

CERN Open Data: Policy to Implementation

J. Boyd obo the CERN OD WG

DPHEP workshop @ CERN

3/10/2024

Many related talks at this workshop:

ALICE: David Chinellato

ATLAS: Zach Marshall

CMS: Julie Hogan

LHCb: Dillon Fitzgerald

LEP: Jacopo Fanini , Dietrich Liko, Matthias Schroeder, Ulrich Schwickerath

OD portal: Pablo Saiz

Introduction

Some brief history:

Open Data Policy Working Group for the LHC expts

- Initiated by Eckhard Elsen (DRC at the time) in Feb 2020
- Working Group consists of 2 representatives each from the large LHC experiments (1 from TOTEM), and a representative from IT and SIS
- WG chaired by Jamie Boyd
- Mandated to draft a policy document relevant to the LHC experiments Open Data



- Policy endorsed by large LHC experiments by end of 2020
- Since then Joachim Mnich (DRC) has asked me to work on expanding the set of experiments that sign up to the policy
 - **First** “small” LHC experiments (TOTEM, LHCf, MoEDAL, FASER, SND@LHC)
 - Now all of these experiments have endorsed the policy!
 - **Second**, non-LHC experiments
 - Discussions with these experiments ongoing....

The ODWG is now a standing WG, with once yearly meetings to discuss the status of the policy implementation.

Policy Strategy

- The policy has been broken down into the 4 levels of data as defined in the DHEP study on data preservation:
 - Level 1 – Scientific publications, and associated additional data
 - Level 2 – Data useful for Education and Outreach
 - Level 3 – Reconstruction level data useful for general physics analysis
 - Level 4 – RAW data
- All large-LHC experiments already released data for L1 and L2 in broadly similar ways, and all agree that L4 is not practically useful
 - The WG discussion therefore strongly focussed on the policy for L3 data
- Decided to write a short general policy document (to be made public) outlining the general principles, and in addition an internal document outlining the implementation of the policy by each experiment
 - Implementation strategy for each experiment cannot be changed without discussing in the WG

Policy document : <https://cds.cern.ch/record/2745133>

Implementation document : <https://cds.cern.ch/record/2745081>



Public Policy Document

- 2 page document:
- Introduction:
 - Motivation and scope
 - Outline DHEP levels of data
- L1: Published results policy:
 - Publish results in Open Access journals
 - Provide additional information through HepData / RIVET etc..
 - No change in current policy
- L2: Outreach and Education
 - Provide rich data samples tailored for outreach/education (easy to use)
- L3: Reconstructed data
 - See next slides...
 - Real change to Collaborations policy
- L4: Raw Data
 - Not useful for external people

The CERN Open Data Policy reflects values that have been enshrined in the CERN Convention for more than sixty years that were reaffirmed in the European Strategy for Particle Physics (2020)¹, and aims to empower the LHC experiments to adopt a consistent approach towards the openness and preservation of experimental data. Making data available responsibly (applying FAIR standards²), at different levels of abstraction and at different points in time, allows the maximum realisation of their scientific potential and the fulfillment of the collective moral and fiduciary responsibility to member states and the broader global scientific community. CERN understands that in order to optimise reuse opportunities, immediate and continued resources are needed. The level of support that CERN and the experiments will be able to provide to external users will depend on available resources.

This policy relates to the data collected by the LHC experiments, for the main physics programme of the LHC — high-energy proton–proton and heavy-ion collision data. The foreseen use cases of the Open Data include reinterpretation and reanalysis of physics results, education and outreach, data analysis for technical and algorithmic developments and physics research. The Open Data will be released through the CERN Open Data Portal which will be supported by CERN for the lifetime of the data. The data will be tailored to the different uses, and will be made available in formats defined by each experiment that afford a range of opportunities for long-term use, reuse and preservation. In general, four levels of complexity of HEP data have been identified by the Data Preservation and Long Term Analysis in High Energy Physics (DPHEP) Study Group³, which serve varying audiences and imply a diversity of openness solutions and practices.

Published Results (Level 1) Policy: Peer-reviewed publications represent the primary scientific output from the experiments. In compliance with the CERN Open Access Policy, all such publications are available with Open Access, and so are available to the public. To maximise the scientific value of their publications, the experiments will make public additional information and data at the time of publication, stored in collaboration with portals such as HEPData,⁴ with selection routines stored in specialised tools. 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 is also made possible through analysis preservation and direct collaboration with external researchers.

Outreach and Education (Level 2) Policy: For the purposes of education and outreach, 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. These data are released with a schedule and scope determined by each experiment. The data are provided in simplified, portable and self-contained formats suitable for educational and public understanding purposes; but are not intended nor adequate for the publication of scientific results. Lightweight environments to allow the easy

also be provided. CERN experiments will make data of such high level of e through the CERN Open Data Portal.⁵

Policy: The LHC experiments will release calibrated reconstructed data for algorithmic, performance and physics studies. The release of these provenance metadata, and by a concurrent release of appropriate ware, reproducible example analysis workflows, and documentation. nts that are compatible with the data and software will be made ovided will be sufficient to allow high-quality analysis of the data llication of the main correction factors and corresponding systematic ations, detector reconstruction and identification. A limited level of vel 3 Open Data will be provided on a best-effort basis by the

periodically, following an appropriate latency period to allow thorough re reconstruction and calibrations, as well as to allow time for the data by the collaboration. The size of the released datasets will be amount of data collected of similar type, with the aim to commence s of the conclusion of the run period. Data may be withheld by an analyses ongoing. Full datasets will be made available at the close of the

¹ European Strategy Group (2020), '2020 Update of the European Strategy for Particle Physics'.

² FAIR Guiding Principles for scientific data management and stewardship. Available at: <https://www.go-fair.org/fair-principles/>.

³ Data management plans are defined by the LHC experiments to address the long-term preservation of internal data products. See: Akopov et al., Status report of the DPHEP Study Group: Towards a global effort for sustainable data preservation in high energy physics. arXiv preprint arXiv:1205.4667 (2012).

⁴ Repository for publication-related High-Energy Physics data: <http://www.hepdata.net>.

The data will be released from the CERN Open Data Portal under the Creative Commons CC0 waiver, and will be identified with persistent data identifiers, and the data must be cited through these identifiers. Similarly, appropriate acknowledgements of the experiment(s) should be included in publications released using such data, and the publications made clearly distinguishable from those released by the collaboration. Any scientific claims in such publications are the responsibility of their authors and not of the experiments. It is expected that scientific results released using Open Data follow best scientific practices. The experiments may impose rules related to the use of the data by members of their respective collaborations.

External authors should be aware that they will not have access to the vast amount of tacit knowledge built up within the LHC collaborations over the decades of design, construction and operation of the experimental apparatus. To allow external scientists to fully benefit from all the data, knowledge and tools, the collaborations may offer appropriate association programmes.

Raw Data (Level 4) Policy: It is not practically possible to make the full raw data-set from the LHC experiments usable in a meaningful way outside the collaborations. This is due to the complexity of the data, metadata and software, the required knowledge of the detector itself and the methods of reconstruction, the extensive computing resources necessary and the access issues for the enormous volume of data stored in archival media. It should be noted that, for these reasons, general direct access to the raw data is not even available to individuals within the collaboration, and that instead the production of reconstructed data (i.e. Level-3 data) is performed centrally. Access to representative subsets of raw data—useful for example for studies in the machine learning domain and beyond—can be released together with Level-3 formats, at the discretion of each experiment.

⁵ CERN Open Data portal: <http://opendata.cern.ch>.

L3 data – discussion (1)

- Any rules relating to publically releasing L3 data need to be approved by each experiments Collaboration Board
- Generally tried to find a good balance between:

1. Making data openly available

And

2. Protecting the collaborations:

- Avoiding collaboration members publishing with open data rather than a Collaboration paper
- Having to deal with wild claims made by external analysts with their data
- Not taking too much resources (both human and computing) from the Collaboration
- After much discussion converged on limiting of the amount of data released after a given time as a tool to make it unattractive for collaboration members to publish with Open Data
 - Here a common approach across experiments is needed since e.g. ATLAS rules could effect a CMS collaborator's behaviour in this regard
 - Exact fraction of data released after what latency experiment specific, but general principle common across experiments

L3 data – discussion (2)

- Summary of each experiments Latency shown below:

	ALICE	ATLAS	CMS	<u>LHCb</u>
Fraction of data released in: 5 yrs (6 yrs for CMS)	10%	25% (but limiting to <20% of the total data at that time)	50% (but limiting to <20% of the total data at that time)	50%
Fraction of data released in: 10 yrs	50%	50% (but limiting to <20% of the total data at that time)	100% (but limiting to <20% of the total data at that time)	100%
End-of-Collaboration	100%	100%	100%	100%

L3 data – discussion (3)

- We should only release L3 Open Data if it can be used for high quality science
 - Release data and simulated samples
 - Release data with best calibrations available
 - Release sufficient information to allow main systematic uncertainties to be applied
 - Release analysis s/w to allow efficient analysis of data
- Data expected to be useable not only for particle physics, but also computing, algorithmic development, big data studies etc..
- Publications using Open Data should:
 - Have appropriate acknowledgements
 - Be clearly identifiable from Collaboration papers
 - Follow best scientific practices
- All experiments to use CERN Open Data Portal to house the released data

Implementation: Data Volumes

- Implementation document showed the expected media resource needs:

	ALICE (TB)	ATLAS (TB)	CMS (TB)	LHCb (TB)	SUM (TB)	Cumulative (PB)
2022	15	0	1089	200	1304	1.3
2023	20	150	872	4600	5642	6.9
2024	105	0	1436	0	1541	8.5
2025	105	0	1768	0	1873	10.4
Total	245	150	5165	4800	10360	

Table updated since, but broadly consistent with previous estimates.

- Currently ~5.5PB of LHC OD stored in the portal
- CERN IT kindly agreed to provide these resources for first 5 years (until end of 2025) - then to be re-discussed
- IT position to develop integrated tape back-end (save costs on storage resources in the long-term)
 - Work ongoing, hopefully to be deployed next year



Level-3 data so far released

7/24

ATLAS releases 65 TB of open data for research

Explore over 75 billion LHC collision events — from home

News ATLAS

4/24

CMS releases 13 TeV proton collision data from 2016

CMS releases 13 TeV proton collision data from 2016

News CMS

12/23

LHCb releases the entire Run I dataset

Today the LHCb collaboration completes the release of the data collected throughout the Run I of the Large Hadron Collider at CERN.

News LHCb

9/23

CMS completes Run-1 heavy ion open data collection

New release of simulations, proton-lead collision data, and proton reference data.

News CMS

12/22

LHCb releases first set of data to the public

The LHCb collaboration has released data from Run 1 of the LHC to the public for the first time, allowing research to be conducted by anyone in the world.

News LHCb

12/21

CMS completes the release of its entire Run-1 proton-proton data

All proton-proton data collected by the CMS experiment during LHC Run-1 (2010-2012) are now available through the CERN Open Data Portal.

News CMS

12/20

First CMS open data from LHC Run 2 released

As the experiments at the Large Hadron Collider (LHC) brace for the start of Run 3 of the accelerator's programme in 2022, the CMS collaboration has released a new research-quality open data recorded by the CMS detector in 2015, the first year of Run 2. The new datasets are now available on the CERN Open Data portal.

News CMS

CMS releases heavy-ion data from 2010 and 2011

CMS releases heavy-ion data from 2010 and 2011

News CMS



Level-3 data so far released

7/24

ATLAS releases 65 TB of open data for research
Explore over 75 billion LHC collision events — from home

News ATLAS

4/24

CMS releases 13 TeV proton collision data from 2016
CMS releases 13 TeV proton collision data from 2016

News CMS

LHCb releases the entire Run-1 dataset

Today the LHCb collaboration co

News LHCb

CMS completes Run-1 heavy-ion data release
New release of simulations, proto

News CMS

LHCb releases first set of data
The LHCb collaboration has rele

News LHCb

CMS completes the release of Run-1 data
All proton-proton data collected b

News CMS

12/21

First CMS open data from LHC Run 2 released

As the experiments at the Large Hadron Collider (LHC) brace for the start of Run 3 of the accelerator's programme in 2022, the CMS collaboration has released a new research-quality open data recorded by the CMS detector in 2015, the first year of Run 2. The new datasets are now available on the CERN Open Data portal.

News CMS

12/20

CMS releases heavy-ion data from 2010 and 2011
CMS releases heavy-ion data from 2010 and 2011

News CMS

No ALICE data release on Open Data portal yet, but will happen VERY soon

A lot of work has been going in ALICE to prepare for this:

- Developed new format to release data in
 - Needed for efficient long-term release of data
- Converted planned datasets to new format
- Established a new software framework for analysing converted data
- Validation of new samples now done
- In final stages of metadata preparation and data transfer to portal, for:
 - pp, p-Pb and Pb-Pb Run 1 data samples



Level 3 Open Data: LHCb

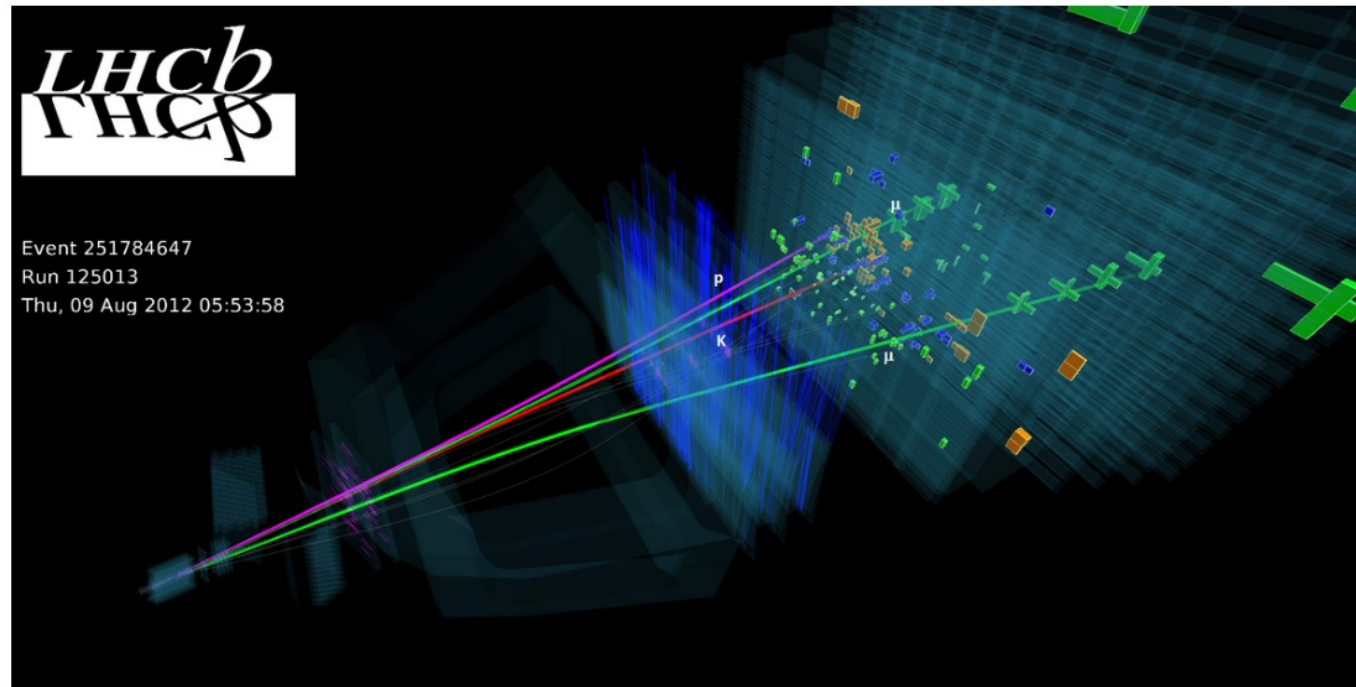
20/12/23

LHCb releases the entire Run I dataset

2023-12-20 by LHCb Collaboration

News

Today the LHCb collaboration completes the release of the data collected throughout the Run I of the Large Hadron Collider at CERN. The sample made available amounts to approximately 800 terabytes (TB) of data. These data, collected by the LHCb experiment in 2011 and 2012, contains information obtained from [proton-proton](#) collisions. The format made available provides pre-filtered data, suitable for a wide range of physics studies. The image below displays an [event](#) recorded during 2012.





Level 3 Open Data: CMS

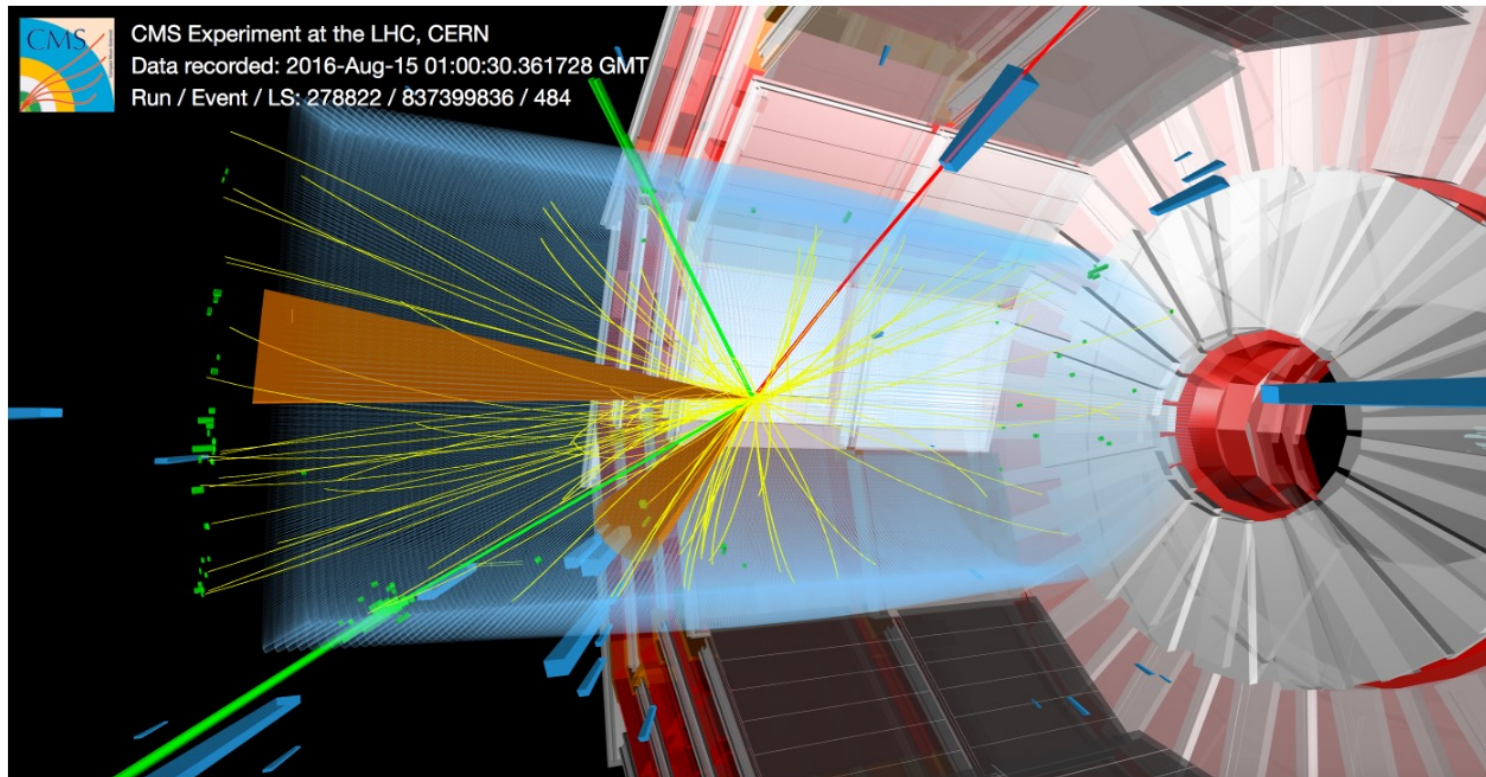
2/4/24

CMS releases 13 TeV proton collision data from 2016

2024-04-02 by CMS Collaboration

News

The CMS experiment at CERN is proud to announce the first release of 13 TeV proton-proton collision data collected in 2016. Over 70 TB of 13 TeV collision data and 830 TB of corresponding simulations are now accessible to the global scientific community and enthusiasts alike through the [CERN Open Data Portal](#).



A candidate event in which a top quark is produced in association with a Z boson. Photo: [CERN Document Server](#)



Level 3 Open Data: ATLAS

1/7/24

ATLAS releases 65 TB of open data for research

2024-07-01 by ATLAS Collaboration

News

Explore over 7 billion LHC collision events – from home

The ATLAS Experiment at CERN has made two years' worth of scientific data available to the public for research purposes. The data include recordings of proton–proton collisions from the Large Hadron Collider (LHC) at a collision energy of 13 TeV. This is the first time that ATLAS has released data on this scale, and it marks a significant milestone in terms of public access and utilisation of LHC data.

“Open access is a core value of CERN and the ATLAS Collaboration,” says Andreas Hoecker, ATLAS Spokesperson. “Since its beginning, ATLAS has strived to make its results fully accessible and reusable through open access archives such as arXiv and HepData. ATLAS has routinely released open data for educational purposes. Now, we’re taking it one step further — inviting everyone to explore the data that led to our discoveries.”

Released under the Creative Commons CC0 waiver, ATLAS has made public all the data collected by the experiment during the 2015 and 2016 proton–proton operation of the LHC. This is approximately 65 TB of data, representing over 7 billion LHC collision events. In addition, ATLAS has released 2 billion events of simulated “Monte Carlo” data, which are essential for carrying out a physics analysis.

External researchers, in particular, are encouraged to explore the ATLAS open data. “Along with the data, we have provided comprehensive documentation on several of our analyses, guiding users through our process step-by-step,” says Zach Marshall, ATLAS Computing co-Coordinator. “These guides provide first-hand experience of working on a real ATLAS result, allowing anyone to test our tools, and evaluate the systematic uncertainties associated with the result for themselves.”

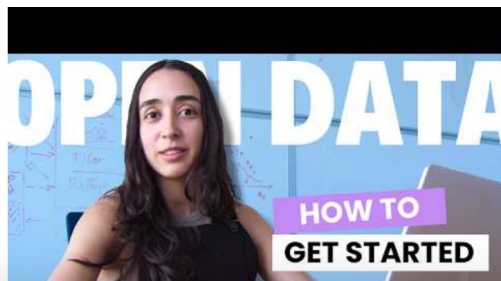
ATLAS traditionally collaborates with non-ATLAS scientists through short-term associations, granting them full access to ATLAS data, internal tools, and information. Through the open data, ATLAS researchers hope to further nurture this dialogue and collaboration. “In particular,” adds Zach, “we’d like to encourage phenomenologists and also computer scientists to explore our datasets, instead of relying on mock-ups.”

Today’s release builds upon previous open data releases for educational use (in 2016 and 2020). “All of our open data releases are now available through the ATLAS open data website,” says Dilia Portillo, ATLAS Outreach and Education co-Coordinator. “The website includes multi-level documentation, video tutorials and online tools aimed at the full-spectrum of users, from high school students to senior particle physics researchers. In addition, the software used to create the education-use open data has been released. This provides a seamless transition from the research open data to all the tutorials for outreach and education, including newly updated Higgs-boson discovery documentation. With a bit of time and dedication, you can go from being a relative novice to carrying out your own analysis.”

The ATLAS open data website also serves as a hub for the community, which includes teachers, students, enthusiasts and, now, scientists. Anyone diving into the open data can also directly engage with ATLAS physicists, who are available to respond to user feedback and take suggestions.

This release marks the start of more to come, with ATLAS’ first release of lead-lead-nuclei collision data data up next. The ATLAS Collaboration, along with the other main LHC experiment collaborations, has committed to making all of its data publicly accessible after a certain time. Openness is deeply ingrained in the culture of high-energy physics, enabling greater accessibility, reproducibility and better science.

Begin your journey with ATLAS open data by following the tutorial below.



Evolving model

- At the time of writing the implementation document in 2020, it became clear that the disk space needs for LHCb L3 Open Data will become prohibitive after Run 1
 - Estimated storage space 10PB (45PB) for Run 2 (3) if using the same model as Run 1
- LHCb have therefore developed a new strategy based on the NtupleWizard for access to LHCb OD
 - Allows a user to skim the full data (stored on LHCb resources) using LHCb CPU resources to produce an ntuple tailored for their use case
 - Smart idea to trade off CPU for storage
 - At advanced stage of development, to be deployed next year
 - More details in later talk by Dillon Fitzgerald
- Current ALICE model will also not scale for Run 3 dataset (would be 2PB/year)
 - Plan to publish skimmed and derived dataset for Run 3 (to be released in 2030+)

Open Data monitoring

- In the context of the CERN Open Science effort a discussion has started about what metrics should be monitored related to CERN Open Data and its usage
 - This is an ongoing discussion, but some current thoughts are shown below
- Some Key Performance Indicators will be made public in the Open Science report, while others will be monitored internally to understand the usage at a more detailed level
- Discussions are ongoing but can think of two sets of KPIs:
 - Related to data stored and accessed
 - Quantitative KPIs can be extracted from the OD portal
 - Dashboard under development in IT
 - Related to publications using CERN OD
 - Quantitative KPIs can be extracted from Inspire (assuming papers reference OD DOIs as they should)
 - This can miss the use of OD which does not lead to publications e.g. for education purposes
 - More difficult to come-up with quantitative KPIs for this, but maybe able to have more qualitative information related to this from e.g. surveying users

Expanding the effort – small-LHC expts



Expanding OD across CERN experiments

- After the successful drafting of the OD policy for the large LHC experiments, I have been asked by the DRC to try to expand this policy to ultimately cover all experiments at CERN (where possible)
- Suggested to start with the small LHC experiment:
 - Some of these experiments are quite different from the large LHC experiments, which may affect the need/use of OD (e.g. emulsion based, passive detectors, only operating during special running conditions)
 - Datasets expected to be small compared to large LHC experiments but experiments have less (human) resources for OD effort

All small-LHC experiments have now agreed to this:

Experiment	Status	Comment
TOTEM	Policy endorsed	Operating since LHC start. Data taking coming to an end this year, OD way towards data preservation?
LHCf	Policy endorsed	Operating since LHC start. Only takes data in special runs. Last running during Run 3.
MoEDAL	Policy endorsed	In Run 1,2 only passive detectors. Installing electronic detectors during Run 3.
FASER	Policy endorsed	Only starting data taking in Run 3 (approved for duration of Run 3).
SND@LHC	Policy endorsed	Electronic detector, and (passive) emulsion detector data.

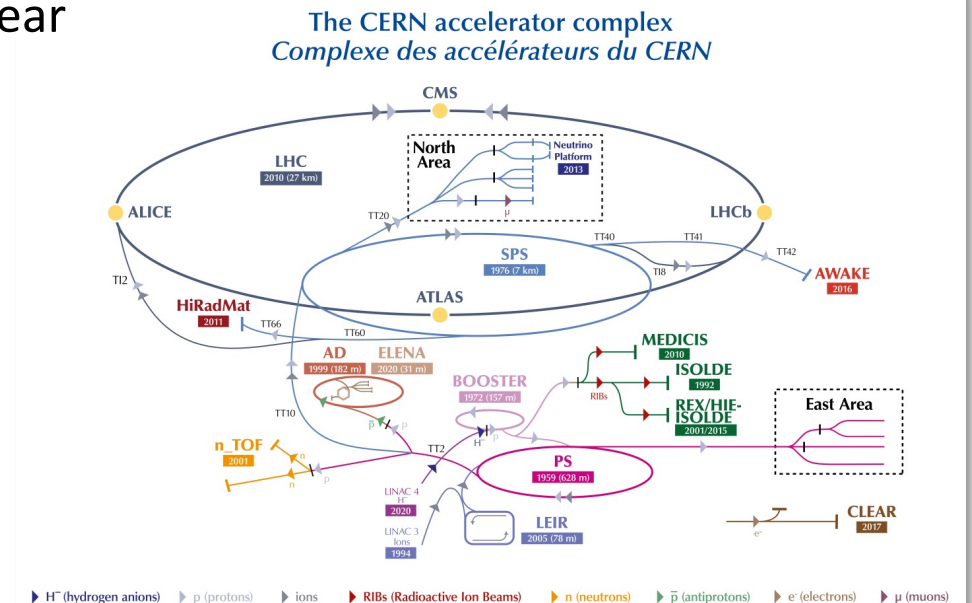


4

Given the 5 long latency before data needs to be released, these experiments do not need to put immediate effort into this yet

Expanding the effort – non-LHC expts

- Have started to discuss with non-LHC experiments, if they could sign-up to the existing LHC policy, or if they should develop a separate policy.
- Meeting organized with experiments Spokespersons on Spring this year
- For experiments associated with some CERN facilities the LHC policy, not well aligned with the procedures from the given physics community:
 - ISOLDE – have developed their own version of the policy:
 - <https://isolde.cern/isolde-open-data-policy>
 - nTOF – will likely develop their own policy, mostly based on existing procedures
 - AD experiments – to discuss if a common AD policy makes sense
- SPS based experiments mostly think LHC policy could work for them, but more discussion needed
 - Further discussion to be organized before the end of the year



Expanding the effort – LEP expts



DELPHI data preservation, re-use, and open access policy

*The DELPHI collaboration,
12 March 2024*

Keywords

DELPHI, OpenData, FAIR, CERN, LEP

The DELPHI experiment conducted e+e- collisions at various energies and produced a unique and irreplaceable data set that has been of great interest to the international scientific community. The experiment ran for a decade between 1989 and 2000, yielding important physics results, especially on electroweak interactions and QCD. However, many results from DELPHI are increasingly difficult to reproduce, as the required detailed knowledge of the scientific community about the detector and the run conditions are vanishing. Therefore, the collaboration strongly discourages attempts to redo high-precision analyses. Despite this, the data still holds potential for further exploration and discovery. Therefore, the data should be preserved and made accessible to the public for various purposes, such as education and citizen science. This document outlines the data preservation, re-use, and open access policy for this valuable data set.

Potential users of the data are encouraged to get in contact with DELPHI scientists to understand the limitations and possibilities. They are also encouraged to register their activity with the DELPHI data preservation board. DELPHI aims to implement the FAIR¹ principle for its data², but its full implementation will be subject to available person power.

Detailed and up to date information about the status and contacts will be made available via the DELPHI web page at <http://delphiweb.cern.ch>.

1 DELPHI Data

The data released by the DELPHI experiment consists of data sets collected during the operation of the experiments between 1989 and 2000. Additionally, various simulated data sets that simulate a large number of physics processes are also being released.

The main format of the data is called SHORT DST, which contains physics information in a compressed format. Reading of this data is supported by various software packages, as described in the documentation. In addition, the original RAW data is also available. Its main purpose is to study individual events with the event server and the display. Currently, the data is available via the CERN EOS storage system, with a backup on CERN's tape archive system, CTA. At a later point in time, the data should also be accessible via the CERN Open Data Portal.

Publications based on the DELPHI data shall give credits to the collaboration and clearly identify the data which has been used, e.g. by quoting an identifier, such as a DOI when available.

Following the convention of CERN, all metadata and data are released under the terms of the CC0 waiver.

2 DELPHI Analysis Software

The DELPHI analysis software, which includes all software for event reconstruction, event simulation, physics analysis and event viewing, is also released under an open-source license and is available via

¹FAIR: Findability, Accessibility, Interoperability, and Reuse

²Guiding Principles for scientific data management and stewardship, currently available at: <https://www.gofair.org/fair-principles>

LEP data represents the highest energy e+e- collision data, very valuable – especially in the context of future e+e- machines (FCC-ee etc...)

ALEPH:

- Open Data available and used in some physics papers
 - (not on CERN OpenData portal though)

DELPHI:

- Recently: s/w ported to latest OS, released Open Data policy
- Data starting to be uploaded to portal

OPEL:

- Renewed effort started recently
- Recently : s/w ported to latest OS, policy being discussed in collaboration

L3:

- No known effort ongoing

https://delphi-www.web.cern.ch/delphi-www/delsec/finalrules/DELPHI_Data_preservation-8.pdf

<https://dphep.web.cern.ch/experiment/aleph>

Summary

- LHC Open Data policy released at end of 2020, now entering implementation period
 - The ODWG standing WG to follow implementation of policy at a high level
 - Should expect some aspects to evolve as we see how this works in practise
 - All large LHC experiments have now (or are very close to) releasing large L3 datasets
 - Will be interesting to see the usage of these dataset
 - Releases broadly in line with expectation from the policy
 - In some cases the model will need to evolve to limit the huge dataset sizes
 - LHCb Ntuple Wizard, ALICE releasing derived/skimmed formats etc...
 - Experiments continuing to release Level-1 and Level-2 Open Data as before
- Expanding the set of experiments with formal Open Data policies
 - All small-LHC experiments have now endorsed the policy
 - Will start to release L3 data before 5 years after the end of Run 3
 - Discussions ongoing with non-LHC experiments
 - Depending on the facility:
 - Plan to write their own policy (tailored for their community)
 - Hope to sign up to the existing LHC policy
- Discussions ongoing about KPIs to monitor Open Data usage:
 - KPIs to be made public as part of the CERN Open Science reports
 - Internal KPIs for more detailed monitoring

Many thanks to the WG members for their valuable input during this process

Backup...



Level 1 Open Data

The HEPData is the tool for storing additional Level-1 data associated to a particular publication. It can store digitized versions of plots, and more detailed information on event selections, efficiencies etc...

← → ↻ 🏠 hepdata.net/record/ins1204284 ☆ 📌 🔄 📄 📱 📂

Popular <https://atlasdqm.cern.ch> [bitly | Basic | a sim...](#) [2011 CERN-Fermit...](#) [Atlas TCT Query](#) [Home - The FASE...](#) [DQ2 Accounting...](#) [Jet/Etmiss Live Pa...](#) [Indico \[LHCC Mee...](#) [LHC Programme C...](#) >> | 📁 All Bookmarks

The [YODA](#) download option now gives the new [YODA2](#) format, with the legacy format still available via the YODA1 download option.

HEPData 🔍 Search HEPData Search

📄 Browse all | 📄 Last updated on 2015-08-25 00:00 | 📄 Accessed 1978 times | 📄 Cite | 📄 JSON

⏪ Hide Publication Information

Constraining R-parity violating Minimal Supergravity with stau₁ LSP in a four lepton final state with missing transverse momentum

The ATLAS collaboration

Conference Paper, 2012.

<https://doi.org/10.17182/hepdata.58712>

INSPIRE Resources

Abstract (data abstract)
CERN-LHC. An interpretation of a search for supersymmetry in final states with four or more leptons (electrons or muons) and missing transverse momentum from proton-proton collisions at a centre-of-mass energy of 7 TeV. The analysis, based on an existing search reported in ATLAS-CONF-2012-001, uses a data sample with total integrated luminosity 2.06 fb⁻¹ recorded in 2011, and finds no significant excess above the expectations from Standard Model processes. Exclusion limits are shown for mSUGRA/CMSSM with m₀=A₀=0, μ>0 and one R-parity violating parameter Lambda₁₂₁=0.032 at the grand unification scale mGUT. This record lists the various limits with acceptances and efficiencies values and gives access to the SLHA files from the analyses.

Table 7 10.17182/hepdata.58712.v1/t7
Data from F 9
Signal acceptance for strong production, gaugino-gaugino production, stau₁-stau₁ production and slepton-slepton production (excluding stau₁-stau₁) as a function of M_{1/2} and Tan(Beta). Statistical fluctuations can be seen, especially where the contributions to the signal region are small (see Figure 11).

cmenergies 7000.0

observables ACC

phrases Exclusive, Proton-Proton Scattering, Jet Production

reactions P P -> .GE.4LEPTONS JETS MM

Showing 50 of 224 values [Show All 224 values](#)

ABS(ETARAP(C=ELECTRON))	< 2.47
ABS(ETARAP(C=ELECTRON,BARREL/END-CAP))	1.37-1.52
ABS(ETARAP(C=JET))	< 2.8
ABS(ETARAP(C=MUON))	< 2.4
E(C=JET)	> 20 GEV
ET(C=ELECTRON)	> 10 GEV
ET(C=ELECTRON,BARREL/END-CAP)	> 15 GEV

Visualize

Level 2 Open Data

alice.physicsmasterclasses.org/MasterClassWebpage.html

Popular https://atlasdqm.c... bitly | Basic | a sim... 2011 CERN-Fermit... Atlas TCT Query CERN-OPEN-2023... Home - The FASE... DQ2 Accounting:...

Einstein in the 21st Century

Main Menu

- Installation
- Support Material
- Students section
- Evaluation
- Instructions for the Institutes
- Description of Exercises
 - English
 - .doc
 - .pdf
 - Deutsch
 - .doc
 - .pdf
 - Français
 - .doc
 - .pdf
 - Italiano
 - .doc
 - .pdf
 - Czech
 - .doc
 - .pdf
 - Portugese
 - .doc
 - .pdf
 - Greek
 - .doc
 - .pdf

As mentioned in the previous section, strange particles do not live long; they decay soon after their production. However, they live long enough to travel some cm distance from the interaction point (IP), where they were produced (primary vertex). Their search is thus based on the identification of their decay products, which must originate from a common secondary vertex.

Neutral strange particles, such as K^0_s and Λ , decay giving a characteristic decay pattern, called V0. The mother particle disappears some cm from the interaction point and two oppositely charged particles appear in its place, which are bent in opposite directions inside the magnetic field of the ALICE solenoid.

In the following red tracks indicate positively charged particles; green tracks indicate negatively charged particles.

The decays we will be looking for are:

$K^0_s \rightarrow \pi^+ \pi^-$

$\Lambda \rightarrow p^+ \pi^-$

$\text{anti } \Lambda \rightarrow p^- \pi^+$

We see that for a pion-pion final state the decay pattern is quasi-symmetric whereas in the pion-proton final state the radius of curvature of the proton is bigger than that of the pion: due to its higher mass the proton carries most of the initial momentum.

We will also be looking for cascade decays of charged strange particles, such as the Ξ^- ; this decays into π^- and Λ ; the Λ then decays into π^- and proton; the initial pion is characterized as a bachelor (single charged track) and is shown in purple.

$\Xi^- \rightarrow \pi^- \Lambda \rightarrow \pi^- p^+ \pi^-$

The search for V0s is based on the decay topology and the identification of the decay products; an additional confirmation of the particle identity is the calculation of its mass; this is done based on the information (mass and momentum) of the decay products as described in the following section.

7. The (invariant) mass calculation

We consider the decay of the neutral kaon to two charged pions, $K^0_s \rightarrow \pi^+ \pi^-$.

Let E , p and m be the total energy, momentum (vector!) and mass of the mother particle (K^0)
 Let E_1 , p_1 and m_1 be the total energy, momentum and mass of the daughter particle number 1 (π^+); and E_2 , p_2 and m_2 the total energy, momentum and mass of the daughter particle number 2 (π^-).

Conservation of energy: $E = E_1 + E_2$ (1)

Example use of real experimental data for education and outreach. In this ALICE masterclass, school students can apply selections to real ALICE data, to emulate a published physics analysis.

The dataset used is openly available for education purposes.

The other large LHC experiments have similar tools for education purposes.



Level 3 Open Data: LHCb

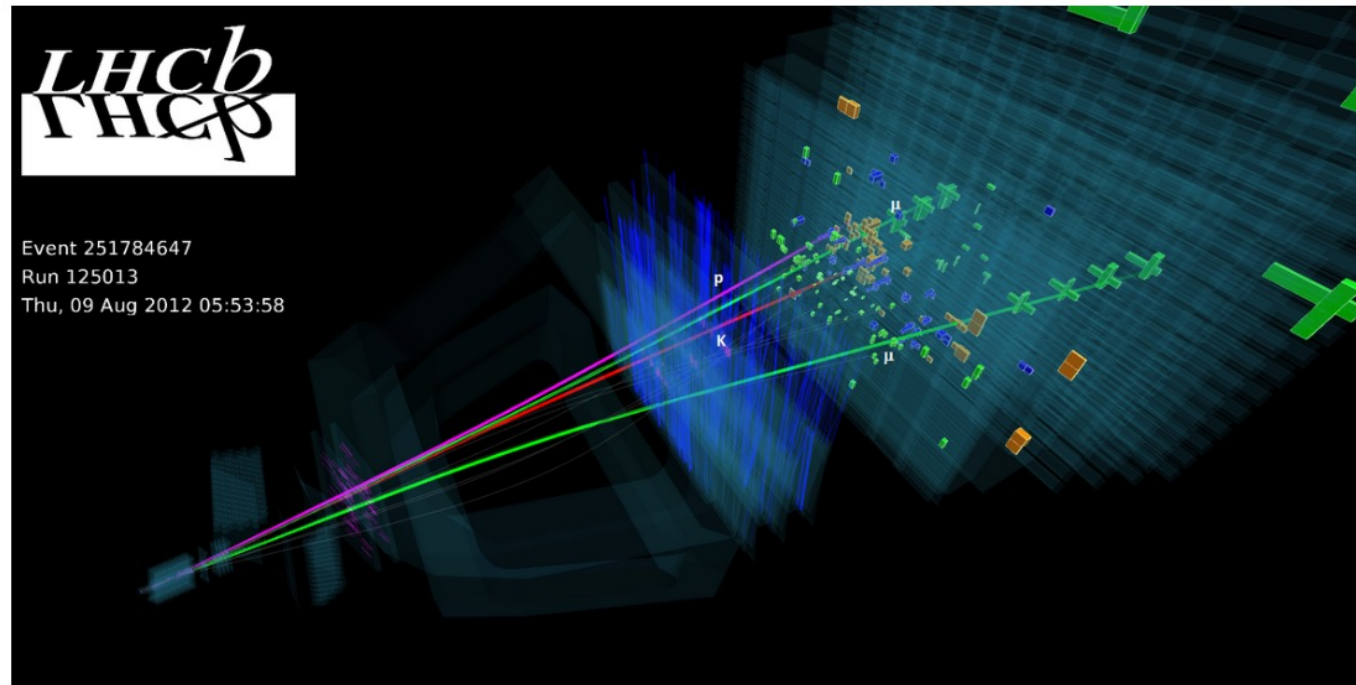
20/12/23

LHCb releases the entire Run I dataset

2023-12-20 by LHCb Collaboration

News

Today the LHCb collaboration completes the release of the data collected throughout the Run I of the Large Hadron Collider at CERN. The sample made available amounts to approximately 800 terabytes (TB) of data. These data, collected by the LHCb experiment in 2011 and 2012, contains information obtained from [proton-proton](#) collisions. The format made available provides pre-filtered data, suitable for a wide range of physics studies. The image below displays an [event](#) recorded during 2012.





Level 3 Open Data: CMS

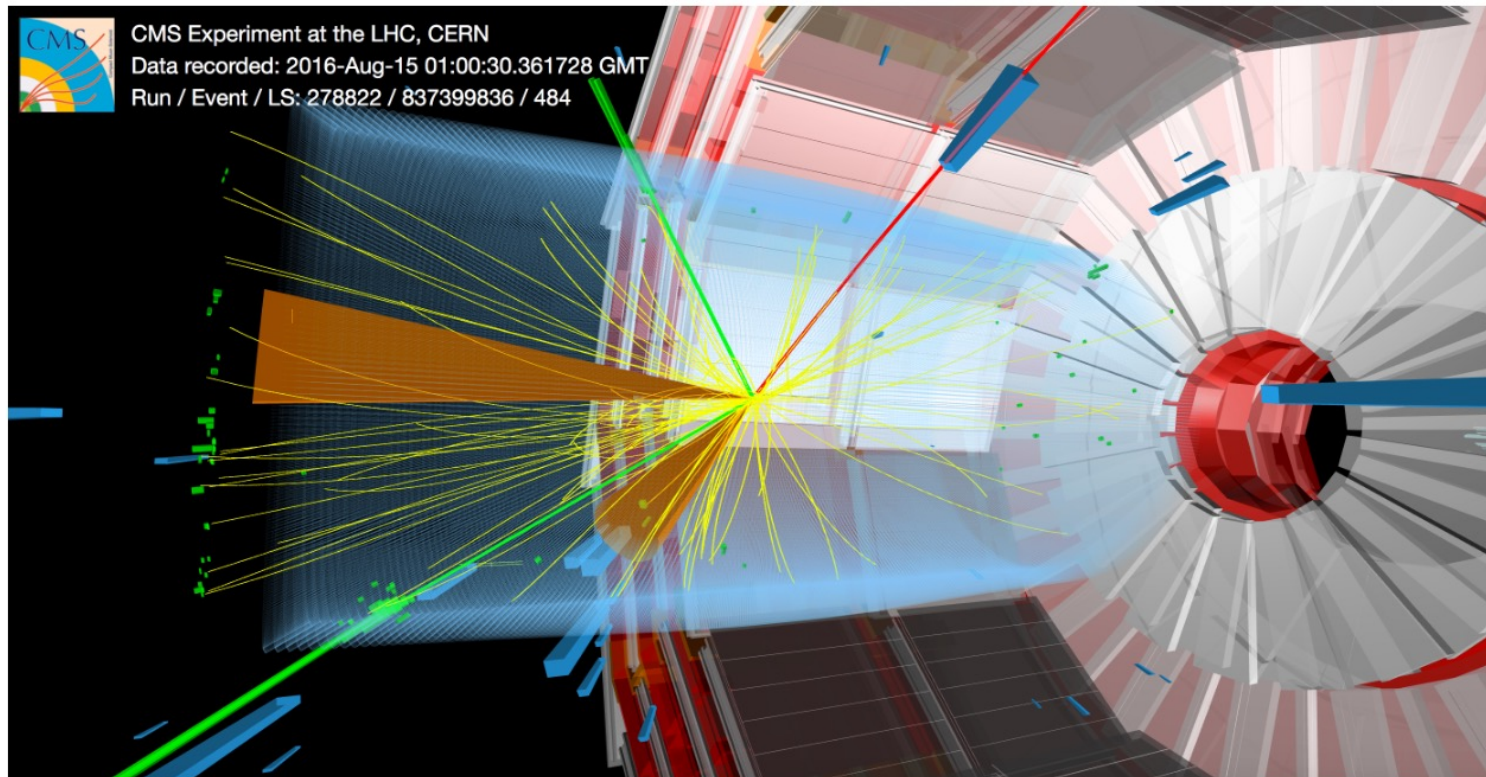
2/4/24

CMS releases 13 TeV proton collision data from 2016

2024-04-02 by CMS Collaboration

News

The CMS experiment at CERN is proud to announce the first release of 13 TeV proton-proton collision data collected in 2016. Over 70 TB of 13 TeV collision data and 830 TB of corresponding simulations are now accessible to the global scientific community and enthusiasts alike through the [CERN Open Data Portal](#).



A candidate event in which a top quark is produced in association with a Z boson. Photo: [CERN Document Server](#)



Level 3 Open Data: ATLAS

1/7/24

ATLAS releases 65 TB of open data for research

2024-07-01 by ATLAS Collaboration

News

Explore over 7 billion LHC collision events – from home

The ATLAS Experiment at CERN has made two years' worth of scientific data available to the public for research purposes. The data include recordings of proton–proton collisions from the Large Hadron Collider (LHC) at a collision energy of 13 TeV. This is the first time that ATLAS has released data on this scale, and it marks a significant milestone in terms of public access and utilisation of LHC data.

“Open access is a core value of CERN and the ATLAS Collaboration,” says Andreas Hoecker, ATLAS Spokesperson. “Since its beginning, ATLAS has strived to make its results fully accessible and reusable through open access archives such as arXiv and HepData. ATLAS has routinely released open data for educational purposes. Now, we’re taking it one step further — inviting everyone to explore the data that led to our discoveries.”

Released under the Creative Commons CC0 waiver, ATLAS has made public all the data collected by the experiment during the 2015 and 2016 proton–proton operation of the LHC. This is approximately 65 TB of data, representing over 7 billion LHC collision events. In addition, ATLAS has released 2 billion events of simulated “Monte Carlo” data, which are essential for carrying out a physics analysis.

External researchers, in particular, are encouraged to explore the ATLAS open data. “Along with the data, we have provided comprehensive documentation on several of our analyses, guiding users through our process step-by-step,” says Zach Marshall, ATLAS Computing co-Coordinator. “These guides provide first-hand experience of working on a real ATLAS result, allowing anyone to test our tools, and evaluate the systematic uncertainties associated with the result for themselves.”

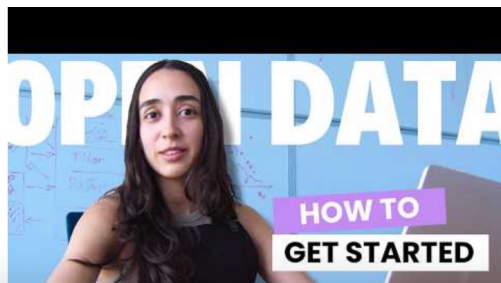
ATLAS traditionally collaborates with non-ATLAS scientists through short-term associations, granting them full access to ATLAS data, internal tools, and information. Through the open data, ATLAS researchers hope to further nurture this dialogue and collaboration. “In particular,” adds Zach, “we’d like to encourage phenomenologists and also computer scientists to explore our datasets, instead of relying on mock-ups.”

Today’s release builds upon previous open data releases for educational use (in 2016 and 2020). “All of our open data releases are now available through the ATLAS open data website,” says Dilia Portillo, ATLAS Outreach and Education co-Coordinator. “The website includes multi-level documentation, video tutorials and online tools aimed at the full-spectrum of users, from high school students to senior particle physics researchers. In addition, the software used to create the education-use open data has been released. This provides a seamless transition from the research open data to all the tutorials for outreach and education, including newly updated Higgs-boson discovery documentation. With a bit of time and dedication, you can go from being a relative novice to carrying out your own analysis.”

The ATLAS open data website also serves as a hub for the community, which includes teachers, students, enthusiasts and, now, scientists. Anyone diving into the open data can also directly engage with ATLAS physicists, who are available to respond to user feedback and take suggestions.

This release marks the start of more to come, with ATLAS’ first release of lead-lead-nuclei collision data data up next. The ATLAS Collaboration, along with the other main LHC experiment collaborations, has committed to making all of its data publicly accessible after a certain time. Openness is deeply ingrained in the culture of high-energy physics, enabling greater accessibility, reproducibility and better science.

Begin your journey with ATLAS open data by following the tutorial below.



WG Original Mandate

The European Commission is driving an Open Data Policy. CERN Council have encouraged the CERN experiments to make their data openly available for external analysis. The CERN Open Data portal is an initial attempt at making the data of the four experiments available through a common interface. So far this has largely been an experimental effort and a common strategy for Open Data has not been established.

This memo calls for forming a small **working group** to explore the conditions for a common Open Data policy from the LHC experiments. The **LHC collaborations are asked to nominate two representatives** each who will be complemented by one representative of RCS-SIS and IT each. The meetings will be convened by Jamie Boyd who reports to the Director of Research and Computing. The charge of the group is to draft a concise document on a common approach by the start of the summer addressing the following open points:

- Who are the expected users of the Open Data and what are their use cases?
 - o What would be the target typical event size for the data that is made available and what would a rough conceptual breakdown of the event content look like? What operations (tightening selections, re-calibrating objects, applying systematic uncertainties etc.) would be possible / not-possible with this event format?
 - o Would simulated data samples be made available with the data? If so which samples, and how would the total simulated data event size compare with the real data? Would there be a mechanism and/or recommendation for users to be able to obtain new simulated samples (e.g. of hypothetical signals)?
 - o What would be the typical total annual data size on the Open Data portal? (note for this to be useful this should be something that can be downloaded to, and processed at a typical university computing installation);
- After what latency period would the collaborations relinquish their exclusive access to experimental data?
- Should there be any rules or recommendations related to who can use the Open data, for example related to members of the experiment that the data is from, or competing experiments?

L3 data – discussion - 2

- Summary of each experiments Latency shown below:

	ALICE	ATLAS	CMS	<u>LHCb</u>
Fraction of data released in: 5 yrs (6 yrs for CMS)	10%	25% (but limiting to <20% of the total data at that time)	50% (but limiting to <20% of the total data at that time)	50%
Fraction of data released in: 10 yrs	50%	50% (but limiting to <20% of the total data at that time)	100% (but limiting to <20% of the total data at that time)	100%
End-of-Collaboration	100%	100%	100%	100%

Resource needs (L3 data)

Expected disk resources for Open Data release (note LHCb have different trigger strategy that increases the size of their dataset. This can be controlled by not releasing some of the exclusive streams if needed).

	ALICE	ATLAS	CMS	LHCb
Run 2	2 PB	0.5 PB	2 PB	10 PB (including Run 1)
Run 3	4 PB	1 PB	4 PB	45 PB
Total	6 PB	1.5 PB	6 PB	55 PB

It is expected that the computing resources will be covered by CERN IT.
The released data will be available through the CERN Open Data Portal.

The experiments will release data in their internal data-formats thus minimizing needed human resources.
The documents have a disclaimer that support for Open Data users will only be provided on a best-effort basis.

Text in public policy document on releasing L3 data

To be released:

- calibrated data
- accompanying Monte Carlo
- Analysis software
- Software environments (VM, containers)

importance of sufficient latency period.

- start 5y after run period
- guideline:
more data → more open data
- full data at close of expt
- escape hatch for ongoing analyses

Reconstructed Data (Level 3) Policy: The LHC experiments will release calibrated reconstructed data with the level of detail useful for algorithmic, performance and physics studies. The release of these data will be accompanied by provenance metadata, and by a concurrent release of appropriate simulated data samples, software, reproducible example analysis workflows, and documentation. Virtual computing environments that are compatible with the data and software will be made available. The information provided will be sufficient to allow high-quality analysis of the data including, where practical, application of the main correction factors and corresponding systematic uncertainties related to calibrations, detector reconstruction and identification. A limited level of support for users of the Level 3 Open Data will be provided on a best-effort basis by the collaborations.

Public data releases will occur periodically, following an appropriate latency period to allow thorough understanding of the data, the reconstruction and calibrations, as well as to allow time for the scientific exploitation of the data by the collaboration. The size of the released datasets will be commensurate with the total amount of data collected of similar type, with the aim to commence data releases within five years of the conclusion of the run period. Data may be withheld by an experiment if there are active analyses ongoing. Full datasets will be made available at the close of the collaboration.

The data will be released from the CERN Open Data Portal under the Creative Commons CC0 waiver, and will be identified with persistent data identifiers, and the data must be cited through these identifiers. Similarly, appropriate acknowledgements of the experiment(s) should be included in publications released using such data, and the publications made clearly distinguishable from those released by the collaboration. Any scientific claims in such publications are the responsibility of their authors and not of the experiments. It is expected that scientific results released using Open Data follow best scientific practices. The experiments may impose rules related to the use of the data by members of their respective collaborations.

External authors should be aware that they will not have access to the vast amount of tacit knowledge built up within the LHC collaborations over the decades of design, construction and operation of the experimental apparatus. To allow external scientists to fully benefit from all the data, knowledge and tools, the collaborations may offer appropriate association programmes.

need for sufficient quality in OD

commit only to best-effort support

data citation (of course) required

but waive responsibility for results

must be clearly marked as external publication

rules for use for own collaborators possible (but not required)

emphasize know-how asymmetry and advertise association programs

Expanding the effort – LEP expts



DELPHI data preservation, re-use, and open access policy

*The DELPHI collaboration,
12 March 2024*

Keywords

DELPHI, OpenData, FAIR, CERN, LEP

The DELPHI experiment conducted e+e- collisions at various energies and produced a unique and irreplaceable data set that has been of great interest to the international scientific community. The experiment ran for a decade between 1989 and 2000, yielding important physics results, especially on electroweak interactions and QCD. However, many results from DELPHI are increasingly difficult to reproduce, as the required detailed knowledge of the scientific community about the detector and the run conditions are vanishing. Therefore, the collaboration strongly discourages attempts to redo high-precision analyses. Despite this, the data still holds potential for further exploration and discovery. Therefore, the data should be preserved and made accessible to the public for various purposes, such as education and citizen science. This document outlines the data preservation, re-use, and open access policy for this valuable data set.

Potential users of the data are encouraged to get in contact with DELPHI scientists to understand the limitations and possibilities. They are also encouraged to register their activity with the DELPHI data preservation board. DELPHI aims to implement the FAIR¹ principle for its data², but its full implementation will be subject to available person power.

Detailed and up to date information about the status and contacts will be made available via the DELPHI web page at <http://delphiweb.cern.ch>.

1 DELPHI Data

The data released by the DELPHI experiment consists of data sets collected during the operation of the experiments between 1989 and 2000. Additionally, various simulated data sets that simulate a large number of physics processes are also being released.

The main format of the data is called SHORT DST, which contains physics information in a compressed format. Reading of this data is supported by various software packages, as described in the documentation. In addition, the original RAW data is also available. Its main purpose is to study individual events with the event server and the display. Currently, the data is available via the CERN EOS storage system, with a backup on CERN's tape archive system, CTA. At a later point in time, the data should also be accessible via the CERN Open Data Portal.

Publications based on the DELPHI data shall give credits to the collaboration and clearly identify the data which has been used, e.g. by quoting an identifier, such as a DOI when available.

Following the convention of CERN, all metadata and data are released under the terms of the CC0 waiver.

2 DELPHI Analysis Software

The DELPHI analysis software, which includes all software for event reconstruction, event simulation, physics analysis and event viewing, is also released under an open-source license and is available via

¹FAIR: Findability, Accessibility, Interoperability, and Reuse

²Guiding Principles for scientific data management and stewardship, currently available at: <https://www.gofair.org/fair-principles>

LEP data represents the highest energy e+e- collision data, very valuable – especially in the context of future e+e- machines (FCC-ee etc...)

ALEPH:

- Open Data available and used in some physics papers
 - (not on CERN OpenData portal though)

DELPHI:

- Recently released Open Data policy
- Currently uploading data to CERN Open Data portal
- Status for L3/OPEL not clear, but probably not happening...

https://delphi-www.web.cern.ch/delphi-www/delsec/finalrules/DELPHI_Data_preservation-8.pdf

<https://dphep.web.cern.ch/experiment/aleph>

CMS Open Data Results

- CMS have been a pioneer of releasing L3 Open Data
- There have been a number of publications on this data covering BSM, QCD, jet-reconstruction and machine learning topics

The screenshot shows the INSPIRE HEP search results page for the query "CMS Open Data". The page features a navigation bar with "Literature", "Authors", "Jobs", "Seminars", "Conferences", and "More...". The search results are displayed in a list format, with each entry including a title, authors, publication information, and citation count. The left sidebar contains filters for "Date of paper", "Number of authors", "Exclude RPP", "Document Type", "Author", and "Subject".

INSPIRE HEP literature find t "CMS Open Data" Q

Literature Authors Jobs Seminars Conferences More...

27 results | cite all Citation Summary Most Recent

Quark-versus-gluon tagging in CMS Open Data with CWoLa and TopicFlow #1
Matthew J. Dolan (U. Melbourne (main)), John Gargalionis (Valencia U., IFIC and Valencia U.), Ayodele Ore (U. Heidelberg, ITP) (Dec 6, 2023)
e-Print: 2312.03434 [hep-ph]
pdf cite claim reference search 1 citation

Search for production of dark fermion candidates in association with heavy neutral gauge boson decaying to dimuon in proton-proton collisions at TeV using CMS open data* #2
Y. Mahmoud (British U. in Egypt), H. Abdallah (Cairo U.), M.T. Hussein (Cairo U.), S. Elgammal (British U. in Egypt) (Apr 19, 2023)
Published in: *Chin.Phys.C* 48 (2024) 4, 043001 · e-Print: 2304.09483 [hep-ex]
pdf DOI cite claim reference search 0 citations

Jet fragmentation properties with CMS open-data #3
Saksevil Arias (CINVESTAV, IPN), Eleazar Cuautele (Mexico U.), Hermes León Vargas (Mexico U., ICN) (Feb 17, 2023)
Published in: *Phys.Scripts* 98 (2023) 3, 035305
pdf DOI cite claim reference search 2 citations

Application of Inferno to a Top Pair Cross Section Measurement with CMS Open Data #4
Lukas Layer (INFN, Padua and Naples U.), Tommaso Dorigo (INFN, Padua and JAEA, Ibaraki), Giles Strong (INFN, Padua) (Jan 24, 2023)
e-Print: 2301.10358 [hep-ex]
pdf cite claim reference search 2 citations

Disentangling quarks and gluons in CMS open data #5
Patrick T. Komiske (MIT, Cambridge, CTP and IAIFI, Cambridge), Serhii Kryhin (MIT, Cambridge, CTP and IAIFI, Cambridge), Jesse Thaler (MIT, Cambridge, CTP and IAIFI, Cambridge) (May 9, 2022)
Published in: *Phys.Rev.D* 106 (2022) 9, 094021 · e-Print: 2205.04459 [hep-ph]
pdf DOI cite claim reference search 19 citations

Inference Aware Neural Optimization for Top Pair Cross-Section Measurements with CMS Open Data #6
Lukas Layer (Naples U.) (Apr 8, 2022)
links cite claim reference search 0 citations

Analyzing N-Point Energy Correlators inside Jets with CMS Open Data #7
Patrick T. Komiske (MIT, Cambridge, CTP), Ian Moutl (Yale U.), Jesse Thaler (MIT, Cambridge, CTP), Hua Xing Zhu (Zhejiang U.) (Jan 19, 2022)
Published in: *Phys.Rev.Lett.* 130 (2023) 5, 051901 · e-Print: 2201.07800 [hep-ph]
pdf DOI cite claim reference search 72 citations

Search for the production of dark matter candidates in association with heavy dimuon resonance using the CMS open data for proton-proton collisions at $\sqrt{s} = 8$ TeV #8
S. Elgammal (British U. in Egypt), M. Louka (British U. in Egypt), A.Y. Ellithi (Cairo U.), M.T. Hussein (Cairo U.) (Sep 23, 2021)
Published in: *Eur.Phys.J.Plus* 138 (2023) 6, 548 · e-Print: 2109.11274 [hep-ex]
pdf DOI cite claim reference search 4 citations

Filters:

- Date of paper:
Number of authors: Single author (3) 10 authors or less (25)
- Exclude RPP: Exclude Review of Particle Physics (27)
- Document Type: article (19) published (17) conference paper (7) thesis (2)
- Author: Jesse Thaler (6) Sergei V. Gleyzer (5) Manfred Paulini (5) Emanuele Usai (4) Bjorn Thomas Burkle (4) Meenakshi Narain (4) Kati Lassila-Perini (3) Mohamed Tarek Hussein (3) Wei Xue (3) Sherif Elgammal (3) [Show 71 more](#)
- Subject: Experiment-HEP (22) Phenomenology-HEP (12) Computing (8) Data Analysis and Statistics (5)