

Reprocessing LHC beam and cosmic ray data with the ATLAS distributed Production System

J.Catmore, K.De, R.Hawkings, A.Hoecker, [A.Klimentov](#),
P.Nevski, A.Read, G.Stewart, A.Vaniachine and R.Walker
ATLAS Collaboration



Outline

- Introduction
 - ATLAS Production System
 - Data processing cycle
 - LHC beam and cosmic ray data
 - Data volume
 - Data distribution
- Data reprocessing
 - Preparations
 - Data staging tests
 - Conditions data and database access
 - Software and Sites validation
 - Statistics and errors analysis
- Final data distribution
- Summary and conclusions



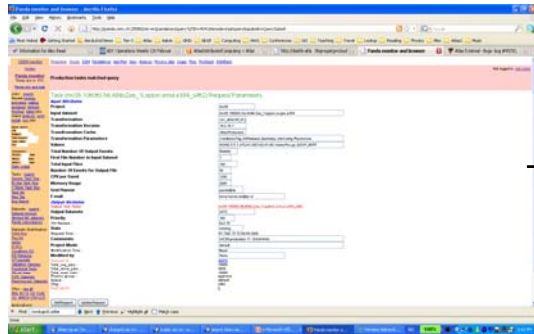
Introduction : ATLAS Production System 1/2

- Manages ATLAS simulation (full chain) and reprocessing jobs on the wLCG
 - Task request interface to define a related group of jobs
 - Input : DQ2 dataset(s) (with the exception of some event generation)
 - Output : DQ2 dataset(s) (the jobs are done only when the output is at the Tier-1)
 - Due to temporary site problems, jobs are allowed several attempts
 - Job definition and attempt state are stored in Production Database (Oracle DB)
 - Jobs are supervised by ATLAS Production System
- Consists of many components
 - DDM/DQ2 for data management
 - PanDA task request interface and job definitions
 - PanDA for job supervision
 - ATLAS Dashboard and PanDA monitor for monitoring
 - Grid Middlewares
 - ATLAS software

DQ2 - ATLAS Distributed Data Management Software
PanDA - the core part of ATLAS Production System



Introduction : ATLAS Production System 2/2



Task request interface

Job brokering is done by the PanDA Service (bamboo) according to input data and site availability

JOBLOGS CLOB(4000)	JOBPARS CLOB(4000)
<pre> <joblog> <stream> <stream> <stream> <fileinto> <filelog> <filelog> <filelog> <filelog> <dataset> <dataset> <dataset> <dataset> </pre>	<pre> <jobpar> <name> <input> <output> <position> <type> <LFN> <type> <metatype> <inputLFN> <metatype> <value> <mc08_100000_McANoZee_1Leptron.simul.EWAT.e384_s462_s1d04016.EWAT.D04034_00001.pool.root<value> <jobactuator> <jobactuator> <name> <outputfile> <name> <position> <type> <LFN> <type> <metatype> <outputLFN> <metatype> <value> <mc08_100000_McANoZee_1Leptron.simul.HITS.e384_s462_s1d04016.HITS.D04016_00001.pool.root<value> <jobactuator> <jobactuator> <name> <inputevents> <name> <position> <type> <natural> <type> <metatype> <plain> <metatype> <value> <50<value> <jobactuator> <jobactuator> <name> <inputevents> <name> <position> <position> <type> <natural> <type> <metatype> <plain> <metatype> <value> <50<value> <jobactuator> <jobactuator> <jobactuator> <name> <randomseed> <name> <position> <type> <natural> <type> <metatype> <plain> <metatype> <value> <1<value> <jobactuator> <jobactuator> <name> <geometry> <name> <position> <position> <position> <type> <string> <type> <metatype> <plain> <metatype> <value> <ATLAS-GE0-02-01-00<value> <jobactuator> <jobactuator> <name> <physlist> <name> <position> <type> <string> <type> <metatype> <plain> <metatype> <value> <QOS.P_BERT<value> <jobactuator> <jobactuator> <name> <jobconfig> <name> <position> <position> <type> <string> <type> <metatype> <plain> <metatype> <value> <VertePos.py<value> <jobactuator> <jobactuator> <name> <release> <name> <position> <position> <type> <LFN> <type> <metatype> <inputLFN> <metatype> <value> <ds0_000001_Atlas Ideal.DBRelease.v00001.DBRelease-5.5.1.tar.gz<value> <jobactuator> <jobactuator> <name> <conditiontag> <name> <position> <position> <type> <string> <type> <metatype> <plain> <metatype> <value> <NONE<value> <jobactuator> <jobactuator> </pre>

Production Database: job definition, job states, metadata



Task states

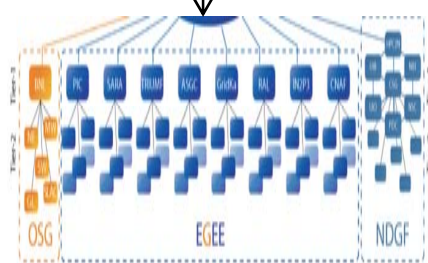


```

[read@charged dulcinea]$ ./taskinputs.py 43044
----- task 43044 -----
Dataset data08_cosmag.00088459.physics_LLCalo.daq.RAW.o4 has 359 files
Incomplete:
Complete : RAL-LCG2_DATATAPE,CERN-PROD_DAQ
No subscriptions

```

Tasks Input: DQ2 datasets



3 Grids/10 Clouds/90+Production Sites

```

[read@charged dulcinea]$ dq2-list-dataset \*tid043169
validl.107406.singlepart_singlelep17.recon.ESD.e380_s513_r634_tid043169
validl.107406.singlepart_singlelep17.recon.RD0.e380_s513_r634_tid043169
validl.107406.singlepart_singlelep17.recon.AOD.e380_s513_r634_tid043169
validl.107406.singlepart_singlelep17.recon.log.e380_s513_r634_tid043169

```

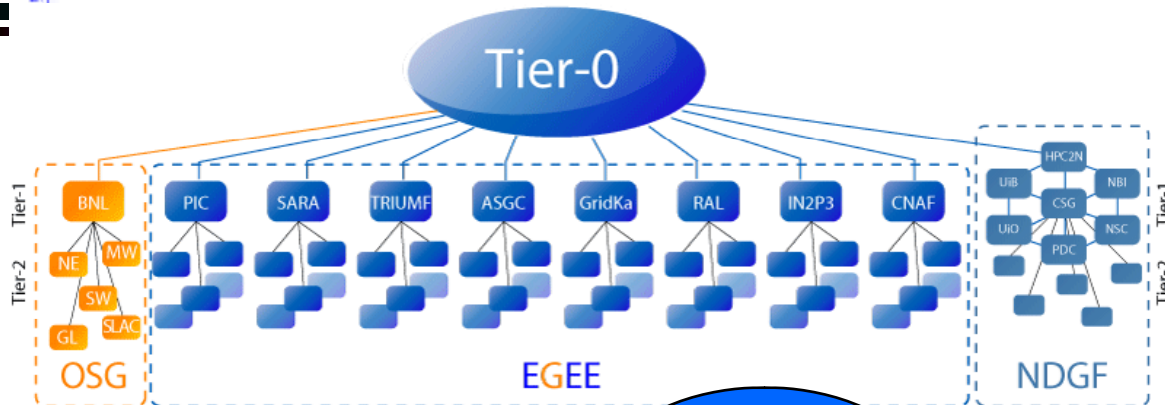
Tasks Output: DQ2 datasets



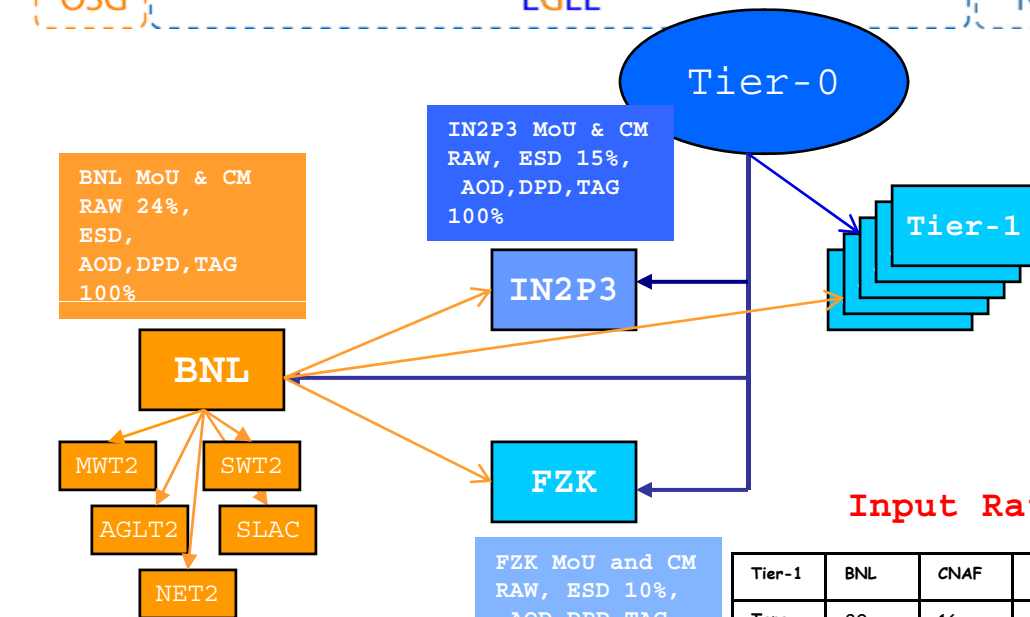
Monitor sites, tasks, jobs



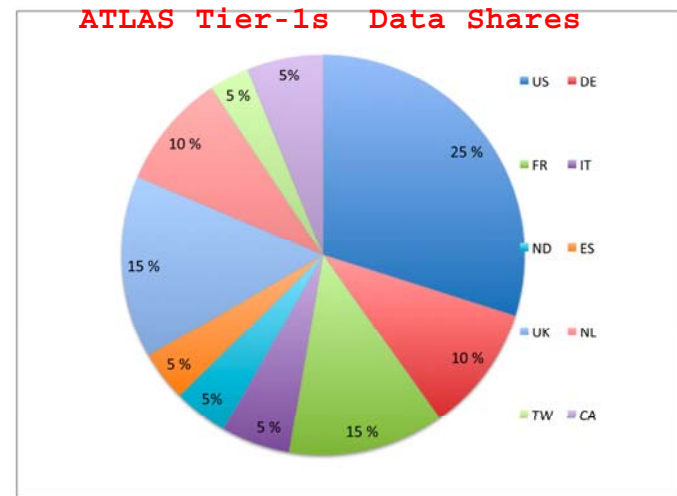
ATLAS Grid Sites and Data Distribution



3 Grids, 10 Tier-1s, ~70 Tier-2(3)s
 Tier-1 and associated Tier-ns form cloud. ATLAS clouds have from 2 to 15 sites. We also have T1-T1 associations.



→ Data export from CERN
 → reProcessed data distribution



Input Rates Estimation (Tier-1s)

Tier-1	BNL	CNAF	FZK	IN2P3	NDGF	PIC	RAL	SARA	TAIWAN	TRIUMF	Summary
Tape (MB/s)	80	16	32	48	16	16	32	48	16	16	320
Disk (MB/s)	240	60	80	100	60	60	80	100	60	60	800
Total (MB/s)	320	76	112	148	76	76	112	148	76	76	1220



Introduction : Data Processing Cycle

- Data processing at CERN (Tier-0 processing)
 - First-pass processing of the primary event stream
 - The derived datasets (ESD, AOD, DPD, TAG) are distributed from the Tier-0 to the Tier-1s
 - RAW data (received from Event Filter Farm) are exported within 24h. This is why first-pass processing can be done by Tier-1s (though this facility was not used during LHC beam and cosmic ray runs)
- Data reprocessing at Tier-1s
 - 10 Tier-1 centers world wide. Each takes a subset of RAW data (Tier-1 shares from 5% to 25%), ATLAS production facilities at CERN can be used in case of emergency.
 - Each Tier-1 reprocessed its share of RAW data. The derived datasets are distributed ATLAS-wide.

Incomplete list of Data Formats:
ESD : Event Summary Data
AOD : Analysis Object Data
DPD : Derived Physics Data
TAG : event meta-information



Reprocessing



- ATLAS collected cosmic ray data in Aug-Nov08 and single beam data in September 2008.
- First reprocessing round was completed in Dec08-Jan09, the second one is just started.



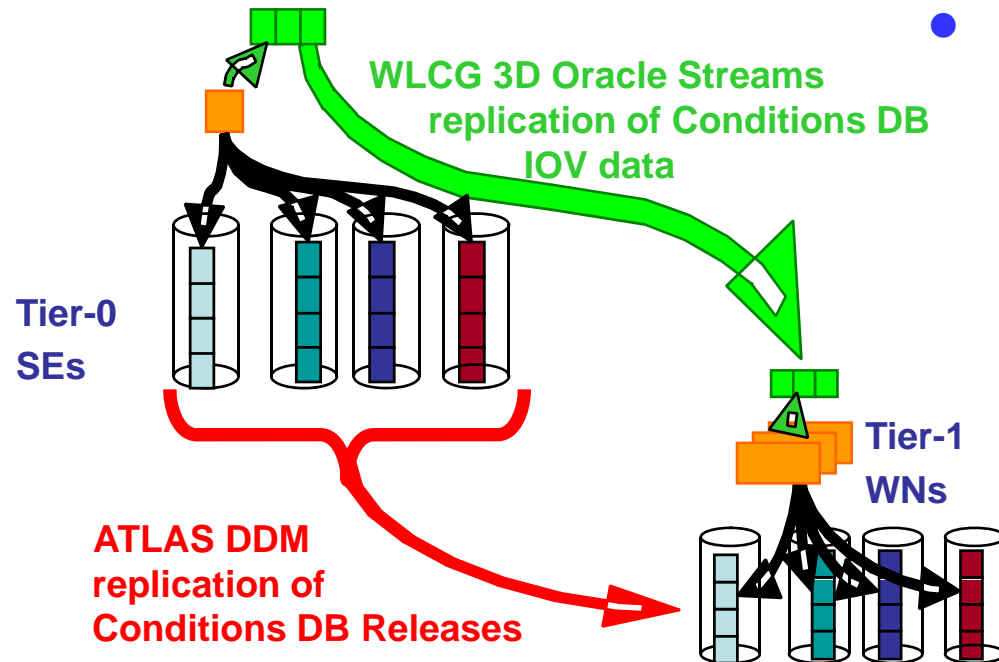
Preparations: Conditions DB Scalability

- In reprocessing on the Grid, instabilities and problems at Tier-1 sites may result in peak database access loads when many jobs are starting at once
 - Peak database loads can be much higher than average access rates
- In preparation for reprocessing, ATLAS Conditions DB scalability tests were increased both in scope and complexity, which allowed the identification and resolution of problems in time for reprocessing.
 - By simulating realistic workflow, ATLAS Conditions DB scalability tests produced Oracle overload conditions at all five Tier-1 sites tested
 - During the overload, the continuous Oracle Streams update of ATLAS Conditions DB data to this Tier-1 site degraded
 - After several hours, this Oracle overload at one Tier-1 site degraded Oracle Streams updates to all other Tier-1 sites
 - This situation has to be avoided
- To assure robust production operations in reprocessing we minimized the number of queries made to Oracle database replicas by taking full advantage of ATLAS technology-independent data access architecture



Scalable Conditions DB Access

- ATLAS reprocessing jobs accessed Conditions DB data in Tier-1 Oracle replicas, in SQLite replicas and in POOL Conditions DB payload files
- Minimization of Oracle access improved robustness of remote database access, which is critical for reprocessing on the distributed NDGF Tier-1 and US ATLAS Tier-2 sites
 - Robust Oracle access effectively doubled the reprocessing capacity at BNL Tier-1

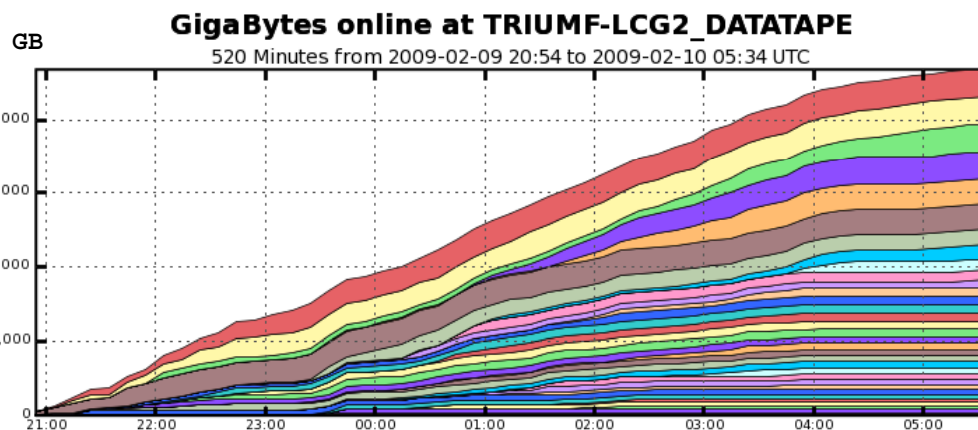


- By taking advantage of the organized nature of scheduled reprocessing, our Conditions DB access strategy leaves the Oracle servers free for 'chaotic' user-driven database-intensive tasks, such as calibration/alignment, detector performance studies and physics analysis



Preparations : Data staging test

- Test bulk recall of data from tape by
 - Using ATLAS Distributed Data Management staging service
 - Ask for 35 datasets comprising 9TB of data in 3k files
 - Target rate for a 10% Tier-1 is 186MB/s



■ fdr08_run2.0052293.physics_Jet.daq.RAW.o3 (677.43)	■ fdr08_run2.0052280.physics_Jet.daq.RAW.o3 (745.86)
■ fdr08_run2.0052290.physics_Muon.daq.RAW.o3 (422.45)	■ fdr08_run2.0052301.physics_Jet.daq.RAW.o3 (750.07)
■ fdr08_run2.0052300.physics_Jet.daq.RAW.o3 (298.08)	■ fdr08_run2.0052293.physics_Bphys.daq.RAW.o3 (232.68)
■ fdr08_run2.0052290.physics_Bphys.daq.RAW.o3 (256.63)	■ fdr08_run2.0052293.physics_Minbias.daq.RAW.o3 (131.23)
■ fdr08_run2.0052304.physics_Jet.daq.RAW.o3 (745.86)	■ fdr08_run2.0052300.physics_Minbias.daq.RAW.o3 (118.21)
■ fdr08_run2.0052283.physics_Jet.daq.RAW.o3 (750.07)	■ fdr08_run2.0052304.physics_Bphys.daq.RAW.o3 (221.80)
■ fdr08_run2.0052301.physics_Muon.daq.RAW.o3 (165.11)	■ fdr08_run2.0052293.physics_Egamma.daq.RAW.o3 (84.40)
■ fdr08_run2.0052290.physics_Jet.daq.RAW.o3 (684.99)	■ fdr08_run2.0052283.physics_Minbias.daq.RAW.o3 (162.84)
■ fdr08_run2.0052280.physics_Bphys.daq.RAW.o3 (221.80)	■ fdr08_run2.0052293.physics_Muon.daq.RAW.o3 (386.78)
■ fdr08_run2.0052301.physics_Egamma.daq.RAW.o3 (98.46)	... plus 16 more

Total: 9,371 , Average Rate: 0.30 /s

- Many Problems Understood
 - Poor performance between SE and MSS systems
 - Stuck tapes leading to files being unavailable
 - Load problems on SRM servers

9371 GB staged in 520 mins -> 300 MB/sec



Before Reprocessing : Site Validation

- Site Validation procedure

- Objective was to validate that all Tier-1 and Tier-2 sites produce identical outputs from the same inputs
 - The computing model envisaged reprocessing at Tier-1s only
 - Some clouds have asked to use their Tier-2s as well, and additionally the Operations team is keen to have spare reprocessing capacity
- Validation procedure :
 - Reconstruct representative files (which are a mixture of streams) at all reprocessing sites
 - Dump the numerical contents of ESD, AOD files into plain text (using Athena)
 - Compare text files and check for line-by-line identity
 - Perform final round of validation (signed off by data quality experts) with a single run/stream/dataset processed exactly as for the reprocessing

The Golden Rule:

Reprocessing may only run at sites which have been validated



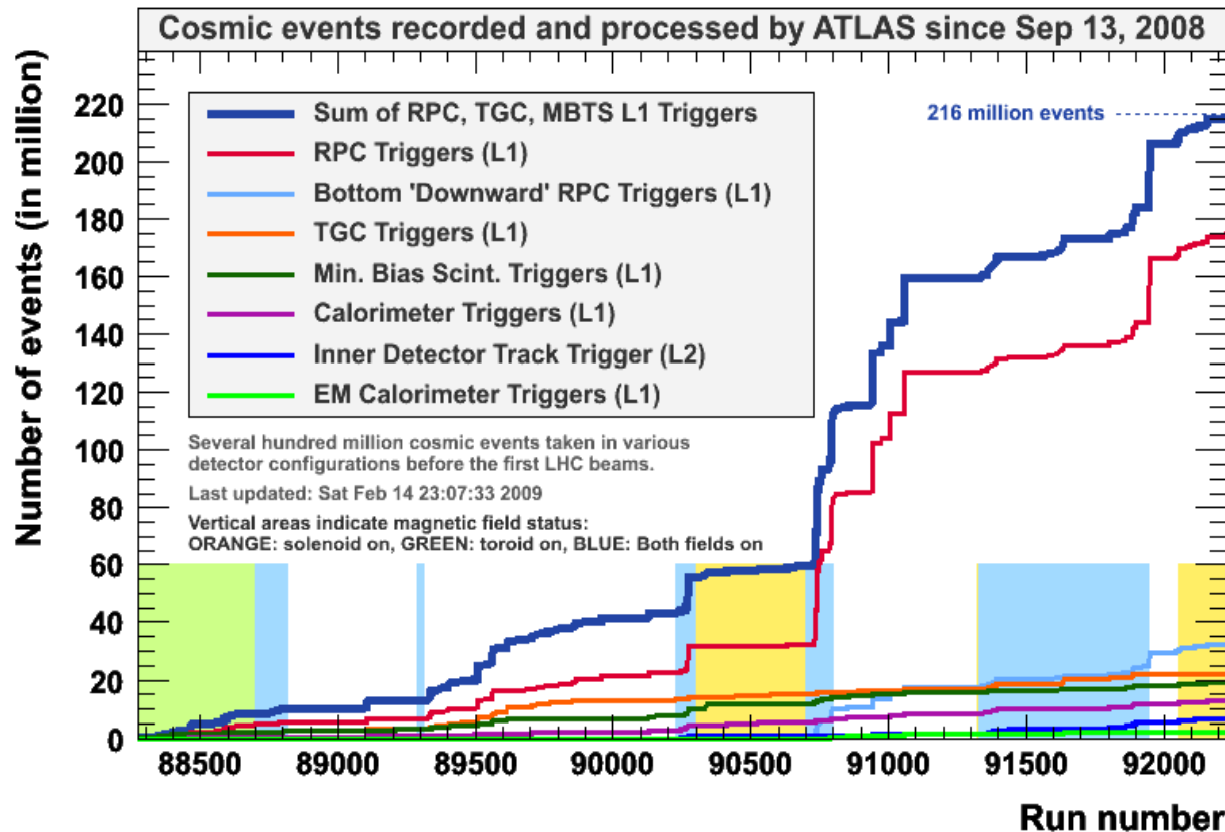
Reprocessing : Run/Stream List and Data Volume



Data volume: **284 million events** in ...

■ ... 127 runs, 11 streams, 1973 datasets, 330559 files → **513 TB of raw data**

■ Output volume: 8321 containers, 2 164 416 files, **110 TB** (all formats, w/o multiple replicas)

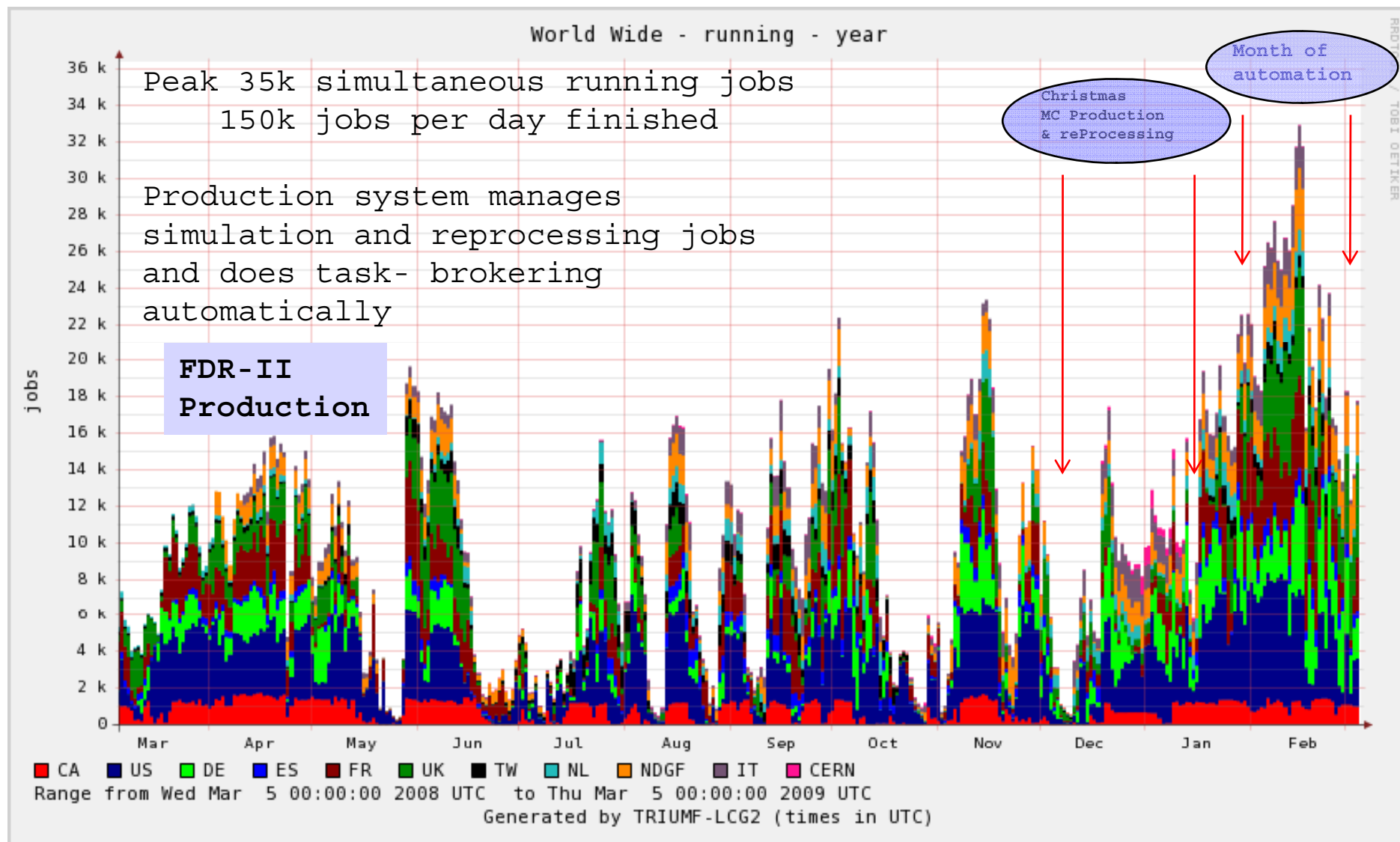


Raw Data distribution between ATLAS Tier-1s, 91% data are disk resident

Tier-1	Datasets (TB)@TAPE	Datasets (TB)@DISK
BNL	648 (89.5)	609 (124.2)
CNAF	85 (21.5)	81 (20.3)
FZK	200 (57.6)	186 (40.3)
IN2P3	269 (61.3)	247 (45.3)
NDGF	93 (14.6)	88 (13.2)
PIC	88 (17.1)	84 (16.4)
RAL	181 (69.0)	159 (52.7)
SARA	237 (47.4)	207 (37.0)
TAIWAN	99 (20.4)	94 (18.5)
TRIUMF	73 (14.6)	2 (0.6)



Production Jobs. Y2008-2009





Reprocessing : Running Jobs



Reprocessing scheme (2060 tasks of types 1 – 3, 630 of type 4)

1. Reprocess RAW data files, produce : EventSummaryData, DataQuality Histograms, Calibration NTUPles and TAG files
2. Produce Analysis Object Data, Derived Physics Data files
3. Merge DQ histograms files
4. Merge Calibration NTUPles files

Job brokering is done by the PanDA Service (bamboo) according to input data and site availability. When a job is defined, it knows which files are on tape and the Production System triggers file pre-staging in these cases.

Job statistic (1 job = 1 file)

T1	CA	CERN	ES	FR	IT	NG	NL	UK	US	Sum
Total Jobs	20707	26348	364	48288	13619	12561	23472	54360	128764	329609
Done Jobs	20150	26015	364	46937	13018	12281	23167	51344	124667	317943
Fraction [%%]	97.3	94.7	100.	97.2	95.6	97.8	98.7	94.5	96.8	96.5
Aborted jobs	557	1459	0	1351	601	280	305	3016	4097	11666
Fraction [%%]	2.7	5.3	0	2.8	4.4	2.2	1.3	5.5	3.2	3.5

Number of attempts per successful job, <av> 1.8

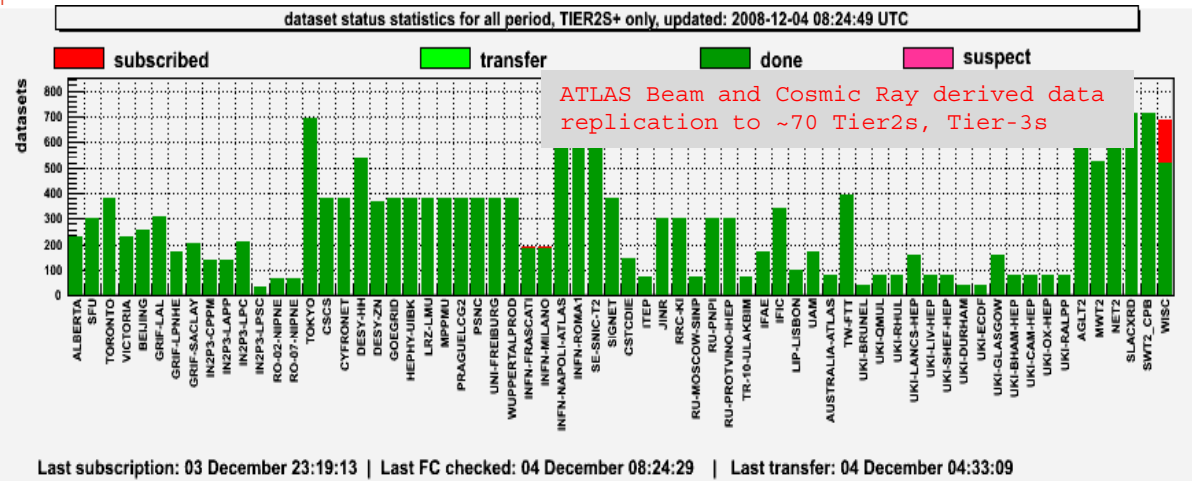
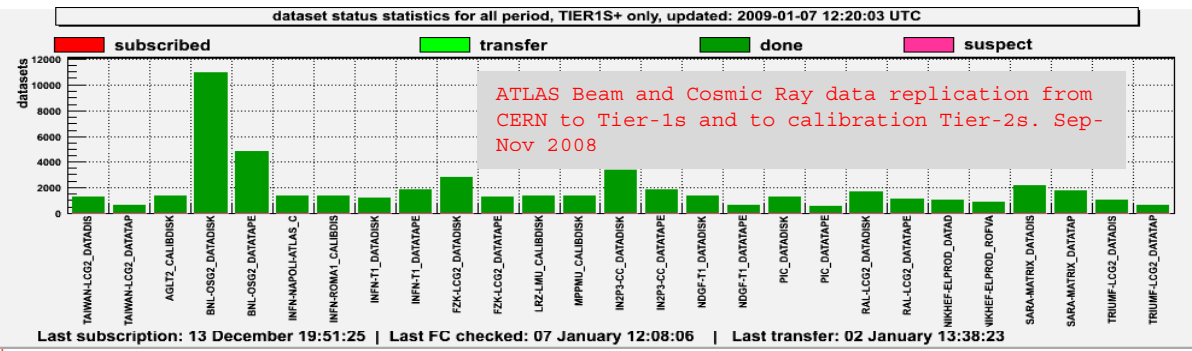
Tier-1	CA	FR	NDGF	CERN	ES	US	IT	UK	NL
#attempts	1.02	1.11	1.14	1.16	1.18	1.39	1.83	2.31	2.85



1 Beam, Cosmic Ray and Reprocessed Data Replication to ATLAS Tiers

- ESD : 2 replicas ATLAS wide (distributed between Tier-1s)
- AOD : 11 replicas ATLAS wide (consolidated at Tier-1s and CERN)
- DPD : 20+ replicas ATLAS wide (consolidated at CERN, Tier-1s and distributed within clouds)
- Calibration datasets replicated to 5 T2 calibration centers
- Data quality HIST datasets replicated to CERN

Cloud	Datasets	Total Files in datasets	Total CpFiles in datasets	Completed	Transfer	Subscribed
TO	8125	80361	79968	8104	0	21
CA	5478	74369	73984	5458	2	18
DE	5478	74369	73977	5458	0	20
ES	5478	74369	74347	5458	20	0
FR	5478	74369	73977	5458	0	20
IT	5478	74369	74347	5458	20	0
NG	5478	74369	74347	5458	20	0
NL	5478	74369	73972	5453	0	25
TW	5478	74369	73977	5458	0	20
UK	5478	74369	73977	5458	0	20
US	8125	80361	79849	8095	0	30



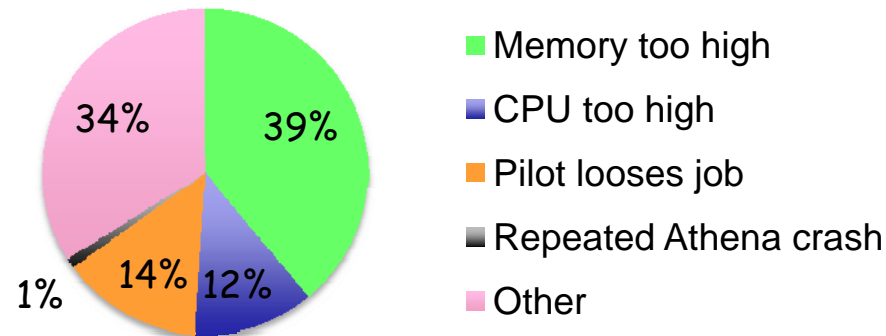
Reprocessed data replication status. 99+% were completely replicated to all Tier-1s



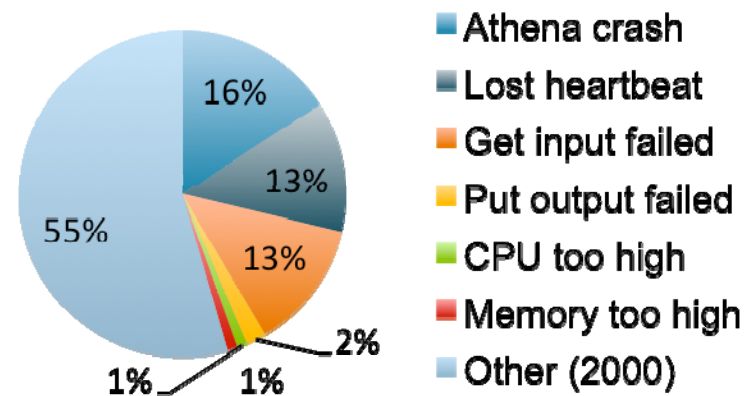
Reprocessing : Brief Error Analysis



Persistent errors – never succeeded (~25% of all errors)



Transient errors – job ultimately succeeded (~75% of all errors)



No single "main reason" but operational issues



Summary



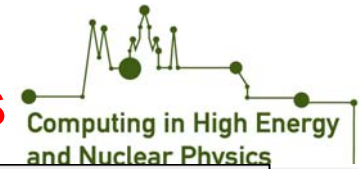
- The ATLAS Production System has been used successfully for LHC beam and Cosmic Ray data (re)processing
- The Production System handled the expected data volume robustly
- ATLAS Distributed Data Management System is robust and detector data as well as reprocessing results are distributed to sites and physics team in a timely manner
- Issues with conditions data and database access were understood and technical solutions found. There is no scalability limit foreseen for database access.
- Data staging was exercised on a 10% scale and reprocessing using bulk (0.5PB) data staging is in progress.
- Grid vs off-Grid data processing issues need more testing
- The second round of reprocessing is started in March and our target is to reprocess 100% of events. We have all machinery to do it.



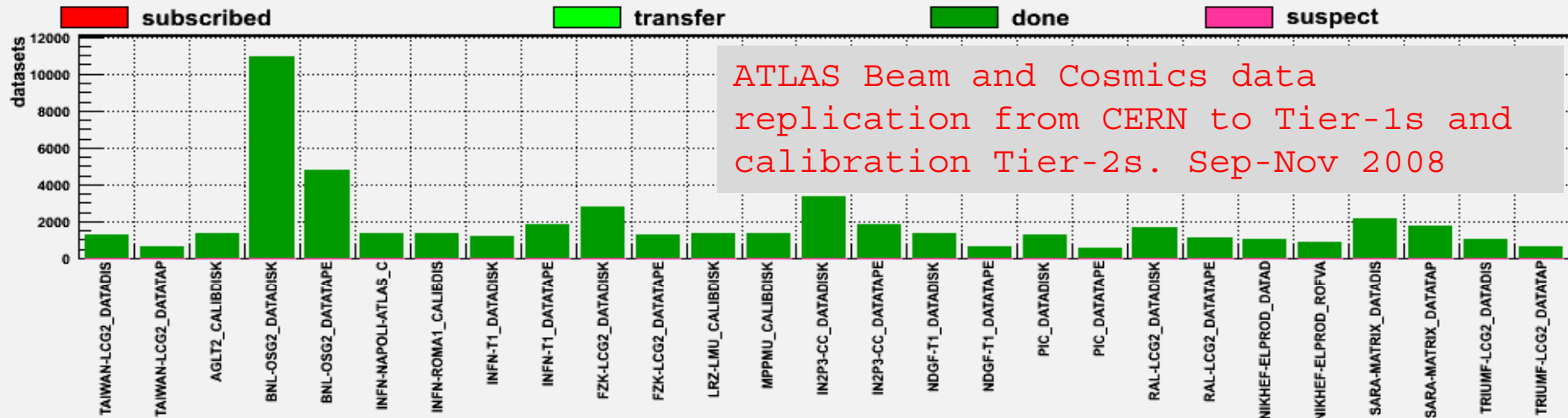
● BACKUP SLIDES



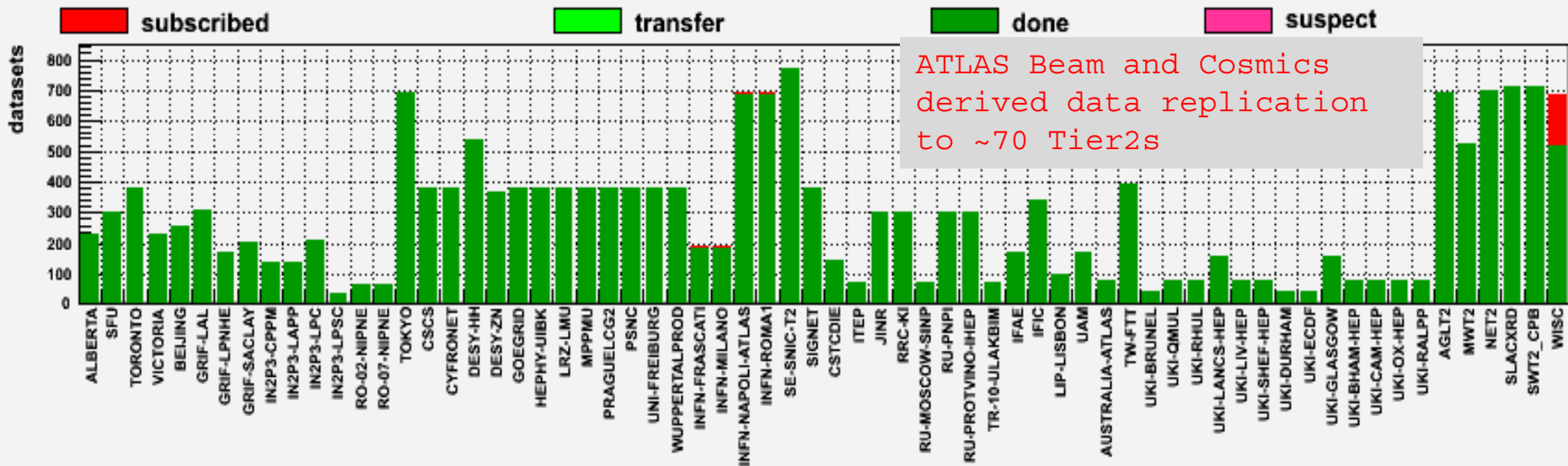
1Beam And Cosmics Data Replication To Tiers



dataset status statistics for all period, TIER1S+ only, updated: 2009-01-07 12:20:03 UTC



Last subscription: 13 December 19:51:25 | Last FC checked: 07 January 12:08:06 | Last transfer: 02 January 13:38:23



Last subscription: 03 December 23:19:13 | Last FC checked: 04 December 08:24:29 | Last transfer: 04 December 04:33:09



Dealing with persistently failing events

- Some events never reprocess
 - 3.5% of all events in last reprocessing
 - 1 failed event = all events in RAW file are not reprocessed = 1 complete luminosity block for that stream not reprocessed (with collisions)
- Generally a failed event will need *new software* to reprocess it
 - After the main campaign, we must re-run all failed files to get to a situation where 100% of events are reprocessed
 - Once finally done these events will be appended to the existing run x stream container, as a final dataset

*Machinery is ready and it will be tested during March09
reprocessing campaign*



Related ATLAS Talks



● Software Components :

- P.Nevski : Knowledge Management System for ATLAS Scalable Task Processing on the Grid
- R.Walker : Advanced Technologies for Scalable ATLAS Conditions Database Access on the Grid

● Grid Middleware and Networking Technologies

- R.Rocha : The ATLAS Distributed Data Management Dashboard
- S.Campana : Experience Commissioning the ATLAS Distributed Data Management system on top of the WLCG Service
- G.Stewart : Migration of ATLAS PanDA to CERN

● Distributed Processing and Analysis

- G.Negri : The ATLAS Tier-0: Overview and Operational Experience
- B.Gaidiouz : Monitoring the ATLAS distributed production