



# Status of U.S. ATLAS Computing Facilities and Distributed Computing

Michael Ernst, BNL  
Brookhaven National Laboratory

U.S. ATLAS Distributed Facilities Meeting

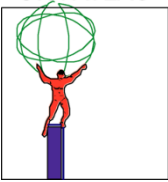
UNL

March 19, 2012



# Outline

- **U.S. ATLAS Facility Organization**
- **Facility Accomplishments and Milestones**
- **Open Science Grid Project Continuation**
- **Facility Contribution & Performance**
- **Facility Capacity Planning**
- **Tier-1 Facility Operations and Performance**
- **Facility Effort and Current Capacities**
- **Summary**



## U.S. ATLAS Facilities under PS&C

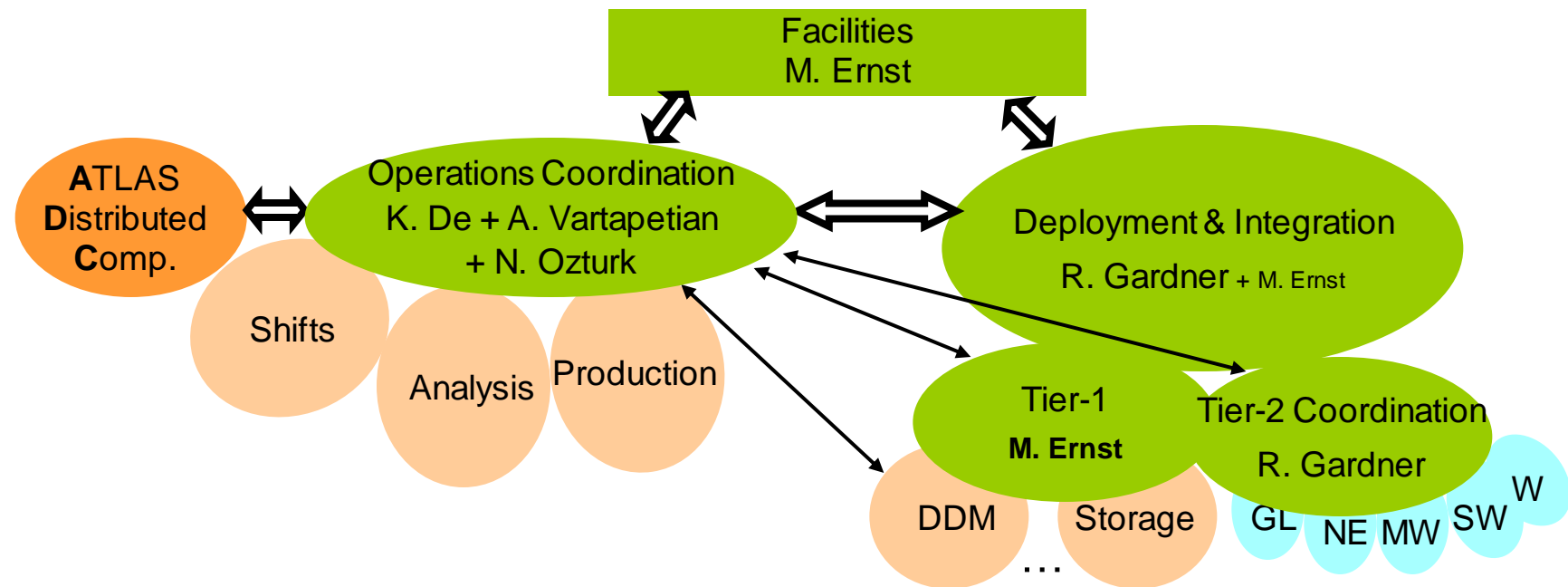
- 2.1, 2.9 Management (Wenaus/Willocq)
- 2.2 Software (Luehring)
  - 2.2.1 Coordination (Luehring)
  - 2.2.2 Core Services (Calafiura)
  - 2.2.3 Data Management (Malon)
  - 2.2.4 Distributed Software (Wenaus)
  - 2.2.5 Application Software (Neubauer)
  - 2.2.6 Infrastructure Support (Undrus)
  - 2.2.7 *Analysis support (retired; redundant)*
  - 2.2.8 Multicore Processing (Calafiura)
- 2.3 Facilities and Distributed Computing (Ernst)
  - 2.3.1 Tier 1 Facilities (Ernst)
  - 2.3.2 Tier 2 Facilities (Gardner)
  - 2.3.3 Wide Area Network (McKee)
  - 2.3.4 Grid Tools and Services (Gardner)
  - 2.3.5 Grid Production (De)
  - 2.3.6 Facility Integration (Gardner)
- 2.4 Analysis Support (Cochran/Yoshida)
  - 2.4.1 Physics/Performance Forums (Black)
  - 2.4.2 Analysis Tools (Cranmer)
  - 2.4.3 Analysis Support Centers (Ma)
  - 2.4.4 Documentation (Luehring)
  - 2.4.5 Tier 3 Coordination (Benjamin)

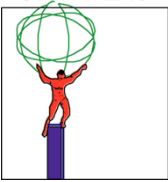


# Facilities Organization

Facility comprising two principal lines

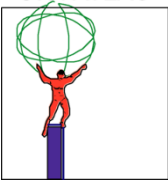
- ❑ Production and Analysis Operation Coordination
- ❑ Facility Deployment, Integration and Operation
- Organization proven to be very effective for >3 years





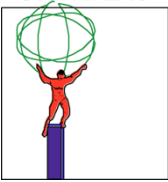
# Future Operations Coordination

- **Observation: US is losing momentum in ADC operations**
  - ◆ In the past had significant influence when Alexei operated DDM
  - ◆ Had significant influence when Kaushik ran Production
  - ◆ Has an impact on our effectiveness and contribution
    - Had several incidents where US cloud was impaired by central or local SW or infrastructure issues w/o corrective actions being taken in a timely fashion
      - Relying on experts outside the US to make sure our resources are adequately filled (CPU and disk) and tasks are running smoothly is likely to result in problems
  - ◆ **Production and Analysis Operations Coordination is part of the Facilities Program (WBS 2.3.5)**
    - There is a team of people working in this area, but their area of activities is either too US-centric or too ATLAS-centric
      - Lacking decent integration of regional and central responsibilities
      - E.g. while DDM is covered to some extent Production at a higher level than job failures is not covered at all
        - Facilities in the “dark” whenever there is a lack of jobs and/or analysis resources are underutilized
  - ◆ **Kaushik et al in the process of defining scope & responsibilities**
    - Crucial to have (pro-)active/agile people in these positions



# Accomplishments – Facilities and Distributed Computing

- **Largest Tier-1 center in ATLAS**
  - ◆ Excellent w.r.t. availability, data and CPU delivered, production and analysis performance
  - ◆ If interventions cannot be done in a transparent way we always schedule outages in the shadow of LHC/ATLAS-wide downtimes
- **U.S. Tier-2s also the best in the ATLAS Tier-2 complex**
  - ◆ 4 out of 5 U.S. Tier-2s in the top 10 (of ~75) sites that do 50% of ATLAS analysis work
- **Tier-3 deployment successfully completed and integrated**
  - ◆ Doug now focusing on distributed analysis support
- **U.S.-led project in the area of Federated Data Stores**
  - ◆ The Tier-1 and all Tier-2s participating
  - ◆ Have created a global namespace across the Facilities in the US
  - ◆ “Demonstrator” is up & running
    - **Functionality verified with ATLAS Analysis jobs (HammerCloud)**
- **Thanks to OSG, a crucial part of the success**
  - ◆ We have prominently contributed to a strong proposal that is funded by DOE & NSF from Apr 1, 2012 for the next 5 years



# Service Integration Activities

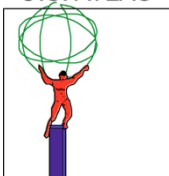
- Besides Facility operations at high performance and high reliability we are/have been working on several service integration activities
  - UIUC joining MWT2 (when done comprising 3 sites: UC, IU, UIUC)
  - ATLAS is moving from “shared area” (NFS, AFS) to virtualized global filesystem based on CVMFS (developed at and supported by CERN) for S/W distribution and Conditions Data access
    - CVMFS-based PanDA sites in production (T1, AGLT2, MWT2, NET2, WT2, SWT2?)
  - Federated Data stores w/ xrootd
    - Focused WG has made excellent progress
    - All tiers participating
    - ATLAS analysis jobs successfully access data through federation
    - Still some work remaining (authentication, performance, ...)
  - Cloud Computing
    - Developing concepts, infrastructure & components applicable to T1, T2 and T3
    - US participating in an ATLAS-wide program of work
  - AutoPyFactory – A new infrastructure for pilot submission
    - Deployment already in progress in the US
    - Pilot factory development & operations now under full control by US Facilities
  - WAN performance optimization and monitoring with perfSONAR-PS
    - Our initiative in collaboration w/ Internet2 has finally paid off
    - LHCOPN monitoring fully implemented at 10 T1s and CERN within only ~3 months, we are in the process of making it a service under WLCG
    - Baseline for LHCONE (LHC Open Network Environment) monitoring
    - Adopted by other regions



# **FY2012 Facility Milestones**

- **January: At the T1 AutoPilot based Pilot submission replaced by AutoPyFactory**
- **February: All U.S. ATLAS Facility Sites using CVMFS**
- **March: Infrastructure deployed to invoke Cloud resources**
- **April: 2012 Pledges installed at T1 and T2s**
- **May: Tier-1 and at least 3 Tier-2s connected to LHCONE**
- **May: UIUC/NCSA becomes integral part of MWT2**
- **June: Federated Data Stores at U.S. Sites in Production**
- **September: 100 GE infrastructure for R&D deployed at T1**





# US ATLAS Installed Capacities

## (Snapshot as of Jan 2012)

2012 Pledged vs Installed Capacities at the US ATLAS Facilities (as of Jan, 2012)

Site	CPU [HEPSpec 2006]		DISK [TB]	
	2011 Pledge	Installed Sep, 2011	2011 Pledge	Installed Sep, 2011
Tier-1	51,980	58,000	5,704	7,100
AGLT2	12,232	36,163	1,654	1,910
MWT2	12,232	36,840	1,654	1,302
NET2	12,232	19,035	1,654	1,100
SWT2	12,232	23,220	1,654	1,260
WT2	12,232	15,816	1,654	1,663
<b>Total</b>	<b>113,140</b>	<b>189,074</b>	<b>13,974</b>	<b>14,335</b>

2011 pledges were made available to ATLAS in April – June (some T2 disk)

2012 pledges are almost met

- CPU installed & operational
- Disk: 100% at T1, 88% at T2
- Tape installed & operational

2012 Pledged capacities at the US ATLAS Facilities

Site	CPU [HEPSpec 2006]		DISK [TB]	
	2012 Pledge	Installed Jan, 2012	2012 Pledge	Installed Jan, 2012
Tier-1	60,000	62,000	6,300	8,100
AGLT2	12,500	36,178	2,200	2,160
MWT2	12,500	36,840	2,200	2,730
NET2	12,500	19,679	1,648	1,200
SWT2	12,500	30,460	2,200	1,439
WT2	12,500	29,376	2,200	2,143
<b>Total</b>	<b>122,500</b>	<b>214,533</b>	<b>16,748</b>	<b>17,772</b>

2013 Planned to be Pledged capacities at the US ATLAS Facilities

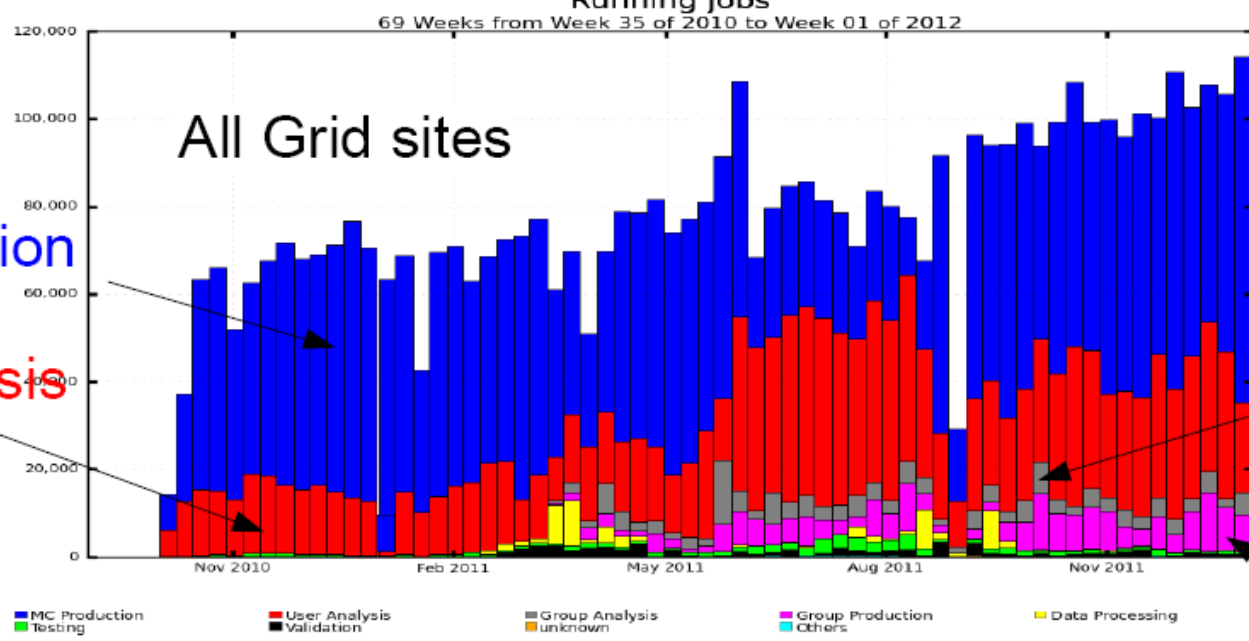
Site	CPU [HEPSpec 2006]		DISK [TB]	
	2013 Pledge	Installed Jan, 2012	2013 Pledge	Installed Jan, 2012
Tier-1	63,000	62,000	7,000	8,100
AGLT2	13,400	36,178	2,500	2,160
MWT2	13,400	36,840	2,500	2,730
NET2	13,400	19,679	2,500	1,200
SWT2	13,400	30,460	2,500	1,439
WT2	13,400	29,376	2,500	2,143
<b>Total</b>	<b>130,000</b>	<b>214,533</b>	<b>19,500</b>	<b>17,772</b>

Site	Tape [TB]		2013 Pledge Installed 1/12	
	2011 Pledge	2012 Pledge		
Tier-1	6,900	8,300	7,590	6,900

CPU Usage  
In 2011

MC production

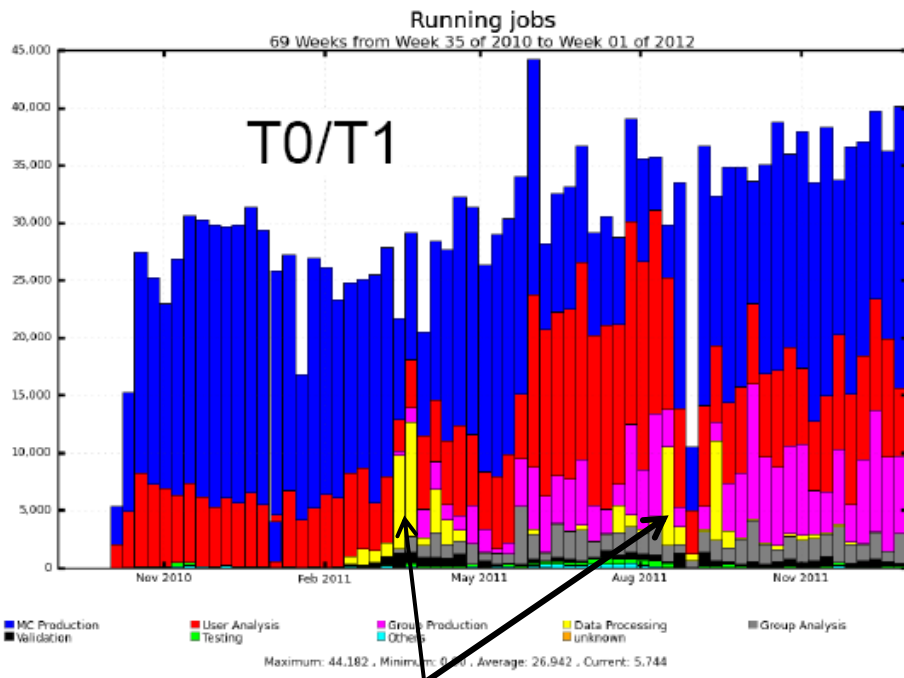
User Analysis



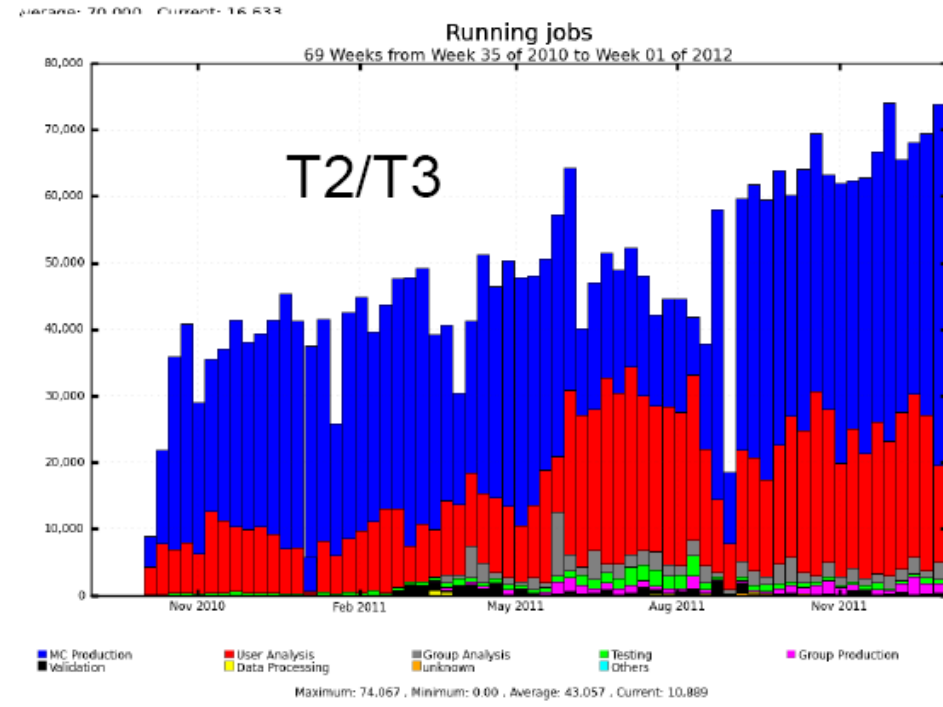
1 week/bin

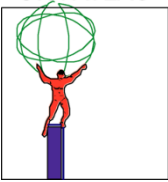
Group Analysis

Group Production



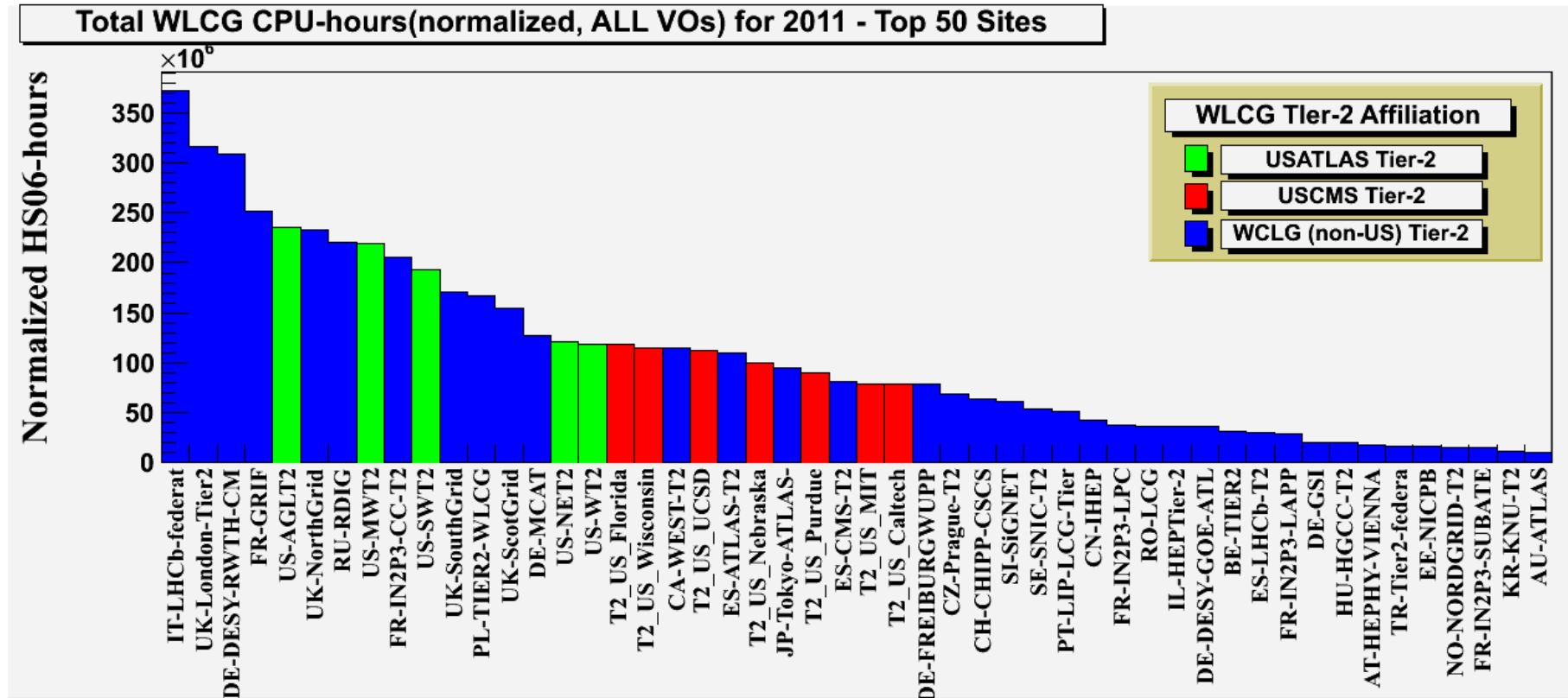
Reprocessing





# WLCG Tier-2 Contribution by Site

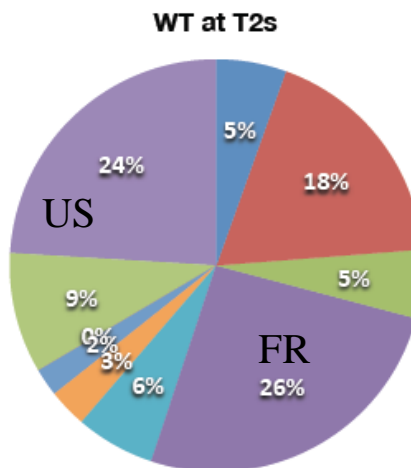
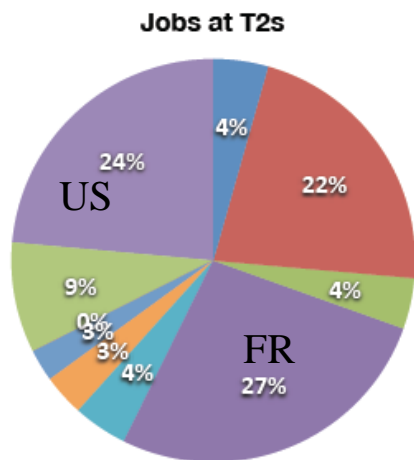
(delivered HEPSpec 2006 \* hours)



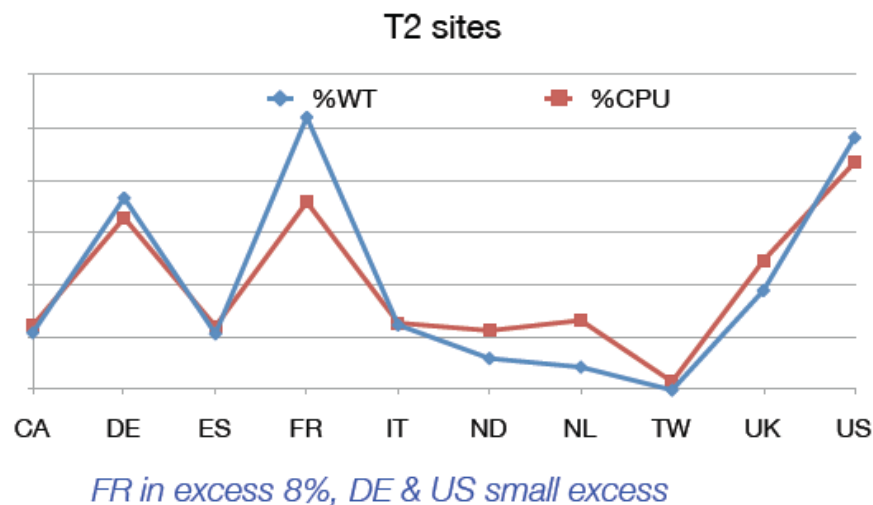
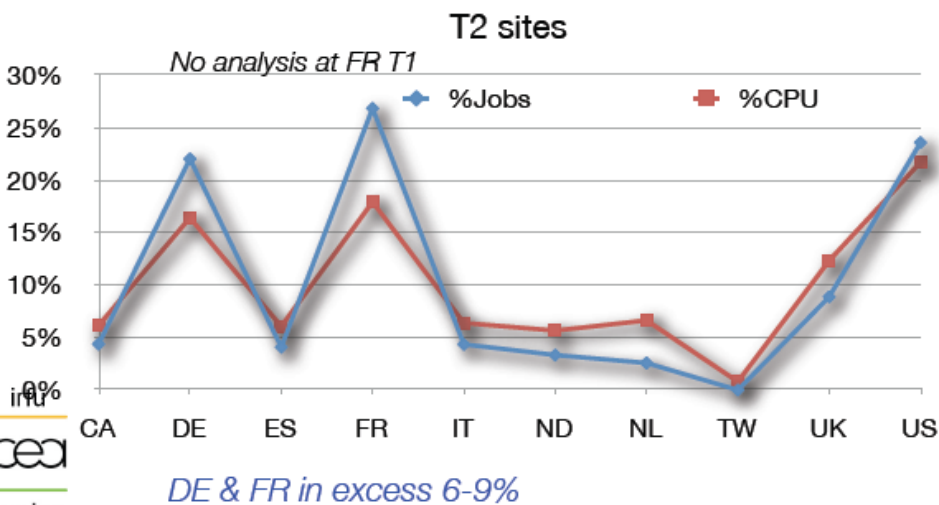
S. McKee

# Analysis jobs at T2s

Presented in July, 2011



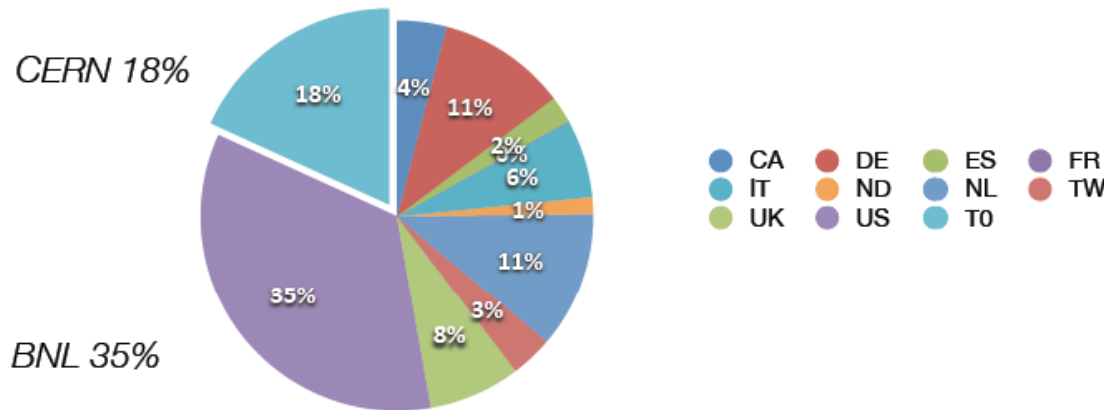
## Relative distribution of analysis jobs compared to CPU pledges



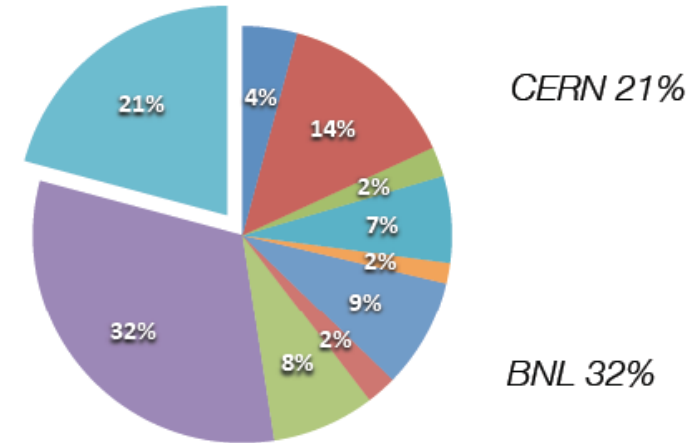
# Analysis jobs at T1s & T0 (CERN)

Slides presented by Chair of the ATLAS International Computing Board in 07/2011

Jobs at T1s & T0



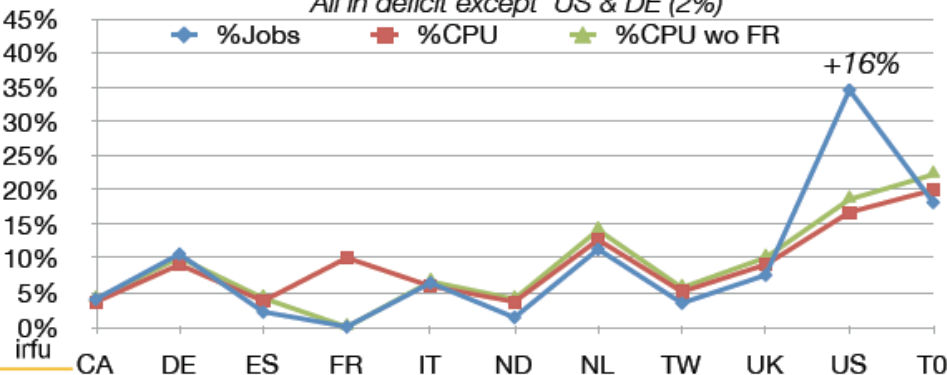
WT at T1s & T0



## Relative distribution of analysis jobs compared to CPU pledges

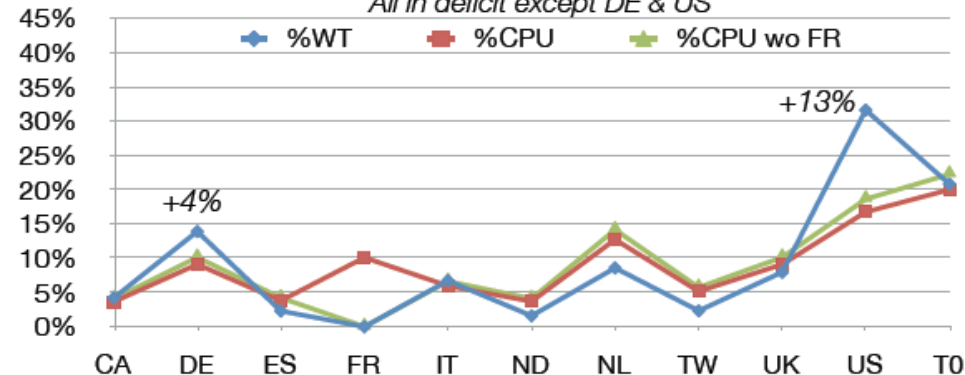
T1 sites

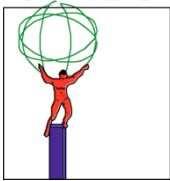
All in deficit except US & DE (2%)



T1 sites

All in deficit except DE & US

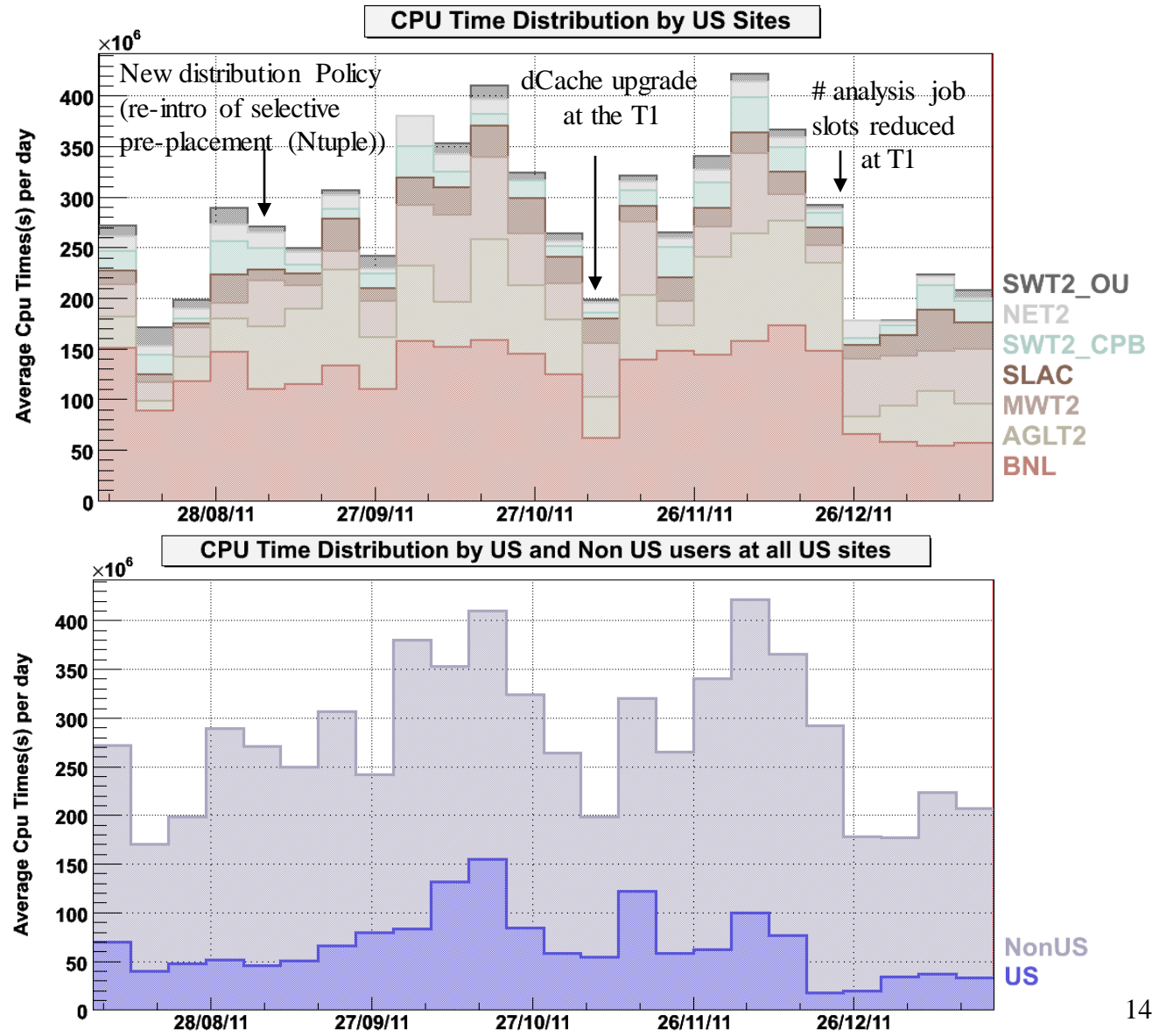




# Contribution to Analysis in the U.S.

(July 2011 – January 2012)

The reduction of analysis slots at Tier-1s to 5% of the total capacity was requested by ATLAS Computing Management in early January in order to speed up high-priority MC production



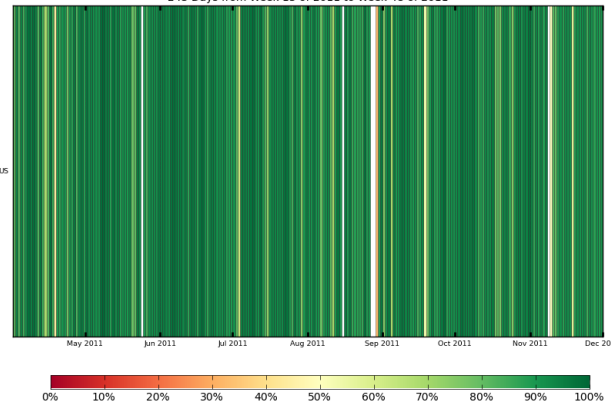
(Plots are stacked)





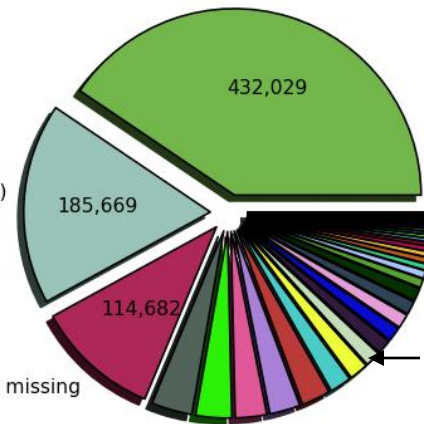
# Analysis Efficiency at US ATLAS T1

Efficiency based on success/all accomplished jobs  
243 Days from Week 13 of 2011 to Week 48 of 2011



Panda Failures by ExitCode (Pie Graph) (Sum: 1,065,666)  
Athena crash - consult log file

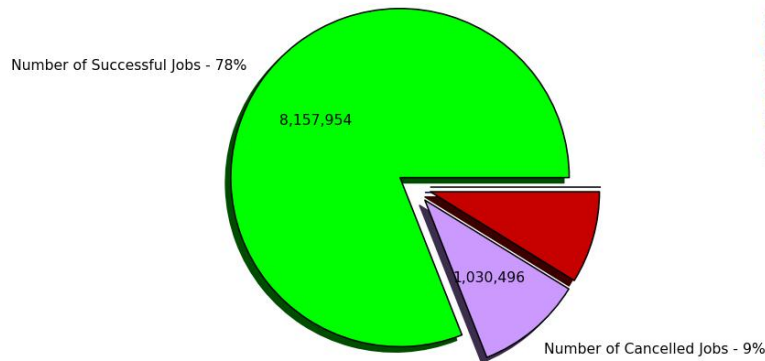
known Problem" (see checklog.txt)



Put error: Local output file missing

First facility related failure (stage-out)

Number of Successful and Failed Jobs (Pie Graph) (Sum: 10,075,481)



Number of Successful Jobs - 78%

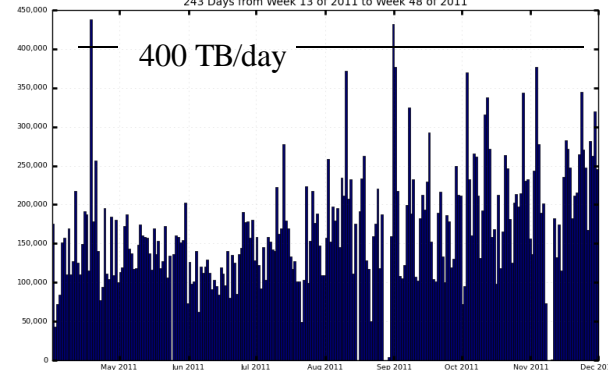
Number of Cancelled Jobs - 9%

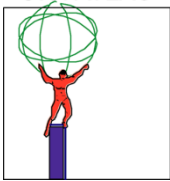
Number of Successful Jobs - 78% (8,157,954)  
Number of Failed Jobs - 11% (887,031)

Number of Cancelled Jobs - 9% (1,030,496)  
Number of Unknown-Status Jobs - 0% (0.00)

- Athena crash - consult log file (432,030)
- Put error: Local output file missing (114,682)
- Expired three days after submission (33,066)
- Athena core dump, or Athena time out, or ConditionsDB exception caught: MySQL error: unable to pin (28,265)
- Not documented, Exitcode: 137 (19,943)
- Put error: Error in copying the file from job workdir to localSE (13,995)
- Not documented, Exitcode: 8 (13,249)
- Put error: LFC registration failed (12,402)
- Adder could not add files to the output datasets (8,130)
- Not documented, Exitcode: 9 (5,540)
- Athena release is not installed in the CE, or trf failed due to "Unknown Problem"
- User work directory too large (37,644)
- trf is not installed in the CE (28,967)
- Athena ran out of memory (17,265)
- Get error: Failed to get LFC replicas (13,467)
- Athena core dump (13,045)
- Get error: Staging input file failed (11,932)
- lost heartbeat (7,038)
- plus 96 more

NBytes Processed in GBs (Time Stacked Bar Graph)  
243 Days from Week 13 of 2011 to Week 48 of 2011





# Comparison w/ Computing Model- CPU

According to Computing Model

Tier-1

Tier-2+3

MC	38% G4
	22% digi+reco
Reproc	20%
Group+User	20%

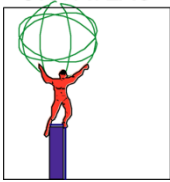
MC	20% G4
	0% digi+reco
Group	20%
User	60%

According to Grid Monitoring

MC	52% G4
	18% digi+reco
Reproc	5%
Group+User	25% (group 10%)

MC	77% ( 67% G4 + other)
	5% digi+reco
Group	2%
User	16%





# ATLAS Resource Usage Report

(Slides presented by B. Kersevan at ICB on 3/15)

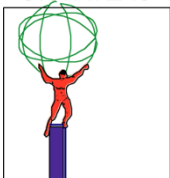
- The ATLAS 2011 Resource Usage report has been completed and submitted on Monday. Some highlights presented here, full document in the twiki:  
[https://twiki.cern.ch/twiki/pub/Atlas/ComputingModel/ATLAS\\_computing\\_usage\\_report\\_2011.pdf](https://twiki.cern.ch/twiki/pub/Atlas/ComputingModel/ATLAS_computing_usage_report_2011.pdf)

<i>LHC and ATLAS parameters for 2011</i>		<i>Predicted p-p</i>	<i>Actual p-p (Jan-Dec 2011)</i>
Rate	events/sec	200	340
Time	Msec total	5.2	4.6
Real data	Bevents total	1.0	1.58
Full Simulation	Bevents total	0.6	2.8
Fast Simulation	Bevents total	0.6	0.7
Real RAW	MB/event	2.8	0.64; 1.2 until July
Real ESD	MB/event	2.2	1.1
Real AOD	MB/event	0.36	0.161
Simulated HITS	MB/event	2	0.8
Simulated ESD	MB/event	3.0	1.9
Simulated AOD	MB/event	0.42	0.3
Full Simulation	HS06sec/event	5100	2700
Fast Simulation	HS06sec/event	400	250
Real Reconstruction	HS06sec/event	200	108
Simulation Recon.	HS06sec/event	340	300

More data and (much) more simulation done in 2011

Lower event sizes due to ATLAS software effort and somewhat lower pile-up

Faster processing speeds due to ATLAS software effort



# Tier-1 Disk

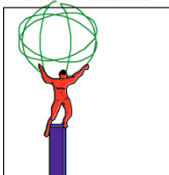
The 0.7 WLCG factor

<i>Tier-1 Disk [PB]</i>	<i>Predicted</i>	<i>Utilizable</i>	<i>Actual (end of 2011)</i>	<i>Comment</i>
Current RAW data	4.3	3.0	8.2	
Real (ESD, AOD, DPD) data	3.8	2.7		
Simulated data	4.4	3.1	11.7	
Calibration and alignment inputs	0.4	0.3	0.3	
Group data	4.5	3.2	1.8	
User data	0.6	0.4	1.3	
Cosmic ray data	0.2	0.1	0.1	
Processing and I/O buffers	3.6	2.5	3.2	3 PB total buffers
<b>Total</b>	<b>22</b>	<b>15.3</b>	<b>26.6</b>	

Big effect of RAW zipping and decreased event sizes

We simulated factors more than expected:  
.. and increased our physics output accordingly!

Consequently we still did use more disk than predicted



# Tier-1 CPU

<i>Tier-1 CPU [kHS06]</i>	<i>Predicted</i>	<i>Actual (average over 2011)</i>
Reprocessing	41	188
Simulation Production	77	
Simulation Reconstruction	28	
Group+User activities	56	44
Misc. other activities	0	12
<b>Total</b>	<b>202</b>	<b>244</b>

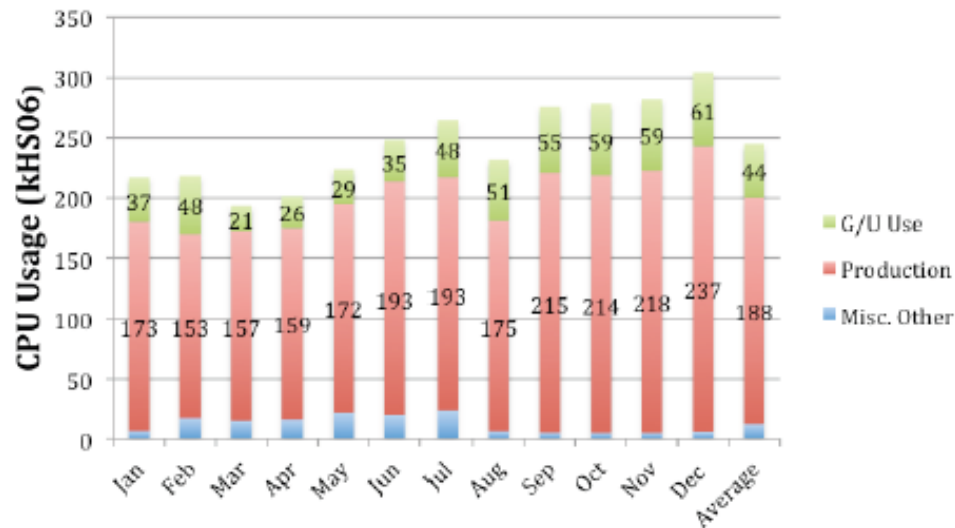
Big impact of simulation production

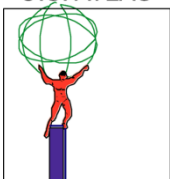
Group+user activities ramped up

We used more CPU than predicted

The Tier-1s as our first-line of service have shown their special role as the providers of capacity for the ATLAS most urgent tasks as Reprocessing, Simulation and Group production, enabling us to deliver the results in time for Summer 2012 conferences, December 2012 Council, Moriond 2013 ....

2011 ATLAS T1 CPU Usage (kHS06)





# Tier-2 Disk

The 0.7 WLCG factor

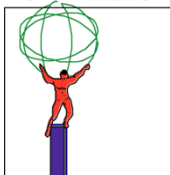
<i>Tier-2 Disk [PB]</i>	<i>Predicted</i>	<i>Utilizable</i>	<i>Actual (end of 2011)</i>	<i>Comment</i>
Real AOD+DPD	13	9	16.4	7.5 PB total cache space provided for dynamic data placement
Simulated data	12	8		
Calibration and alignment output	0.3	0.2	0.3	
Group data	7	5	2.9	
User data	2	1.4	1.4	
Processing and I/O buffers	1	0.7	0.4	4.4 PB total buffers
<b>Total</b>	<b>35</b>	<b>24</b>	<b>22</b>	

Optimizing our replica placement by using PD2P according to popularity matrix

~4 PB of more buffering space needed (and deployed).

- By the end of 2011 we efficiently used practically all of Tier-2 disk space, as predicted! In the first months of 2012 the volumes increased further in preparation for Moriond and other winter conferences.

Statement is only true if we accept 70% Efficiency factor – wasteful and does not reflect reality



# Tier-2 CPU

<i>Tier-2 CPU [kHS06]</i>	<i>Predicted</i>	<i>Actual (average over 2011)</i>
Simulation Production	56	300
Group activities	56	84
User activities	163	
Other Misc. activities	0	21
<b>Total</b>	<b>275</b>	<b>405</b>

Big impact of simulation production

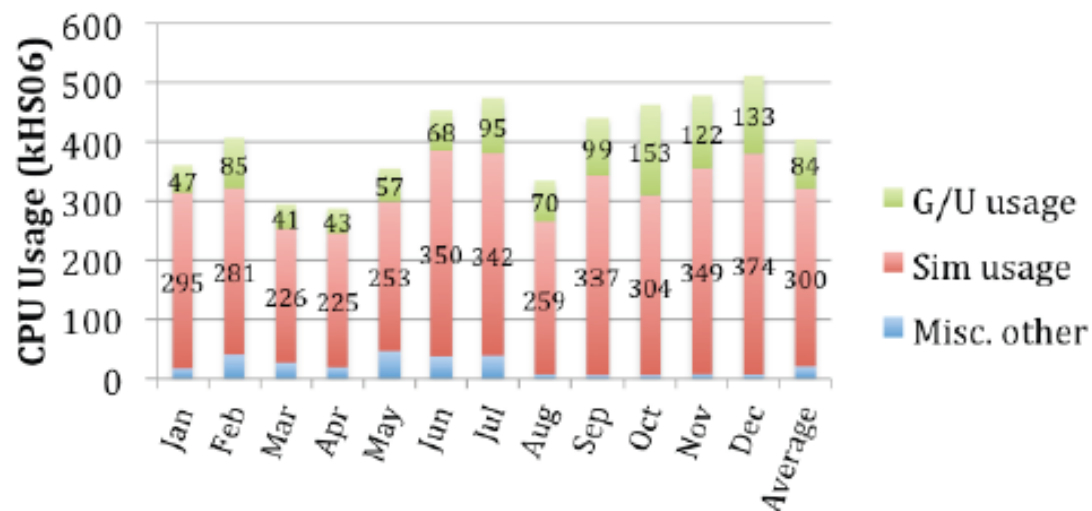
Over-estimated group+user activity:  
Success of group production

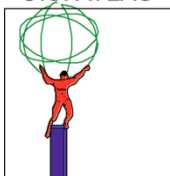
Eventual use much more than predicted

**The Tier-2s, providing the needed extra capacity for simulation production got us out of the tight spots!**

To repeat the argument: we did not do simulation production to fill the capacity but to get more and better physics results! ATLAS made too moderate predictions of what we would need - and we are updating our resource requests accordingly.

**2011 ATLAS T2 CPU Usage (kHS06)**



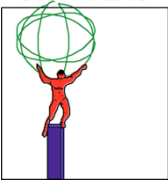


# Resource Request 2012 - 2014

- The resource request is still in preparation ( yes, we are late ), so to present detailed projections is too early.
- The main parameters are:
  - Due to increased pileup and energy ~ 24 collisions on average @ 8 TeV:
    - The event sizes will increase by a considerable factor.
    - The CPU processing times will increase by a considerable factor.
  - ATLAS software is of course fighting this with alacrity.
  - We will need to also optimize our replication policy:
    - Reduce the number of replicas if needed
    - Rely on dynamic data placement to a greater level

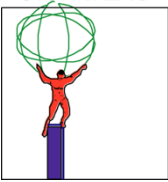
LHC and data taking parameters		2011 pp actual	2012 pp
Rate [Hz]	Hz	340	400
Time [sec]	MSeconds	4.6	3.8
Real data	B Events	1.6	1.5
<b>Simulated data</b>			
Full Simulation	B Events	2.8	1.7
Fast Simulation	B Events	1.6	2
<b>Event sizes</b>			
Real RAW	MB	0.64	0.8
Real ESD	MB	1.1	2.8
Real AOD	MB	0.16	0.42
Sim HITS	MB	0.8	1
Sim ESD	MB	1.9	3.8
Sim AOD	MB	0.26	0.65
<b>CPU times per event</b>			
Full sim	HS06 sec	2700	3300
Fast sim	HS06 sec	250	310
Real recon	HS06 sec	108	230
Sim recon	HS06 sec	300	830
Group analysis	HS06 sec	20	20
User analysis	HS06 sec	0.4	0.4





# Observations & Concerns

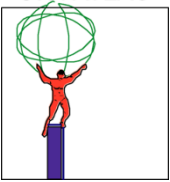
- **Usage of Resources at T1s vs. T2s seems suboptimal in 2011**
  - ◆ T1: small fraction of (CPU) resources were used for T1-specific tasks (reprocessing, group production and group analysis); simulation and user analysis dominated workload
  - ◆ T2: disk resources were underutilized while T1 disk was always full
  - ◆ Resources at T1s and T2s were underutilized during extended periods
  - ◆ Lack of Transparency (for Facilities) as to what the resources were/are actually used for
    - ◆ Processing tasks
    - ◆ Data Types



# Resource Usage Optimization

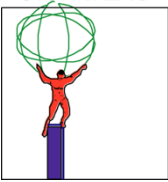
- Computing at T1s is in general significantly more expensive than at T2s
- whenever possible, T2 resources should be utilized for tasks that don't require the level of performance and reliability T1s are required to provide for services across the entire (T1) facility
- as was presented in recent reports and ADC meetings simulation dominates resource usage at T1s and T2s, while T1-specific tasks (i.e. reprocessing, which is also under discussion for execution at T2s, at least for MC) were rarely executed by ATLAS
- there is no compelling reason for providing resources for AOD and D3PD/ntuple (user) analysis at T1s. Analysis resources there should be restricted to group analysis, and for such analyses that require access to data categories that are not accessible at other Tiers.





# Resource Usage Optimization

- **Proposal: Regions should be given some flexibility as to how/where they provision computing resources to ATLAS**
  - ◆ shifting resources between Tiers should be an option as long as the aggregate capacity across a region equals the amount requested by ATLAS, and the performance and reliability of the overall regional facility complex is at the required level.
  - ◆ regional flexibility should be allowed regarding the distribution of resources at T1 vs. T2s, in particular CPU but also disk
- **US ATLAS pledging resources at ~23% of total, but deliver**
  - ◆ 120% of pledged Tier-1 capacity (CPU and Disk)
  - ◆ 180% of pledged Tier-2 capacity (CPU) and 100% Disk
  - ◆ Computing in the US is run as a program with the T1 and all T2s integrated and managed under the US ATLAS Operations Program/Facilities



# Evolution in LHC Computing

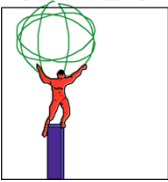
- WLCG TEGs have almost completed their work
  - ◆ Summary by Simone at <https://indico.cern.ch/getFile.py/access?contribId=77&sessionId=0&resId=1&materialId=slides&confId=169695>
- SW&C workshop last week at CERN focused on development activities in the area of distributed computing
- Despite popular opinion, the Grid is doing us very very well.
  - ◆ 1000's of users can process petabytes of data with millions (billions?) of jobs
- At the same time, we are starting to hit some limits:
  - ◆ Scaling up, elastic resource usage, global access to data
- What can we learn from external innovations? (without disrupting operations!)
- Various R&D Projects and Task Forces were formed one year ago
  - ◆ NoSQL databases R&D
  - ◆ Cloud Computing R&D
  - ◆ XROOTD Federation and File level Caching Task force
  - ◆ Event Level Caching R&D
  - ◆ Tier3 Monitoring Task Force
  - ◆ CVMFS Task Force
  - ◆ Multicores Task Force
  - ◆ Also *Network Monitoring...*
- <https://twiki.cern.ch/twiki/bin/viewauth/Atlas/TaskForcesAndRnD>



# NoSQL: “Big data processing” in DDM

## DB SCALING

- Evaluated several products, finally chose the Hadoop ecosystem
  - ◆ distributed filesystem, non-relational database, SQL-like and non-SQL-like data processing languages
  - ◆ can also distribute any application that uses stdin/stdout (e.g. grep, awk, python, ....)
  - ◆ 12 nodes managed by Puppet, 30TB available space
    - Puppet configuration can be shared!
    - encrypted hourly backup of cluster state to Dropbox
  - ◆ stable and efficient backend
    - verified recovery under hardware and software failures
    - zero (point one) maintenance effort
    - upgraded cluster from SLC5 to SLC6 in place with only 5 minute downtime, got IO boost of factor 3 due to new async-IO kernel
  - ◆ migrated one application from Cassandra to Hadoop within 2 days
  - ◆ Excellent foundation for data-related statistics plots
    - In the US we are planning to use the infrastructure to show what our disk resources are used for



# BNL Cassandra: Current Status

- Test cluster at BNL has been in operation for many months with no major problems, is largely automated and demonstrated resilience against hardware failures in a few incidents
- Cassandra DB is exposed via a Web service (JSON format for data)
- Not a replacement for Oracle but is to complement it and reduce its load
- We'll need both job and file data (so far only job data has been loaded)
- Existing 1TB SSD capacity per node allows for roughly 6 months worth of monitoring data to be stored (which includes operations overhead)
  - Hiro has managed to create a variety of statistics plots very efficiently

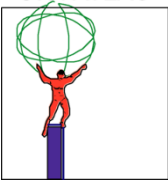


# Virtualization and Cloud R&D

EFFICIENCY, ELASTICITY

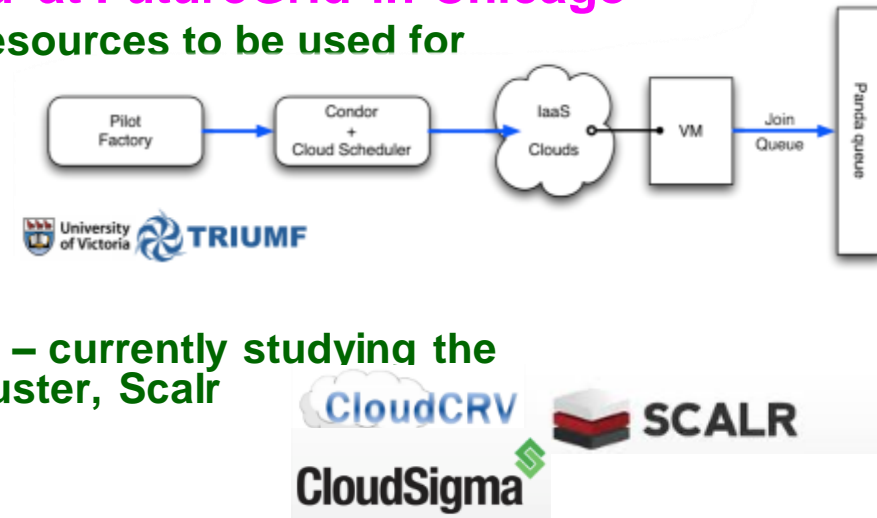
- Active participation, almost 10 persons working part time on various topics
  - ◆ <https://twiki.cern.ch/twiki/bin/view/Atlas/CloudcomputingRnD>
  - ◆ In the US activities at the BNL, LBNL, Fresno and by Doug
- Data Processing
  - ◆ Panda Queues in the Cloud
    - Centrally managed, non-trivial deployment but scalable
    - Benefits ATLAS & sites, transparent to users
  - ◆ Tier3 Analysis Clusters (Instant cloud site)
    - User/Institute Managed, Low/Medium Complexity
  - ◆ Personal Analysis Queue (~One click, run my jobs)
    - User Managed, Low Complexity (almost transparent)
- Data Storage
  - ◆ Short term data caching to accelerate above data processing use cases
    - Transient data
  - ◆ Long term data storage in the cloud
    - Integrate with DDM

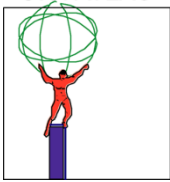
Slides by  
Dan



# Cloud: Achievements

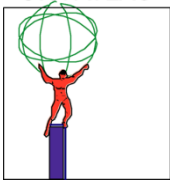
- We saw some good achievements presented Tuesday:
  - ◆ UVic Cloud Scheduler demonstrated at FutureGrid in Chicago
    - Will be good for quickly adding new resources to be used for analysis or production
    - Nearly in production, I/O to be tuned
  - ◆ Proof/Batch Analysis Clusters
    - Good performance evaluations
    - The “instant cluster” use case is clear – currently studying the management tools: CloudCRV, StarCluster, Scalr
  - ◆ MC Production in the Cloud
    - Demonstrated already at CloudSigma. Possible to benefit from above tools?
- Plus other activities: PanDA data in Cloud, APF, OpenStack, ...





# Cloud Futures

- Many of the activities are reaching a point where we can start getting feedback from users
  - ◆ Should focus in the next months on eliminating options, and determine what we can deliver in production
- Cloud Storage
  - ◆ This is the hard part. Some free S3 endpoints are just coming online, so effective R&D is only starting now.
  - ◆ Looking forward to good progress in caching (Xrootd/HDFS in cloud) and DDM S3 Evaluation (Rucio incubator proposal)
- Support the grid sites who want to offer private cloud resources
  - ◆ Develop guidelines, best practices
  - ◆ Good examples already, e.g. LxCloud, PIC, BNL, and others.



# Data Federations

GLOBAL DATA ACCESS

- Since Sep 2011 ~10 sites reporting to global federation
- Performance studies with various TTreeCache options under study
- Adoption if decent WAN performance achievable



USATLAS Federated Xrootd Status - 2012-03-15 08:30:33

Frequently Asked Questions

Host: atl-prod09.slac.stanford.edu (atl-prod09.slac.stanford.edu)

Metric	Last Executed	Enabled?	Next Run Time	Status
org.atlas.xrootd.grid.vxrep-compare	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.grid.vxrep-direct	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.grid.vxrep-fax	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.ping	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	OK

Host: atlas29.hep.anl.gov (atlas29.hep.anl.gov)

Metric	Last Executed	Enabled?	Next Run Time	Status
org.atlas.xrootd.ping	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-compare	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	CRITICAL
org.atlas.xrootd.vxrep-direct	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	CRITICAL
org.atlas.xrootd.vxrep-fax	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	CRITICAL

Host: atlgirdtp01.phy.duke.edu (atlgirdtp01.phy.duke.edu)

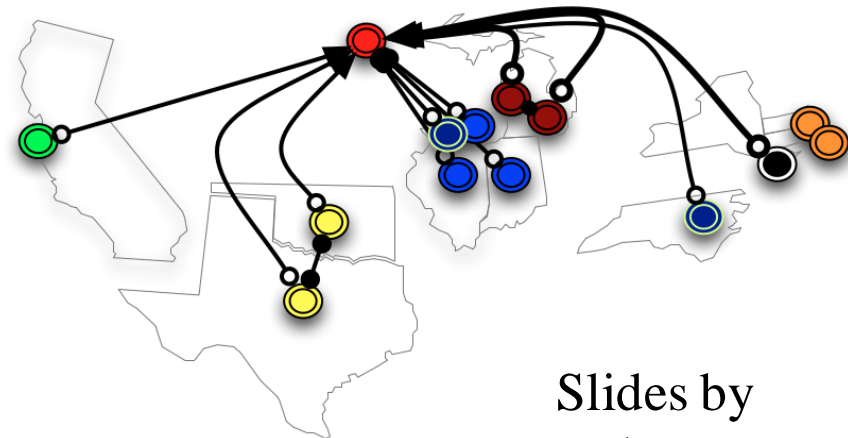
Metric	Last Executed	Enabled?	Next Run Time	Status
org.atlas.xrootd.ping	2012-03-15 08:20:03 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-compare	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-direct	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-fax	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	OK

Host: dedoor09.asatlas.hnl.gov (dedoor09.asatlas.hnl.gov)

Metric	Last Executed	Enabled?	Next Run Time	Status
org.atlas.xrootd.ping	2012-03-15 08:20:03 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-compare	2012-03-15 08:20:01 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-direct	2012-03-15 08:20:03 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.vxrep-fax	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK

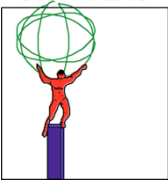
Host: dedoor10.asatlas.hnl.gov (dedoor10.asatlas.hnl.gov)

Metric	Last Executed	Enabled?	Next Run Time	Status
org.atlas.xrootd.grid.vxrep-compare	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.grid.vxrep-direct	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.grid.vxrep-fax	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK
org.atlas.xrootd.ping	2012-03-15 08:20:02 CDT	YES	2012-03-15 08:35:00 CDT	OK



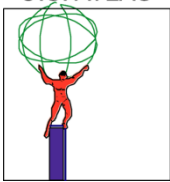
Slides by  
Rob





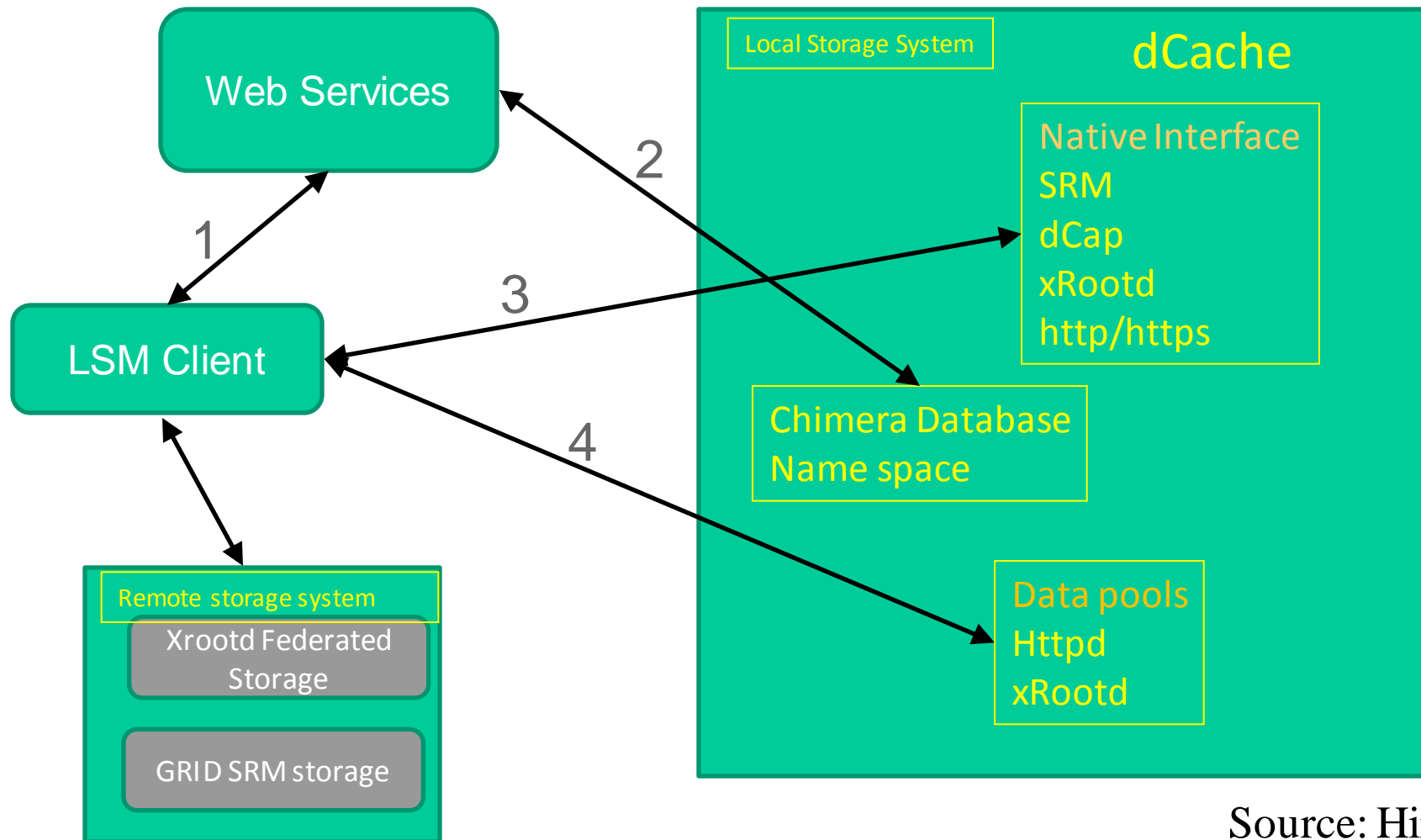
# Xrootd Federation R&D to Production

- Current set of sites regularly testing at significant analysis job scale (in HammerCloud)
- Provide redirector of highly performing data sources
- More experience with TTreeCache settings with well defined examples for users
- More monitoring of IO
- Activate X509 on all sites
- Explore augmenting current ANALY workflow to use FAX when problems with local SE or missing files (Athena or lsm, eg.)
  - ◆ or to expand number of queues available to users (no local input dataset requirement)
- Other regions express interest in trying

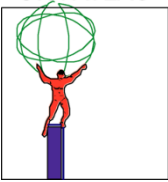


# Resilient data access with Local Site Mover (LSM) w/ multi-protocol support

Goal: Prevent jobs from failing because files in the local storage system are not accessible, i.e. in case of transient failures



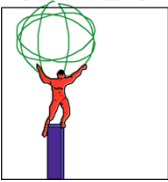
Source: Hiro



# Tier3 Site Monitoring

UNDERSTAND USAGE

- **Aim** – enable monitoring of off-grid Tier3 sites
- **Monitoring should include**
  - ◆ Infrastructure monitoring
  - ◆ Storage monitoring
  - ◆ Task processing monitoring
  - ◆ Site summaries presentation at Dashboard's historical view, DQ2 popularity
- **Methods**
  - ◆ Ganglia as site level monitoring provider
  - ◆ Storage and task processing software Ganglia add-ons development
  - ◆ Development of Ganglia add-ons to data transmission to Dashboard, DQ2 popularity
- **Status**
  - ◆ Proof, XRootD, Lustre monitoring solutions for sites are ready
  - ◆ Wiki page with installation/configuration instructions: <https://svnweb.cern.ch/trac/t3mon/wiki/T3MONHome>
  - ◆ Repo: <http://t3mon-build.cern.ch/t3mon/>
  - ◆ Several installations on sites in Russia, USA
- **Plan**
  - ◆ Completion of data transmission add-ons development
  - ◆ Involve more sites into testing



# XRootD Federation Monitoring

- **Goal – Monitor data transfer between sites in federation**
  - ◆ **Collector is ready**
  - ◆ **Site level interface for Ganglia is in development**
- **Summaries of data transfers will be presented at federation level interface**
  - ◆ **ActiveMQ transmission module is ready**
  - ◆ **Federation transfer interface is in system design stage**

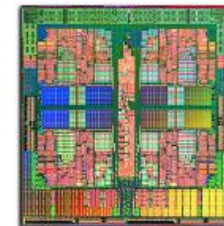


# Multi-Core

Slides by  
R. Walker

CPU SCALING

- **Multicore or whole-node scheduling coming**
  - ◆ 64 bit reco memory footprint & AthenaMP
  - ◆ local batch system scaling - otherwise #slots blows up
  - ◆ cloud computing - whole-node (VM) by definition
- **Readiness to use**
  - ◆ running N serial job on whole-node slot is inefficient - hold node until last job finishes
  - ◆ AthenaMP for Reco and G4sim, but not digi (prevents MC reco usage)
    - no memory requirement for G4sim, but still want MP
    - fewer and bigger output files (helps DDM and no need to merge)
    - hope AthenaMP for G4sim can be validated and used for part of mc12
  - ◆ use resources only available as whole-node (cloud, hpc site)
    - need to get more experience using sites offering these (CERN and several T1/2s have multi-core queues)



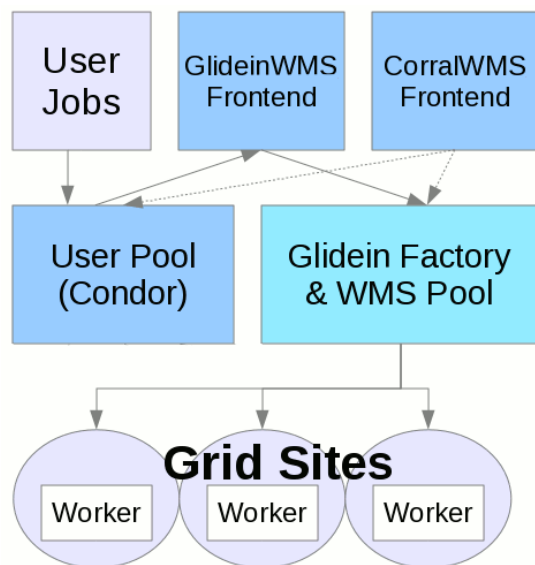
See the WLCG WM TEG Report on Multicore:

<https://twiki.cern.ch/twiki/bin/view/LCG/WMTEGMulticore>

Propose to add min/max CPU requirements to the CE JDL.



# GlideInWMS



- **Motivation**

- ♦ stop running user jobs under same uid on WN, or risk security incident
- ♦ common solution with CMS to use glxexec

- **Overview**

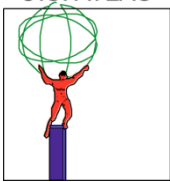
- ♦ glideinWMS builds distributed Condor pool (popular local batch system)
- ♦ user-specific panda pilots submitted to this, with user proxy (via MyProxy) delegated to WN

- **Status**

- ♦ partially used in conjunction with new pilot factory (APF)
- ♦ no user proxies used during test phase
- ♦ plan to scale up to O(10k) on few test sites
- ♦ use also for production pilots (to increase scale)
- ♦ transparent to user, apart from MyProxy store requirement from client
- ♦ potential benefit (fairshare transparency, ssh-to-job) but also risks - testing phase
- ♦ no immediate requirement for glxexec from all sites, but volunteer sites welcome
  - WLCG likely to require sites to install glxexec in the mid-term

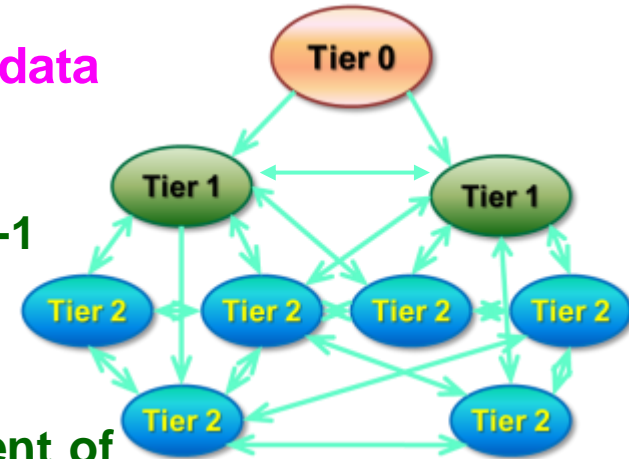
- **Risk**

- ♦ Overlay system w/ PanDA on top of glideinWMS increases complexity
  - Hard to debug in case of problems (why are jobs not flowing?)
- ♦ Potential ownership/permission-related DDM issues (on SEs and in LFC)

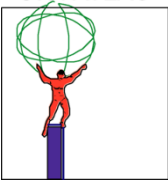


# The ATLAS Data Model has changed

- Moved away from the strict MONARC model
- 4 recurring themes:
  - ◆ Flat(ter) hierarchy: Any site can replicate data from any other site
  - ◆ Multi Cloud Production
    - Need to replicate output files to remote Tier-1
  - ◆ Dynamic data caching: Analysis sites receive datasets from any other site “on demand” based on usage pattern
    - Possibly in combination with pre-placement of data sets by centrally managed replication of whole datasets
  - ◆ Remote data access: local jobs accessing data stored at remote sites
- ATLAS is now heavily relying on multi-domain networks and needs decent e2e network monitoring



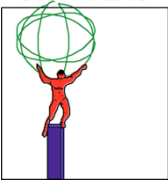
USE THE NETWORK



# Network Monitoring based on perfSONAR

- perfSONAR-PS Deployments in US, Italian and Canadian clouds.
- US completely covers ATLAS Tier-1 and Tier-2 sites
- Italian deployment targets all LHC experiments using one set of perfSONAR-PS instances.
- Canadian cloud chose Dell R610s with 10GE for bandwidth nodes.
- LHCONE deployment targeted at measuring initial baseline (before most sites transition to LHCONE) and then tracking the transition to verify the impact of LHCONE.
  - ◆ This deployment is temporary. Long-term LHCONE monitoring plan is under discussion.
  - ◆ 16 “early adopter” sites selected for monitoring
  - ◆ Details at <https://twiki.cern.ch/twiki/bin/view/LHCONE/SiteList>
  - ◆ Monitoring at <https://130.199.185.78:8443/exda/?page=25&cloudName=LHCONE>
  - ◆ Still 2 sites without instances setup: DESY-HH (h/w arrived) and GRIF/LAL
  - ◆ Some issues (configuration problems, firewalls, etc) being looked at
- This is a vital piece of distributed facilities infrastructure, therefore it really needs to become a WLCG core service
  - Was brought to the attention of the GDB and the WLCG Project Leader
  - Should be operated as a community effort





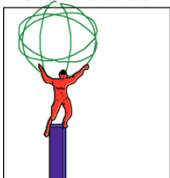
# Summary

- Computing was a big success as enabler for physics, on its own metrics but also on the ultimate metric of timely physics output
- The Facilities, the Tier-1 and the Tier-2 centers, have performed well in 2011 LHC data taking, managed data processing and user analysis
  - ◆ We have a very effective Integration Program in place to ensure readiness of Facility services
- The U.S. ATLAS Computing Facilities need sufficient funding to be on track to meet the ATLAS performance and capacity requirements, long-term and independent from the long shutdown
  - ◆ Tier-2 funding now secured until 2016 (waiting on numbers from NSF)
  - ◆ Facilities, in collaboration with ATLAS Computing Management, needs to work on resource usage optimization to make computing more cost-effective
- Toward the end of 2011 Tier-3 deployment in the U.S. successfully completed
- OSG proposal submitted last year is now funded - continuation secured for 5 more years
  - ◆ DOE/OHEP funding contribution focusing on LHC computing (slight relief for T1 operating funds)
  - ◆ Exact NSF funding level still unclear, ~24% cut in DOE lab funding (down to \$8.250 M from 11M)
  - ◆ In the process of working on “Satellites” to reduce impact from cuts
- Leadership or strong participation of Facilities in ATLAS Computing R&D activities, e.g.
  - ◆ Cloud Computing
  - ◆ Federated Data Stores
- Overall, the Facilities in the U.S. have performed very well during the 2011 run, and I am (still) looking with confidence into the future



# Additional Material

---

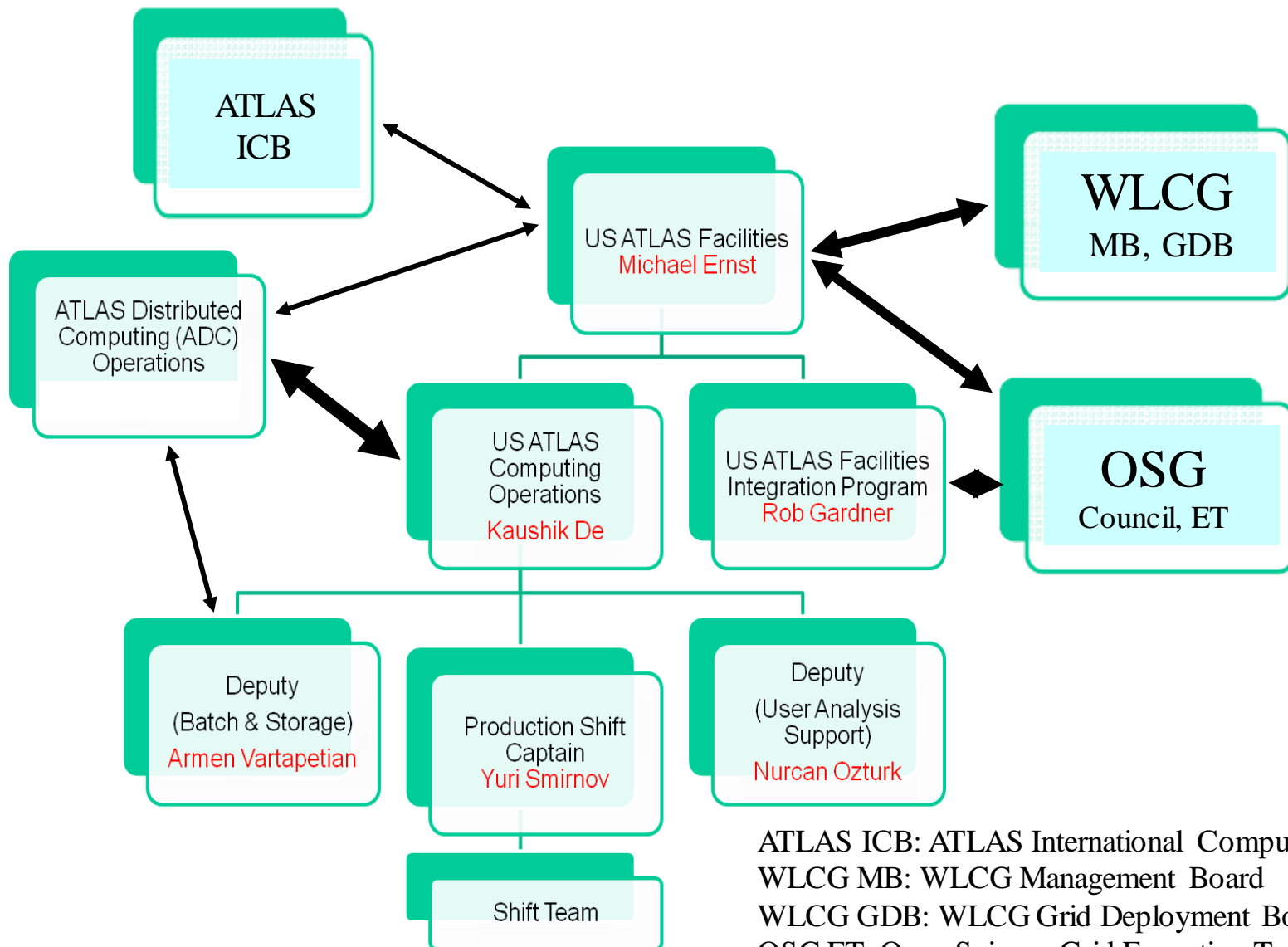


# MC Numbers for 2012

- **Disk consumption, assuming  $\mu=30$ :**
  - **HITS: 1MB/event** (guess based on 0.8 MB/event in 2011)
  - **RDO: 2.4 MB/event** (from ttbar sample with 0.7 factor; was 1.8 MB/event in 2011)
  - **AOD: 0.45 MB/event:** (was 0.25 MB/event in 2011)
  - **ESD: 2.6 MB/event:** (was 1.6 MB/event in 2011)
    - tests from T. Kittelmann:
      - [https://atlas-pmb.web.cern.ch/atlas-pmb/temp/pileup\\_rdoreco\\_jan2012/](https://atlas-pmb.web.cern.ch/atlas-pmb/temp/pileup_rdoreco_jan2012/)
    - using factor 0.7 to get from ttbar to the average (derived from MC11)
  - We need input from Phys. Coord. on the expected total numbers of events.
    - e.g. 5B events, AOD size 0.45M = 2.25 PB single replica (!).
- **CPU consumption:**
  - **Fullsim G4: 360 sec/event (x 9 = 3300 HSo6)**
    - taken from MC11 300 sec/event times 1.2 factor (guess based on energy increase).
  - **Fastsim G4: 34 sec/event ( 310 HSo6)**
    - taken from MC11 28 sec/event times 1.2 factor.
  - **Digi+reco:** based on  $\mu=30$  current estimates:
    - Digi  $\mu=30$  for ttbar (grid jobs): 50 sec/event so multiply by 0.7 gives 35 sec/event for the average.
    - Reco  $\mu=30$  from (tests from T. Kittelmann on above RDOs):
      - [https://atlas-pmb.web.cern.ch/atlas-pmb/temp/pileup\\_rdoreco\\_jan2012/](https://atlas-pmb.web.cern.ch/atlas-pmb/temp/pileup_rdoreco_jan2012/)
    - gives 42 sec/event, again multiply by 0.7 gives 30 sec/event for the average.
    - So the sum would give **65 sec/event ( 590 HSo6)**, to be compared to 33 sec/event from MC11.



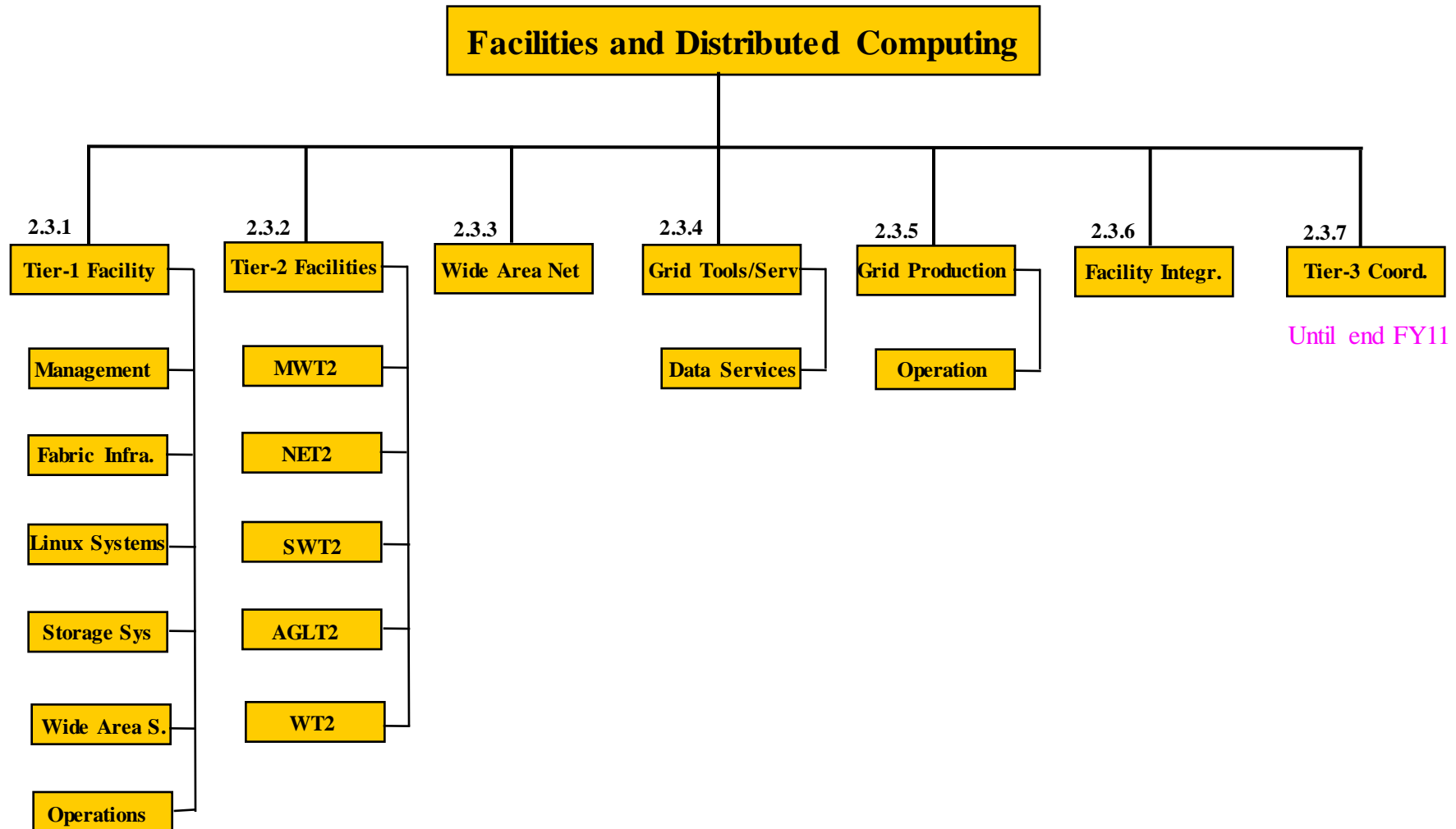
# U.S. ATLAS Facilities Operations - Present and Future

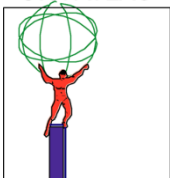


ATLAS ICB: ATLAS International Computing Board  
WLCG MB: WLCG Management Board  
WLCG GDB: WLCG Grid Deployment Board  
OSG ET: Open Science Grid Executive Team



# WBS 2.3 - Facilities Organization

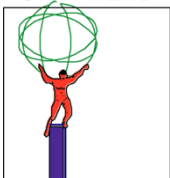




# Physics drives Computing

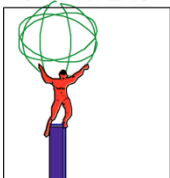
(3 slides from J. Shank at September LHCC)

- Our Computing Model and resource usage is driven by our physics goals
  - We don't want computing to be a bottleneck to physics publications
    - We have achieved this so far in the LHC data taking era
      - For example, some recent results presented at the EPS conference were using data that had been taken three weeks before the conference.
- Changes to our data distribution plan made better use of our facilities and allowed us to take data at a higher trigger rate.
- One ingredient allowing us to achieve this is a lot of hard work that went into improving our reconstruction CPU time/event and event sizes in the face of increasing pileup
- Work continues on these improvements and we also are improving our simulation time and our fast Monte Carlo, which now reproduces the data well enough that physics groups are using it for publications in preparation now. (up to 10 times faster than our Geant4 simulation)



# New Data Distribution Policy

- We have eliminated bulk ESD as an analysis format.
  - Just maintain a rolling buffer at T1's of a small fraction of the ESD for special studies
  - Helped successfully guide users to more space-efficient derived physics data.
- RAW (compression)
  - We place one full copy of the RAW data on disk distributed over the ensemble of Tier 1s, but we now compress the RAW, making the data 60% of the uncompressed size.
- Physics analysis is done from AOD, dESD and Ntuples produced by physics groups.
  - Distributed to T1s, and dynamically to T2
- Overall, our data distribution has been flexible and we adapt it in order to optimize our physics output and make maximal use of available CPU and disk resources.



# Dynamic Data Placement

- In use in ATLAS for over 1 year now.
- Derived physics datasets are copied to T2 disk only when there is a demand
  - The first time a user submits a job using a dataset as input, that data is copied to a T2 (selected by brokering algorithm)
    - Subsequent jobs will go to that T2 up to a set number
    - Further jobs submitted using that dataset trigger more replicas to other T2's
- Some T2s were being under-utilized (both CPU and Disk)
  - Algorithms (both brokering and data placement) are being tuned to fix this problem
- AOD, Ntuple and DESD are popular → we now are doing some pre-placement of these to level CPU usage across sites
- We are still in the era where physics group's Ntuple samples are small enough (~ few TB) that users can copy them to local resources to do analysis.
  - We expect this to not be possible soon (2012) and users will have to use the T2s for this part of analysis also, so user activity will increase
- Current (20-09-2011) T2 disk usage: 20 PB
  - 35 PB pledged until April, 2012.
    - We expect to fully utilize this space by then





# mc11 sample size



- Physics / CP group needs for MC sample size are unprecedented:

## MC11 simulation (full+fast)

NEvents Processed in MEvents (Million Events) (Cumulative Graph)

4607 Hours from Week 26 of 2011 to Week 02 of 2012 UTC

Simulation (G4+fast)

Aug – Dec 2011

Legend for MC11 simulation (full+fast):  
Egamma (100.0%), Tau (100.0%), Top (100.0%), Higgs (100.0%), Exotics (100.0%), BPhys (100.0%), SUSY (100.0%), Egamma (100.0%), Tau (100.0%), Top (100.0%), Higgs (100.0%), Exotics (100.0%), BPhys (100.0%), SUSY (100.0%)

## MC11 a+b+c reconstruction (full+fast)

NEvents Processed in MEvents (Million Events) (Cumulative Graph)

4607 Hours from Week 26 of 2011 to Week 02 of 2012 UTC

digi+reco

Legend for MC11 a+b+c reconstruction (full+fast):  
Egamma (100.0%), Tau (100.0%), Top (100.0%), Higgs (100.0%), Exotics (100.0%), BPhys (100.0%), SUSY (100.0%), Egamma (100.0%), Tau (100.0%), Top (100.0%), Higgs (100.0%), Exotics (100.0%), BPhys (100.0%), SUSY (100.0%)

Total: 2,536 - Average Rate: 0.00/s

## Current MC11b requests

### HIGH PRIORITY FULLSIM

Group Tau stats 40.17M  
Group StandardModel stats 151.91M  
Group Exotics stats 26.852M  
Group Higgs stats 217.10M  
Group Top stats 53.49M  
Group FlavourTag stats 37.1M  
Group JetEtmis stats 174.55M  
Group BPhys stats 51.0M  
Group SUSY stats 56.34M  
Group Egamma stats 145.9M  
Group physics stats 81.98M  
total 1036.39M

### NORMAL PRIORITY FULLSIM

Group Tau stats 16.5M  
Group StandardModel stats 93.3M  
Group Exotics stats 34.6M  
Group Higgs stats 37.84M  
Group Top stats 7.65M  
Group FlavourTag stats 18.9M  
Group JetEtmis stats 203.6M  
Group BPhys stats 5.8M  
Group SUSY stats 23.0826M  
Group Egamma stats 47.65M  
Group physics stats 24.6M  
513.52 M

### HIGH PRIORITY Atlfast-II

Group Tau stats 0.5M  
Group StandardModel stats 116M  
Group Exotics stats 5.125M  
Group Higgs stats 53.16M  
Group Top stats 228.57M  
Group JetEtmis stats 16.2M  
Group BPhys stats 31.5M  
Group SUSY stats 23.48M  
Group Egamma stats 98.0M  
total 572.57M

TOTAL HIGH PRIO: 1665M  
(including single particles)

### NORMAL PRIORITY Atlfast-II

Group SUSY stats 87.9M  
Group Top stats 114.48M  
Group StandardModel stats 56.1M  
Group BPhys stats 106.8M  
Group Exotics stats 6.655M  
total 372 M

TOTAL NORMAL PRIO: 907 M  
(including single particles)

Info from Borut Kersevan

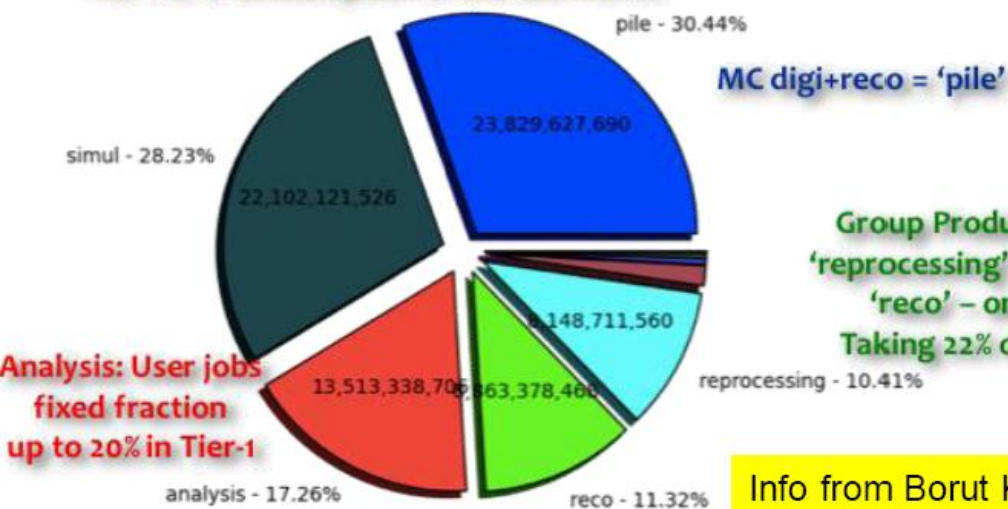


# Monte Carlo production



- Huge effort from production team, ADC, computing sites to ensure high availability of resources over Christmas (also using Tier-0...)
  - Achieved rates: digi+reco 27M events/day, G4 sim 11M/day, AtlFast II sim 6M/day
  - digi+reco is bottleneck, did not reach hoped for 50M/day
    - Huge unanticipated load from **group production** (DPDs) taking 22% of Tier-1 CPU – had to be manually throttled to avoid blocking MC production
    - Huge pileup sample disk space requirements (8 TB), only few Tier-2s can run digi+reco

Tier-1 CPU consumption in the last month



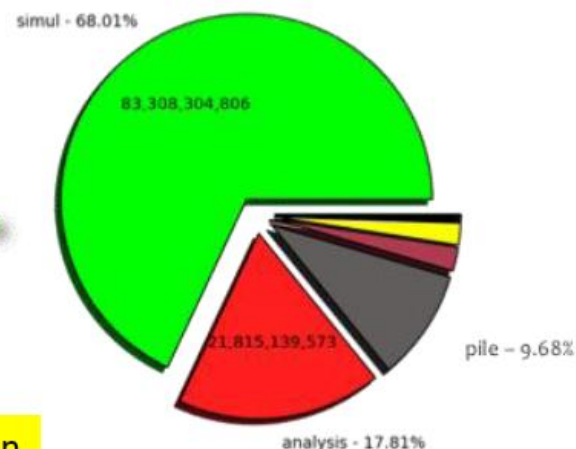
Analysis: User jobs  
fixed fraction  
up to 20% in Tier-1

10th January 2012

Info from Borut Kersevan

Richard Hawkings / Kevin Einsweiler

Tier-2 CPU consumption in the last month



7

In early January CREM decided to reduce analysis share at T1s to 5%

- Slightly higher in the US (~15%) due to non-pledged resources



# MC Disk Space Consumption

- Average event sizes from DDM (thanks to Cedric for the lists!) for MC11:
  - Fullsim HITS: **0.8 MB/event** (was 1MB/event in MC10)
  - Fastsim HITS: **0.6 MB/event**
  - AOD: **0.26MB/event** (was 0.4MB/event in MC10)
    - Average number of replicas : **4** (2 in Tier-1 + 2 in Tier-2).
  - ESD: **1.9 MB/event** (was up to 3MB/event in MC10)
  - RDO: **1.8 MB/event** (was 2MB/event in MC10)
- From Stephane's monitoring and DDM the current occupancies on DATADISK are for MC (**numbers are still increasing!**):

## Tier-1: only primary replicas

Tier-1	AOD	ESD	RDO	EVNT	D3PD/NTUP	Other	Sum
MC10b	951.5	269.4	182.8	199.4	101.4	243.7	1948.2
MC11a	753.8	406.7	86.9	136.8	210.4	146.9	2656.8
MC11b/c	517.4	397.9					
<b>Total</b>	<b>2459.1</b>	<b>1074</b>	<b>269.7</b>	<b>336.2</b>	<b>311.8</b>	<b>390.6</b>	<b>4841.4</b>
<b>Comp. Model</b>	<b>1080</b>	<b>none</b>	<b>171</b>	<b>none</b>	<b>139 (as legacy)</b>	<b>none</b>	<b>1390</b>

## Tier-2: primary (MC10) + secondary replicas

Tier-2	AOD	ESD	RDO	EVNT	D3PD/NTUP	Other	Sum
MC10	1117.2	340.4	none	11.6	321.3	0.3	2609.2
MC10b	516.8	301.6					
MC11a	927.9	0	none	9.4	1580.5	0.0	3037.4
MC11b/c	519.6	0					
<b>Total</b>	<b>3081.5</b>	<b>642</b>	<b>none</b>	<b>21</b>	<b>1901.8</b>	<b>0.3</b>	<b>5646.4</b>
<b>Comp. Model</b>	<b>7300 (8 replicas!)</b>	<b>none</b>	<b>none</b>	<b>none</b>	<b>1100 (as legacy)</b>	<b>none</b>	<b>8400</b>





# DATADISK Situation

(Borut Kersevan at 26 Jan CREM)

## How much do we need for MC12?

- In 2012 we should get additional 1.2 PB for MC in Tier-1 and 3 PB in Tier-2.
- Guessing the volume a problem (being worked on):
  - e.g. 5B events, AOD size 0.3M = 1.5 PB single replica (!).

## What is the legacy we need to keep? Tier-1 is the concern (move to Tier-2?):

### MC11 a/b/c:

- MC11 b/c: one AOD replica will be on the order of 800 TB, so assuming 2 replicas 1.6 PB
- MC11a: two replicas about 500 TB
- MC10b: keeping 2 AOD replicas is about 1 PB.
- It would probably make sense to **keep the EVNT** (in 2 replicas?), so order 200TB.
- Keep D3PDs: about 400 TB is the guess.
- Various legacy/special samples (MCXY\_14TeV, 2TeV, 900 GeV): about 1PB.
- Pileup HITS: about 400 TB (can reduce maybe to a few replicas?)
- logfiles: order 300 TB - we cannot move them (to tape) I think
- In total: 5.4 PB

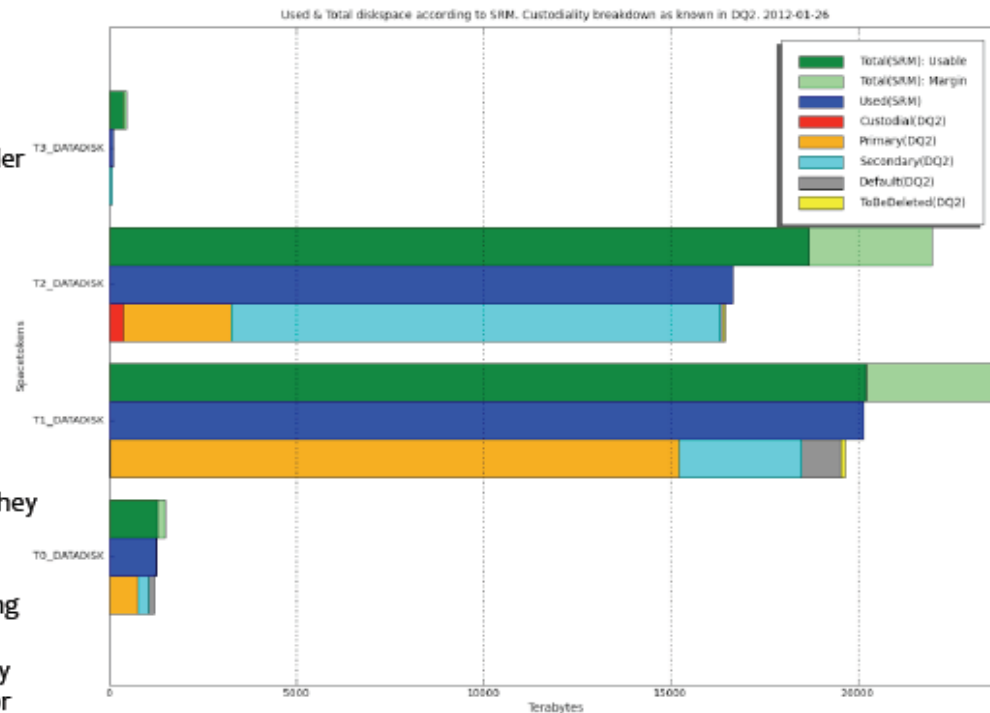
## What we need to clean:

- Transient datasets (gain ~1 PB) - this I should do.
- We clean the ESD and RDO (gain 1.3 PB)? Or move them to Tier-2s? (they are still used, so groups might want to put them on groupdisk)

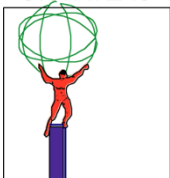
## Proposal:

- Move the legacy/special samples to Tier-2? (about 1PB), this would bring us to 4.4PB in Tier-1 (matching the model).
- Consolidate AODs in Tier-2s, if needed we can then remove the primary copies from Tier-1s and set the AODs in Tier-2s as primary (as we did for MC10).
- We will probably need to keep all of the MC10b/MC11a AODs only on Tier-2s very soon.
- We should also detail the Comp. Model that we can have  $\geq 4$  'active Campaigns' (samples):
  - In 2012: MC10b, MC11a/b/c, MC12(a/b...)
  - MC10 AODs could be retired to tape copies.
- Maybe change the Comp. Model to keep 1 primary copy of ESD and RDO at Tier-2s?

Total MC [TB]	From Monitoring	Comp. Model
Tier-1 (primary)	8125	4400
Tier-2	6572	11000



DATADISK [TB]	Full	Free
Tier-1 (primary)	20147	3652
Tier-2	16656	5317



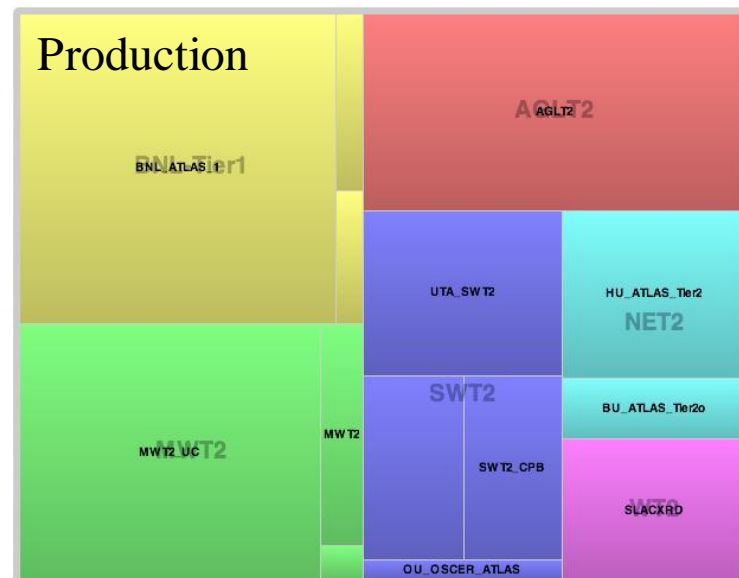
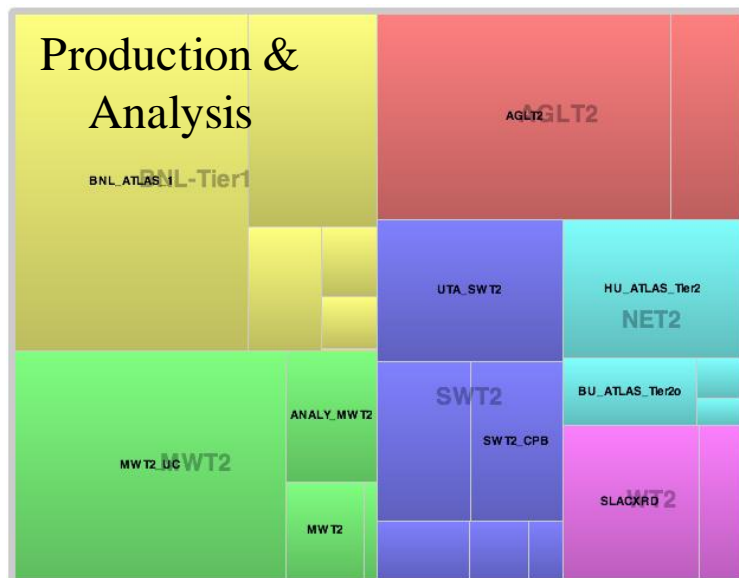
# Resource Planning for 2012 – 2014

## (Borut Kersevan at Mar 1 CREM)

- Some assumptions that came from the discussions (detailed numbers to follow):
  - **We cannot keep more than one MC or data campaign in Tier-1s (it seems):**
    - Change of policy ( preliminaries?):
      - **For both data and MC:**
        - 2 primary AOD copies in Tier-1s.
        - 2 primary DESD copies in Tier-1s.
        - 2 secondary AOD copies in Tier-2s.
        - 1 secondary DESD copy in Tier-2s.
      - Once the reprocessing is done, the Tier-2 secondaries become primary, the Tier-1 copies get cleaned (secondary).
        - Similar as to what we have decided for 2011 MC.
  - **For 2012, assume 1.5 campaigns effective (data and MC):**
    - A preliminary reprocessing ( MC12a, beginning of summer): ~0.25 reprocessings.
    - Full 2012 reprocessing.
    - Reprocessing important searches MC and data from 2011 (i.e. Higgs): ~0.25 reprocessings.
  - **For 2013 and 2014:**
    - Res-simulate all MC 2010-2012.
    - Each year full reconstruction campaign (data and MC)
    - in 2014 order 1B (500 M fastsim) MC @ 13 TeV.
  - **Group+user production rates and data volume equal to what was in 2011 ( defendable?).**
    - Once I do the sums, we might reduce the group space requests somewhat?
  - **Volume of the ‘contingency’ disk space for secondary replicas in Tier-1s and Tier-2s ?**
    - What would be the volume to have? 30% of the total volume? 40% ?



# Relative Contribution by Tier (US only)



By Application

