

CMS Open data - process, challenges and plans

HSF Data Analysis Working Group - 3 July 2023



Kati Lassila-Perini
Helsinki Institute of Physics - Finland
CMS Data preservation and open access coordinator

1

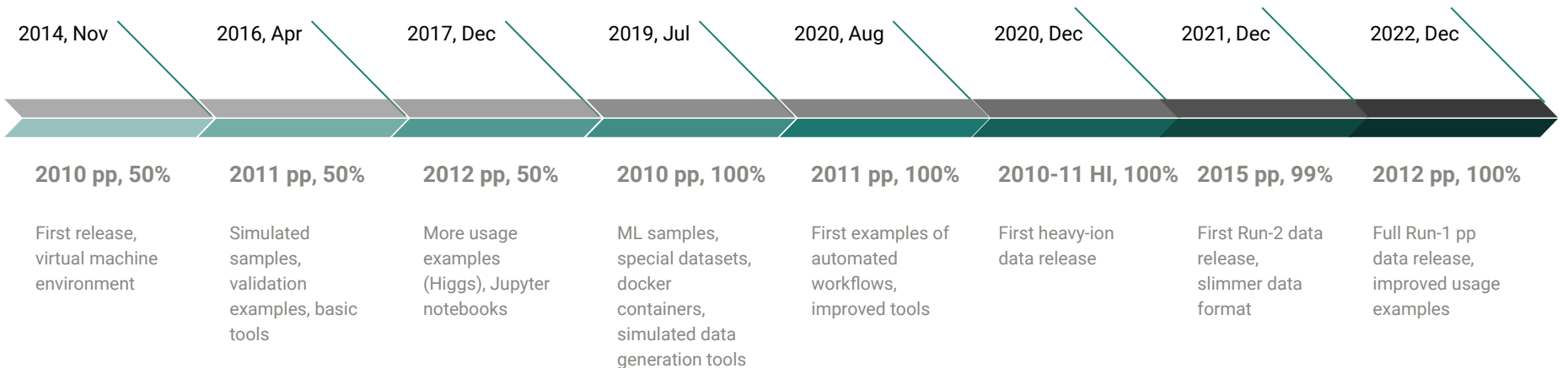
CMS Open data

Continuous releases since 2014

- open data appreciated and in use



Release timeline



For details, see “[CMS Open data](#)” at a recent workshop by Julie Hogan, DPOA co-convener

INSPIRE HEP literature references.reference.doi:10.7483/OPENDATA.CMS*

74 results | cite all Citation Summary lit Most Recent

Date of paper

2015 2023

Number of authors

Single author 12

10 authors or less 64

Exclude RPP

Exclude Review of Particle Physics 74

Document Type

article 48

published 33

conference paper 22

thesis 4

Potential of the Julia programming language for high energy physics computing #1

J. Eschle (U. Zurich (main)), T. Gal (Erlangen - Nuremberg U., Theorie III), M. Giordano (Imperial Coll., London), P. Gras (IRFU, Saclay), B. Hegner (CERN) et al. (Jun 6, 2023)

e-Print: 2306.03675 [hep-ph]

pdf cite claim reference search 0 citations

Baler -- Machine Learning Based Compression of Scientific Data #2

Fritjof Bengtsson (Lund U. (main)), Caterina Doglioni (Manchester U.), Per Alexander Ekman (Lund U. (main)), Axel Gallén (Lund U. (main)), Pratik Jawahar (Manchester U.) et al. (May 3, 2023)

e-Print: 2305.02283 [physics.comp-ph]

pdf cite claim reference search 0 citations

Quantum Generative Adversarial Networks For Anomaly Detection In High Energy Physics #3

Elie BERMOT (IBM, Zurich and ETH, Zurich (main)), Christa Zoufal (IBM, Zurich), Michele Grossi (CERN), Julian Schuhmacher (IBM, Zurich), Francesco Tacchino (IBM, Zurich) et al. (Apr 27, 2023)

e-Print: 2304.14439 [quant-ph]

pdf cite claim reference search 0 citations

Search for the production of dark fermion candidates in association with heavy neutral gauge boson decaying to dimuon in proton-proton collisions at $\sqrt{s} = 8$ TeV using the CMS open data #4

CMS open data are actively in use - no problems have ever arisen!
Search inspire

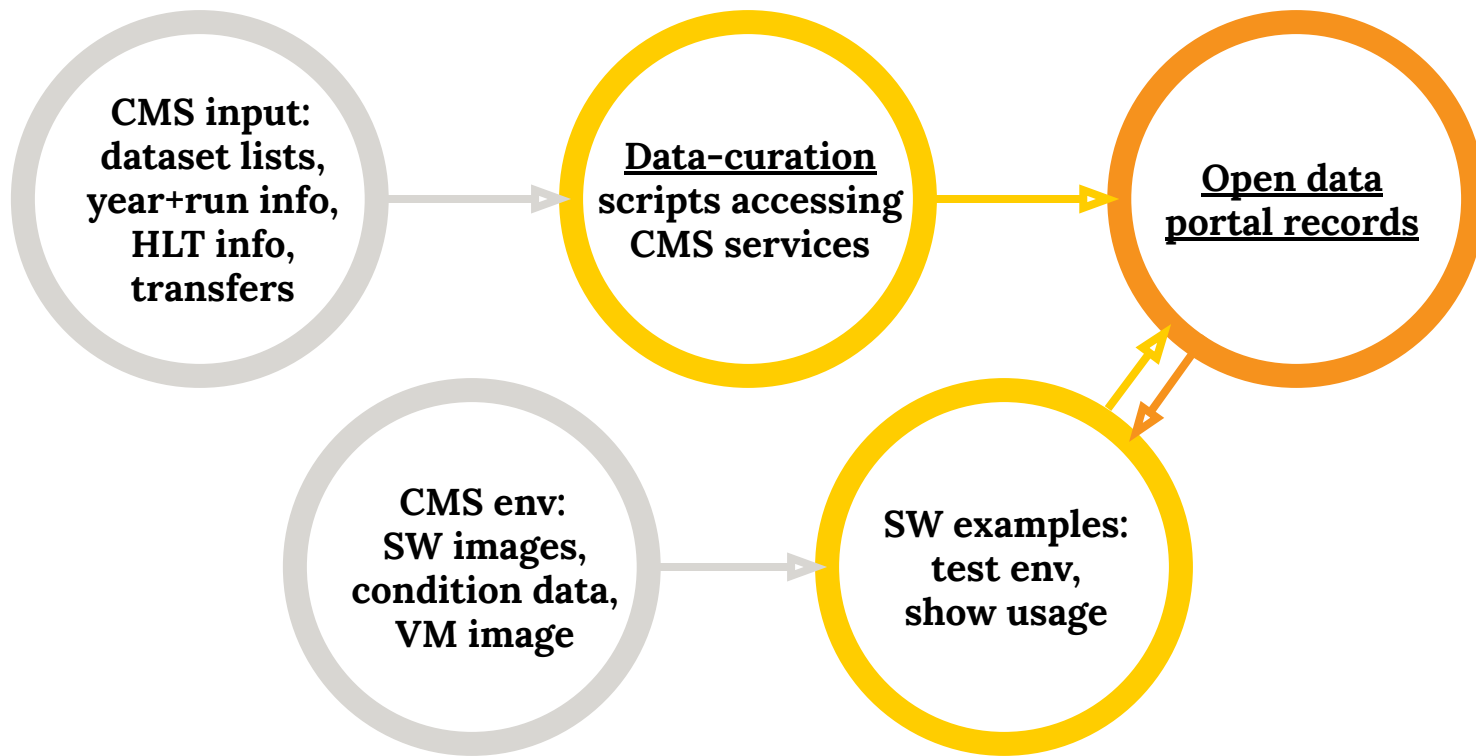


2

The process

What CMS needs → what external users need

Release procedure



**Contextual
metadata:
HOW TO
USE?**

Dataset characteristics

76523854 events, 1607 files, 1.5 **TB** in total.

System details

Recommended [global tag](#) for analysis: 76X_dataRun2_16Dec2015_v0

Recommended release for analysis: CMSSW_7_6_7

How were these data selected?

Events stored in this primary dataset were selected because of the presence of a high scalar sum of the jet transverse momenta (HT); or at least one or two energetic jets.

Data taking / [HLT](#)

The collision data were assigned to different RAW datasets using the following [HLT configuration](#).

Data processing / [RECO](#)

This primary MINIAOD dataset was processed from the RAW dataset by the following step:

Step: RECO

Release: CMSSW_7_6_3

Global tag: 76X_dataRun2_v15

[Configuration file for RECO step reco_2015D_JetHT](#)

[HLT trigger paths](#)

The possible [HLT trigger paths](#) in this dataset are:

[HLT_AK8DIPFJet250_200_TrimMass30_BTagCSV0p45](#)

[HLT_AK8DIPFJet280_200_TrimMass30_BTagCSV0p45](#)

[HLT_AK8PFHT600_TrimR0p1PT0p03Mass50_BTagCSV0p45](#)

**Content
metadata:
WHAT?**

**Provenance
metadata:
FROM
WHERE?**



**Context:
environment
software**

Dataset characteristics

76523854 events. 1607 files. 1.5 TB in total.

System details

Recommended [global tag](#) for analysis: 76X_dataRun2_16Dec2015_v0

Recommended release for analysis: CMSSW_7_6_7

How were these data selected?

Events stored in this primary dataset were selected because of the presence of two energetic jets.

Data taking / HLT

The collision data were assigned to different RAW datasets using the following

Data processing / RECO

This primary MINIAOD dataset was processed from the RAW dataset by the following

Step: RECO

Release: CMSSW_7_6_3

Global tag: 76X_dataRun2_v15

[Configuration file for RECO step reco_2015D_JetHT](#)

HLT trigger paths

The possible HLT trigger paths in this dataset are:

[HLT_AK8DIPFJet250_200_TrimMass](#)

[HLT_AK8DIPFJet280_200_TrimMass](#)

[HLT_AK8PFHT600_TrimRop1PT0p03](#)

**Context:
trigger path**

**Context:
validated
data selection**

JetHT primary dataset in MINIAOD format from RunD of 2015 (/J.../16Dec2015-v1/MINIAOD)

/JetHT/Run2015D-16Dec2015-v1/MINIAOD, CMS collaboration

Cite as: CMS collaboration (2021). JetHT primary dataset in MINIAOD format from RunD of 2015. [Open Data Portal](#). DOI:10.7483/OPENDATA.CMS.IDN0.S11Z

[Dataset](#) [Collision](#) [CMS](#) [13TeV](#) [CERN-LHC](#)

Description

JetHT primary dataset in MINIAOD format from RunD of 2015. Run period from run number 26030 to 260627.

The list of validated runs, which must be applied to all analyses, either with the full validation or for an analysis requiring

[Validated runs, full validation](#)

**Context:
usage
instructions**

How can you use these data?

You can access these data through the CMS Open Data container or the CMS Virtual Machine. See the instructions for setting up one of the two alternative environments and getting started in

[Running CMS analysis code using Docker](#)

[How to install the CMS Virtual Machine](#)

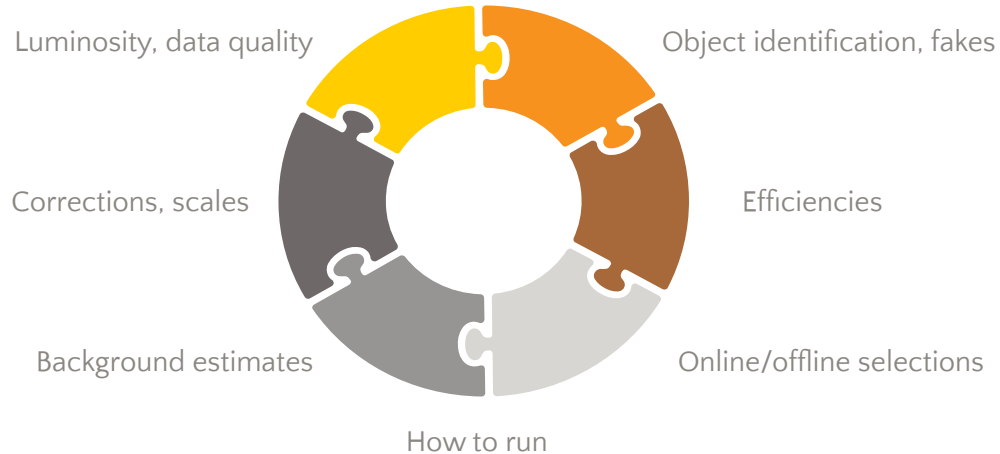
[Getting started with CMS open data](#)

File Indexes

Filename	Size	
CMS_Run2015D_JetHT_MINIAOD_16Dec2015-v1_00000_file_index.txt	2.4 kB	List Files Download
CMS_Run2015D_JetHT_MINIAOD_16Dec2015-v1_50000_file_index.txt	123.9 kB	List Files Download

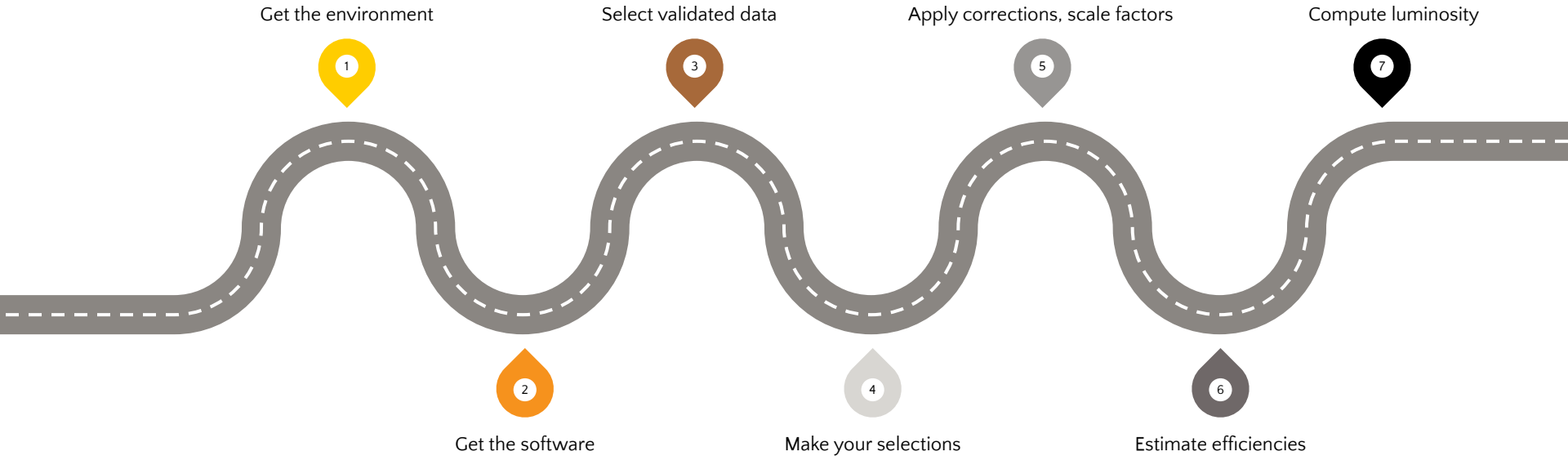


Contextual metadata - should cover a lot!





How to put this together?





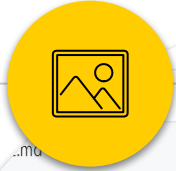
Providing contextual metadata in a useful way

- Best through example workflows
 - automated, machine and human-readable
- Not easy to define (even with >1200 papers...)
 - partly because
 - analysis processes are complex
 - CMS data format supports a wide range of use cases
 - but also because we, as a community, have undervalued:
 - documentation
 - common tools
 - analysis code reuse.

3

Support to external analysts

- [Open data portal information](#)
- [CMS Open data guide](#)
- [CERN open data forum](#)
- [CMS open data workshops](#)



tidy up a bit to align with 2012 version struc...

add trigger analyzer

Update README.md (#102)

README.md

Physics Objects Extractor (PhysObjectExtractor) for 2015MiniAOD data

Description

The `PhysObjectExtractor` package is the heart of the POET repository. It contains a collection of `QAnalyzers` that extract information from different physics objects into a `ROOT` file called `output.root`. These have been written separately for clarity and can be executed modularly using the configuration file called `poet_cfg.py`.

We need a logo!!!

THE example code: **POET**

Physics Objects Extractor Tool

Put together by many people in the CMS open data group from various sources within CMS.

Covers Run1 AOD and Run2 MiniAOD
Not a negligible effort!




CMS Open Data workshop 2022!



Yearly CMS Open data workshops
Next: 11-14 July!!!



Watch on  YouTube



4

Challenges

Open data have value only when in use but usability does not come for free

A little interlude to what it takes to open science to happen

4.1

Roles & responsibilities



Open science - what it takes to make it happen



4.2

Best practices

To preserve the knowledge at the time of active analysis



Efforts are needed

Manage your code in versioned code repositories

Capture the steps in your analysis

Package your analysis environment in software containers



Document everything from the start

Define easily reusable workflows

Use continuous, automated testing

Best practices require time but they will pay off:

for the individual, for the group, and eventually, for open science!

5

Plans

and some topics for discussion...



Approximate plans

	2023				2024			
	Q1	Q2	Q3	Q4	Q1	Q2	Q3	Q4
2013 Heavy-ion data release preparations	High	High	High	Low	Low	Low	Low	Low
2023 CMS Open data workshop preparations	Low	Medium	Medium	Low	Low	Low	Low	Low
2016 pp data release preparations (includes NanoAOD)	Low	Medium	High	Medium	Low	Low	Low	Low
2024 CMS open data workshop preparations (examples with NanoAOD)	Low	Low	Low	High	High	High	High	Low
Update CMS open data guide	Low	Low	Low	Medium	Medium	Medium	Medium	Low
2017 pp data release preparations (to be discussed and approved)	Low	Low	Low	Low	Low	Low	Medium	High
Improve metadata for workflows	Low	Low	High	High	High	High	Low	Low
Benchmarking workflows on public cloud resources	Low	Low	Low	Medium	Medium	Medium	Medium	Medium



Topics for discussion

Person-power

**Not everyone
agrees on open
data**

**Usage
patterns:
our
expectations vs
reality**

**Value to OD
efforts by big
OD/OS
projects?**

**Goal of open
research data:
research
(by others),
not outreach
(by us)**



Thank you!

Any **questions** ?

And thanks to [SlidesCarnival](#) for this free presentation template