



---

# CERN Open Data and Data Analysis Knowledge Preservation

Tibor Šimko



Digital Library 2015 · 21–23 April 2015 · Jasná, Slovakia

1

# Invenio

# What is Invenio?

- digital library and document repository software
  - mature platform: first public release in 2002
  - rich data: articles, books, notes, photos, videos, *software*, *data*
- some Invenio-based services at CERN:



- co-developed by an international collaboration



- participating in EU FP7/H2020 projects



Welcome to INSPIRE, the High Energy Physics information system. Please direct questions, comments or concerns to [feedback@inspirehep.net](mailto:feedback@inspirehep.net).

HEP :: HEPNAMES :: INSTITUTIONS :: CONFERENCES :: JOBS :: EXPERIMENTS :: JOURNALS :: HELP

Information References Citations Filter Plots HepData

Search for new phenomena in final states with large jet multiplicities and missing transverse momentum at  $\sqrt{s} = 8 \text{ TeV}$  proton-proton collisions using the ATLAS experiment - ATLAS Collaboration (Aad, Georges *et al.*) JHEP 1310 (2013) 130 arXiv:1308.1841 [hep-ex] CERN-PH-EP-2013-110

THIS DATA COMES FROM DURHAM HEPDATA PROJECT

DATASETS:

Description: MET/sqrt(HT) distributions for the multi-jet + flavour stream with PTmin=50 GeV, and exactly eight jets, with the signal region selection, other than that on MET/sqrt(HT) itself. A selection of zero b-jets is applied.

Go to the record

Plain

$|ETARAP(B^-jet)| < 2.5$   
 $N - BJETS(p_T > 40 \text{ GeV}) = 0.0$   
 $N - JETS(p_T > 50 \text{ GeV}) = 8.0$   
 $p_T \rightarrow JETS MM$   
 $= dATA$        $= MC$   
 $|ETARAP(jet)| < 2.0$

$ET_{MISSING}/\sqrt{HT} (\text{GeV}^{0.5})$        $EVENTS/4 \text{ GeV}^{0.5}$

Events per bin

ET(C=MISSING)/SQRT(HT) IN GeV

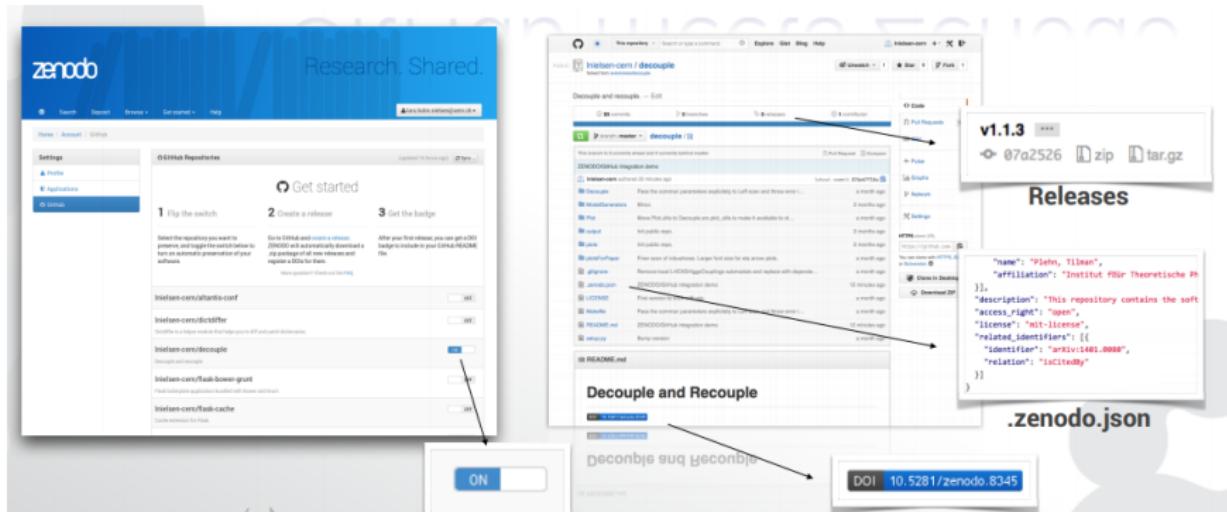
DP:1 DP:2

Collapsing

Bin Range	DP:1	DP:2
0.25	20504.0	21017.74 ± 0.25
0.75	42632.0	42806.83 ± 0.00
1.25	35848.0	35159.43 ± 0.00
1.75	19376.0	18926.27 ± 0.00
2.25	7872.0	7742.63 ± 0.00
2.75	2720.0	2686.48 ± 1.43
3.25	792.0	880.64 ± 4.83
3.75	120.0	120.00 ± 0.00

TS4 GEV<sup>0.5</sup>

- automated GitHub ↔ Zenodo bridge
- push new release to GitHub → automatic archival on Zenodo
- software preserved, minted with a DOI, made citable



<https://guides.github.com/activities/citable-code>

- link data (DATAVERSE) to code (ZENODO) to papers (INSPIRE)
- example: hep-ex/0011057, arXiv:1401.0080

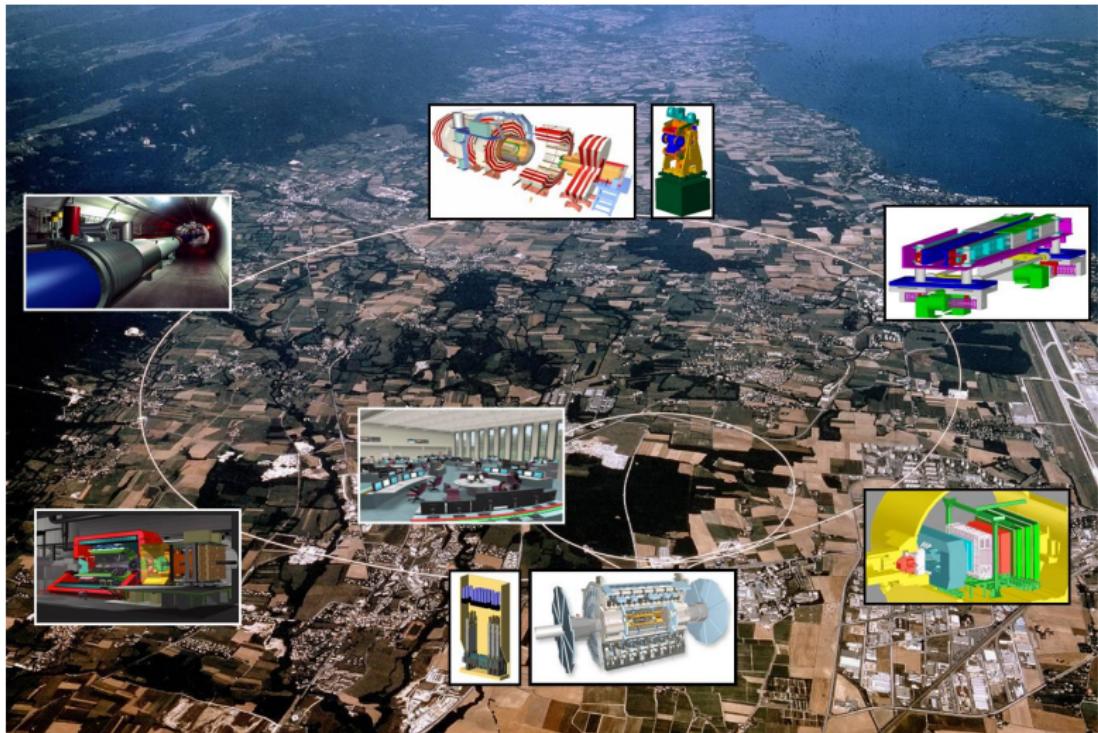
The screenshot displays three interconnected web pages:

- Zenodo Page:** Shows a search result for "decouple software associated to arXiv:1401.0080". It includes a table with columns: Name, Date, and Version. A "View on Zenodo" button is present.
- INSPIRE Home Page:** Features the INSPIRE logo and navigation links like Home, HomeNames, Institutions, Conferences, Jobs, Experiments, Journals, and Help. The URL is <http://inspirehep.net/>.
- INSPIRE Record Page:** For arXiv:1401.0080. It shows the title "A Novel Approach to Higgs Coupling Measurements", authors (Cranmer, Kyle; Kraiss, Sven (New York University)), and the arXiv ID. It includes sections for Description, Abstract, Notes, References (4), Citations (0), Plots, and Authors. The URL is <http://inspirehep.net/search?p=find+arXiv:1401.0080>.

2

# Data Analysis

# CERN LHC Experiments

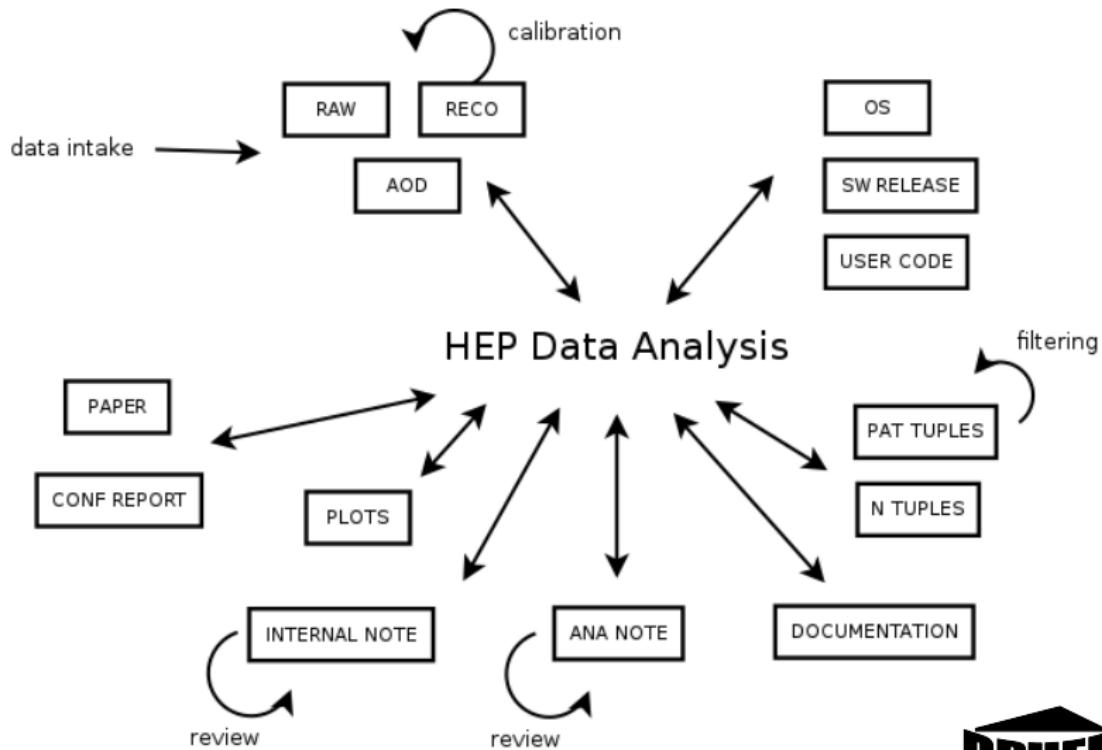


# Large Scale Solutions



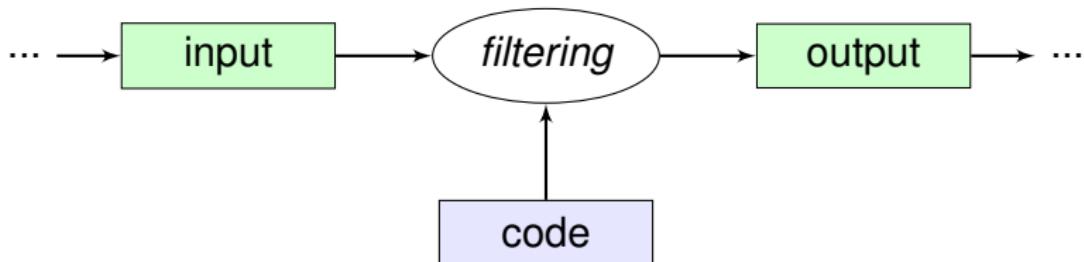
Primary site: 100k cores (10k nodes), 100k disks (50 PB), 21k NIC  
Grid: 13 Tier-1 sites, 155 Tier-2 sites, 10 Gbps optical fibre

# Preserve an Analysis?



# Big Data?

<b>data</b>	<b>scale</b>	<b>knowledge</b>
raw	~GB / sec	calibration, conditioning
reconstructed	~PB / year	filtering, selection
reduced	~TB / analysis	user code, physics objects
publication	~GB / analysis	correlation, data behind plots



Analysis Train

**Final selection step**

OS: SLC 5 x

Analysis software: CMSSW 5\_3\_x

User code: Please enter URL to your code Tag? e.g v2.3

Example of supported repositories:  
git@github.com:johndoe/myrepo.git  
svn@svnweb.cern.ch:cern/wsvn/myrepo

Harvest?  yes, harvest user code  
 no, create link only

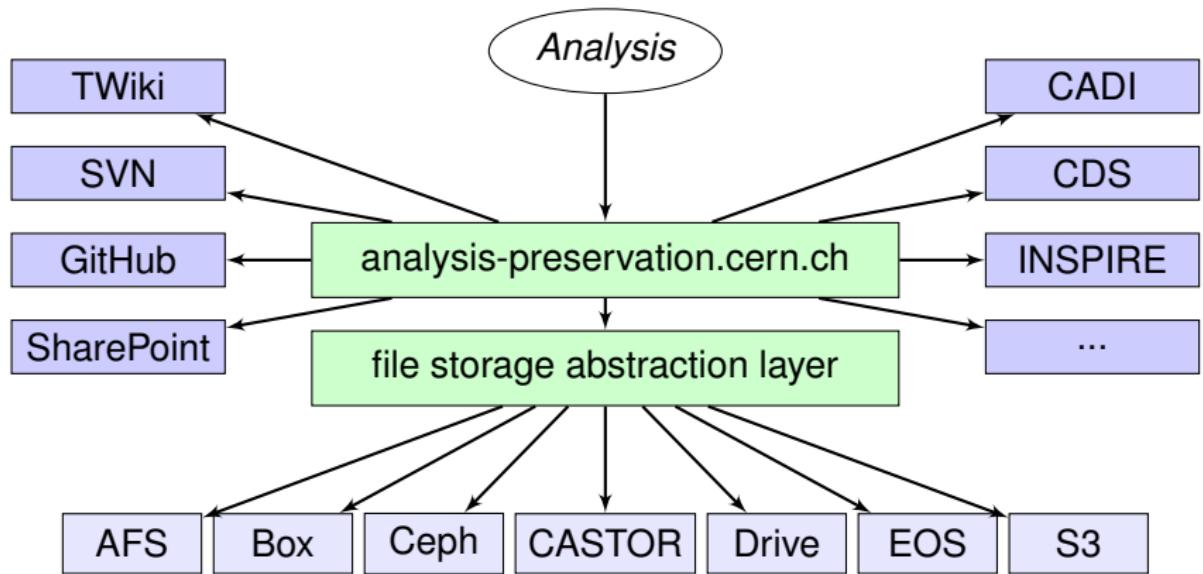
Input Data Files:  taken from output of pre-selection step  
 taken from output of custom analysis step

Output Data Files: Please enter path to data files used + Add another file

Example of supported protocols:  
xroot://castorpublic.cern.ch/castor/cern.ch/user/johndoe/mydir/myfile.root  
root://eospublic.cern.ch//eos/lhcb//myfile.root  
file:///tmp/myfile.root  
http://john.doe.example.org/myfile.root

Harvest?  yes, harvest files  no, create link only

# System Architecture



- record format: extended MARC21

- “technical” metadata: beyond bytes

- e.g. 256 “computer file characteristics”

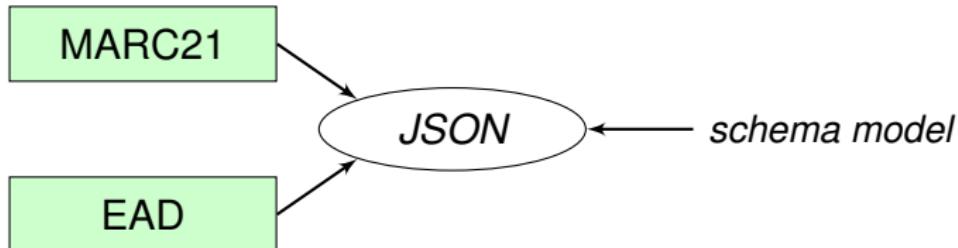
\$a characteristics	\$e events	\$t text
\$b bytes	\$f files	...

- “knowledge” metadata: semantics

- e.g. 505 “formatted contents note” CSV column information

\$t title	\$g miscellaneous
-----------	-------------------

- internal format: JSON



3

# Open Data

## Data policies:

- restricted → embargo period → open

*[...] Data with high abstraction, such as AOD, will be conditionally made publicly available after an embargo period of 5 years after publication for 10% of the data and 10 years for 100% of the data [...]” —ALICE Data Policy*

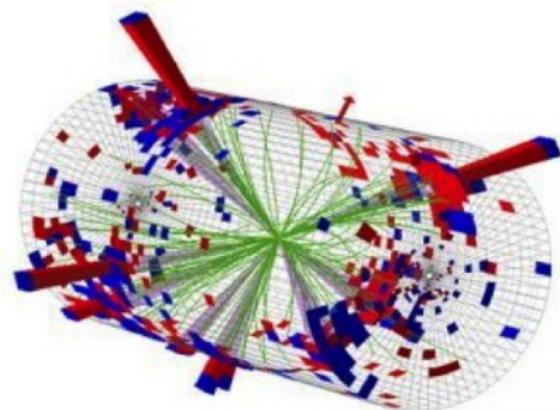
## Challenges:

- audience:

- data miners
  - citizen scientists
  - high-school students
  - general public

- computing:

- exploring in the browser
  - specialised VMs



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](#)

© 2014 CERN Open Data Terms of Use Privacy Policy Contact

@tiborsimko · @inveniosoftware

17 / 26

## Education



The CMS (Compact Muon Solenoid) experiment is one of two large general-purpose detectors built on the Large Hadron Collider (LHC). Its goal is to investigate a wide range of physics such as the characteristics of the Higgs boson, extra dimensions or dark matter.

[Explore CMS >](#)



ALICE (A Large Ion Collider Experiment) is a heavy-ion detector designed to study the physics of strongly interacting matter at extreme energy densities, where a phase of matter called quark-gluon plasma forms. More than 1000 scientists are part of the collaboration.

[Explore ALICE >](#)



The ATLAS (A Toroidal LHC Apparatus) experiment is a general purpose detector exploring topics like the properties of the Higgs-like particle, extra dimensions of space, unification of fundamental forces, and evidence for dark matter candidates in the Universe.

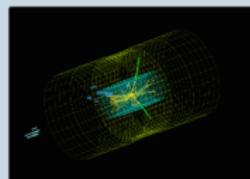
[Explore ATLAS >](#)



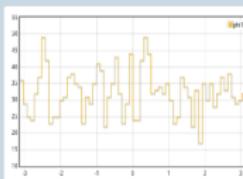
The LHCb (Large Hadron Collider beauty) experiment aims to record the decay of particles containing b and anti-b quarks, known as B mesons. The detector is designed to gather information about the identity, trajectory, momentum and energy of each particle.

[Explore LHCb >](#)

For education purposes, the complex primary data need to be processed into a format (examples below) that is good for simple applications. Get in touch if you wish to build your own applications similar to those shown here.



[Visualise events >](#)



[Visualise histograms >](#)



[Learning Resources >](#)

opendata

ABOUT SEARCH EDUCATION RESEARCH

Home Education Visualise Events CMS

Explore CMS open data and visualise detector events

Need HELP?

/Mu.Ig:Events/Run\_146436/Event\_90626440

Detector Model

- Tracker Barrels
- Tracker Endcaps
- ECAL Barrel
- ECAL Endcaps
- ECAL Preshower
- HCal Barrel
- HCal Endcaps
- HCal Outer
- HCal Forward
- Drift Tubes (muon)
- Cathode Strip Chambers (muon)
- Resistive Plate Chambers (muon)

Tracking

- Tracks Reco.

ECAL

- Barrel Rec. Hits
- Endcap Rec. Hits
- Preshower Rec. Hits

HCal

- Barrel Rec. Hits
- Endcap Rec. Hits
- Forward Rec. Hits
- Outer Rec. Hits

Muon

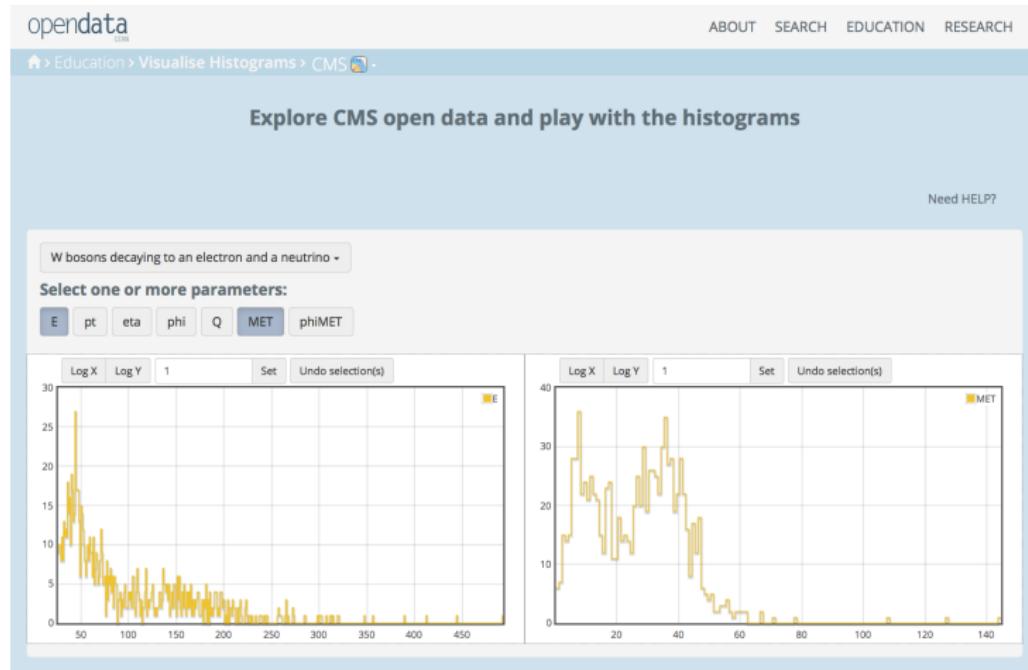
- Matching muon chambers

Physics Objects

- Electron Tracks (GSF)
- Tracker Muons (Reco)
- Stand-alone Muons (Reco)
- Global Muons (Reco)
- Calorimeter Energy Towers
- Jets
- Missing Et (Bacon)

A 3D visualization of a particle collision event in the CMS detector. The event shows a central collision vertex with outgoing particles, reconstructed tracks, and energy deposits in calorimeters. A coordinate system is shown at the bottom right.

# Basic histogramming



## Research



To analyse CMS data, a Virtual Machine with the CMS analysis environment is provided. The data can be accessed directly through the VM. In the primary datasets, no selection nor identification criteria have been applied. For this release, no simulated Monte Carlo datasets are provided.

[Explore CMS >](#)



According to the ALICE data preservation strategy, reconstructed data and Monte Carlo data as well as the analysis software and documentation needed to process them will be made available on a time scale of 5 years (for 10% of the data). Thus, the first release of ALICE research data will happen in 2018.



According to the ATLAS Data Access Policy, reconstructed data and accompanying tools will be released after reasonable embargo periods.



According to the LHCb External Data Access Policy, reconstructed data and accompanying tools will be released after reasonable embargo periods.

For research purposes, specific software environments and tools need to be deployed to analyse these complex primary data. In addition to the data below, you will find instructions for setting up your working environments [here](#).



[Install your Virtual Machine >](#)



[Start analysing the data >](#)

## Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21ReReco-v1/AOD) 2014

/Mu/Run2010B-Apr21ReReco-v1/AOD

CMS collaboration

Cite as: CMS collaboration (2014). Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21ReReco-v1/AOD). CERN Open Data Portal. DOI:  
[10.7483/OPENDATA.CMS.B8MR.C4A2](https://doi.org/10.7483/OPENDATA.CMS.B8MR.C4A2)

Collection

CMS Primary Datasets

Collision Energy

7TeV

Accelerator

CERN-LHC

Experiment

CMS

## Description

Mu primary dataset in AOD format from RunB of 2010

## Characteristics

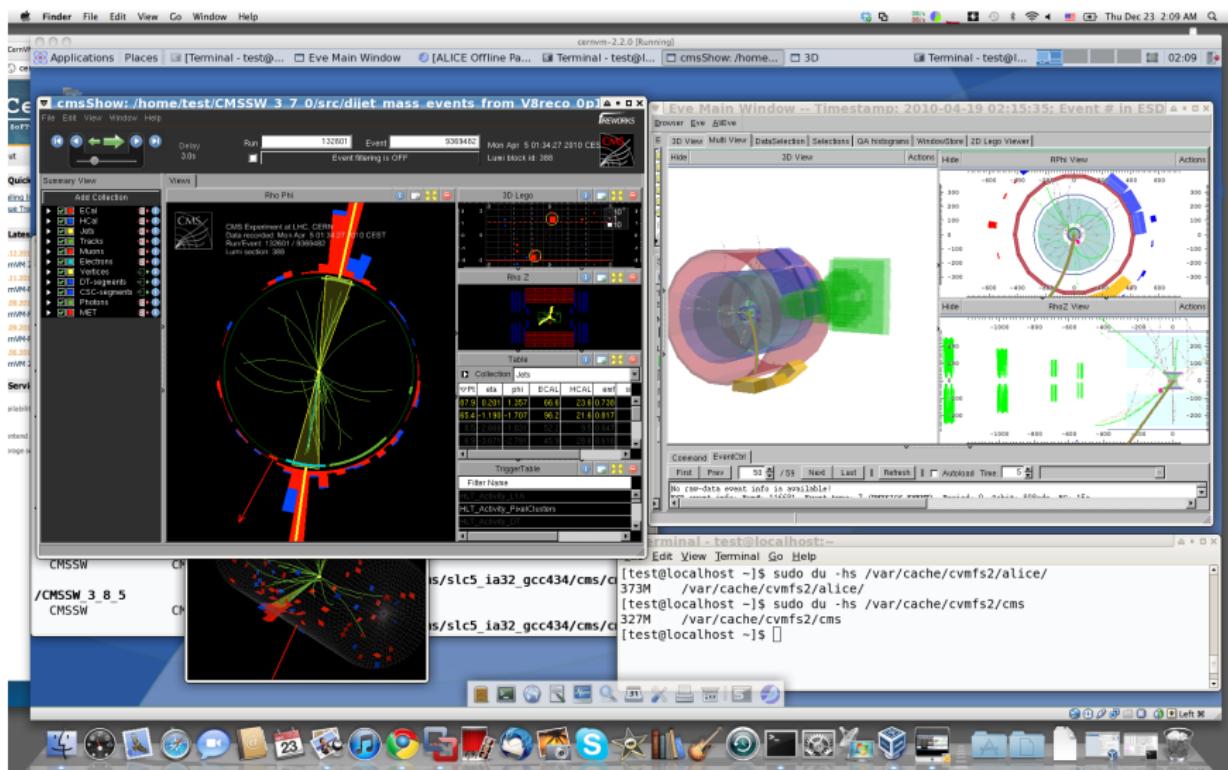
Dataset: 32376291 events 2979 files 3.2 TB in total

## System Details

Software release: CMSSW\_4\_2\_1\_patch1

## Indexes

# CernVM Virtual Machine



# Open Data? Who cares?

**NewScientist** Physics & Math

Home News In-Depth Articles Opinion CultureLab Galleries Topic Guides Last Word Subscribe Dating

SPACE TECH ENVIRONMENT HEALTH LIFE PHYSICS&MATH SCIENCE IN SOCIETY

Home | Physics & Math | News

**Run your own experiment using CERN's public LHC data**

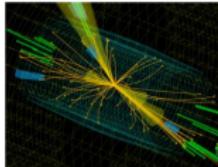
18:30 24 November 2014 by Jacob Aron  
For similar stories, visit the [The Large Hadron Collider](#) Topic Guide

Why build your own particle accelerator when you can borrow CERN's? The home of the Large Hadron Collider near Geneva, Switzerland, has started putting data from its experiments online for anyone to use. They hope it could fuel education, art and perhaps even physics discoveries.

"It's very important that we keep this data open and usable," says Kati Lassila-Perini of the CMS experiment at the LHC, which has uploaded 27 terabytes of data to the new [CERN Open Data Portal](#). A web interface lets you visualise the paths of particles created by collisions at the LHC, or you can work directly with the data for more serious analysis.

Other LHC experiments have uploaded smaller data sets for educational purposes, allowing budding particle physicists to try their hand at massive-scale physics. The data could also be turned into art or music, as was done previously at a [CERN arts festival](#).

It's possible that physicists might search through the data and make discoveries, but most people with the necessary skills already work with similar experiments, says Lassila-Perini. "There are not so many spare physicists around."



Maybe you'll be able to spot a particle in data from the LHC (Image: CERN)

ADVERTISEMENT

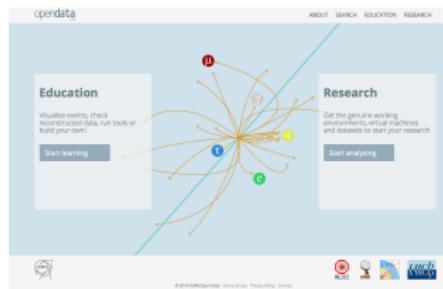
Like 1021 Tweet 299 Share 294



- 82,000 distinct users visited the site
- 21,000 distinct users viewed data records
- 16,000 distinct users used event display
- 3,000 distinct users used histogramming

7

# Conclusions



*Capturing and disseminating knowledge  
of data, code, platform, processes  
to enable future data reuse*

## (Open) Data Analysis Preservation Framework

<http://opendata.cern.ch/>

**CERN IT** J. Cowton, P. Fokianos, J. Kunčar, T. Smith, T. Šimko

**CERN Library** S. Dallmeier-Tiessen, P. Herterich, L. Rueda

**ALICE** M. Gheata, C. Grigoras

**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, B. Couturier, A. Trisovic

**CERN CernVM** J. Blomer

**CERN EOS** L. Mascetti

**DASPOS** M. Hildreth

**DPHEP** F. Berghaus