



**Argonne**  
NATIONAL  
LABORATORY

*... for a brighter future*



U.S. Department  
of Energy

UChicago ►  
Argonne<sub>LLC</sub>



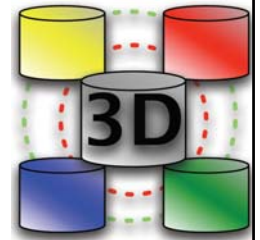
**Office of  
Science**  
U.S. DEPARTMENT OF ENERGY

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

## *Joint Task Force on ATLAS DB Performance*

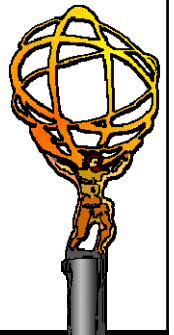


*WLCG Distributed Database Operations Workshop  
CERN, Geneva, Switzerland  
November 11-12, 2008  
Alexandre Vaniachine (ANL)*

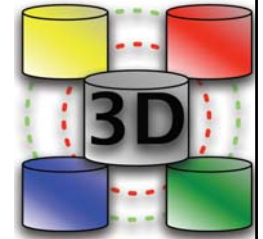


## Conditions DB Development and Deployment Cycle

- For access to Conditions DB data ATLAS experiment adopted
  - Common LHC software: COOL/CORAL
    - *The COOL/CORAL development cycle takes long time (months)*
  - Common WLCG 3D Services: Oracle Streams at Tier-1 sites
    - *Changes in Oracle RAC hardware configurations requires careful studies prior to approval*
    - *Oracle RAC hardware deployment at the Tier-1 sites takes more than six months*
- Given long lead times, ATLAS must detect any DB limitations *in advance*
  - Detecting limitations in software, server configuration, and hardware is a main purpose of ATLAS database stress tests
    - *Database stress tests on the Grid provide most valuable results*
- The goal of database stress tests is to detect scalability limits of the hardware deployed at the Tier-1 sites,
  - so that the server overload conditions can be safely avoided in a production environment



# ATLAS DB Performance Studies



- The very first ATLAS DB stress tests detected limitations in software
  - based on ATLAS feedback, the next COOL version was optimized
- Next stress tests found that the server memory may become a bottleneck

## Proposed Targets for 3D Milestones in CCRC08-2

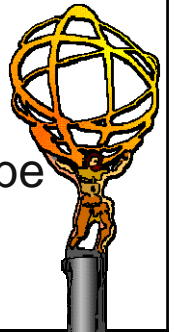
- Last LCG Management Board requested to collect any remaining ATLAS milestones for the 3D setup
- Reprocessing during CCRC08-2 in May provides an opportunity to validate that the 3D deployment is ready for LHC data taking

Tier-1 site	Sessions/Node
TRIUMF	40
FZK	80
IN2P3	60
CNAF	40
SARA	140
NDGF	100
ASGC	80
RAL	90
BNL	220
PIC	40

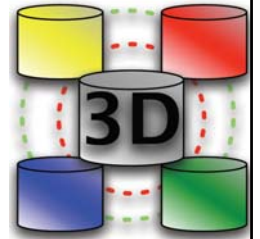
- Robust Oracle RAC operation at the Session/Node target demonstrates that 3D setup at the site is ready for ATLAS data taking
  - Achieving 50% of the target shows that site is 50% ready
- In addition, during ATLAS reprocessing week DBAs at the Tier-1 sites should collect Oracle performance metrics: CPU, IO and cache hit ratios
  - Details will be presented to DBAs at the Database Track



- ATLAS memory requirements became a target for the remaining 3D milestone during CCRC'08
- During the week of ATLAS data reprocessing the cumulative ATLAS 3D milestone target for CCRC'08 was exceeded by 30%
- CCRC'08 did not probe data-intensive COOL access to DCS data

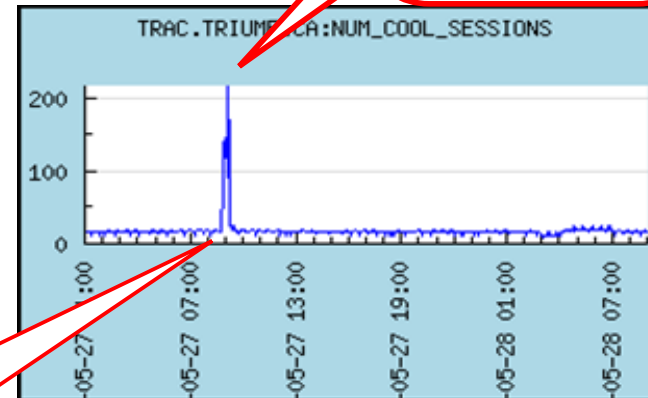
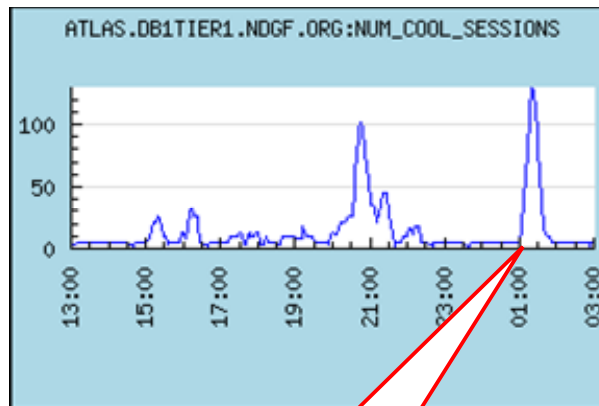


# ATLAS CCRC'08 Experience



- At the CCRC'08 Post Mortem Workshop ATLAS expressed first concerns on peak loads typical for a Grid computing environment:

## Monitoring Burst Load at Tier-1s

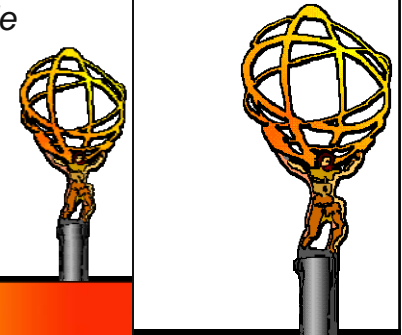


TRIUMF  
CPU load at peak:  
■ 98% peak for 1 min  
■ 70% average

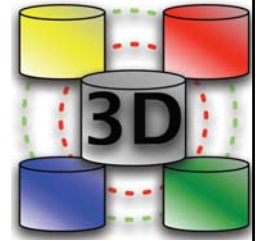
Short burst periods correspond to submission of limited number of test jobs for data reprocessing tasks

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

- ATLAS does not have spare WLCG 3D Oracle capacities deployed to support burst loads when many data reprocessing jobs start at once on an empty cluster

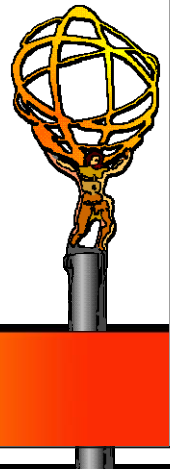


## ATLAS DB Stress Tests Announcement in June

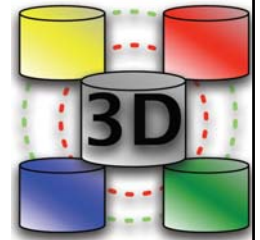


- During CCRC'08 we did not obtain requested Oracle performance metrics (CPU, IO and cache hit ratios) during reprocessing, thus, to find
  - What are the scalability limits and database access patterns of ATLAS reprocessing jobs under realistic conditions?at the WLCG CCRC'08 Post-Mortem Workshop ATLAS announced that Oracle stress tests will be added to reprocessing:

- Since there is little room for an error in our estimates of 3D capacities required for ATLAS
  - We have to validate these estimates using latest ATLAS software and Computing Model
- ATLAS validation of 3D services will continue with further data reprocessing exercises increasing in scope and complexity:
  - We plan to complete testing using latest FDR-2 data samples
  - We will test all Tier-1 sites using access to the DCS data
    - *Expect higher load from new DCS tests at all Tier-1*

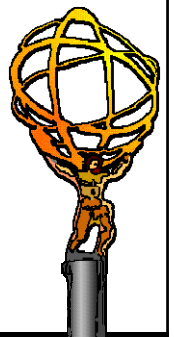




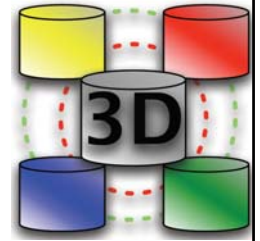


## DB Stress Tests during FDR-2 Reprocessing Exercise

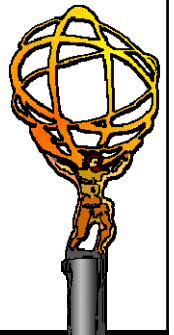
- In our first exercise called “Bulk Reprocessing” we processed the FDR-2 data at most of the Tier-1 sites
  - We produced ESD, AOD, DPD, and later TAGs and merged AODs
    - *Each ATLAS data processing step involves database access*
- Database infrastructure worked well during first ATLAS exercise in Bulk Reprocessing of the FDR-2 data
  - The proposed Conditions DB Release technology (for COOL payload files distribution) has been successfully proven in a large scale test
  - 3D streaming worked very well - all reprocessing jobs running world-wide found their Conditions DB data at the Tier-1 Oracle
  - Due to a chaotic nature of Grid computing the peak database load can be much higher than the required average
    - *During “Bulk Reprocessing” such overload observed at one Tier-1 where concurrent jobs peaked at 11 times above the average load*
- Despite immediate steps taken to upgrade the site hardware as well as ATLAS software that were made to reduce probability of such condition
  - ATLAS decided to develop technology to avoid peak database load



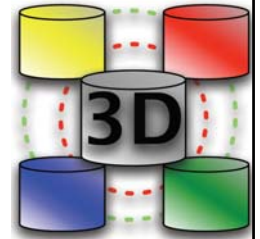
# Database Stress Tests during “Realistic Reprocessing”



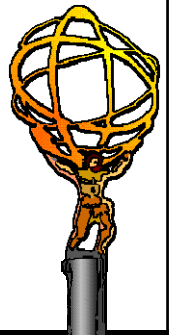
- Next ATLAS larger-scale exercise was called “Realistic Reprocessing”
  - Larger volumes of simulated FDR-2 data were used for reprocessing at three Tier-1 sites: NDGF, SARA, TRIUMF
- Instead of been “Realistic” this was a “Worst Case Scenario” test
  - Our added database stress test put too much load for Oracle servers
    - *ATLAS now developed an improved Realistic Test*
  - Athena software has not been updated and was using another “worst case scenario” database access pattern
    - *With a new improved Athena database access strategy (faster sessions and tunable 10-minute in-memory Athena caches for COOL DCS data) this problem is now gone*
- Launch of database stress tests has not been coordinated with the DBAs
  - Proper monitoring procedures are put in place to alert the DBAs
- Despite deficiencies the DB stress tests provided critical knowledge:
  - In ATLAS workflow the IO load is more important than the CPU load
    - *There is no longer a need to add database stress tests to ATLAS jobs, since jobs now read realistic Conditions DB data volumes*



## Database Performance Limitations and Actions Taken

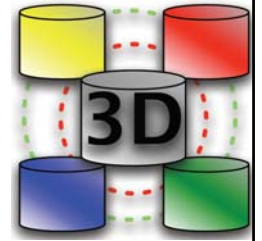


- During ATLAS database stress tests Oracle scalability limits were detected at all five Tier-1 sites tested: SARA, NDGF, GridKa, TRIUMF and PIC
  - In all cases IO limits were reached at peak loads for database access at the number of concurrent jobs that were 5-10 times higher than the required average access rates for each of these sites
    - *In few hours IO saturation at one Tier-1 site may degrade Oracle Streams updates to all other Tier-1 sites, thus, it must be avoided*
- Database peak load avoidance technology has been prototyped and tested
  - Our approach based upon the idea for “pilot” job submission on the Grid:
    - *Instead of the actual query ATLAS will send the “pilot” query first*
- To avoid limitations due to the ORA-01555 ‘snapshot too old’ error an improved database access strategy is now implemented in new ATLAS software version
- Changes in WLCG 3D Services configuration setting were proposed
  - avoid 15-min limit for duration of queries (by G. Dimitrov and F.

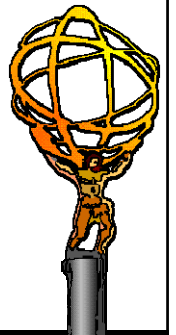




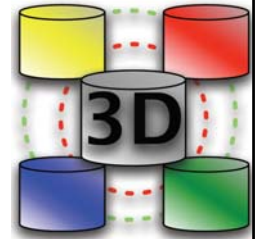
## Avoiding Database Peak Loads in Grid Computing



- In distributed data processing scenarios one must take into account the chaotic nature of Grid computing
  - To validate database performance at peak loads, we stress tested databases at concurrent jobs rates that are much higher than required average rates
    - *This has been achieved through coordinated database stress tests performed in a series of ATLAS reprocessing exercises at Tier-1s*
- Joint task force analysis of server performance under these stress tests conditions found that Conditions DB data access is limited by the IO
  - An unacceptable side-effect of the IO saturation is a degradation of the WLCG 3D Service that updates ATLAS Conditions DB data at all ten ATLAS Tier-1 sites using Oracle Streams technology
    - *Worldwide degradation of 3D Oracle Streams must be avoided*
- To avoid such bottlenecks ATLAS prototyped and tested novel database peak load avoidance approach for Grid computing
  - The foundation of the ATLAS approach is based upon the proven idea of “pilot” job submission on the Grid:
    - *instead of the actual query send a “pilot” query first*

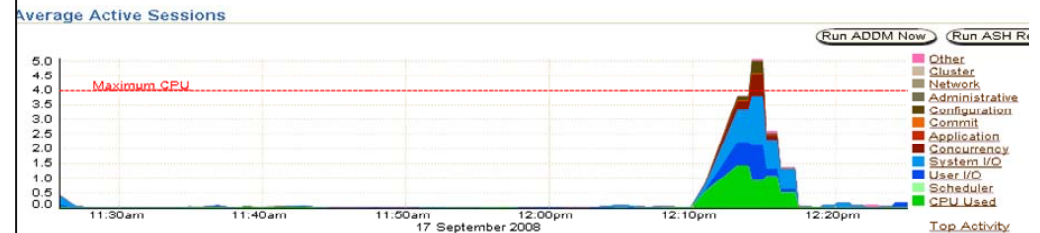


# ATLAS Pilot Query Approach



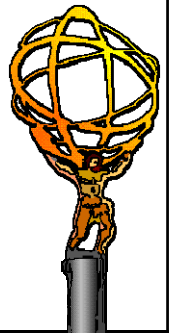
```
CREATE OR REPLACE FUNCTION "PILOT_QUERY" return
  number
authid definer
is
  nload number;
  threshold number;
  intervalsec number;
begin
  threshold:=2;
  intervalsec:=15;
  select round(count(*)/intervalsec,1) a
    into nload
  from v$active_session_history
  where sample_time>sysdate-intervalsec/(24*60*60)
  return nload;
end;
```

Query returns load similar to the OMS « Top Activity » plot number:



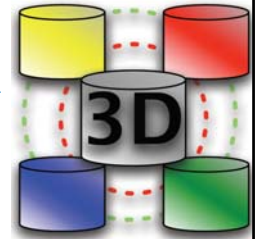
Would be integrated in production system for throttling jobs when load is high  
Threshold should be personalized according to T1 site  
Timings of function call and time interval average must be tuned to work

- ATLAS pilot query approach requires WLCG 3D Services approval and DBA privileges for installation at the Tier-1 site servers
  - The 3D Services alternative to the proposed ATLAS pilot query approach is to buy more disks to support peak database loads

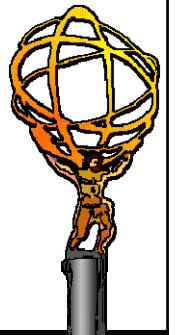


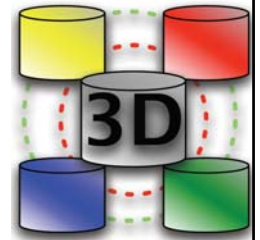
*Slides by F. Viegas*

## Performance of New Software for Conditions DB Access



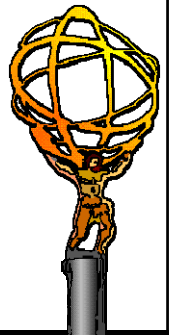
- New database access technologies were successfully validated during next round of ATLAS preparations for data reprocessing at Tier-1 sites
- On request of ATLAS reprocessing coordinator we tried reprocessing of ten cosmics and four FDR datasets at RAL, CERN, TRIUMF, BNL and NDGF with more than fifteen thousand jobs finished
- A largest task with more than six thousand jobs was processed at RAL without reaching any Oracle server limitations
- Due to Oracle monitoring deficiencies in the 3D database dashboard a smaller 344-jobs task had database access problems at NDGF
  - Task completed successfully after NDGF server restart and recommended 3D configuration settings were applied
- Recommendations based upon ATLAS experience during this rare 3D Services incident are presented at the next slide



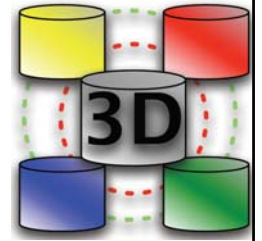


## Requested Improvements for 3D Services Procedures

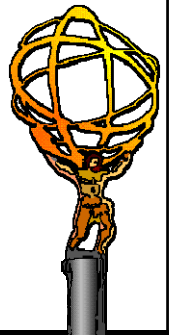
- WLCG 3D Services incidents reports should be prepared jointly with experiments' contacts for databases
  - Otherwise WLCG reports may contain incomplete information
    - *WLCG 3D Services incidents follow-up procedure should also involve experiments' contacts for databases*
- WLCG 3D Services dashboard for ATLAS experiment should not be all green when the 3D Oracle Streams are blocked
  - This prevents alerting Tier-1 DBAs or other ATLAS actions
    - *3D Services dashboard is part of database monitoring display for shifters in ATLAS world-wide computing operations control room*
- Procedures to alert Tier-1 DBAs should be documented and validated
  - Experiments' contacts for databases should be able to alert DBAs
    - *Some Tier-1 sites do not accept e-mails with ATLAS alerts*
- WLCG problem reporting procedures should be documented at
  - <https://twiki.cern.ch/twiki/bin/view/PSSGroup/ProblemReporting>
    - *Currently, this web page is "under construction"*



## Joint Task Force of ATLAS and CERN IT-DM

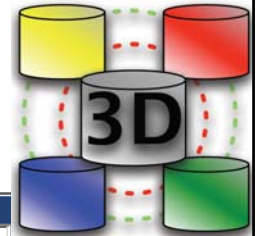


- Experience with database stress tests on the Grid found that test conditions are difficult to reproduce or predict
  - To conduct further stress tests in a controlled environment a joint task force of ATLAS and CERN IT-DM created a testbed for database performance studies at the dedicated CERN IT Oracle RAC server
- Using new ATLAS software the joint task force were unable to detect Oracle scalability limitations at the CERN IT server dedicated for these database performance studies
  - during ATLAS cosmics reprocessing task with 3,000-jobs L. Canali and F. Viegas found that the server load was low
    - *Work continues to systematically study Conditions DB access via COOL/CORAL (concurrent reconstruction jobs) without and with Oracle Streams*
- Our priority is confirmation of the COOL 2.6 performance improvements thanks to a removal of the unnecessary SELECT COUNT(\*) queries
  - Use the script for test jobs with realistic DCS access patterns
    - *Validation with actual reprocessing jobs can be done months later, which may be late for feedback cycle before 2009 LHC data taking*





# Joint Task Force: What We Do and Who we Are?



CHEP 2009



Search

21-27 March 2009

Prague

[Home](#) > [Call for Abstracts](#) > [View my abstracts](#) > [Abstract details](#)

modify

withdraw

## Advanced Technologies for Scalable ATLAS Conditions Database Access on the Grid

Abstract ID:81

### Content:

During massive data reprocessing operations an ATLAS Conditions Database application must support concurrent access from numerous ATLAS data processing jobs running on the Grid. By simulating realistic workflow, ATLAS database scalability tests provided feedback for Conditions DB software optimization and allowed precise determination of required distributed database resources. 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 average access rates. To validate database performance at peak loads, we tested database scalability at very high concurrent jobs rates. This has been achieved through coordinated database stress tests performed in series of ATLAS reprocessing exercises at the Tier-1 sites. The goal of database stress tests is to detect scalability limits of the hardware deployed at the Tier-1 sites, so that the server overload conditions can be safely avoided in a production environment. Our analysis of server performance under stress tests indicates that Conditions DB data access is limited by the disk I/O throughput. An unacceptable side-effect of the disk I/O saturation is a degradation of the WLCG 3D Services that update Conditions DB data at all ten ATLAS Tier-1 sites using the technology of Oracle Streams. To avoid such bottlenecks we prototyped and tested novel approach for database peak load avoidance in Grid computing. Our approach is based upon the proven idea of "pilot" job submission on the Grid: instead of the actual query ATLAS utility library sends to the database server a "pilot" query first.

### Summary:

We present detail results of our database access scalability studies as well as technologies developed to eliminate database access bottlenecks in a Grid computing environment.

**Primary authors:** VANIACHINE, Alexandre (Argonne)

DIMITROV, Gancho (LBNL)

NEVSKI, Pavel (BNL)

BASSET, Romain (CERN)

CANALI, Luca (CERN)

GIRONE, Maria (CERN)

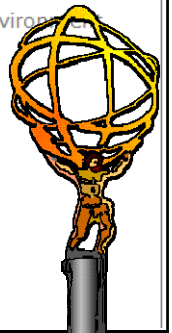
HAWKINGS, Richard (CERN)

VALASSI, Andrea (CERN)

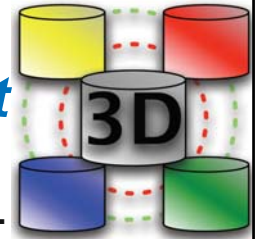
VIEGAS, Florbela (CERN)

WALKER, Rodney (LMU Munich)

WONG, Andrew (TRIUMF)



# Latest Results Confirmed that IO Load is More Important



- Presented at the ATLAS Software and Computing Workshop last week:

## DM Scalability data

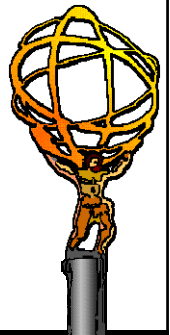
CERN IT Department

- Characterized DCS and reconstruction jobs
  - CPU load is about 10-20% of a CPU core per job
  - IO load more important
  - Depends on the cache status (i.e. second execution of a job is faster)
- Typical IO load
  - Measured with flushed cache
  - It's mostly RANDOM IO
  - Reconstruction: ~5k blocks in ~100 sec -> 40 MB -> 50 IOPS
  - DCS: ~40k blocks -> 300 MB -> ~400 IOPS

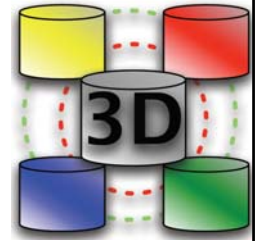
CERN IT Department  
CH-1211 Genève 23  
Switzerland  
[www.cern.ch/it](http://www.cern.ch/it)

Maria Girone

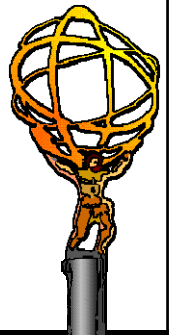
Physics Database Services 1



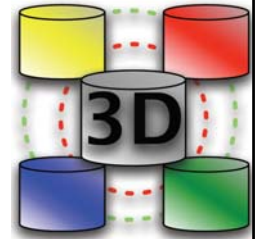
## Task Force Input to ATLAS DB Resource Requirements



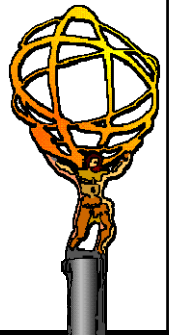
- Our joint task force on database performance studies provided key input for ATLAS database resource requirements for 2009
  - Assuming (pessimistic) 50 MB of random disk IO per job, the required IO rates at Tier-1 sites range from 200 to 1300 IOPS depending on the amount of CPU cores at the Tier-1 and its associated Tier-2 sites
    - *Sites will need IO upgrades only if they have to support peak loads*
  - For comparison ATLAS Tier-0 ATLR server IO rate was at the level of 2000 IOPS during 2008 operations
- In these *conservative* estimates Tier-2 Conditions DB access pattern assumed to be the same as for Tier-1 sites:
  - The seven-hour jobs spends five minutes accessing the Oracle RAC
- 2009 pledges for doubling of CPU deployment are not accounted for
- ATLAS targets January 2009 for a decision on Conditions DB access model (deployment of FroNTier)
  - New database access model may lower our DB resource requirements

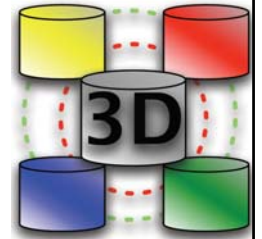


## ATLAS Schedule Relevant to 3D Services



- November-December:
  - Preparation for reprocessing
    - *Intermittent load on Oracle servers at Tier-1s during preparation*
  
- 20-21 January 2009:
  - ATLAS “2008 run post-mortem” workshop is targeted for a decision on Conditions DB access technology
    - *Should ATLAS deploy FroNTier for DB access at Tier-1/2/3?*
  
- January-March:
  - Reprocessing
    - *Continuous load on Oracle servers at Tier-1s during reprocessing*





## Conclusions

- During massive bulk data reprocessing operations at the Tier-1 sites a major ATLAS database application – the Conditions DB – is required to deliver robust concurrent access from numerous ATLAS data processing jobs running at the Tier-1 sites
- To provide input to upgrades WLCG 3D hardware resources, ATLAS database stress tests simulate realistic Conditions DB workflow
  - The goal of database stress tests is to overload the database cluster by launching many data processing jobs in parallel
    - *Database stress tests detect scalability limits of the hardware deployed at the Tier-1 sites, so that the server overload conditions can be safely avoided in a production environment*
- Database stress tests coordinated by the joint DB performance task force
  - Provide feedback to WLCG AA developers of the COOL/CORAL software used for access to Conditions DB data
  - Provide input for server reconfiguration requests to WLCG 3D Service
  - Provide input for database capacities requests to WLCG MB
    - *Our stress tests allowed precise determination of the actual ATLAS requirements for distributed database resources at the Tier-1 sites*

