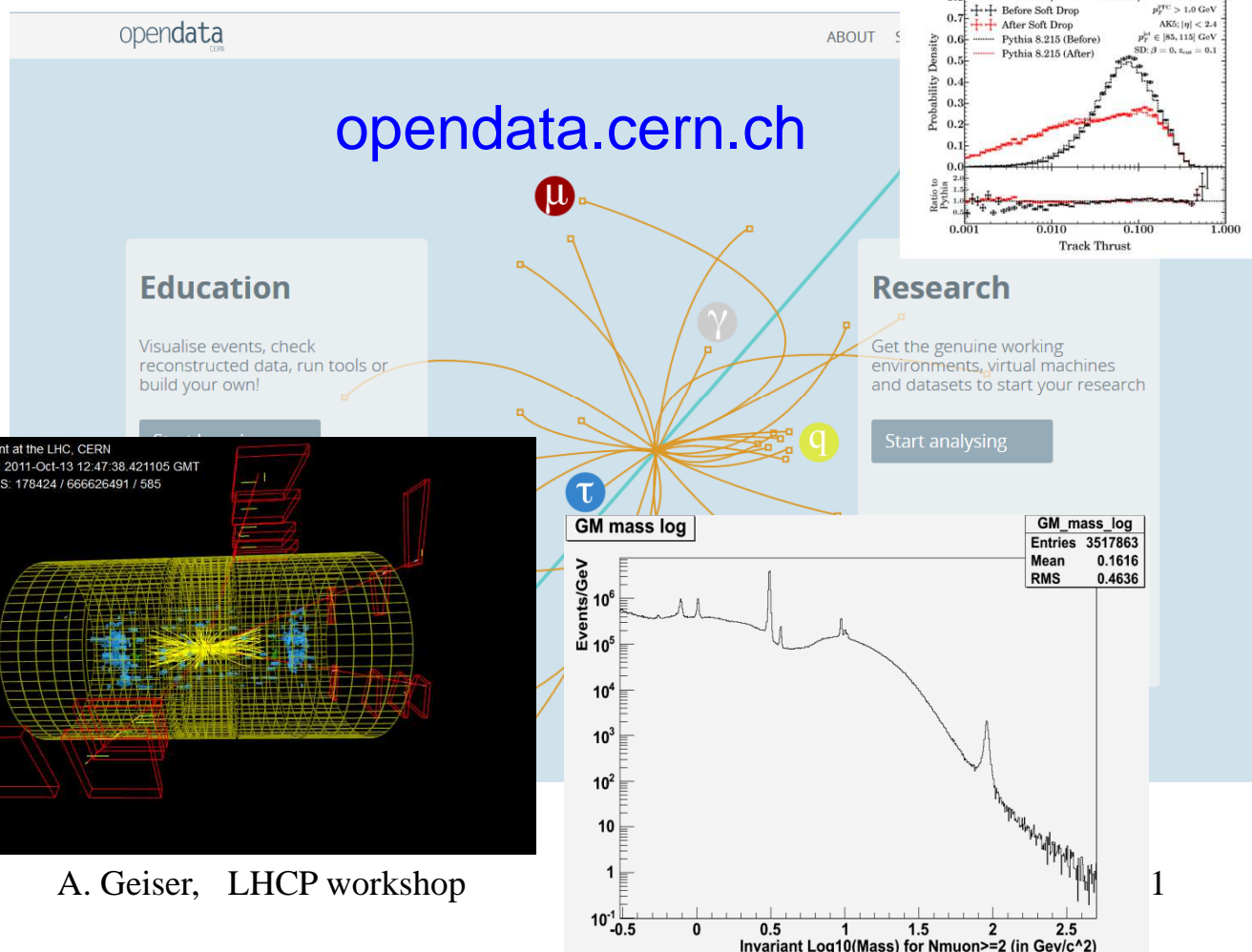


CMS Open Data in Research

Achim Geiser for the CMS collaboration

Outreach session of LHCP workshop, Shanghai, China, 15.05.2017

- The vision
- The implementation
- CMS Open Data for Research
- Status, results and prospects
- Conclusions



15.05.17

A. Geiser, LHCP workshop

1

LHC plans for open data future

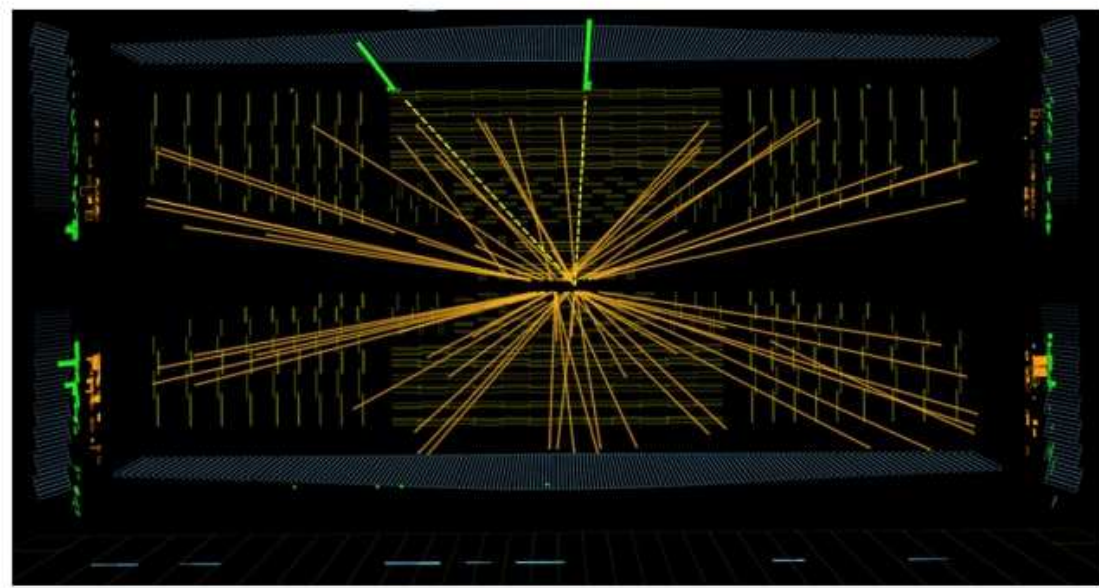
Researchers share results to keep them accessible.

Elizabeth Gibney

26 November 2013

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

PDF Rights & Permissions



Thomas McAuley/Lucas Taylor/CMS Collection/CERN

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

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

The implementation: Open sharing

- **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**

The screenshot shows the CERN Open Data Portal interface. At the top, the logo 'opendata CERN' is on the left, and navigation links 'ABOUT SEARCH EDUCATION RESEARCH' are on the right. The main content area is split into two columns: 'Education' on the left and 'Research' on the right. The 'Education' section includes the text 'Visualise events, check reconstructed data, run tools or build your own!' and a 'Start learning' button. The 'Research' section includes 'Get the genuine working environments, virtual machines and datasets to start your research' and a 'Start analysing' button. A central network diagram features a central hub with numerous lines radiating outwards to various nodes, some of which are labeled with Greek letters: μ , γ , τ , e , and q . A diagonal teal line runs across the diagram. At the bottom of the screenshot, the text 'A. Geiser, LHCP workshop' is visible.


this talk:

**focus on
Research
applications**

(many educational applications available from all four experiments)


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.

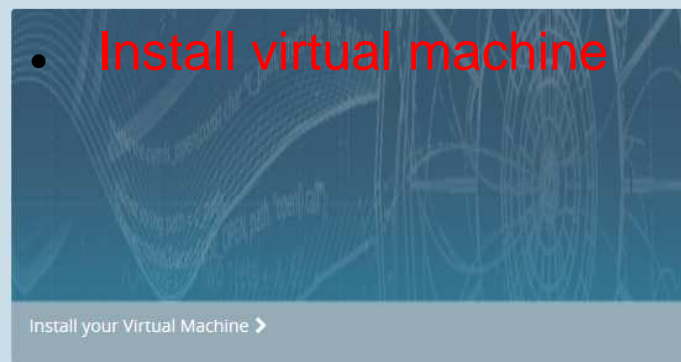


According to the LHCb External Data Access Policy, reconstructed data and accompanying tools will be released after reasonable embargo periods.

• 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



- Install virtual machine

Install your Virtual Machine >



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

Start analysing the data >

~15 min to set up

so far:

only CMS
released
Research
level data

-> pioneer

Open Data in CMS

- **CERN Open Data Portal:** <http://opendata.cern.ch/about/CMS>
- **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:

- 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 >200 TB of corresponding MC data in April 2016 **~ half the respective full datasets**
- 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 still useful and usable**
- **information is missing for an analysis to be completed ?**
(e.g. currently luminosity values for collision data and cross sections for MC)
-> **we are more than happy to take the feedback at opendata.support@cern.ch and provide them** (as long as we have them available ourselves)
- **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 are full of good will but 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. No grid jobs, no batch jobs on farm, no CMS account needed.

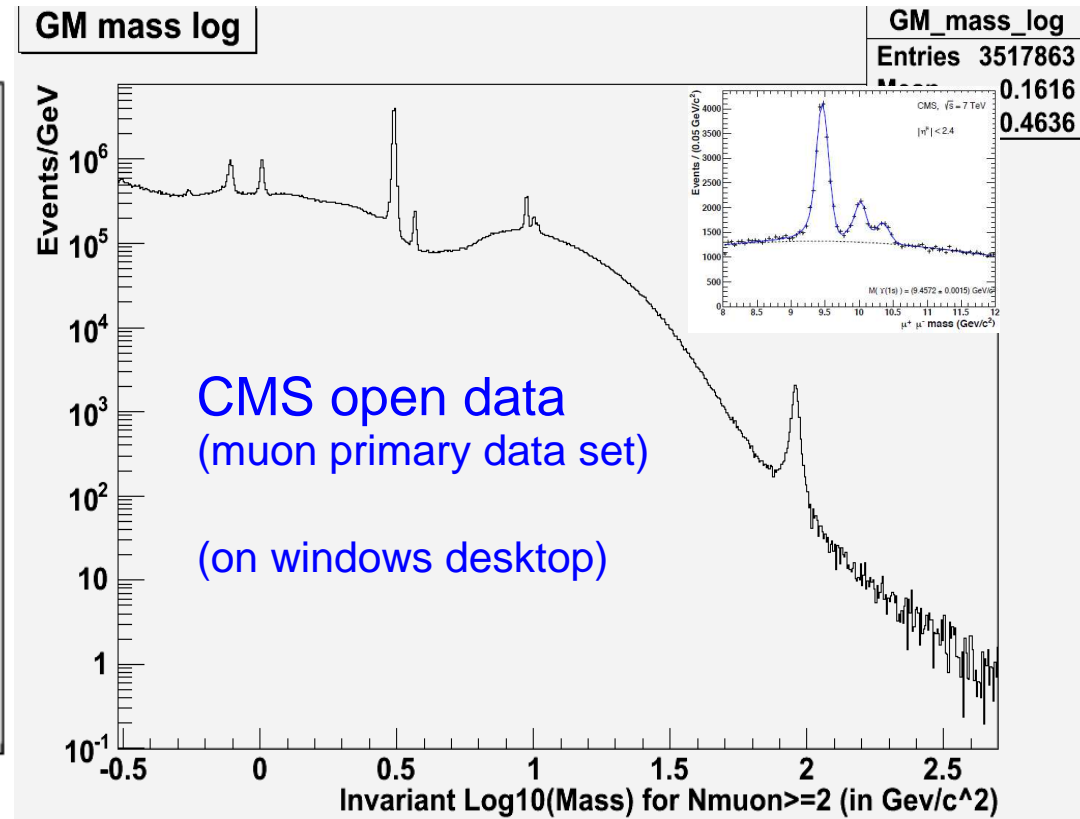
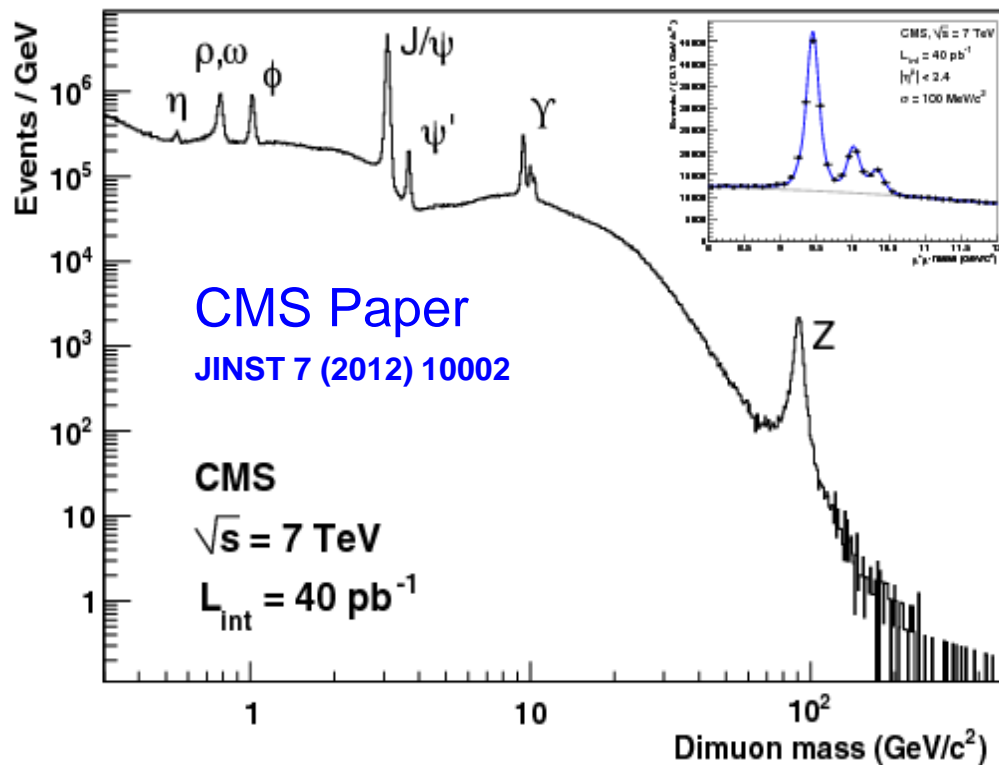
Many obtained by **undergraduate students** supervised by experienced physicists
-> **excellent training opportunities!**

Provide references: validation/benchmarking/analysis examples

Open release of 2010 data in fall 2014

Using open data portal: <http://opendata.cern.ch/about/CMS>

Dimuon invariant mass distribution



Open Data benchmark analysis: “Ridge”

A. Nassirpour, summer student 2016

<https://indico.desy.de/getFile.py/access?contribId=4&resId=0&materialId=slides&confId=15932>

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

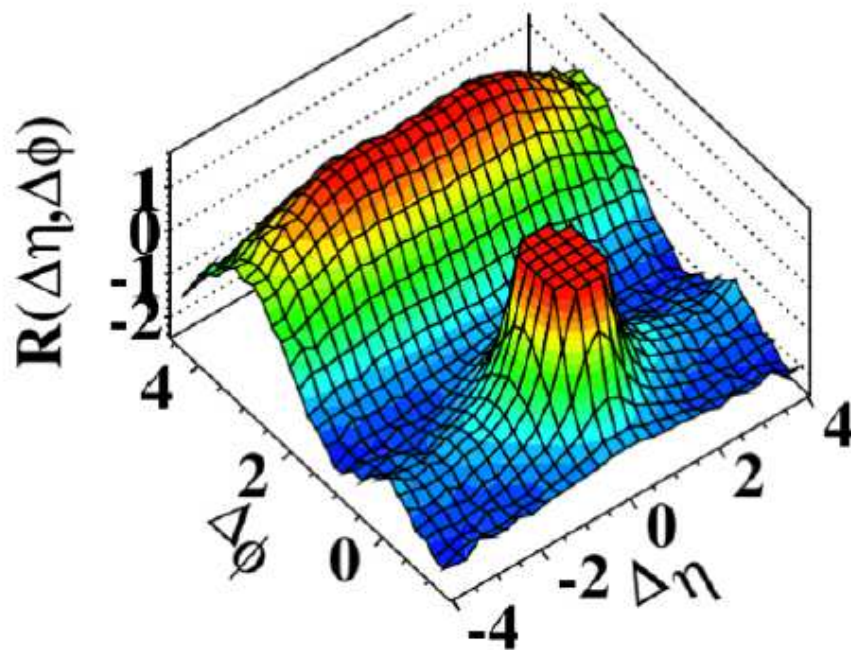
CMS Paper

JHEP 1009 (2010) 091

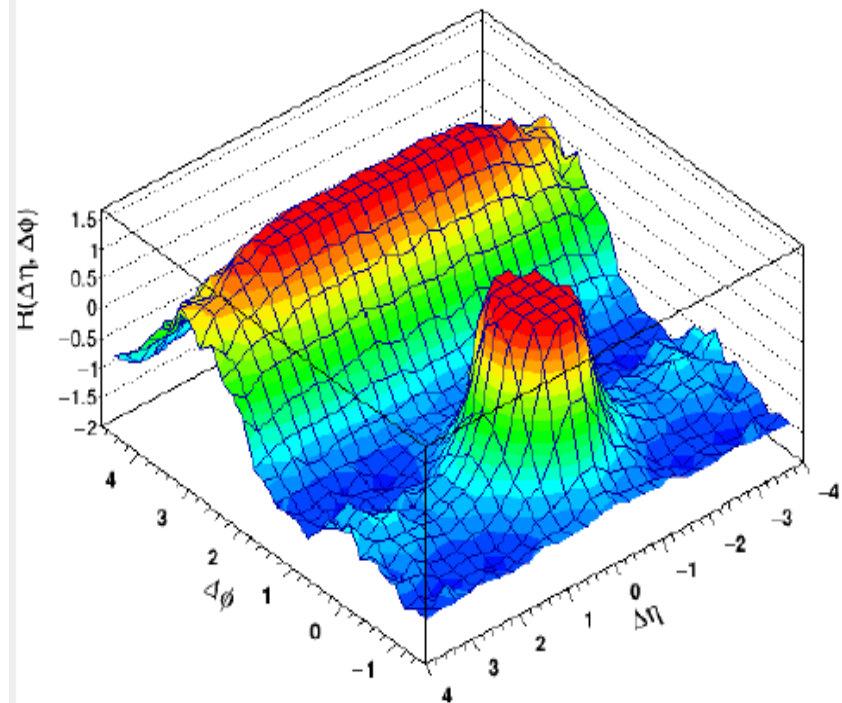
CMS open data

(summer student on office desktop)

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



2-Particle Correlation Function, $N_{ch}^{offline} > 110$, $1.0\text{GeV}/c < p_T < 3.0\text{GeV}/c$



Mine data to test new (aspects of) theories

Exposing the QCD Splitting Function with CMS Open Data

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

+ some CMS support (S. Rappoccio)

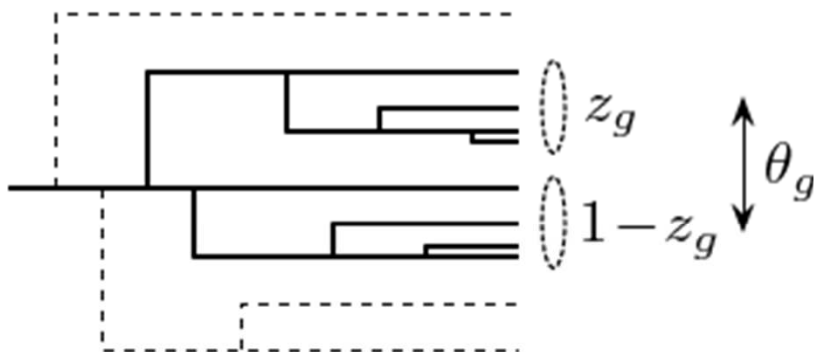
Apr 17, 2017 - 7 pages

MIT-CTP-4891

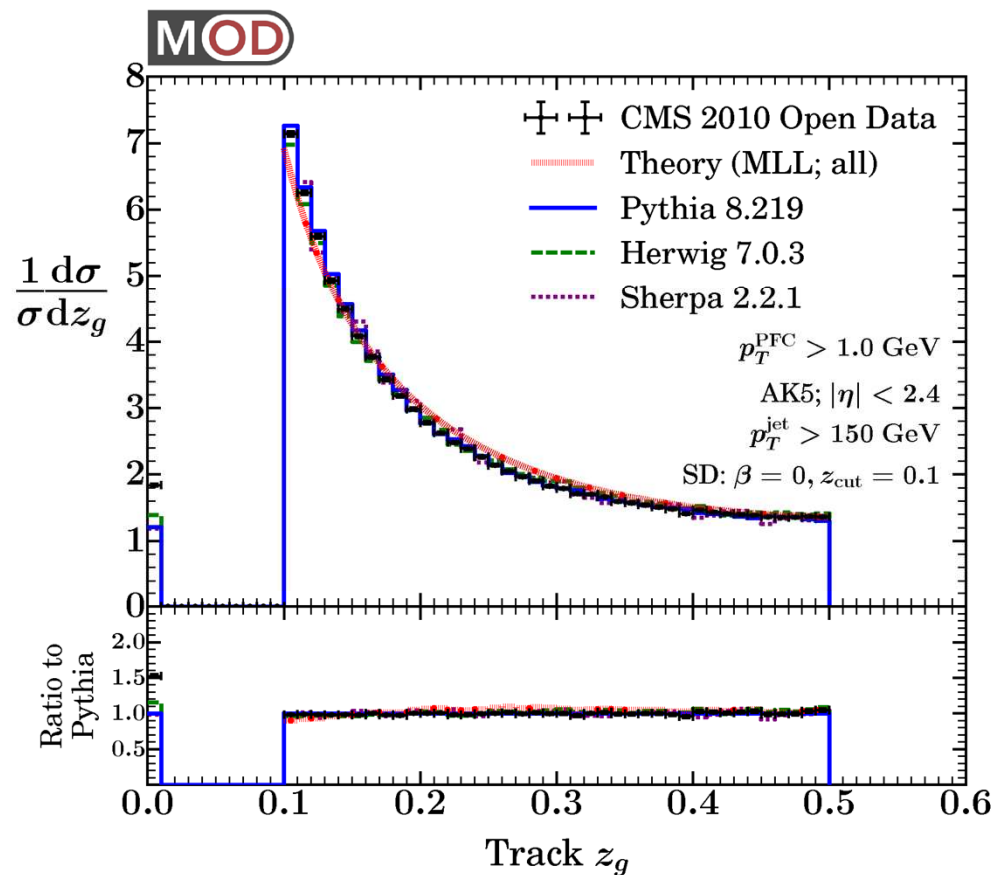
e-Print: [arXiv:1704.05066](https://arxiv.org/abs/1704.05066) [hep-ph] | [PDF](#)

(soon) **first ever published**
CMS Open Data result

using track (and PF) candidates
from CMS 2010 pp jet data



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



Mine data to test new (aspects of) theories

Jet Substructure Studies with CMS Open Data

Aashish Tripathy, Wei Xue, Andrew Larkoski, Simone Marzani, Jesse Thaler + some CMS support (S. Rappoccio)

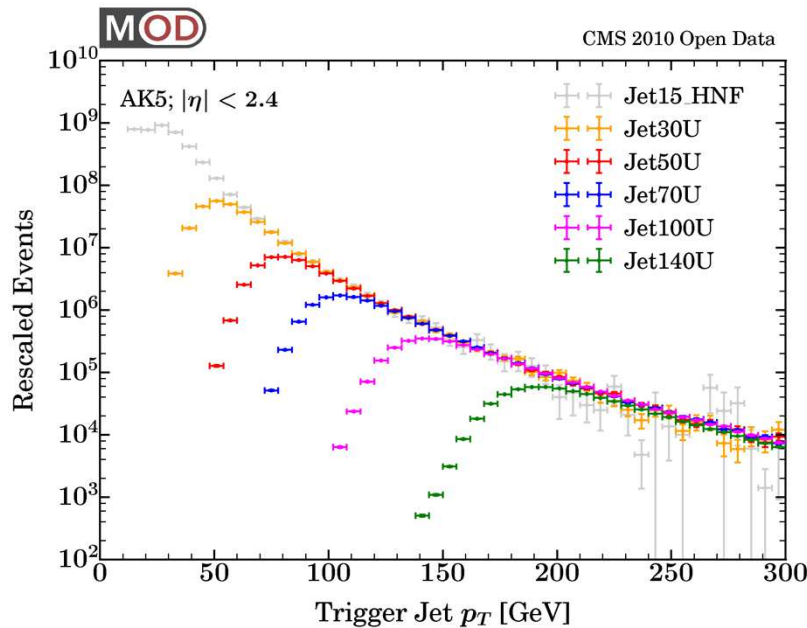
Apr 19, 2017 - 35 pages

MIT-CTP-4890

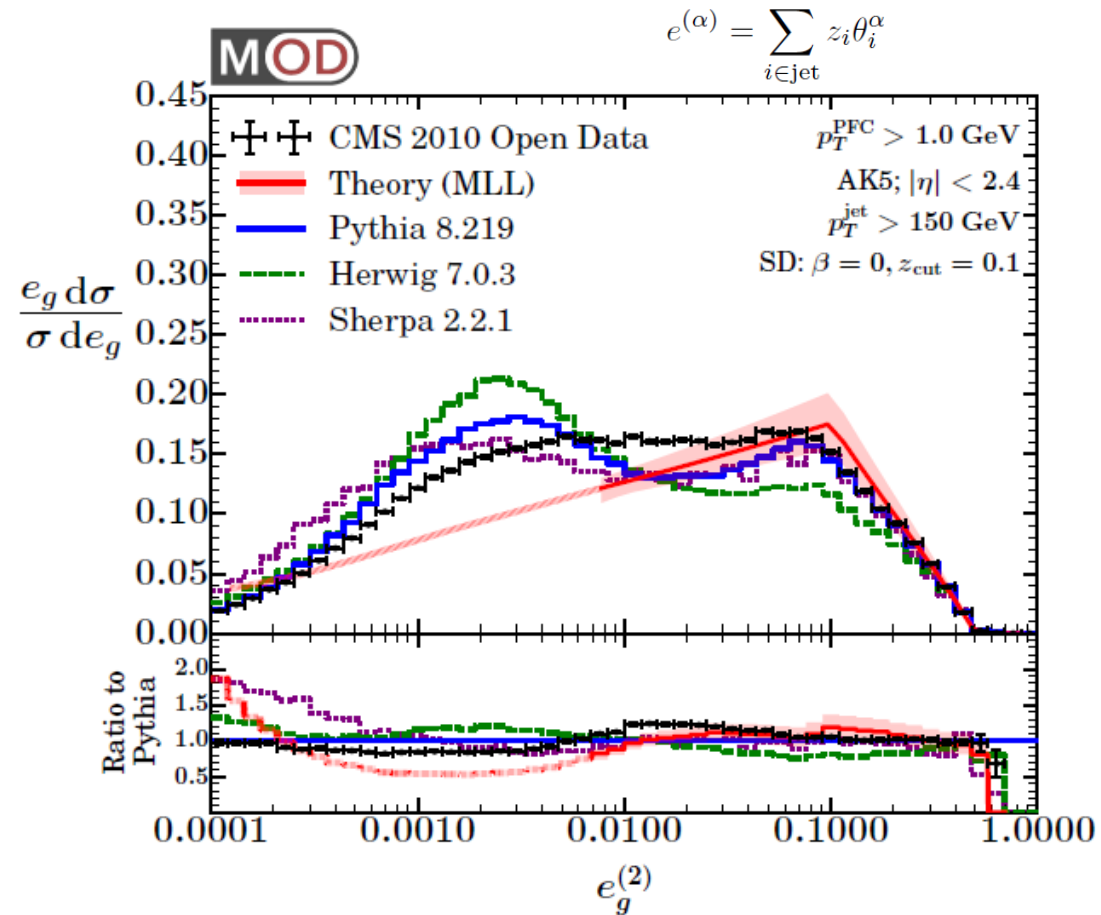
e-Print: [arXiv:1704.05842](https://arxiv.org/abs/1704.05842) [hep-ph] | [PDF](#)

combine jet data from different trigger thresholds -> challenge

using particle flow (PF) candidates from CMS 2010 pp jet data



expected kink structure observed, higher p_T jet measurements needed to disentangle perturbative from nonperturbative effects



Feedback to community

Jet Substructure Studies with CMS Open Data

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

Apr 19, 2017 - 35 pages

MIT-CTP-4890

e-Print: [arXiv:1704.05842](https://arxiv.org/abs/1704.05842) [hep-ph] | [PDF](#)

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

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

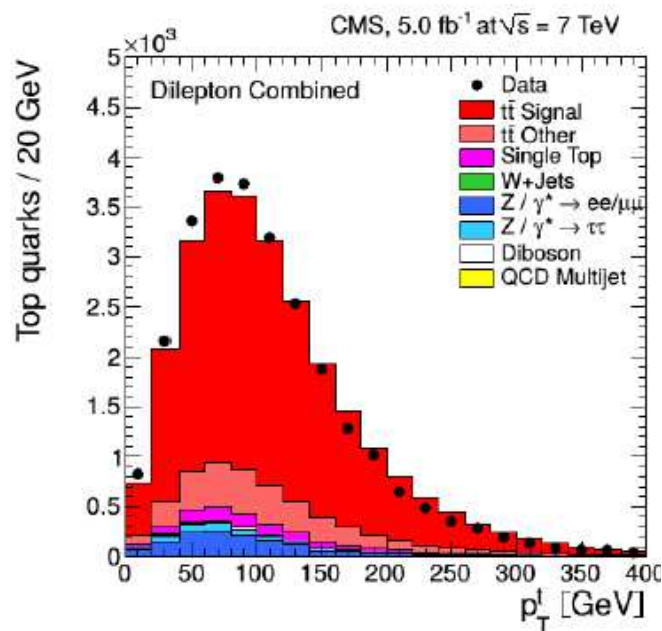
Conclusions: “We hope our experience motivates the LHC collaborations to further their investment in public data release and encourages the particle physics community to exploit the scientific potential of open datasets”

Open data benchmark analysis: top production

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

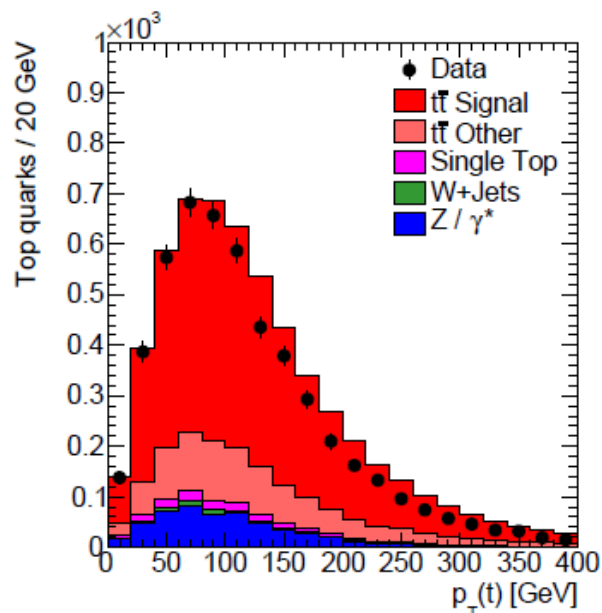
CMS Paper

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

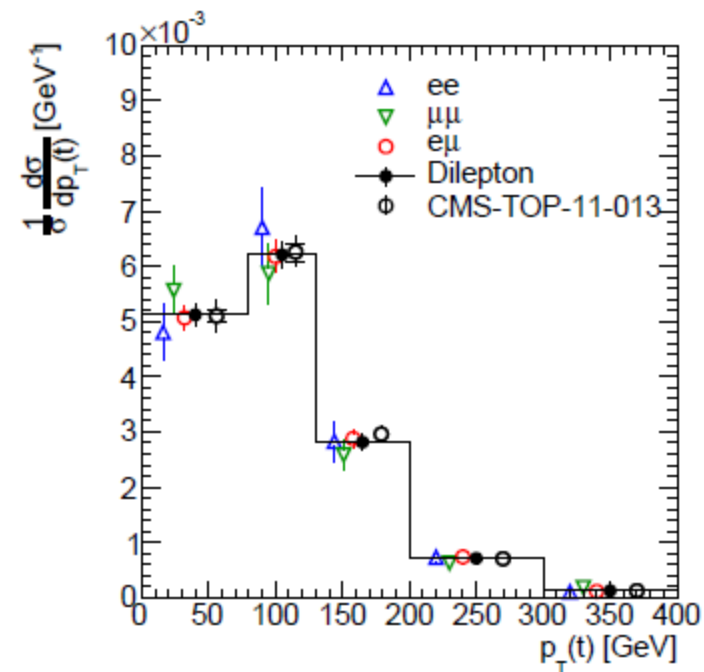


CMS Open Data

(O. Zenaiev, work in progress)

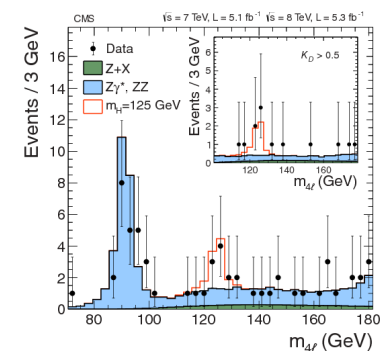
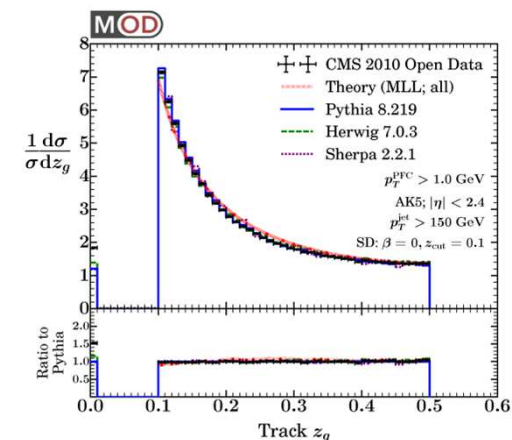
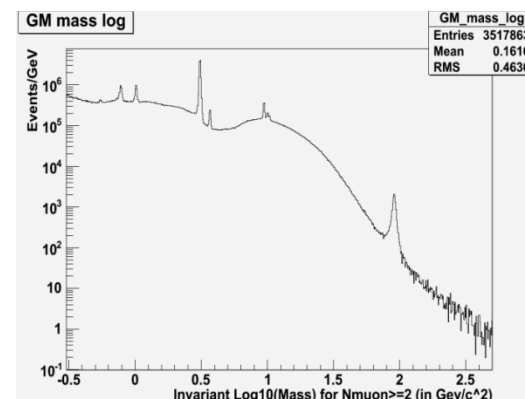


comparison of
norm. cross sections
(O. Zenaiev, work in progress)



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
- 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 (Yandex, see backup)
- upcoming 2012 data release: **on the way towards public reconstruction of the Higgs discovery**



Backup

Extended Vision (for discussion)

My personal extension of initial vision:

(for discussion, **not** a collaboration statement)

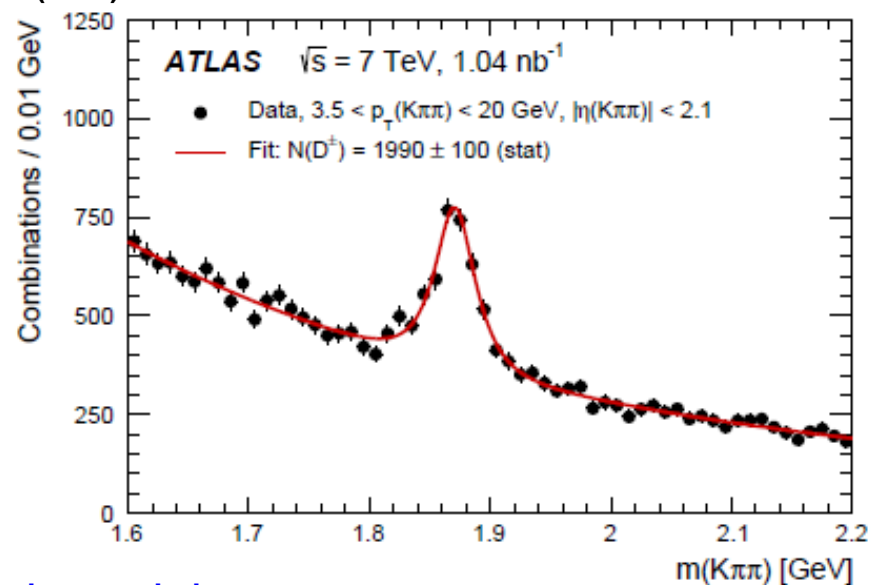
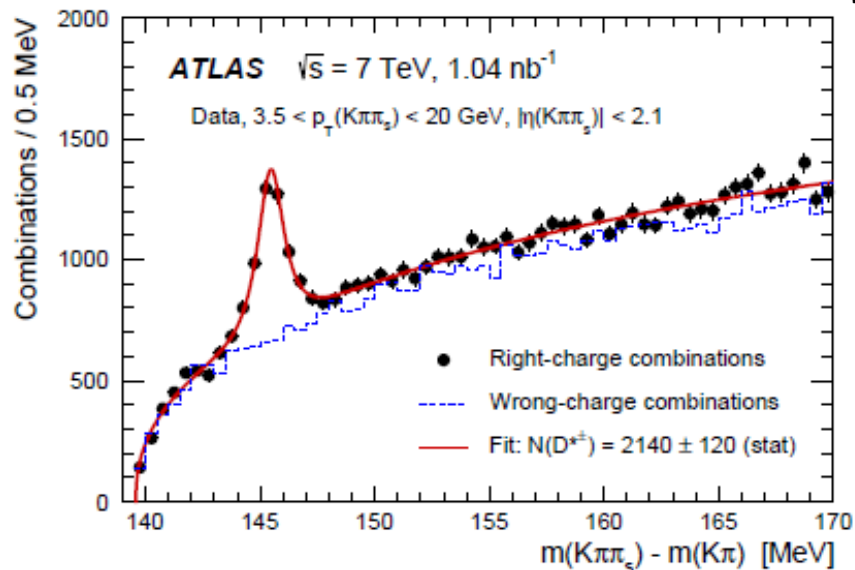
With ~1% of additional resources, aim to achieve

~10% additional scientific output (physics papers)

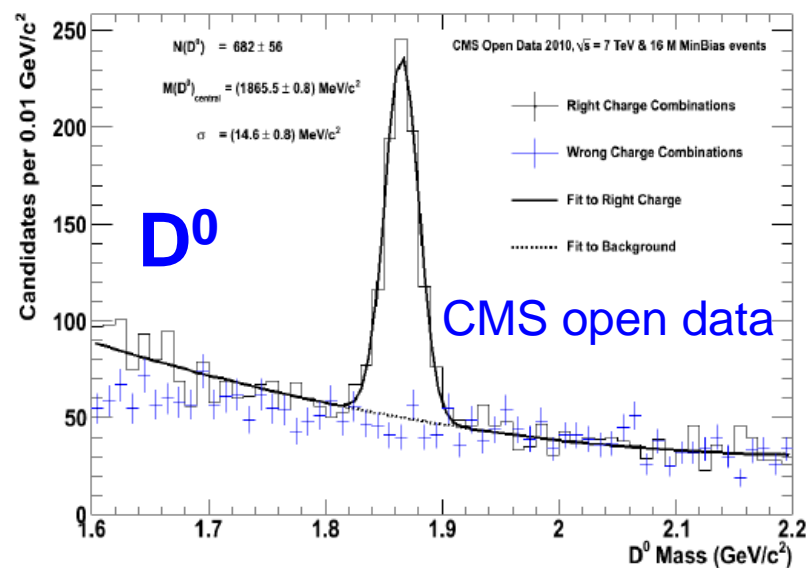
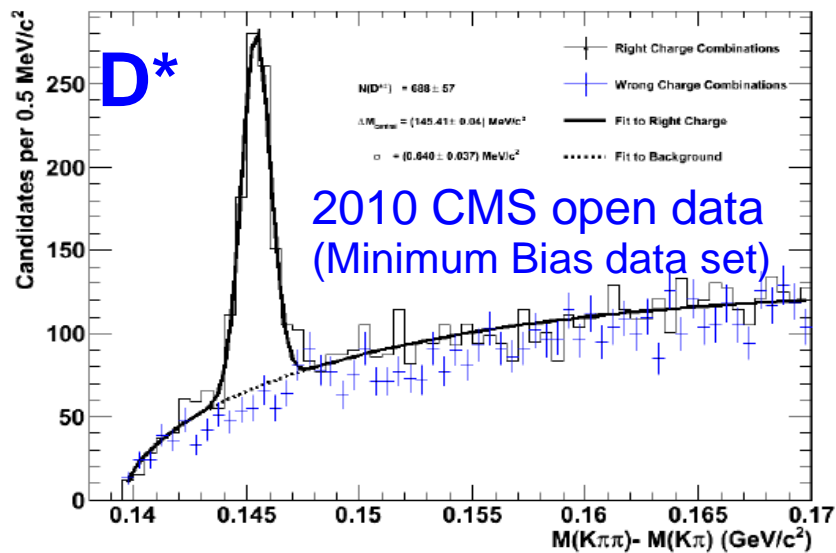
from both external and internal use of preserved data
over lifetime of experiment + 10-20 years

“Recreate” ATLAS result from CMS data: Low p_T D^* production (new for CMS)

Nucl. Phys. B907 (2016) 717



B. Sheeran, N. Stefaniuk, work in progress



Machine Learning with CMS Open Data: Yandex

Problem 1: Data Certification (CMS)

From seminar at DESY, 14.2.17

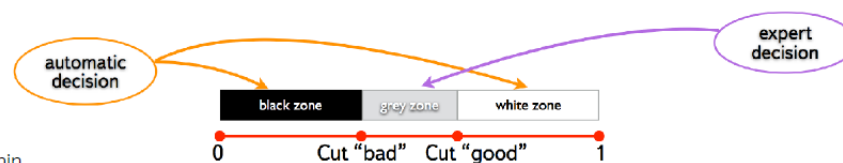
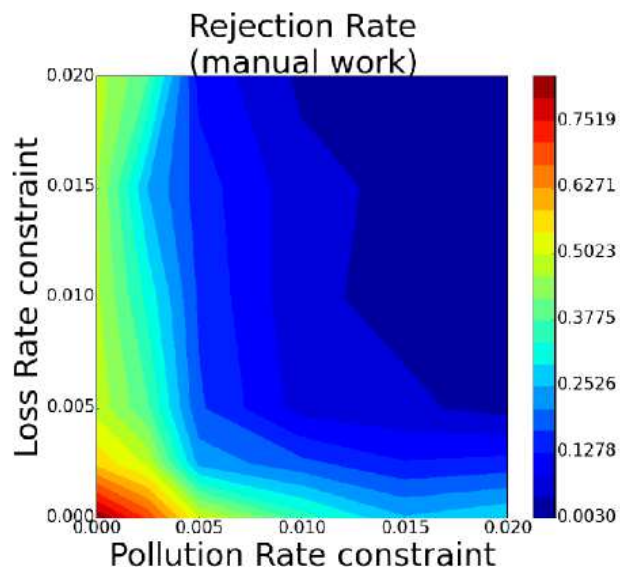
- Traditionally, quality of the data at CERN CMS experiment is determined manually which requires considerable amount of human efforts;
- ML can save some of those efforts;
- Data: CMS 2010B run open data;
- Aim: automated classification of Lumisections as “good” or “bad”;
- Features: particle flow jets, Calorimeter Jets, Photons, Muons;
- The dataset was flagged by experts (3 FTE).

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

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

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

Results



Andrey Ustuzhanin

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 Ustuzhanin

15.05.17

13

A. Geiser, LHCP workshop

18