

FAIR Data for Complex HEP Datasets

FAIROS-HEP Kick-Off Workshop - 8-10 Feb 2023



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

1

CMS Open data - FAIR?

Findable - Accessible - Interoperable - Reusable



FAIR? My interpretation... (see also Tibor's talk)

FINDABLE

Do you know where to look for them?

Can you find what you need?

F

A

ACCESSIBLE

Can you download them?

Are they in some common format?

Do you have the tools to
open the data files?

INTEROPERABLE

I

R

Do you know how to use?

Can you make new research with
them?

REUSABLE



FAIR? My interpretation...

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

F

ACCESSIBLE

Can you download them?

Yes, or they can be streamed with XRootD

A

Are they in some common format?

No, not really...

Do you have the tools to open the data files?
No, but they are provided

INTEROPERABLE

I

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

R

REUSABLE



Metadata for complex research data

**Content
metadata
WHAT?**

**Provenance
metadata
FROM
WHERE?**

**Contextual
metadata
HOW TO
USE?**

**Contextual:
HOW TO
USE?**

Dataset characteristics

76523854 events, 1607 files, 1.5 **TB** in total.

System details

Recommended global tag for analysis: 76X_dataRun2_16Dec2015_v0

Recommended release for analysis: CMSSW_7_6_7

How were these data selected?

Events stored in this primary dataset were selected because of the presence of a high scalar sum of the jet transverse momenta (HT); or at least one or two energetic jets.

Data taking / HLT

The collision data were assigned to different RAW datasets using the following HLT configuration.

Data processing / RECO

This primary MINIAOD dataset was processed from the RAW dataset by the following step:

Step: RECO

Release: CMSSW_7_6_3

Global tag: 76X_dataRun2_v15

Configuration file for RECO step reco_2015D_JetHT

HLT trigger paths

The possible HLT trigger paths in this dataset are:

HLT_AK8DIPFJet250_200_TrimMass30_BTagCSV0p45

HLT_AK8DIPFJet280_200_TrimMass30_BTagCSV0p45

HLT_AK8PFHT600_TrimR0p1PT0p03Mass50_BTagCSV0p45

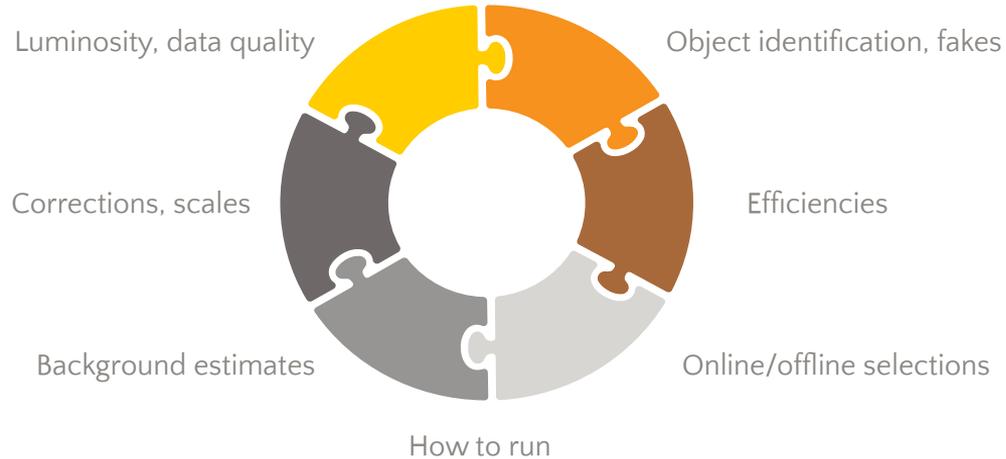
**Content:
WHAT?**

**Provenance:
FROM
WHERE?**



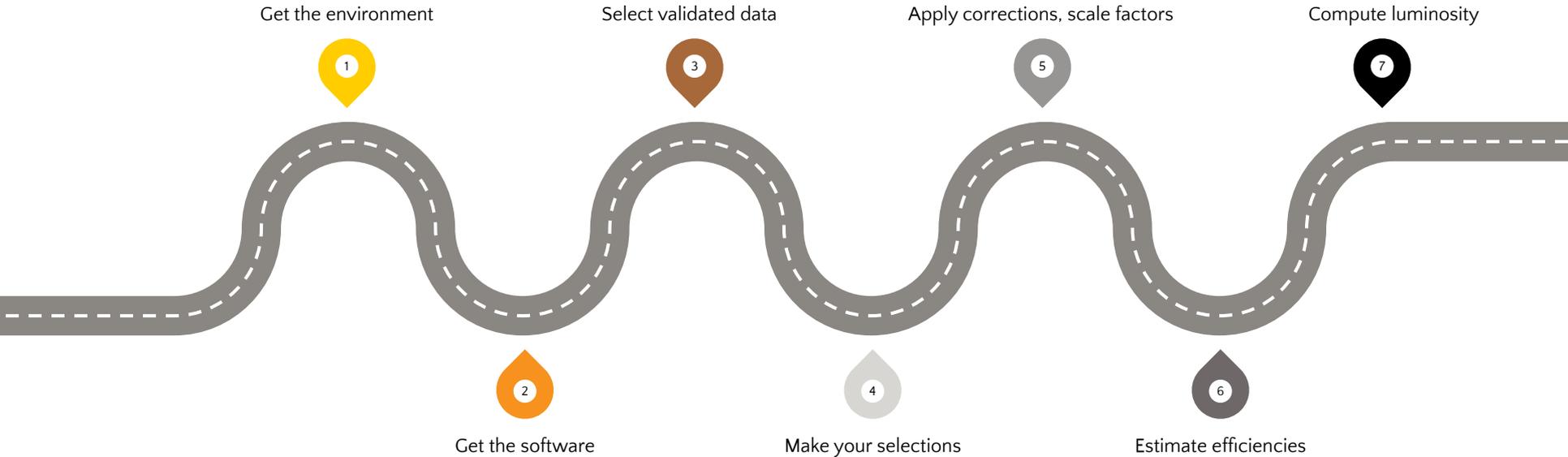


Contextual metadata - should cover a lot!





How to put this together?





Providing contextual metadata in a useful way

- Best through example workflows
 - automated, machine and human-readable
- Not easy to define (even with >1100 papers...)
 - partly because
 - analysis processes are complex
 - CMS data format supports a wide range of usecases
 - but also because we, as a community, have undervalued:
 - documentation
 - common tools
 - analysis code reuse.



tidy up a bit to align with 2012 version struc...

add trigger analyzer

Update README.md (#102)

README.md

Physics Objects Extractor (PhysObjectExtractor) for 2015MiniAOD data

Description

The `PhysObjectExtractor` package is the heart of the POET repository. It contains a collection of `Analizers` that extract information from different physics objects into a `ROOT` file called `output.root`. Each `Analyzer` has been written separately for clarity and can be executed modularly using the configuration file called `poet_cfg.py`.

We need a logo!!!

THE example code: **POET**

Physics Objects Extractor Tool

Put together by many people in the CMS open data group from various sources within CMS.

See Julie's talk on Wednesday.

Not a negligible effort!

2

Open science - how to?

Example workflows do not come out of the blue - it has been a major effort to put them together for CMS open data.

A little interlude to what it takes to open science to happen

2.1

Roles & responsibilities



Open science - what it takes to make it happen



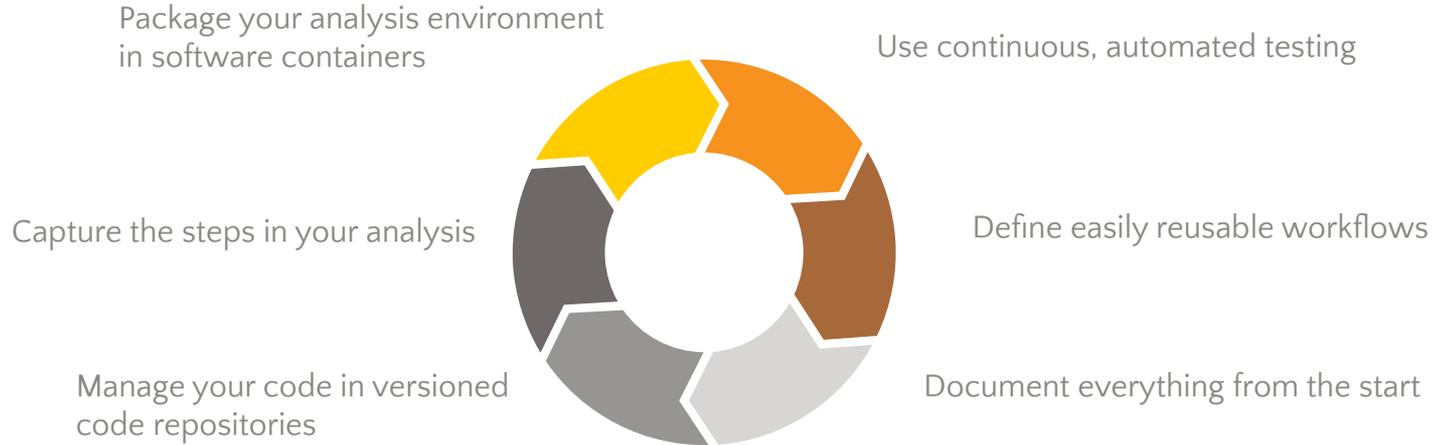
2.2

Best practices

To preserve the knowledge at the time of active analysis



Efforts are needed

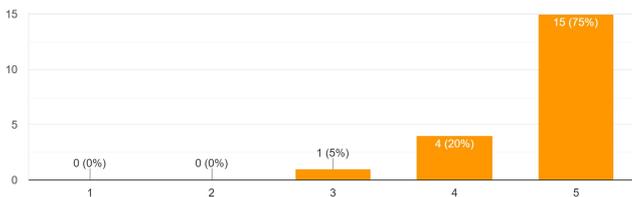


Best practices require time but they will pay off:

for the individual, for the group, and eventually, for open science!

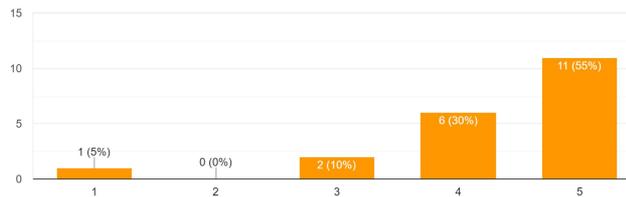
Do you think that following best practices in software development in the data analysis work is important?

20 responses



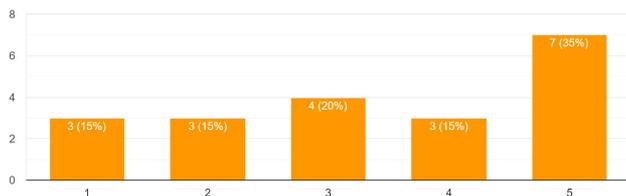
Did this workshop motivate you to improve your software and code management practices to better preserve your analysis work?

20 responses



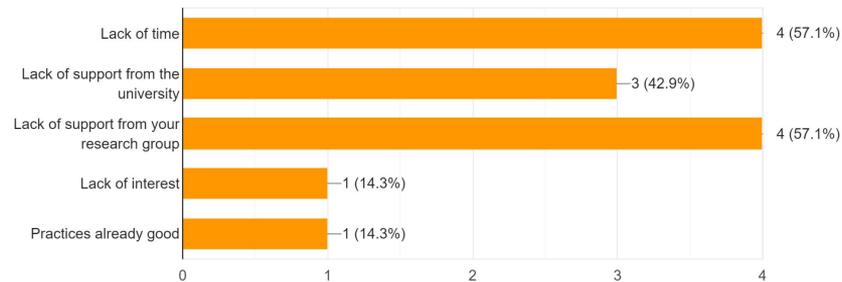
Do you think that you actually have the possibility to improve your software and code management practices and apply them to your work?

20 responses



If you answered 1 or 2 in the above what are the factors that make it difficult?

7 responses



Feedback from a recent Open data -workshop for PhD students in physics (outside HEP)





All parts count

It is the researcher that has the key role in preserving the knowledge and needs support in doing so!

3

Assessing FAIRness with workflows

Automated workflows that concretely access, use and test the CMS open data metadata in the research context.



CMS Open Data workshop 2022!

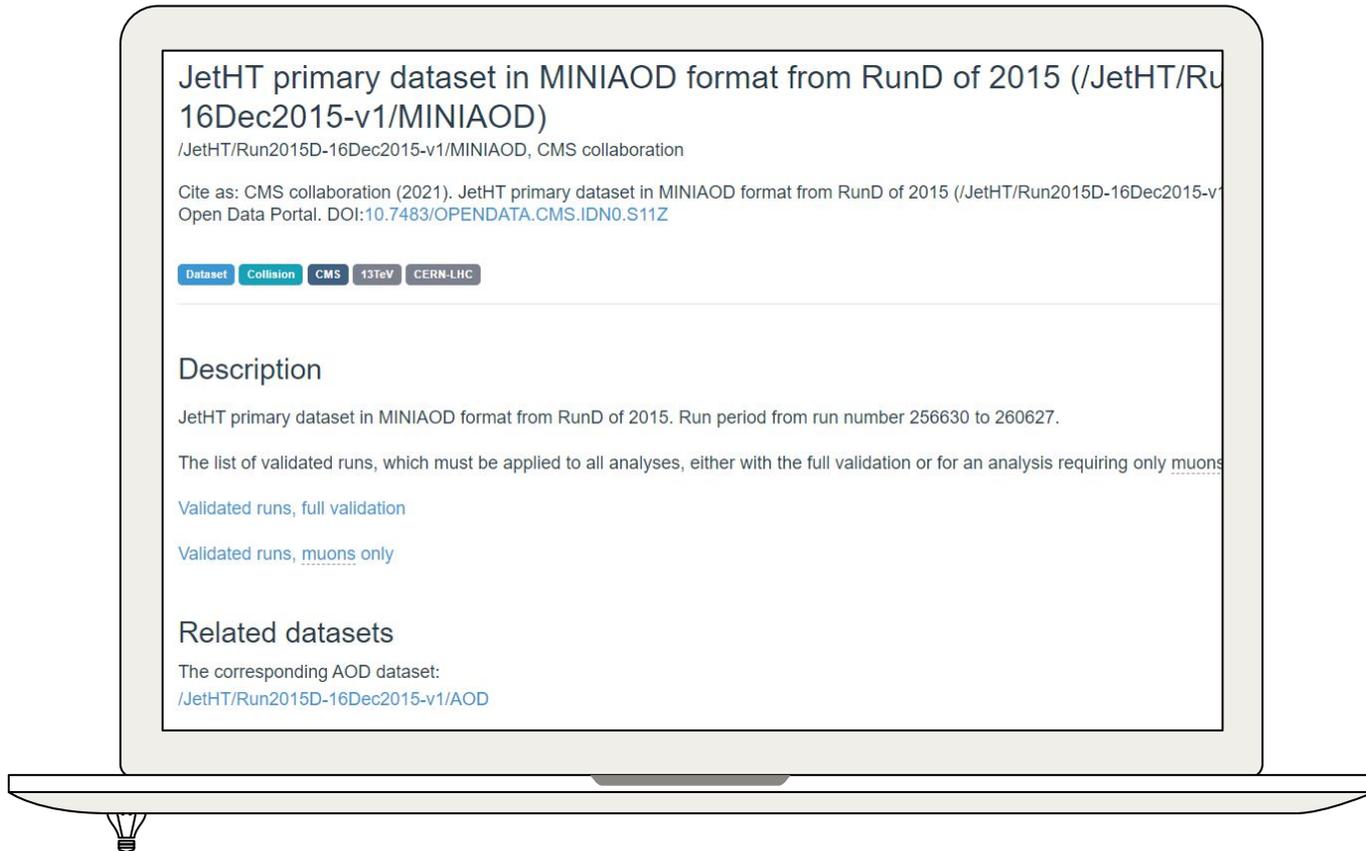


**The code for the workflow?
Assembled for the CMS OD workshops!**



Watch on  YouTube





Open dataset record **display**



```
{ } metadata.json U ×
{} metadata.json > ...
1  {
2  |   "$schema": "http://opendata.cern.ch/schema/records/record-v1.0.0.json",
3  |   "abstract": {
4  |     |   "description": "<p>JetHT primary dataset in MINIAOD format from RunD of 2015. Run period from run number 25663
5  |     |   |   p><p>The list of validated runs, which must be applied to all analyses, either with the full validation or for
6  |     |   |   |   requiring only muons, can be found in:</p>",
7  |     |   |   |   "links": [
8  |     |   |   |   |   {
9  |     |   |   |   |   |   "description": "Validated runs, full validation",
10 |     |   |   |   |   |   "recid": "14210"
11 |     |   |   |   |   |   },
12 |     |   |   |   |   |   {
13 |     |   |   |   |   |   |   "description": "Validated runs, muons only",
14 |     |   |   |   |   |   |   "recid": "14211"
15 |     |   |   |   |   |   }
16 |     |   |   |   |   }
17 |     |   |   |   }
18 |     |   |   }
19 |     |   }
20 |     "accelerator": "CERN-LHC",
21 |     "collaboration": {
22 |       |   "name": "CMS collaboration"
23 |     },
24 |     "collections": [
25 |       |   "CMS-Primary-Datasets"
26 |     ],
27 |     "collision_information": {
28 |       |   "energy": "13TeV",
29 |       |   "type": "pp"
30 |     },
31 |     "control_number": "24124",
32 |     "date_created": [
```

Underlying metadata

**Context:
environment
software**

Dataset characteristics

76523854 events. 1607 files. 1.5 TB in total.

System details

Recommended [global tag](#) for analysis: 76X_dataRun2_16Dec2015_v0

Recommended release for analysis: CMSSW_7_6_7

How were these data selected?

Events stored in this primary dataset were selected because of the presence of two energetic jets.

Data taking / HLT

The collision data were assigned to different RAW datasets using the following

Data processing / RECO

This primary MINIAOD dataset was processed from the RAW dataset by the following

Step: RECO

Release: CMSSW_7_6_3

Global tag: 76X_dataRun2_v15

[Configuration file for RECO step reco_2015D_JetHT](#)

HLT trigger paths

The possible HLT trigger paths in this dataset are:

[HLT_AK8DIPFJet250_200_TrimMass](#)

[HLT_AK8DIPFJet280_200_TrimMass](#)

[HLT_AK8PFHT600_TrimRop1PT0p03](#)

**Context:
trigger path**

**Context:
validated
data selection**

JetHT primary dataset in MINIAOD format from RunD of 2015 (/J.../16Dec2015-v1/MINIAOD)

/JetHT/Run2015D-16Dec2015-v1/MINIAOD, CMS collaboration

Cite as: CMS collaboration (2021). JetHT primary dataset in MINIAOD format from RunD of 2015. [Open Data Portal](#). DOI:10.7483/OPENDATA.CMS.IDN0.S11Z

[Dataset](#) [Collision](#) [CMS](#) [13TeV](#) [CERN-LHC](#)

Description

JetHT primary dataset in MINIAOD format from RunD of 2015. Run period from run number 26030 to 260627.

The list of validated runs, which must be applied to all analyses, either with the full validation or for an analysis requiring

[Validated runs, full validation](#)

**Context:
usage
instructions**

How can you use these data?

You can access these data through the CMS Open Data container or the CMS Virtual Machine. See the instructions for setting up one of the two alternative environments and getting started in

[Running CMS analysis code using Docker](#)

[How to install the CMS Virtual Machine](#)

[Getting started with CMS open data](#)

File Indexes

Filename	Size	
CMS_Run2015D_JetHT_MINIAOD_16Dec2015-v1_00000_file_index.txt	2.4 kB	List Files Download
CMS_Run2015D_JetHT_MINIAOD_16Dec2015-v1_50000_file_index.txt	123.9 kB	List Files Download



First thoughts: to-do list

- The CODP metadata is display-oriented and some fields of interests contain HTML
 - Not easily machine-readable
- ⇒ Modify the existing CODP metadata structure

- Some contextual metadata is not exposed but embedded in the software or in condition data
 - Trigger prescales, corrections, scales etc
- ⇒ Will be addressed with nanoAOD...



First concrete attempts...

Container image:
cernopendata-client
Task:
Get metadata

1

Container image:
cernopendata-client
Task:
Get validated runs list

3

Container image:
cmsopendata/cmssw_7_6_7-...
Task:
Run the job(s)

5

Container image:
root
Task:
Make reference plots/output

7

2

Container image:
python/bash
Task:
Find needed values from metadata

4

Container image:
python/bash
Task:
Build the job configuration

6

Container image:
gitlab-registry.cern.ch/cms-cloud/brilws-docker
Task:
Get the luminosity

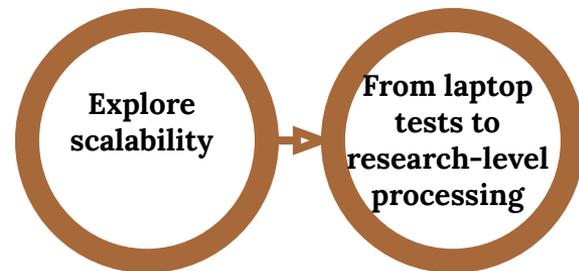
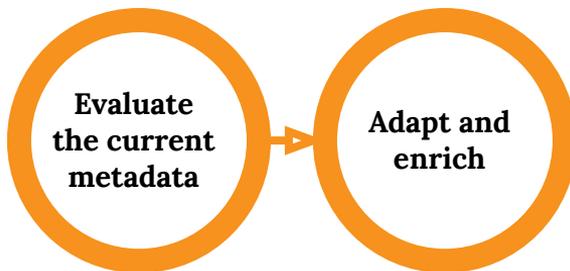
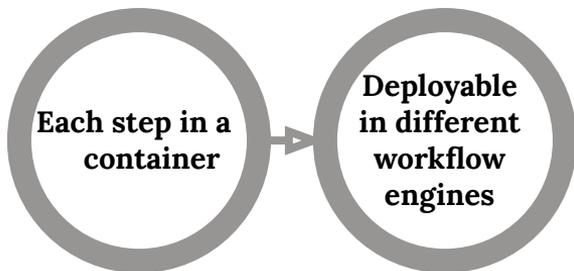
4

Conclusions

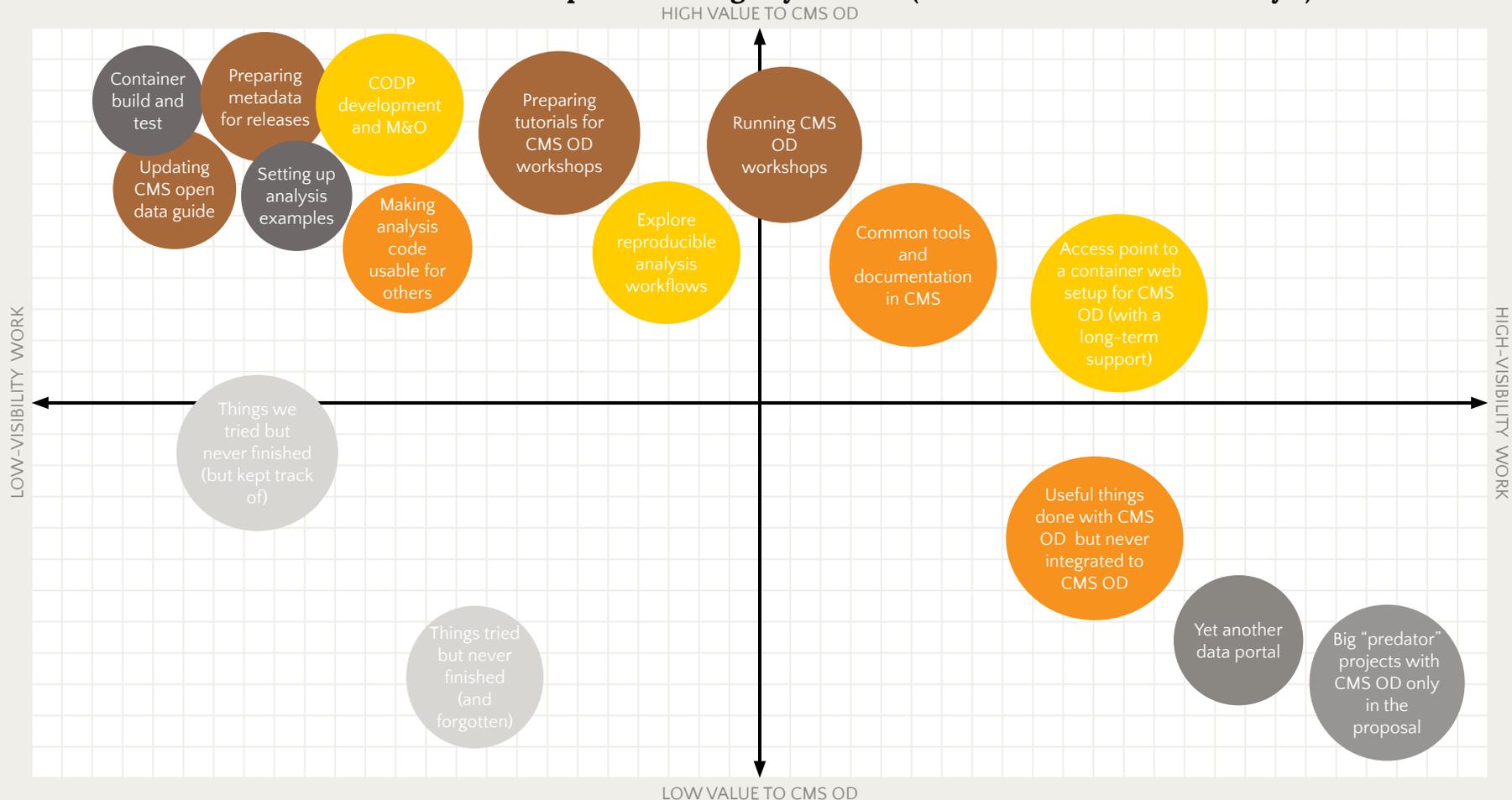
and some personal considerations...



Proposal: CMS open data as a testbed



Personal considerations: Value to CMS open data vs “glory” Matrix (not to be taken too seriously...)





Thank you!

Any **questions** ?

And thanks to [SlidesCarnival](#) for this free presentation template