

Data ex Machina

Machine Learning with Jets in CMS Open Data

Machine Learning for Jet Physics 2020

Eric M. Metodiev

Center for Theoretical Physics

Massachusetts Institute of Technology



Patrick
Komiske



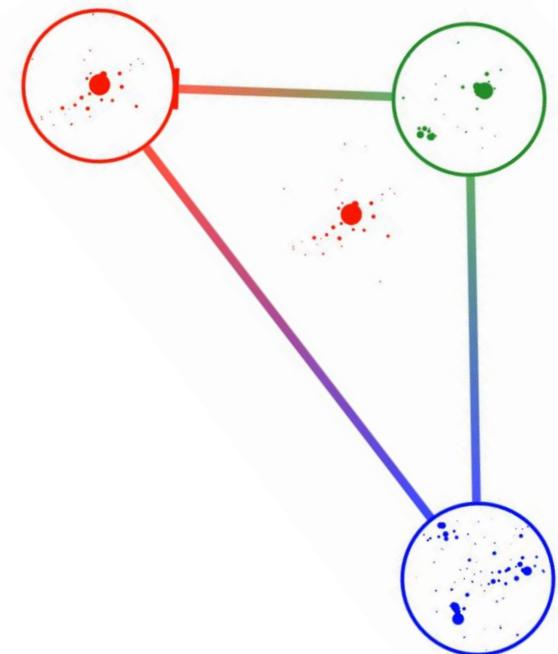
Radha
Mastandrea

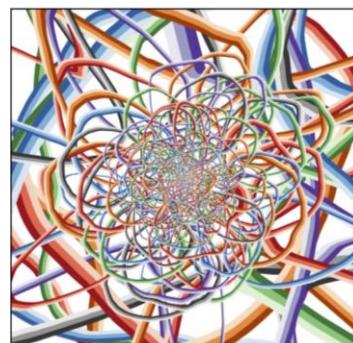
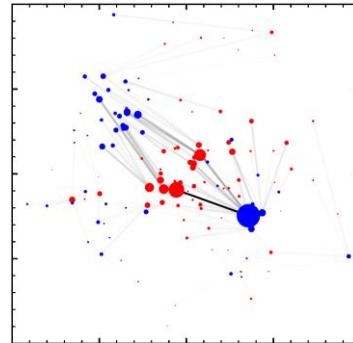
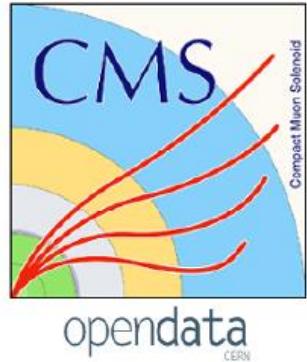


Preksha
Naik



Jesse
Thaler





CMS Open Data

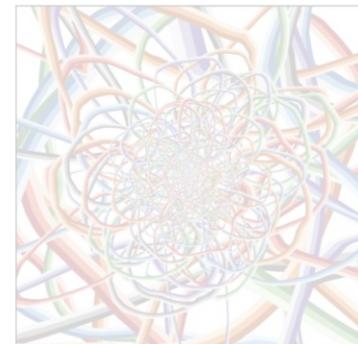
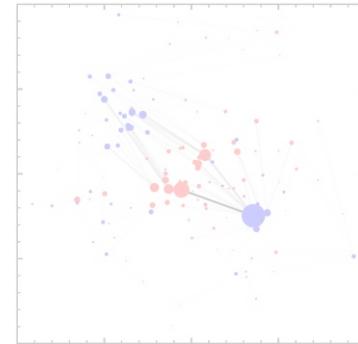
A new public dataset for jet studies

Unsupervised Learning

A metric for collider events

Supervised Learning

Training directly on collider data



CMS Open Data

A new public dataset for jet studies

Unsupervised Learning
A metric for collider events

Supervised Learning
Training directly on collider data

opendata.cern.ch

Explore more than **two petabytes**
of open data from particle physics!

jet primary dataset

search examples: [collision datasets](#), [keywords:education](#), [energy:7TeV](#)

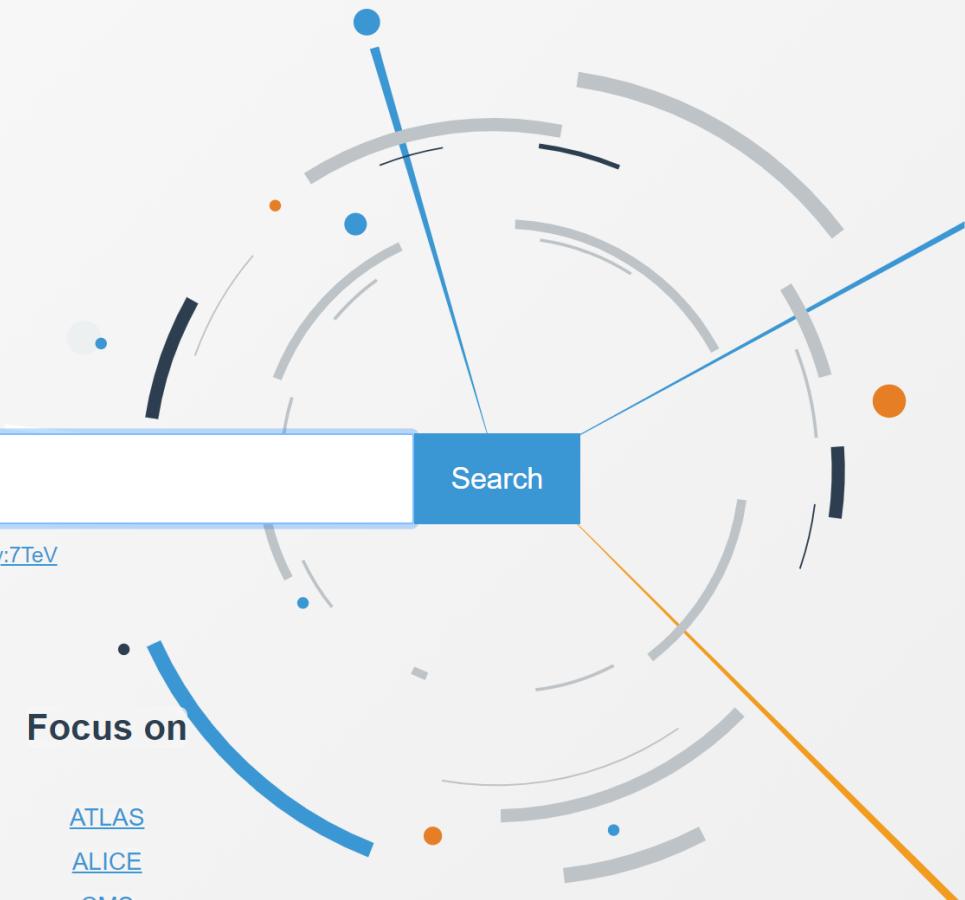
Explore

[datasets](#)
[software](#)
[environments](#)
[documentation](#)

Focus on

[ATLAS](#)
[ALICE](#)
[CMS](#)
[LHCb](#)
[OPERA](#)
[Data Science](#)

▼ Get started ▼



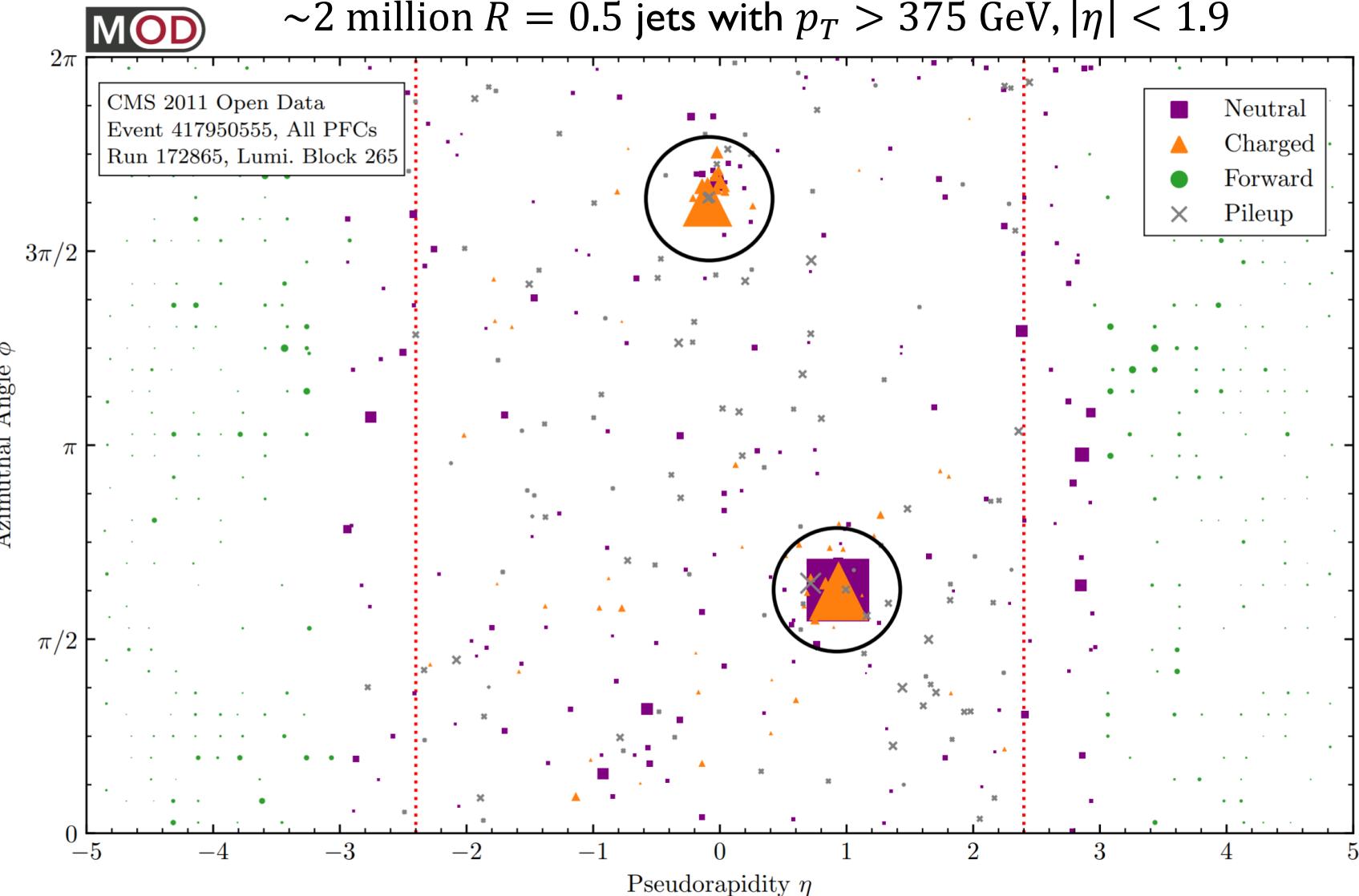


CMS 2011 A Jet Primary Dataset (+ Simulation)

2.3 fb^{-1} of 7 TeV proton-proton collision data. [\[link\]](#)



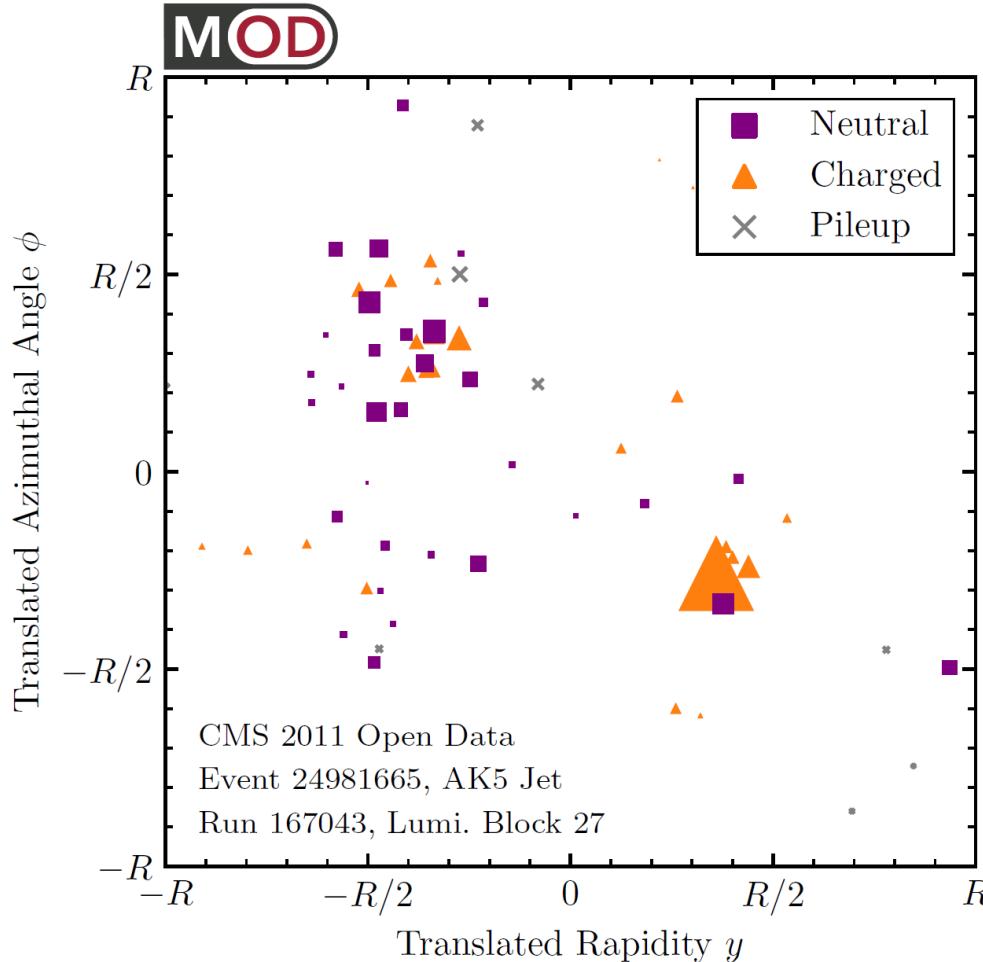
~2 million $R = 0.5$ jets with $p_T > 375 \text{ GeV}$, $|\eta| < 1.9$





CMS 2011 A Jet Primary Dataset (+ Simulation)

2.3 fb^{-1} of 7 TeV proton-proton collision data.



Our processed jet dataset is public!

[\[Komiske, Mastandrea, EMM, Naik, Thaler, 1908.08542\]](#)

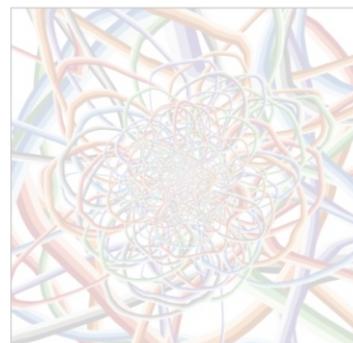
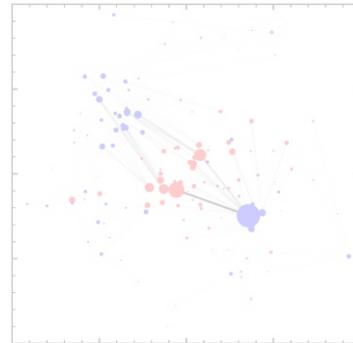
~2 million $R = 0.5$ anti- $k\text{T}$ jets recorded by CMS
 $p_T > 375 \text{ GeV}, |\eta| < 1.9$

Jets as lists of particle flow candidates:
 $[p_T, y, \phi, \text{ID}, \text{vertex}]$

Plus additional information:
Jet energy correction factors
Monte Carlo samples
Detector simulation

[Zenodo record](#)

[Binder demo](#) to download and read the dataset



CMS Open Data

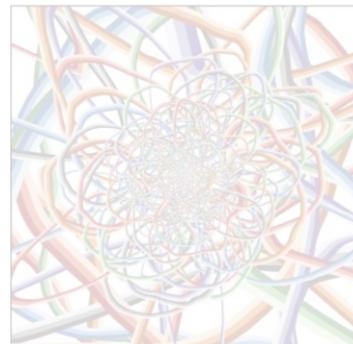
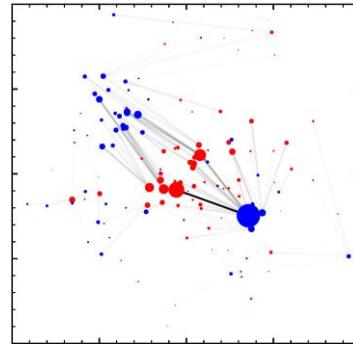
A new public dataset for jet studies

Unsupervised Learning

A metric for collider events

Supervised Learning

Training directly on collider data



CMS Open Data

A new public dataset for jet studies

Unsupervised Learning

A metric for collider events

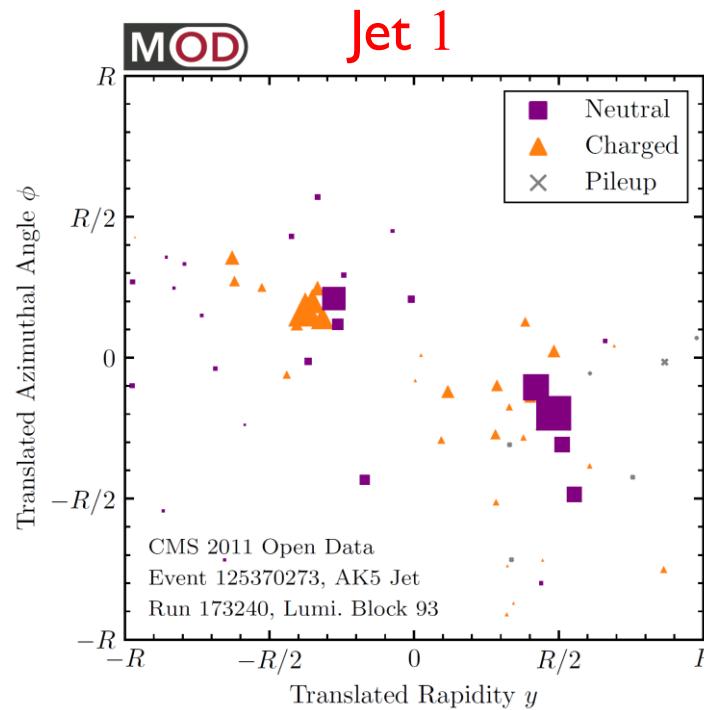
Supervised Learning

Training directly on collider data

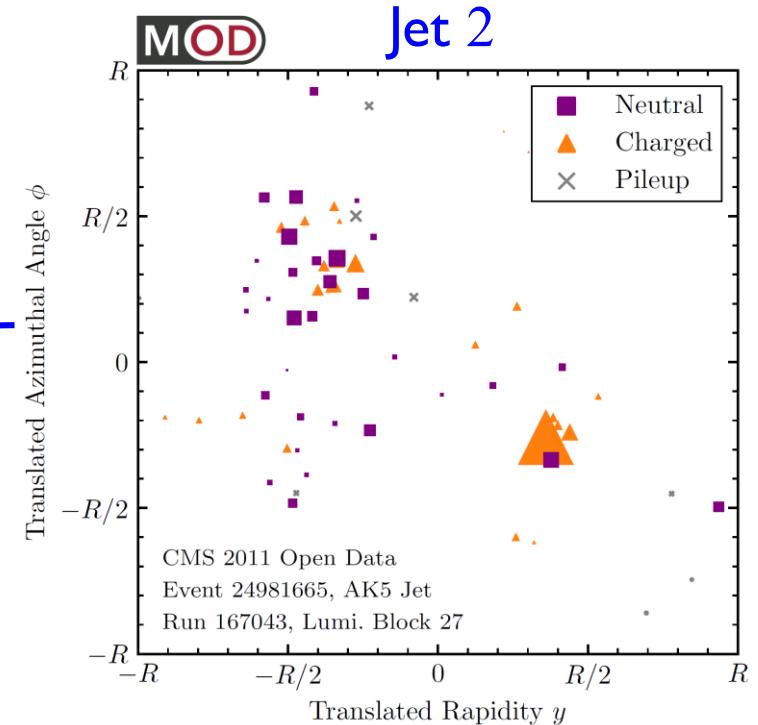
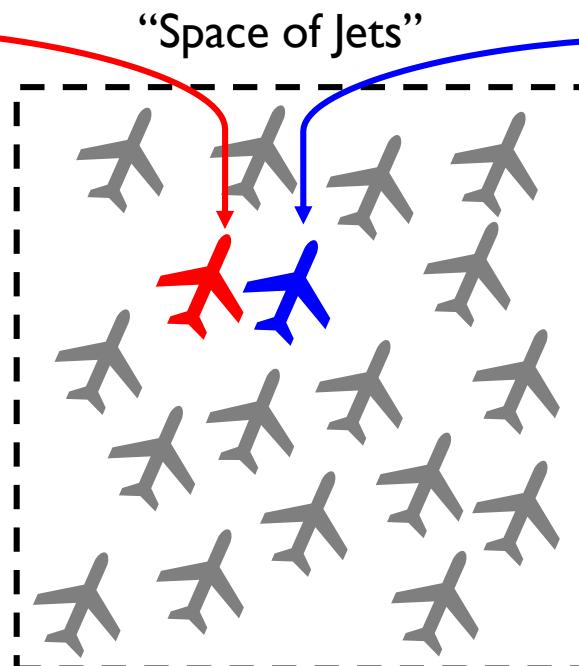
When are two jets similar?

Many unsupervised methods rely on a **distance matrix**. Need a physically-sensible **metric** between jets!

These two jets “look” similar, but have different numbers of particles, flavors, and locations.

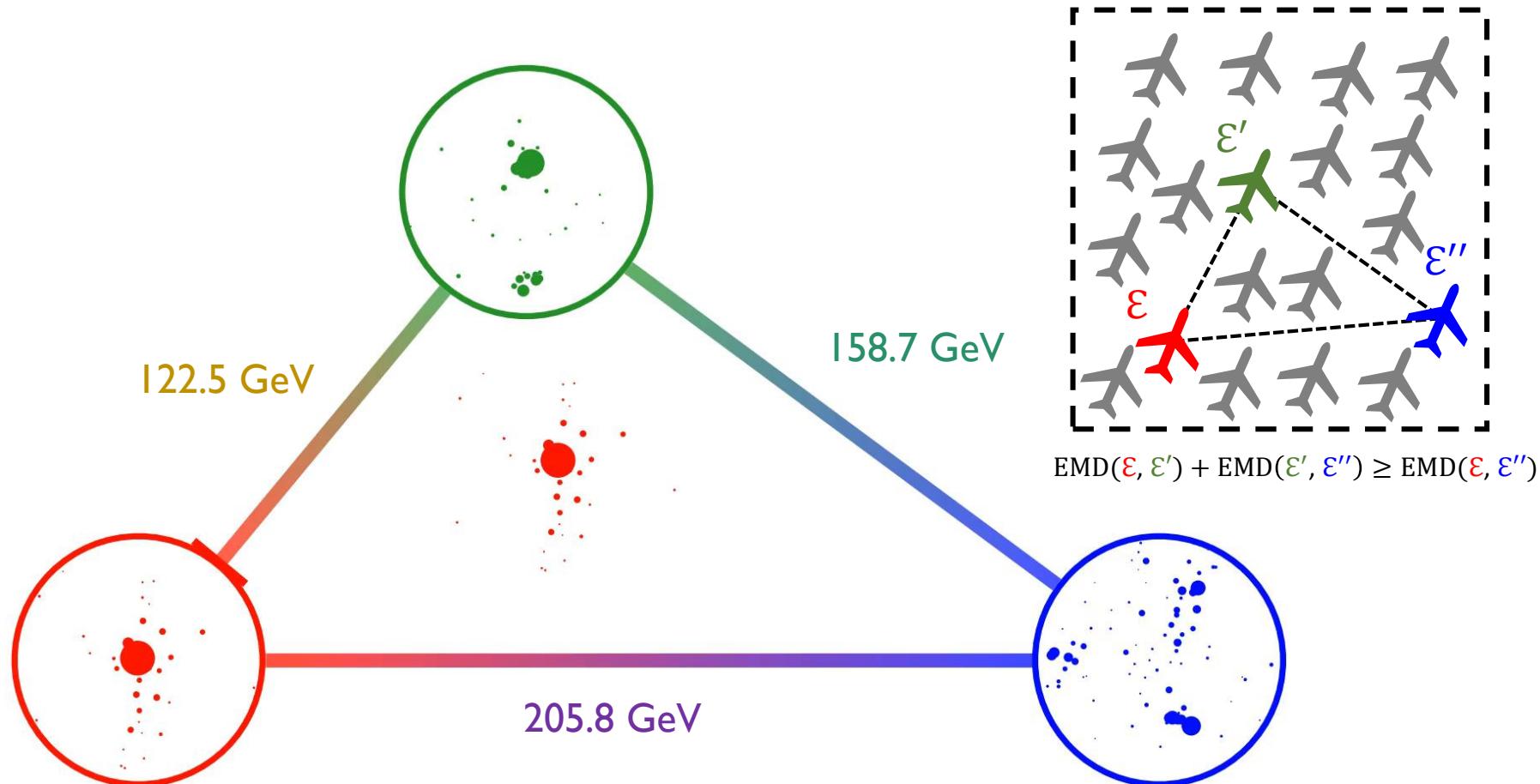


How do we quantify this?



The Energy Mover's Distance

“Energy” Mover’s Distance: the minimum “work” (**energy** \times angle) to rearrange one jet (pile of energy) into another

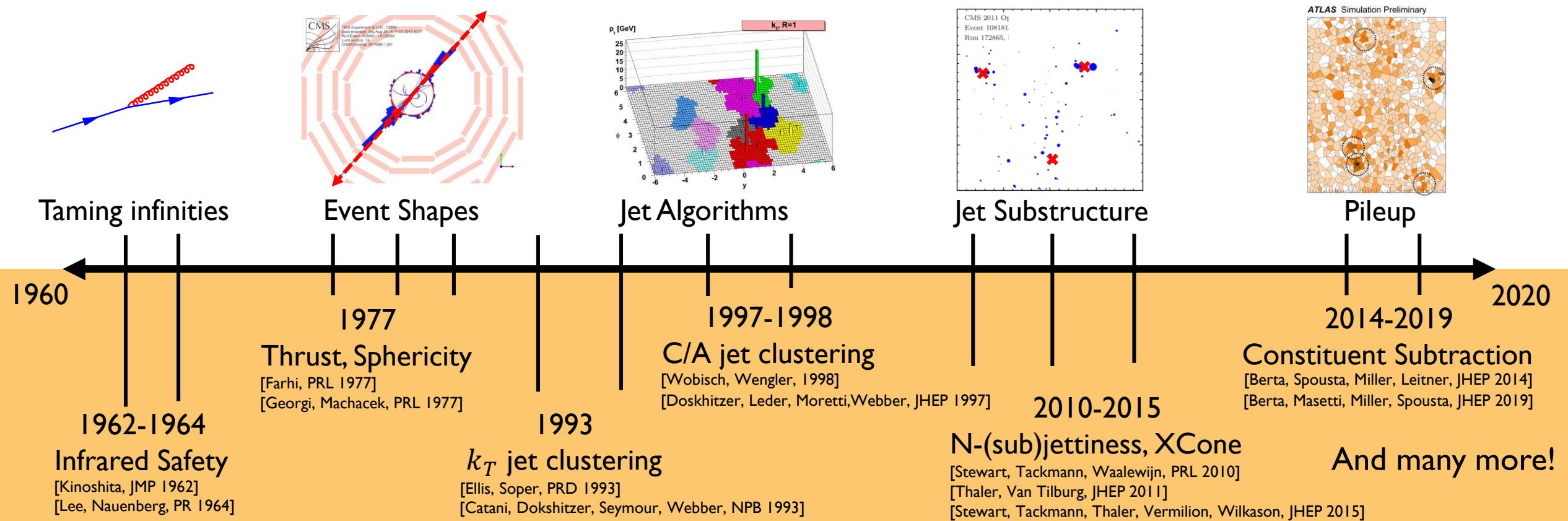


[Komiske, EMM, Thaler, 1902.02346]

Earth Mover's Distance: [Peleg, Werman, Rom] [Rubner, Tomasi, Guibas]

[See Today's Talk on Optimal Transport](#)

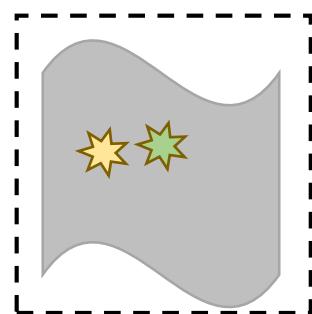
Six Decades of Collider Techniques



Six Decades of Collider Techniques as Optimal Transport!

[Komiske, EMM, Thaler, to appear]

Smooth function of energy distribution are finite in QFT



$$\text{EMD}(\mathcal{E}, \mathcal{E}') < \delta \\ \rightarrow |\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')| < \epsilon$$

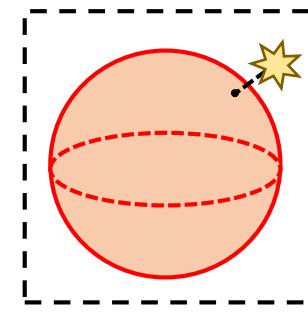
Taming infinities

1960

1962-1964

Infrared Safety
[Kinoshita, JMP 1962]
[Lee, Nauenberg, PR 1964]

Event shapes as distances to the 2-particle manifold



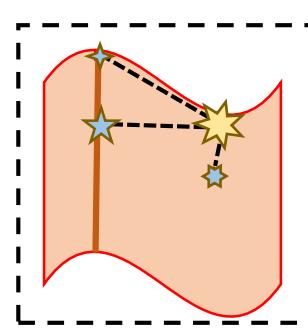
$$t(\mathcal{E}) = \min_{|\mathcal{E}'|=2} \text{EMD}(\mathcal{E}, \mathcal{E}')$$

Event Shapes

1977

Thrust, Sphericity
[Farhi, PRL 1977]
[Georgi, Machacek, PRL 1977]

Jets are N-particle event approximations

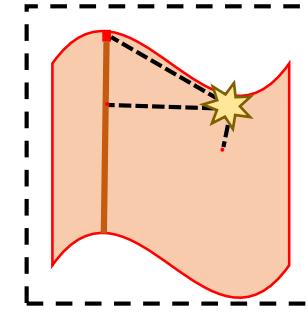


$$\mathcal{J}(\mathcal{E}) = \operatorname{argmin}_{|\mathcal{E}'|=N} \text{EMD}(\mathcal{E}, \mathcal{E}')$$

Jet Algorithms

k_T jet clustering
[Ellis, Soper, PRD 1993]
[Catani, Dokshitzer, Seymour, Webber, NPB 1993]

Subtract a pileup as a uniform distribution



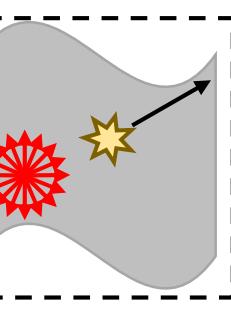
$$\mathcal{E} - \mathcal{U}$$

Jet Substructure

2010-2015

N-(sub)jettiness, XCone
[Stewart, Tackmann, Waalewijn, PRL 2010]
[Thaler, Van Tilburg, JHEP 2011]
[Stewart, Tackmann, Thaler, Vermilion, Wilkason, JHEP 2015]

And many more!



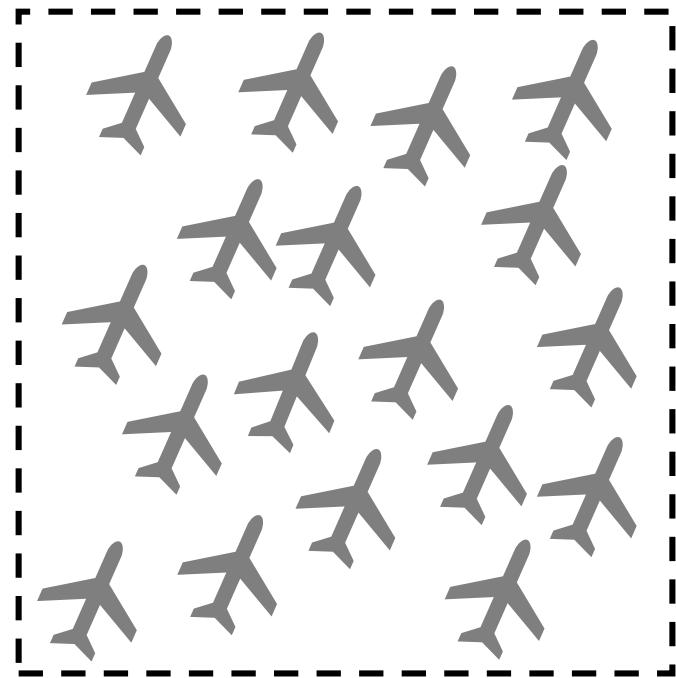
Pileup

2014-2019

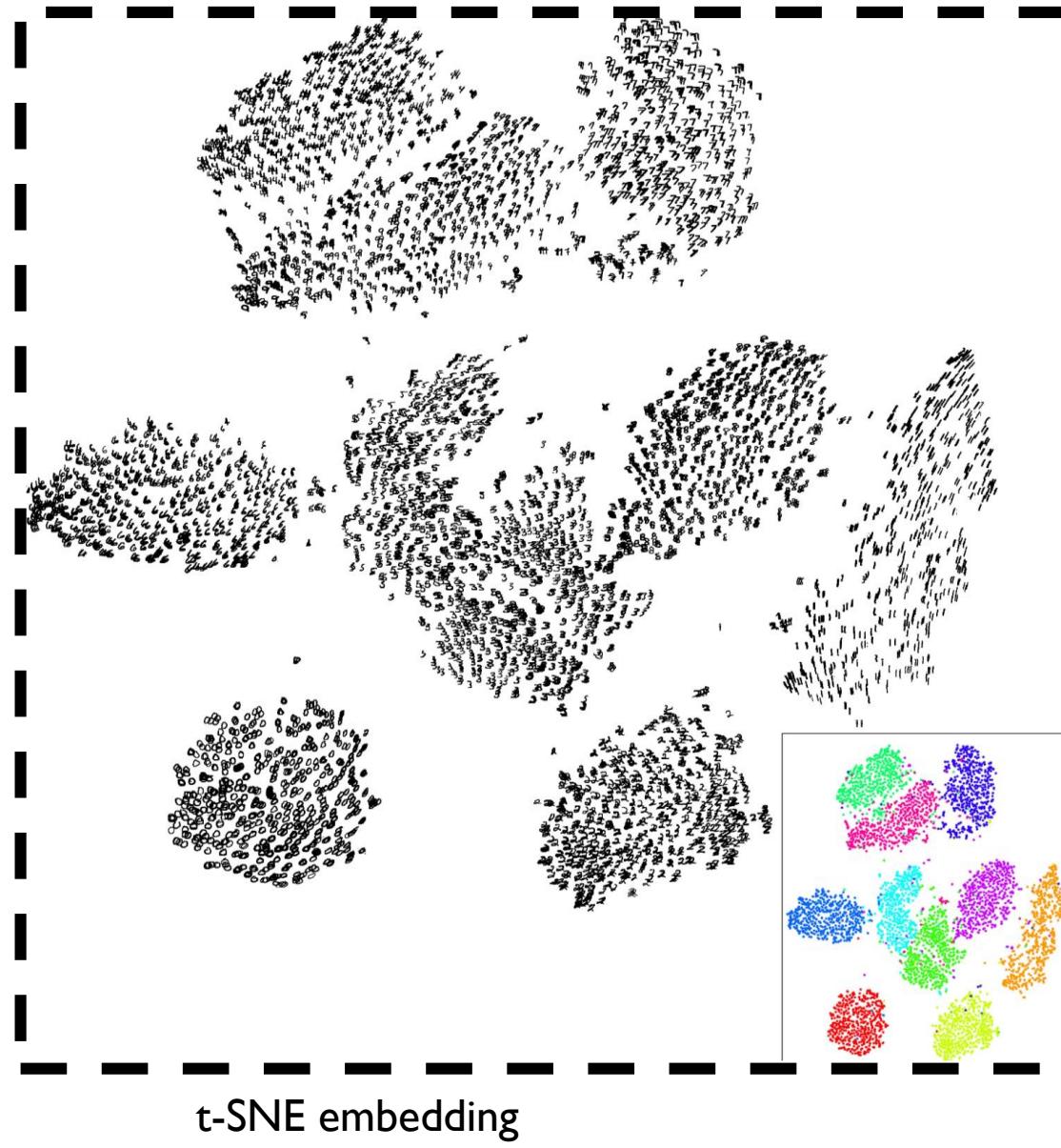
Constituent Subtraction
[Berta, Spousta, Miller, Leitner, JHEP 2014]
[Berta, Masetti, Miller, Spousta, JHEP 2019]

Visualizing the Manifold

What does the space of jets look like?

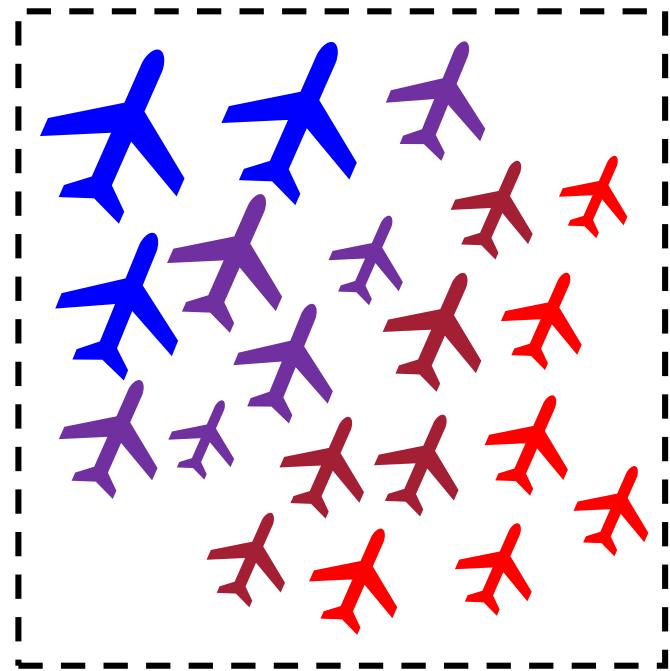


[van der Maaten, Hinton, JMLR 2008]



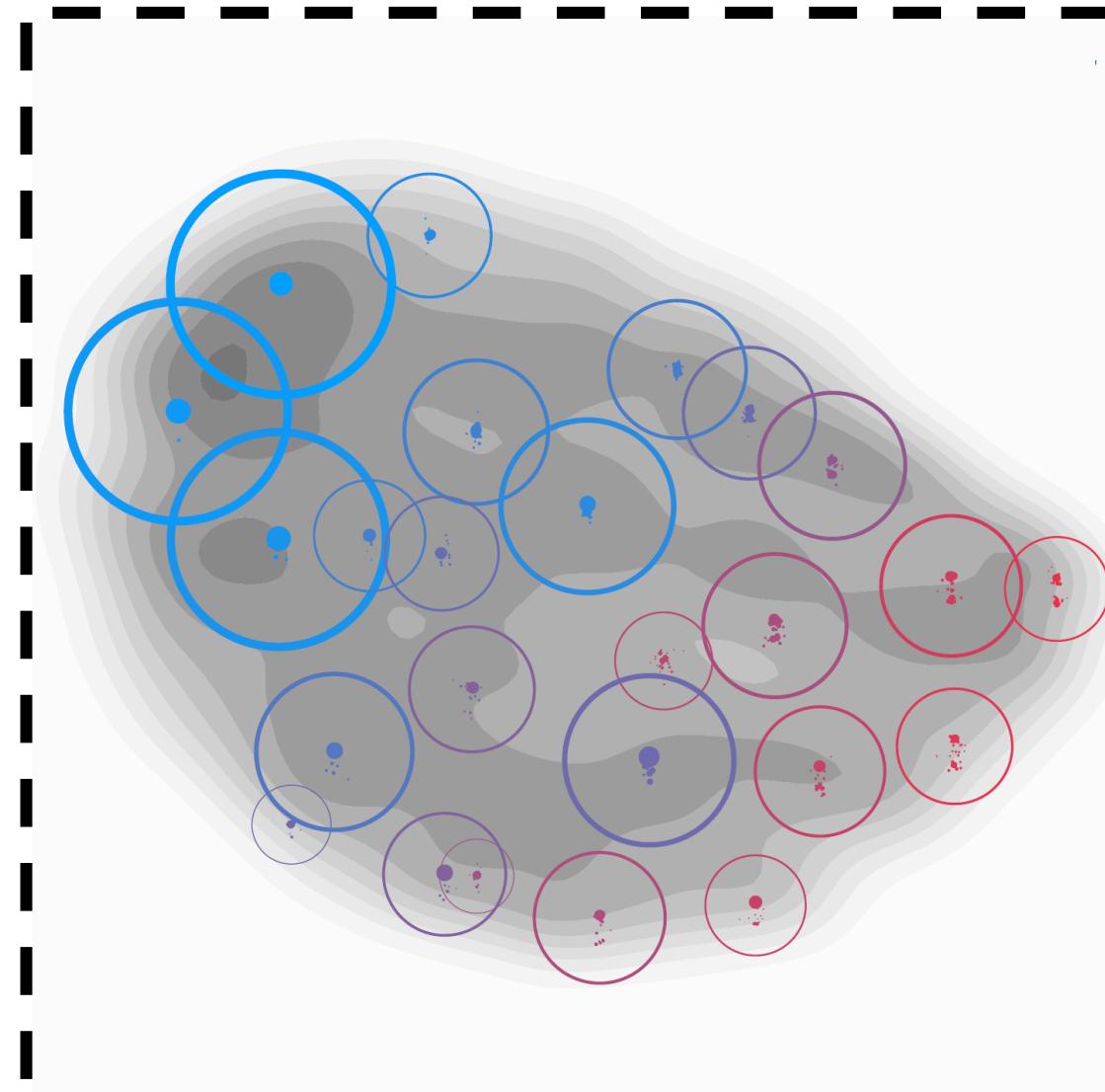
Visualizing the Manifold

What does the space of jets look like?



[van der Maaten, Hinton, JMLR 2008]

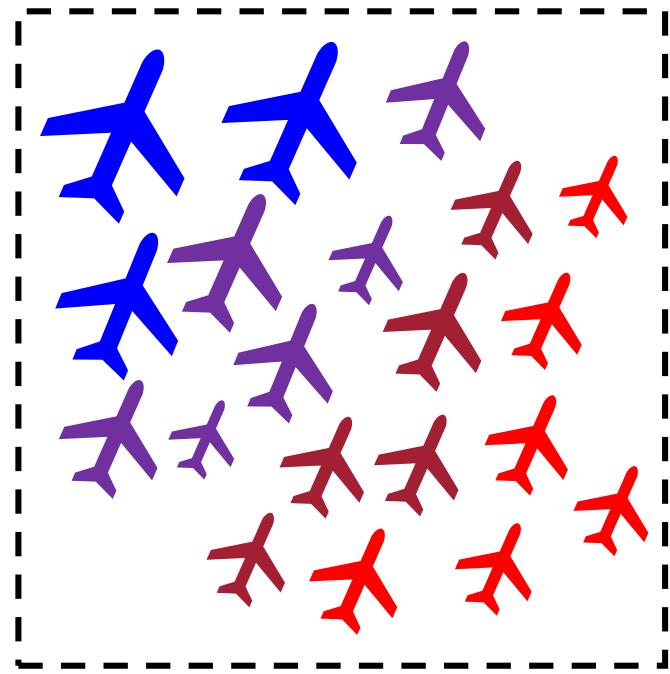
[\[Komiske, Mastandrea, EMM, Naik, Thaler, 1908.08542\]](#)



t-SNE embedding: 25-medoid jets shown

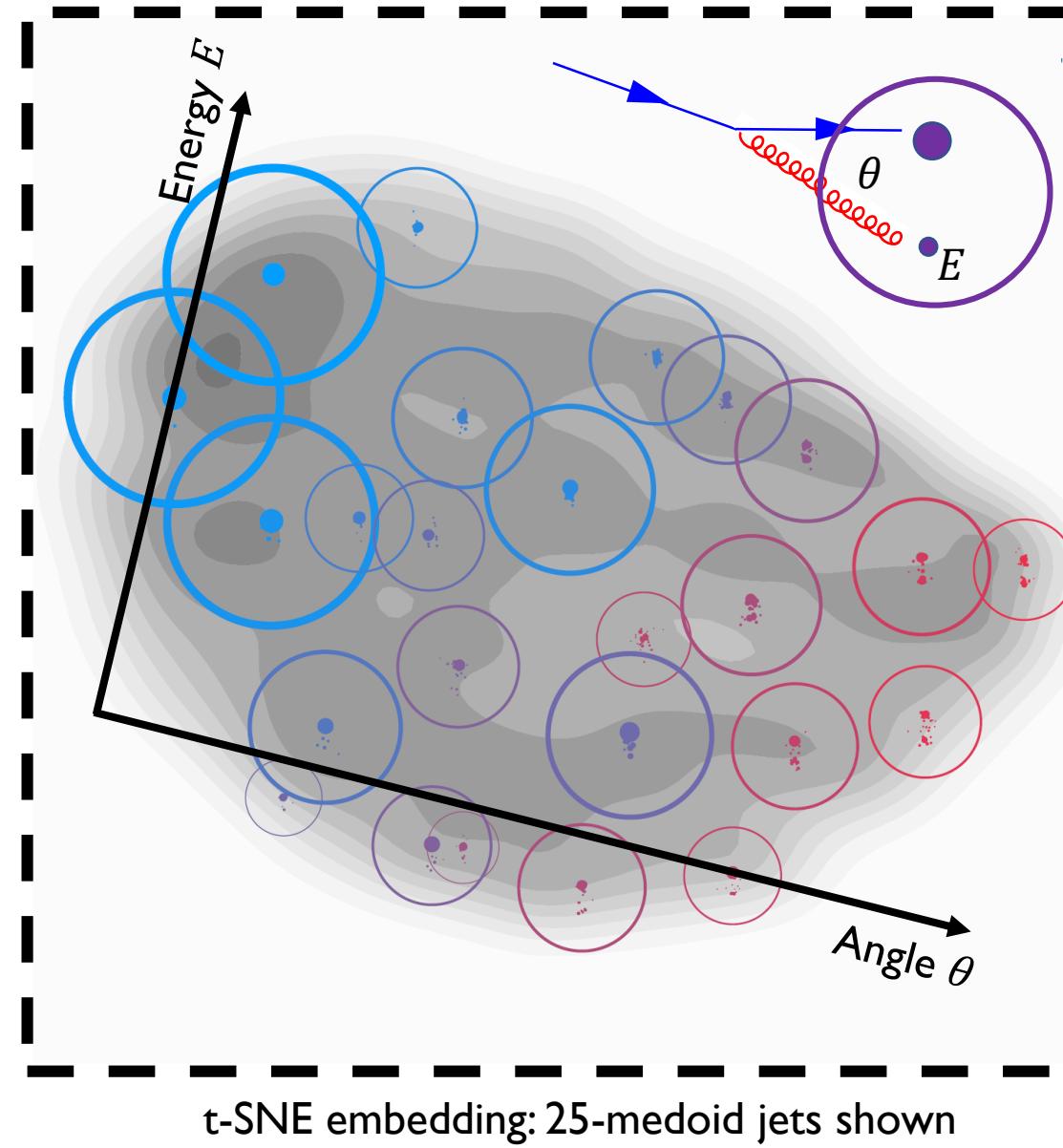
Visualizing the Manifold

What does the space of jets look like?



[van der Maaten, Hinton, JMLR 2008]

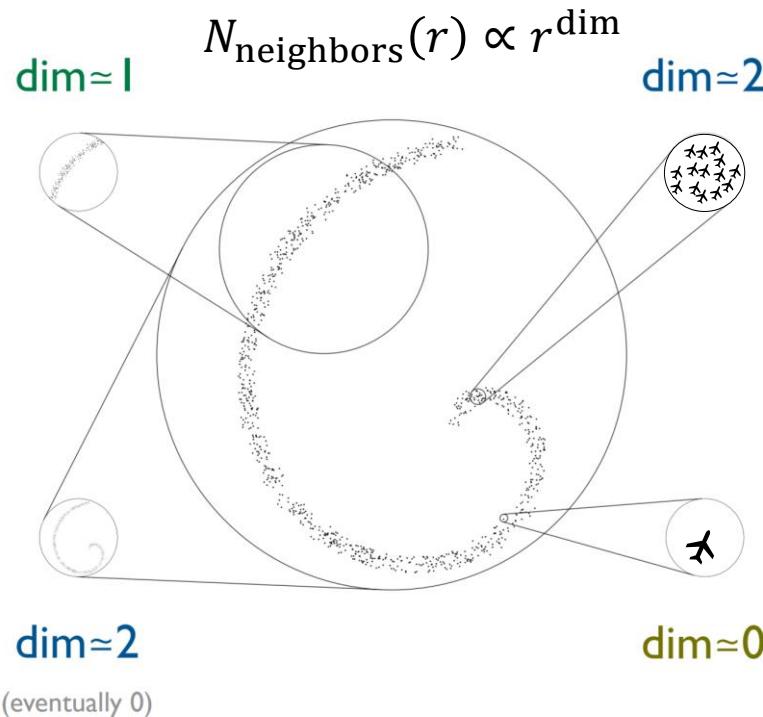
[Komiske, Mastandrea, EMM, Naik, Thaler, 1908.08542]



t-SNE embedding: 25-medoid jets shown

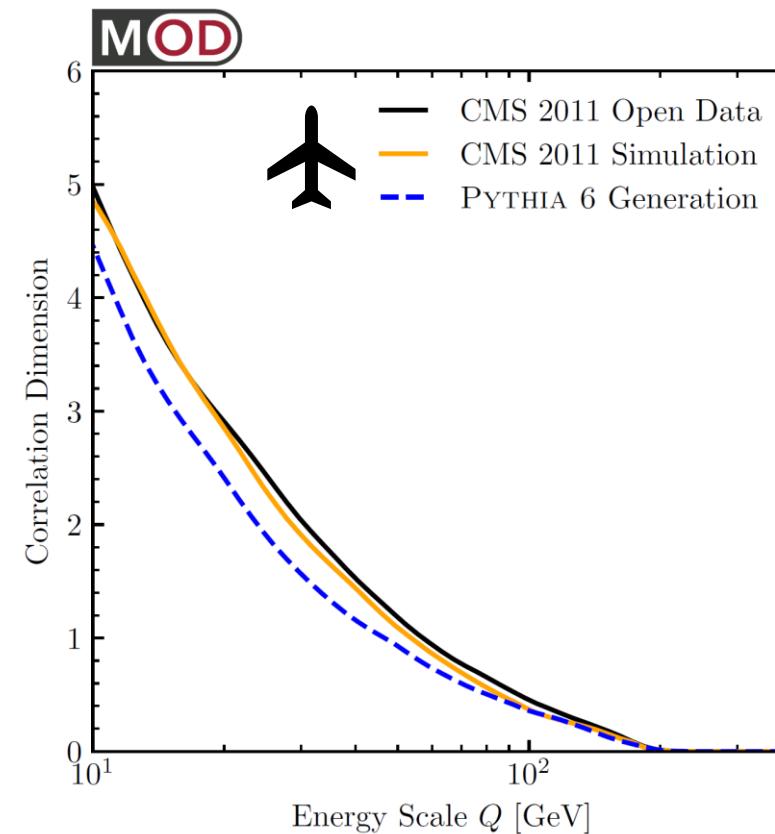
Correlation Dimension

Conceptual Idea



[Grassberger, Procaccia, PRL 1983] [Kegl, NeurIPS 2002]

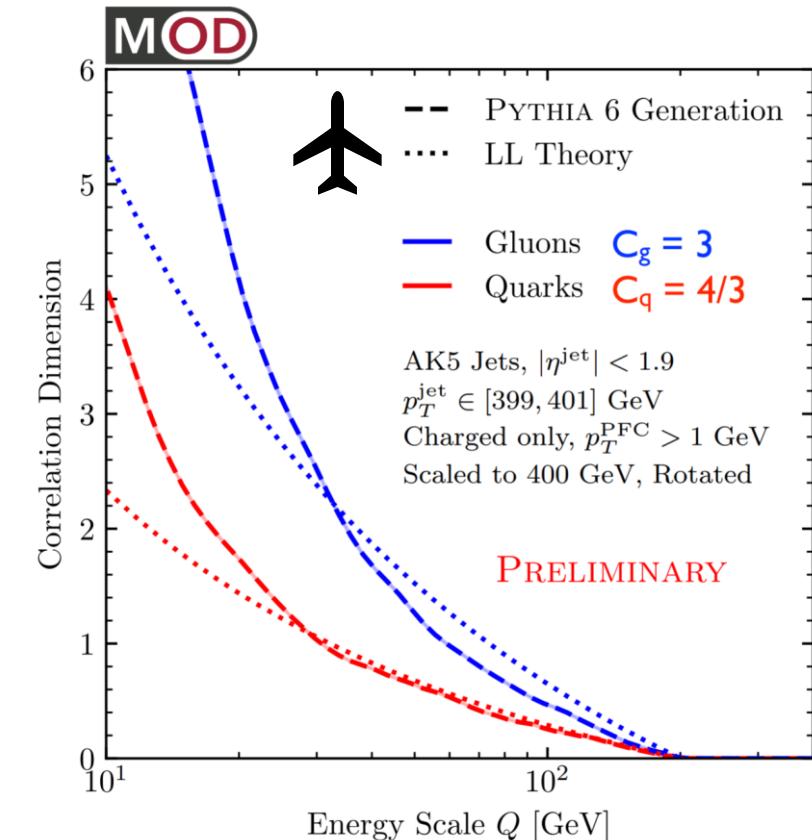
Experimental Data



Dimension blows up at low energies.

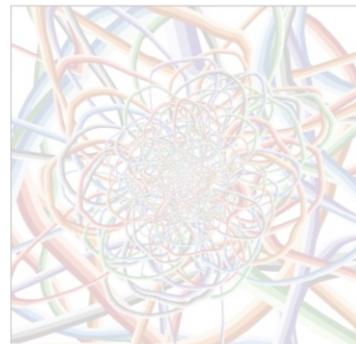
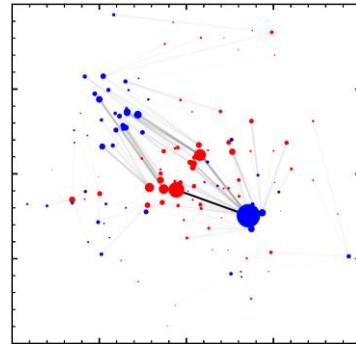
[Komiske, Mastandrea, EMM, Naik, Thaler, 1908.08542]

Theoretical Calculation



See extra slides for calculation sketch.

[See Jack's Talk](#) for more on dimensionality



CMS Open Data

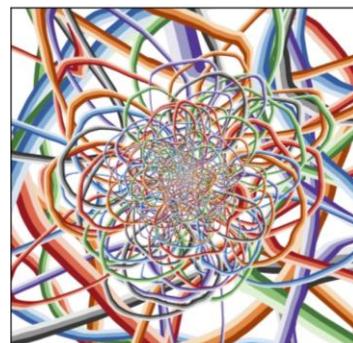
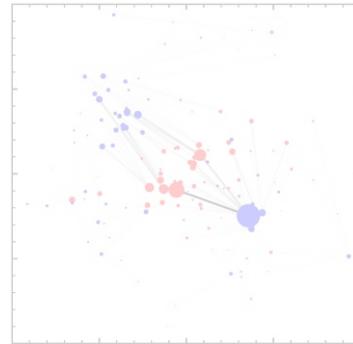
A new public dataset for jet studies

Unsupervised Learning

A metric for collider events

Supervised Learning

Training directly on collider data



CMS Open Data

A new public dataset for jet studies

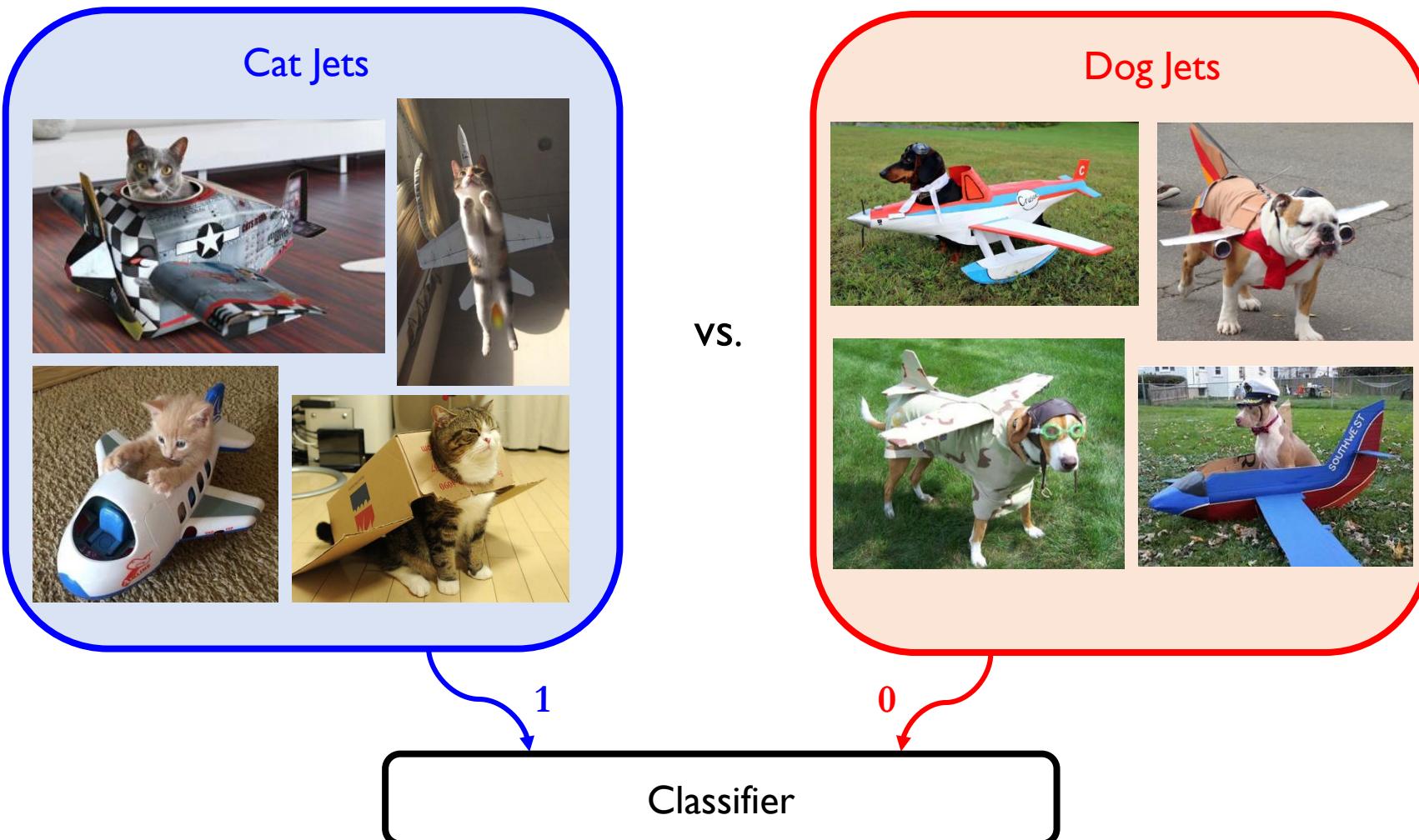
Unsupervised Learning

A metric for collider events

Supervised Learning

Training directly on collider data

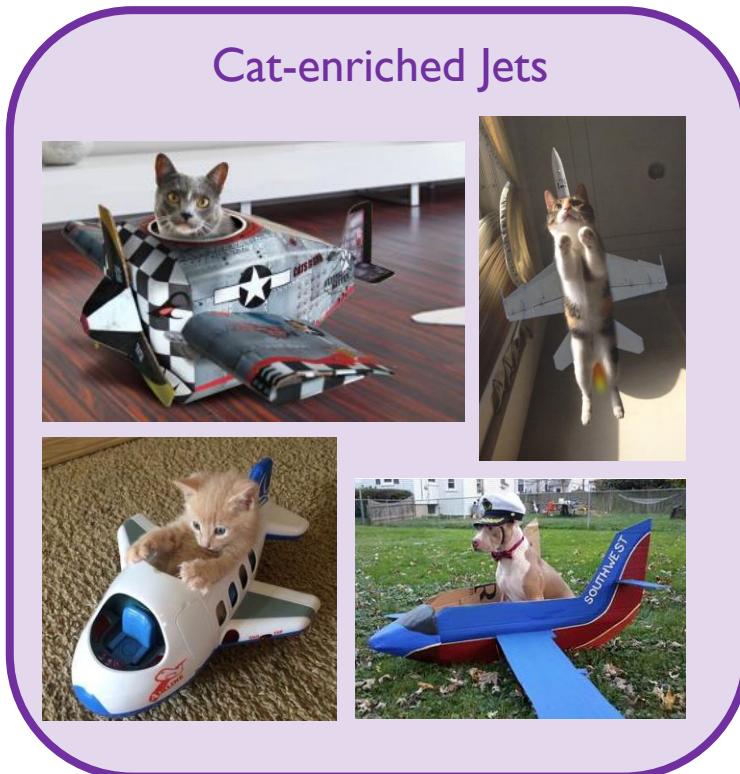
Training on pure samples: Cat jets vs. Dog jets



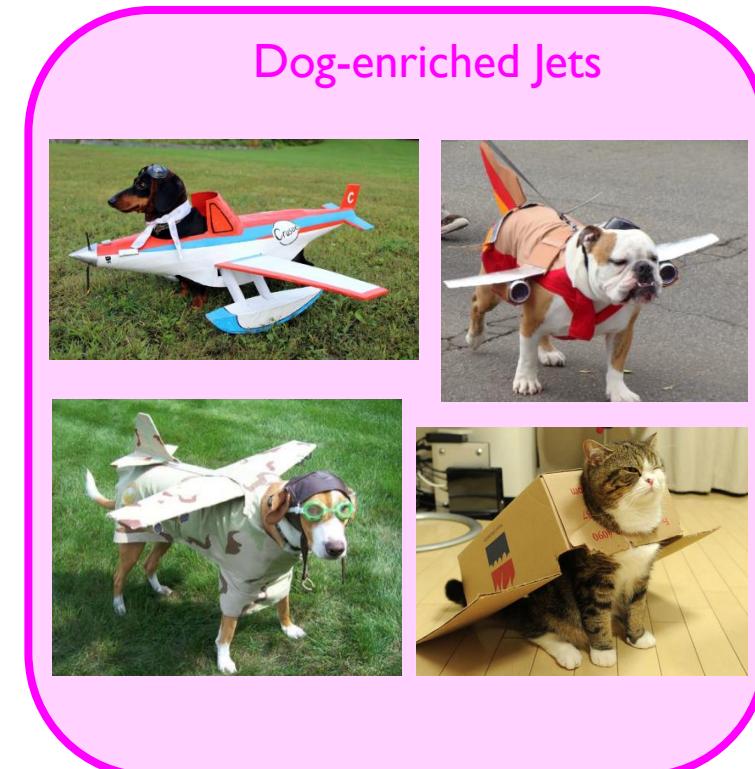
Training on mixed samples: Cat jets vs. Dog jets



Classification
Without Labels
(CWoLa)



vs.



1

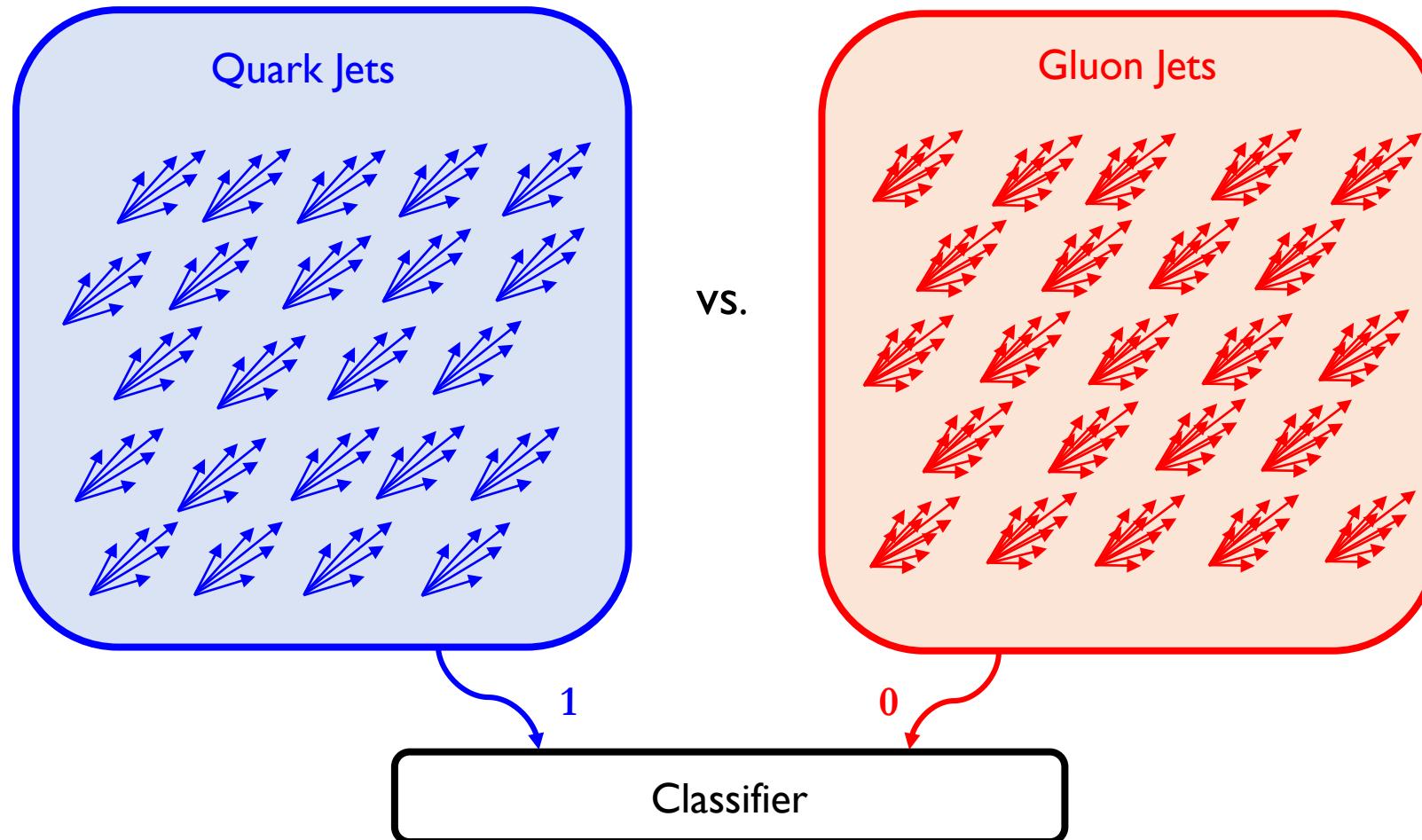
0

Classifier

Used by CMS for ttbb! [\[CMS 1909.05306\]](#)

This defines an equivalent classifier to the pure case!

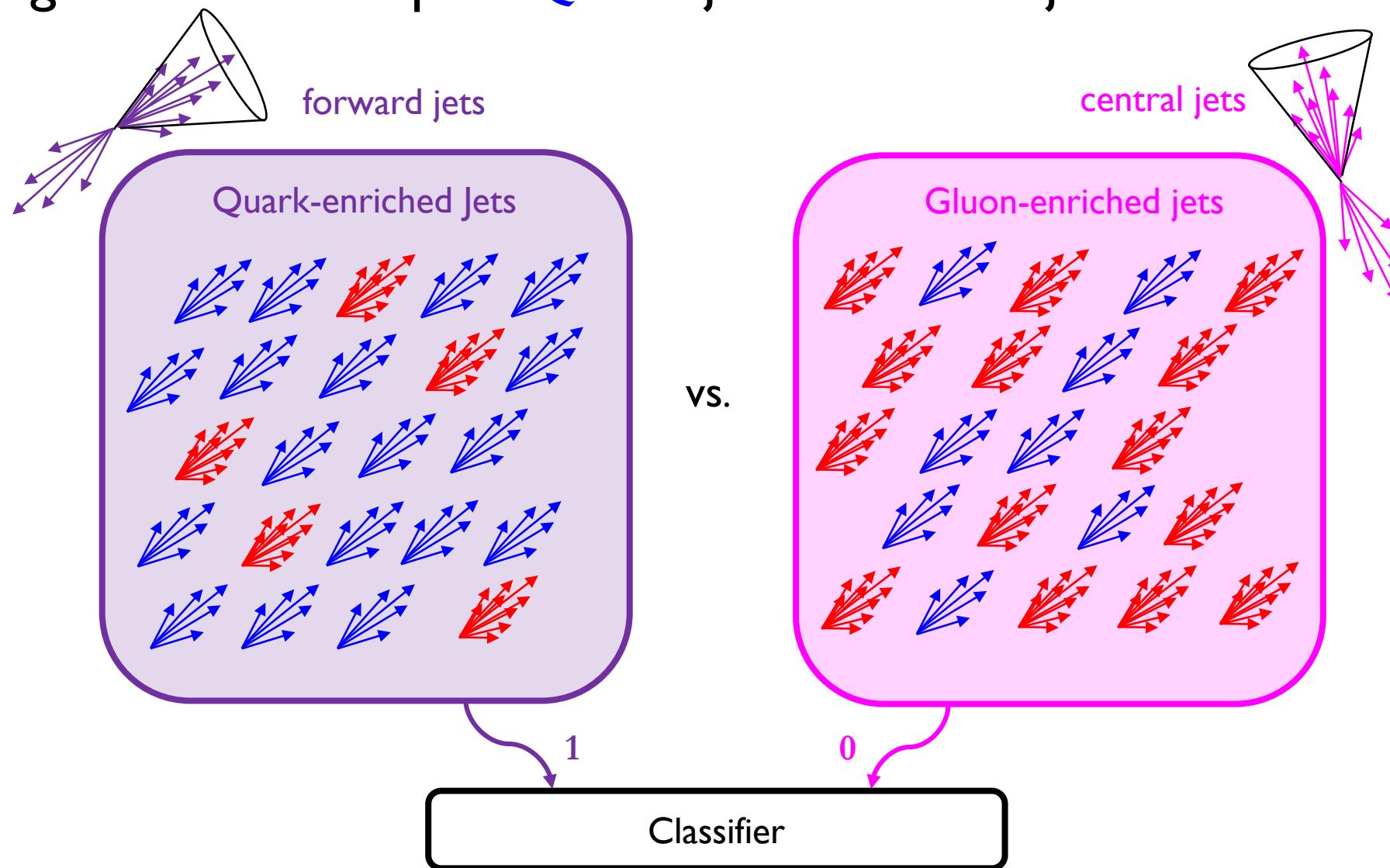
Training on pure samples: Quark jets vs. Gluon jets



Training on mixed samples: Quark jets vs. Gluon jets



Classification
Without Labels
(CWoLa)



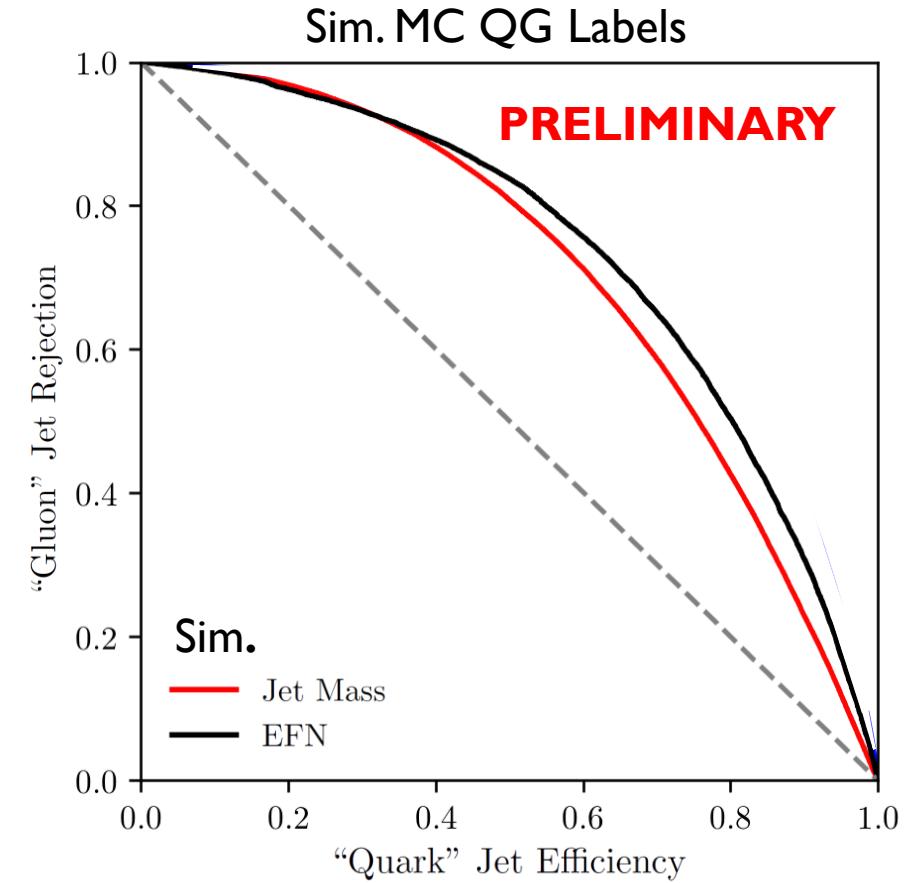
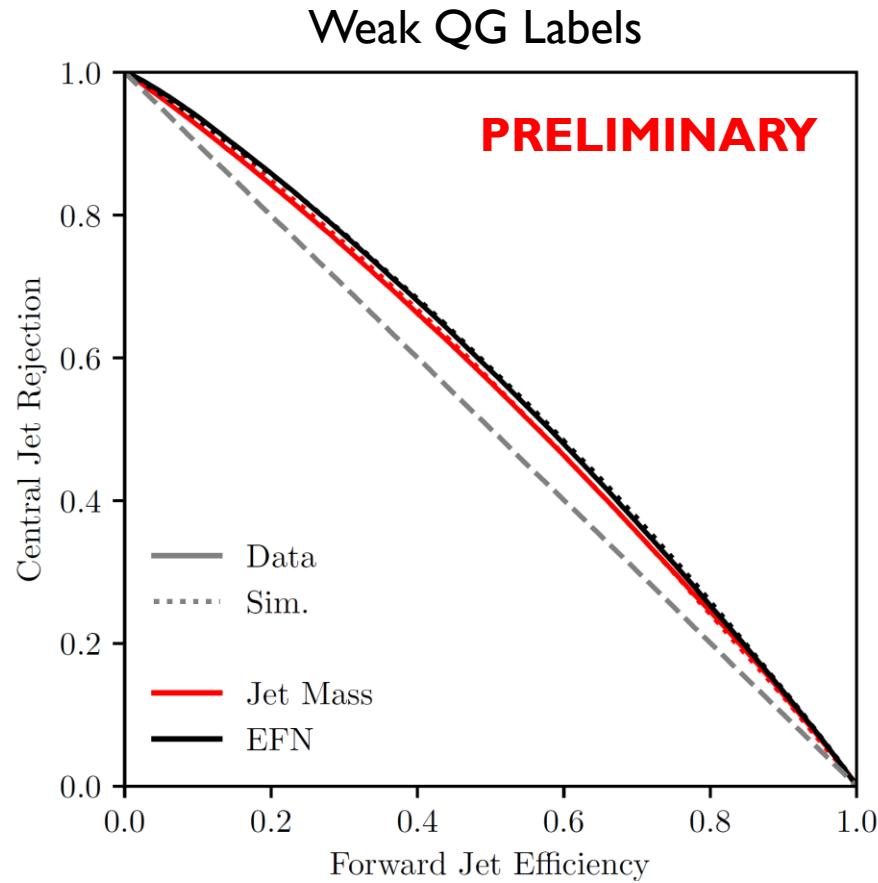
[EMM, B. Nachman, J. Thaler, 1708.02949]

[P.T. Komiske, EMM, B. Nachman, M.D. Schwartz, 1801.10158]

[L. Dery, B. Nachman, F. Rubbo, A. Schwartzman, 1702.00414] [T. Cohen, M. Freytsis, B. Ostdiek, 1706.09451]

Training on Data!

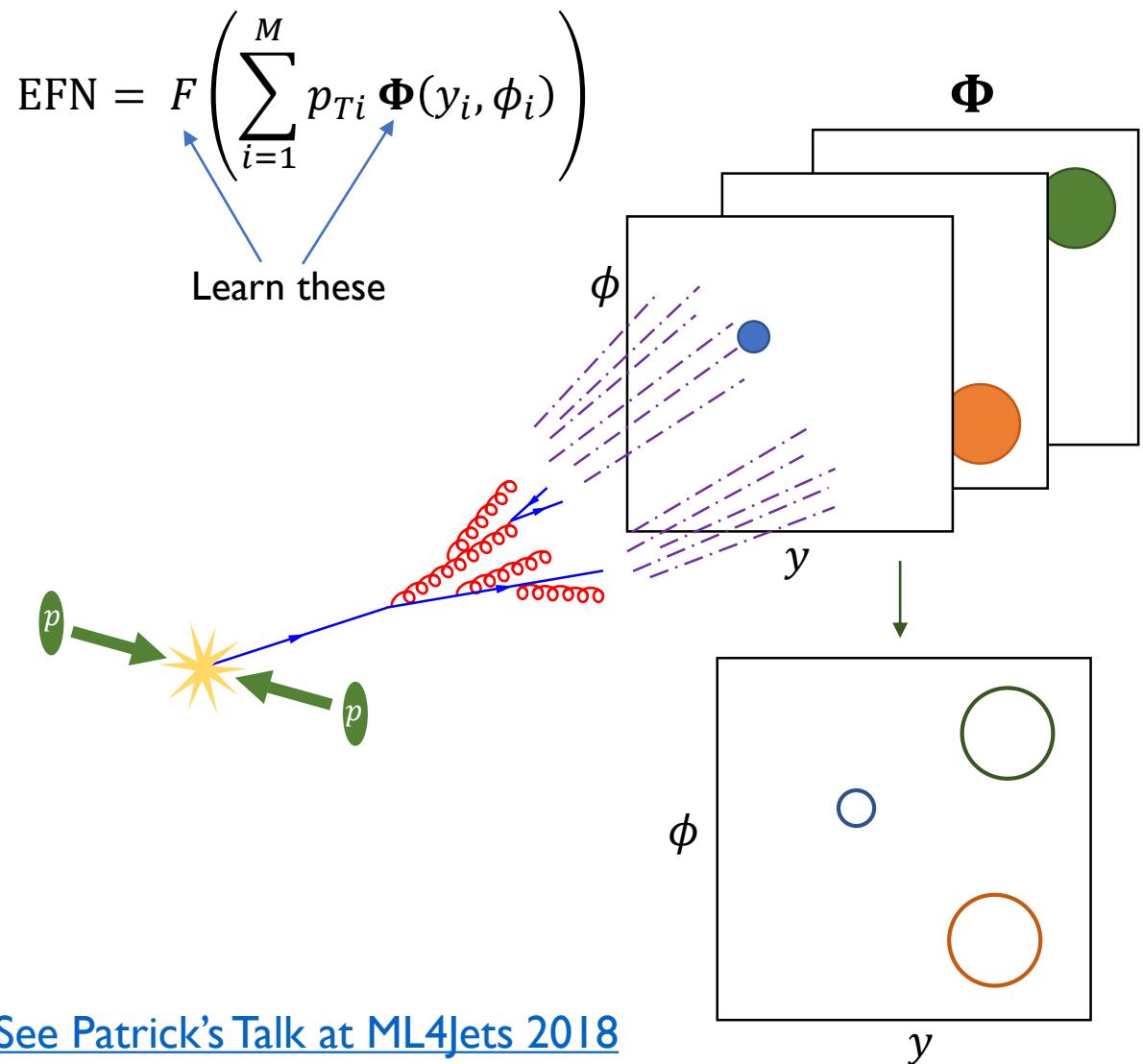
Central Jets ($|\eta^{\text{jet}}| < 0.7$): ~45% quark jets
Forward Jets ($|\eta^{\text{jet}}| > 0.7$): ~65% quark jets



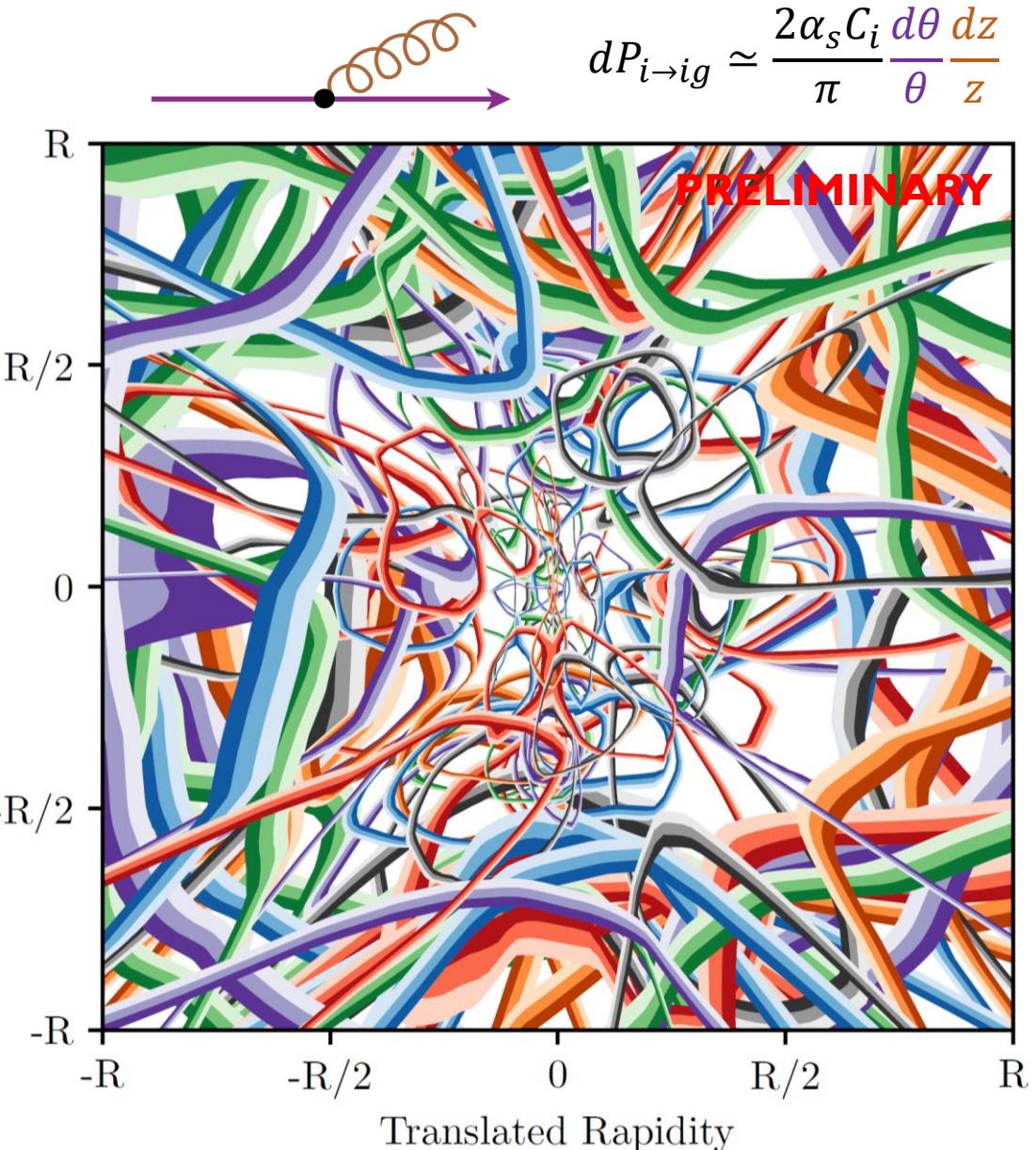
To reduce sample dependence, we train an EFN on tracks with $p_T^{\text{PFC}} > 1 \text{ GeV}$ and remove pileup.

Or high-dimensional unfolding? [See Patrick’s Talk](#)

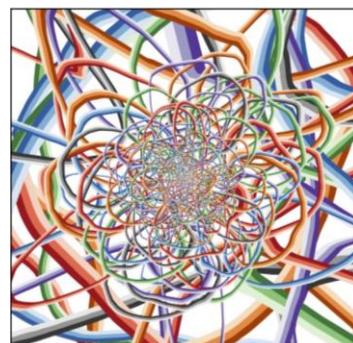
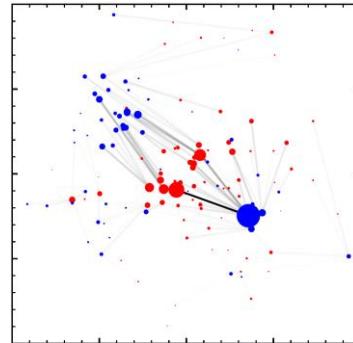
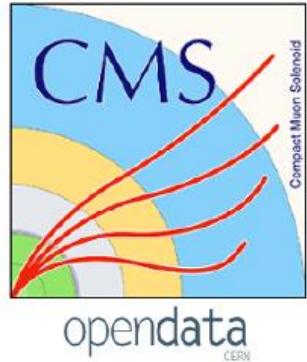
What is the model learning?



[See Patrick's Talk at ML4Jets 2018](#)



Visualizing 256 filters for EFN (weakly) trained on data



CMS Open Data

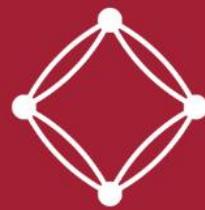
A new public dataset for jet studies

Unsupervised Learning

A metric for collider events

Supervised Learning

Training directly on collider data



EnergyFlow

Search docs

Home

Welcome to EnergyFlow

References

Copyright

Getting Started

Installation

Demos

Examples

FAQs

Release Notes

News

Documentation

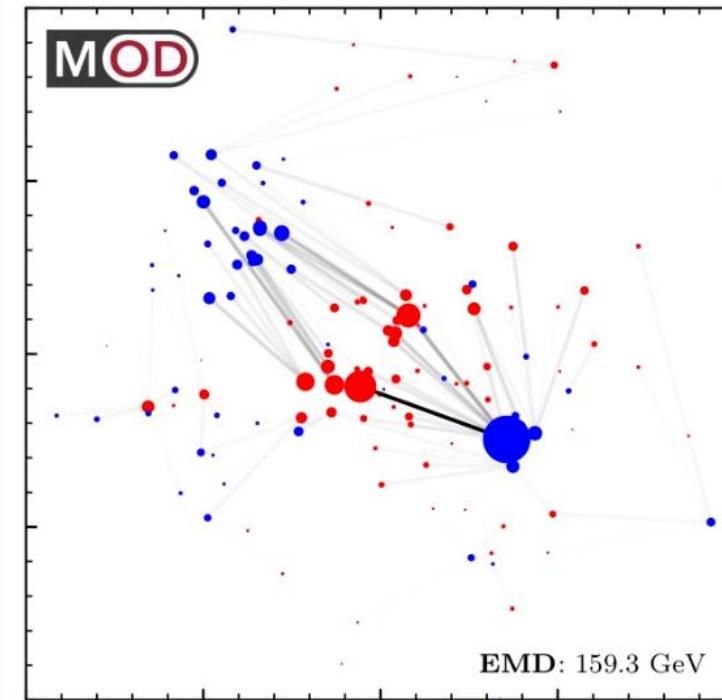
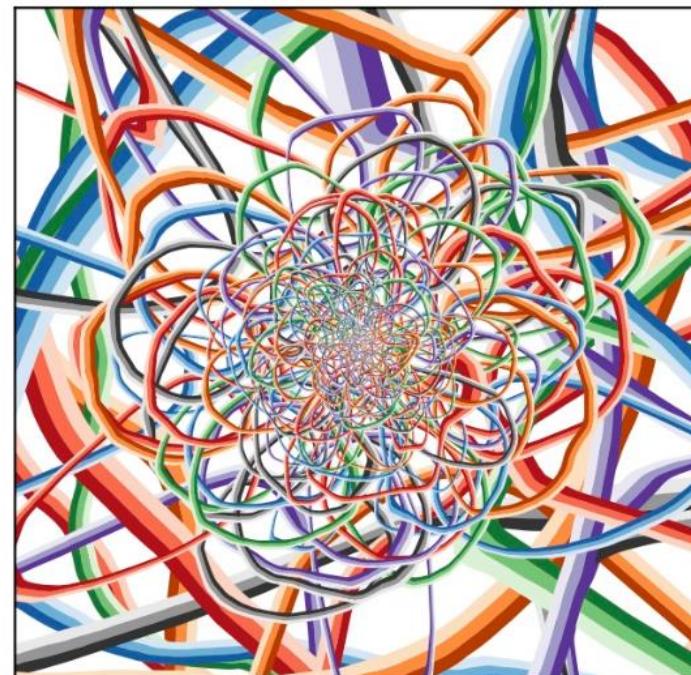
Architectures

Docs » Home

<https://energyflow.network>

pip install energyflow

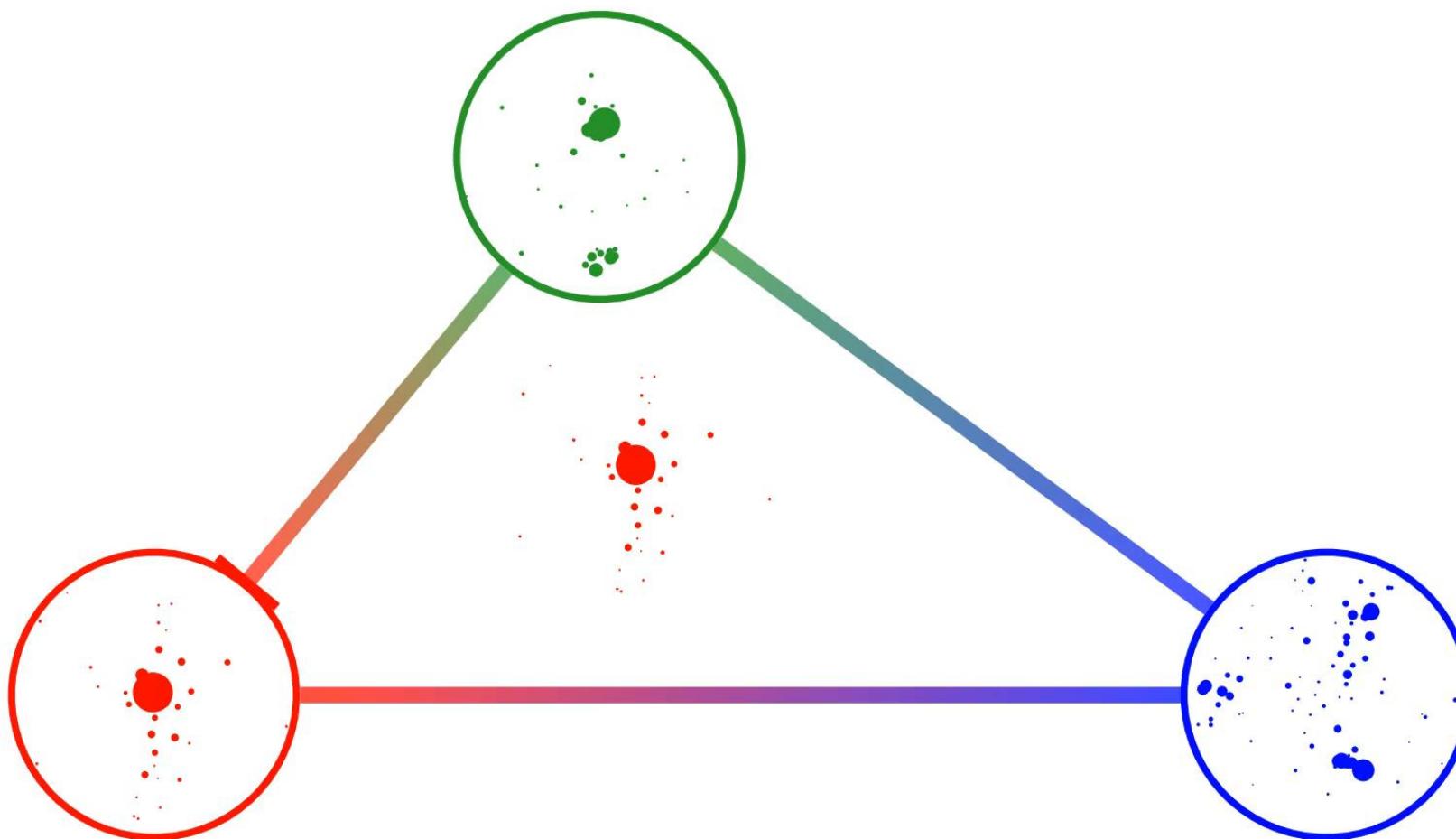
Welcome to EnergyFlow



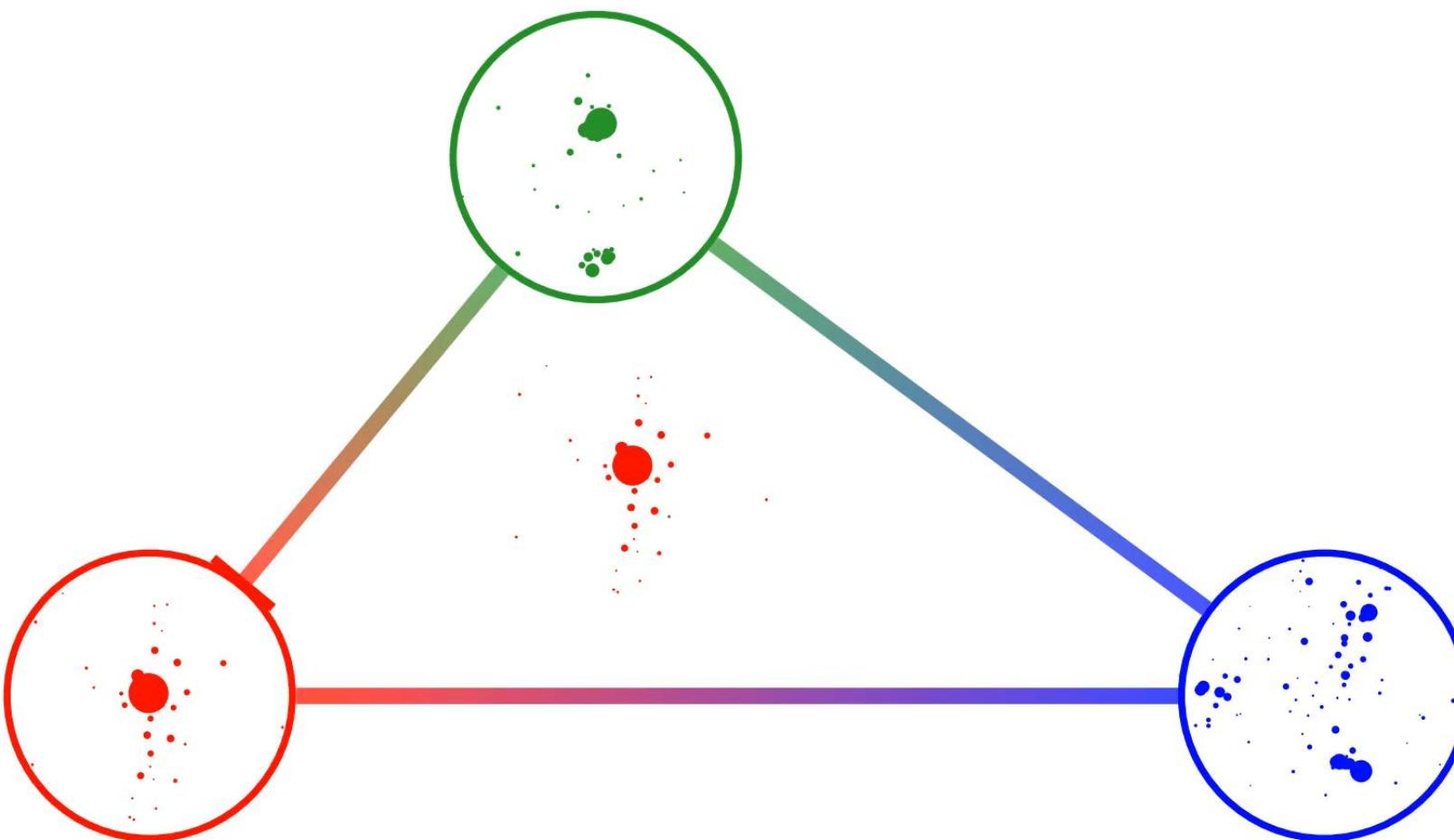
EnergyFlow is a Python package containing a suite of particle physics tools:

- **Energy Flow Polynomials:** EFPs are a collection of jet substructure observables which form a complete linear basis of IRC-safe observables. EnergyFlow provides tools to compute EFPs on

The End
Thank you!



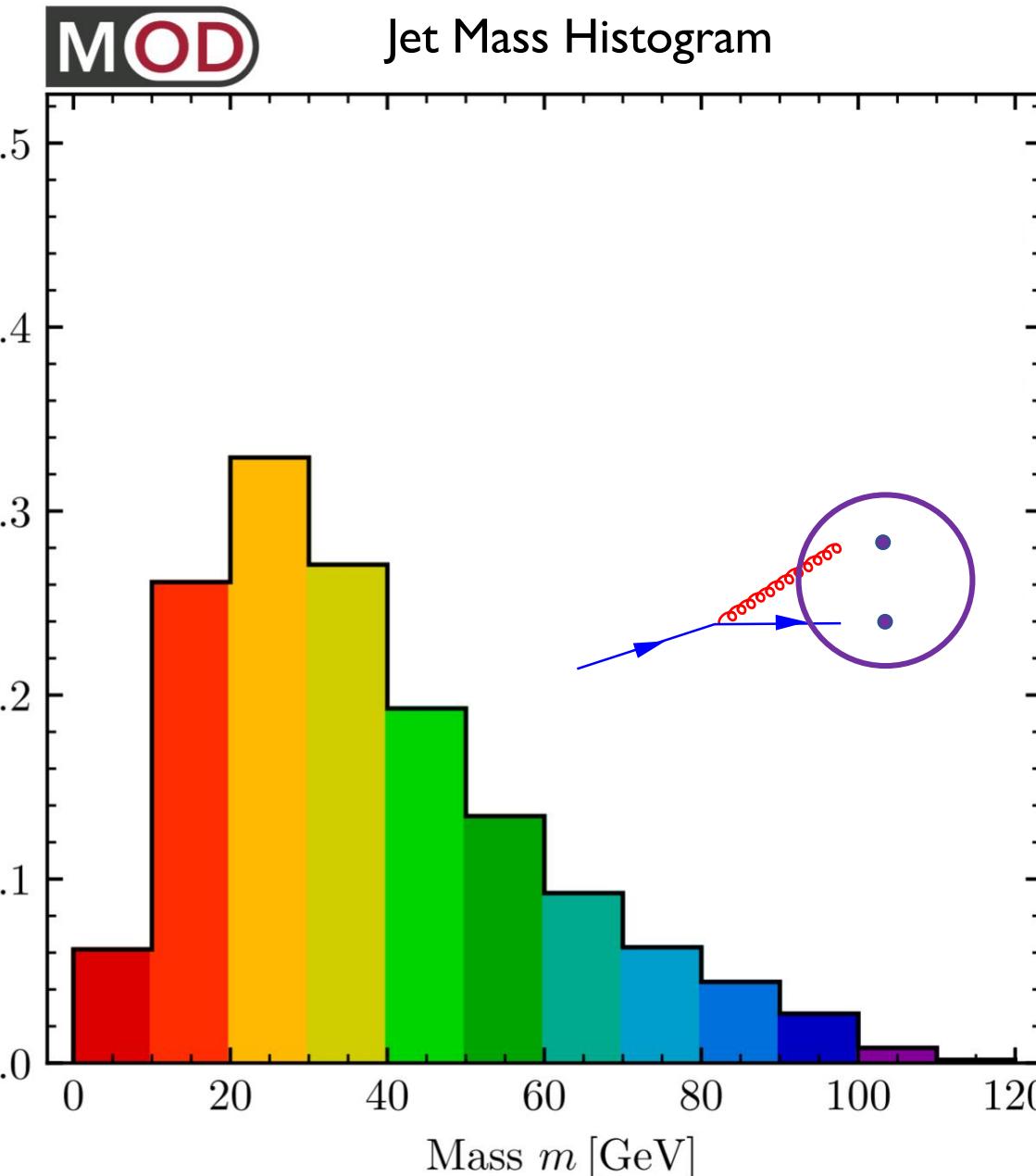
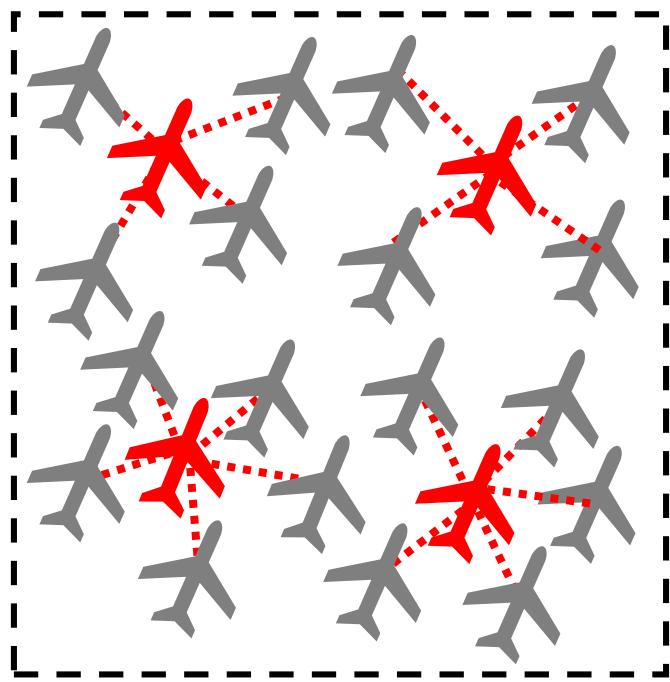
Extra Slides



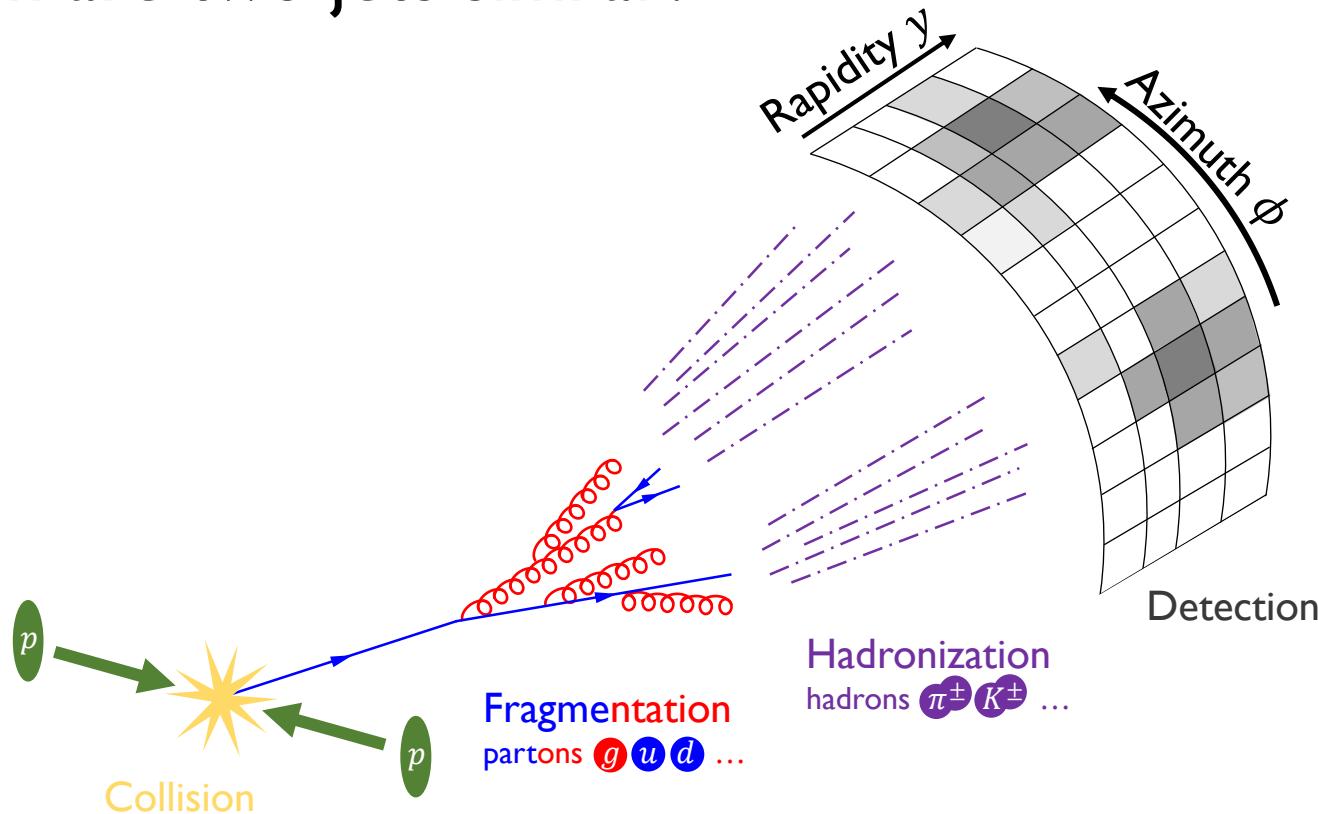
Most Representative Jets

$$\text{Jet Mass: } m = \left(\sum_{i=1}^M p_i^\mu \right)^2$$

Measures how “wide” the jet is.



When are two jets similar?

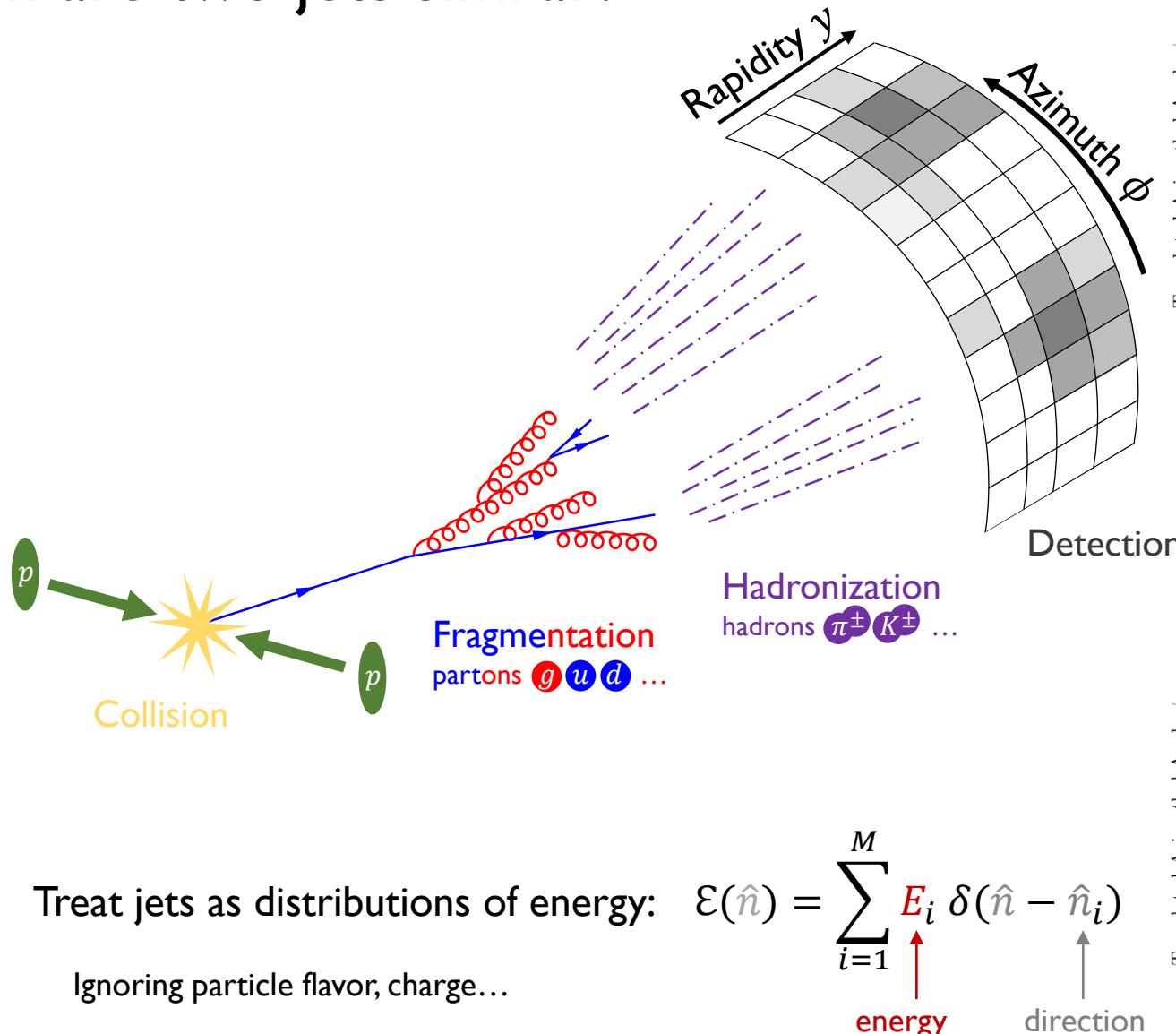


The energy flow (distribution of energy) is the information that is robust to:
fragmentation, hadronization, detector effects, ...

[\[N.A. Sveshnikov, F.V. Tkachov, 9512370\]](#)
[\[F.V. Tkachov, 9601308\]](#)
[\[P.S. Cherzor, N.A. Sveshnikov, 9710349\]](#)

Energy flow \Leftrightarrow Infrared and Collinear (IRC) Safe information

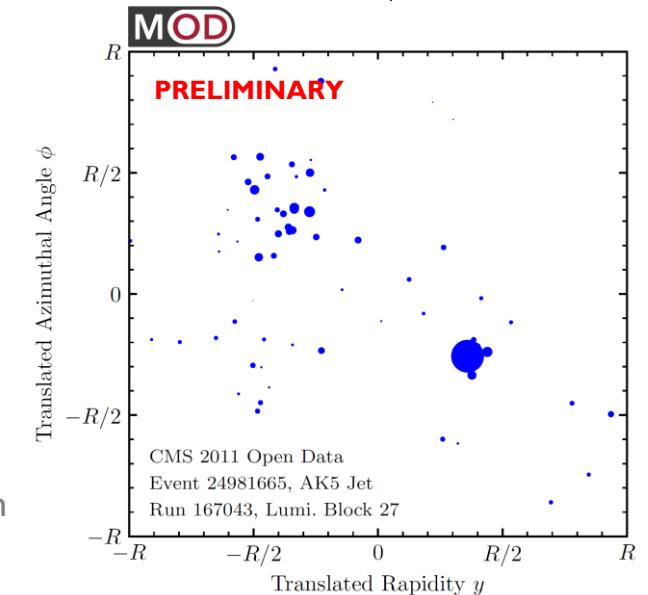
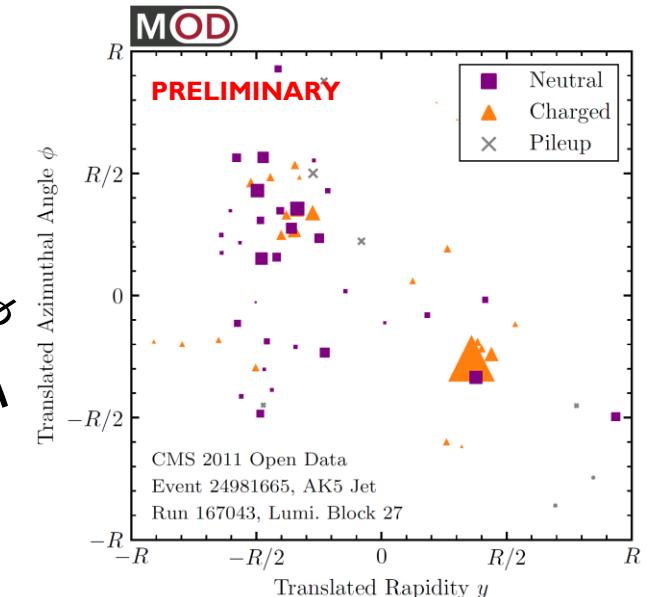
When are two jets similar?



Treat jets as distributions of energy: $\mathcal{E}(\hat{n}) = \sum_{i=1}^M E_i \delta(\hat{n} - \hat{n}_i)$

Ignoring particle flavor, charge...

↑
energy
↑
direction



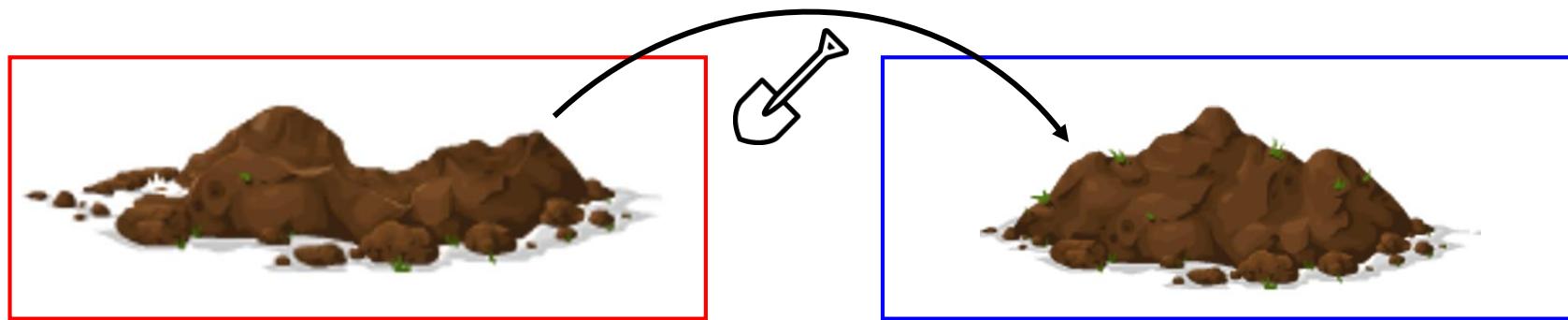
The Energy Mover's Distance

Review: *The Earth Mover's Distance*

Earth Mover's Distance: the minimum “work” (stuff \times distance) to rearrange one pile of dirt into another

[Peleg, Werman, Rom]

[Rubner, Tomasi, Guibas]



Metric on the space of (normalized) distributions: *symmetric, non-negative, triangle inequality*

Distributions are close in EMD \Leftrightarrow their expectation values are close.

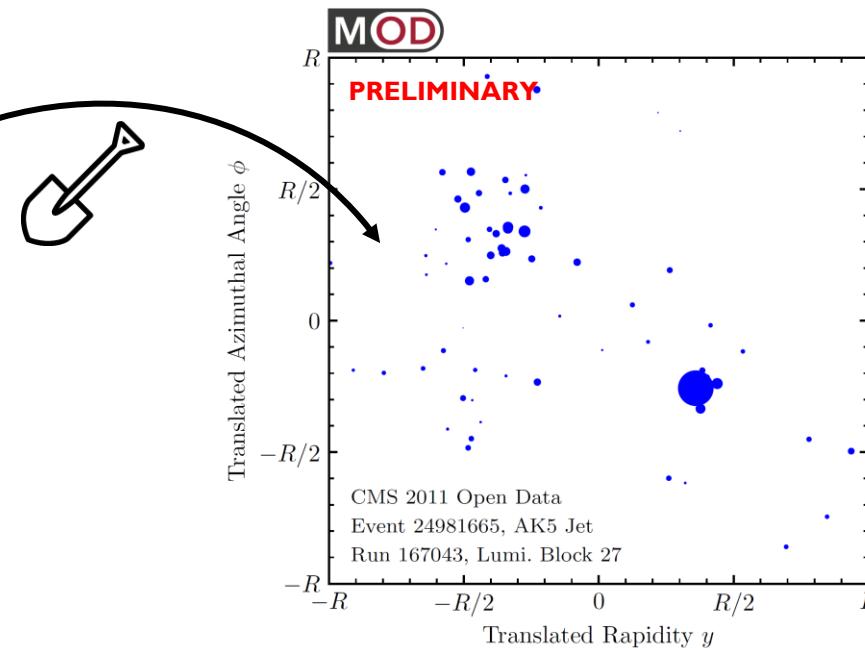
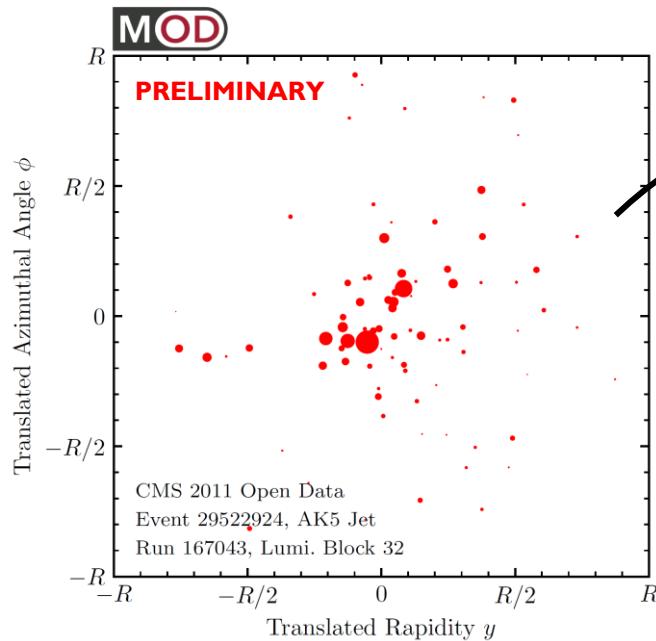
Also known as the 1-Wasserstein metric.

The Energy Mover's Distance

From Earth to Energy

Energy Mover's Distance: the minimum “work” (**energy** \times angle) to rearrange one event (pile of energy) into another

[Komiske, EMM, Thaler, 1902.02346]



$$\text{EMD}(\mathcal{E}, \mathcal{E}') = \min_{\{f\}} \sum_{i=1}^M \sum_{j=1}^{M'} f_{ij} \frac{\theta_{ij}}{R} + \left| \sum_{i=1}^M E_i - \sum_{j=1}^{M'} E'_j \right|$$

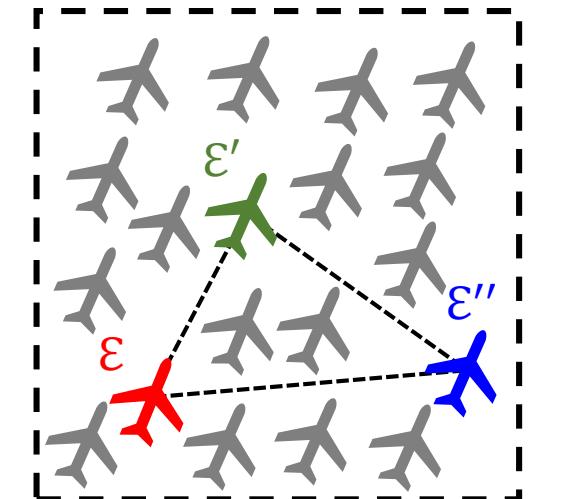
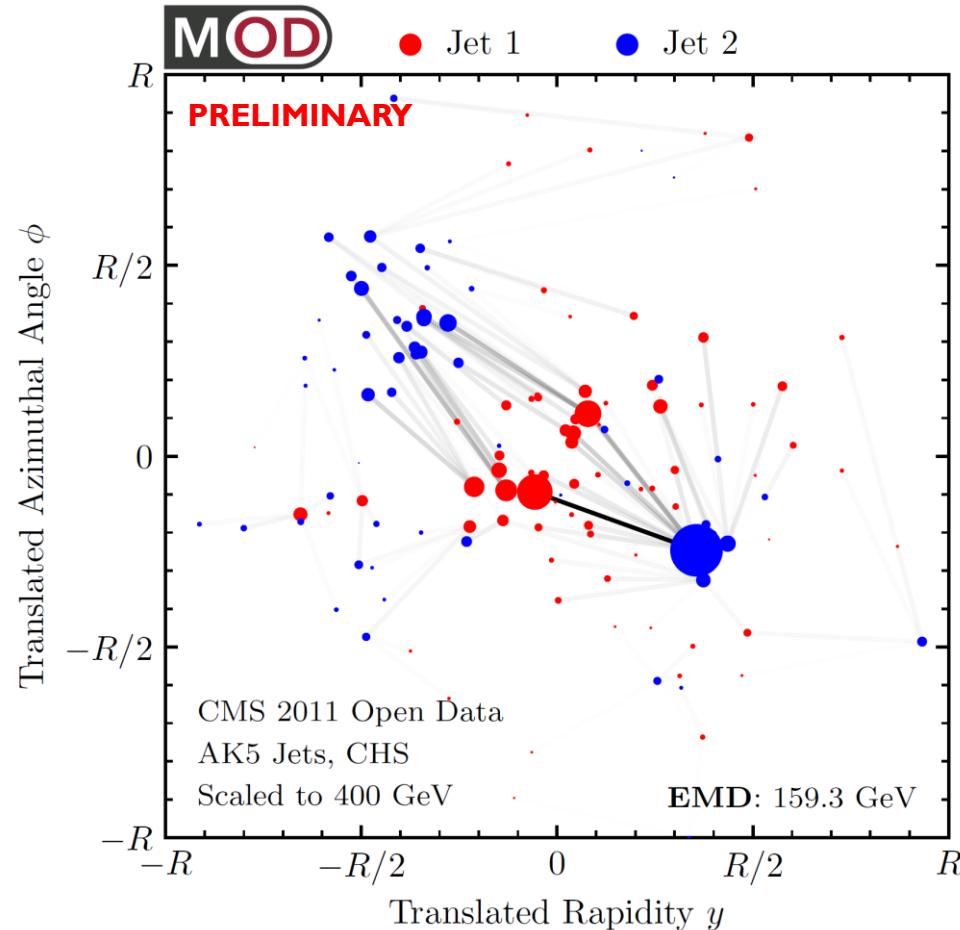
Difference in radiation pattern Difference in total energy

The Energy Mover's Distance

From Earth to Energy

Energy Mover's Distance: the minimum “work” (**energy** \times angle) to rearrange one event (pile of energy) into another

[Komiske, EMM, Thaler, 1902.02346]



EMD has dimensions of energy

True metric as long as $R \geq \frac{1}{2}\theta_{\max}$
 $R \geq$ the jet radius, for conical jets

Solvable via Optimal Transport problem.

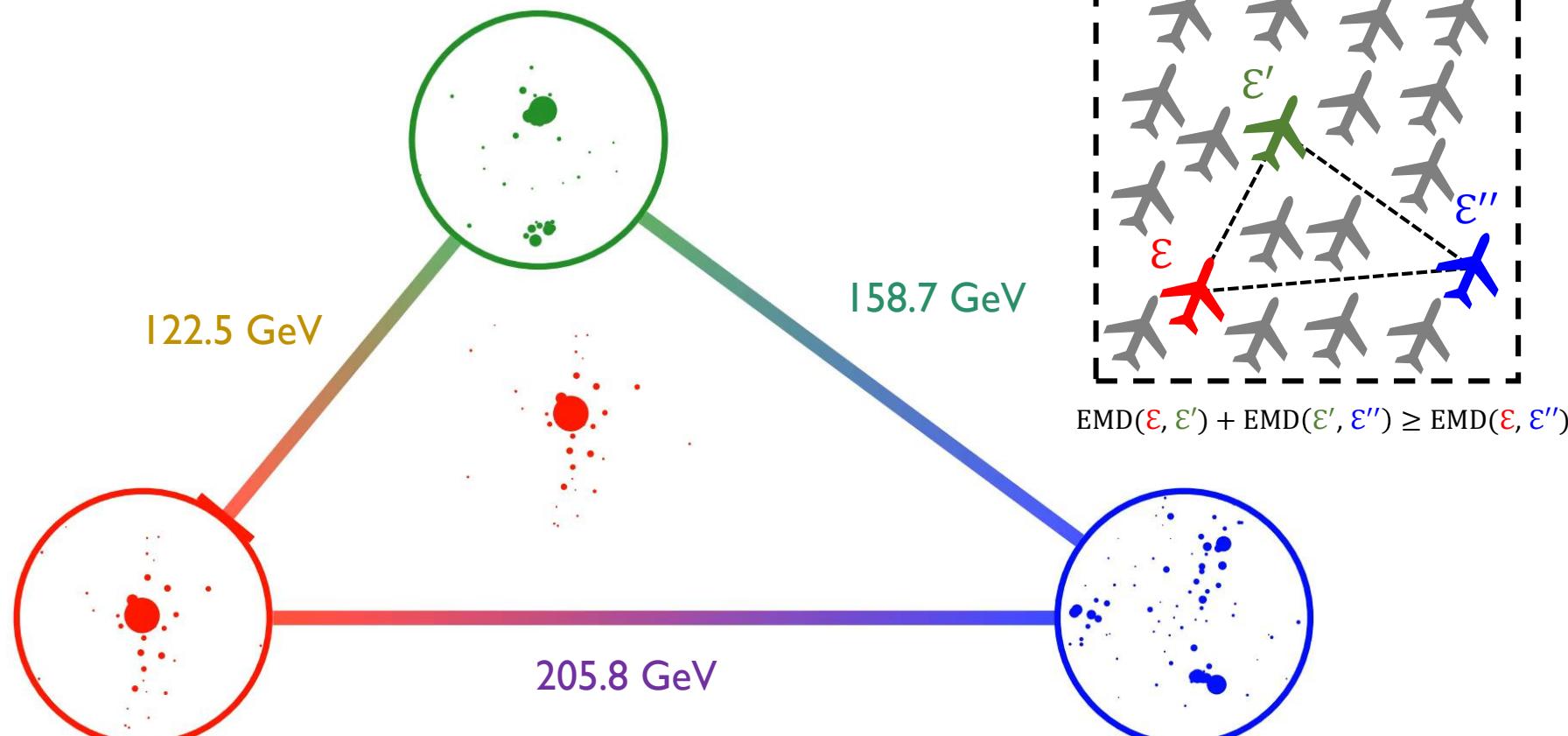
~1ms to compute EMD for two jets with 100 particles.

The Energy Mover's Distance

From Earth to Energy

Energy Mover's Distance: the minimum “work” (**energy** \times angle) to rearrange one event (pile of energy) into another

[Komiske, EMM, Thaler, 1902.02346]



<https://energyflow.network>

Energy Moving and IRC Safety

Events close in EMD are close in any infrared and collinear safe observable!

Additive IRC-safe observables:

Energy Mover's Distance

$$\text{EMD}(\mathcal{E}, \mathcal{E}') \geq \frac{1}{RL} |\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')|$$

Difference in observable values

“Lipschitz constant” of Φ
i.e. bound on its derivative

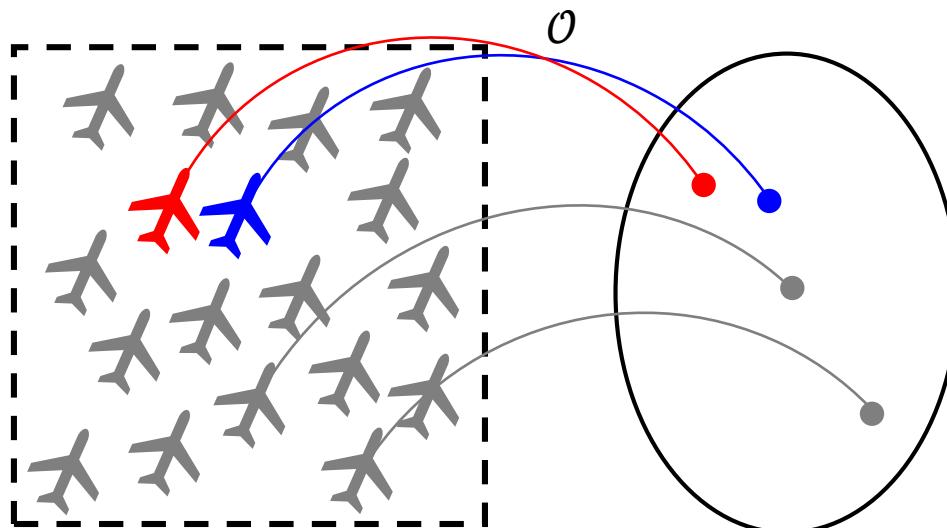
$$\mathcal{O}(\mathcal{E}) = \sum_{i=1}^M \mathbf{E}_i \Phi(\hat{n}_i)$$

e.g. $\beta \geq 1$ jet angularities:

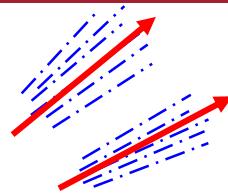
[\[Berger, Kucs, Sterman, 0303051\]](#)

[\[Larkoski, Thaler, Waalewijn, 1408.3122\]](#)

$$|\lambda^{(\beta)}(\mathcal{E}) - \lambda^{(\beta)}(\mathcal{E}')| \leq \beta \text{EMD}(\mathcal{E}, \mathcal{E}')$$



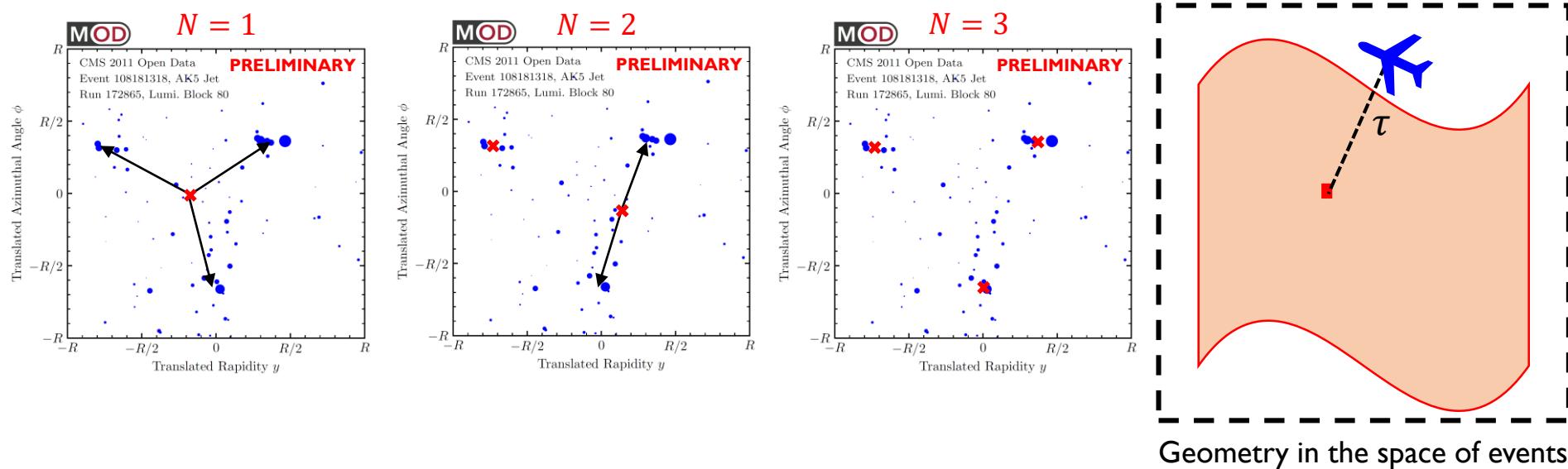
Old Observables in a New Language



N -subjettiness is the EMD between the event and the closest N -particle event.

$$\tau_N^{(\beta)}(\mathcal{E}) = \min_{N \text{ axes}} \sum_{i=1}^M E_i \min\{\theta_{1,i}^\beta, \theta_{2,i}^\beta, \dots, \theta_{N,i}^\beta\} \longrightarrow \tau_N(\mathcal{E}) = \min_{|\mathcal{E}'|=N} \text{EMD}(\mathcal{E}, \mathcal{E}').$$

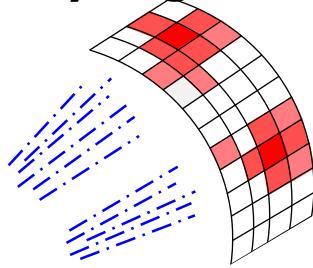
$\beta \geq 1$ is p-Wasserstein distance with $p = \beta$.



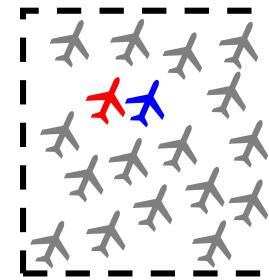
Thrust is the EMD between the event and two back-to-back particles.

$$t(\mathcal{E}) = E - \max_{\hat{n}} \sum_i |\vec{p}_i \cdot \hat{n}| \longrightarrow t(\mathcal{E}) = \min_{|\mathcal{E}'|=2} \text{EMD}(\mathcal{E}, \mathcal{E}') \quad \text{with } \theta_{ij} = \hat{n}_i \cdot \hat{n}_j, \quad \hat{n} = \vec{p}/E$$

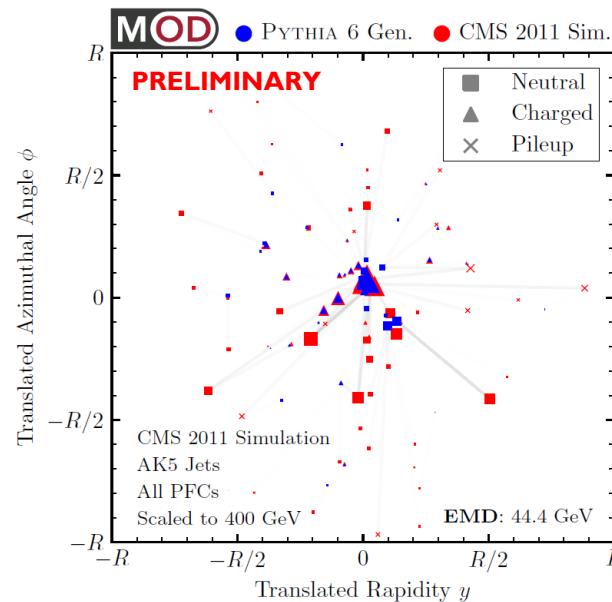
Quantifying Pileup and Detector Effects with EMD



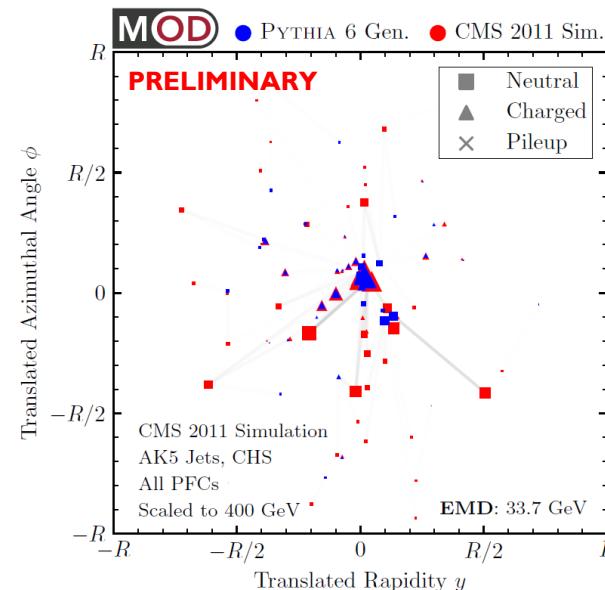
Gen./Sim. EMD universally quantifies pileup and detector effects.



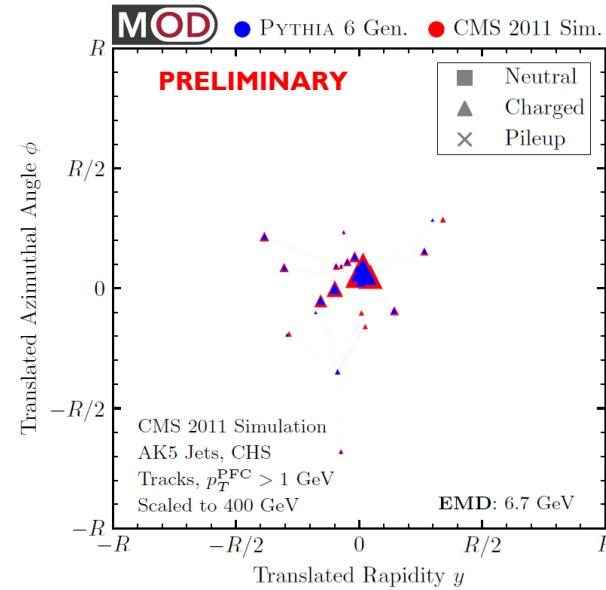
Gen./Sim. EMD: 44.4 GeV



Gen./Sim. EMD: 33.7 GeV



Gen./Sim. EMD: 6.7 GeV



+ charged hadron subtraction

+ Tracks only, $p_T^{\text{PFC}} > 1 \text{ GeV}$ cut

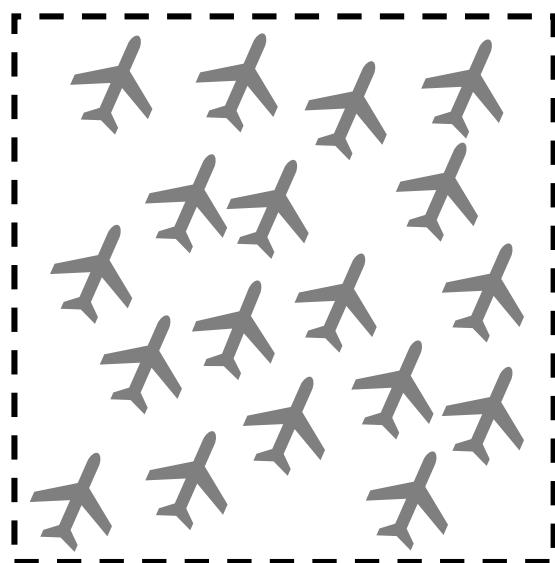
See extra slides for histograms. Can also quantify hadronization effects this way.

Exploring the Space of Jets: Visualizing the Manifold

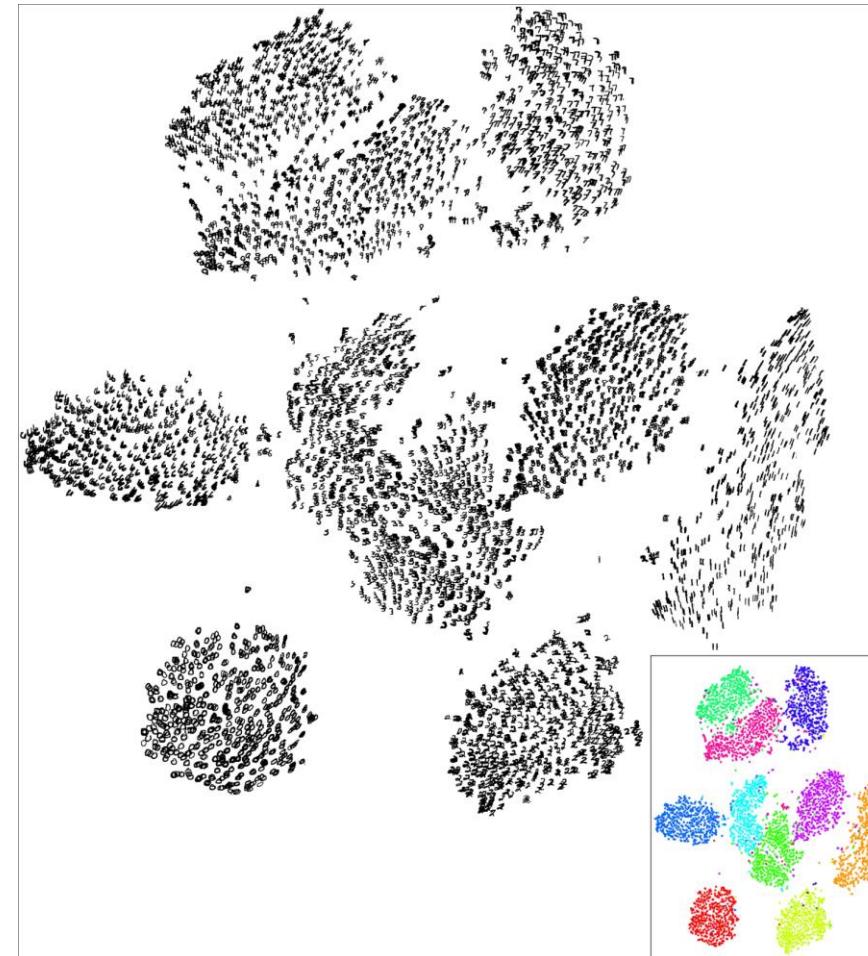
Visualize the space of events with t-Distributed Stochastic Neighbor Embedding (t-SNE).

[\[L. van der Maaten, G. Hinton\]](#)

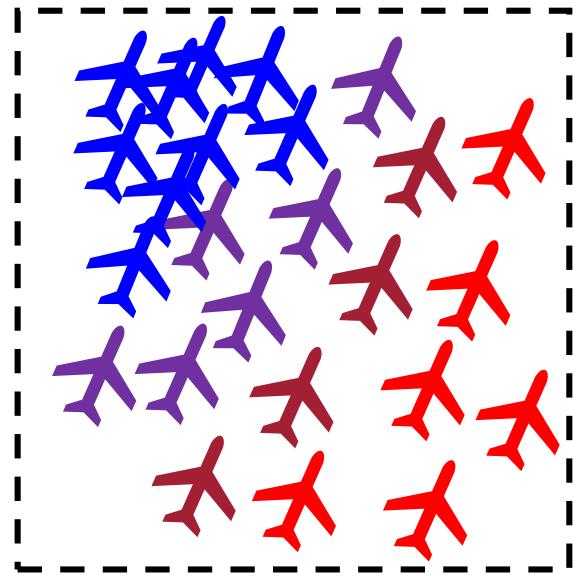
Finds an embedding into a low-dimensional manifold that respects distances.



What does the space
of jets look like?

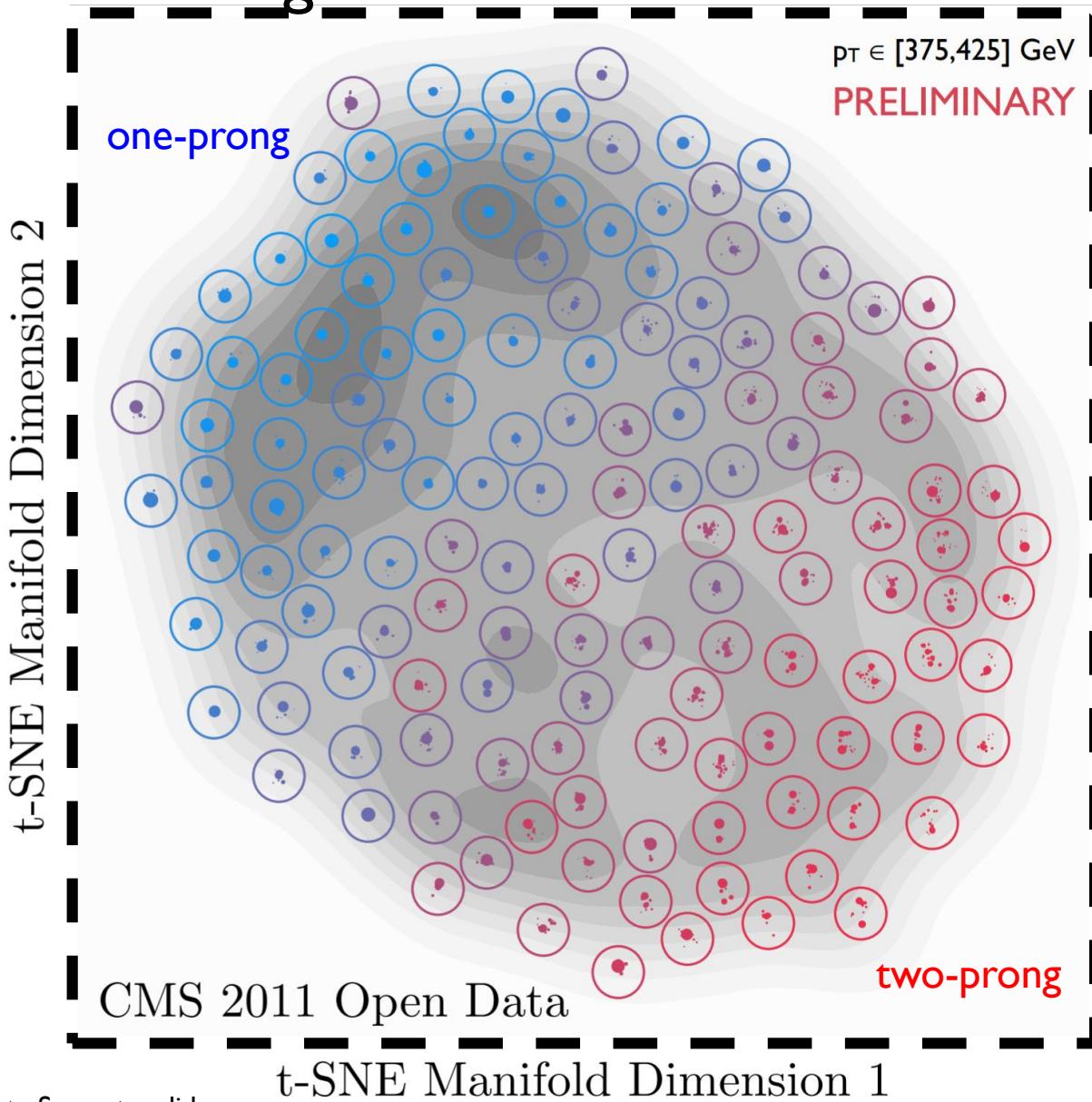


Exploring the Space of Jets: Visualizing the Manifold



What does the space
of jets look like?

Quantify (and calculate?) dimension of the space of jets. See extra slides.



Exploring the Space of Jets: Correlation Dimension

VOLUME 50, NUMBER 5

PHYSICAL REVIEW LETTERS

31 JANUARY 1983

Characterization of Strange Attractors

Peter Grassberger^(a) and Itamar Procaccia

Chemical Physics Department, Weizmann Institute of Science, Rehovot 76100, Israel

(Received 7 September 1982)

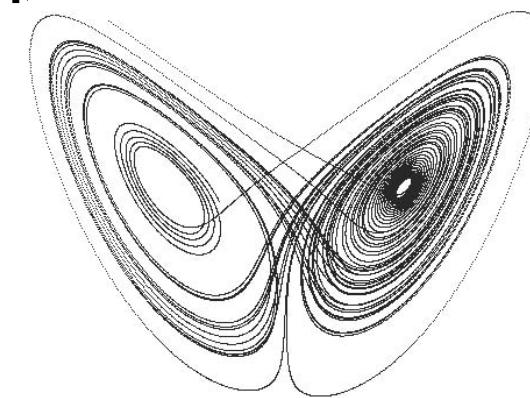
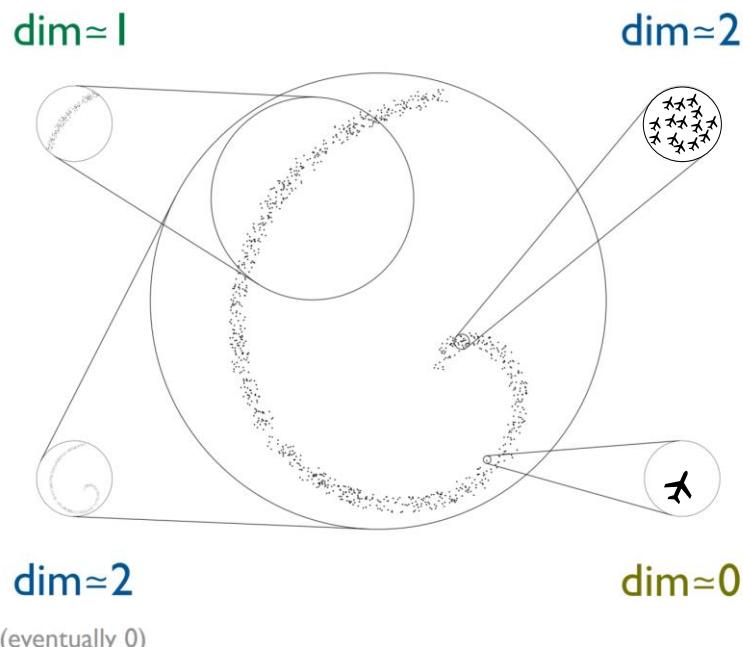
A new measure of strange attractors is introduced which offers a practical algorithm to determine their character from the time series of a single observable. The relation of this new measure to fractal dimension and information-theoretic entropy is discussed.

Intuition:

$$N_{\text{neighboring}}(r) \propto r^{\dim}$$

points

$$\longrightarrow \dim(r) = r \frac{\partial}{\partial r} \ln N_{\text{neighbors}}(r)$$



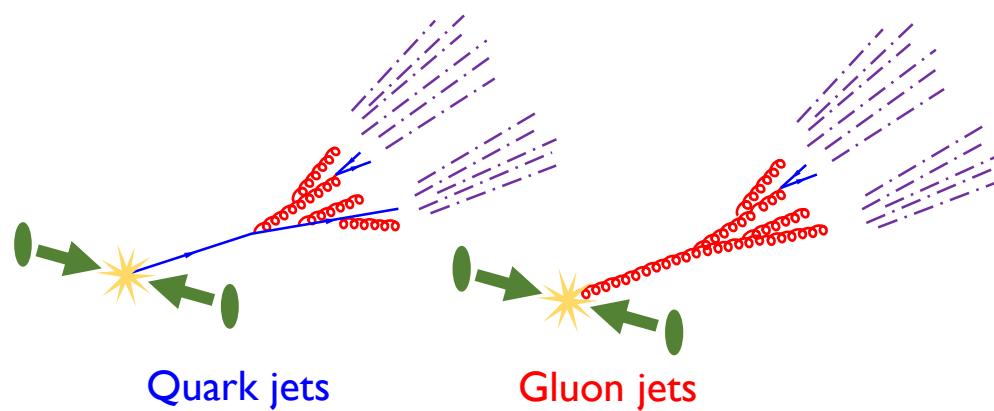
Correlation dimension:

$$\dim(Q) = Q \frac{\partial}{\partial Q} \ln \sum_{i=1}^N \sum_{j=1}^N \Theta[\text{EMD}(\varepsilon_i, \varepsilon_j) < Q]$$

Energy scale Q
dependence

Count neighbors in
ball of radius Q

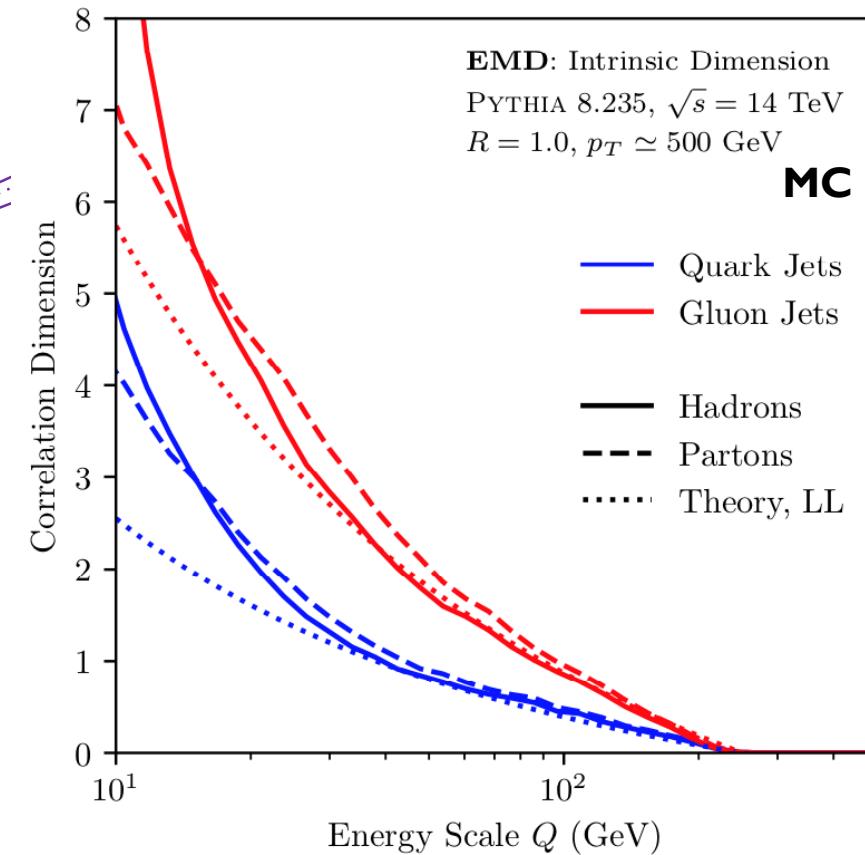
Exploring the Space of Jets: Correlation Dimension



$$\text{At LL: } \dim_{q/g}(Q) = -\frac{8\alpha_s C_{q/g}}{\pi} \ln \frac{Q}{p_T/2}$$

+ 1-loop running of α_s

$$C_q = C_F = \frac{4}{3}$$
$$C_g = C_A = 3$$



EMD: Intrinsic Dimension
PYTHIA 8.235, $\sqrt{s} = 14$ TeV
 $R = 1.0, p_T \simeq 500$ GeV

MC

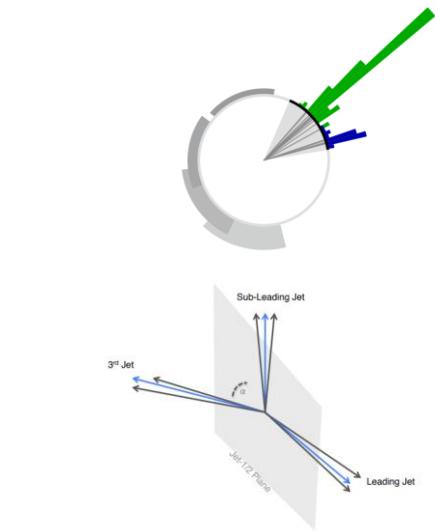
- Quark Jets
- Gluon Jets
- Hadrons
- - Partons
- Theory, LL

Dimension blows up at low energies.

Jets are “more than fractal”

CMS Open Data

Many exciting physics applications with the CMS Open Data already.



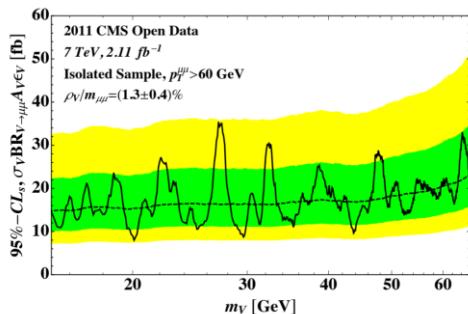
Exposing the QCD splitting function

[\[Tripathee, Xue, Larkoski, Marzani, Thaler, 1704.05842\]](#)

[\[Larkoski, Marzani, Thaler, Tripathee, Xue, 1704.05066\]](#)

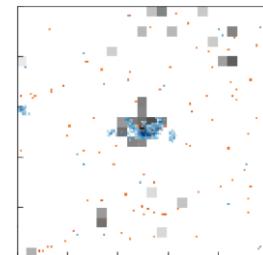
Looking for parity violation in jets

[\[Lester, Schott, 1904.11195\]](#)



Searching for dimuon resonances

[\[Cesarotti, Soreq, Strassler, Thaler, Xue, 1902.04222\]](#)



Analyzing collision data with deep learning techniques

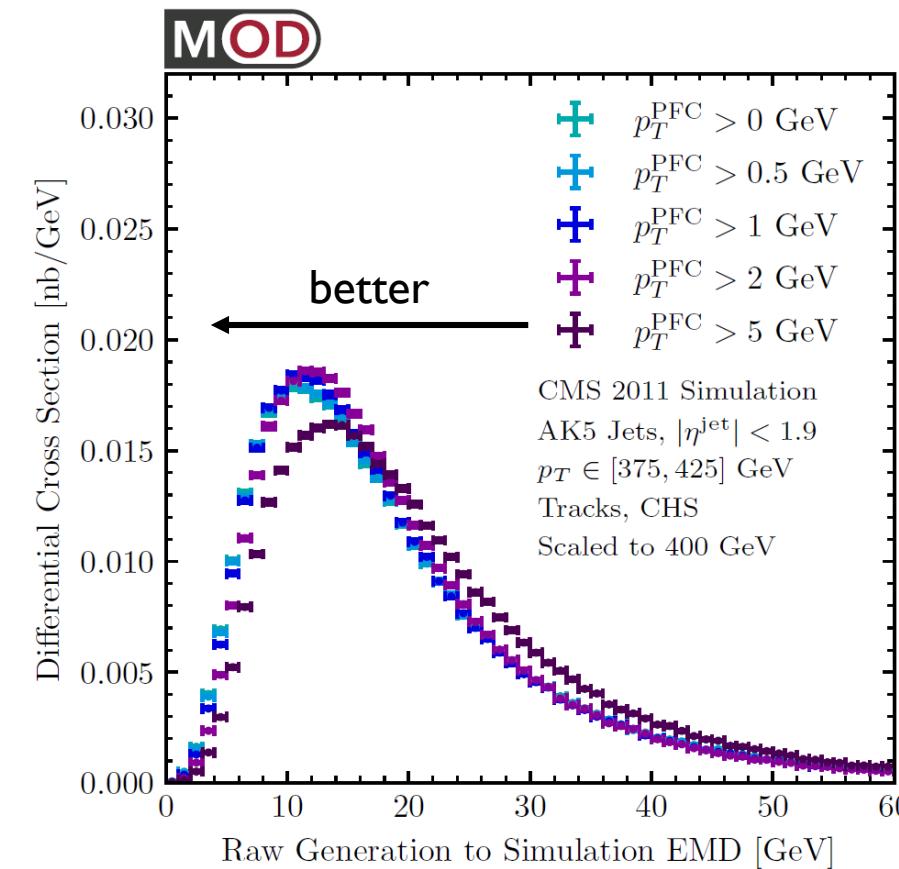
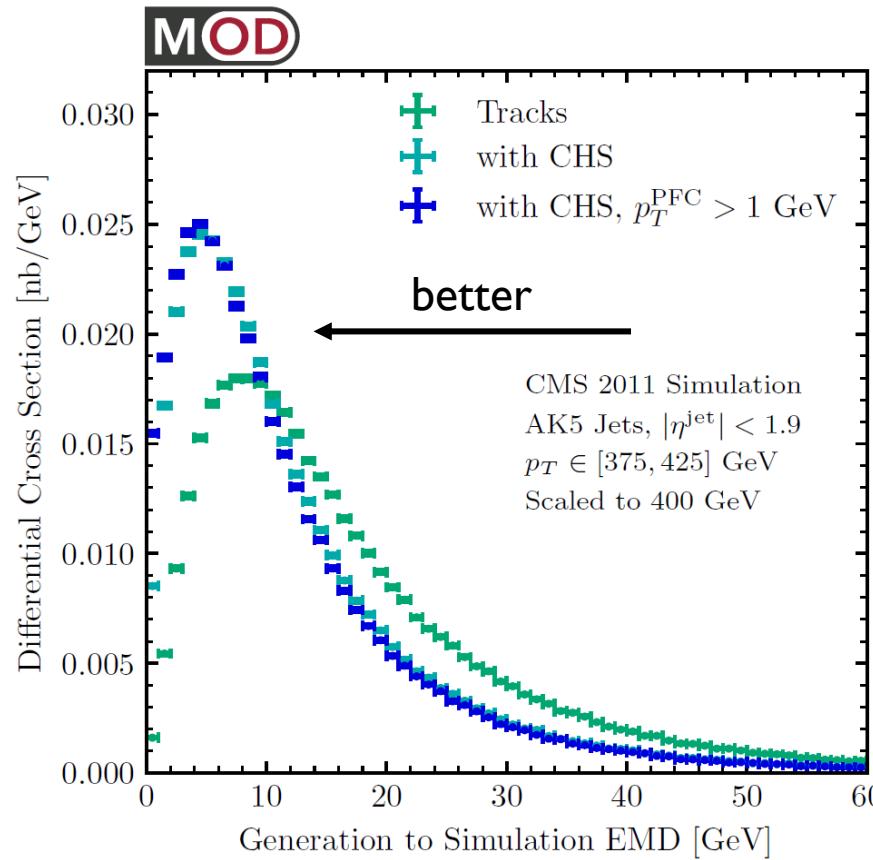
[\[Madrazo, Cacha, Iglesias, de Lucas, 1708.07034\]](#)

[\[Andrews, Paulini, Gleyzer, Poczos, 1807.11916\]](#)

[\[Andrews, et al., 1902.08276\]](#)

Quantifying Pileup and Detector Effects with EMD

Gen./Sim. EMD universally quantifies pileup mitigation and detector effects.



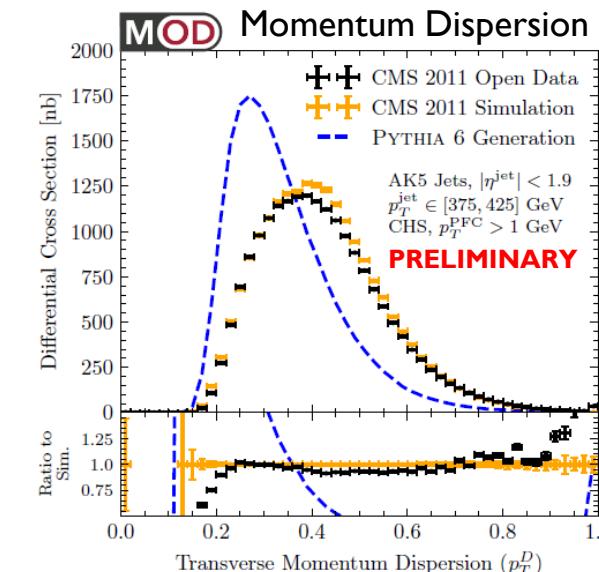
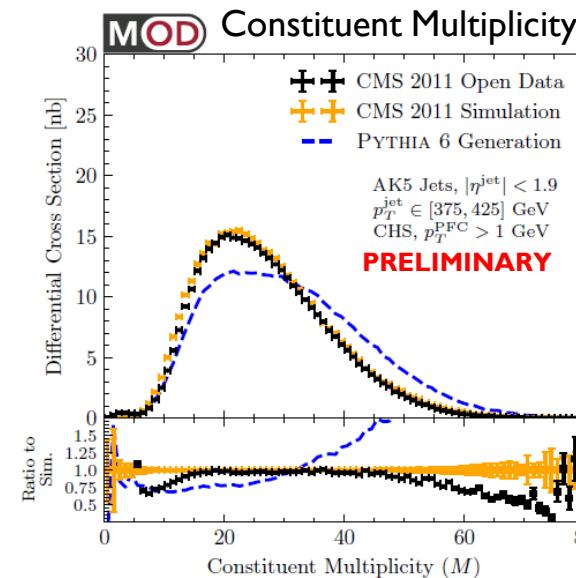
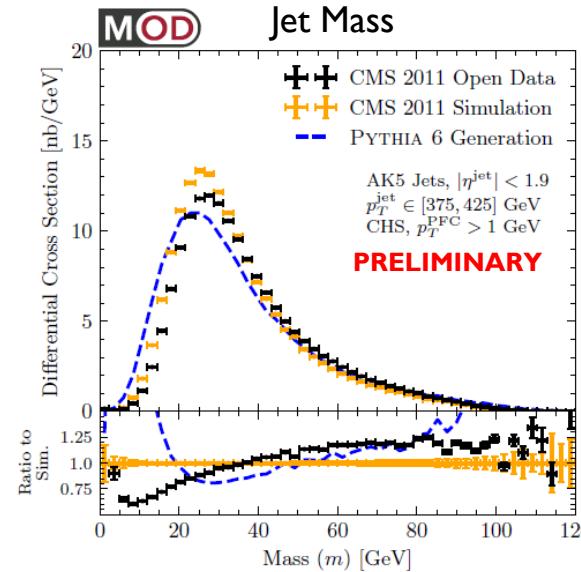
Jet Substructure Observables

Study jet substructure at truth and detector level.

$$m^2 = \left(\sum_{i \in \text{Jet}} p_i^\mu \right)^2$$

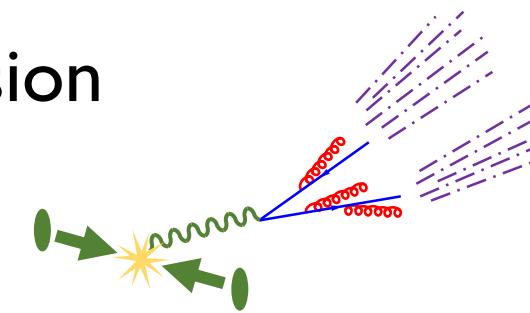
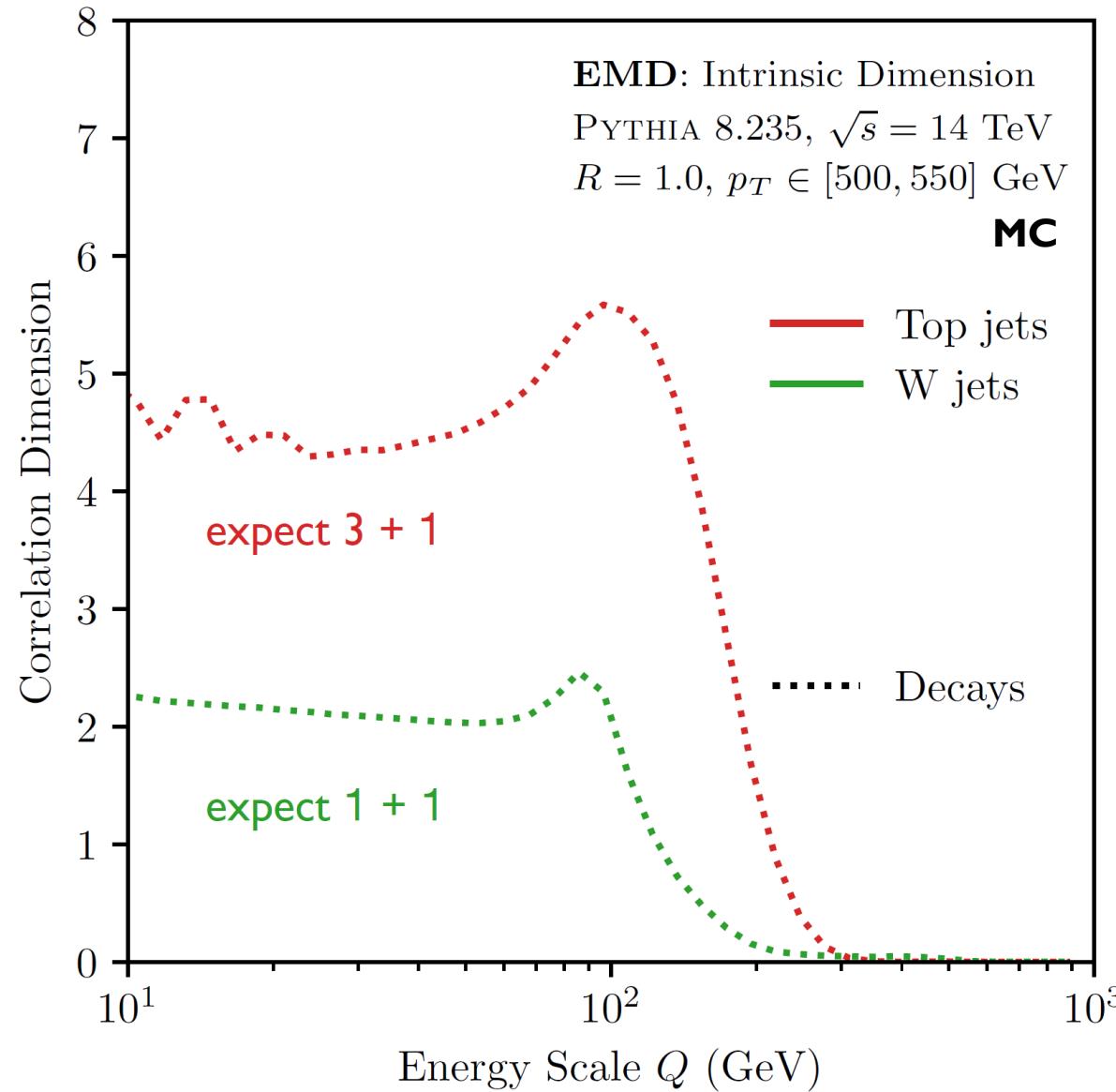
$$M = \sum_{i \in \text{Jet}} 1$$

$$p_T^D = \frac{\sum_{i \in \text{Jet}} p_{T,i}^2}{\left(\sum_{i \in \text{Jet}} p_{T,i} \right)^2}$$



Similar to: [\[Larkoski, Marzani, Thaler, Tripathee, Xue, 1704.05066\]](#)

Exploring the Space of Jets: Correlation Dimension



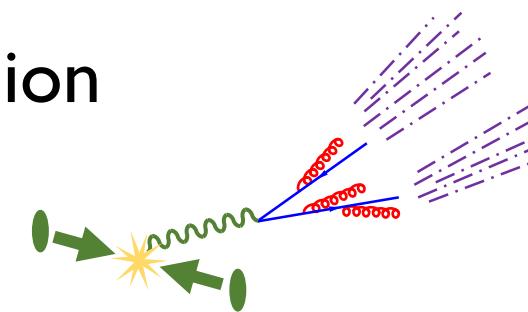
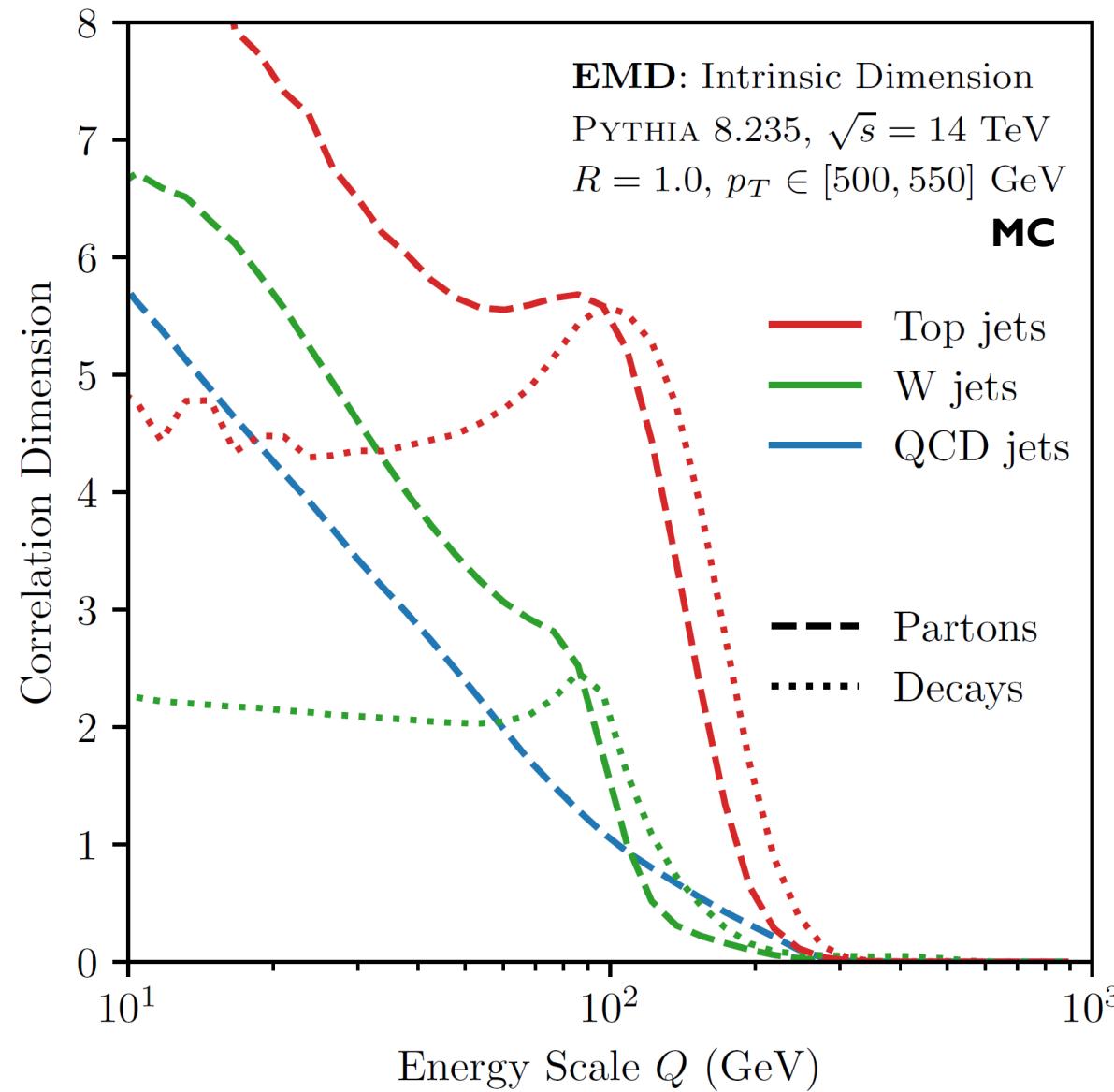
QCD jets are simplest.

W jets are more complicated.

Top jets are most complex.

“Decays” have \sim constant dimension.

Exploring the Space of Jets: Correlation Dimension



QCD jets are simplest.

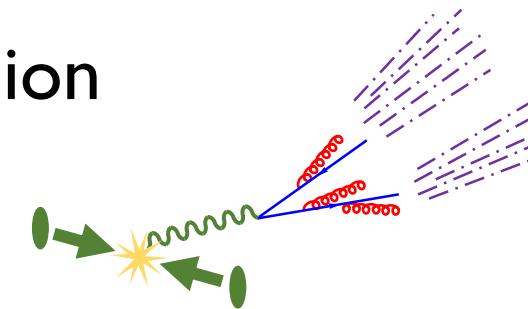
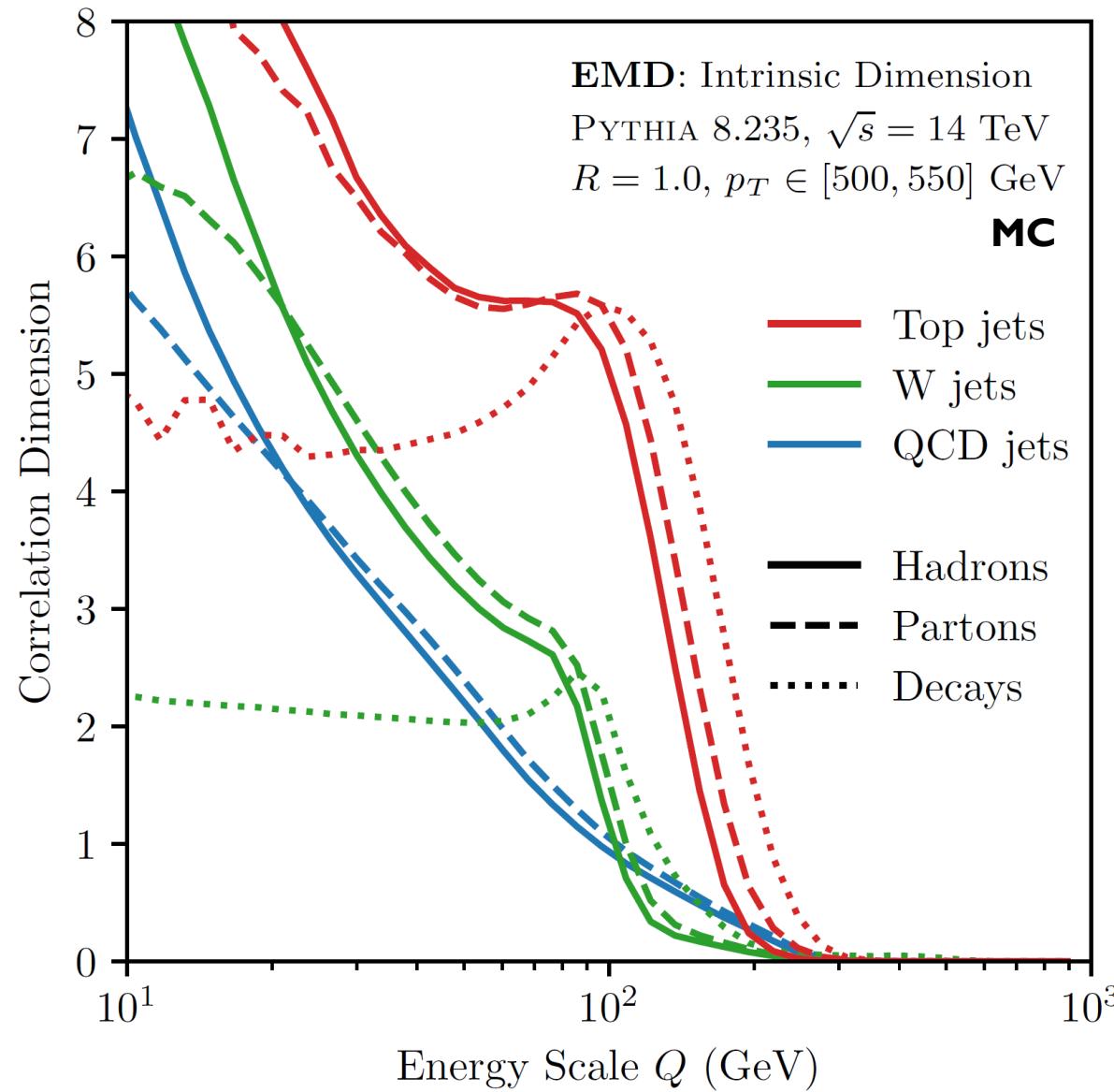
W jets are more complicated.

Top jets are most complex.

“Decays” have \sim constant dimension.

Fragmentation becomes more complex at lower energy scales.

Exploring the Space of Jets: Correlation Dimension



QCD jets are simplest.

W jets are more complicated.

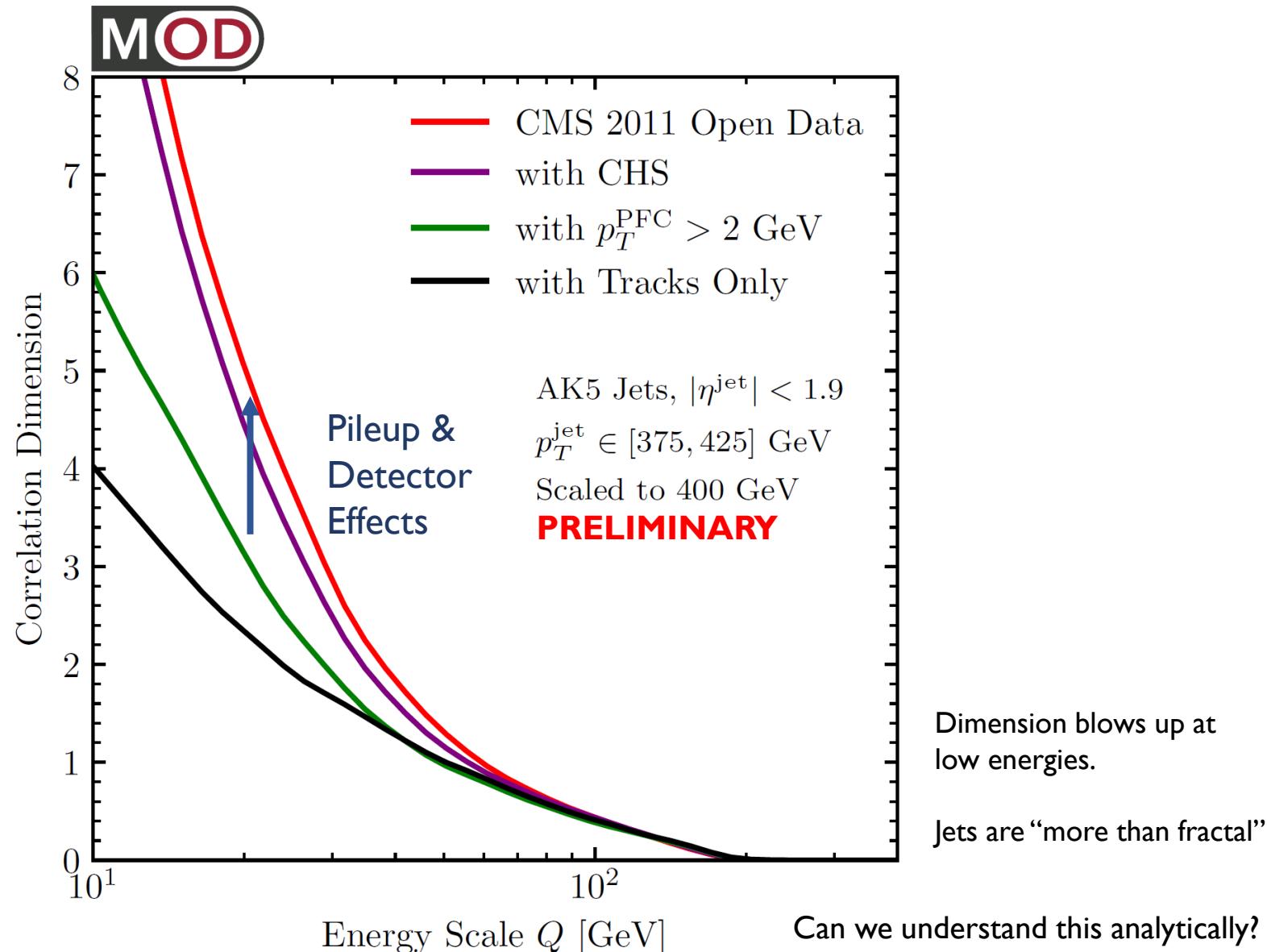
Top jets are most complex.

“Decays” have \sim constant dimension.

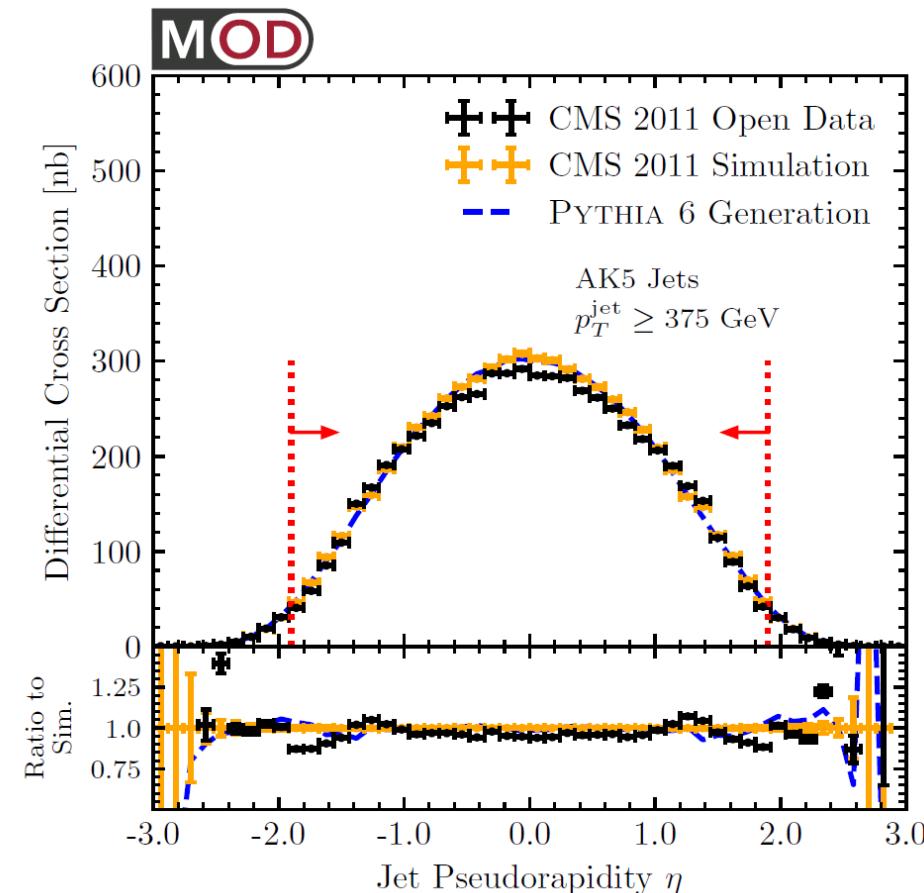
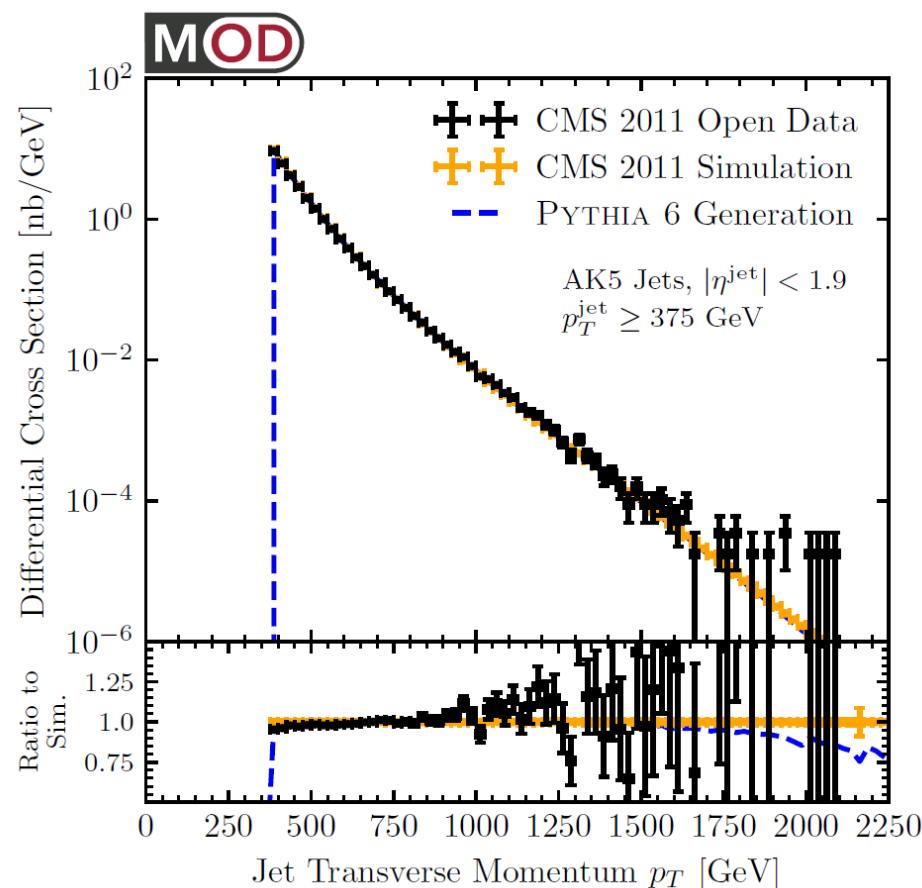
Fragmentation becomes more complex at lower energy scales.

Hadronization becomes relevant at scales around 20 GeV.

Exploring the Space of Jets: Correlation Dimension

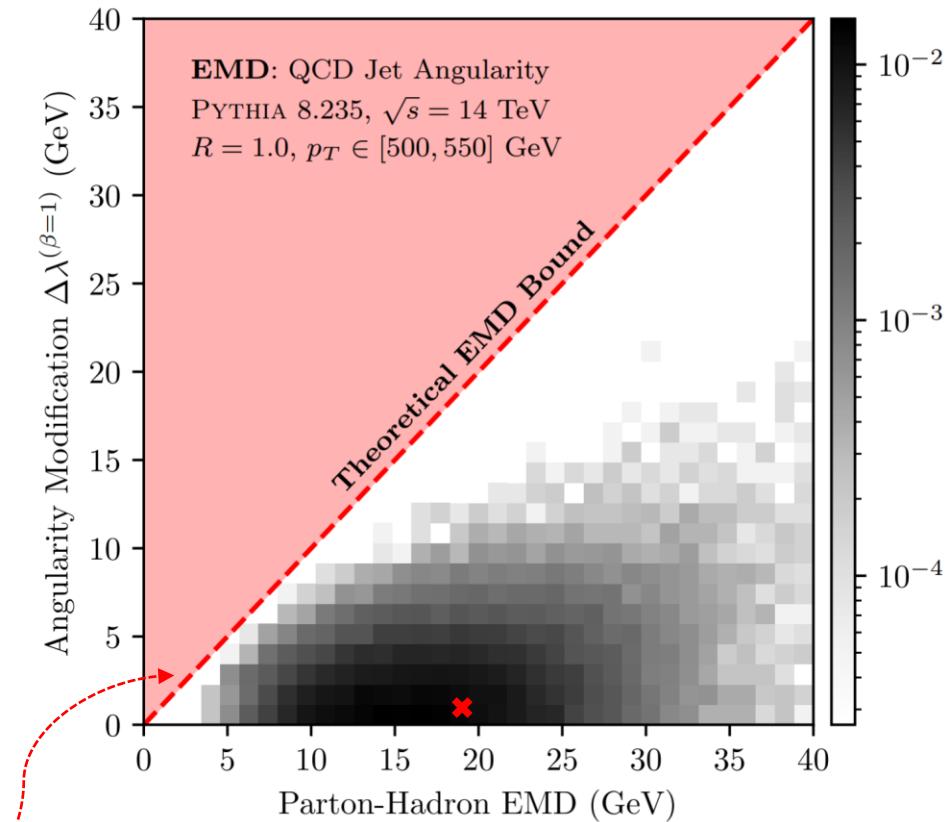
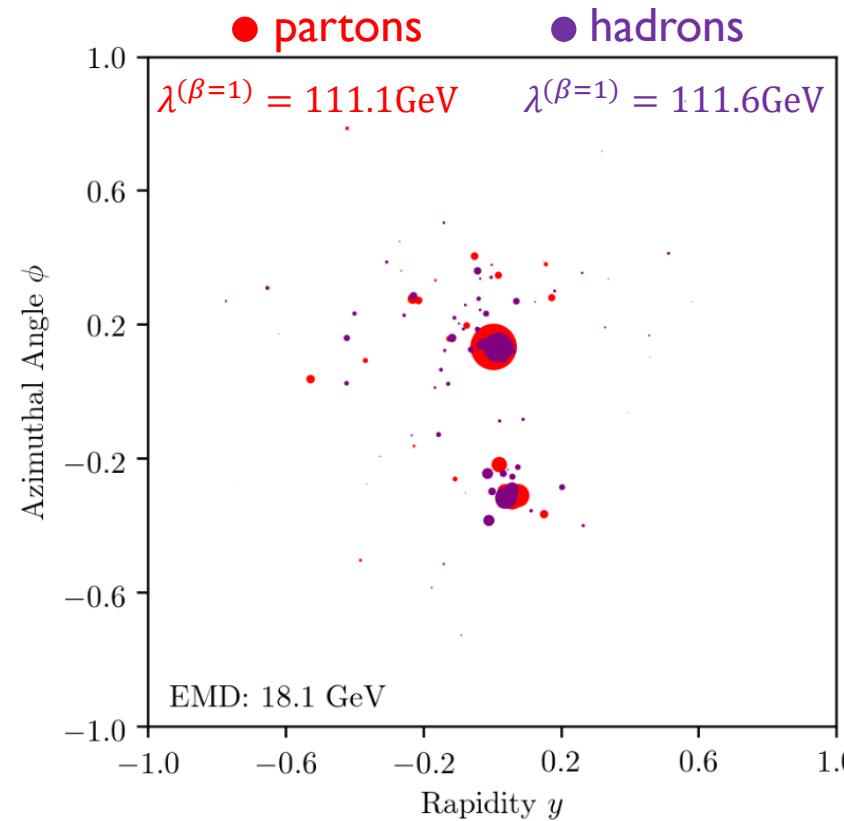


Jet Kinematic Distributions



Quantifying event modifications: Hadronization

$$\lambda^{(\beta=1)} = \sum_{i=1}^M E_i \theta_i$$



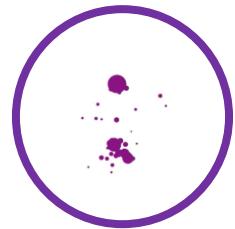
$$\varepsilon = \varepsilon_{\text{partons}}$$

$$\varepsilon' = \varepsilon_{\text{hadrons}}$$

$$|\lambda^{(\beta=1)}(\varepsilon) - \lambda^{(\beta=1)}(\varepsilon')| \leq \text{EMD}(\varepsilon, \varepsilon')$$

Exploring the Space of Events: Jet Classification

Classify W jets vs. QCD jets



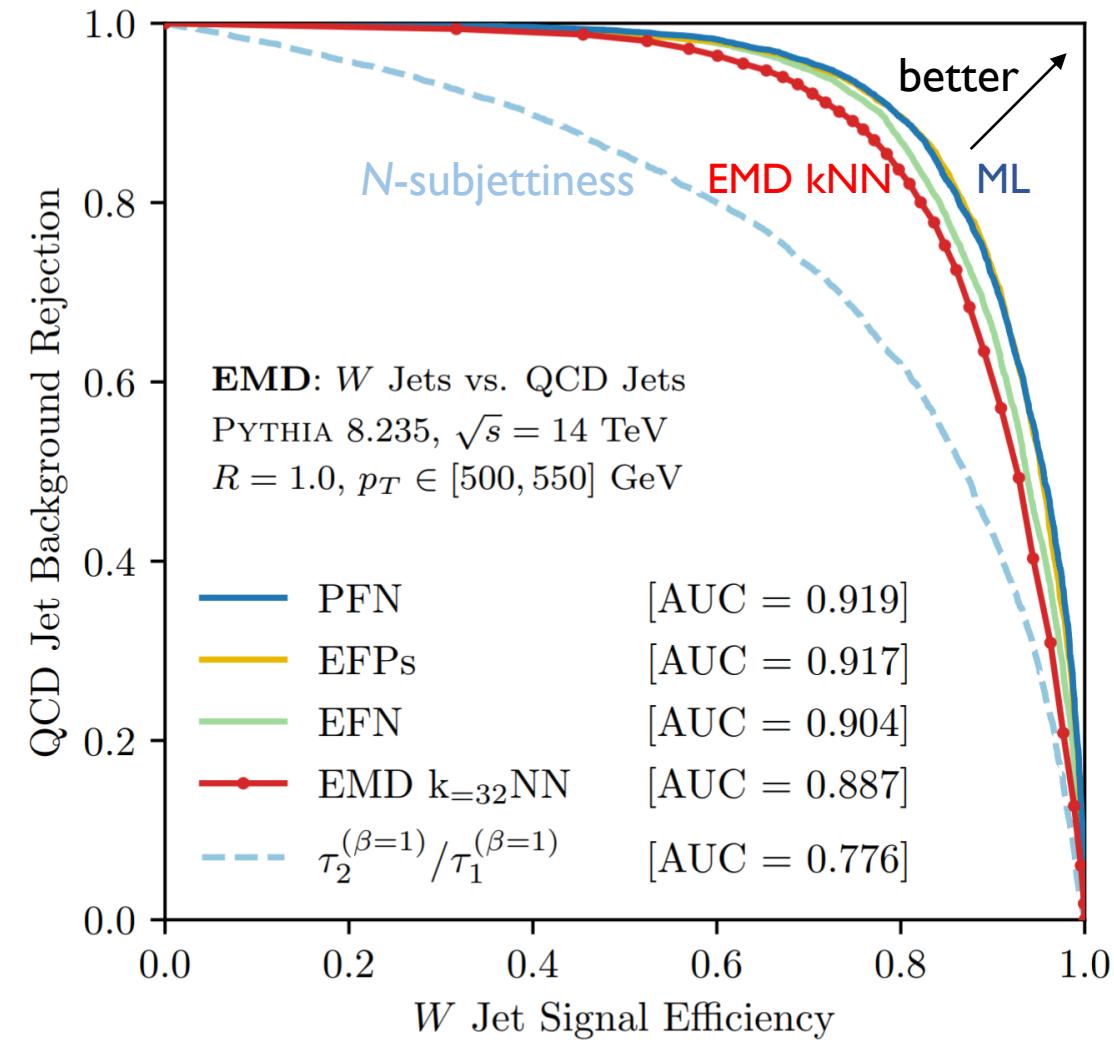
vs.



Look at a jet's nearest neighbors (kNN) to predict its class.

Optimal IRC-safe classifier with enough data.

Nearing performance of ML.



Exploring the Space of Events

Use EMD as a measure of event similarity

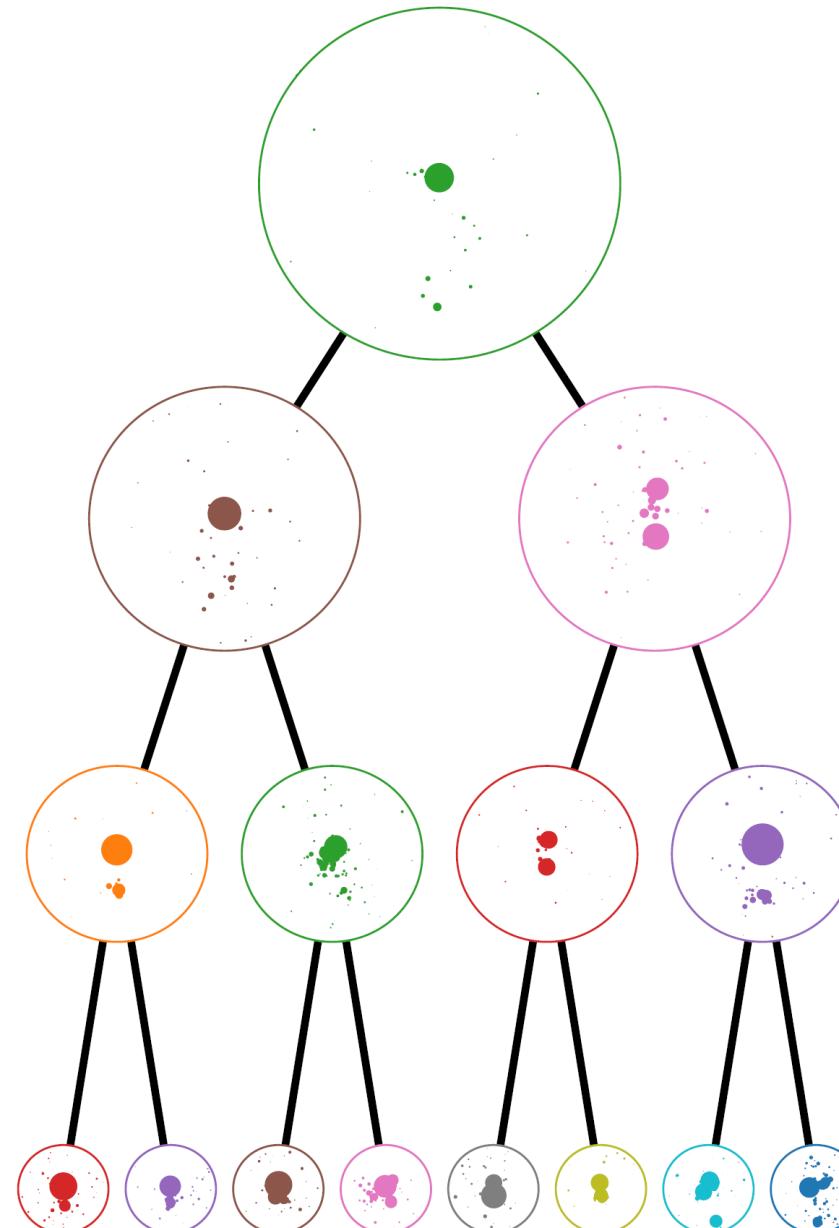
Unsupervised clustering algorithms can be used to cluster events

Jets are clusters of particles
???? are clusters of jets

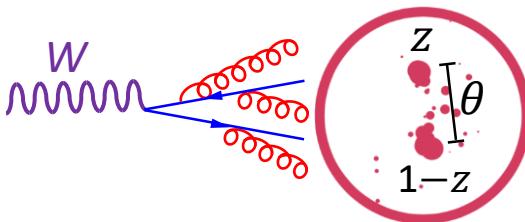
VP Tree: $O(\log(N))$ neighbor query time

Much more to explore.

Vantage Point (VP) Tree



Exploring the Space of Events: W jets



W jets are 2-pronged:

z : Energy Sharing of Prongs

θ : Angle between Prongs

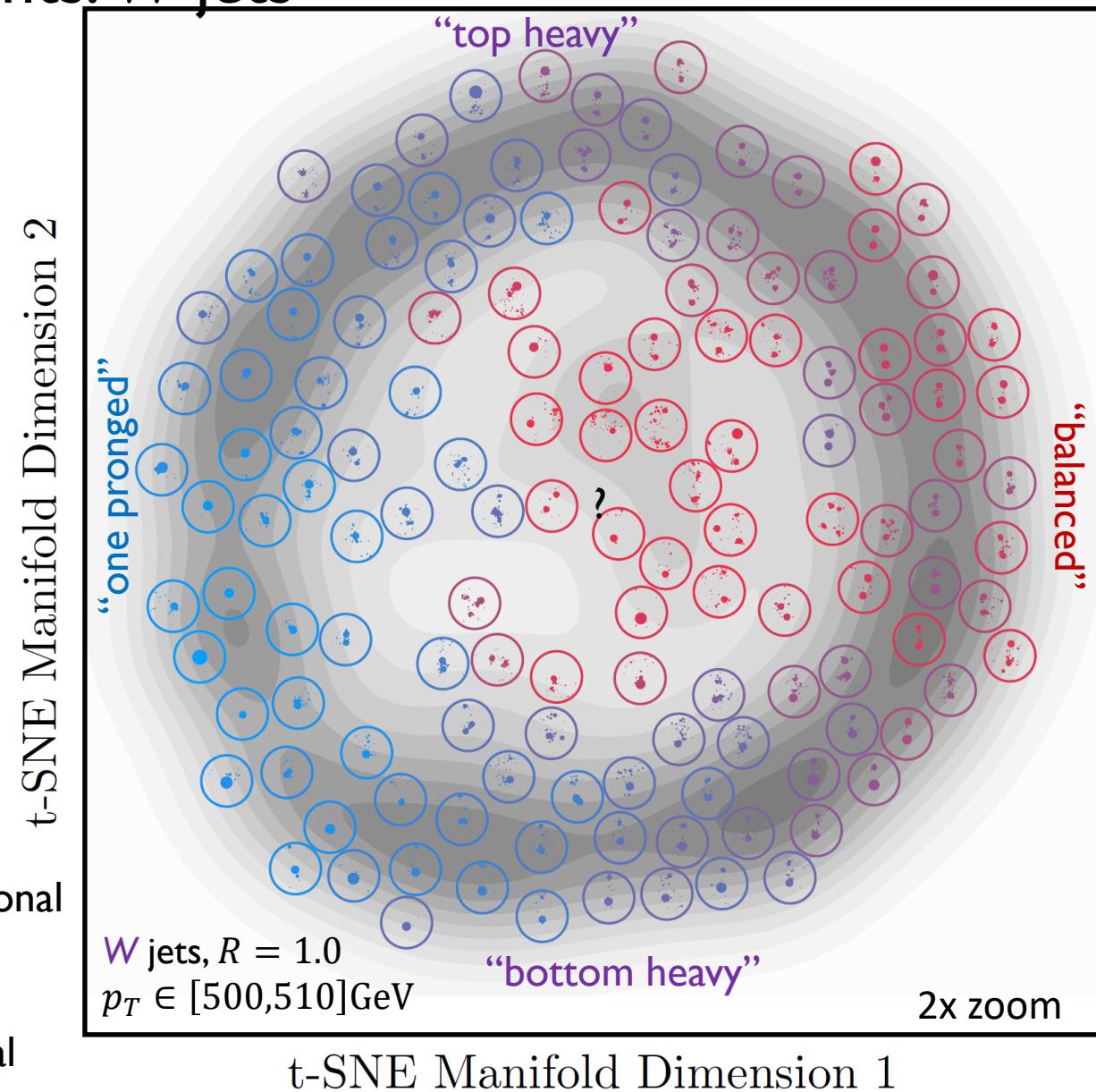
φ : Azimuthal orientation

Constrained by W mass:

$$z(1 - z)\theta^2 = \frac{p_{\mu J}^2}{p_T^2} = \frac{m_W^2}{p_T^2}$$

Hence we expect a **two-dimensional** space of W jets.

After φ rotation: **one-dimensional**



Exploring the Space of Jets: Correlation Dimension

Sketch of leading log (one emission) calculation:

$$\dim_{q/g}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_{i=1}^N \sum_{j=1}^N \Theta[\text{EMD}(\varepsilon_i, \varepsilon_j) < Q]$$

$$= Q \frac{\partial}{\partial Q} \ln \Pr [\text{EMD} < Q]$$

$$= Q \frac{\partial}{\partial Q} \ln \Pr [\lambda^{(\beta=1)} < Q; C_{q/g} \rightarrow 2 C_{q/g}]$$

$$= Q \frac{\partial}{\partial Q} \ln \exp \left(- \frac{4\alpha_s C_{q/g}}{\pi} \ln^2 \frac{Q}{p_T/2} \right)$$

$$= - \frac{8\alpha_s C_{q/g}}{\pi} \ln \frac{Q}{p_T/2}$$

+ 1-loop running of α_s

$$C_q = C_F = \frac{4}{3}$$

$$C_g = C_A = 3$$

