



Designing and implementing the LHC Open Data policy at CERN

Jamie Boyd (CERN)

14 December 2023

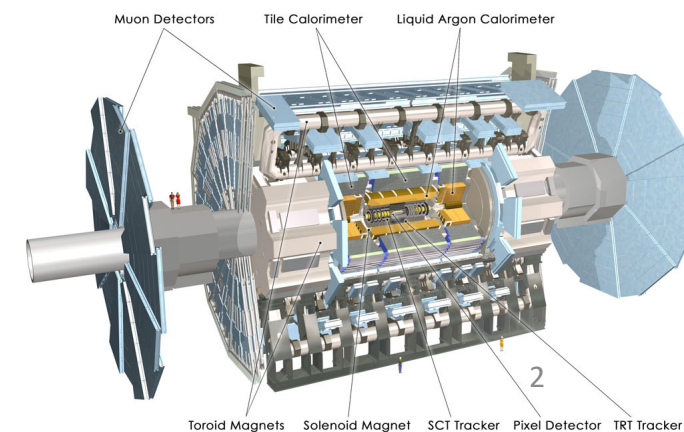
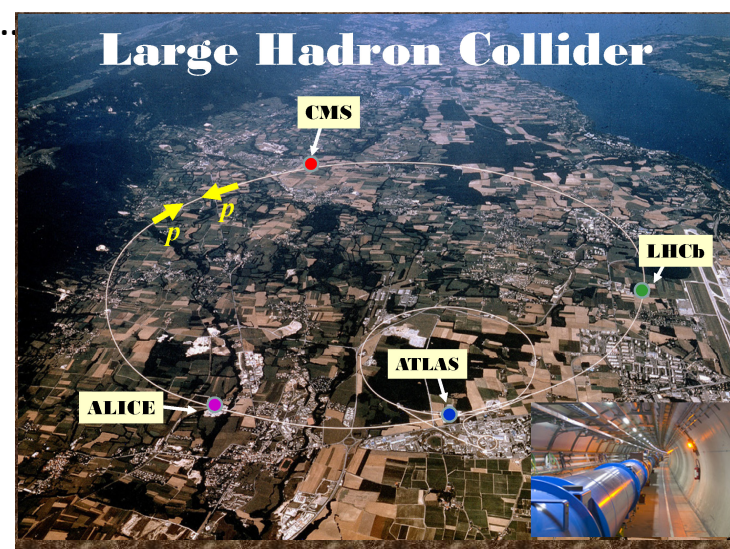
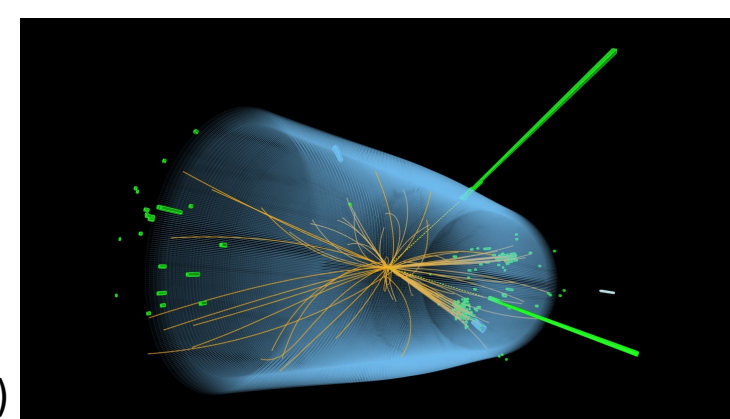
Policy developed along with the LHC experiments:





Context

- LHC physics
 - Study high energy frontier (discovery of Higgs boson, search for new particles/forces etc...)
 - Complex detectors (millions of channels), large data volumes ($O(50)$ PB/year raw data) etc.
- LHC community
 - 4 Large international collaborations (up to 3000 scientists / collaboration)
 - Collaboration lifetimes several decades
 - CERN as host laboratory
 - Collaborations have their own governance
- Increasing importance of Open Data
 - European Commission (relevant for funding applications)
 - European Strategy for Particle Physics (2020 update)
- CERN management mandated working group to explore a common open data policy across the collaborations
 - To be formally endorsed by the Collaborations
- Expected community to use LHC Open Data
 - Professional physicists
 - Non-physics scientists (e.g. computer scientists: machine learning, reconstruction algorithms...)
 - Interested non-scientists





Challenges

- Concerns from the collaborations:
 - **Ownership**
 - Scientists who worked to design (R&D), build, operate the detectors unwilling to lose “ownership” of their data
 - **Effort**
 - Concern from experiment management that could lose effort to operate the experiment if people can analyze the open data without contributing to the experiment
 - Open Data policy of one LHC collaboration can effect other LHC collaborations => push for common policy!
 - **Scientific rigour**
 - Worry about lack of scientific rigour in analysis of open data (spurious claims)
 - **Resources**
 - Required resources within experiments (preparing open data datasets, documentation, storage space etc...), and from CERN side (person power and computing resources)



LHC data levels

- HEP Data Preservation (DPHEP) study group, divided particle physics data into 4 levels:

Small

- Level 1: scientific papers and associated auxiliary data
- Level 2: data tailored for outreach and education purposes

Large

- Level 3: output of data reconstruction. The input for physics analysis.
- Level 4: the raw data from the experiment

- Prior to the OD policy, all collaborations released Level 1, 2 data, and all agree that level 4 data is not useful for external bodies
 - Nearly all discussion was on Level 3 data



Main points in the policy

- **Level 1 data:**
 - Continue to release, including as much auxiliary data as possible to allow re-interpretation of the results (HepData database)
- **Level 2 data:**
 - Continue to release in appropriate formats/schedule
- **Level 3 data:**
 - Next slide
- **Level 4 data:**
 - Not useful, will not be released



After ~1 year of discussions the policy was endorsed by the large LHC experiments in late 2020.
Policy document: <https://cds.cern.ch/record/2745133>



Main points in the policy

- Level 3 data:
 - Release data within 5 years after end of running period
 - *Latency key to counter resistance from within the collaborations*
 - Collaboration can withhold releasing data in special circumstances (unfinished high profile analysis ongoing)
 - Exact format determined by collaboration
 - Same format as used internally in the collaboration for physics
 - Also release analysis software and simulated data samples
 - Needed to allow *meaningful scientific study of the data*
 - Documentation / support offered on best effort basis
 - Data released via CERNs OpenData portal
 - Storage media supplied by CERN (may not be long term solution, but for first 5 year period)
- Open Data policy important for preparing experiments for long term data preservation needs
 - Many aspects of Open Data relevant for Data Presevation



After ~1 year of discussions the policy was endorsed by the large LHC experiments in late 2020.

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



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-Fermilab](#) [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

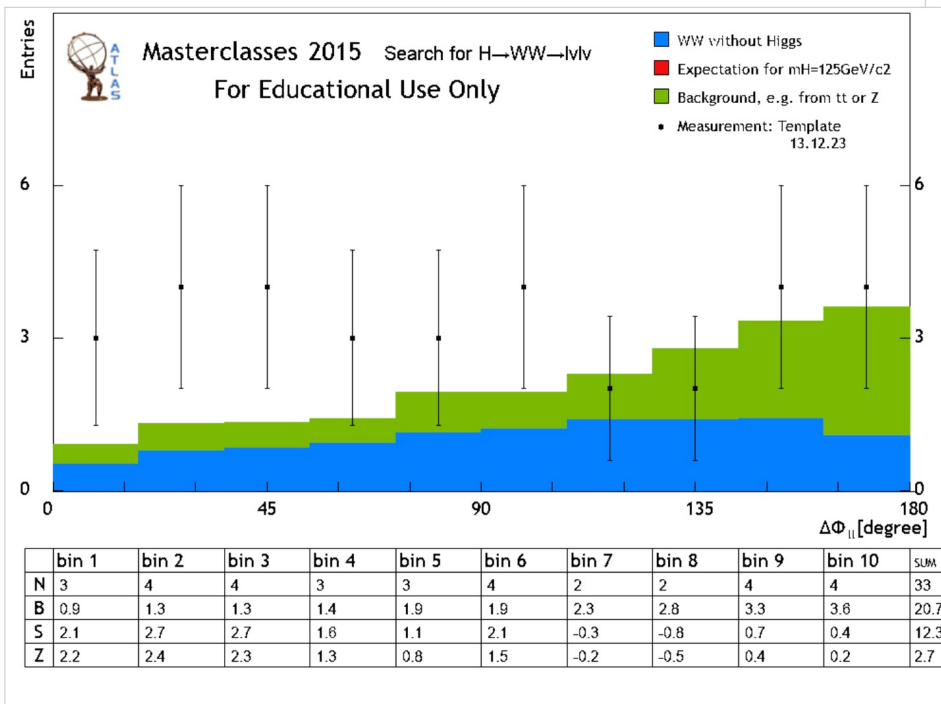
Visualize

<https://www.hepdata.net/>

Level 2 Open Data



WW-Histogramme



Inhalt

- Aims/Tasks
- Identifying particles
- Identifying events
- Measurement
- Analysis**
- Please choose your institute...

Example use of real experimental data for education and outreach. In this ATLAS masterclass, school students can apply selections to real ATLAS 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.

number of bins [1 ... 20] maximum of y axis

standardization 1 2 3 cut on bin number

Higgs contribution



Level 3 Open Data

The CERN Open Data Portal is where the Level 3 Open Data is stored and can be accessed from. Typically a few PB of Level 3 Open data will be released through the portal per year.

opendata.cern.ch/search?page=1&size=20&experiment=CMS

Search

Sort by: **Best match** | asc | Display: **detailed** | 20 results

Found 13155 results.

CMS completes Run-1 heavy ion open data collection
New release of simulations, proton-lead collision data, and proton reference data.
300TB dataset

Getting Started with CMS 2011 and 2012 Open Data
To analyse CMS data collected in 2011 and 2012, you need version 5.3.32 of CMSSW, supported only on Scientific Linux 6. If you are unfamiliar with Linux, take a look at this short introduction to Linux or try this...

Running CMS analysis code using Docker
As an alternative to using a virtual machine, you can run CMS analysis code in a Docker container. If you have not already installed Docker...

CMS Guide to research use of CMS Open Data
If you are interested in finding hints, tips and guidance for conducting a research-oriented analysis using CMS Open Data, please see our notes on this page.

A 300TB dataset

Useful information/
documentation



Level 4 ~~Open~~ Data

A tiny fraction of 1 RAW event from the ATLAS experiment:

0x00000015	0x20000e3f	536874559	lvl1 trigger info[0]	L1 Trigger Bits Before Prescale
0x00000016	0x100000c0	268435648	lvl1 trigger info[1]	
0x00000017	0x8000043f	2147484735	lvl1 trigger info[2]	
0x00000018	0x00021007	135175	lvl1 trigger info[3]	
0x00000019	0x00000e10	3600	lvl1 trigger info[4]	
0x0000001a	0x00080000	524288	lvl1 trigger info[5]	
0x0000001b	0x02c00400	46138368	lvl1 trigger info[6]	
0x0000001c	0x00020001	131073	lvl1 trigger info[7]	
0x0000001d	0x00000816	2070	lvl1 trigger info[8]	L1 Trigger Bits After Prescale
0x0000001e	0x100000c0	268435648	lvl1 trigger info[9]	
0x0000001f	0x80000018	2147483672	lvl1 trigger info[10]	
0x00000020	0x00021001	135169	lvl1 trigger info[11]	
0x00000021	0x00000e10	3600	lvl1 trigger info[12]	
0x00000022	0x00000000	0	lvl1 trigger info[13]	
0x00000023	0x02c00400	46138368	lvl1 trigger info[14]	
0x00000024	0x00020000	131072	lvl1 trigger info[15]	
0x00000025	0x00000010	16	lvl1 trigger info[16]	L1 Trigger Bits After Veto
0x00000026	0x00000000	0	lvl1 trigger info[17]	
0x00000027	0x00000008	8	lvl1 trigger info[18]	
0x00000028	0x00000000	0	lvl1 trigger info[19]	
0x00000029	0x00000810	2064	lvl1 trigger info[20]	
0x0000002a	0x00000000	0	lvl1 trigger info[21]	
0x0000002b	0x00000400	1024	lvl1 trigger info[22]	
0x0000002c	0x00000000	0	lvl1 trigger info[23]	

The LHC experiments are incredibly complex, with millions of readout channels covering many detector technologies.

The programs to reconstruct the raw data are hugely complex and need a lot of expertise in the experiment to be able to use in a useful way.

It is therefore not realistic for people from outside of the experiment to be able to extract meaningful results from the raw data.



Level 3 Open Data Use

The policy discussed in this talk started at the beginning of 2021, and with up to 5 years latency. Therefore, not so much data has been released through the policy to date. However, the CMS experiment has been releasing Open Data for several years.

INSPIRE HEP

literature

Literature Authors Jobs Seminars Conferences More...

80 results | cite all Citation Summary Most Recent

Quark-versus-gluon tagging in CMS Open Data with CWoLa and TopicFlow #1
Matthew J. Dolan, John Gargalionis, Ayodele Ore (Dec 6, 2023)
e-Print: [2312.03434](#) [hep-ph]
pdf cite claim reference search 0 citations

Jet Energy Calibration with Deep Learning as a Kubeflow Pipeline #2
Daniel Holmberg (U. Helsinki (main)), Dejan Golubovic (CERN), Henning Kirschenmann (Helsinki Inst. of Phys.) (Aug 23, 2023)
Published in: *Comput.Softw.Big Sci.* 7 (2023) 1, 9 • e-Print: [2308.12724](#) [hep-ex]
pdf DOI cite claim reference search 0 citations

Potential of the Julia Programming Language for High Energy Physics Computing #3
Jonas Eschle (U. Zurich (main)), Tamás Gál (Erlangen - Nuremberg U., Theorie III), Mosè Giordano (Imperial Coll., London), Philippe Gras (IRFU, Saclay), Benedikt Hegner (CERN) et al. (Jun 6, 2023)
Published in: *Comput.Softw.Big Sci.* 7 (2023) 1, 10 • e-Print: [2306.03675](#) [hep-ph]
pdf DOI cite claim reference search 2 citations

Baler -- Machine Learning Based Compression of Scientific Data #4
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

Leveraging an open source serverless framework for high energy physics computing #5
Vincenzo Eduardo Padulano (Valencia, Polytechnic U. and CERN), Pablo Oliver Cortés (Valencia, Polytechnic U.), Pedro Alonso-Jordá (Valencia, Polytechnic U.), Enric Tejedor Saavedra (CERN), Sebastián Risco (Valencia, Polytechnic U.) et al. (May 1, 2023)
Published in: *J.Supercomput.* 79 (2023) 8, 8940-8965

Filters:

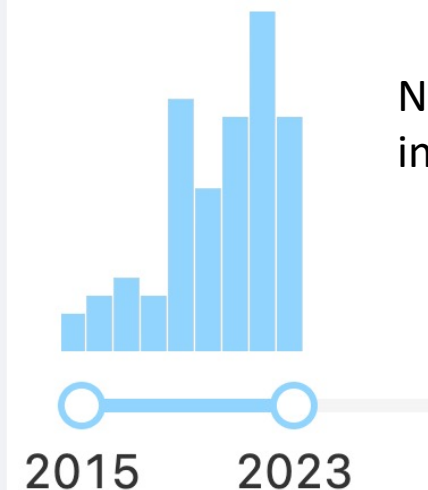
- Date of paper:
- Number of authors: Single author (14) 10 authors or less (70)
- Exclude RPP: Exclude Review of Particle Physics (80)
- Document Type: article (52) published (39) conference paper (22) thesis (6) review (1)
- Author: Jesse Thaler (11) Kati Lassila-Perini (6)

80 papers have been released using that data.

Mostly covering:

- Pure physics research
- Using Machine Learning techniques to improve event reconstruction

Date of paper

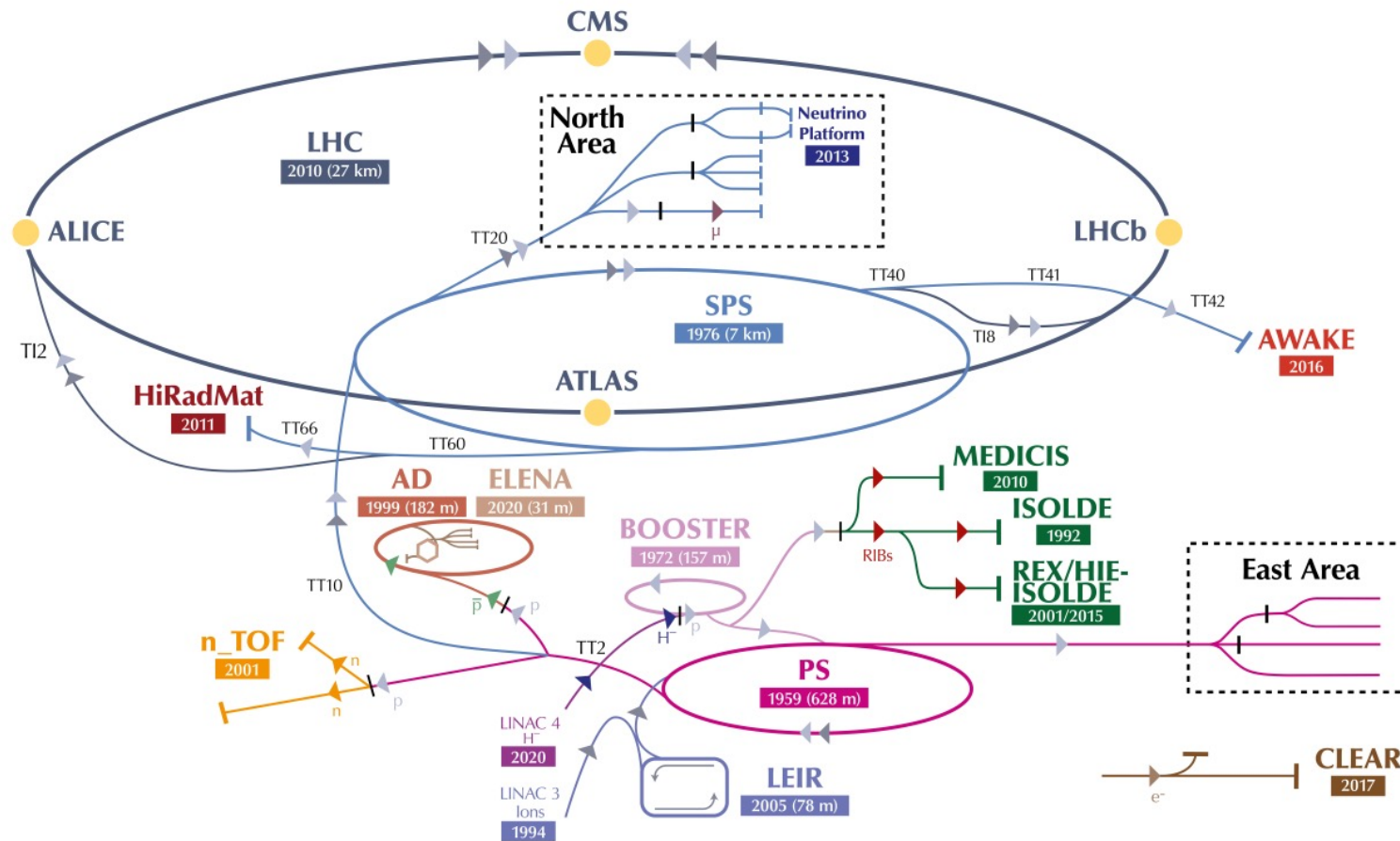


Number of papers increasing per year



Future Plans

The CERN accelerator complex *Complexe des accélérateurs du CERN*



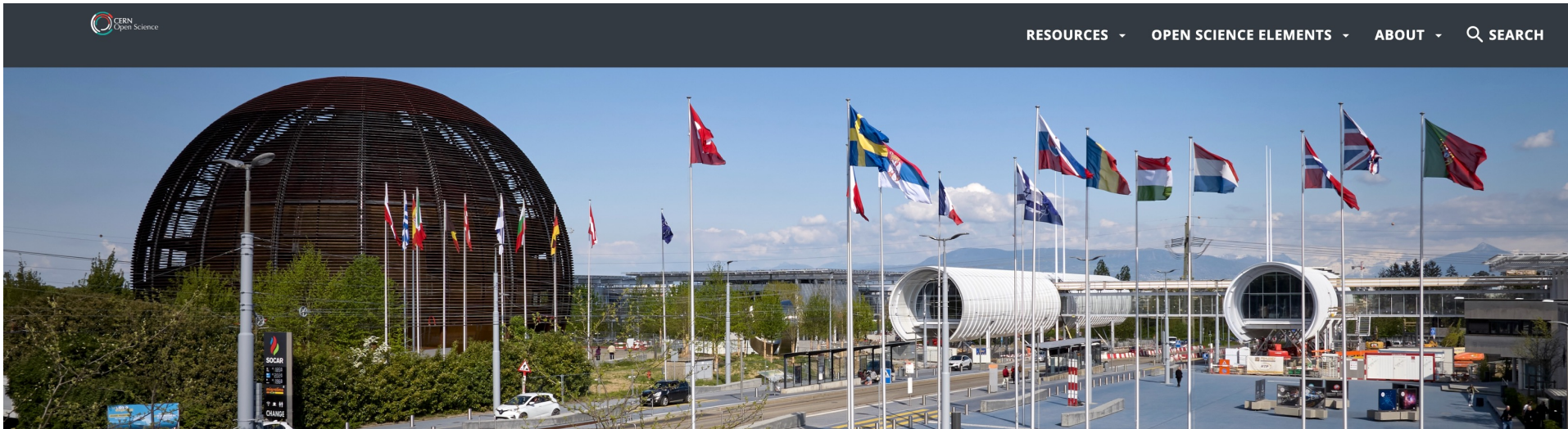
- CERN has a diverse accelerator complex feeding many experiments, not just the large-LHC experiments discussed so far.
- CERN management has asked us to start the process of getting the smaller CERN experiments to sign up to the Open Data policy.
- So far the initial response to this has been generally positive, but there are several more steps to go through.
- In my experience the smaller experiments are less worried about issues related to “effort” and “ownership” than the large LHC experiments, but are more worried about “resources” since they have less people in the experiment.



Open Science at CERN

Since the Open Data policy was released CERN has produced its broader Open Science policy (released in 2022). This covers the following areas:

- Open Access
- Open Data +
- Open Source Software
- Open Source Hardware
- Research Integrity
- Open Infrastructure
- Research Assessment
- Training & Outreach
- Citizen Science



Welcome to CERN Open Science

At CERN, we believe that the practice of open science is key to delivering on our organizational [mission](#).

CERN Open Science describes an evolving ecosystem of policies, initiatives, services, and technologies, driven by people from across the organization with the goal to maximize the global impact of research conducted at CERN.

NEWS





Summary

- The CERN LHC Open Data policy was released at the end of 2020
- It was crafted to find the balance between the wish to open the data, with the constraints/concerns of the experiments
 - We believe a good balance was found in the final policy document
- Following this, Open Data has been released and is being used
- We are now trying to enlarge the policy to cover all experiments at CERN
 - Smaller experiments have different challenges which need to be addressed
- The effort is now embedded in the recent CERN Open Science effort