# LHC experiments and their Open Data

## LHCP 2021

Edgar F. Carrera
On behalf of the ALICE, ATLAS, CMS, and LHCb Collaborations
Universidad San Francisco de Quito, Ecuador

UNIVERSIDAD
SAN FRANCISCO

June 9, 2021

# CERN Open Data Policy for LHC experiments

**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

Data storage solutions at the CERN data centre (Image: CERN)

Geneva, 11 December 2020. The four main LHC collaborations (ALICE, ATLAS, CMS and LHCb) have unanimously endorsed a new open data policy for scientific experiments at the Large Hadron Collider (LHC), which was presented to the CERN Council today. The policy commits to publicly releasing so-called level 3 scientific data, the type required to make scientific studies, collected by the LHC experiments. Data will start to be released approximately five years after collection, and the aim is for the full dataset to be publicly available by the close of the experiment concerned. The policy addresses the growing movement of open science, which aims to make scientific research more reproducible, accessible, and collaborative.

- CERN committed to openness and preservation for a long time
- Policy relates to data collected by LHC experiments
- Endorsed by all LHC experiments
- Different levels of abstraction:

- Level 1 (published results and numerical information (likelihoods))
- Level 2 (outreach and education)
  - simplified format
- Level 3 (reconstructed data)
  - research quality
  - latency/embargo periods apply
  - restrictions apply
- Level 4 (raw data)
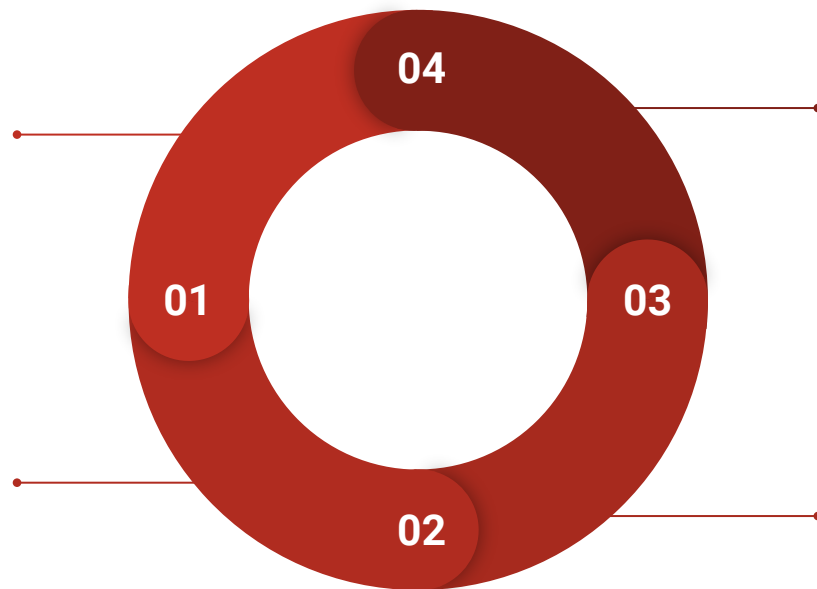
# Level 3 (reconstructed data)

- All LHC experiments share common strategies for Level 1, 2, and 4.
- Level 3 strategies differ, within a common ground.
- LHC experiments will release:

**Calibrated reconstructed data**
- Including provenance data, simulated samples, example workflows and documentation

**Periodically**
- With an appropriate latency period.
- Data releases within 5 years of the conclusion of run period.
- Timeline defined by experiments but full datasets will be made public at the end of the collaborations
- Some data may be withheld if active analyses ongoing

**And may offer association programs**
- External authors may not get access to cumulative knowledge within the collaborations
- Association programs may be provided

**Using the CERN Open Data Portal (CODP)**

01  02  03  04

# CERN Open Data Portal (CODP)



https://opendata.cern.ch/

- Access point to data
- Also software and documentation

- Products shared under open licenses
- Issued with a Digital Object Identifier (DOI) to make them citable

# ALICE open data for RESEARCH

## Release Policy:
- Public data releases expected periodically
- Needed appropriate latency period to allow (**true for all experiments**):
  - thorough understanding of the data
  - reconstruction and calibrations
  - scientific exploitation of the data
- Aim to commence data releases within five years of the conclusion of the run period
- Size of the released datasets commensurate with the amount of data collected

## Data Products and computing environment:
- Conversion of Run 1 and Run 2 ESD and AOD into new AOD format (based on ALICE O2 Project)
- Also derived nanoAOD considered
- Real (HI) and simulated datasets in new AOD format expected
- Related software via VM and/or containers

## Current status and timeline:
- 5% (7%) of Pb-Pb (pp) 2010 ESD datasets released
- Conversion into new AOD format expected by the end of 2021
- Expect to populate CODP with new AOD from 2022

# ALICE open data for OUTREACH AND EDUCATION



- ESD files in COPD have also been used for outreach and education
- new version of the ALICE - strange particles masterclass exercise is web-based: https://alice-web-masterclass.web.cern.ch/ (all data sets needed are provided by the web server)

# ATLAS open data for Research

## Release Policy:
- Preparing periodic releases in alignment with new CERN-LHC policy
- Appropriate latency period needed to understand and scientifically exploit the data

## Data Products and computing environment:
- Full likelihoods released (level 1 data) for reinterpretation
- Real and simulated datasets in PHYSLITE format (calibrated objects and information to compute systematics)
- Associated software (containers)

## Current status and timeline:
- Only datasets for education purposes released so far.
- Association programs established (case-by-case for external proposals) for specific analyses
- Special datasets may also be approved for release



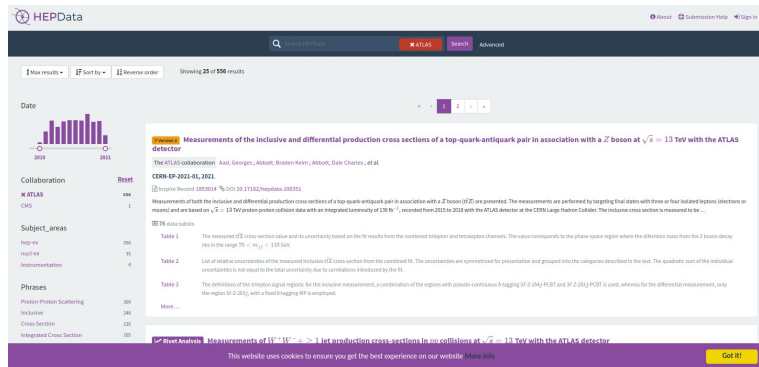*ATLAS* approach to releasing likelihoods for reinterpretations

ReINPS2021

Eric Schanet

on behalf of the *ATLAS* collaboration

February 15, 2021

Luxembourg National Research Fund

LMU LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN

ATLAS EXPERIMENT

**Publishing** likelihoods

# ATLAS open data for Outreach and Education



ATLAS Open Data provides open access to proton-proton collision data at the LHC for educational purposes. Designed in collaboration with students and teachers, ATLAS Open Data resources are ideal for high-school, undergraduate and postgraduate students – or even enthusiastic self-learners! So whether you have an hour or a semester, try your hand at analysing the 13 TeV proton-... Experiment. Create a simple histogram, write and visualise a particle-physics analysis directly from y... resources to take a deep-dive into the ATLAS analysis framework and re-discover the Higgs boson!

## Online Open Data Analysis

Explore ATLAS open datasets and physics analyses directly from your browser with the help of our cloud computing resources. These "Jupyter notebooks" allow you to easily interact with the data without downloading files or writing code. If you have a CERN computing account, you can also explore these notebooks through the SWAN
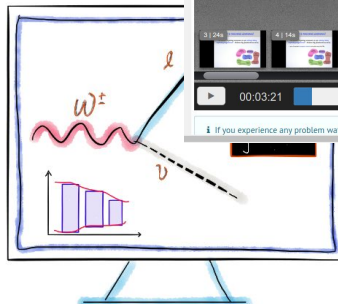
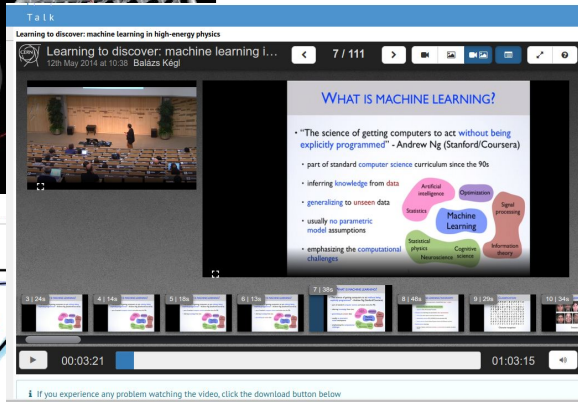A look inside & around the ATLAS detector

**Histogram Analyser: Real & Simulated Data**

m-analyser-02/

Perform real HEP analysis with your mouse

**Jupyter Notebooks**

# CMS open data for Research

## Release Policy:
- 50% after 6 years and
- 100% within 10 years of data taking
- Limit of 20% if data taking still planned
- 100% at the end of collaboration
- CB approval always needed ( which can modify the dates and sizes)

## Data Products and computing environment:
- Real and simulated datasets in AOD (Run 1) format
- MiniAOD and nanoAOD (Run 2 and beyond)
- VM and Docker containers with CMSSW software
- Trigger info and conditions database (align., calib, etc.) through CVMFS
- Data quality and luminosity information tools
- Additional level-1 (numerical info) through Inspire and HEPData.

## Current status and timeline:
- Most pp data from Run 1 released (2010-2012)
- First batch of HI Run 1 data released in 2020
- Few other datasets released for, e.g., machine learning studies (including Run 2 simulations)
- First Run 2 release expected by the end of 2021 or early 2022

## Usage:
- About 30 articles* produced (SM, BSM, computing)
- About 15 published in indexed journals

Journal of Instrumentation

Opportunities and challenges of Standard Model production cross section measurements in proton-proton collisions at √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]
Published 14 January 2020 • © 2020 IOP Publishing Ltd and Sissa Medialab
Journal of Instrumentation, Volume 15, January 2020
Citation A. Apyan et al 2020 JINST 15 P01009

Article PDF

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

Highlights   Recent   Accepted   Collections   Authors   Referees   Search   Press

Open Access

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

Cari Cesarotti, Yotam Soreq, Matthew J. Strassler, Jesse Thaler, and Wei Xue
Phys. Rev. D 100, 015021 – Published 16 July 2019

Article   References   Citing Articles (5)   PDF   HTML   Export Citation

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

Highlights   Recent   Accepted   Collections   Authors   Referees   Search   Press   About

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

Article   References   Citing Articles (22)   PDF   HTML   Export Citation

Springer Link

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

The Standard Model violates parity, but only by mechanisms which are invisible to Large Hadron Collider (LHC) experiments (on account of the lack of initial state polarisation or spin-sensitivity in the detectors). Nonetheless, new physical processes could potentially violate parity in ways which are detectable by those same experiments. If those sources of new physics occur only at LHC energies, they are untested by direct searches. We probe the feasibility of such measurements using approximately 0.2 fb⁻¹ of data which was recorded in 2012 by the CMS collaboration and made public within the CMS Open Data initiative. In particular, we test an inclusive three-jet event selection which is primarily sensitive to non-standard parity violating effects in quark-gluon interactions. Within our measurements, no significant deviation from the Standard Model is seen and no obvious experimental limitations have been found. We discuss other ways that searches for non-standard parity violation could

*link not exact but a reference

# CMS open data for OUTREACH AND EDUCATION

# LHCb open data for Research, outreach and education

**Release Policy:**
- 50% after 5 years and
- 100% after 10 years (from the end of LHC running period)
- Some restrictions may apply but 100% at the end of collaboration.

**Data Products and computing environment:**
- Real and simulated data in DST and microDST formats
- Associated software (containers)
- Calibration tools (at user's responsibility)

**Current status and timeline:**
- Small sample of simulated events released in 2020 for quantum-inspired machine learning techniques (led to paper pre-print)
- Otherwise, releases currently limited to outreach and education

Apr 2020

### Quantum-inspired Machine Learning on high-energy physics data

Marco Trenti,[1] Lorenzo Sestini,[2] Alessio Gianelle,[2] Davide Zuliani,[1,2]
Timo Felser,[1,2,3] Donatella Lucchesi,[1,2] and Simone Montangero[1,2]

[1] *Dipartimento di Fisica e Astronomia "G. Galilei", Università di Padova, I-35131 Padova, Italy*
[2] *INFN, Sezione di Padova, I-35131 Padova, Italy.*
[3] *Theoretische Physik, Universität des Saarlandes, D-66123 Saarbrücken, Germany.*
(Dated: April 30, 2020)

One of the most challenging big data problems in high energy physics is the analysis and classification of the data produced by the Large Hadron Collider at CERN. Recently, machine learning techniques have been employed to tackle such challenges, which, despite being very effective, rely on classification schemes that are hard to interpret. Here, we introduce and apply a quantum-inspired machine learning technique and, exploiting tree tensor networks, we show how to efficiently classify b-jet events in proton-proton collisions at LHCb and to interpret the classification results. In particular, we show how to select important features and adapt the network geometry based on information acquired in the learning process. Moreover, the tree tensor network can be adapted for optimal precision or fast response in time without the need of repeating the learning process. This paves the way to high-frequency real-time applications as needed for current and future LHC event classification to trigger events at the tens of MHz scale.

# Summary

- LHC experiments at CERN have recently (Dec, 2020) announced a new, revitalized policy on open data
- The new policy puts emphasis on the release of Level 3 (research quality) scientific data
- LHC experiments have agreed to start releasing Level 3 data as soon as 5 years after data taking periods end, and to make 100% of such data available after the end of each collaboration.
- While strategies differ and some restrictions are imposed, all the LHC experiments are committed to open data for research.
- All LHC experiments have made important contributions to outreach and education using their open data.

# Backup slides

# Summary of common L1, L2 and L4 strategies

- **Level 1 (published results):**
  - Peer-reviewed publications available with Open access
  - Make public additional information and data at the time of publication (in portals such as HEPData)
  - The data made available may include simplified or full binned likelihoods, as well as unbinned likelihoods based on datasets of event-level observables extracted by the analyses.
  - Reinterpretation of published results possible through analysis preservation and direct collaboration with external researchers.
- **Level 2 (outreach and education):**
  - Dedicated subsets of data are used, selected and formatted to provide rich samples to maximise their educational impact, and to facilitate the easy use of the data.
  - Data are released with a schedule and scope determined by each experiment.
  - Data are provided in simplified, portable and self-contained formats suitable for educational and public understanding purposes.
  - Not intended or adequate for the publication of scientific results.
  - Lightweight environments to allow the easy exploration of these data may be provided.
  - Data accessible through the CERN Open Data Portal.
- **Level 4 (raw data):**
  - Not intended for release
  - Small samples may be approved for release if useful for, e.g., the development of new reconstruction algorithms.