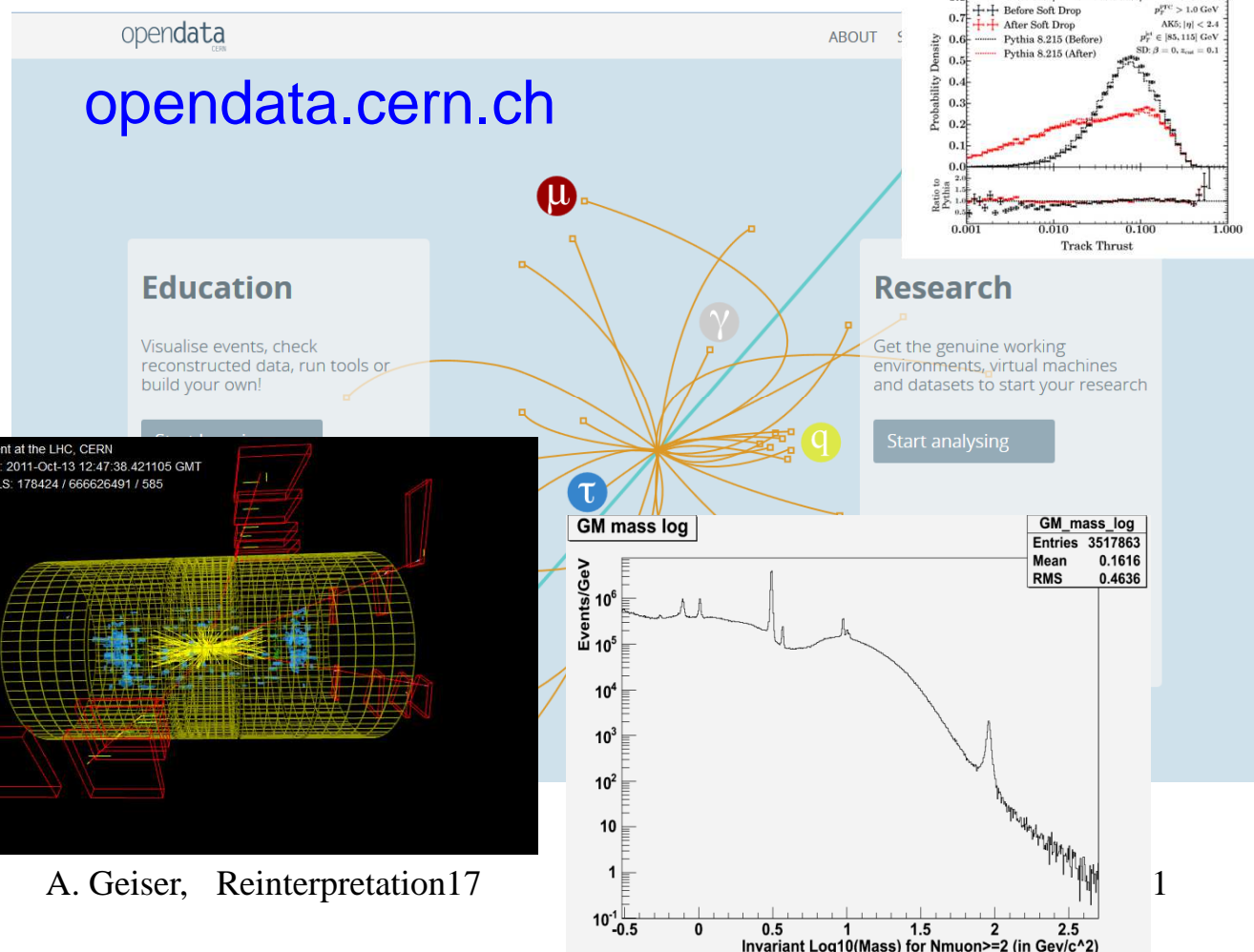# CMS Open Data in Research

Achim Geiser for the CMS collaboration (Achim.Geiser at desy.de)

Reinterpretation17 workshop, Fermilab, Batavia, IL, USA, 17.10.2017

- The vision
- The implementation opendata.cern.ch
- What it is (not)
- CMS Open Data for Research
- Status, results and prospects
- (slide) Tutorial
- Conclusions

## The Vision

- **Preserve data and knowledge (metadata)**

- **Open sharing** – data and knowledge more likely to survive if constantly used
  -> enlightened self-interest

- **Make data available to school pupils and researchers alike** - allow them e.g. to reconstruct the Higgs discovery

- (Allow CMS physicists to **recreate results** from ATLAS and vice versa -> backup)

- **Mine data to test new theories and provide crucial references**

- **Contain cost** to ~1% of operating costs -> worth the effort

**NATURE | NEWS**

عربي

# LHC plans for open data future

Researchers share results to keep them accessible.

Elizabeth Gibney

26 November 2013

PDF  Rights & Permissions

statements by
C. Diaconu (DPHEP)
M. Hildreth (DASPOS)
K. Lassila-Perini (CMS)
J. Shiers (CERN,DPHEP)
D. South (DESY, HERA)

Thomas McCauley/Lucas Taylor/CMS Collection/CERN

Data from the Large Hadron Collider, such as this decay of a Higgs boson, could be made publicly available.

- **CERN Open Data Portal:**   **opendata.cern.ch**

- Access point to growing range of data produced through research at CERN. Disseminates preserved output from various research activities, including accompanying software and documentation needed to understand and analyze the data being shared.

- Adheres to established global standards in data preservation and Open Science: the products are shared under open licenses; issued with a digital object identifier (DOI) to make them citable objects in the scientific discourse.

- Close collaboration between experiments, CERN IT and scientific information services

opendata                                          ABOUT   SEARCH   EDUCATION   RESEARCH

**Education**

Visualise events, check
reconstructed data, run tools or
build your own!

Start learning

**Research**

Get the genuine working
environments, virtual machines
and datasets to start your research

Start analysing

**this talk:**

**focus on
Research
applications**

(many educational
applications available
from all four experiments)

17.10.2017                         A. Geiser,   Reinterpretation17                              3

# The implementation:   Data and knowledge

**Research**

To analyse CMS data, a Virtual Machine with the CMS analysis environment is provided. The data can be accessed directly through the VM. In the primary datasets, no selection nor identification criteria have been applied. The 2011 data release includes simulated Monte Carlo datasets, but no simulated datasets are provided for the 2010 release.

Explore CMS >

According to the ALICE data preservation strategy, reconstructed data and Monte Carlo data as well as the analysis software and documentation needed to process them will be made available on a time scale of 5 years (for 10% of the data). Thus, the first release of ALICE research data will happen in 2018.

According to the ATLAS Data Access Policy, reconstructed data and accompanying tools will be released after reasonable embargo periods.

According to the LHCb External Data Access Policy, reconstructed data and accompanying tools will be released after reasonable embargo periods.

- **CERN Open Data Portal:**  **opendata.cern.ch**

For research purposes, specific software environments and tools need to be deployed to analyse these complex primary data. In addition to the data below, you will find instructions for setting up your working environments here

**~15 min to set up**

- **Install virtual machine**

Install your Virtual Machine >

**so far:**

- **Install CMS software**

(data in AOD format, same as used by CMS physicists)

Start analysing the data >

**only CMS released Research level data**

**-> pioneer**

# What it is not:   (in the context of this workshop)

- **not a tool to browse existing published CMS results**

  **-> use e.g.  INSpire,  arXiv,  …**


- **not a tool to (re)interpret published results by comparing with theory**

  **-> use e.g.  HEPdata,  Rivet,  …**


- **not a toolbox to recast published results into a different form**

  **-> use recasting tools  (see preceding and later contributions)**


- **not  dedicated to BSM applications (scope is general, so far dominated by SM applications, but BSM use possible and encouraged)**

# What it is: (in the context of this workshop)

- **a setup to do whatever a CMS member did, could have done or could still do with the CMS data, without any formal constraint for non-CMS members**

- **e.g. frequent theorist complaint/request:**

  **paper X does not present the results in the way I need them for my purposes, recasting is not possible for reason Y, could you please change the results?** (or the way they are presented)

- **alternative solution: stop complaining, use Open Data and change them yourself !**

  **-> (approximately) reproduce the results, or produce new ones**
  **-> modify whatever you want to modify**
  **-> compare to your favorite hypothesis**

  **real published example:   next talk by J. Thaler**

**drawback:**

- **can only be done on already released datasets (embargo period 3-4 years)**

- **will probably need a similar effort as if a CMS person or group would have done it (no magic)**

# Information about CMS Open Data

- **CERN Open Data Portal:   http://opendata.cern.ch/about/CMS**
  (see also https://twiki.cern.ch/twiki/bin/view/CMSPublic/CMSPublicData )

- **CMS data preservation, re-use and open access policy**
  **http://opendata.cern.ch/record/411**
  defines approach to data access at various levels:

- **CMS (DPHEP) Open Data levels:**

  - **Level 1 – Open access publication and additional numerical data          INSPIRE**

  - **Level 2 – Simplified data for Outreach and Education          Open Data - Education**

  - **Level 3 – Reconstructed data and the software to analyze them     Open Data - Research**

  - **Level 4 – Raw data, and the software to reconstruct and analyze them**

**CMS Open Data for Research:  AOD format** (CMS root)

talk J. Thaler

- 1st release of 28 TB of reconstructed 2010 **7 TeV pp collision data** in Nov. 2014

- 2nd release of 130 TB  of  2011 **7 TeV pp collision data** and          **~ half the respective full datasets**
  >200 TB of corresponding **MC data**     in April 2016

- 3rd release of  **8 TeV pp data + MC** (~2 PB) approved for later this year

# The challenge:   knowledge preservation

**HEP doing well with "immediate" metadata**,  **such as**

- beam conditions,  event and run numbers,  provenance information (processing and reconstruction chain, software versions) recorded together with data at time of data set creation

**doing poorly with "context" metadata,**  **such as**

- how to pick up the right objects in the data and their documentation

- how to know if there are additional selections, corrections, …

- in general, practical information needed to put data in context and analyze them: information readily available and even obvious at time of immediate data analysis, but then easily forgotten

- **Open Data helps/forces us to meet this challenge**

**Information must be collected and released together with the data**

# How we (try to) meet the challenge

- **information provided is not perfect** (and will not be) **but useful and usable**

- **information is missing for an analysis to be completed ?**

  **-> we are more than happy to take the feedback at** opendata.support@cern.ch **and provide it** (as long as we have it available ourselves)

- e.g. luminosity values for collision data recently added; cross sections for MC being added

- **being done for the first time** (in HEP) **-> learning process for everyone,** for users to learn to use these data, for us to gather and provide the necessary information from internal sources

- **we have plenty of good will but are very low on resources** -> be patient

most results presented in next slides obtained starting from scratch on CERNVM virtual machines, using windows or linux office desktop or laptop computers (can do it "from your kitchen"!), using publicly available documentation of CMS software (**open source**!). No grid jobs, no batch jobs on farm, no CMS account.

many obtained by **undergraduate students** supervised by experienced physicists -> **excellent training opportunities**!

## Open release of 2010 data in fall 2014
Using open data portal:  http://opendata.cern.ch/about/CMS

## Dimuon invariant mass distribution



CMS Paper

**JINST 7 (2012) 10002**

CMS open data
(Mu primary data set)

(on windows desktop)

~few weeks real time to set  up public example
~1 day real time to run it

Unexpected „Ridge" was observed in 2010 pp data,
**JHEP 1009 (2010) 091 (topcite 500+)**

Can be ~reproduced by selecting high multiplicity triggers in Minimum Bias dataset of 2010 open data

## CMS Paper

**JHEP 1009 (2010) 091**

## CMS open data
(summer student on office desktop)
~few weeks real time

Open Data



(d) CMS $N \geq 110$, $1.0 \mathrm{GeV/c} < p_T < 3.0 \mathrm{GeV/c}$

$R(\Delta\eta, \Delta\phi)$

2-Particle Correlation Function, $N_{sh}^{offline} > 110$, $1.0 \mathrm{GeV/c} < p\_T < 3.0 \mathrm{GeV/c}$

$H(\Delta\eta, \Delta\phi)$

use 2011 pp Open Data (2.5 fb$^{-1}$) + MC,
no usage of advanced CMS tools, simplified acceptance correction

CMS Paper

**CMS-TOP-11-013,
EPJ C73 (2013) 2339**

CMS Open Data
(O. Zenaiev)
~two months

comparison of
norm. cross sections
(O. Zenaiev)

## Problem 1: Data Certification (CMS)

From seminar at DESY, 14.2.17
A. Ustyuzhanin

Traditionally, quality of the data at CERN CMS experiment is
determined manually which requires considerable amount of
human efforts;

ML can save some of those efforts;
Data: CMS 2010B run open data;
Aim: automated classification of
Lumisections as "good" or "bad";
Features: particle flow jets, Calorimeter Jets, Photons, Muons;
The dataset was flagged by experts (3 FTE).

$$\text{Rejection Rate} = \frac{\text{Rejected}}{\text{Total quantity of samples}} \to \min;$$

$$\text{Pollution Rate} = \frac{\text{False Positive}}{\text{True Positive} + \text{False Positive}} \leq \text{const};$$

$$\text{Loss Rate} = \frac{\text{False Negative}}{\text{True Positive} + \text{False Negative}} \leq \text{const}.$$

## Results

Rejection Rate (manual work)



Loss Rate constraint

Pollution Rate constraint

expert decision

automatic decision

black zone   grey zone   white zone

0     Cut "bad"   Cut "good"     1

Andrey Ustyuzhanin                                                    12

The aim is to minimise the Manual work with
low Loss Rate ("good" classified as "bad") and
Pollution Rate ("bad" classified as "good");

~80% saving on manual work is feasible for
Pollution & Loss rate of 0.5%.
Next steps: adopt technique for 2016 data &
run in production

http://bit.ly/2I0MLiN

Andrey Ustyuzhanin                                      13

# Mine data to test new (aspects of) theories

## Exposing the QCD Splitting Function with CMS Open Data

Andrew Larkoski, Simone Marzani, Jesse Thaler, Aashish Tripathee, Wei Xue

+ some CMS support  (S. Rappoccio)

**Phys Rev Lett 119 (2017)  132003**

Apr 17, 2017 - 7 pages

MIT-CTP-4891
e-Print: arXiv:1704.05066 [hep-ph] | PDF

**first ever published
CMS Open Data results**

## Jet Substructure Studies with CMS Open Data

Aashish Tripathee, Wei Xue, Andrew Larkoski, Simone Marzani, Jesse Thaler

Apr 19, 2017 - 35 pages

**Phys Rev D96 (2017)  074003**

MIT-CTP-4890
e-Print: arXiv:1704.05842 [hep-ph] | PDF

**observed jet substructure
agrees with predictions
from first principles using
QCD splitting functions**

**-> next talk**

# Mine data to test new (aspects of) theories

**Open Data analysis example in preparation:**

**Search in leptonic channels for heavy resonances decaying to long-lived neutral particles**

**JHEP 1302 (2013) 085
CMS-EXO-11-101**

Theory: "Hidden Valley" Search

in practice: Search for leptons
with "long-distance" displacement
from primary vertex, originating from
decay of X particles

**-> get limits e.g. on H->XX
(details see paper)**

**some practical aspects below**

# Open Data Tutorial

Slides, not online (would take too long), but written such that it can be tried immediately

Focus on:

Dimuon mass spectrum

(simple, exists, works)

only prerequisite: know a bit of Linux and ROOT

Displaced Lepton Search

(conceptual, in preparation,

expose challenges)

# Open Data Tutorial

start your favourite laptop, desktop, …    windows, linux or MacOS
(at least **2 GB memory**, administrator rights or VirtualBox preinstalled)

**with any web browser:**    **opendata.cern.ch**
(see also https://twiki.cern.ch/twiki/bin/view/CMSPublic/CMSPublicData )

Portal appearance
might change soon.
Content will
stay/be extended.



choose Research

("Start analysing")

side remarks:

VM is faster on windows than on linux!

Tutorial will work (almost ) anywhere
(except on the Fermilab guest  network ..)

# Open Data Tutorial



choose

"Install your Virtual Machine"

# Open Data Tutorial



choose

"CMS Virtual Machines"

# Open Data Tutorial



different Virtual Machines for 2010 and 2011 data !

-> you launch the Virtual Machine & graphical user interface

**1** choose "2010"

**2** if not yet done:

download /install VirtualBox

(see FAQ at bottom of page)

**3** download CMS VM Image

& double click it

# Open Data Tutorial



further down the same page

**How to Test & Validate?**

The validation procedure tests that the CMS environment is installed and operational on your virtual machine, and that you have access to the ROOT files. You may skip this step if you want, and head straight to Getting Started with CMS data. However, these steps give you a quick introduction to the CMS environment.

**2**

Set up the CMS environment and run a demo analyzer

Open a terminal with the X terminal emulator (an icon bottom-left of the VM screen)
Execute the following command; this command builds the local release area (the directory structure) for CMSSW, and only needs to be run once:

**1**
```
cmsrel CMSSW_4_2_8
```

Change to the `CMSSW_4_2_8/src/` directory:
```
cd CMSSW_4_2_8/src/
```

Then, run the following command to create the CMS runtime variables:
```
cmsenv
```

Create a working directory for the demo analyzer, change to that directory and create a "skeleton" for the analyzer:
```
mkdir Demo
cd Demo
mkedanlzr DemoAnalyzer
```

Compile the code:
```
cd DemoAnalyzer
scram b
```

Change the file name in the configuration file `demoanalyzer_cfg.py` in the DemoAnalyzer directory: i.e. replace `file:myfile.root` with `root://eospublic.cern.ch` `//eos/opendata/cms/Run2010B/Mu/AOD/Apr21ReReco-v1/0000/00459D48-EB70-E011-AF09-90E6BA19A252.root`
Change the max number of events to 10 (i.e change -1 to 10 in `process.maxEvents = cms.untracked.PSet( input = cms.untracked.int32(-1)` )
Move two directories back using:
```
cd ../..
```

And then run
```
cmsRun Demo/DemoAnalyzer/demoanalyzer_cfg.py
```

Test & Validate:
(do not skip)

**1**

install relevant CMS software
(CMSSW_4_2_8 for 2010 data)

follow instructions

-> you will have technically run your first CMS job ☺
-> continue with **2** "Getting started with CMS data"

# Open Data Tutorial



you will already have done this

manually inspect the content of a CMS AOD ROOT file (located at CERN) in order to get a "feel"

(follow instructions)

# Open Data Tutorial
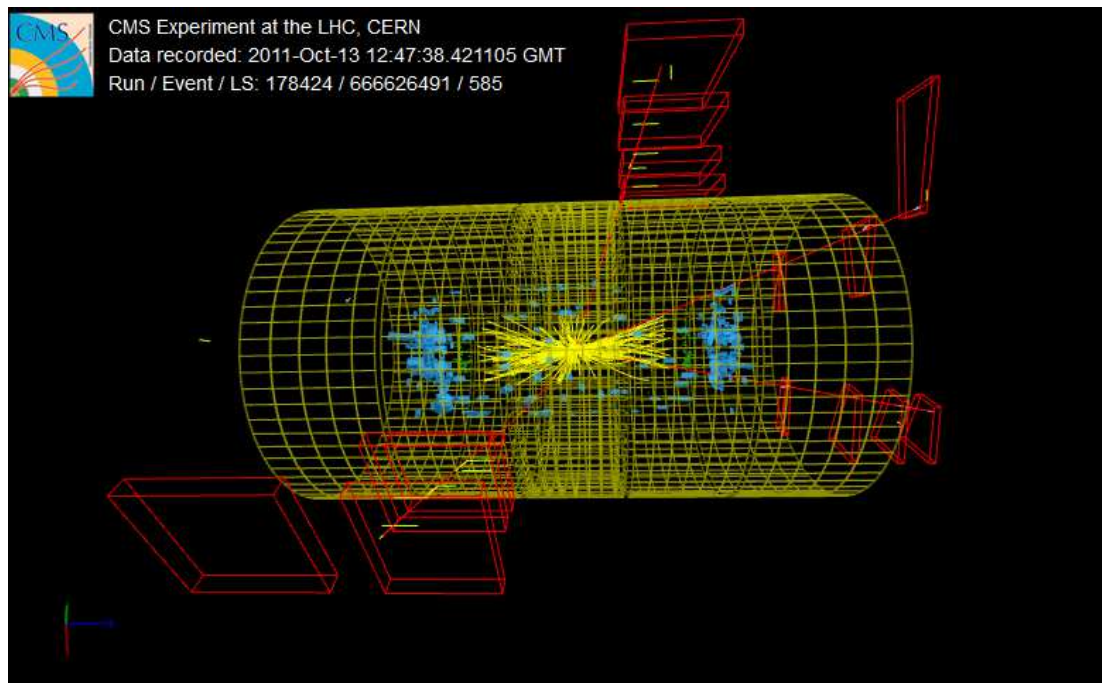


start the ROOT browser

(previous experience with ROOT is helpful)

inspect the variables

# Open Data Tutorial

**Intermezzo:** at this stage (or at any other time), might want to looki at some CMS **event displays** on your standard browser (**no VM needed**).

e.g. real Higgs -> 4 muon candidate



CMS Experiment at the LHC, CERN
Data recorded: 2011-Oct-13 12:47:38.421105 GMT
Run / Event / LS: 178424 / 666626491 / 585

**in browser:**

go to **education** part

-> "visualize events"

and follow instructions

consider also https://twiki.cern.ch/twiki/bin/view/CMSPublic/CMSPublicData

# Open Data Tutorial

continue further down the "Getting Started" page



for this tutorial:

choose Option A

# Open Data Tutorial



you have now reached the dimuon mass spectrum example
(works technically like "Test & Validate")

follow instructions to download and run it

-> inspect histograms on resulting root file
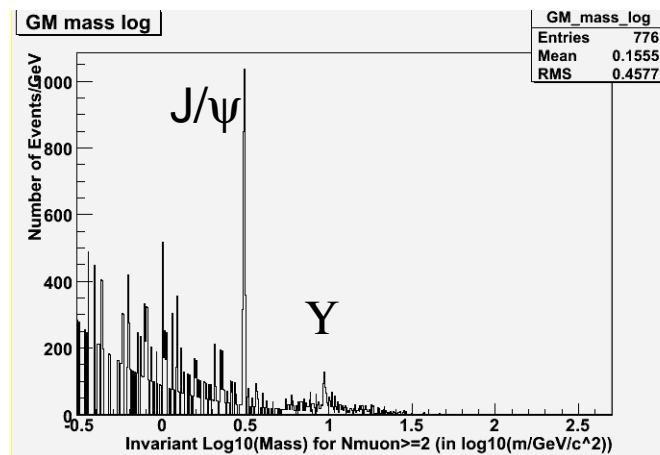
# Open Data Tutorial

-> inspect histograms on resulting ROOT file

with 10000
input events

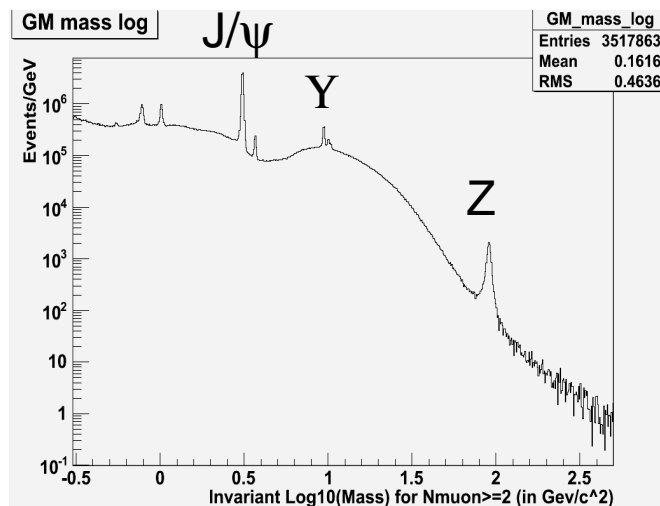(default,
runs ~few minutes)



there you are!

now ready to
edit example
and
add/change
what you like

with full 2010
Mu dataset,
set logy option

(I/O to CERN
 takes ~1 day
 at ~10 Mbit/s)

# Open Data Tutorial: conceptual extension

e.g. for **"displaced lepton" analysis**:      (example in preparation)

repeat previous exercise (or other) with **2011 data**  (fully documented)

learn how to **find, access and treat physics objects** you are interested in
(e.g. select non-vertex-associated leptons, possibly do revertexing (tools exist!))

learn how to **select the most relevant dataset(s)**
and how to **identify and treat the most relevant triggers**

> we can (try to) help in all these steps

check availability of **suitable MC sets**   (among the already released ones),
possibly **reweight for your application**;
if necessary, **try to generate your own exotic signal set** (not at all documented
yet, being tried; if successful, might become possible with full detector simulation, but
will remain a challenge! Or try to use external simplified simulation tools)

possibly extract your **personal reduced data set**  (not necessary, but allowed)

**do your analysis** and
**publish** with your (non-CMS) name on it, cite DOIs of CMS open datasets used

# Extended Vision (for discussion)

my **personal extension of initial vision:**

(for discussion, not a collaboration statement)


**with** **~1% of additional resources** **aim to achieve**
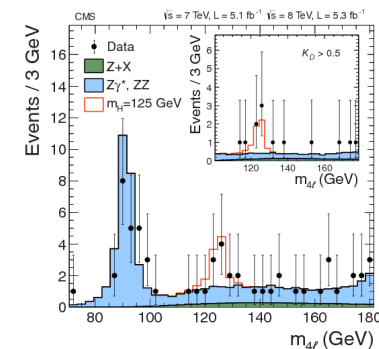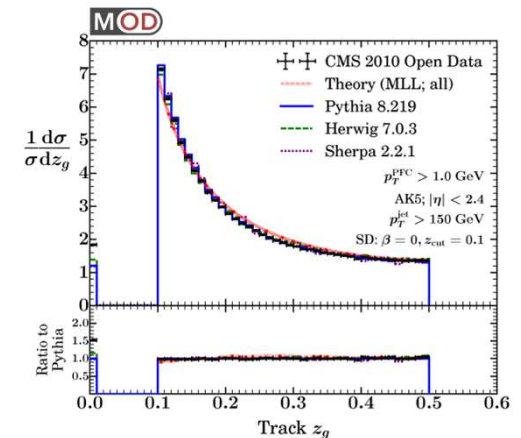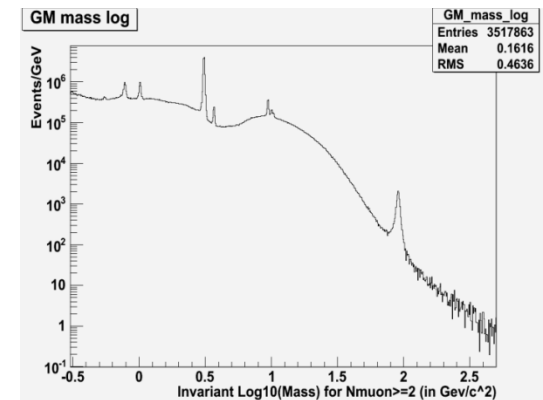
**~10% additional scientific output** (physics papers)

from both external and internal use of preserved/open data

over lifetime of experiment + 10-20 years

# Conclusions

- Open Data release of full CMS 2010 run B data and 2011 run A **data + MC** available on http://opendata.cern.ch

- well prepared by CERN and CMS IT and Open Data teams: anybody can use it and it works

- contains benchmark physics analysis and validation examples (more coming)

- involves nontrivial challenges being worked on

- first physics results from 2010 open data just published by group of theorists from MIT
  **-> hopefully start of long and fruitful future of full exploitation of High Energy Physics data beyond actual collaborations**

- also used for machine learning

- upcoming 2012 data release: on the way towards public reconstruction of the Higgs discovery

- **feel encouraged to use it for your purposes!**

17.10.2017                      A. Geiser,   Reinterpretation17                      30

# Backup

# Feedback to community

**Jet Substructure Studies with CMS Open Data**

Aashish Tripathee, Wei Xue, Andrew Larkoski, Simone Marzani, Jesse Thaler

Apr 19, 2017 - 35 pages

MIT-CTP-4890

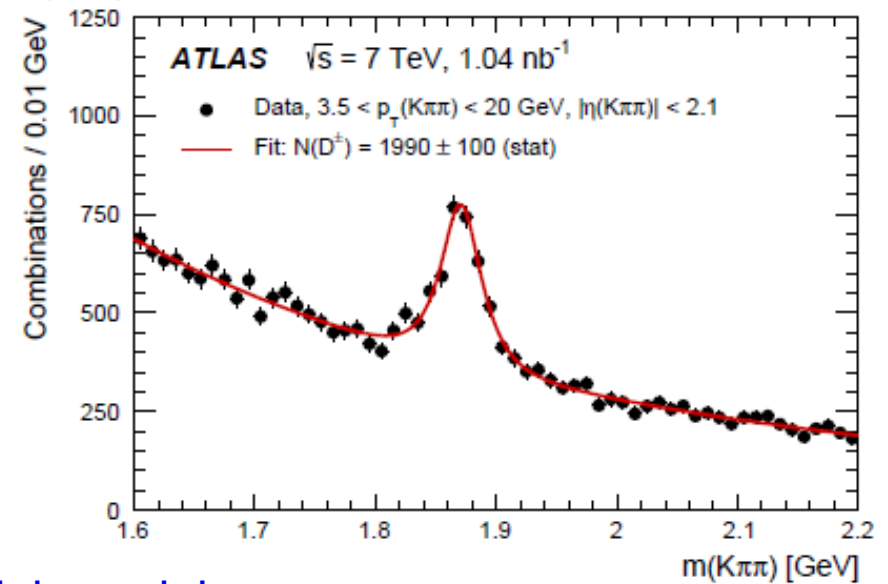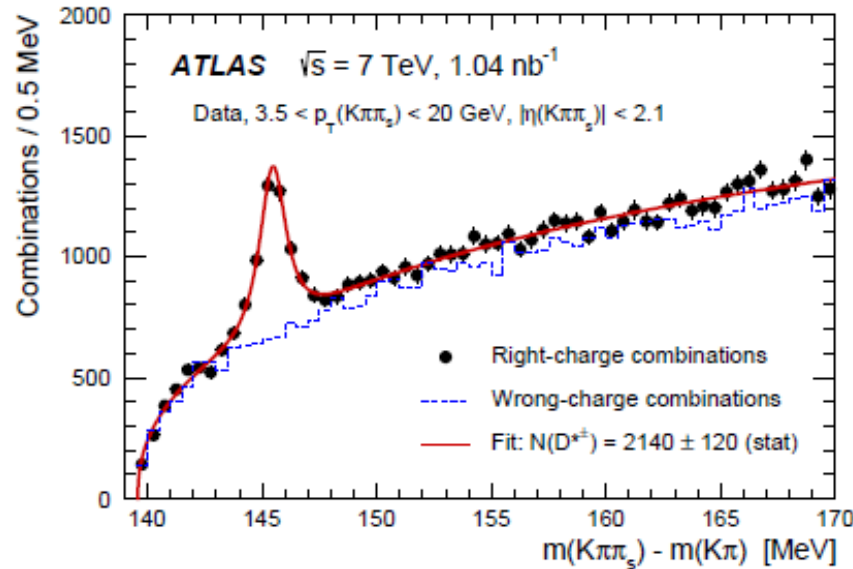e-Print: arXiv:1704.05842 [hep-ph] | PDF

see also talk J. Thaler
+ discussion

Contains section with Advice to community, Challenges, and Recommendations
(see there)

Releases of 2011 CMS data+MC "exciting"
-> properly evaluate detector systematics

Conclusions: "We hope our experience motivates the LHC
collaborations to further their investment in public data release and
encourages the particle physics community to exploit the scientific
potential of open datasets"

Nucl. Phys. B907 (2016) 717



B. Sheeran, N. Stefaniuk, work in progress