# Open Science at CERN
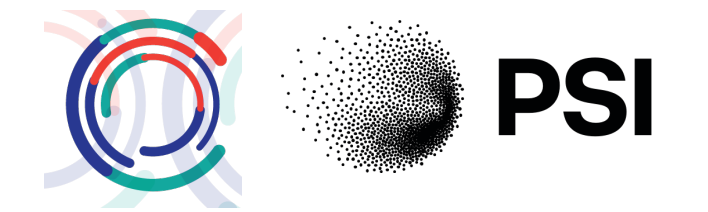
... using the CMS experiment as leading example

openscience.cern

**Clemens Lange (Paul Scherrer Institute PSI)**
EuroLabs Advanced Training: Open Science and Data Management
26th November 2024

# Introducing the CMS Experiment at CERN

Video: https://www.youtube.com/watch?v=EB5eZIR3AoM

# Hi, I'm Clemens :-)



I'm a particle physicist at Paul Scherrer Institute, working on the CMS experiment at the Large Hadron Collider (LHC) at CERN
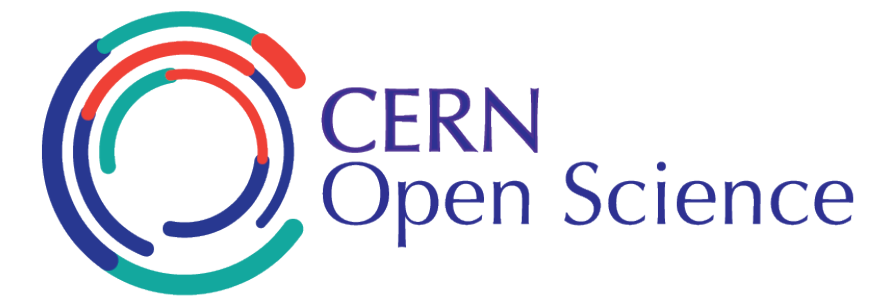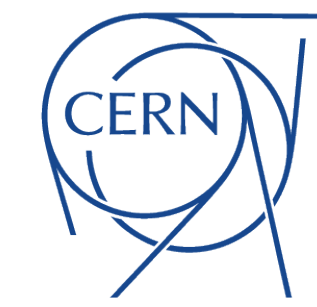
I roughly share my work time as follows

> Analysing the particle collisions provided by the LHC, recorded by the CMS detector (30-50%)

> Building and testing new pixel detectors for the upgrade of the CMS experiment (40-60%)

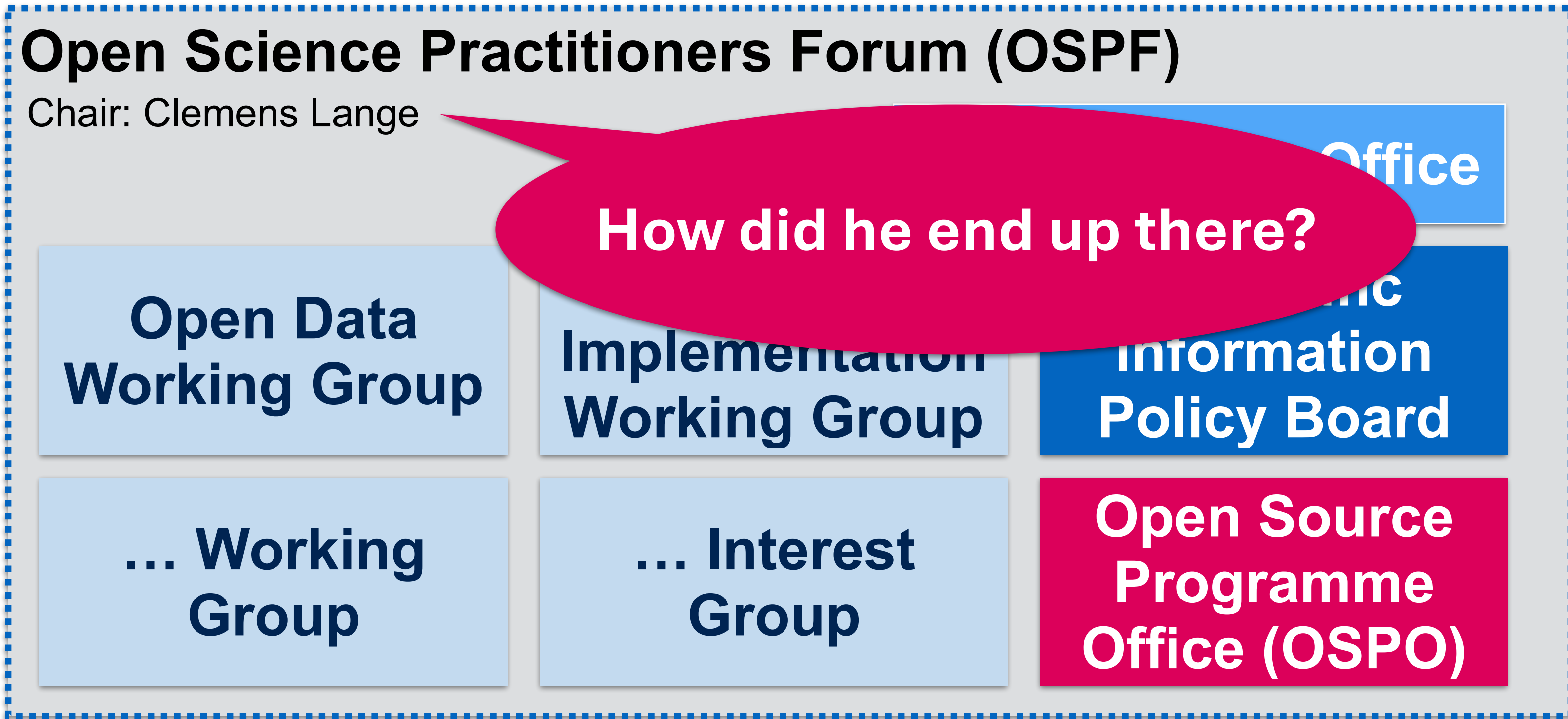> Detector operations (5%)

> Open data/computing (~10%)

**???**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN          26.11.2024

# Open Science Governance at CERN

**Director General | Director of Research and Computing**

**Open Science Steering Board**

**Open Science Practitioners Forum (OSPF)**

Chair: Clemens Lange

**How did he end up there?**

Office

**Open Data Working Group**

**Implementation Working Group**

Information Policy Board

**… Working Group**

**… Interest Group**

**Open Source Programme Office (OSPO)**

openscience.cern
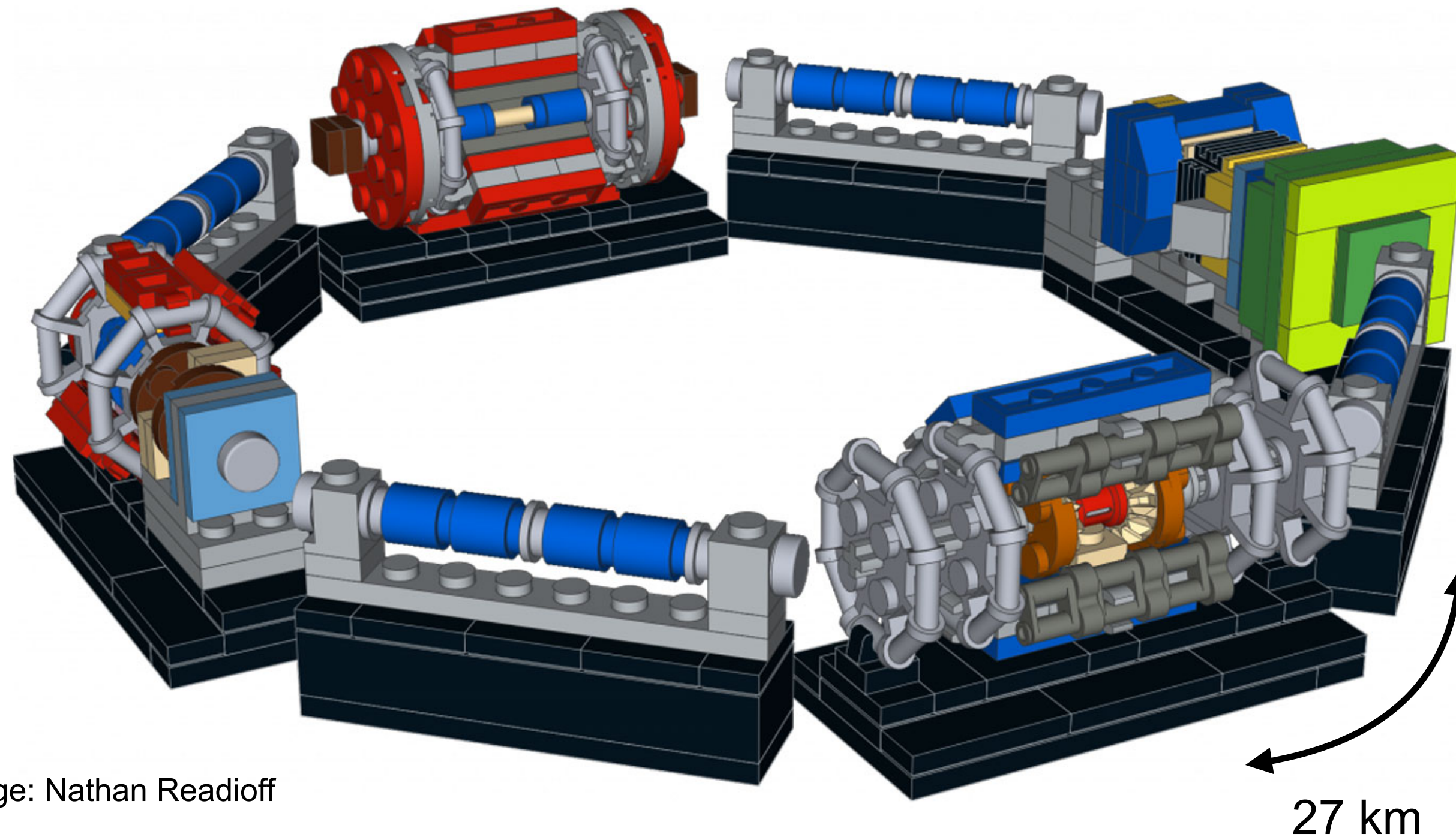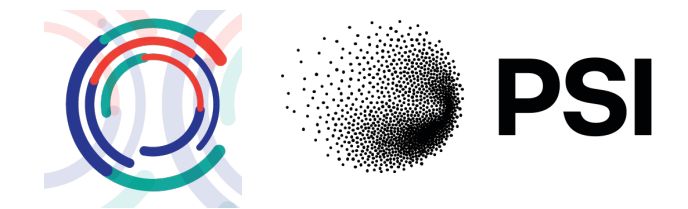
# High-Energy Physics at the CERN LHC



Image: Nathan Readioff

27 km

Four large experiments:

> ATLAS (5500 members of which almost 3000 scientific authors)

> ALICE (almost 2000 members)

> CMS (4000 particle physicists, engineers, computer scientists, technicians and students)

> LHCb (about 1700 scientists, engineers and technicians)

... plus several smaller ones

**Today: more than 13,000 people involved in the experiments**

26.11.2024

# The CMS Experiment

**CMS DETECTOR**

Total weight : 14,000 tonnes
Overall diameter : 15.0 m
Overall length : 28.7 m
Magnetic field : 3.8 T

**STEEL RETURN YOKE**
12,500 tonnes

**SILICON TRACKERS**
Pixel (100x150 μm) ~1m² ~124M channels
Microstrips (80x180 μm) ~200m² ~9.6M channels

**SUPERCONDUCTING SOLENOID**
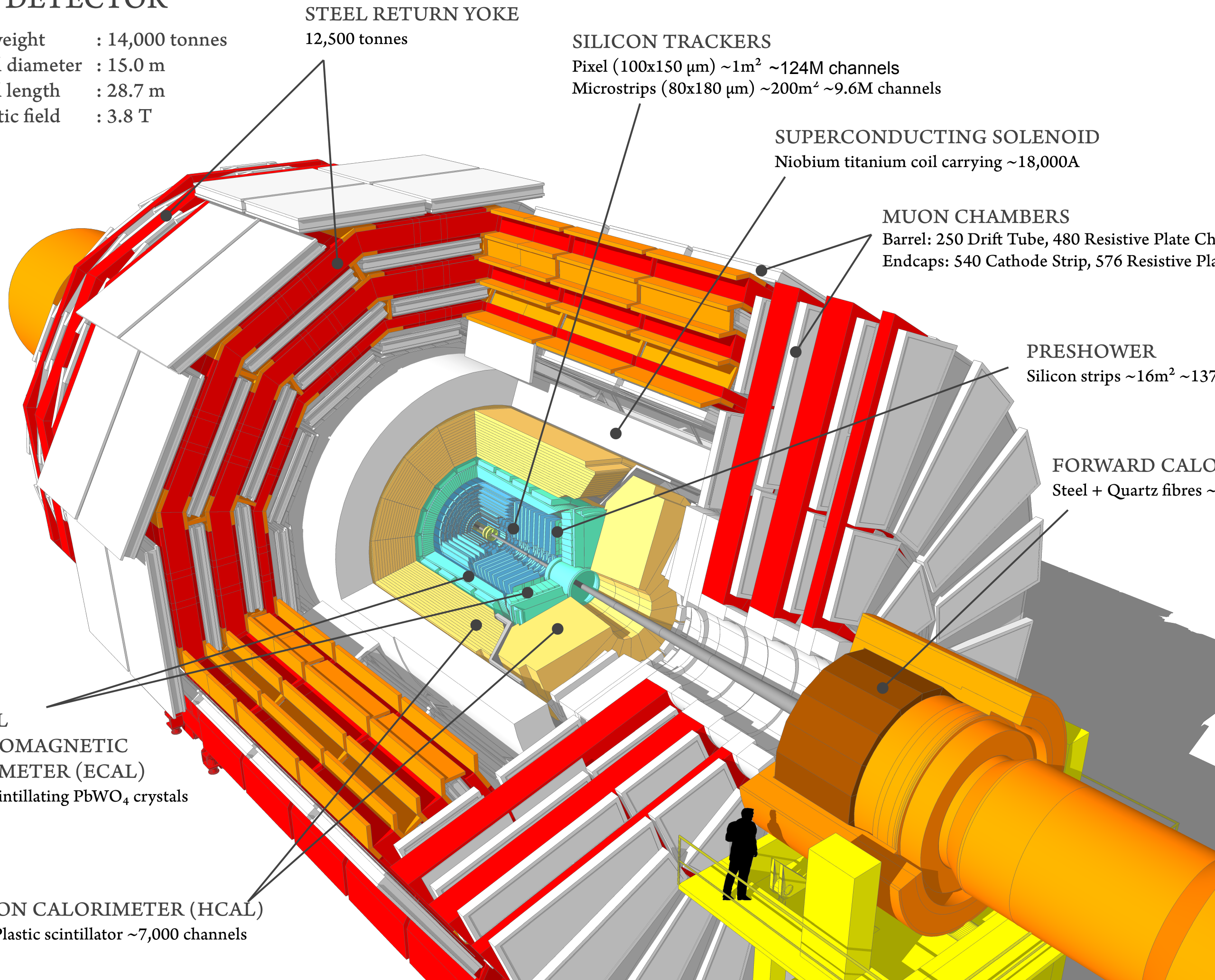Niobium titanium coil carrying ~18,000A

**MUON CHAMBERS**
Barrel: 250 Drift Tube, 480 Resistive Plate Chambers
Endcaps: 540 Cathode Strip, 576 Resistive Plate Chambers

**PRESHOWER**
Silicon strips ~16m² ~137,000 channels

**FORWARD CALORIMETER**
Steel + Quartz fibres ~2,000 Channels

**CRYSTAL ELECTROMAGNETIC CALORIMETER (ECAL)**
~76,000 scintillating $PbWO_4$ crystals

**HADRON CALORIMETER (HCAL)**
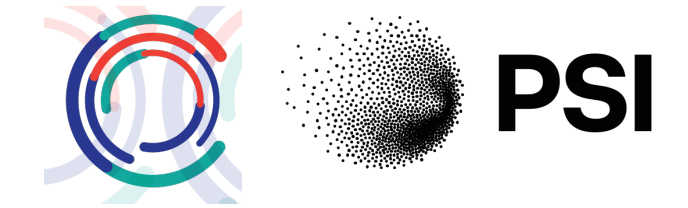Brass + Plastic scintillator ~7,000 channels

Record up to **40,000,000 events** of the LHC collisions **per second**, 24/7 (almost) all year long

Goal: understand the smallest building blocks of matter

**~134 million readout channels** — extraordinary levels of technical sophistication

**These data are unique, e.g. the Higgs boson can only be measured at the LHC**

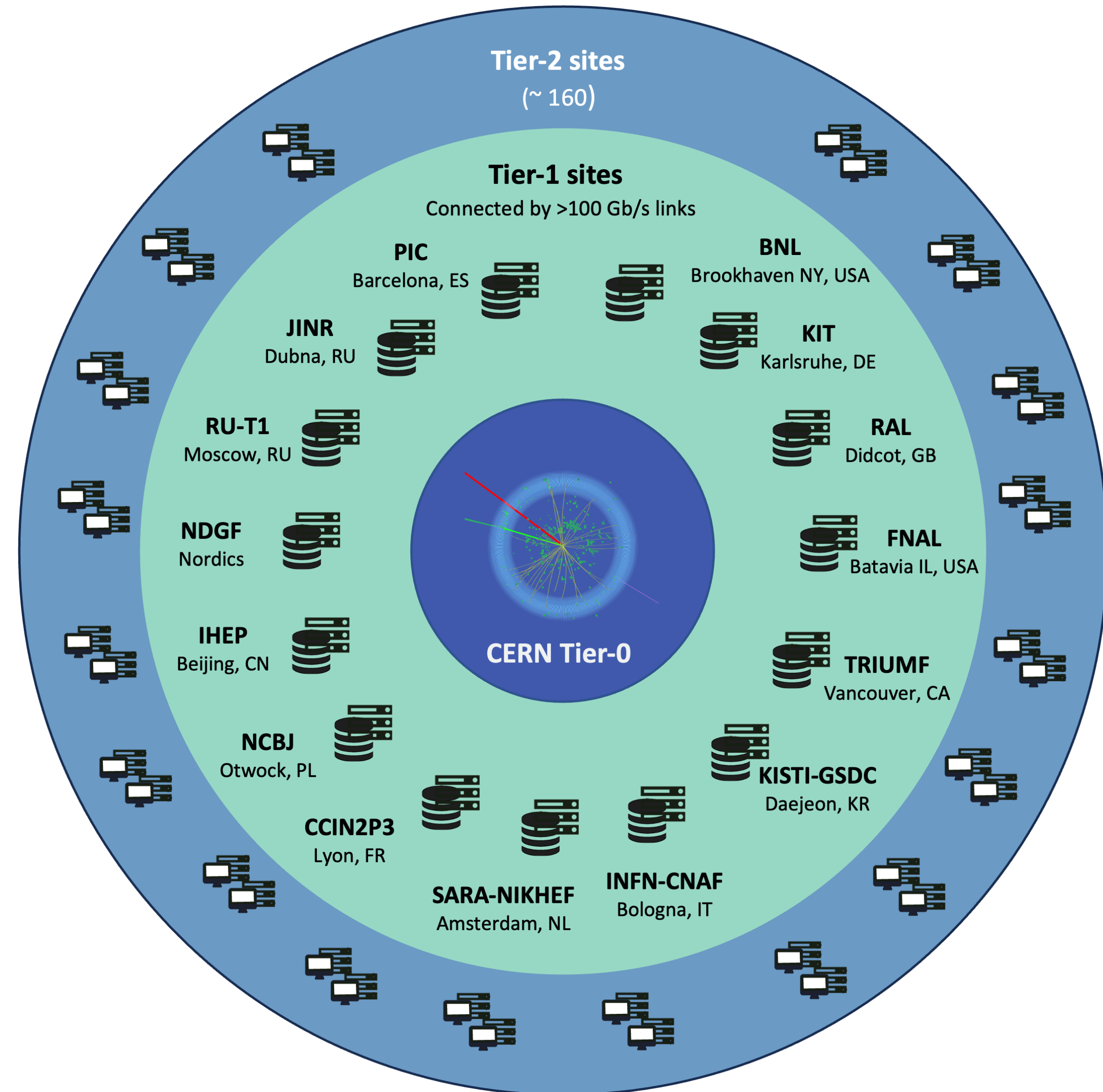# Data Processing using the Worldwide LHC Computing Grid

Centralised processing takes place in an **automated** way on the Worldwide LHC Computing Grid (> 170 computer centres, > 1 million computer cores, 2 exabytes of storage)
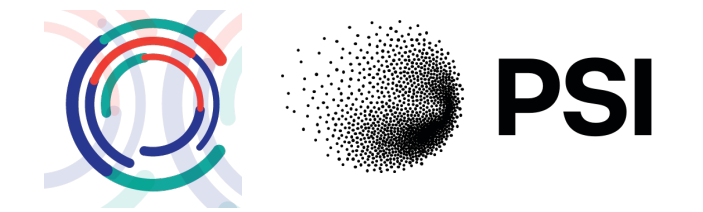
Each experiment has its own computing model

CMS: container images containing base Linux operating system with experiment software served through CernVM file system

**This is what makes things easy**



Source: https://wlcg-public.web.cern.ch/tiers

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

26.11.2024

# CMS Publications

Show all | Total | Exotica | Standard Model | Supersymmetry | Higgs | Top | Heavy Ions

B and Quarkonia | Forward and Soft QCD | Beyond 2 Generations | Detector Performance

**1332 collider data papers submitted as of 2024-11-05**

Discovery of the Higgs boson (No. 183)

Interactive version at http://cms-results.web.cern.ch/cms-results/public-results/publications-vs-time/

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

26.11.2024

# Current State of Analysis Preservation in CMS (1)

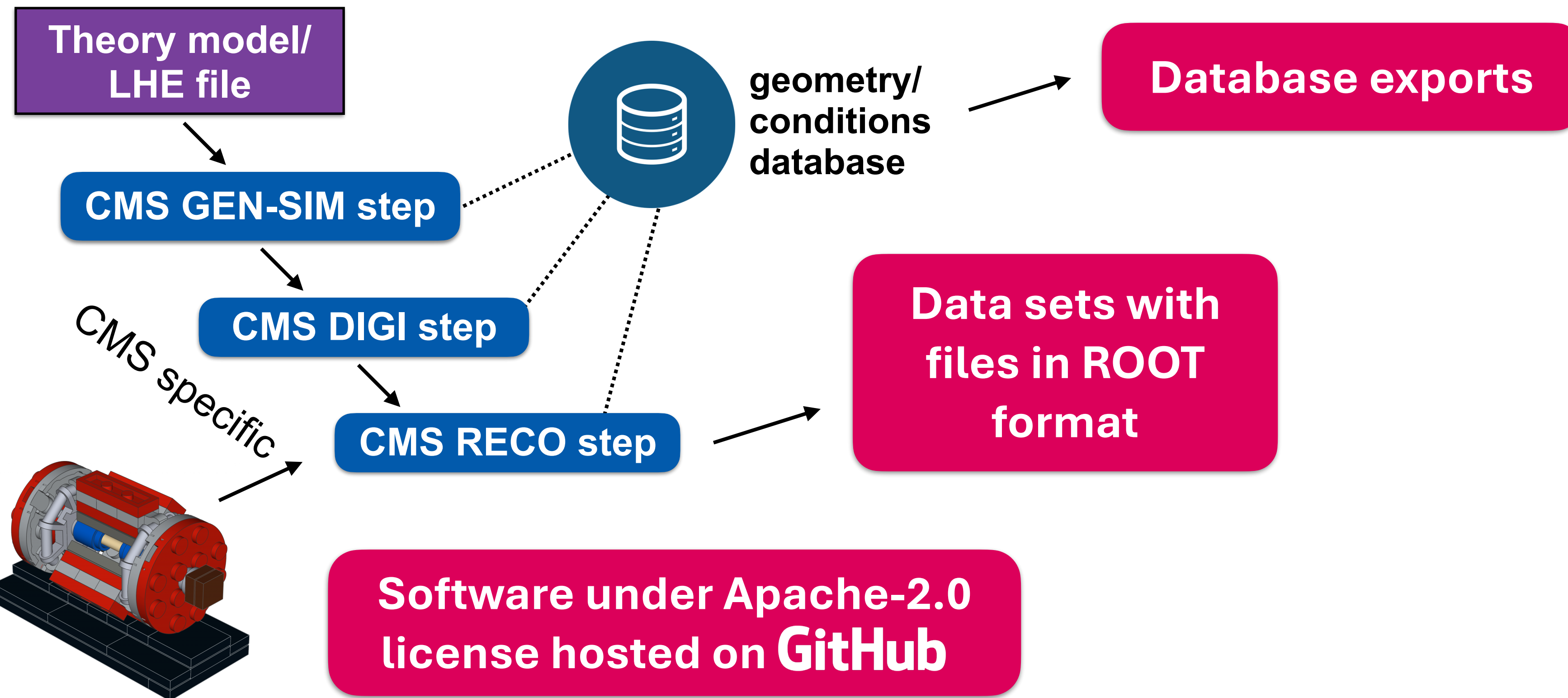Event generation+simulation as well as reconstruction (both data and MC) centralised

> Software and database tags preserved and archived



**Theory model/ LHE file**

**CMS GEN-SIM step**

geometry/ conditions database

**Database exports**

CMS specific

**CMS DIGI step**

**Data sets with files in ROOT format**

**CMS RECO step**

**Software under Apache-2.0 license hosted on GitHub**

Experiment data

**Reproducible**

# CMS Data Aggregation Service (internal)



All data sets created using grid resources are registered

> Naming is important!

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN        26.11.2024

# From Internal CMS Data Sets to Open Data

A lot of information encoded in a data set:

$m_{Higgs}$

centre-of-mass energy

process

underlying event tune

event generator(s)

Processing campaign

/ttHTobb_M125_TuneCP5_13TeV-powheg-pythia8/

RunIISummer20UL16NanoAODv9-106X_mcRun2_asymptotic_v17-v2/

NANOAODSIM

Data tier

**This is the hard part**

There is **a lot more metadata that is not visible** though, added to public records such as https://opendata.cern.ch/record/67645 — **can trace/reproduce entire chain** — often used by CMS collaborators now!

Open Data are eventually transferred to public EOS instance at CERN

Corresponding records published on CERN Open Data Portal (with DOIs)

# 10 Years of CMS Open Data and the CERN Open Data Portal

Since 2014, CMS has released ~4.5 petabytes of open data available on the <u>CERN Open Data Portal</u>

> Both collision and simulation data sets

> Entire Run-1 (2010-2012) + 2015 data sets, fraction of 2016

More information: <u>https://cms.cern/news/cms-celebrates-decade-open-data</u>

Clemens Lange (<u>clemens.lange@cern.ch</u>) — Open Science at CERN          26.11.2024

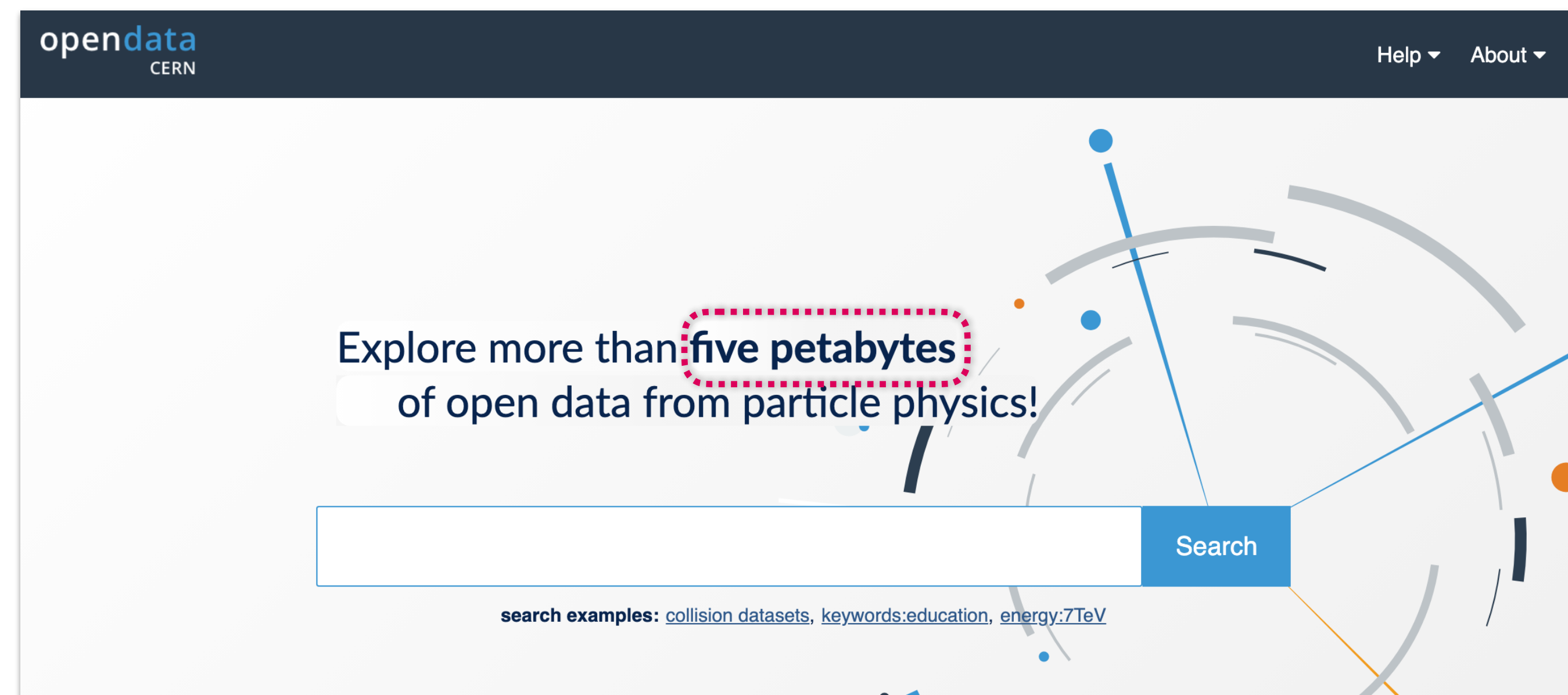# From CMS to CERN Open Data



At the end of 2020, all large LHC experimental collaborations have endorsed a <u>new open data policy</u>

> Following existing CMS policy

Commit to publicly **releasing data required to ma**

Data and simula released approximately five years after collection (50%)

> Released und Commons CC

> Full dataset by the close of the experiment

**ATLAS and LHCb (0.8 PB) recently released their first research-level Open Data**

**In 2022, developed CERN Open Science Policy and its Implementation Plan**

> Level 1: Open access publication and additional

utreach and Education

and the software to analyse them

> Level 4: Raw data, and the software to reconstruct and analyse them

higher comp

# Challenges of Providing Open Data

## Data: available ≠ usable

Open Data needs to be FAIR:

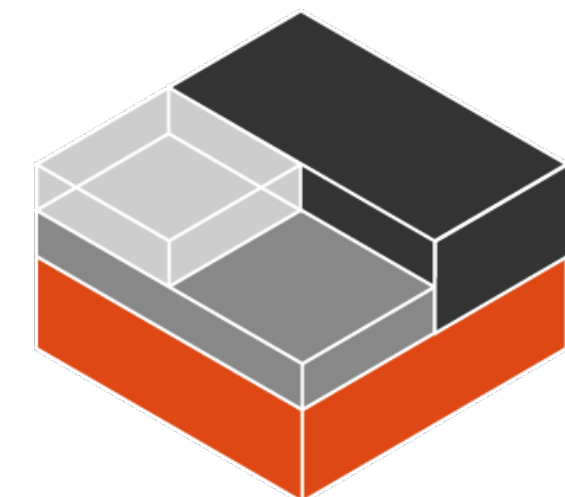**F**indable ➡ CERN Open Data Portal records

**A**ccessible ➡ reliable storage and access technology

**I**nteroperable ➡ provide good documentation, avoid jargon

**R**eusable ➡ **preserve software** (and hardware to run it if needed), data provenance, workflows

**Building CMSSW container (Docker) images got me involved in CMS Open Data**

26.11.2024

# Getting Others to Use Your Open Data

Beyond the data sets available on the CERN Open Data Portal, we provide:

Analysis examples with different levels of complexity (scientific and education)

The required software

A separate CMS Open Data Guide

> In particular, trying to explain **how to use** the data and **what to do** with them in addition to **what is** in the data

Workshops with Software Carpentry style tutorials:

> 2020 CMS Open Data Workshop for Theorists
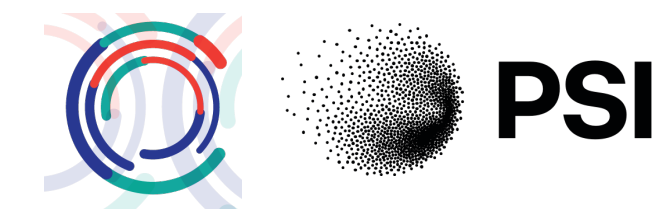
> 2021 CMS Open Data Workshop

> 2022 CMS Open Data Workshop at CERN

> 2023 CMS Open Data Workshop at Fermilab LPC

> 2024 CMS Open Data Workshop & Hackathon at CERN



We've studied this data...

Now it's your turn to explore!

Photo Illustration: CNN/Adobe Stock/Universal Pictures/Warner Bros. Pictures

CMS Open Data Workshop & Hackathon

July 29th - Aug 1st, 2024    CERN IdeaSquare

# CMS Publications



Discovery of the Higgs boson (No. 183)

OD results: **Equivalent of a new working group** (but we are not authors)

Interactive version at http://cms-results.web.cern.ch/cms-results/public-results/publications-vs-time/

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

26.11.2024

# Beyond Open Data

# Current State of Analysis Preservation in CMS (2)

**PSI**

Event generation+simulation as well as reconstruction (both data and MC) centralised

**>** Software and database tags preserved and archived

Internal documentation (analysis notes) preserved

Paper publications **open access** (since 2014 under SCOAP³), Preprints available on arXiv



Theory model/ LHE file

geometry/ conditions database

CMS GEN-SIM step

CMS specific

CMS DIGI step

CMS RECO step

The physics analysis™

Internal documentation

**HEPData**

Experiment data

**Reproducible**

**Preserved**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN                    26.11.2024

# HEPData

Paper publications typically in PDF format — nice to read and print, but challenging to extract actual data

Goal: provide figures and tables in tabulated, machine-readable format

Can be used for comparisons, reinterpretations etc.



doi.org/10.17182/hepdata.102646

26.11.2024

# Increasing HEPData Adoption

LHC publications with HEPData records (2022-11-25)

Source

CMS has HEPData records for almost all analyses since 2022

How?

In 2018, hepdata_lib was born: "[Python] Library for getting your data into HEPData" → provide users with **good tooling**, show that it is **easy**, **make it mandatory** in the collaboration (2021)

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN
26.11.2024

# Open Software and Hardware

PSI

# Side Note: Open-Sourcing Software and Hardware

CERN has recently established an Open Source Programme Office (OSPO)

Goal is to support the CERN community in the process of **making internal projects public**, e.g.:

> Identify and apply suitable license

> Guidelines for project maintenance and support

Also, provide a public catalogue of CERN's open source projects

More information: opensource.cern

Adapted from https://xkcd.com/2347/

Adapted from https://xkcd.com/2347/

# Preserving Physics Analyses

PSI

# Current State of Analysis Preservation in CMS (2)

Event generation+simulation as well as reconstruction (both data and MC) centralised

**>**Software and database tags preserved and archived

Internal documentation (analysis notes) preserved

Paper publications **open access** (since 2014 under SCOAP³), Preprints available on arXiv

**Theory model/ LHE file**

**geometry/ conditions database**

**CMS GEN-SIM step**

CMS specific

**CMS DIGI step**

**CMS RECO step**

Experiment data

**The physics analysis™**

**Details unknown**

**Reproducible**

Internal documentation

PUBLISHED

HEPData

**Preserved**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN                26.11.2024

# Helping Your Future You

PSI

*"Your closest collaborator is you six months ago… and your younger self doesn't reply to emails"* → **preserving your analysis pipelines will help you in your immediate future**

Theory model/ LHE file

geometry/ conditions database

**Your work of the past N years?**

CMS GEN-SIM step

CMS specific

CMS DIGI step

CMS RECO step

Experiment data

ntuplisation/selection

*Analysis specific*

Statistical analysis

Internal documentation

HEPData

PUBLISHED

**Reproducible**

**Currently not preserved**

**Preserved**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

26.11.2024

# Steps towards Reusable Analyses

| 1. Capture software | 2. Capture commands | 3. Capture workflow |
|---|---|---|
| Individual analysis stages in an executable way (including all dependencies) | How to run the captured software? | How to connect the individual analysis steps? |

Capturing analysis code almost trivial today

Requires e.g. two additional files in a GitLab repository → something you will learn this week

# CMS Open Data simplified analysis example



https://opendata.cern.ch/record/5500

# Steps towards Reusable Analyses

**1. Capture software**

Individual analysis stages in an executable way (including all dependencies)

**2. Capture commands**

How to run the captured software?

**3. Capture workflow**

How to connect the individual analysis steps?

Capturing analysis code almost trivial today

Requires e.g. two additional files in a GitLab repository → something you will learn this week

# Why Did the Demo not Work Everywhere?



**LHC / HL-LHC Plan**

Today

| | LHC | | HL-LHC |

| Run 1 | Run 2 | Run 3 | Run 4 - 5... |

LS1 — splice consolidation, button collimators, R2E project
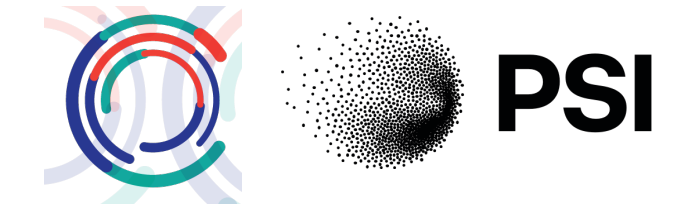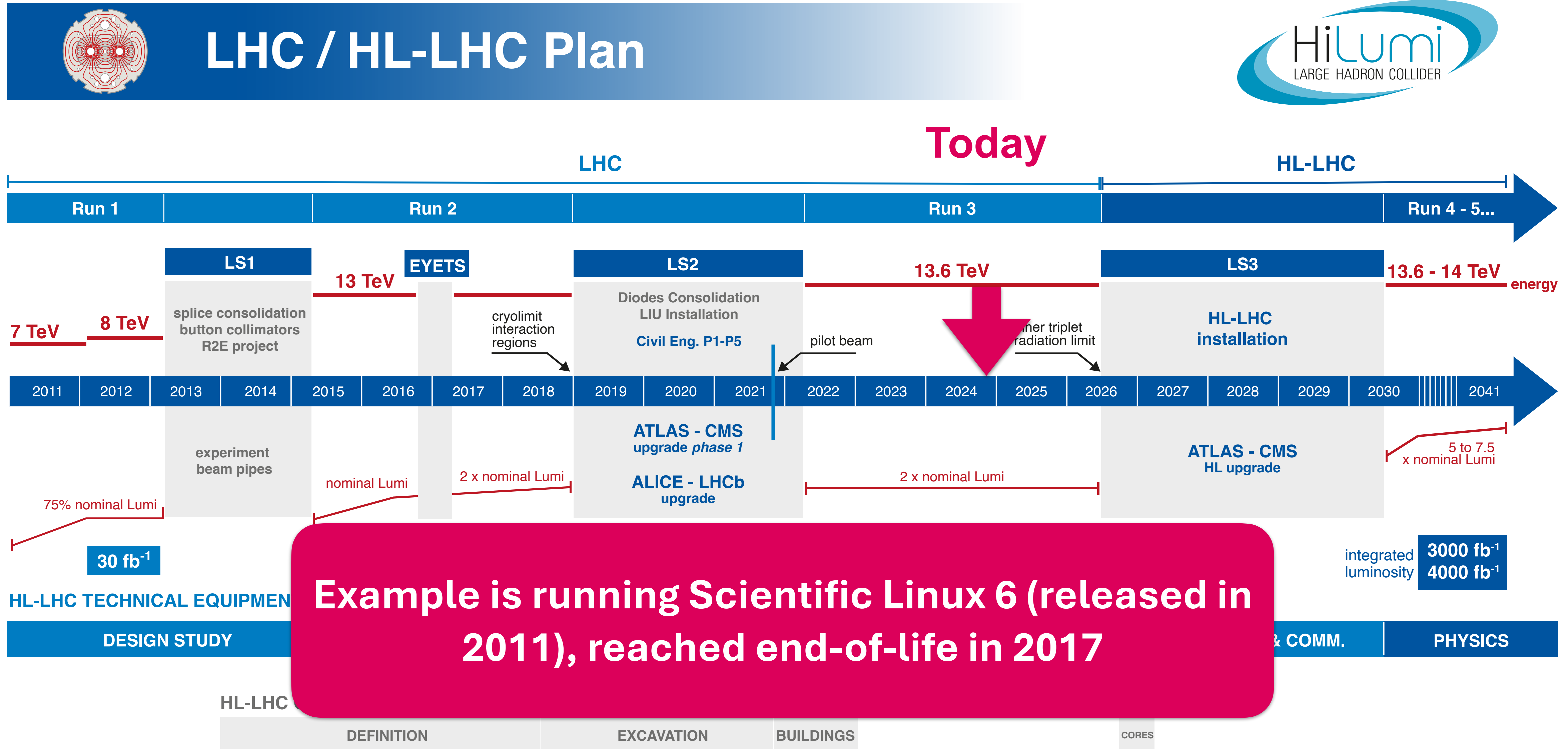EYETS
LS2 — Diodes Consolidation, LIU Installation, Civil Eng. P1-P5
13 TeV
13.6 TeV
LS3 — HL-LHC installation
13.6 - 14 TeV energy

7 TeV
8 TeV

cryolimit interaction regions
pilot beam
inner triplet radiation limit

| 2011 | 2012 | 2013 | 2014 | 2015 | 2016 | 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 | 2025 | 2026 | 2027 | 2028 | 2029 | 2030 | 2041 |

experiment beam pipes

ATLAS - CMS upgrade *phase 1*
ALICE - LHCb upgrade

ATLAS - CMS HL upgrade

75% nominal Lumi
nominal Lumi
2 x nominal Lumi
2 x nominal Lumi
5 to 7.5 x nominal Lumi

30 fb$^{-1}$
integrated luminosity: 3000 fb$^{-1}$ / 4000 fb$^{-1}$

HL-LHC TECHNICAL EQUIPMENT

DESIGN STUDY

**Example is running Scientific Linux 6 (released in 2011), reached end-of-life in 2017**

& COMM.

PHYSICS

HL-LHC

DEFINITION | EXCAVATION | BUILDINGS | CORES

Source: https://project-hl-lhc-industry.web.cern.ch/content/project-schedule

31    Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN                26.11.2024

# Tooling: Software Containers

Software containers enable portability of (compiled) code

They allow e.g. to compile and run old and recent CMSSW versions on today's operating systems and processor architectures (but also your analysis code from last year)

> "Works on my *and your* machines" — from laptop to batch/grid/cloud



Advantage: **You know exactly which version of your code is running**

> Ideally built automatically using continuous integration (e.g. GitHub/GitLab)

Also useful for analysis development in general (or e.g. DAQ software, machine learning, …)

# Steps towards Reusable Analyses

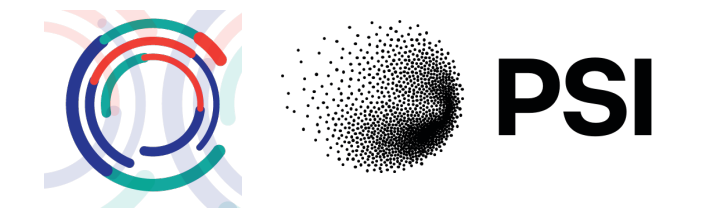| 1. Capture software | 2. Capture commands | 3. Capture workflow |
|---|---|---|
| Individual analysis stages in an executable way (including all dependencies) | How to run the captured software? | How to connect the individual analysis steps? |

Capturing analysis code almost trivial today

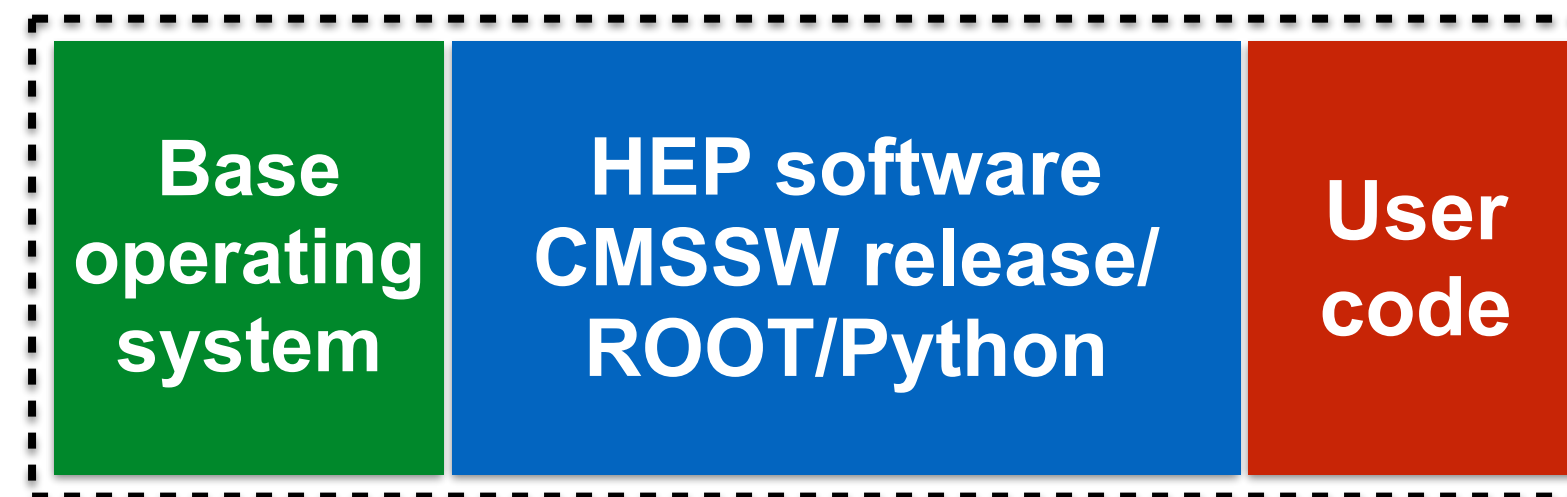Requires e.g. two additional files in a GitLab repository → something you will learn this week

**Once commands have been captured, can run individual analysis steps**

# Steps towards Reusable Analyses

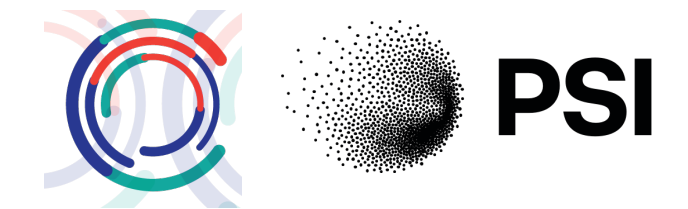| 1.Capture software | 2. Capture commands | 3. Capture workflow |
|---|---|---|
| Individual analysis stages in an executable way (including all dependencies) | How to run the captured software? | How to connect the individual analysis steps? |

Capturing analysis code almost trivial today

Requires e.g. two additional files in a GitLab repository → something you will learn this week

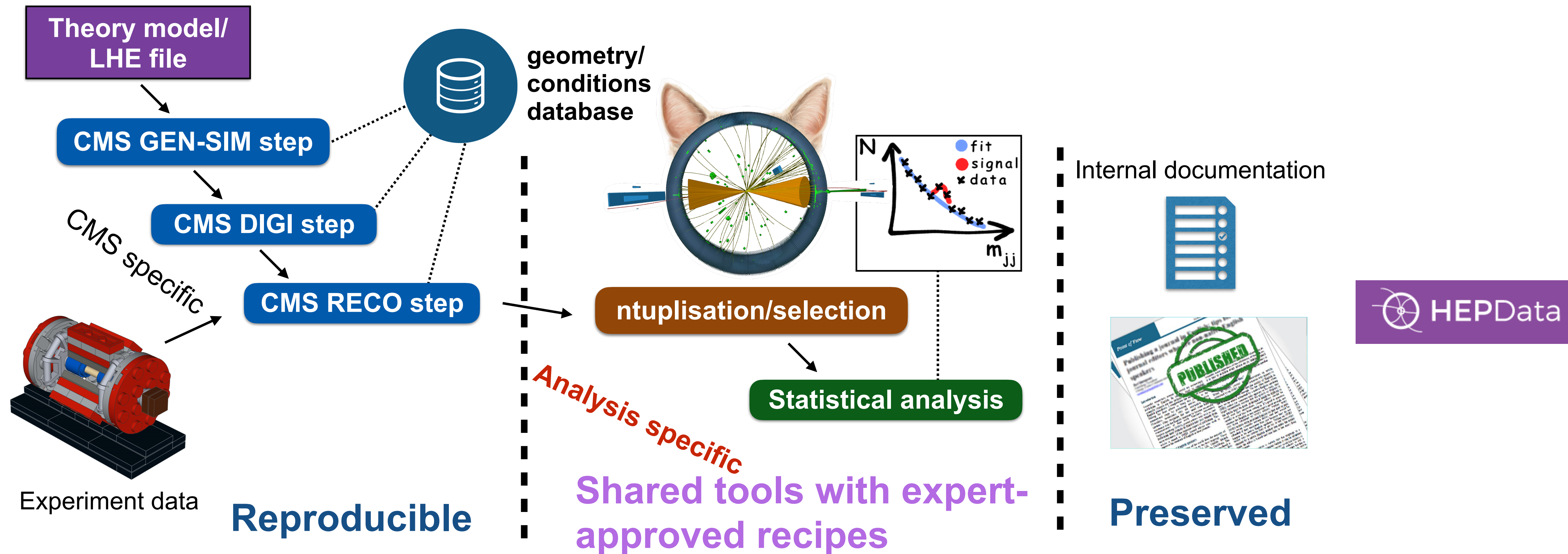Once commands have been captured, can run individual analysis steps

**Capturing the workflow can be achieved in various ways**

**Several tools exist, e.g. SnakeMake (available in REANA)**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN          26.11.2024
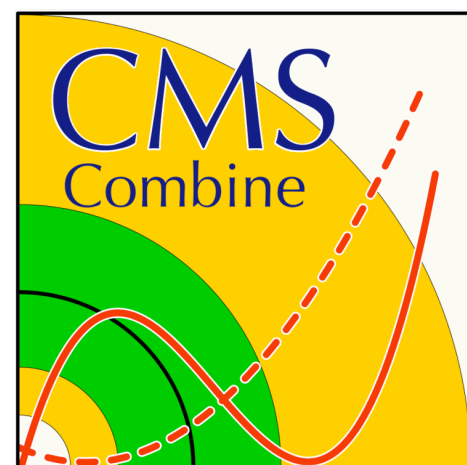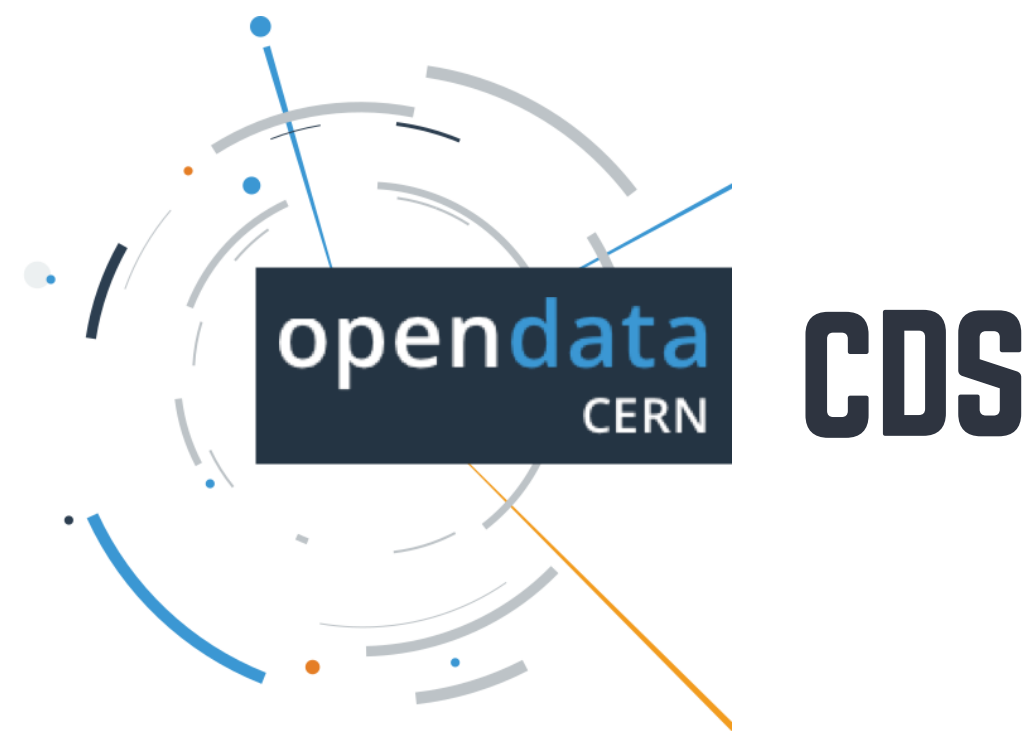
# Towards Common Analysis Tools in CMS

CMS established a new Common Analysis Tools (CAT) group at the end of 2022

This group is now working with various groups in CMS towards improved data processing tools, analysis workflows and their preservation as well as statistical inference tools (and much more)

**Theory model/ LHE file**

**geometry/ conditions database**

**CMS GEN-SIM step**

CMS specific

**CMS DIGI step**

**CMS RECO step**

N

fit
signal
data

$m_{jj}$

**ntuplisation/selection**

Analysis specific

**Statistical analysis**

Internal documentation

PUBLISHED

HEPData

Experiment data

**Reproducible**

**Shared tools with expert-approved recipes**

**Preserved**

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN
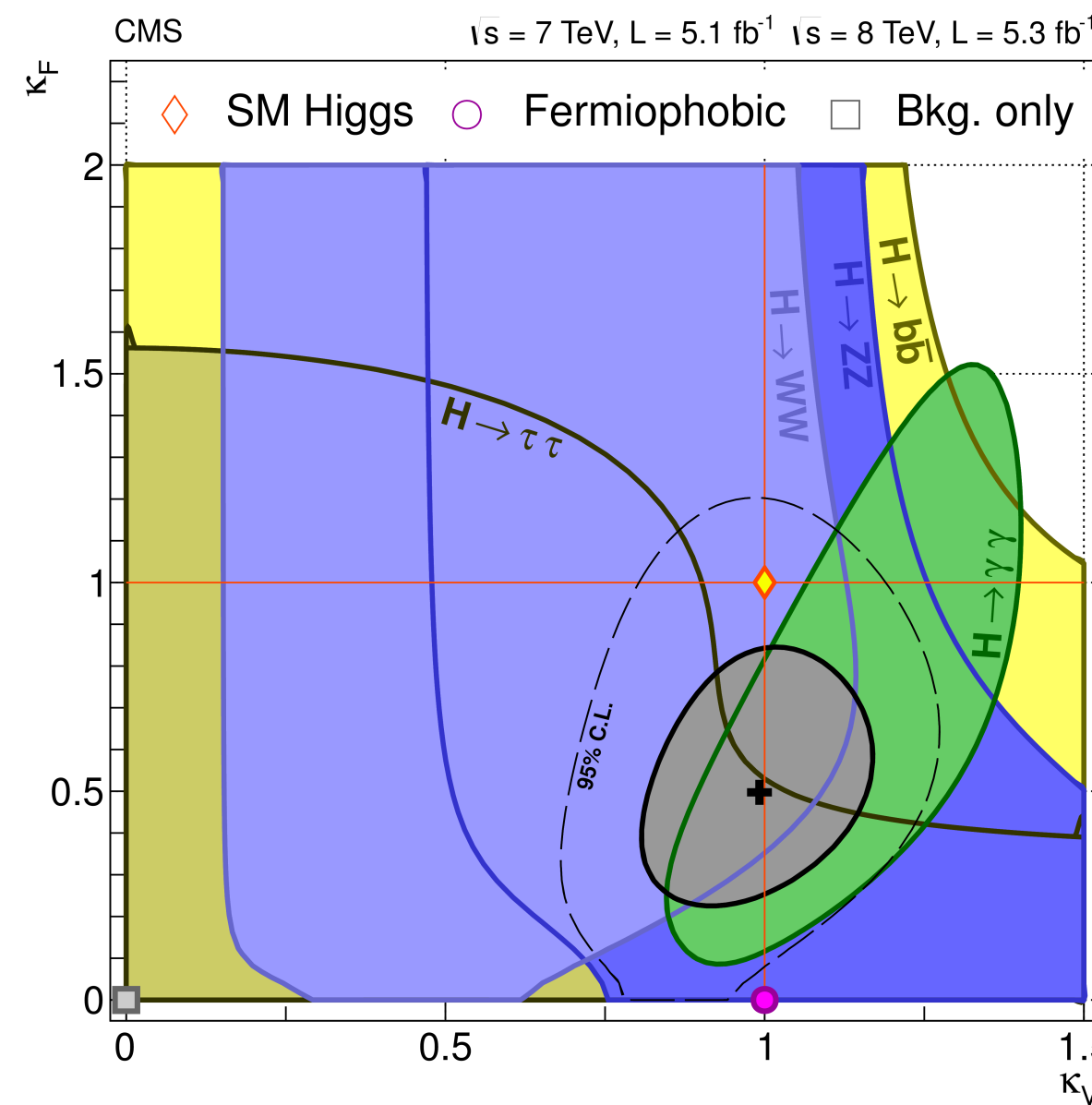
26.11.2024

# Open Science in Action

The CMS Collaboration recently released the full statistical model ("set of likelihood functions") of the measurements that contributed to establishing the **discovery of the Higgs boson** in 2012 including the required software
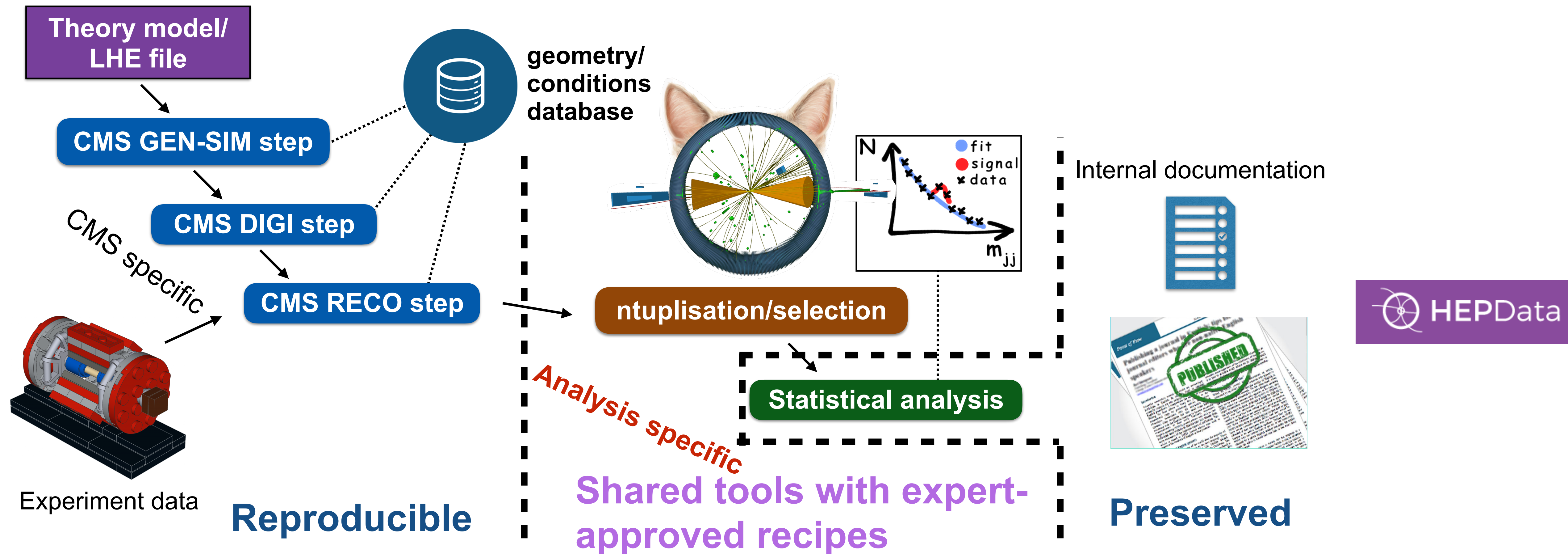
doi.org/10.17181/c2948-e8875



Open-Access Publications

Reusable physics analyses

Already being used outside CMS!

FOSS software

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

The new CMS CAT group works towards closing the gap in analysis preservation and reusability

However, **analysts need to be part of this change**



Theory model/LHE file

geometry/conditions database

CMS GEN-SIM step

CMS DIGI step

CMS specific

CMS RECO step

ntuplisation/selection

Analysis specific

Statistical analysis

Internal documentation

Experiment data

**Reproducible**

**Shared tools with expert-approved recipes**

**Preserved**

HEPData

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN

26.11.2024

Large collaborations move slowly, **policies help enforce standards — grassroots initiatives can make a difference**

Try to use **clear and understandable naming**, e.g. for data sets

Make **recurring/useful workflows** a **routine**/automate them

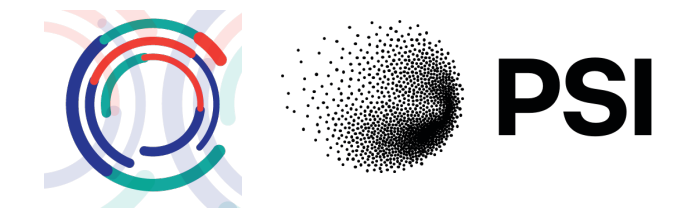Use **version control** — for collaboration with others ideally also have **tests**

Your computing environment might change within a few months time — **software containers provide portability**



HOW LONG CAN YOU WORK ON MAKING A ROUTINE TASK MORE EFFICIENT BEFORE YOU'RE SPENDING MORE TIME THAN YOU SAVE? (ACROSS FIVE YEARS)

|  | HOW OFTEN YOU DO THE TASK | | | | | |
|---|---|---|---|---|---|---|
|  | 50/DAY | 5/DAY | DAILY | WEEKLY | MONTHLY | YEARLY |
| 1 SECOND | 1 DAY | 2 HOURS | 30 MINUTES | 4 MINUTES | 1 MINUTE | 5 SECONDS |
| 5 SECONDS | 5 DAYS | 12 HOURS | 2 HOURS | 21 MINUTES | 5 MINUTES | 25 SECONDS |
| 30 SECONDS | 4 WEEKS | 3 DAYS | 12 HOURS | 2 HOURS | 30 MINUTES | 2 MINUTES |
| 1 MINUTE | 8 WEEKS | 6 DAYS | 1 DAY | 4 HOURS | 1 HOUR | 5 MINUTES |
| 5 MINUTES | 9 MONTHS | 4 WEEKS | 6 DAYS | 21 HOURS | 5 HOURS | 25 MINUTES |
| 30 MINUTES |  | 6 MONTHS | 5 WEEKS | 5 DAYS | 1 DAY | 2 HOURS |
| 1 HOUR |  | 10 MONTHS | 2 MONTHS | 10 DAYS | 2 DAYS | 5 HOURS |
| 6 HOURS |  |  |  | 2 MONTHS | 2 WEEKS | 1 DAY |
| 1 DAY |  |  |  |  | 8 WEEKS | 5 DAYS |

HOW MUCH TIME YOU SHAVE OFF

Source: https://xkcd.com/1205/

# On to You!

CERN has its first Open Science Policy plus an Implementation Plan

> Openly available and meant as "inspiration" for other institutions

The LHC experiments are making an effort to preserve larger parts of the physics analysis chain

> Whether this is successful will depend a lot on the analysts themselves

> There are several examples how Open Science made even internal collaboration better

I hope you will see the advantages of a more structured/systematic approach — this week's training will provide with the skills required
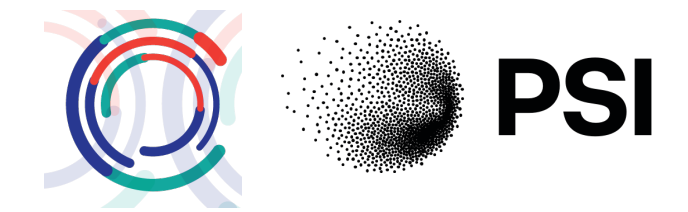
> Your future self will probably thank you

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN          26.11.2024

# New: CERN Open Science Policy

Captures current practice and states vision across multiple Open Science domains:

Open Access to Publications

Open Research Data

Open Software

Open Hardware

Citizen Science

Research Integrity, Reuse & Reproducibility

Infrastructure for Open Science

Research Assessment & Evaluation

Education, Training & Outreach

v1.0 released Oct 2022: https://cds.cern.ch/record/2835057

For more information, see https://openscience.cern/

**>Have a look at the implementation plan!**

> Data Management Plan template: https://openscience.cern/index.php/DMP

The **likelihood function** is a particularly special data product

> Small, information-dense, overall summary of the analysis

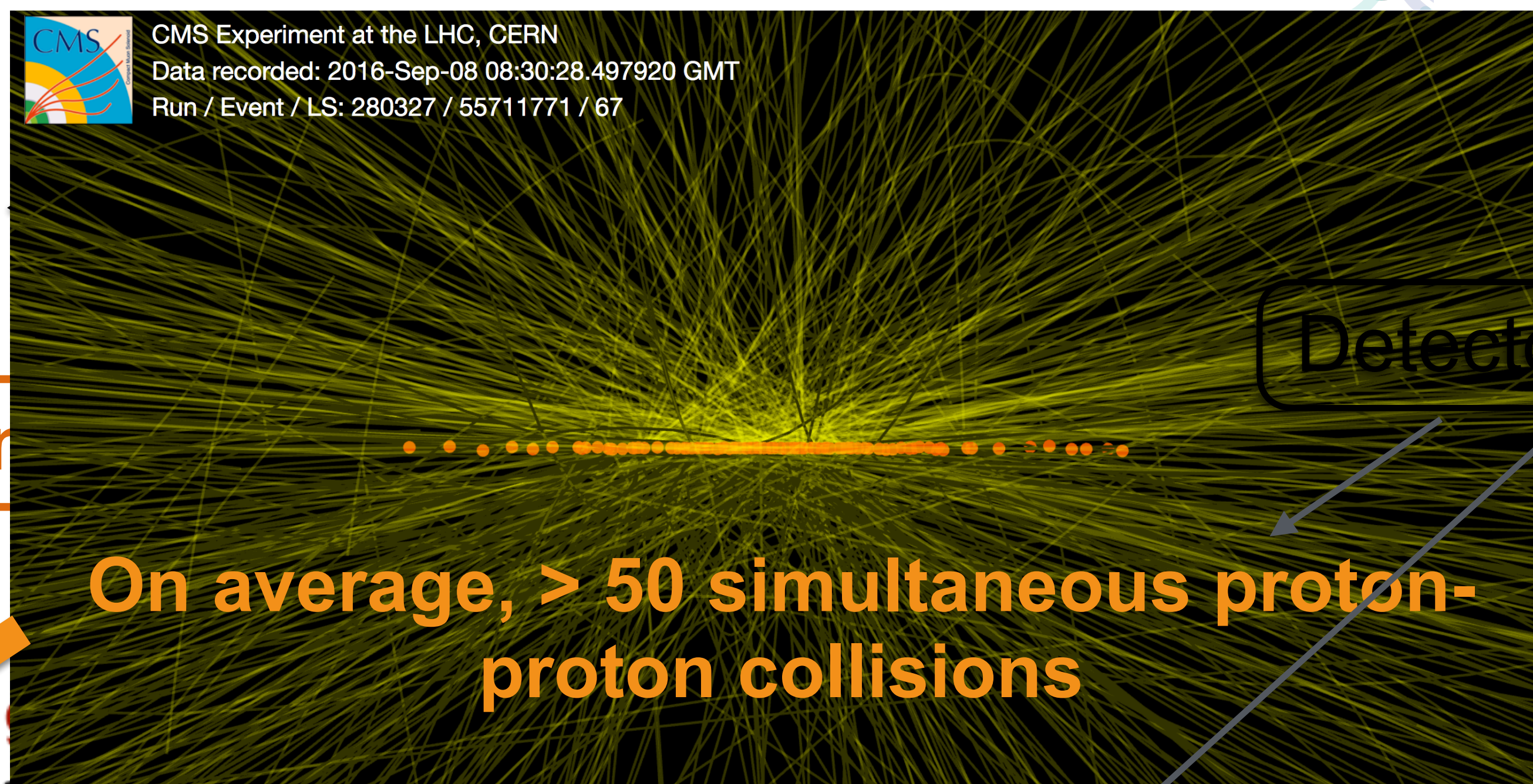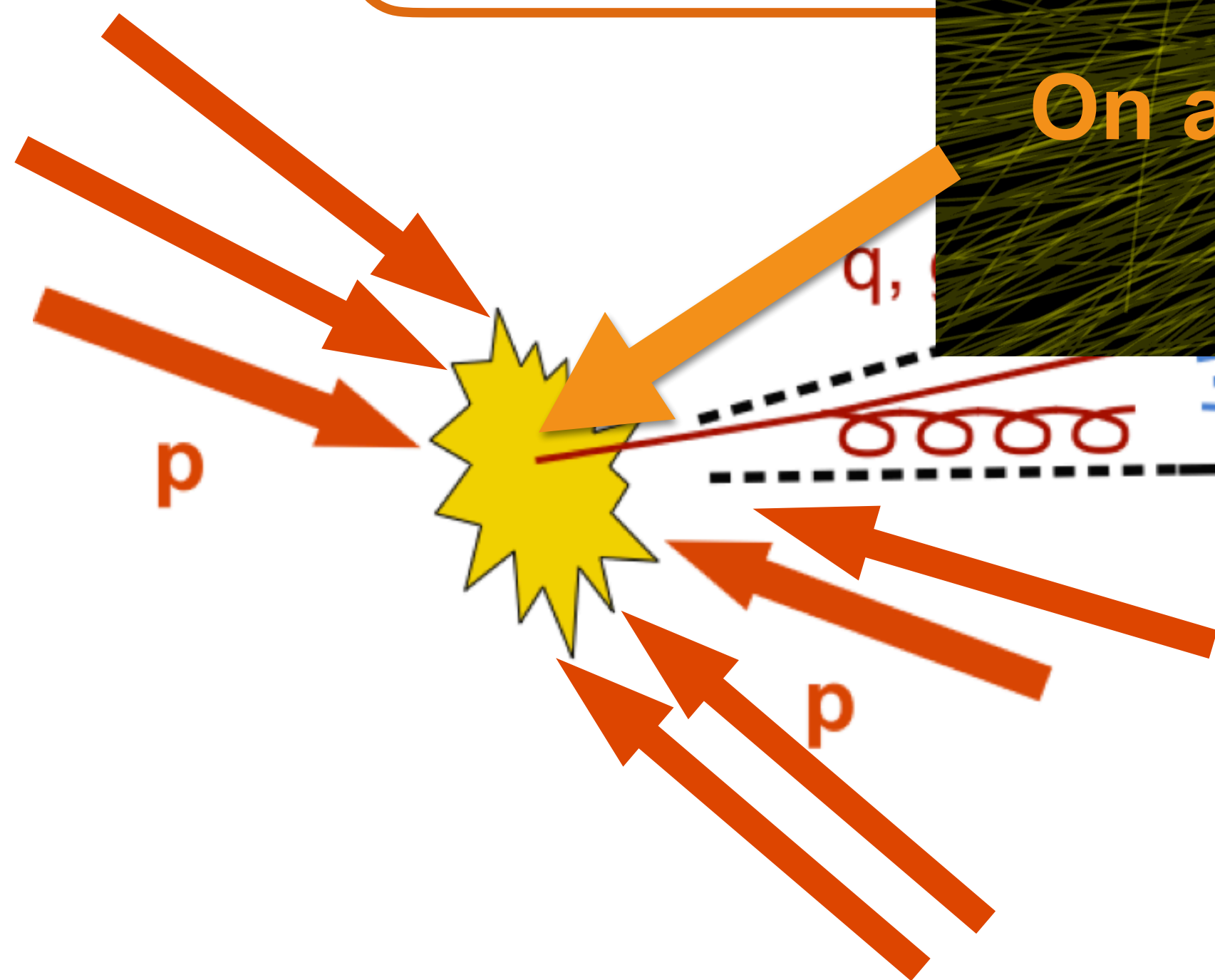> Almost every analysis decision is reflected in the likelihood



Expect to see more full likelihoods from CMS in the next few months

**Theory**
**(perturbation theory)**
**/ LHC pp collisions**

**Pileup+U...**

Detector noise

CMS Experiment at the LHC, CERN
Data recorded: 2016-Sep-08 08:30:28.497920 GMT
Run / Event / LS: 280327 / 55711771 / 67

**On average, > 50 simultaneous proton-proton collisions**

q,

p

p

π, K, ...

# Further challenges in high-energy physics data analysis



Collaborations make huge internal review effort (months to years) to **ensure accurate interpretation of the data**

> False claims (also from OD users) could risk erosion of public trust

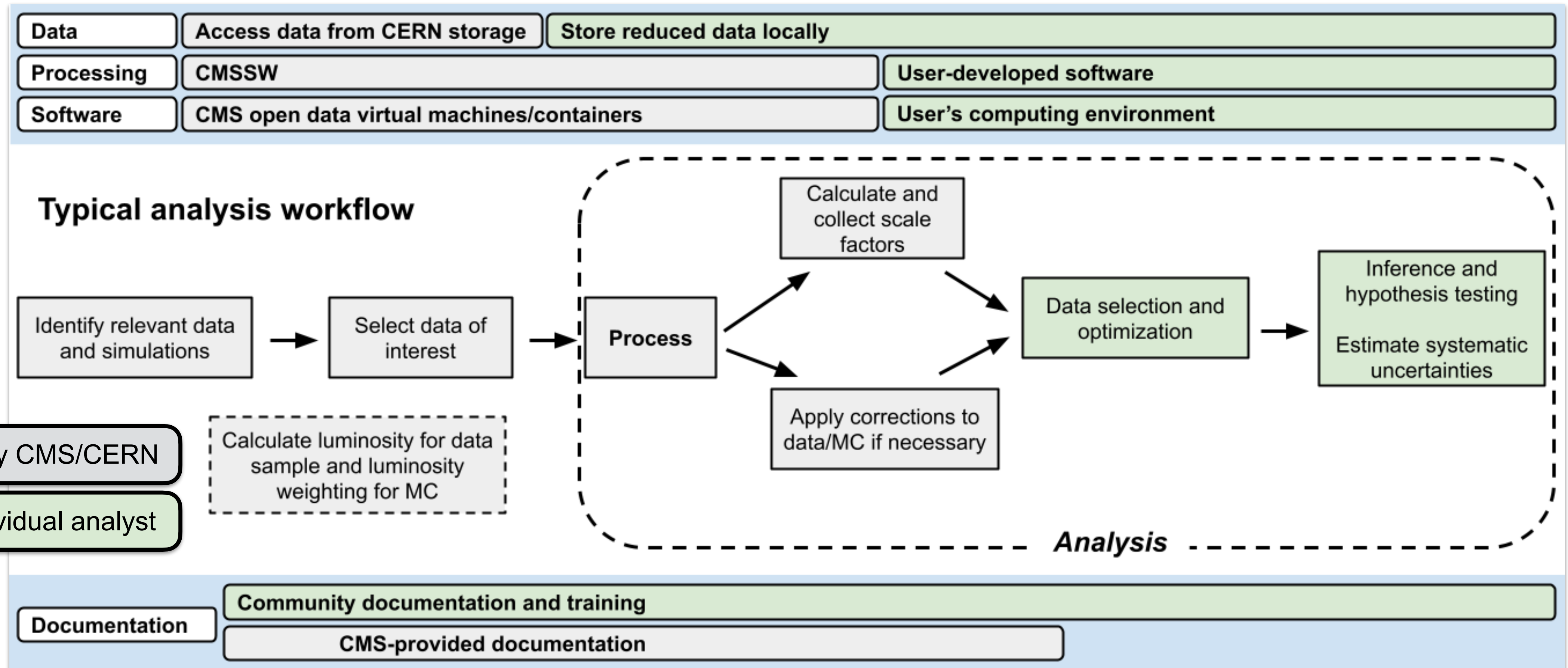**Small deviations can make a big difference**

> A few events could mean a discovery

Physics objects definitions are analysis-dependent

> An electron in one analysis might not be one in another due to different reconstruction algorithms used

Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN                    26.11.2024

The analysis part usually takes a lot of iterations

# Analysing Collider Data is very Challenging

We can **only store 0.05‰ of the collisions** (1 in 20,000 events or 2,000 events per second)

> A multi-stage trigger system selects events of interest — this bias needs to be taken into account when performing an analysis

A raw event has the size of about 2 megabytes

> We have recorded tens of billions of events, and simulated even more

> **Size can be reduced at the cost of information loss** — expertise required

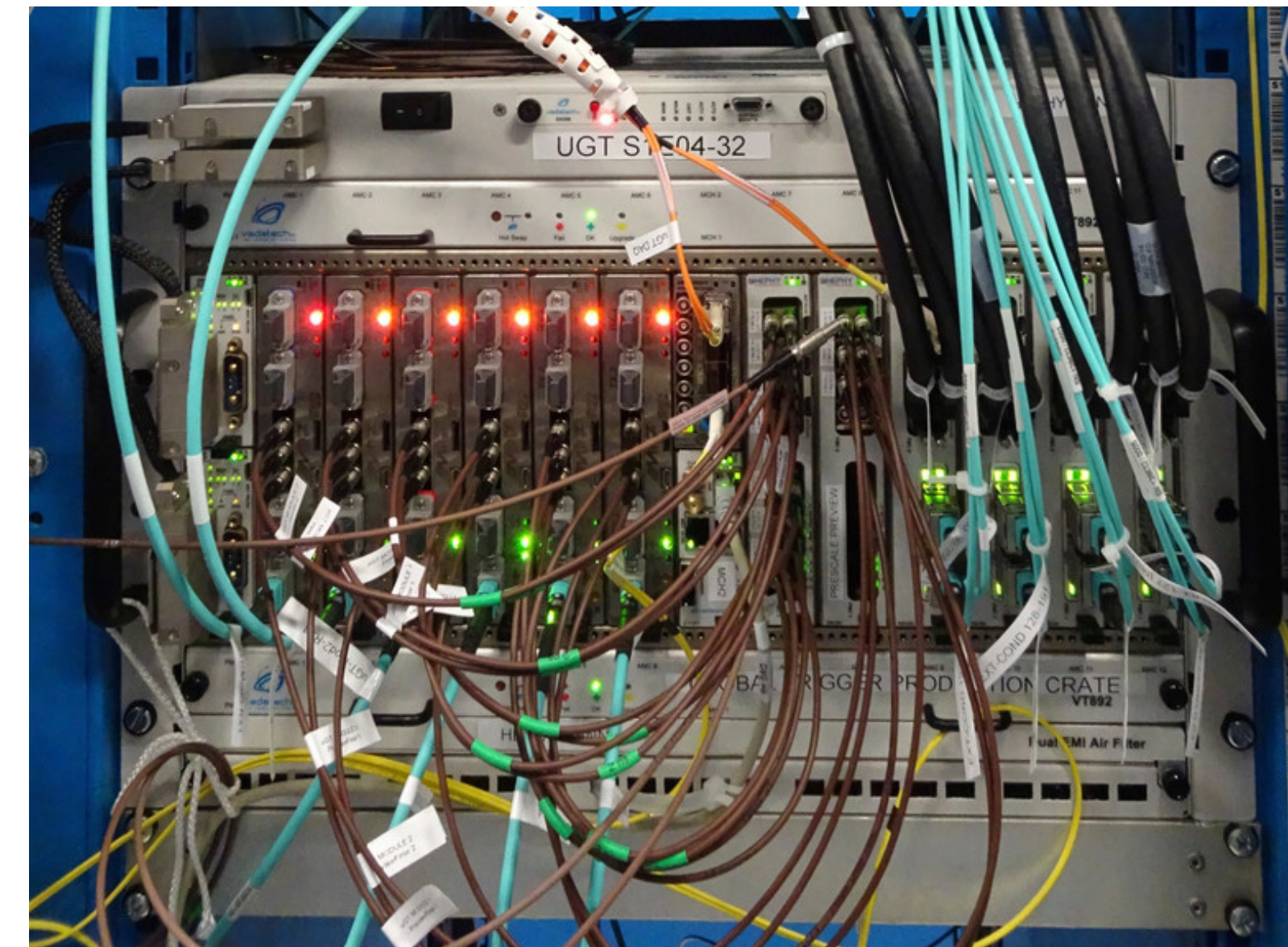> For Run 1, we have largely released "Analysis object data" (500 kB/event)

> For Run 2 (2015+), we release MiniAODs and NanoAODs (2 kB/event)

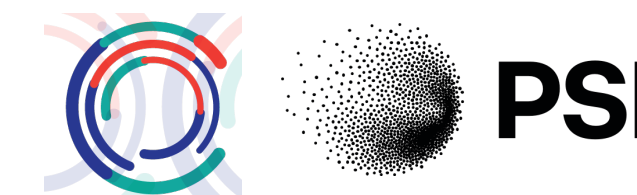Billions of events need **significant computing power** for processing

A complete physics analysis needs to take **dozens of systematic uncertainties** into account

> Understanding the relevance of individual uncertainties needs expertise

**Statistical interpretation** needs particular care



Clemens Lange (clemens.lange@cern.ch) — Open Science at CERN                26.11.2024

# Computational Challenges

We provide simplified analysis examples to lower the threshold to get started

> Pro: users can obtain a result/plot rather quickly

> Contra: these are usually far from realistic

At least the first step of the analysis chain requires substantial computing resources, ideally high-throughput batch processing systems
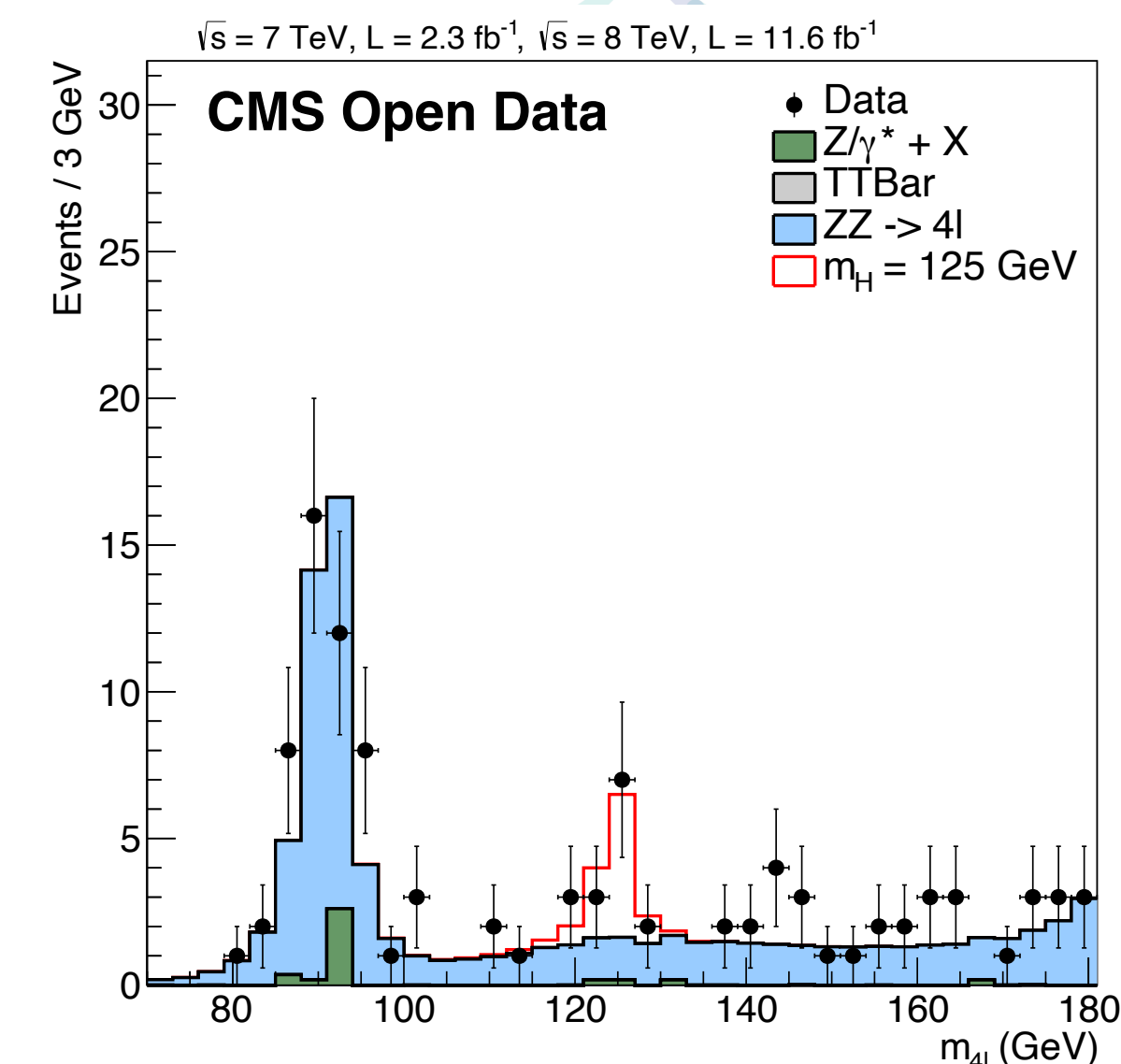
> Data sets can be processed in an "embarrassingly parallel" way

> We provide examples/tutorials on using public cloud resources

**kubernetes**

Simulation of new processes needs CMSSW

> Parts of the software are more than a decade old ➜ interfacing can be difficult



$\sqrt{s} = 7$ TeV, L = 2.3 fb$^{-1}$, $\sqrt{s} = 8$ TeV, L = 11.6 fb$^{-1}$

CMS Open Data

- Data
- $Z/\gamma^* + X$
- TTBar
- ZZ -> 4l
- $m_H = 125$ GeV

Events / 3 GeV

$m_{4l}$ (GeV)

DOI:10.7483/OPENDATA.CMS.JKB8.RR42



```
[15:00:29] cmsusr@989a8697067a ~/CMSSW_4_4_7/src $ root -b
*******************************************
*                                         *
*        W E L C O M E   to   R O O T      *
*                                         *
*   Version  5.27/06b    5 November 2010   *
*                                         *
*   You are welcome to visit our Web site  *
*          http://root.cern.ch             *
*                                         *
*******************************************

ROOT 5.27/06b (branches/v5-27-06-patches@36515, Nov 05 2010,
  15:46:56 on linuxx8664gcc)

CINT/ROOT C/C++ Interpreter version 5.18.00, July 2, 2010
Type ? for help. Commands must be C++ statements.
Enclose multiple statements between { }.
root [0]
```
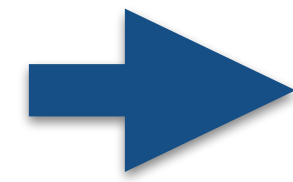
# Keeping up

When developing examples, we now aim to use open tools combined with container technologies for automatic and regular validation
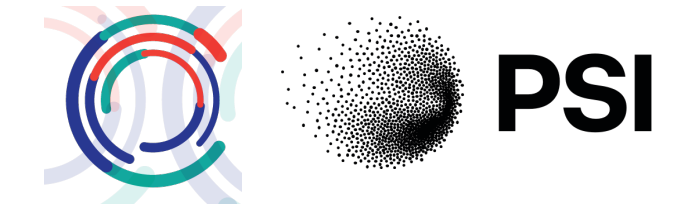
>Continuous integration using CERN's GitLab installation

>Simpler examples also run as GitHub actions **GitHub**

For easier usability, we provide examples on how get out of the HEP-specific software tool chain to industry standard tools
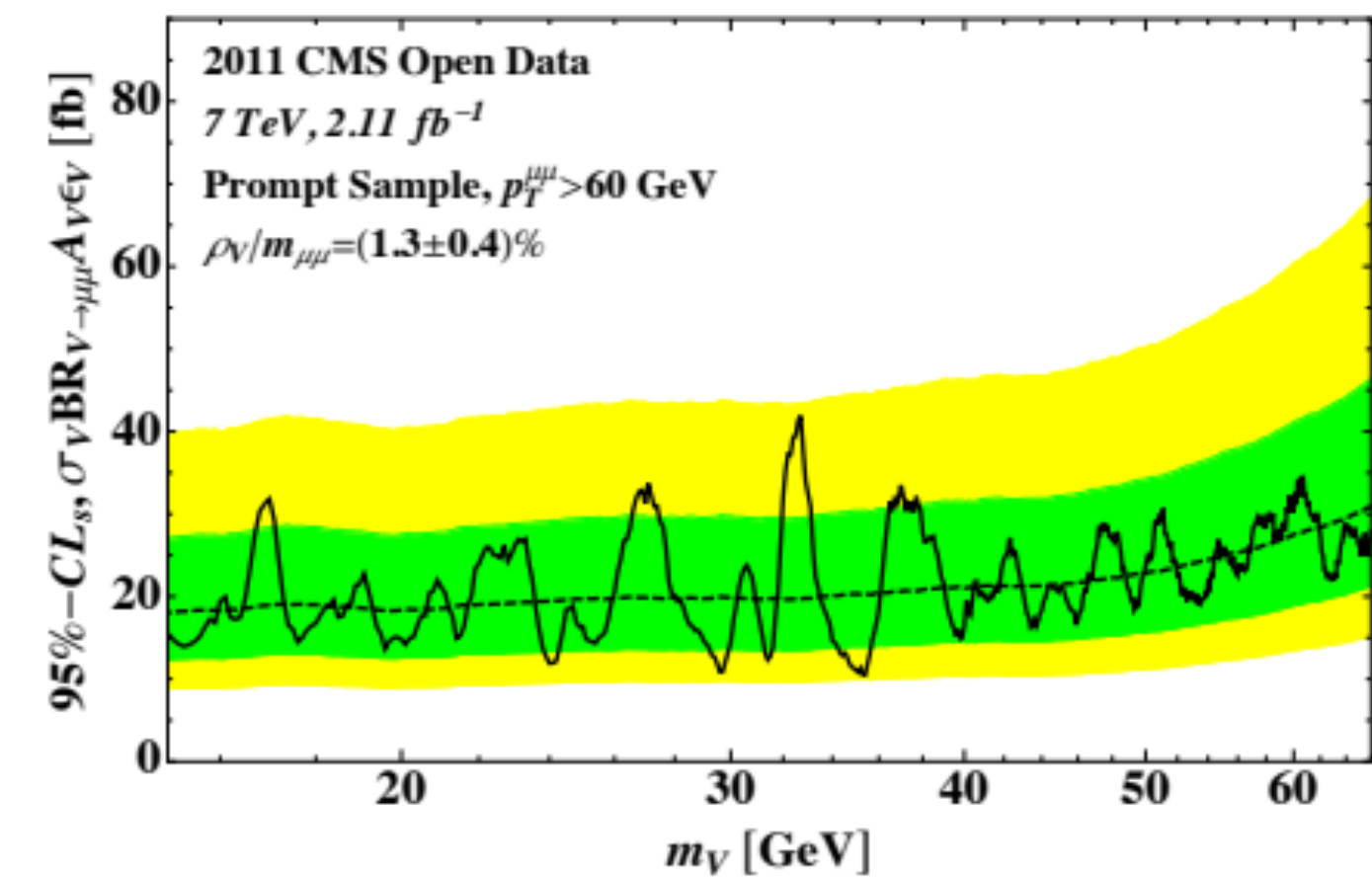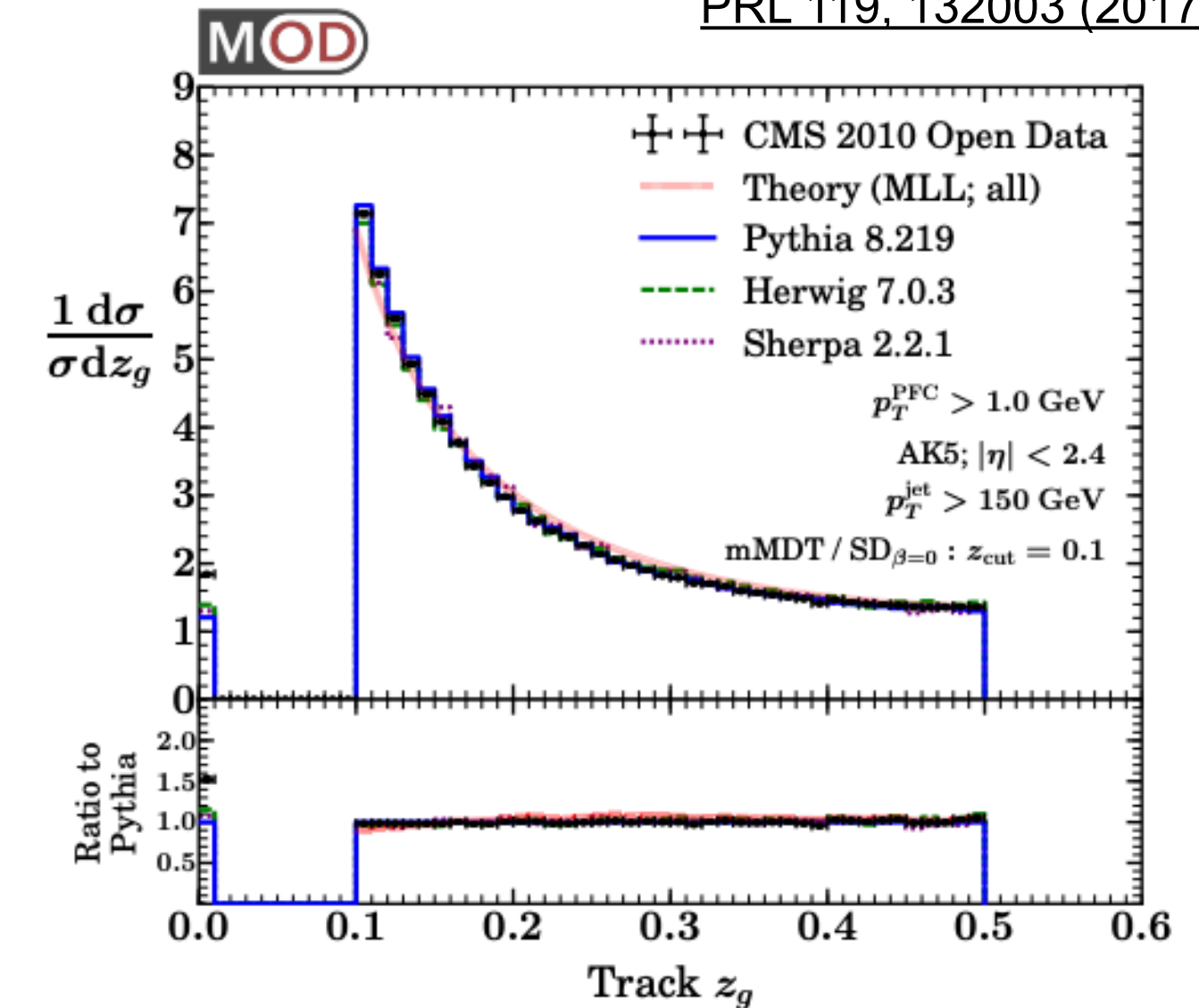
By now, CMS Open Data have been used for both actual physics results and also several computing-related projects

**Eventually, the data might be used to unveil hidden physics!**

# CERN Code of Conduct: Creativity & Professionalism

As CERN contributors, we

> Follow developments within our domain.

> Use our professional experience in a constructive manner.

> Contribute to the evolution of CERN by committing to **sharing our knowledge**.

> **Share** with internal parties **any information that could benefit them in their work**.

> Are open to new ideas and approaches.

> Adopt alternative outlooks in order to generate new thoughts and concepts.

> Conduct our work in a **structured way** to **enhance knowledge transfer and continuity**.