# Theory Perspective on Machine Learning for Jets

## Jesse Thaler

MIT

# The NSF AI Institute for Artificial Intelligence and Fundamental Interactions (IAIFI) *"eye-phi"*



*Advance physics knowledge — from the smallest building blocks of nature to the largest structures in the universe — and galvanize AI research innovation*
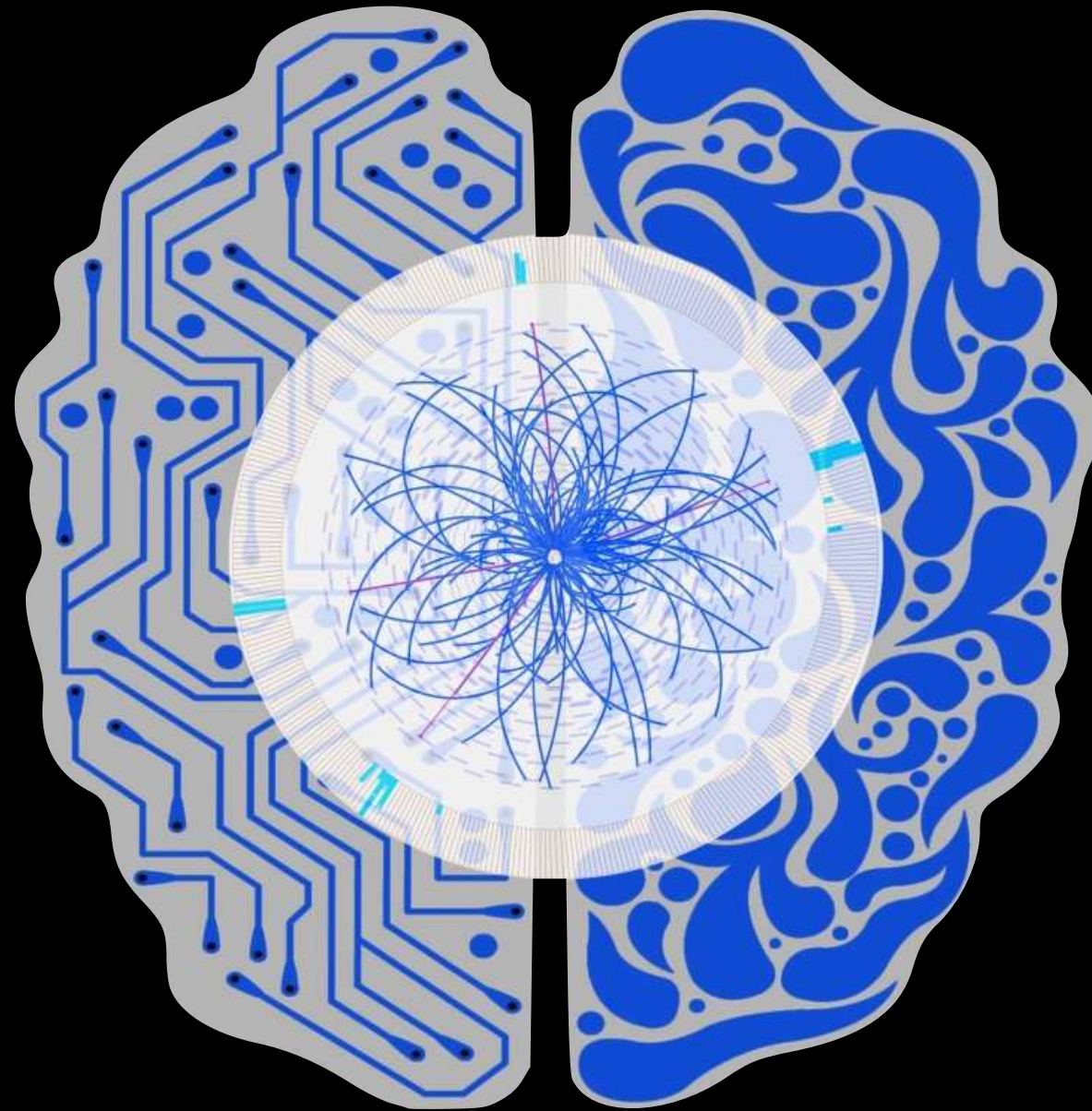
[http://iaifi.org, MIT News Announcement]

# AI²: Ab Initio Artificial Intelligence

*Machine learning* that incorporates
*first principles, best practices, and domain knowledge*
*from fundamental physics*

Symmetries, conservation laws, scaling relations, limiting behaviors, locality, causality,
unitarity, gauge invariance, entropy, least action, factorization, unit tests,
exactness, systematic uncertainties, reproducibility, verifiability, …

# Confronting the Black Box



*How do we develop robust machine learning for jet analyses?*

# Likelihood Ratio Trick

*Key example of simulation-based inference*

Goal:   Estimate p(x) / q(x)

Training Data:   Finite samples P and Q

Learnable Function:   f(x) parametrized by, e.g., neural networks

Loss Function(al):   $L = -\Big\langle \log f(x) \Big\rangle_P + \Big\langle f(x) - 1 \Big\rangle_Q$

Asymptotically:   $\underset{f(x)}{\arg\min} \, L = \dfrac{p(x)}{q(x)}$   *Likelihood ratio*

$-\underset{f(x)}{\min} \, L = \displaystyle\int \mathrm{d}x \, p(x) \log \frac{p(x)}{q(x)}$   *Kullback–Leibler divergence*

[see e.g. Cranmer, Pavez, Louppe, arXiv 2015; D'Agnolo, Wulzer, PRD 2019; simulation-based inference in Cranmer, Brehmer, Louppe, PNAS 2020; relation to f-divergences in Nguyen, Wainwright, Jordan, AoS 2009; Nachman, JDT, arXiv 2021]

# Likelihood Ratio Trick

*Key example of simulation-based inference*

Asymptotically, same structure as Lagrangian mechanics!

Action: $L = \int \mathrm{d}x\, \mathcal{L}(x)$

Lagrangian: $\mathcal{L}(x) = -p(x)\log f(x) + q(x)\big(f(x) - 1\big)$

Euler-Lagrange: $\dfrac{\partial \mathcal{L}}{\partial f} = 0$

Solution: $f(x) = \dfrac{p(x)}{q(x)}$

*Requires shift in theoretical focus from solving problems to specifying problems*

[see e.g. Cranmer, Pavez, Louppe, arXiv 2015; D'Agnolo, Wulzer, PRD 2019; simulation-based inference in Cranmer, Brehmer, Louppe, PNAS 2020; relation to f-divergences in Nguyen, Wainwright, Jordan, AoS 2009; Nachman, JDT, arXiv 2021]

*"What is the machine learning?"*

For this loss function, an estimate of the likelihood ratio derived from sampled data and regularized by the network architecture and training paradigm
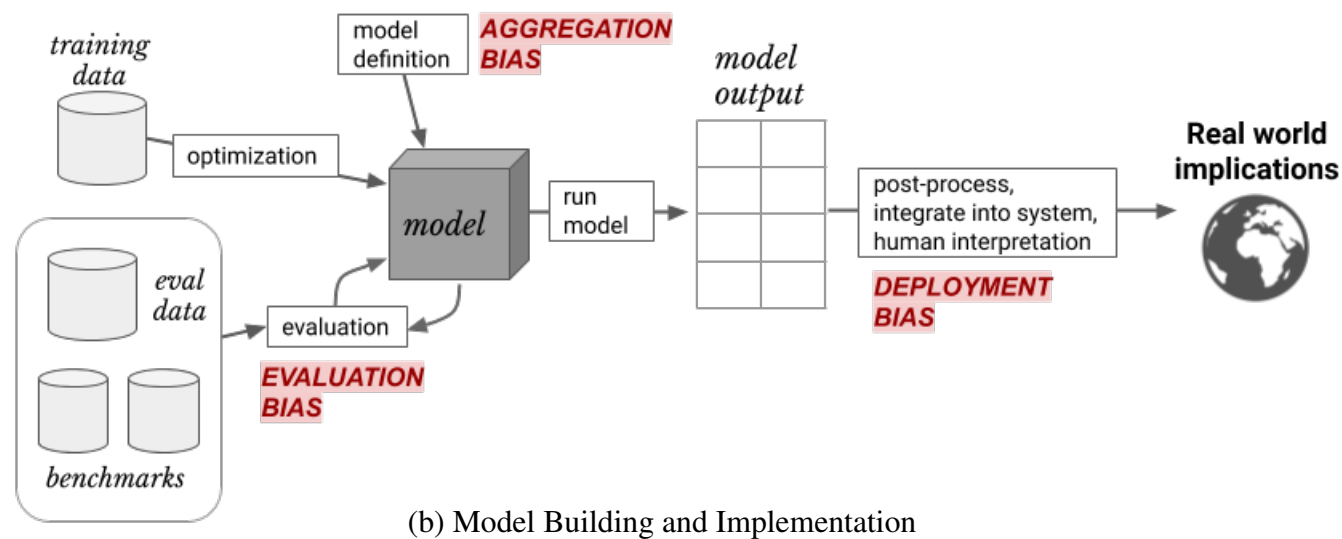
*"But where's the physics?!"*
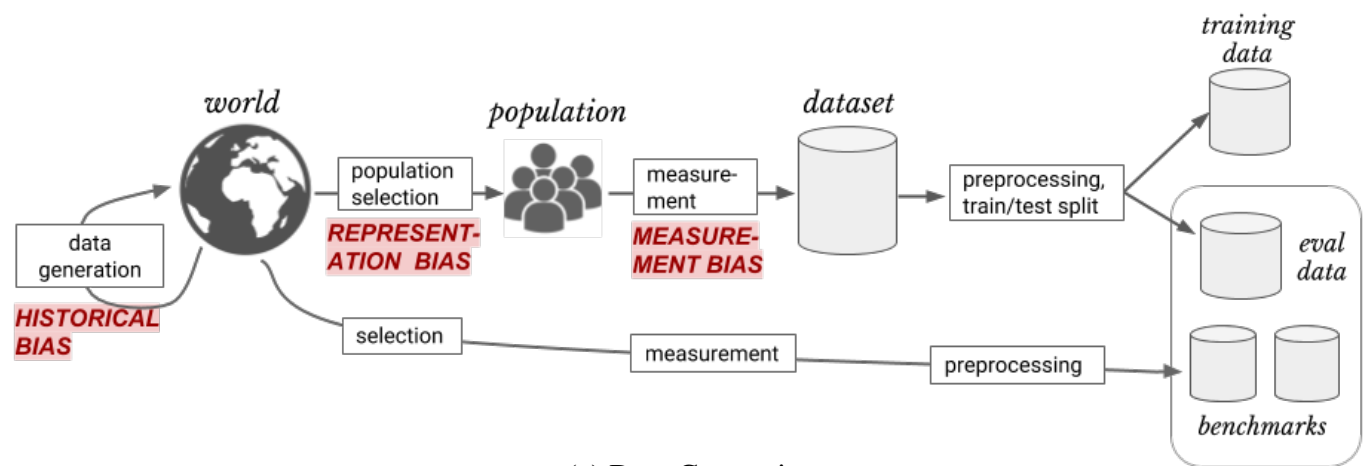
In the choice of loss function, data samples, network architecture, and training paradigm

*" . . . "*

# Many Reasons to be Wary!

*"A Framework for Understanding Unintended Consequences of Machine Learning"*



(a) Data Generation



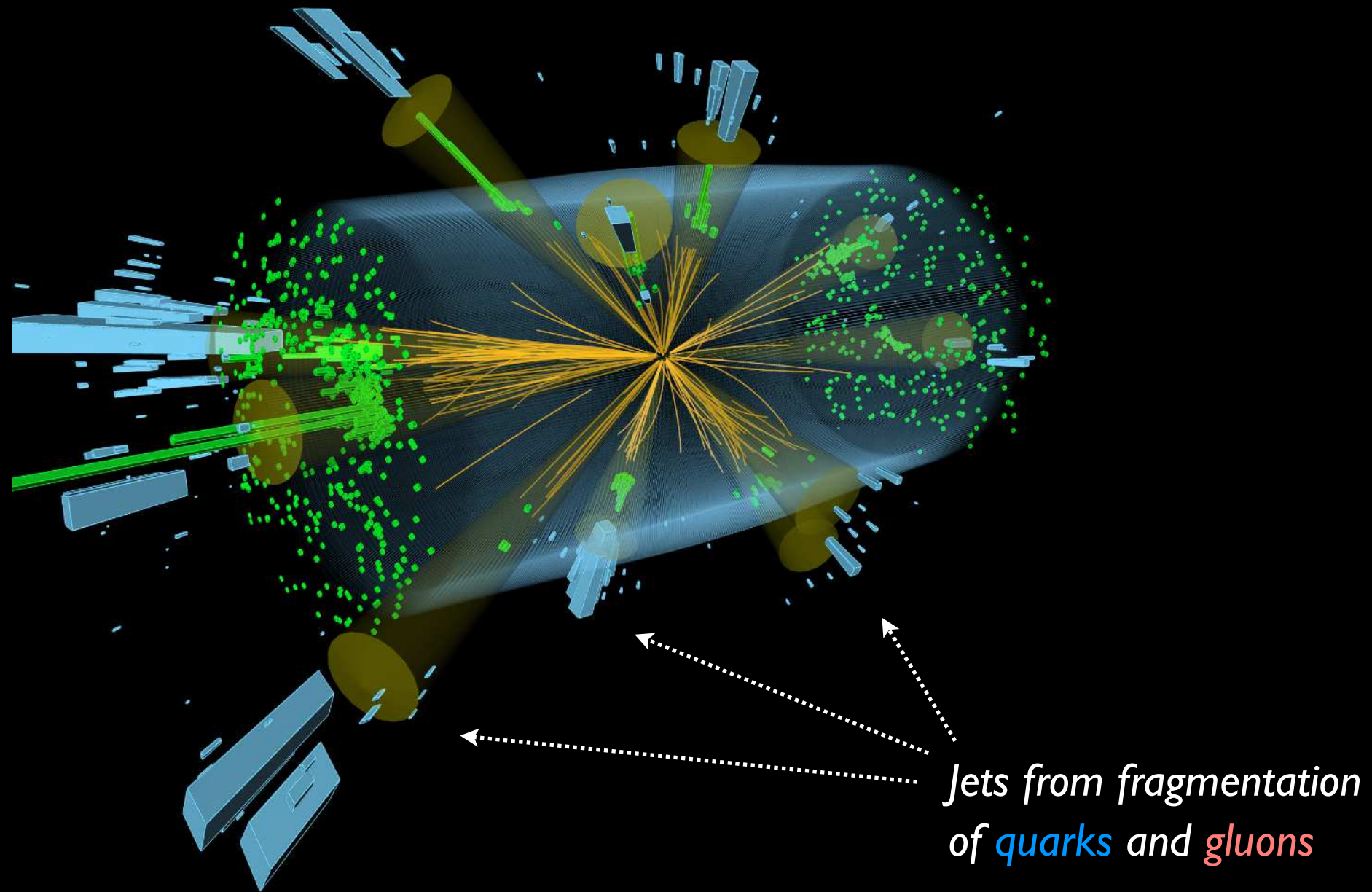(b) Model Building and Implementation

1. **Historical bias** arises when there is a misalignment between world as it is and the values or objectives to be encoded and propagated in a model. It is a normative concern with the state of the world, and exists even given perfect sampling and feature selection.

2. **Representation bias** arises while defining and sampling a development population. It occurs when the development population under-represents, and subsequently fails to generalize well, for some part of the use population.

3. **Measurement Bias** arises when choosing and measuring features and labels to use; these are often proxies for the desired quantities. The chosen set of features and labels may leave out important factors or introduce group- or input-dependent noise that leads to differential performance.

4. **Aggregation bias** arises during model construction, when distinct populations are inappropriately combined. In many applications, the population of interest is heterogeneous and a single model is unlikely to suit all subgroups.

5. **Evaluation bias** occurs during model iteration and evaluation. It can arise when the testing or external benchmark populations do not equally represent the various parts of the use population. Evaluation bias can also arise from the use of performance metrics that are not appropriate for the way in which the model will be used.

6. **Deployment Bias** occurs after model deployment, when a system is used or interpreted in inapppropriate ways.

For collider physics, "bias" ≈ "systematic uncertainty"

[h/t David Kaiser, MIT SERC; Suresh, Guttag, arXiv 2019]

# Machine Learning for Jet Substructure



Jets from fragmentation
of *quarks* and *gluons*

*How can we leverage theory to advance machine learning for jets?*

# My (Evolving) Perspective

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
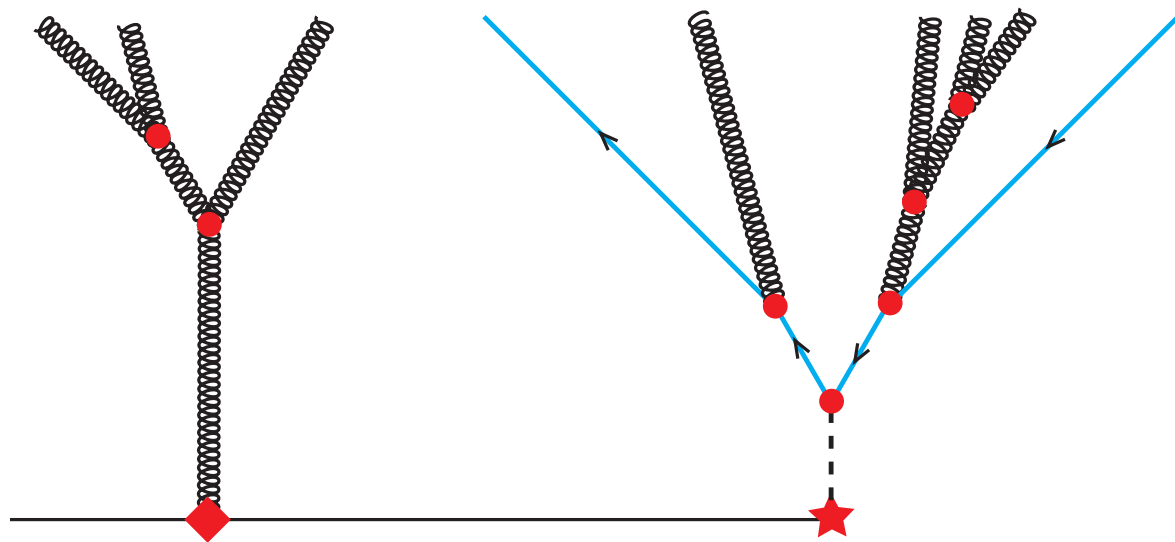structures and symmetries of our problem

Machine learning methods can only be as robust and reliable
as the data samples used for training

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms

*Examples below are representative, not exhaustive; apologies!*

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
structures and symmetries of our problem

Machine learning methods can only be as robust and reliable
as the data samples used for training

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms

# Compute Likelihood Ratios Directly?

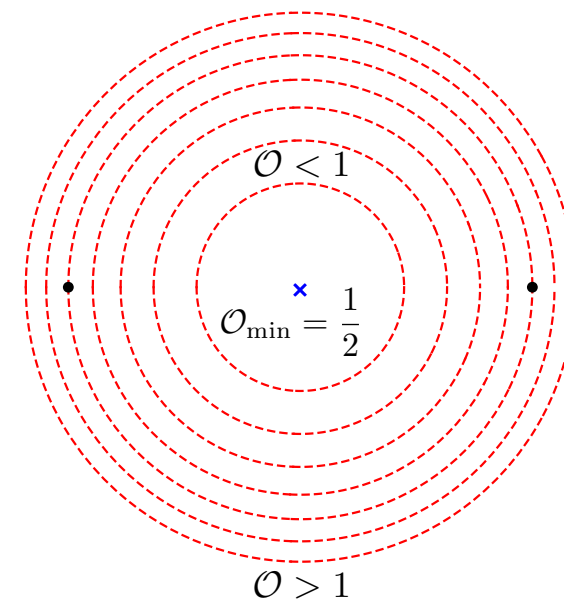*Yes, you can! And if you have enough calculational power, you should!*

## Shower Deconstruction



$$\chi(\{p,t\}_N) = \frac{P(\{p,t\}_N | \mathrm{S})}{P(\{p,t\}_N | \mathrm{B})}$$

[Soper, Spannowsky, PRD 2011]
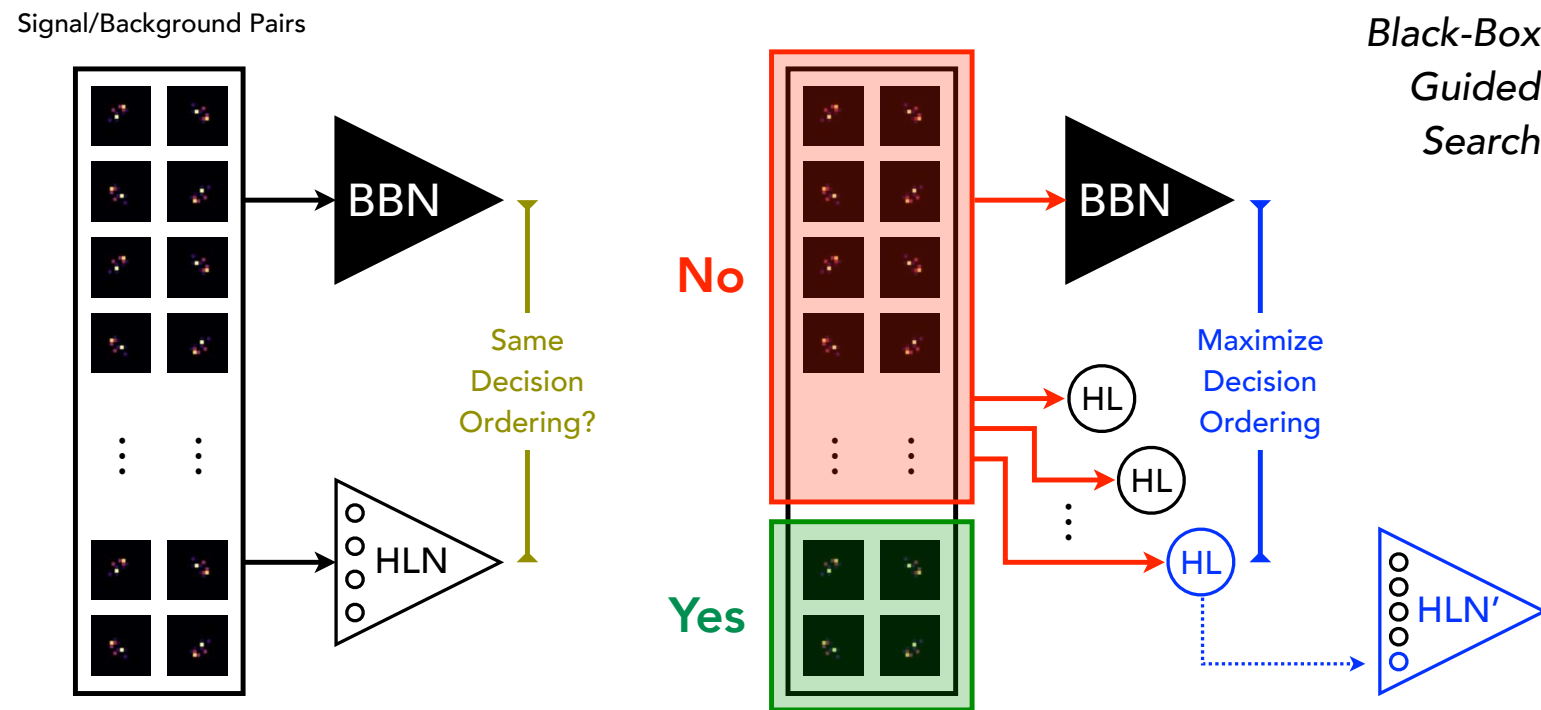
## Color Singlet Identification



$$\frac{|\mathcal{M}_\mathcal{B}|^2}{|\mathcal{M}_\mathcal{S}|^2} \simeq \frac{1 - \cos\theta_{ak} + 1 - \cos\theta_{bk}}{1 - \cos\theta_{ab}}$$
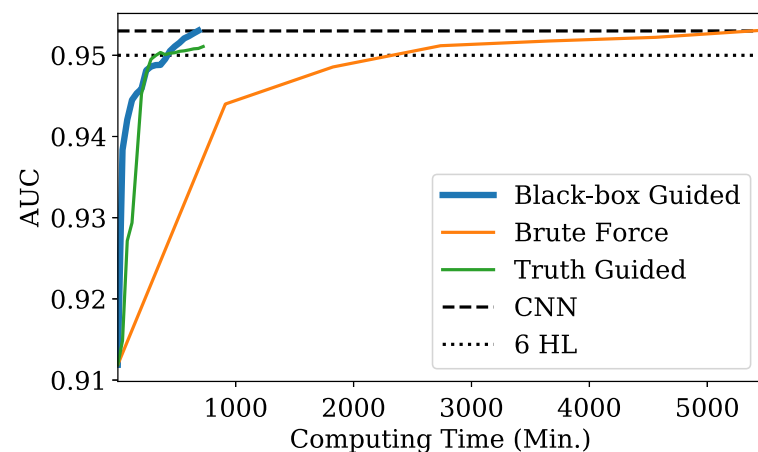
[Buckley, Callea, Larkoski, Marzani, SciPost 2020]

Challenge is that in most cases, best estimate of likelihood ratio comes from complex simulations with no closed form expression
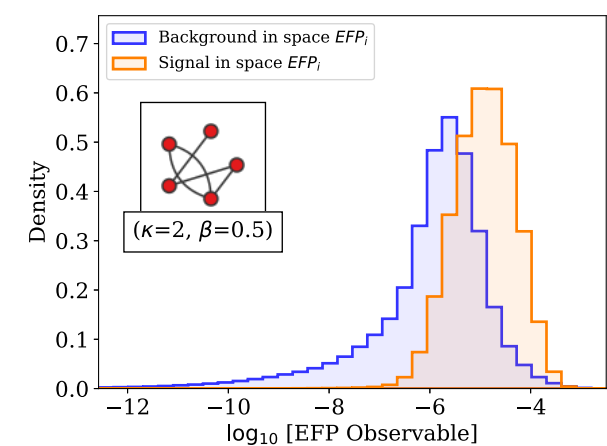
# Identifying Novel Jet Observables

*When **likelihood ratio** is not a function of standard high-level observables*



Signal/Background Pairs

BBN

HLN

Same Decision Ordering?

No

Yes

*Black-Box Guided Search*

BBN

HL

HL

HL

Maximize Decision Ordering

HLN'

A glimpse at an **alternative history** for field of jet substructure



| Iteration $(n)$ | EFP | $\kappa$ | $\beta$ | Chrom # |
|---|---|---|---|---|
| 0 | $M_{\text{jet}} + p_{\text{T}}$ | – | – | – |
| 1 | | 2 | $\frac{1}{2}$ | 2 |
| 2 | | 0 | 2 | 2 |
| 3 | | 0 | – | 1 |

Background in space $EFP_i$
Signal in space $EFP_i$
$(\kappa=2, \beta=0.5)$

Density

$\log_{10}$ [EFP Observable]

AUC

Computing Time (Min.)

Black-box Guided
Brute Force
Truth Guided
CNN
6 HL

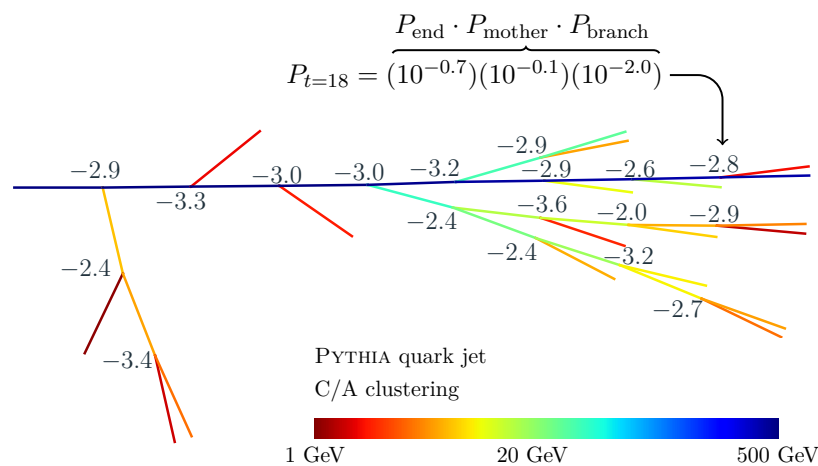[Faucett, JDT, Whiteson, PRD 2021; using Komiske, Metodiev, JDT, JHEP 2018]

# Theory-Inspired Likelihood Parametrizations

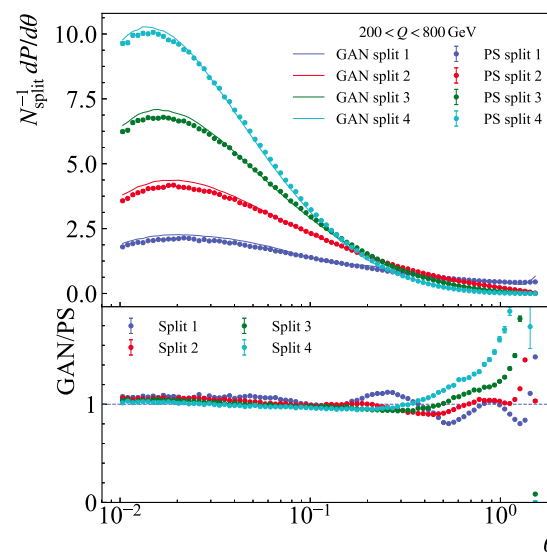*Flexible frameworks for parton-shower-like modeling*

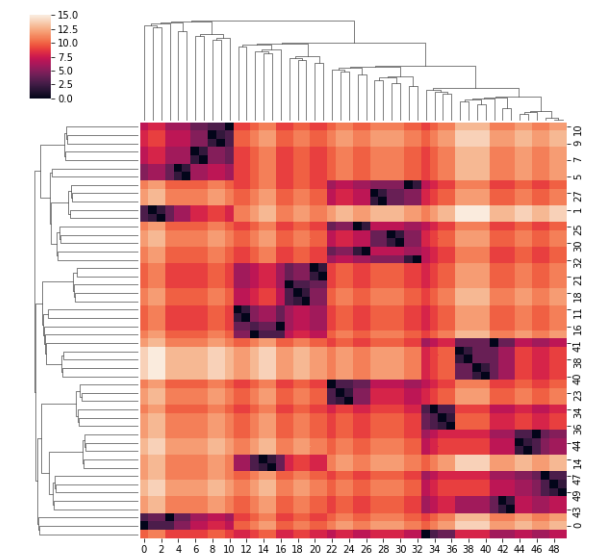## JUNIPR



[Andreassen, Fei...]

## DGLAP White Box



[Lai, Neill, Płoskoń, Ringer, arXiv 2020]

## Ginkgo



[Cranmer, Drnevich, Macaluso, Pappadopulo, arXiv 2021]

...s are based on constrained generative models,
...generalizing better than generic methods

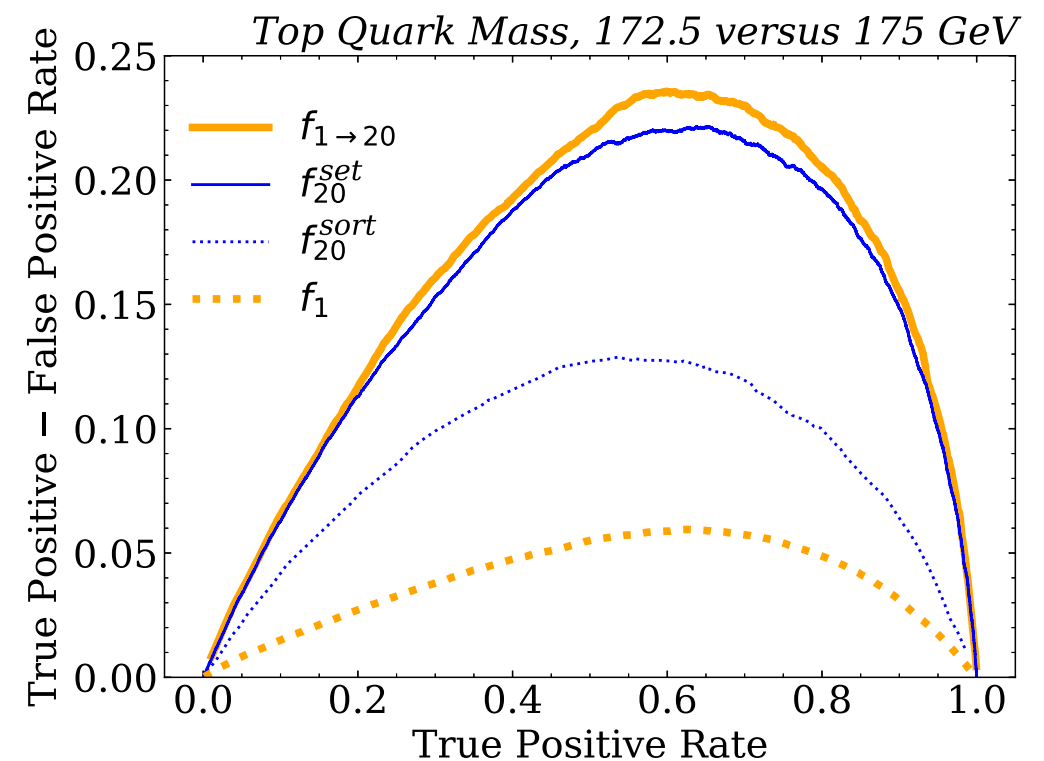Exhibit close relationship between generation and inference

# Sidestepping Per-Event Likelihood Ratios?

*Could be helpful, though per-event information is usually complete*

Collider events are independent and identically distributed…

$$\prod_{i=1}^{N} \frac{p(x_i|\theta_A)}{p(x_i|\theta_B)} = \frac{p(\{x_1,\ldots,x_N\}|\theta_A)}{p(\{x_1,\ldots,x_N\}|\theta_B)}$$

…therefore, (parametrized) per-event binary classifiers can be used to construct asymptotically optimal per-ensemble inference tools



*Top Quark Mass, 172.5 versus 175 GeV*

Legend:
- $f_{1\to20}$
- $f_{20}^{set}$
- $f_{20}^{sort}$
- $f_1$

y-axis: True Positive − False Positive Rate
x-axis: True Positive Rate

[Nachman, JDT, arXiv 2021; see mixed sample discussion in Metodiev, Nachman, JDT, JHEP 2017]

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
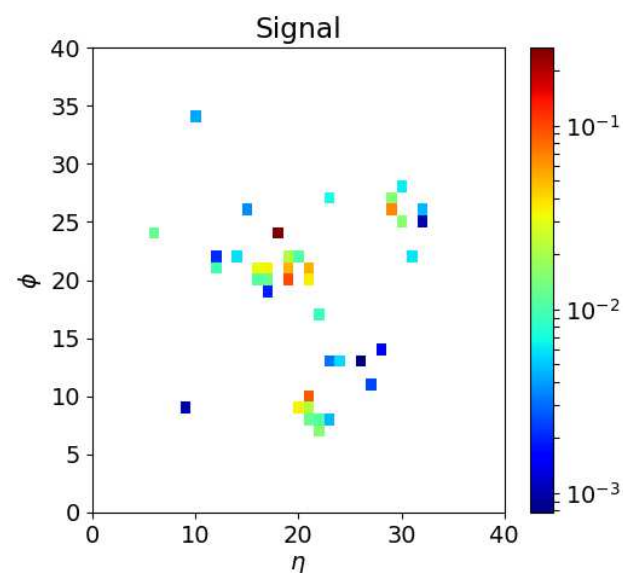structures and symmetries of our problem

Machine learning methods can only be as robust and reliable
as the data samples used for training

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms
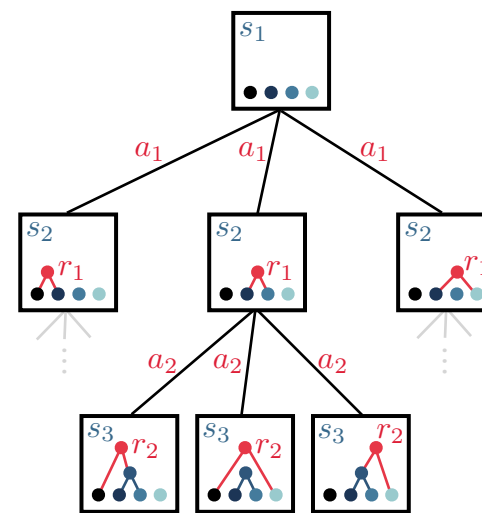
# Jet Representations

## Pixelized Image

*Calorimetry*



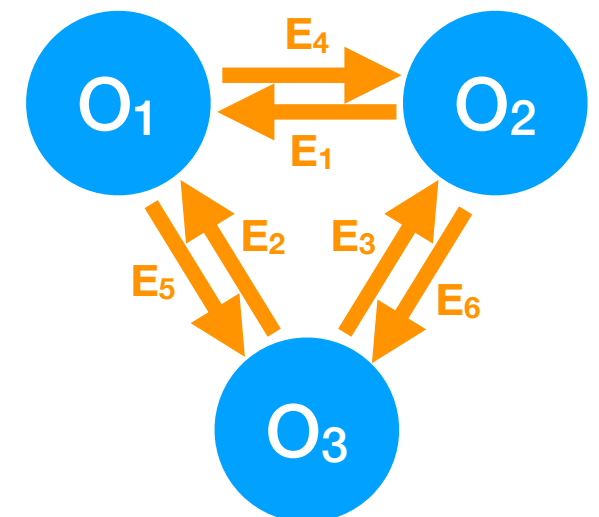[review in Kagan, arXiv 2020]

## Hierarchical Tree

*Binary Splittings*



[e.g. Brehmer, Macaluso, Pappadopulo, Cranmer, NeurIPS 2020]
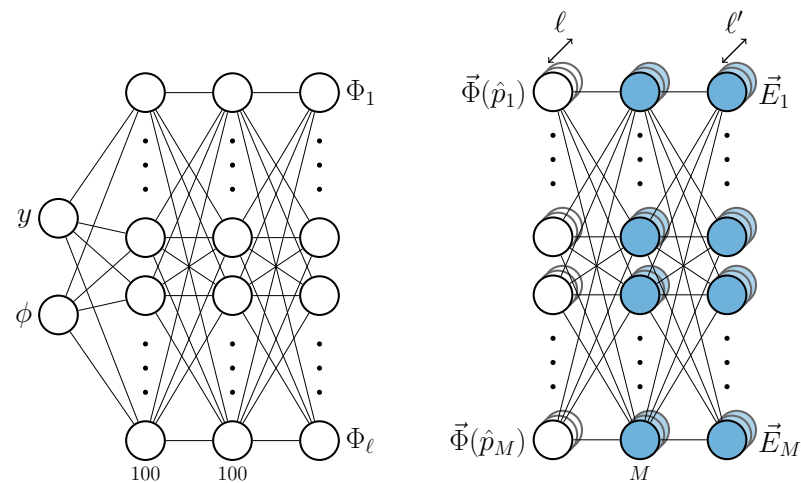
## Graphs

*Pairwise Interactions*



[e.g. Moreno, Cerri, Duarte, Newman, Nguyen, Periwal, Pierini, Serikova, Spiropulu, Vlimant, EPJC 2020]

*Imposes implicit theoretical prior (typically a good thing!)*
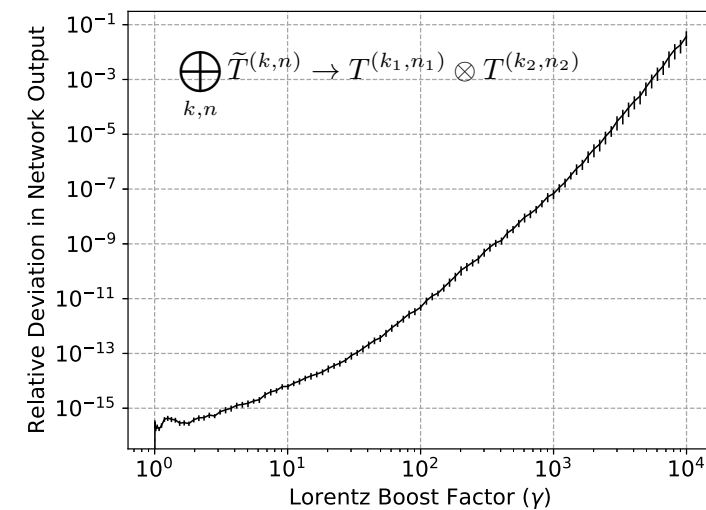*Influences choice of network architecture*

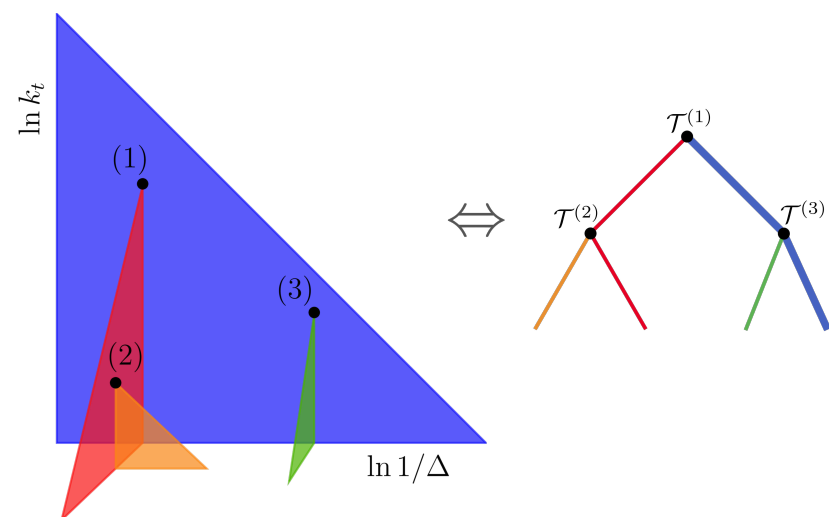# From Principles to Network Architectures

## Permutation Equivariance



[Dolan, Ore, PRD 2021]

## Lorentz Equivariance



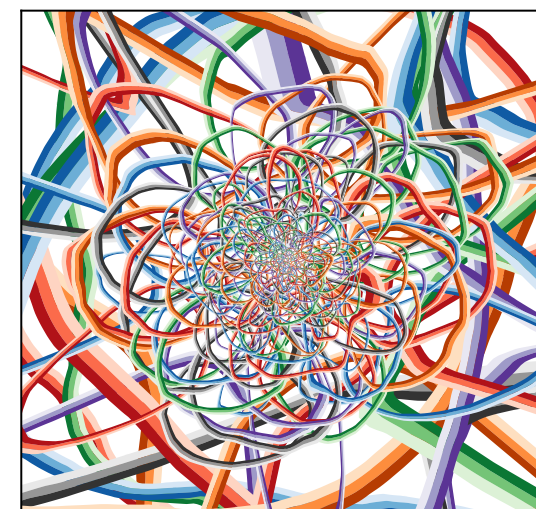$$\bigoplus_{k,n} \widetilde{T}^{(k,n)} \to T^{(k_1,n_1)} \otimes T^{(k_2,n_2)}$$

[Bogatskiy, Anderson, Offermann, Roussi, Miller, Kondor, arXiv 2020]

## Lund Plane Emissions



[Dreyer, Qu, JHEP 2021]

## Infrared and Collinear Safety



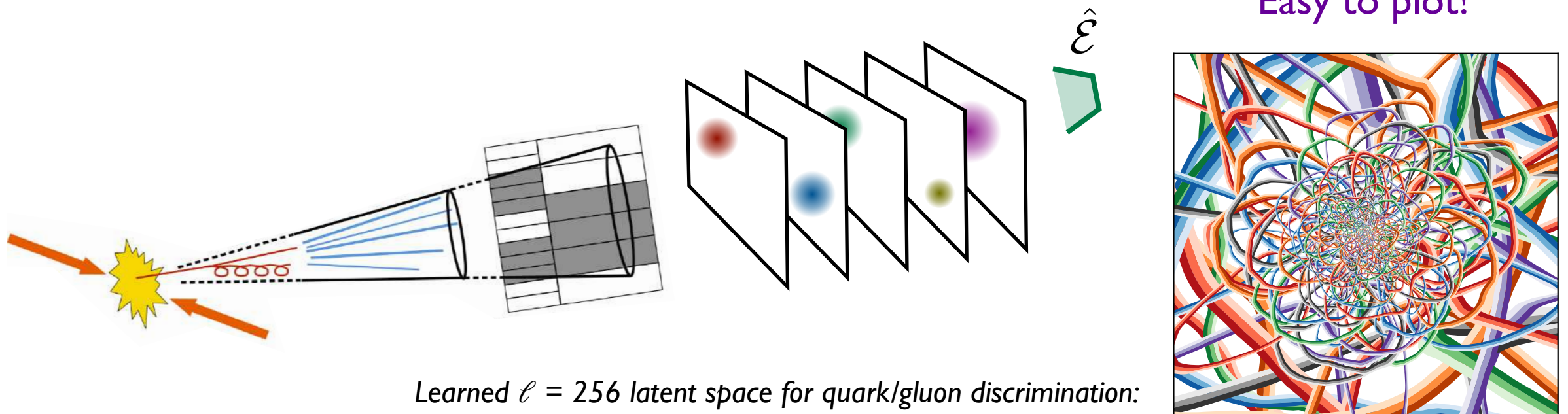[Komiske, Metodiev, JDT, JHEP 2019]

# Energy Flow Networks

*Architecture designed around symmetries and interpretability*

Latent space of dim $\ell$

Permutation invariant

Linear weights (IRC safe)

$$S(\mathcal{J}) = F(V_1, V_2, \ldots, V_\ell) \qquad V_a(\mathcal{J}) = \sum_{i \in \mathcal{J}} E_i \, \Phi_a(\hat{n}_i)$$
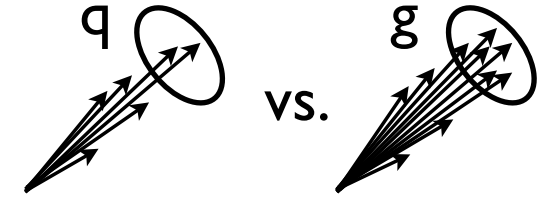
Easy to plot!

$\hat{\mathcal{E}}$



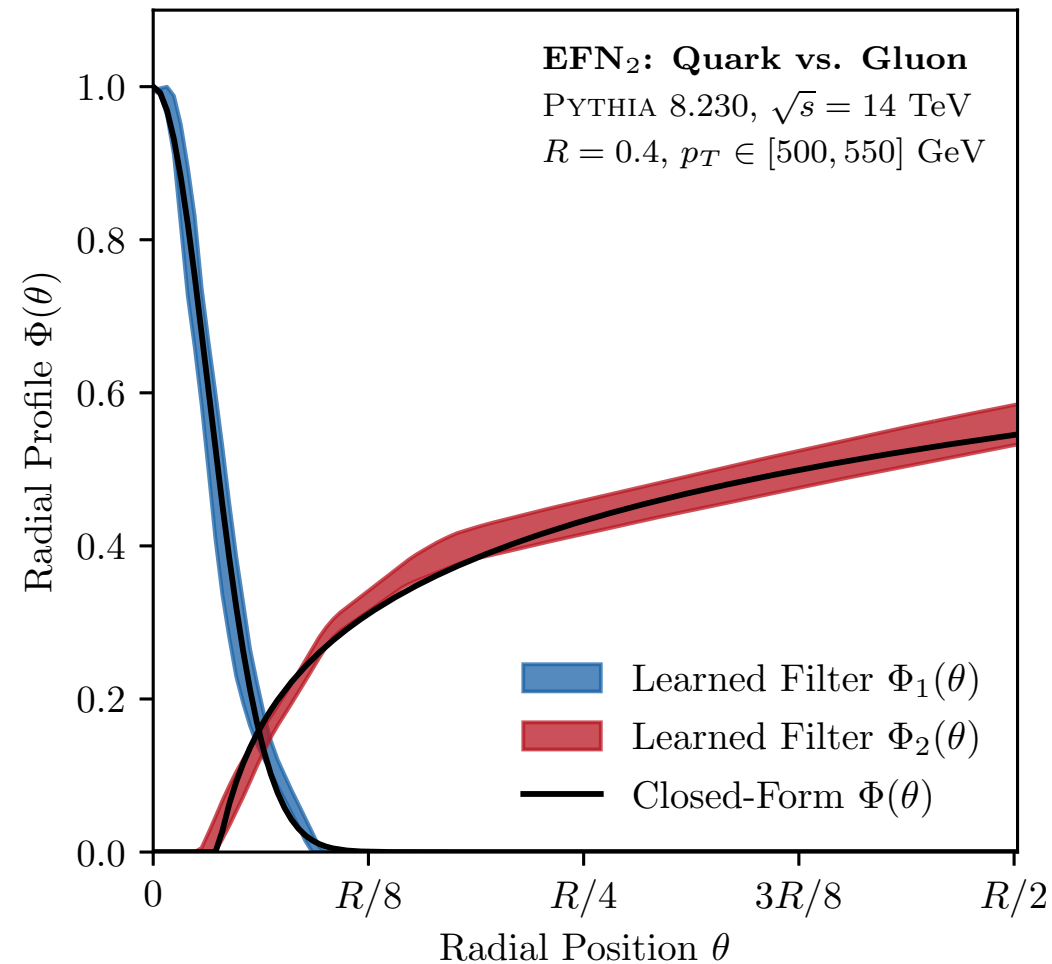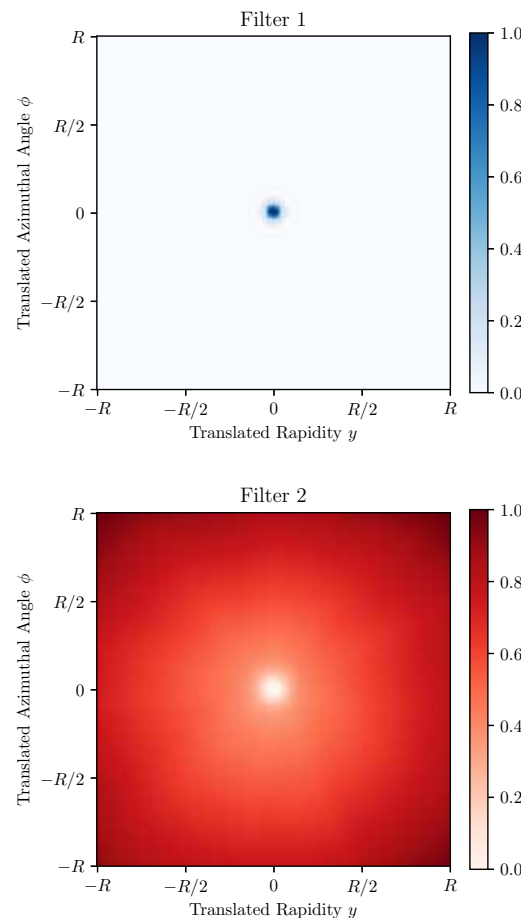*Learned $\ell = 256$ latent space for quark/gluon discrimination:*

[Komiske, Metodiev, JDT, JHEP 2019; see also Komiske, Metodiev, JDT, JHEP 2018; code at energyflow.network; special case of Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, NIPS 2017; histogram pooling in Cranmer, Kreisch, Pisani, Villaescusa-Navarro, Spergel, Ho, ICLR SimDL 2021]
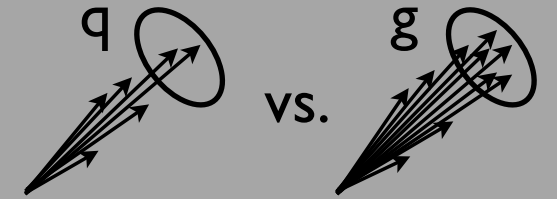
# Learning from the Machine

For $\ell = 2$, EFN learns radial moments: $\displaystyle\sum_{i\in\text{jet}} z_i\, f(\theta_i)$    cf. Angularities: $f(\theta) = \theta^\beta$



Filter 1

Filter 2

EFN$_2$: Quark vs. Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV

Learned Filter $\Phi_1(\theta)$
Learned Filter $\Phi_2(\theta)$
Closed-Form $\Phi(\theta)$

Radial Profile $\Phi(\theta)$

Radial Position $\theta$
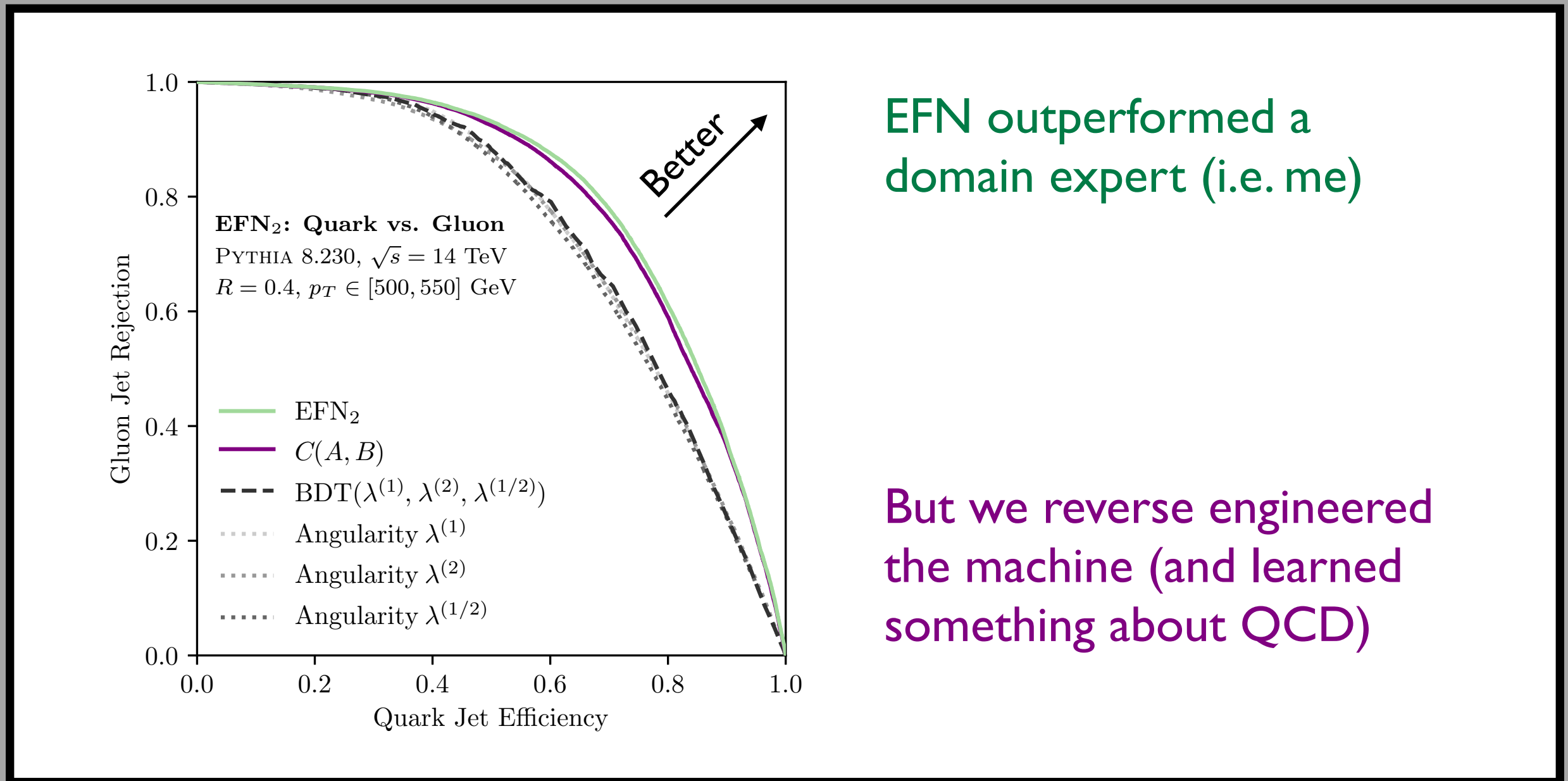
Traditional QCD observables emphasize homogeneous angular scaling
But EFN reveals that likelihood ratio exhibits collinear/wide-angle separation

[Komiske, Metodiev, JDT, JHEP 2019;
cf. Larkoski, JDT, Waalewijn, JHEP 2014; using Berger, Kucs, Sterman, PRD 2003; Ellis, Vermilion, Walsh, Hornig, Lee, JHEP 2010]

# Learning from the Machine

q vs. g

For $\ell = 2$, EFN learns radial moments: $\displaystyle\sum_{i \in \text{jet}} z_i f(\theta_i)$     cf. Angularities: $f(\theta) = \theta^\beta$



EFN$_2$: Quark vs. Gluon
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV

Legend:
- EFN$_2$
- $C(A, B)$
- BDT($\lambda^{(1)}, \lambda^{(2)}, \lambda^{(1/2)}$)
- Angularity $\lambda^{(1)}$
- Angularity $\lambda^{(2)}$
- Angularity $\lambda^{(1/2)}$

Axes: Gluon Jet Rejection vs. Quark Jet Efficiency. "Better" arrow.

EFN outperformed a domain expert (i.e. me)

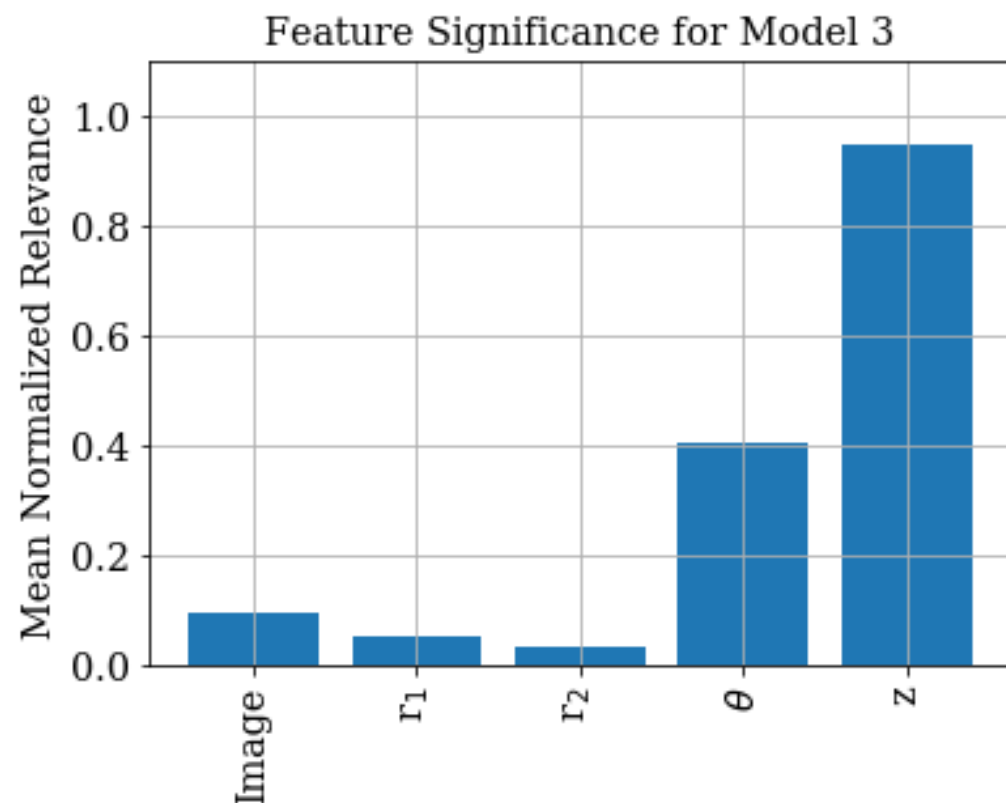But we reverse engineered the machine (and learned something about QCD)

[Komiske, Metodiev, JDT, JHEP 2019;
cf. Larkoski, JDT, Waalewijn, JHEP 2014; using Berger, Kucs, Sterman, PRD 2003; Ellis, Vermilion, Walsh, Hornig, Lee, JHEP 2010]
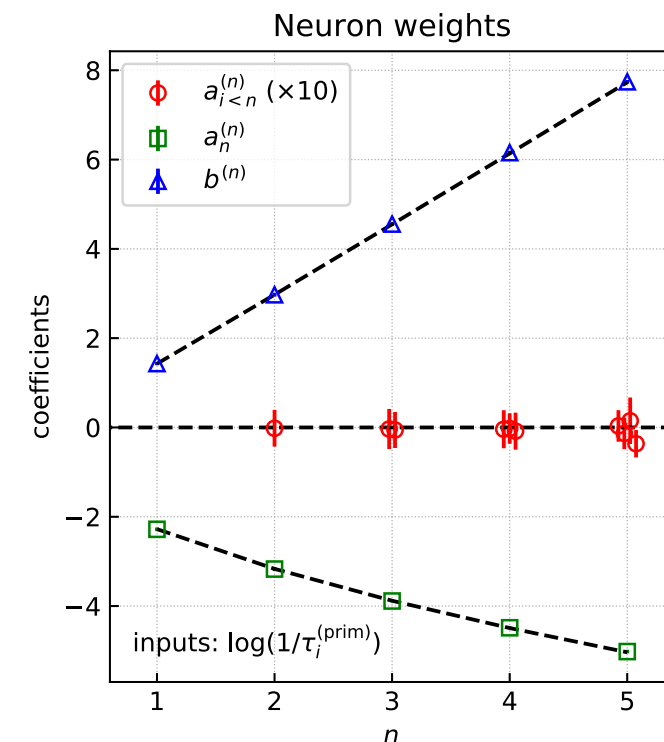
# More Network Architectures for Interpretability

*Rendering the black box more transparent*

## Input Feature Relevance



Feature Significance for Model 3

[Agarwal, Hay, Iashvili, Mannix, McLean, Morris, Rappoccio, Schubert, JHEP 2021]

## Analytic Calculations



Neuron weights

[Kasieczka, Marzani, Soyez, Stagnitto, JHEP 2020]

Imposing specific theoretical structures might reduce performance
but might also yield better robustness/generalizability

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
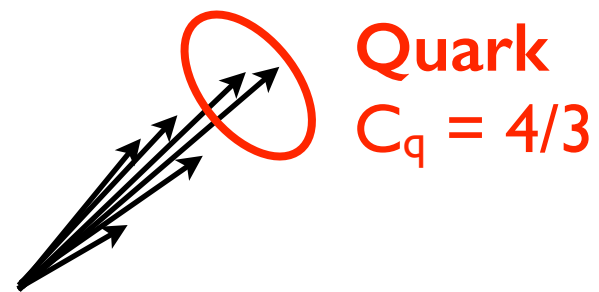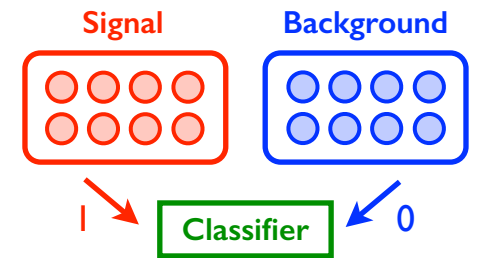approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
structures and symmetries of our problem

**Machine learning methods can only be as robust and reliable
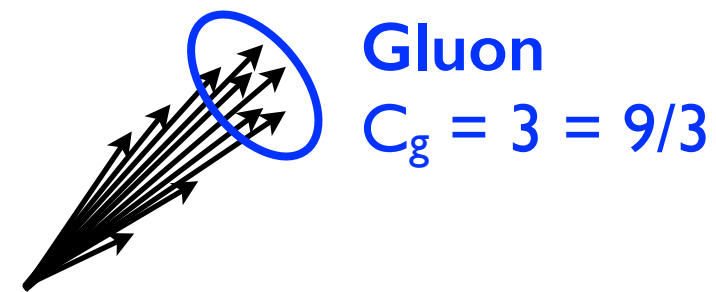as the data samples used for training**

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms

# Quark/Gluon Classification

*"Hello, World!" of Jet Physics*

**Quark**
$C_q = 4/3$

**vs.**

**Gluon**
$C_g = 3 = 9/3$

Find $h\left(\ \right)$ such that

$$h(\text{Quark}) = 1$$
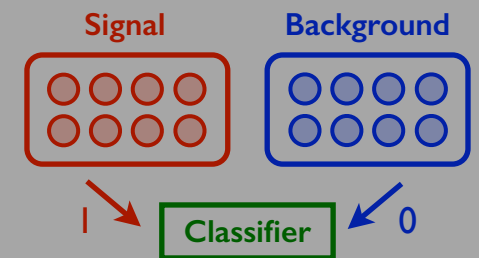$$h(\text{Gluon}) = 0$$

**Best you can do:** $h(\mathcal{J}) = \left(1 + \dfrac{p(\mathcal{J}|\text{G})}{p(\mathcal{J}|\text{Q})}\right)^{-1}$

(Neyman-Pearson lemma)

*Likelihood ratio yields optimal binary classifier (and vice versa)*

[see e.g. Gras, Höche, Kar, Larkoski, Lönnblad, Plätzer, Siódmok, Skands, Soyez, JDT, JHEP 2017; Komiske, Metodiev, Schwartz, JHEP 2017; Komiske, Metodiev, JDT, JHEP 2018]

# Quark/Gluon Classification

*"Hello, World!" of Jet Physics*

**Quark**
$C_q = 4/3$

**vs.**

**Gluon**
$C_g = 3 = 9/3$

*What do you mean by "quark" and "gluon"?*

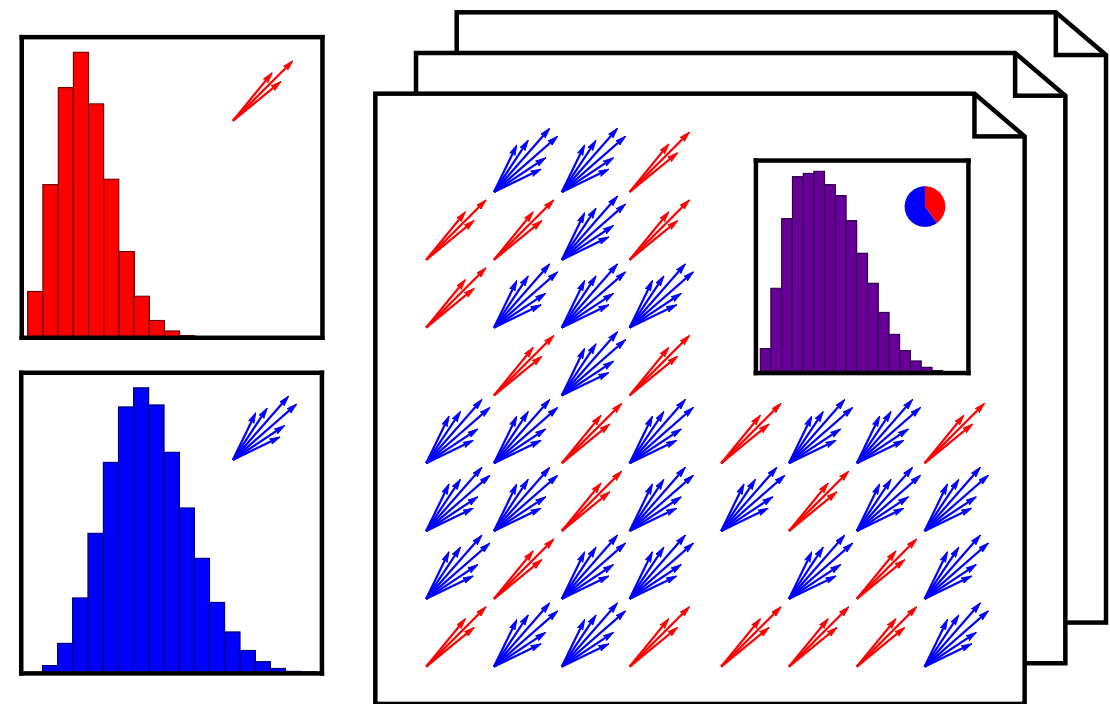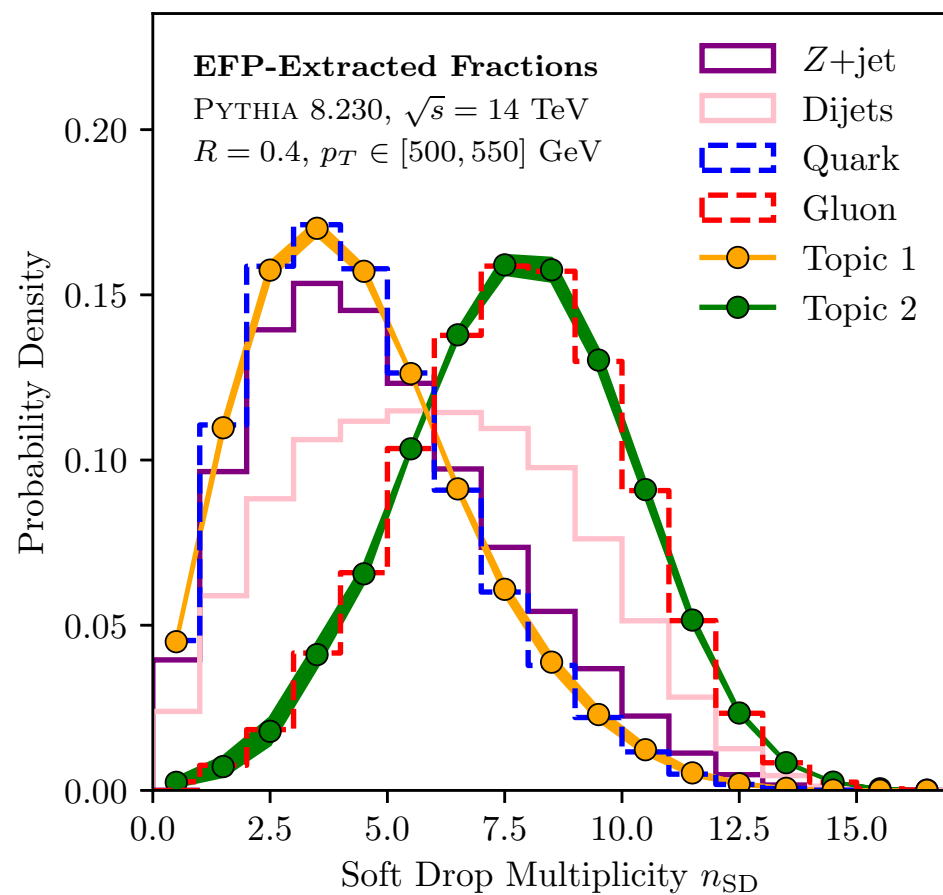*Jets are clusters of colorless hadrons!*

*Parton shower "truth" is but a (useful) fiction!*

*Likelihood ratio yields optimal binary classifier (and vice versa)*

[see e.g. Gras, Höche, Kar, Larkoski, Lönnblad, Plätzer, Siódmok, Skands, Soyez, JDT, JHEP 2017; Komiske, Metodiev, Schwartz, JHEP 2017; Komiske, Metodiev, JDT, JHEP 2018]

# Topic Modeling to Disentangle Data Samples

While you can't unambiguously label individual jets, you can
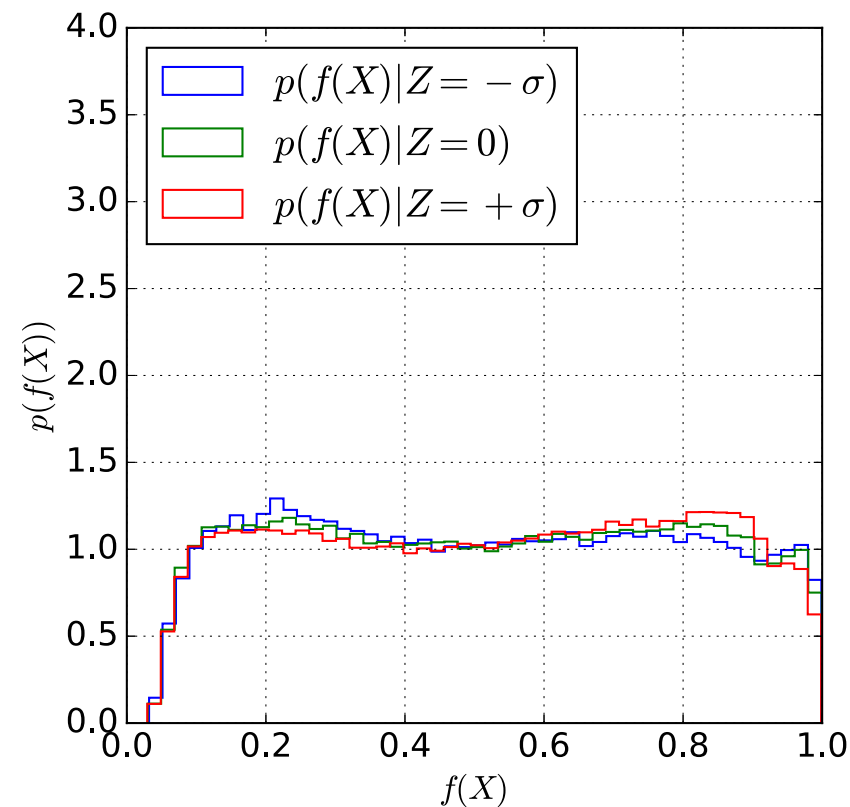extract quark and gluon distributions from hadron-level measurements



EFP-Extracted Fractions
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV

Legend:
- $Z$+jet
- Dijets
- Quark
- Gluon
- Topic 1
- Topic 2

Probability Density vs Soft Drop Multiplicity $n_{\mathrm{SD}}$

Key concept from natural language
processing: "anchor words"

[Komiske, Metodiev, JDT, JHEP 2018; cf. ATLAS, PRD 2019; using Metodiev, Nachman, JDT, JHEP 2017; Metodiev, JDT, PRL 2018]
see also Blanchard, Flaska, Handy, Pozzi, Scott, PLMR 2013; Katz-Samuels, Blanchard, Scott, JMLR 2016]
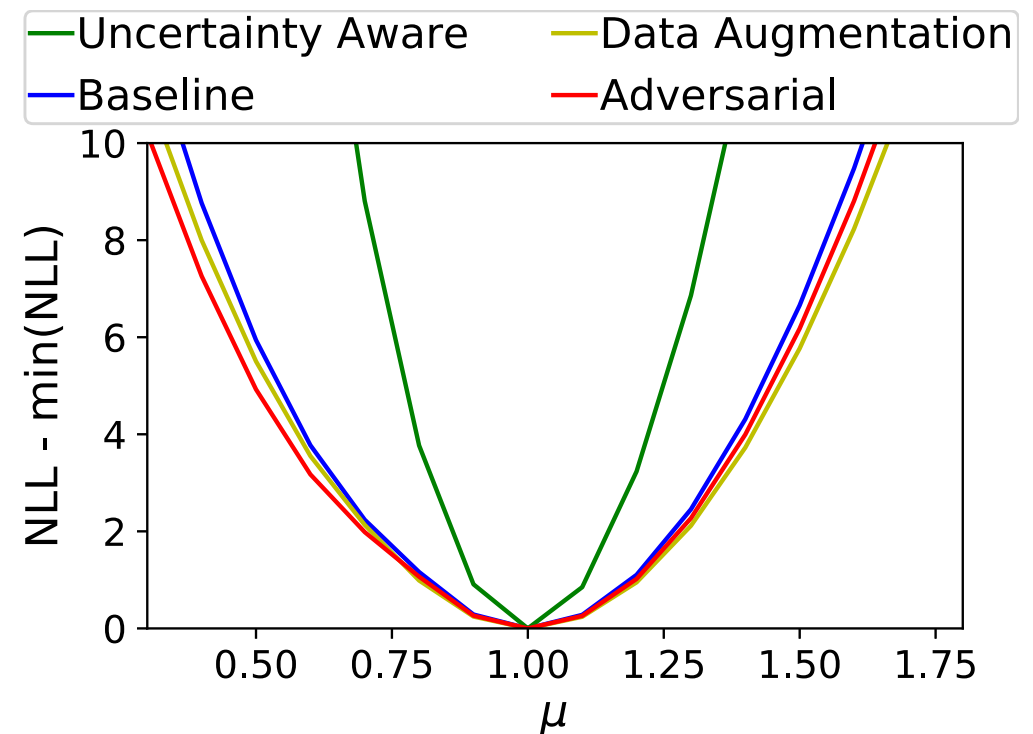
# Parameterized Data Samples

*Incorporating nuisance parameters into training*

## Learn to Pivot…                     …or Learn to Profile?



Legend (left plot):
- $p(f(X)|Z=-\sigma)$
- $p(f(X)|Z=0)$
- $p(f(X)|Z=+\sigma)$

[Louppe, Kagan, Cranmer, NeurIPS 2017]



Legend (right plot):
- Uncertainty Aware
- Data Augmentation
- Baseline
- Adversarial

[Ghosh, Nachman, Whiteson, arXiv 2021]

## Regardless of the approach, inspires renewed theoretical focus on uncertainty modeling for Monte Carlo generation

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
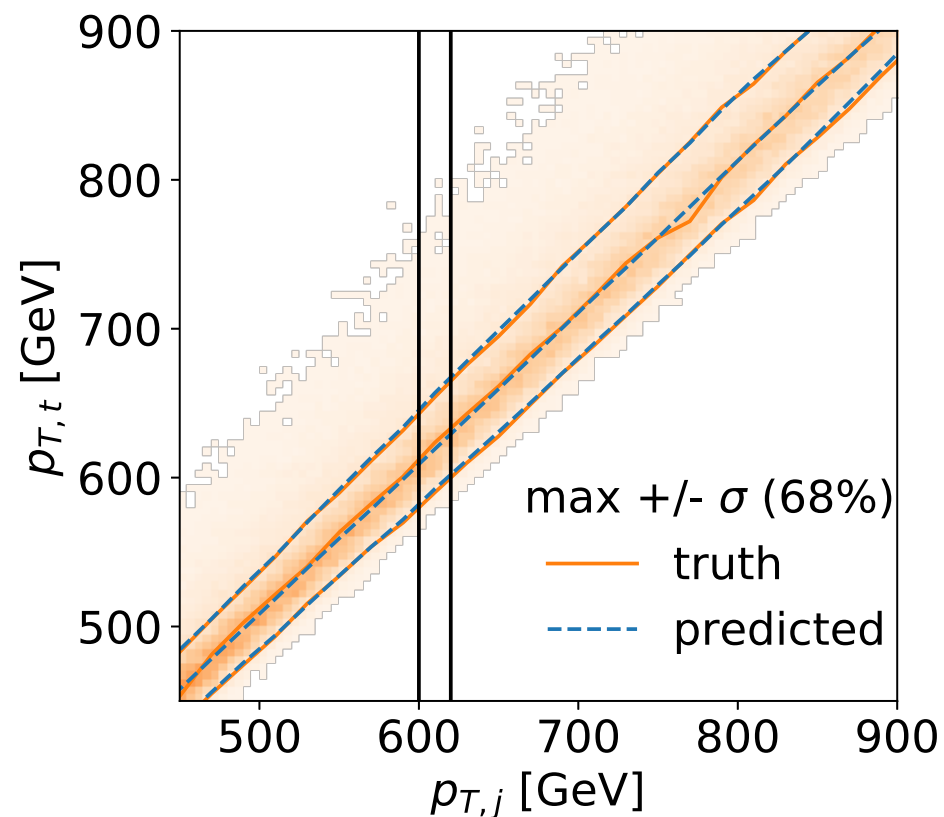structures and symmetries of our problem

Machine learning methods can only be as robust and reliable
as the data samples used for training

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms
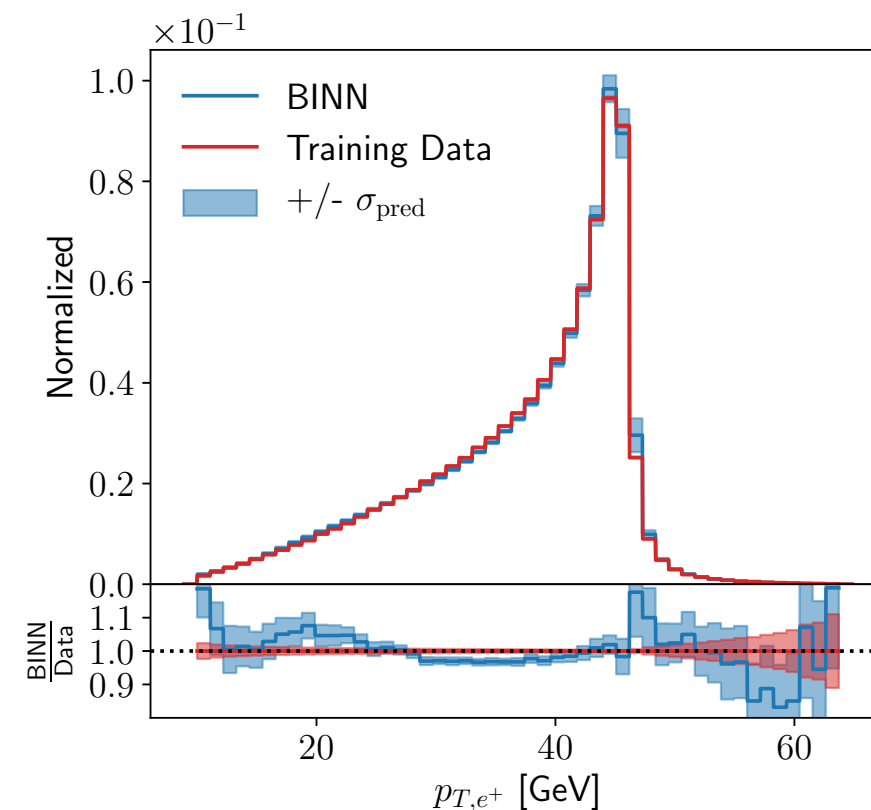
# Loss Function for Bayesian Analyses

*Treating network parameters as having a prior distribution*

## Jet Energy Scale



[Kasieczka, Luchmann, Otterpohl, Plehn, SciPost 2020]

## Event Generation



[Bellagente, Haußmann, Luchmann, Plehn, arXiv 2021]

Use ELBO loss to capture both statistical and systematic uncertainties
Worth developing a frequentist version of this approach?

# Loss Function for Bayesian Analyses

*Treating network parameters as having a prior distribution*

Nothing intrinsically wrong with Bayesian analyses,
but have to be aware of cases with strong prior dependence

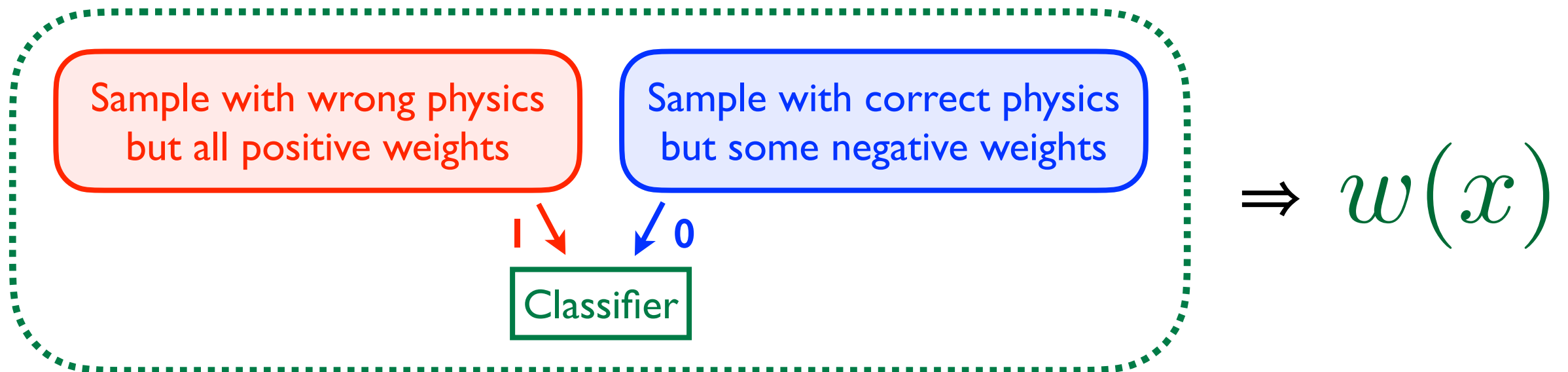$$\int d\theta \, \theta \, p(\theta|\text{data}) \neq \max_{\theta} p(\text{data}|\theta)$$

If needed, can treat neural networks as static objects
and calibrate them as if they were ordinary observables

[see further discussion in Cranmer, Pavez, Louppe, arXiv 2015; Nachman, SciPost 2020]

Use ELBO loss to capture both statistical and systematic uncertainties
Worth developing a frequentist version of this approach?

# Training Paradigm for Preserved Uncertainties

*When the goal is to maintain statistical properties*

Sample with wrong physics but all positive weights

Sample with correct physics but some negative weights

$\Rightarrow w(x)$

1 → ← 0

Classifier

**Plain reweighting** yields all positive weights with correct asymptotic probability density

$$p_{\text{wrong}}(x) \times w(x) = p_{\text{correct}}(x)$$

**Improved resampling** through auxiliary neural network yields correct statistical uncertainties
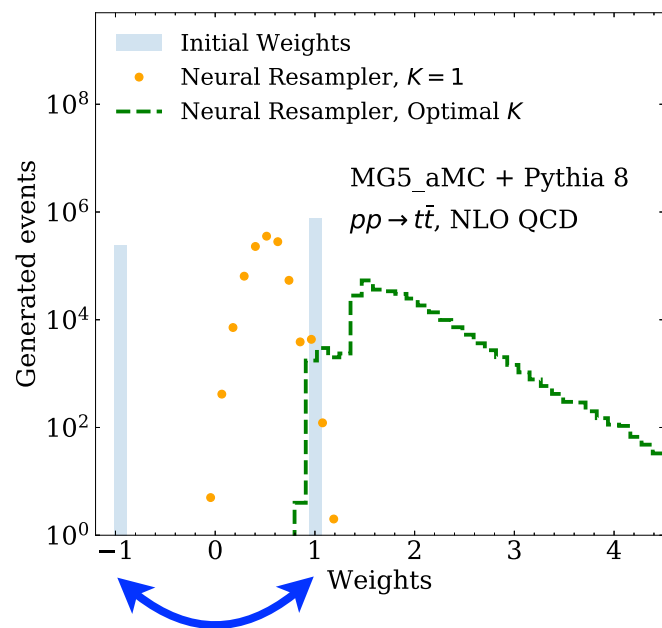
$$\left(\frac{\delta p}{p}\right)^2 = \frac{\langle w^2 \rangle}{\langle w \rangle^2}$$

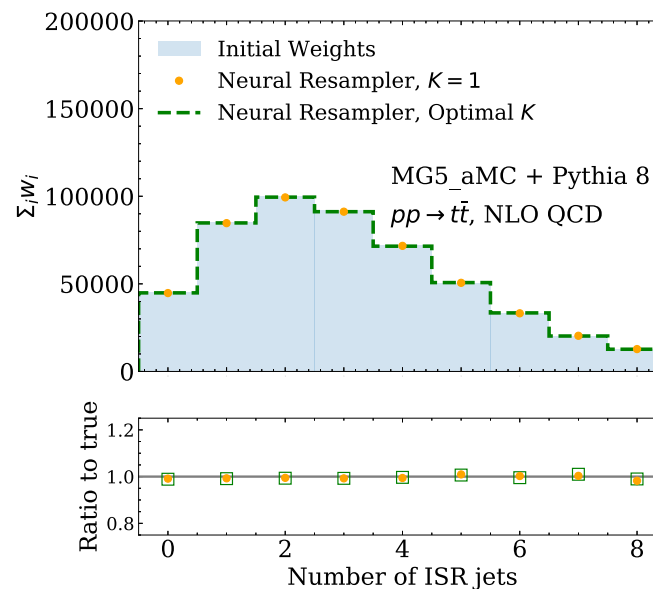[Nachman, JDT, PRD 2020; building on Andersen, Gutschow, Maier, Prestel, EPJC 2020]

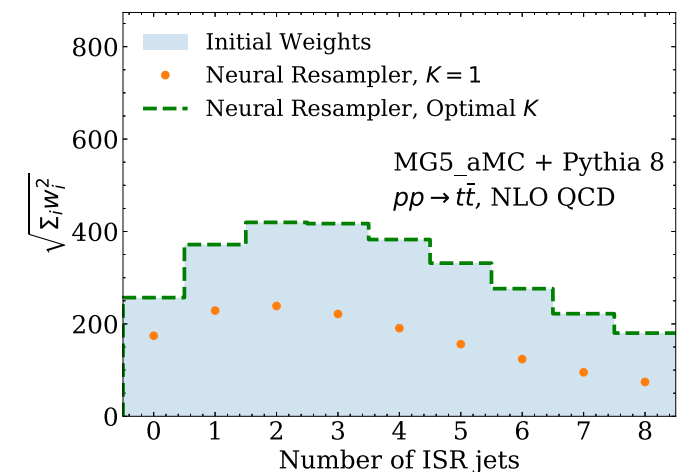## Case Study in Jet Physics at Large Hadron Collider

**Original sample: large weight cancellations**



**Reweighting recovers desired distribution**



**Resampling recovers desired uncertainties**



**Improved resampling** through auxiliary neural network yields correct statistical uncertainties

$$\left(\frac{\delta p}{p}\right)^2 = \frac{\langle w^2 \rangle}{\langle w \rangle^2}$$

[Nachman, JDT, PRD 2020; building on Andersen, Gutschow, Maier, Prestel, EPJC 2020]

# Theory Perspective on ML for Jets

When striving for "interpretable machine learning"
we are essentially hoping that likelihood ratios can be
approximated via theoretically well-motivated forms

We can impose theoretical priors by judicious choice of
network architecture that captures the underlying
structures and symmetries of our problem

Machine learning methods can only be as robust and reliable
as the data samples used for training

We are making progress towards uncertainty quantification,
using more elaborate loss functions and training paradigms

*Looking forward to your thoughts and discussion!*

# *Backup Slides*

# Using All *Five* LHC Interaction Points



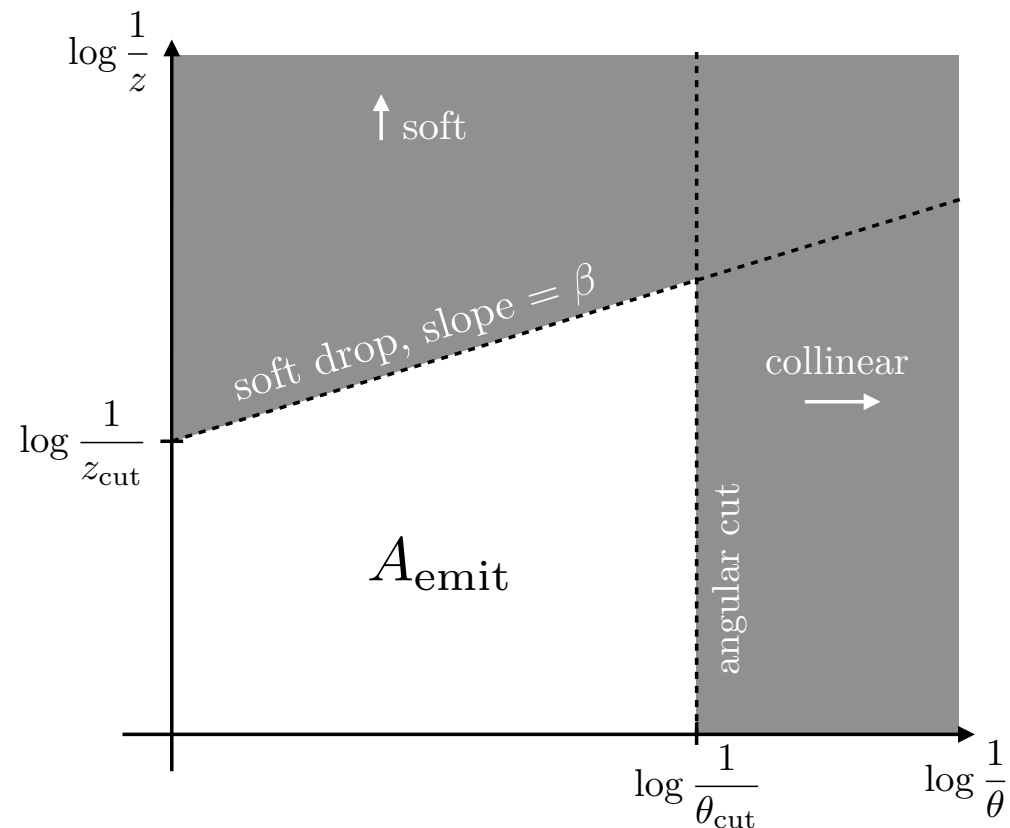P1          P2          P5          P8          R1
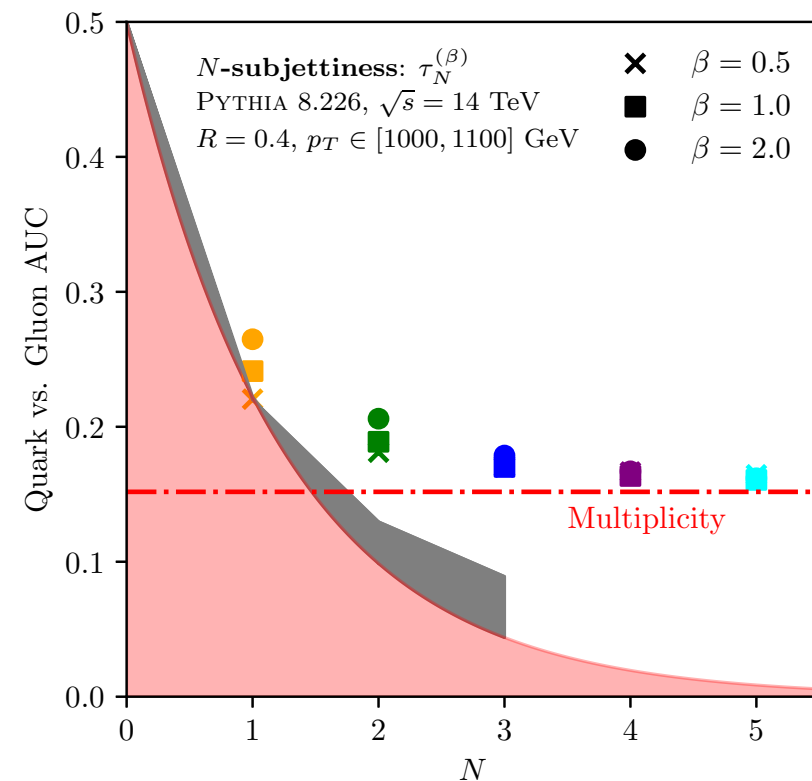
# E.g. Quark/Gluon Classification

*In various limits, likelihood ratio is monotonically related to…*

## Soft-dropped Multiplicity



[Frye, Larkoski, JDT, Zhou, JHEP 2017]

## IRC Safe Multiplicity



$N$-subjettiness: $\tau_N^{(\beta)}$
PYTHIA 8.226, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [1000, 1100]$ GeV

$\beta = 0.5$
$\beta = 1.0$
$\beta = 2.0$
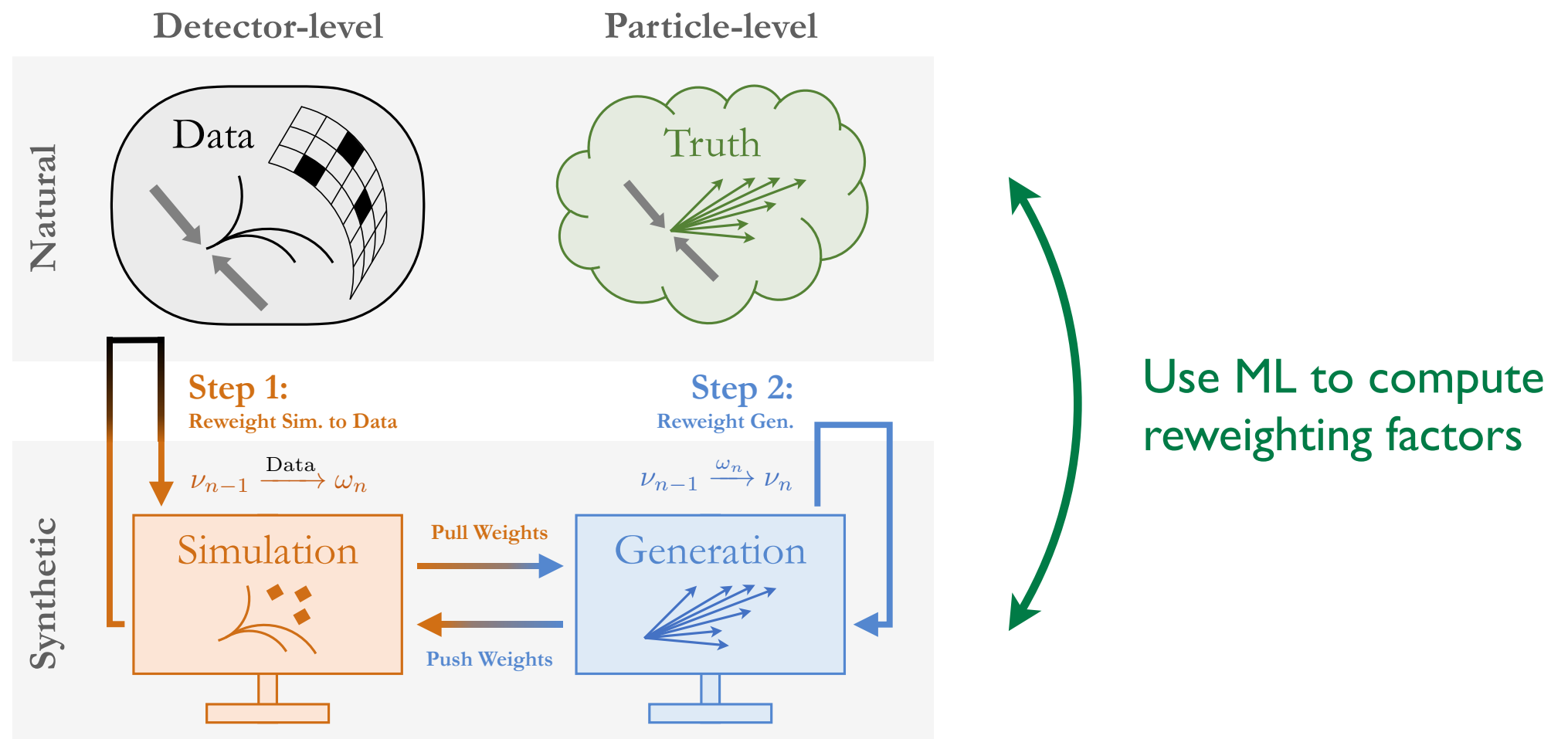
Multiplicity

[Larkoski, Metodiev, JHEP 2019]

Away from these limits, the likelihood ratio is
**not typically simple, elegant, or interpretable** (but we can hope!)
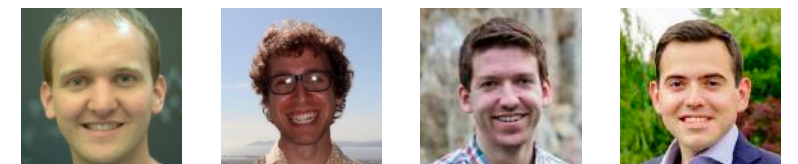
# E.g. Detector Unfolding

**OmniFold**

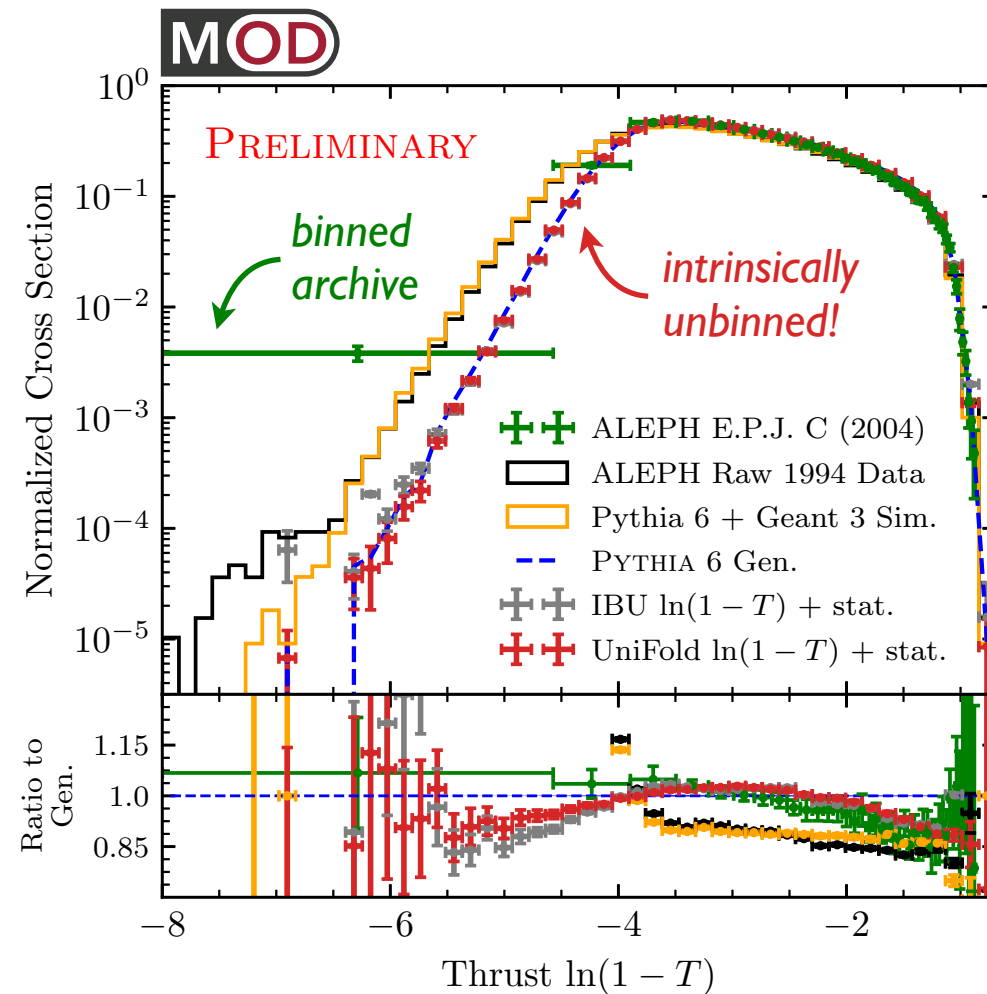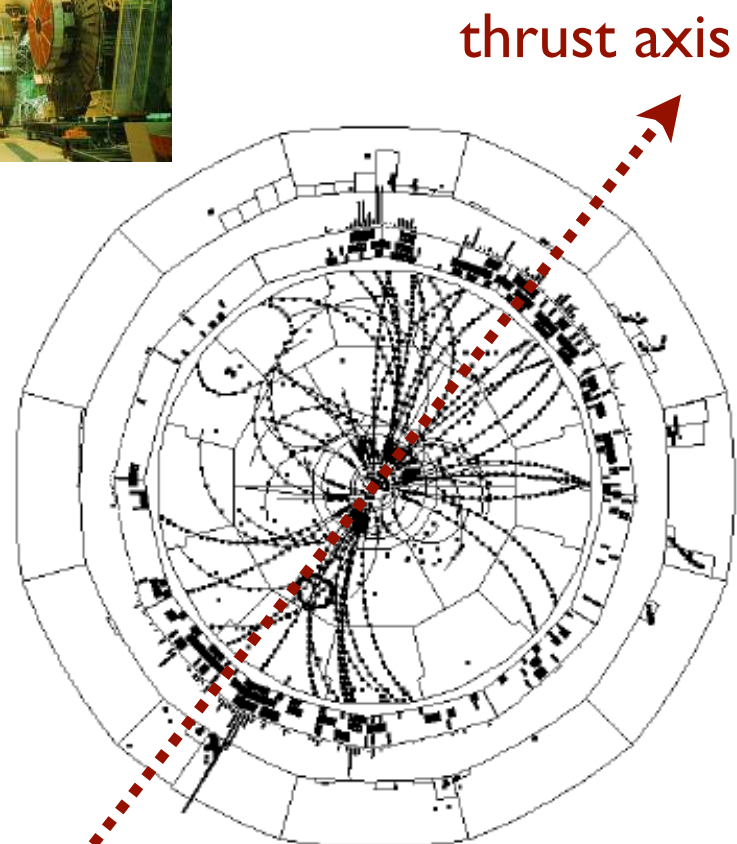*Multi-dimensional unbinned detector corrections via iterated application of likelihood ratio trick*



Detector-level    Particle-level

Natural

Data    Truth

Step 1:
Reweight Sim. to Data

Step 2:
Reweight Gen.

$\nu_{n-1} \xrightarrow{\text{Data}} \omega_n$    $\nu_{n-1} \xrightarrow{\omega_n} \nu_n$

Synthetic

Simulation    Generation

Pull Weights

Push Weights

Use ML to compute reweighting factors

[Andreassen, Komiske, Metodiev, Nachman, JDT, PRL 2020; + Suresh, ICLR SimDL 2021;
Komiske, McCormack, Nachman, arXiv 2021; see unfolding comparison in Petr Baron, arXiv 2021]

# E.g. Detector Unfolding
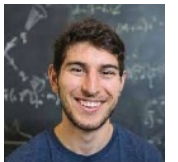


**Back to the Future with ALEPH Archival Data**

thrust axis

PRELIMINARY

*binned archive*

*intrinsically unbinned!*

ALEPH E.P.J. C (2004)
ALEPH Raw 1994 Data
Pythia 6 + Geant 3 Sim.
PYTHIA 6 Gen.
IBU $\ln(1-T)$ + stat.
UniFold $\ln(1-T)$ + stat.

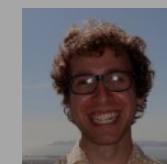Normalized Cross Section

Ratio to Gen.

Thrust $\ln(1-T)$

[talk by Badea, ICHEP 2020; cf. ALEPH, EPJC 2004]
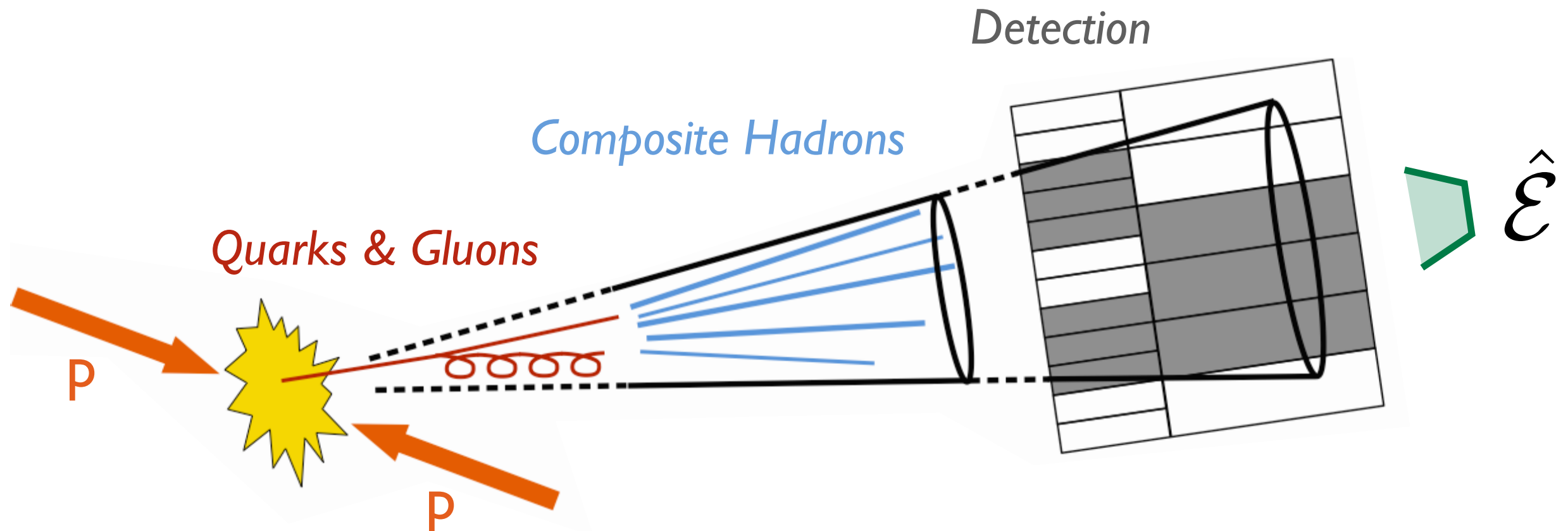[see also Badea, Baty, Chang, Innocenti, Maggi, McGinn, Peters, Sheng, JDT, Lee, PRL 2019; H1, DIS2021]

[Andreassen, Komiske, Metodiev, Nachman, JDT, PRL 2020; + Suresh, ICLR SimDL 2021;
Komiske, McCormack, Nachman, arXiv 2021; see unfolding comparison in Petr Baron, arXiv 2021]

# Energy Flow Representation

*Emphasizes* *infrared and collinear safety*

*Theory*

*Detection*

*Composite Hadrons*

*Quarks & Gluons*

P

P

$\hat{\mathcal{E}}$

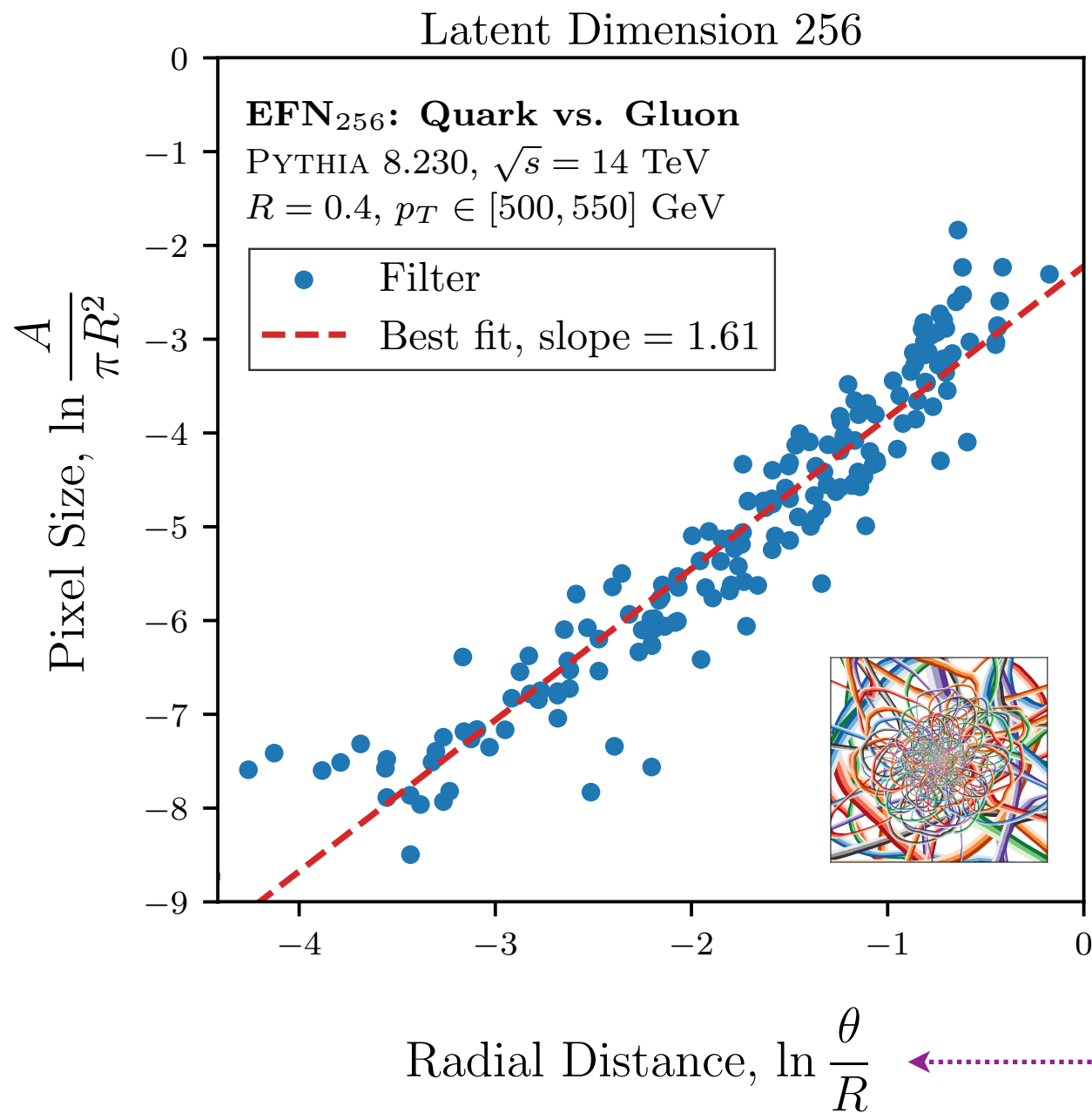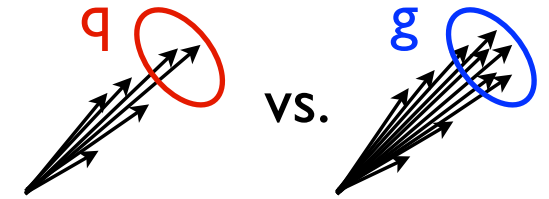## Energy Flow:

Robust to hadronization and detector effects
Well-defined for massless gauge theories

$$\hat{\mathcal{E}} \simeq \lim_{t \to \infty} \hat{n}_i T^{0i}(t, vt\hat{n})$$

[see e.g. Sveshnikov, Tkachov, <u>PLB 1996</u>; Hofman, Maldacena, <u>JHEP 2008</u>; Mateu, Stewart, JDT, <u>PRD 2013</u>;
Belitsky, Hohenegger, Korchemsky, Sokatchev, Zhiboedov, <u>PRL 2014</u>; Chen, Moult, Zhang, Zhu, <u>PRD 2020</u>]
[complementary perspective on IRC unsafe information in Chakraborty, Lim, Nojiri, Takeuchi, <u>JHEP 2020</u>]

# Machine Learning Collinear QCD
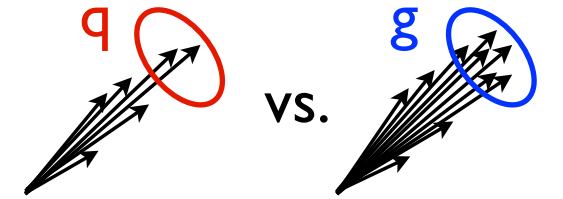


Latent Dimension 256

**EFN$_{256}$: Quark vs. Gluon**
PYTHIA 8.230, $\sqrt{s} = 14$ TeV
$R = 0.4$, $p_T \in [500, 550]$ GeV

- Filter
- - - Best fit, slope = 1.61

Pixel Size, $\ln \frac{A}{\pi R^2}$

Radial Distance, $\ln \frac{\theta}{R}$

$C_q = 4/3$
$C_g = 3$

$\theta$
$z$

$$\mathrm{d}P_{i \to ig} \simeq \frac{2\alpha_s}{\pi} C_i \frac{\mathrm{d}\theta}{\theta} \frac{\mathrm{d}z}{z}$$
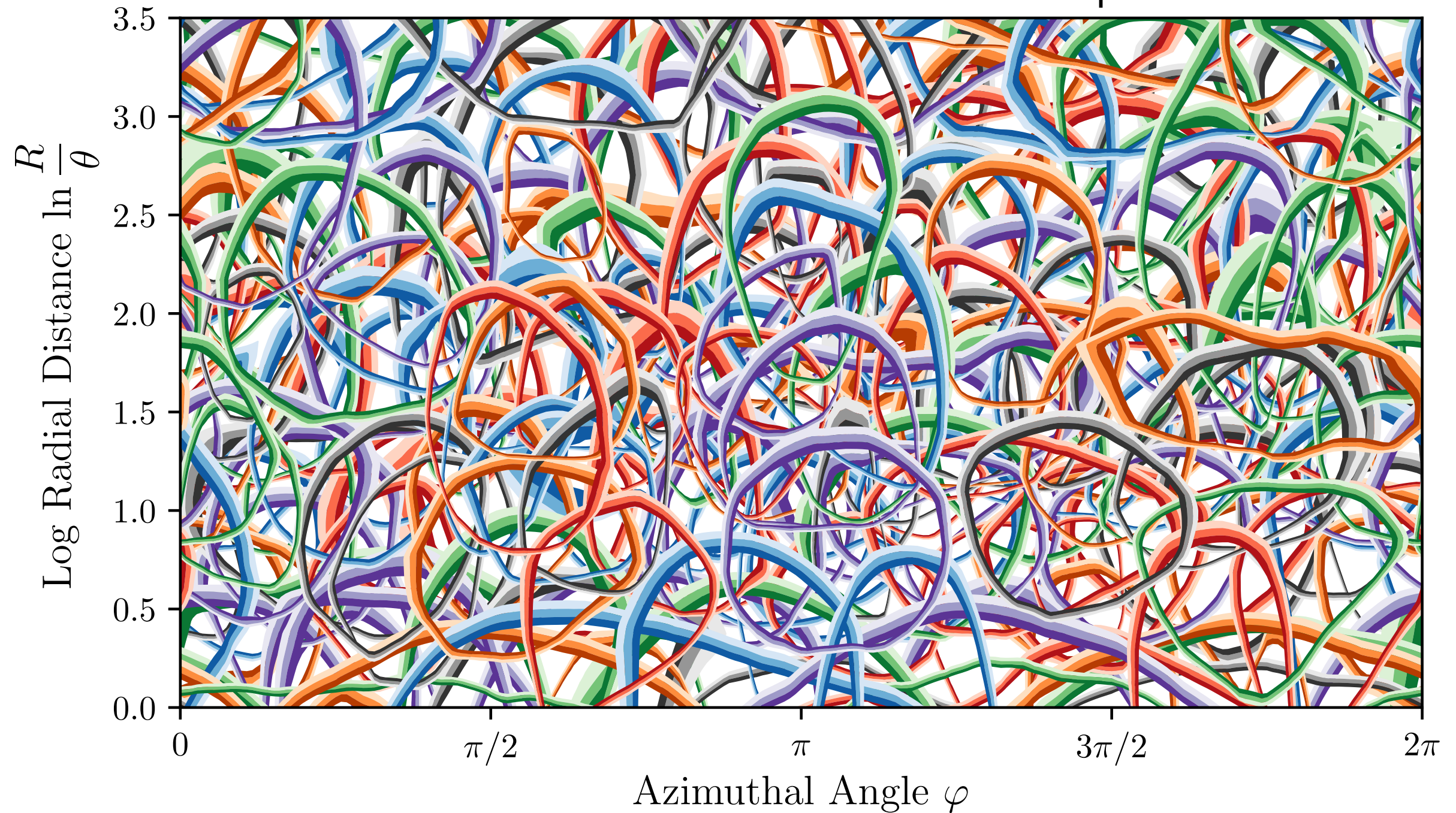
Collinear  Soft

q  vs.  g

[Komiske, Metodiev, JDT, JHEP 2019]

# En Route to the Lund Plane



Coordinate transformation to the emission plane

[Komiske, Metodiev, JDT, JHEP 2019; see also Dreyer, Salam, Soyez, JHEP 2018]