-1800 job failures).
- **Solution:** If this continues to be the practice among users, he needs to add more changes on the local site mover and add some extra 'spices' to be able to scale the whole system.(except the data availability since for that lib files would need to go to HOTDISK by default or some new HOTUSERDISK token).

# User File Access Pattern (2)



- I looked at the user jobs Pedro gave as examples.

user	gregor	mdavie	csandova	zmeng
jobsetID	2354	3909	639	3879
# subjobs	4999	2295	3884	1390
#input file/job	4 ESD	1 D3PD	1 ESD	1 ESD
average run time/job	3h	3' to 37'	13' to 3h	14' to 50'
average input file size	800 MB	200 MB (some 6 MB)	3 GB	3 GB
average output file size	60 KB	300-900 KB	200 MB	130 MB

# User File Access Pattern (3)



- **Conclusion:** Users are splitting into many subjobs and their output file size is small, then they have difficulty to download them by dq2-get. As discussed in DAST help list last week, the file-look up can be as long as 30s in DDM system, so users have been advised to merge their output files at the source and then download.
- The default limit on the input file size is 14GB in the pilot. So these users could use to run on more input files in one job. They are probably splitting more to be on the safe side.
- **How could we advise users?** Can pathena print a message when the job is splitted into too many subjobs?
- **This needs to be further discussed.**



# Storage space alerts for SCRATCHDISKs



DAST receives notifications from DQ2 system. An example follows from today. NET2 is often in the list. Is there a problem with cleaning the SCRATCHDISK?

Date: Tue, 12 Oct 2010 10:06:25 +0200

From: ddmusr01@cern.ch

To: atlas-project-adc-operations-analysis-shifts@cern.ch,  
fernando.harald.barreiro.megino@cern.ch,  
alessandro.di.girolamo@cern.ch, dan@vanderster.com

Subject: [DQ2 notification] Storage space alerts for SCRATCHDISKs

Site	Free(TB)	Total(TB)
IN2P3-LPSC_SCRATCHDISK	0.359 (16%)	2.199
NET2_SCRATCHDISK	0.047 (8%)	0.537

For questions related to this alarm service: atlas-dq2-dev@cern.ch

# How US Cloud Provide Support to DAST



- Site issues are followed by GGUS tickets:
  - Good response from site admins
  - Some site admins also watch for issues reported in the DA help forum for their sites, thus prompt support.
- General questions:
  - [atlas-support-cloud-us@cern.ch](mailto:atlas-support-cloud-us@cern.ch)
  - Good response, Yuri often keeps an eye on this list.
- Release installation issues at US sites:
  - Xin Xhao
- Also private communications with Paul Nilsson, Torre Weanus on the pilot and Panda monitor issues (not only for US sites of course).

# Summary



- US cloud continues to be a successful cloud in ATLAS distributed user analysis.
- Good US site performances, one of the least problematic analysis sites in ATLAS.
- User analysis jobs fail with different errors. The dominant ones are classified as “panda” and “athena” on Panda monitor. “athena” errors often refer to user code problems. “panda” errors vary, difficult to observe a pattern or classify as common problem.
- Good US cloud support to DAST. No operational issues to report in the US Facility, in particular from the DAST point of view.
- User access pattern needs to be further discussed for better site/storage performances.
- I have not mentioned here the new Panda Dynamic Data Placement (PD2P) mechanism which helped greatly users (long waiting period in the queue) and sites (storage). I expect Kaushik will cover this topic.