

Machine Learning

An Essential Toolkit for Particle Physics

Patrick T. Komiske III

Massachusetts Institute of Technology
Center for Theoretical Physics

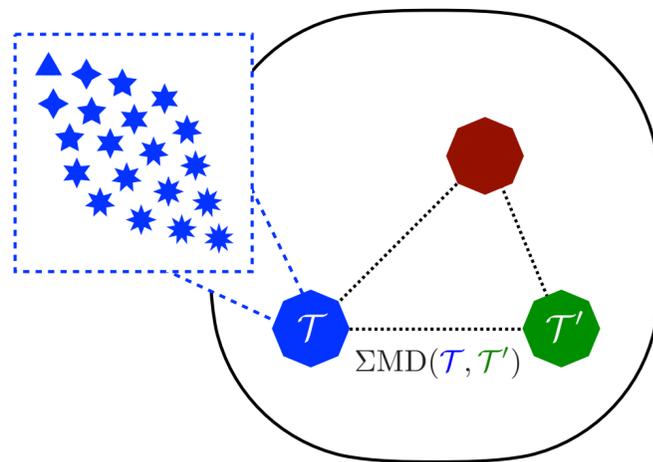
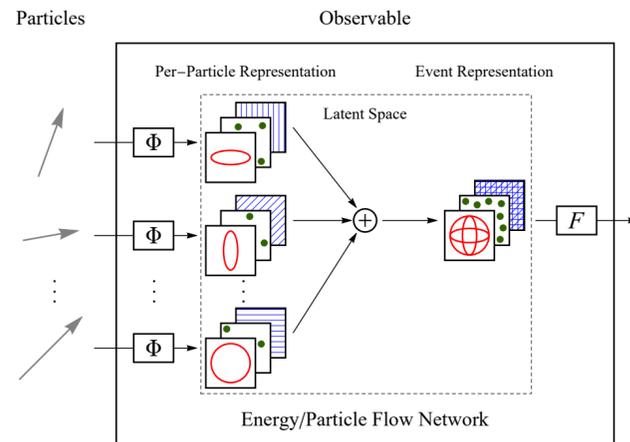
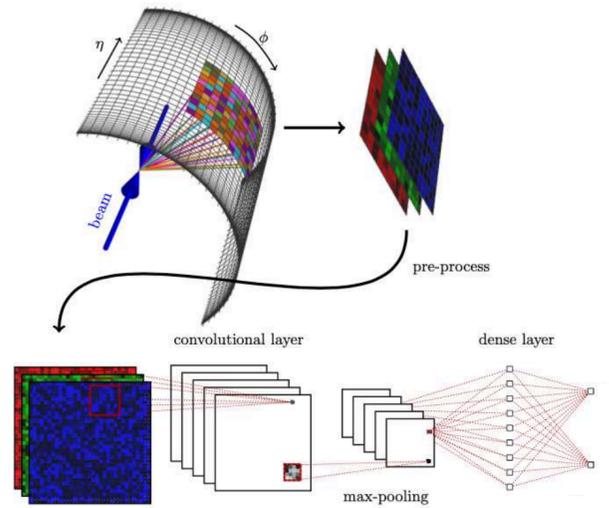
[Snowmass Computational Frontier Workshop](#) – ML Subgroup

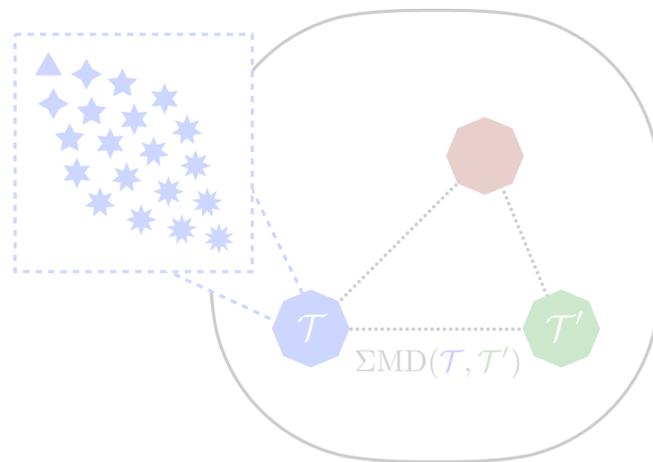
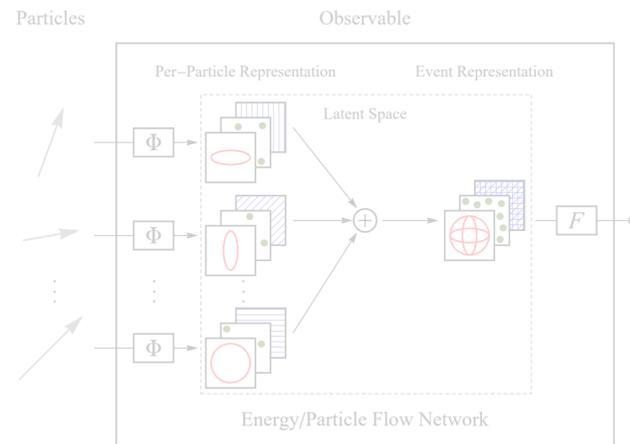
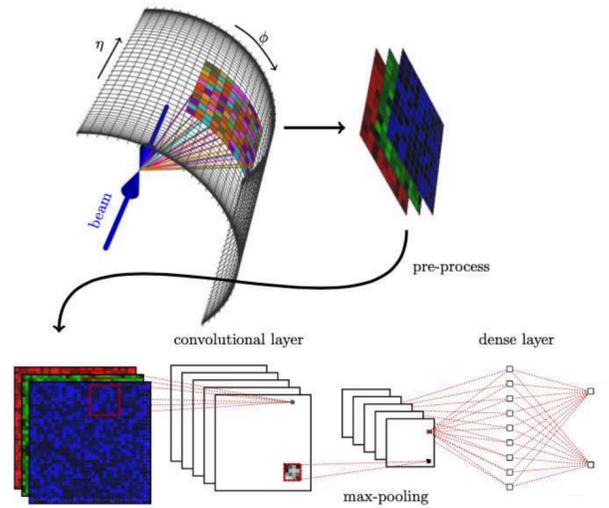
August 10, 2020

Ubiquity of ML in HEP

Lightning Review

Future Directions





Ubiquity of ML in HEP

Lightning Review

Future Directions

Machine Learning Permeates High-Energy Physics



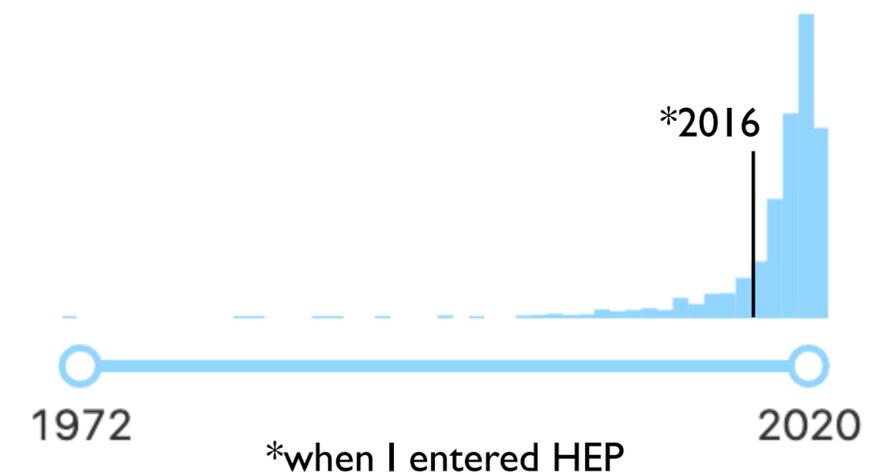
arXiv Category

<input type="checkbox"/> astro-ph.IM	316
<input type="checkbox"/> astro-ph.CO	262
<input type="checkbox"/> hep-ex	247
<input type="checkbox"/> hep-ph	245
<input type="checkbox"/> cs.LG	207
<input type="checkbox"/> physics.data-an	183
<input type="checkbox"/> astro-ph.HE	152
<input type="checkbox"/> stat.ML	144
<input type="checkbox"/> physics.ins-det	112
<input type="checkbox"/> astro-ph.GA	110
<input type="checkbox"/> physics.comp-ph	94
<input type="checkbox"/> hep-th	87
<input type="checkbox"/> gr-qc	85
<input type="checkbox"/> cs.CV	44
<input type="checkbox"/> nucl-th	44
<input type="checkbox"/> cond-mat.dis-nn	41
<input type="checkbox"/> cond-mat.stat-mech	41
<input type="checkbox"/> hep-lat	40
<input type="checkbox"/> quant-ph	39
<input type="checkbox"/> nucl-ex	37

Subject

<input type="checkbox"/> Astrophysics	556
<input type="checkbox"/> Computing	473
<input type="checkbox"/> Instrumentation	470
<input type="checkbox"/> Experiment-HEP	377
<input type="checkbox"/> Phenomenology-HEP	255
<input type="checkbox"/> Data Analysis and Statistics	185
<input type="checkbox"/> Other	172
<input type="checkbox"/> General Physics	138
<input type="checkbox"/> Gravitation and Cosmology	99
<input type="checkbox"/> Theory-HEP	80

Date of paper



ML in HEP is a young, vibrant, growing field with exciting potential!

Machine Learning Fundamentals

The Power of ML

Comes at a cost

▶ *Interpolation in high-dimensional spaces*

Combats the curse of dimensionality

Loses analytic understandability/tractability

▶ *Automatic feature extraction*

Ensures relevant features are not missed

Cannot easily convey what features are used

▶ *Asymptotically optimizes performance*

Provides useful/practical statistical power

Training is difficult with few global guarantees

Responsible ML Considerations

▶ *Available data*

Data source, number of samples, labels, reliability

▶ *Learning paradigm*

Fully/weakly/un-supervised, classification/regression/generation

▶ *Inputs and outputs*

Size/shape, symmetries, dimensionality

▶ *Model architecture*

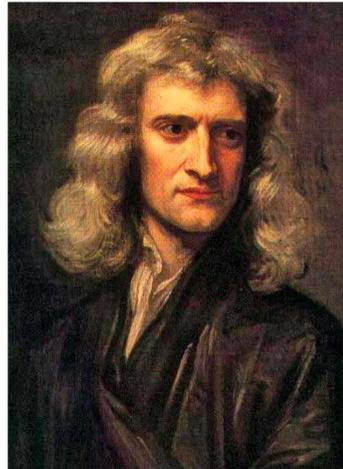
Expressibility, loss function, hyperparameters, validation/testing

▶ *Deployment strategy*

Model implementation, training/evaluation speed, uncertainties

My Perspective on ML in HEP

Machine learning is here to stay in high-energy physics

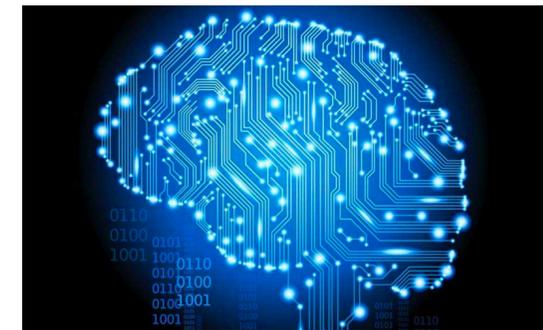


Calculus

- ▶ Fundamental to quantitative study of functions
- ▶ Taught to all undergraduate physics students
- ▶ Essential for understanding modern physics

Machine Learning

- ▶ Fundamental to statistical data analysis
- ▶ Increasingly taught in undergrad science programs
- ▶ Increasingly essential for modern physics/science



Machine learning straddles theory/experiment/computation divide

Close coordination between theory and experiment is essential in the current era of uncertainty in particle physics

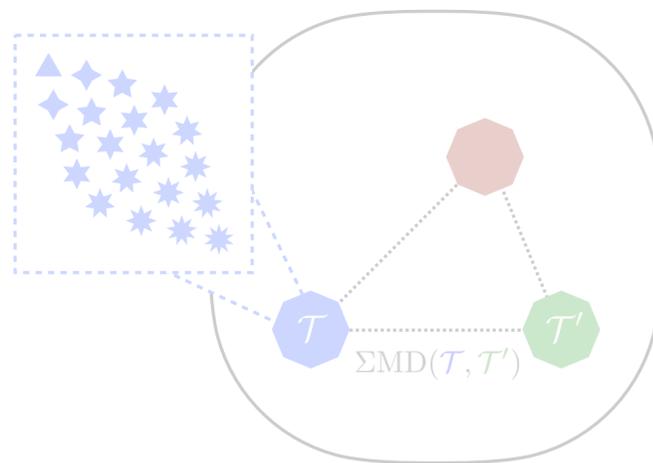
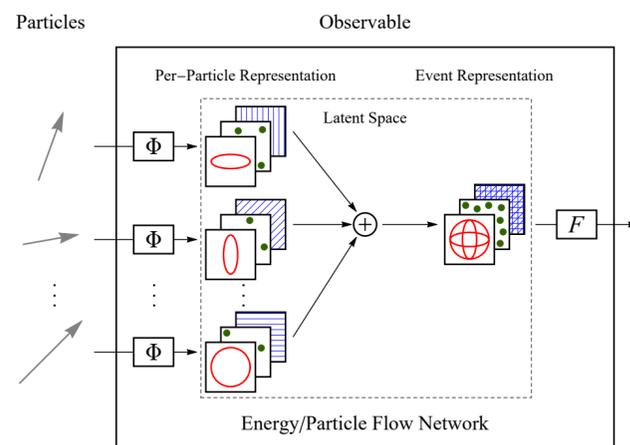
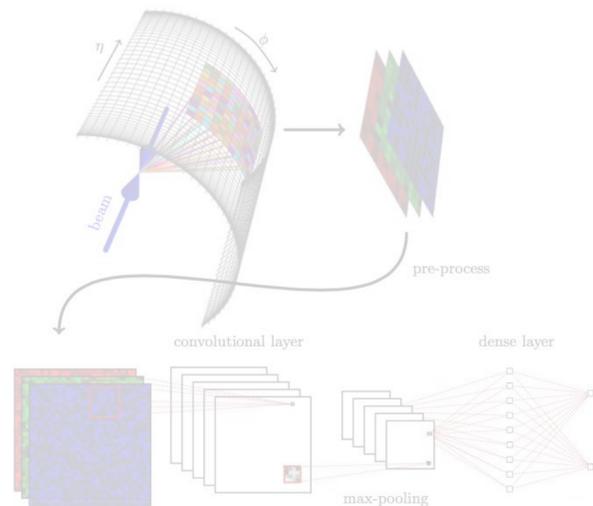
▶ *Where has ML already had an impact in high-energy physics?*

Part II of this talk

Key Questions During the Snowmass Process

▶ *What is the role of ML in particle physics (and vice-versa) in the future?*

Part III of this talk



Ubiquity of ML in HEP

Lightning Review

Disclaimers

- ▶ Examples biased towards collider physics
- ▶ References are not exhaustive

Future Directions

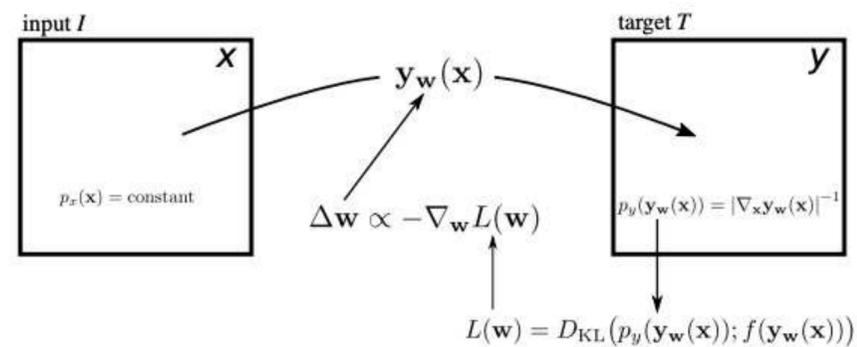
“Uncontroversial” Applications of Machine Learning in HEP

ML can provide a uniformly improved, drop-in replacement for some mathematical tasks

Efficient Monte Carlo Integration

[Bendavid, [1707.00028](#); Klimek, Perelstein, [1810.11509](#)]

“A continuum implementation of the VEGAS algorithm”

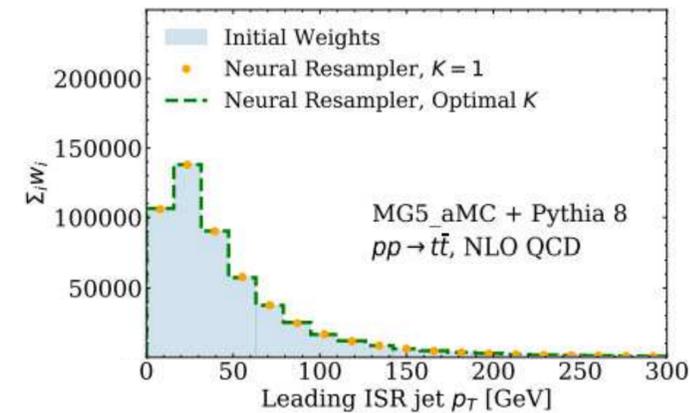


Algorithm	# of Func. Evals	$\sigma_w / \langle w \rangle$	σ_I / I (2e6 add. evts)
VEGAS	300,000	2.820	$\pm 2.0 \times 10^{-3}$
Foam	3,855,289	0.319	$\pm 2.3 \times 10^{-4}$
Generative BDT	300,000	0.082	$\pm 5.8 \times 10^{-5}$
Generative BDT (staged)	300,000	0.077	$\pm 5.4 \times 10^{-5}$
Generative DNN	294,912	0.083	$\pm 5.9 \times 10^{-5}$
Generative DNN (staged)	294,912	0.030	$\pm 2.1 \times 10^{-5}$

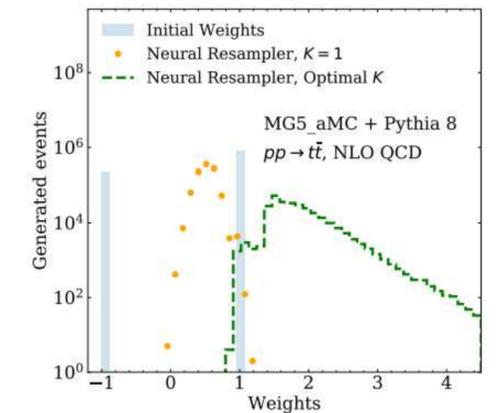
DNN has 10-100x smaller error with the same or fewer samples!

Positive MC-Weight Resampling

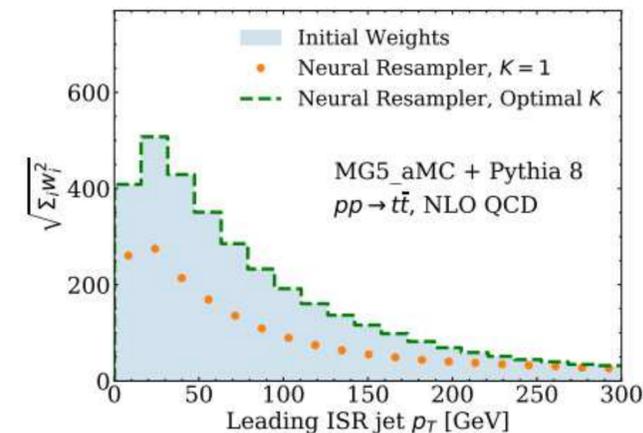
[Andersen, Gutschow, Maier, Prestel, [2005.09375](#);
Nachman, Thaler, [2007.11586](#)]



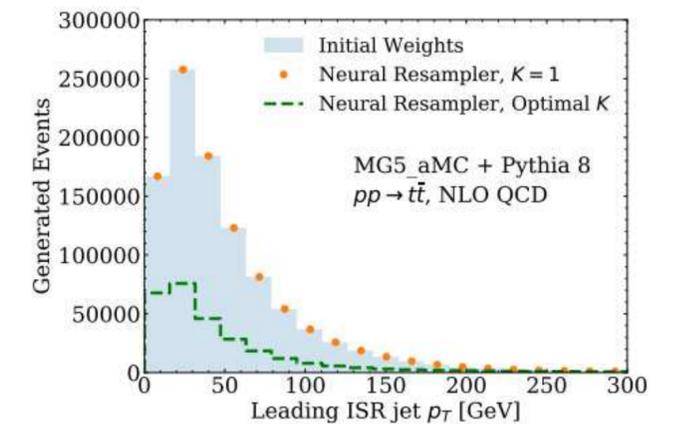
Central value preserved ...



using only positive weights ...



while preserving uncertainty ...



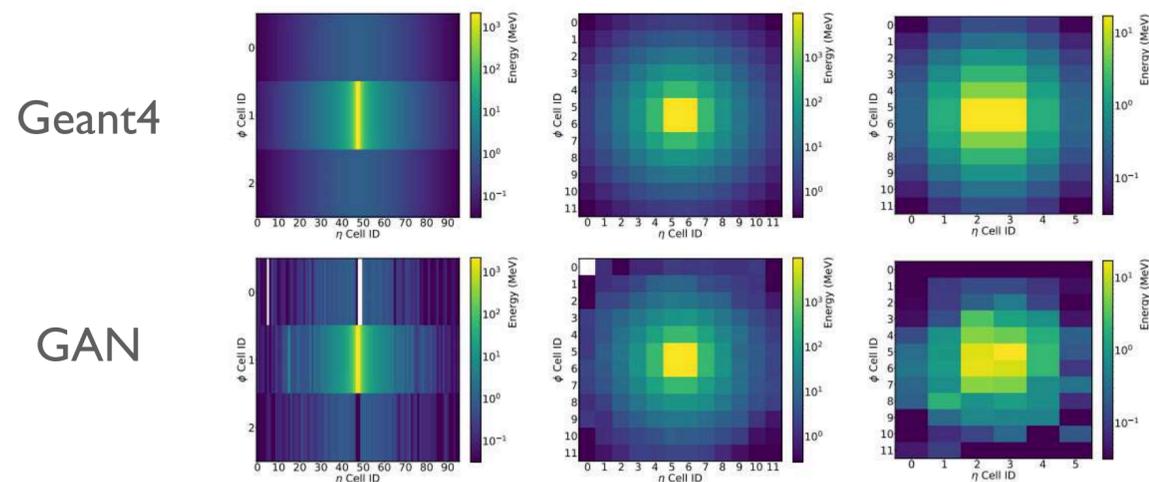
and requiring fewer events!

Improvement to Computational Speed/Efficiency with ML

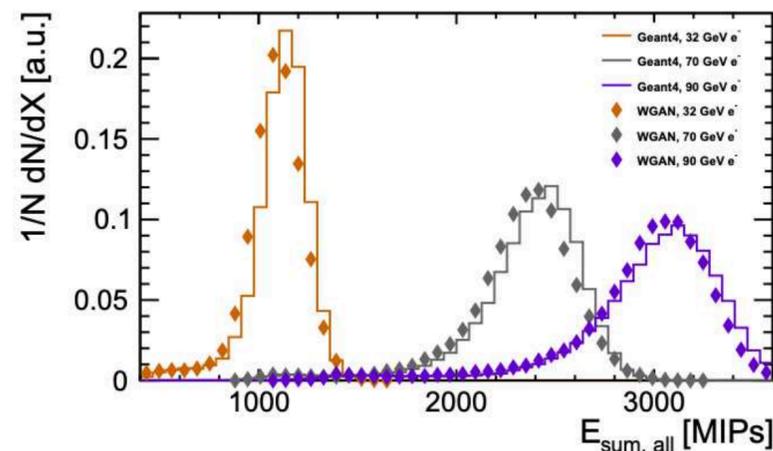
Fast Electromagnetic Calorimeter Simulation

[Paganini, de Oliveira, Nachman, [1705.02355](#) [1712.10321](#);
Erdmann, Glombitza, Quast, [1807.01954](#)]

Generative models can be $O(10^2-10^5)$ x
quicker than full detector-simulation



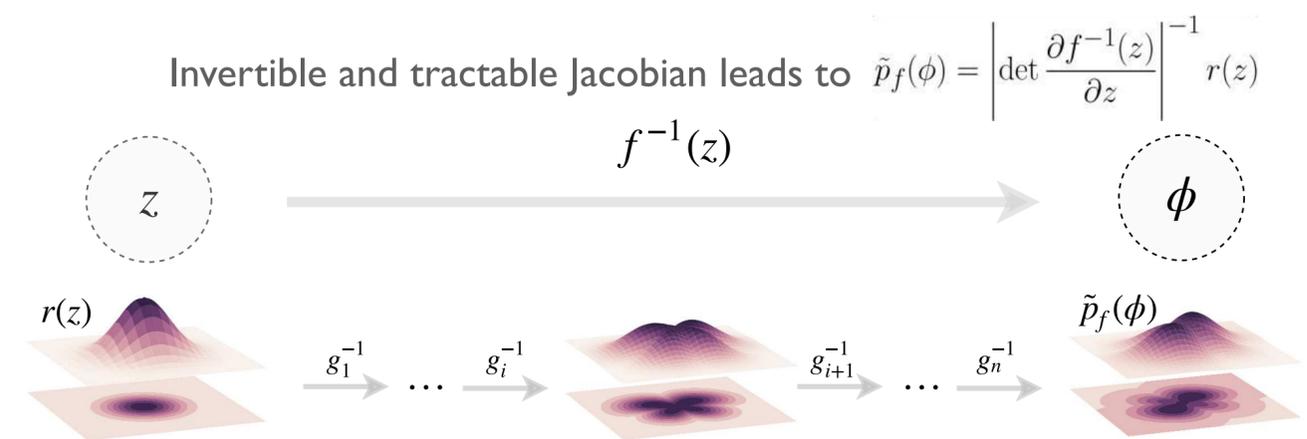
Wasserstein GANs
have stabler training
and good agreement



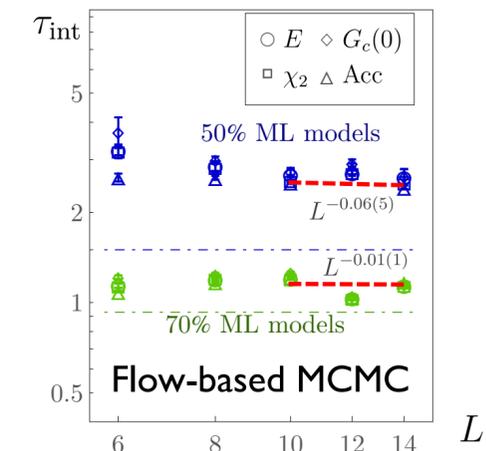
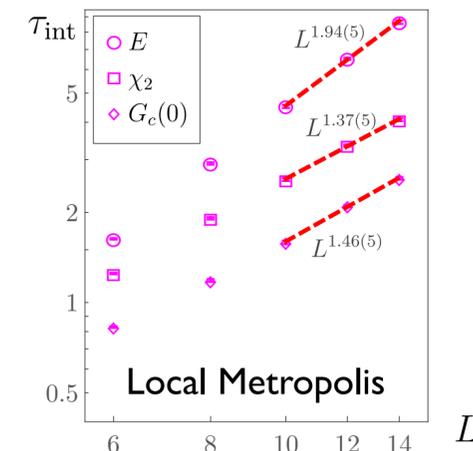
Improved MCMC for Lattice Field Theory

[Albergo, Kanwar, Shanahan, [1904.12072](#); Talk by G. Kanwar]

Normalizing flows can be used to
sample from complicated distributions



Power law growth of autocorrelation time avoided with ML

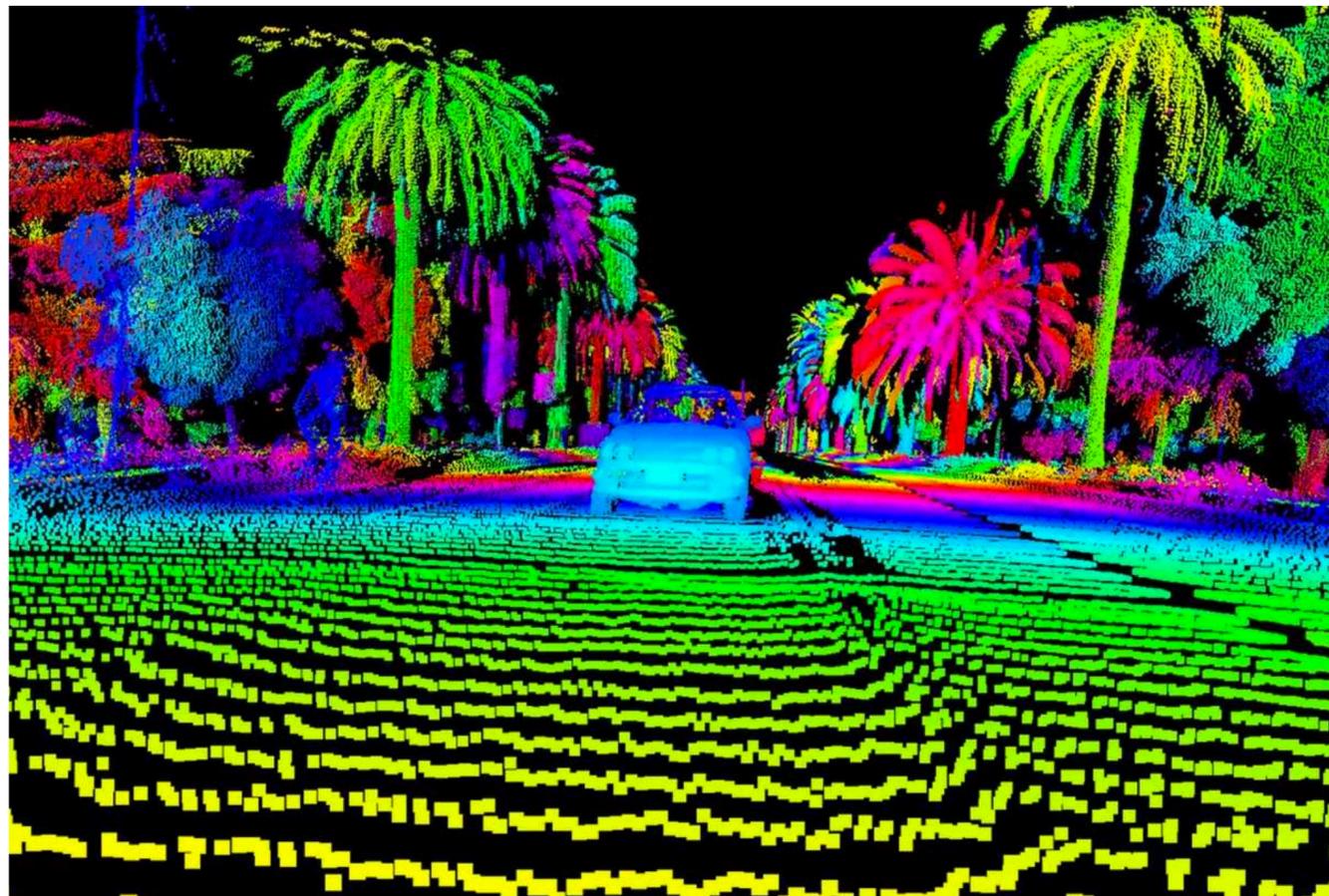


Neural Network Architectures for Particle Physics

Maximally appropriate ML architectures respect symmetries of the underlying data

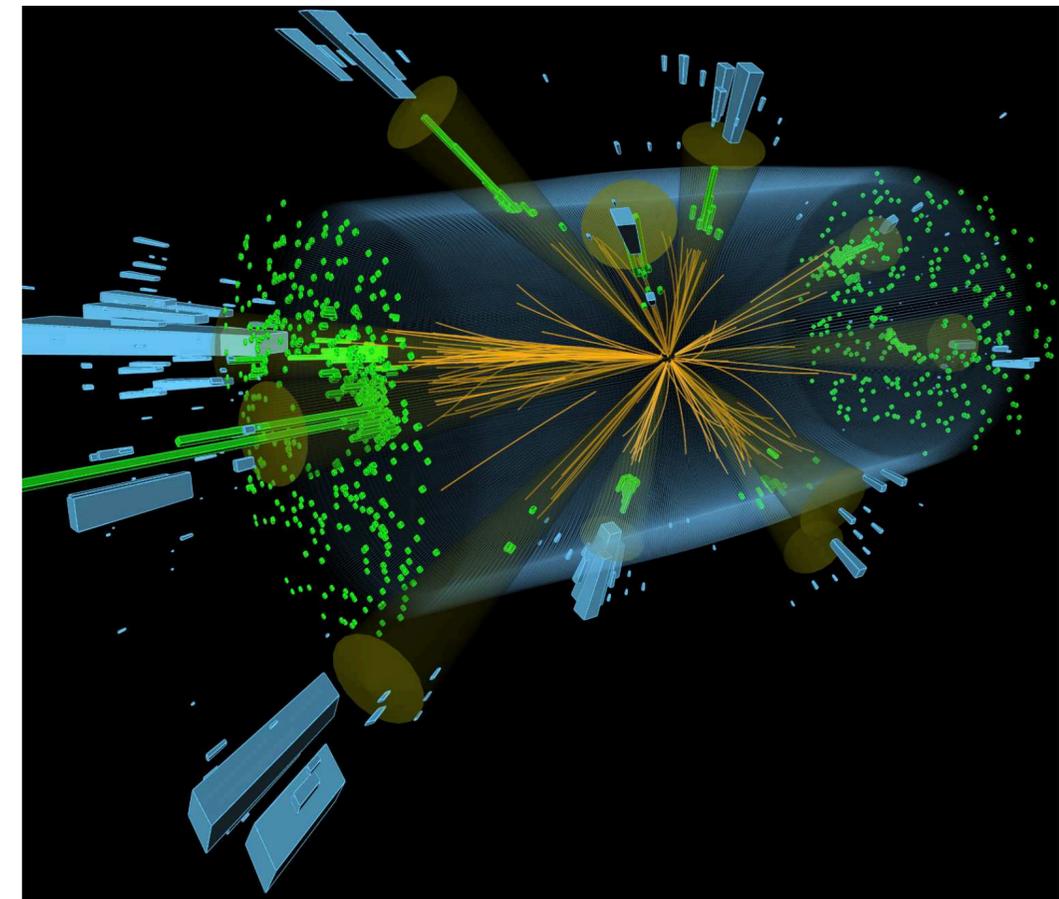
Particle physics events are naturally point clouds (alternatively, images e.g. calorimeters)

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



LIDAR data from self-driving car sensor

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

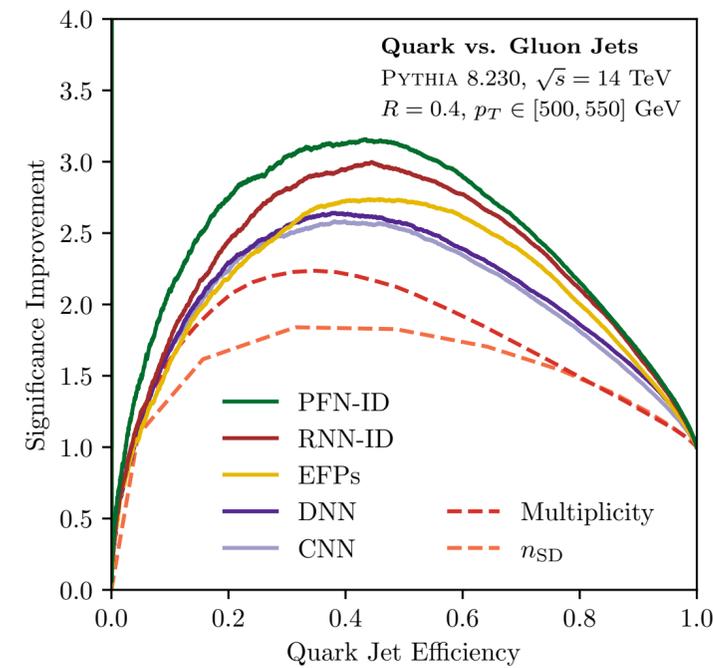
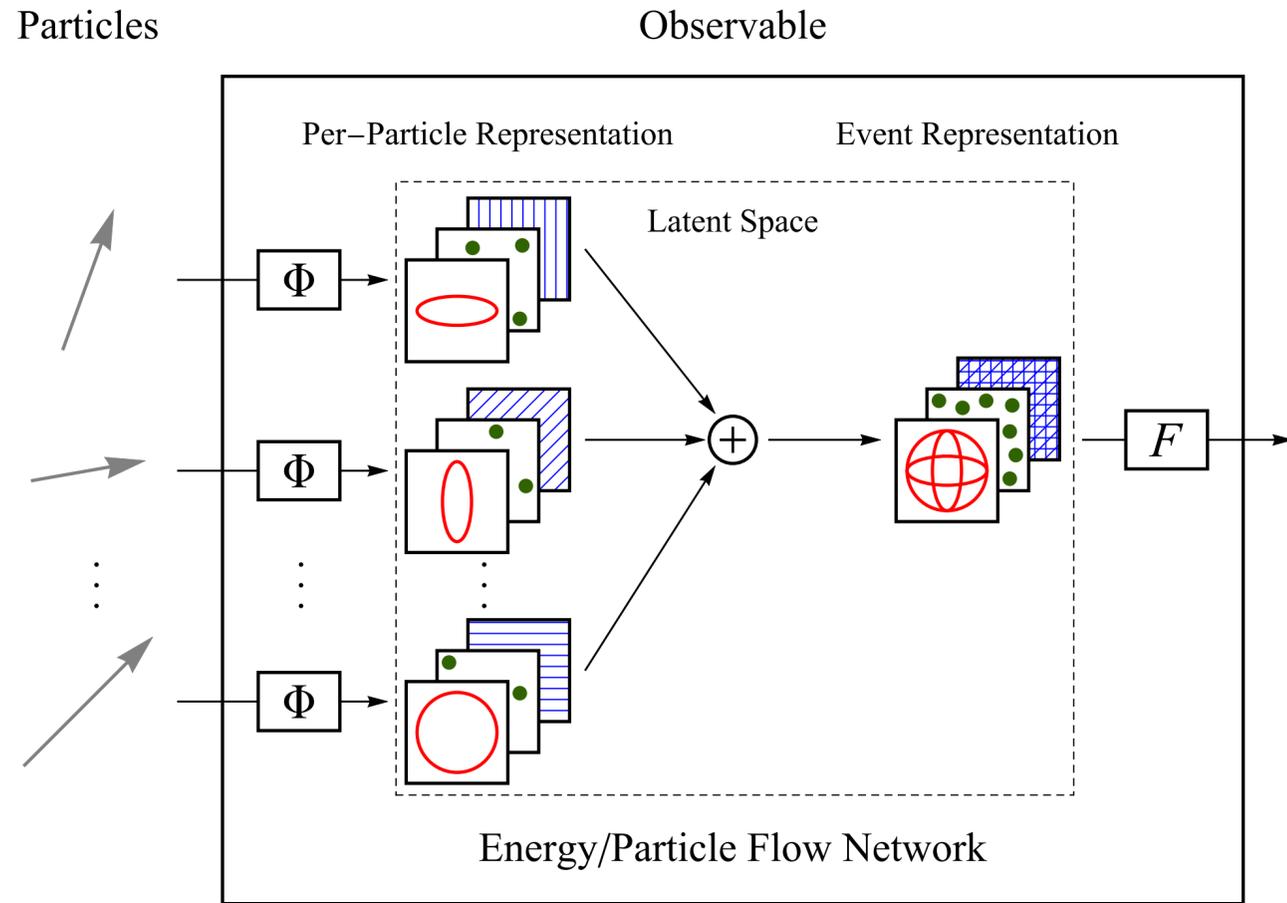


Multi-jet event at CMS

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

Deep Sets for Particle Jets

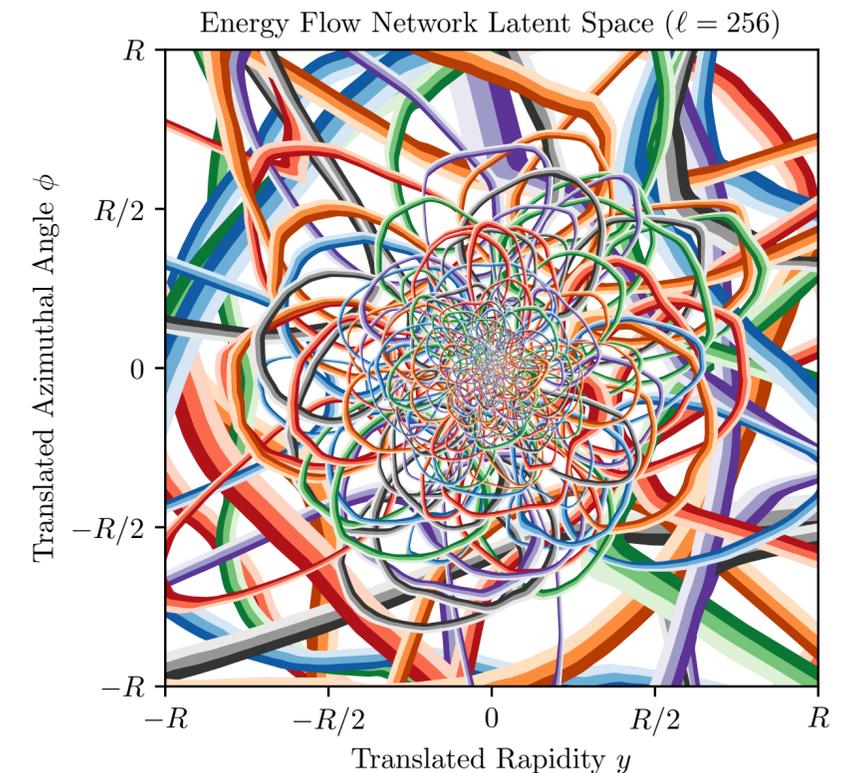
[Zaheer, Kottur, Ravanbakhsh, Póczos, Salakhutdinov, Smola, [1703.06114](#);
PTK, Metodiev, Thaler, [1810.05165](#);
[EnergyFlow Python Package](#)]



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)

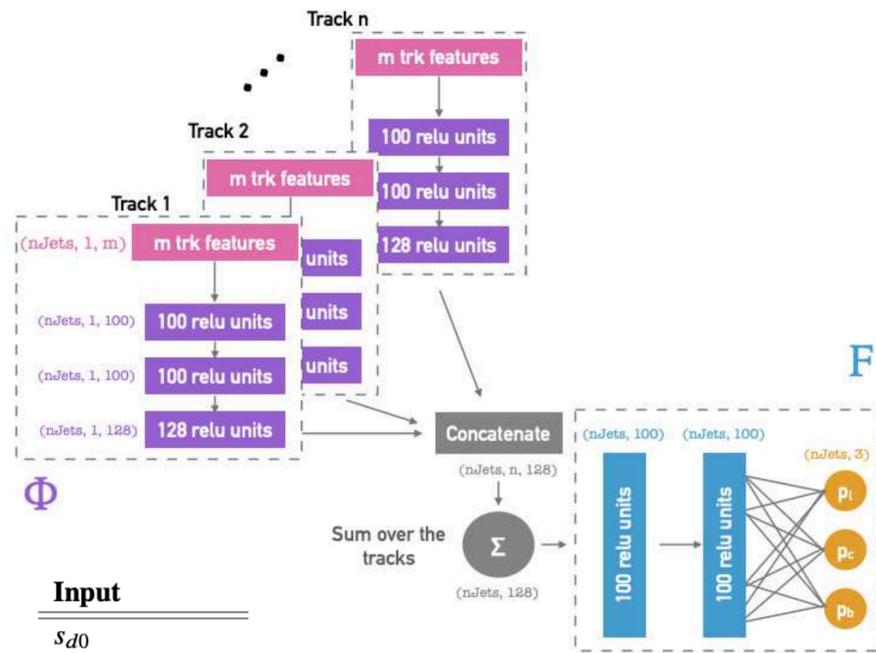
$$\text{EFN}(\{p_1^\mu, \dots, p_M^\mu\}) = F \left(\sum_{i=1}^M z_i \Phi(\hat{p}_i) \right)$$

IRC-safe latent space

Other Physics-Inspired Architectures

Deep Impact Parameter Sets

[ATL-PHYS-PUB-2020-014]



Input

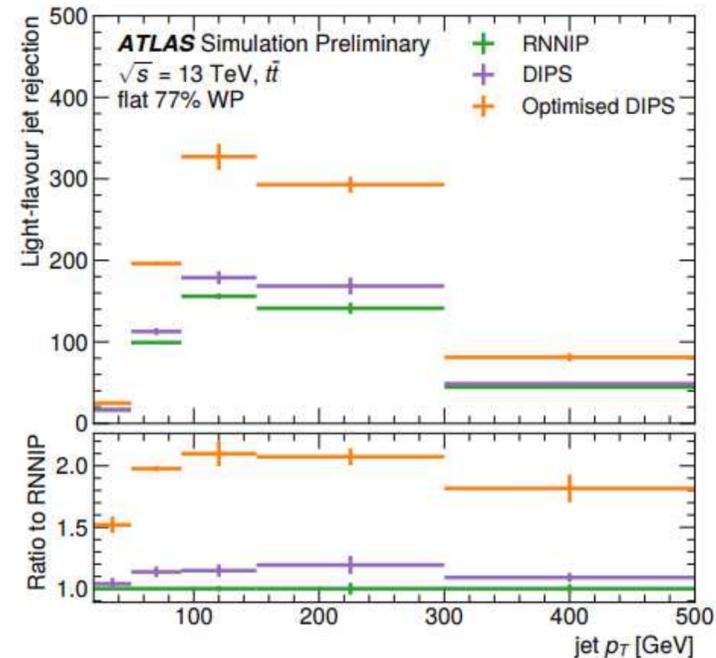
- s_{d0}
- s_{z0}
- $\log p_T^{frac}$
- $\log \Delta R$
- IBL hits
- PIX1 hits
- shared IBL hits
- split IBL hits
- nPixHits
- shared pixel hits
- split pixel hits
- nSCTHits
- shared SCT hits

Achieves 2x better flavor tagging than an RNN

“Optimised” includes additional features per track

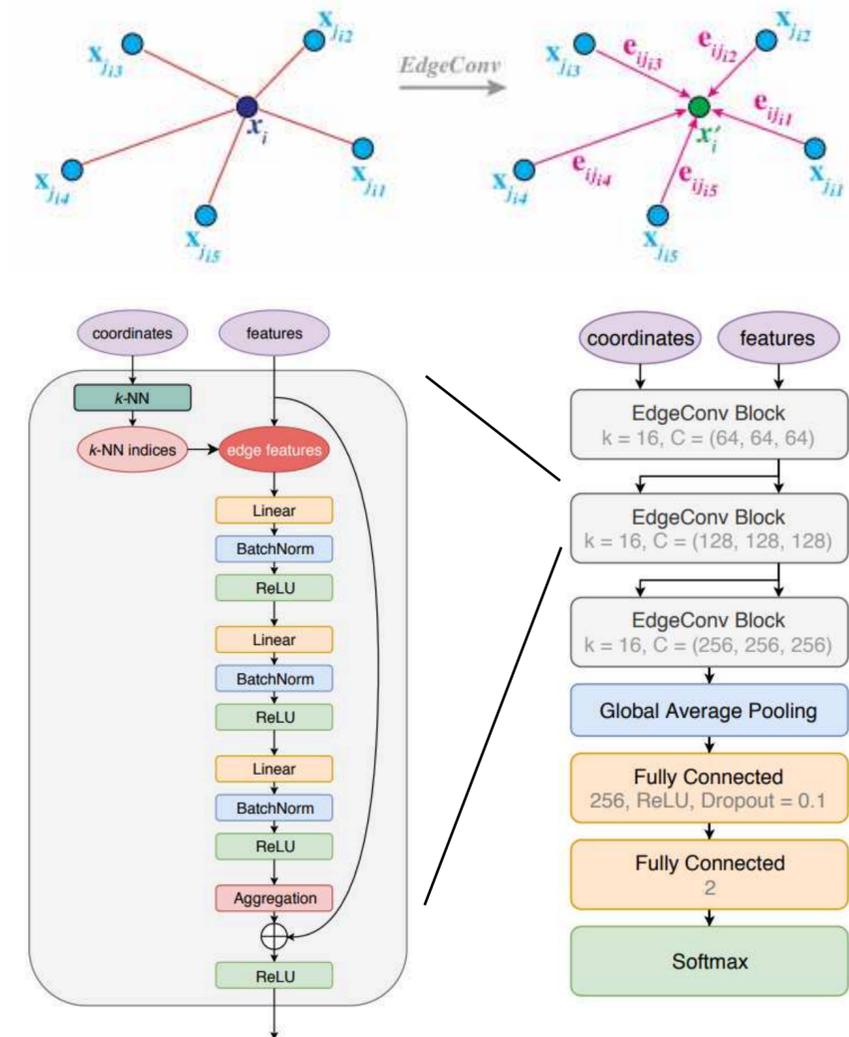
Lorentz Group Equivariant NN

[Bogatskiy, Anderson, Offermann, Roussi, Miller, Kondor, 2006.04780]



Dynamic Graph CNNs (e.g. Particle Net)

[Wang, Sun, Liu, Sarma, Bronstein, Solomon, 1801.07829; Qu, Gouskos, 1902.08570]



Preserves permutation symmetry while prioritizing relationships between inputs

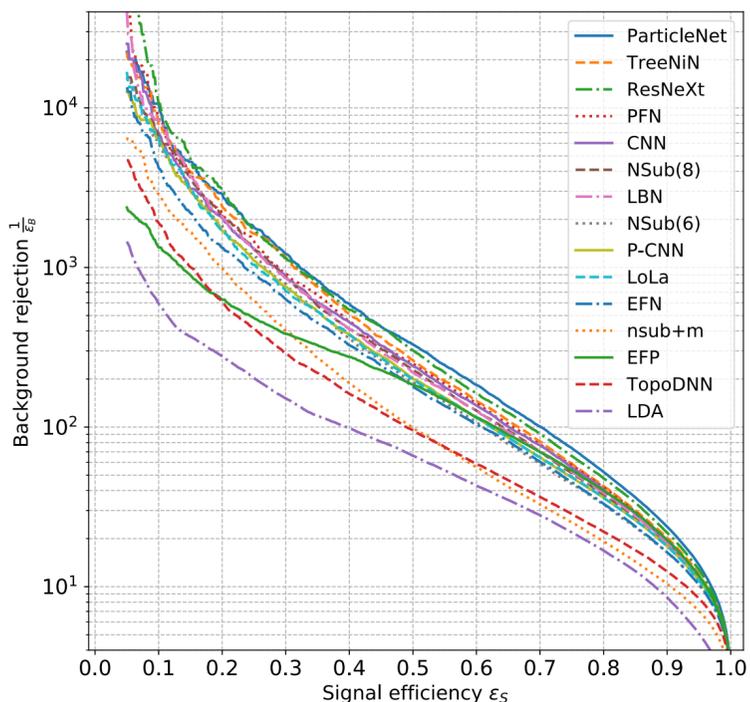
Improved Performance on Classic HEP Tasks with ML

ML optimizes performance

(Jet) Classification

[Kasieczka, Plehn, et al., [1902.09914](#);
many many other references...]

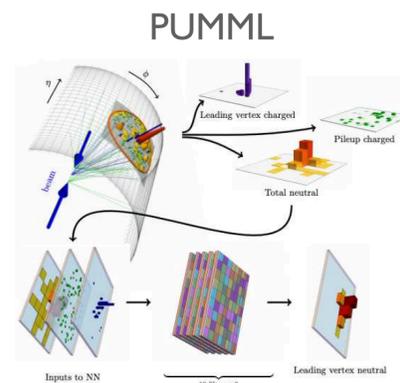
Community top-tagging comparison



Regression (e.g. to remove pileup)

[PTK, Metodiev, Nachman, Schwartz, [1707.08600](#);
[ATL-PHYS-PUB-2019-028](#);

see also Arjona Martínez, Cerri, Spiropulu, Vlimant, Pierini [1810.07988](#)]



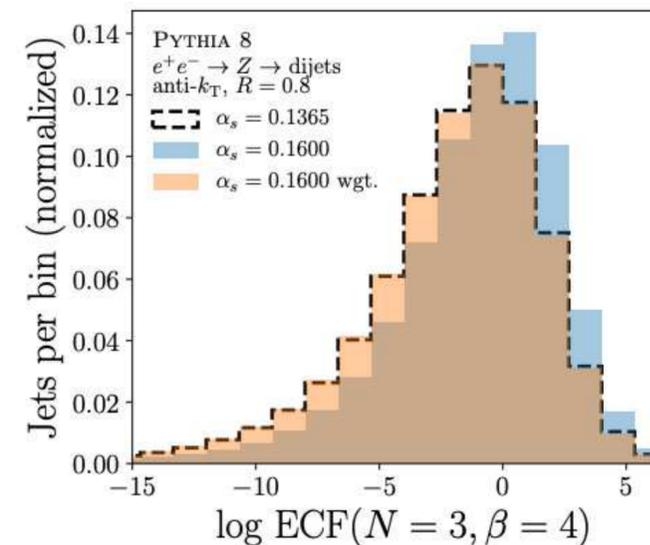
ML can enable an unbinned version of a binned method

(Full) Phase-Space Reweighting

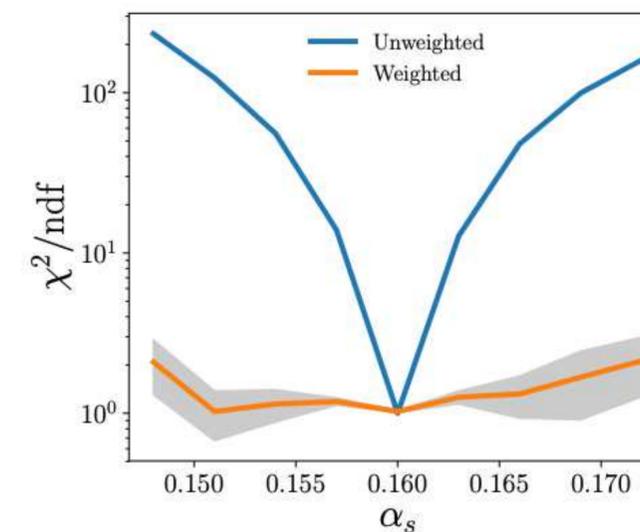
[Cranmer, Pavez, Louppe, [1506.02169](#);
Andreassen, Nachman, [1907.08209](#)]

Likelihood-free inference via classification and Neyman-Pearson

DCTR uses a single high-dimensional reweighting...



to match ECF distributions ...



and multiplicity distributions!

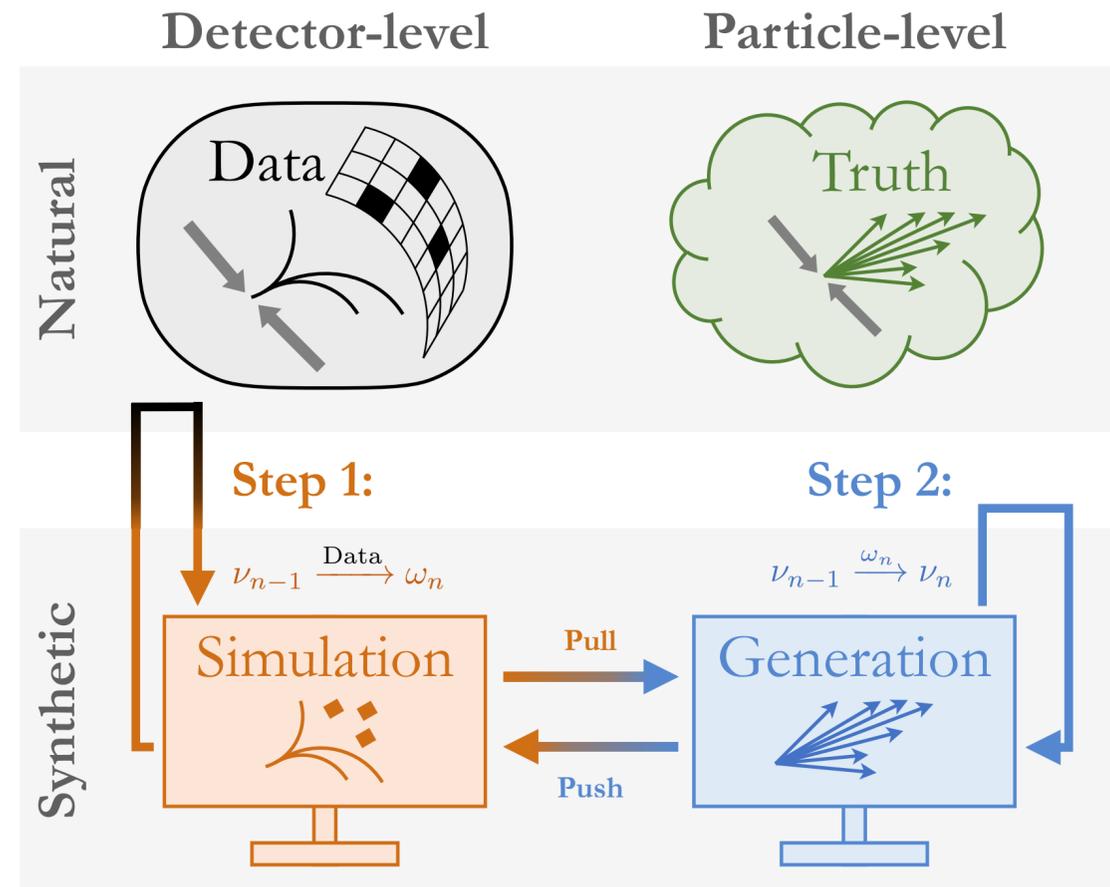
High-dimensional reweighting is useful for many tasks ...

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_{n-1} to data, pulls weights back to particle-level Gen_{n-1}

Step 2 – Reweights Gen_{n-1} to (step 1)-weighted gen_{n-1} , pushes weights to detector-level Sim_n

OmniFold – i.e. continuous IBU

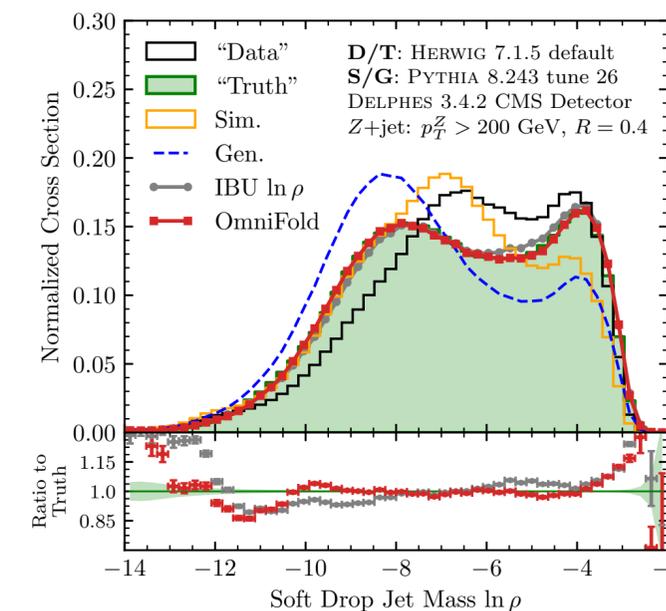
$$\text{Step 1} - \omega_n(m) = \nu_{n-1}^{\text{push}} \times L[(1, \text{Data}), (\nu_{n-1}^{\text{push}}, \text{Sim})](m)$$

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

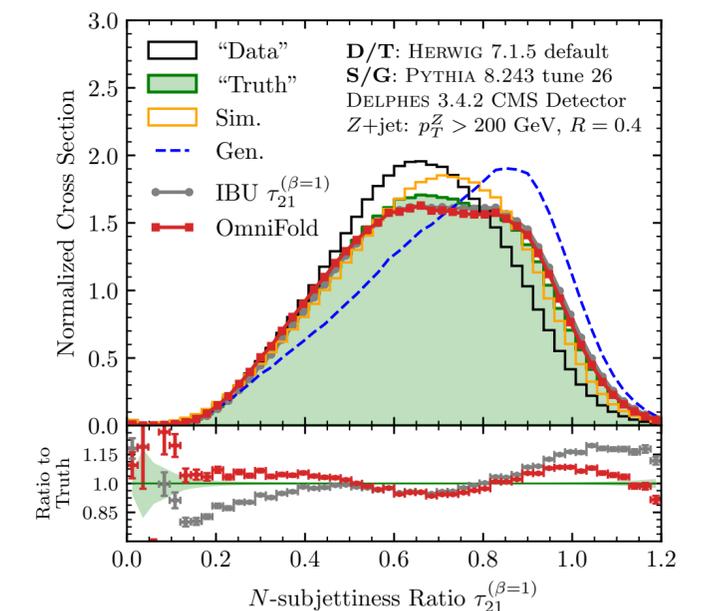
Unfold any* observable $p_{\text{Gen}}(t)$ using universal weights $\nu_n(t)$

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

*Observables should be chosen responsibly



IRC safe



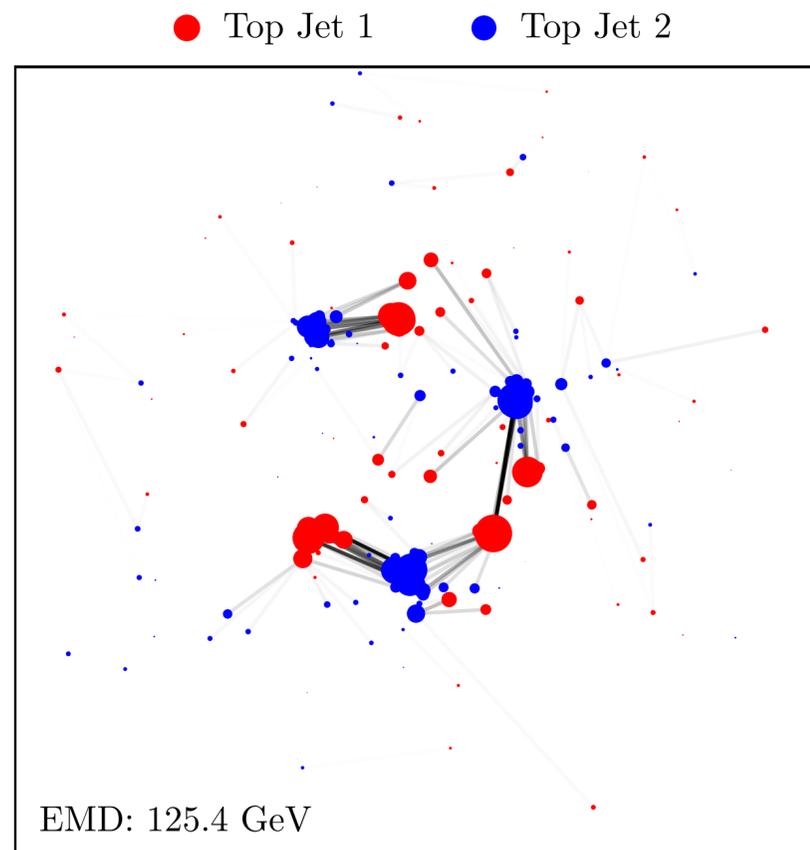
Sudakov safe

Non-Parametric ML – The Energy Mover’s Distance (EMD)

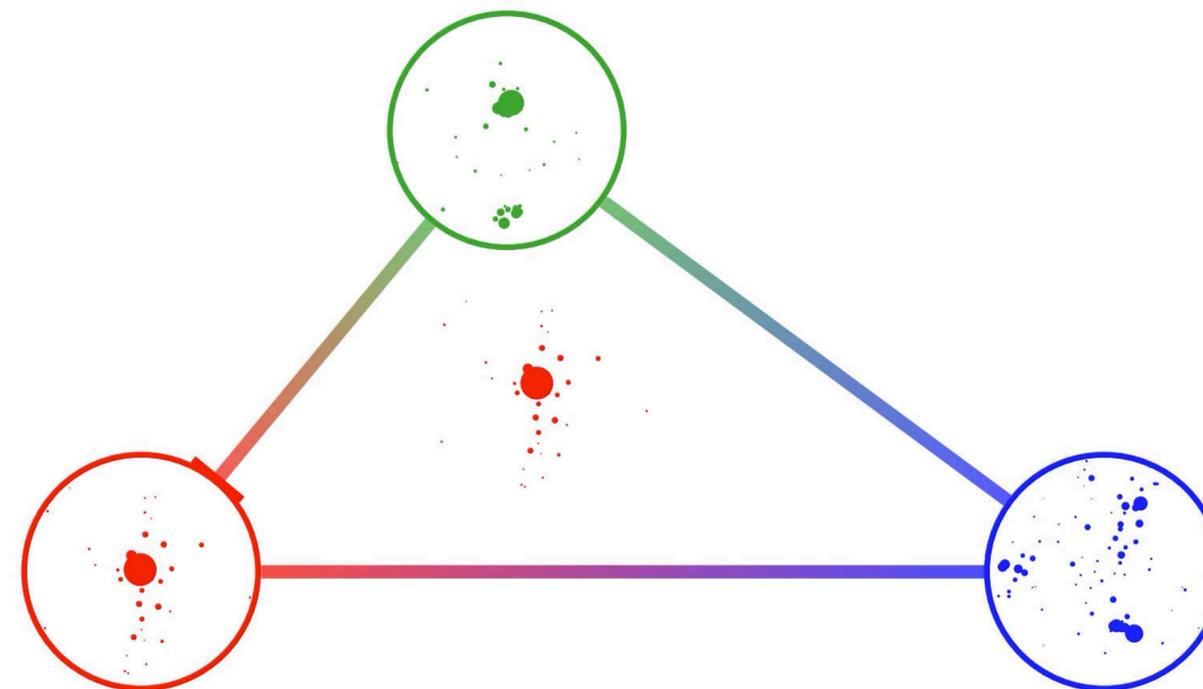
[PTK, Metodiev, Thaler, [PRL 2019](#);

applied on CMS Open Data: PTK, Mastandrea, Metodiev, Naik, Thaler, [1908.08542](#)]

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

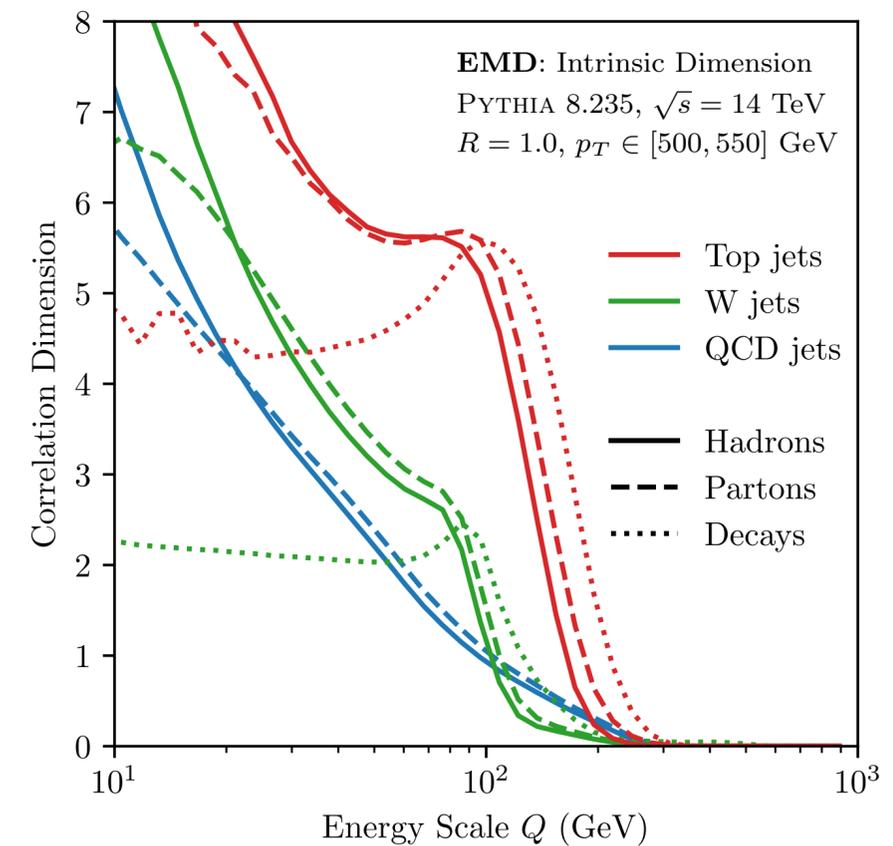


EMD is the work required to rearrange one event into another



Triangle inequality

$$0 \leq \text{EMD}(\mathcal{E}, \mathcal{E}') \leq \text{EMD}(\mathcal{E}, \mathcal{E}'') + \text{EMD}(\mathcal{E}'', \mathcal{E}')$$

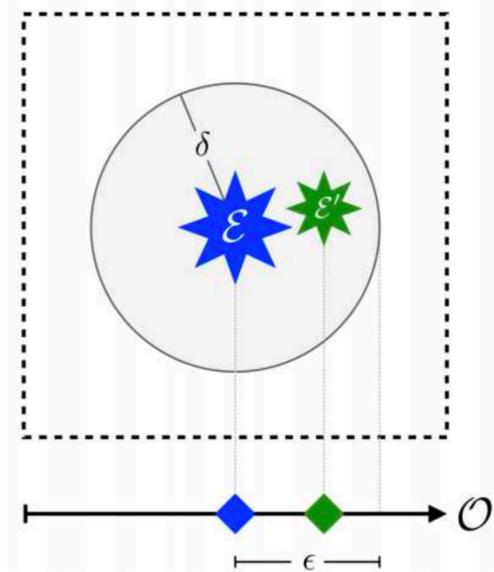


Intrinsic dimensionality of dataset highlights physics at all scales

ML Enables New Theoretical Paradigms – Six Decades of Collider Techniques as Geometry

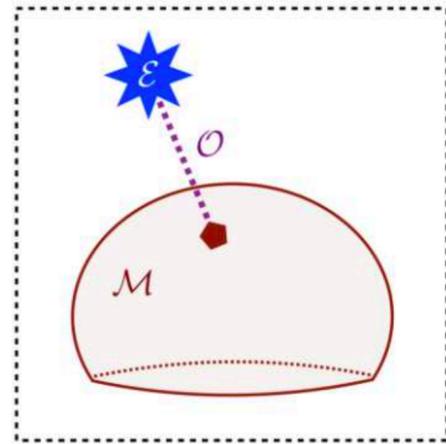
[PTK, Metodiev, Thaler, 2004.04159]

IRC Safety is smoothness in the space of events



Taming infinities

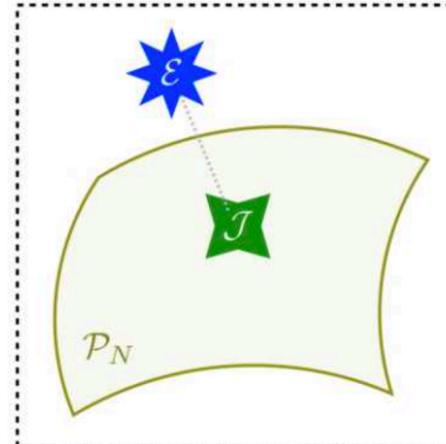
Event shapes are distances from events to manifolds.



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

Event Shapes

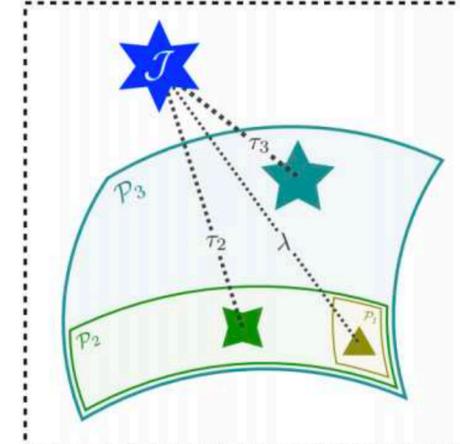
Jets are projections to few-particle manifolds.



$$J = \operatorname{argmin}_{\mathcal{E}' \in \mathcal{P}_N} \text{EMD}_{\beta, R}(\mathcal{E}, \mathcal{E}')$$

Jet Algorithms

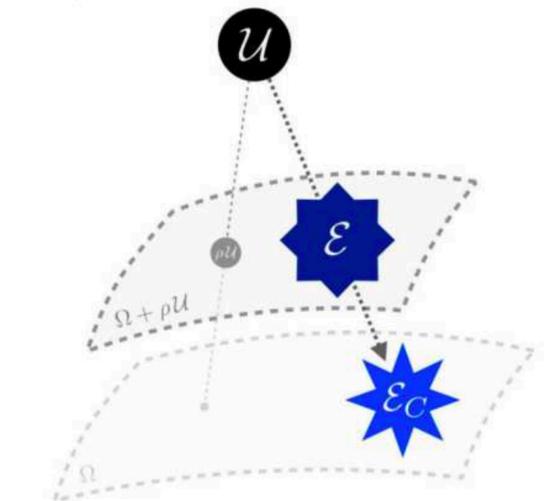
Substructure resolves emissions within the jet.



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

Jet Substructure

Pileup mitigation moves away from uniform radiation.



$$\mathcal{E}_C = \operatorname{argmin}_{\mathcal{E}'} \text{EMD}(\mathcal{E}, \mathcal{E}' + \rho \mathcal{U}).$$

Pileup



Handling Uncertainties

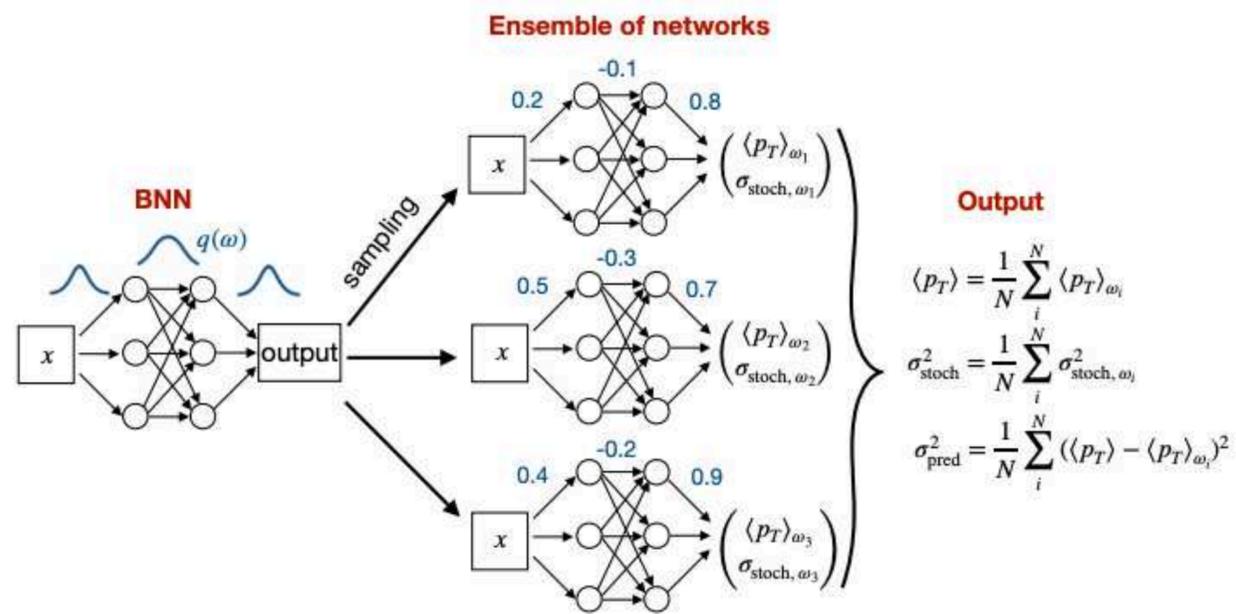
ML can assist in handling uncertainties (but can also create its own difficulties)

Adversarial training to reduce dependence on uncertain regions

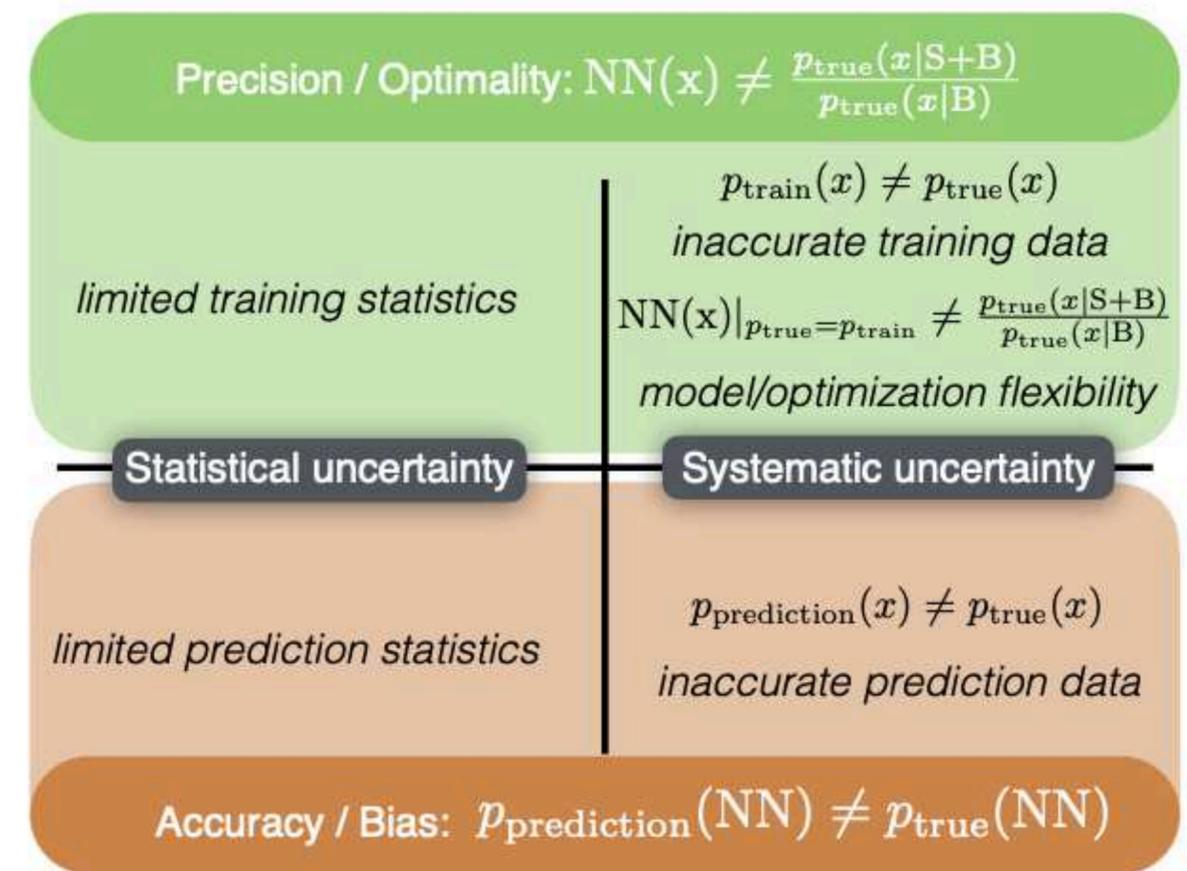
[Englert, Galler, Harris, Spannowsky, [1807.08763](#)]

Bayesian neural networks can estimate some uncertainties

[Bollweg, Haußmann, Kasieczka, Luchmann, Plehn, Thompson, [1904.10004](#);
Kasieczka, Luchmann, Otterpohl, Plehn, [2003.11099](#)]



Sources of uncertainty in a statistical analysis



[Nachman, [1909.03081](#)]

Parametrized models could enable efficient profiling to handle systematic uncertainties

[similar to Baldi, Cranmer, Faucett, Sadowski, Whiteson, [1601.07913](#)]

... And So Much More HEP-Specific ML

Weakly Supervised Learning

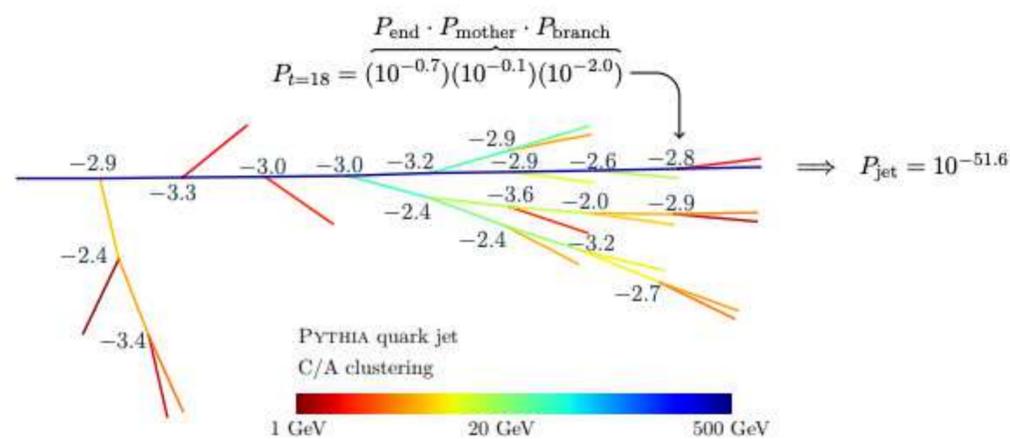
e.g. LLP, CWoLa

[Dery, Nachman, Rubbo, Schwartzman, [1702.00414](#);
 Metodiev, Nachman, Thaler, [1708.02949](#);
 Cohen, Freytsis, Ostdiek, [1706.09451](#);
 PTK, Metodiev, Nachman, Schwartz, [1801.10158](#)]

Unsupervised Learning

e.g. JUNIPR

[Andreassen, Feige, Frye, Schwartz, [1804.09720](#), [1906.10137](#)]

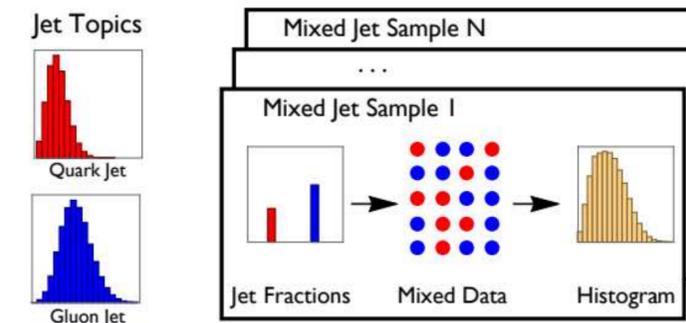


Probabilistic formation of a Pythia jet

Topic Modeling

e.g. Jet Topics

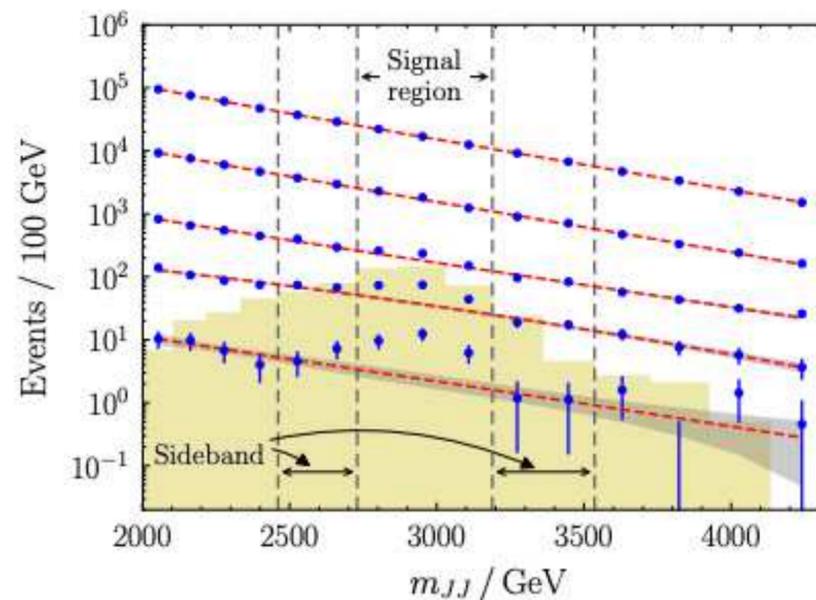
[Metodiev, Thaler, [1802.00008](#);
 PTK, Metodiev, Thaler, [1809.01140](#)]



Anomaly Detection

tons of recent work

[Collins, Howe, Nachman, [1805.02664](#);
 Farina, Nakai, Shih, [1808.08992](#);
 Heibel, Kasieczka, Plehn, Thompson, [1808.08979](#);
 Cerri, Nguyen, Pierini, Spiropulu, Vlimant, [1811.10276](#);
 Blance, Spannowsky, Waite, [1905.10384](#);
 See LHC Olympics 2020 anomaly detection workshop]

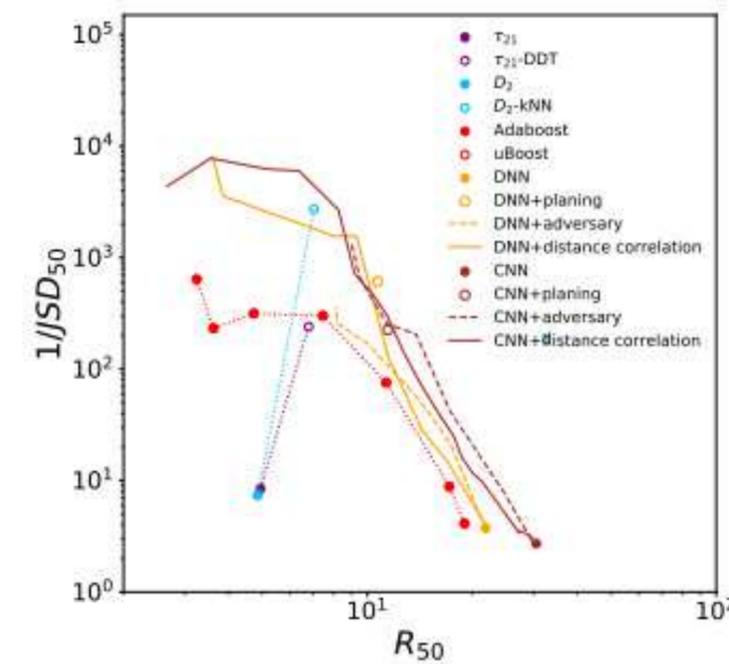


CWoLa hunting to enhance a resonance search

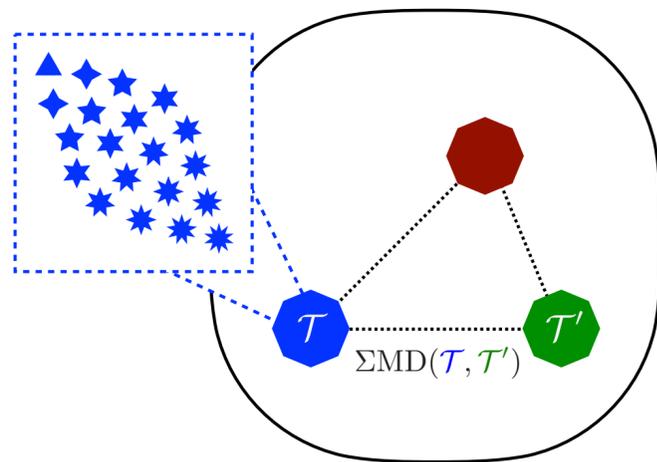
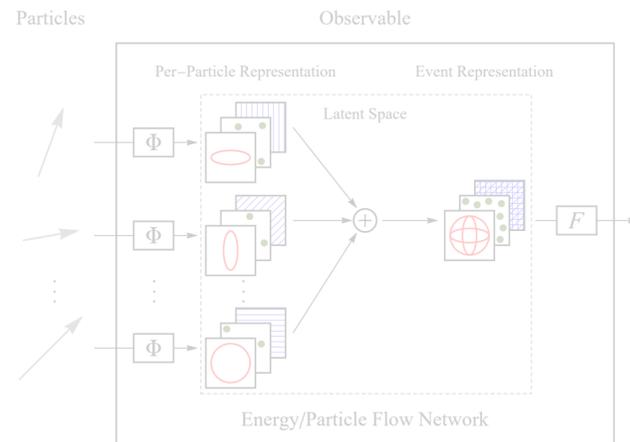
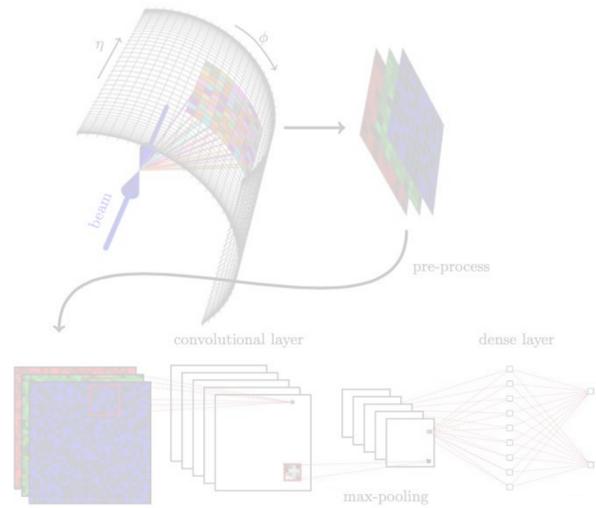
Decorrelation Methods

e.g. DDT, Planing, DisCo

[Dolen, Harris, Marzani, Rappoccio, Tran, [1603.00027](#);
 Chang, Cohen, Ostdiek, [1709.10106](#);
 Kasieczka, Shih, [2001.05310](#);
 Kasieczka, Nachman, Schwartz, Shih, [2007.14400](#)]



Decorrelating ML taggers from a key observable



Ubiquity of ML in HEP

Lightning Review

Future Directions

Thoughts on the Future

▶ *Machine learning will be essential in maximizing HEP potential*

We should capitalize on the opportunity to optimize

ML is both a computational tool and a useful formalism/language

▶ *Diversity of applications is impressive and will become more so*

Rapidly advancing beyond “hammer and nail” approach for ML in HEP

New observables, architectures, paradigms have been developed

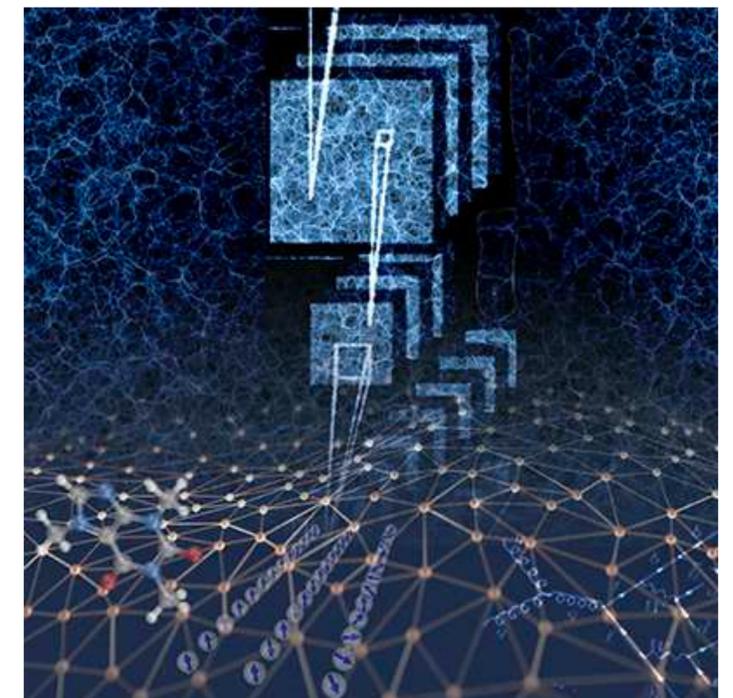
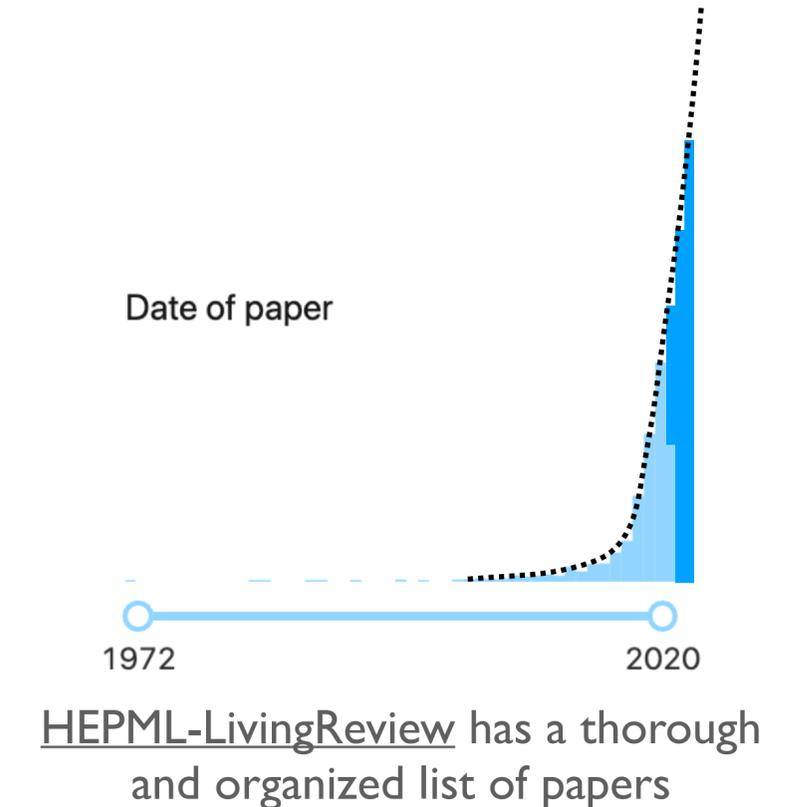
▶ *Collaboration across traditional lines will enable success*

Learn from and contribute to the highly-vibrant ML community

ML in HEP needs insight from theory and experiment

▶ *ML strongly overlaps with other computational frontier areas*

Software workflows, reproducible analyses, public datasets are critical



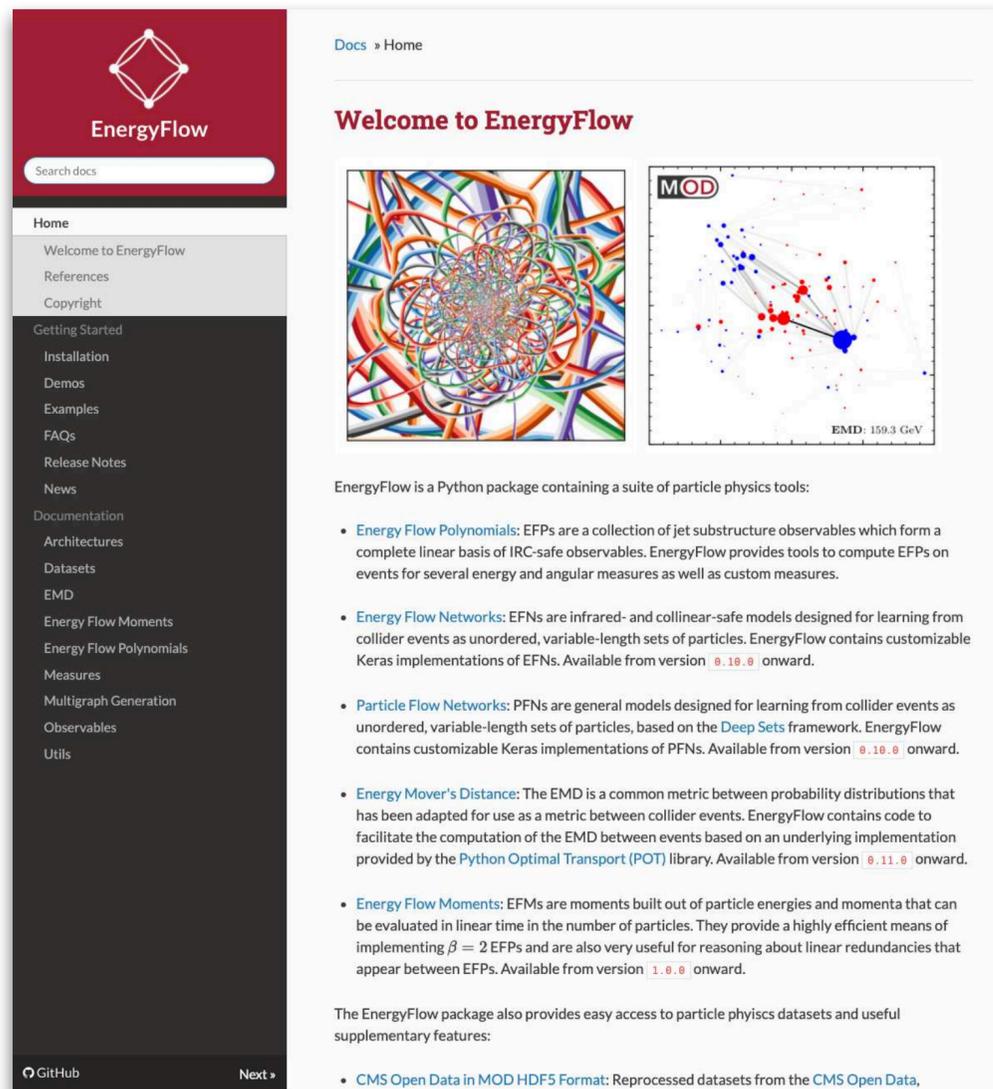
[Reviews of Modern Physics Cover December 2019
from Machine learning and the physical sciences]

Thank You!

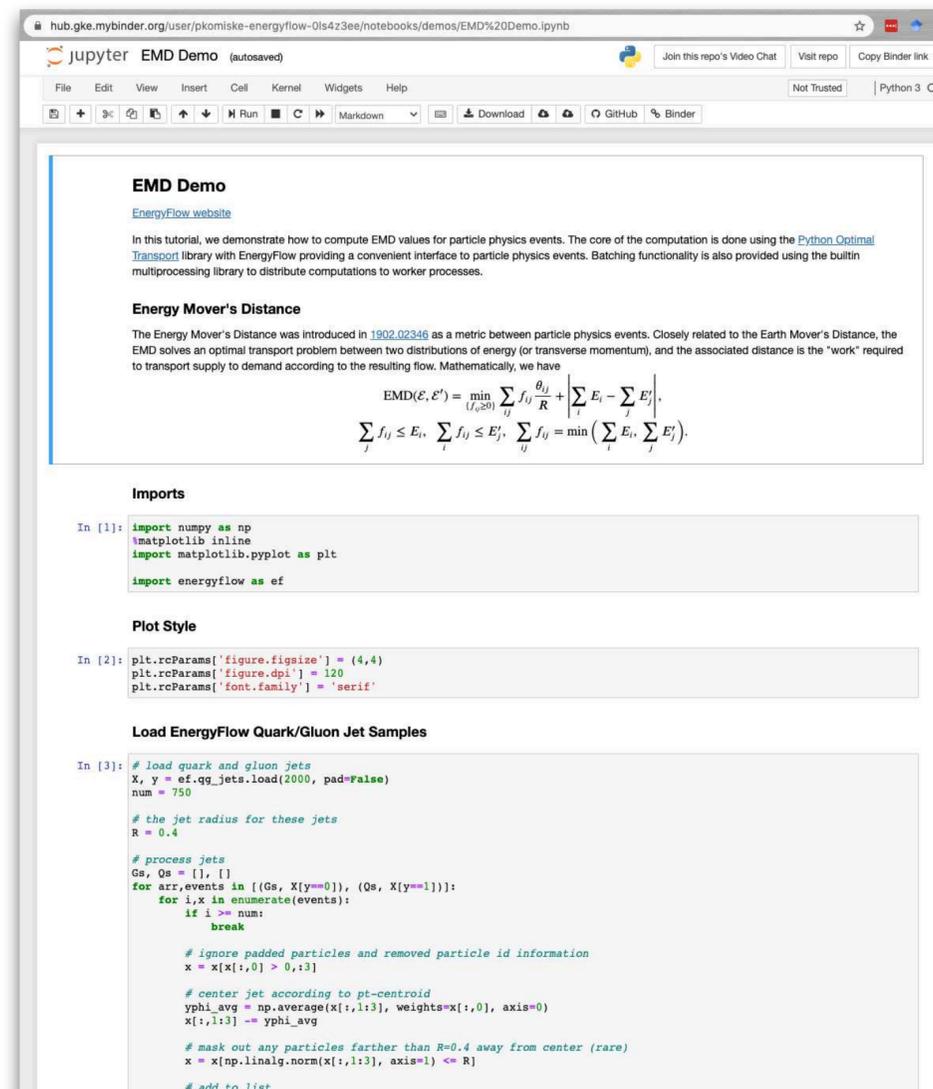
EnergyFlow Python Package

`pip3 install energyflow`

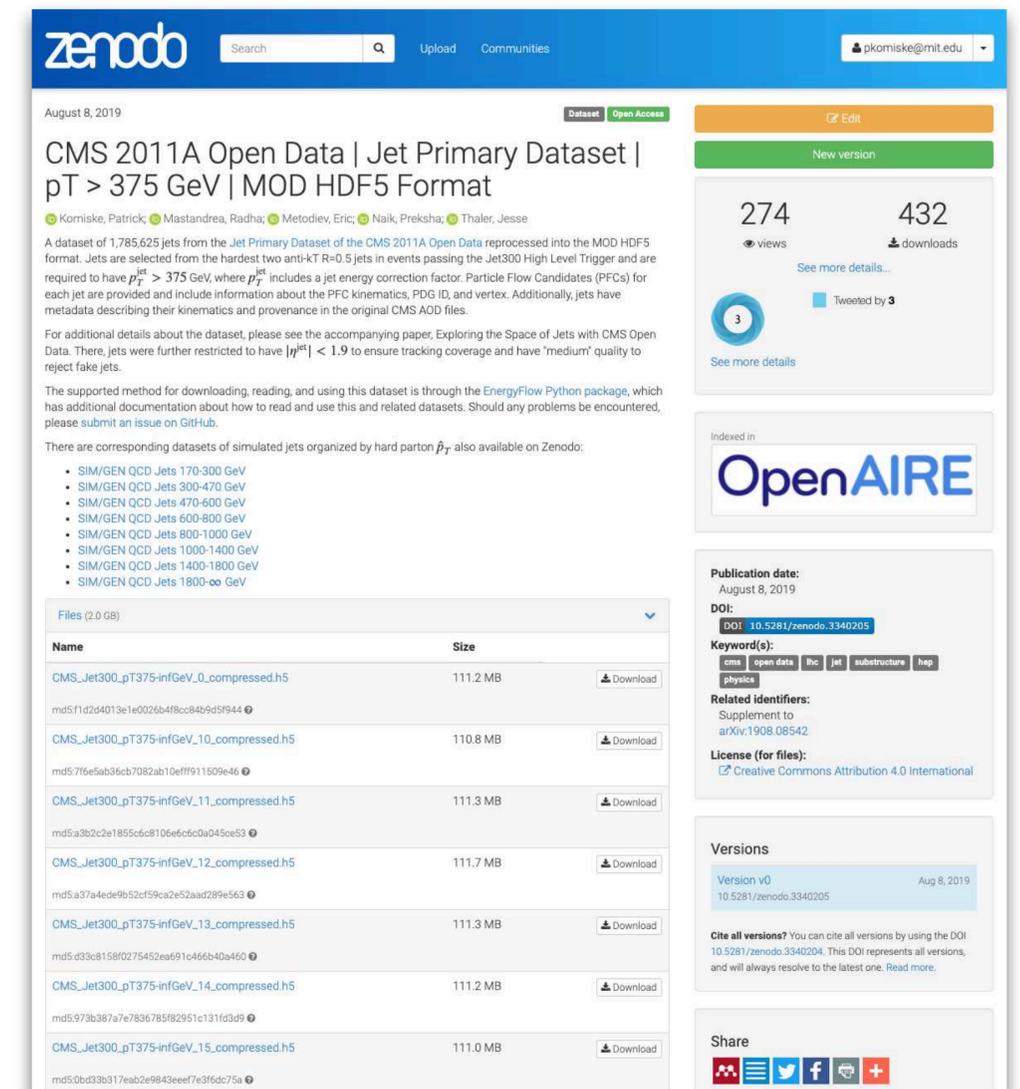
Keras/Tensorflow implementations of Energy/Particle Flow Networks
Interfaces with reprocessed [CMS 2011A Jet Primary Dataset](#) hosted on [Zenodo](#)
Detailed [examples](#), [demos](#), and [documentation](#)



The screenshot shows the EnergyFlow website. The header features the EnergyFlow logo and a search bar. A sidebar on the left contains navigation links: Home, Welcome to EnergyFlow, References, Copyright, 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 is titled "Welcome to EnergyFlow" and includes two images: a colorful particle flow network and a scatter plot labeled "MOD" with "EMD: 159.3 GeV". Below the images, a paragraph states "EnergyFlow is a Python package containing a suite of particle physics tools:". A list of features follows, including Energy Flow Polynomials, Energy Flow Networks, Particle Flow Networks, Energy Mover's Distance, and Energy Flow Moments. At the bottom, it mentions "CMS Open Data in MOD HDF5 Format: Reprocessed datasets from the CMS Open Data."



The screenshot shows a Jupyter Notebook titled "EMD Demo" in a browser window. The notebook content includes a title "EMD Demo" with a link to the EnergyFlow website. The text explains that the tutorial demonstrates how to compute EMD values for particle physics events using the Python Optimal Transport library. It defines the Energy Mover's Distance and provides the mathematical formula:
$$\text{EMD}(E, E') = \min_{(f_{ij})} \sum_{ij} f_{ij} \frac{\theta_{ij}}{R} + \left| \sum_i E_i - \sum_j E'_j \right|$$
 with constraints $f_{ij} \leq E_i$, $f_{ij} \leq E'_j$, and $\sum_{ij} f_{ij} = \min(\sum_i E_i, \sum_j E'_j)$. The notebook also shows code for imports, plot style settings, and loading quark and gluon jets. The code includes comments for jet radius, processing jets, ignoring padded particles, and centering jets.



The screenshot shows the Zenodo dataset page for "CMS 2011A Open Data | Jet Primary Dataset | pT > 375 GeV | MOD HDF5 Format". The page includes a search bar, user profile, and dataset statistics: 274 views and 432 downloads. It lists the authors: Komiske, Patrick; Mastandrea, Radha; Metodiev, Eric; Naik, Preksha; Thaler, Jesse. A list of 15 files is provided, each with a name, size, and download link. The files are organized by hard parton pT. The page also features OpenAIRE indexing, publication date (August 8, 2019), DOI (10.5281/zenodo.3340205), and related identifiers.

Iterated Bayesian Unfolding (IBU)

Histogram-based unfolding method for a small number of observables

Choose observable(s) and binning at **detector-level** and **particle-level**

measured distribution: $m_i = \text{Pr}(\text{measure } i)$ true distribution: $t_j^{(0)} = \text{Pr}(\text{truth is } j)$

Calculate *response matrix* R_{ij} from **generated/simulated** pairs of events

$R_{ij} = \text{Pr}(\text{measure } i \mid \text{truth is } j)$ R is in general non-square and non-invertible

Calculate new particle-level distribution using Bayes' theorem

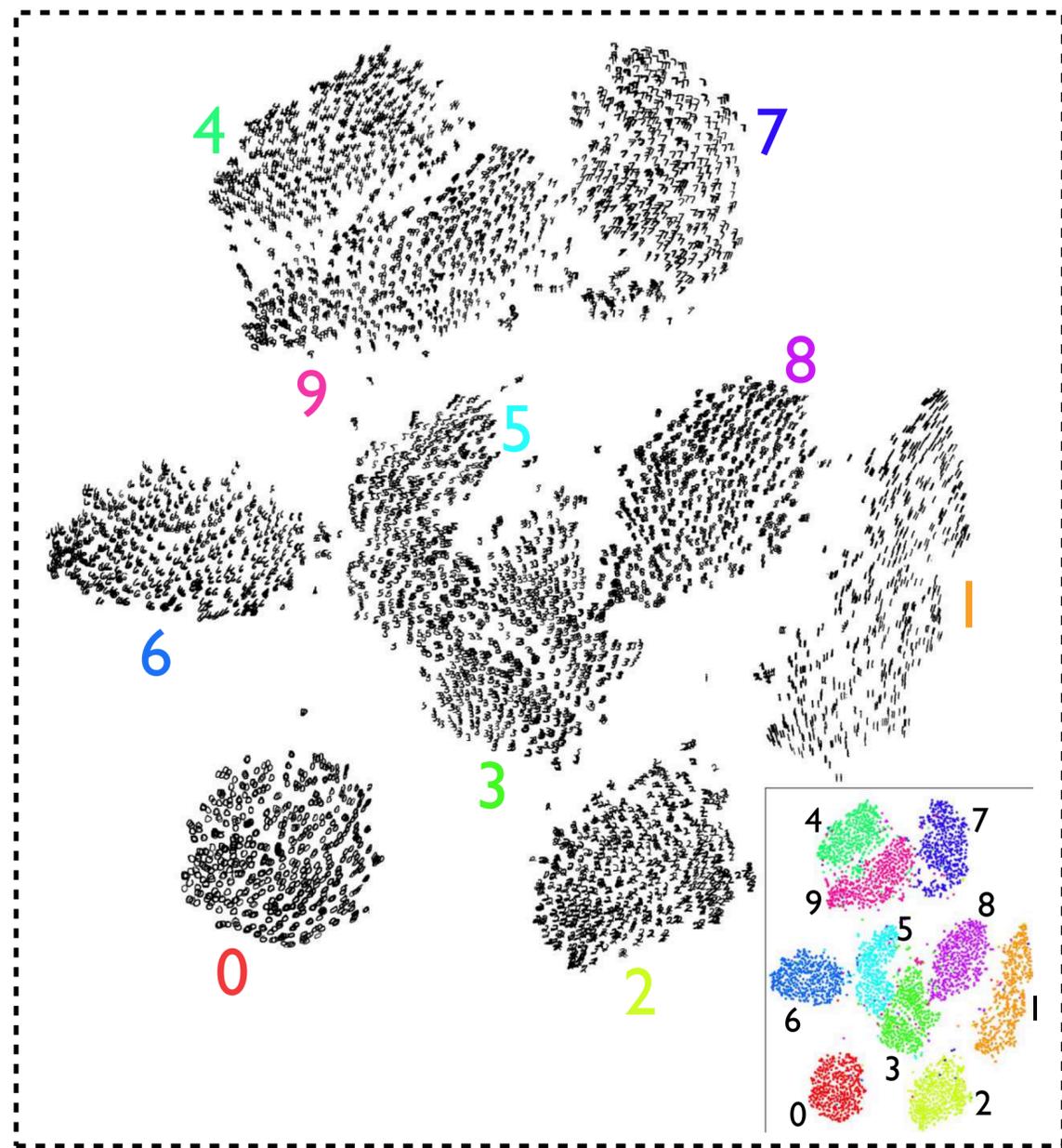
$$t_j^{(n)} = \sum_i \text{Pr}(\text{truth}_{n-1} \text{ is } j \mid \text{measure } i) \times \text{Pr}(\text{measure } i) = \sum_i \frac{R_{ij} t_j^{(n-1)}}{\sum_k R_{ik} t_k^{(n-1)}} \times m_i$$

Iterate procedure to remove dependence on prior

[Richardson, 1972; Lucy, 1974; D'Agostini, 1995]

Visualizing Geometry in the Space of Events

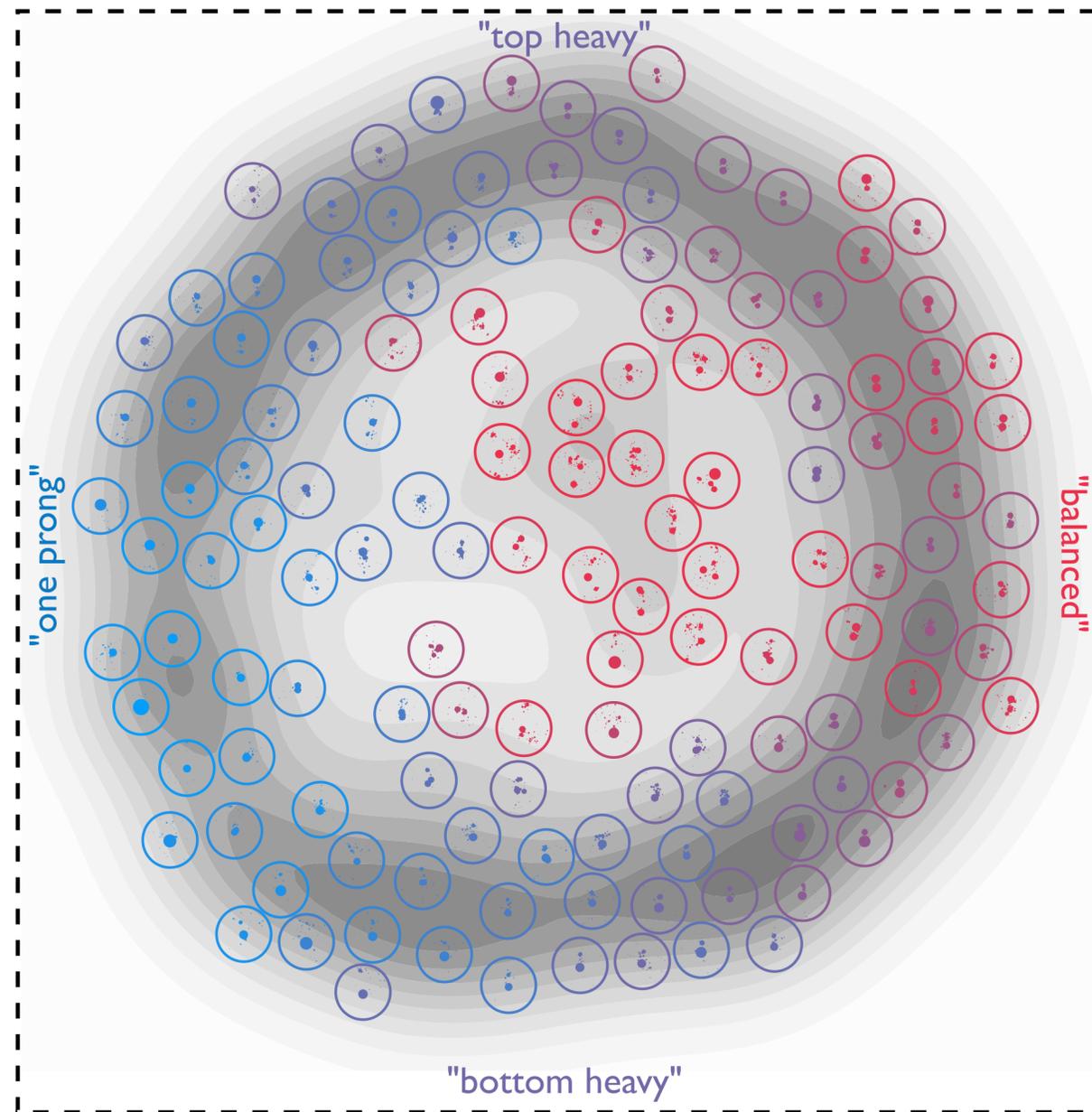
t-Distributed Stochastic Neighbor Embedding (t-SNE)
MNIST handwritten digits



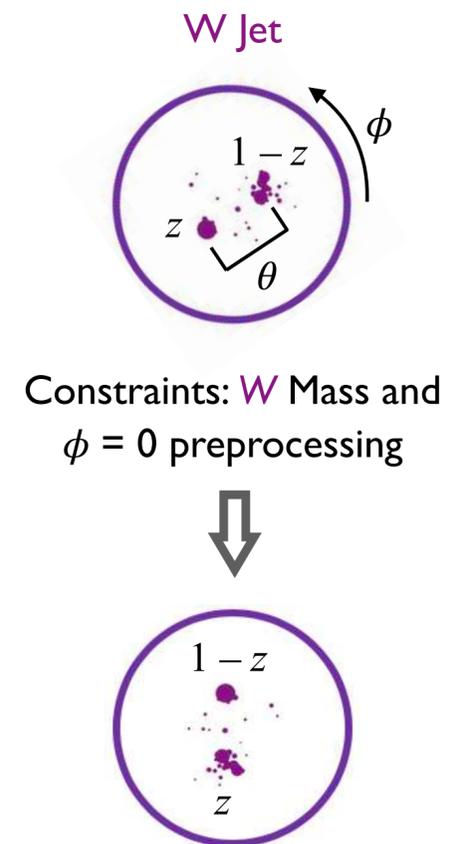
[L. van der Maaten, G. Hinton, JMLR 2008]

[PTK, Metodiev, Thaler, PRL 2019]

Geometric space of W jets

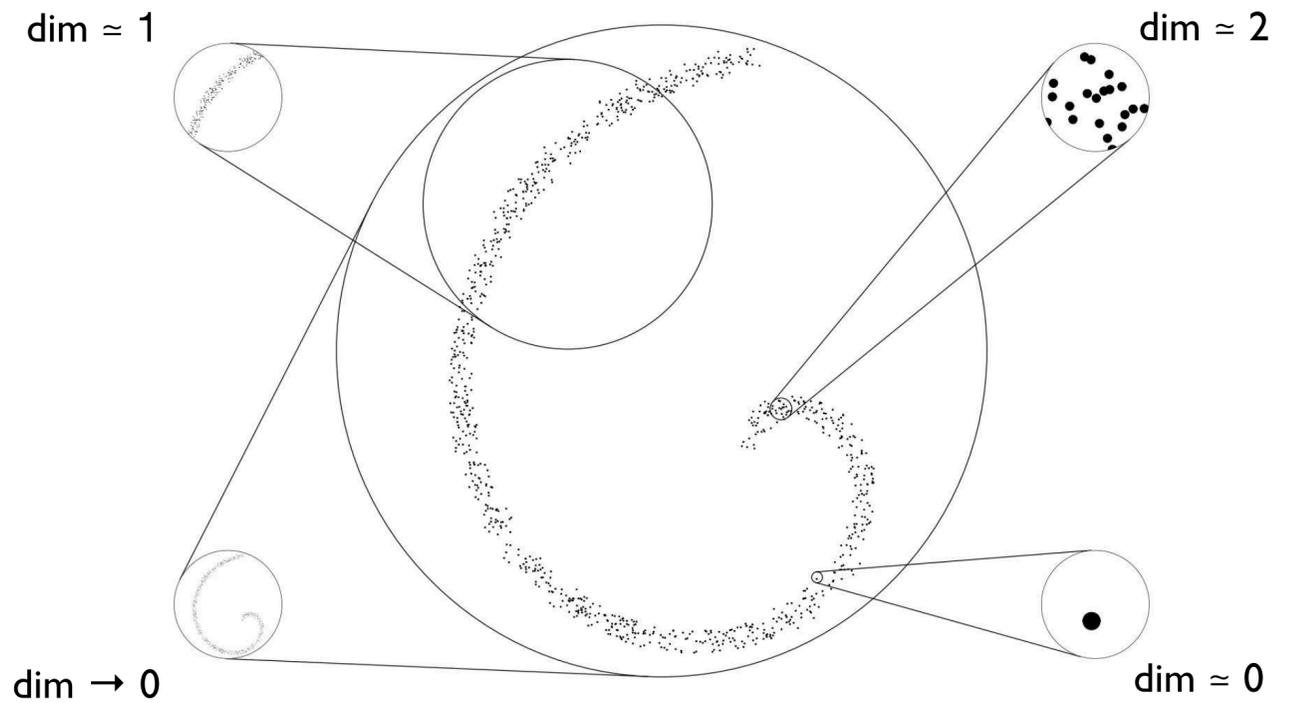


Gray contours represent the density of jets
Each circle is a particular W jet



Quantifying Event-Space Manifolds

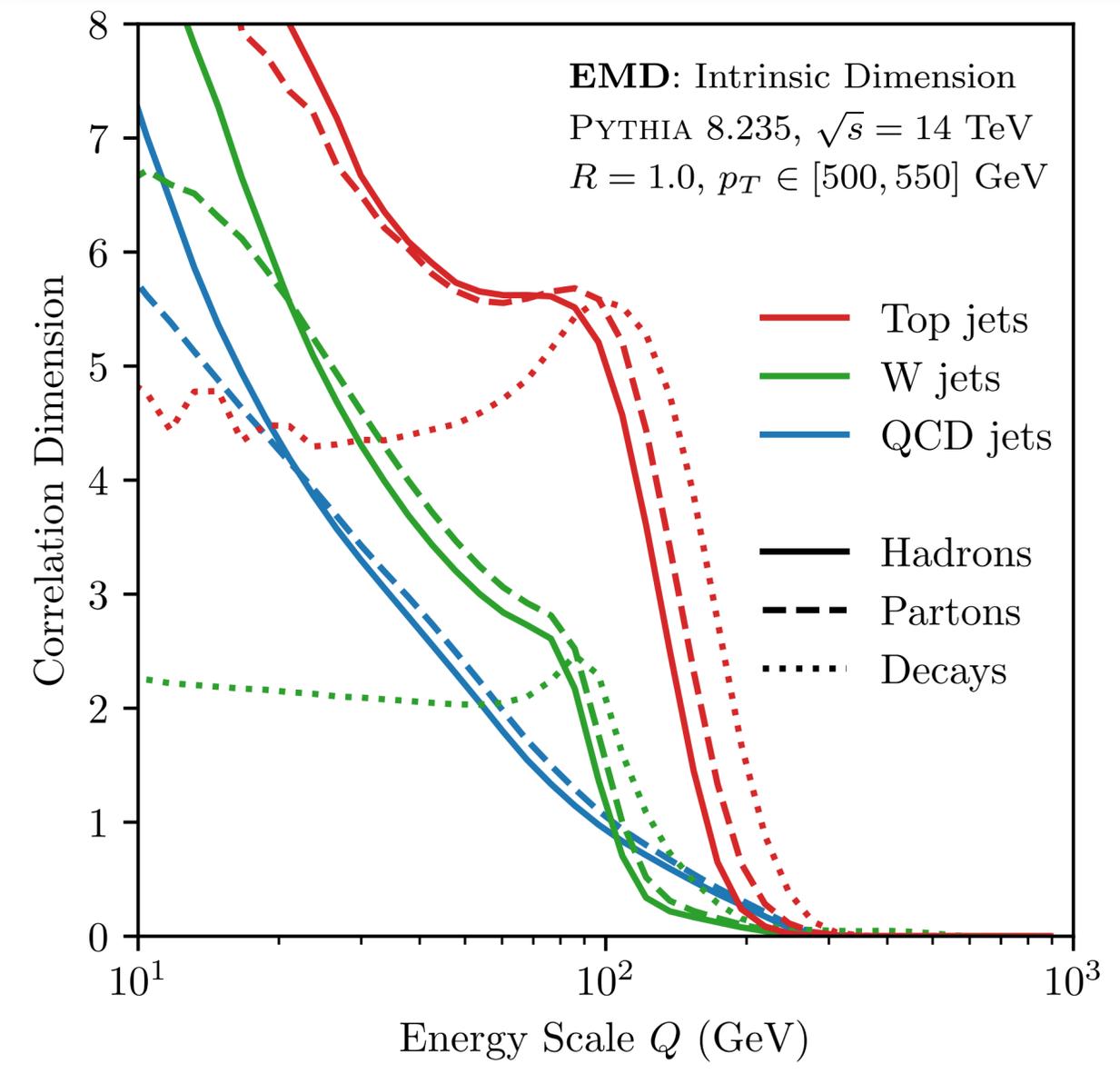
Correlation dimension: *how does the # of elements within a ball of size Q change?*



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

- 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

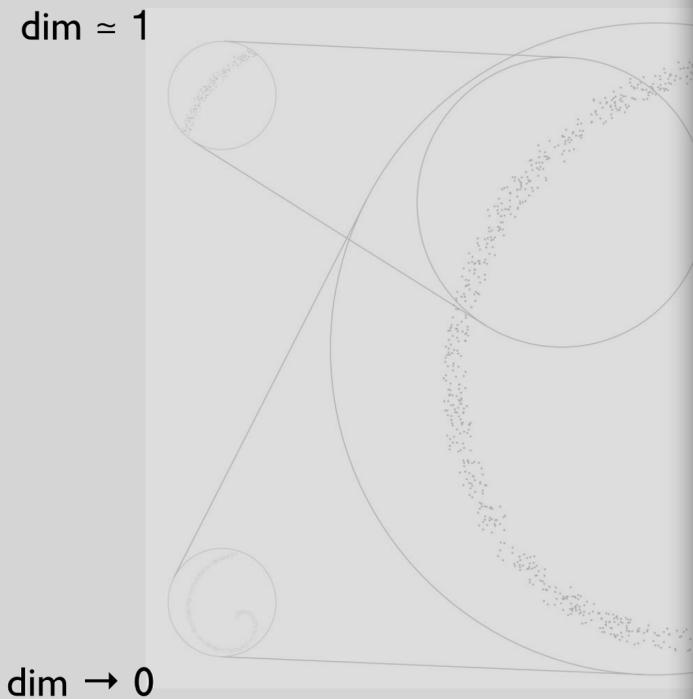
$$\text{dim}(Q) = Q \frac{\partial}{\partial Q} \ln \sum_i \sum_j \Theta(\text{EMD}(\mathcal{E}_i, \mathcal{E}'_j) < Q)$$



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

Quantifying Event-Space Manifolds

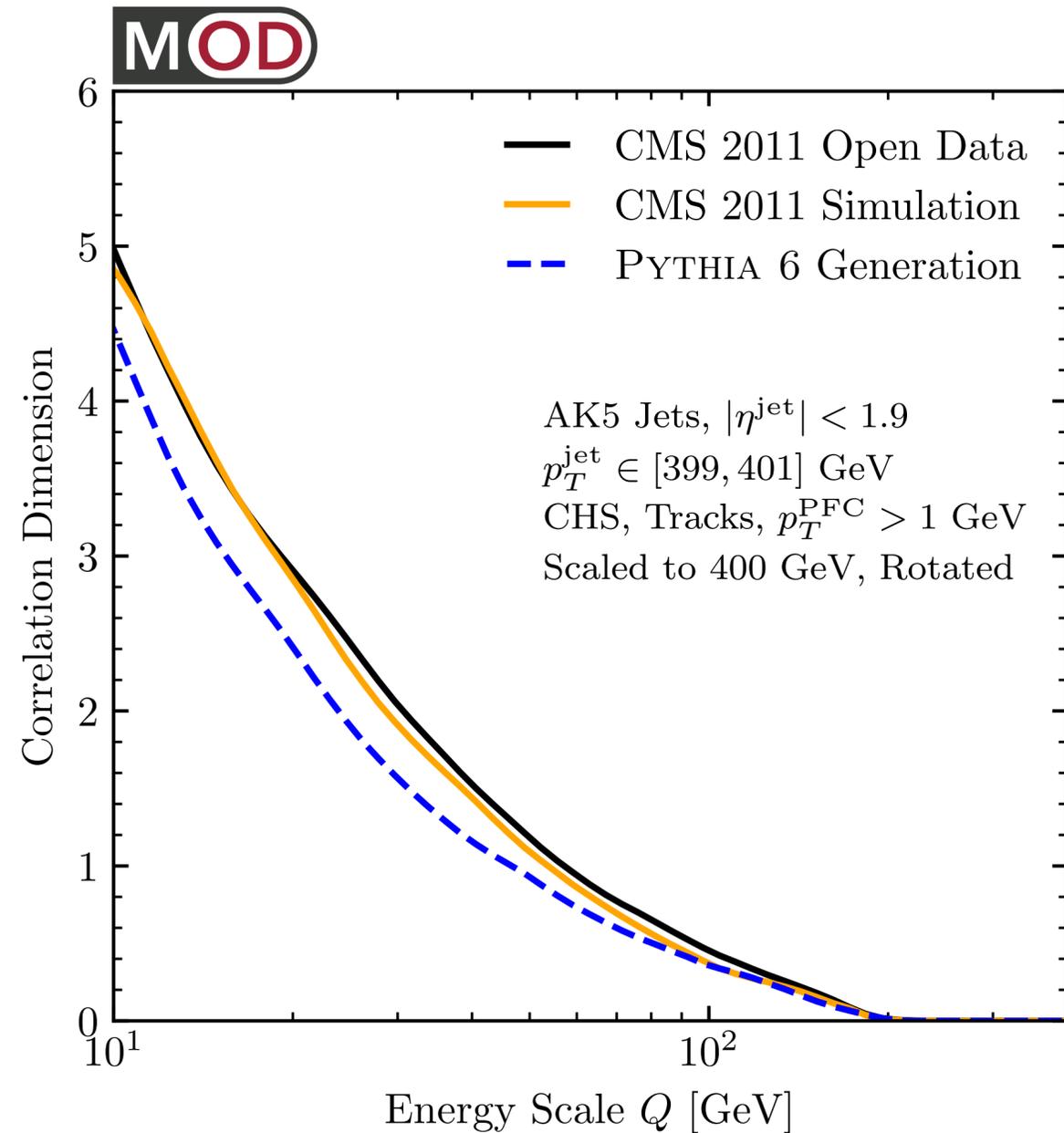
Correlation dimension
elements within a ball



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

Correlation dimension
Decays are "constant"
Complexity hierarchy
Fragmentation increases
Hadronization important

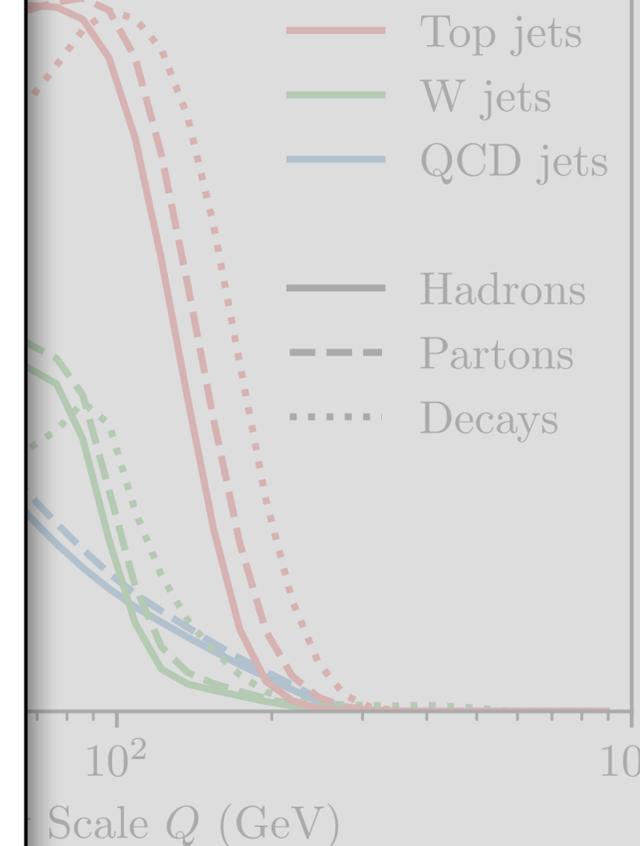
... in CMS Open Data



*More in backup

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

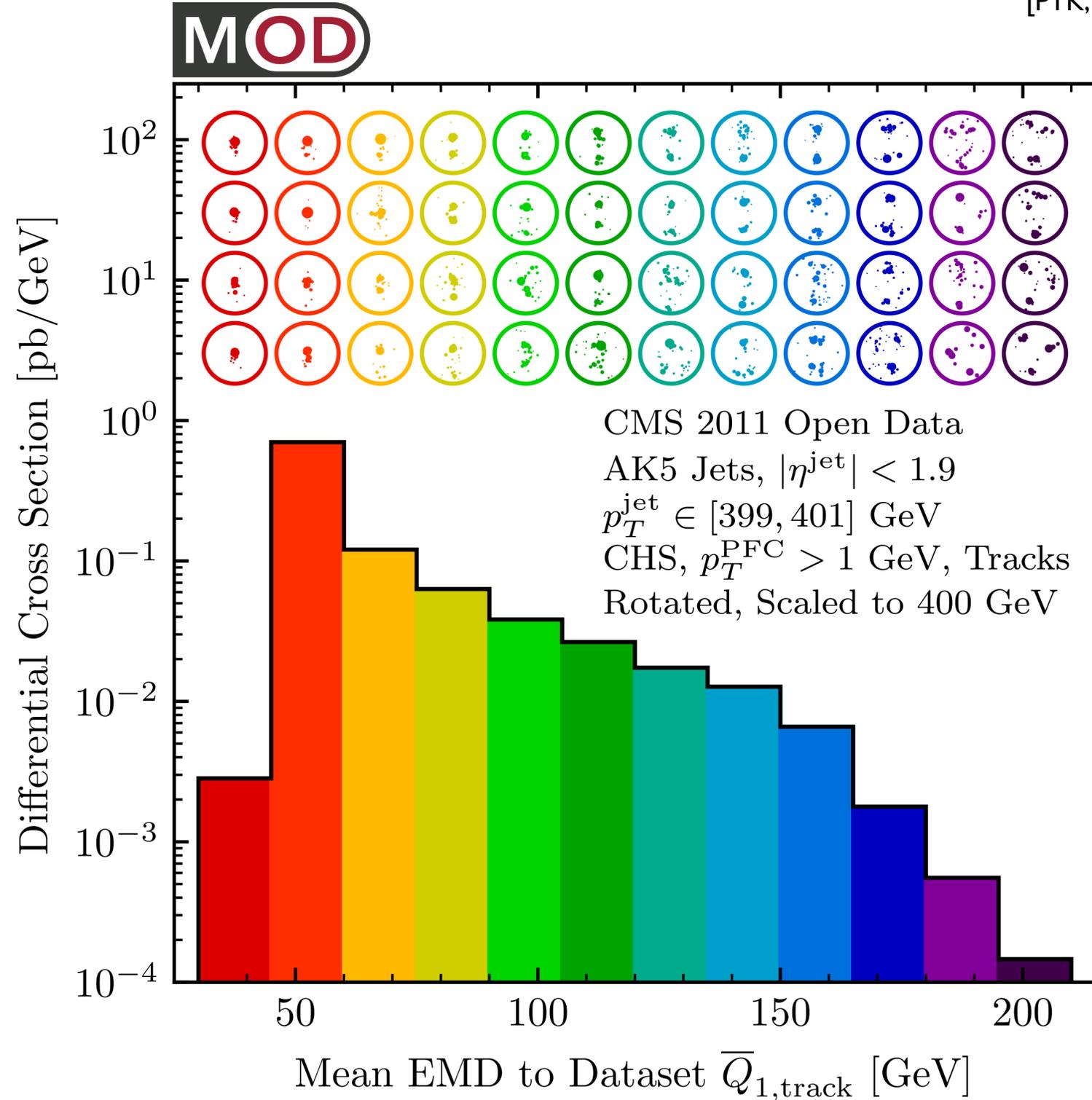
EMD: Intrinsic Dimension
PYTHIA 8.235, $\sqrt{s} = 14$ TeV
 $R = 1.0, p_T \in [500, 550]$ GeV



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

Visualizing Geometry in CMS Open Data

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



EMD for anomaly detection

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

n^{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?

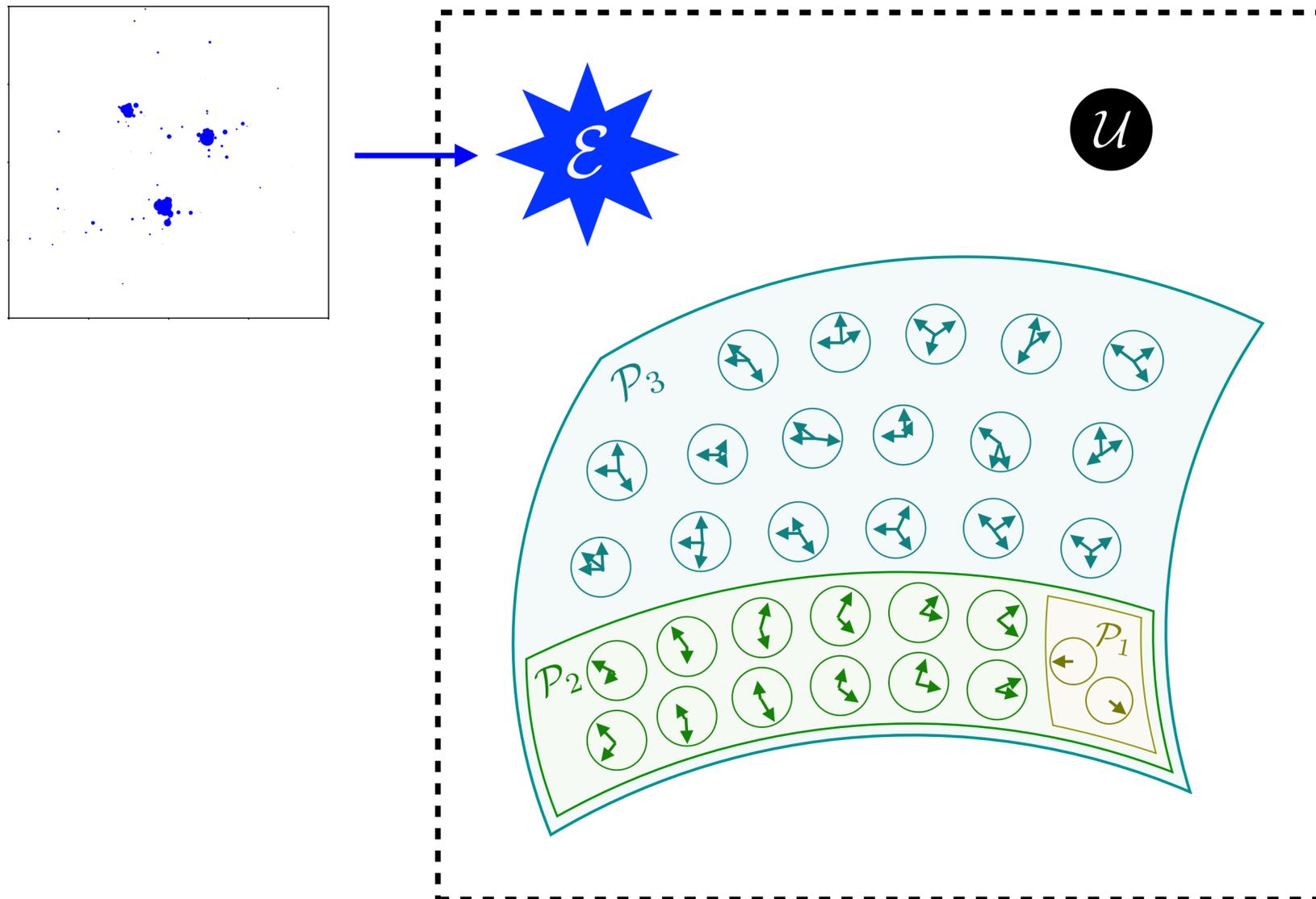
jet from dataset

medoids

N-particle Manifolds in the Space of Events

[PTK, Metodiev, Thaler, 2004.04159]

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



\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

⋮

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

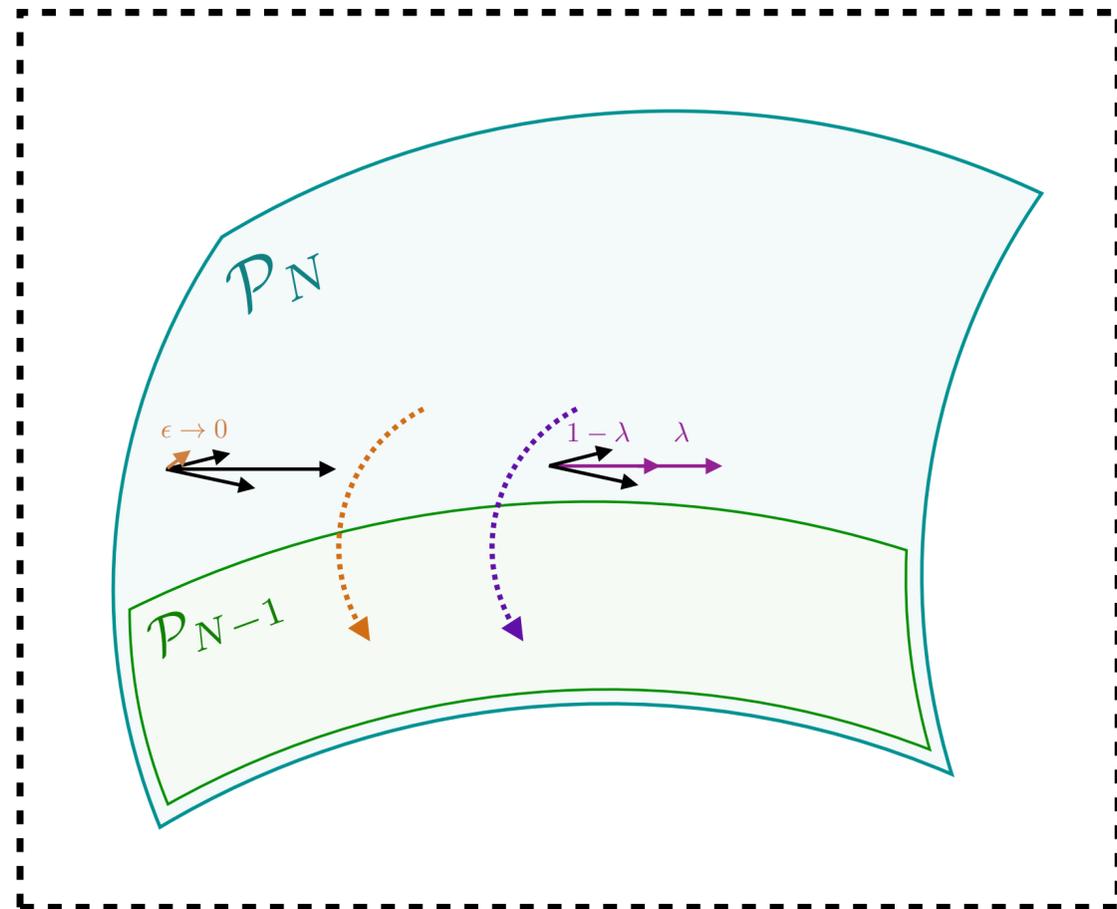
by **soft** and **collinear** limits

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

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

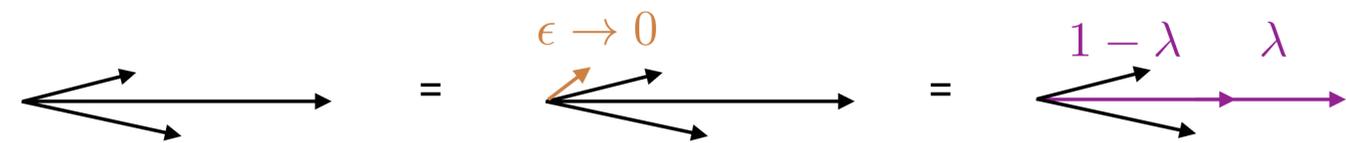
[PTK, Metodiev, Thaler, 2004.04159]

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



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

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

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

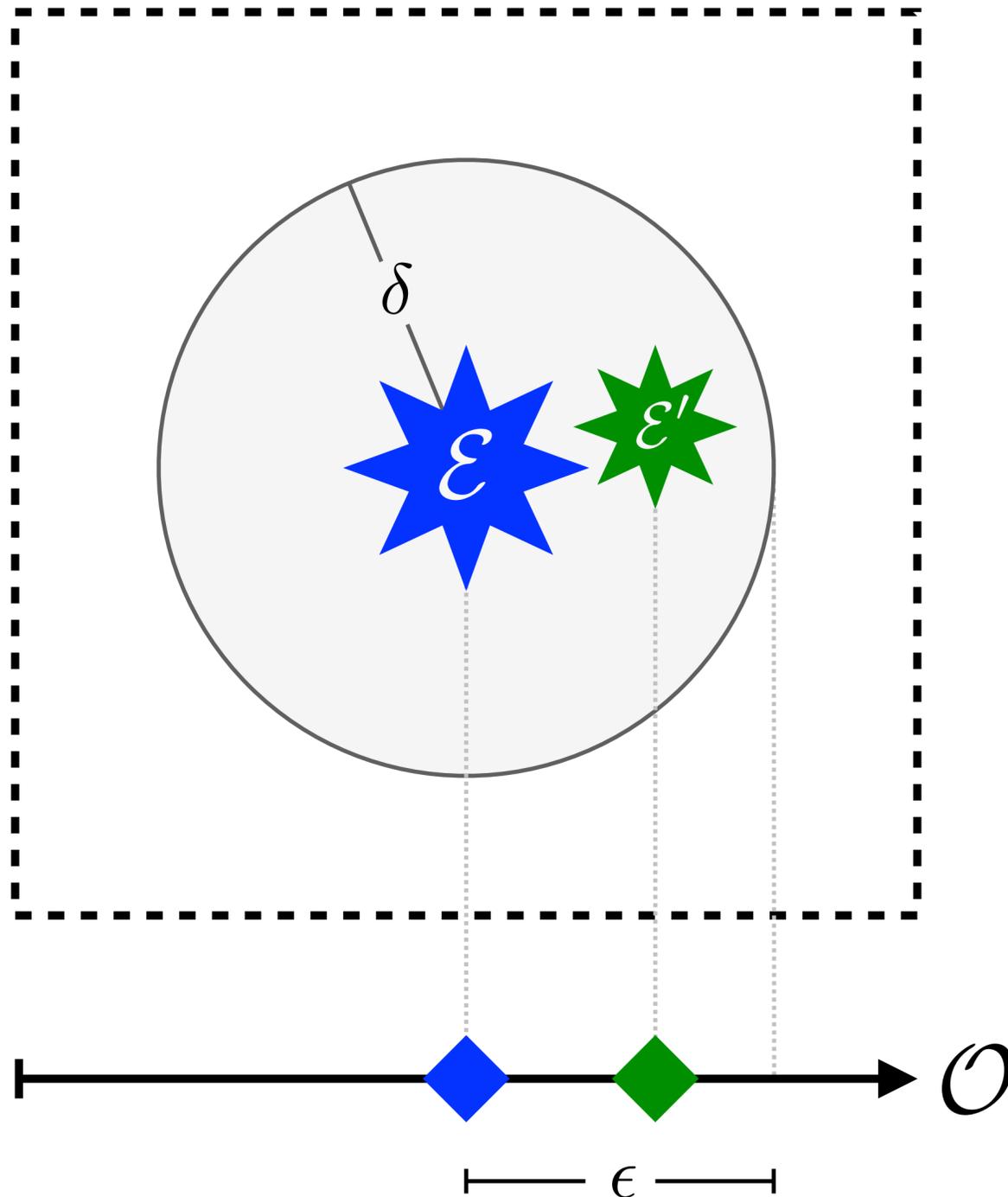
$$\mathcal{O}(p_1^\mu, \dots, p_M^\mu) = \lim_{p_0^\mu \rightarrow p_1^\mu} \mathcal{O}(\lambda p_0^\mu, (1 - \lambda)p_1^\mu, \dots, p_M^\mu)$$

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‡ of events

‡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]

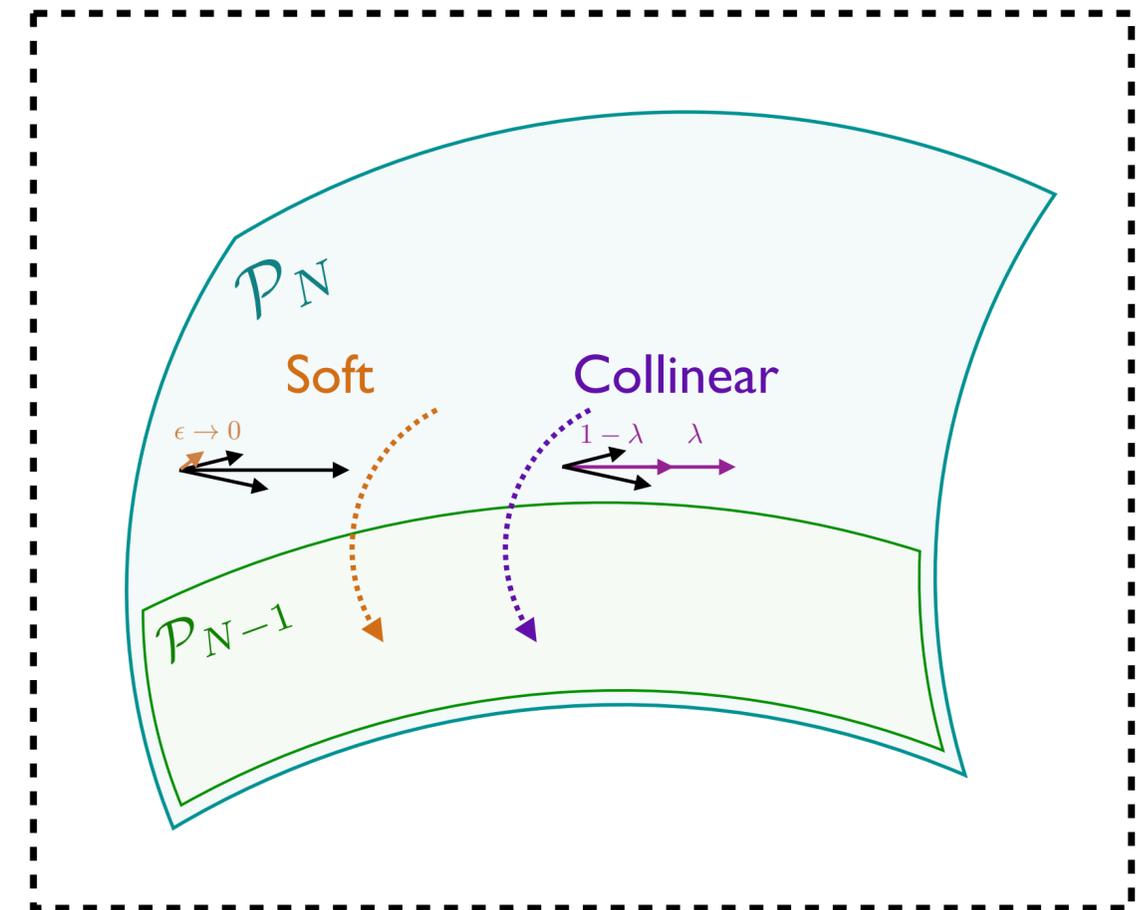
Is a statement of integrability on each P_N

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

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

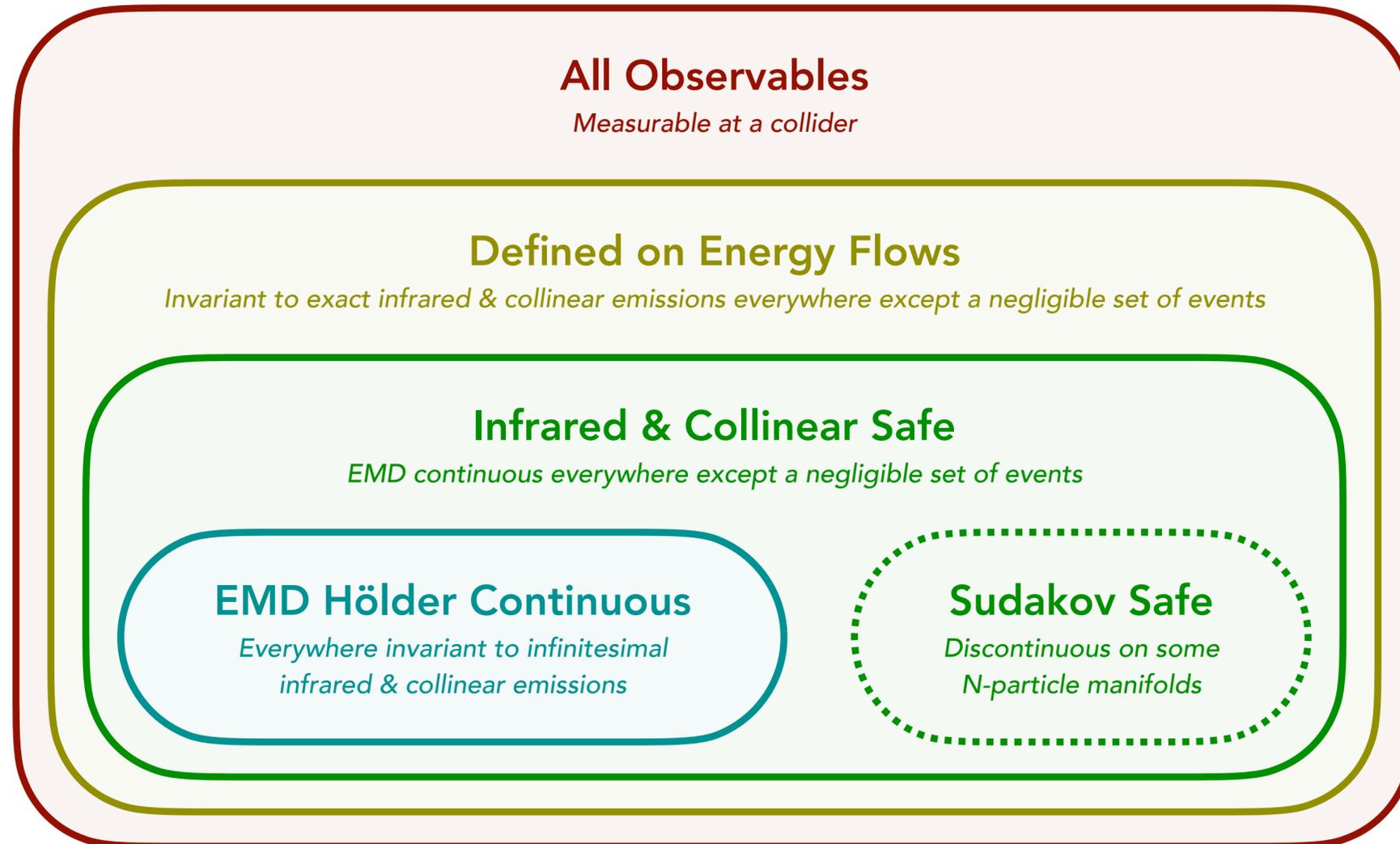
Example: $V(\mathcal{E}) = \mathcal{T}_2(\mathcal{E}) \left(1 + \frac{1}{\ln E(\mathcal{E})/\mathcal{T}_3(\mathcal{E})} \right)$ is EMD continuous but not EMD Hölder continuous (it is Sudakov safe)

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	
<i>Sudakov Safe</i>	
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] (τ_{N+1}/τ_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]	