

# The (Metric) Space of Collider Events

Patrick T. Komiske III

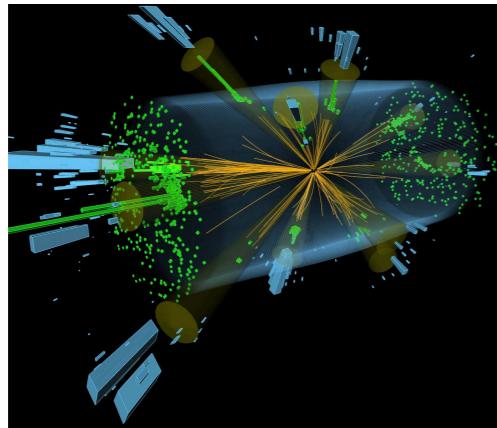
Massachusetts Institute of Technology  
Center for Theoretical Physics

*with Eric Metodiev and Jesse Thaler, [1902.02346](#)*

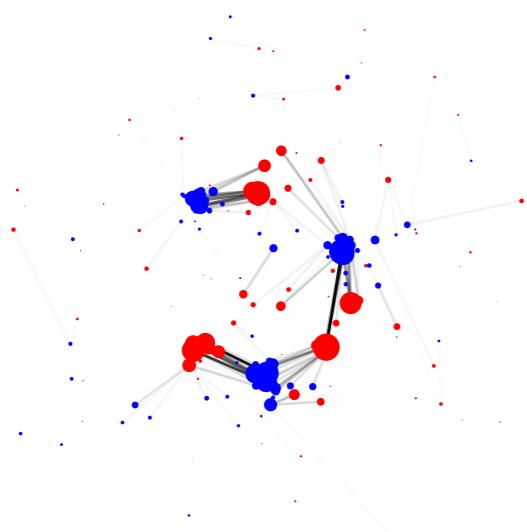
Elementary Particle Theory Seminar – Maryland Center for Fundamental Physics

University of Maryland, College Park

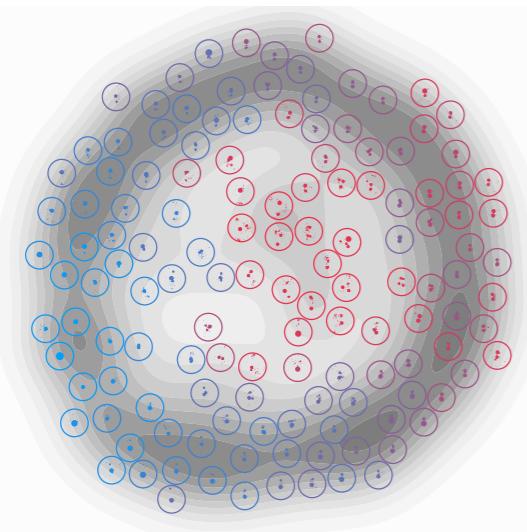
March 25, 2019



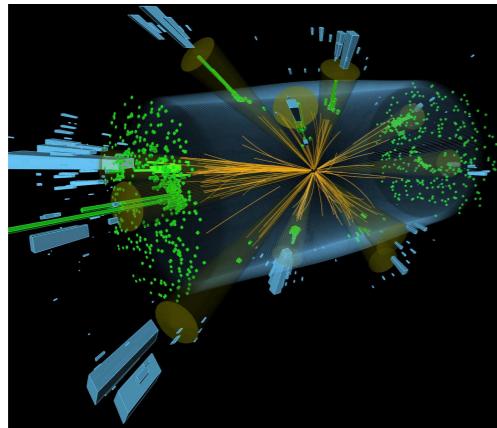
# Collider Event Foundations



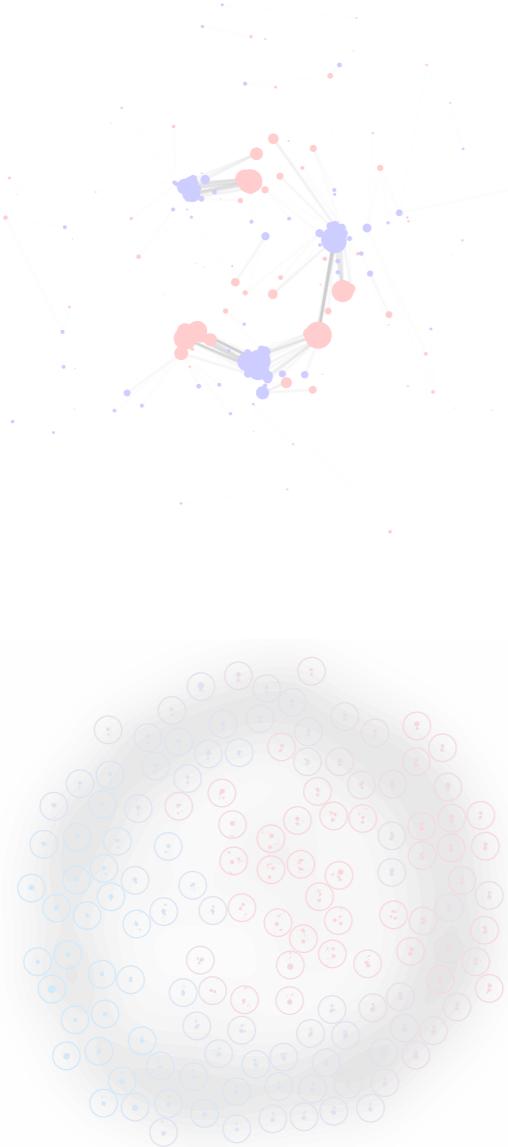
## The Energy Mover's Distance



## Particle Physics Applications



# Collider Event Foundations

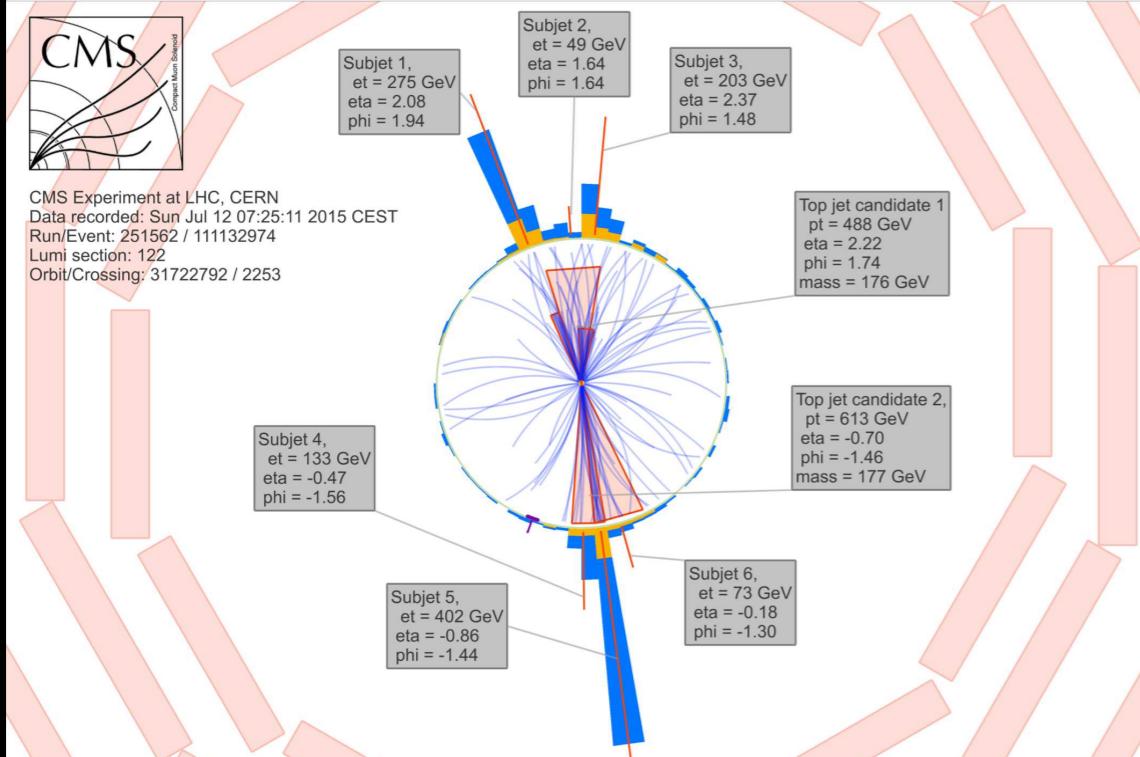
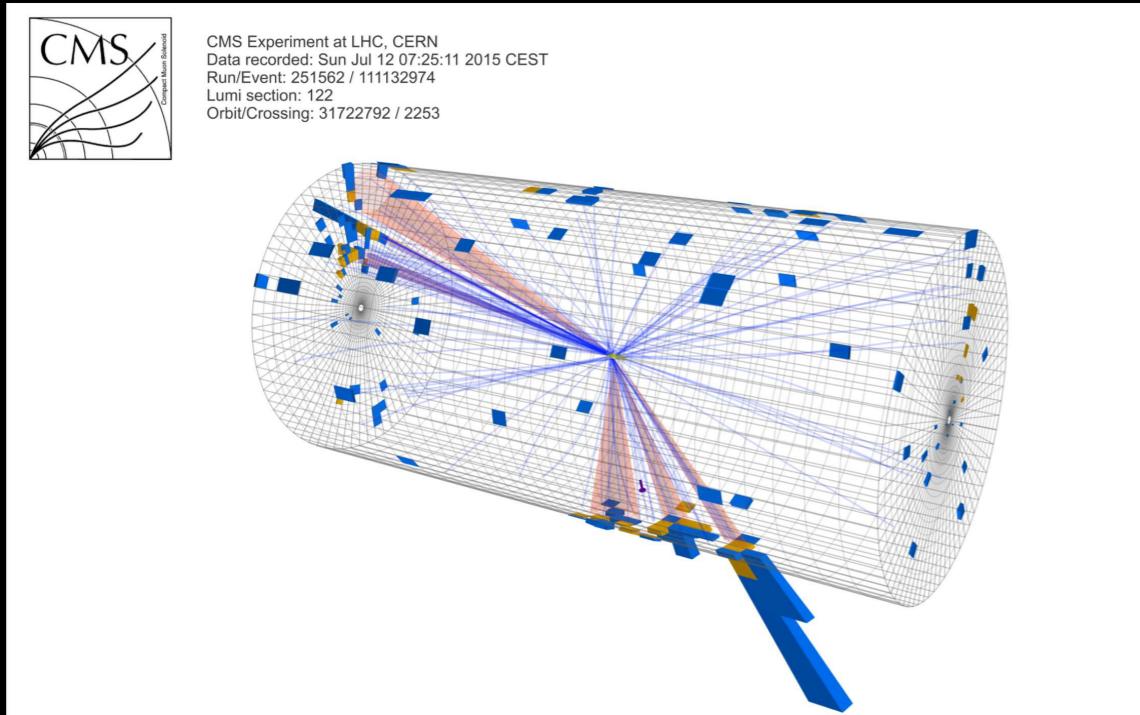


The Energy Mover's Distance

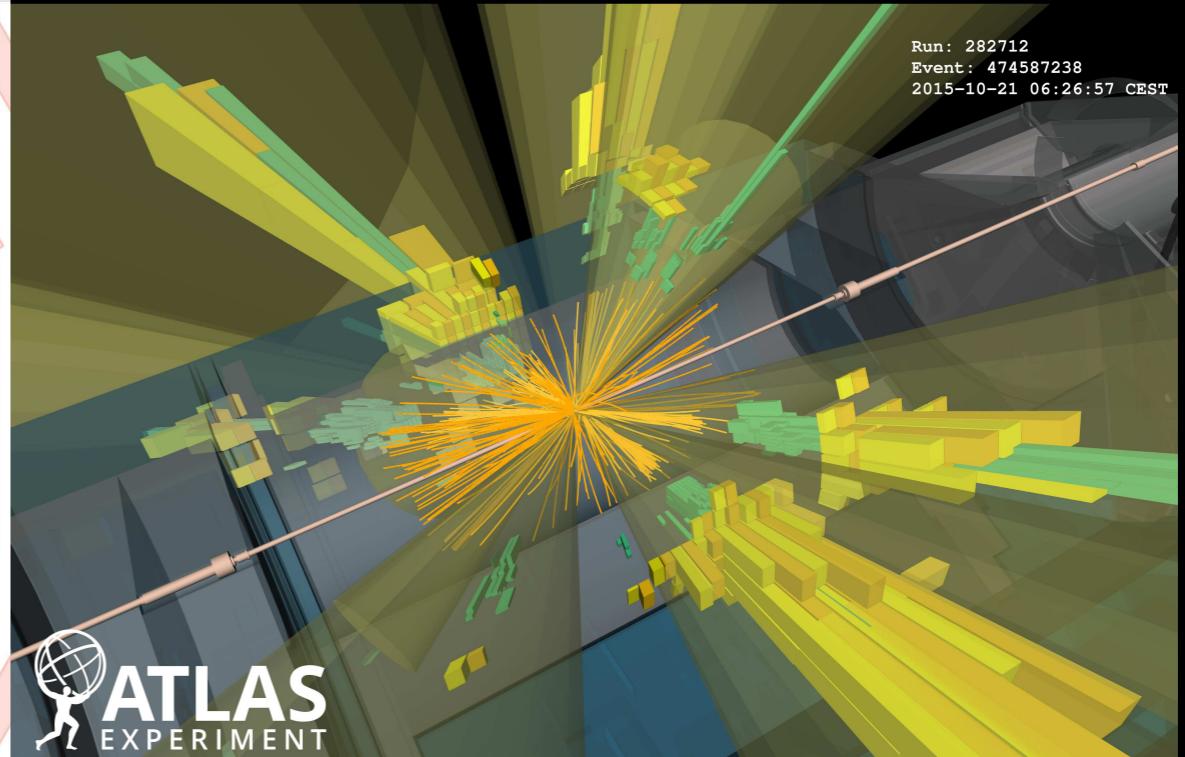
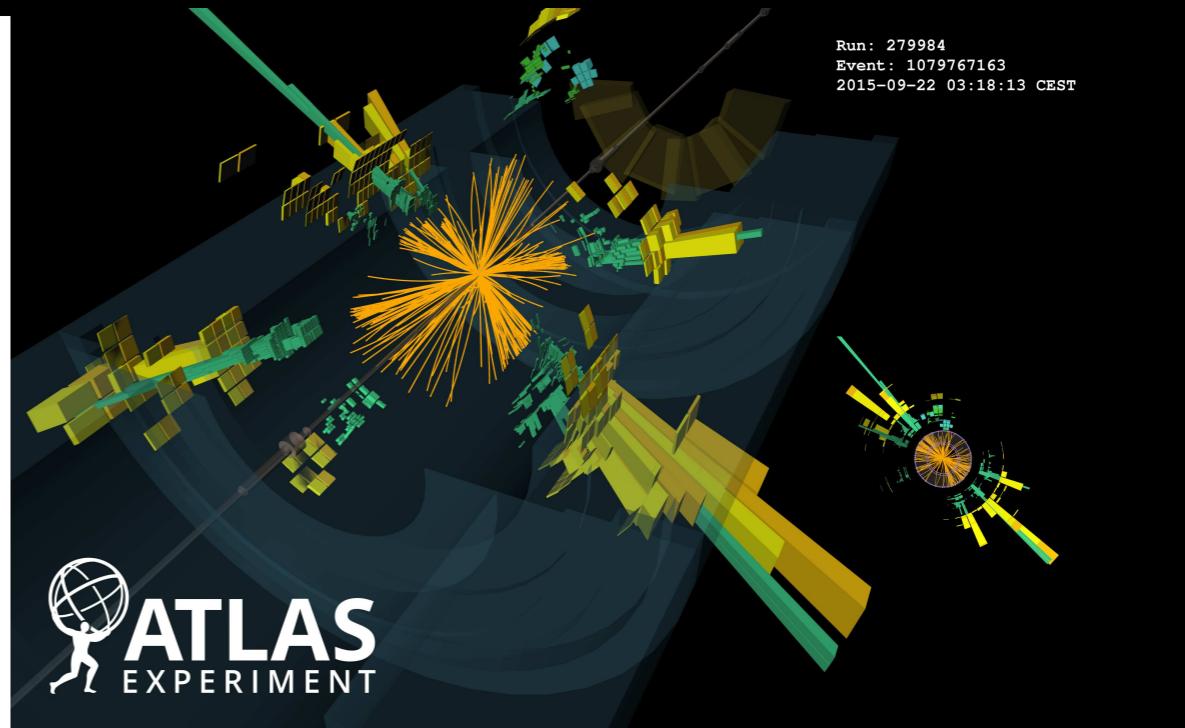
Particle Physics Applications

# Fascinating Event Topologies at the LHC

New physics searches involve complicated final states including jets (collimated sprays of hadrons)



CMS hadronic  $t\bar{t}$  event



ATLAS high jet multiplicity events

# Jet Formation in Theory

## Hard collision

*Good understanding via perturbation theory*

## Fragmentation

*Semi-classical parton shower, effective field theory*

## Hadronization

*Poorly understood (non-perturbative), modeled empirically*

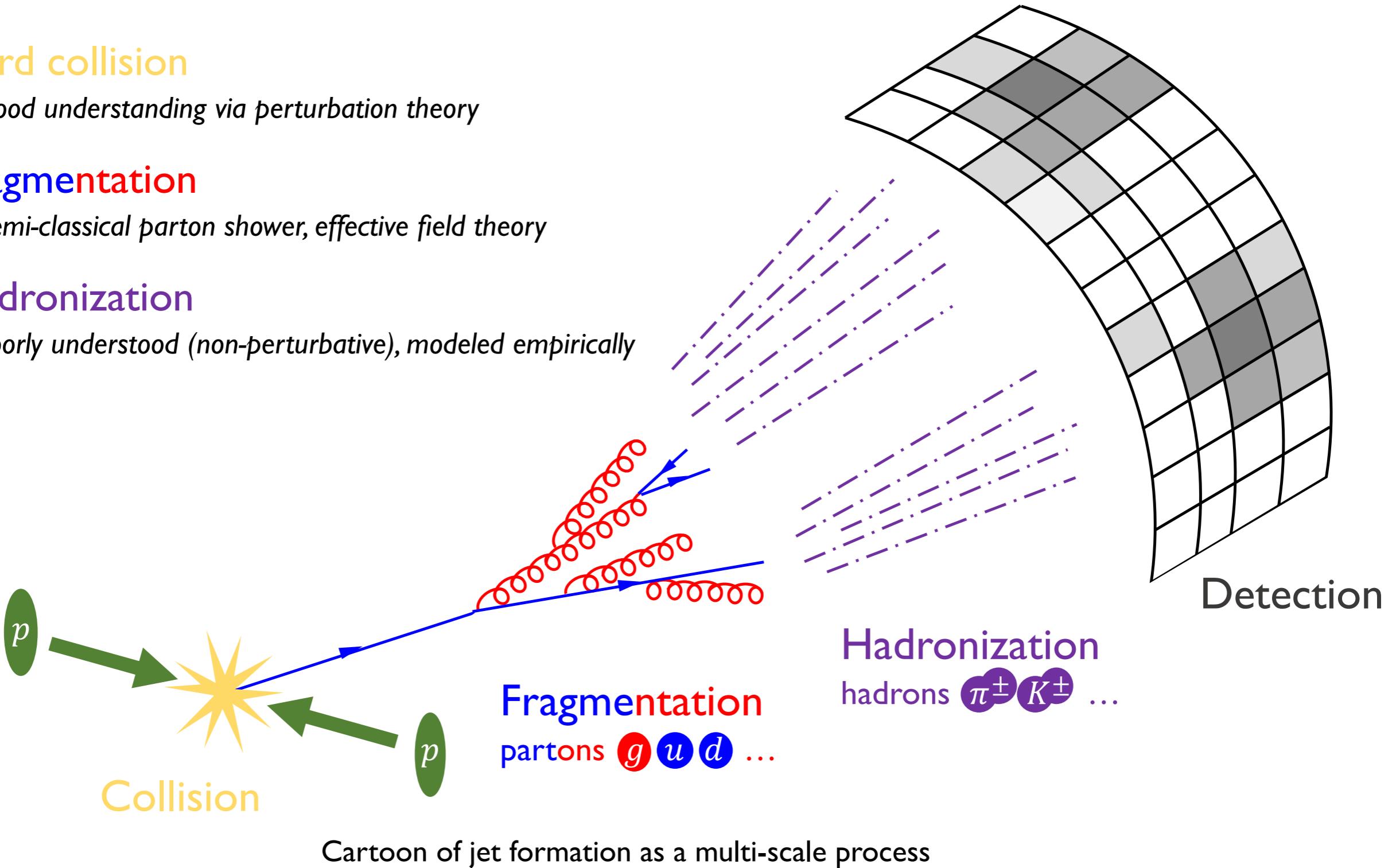
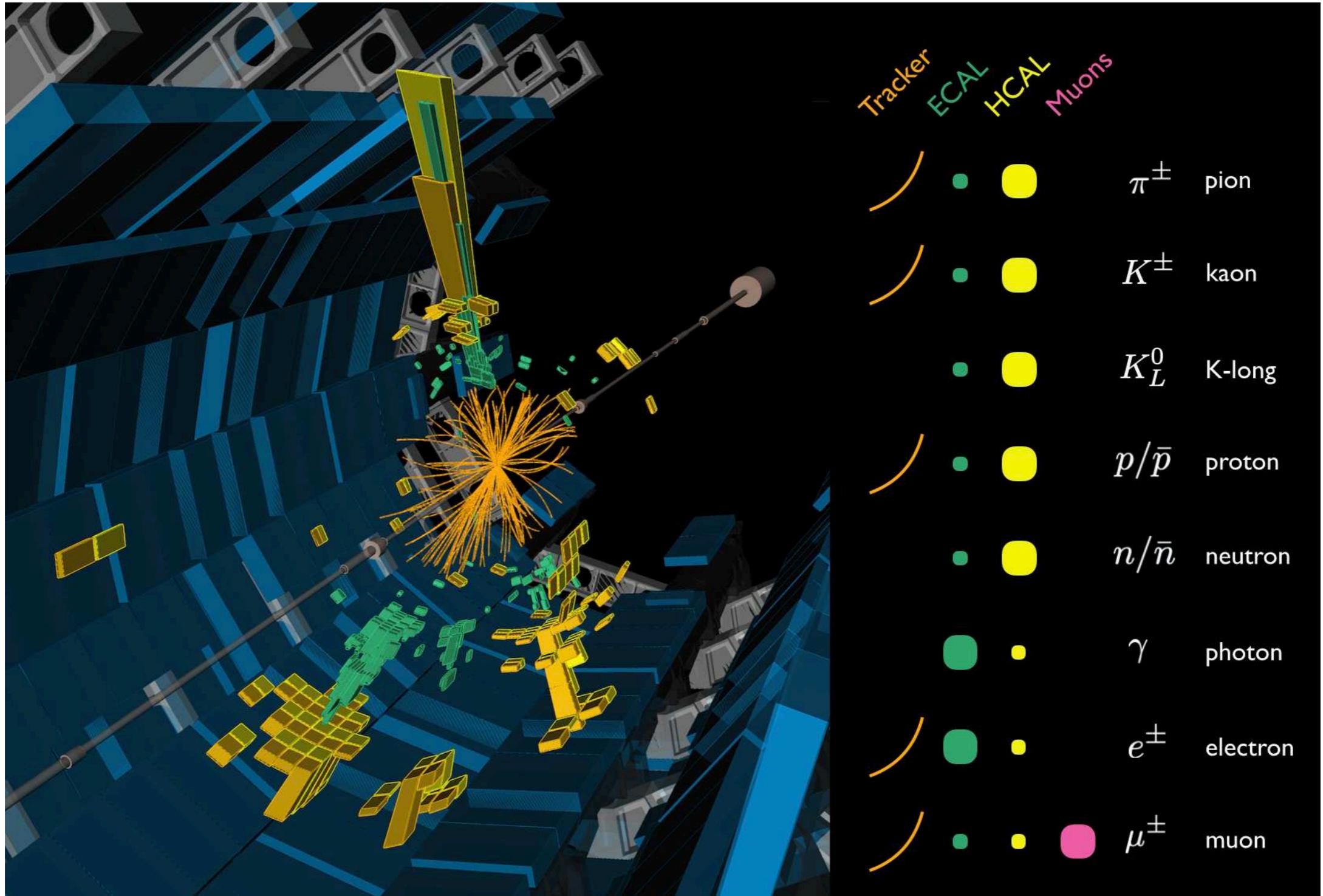


Diagram by Eric Metodiev

# Jet Detection in Experiment



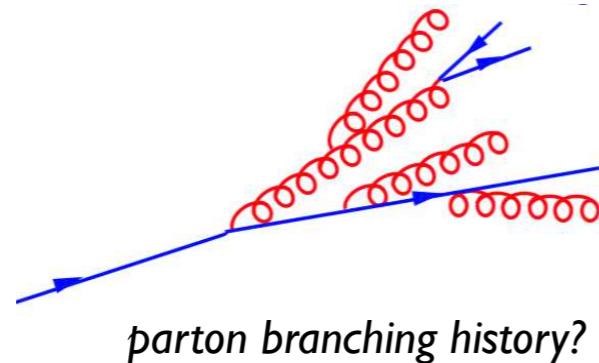
*What information is both theoretically and experimentally robust?*

# Energy Flow

## Events, Theoretically

$$|\mathcal{E}\rangle = |(p_1^\mu, \vec{q}_1); (p_2^\mu, \vec{q}_2); \dots\rangle$$

*quantum state?*

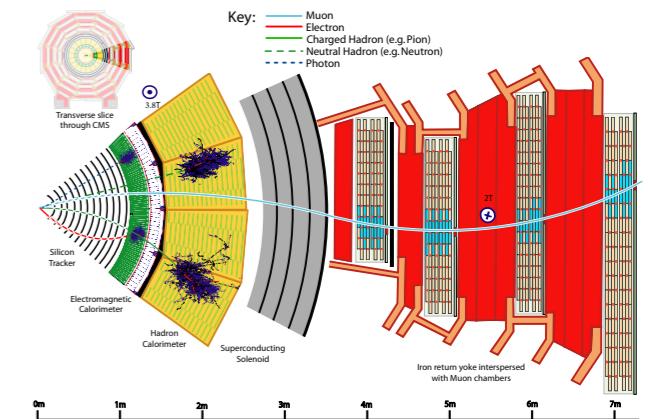


*The energy flow (distribution of energy) is robust to fragmentation, hadronization, detector effects*

## Events, Experimentally



$O(10 \text{ million})$  electrical signals?

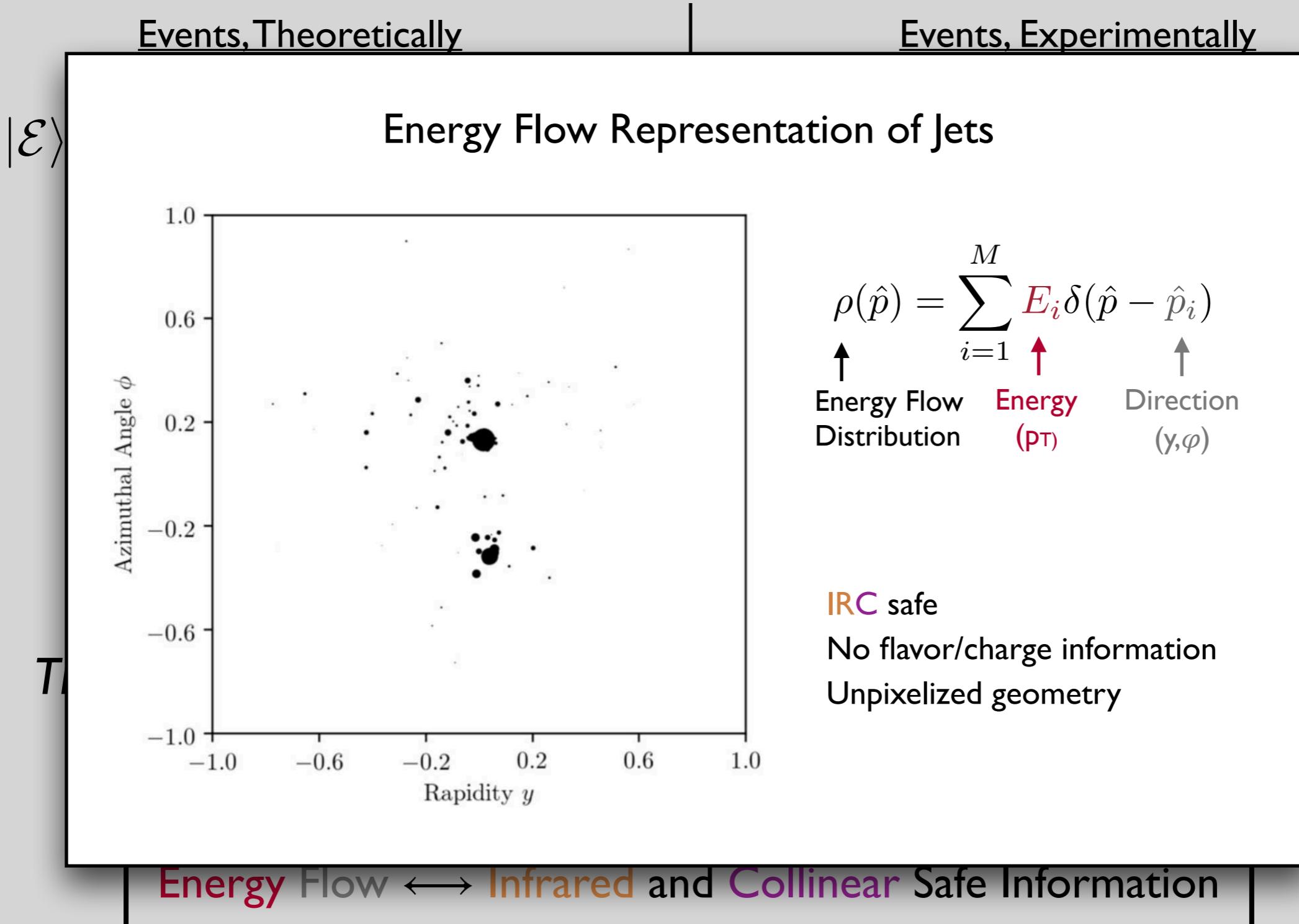


*set PF candidates?*

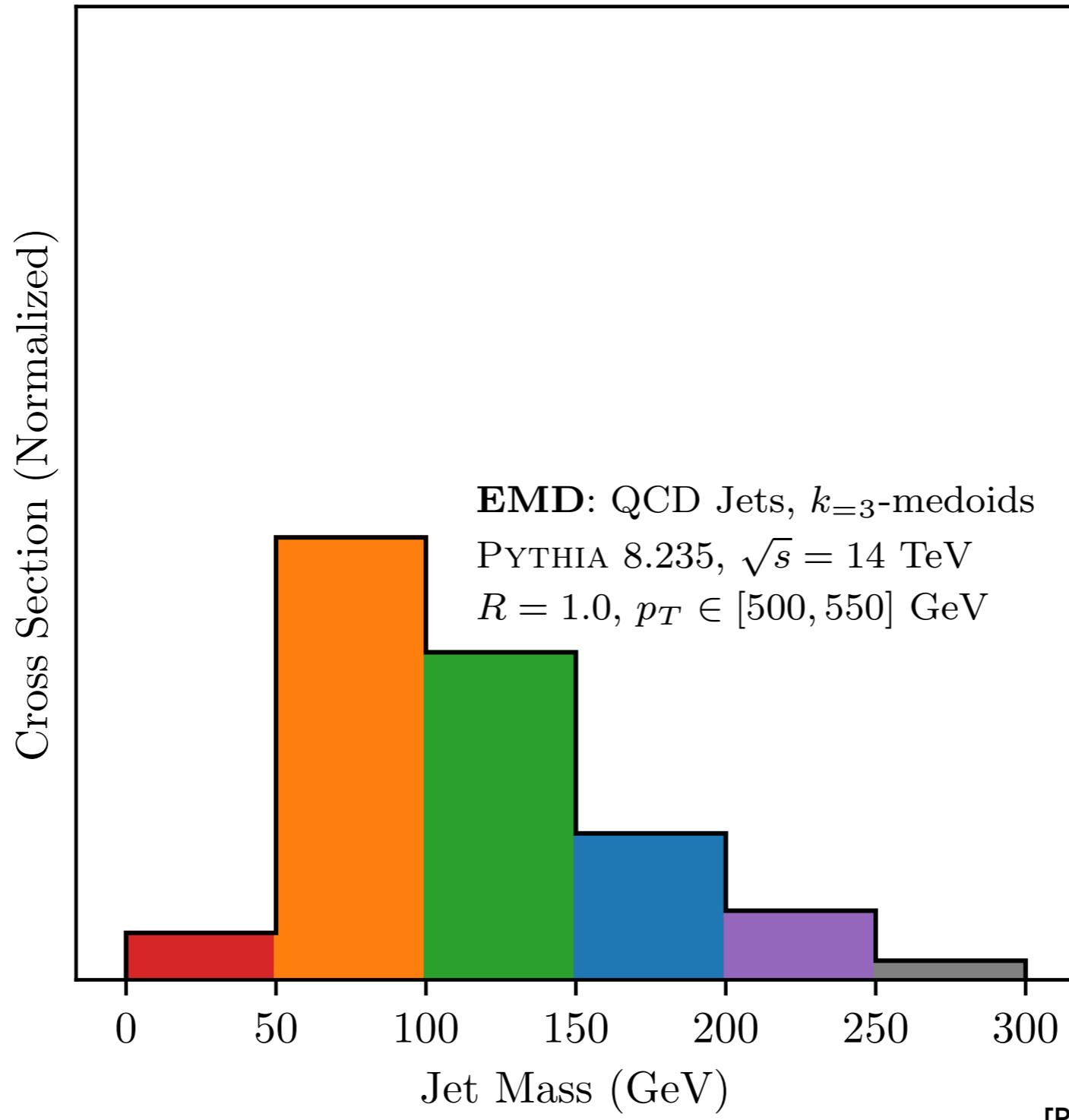


**Energy Flow  $\longleftrightarrow$  Infrared and Collinear Safe Information**

# Energy Flow

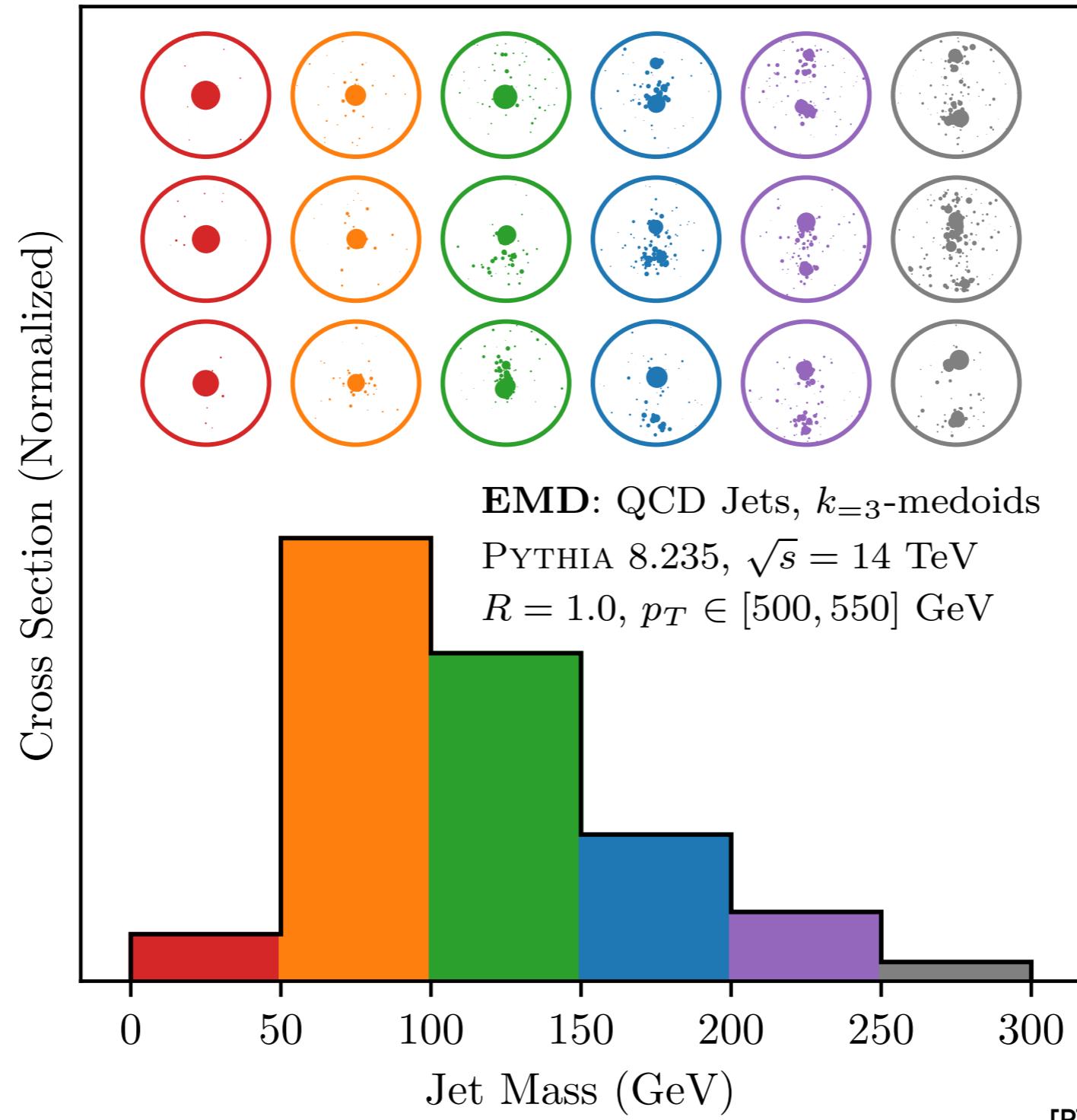


# Particle Physics Histograms



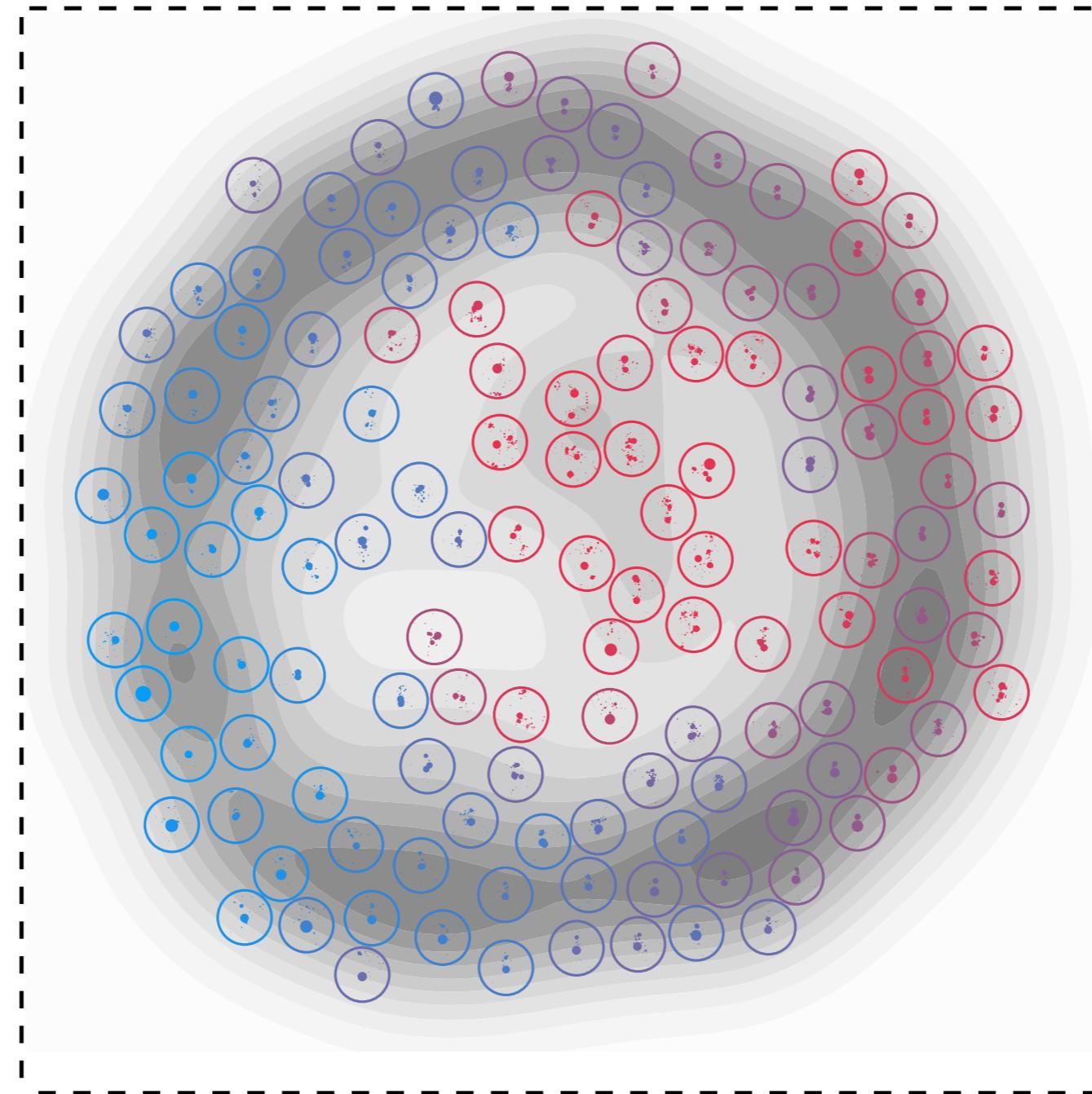
[PTK, Metodiev, Thaler, [1902.02346](#)]

# Particle Physics Histograms



[PTK, Metodiev, Thaler, [1902.02346](#)]

# Boosted W Jets

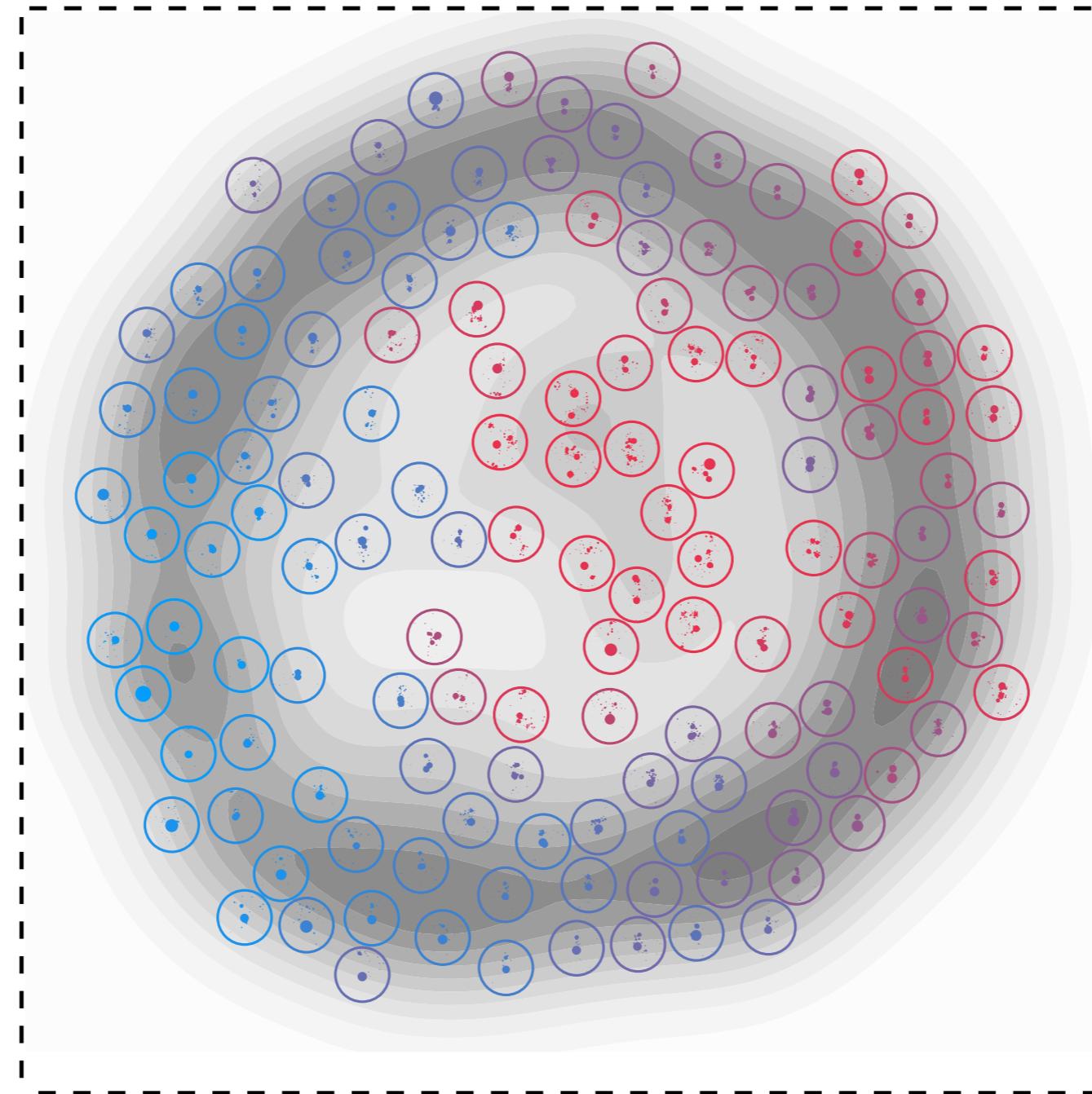


*Abstract space of W jets*

[PTK, Metodiev, Thaler, [1902.02346](#)]

# Boosted W Jets

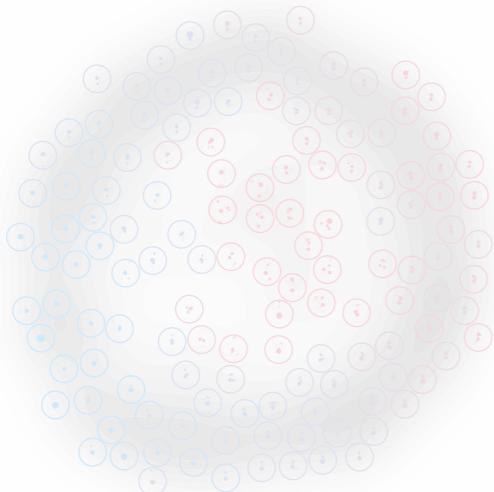
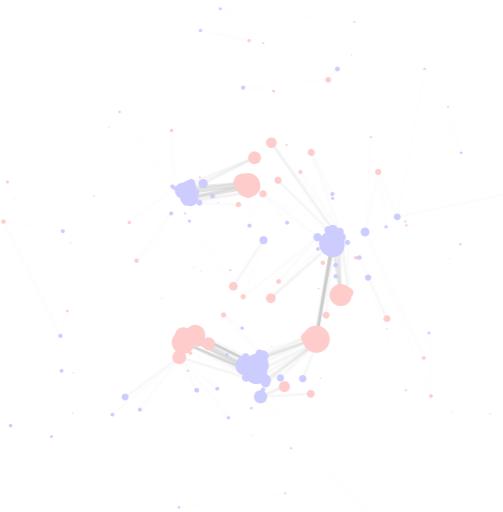
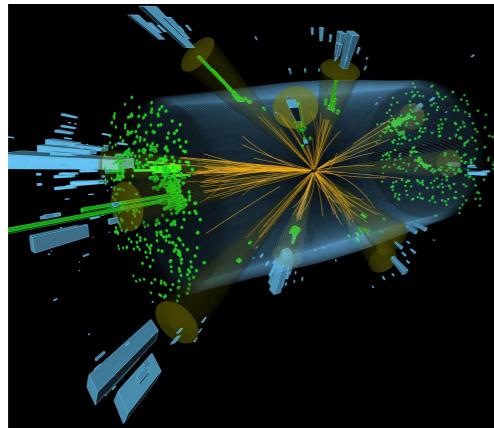
Gray contours represent the density of jets



Each circle is a particular  $W$  jet

*Abstract space of  $W$  jets*

[PTK, Metodiev, Thaler, [1902.02346](#)]

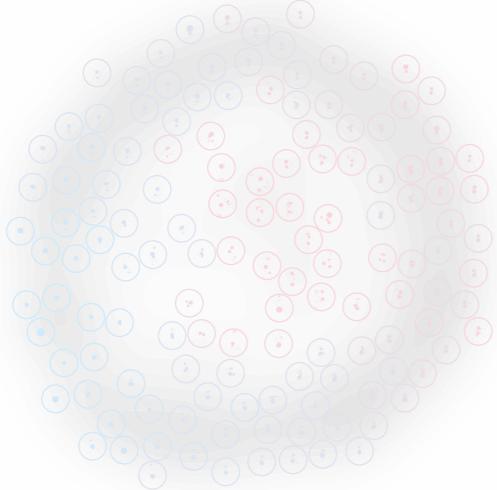
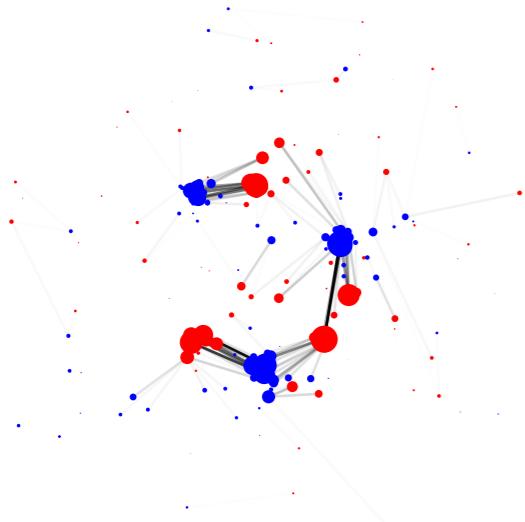
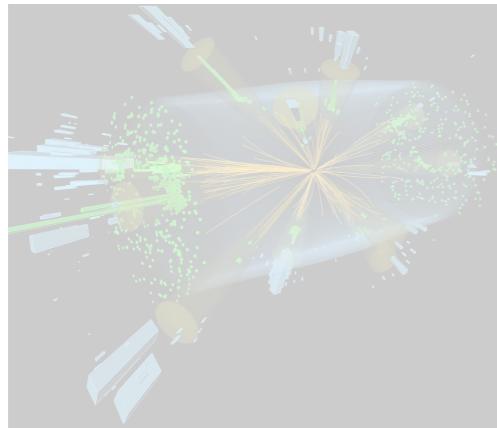


# Collider Event Foundations

*IRC-safe energy flow is theoretically and experimentally robust*

The Energy Mover's Distance

Particle Physics Applications



## Collider Event Foundations

*IRC-safe energy flow is theoretically and experimentally robust*

## The Energy Mover's Distance

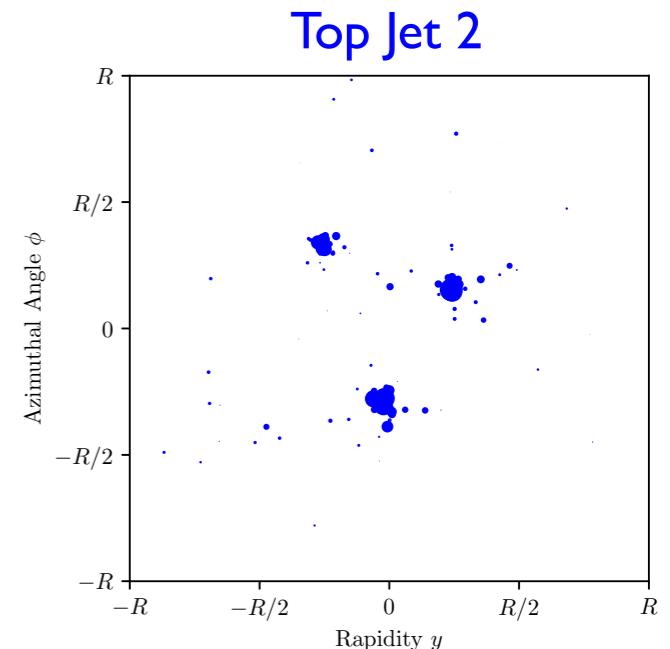
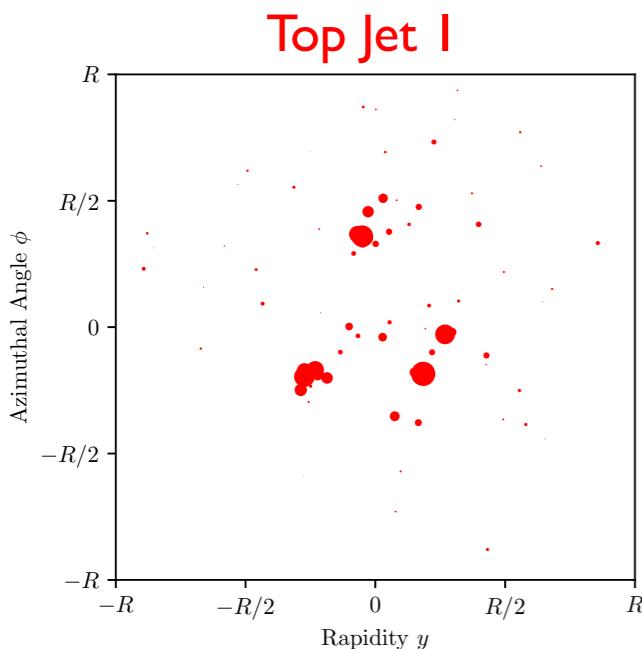
## Particle Physics Applications

# The Earth Mover's Distance

A metric on normalized distributions in a space with a ground distance measure

↳ symmetric, non-negative, triangle inequality, zero iff identical

*The minimum "work" (stuff x distance) required to transport supply to demand*



Related to *optimal transport* theory – commonly used as a metric on the space of images

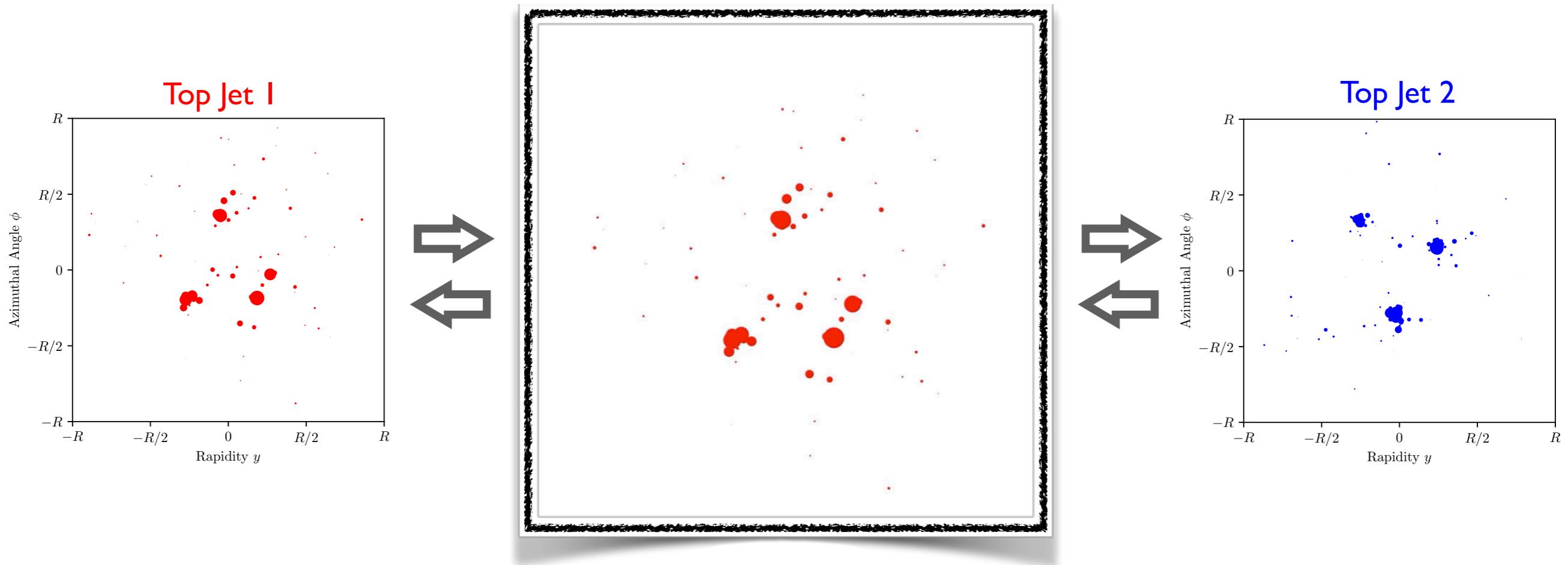
[Peleg, Werman, Rom, [IEEE 1989](#); Rubner, Tomasi, Guibas, [ICCV 1998](#), [ICJV 2000](#); Pele, Werman, [ECCV 2008](#); Pele, Taskar, [GSI 2013](#)]

# The Earth Mover's Distance

A metric on **normalized distributions** in a space with a **ground distance measure**

↳ symmetric, non-negative, triangle inequality, zero iff identical

*The minimum "work" (**stuff** x **distance**) required to transport **supply** to **demand***



Related to **optimal transport** theory – commonly used as a metric on the space of images

[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

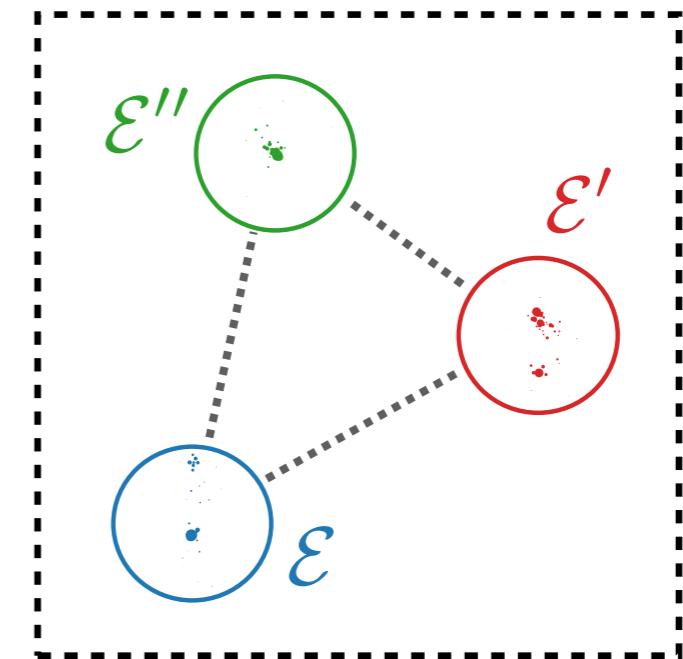
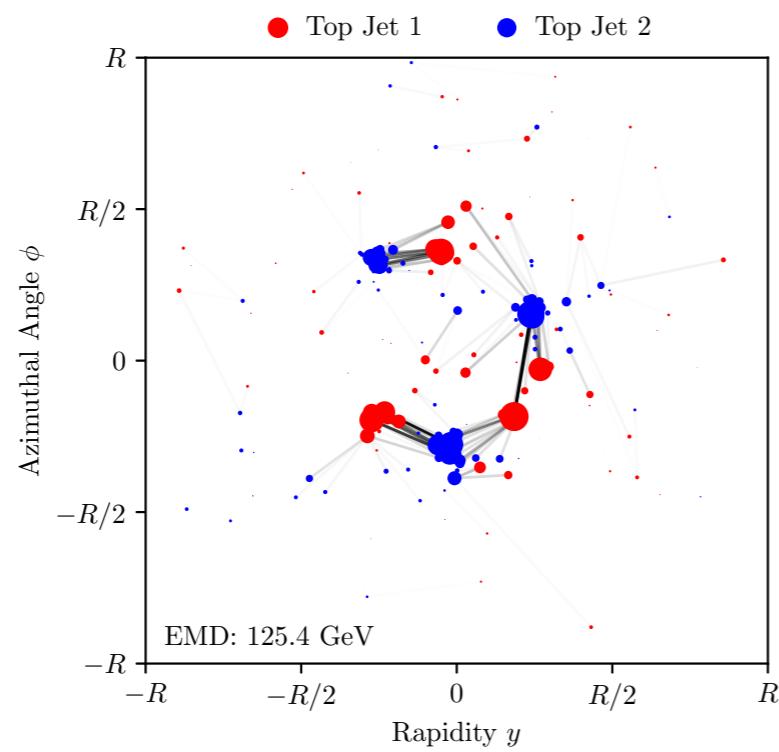
[PTK, Metodiev, Thaler, 1902.02346]

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

$$\text{EMD}(\mathcal{E}, \mathcal{E}') = \min_{\{f_{ij} \geq 0\}} \sum_i \sum_j f_{ij} \frac{\theta_{ij}}{R} + \left| \sum_i E_i - \sum_j E'_j \right|$$

Cost of optimal transport      Cost of energy creation

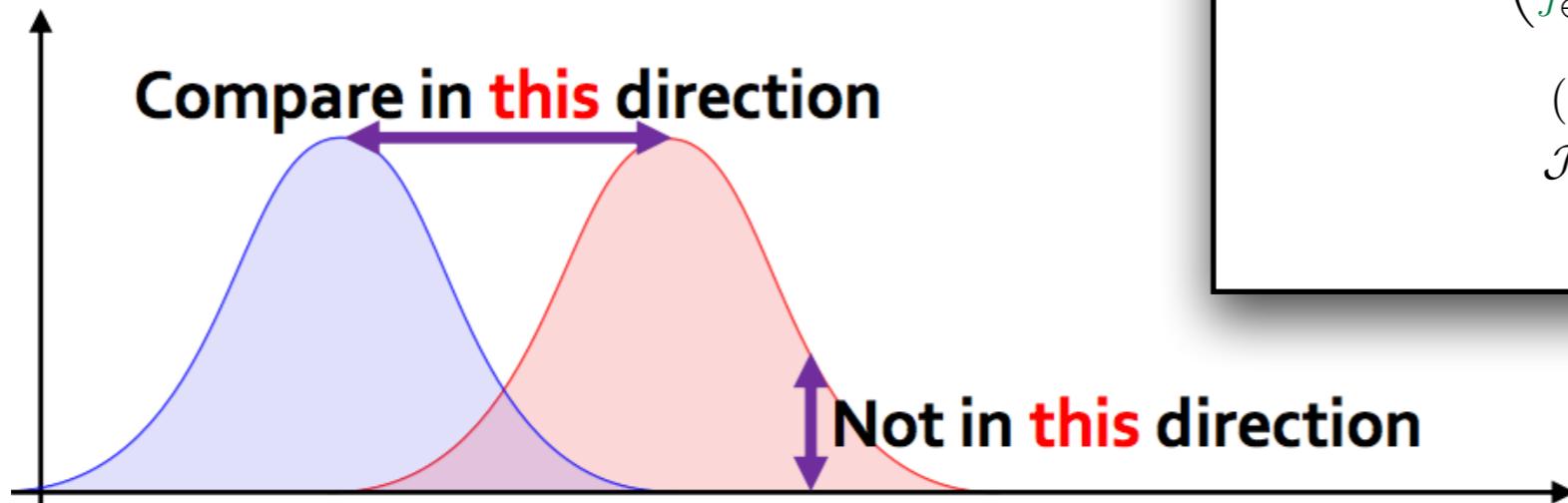
$$\sum_j f_{ij} \leq E_i, \quad \sum_i f_{ij} \leq E'_j, \quad \sum_{ij} f_{ij} = \min \left( \sum_i E_i, \sum_j E'_j \right)$$



Triangle inequality satisfied for  $R \geq d_{\max}/2$   
 $0 \leq \text{EMD}(\mathcal{E}, \mathcal{E}') \leq \text{EMD}(\mathcal{E}, \mathcal{E}'') + \text{EMD}(\mathcal{E}'', \mathcal{E}')$

# Mathematics of the Earth Mover's Distance

p-Wasserstein distance is a metric on probability distributions



[figure from Kun, [Math n Programming](#)]

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

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

Earth mover's distance is **1-Wasserstein** metric on discrete distributions

## Recent use in Machine Learning

### *Wasserstein Generative Adversarial Networks*

[Arjovsky, Chintala, Bottou, [1701.07875](#);

in particle physics:

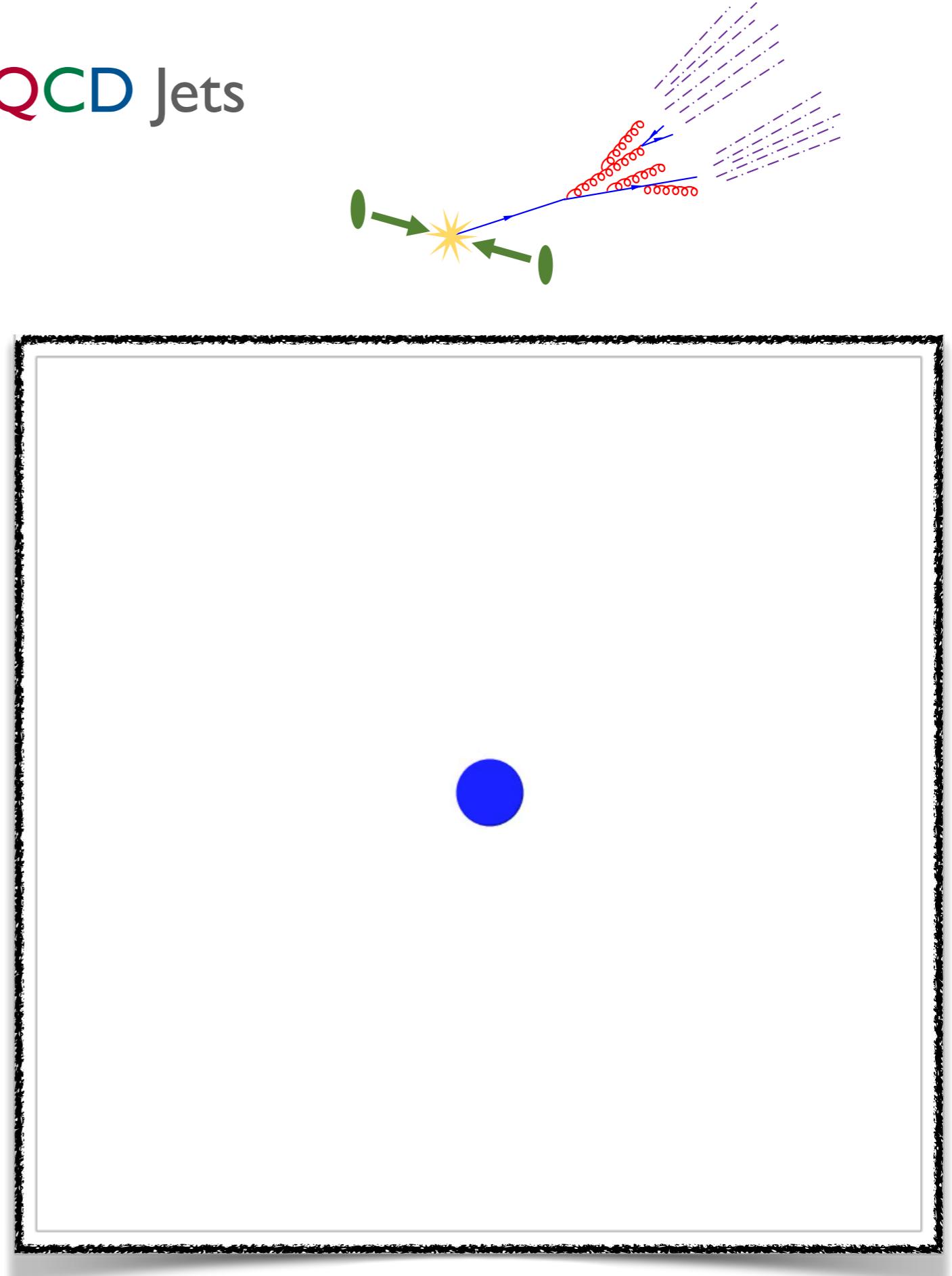
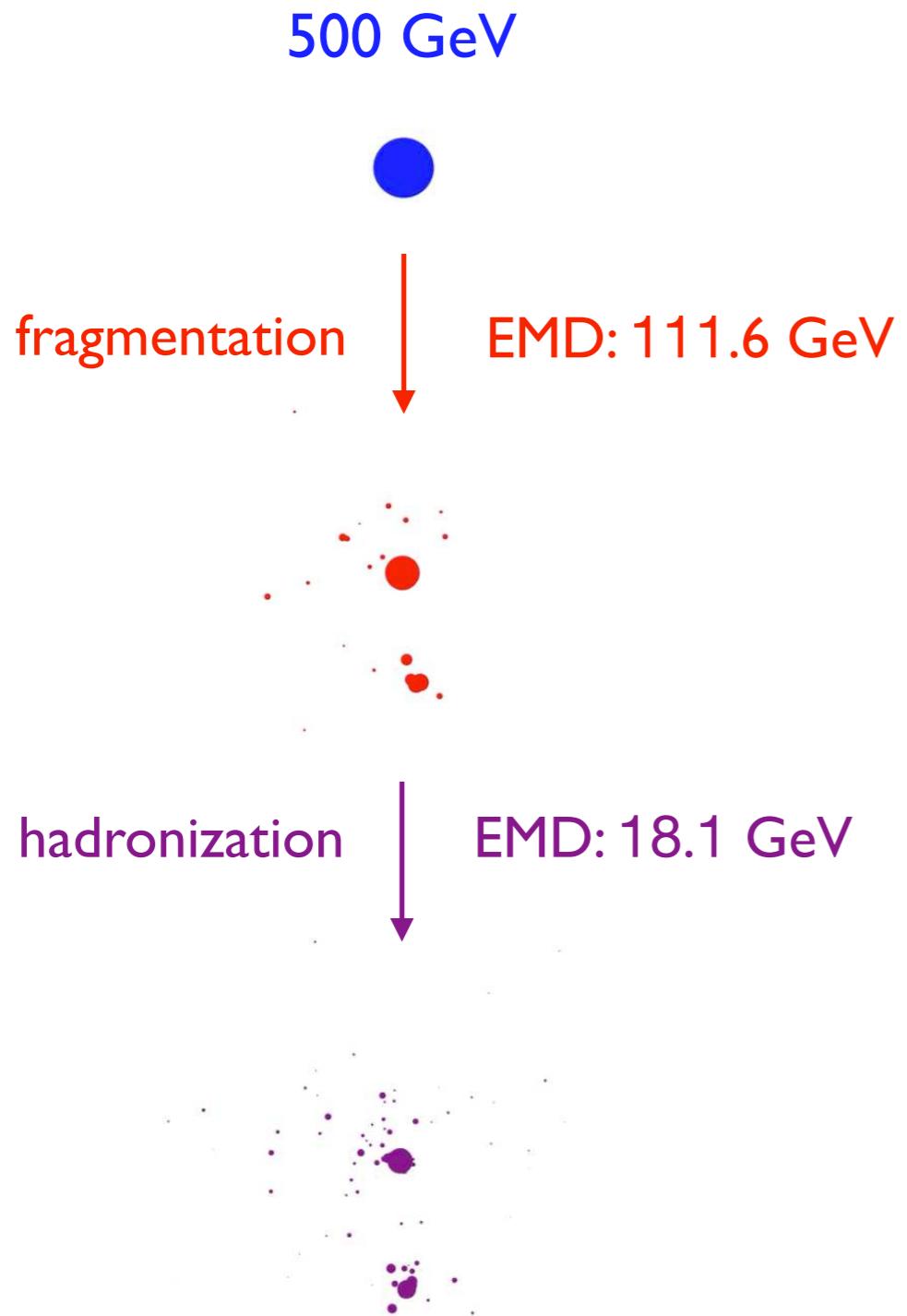
- Erdmann, Geiger, Glombitza, Schmidt, [1802.03325](#)
- Erdmann, Glombitza, Quast, [1807.01954](#)

### *Wasserstein(-Wasserstein) Autoencoders*

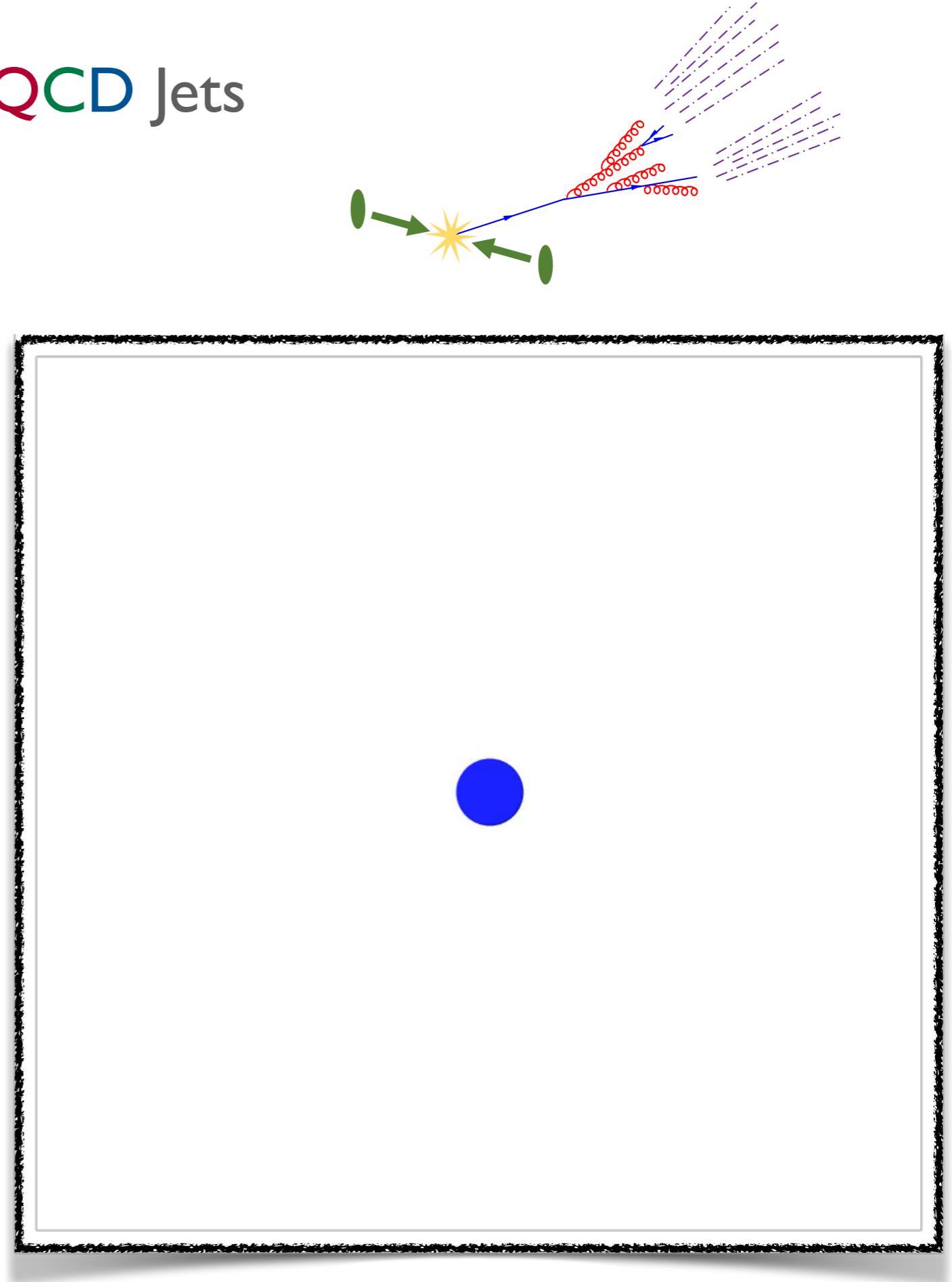
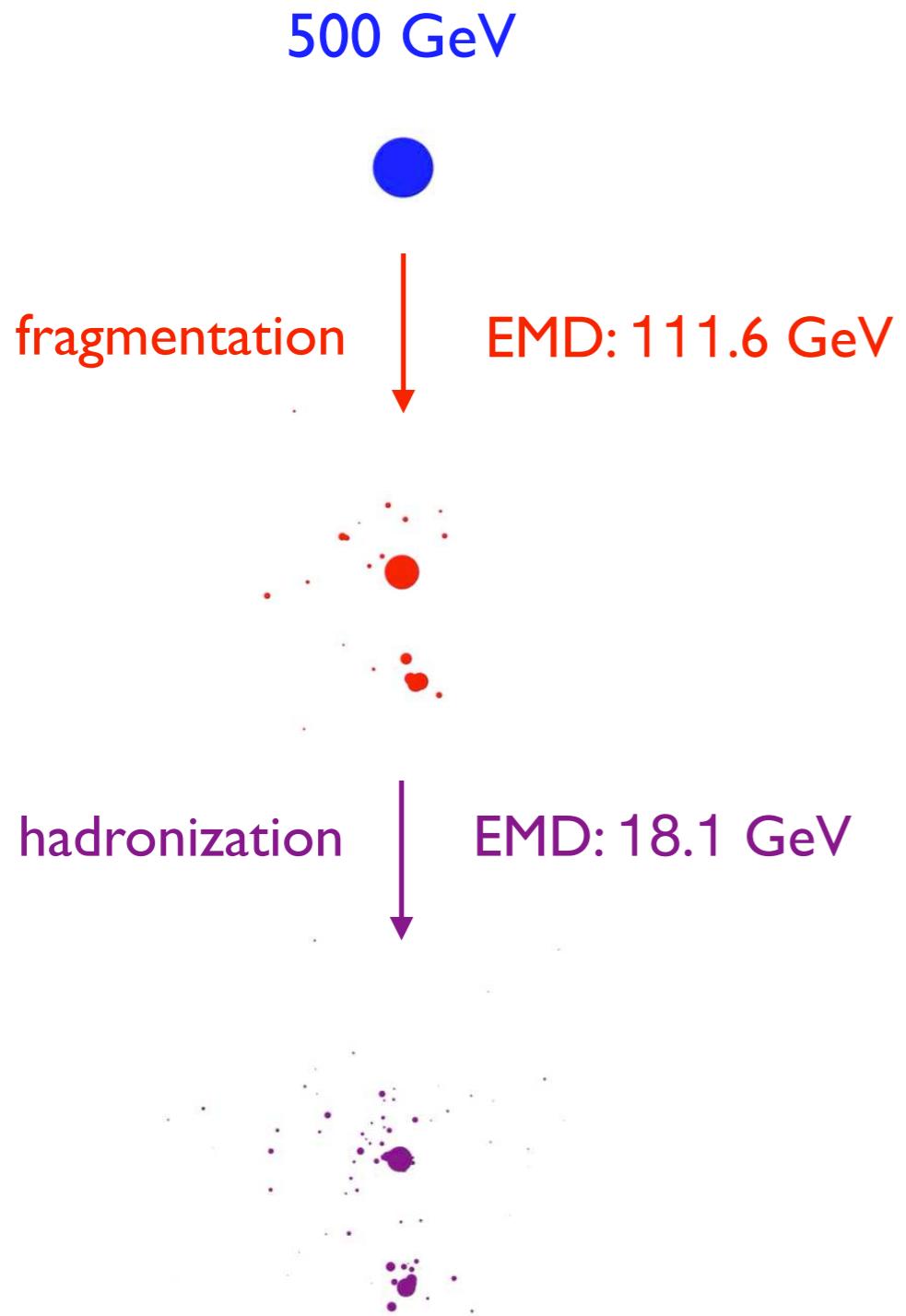
[Tolstikhin, Bousquet, Gelly, Schölkopf, [1711.01558](#)]

[Zhang, Gao, Jiao, Liu, Wang, Yang, [1902.09323](#)]

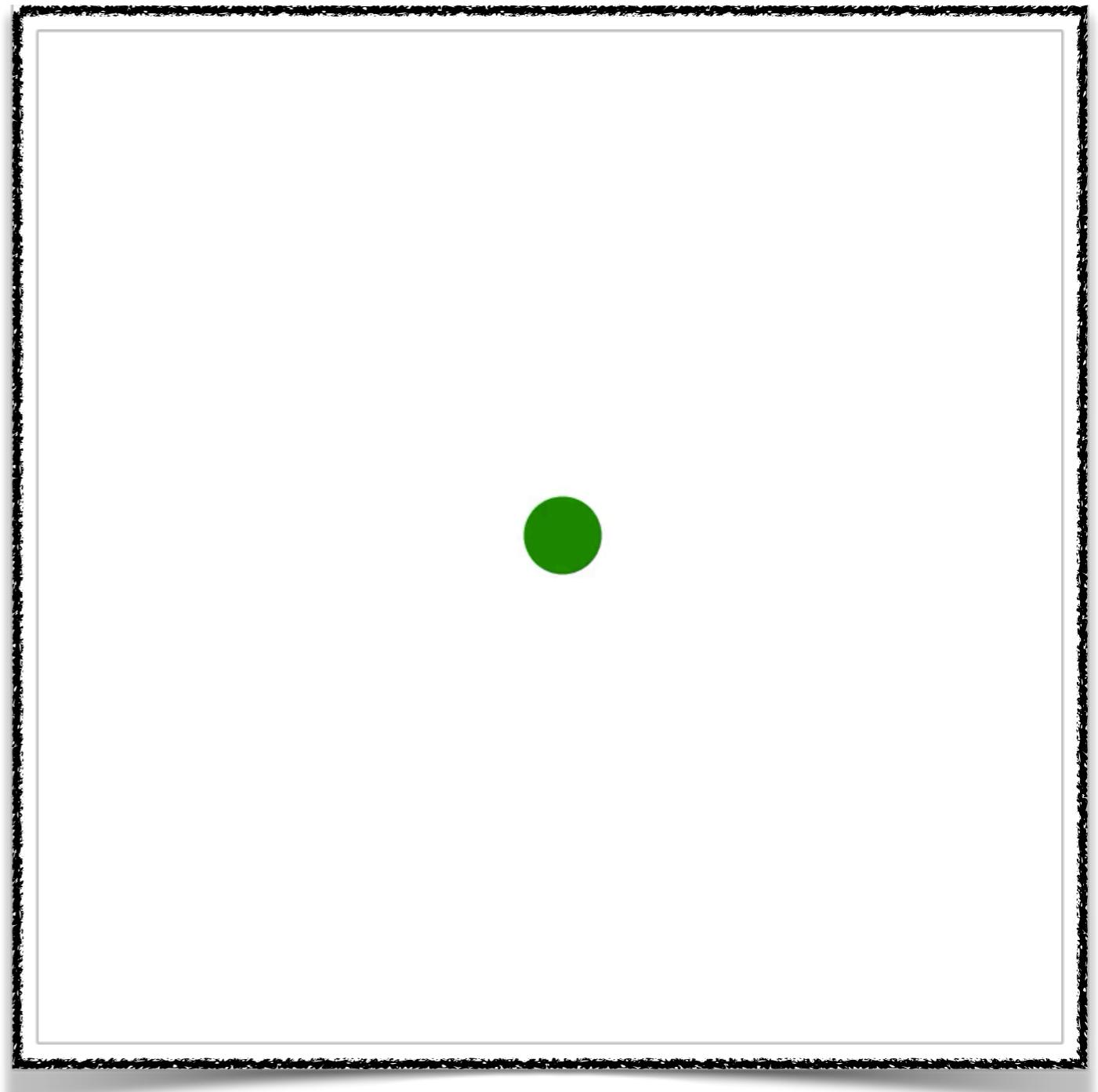
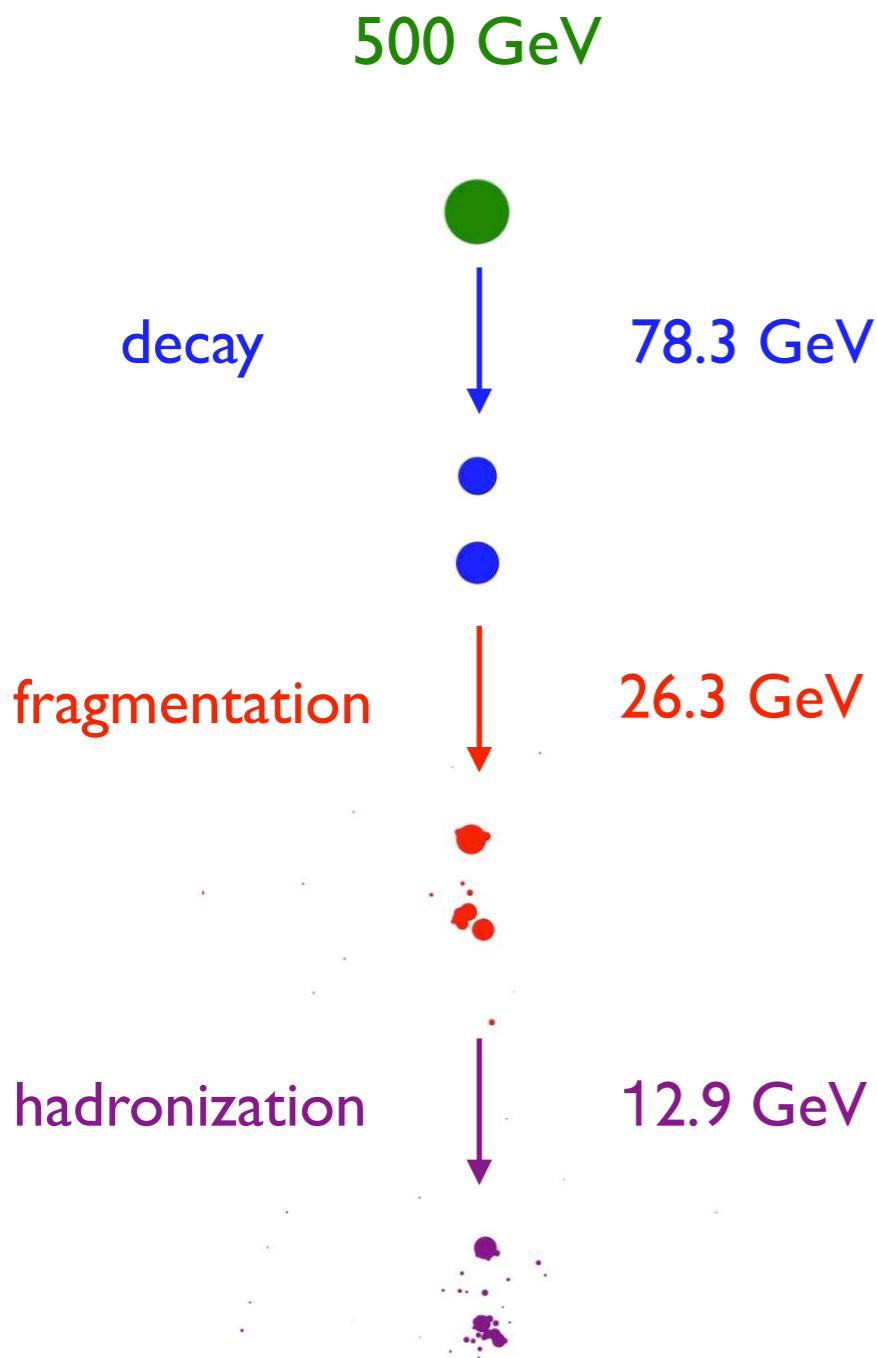
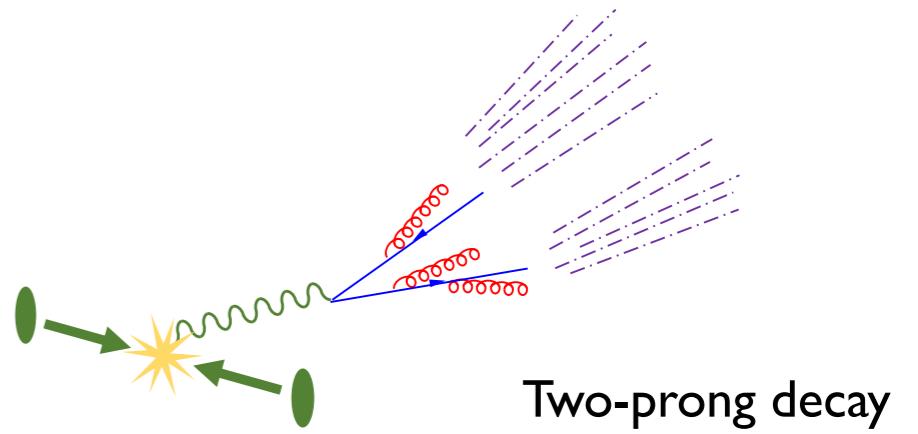
# Visualizing Jet Formation – QCD Jets



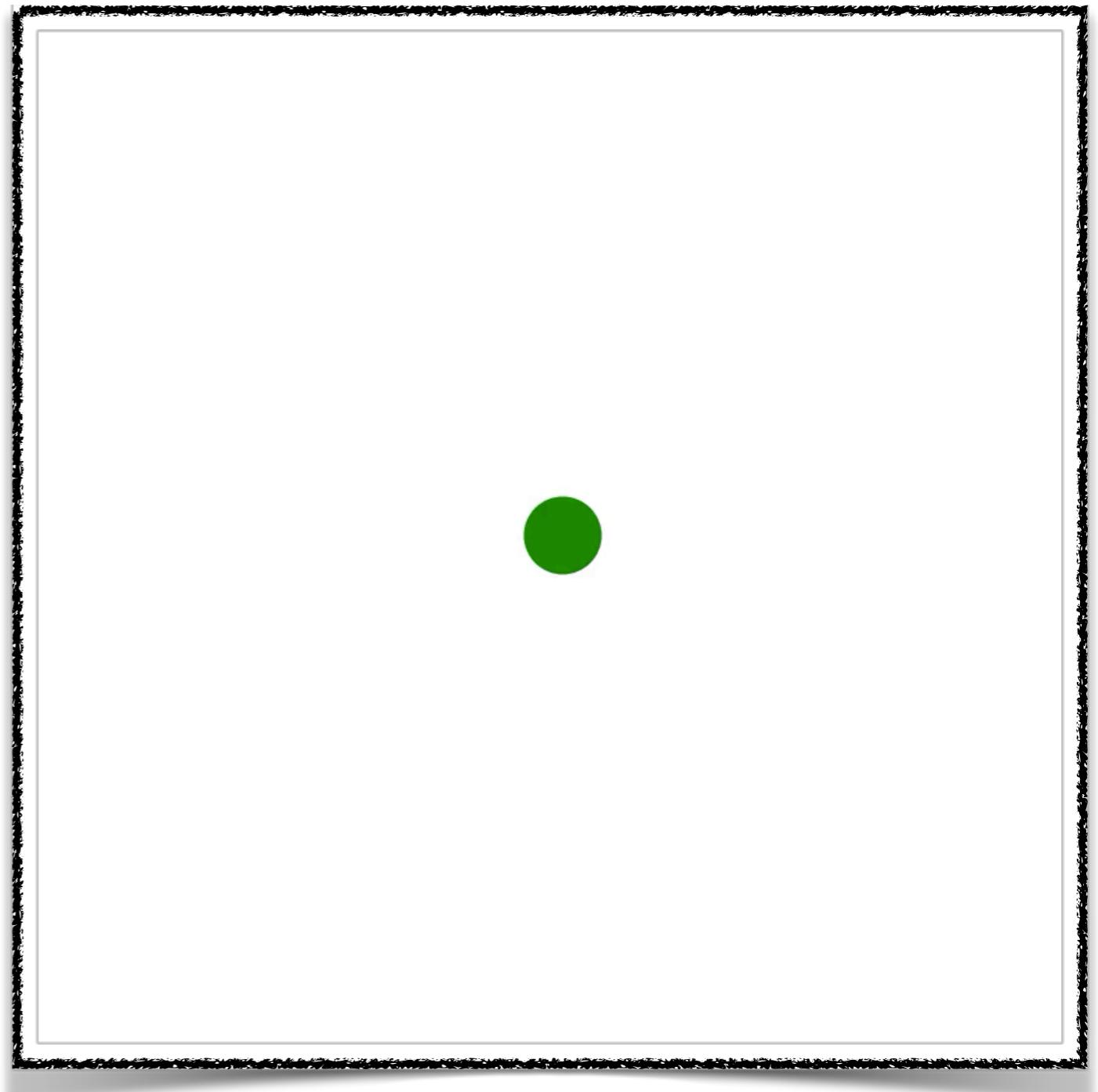
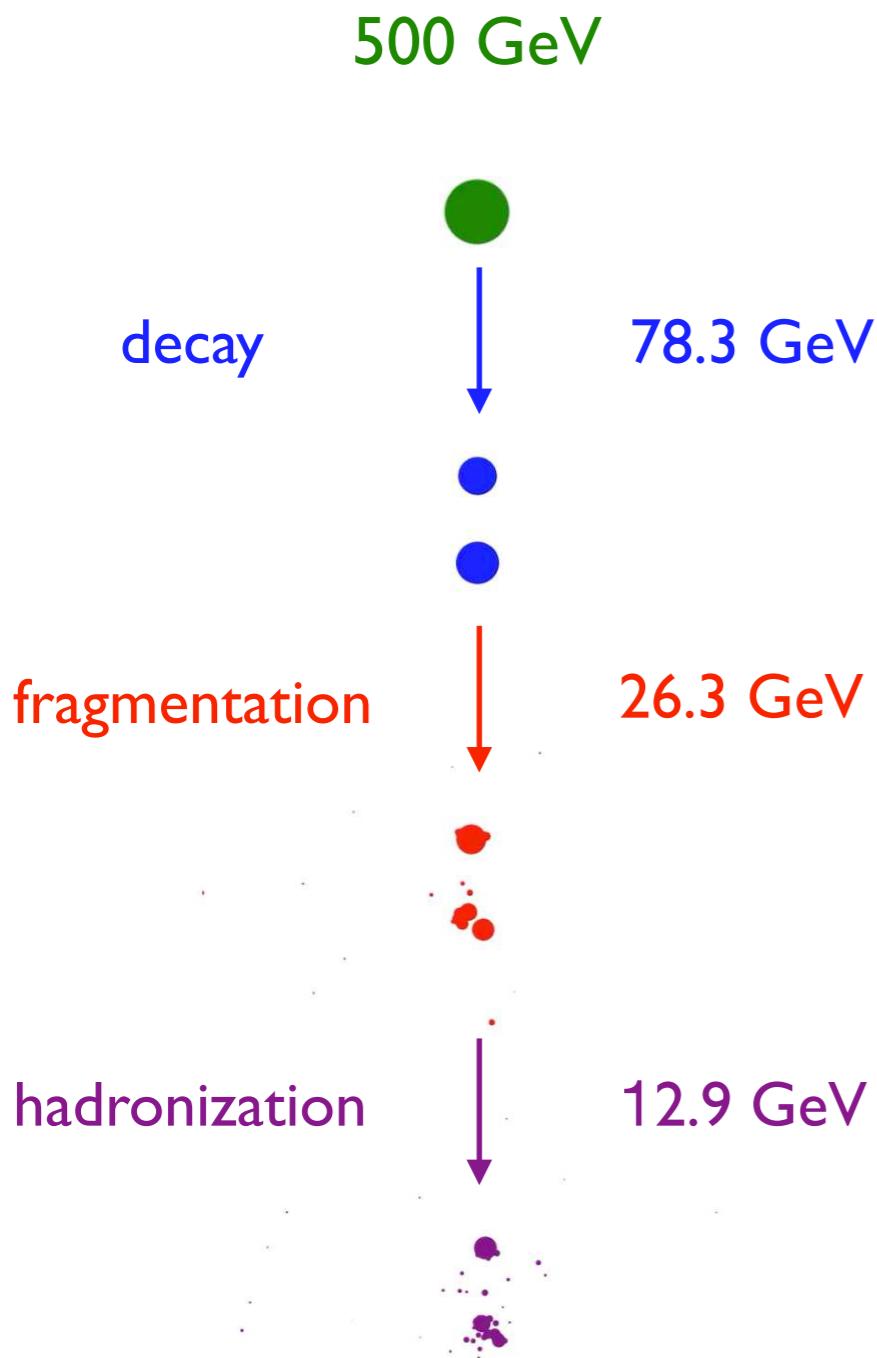
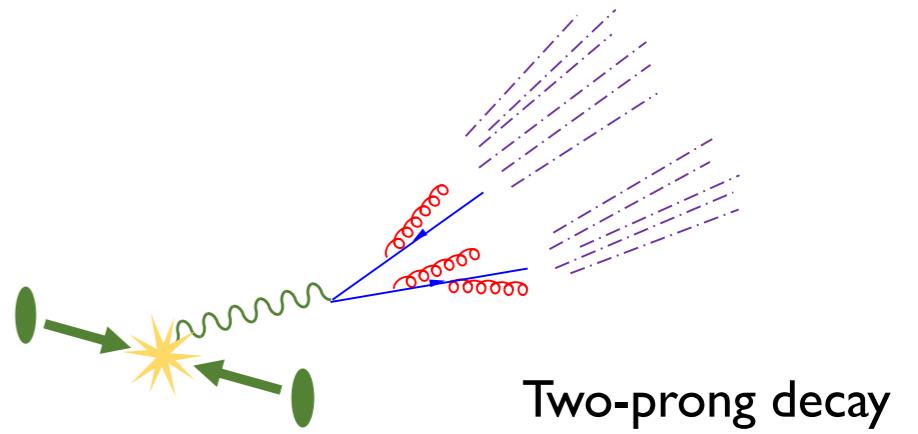
# Visualizing Jet Formation – QCD Jets



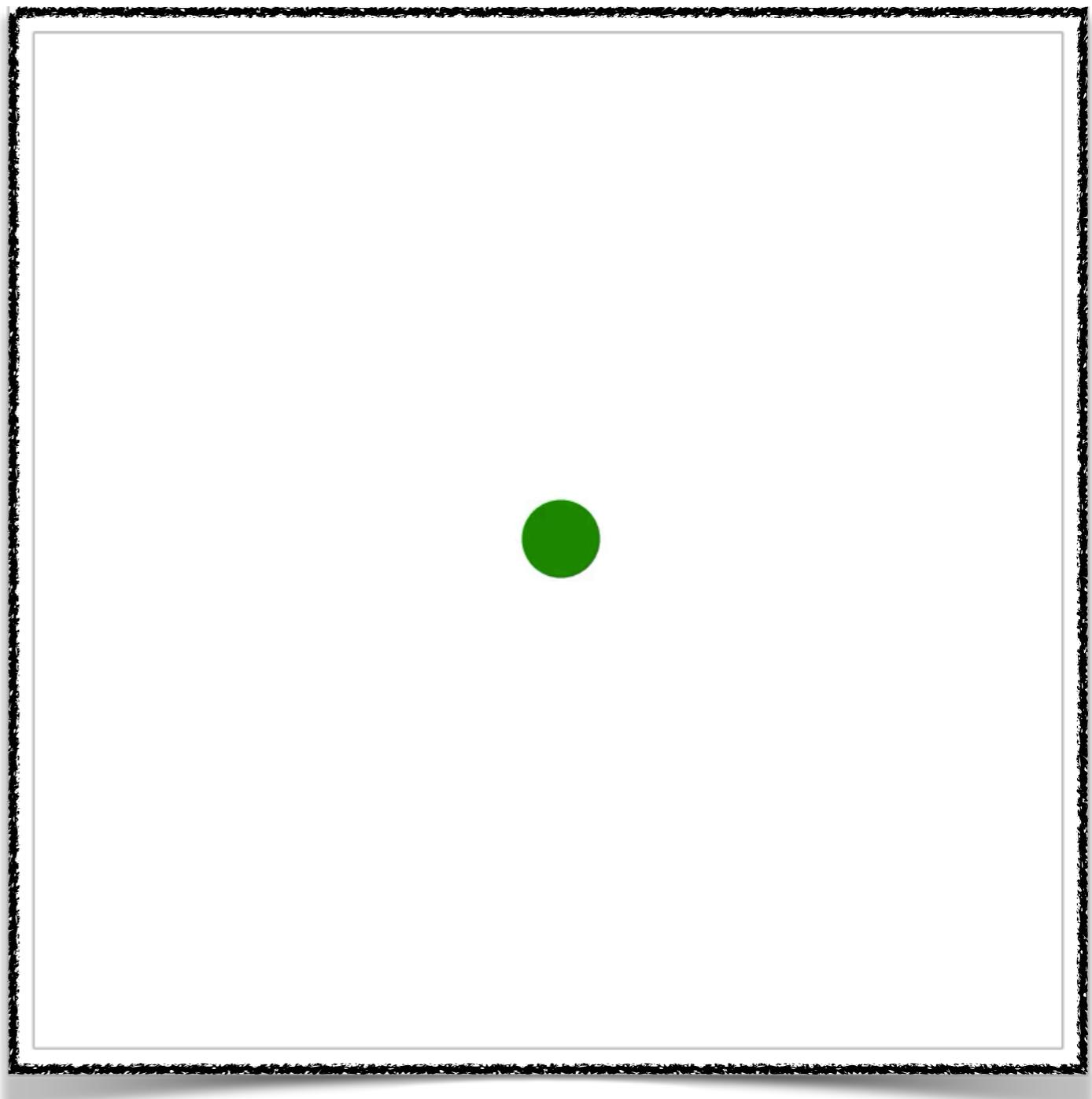
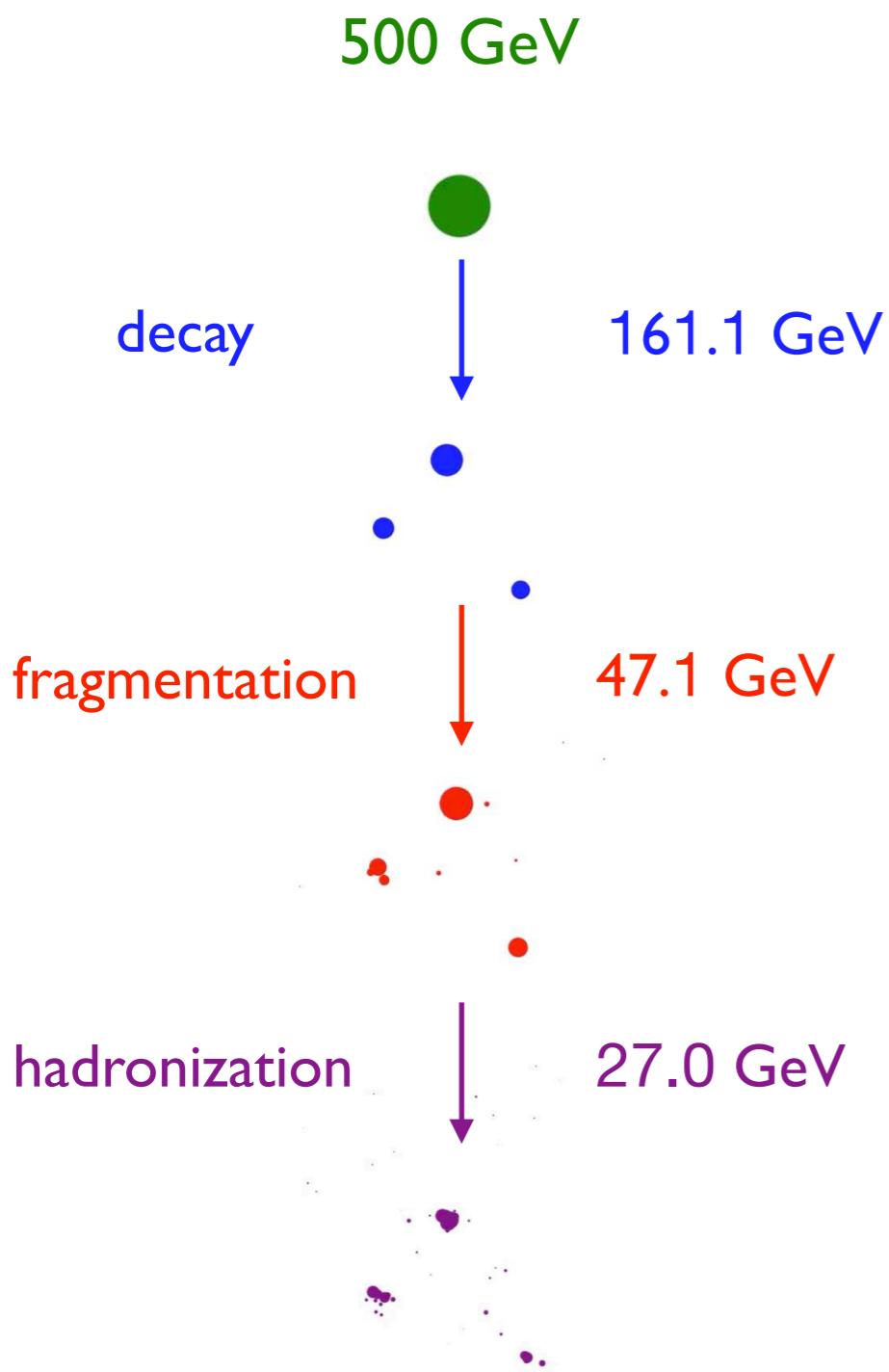
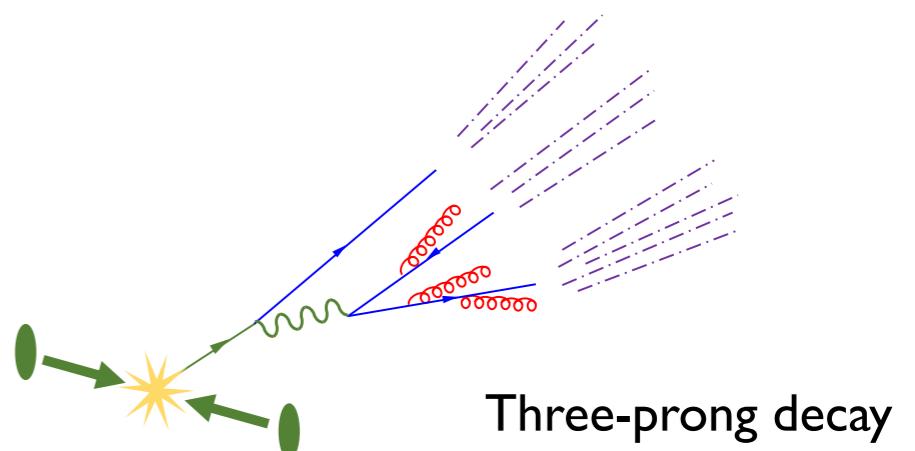
# Visualizing Jet Formation – W Jets



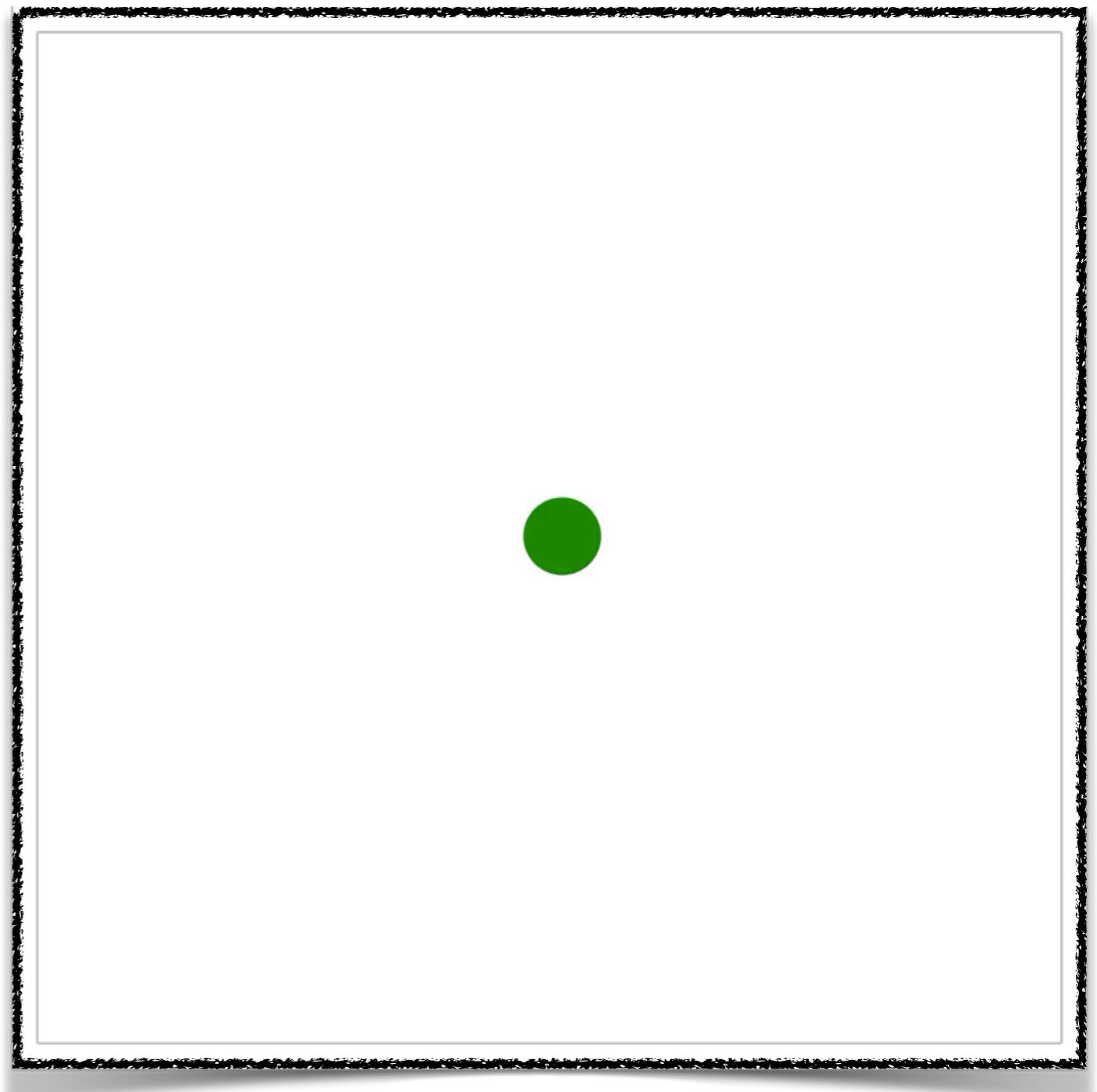
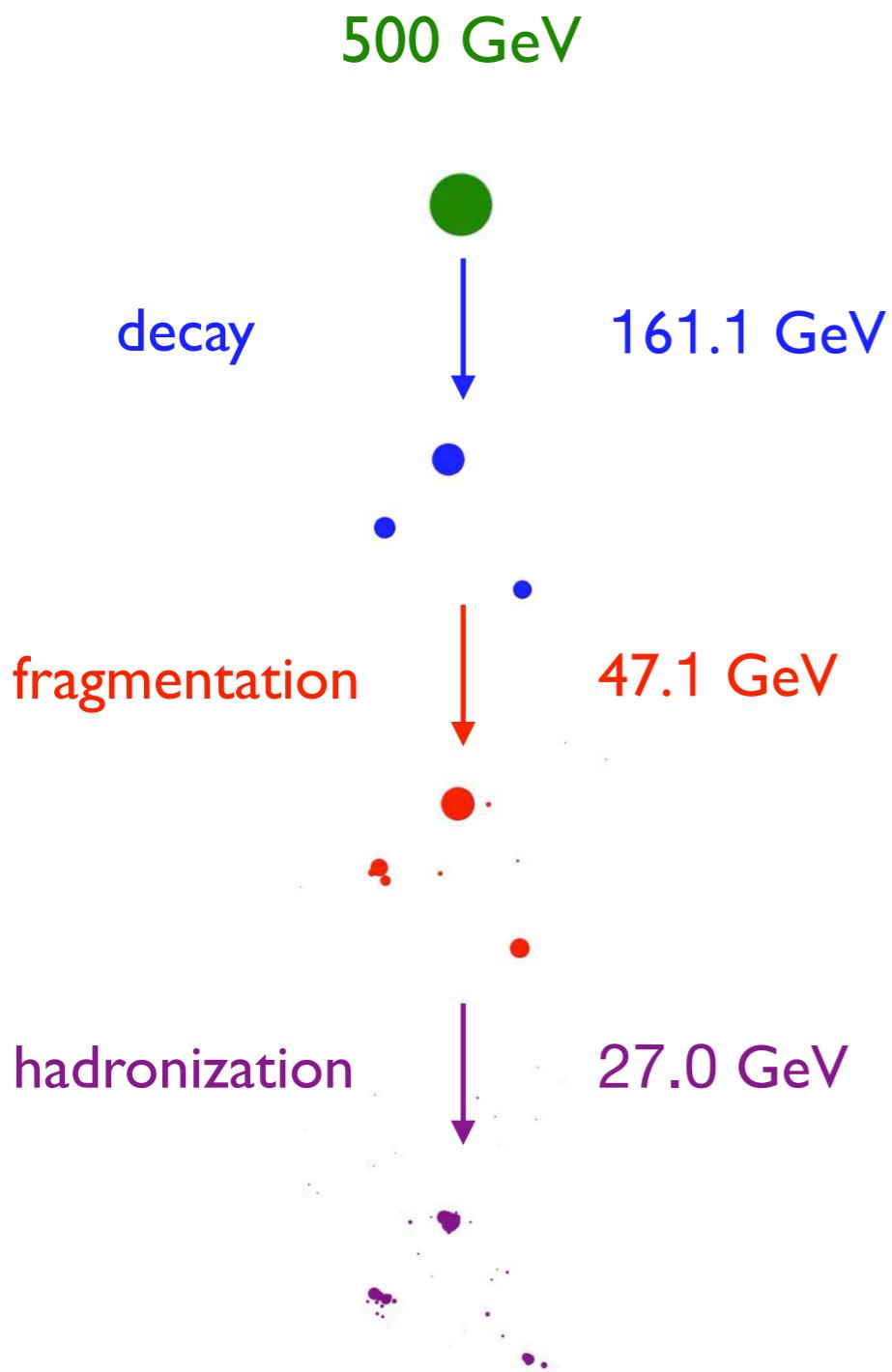
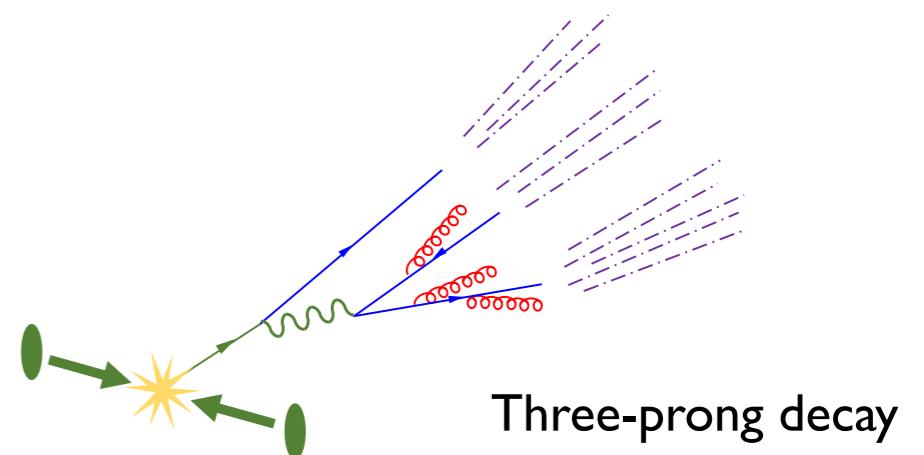
# Visualizing Jet Formation – W Jets

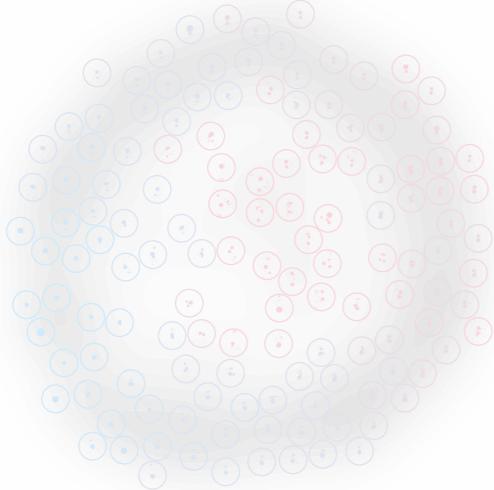
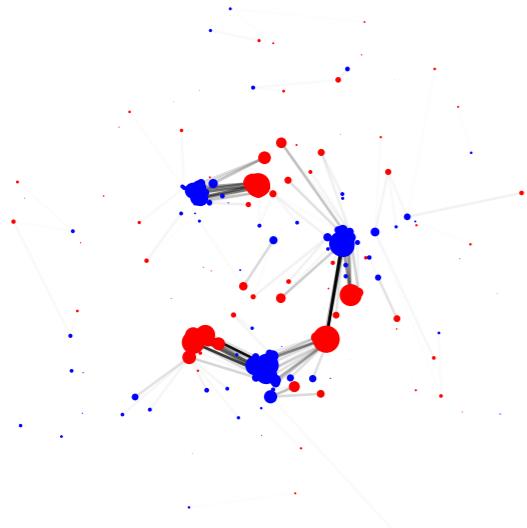
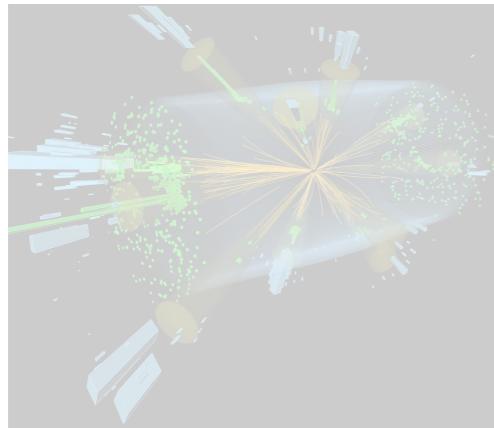


# Visualizing Jet Formation – Top Jets



# Visualizing Jet Formation – Top Jets





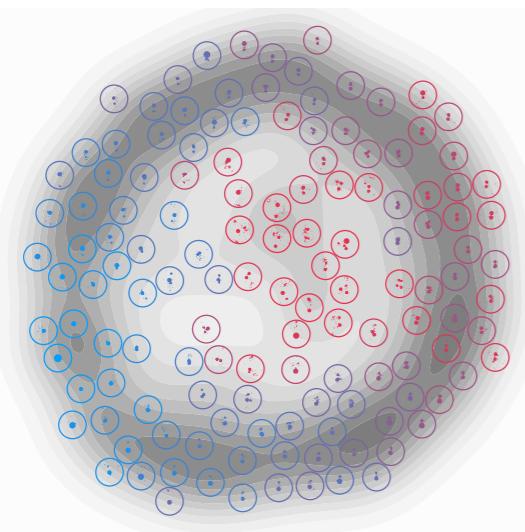
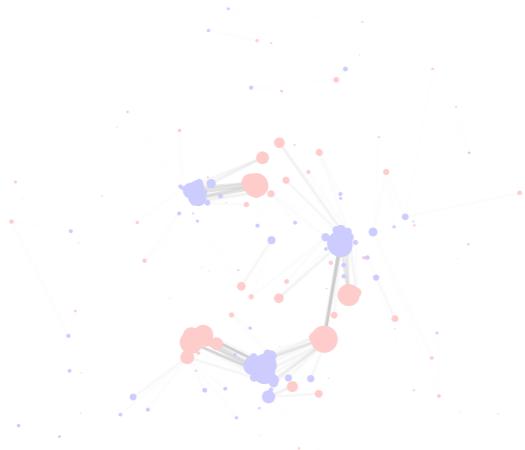
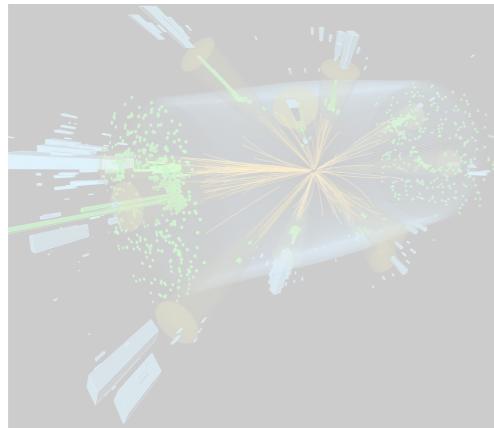
## Collider Event Foundations

*IRC-safe energy flow is theoretically and experimentally robust*

## The Energy Mover's Distance

*Quantifies the difference in energy flow between events*

## Particle Physics Applications



## Collider Event Foundations

*IRC-safe energy flow is theoretically and experimentally robust*

## The Energy Mover's Distance

*Quantifies the difference in energy flow between events*

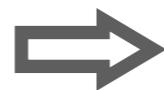
## Particle Physics Applications

# Quantifying Event Modifications

[PTK, Metodiev, Thaler, [1902.02346](#)]

## Mathematics

$\text{I-Wasserstein metric bounds the difference in expectation values between distributions}$



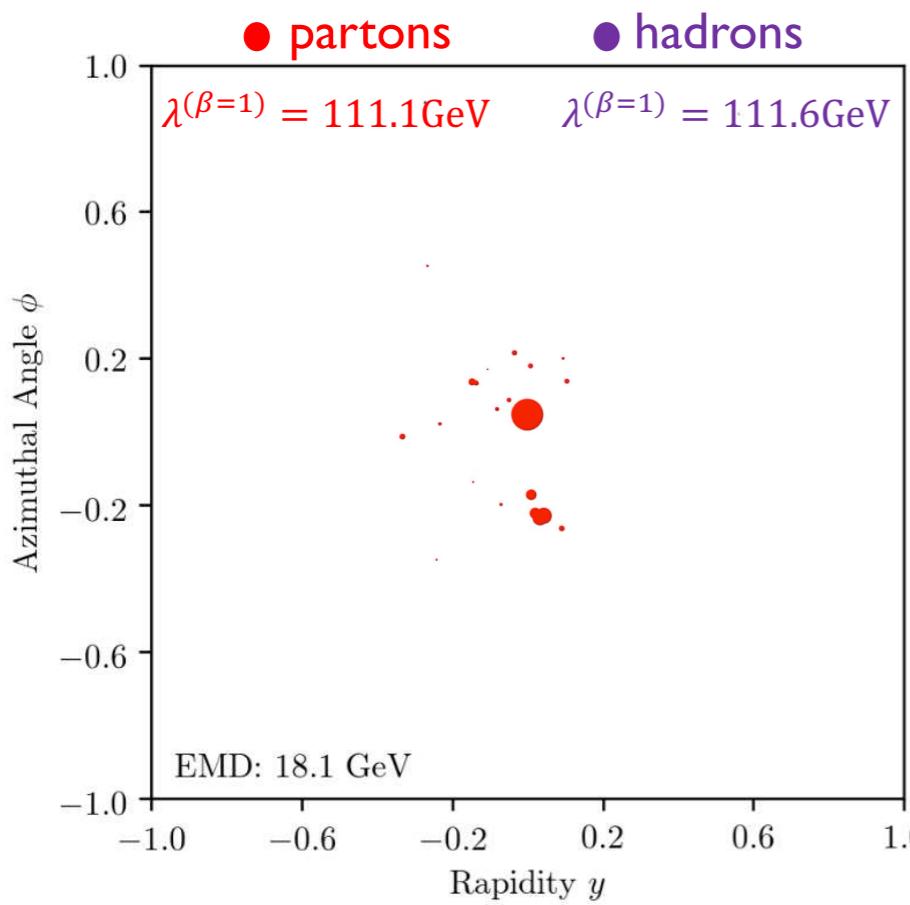
## Physics

$\text{Events close in EMD are close according to } \text{IRC-safe observables}$

$$\text{EMD}(\mathcal{E}, \mathcal{E}') \geq \frac{1}{RL} \left| \sum_i E_i \Phi(\hat{p}_i) - \sum_j E'_j \Phi(\hat{p}'_j) \right| = \frac{1}{RL} |\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')|$$

via Kantorovich-Rubinstein duality

Additive **IRC**-safe observable

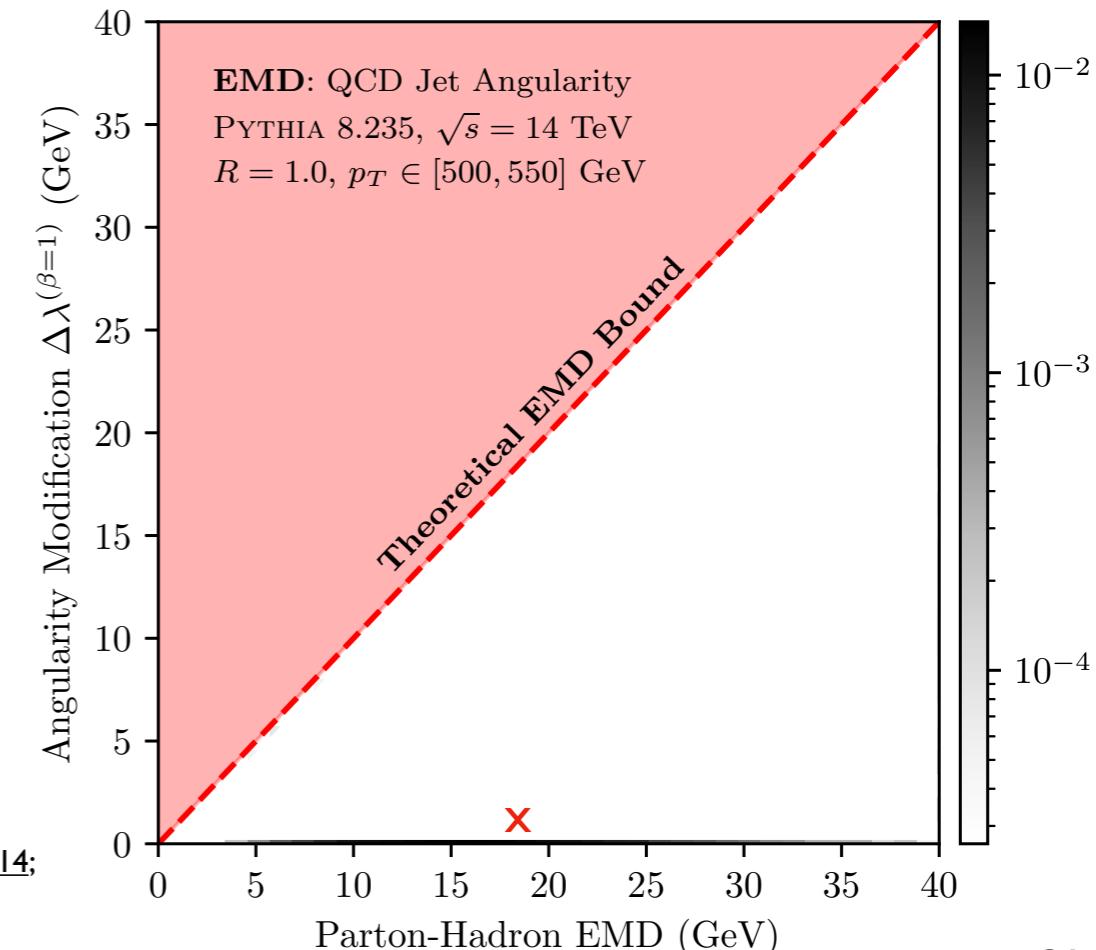


e.g. bounding **IRC**-safe angularities

$$\lambda^{(\beta)}(\mathcal{E}) = \sum_i E_i \theta_i^\beta$$

Can do the same for pileup, detector effects

[Berger, Kucs, Sterman, [hep-ph/0303051](#);  
Ellis, Vermilion, Walsh, Hornig, Lee, [1001.0014](#);  
Larkoski, Thaler, Waalewijn, [1408.3122](#)]

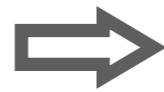


# Quantifying Event Modifications

[PTK, Metodiev, Thaler, [1902.02346](#)]

## Mathematics

$\text{I-Wasserstein metric bounds the difference in expectation values between distributions}$



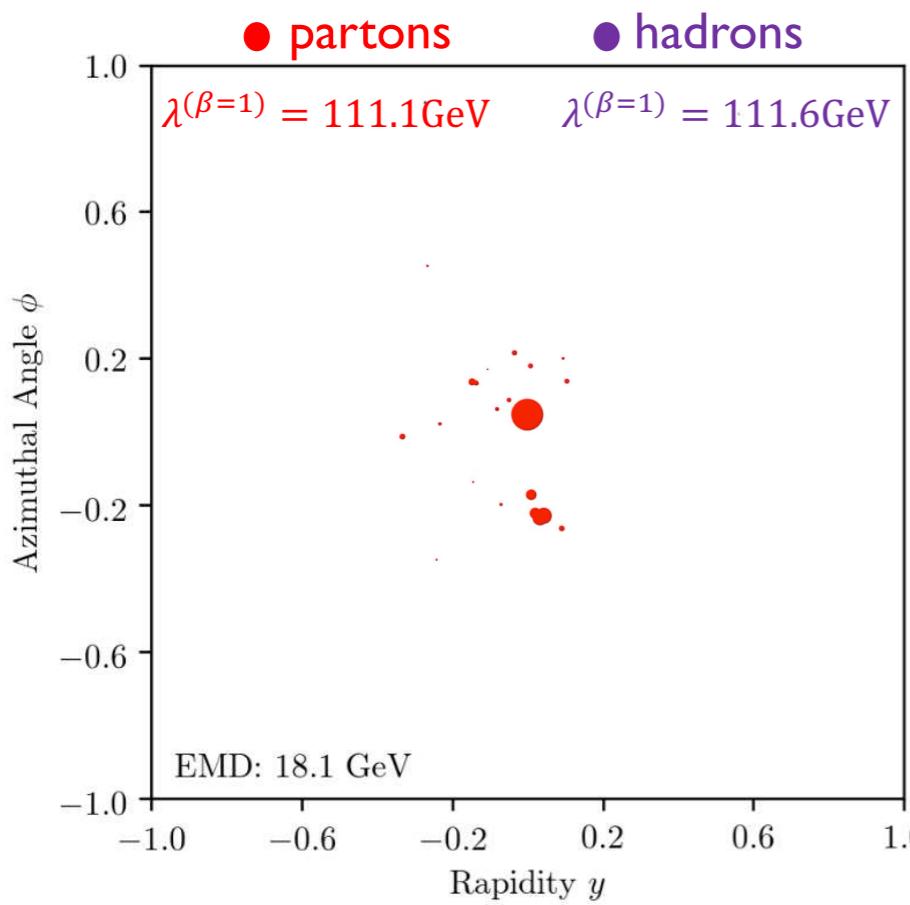
## Physics

$\text{Events close in EMD are close according to } \text{IRC-safe observables}$

$$\text{EMD}(\mathcal{E}, \mathcal{E}') \geq \frac{1}{RL} \left| \sum_i E_i \Phi(\hat{p}_i) - \sum_j E'_j \Phi(\hat{p}'_j) \right| = \frac{1}{RL} |\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')|$$

via Kantorovich-Rubinstein duality

Additive **IRC**-safe observable

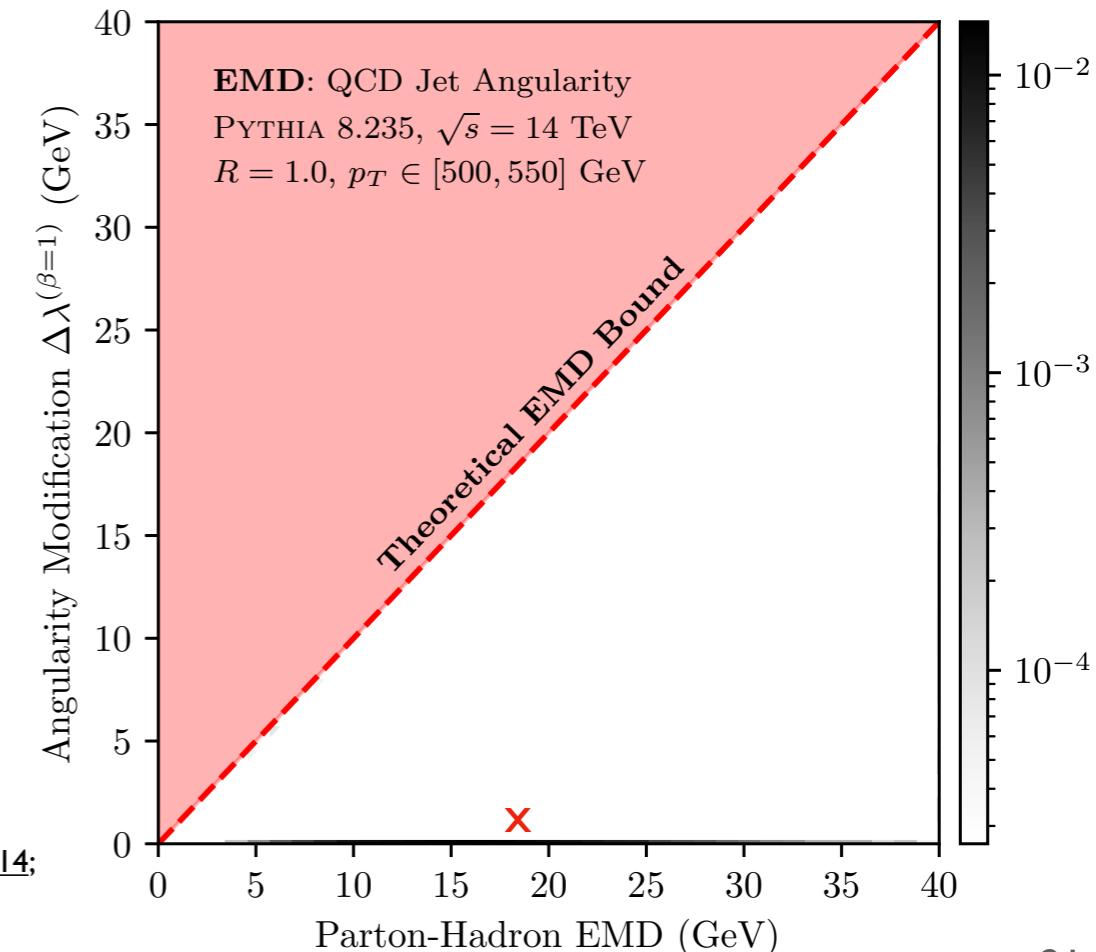


e.g. bounding **IRC**-safe angularities

$$\lambda^{(\beta)}(\mathcal{E}) = \sum_i E_i \theta_i^\beta$$

Can do the same for pileup, detector effects

[Berger, Kucs, Sterman, [hep-ph/0303051](#);  
Ellis, Vermilion, Walsh, Hornig, Lee, [1001.0014](#);  
Larkoski, Thaler, Waalewijn, [1408.3122](#)]

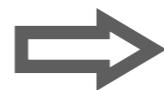


# Quantifying Event Modifications

[PTK, Metodiev, Thaler, [1902.02346](#)]

## Mathematics

$\text{I-Wasserstein metric bounds the difference in expectation values between distributions}$



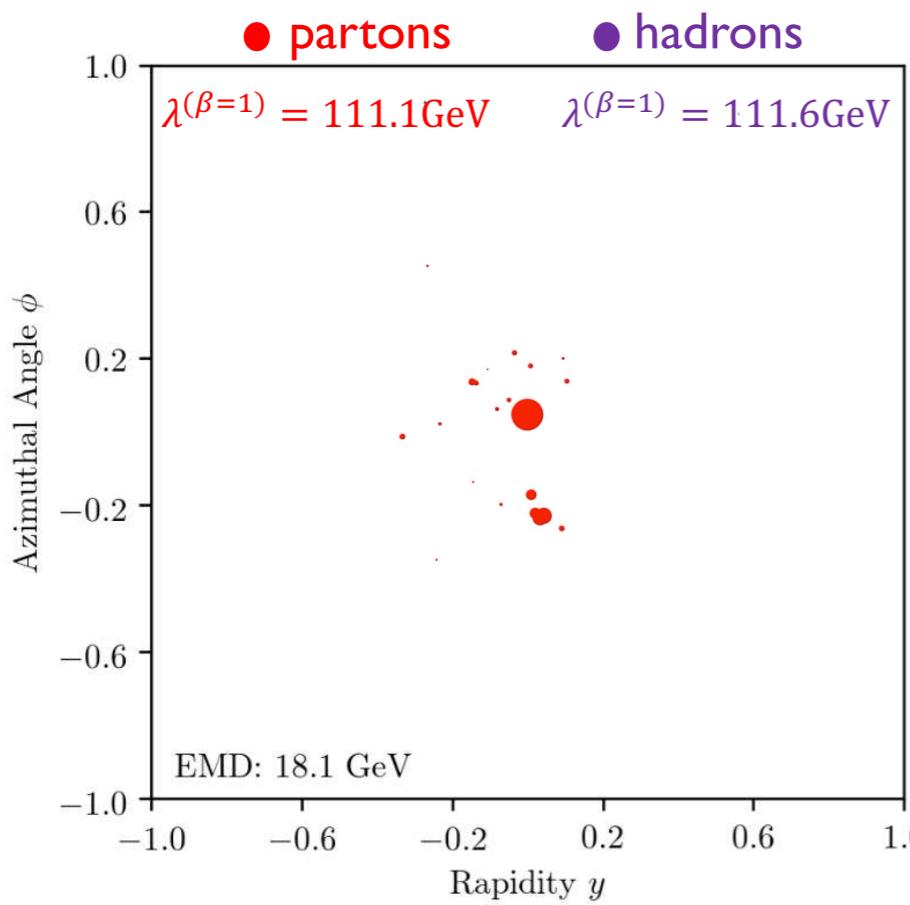
## Physics

$\text{Events close in EMD are close according to } \text{IRC-safe observables}$

$$\text{EMD}(\mathcal{E}, \mathcal{E}') \geq \frac{1}{RL} \left| \sum_i E_i \Phi(\hat{p}_i) - \sum_j E'_j \Phi(\hat{p}'_j) \right| = \frac{1}{RL} |\mathcal{O}(\mathcal{E}) - \mathcal{O}(\mathcal{E}')|$$

via Kantorovich-Rubinstein duality

Additive **IRC**-safe observable

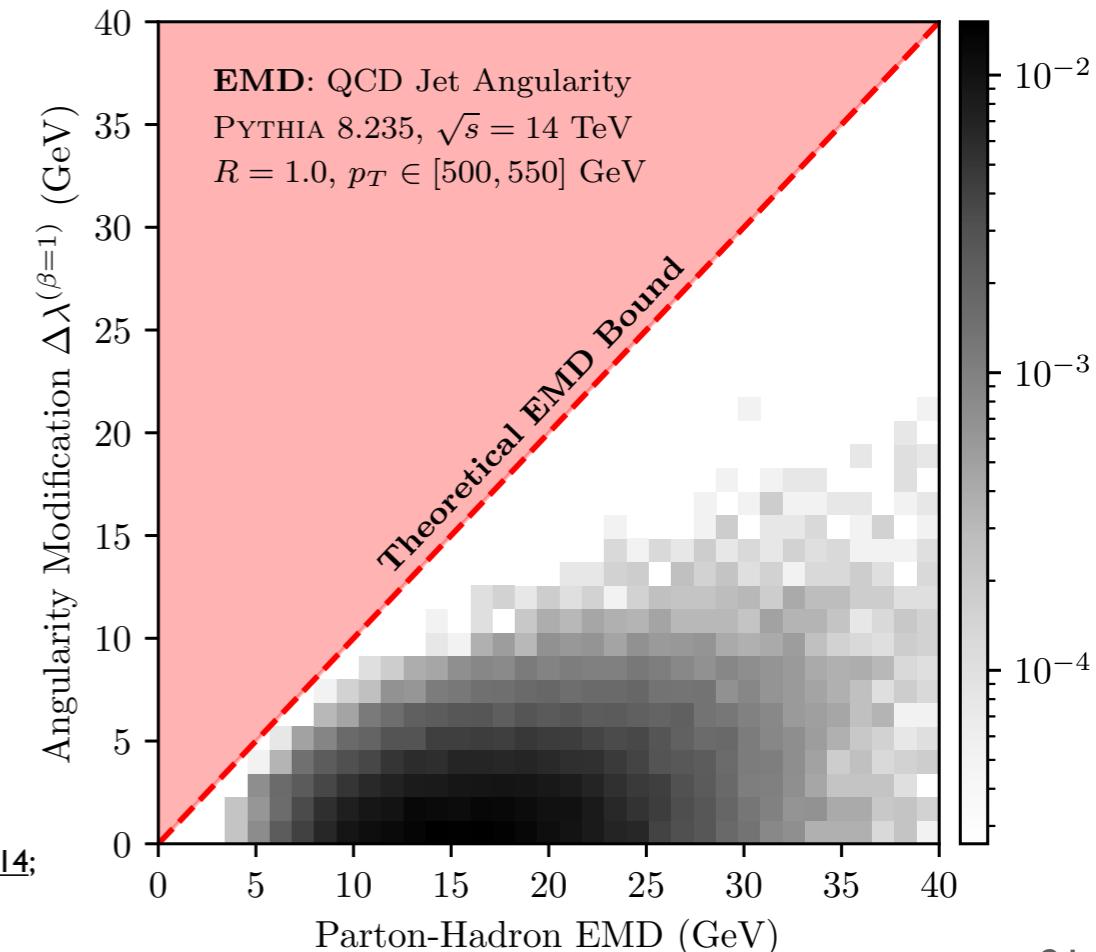


e.g. bounding **IRC**-safe angularities

$$\lambda^{(\beta)}(\mathcal{E}) = \sum_i E_i \theta_i^\beta$$

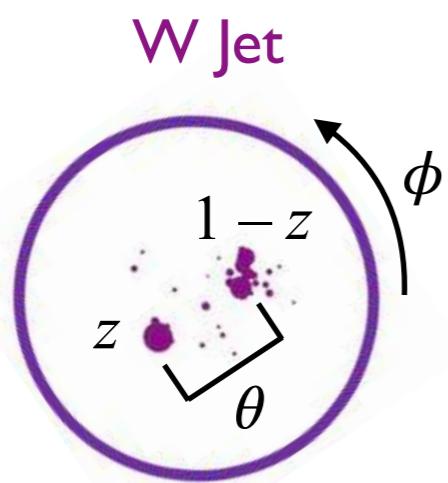
Can do the same for pileup, detector effects

[Berger, Kucs, Sterman, [hep-ph/0303051](#);  
Ellis, Vermilion, Walsh, Hornig, Lee, [1001.0014](#);  
Larkoski, Thaler, Waalewijn, [1408.3122](#)]

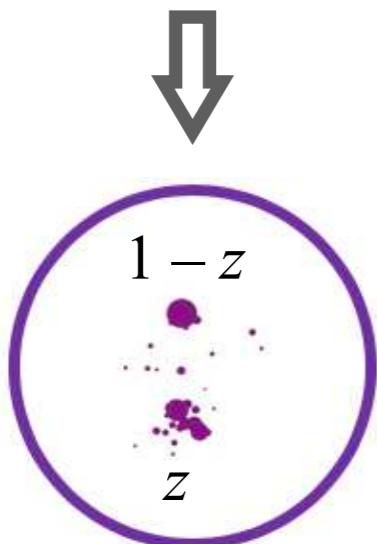


# Visualizing the Metric Space of $W$ Jets

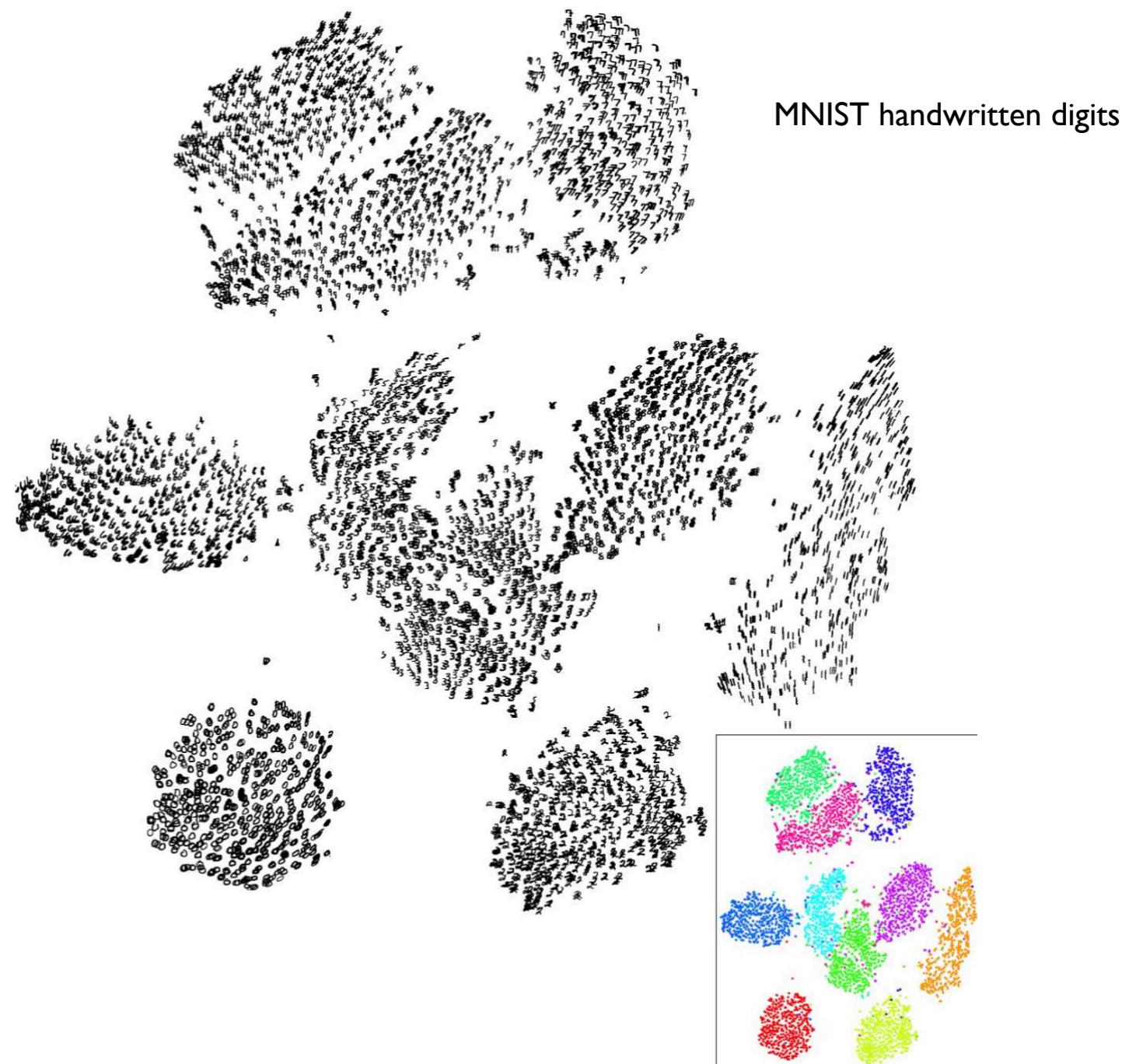
*Embed high-dimension manifold in low-dimensional space?*



Constraints:  $W$  Mass and  $\phi = 0$  preprocessing



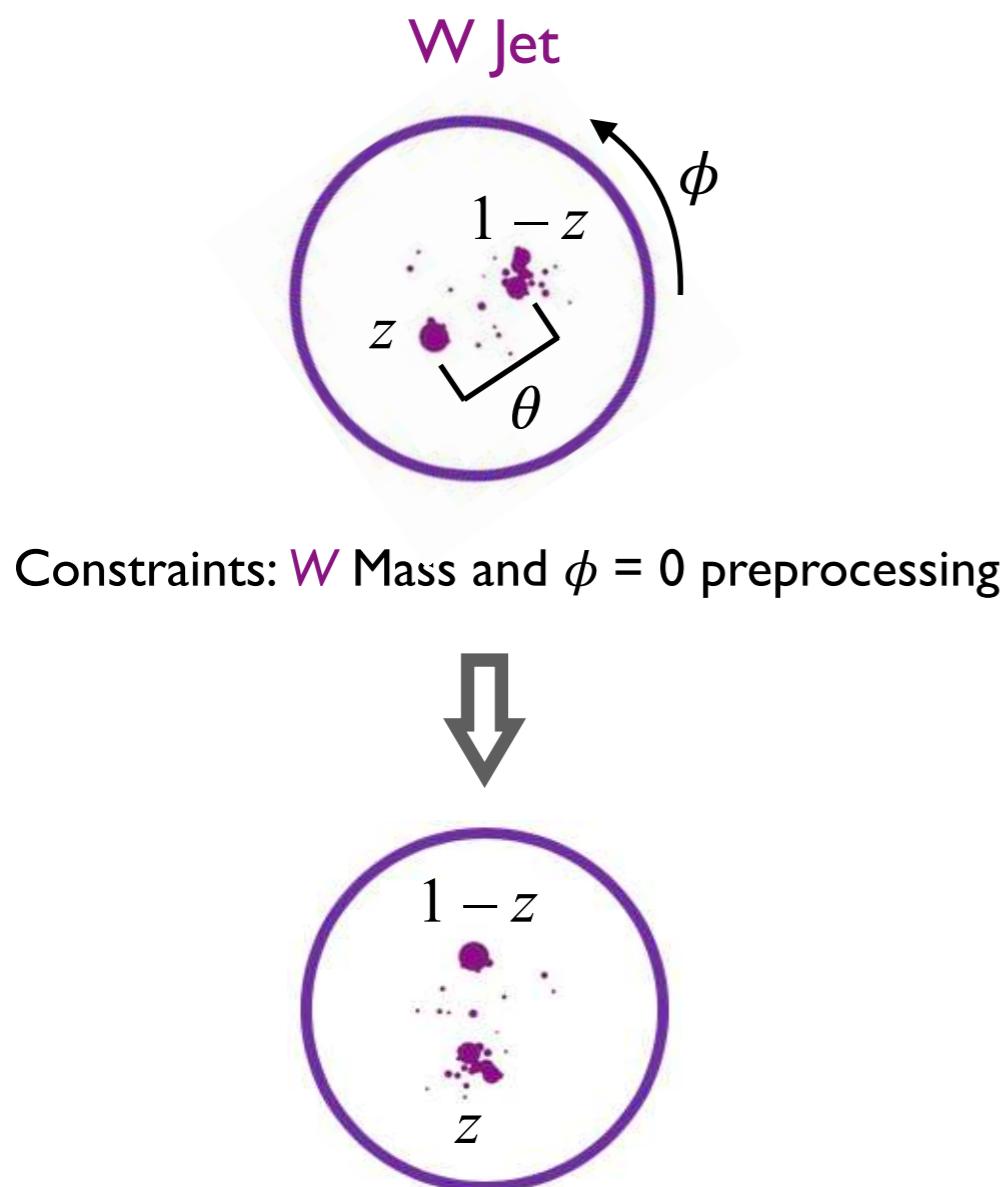
t-Distributed Stochastic Neighbor Embedding



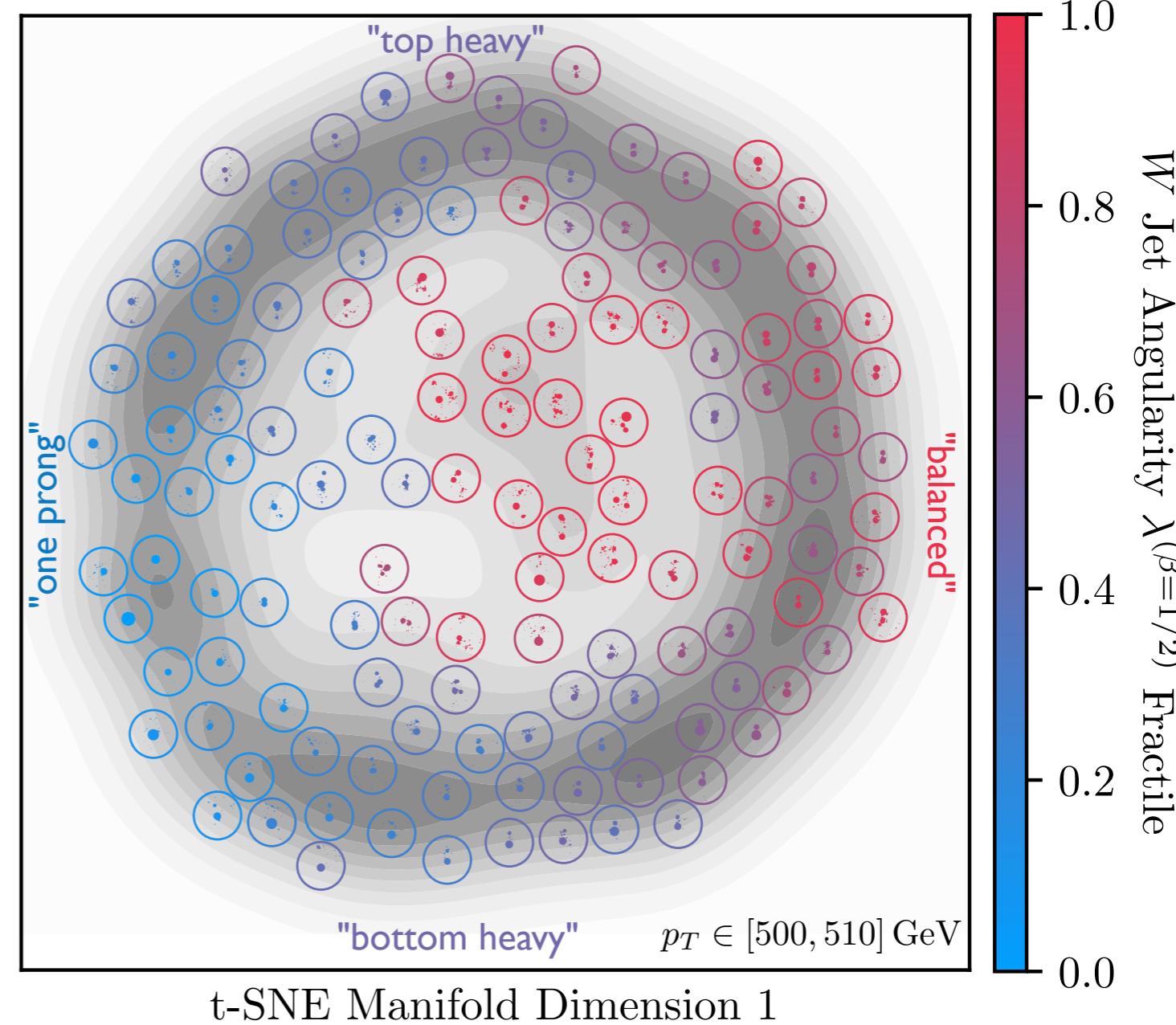
[L. van der Maaten, G. Hinton, JMLR 2008 ]

# Visualizing the Metric Space of W Jets

*Embed high-dimension manifold in low-dimensional space?*

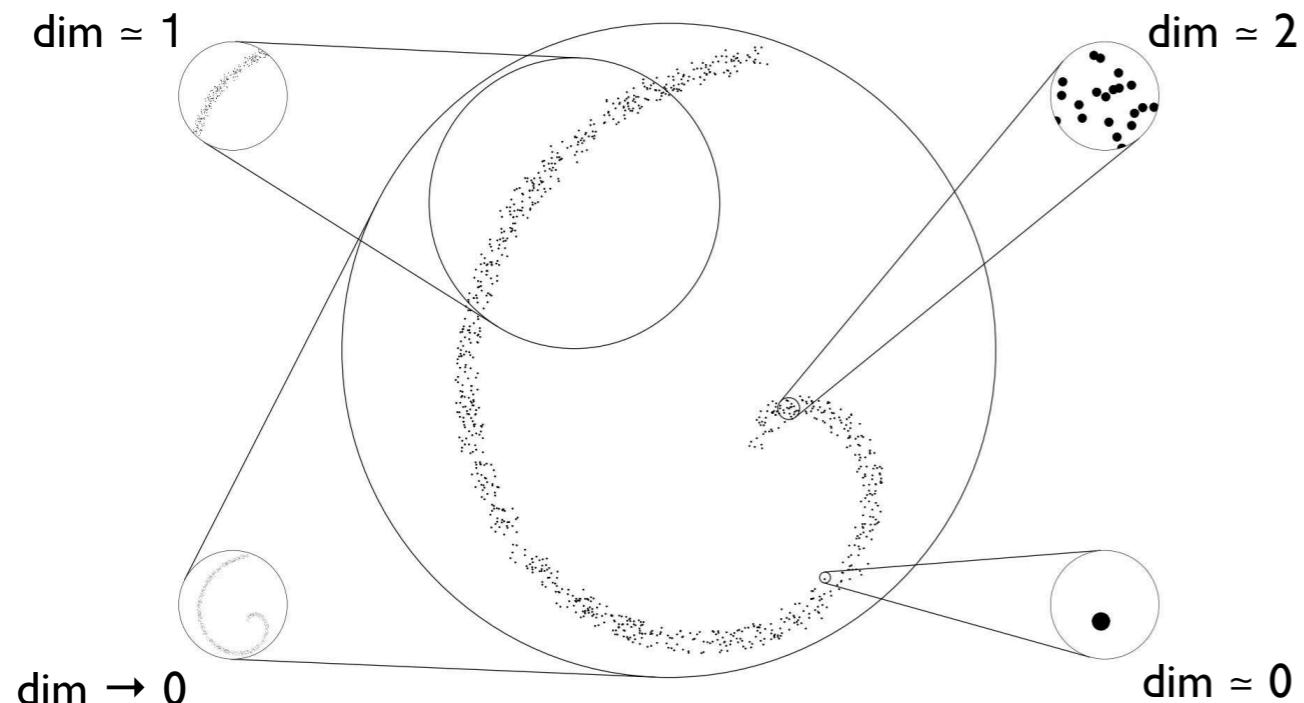


t-Distributed Stochastic Neighbor Embedding



# Manifold Dimensions of Event Space

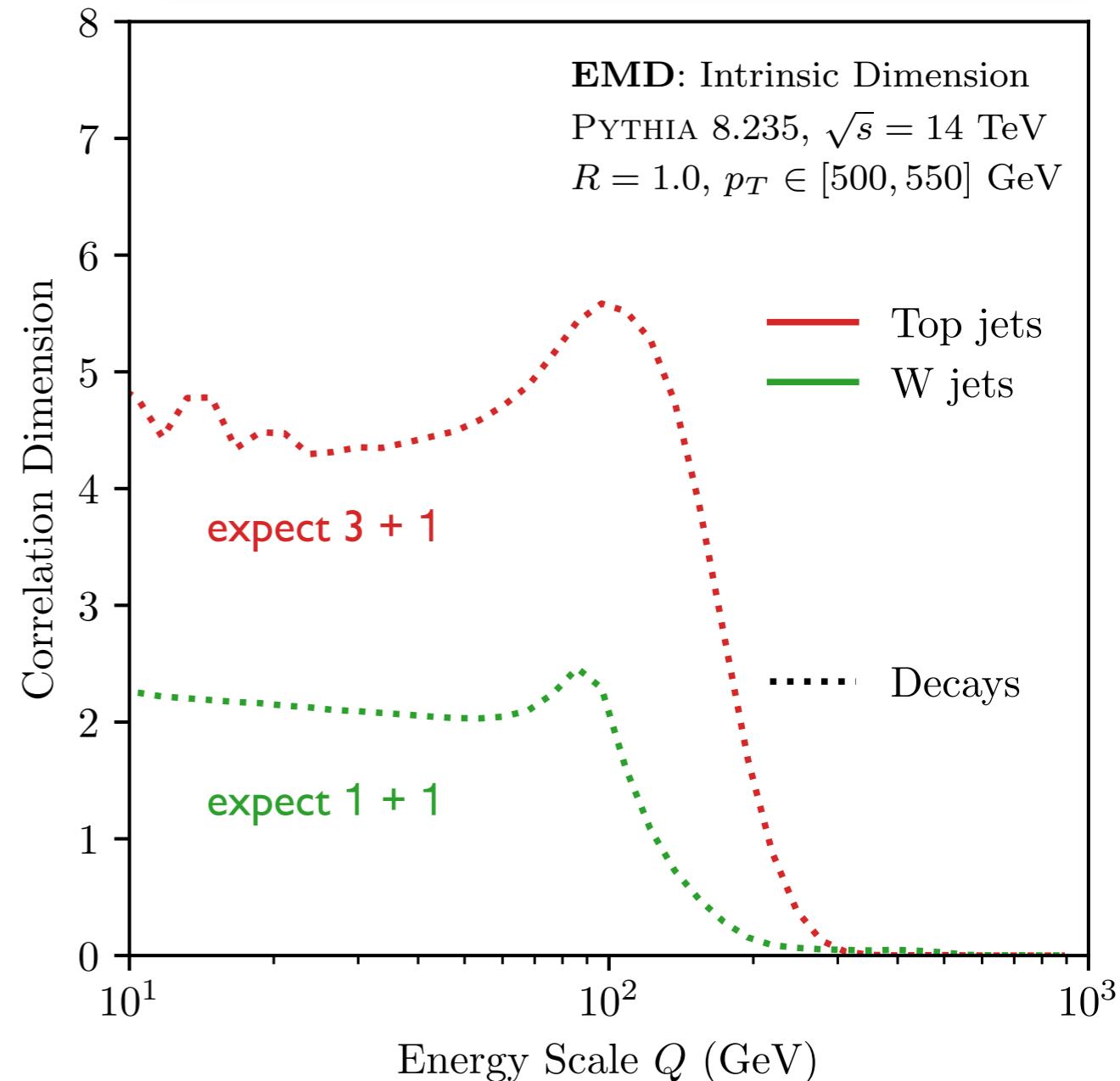
**Correlation dimension: how does the # of elements within a ball of size  $Q$  change?**



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

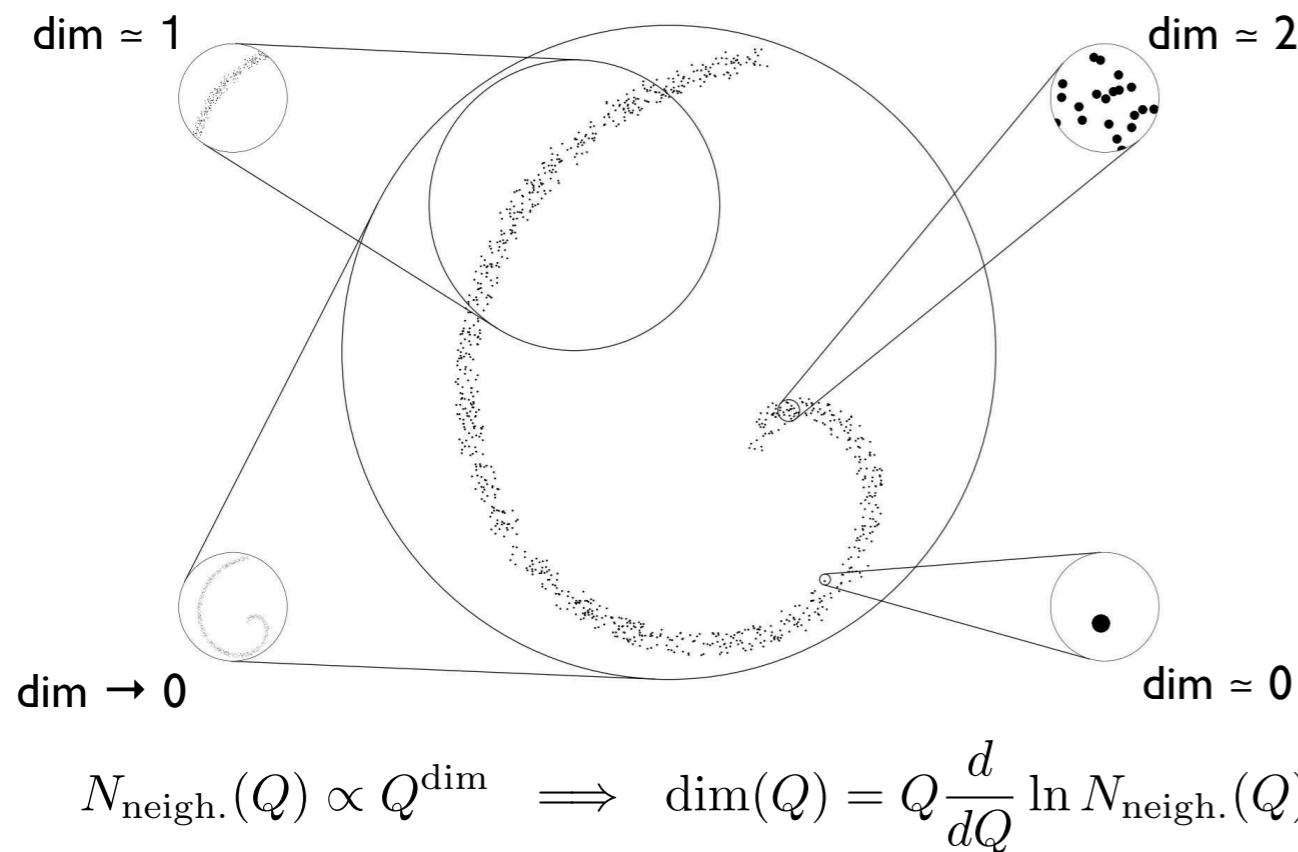
**Correlation dimension lessons:**  
Decays are "constant" dim. at low  $Q$

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



# Manifold Dimensions of Event Space

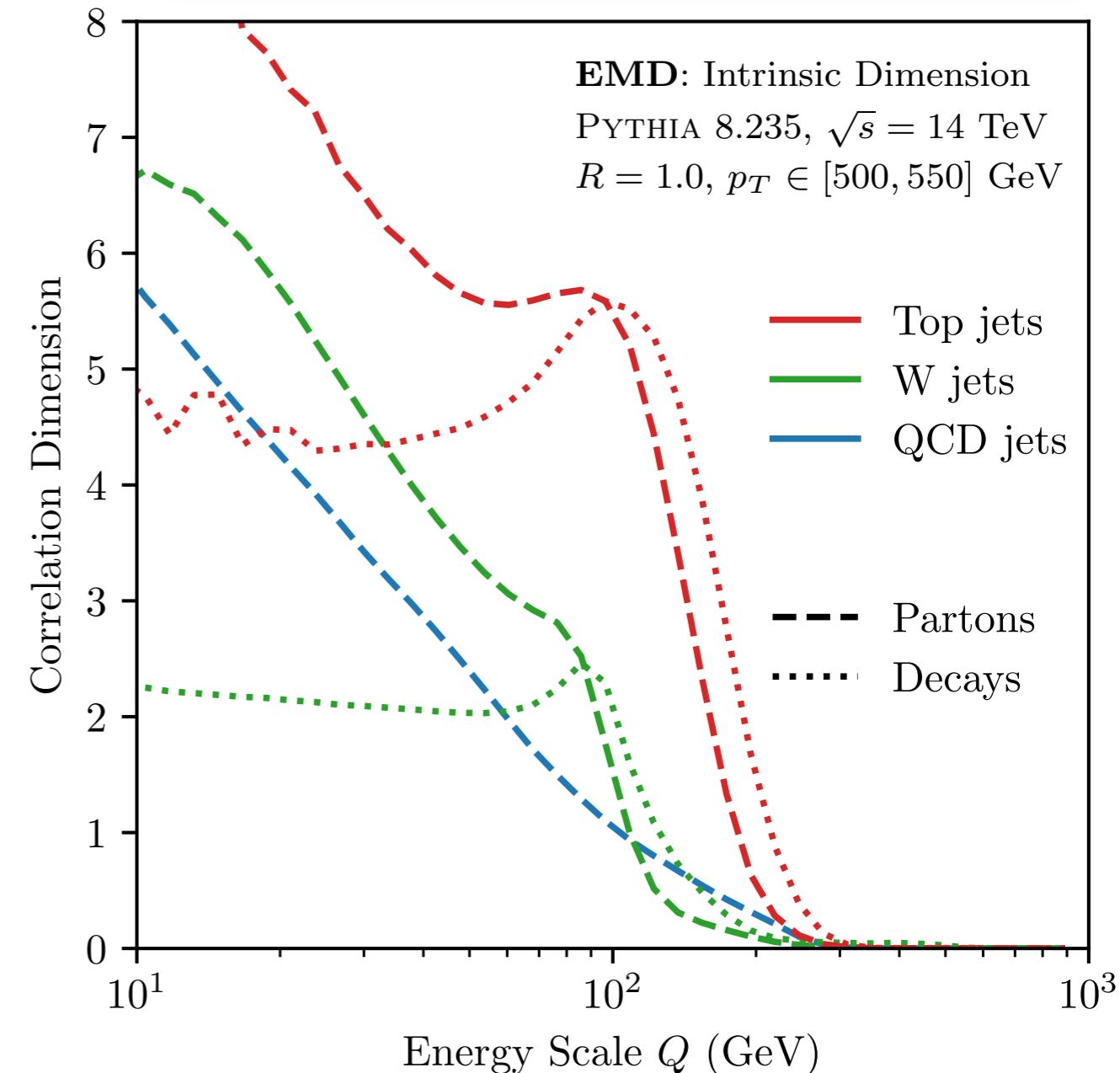
Correlation dimension: how does the # of elements within a ball of size  $Q$  change?



**Correlation dimension lessons:**

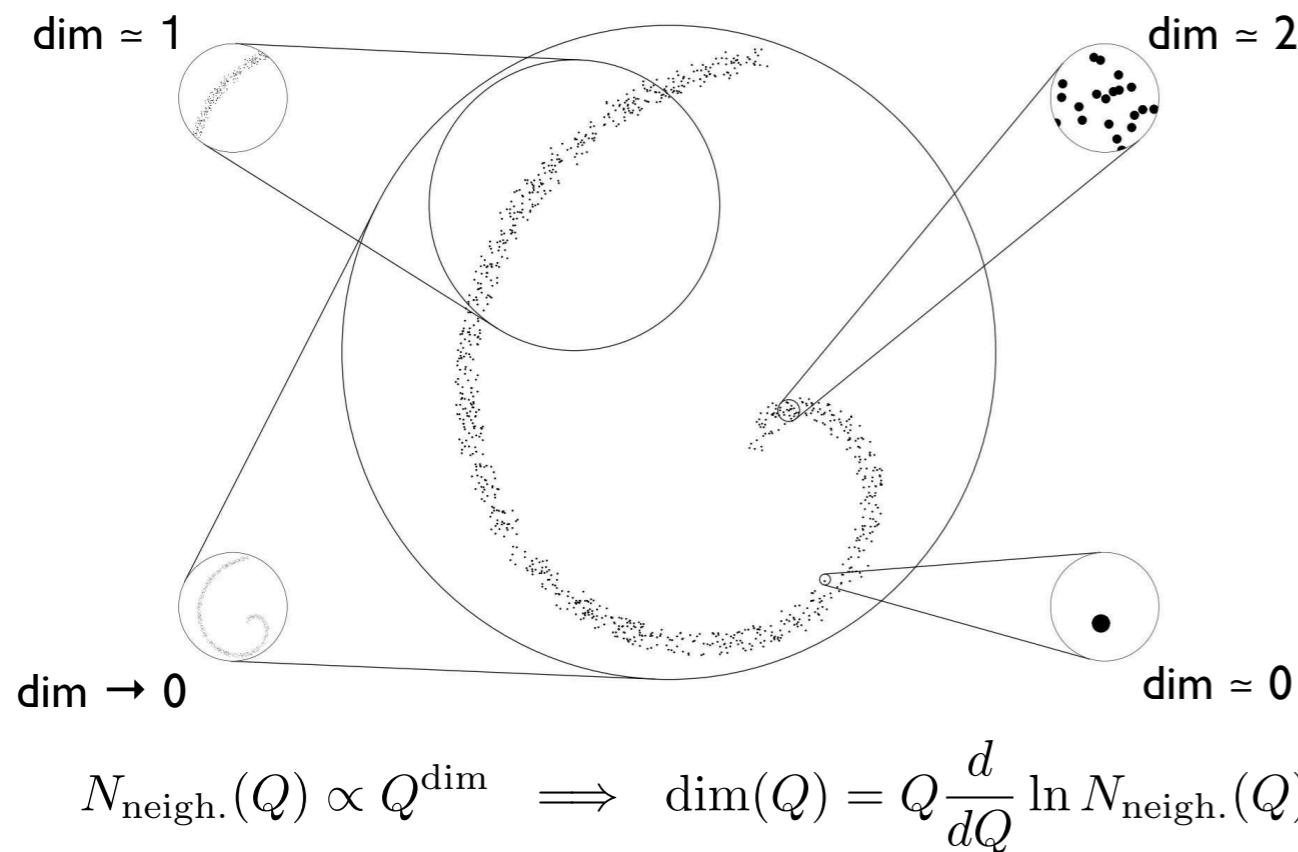
- Decays are "constant" dim. at low  $Q$
- Complexity hierarchy: QCD < W < Top
- Fragmentation increases dim. at smaller scales

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



# Manifold Dimensions of Event Space

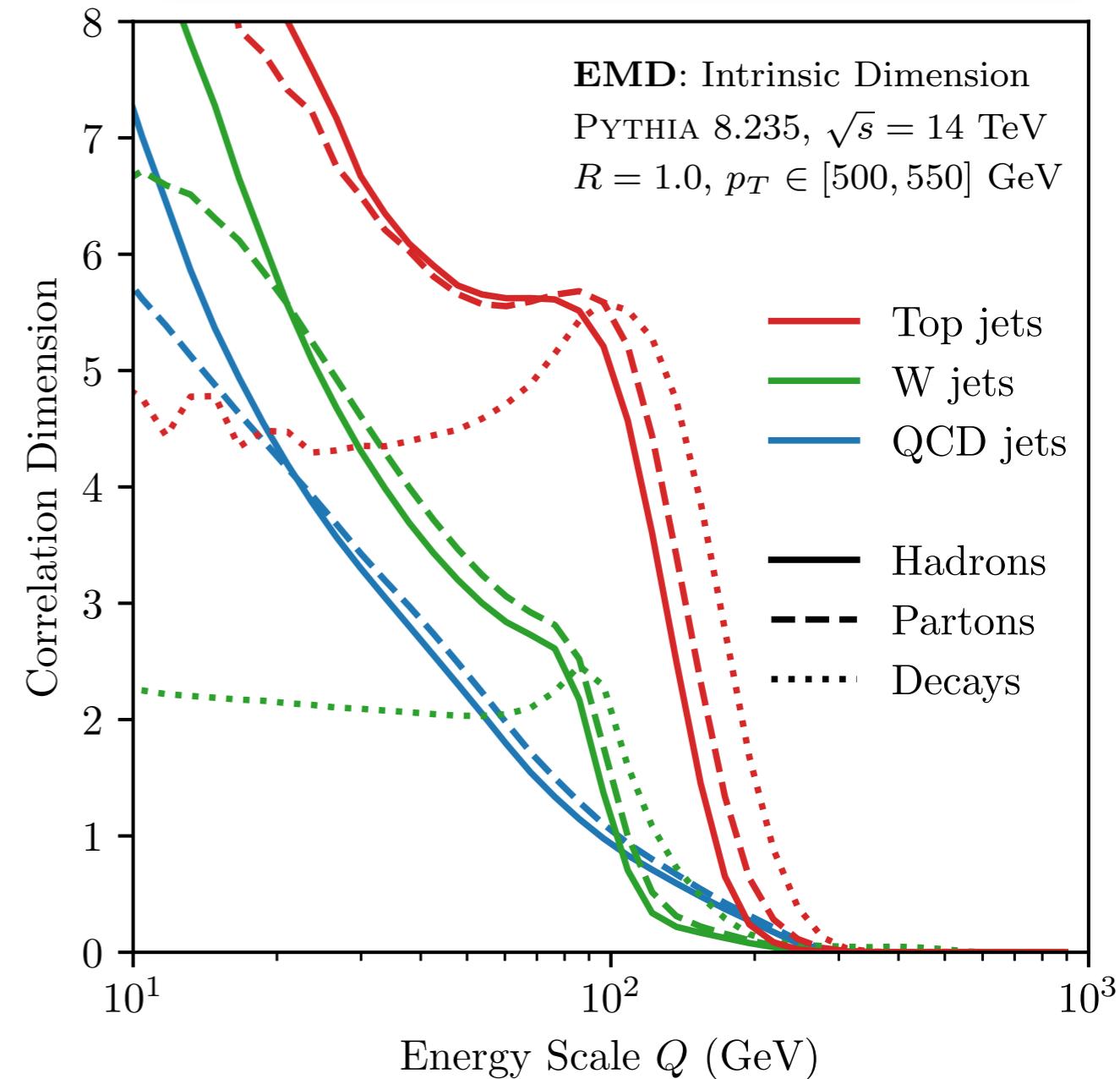
Correlation dimension: how does the # of elements within a ball of size  $Q$  change?



**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

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



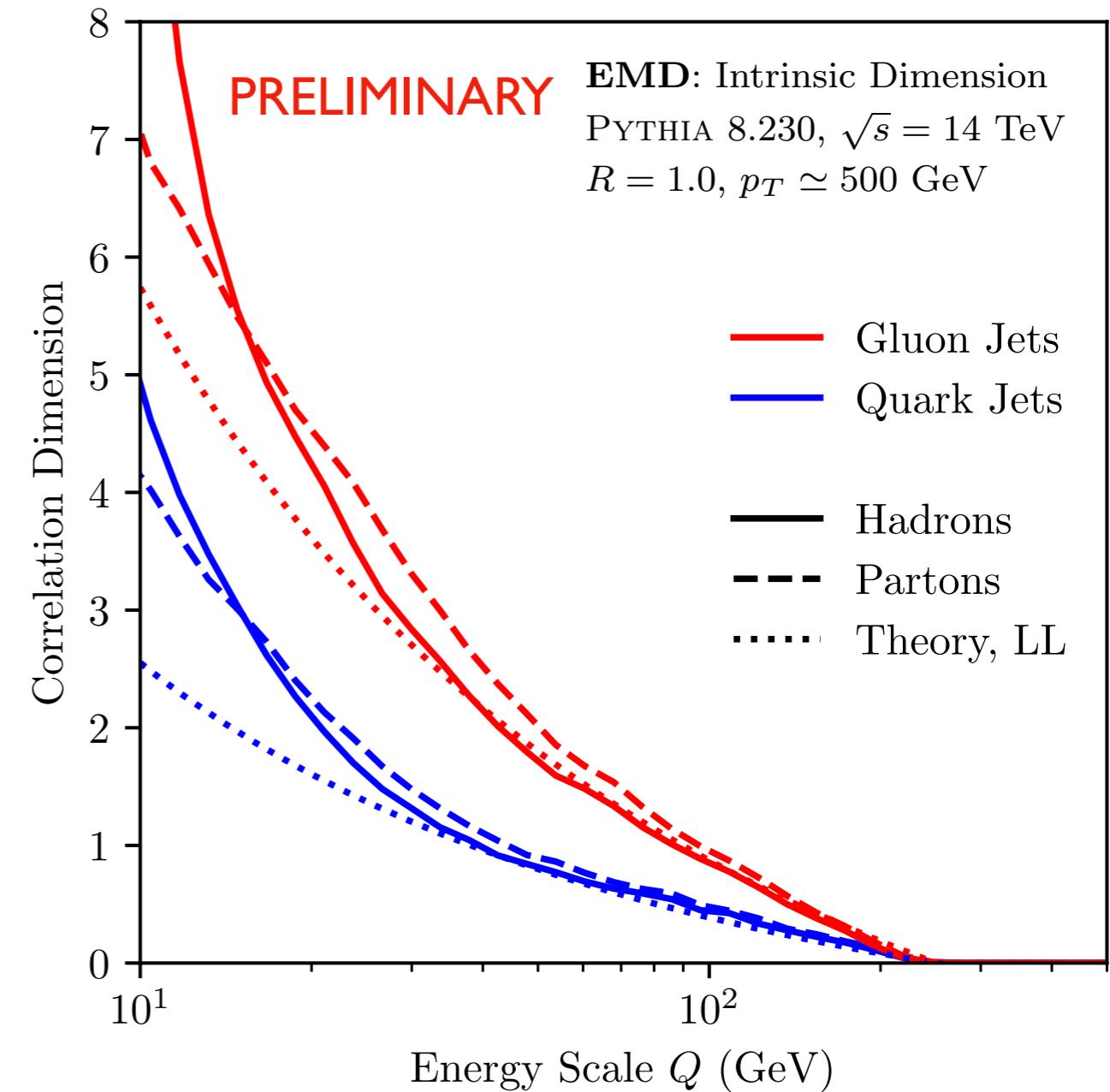
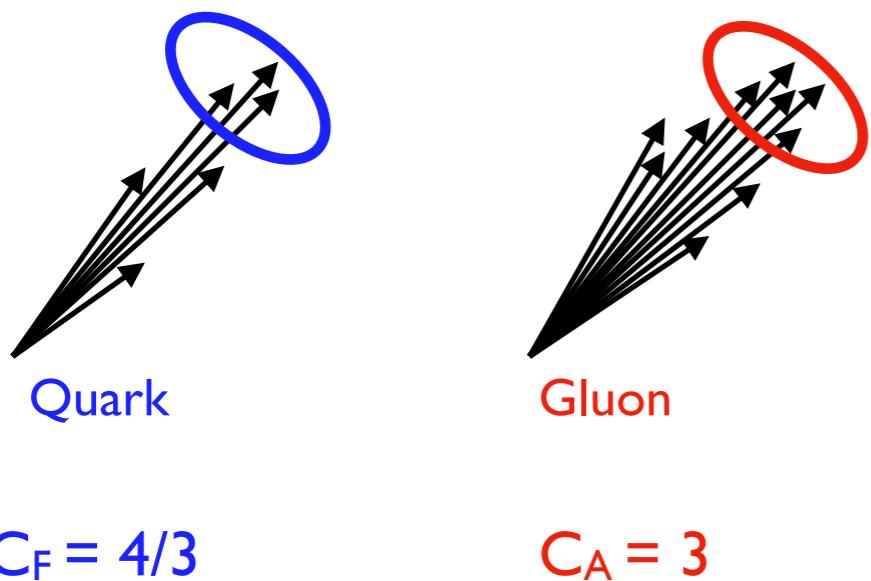
# Quark and Gluon Correlation Dimensions

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

Leading log calculation:

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

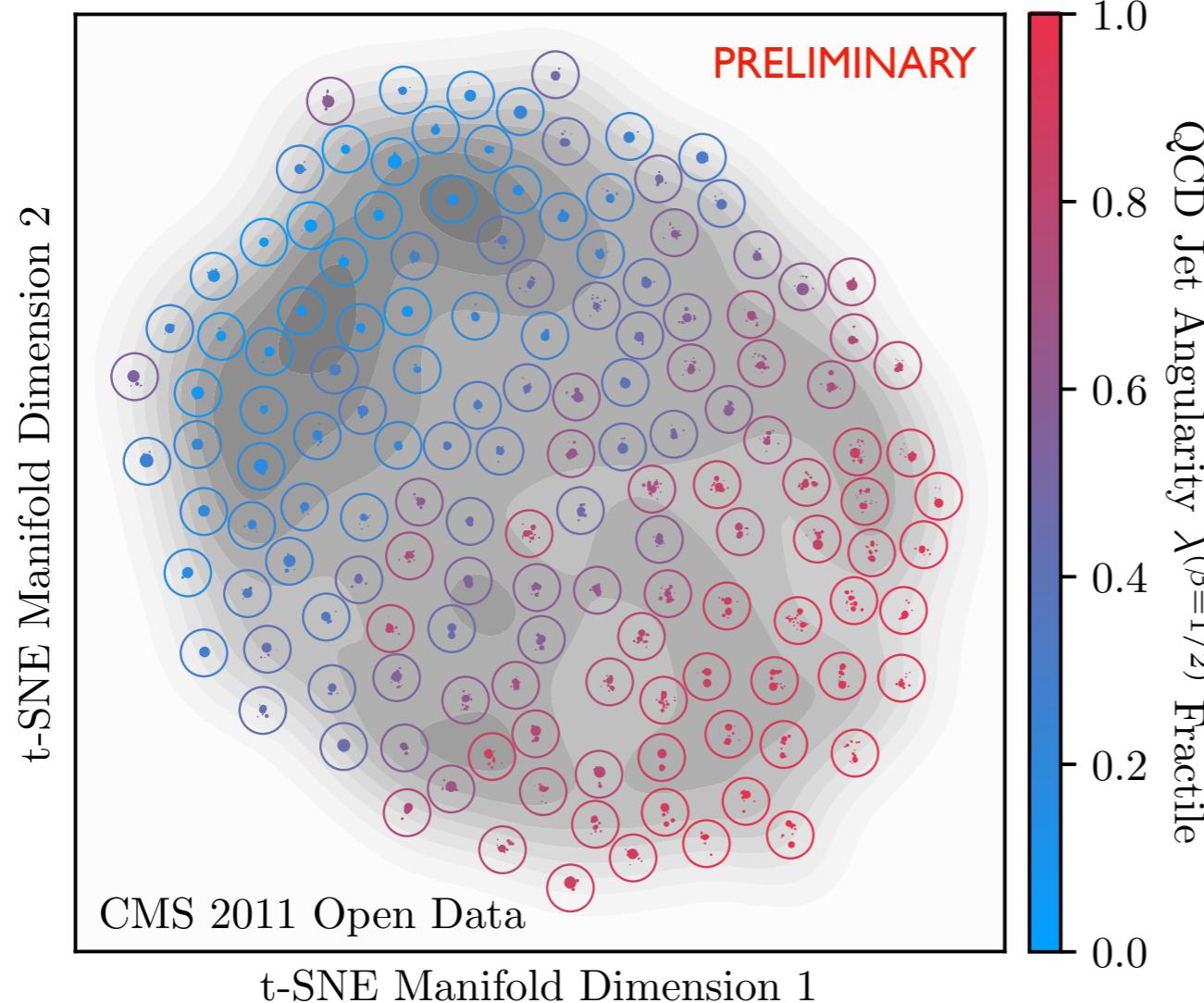
↑  
color factor



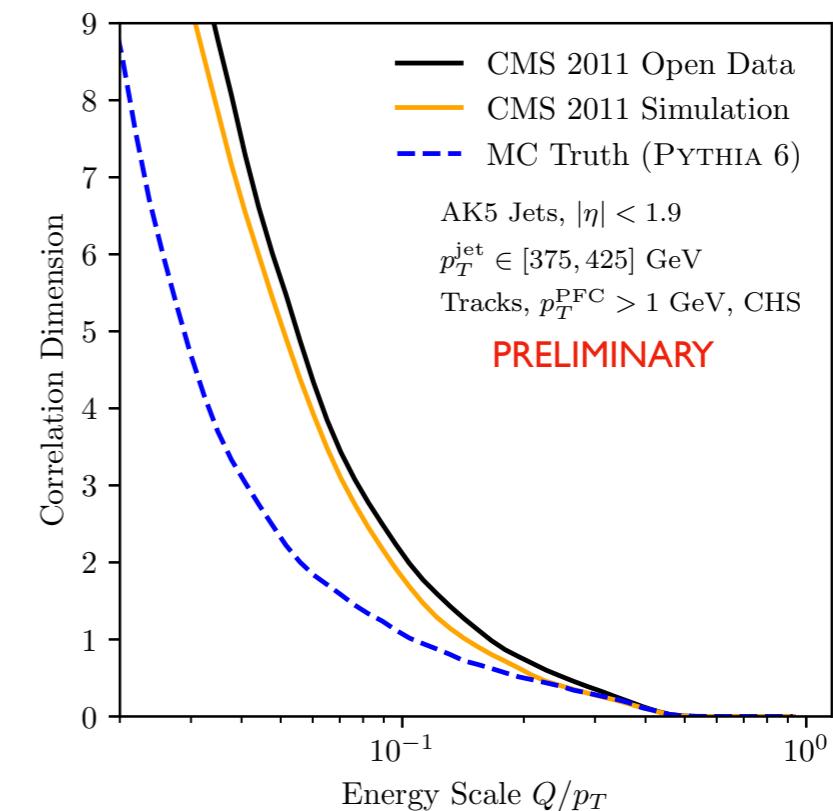
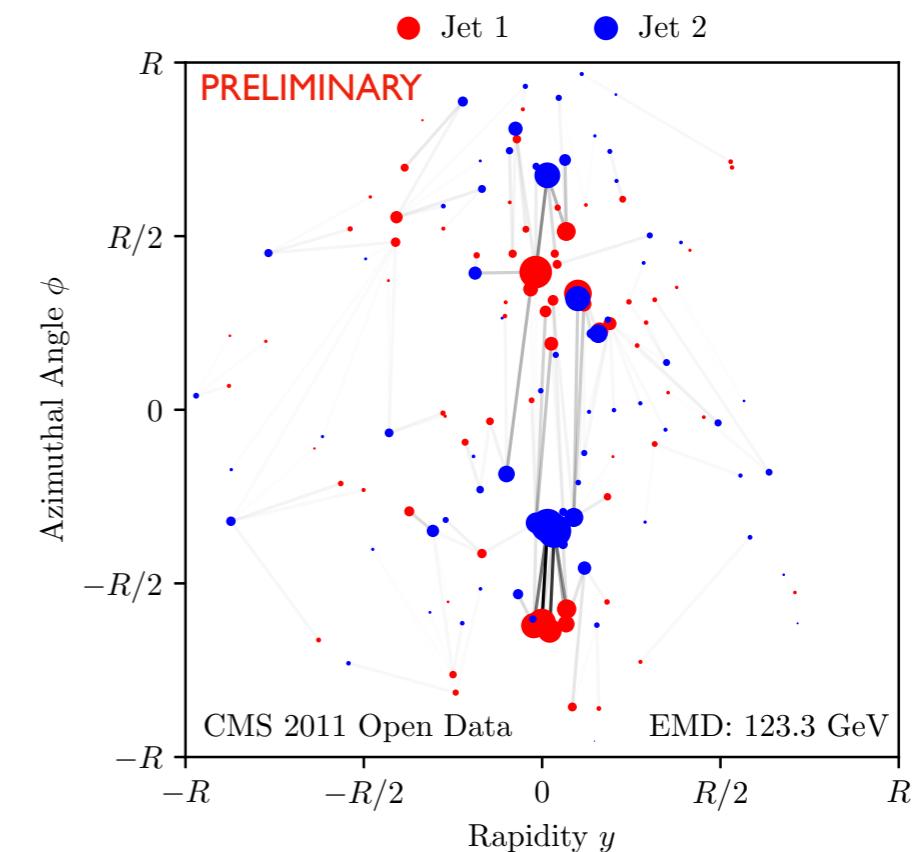
# Visualizing Jets with CMS Open Data



*Application of EMD techniques to the jet primary dataset in CMS 2011 Open Data*



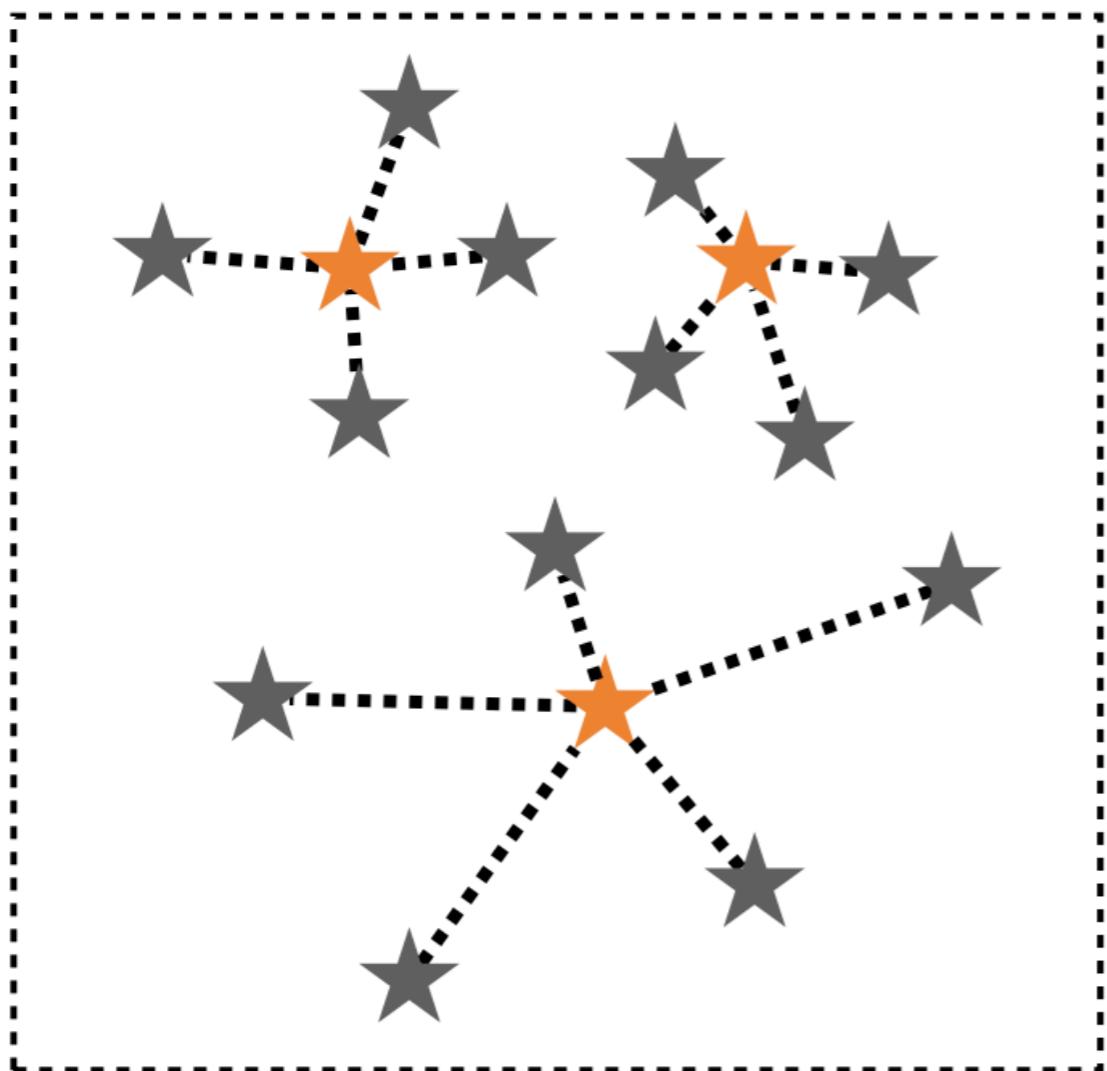
[Mastandrea, Naik, PTK, Metodiev, Thaler, *in progress*]



# Identifying Representative Jets

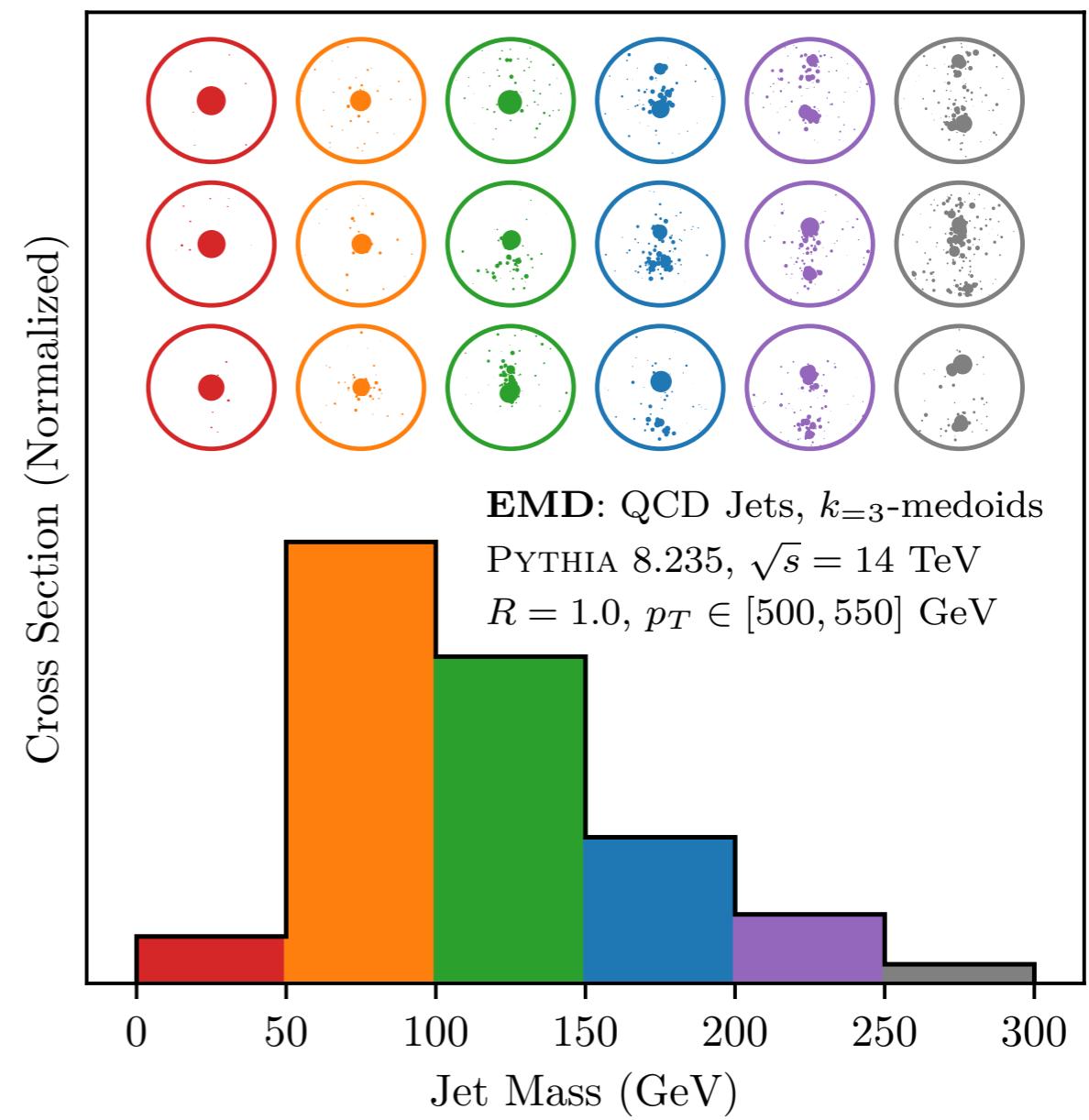
**medoid**: element selected to best represent a set of elements

**k-medoids**: k clusters to minimize total distance of points to medoids



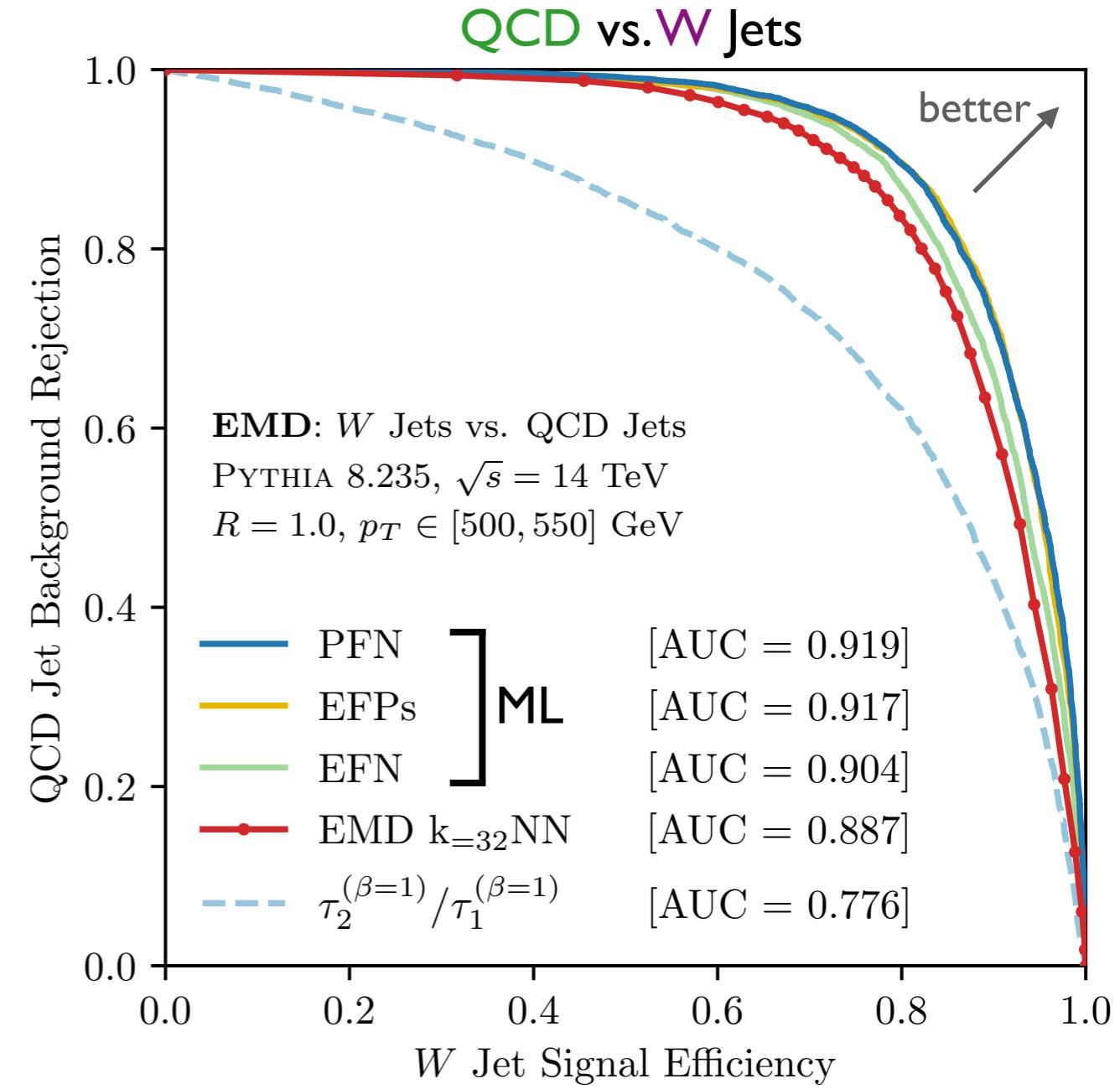
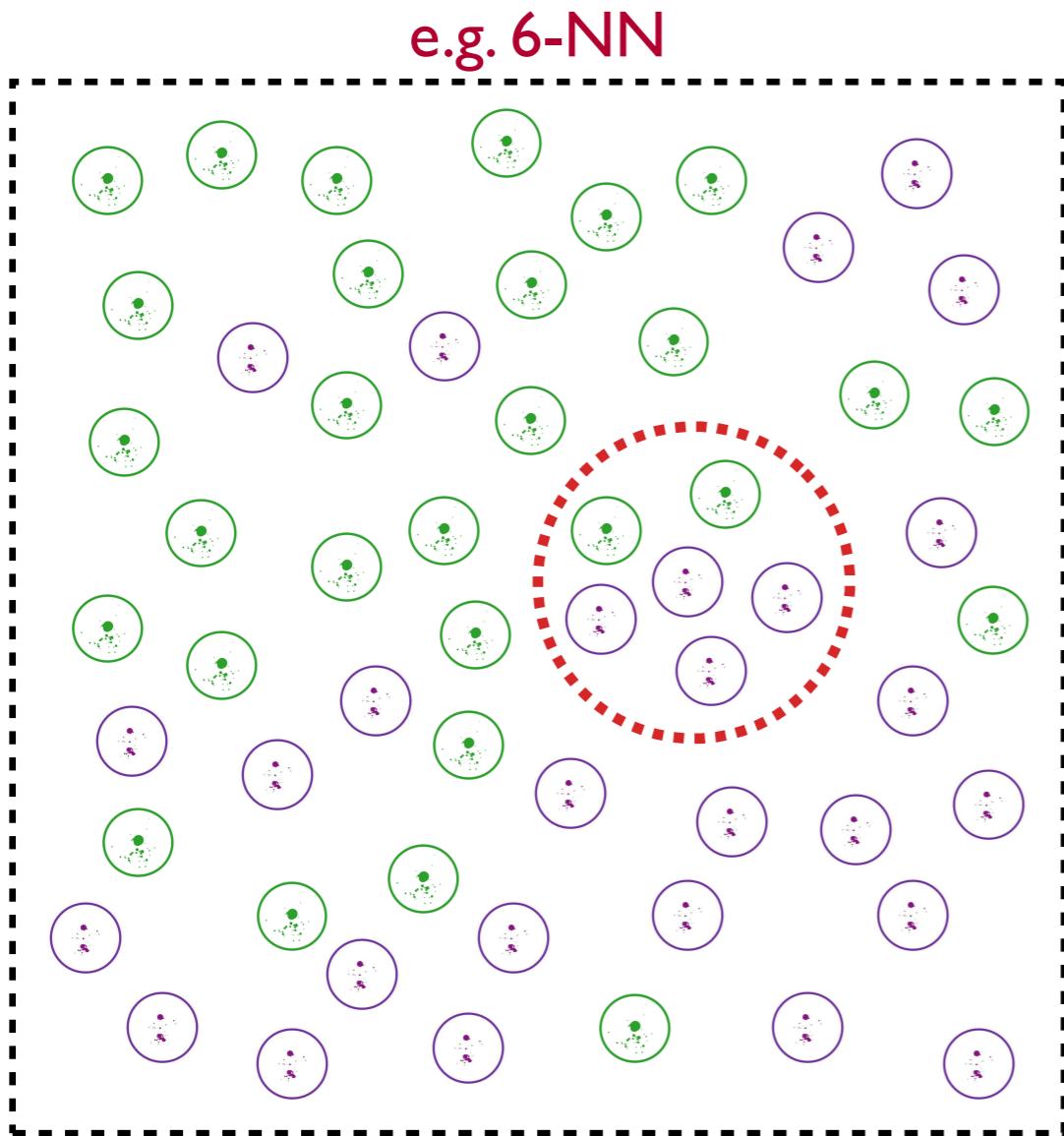
**3-medoid**

[diagram by Jesse Thaler]



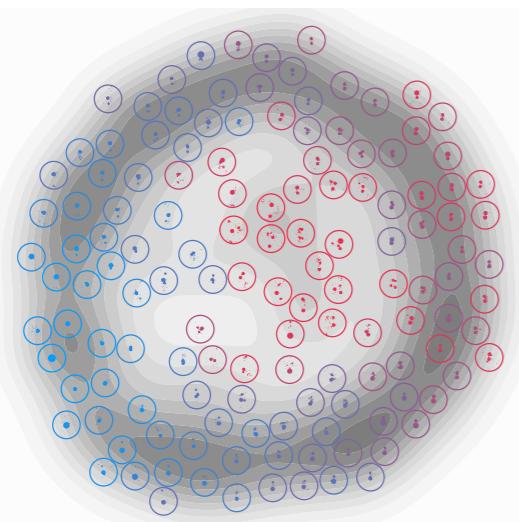
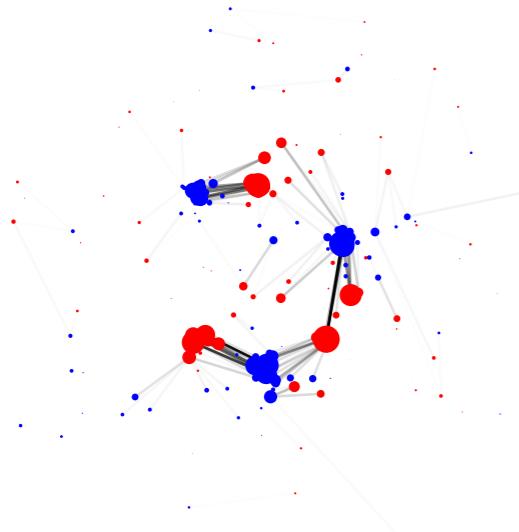
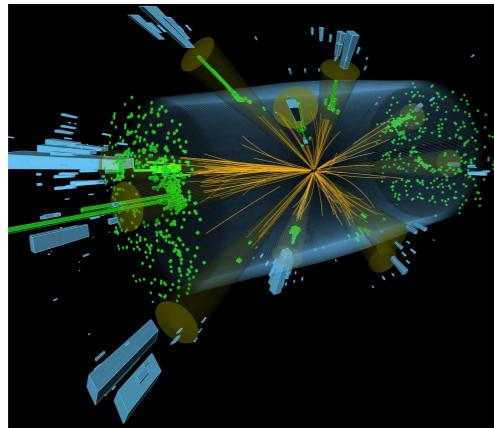
[PTK, Metodiev, Thaler, [1902.02346](#)]

# Jet Classification via Nearest-Neighbor Density Estimation



[PTK, Metodiev, Thaler, [1902.02346](#);

comparison to Thaler, Van Tilburg, [1011.2268](#), [1108.2701](#); PTK, Metodiev, Thaler, [1712.07124](#), [1810.05165](#);



## Collider Event Foundations

*IRC-safe energy flow is theoretically and experimentally robust*

## The Energy Mover's Distance

*Quantifies the difference in energy flow between events*

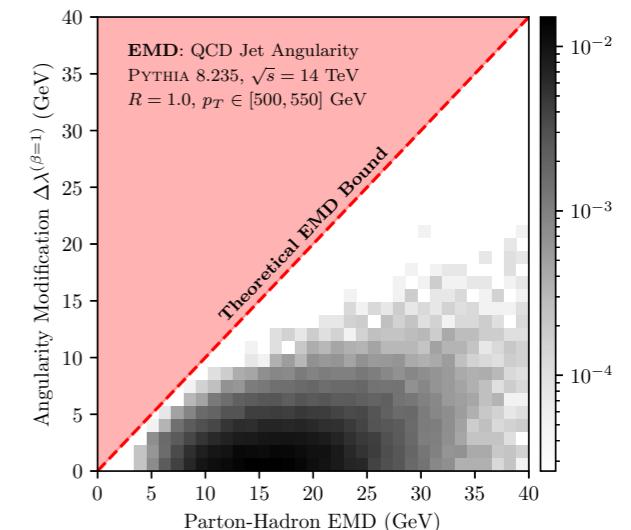
## Particle Physics Applications

*Quantifying modifications, visualizing and exploring event space*

# Further Directions

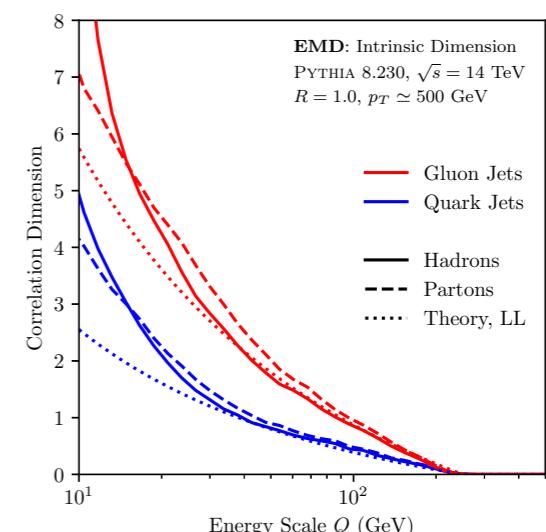
## Experimental

- Quantify (or even mitigate?) pileup/detector effects
- Non-parametric density estimates (unfolding?)
- Automated data compression (triggering?)



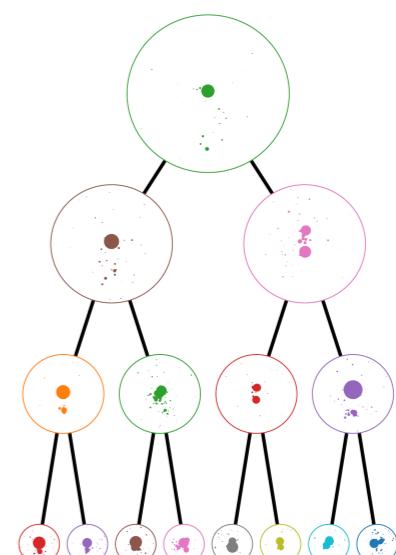
## Theoretical

- Define new observables with EMD?
- Precision **QCD** calculations of event space geometry?
- Event Mover's Distance between ensembles?



## Algorithmic

- Loss function for modern ML in particle physics?
- Metric trees to turn  $O(N^2)$  into  $O(N \log N)$ ?





# BOOST 2019

[BOOST 2019, July 22-26, MIT]

Phenomenology | Reconstruction | Searches | Algorithms | Measurements | Calculations  
Modeling | Machine Learning | Pileup Mitigation | Heavy-Ion Collisions | Future Colliders

# EnergyFlow Python Package

<https://energyflow.network>

Parallelized EMD calculations via the Python Optimal Transport library

Keras implementations of EFNs, PFNs, DNNs, CNNs, efficient EFP computation

Several detailed [examples](#) and [demos](#) for common use cases and visualization procedures

The screenshot shows the EnergyFlow documentation website. The header features a red logo with a white diamond shape and the word "EnergyFlow". Below it is a search bar labeled "Search docs". The main navigation menu includes links for "Home", "Welcome to EnergyFlow", "References", "Copyright", "Getting Started", "Installation", "Demo", "Examples", "FAQs", "Documentation", "Energy Flow Polynomials", "Architectures", "EMD", "Measures", "Generation", "Utils", and "Datasets". On the right, there's a "Docs » Home" breadcrumb link and a "Welcome to EnergyFlow" title. Below the title are three images: a colorful tangle of lines, a diagram of a "Per-Particle Representation" block, and a scatter plot titled "EMD: 125.4 GeV". The main content area describes EnergyFlow as a Python package for particle physics, mentioning its evolution from EFPs to EFNs and PFNs, and its implementation of the Energy Mover's Distance (EMD). It lists several key features:

- **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 events for several energy and angular measures as well as custom measures.
- **Energy Flow Networks:** EFNs are infrared- and collinear-safe models designed for learning from collider events as unordered, variable-length sets of particles. EnergyFlow contains customizable Keras implementations of EFNs.
- **Particle Flow Networks:** PFNs are general models designed for learning from collider events as unordered, variable-length sets of particles, based on the Deep Sets framework. EnergyFlow contains customizable Keras implementations of PFNs.
- **Energy Mover's Distance:** The EMD is a common metric between probability distributions that has been adapted for use as a metric between collider events. EnergyFlow contains code to

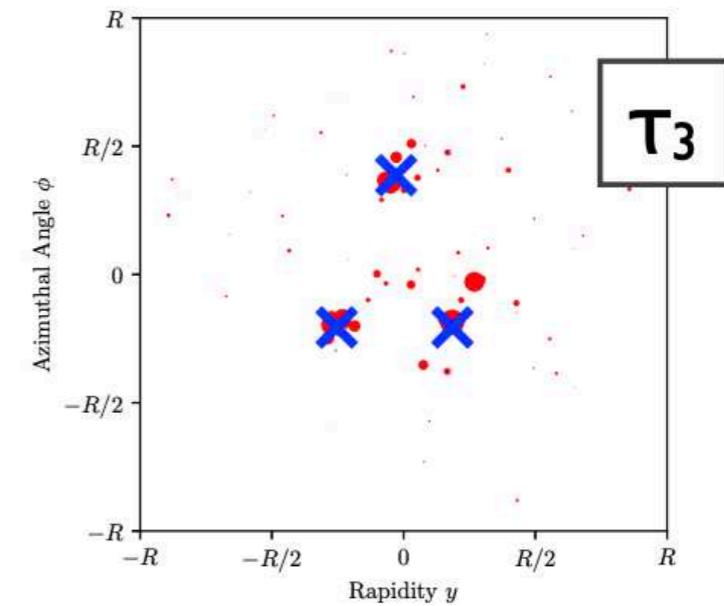
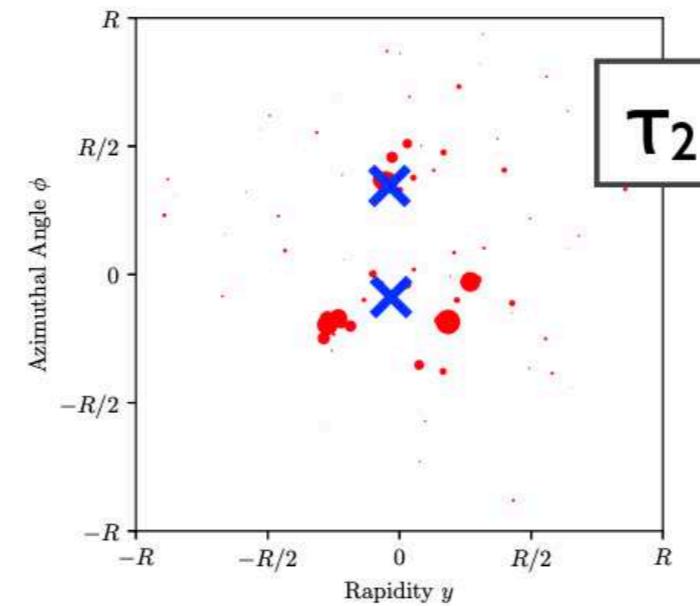
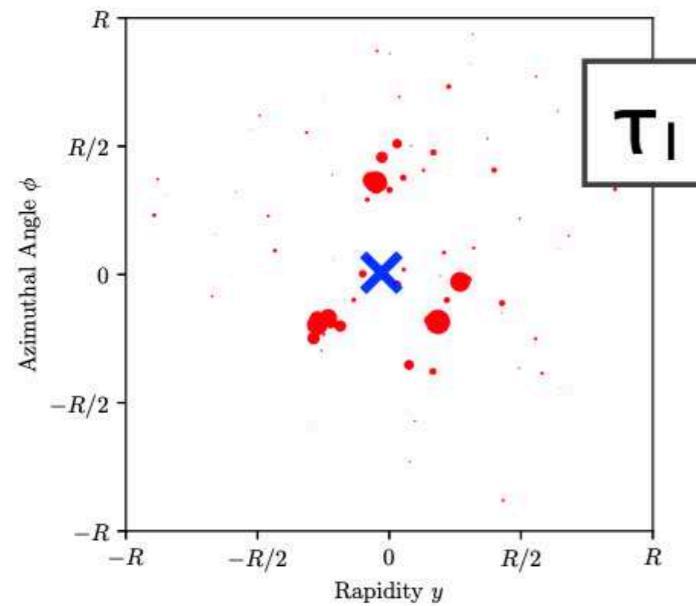
# Backup Slides

# Connection to N-subjettiness

$$\tau_N^{(\beta)}(\mathcal{E}) = \min_{N \text{ axes}} \sum_i E_i \min \left\{ \theta_{1,i}^\beta, \theta_{2,i}^\beta, \dots, \theta_{N,i}^\beta \right\}$$

↑ kind of arbitrary

IRC safe



$$\tau_N(\mathcal{E}) = \min_{|\mathcal{E}'|=N} \text{EMD}(\mathcal{E}, \mathcal{E}') \quad \text{for } \beta = 1$$

↑ very satisfying

Related to p-Wasserstein metric for  $p = \beta > 1$

[slide from talk by J.Thaler]

[JDT, Van Tilburg, [1011.2268](#), [1108.2701](#);  
based on Brandt, Dahmen, [ZPC 1979](#); Stewart, Tackmann, Waalewijn, [1004.2489](#)]

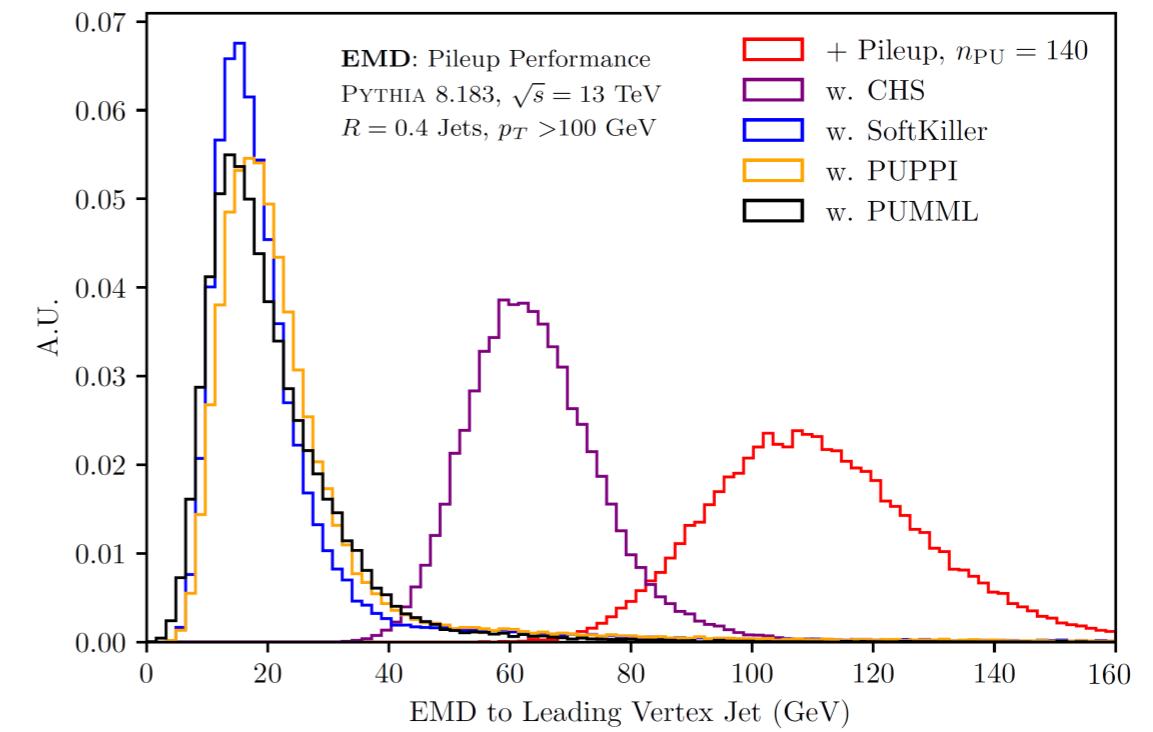
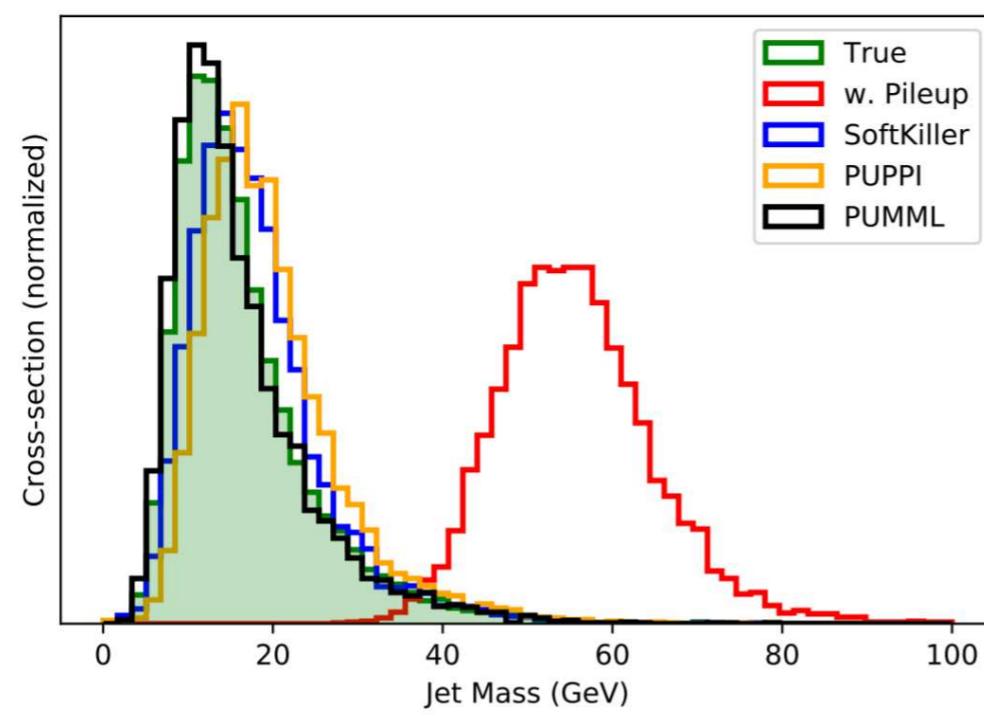
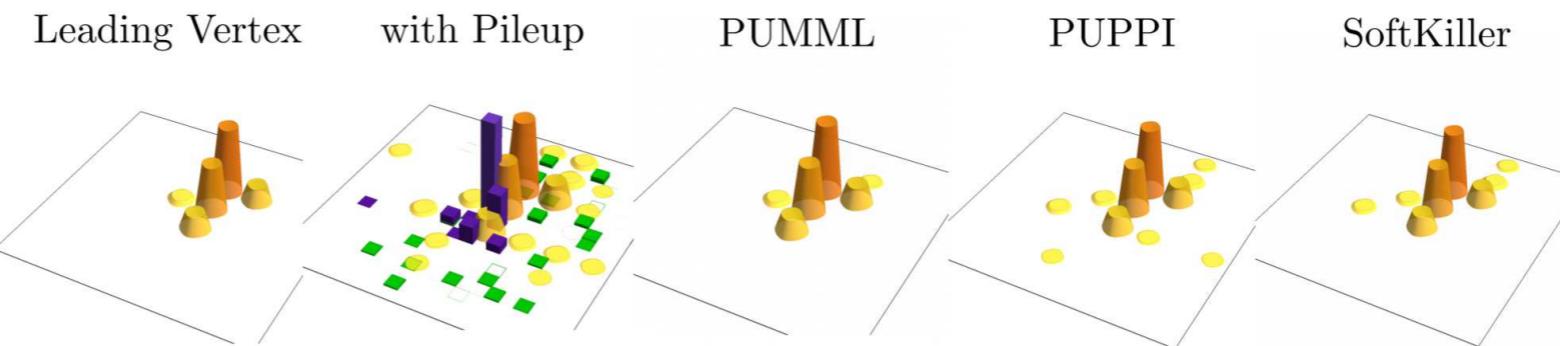
# Pileup Removal with Machine Learning (PUMML) and EMD

## PUMML with jet images

- pixel-based loss function
- compared specific IRC-safe observables

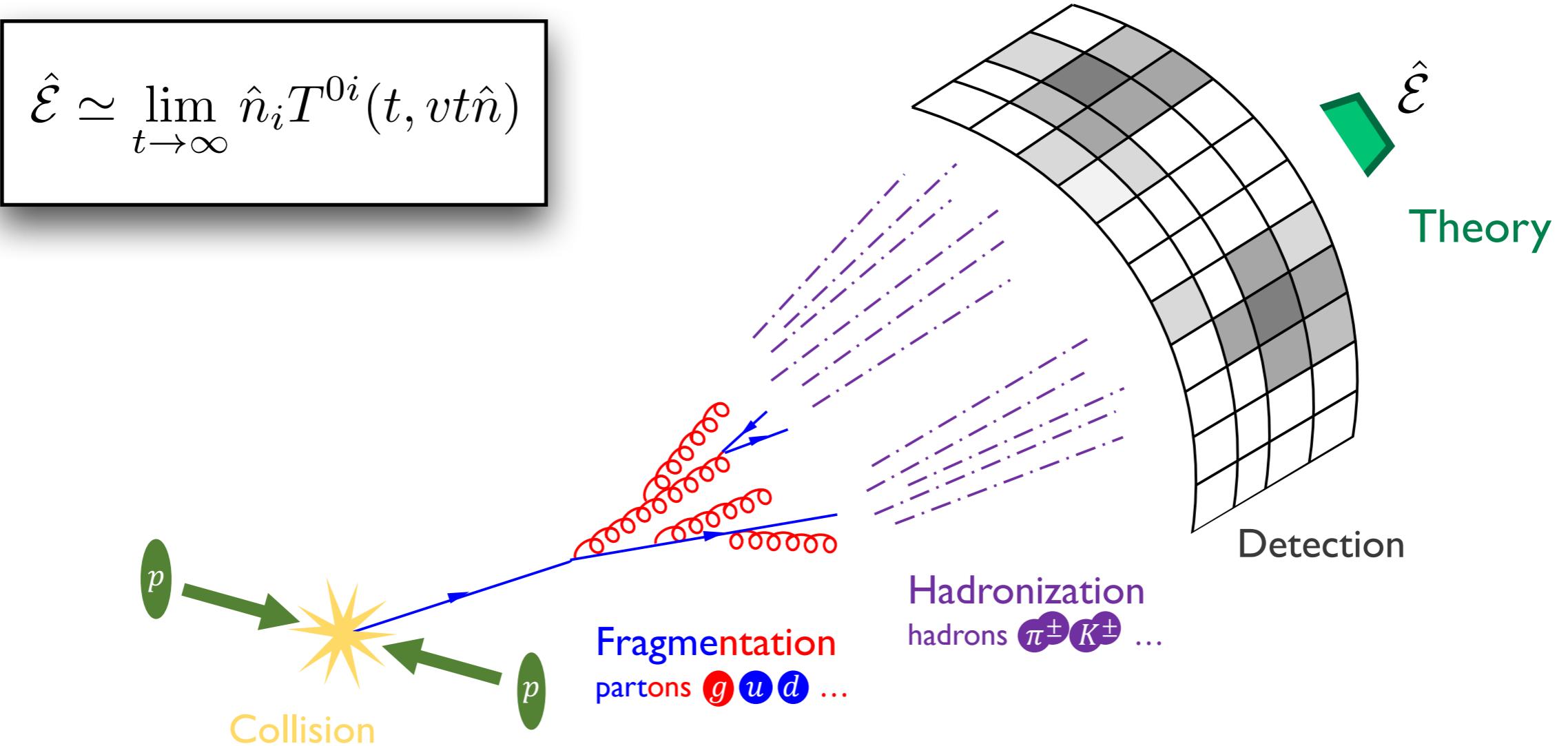
## PUMML with EMD

- no pixelation
- related to all IRC-safe observables



[PTK, Metodiev, Nachman, Schwartz, [J707.08600](#)]

# Stress-Energy Flow Operator



Stress-energy flow – measure of event/jet structure that is robust to non-perturbative and detector effects

[Sveshnikov, Tkachov, hep-ph/9512370; Hofman, Maldacena, 0803.1467; Mateu, Stewart, Thaler, 1209.3781; PTK, Metodiev, Thaler, 1712.07124, 1810.05165]