



# Data preservation at LHC

- November 29, 2013 -

**S. Amerio**

(University of Padova, INFN)

Data preservation activities are becoming more relevant in HEP experiments.  
Long term data preservation (and open access) matters for

- **social reasons:** our experiments are paid by the community, data should be preserved and made available to general public.
- **scientific reasons:**
  - old data can always be re-analyzed in search of new signals or to improve measurements.
  - scientific results should be reproducible.

All LHC experiments are now actively working on data preservation issues.

- Dedicated task forces
- Dedicated sessions/talks at collaboration meetings/workshops
- DP activities are considered as service tasks to the experiment

### ***Common areas of work:***

- Open access and outreach
- Bit and software preservation
- Validation systems
- Documentation/Analysis preservation

CMS open access policy official since March 2012, LHCb since March 2013  
 Atlas/Alice: drafts under discussion within their collaborations.

CMS/LHCb policy (based on the DPHEP levels of data preservation):

<b>Level 1 (published data)</b>	All scientific results are public. Additional data associated with the results will also be made available (e.g. Histograms data); additional info archived in Inspire or HEPdata.
<b>Level 2 (samples for educational purposes)</b>	Data samples in simplified format for event displays and masterclass exercises.
<b>Level 3 (reconstructed data)</b>	CMS will release <b>50% of their data 3 years after data taking</b> LHCb will release <b>50% 5 years after data is taken, 100% after 10 years.</b>
<b>Level 4 (raw data)</b>	Due to the complexity of the raw data processing stage, the extensive computing resources required and enormous access to tape resources, direct access to raw data is not permitted to individuals within the collaboration. Raw data processing is performed centrally. Due to this, CMS/LHCb are currently <b>not planning to allow open access to raw data</b>

CMS open access policy official since March 2012, LHCb since March 2013  
 Atlas/CMS: drafts under discussion within their collaborations.

CMS/LHCb policy (based on the DPHEP levels of data preservation):

<b>Level 1 (published data)</b>	All scientific results are public. Additional data associated with the results will also be made available (e.g. Histograms data); additional info archived in Inspire or HEPdata.
<b>Level 2 (samples for educational purposes)</b>	Data samples in simplified format for event displays and masterclass exercises.
<b>Level 3 (reconstructed data)</b>	CMS will release <b>50% of their data 3 years after data taking</b> LHCb will release <b>50% 5 years after data is taken, 100% after 10 years.</b>
<b>Level 4 (raw data)</b>	Due to the complexity of the raw data processing stage, the extensive computing resources required and enormous access to tape resources, direct access to raw data is not permitted to individuals within the collaboration. It means we need to provide centrally. Due to this, CMS/LHCb are providing external users with <b>open access to raw data</b>

- Access to data
- Software to analyse data
- Documentation

- Preparing for a public release of part of 2010 collision and simulated data in AOD format, appropriate for physics analysis, and accompanied by stable, open source software needed for a number of example analysis and suitable documentation.
- A virtual machine (VM) image has been prepared, in a format usable by the freely-available VirtualBox application.
- For the access to the data, the initial work flow is kept as close to the standard one as possible, which uses xroot. An xrootd server has been commissioned with anonymous read-only access, further limited by firewall to include only those sites involved in the testing phase.
- Preparing high-school level classroom applications chosen as a pilot use-case[6] - tools to process the data and documentation into an appropriate format, studies on:
  - a common definition of data contents for most HEP experiments
  - a common HEP ontology with Linked Data methods.

From Kati Lassili-Perini poster at CHEP2013

Open data project funded by Finnish government in collaboration with IT Center for Science. The project is part of a larger plan of Finnish Ministry of Education to put scientific data in use, and a pilot project for the platform of the open research material at CSC.  
<https://twiki.cern.ch/twiki/bin/view/HIPCMSExperiment/CMSOpenDataProject>

# Other open access & outreach projects

All LHC experiments have released samples of data in simplified format for educational purposes, to be used in event displays and masterclass exercises.

## Public Data

- Data Samples
- Analysis Tools & Formats
- CMS Masterclass
- Useful links

### CMS Public Data

The CMS experiment at the LHC has released a portion of its data to the public for use in education and outreach. Explore this page to find out more about the data and how to analyse it yourself.

Try the online event display below.

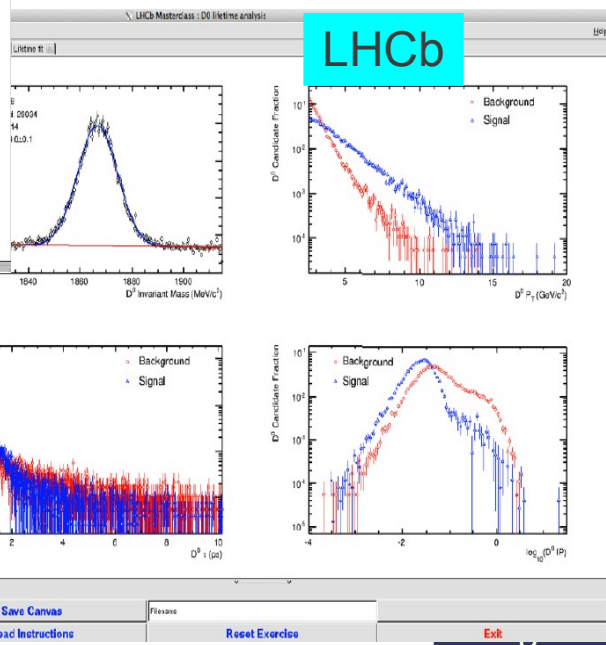
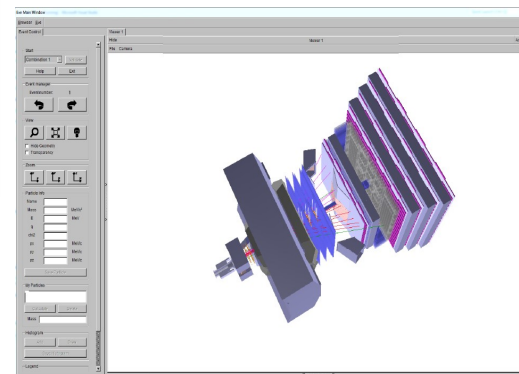
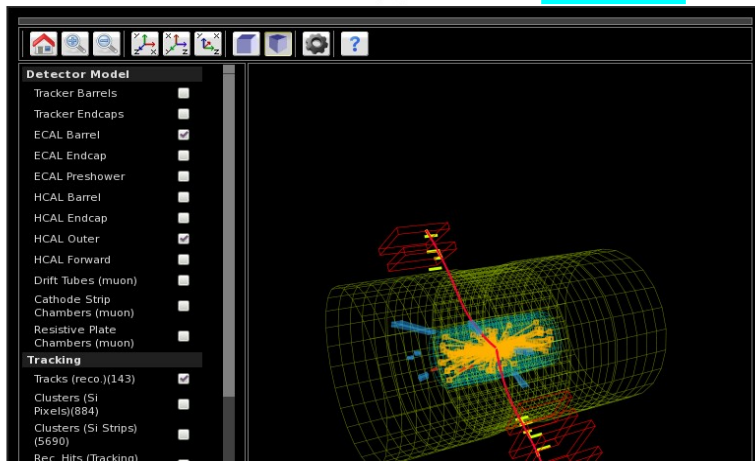
Use the Mouse to rotate.

Ctrl+Mouse or Ctrl +  $\leftarrow$   $\rightarrow$   $\uparrow$   $\downarrow$  to pan x/y.

Shift+Mouse or Shift +  $\leftarrow$   $\rightarrow$  to zoom.

Follow the links at left to access the full version of the display and thousands of ev

CMS



Alice

The calculator pops DP @ LHC

## CMS/LHCb

- All raw data preserved
- Define legacy datasets / software releases: all Run I data reprocessed into a legacy dataset with a single software version
- Use virtualization techniques to run legacy software
- Backup policies
- Data integrity checks
  
- Analysis level software preserved in dedicated repositories

Alice/Atlas developinig similar plans

An efficient validation system is a fundamental ingredient of any long term data preservation plan.

*Archived data and software need to be regularly checked against data loss and/or corruption, new operating systems and hardware.*

LHC experiments are building long term validation systems *on top of their current validation tools.*

**CMS:** a set of reference plots being defined covering physics results: from individual objects (muons, electrons,...) to high level physics signatures involving basic analysis selections.

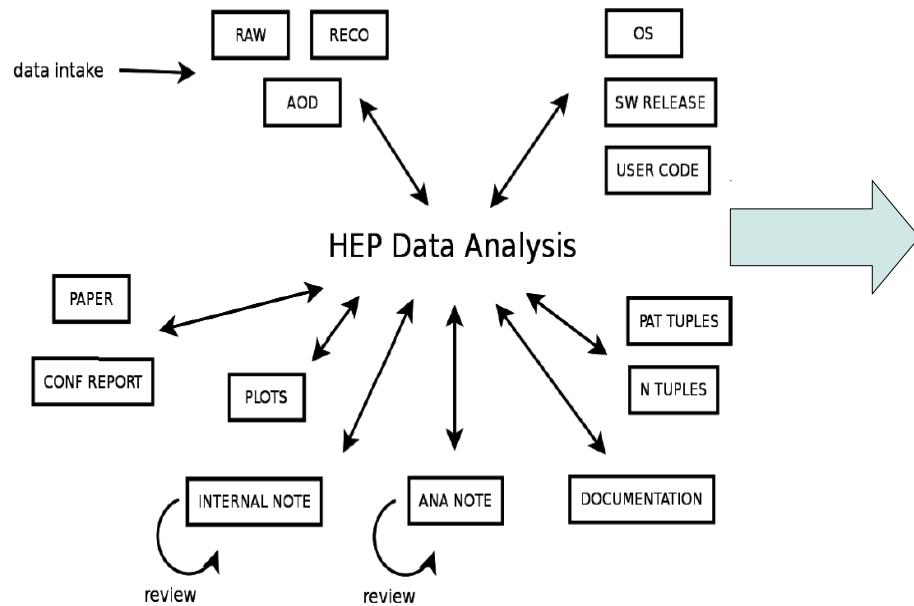
**LHCb:** for each step of data processing define long term future validation “analysis” in the current LHCb Performance & Regression tool

The screenshot shows the LHCb Performance & Regression tool interface. The top header includes the LHCb Computing logo and the text "LHCb PERFORMANCE & REGRESSION". Below the header is a navigation menu with "Home", "Jobs", "Job Descriptions", "Analyse", and "help". Under "Analyse", there are sub-menus for "BRUNEL", "DAVINCI", "GAUSS", and "MOORE". The "GAUSS" menu is currently selected. The main content area displays "Successful handlers for this application:" followed by a list of handlers: "TimingHandler" and "gaussGenerator". Below this, there is a section titled "Choose type of analysis:" with six buttons: "BASIC", "HISTOGRAMS", "OVERVIEW", "TIMING", "TIMING\_COMPARISON", and "TREND".



## CMS/LHCb/Cern-IT

The Invenio team at CERN, with input from CMS and LHCb, will setup a prototype of a tool to make recording and documenting of the workflow and intermediate data easy.



The screenshot shows the CERN Document Server interface. At the top, it says 'CERN Accelerating science' with 'Sign in' and 'Directory' links. Below that is the 'CERN Document Server' header with the tagline 'Access articles, reports and multimedia content in HEP'. A search bar is visible with the 'INSPIRE' logo. A welcome message for INSPIRE is displayed. Below the search bar is a Zenodo banner with the text 'zenodo Research. Shared.' and navigation links for 'Search', 'Collections', 'Upload', and 'Get started'. A search results section shows 'Search 389 records for:' and 'Showing records 1 to 10 out of 389 results.' The first result is a conference paper from 29 April 2013 titled 'Considering formal assessment in learning analytics within a PLE: the HOU2LEARN case' by Koufocheri, Eleni; Xenos, Michalis. The authors' names and counts are listed on the left.

## Atlas

RECAST: framework developed to extend impact of existing analysis; it allows an existing analysis to be reinterpreted under an alternate model hypothesis --> complete information from original analysis needs to be preserved.

LHC experiments are actively working on data preservation.

They are working within **DPHEP**; collaboration is encouraged but not organized so far.

Same problems to face --> a lot of room for collaboration/common projects.  
At the moment not discussed yet within WLCG.

The final goal should be the development of a common framework for long term preservation of HEP data.

- Backup -

Study group started in 2009

Blueprint released in May 2012 (available at [dphep.org](http://dphep.org))

Priorities:

- *Secure data in all experiments*
- *Consolidate the on-going international collaboration*
- *Promote common multi-experiment projects and/within interdisciplinary cooperation*

DPHEP is now moving from a study group to an organization (collaboration)

- **Full time project manager** (J.Shiers) to coordinate the activities
- Definition of a **collaboration agreement** to be signed by laboratories and experiments → already signed by Cern.
- Regular meetings with representatives from all lab/experiments to discuss ongoing DP projects.

DPHEP vision:

By 2020 all archived data easily findable, fully usable by designated communities with clear (open) access policies.

Best practices, tools and services well run-in, fully documented and sustainable; built in common with other disciplines, based on standards.