# Start your day with a cup of CMS open data

# How do I take my cup of CMS Open Data?

Rikab Gambhir

Available at a computer near you!

Photo by Kelly Sikkema on Unsplash

### I like my CMS Open Data like I like my coffee ....

# Start your day with a cup of CMS open data

#### Available at a computer near you!

Photo by Kelly Sikkema on Unsplasi

# I like my CMS Open Data like I like my coffee ...

- Very easily accessible anywhere I am
- Takes only a few seconds to minutes to set up
- Highly preprocessed and prepackaged
- Don't have to understand all the details of how it was made
- Helps me make plots
- Can order online
- Made by somebody else
- Contains flavor information



#### Available at a computer near you!

Photo by Kelly Sikkema on Unsplash

•

Admittedly, the last few are a stretch

# like my CMSThis Talklike I

# Start your day with a cup of

CMS Open Data, who uses it, and how it's being used



Don't have to understand all

Δ

My own experiences and anecdotes with CMS Open Data





Photo by Kelly Sikkema on Unsplas

This Talk

CMS Open Data, who uses it, and how it's being used



My own experiences and anecdotes with **CMS Open Data** 



#### **CMS Open Data**

Google

cms open data				× 🔅 Q	
News	For medicare	Images	Open Payments	Payments Search	Payment

# According to Google...

About 743,000,000 results (0.44 seconds)

CMS (.gov) https://data.cms.gov :

#### CMS data

Official site of the Center's for Medicare & Medicaid Services (CMS) data. Find CMS program datasets, tools, and more.

Explore Data · Medicare Fee-For-Service... · Provider Data Catalog · About Us

https://openpaymentsdata.cms.gov

#### **Open Payments - CMS**

The **Open** Payments Search Tool is used to search payments made by drug and medical device companies to physicians, physician assistants, advanced practice nurses ...

https://www.cms.gov > openpayments > data

#### **Open Payments Data Overview**

7 days ago — Open Payments data is publicly accessible information about payments and transfers of value that reporting entities make to covered ...

https://www.cms.gov > newsroom > data

#### Data

This data tool lets you filter publicly available data sets by geography, health care setting, and document types. You can also sign up for email updates to ...

#### CERN

https://opendata.cern.ch > docs > about-cms

#### About CMS

All CMS publications are open access. Some of the papers also include open data in the form of additional tables, plots, graphs and Rivet packages. Policies.

#### **CMS Open Data**

7

Google	cms open data X
	News         For medicare         Images         Open Payments         Payments         Payments
	About 743,000,000 results (0.44 seconds)
	CMS (.gov) https://data.cms.gov :: Not this! CMS data Official site of the Center's for Medicare & Medicaid Services (CMS) data. Find CMS program datasets, tools, and more. Explore Data · Medicare Fee-For-Service · Provider Data Catalog · About Us https://openpaymentsdata.cms.gov :: Open Payments - CMS The Open Payments Search Tool is used to search payments made by drug and medical device companies to physicians, physician assistants, advanced practice nurses https://www.cms.gov · openpayments · data :: Open Payments Data Overview 7 days ago — Open Payments data is publicly accessible information about payments and
	transfers of value that reporting entities make to covered         https://www.cms.gov > newsroom > data ::         Data         This data tool lets you filter publicly available data sets by geography, health care setting, and document types. You can also sign up for email updates to         CERN         https://opendata.cern.ch > docs > about-cms :         About CMS         All CMS publications are open access. Some of the papers also include open data in the form of additional tables, plots, graphs and Rivet packages. Policies.

# According to Google...

Х

#### **CMS Open Data**

#### http://opendata.cern.ch/ opendata Help About • CERN Explore more than three petabytes of open data from particle physics! Search Start typing... search examples: collision datasets, keywords:education, energy:7TeV **Explore** Focus on **ATLAS** datasets software ALICE CMS environments documentation LHCb **OPERA**

#### In 2020...



[Thaler, Adventures with Public Collider Data (2020)]



[Thaler, Adventures with Public Collider Data (2020)]

Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

10



[Thaler, Adventures with Public Collider Data (2020)]

Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

11

# **Some Fun Recent Highlights**

\*With a strong bias towards research done my members of my group/home institution





$$\operatorname{ENC}(R_L) = \left(\prod_{k=1}^N \int d\Omega_{\vec{n}_k}\right) \delta(R_L - \Delta \hat{R}_L) \\ \cdot \frac{1}{(E_{\text{jet}})^N} \left\langle \mathcal{E}(\vec{n}_1) \mathcal{E}(\vec{n}_2) \dots \mathcal{E}(\vec{n}_N) \right\rangle$$

13

**Energy-Energy** 

**Correlators** (*EEC*s (and  $E^NCs$ )) let us explore different aspects of QCD, including scaling behavior, collinear structure, phase transitions, and more

Explored in CMS Open Data!



Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

14



15

Different length scales probe different regimes of QCD!

Can see it all within Open Data!



# **Example 2: Jet Topics**

Slight difference in detector response for forward vs. central **quark** vs. **gluon** jets



A jet *x* is never purely a **quark** jet or a **gluon** jet, but rather a mixture:

$$p_{\text{mixed}}(\vec{x}) = f_q \, p_{\text{quark}}(\vec{x}) + (1 - f_q) \, p_{\text{gluon}}(\vec{x})$$

#### **Example 2: Jet Topics**

Can turn these quark/gluon distributions into measurements of **fundamental constants** of QCD in CMS Open Data!



Correlation dimensions are defined using Wasserstein geometry, ask me about it later!

# This Talk

CMS Open Data, who uses it, and how it's being used



My own experiences and anecdotes with CMS Open Data



19

# Open Data as a teaching tool

Me as an undergrad in 2018 joining the Rutgers CMS B2G Group



20

#### My favorite dataset: CMS2011AJets

Jet Data collected in 2011 Run A

Applied HLT Jet300 single-jet trigger

AK5 Jets with  $p_{\tau}$  > 375 GeV

AOD files located at Record 21, with associated MC (in both SIM/GEN varieties) at Records 1364 - 1369

Perfect for QCD & Jet studies!

#### http://opendata.cern.ch/record/21

Jet primary dataset in AOD format from RunA of 2011 (/Jet/Run2011A-12Oct2013v1/AOD)

/Jet/Run2011A-12Oct2013-v1/AOD, CMS collaboration

Cite as: CMS collaboration (2016). Jet primary dataset in AOD format from RunA of 2011 (/Jet/Run2011A-12Oct2013-v1/AOD). CERN Open Data Portal. DOI:10.7483/OPENDATA.CMS.UP77.P6PQ

opendata

CERN



Pictured: Dijet mass of QCD samples from CMS Open Sim at truth and detector level, [RG, Nachman, Thaler, 2205.05084]

https://energyflow.network/docs/datasets/ [Komiske, Mastandrea, Metodiev, Naik, Thaler, 1908.08542] [Tripathee, Xue, Larkoski, Naik, Marzani, 1704.05842]

#### My favorite way to access open data: opendata The MIT Open Data (MOD) Format

Processed AOD files into manageable "MOD HDF5" text files hosted at https://zenodo.org/record/3340205

Very easy to access – no CMSSW, no virtual machines, no ROOT, no complicated AODs ...

Can easily download *anywhere* on *any machine* with *energyflow*:

```
import energyflow as ef
# Load data
specs = [f'{500} <= corr_jet_pts <= {1000}', f'abs_jet_eta < {1.9}', f'quality >= {2}']
sim = ef.mod.load(*specs, dataset='cms')
```



#### Try pip install energyflow

Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

CERN

https://energyflow.network/docs/datasets/ [Komiske, Mastandrea, Metodiev, Naik, Thaler, <u>1908.08542</u>] [Tripathee, Xue, Larkoski, Naik, Marzani, <u>1704.05842</u>]

# My favorite way to access open data: Opendata The MIT Open Data (MOD) Format

Processed AOD files into manageable "MOD HDF5" text files hosted at <u>https://zenodo.org/record/3340205</u>
Very easy to access – no CMSSW, no virtual machines, no ROOT, no complicated AODs …
Can easily download anywhere on any machine with energy flow:
This is the reason why CMS2011AJets is my favorite dataset – it's the easiest one to access!

# import energyflow as ef # Load data specs = [f'{500} <= corr\_jet\_pts <= {1000}', f'abs\_jet\_eta < {1.9}', f'quality >= {2}'] sim = ef.mod.load(\*specs, dataset='cms')

23

Try pip install energyflow

This code follows the example analysis I made at https://github.com/rikab/SHAPER/blob/main/examples/example.ipynb [Ba, Dogra, RG, Tasissa, Thaler, 2302.12266] **Try** pip install pyshaper

# My typical workflow:

#### **Step 1**: Download CMS Open Data!

Pictured: An AK5 Jet measured during Run A in 2011



does minor preprocessing, and converts to *np* arrays

On a fresh machine, takes only 5 minutes to download a 100,000 jet sample



Ease of download makes open data great as an example data set (especially for tutorials)! I don't have to worry about Pythia. Geant, etc ...

-0.25

0.00

Rapidity

0.25

0.50

This code follows the example analysis I made at <u>https://github.com/rikab/SHAPER/blob/main/examples/example.ipynb</u> [Ba, Dogra, **RG**, Tasissa, Thaler, <u>2302.12266</u>] Try pip install pyshaper

### My typical workflow:

#### Step 2: Set up calculations, e.g.

<pre># Sample from a normalized uniform distribution def uniform_sampler(N, param_dict):     points = torch.FloatTensor(N, 2).uniform_(-R, R).to(device)     zs = torch.ones((N,)).to(device) / N     return (points, zs)</pre>
_isotropy = Observable({}, uniform_sampler, beta = beta, R = R)
*****************
##### N-Point-Ellipsiness ##### ##############################
<pre># Sample points from N uniform ellipses plus weighted points at their center def point_ellipse_sampler(N, param_dict):</pre>
<pre>centers = param_dict["Points"].params num = param_dict["Points"].N radiil = param_dict["Radius1"].params radii2 = param_dict["Radius2"].params angles = param_dict["Angles"].params weights = param_dict["Weights"].params</pre>
<pre>phi = 2 * np.pi * torch.rand(num, N).to(device) r = torch.sqrt(torch.rand(num, N)).to(device) points = torch.stack([radii1[:, None] * torch.cos(phi + angles[:, None]), radii2[:, None] * torch.sin(phi + angles[:, None])] points = torch.cat([point for point in points], dim=1)</pre>
<pre># Concatenate and reweight e = torch.cat([centers, points.T], dim=0) z1 = torch.cat([weights[i] * torch.ones((1,), device=device) for i in range(num)], dim=0) z2 = torch.cat([weights[num + i] * torch.ones((N,), device=device) / N for i in range(num)], dim=0) z = torch.cat([z1, z2], dim=0) return (e, z)</pre>
<pre>3pointellipsiness = Observable({"Points": Coordinates2D(3), "Weights": Simplex(2*3), "Radius1": PositiveReals(3, 0), "Radius2":  </pre>

For me, this usually involves defining QCD observables or building ML tools to act on the data – This is where all the physics happens!

This code follows the example analysis I made at <u>https://github.com/rikab/SHAPER/blob/main/examples/example.ipynb</u> [Ba, Dogra, **RG**, Tasissa, Thaler, <u>2302.12266</u>] Try pip install pyshaper

# My typical workflow:

Step 3: Run all calculations on the data!

```
plot_dictionary = {
    "plot_directory" : "Plots/Test",
    "gif_directory" : "Plots/Test/gifs",
    "extension" : "png",
    "title" : "CMS Jets"
}
# Initialize SHAPER
shaper = Shaper(observables, device)
shaper.to(device)
emds, params = shaper.calculate(dataset, epochs = 500, verbose=True, lr = 0.01, N = 100, scaling = 0.9, epsilon = 0.001)
```

(Often done on a big cluster rather than a Jupyter notebook ...)

This code follows the example analysis I made at https://github.com/rikab/SHAPER/blob/main/examples/example.ipynb [Ba, Dogra, RG, Tasissa, Thaler, 2302.12266] [RG, Thaler, Wu, WIP]

10

10

Data-Based Dijet Example

BSM, Delphes

BSM, Geant

BSM. Delphe

Data-Based Dijet Example

10-

QCD, Delphes

QCD, Geant

OCD, Delphes + O

### My typical workflow:

#### Step 4: Plots Plots Plots Plots Plots!



Try pip install GaussianAnsatz

#### **CMS Open Sim for Calibration**



[RG, Nachman, Thaler, <u>2205.05084</u> [RG, Nachman, Thaler, <u>2205.05084</u>] [RG, Nachman, Thaler, WIP] Try pip install GaussianAnsatz

#### CMS Open Sim for Calibration

29



[RG, Nachman, Thaler, 2205.03413]

Try pip install GaussianAnsatz

#### **CMS Open Sim for Uncertainty Estimation**

![](_page_29_Figure_3.jpeg)

... Using Open Data to understand **detector efforts** and quantify **uncertainties and correlations** with **ML**!

![](_page_29_Picture_5.jpeg)

### **Hearing the Shapes of Jets**

![](_page_30_Figure_2.jpeg)

Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

31

#### How wide are QCD jets?

![](_page_31_Figure_2.jpeg)

Determining the radius distribution of q/g jets in data with an **MIT Summer Research Program undergrad** (Xinyue Wu)! From zero to this in a few weeks!

![](_page_31_Picture_4.jpeg)

#### **Prototyping new metrics**

![](_page_32_Figure_2.jpeg)

## How do I take my Conclusion pen

# Start your day with a cup of

But it's good to have some variety in coffee!

How can we enable more datasets to be made easily accessible and useable?

![](_page_33_Picture_4.jpeg)

Photo by Kelly Sikkema on Unsplas

Admittedly, the last few are a stretch

Rikab Gambhir – CMS Open Data – FNAL – 11 July 2023

34