# Open data provenance and reproducibility: a case study from publishing CMS open data
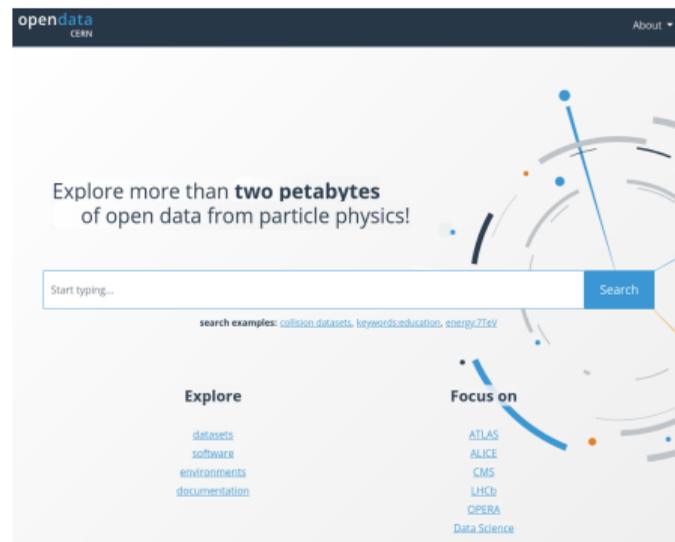
Tibor Šimko[1]    Heitor Pascoal de Bittencourt[2]    Edgar Carrera[2]    Clemens Lange[2]
Kati Lassila-Perini[2]    Lara Lloret[2]    Tom McCauley[2]    Jan Okraska[1]
Daniel Prelipcean[1]    Mantas Savaniakas[2]
on behalf of the CERN Open Data team and the CMS Collaboration

[1]CERN Open Data team   [2]CMS Collaboration

*24th International Conference on Computing in High Energy and Nuclear Physics (CHEP)*
*Adelaide, Australia, 4–8 November 2019*

# CERN Open Data

- launched in November 2014
- rich content
  - collision and simulated datasets for research
  - derived datasets for education
  - configuration files and documentation
  - virtual machines and container images
  - software tools and analysis examples
- total size in November 2019
  - over 7'000 bibliographic records
  - over 800'000 files
  - over 2 petabytes
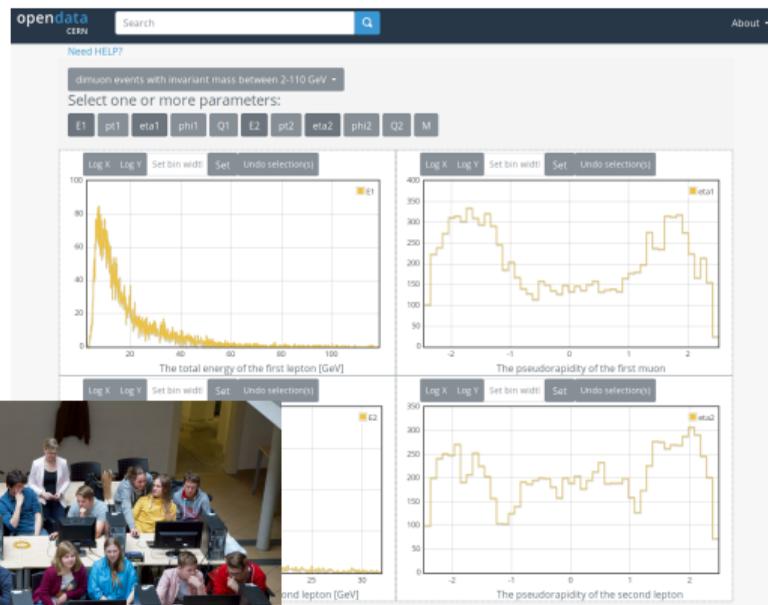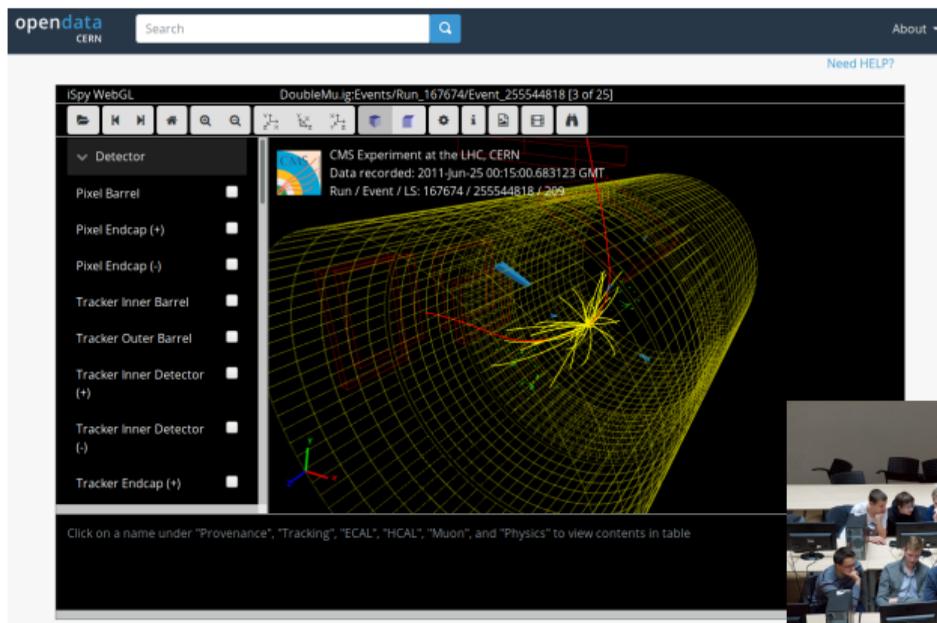


`http://opendata.cern.ch`

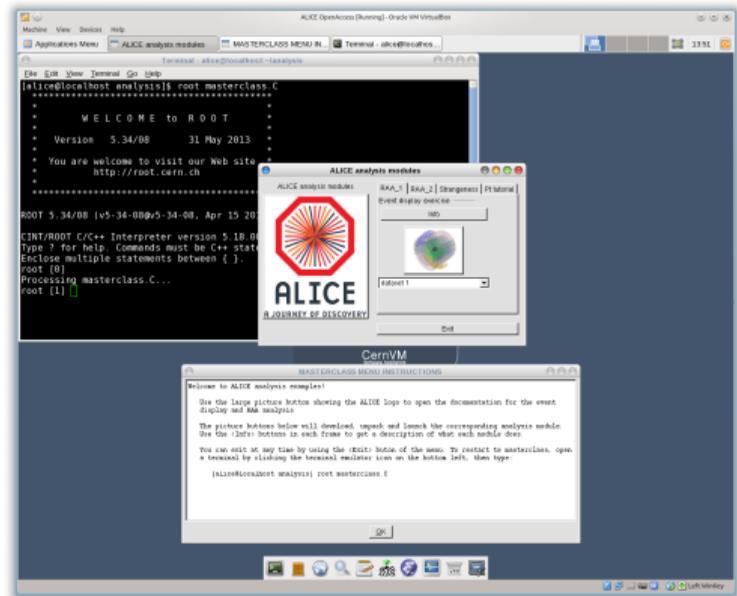Developed by CERN-IT in close collaboration with Experiments

# Education-oriented use cases



Interactive event display and histogramming for derived datasets

# Research-oriented use cases



Run CernVM Virtual Machines



Run realistic physics analysis examples

# Enables independent theoretical research



arXiv:1704.05066    arXiv:1807.11916    arXiv:1902.04222

Searches, QCD jet studies, Machine Learning. . .

Over twenty papers citing CMS open data

. . . that the CMS collaboration start to cite!

[33] A. Larkoski et al., "Exposing the QCD splitting function with CMS Open Data", *Phys. Rev. Lett.* **119** (2017) 132003, doi:10.1103/PhysRevLett.119.132003, arXiv:1704.05066.

[34] A. Tripathee et al., "Jet Substructure Studies with CMS Open Data", *Phys. Rev. D* **96** (2017) 074003, doi:10.1103/PhysRevD.96.074003, arXiv:1704.05842.

# New CMS open data release



**Release highlights**

- This is the fourth release of high-level CMS open data, following release of around 50% of data from the LHC's Run 1: 2010 data in 2014, 2011 data in 2016, and 2012 data in 2017. This brings the volume of CMS open data to more than 2 PB.
- The release includes datasets prepared specifically for use in Machine Learning or in data science.
  - A dataset derived from Run 2 simulation data is devoted to the challenge of event and object tagging in events with two b quarks produced from the decay of a Higgs boson. It is particularly difficult to distinguish this Higgs signal channel from the background.
  - Further datasets derived from Run 1 and Run 2 simulated data are devoted to identifying top quarks produced in events and to studying the flavour content of jets.
  - Another dataset is devoted to the challenge of particle tracking in the future era of high-luminosity collisions and is derived from simulations of collisions in the tracker after Phase 2 upgrades.
- The parent datasets and production workflows for the ML samples also available for full reproducibility.
  - These include the first simulation samples in the "MiniAODSIM" format in use in Run 2 data analysis. Small samples of raw data are released, useful for testing of the data-processing chain and eventually reconstruction-algorithm development.
- Instructions are now available on how to generate simulated events in the open data environment.
- The release completes the 2010 data release with now all proton-proton data available publicly and adds some simulated data also for 2010 data taking.
- Contains datasets from early commissioning runs used in studies with CASTOR calorimeter and corresponding simulations.
- In addition to already available 2012 simulation data, large amount of 2012 simulation data of rarely used processes is now available on demand.
- Search functionality for simulation data is now available based on physics processes.

Latest batch of CMS open data was released in Summer 2019

# Example 1: Data provenance of simulated datasets



Simulated dataset BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph in MINIAODSIM format for 2016 collision data

/BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIISummer16MiniAODv2-PUMoriond17_80X_mcRun2_asymptotic_2016_TrancheIV_v6-v1/MINIAODSIM, CMS Collaboration

Cite as: CMS Collaboration (2019). Simulated dataset BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph in MINIAODSIM format for 2016 collision data. CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.7N4X.Z7FA

`Dataset` `Simulated` `Exotica` `Gravitons` `CMS` `13TeV` `CERN-LHC`

## How were these data generated?

These data were generated in several steps (see also CMS Monte Carlo production overview):

**Step LHE**
Release: CMSSW_7_1_16
Output dataset: /BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIIWinter15wmLHE-MCRUN2_71_V1-v1/LHE
Note: To get the exact generator parameters, please see Finding the generator parameters.

**Step SIM**
Release: CMSSW_7_1_20
📄 Configuration file for SIM (link)
Output dataset: /BulkGravTohhTohbbhbb_narrow_M-4500_13TeV-madgraph/RunIISummer15GS-MCRUN2_71_V1-v1/GEN-SIM

**Step HLT RECO**
Release: CMSSW_8_0_21

▶ full capture of data generation steps
▶ full capture of compute environments
▶ full capture of configuration files
▶ full capture of production scripts

Data records come with full provenance information

# Capturing data provenance via ad-hoc curation techniques



Dedicated data curation scripts



CMS DAS          CMS McM

Mining several CMS collaboration sources

# Harmonising year-dependent sources

```
"methodology": {
  "description": "<p>These data were generated in several steps (see also <a href=\"/docs/cms-mc-production-overview\">CMS M
  "steps": [
    {
      "configuration_files": [
        {
          "script": "#!/bin/bash\nsource /cvmfs/cms.cern.ch/cmsset_default.sh\nexport SCRAM_ARCH=slc5_amd64_gcc462\nif [ -r
          "title": "Production script"
        },
        {
          "title": "Generator parameters",
          "url": "https://cms-pdmv.cern.ch/mcm/public/restapi/requests/get_fragment/HIG-Summer12-02276"
        },
        {
          "cms_confdb_id": "a97a2f6c22dfba999c0131657a81ecfd",
          "process": "SIM",
          "title": "Configuration file"
        }
      ],
      "generators": [
        "pythia6"
      ],
      "global_tag": "START53_V7C::All",
      "output_dataset": "/BBH_HToTauTau_M_125_TuneZ2star_8TeV_pythia6_tauola/Summer12-START53_V7C-v1/GEN-SIM",
      "release": "CMSSW_5_3_13",
      "type": "SIM"
    },
    {
      "configuration_files": [
```

From year-dependent DAS/McM information to year-independent Open Data JSON schema

# Example 2: Raw data samples for 2010–2012 data

## SingleElectron primary dataset sample in RAW format from RunA of 2011 (from /SingleElectron/Run2011A-v1/RAW)

/SingleElectron/Run2011A-v1/RAW, CMS collaboration

Cite as: CMS collaboration (2019). SingleElectron primary dataset sample in RAW format from RunA of 2011 (from /SingleElectron/Run2011A-v1/RAW). CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.6O84.WLN8

`Dataset` `Collision` `CMS` `7TeV` `CERN-LHC`

### Description

A sample from SingleElectron primary dataset in RAW format from RunA of 2011. Run range [161224,163286].

This dataset contains selected runs from 2011 RunA. The list of validated lumi sections, which must be applied to all analyses on events reconstructed from these data, can be found in

CMS list of validated runs Cert_160404-180252_7TeV_ReRecoNov08_Collisions11_JSON.txt

### Dataset characteristics

**2064298** events. **116** files. **424.3 GB** in total.

### How can you use these data?

These data are in RAW format and not directly usable in analysis. The reconstructed data reprocessed from these RAW data are included in the data of this record. The reconstruction step can be repeated with the configuration file below and the resulting AOD has been confirmed to be identical with the original one with comparison code available in Validation code to plot basic physics objects from AOD

## SingleElectron primary dataset in AOD format from RunA of 2011 (/SingleElectron/Run2011A-12Oct2013-v1/AOD)

/SingleElectron/Run2011A-12Oct2013-v1/AOD, CMS collaboration

Cite as: CMS collaboration (2016). SingleElectron primary dataset in AOD format from RunA of 2011 (/SingleElectron/Run2011A-12Oct2013-v1/AOD). CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.P87Z.TXTV

`Dataset` `Collision` `CMS` `7TeV` `CERN-LHC`

### Description

SingleElectron primary dataset in AOD format from RunA of 2011. Run period from run number 160404 to 173692.

This dataset contains all runs from 2011 RunA. The list of validated runs, which must be applied to all analyses, can be found in

CMS list of validated runs Cert_160404-180252_7TeV_ReRecoNov08_Collisions11_JSON.txt

### Dataset characteristics

**41709195** events. **1542** files. **5.8 TB** in total.

### How were these data selected?

Events stored in this primary dataset were selected because of the presence of at least one high-energy electron in the event.

**Data taking / HLT**
The collision data were assigned to different RAW datasets using the following HLT configuration.

**Data processing / RECO**
This primary AOD dataset was processed from the RAW dataset by the following step:
Step: RECO
Release: CMSSW_5_3_12_patch1
Global tag: FT_R_53_LV5::All
Configuration file for RECO step reco_2011A_SingleElectron

RAW                                       AOD

# Can we reprocess raw data samples from 2010-2012?

**3. Workflow**

The workflow can be logically divided into several parts:

0. *Upload all files.*
   Some files cannot be generated at run time and need to be uploaded.

```
inputs:
  files:
    - src/PhysicsObjectsHistos.cc
    - BuildFile.xml
    - demoanalyzer_cfg.py
```

1. *Fix the CMS SW environment variables manually.*
   First, we have to set up the environment variables accordingly for the CMS SW. Although this is done in the docker image, reana overrides them and they need to be reset. This is done by invoking the cms entrypoint.sh script commands.

   See also this issue.

```
$ source /opt/cms/cmsset_default.sh
$ scramv1 project CMSSW CMSSW_5_3_32
$ cd CMSSW_5_3_32/src
$ eval `scramv1 runtime -sh`
```

2. *Create the specific CMS path.*
   CMS specific data analysis framework requires two directory levels. See also this issue.

```
$ mkdir Reconstruction && cd Reconstruction
$ mkdir Validation && cd Validation
```

3. *Create the reconstruction file.*
   See also this repo.

```
$ cmsDriver.py reco -s RAW2DIGI,L1Reco,RECO,USER:EventFilter/HcalRawToDigi/hcallaserhbhehffilter2012_cf
```

4. *Adjust the reconstruction file to the specific data file.*
   Although generated using parameters, the reconstruction file still requires changes.

```
$ sed -i 's/from Configuration.AlCa.GlobalTag import GlobalTag/process.GlobalTag.connect = cms.string("
$ sed -i 's/# Other statements/from Configuration.AlCa.GlobalTag import GlobalTag/g' reco_cmsdriver.py
$ sed -i 's/process.GlobalTag = GlobalTag(process.GlobalTag, 'FT_53_LV5_AN1::All', ''/process.GlobalTa
```

5. *Link the CVMFS files.*
   The ls -l commands are explicitly needed to make sure that the cms-opendata-conddb.cern.ch directory has actually expanded in the image, according to this guide. See also this issue.

```
$ ln -sf /cvmfs/cms-opendata-conddb.cern.ch/FT_53_LV5_AN1_RUNA FT_53_LV5_AN1
$ ln -sf /cvmfs/cms-opendata-conddb.cern.ch/FT_53_LV5_AN1_RUNA.db FT_53_LV5_AN1_RUNA.db
$ ls -l
$ ls -l /cvmfs/
```

6. *Run the reconstruction.*
   At this point all environment variables and files should be proper.

```
$ cmsRun reco_cmsdriver.py
```

7. *Adjust project structure for validation*
   Copy the required files for the next steps.

```
$ mkdir src
$ scp ../../../../src/PhysicsObjectsHistos.cc ./src
$ scp ../../../../BuildFile.xml .
$ scp ../../../../demoanalyzer_cfg.py .
```

8. *Run CMS scram command to fix libraries.*
   Most importantly, the *BuildFile.xml* has to be inside the directory where the *scram* command is executed.
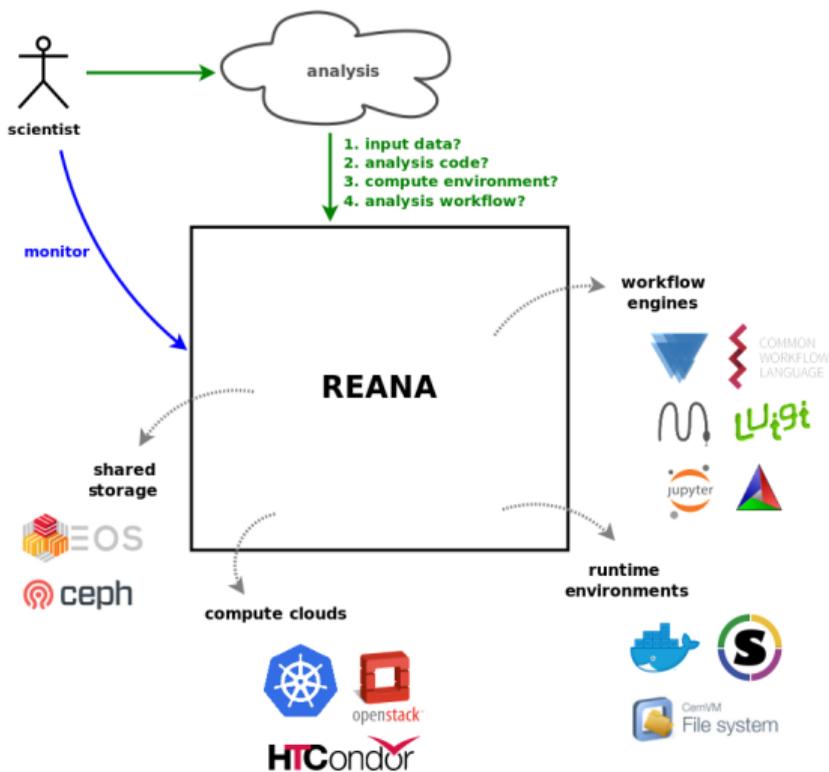
```
$ scram b
```

9. *Run the validation file.*
   See also this repo

```
$ cmsRun demoanalyzer_cfg.py
```

Workflow steps to run CMS reconstruction in CMSSW environment

# Running scientific workflows on containerised clouds



- ▶ REANA reproducible analysis platform
  `http://www.reana.io`
- ▶ multiple workflow systems
  (CWL, Serial, Yadage)
- ▶ multiple compute backends
  (Kubernetes, HTCondor, Slurm)
- ▶ multiple shared storage
  (Ceph, EOS, NFS)

reproducibility

code + data + environment + workflow

# Preserving CMS software stack environment



CMSSW docker image with "embedded" CVMFS



Condition data for open data analyses
are available on "live" CVMFS

# Automated reconstruction workflows

```
dataset=Jet
year=2011A
```

1 input parameters

↓

2 workflow factory

→

3 reana.yaml

→

4 run by REANA platform

5 serving open data files

↓

6 output histograms

Parametrised workflow runnable on REANA reproducible analysis platform

# Conclusions

CMS open data now contains detailed provenance information

- ▶ knowing "how the data came about" enhances current knowledge and future reuse
- ▶ capturing data provenance requires non-trivial information hunt and harmonisation
- ▶ *a posteriori* approach: running after $\sim$5 year old data and procedures
- ▶ *a priori* approach: ultra legacy run to generate preservation-friendly assets?

Successful RAW to AOD reconstruction tests on open data

- ▶ AOD reconstruction and histogram verification permitted to validate approach
- ▶ using non-production compute environment ensures reproducibility



`http://opendata.cern.ch`