



# Distributed processing and analysis of ATLAS experimental data

Dario Barberis

(Genoa University/INFN)

On behalf of the ATLAS Collaboration

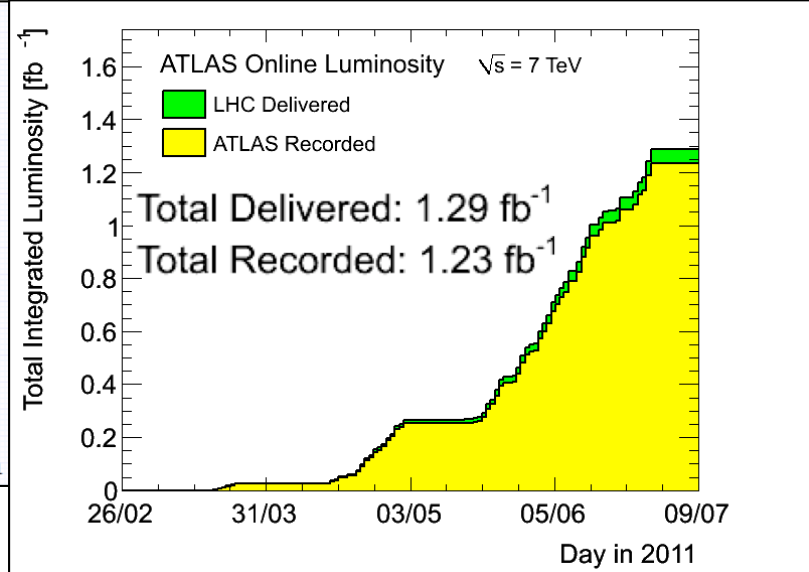
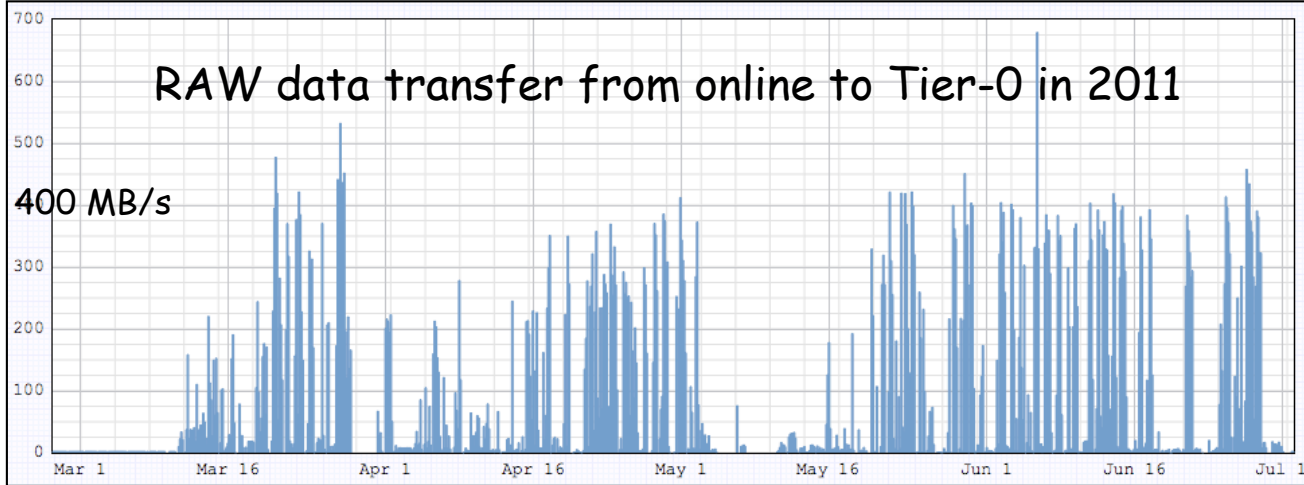


# Overview

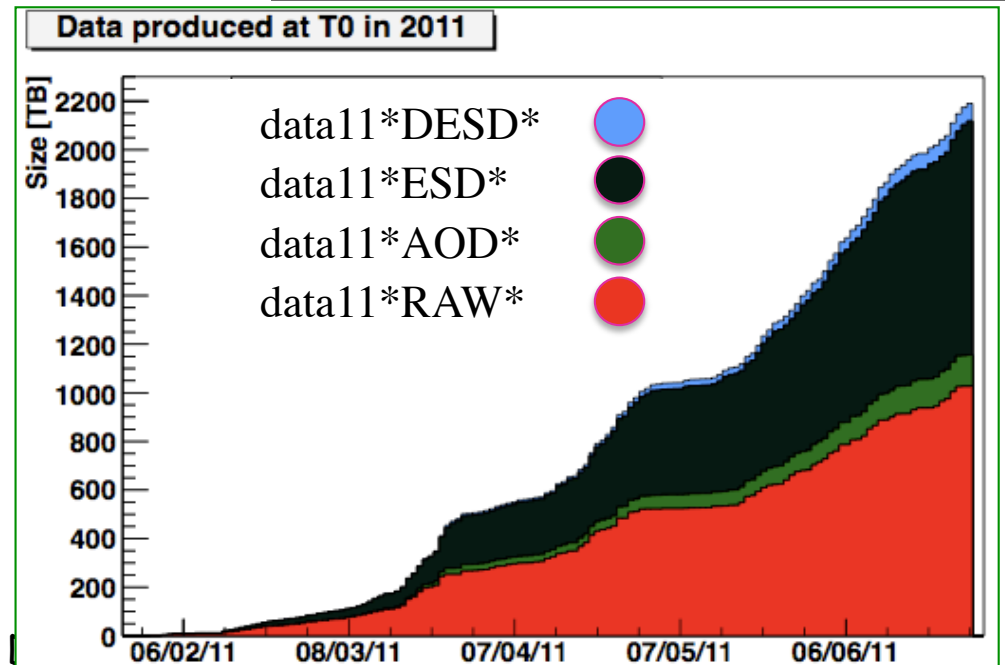
- Data collection
- Tier-0 processing
  - Fast calibration loop
  - Event reconstruction
  - Data export
- Some of the key distributed computing (Grid) technologies:
  - Data management and distribution with DDM/DQ2
  - Workload management with Panda
    - Re-processing campaigns and simulation production
    - Distributed analysis
  - Conditions Databases with Frontier
- Evolution of the computing model
- Outlook

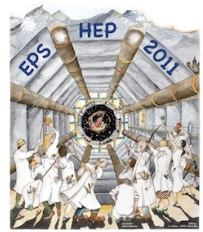


# Data taking in 2011



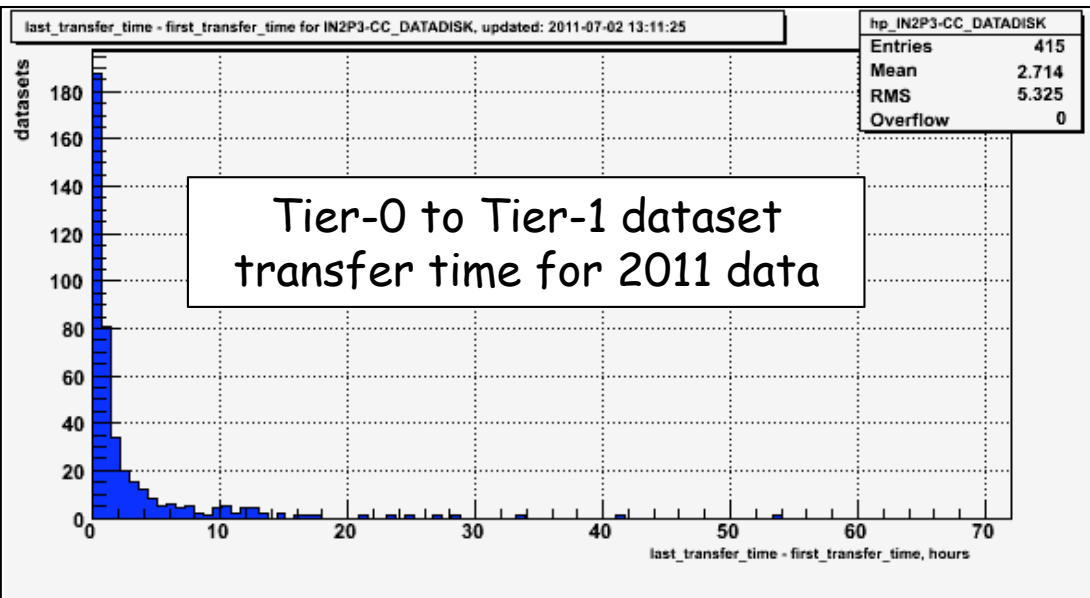
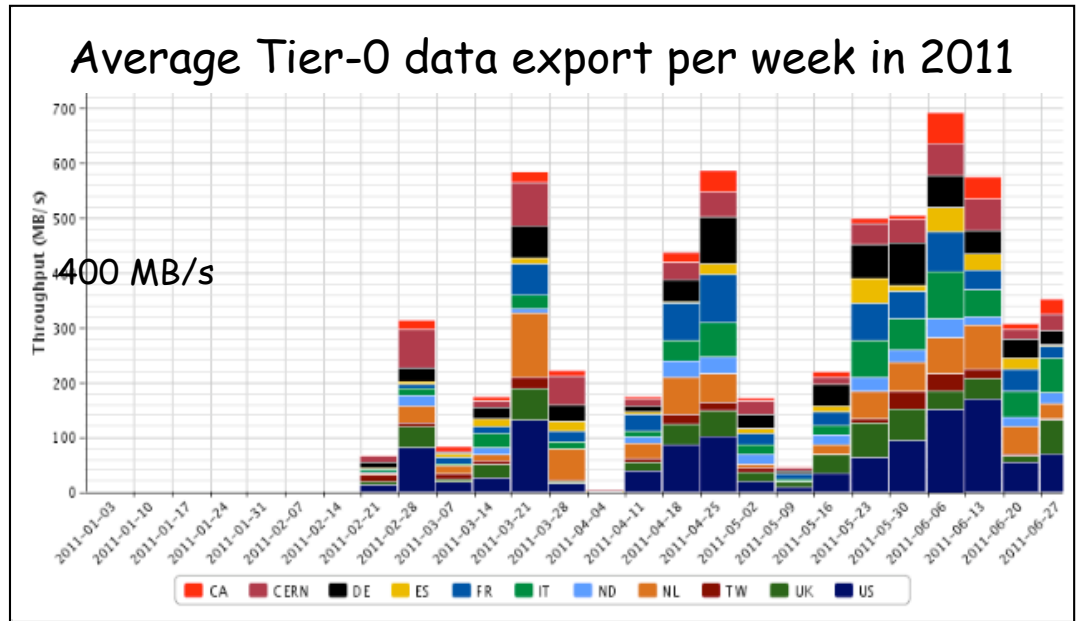
- We took until end June 2011 1 PB of RAW data. All data were:
  - Calibrated in real time (within 36 hours)
  - Reconstructed at Tier-0
  - Distributed on the Grid to 10 Tier-1 and many Tier-2 sites
- In total we produced 2.2 PB of distributed data





# Data distribution on the Grid

- Data export from Tier-0 to Tier-1s:
  - RAW: 1 primary copy (on disk) + 1 custodial copy (on tape)
  - ESD: 1 primary + 1 secondary copy (both on disk at different sites)
  - DESD: 2 primary copies
  - AOD: 2 primary + 1 secondary copy

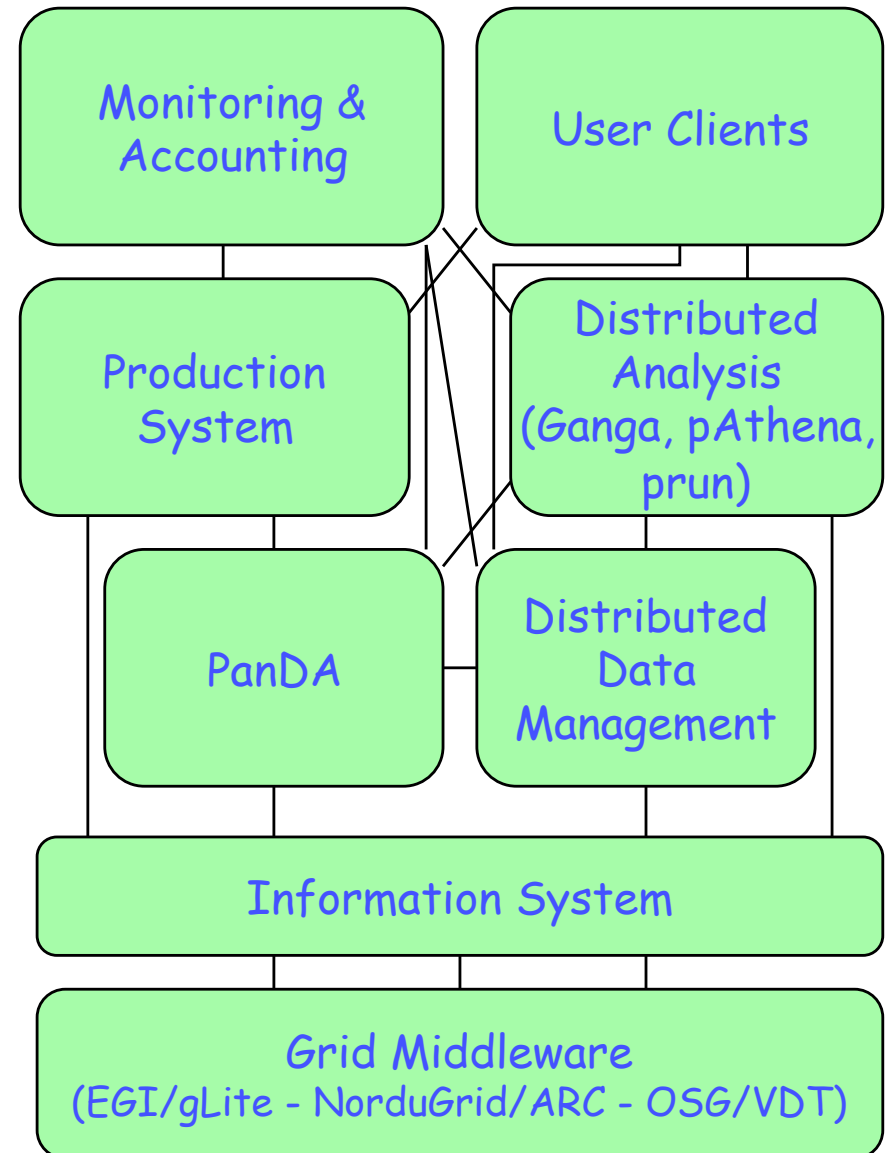


- Data are available for analysis in "almost-real" time. Example:
  - data11\_7TeV AOD distribution (to one specific Tier-1 but they are all similar):
  - on average 2.7 hours to complete the dataset



# ATLAS Grid Architecture

- ATLAS runs on 3 middleware suites:
  - gLite in most of Europe and several other countries (including all A-P countries)
  - ARC in Scandinavia and a few other small European countries
  - VDT in the USA
- ATLAS Grid tools interface with the middleware and shield the users from it
  - They also add a lot of functionality that is ATLAS specific
- The ATLAS Grid architecture is based on few main components:
  - Information system
  - Distributed data management (DDM)
  - Distributed production and analysis job management system (PanDA)
  - Distributed production (ProdSys) and analysis (Ganga/pAthena/prun) interfaces
  - Monitoring and Accounting tools
- DDM is the central link between all components
  - As data access is needed for any processing and analysis step!





# Distributed data management: DDM/DQ2

- The Distributed Data Management (DDM) architecture is implemented in the DQ2 tools and additional services
- The unit of storage and transfer is the dataset:
  - A dataset contains all files with statistically equivalent events
- DDM takes care of:
  - Distributing data produced by Tier-0 to Tier-1s and Tier-2s
  - Distributing simulated and reprocessed data produced by Tier-1/2s
  - Distributing user and group datasets as requested
  - Managing data movement generated by production activities
  - Cataloguing datasets (files, sizes, locations etc.)
  - Providing usage information for each dataset replica
  - Deleting obsolete or unnecessary replicas of datasets from disk when disks are full
  - Providing end-users with client tools to operate on datasets (import/export/move etc)





# Distributed data management: DDM/DQ2

- Data are transferred around the world steadily at high rates (the Grid never sleeps!)

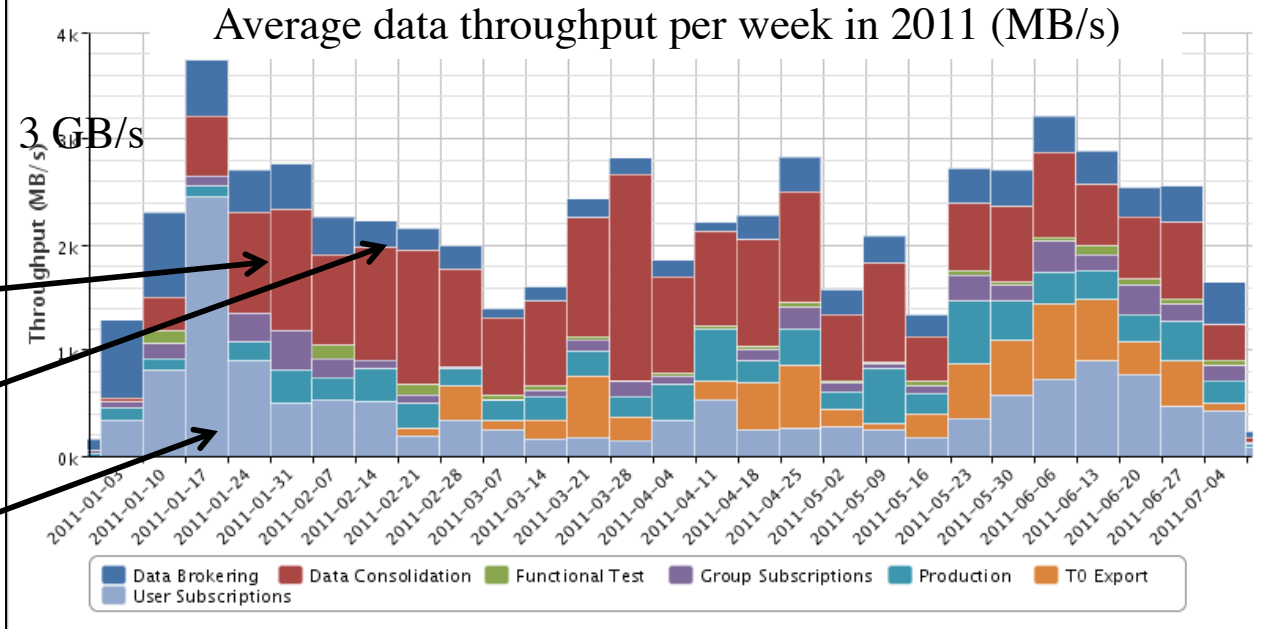
- Delicate balance between

■ Pre-placement

■ Dynamic data placement

➤ With automatic caching and cleaning

■ User requests



	TOTAL-	CA+	CERN+	DE+	ES+	FR+	IT+	ND+	NL+	TW+	UK+	US+
TOTAL-	93 % 2 GB/s	95 % 111 MB/s	92 % 372 MB/s	92 % 330 MB/s	91 % 98 MB/s	92 % 271 MB/s	90 % 118 MB/s	92 % 106 MB/s	88 % 169 MB/s	88 % 53 MB/s	93 % 211 MB/s	96 % 446 MB/s
CA+	97 % 118 MB/s	97 % 65 MB/s	95 % 15 MB/s	98 % 5 MB/s	95 % 2 MB/s	98 % 4 MB/s	99 % 3 MB/s	96 % 2 MB/s	98 % 5 MB/s	90 % 1 MB/s	95 % 4 MB/s	99 % 13 MB/s
CERN+	90 % 185 MB/s	90 % 6 MB/s	88 % 71 MB/s	90 % 20 MB/s	79 % 4 MB/s	95 % 15 MB/s	85 % 6 MB/s	90 % 9 MB/s	85 % 9 MB/s	91 % 7 MB/s	87 % 13 MB/s	95 % 27 MB/s
DE+	91 % 373 MB/s	88 % 7 MB/s	93 % 38 MB/s	93 % 218 MB/s	93 % 7 MB/s	94 % 15 MB/s	94 % 13 MB/s	90 % 8 MB/s	94 % 16 MB/s	89 % 5 MB/s	90 % 11 MB/s	93 % 34 MB/s
ES+	93 % 111 MB/s	94 % 2 MB/s	90 % 13 MB/s	93 % 6 MB/s	93 % 61 MB/s	92 % 5 MB/s	94 % 3 MB/s	95 % 3 MB/s	88 % 4 MB/s	85 % 2 MB/s	91 % 5 MB/s	94 % 8 MB/s
FR+	92 % 319 MB/s	86 % 6 MB/s	96 % 45 MB/s	93 % 17 MB/s	91 % 5 MB/s	91 % 164 MB/s	94 % 8 MB/s	97 % 9 MB/s	91 % 14 MB/s	88 % 5 MB/s	88 % 15 MB/s	96 % 31 MB/s
IT+	91 % 127 MB/s	92 % 3 MB/s	94 % 29 MB/s	95 % 7 MB/s	92 % 2 MB/s	91 % 6 MB/s	90 % 62 MB/s	96 % 3 MB/s	92 % 3 MB/s	79 % 2 MB/s	87 % 3 MB/s	94 % 8 MB/s
ND+	93 % 91 MB/s	97 % 1 MB/s	94 % 18 MB/s	97 % 4 MB/s	95 % 1 MB/s	95 % 3 MB/s	94 % 2 MB/s	92 % 46 MB/s	91 % 2 MB/s	94 % 722 kB/s	95 % 3 MB/s	96 % 9 MB/s
NL+	88 % 182 MB/s	81 % 3 MB/s	95 % 34 MB/s	91 % 12 MB/s	69 % 3 MB/s	90 % 7 MB/s	88 % 4 MB/s	89 % 5 MB/s	87 % 92 MB/s	88 % 2 MB/s	91 % 6 MB/s	90 % 14 MB/s
TW+	92 % 54 MB/s	93 % 1 MB/s	91 % 11 MB/s	94 % 4 MB/s	852 kB/s	97 % 4 MB/s	88 % 1 MB/s	91 % 2 MB/s	94 % 1 MB/s	90 % 22 MB/s	94 % 2 MB/s	92 % 4 MB/s
UK+	93 % 206 MB/s	86 % 3 MB/s	89 % 29 MB/s	91 % 7 MB/s	68 % 3 MB/s	91 % 8 MB/s	87 % 4 MB/s	91 % 3 MB/s	88 % 5 MB/s	90 % 2 MB/s	95 % 130 MB/s	91 % 14 MB/s
US+	95 % 517 MB/s	91 % 14 MB/s	95 % 68 MB/s	94 % 30 MB/s	90 % 9 MB/s	94 % 39 MB/s	94 % 14 MB/s	91 % 17 MB/s	91 % 17 MB/s	81 % 4 MB/s	93 % 21 MB/s	96 % 285 MB/s

- Excellent data transfer efficiency achieved overall

■ 93% average success rate in 2011

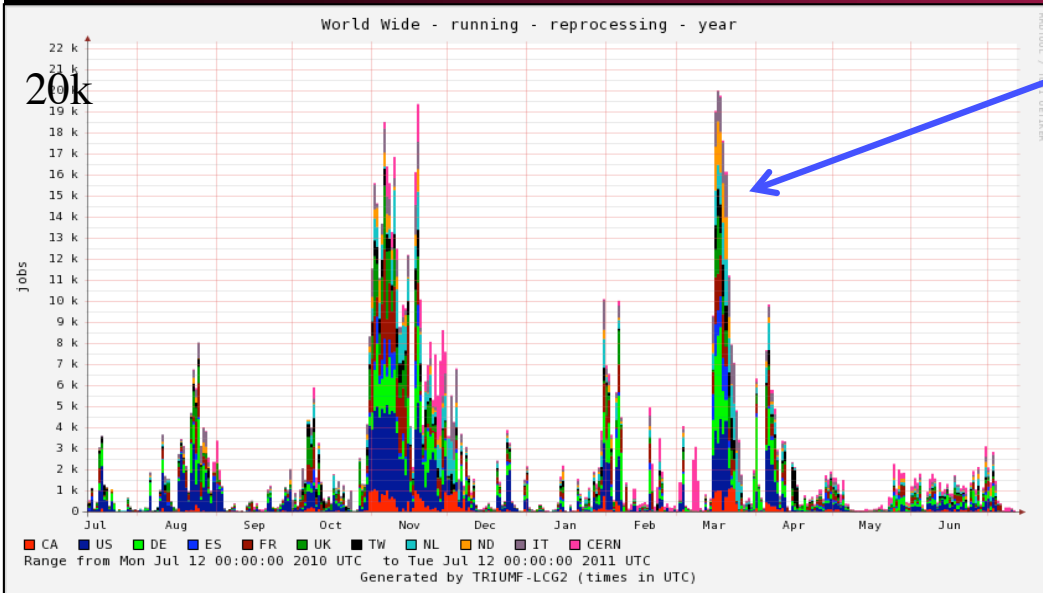
➤ First retry always succeeds

- Users can direct the outputs of their analysis jobs to their "home" on the Grid

■ Asynchronous transfer (plus retry) guarantees success in shortest time

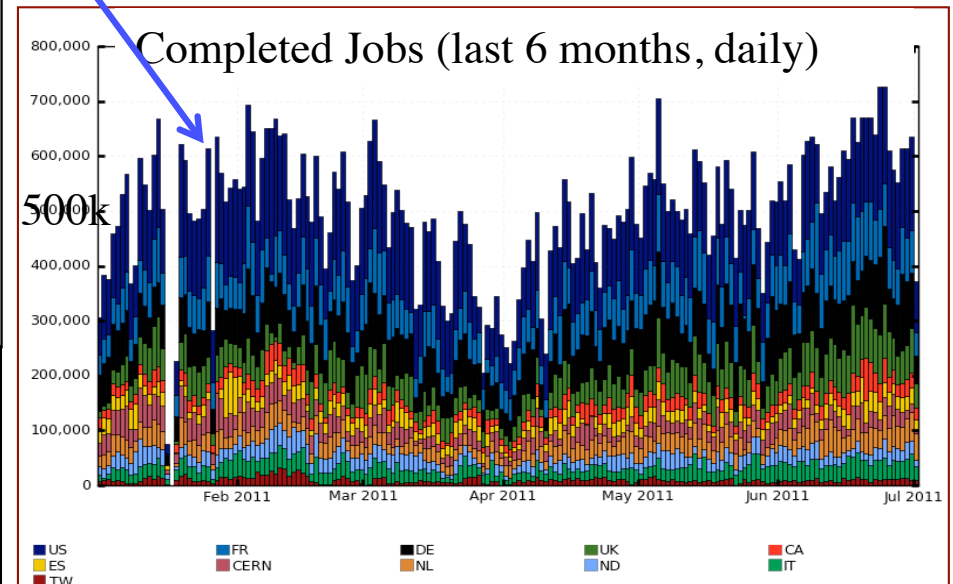


# Reprocessing and simulation production



One year of reprocessing campaigns

Over 500,000 simulation production and data reprocessing jobs/day on the Grid

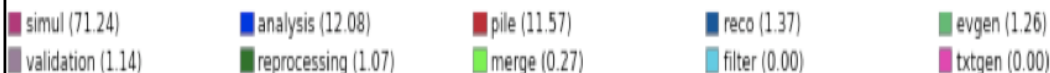
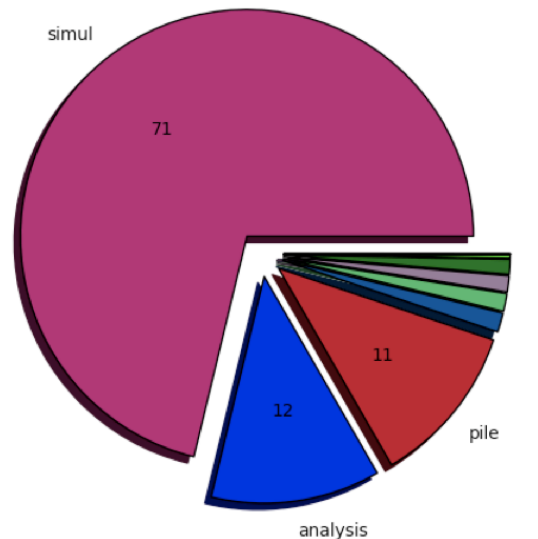


~80k jobs running simultaneously

Analysis tasks are 50% of the jobs but use 12% of total available CPU time

Re-run frequently to produce newer n-tuples

CPU share by activity  
(April-July 2011):  
86% simulation  
12% analysis  
1% validation  
1% reprocessing

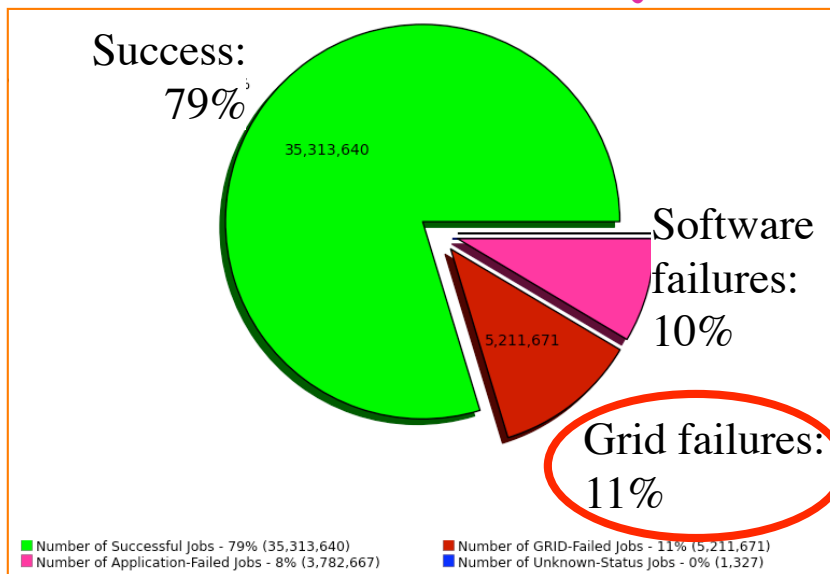
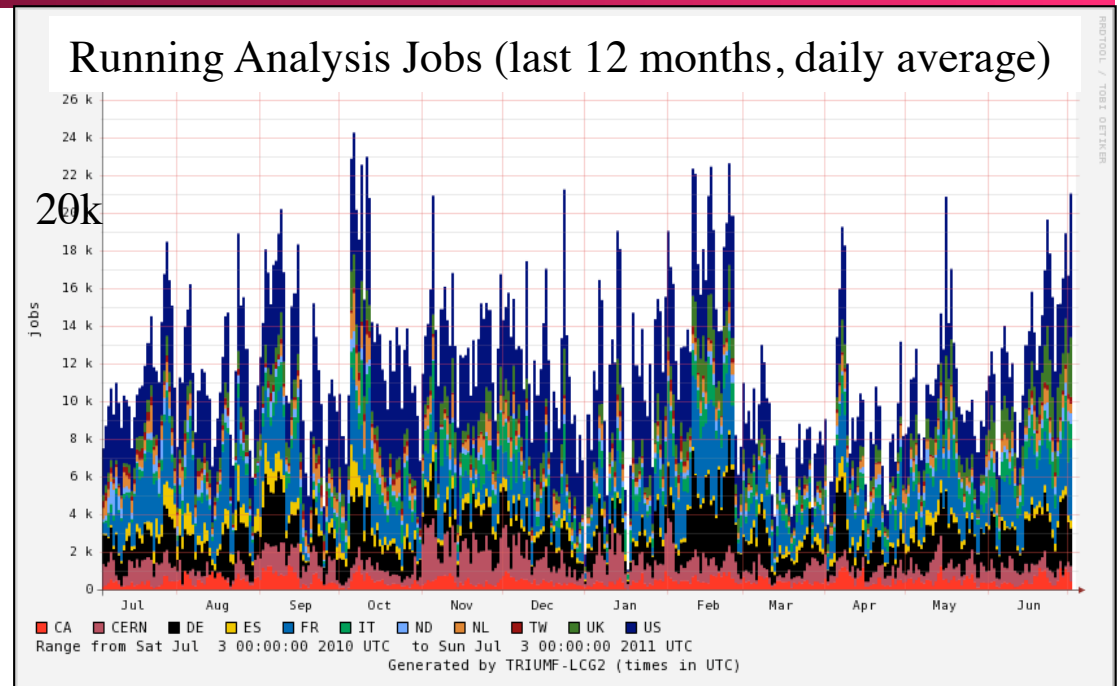






# Distributed analysis on the Grid

- Analysis jobs run world-wide
  - Jobs go to the data as much as possible
- Grid reliability issues...
  - automatic exclusion (and re-inclusion) of analysis queues that do not perform well, measured via automatic HammerCloud test jobs



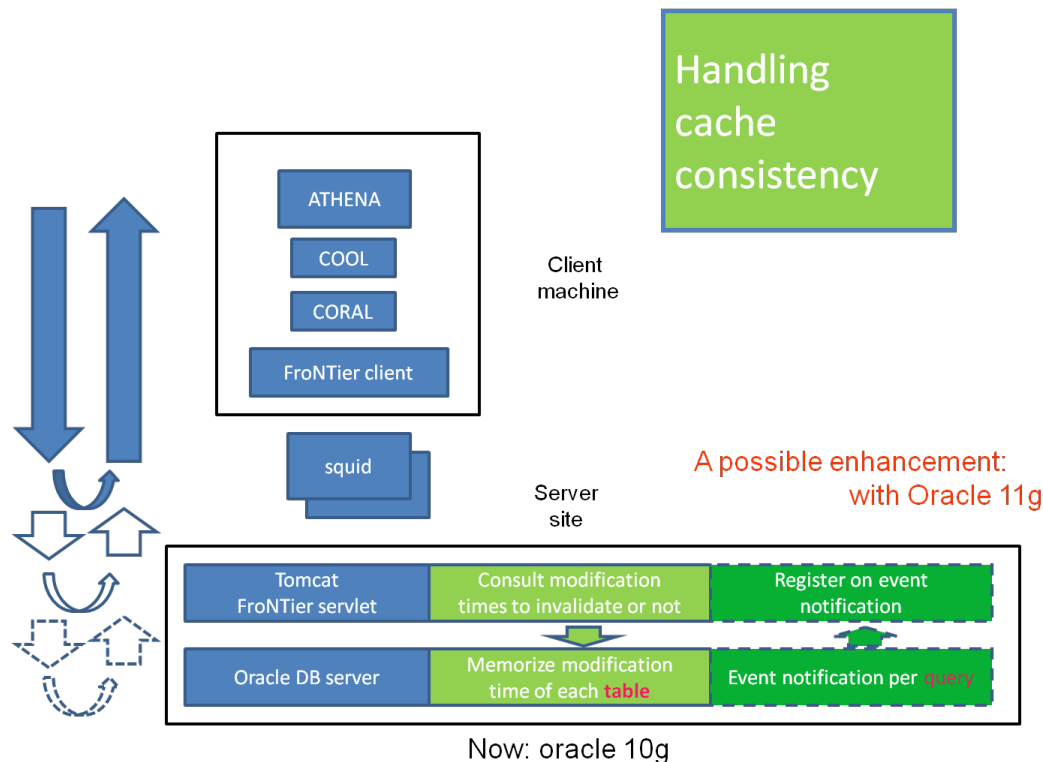
- Work in progress to improve task efficiency (and user happiness)
  - Merging of output files
  - Automatic retrieval of jobs that fail for well-defined Grid-related reasons
  - Improved analysis tasks book-keeping, to better keep track of the whole workflow



# Conditions Databases

- Frontier deployed in 2009 to enable distributed access to the conditions DB
- Flow of database data:
  - Oracle: CERN online -> CERN offline -> 3D (BNL, TRIUMF, RAL, KIT, IN2P3-CC)
  - Frontier server at each of the above sites connects to local Oracle database
  - Local Squid contacts nearest Frontier server
    - With failover to next-to-nearest

## Map of installed Squids



- Frontier reduces considerably the access time to DB data from remote sites
- It is particularly important for sites with low bandwidth and high latency towards Oracle servers



# Evolution of the Computing Model in 2011

- ✓ Break the cloud\* boundaries
  - Introduce flexibility in data distribution and job assignment
- ✓ Allow inter-cloud direct Tier-1 $\leftrightarrow$ Tier-2 and Tier-2 $\leftrightarrow$ Tier-2 transfers according to network connectivity
  - For data placement, user subscriptions and job I/O
- ✓ Allow job distribution from Tier-1s to Tier-2s in other clouds
  - Output files are then collected back to the original Tier-1 (of course)
- ✓ Reduce the number of data replicas to have more data on disk
- ✓ Introduce dynamic data replication and deletion based on dataset popularity
- ✓ Reduce the multiplicity of Oracle database servers and equip all remaining ones with Frontier web servers
- Integrate all 11 LFCs into a single catalogue at CERN (work in progress)
  - No longer one catalogue for each cloud
- Move towards using CVMFS (web-based file system) for software release and conditions data files distribution (tests in progress)
- ★ (An ATLAS Grid cloud includes a Tier-1 and all associated Tier-2/3s)



# Summary and Outlook

- The ATLAS Distributed Computing infrastructure is working thanks to many efforts in preparation and many people working in operations
- We are able to
  - Process, distribute, and reprocess the data
  - Analyse the data
  - Provide support to our large community
- As we get experience with *reality* we are looking at the evolution of the model and our implementations, e.g.
  - Less-strict cloud model
  - Better data distribution for analysis
  - Improved support for analysis