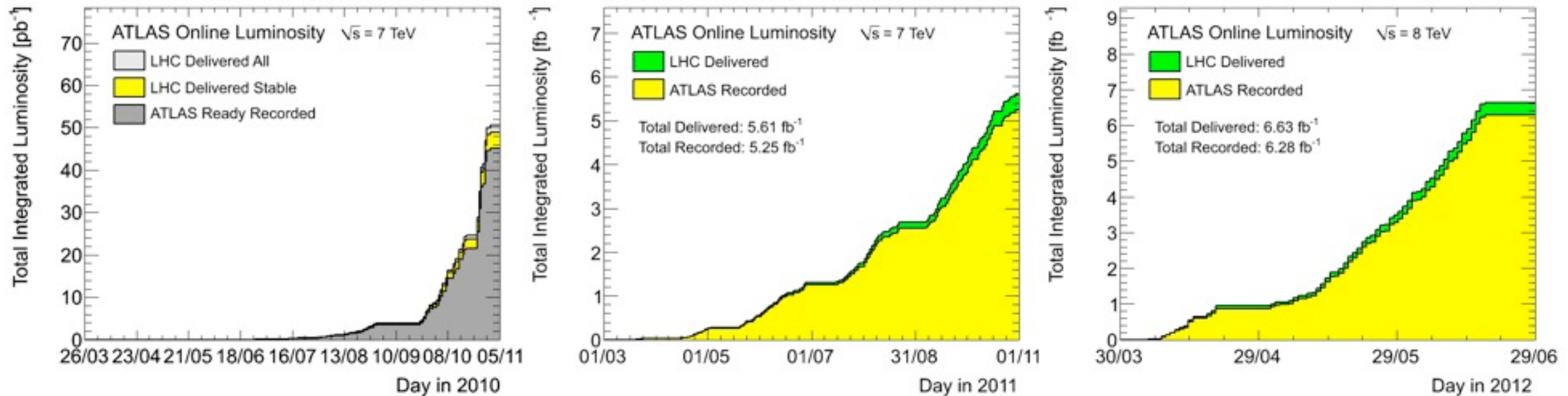


# The Present and Future Challenges of Distributed Computing in the ATLAS experiment

Ueda I.  
on behalf of the ATLAS Collaboration

# ATLAS has collected a large amount of data



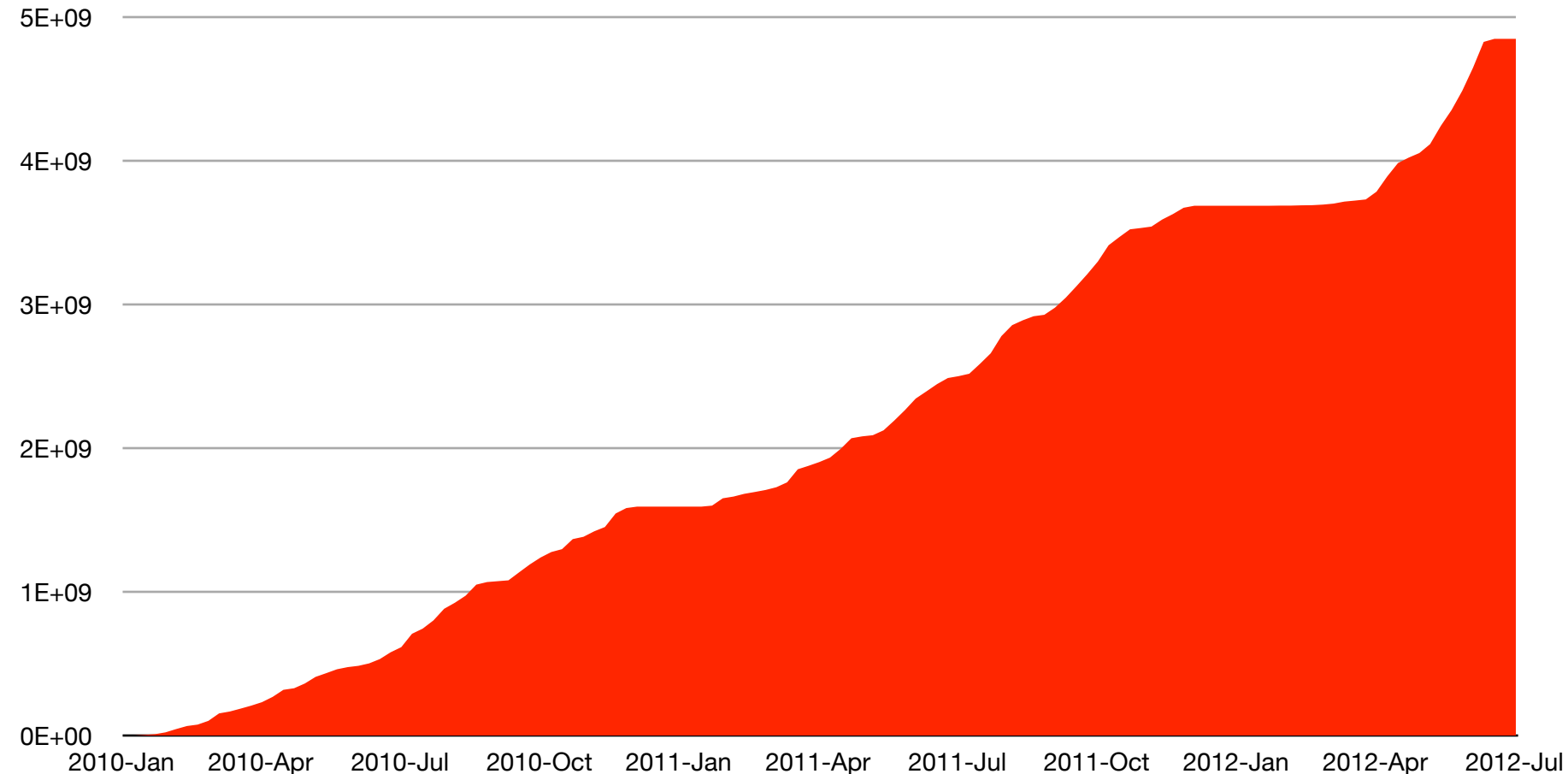
ATLAS RAW Events Registered at Tier-0

More than 5 fb<sup>-1</sup> at 7TeV  
and more than 6 fb<sup>-1</sup> at 8  
TeV in pp collisions

- and nearly 170 μb<sup>-1</sup> of  
heavy-ion collisions

Corresponding to 4.4  
(pp) and 0.4 (HI) billion  
events up to now since  
2010

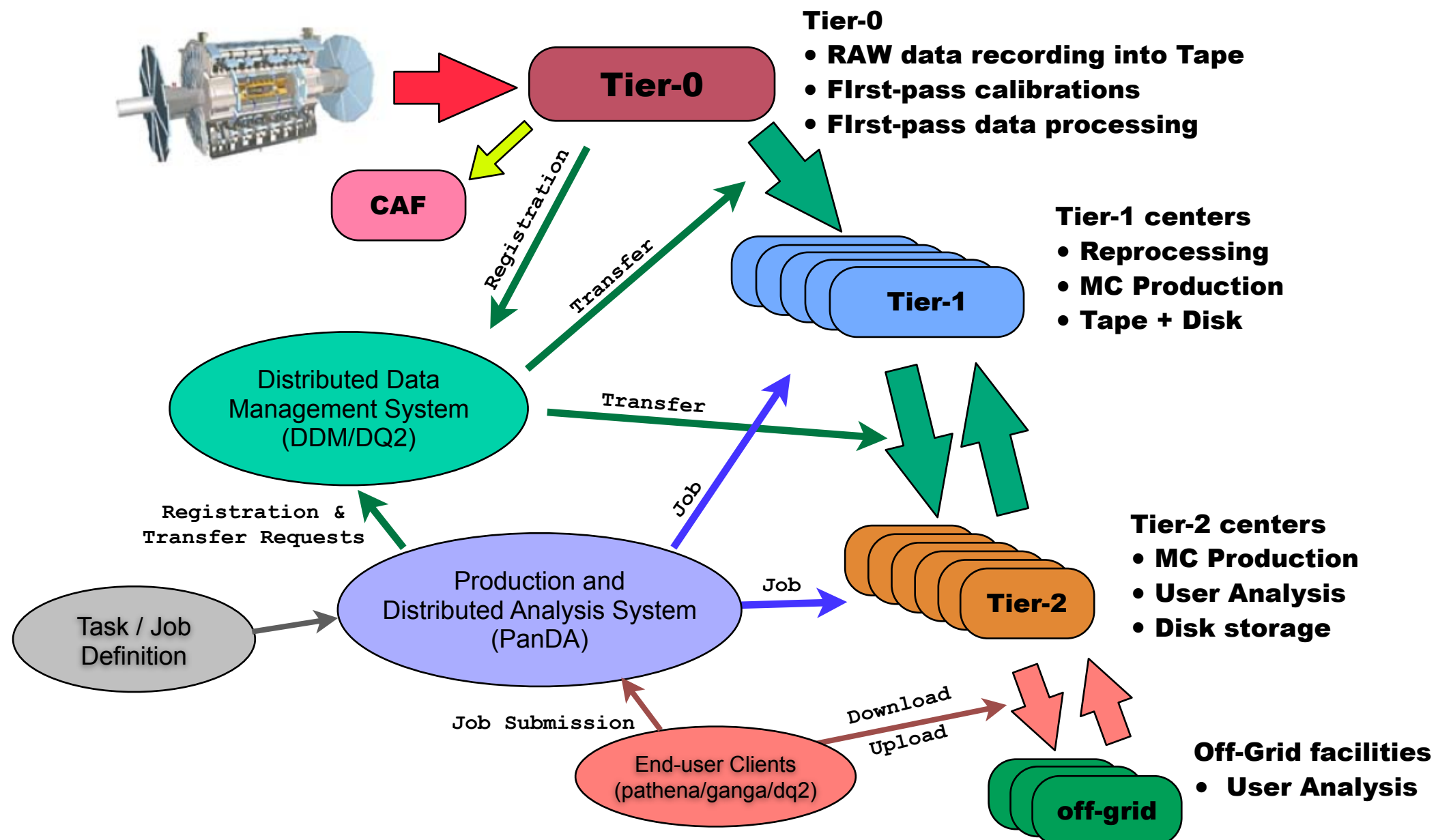
Recorded event rate is set  
by the trigger (up to 400 Hz)



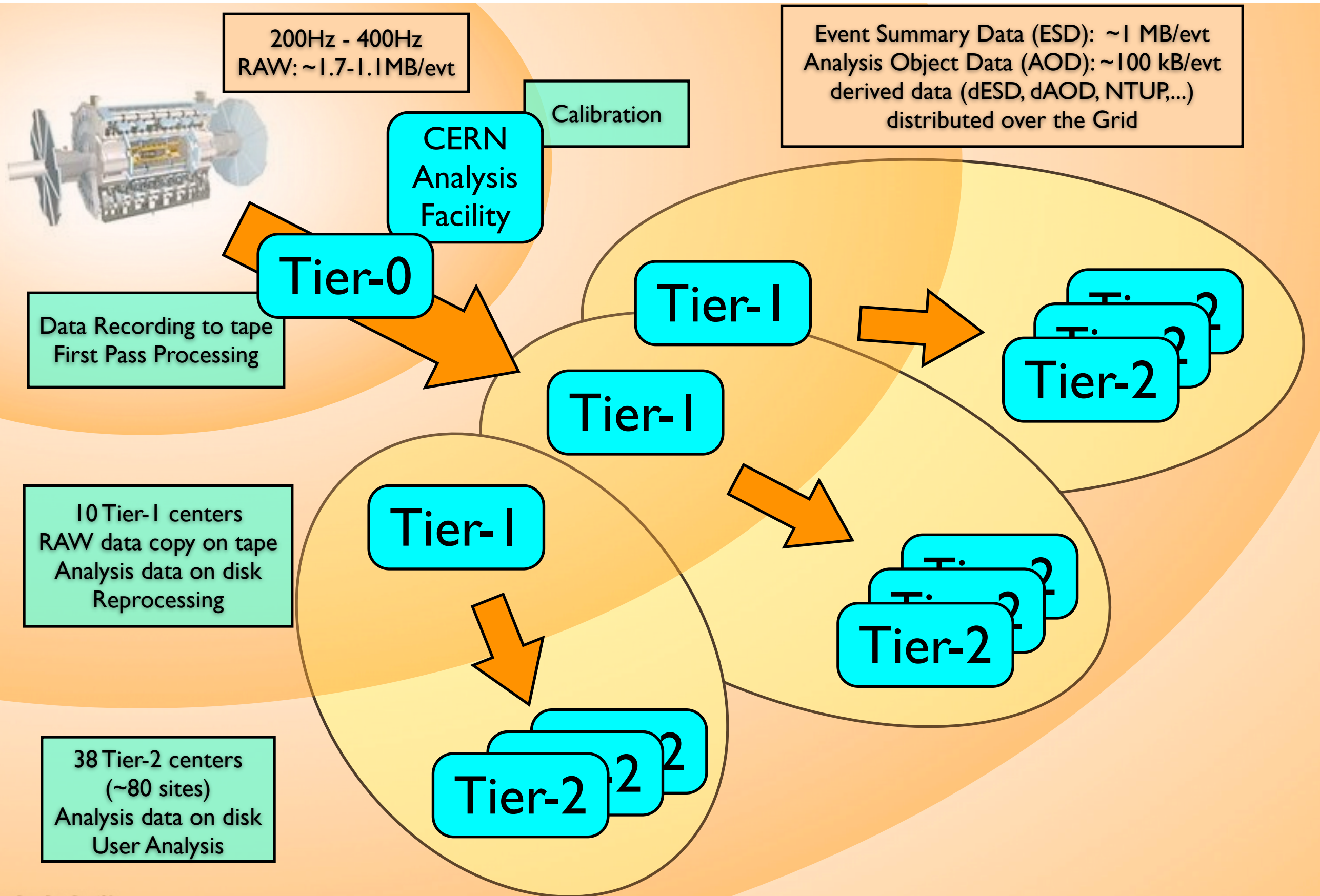
# ATLAS Computing System

The ATLAS Distributed Computing System manages world-wide

- Data processing, MC production and analysis jobs running on over 120k computing cores
- Data transfers to and accesses from ~ 50 PB disk space (and ~ 30 PB tape)
- On over more than 100 sites set-up with the “LHC Computing Grid” middleware

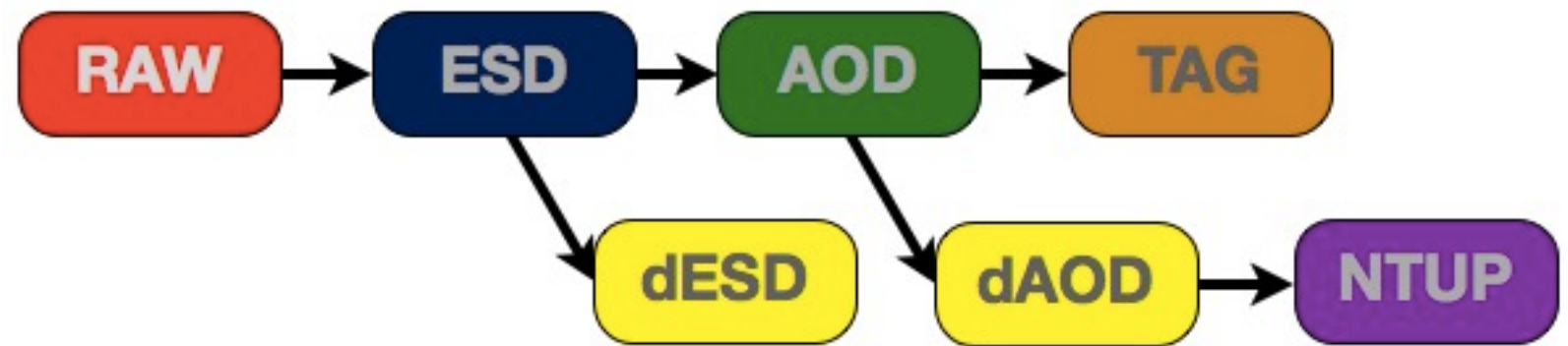


# ATLAS Computing Model





# ATLAS Computing Model



## Data Types:

- RAW -- raw data from the detector -- need processing before analysis
- ESD (Event Summary Data) -- output of event reconstruction
- AOD (Analysis Object Data) -- event representation with reduced information for physics analysis (ATLAS-wide format)
- DPD (Derived Physics Data) -- representation for end-user analysis. Produced for working groups or individual end-users (group-specific format)
  - dESD (performance groups), dAOD (physics groups), NTUP (physics groups and end-users)
- TAG -- event-level metadata (event tags), short event summaries primarily for event selection

## LHC Computing Grid: computing facilities distributed world-wide organized in “tiers”

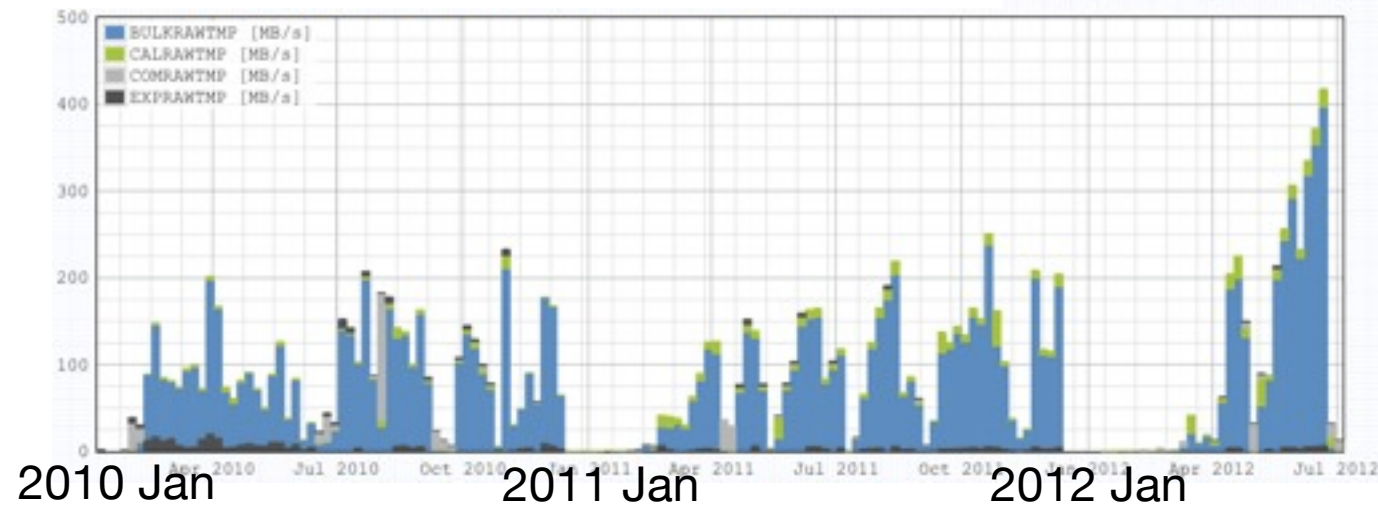
- Tier-0 (CERN) : record RAW data on tape, first-pass processing
  - CERN Analysis Facility: mainly for calibration
- Tier-1 : store replicas of RAW on tape, reprocessing, conditions database (originally\*)
- Tier-2 : group and end-user analysis (originally\*)
- MC simulation : wherever possible (and capable)
- Data distribution / replication over the Grid
  - for redundancy -- to secure data with replicas
  - for accessibility -- more replicas for data used frequently
- ATLAS software releases installed at sites with special Grid jobs (originally\*)

\*: the original model has been revised  
(next slide)

# ATLAS Computing Model revised



RAW data rate from online into T0 system



	2010	2011	2012
Trigger Rate	< 200 Hz	< 400 Hz	< 400 Hz
RAW	1.7MB	1.1MB	(1.2MB)
compressed		0.66 MB	0.73 MB
ESD	1.05 MB	1.21 MB	1.86 MB
AOD	0.09 MB	0.16 MB	0.19 MB

Adjusting the model designed from first principles to varying real conditions during the first years of data-taking

- Unexpected usage pattern on data types
  - data distribution plans revised
  - dynamic data placement following the usage
- Higher event recording rate up to 400 Hz and RAW data on DISK since 2011
  - RAW data compression (factor ~2)
  - ESD with limited lifetime -- allowing prompt detector studies but removed after several weeks to cope with higher trigger rate and RAW on disk within the disk capacity

Monitoring is the “key” -- a lot of improvements and new monitoring + automatic control tools have been made

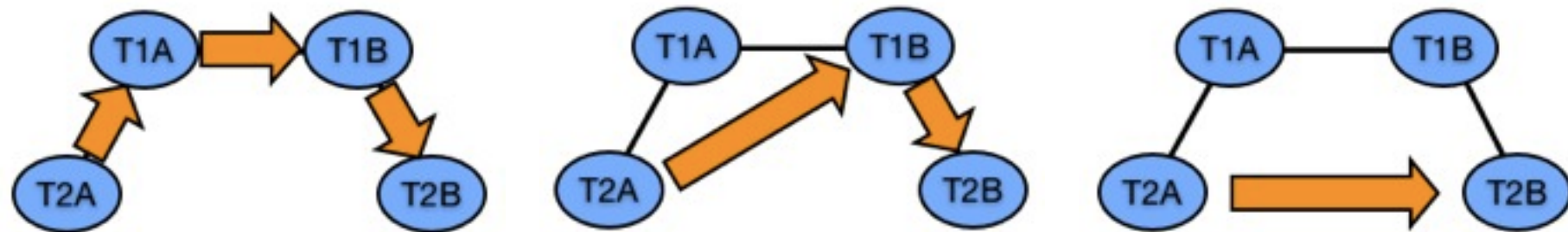
- to understand the activities and needs
- for stable operations and to minimize the manual intervention

# ATLAS Computing Model revised



## Data transfer path revisited to enable direct transfers

- Faster transfer path is chosen for efficiency, less load to the system, and **end-user convenience**
- Higher network bandwidth available than anticipated, less hierarchy needed
  - from “tree” like topology with Tier-1 - Tier-2 association
  - to “mesh” of any Tier-\* combination **based on the measured network** performance



## New technologies arising

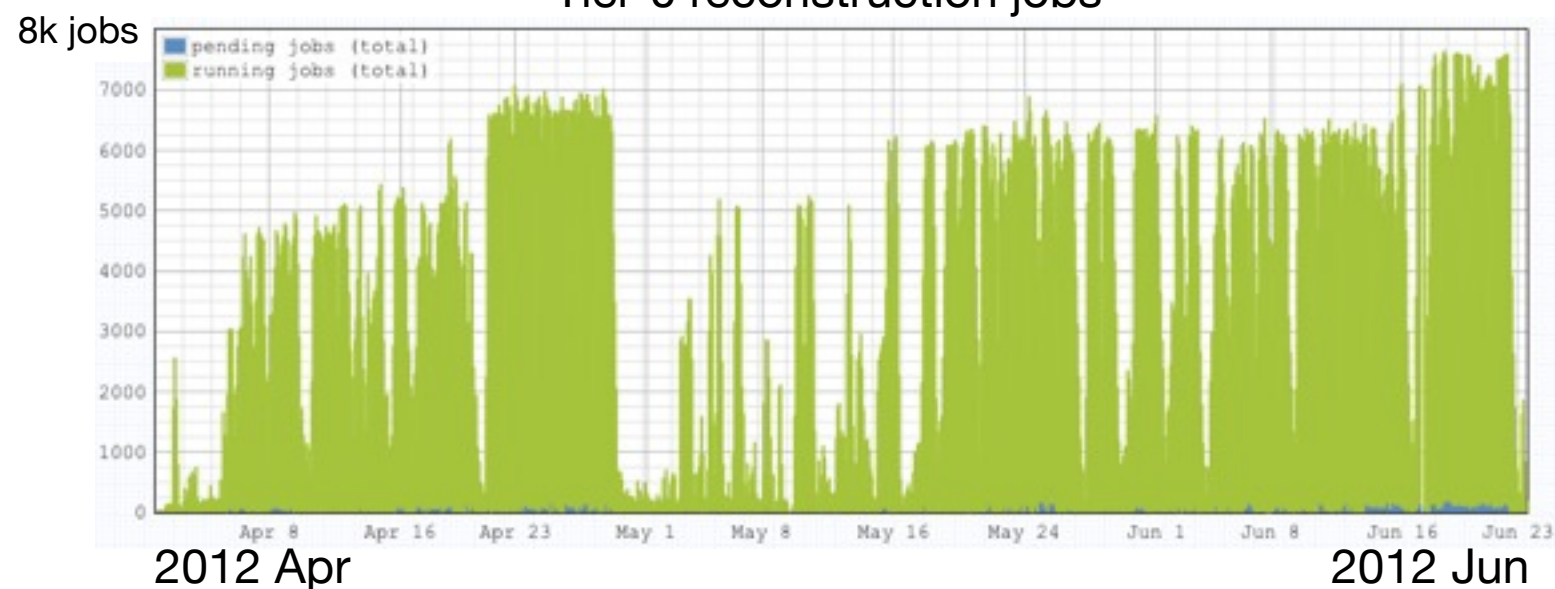
- Frontier/Squid : remote access to the database at Tier-1s from any site via http
  - Conditions data access from any site, not only at Tier-1s
- CernVM-FS : network file system with http
  - No more need to pre-install ATLAS software releases at sites (can avoid load on shared file system)

# ATLAS Tier-0 System

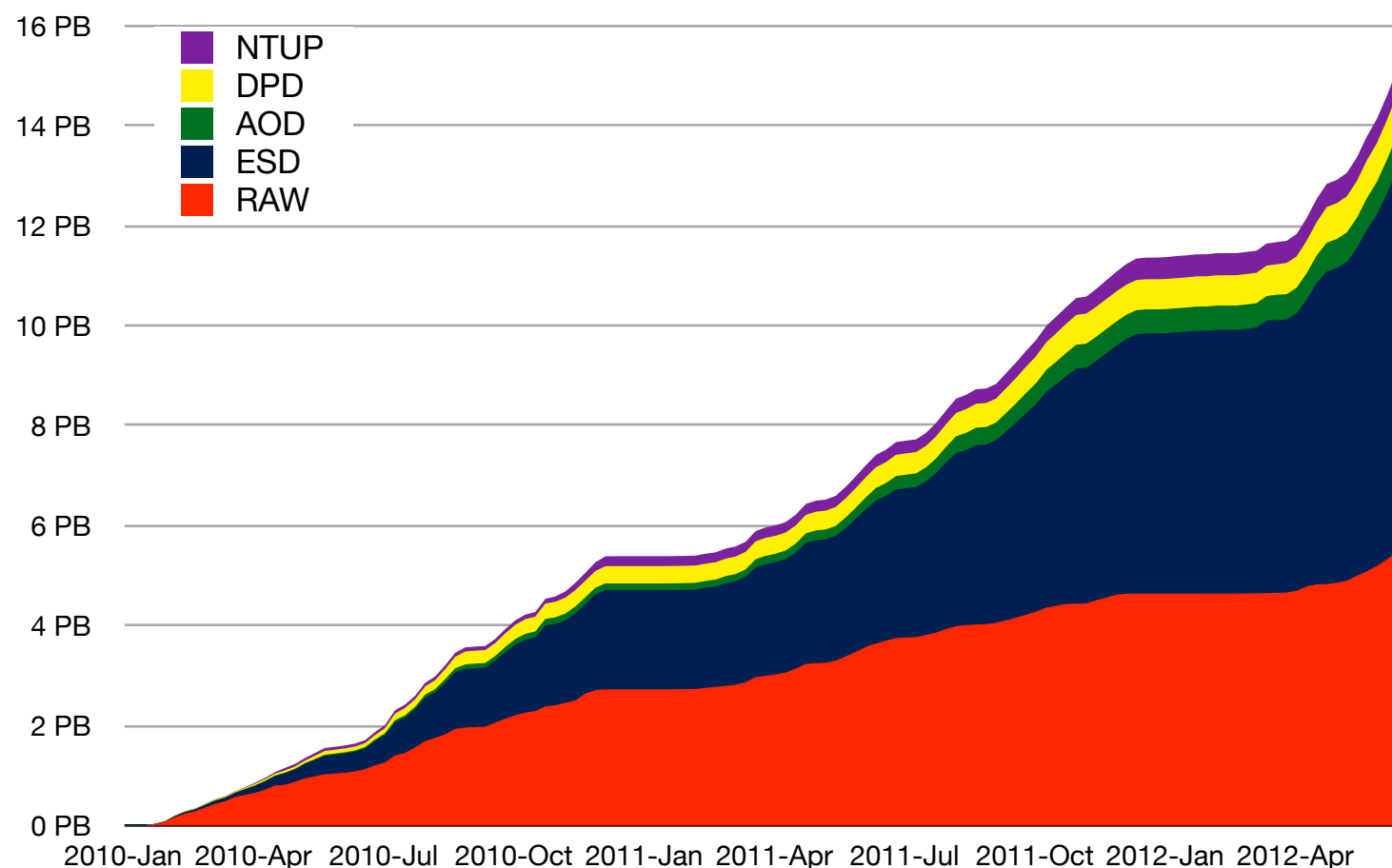
ATLAS Tier-0 system has been running reliably, stably, successfully

- First-pass data processing has kept up with LHC performance
  - Extension of dedicated resources each year
  - Flexible use of non-dedicated resources – up to 7.5k job slots
- High quality data reconstruction already from the first-pass processing
  - Express stream processing and calibration loop before bulk-processing
  - Most 2012 data are used in physics analysis for official physics results without reprocessing
- A comprehensive monitoring suite helping operations
  - Fast and flexible reaction from development team to requests

Tier-0 reconstruction jobs



ATLAS Data Registered at Tier-0





# ATLAS Distributed Computing

## Data processing on the Grid sites

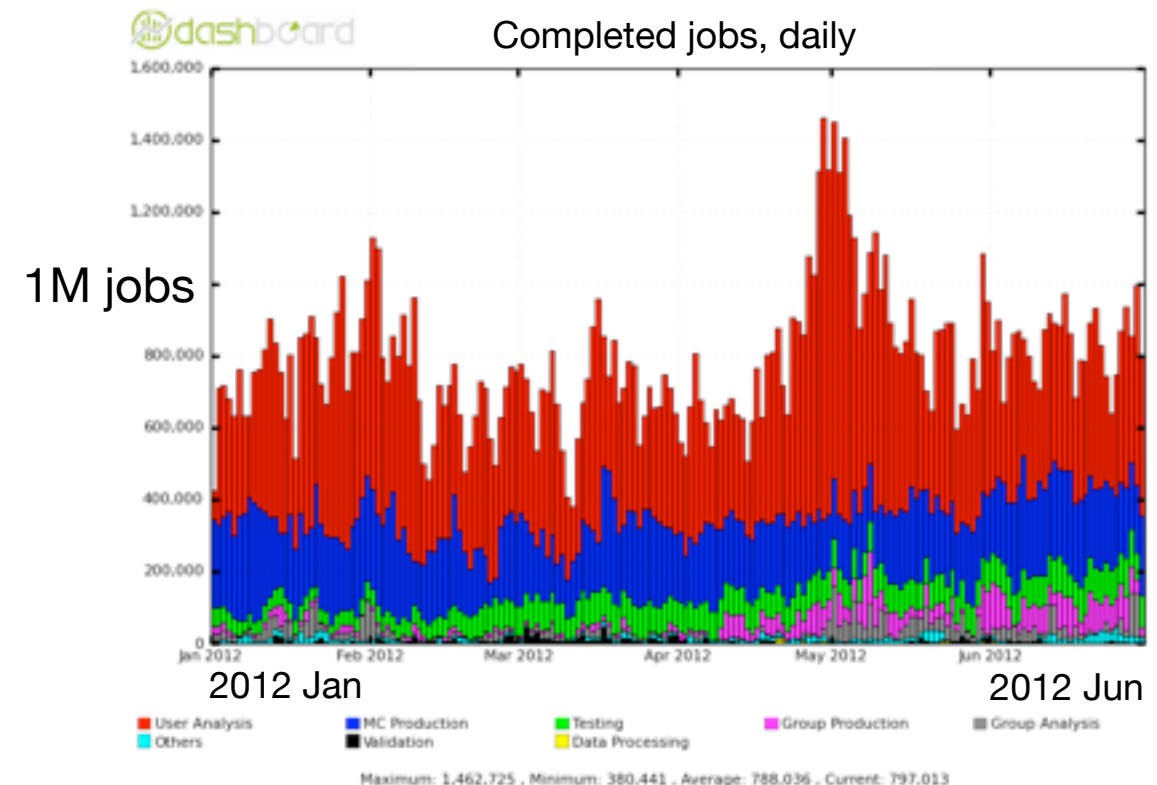
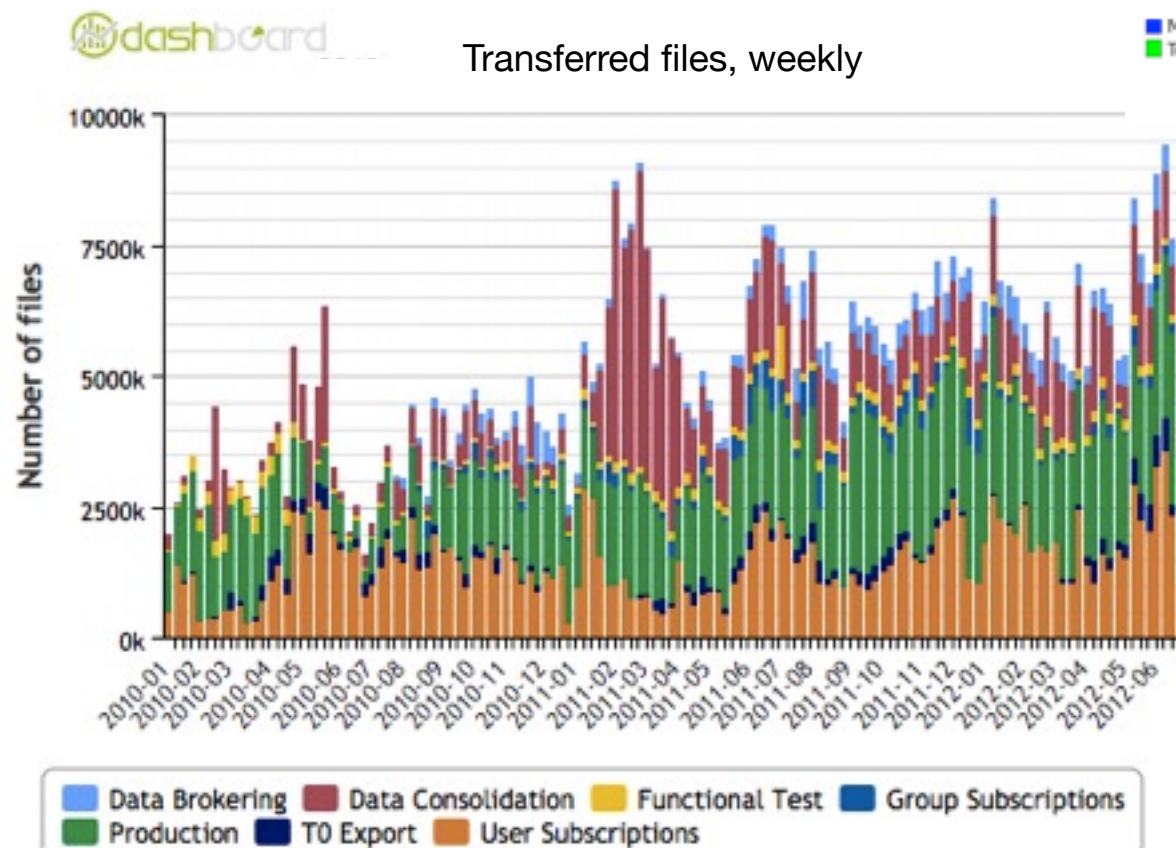
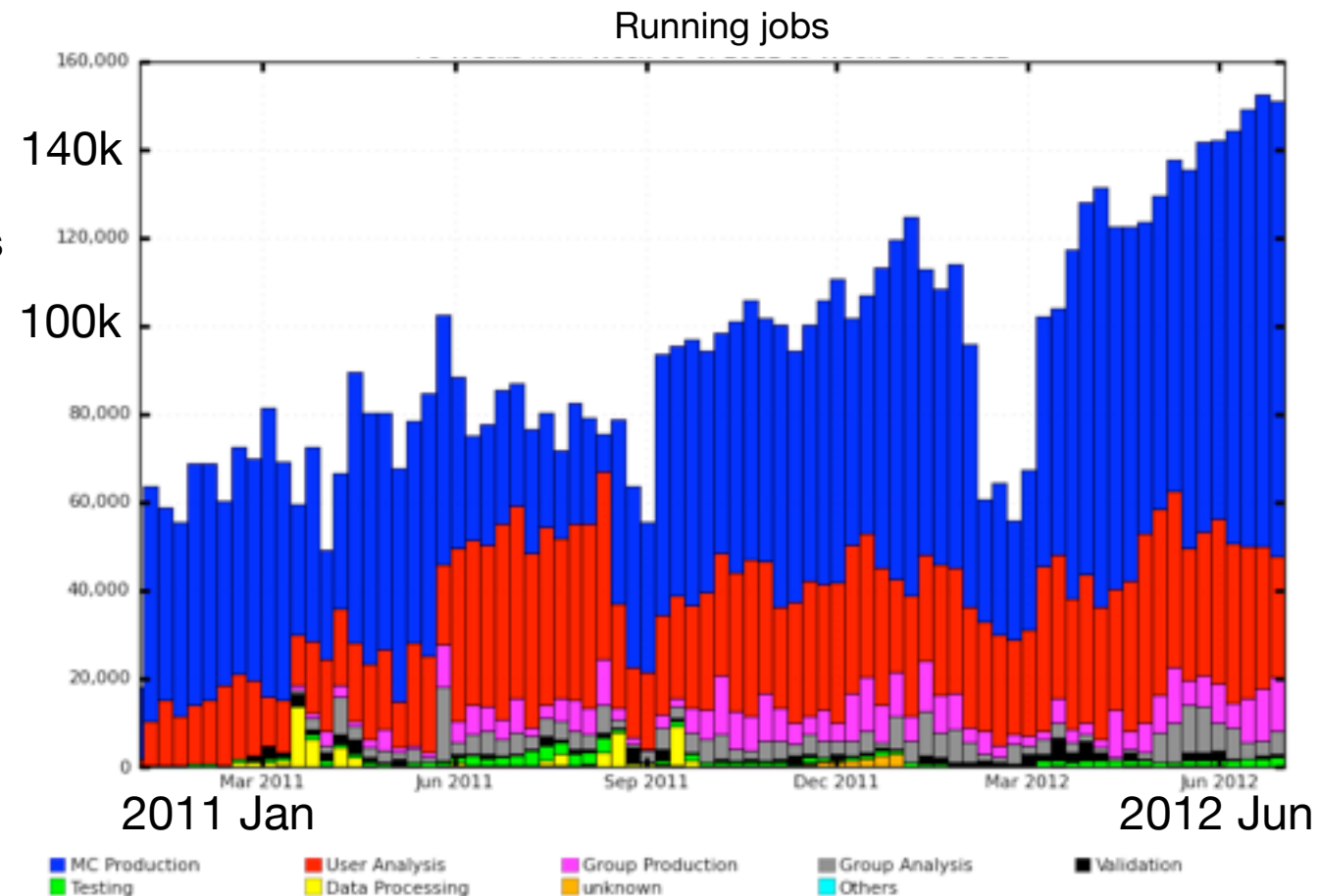
- More than 100k computing cores (doubled in the last 18 months)
- Data processing and analysis of about a million jobs daily

## Data distribution to the Grid sites

- Transferring about a million files daily
- Input and output of production and analysis jobs

## End-user analysis

- can be run with a simple command
- Output can be transferred to their “home” site on Grid (manually or automatically) as well as downloaded to off-Grid computers (not monitored)



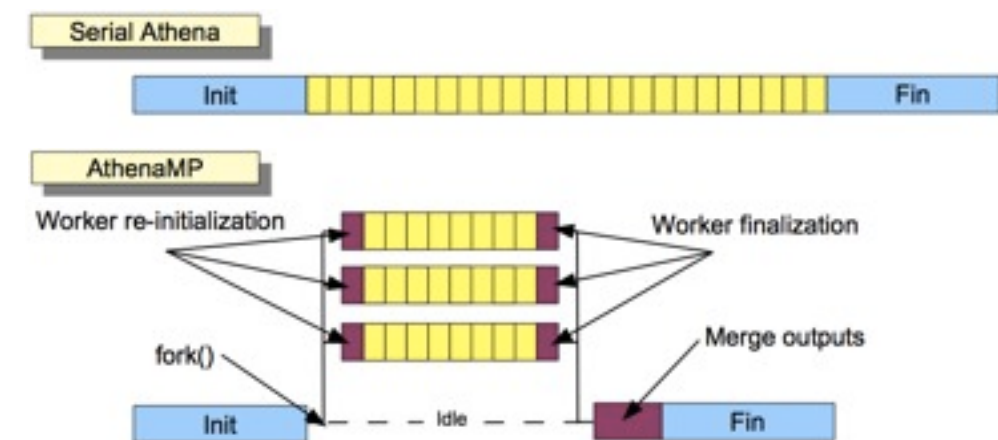
# New Challenges

## The current system

- “Data Grid” concept = Jobs go to Data
- Sites set up with the Grid middleware
  - on top of classical “standard” batch systems
  - job slot = computing (physical / logical) core

## On-going efforts

- AthenaMP -- event-parallel version of the ATLAS framework to optimize resource/memory utilization
  - job slot = multi-core
- Storage Federation
  - WAN access to data (jobs access remote data)
- Virtualization and Cloud Computing
  - utilizing the new “standard” == cloud
- Evolution / new development of the ATLAS computing systems



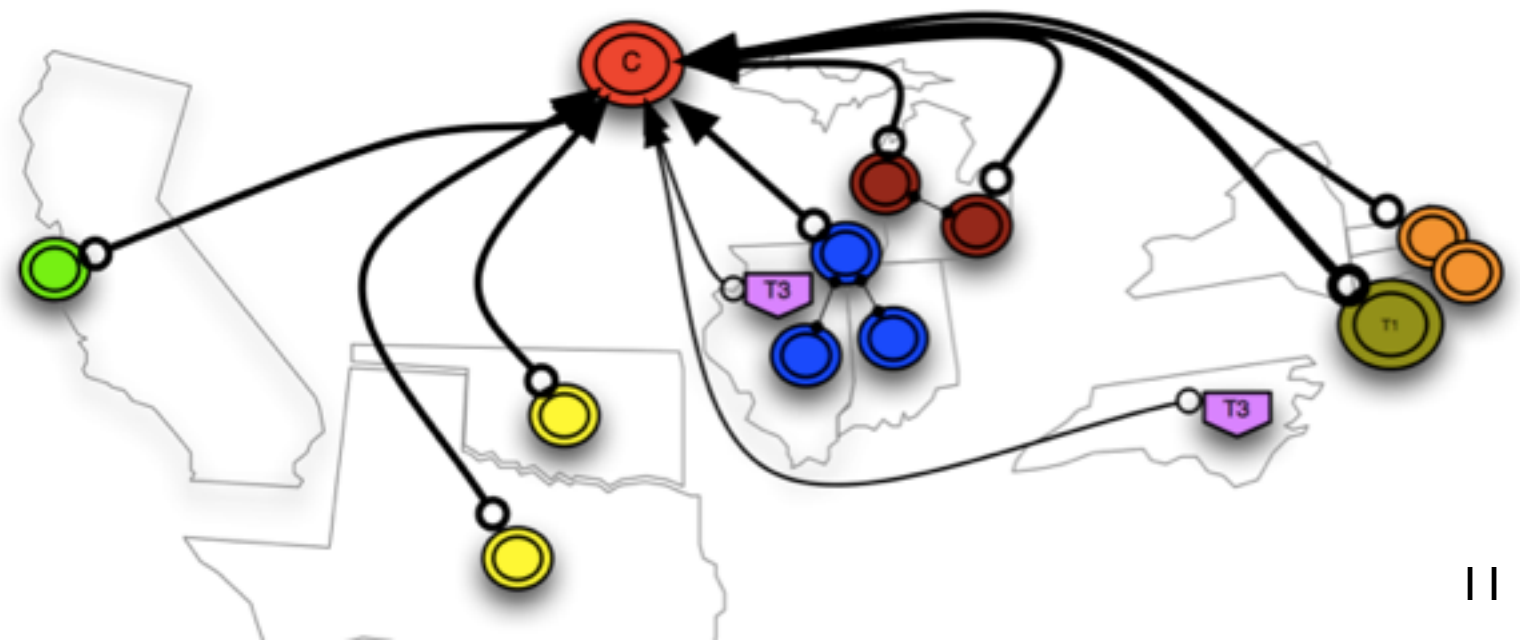
# Storage Federation

The current system is based on the “Data Grid” concept

- Jobs go to data -- access via LAN
- Replicate data for higher accessibility
  - transfer the whole dataset
- Jobs to be re-assigned when the data there is not available

“Storage Federation” provides new access modes & redundancy

- Jobs access data on shared storage resources via WAN
- Analysis jobs may not need all the information / all the files
  - Transfer a part of the dataset
  - File and Event Level Caching
- System of Xrootd ‘redirectors’ is the possible working solution today
  - Work in past year within US ATLAS Computing Facility to develop the concept and test performance
  - Test being extended from regional to global



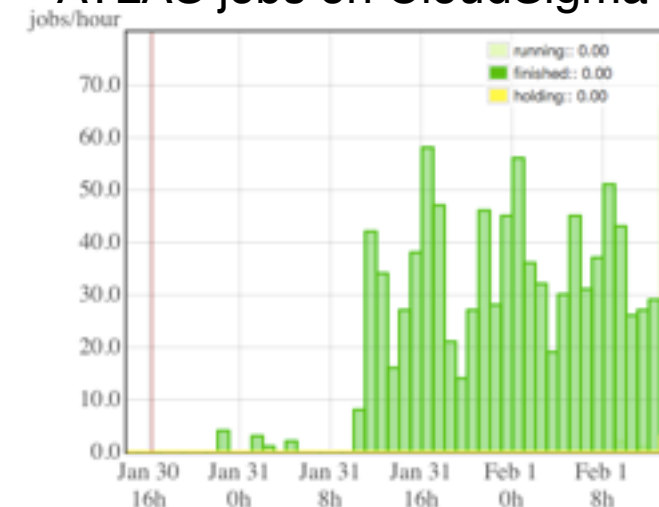


# Virtualization and Cloud Computing

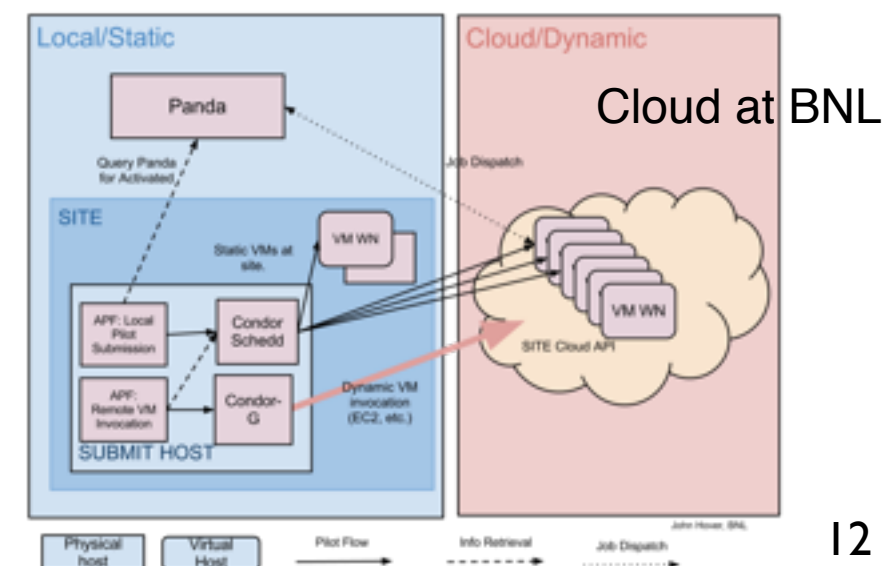
**Goal:** Integrate academic/commercial cloud resources into the ATLAS computing system

- Clouds in the ATLAS production/analysis system
  - 'Helix Nebula - the Science Cloud' -- a European cloud computing platform.
    - **First step completed:** ATLAS simulation jobs running on the contributing commercial cloud providers: CloudSigma, TSystems and ATOS
  - 'Cloud Scheduler'
    - University of Victoria, NRC
    - Tens of thousands of jobs have been executed at FutureGrid (Chicago) and Synnefo cloud (UVic)
- Institutional analysis clusters in the cloud
- Personal analysis in the Cloud via ATLAS system
  - Performance comparison LxCloud vs. LxBatch
- Cloud facilities within Grid sites
  - To be more generic to support various experiments
- Data storage
  - Performance studies
  - Integration into the ATLAS data management system

ATLAS jobs on CloudSigma



ATLAS jobs on Cloud Scheduler





# Evolution of the ATLAS computing systems



## New developments

- The current system has been working well, but we would like more features, flexibility and less manual intervention
- Plans based on the current experiences and new technologies to sustain the load in the coming years. To be ready for LHC restart (2014) ...
- With some scalability re-consideration
- Rucio: the next ATLAS Distributed Data Management System
  - Global management of the space, rather than per-site
  - Reflect the reality of the LHC Computing Grid and non-LCG world
  - Utilize not only the currently used relational DB (Oracle), but also non-relational structured storage (Hadoop)
- JEDI: the next ATLAS Production (Job Definition) System
  - Dynamic job definition with automatic adjustments of job specifications
  - Faster completion of tasks
    - Adjusting priorities and resource requirement for failed jobs -- may need more CPU time, memory or disk space
    - Re-running jobs to recover processed data lost due to hardware failures
  - and some ideas of new workflow (eg. master job & worker jobs)

# Summary

ATLAS Computing System has been running extremely well for the large-scale data processing, distribution and analysis

The Model and the System continue to evolve

- Flexibility and efficiency: from static to dynamic
  - data placement, conditions data distribution, software installation, data transfer path
- Network is the key for more flexible and efficient usage of the resources
- Trends of technologies in the computing industry are followed closely and used whenever possible