# ICFA Data Lifecycle Panel: Introduction

ICFA Data Lifecycle Panel meeting - April 15, 2024

Kati Lassila-Perini
Helsinki Institute of Physics - Finland
CMS Data preservation and open access coordinator

# 1 Who am I?

My background and reasons for agreeing to chair the panel.

# Kati Lassila-Perini

Experimental particle physicist, Ph.D

Research coordinator at Helsinki Institute of Physics (HIP)

CMS data preservation and open access (DPOA) coordinator



kati.lassila-perini @ cern.ch

**Decade of <u>CMS open data</u>**
– **with a small dedicated team**

CMS open data <u>in use</u>

4

Jun 22, 2014 – Apr 8, 2024

cernopendata / opendata.cern.ch

**More hands–on than "coordination"**

cernopendata / data-curation

Contributions to master, excluding merge commits

tiborsimko #1
844 commits 13,410,297 ++ 9,864,205 --

pamfilos #2
459 commits 113,278 ++ 51,965 --

katilp #3
245 commits 168,068 ++ 23,896 --

ioannistsanaktsidis #4
132 commits 27,305 ++ 19,069 --

tiborsimko #1
214 commits 16,202,028 ++ 37,564 --

heitorPB #2
38 commits 73,803 ++ 39,656 --

katilp #3
37 commits 65,152 ++ 6,000 --

OsamaMomani #4
21 commits 16,784 ++ 100 --

CERN OD team

Me

CMS students

"CV"

5

# Are CMS open data **FAIR**?

**FINDABLE**

Do you know where to look for them?

Can you find what you need?

**ACCESSIBLE**

Can you download them?

**F** **A**

**I** **R**

Are they in some common format?

Do you have the tools to open the data files?

**INTEROPERABLE**

Do you know how to use?

Can you make new research with them?

**REUSABLE**

# FAIR? Yes...

**FINDABLE**

Do you know where to look for them?
*Yes, CERN open data portal*

Can you find what you need?

*Yes, there are search functions*

**ACCESSIBLE**

Can you download them?

*Yes, or they can be streamed with XRootD*

Are they in some common format?
*Partly yes, partly no...*

Do you have the tools to open
the data files? *Yes or no, if no, they are provided*

**INTEROPERABLE**

Do you know how to use?
*There are instructions to get started...*

Can you make new research with
them?

*Yes, it takes at least as much as for the CMS people*

**REUSABLE**

F A I R

...but, some years ago...

**To make the process FAIR...**

**Private notes**
**GitHub issues**
**E-mails**
**...**

**Scripts here**
**Recipes there**
**Command history**
**...**

CMS open data release guide

# CMS open data release guide

This document is work in progress to describe the different steps needed in the preparations of CMS public data releases.

CMS open data are accompanied with rich metadata such as

- the type and the size (in terms of events and volume) of each dataset
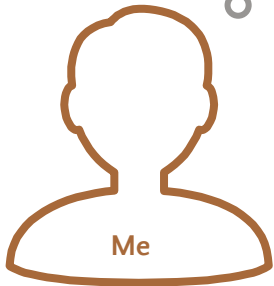- the provenance information i.e. how these data were collected or generated and then reprocessed
- the context metadata, i.e. in which environment and with which additional assets these data should be used.

The release preparation consists of two parts: first, identifying the information sources internal to CMS (databases and servers, or documentation pages) and, then, building the open data records based on this information.

The scripts for the latter part, i.e. preparing the open data records, are stored and available in the data-curation repository of the CERN Open data portal GitHub repository.

The goal of this document is to describe where and how the inputs to these data-curation scripts are taken from the CMS internal resources. In addition, it aims to cover the full knowledge needed for the data release including the administrative procedures for the release approvals and storage requests.

These pages are deployed as mkdocs from https://gitlab.cern.ch/cms-opendata/cms-opendata-releaseguide. The access is restricted to the `cms-web-access` list. Additional access rights can be added in the `dpoa-cms-data-release-docs-access` e-group (by cms-dpoa-coordinators, admins of that group).

**... I needed to change how I work.**

## 2  Data Lifecycle panel

The mandate is broad, even overwhelming…

My expertise is CMS/CERN/Open data/FAIR–centric

Diverse range of expertise within the panel – looking forward!

Practicalities

# ICFA statement on the Data Lifecycle Panel
# Mandate of the Data Lifecycle Panel

"

# Rely on broad expertise

DPHEP

Panel members

SCIC

# **Practicalities**

Frequency of the meetings?

- ◉ Monthly? Set a regular day and time?

Privacy of the agendas?

- ◉ <u>Agendas</u> public, but attendance and recordings members only?

Communications:

- ◉ Mailing list: ICFA-Data-Lifecycle-panel at cern.ch
- ◉ Collect input from the members through "surveys" for each meeting.
- ◉ Other channels?

Web site

- ◉ <u>https://icfa.hep.net/icfa-panel-on-the-data-lifecycle/</u>

# 3 Stakeholders

Individuals researchers, within the collaborations, carry out the work.

The surrounding stakeholders may either empower or restrict the researchers' ability to adopt best practices for the full "data lifecycle".

# **Stakeholders**

Infrastructure
Resources
*Open science guidelines*

Host
laboratory

Researcher

Home institute

Task definition, share of time on them
Resources (salary, travel, personal tools)
Education, training
*Research/data management guidelines*

Experiments
Working group

Responsibilities
Resources (computing, software)
Practices ("how do we work")
*Policies*

15

*Anecdotic feedback from an Open data –workshop for PhD students in physics (outside HEP)*

**4** Guiding principle: Bridge the gap

between words and actions
between prototyping and implementation

# Bridge the gap – between words and action



Task definition, share of time on them
Resources (salary, travel, personal tools)
Education, training
*Research/data management guidelines*

Home institute

Infrastructure
Resources
*Open science guidelines*

Host
laboratory

Researcher

Experiments
Working group

Responsibilities
Resources (computing, software)
Practices ("how do we work")
*Policies*

# **Bridge the gap** – between words and action



Task definition, share of time on them
Resources (salary, travel, personal tools)
Education, training
*Research/data management guidelines*

Infrastructure
Resources
*Open science guidelines*

Home institute

Host
laboratory

Researcher

Experiments
Working group

Responsibilities
Resources (computing, software)
Practices ("how do we work")
*Policies*

*Words are here*

# **Bridge the gap** – between words and action



Task definition, share of time on them
Resources (salary, travel, personal tools)
Education, training
*Research/data management guidelines*

Infrastructure
Resources
*Open science guidelines*

Home institute

Host laboratory

Researcher

**Action is here…**
**…and to make it FAIR needs:**

Knowledge preservation
Data and software skills
Tools, Time

Experiments
Working group

Responsibilities
Resources (computing, software)
Practices ("how do we work")
*Policies*

*Words are here*

20

# How to assess panel's actions on open science and FAIR practices?

Panel has "a focus on open science and FAIR practices"

Panel's actions are effective if they help:

Individuals carry out the work, not panels or committees.

integrate open science practices into researchers' daily work

identify factors that make FAIR practices difficult in daily work

# **Bridge the gap** – between prototyping and implementation

Prototyping

can be

Cross experiments

Short-term project

can have

Researcher

Cross disciplines

Separate funding

Industry involvement

Project specific goals

# Bridge the gap – between prototyping and implementation

# **Bridge the gap** – between prototyping and implementation



Home institute

Host laboratory

Researcher

Experiments Working group

Prototyping

**To make this reusable:**

Knowledge preservation
Data and software skills
Tools, Time
Other success factors?

Implementation

can be

Cross experiments

Short-term project

can have

Researcher

Cross disciplines

Separate funding

Industry involvement

Project specific goals

# 5 Your input

Thanks to everyone having responded the survey!

**Panel's actions are effective if they help:**

**integrate open science practices into researchers' daily work**

**identify factors that make FAIR practices difficult in daily work**
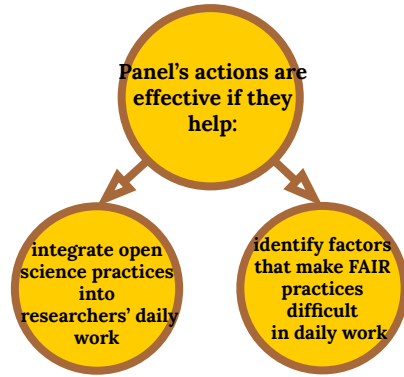
**Panel should identify factors for successful implementation of the outcomes of prototyping / exploring**

Left chart:
- Strongly agree: 3
- Agree: 6
- Neither agree nor disagree: 2
- Disagree: 0
- Strongly disagree: 0

Responses

Right chart:
- Strongly agree: 6
- Agree: 3
- Neither agree nor disagree: 2
- Disagree: 0
- Strongly disagree: 0

Responses

I think developing a way to assess the panel's positive impacts is very important, though it is challenging. I think we can do a lot of good work in helping to spread best practices from where they are developed to the wider community.

This panel covers so much ground that I don't know exactly where to start. I think it would be better to start with the definitions of the words so that all panel members have the same understanding in the discussion (e.g. What is the data lifecycle? the data lifecycle may also vary depending on the scale of the experiment...)

I suggest we do analyze the mandate and establish a 2-3 year plan of action in some of the directions (we will not be able to follow everything at the same level of involvement). I believe that the guiding principle for the panel's action is to have a strategic approach for the longer term (as opposed to already well covered aspects in the community).

I hope that this body can be an advocate for common tools and services, because the experiments themselves can not do that. They have a hard enough time advocating for funding of their own operations. In the end both sides are needed.

We need to define a multi-branched structure that combines advanced technology initiatives, pathways from prototyping into production and making clear to the community technology trends and how they are related to the coming "data intensive" and analysis challenges, and possible solutions following in-depth prototyping and integration with some of the major data management and analysis toolsets of the collaborations. We also need to propagate knowledge and trends in advanced networking and how these can be exploited to solve some of the challenges of the HL LHC era.
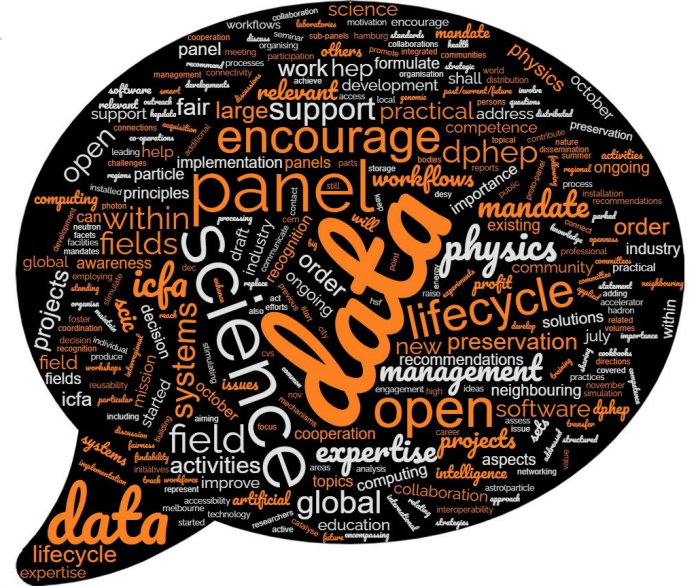
"

**Outlook**

**Looking forward to the next panel meeting in May (tbc).**

Discuss the plan of action
Panel substructure?

Report at the ICFA meeting in July

# Thank you!

## Questions?

And thanks to SlidesCarnival for this free presentation template