

Published on *CERN openlab* (<http://openlab-archive-phases-iv-v.web.cern.ch>)

[Home](#) > Database - Activities Covering 2012

Database - Activities Covering 2012 ^[1]

In January 2012, the start of the fourth phase of CERN openlab saw an increase in the number of Oracle sponsored fellows. In addition to expanding existing openlab studies of database technology, database monitoring and replication and virtualisation, two new work areas were added to the CERN openlab Database Competence Centre (DCC) programme of work in openlab IV: an investigation of the feasibility of using a relational database for physics data analysis and a study of how CERN could take advantage of Oracle analytics capabilities to extract knowledge from and gain further insight into the large amounts of data collected on the LHC operation and other CERN systems.

Discovering the Higgs boson in an Oracle database

Maaïke Limper, an experienced physicist from the ATLAS experiment, joined the DCC team at the beginning of 2012 to investigate the feasibility of using a relational database for physics data analysis, the first of the two new topics. She quickly made progress porting a C++ based analysis using CERN's ROOT framework (<http://root.cern.ch> ^[2]) to a mixture of SQL (Structured Query Language) and PL/SQL (Procedural Language / Structured Query Language) operating on data stored in Oracle. The analysis chosen was a search for events where a standard model Higgs boson is produced in association with a Z0 boson, a particularly interesting Higgs decay mode.

The first step, of course, consisted in making sure that the analysis in an Oracle database produces results that are absolutely identical to those produced with the ROOT based analysis. Maaïke successfully demonstrated that the results are actually identical, at least for a simple analysis, as can be seen from the plots on page 25 and as highlighted and presented during the DCC team visit to Oracle's headquarters in March.

Further time and effort was then needed to understand how to efficiently invoke C++ algorithms that were too complex for rewriting in PL/SQL, and to optimise the internal database processing. By the end of the year, although the database version of the analysis was slower when running on a Monte Carlo generated set of 30,000 signal events, it was

more efficient when run against a set of some 1,650,000 background events. The reason for this is that the background sample has fewer interesting events and these are selected very efficiently by the database.

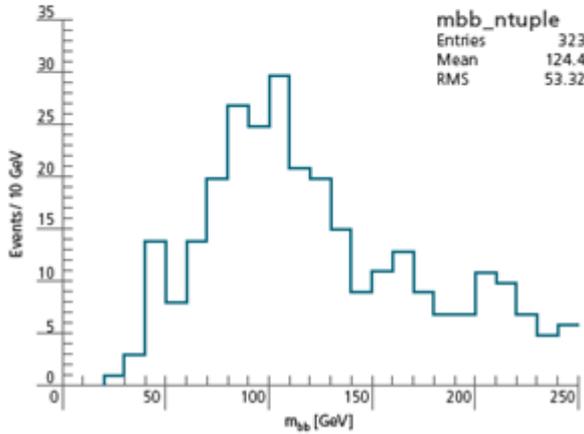
In principle, this is good news for the database version since a 'real-world' analysis of the data would be over a much larger sample of events with an even smaller percentage of events tagged as containing a $Z^0 \rightarrow ll$ decay. However, in the real world, the ROOT based analyses are performed in parallel around the world at sites participating in the Worldwide LHC Computing Grid (WLCG), with many thousands of analyses running at any one time. As foreseen when preparing the programme of work, prior to the start of the fourth phase of CERN openlab, access to Exadata systems is now required to see how the database would support multiple different analyses in parallel. The DCC team looks forward to being able to make these tests in 2013.

Extracting knowledge from operations data of LHC with Oracle Analytics

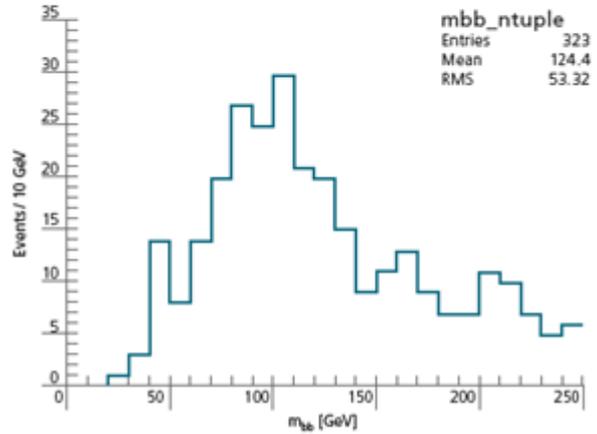
Whilst the physics output of the LHC experiments are stored outside the relational databases, Oracle databases nevertheless play a key role in the operation of the LHC accelerator. Logging data from the thousands of different LHC components has grown from 60 terabytes (TB) and 2.2 trillion records at the time of the last openlab annual report to 165 TB and 4.8 trillion records today. Greg Doherty, the DCC Development Executive Sponsor at Oracle has long commented that active analysis and mining of this huge dataset might help improve the operation of the accelerator by identifying patterns of events that indicate potential problems.

Plots showing results from an Oracle-based analysis (left) identical to those from a ROOT based analysis (right)

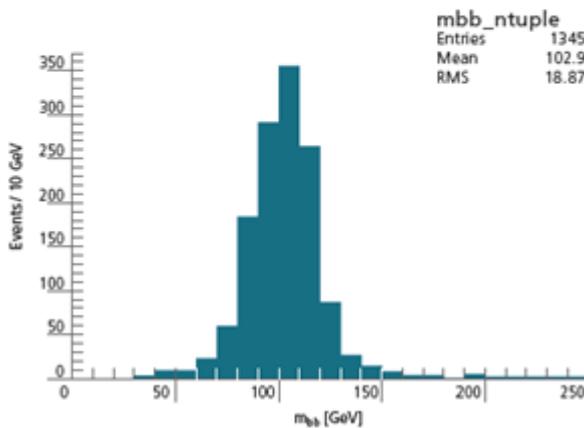
Invariant di-bjet mass



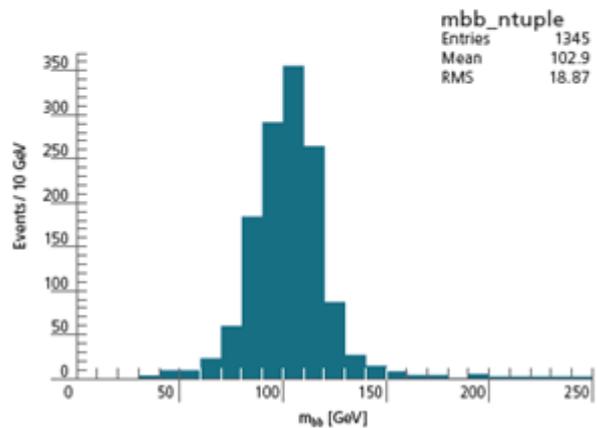
Invariant di-bjet mass



Invariant di-bjet mass



Invariant di-bjet mass



[3]

Following up on these ideas was the reason for the creation of the second new work area in the DCC for openlab IV. Manuel Martin Marquez joined the team over the summer and was responsible for organising a successful Data Analytics workshop at CERN in November. The workshop brought together speakers from CERN, Oracle and another openlab partner, Siemens, to address specific use cases at CERN. Presentations on Oracle's Real-Time Decisions (RTD) and Oracle Endeca products attracted much interest and a regular Forum has been established to continue work in this area.

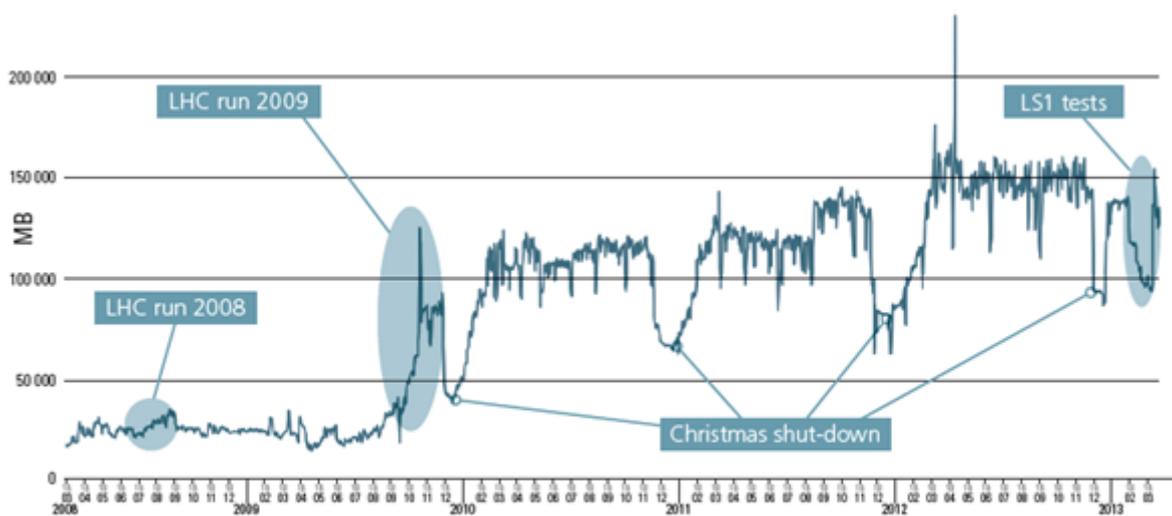
Database technology, replication and virtualisation

The well-established lines of the collaboration in CERN openlab for core database technology, database replication, and in the monitoring and virtualisation areas, have of course continued in openlab IV. Two of these areas, database technology and virtualisation, converge with the development of multi-tenancy database technology for the

next generation Oracle database. This potentially offers a very attractive and efficient option to replace databases in virtual machines and the database team at CERN, including Andrei Dumitru from CERN openlab, extensively tested this feature in the framework of the Oracle Database 12c Beta Testing programme. The next step will be to exploit pluggable databases, and other 12c features, in a production environment.

Deploying technology in production which had been tested first in the CERN openlab framework is exactly what has been done in the area of database replication. Oracle Active Data Guard was evaluated in openlab during the early stages of the 11g database release. When CERN databases migrated to 11g in early 2012, Active Data Guard was deployed to replace the Streams based replication of the databases supporting the online operations of the LHC experiments. Active Data Guard replicas have also been added to offload some primary databases, notably in the accelerator operations area where the secondary databases can be used to support the data analytics work mentioned earlier without any risk of affecting production accelerator operations.

CERN Accelerator Logging Service daily storage



[4]

The logging service stores data using Oracle RAC databases, of close to one million pre-defined signals coming from heterogeneous sources, and it provides access to logged data for more than 700 registered individuals, more than 100 registered custom applications from around CERN, and even offsite access for purposes such as the CNGS experiments in Gran Sasso Italy.

In the replication area, the DCC welcomed a new team member, Lorena Lobato Pardavila, who has been working on various GoldenGate features, notably an improved monitoring plugin for Enterprise Manager and technology that has the potential to address issues identified in earlier openlab work. The DCC team hopes to join a forthcoming beta programme and to invest more effort in this area during 2013.

The final new team member, Ignacio Coterillo Coz, has been working in the virtualisation area, continuing the DCC work with the Oracle VM team, as well as the joint openlab project with Intel to evaluate the impact of Single-Root I/O Virtualisation (SR-IOV) for virtualised database workloads. Here, the ability of Oracle's Real Application Testing technology to record and replay workloads means that the DCC team can measure directly the impact of SR-IOV on the performance of CERN's applications rather than having to make assumptions with we extrapolations from test loads. These tests show that use of SR-IOV can reduce the time taken to execute I/O intensive queries significantly (by up to a factor of three

in some cases) which is very positive as the team looks to further exploit virtualisation to improve service flexibility and overall reliability.

In summary, it has been another fruitful year for the DCC. Oracle's extended support in openlab IV has enabled the team to launch interesting investigations for physics data analysis and data mining, as well as to continue evaluating exciting new developments, giving direct feedback to the Oracle developers. Although the work in the latter area cannot be disclosed, the solutions tested by the DCC team are promising and very much in line with CERN's ideas for future service needs. Working on these in the openlab context has been again extremely beneficial to both CERN and Oracle and more will come in 2013.

- [Visit Us](#)
- [RSS Feeds](#)

DISCLAIMER: This Web page contains pointers to material related to the management of CERN openlab in the Information Technology Department at the European Organization for Nuclear Research (CERN). Their use and distribution are regulated by the [CERN copyright notice](#).



Source URL: <http://openlab-archive-phases-iv-v.web.cern.ch/database-activities-covering-2012>

Links

[1] <http://openlab-archive-phases-iv-v.web.cern.ch/database-activities-covering-2012>

[2] <http://root.cern.ch/>

[3] <http://openlab-archive-phases-iv-v.web.cern.ch/sites/openlab-archive-phases-iv-v.web.cern.ch/files/Plots%20Showing%20results.png>

[4] <http://openlab-archive-phases-iv-v.web.cern.ch/sites/openlab-archive-phases-iv-v.web.cern.ch/files/Cern%20Accelerator%20Loggin%20Service%20daily%20Storage.png>