



Deutsches Elektronen-Synchrotron DESY
CMS Group

nanoAOD(plus) validation from comparison to 2010 MuMonitor and MuOnia Open Data examples

Summer Student Report of
Fabian Stäger

Supervisors:
Achim Geiser
Josry Metwally

4 September 2019

Contents

1. Introduction	2
1.1. Data formats at the CMS experiment	2
1.2. nanoAOD(plus)	2
1.3. Validation of nanoAOD(plus)	3
2. Open Data examples	3
2.1. CMS Open Data	3
2.2. MuMonitor example	3
2.3. MuOnia example	4
3. Validation of nanoAOD(plus)	5
3.1. The procedure	5
3.2. Improvement of the nanoAOD(plus) ntuple	5
3.3. Muon and MuOnia datasets	8
3.4. Electron dataset	10
4. Conclusion	10
5. References	11
A. Appendix	12
A.1. Changes to nanoAOD(plus) ntuple by version	12
A.2. Validated nanoAOD(plus) variables	13
A.3. Additional histograms	15

1. Introduction

1.1. Data formats at the CMS experiment

Data from the CMS experiment comes in different data formats named RAW, RECO, AOD, miniAOD, and nanoAOD. RAW contains the raw data collected by the detector. RECO contains all the reconstructed physics object information. AOD (Analysis Object Data) is a subset of RECO, containing sufficient information for most physics analyses [1]. miniAOD was introduced for Run 2 in 2014. This format is targeted at being approximately 10% the size of the AOD format and to have sufficient information for about 80% of CMS analyses [2]. nanoAOD is a much smaller format targeted at being sufficient for at least 50% of analyses [3]. It is currently only available for Run 2 data from 2016 onwards. RECO, AOD, and miniAOD are EDM (event data model) files and are accessible through the CMSSW framework. nanoAOD is a flat ntuple format that is readable with plain ROOT. Table 1.1 gives a summary of the data formats of reconstructed objects and their availability for Run 1 and Run 2.

data format	file type	size [kB/event]	available for Run 2	available for Run 1
RECO	EDM	1400	no	yes
AOD	EDM	480	yes, partially	yes
miniAOD	EDM	40	yes	no
nanoAOD	flat ntuple	1	yes, 2016 onward	no

Table 1.1: Data formats of reconstructed objects at the CMS experiment. [1][2][3]

1.2. nanoAOD(plus)

Data in the nanoAOD format is much smaller and more easily accessible than data in the miniAOD, AOD, and RECO formats. Even complex analyses on big datasets can be computed in the order of minutes to hours on nanoAOD. The flat ntuple format makes it easier to read the data and removes dependence on CMSSW versions and a CMS environment. Especially on Run 1 data, where the older CMSSW versions require access through a virtual machine, being able to run analyses on a flat ntuple would be very convenient. Furthermore, the absence of nanoAOD and miniAOD for Run 1 data makes it unnecessarily challenging to adapt Run 2 analyses based on these formats for Run 1.

For these reasons, a project to create a nanoAOD-like format for datasets for which nanoAOD does not exist is currently ongoing. This format should have the same structure and content as the existing nanoAOD where possible, and some additional content where useful, potentially even replacing AOD for Run 1 in the long term [4]. In reference to this additional content, the project is referred to as nanoAOD(plus) in this report.

1.3. Validation of nanoAOD(plus)

nanoAOD for Run 2 data is created from miniAOD. As miniAOD does not exist for Run 1 data, all Run 2 algorithms for the nanoAOD content have to be rewritten to work directly on basic Run 1 AOD variables. As adapting these algorithms to work on AOD is not always straightforward, validation of nanoAOD(plus) variables is essential. Validation can be done either directly or indirectly. Direct validation is done by comparing technical distributions, i.e. directly comparing variable distributions from the existing nanoAOD and nanoAOD(plus). For this purpose nanoAOD(plus) ntuples have to be created for Run 2 datasets. Indirect validation is done by reproducing known physics distributions from Run 1 data using nanoAOD(plus). In my summer student project my task was to indirectly validate nanoAOD(plus) variables by reproducing histograms from Open Data examples on 2010 datasets using a nanoAOD(plus) ntuple.

2. Open Data examples

2.1. CMS Open Data

CMS Open Data is original data from the CMS experiment that has been made available to the public via the CERN Open Data Portal (<http://opendata.cern.ch>). It can be used by researchers outside the CMS collaboration to perform independent research. For scientific outreach and educational purposes, simplified data can be used by students or other interested persons [5]. CMS Open Data does not only contain datasets, but also analysis and validation examples. For the purpose of indirect validation of nanoAOD(plus) variables, I used the Open Data MuMonitor and MuOnia examples together with the 2010 Muon, MuOnia, and Electron datasets.

2.2. MuMonitor example

The MuMonitor example [6] is a validation example for the 2010 Muon and MuMonitor datasets, based on the dimuon invariant mass spectrum from global muons. Figure 2.1 shows the official CMS plot of the dimuon invariant mass spectrum and the corresponding plot from the MuMonitor example with the Muon dataset.

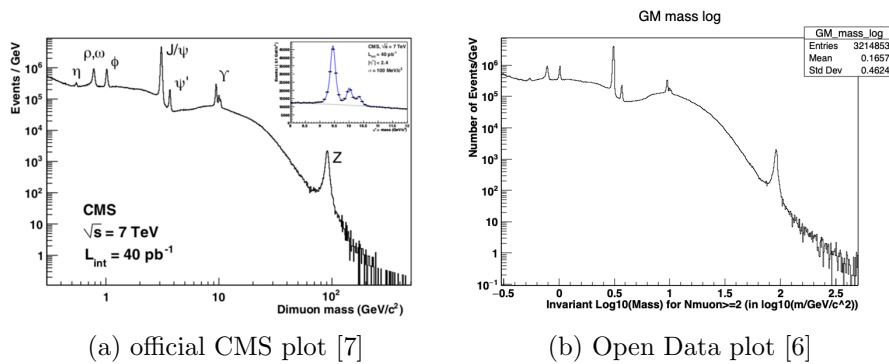


Figure 2.1: Dimuon invariant mass spectrum in CMS Run 1 data from 2010.

2.3. MuOnia example

The MuOnia example [8] is a not yet publicly available validation example for the 2010 Muon, MuOnia, and Electron datasets, based on (di)muon, (di)electron, and raw track, vertex, and jet spectra. Figures 2.2, 2.3, and 2.4 show some of the histograms in the MuOnia example together with the official CMS plots they were inspired by.

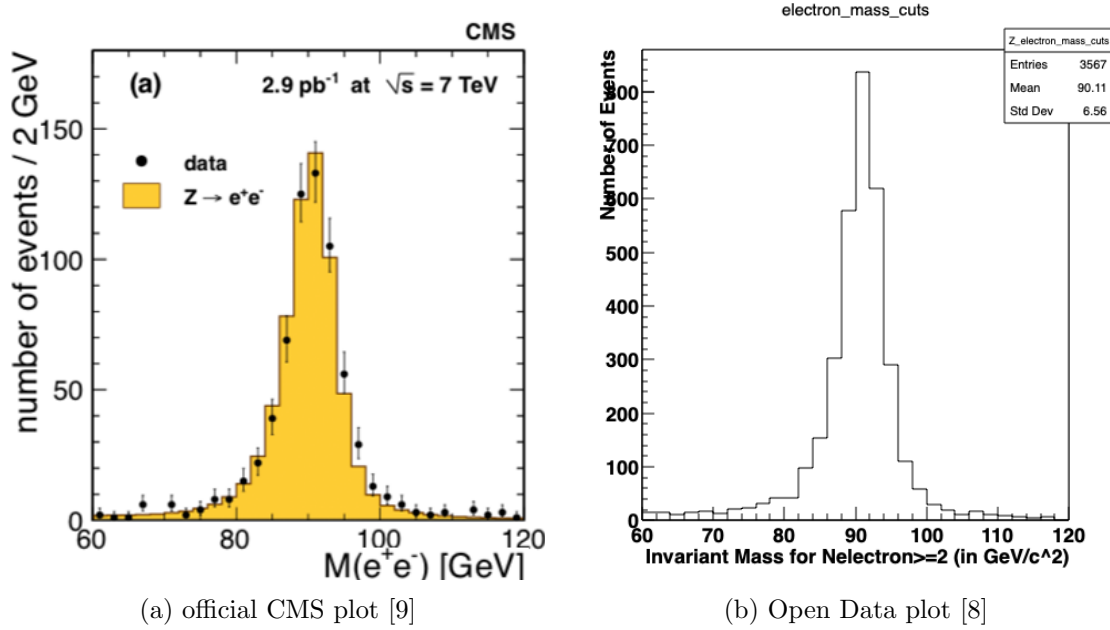


Figure 2.2: $Z \rightarrow e^+e^-$ resonance in CMS Run 1 data from 2010.

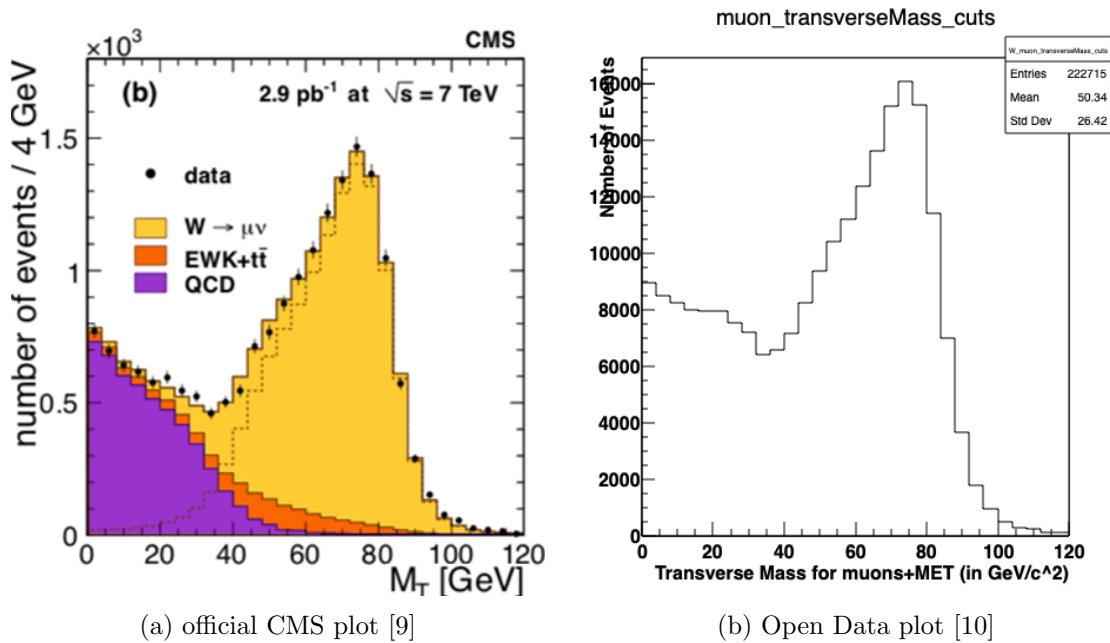


Figure 2.3: Transverse mass in $W \rightarrow \mu\nu_\mu$ decays in CMS Run 1 data from 2010.

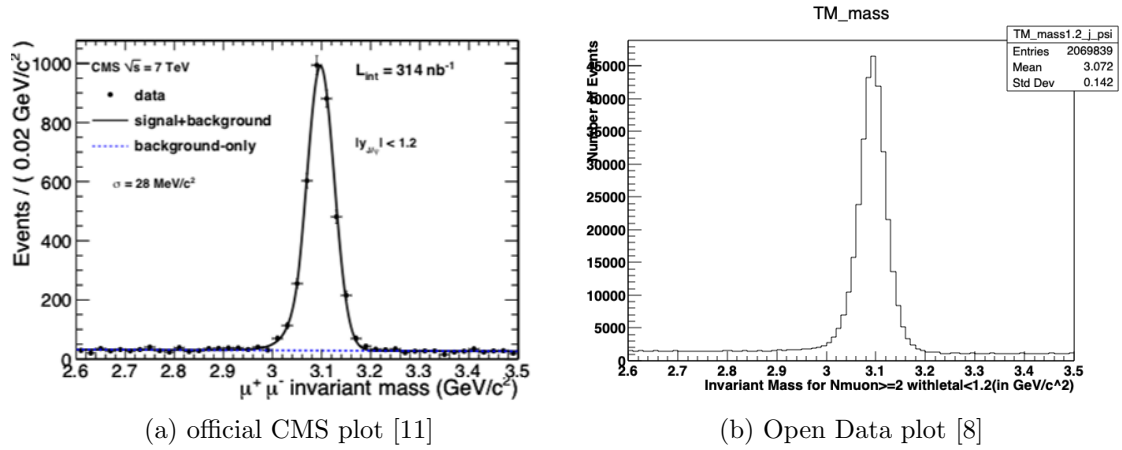


Figure 2.4: $J/\psi \rightarrow \mu^+\mu^-$ resonance with $|\text{rap}| < 1.2$ in CMS Run 1 data from 2010.

3. Validation of nanoAOD(plus)

3.1. The procedure

Indirect validation of nanoAOD(plus) variables from comparison to Open Data examples is done using the following procedure: First, a histogram from an Open Data example is picked and the cuts applied on AOD variables are looked up in the code. The same cuts are then applied on nanoAOD(plus) variables in a code that produces histograms from the nanoAOD(plus) ntuple. Finally, the reproduced histogram from the nanoAOD(plus) ntuple and the original histogram from the Open Data example are compared. If the two histograms are not identical, the source of the difference has to be identified and the procedure is repeated with an adapted nanoAOD(plus) ntuple.

3.2. Improvement of the nanoAOD(plus) ntuple

During the course of my summer student project, validation on Open Data examples has contributed to the improvement of the nanoAOD(plus) ntuple in several iterations. In this section, some of the changes on the nanoAOD(plus) ntuple relevant to my project are explained on the example of the $\Upsilon \rightarrow \mu^+\mu^-$ resonance from the MuOnia example with the Muon 2010 dataset. Figure 3.1 shows the $\Upsilon \rightarrow \mu^+\mu^-$ histogram from the Open Data example with the histogram reproduced using the nanoAOD(plus) ntuple in the version from November 2018 (version name .foroptmu) plotted on top of it. The selection requires two muons of opposite charge with an upper limit on the pseudorapidity ($|\eta| < 1.0$).

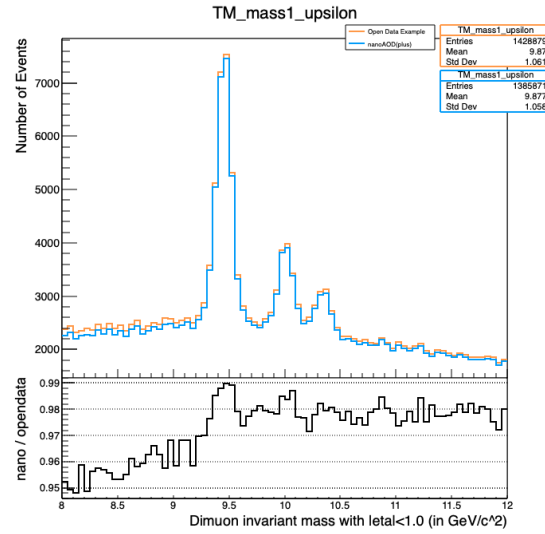


Figure 3.1: $\Upsilon \rightarrow \mu^+ \mu^-$ resonance. Open Data example histogram (orange) [8] and reproduced histogram using nanoAOD(plus) Muon ntuple version .foroptmu (blue).

The lower number of entries in the reproduced histogram ($\sim 3\%$) is due to preselection cuts on low- p_T muons in the production of the nanoAOD(plus) ntuple. These preselection cuts are removed in the version from 7 August 2019 (version name .zerobias). Figure 3.2 shows the Open Data example histogram and the histogram reproduced from this second version of the nanoAOD(plus) ntuple.

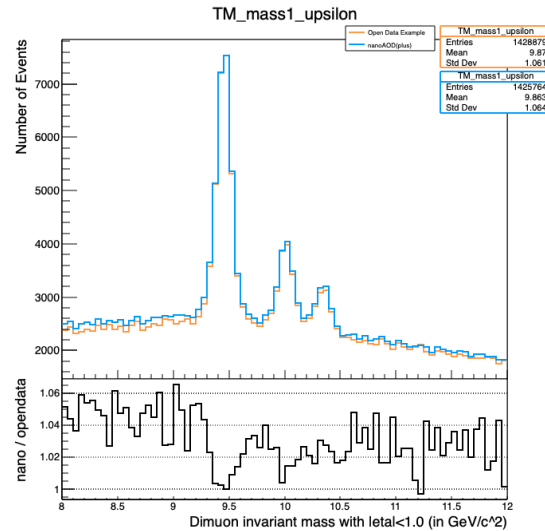


Figure 3.2: $\Upsilon \rightarrow \mu^+ \mu^-$ resonance. Open Data example histogram (orange) [8] and reproduced histogram using nanoAOD(plus) Muon ntuple version .zerobias (blue).

With the preselection cuts removed, the number of entries in the two histograms is now almost identical ($\sim 0.2\%$ difference), but the ratio plot shows a shift away from the $\Upsilon(nS)$ peaks in the reproduced histogram. A closer investigation revealed a bug in the sorting of the `Muon_phi` variable, resulting in some muons being assigned a wrong value for the azimuthal angle ϕ . This leads to miscalculations of the invariant mass of the dimuon system. This bug is fixed in the nanoAOD(plus) version from 19 August 2019 (version name `.zerobias2`). Figure 3.3 shows the Open Data example histogram and the histogram reproduced from this third version of the nanoAOD(plus) ntuple.

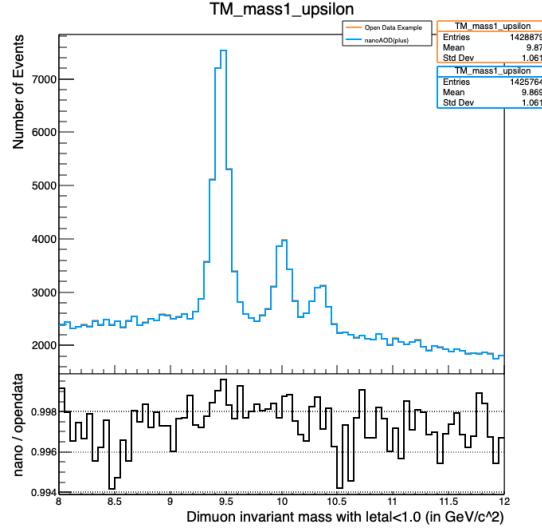


Figure 3.3: $\Upsilon \rightarrow \mu^+ \mu^-$ resonance. Open Data example histogram (orange) [8] and reproduced histogram using nanoAOD(plus) Muon ntuple version `.zerobias2` (blue).

With the bug removed, the shape of the two histograms now looks identical. However, there is still a small difference in the number of entries of about 0.2% . This is solved by extending the length of the vectors in the nanoAOD(plus) ntuple from 32 to 128 (i.e. allowing to store up to 128 muons per event) in the nanoAOD(plus) version from 23 August 2019 (version name `.zerobias3`). With this change the histogram is now exactly reproduced from the nanoAOD(plus) ntuple, except for a small fluctuation in the $\Upsilon(1S)$ peak which is likely due to a small number of events being placed in the wrong bin due to a rounding error. Figure 3.4 shows the Open Data example histogram and the histogram reproduced from this fourth version of the nanoAOD(plus) ntuple. The complete list of changes relevant to this project in each version of the nanoAOD(plus) ntuple are listed in table A.1 in the appendix.

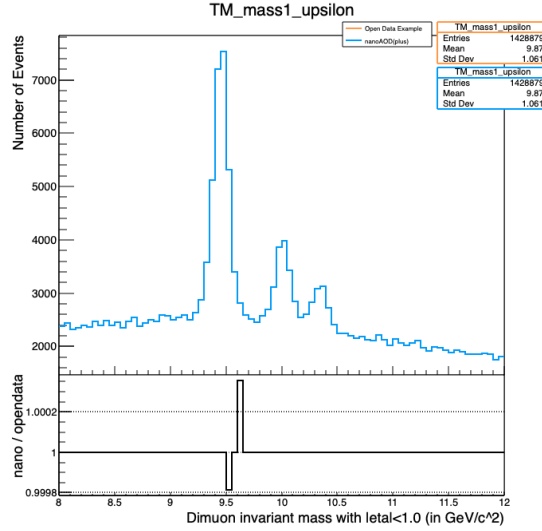


Figure 3.4: $\Upsilon \rightarrow \mu^+ \mu^-$ resonance. Open Data example histogram (orange) [8] and reproduced histogram using nanoAOD(plus) Muon ntuple version .zerobias3 (blue).

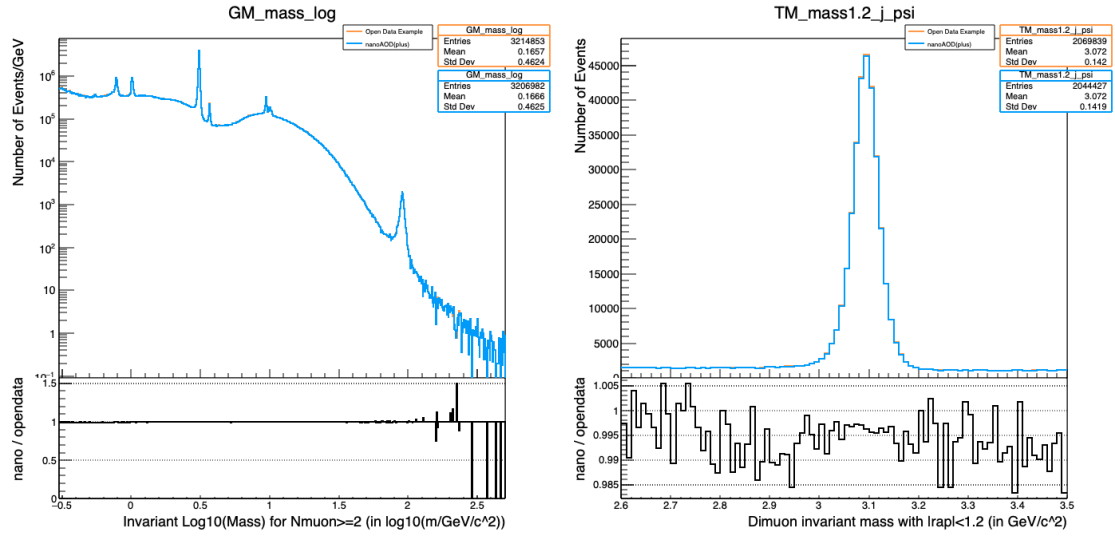
3.3. Muon and MuOnia datasets

While many of the muon distributions from the Open Data MuOnia and MuMonitor examples have been reproduced exactly with the latest version of the nanoAOD(plus) ntuple, in some histograms there are still small differences remaining. Figure 3.5 shows four of the muon histograms that are not yet exactly reproduced.

In the case of the dimuon invariant mass spectrum from the Muon dataset (figure 3.5a), the reproduced histogram from the nanoAOD(plus) ntuple has $\sim 0.2\%$ fewer entries. This deviation could possibly be due to a cut on the number of valid hits in the tracker that for global muons can not yet be exactly reproduced on nanoAOD(plus) variables.

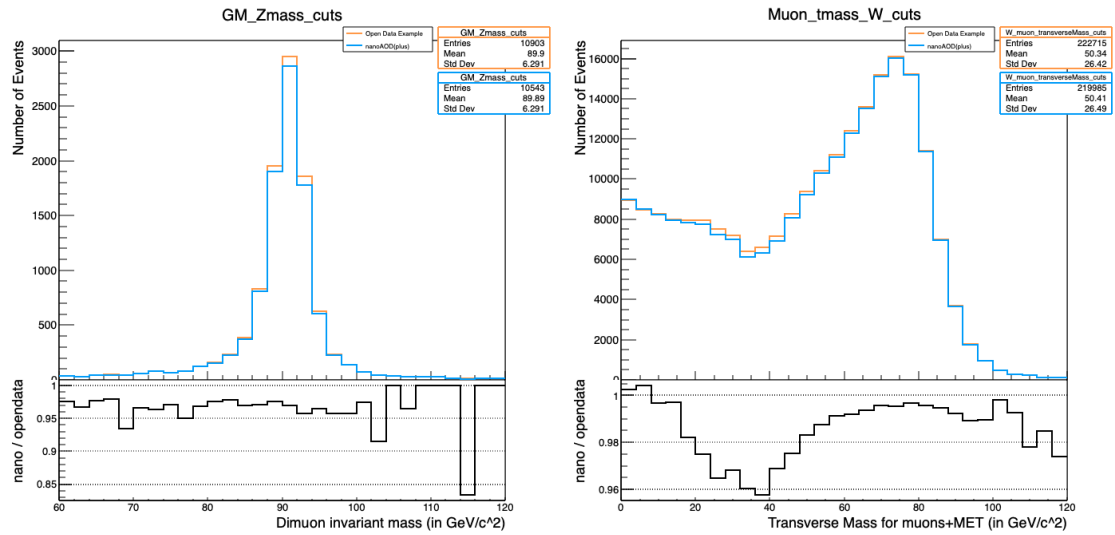
For the $J/\psi \rightarrow \mu^+ \mu^-$ resonance from the MuOnia dataset (figure 3.5b), the $\sim 1\%$ deviation has to be further investigated. A possible reason is a bug in the code that produces the histogram in the Open Data example. The code contains an upper limit on the vertex distance of the two muons. However, it does not take the absolute value of the distance, which means that the cut depends on the sign of the distance. It was attempted to reproduce this bug with the nanoAOD(plus) ntuple, but it can only be reproduced exactly if the ordering of the muons is identical in both data formats, which might not be true in all cases.

The reproduced histogram of the $Z \rightarrow \mu^+ \mu^-$ resonance (figure 3.5c) from the Muon dataset has $\sim 3\%$ fewer entries than the corresponding histogram from the Open Data example. This deviation could again come from a cut on the number of valid tracker hits. Another possible contribution comes from a cut on dxy. The Open Data example calculates the value of dxy using the global muon fit, while the nanoAOD(plus) variable cut on is calculated using the tracker muon fit.



(a) dimuon invariant mass spectrum [6],
MuMonitor example, Muon dataset

(b) $J/\psi \rightarrow \mu^+ \mu^-$ resonance [8],
MuOnia example, MuOnia dataset



(c) $Z \rightarrow \mu^+ \mu^-$ resonance [8],
MuOnia example, Muon dataset

(d) $W \rightarrow \mu \nu_\mu$ transverse mass [10],
MuOnia example, Muon dataset

Figure 3.5: Histograms from the Muon and MuOnia 2010 datasets. Open Data example histograms (orange) and reproduced histograms using nanoAOD(plus) ntuple (blue).

The reasons for the deviation of $\sim 1\%$ in the $W \rightarrow \mu \nu_\mu$ transverse mass histogram (figure 3.5d) could be the same as in the $Z \rightarrow \mu^+ \mu^-$ case. In addition, the code that produces the histogram in the MuOnia example contains a bug that again can not be exactly reproduced with the nanoAOD(plus) ntuple.

3.4. Electron dataset

All the electron histograms in the MuOnia Open Data example have been reproduced exactly from nanoAOD(plus), meaning that all tested electron variables have been validated. Figure 3.6 shows the Open Data example histogram and the reproduced histogram from nanoAOD(plus) of the $Z \rightarrow e^+e^-$ resonance. The small fluctuation in two bins around the peak is likely due to a rounding error.

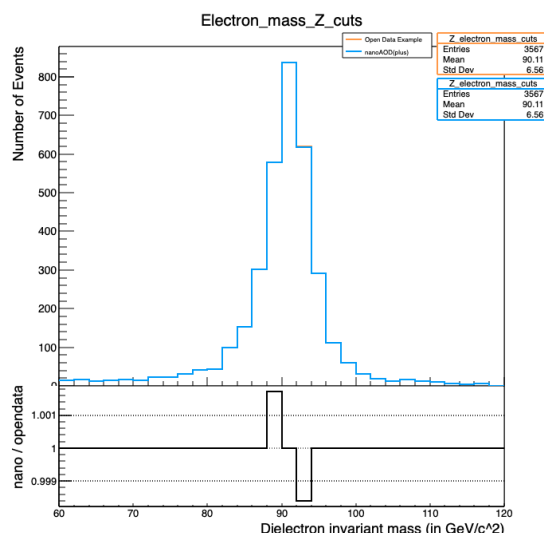


Figure 3.6: $Z \rightarrow e^+e^-$ resonance from the Electron 2010 dataset. Open Data example histogram (orange) [8] and reproduced histogram using nanoAOD(plus) Electron ntuple (blue).

4. Conclusion

All electron histograms in the MuOnia Open Data example have been reproduced exactly from nanoAOD(plus), meaning that all tested electron variables have been validated for 2010 data. The muon histograms in the MuOnia and MuMonitor Open Data examples have been reproduced either exactly or very close, so all tested muon variables have been at least partially validated for 2010 data. However, another iteration is needed to eliminate the remaining differences. A list of all validated variables with references to the corresponding histograms is given in table A.2.

Additionally, some minor bugs have been found in the code that produces the MuOnia Open Data example. Fixes to these bugs have been proposed and will be implemented before the MuOnia example goes public.

In a next step, it should be attempted to reproduce some of the invariant mass and transverse mass plots as close as possible using only official nanoAOD variables. This would give an idea about which of the introduced *plus* variables add any crucial information. This would bring us a step closer to the goal of being able to analyze Run 1 and Run 2 data using the same nanoAOD algorithms.

5. References

- [1] *CMS WorkBookDataFormats*. URL: <https://twiki.cern.ch/twiki/bin/view/CMSPublic/WorkBookDataFormats>.
- [2] G. Petrucciani, A. Rizzi, and C. Vuosalo. *Mini-AOD: A New Analysis Data Format for CMS*. arXiv:1702.04685 [physics.ins-det], 2015.
- [3] A. Rizzi and G. Petrucciani. *A further reduction in CMS event data for analysis: the NANO AOD format*. CHEP 2018 Conference, Sofia, Bulgaria. URL: <https://indico.cern.ch/event/587955/contributions/2937531/>.
- [4] A. Geiser. *nanoAOD-like*. URL: <https://twiki.cern.ch/twiki/bin/view/Main/NanoAODlike>.
- [5] *CMS Open Data*. URL: <http://opendata.cern.ch/docs/about-cms>.
- [6] A. Geiser et al. *Validation code for 2010 Mu and MuMonitor datasets, based on dimuon mass spectrum*. DOI:10.7483/OPENDATA.CMS.A8CP.HBJQ.
- [7] CMS Collaboration. *Performance of CMS muon reconstruction in pp collision events at $\sqrt{s} = 7$ TeV*. CMS-MUO-10-004. arXiv:1206.4071 [physics.ins-det], 2012.
- [8] L. Thulasidharan. *J/ψ , Υ and Z studies with CERN Open Data*. DESY Summer Student Report, 2018.
- [9] CMS Collaboration. *Measurements of Inclusive W and Z Cross Sections in pp Collisions at $\sqrt{s} = 7$ TeV*. CMS-EWK-10-002. arXiv:1012.2466 [hep-ex], 2010.
- [10] *private communications*.
- [11] CMS Collaboration. *Prompt and non-prompt J/ψ production in pp collisions at $\sqrt{s} = 7$ TeV*. CMS-BPH-10-002. arXiv:1011.4193 [hep-ex], 2010.

A. Appendix

A.1. Changes to nanoAOD(plus) ntuple by version

change	corresponding AOD variable
.zerobias (7 August 2019)	
remove muon preselection cuts	
new variable Muon_gChi2	muon->globalTrack()->normalizedChi2()
new variable Muon_gpt	muon->globalTrack()->pt()
new variable Muon_geta	muon->globalTrack()->eta()
new variable Muon_gphi	muon->globalTrack()->phi()
new variable Muon_nValidMu	muon->globalTrack()->hitPattern(). numberOfValidMuonHits()
new variable Muon_isGood	muon::TMOneStationTight
new variable Muon_isGoodLast	muon::TMLastStationTight
new variable MET_pt	pfmet->begin()->pt()
new variable MET_phi	pfmet->begin()->phi()
new variable CaloMET_pt	calomet->begin()->pt()
.zerobias2 (19 August)	
bug fix to Muon_phi	
new variable Muon_gnValid	muon->globalTrack()->hitPattern(). numberOfValidTrackerHits()
new variable Muon_gnPix	muon->globalTrack()->hitPattern(). numberOfValidPixelHits()
.zerobias3 (23 August)	
increase length of vectors to 128	
remove 5 GeV electron p_T cut	
new variable Muon_Chi2	muon->innerTrack()->normalizedChi2()
new variable Muon_isGoodAng	muon::TMLastStationAngTight
new variable Muon_isArbitrated	muon::TrackerMuonArbitrated
.json2 (electrons only, 30 August)	
new variable Electron_deltaEtaSCTR	electron-> deltaEtaSuperClusterTrackAtVtx()
new variable Electron_deltaPhiSCTR	electron-> deltaPhiSuperClusterTrackAtVtx()

Table A.1: List of changes to nanoAOD(plus) ntuple by version.

A.2. Validated nanoAOD(plus) variables

nanoAOD variable	corresponding AOD variable	validated	reference
nElectron	size()	yes	figure A.1d
Electron_charge	electron->charge()	yes	figure 3.6
Electron_deltaEtaSCTR	electron->deltaEtaSuperClusterTrackAtVtx()	yes	figure A.2a
Electron_dr03EcalRecHitSumEt	electron->dr03EcalRecHitSumEt()	yes	figure A.2d
Electron_dr03TKSumPt	electron->dr03TkSumPt()	yes	figure A.2f
Electron_eta	electron->eta()	yes	figure A.1b
Electron_phi	electron->phi()	yes	figure A.1c
Electron_pt	electron->pt()	yes	figure A.1a
Electron_hoe	electron->hcalOverEcal()	yes	figure 3.6
Electron_lostHits	electron->gsfTrack()->trackerExpectedHitsInner().numberOfHits()	yes	figure A.1e
Electron_mass	electron->mass()	no	
Electron_sieie	electron->sigmaIetaIeta()	yes	figure A.2c
Electron_convDist	electron->convDist()	yes	figure A.1h
Electron_convDcot	electron->convDcot	yes	figure A.1g
Electron_deltaPhiSCTR	electron->deltaPhiSuperClusterTrackAtVtx()	yes	figure A.2b
Electron_dr03HcalTowerSumEt	electron->dr03HcalTowerSumEt()	yes	figure A.2e
Electron_isEB	electron->isEB()	yes	figure A.2b
Electron_isEE	electron->isEE()	yes	figure A.2d
Electron_SCeta	electron->superCluster->eta()	yes	figure A.1f
nMuon	size()	yes	figure A.3a
Muon_charge	muon->charge()	yes	figure 3.4
Muon_dxy	muon->track->dxy()	partially	figure A.6d
Muon_eta	muon->eta()	yes	figure A.4b
Muon_phi	muon->phi()	yes	figure A.4c
Muon_pt	muon->pt()	yes	figure A.4a
Muon_isGlobal	muon->isGlobalMuon()	partially	figure 3.5a
Muon_isTracker	muon->isTrackerMuon()	partially	figure 3.5b
Muon_mass	0.105658 GeV	yes	figure 3.4
Muon_pfRelIso03_all	(muon->isolationR03().sumPt() + muon->isolationR03().hadEt() + muon->isolationR03().emEt()) / muon->pt()	partially	figure 3.5c
Muon_x	muon->vx()	yes	figure A.6a
Muon_y	muon->vy()	yes	figure A.6b
Muon_z	muon->vz()	yes	figure A.6c

Muon_nValid	muon->innerTrack()-> hitPattern(). numberOfValidTrackerHits()	yes	figure A.6g
Muon_nPix	muon->innerTrack()-> hitPattern(). numberOfValidPixelHits()	yes	figure A.6h
Muon_trkIdx	Muon Track Index	yes	figure A.6g
Muon_gChi2	muon->globalTrack()-> normalizedChi2()	partially	figure A.3c
Muon_geta	muon->globalTrack()->eta()	partially	figure A.4e
Muon_gphi	muon->globalTrack()->phi()	partially	figure A.4f
Muon_gpt	muon->globalTrack()->pt()	partially	figure A.4d
Muon_nValidMu	muon->globalTrack()-> hitPattern(). numberOfValidMuonHits()	partially	figure A.6e
Muon_isGood	TMOneStationTight	partially	figure A.6i
Muon_isGoodLast	TMLastStationTight	partially	figure 3.5c
Muon_gnValid	muon->globalTrack()-> hitPattern(). numberOfValidTrackerHits()	partially	figure A.6e
Muon_gnPix	muon->globalTrack()-> hitPattern(). numberOfValidPixelHits()	yes	figure A.6f
Muon_Chi2	muon->innerTrack()-> normalizedChi2()	yes	figure A.3b
Muon_isGoodAng	TMLastStationAngTight	partially	figure 3.5b
Muon_isArbitrated	TrackerMuonArbitrated	partially	figure 3.5b
MET_pt	pfmet->begin()->pt()	yes	figure A.5a
MET_phi	pfmet->begin()->phi()	partially	figure 3.5d
CaloMET_pt	calomet->begin()->pt()	yes	figure A.5b

Table A.2: List of validated or partially validated nanoAOD(plus) variables with reference plots and corresponding AOD variables.

A.3. Additional histograms

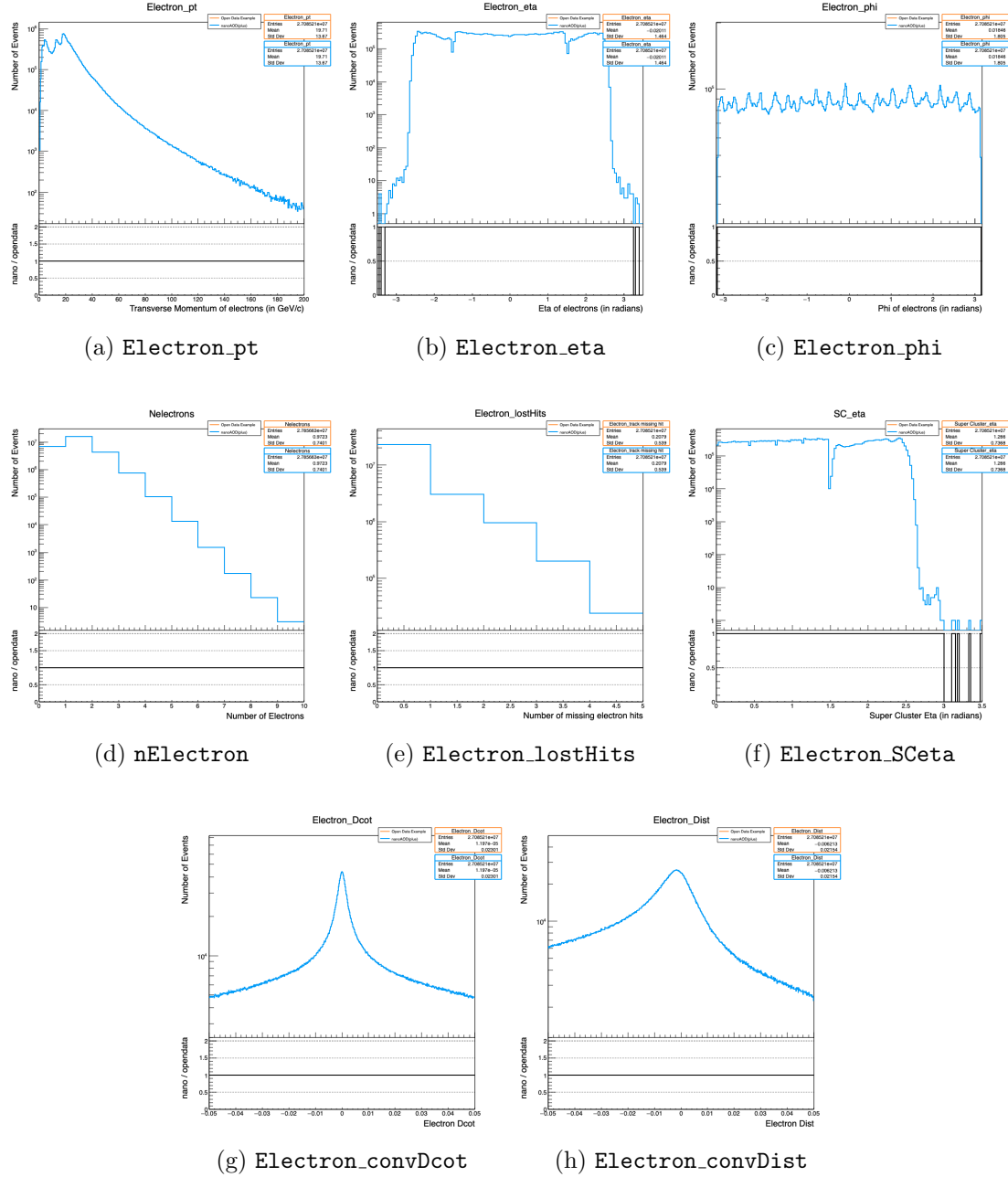


Figure A.1: Electron variables. Open Data example histograms (orange) [10] and reproduced histograms using nanoAOD(plus) Electron ntuple (blue).

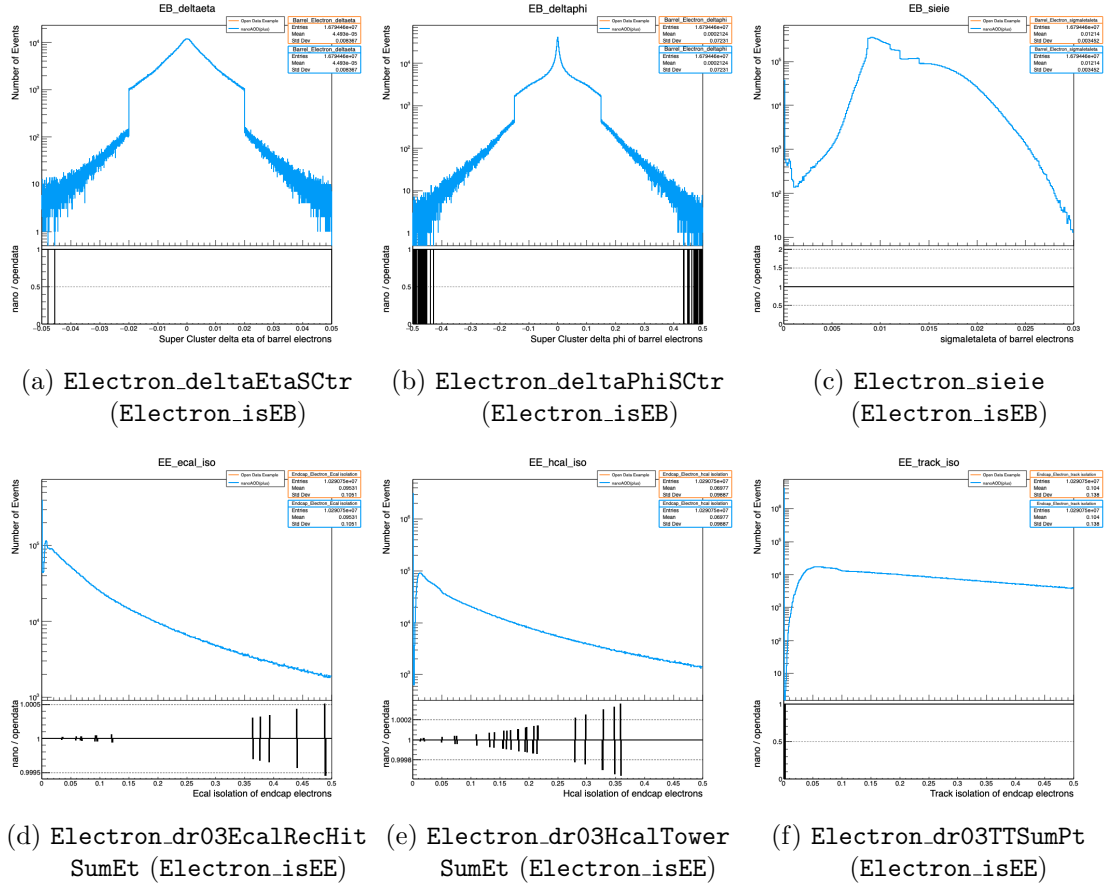


Figure A.2: Electron variables for barrel (EB) and endcap (EE). Open Data example histograms (orange) [10] and reproduced histograms using nanoAOD(plus) Electron ntuple (blue).

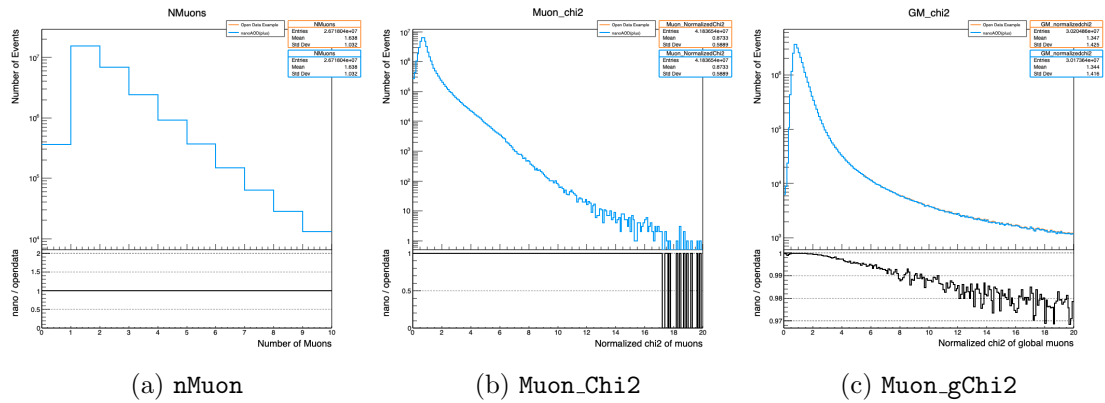


Figure A.3: Muon variables. Open Data example histograms (orange) [10] and reproduced histograms using nanoAOD(plus) Muon ntuple (blue).

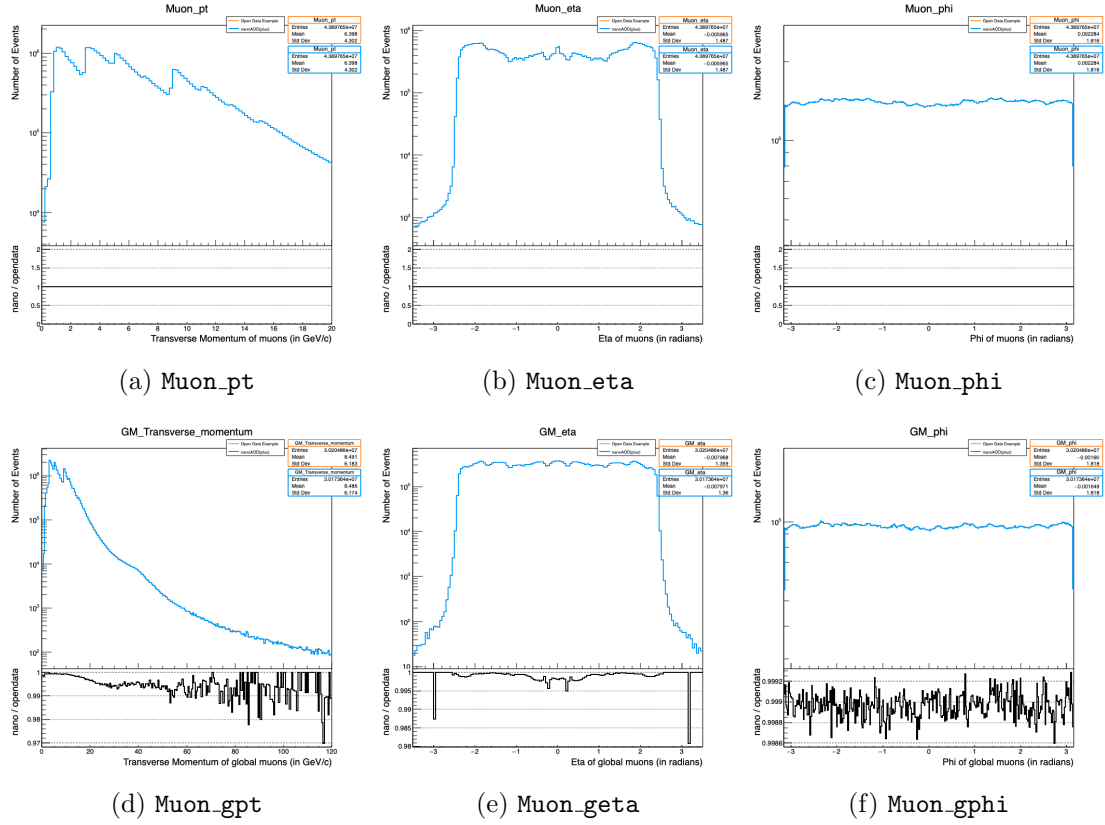


Figure A.4: Muon kinematic variables. Open Data example histograms (orange) [10] and reproduced histograms using nanoAOD(plus) Muon ntuple (blue).

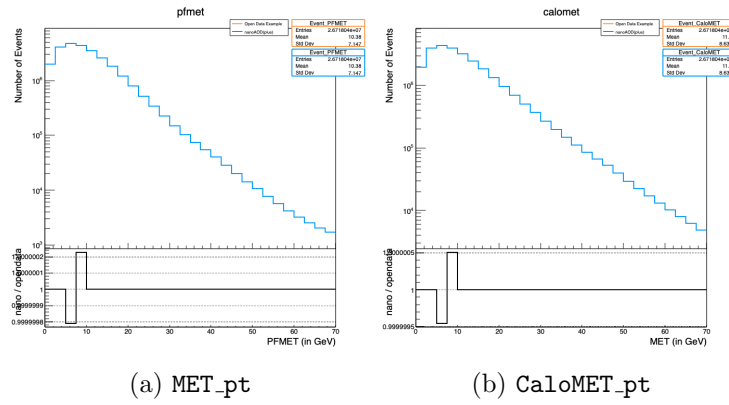


Figure A.5: Missing transverse energy variables. Open Data example histograms (orange) [10] and reproduced histograms using nanoAOD(plus) Muon ntuple (blue).

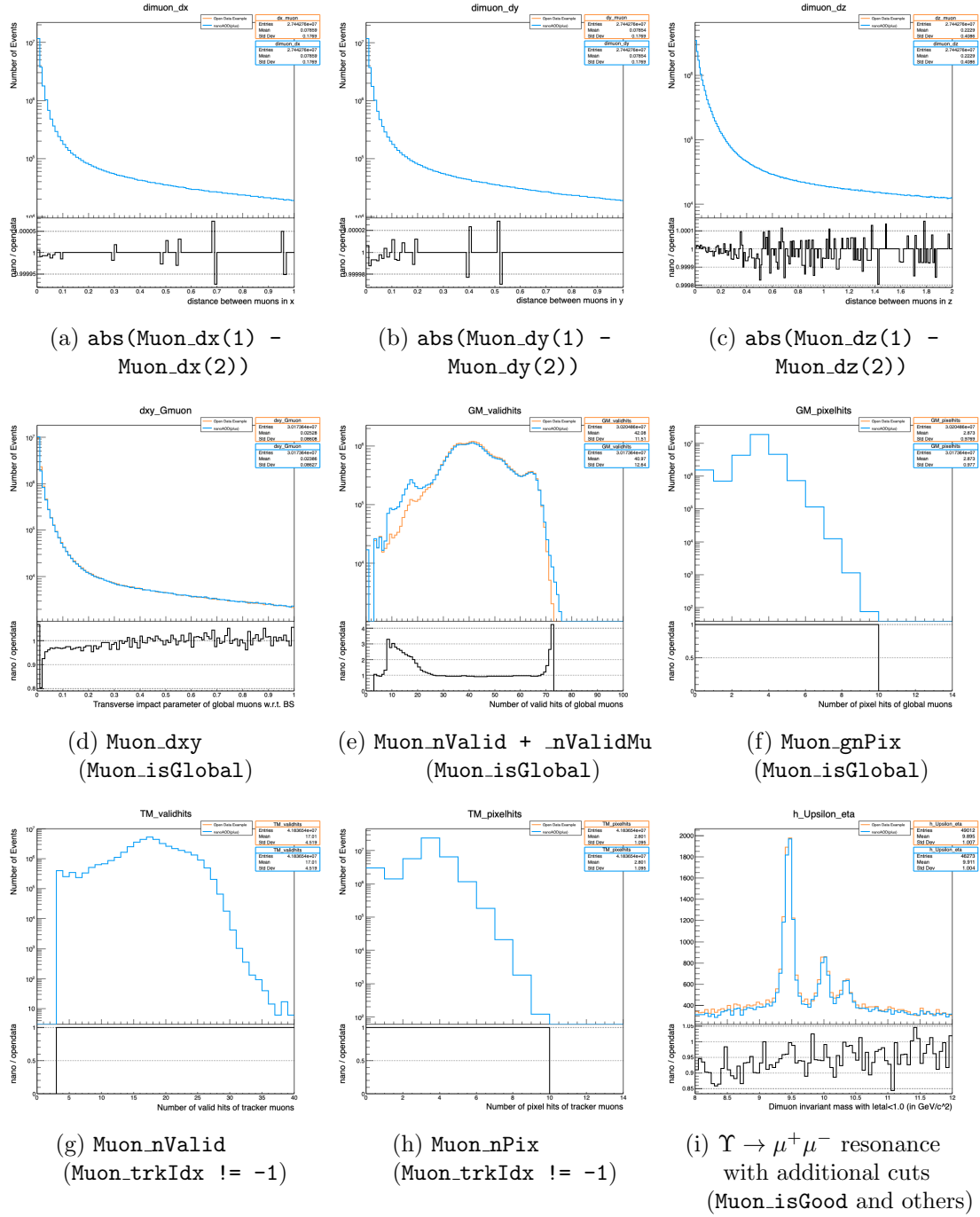


Figure A.6: Muon variables. Open Data example histograms (orange) [10][8] and reproduced histograms using nanoAOD(plus) Muon ntuple (blue).