# Experiments' updates: new DB applications and deployment plans

- **ATLAS: Florbela Viegas, Gancho Dimitrov, Alexander Vaniachine, Solveig Albrand, Pedro Salgado**
- **LHCb: Marco Clemencic**
- **CMS: Lee Lueking**

- CCRC'02 exercise

- CCRC Metrics and statistics collection

- Monitoring and ADC shift work integration

- MDT Calibration dataflow  and replication

- PVSS activities and future plans

- AMI replication status and plans

- DQ2 criticality and disaster recovery options

- LHCb activities

- CMS activities

# CCRC02  activity and plans

- Sasha's slides

Additionally to this metric (sessions per node) we want to  collect statistics for:

1. Resource usage per process:

   - Physical and Logical I/O used
   - CPU used
   - Cache hit ratio for the session

   setting the standard session audit is enough to extract them.

1. Overall machine/database behavior for I/O and CPU:

   - using the standard OS metrics collected by Oracle this can be easily obtained, might need some saving to tables periodically, as these may be overwritten. So the settings on these statistics gathering must be done by all T1s.

2. Statistics for SQL activity from the sessions:

   - For getting statistics of selects, nature of selects and other more the only way is to set the whole database on trace, and process the resulting files. This adds a bit of CPU overhead to the database, but it might be worth setting for a limited time, as it gives a very good picture of "typical" usage per process.

Some of these statistics were already collected ad-hoc in the last tests at BNL and TRIUMF.

# Streams monitor – graphs

- **The ADC shift model:**

- Operations expert on-call has to determine source of problem and escalation

- There should be a simple way to monitor the « health » of the conditions databases at CERN and T1s

- The diagnosys checklist, and visible indicators of health, should include:

  - Number of processes resource limit: Are the jobs getting errors because they can't connect?)

  - Load on the database (Are the jobs too slow? Are they accumulating in the database active sessions)

  - Activity on database ( Is there I/O contention? Are there locks causing hung jobs?)

  - Is the replication within the expected delay ? (Are jobs not getting new enough data?)

- Should alarms be setup for these events, in a homogeneous fashion for the whole T0/T1 operations? Using nagios, lemon, cron?

# SLS for atlas_coolprod service



Browser window showing: https://lemonweb.cern.ch/sls/service.php?id=phydb_atlas_coolprod

**SLS Service Level Status overview**

Home | Search | KPIs | Tags | Admin | Documentation

## Oracle RAC for atlas_coolprod services
24 Apr 2008 Thu

### Service information
full name: **Oracle RAC for atlas_coolprod services**
short name: atlas_coolprod
group: IT/PSS
site: CERN

email: **phydb.support@cern.ch**
web site: http://cern.ch/phydb

### Service availability (more)
availability:
percentage: 100%
availability info: Service fully available
status: **available**

last update: 09:53:44, 24 Apr 2008
(4 minutes ago)
expires after: 60 minutes

rss feed with status changes

### Part of (subservice of):
ATLAS RAC database

### Subservices
none / not declared

### Clusters, subclusters and nodes
none / not declared

### Depends on
none / not declared

### Depended on by
none / not declared

- A new replication model: MDT calibration
- Dataflow:

- Full MDT calibration data replication cycle must be complete in 24 hours.

- For T0-T1 replication, the already setup mechanism is fine, MDT calibration is one more schema to replicate.

- From Rome, Michigan and Munich to CERN, replication is monitored and maintained by site DBAs. A Service Level Agreement should be put in place.

- The jobs at these sites should get their data from the nearest T1 (INFN, BNL and Gridka.

- The other replication components and processes have to be monitored as part of the MDT calibration data replication service, so some tool has to be put in place for the ADC shifters monitoring – Gridmap?

- **Test replica from INTR to INT8R has been setup**



CAPTURING 2782.49 LCRs/s
PROPAGATING 2778.32 LCRs/s
APPLYING 2751.62 LCRs/s

INTR.CERN.CH(CERN) → INT8R.CERN.CH(CERN)

- **Achieved improvement with almost a factor of 2 in the PVSS replication throughput since the last tests done in February**

  **From 2000-2200 LCR/sec before to 3500-3800 LCRs/sec now (maximum steady rate without accumulation of latency)**

- **Special thanks to Luca Canali and Eva Dafonte.**

- Instead of the standard heap organized PVSS 'EVENTHISTORY_xxx' tables + PK, we create them as Index-Organized ones (IOT), thus reducing the I/Os and the disk space usage.

- Two hidden parameters, found from Luca, determine when a flow control to happen (the CAPTURE default is 15000 unbrowsed mesagges is the queue). Setting it to higher values caused less frequent flow control to happen.

  _buffered_publisher_flow_control_threshold

  and

  _capture_publisher_flow_control_threshold

**LCRs Applied**
**STRMADMIN_APPLY@INT8R.CERN.CH**
Generated on Friday 18th of April 2008 01:06:49 PM
Time range 10:30:26 18-04-2008 <-> 13:06:13 18-04-2008

A short period, when the clients were inserting ~ 4000-4200 rows/sec

**Apply Latency**
**STRMADMIN_APPLY@INT8R.CERN.CH**
Generated on Friday 18th of April 2008 01:07:53 PM
Time range 10:30:46 18-04-2008 <-> 13:07:33 18-04-2008

# Is the achieved throughput enough?

- The results are promising and we hope that the achieved throughput will be enough for dealing with the 'normal' PVSS insertion plus cases of burst of inserts on 'start of run' and other events.

- So far the overall INSERT activity hasn't been more than 500 rows/sec


PVSS Overall Activity

# Oracle streams replication of the COOL data



CAPTURING 238.09 LCRs/s
PROPAGATING 212.02 LCRs/s
APPLYING 130.7 LCRs/s
ASGC(TAIWAN)

CAPTURING 238.09 LCRs/s
PROPAGATING 237.45 LCRs/s
APPLYING 131.74 LCRs/s
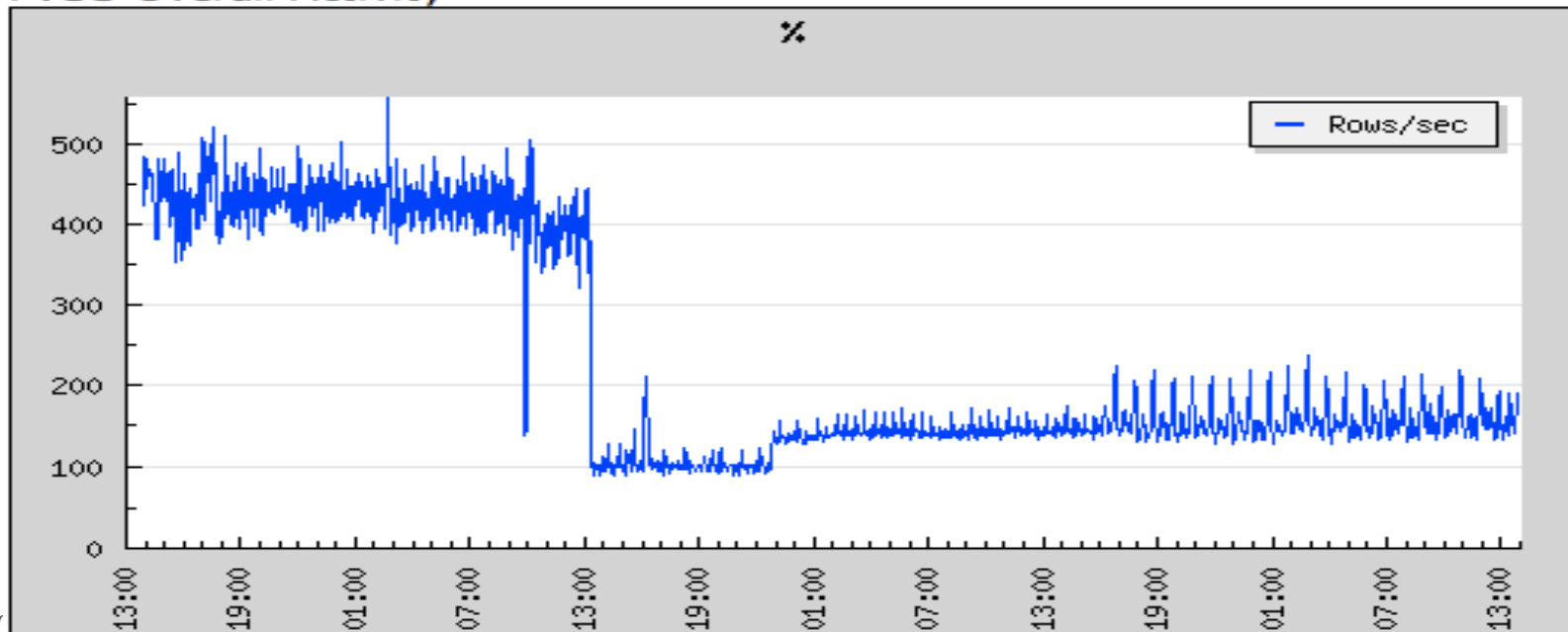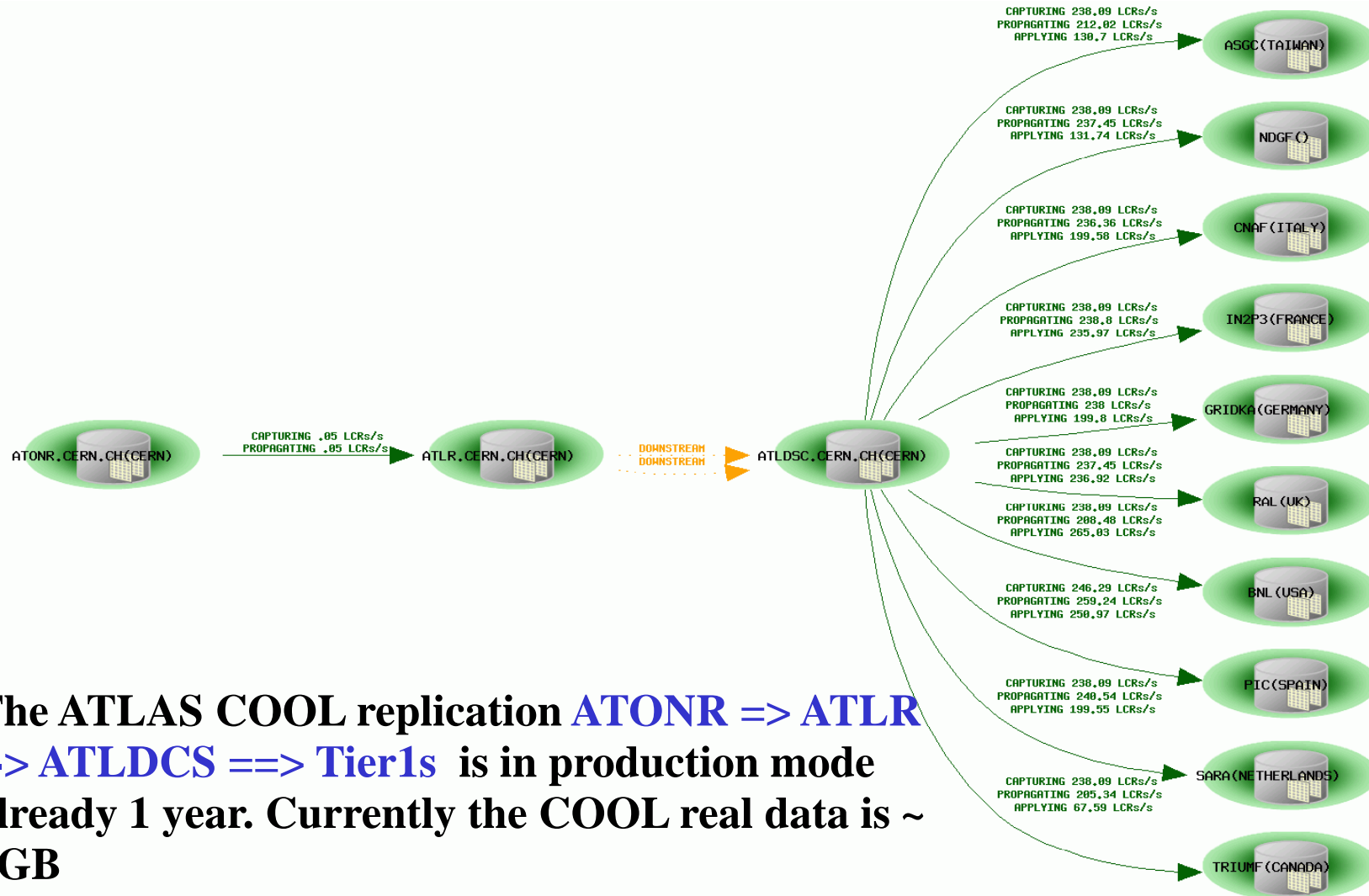NDGF()

CAPTURING 238.09 LCRs/s
PROPAGATING 236.36 LCRs/s
APPLYING 199.58 LCRs/s
CNAF(ITALY)

CAPTURING 238.09 LCRs/s
PROPAGATING 238.8 LCRs/s
APPLYING 235.97 LCRs/s
IN2P3(FRANCE)

CAPTURING 238.09 LCRs/s
PROPAGATING 238 LCRs/s
APPLYING 199.8 LCRs/s
GRIDKA(GERMANY)

CAPTURING 238.09 LCRs/s
PROPAGATING 237.45 LCRs/s
APPLYING 236.92 LCRs/s

CAPTURING 238.09 LCRs/s
PROPAGATING 208.48 LCRs/s
APPLYING 265.03 LCRs/s
RAL(UK)

CAPTURING 246.29 LCRs/s
PROPAGATING 259.24 LCRs/s
APPLYING 250.97 LCRs/s
BNL(USA)

CAPTURING 238.09 LCRs/s
PROPAGATING 240.54 LCRs/s
APPLYING 199.55 LCRs/s
PIC(SPAIN)

CAPTURING 238.09 LCRs/s
PROPAGATING 205.34 LCRs/s
APPLYING 67.59 LCRs/s
SARA(NETHERLANDS)

TRIUMF(CANADA)

ATONR.CERN.CH(CERN)

CAPTURING .05 LCRs/s
PROPAGATING .05 LCRs/s

ATLR.CERN.CH(CERN)

DOWNSTREAM
DOWNSTREAM

ATLDSC.CERN.CH(CERN)

**The ATLAS COOL replication ATONR => ATLR --> ATLDCS ==> Tier1s is in production mode already 1 year. Currently the COOL real data is ~ 4GB**

# The detector conditions data and its replication

- The detector conditions data will be the biggest fraction of data replicated to the Tier1s. This data resides into the ATLAS_COOLOFL_DSC account on the ATLAS 'offline' DB 'ATLR'.

- For the purpose of not flooding the replication, the schema is temporary taken out of the flow.

- After processing the PVSS data gotten in the last few months from all 11 sub-detectors, the above schema will be instantiated on all destination databases using transportable tablespace method ( expected to be 10s of GBs )

- **The Oracle replica is active for**
  - 16 COOLONL_xxx schemas
  - 16 COOLOFL_xxx schemas
  - CONF schemas (TRIGGER, PIXEL, TGC …)
  - OKS_TDAQ, ATLAS_RUN_NUMBER
  - ATLOG (atlas_logbook), COCA(Collection and Cache), MDA (Monitoring Data Archiving)
- **To be added**
  - PVSS archive (so far, 11 schemas )
  - PVSS CONF (configuration data )
  - ATLAS_MDT_DCS

- AMI Status – Solveig's slides

- DQ2 – slides by Pedro

- Darios's quote:

- " Data taking may be in trouble after a few hours if no new dataset registration can proceed and therefore no data is moved and no jobs can be submitted (including detector monitoring and calibration), […] consider 4 hours as a kind of absolute maximum."

- slides by Marco Clemencic

- **Slides by Lee Lueking**