

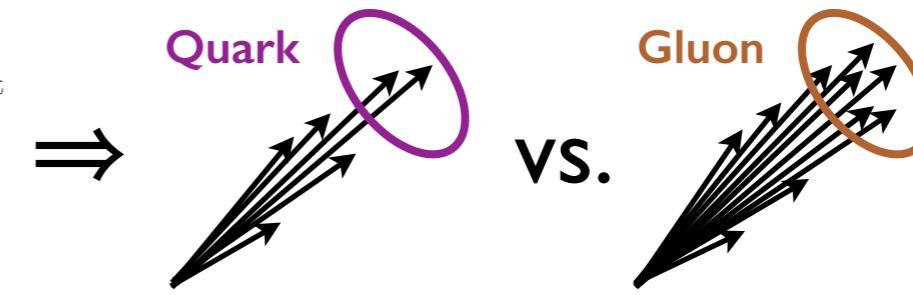
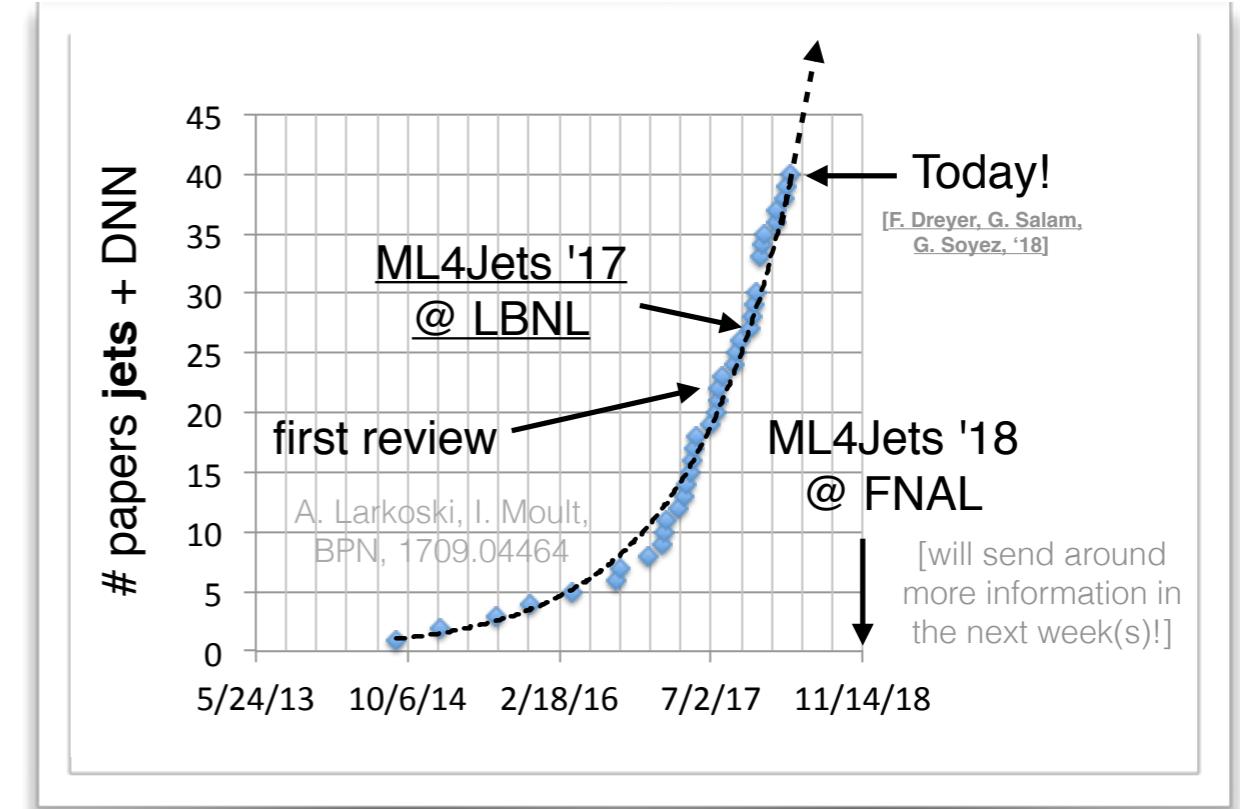
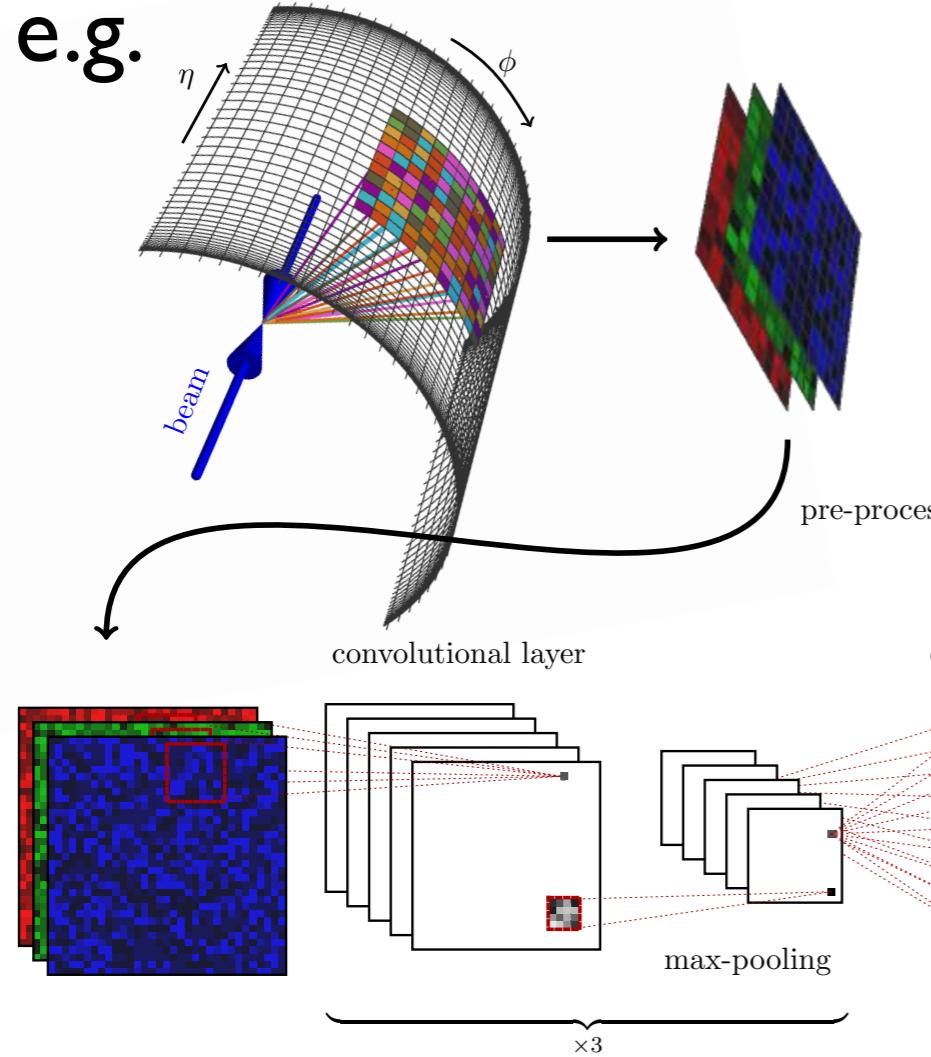
# Deep Sets for Particle Jets

Jesse Thaler



LCTP Particle Theory Seminar, U. Michigan — November 9, 2018

# The Rise of Machine Learning for Jets



[e.g. Komiske, Metodiev, Schwartz, 1612.01551; Nachman, Boost 2018 Talk, July 20, 2018; reviews in Larkoski, Moult, Nachman, 1709.04464; Guest, Cranmer, Whiteson, 1806.11484]

# My Perspective c. 2016



“Deep Learning”      vs.      “Deep Thinking”

# My Perspective c. 2018

BOOST 2018

10<sup>th</sup> International Workshop on Boosted Objects  
Phenomenology, Reconstruction and Searches

“Deep Learning”

&

~~vs.~~

“Deep Thinking”

*New first-principles studies of QCD  
facilitated by advances in  
mathematics, statistics, and computer science*

Desired Outcomes  $\Leftrightarrow$  Algorithms/Observables

# Proximate Reasons for My Conversion



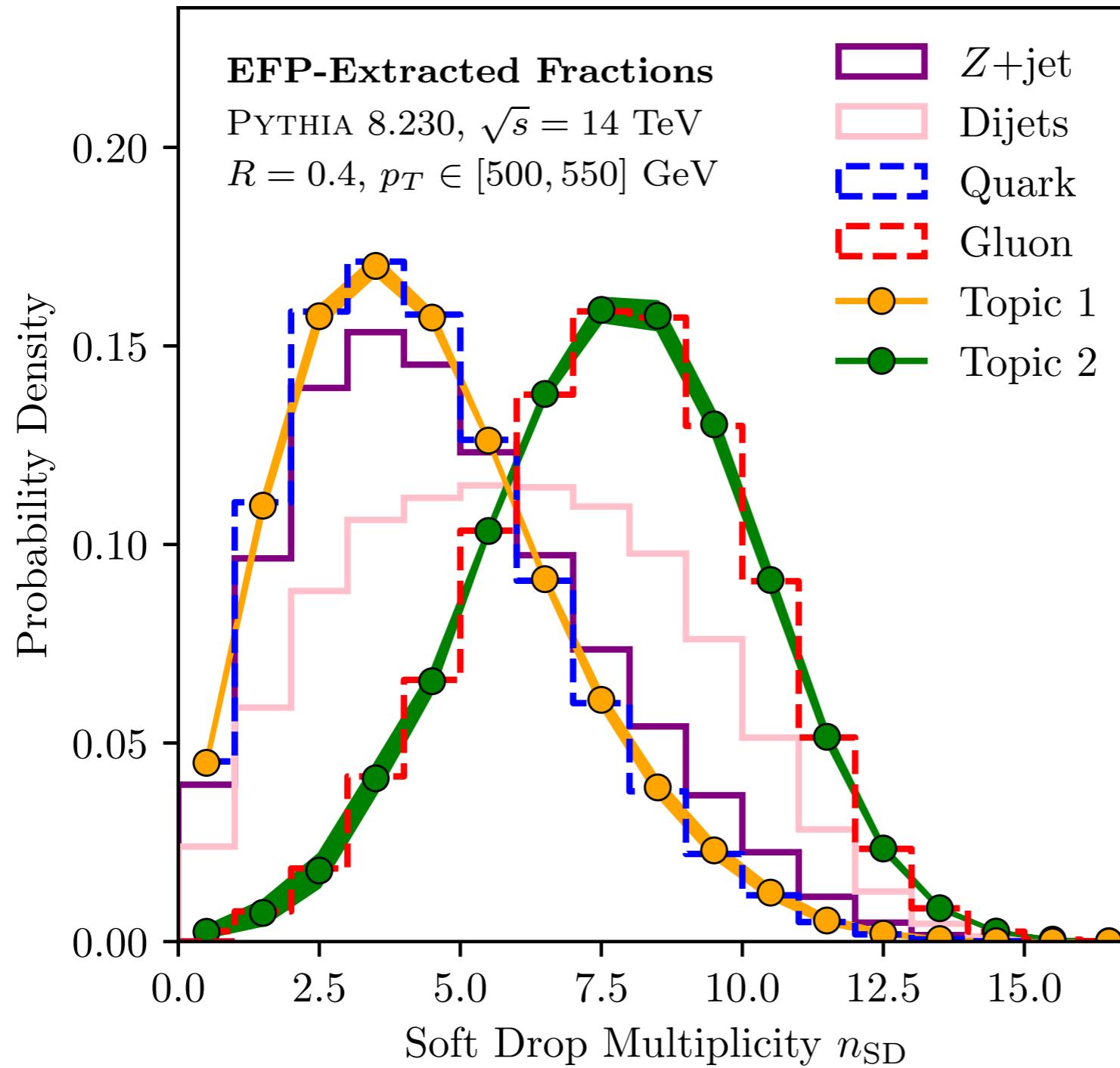
Patrick Komiske



Eric Metodiev

*plus Ben Nachman, Kyle Cranmer, Daniel Whiteson, Mike Williams, Matt Schwartz, Dan Roberts, Phiala Shanahan, ...*

# Physics Reasons for My Conversion (I)



For Offline Discussions

First-principles QCD meets blind source separation with same underlying structure as...

[Komiske, Metodiev, JDT, 1809.01140;  
see also Metodiev, JDT, 1802.00008; Metodiev, Nachman, JDT, 1708.02949]

# Topic Modeling

## Topics

gene	0.04
dna	0.02
genetic	0.01
...	

life	0.02
evolve	0.01
organism	0.01
...	

brain	0.04
neuron	0.02
nerve	0.01
...	

data	0.02
number	0.02
computer	0.01
...	

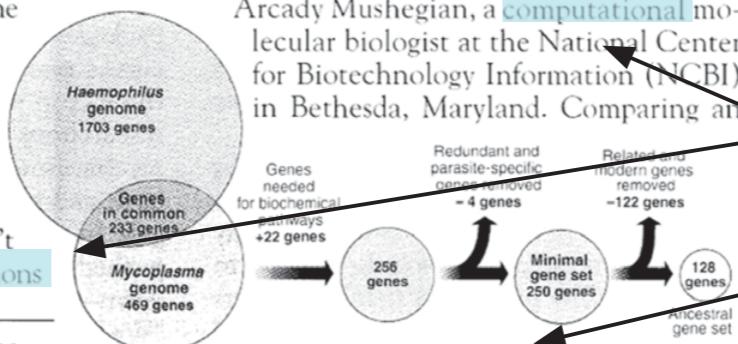
## Documents

### Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

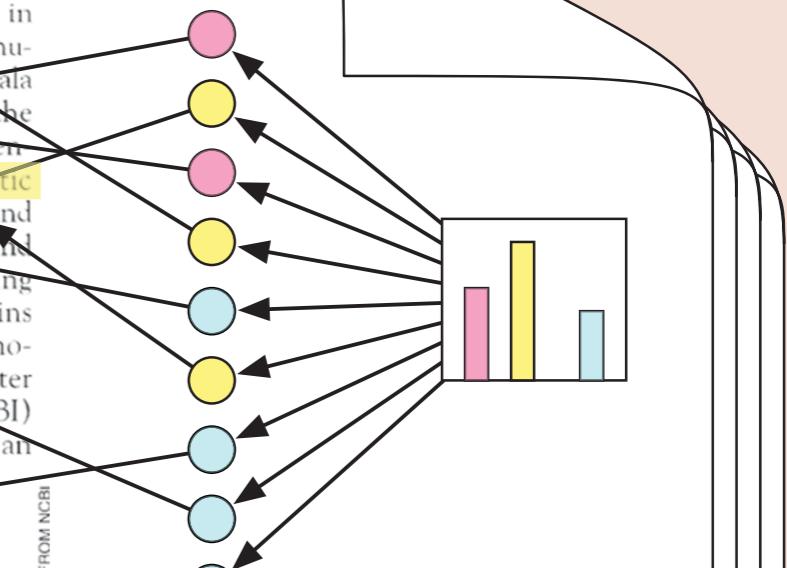
"are not all that far apart," especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. "It may be a way of organizing any newly sequenced genome," explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

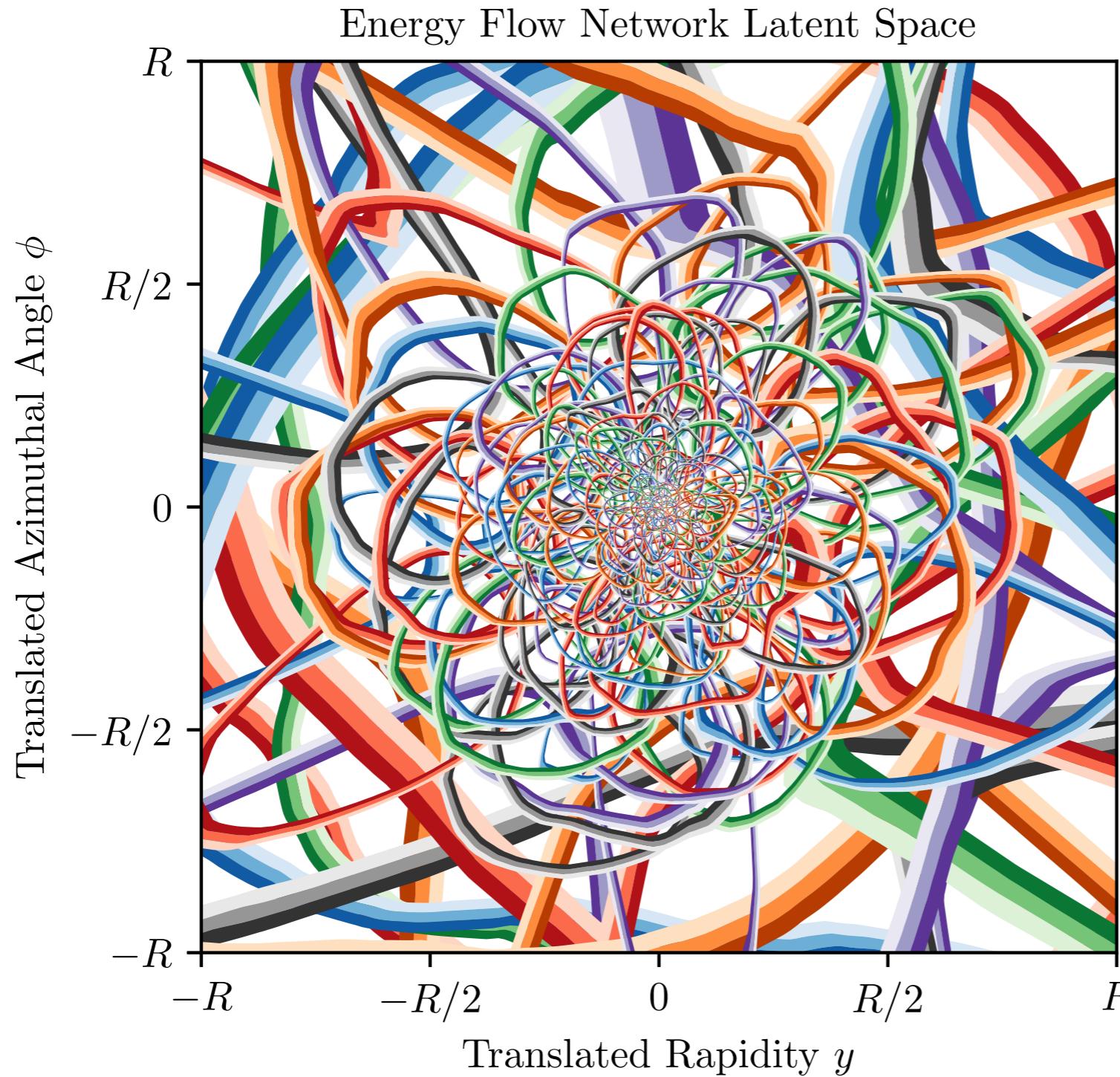
SCIENCE • VOL. 272 • 24 MAY 1996

## Topic proportions and assignments



[Blei, 2012]

# Physics Reasons for My Conversion (2)



Today's Talk

First-principles QCD  
meets neural networks  
with same underlying  
symmetries as...

[Komiske, Metodiev, JDT, 1810.05165;  
see also Komiske, Metodiev, JDT, 1712.07124]

# Point Clouds

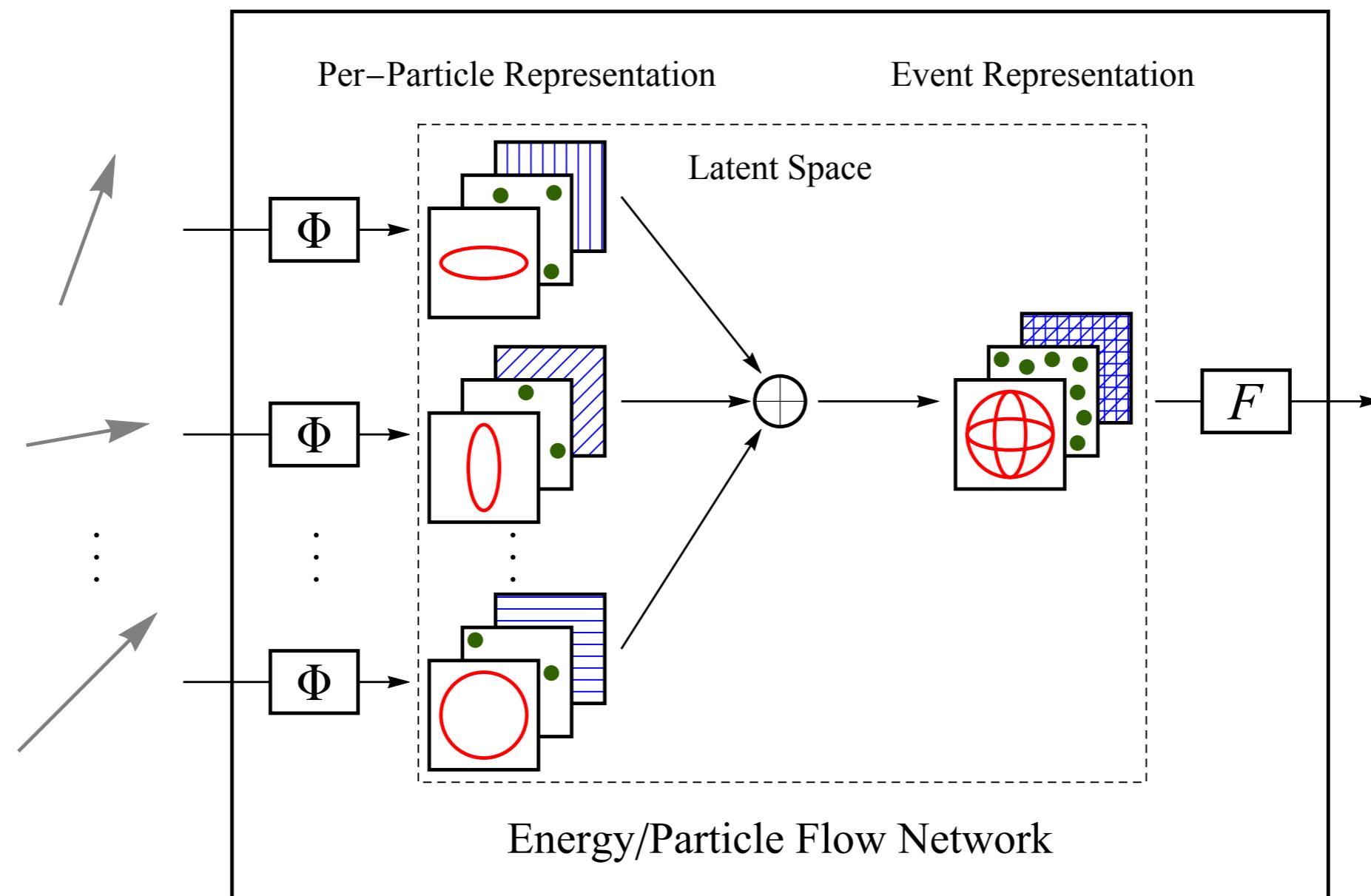


[Popular Science, 2013]

# Introducing Energy Flow Networks

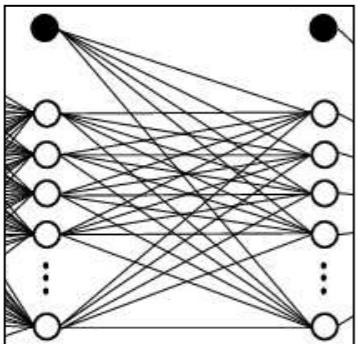
Particles

Observable

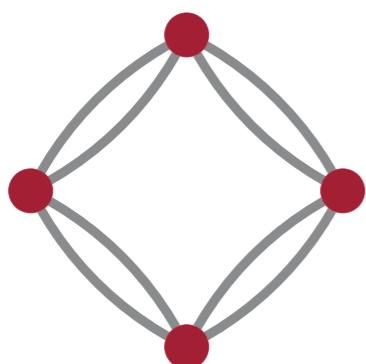


[Komiske, Metodiev, JDT, 1810.05165]

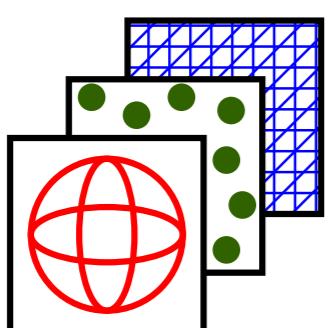
# Outline



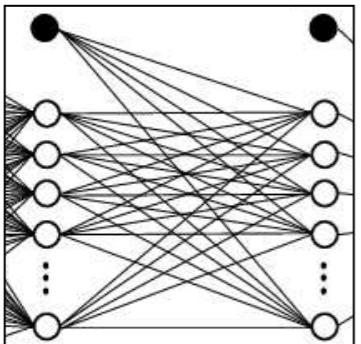
Into the Network



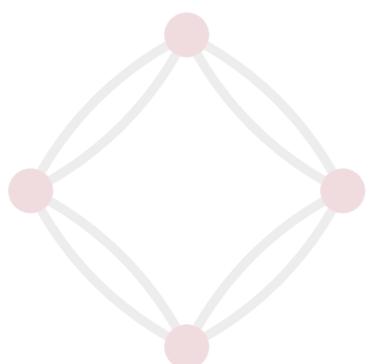
Symmetries & Safety



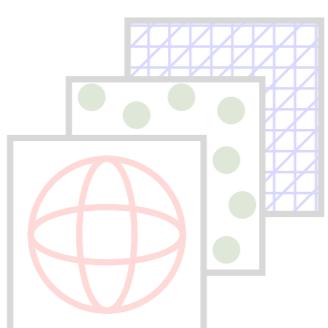
Deep Sets for Particle Jets



## Into the Network



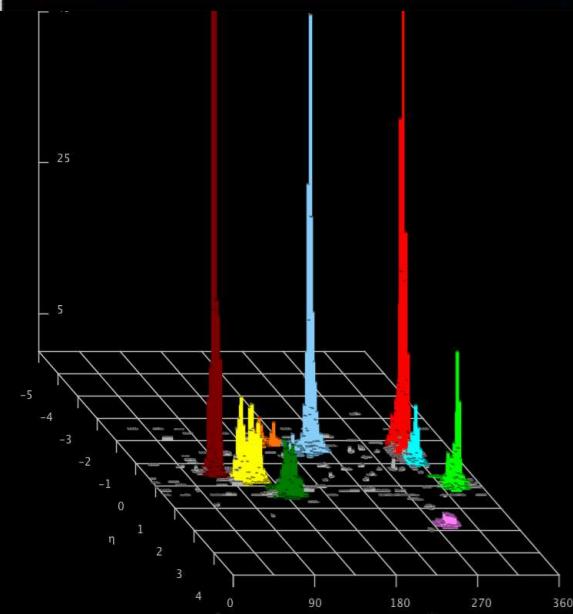
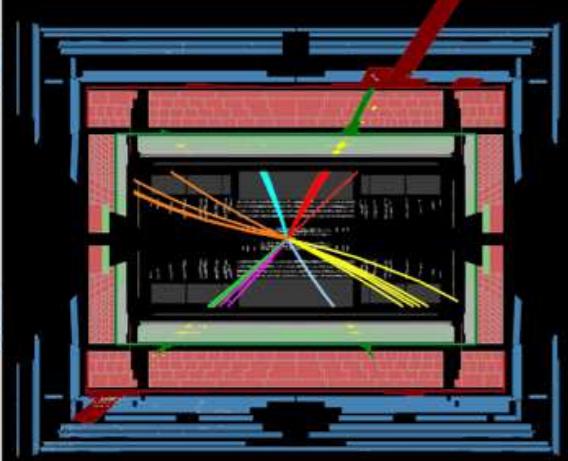
## Symmetries & Safety



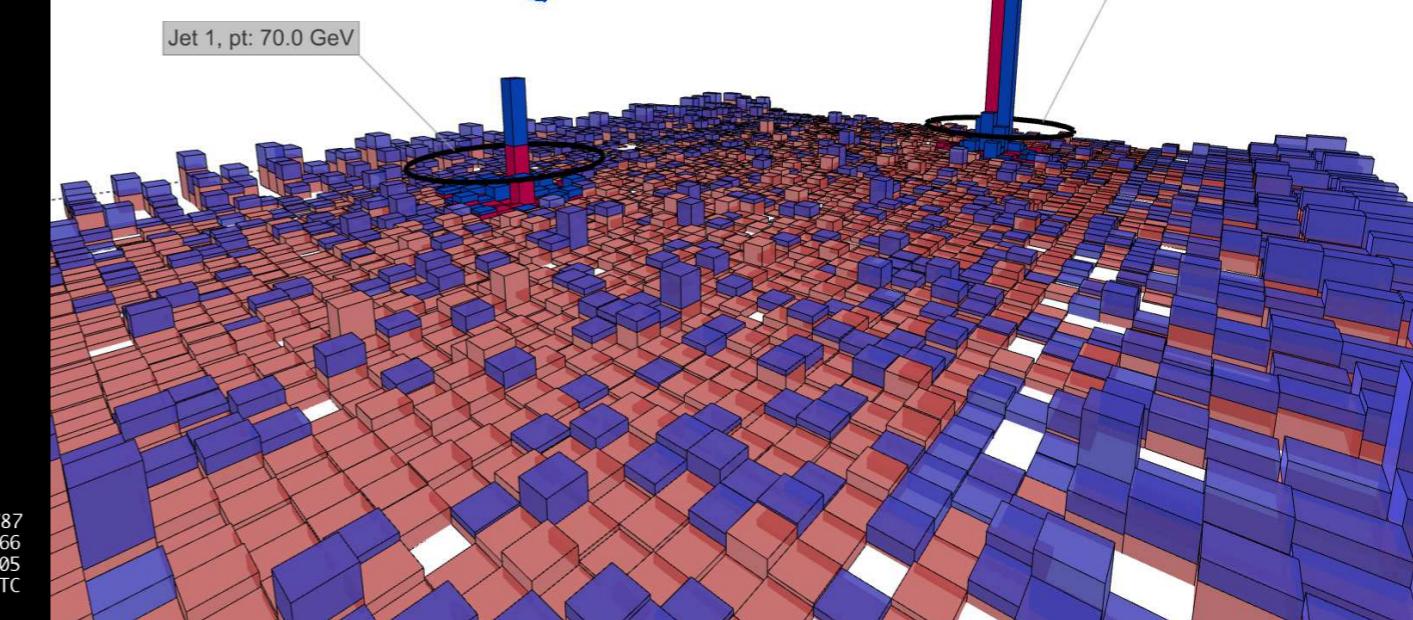
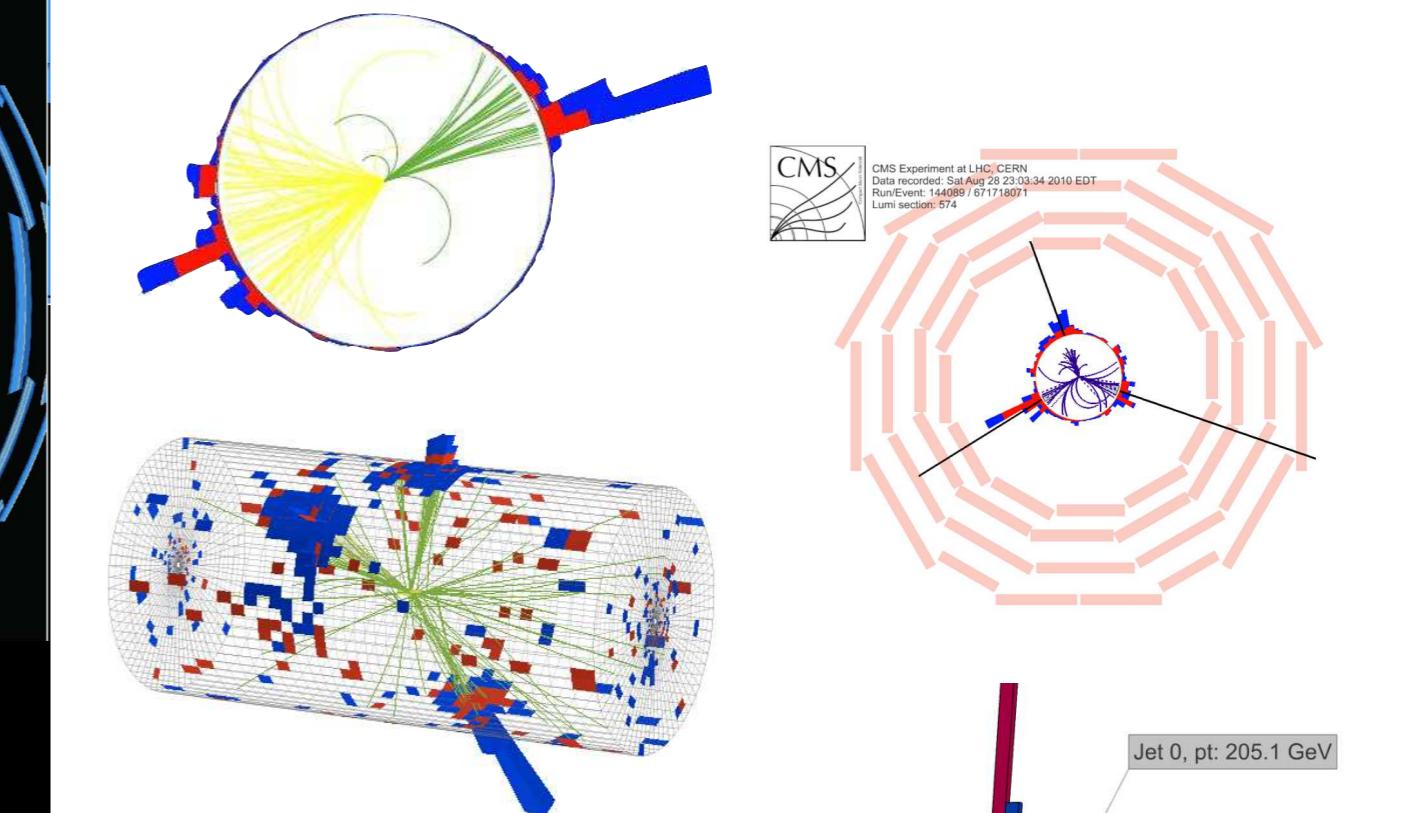
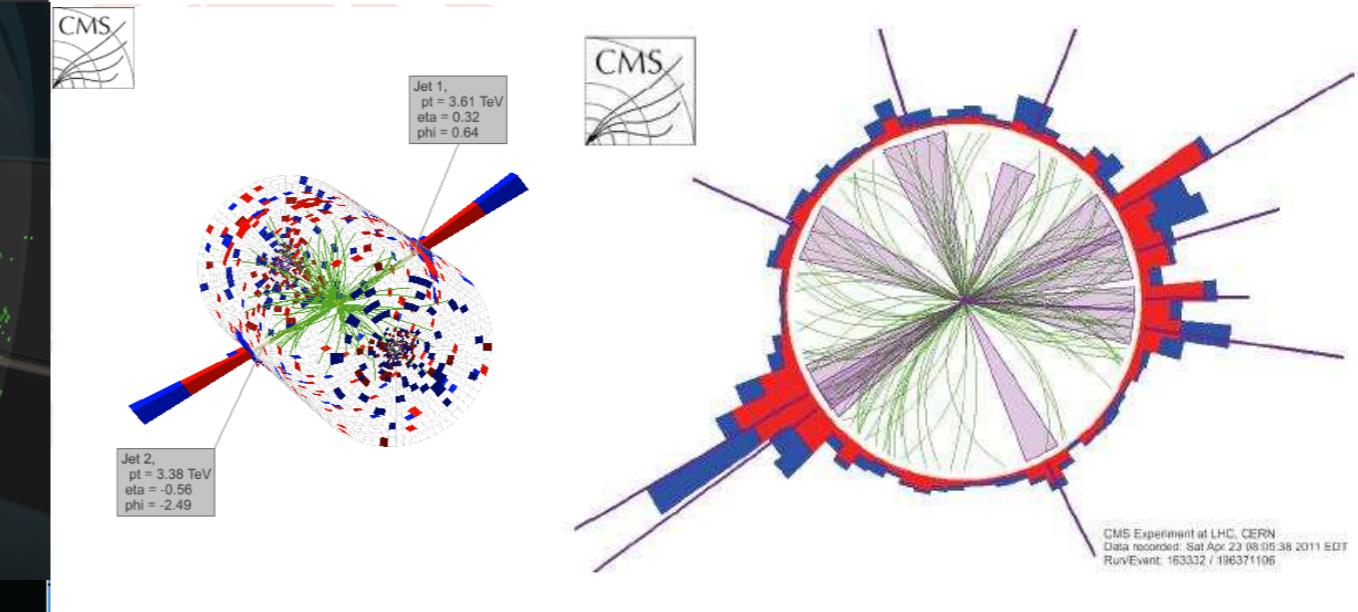
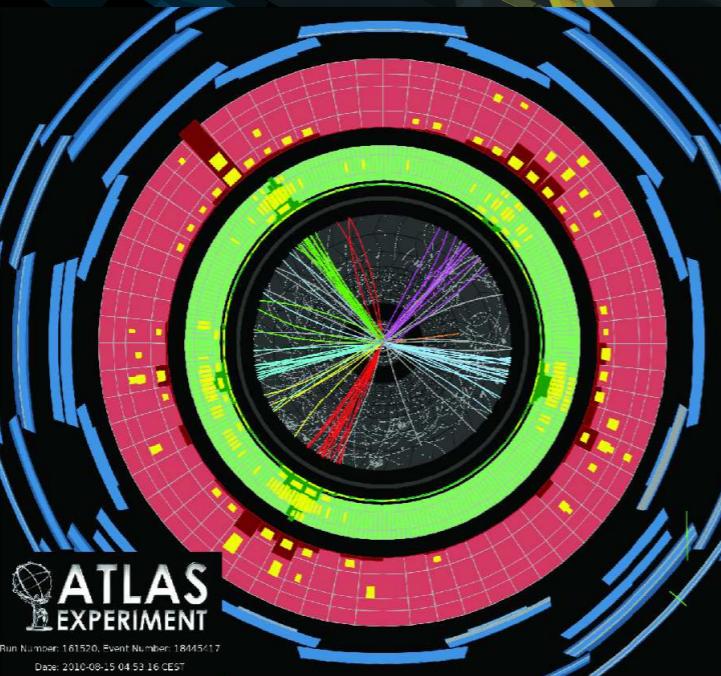
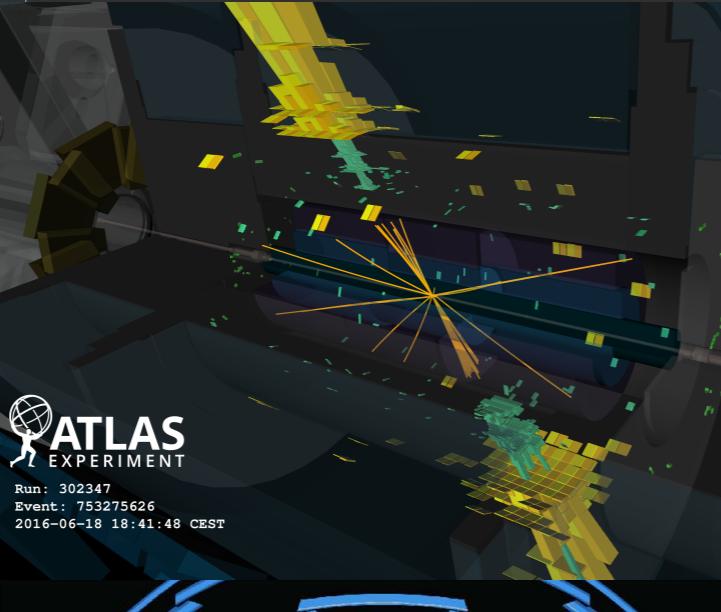
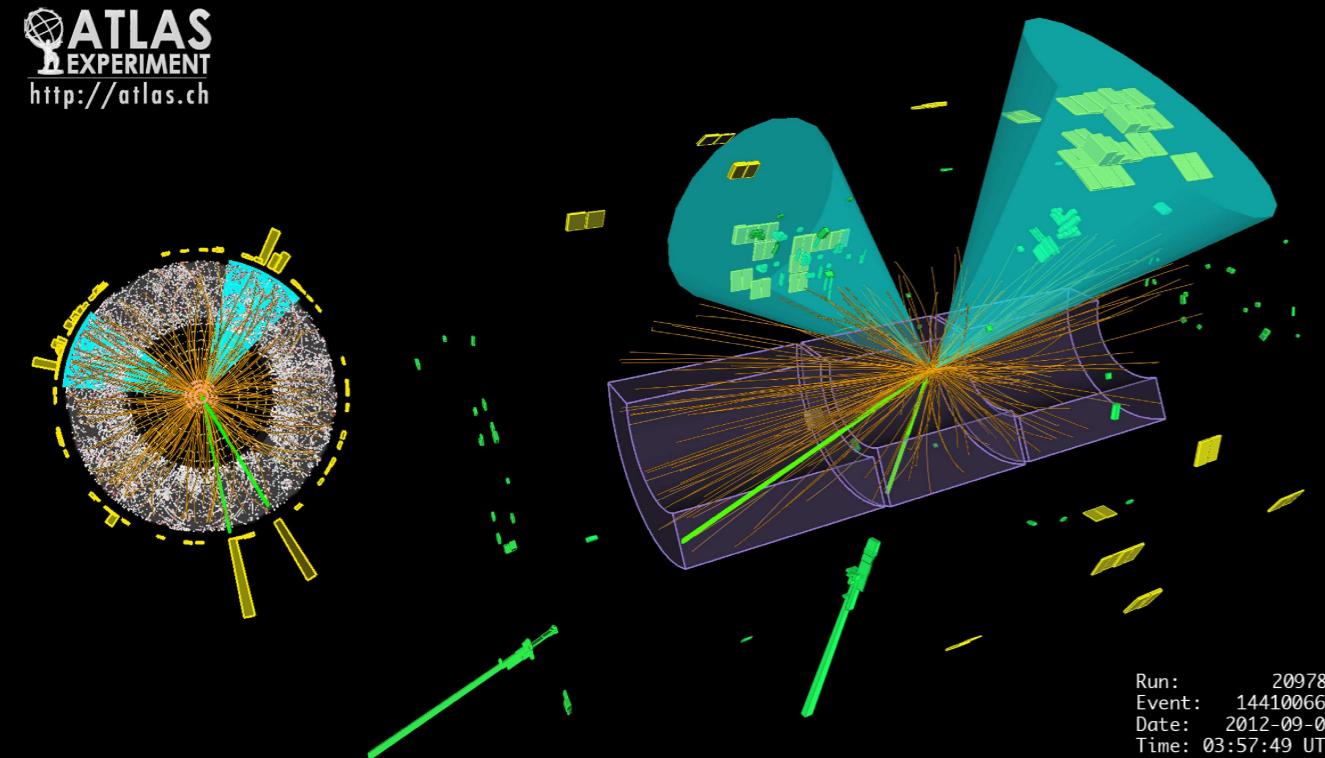
## Deep Sets for Particle Jets

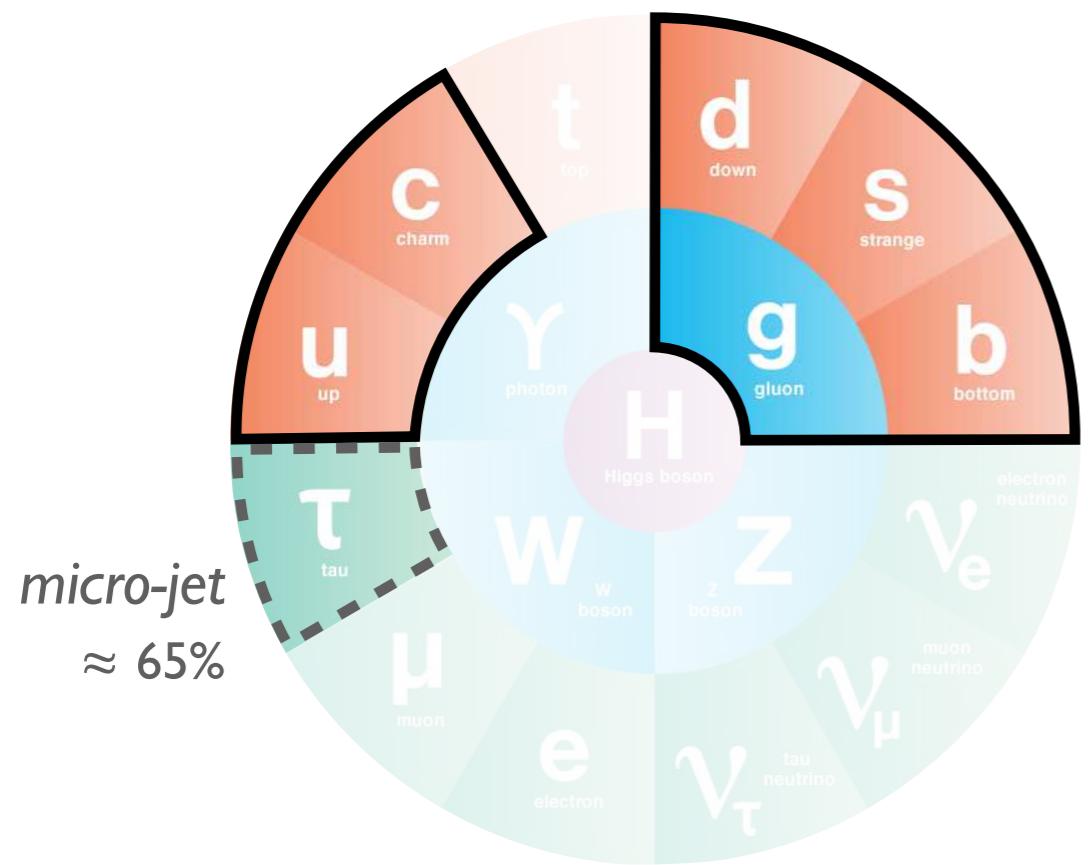
Run Number: 159224, Event Number: 3533152

Date: 2010-07-18 11:05:54 CEST



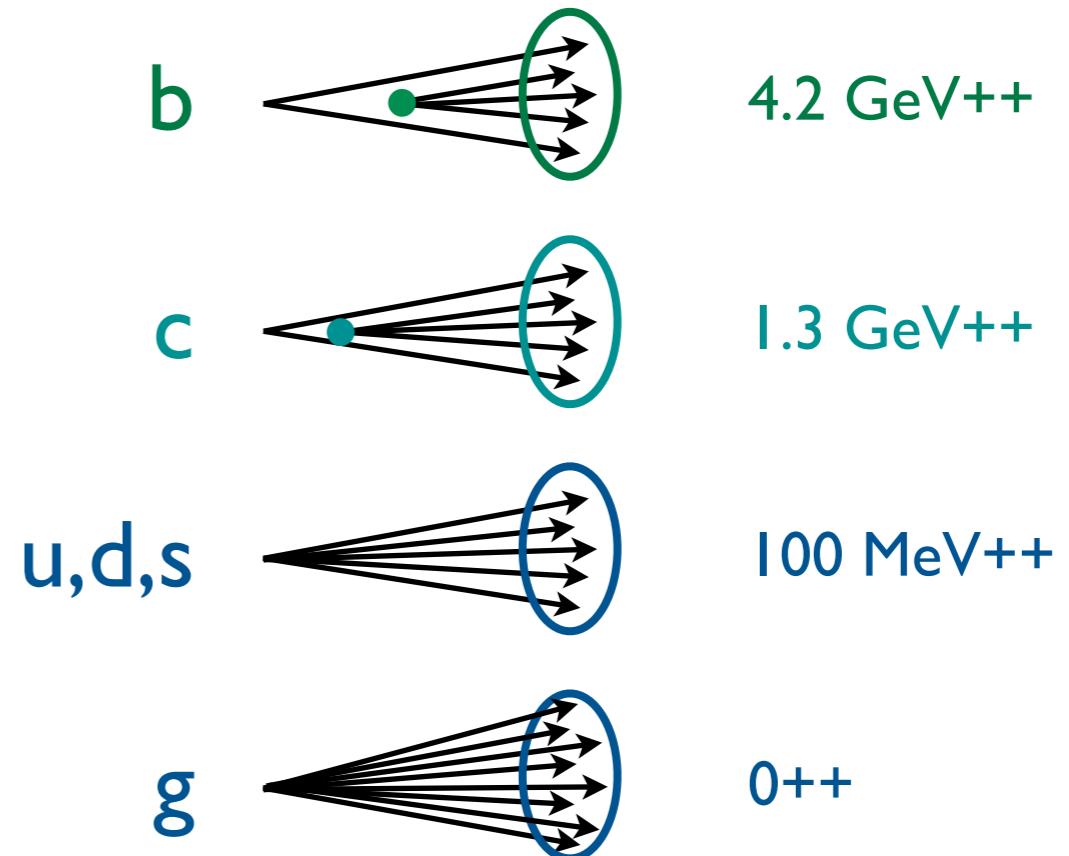
**ATLAS**  
EXPERIMENT  
<http://atlas.ch>

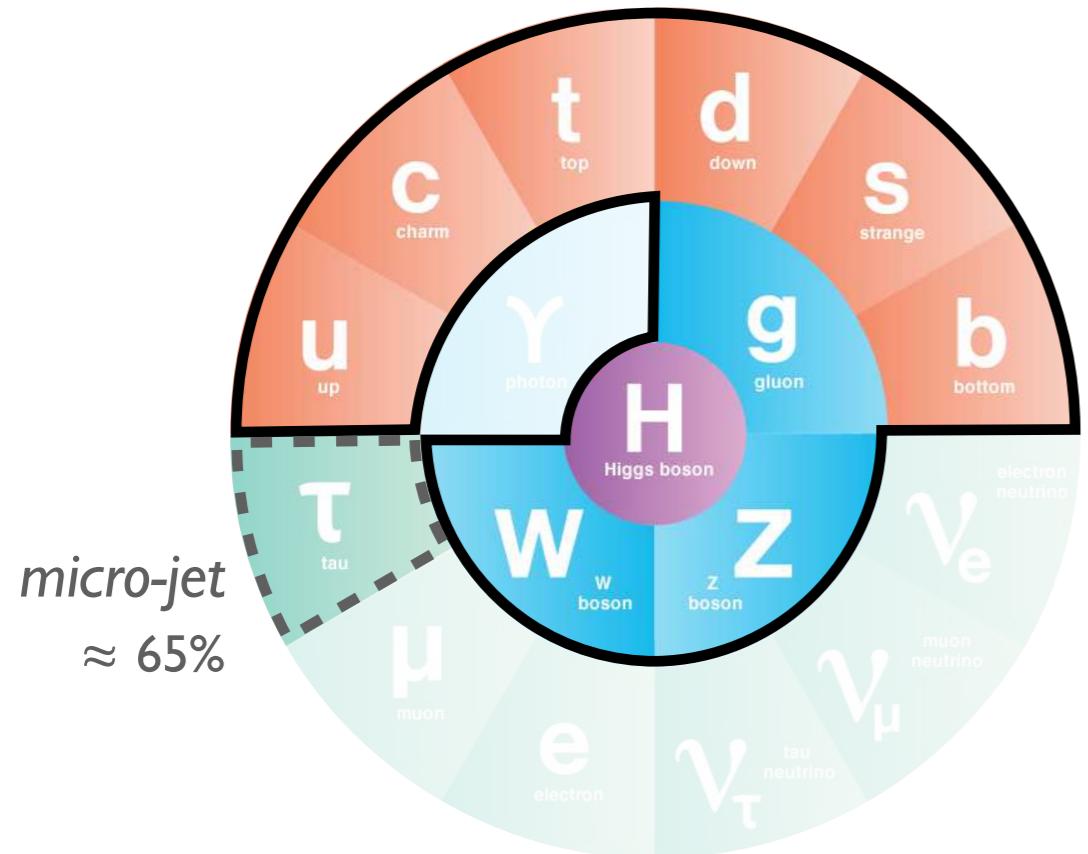




## *Jets from the Standard Model*

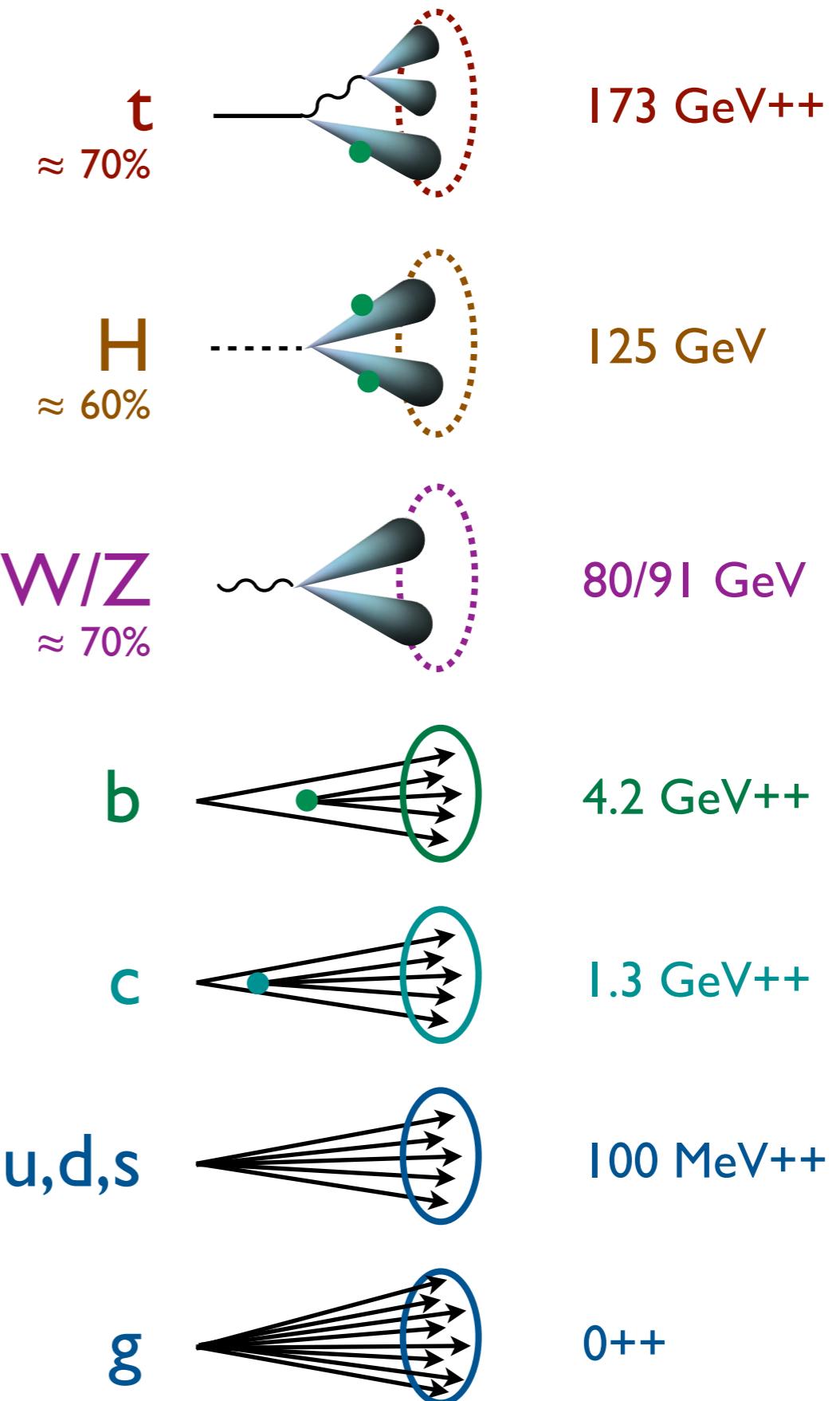
$\text{++}$  = Mass from QCD Radiation

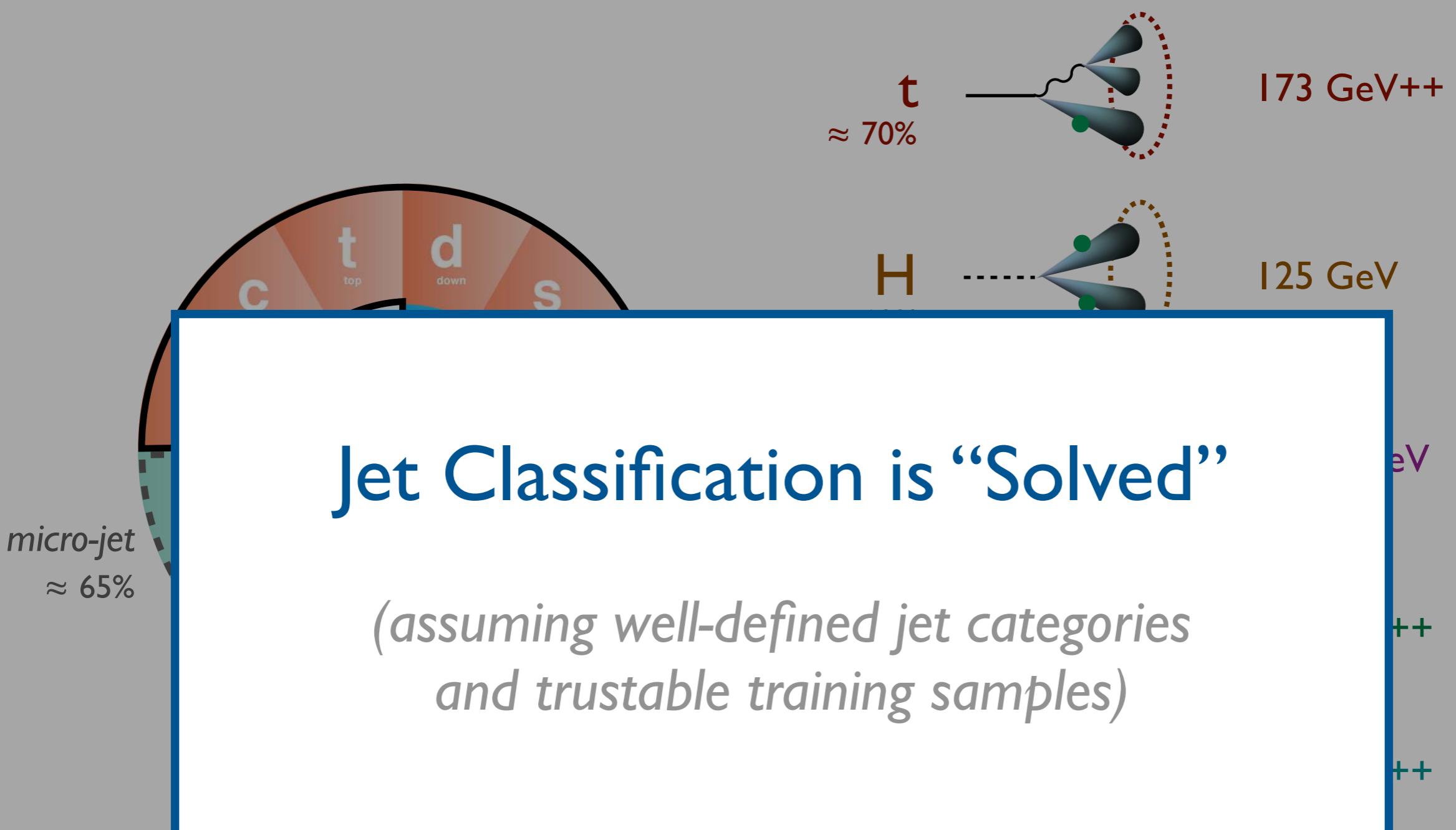




## *Jets from the Standard Model*

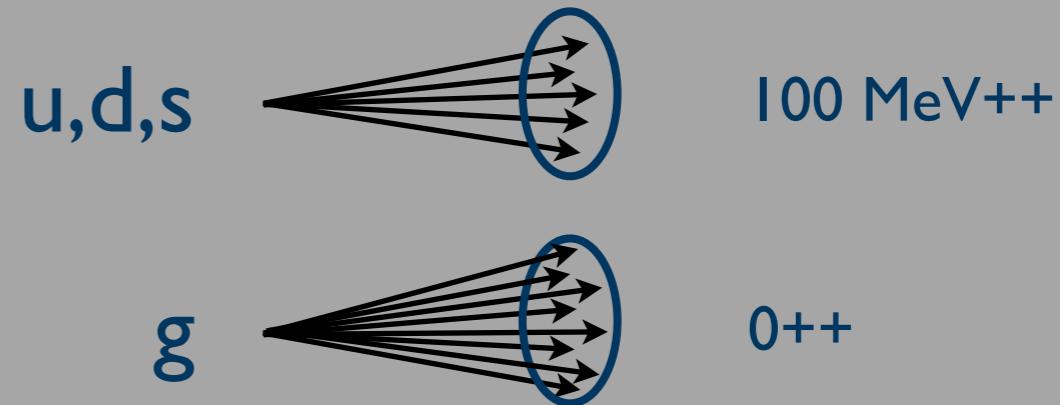
++ = Mass from QCD Radiation

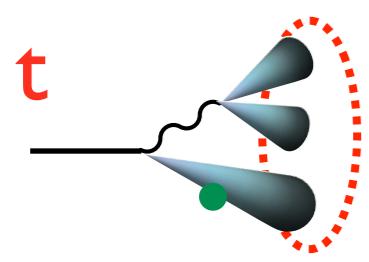




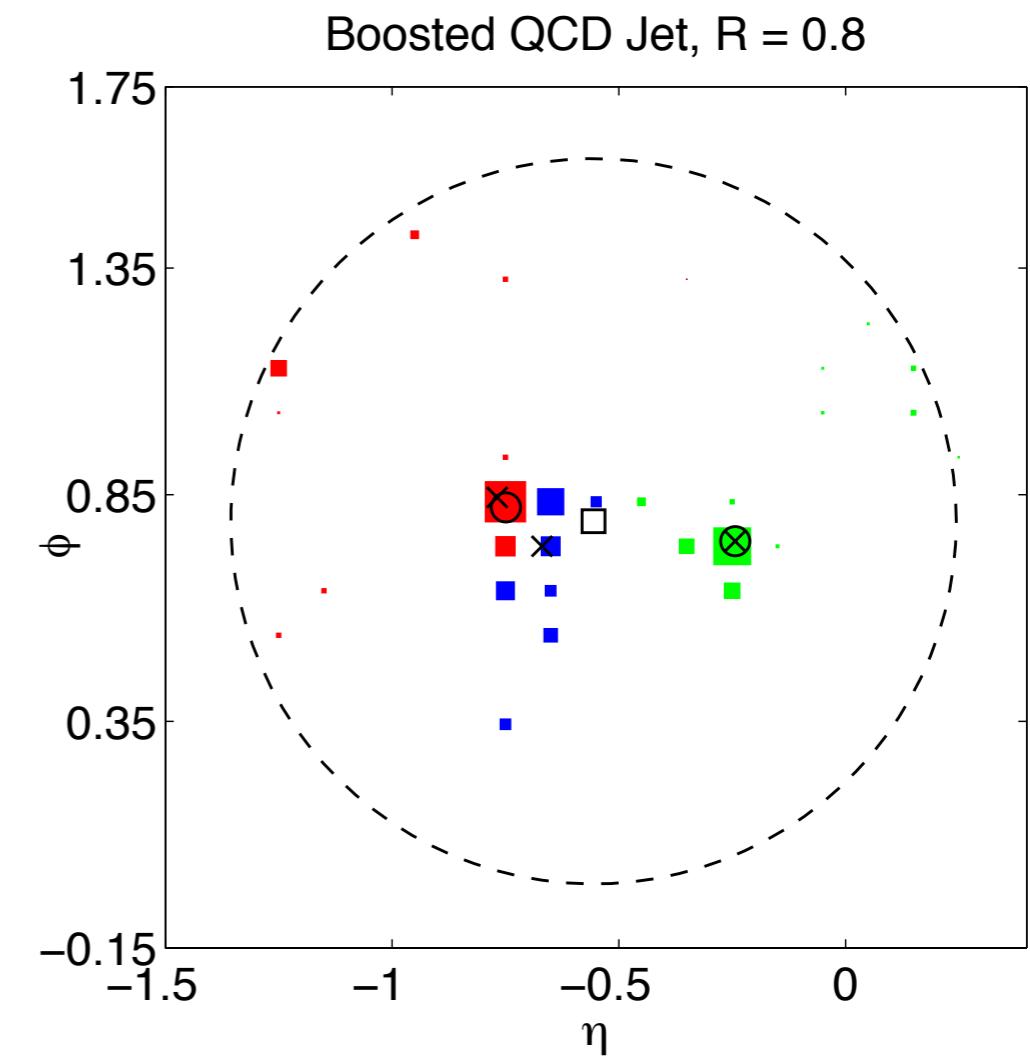
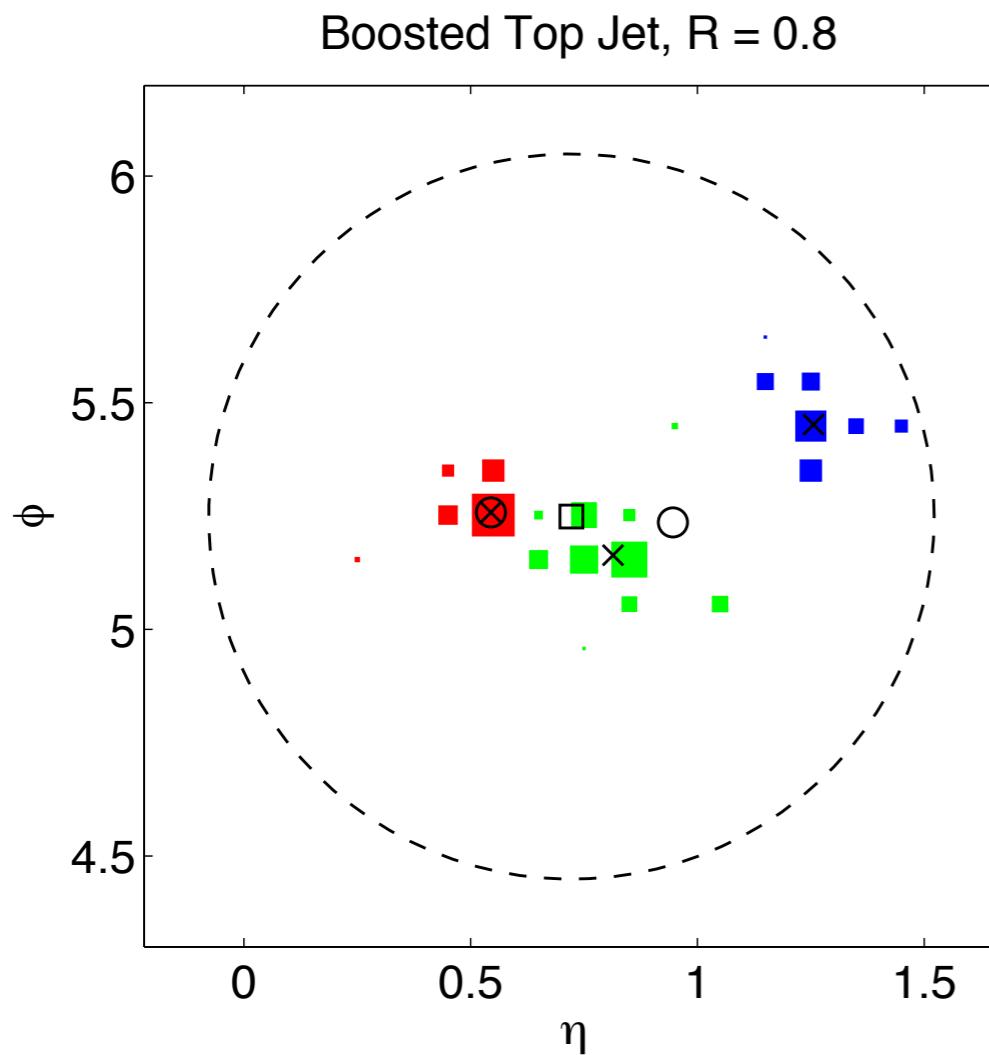
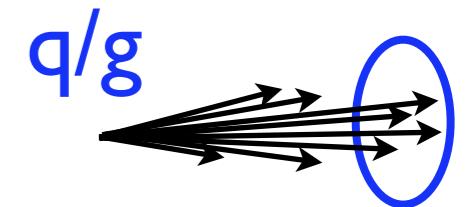
## Standard Model

++ = Mass from QCD Radiation

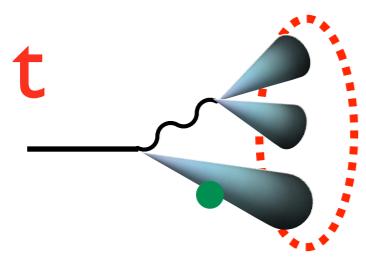




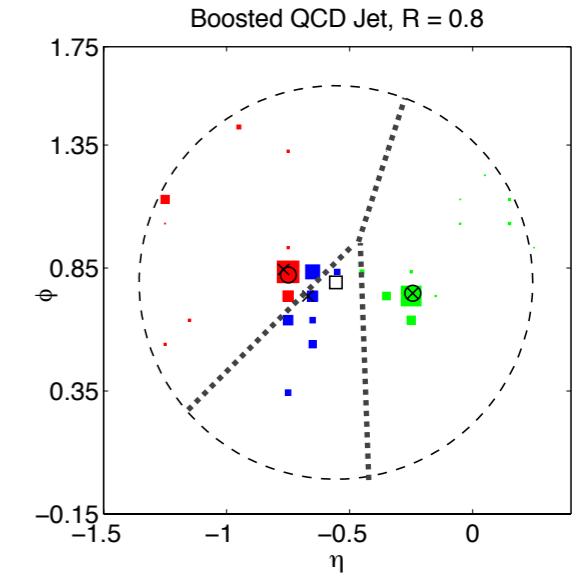
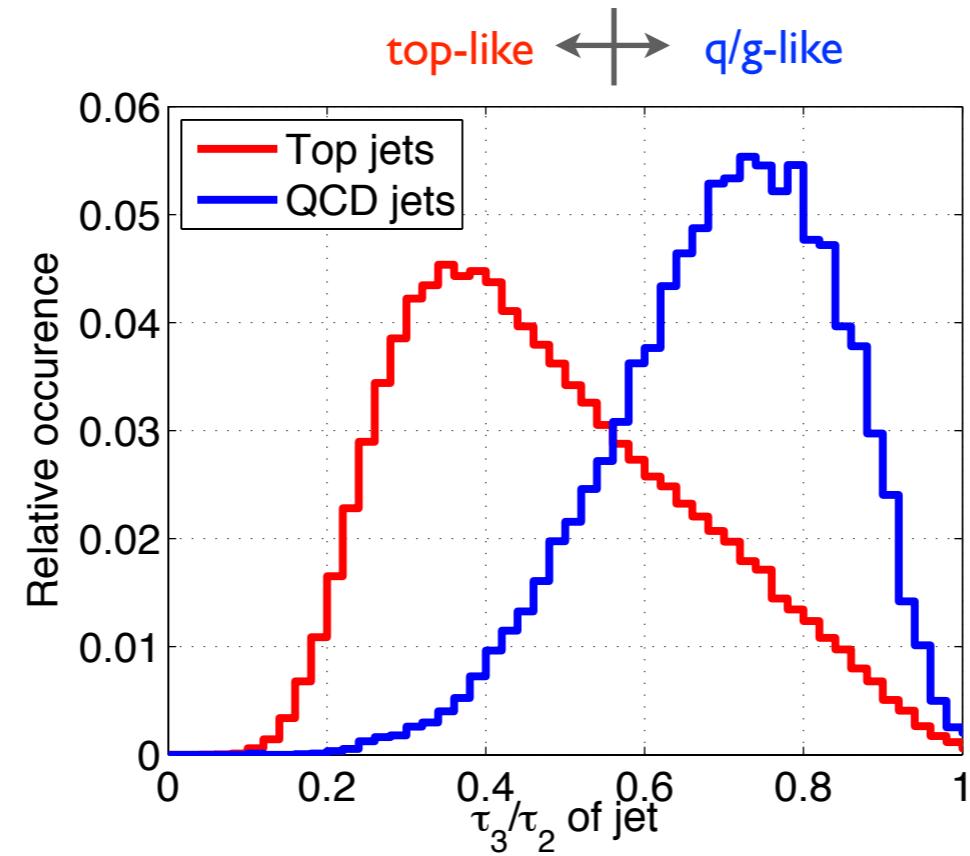
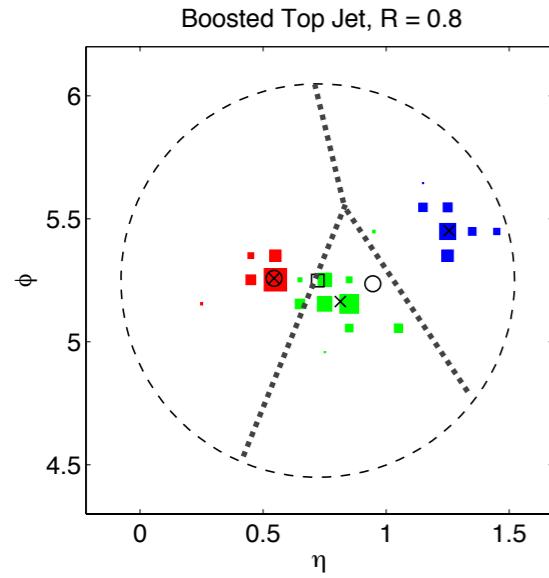
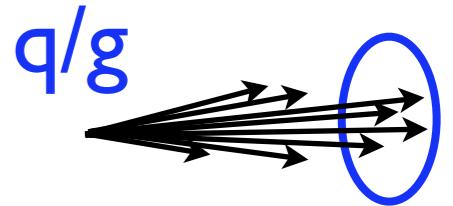
# 3-Prong vs. 1-Prong



*If your eyes can do it...*



# 3-Prong vs. 1-Prong

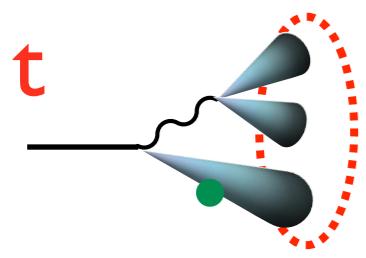


## *N*-subjettiness

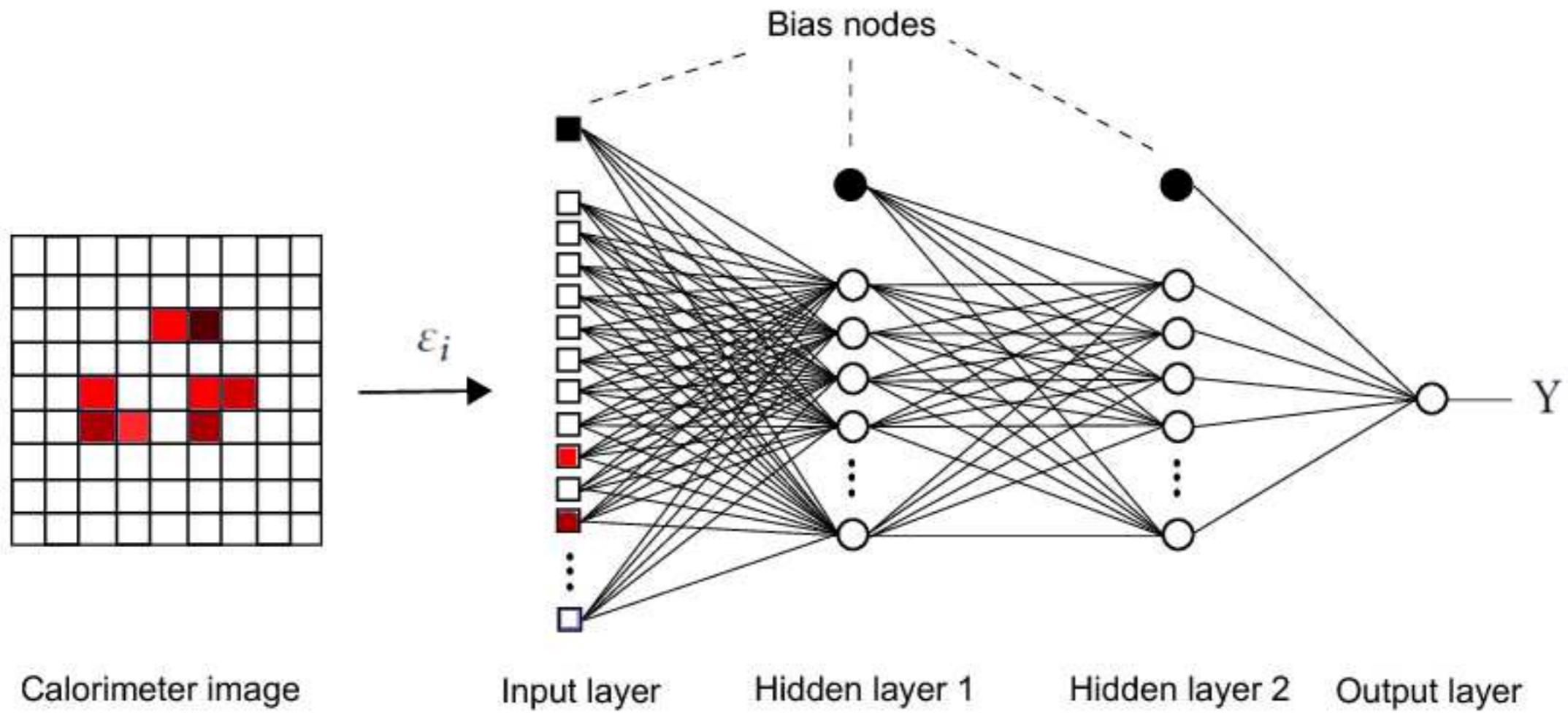
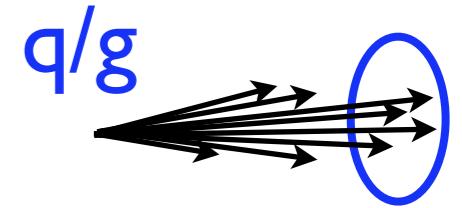
“Deep Thinking”:

$$\tau_N = \sum_k p_{T,k} \min \{ \Delta R_{k,1}, \Delta R_{k,2}, \dots, \Delta R_{k,N} \}$$

[e.g. JDT, Van Tilburg, 1011.2268, 1108.2701]



## 3-Prong vs. 1-Prong



“Deep Learning”: BDTs, FLDs, DNNs, CNNs, RNNs, ...

[e.g. Almeida, Backović, Cliche, Lee, Perelstein, 1501.05968]

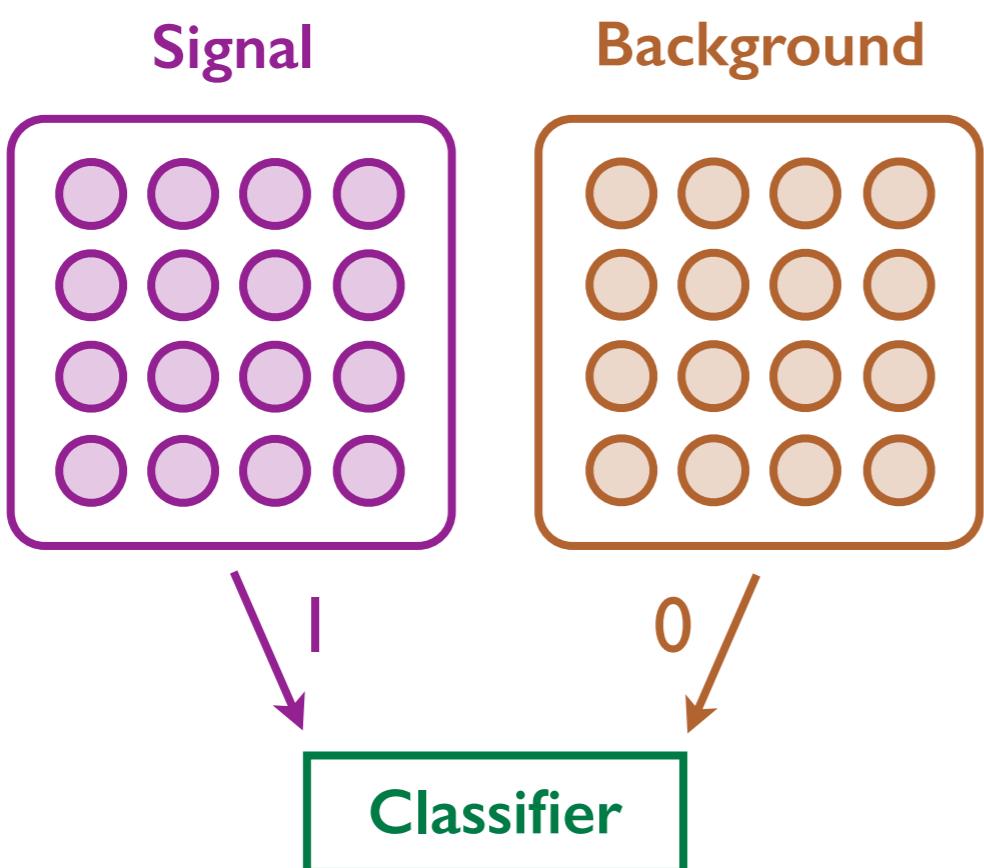
# A Cartoon of Machine Learning

For fully-supervised jet classification

(see backup for regression, generation, modeling)

$$\ell_{\text{MSE}} = \left\langle (\textcolor{violet}{h}(\vec{x}) - 1)^2 \right\rangle_{\text{signal}} + \left\langle (\textcolor{violet}{h}(\vec{x}) - 0)^2 \right\rangle_{\text{background}}$$

Classifier    Inputs



## Minimize Loss Function

(assuming infinite training sets,  
and flexible enough functional form)

$$h(\vec{x}) = \frac{p_{\text{sig}}(\vec{x})}{p_{\text{sig}}(\vec{x}) + p_{\text{bkgd}}(\vec{x})}$$

Optimal Classifier (Neyman–Pearson)

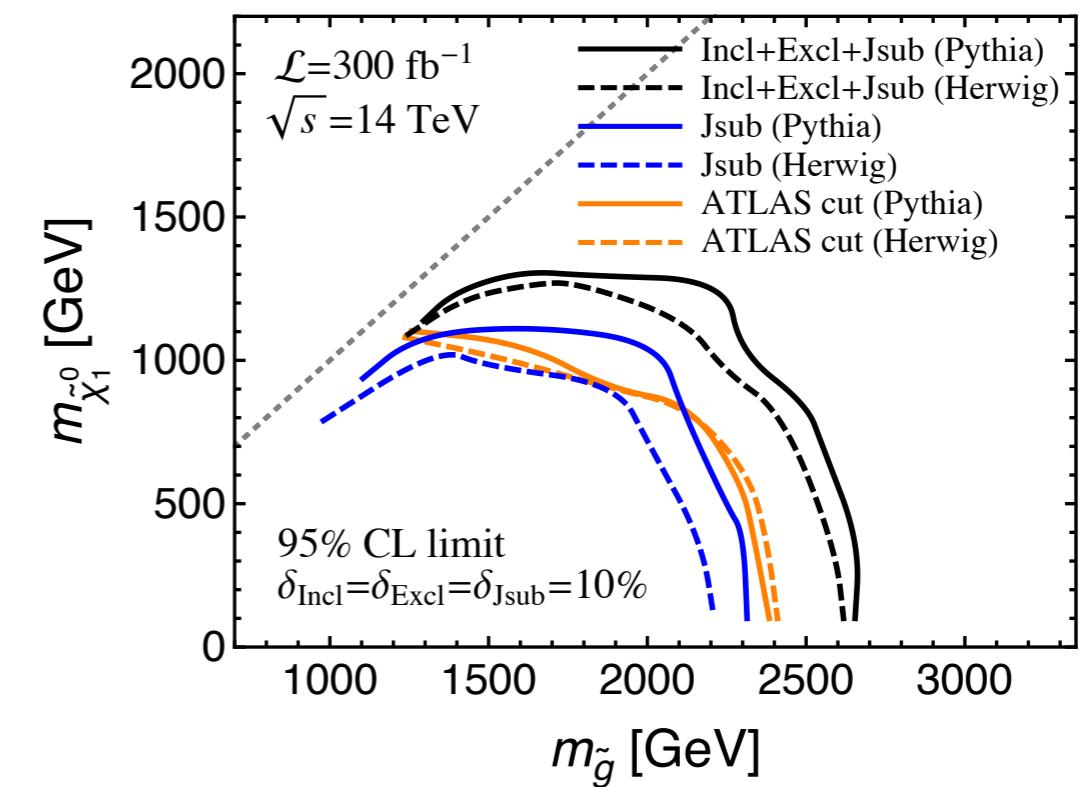
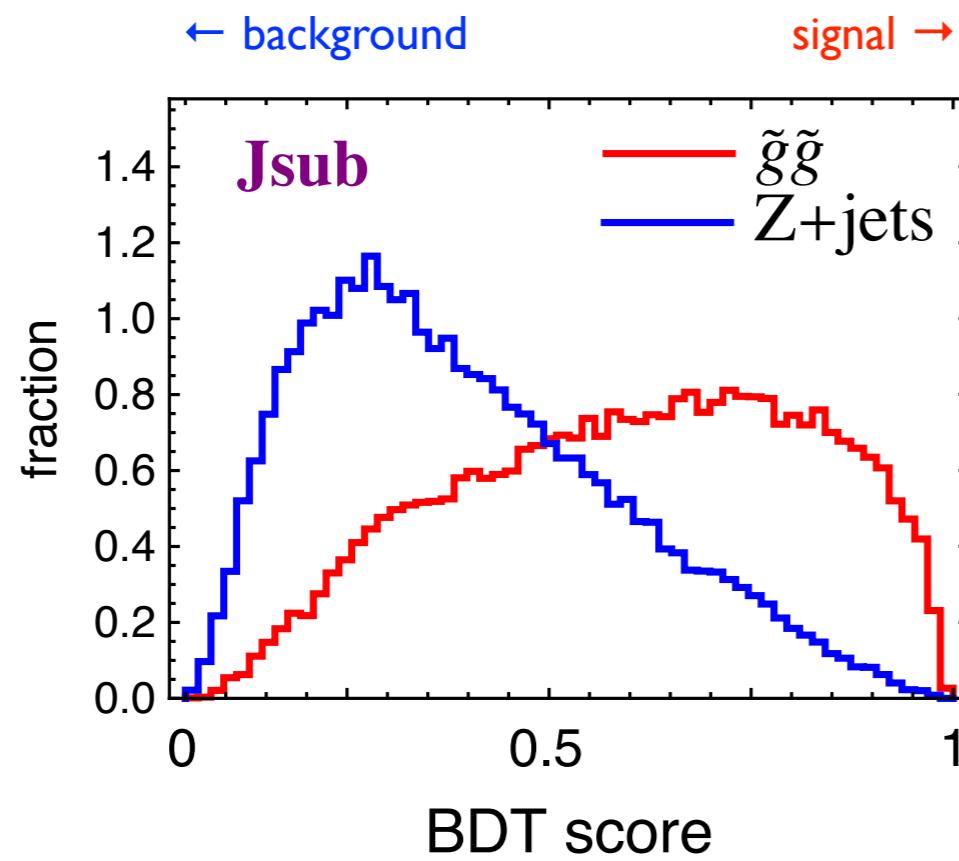
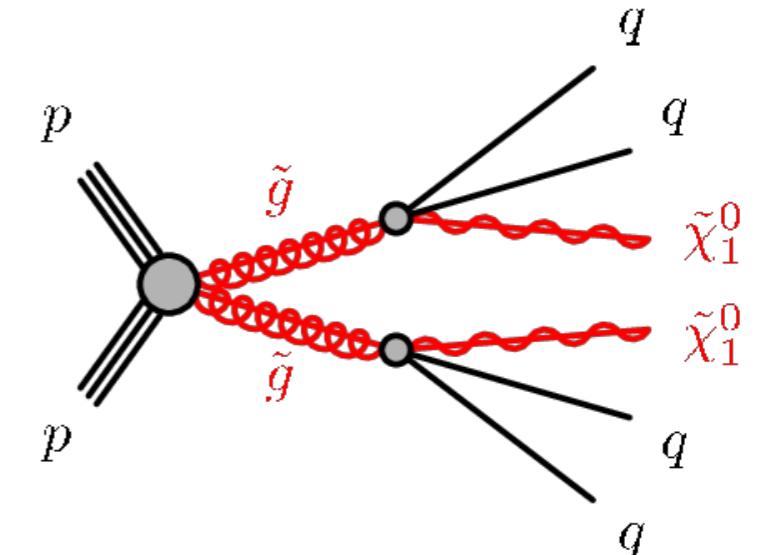
# E.g. SUSY Search for Gluino Pairs

Classifier: Boosted decision tree (for each of 4 jets)

Inputs: Jet mass, width, track multiplicity

Signal: Quark enriched ( $C_F = 4/3$ )

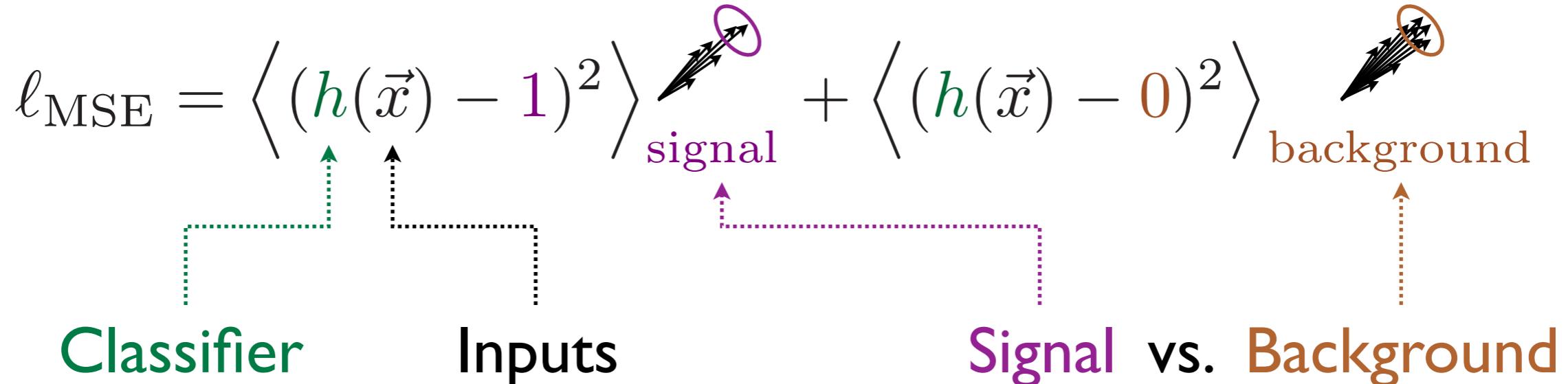
Background: Gluon enriched ( $C_A = 3$ )



[Bhattacherjee, Mukhopadhyay, Nojiri, Sakakie, Webber, 1609.08781]

# Jet Classification Studies

*Mix and match*



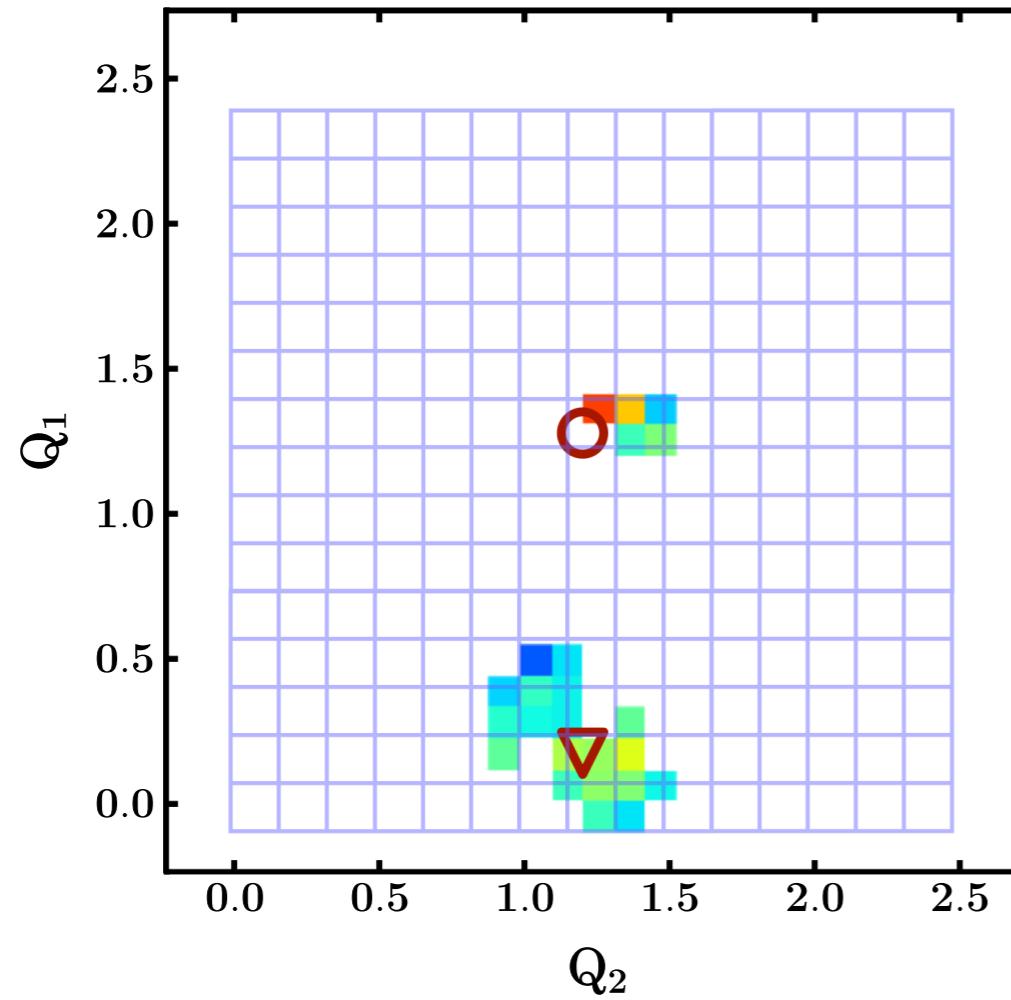
[Lönnblad, Peterson, Rögnvaldsson, 1990, ..., Cogan, Kagan, Strauss, Schwartzman, 1407.5675; Almeida, Backović, Cliche, Lee, Perelstein, 1501.05968; de Oliveira, Kagan, Mackey, Nachman, Schwartzman, 1511.05190; Baldi, Bauer, Eng, Sadowski, Whiteson, 1603.09349; Conway, Bhaskar, Erbacher, Pilot, 1606.06859; Guest, Collado, Baldi, Hsu, Urban, Whiteson, 1607.08633; Barnard, Dawe, Dolan, Rajcic, 1609.00607; Komiske, Metodiev, Schwartz, 1612.01551; Kasieczka, Plehn, Russell, Schell, 1701.08784; Loupe, Cho, Becot, Cranmer, 1702.00748; Pearkes, Fedorko, Lister, Gay, 1704.02124; Datta, Larkoski, 1704.08249, 1710.01305; Butter, Kasieczka, Plehn, Russell, 1707.08966; Fernández Madrazo, Heredia Cacha, Lloret Iglesias, Marco de Lucas, 1708.07034; Aguilar Saavedra, Collin, Mishra, 1709.01087; Cheng, 1711.02633; Luo, Luo, Wang, Xu, Zhu, 1712.03634; Komiske, Metodiev, JDT, 1712.07124; Macaluso, Shih, 1803.00107; Fraser, Schwartz, 1803.08066; Choi, Lee, Perelstein, 1806.01263; Lim, Nojiri, 1807.03312; Dreyer, Salam, Soyez, 1807.04758; Moore, Nordström, Varma, Fairbairn, 1807.04769; plus my friends who will scold me for forgetting their paper; plus many ATLAS/CMS performance studies]

# Jet Classification Studies

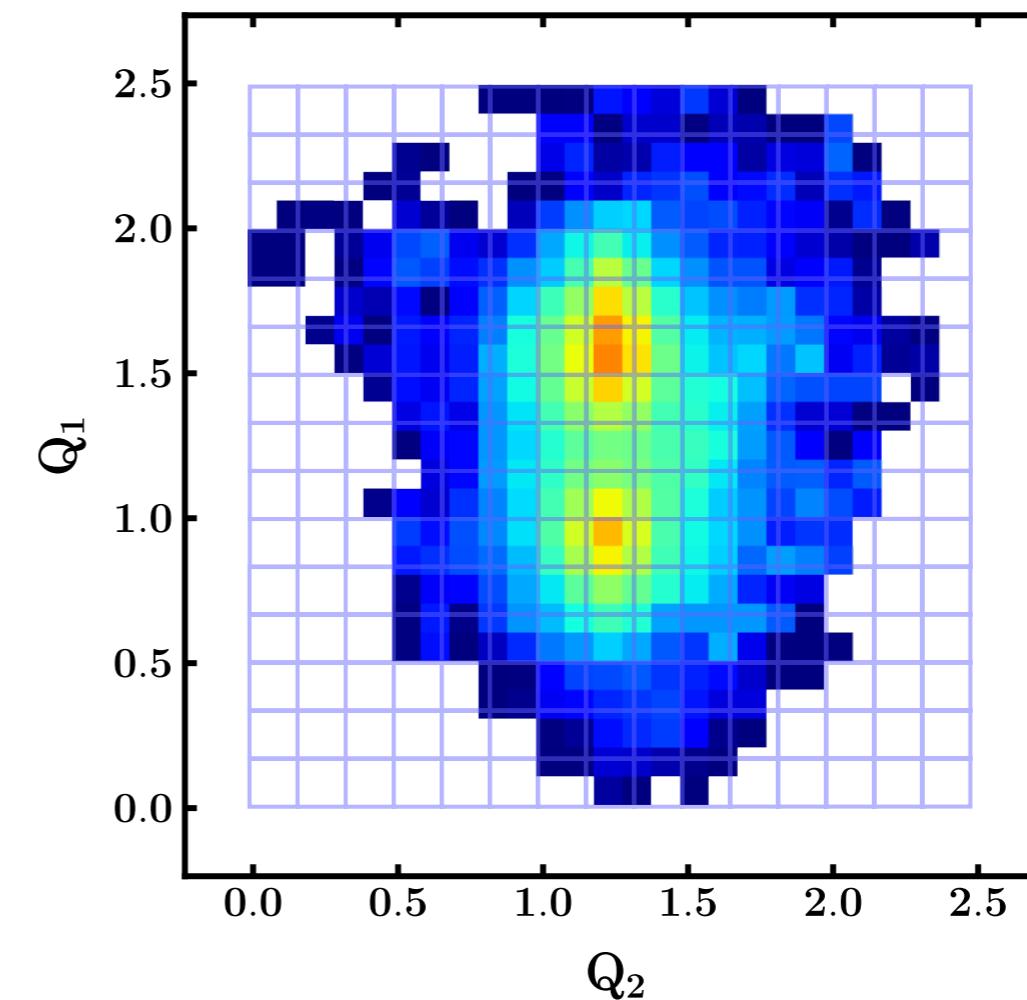
*Mix and match*

Standard CNN input: Jet images

Individual W jet



Ensemble average

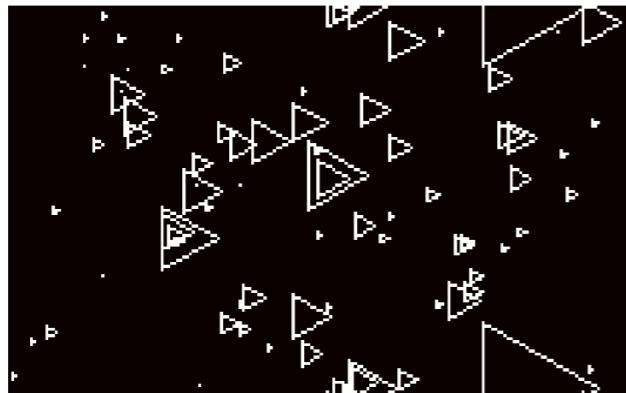


[Cogan, Kagan, Strauss, Schwartzman, 1407.5675]

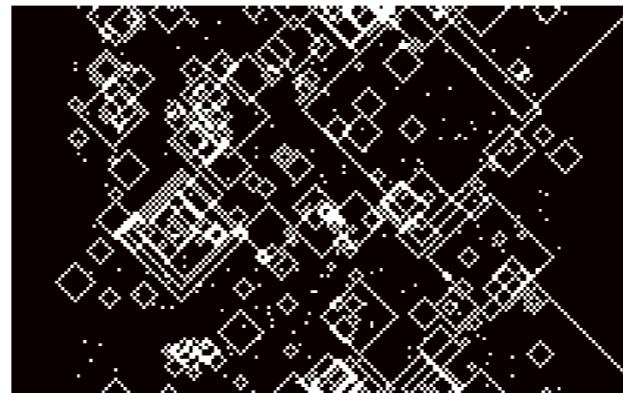
# Jet Classification Studies

*Mix and match*

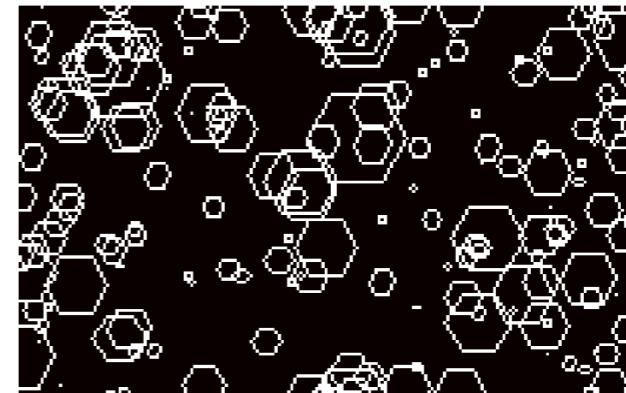
## Novel CNN input: Abstract event images



(a) Photons



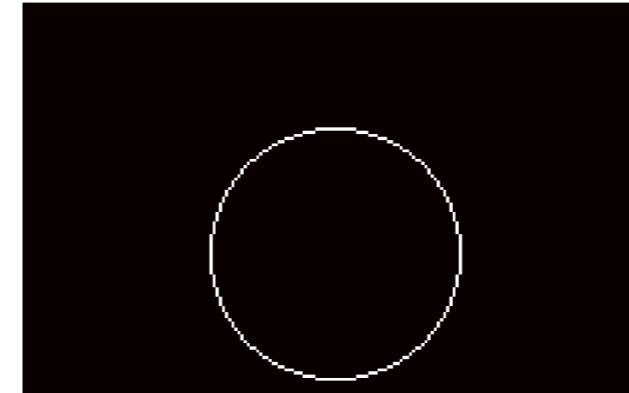
(b) Charged Particles



(c) Neutral Hadrons



(d) Lepton

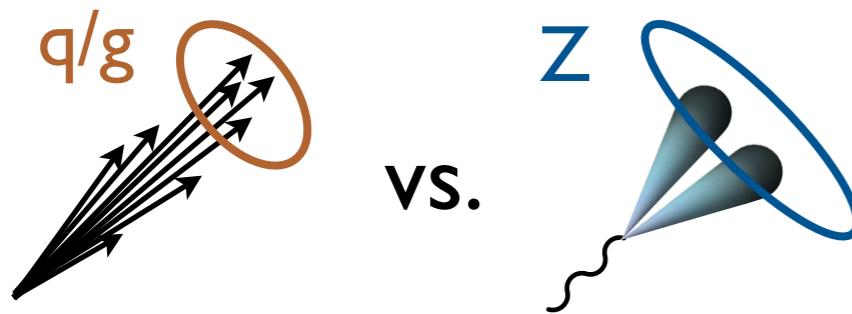


(e)  $E_T^{\text{miss}}$

Addresses sparsity problem of standard energy-to-intensity mapping

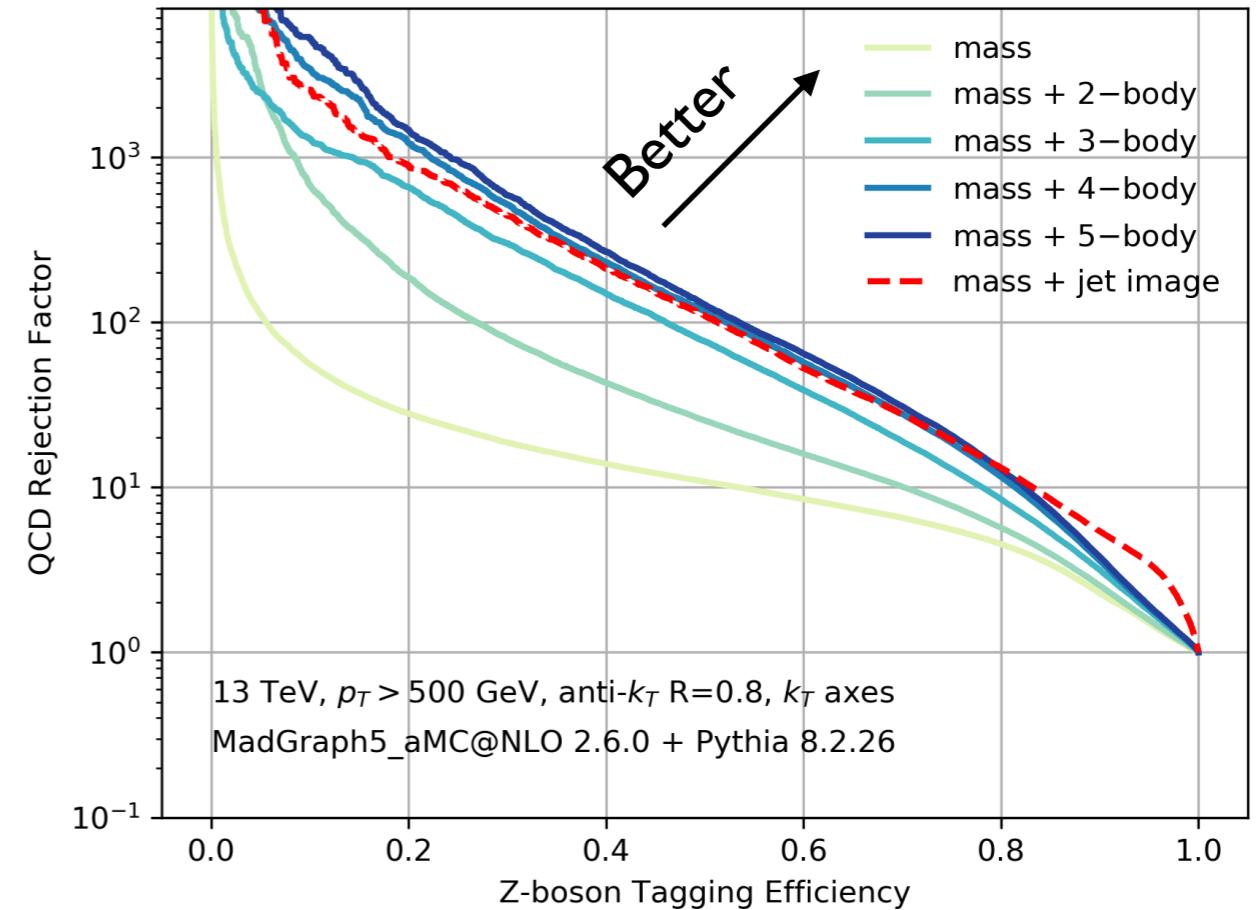
[Nguyen, Weitekamp, Anderson, Castello, Cerri, Pierini, Spiropulu, Vlimant, 1807.00083;  
using Fernández Madrazo, Heredia Cacha, Lloret Iglesias, Marco de Lucas, 1708.07034]

# Evidence for Performance Saturation



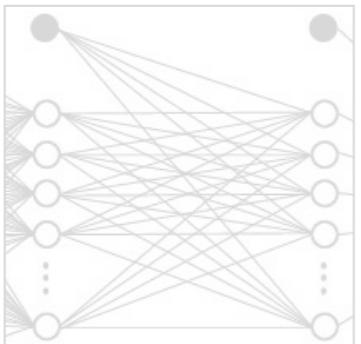
*“Any sufficiently advanced technology  
is indistinguishable from magic”*

Jet Images CNN  $\approx$  “Expert” BDT

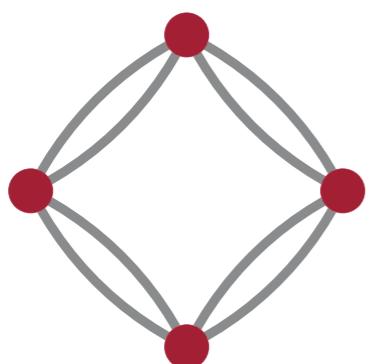


*Next frontier is robustness, versatility & transparency*

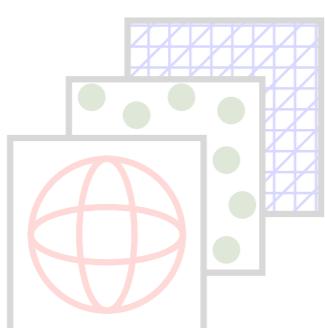
[plot from Moore, Nordström, Varma, Fairbairn, 1807.04769; see also Datta, Larkoski, 1704.08249]



Into the Network



Symmetries & Safety



Deep Sets for Particle Jets

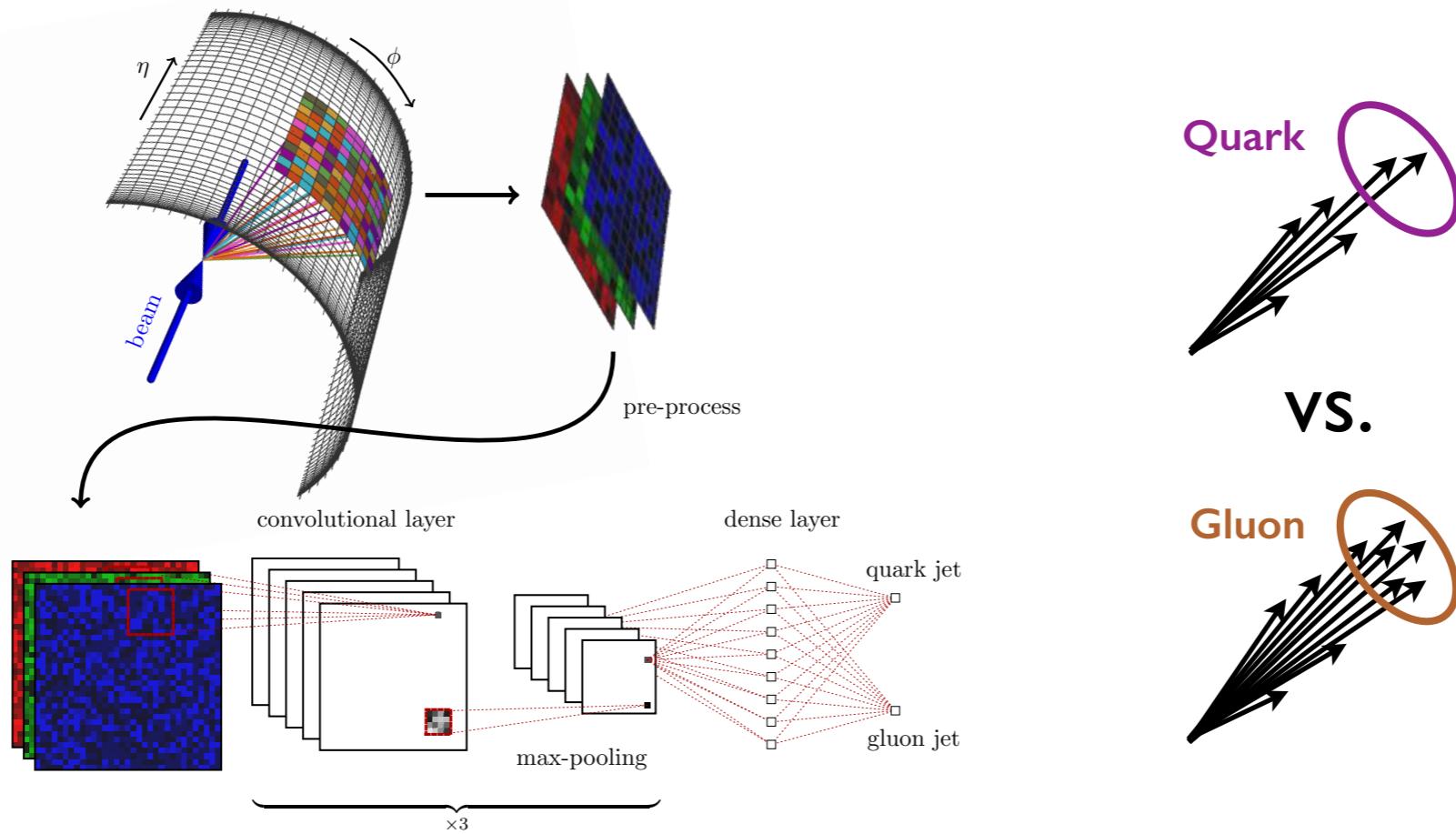


Patrick Komiske

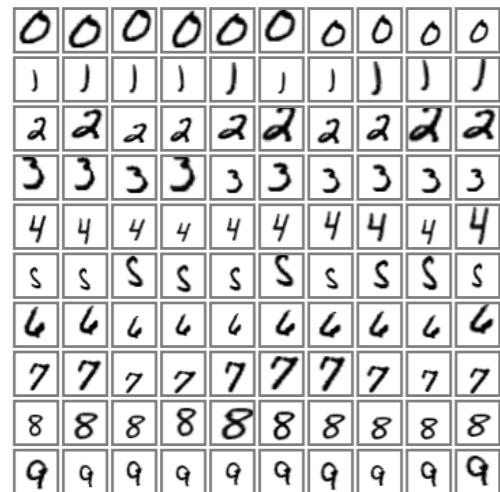
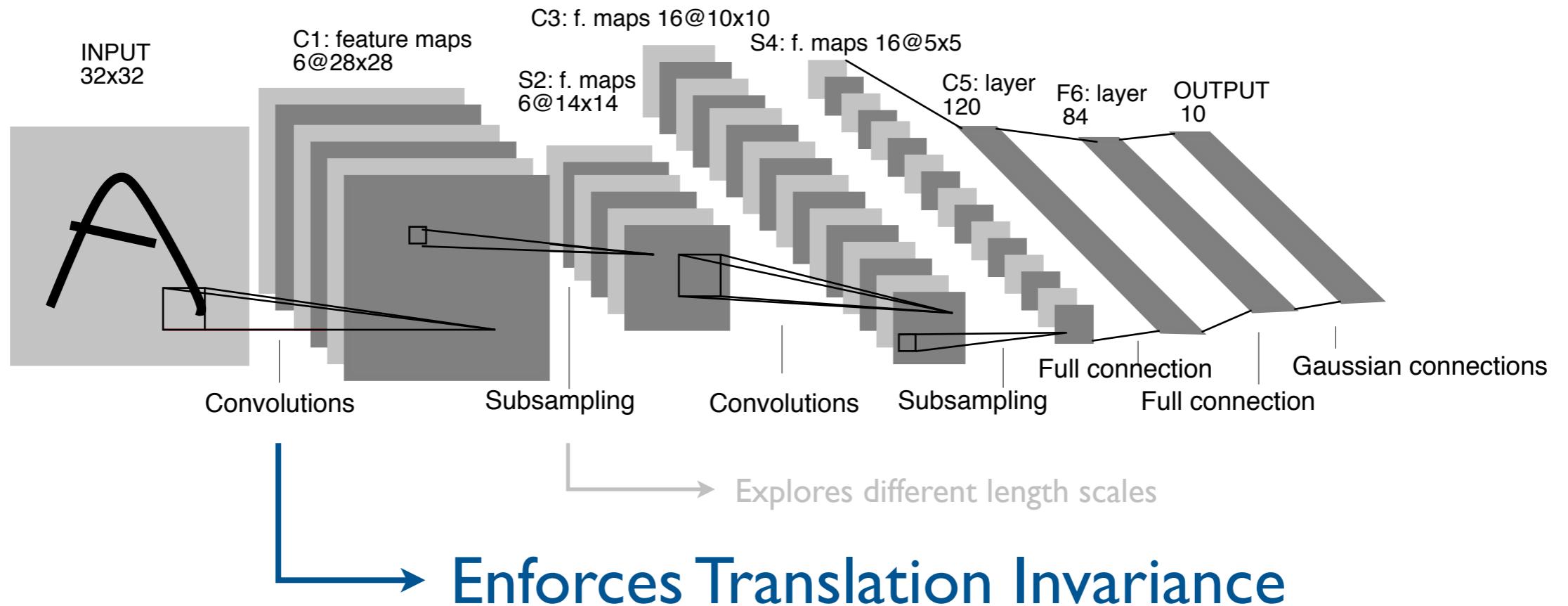


Eric Metodiev

*Two grad students walk into my office with their CNN...*



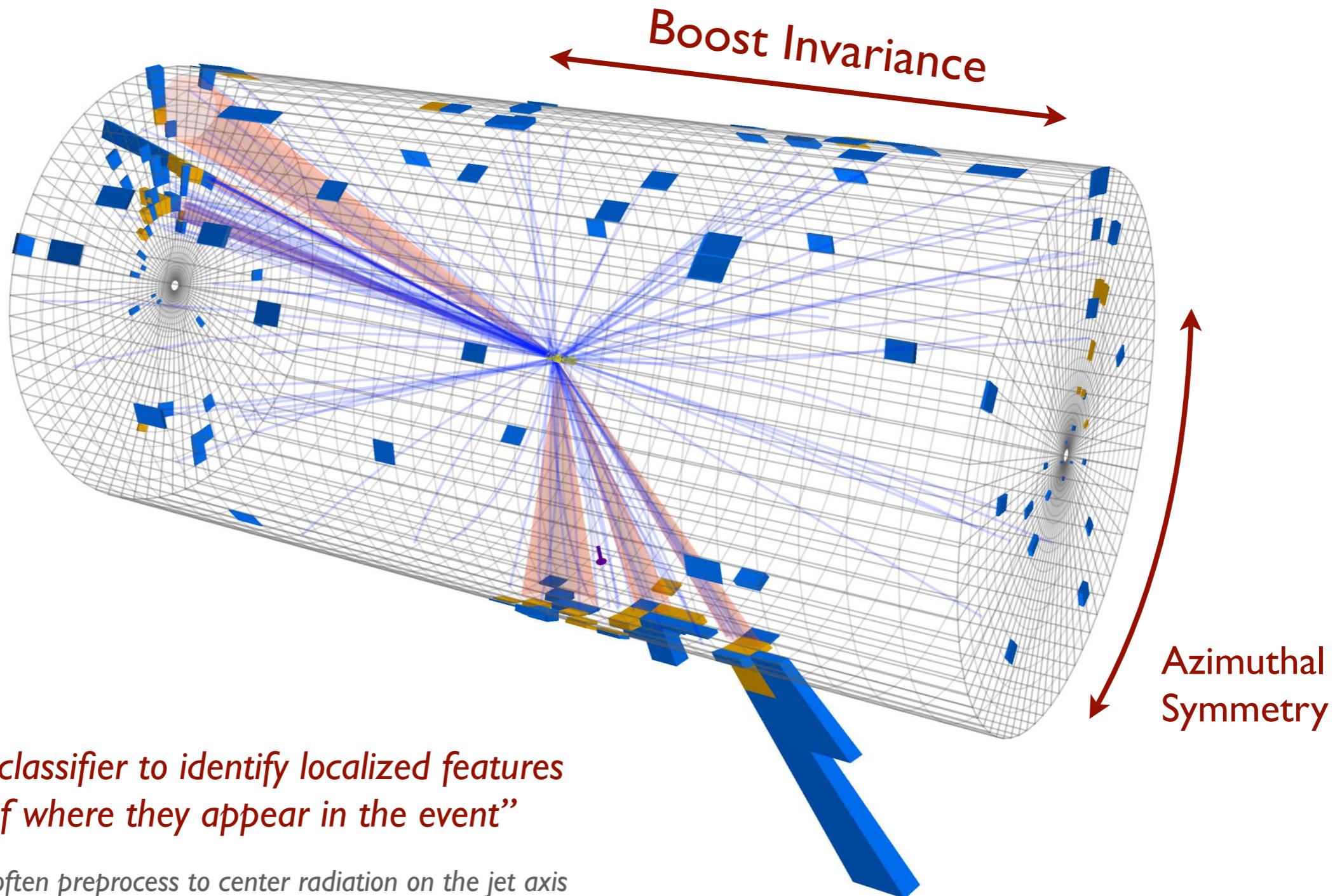
# Symmetries of a CNN



*"I want my classifier to identify localized features regardless of where they appear in the image"*

[image from LeCun, Bottou, Bengio, Haffner, 1998]

# Symmetries of Collision Events



[image from CMS, 2015]

# The Physics-First Approach

Underlying Physics



“Deep Thinking”

Natural Data Representation



“Deep Learning”

Suitable Algorithm

# The Buzzword-First Approach

Questionable Physics



“Wishful Thinking”

Unnatural Data Representation



“Shoehorning”

Cool-Sounding Algorithm

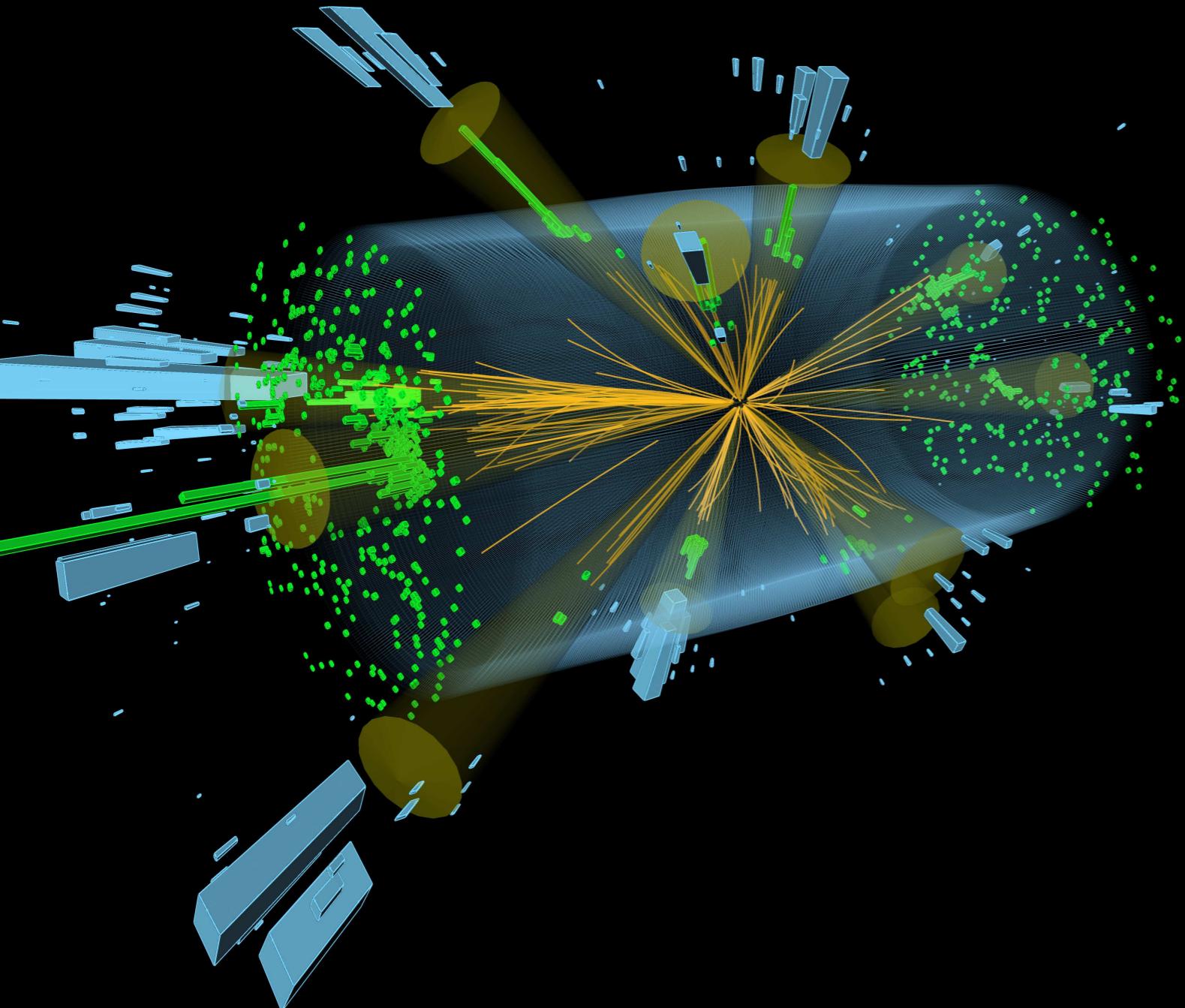
# Why CNNs Aren't Ideal

Underlying Physics  
Natural Data Representation  
Suitable Algorithm

## Translations/Boost Invariance



# Debris Taxonomy



T E H M



$\gamma$

photon



$e^+$

electron



$\mu^+$

muon



$\pi^+$

pion



$K^+$

kaon



$K_L^0$

K-long



$p/\bar{p}$

proton



$n/\bar{n}$

neutron

elementary

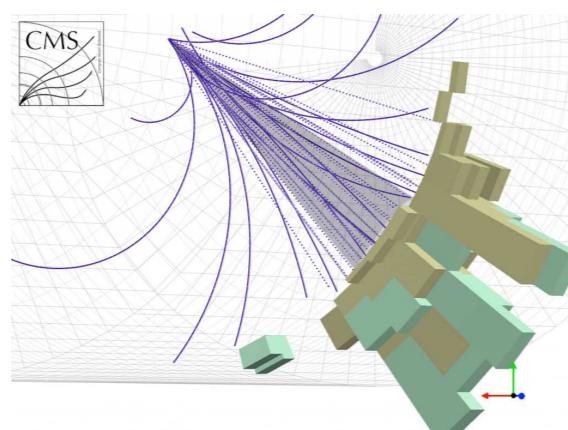
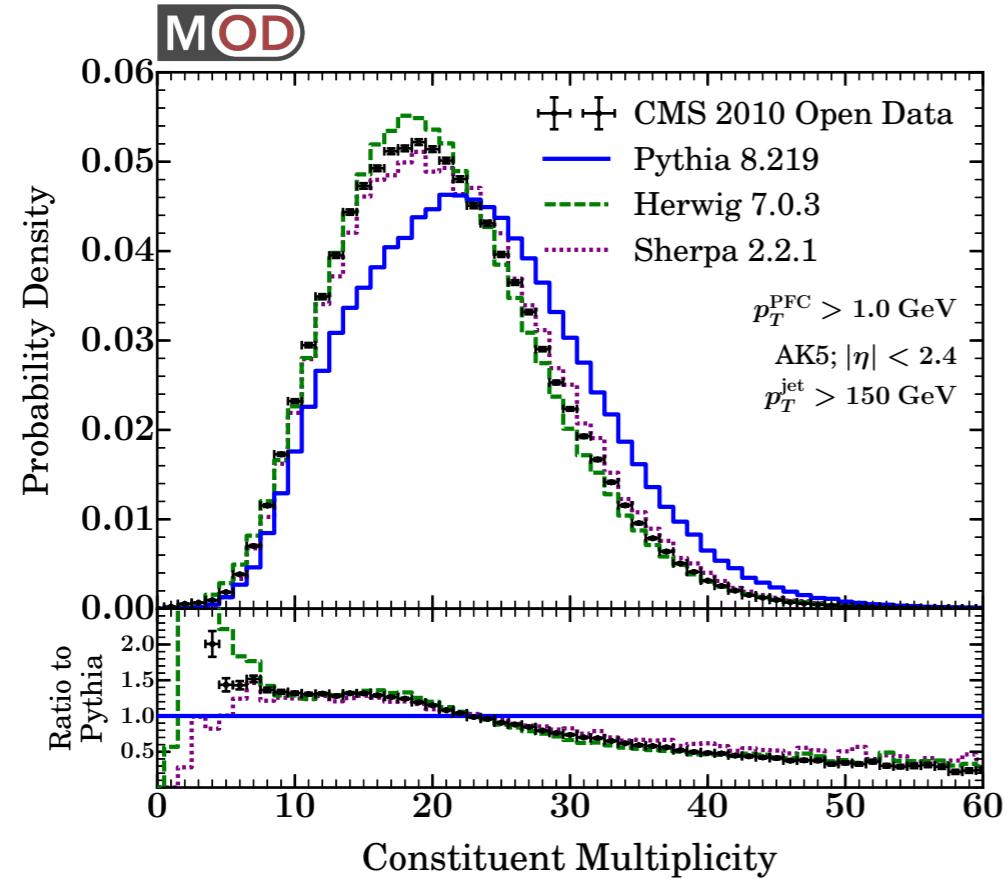
composite

# Point Clouds



[Popular Science, 2013]

# Key Fact #1



Jet constituents:  
*Particle-like objects*  
*Variable-length*  
*Unordered set*

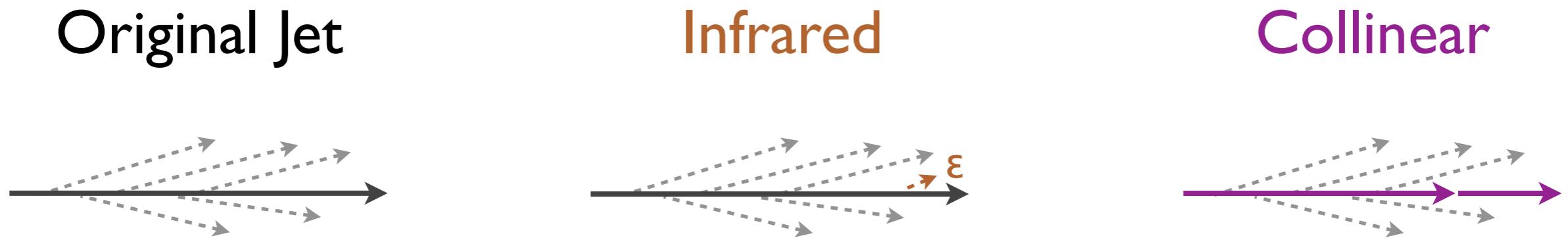
Per particle:  
 $\{E, p_x, p_y, p_z\}$  or  $\{p_T, \gamma, \Phi, m\}$   
*Flavor/charge labels*  
*Vertex information*  
*Quality criteria, etc.*

[plot from Tripathee, Xue, Larkoski, Marzani, JDT, 1704.05842]

# Key Fact #2

Wide range of interesting observables are “safe”

*Interesting  $\approx$  Calculable in fixed-order perturbation theory*



**IRC Safe Observable:** Insensitive to **IR** or **C** emissions

Enforces smooth interpolation between  
variable-length inputs (i.e.  $N \rightarrow N-1$ )

# In the Backup (2017)

Underlying Physics  
Natural Data Representation  
Suitable Algorithm

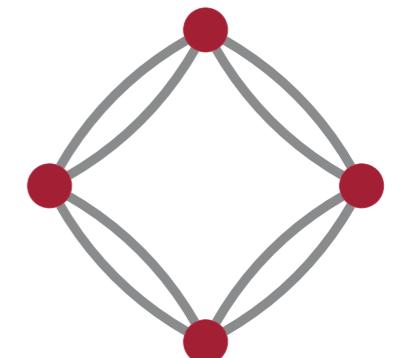
Infrared/Collinear Safety



Energy Flow Polynomials



Linear Regression



[Komiske, Metodiev, JDT, 1712.07124;  
<https://energyflow.network>]

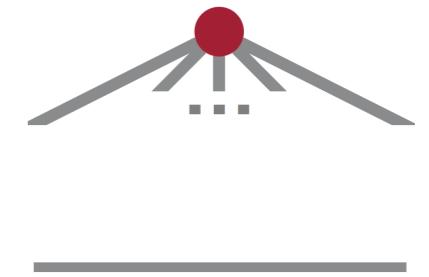
# In the Backup (2019?)

Underlying Physics  
Natural Data Representation  
Suitable Algorithm

## Infrared/Collinear Safety



## Energy Flow Moments



## Linear Regression + Linear Runtime

[Komiske, Metodiev, JDT, we've been promising this paper for 9 months]

# Today's Talk (2018)

Underlying Physics  
Natural Data Representation  
Suitable Algorithm

## Quantum-Mechanical Indistinguishability

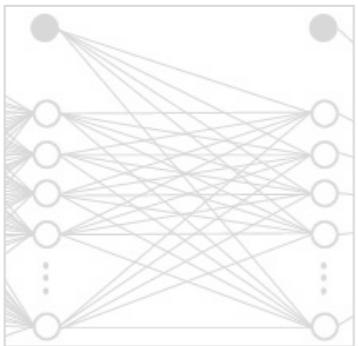


## Variable-Length Unordered Sets of Particles

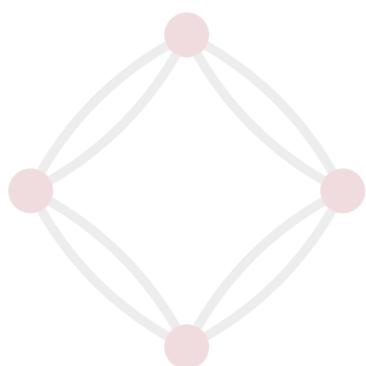


???

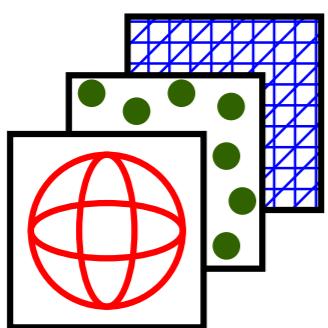
[Komiske, Metodiev, JDT, 1810.05165]



Into the Network



Symmetries & Safety



Deep Sets for Particle Jets

# The Power of Addition

**Additive Observable:**  $\mathcal{O} = \sum_{i \in \text{jet}} \Phi(E_i, \vec{p}_i, \dots)$

(often comes up in the context of  
SCET factorization theorems)

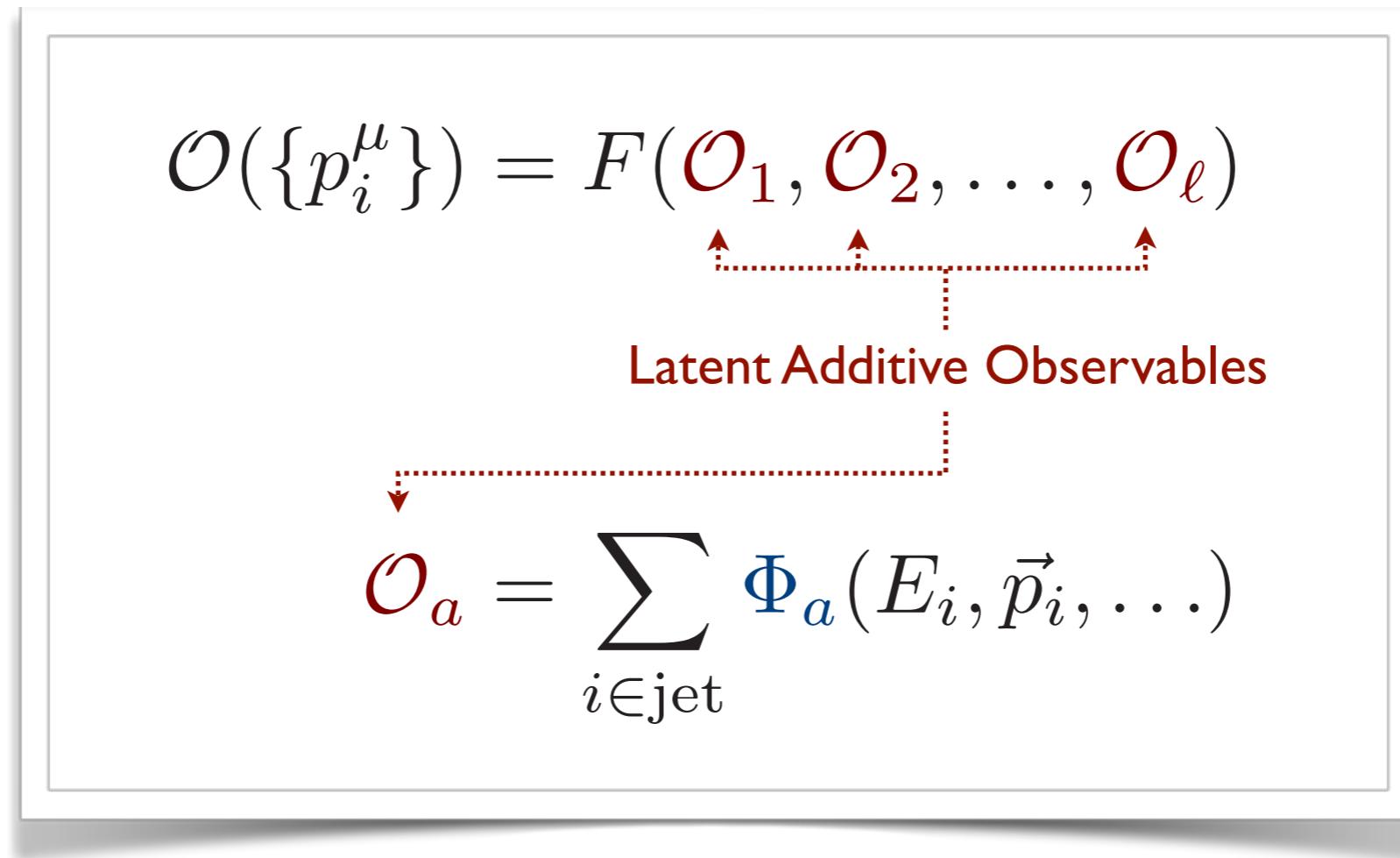
- Permutation invariant by construction
- Easily adapts to variable-length inputs
- Can approximate  $\Phi$  with neural networks
- Can incorporate additional particle properties
- Linear runtime in number of particles

**Additive Safe Observable:**  $\mathcal{O} = \sum_{i \in \text{jet}} E_i \Phi(\hat{p}_i) \quad \hat{p}_i = \frac{\vec{p}_i}{E_i}$

IRC safety guaranteed by energy weighting

# Conjectured Generalization

*Arbitrary permutation-symmetric observable?*



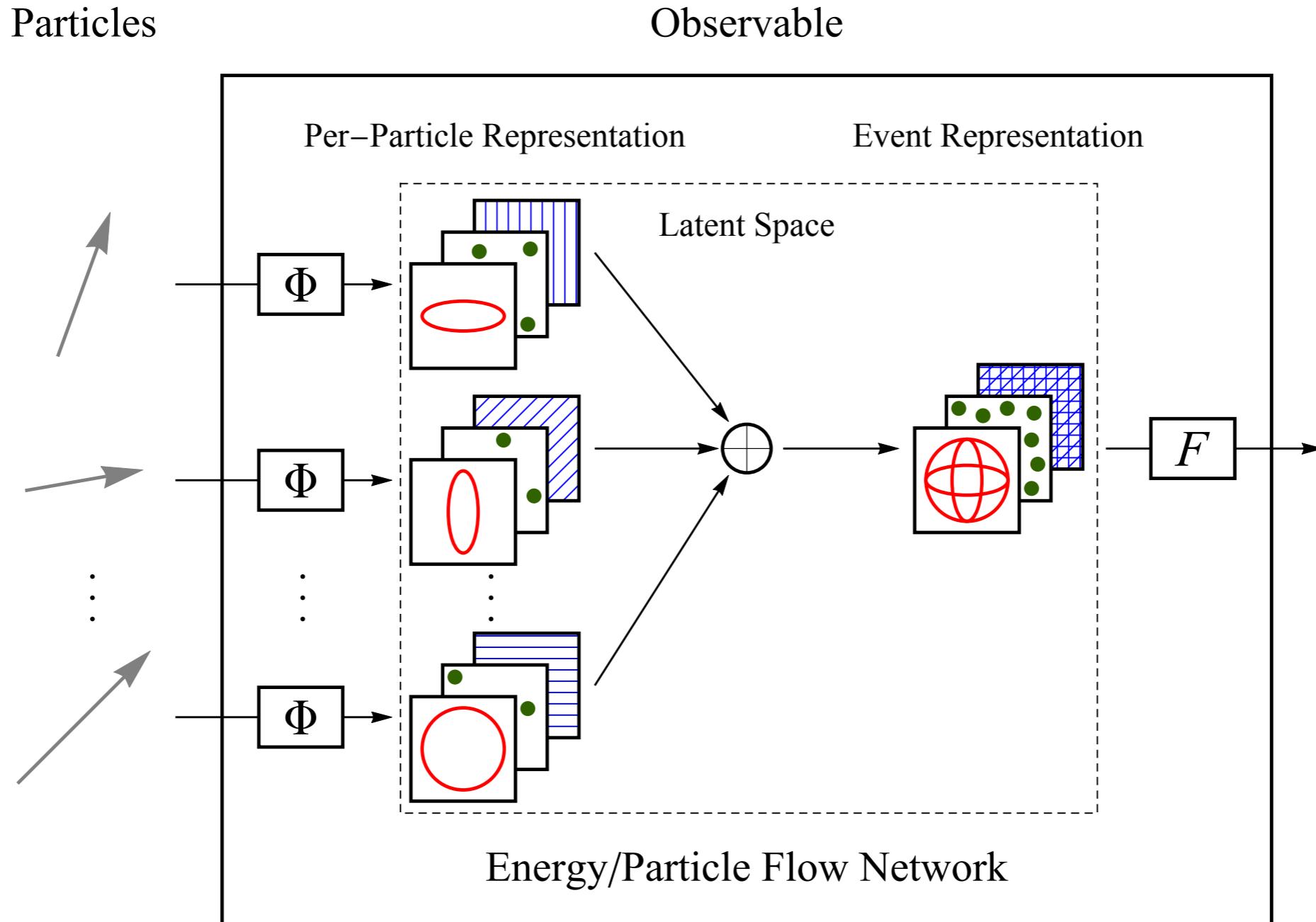
## Energy / Particle Flow Networks

IRC-safe  $\Phi$

General  $\Phi$

[Komiske, Metodiev, JDT, 1810.05165]

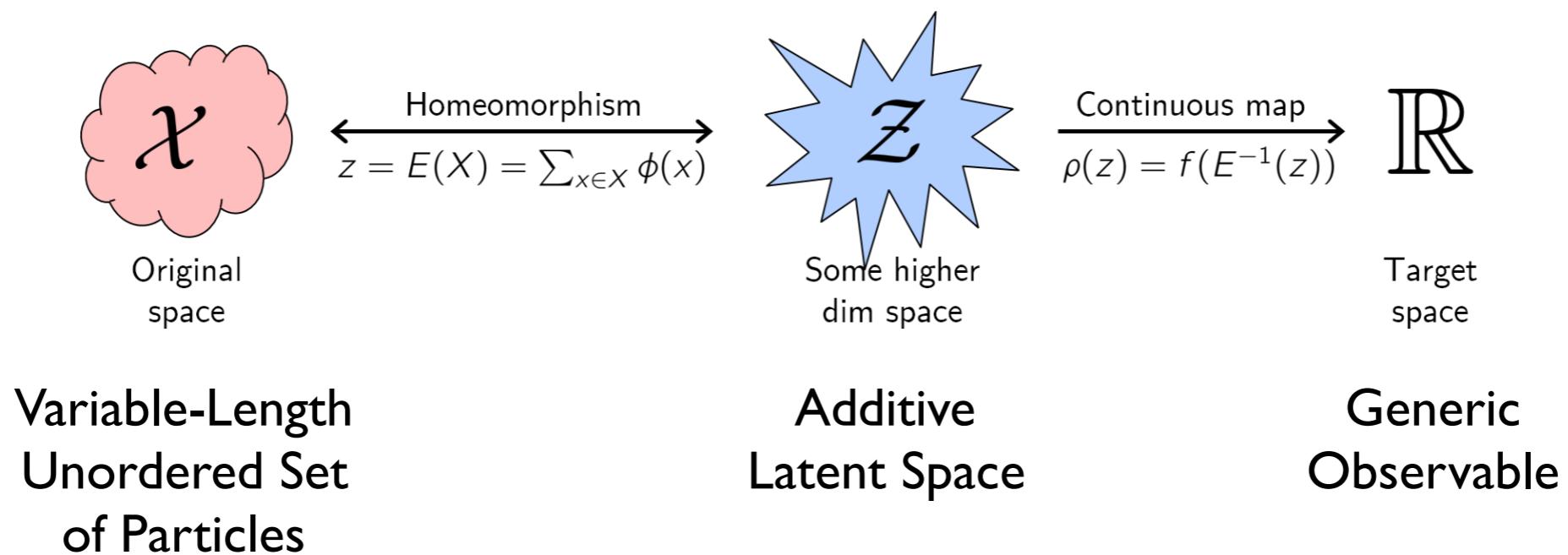
# Conjectured Generalization



# Meanwhile in ML-Land: Deep Sets

**Theorem 2** A function  $f(X)$  operating on a set  $X$  having elements from a countable universe, is a valid set function, i.e., **invariant** to the permutation of instances in  $X$ , iff it can be decomposed in the form  $\rho \left( \sum_{x \in X} \phi(x) \right)$ , for suitable transformations  $\phi$  and  $\rho$ .

↑  
(!)



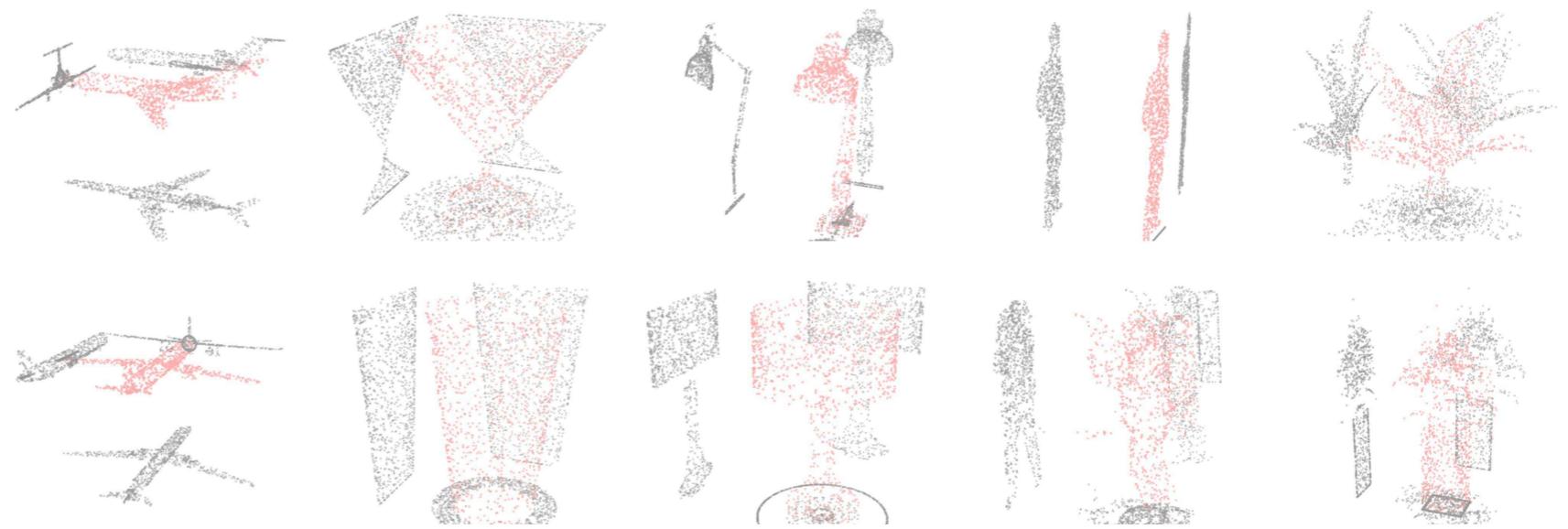
[Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, 1703.06114]

# Deep Sets for...

## Celebrity Face Anomaly Detection



## Point Cloud Classification



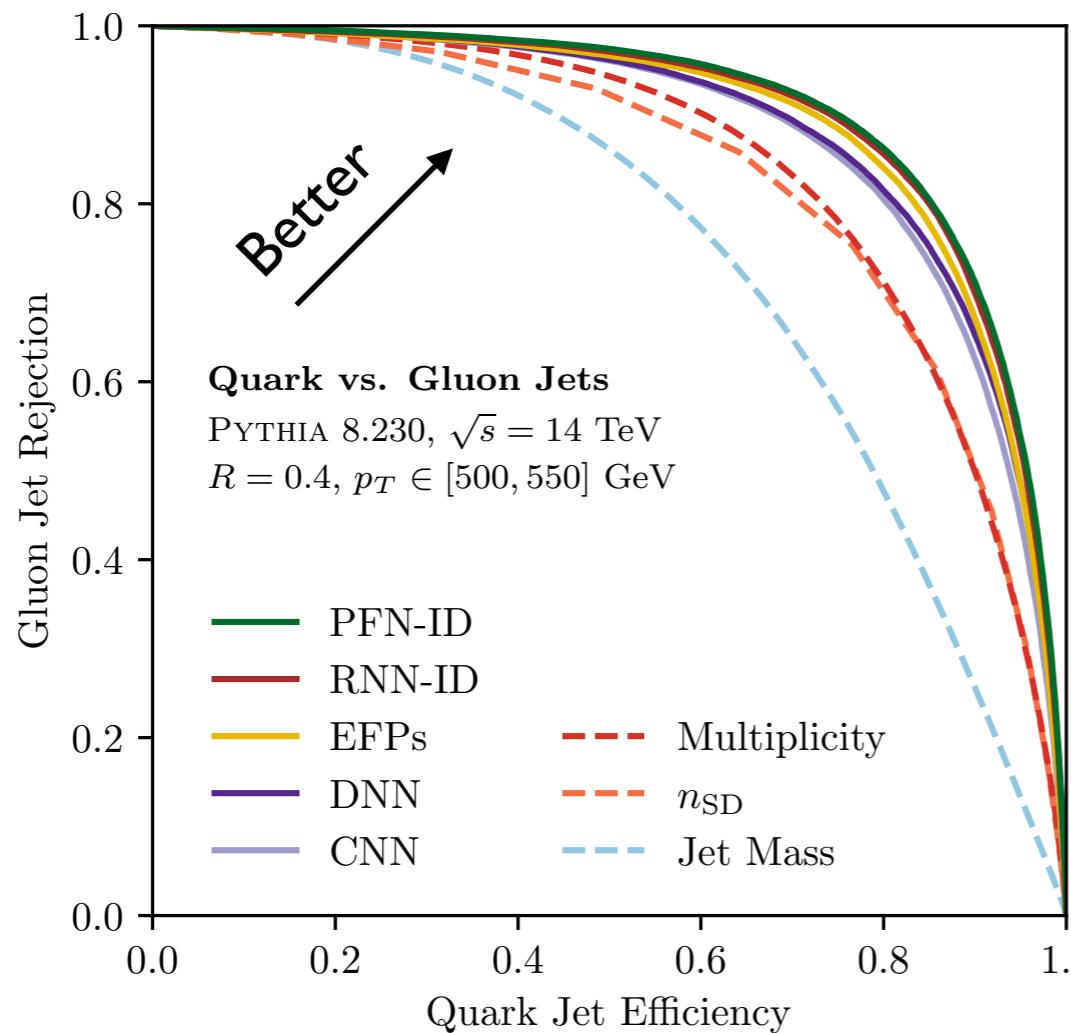
[Zaheer, Kottur, Ravanbakhsh, Poczos, Salakhutdinov, Smola, 1703.06114]

# Deep Sets for Q vs. G

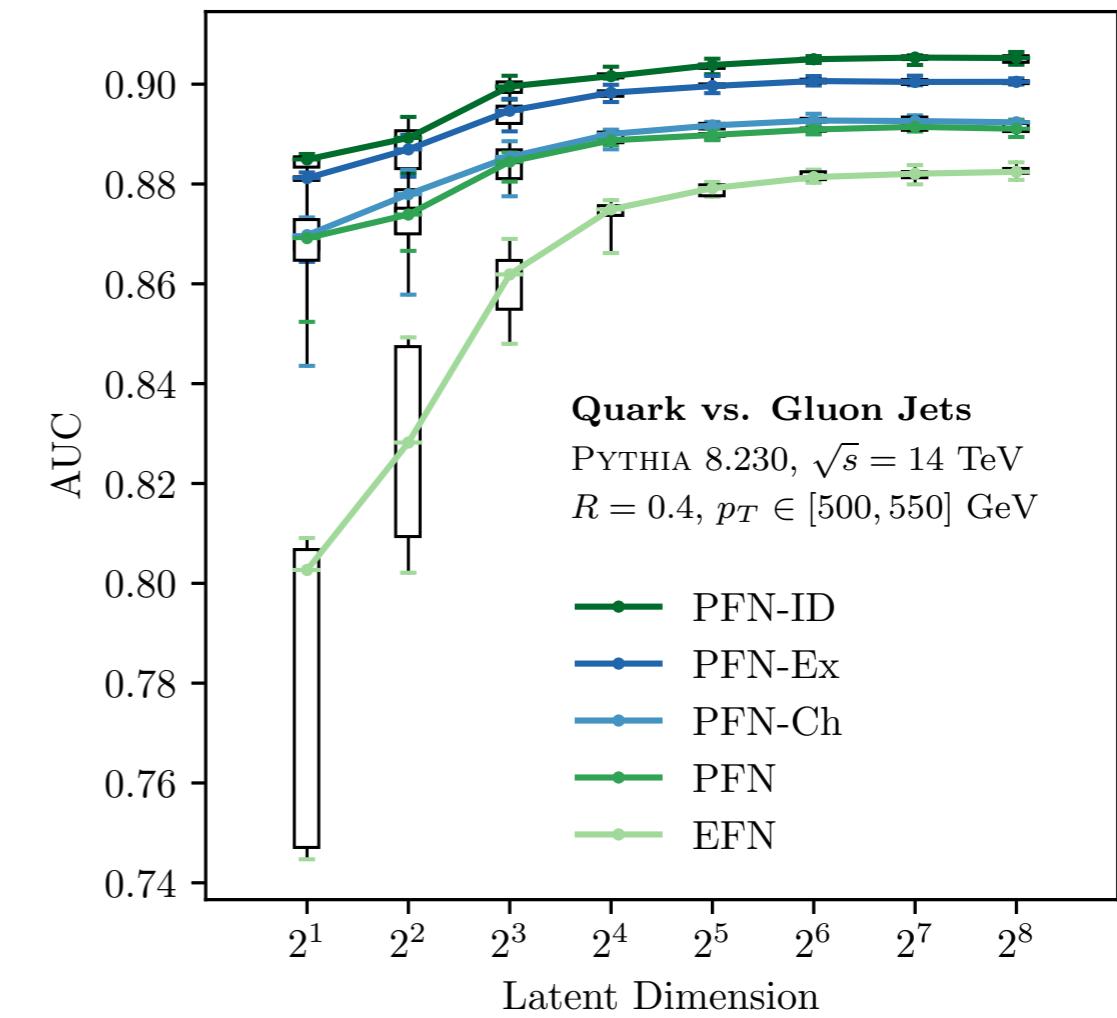
The “Hello, World!” of jet classification



Competitive with  
previous methods



Performance improves  
with more information

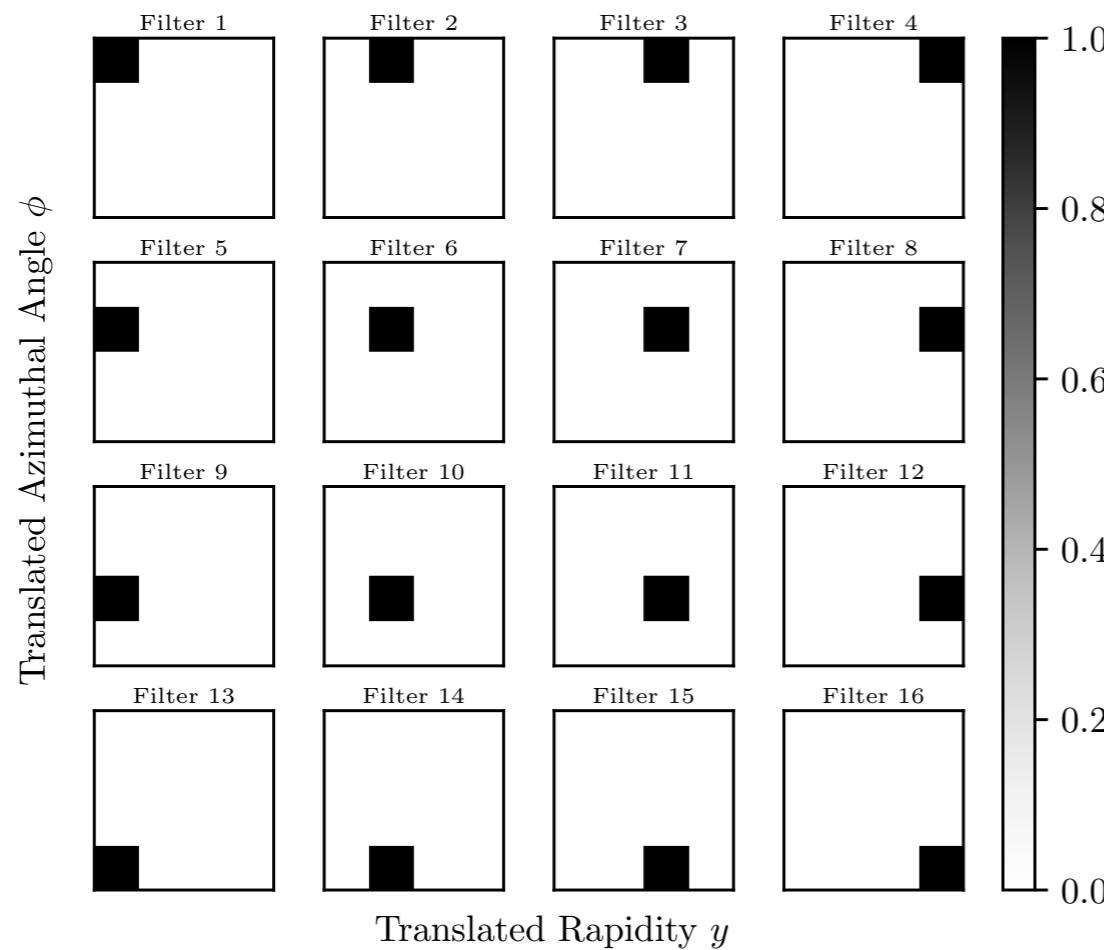


# Latent Space Visualization

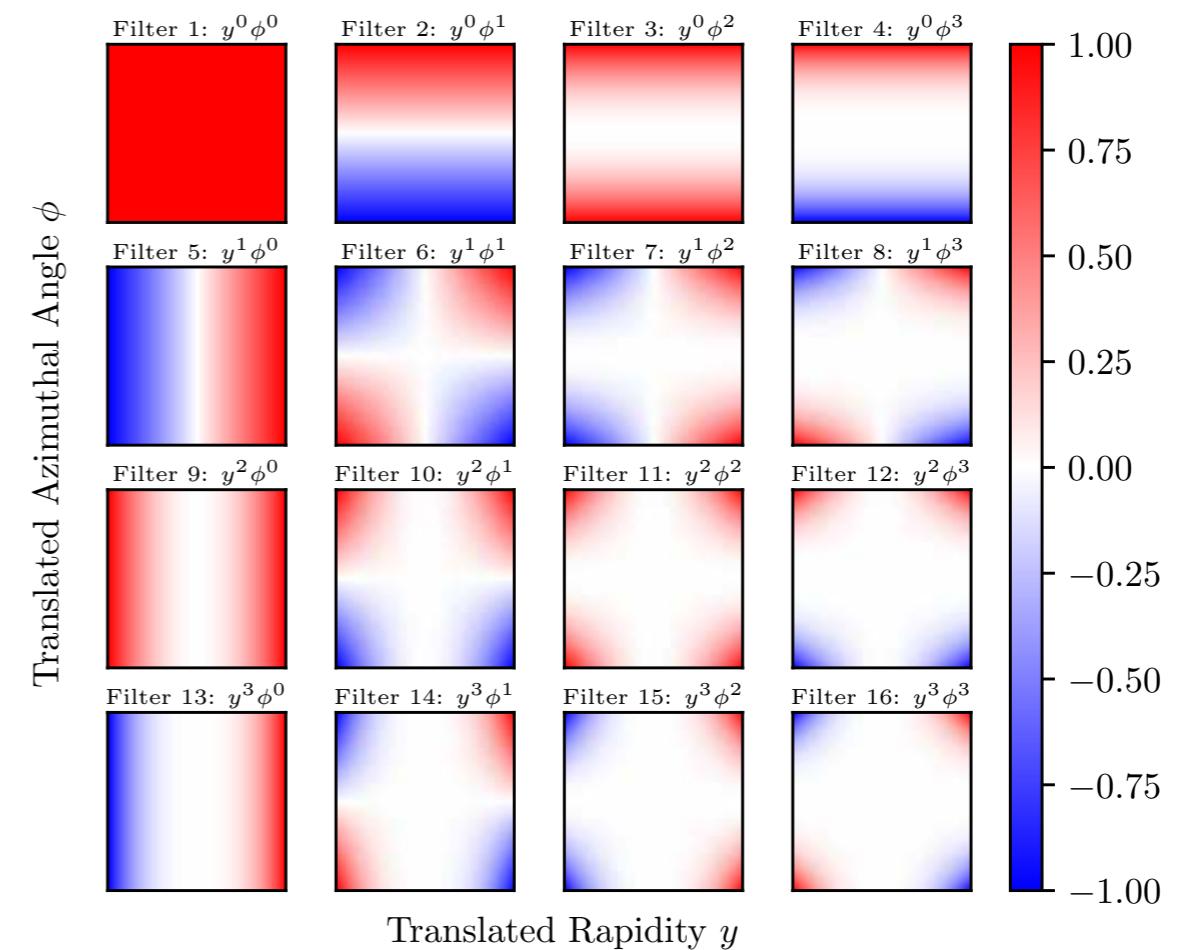
$$\mathcal{O}(\{p_i^\mu\}) = F(\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_\ell)$$

**IRC-safe:**  $\mathcal{O}_a = \sum_{i \in \text{jet}} p_{Ti} \Phi_a(\phi_i, y_i)$

## Calorimeter Pixels



## Radiation Moments

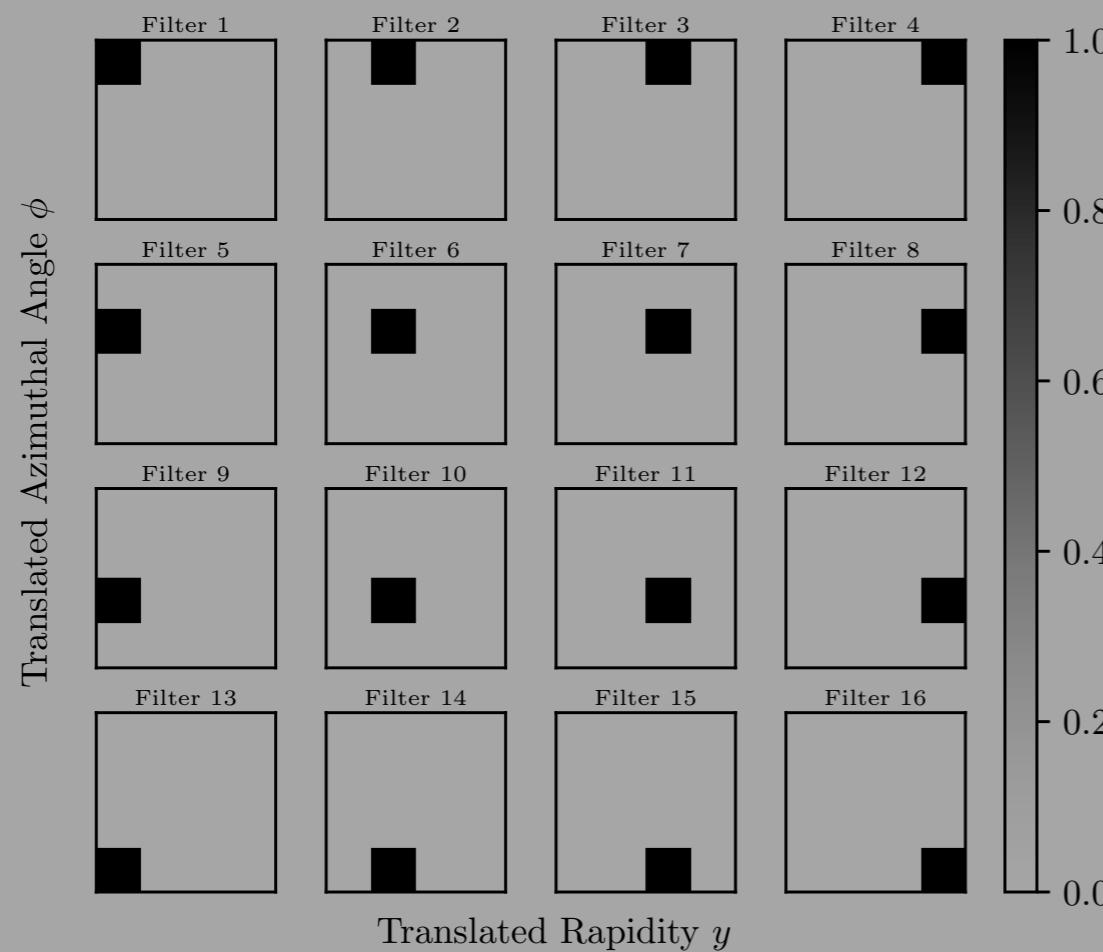


# Latent Space Visualization

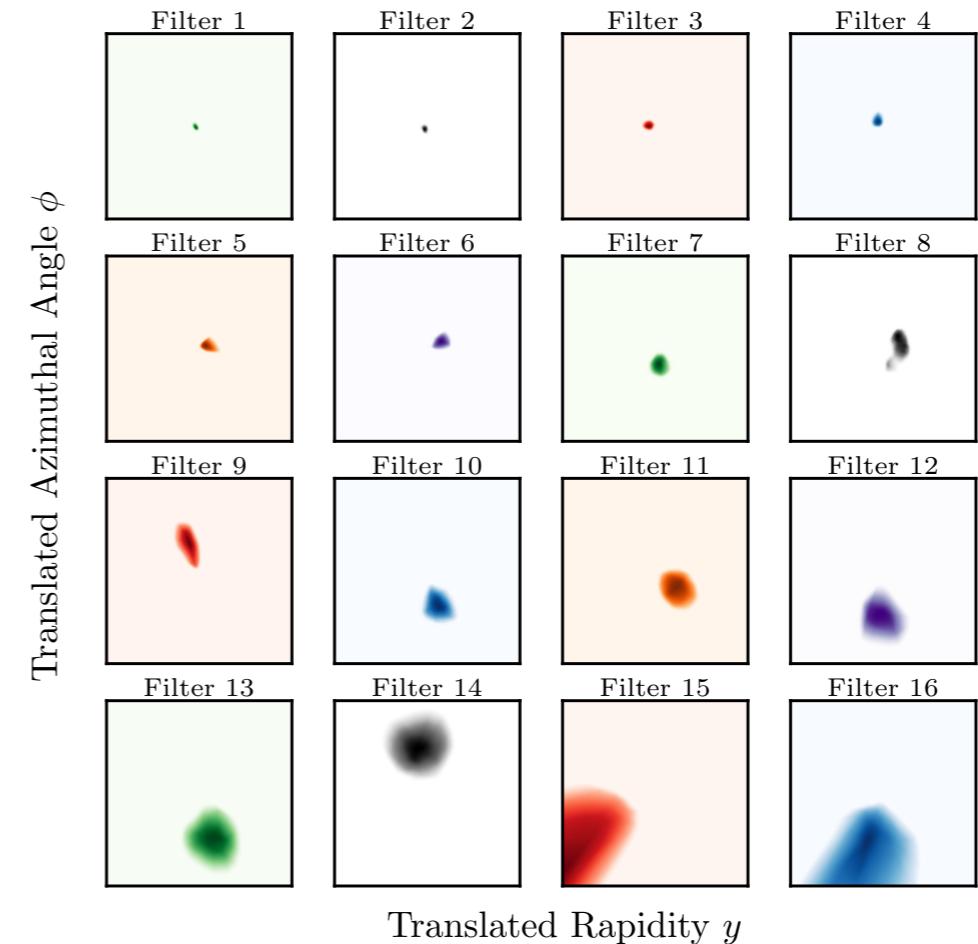
$$\mathcal{O}(\{p_i^\mu\}) = F(\mathcal{O}_1, \mathcal{O}_2, \dots, \mathcal{O}_\ell)$$

**IRC-safe:**  $\mathcal{O}_a = \sum_{i \in \text{jet}} p_{Ti} \Phi_a(\phi_i, y_i)$

## Calorimeter Pixels

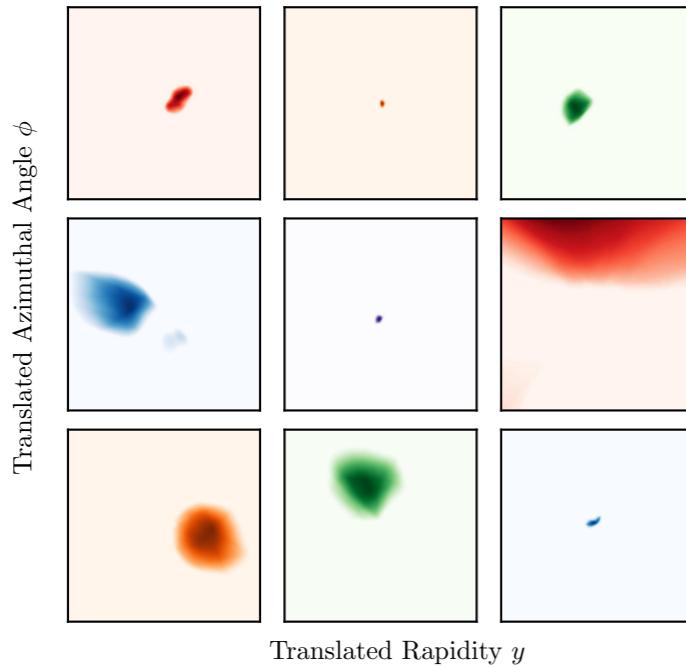


## EFNs: Dynamic Pixilation

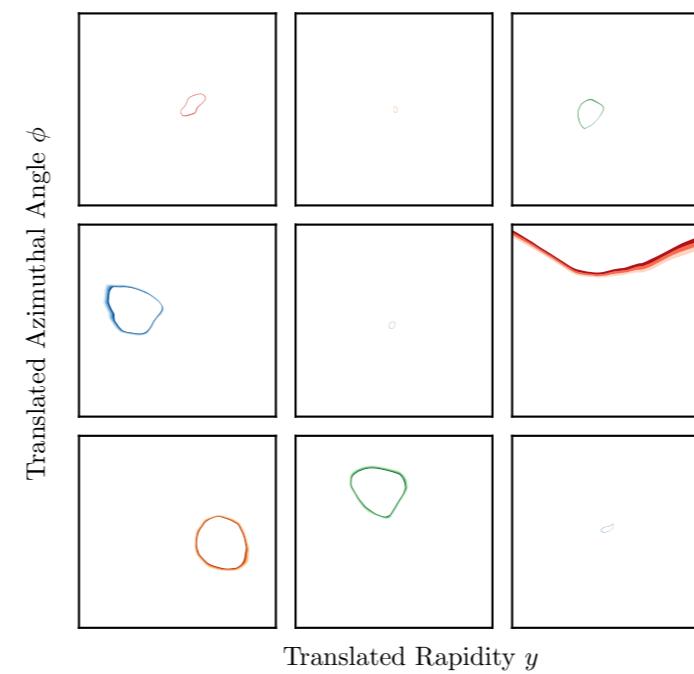


# Psychedelic Visualization

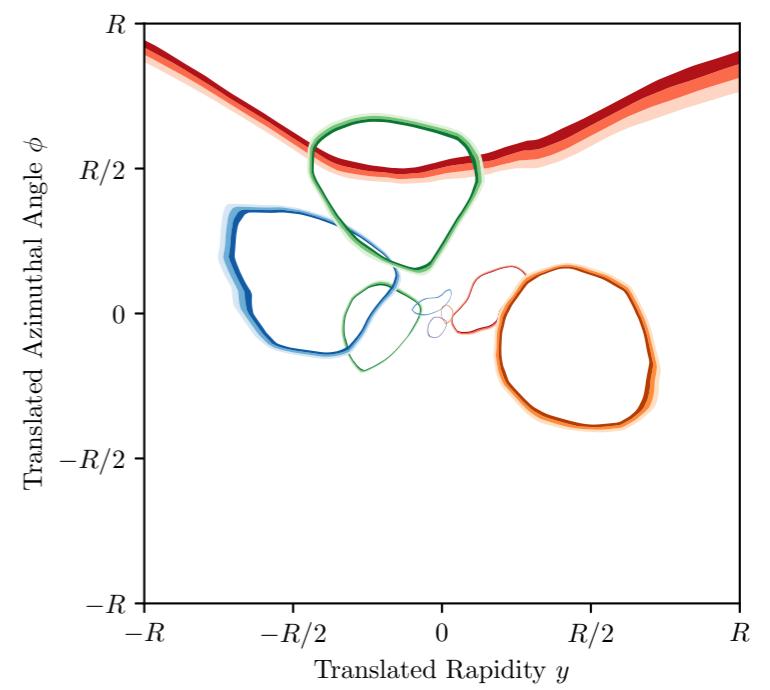
## Latent Filters



## 50% Contours

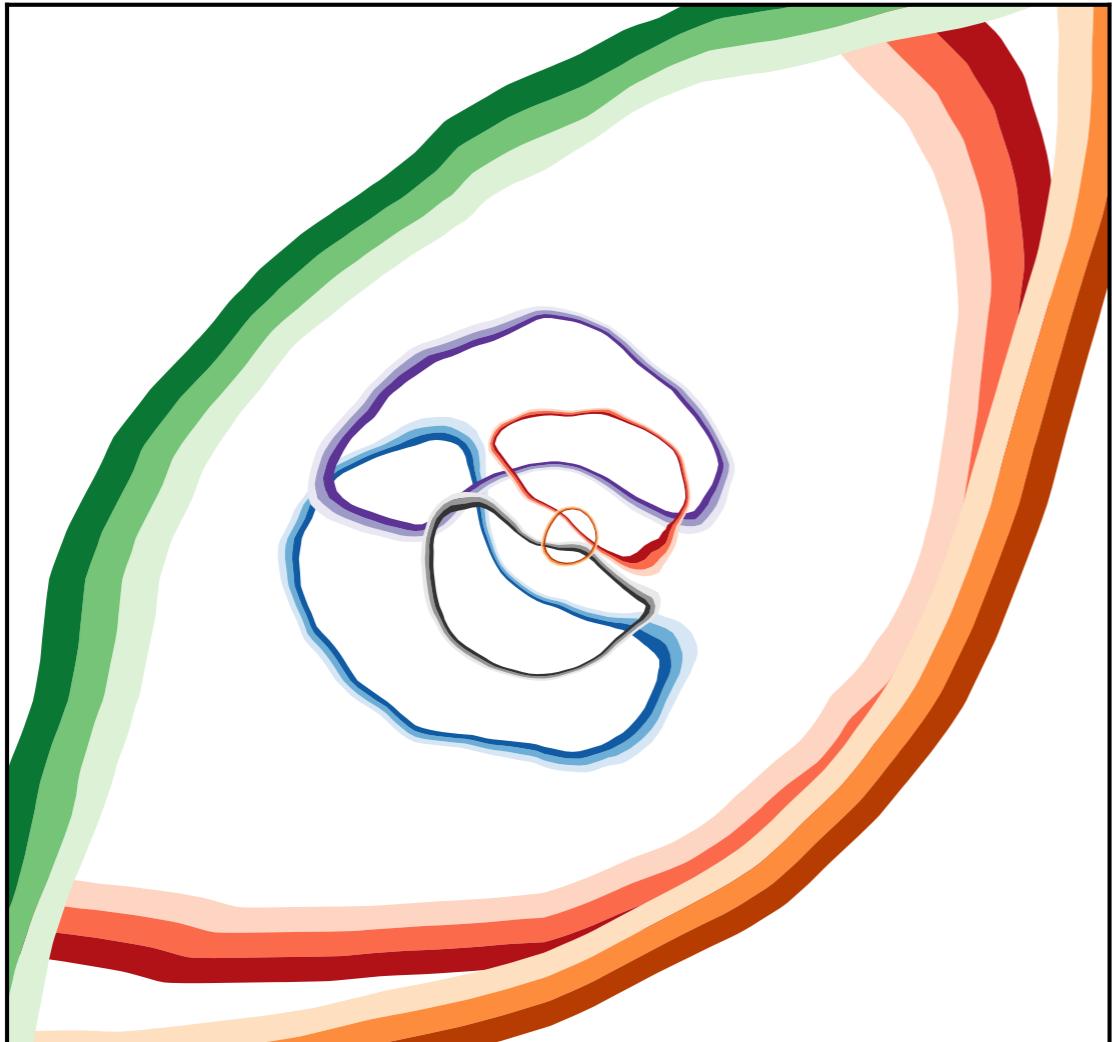


## Overlay

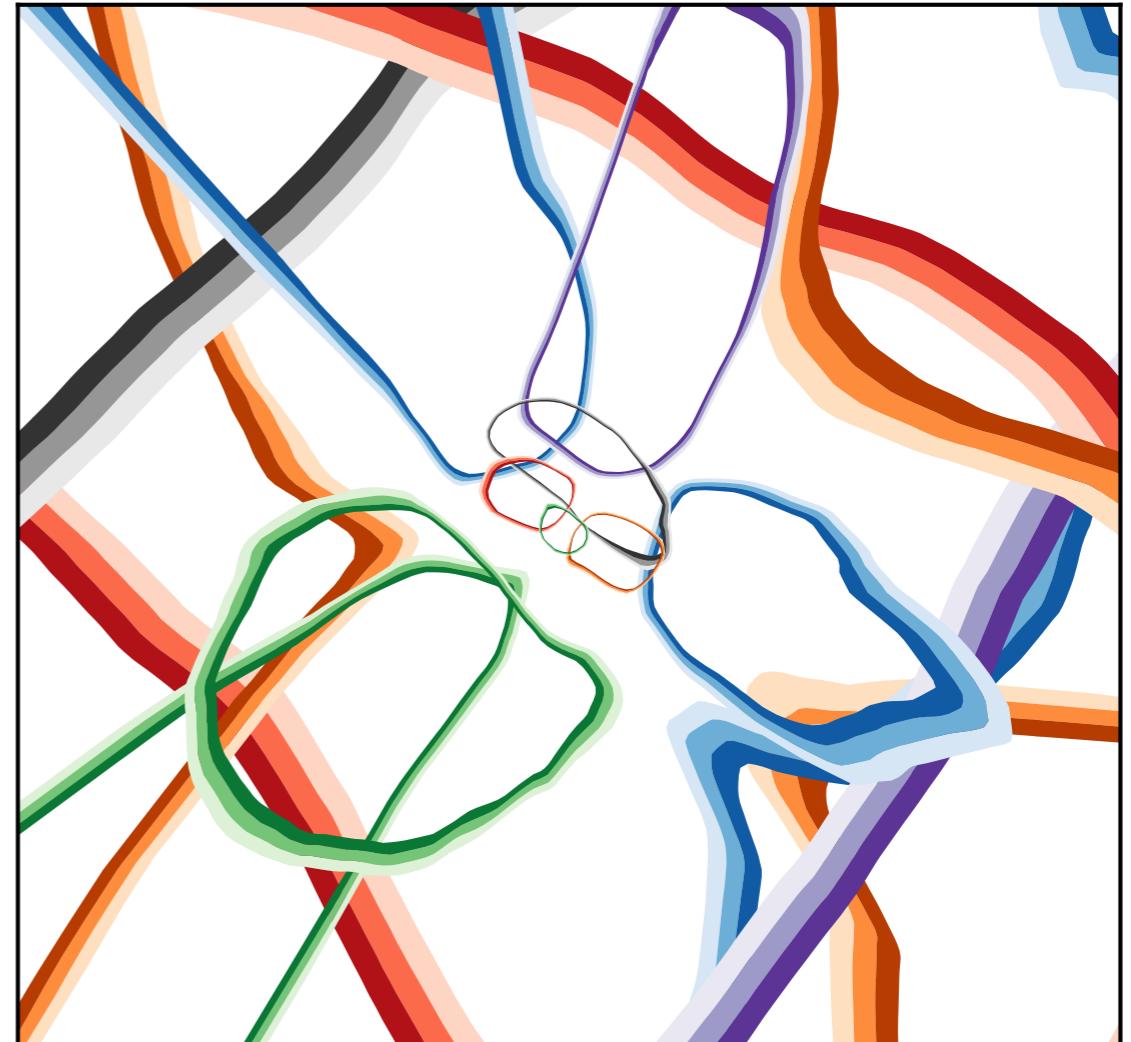


# Psychedelic Visualization

Latent Dimension 8

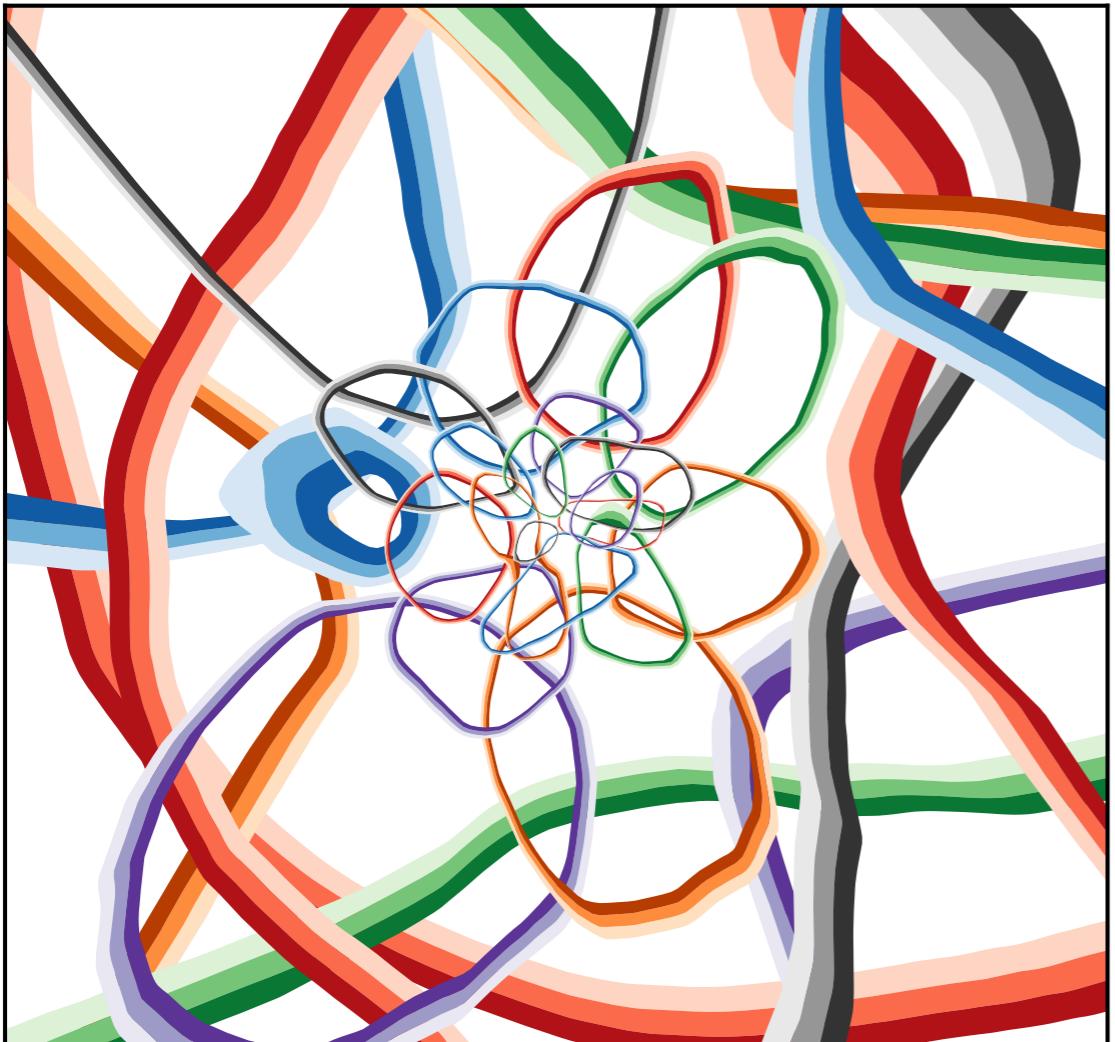


Latent Dimension 16

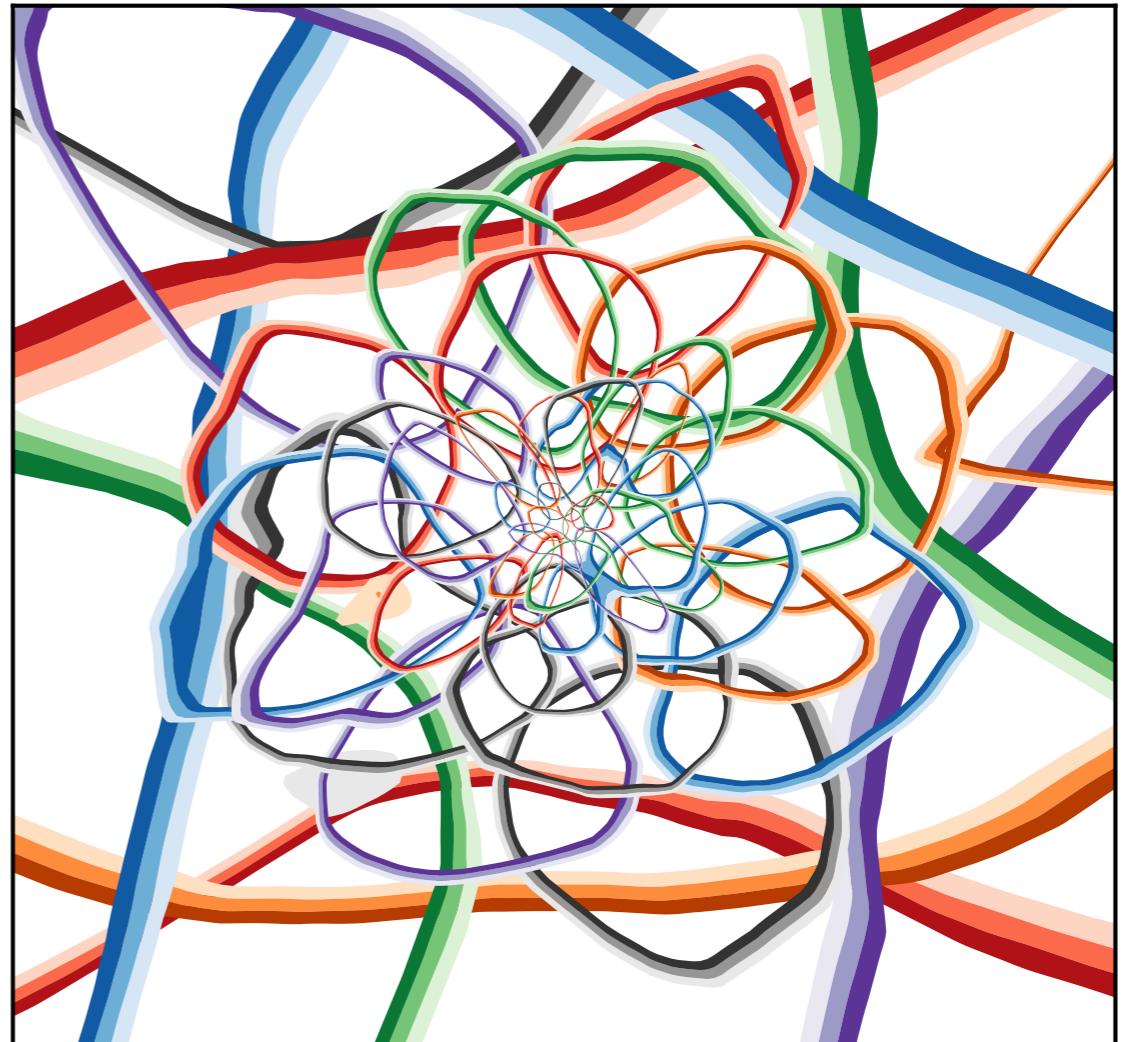


# Psychedelic Visualization

Latent Dimension 32

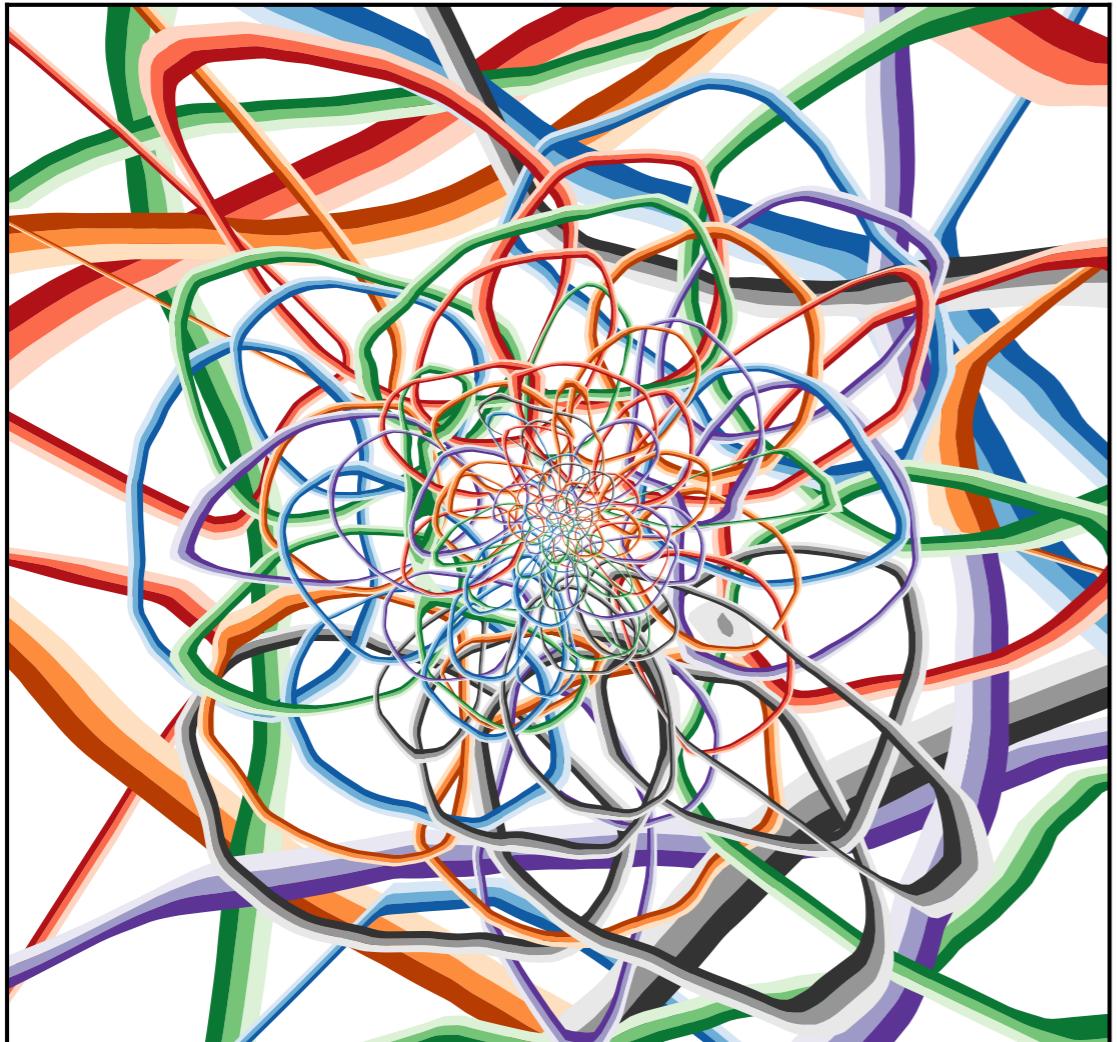


Latent Dimension 64

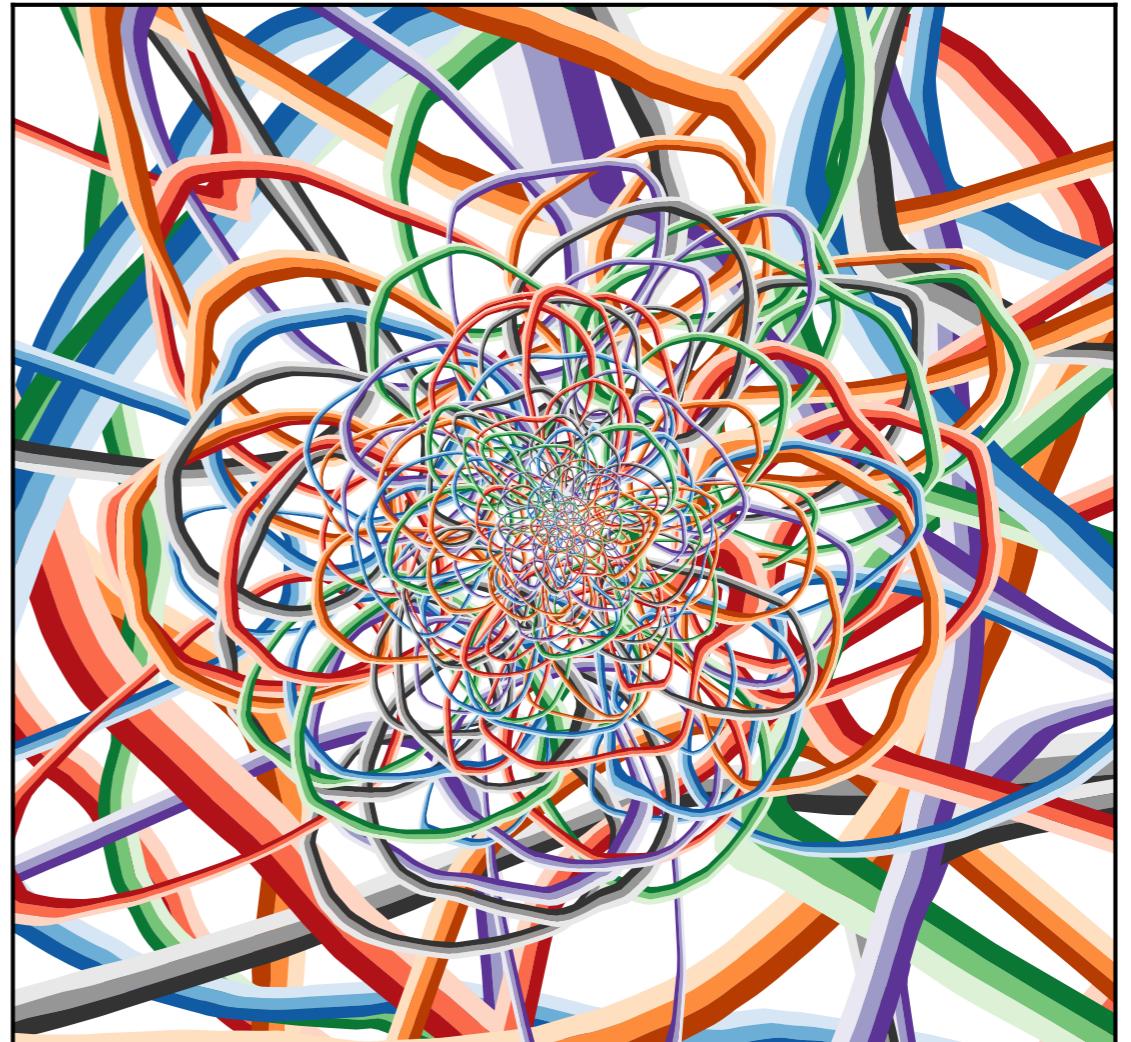


# Psychedelic Visualization

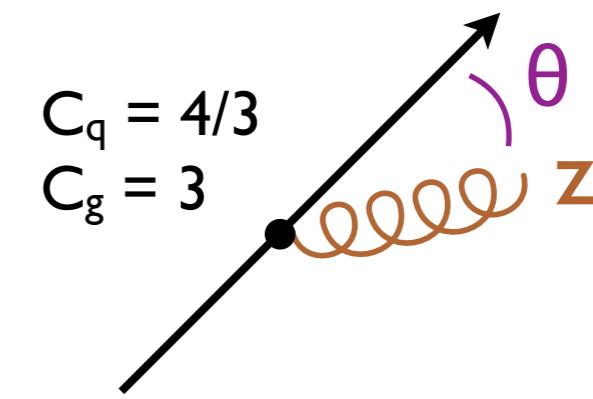
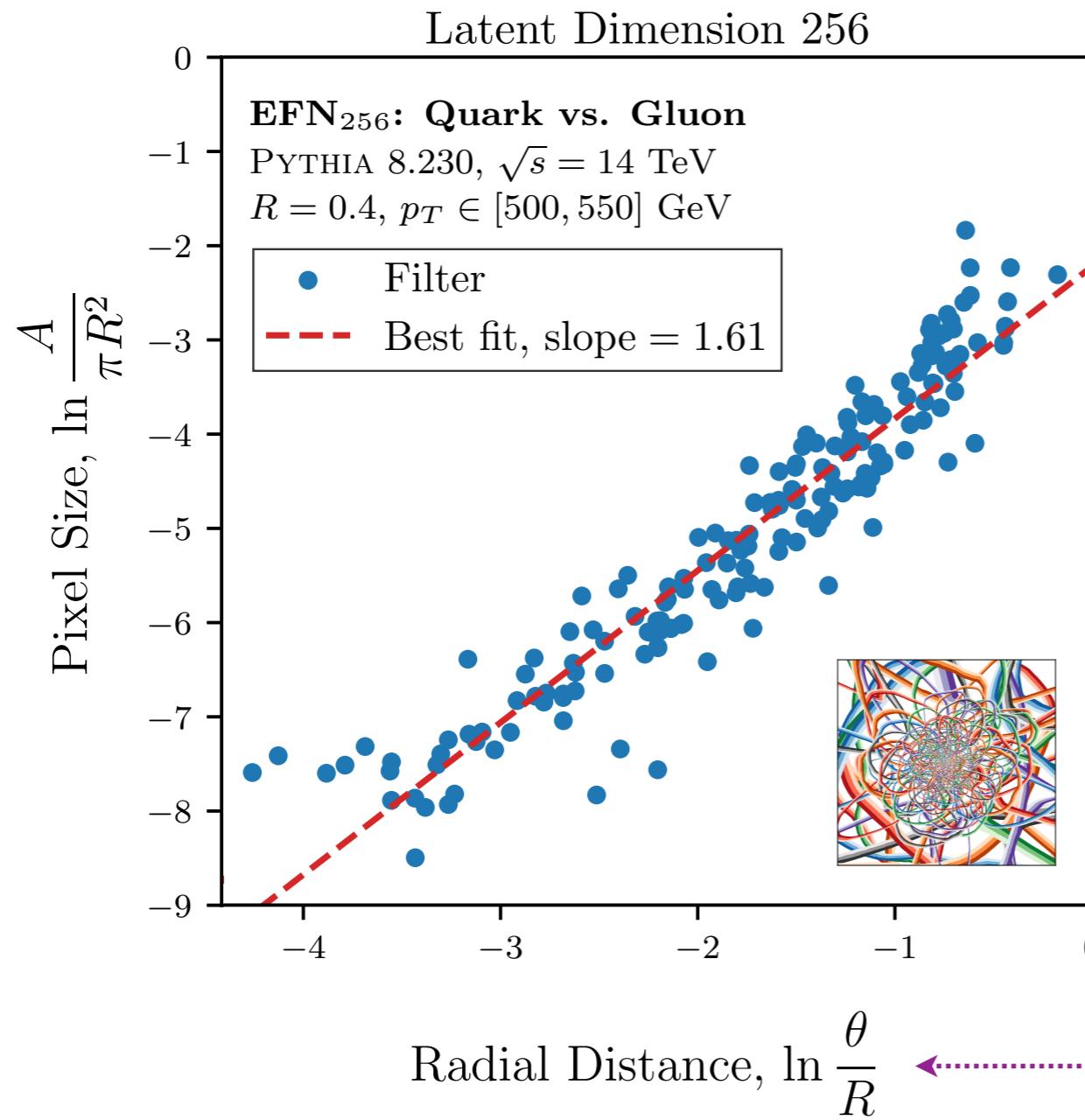
Latent Dimension 128



Latent Dimension 256



# Learning the Singularity Structure of QCD

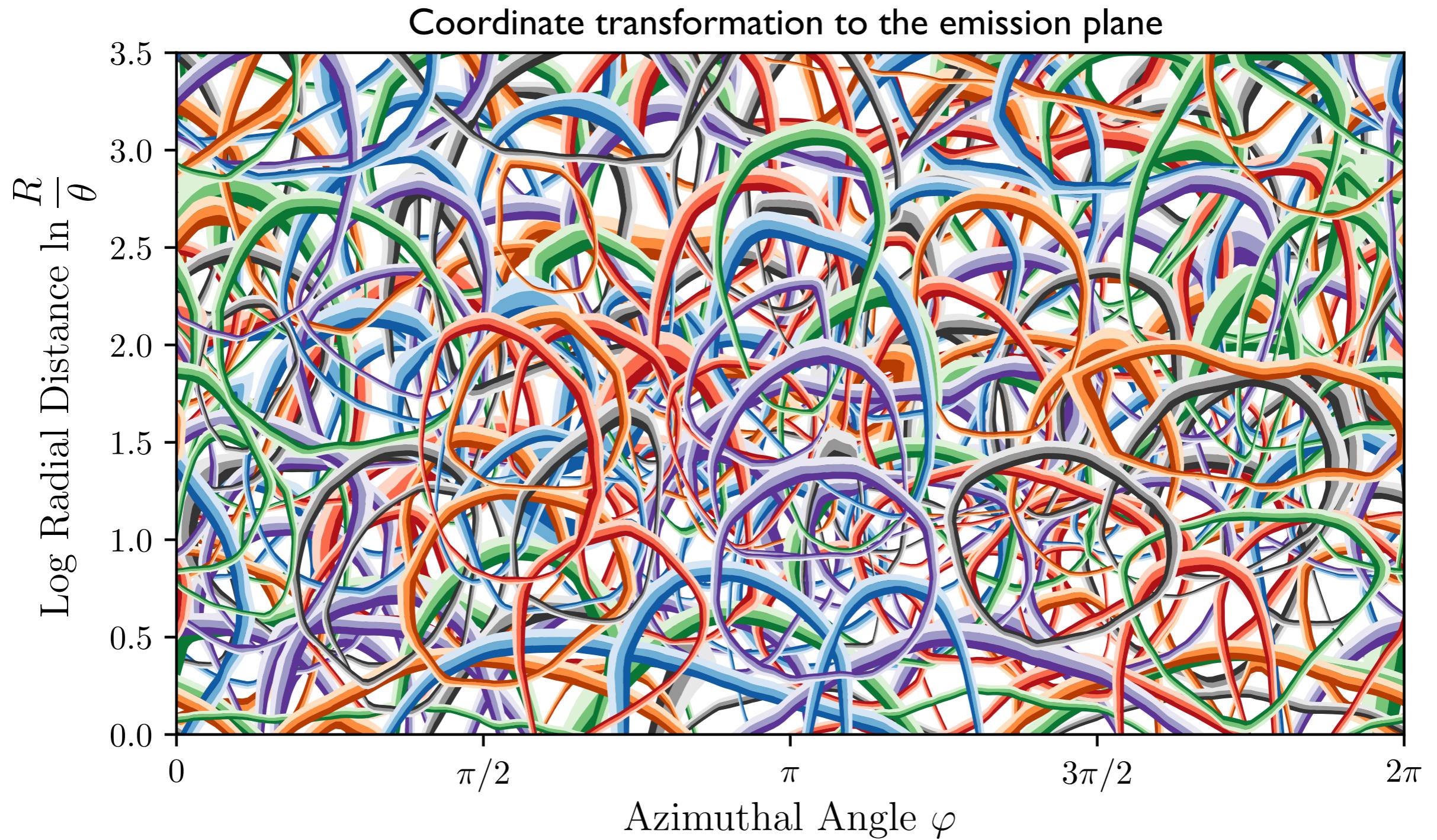


$$dP_{i \rightarrow ig} \sim \frac{2\alpha_s}{\pi} C_i \frac{d\theta}{\theta} \frac{dz}{z}$$

Collinear      Soft

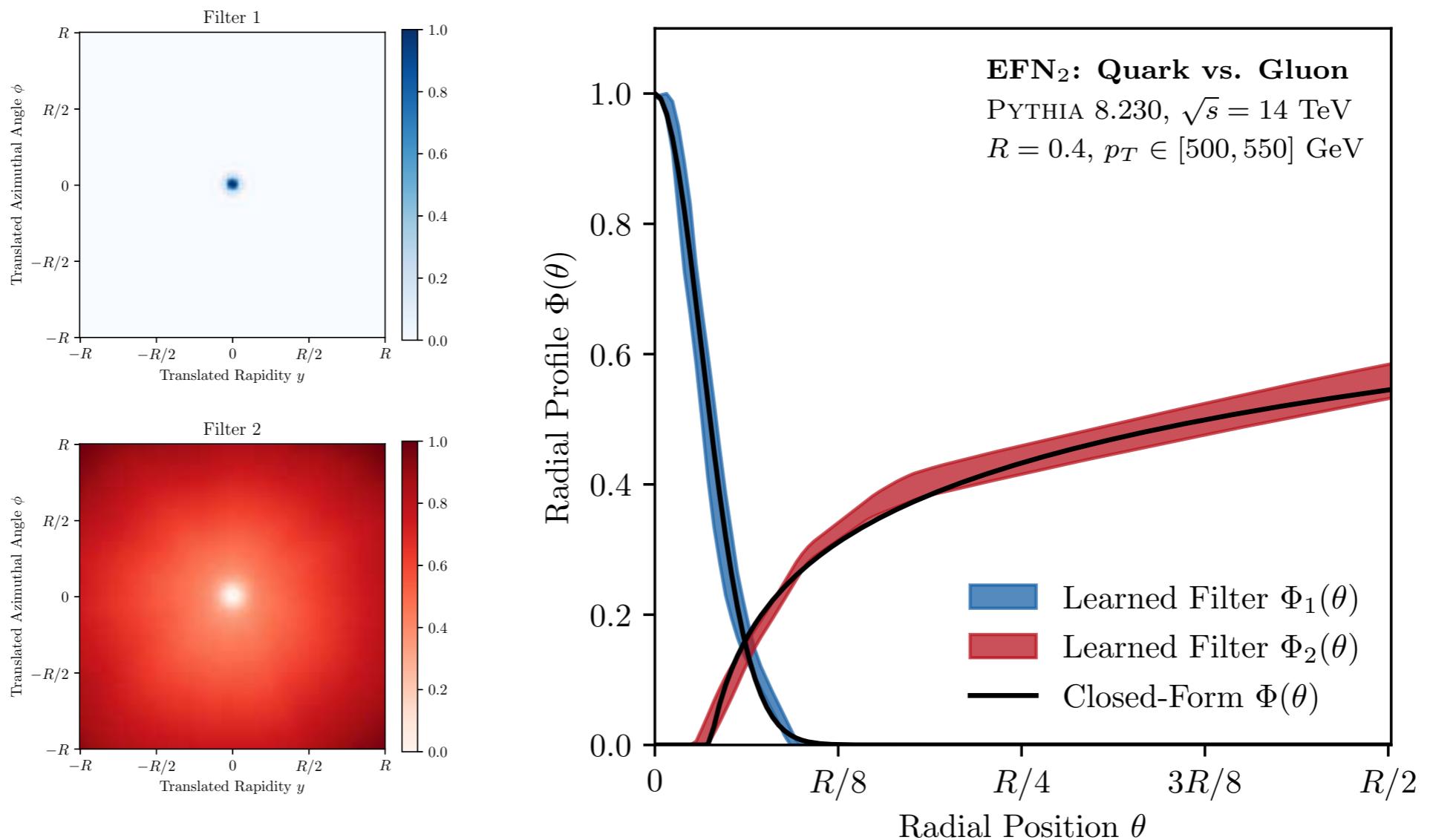
A dotted purple arrow points from the text "Collinear" to the purple bracket under "dθ/θ". Another dotted purple arrow points from the text "Soft" to the orange bracket under "dz/z".

# Ready for the MoMA



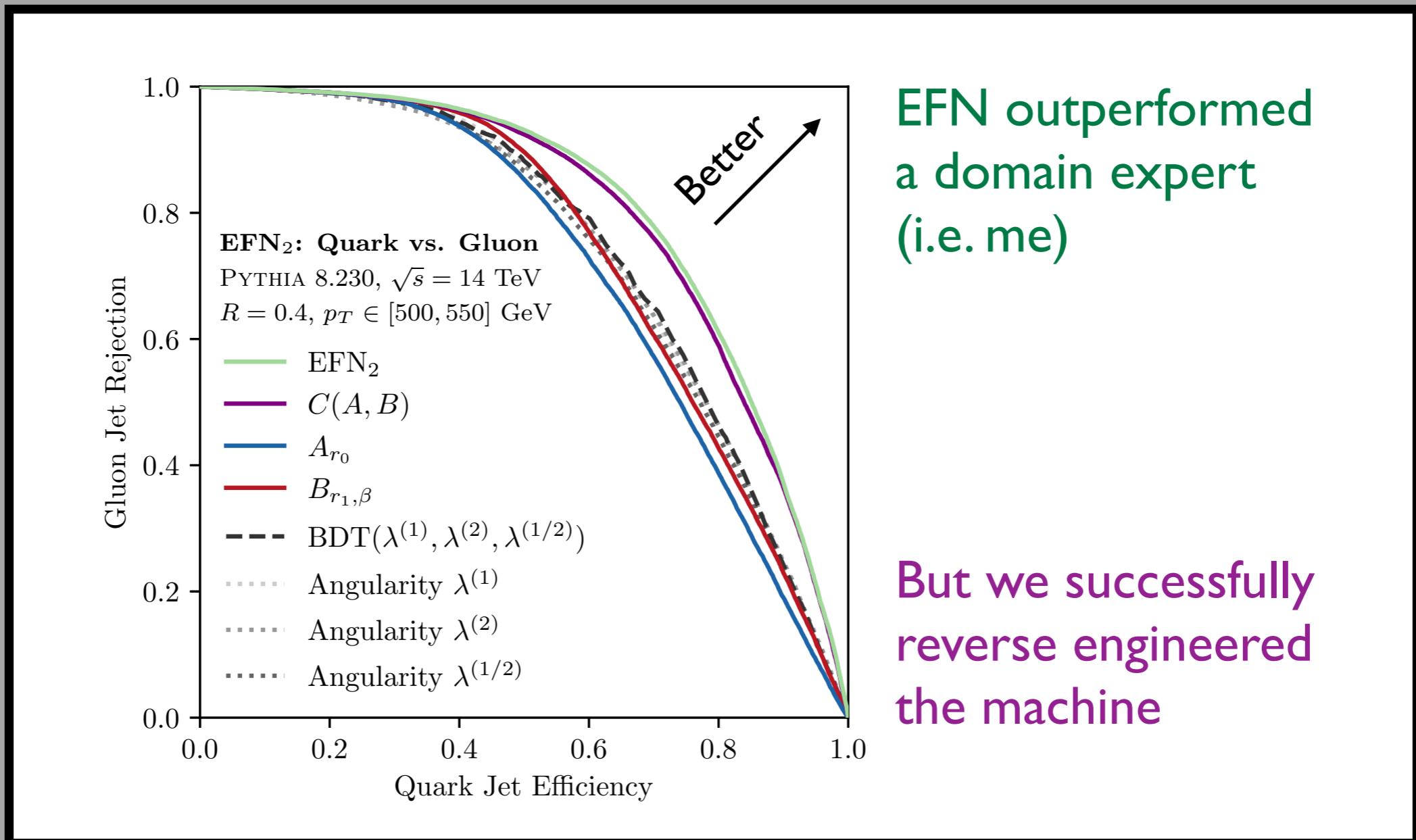
# “What is the Machine Learning?”

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



# “What is the Machine Learning?”

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



# The Broader Lesson

“Deep Learning”

&

~~vs.~~

“Deep Thinking”

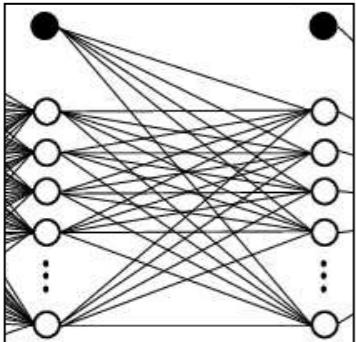
*Advances in mathematics and computer science (Deep Sets)*



*Advances in collider physics (EFN/PFN)*

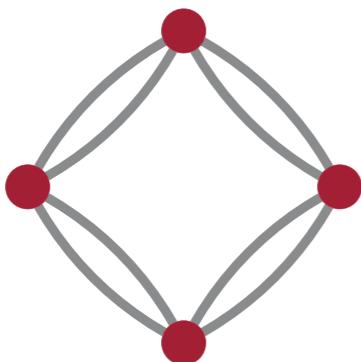
*Weighted Point Sets*    $\Leftrightarrow$    *IRC Safety*

# Summary



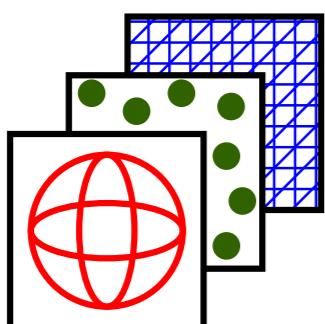
## Into the Network

*Embracing the rise of machine learning for collider physics*



## Symmetries & Safety

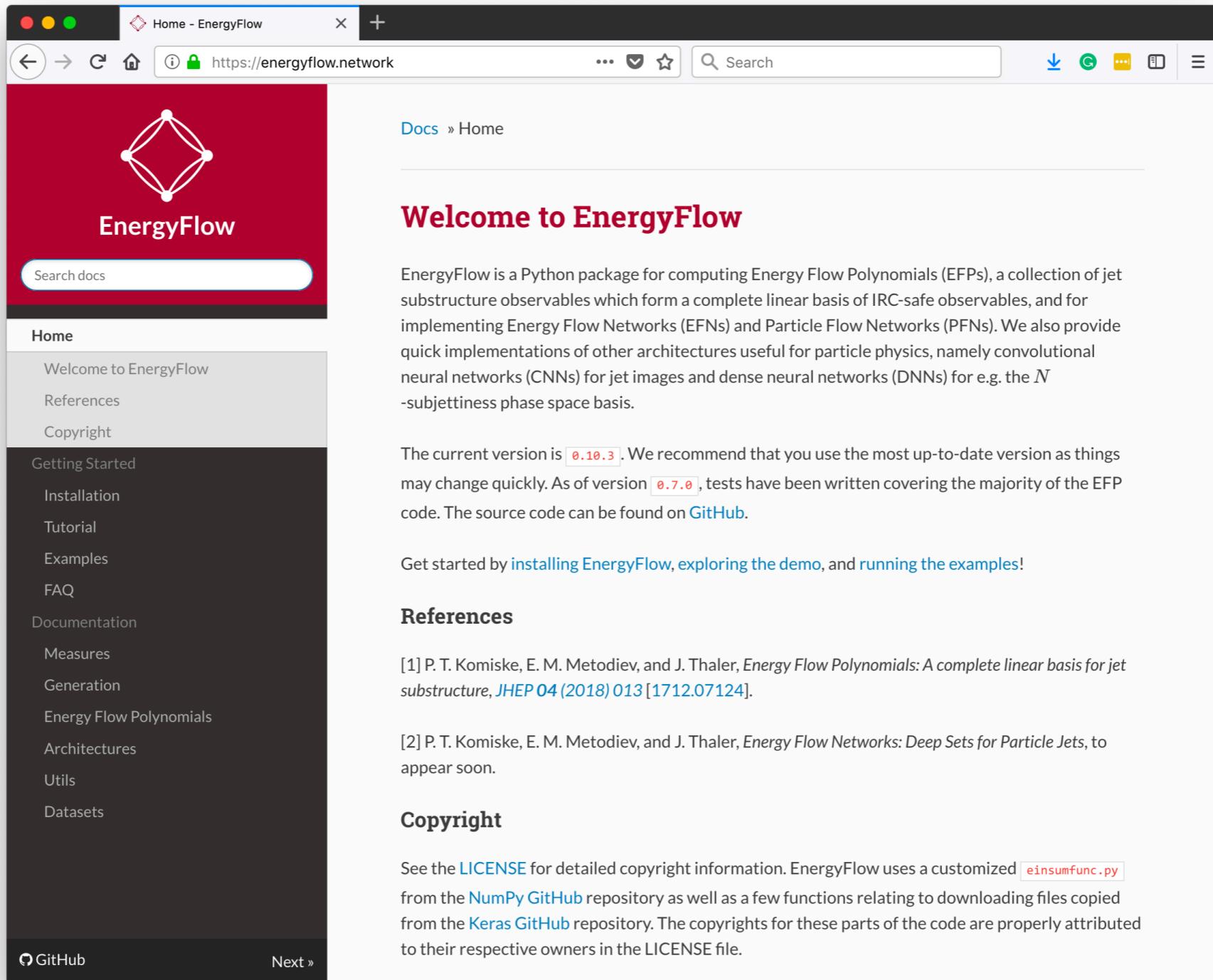
*Importance of indistinguishability and energy weighting*



## Deep Sets for Particle Jets

*EFN/PFNs as a step towards robustness, versatility, and transparency*

# energyflow.network



The screenshot shows a web browser displaying the [EnergyFlow](https://energyflow.network) documentation. The page has a dark theme with a red header. The header features the EnergyFlow logo (a diamond shape with internal lines) and the text "EnergyFlow". Below the header is a search bar labeled "Search docs". The main content area has a white background. At the top of the content area, there is a breadcrumb navigation showing "Docs » Home". The main title is "Welcome to EnergyFlow". The text describes EnergyFlow as a Python package for computing Energy Flow Polynomials (EFPs), which form a complete linear basis of IRC-safe observables, and for implementing Energy Flow Networks (EFNs) and Particle Flow Networks (PFNs). It also mentions quick implementations of convolutional neural networks (CNNs) for jet images and dense neural networks (DNNs) for e.g. the  $N$ -subjettiness phase space basis. A note indicates the current version is 0.10.3, and it recommends using the most up-to-date version due to rapid changes. It also mentions version 0.7.0 where tests have been written covering the majority of the EFP code, and that the source code can be found on [GitHub](#). Below this, there is a call to action: "Get started by [installing EnergyFlow](#), [exploring the demo](#), and [running the examples](#)!". There are sections for "References" and "Copyright". The "References" section lists two papers: [1] P. T. Komiske, E. M. Metodiev, and J. Thaler, *Energy Flow Polynomials: A complete linear basis for jet substructure*, *JHEP* 04 (2018) 013 [[1712.07124](#)]. and [2] P. T. Komiske, E. M. Metodiev, and J. Thaler, *Energy Flow Networks: Deep Sets for Particle Jets*, to appear soon. The "Copyright" section explains that EnergyFlow uses a customized `einsumfunc.py` from the [NumPy GitHub](#) repository and a few functions from the [Keras GitHub](#) repository, with copyrights properly attributed in the LICENSE file.

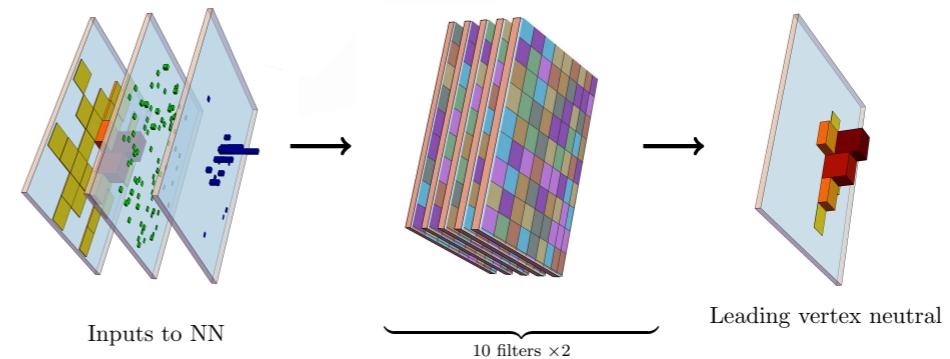
# *Backup Slides*

# Beyond Classification

## PUMML

### *Pileup Mitigation*

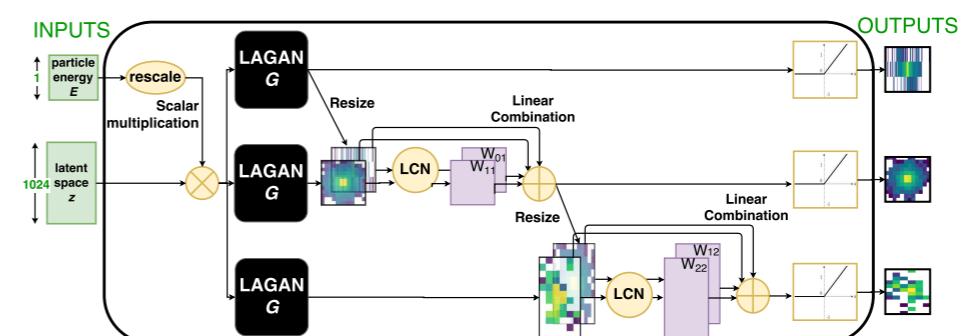
[Komiske, Metodiev, Nachman, Schwartz, 1707.08600]



## CaloGAN

### *Fast Detector Simulation*

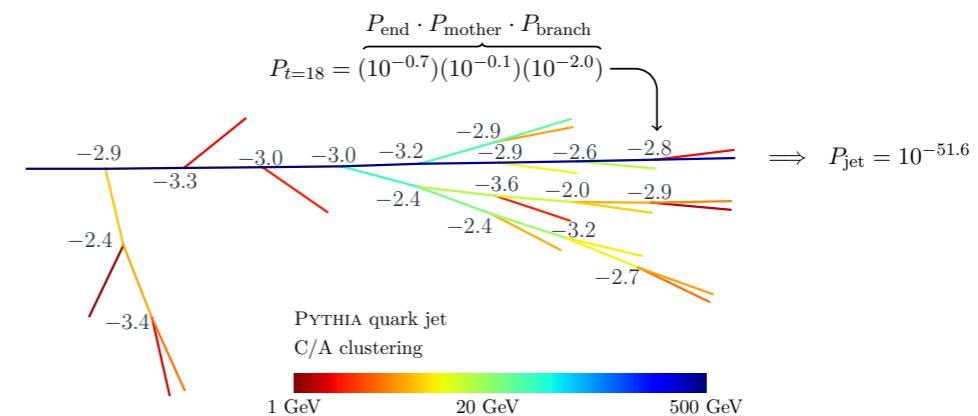
Paganini, de Oliveira, Nachman, 1705.02355, 1712.10321;  
see also de Oliveira, Michela Paganini, Nachman, 1701.05927]



## JUNIPR

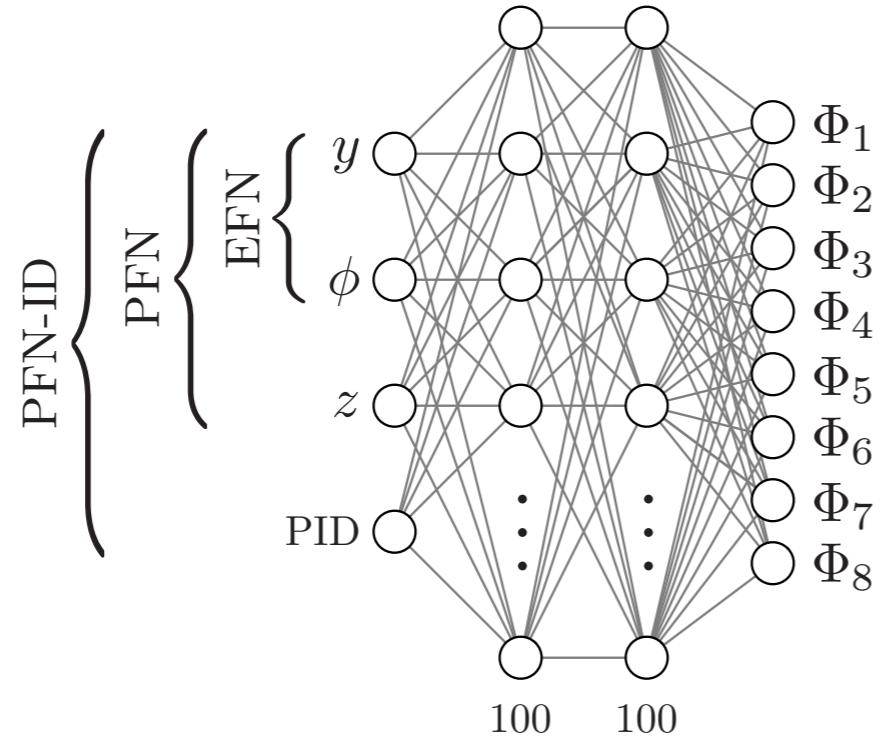
### *Probability Modeling*

[Andreassen, Feige, Frye, Schwartz, 1804.09720;  
see also Monk, 1807.03685]

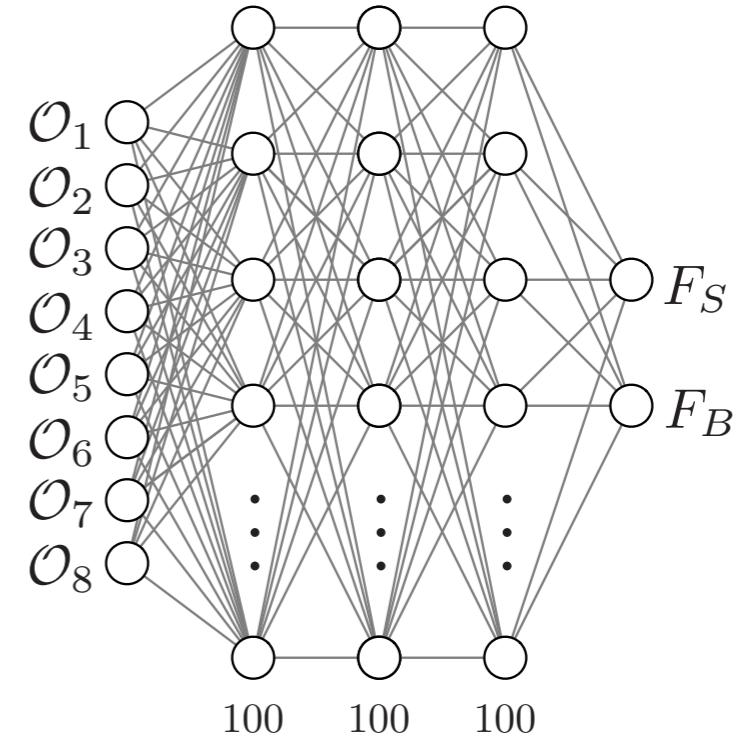


# Architecture Details

## Per-Particle:



## Latent Combiner: F



## Per-Jet (Latent Space):

$$\text{EFN: } \mathcal{O}_a = \sum_{i \in \text{jet}} z_i \Phi_a(y_i, \phi_i)$$

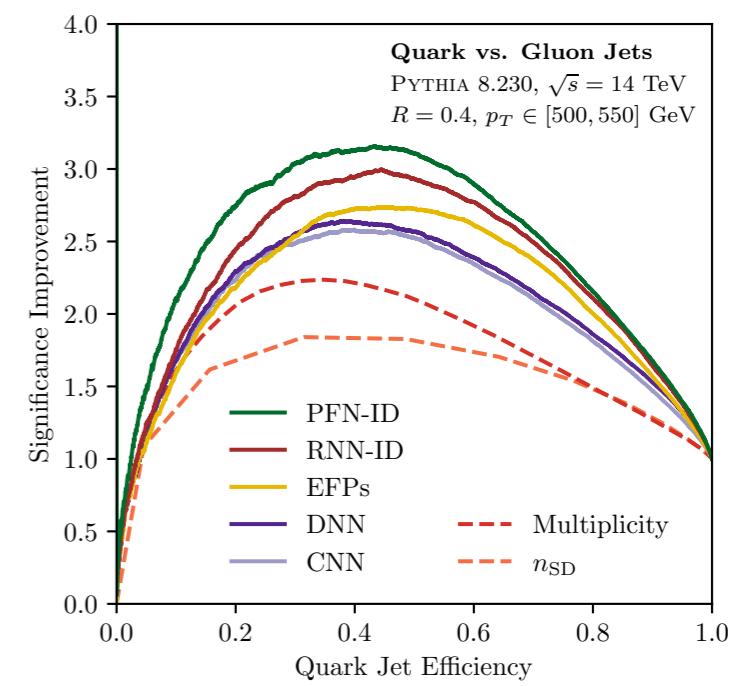
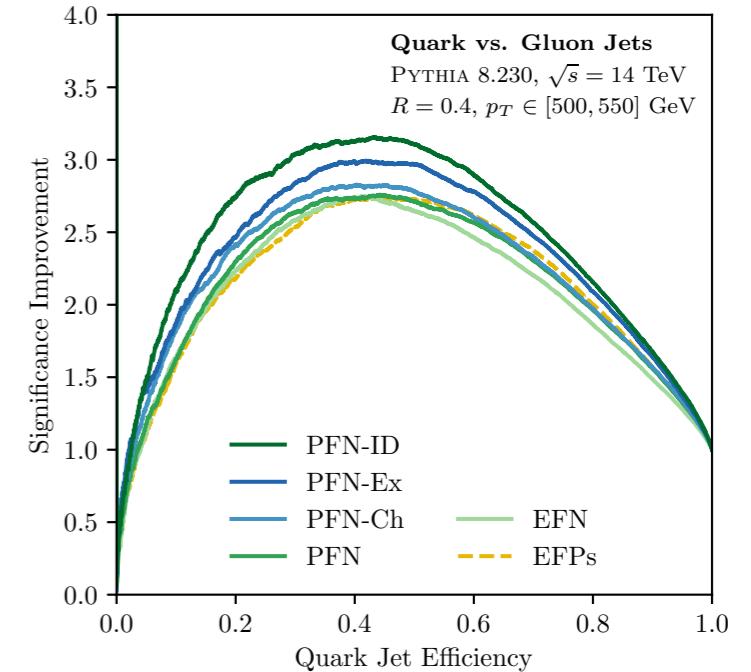
$$\text{PFN: } \mathcal{O}_a = \sum_{i \in \text{jet}} \Phi_a(y_i, \phi_i, z_i, \text{PID}_i)$$

## Final Discriminant:

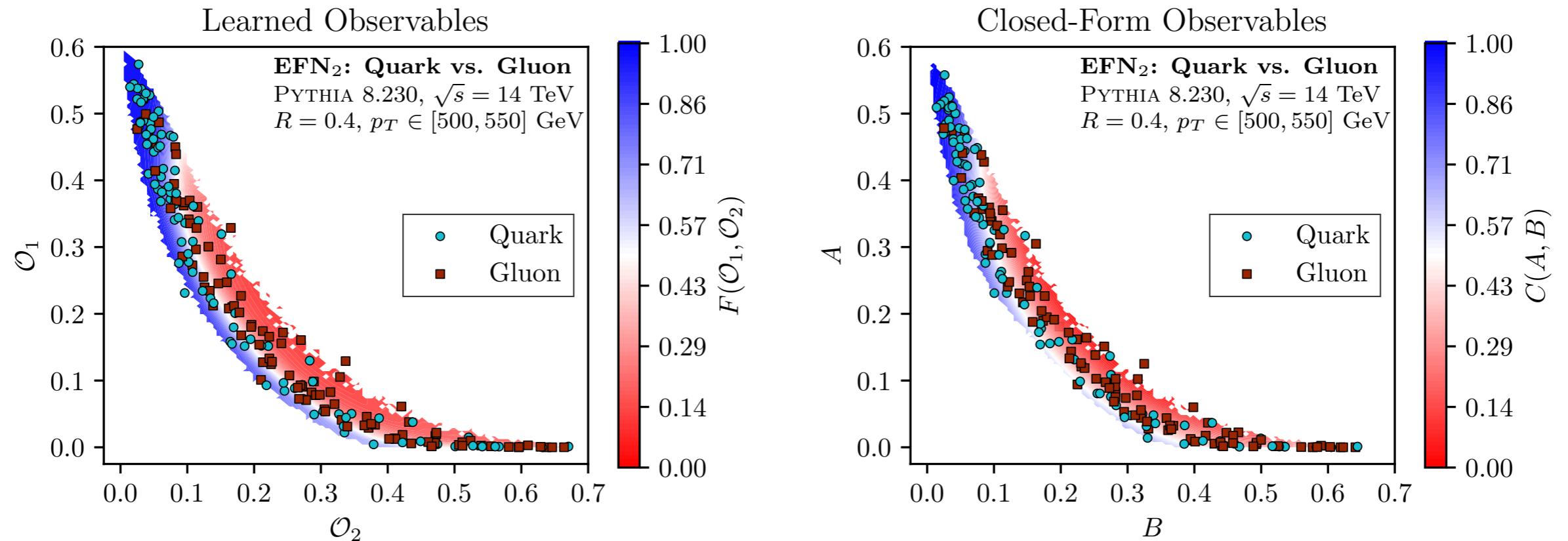
$$\text{softmax}(F_S, F_B)$$

# More Quark/Gluon Performance

Model	AUC	$1/\varepsilon_g$ at $\varepsilon_q = 50\%$
PFN-ID	<b>0.9052</b> $\pm 0.0007$	<b>37.4</b> $\pm 0.7$
PFN-Ex	0.9005 $\pm 0.0003$	34.7 $\pm 0.4$
PFN-Ch	0.8924 $\pm 0.0001$	31.2 $\pm 0.3$
PFN	0.8911 $\pm 0.0008$	30.8 $\pm 0.4$
EFN	0.8824 $\pm 0.0005$	28.6 $\pm 0.3$
RNN-ID	0.9010	34.4
RNN	0.8899	30.5
EFP	0.8919	29.7
DNN	0.8849	26.4
CNN	0.8781	25.5
$M$	0.8401	19.0
$n_{SD}$	0.8297	14.2
$m$	0.7401	7.2



# Reverse Engineering the Machine



Fascinating QCD question about why this is a better strategy than traditional angularities

# *Energy Flow Polynomials*

- Underlying Physics
- Natural Data Representation
- Suitable Algorithm

What is the space of *all*  
IRC-safe observables?

# Examples from Jet Substructure

→ Underlying Physics  
Natural Data Representation  
Suitable Algorithm

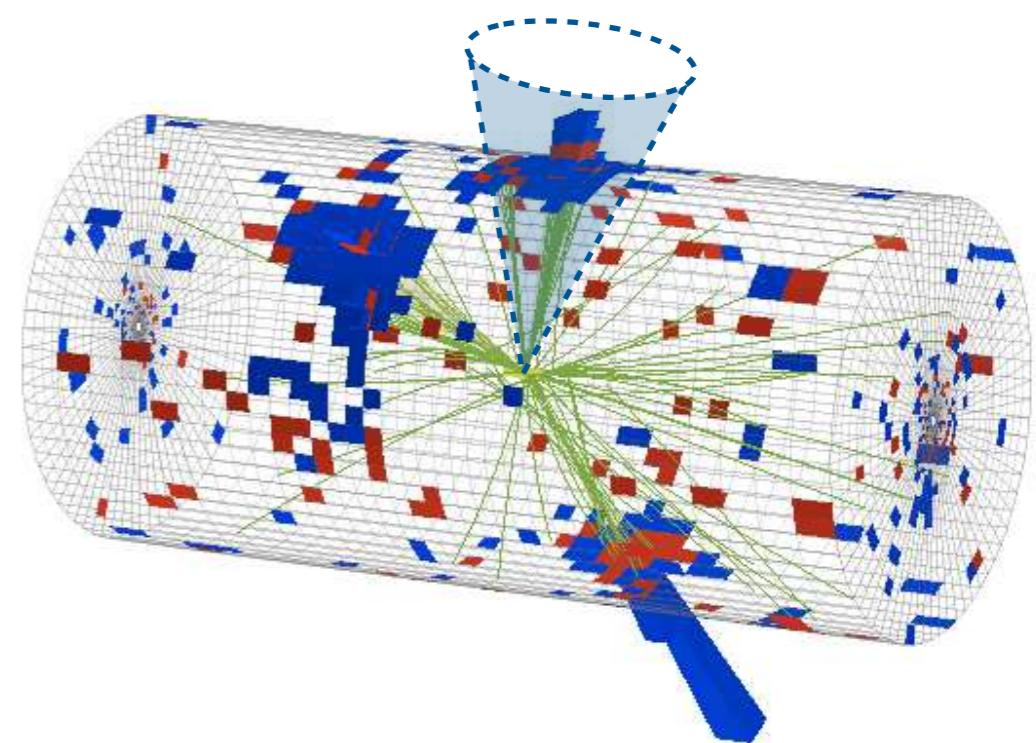
Jet pt:  $\sum_{i \in \text{jet}} p_{T,i}$  **IRC Safe**

$p_T^D$ :  
[CMS HIG-11-027]  $\sum_{i \in \text{jet}} \frac{p_{T,i}^2}{p_{T\text{jet}}^2}$  **IR Safe**  
**C Unsafe**

Multiplicity:  $\sum_{i \in \text{jet}} 1$  **IRC Unsafe**

Jet Mass:  $\sum_{i,j \in \text{jet}} p_i \cdot p_j$  **IRC Safe**

N-subjettiness:  
[JDT, Van Tilburg,  
1011.2268, 1108.2701]  $\sum_{i \in \text{jet}} p_{T,i} \min \{ \Delta R_{i,1}, \Delta R_{i,2}, \dots, \Delta R_{i,N} \}^\beta$  **IRC Safe**  
But ratios are only “Sudakov safe”!



# A Systematic Expansion

Underlying Physics  
→ Natural Data Representation  
Suitable Algorithm

Expand\* any IRC safe observable in small energy limit

$$\begin{aligned} \mathcal{S} = & \sum_i E_i f_1^{\mathcal{S}}(\hat{n}_i) + \sum_{ij} E_i E_j f_2^{\mathcal{S}}(\hat{n}_i, \hat{n}_j) \\ & + \sum_{ijk} E_i E_j E_k f_3^{\mathcal{S}}(\hat{n}_i, \hat{n}_j, \hat{n}_k) + \dots \end{aligned}$$

Form enforced by:	Particle Relabeling	Infrared Safety	Collinear Safety
-------------------	---------------------	-----------------	------------------

Further expand\* each angular function in pairwise angles

$$z_i = \frac{E_i}{E_{\text{jet}}} \quad \cos \theta_{ij} = \hat{n}_i \cdot \hat{n}_j$$

[Komiske, Metodiev, JDT, 1712.07124; see also Tkachov, hep-ph/9601308]

# The Energy Flow Polynomials

Underlying Physics  
→ Natural Data Representation  
Suitable Algorithm

$$\text{EFP}_G = \sum_{i_1=1}^M \cdots \sum_{i_N=1}^M z_{i_1} \cdots z_{i_N} \prod_{(k,\ell) \in G} \theta_{i_k i_\ell}$$

*Multigraph* *Angular Scaling*

$\downarrow$   $\downarrow \beta$

All N-tuples N Energy Fractions

Polynomial in Pairwise Angles

## A Linear Basis for Jet Substructure (!)

[Komiske, Metodiev, JDT, 1712.07124]

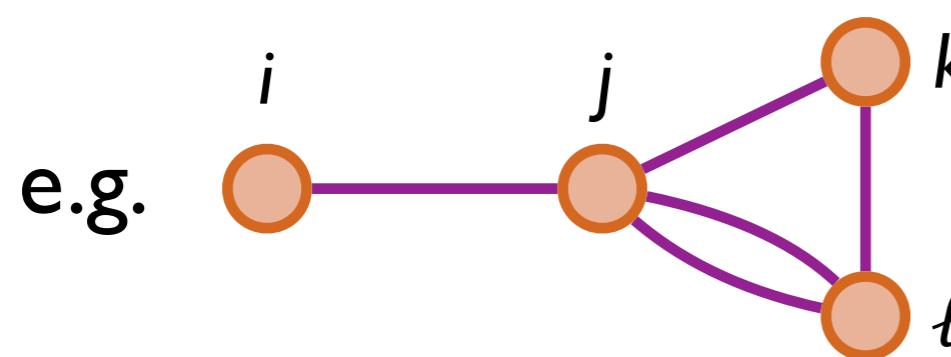
# The Energy Flow Polynomials

Underlying Physics  
 → Natural Data Representation  
 Suitable Algorithm

$$\begin{array}{c}
 \text{Multigraph} \\
 \downarrow \\
 \text{EFP}_G = \sum_{i_1=1}^M \cdots \sum_{i_N=1}^M z_{i_1} \cdots z_{i_N} \prod_{(k,\ell) \in G} \theta_{i_k i_\ell}^\beta
 \end{array}$$

$\text{All N-tuples}$	$\text{N Energy Fractions}$	$\text{Polynomial in Pairwise Angles}$
-----------------------	-----------------------------	--

e.g.



$$= \sum_{ijkl} z_i z_j z_k z_l \theta_{ij} \theta_{jk} \theta_{jl}^2 \theta_{kl}$$

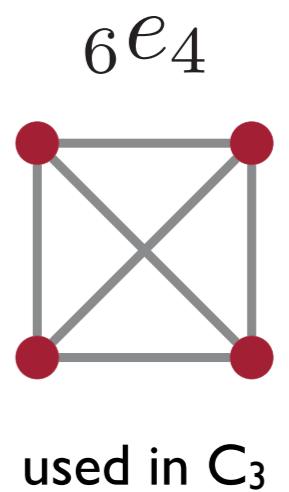
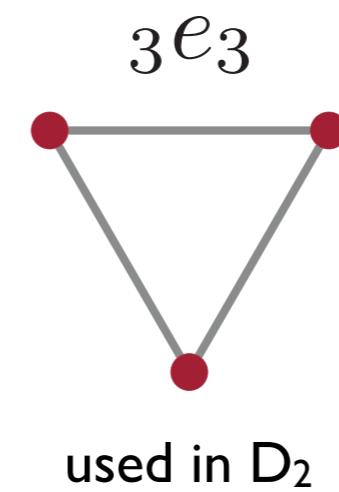
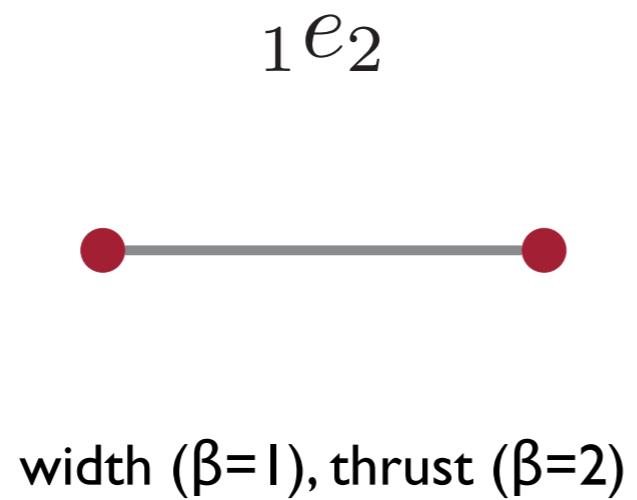
## A Linear Basis for Jet Substructure (!)

[Komiske, Metodiev, JDT, 1712.07124]

# Down the Rabbit Hole

Underlying Physics  
→ Natural Data Representation  
Suitable Algorithm

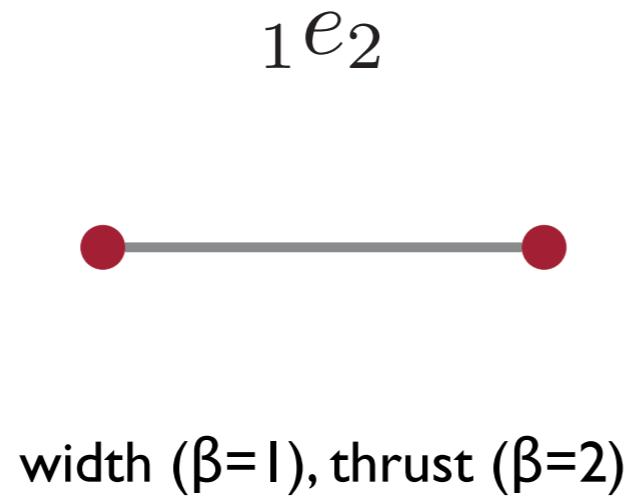
Known Structures:



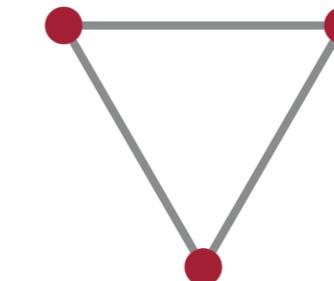
# Down the Rabbit Hole

Underlying Physics  
→ Natural Data Representation  
Suitable Algorithm

Known Structures:

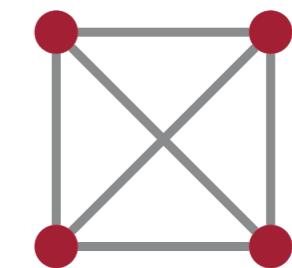


$3e_3$



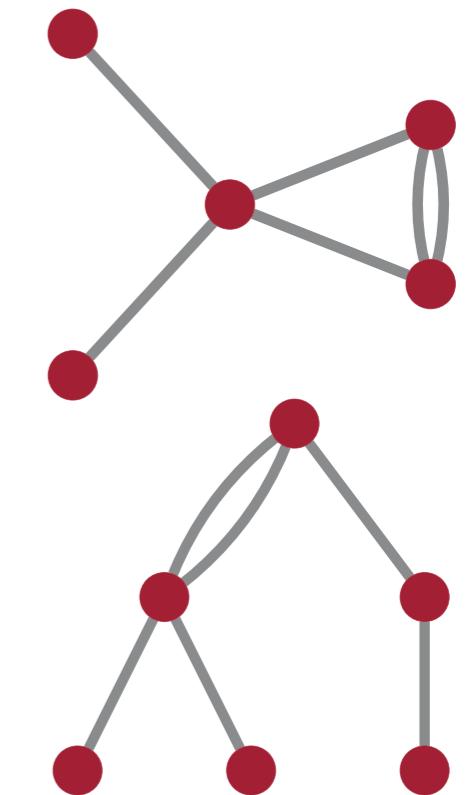
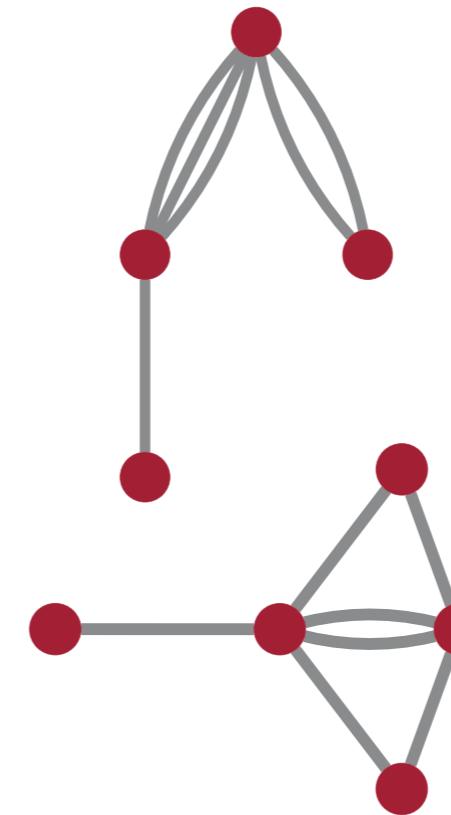
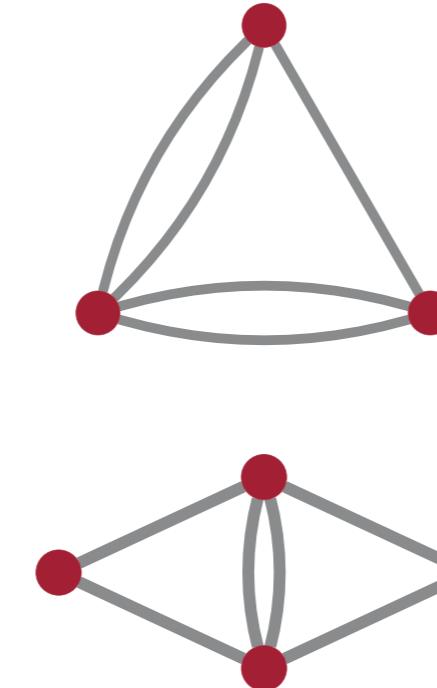
used in  $D_2$

$6e_4$



used in  $C_3$

No Idea:

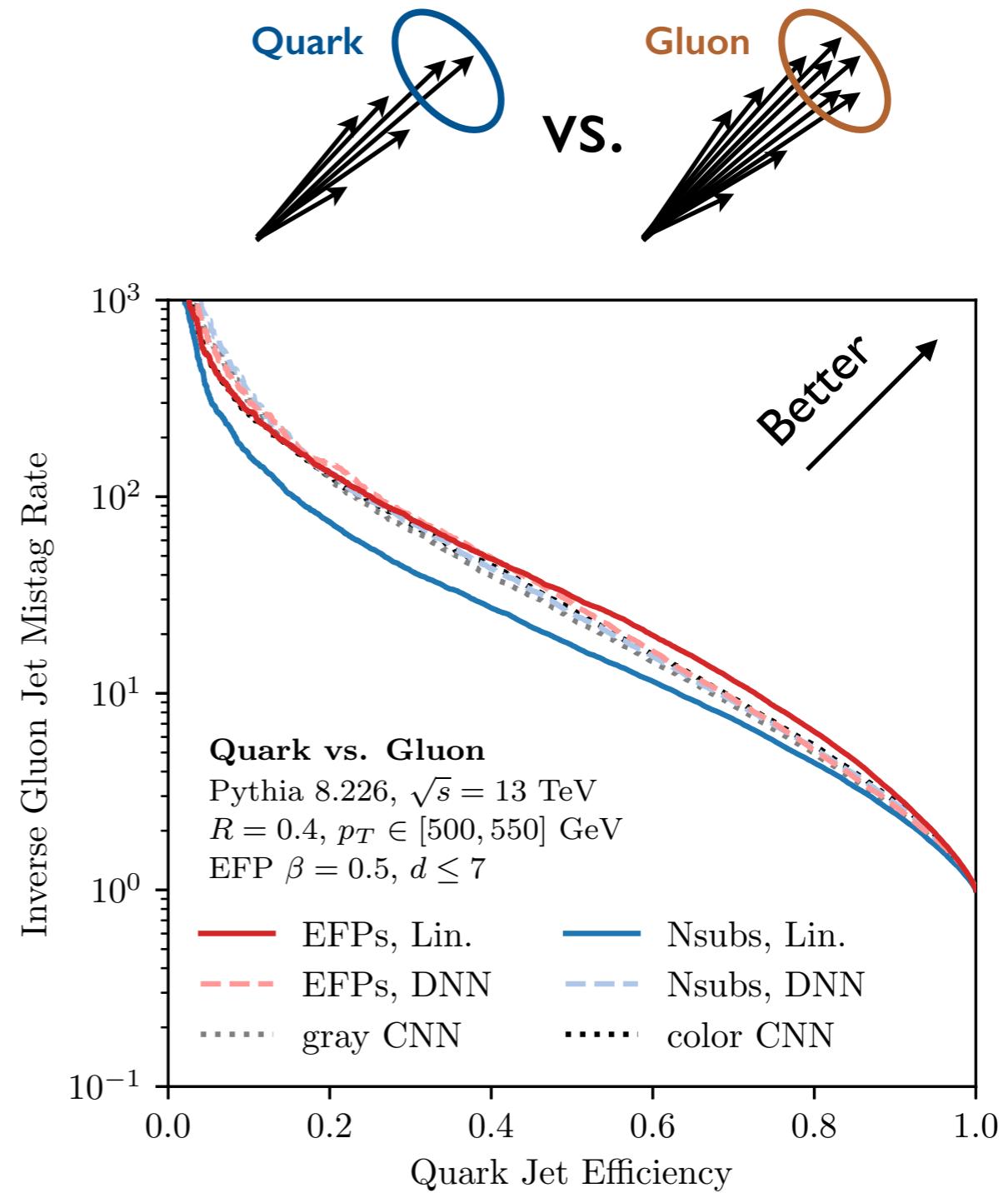
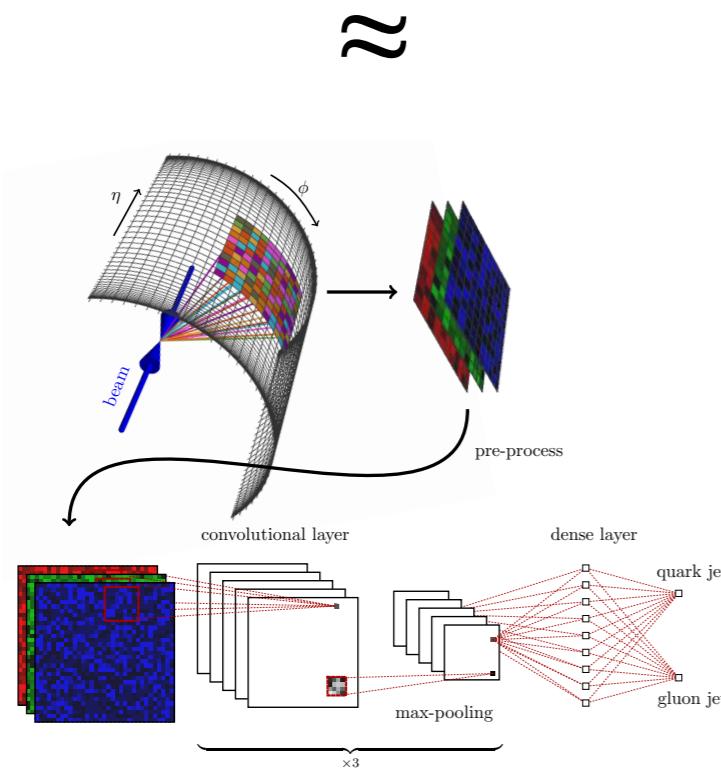


# Linear Regression $\approx$ CNN

Underlying Physics  
 Natural Data Representation  
 → Suitable Algorithm

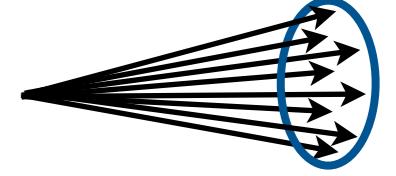
*“...indistinguishable from magic”*

$$\mathcal{S} = \sum_G s_G \text{EFP}_G$$



[Komiske, Metodiev, JDT, 1712.07124; Komiske, Metodiev, Schwartz, 1612.01551]

# Comparing Data Representations



Original 4-Vectors:  $\{p_1^\mu, p_2^\mu, \dots, p_N^\mu\}$

*Variable-length, unordered set*

EFP Basis:

$$\text{EFP}_G = \sum_{i_1=1}^M \cdots \sum_{i_N=1}^M z_{i_1} \cdots z_{i_N} \prod_{(k,\ell) \in G} \theta_{i_k i_\ell}$$

*Automatically permutation invariant  
Linear spanning basis (over-complete)  
Computational nightmare of  $O(M^N)$ ?*



Too naive, can use  
variable elimination

# The Physics-Meets-Computation Approach

Underlying Physics



Natural & Efficient Data Representation



Desired Computational Property

- Underlying Physics
- Natural Data Representation
- Suitable Algorithm

What is the space of *all*  
linearly-computable  
permutation-invariant  
IRC-safe observables?

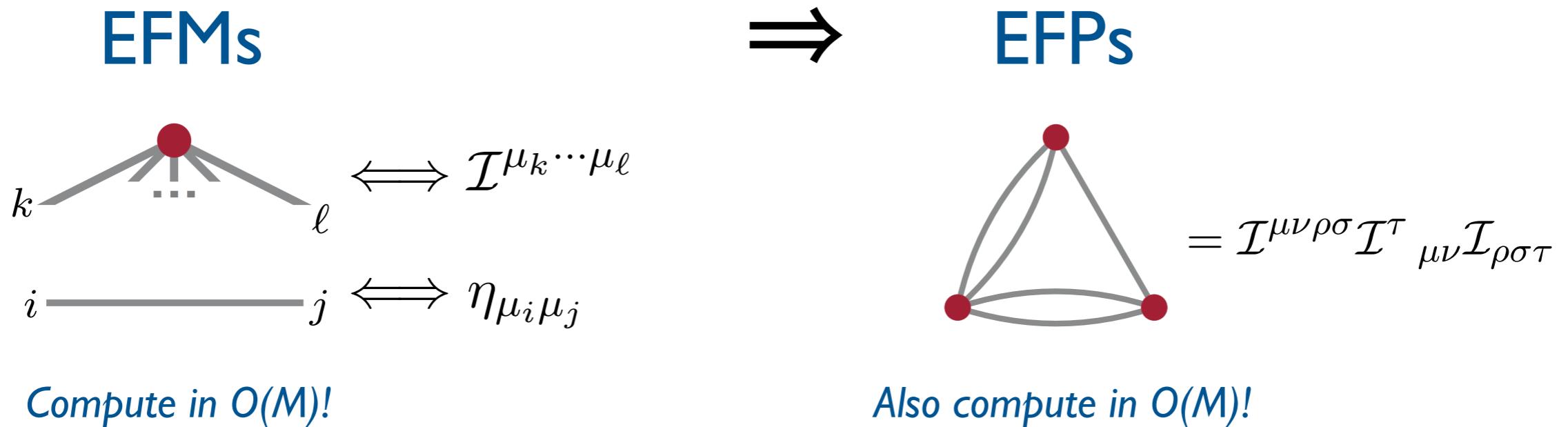
# The Energy Flow Moments

Underlying Physics  
 → Natural Data Representation  
 Suitable Algorithm

$$\mathcal{I}^{\mu_1 \mu_2 \cdots \mu_v} = \sum_{i=1}^M E_i \hat{p}^{\mu_1} \hat{p}^{\mu_2} \cdots \hat{p}^{\mu_v}$$

Particle  
Relabeling      Infrared  
Safety

Special Choice  
of Angle       $\theta_{ij} = 2 \eta_{\mu\nu} \hat{p}_i^\mu \hat{p}_j^\nu$



[Komiske, Metodiev, JDT, we've been promising this paper for 9 months]