



Argonne
NATIONAL
LABORATORY

... for a brighter future



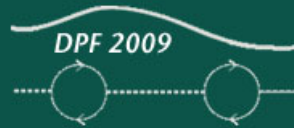
U.S. Department
of Energy

UChicago ►
Argonne_{LLC}



**Office of
Science**
U.S. DEPARTMENT OF ENERGY

A U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC



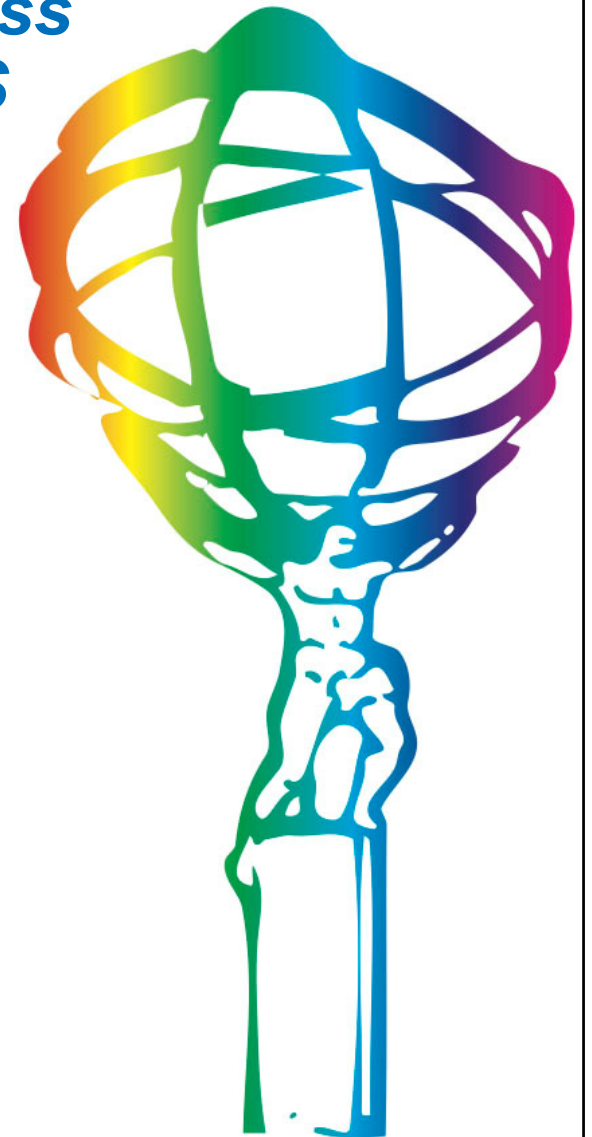
DPF 2009

2009 Meeting of the Division of Particles and
Fields of the American Physical Society (DPF 2009)

26-31 JULY 2009

Wayne State University, Detroit, MI

Scalable Database Access Technologies for ATLAS Distributed Computing



DPF2009, Wayne State University

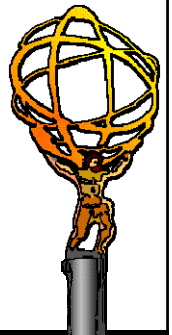
Detroit, Michigan, July 26-31

Alexandre Vaniachine

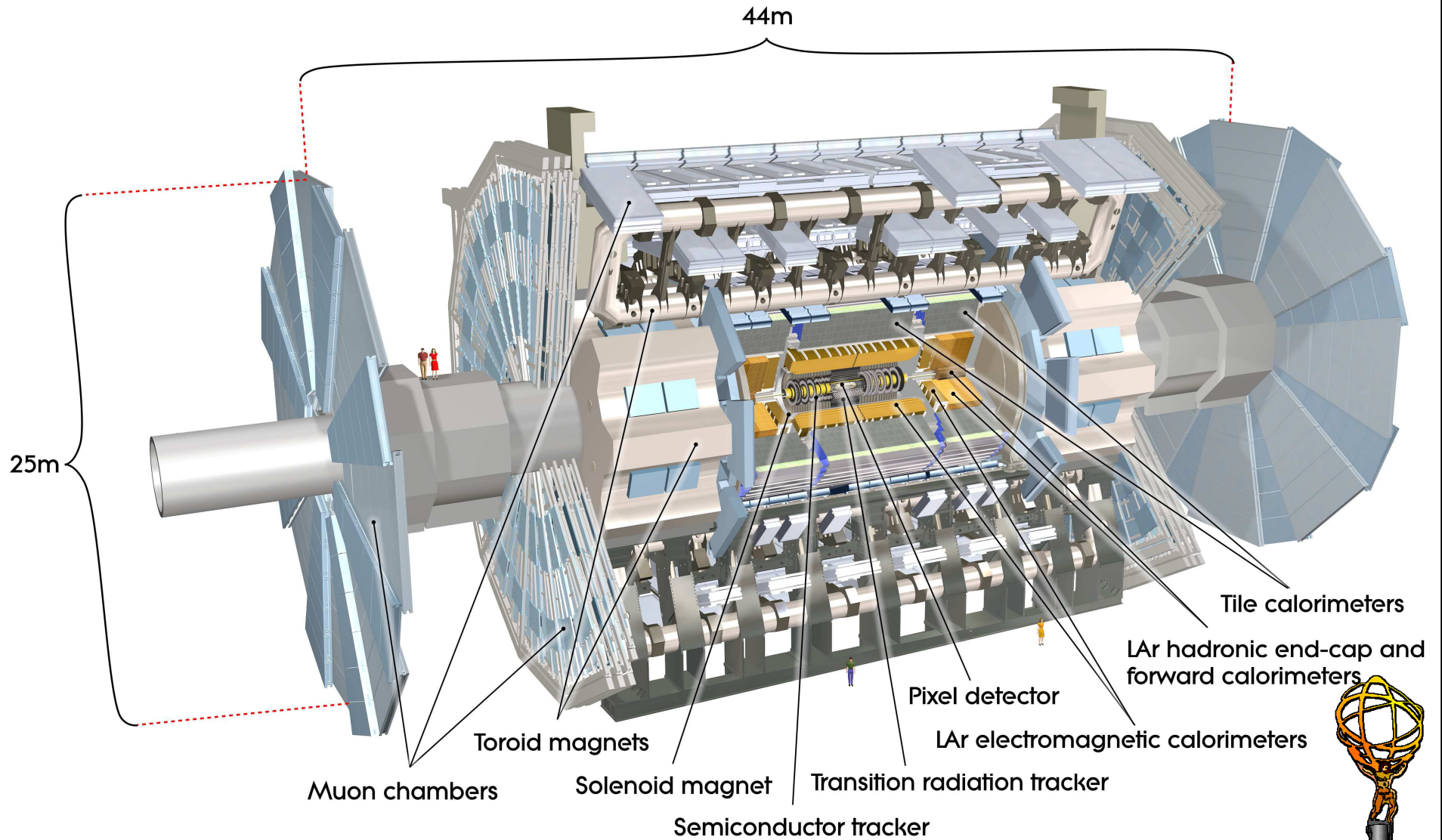
on behalf of the ATLAS Collaboration

Outline

- Complexity of the ATLAS detector is mirrored in our Conditions DB
- Data reconstruction – a starting point for any ATLAS data analysis
- Database access in data reconstruction
- Redundant database deployment infrastructure for Conditions DB
- Database access for user analysis
- Database access for Monte Carlo simulations
- Conclusions

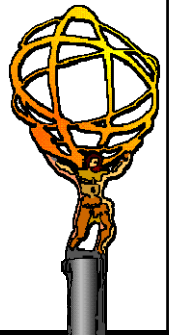


ATLAS Detector is Complex - Many Subdetectors

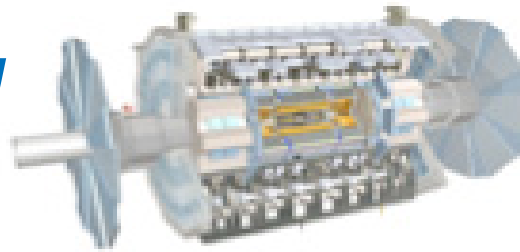


Managing Complexity

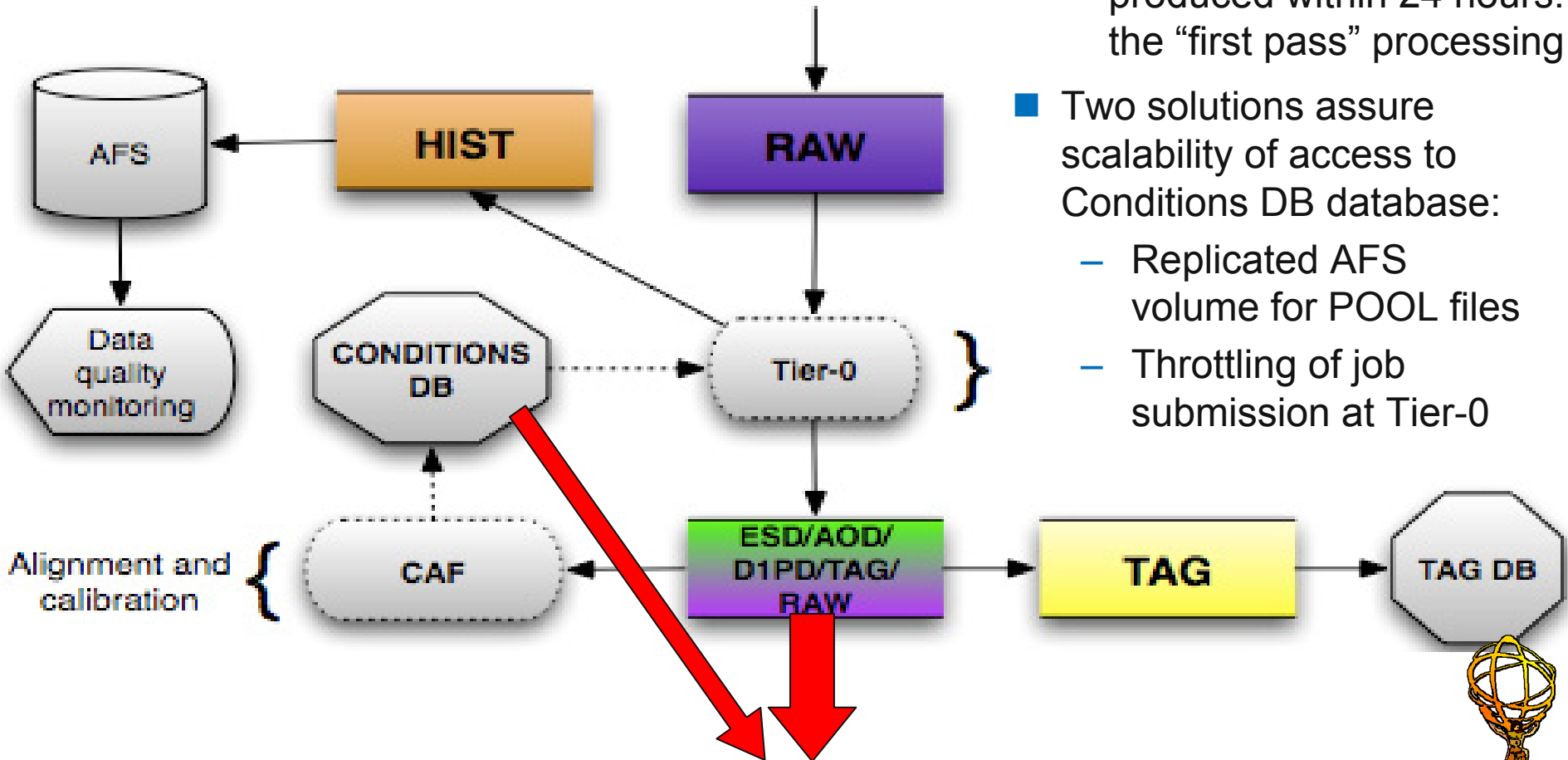
- Driven by the complexity of the detector the Conditions DB is complex:
 - It contains both database-resident information and external data in separate files, which are referenced by the database-resident data
 - *These files are in a common LHC format called POOL*
- ATLAS database-resident information exists in its entirety in Oracle but can be distributed in smaller slices of data using SQLite
 - a lightweight file-based technology
- Latest database access statistics provides some examples:
- These Conditions DB data are organized in 16 database schemas:
 - Two GLOBAL schemas (ONL/OFL) plus 1 or 2 per each subdetector
 - *Total of 747 tables organized in 122 folders plus system tables*
- 35 distinct database-resident payloads from 32 bit to 16 MB in size
 - Referencing 64 external POOL files in total
- To process a 2 GB file with 1K raw events a typical reconstruction job makes ~2K queries to read ~40 MB of database-resident data
 - Some jobs read tens of MB extra
 - Plus about the same volume of data is read from external POOL files



Offline Data Processing at CERN Tier-0



- Conditions DB is critical for data reconstruction at CERN using alignment and calibration constants produced within 24 hours: the “first pass” processing



- Two solutions assure scalability of access to Conditions DB database:
 - Replicated AFS volume for POOL files
 - Throttling of job submission at Tier-0

- To distributed computing facilities (the Grid) – see next slide



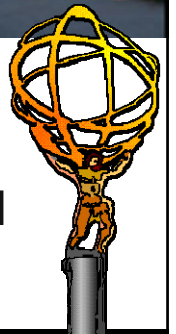
Where the LHC Findings Will Likely Be?

- “I’m willing to bet that when we do the first pass at the Tier-0, we won’t find anything definitive, not only because there will be little time but also because the calibrations and even algorithms will not be fully tuned”
- “The findings will likely be at the Tier-1s on the reprocessed data with refined calibrations and algorithms, and from analyses performed primarily at the Tier-2s”



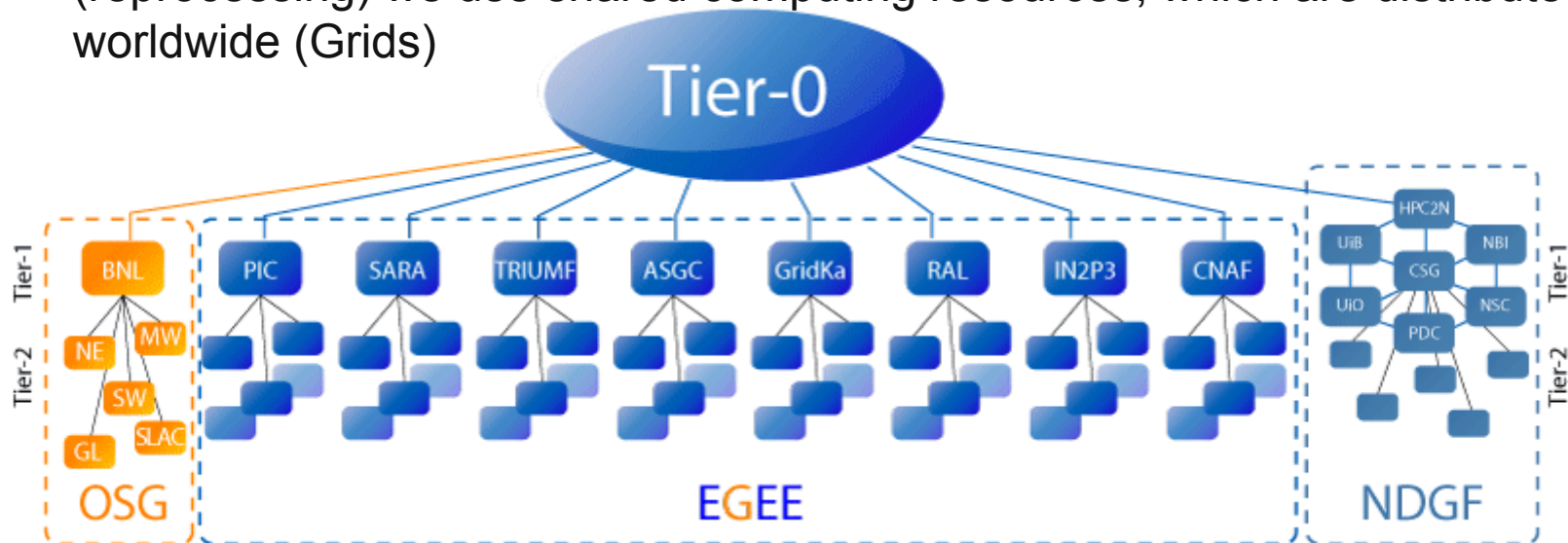
- Picture and Quote by Jamie Shiers

- Leader of the Worldwide LHC Computing Grid Support Group, CERN
 - <http://www.isgtw.org/?pid=1001318>

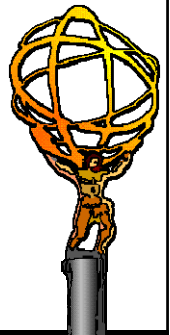


Distributed Computing Infrastructure

- Described in detail in the A. Stradling and A. Farbin talk in this session
- For data processing with improved alignment and calibration constants (reprocessing) we use shared computing resources, which are distributed worldwide (Grids)



- ATLAS uses three Grids (each with a different interface) split in ten “clouds” organized as large computing centers with tape data storage (Tier-1 sites) each associated with 5-6 smaller computing centers (Tier-2 sites). Plus more than a hundred of Tier-3 sites – this is a physicist’s own computing facility at the university or the department
 - ATLAS distributed computing power is six times higher than at Tier-0



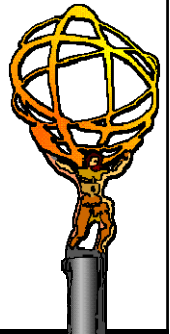
Scalable Database Access is Critical for Reprocessing

- The reprocessing at Tier-1 sites uses specific refined alignment and calibrations which are collected from subdetector groups, certified and versioned collectively
- Reprocessing improves the particle identification and measurements over the “first pass” processing at CERN

Both the Software Release build and the Database Release preparation are on a critical path in ATLAS reprocessing workflow shown on the chart

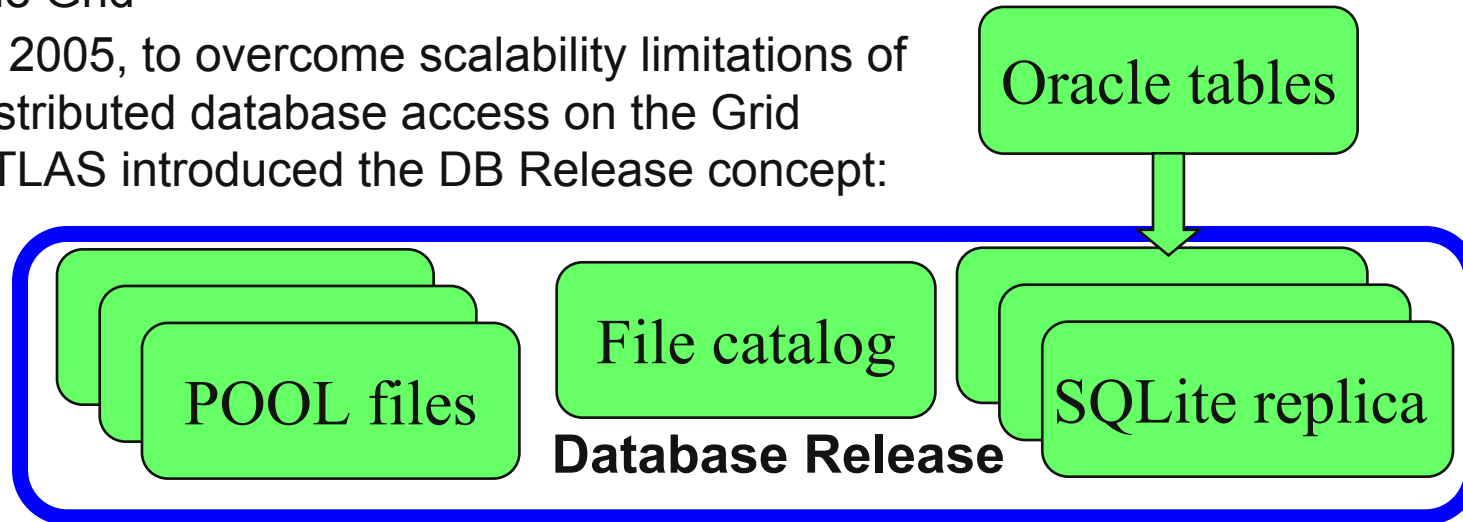


- ATLAS Database Release approach fully satisfies reprocessing requirements, which has been proven on a scale of one billion database queries during two reprocessing campaigns of 0.5 PB of single-beam and cosmics data on the Grid

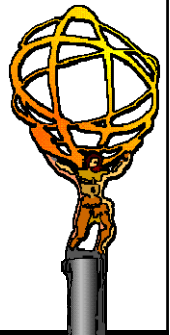


What is an ATLAS Database Release

- None of Tier-0 solutions for scalable database access are available on the Grid
- In 2005, to overcome scalability limitations of distributed database access on the Grid ATLAS introduced the DB Release concept:

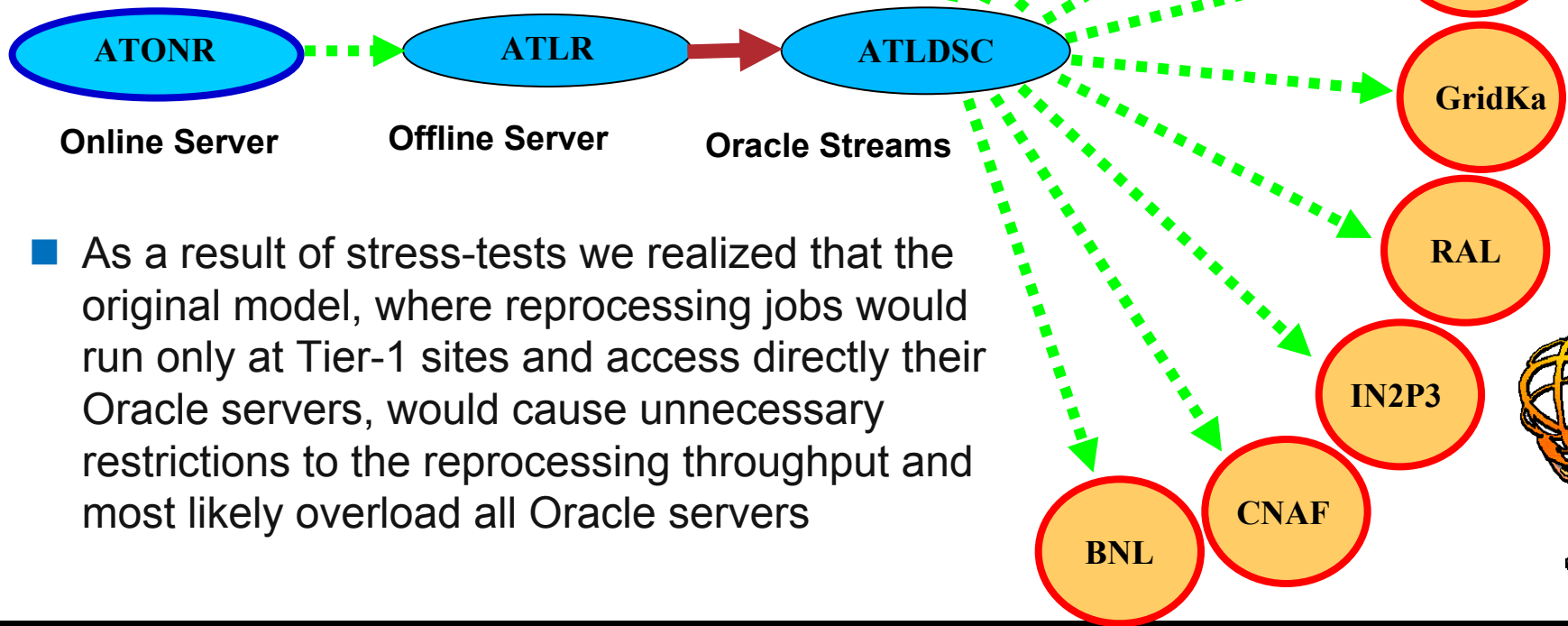


- Years of experience resulted in continuous improvements in the DB Release approach, which now provides solid foundation for ATLAS Monte Carlo simulation in production
- In 2007 the DB Release approach was proposed as a backup for database access in reprocessing at Tier-1s
- In a recent fast reprocessing campaign the DB Release encapsulated in a single dataset a 1 GB slice of the Conditions DB data from a two-week summer data taking period. The dataset was “frozen” to guarantee reproducibility of the reprocessing results.



Conditions DB Distribution

- In addition to Database Releases, ATLAS Conditions DB data are delivered to all ten Tier-1 sites via continuous updates using Oracle Streams technology
- To assure scalable database access during reprocessing ATLAS conducted Oracle stress-testing at the Tier-1 sites

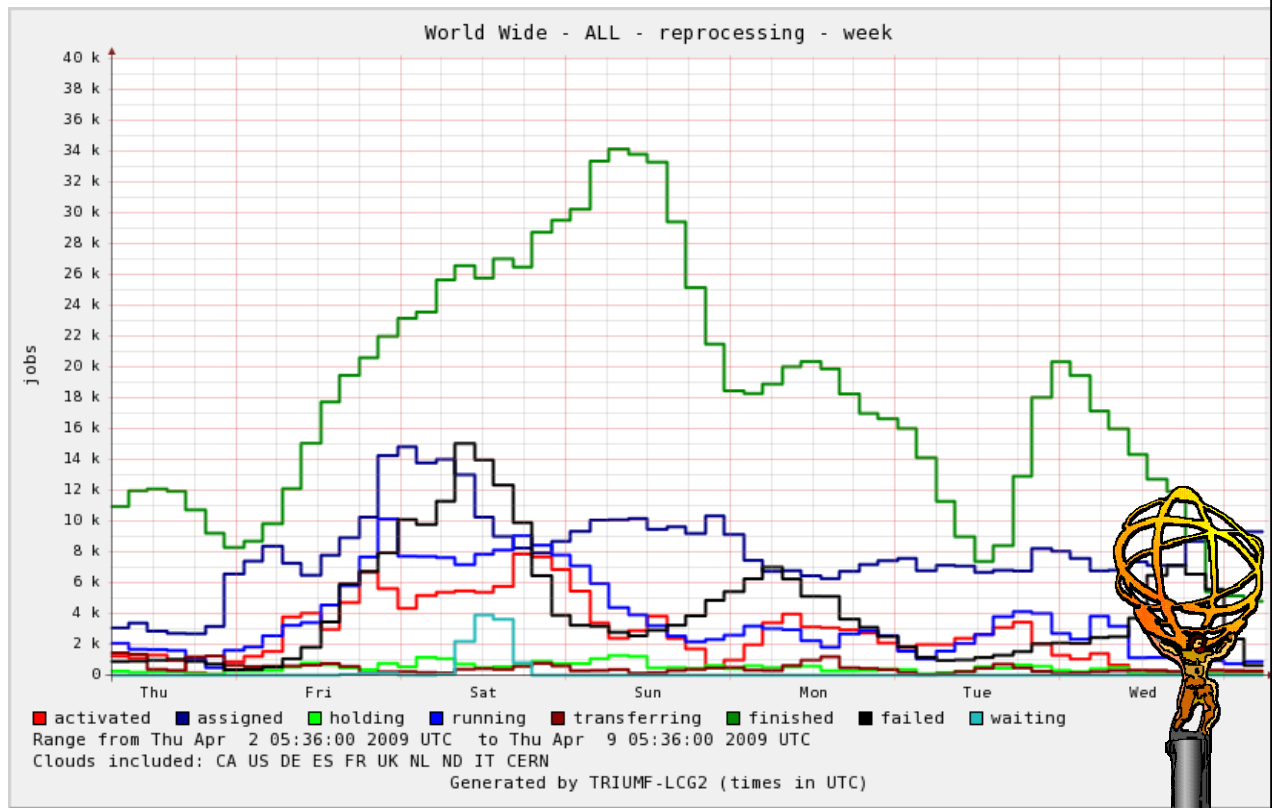


- As a result of stress-tests we realized that the original model, where reprocessing jobs would run only at Tier-1 sites and access directly their Oracle servers, would cause unnecessary restrictions to the reprocessing throughput and most likely overload all Oracle servers

Conditions DB Scalability Challenges in Reprocessing

Additional challenges exacerbated the main problem with Oracle overload:

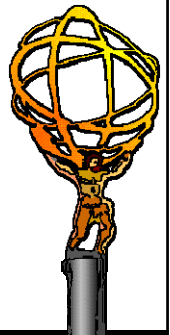
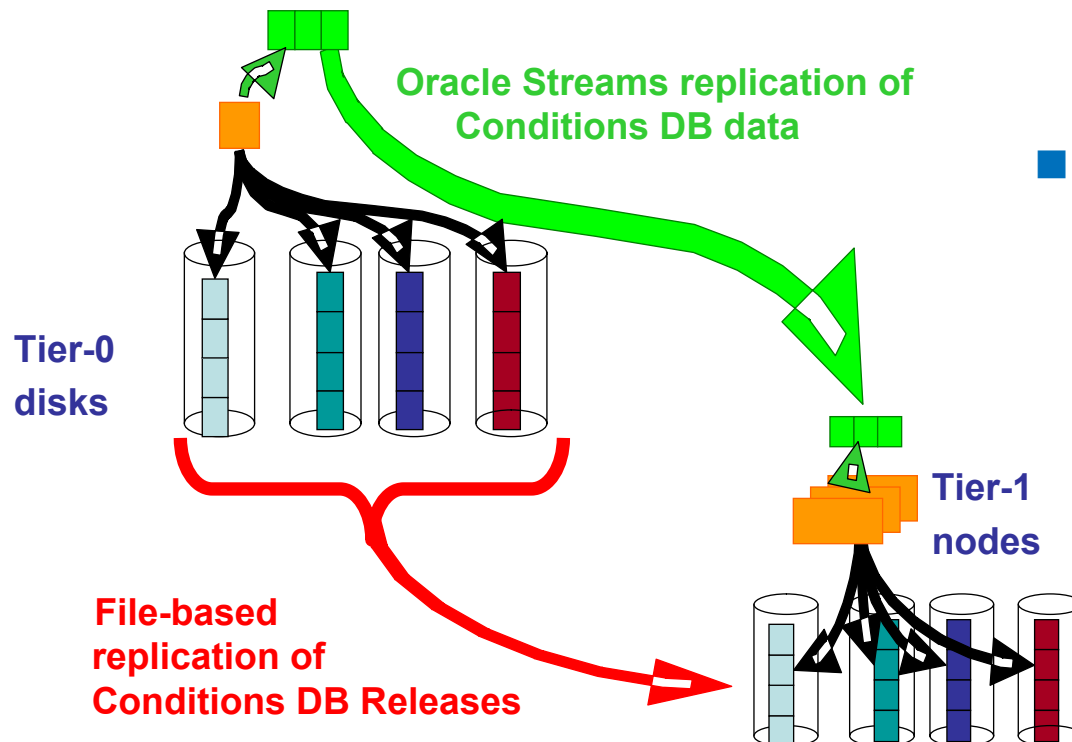
- Reprocessing jobs for the cosmics data are five time shorter than the baseline jobs reconstructing the LHC collision data
 - resulting in a fivefold increase in the Oracle load
- Having data on Tier-1s disks increased Oracle load six fold
 - in contrast to the original model of reprocessing data from tapes
- Combined with other limitations these factors required increase in scalability by orders of magnitude
- Thus, the DB Release approach, developed as a backup, was selected as a baseline



ATLAS Strategic Decisions for Reprocessing

- Read most of database-resident data from SQLite
- Optimize SQLite access and reduce volume of SQLite replicas
- Maintain access to Oracle
 - to assure a working backup technology, when required
- As a result of these decisions we overcome the Conditions DB scalability challenges in ATLAS reprocessing

- For the reprocessing we now have a robust but flexible technology for Conditions DB access
- By enabling reprocessing at the Tier-2 sites, the Conditions DB Release approach effectively doubled CPU capacities at the BNL Tier-1 site during Christmas reprocessing campaign

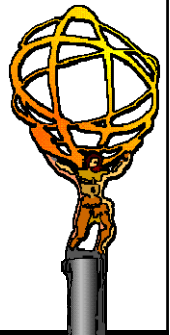


Redundant Database Deployment Infrastructure

- Since Conditions DB is critical for operations with LHC data, we are developing the system where a different technology can be used as a redundant backup, in case of problems with a baseline technology

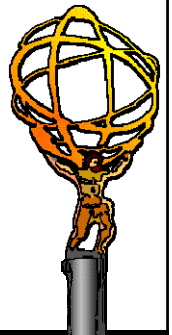
Use Case	Baseline	Backup
Reprocessing	SQLite	Oracle
Late-coming components	<i>db-on-demand</i>	Pilot Query
User Analysis	FroNTier	Oracle
Late-coming components and/or improvements	DoubleCheck, Software for deployment	Software for fast remote access

- Status of some late-coming components is reported in the following slides
 - *db-on-demand* is a system for automated Conditions DB Release packaging and validation
 - Pilot Query is a system for throttling job submission on the Grid
 - DoubleCheck is the FroNTier cache consistency solution for ATLAS
 - FroNTier is a web data caching system for scalable database access
 - *Initial implementation did not maintain its cache consistency*



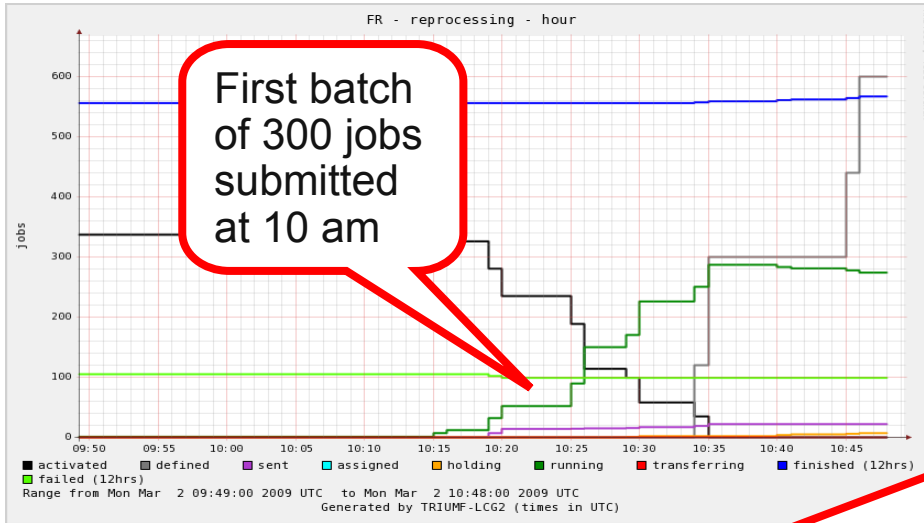
DoubleCheck: FroNTier Cache Consistency for ATLAS

- Piggybacking on recent CMS progress addressing the cache consistency problem, ATLAS resumed FroNTier development and testing in 2008
- In CMS case the cache consistency is checked for a single table at a time
 - This does not work for ATLAS, as most our queries are for two tables
 - *Hence the name DoubleCheck is chosen for the ATLAS solution*
- A major milestone in DoubleCheck development was achieved in July:
 - The proof-of-principle test demonstrated that CERN cache consistency solution for CMS can be extended to work for ATLAS
- With no showstoppers in sight, FroNTier development and testing in ATLAS continues increasing in scope and complexity
 - Details presented in the ATLAS talk by S. McKee in this session
- FroNTier/Squid deployment in U.S. ATLAS:
 - For better performance FroNTier at BNL is installed on two nodes
 - *with both Squids used as accelerators*
 - Redundant nodes established and working at Michigan and Chicago
 - Established, initial testing done at SLAC and Indiana
 - Initial phases at Boston, Harvard and UT Arlington



Pilot Query: Proof-of-principle Demonstrated

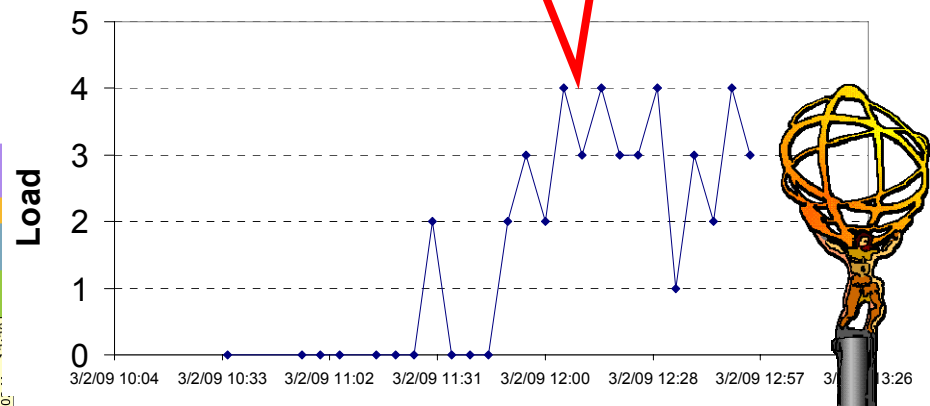
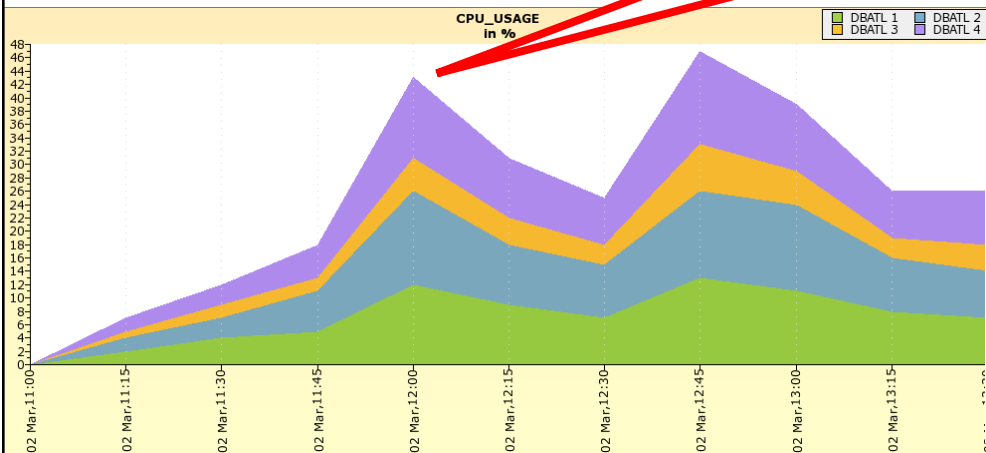
- Throttling Oracle server load on the Grid (at the Tier-1 site in Lyon)



- Development of the next generation Pilot Query system is complete and ready for testing

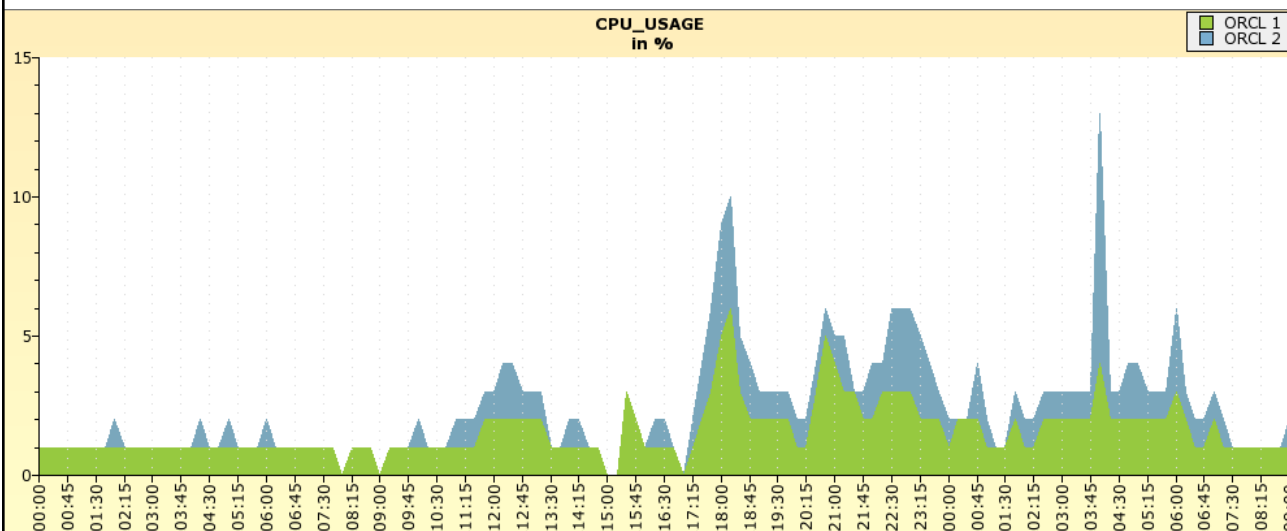
Monitoring shows Oracle load limited by the Pilot Query technology

Because we set ATLAS application-specific Oracle load limit at 4

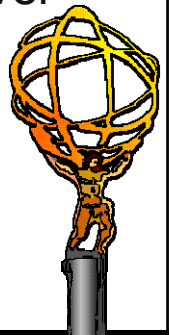


Use of Oracle Access in ATLAS Distributed Computing

- For years ATLAS Monte Carlo simulation jobs used SQLite replicas for access to simulated Conditions DB data
 - Recently, led by U.S. ATLAS efforts, Monte Carlo simulations are becoming more realistic by using access to real Conditions DB data
 - *this new type of simulation jobs requires access to Oracle servers*
- Number of jobs of this type that run at all U.S. ATLAS Tier-2s in June:
 - Michigan:[50](#) Boston:[20](#) Harvard:[9](#) IU_OSG:[10](#) MWT2_IU:[10](#)
Chicago:[20](#) Oklahoma:[38](#) SLAC:[23](#) UT Arlington:[10](#) UT Dallas:[10](#)
- All jobs finished successfully after accessing BNL Oracle server

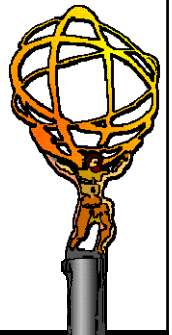


- The monitoring plot shows increase in load at the BNL Tier-1 Oracle server caused by these simulation jobs



Conclusions

- ATLAS has a well-defined strategy for redundant deployment of critical database-resident data
 - This strategy is based on the usage of the most suited technology for each use case
 - *ATLAS experience demonstrated that this strategy worked well as new unanticipated requirements emerged*
- ATLAS database deployment strategy scales well for reprocessing
 - The redundant database deployment infrastructure fully satisfies both Full Reprocessing and Fast Reprocessing requirements
 - *Steps being taken to assure that Oracle can be used as a backup in case of unexpected problems with the baseline approach*
- For scalable database access in user analysis the FroNTier technology is undergoing development in collaboration with U.S. CMS
- Each major ATLAS use cases is functionally covered by more than one of the available technologies, so that we can achieve a redundant and robust data access system, ready for the challenge of the first impact with LHC collision data





Argonne
NATIONAL
LABORATORY

... for a brighter future



U.S. Department
of Energy

UChicago ►
Argonne_{LLC}



**Office of
Science**

U.S. DEPARTMENT OF ENERGY

A U.S. Department of Energy laboratory
managed by UChicago Argonne, LLC

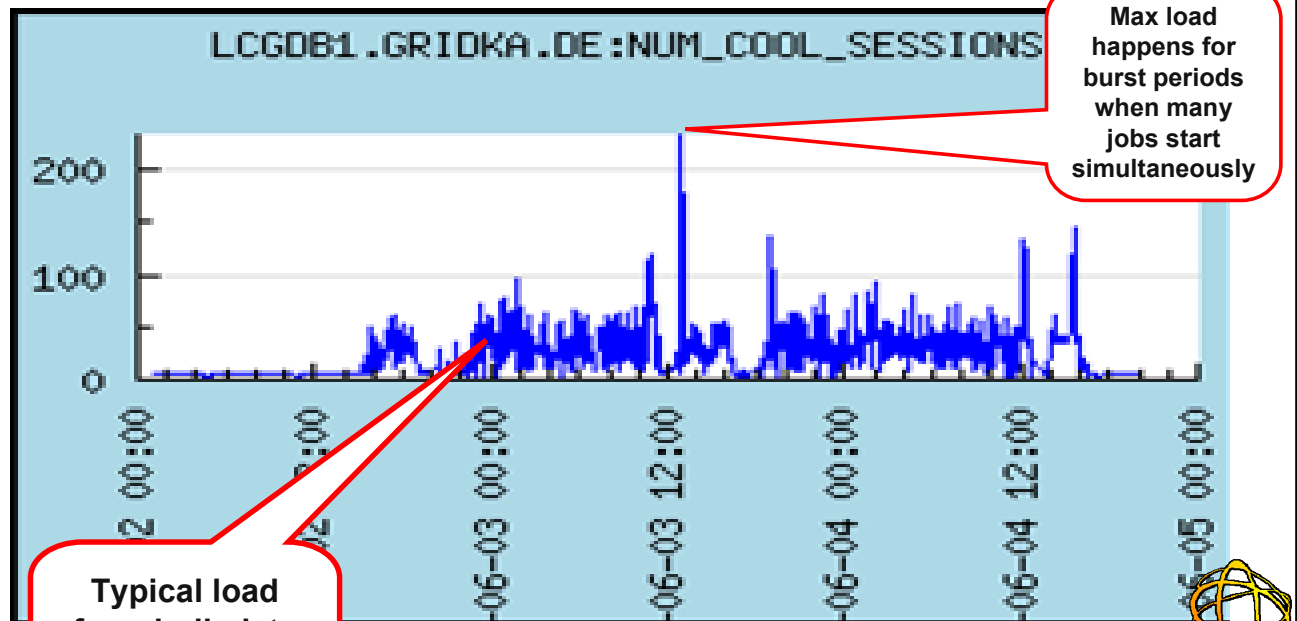
Backup



Peak Loads are Typical in Database Access on the Grid

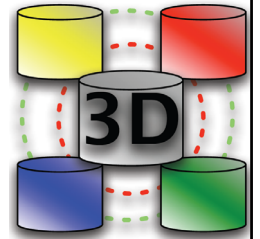
- Capacities supporting nominal throughput are not sufficient on the Grid
- In distributed data processing one must take into account the chaotic nature of Grid computing characterized by peak loads, which can be much higher than nominal access rates
 - DPF2004
- Instabilities at Tier-1 sites may result in peak database access loads when many jobs are starting at once
- This may create overload of Oracle servers and degrade Oracle Streams replication worldwide

Monitoring CCRC'08 Bulk Data Reprocessing at Tier-1

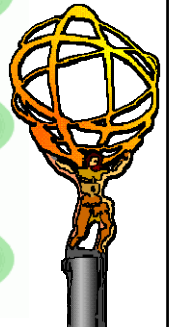
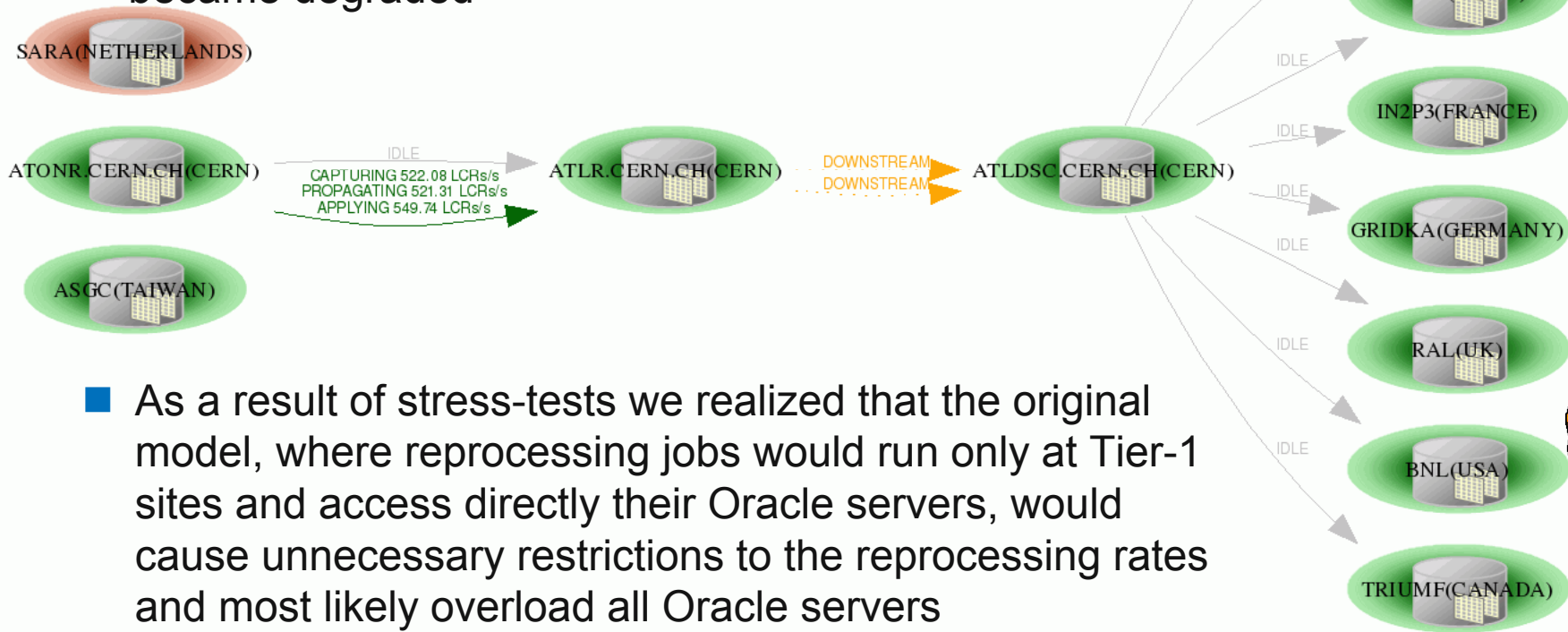


Note: Summary load on both Oracle RAC nodes at FZK Tier-1

Problem: Escalation of WLCG 3D Incidents



- In 2008, Oracle overload was experienced at all five Tier-1 sites tested
- During overload, Oracle Streams updates of Conditions DB data to this Tier-1 site are degraded for hours
- After several hours of Oracle overload at one Tier-1 site, Conditions DB updates to **all** other Tier-1 sites became degraded



- As a result of stress-tests we realized that the original model, where reprocessing jobs would run only at Tier-1 sites and access directly their Oracle servers, would cause unnecessary restrictions to the reprocessing rates and most likely overload all Oracle servers

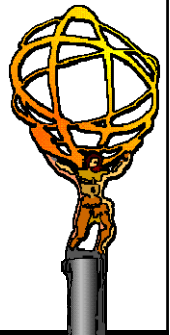
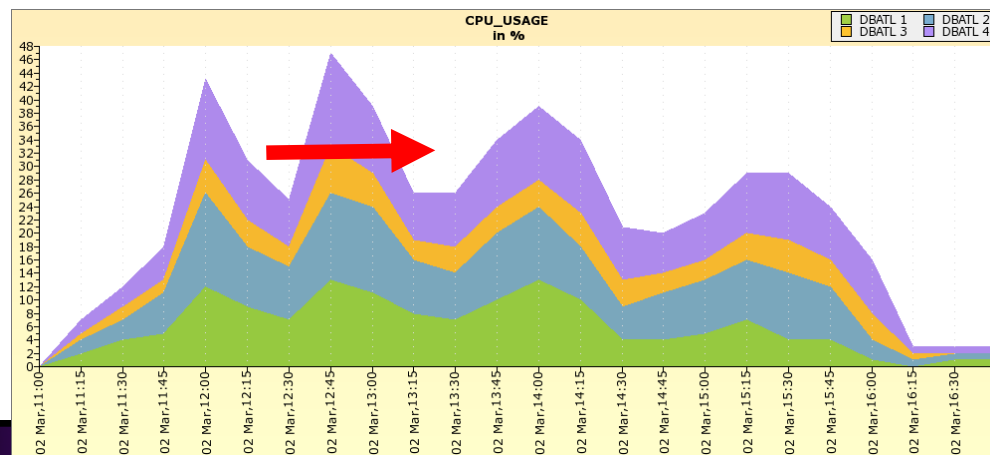
How Does ATLAS Pilot Query Work?

- Sample output from the finished job:
running on ccwl0613 on Mon Mar 2 13:06:01 2009
Database operations pilot at LYON
pilot detected status GO - Load: 03.10 Sessions: 813 Threshold:4

- An example of a job held at a lower threshold:

Database operations pilot at LYON

pilot detected status NOGO - Load: 03.00 Sessions: 408 Threshold:2
Mon, 02 Mar 2009 12:21:35 avoiding load of 03.00 at 408 concurrent COOL sessions
pilot detected status NOGO - Load: 02.10 Sessions: 477 Threshold:2
Mon, 02 Mar 2009 12:28:23 avoiding load of 02.10 at 477 concurrent COOL sessions
pilot detected status NOGO - Load: 02.70 Sessions: 483 Threshold:2
Mon, 02 Mar 2009 12:44:17 avoiding load of 02.70 at 483 concurrent COOL sessions
pilot detected status NOGO - Load: 02.90 Sessions: 780 Threshold:2
Mon, 02 Mar 2009 13:08:27 avoiding load of 02.90 at 780 concurrent COOL sessions
pilot detected status GO - Load: 01.90 Sessions: 673 Threshold:2



Pilot Query: Throttling Jobs Submission at Tier-1s

- Nominal throughput is not enough
 - Instabilities at Tier-1 sites result in peak Oracle loads when many jobs are starting at once
 - *peak loads can be much higher than the nominal load*
- Oracle overload at one site may result in a worldwide degradation of ATLAS data distribution via Oracle Streams
 - To prevent that from happening we must throttle jobs submission at Tier-1 sites using Pilot Query
- Development of the next generation Pilot Query system is complete and ready for testing

ATLAS COOL Pilot Query Monitoring Page

Information about the Pilot Query: [CoolReprocessingTests TWIKI](#)

you are connected to: **PIC**

Select Time Interval for graph generation

From Date: Day: 1 Month: 6 Year: 2009 Hour: 9 Min: 2

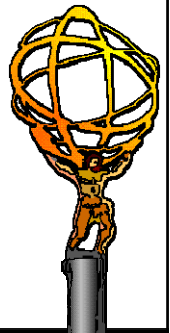
To Date: Day: 8 Month: 6 Year: 2009 Hour: 9 Min: 2

Overview	Min, Max, Mean
# tickets waiting 0	Mean time requested-called in last hour no info
# tickets completed 8	Min time requested-called in last hour no info
# tickets completed in last hour 0	Max time requested-called in last hour no info
# tickets completed in last 24h 0	Mean time requested-called in last 24h no info
	Min time requested-called in last 24h no info
	Max time requested-called in last 24h no info
List of tickets waiting	Mean time requested-called last 10 tickets 56
List of tickets completed	Min time requested-called last 10 tickets 45
	Max time requested-called last 10 tickets 90

Clients that left	Current DB load
Clients having left (in total) 4	load as in gv\$osstat .738
Clients having left in the last hour 0	
Clients having left in the last 24h 0	
List of clients having left	

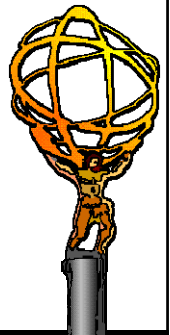
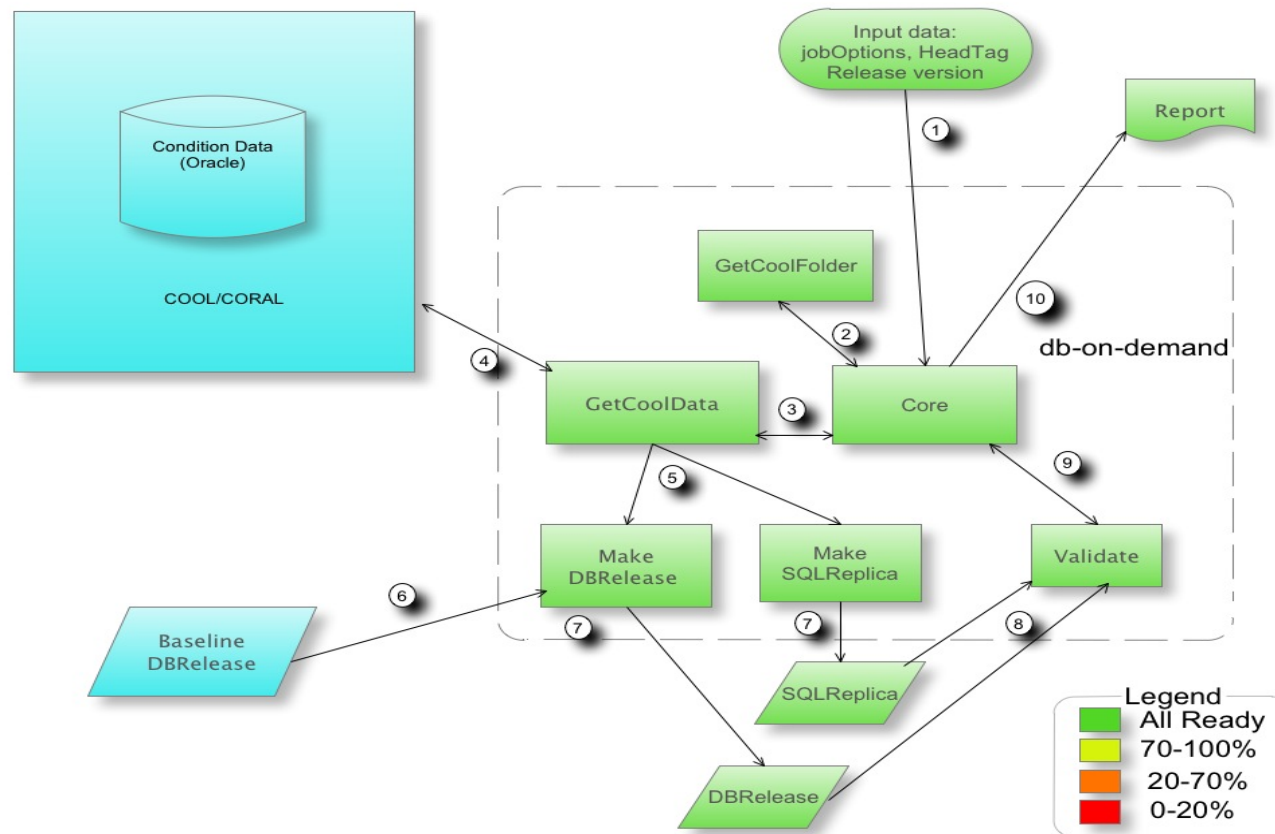
Current Pilot Query Parameters:

PARAMETER	VALUE
RUN	1



Progress in db-on-demand Development

- Integration with ATLAS production system started
 - New use cases have been indentified and tested
 - *System was used during Fast Reprocessing in July*



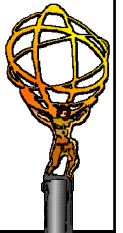
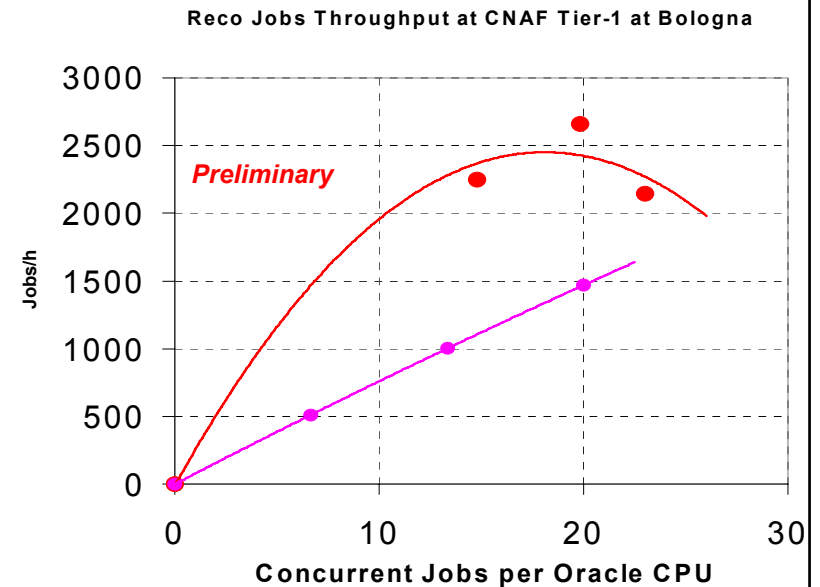
ATLAS Oracle Scalability Testing

- The goal of database scalability testing is to detect hardware limits of Oracle servers deployed at the Tier-1 sites, so that the server overload conditions can be safely avoided in a production environment
- First tests showed that Oracle capacities are sufficient for expected nominal jobs throughput
- Recent tests and operational experience in 2009 confirmed our expectations

First Oracle Scalability Tests in 2007

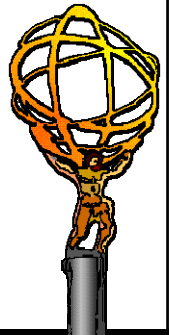
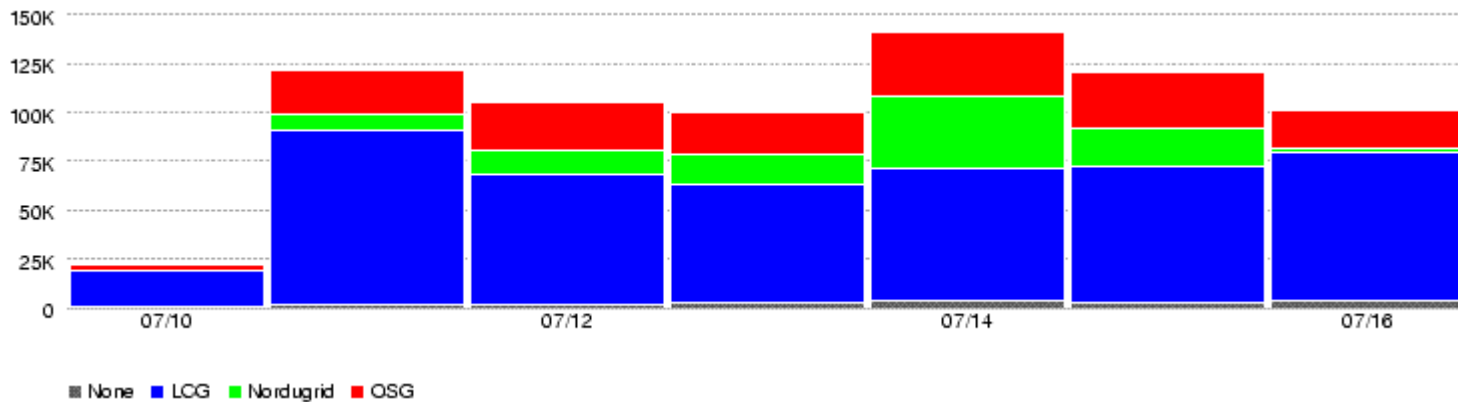
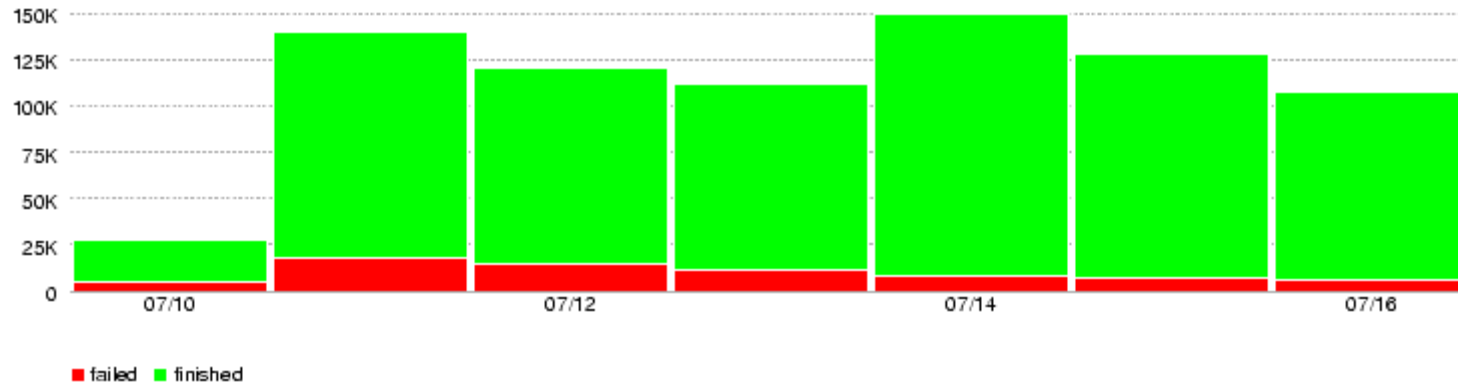
- Test jobs read realistic Conditions DB data workload at random
- We estimate that ATLAS daily reconstruction and/or analysis jobs rates will be in the range from 100,000 to 1,000,000 jobs/day
- For each of ten Tier-1 centers that corresponds to the Conditions DB access rates of 400 to 4,000 jobs/hour

- Thus, preliminary results from the first scalability test were promising
 - We got initial confirmation that ATLAS capacities request to WLCG (3-node clusters at all Tier-1s) is close to what will be needed for reprocessing in the first year of ATLAS operations



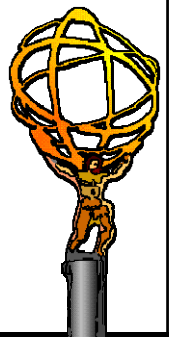
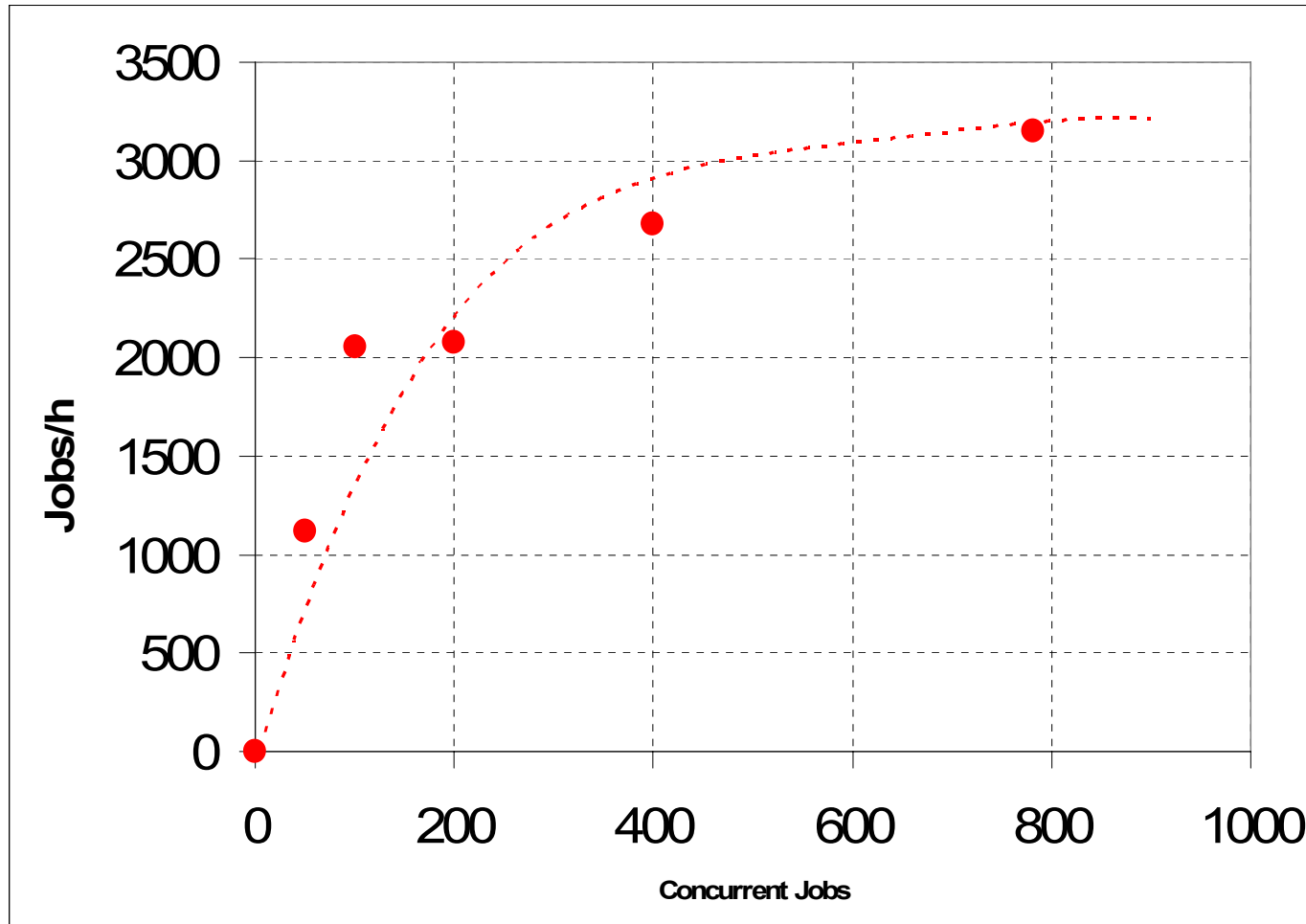
ATLAS Production Rates

- In agreement with our 2007 expectations, current ATLAS production rates reached levels above 100,000 jobs/day



Latest ATLAS Scalability Test Results at PIC

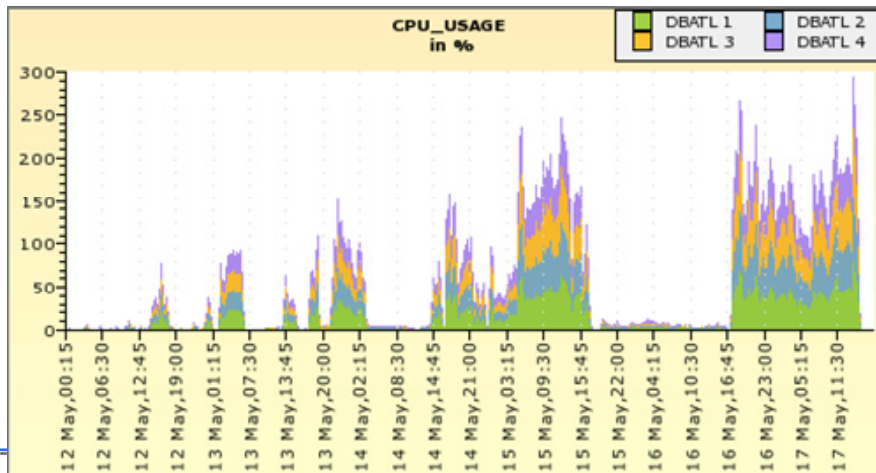
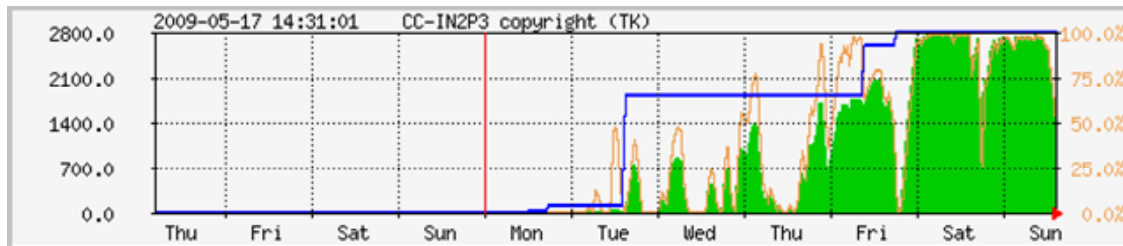
- 2009 scalability test results are in agreement with our previous findings:
 - Oracle capacities are sufficient for expected nominal jobs throughput



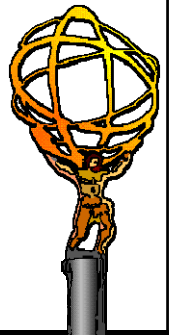
Latest Nominal Throughput Test at CC-IN2P3

- In addition to scalability tests, a comprehensive test to validate Oracle capacities deployed at Lyon Tier-1 site has been done

Summary of the week



- Lyon test confirmed that ATLAS used correct projections in our WLCG request for Oracle capacities deployed at the Tier-1s
 - Tier-1 Oracle capacities are correctly provisioned for projected a nominal jobs throughput

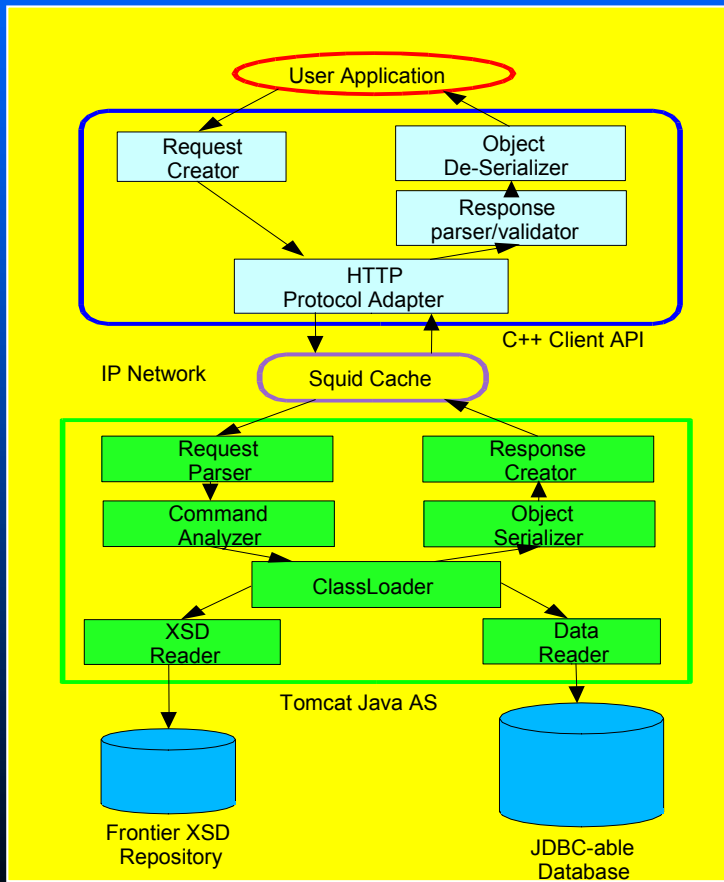


FroNTier Evaluation in ATLAS



Frontier & POOL (Simplified)

Sergey Kosyakov



- In collaboration with CMS, ATLAS started an evaluation of the promising FroNTier technology in 2006
 - Our tests found good performance gains of cached data
- Later tests found that because FroNTier/squid itself does not maintain its cache consistency, considerable efforts must be spent to assure that ATLAS applications obtain stable results in the case of ongoing changes to the Conditions DB
- Piggybacking on recent developments addressing the cache consistency problem, ATLAS resumed FroNTier development and testing in 2008

Slide by L. Lueking, 3D Workshop, December 14, 2004

New Approach to FroNTier Cache Coherency Problem

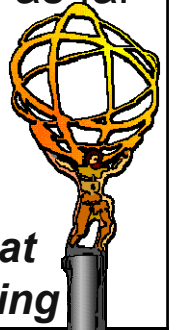
- Piggybacking on recent CMS developments addressing the cache consistency problem, ATLAS resumed FroNTier development and testing
 - Facilitating that, our U.S. CMS colleagues share their FroNTier experience, participate in common discussions, etc.



Cache coherency problem

- Need to get changes propagated to all sites
- First CMS solution: “short” queries (for data that could change) & “long” queries (for data that wasn't supposed to change), didn't work well:
 - too long: ~day for “short”, ~year for “long”
 - much shorter would overwhelm servers
 - sometimes data in “long” queries changed even though it wasn't supposed to, forcing us to flush all caches
- New solution based on http's Last-Modified and If-Modified-Since headers

- Recently, CERN IT/DB found a workaround for the Oracle bug affecting the proposed CMS solution:
 - Oracle PL/SQL script reduces probability that job gets stale data from the squid web cache
- Our U.S. CMS colleagues have been very open as far as their development and are a good source of advice



Slide from D. Dykstra talk at the ATLAS Database Meeting

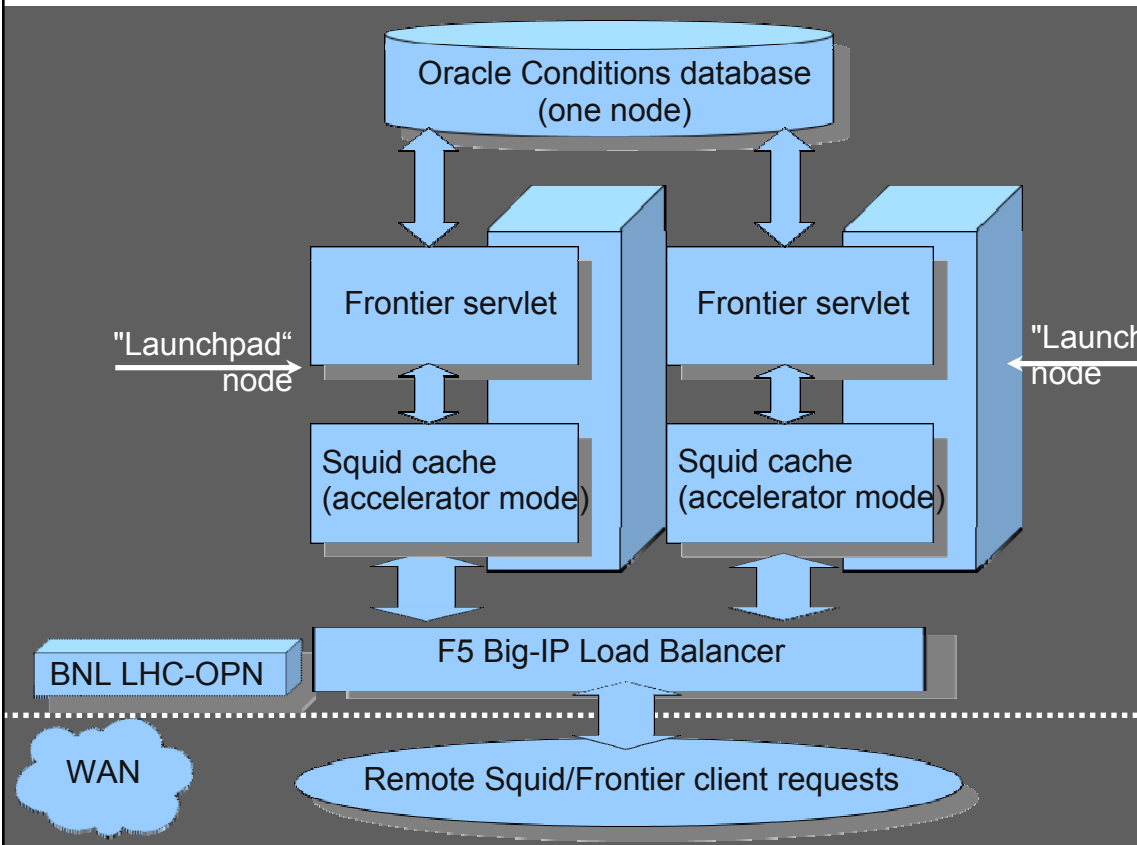
03/03/09

Frontier new cache coherency

4

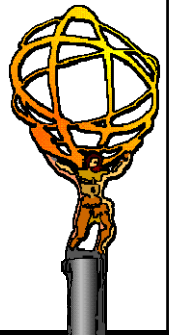
U.S. ATLAS Project on FroNTier Evaluation

- First T1/T2 tests executed at BNL/AGLT2 in the context of reprocessing tasks spotted initial problems and identified various ways for improvement
 - For better performance FroNTier at BNL was installed on both nodes,
 - *with both Squids used as accelerators:*



Squid deployment at U.S. T2s:

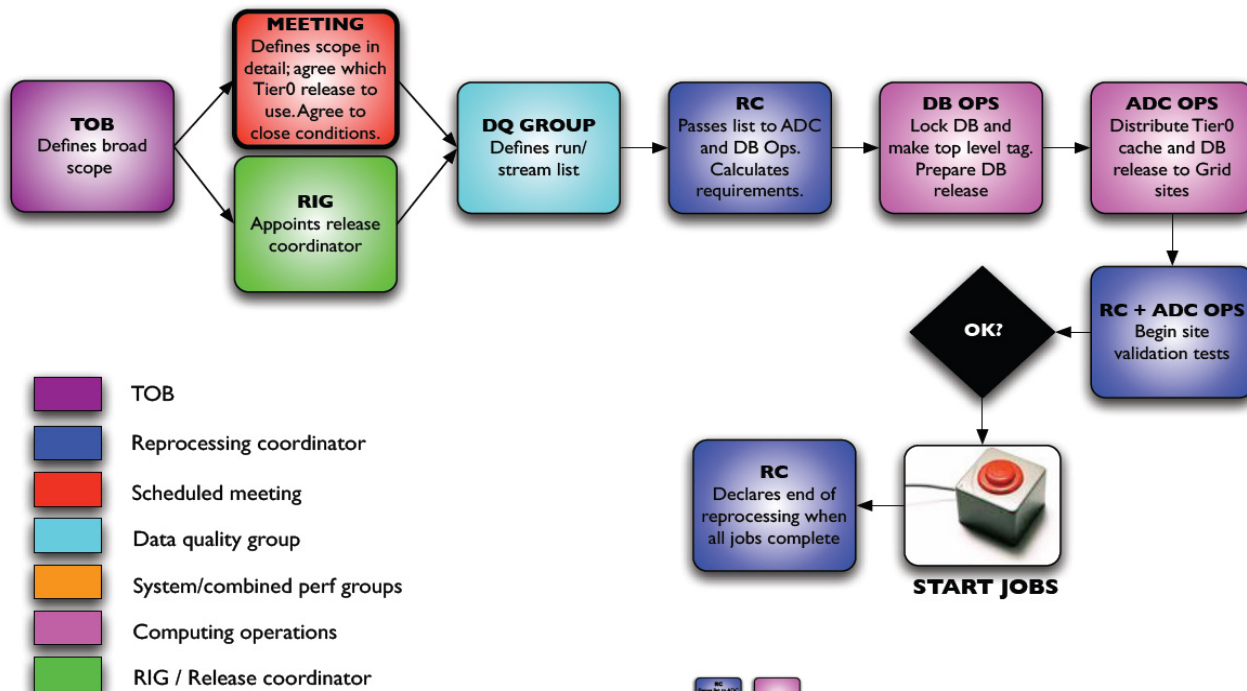
- Redundant nodes established, working, and tested at AGLT2 (Michigan), MWT2 (Chicago)
- Established, initial testing done: WT2 (SLAC), MWT2 (Indiana)
- Initial phases at NET2 (Boston, Harvard), SWT2 (UTA)
- FroNTier testing is a very active area of rapid developments



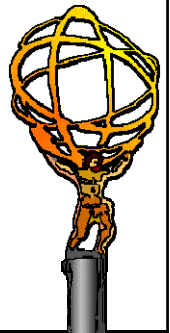
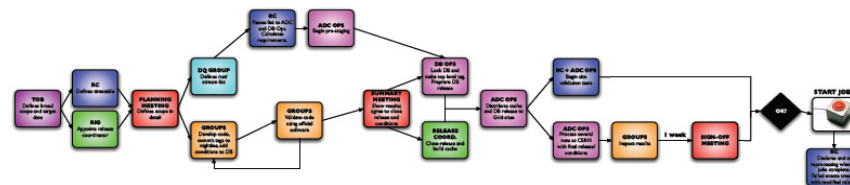
Fast Reprocessing

- As well as reprocessing very large quantities of data (Full Reprocessing) we are doing reprocessing of smaller amounts of data, much quicker, to give fast feedback to systems and groups
- In July we exercised Fast Reprocessing of cosmics data taken in 2009

- Fast Reprocessing started within several days after the end of a two-week data taking period
- 0.3 PB of data were reprocessed on the Grid within one week



for comparison - full procedure:



Summary: Getting Ready for LHC Data Taking

Reprocessing:

- SQLite:
 - *db-on-demand* undergoes integration with production system
- Oracle:
 - Oracle capacities deployed at the Tier-1 sites are been validated for the nominal throughput
 - Next generation Pilot Query system to prevent Oracle overload is ready for testing

User Analysis:

- FroNTier:
 - The proof-of-principle test of demonstrated that ATLAS DoubleCheck cache consistency solution for FroNTier works
 - Software to support FroNTier is on track for delivery for LHC data
- Oracle:
 - Good prospects for remote access performance improvements

Simulations:

- Both SQLite and Oracle are in use

