

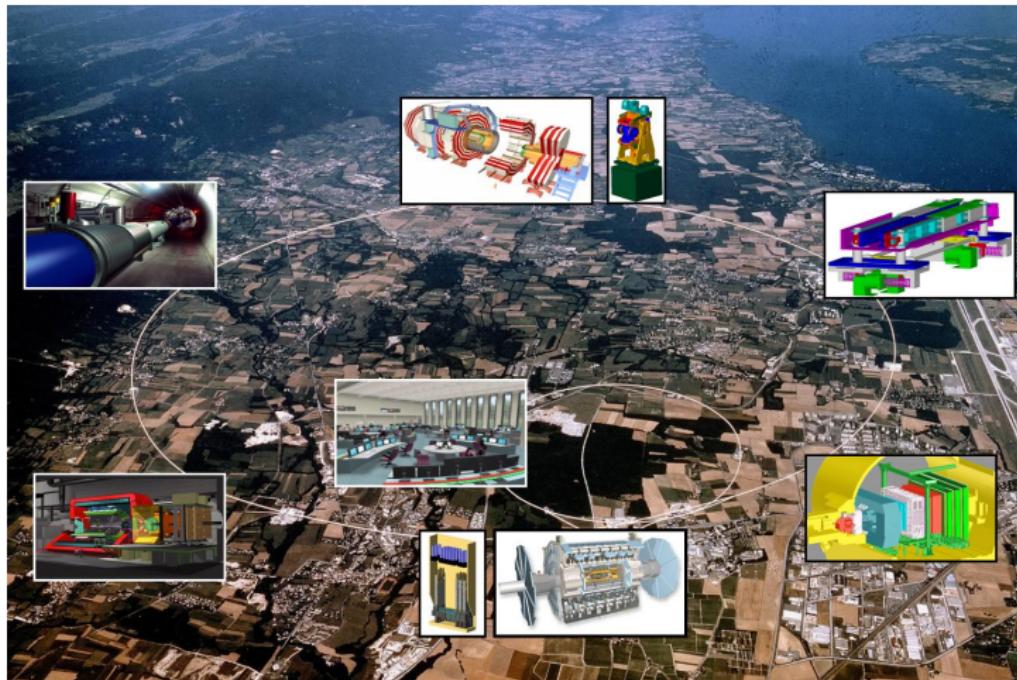


Open is not enough: fostering reproducible and reusable particle physics research

Tibor Šimko

CERN

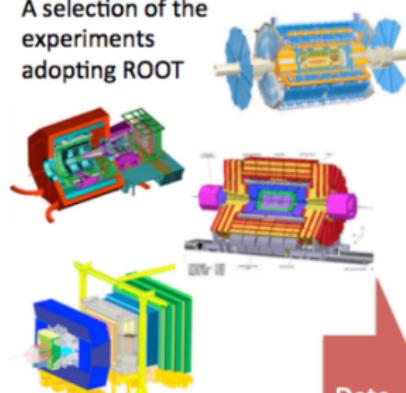
CERN Large Hadron Collider



Large Hadron Collider with ATLAS, ALICE, CMS, LHCb detectors

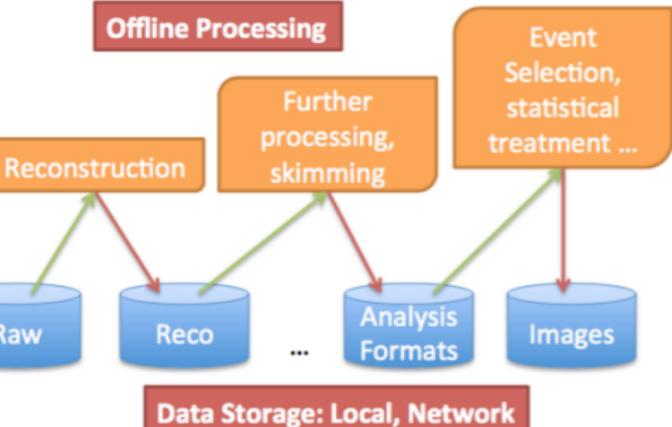
HEP data analyses

A selection of the experiments adopting ROOT



Event Filtering

Data



D. Krücker *et al* <https://indico.desy.de/indico/event/18343>

Typical data chain from acquisition through analysis to plots

LHC data pyramid

~KB/paper

~GB/analysis

~TB/analysis

~PB/year

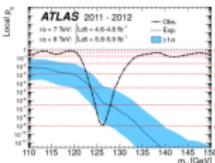
~GB/sec

1. Provision of additional documentation for the published results

2. Simplified data formats for analysis in outreach and training exercises

3. Reconstructed data and simulations as well as the analysis level software to allow a full scientific analysis

4. Basic raw level data (if not yet covered as level 3 data) and their associated software which allows access to the full potential of the experimental data



↑
analysis

↑



Four data levels for capture, preservation and opening

CERN Open Data portal

Explore more than **1 petabyte** of open data from particle physics!

Start typing...

search examples: collision datasets, keywords education, energy/TeV

Explore

- datasets
- software
- environments
- documentation

Focus on

- ATLAS
- ALICE
- CMS
- LHCb

Get started

<http://opendata.cern.ch/>

LHC collaboration data policies

Alice data preservation strategy

Guido, October 6, 2014

The Alice Experiment at the LHC experiments on how to proceed in the future, the basic constraints of preservation to honour and fulfil societal interests by the international community. This document aims to provide a general overview of the current situation, and to highlight the main issues, long term preservation must be an essential aspect of the data processing framework and will be the focus of the data analysis group. This document also aims to provide a general overview of the Alice data preservation strategy and policies. Recommendation, long term preservation of data is a key concern of the LHC data preservation strategy and policies. Recommendation, long term preservation of data is a key concern of the LHC data preservation strategy and policies. Recommendation, long term preservation of data is a key concern of the LHC data preservation strategy and policies.

The answer documents also discuss the basic principles and rule guide the policies adopted by the SGDD data preservation policy.

Alice data formats

The Alice experiment performs many step of the data processing chain starting from basic raw data delivered by the detector or the trigger system, making it physics analysis-ready data and ending with physics analysis results. The data processing chain consists of several stages, each stage producing a specific type of data and leaving conditions messages are needed to trackback raw data information into physics interlock, calibration, simulation and reconstruction.

- 1. Raw data recorded by the detectors along with the associated status discriminating between different types of raw data, such as triggers, calibration and reconstruction experiments. These provide the input of the reconstruction algorithm, together with the calibration data.
- 2. Measured raw data, including rates of the most generic level (RAW mode) and data resulting from the raw data processing chain. The raw data are produced by the reconstruction algorithm.
- 3. Event summary (ES) data produced by the reconstruction algorithm, for both Monte Carlo and real data, containing information on the event topology and particle flow.
- 4. Event summary (ES) data produced by the reconstruction algorithm, for both Monte Carlo and real data, containing information on the event topology and particle flow.
- 5. Event analysis object files, saved individually or together with the generic purpose AMF for specific analysis.
- 6. Published physics results.

These different formats of the Alice data serve as a specific volume for data preservation. While formats can differ with time, the software provides sufficient metadata to record and process any format. An alternative to raw data preservation is to keep the raw data and produce a compressed version of the raw data, which is better suited for the final publication of the results, as mentioned in the answer document.

The answer documents also discuss the data preservation strategy of the Alice experiment. There are no specific rules for the final publication of the results, but they can be agreed upon by the members of the Alice collaboration upon approval to the Alice Physics Board. The original datasets used to produce published results, together with the relevant software source code and services are called to being open to preservation.



Restricted data → embargo period (~5 years) → open data

Information organisation

The screenshot shows a detailed view of a dataset page. At the top, there's a navigation bar with 'opendata' and 'cern' on the left, a search bar with a magnifying glass icon in the center, and a 'About' link on the right. Below the header, the main content area has a title 'Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21ReReco-v1/AOD)' and a subtitle '(/Mu/Run2010B-Apr21ReReco-v1/AOD, cms collaboration)'. It includes a citation: 'Cite as: CMS collaboration (2014). Mu primary dataset in AOD format from RunB of 2010 (/Mu/Run2010B-Apr21ReReco-v1/AOD). CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.B8MR.C4A2'. Below the citation are several tabs: 'Dataset' (selected), 'Citation', 'CMS', 'Collision energy 7TeV', 'Accelerator LHCnLHC', and 'Parent Dataset: /Mu/Run2010B-v1/RAW'. The main content area is divided into sections: 'Description', 'Notes', 'Related Datasets', 'Characteristics', and 'System Details'. Each section contains specific details about the dataset, such as the number of events, global tag, and recommended release for analysis.

How were these data selected?

There are four categories of triggers in the Mu dataset (with significant overlaps):

- 70% inclusive single muon triggers with varying trigger pt threshold 3.5,7,9,11,13,15,17,19,21 GeV plus a few with loosened quality cuts.
- 20% isolated single muon triggers with varying trigger pt threshold 9,11,13,15,17 GeV.
- 10% inclusive dimuon triggers with varying trigger pt threshold 3.5 GeV plus one Z->mu mu trigger with loosened quality cuts.
- 20% combinations of muon triggers with various pt thresholds 3.5,7,8,9,11 GeV with some EM/e/gamma or hadronic/jet energy deposit with thresholds 6-100 GeV.

How were these data validated?

During data taking all the runs recorded by CMS are certified as good for physics analysis if all subdetectors, trigger, lumi and physics objects (tracking, electron, muon, photon, jet and MET) show the expected performance. Certification is based first on the offline shifters evaluation and later on the feedback provided by detector and Physics Object Group experts. Based on the above information, which is stored in a specific database called Run Registry, the Data Quality Monitoring group verifies the consistency of the certification and prepares a json file of certified runs to be used for physics analysis. For each reprocessing of the raw data, the above mentioned steps are repeated. For more information see:

CMS data quality monitoring: Systems and experiences

The CMS Data Quality Monitoring software experience and future improvements

The CMS data quality monitoring software: experience and future prospects

How can you use these data?

You can access these data through the CMS Virtual Machine. See the instructions for setting up the Virtual Machine and getting started in

How to install the CMS Virtual Machine

Getting started with CMS open data

Curated context information about data selection, validation, use

Information discovery

opendata CERN

Search

About ▾

Filter by type

<input type="checkbox"/> Dataset	997
<input type="checkbox"/> Collision	100
<input type="checkbox"/> Derived	173
<input type="checkbox"/> Simulated	723
<input type="checkbox"/> Documentation	56
<input type="checkbox"/> About	8
<input type="checkbox"/> Activities	19
<input type="checkbox"/> Authors	3
<input type="checkbox"/> Guide	16
<input type="checkbox"/> Help	2
<input type="checkbox"/> Policy	4
<input type="checkbox"/> Report	1
<input type="checkbox"/> Environment	19
<input type="checkbox"/> Condition	5
<input type="checkbox"/> VM	11
<input type="checkbox"/> Validation	3
<input type="checkbox"/> Glossary	22
<input type="checkbox"/> News	9
<input type="checkbox"/> Software	33
<input type="checkbox"/> Analysis	16
<input type="checkbox"/> Framework	4
<input type="checkbox"/> Tool	8
<input type="checkbox"/> Validation	5
<input type="checkbox"/> Suplementaries	2642
<input type="checkbox"/> Configuration	917
<input type="checkbox"/> Luminosity	3
<input type="checkbox"/> Trigger	1722

Sort by: Most recent ▾ asc. ▾

Display: detailed ▾ 20 results ▾

Found 3778 results. < 1 2 3 4 5 6 7 8 9 >

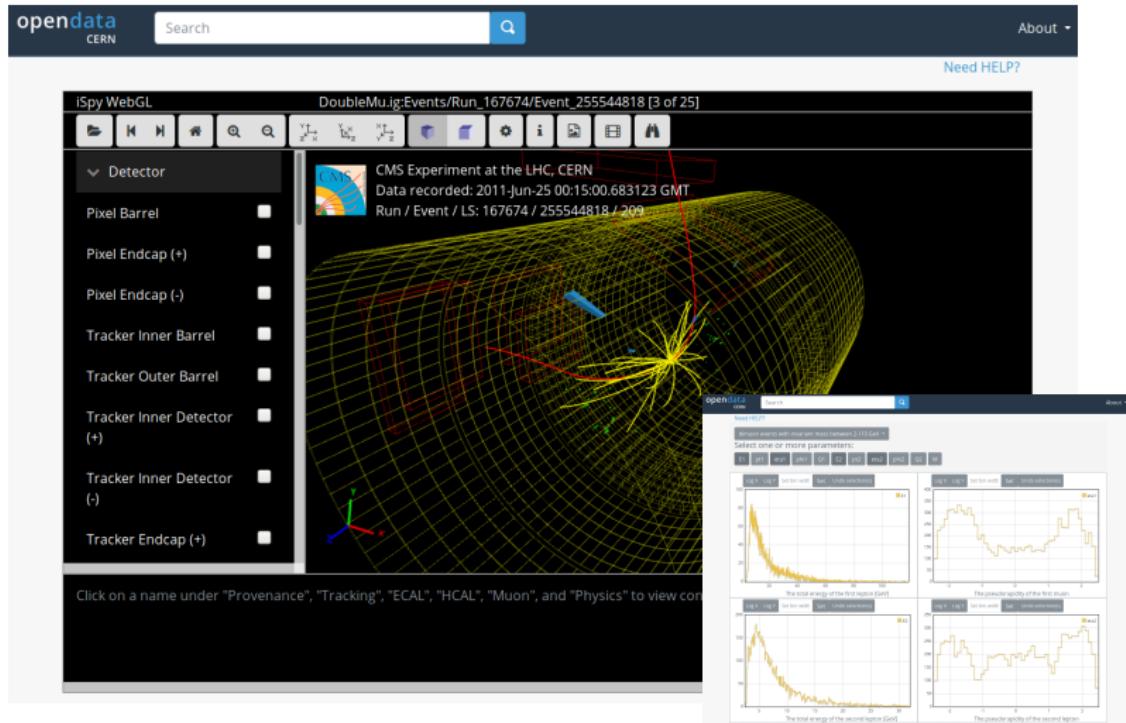
/TTJets_MSDecays_scaleup_mt172_5_7TeV-madgraph-tauola/Summer11LegDR-PU_S13_START53_LV6-v1/AODSIM
Simulated dataset TTJets_MSDecays_scaleup_mt172_5_7TeV-madgraph-tauola in AODSIM format for 2011 collision data (SM Systematic Variations)
See the description of the sim...
[Dataset](#) [Simulated](#) [CMS](#)

/Vector1MToZZTo4L_M-125p6_7TeV-jHUGenV3-pythia6/Summer11LegDR-PU_S13_START53_LV6-v1/AODSIM
Simulated dataset Vector1MToZZTo4L_M-125p6_7TeV-jHUGenV3-pythia6 in AODSIM format for 2011 collision data (SM Inclusive)
See the description of the simulated dataset nam...
[Dataset](#) [Simulated](#) [CMS](#)

/VBFHiggs0PToGG_M-125p6_7TeV-jHUGenV4-pythia6-tauola/Summer11LegDR-PU_S13_START53_LV6-v1/AODSIM

Explore a variety of data, software, VMs, supplementary material

Education use cases



Use interactive event display and basic histogramming

Analysis examples

opendata CERN Search About ▾

Higgs-to-four-lepton analysis example using 2011-2012 data

Jomhari, Nur Zulaiha; Geiser, Achim; Bin Anuar, Afiq Aizuddin; (2017). Higgs-to-four-lepton analysis example using 2011-2012 data. CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.JKB8.RR42

Software Analytics CMS Accelerator CERN/LHC

Description

This research level example is a strongly simplified reimplementation of parts of the original CMS Higgs to four lepton analysis published in *Phys.Lett. B716* (2012) 30-61, arXiv:1207.7235.

The published reference plot which is being approximated in this example is https://inspirehep.net/record/1124338/files/H4l_mass_3.png. Other Higgs final states (e.g. Higgs to two photons), which were also part of the same CMS paper and strongly contributed to the Higgs boson discovery, are not covered by this example.

The example consists of different levels of complexity. The highest level minimal understanding of the content of this paper and of the meaning educational exercises. The lower levels might also be interesting for ed with the linux operating system and the ROOT analysis tool.

Use with

The example uses legacy versions of the original CMS datasets in the A publication due to improved calibrations. It also uses legacy versions o but not identical to, the ones in the original publication. These legacy d in many later CMS publications.

/DoubleElectron/Run2011A-12Oct2013-v1/AOD
/DoubleMu/Run2011A-12Oct2013-v1/AOD

CMS Preliminary $\bar{G} = 7 \text{ TeV}, L = 5.05 \text{ fb}^{-1}$, $\bar{G} = 8 \text{ TeV}, L = 5.26 \text{ fb}^{-1}$
CMS Open Data $\bar{G} = 7 \text{ TeV}, L = 2.3 \text{ fb}^{-1}$, $\bar{G} = 8 \text{ TeV}, L = 11.6 \text{ fb}^{-1}$

Events / 3 GeV

Events / 3 GeV

Legend: Data, Z+X, ZZ, ZZ -> 4l, m_h = 126 GeV

Legend: Data, Z+X, ZZ, ZZ -> 4l, m_h = 126 GeV

Install virtual machines and run realistic physics analysis examples

Research use cases

PRL 119, 132005 (2017)

PHYSICAL REVIEW LETTERS

work ending
28 SEPTEMBER 2017

PHYSICAL REVIEW D 96, 074001 (2017)

Exposing the QCD Splitting Function with CMS Open Data

Andrew Larkoski,^{1,*} Steven Maitra,¹ Jesse Thaler,^{1,2} Anatoli Trifunovic,¹ and Wen Xu³

¹Center for Theoretical Physics, Research School of Physics and Engineering, The Australian National University, Canberra ACT 0485, Australia

²Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

³Department of Physics, University of California, Berkeley, Berkeley, California 94720, USA

The splitting function is a universal property of quantum chromodynamics (QCD) which describes how energy is shared between partons. Despite its ubiquitous appearance in many QCD calculations, the splitting function cannot be measured directly, since it always appears multiplied by a collinear singularity like $\delta(\hat{z})$. In this Letter, we show how to measure the splitting function directly from the jet splitting function for sufficiently high jet energies. This provides a way to expose the splitting function through jet substructure measurements at the Large Hadron Collider. In this Letter, we use public data sets from the CMS experiment to study the two-prong substructure of jets and test the $1 \rightarrow 2$ splitting function of QCD. To our knowledge, this is the first-ever physics analysis based on the CMS Open Data.

DOI: 10.1103/PhysRevLett.119.132005

Quantum chromodynamics (QCD), like any weakly coupled gauge theory, exhibits universal behavior in the small-angle limit. When two partons become collinear in QCD, they split into two new partons whose momenta factorizes into a $2 \rightarrow 2 = 1$ -scattering cross section multiplied by a universal $1 - 2$ splitting probability, which conserves energy and momentum. Universality of the splitting function is a fundamental property of QCD and appears in many applications, most famously in duality [1–3]. Duality is also a key element of the evolution equation [1–3] (see also [4–13]), and it is at the heart of the factorization theorem in hadron-hadron collisions [14–17]. In addition, parton shower programs have adopted the splitting function approach [18–20], and fixed-order subtraction schemes enforce the $1 \rightarrow 2$ splitting function [19–21], while the jet clustering series is based on the splitting function [22–24]. The splitting function can be extended to multi-parton splittings at tree level and beyond [25–41]; however, in all orders of perturbation theory [42–45] it is typically the $1 \rightarrow 2$ splitting function that has been measured recently, by subtraction techniques [46–52], because it is distinguished by a $1 \rightarrow 2$ cross section of heavy particles (the “jet”) that has appeared in many jet substructure studies (including those of the Large Hadron Collider (LHC)) [53–56].

Due to the oblique nature of the $1 \rightarrow 2$ splitting function, it cannot be directly measured at a collider, since collinear universality is inseparable from the existence of color singularities and closely related nonperturbative configurations of gluons. Specifically, when two partons separated by an angle θ , the $1 \rightarrow 2$ splitting function takes the form

$$dP_{1 \rightarrow 2} = \frac{d\theta}{\theta} dP_{\text{coll}}(\theta), \quad (1)$$

0551-9600/17/11132005/13 \$15.00 © 2017 American Physical Society

132005-1 © 2017 American Physical Society

where the P_{coll} are the Altarelli-Patrizzi QCD splitting functions [5] which depend on the momentum fraction z and the parton flavors i, j , and i, j . Crucially, this expression is manifestly $1 \rightarrow 2$ symmetric, so it is well suited to cancel corresponding virtual singularities from loop diagrams. In this sense, there is no way to directly measure the splitting function, since it is always accompanied by a coarse, overwhelming indirect evidence that $P_{\text{coll}}(z)$ is a universal function from the many sources of QCD at different high-energy scattering (see e.g. [17–19]).

In this Letter, we show how to directly test the $1 \rightarrow 2$ splitting function in QCD by studying the two-prong substructure of jets. Our method is based on soft drop [26–28], a technique that removes soft radiation from a jet and build two-prong substructure is found. When applied to ordinary quarks and gluons, soft drop is a good approximation to the $1 \rightarrow 2$ splitting function that exposes the collinear core of the jet. As shown in Ref. [21], the massless exchange between the two prongs (denoted τ_{coll}) is proportional to the splitting function in Eq. (1), and the cross section for z_1 asymptotes to the QCD splitting function in the high-energy limit. While various forms of τ_{coll} have appeared in many jet substructure studies (including those of the LHC) [27,28], to the best of our knowledge, no published τ_{coll} distribution has ever been presented using actual collider data, though there have been several theoretical predictions [29–32] and ATLAS and ALICE [33] Collaborations. Here, we present the first analysis of τ_{coll} using LHC data, taking advantage for the first time of the CMS Open Data.

The CMS Open Data are derived from the 2016 center-of-mass proton-proton collisions recorded in 2016 and released to the public via the CERN Open Data Portal in November 2017 [34]. The data are presented as analysis object data (AOD) format, which is a CMS-specific data schema based

Jet substructure studies with CMS open data

Anatoli Trifunovic,^{1,3} Wei Xu,^{1,3} Andrew Larkoski,^{1,3} Steven Maitra,^{1,3} and Jesse Thaler,^{1,3}

¹Center for Theoretical Physics, Massachusetts Institute of Technology, Cambridge, Massachusetts 02139, USA

²Physics Department, University of Michigan, Ann Arbor, Michigan 48109, USA

³University of Buffalo, The State University of New York, Buffalo, New York 14260-1500, USA

(Received 4 May 2017; published 3 October 2017)

We use public data from the CMS experiment to study the two-prong substructure of jets. The CMS open data are presented as analysis object data (AOD) format, which is a CMS-specific data schema based on the 2016 center-of-mass proton-proton collisions at the Large Hadron Collider. In this Letter, we use public data sets from the CMS open data and data obtained from particle shower generators, and we also compare to analytic jet substructure calculations performed in modified leading-logarithmic order. Although the CMS open data do not include simulation data to help estimate systematic uncertainties, we use track-only observables to validate these substructure studies.

DOI: 10.1103/PhysRevD.96.074001

INTRODUCTION

In November 2016, the CMS experiment at the Large Hadron Collider (LHC) released the CMS Open Data project [1]. To our knowledge, this is the first time in the history of particle physics that research-grade collision data has been made publicly available by a major particle physics official experimental collaboration. The CMS open data were reconstructed from 7 TeV proton-proton collisions in the 2016 LHC run, and are intended for use by the physics community when pileup contamination was minimal and trigger thresholds were relatively low. The CMS open data provide an excellent opportunity to the particle physics community to study the properties of jets, which are more difficult at higher luminosities and for demonstrating the scientific value of open data releases.

In this Letter, we use the CMS open data to study the substructure of jets. Jets are collimated sprays of particles that are copiously produced in LHC collisions, and by studying their internal structure, one can gain valuable information about their parentage [2–10]. A key application of jet substructure is tagging boosted heavy objects like top quarks and Higgs bosons [11–14], which are often produced with soft drops [15–20]—that is, jets composed of many subjets with large angular separation. We study these various two-prong substructure observables. In addition to comparing the CMS open data to particle shower generators, we also compare to analytic jet substructure calculations using recently developed analytic techniques [16–20,26,28]. In a companion paper, we will add to drop to the CMS open data to study the two-prong substructure of jets [20]. A similar strategy was used in preliminary CMS [16], STAR [28], and ALICE [20] heavy ion studies to test for dense QCD radiation and the splitting function from the dense QCD medium [20,21].

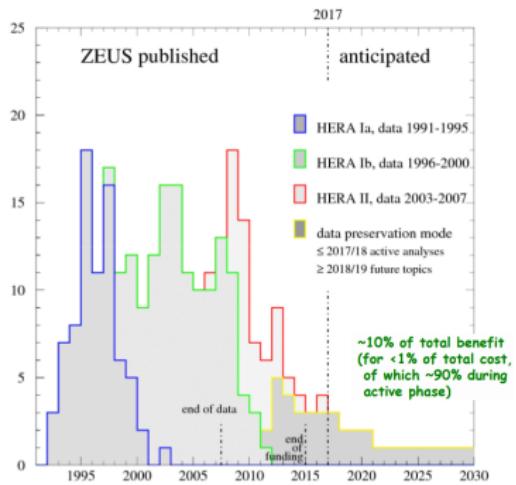
For studying jet substructure, the key feature of the CMS open data is that they contain full information about particle

*To highlight the vibrancy of the field, we have attempted to list all published jet substructure measurements from ATLAS and CMS [2–10].

¹The original mass drop trigger [24] was a pioneering trigger in jet substructure; see also precursor work in Refs. [2–5].

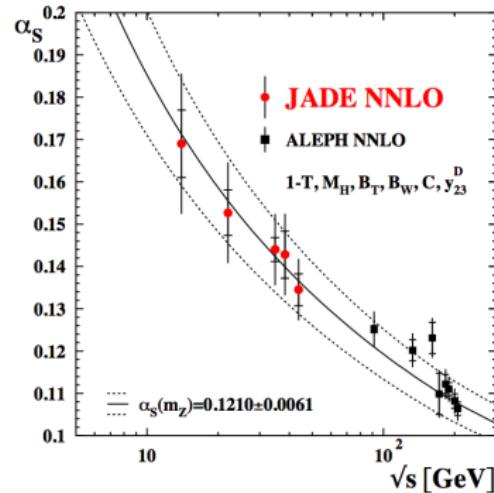
Independent analyses by theorists (Jesse Thaler *et al*, MIT)

Long-term value of data!



Achim Geiser <https://indico.cern.ch/event/588219>

Collaborations publish papers even ~ 15 years after data taking ends



DPHEP <https://arxiv.org/abs/1205.4667>

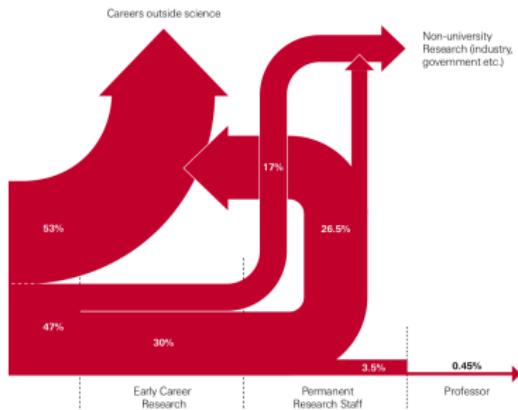
JADE data (1979–1986) still unique even ~ 35 years later

Long-term value of knowledge?



CMS collaboration

Experimental physics done by groups of \sim 3000 physicists

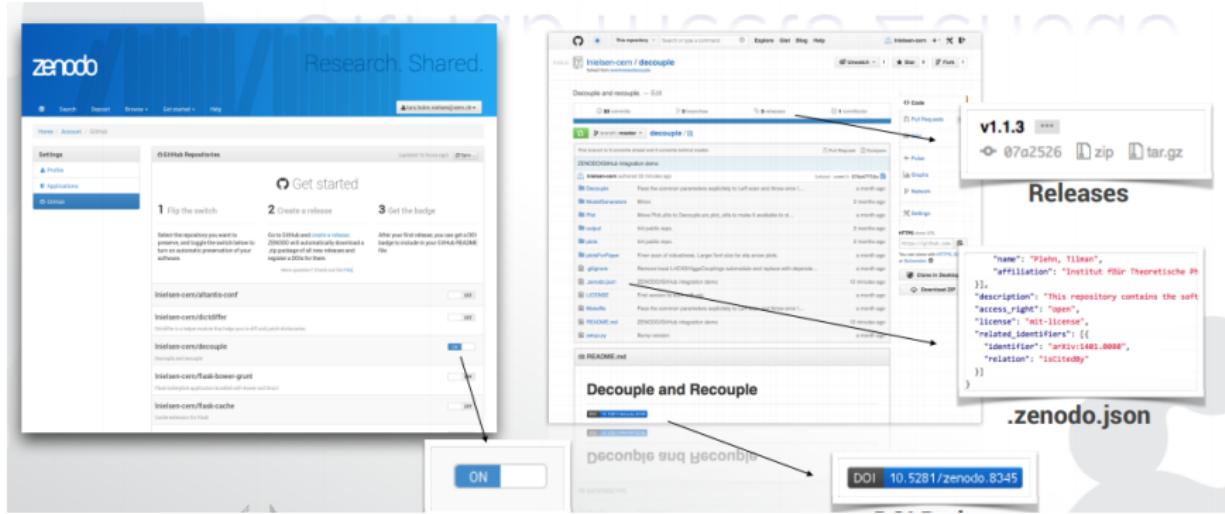


Career after PhD

THE ROYAL SOCIETY

High turnover of young researchers

Preserving scientific code



<https://guides.github.com/activities/citable-code>

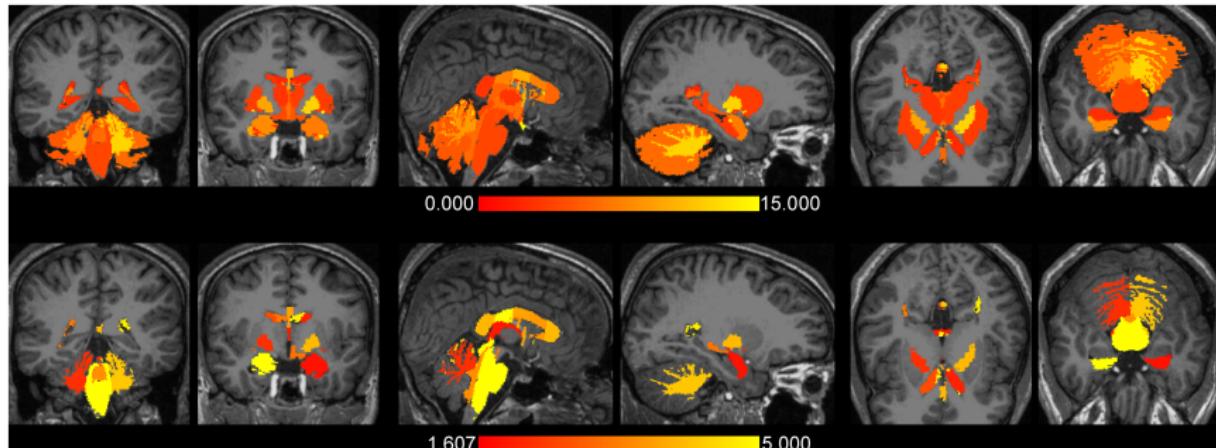
GitHub ↔ Zenodo bridge to automatically preserve releases

Code is not enough

The Effects of FreeSurfer Version, Workstation Type, and Macintosh Operating System Version on Anatomical Volume and Cortical Thickness Measurements

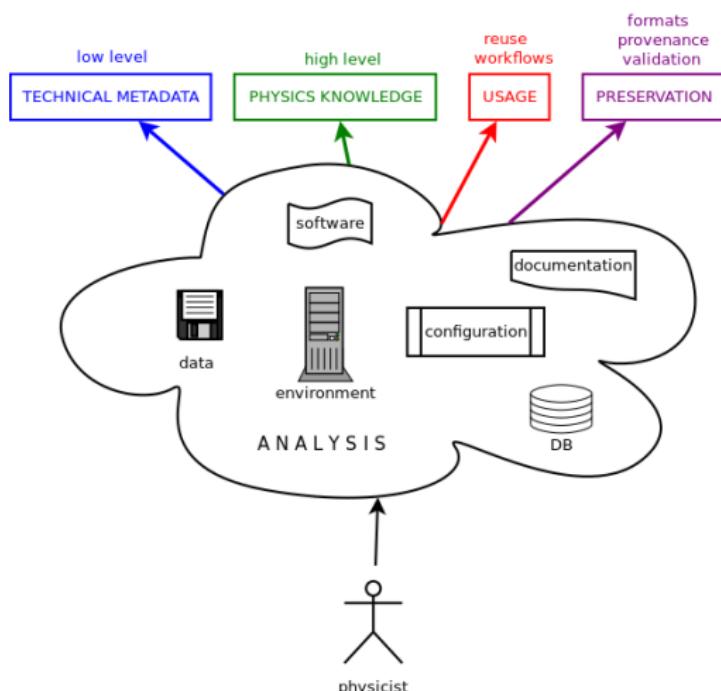
Ed H. B. M. Gronenschild , Petra Habets, Heidi I. L. Jacobs, Ron Mengelers, Nico Rozendaal, Jim van Os, Machteld Marcelis

Published: June 1, 2012 • DOI: 10.1371/journal.pone.0038234



$8.8 \pm 6.6\%$ (volume) and $2.8 \pm 1.3\%$ (cortical thickness)

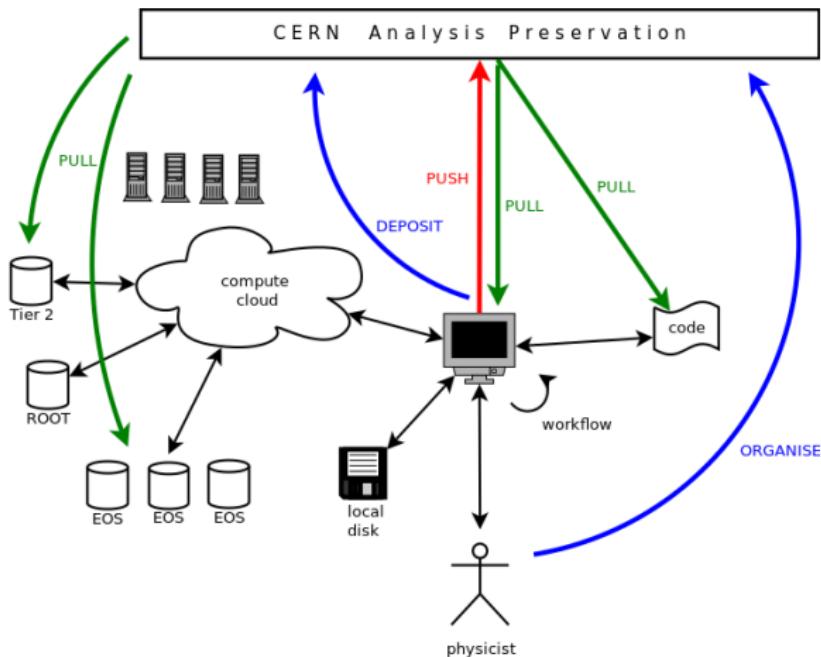
Preserving analyses 1/2



- JSON Schema
- W3C DCAT
- domain-specific fields

Structuring knowledge behind research data analysis

Preserving analyses 2/2



INVENIO)

- datasets:
local storage,
cloud storage
- software:
Git, SVN
- information:
DBs, TWiki,
SharePoint
- protocols:
HTTP, XRootD

Taking consistent snapshot of analysis assets at a certain time

CERN Analysis Preservation

The screenshot shows the CERN Analysis Preservation web interface. At the top, there's a navigation bar with 'CERN Analysis Preservation' and various icons for search, create, search, and more. Below the header, a 'Submission Form' section titled 'Preserve your analysis' asks for an 'Analysis Name' (input field) and has a large blue 'Start Preserving' button. To the right, several sections are listed with arrows: 'Basic Information' (with a note about providing information for all parts), 'Stripping/Turbo selections [0 items]', 'ntuple/userDST-production [0 items]', and 'User Analysis'. The 'User Analysis' section is expanded, showing a terminal window with the following commands:

```
$ pip install cap-client
$ export CAP_SERVER_URL=https://analysispreservation.cern.ch/
$ export CAP_ACCESS_TOKEN=<your generated access token from server>
$ cap-client files upload <file path> --pid/-p <existing pid>
$ cap-client files upload file.json -p 89b593c498874ec8bcacf88944c458a7
File uploaded successfully.
```

Web based and command-line based deposit workflows

How FAIR is scientific data?

■ Findable

- persistent identifiers
- rich metadata
- indexed and searchable

■ Accessible

- retrievable by identifiers
- standard protocols
- metadata vs data accessibility

■ Interoperable

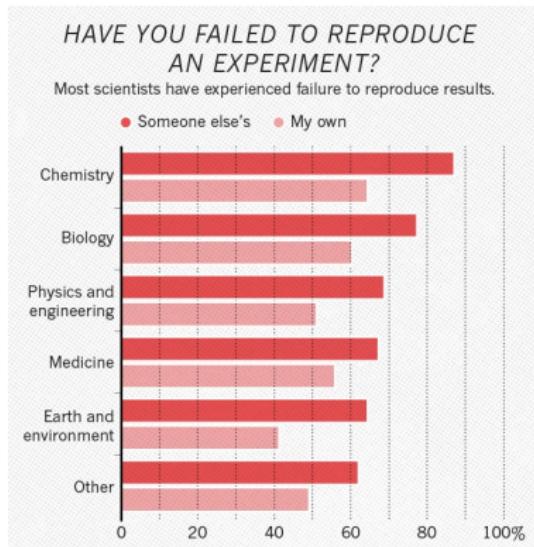
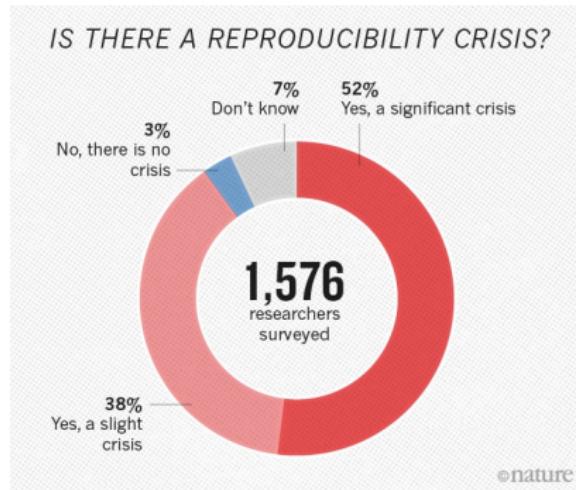
- knowledge representation language
- common vocabularies
- references to other metadata and data

■ Reusable

- domain-relevant attributes and community standards
- clear licensing
- provenance tracking

<https://www.nature.com/articles/sdata201618>

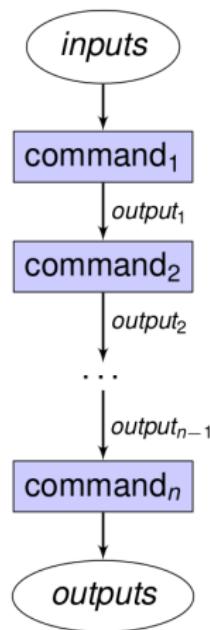
Reusable and reproducible?



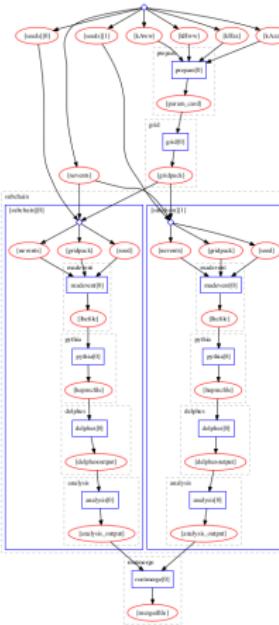
<https://www.nature.com/news/1-500-scientists-lift-the-lid-on-reproducibility-1.19970>

Half of researchers cannot reproduce their own experimental results

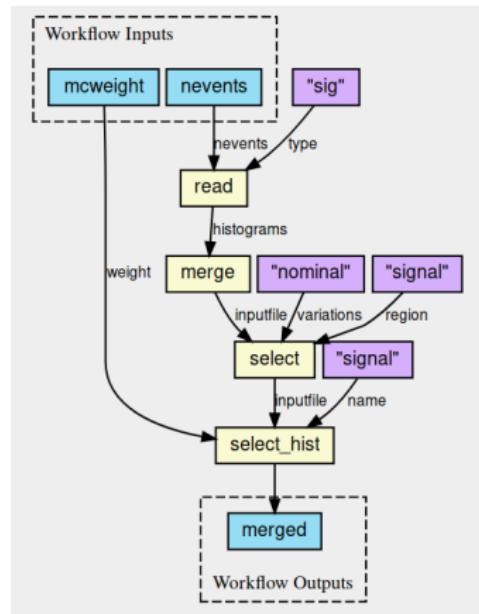
Computational workflows



Serial



Yadage



CWL

REANA Reusable Analyses



Reproducible research data analysis platform

Flexible

Run many computational workflow engines.



Scalable

Support for remote compute clouds.



Reusable

Containerise once, reuse elsewhere. Cloud-native.



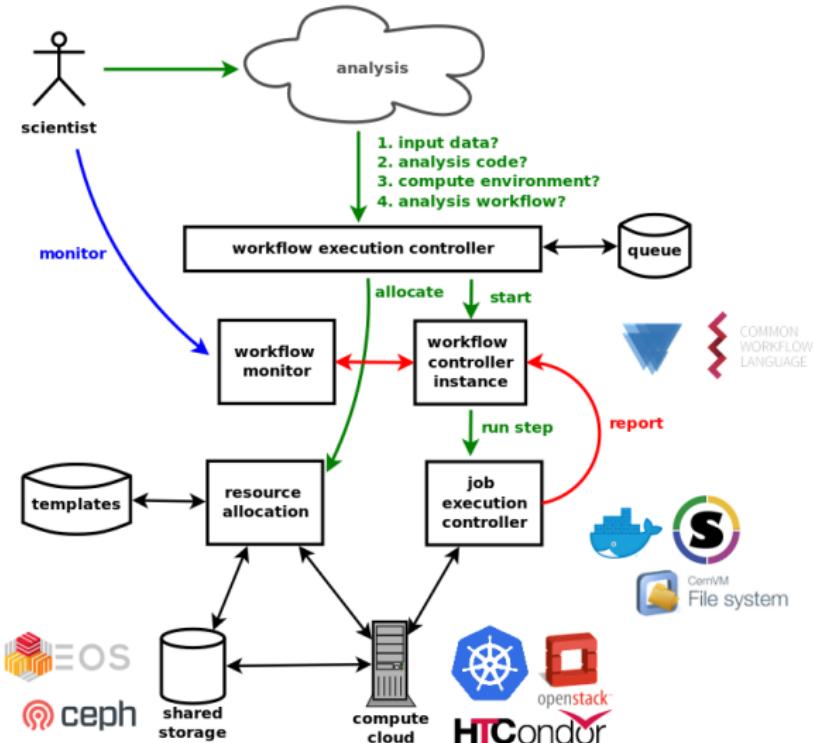
Free

Free Software. MIT licence.
Made with ❤ at CERN.



<http://www.reana.io/>

REANA architecture



Data production examples

The screenshot shows a web page from the "opendata CERN" site. At the top, there's a search bar with a magnifying glass icon and an "About" link. Below the header, a title reads "Validation code for reprocessing AOD from 2011 MinimumBias RAW sample". It includes author information ("Lassila-Perini, Kati") and a citation ("Cite as: Lassila-Perini, Kati; (2017). Validation code for reprocessing AOD from 2011 MinimumBias RAW sample. CERN Open Data Portal. DOI:10.1461/OPENDATA.CMS:399.BGv4"). There are tabs for "Software", "Analysis", "IMF", and "Associated datasets". A "Description" section follows, containing text about the code's purpose and usage. Below it is a "Use with" section with a dataset link. A "Characteristics" section shows a dataset of 1 file (14.6 kB). Finally, a "System Details" section provides instructions for using the CMS Open Data VM environment.

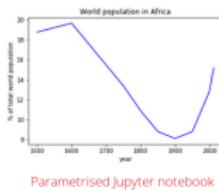
This screenshot shows another page from the "opendata CERN" site. The title is "Example code for production of flat jet tuple using 2011 data". It includes author information ("Zenayev, Olegandr; Haapalehto, Matias") and a citation ("Cite as: Zenayev, Olegandr; Haapalehto, Matias; (2017). Example code for production of flat jet tuple using 2011 data. CERN Open Data Portal. DOI:10.1461/OPENDATA.CMS:1793.120"). There are tabs for "Software", "Analysis", "IMF", and "Associated datasets". A "Description" section explains the purpose of the code. Below it is a "Use with" section with a dataset link. A "Characteristics" section shows a dataset of 1 file (17.0 kB). A "System Details" section provides instructions for using the CMS Open Data VM environment.

Reconstruction and flat jet tuple production from CMS open data

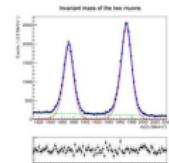
Data analysis examples

```
$ cat inputs/names.txt  
Jane Doe  
Joe Bloggs  
$ reana-client start  
$ reana-client download  
$ cat outputs/greetings.txt  
Hello Jane Doe!  
Hello Joe Bloggs!
```

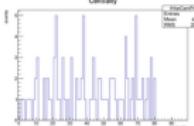
"Hello world"



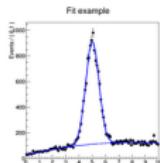
Parametrised Jupyter notebook



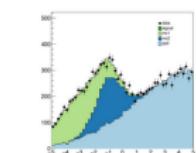
LHCb rare charm decay search



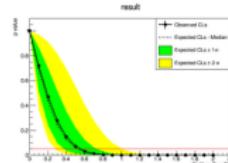
ALICE LEGO train test run



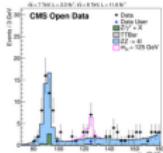
ROOT/RooFit physics analysis



ATLAS BSM search



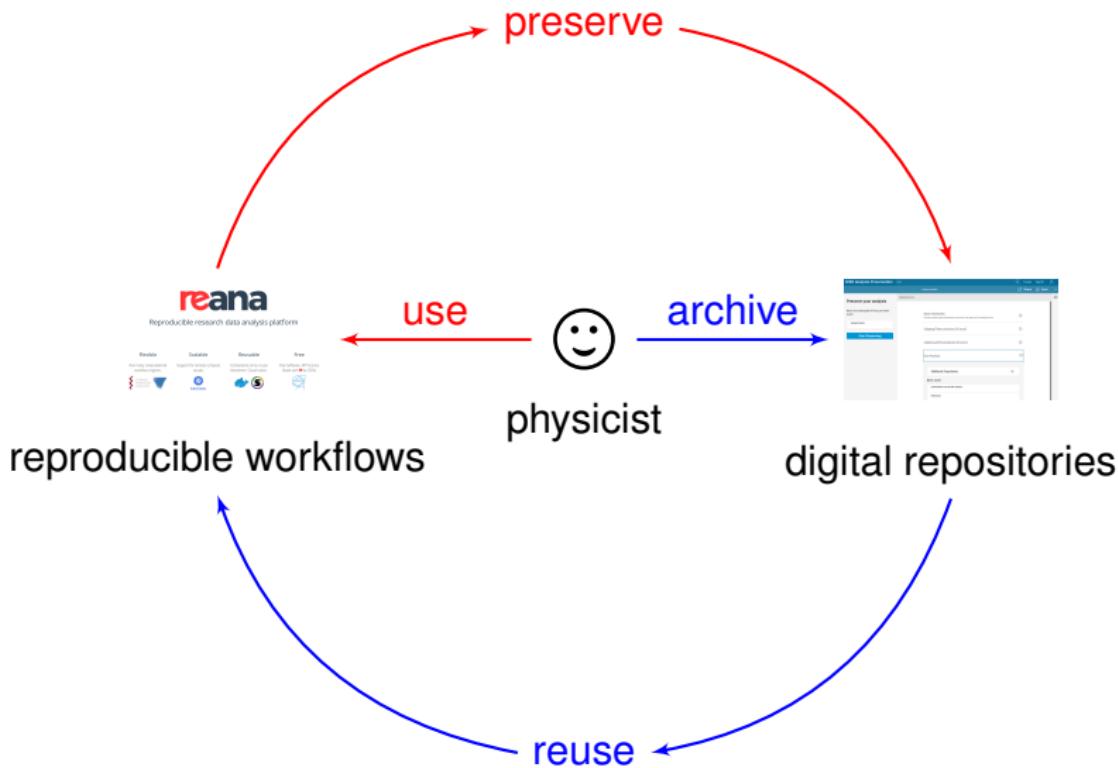
ATLAS RECAST



CMS Higgs-to-four-leptons

Several physics analysis examples

Reproducibility \rightleftharpoons Preservation



Fostering FAIR science



CERN Open Data



CERN Analysis Preservation

[Log in with CERN](#)

reana

Reproducible research data analysis platform

Flexible

Run many computational workflow engines.



Scalable

Support for remote computer clouds.



Reusable

Convenience once, reuse elsewhere. Cloud native.



Free

Free Software. MIT license. Made with ❤ at CERN.



REANA Reusable Analyses

“Capturing, preserving and sharing FAIR data and actionable knowledge behind particle physics data analyses in order to facilitate future data reuse”

CERN IT R. Maciulaitis, J. Okraska, D. Rodriguez, T. Šimko · **CERN SIS** S. Feger, P. Fokianos, A. Lavasa, S. van de Sandt, A. Trzcińska · **ALICE** Y. Foka, M. Gheata, C. Grigoras, M. Zimmermann · **ATLAS** K. Cranmer, L. Heinrich, A. Sanchez Pineda, D. Rousseau, F. Socher · **CMS** H. Bittencourt, A. Calderon, E. Carrera, A. Geiser, A. Huffman, C. Lange, K. Lassila-Perini, L. Lloret, T. McCauley, A. Rao, A. Rodriguez Marrero · **LHCb** S. Amerio, C. Burr, B. Couturier, S. Neubert, C. Parkes, S. Roiser, A. Trisovic · **OPERA** G. De Lellis, S. Dmitrievsky · **CERN CernVM** J. Blomer · **CERN EOS** L. Mascetti, H. Rousseau · **CERN Kubernetes** R. Rocha · **CERN OpenShift** A. Lossent

Conclusions

- opening complex LHC physics data beyond education use cases
 - enables independent physics research
 - validates and enriches knowledge preservation techniques
- open data is not enough
 - code, environments, runnable recipes and analysis workflows
 - science is born naturally ‘closed’; FAIR principles apply early
- “top-down” approach
 - funding agency requirements for preservation and openness
 - collaboration best practices for better reproducibility
- “bottom-up” activities
 - building tools to fit scientists’ daily workflow
 - pre-producibility rather than re-producibility
- time is ripe for driving the change
 - technological progress meets sociological challenges

Further reading 1: CERN Courier



CERN Courier March/April 2019

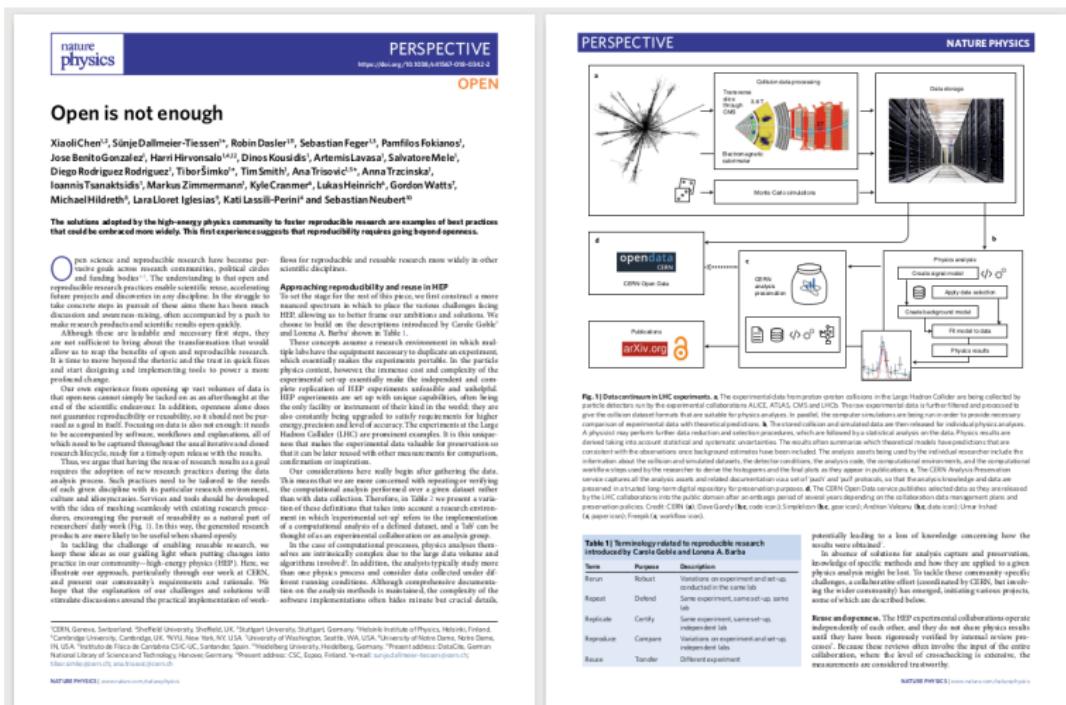
“The rise of Open Science”

Features:

- open access
- open data
- open source
- open science

<https://cerncourier.com/>

Further reading 2: Nature Physics



<https://www.nature.com/articles/s41567-018-0342-2.pdf>

References



CERN Open Data

- <http://opendata.cern.ch>
- <http://github.com/cernopendata>
- [cernopendata](#)



CERN Analysis Preservation

- <http://analysispreservation.cern.ch>
- <http://github.com/cernanalysispreservation>
- [analysispreserv](#)



REANA

- <http://www.reanahub.io>
- <http://github.com/reanahub>
- [reanahub](#)



Invenio

- <http://inveniosoftware.org>
- <http://github.com/inveniosoftware>
- [inveniosoftware](#)



Zenodo

- <https://zenodo.org>
- <http://github.com/zenodo>
- [zenodo_org](#)