

Erweitertes Forschungsdatenmanagement: Langzeitar Archivierung von Forschungsdaten in der Hochenergiephysik

Patricia Herterich

@pherterich

Humboldt-Universität zu Berlin & CERN Scientific Information Service

6. Bibliothekskongress Leipzig 2016

14.03.2016

Vision: Open & Reproducible Research

- Open Research or Open Science aims at making research as transparent and reproducible as possible
 - Open Access to publications
 - Access to data produced or used
 - Access to code and methodology applied
 - Use of Open Source software instead of proprietary solutions whenever possible

Policy framework - CERN

The CERN Convention (1953) contains what is effectively an early Open Research manifesto:

“... the results of its experimental and theoretical work shall be published or otherwise made generally available”

Policy framework – LHC Collaborations

Approved CB 20th June 2014

ATLAS Data Access Policy

May 21st 2014

Introduction

ATLAS has fully supported the principle of open access in its publication policy. This document outlines the policy of ATLAS as regards open access to data at different levels as described in the DPHEP [1] model. The main objective is to make the data available in a usable way to people external to the ATLAS collaboration.

The ATLAS policy for data preservation is described in a separate document. The collaboration's need to preserve data for its own use shares some requirements with making them open access. To support open access to data additional resources will be required to develop and support the tools to make the data available.

Policies for LHCb External Data Access Policy

Open access to levels of increa

ur
is

LHCb Public Note

Issue: 1
Revision: 1

Reference: LHCb-PUB-2013-003
Created: 22nd April 2013
Last modified: 22nd April 2013

ALICE data preservation strategy

Sunday, October 6, 2013

The data harvested by the ALICE Experiment up to now and to be harvested in the future constitute the return of investment in human and financial resources by the international community. These data embed unique scientific information for the in depth understanding of the profound nature and origin of matter. Because of their uniqueness, long term preservation must be an essential objective of the data processing framework and will lay the foundations of the ALICE Collaboration legacy to the scientific community as well as to the general public. These considerations call for a detailed assessment of the ALICE data preservation strategy and policy. Documentation, long term preservation at various levels of abstraction, data access and analysis policy and software availability constitute the key elements of such a data preservation strategy allowing future collaborators, the wider scientific community and the general public to analyze data for educational purpose and for eventual reassessment of the published results. The present document describes the basic principles that will guide the redaction addressed by the ALICE data preservation policy.

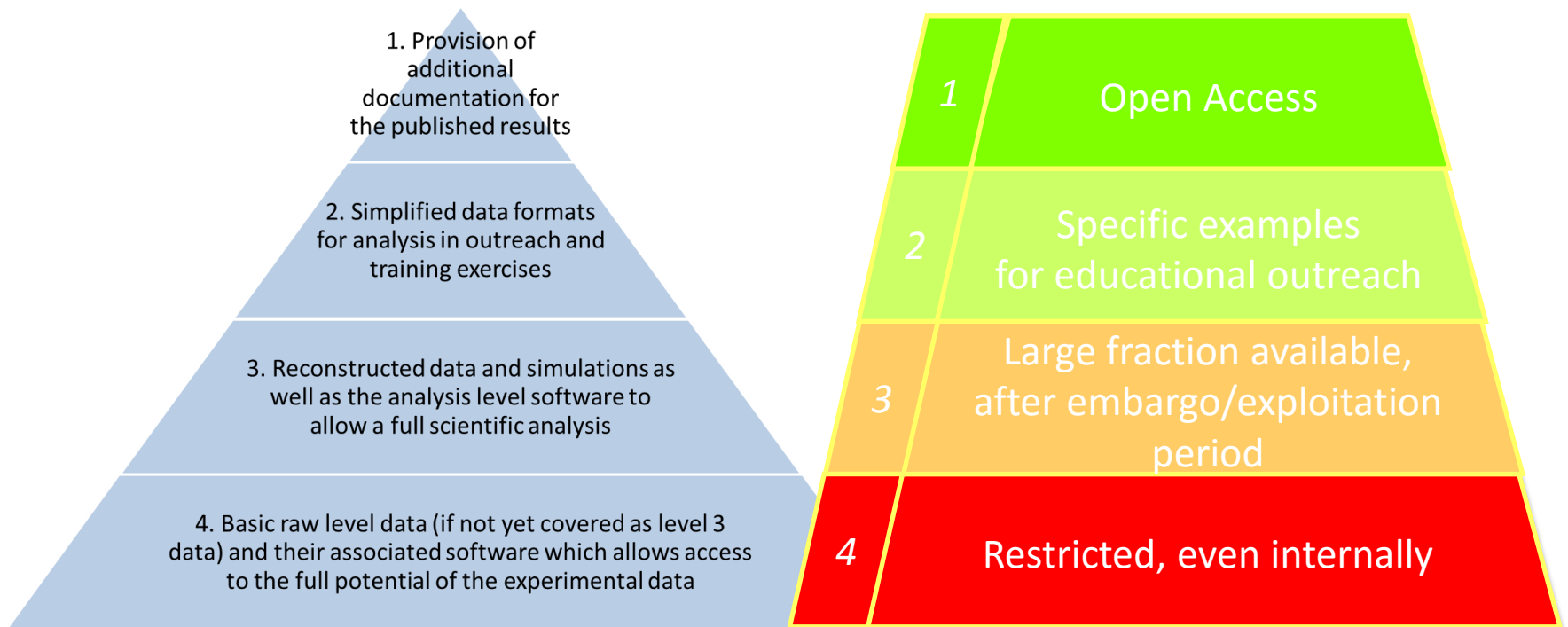
CMS data preservation, re-use and open access policy

CMS data are unique and are the result of vast and long-term moral, human and financial investment by the international community. There is unique scientific opportunity in re-using these data, at different level of abstraction and at different points in time¹. This opportunity calls for our collective responsibility, and poses unprecedented challenges as no data sample of this complexity and value has ever been preserved or made available for later re-use.

The CMS collaboration is committed to preserve its data, at different levels of complexity, and to allow their re-use by a wide community including: collaboration members long after the data are taken, experimental and theoretical HEP scientists who were not members of the collaboration, educational and outreach initiatives, and citizen scientists in the general public.

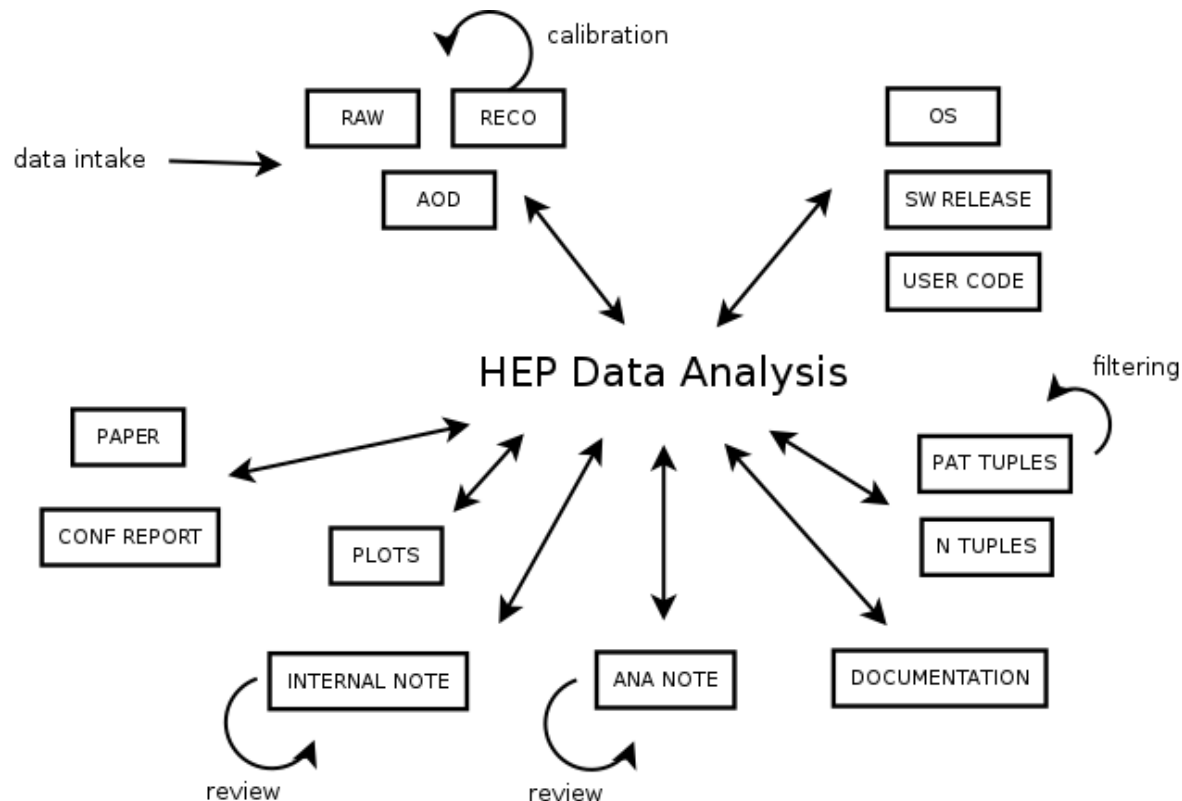
CMS upholds the principle that open access to the data will, in the long term, allow the maximum realization of their scientific potential. To that extent, CMS will provide open access to its data after a suitable but relatively short embargo period, allowing CMS collaborators to fully exploit their scientific potential.

Policy framework – LHC Collaborations



More than research data

“reconstructed data and simulations **as well as the analysis level software to allow a full scientific analysis**”



CERN Open Data Portal

opendata
CERN

ABOUT SEARCH EDUCATION RESEARCH

Education

Visualise events, check reconstructed data, run tools or build your own!

Start learning

Research

Get the genuine working environments, virtual machines and datasets to start your research

Start analysing

CERN Open Data Portal

- Launched Nov 2014: <http://opendata.cern.ch/>
- Public access point to data (incl. software and documentation) produced at CERN
- Upon release access to 27 TB of CMS data + educational data from all 4 LHC experiments
- Github repo: <https://github.com/cernopendata/>
- Based on Digital Library Software Invenio 2.0
- “preservation in the open”

Data

The screenshot displays the OpenData CERN website interface. At the top, the 'opendata CERN' logo is on the left, and navigation links for 'ABOUT', 'SEARCH', 'EDUCATION', and 'RESEARCH' are on the right. A search bar is located below the navigation. The main content area shows a breadcrumb trail: 'Home > CMS > CMS Primary Datasets'. The dataset title is 'Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21ReReco-v1/AOD) 2014'. Below the title, the URL '/Mu/Run2010B-Apr21ReReco-v1/AOD' and 'CMS collaboration' are listed. A 'Cite as' section provides the citation: 'CMS collaboration (2014). Mu primary dataset in AOD format from RunB of 2010. doi:10.7483/OPENDATA.CMS.B8MR.C4A2'. A filter bar includes 'Collection' (selected as 'CMS Primary Datasets'), 'Collision Energy' (7TeV), and 'Accelerator'. The left sidebar contains sections for 'Description' (Mu primary dataset in AOD format from RunB of 2010) and 'Characteristics' (Dataset: 32376291 events, 2979 files, 3.2 TB in total). The right sidebar contains a detailed information panel with sections: 'How can you use these data?' (with links to CMS data quality monitoring software, CMS Virtual Machine instructions, and getting started with CMS open data), 'Issues & Limitations' (linking to a list of validated runs), and 'Disclaimer' (stating the data is released under a Creative Commons CC0 waiver). A 'Public Domain' logo is present, and an 'Export MARCXML' button is at the bottom right of the panel. The footer includes the CERN logo, copyright information '© 2014 CERN Open Data', and logos for ALICE, ATLAS, CMS, and LHCb.

Software & documentation

🏠 > CMS > CMS Tools

Software to preprocess the Muon and Electron datasets for the two-lepton/four-lepton analysis example

Rodriguez Marrero, Ana

Cite as: Rodriguez Marrero, A. (2014). Software to preprocess the Muon and Electron datasets for the two-lepton/four-lepton analysis example. CERN Open Data Portal. DOI: [10.7483/OPENDATA.CMS.GS6N.54B9.2](https://doi.org/10.7483/OPENDATA.CMS.GS6N.54B9.2)

Collection CMS Tools Accelerator CERN-LHC Experiment CMS

Description

Software to produce the intermediate data files, derived from the [primary datasets](#), for a simple analysis on Z decays to two leptons and the ZZ decays to four leptons.

URLs

<https://github.com/ayrodrig/pattuples2010>

How can you use this?

See the run instructions in [GitHub](#)

Use with..

Use this with the following datasets:
Electron primary dataset in AOD format from RunB of 2010 (/Electron/Run2010B-Apr21Reco-v1/)
Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21Reco-v1/)

CMS Virtual Machines: How to install

The CMS-specific VM includes the [ROOT framework](#) and [CMSSW](#). Follow the instructions below to setup a CERN Virtual Machine on your computer with CMS data

1. How to install a CERN VM
2. Issues & Limitations

How to install a CERN Virtual Machine

Step 1: Installing VirtualBox

VirtualBox is a free, open source and multiplatform application to run virtual machines; you can [download](#) the package for your platform from the [VirtualBox website](#). You will need administrative ("root") privileges on every platform to perform the installation of VirtualBox.

Note: the latest tested version of VirtualBox working with CernVM is 4.3.14. If you have troubles with the latest version of VirtualBox, pick that one: the previous version is available [on a different page](#).

Forthcoming: new data releases

270 TB primary and simulated data



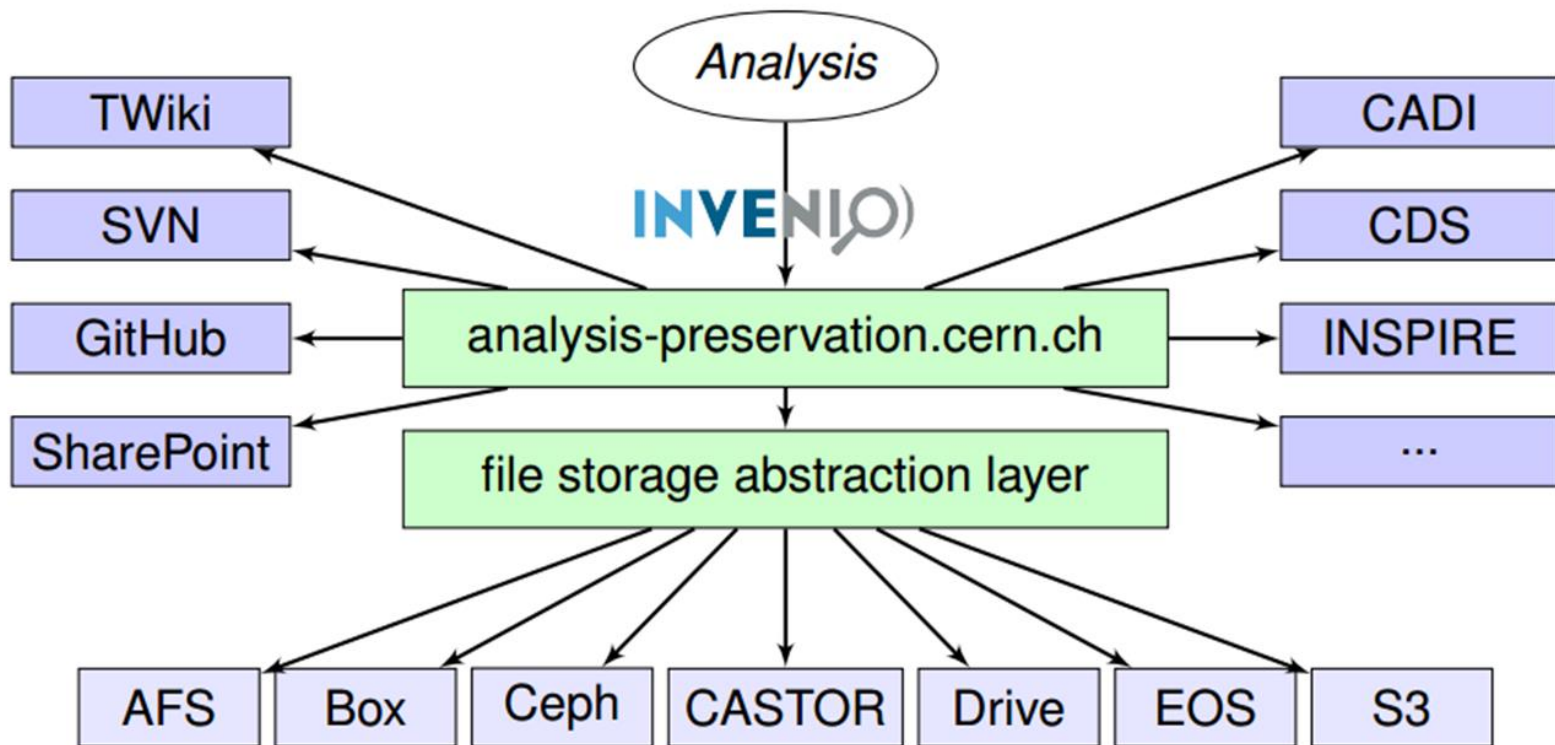
CERN Analysis Preservation

“**closed** counterpart” to CERN Open Data that captures the complexity of

- The data
- The processing steps
- Code involved
- Documentation, Physics information
- Peer review, QA

i.e. all the information contributing to the research claim/presentation/publication to enable future reuse

CERN Analysis Preservation



Work in progress



DEMO

 patricia.herterich@cern.ch [LOG OUT](#)

 Start typing

in [All Collections](#)

2 records found.



ALICE



ATLAS



CMS




LHCb

Capturing what is relevant

The image shows two overlapping web application windows. The top window is titled "Post-AOD Processing" and contains several form sections: "OS" with "Name" and "Version" fields; "Analysis Software" with "Software Used", "Version", and "Tag" fields; "User Code" with a "URL" field, a "Tag" field, and "YES" and "NO" radio buttons; and "Input Data Files" with two empty input fields. The bottom window is titled "Physics Information" and contains a "Final State Particles" section with a "Select Particle" dropdown, "Number", "PT Cut", and "ETA Cut" fields. Below this is a "Cuts" section with a "+ Add New Item" button. The "Vetos" section contains a table with columns for "Particle", "Number", "PT", and "ETA", and a "Select Particle" dropdown menu. The dropdown menu is open, showing a list of options: "Select Particle", "electron", "muon" (highlighted in blue), and "bjet". A large red stamp with the word "DEMO" is overlaid on the right side of the image.

DEMO

Integrating other resources



Documentations

CADI ID

URL

Keyword

Comment

Internal Discussions

URL

Presentations

URL

Publications

Journal Title

Journal Year

Journal Volume

Journal Issue

Journal Page

Identifiers [+ Add New Item](#)




Permissions



Basic Info JSON Permissions

Type email to FILTER or ADD access rights

Private

User	Index	Read	Write
patricia.herterich@cern.ch			

Open challenges

- Search syntax for CERN Analysis Preservation
- Integration of CERN Analysis Preservation in publication approval processes, incl. rights and permissions for searching, editing, reviewing
- DOI minting and versioning for content in the CERN Analysis Preservation system
- Publication of preserved analyses, e.g. in the CERN Open Data Portal

Lessons learnt

- Successful services are a collaborative effort between library, IT and the researchers
- Services need to be easy to use for researchers in their everyday workflows and re-use a lot of existing sources
- Don't be afraid of testing the service often and doing many iterations!

THANKS TO

CERN IT J. Cowton, J. Delgado, J. Kunčar, M. Neumann, T. Smith, T. Šimko

CERN SIS S.Dallmeier-Tiessen, A. Dani, P. Fokianos, L. Rueda

ALICE M. Gheata, C. Grigoras, M. Zimmermann

ATLAS K. Cranmer, L. Heinrich, D. Rousseau, F. Socher

CMS A. Calderon, A. Huffman, K. Lassila-Perini, T. McCauley, A. Rao, A. Rodriguez Marrero

LHCb S. Amerio, M. Bettler, B. Couturier, T. Head, A. Trisovic, A. Ustyuzhanin

CERN CernVM J. Blomer

CERN EOS L. Mascetti

DASPOS M. Hildreth, C. Vardeman, G. Watts

DPHEP F. Berghaus, J. Shiers

Work sponsored by the Wolfgang Gentner Programme of the Federal of Education and Research

SPONSORED BY THE



Federal Ministry
of Education
and Research

Sources

- CERN convention:
<http://council.web.cern.ch/council/en/governance/convention.html>
- Data policies:
 - ALICE Collaboration (2013). ALICE data preservation strategy. CERN Open Data Portal. <http://doi.org/10.7483/OPENDATA.ALICE.54NE.X2EA>
 - ATLAS Collaboration (2014). ATLAS Data Access Policy. CERN Open Data Portal. <http://doi.org/10.7483/OPENDATA.ATLAS.T9YR.Y7MZ>
 - CMS Collaboration (2012). CMS data preservation, re-use and open access policy. CERN Open Data Portal. <http://doi.org/10.7483/OPENDATA.CMS.UDBF.JKR9>
 - LHCb Collaboration (2013). LHCb External Data Access Policy. CERN Open Data Portal. <http://doi.org/10.7483/OPENDATA.LHCb.HKJW.TWSZ>