

The LHCb NTuple Wizard

Sebastian Neubert¹

¹HISKP Bonn

Mini Workshop Best Practices in
Model preservation
14.12.23



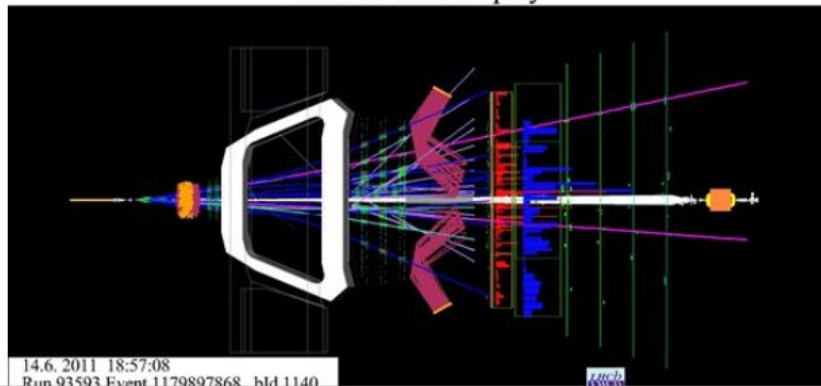
News › News › Topic: Knowledge sharing

LHCb releases first set of data to the public

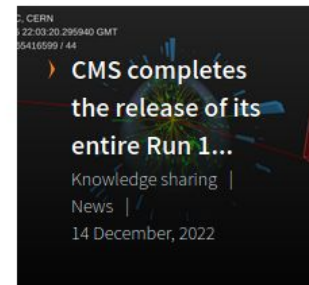
The LHCb collaboration has released data from Run 1 of the LHC to the public for the first time, allowing research to be conducted by anyone in the world

8 DECEMBER, 2022 | By LHCb collaboration

LHCb Event Display



Related Articles



Dataset × Collision × LHCb ×

Include on-demand datasets

Filter by type

- ▼ Dataset 31
 - Collision 28
 - Derived 3
 - ▶ Documentation 2149
 - ▶ Environment 1
 - News 1
 - ▶ Software 1

Filter by experiment

- ALICE 14
- CMS 224
- LHCb 28

Filter by year

- 2011 13
- 2012 14

Filter by file type

- DST 15
- MDST 12
- root 1

Filter by collision type

- pp 27

Filter by collision energy

Sort by: Title A-Z ▾ asc. ▾

Display: detailed ▾ 20 results ▾

Found 28 results.

LHCb 2011 Beam3500GeV MagDown EW Stream Stripping21r1

proton-proton (pp) collision data collected by the LHCb experiment in the year 2011 of Run1 of the LHC....

[Dataset](#) [Collision](#) [LHCb](#)

LHCb 2011 Beam3500GeV MagDown EW Stream Stripping21r1p1

proton-proton (pp) collision data collected by the LHCb experiment in the year 2011 of Run1 of the LHC....

[Dataset](#) [Collision](#) [LHCb](#)

LHCb 2011 Beam3500GeV MagDown EW Stream Stripping21r1p2

proton-proton (pp) collision data collected by the LHCb experiment in the year 2011 of Run1 of the LHC....

[Dataset](#) [Collision](#) [LHCb](#)

LHCb 2011 Beam3500GeV MagDown LEPTONIC Stream Stripping21r1

proton-proton (pp) collision data collected by the LHCb experiment in the year 2011 of Run1 of the LHC....

[Dataset](#) [Collision](#) [LHCb](#)

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

DPHEP Levels of Data Complexity

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

1. Published results + additional information

- supplemental data tables, ntuples
- HEPData entries, rivet plugins
- notes, technical information
- documentation, slides
- analysis code, jupyter notebooks

2. Education and Outreach

- simplified data formats, e.g. highly preprocessed ntuples

3. Reconstructed data + analysis level software

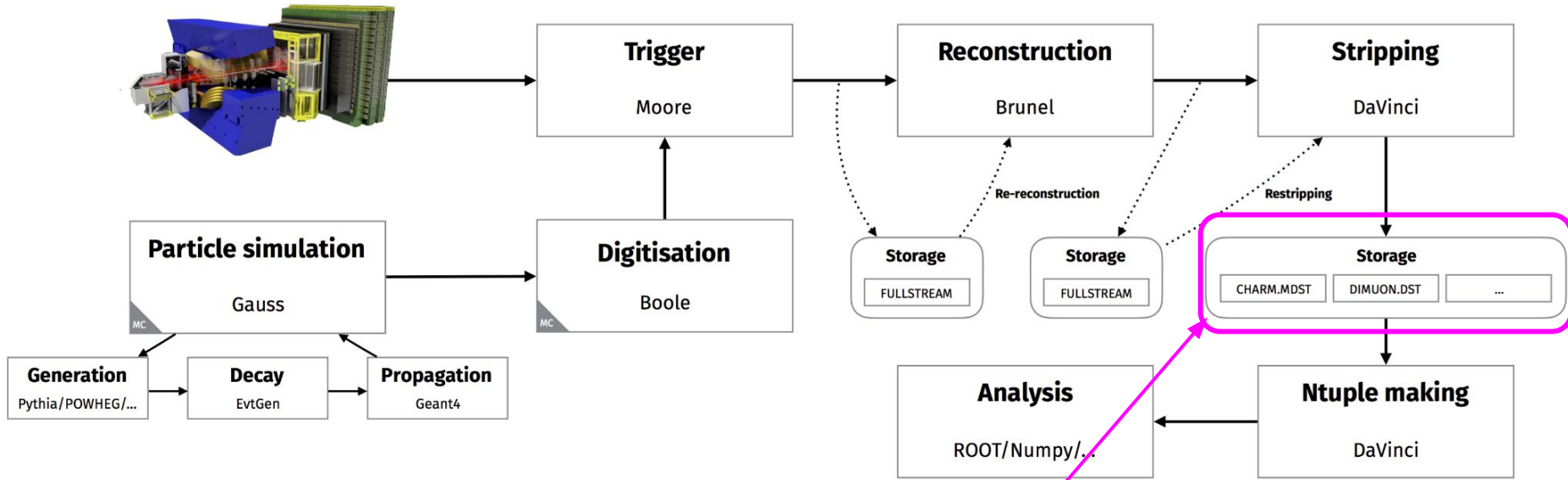
- Calibrated reconstructed data with the level of detail useful for algorithmic, performance and physics studies
- preservation of analysis level experiment-specific software

4. Raw data + reconstruction software

- Not released for LHC data

LHCb Level 3 Data

Release policy: 50% @ 5yrs, 100% @ 10yrs
after end of running period



- Level 3 data in LHCb **defined as the output of the stripping**
- Same level of abstraction accessed by LHCb members
- Organized in ~10 streams, according to physics signature
- Software needed to access data (DaVinci) [is open source](#), available via CVMFS (or container)
- Documentation: [LHCb Starterkit](#) openly available

Data to be released very soon: full 2011/12 Stripping Output

Mindaugas Sarpis, Bonn



BHADRON.MDST
BHADRONCOMPLETEEVENT.DST
CHARM.MDST
CHARMCOMPLETEEVENT.DST
DIMUON.DST
EW.DST
LEPTONIC.MDST
RADIATIVE.DST
SEMILEPTONIC.DST

} Already
released
~ 200TB

} Ready to be released
(preparing press announcement)
~ 800TB

All this data comes with meta-data and documentation.

The data we are releasing now contains almost 7500 Pre-Selections (stripping lines)

- This includes common particles used to build more complex decays
- DST files contain the full events - data mining possible
- MDST files only contain the particles selected by the stripping lines (rest of event is discarded)

But it still is in DST/MDST Format and requires the LHCb Software to read.

LHCb Preselections

All the info is there but you need to know how to look for it. E.G.:

The decay $B^0 \rightarrow D \pi$
with $D \rightarrow K_s h$
where the K_s is reconstructed from two downstream tracks

is selected by the

B02DPiD2KSHDDBeauty2CharmLine

The screenshot displays the LHCb Preselections search interface. At the top, there are tabs for 'Documentation', 'Stripping', and 'LHCb'. Below the tabs, there are several filter sections:

- Include on-demand datasets:**
- Filter by type:** Documentation (2533), Stripping (2533)
- Filter by experiment:** LHCb (2533)
- Filter by year:** 2011 (956), 2012 (1577)
- Filter by stripping stream:** COMMONPARTICLES (838), EW (382), LEPTONIC (898), RADIATIVE (16), bhadron (2533), bhadroncompleteevent (239), charm (640), charmcompleteevent (87), commonparticles (554), dimuon (106)
- Filter by stripping version:** stripping21 (854), stripping21r0p1 (274), stripping21r0p2 (449), stripping21r1 (848), stripping21r1p1 (108)
- Filter by event number:** 0-999 (0), 1000-9999 (0), 10000-99999 (0), 100000-999999 (0), 1000000-9999999 (0), 10000000- (0)

At the top right, there are options for 'Sort by: Best match', 'asc.', 'Display: detailed', and '20 results'. Below this, it says 'Found 2533 results.' The search results are listed in a table with the following entries:

Search Result	Count
LHCb Stripping V21 BHADRON Stream B02DKWSD2PI0HHHRESOLVEDBEAUTY2CHARM Line	2533
LHCb Stripping V21 BHADRON Stream B02DPID2KSHDDBEAUTY2CHARM Line	1577
LHCb Stripping V21 BHADRON Stream B02DPIPIWSD2HHHPIDBEAUTY2CHARM Line	838
LHCb Stripping V21 BHADRON Stream B02DPIWSD2PIOHHHMERGEDBEAUTY2CHARM Line	640
LHCb Stripping V21 BHADRON Stream B02DPIWSNOIPDS2HHHPIDBEAUTY2CHARM Line	274
LHCb Stripping V21 BHADRON Stream B02DSTARKDST2D0PI-D2KSHHDDBEAUTY2CHARM Line	848
LHCb Stripping V21 BHADRON Stream B02DSTARSKDDST2D0PIBEAUTY2CHARM Line	848
LHCb Stripping V21 BHADRON Stream B02DSTARSPIDDDST2D0PIBEAUTY2CHARM Line	108

Level 3 Data - Resource Projections

2020 projections

	ALICE	ATLAS	CMS	LHCb
Run 2	2 PB	0.5 PB	2 PB	10 PB (including Run 1)
Run 3	4 PB	1 PB	4 PB	45 PB
Total	6 PB	1.5 PB	6 PB	55 PB

Mitigation Strategies:

- Provide protected access to existing copies of stripping/turbo output via WG-production slots. Needs “ntupling wizard”
- Provide direct access to data on grid storage

FAIR Data Principles

[The FAIR Guiding Principles for scientific data management and stewardship. *Nature Sci Data* 3, 160018 (2016). <https://doi.org/10.1038/sdata.2016.18>]



Findable: Metadata and data should be easy to find for both humans and computers.



Accessible: The exact conditions under which the data is accessible should be provided in such a way that humans and machines can understand them.



Interoperable: The (meta)data should be based on standardized vocabularies, ontologies, thesauri etc. so that it integrates with existing applications or workflows.



Reusable: Metadata and data should be well-described so that they can be replicated and/or combined in different research settings.

<https://go-fair.org>

Solved by

<https://opendata.cern.ch>

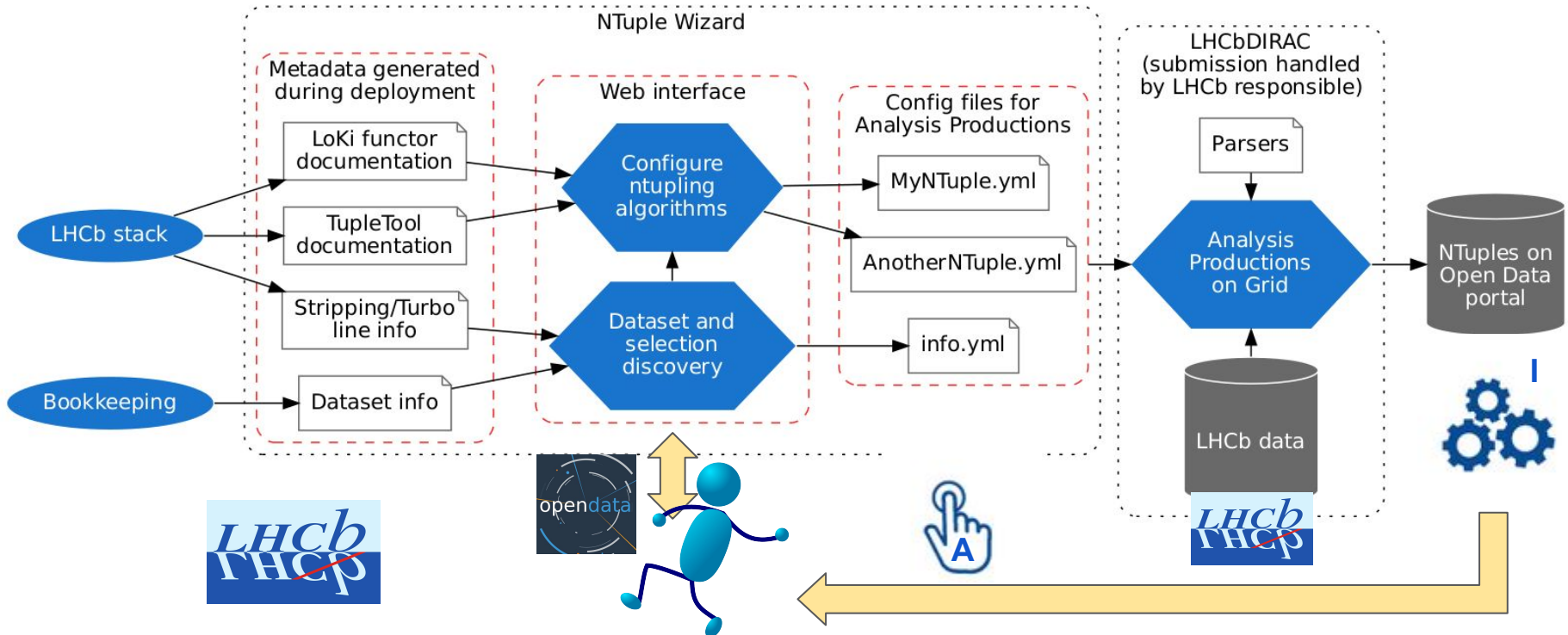
needs continuous curation by the collaborations

Needs dedicated work by the experimental collaborations (link to NFDI, PUNCH4NFDI)

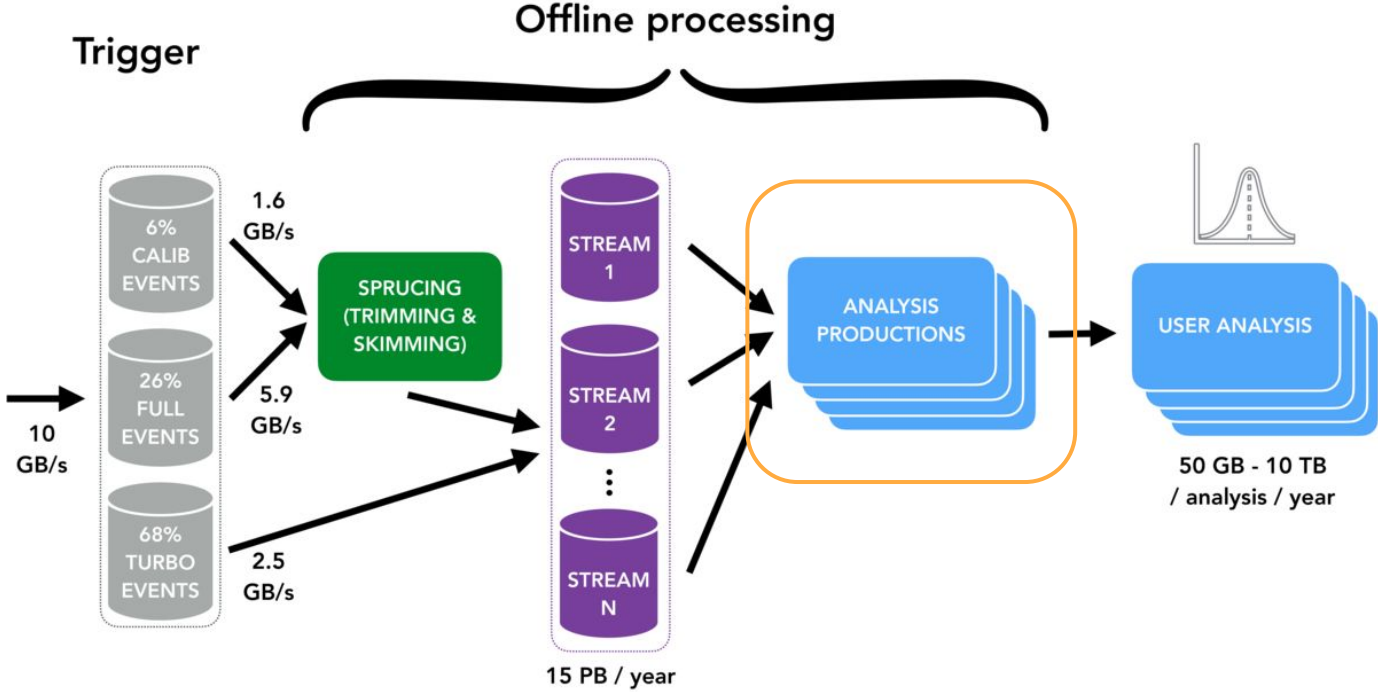
NtupleWizard overview

[[Comput Softw Big Sci 7,\(2023\) 6](#)]

- Can be used to create AP job scripts for LHCb-internal analysts



Analysis productions in LHCb



- Goals:**
- Allow analysts to creatively solve their analysis.
 - Flexible choice of tools and methods.
 - Preserve ingredients needed for **interpretation of data**

centrally managed and preserved

managed by proponents / PWG

data preparation

data interpretation

Decay search

Head (exactly): B^0 | Contains (all of): $J/\psi(1S)$ | Show only selected:

Tags (none of): $undefined-unstable$ $charge-violating$ | Stripping line

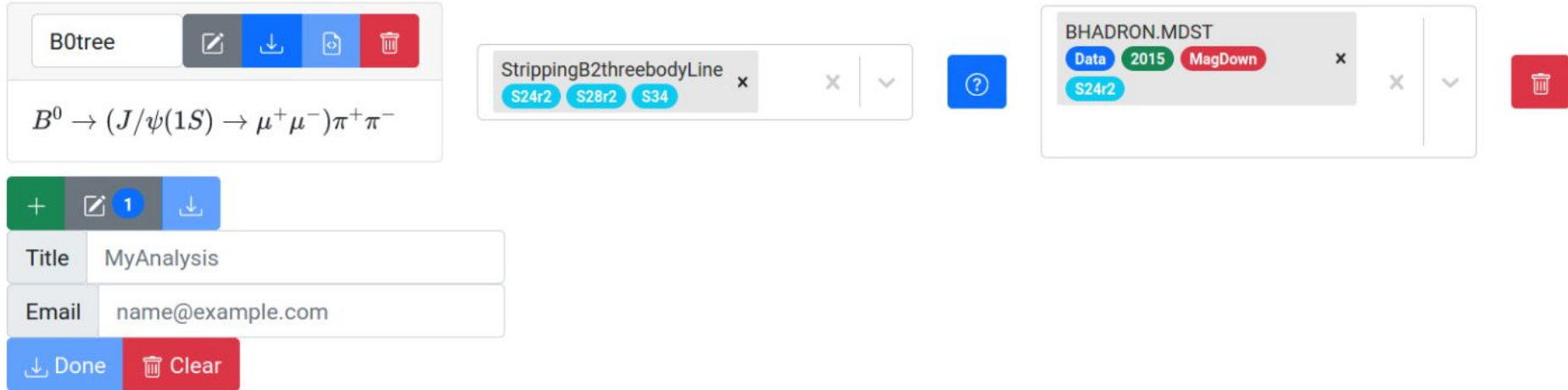
- $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ \mu^-) \pi^0$
4 Stripping lines
- $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ \mu^-) \pi^+ \pi^-$
1 Stripping line
- $B^0 \rightarrow (J/\psi(1S) \rightarrow e^+ e^-) (K_S^0 \rightarrow \pi^+ \pi^-)$
4 Stripping lines
- $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ e^-) (K_S^0 \rightarrow \pi^+ \pi^-)$
1 Stripping line **lepton-flavour-violating**
- $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ \mu^-) (K_S^0 \rightarrow \mu^+ \mu^-)$
2 Stripping lines
- $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ \mu^-) (K_S^0 \rightarrow \pi^+ \pi^-)$
8 Stripping lines

Fig. 3 Example of the decay candidate search function of the Ntuple Wizard.

Ntuple Wizard



Production configuration



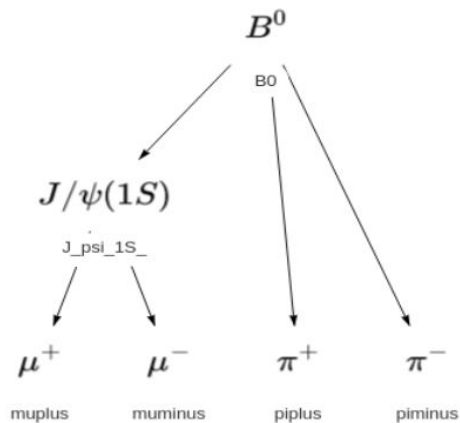
The screenshot displays the production configuration interface of the Ntuple Wizard. It features several components:

- Process Selection:** A dropdown menu showing "B0tree" with edit, download, share, and delete icons.
- Process Definition:** A text box containing the decay equation: $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+ \mu^-) \pi^+ \pi^-$.
- Stripper Selection:** A dropdown menu showing "StrippingB2threebodyLine" with a count of "x" and a dropdown arrow. Below it are three selected options: "S24r2", "S28r2", and "S34".
- Help:** A blue question mark icon.
- Dataset Selection:** A dropdown menu showing "BHADRON.MDST" with a count of "x" and a dropdown arrow. Below it are three selected options: "Data", "2015", and "MagDown".
- Form Fields:** A "Title" field with the value "MyAnalysis" and an "Email" field with the value "name@example.com".
- Buttons:** A green "+" button, a blue "1" button, a blue download icon, a blue "Done" button, and a red "Clear" button.
- Tray:** A red trash icon on the far right.

Fig. 4 Example of the data set selection and production configuration step of the Ntuple Wizard.

② Configure $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+\mu^-)\pi^+\pi^-$

8Q



Select by category

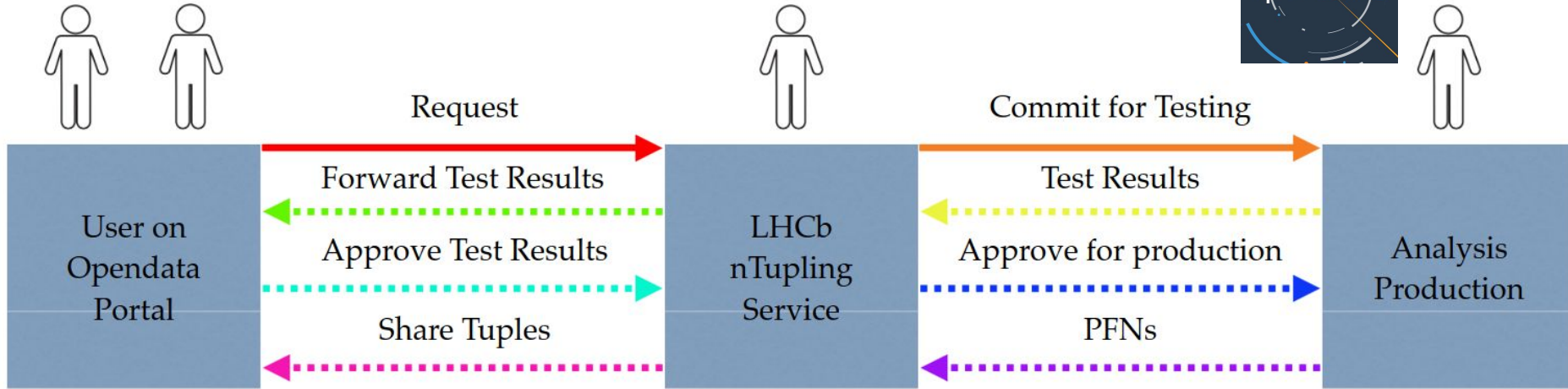
Hadron Meson Lepton X0 X+ X- Down Beauty Charm Up LongLived ShortLived Stable StableCharged Scalar Vector Spinor

Current selection: $B^0 \rightarrow (J/\psi(1S) \rightarrow \mu^+\mu^-)\pi^+\pi^-$

5 TupleTools		+
TupleToolANNPID	<input type="checkbox"/>	<input type="checkbox"/>
TupleToolEventInfo	<input type="checkbox"/>	<input type="checkbox"/>
TupleToolGeometry	<input type="checkbox"/>	<input type="checkbox"/>
TupleToolKinematic	<input type="checkbox"/>	<input type="checkbox"/>
TupleToolPid	<input type="checkbox"/>	<input type="checkbox"/>

Wizard Integration with OD Portal - WIP

Output: Plain Root NTuples



- NTuples will be cleaned up after 30* days
- Can be promoted to permanent records with DOI

- Requests are processed on LHCb internal analysis productions system on the grid

Improving documentation of the NTuples

Problem:

- Unclear what quantities are contained within the NTuple
- Entries in the NTuple can be quite cryptic and variable naming is customizable in many cases
- Relations between quantities are lost, e.g. Decay Mother -> Daughter relation



Interoperable: The (meta)data should be based on standardized vocabularies, ontologies, thesauri etc. so that it integrates with existing applications or workflows.



Reusable: Metadata and data should be well-described so that they can be replicated and/or combined in different research settings.

Well known problems, also in other fields of data science!

“**Semantic Web**”

Knowledge Graphs

[illustration from ontotext.com]

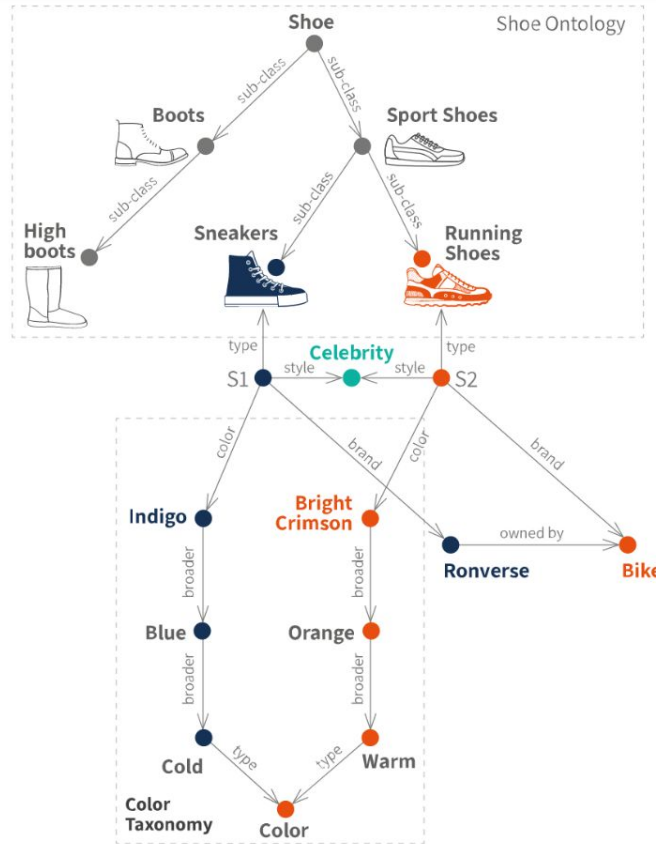
Plain Graph



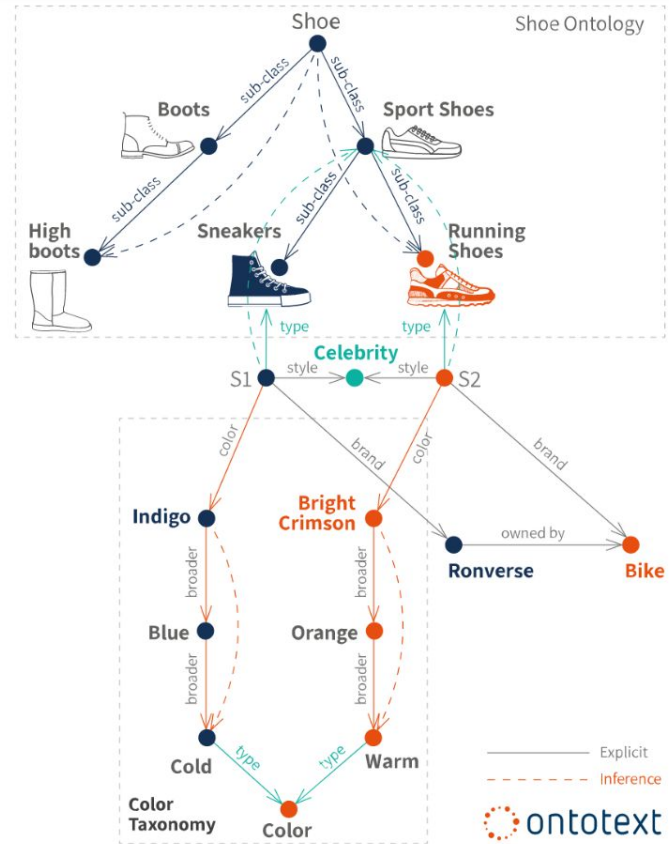
Standards for building a knowledge graph are available, which are universally digestible

Domain knowledge (Particle Physics, LHCb, Spectroscopy...) represented in Ontology

Knowledge Graph



Knowledge Graph with Inference



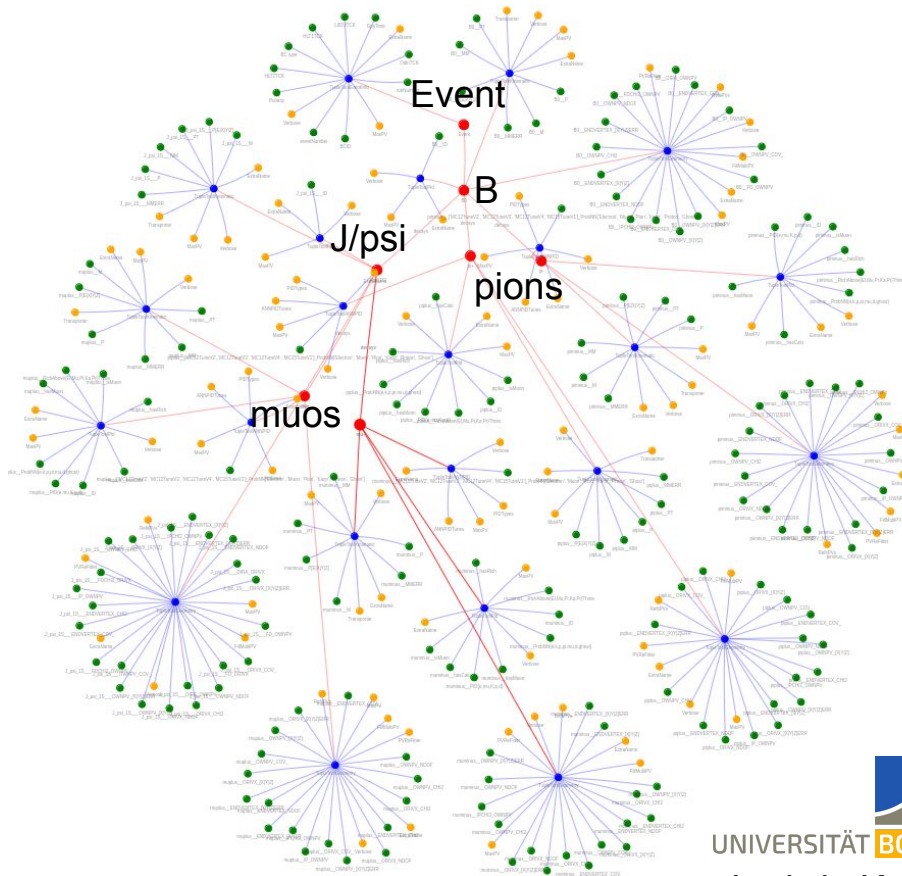
Frist Graph representation of an LHCb NTuple

Beispiel hier $B \rightarrow J/\psi \pi^+ \pi^-$

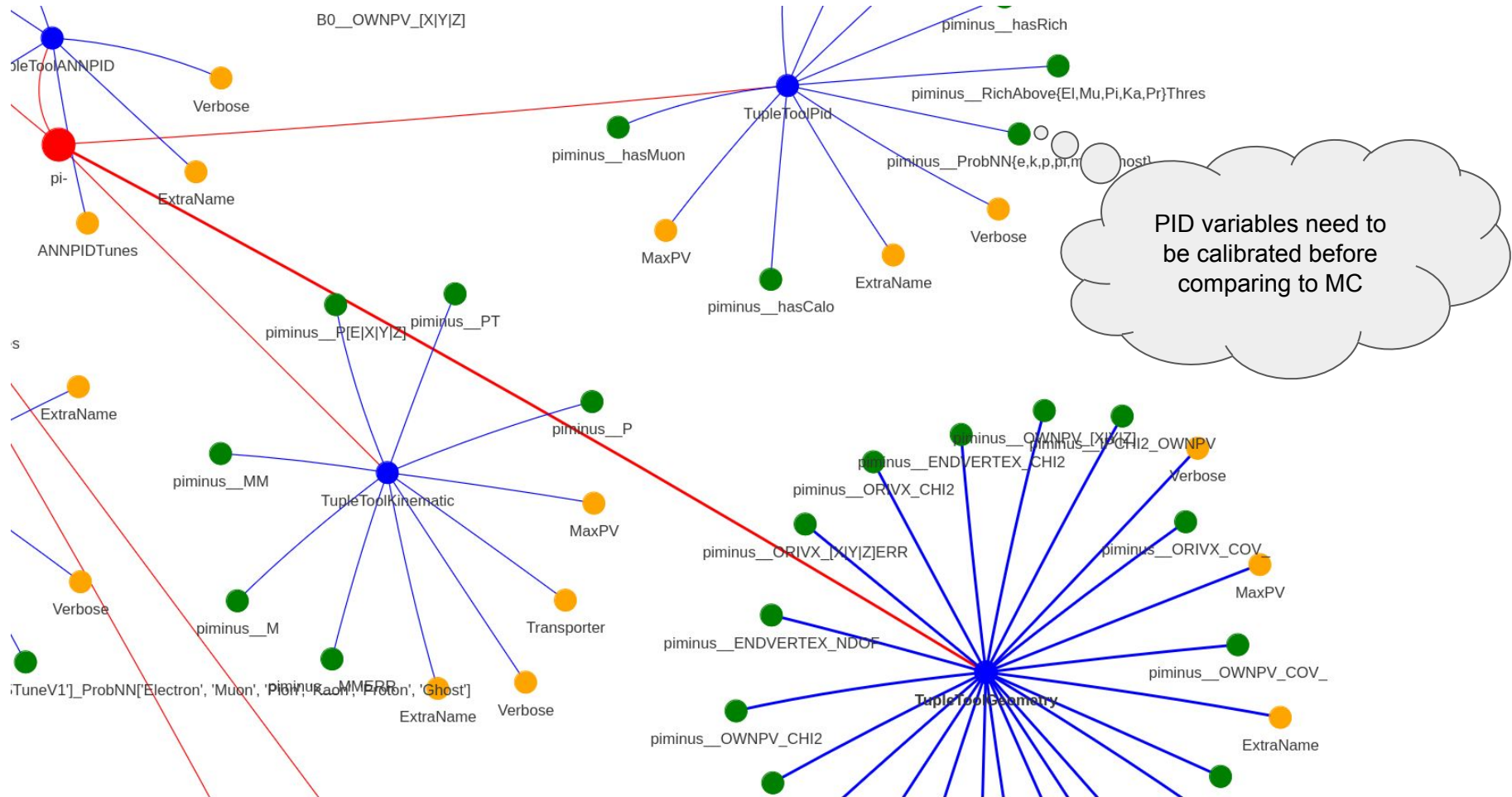
Preserving the links between the **variables** and the **tools** will allow to produce documentation for each one.

Nontrivial applications:

- Assistant to discover useful tools within Wizard
- Discover potential postprocessing/calibration steps implied by some variables

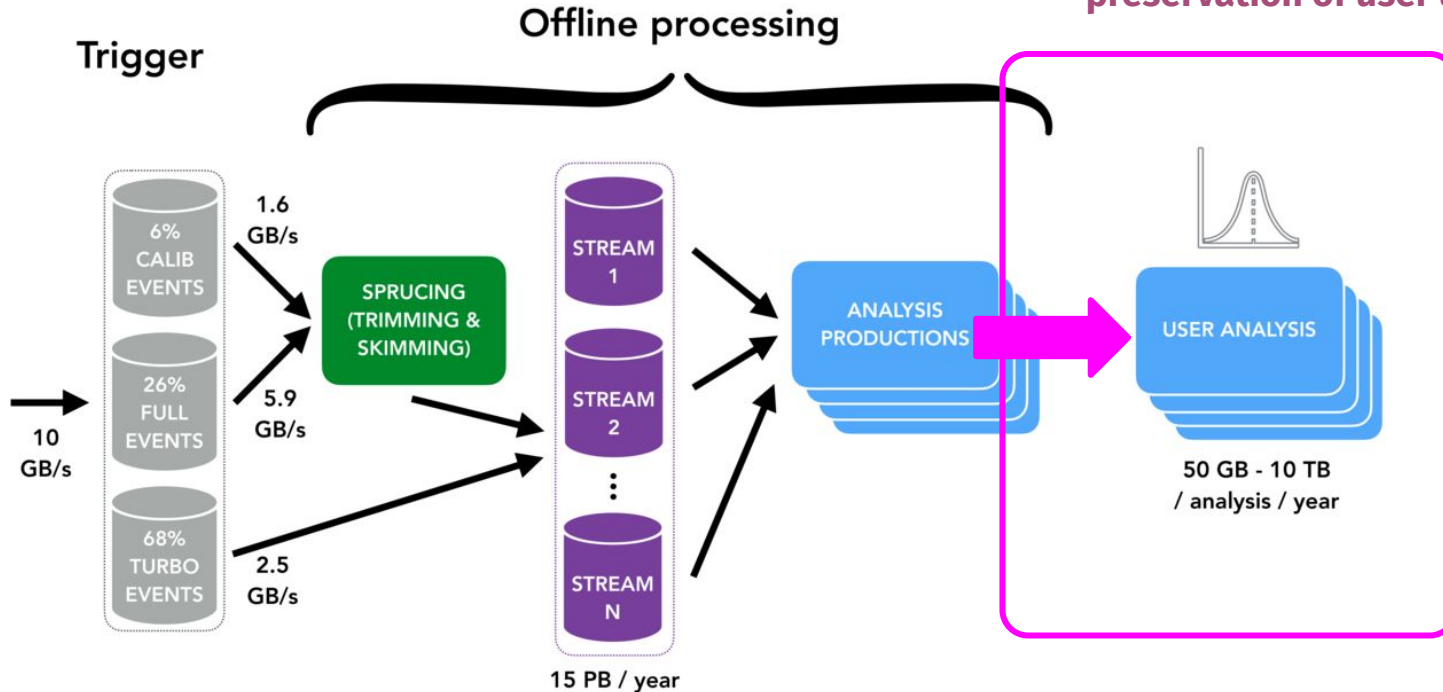


Frist Graph representation of LHCb NTuple



Preservation of User Analyses

This talk:
preservation of user analyses



Goals:

Allow analysts to creatively solve their analysis.

Flexible choice of tools and methods.

Preserve ingredients needed for **interpretation of data**

centrally managed and preserved

managed by proponents / PWG

data preparation

data interpretation

Snakemake workflow description

Set of analysis scripts, input data, and parameters + tacit knowledge how and in what order to run them

Machine readable description of workflow (similar to Makefile for software build)

- Snakemake selected as top recommendation after comparative review in 2017 (see LHCb-INT-2017-021)
- Wide use inside collaboration
- Feature complete
- Easy to get started
- Supported by CERN REANA

Snakemake is very well documented

<https://snakemake.readthedocs.io/en/stable/>

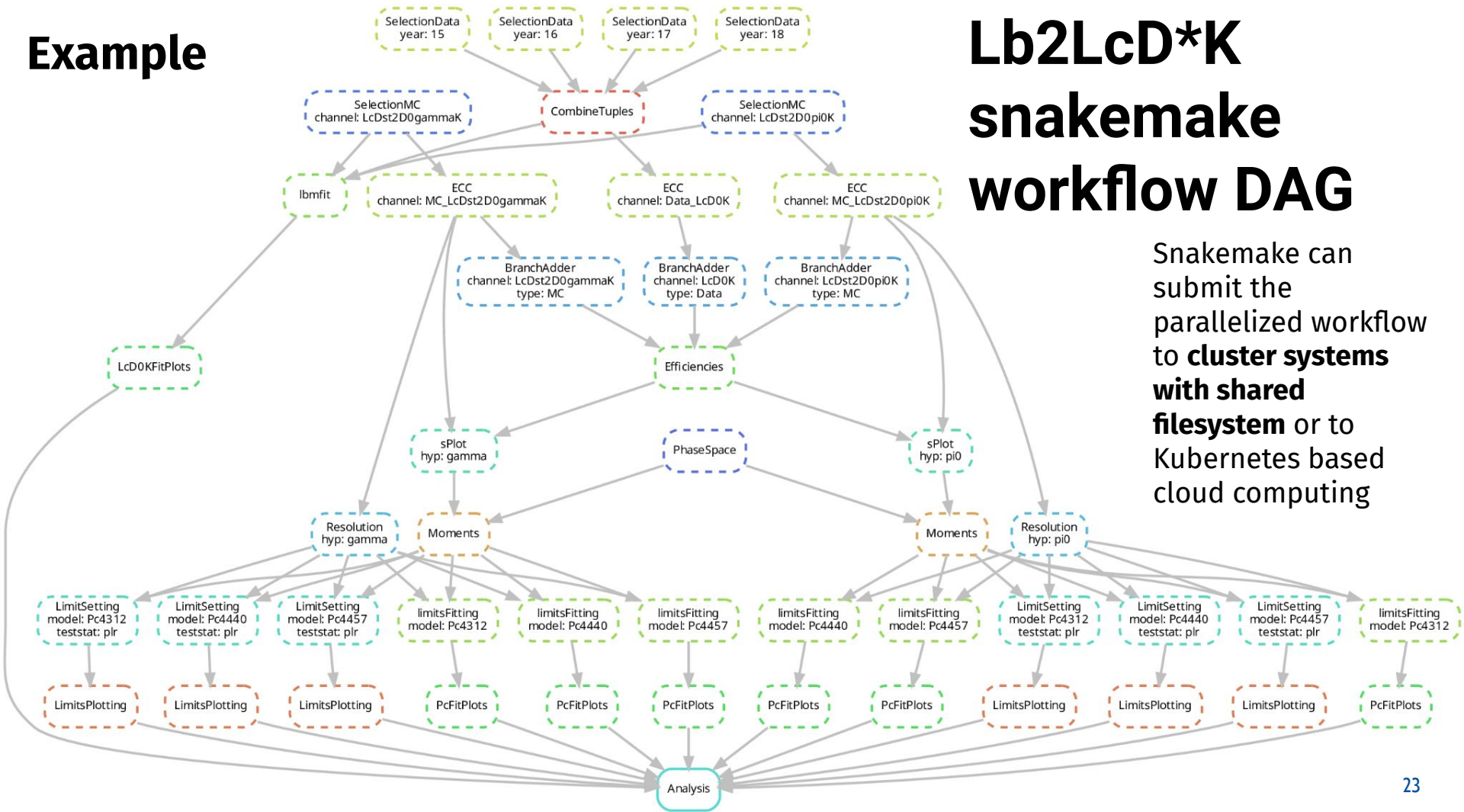
https://snakemake.readthedocs.io/en/stable/snakefiles/best_practices.html

<https://hsf-training.github.io/analysis-essentials/snakemake/README.html>

Example

Lb2LcD*K snakemake workflow DAG

Snakemake can submit the parallelized workflow to **cluster systems with shared filesystem** or to **Kubernetes based cloud computing**



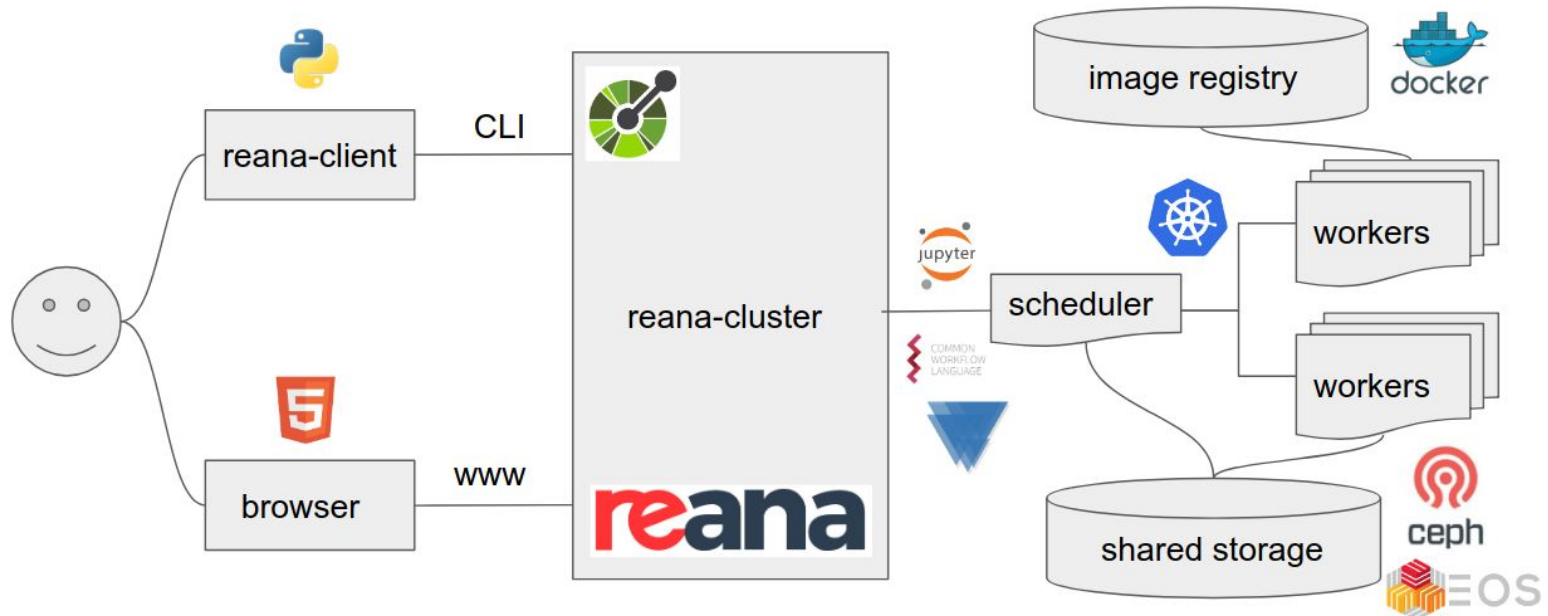
REANA: <https://docs.reana.io/>

Analysis facility dedicated to developing reproducible analyses

REANA Mattermost

Supported:
Snakemake, Common Workflow Language, Yadage

Kubernetes, HTCondor, Slurm









Deploy analysis to REANA via gitlab

<https://reana.cern.ch/profile>

reana.yaml 415 Bytes




```
1 # reana-snakemake.yaml
2 version: 0.1.0
3 inputs:
4   files:
5     - scripts/filter_tree.C
6     - scripts/helloworld.C
7   directories:
8     - workflow/snakemake
9     - output/logs
10  parameters:
11    input: workflow/snakemake/config.yaml
12  workflow:
13    type: snakemake
14    file: workflow/snakemake/Snakefile
15    resources:
16      cvmfs:
17        - lhcb.cern.ch
18        - lhcbdev.cern.ch
19        - lhcb-condb.cern.ch
20  outputs:
21    files:
22      - x.pdf
```

Your GitLab projects

-  couchWGDB
sneubert/couchWGDB
-  cookietest
sneubert/cookieetest
-  PhiKKmatrix
sneubert/phiKKmatrix
-  Analysis Workflow Template
sneubert/analysis-workflow-template
-  Snakemake Best Practices
sneubert/snakemake-best-practices
-  IR3Detector
sneubert/ir3detector

Switch on to deploy to REANA with next gitlab-ci job

Pipeline Needs Jobs 2 Tests 0

Run	External
 analysis 	 default

REANA Webinterface <https://reana.cern.ch>

✓ Analysis Workflow Template #17
Finished 5 days ago

finished in 2 min 42 sec
step 2/2

>_ Logs Workspace Specification

Step hello_world

finished

Kubernetes

gitlab-registry.cern.ch/lhcb-docker/os...

\$ ZSH_VERSION= VIRTUAL_ENV= PYT...

```
job: :  
-----  
| Welcome to ROOT 6.26/00                https://root.cern |  
| (c) 1995-2021, The ROOT Team; conception: R. Brun, F. Rademakers |  
| Built for linuxx8664gcc on Mar 05 2022, 12:03:00           |  
| From tag , 3 March 2022                                   |  
| With                                                       |  
| Try '.help', '.demo', '.license', '.credits', '.quit'/'.q'|  
-----
```

```
Processing scripts/helloworld.C("filtered_tree.root")...  
RooRealVar::Hello World from Sebastian Neubert = 0 L(-42 - 42)  
TFile**      filtered_tree.root  
TFile*       filtered_tree.root  
KEY: TTree  tree;1 tree  
Info in <TCanvas::Print>: pdf file x.pdf has been created  
(int) 0
```

Open Jupyter Notebook

Delete workflow

Interactive session
possible via Jupyter

Summary

- LHCb is completing its Run I Open data release
 - Data taken in 2011/12
 - ~ 7500 Preselections
 - ~ 800TB of data
 - MC Samples on demand
- Custom (M)DST format
- Available data \neq accessible, useable data
- NTuple Wizard will provide
 - Better access to the data, easier to use output format
 - Avoid dedicated open-data replicas of the large data sets
 - Machine readable documentation - knowledge graphs
- Preserving functional objects - such as unbinned Likelihoods - requires preserving runnable code - REANA

Backup

Policies the CERN experiments have given themselves

[CERN Open Data Policy 2020](#)

Initiated beginning 2020 by the chair of the European Commission

CERN director of research: Mandate for a working group to draft a common policy for all LHC experiments

Endorsed by the Collaboration Boards of ALICE, ATLAS, CMS and LHCb

[CERN Open Science Policy 2022](#)

Includes all experiments at CERN

<https://openscience.cern/>

Includes a wider scope of topics:

- Open access, open data, open source, open hardware
- Research integrity, research assessment
- Open infrastructure
- Training and outreach, citizen science

New Open Science Steering Board to be instantiated at CERN (S.N. LHCb delegate)

Open data policy: Level 3 data releases

Reconstructed Data (Level 3) Policy: The LHC experiments will release calibrated reconstructed data with the level of detail useful for algorithmic, performance and physics studies. The release of these data will be accompanied by provenance metadata, and by a concurrent release of appropriate simulated data samples, software, reproducible example analysis workflows, and documentation. Virtual computing environments that are compatible with the data and software will be made available. The information provided will be sufficient to allow high-quality analysis of the data including, where practical, application of the main correction factors and corresponding systematic uncertainties related to calibrations, detector reconstruction and identification. A limited level of support for users of the Level 3 Open Data will be provided on a best-effort basis by the collaborations.

Level 3 data is addressed at professional researchers

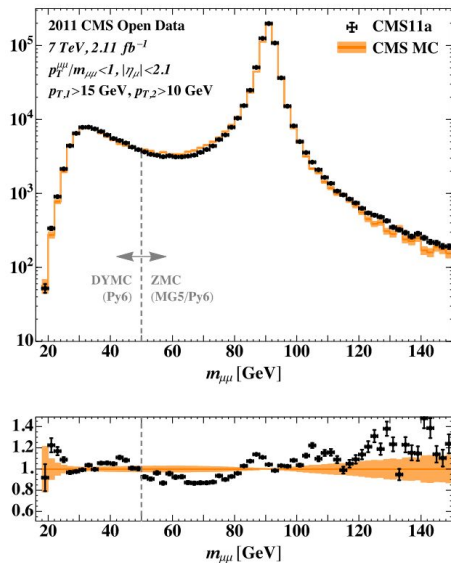
How is LHC Open Data going to be used?

No experience for LHCb data, yet.

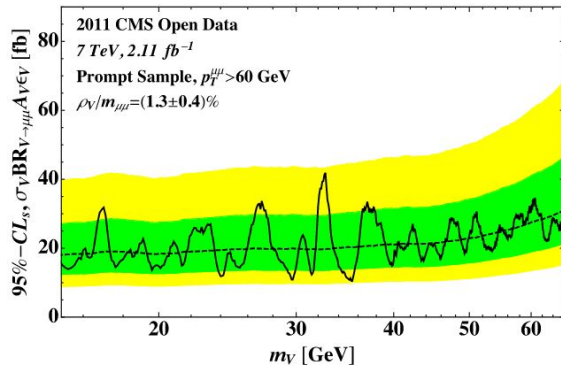
But various studies done on CMS open data.

Overview: [arXiv:2106.05726](https://arxiv.org/abs/2106.05726)

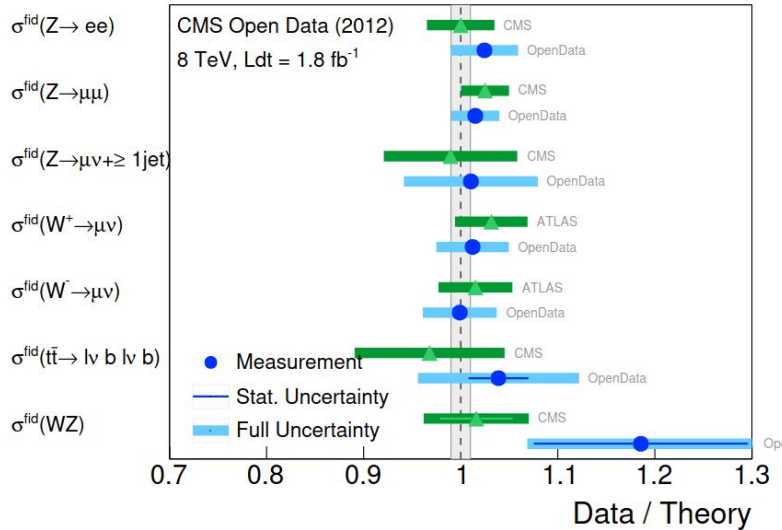
Dimuon spectrum [PRD100(2019)015021]:



Search for narrow dimuon resonances



SM cross section measurements on CMS open data [1907.08197]

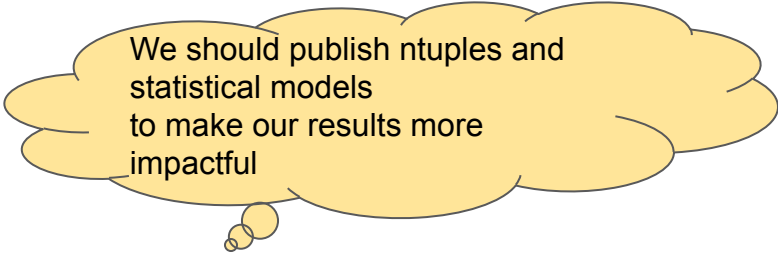


Going beyond level 3 data

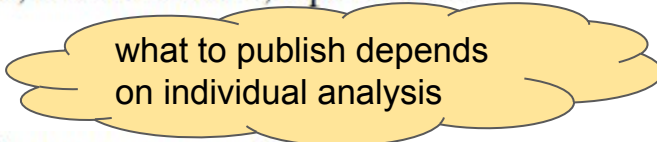
Open science and Open data policies:

5. Research integrity, reuse and reproducibility

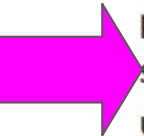
CERN is committed to ensuring the integrity of research. In order to facilitate the reuse of its research products, CERN provides infrastructures to accommodate the scale and complexity of its research outputs. Reuse and reproducibility are facilitated by practising comprehensive analysis preservation to capture relevant research objects, such as research data releases with supporting metadata, auxiliary data, linked software, reproducible analysis workflows, documentation, etc.



We should publish ntuples and statistical models to make our results more impactful



what to publish depends on individual analysis



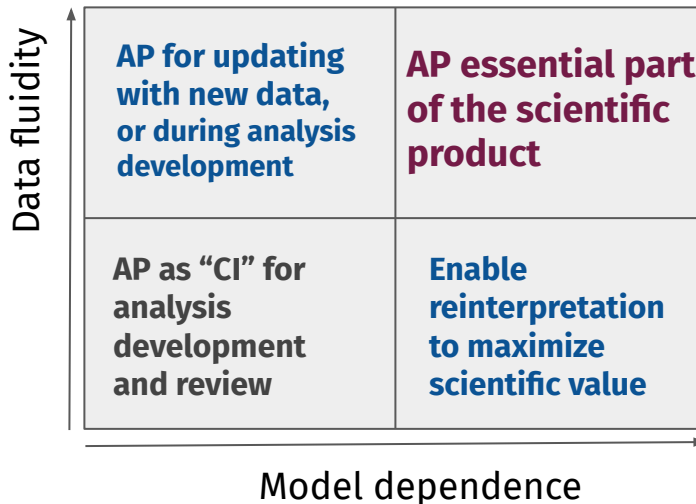
Published Results (Level 1) Policy: Peer-reviewed publications represent the primary scientific output from the experiments. In compliance with the CERN Open Access Policy, all such publications are available with Open Access, and so are available to the public. To maximise the scientific value of their publications, the experiments will make public additional information and data at the time of publication, stored in collaboration with portals such as HEPData,⁴ with selection routines stored in specialised tools. The data made available may include simplified or full binned likelihoods, as well as unbinned likelihoods based on datasets of event-level observables extracted by the analyses. Reinterpretation of published results is also made possible through analysis preservation and direct collaboration with external researchers.

Who should use full analysis preservation (AP)?

Data fluidity

- updating analysis with new data
 - **e.g. early measurements**
- control channels and their analysis for calibrations and efficiencies
 - during commissioning
 - precision measurements
- **combining measurements**

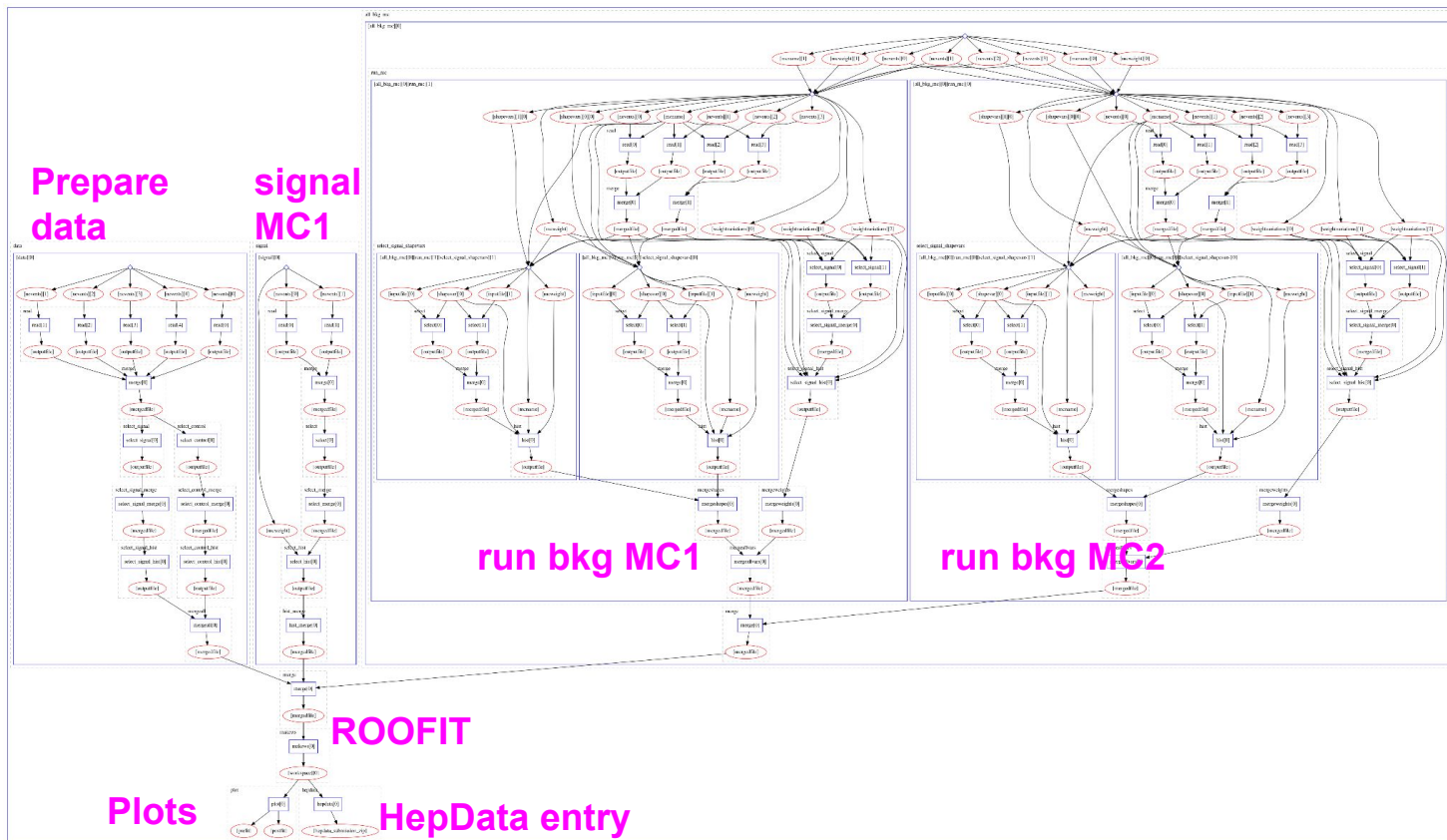
The level of detail of analysis preservation and published research products need to be decided case-by-case



Model dependence

- significant phenomenology input
 - **amplitude analyses!**
- choice of observables based on theory input
- **auxiliary inputs:** e.g. Formfactors
- MC generators / samples
- statistical methodology

<https://github.com/reanahub/reana-demo-bsm-search>



Analysis Productions

Starterkit Lesson

- Ntuple production **metadata preserved automatically!**
- Not yet supporting Run 3 DaVinci
 - Conversion will be simple from

```
lb-run DaVinci/vXrY \  
    gaudirun.py my_options.py
```

- Need to maintain link between ntuple production and analysis
- Will be able to use [apd](#) (Analysis Production Data)
- Provides PFN(s) for datasets
 - Designed to allow analyses to be rerunnable long-term

```
1 import apd
2 datasets = apd.AnalysisData("MyWG", "MyAnalysis")
3
4 rule train_bdt:
5     input:
6         data = datasets(datatype="2022", mc=False),
7         mc = datasets(datatype="2022", mc=True)
8     output:
9         fn = "classifier.pkl"
10    shell:
11        "scripts/train_bdt.py --data {' '.join(input.data)} --mc {' '.join(input.mc)}"
```

Concluding remarks: Curating Research Products

- Different scientific questions require different levels of detail in the empirical evidence.
 - **Level of model-dependence** will influence how much the experimental data can be “compressed” into a few numbers.
 - Techniques that allow reinterpretation of the data are the same as those needed to adapt to a fluid dataset
- Decisions on the level of detail of analysis preservation have to be tuned to the individual study - IMHO: avoid one-fits-all solutions
 - It is possible to support this with a small number of **generic tools, practices, and standards**
- This **data curation** requires dedicated resources.
 - Maximizing scientific value is not for free
- The technologies used to support the effort are **very useful beyond fundamental science**. Come join us!

The Open Science Philosophy (at CERN)

Recognize the **universal importance of the fundamental scientific knowledge** produced at CERN and the key role of openness in the pursuit of CERN organisational mission.

Commits to the **advancement of science** and wide dissemination of knowledge by adopting practices to make scientific research more open, global, collaborative and responsive to societal changes.

In fulfilment of the **collective moral and fiduciary responsibility** to member states and the broader global scientific community

Data collected at the LHC is a heritage to humanity.

It has been obtained through collaborative work using public funds.

Therefore, CERN is committed to preserve, curate, steward and share the data with the public.

Goals of Open Data - Maximizing Scientific Value

- Validation / reproduction of published results
- Reinterpretation of the data
 - test future theories
 - refine phenomenological models
 - use different statistical tools
- Reuse of data sets
 - Combined analyses
 - Use collected data as input for future studies
 - Algorithm development (e.g. machine learning community)
- Data mining
 - search for interesting physics in unexplored parts of the data
 - use new techniques to (re-)select data

We cannot anticipate the questions future generations might ask of this data.

require different levels
of data complexity

Open Science Landscape - Recent Trends

- Funding agencies: requests for data management plans
- Publishers: requests for data products allowing to
 - validate / reproduce results
 - reuse data for further studies

Science Community: “Data is not enough”:

- Papers with code <https://paperswithcode.com/>
- Interactive publications
- Federated infrastructures and computing/science portals (e.g. NFDI)
- Not a new realization (see e.g. DPHEP study group [2013 status report](#)) but technology (esp cloud computing, containerization) has made progress!
- Development driven especially through bioinformatics and machine learning / AI community