# Using CMS Open Data in research — challenges and directions

## CHEP 2021

Kati Lassila-Perini[1], Clemens Lange[2], Edgar Carrera Jarrin[3], and Mathew Bellis[4].

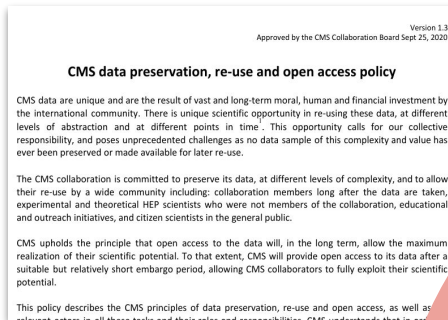[1]Helsinki Institute of Physics, Finland
[2]CERN, Switzerland
[3]Universidad San Francisco de Quito, Ecuador
[4]Siena College, USA

May 20, 2021

# Introduction

**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 levels of abstraction and at different points in time. 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.

This policy describes the CMS principles of data preservation, re-use and open access, as well as relevant actors in all these tasks and their roles and responsibilities. CMS understands that in ...

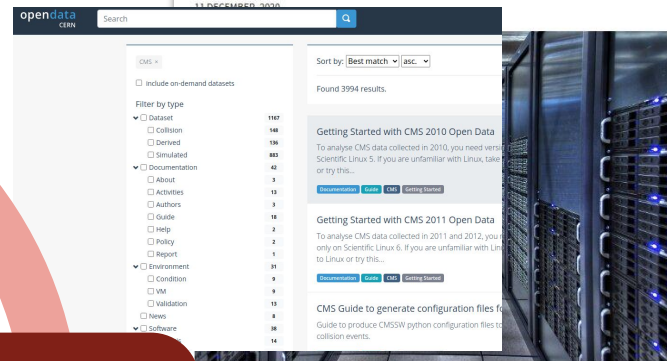## CMS Open Data releases

- Leading CERN's and LHC's efforts in open science

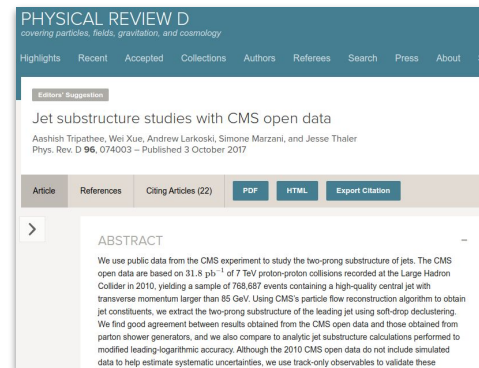- > 2PB of data released since 2014 using the CERN Open Data Portal (CODP)

**CERN announces new open data policy in support of open science**

A new open data policy for scientific experiments at the Large Hadron Collider (LHC) will make scientific research more reproducible, accessible, and collaborative

11 DECEMBER 2020

## Feedback

- CERN Open Data Forum
- From comments in published articles
- From support email
- From Informal discussions

## Usage in Research*

Topics include: standard model (SM) novel studies, re-measurements, searches beyond the SM, new methods and techniques (e.g., machine learning algos)
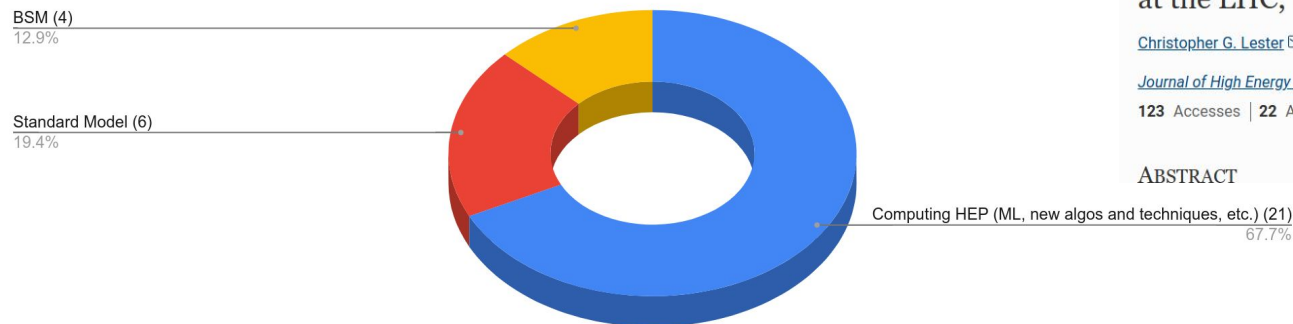
**PHYSICAL REVIEW D**
*covering particles, fields, gravitation, and cosmology*

Editors' Suggestion

Jet substructure studies with CMS open data

Aashish Tripathee, Wei Xue, Andrew Larkoski, Simone Marzani, and Jesse Thaler
Phys. Rev. D **96**, 074003 – Published 3 October 2017

ABSTRACT

We use public data from the CMS experiment to study the two-prong substructure of jets. The CMS open data are based on $31.8~pb^{-1}$ of 7 TeV proton-proton collisions recorded at the CMS Collider in 2010, yielding a sample of 768,687 events containing a high-quality central jet with transverse momentum larger than 85 GeV. Using CMS's particle flow reconstruction algorithm to obtain jet constituents, we extract the two-prong substructure of the leading jet using soft-drop declustering. We find good agreement between results obtained from the CMS open data and those obtained from parton shower generators, and we also compare to analytic jet substructure calculations performed to modified leading-logarithmic accuracy. Although the 2010 CMS open data do not include simulated data to help estimate systematic uncertainties, we use track-only observables to validate these

*Not an exact search link but a reference

# Articles using CMS open data*

| AREA | NUMBER OF ARTICLES (as of may 13, 2021) | | | |
|---|---|---|---|---|
| | Published in a journal | Arxiv only or submitted to a journal | Contribution to conferences | TOTAL |
| Computing HEP (ML, new algos and techniques, etc.) | 8 | 5 | 8 | 21 |
| Standard Model | 5 | | 1 | 6 |
| BSM | 2 | 2 | | 4 |
| TOTAL | 15 | 7 | 9 | 31 |

**ARTICLES USING CMS OPEN DATA**



BSM (4)
12.9%

Standard Model (6)
19.4%

Computing HEP (ML, new algos and techniques, etc.) (21)
67.7%

Journal of Instrumentation

**Opportunities and challenges of Standard Model production cross section measurements in proton-proton collisions at $\sqrt{s}$=8 TeV using CMS Open Data**

A. Apyan[1], W. Cuozzo[2], M. Klute[2], Y. Saito[2], M. Schott[2,3] and B. Sintayehu[2]

PHYSICAL REVIEW D
*covering particles, fields, gravitation, and cosmology*

Highlights    Recent    Accepted    Collections    Authors    Referees    Search    Press    About

Open Access

Searching in CMS open data for dimuon resonances with substantial transverse momentum

Cari Ces
Phys. Re    SpringerLink

Regular Article - Experimental Physics | Open Access | Published: 16 December 2019

**Testing non-standard sources of parity violation in jets at the LHC, trialled with CMS Open Data**

Christopher G. Lester ✉ & Matthias Schott

*Journal of High Energy Physics* **2019**, Article number: 120 (2019) | Cite this article

**123** Accesses | **22** Altmetric | Metrics

ABSTRACT

* Link is not exact but just a reference for easier search.

# CMS open data

## Release Policy

- Some embargo time and restrictions apply.
- Start of release after a few years of the end of data taking
- Most of Run 1 data released

## Data Format

- Analysis Object Data (AOD) format
- Based on ROOT and CMSSW
- Research quality
- Slimmer miniAOD and nanoAOD used and foreseen in Run 2.
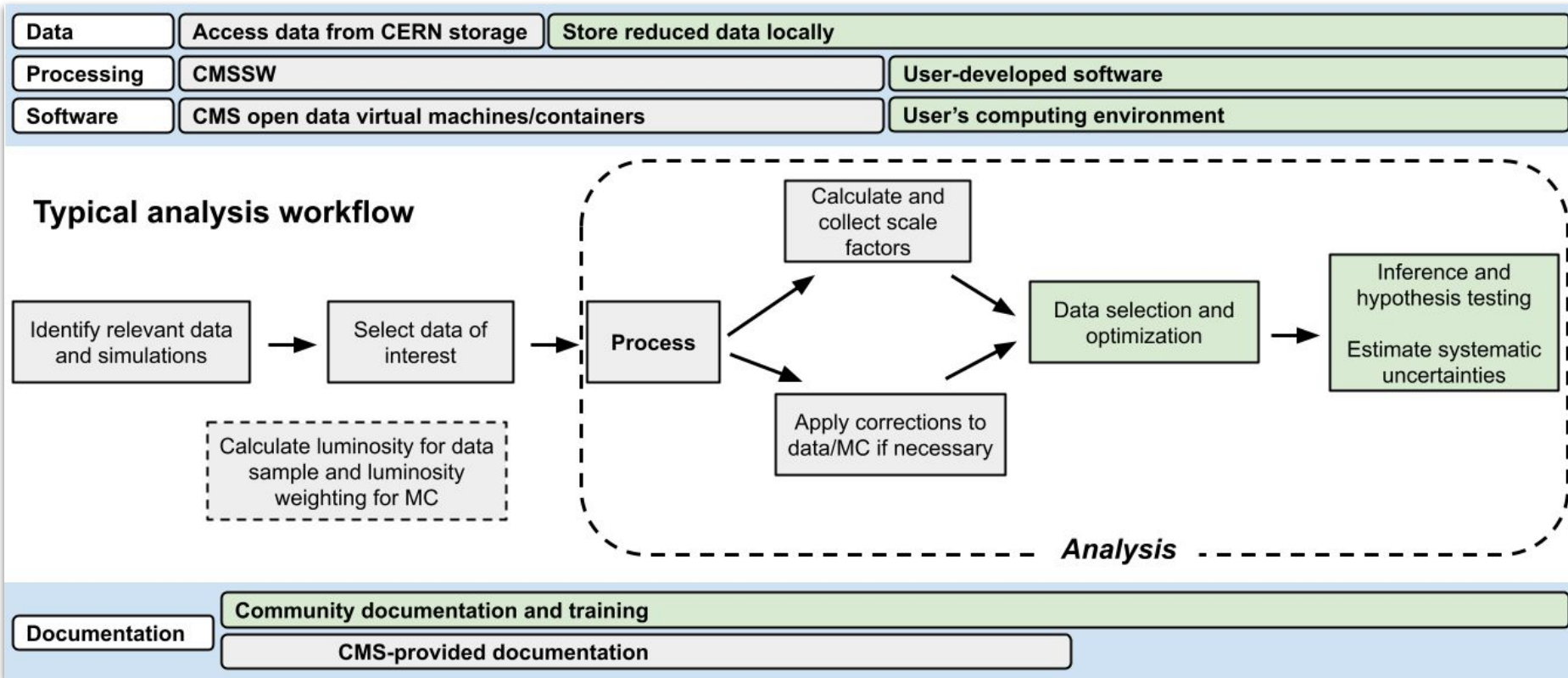
## Data products
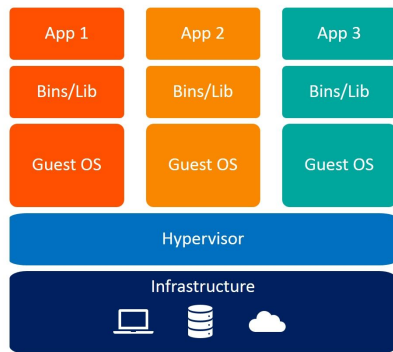
- Collision Data
- Simulated Data (MC)

## Software and Associated products

- CMSSW
- Data Quality
- Conditions database (alignment, calibration, etc)
- Luminosity information
- Examples (some with automated workflows) and topical guide pages
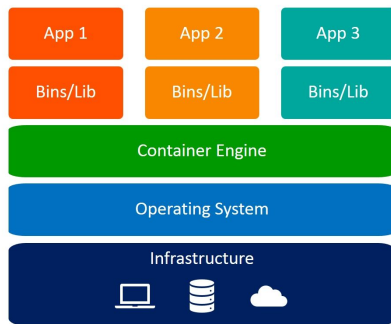
# Using CMS open data

| Data | Access data from CERN storage | Store reduced data locally | |
|---|---|---|---|
| Processing | CMSSW | | User-developed software |
| Software | CMS open data virtual machines/containers | | User's computing environment |

**Typical analysis workflow**

Identify relevant data and simulations → Select data of interest → **Process**

Calculate and collect scale factors

Apply corrections to data/MC if necessary

Data selection and optimization

Inference and hypothesis testing

Estimate systematic uncertainties

Calculate luminosity for data sample and luminosity weighting for MC

*Analysis*

| Documentation | Community documentation and training |
|---|---|
| | CMS-provided documentation |

# Using CMS open data



Virtual Machines

Containers

- Based on CernVM
- Quick first access
- CMSSW dependencies accessed through CVMFS
- Challenging to scale up for full analysis

- Different kind of images provided
- CVMFS may be dispensed with
- Layered file system makes it easier to maintain and preserve (Gitlab CI).
- Easier to use in batch compute systems.

Data Access: **XRootD** (remote) or http/XRootD (download).

VM or Docker container (CMSSW)

Condition Data (alignment, calibration, etc.) through CVMFS.

# User feedback and challenges

## Data Complexity

- Objects defined in CMSSW (C++) classes (AOD)
- Multiple definitions of physics objects
- Pile up
- Selection efficiencies, fake rates, calibrations, corrections
- Triggers, datasets, duplicate events
- Info overload and superfluous data

## Software Complexity

- Complexity of CMSSW (object properties are C++ classes)
- Difficult to navigate
- ROOT structures
- Difficulty to deal with legacy versions
- Procedures very analysis-specific

## Documentation And Extended Examples

## Scalability

- Order of TB datasets
- Batch/parallel needed
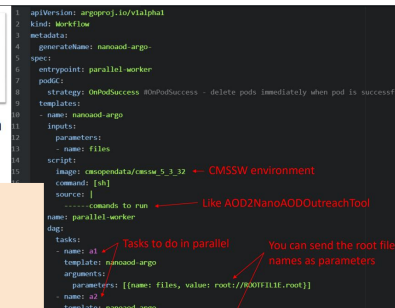- Slow development cycle

## Long-term usability

- CMS Run 1 open data fully dependent on CMSSW-compatible environment
- VMs and containers not completely independent of computing progress: require maintenance
- nanoAOD format (slimmer and better for long term) not yet available for Run 1.

# Measures to improve usability



- Easy way to manage a workflow inside Kubernetes clusters

## Preserved workflows as examples

- Workflows linked from CODP
- Different levels of complexity
- Github and/or GitLab CI/CD
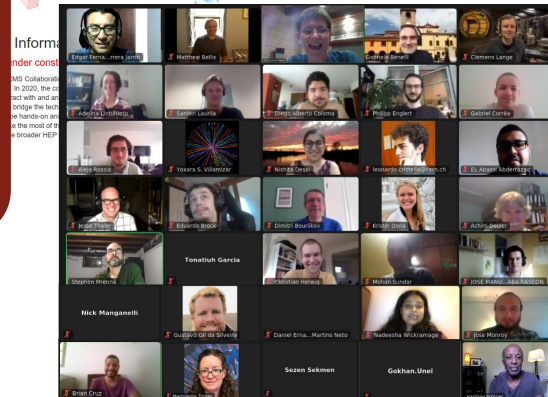- Usage of Kubernetes/minikube
- Some workflows executed in REANA

## Data analysis in the cloud

- Usage of scalable commercial clusters
- Simplify and test workflows on Kubernetes cloud clusters
- Test cloud local data download (via XRootD) and local cached CVMFS server

CMS Open Data Workshop 2021

CERN (Virtual)

Online (link)

Jul 19-22, 2021

2:30 pm - 6:40 pm (CET)

**Instructors:** Matt Bellis, Edgar Carrera Jarrin, Julie Hogan, Clemens Lange, Kati Lassila-Perini

**Helpers:** helper one, helper two

## Documentation and Training

- Twiki pages, CODP records, Github code available but difficult
- Analysis examples in CODP hard to navigate
- CMS Open Data Guide
- Since 2020, regular hands-on workshops and training events

### CMS Open Data Guide

⚠️ **Warning**

This page is under construction.

Welcome to the official guide for CMS open data. This page is still under construction. We appreciate your feedback and/or your help building this guide.

**How to use this site**

There are three main tabs to help you navigate the site. It starts with the **Computing Tools** most likely needed to deal with CMS open data. Then, there is a little review of **CMSSW**, which is the software used by CMS. Finally the **Analysis** section guides you through the different steps (in the most general order) that you need to follow for performing a particle physics analysis with CMS open data.

# Summary and Outlook

- CMS has spearheaded CERN's open data efforts
- CMS open data have been successfully used in original scientific research (jet physics, sm measurements, new methods and algorithms)
- Usage of these data has opened opportunities for deeper collaboration between theorists and experimentalists
- Limitations and challenges have been identified from user feedback and self assessment:
    - Information extraction from data, procedures, documentation and examples.
- Measures (within limited person-power capabilities) are being taken in order to improve usability:
    - Better (automated) implementation of workflows, cloud computing, CMS Open Data Guide and training events.
- CMS will maintain its commitment to open data and open science
- Next CMS Open Data Workshop, July 19-22, 2021. Registrations open: https://indico.cern.ch/e/CmsODW2021