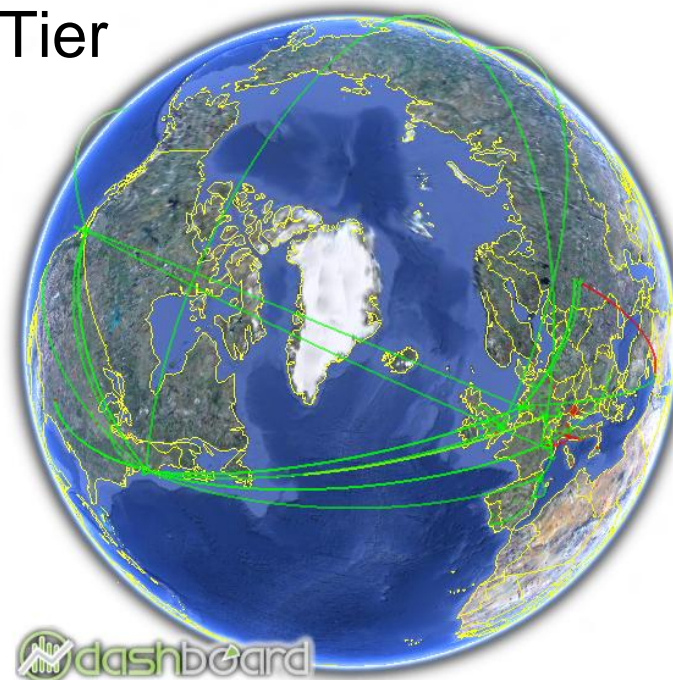


ATLAS Distributed Computing: First Experiences with Real Data

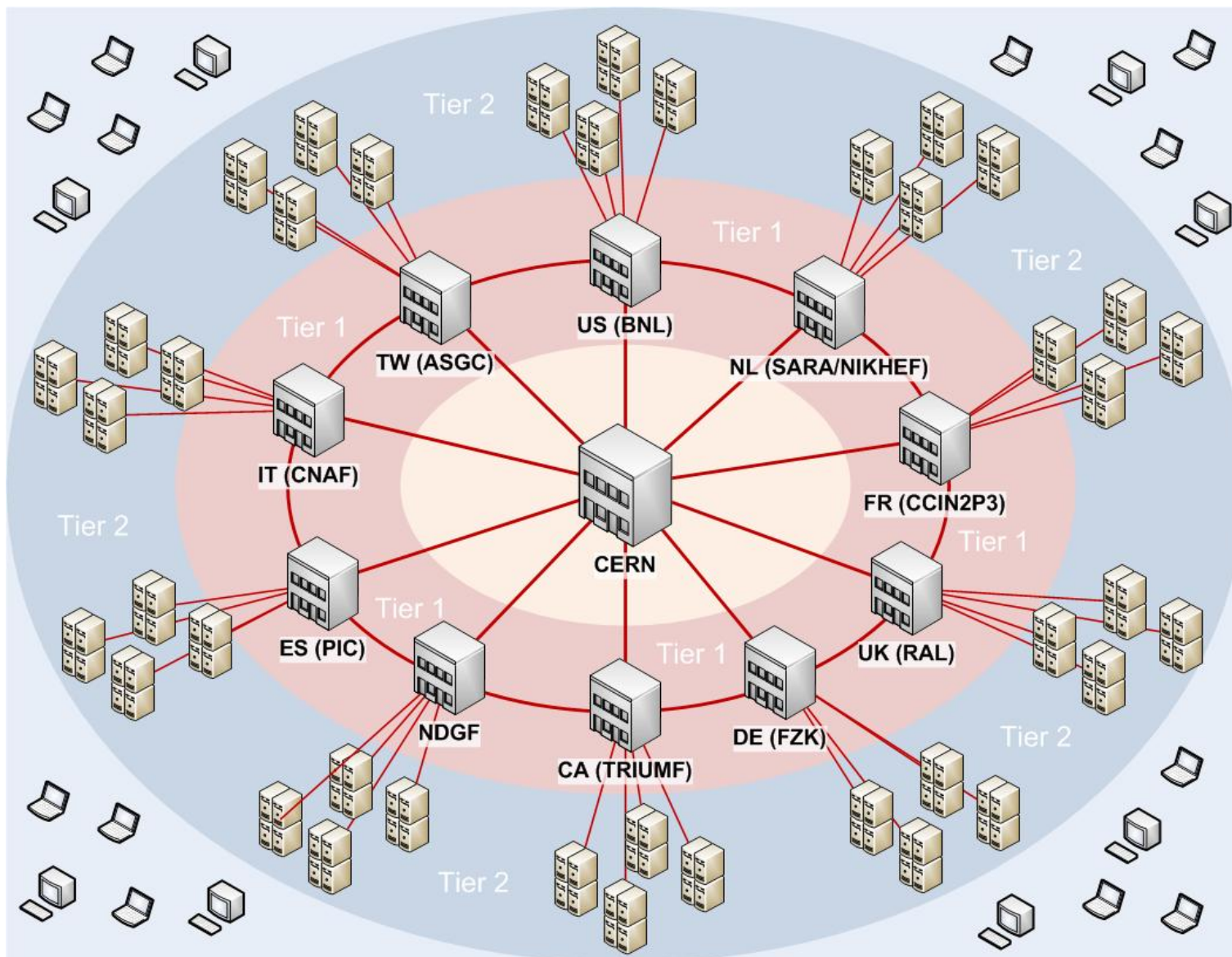
Dan van der Ster – CERN IT Experiment Support
on behalf of the ATLAS Collaboration

*12th Topical Seminar on Innovative Particle and Radiation
Detectors (IPRD10), 7-10 June 2010, Siena, Italy*

- Summary of the ATLAS Computing Model
- Some of the key technologies:
 - Data management and distribution with DQ2
 - Workload management with Panda
 - Conditions Databases with FroNTier
- Some of the key activities:
 - Grid Data Processing
 - Distributed User Analysis
 - User Support
 - Tier 3

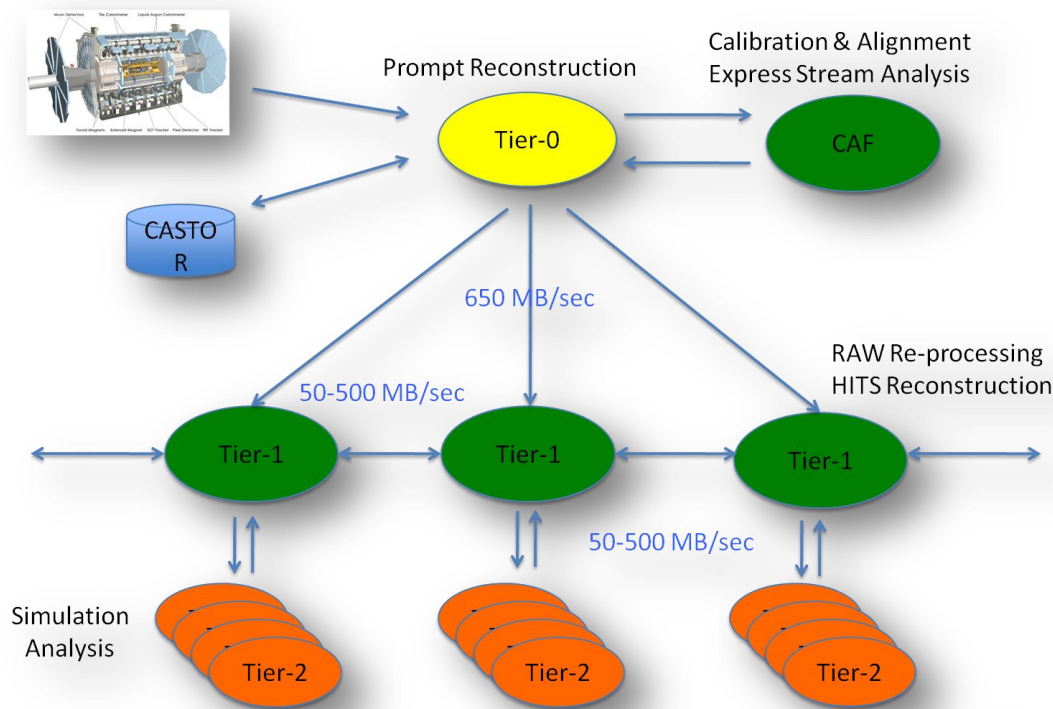


Acknowledgements: D.Barberis,
K.Bos, S.Campana , D.Front,
A.Klimentov, M.Lamanna, D.Smith



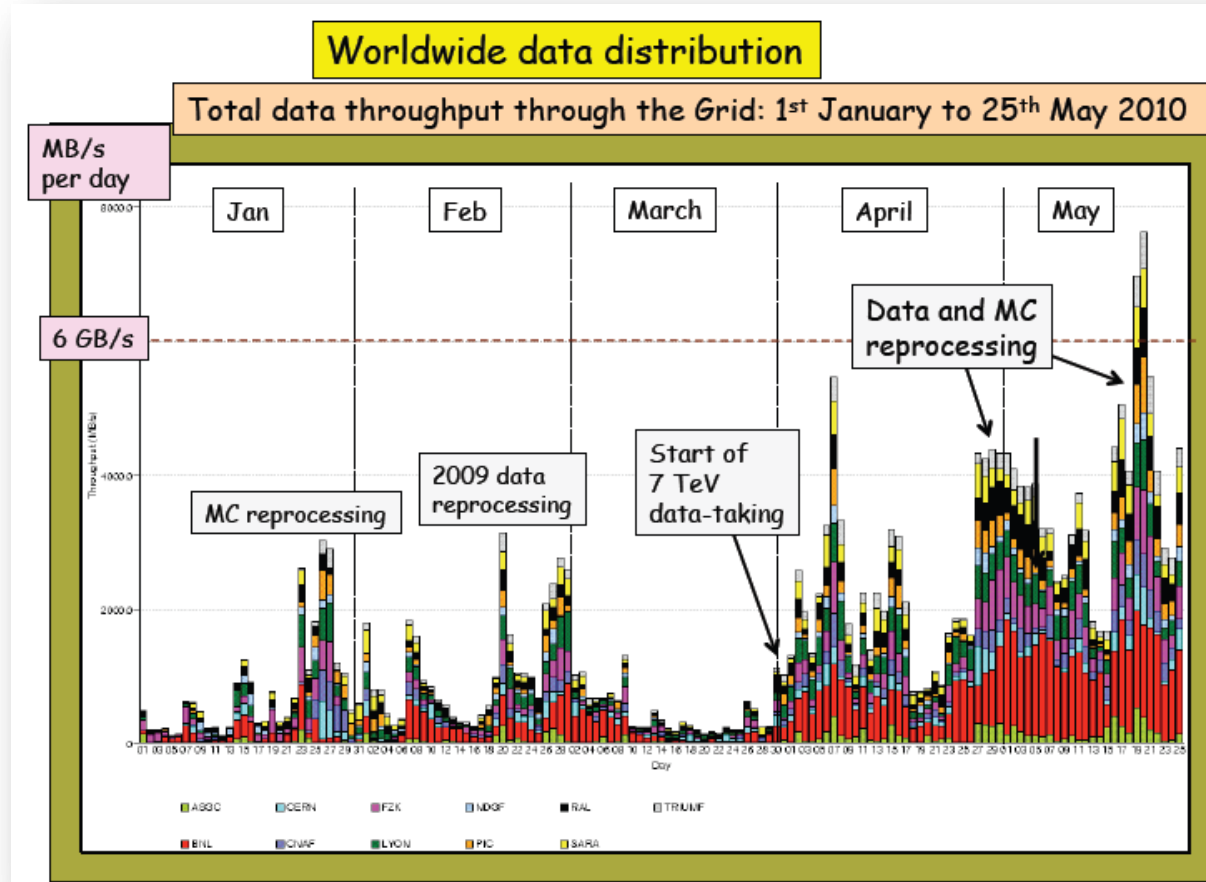
- Tier-0 (CERN)
 - RAW Detector Data Acquisition and archive to tape
 - Calibration and Alignment
 - First processing
 - Data distribution to Tier-1's
- Tier-1's (10 big Computer Centers)
 - One Tier-1 at the head of each *cloud*
 - Archive a share of the RAW Detector Data to tape (2nd copy)
 - Re-process those data when needed (new software, new calibration)
 - Archive Simulated data to tape and reconstruct when needed
 - Bulk analysis jobs but also user analysis in some cases
 - Data distribution to Tier-2's
- Tier-2's (60 mid size computer centers)
 - Many attached to a Tier-1 to form a cloud
 - Simulation Production
 - User analysis
- Tier-3's (100 (?) home institutes, faculty facilities)
 - End user analysis
 - None pledged resources; Not under ATLAS control

- RAW data master copy stored at CERN
- RAW data distributed over all Tier-1's
 - Tier-1 is responsible for preserving data on tape
 - And recall it for re-processing
- Cloud independence: All derived data available in each cloud
 - Generally, there should be a cloud with free CPU's
 - Generally, data should not have to move between clouds
- All data is pre-placed in each cloud
 - For controlled processing in Tier-1's
 - For user analysis in Tier-2's
- New data produced in a cloud should be archived there
 - Only Tier-1's are required to have tape archives
 - Also true for the Tier-0 (CERN)

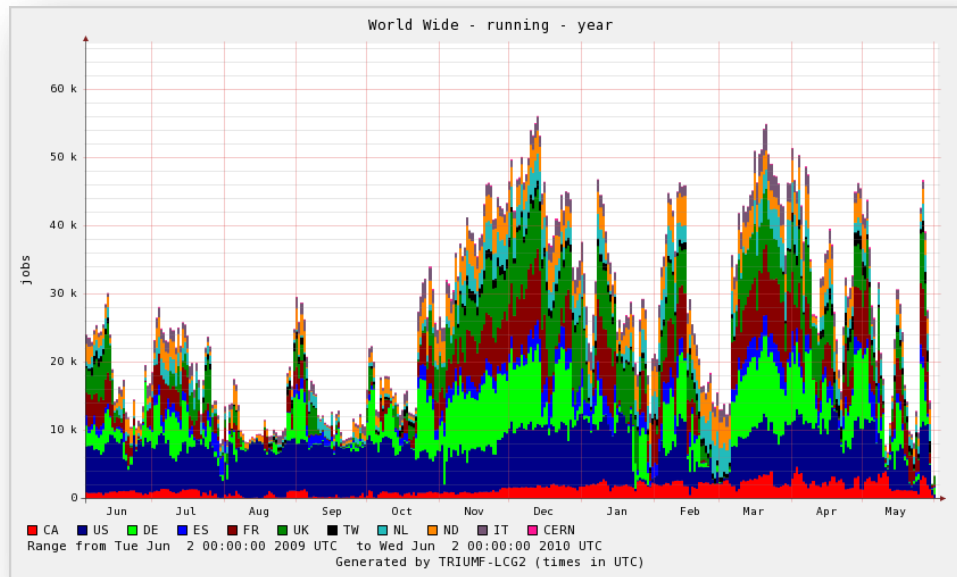


3

- Three sources of data: Data Taking, Monte Carlo generation and Reprocessing, and Data Reprocessing



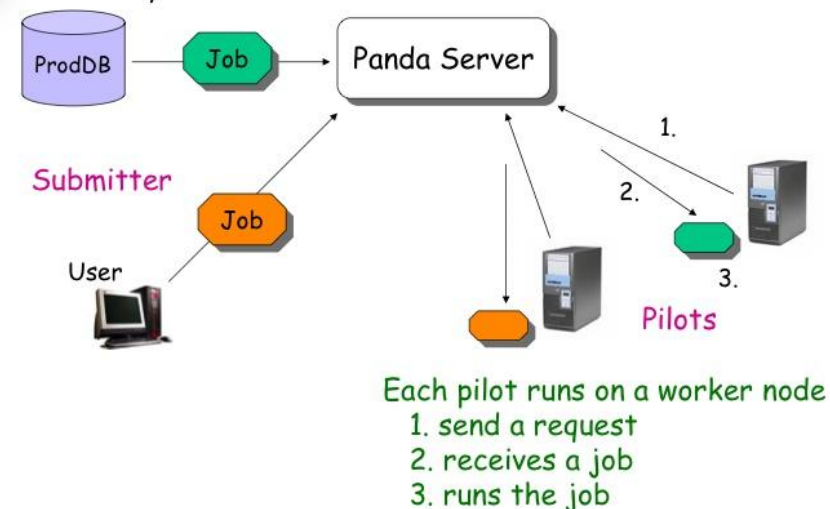
- Possible evolutions of the model:
- Cross-cloud T2->T2 currently goes through a T1
 - Could better exploit networks if some direct transfers were allowed
 - Need to work with middleware folks
- Data pre-placement might not be ideal for various activities.
 - Presently evaluating on-demand data placement and caching.



- PanDA is used to run all MC and Reprocessing, and ~75% of the user analysis worldwide
- PanDA@CERN deployed >1 year ago and is running successfully.
- The service was well prepared thanks to pre-exercises such as STEP'09

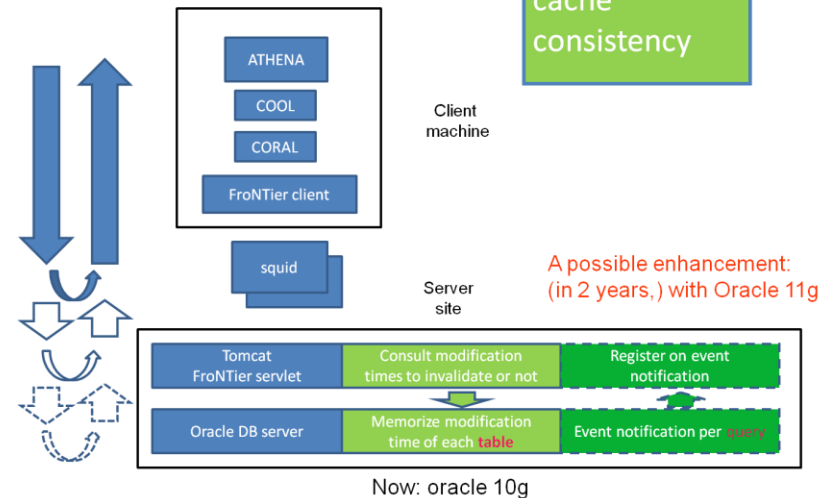
- Panda load depends more on the number of resources (~70 sites), and less so with the amount of data

Production system

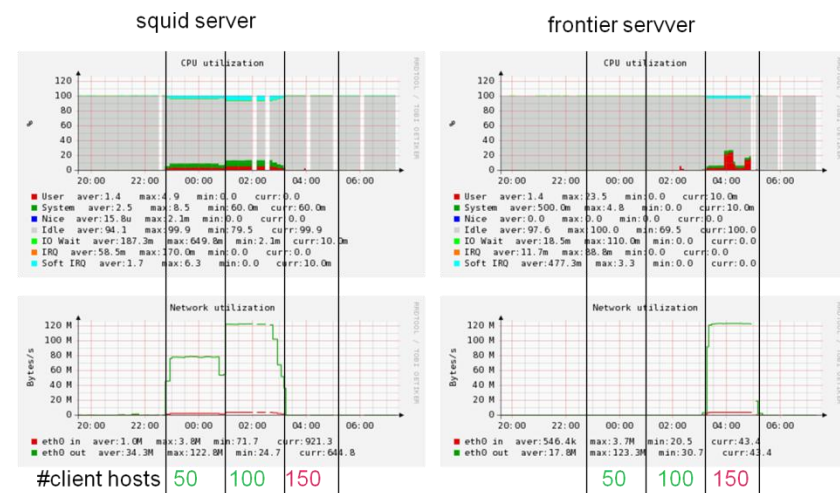


- FroNTier deployed to enable distributed access to the conditions DB
- Working toward making it more transparent to the end users

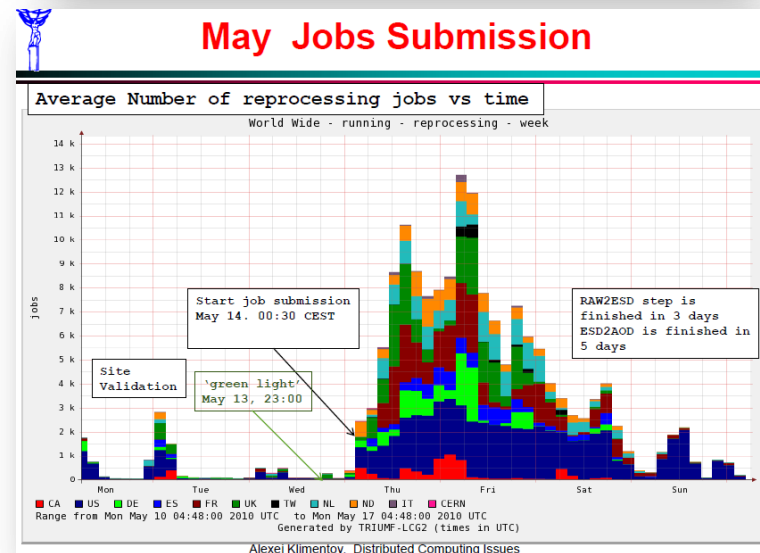
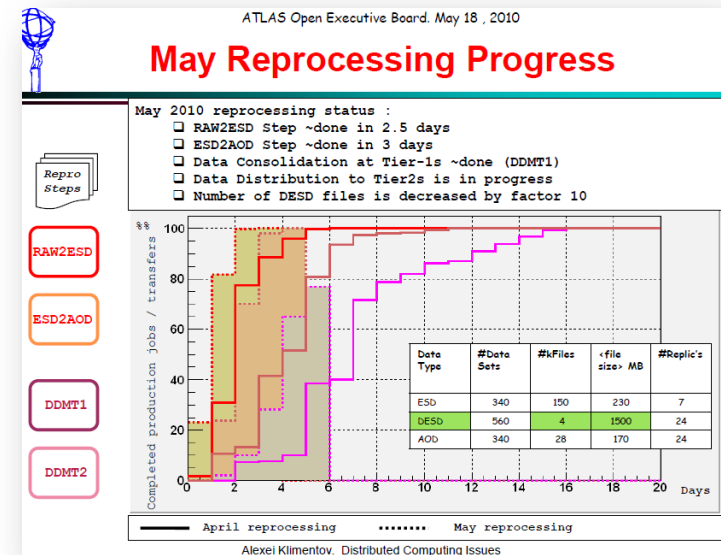
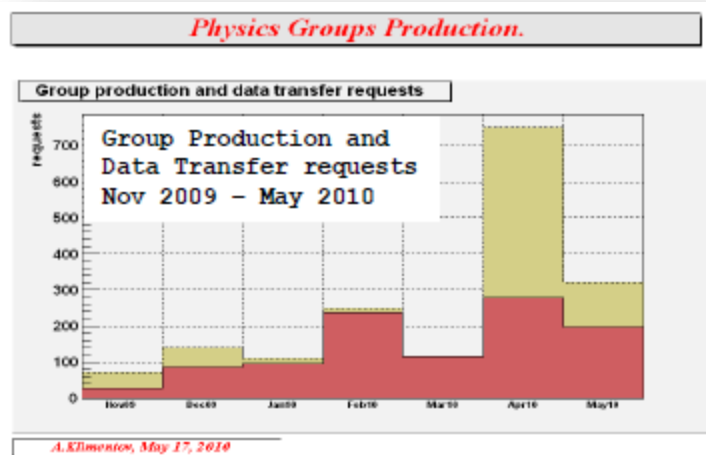
Handling
cache
consistency

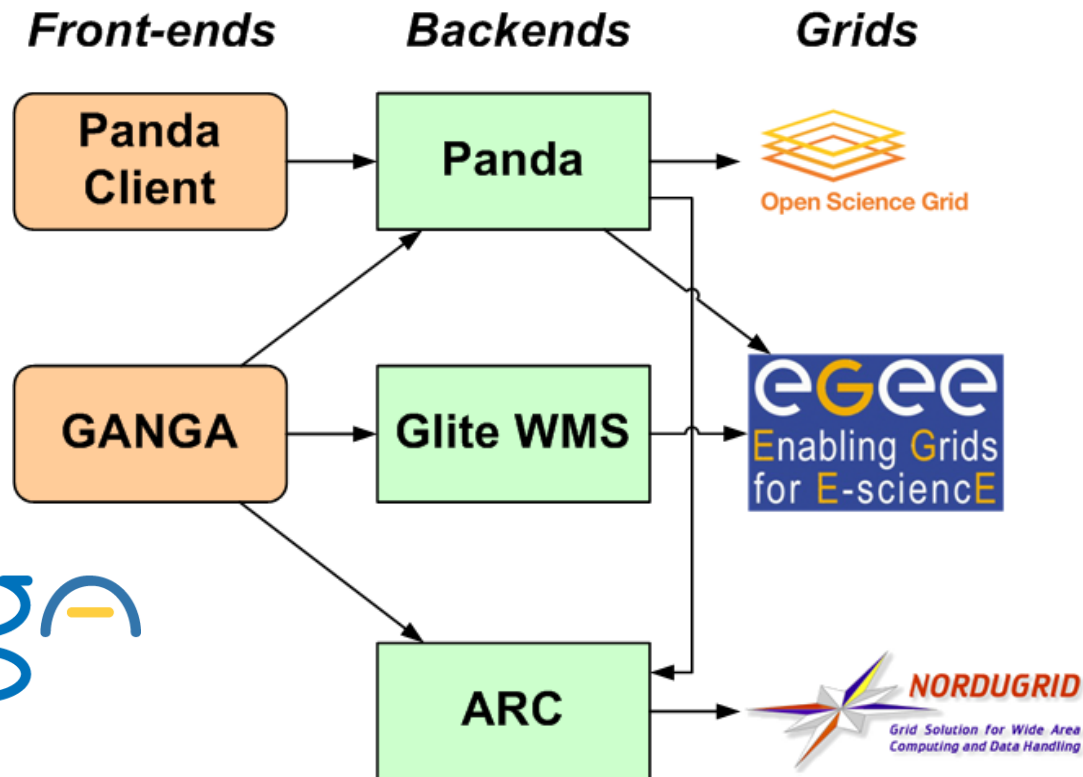


Map of installed Squids



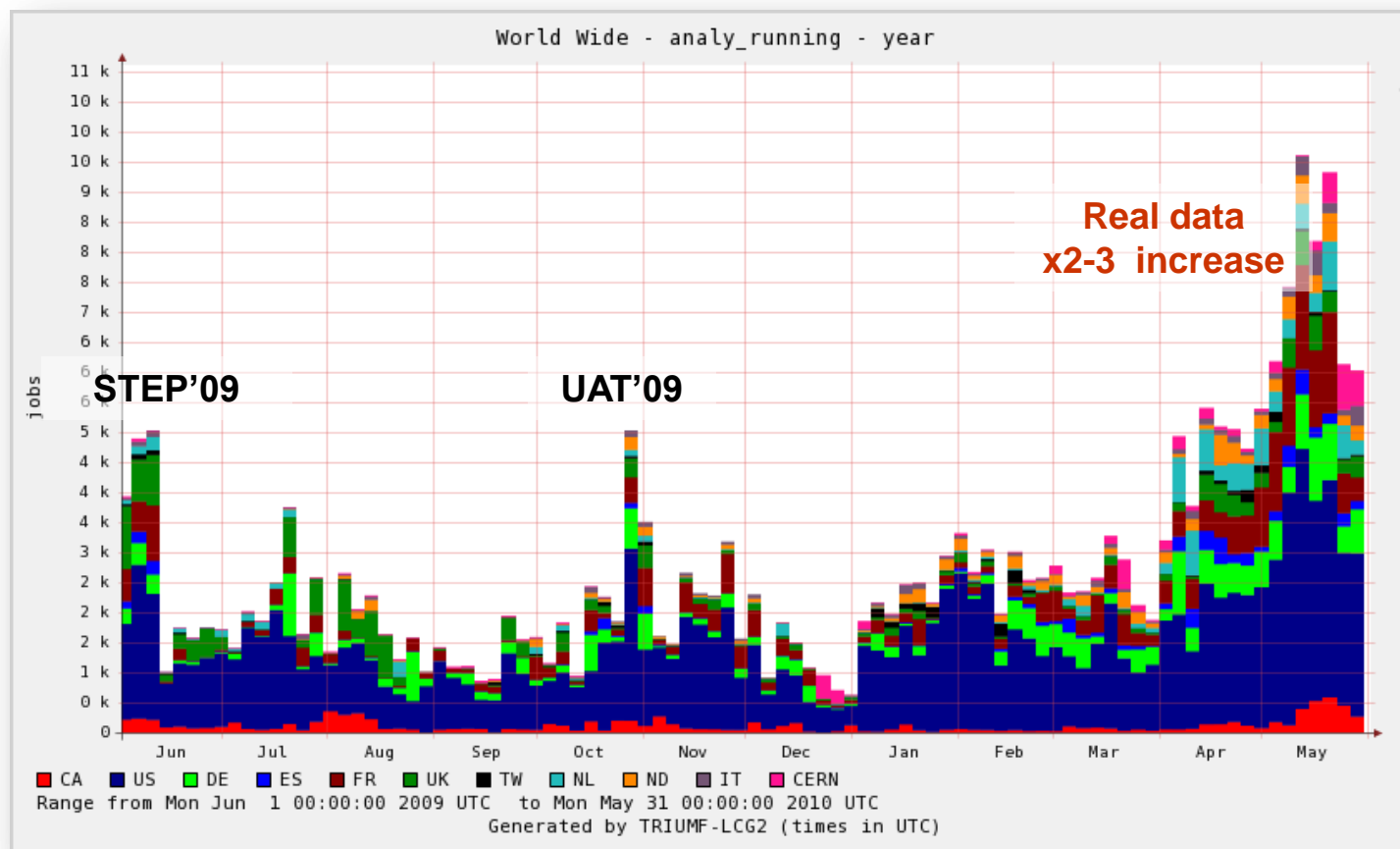
- GDP covers MC Production, Reprocessing, and Group Production
- Four Reprocessing campaigns since April
 - All requirements are fulfilled
 - Some site and DDM issues solved
 - Some long tails have been tackled by improved shifter procedures
- Group Production (D3PD production via production system) activity started trials in October and has ramped up lately

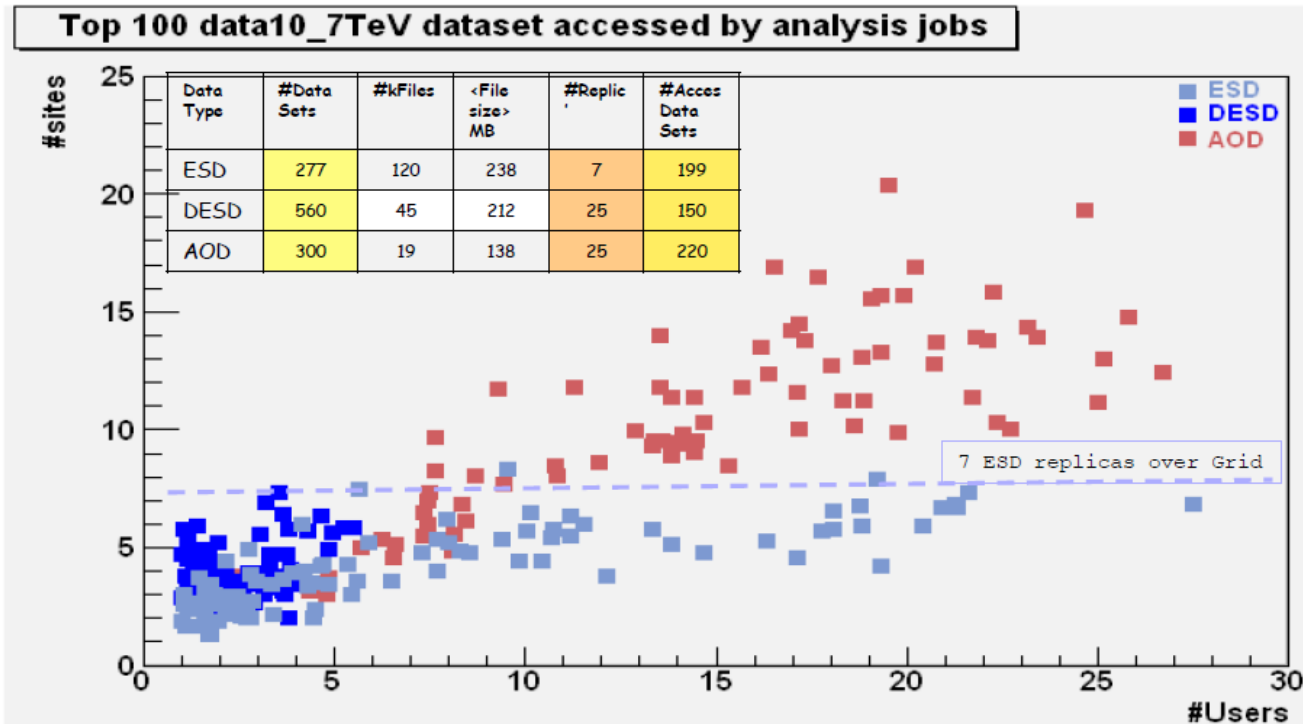




- Basic model: Data is pre-distributed to the sites, jobs are brokered to a site having the data
- Large dataset containers are distributed across clouds, so the front-ends do not restrict jobs to a cloud. i.e. DA jobs run anywhere in the world.

- In April and May: >900 users, 6.1 million successful jobs (+30% more failed)
- Reality is 2-3x larger than our average pre-data experiences, though previous tests were essential preparation.

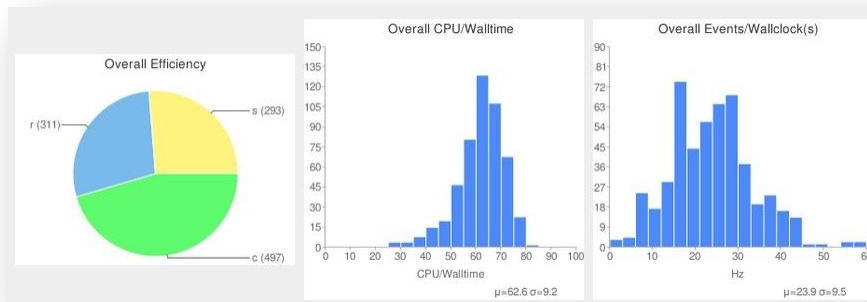




A.Klimentov, May 16, 2010

- 70+ sites get data, but ~30-40 are actively used. DESD's distributed widely, but not widely accessed.
 - Investigating popularity-based data movement and caching
- Small output file size – distributing millions of tiny files.

- We pre-validate sites for distributed analysis with Functional and Stress tests:
 - GangaRobot is running a continuous stream of short user analysis jobs at all grid sites
 - Results fed into SAM
 - Manual or automatic blacklisting
 - HammerCloud is used for on-demand stress tests spanning one or many sites
 - Used to commission new sites, tune the performance at existing sites, and to benchmark sites to make comparisons
- HammerCloud
 - Invested ~200k CPU-days of stress testing jobs since late 2008.



- We have ~1000 active distributed analysis users
 - They should not need to be distributed computing experts – The Grid is a black box that should just work
 - Grid workflows are still being tuned – not everything is 100% naïve user-proof
 - Supporting the users to get real work done is critical (it will stay like this!)
- ATLAS introduced a team of expert user support shifters in fall 2008.
- DAST: Distributed Analysis Support Team
 - Class 2 (off-site) ATLAS shifts; week-long shifts in EU and NA time zones (Asia-Pacific shifters wanted...)
 - 1st and 2nd-level support: better incorporate new shifters and shares the load in times of high demand
 - DAST is a ~15 member team; each takes a shift every 4-8 weeks.
- Users discuss all problems on a single “DA Help” eGroup
 - Discussion about all grid tools, workflows, problems
 - Not just DA – also data management questions

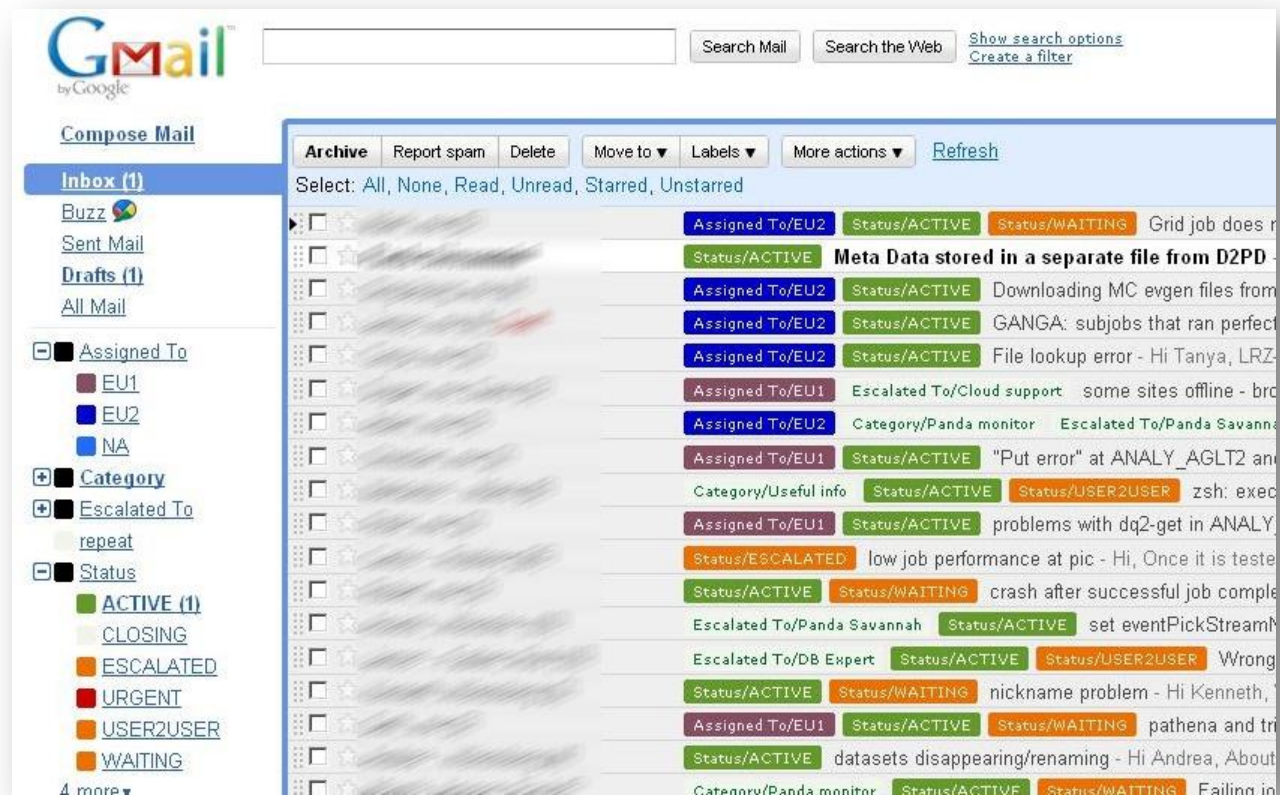
- Two main requirements of the DA User Support discussions:
 - R1: Enable (and encourage) user-to-user support
 - R2: Enable shifters to track issues, shifter assignments, escalations to experts, and moderate threads
- R1 points at a mailing list or eGroup; R2 points at a trouble ticketing system (e.g. RT).

Gmail has worked very well as a compromise solution

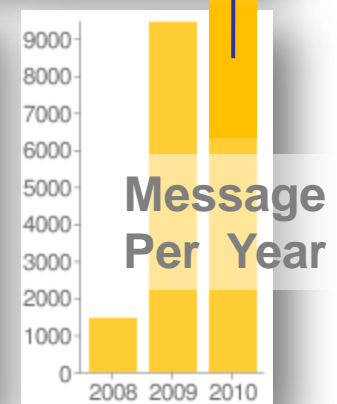
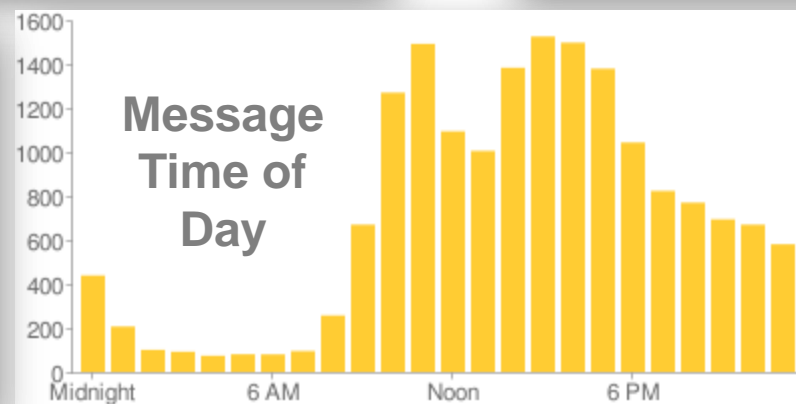
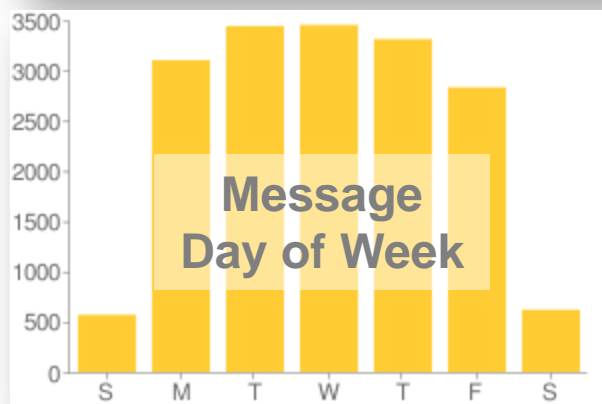
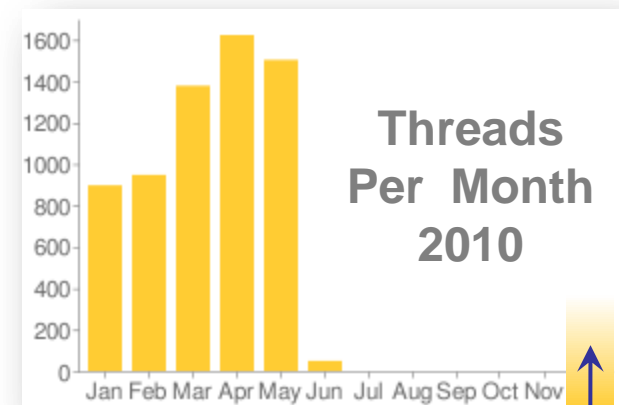
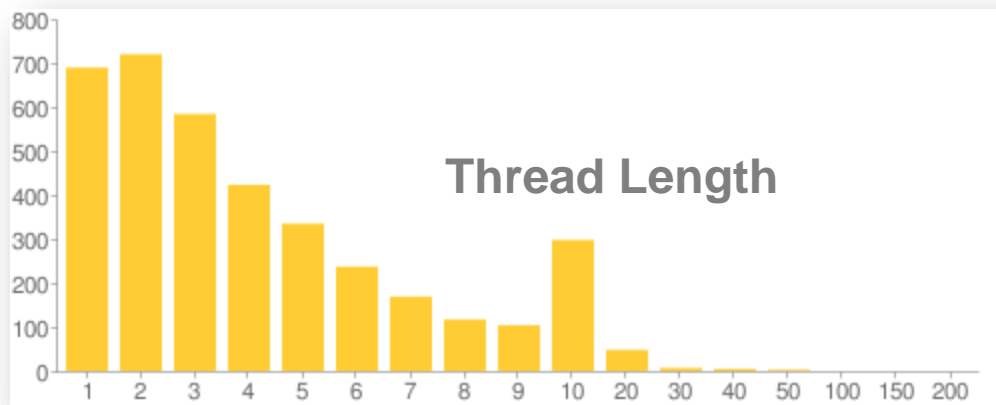
Shared account for the DAST shifters

Threading keeps issues separated

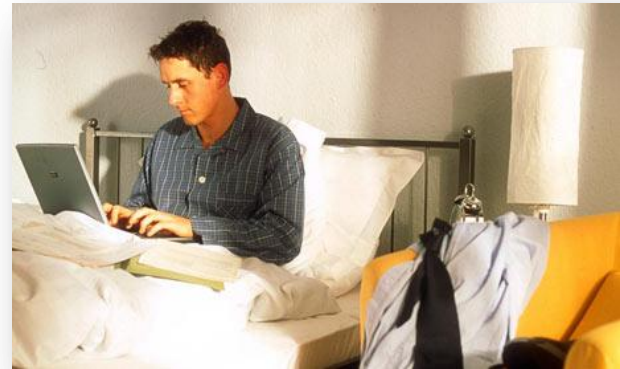
Track and assign issues with thread labels



- Some interesting measures of the DA Help messages and threads...



- Enabling Tier 3 activity is the next essential step in ATLAS Computing
- Formed working groups in February 2010 to study:
 - Distributed storage (Lustre/Xrootd/GPFS)
 - DDM-Tier3 link
 - Tier 3 Support
 - PROOF
 - Software / Conditions Data
 - Virtualization
- WG's are wrapping up now. Sites are starting to get connected.



- The ATLAS Distributed Computing infrastructure is working thanks to many efforts in preparation
- We are able to
 - process, distribute, and reprocess the data
 - analyse the data
 - provide support to our large community
- and we are tackling the next frontier: Tier 3
- 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?