

# The Hidden Geometry of Particle Collisions

Patrick T. Komiske III

Massachusetts Institute of Technology  
Center for Theoretical Physics

*Based on work with Eric Metodiev and Jesse Thaler*

[1902.02346](#) (PRL)

[2004.04159](#) (JHEP)

BSM Pandemic Double Feature

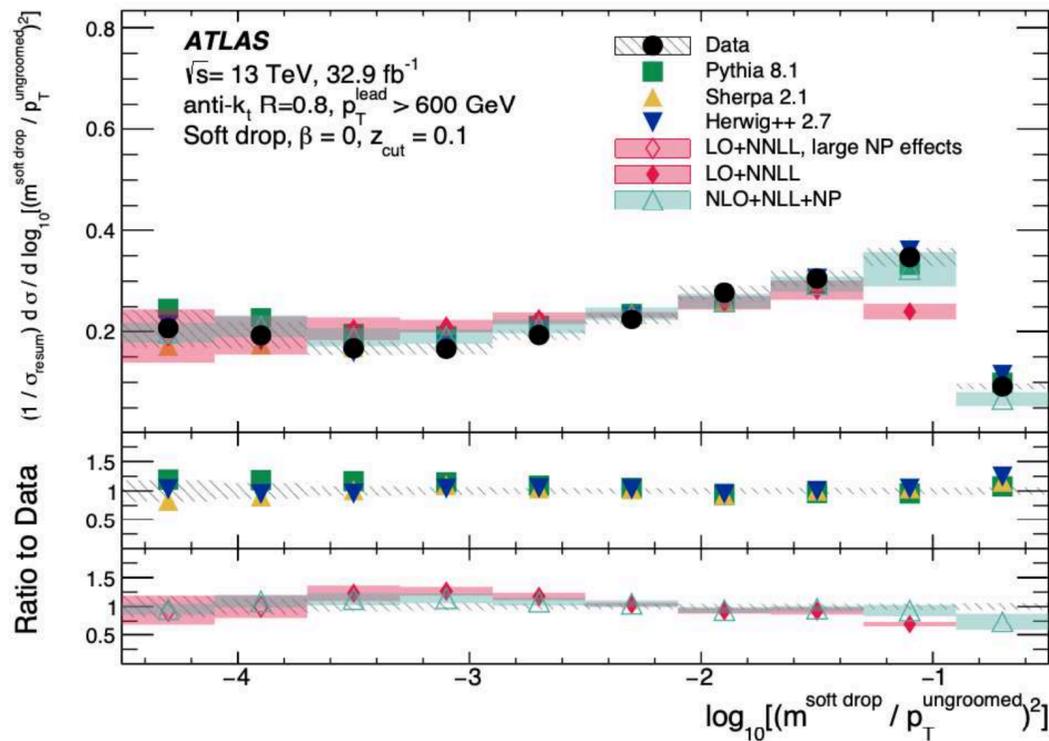
December 1, 2020

# Developing New Analysis Frameworks

Particle theory *makes predictions* and *invents new models*

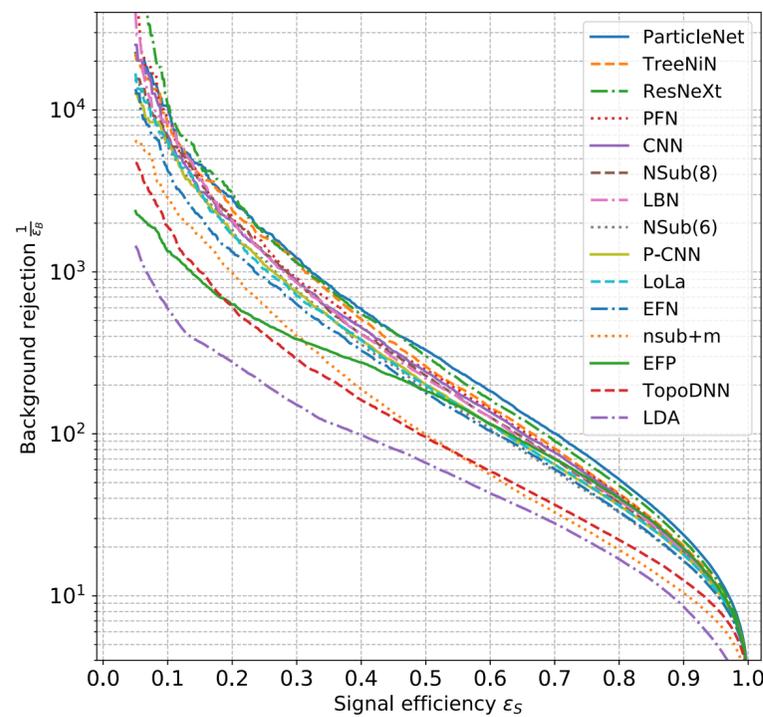
Particle experiment *conducts measurements* and *searches for new physics*

Soft Drop jet mass measurement



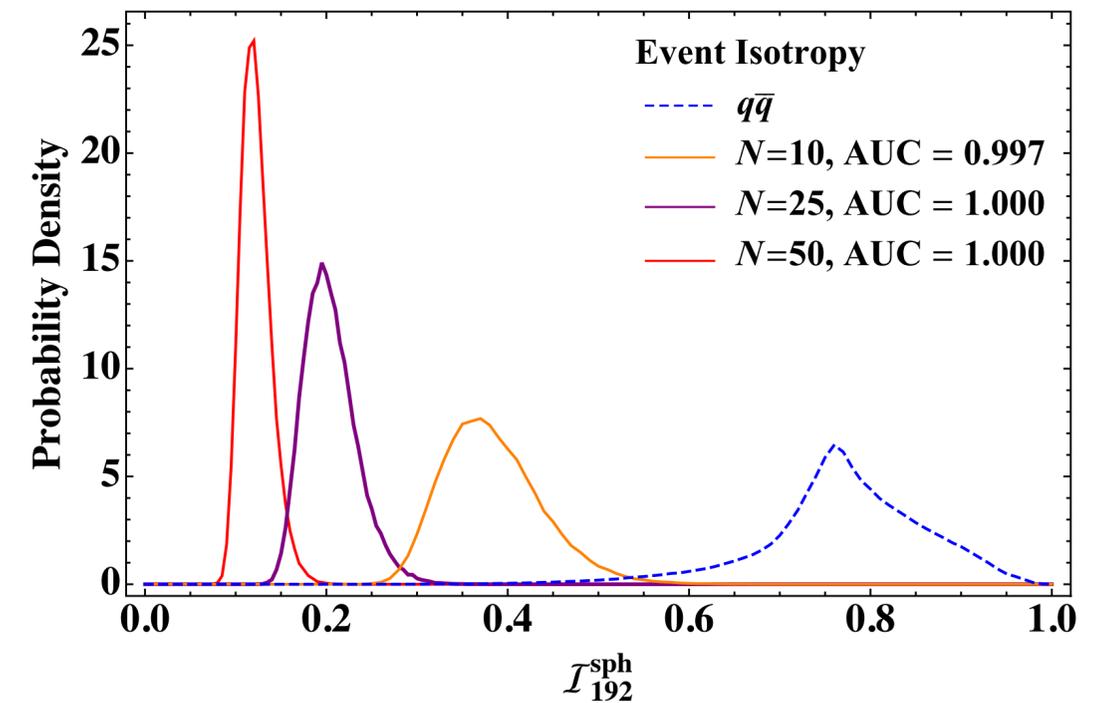
[Larkoski, Marzani, Soyez, Thaler, JHEP 2014;  
 ATLAS, PRL 2018]

Comparison of ML top taggers



[Kasieczka, Plehn, et al., 1902.09914;  
 using PTK, Metodiev, Thaler, 1810.05165, and others]

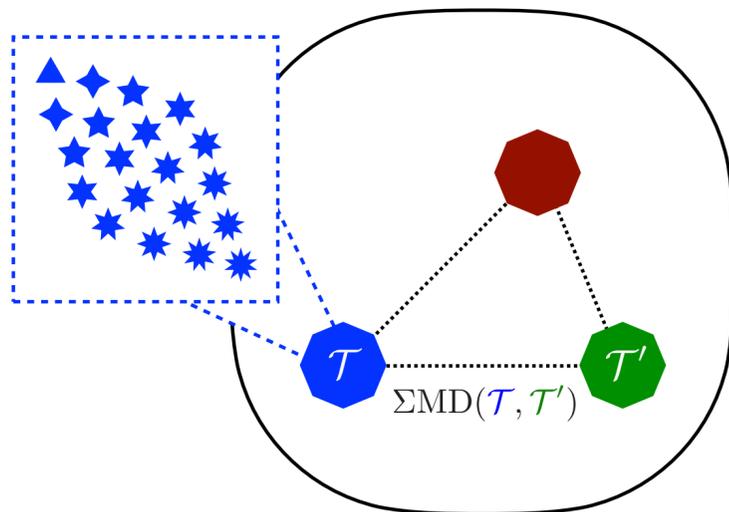
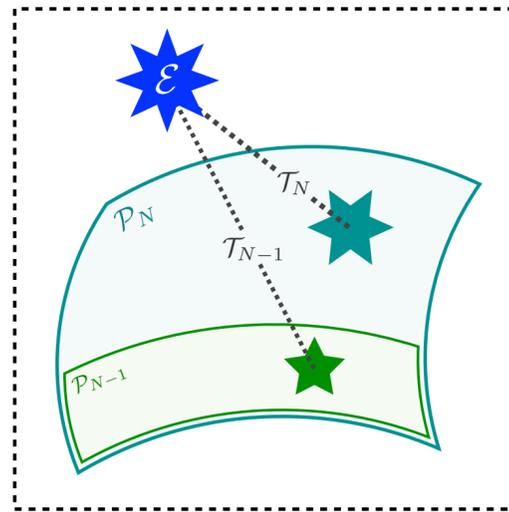
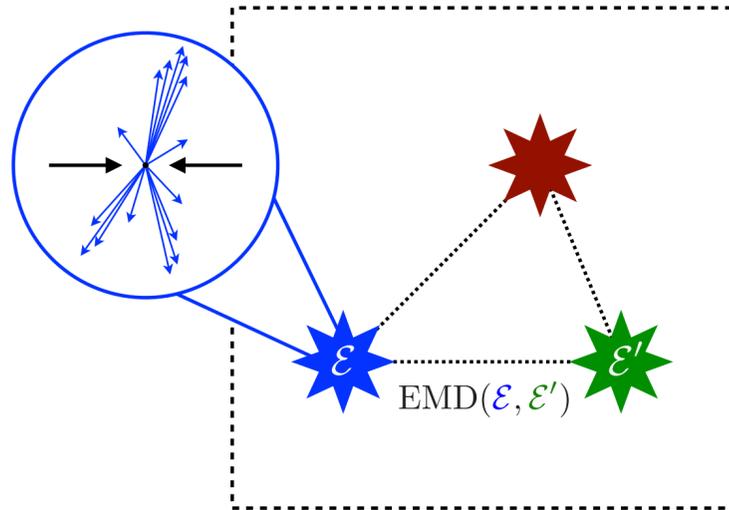
New observables sensitive to new signals



[Cesarotti, Thaler, JHEP 2020;  
 utilizing PTK, Metodiev, Thaler, PRL 2019]

Twenty more years of the LHC + any future collider(s)

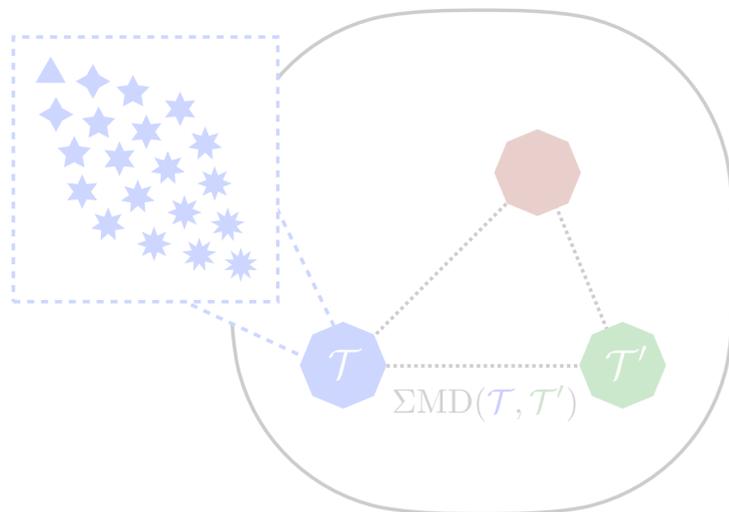
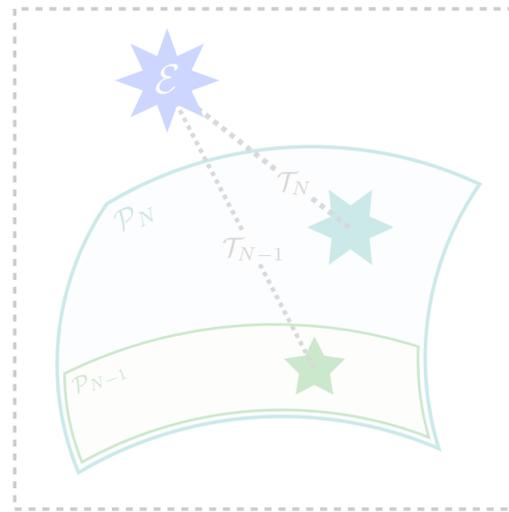
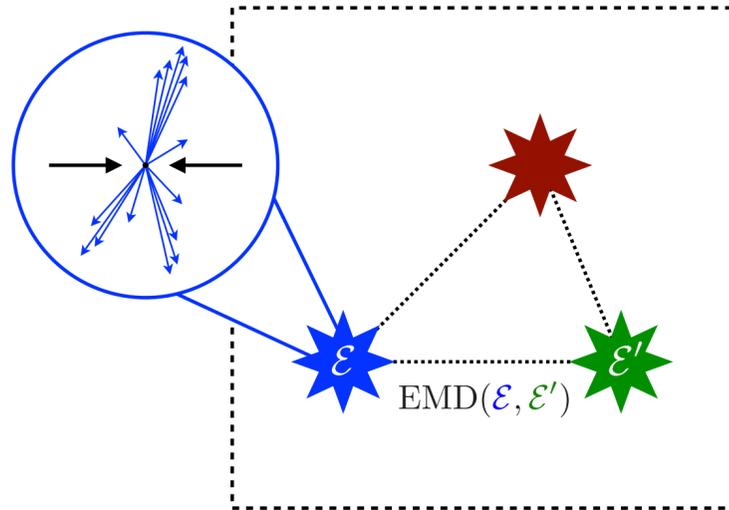
Maximizing physics potential will require insights from data science and ML



The (Metric) Space of Events

Revealing Hidden Geometry

[Theory Space]



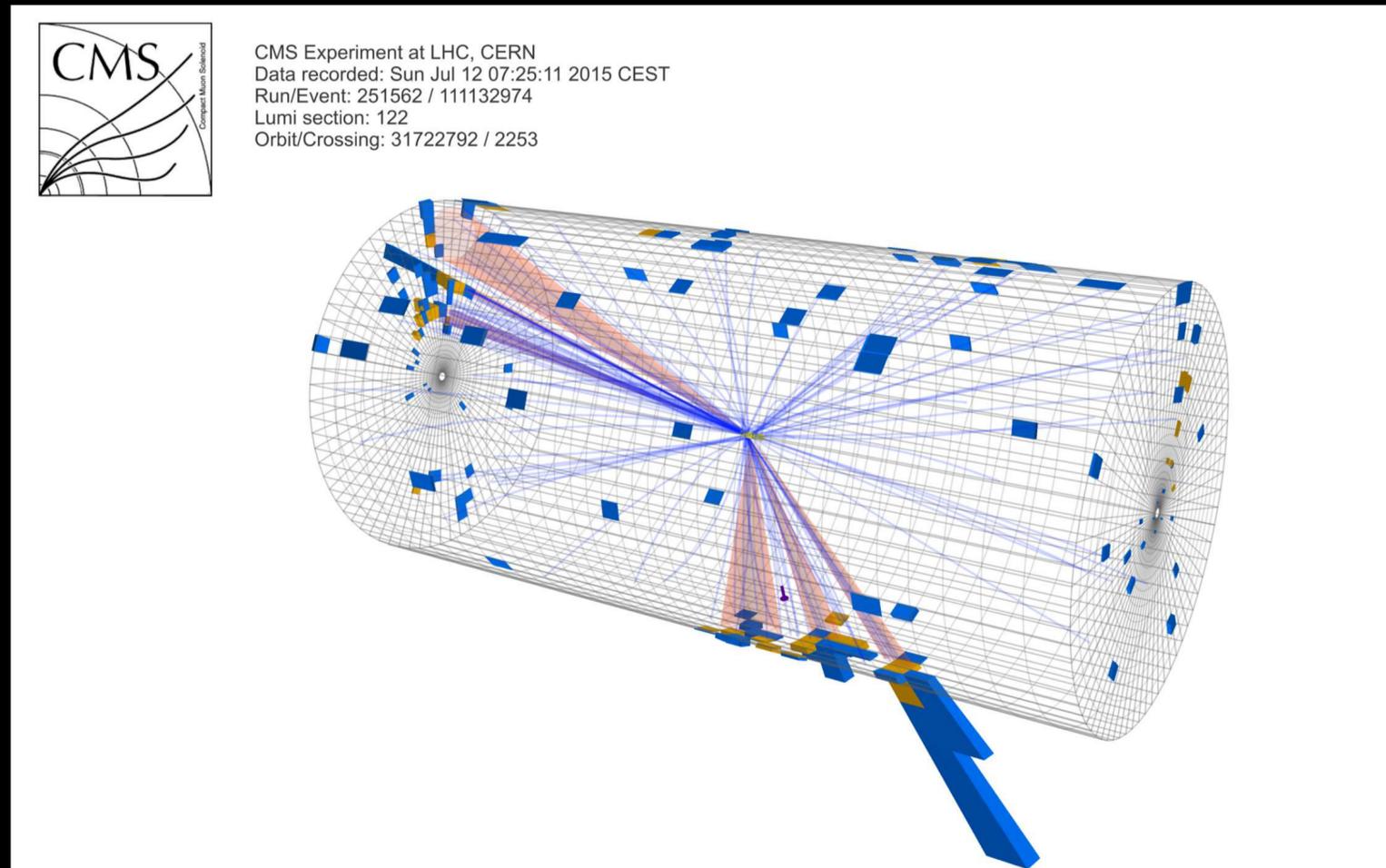
# The (Metric) Space of Events

## Revealing Hidden Geometry

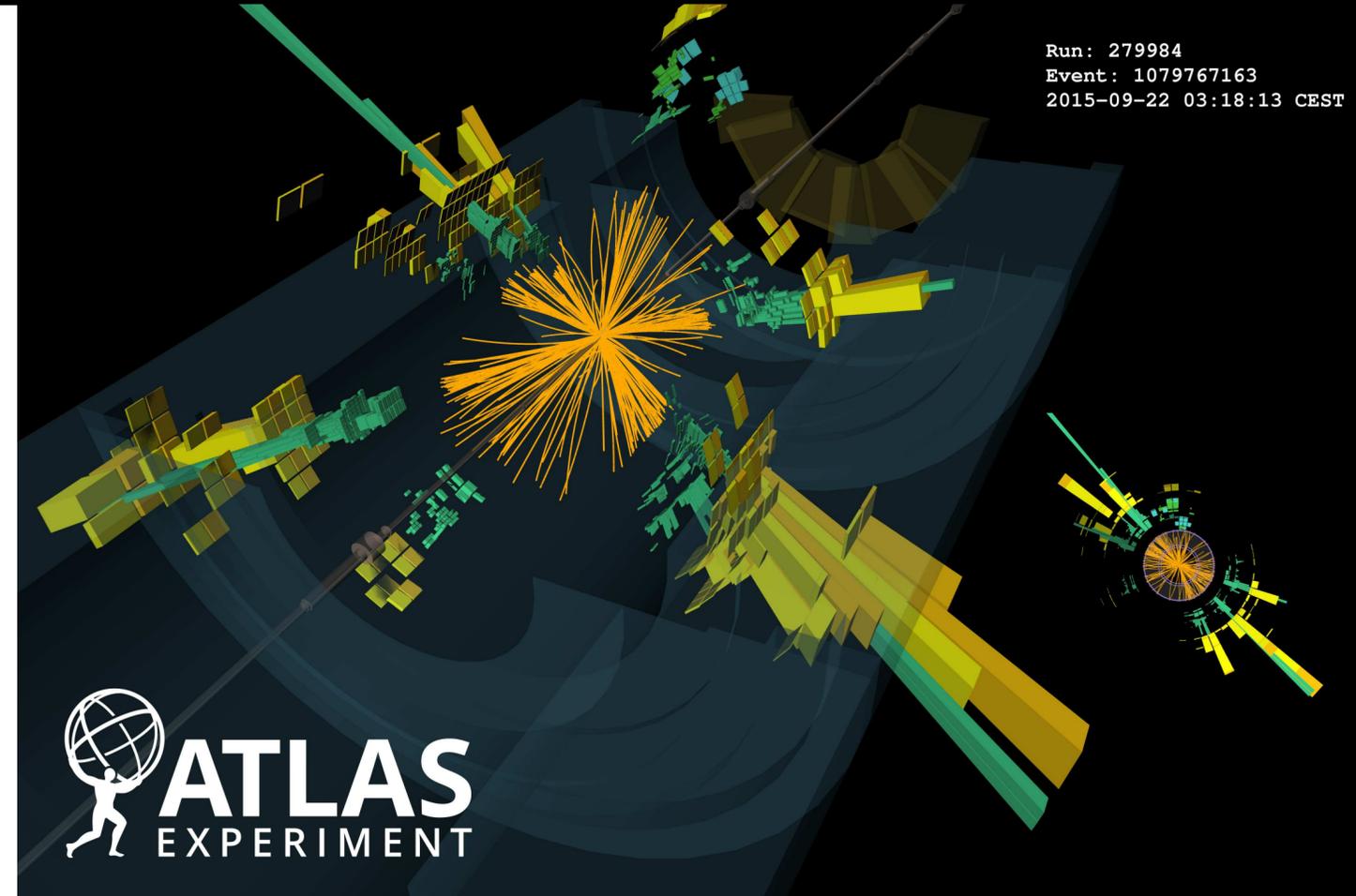
### [Theory Space]

# Explicit Geometry – Individual Events at the LHC

High-energy collisions produce final state particles with *energy*, *direction*, *charge*, *flavor*, and *other quantum numbers*



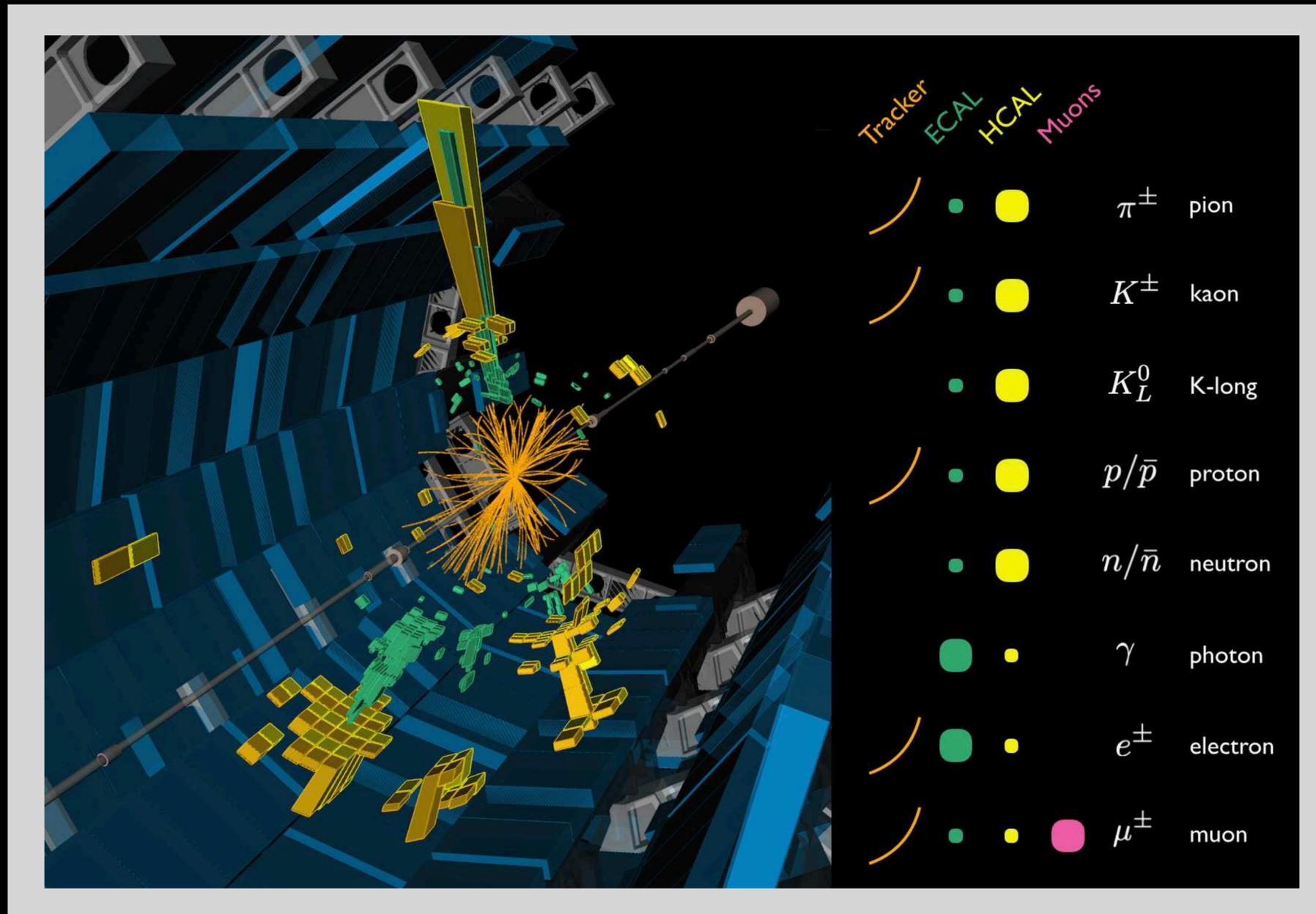
CMS hadronic  $t\bar{t}$  event



ATLAS high jet-multiplicity event

# Explicit Geometry – Individual Events at the LHC

High-energy collisions produce final state particles with *energy*, *direction*, *charge*, *flavor*, and *other quantum numbers*

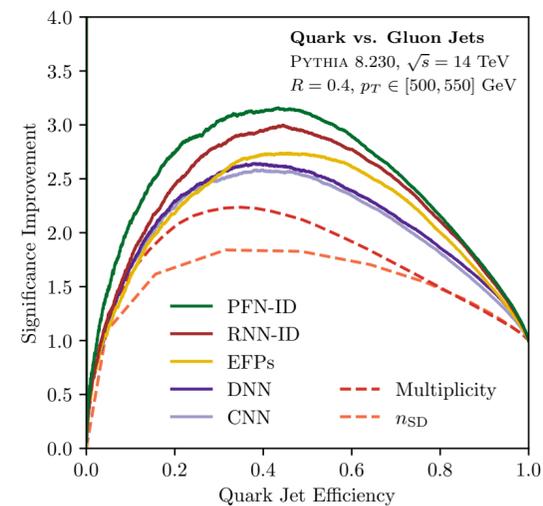
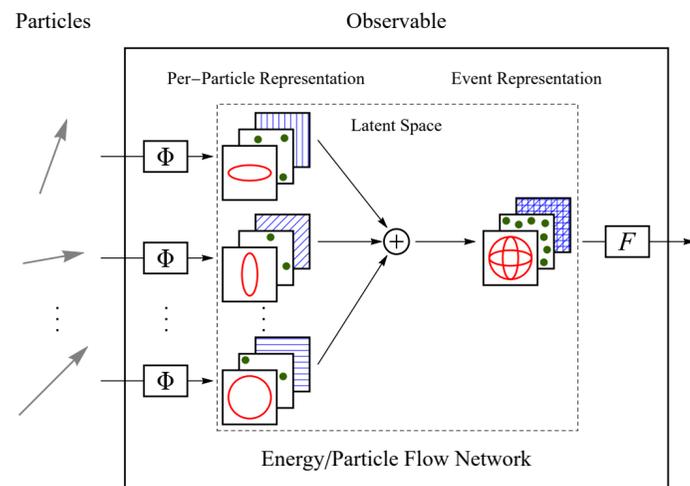


# Machine-Learning-Inspired Methods for Particle Physics

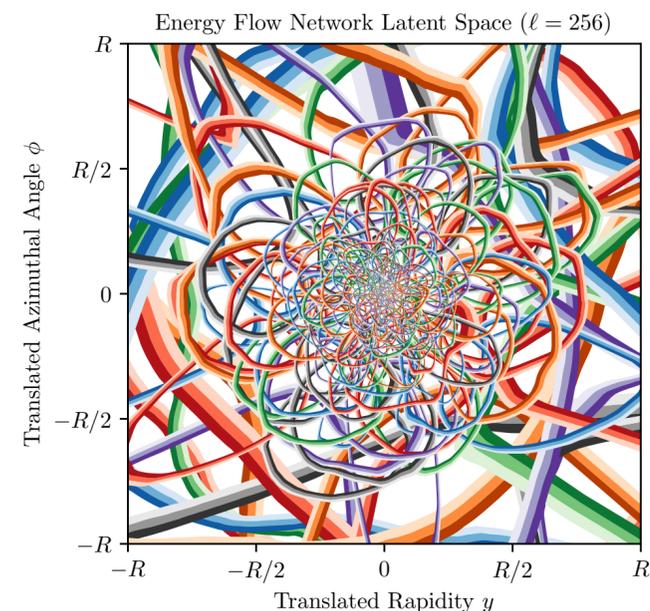
## Energy/Particle Flow Networks (EFNs/PFNs)

[PTK, Metodiev, Thaler, JHEP 2019]

Permutation symmetric neural network architecture for events with variable numbers of particles



Can be used to build powerful taggers



Latent space visualization reveals what the network has learned

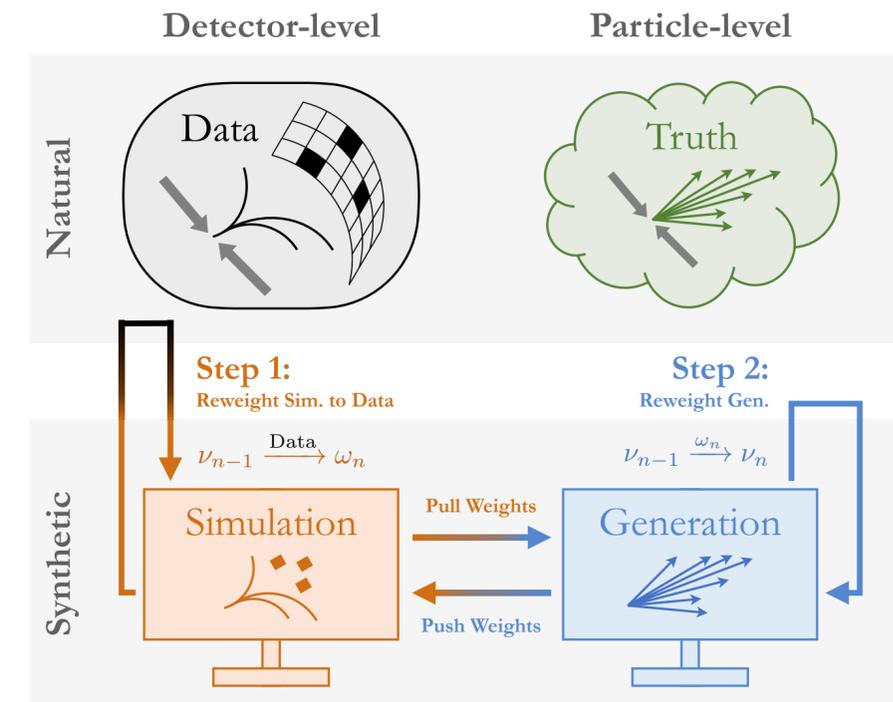
Dynamic pixel sizing related to collinear singularity of QCD!

## OmniFold

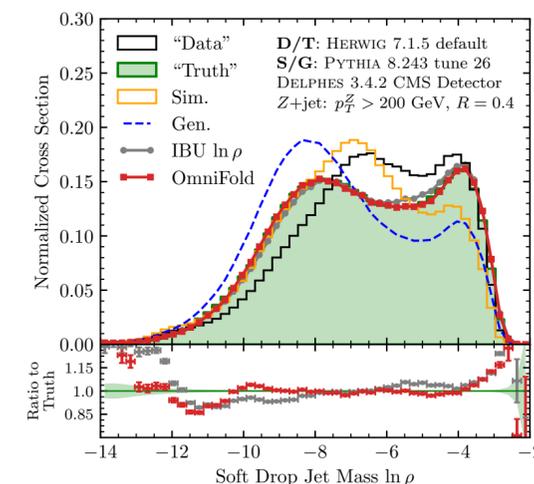


[Andreassen, PTK, Metodiev, Nachman, Thaler, PRL 2020]

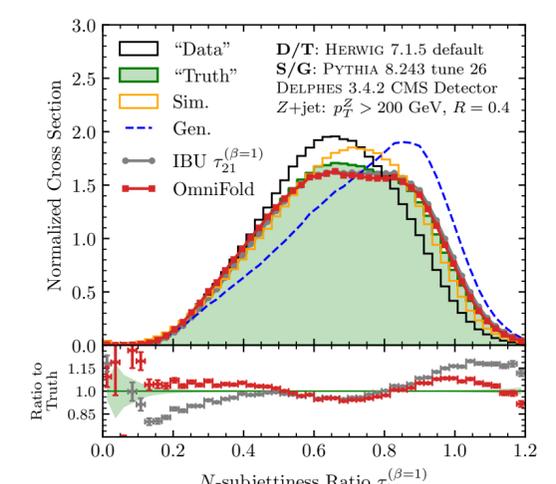
Unbinned, full-phase space unfolding of all observables simultaneously



Single application of OmniFold succeeds where multiple IBUs are required!

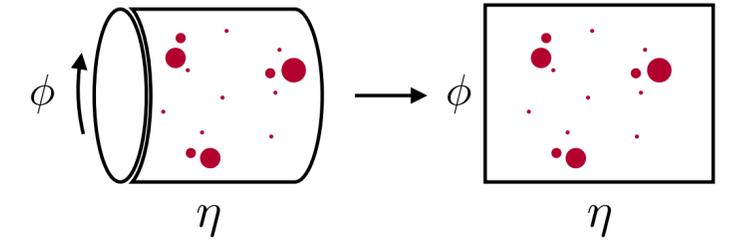


IRC safe



Sudakov safe

# Back to Explicit Geometry – Events as Distributions of Energy



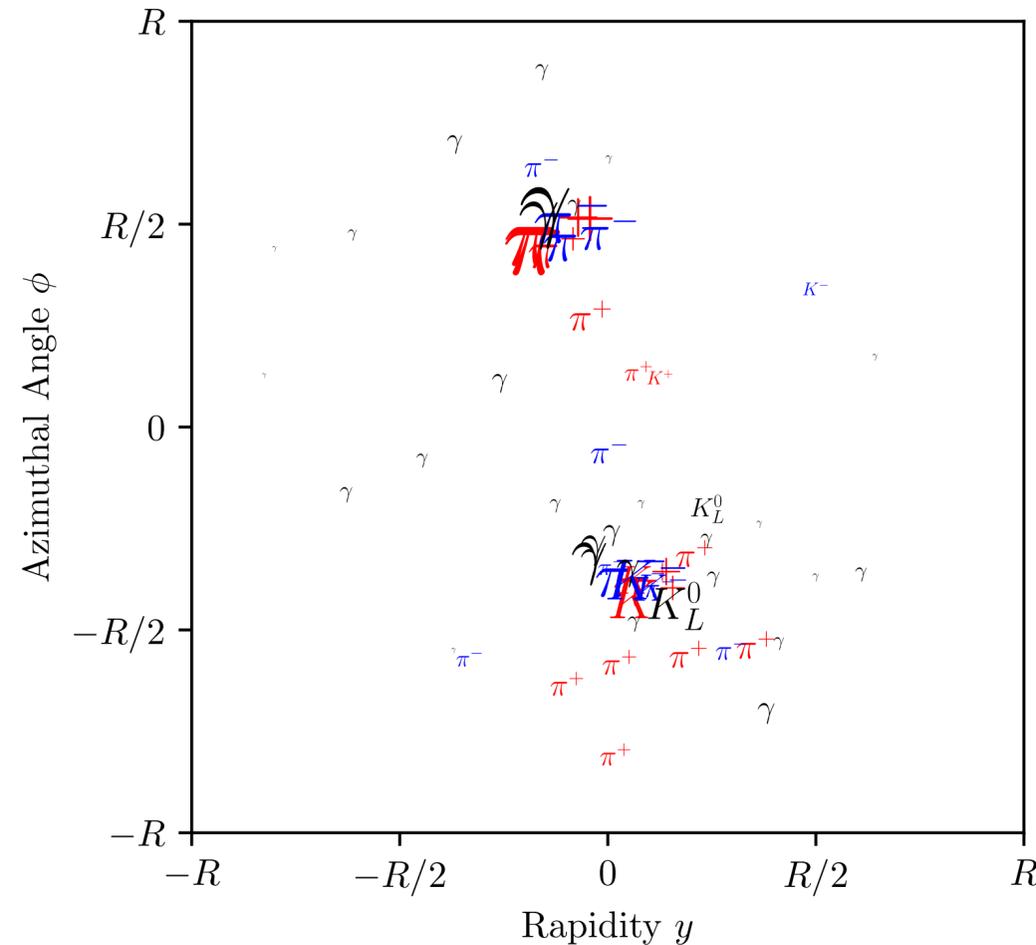
[PTK, Metodiev, Thaler, [JHEP 2019](#); PTK, Metodiev, Thaler, [JHEP 2020](#)]

**Energy flow distribution fully captures IRC-safe information**

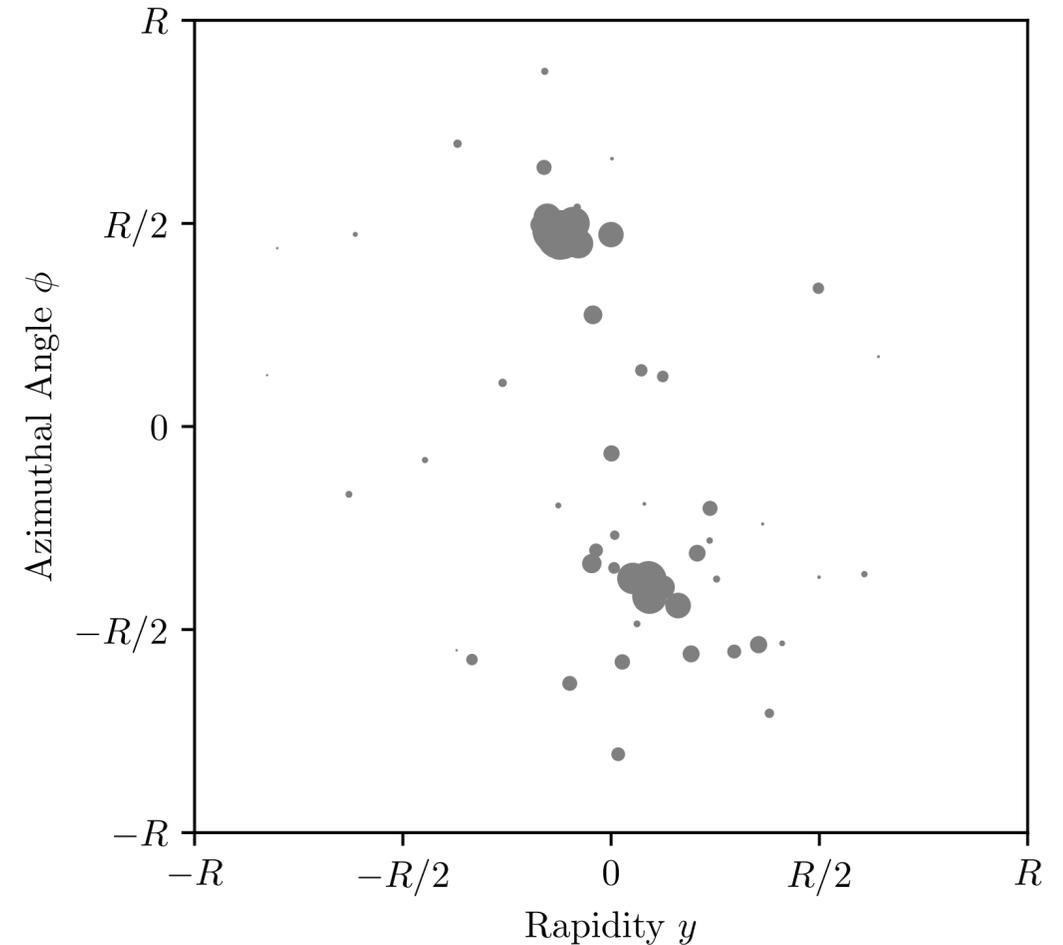
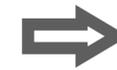
$$\mathcal{E}(\hat{n}) = \sum_{i=1}^M E_i \delta(\hat{n} - \hat{n}_i)$$

↑ Energy Flow Distribution    
 ↑ Energy (pT)    
 ↑ Direction (y, φ)

Full event is a set of particles having momentum and charge/ flavor

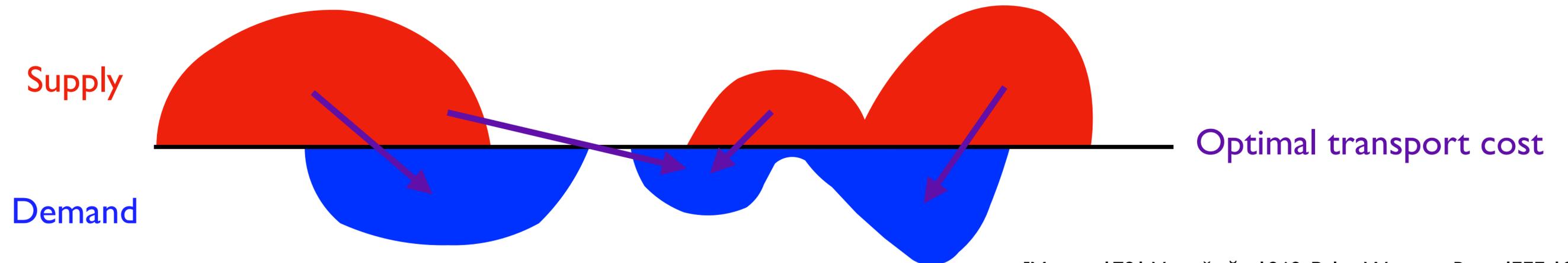
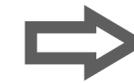
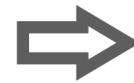


The **energy** flow is unpixelized and ignores charge/ flavor information



# Very Basic Question – When are Two Distributions Similar?

Optimal transport minimizes the “work” (stuff x distance) required to transport supply to demand



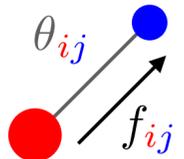
[Monge, 1781; Vaserštejn, 1969; Peleg, Werman, Rom, IEEE 1989; Rubner, Tomasi, Guibas, ICCV 1998, ICJV 2000; Pele, Werman, ECCV 2008; Pele, Taskar, GSI 2013]

# The Energy Mover's Distance (EMD)

[PTK, Metodiev, Thaler, [PRL 2019](#)]

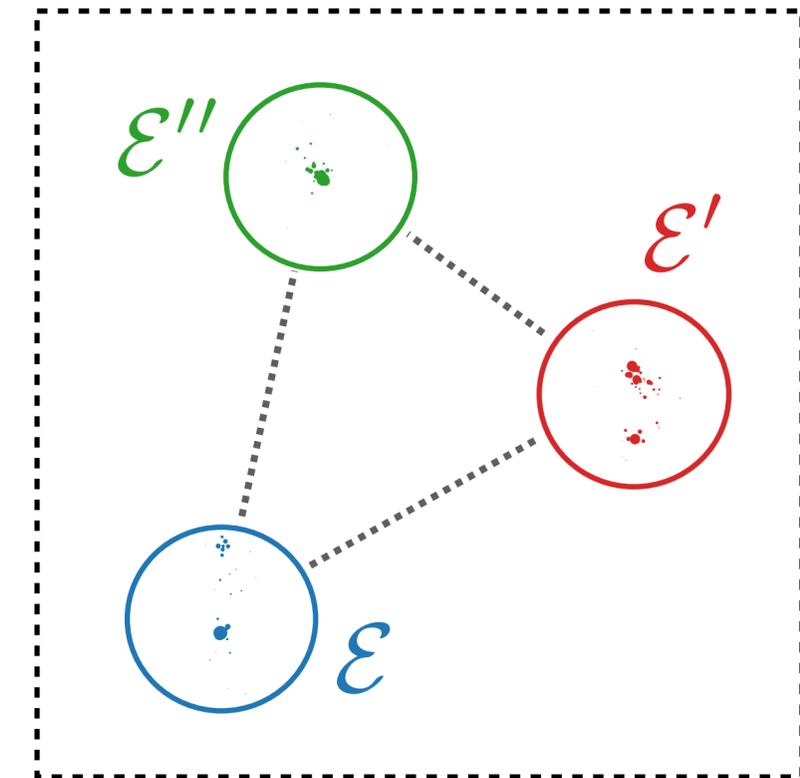
EMD between *energy flows* defines a *metric* on the space of events

$$\text{EMD}_{\beta,R}(\mathcal{E}, \mathcal{E}') = \underbrace{\min_{\{f_{ij} \geq 0\}} \sum_i \sum_j f_{ij} \left( \frac{\theta_{ij}}{R} \right)^\beta}_{\text{Cost of optimal transport}} + \underbrace{\left| \sum_i E_i - \sum_j E'_j \right|}_{\text{Cost of energy creation}}$$

$$\underbrace{\sum_j f_{ij} \leq E_i, \quad \sum_i f_{ij} \leq E'_j, \quad \sum_{ij} f_{ij} = \min \left( \sum_i E_i, \sum_j E'_j \right)}_{\text{Capacity constraints to ensure proper transport}}$$


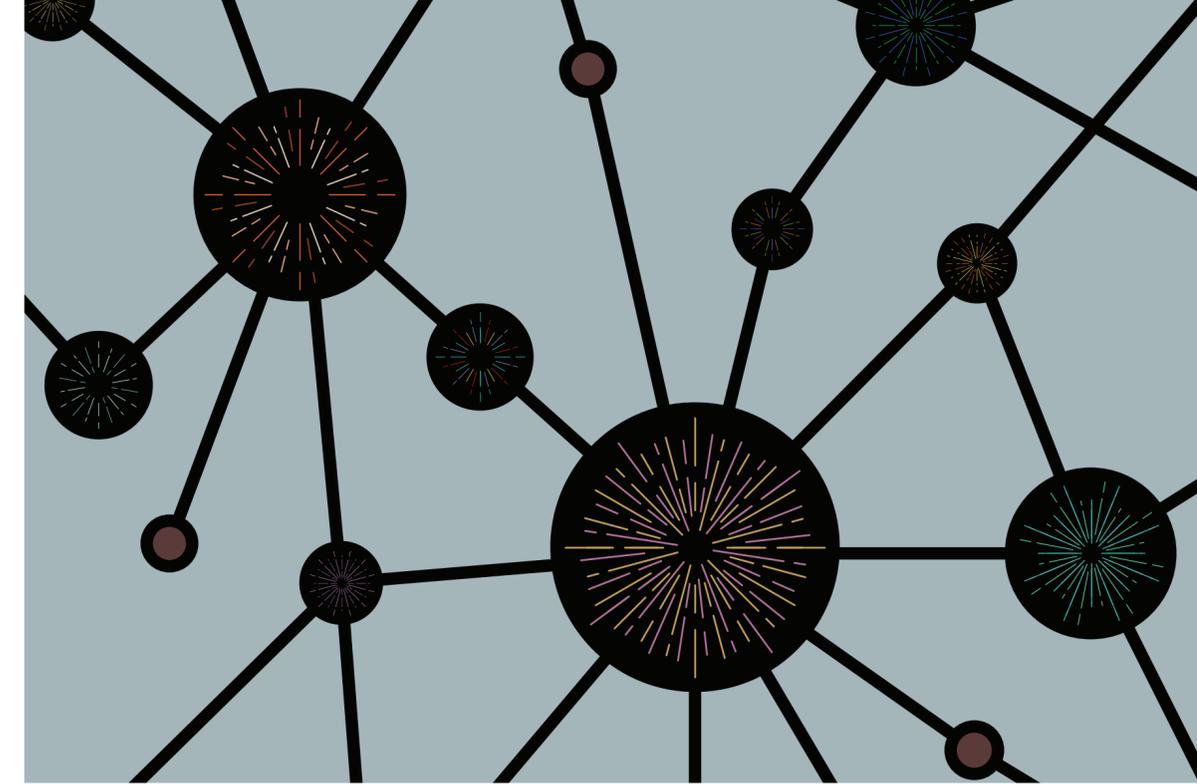
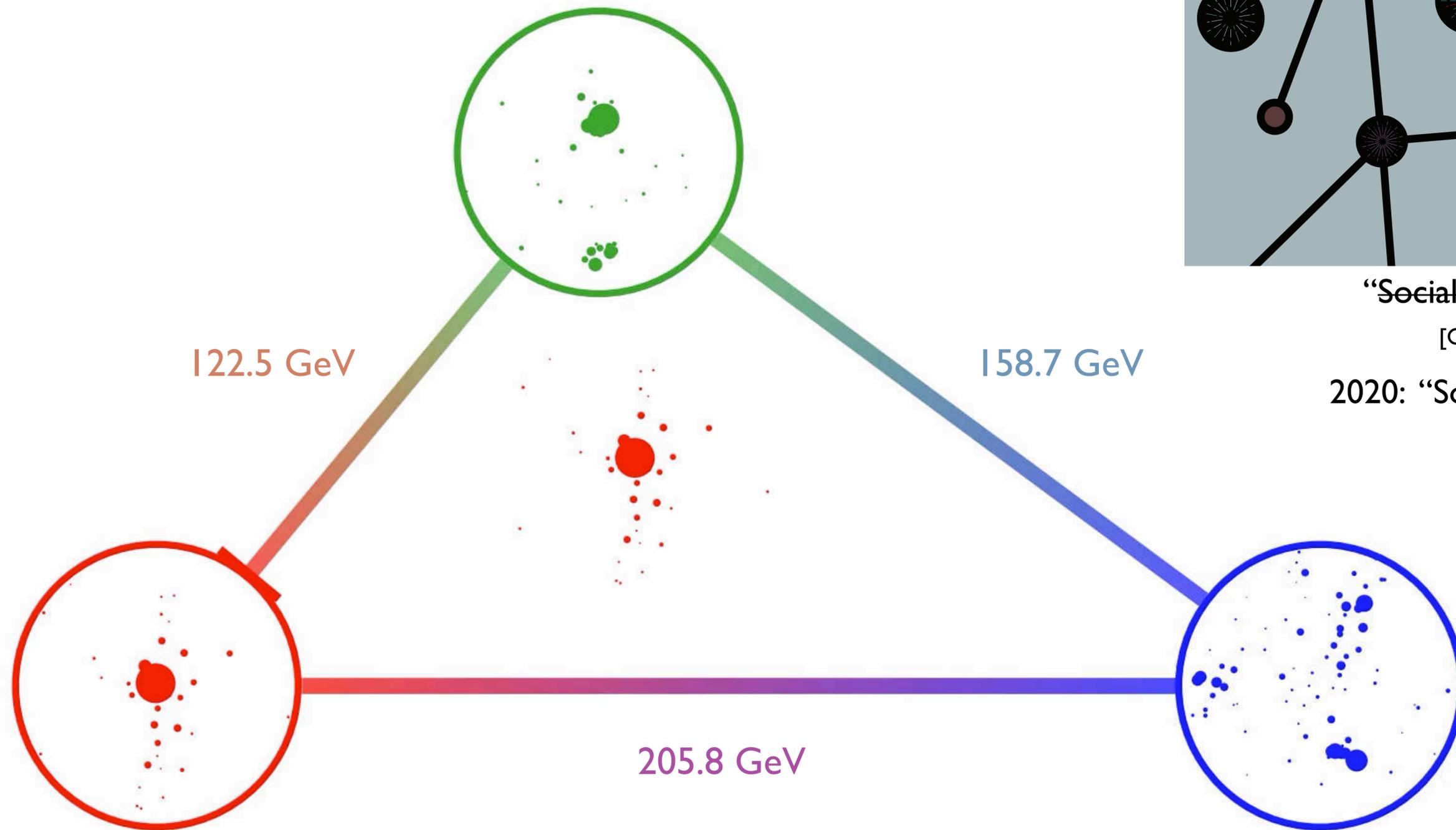
$R$ : controls cost of transporting energy vs. destroying/creating it

$\beta$ : angular weighting exponent



Triangle inequality satisfied for  $R \geq d_{\max}/2$   
 $0 \leq \text{EMD}(\mathcal{E}, \mathcal{E}') \leq \text{EMD}(\mathcal{E}, \mathcal{E}'') + \text{EMD}(\mathcal{E}'', \mathcal{E}')$   
 i.e.  $R \geq$  jet radius for conical jets

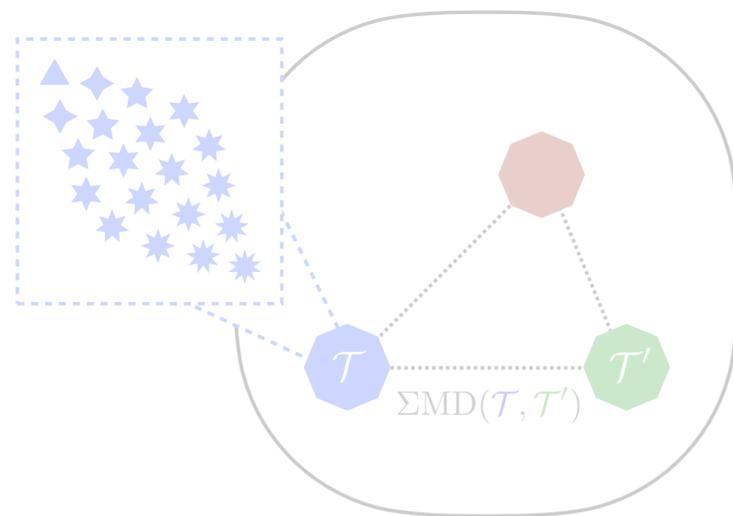
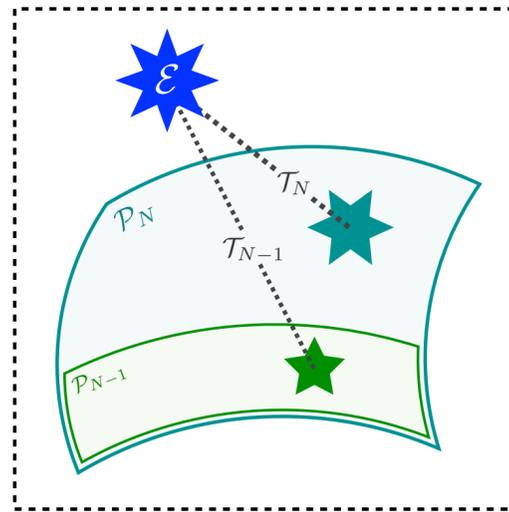
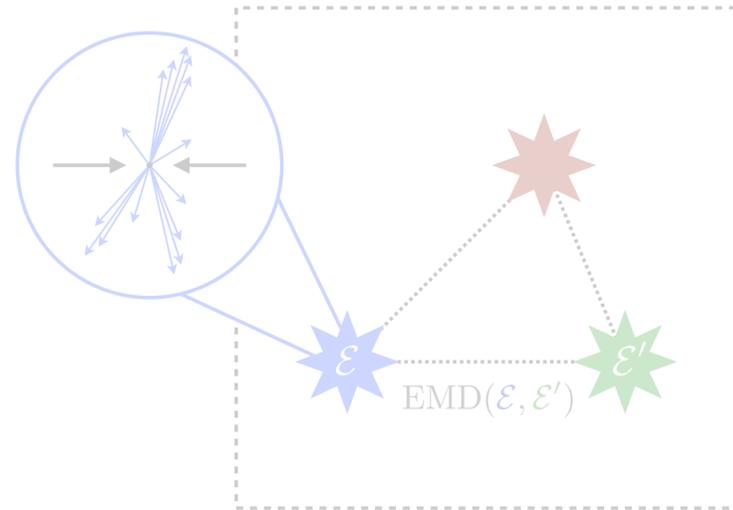
# Geodesics in the Space of Events



“Social networking of jets”

[Chu, MIT News 2019]

2020: “Social distancing of jets”



The (Metric) Space of Events

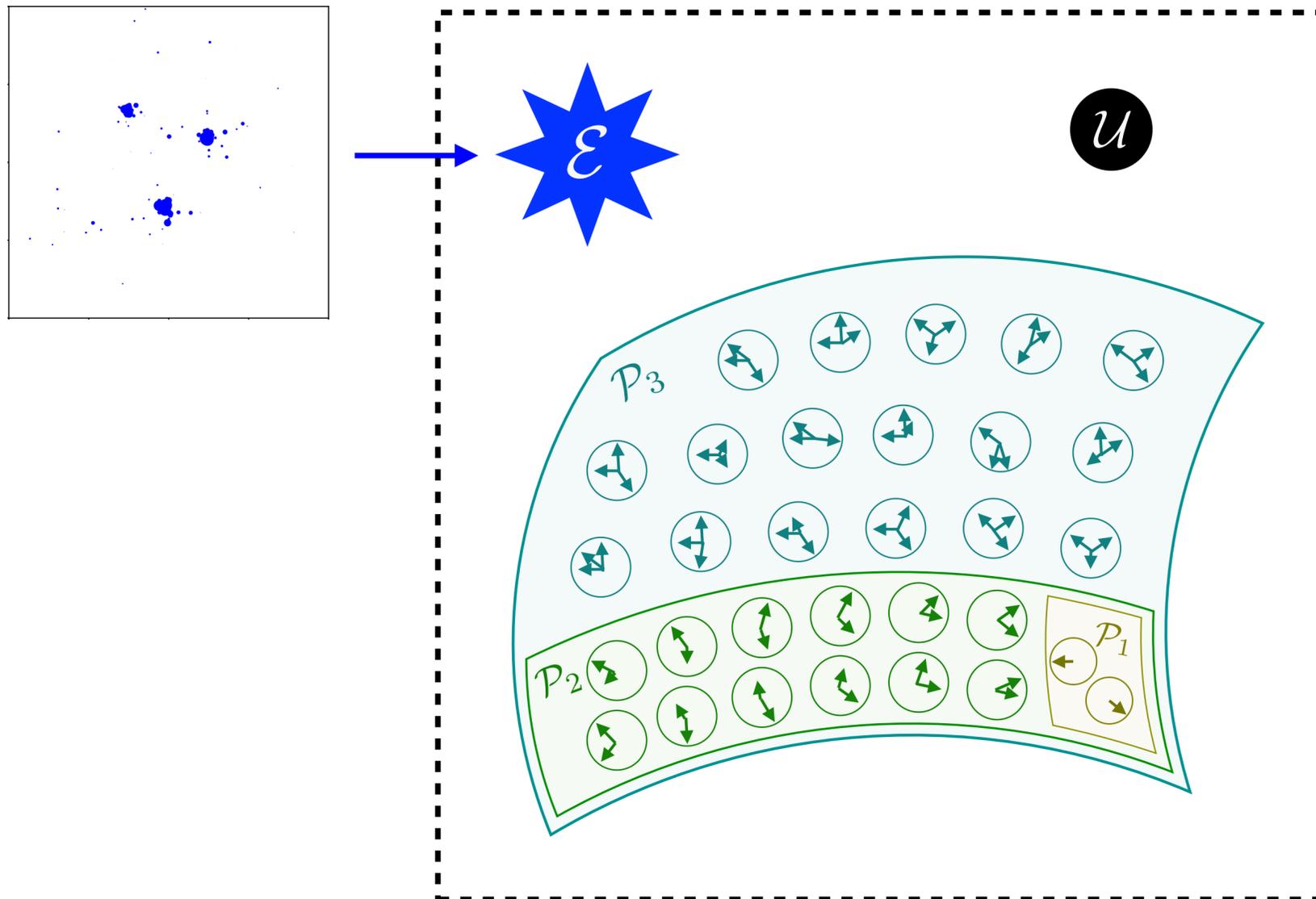
Revealing Hidden Geometry

[Theory Space]

# N-particle Manifolds in the Space of Events

[PTK, Metodiev, Thaler, JHEP 2020]

$$\mathcal{P}_N = \text{set of all } N\text{-particle configurations} = \left\{ \sum_{i=1}^N E_i \delta(\hat{n} - \hat{n}_i) \mid E_i \geq 0 \right\}$$



$\mathcal{P}_1$ : manifold of events with one particle

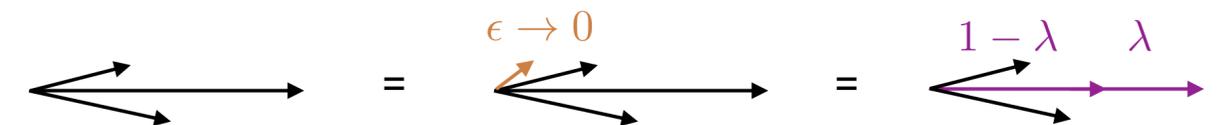
$\mathcal{P}_2$ : manifold of events with two particles

$\mathcal{P}_3$ : manifold of events with three particles

⋮

$$\mathcal{P}_N \supset \mathcal{P}_{N-1} \supset \cdots \supset \mathcal{P}_3 \supset \mathcal{P}_2 \supset \mathcal{P}_1$$

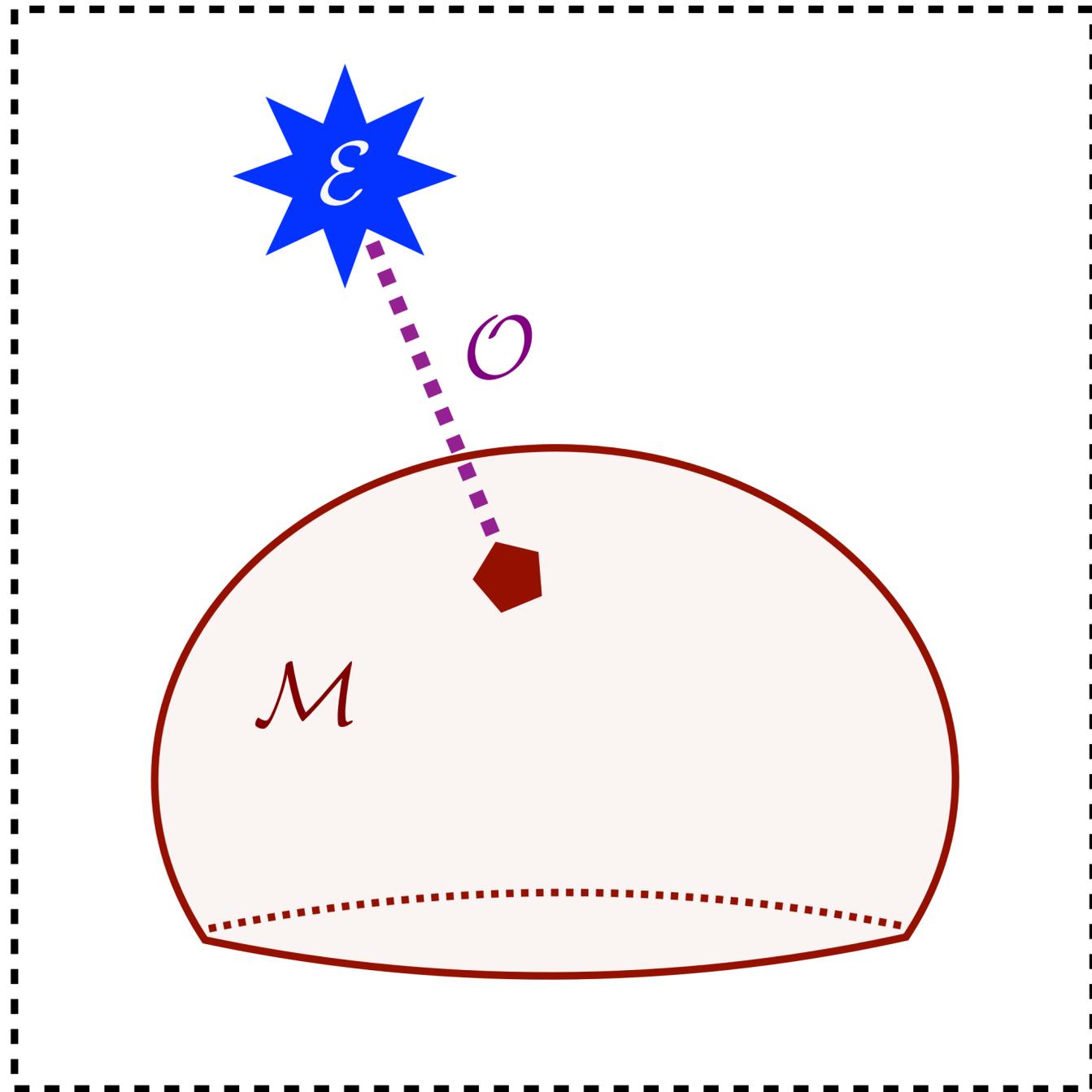
by **soft** and **collinear** limits



$\mathcal{U}$  Uniform event, not contained in any  $\mathcal{P}_N$

# Defining Observables via Event Space Geometry

[PTK, Metodiev, Thaler, JHEP 2020]



Many common *observables* are distance of closest approach from event to a specific *manifold*

$$\mathcal{O}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{M}} \text{EMD}_{\beta, R}(\mathcal{E}, \mathcal{E}')$$

EMD variant for equal-energy events

$$\text{EMD}_{\beta}(\mathcal{E}, \mathcal{E}') = \lim_{R \rightarrow \infty} R^{\beta} \text{EMD}_{\beta, R}(\mathcal{E}, \mathcal{E}') = \min_{\{f_{ij} \geq 0\}} \sum_{i=1}^M \sum_{j=1}^{M'} f_{ij} \theta_{ij}^{\beta}$$

Enforces equal energy (else infinity)      on equal-energy events

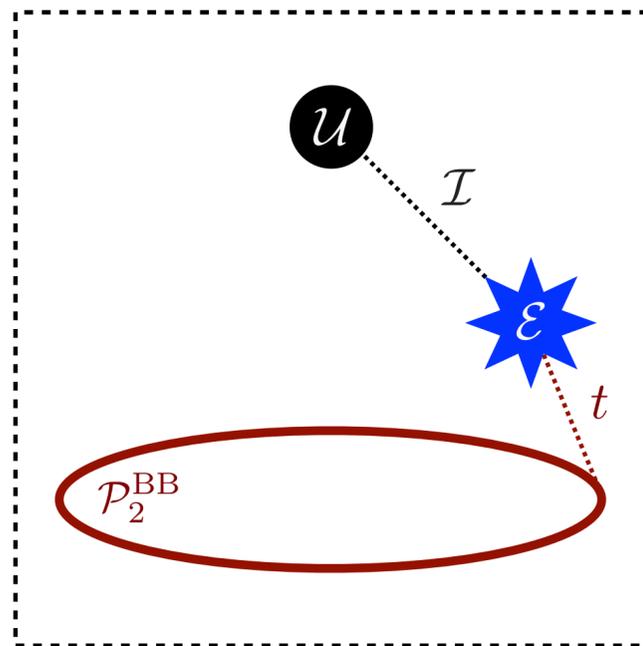
# Defining Observables via Event Space Geometry

$$\mathcal{O}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{M}} \text{EMD}_{\beta,R}(\mathcal{E}, \mathcal{E}')$$

[PTK, Metodiev, Thaler, JHEP 2020]

## Thrust, sphericity, isotropy\*

Distance of closest approach to a specific manifold



$$t(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{P}_2^{\text{BB}}} \text{EMD}_2(\mathcal{E}, \mathcal{E}')$$

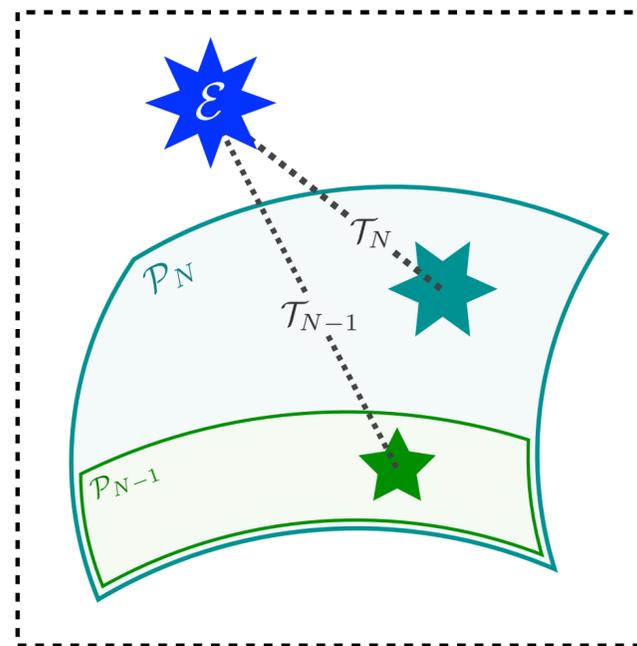
$$\sqrt{s(\mathcal{E})} = \min_{\mathcal{E}' \in \mathcal{P}_2^{\text{BB}}} \text{EMD}_1(\mathcal{E}, \mathcal{E}')$$

$$\mathcal{I}^{(\beta)}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{M}_{\mathcal{U}}} \text{EMD}_{\beta}(\mathcal{E}, \mathcal{E}')$$

[Farhi, PRL 1977; Georgi, Machacek, PRL 1977]  
\*New! [Cesarotti, Thaler, 2004.06125]

## N-jettiness

Minimum distance from event to N-particle manifold



without beam region

$$\mathcal{T}_N^{(\beta)}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{P}_N} \text{EMD}_{\beta}(\mathcal{E}, \mathcal{E}')$$

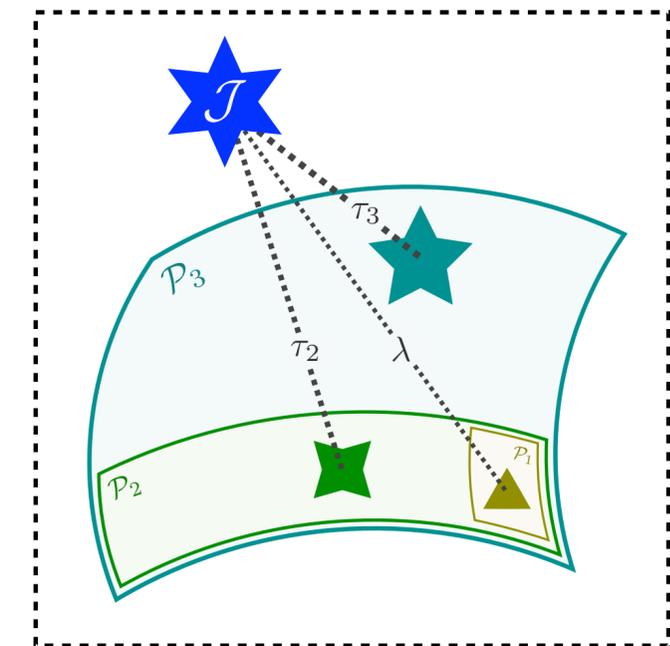
with constant beam distance  $R^{\beta}$

$$\mathcal{T}_N^{(\beta,R)}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{P}_N} \text{EMD}_{\beta,R}(\mathcal{E}, \mathcal{E}')$$

[Brandt, Dahmen, Z. Phys 1979;  
Stewart, Tackmann, Waalewijn, PRL 2010]

## N-subjettiness, angularities

Smallest distance from jet to N-particle manifold



for recoil-free angularity

$$\lambda_{\beta}(\mathcal{J}) = \min_{\mathcal{J}' \in \mathcal{P}_1} \text{EMD}_{\beta}(\mathcal{J}, \mathcal{J}')$$

$$\tau_N^{(\beta)}(\mathcal{J}) = \min_{\mathcal{J}' \in \mathcal{P}_N} \text{EMD}_{\beta}(\mathcal{J}, \mathcal{J}')$$

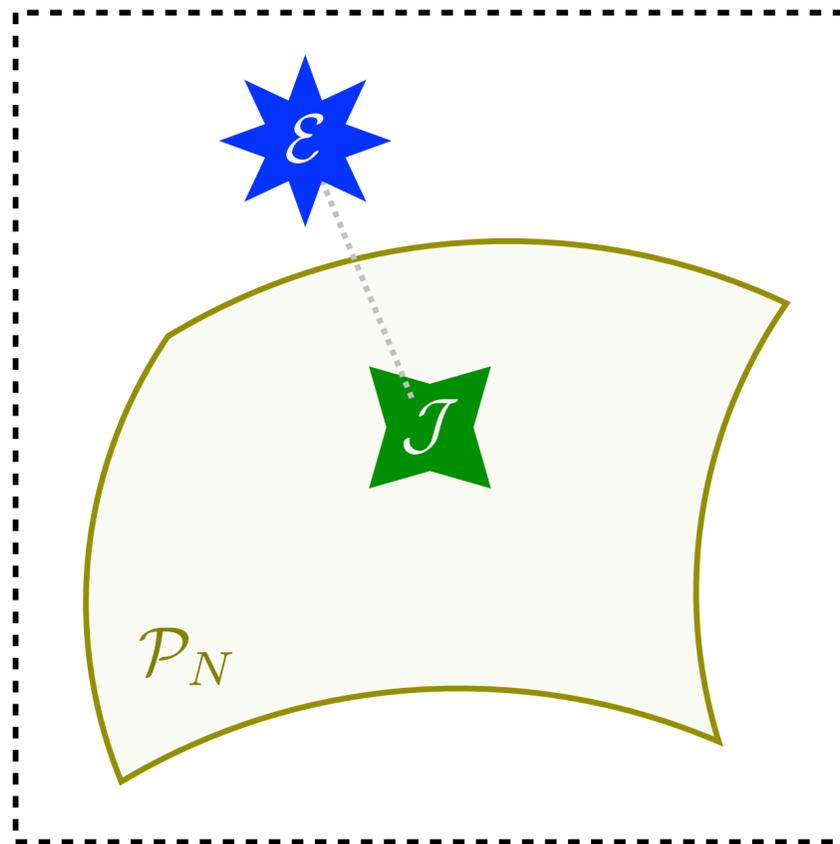
[Ellis, Vermilion, Walsh, Hornig, Lee, JHEP 2010;  
Thaler, Van Tilburg, JHEP 2011, JHEP 2012]

# Jets in the Space of Events – The Closest $N$ -particle Description of an $M$ -particle Event

[PTK, Metodiev, Thaler, [JHEP 2020](#)]

## Exclusive cone finding

*X*Cone finds  $N$  jets by minimizing  $N$ -jettiness

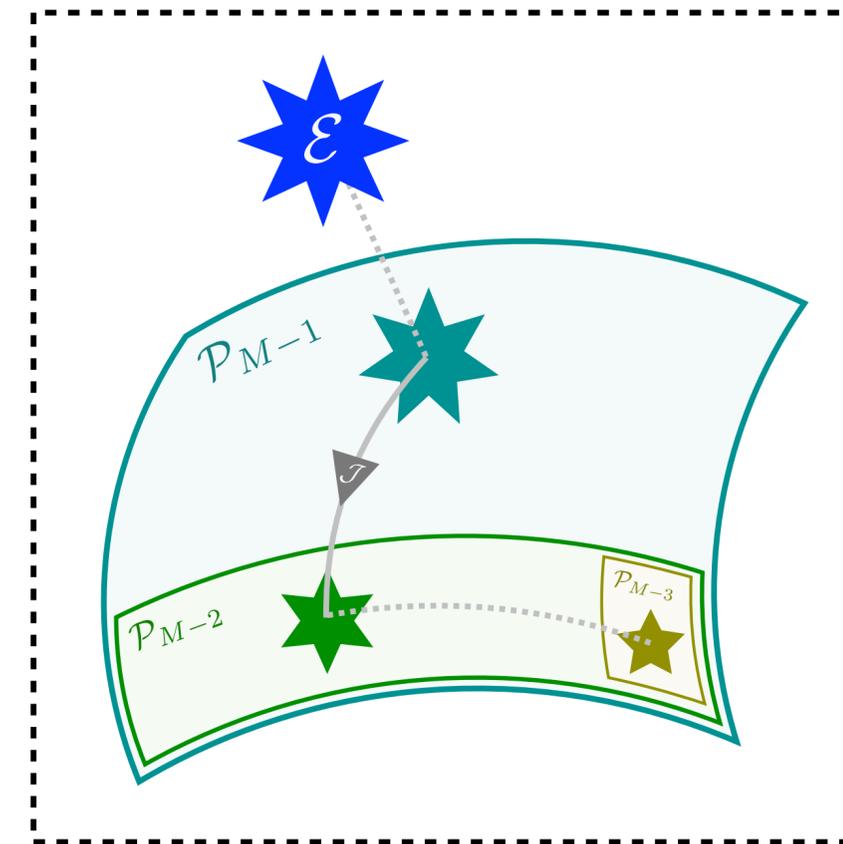


$$\mathcal{J}_{N,\beta,R}^{\text{XCone}}(\mathcal{E}) = \arg \min_{\mathcal{J} \in \mathcal{P}_N} \text{EMD}_{\beta,R}(\mathcal{E}, \mathcal{J})$$

[Stewart, Tackmann, Thaler, Vermilion, Wilkason, [JHEP 2015](#);  
Thaler, Wilkason, [JHEP 2015](#)]

## Sequential recombination

Iteratively merges particles or identifies a jet



“destroying” energy corresponds to identifying a jet

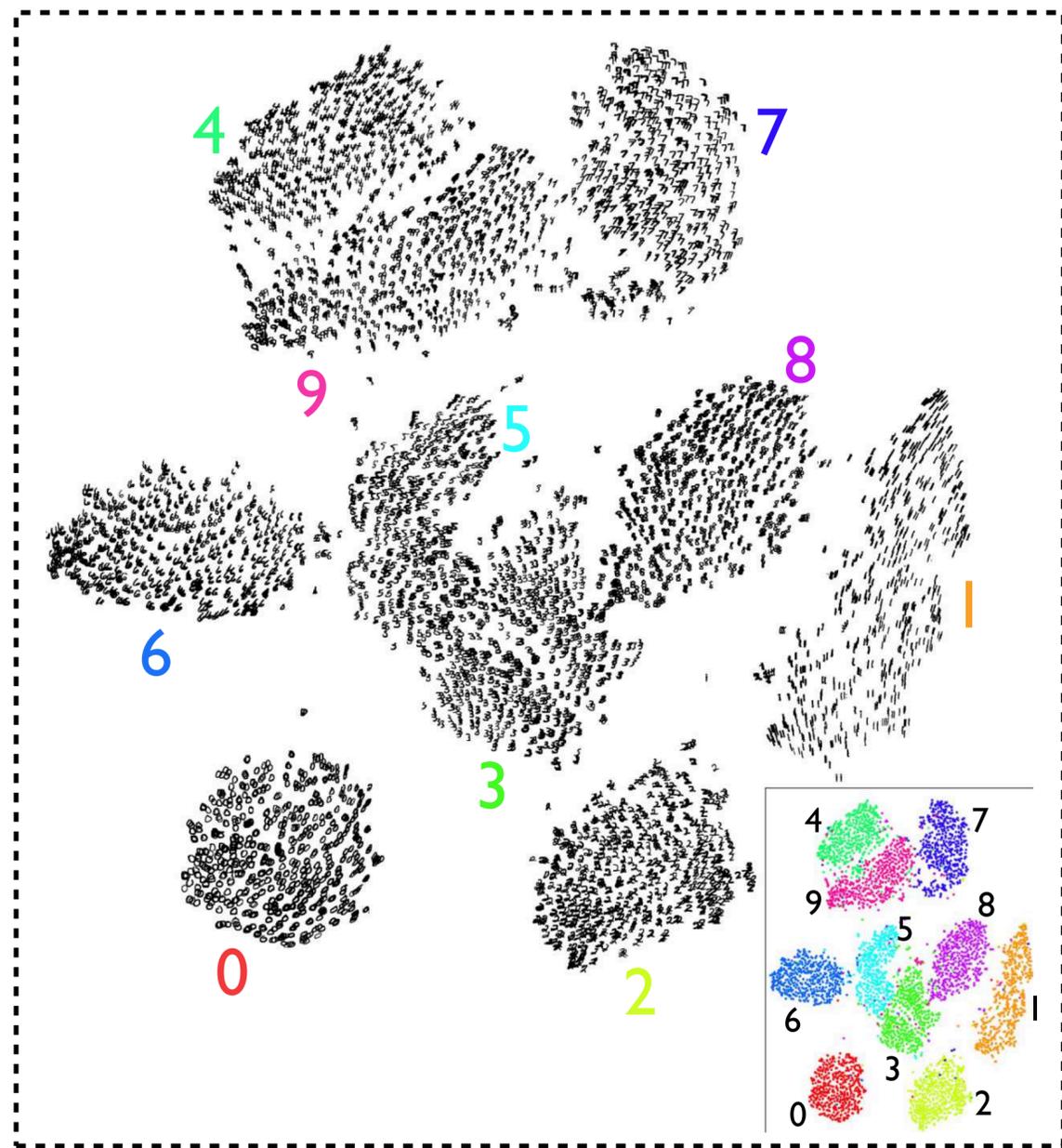
event with one fewer particle after one step

$$\mathcal{E}_{M-1}^{(\beta,R)}(\mathcal{E}_M) = \arg \min_{\mathcal{E}'_{M-1} \in \mathcal{P}_{M-1}} \text{EMD}_{\beta,R}(\mathcal{E}_M, \mathcal{E}'_{M-1})$$

[Catani, Dokshitzer, Seymour, Webber, [Nucl. Phys. B 1993](#);  
Ellis, Soper, [PRD 1993](#);  
Dokshitzer, Leder, Moretti, Webber, [JHEP 1997](#);  
Cacciari, Salam, Soyez, [JHEP 2008](#)]

# Visualizing Geometry in the Space of Events

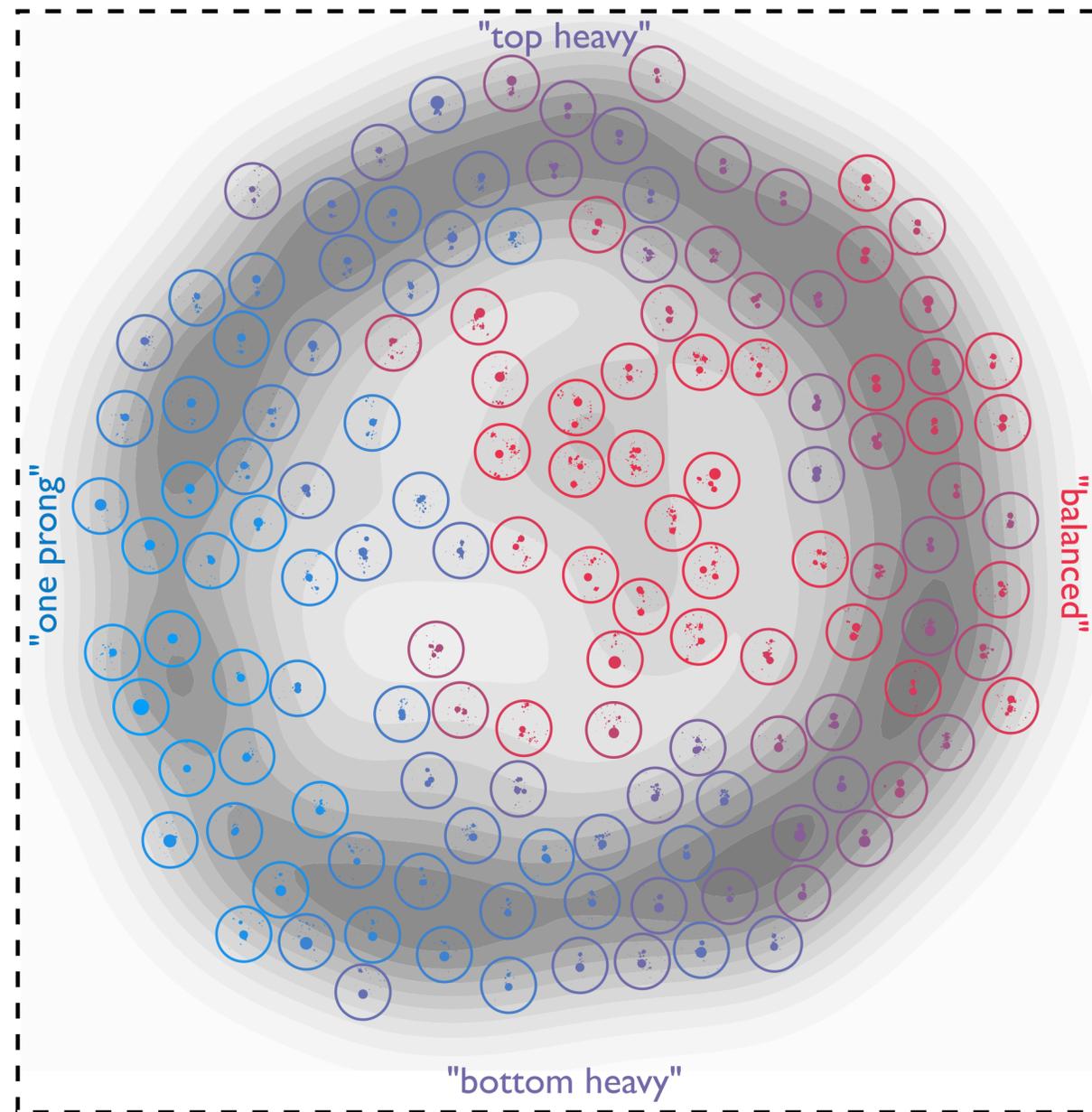
t-Distributed Stochastic Neighbor Embedding (t-SNE)  
MNIST handwritten digits



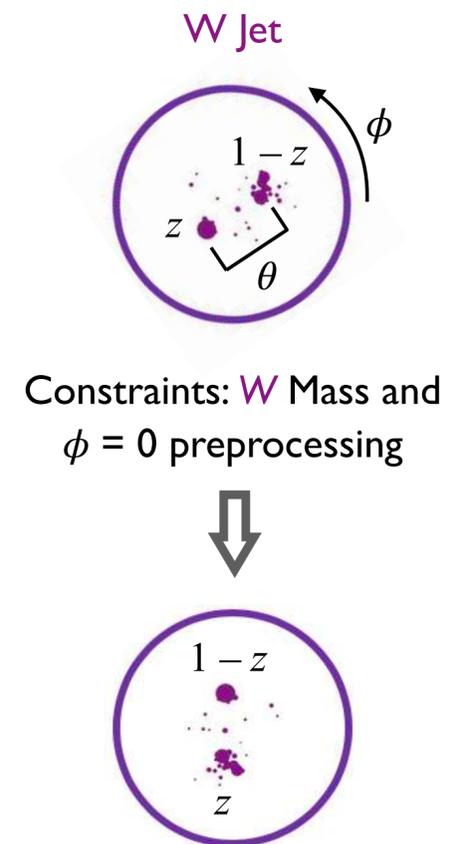
[L. van der Maaten, G. Hinton, JMLR 2008 ]

[PTK, Metodiev, Thaler, PRL 2019]

Geometric space of  $W$  jets

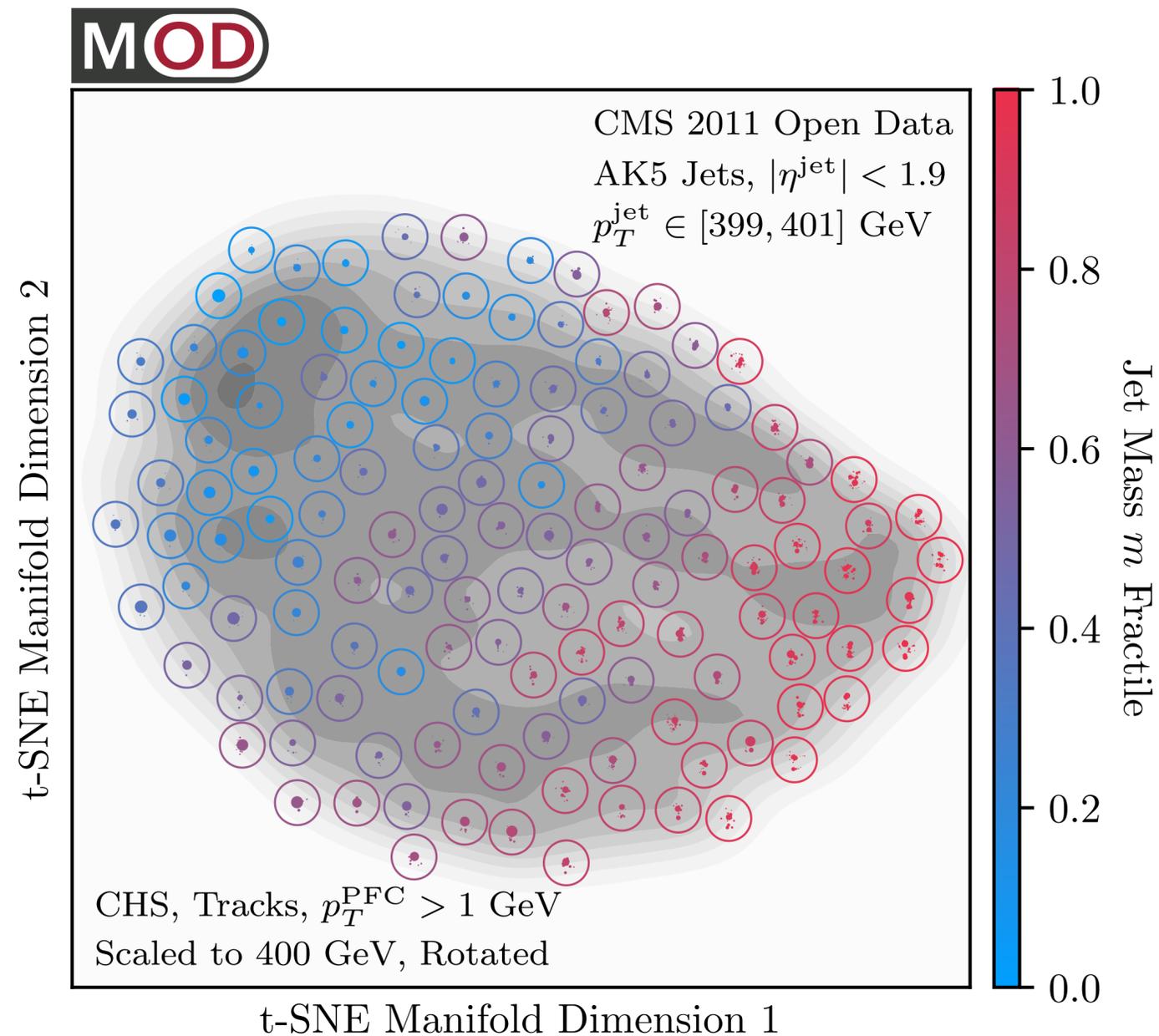


Gray contours represent the density of jets  
Each circle is a particular  $W$  jet

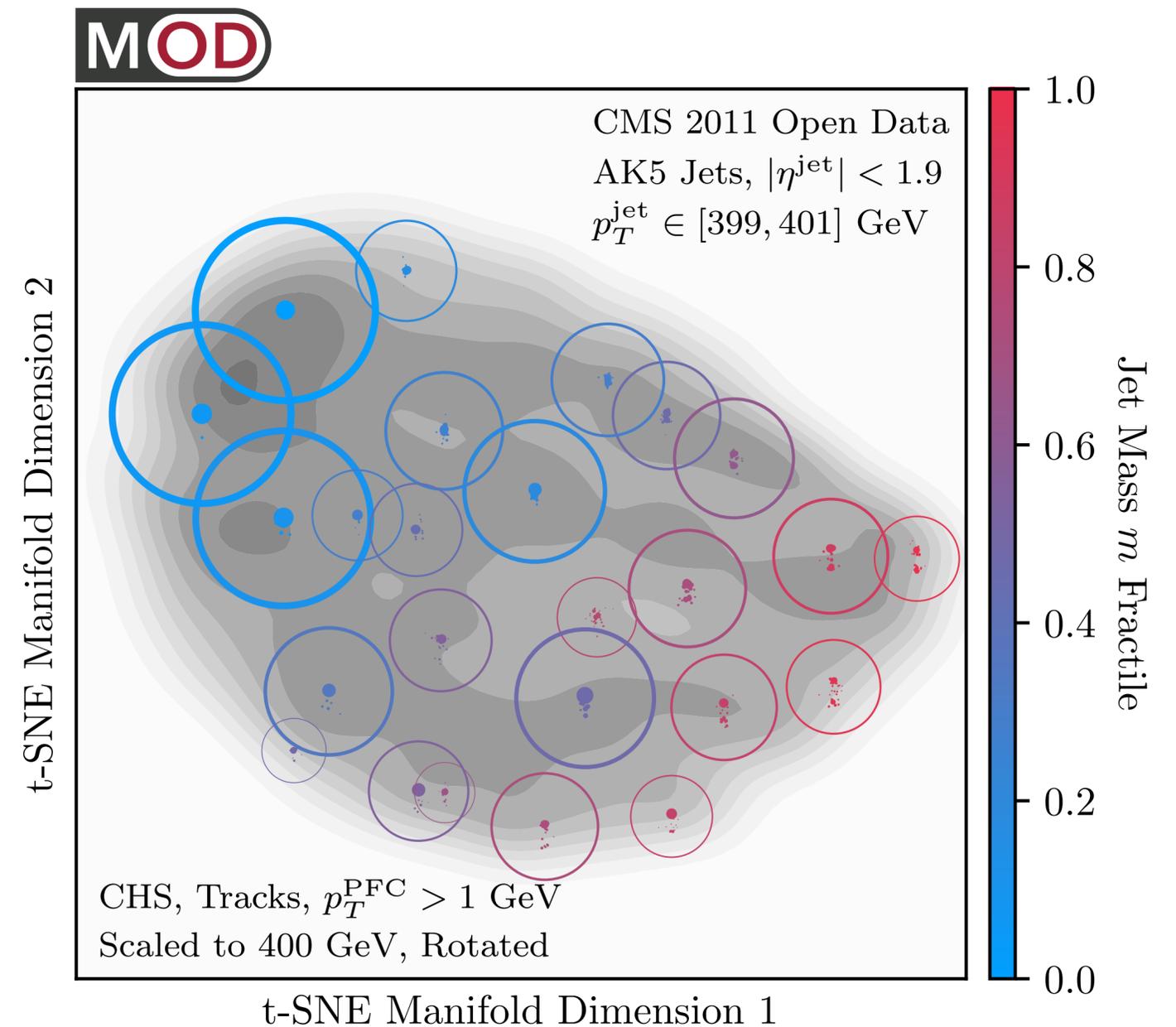


# Visualizing Geometry in CMS Open Data

[PTK, Mastandrea, Metodiev, Naik, Thaler, [PRD 2019](#); code and datasets at [energyflow.network](#)]



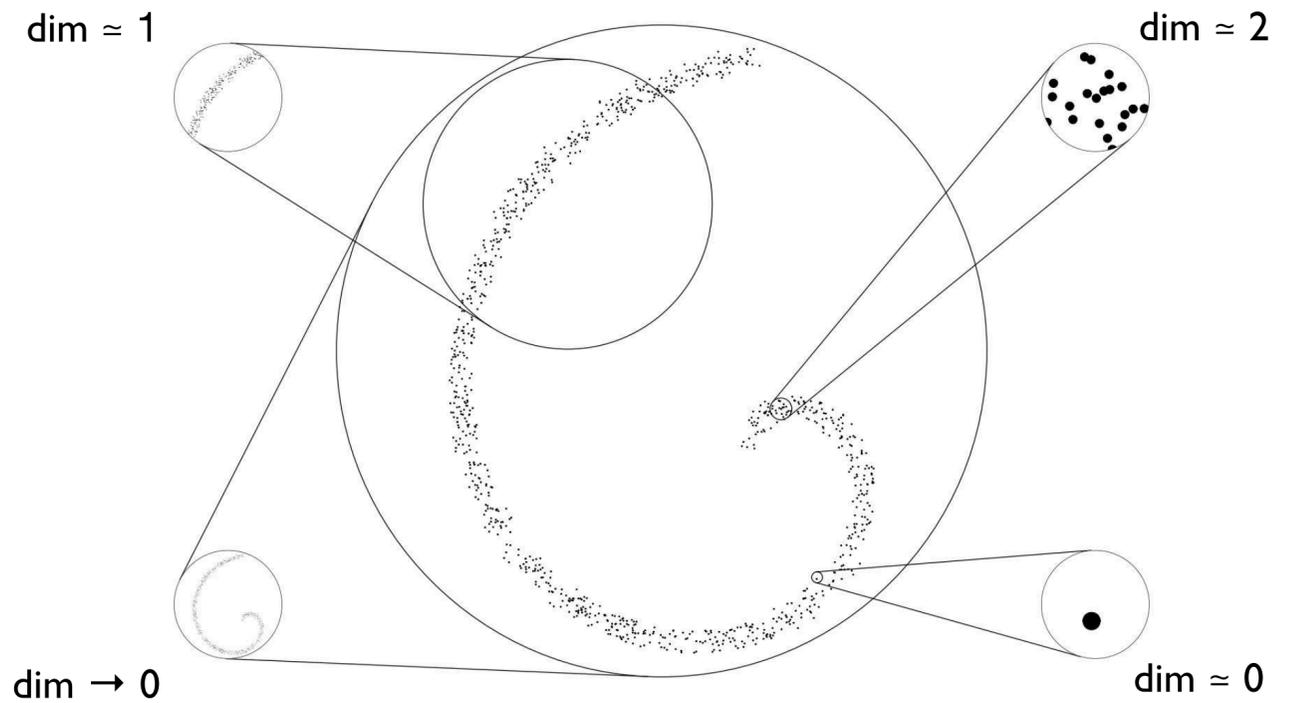
Example jets sprinkled throughout



25 most representative jets (“medoids”)  
Size is proportional to number of jets associated to that medoid

# Quantifying Event-Space Manifolds

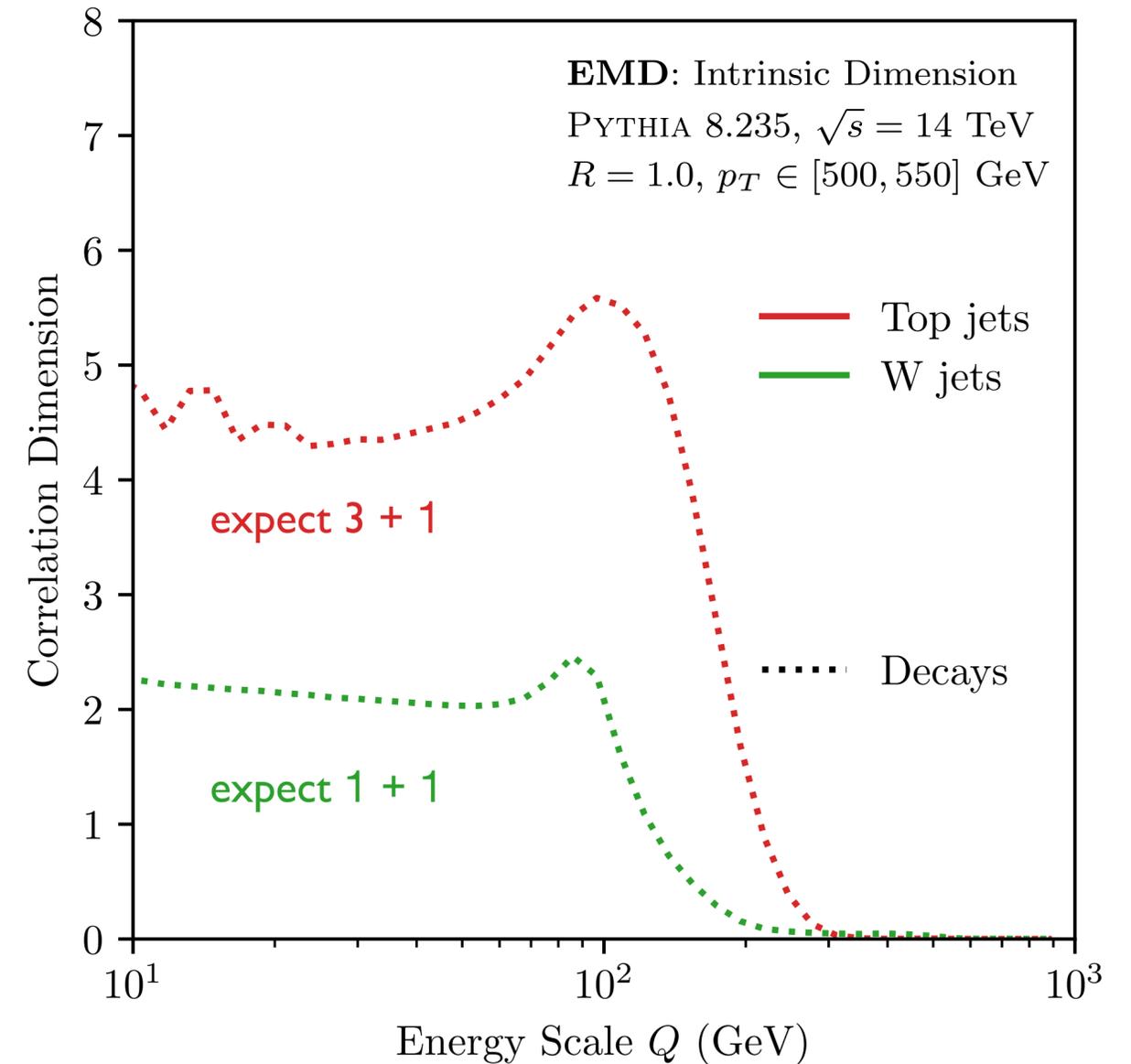
Correlation dimension: *how does the # of elements within a ball of size  $Q$  change?*



$$N_{\text{neigh.}}(Q) \propto Q^{\text{dim}} \implies \text{dim}(Q) = Q \frac{d}{dQ} \ln N_{\text{neigh.}}(Q)$$

Correlation dimension lessons:  
Decays are "constant" dim. at low  $Q$

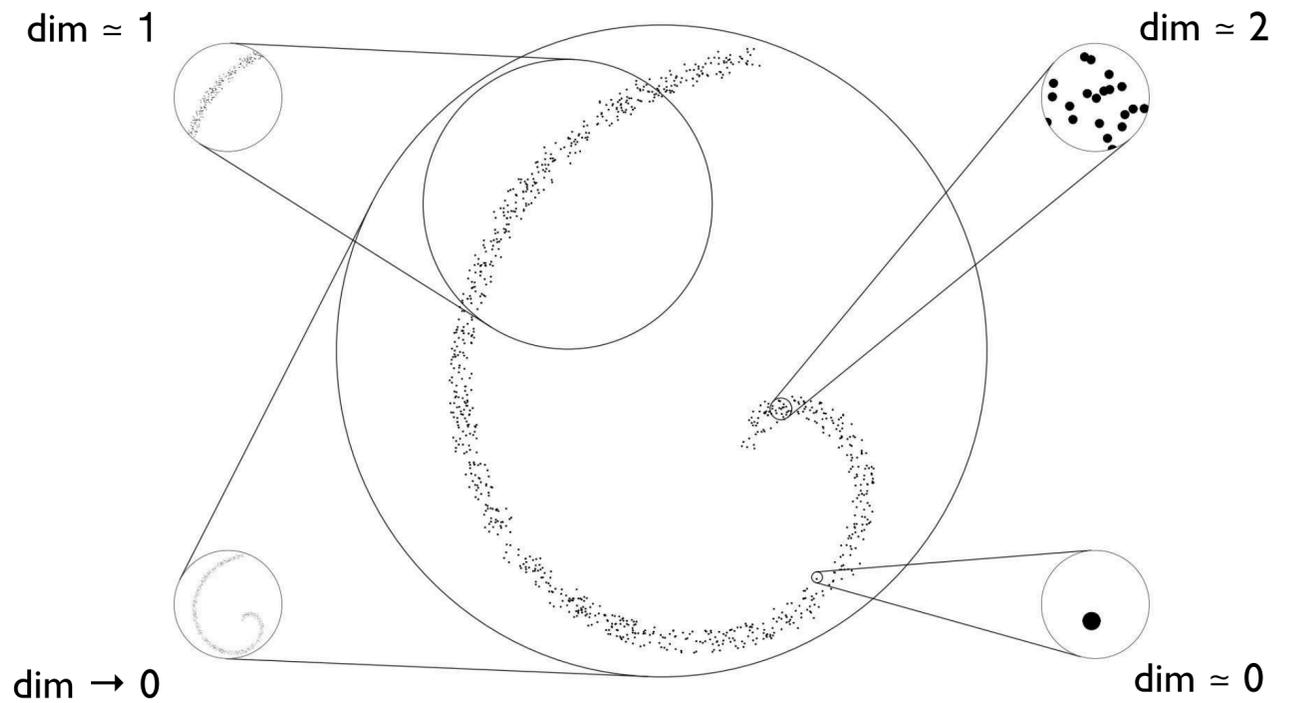
$$\text{dim}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$



[Grassberger, Procaccia, PRL 1983; PTK, Metodiev, Thaler, PRL 2019]

# Quantifying Event-Space Manifolds

Correlation dimension: *how does the # of elements within a ball of size  $Q$  change?*



$$N_{\text{neigh.}}(Q) \propto Q^{\text{dim}} \implies \text{dim}(Q) = Q \frac{d}{dQ} \ln N_{\text{neigh.}}(Q)$$

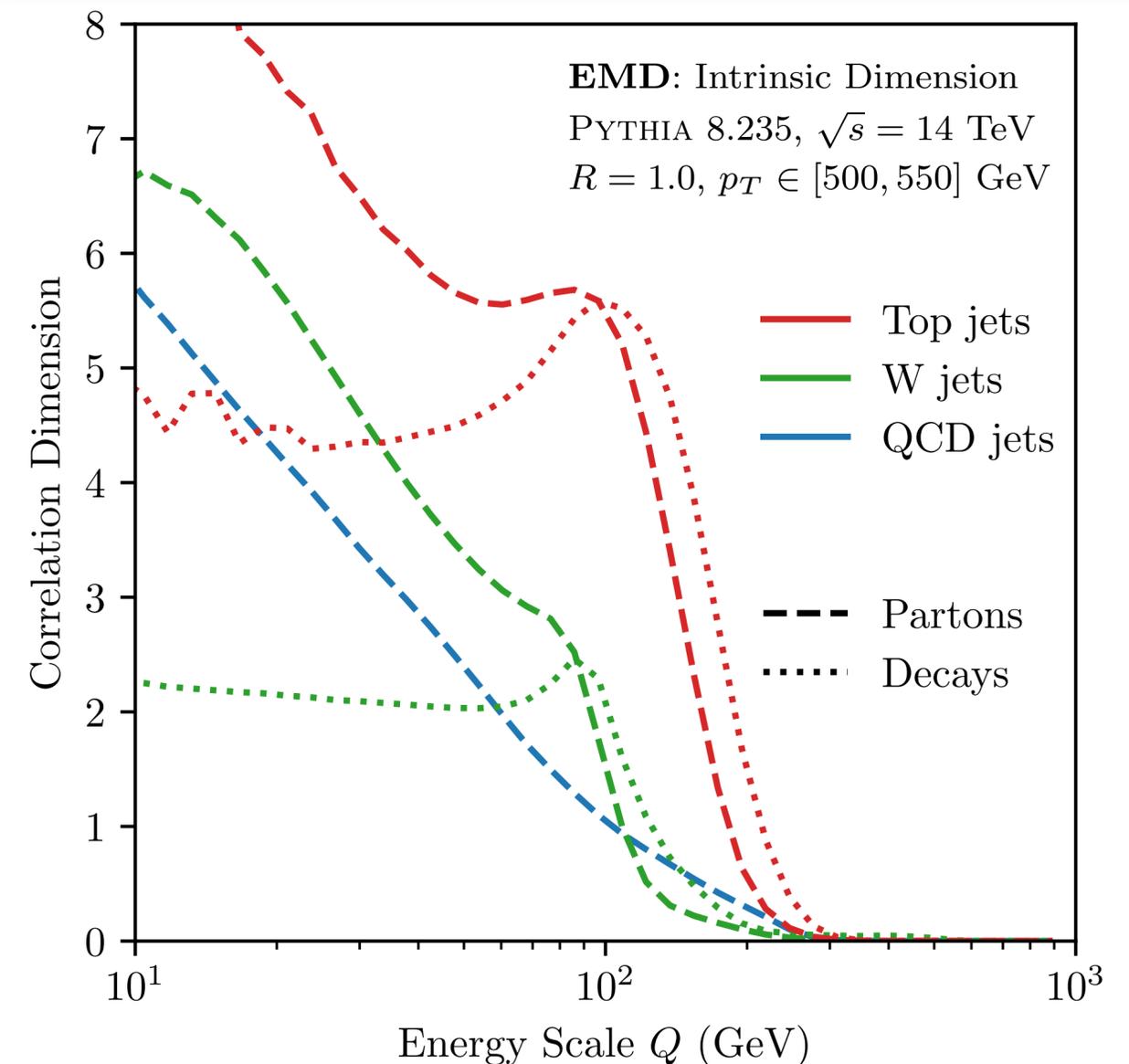
Correlation dimension lessons:

Decays are "constant" dim. at low  $Q$

Complexity hierarchy: QCD < W < Top

Fragmentation increases dim. at smaller scales

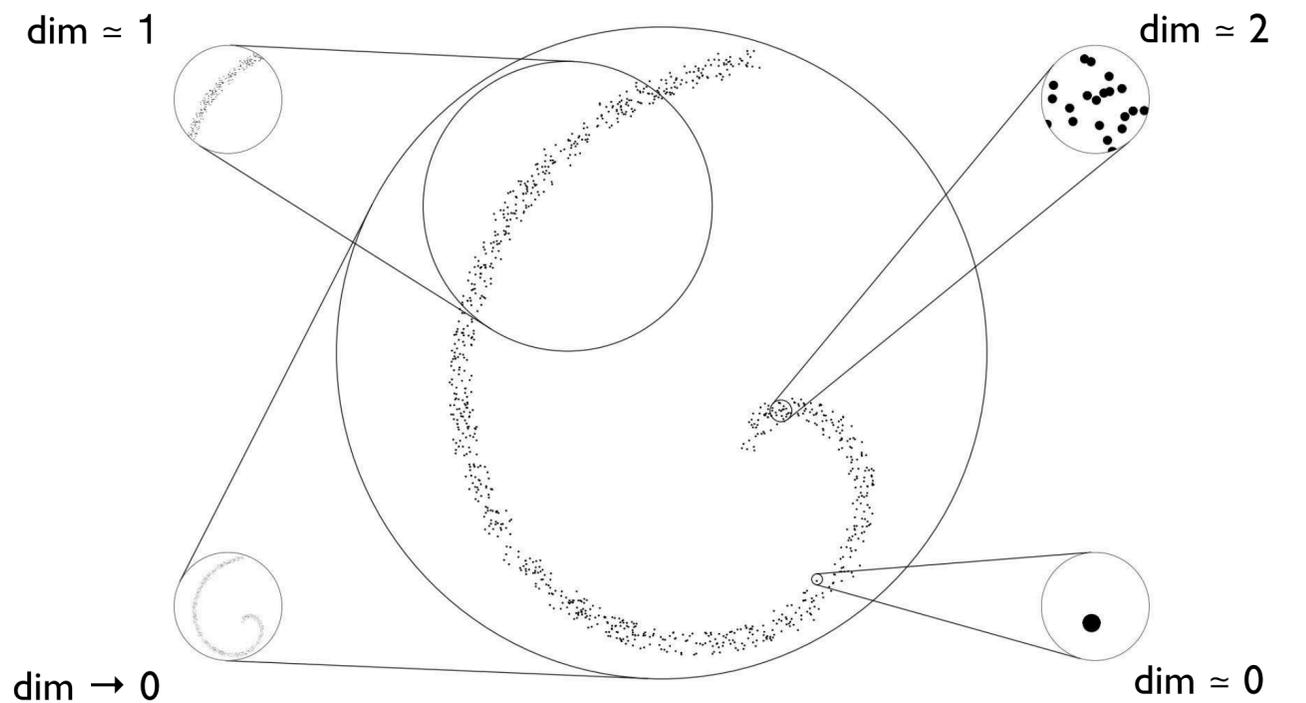
$$\text{dim}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$



[Grassberger, Procaccia, PRL 1983; PTK, Metodiev, Thaler, PRL 2019]

# Quantifying Event-Space Manifolds

Correlation dimension: *how does the # of elements within a ball of size  $Q$  change?*

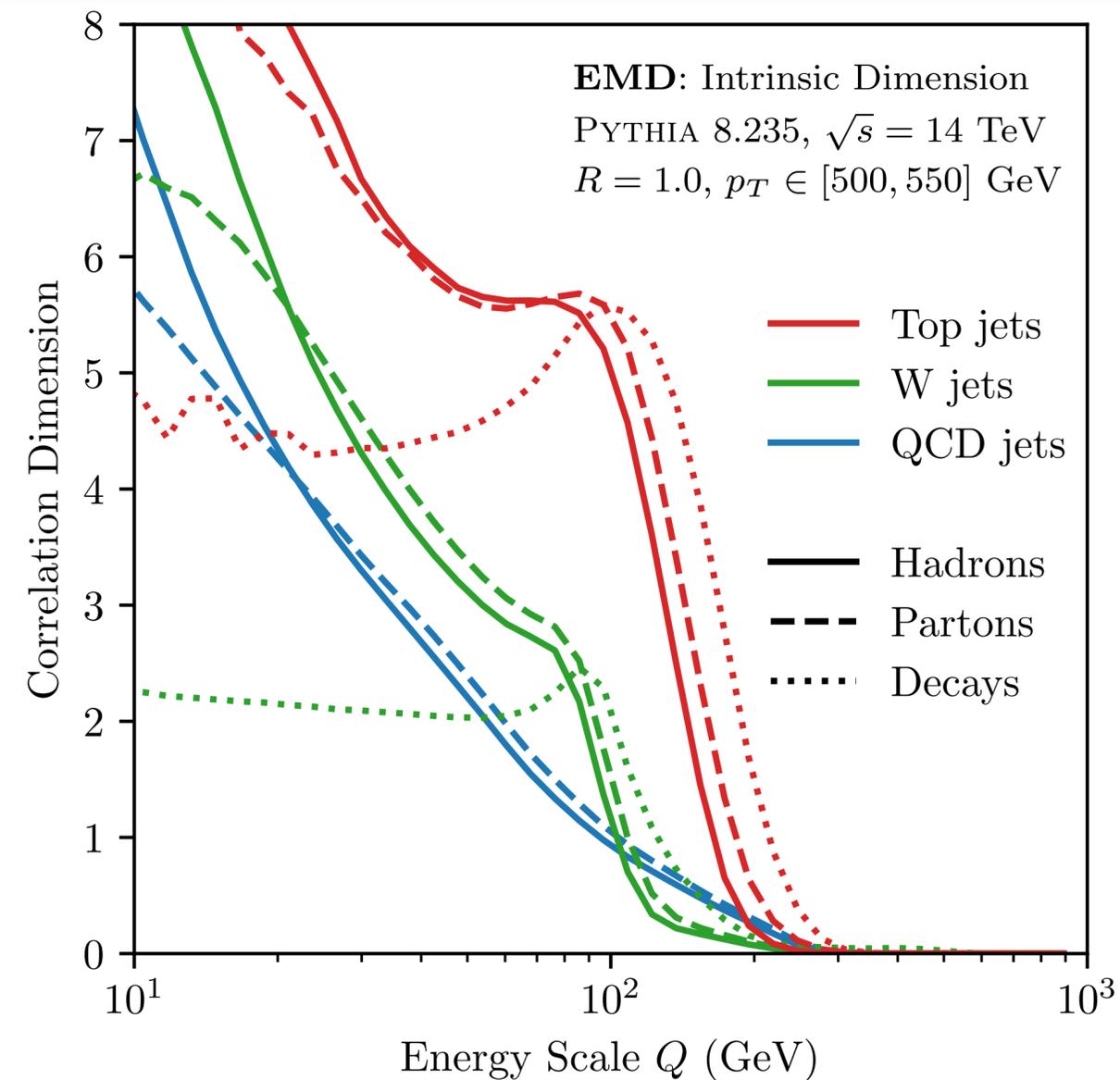


$$N_{\text{neigh.}}(Q) \propto Q^{\text{dim}} \implies \text{dim}(Q) = Q \frac{d}{dQ} \ln N_{\text{neigh.}}(Q)$$

**Correlation dimension lessons:**

- Decays are "constant" dim. at low  $Q$
- Complexity hierarchy: QCD < W < Top
- Fragmentation increases dim. at smaller scales
- Hadronization important around 20-30 GeV

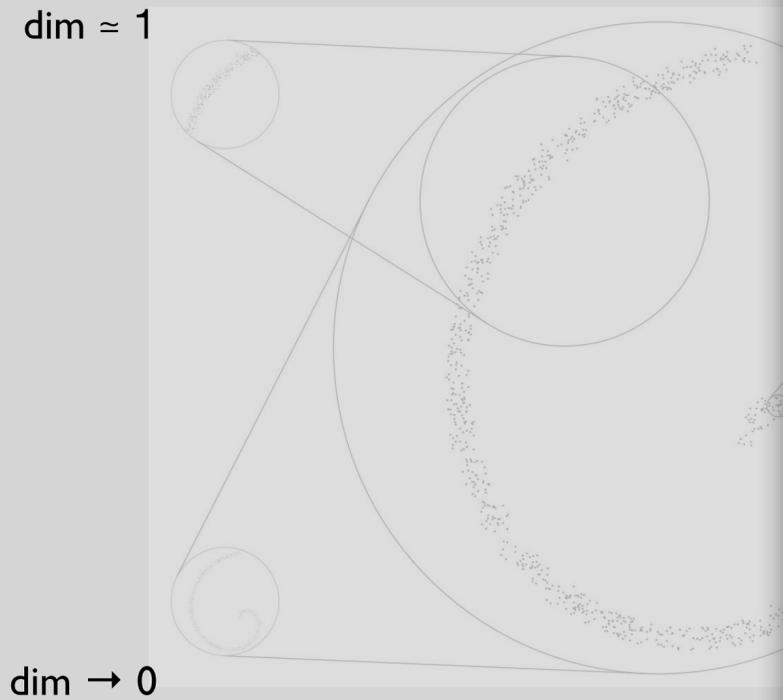
$$\text{dim}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$



[Grassberger, Procaccia, PRL 1983; PTK, Metodiev, Thaler, PRL 2019]

# Quantifying Event-Space Manifolds

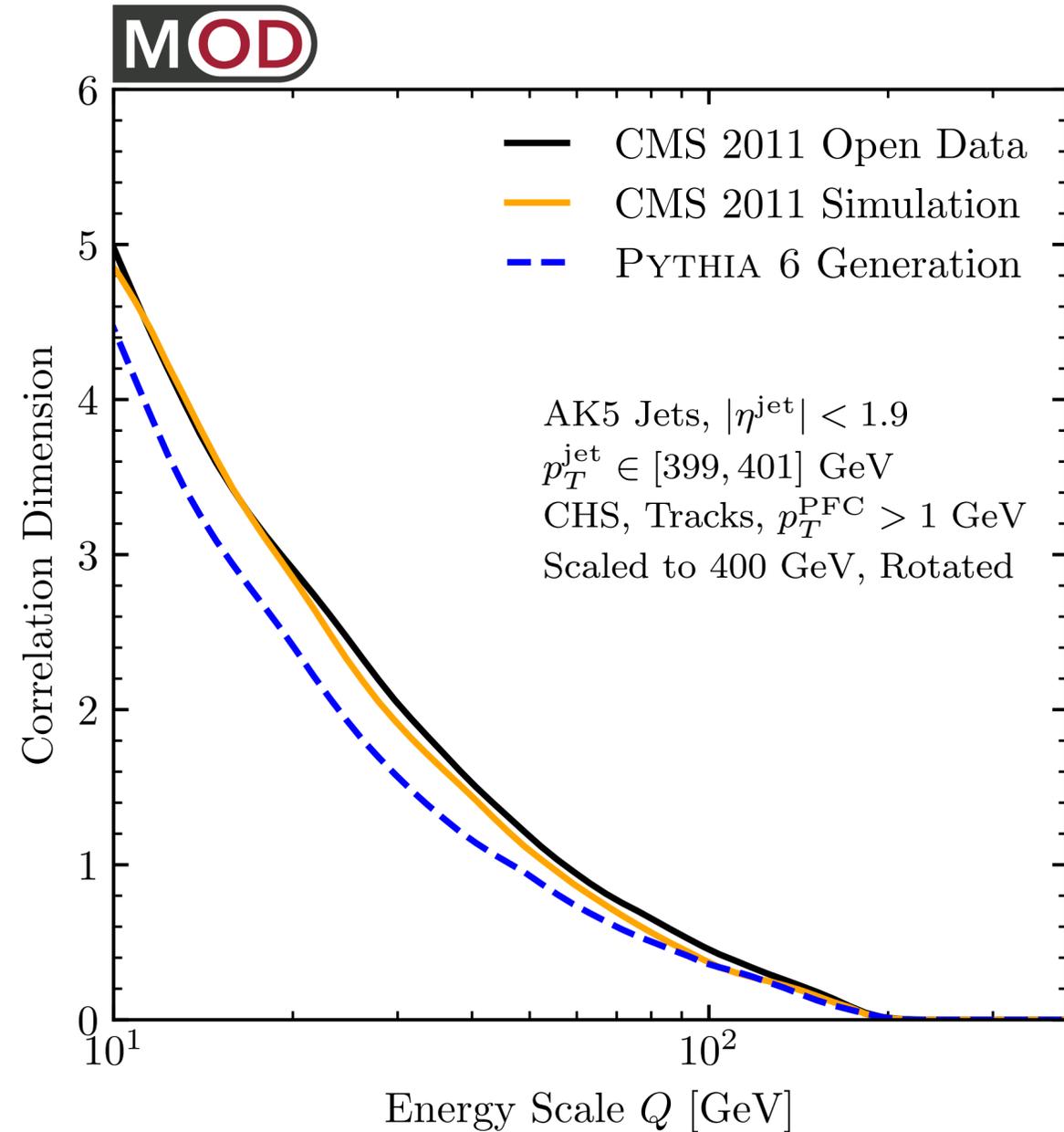
Correlation dimension: how many elements within a ball of size  $Q$



$$N_{\text{neigh.}}(Q) \propto Q^{\text{dim}} \implies \text{dim}(Q)$$

Correlation dimension  
 Decays are "constant" dimension  
 Complexity hierarchy: QCD jets > Partons > Hadrons  
 Fragmentation increases dimension  
 Hadronization important

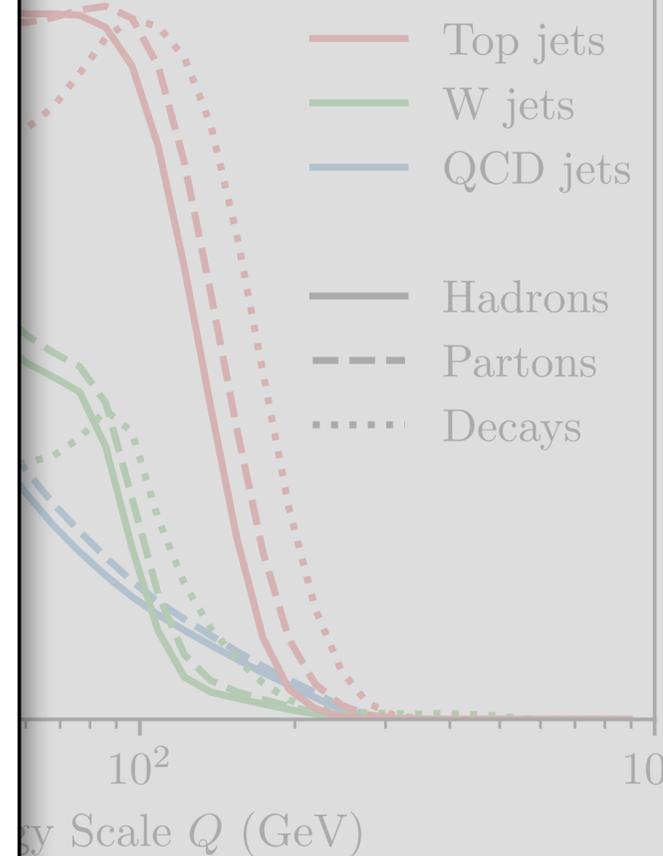
... in CMS Open Data



\*More in backup

$$\sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$

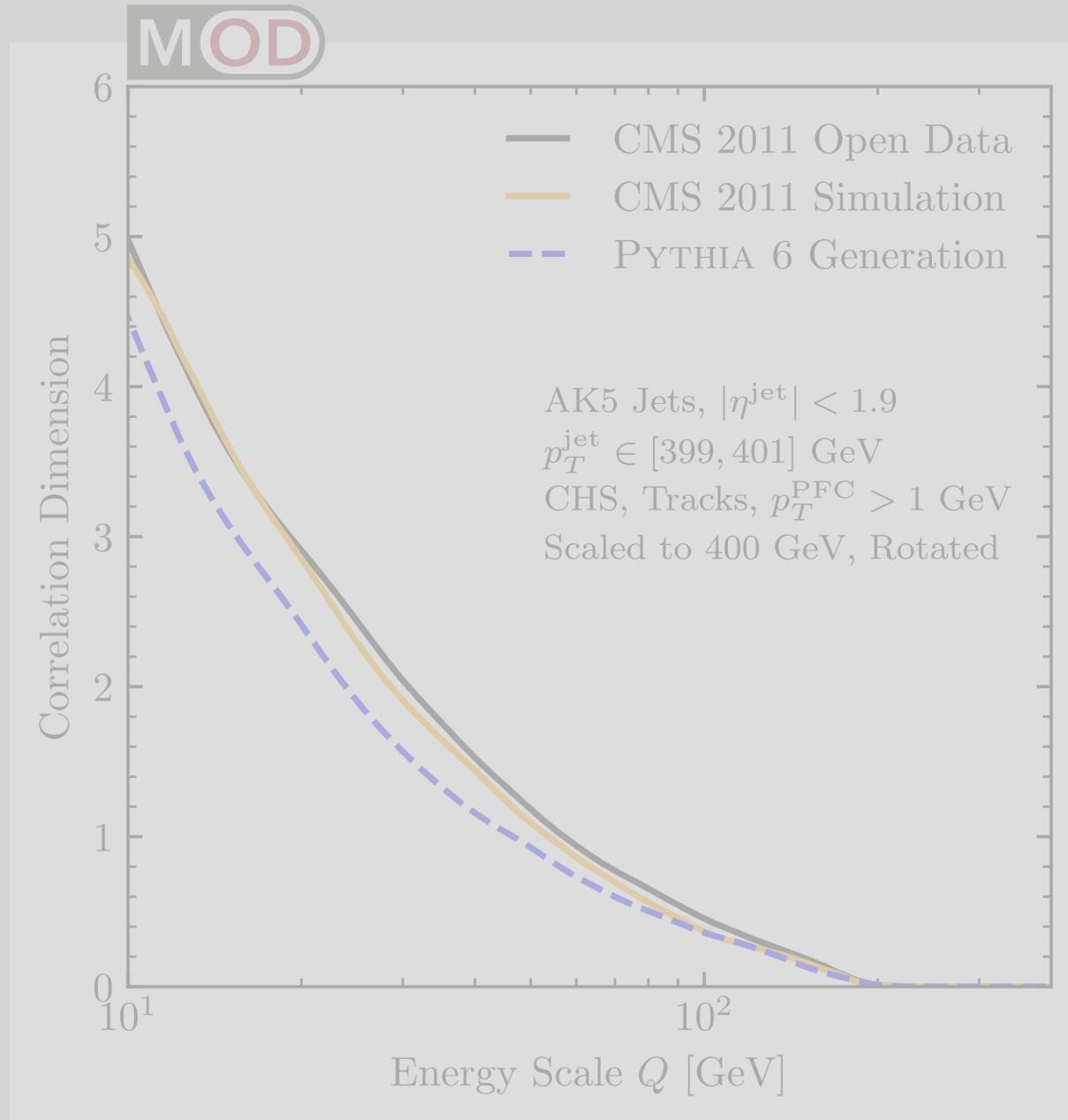
EMD: Intrinsic Dimension  
 PYTHIA 8.235,  $\sqrt{s} = 14$  TeV  
 $R = 1.0, p_T \in [500, 550]$  GeV



Procaccia, PRL 1983; PTK, Metodiev, Thaler, PRL 2019]

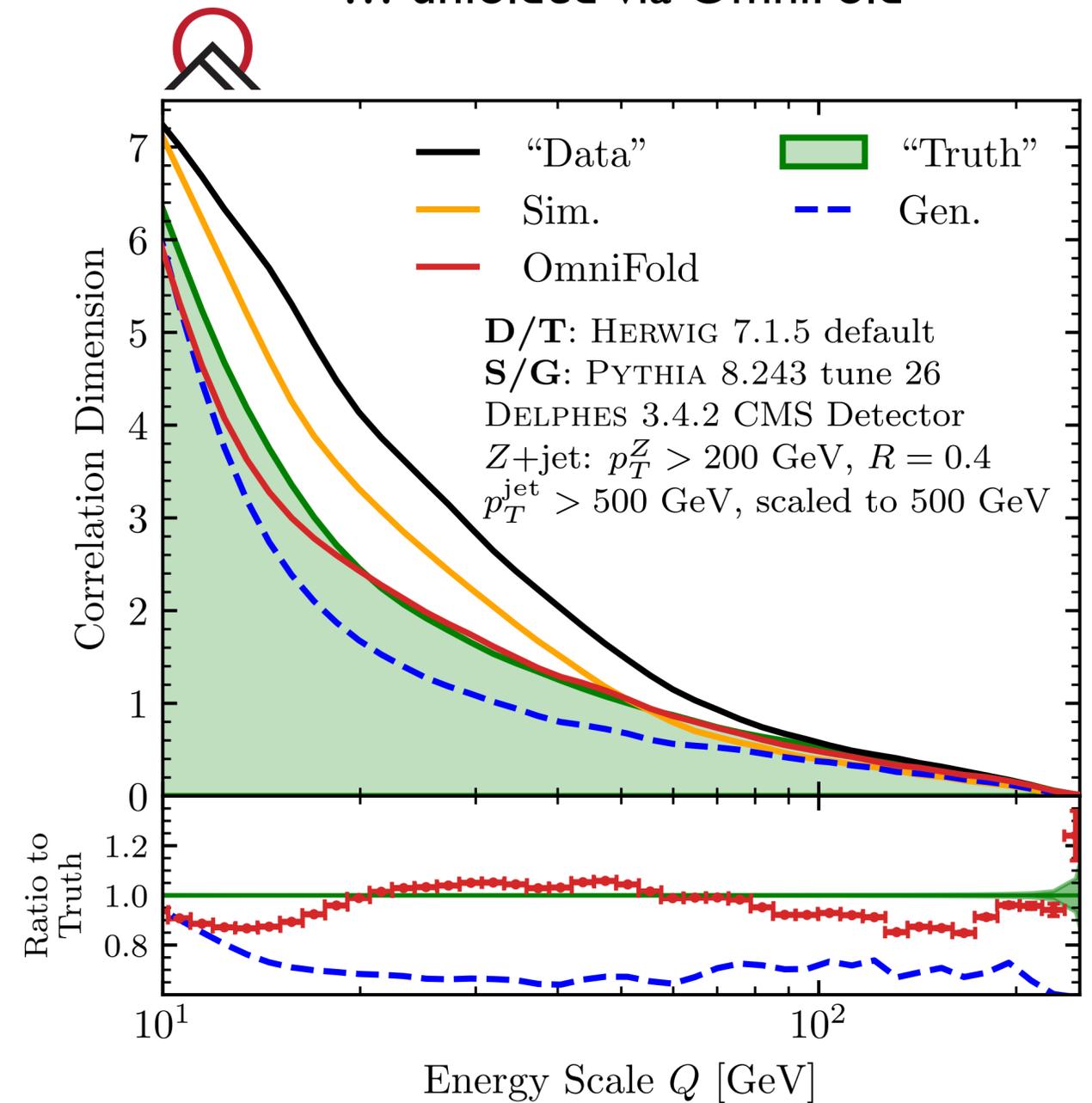
# Quantifying Event-Space Manifolds

... in CMS Open Data

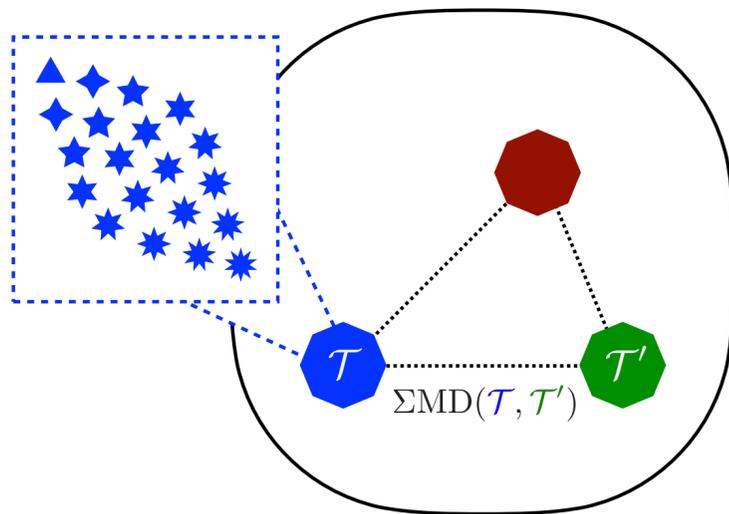
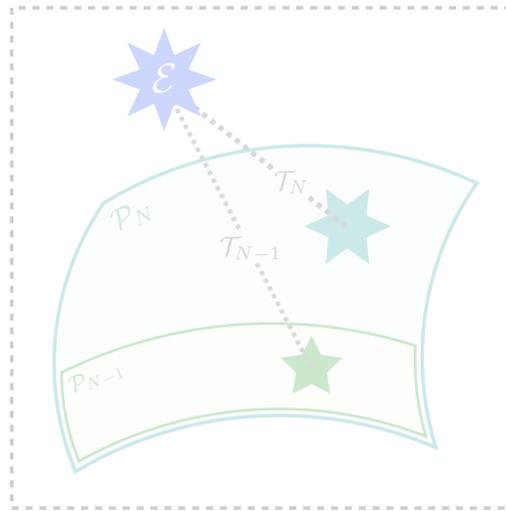
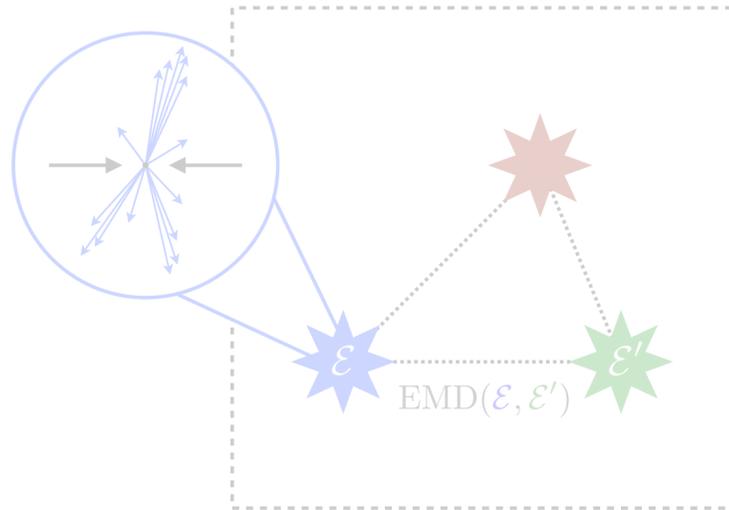


\*More in backup

... unfolded via OmniFold



Herwig as “data” and Pythia as MC



The (Metric) Space of Events

Revealing Hidden Geometry

[Theory Space]

# Templated Metric Construction

## Inputs

- “Points” that live in the ground metric space
- “Ground metric” that measures point distances
- “Weights” associated to each point

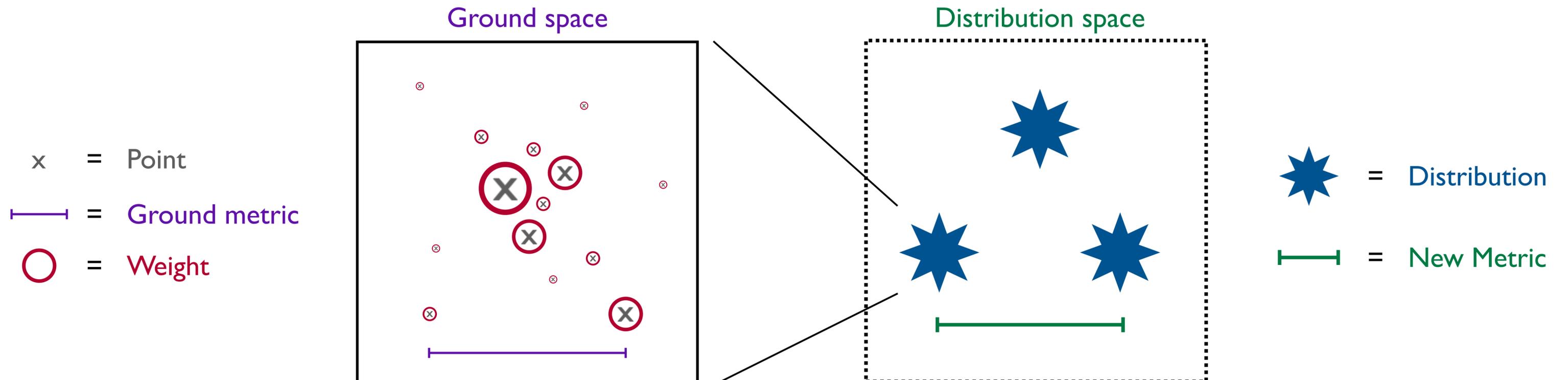
## Output

- A new metric for collections of weighted points
- A metric space where these distributions live

p-Wasserstein metric from optimal transport theory

$$W_p(\mu, \nu) = \left( \inf_{J \in \mathcal{J}(\mu, \nu)} \int_{M \times M} d(x, y)^p dJ(x, y) \right)^{1/p}$$

$(M, d)$ , metric space  
 $\mathcal{J}(\mu, \nu)$ , space of joint distributions with marginals  $\mu, \nu$



# Templated Metric Construction – Energy Mover’s Distance

[PTK, Metodiev, Thaler, [PRL 2019](#)]

## Inputs

- “Points” that live in the ground metric space
- “Ground metric” that measures point distances
- “Weights” associated to each point

## Output

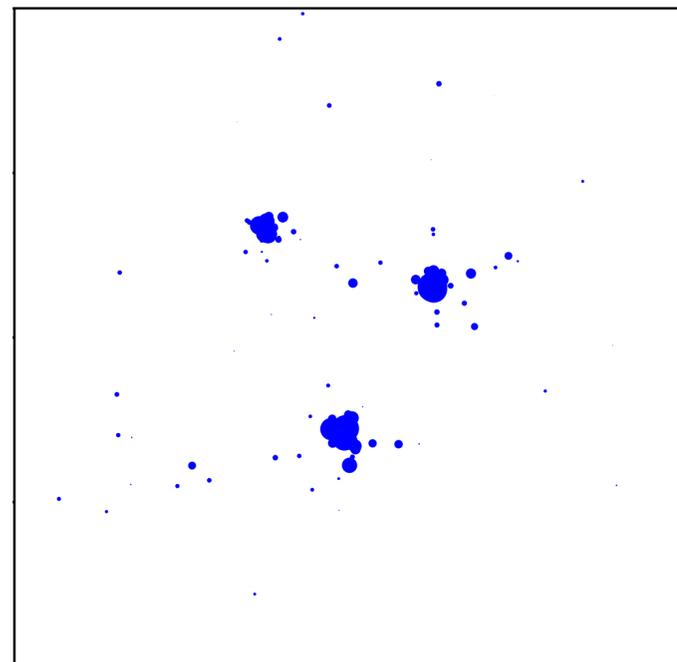
- A new metric for collections of weighted points
- A metric space where these distributions live

p-Wasserstein metric from optimal transport theory

$$W_p(\mu, \nu) = \left( \inf_{J \in \mathcal{J}(\mu, \nu)} \int_{M \times M} d(x, y)^p dJ(x, y) \right)^{1/p}$$

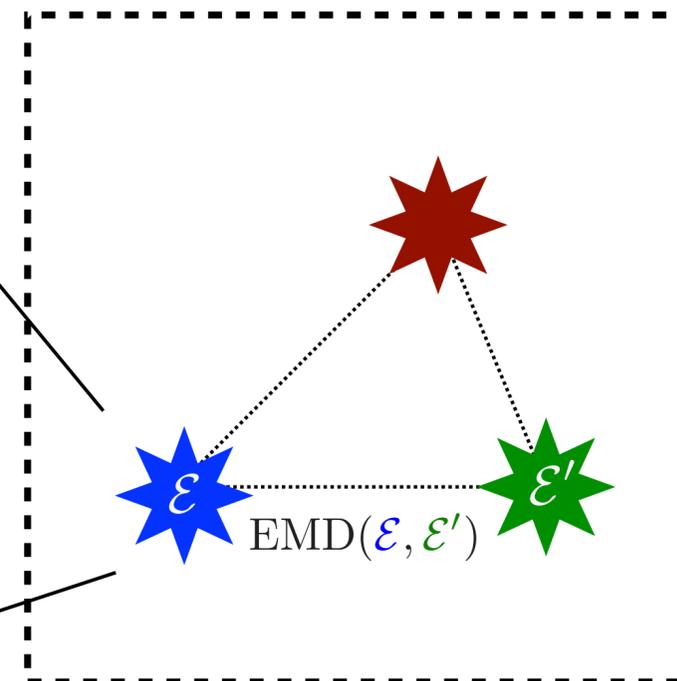
$(M, d)$ , metric space  
 $\mathcal{J}(\mu, \nu)$ , space of joint distributions with marginals  $\mu, \nu$

Rapidity-azimuth plane



- $x$  = Direction  $(y, \phi)$
- $\text{—|—}$  = Angular distance  $\Delta R$
- $\bigcirc$  = Energy

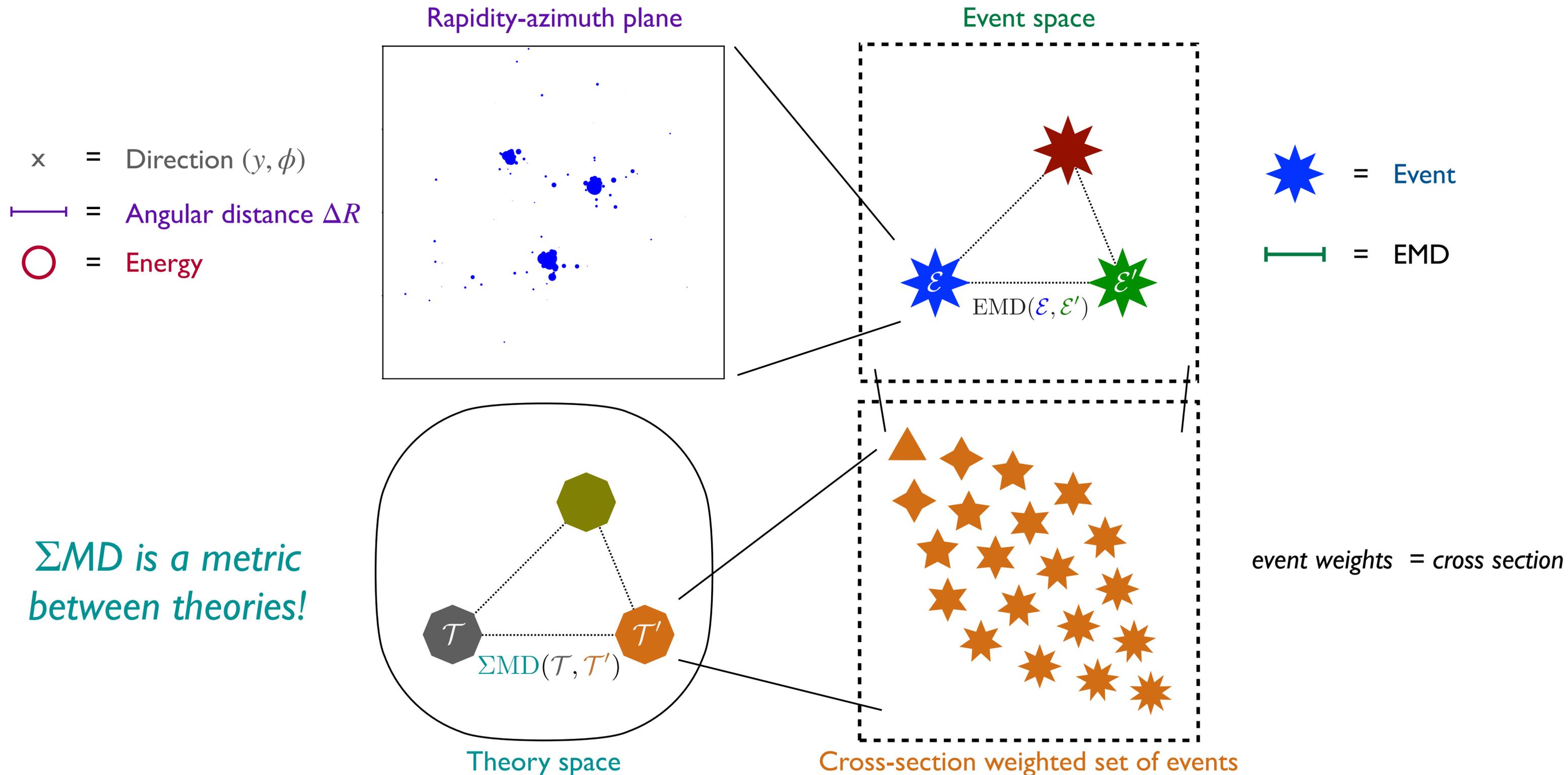
Event space



- $\star$  = Event
- $\text{—|—}$  = EMD

# Bootstrapping to the Cross-Section Mover's Distance ( $\Sigma MD$ )

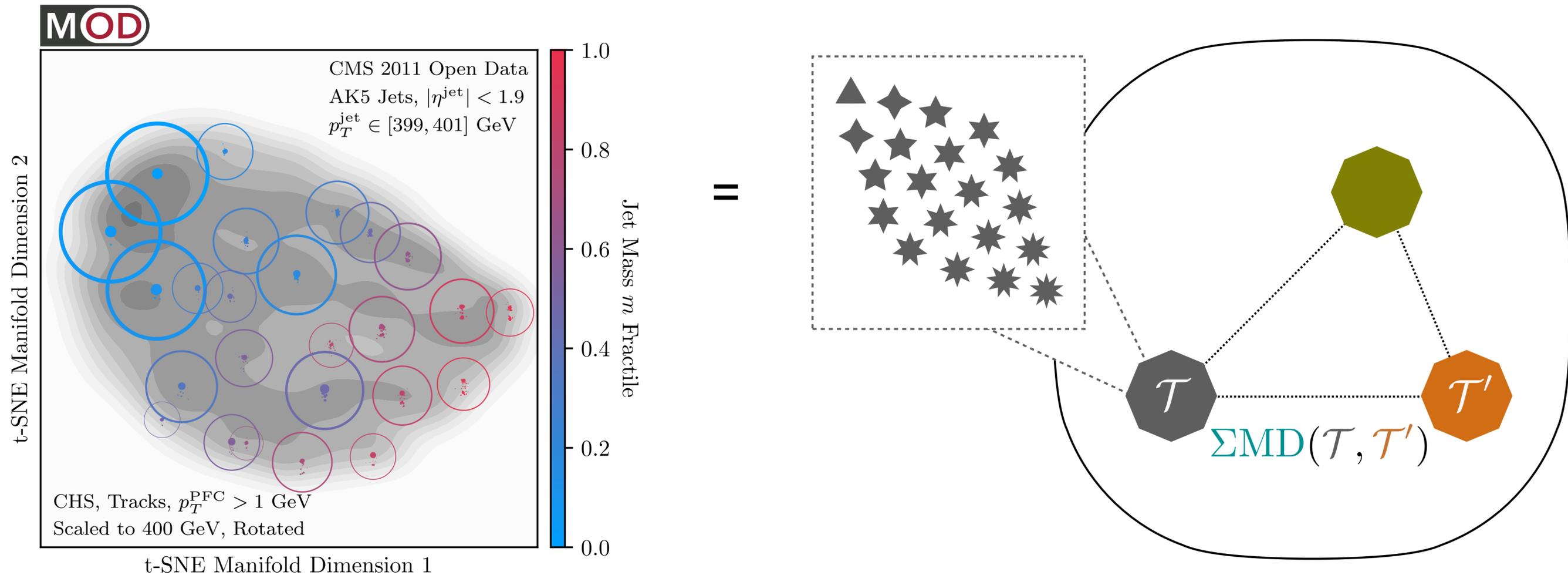
[PTK, Metodiev, Thaler, 2004.04159]



# The Space of Theories

[PTK, Metodiev, Thaler, 2004.04159]

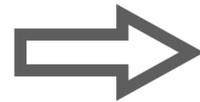
$\Sigma$ MMD provides a rigorous construction of theory space



\*Theories are distinguished by their **energy** flows only

# Applications of $\Sigma$ MMD and the Space of Theories

$N$ -(sub)jettiness

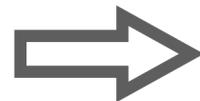


$k$ -eventiness defined

$$\mathcal{V}_k^{(\gamma)}(\{\sigma_i, \mathcal{E}_i\}) = \min_{\mathcal{K}_1, \dots, \mathcal{K}_k} \sum_{i=1}^N \sigma_i \min \{ \text{EMD}(\mathcal{E}_i, \mathcal{K}_1), \dots, \text{EMD}(\mathcal{E}_i, \mathcal{K}_k) \}^\gamma$$

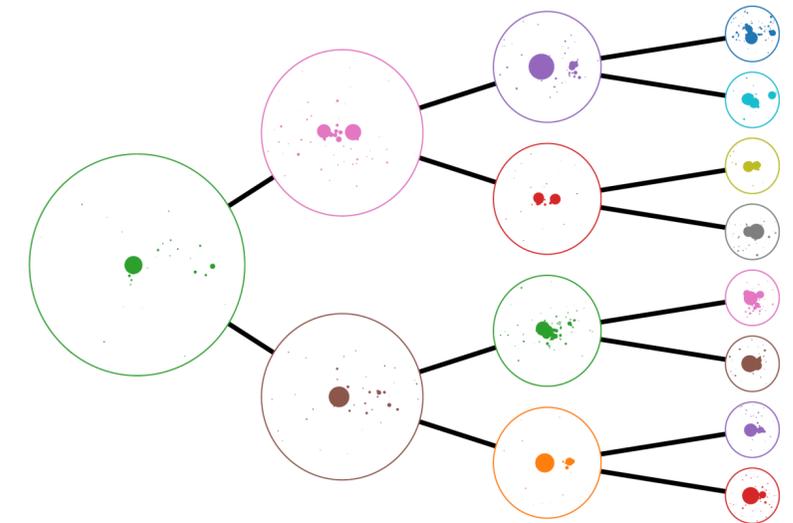
$$\mathcal{V}_k^{(\gamma)}(\mathcal{T}) = \min_{|\mathcal{T}'|=k} \Sigma\text{MMD}_\gamma(\mathcal{T}, \mathcal{T}')$$

Jet clustering



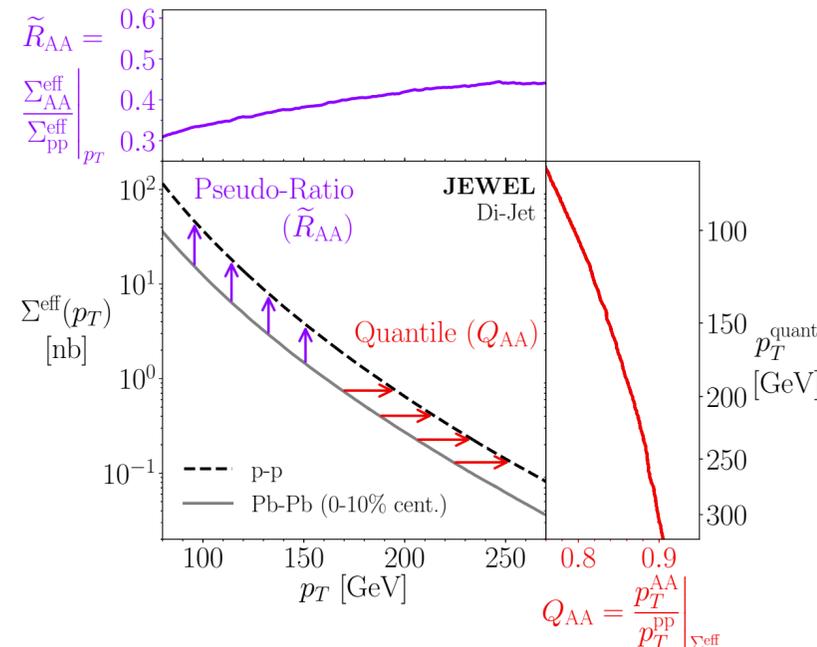
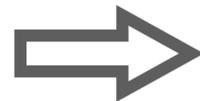
Event clustering enabled

- Exclusive cone finding
- Sequential recombination



Jet quenching in HI collisions

[Brewer, Milhano, Thaler, PRL 2019]



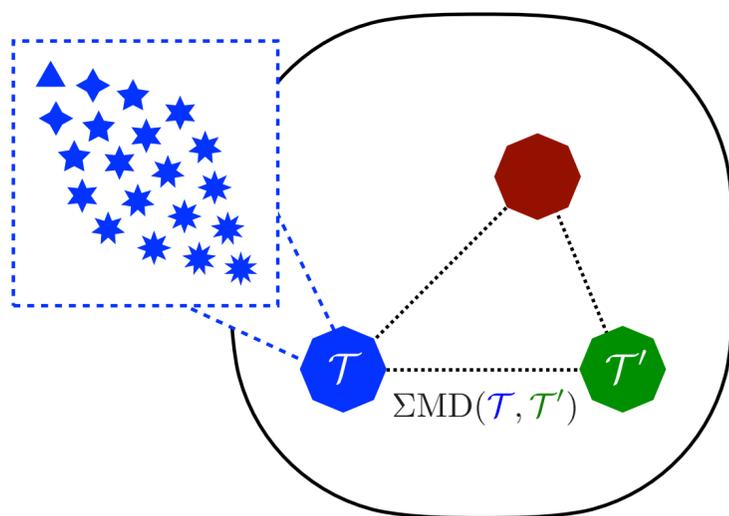
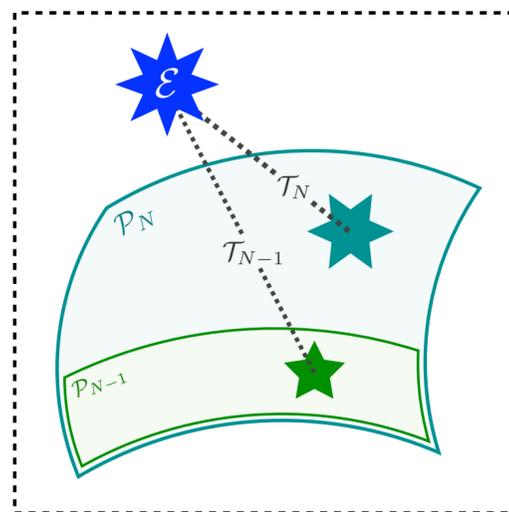
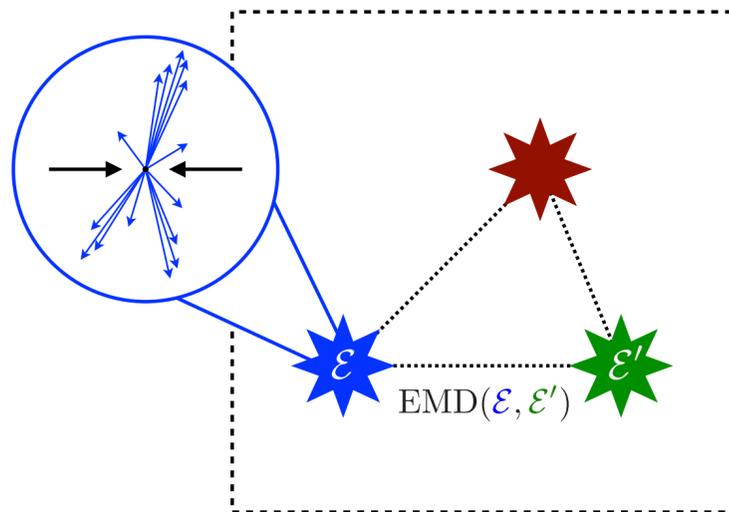
Quantile matching:

$$\Sigma_{pp}^{\text{eff}}(p_T^{\text{quant}}) \equiv \Sigma_{AA}^{\text{eff}}(p_T^{AA})$$

...is exactly a theory moving problem!

$$p_T^{\text{quant}} = \text{TM}(\mathcal{T}_{AA}, \mathcal{T}_{pp})[p_T^{AA}]$$

↑  
optimal  $p_T$ -only theory movement



## The (Metric) Space of Events

- Energy flow is theoretically and experimentally robust
- EMD metrizes the space of energy flows (events)
- Manifolds in the space of events can be visualized and quantified

## Revealing Hidden Geometry

- Event space exhibits a rich geometry that can be probed using the EMD
- Decades worth of collider techniques are naturally described in this geometry
- Many new techniques are suggested, and new light is shed on old ones

## [Theory Space]

- Is rigorously constructed using the cross-section mover's distance  $\Sigma MD$
- $\Sigma MD$  uses the EMD as ground metric and cross sections as weights
- Theories can be explored with tools developed for events (e.g.  $k$ -eventiness)

# EnergyFlow Python Package

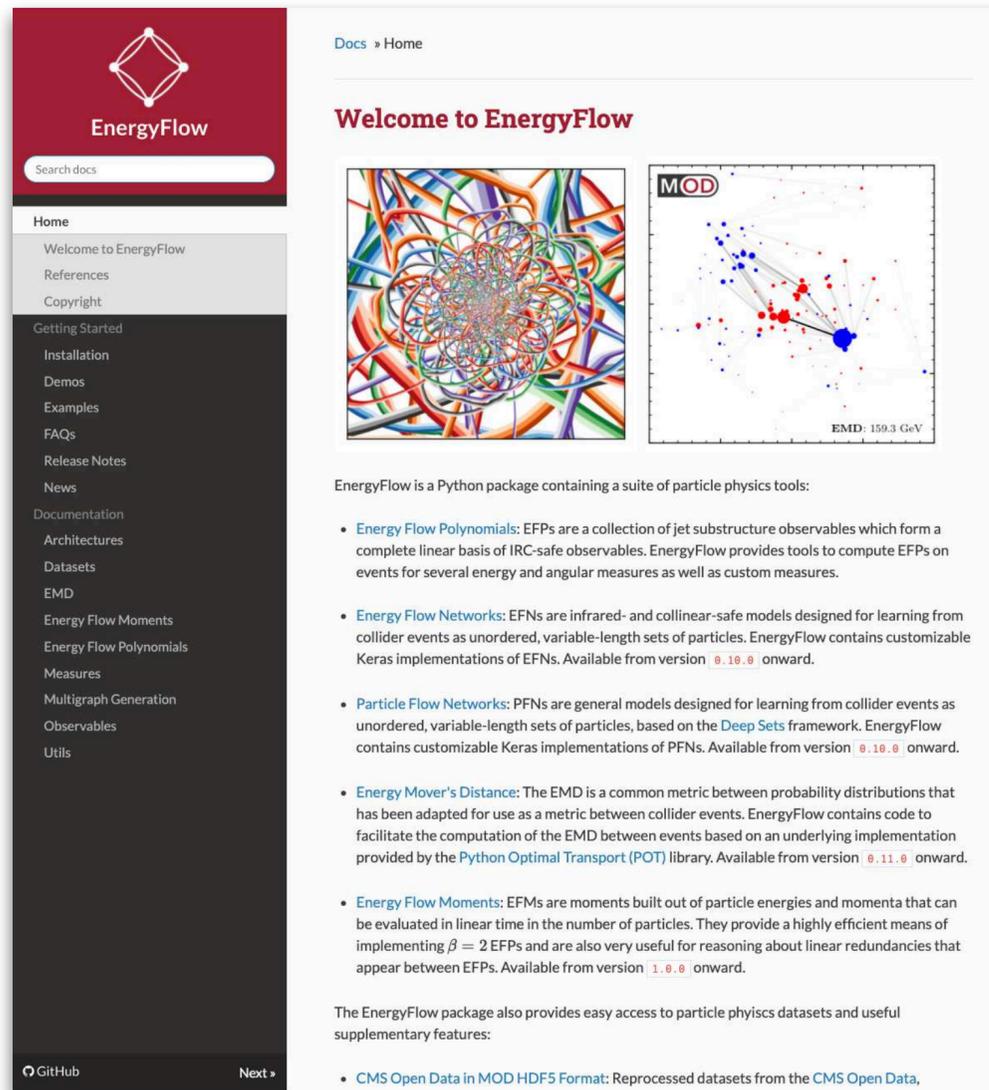
`pip3 install energyflow wasserstein`

Parallelized **EMD** calculations via the Wasserstein library

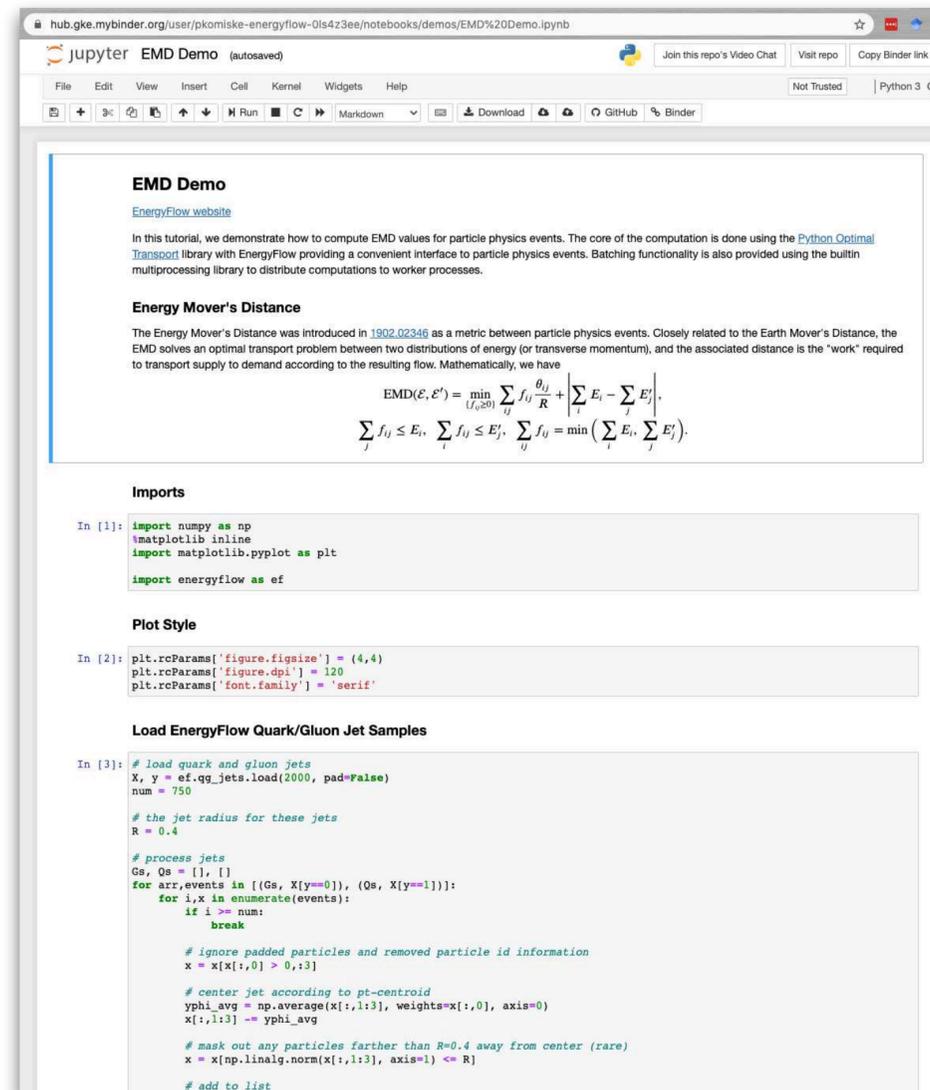
**EFN/PFN** implementations in TensorFlow/Keras

Detailed **examples, demos, and documentation**

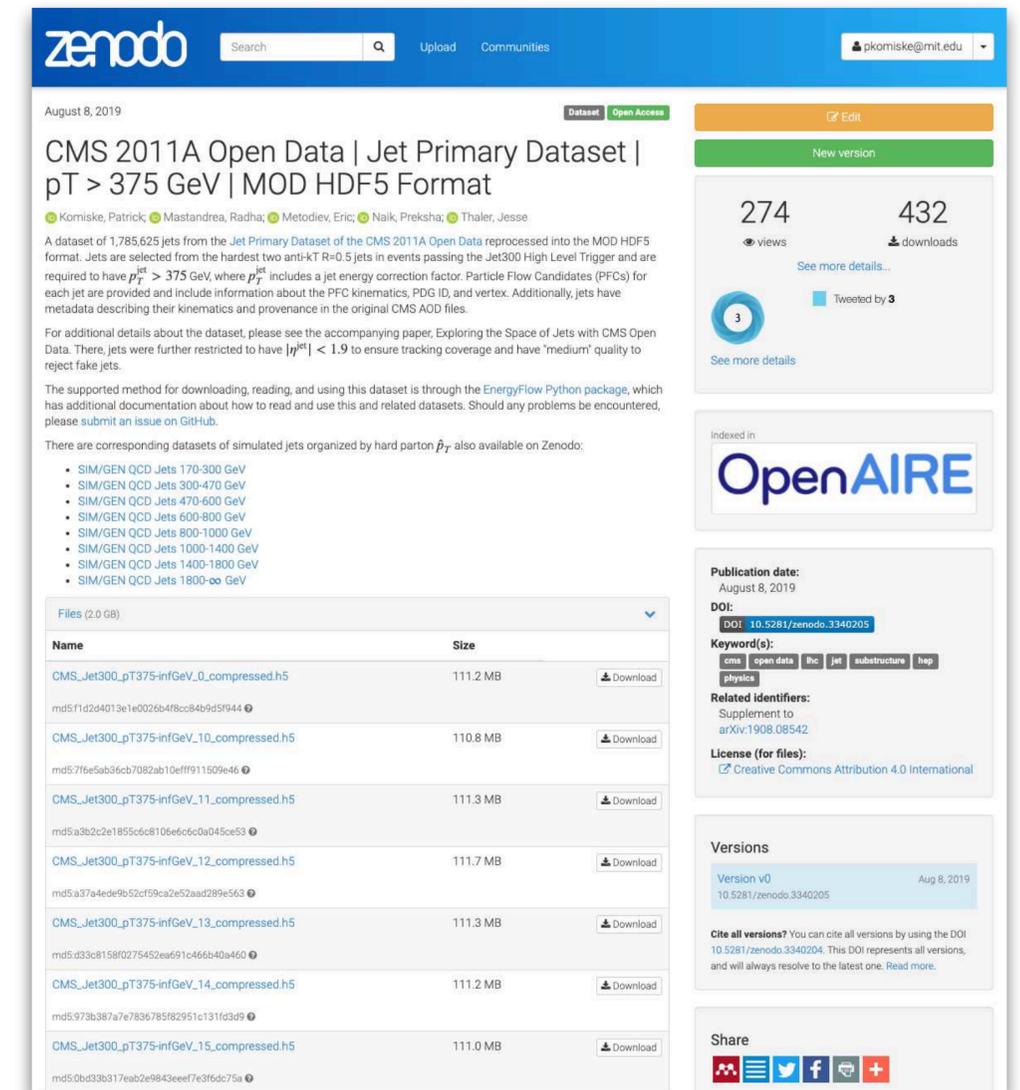
Interfaces with **CMS 2011A Jet Primary Dataset** hosted on **Zenodo**



The screenshot shows the EnergyFlow website. On the left is a dark sidebar with a search bar and a navigation menu including Home, Getting Started, Installation, Demos, Examples, FAQs, Release Notes, News, Documentation, Architectures, Datasets, EMD, Energy Flow Moments, Energy Flow Polynomials, Measures, Multigraph Generation, Observables, and Utils. The main content area has a red header with the EnergyFlow logo and a search bar. Below the header, it says "Welcome to EnergyFlow" and features two images: a colorful particle flow diagram and a scatter plot labeled "MOD" with "EMD: 159.3 GeV". The text below explains that EnergyFlow is a Python package for particle physics tools. It lists several key features: Energy Flow Polynomials (EFPs), Energy Flow Networks (EFNs), Particle Flow Networks (PFNs), Energy Mover's Distance (EMD), and Energy Flow Moments (EFMs). At the bottom, it mentions access to CMS Open Data in MOD HDF5 Format.



The screenshot shows a Jupyter Notebook titled "EMD Demo" in a browser. The notebook content includes: a title "EMD Demo" with a link to the EnergyFlow website; an introductory paragraph about computing EMD values; a section on "Energy Mover's Distance" with a mathematical definition: 
$$\text{EMD}(\mathcal{E}, \mathcal{E}') = \min_{(f_{ij})} \sum_{ij} f_{ij} \frac{\theta_{ij}}{R} + \left| \sum_i E_i - \sum_j E'_j \right|$$
 and constraints  $\sum_j f_{ij} \leq E_i$ ,  $\sum_i f_{ij} \leq E'_j$ ,  $\sum_{ij} f_{ij} = \min(\sum_i E_i, \sum_j E'_j)$ ; an "Imports" section with code for numpy, matplotlib, and energyflow; a "Plot Style" section with code for plt.rcParams; and a "Load EnergyFlow Quark/Gluon Jet Samples" section with code for loading and processing jets.



The screenshot shows the Zenodo dataset page for "CMS 2011A Open Data | Jet Primary Dataset | pT > 375 GeV | MOD HDF5 Format". It includes the dataset title, authors (Komiske, Patrick; Mastandrea, Radha; Metodiev, Eric; Naik, Preksha; Thaler, Jesse), and statistics: 274 views and 432 downloads. It lists 15 files for download, each with its name and size (e.g., CMS\_Jet300\_pt375-infGeV\_0\_compressed.h5, 111.2 MB). The page also shows the publication date (August 8, 2019), DOI (10.5281/zenodo.3340205), and related identifiers.

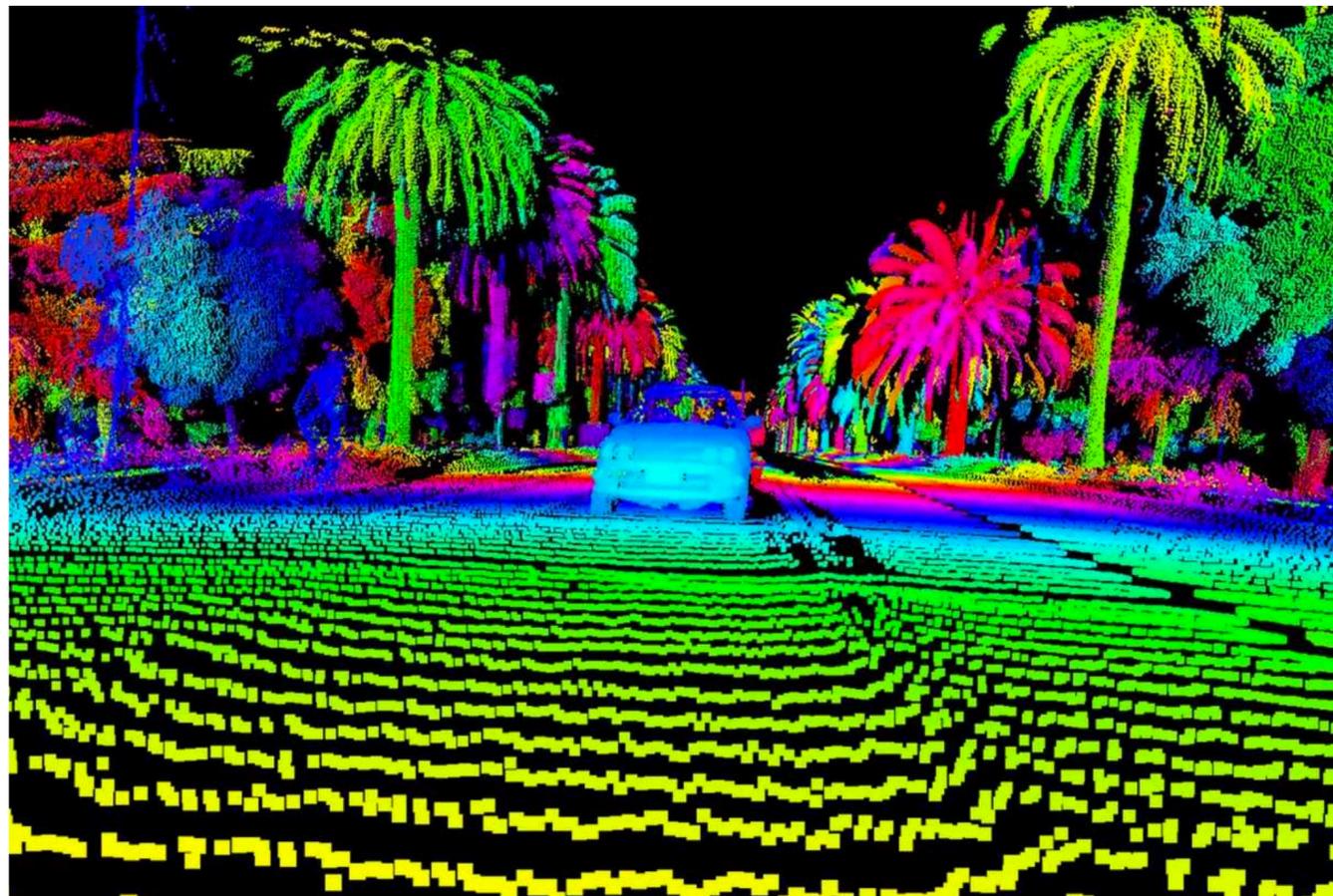
# Additional Slides

# Neural Network Architectures for Particle Physics

*Maximally appropriate ML architectures respect symmetries of the underlying data*

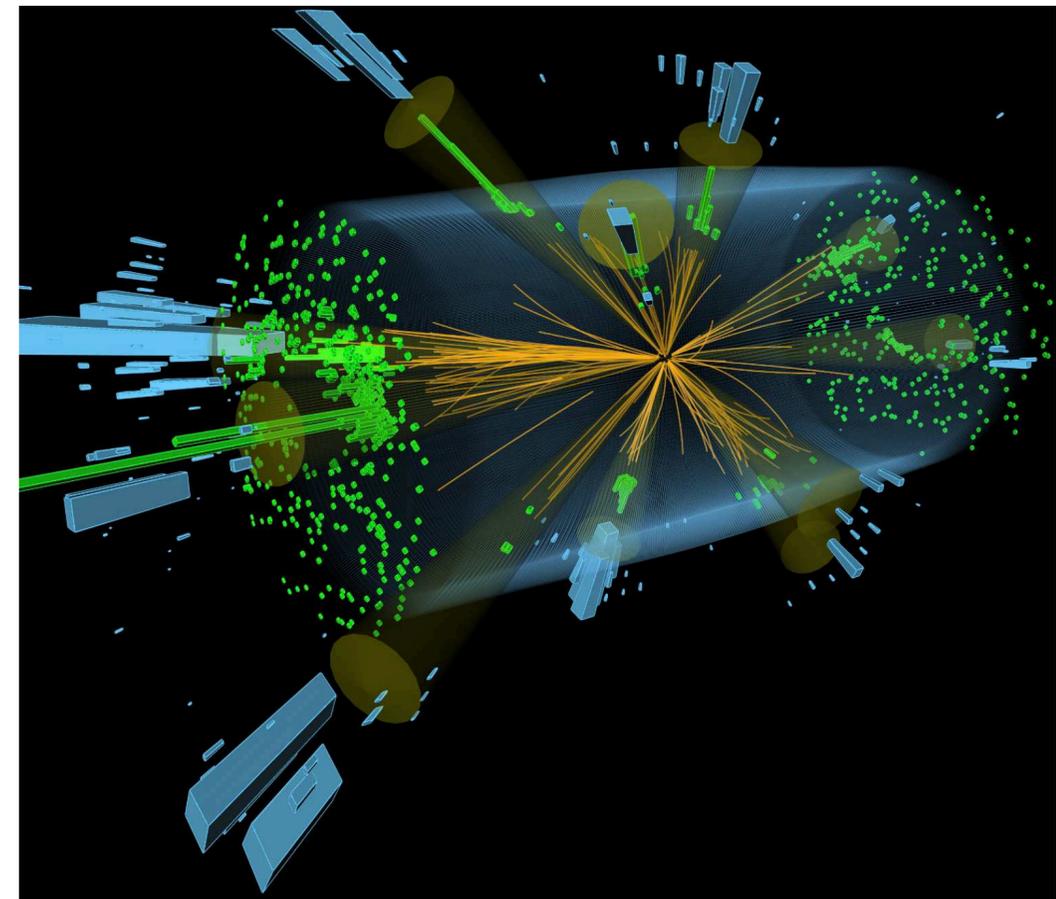
Particle physics events are naturally point clouds (alternatively, images e.g. calorimeters)

Point cloud: "A set of data points in space" –Wikipedia



LIDAR data from self-driving car sensor

An **unordered**, **variable length** collection of particles

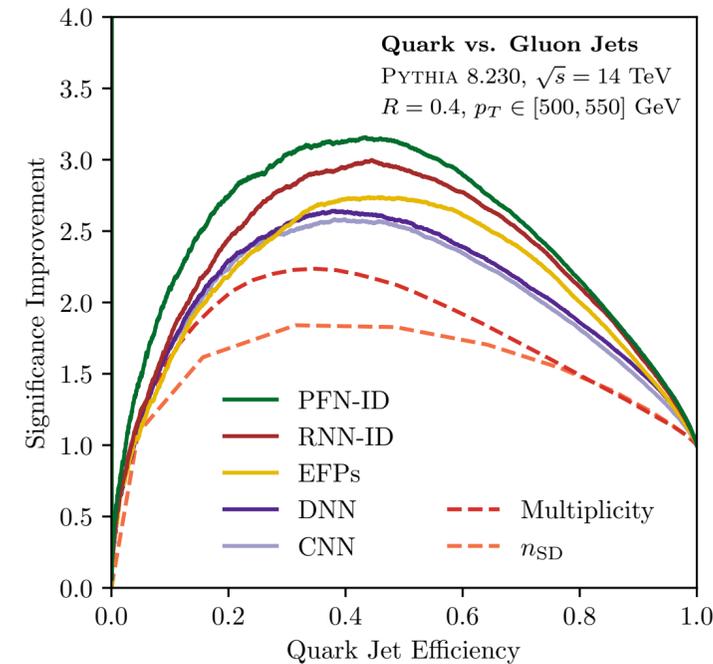
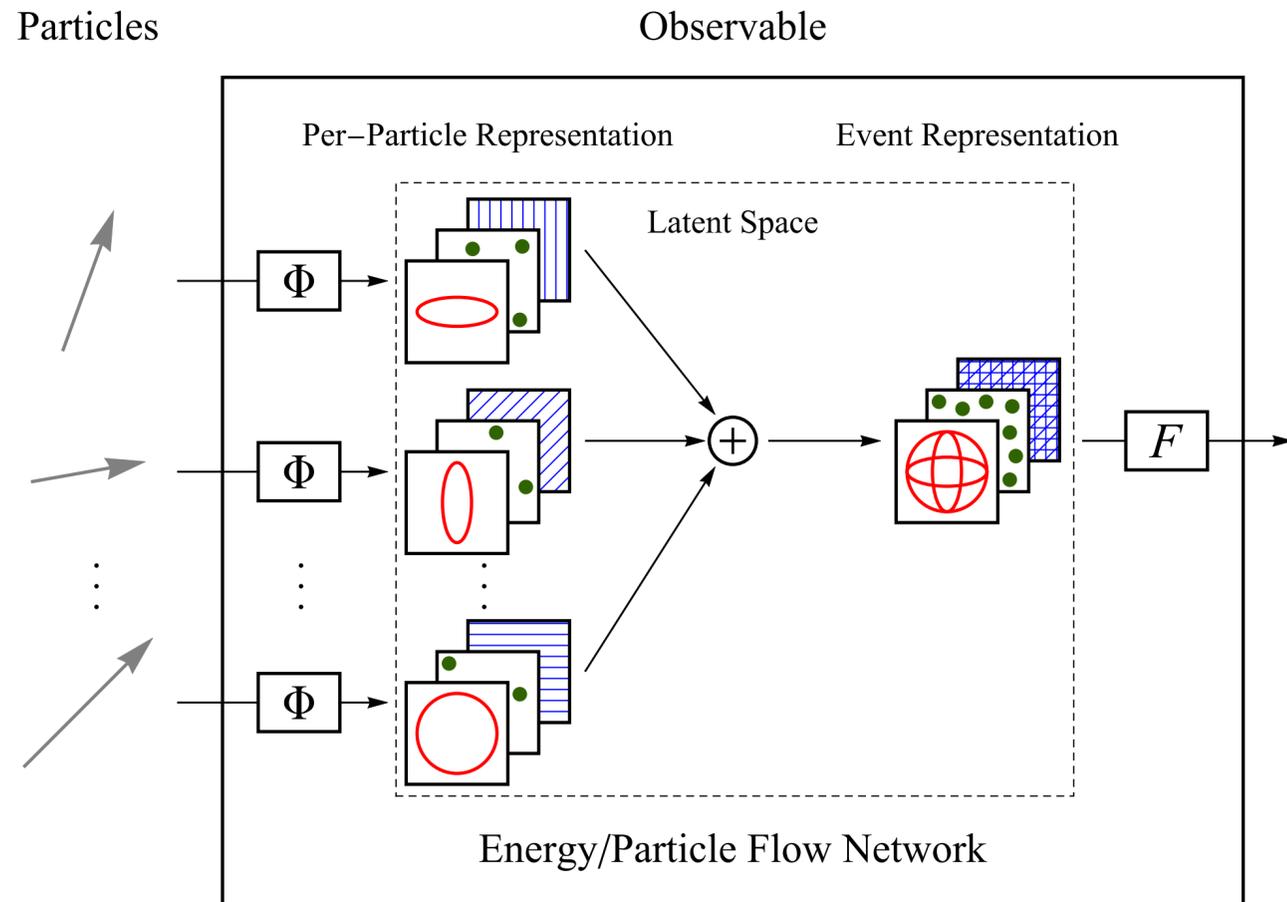


Multi-jet event at CMS

Due to quantum-mechanical indistinguishability  
Due to probabilistic nature of event formation

# Energy/Particle Flow Networks (EFNs/PFNs)

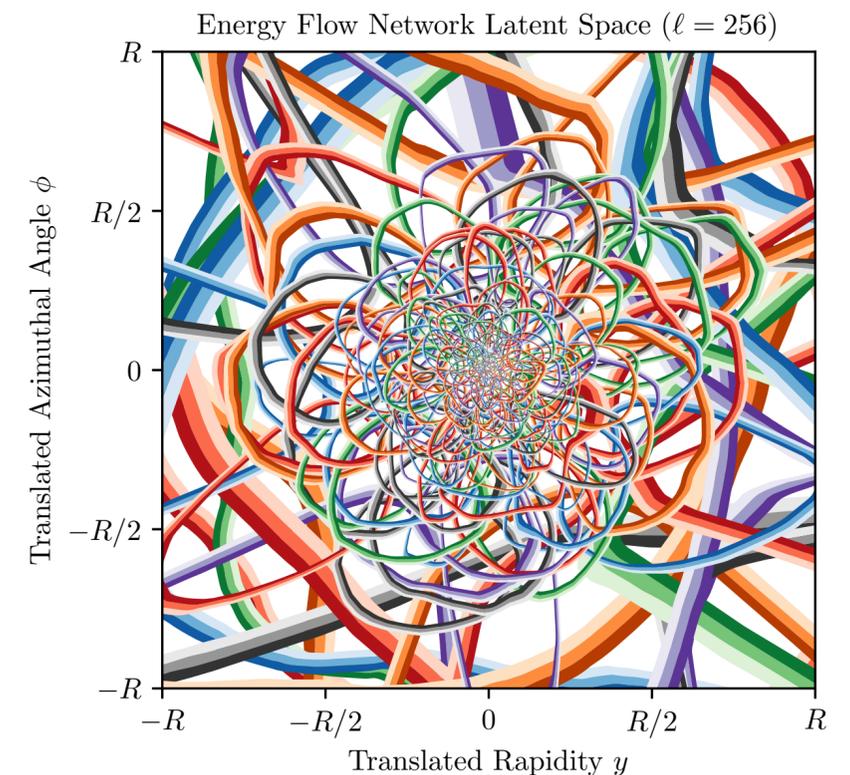
[Zaheer, Kottur, Ravanbakhsh, Póczos, Salakhutdinov, Smola, [1703.06114](#);  
PTK, Metodiev, Thaler, [1810.05165](#);  
[EnergyFlow Python Package](#)]



*Improved performance (and training) compared to RNN and CNN*

*Latent space visualization reveals what the network has learned*

Dynamic pixel sizing related to collinear singularity of QCD!



**Particle Flow Network (PFN)**

$$\text{PFN}(\{p_1^\mu, \dots, p_M^\mu\}) = F \left( \sum_{i=1}^M \Phi(p_i^\mu) \right)$$

Fully general latent space

**Energy Flow Network (EFN)**

$$\text{EFN}(\{p_1^\mu, \dots, p_M^\mu\}) = F \left( \sum_{i=1}^M z_i \Phi(\hat{p}_i) \right)$$

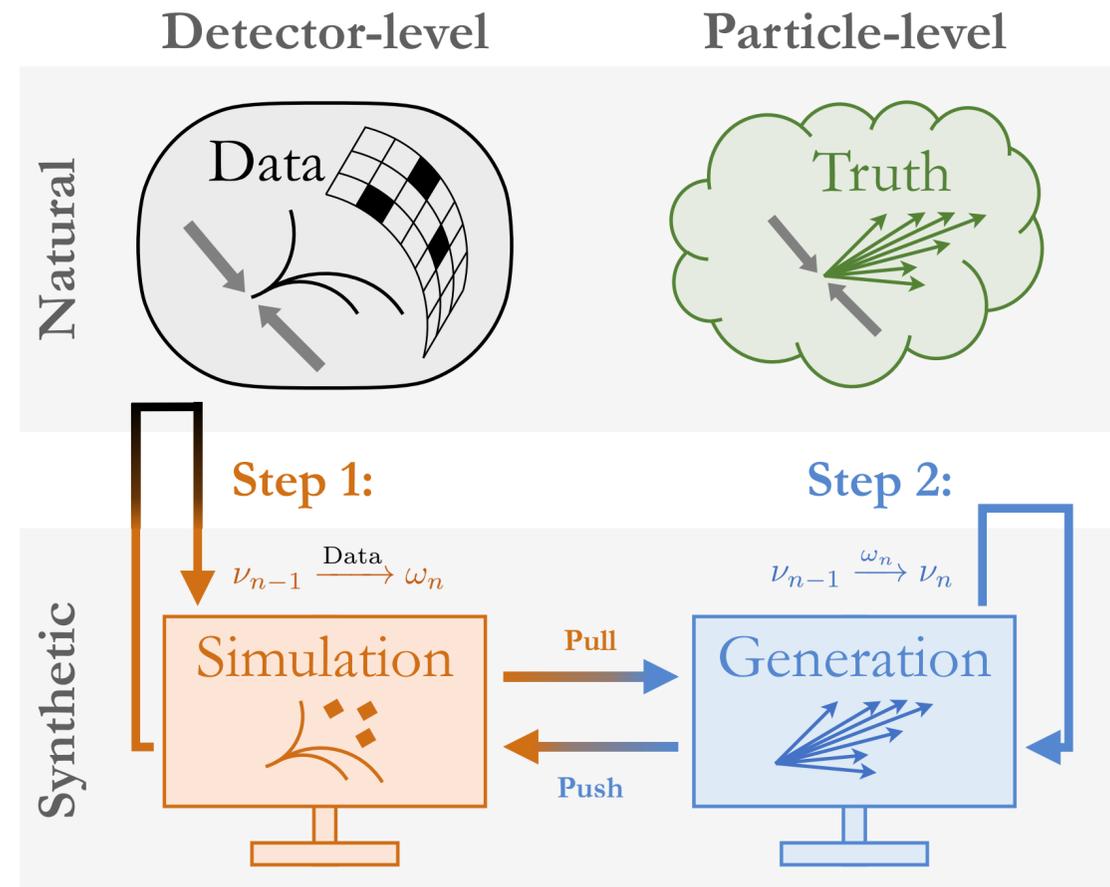
IRC-safe latent space

# OmniFold – Unbinned, Full Phase-Space Unfolding



*OmniFold* weights particle-level **Gen** to be consistent with **Data** once passed through the detector

[Andreassen, PTK, Metodiev, Nachman, Thaler, [1911.09107](#); PTK talk at ML4Jets 2020]



Step 1 – Reweights **Sim<sub>n-1</sub>** to data, pulls weights back to particle-level **Gen<sub>n-1</sub>**

Step 2 – Reweights **Gen<sub>n-1</sub>** to (step 1)-weighted **gen<sub>n-1</sub>**, pushes weights to detector-level **Sim<sub>n</sub>**

## OmniFold – i.e. continuous IBU

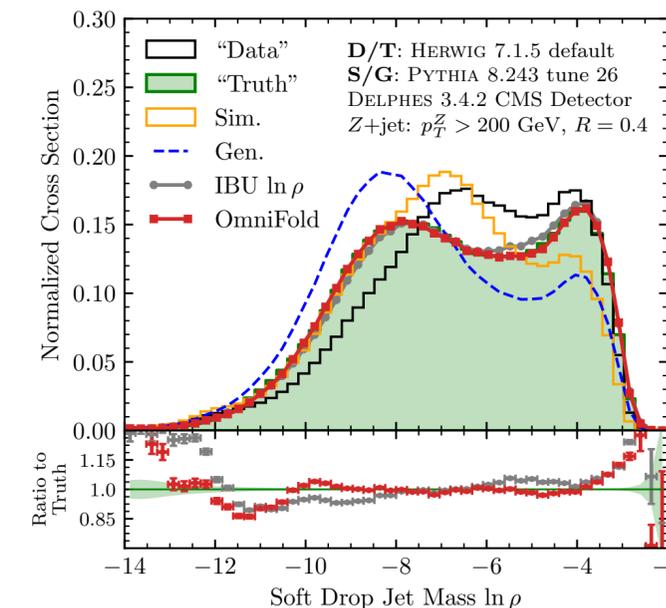
$$\text{Step 1} - \omega_n(m) = \nu_{n-1}^{\text{push}} \times L[(1, \text{Data}), (\nu_{n-1}^{\text{push}}, \text{Sim})](m)$$

$$\text{Step 2} - \nu_n(t) = \nu_{n-1}(t) \times L[(\omega_n^{\text{pull}}, \text{Gen}), (\nu_{n-1}, \text{Gen})](t)$$

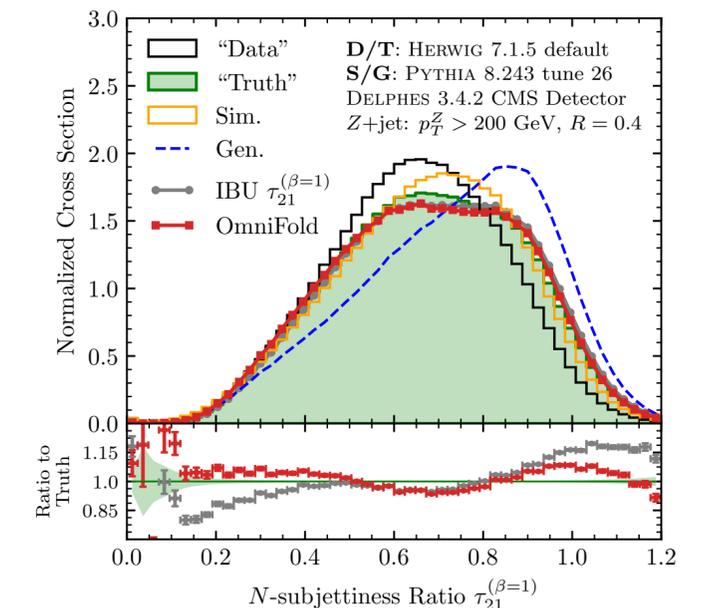
Unfold any\* observable  $p_{\text{Gen}}(t)$  using universal weights  $\nu_n(t)$

$$p_{\text{unfolded}}^{(n)}(t) = \nu_n(t) \times p_{\text{Gen}}(t)$$

\*Observables should be chosen responsibly



IRC safe



Sudakov safe

# Explicit Geometry – Individual Events in Theory

## Hard collision

Good understanding via perturbation theory

## Fragmentation

Semi-classical parton shower, effective field theory

## Hadronization

Poorly understood (non-perturbative), modeled empirically

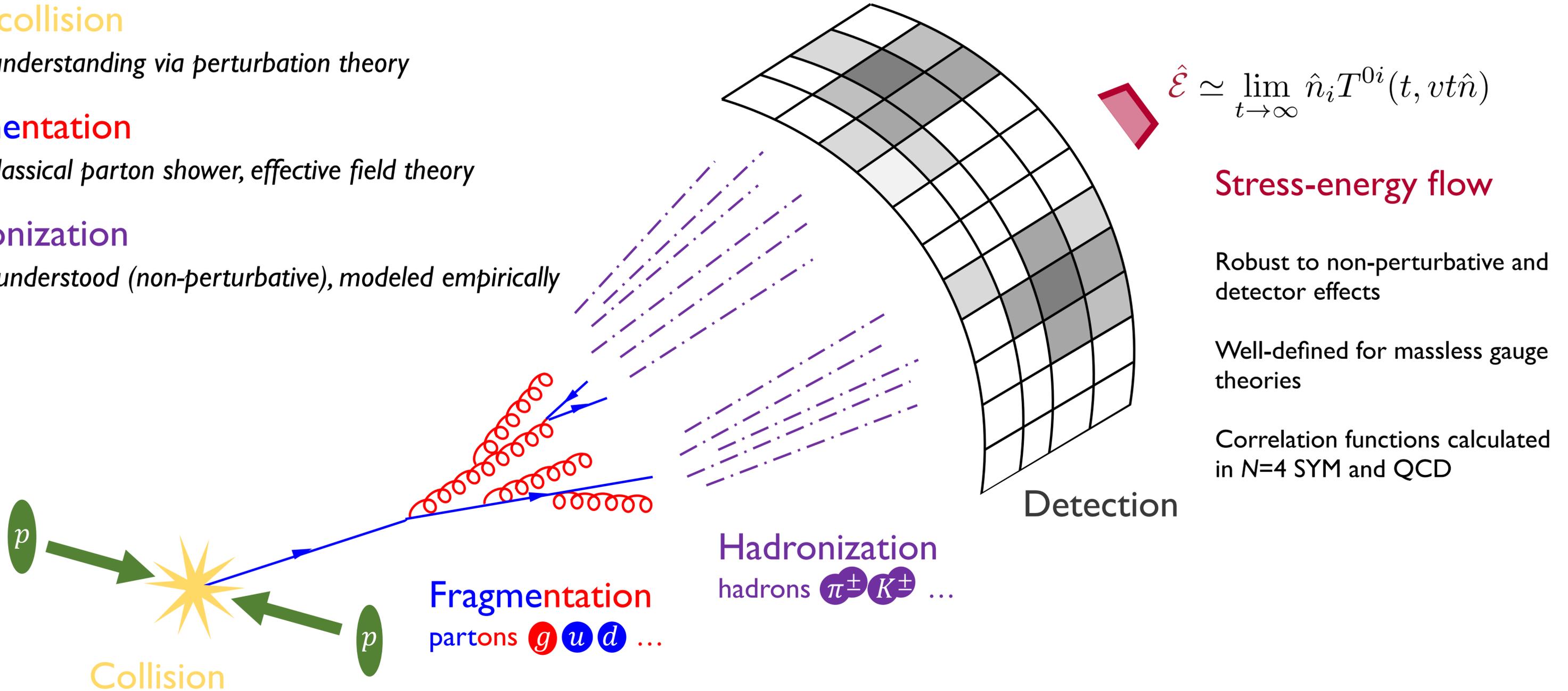


Diagram by Eric Metodiev

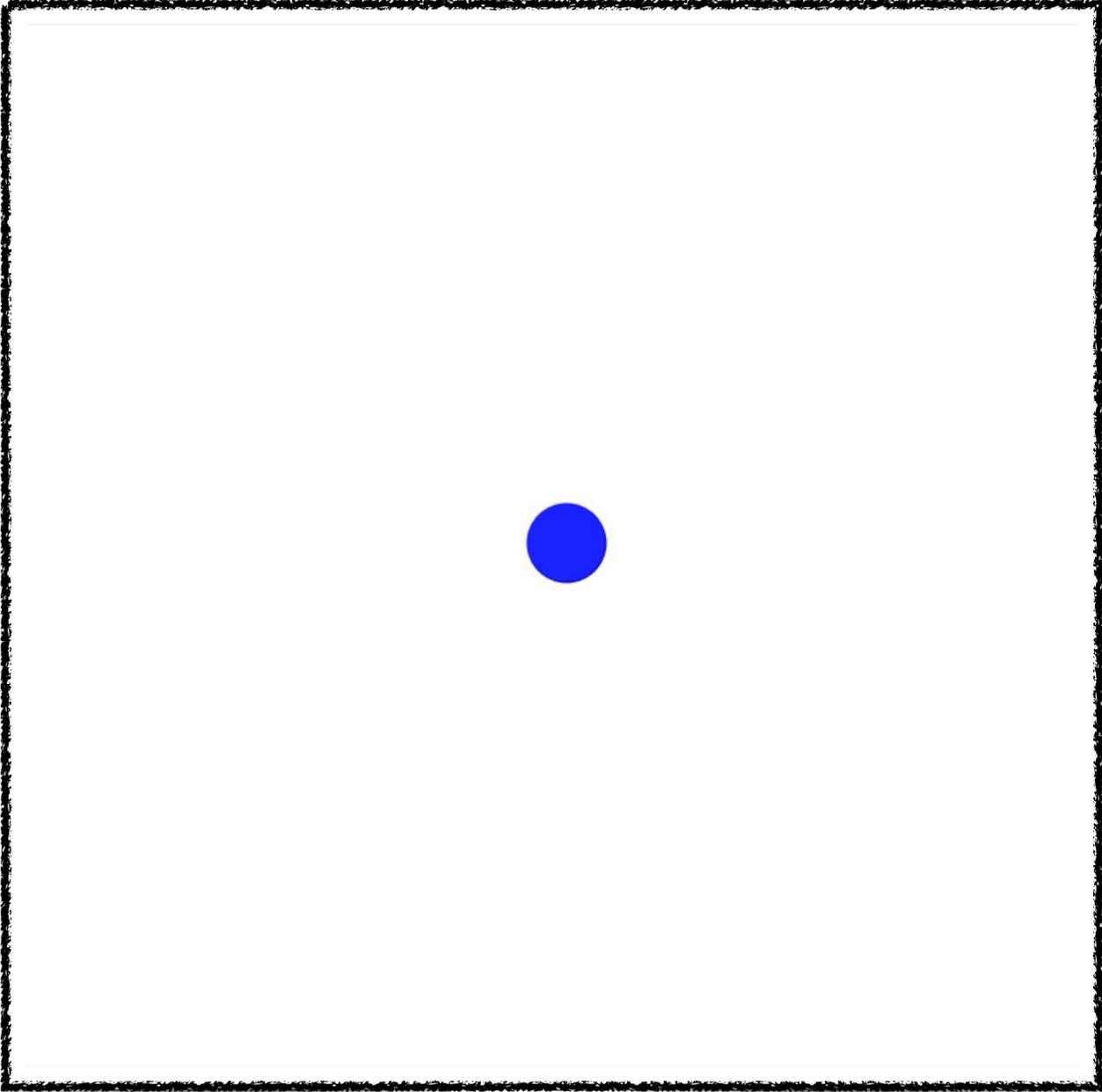
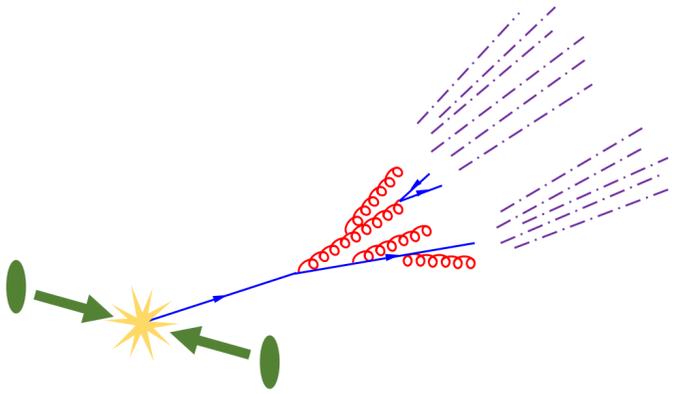
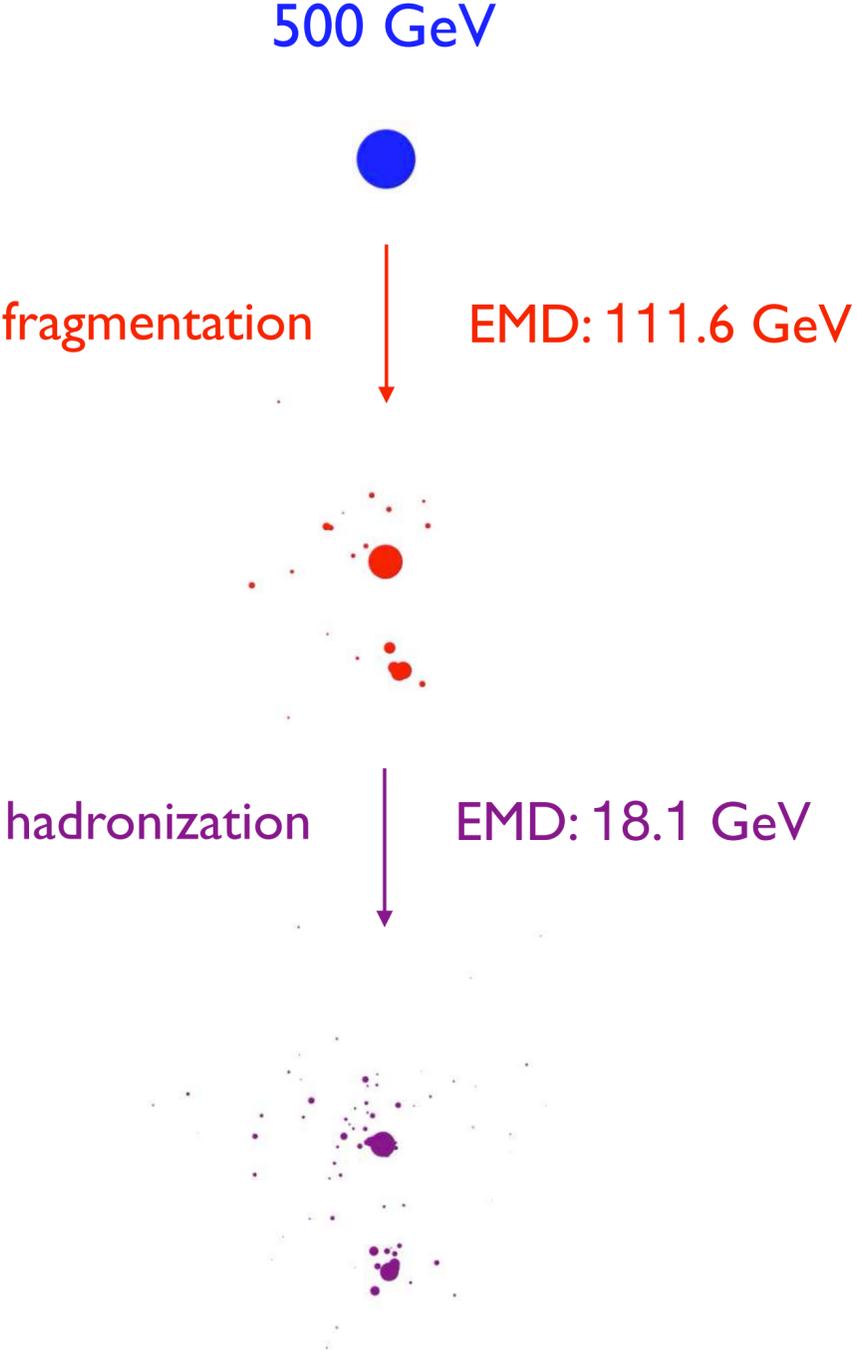
[Sveshnikov, Tkachov, [PLB 1996](#); Hofman, Maldacena, [JHEP 2008](#); Mateu, Stewart, Thaler, [PRD 2013](#); Belitsky, Hohenegger, Korchemsky, Sokatchev, Zhiboedov, [PRL 2014](#); Chen, Moulton, Zhang, Zhu, [2004.11381](#); Dixon, PTK, Moulton, Thaler, Zhu, [to appear soon, see more here](#)]

# Table of Observables Defined via Event Space Geometry

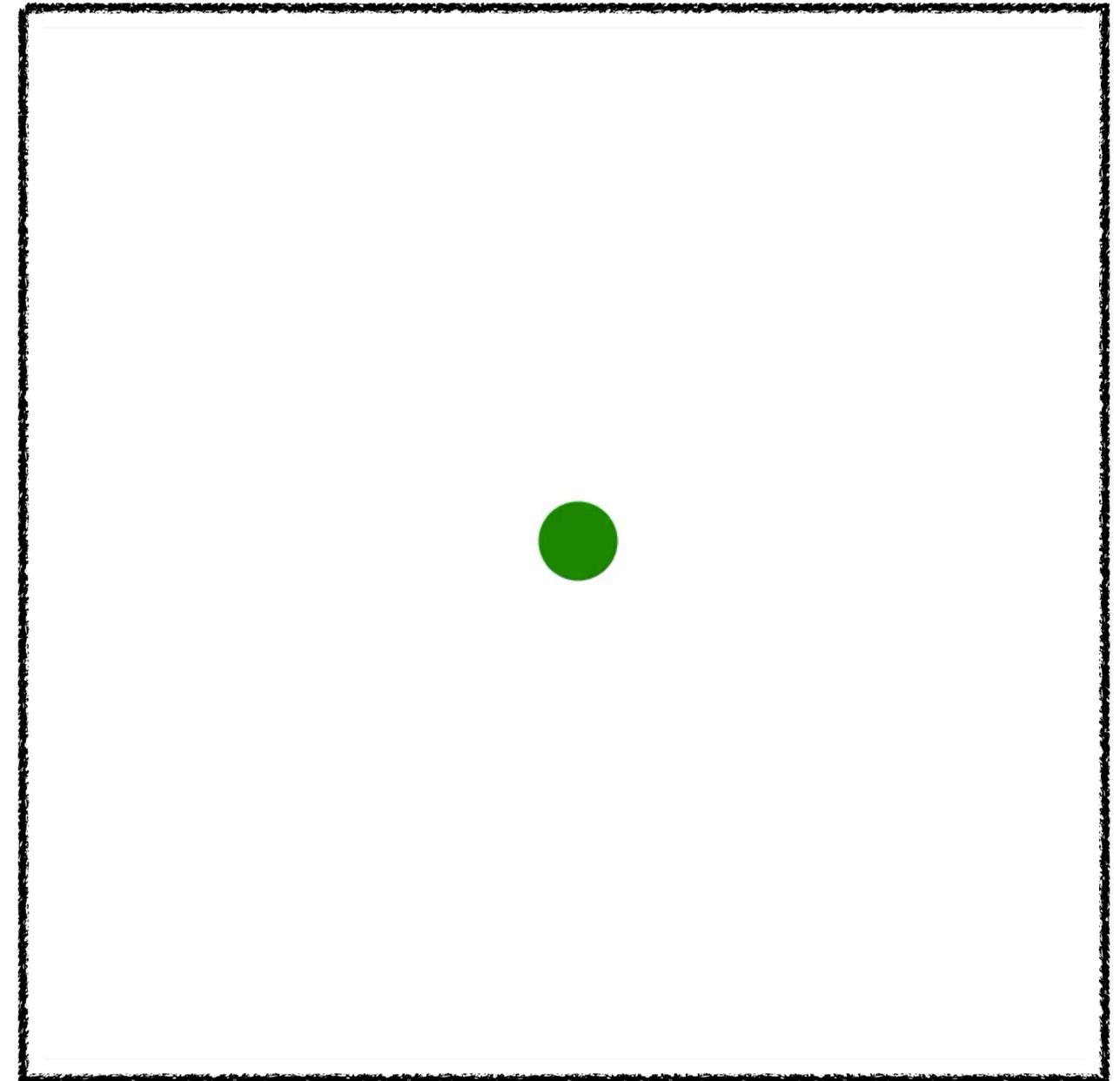
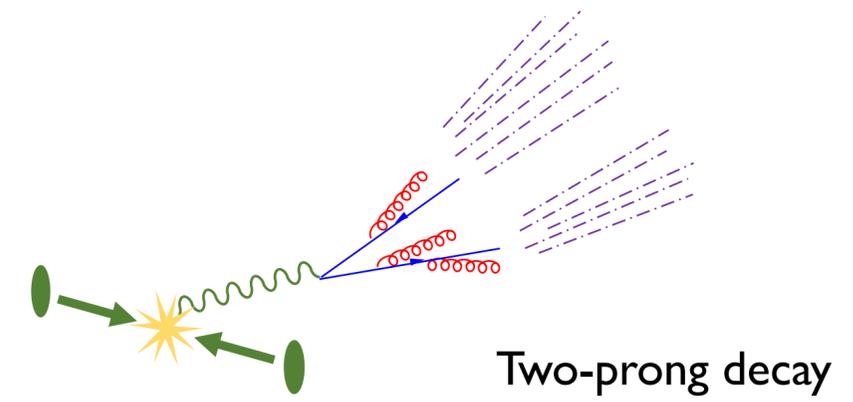
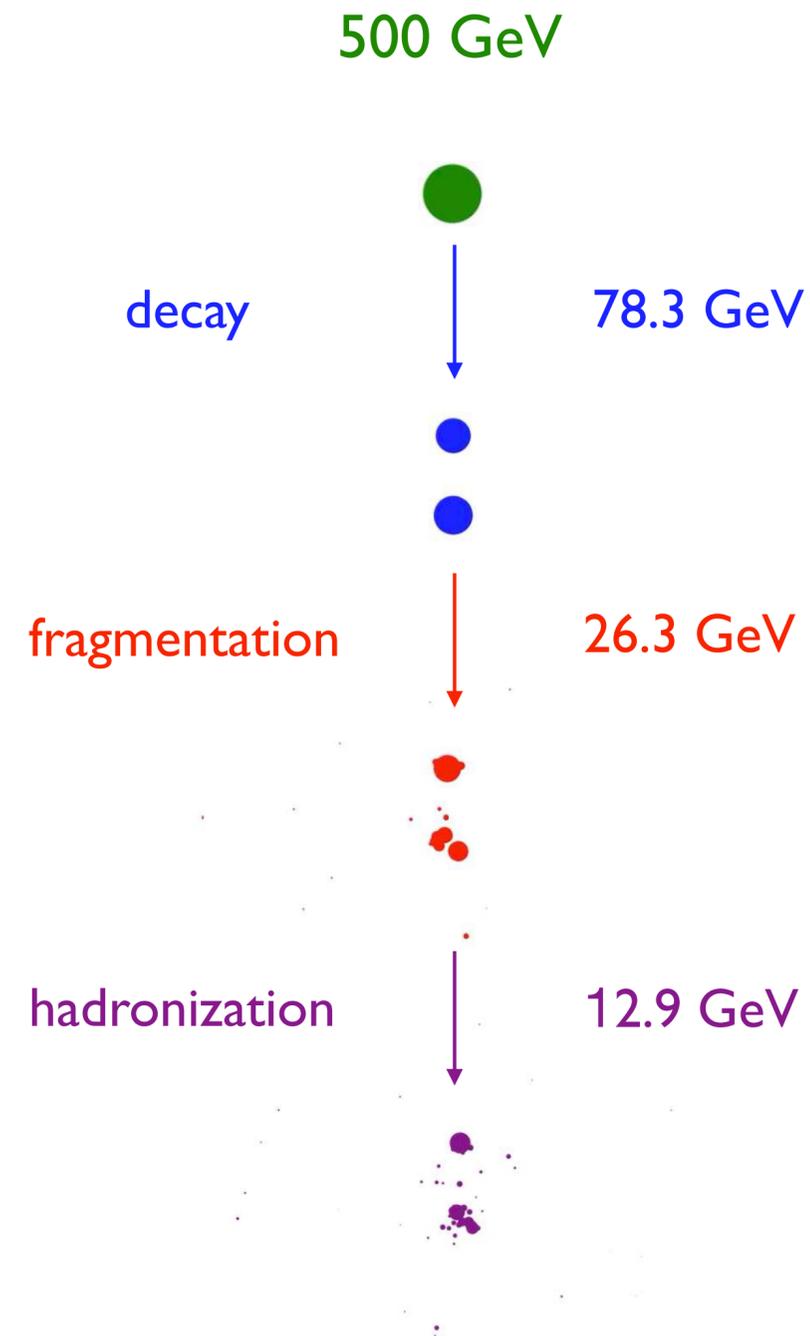
$$\mathcal{O}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{M}} \text{EMD}_{\beta, R}(\mathcal{E}, \mathcal{E}')$$

$\mathcal{O}(\mathcal{E}) = \min_{\mathcal{E}' \in \mathcal{M}} \text{EMD}_{\beta}(\mathcal{E}, \mathcal{E}')$			
Name		$\beta$	Manifold $\mathcal{M}$
Thrust	$t(\mathcal{E})$	2	$\mathcal{P}_2^{\text{BB}}$ : 2-particle events, back to back
Spherocity	$\sqrt{s(\mathcal{E})}$	1	$\mathcal{P}_2^{\text{BB}}$ : 2-particle events, back to back
Broadening	$b(\mathcal{E})$	1	$\mathcal{P}_2$ : 2-particle events
$N$ -jettiness	$\mathcal{T}_N^{(\beta)}(\mathcal{E})$	$\beta$	$\mathcal{P}_N$ : $N$ -particle events
Isotropy	$\mathcal{I}^{(\beta)}(\mathcal{E})$	$\beta$	$\mathcal{M}_{\mathcal{U}}$ : Uniform events
Jet Angularities	$\lambda_{\beta}(\mathcal{J})$	$\beta$	$\mathcal{P}_1$ : 1-particle jets
$N$ -subjettiness	$\tau_N^{(\beta)}(\mathcal{J})$	$\beta$	$\mathcal{P}_N$ : $N$ -particle jets

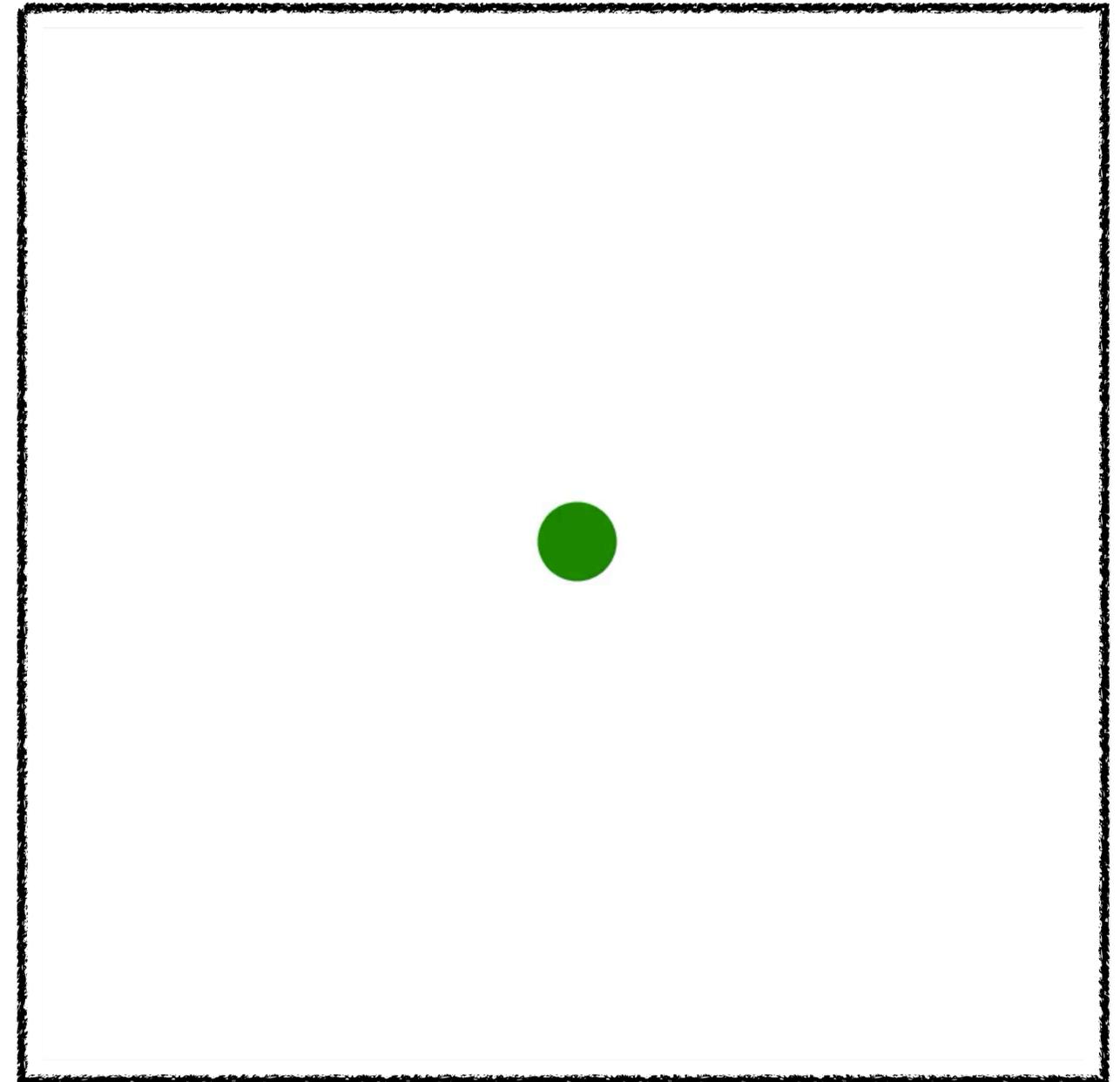
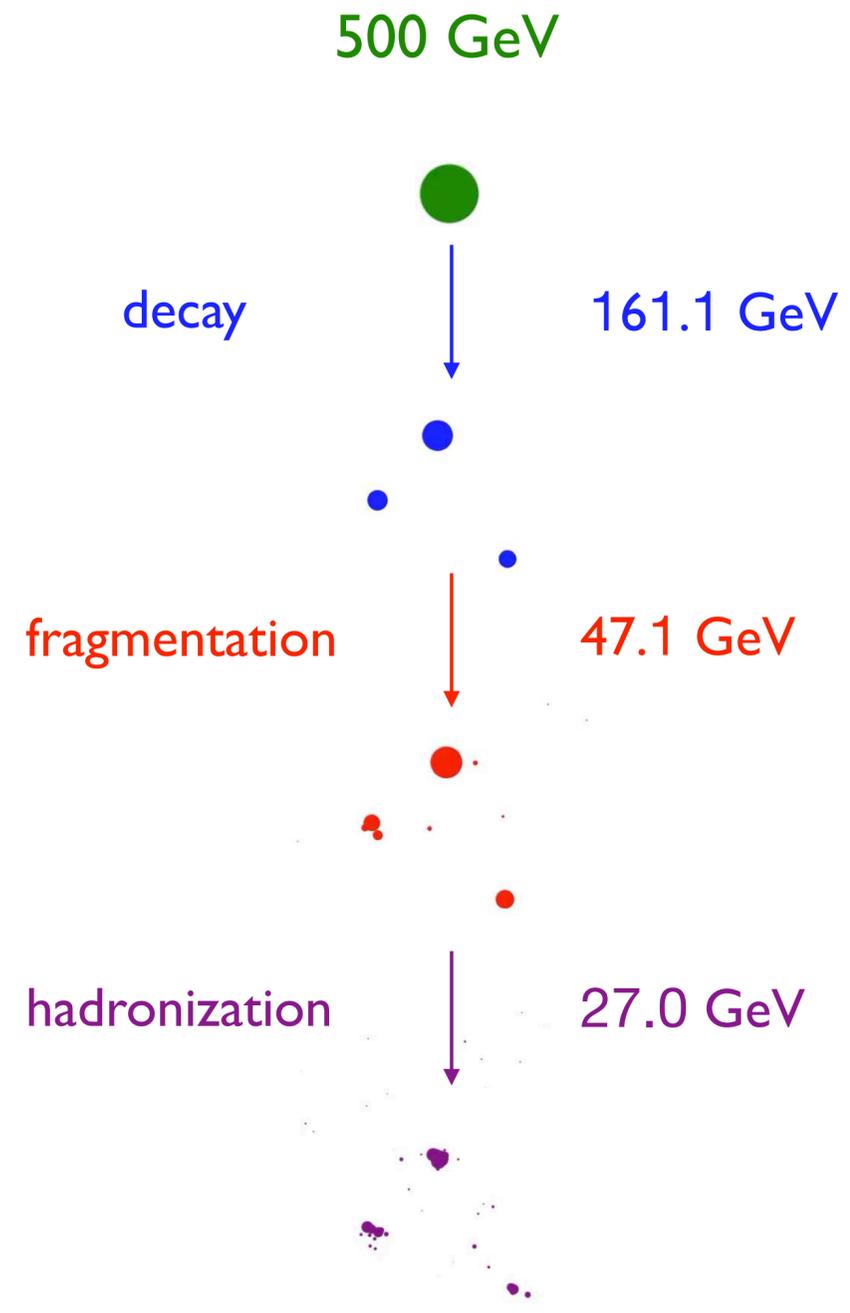
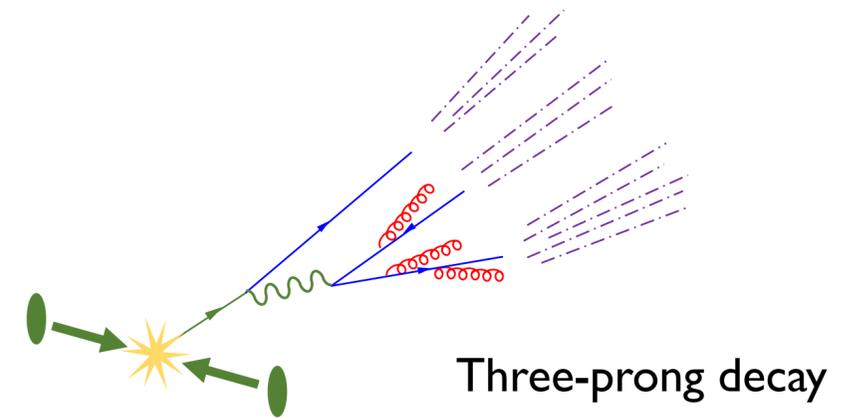
# Visualizing Jet Formation – QCD Jets



# Visualizing Jet Formation – W Jets



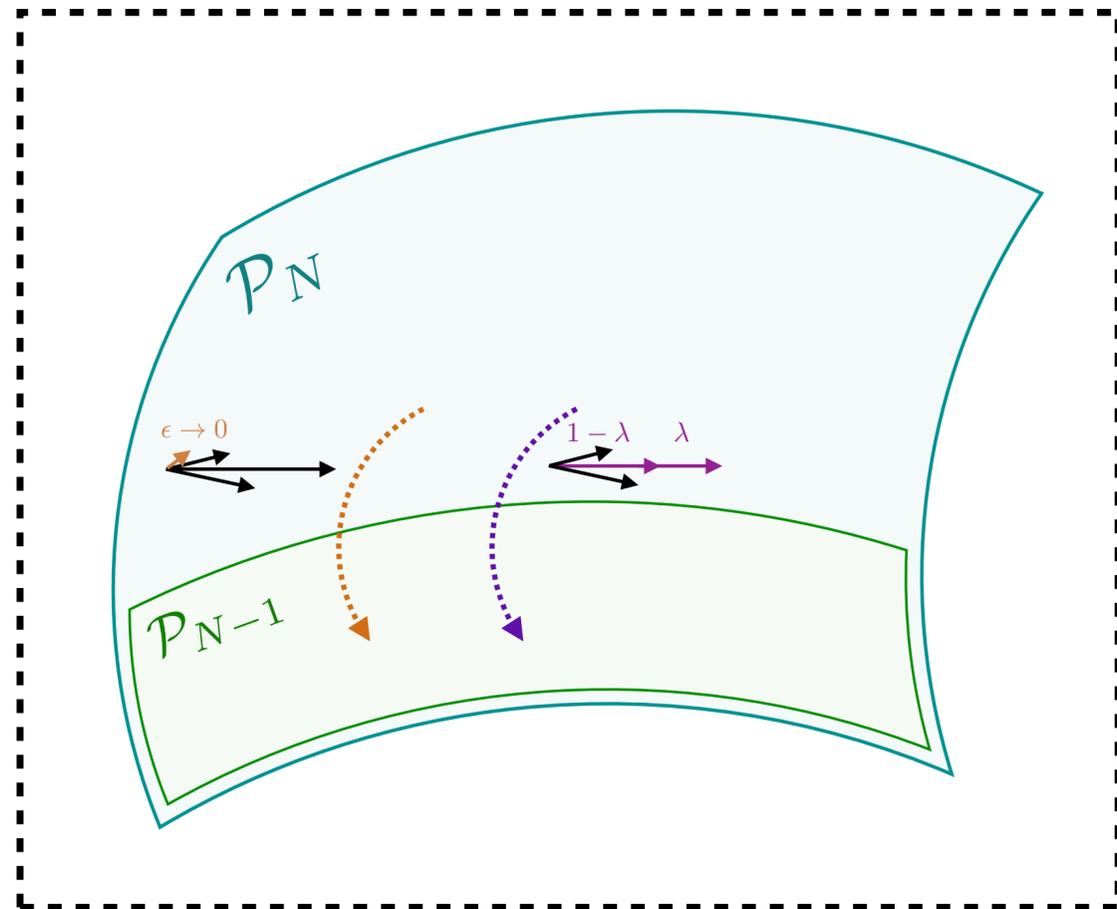
# Visualizing Jet Formation – Top Jets



# N-particle Manifolds in the Space of Events – Infrared Divergences

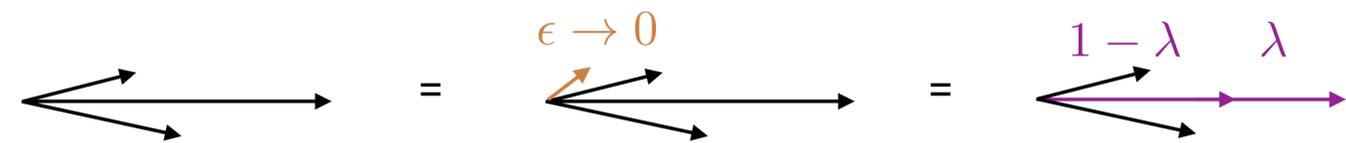
[PTK, Metodiev, Thaler, JHEP 2020]

$$\mathcal{P}_N = \text{set of all } N\text{-particle configurations} = \left\{ \sum_{i=1}^N E_i \delta(\hat{n} - \hat{n}_i) \mid E_i \geq 0 \right\}$$



$$dP_{i \rightarrow ig} \simeq \frac{2\alpha_s}{\pi} C_a \frac{d\theta}{\theta} \frac{dz}{z}$$

Energy flow is unchanged by exact soft/collinear emissions



Functions of energy flow automatically satisfy exact IRC invariance!

Real and virtual divergences appear naturally together

$$\mathcal{P}_N \supset \mathcal{P}_{N-1} \supset \cdots \supset \mathcal{P}_3 \supset \mathcal{P}_2 \supset \mathcal{P}_1$$

by soft and collinear limits

# Defining IRC Safety Precisely

[Sterman, Weinberg, [PRL 1997](#); Sterman, [PRD 1978](#); Banfi, Salam, Zanderighi, [JHEP 2005](#)]

*Infrared and collinear safety is a proxy for perturbative calculability of an observable*

## Exact IRC invariance

$$\mathcal{O}(p_1^\mu, \dots, p_M^\mu) = \mathcal{O}(0p_0^\mu, p_1^\mu, \dots, p_M^\mu)$$

$$\mathcal{O}(p_1^\mu, \dots, p_M^\mu) = \mathcal{O}(\lambda p_1^\mu, (1 - \lambda)p_1^\mu, \dots, p_M^\mu)$$

Guarantees observable is well-defined on **energy** flows

Allows for pathological observables, e.g. pseudo-multiplicity

## Smooth IRC invariance

$$\mathcal{O}(p_1^\mu, \dots, p_M^\mu) = \lim_{\epsilon \rightarrow 0} \mathcal{O}(\epsilon p_0^\mu, p_1^\mu, \dots, p_M^\mu)$$

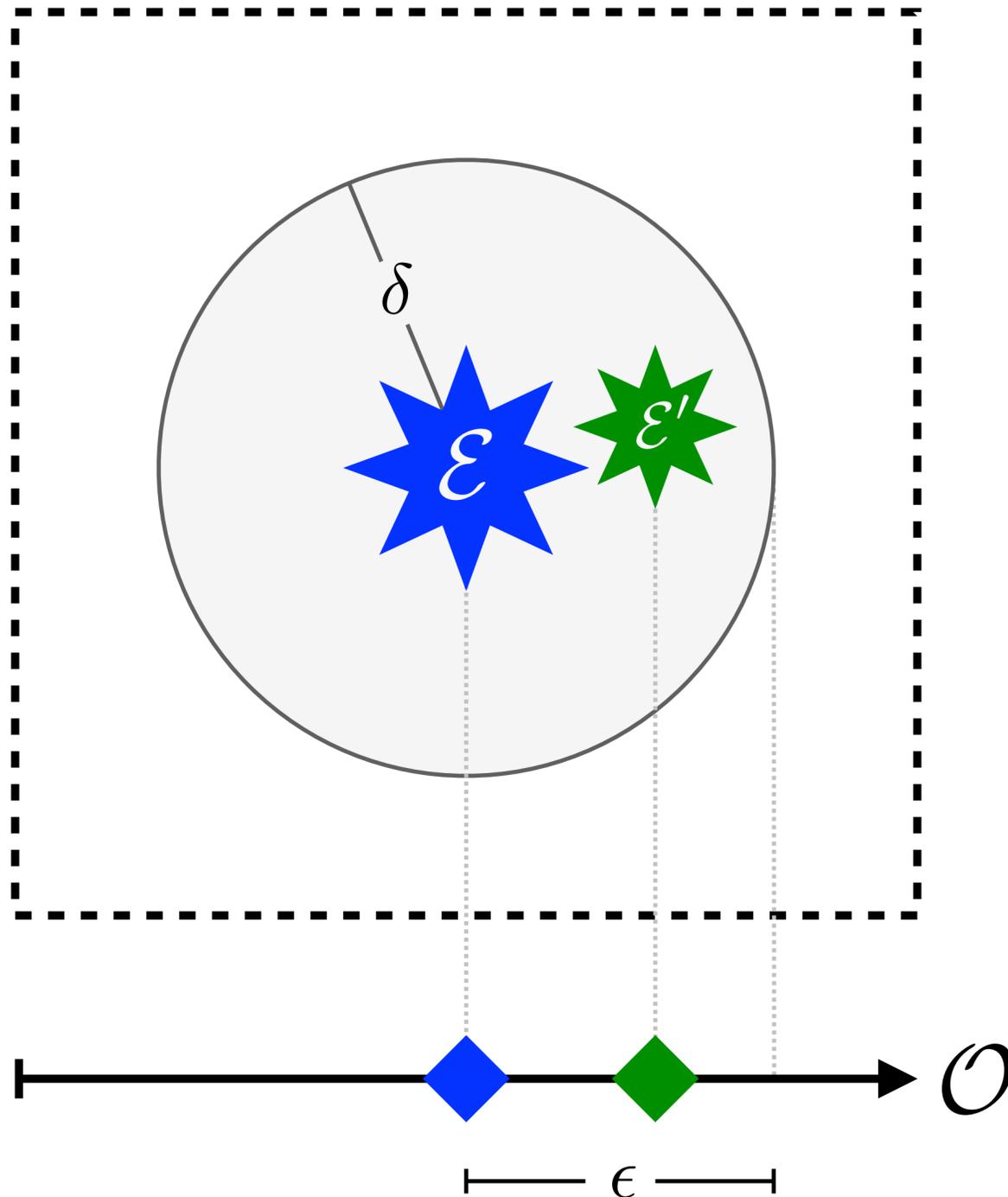
$$\mathcal{O}(p_1^\mu, \dots, p_M^\mu) = \lim_{p_0^\mu \rightarrow p_1^\mu} \mathcal{O}(\lambda p_0^\mu, (1 - \lambda)p_1^\mu, \dots, p_M^\mu)$$

Eliminates common observables with hard boundaries

All Observables	Comments
Multiplicity ( $\sum_i 1$ )	IR unsafe and C unsafe
Momentum Dispersion [65] ( $\sum_i E_i^2$ )	IR safe but C unsafe
Sphericity Tensor [66] ( $\sum_i p_i^\mu p_i^\nu$ )	IR safe but C unsafe
Number of Non-Zero Calorimeter Deposits	C safe but IR unsafe
Defined on Energy Flows	
Pseudo-Multiplicity ( $\min\{N \mid \mathcal{T}_N = 0\}$ )	Robust to exact IR or C emissions
Infrared & Collinear Safe	
Jet Energy ( $\sum_i E_i$ )	Disc. at jet boundary
Heavy Jet Mass [67]	Disc. at hemisphere boundary
Soft-Dropped Jet Mass [38, 68]	Disc. at grooming threshold
Calorimeter Activity [69] ( $N_{95}$ )	Disc. at cell boundary

# More EMD Geometry – Continuity in the Space of Events

[PTK, Metodiev, Thaler, 2004.04159]



## Classic $\epsilon - \delta$ definition of continuity in a metric space

An observable  $\mathcal{O}$  is **EMD continuous** at an event  $\mathcal{E}$  if, for any  $\epsilon > 0$ , there exists a  $\delta > 0$  such that for all events  $\mathcal{E}'$ :

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

## Towards a geometric definition of IRC Safety

**IRC Safety = EMD Continuity\***

\*on all but a negligible set<sup>‡</sup> of events

<sup>‡</sup>a negligible set is one that contains no positive-radius EMD-ball

⋮

# Perturbation Theory in the Space of Events

[PTK, Metodiev, Thaler, 2004.04159]

## Sudakov safety

[Larkoski, Thaler, JHEP 2014; Larkoski, Marzani, Thaler, PRD 2015]

Some observables have discontinuities on  $P_N$  for some  $N$

A resummed IRC-safe companion can mitigate the divergences

$$p(\mathcal{O}_{\text{Sudakov}}) = \int d\mathcal{O}_{\text{Comp.}} p(\mathcal{O}_{\text{Sudakov}} | \mathcal{O}_{\text{Comp.}}) p(\mathcal{O}_{\text{Comp.}})$$

Event geometry suggests  $N$ -(sub)jettiness as universal companion

## Fixed-order calculability

[Sterman, PRD 1979; Banfi, Salam, Zanderighi, JHEP 2005]

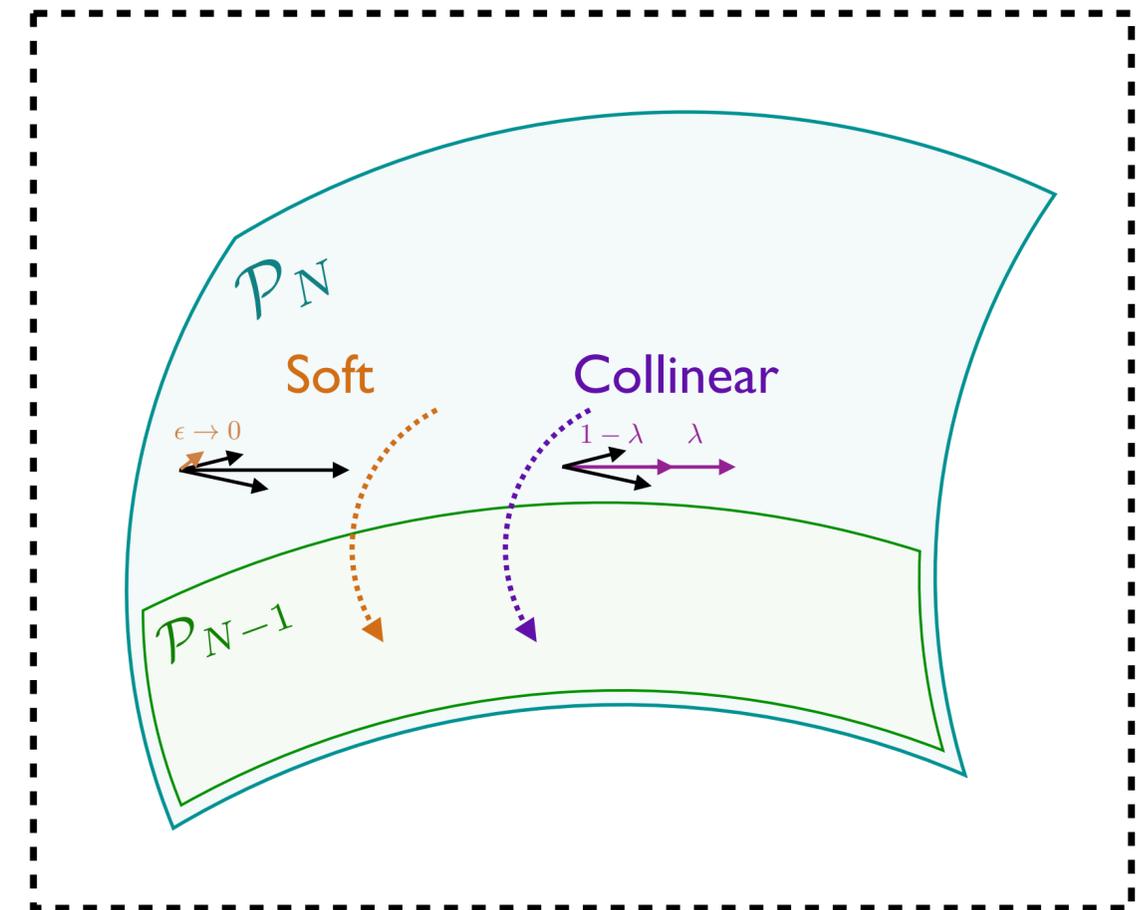
Is a statement of integrability on each  $P_N$

EMD continuity must be upgraded to EMD-Hölder continuity on each  $P_N$

$$\lim_{\mathcal{E} \rightarrow \mathcal{E}'} \frac{\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')}{\text{EMD}(\mathcal{E}, \mathcal{E}')^c} = 0, \quad c > 0$$

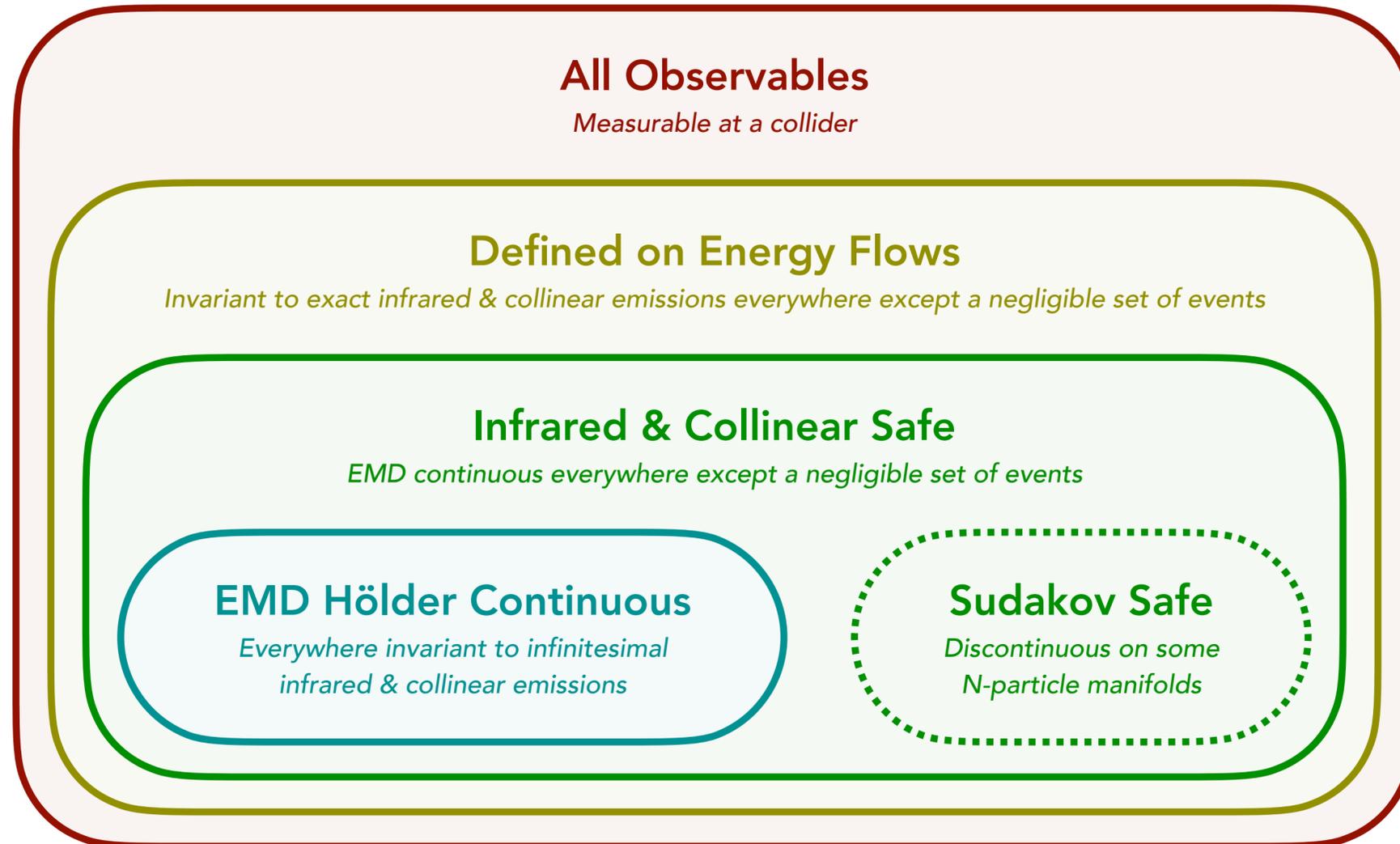
Example:  $V(\mathcal{E}) = \mathcal{T}_2(\mathcal{E}) \left( 1 + \frac{1}{\ln E(\mathcal{E})/\mathcal{T}_3(\mathcal{E})} \right)$  is EMD continuous but not EMD Hölder continuous (it is Sudakov safe)

Infrared singularities of massless gauge theories appear on each  $P_N$



# Hierarchy of IRC Safety Definitions

[PTK, Metodiev, Thaler, 2004.04159]



All Observables	Comments
Multiplicity ( $\sum_i 1$ )	IR unsafe and C unsafe
Momentum Dispersion [65] ( $\sum_i E_i^2$ )	IR safe but C unsafe
Sphericity Tensor [66] ( $\sum_i p_i^\mu p_i^\nu$ )	IR safe but C unsafe
Number of Non-Zero Calorimeter Deposits	C safe but IR unsafe

Defined on Energy Flows	
Pseudo-Multiplicity ( $\min\{N \mid \mathcal{T}_N = 0\}$ )	Robust to exact IR or C emissions

Infrared & Collinear Safe	
Jet Energy ( $\sum_i E_i$ )	Disc. at jet boundary
Heavy Jet Mass [67]	Disc. at hemisphere boundary
Soft-Dropped Jet Mass [38, 68]	Disc. at grooming threshold
Calorimeter Activity [69] ( $N_{95}$ )	Disc. at cell boundary

Sudakov Safe	
Groomed Momentum Fraction [39] ( $z_g$ )	Disc. on 1-particle manifold
Jet Angularity Ratios [37]	Disc. on 1-particle manifold
$N$ -subjettiness Ratios [47, 48] ( $\tau_{N+1}/\tau_N$ )	Disc. on $N$ -particle manifold
$V$ parameter [36] (Eq. (2.11))	Hölder disc. on 3-particle manifold

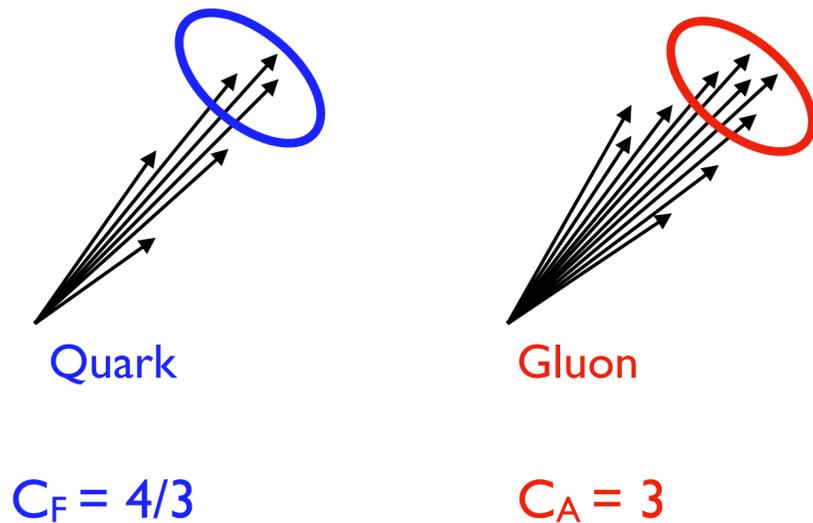
EMD Hölder Continuous Everywhere	
Thrust [40, 41]	
Sphericity [42]	
Angularities [70]	
$N$ -jettiness [44] ( $\mathcal{T}_N$ )	
$C$ parameter [71–74]	Resummation beneficial at $C = \frac{3}{4}$
Linear Sphericity [72] ( $\sum_i E_i n_i^\mu n_i^\nu$ )	
Energy Correlators [36, 75–77]	
Energy Flow Polynomials [15, 17]	

# Quark and Gluon Correlation Dimensions

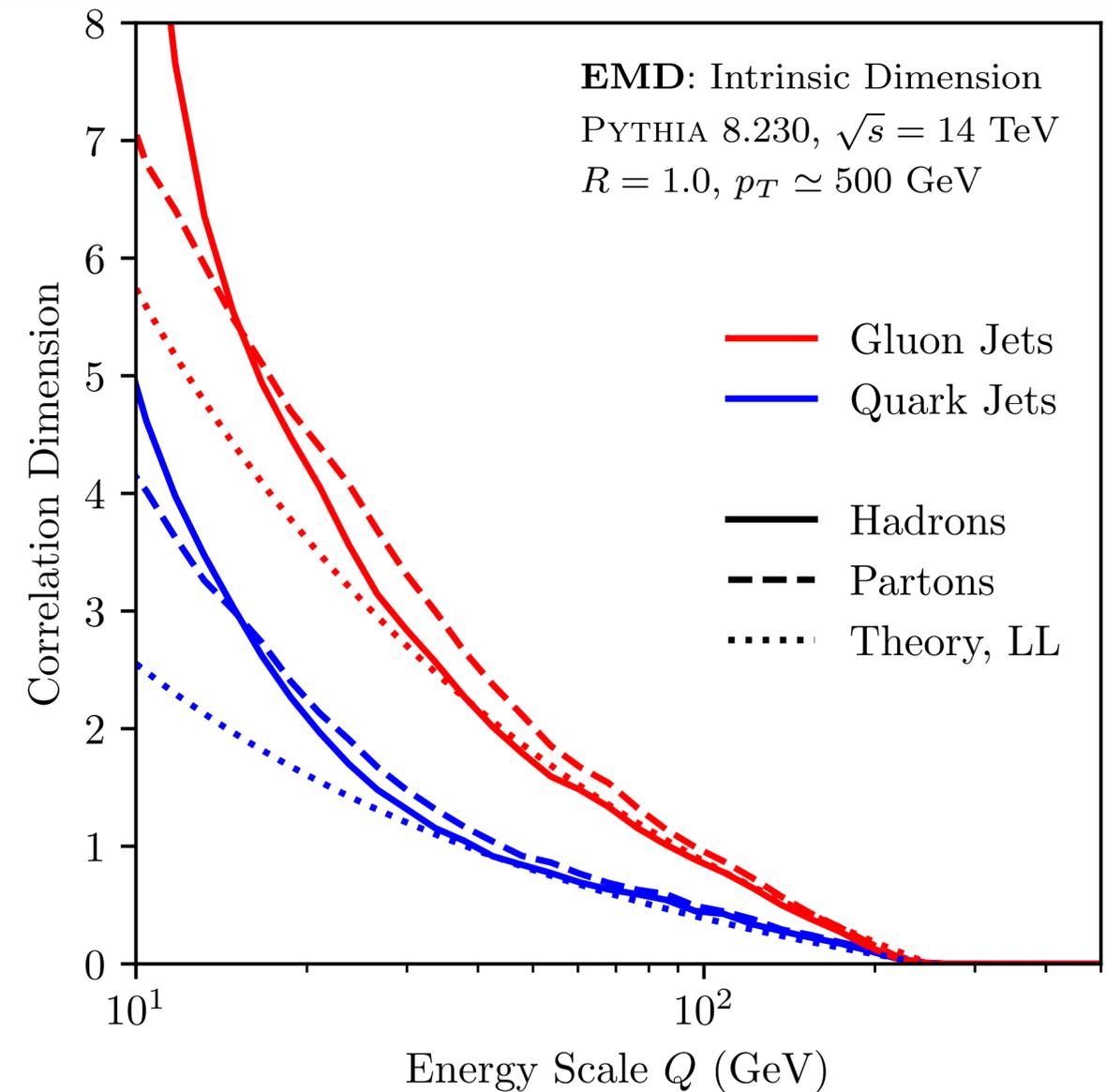
Leading log (single emission) calculation:

$$\text{dim}_i(Q) \simeq -\frac{8\alpha_s}{\pi} C_i \ln \frac{Q}{p_T/2}$$

↑  
color factor



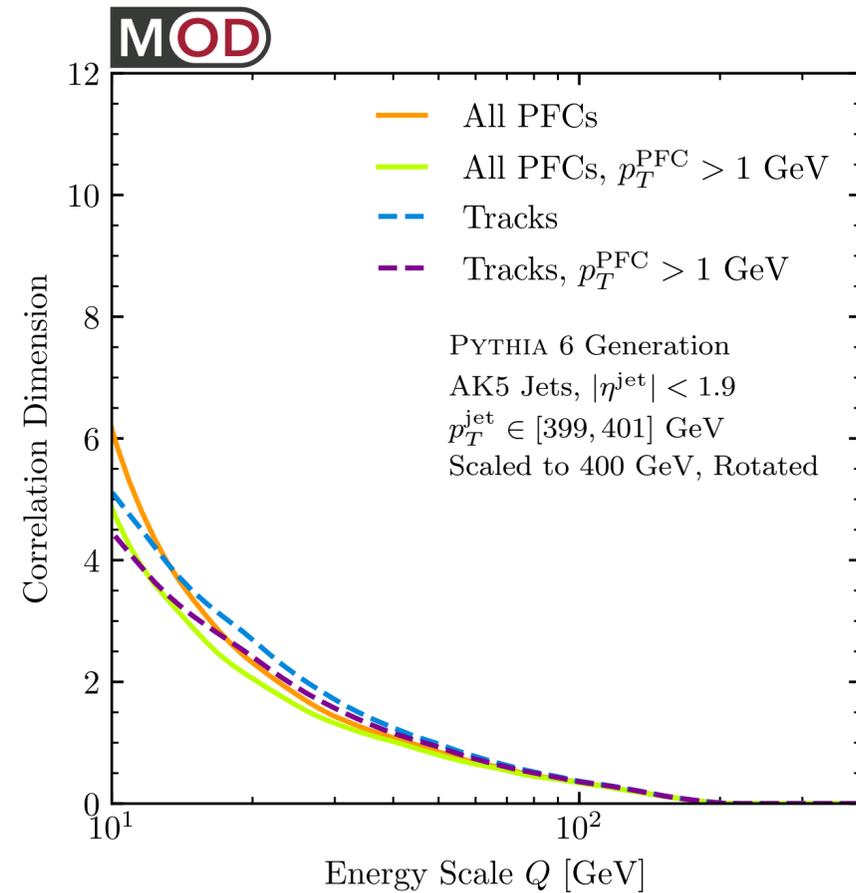
$$\text{dim}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$



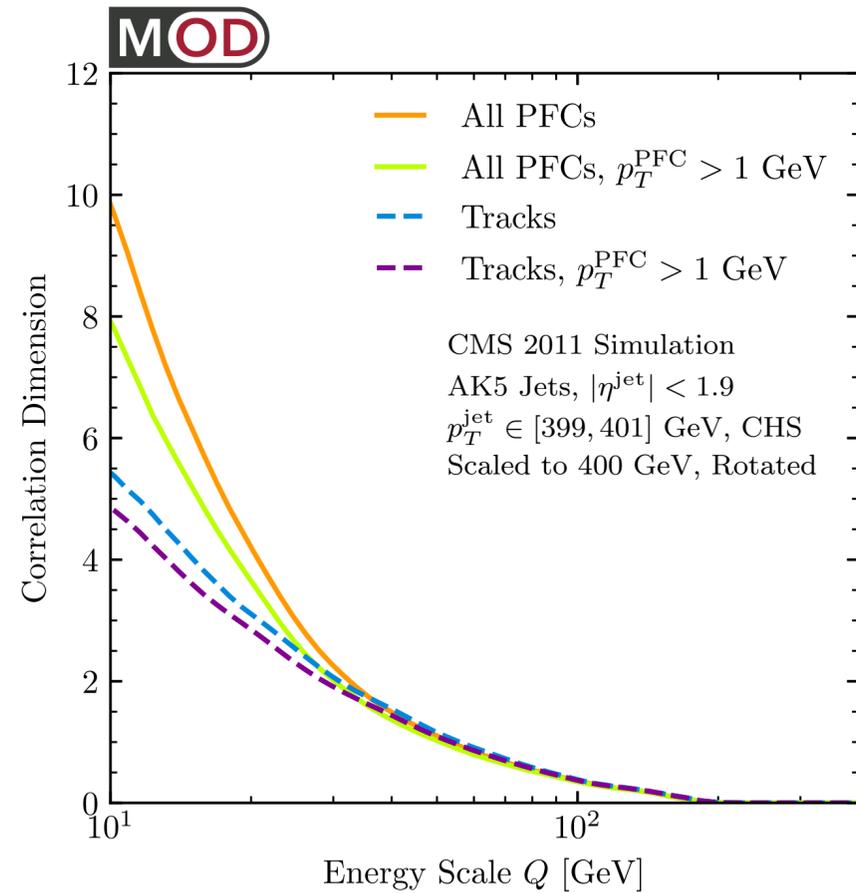
[PTK, Metodiev, Thaler, to appear soon]

# Correlation Dimension at Particle and Detector Levels

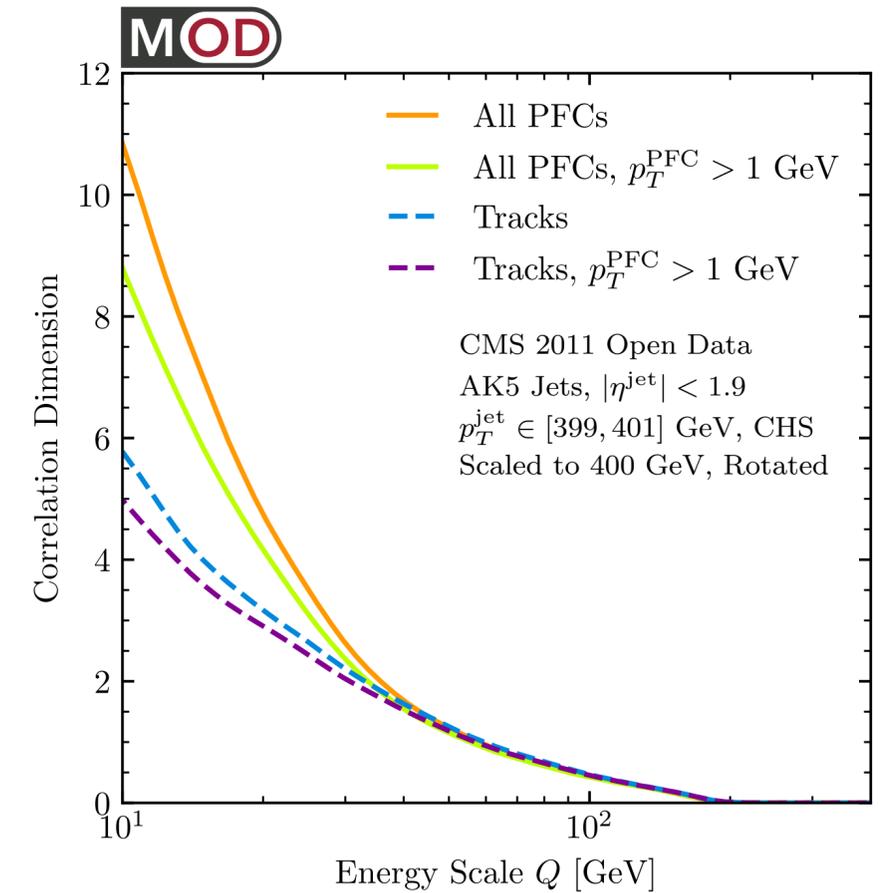
## Particle-level (PYTHIA)



## Detector-level (PYTHIA + GEANT 4)

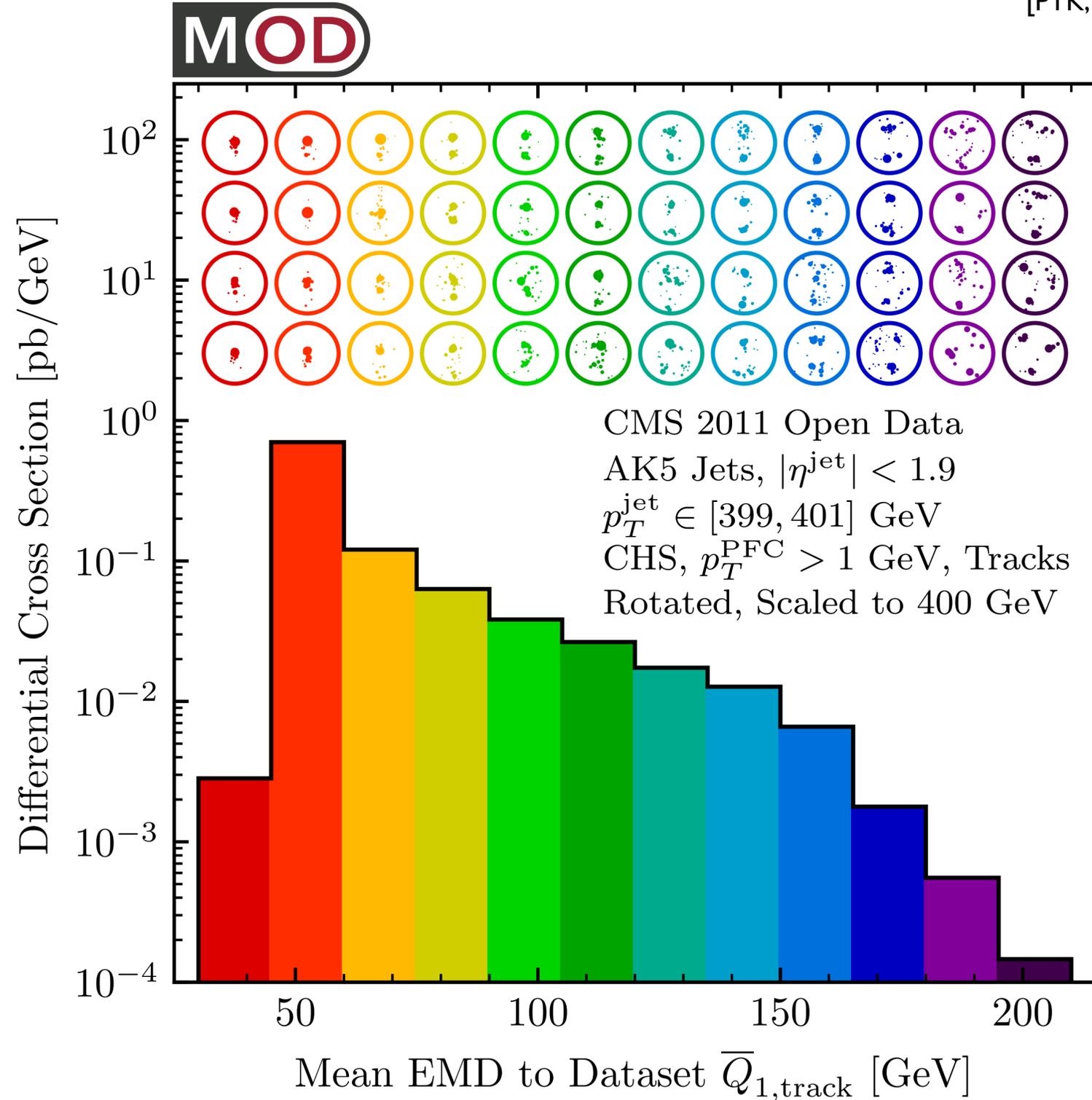


## CMS Open Data



# Visualizing Geometry in CMS Open Data

[PTK, Mastandrea, Metodiev, Naik, Thaler, PRD 2019; code and datasets at [energyflow.network](http://energyflow.network)]



## EMD for anomaly detection

← 4 medoids in each bin of anomaliness  $\bar{Q}_1$

$n^{\text{th}}$  moment of EMD distribution for a dataset

$$\bar{Q}_n(\mathcal{I}) = \sqrt[n]{\frac{1}{N} \sum_{k=1}^N (\text{EMD}(\mathcal{I}, \mathcal{J}_k))^n}$$

How far does this go?

$$\mathcal{V}_k = \frac{1}{N} \sum_{i=1}^N \min \{ \text{EMD}(\mathcal{J}_i, \mathcal{K}_1), \dots, \text{EMD}(\mathcal{J}_i, \mathcal{K}_k) \}$$

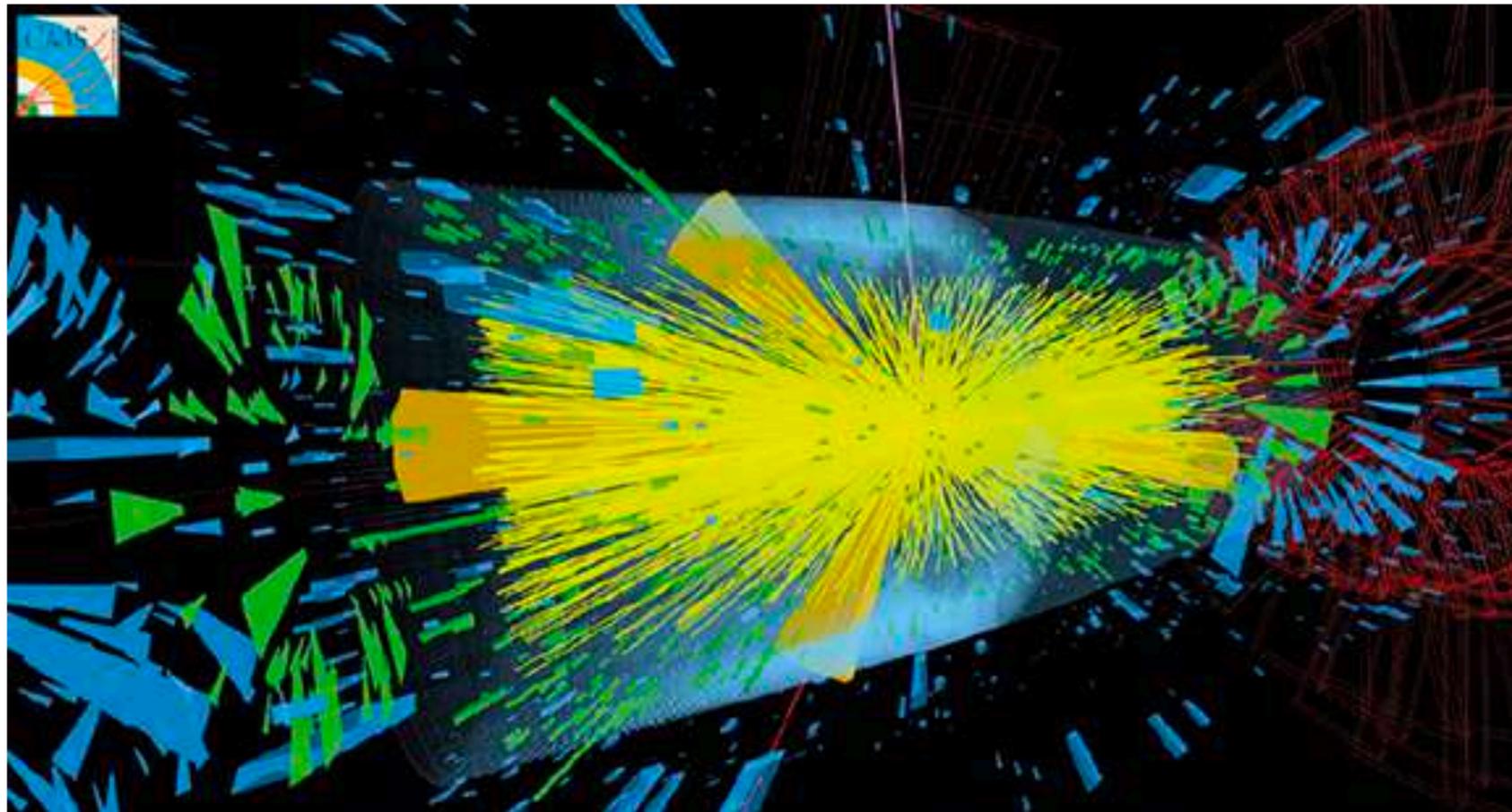
↑  
k-eventiness?

↑  
jet from dataset

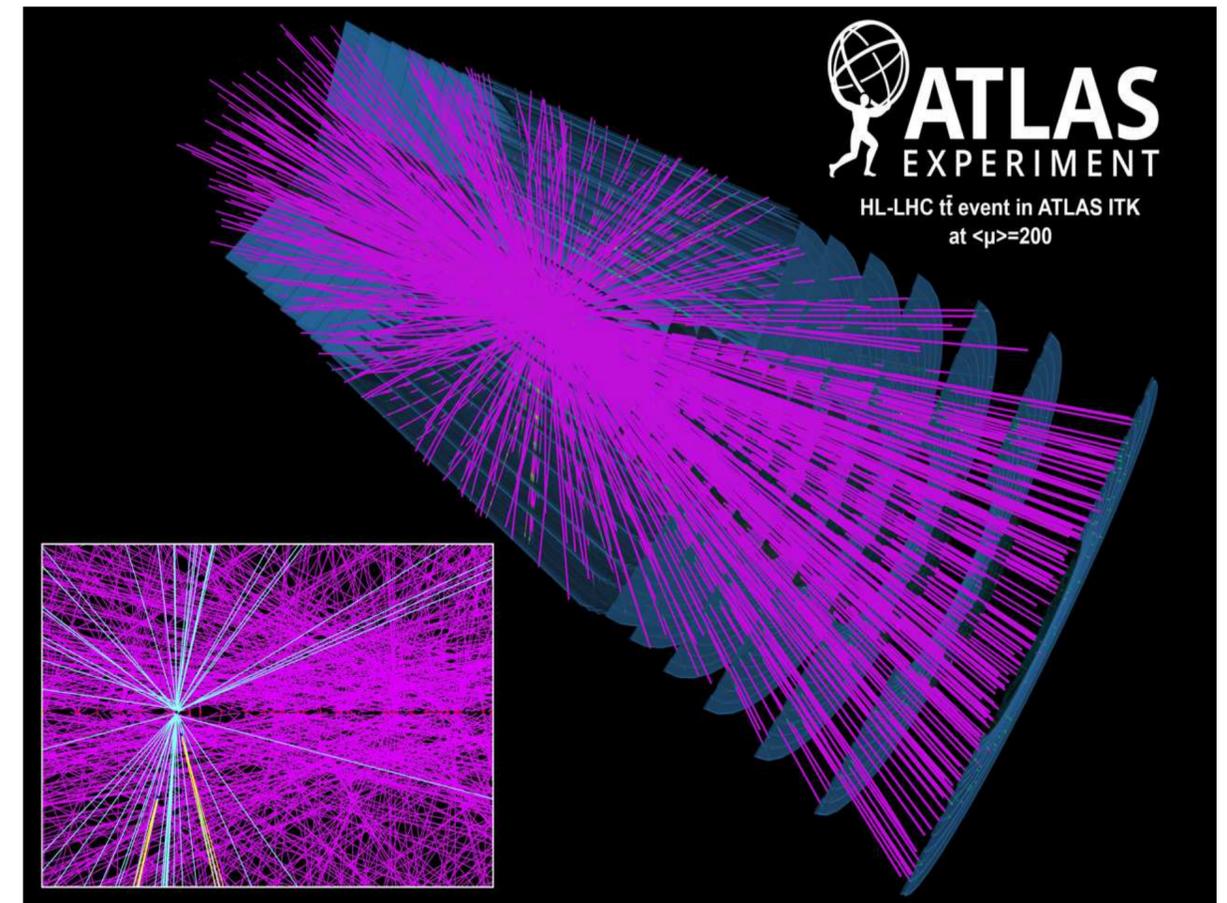
↑  
medoids

# Pileup at the (HL-)LHC

*Pileup is uniform (on average) radiation from additional proton-proton collisions*



VBF Higgs + 200 pileup vertices

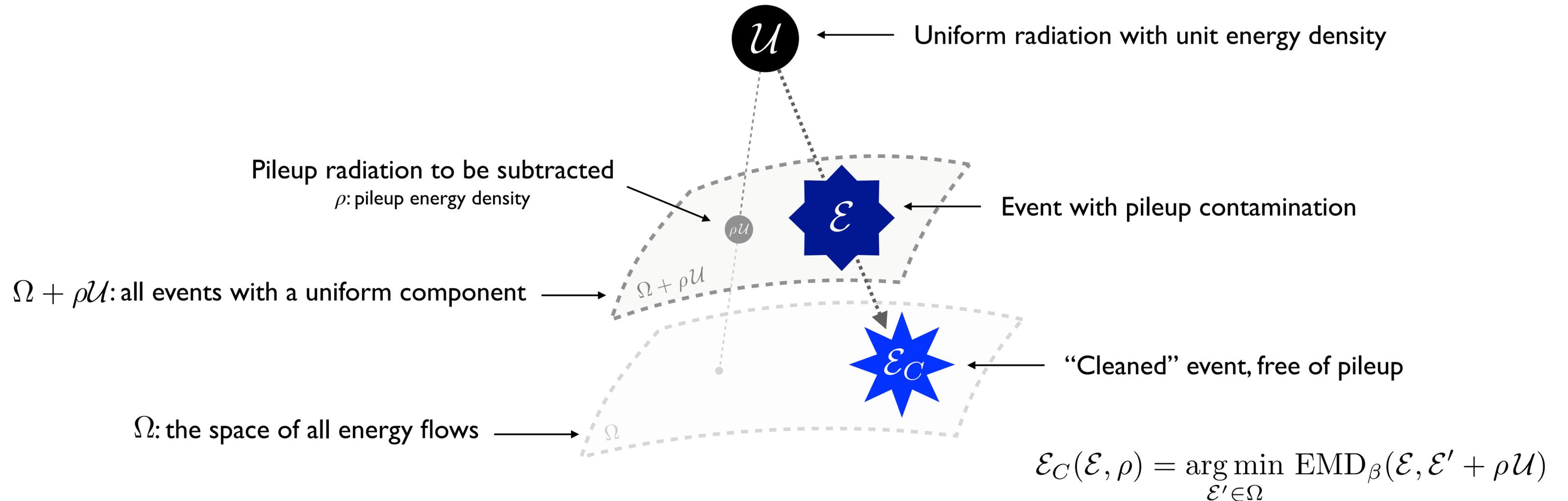


$t\bar{t}$  + 200 pileup vertices

# Pileup Mitigation in Event Space

Pileup: uniform (on average) radiation from additional proton-proton collisions

Pileup mitigation: “moving away” from the uniform event

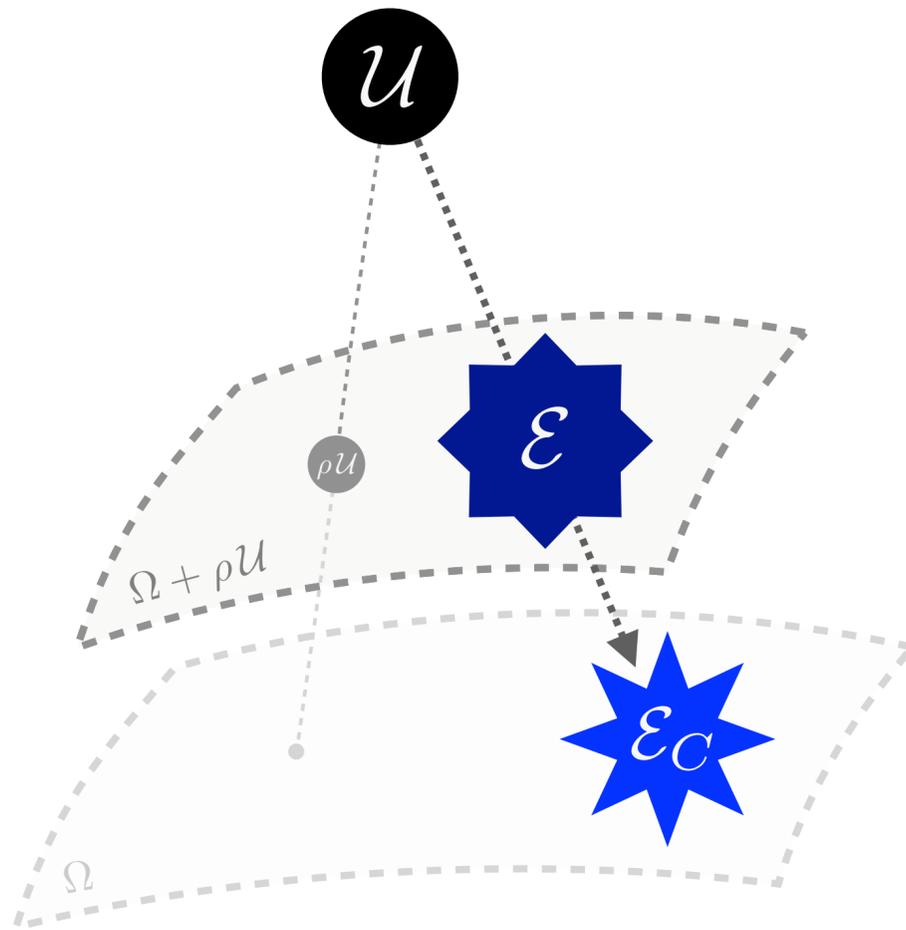


# Pileup Mitigation in Event Space

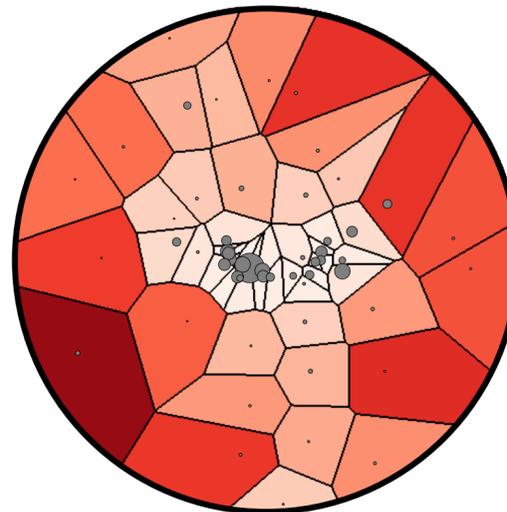
## Area subtraction

Particle energies corrected proportional to area of associated region

$$\mathcal{E}_C(\mathcal{E}, \rho) = \arg \min_{\mathcal{E}' \in \Omega} \text{EMD}_\beta(\mathcal{E}, \mathcal{E}' + \rho\mathcal{U})$$



### Voronoi

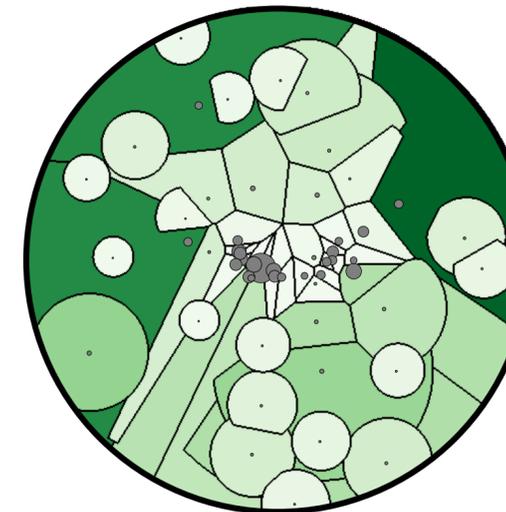


[Cacciari, Salam, Soyez, JHEP 2008]

Voronoi regions IRC unsafe

Sensitive to small modifications

### Constituent subtraction



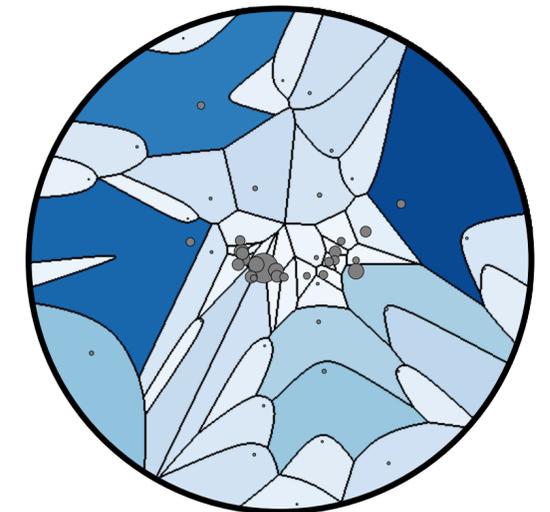
[Berta, Spouta, Miller, Leitner, JHEP 2008]

Lays down grid of “ghost” particles

Ghosts associate to nearest particle

Vanished particles don’t attract ghosts

### Apollonius



[PTK, Metodiev, Thaler, 2004.04159]

Ghosts are optimally assigned to particles by minimizing EMD

Apollonius regions have an understood continuum limit

# Beyond Observables via Weighted Cross Sections

## Standard observable (e.g. EFPs)

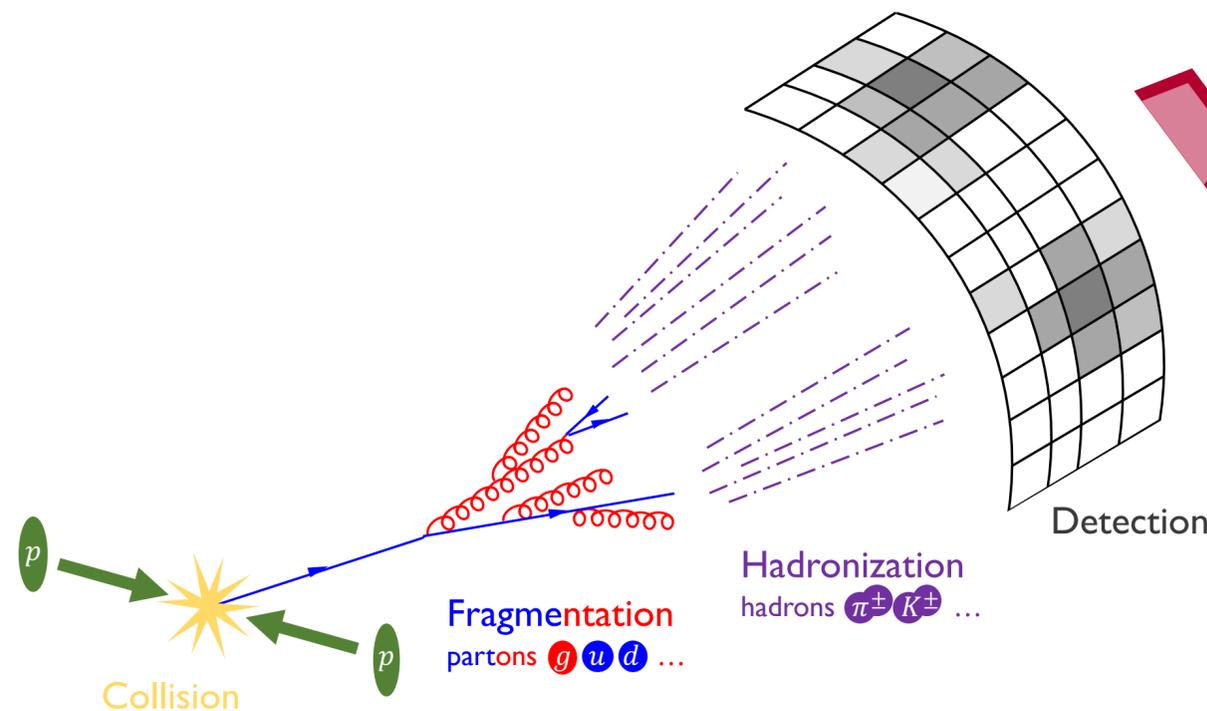
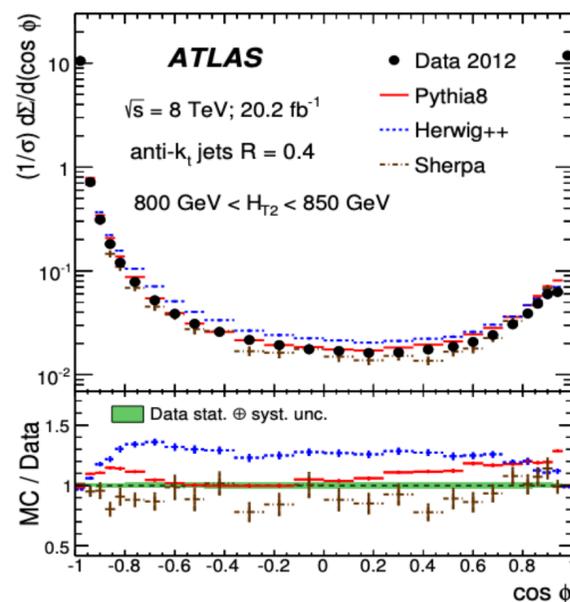
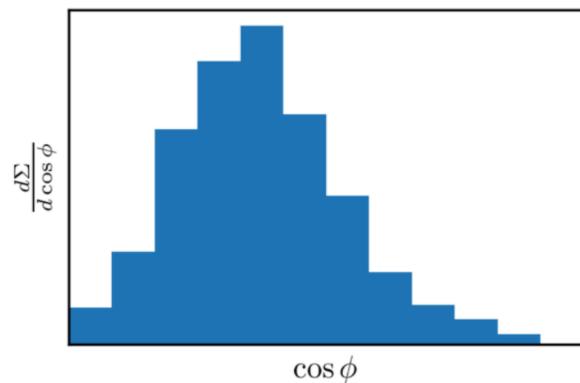
Calculate a single number for each jet/event and study distribution of values

## Weighted cross section

Calculate a distributional quantity per event and study the mean distribution

e.g. energy-energy correlator (EEC)

$$\frac{d\Sigma}{d \cos \phi} = \sum_{i,j} \int d\sigma \frac{E_i E_j}{Q^2} \delta(\cos \theta_{ij} - \cos \phi)$$



$$\mathcal{E}(\hat{n}) = \int_0^\infty dt \lim_{r \rightarrow \infty} r^2 n^i T_{0i}(t, r \hat{n})$$

Stress-energy flow

$$\frac{1}{\sigma_{\text{tot}}} \frac{d\sigma}{d\hat{n}_1 \cdots d\hat{n}_N} = \frac{\langle \mathcal{O} \mathcal{E}(\hat{n}_1) \cdots \mathcal{E}(\hat{n}_N) \mathcal{O}^\dagger \rangle}{\langle \mathcal{O} \mathcal{O}^\dagger \rangle}$$

Correlations of energy flow operators can be directly studied!

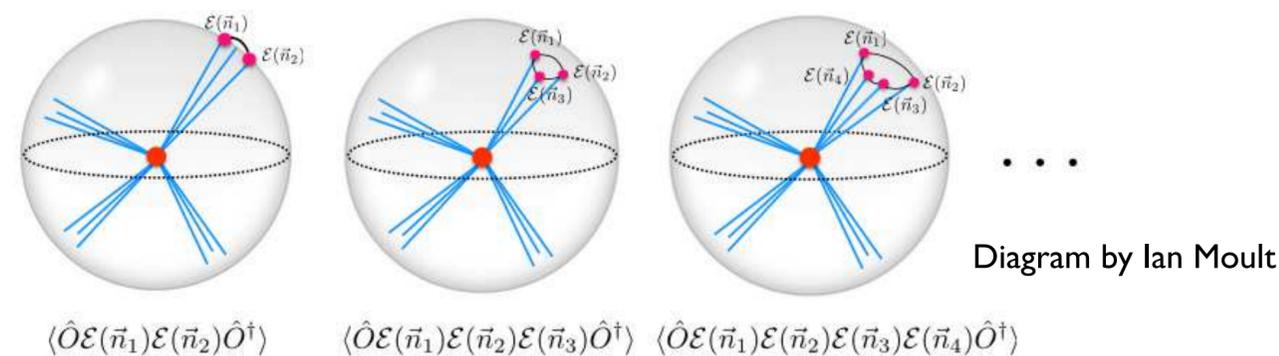


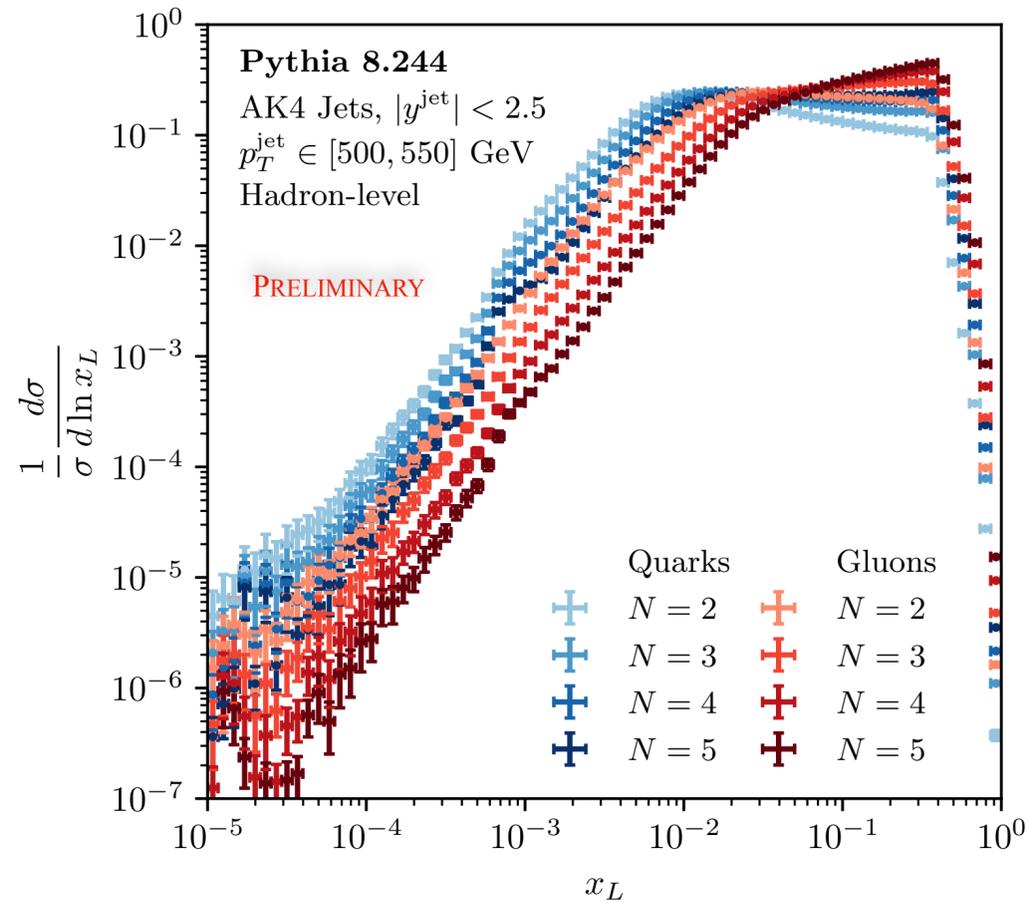
Diagram by Ian Moutl

# Energy-Energy Correlators – Projection to Longest Side

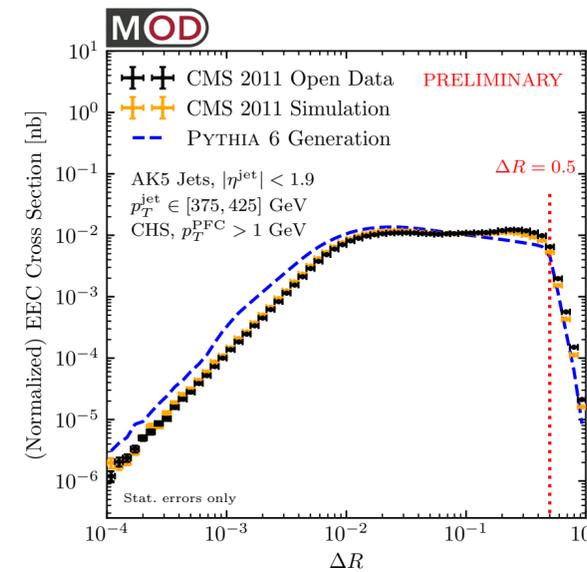
[PTK, Moul, Thaler, Zhu, to appear soon]

*Integrate out shape dependence but keep overall size dependence*

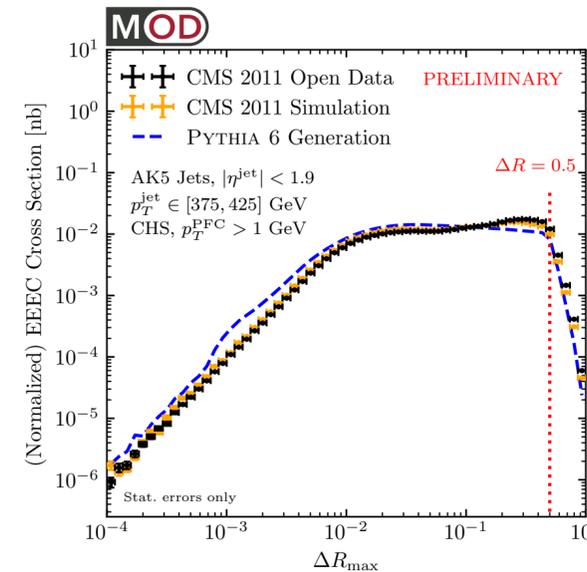
$$\frac{d\Sigma^{[N]}}{dx_L} = \sum_n \sum_{1 \leq i_1 \leq \dots \leq i_N \leq n} \int d\sigma_n \frac{E_{i_1} \cdots E_{i_N}}{Q^N} \delta(x_L - \max_{1 \leq j < k \leq N} \{\theta_{i_j i_k}\})$$



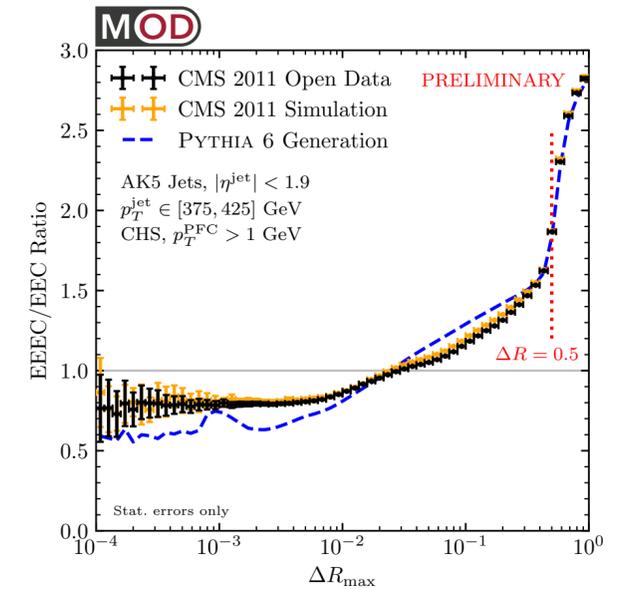
$N = 2$



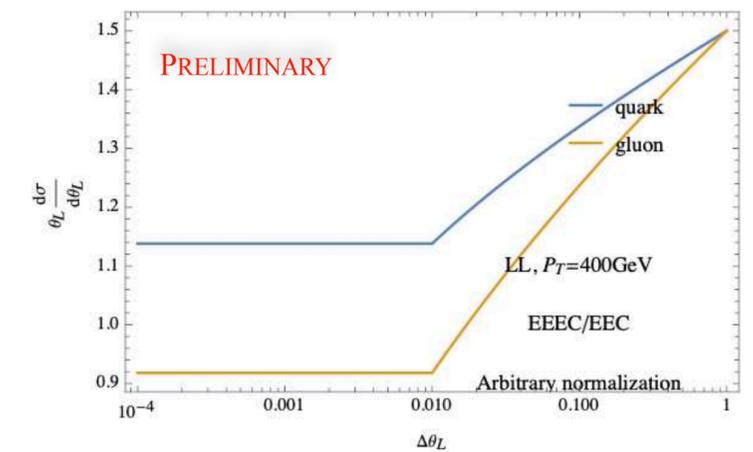
$N = 3$



EEEC/EEC Ratio



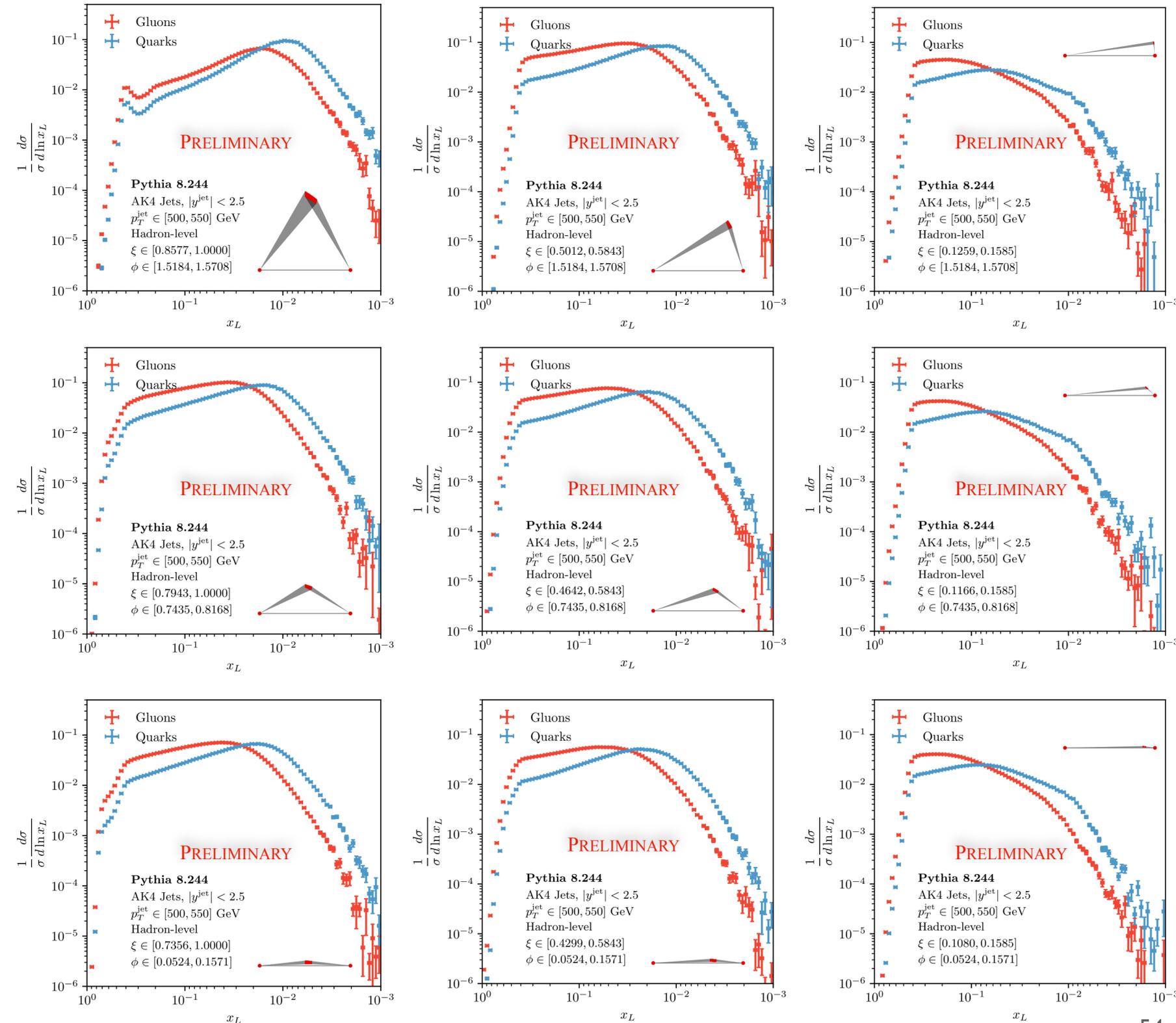
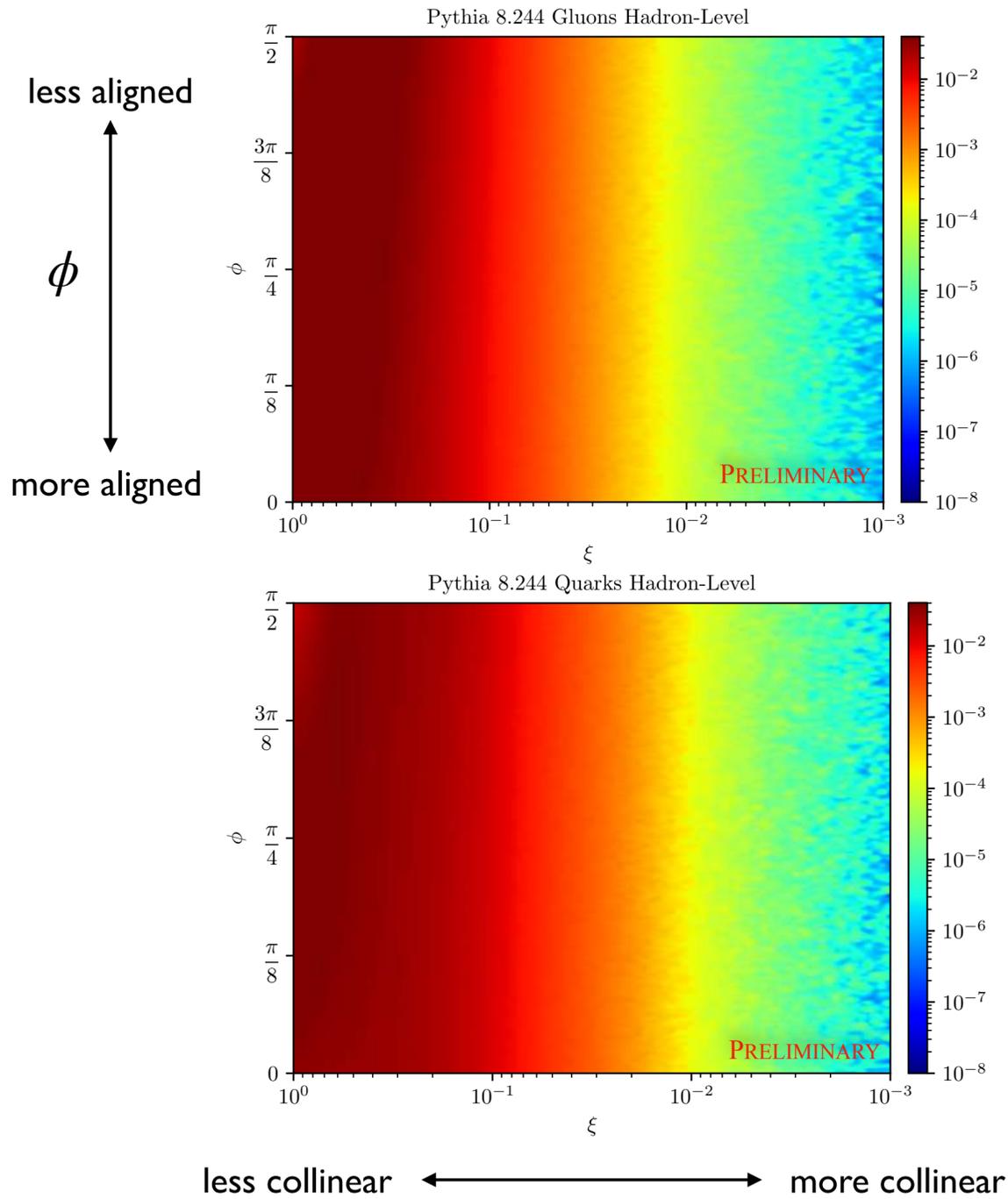
LL prediction of ratio



# EEEC – Full Shape Dependence

[PTK, Moutl, Thaler, Zhu, to appear soon]

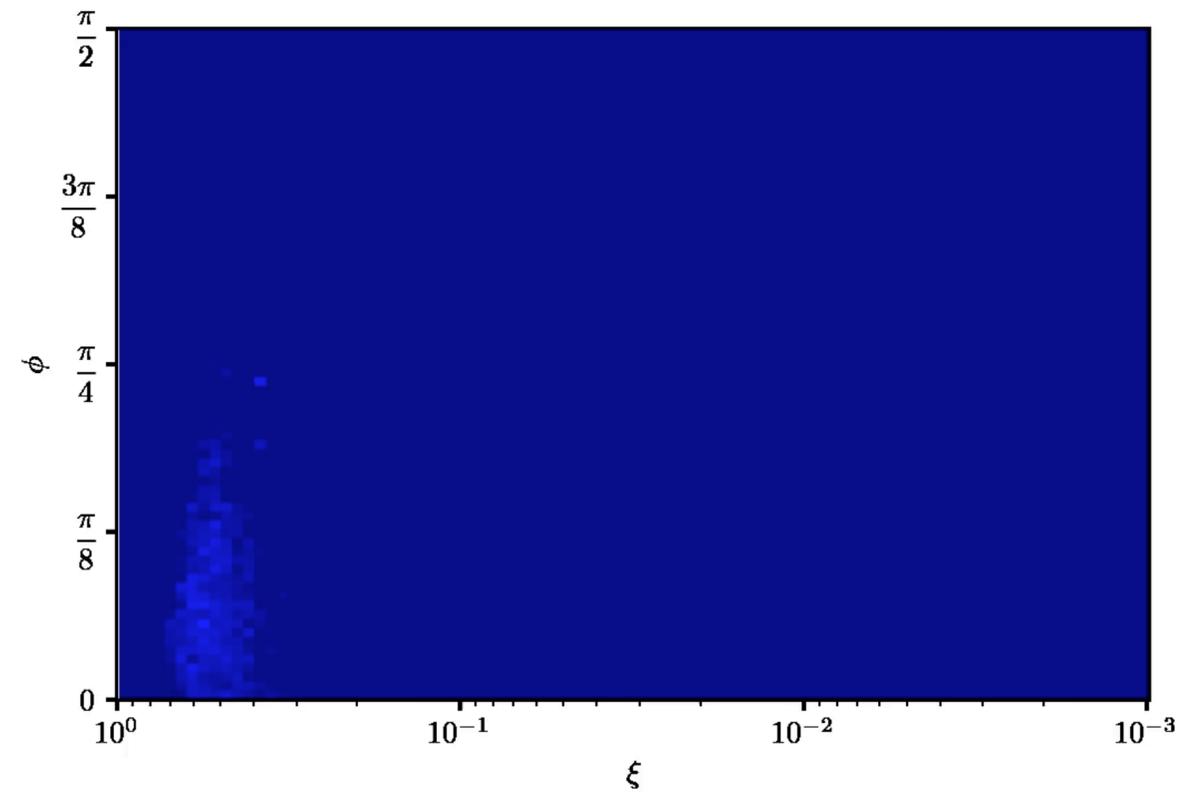
For  $x_L \sim 0.01$



# Visualizing the 3D EEEC

[PTK, Mout, Thaler, Zhu, to appear soon]

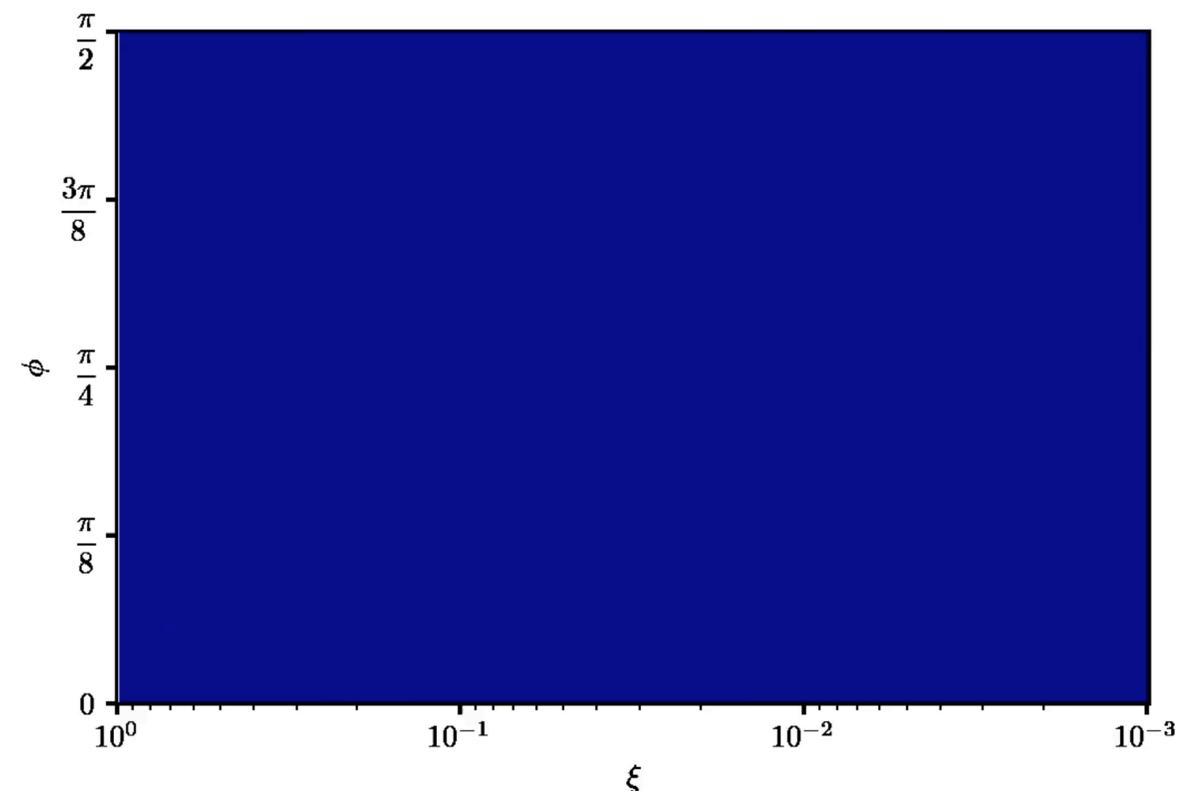
Pythia **Gluon** Jets  
 $p_T \in [500, 550]$  GeV



Time in the videos corresponds to  
 $\ln x_L$  going from 0 to  $-\infty$

Color corresponds to log of EEC  
(red is large, blue is small)

Pythia **Quark** Jets  
 $p_T \in [500, 550]$  GeV



Uniformly persistent **red** is roughly the  
perturbatively accessible region