

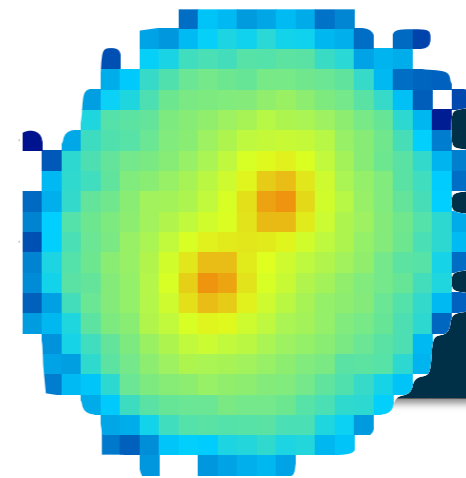
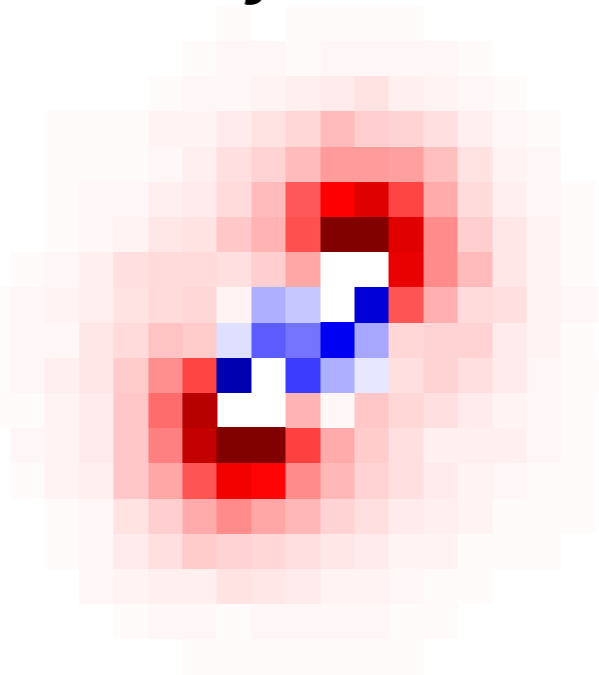
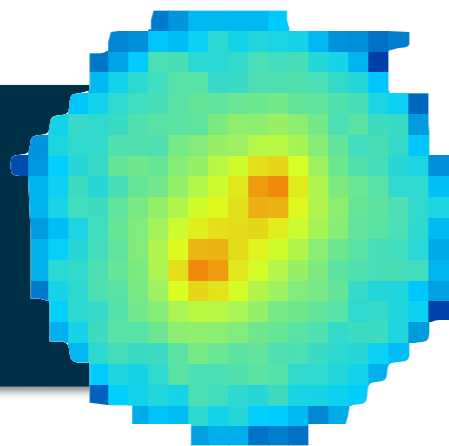
Modern Machine Learning

for Classification, Regression,
and Generation in Jet Physics



Benjamin Nachman

Lawrence Berkeley National Laboratory



CERN Data Science Seminar, November 14, 2017

High Energy Physics at the LHC

Center-of-mass energy = 13 TeV



Run: 302347

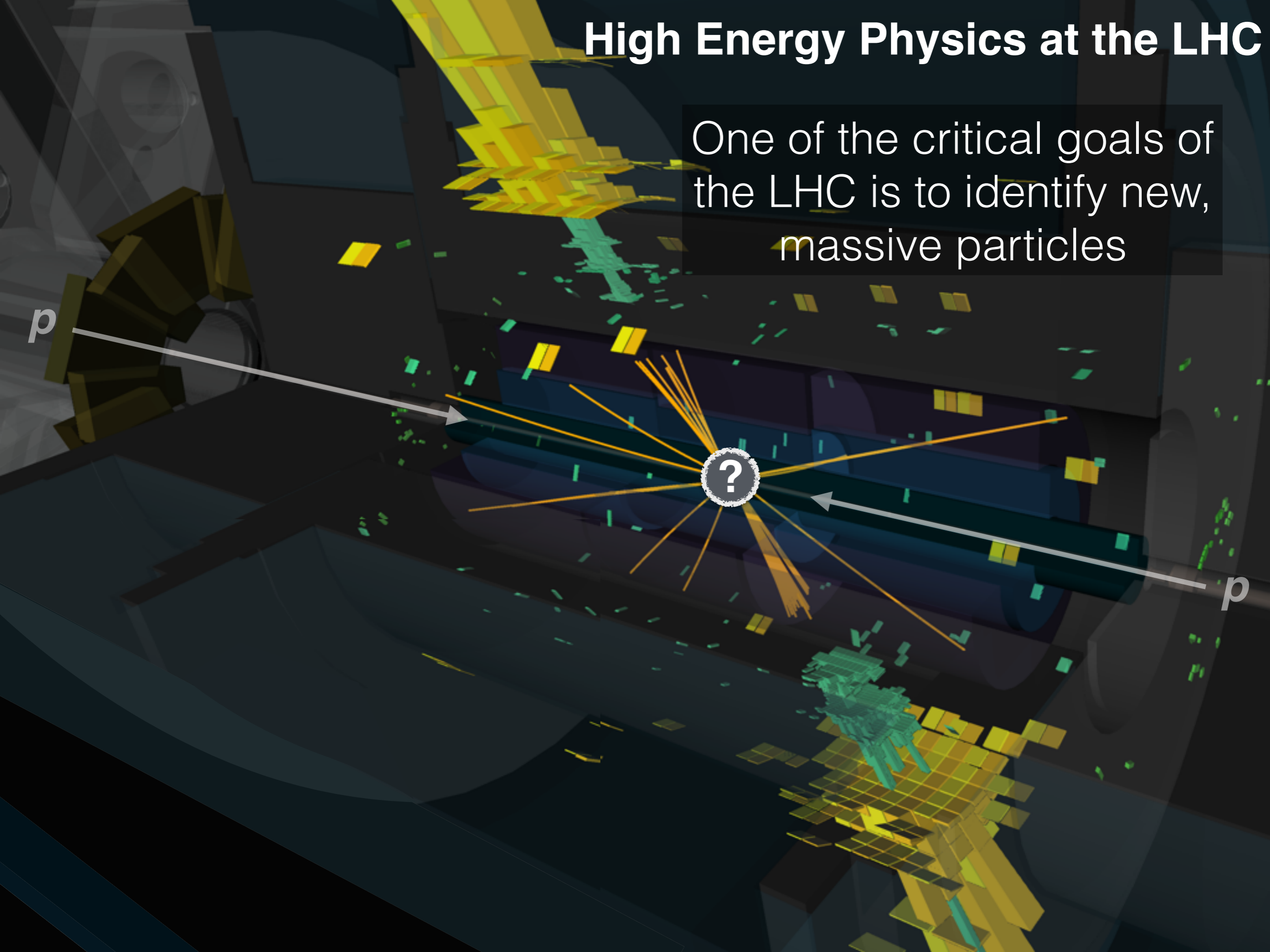
Event: 753275626

2016-06-18 18:41:48 CEST

Credit: All collision event displays from the **ATLAS** Collaboration

High Energy Physics at the LHC

One of the critical goals of the LHC is to identify new, massive particles



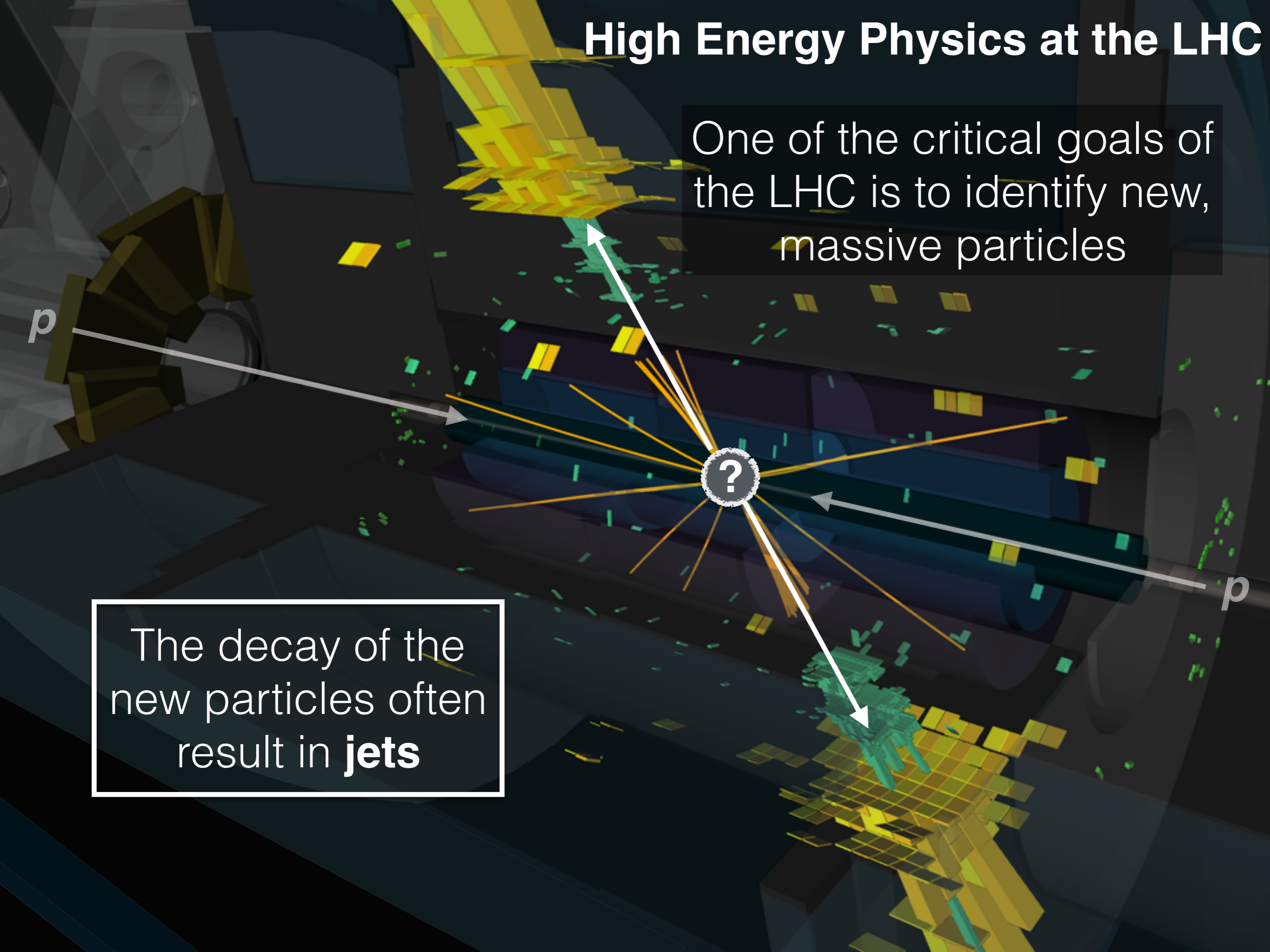
High Energy Physics at the LHC

One of the critical goals of the LHC is to identify new, massive particles

p

p

The decay of the new particles often result in **jets**



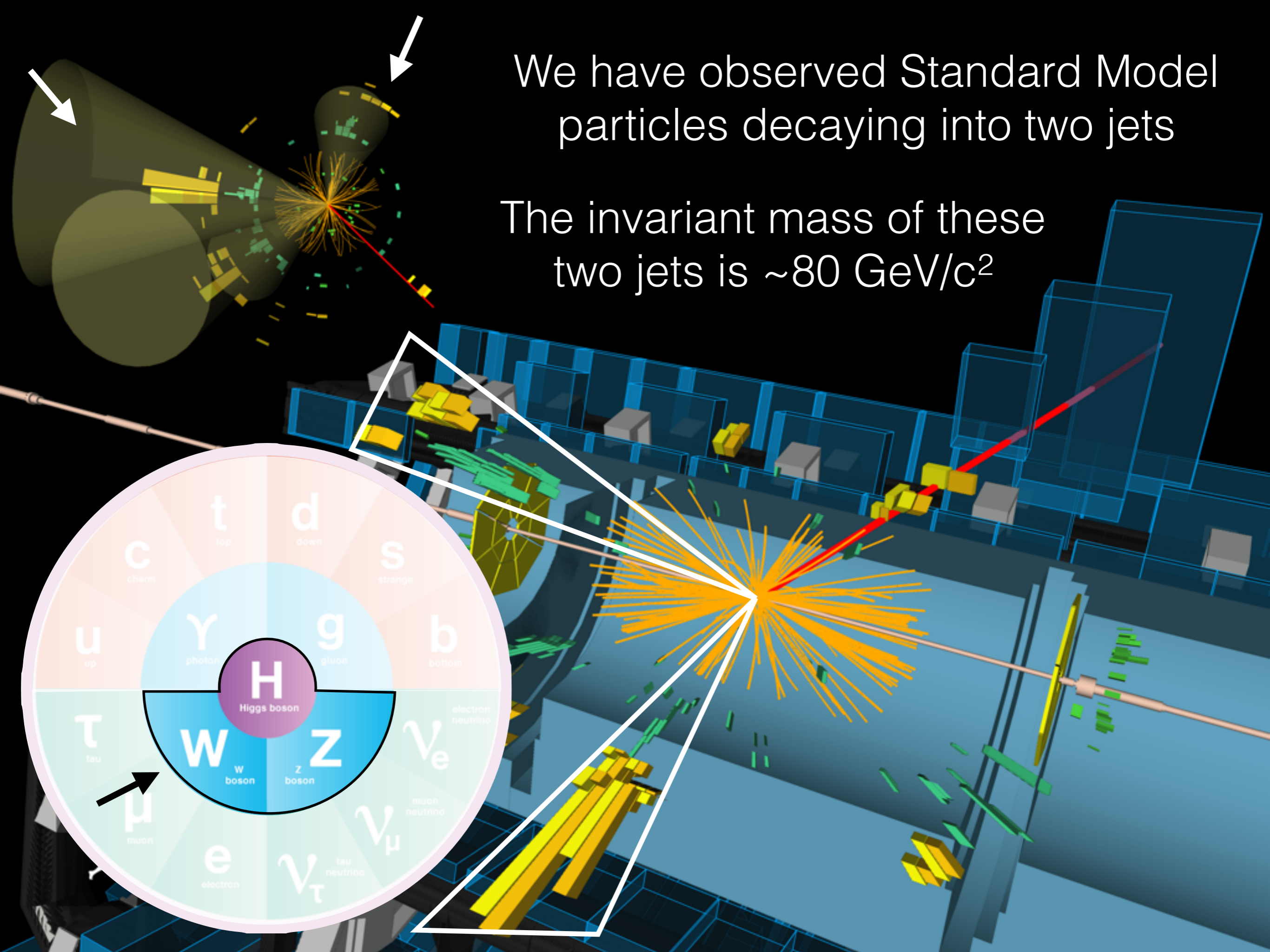
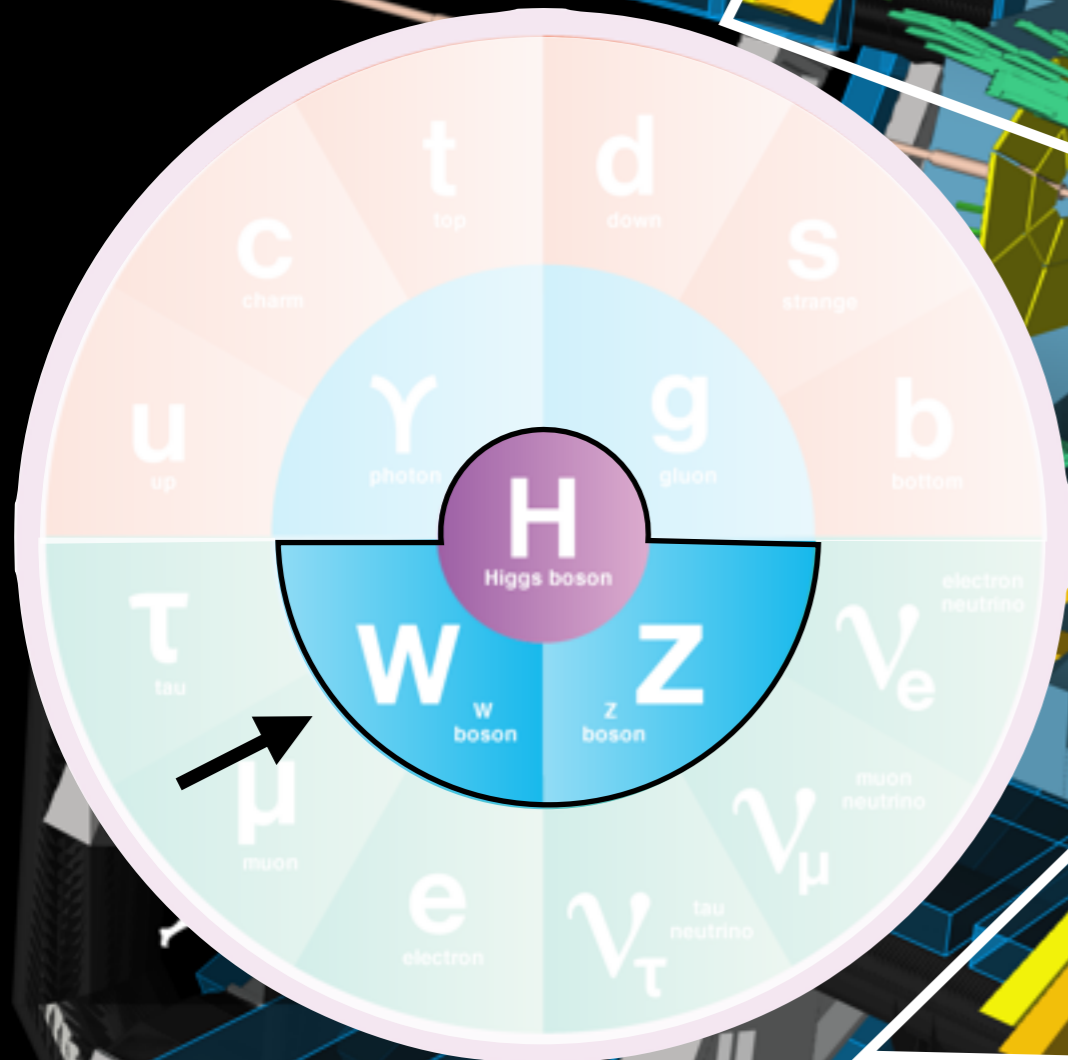


We have observed Standard Model particles decaying into two jets

The invariant mass of these two jets is $\sim 80 \text{ GeV}/c^2$

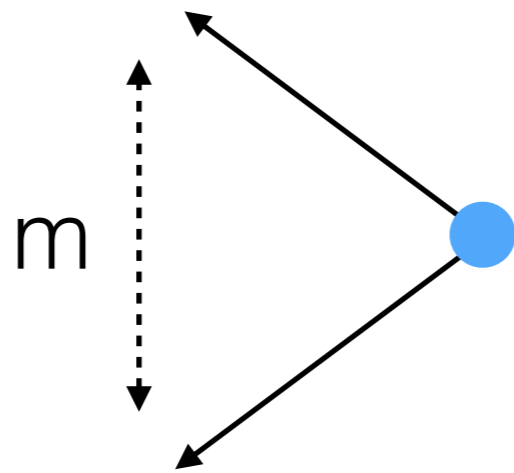
We have observed Standard Model particles decaying into two jets

The invariant mass of these two jets is $\sim 80 \text{ GeV}/c^2$

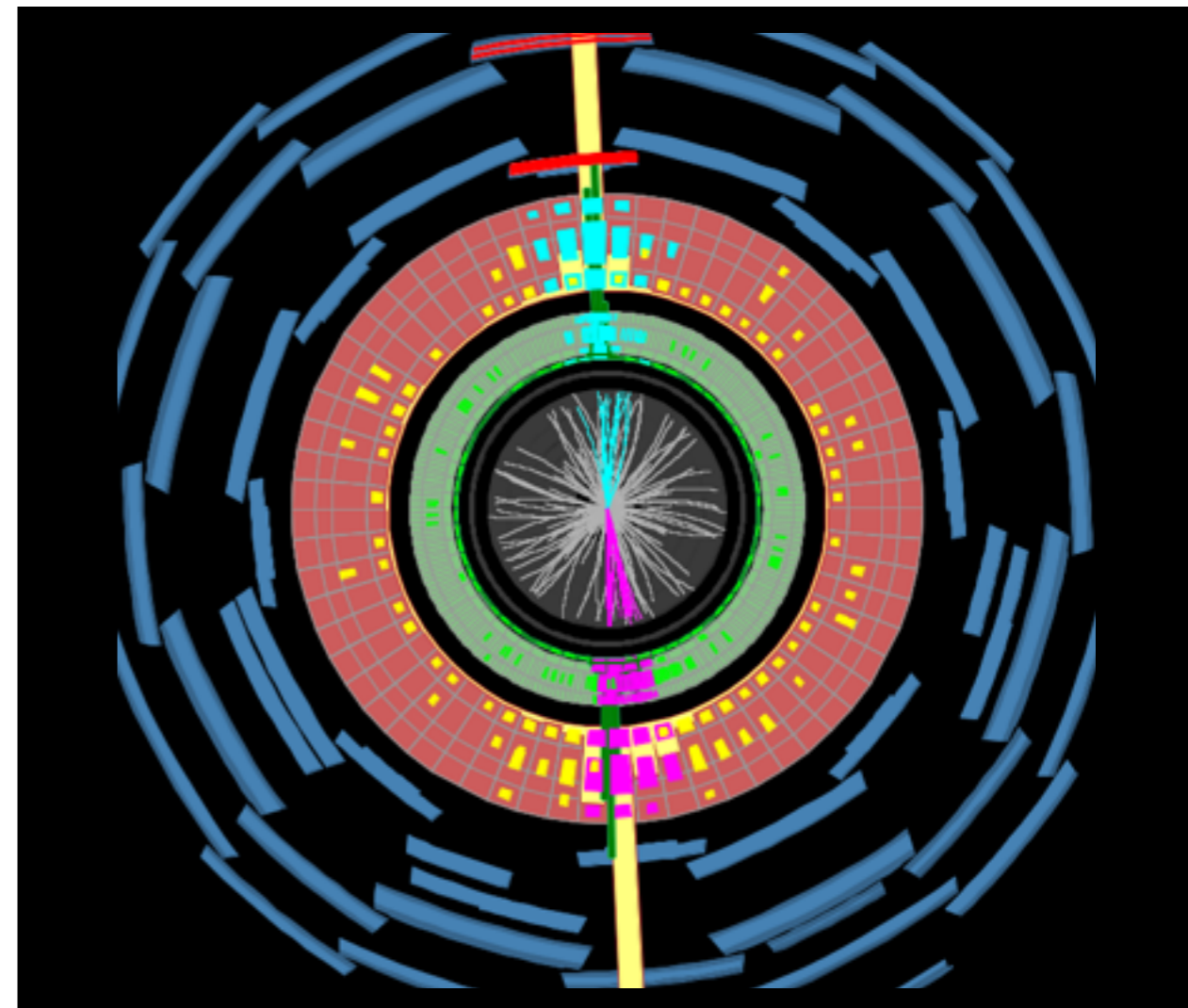
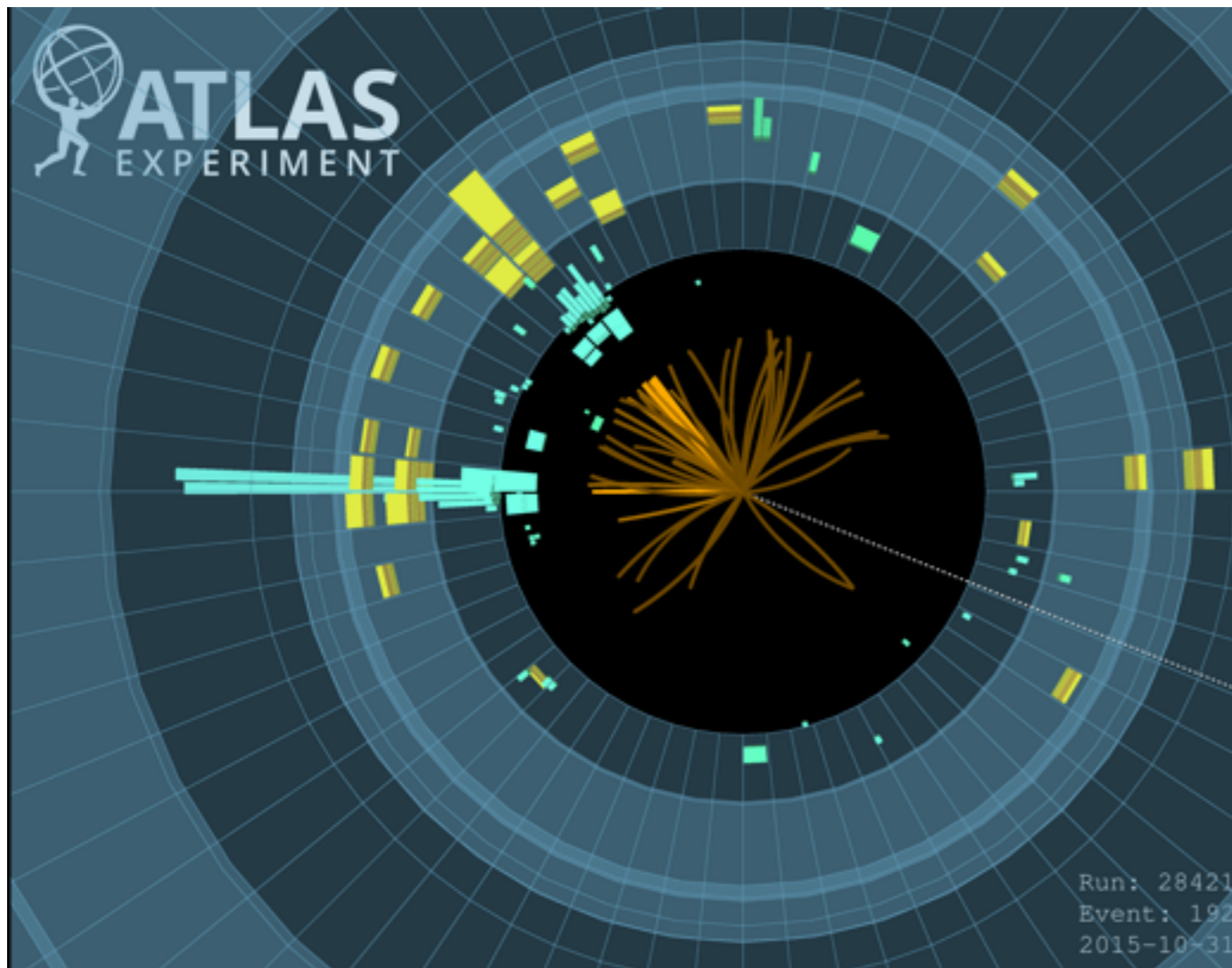
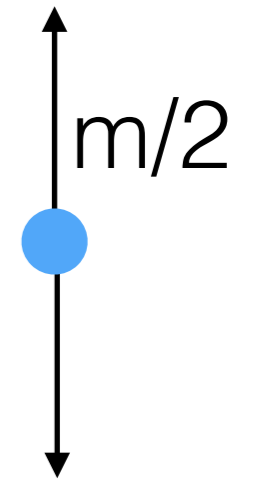


What if you take one of those SM dijet resonances and Lorentz boost it?

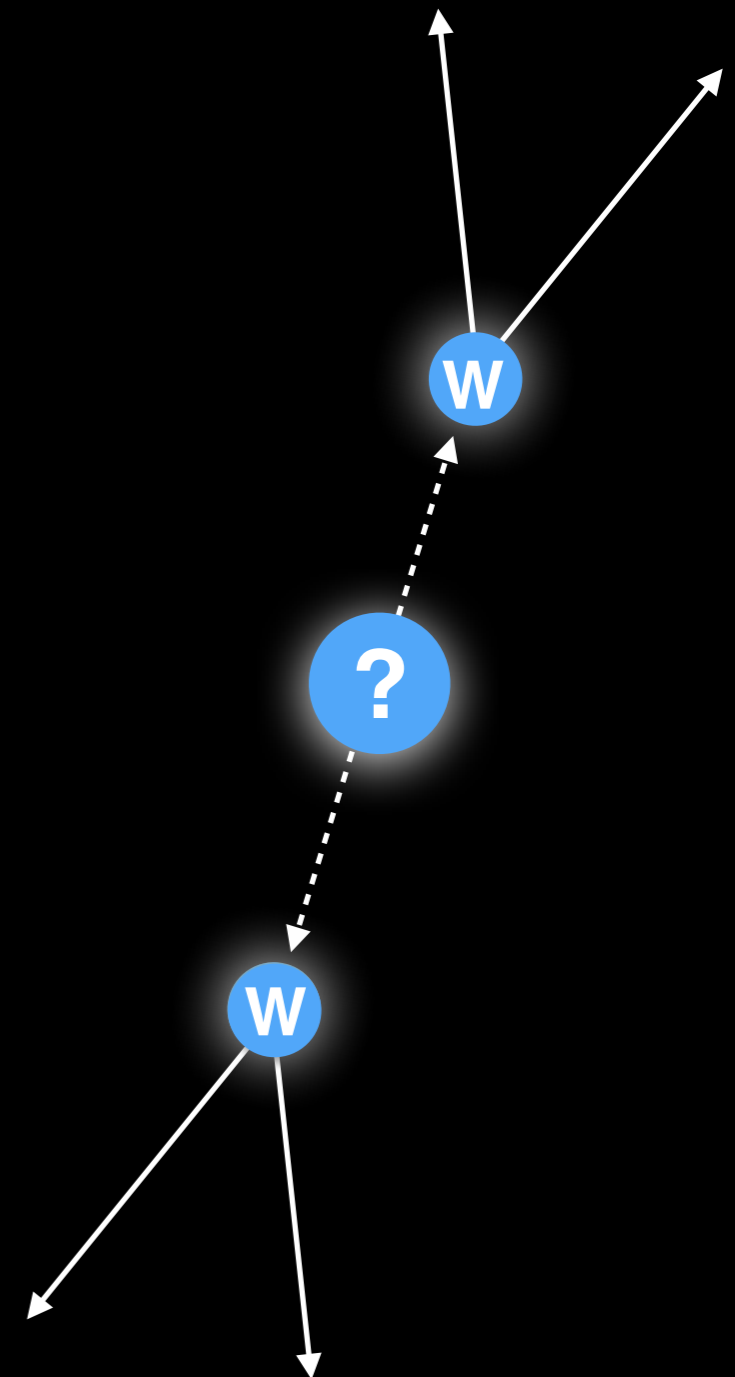
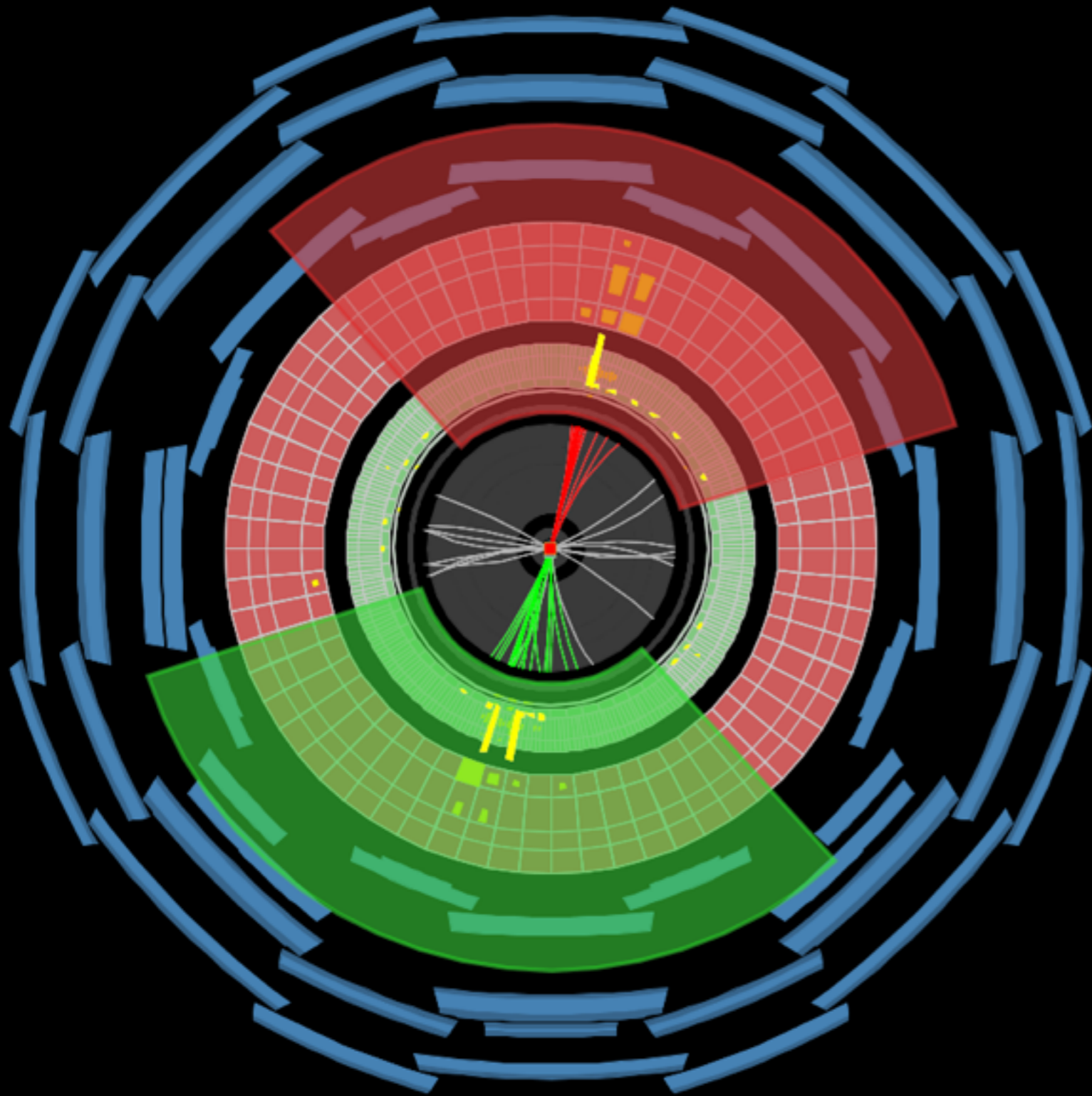
$$\phi \sim 1/\gamma = m/E$$



$$\gamma = E/m$$

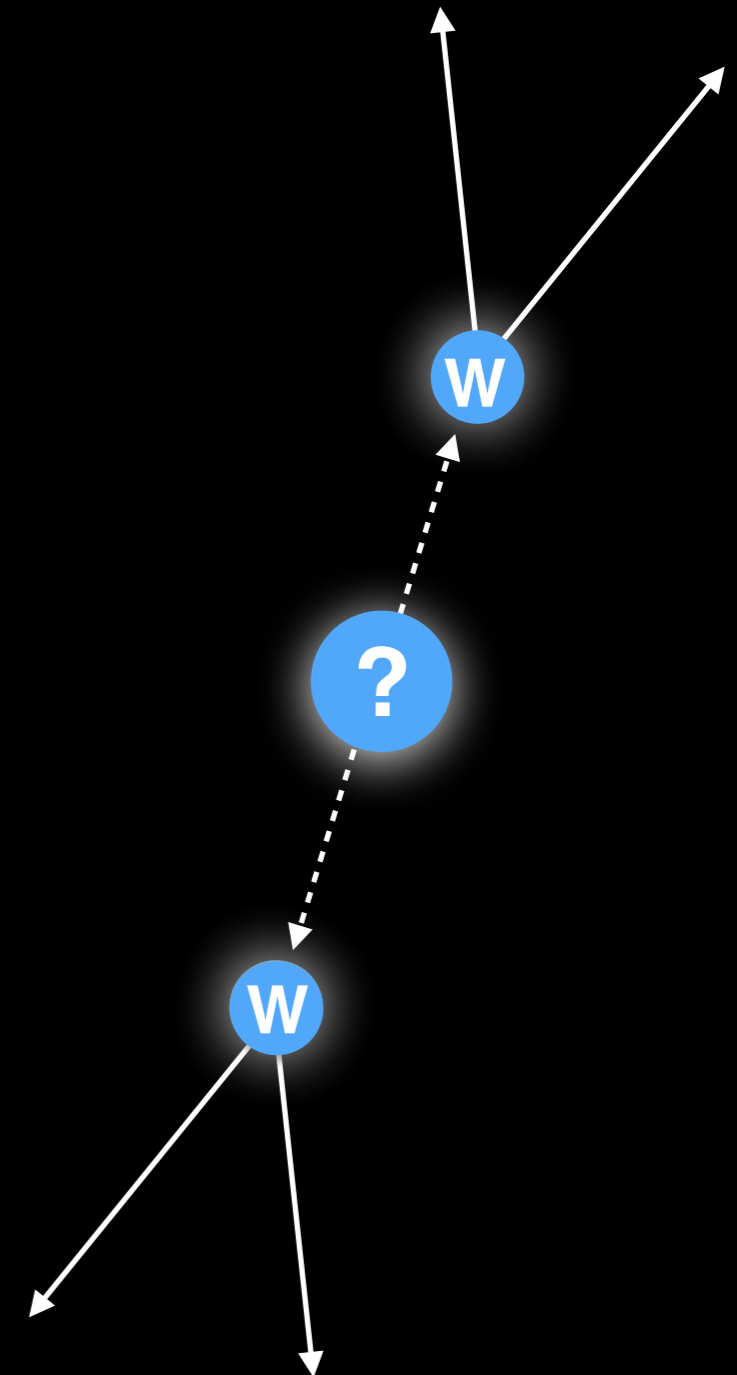


W bosons are naturally boosted if they result from the decay of something even heavier



W bosons are naturally boosted if they result from the decay of something even heavier

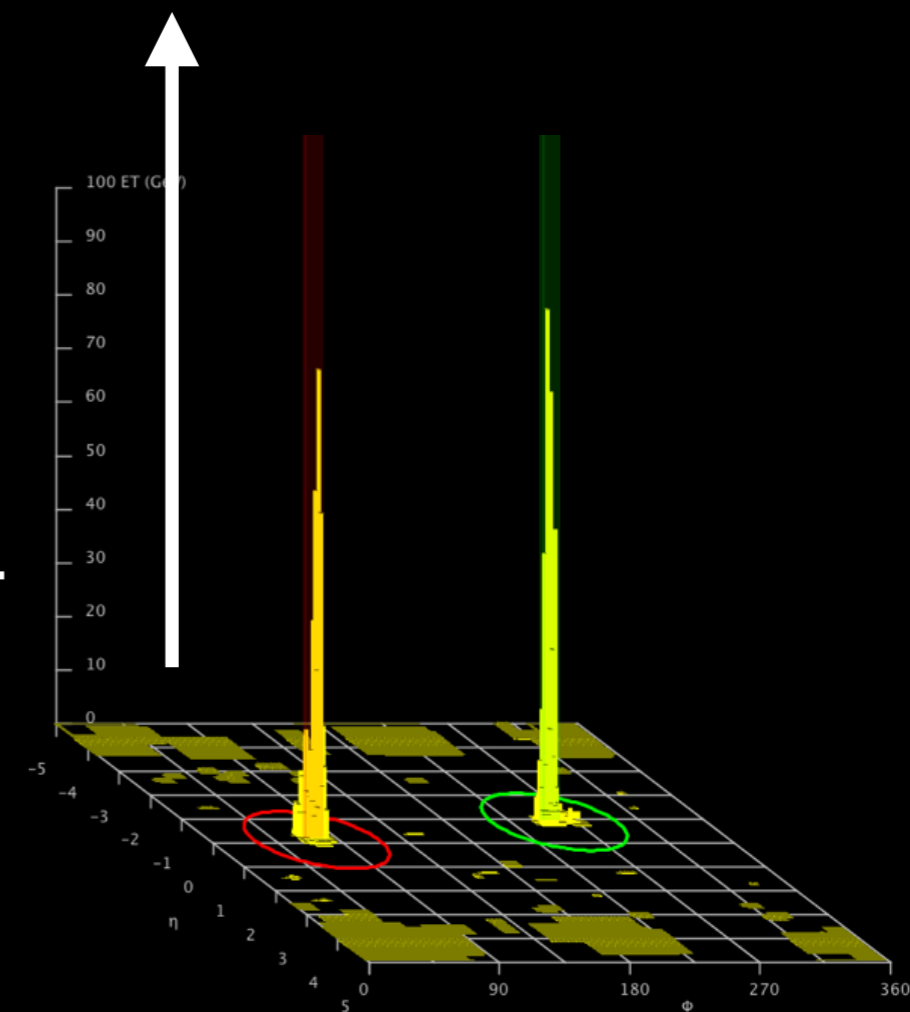
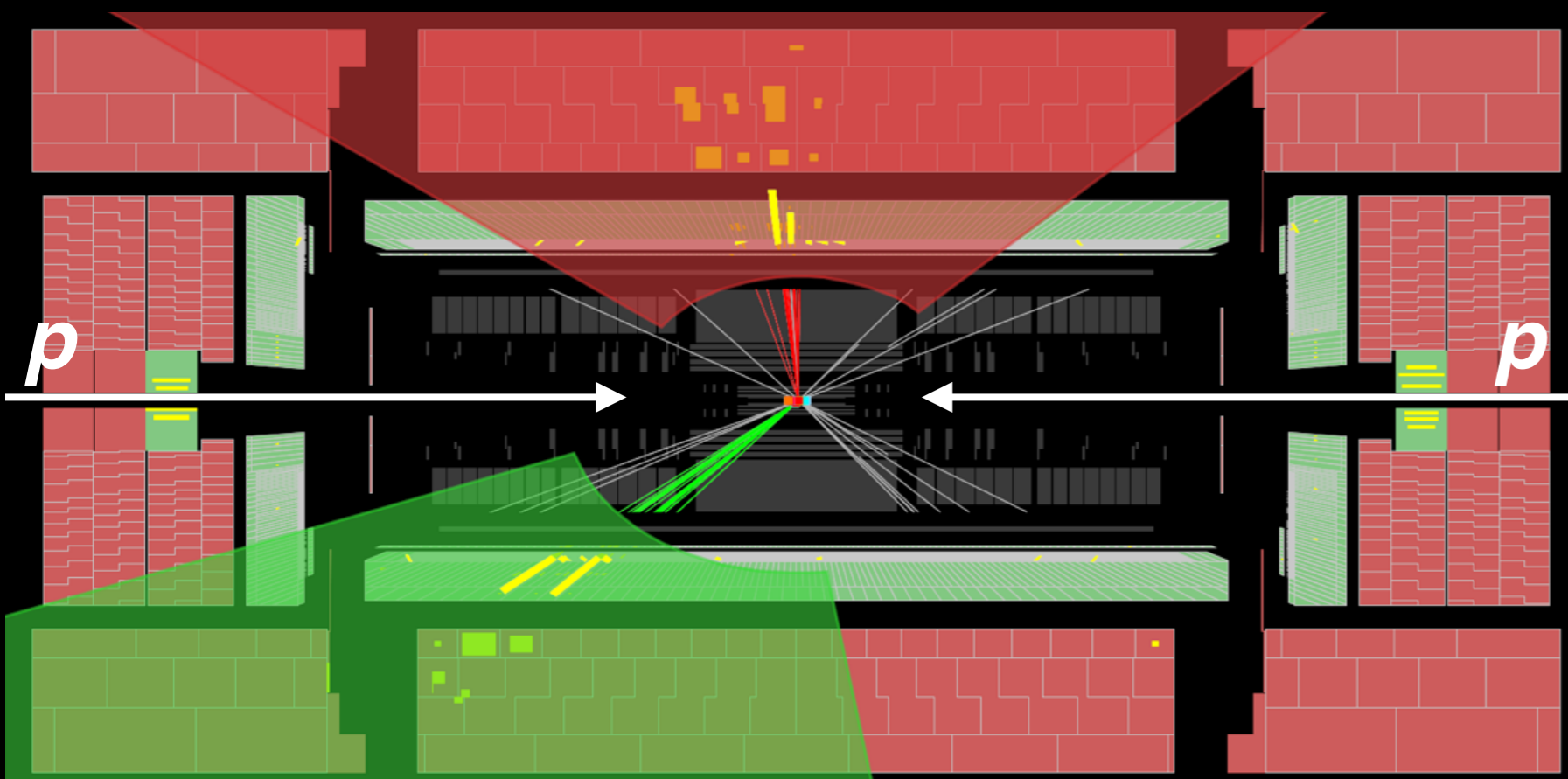
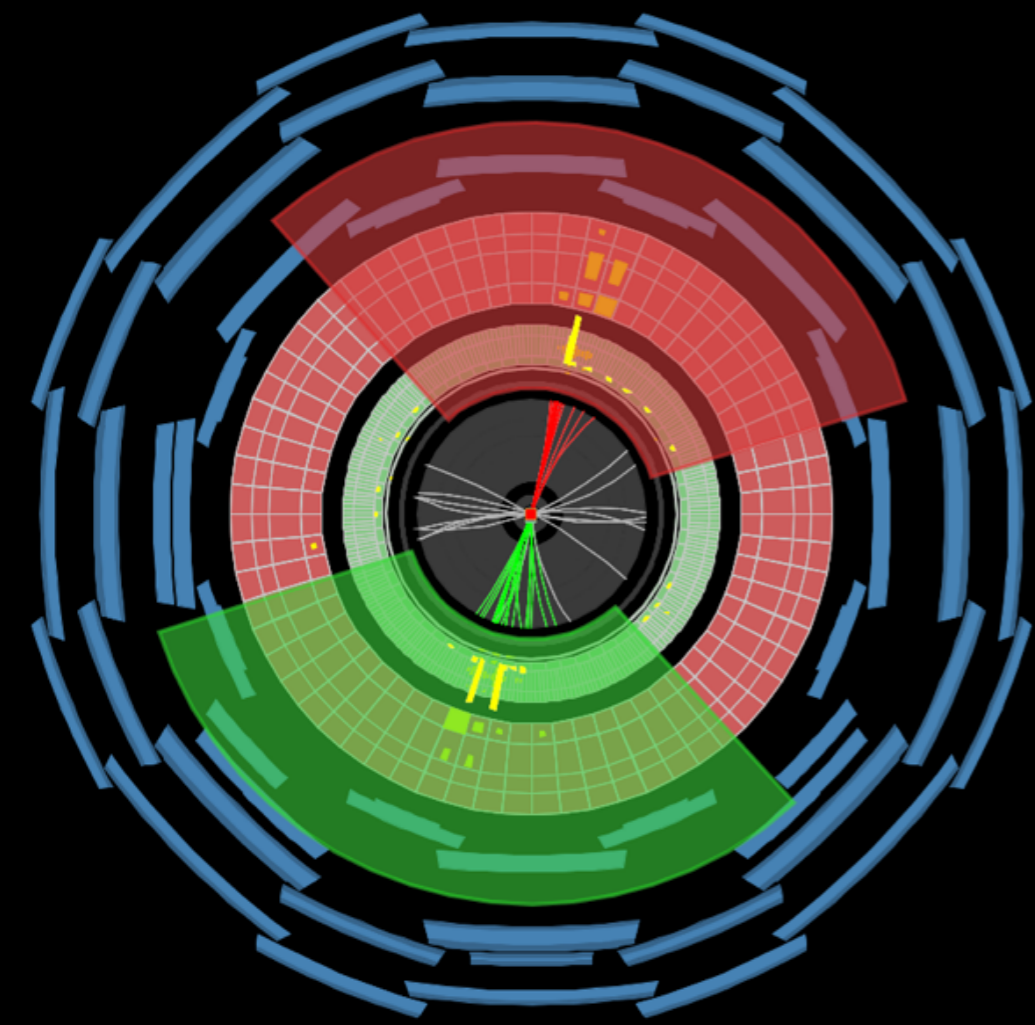
Goal: Find W jets in an enormous sea of generic q/g jets



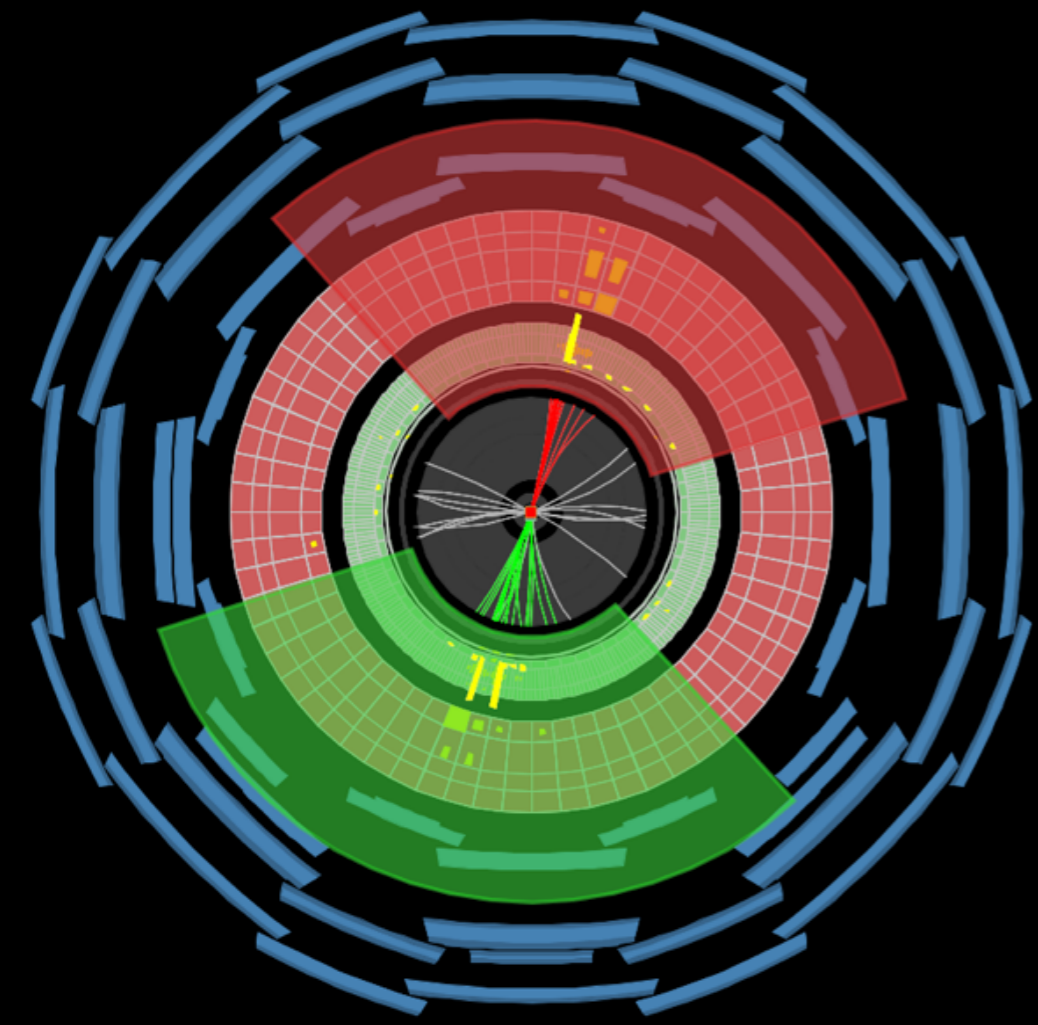
These jets have a non-trivial structure!

Searching for new particles decaying into boosted W bosons requires **looking at the radiation pattern inside jets**

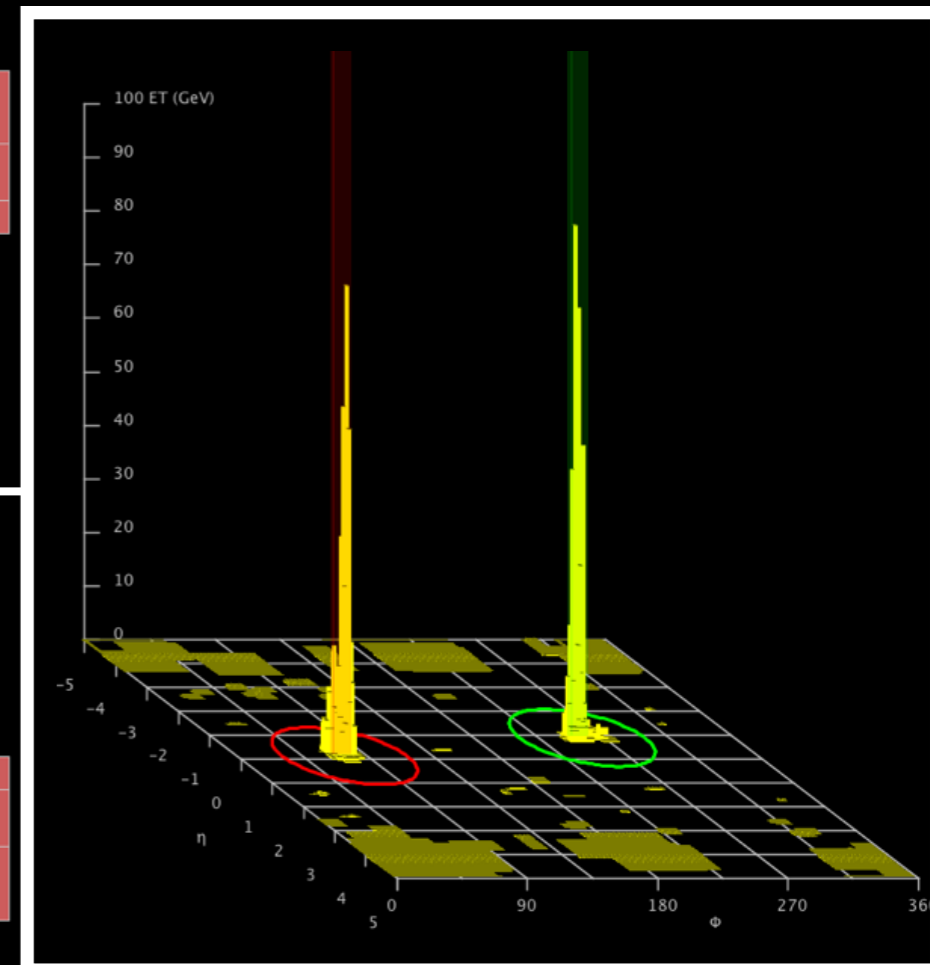
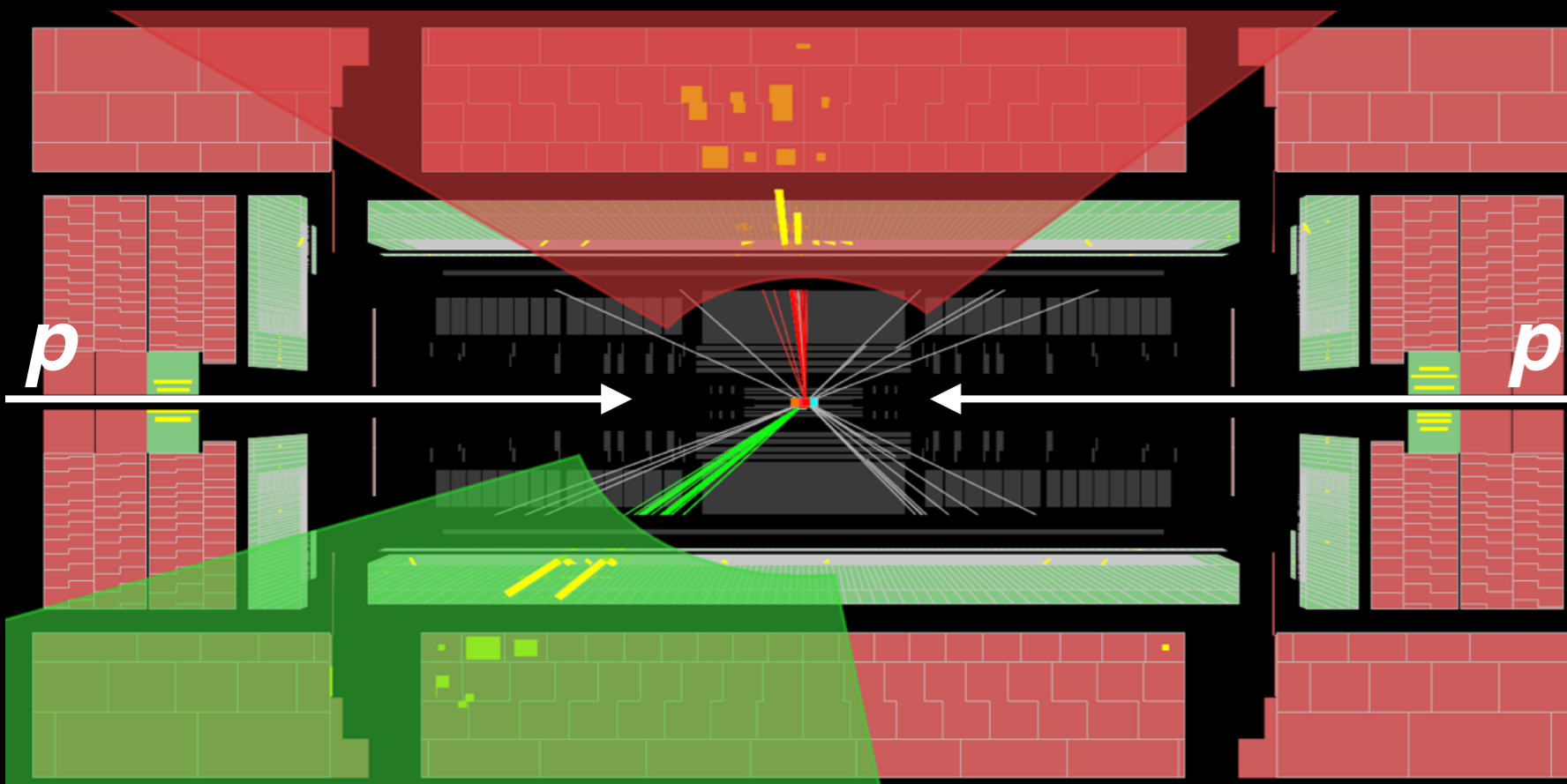
momentum transverse to the beam (p_T)



Up next: jet images



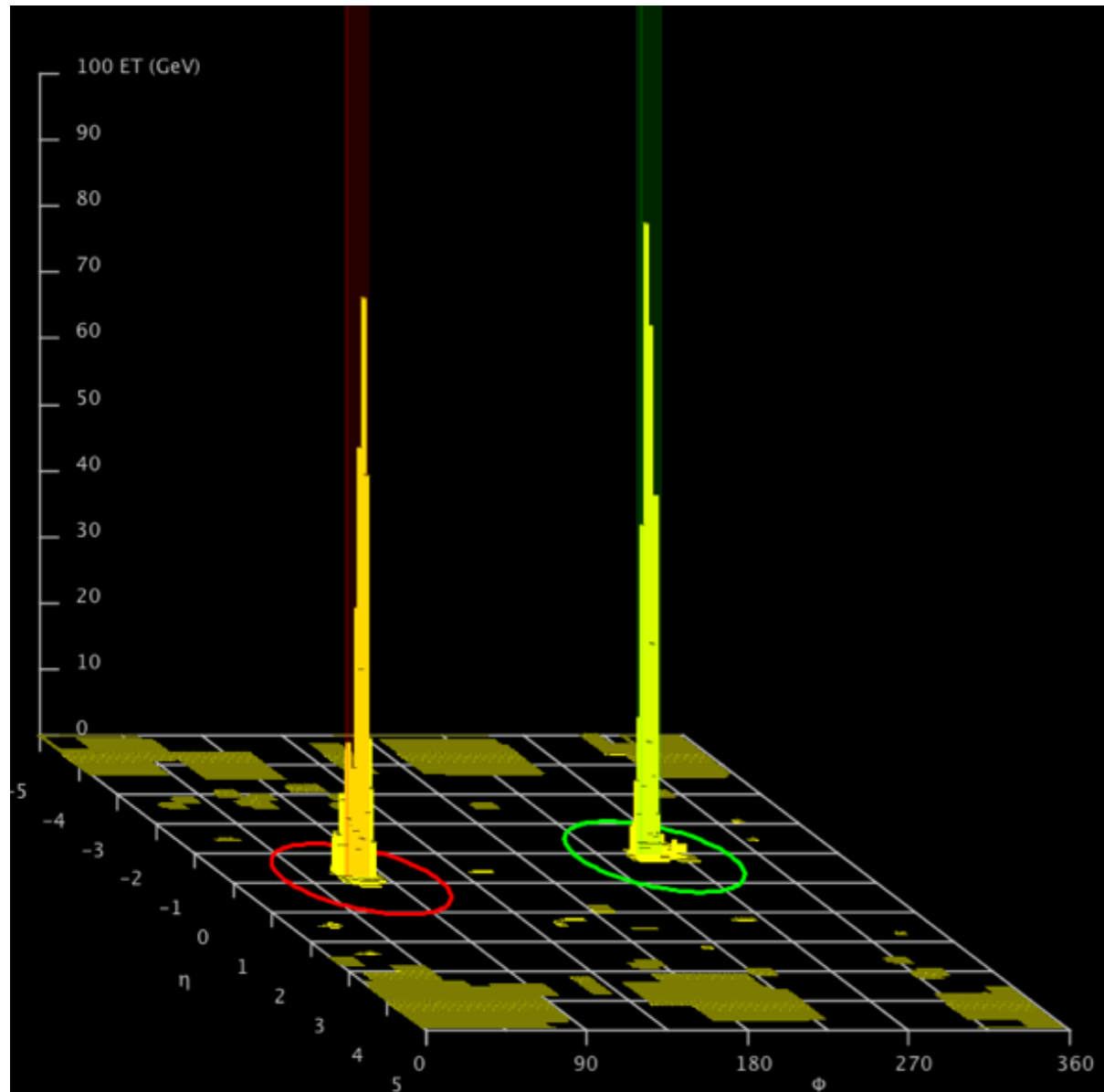
like a digital image!



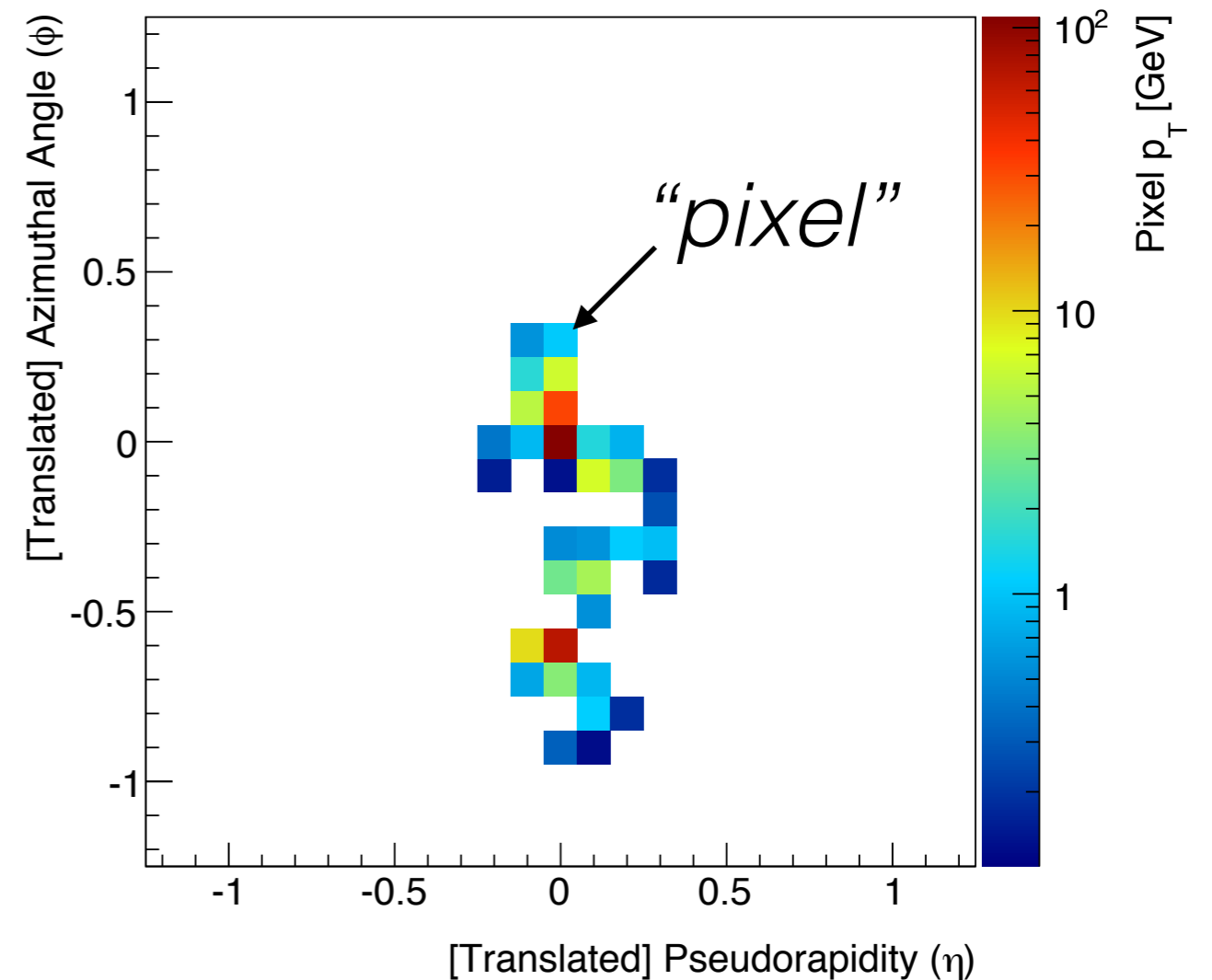
N.B. this is not the only way to represent a jet - more on that later

the Jet Image

J. Cogan et al. JHEP 02 (2015) 118



Boosted W

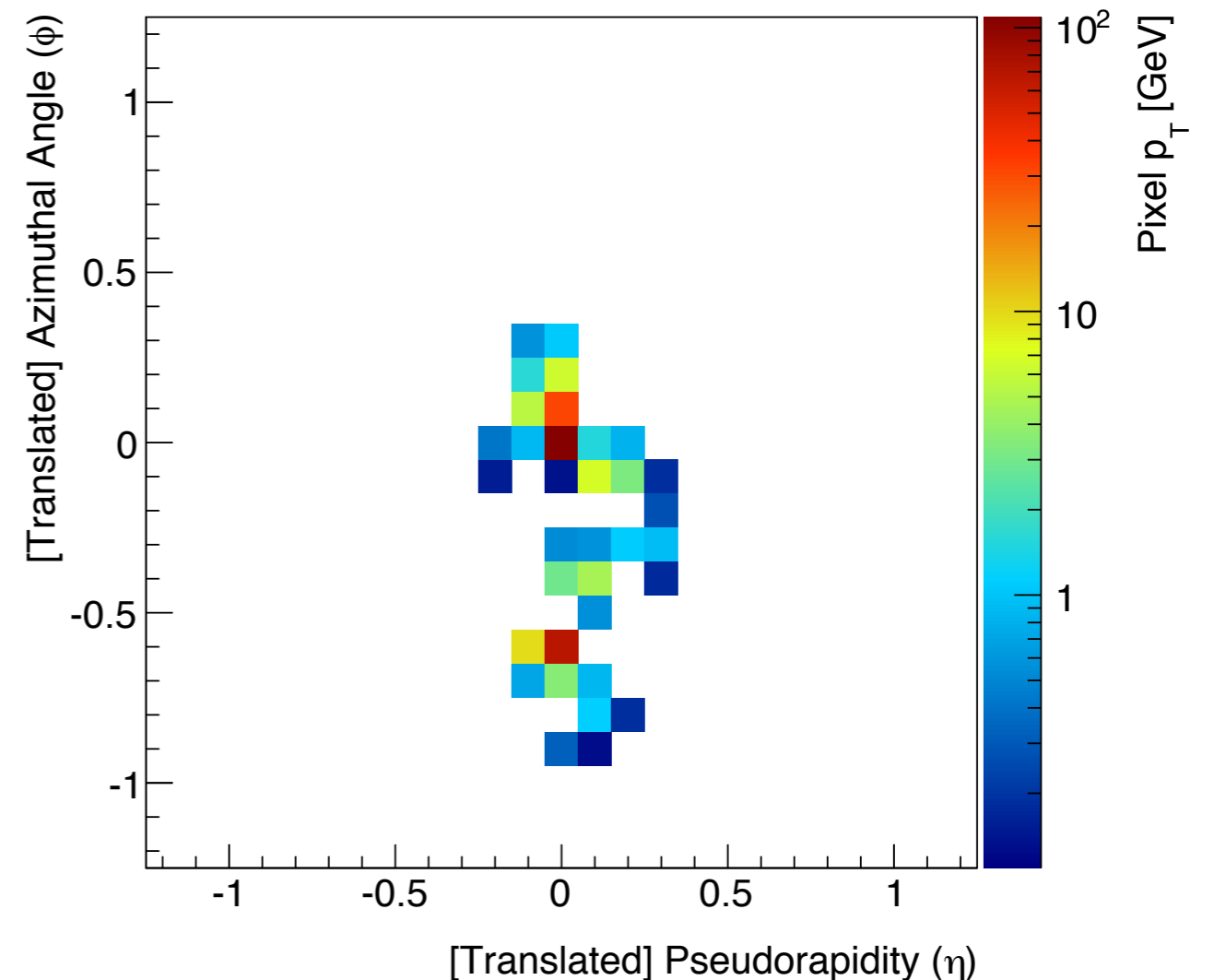


*nothing like a
'natural' image!*

the Jet Image

J. Cogan et al. JHEP 02 (2015) 118

Boosted W



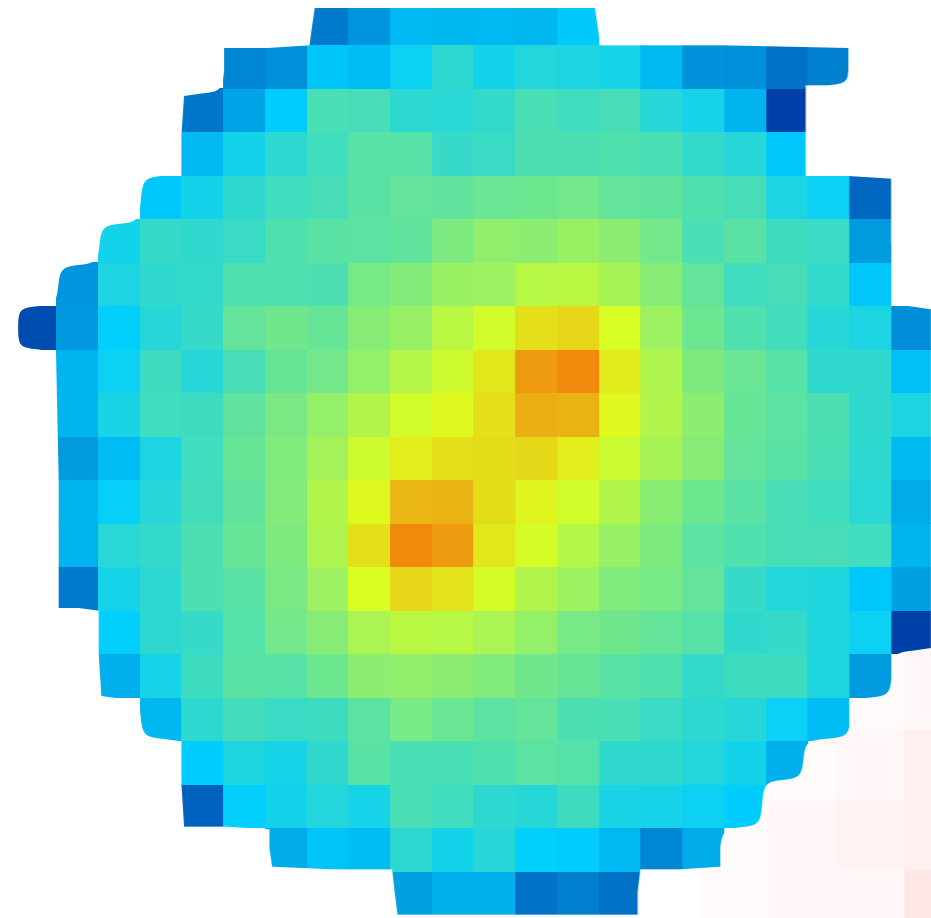
Credit: Peter G Trimming (Wikipedia)

*no smooth edges, clear features, low
occupancy (number of hit pixels)*

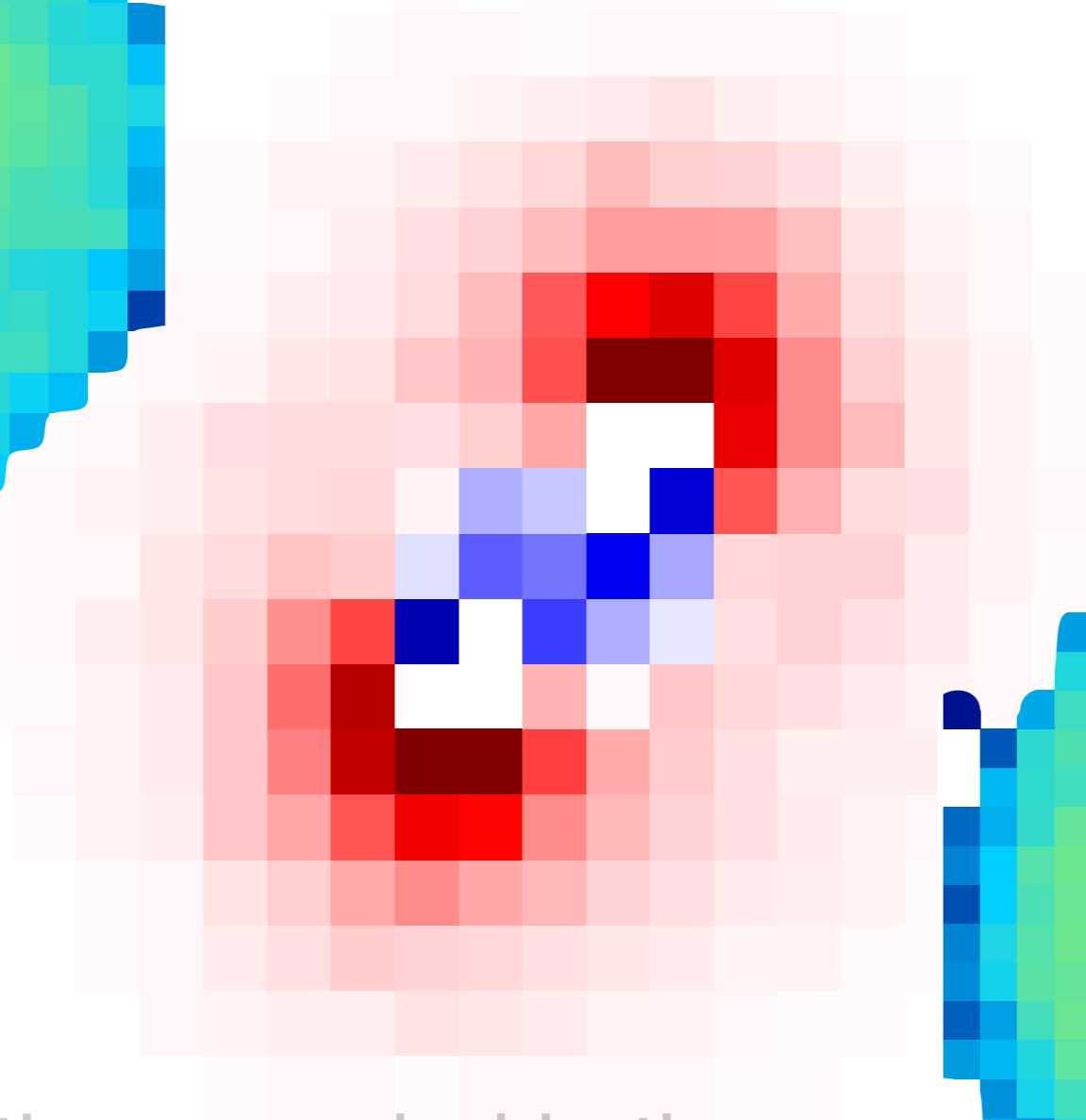
Why images?

Can directly visualize physics

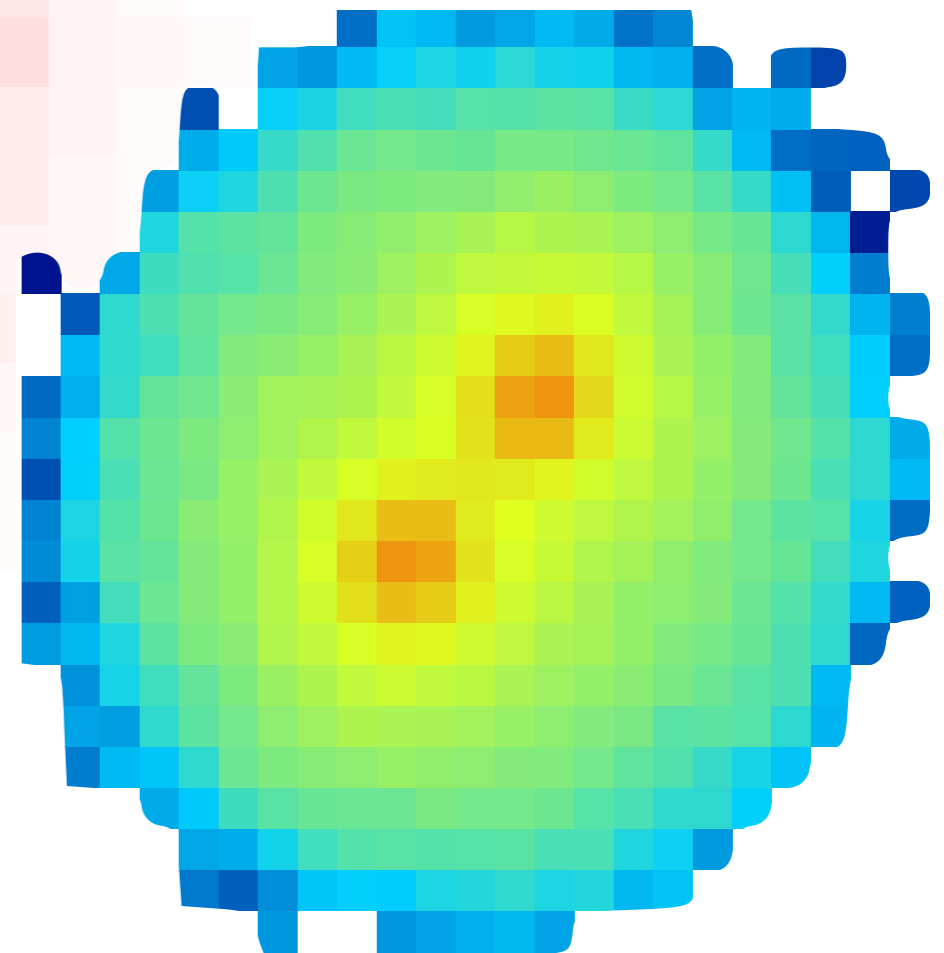
and we can benefit from the extensive image processing literature



$W \rightarrow q\bar{q}$



$g \rightarrow q\bar{q}$



there is information encoded in the physical distance between pixels

Why images?

Can directly visualize physics

and we can benefit from the
extensive image processing literature

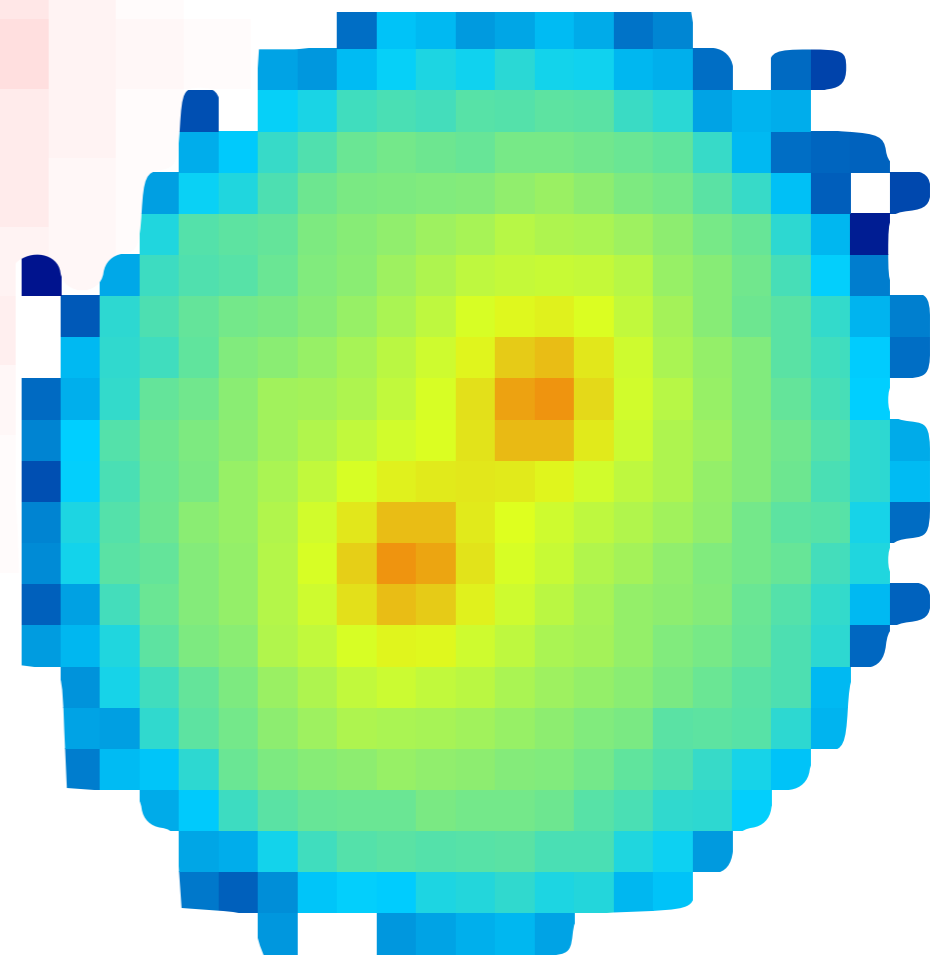
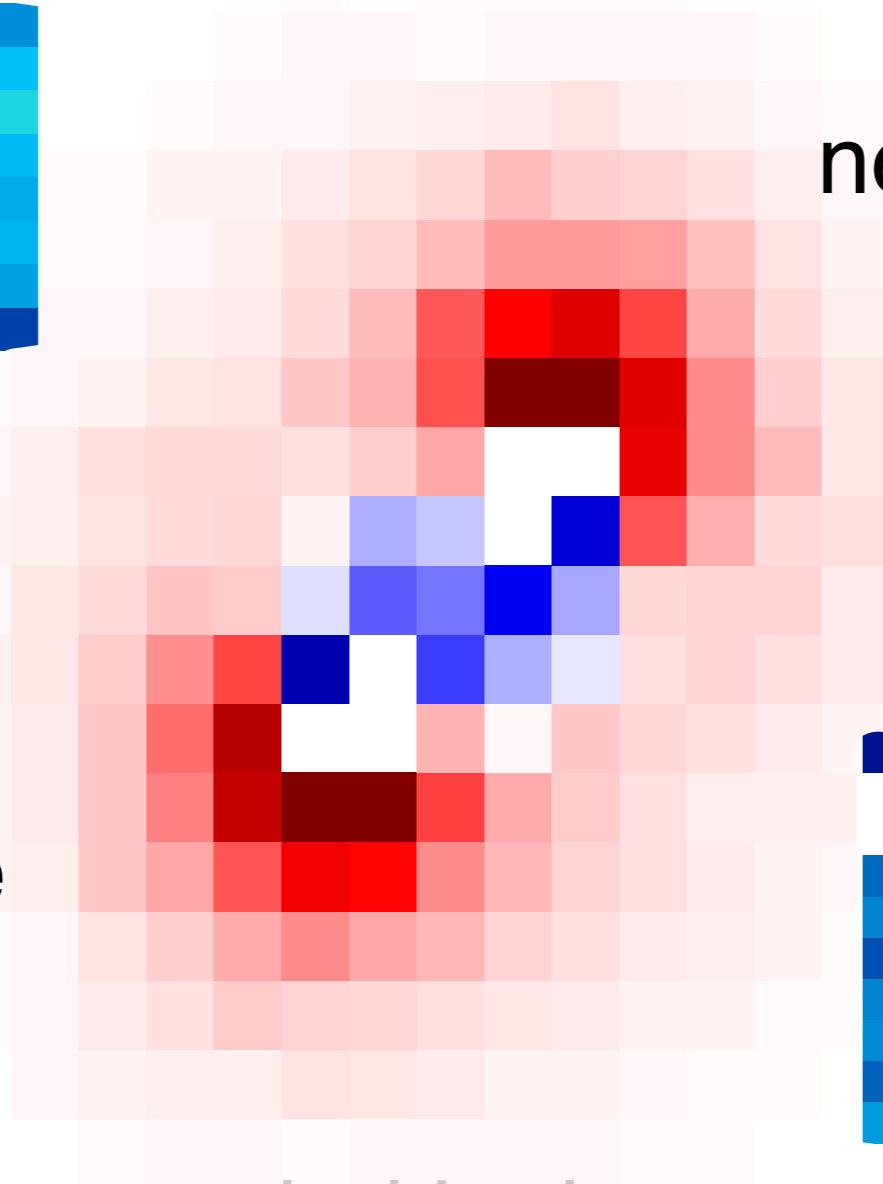
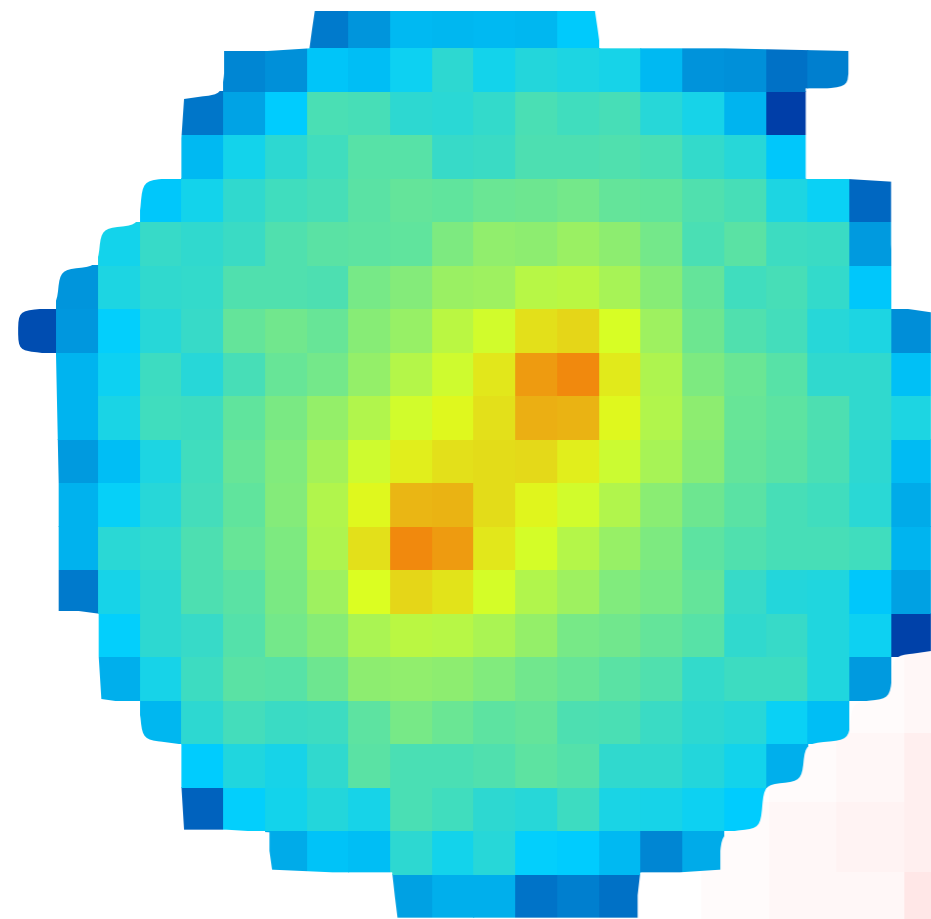
net strong-force charge

$g \dashrightarrow q\bar{q}$

$W \dashrightarrow q\bar{q}$

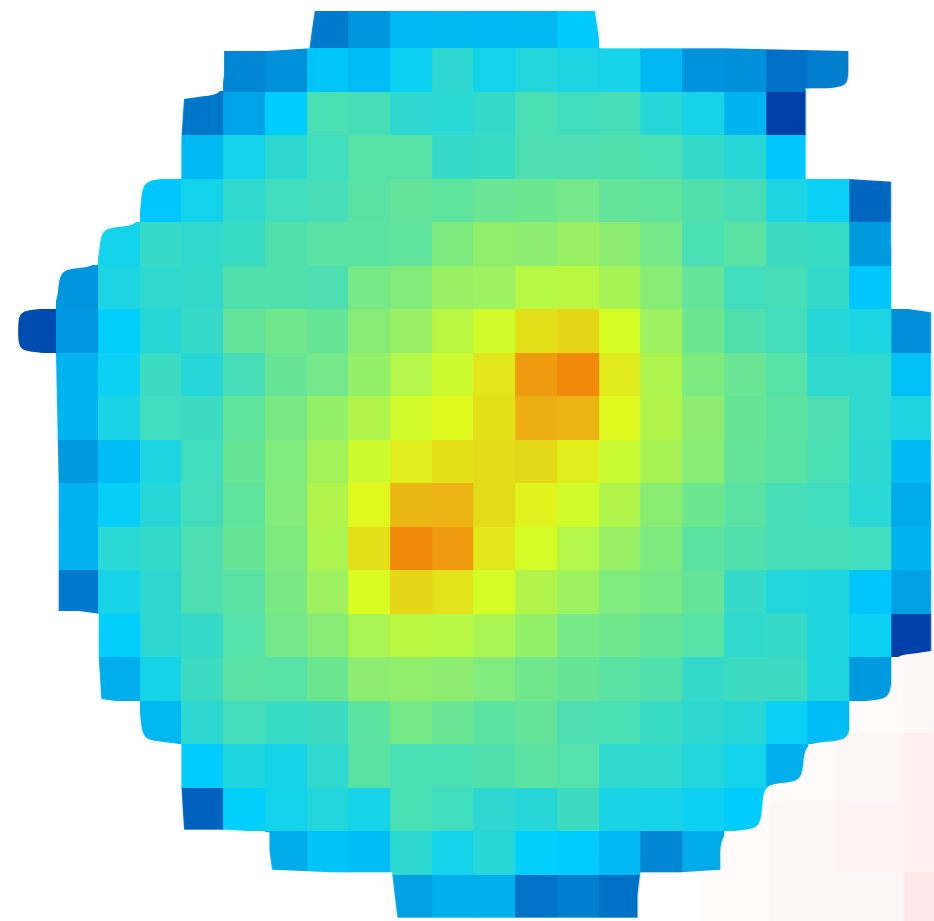
radiates like a dipole
(no net charge)

there is information encoded in the
physical distance between pixels

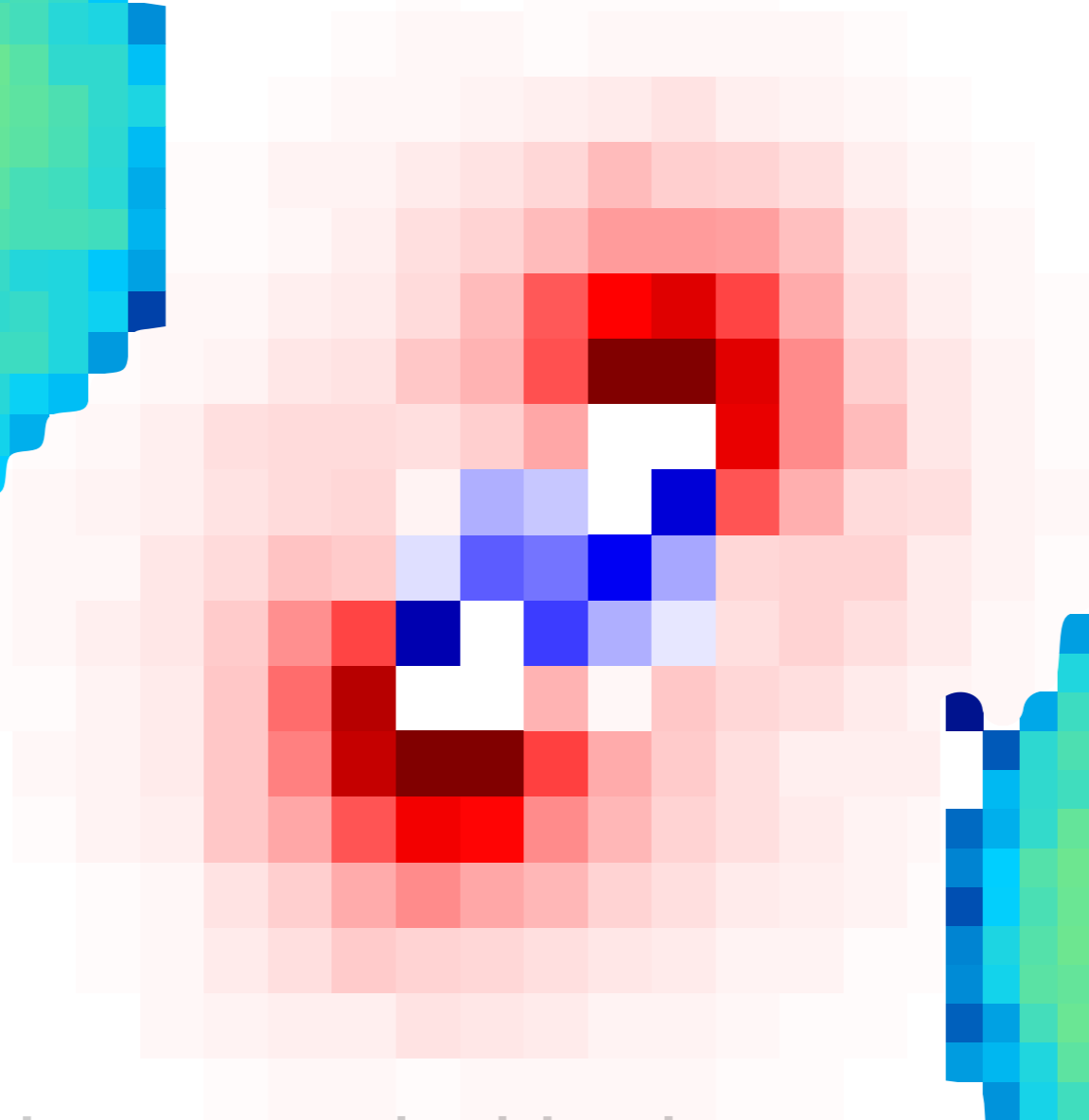


Why images?

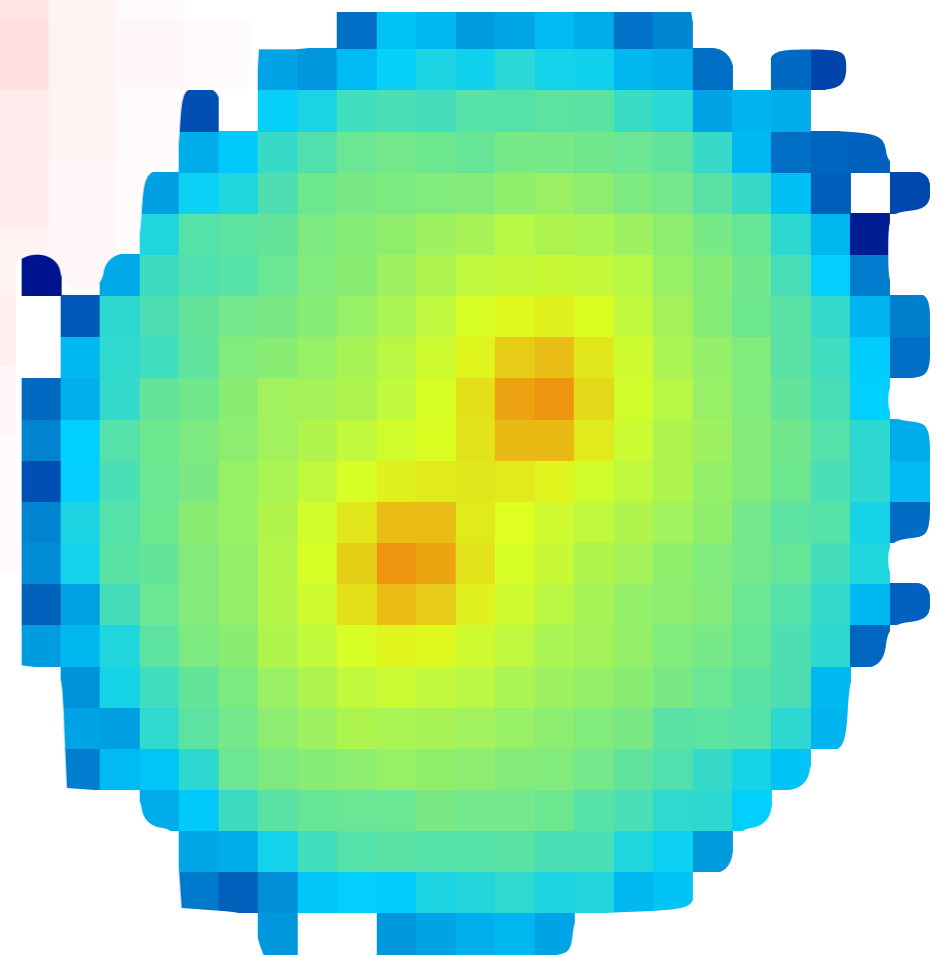
Can directly visualize physics
and we can benefit from the
extensive image processing literature



$W \rightarrow q\bar{q}$



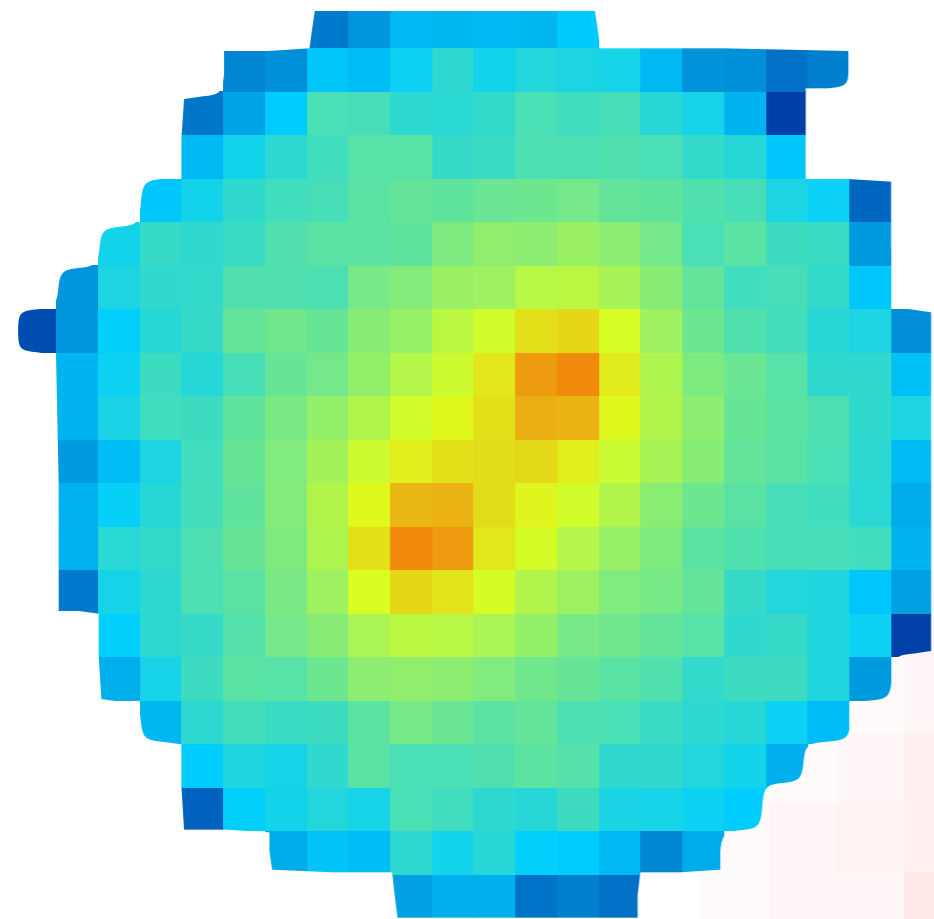
$g \rightarrow q\bar{q}$



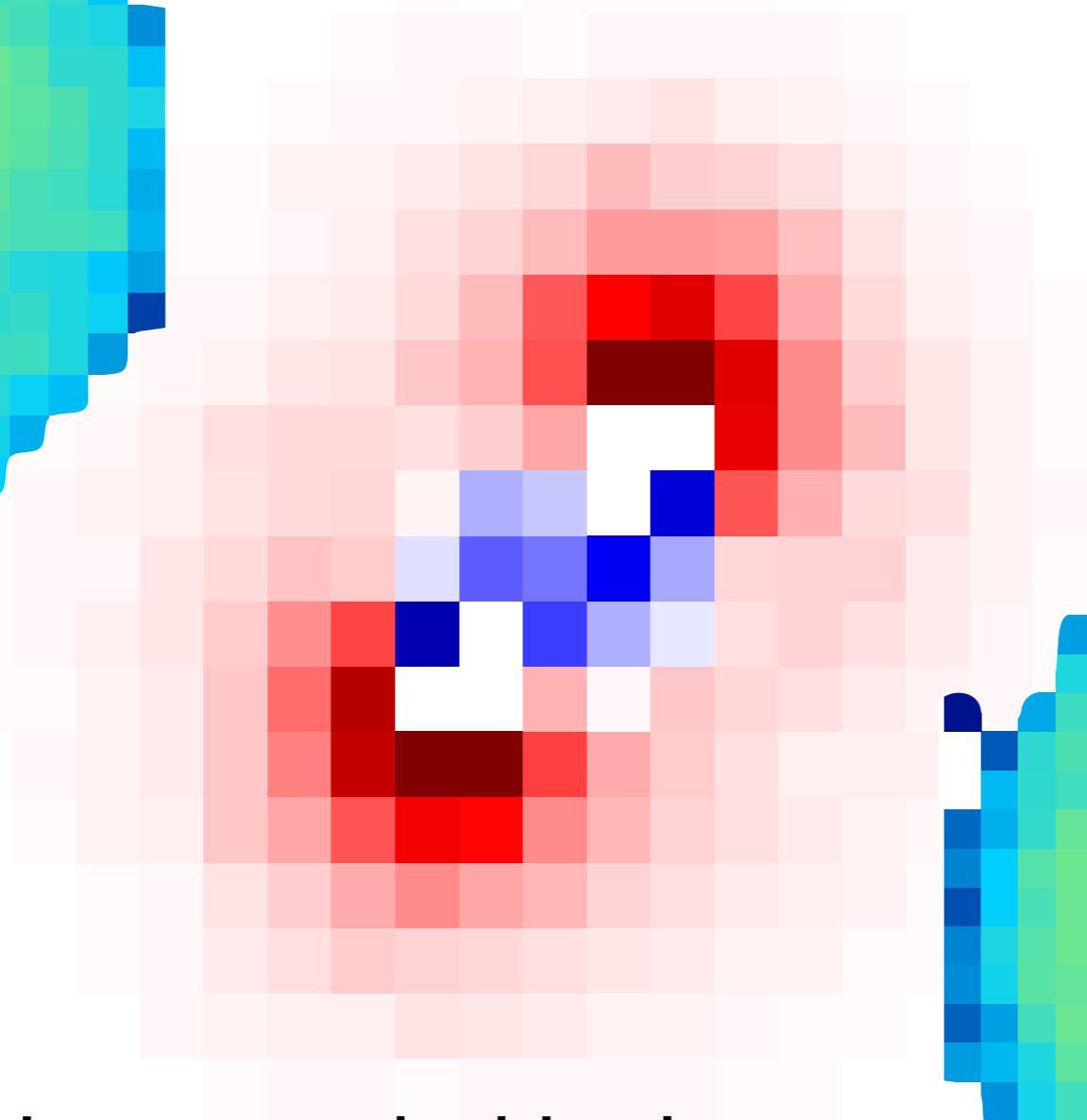
there is information encoded in the
physical distance between pixels

Why images?

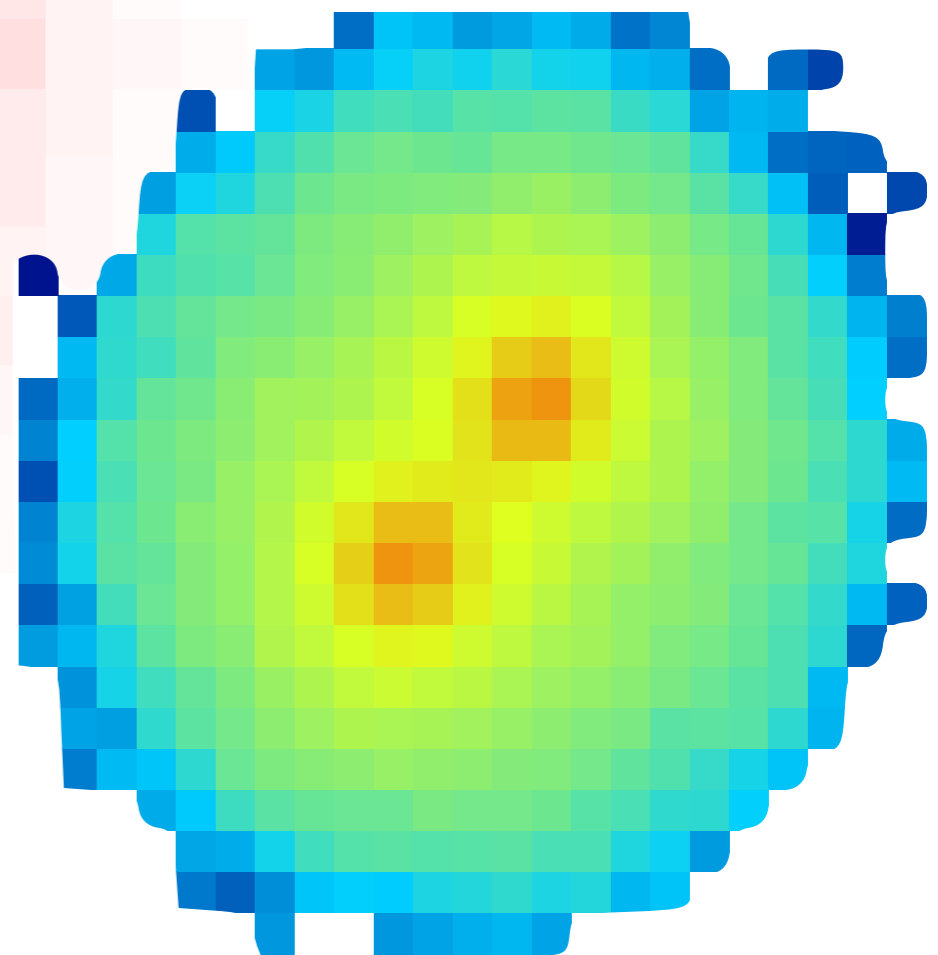
Can directly visualize physics
and we can benefit from the
extensive image processing literature



$W \rightarrow q\bar{q}$



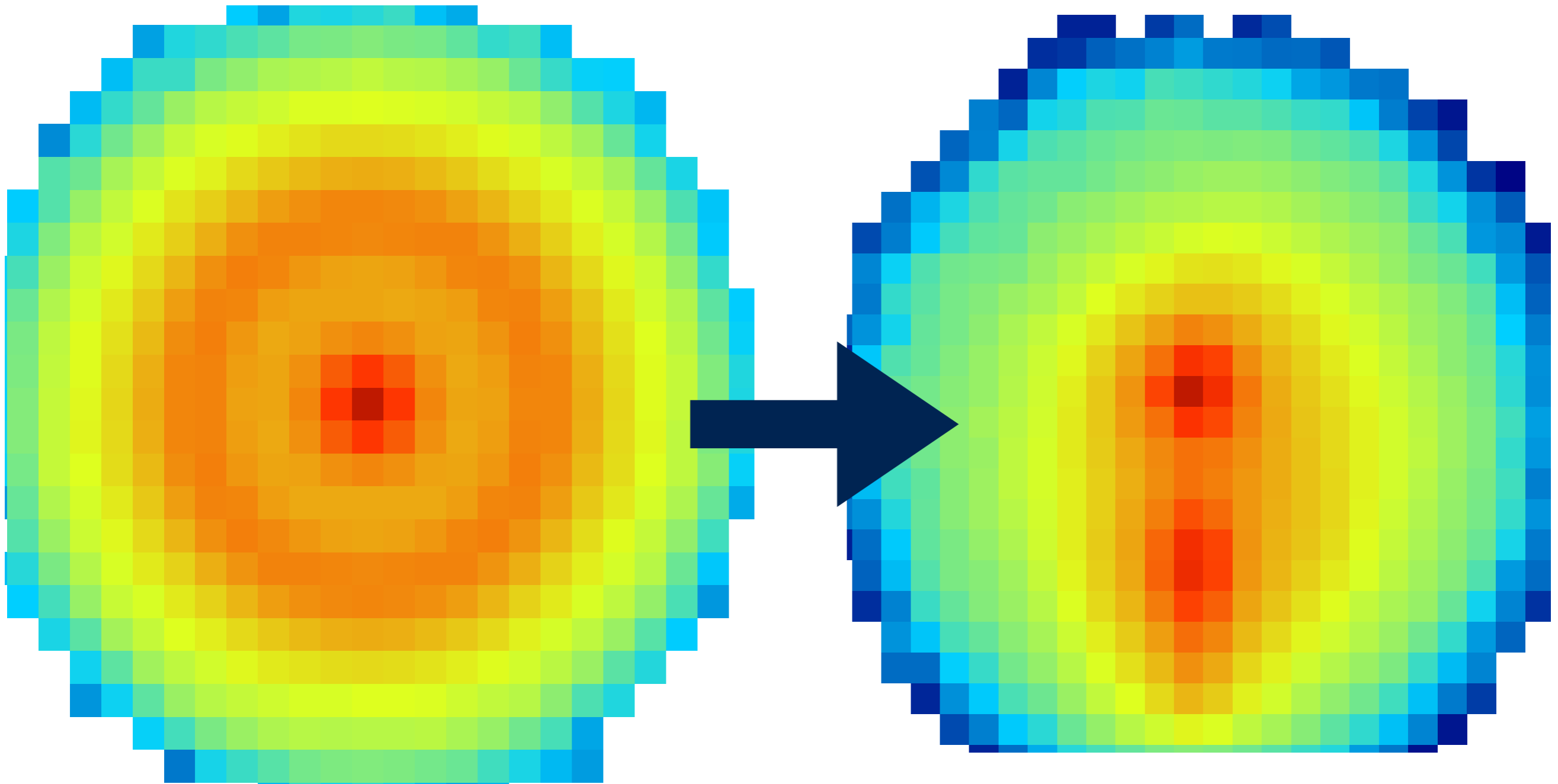
$g \rightarrow q\bar{q}$



there is information encoded in the
physical distance between pixels

Pre-processing & spacetime symmetries

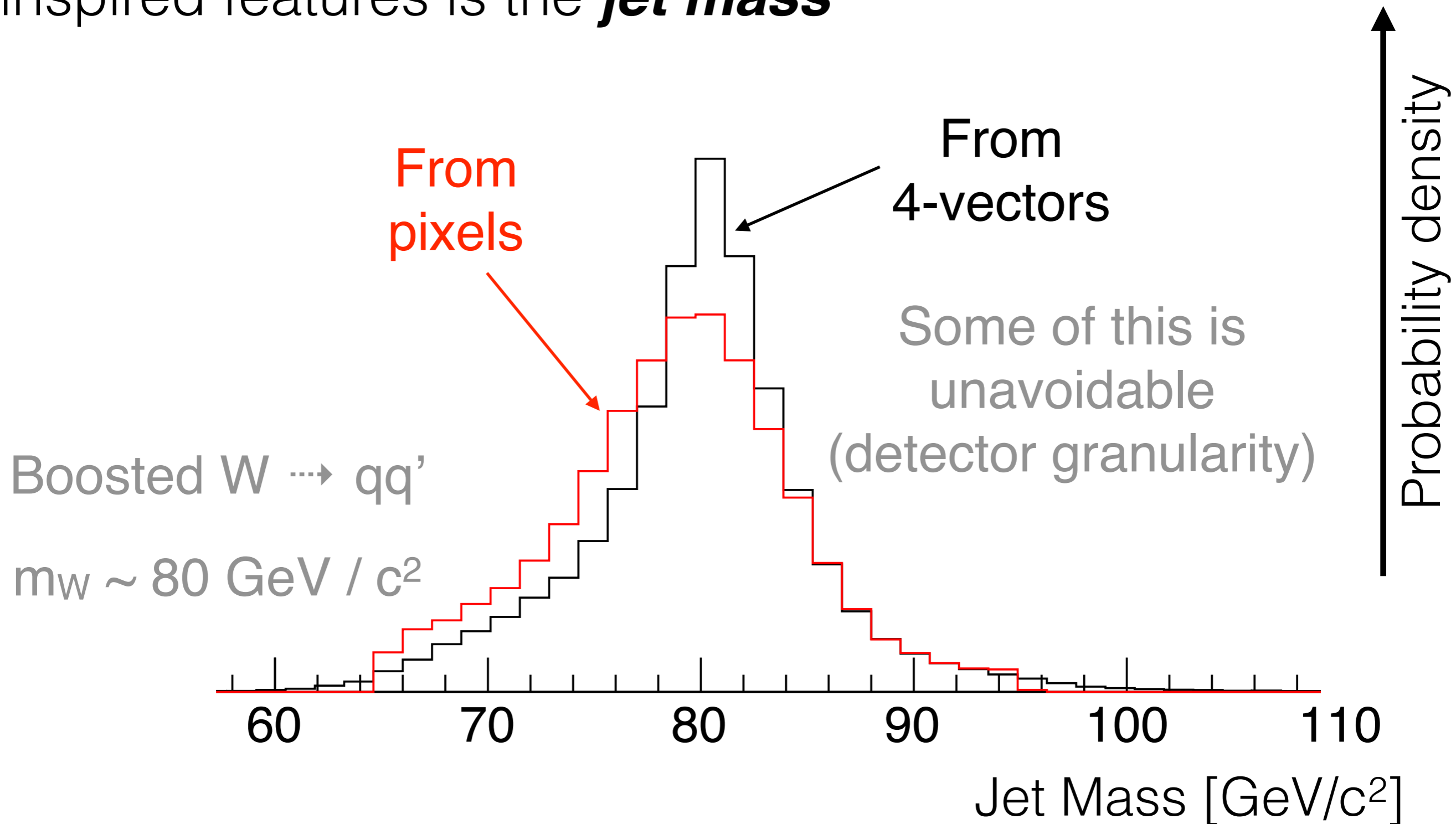
One of the first typical steps is pre-processing



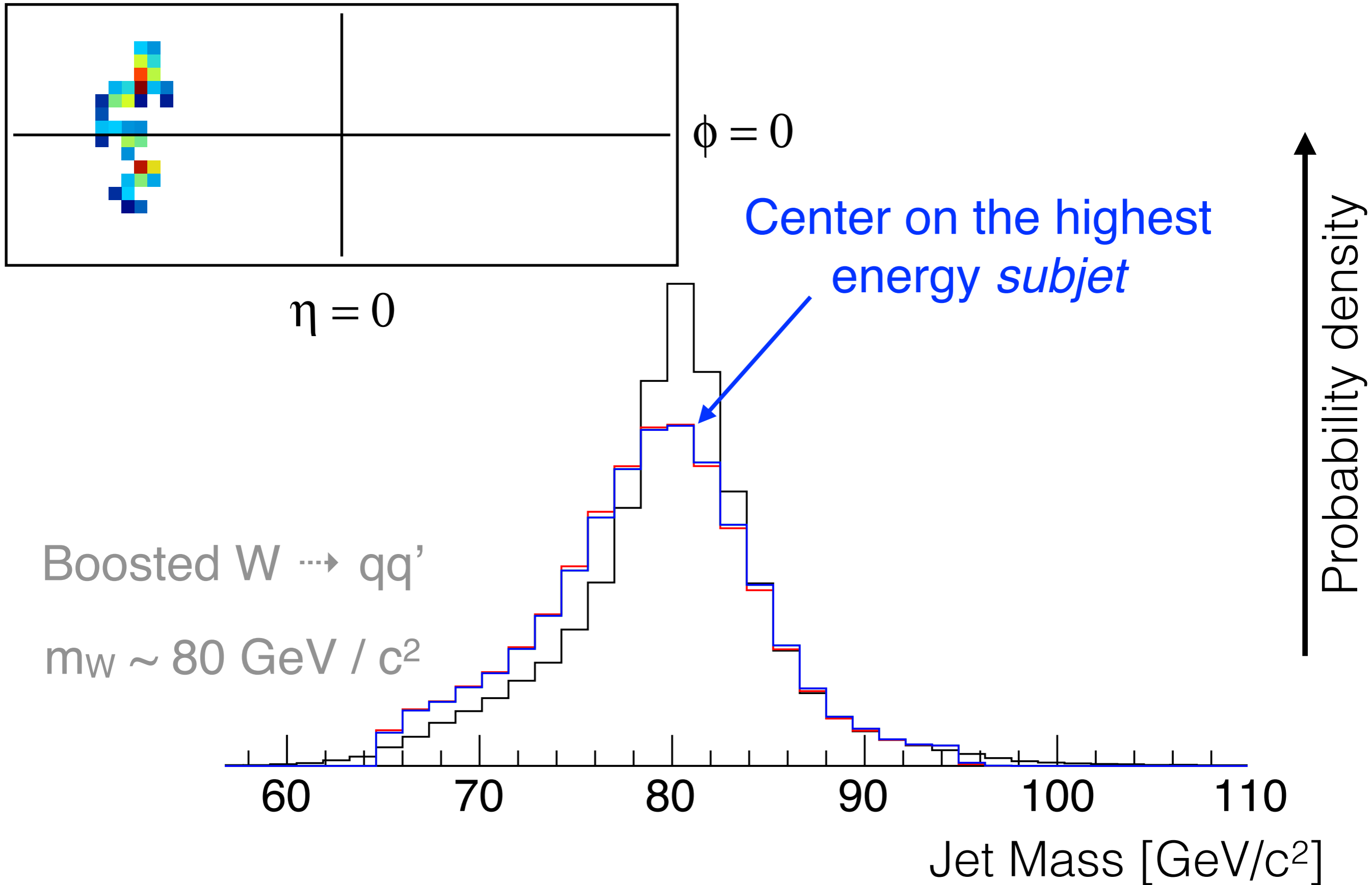
Can help to learn faster & smarter; but must be careful!

Pre-processing & spacetime symmetries

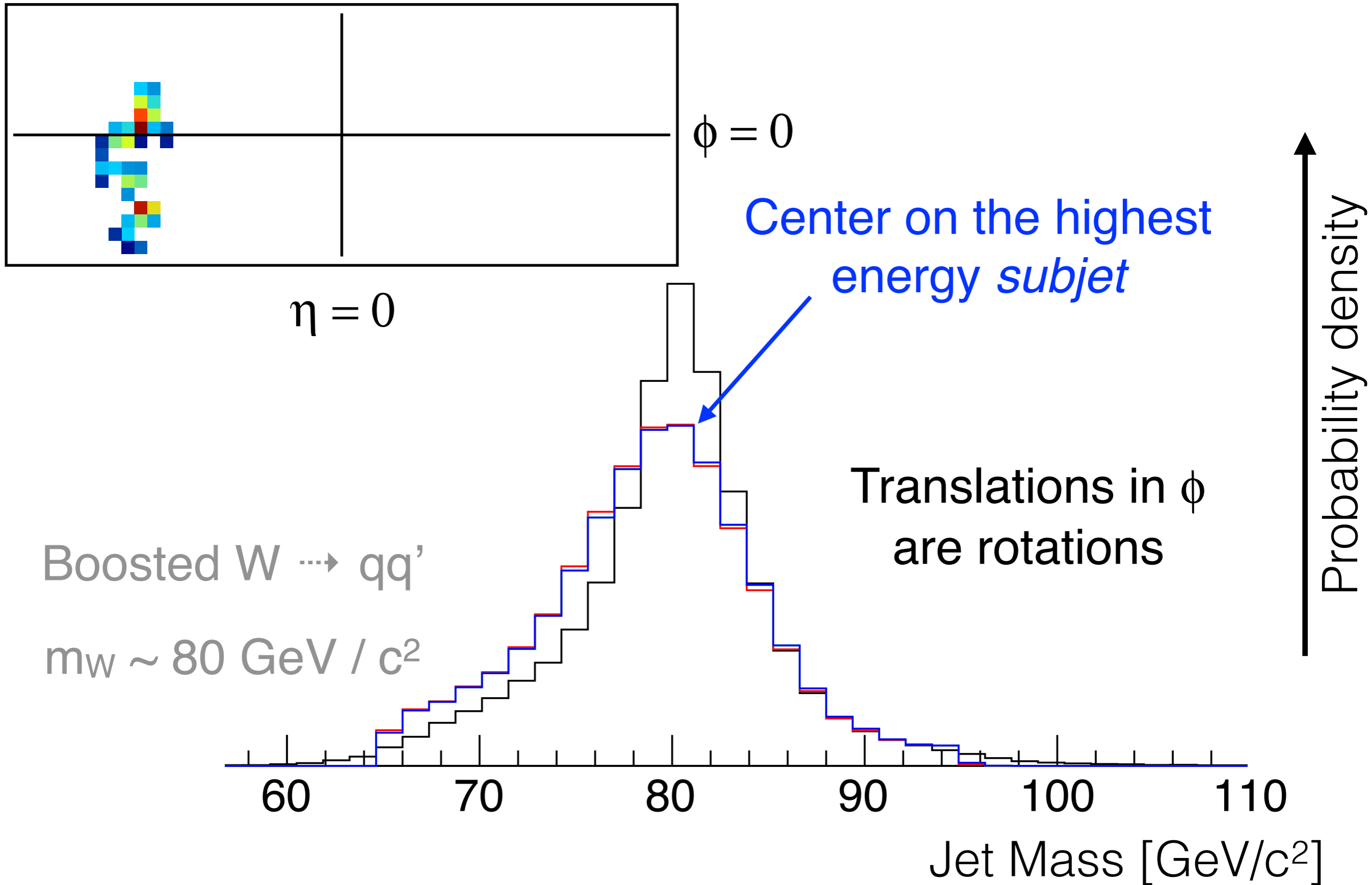
One of the most useful physics-inspired features is the **jet mass**



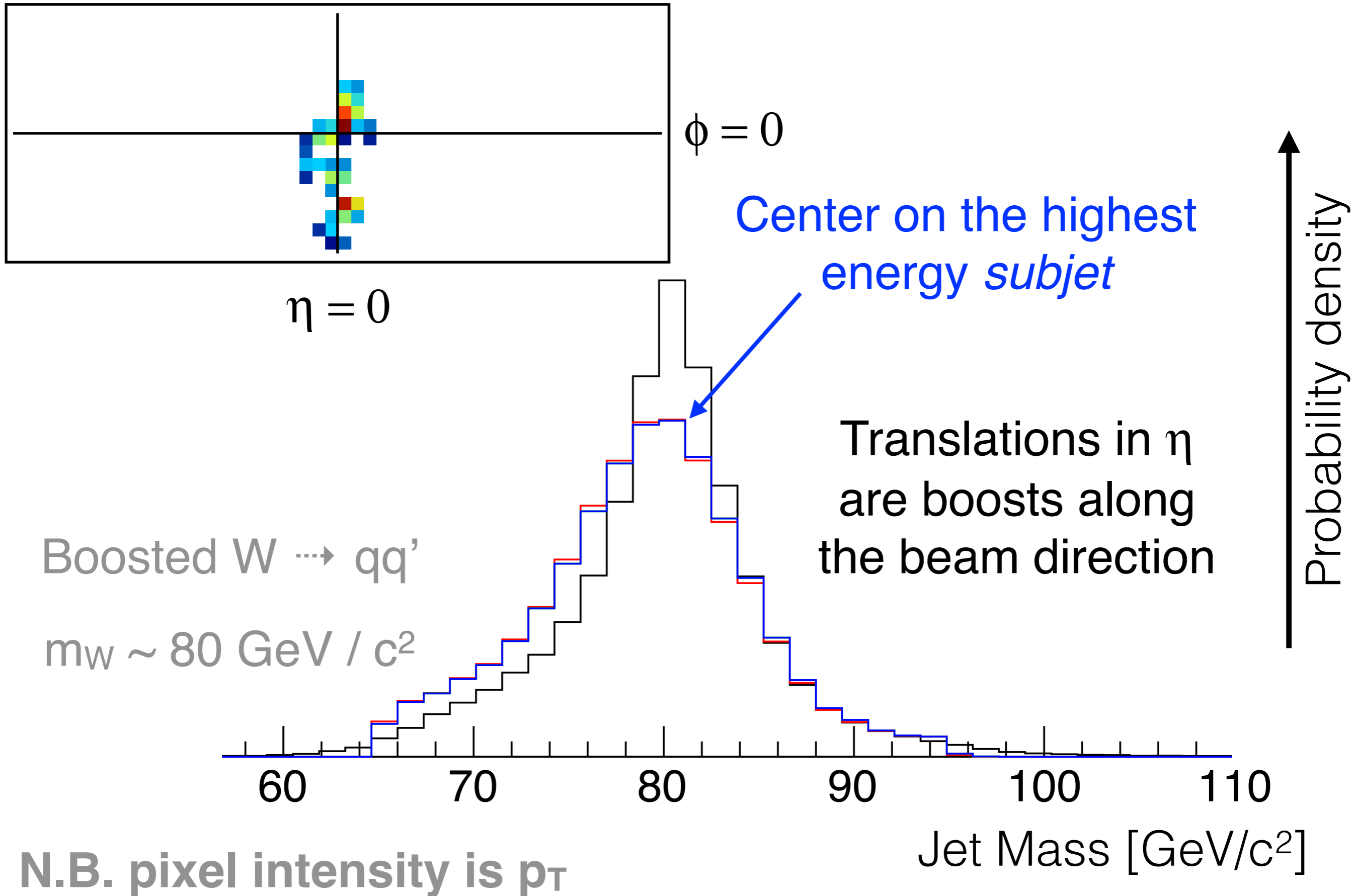
Pre-processing & spacetime symmetries



Pre-processing & spacetime symmetries

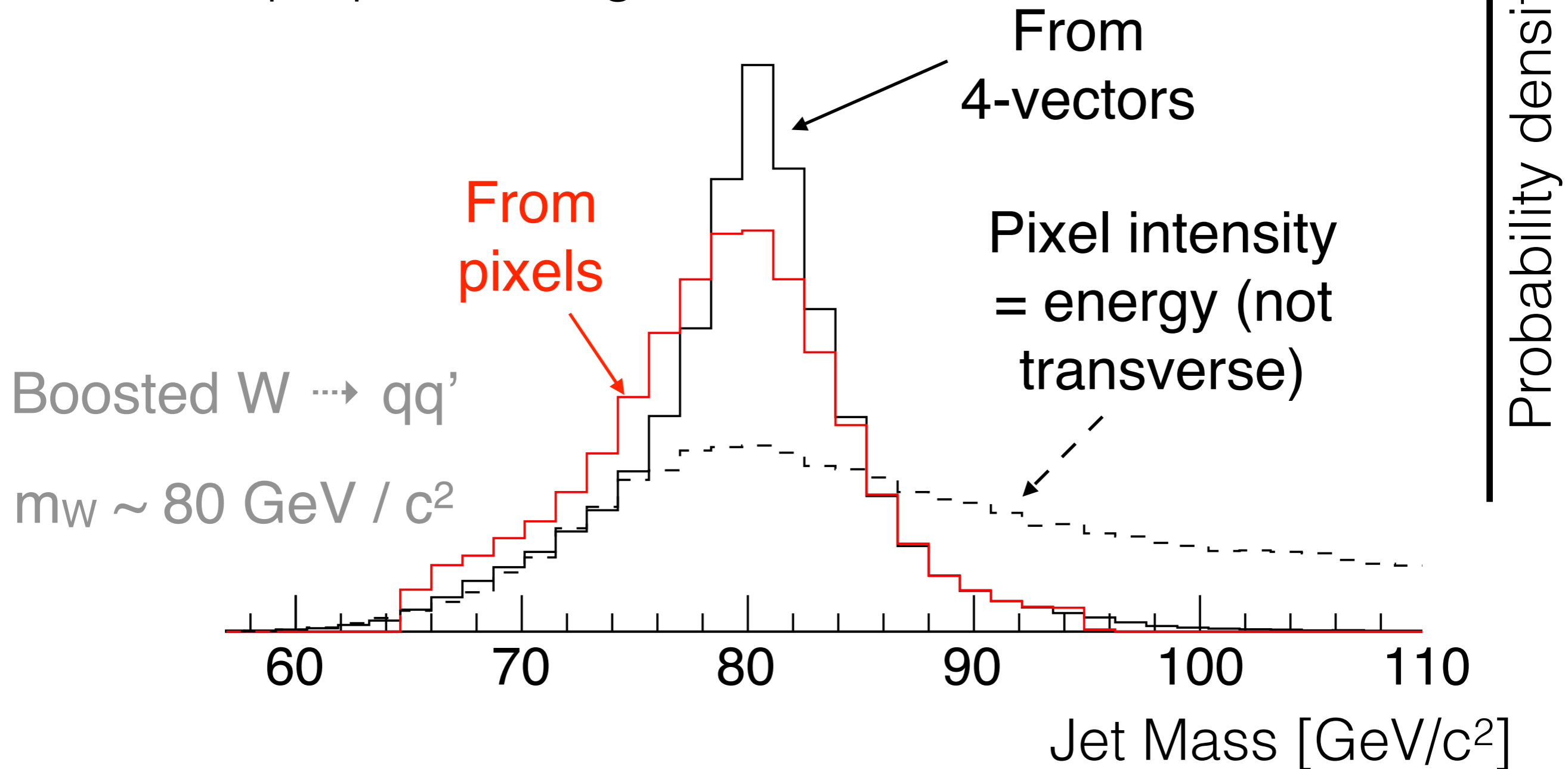


Pre-processing & spacetime symmetries



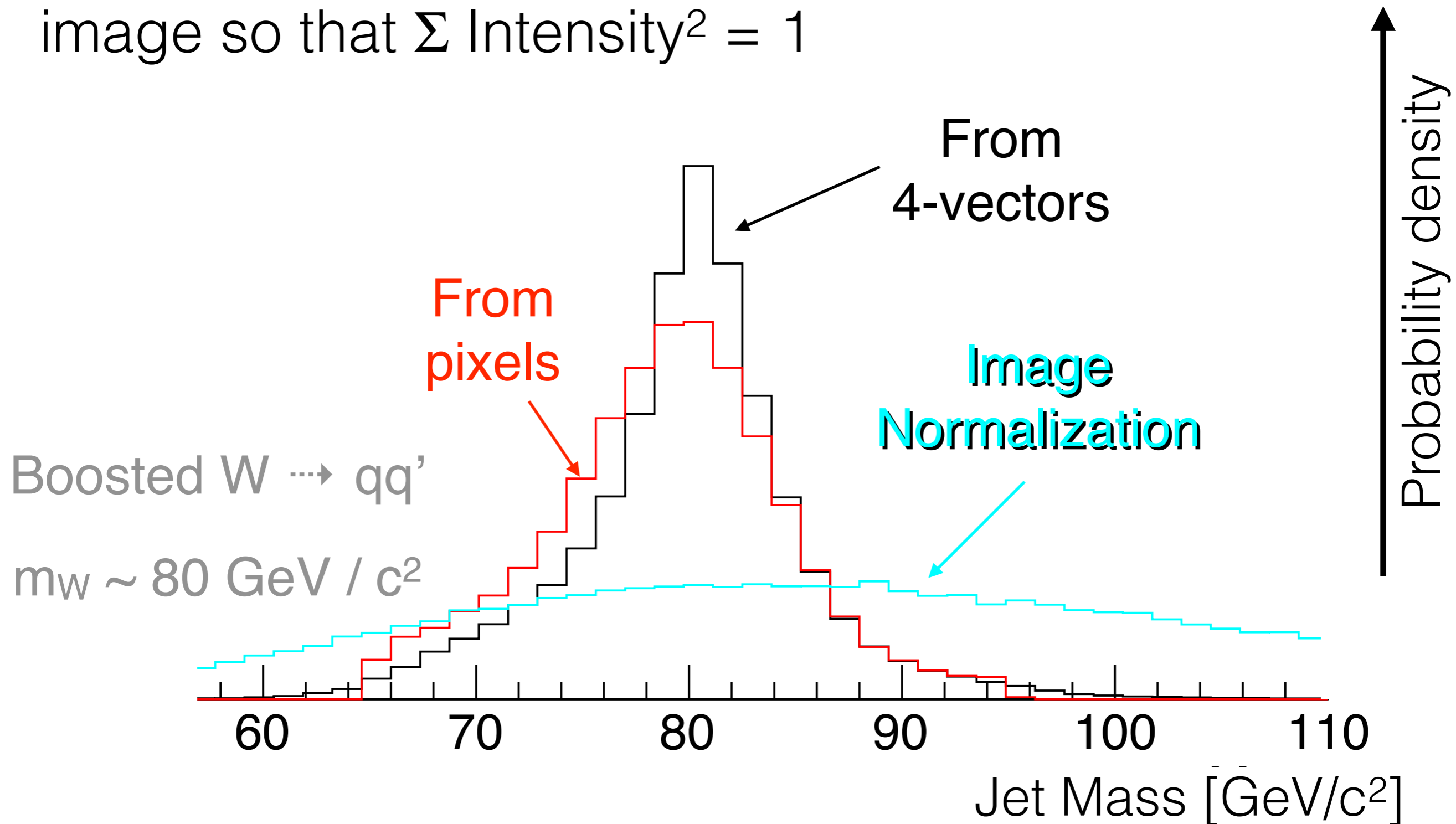
Pre-processing & spacetime symmetries

Information can be **washed out** without care in preprocessing

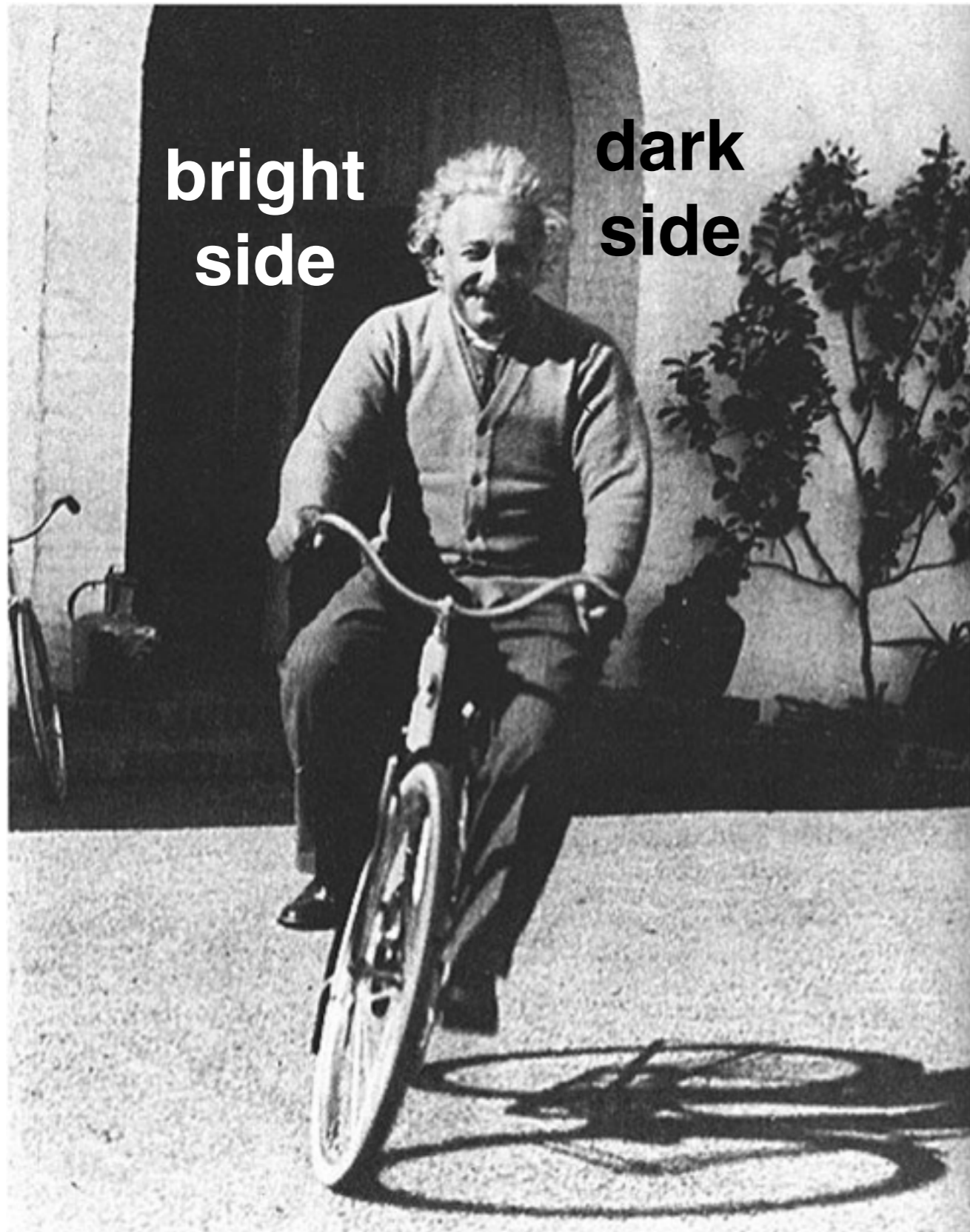


Pre-processing & spacetime symmetries

It is common to normalize each image so that $\sum \text{Intensity}^2 = 1$



Intuition via analogy *why normalization can hurt*



In both pictures, total intensity of Einstein's face is about the same.



However, his face's **image mass** is quite different!

Intuition via analogy *why normalization can hurt*



**bright
side**

**dark
side**

In standard computer vision, you likely don't want to be sensitive to this! ...not the case for jet images!

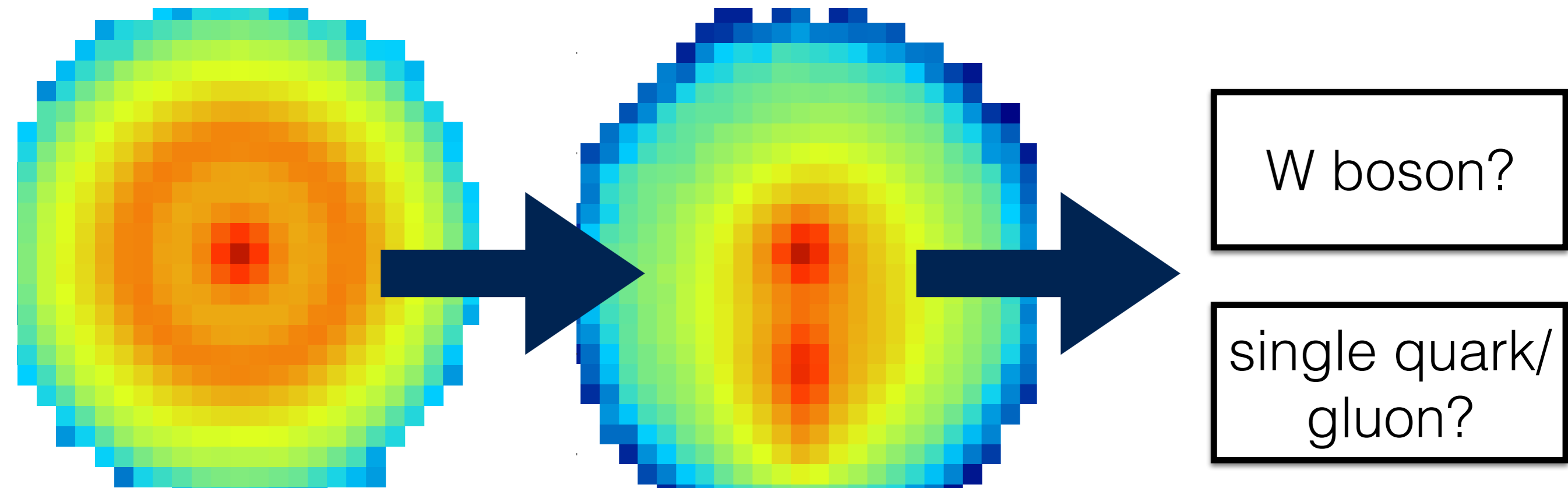
In both pictures, total intensity of Einstein's face is about the same.



**uniform moderate
intensity**

However, his face's **image mass** is quite different!

Now, with a carefully processed image, we can ask: where did this jet come from?

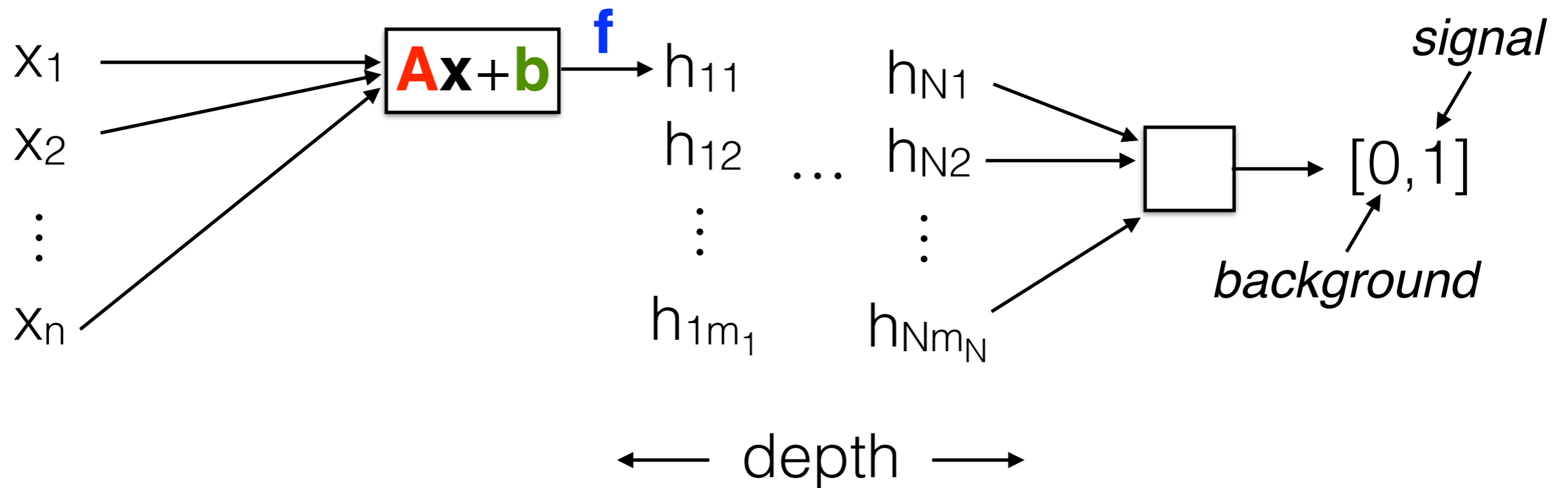


ultimate classification is achieved with modern machine learning using **all pixels as input!**

Modern Deep NN's for Classification

Neural Network: composition of functions $f(\mathbf{Ax}+\mathbf{b})$ for inputs \mathbf{x} (features) matrix \mathbf{A} (weights), bias \mathbf{b} , non-linearity f .

N.B. I'm not mentioning biology - there may be a vague resemblance to parts of the brain, but that is not what modern NN's are about.



Modern Deep NN's for Classification

Neural Network: composition of functions $f(\mathbf{Ax} + \mathbf{b})$ for inputs \mathbf{x} (features) matrix \mathbf{A} (weights), bias \mathbf{b} , non-linearity f .

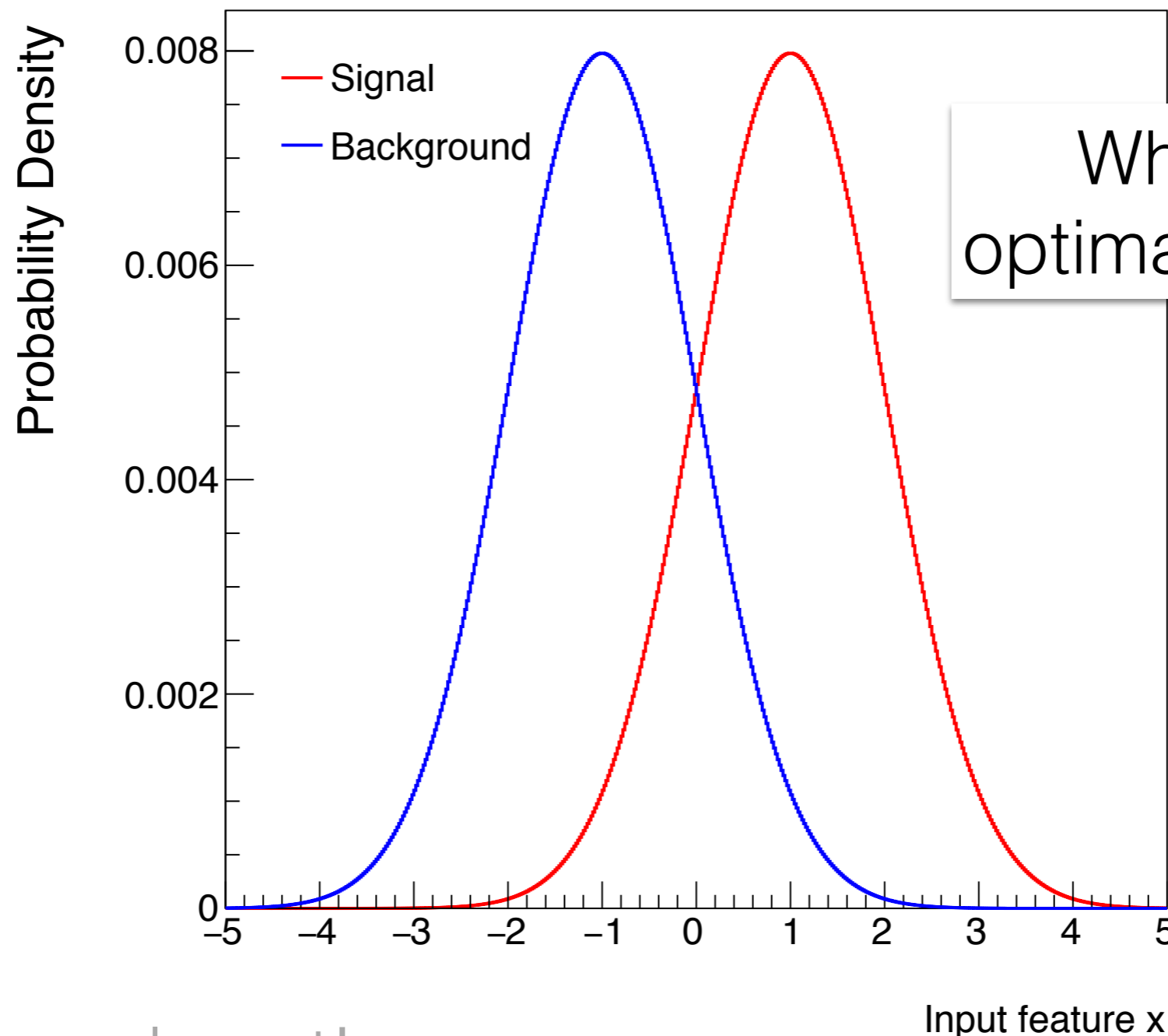
Fact: NN's can approximate "any" function.

*Why
useful?*

For classification, there is an optimal function to learn: the likelihood ratio, $LL(x) = p_S(x) / p_B(x)$.

Getting into the machine's mind

Let's consider an important special case:
binary classification in 1D



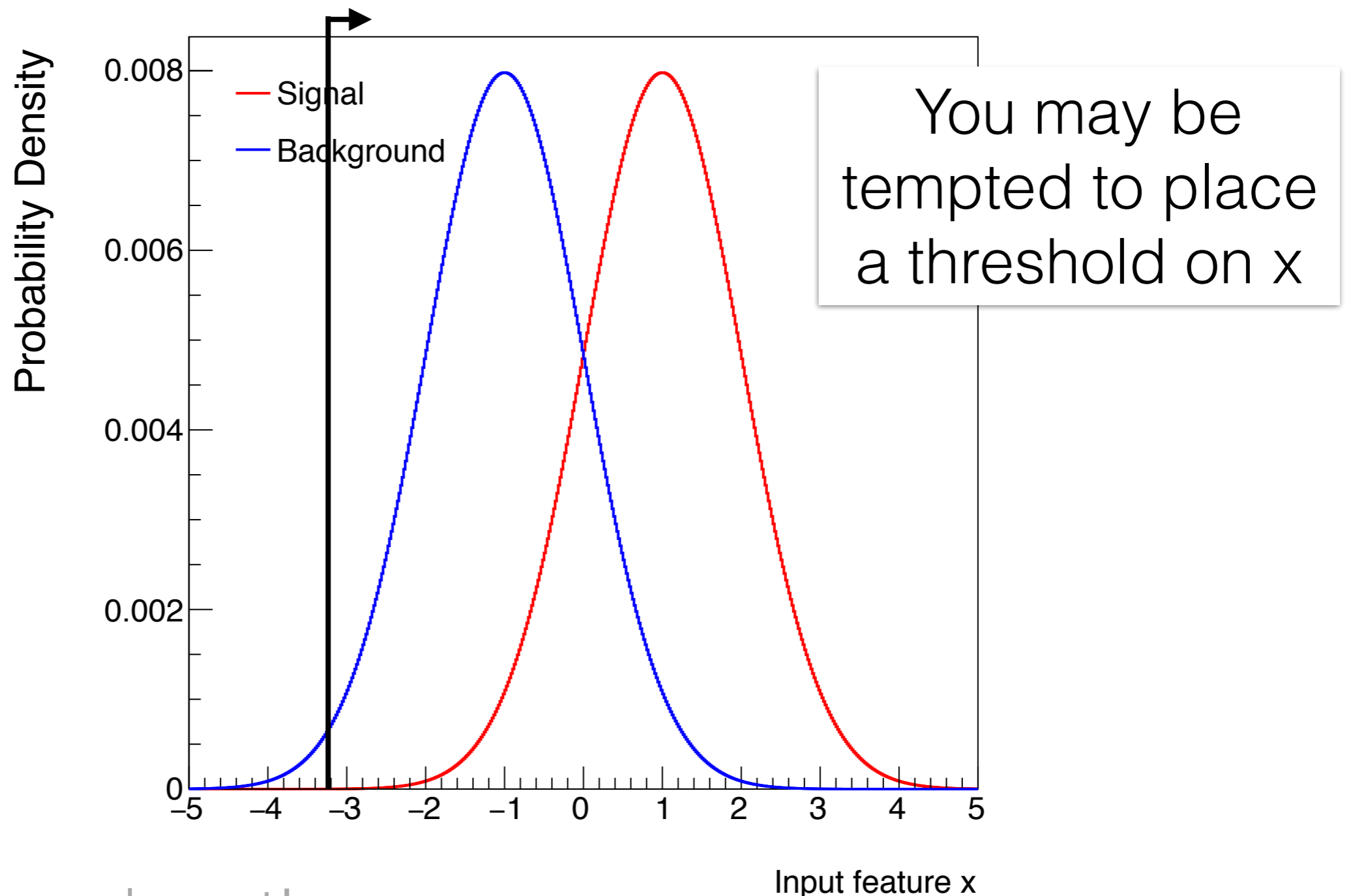
What is the optimal classifier?

could be e.g.
the jet mass

⇒ Try this example out!

Getting into the machine's mind

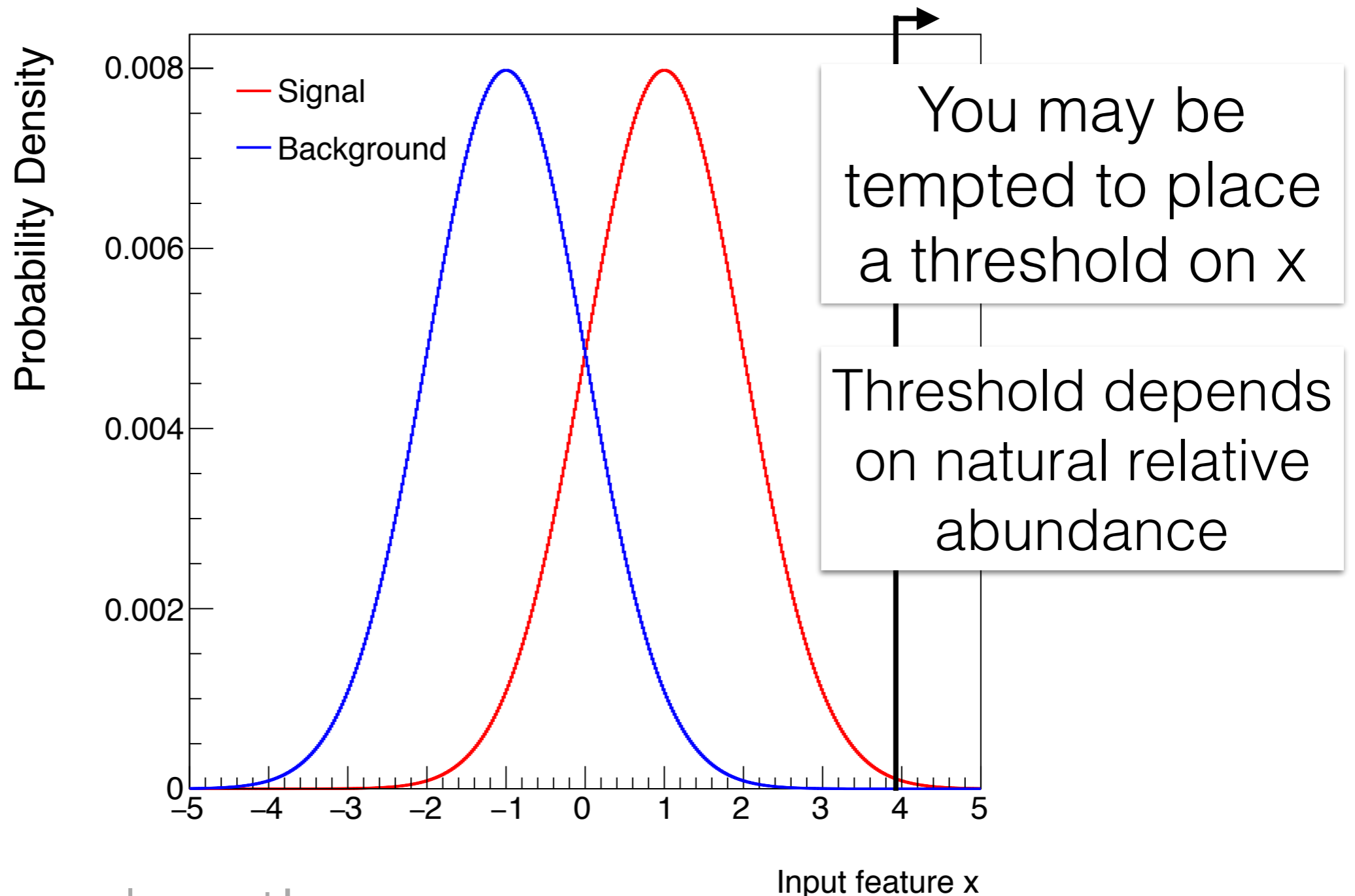
Let's consider an important special case:
binary classification in 1D



⇒ Try this example out!

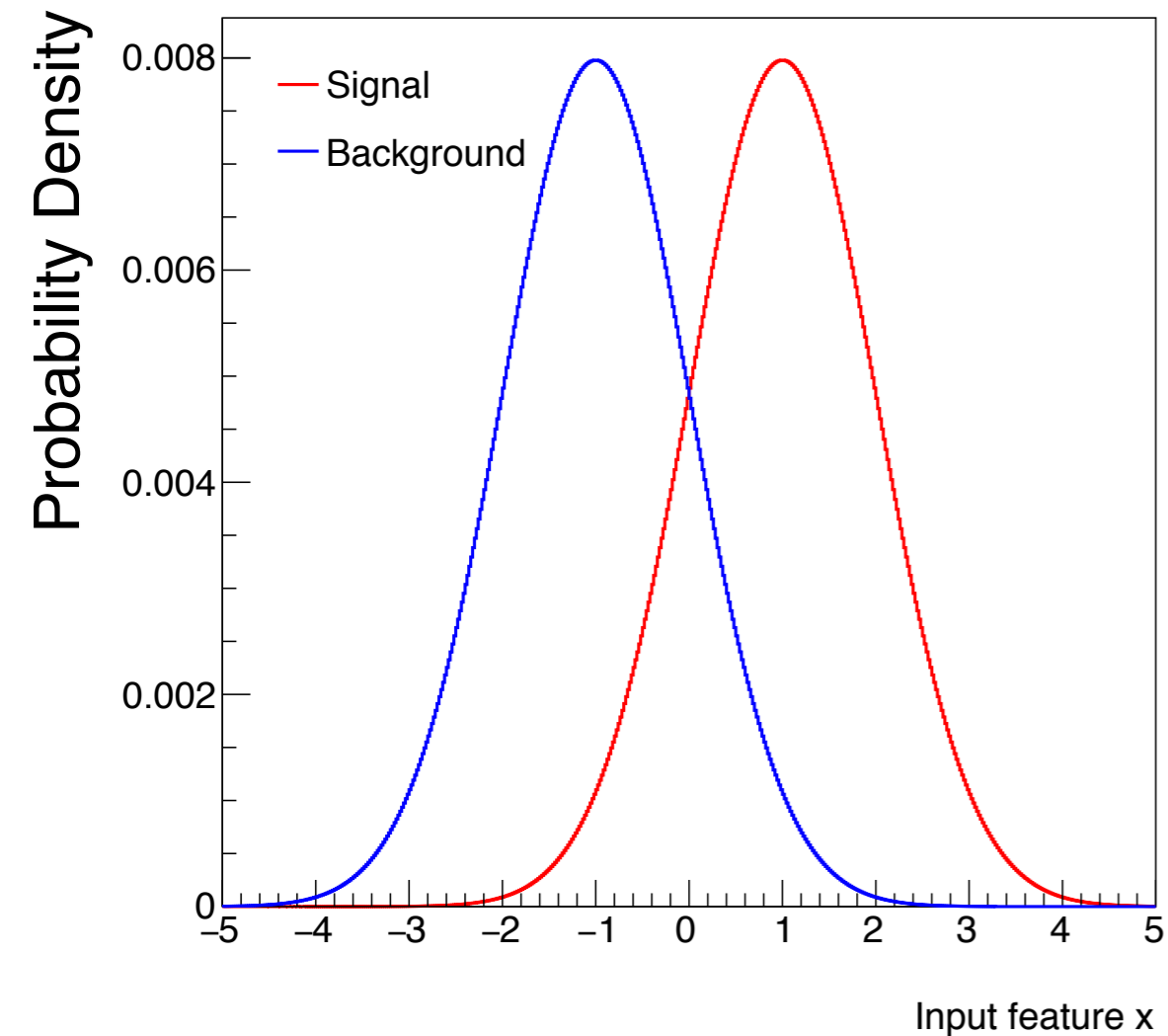
Getting into the machine's mind

Let's consider an important special case:
binary classification in 1D

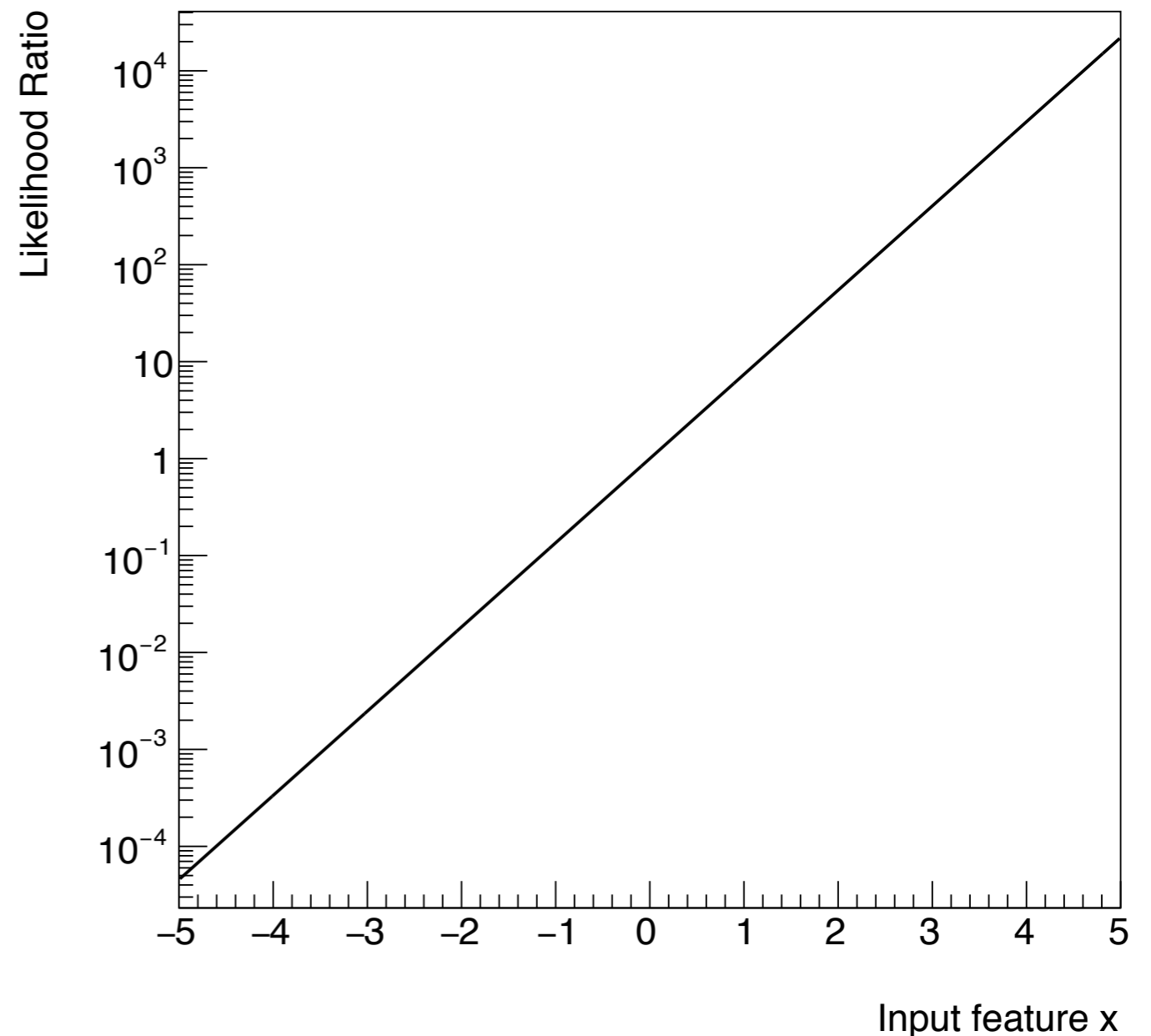


⇒ Try this example out!

Getting into the machine's mind



Is the simple threshold cut optimal?

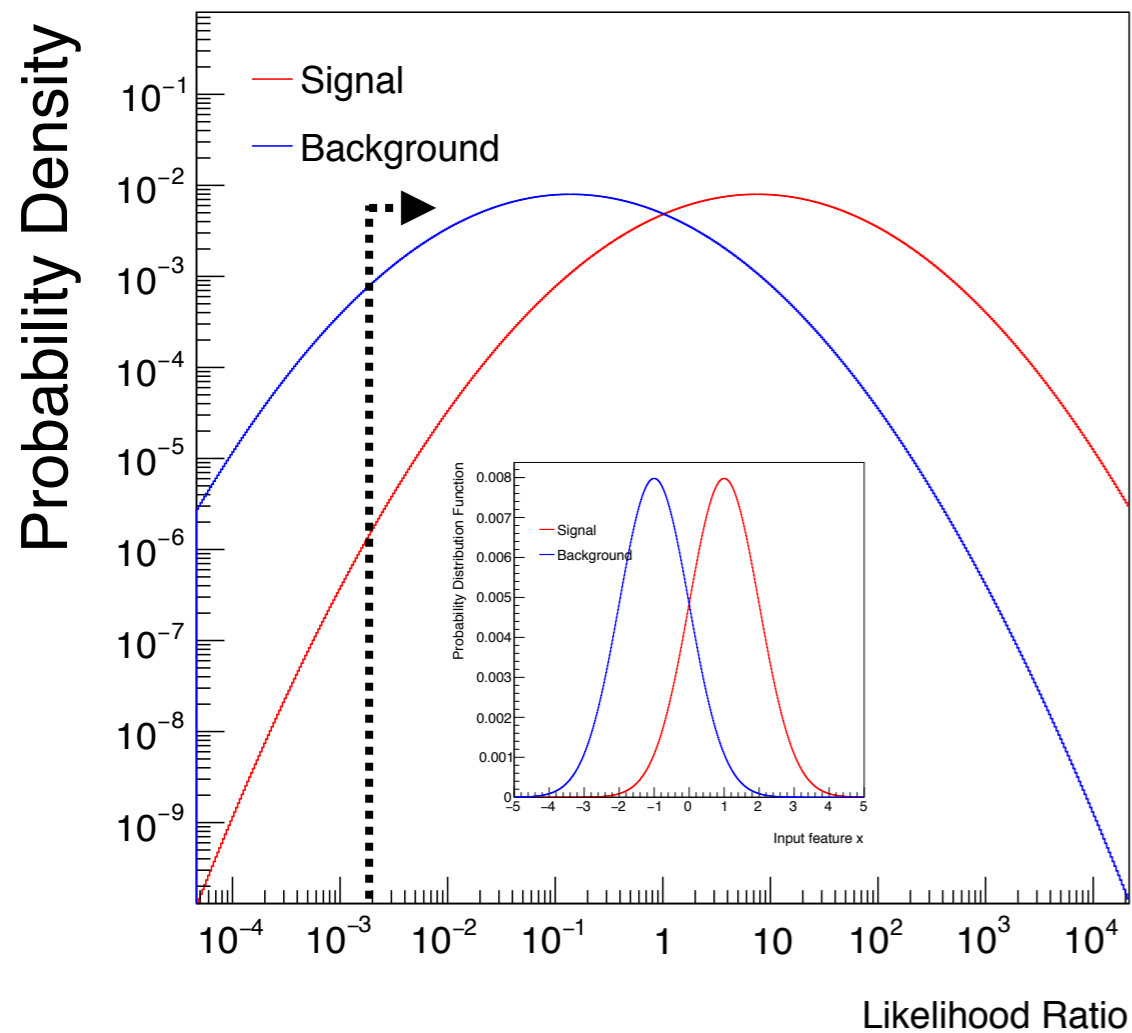


In this simple case, the log LL is proportional to x:

no need for non-linearities!

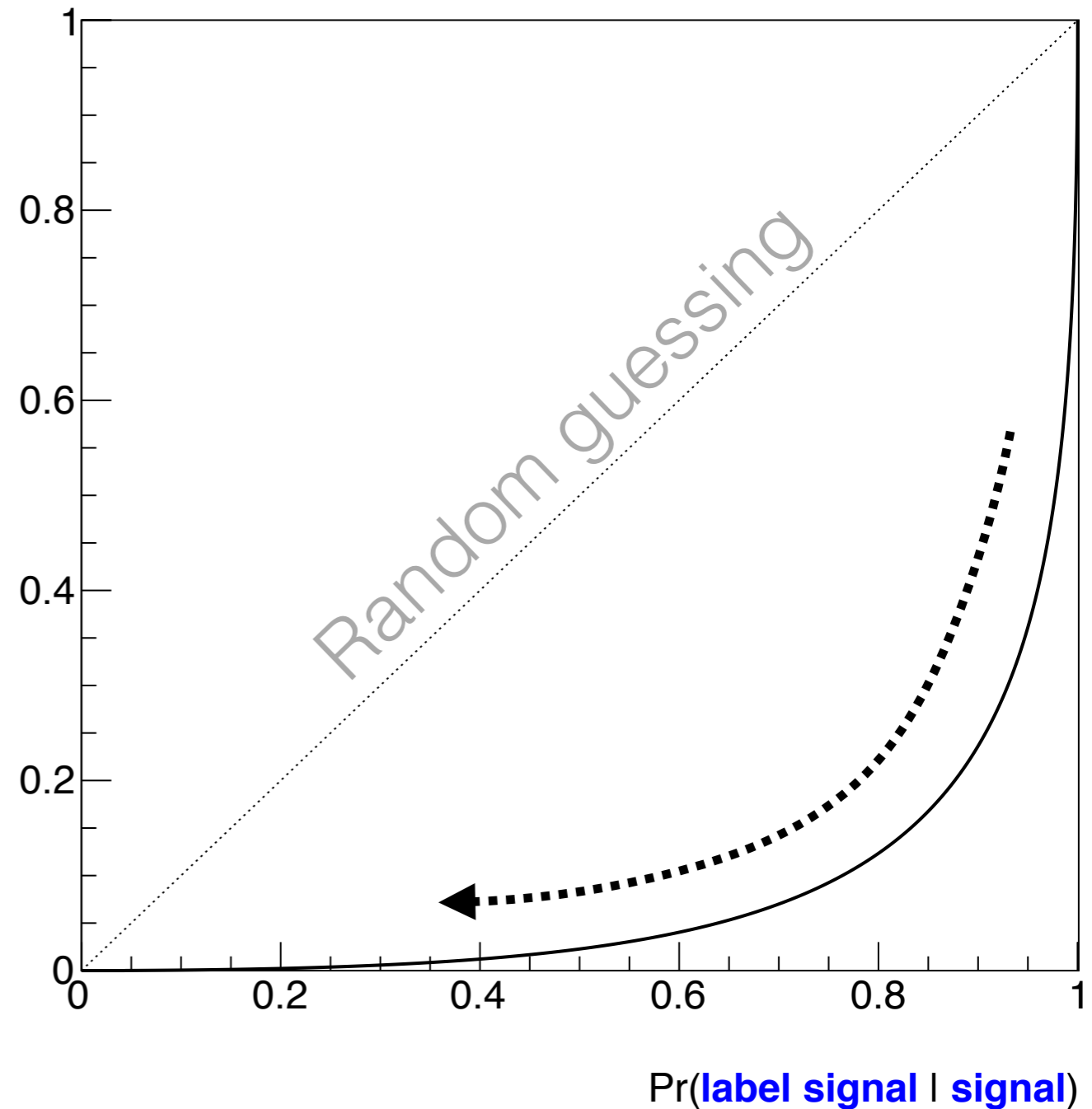
Threshold cut is optimal

Getting into the machine's mind



The optimal procedure is a threshold on the LL

$\Pr(\text{label signal} \mid \text{background})$



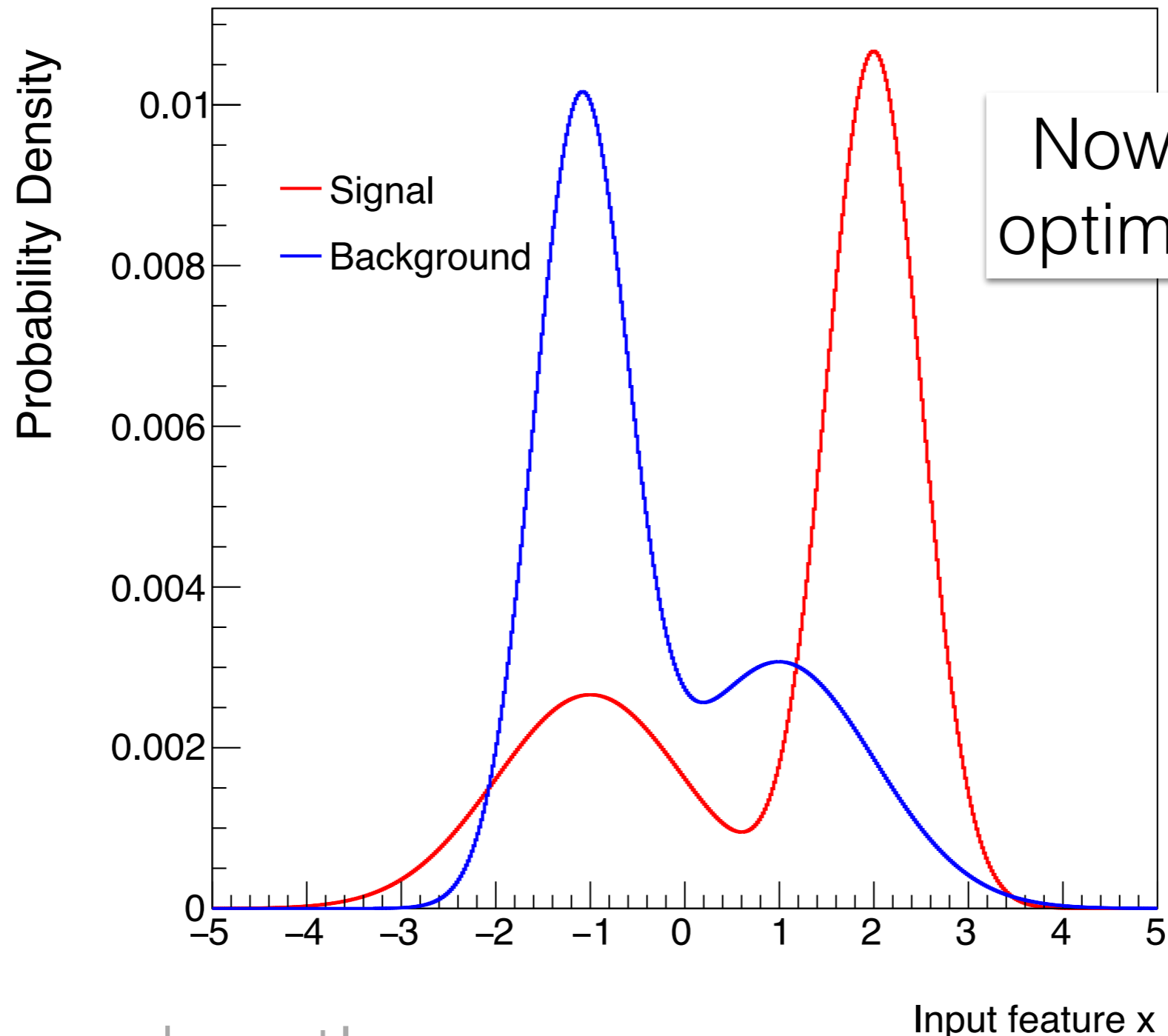
“Receiver Operating Characteristic” (**ROC**) Curve

⇒ Try this example out!

Getting into the machine's mind

What if the distribution of x is complicated?

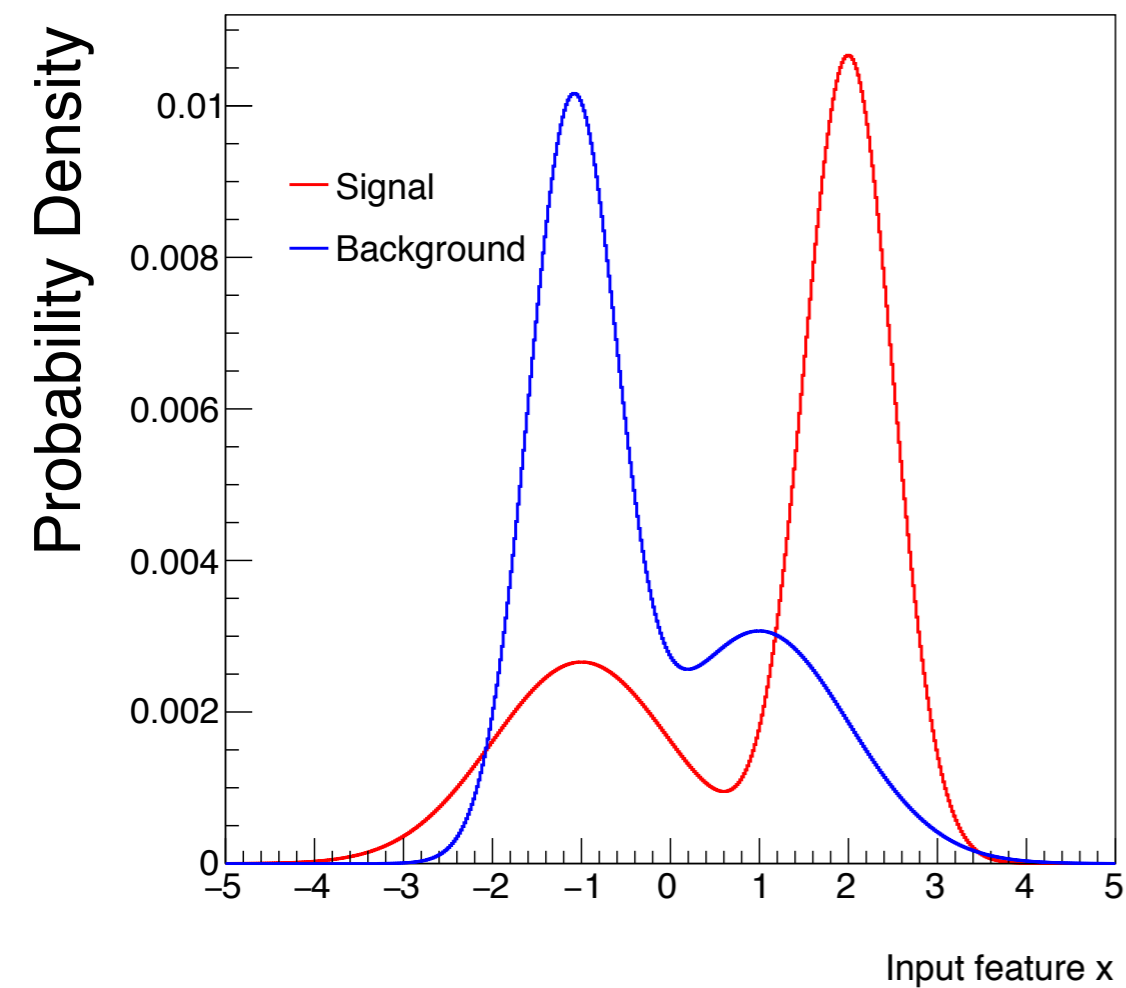
Real life is complicated!



Now what is the optimal classifier?

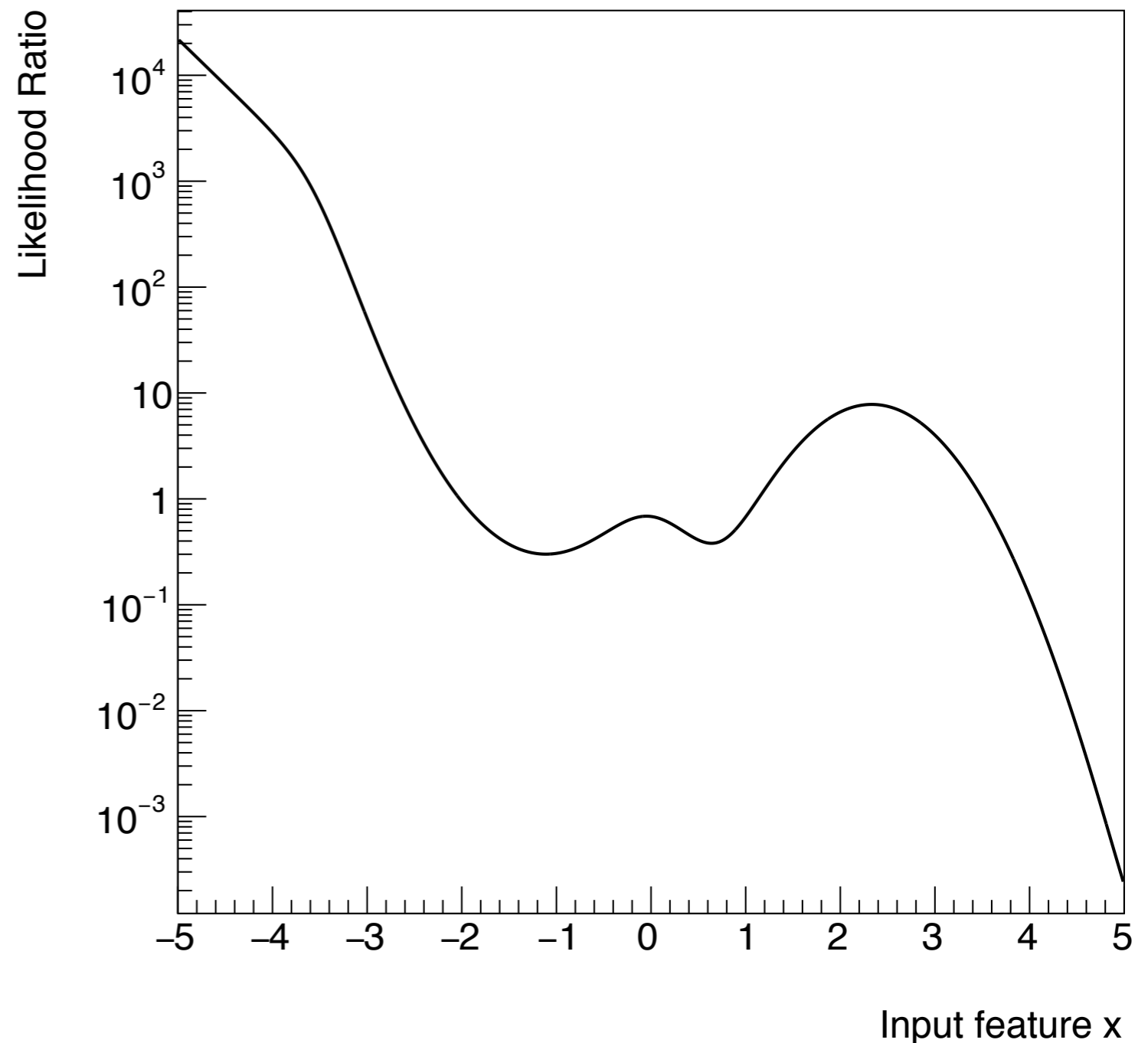
⇒ Try this example out!

Getting into the machine's mind



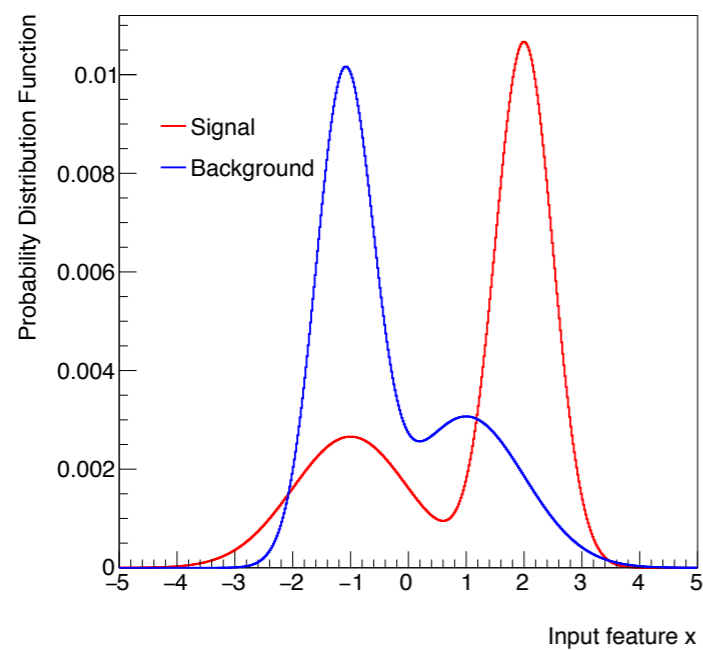
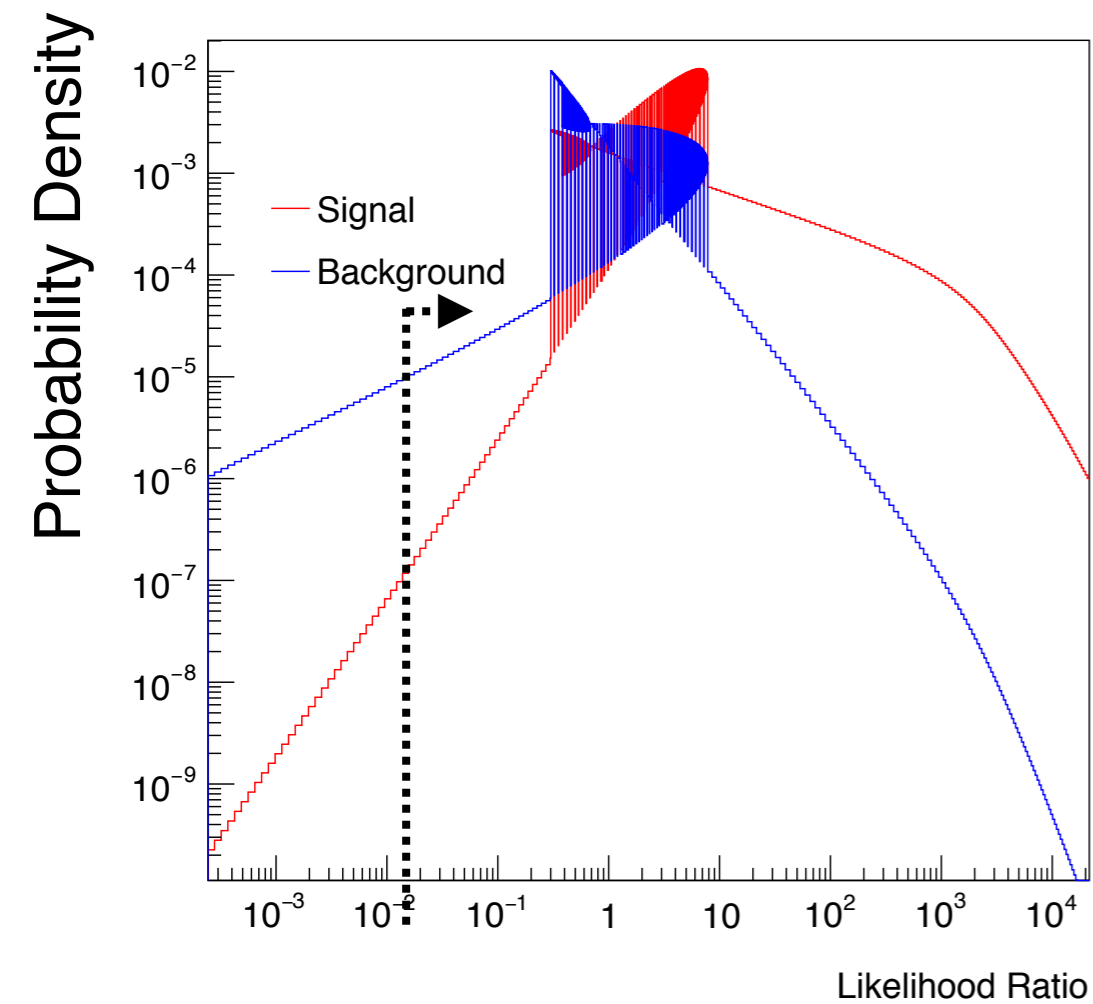
A threshold on x
would be sub-optimal

In this case, LL is highly
non-linear function of x

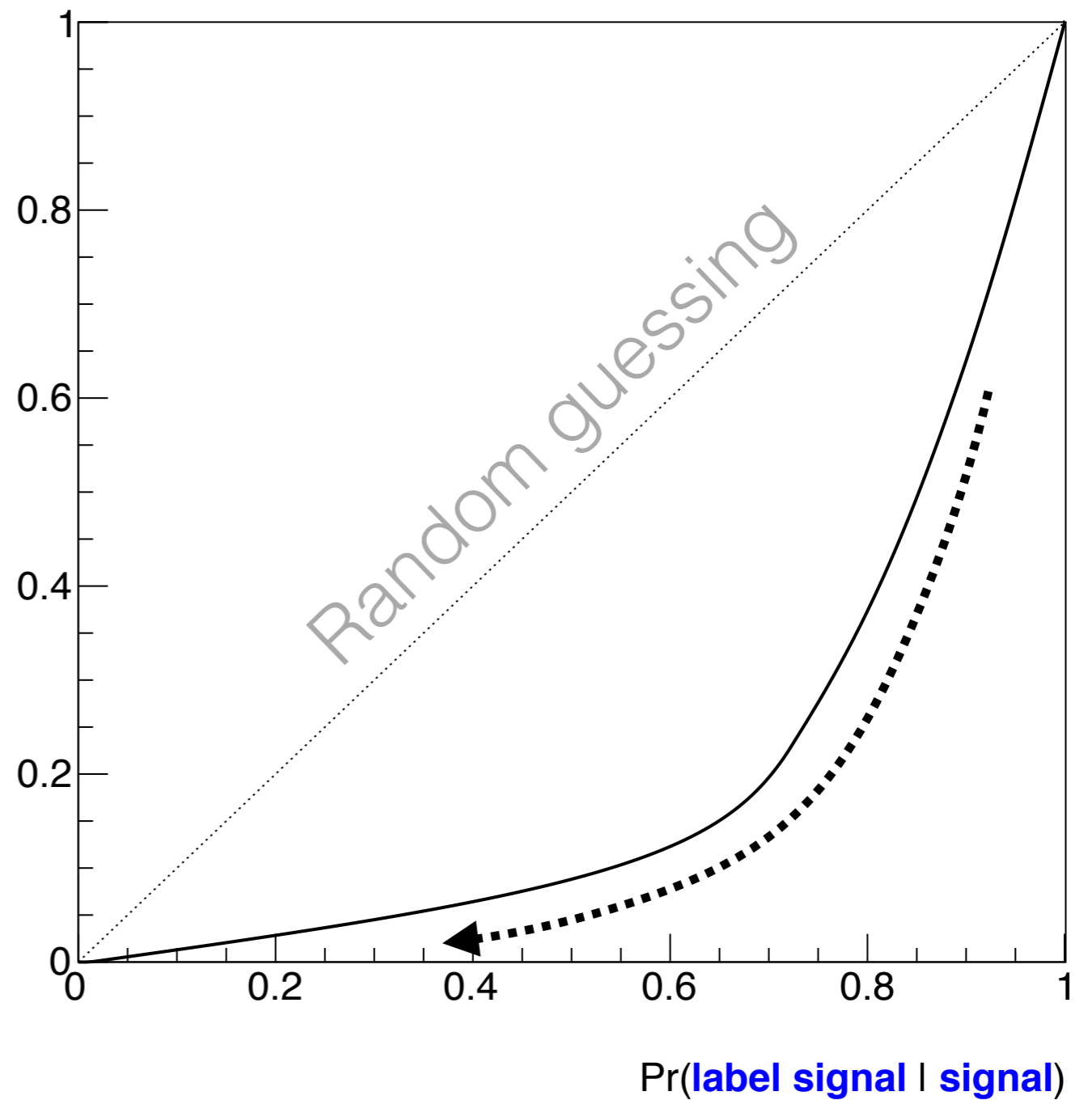


⇒ Try this example out!

Getting into the machine's mind



$\Pr(\text{label signal} \mid \text{background})$

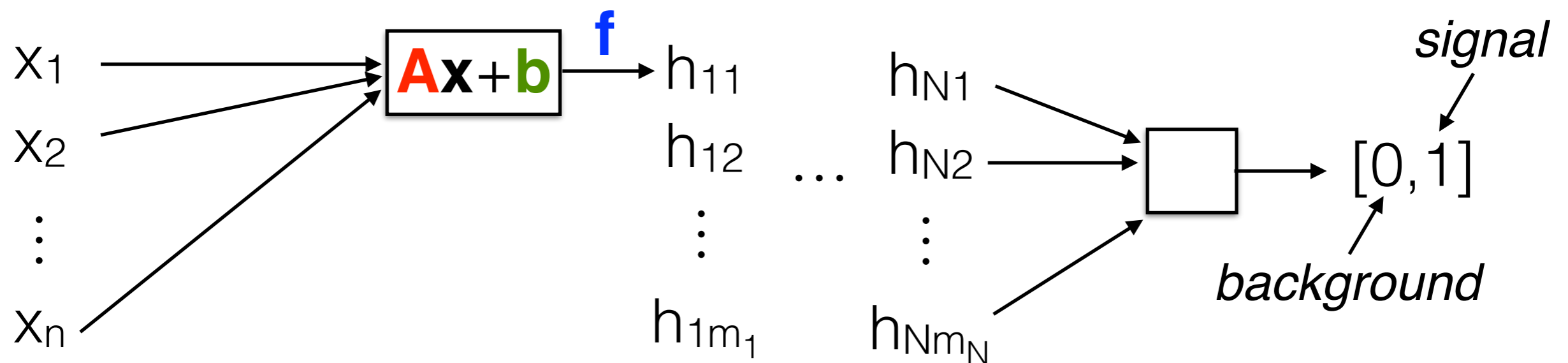


The curse of dimensionality

In principle, you can do the same thing in $N > 1$ dimensions. However, it very quickly gets out of hand!

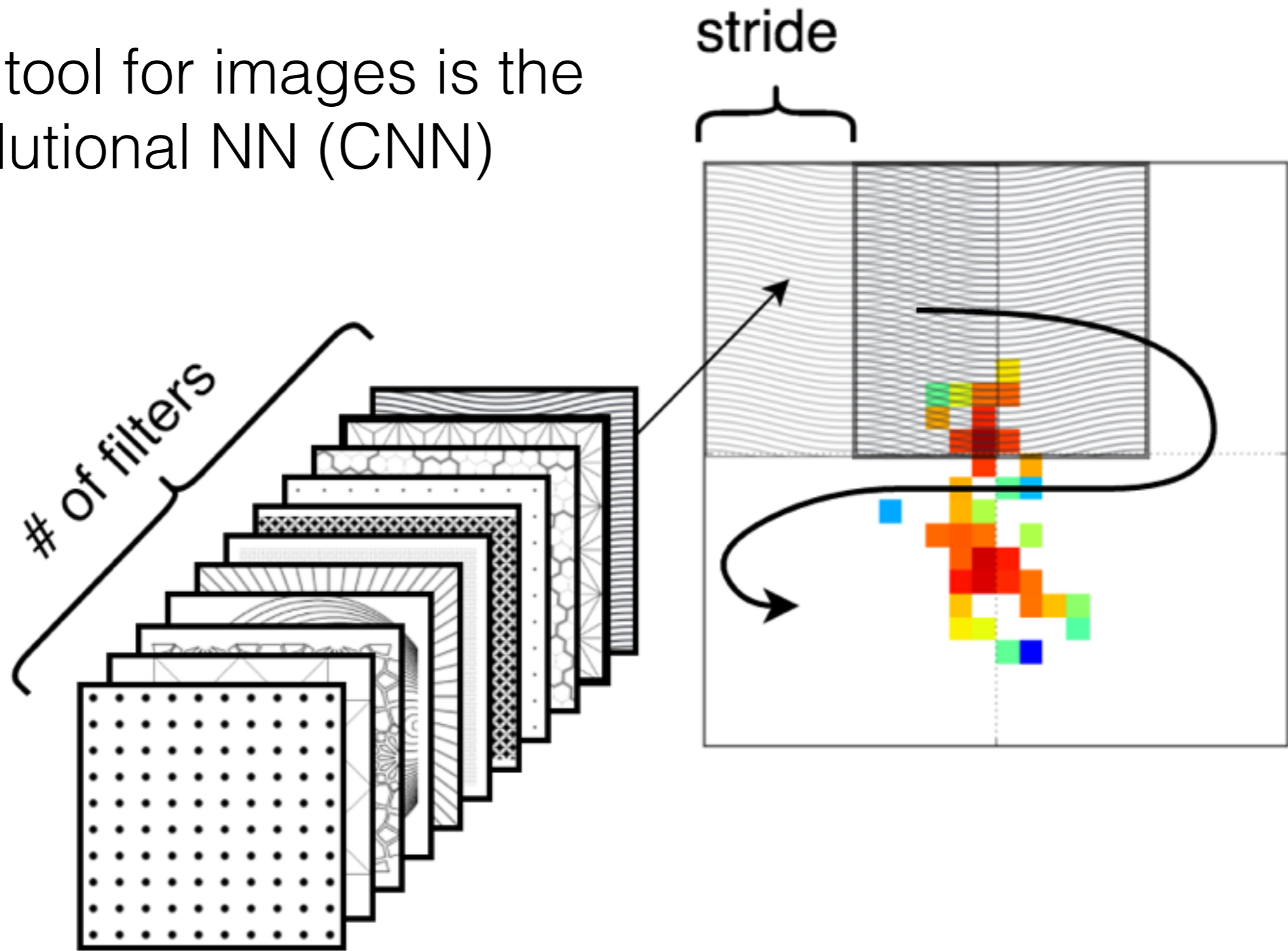
That is where NN's come in.

Image ~ 1000 dimensional



Let's see how we can use DNN's for jet image classification

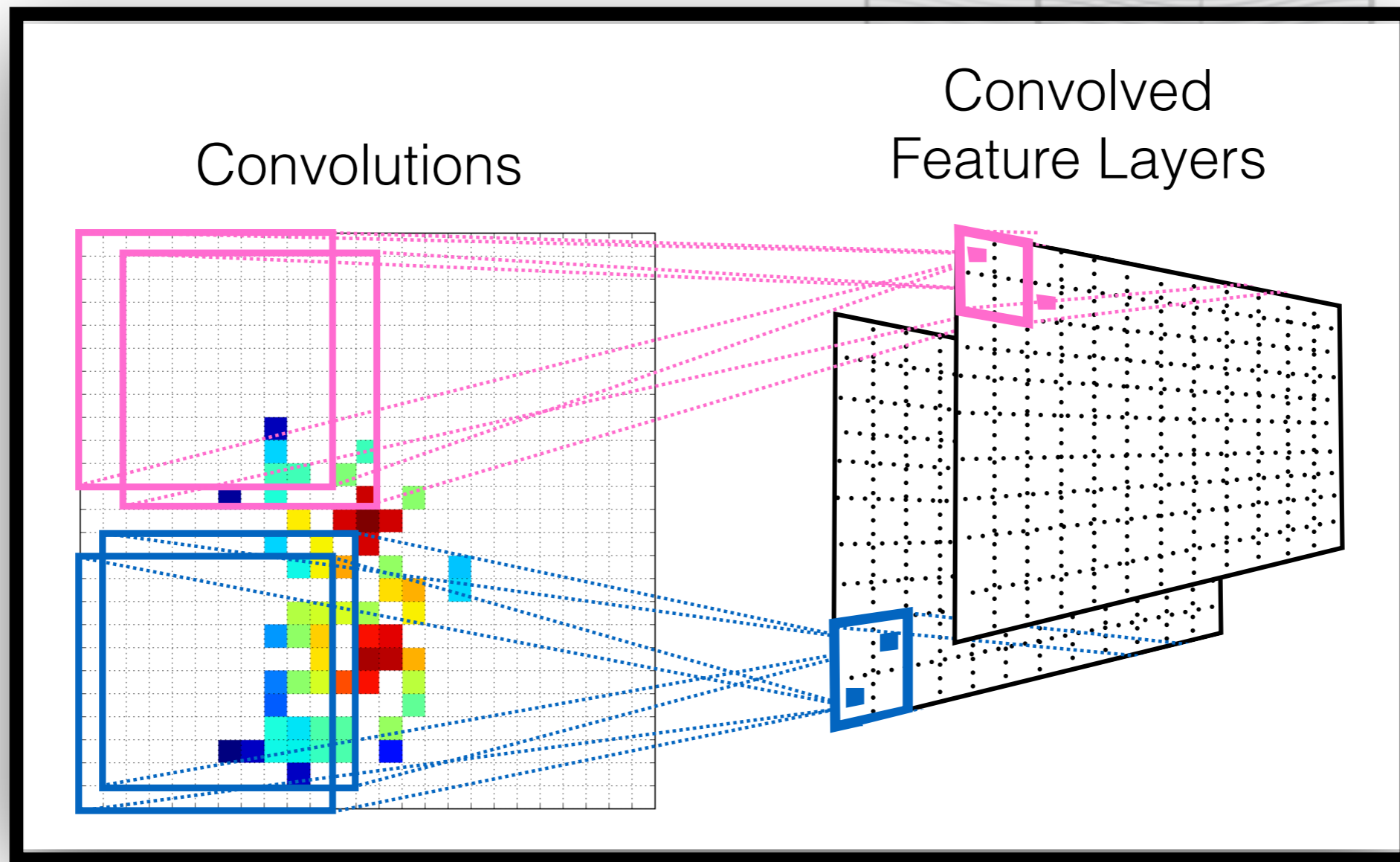
Common tool for images is the convolutional NN (CNN)



The filter is like the **A**, only the dimensionality is now the filter size ($\ll n$) and not the image size (n).

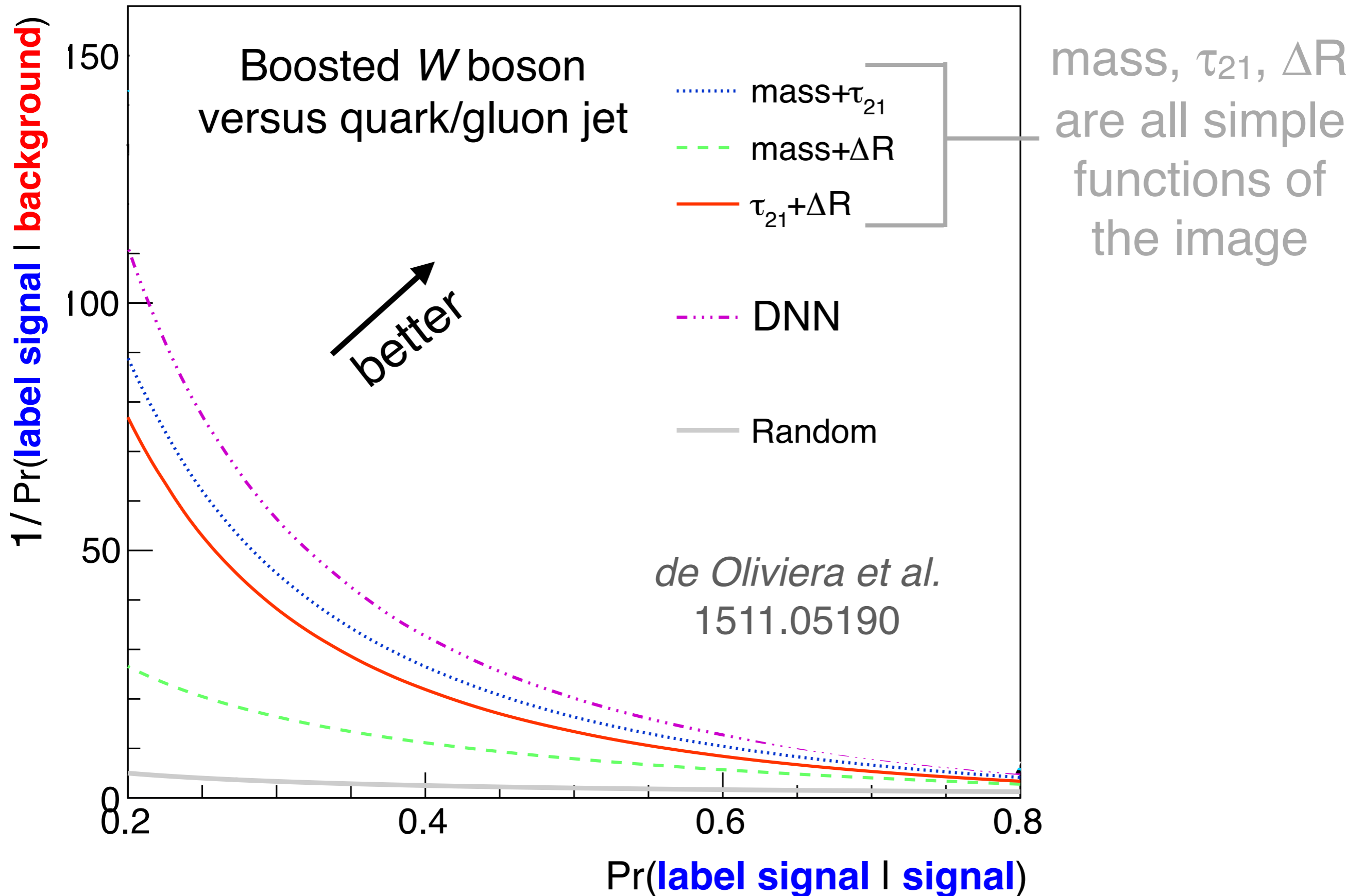
Common tool for images is the
convolutional NN (CNN)

stride

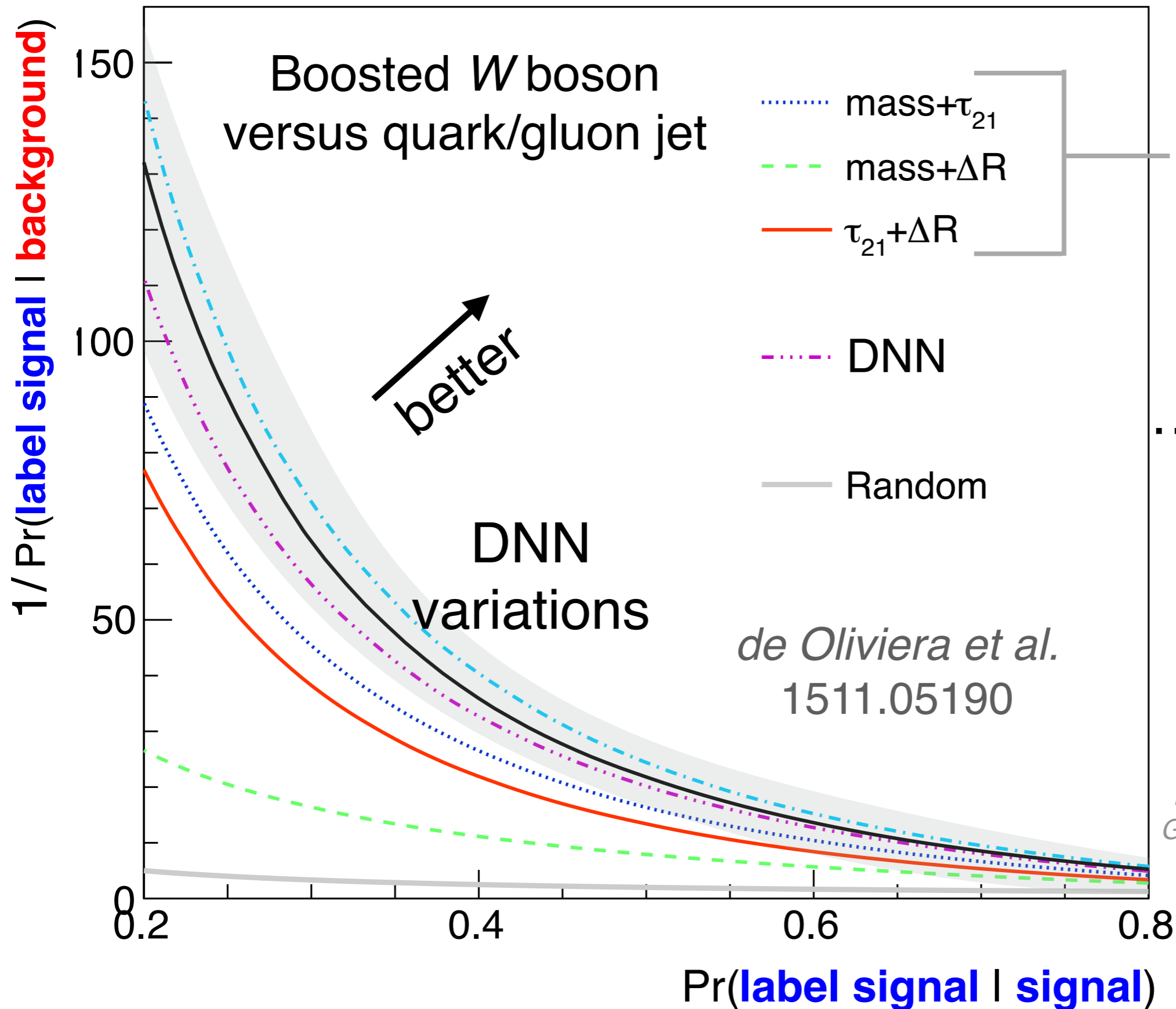


The filter is like the **A**, only the dimensionality is now
the filter size ($\ll n$) and not the image size (n).

Modern Deep NN's for Classification



Modern Deep NN's for Classification



mass, τ_{21} , ΔR are all simple functions of the image

...what the DNN is learning is active R&D!

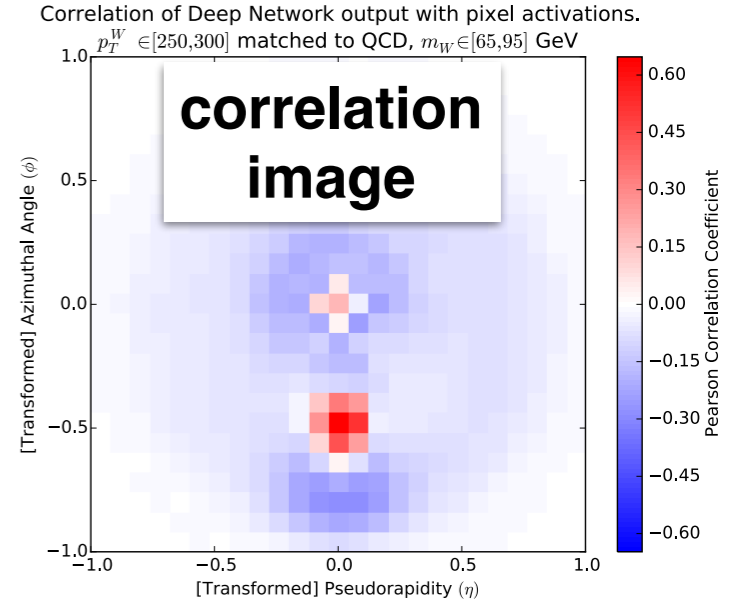
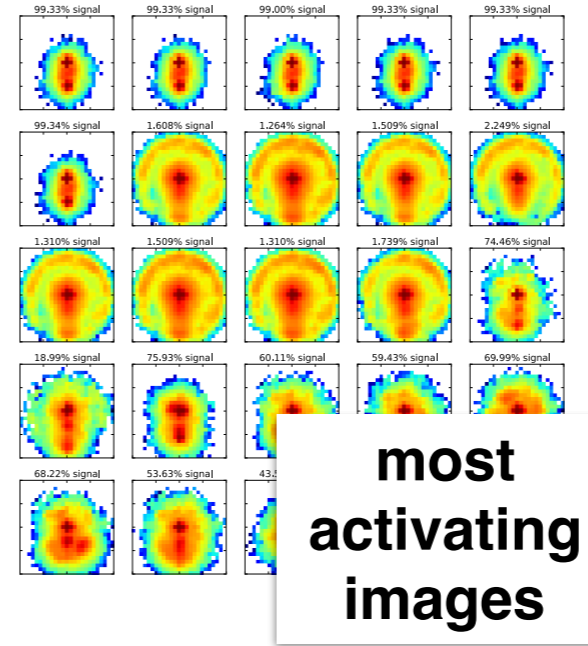
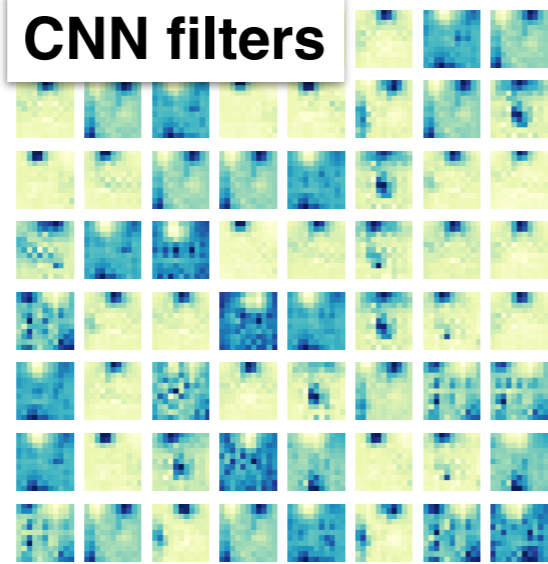
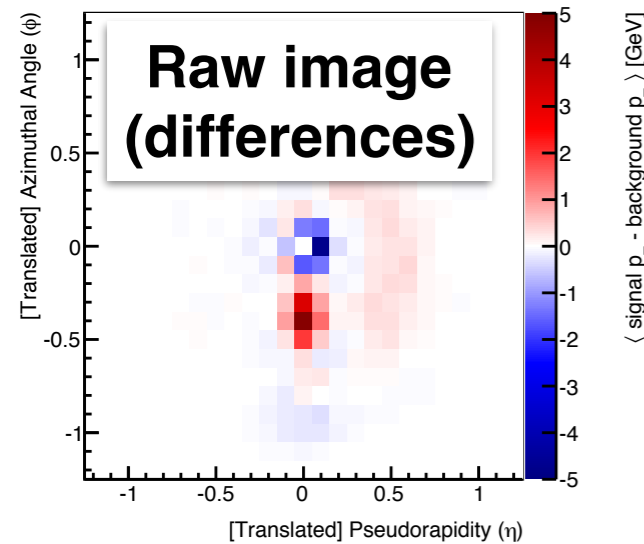
See also

- L. Almeida et al. 1501.05968*
- Baldi et al. 1603.09349*
- J. Barnard et al. 1609.00607*
- P. Komiske et al. 1612.01551*
- G. Kasieczka et al. 1701.08784*
- W. Bhimji et al. 1711.03573*

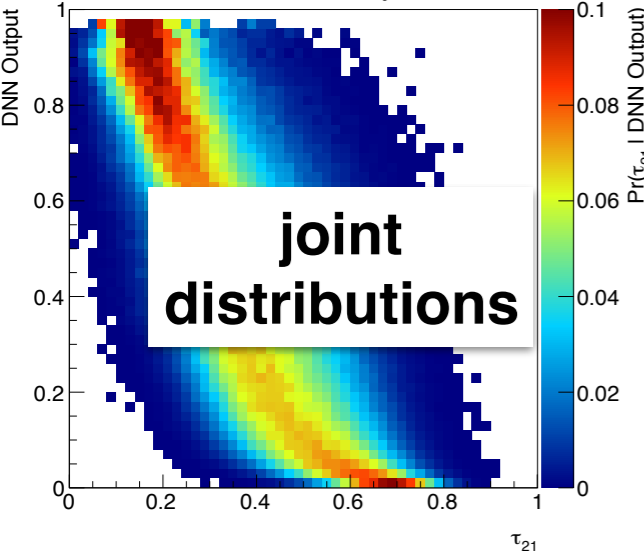
Learning about Learning

Opening the **box** is critical for improving robustness

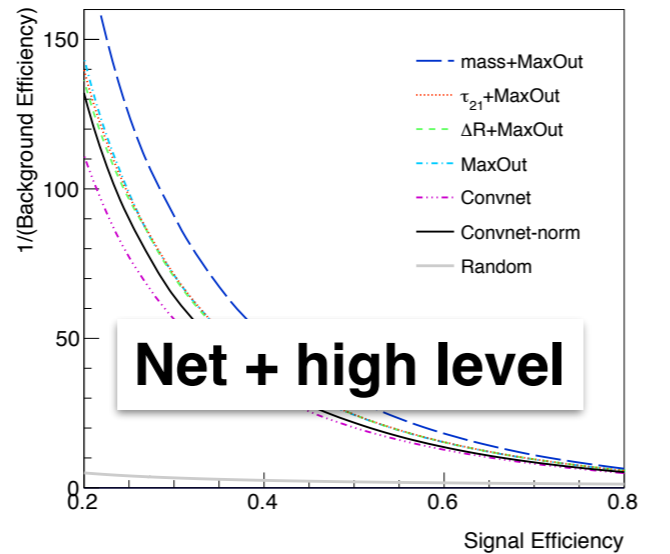
$250 < p_T/\text{GeV} < 260 \text{ GeV}$, $0.59 < \tau_{21} < 0.61$, $79 < \text{mass}/\text{GeV} < 81$
 $\sqrt{s} = 13 \text{ TeV}$, Pythia 8



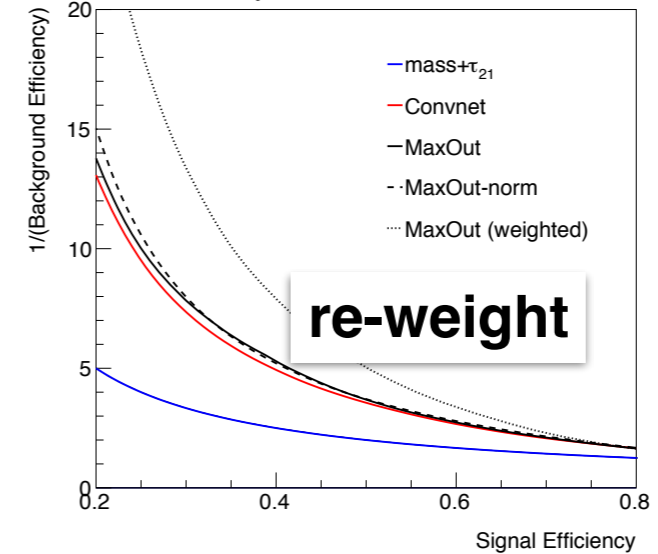
$250 < p_T/\text{GeV} < 300 \text{ GeV}$, $65 < \text{mass}/\text{GeV} < 95$
 QCD, $\sqrt{s} = 13 \text{ TeV}$, Pythia 8



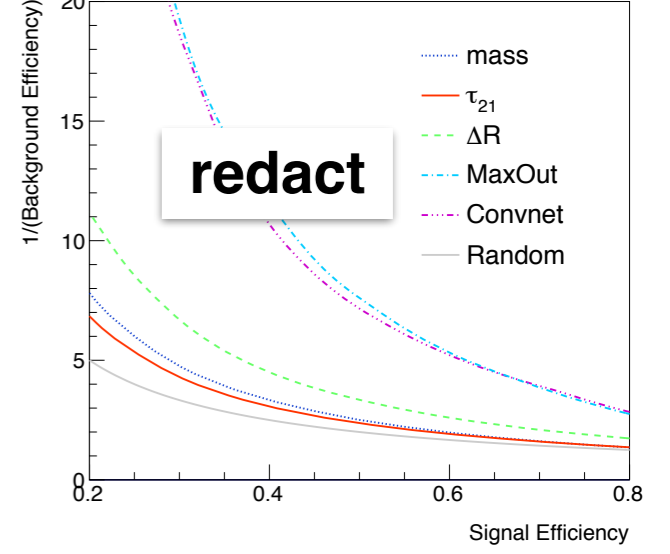
$240 < p_T/\text{GeV} < 260 \text{ GeV}$, $0.19 < \tau_{21} < 0.21$, $79 < \text{mass}/\text{GeV} < 81$
 $\sqrt{s} = 13 \text{ TeV}$, Pythia 8



$250 < p_T/\text{GeV} < 300 \text{ GeV}$, $0.2 < \tau_{21} < 0.8$, $65 < \text{mass}/\text{GeV} < 95$
 Pythia 8, $\sqrt{s} = 13 \text{ TeV}$

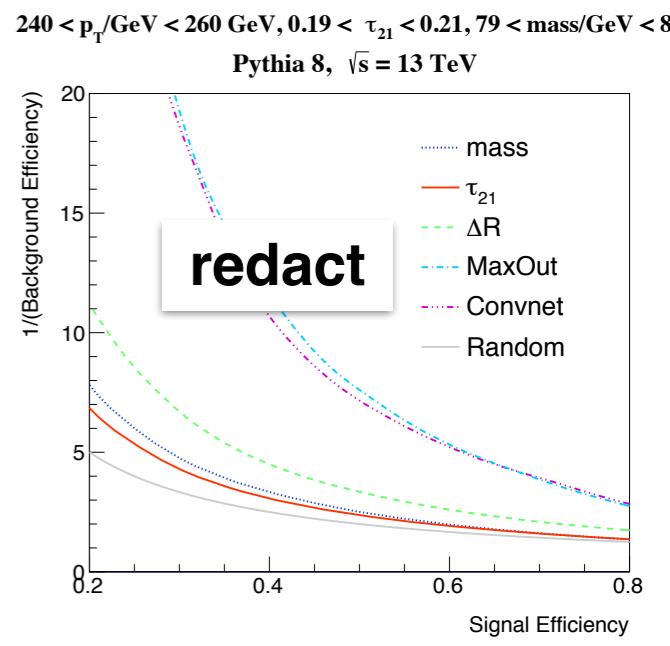
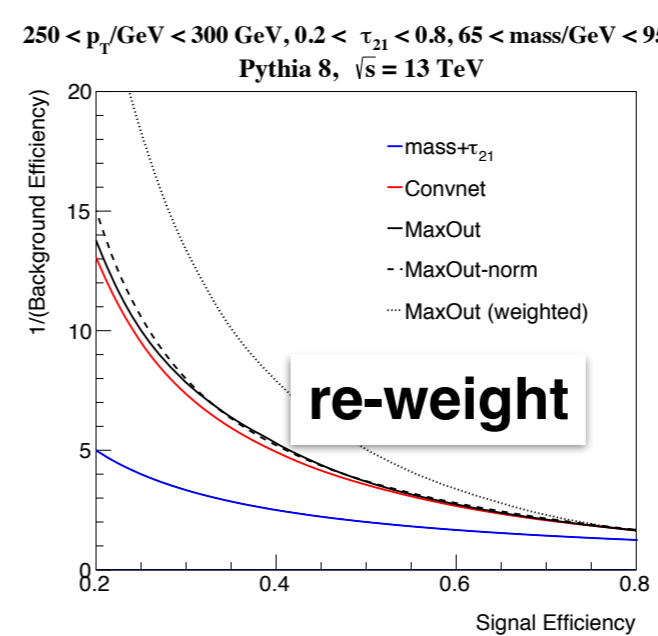
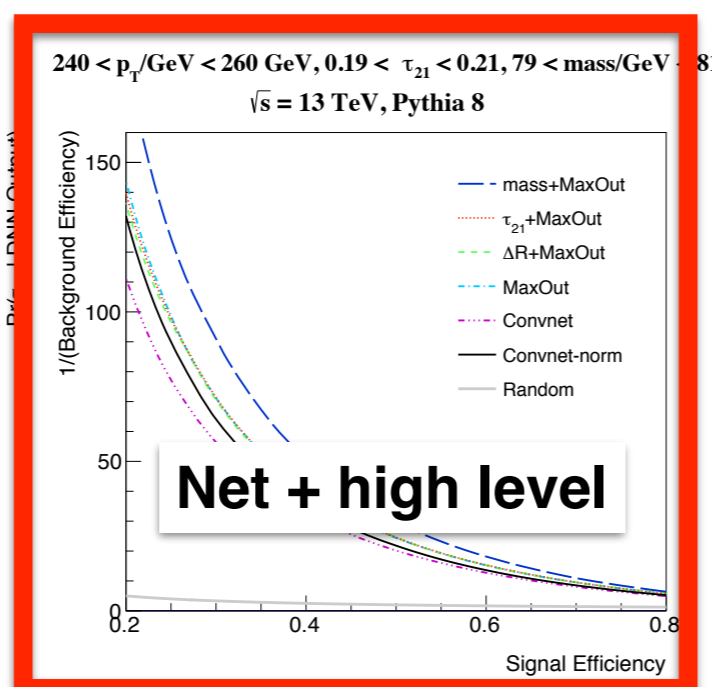
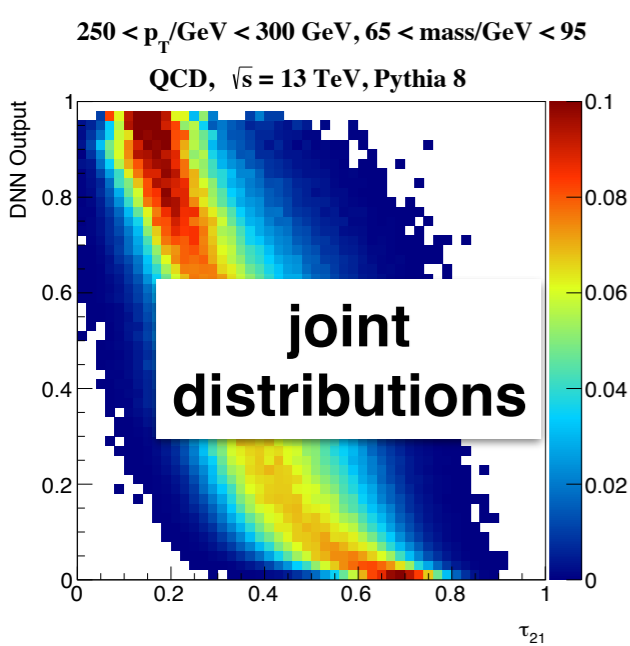
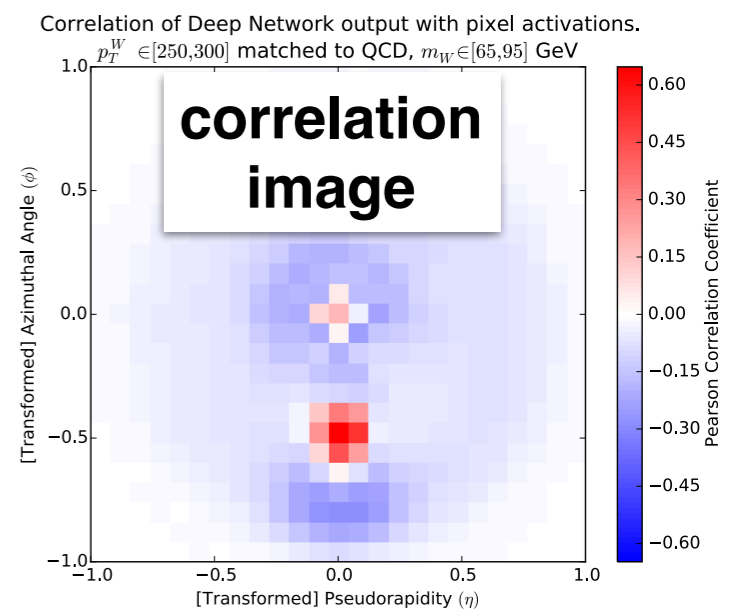
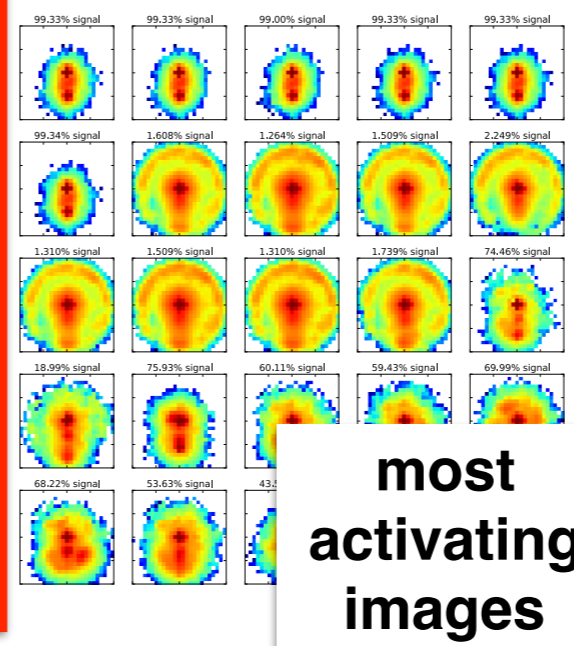
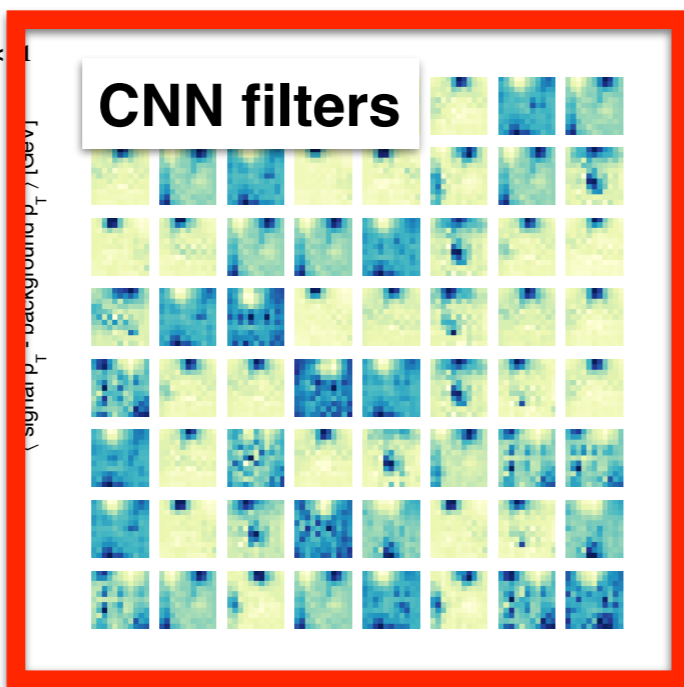
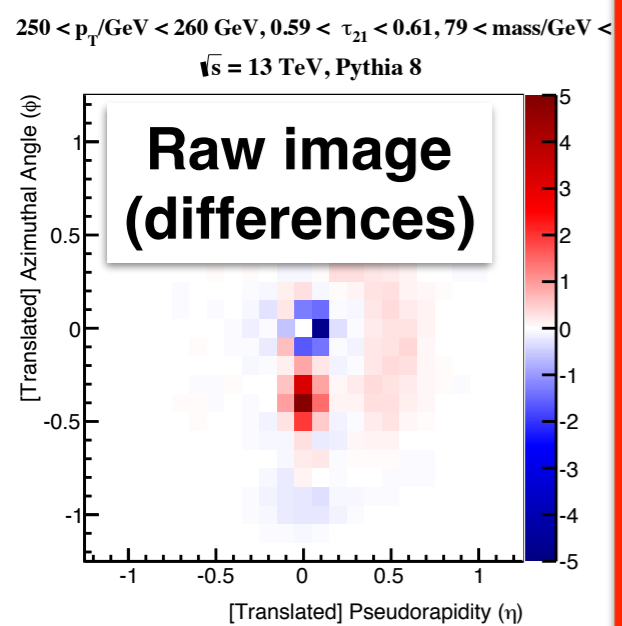


$240 < p_T/\text{GeV} < 260 \text{ GeV}$, $0.19 < \tau_{21} < 0.21$, $79 < \text{mass}/\text{GeV} < 81$
 Pythia 8, $\sqrt{s} = 13 \text{ TeV}$



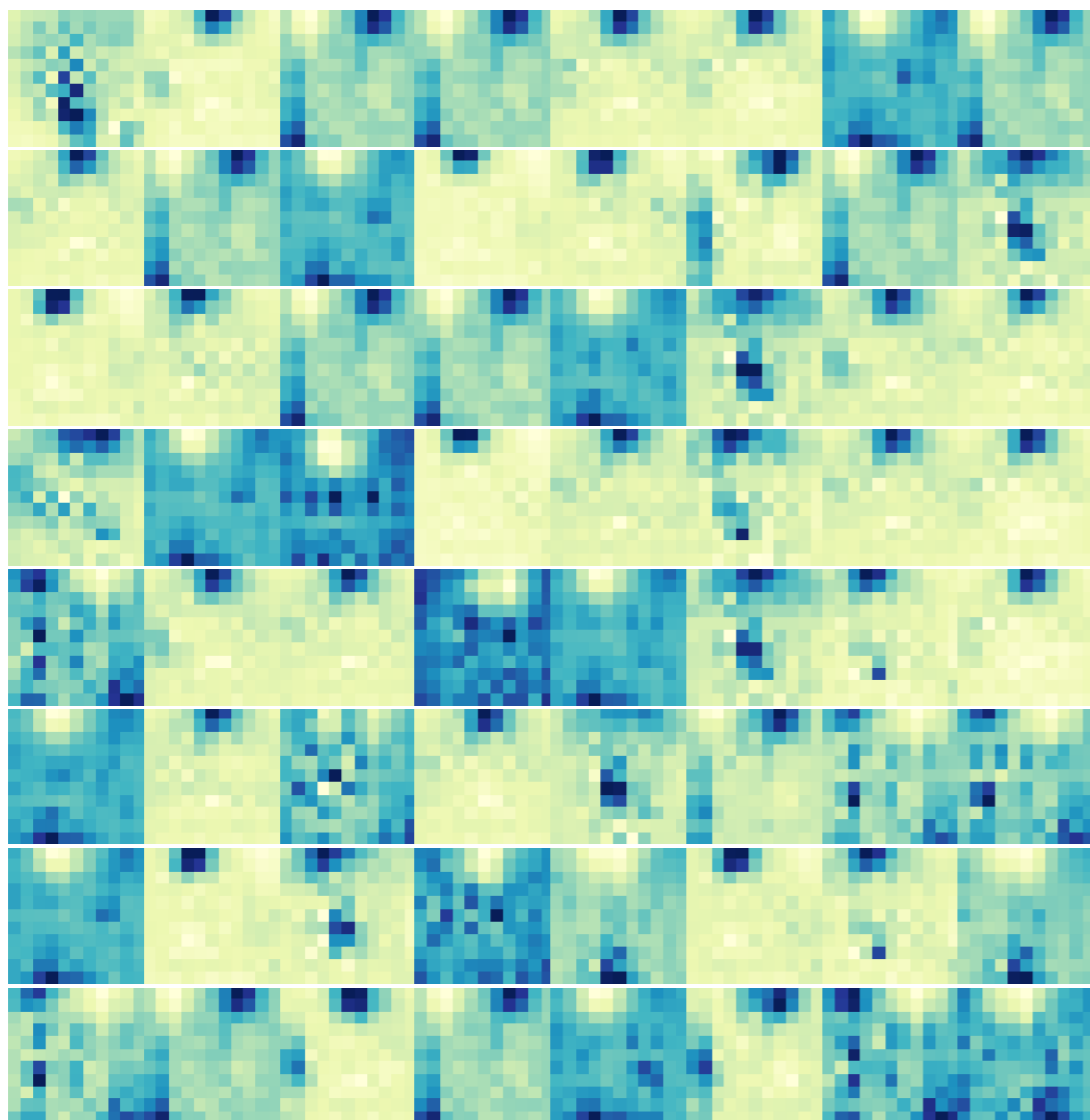
Learning about Learning

Opening the **box** is critical for improving robustness



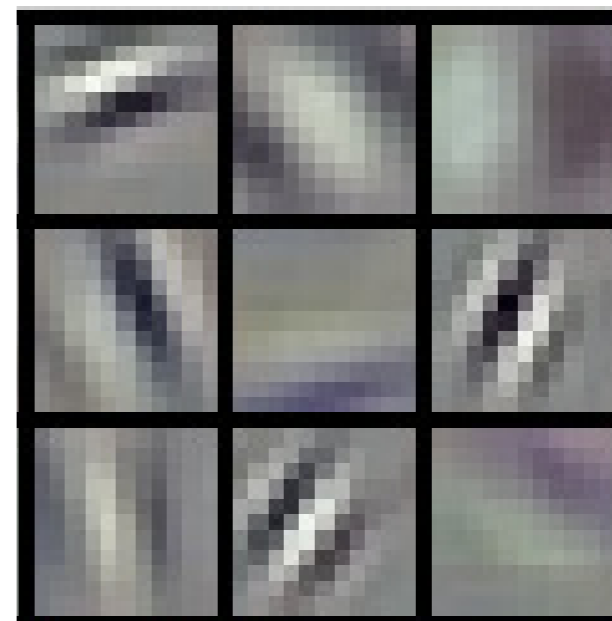
Convolutional Filters

Filters are images! Can visualize 'higher-level features' learned by the network



Jet Images

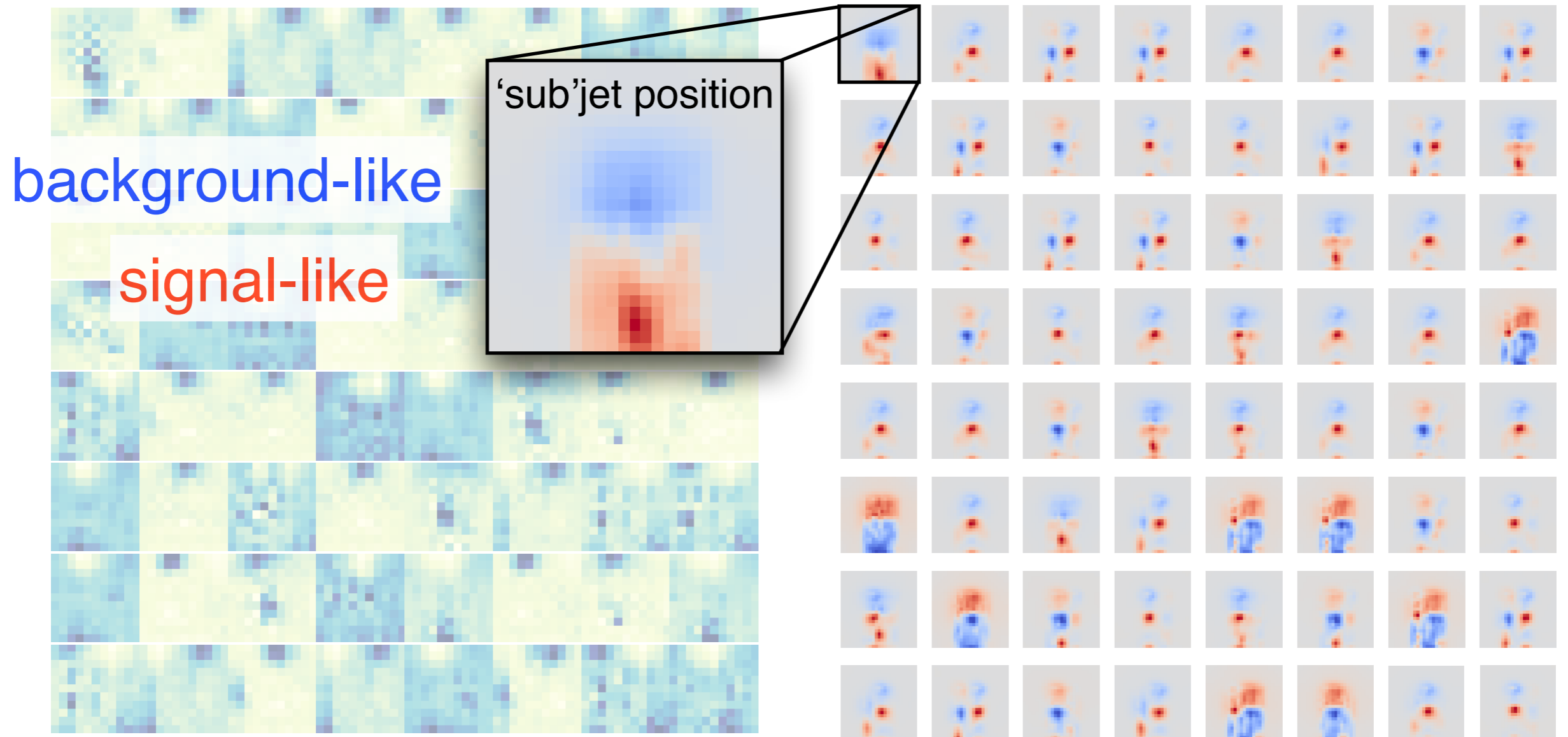
learned edge detection



“Natural” Images

Convolutional Filters

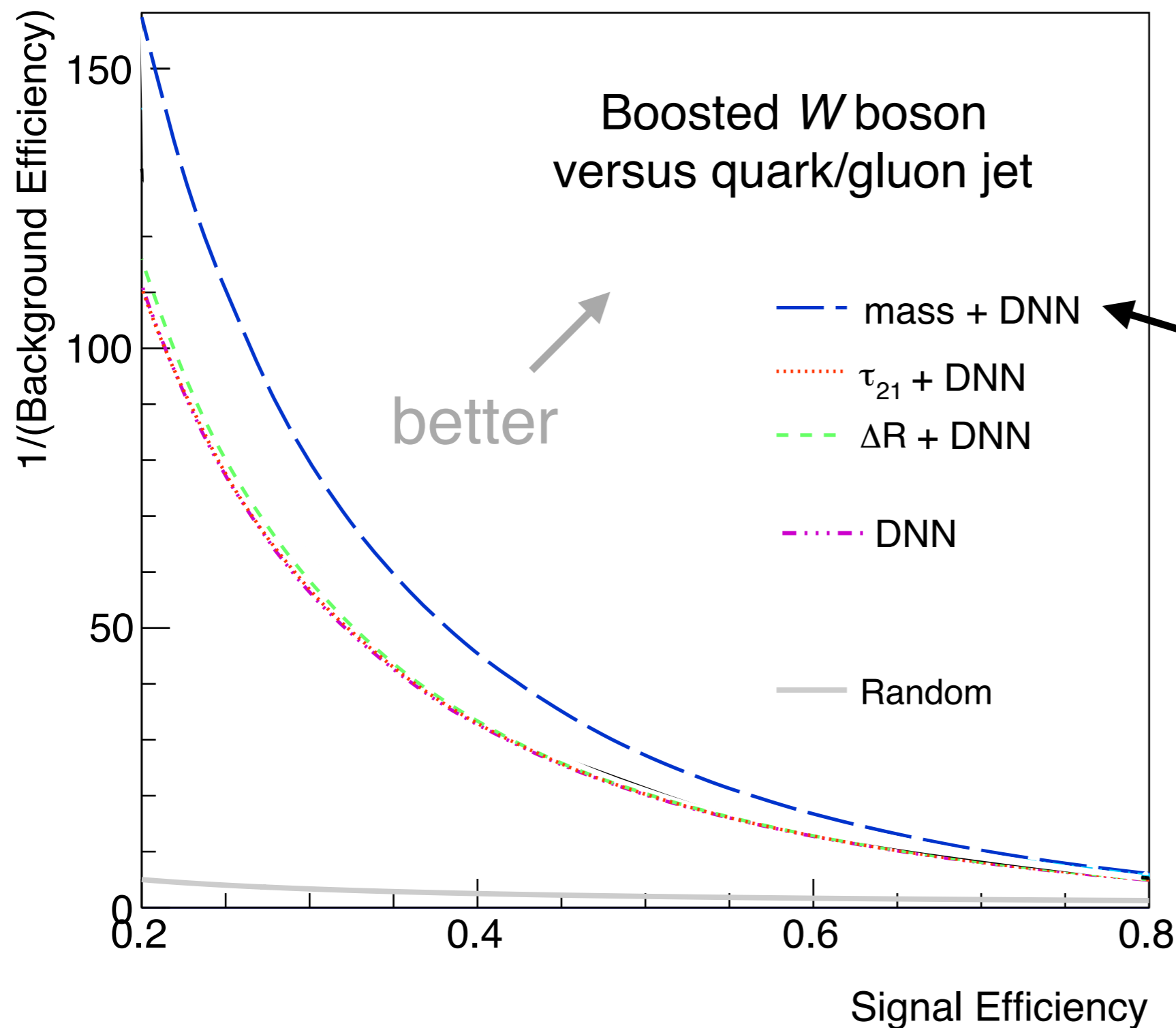
Filters are images! Can visualize 'higher-level features' learned by the network



Jet Images Layer 1 Filters

Filters convolved with
signal - background

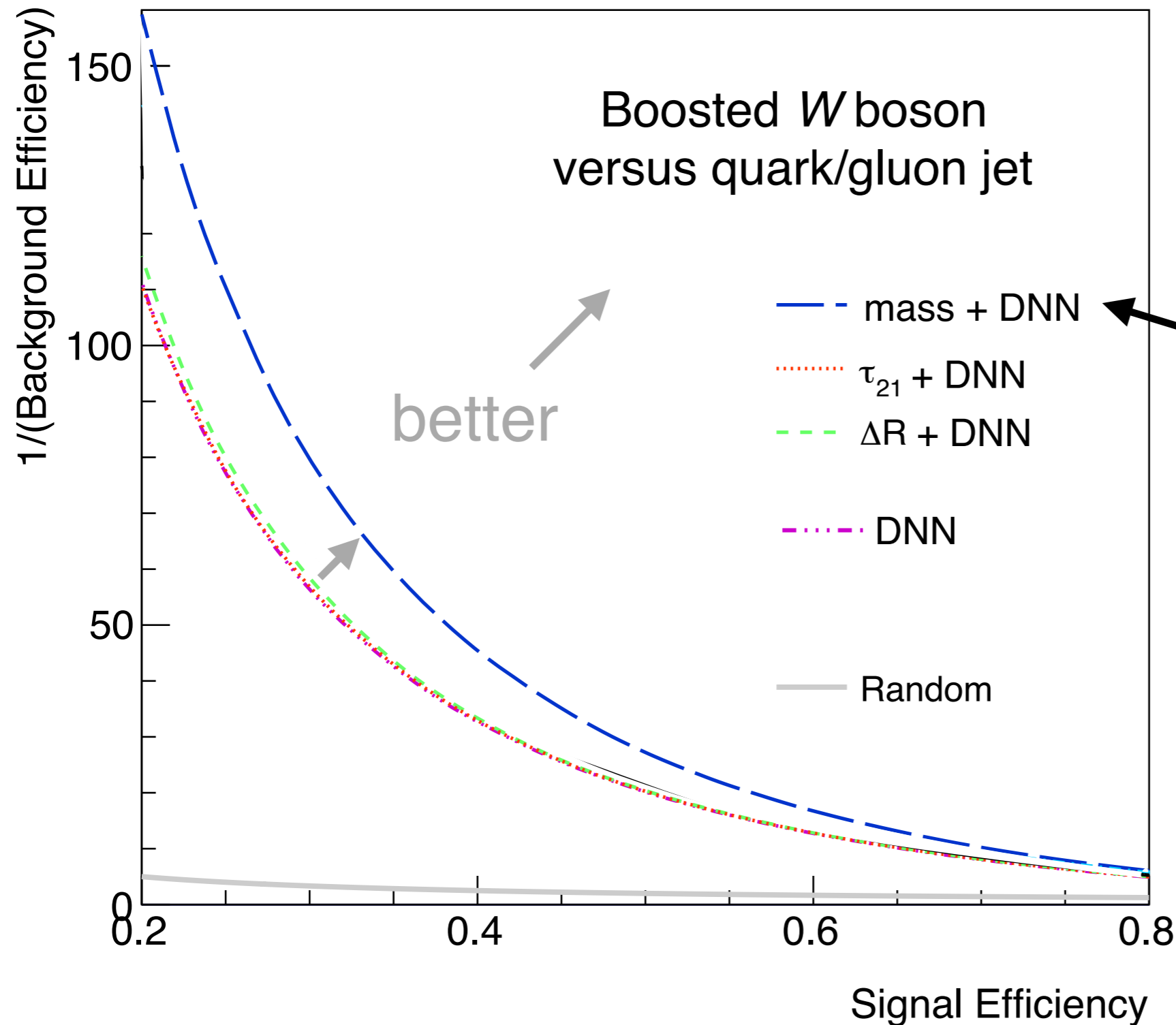
Is there more to learn that we know about?



Idea: explicitly
combine NN with
a known feature

Has learned
image mass?

Is there more to learn that we know about?

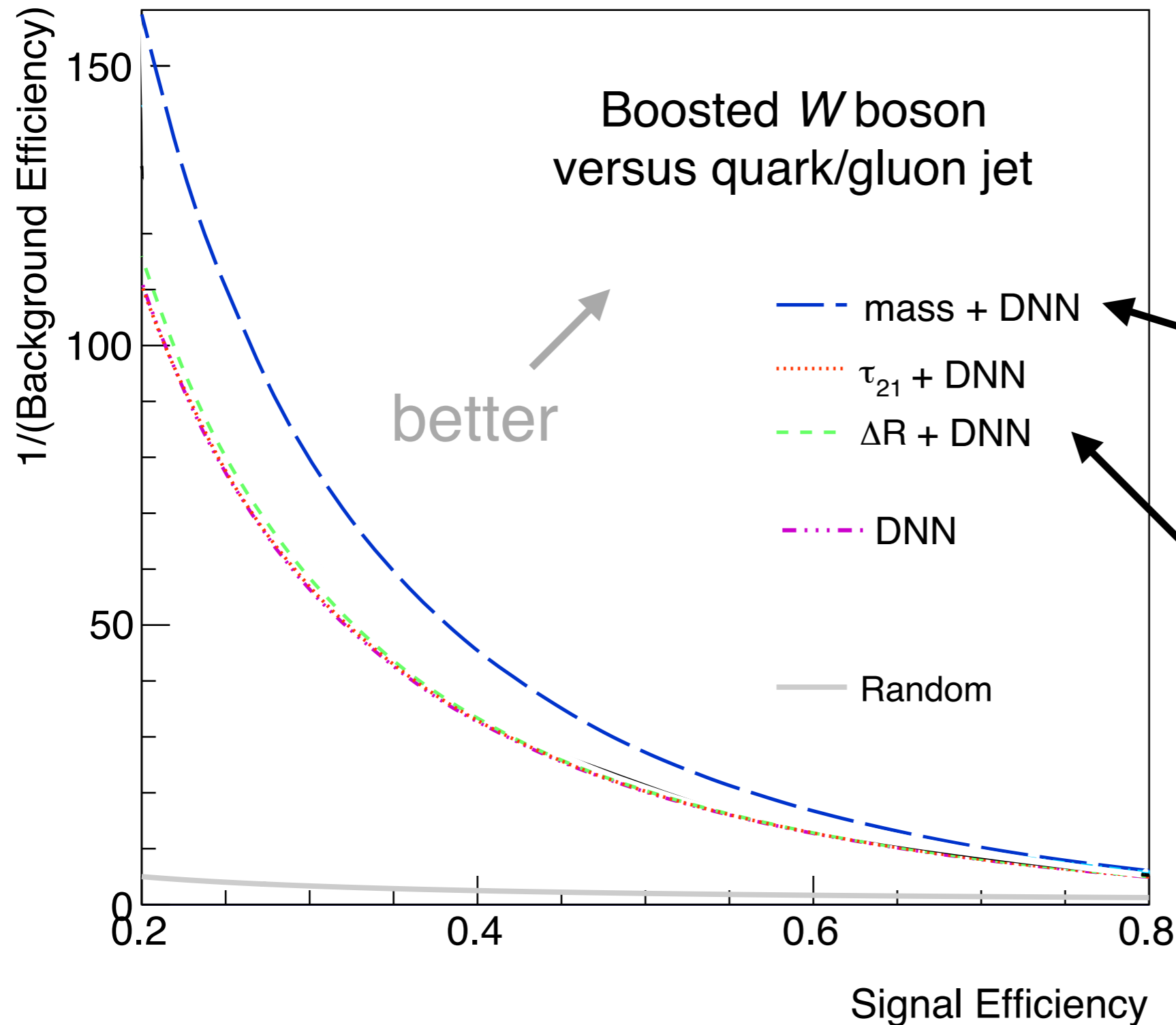


Idea: explicitly combine NN with a known feature

Has learned image mass?

not fully!

Is there more to learn that we know about?



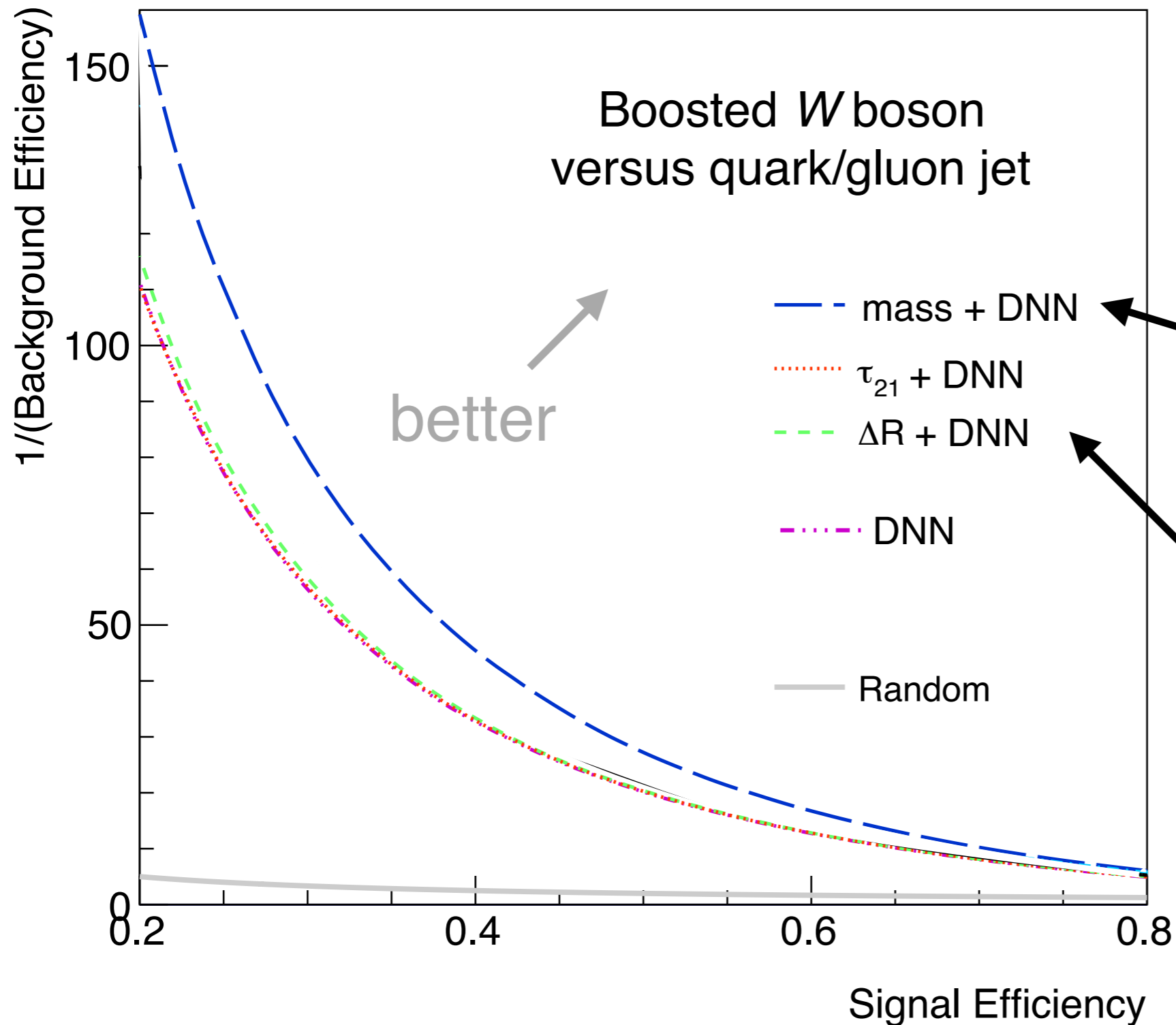
Idea: explicitly
combine NN with
a known feature

Has learned
image mass?

not fully!

What about the
distance between
subjects?

Is there more to learn that we know about?



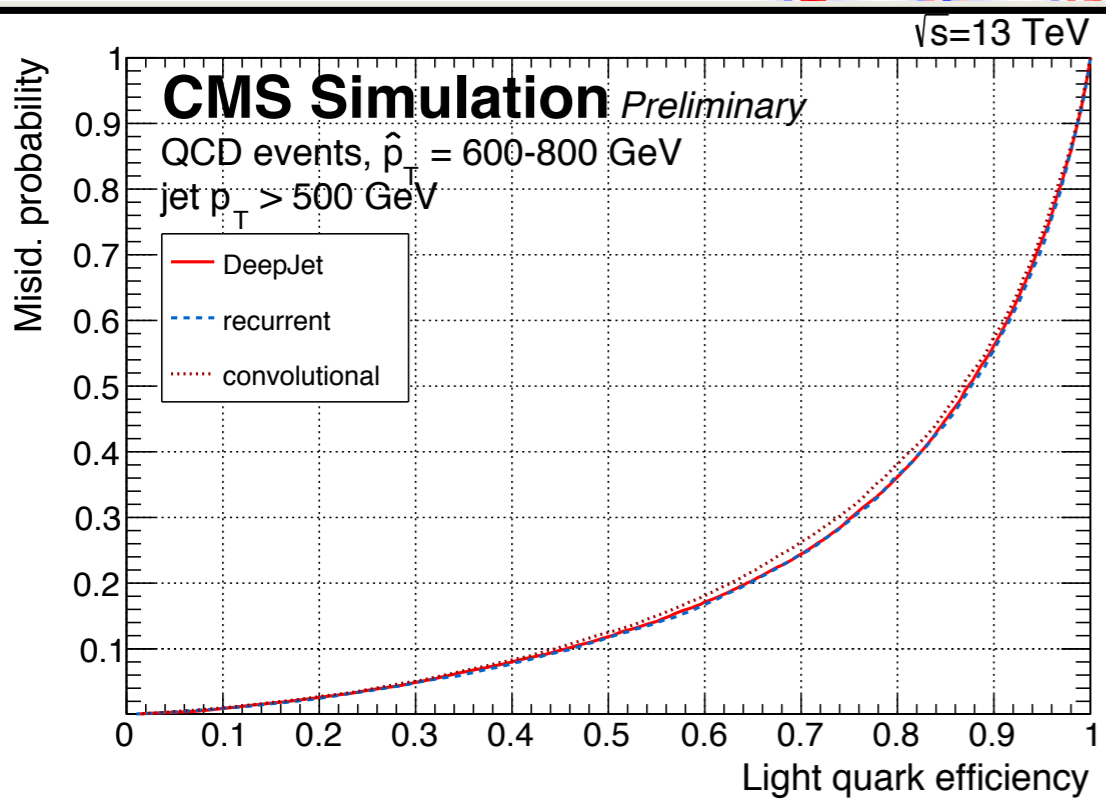
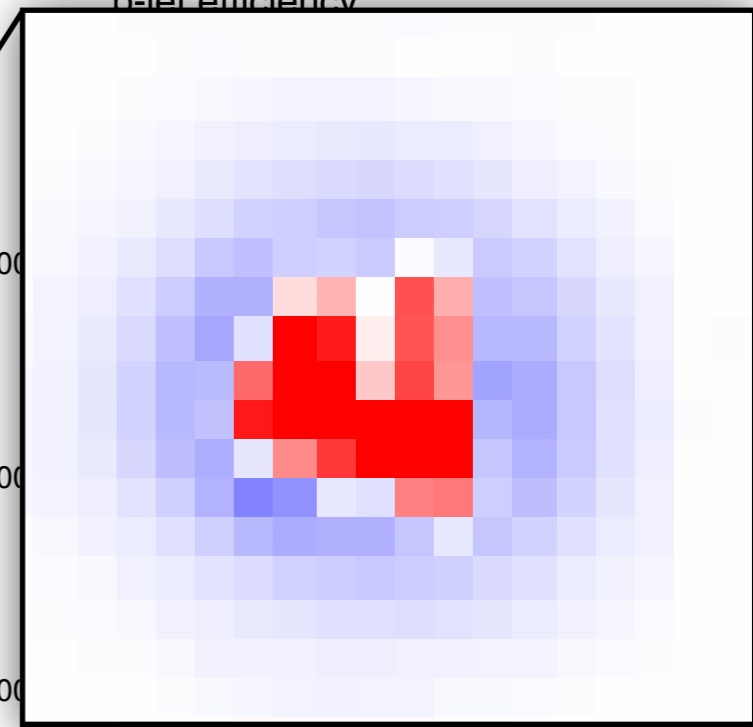
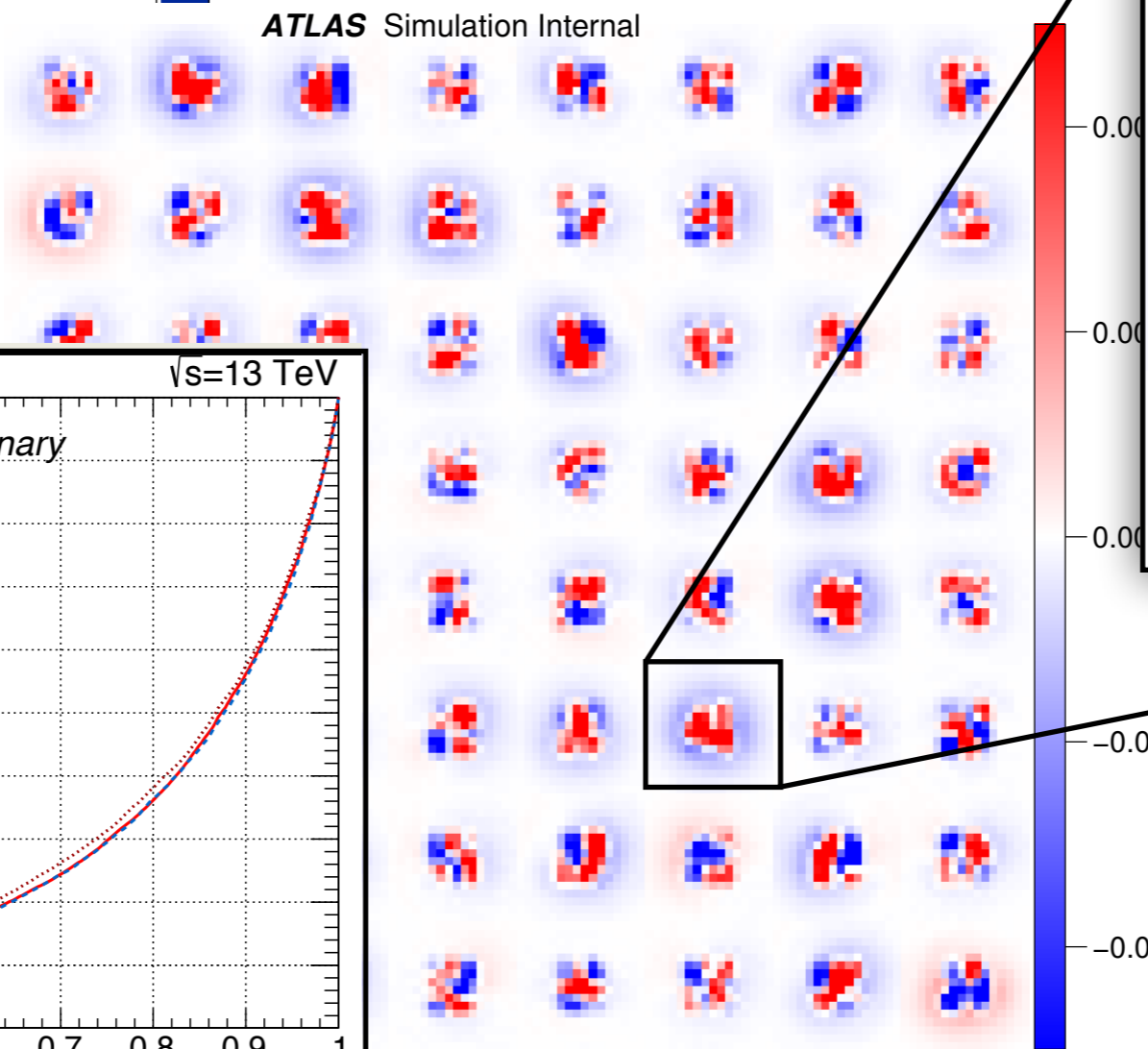
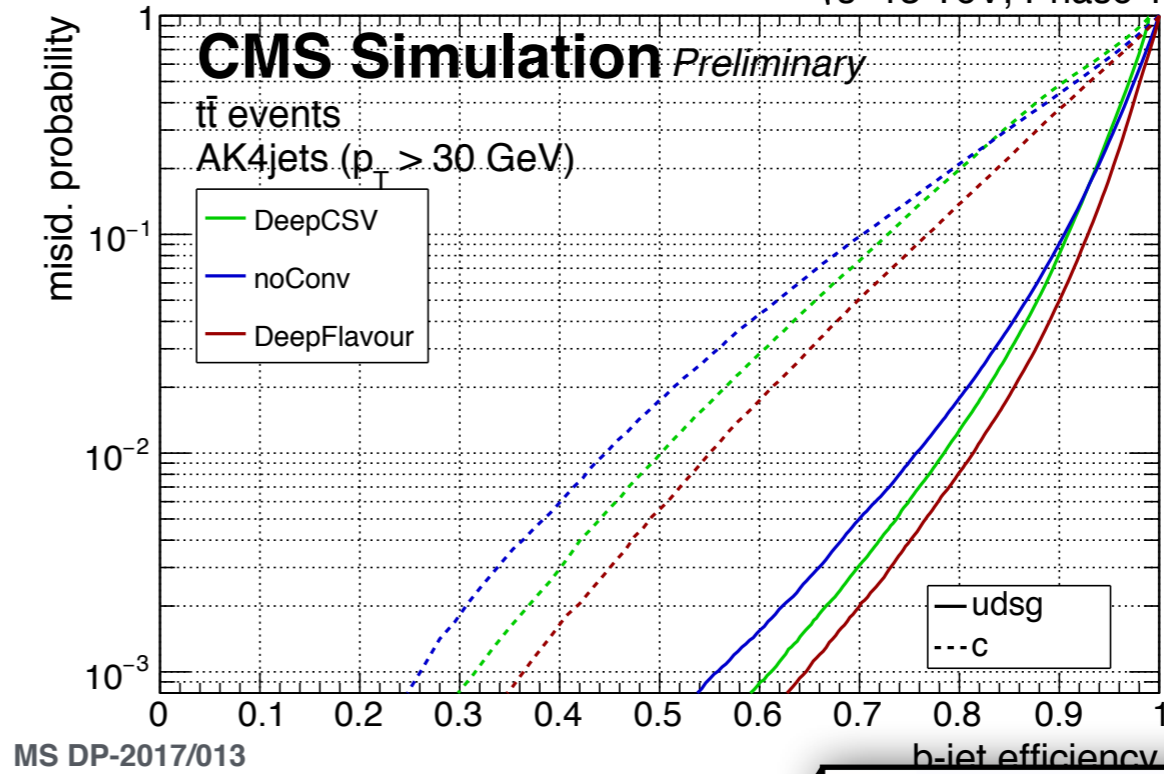
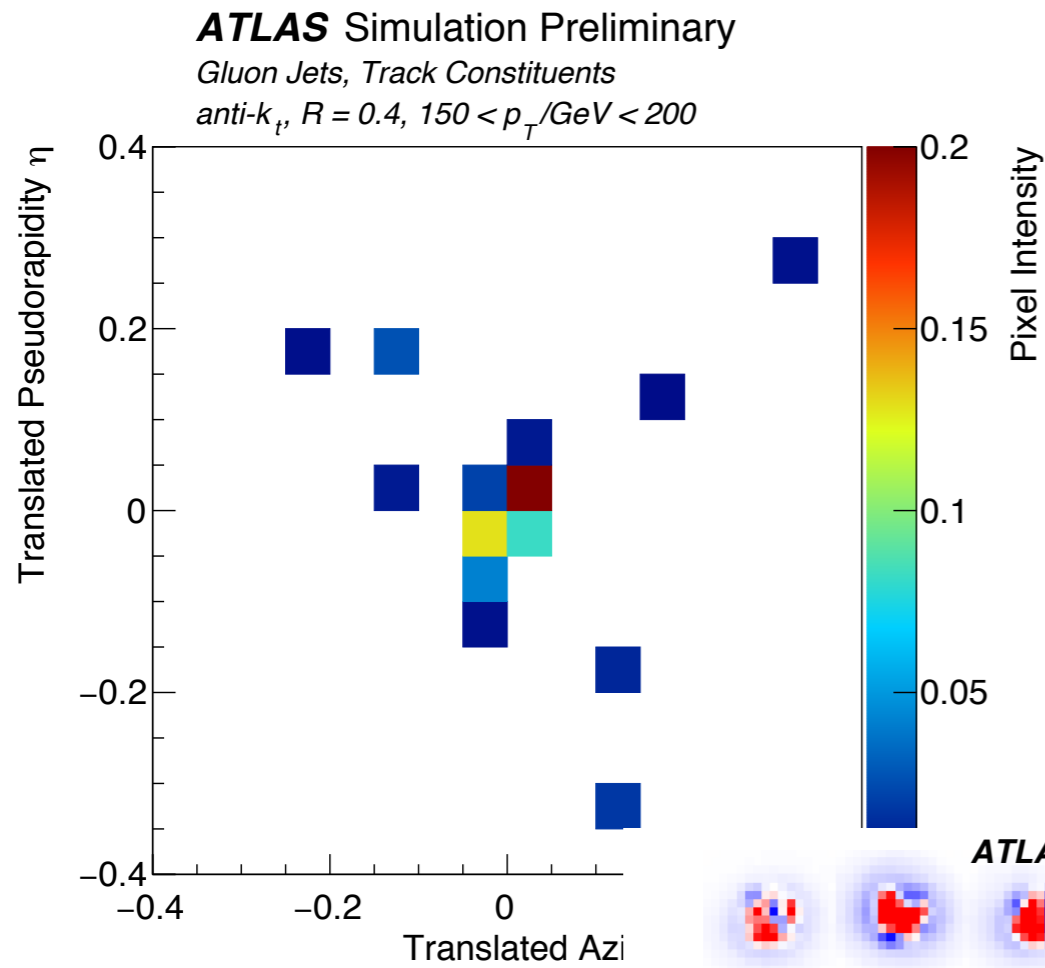
Idea: explicitly combine NN with a known feature

Has learned image mass?

not fully!

What about the distance between subjects?

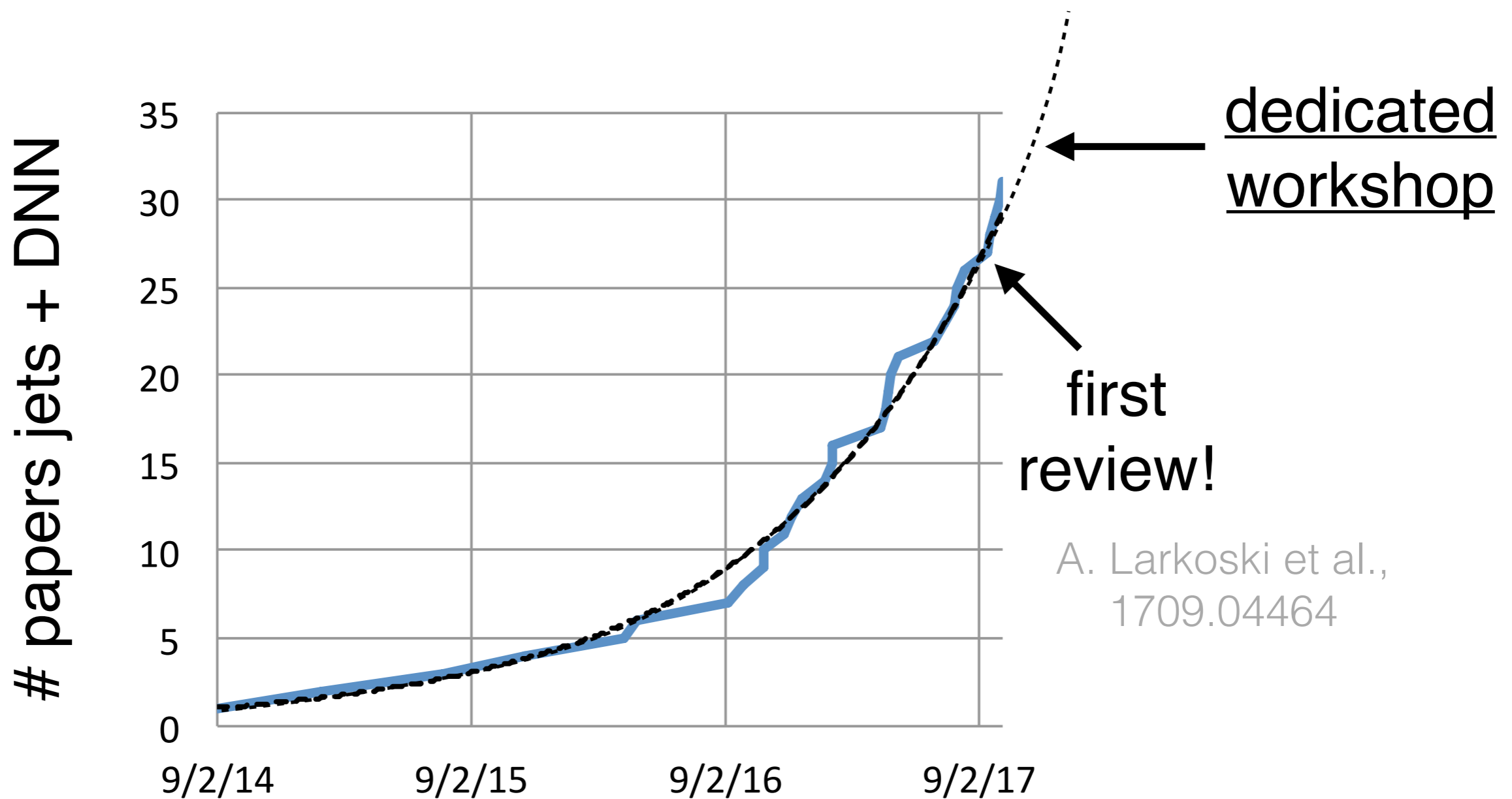
yes!



See also
 Markus'
 seminar

Exciting New Directions

So far only scratches the surface
...this is a very active field of research!

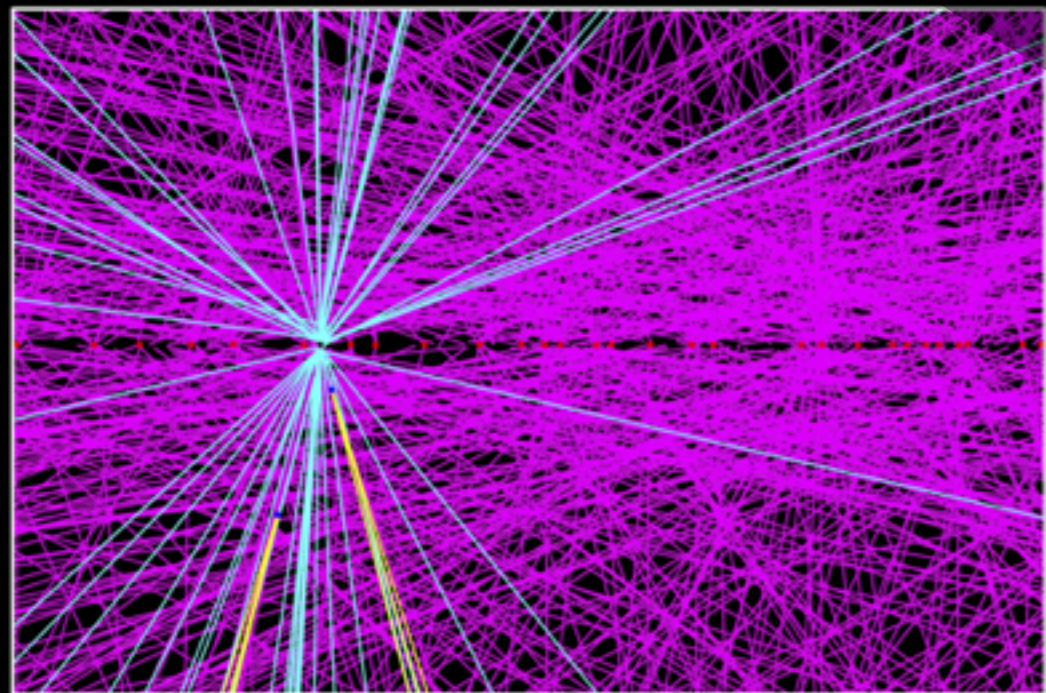


Exciting New Directions I: Removing Noise

pp collisions at the LHC
don't happen one at a time!

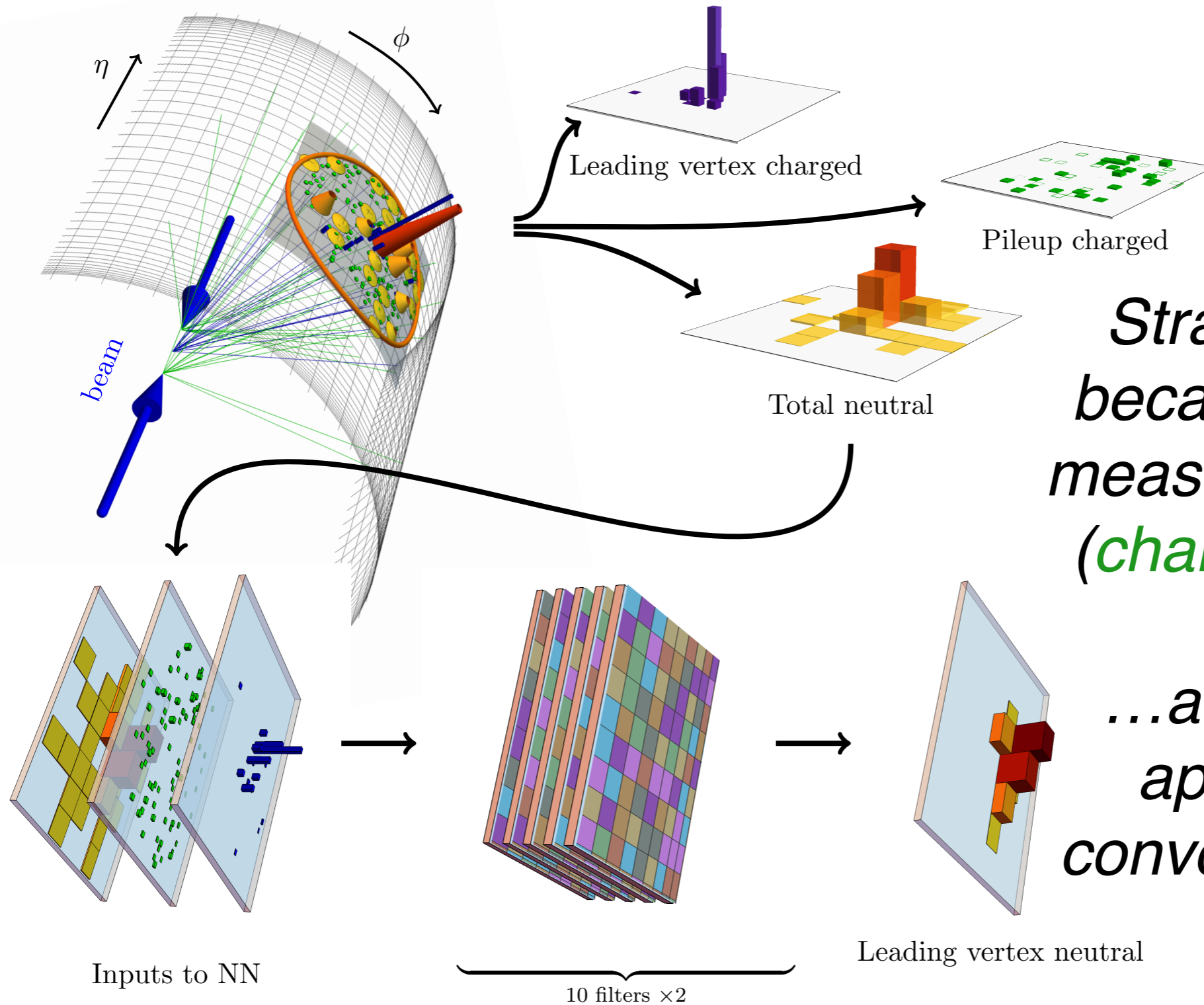


the extra collisions are called **pileup**
and add soft radiation on top of our jets



this is akin to image
de-noising - we can
use ML for that!

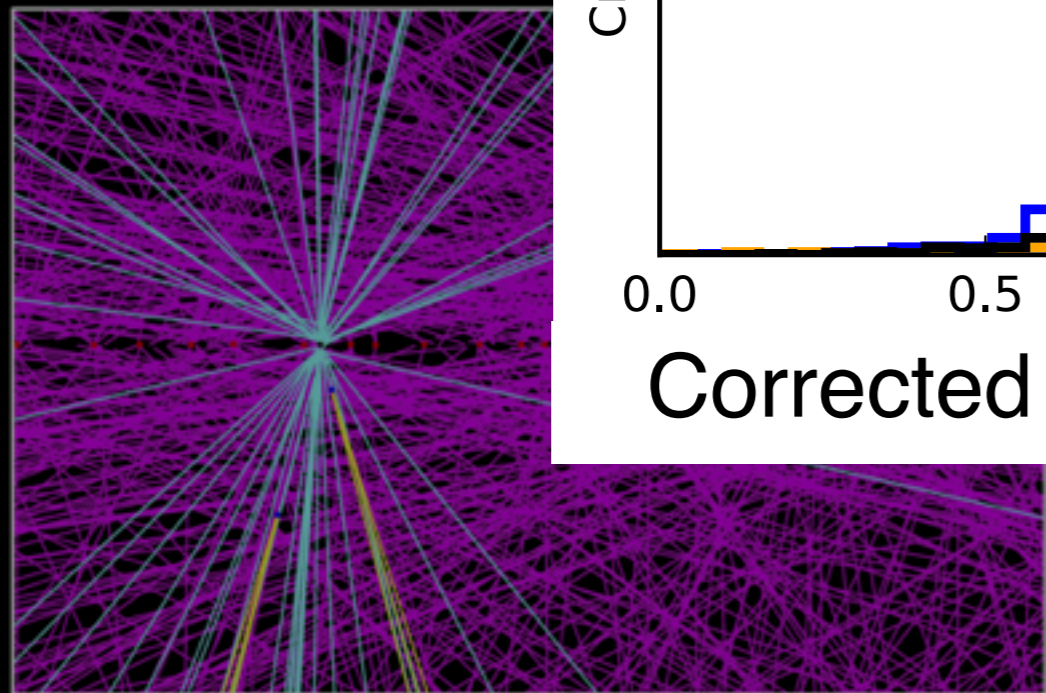
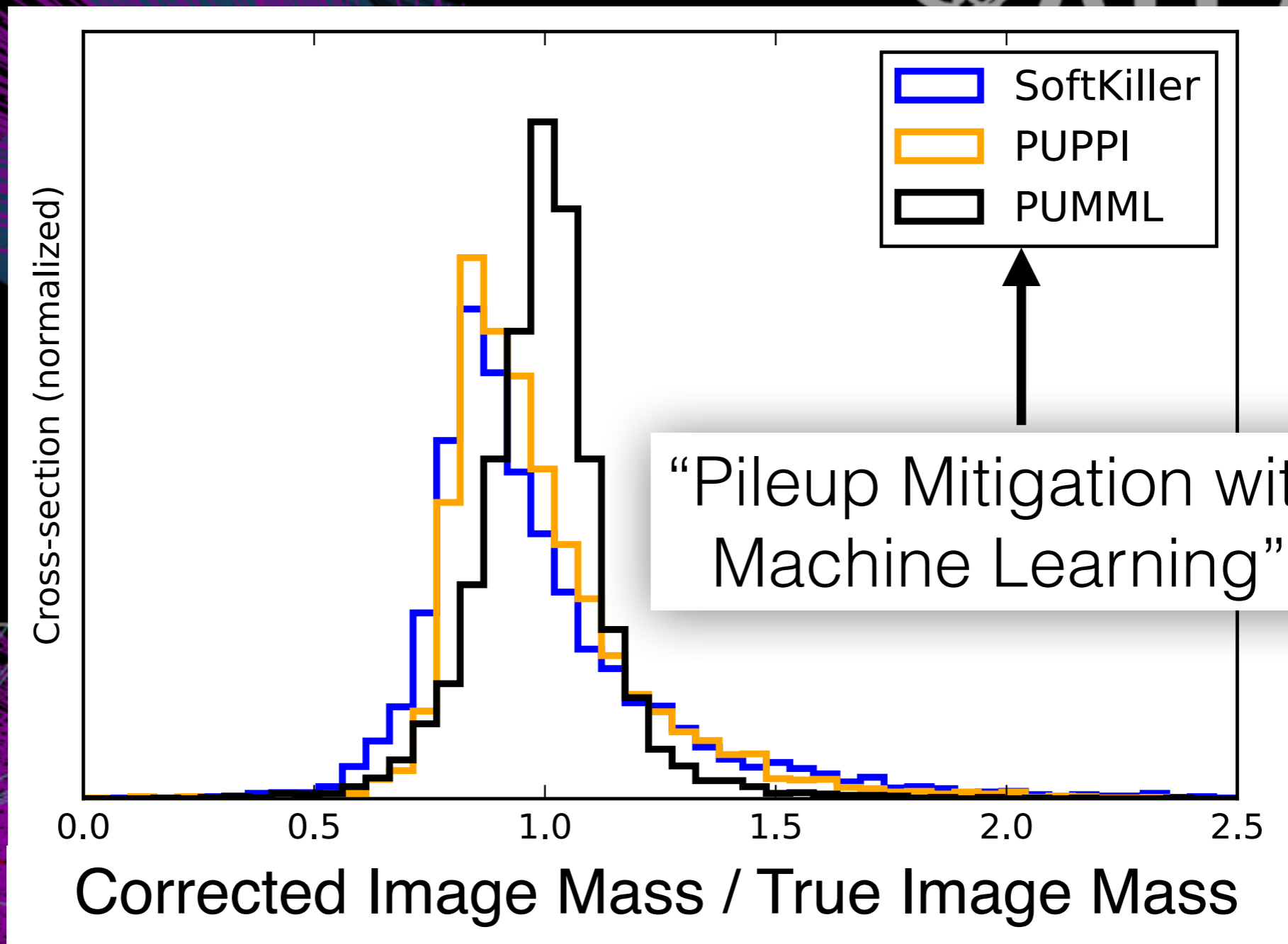
Exciting New Directions I: Removing Noise



*Strange noise
because we can
measure $\sim 2/3$ of it
(charged pileup)*

*...also a natural
application of
convolutional NNs!*

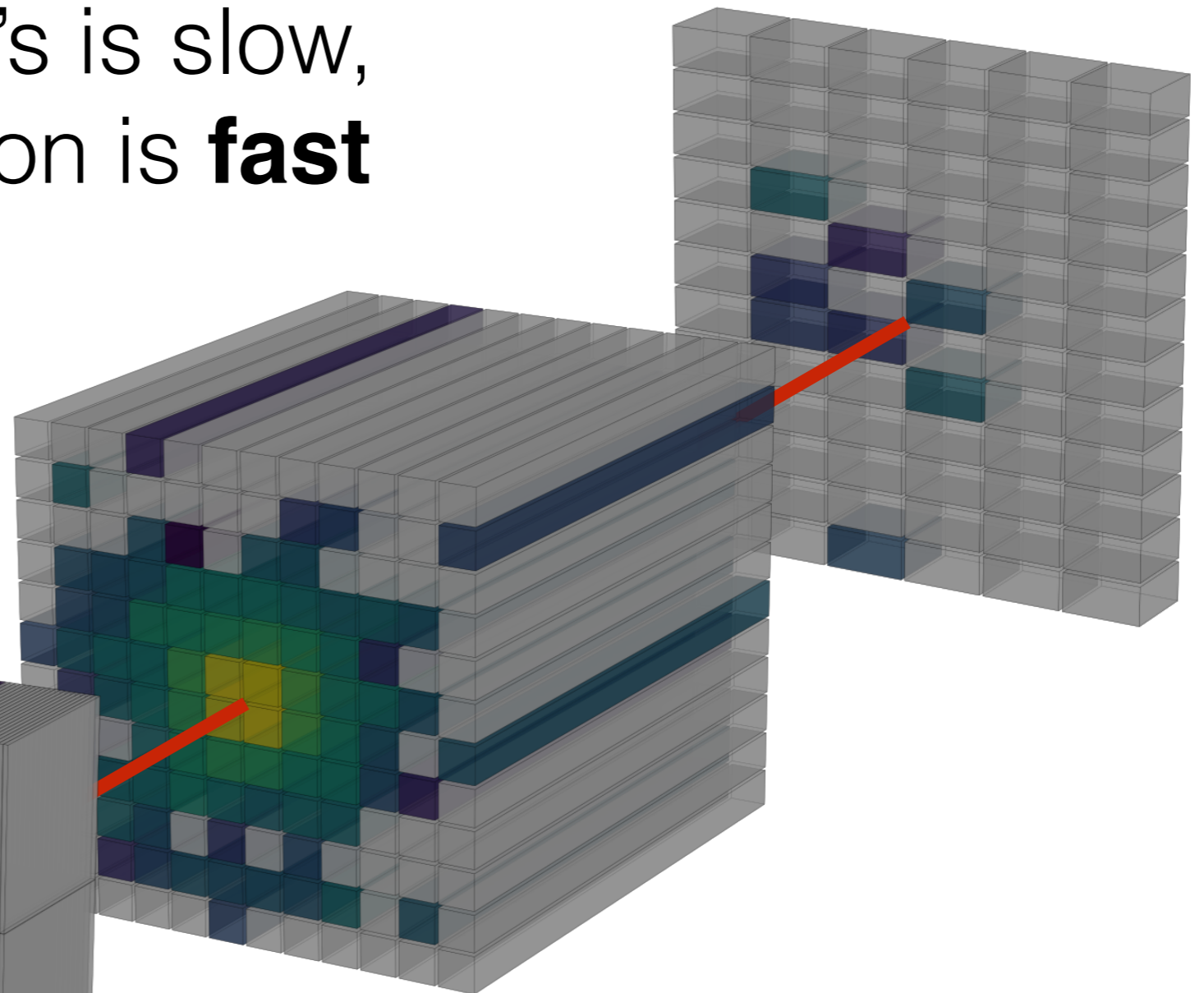
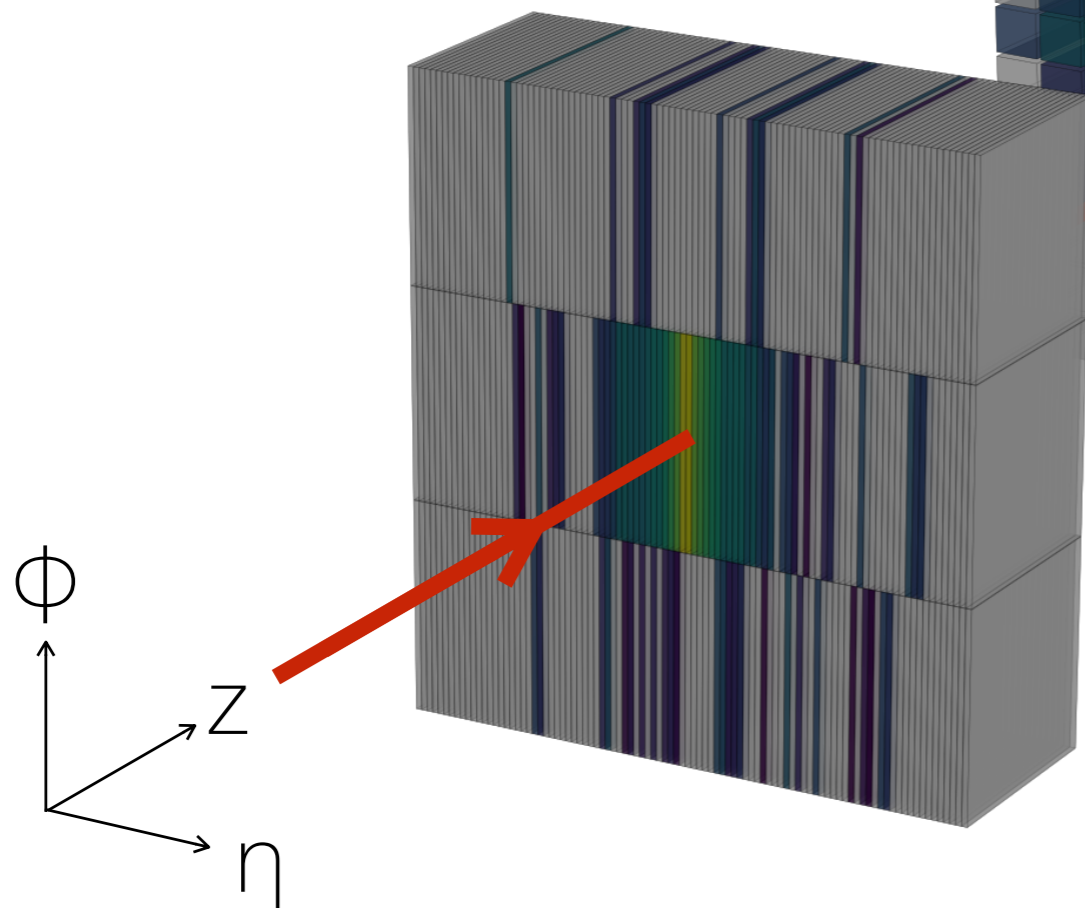
Exciting New Directions I: Removing Noise



Exciting New Directions II: Simulation NN

Training NN's is slow,
but evaluation is **fast**

Physics-based
simulations of
jets are **slow**



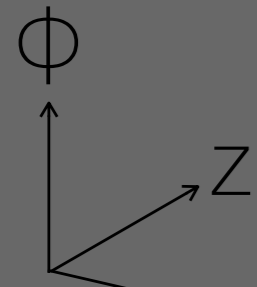
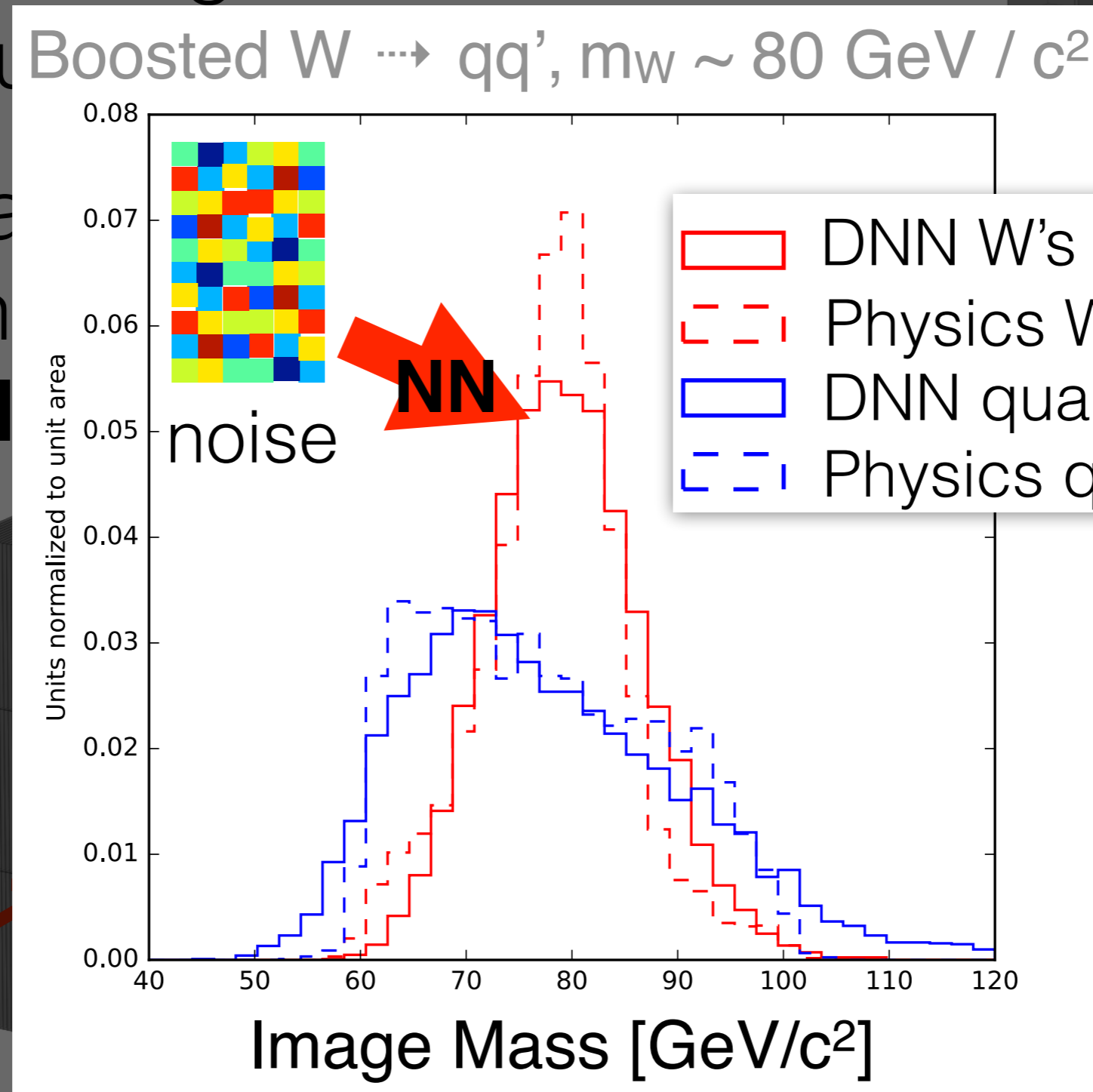
What if we can learn to
simulate jets with a NN?

Exciting New Directions II: Simulation NN

Training NN's is slow,

but Boosted $W \rightarrow qq'$, $m_W \sim 80 \text{ GeV} / c^2$

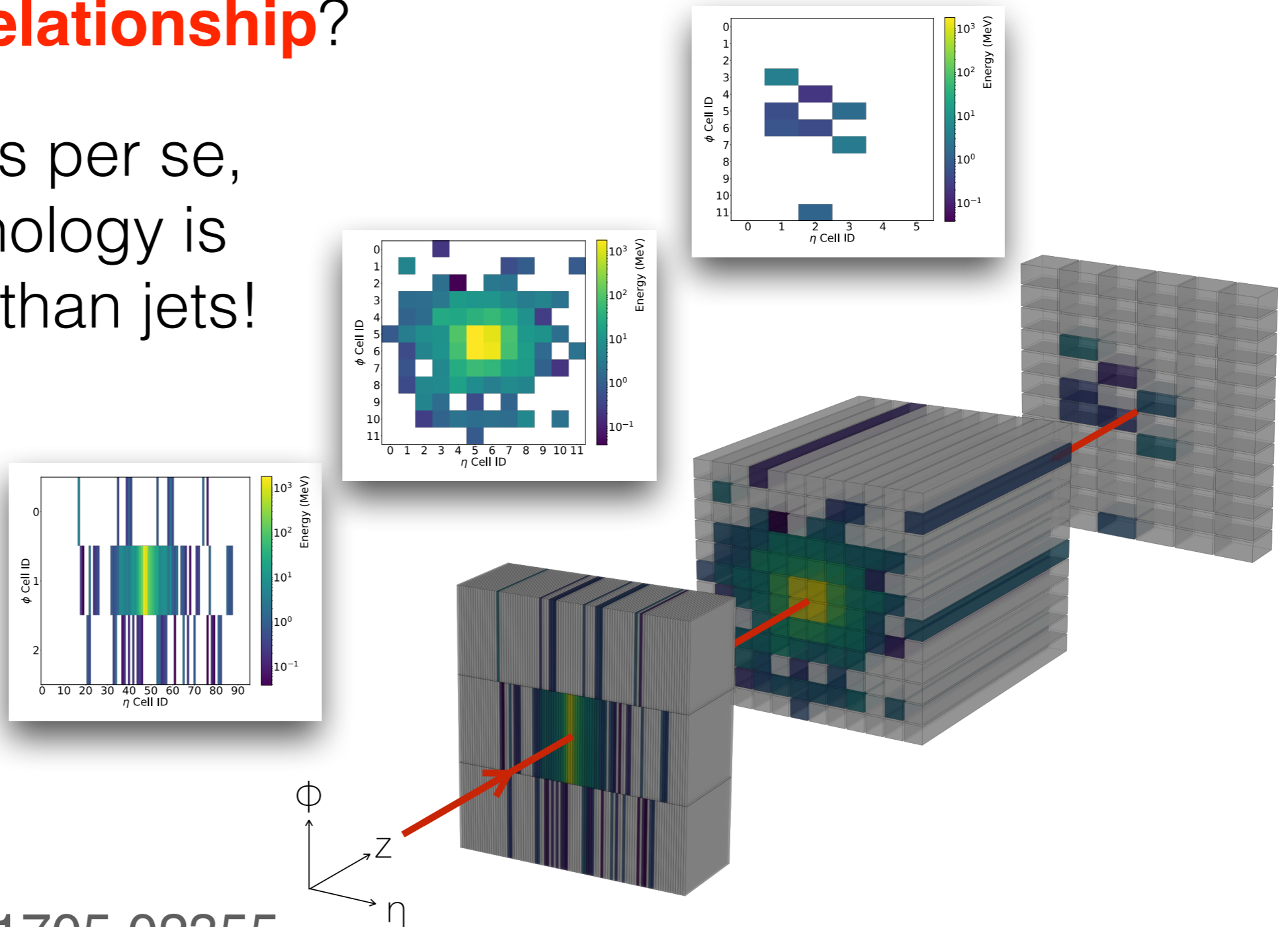
Physics-based simulation jets are slow



+ More Layers for Generation

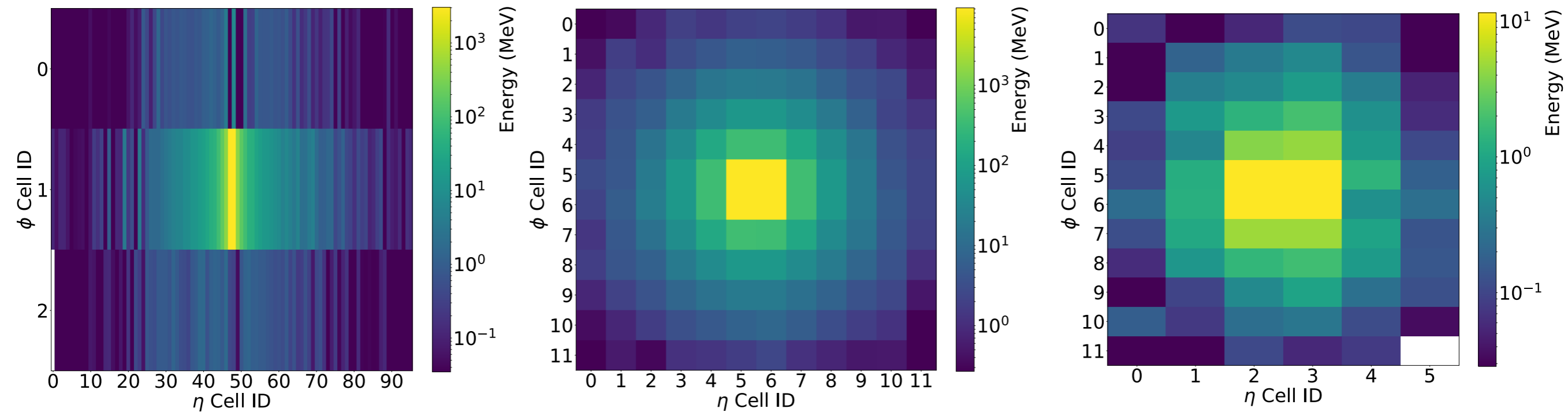
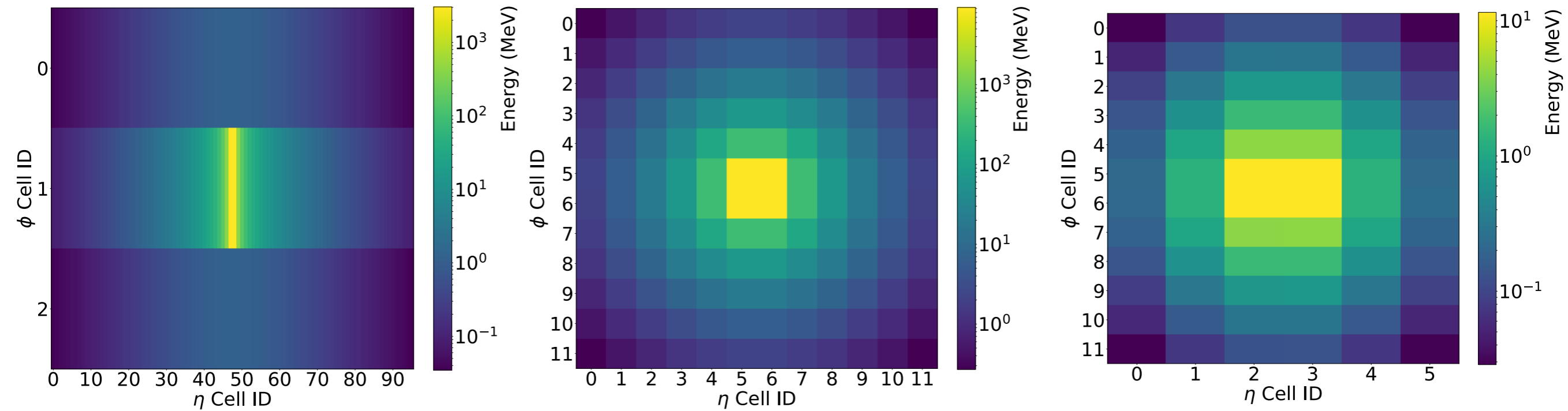
What about **multiple layers** with **non-uniform granularity** and a **causal relationship**?

Not jet images per se,
but the technology is
more general than jets!



Average Images

Geant4



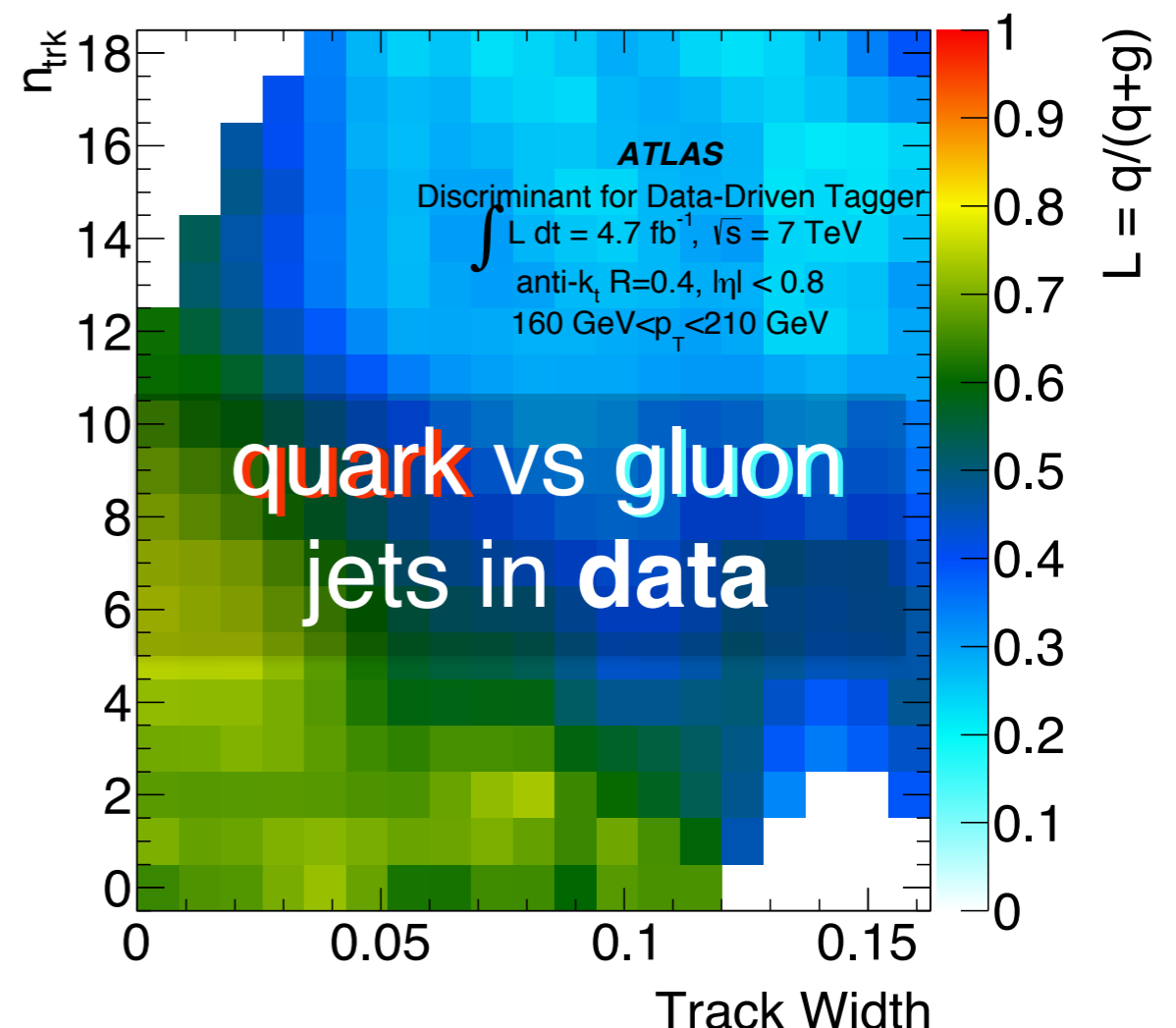
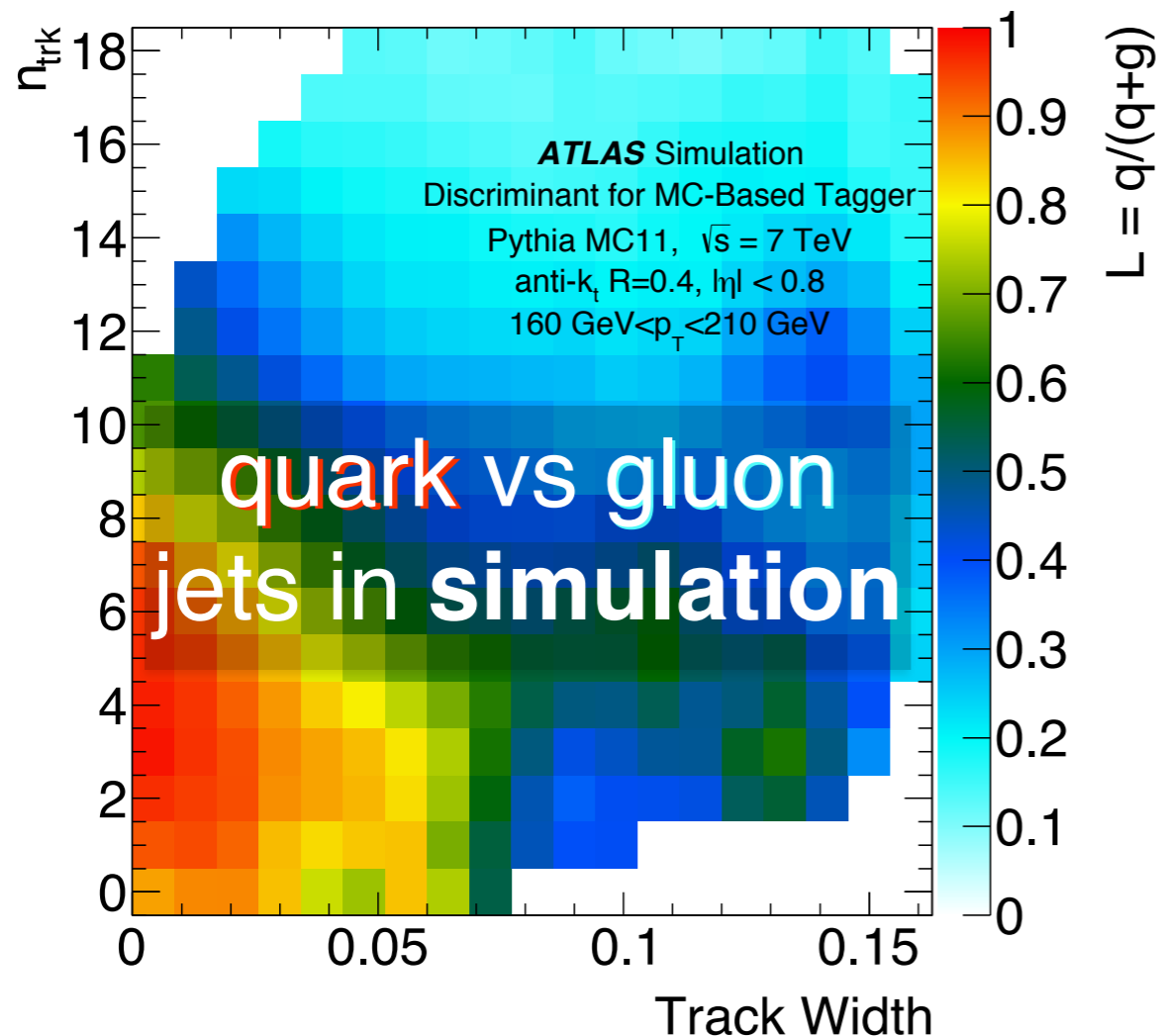
M. Paganini et al., 1705.02355

Generation Method	Hardware	Batch Size	milliseconds/shower
GEANT4	CPU	N/A	1772 ←
CALOGAN	CPU	1	13.1
		10	5.11
		128	2.19
		1024	2.03
	GPU	1	14.5
		4	3.68
		128	0.021
		512	0.014
		1024	0.012 →

See also [S. Vallecorsa et al. \(GeantV\)](#), [C. Guthrie et al. \(NYU\)](#), [W. Wei et al. \(LCD dataset group\)](#), [D. Salamani et al. \(Geneva\)](#), [D. Rousseau et al. \(Orsay\)](#), [L. de Oliveira et al. \(Berkeley\)](#)

Where next III: Learning directly from data

For supervised learning, we depend on labels
labels usually come from simulation

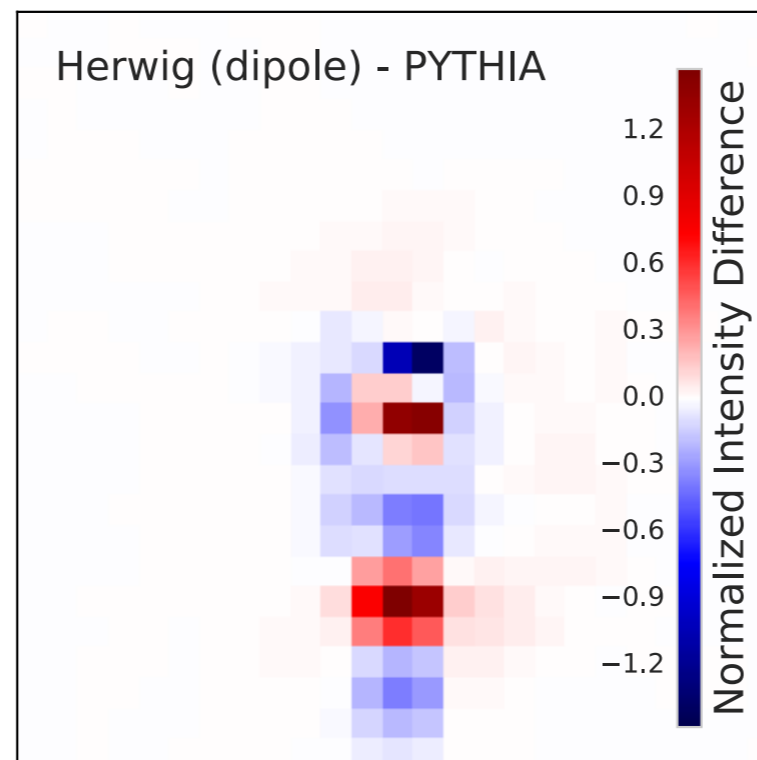
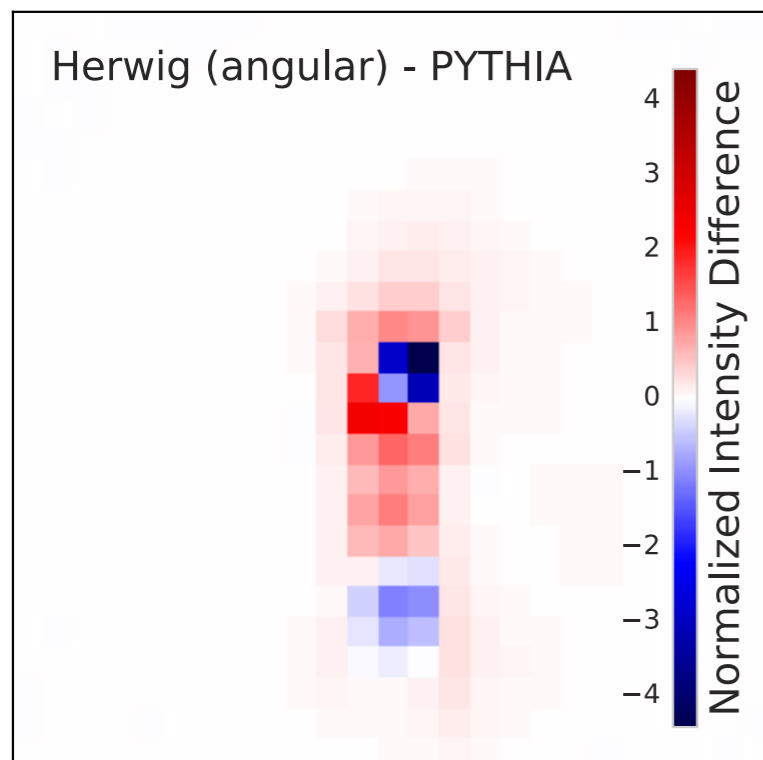
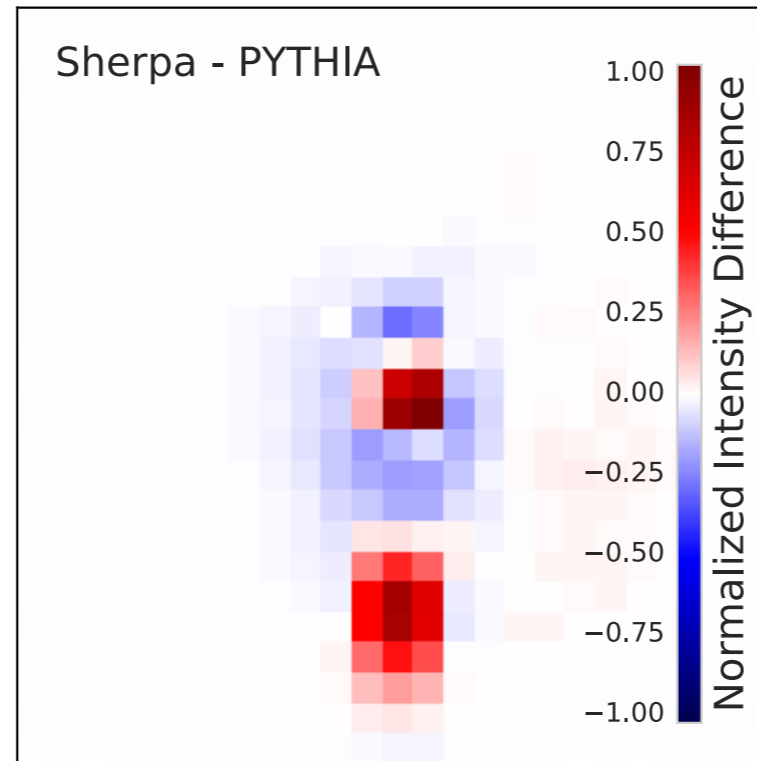
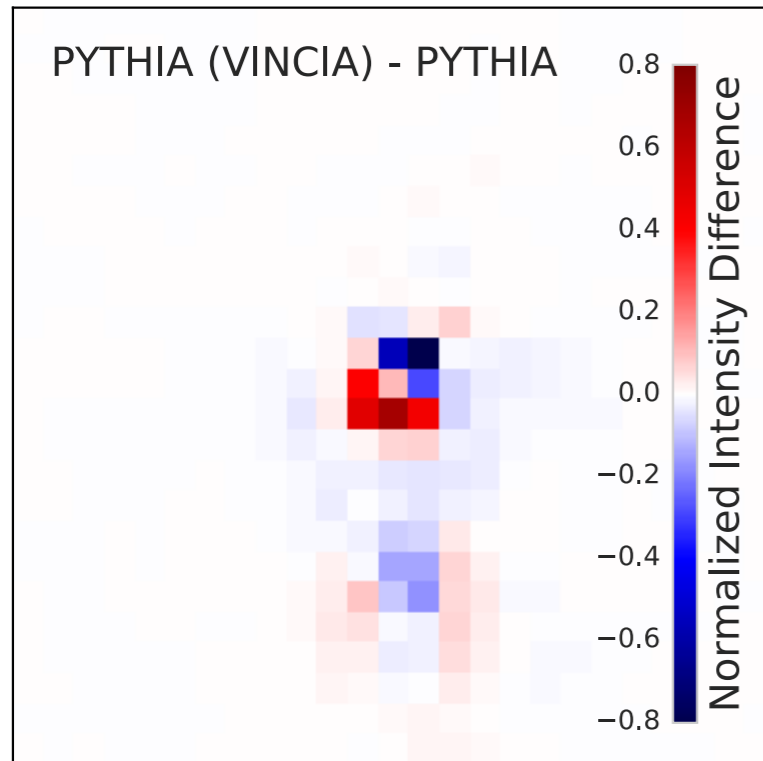


What if data and simulation are very different?

...your classifier will be sub-optimal

Where next III: Learning directly from data

Boosted W boson jets



J. Barnard et al.
Phys. Rev. D 95, 014018 (2017)

DNN classifiers
can **exploit**
subtle features

subtle features are
hard to model !

we need to be
careful about which
models we use -
only data is correct

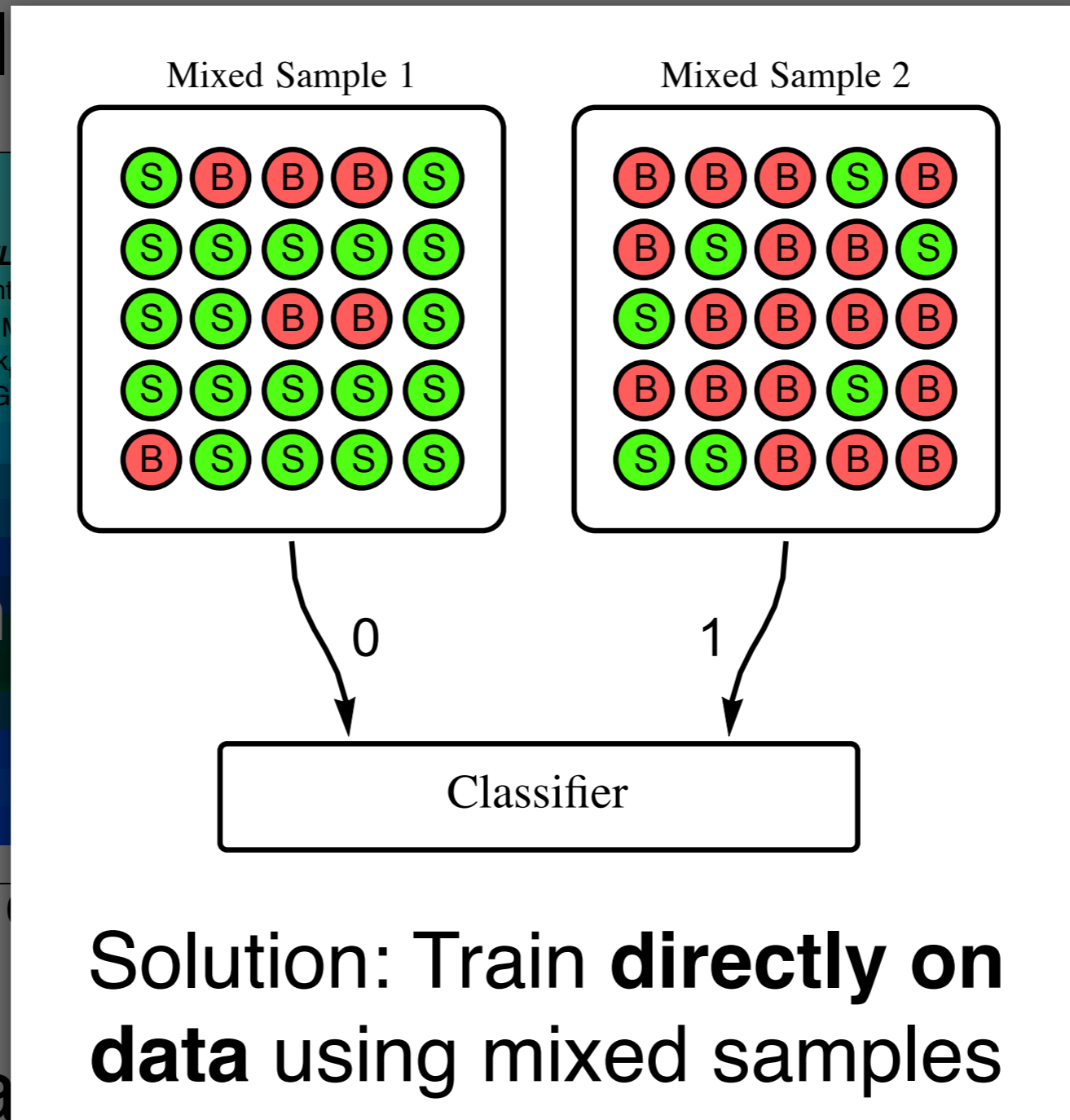
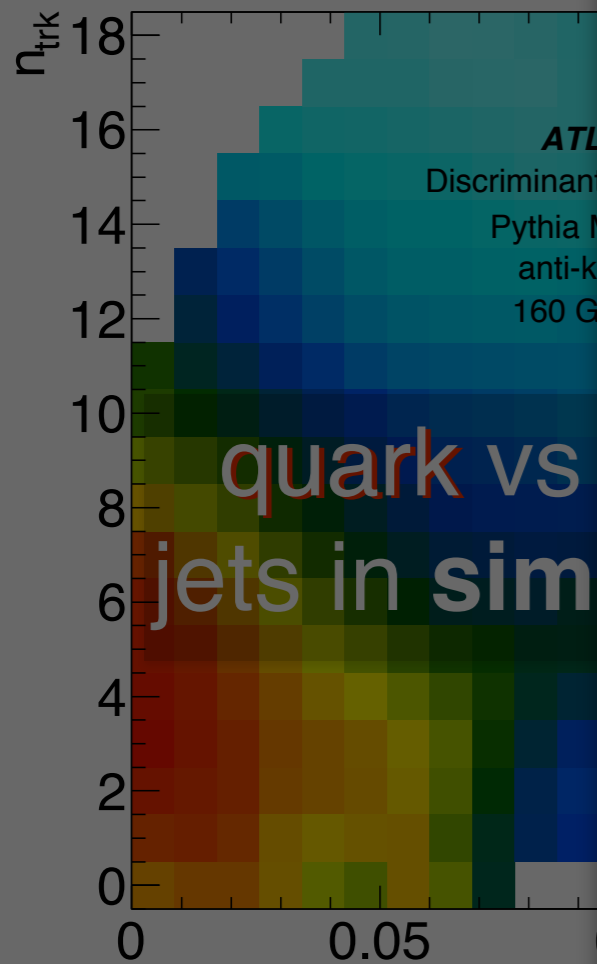
For a mixed approach, see
[G. Louppe et al.](#)

N.B. not all of these have been tuned to the same data

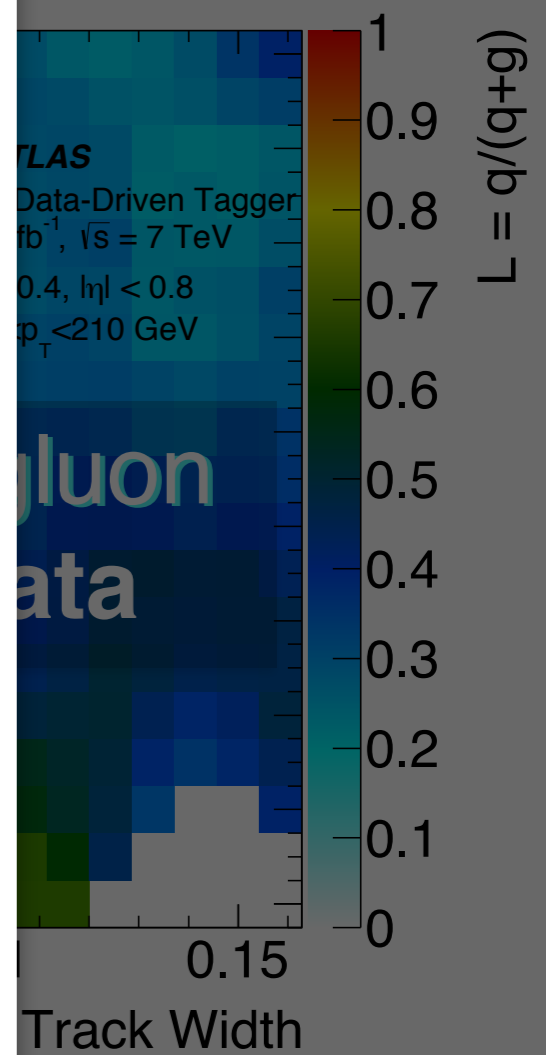
Where next III: Learning directly from data

For supervised learning, we depend on labels

label



ion



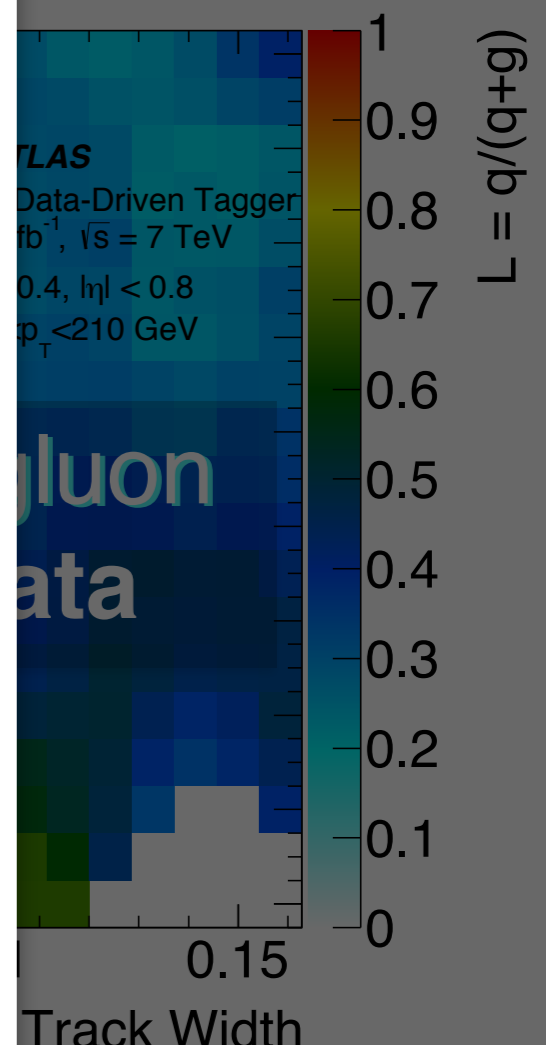
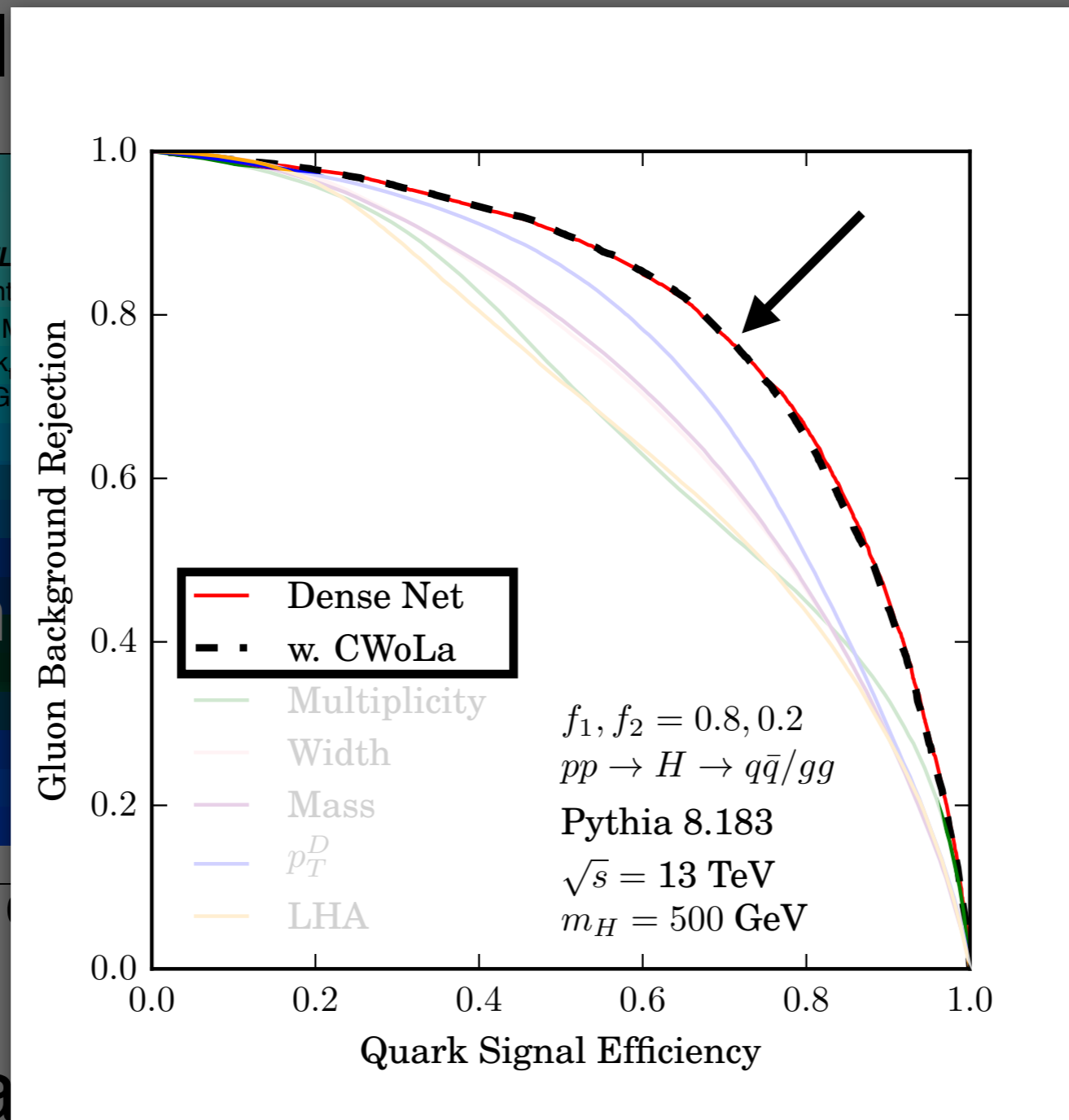
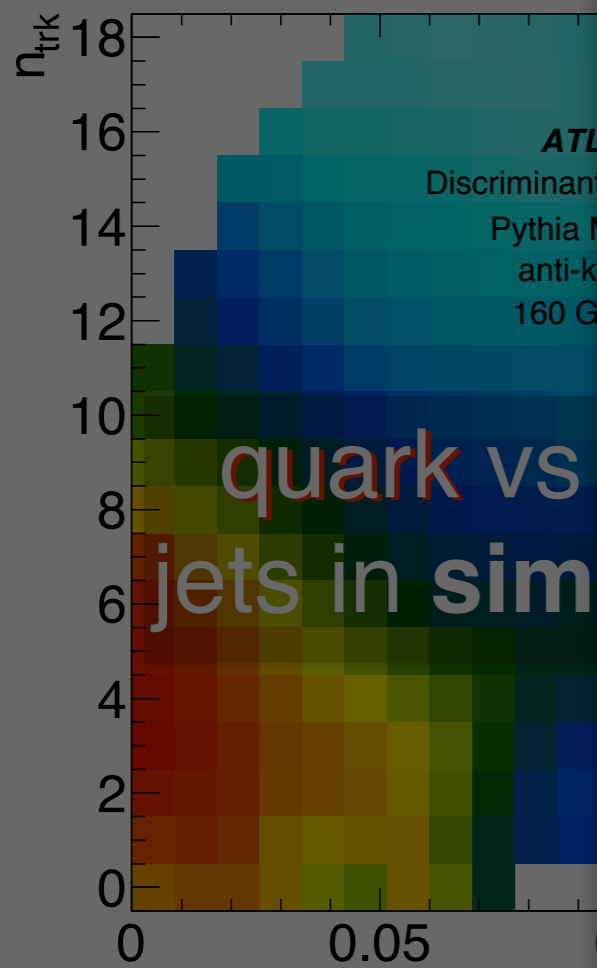
What if da

ifferent?

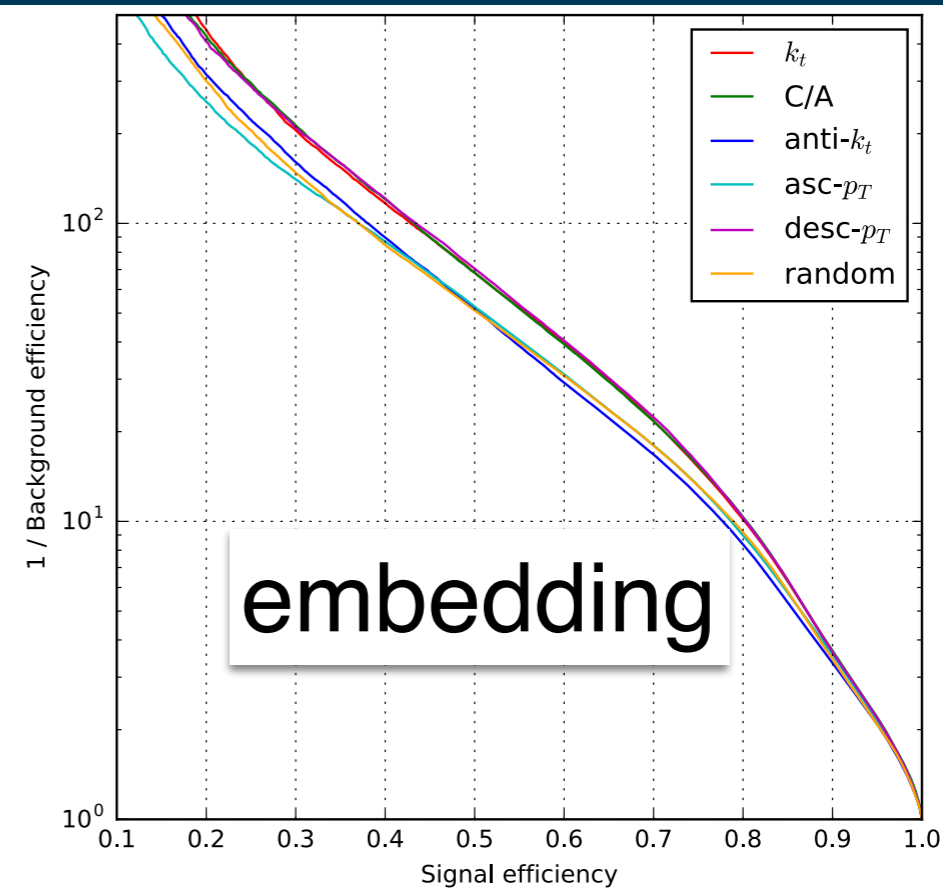
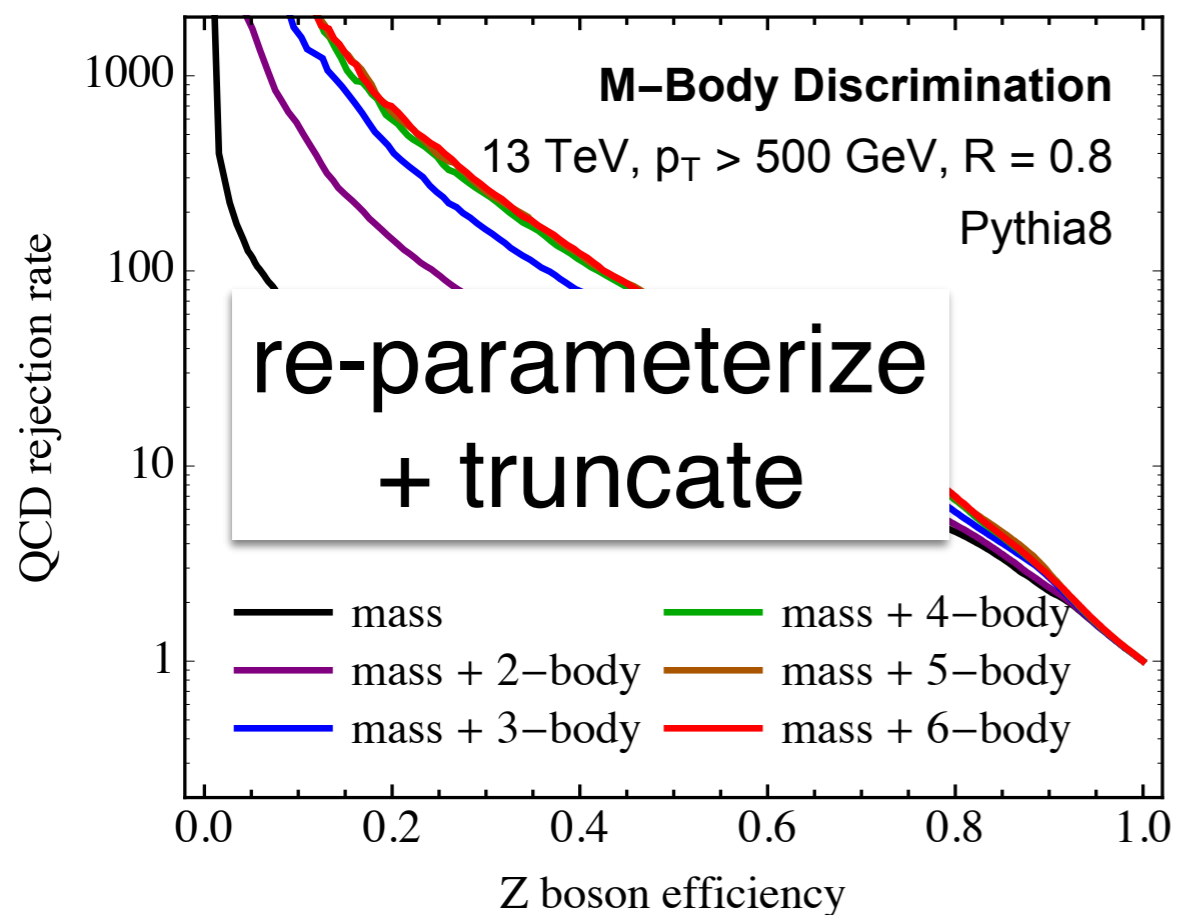
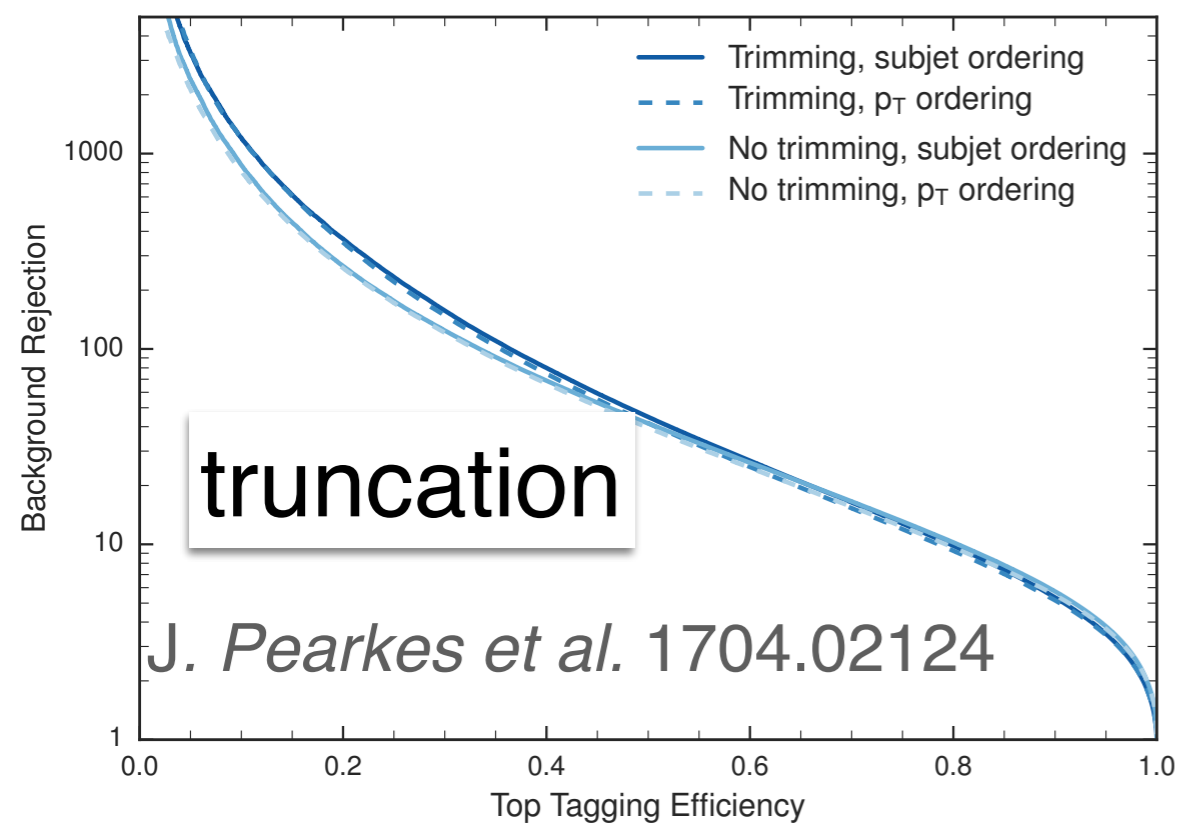
...your classifier will be sub-optimal

Where next III: Learning directly from data

For supervised learning, we depend on labels
label



Beyond Images



A. Butter et al. 1707.08966
(truncate + augment/embed)

K. Datta et al. 1710.01305 (re-param)

T. Cheng 1711.02633 (RNN)

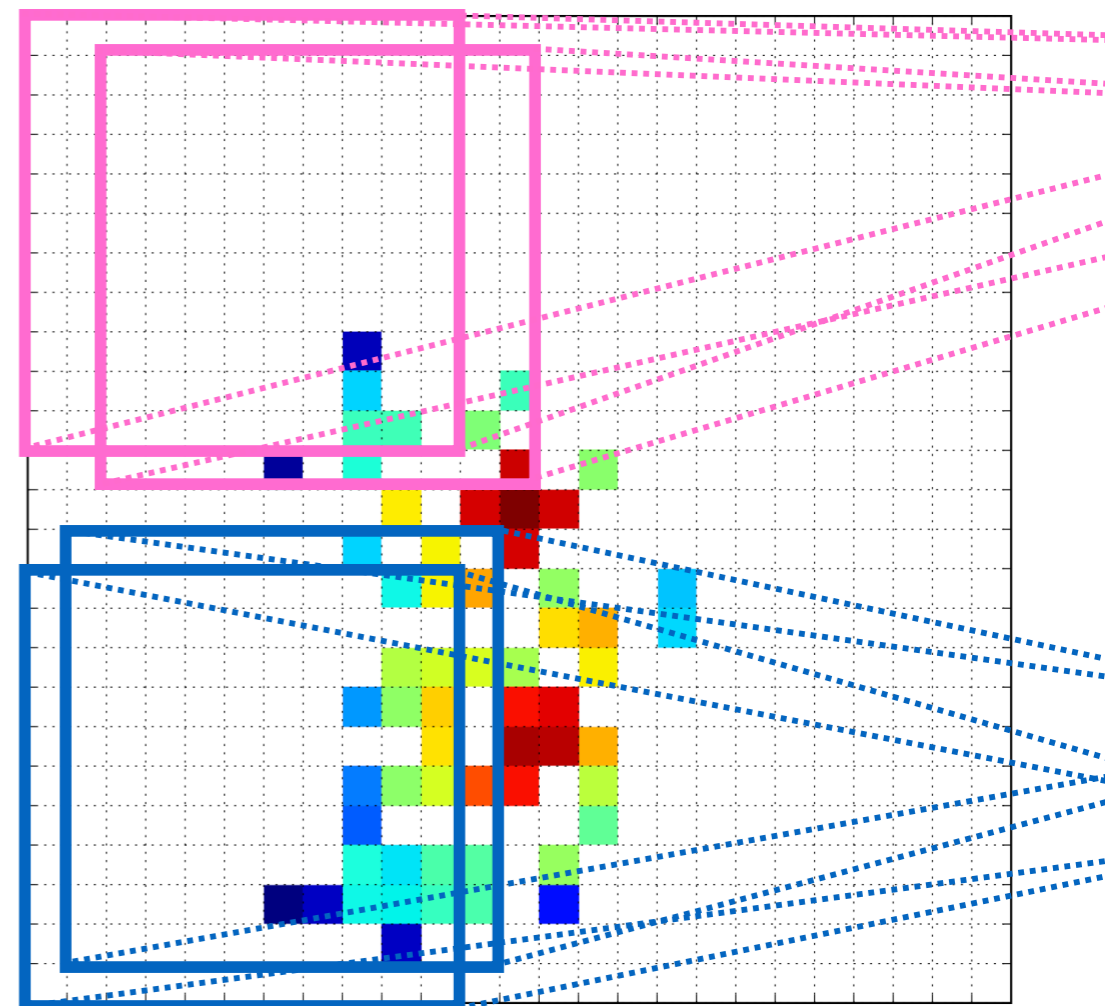
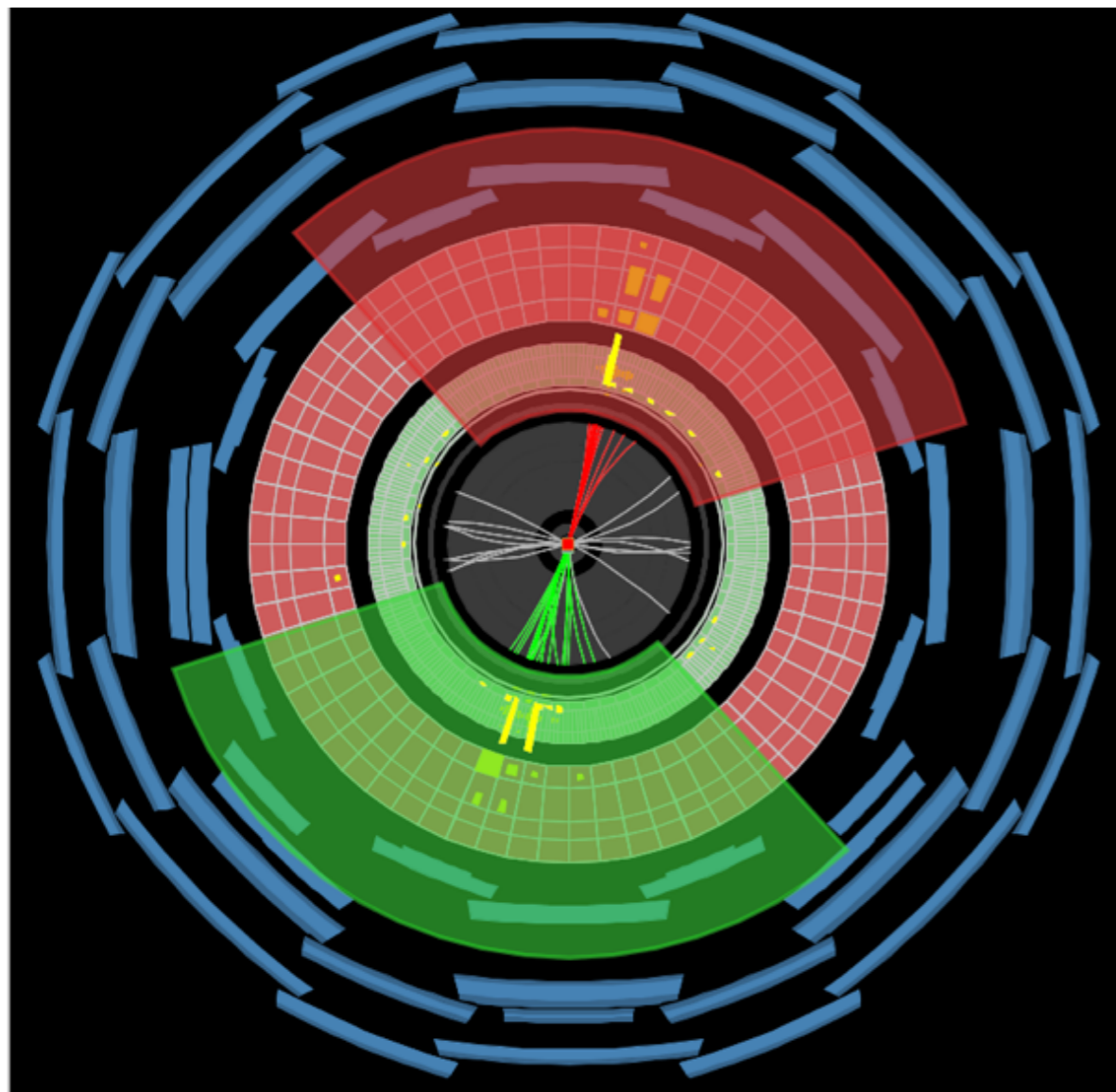
J. A. Aguilar-Saavedra et al.
1709.01087 (re-param)

+ flavor tagging (see backup)

+ many more results at the
dedicated workshop next month!

Conclusions and outlook

(Jet) image-based NN classification, regression, and generation are powerful tools for fully exploiting the physics program at the LHC



The key to robustness is to study what is being learned; this may even help us to learn something new about nature!

Collaborators



Lucio
Dery

Stanford



Michela
Paganini

Yale



Eric
Metodiev

MIT



Patrick
Komiske

MIT



Zihao
Jiang

Stanford



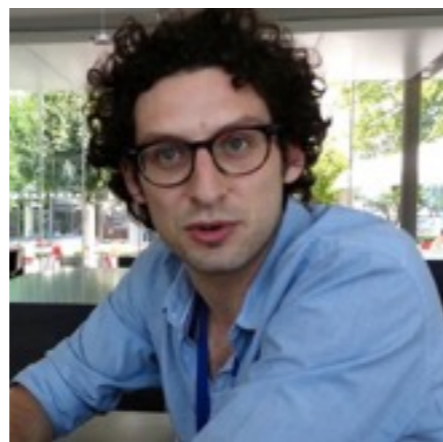
Francesco
Rubbo

SLAC



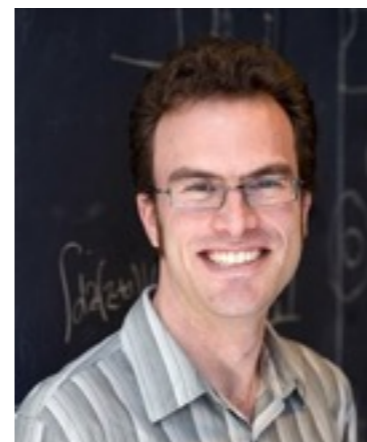
Luke
de Oliveira

VAI tech.



Michael
Kagan

SLAC



Jesse
Thaler

MIT



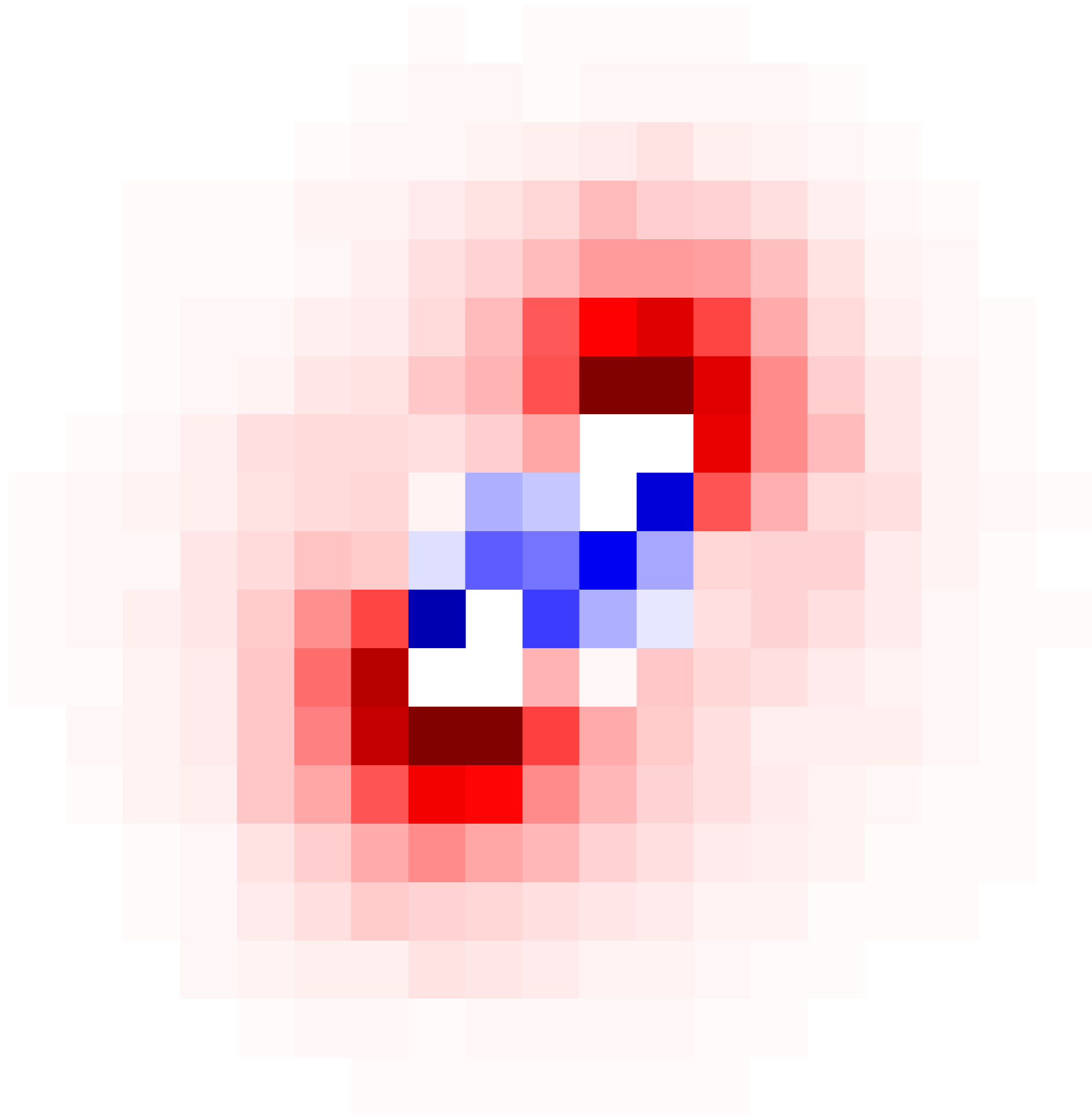
Matt
Schwartz

Harvard



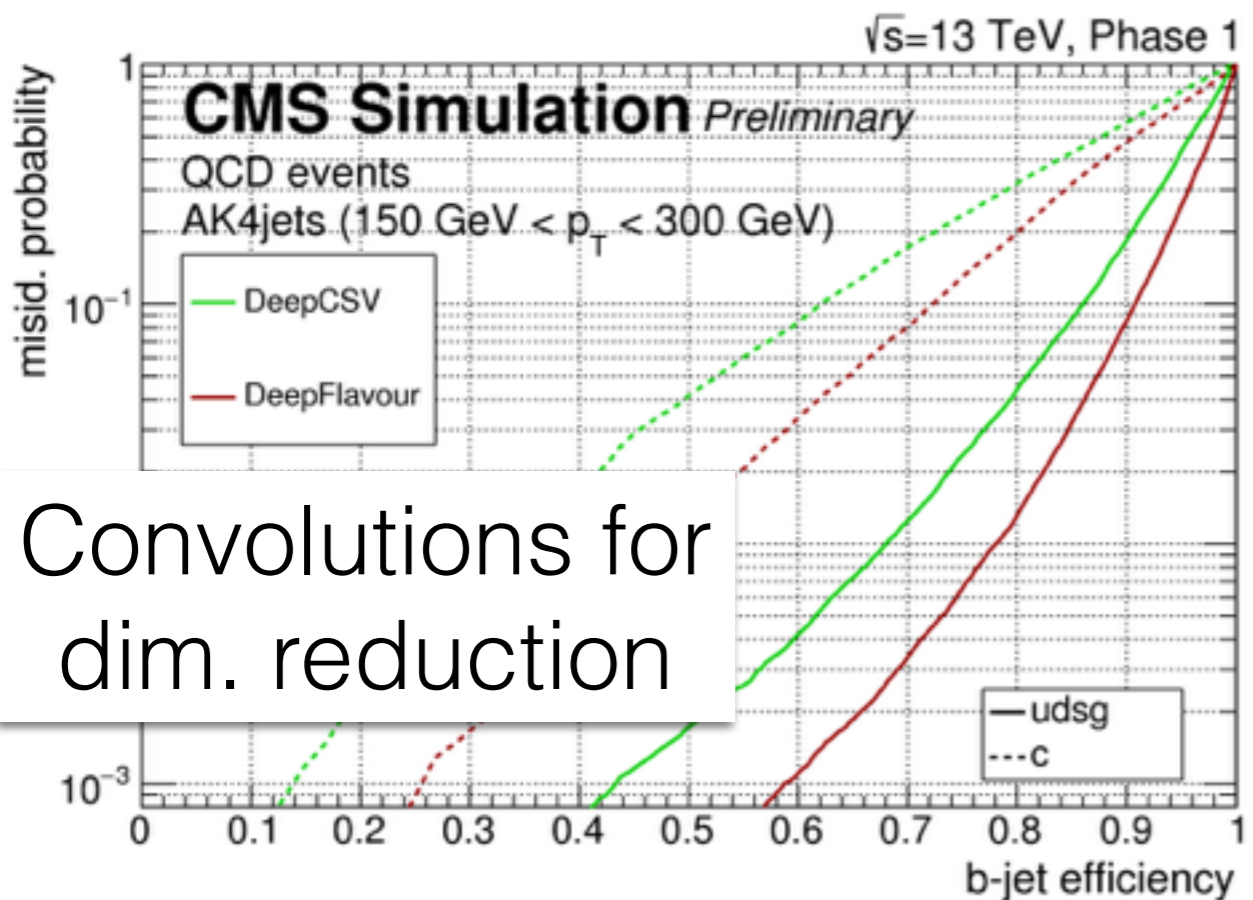
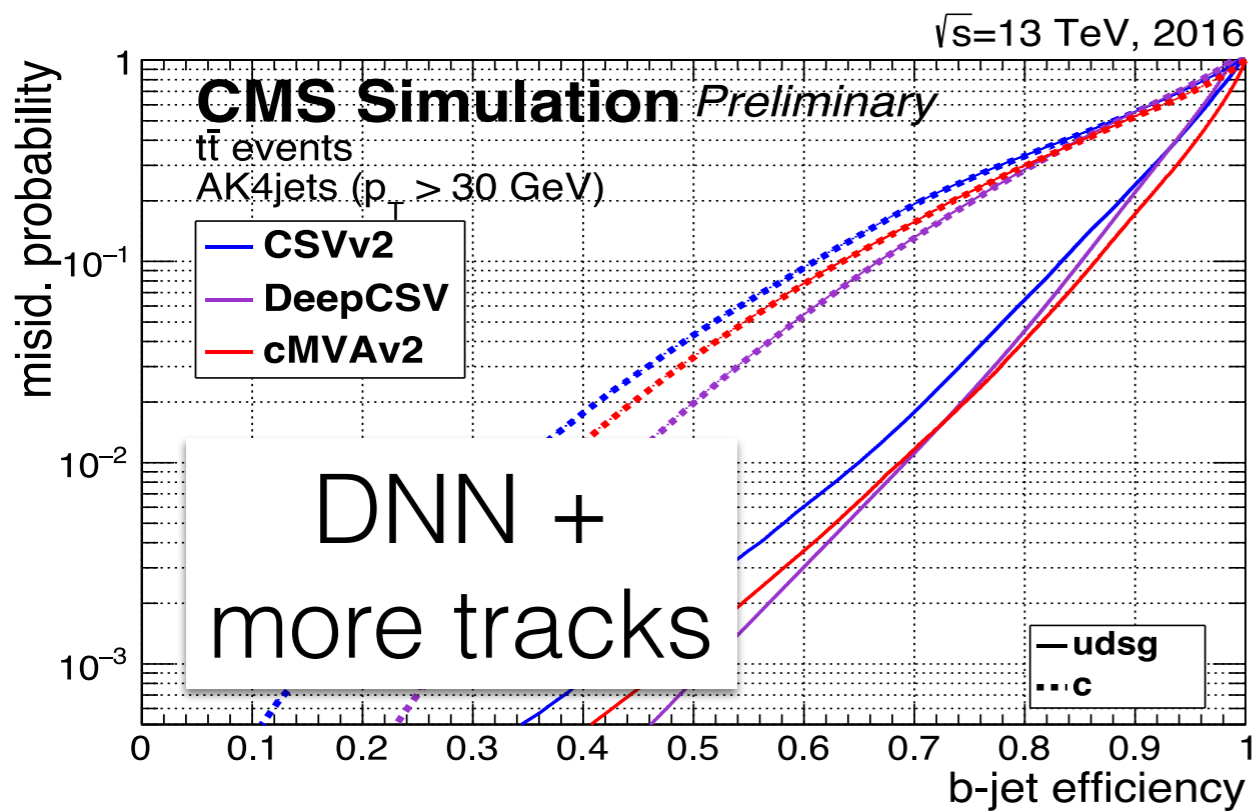
Ariel
Schwartzman

SLAC



Fin.

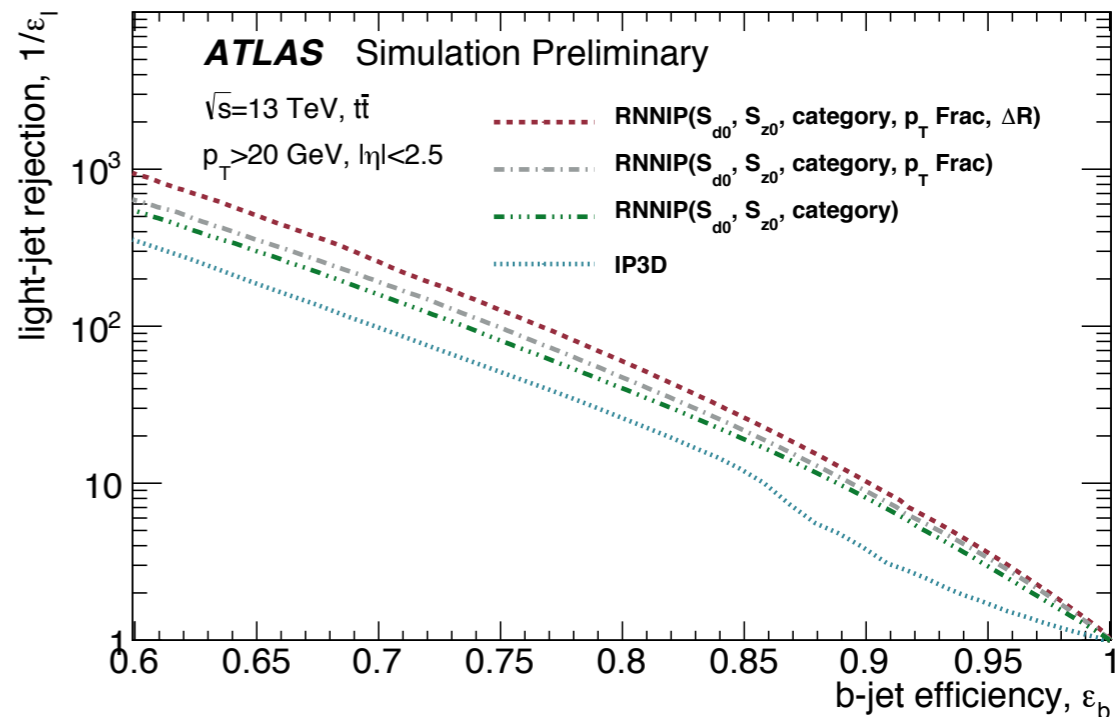
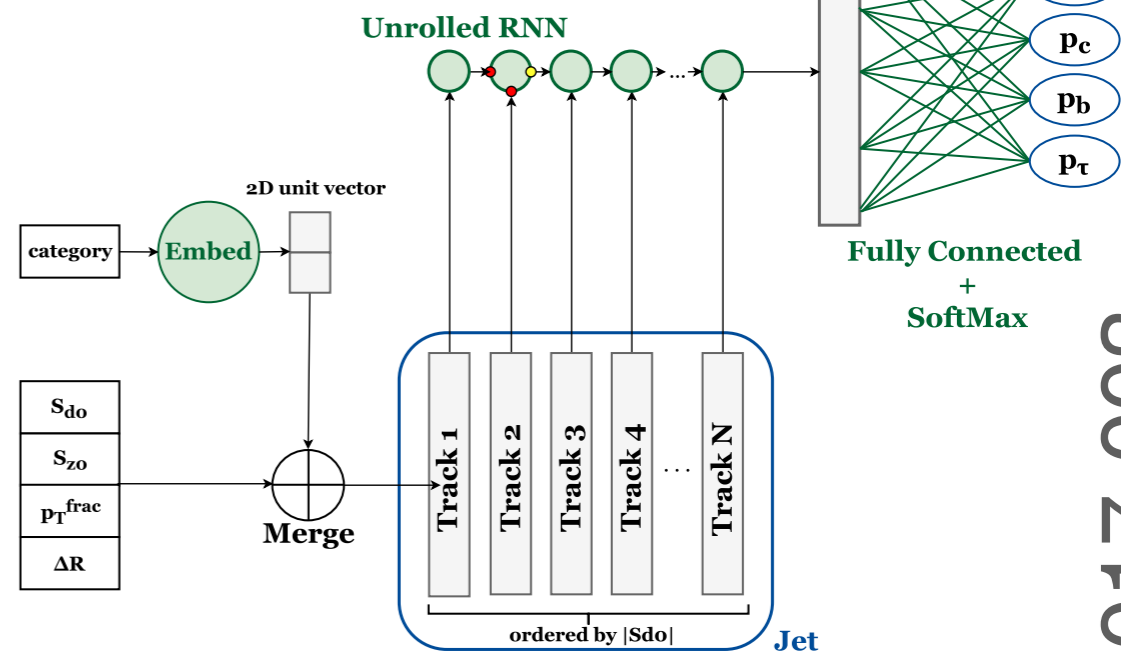
Backup



CMS-DP-2017-005

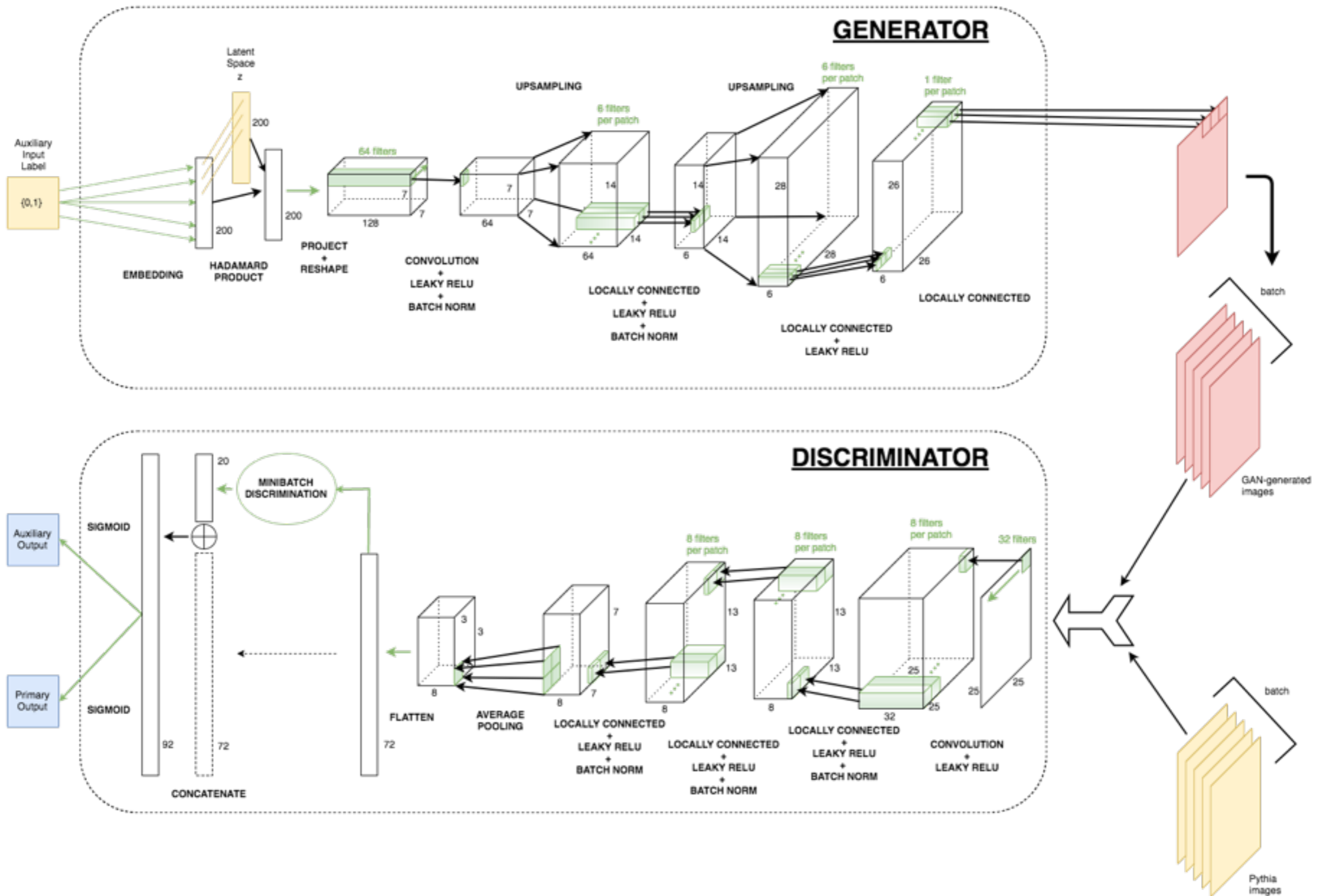
CMS-DP-2017-013

Sequence of tracks
 → Recurrent NNs

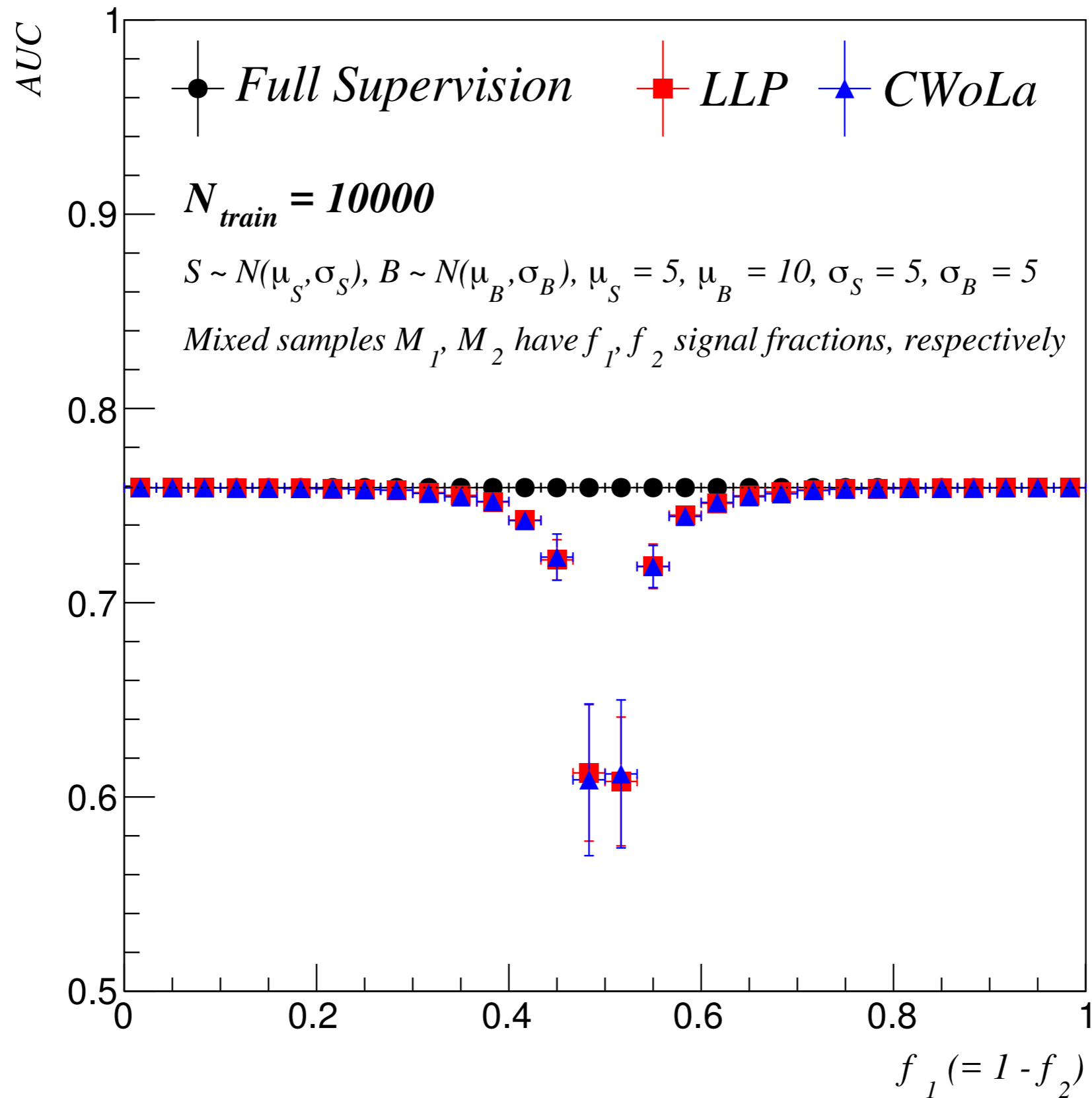


ATL-PHYS-PUB-2017-003

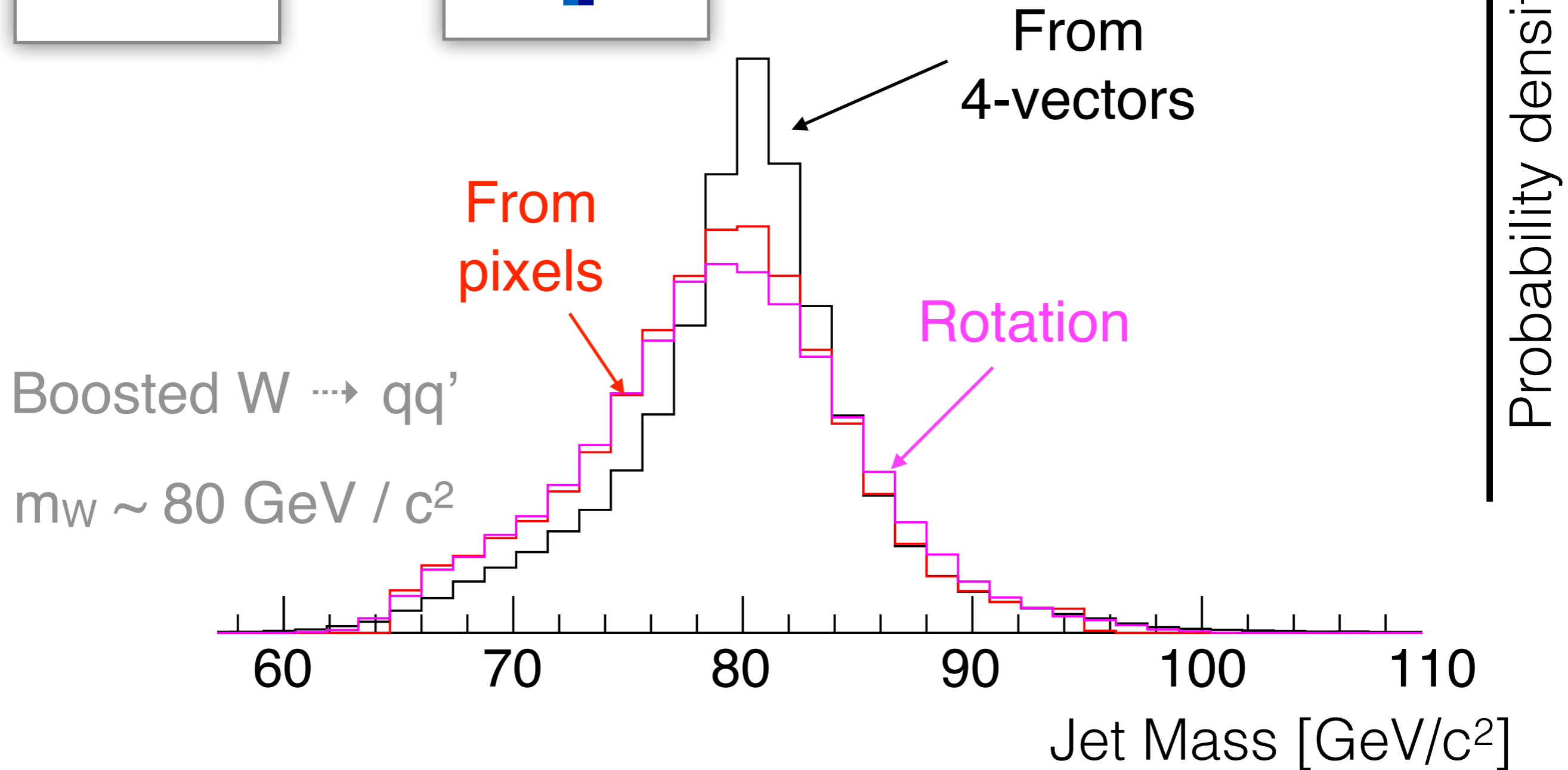
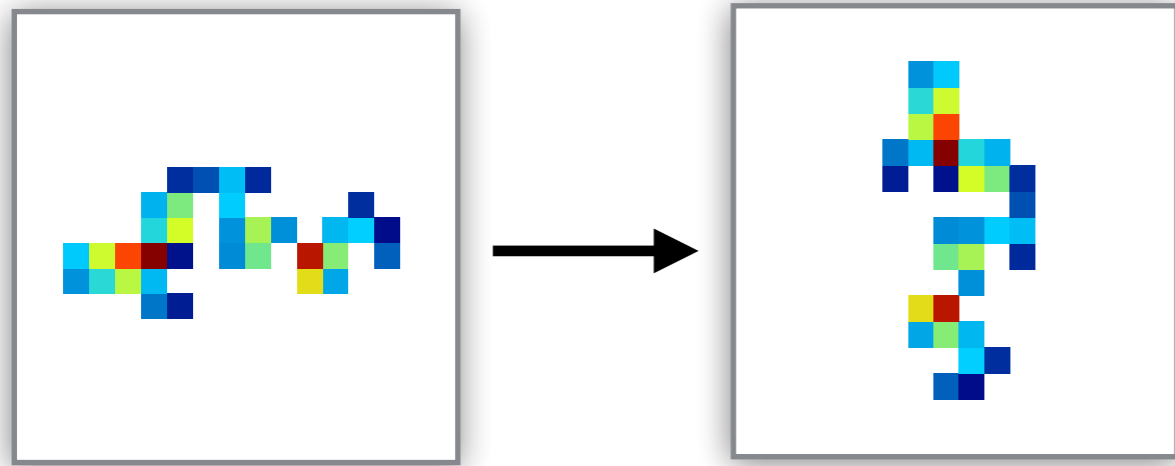
Locally Aware GAN (LAGAN)



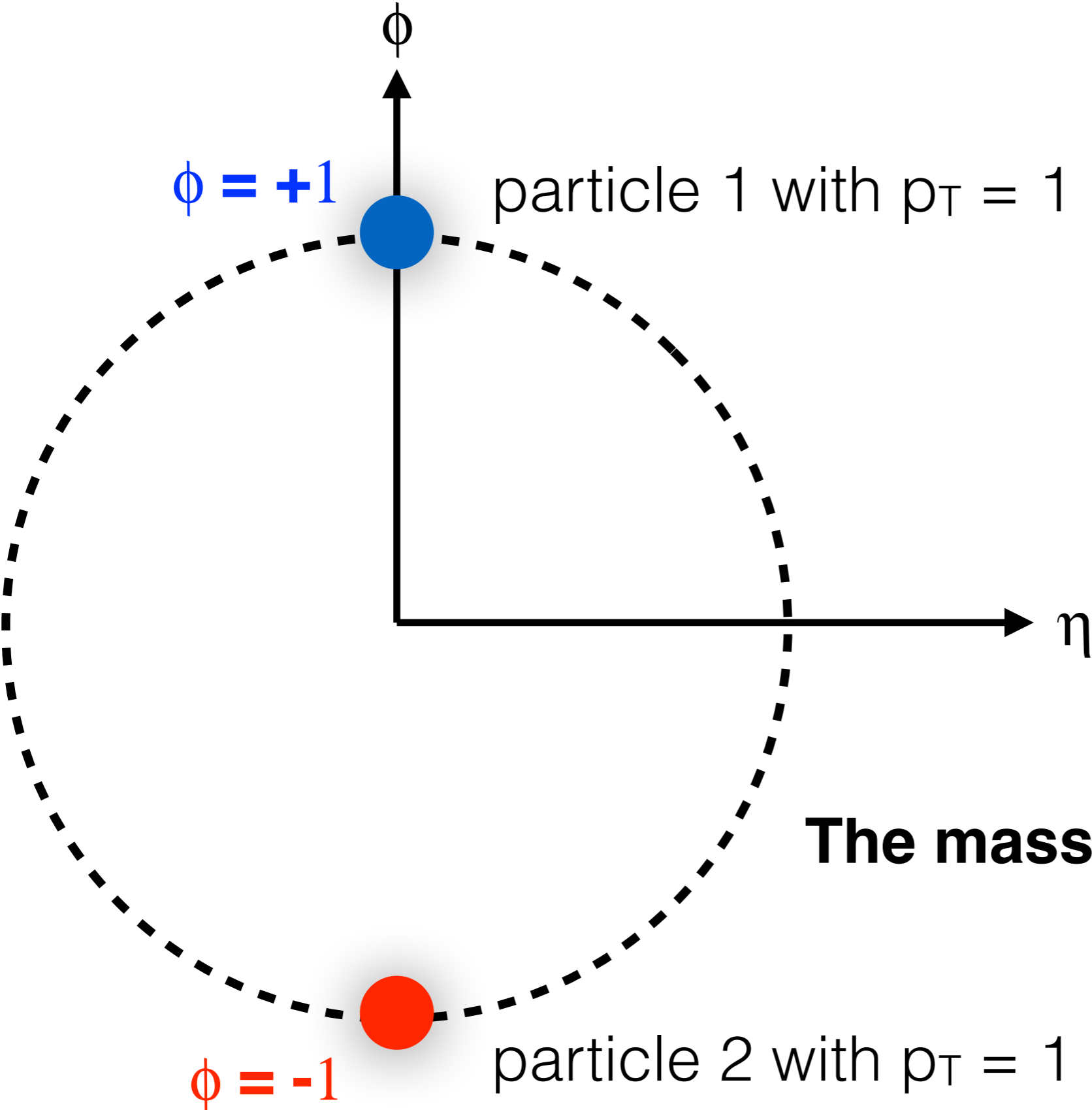
Learning when you know (almost) nothing



Pre-processing & spacetime symmetries

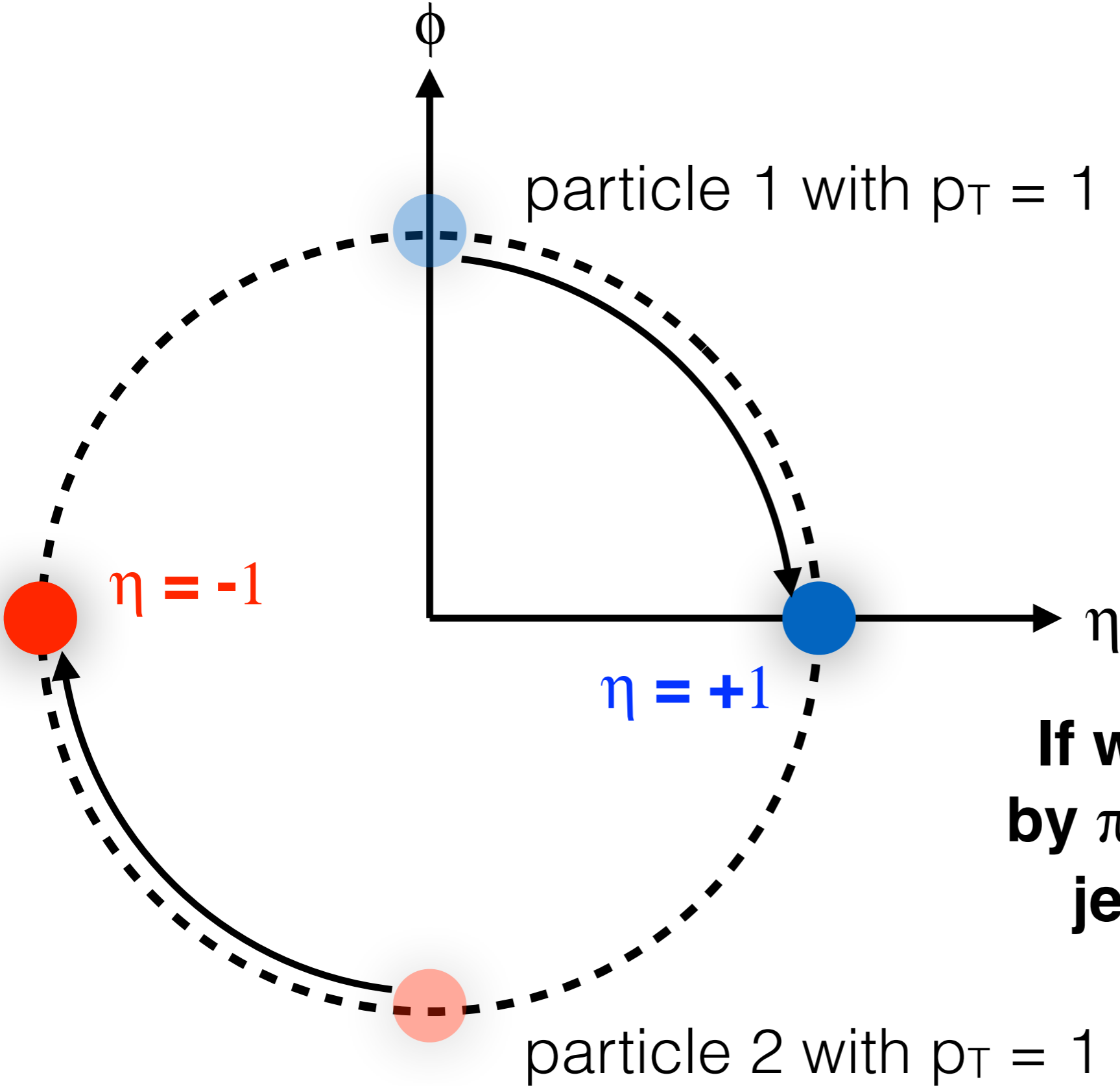


Pre-processing & spacetime symmetries



The mass of this 'jet' is ~ 1.7

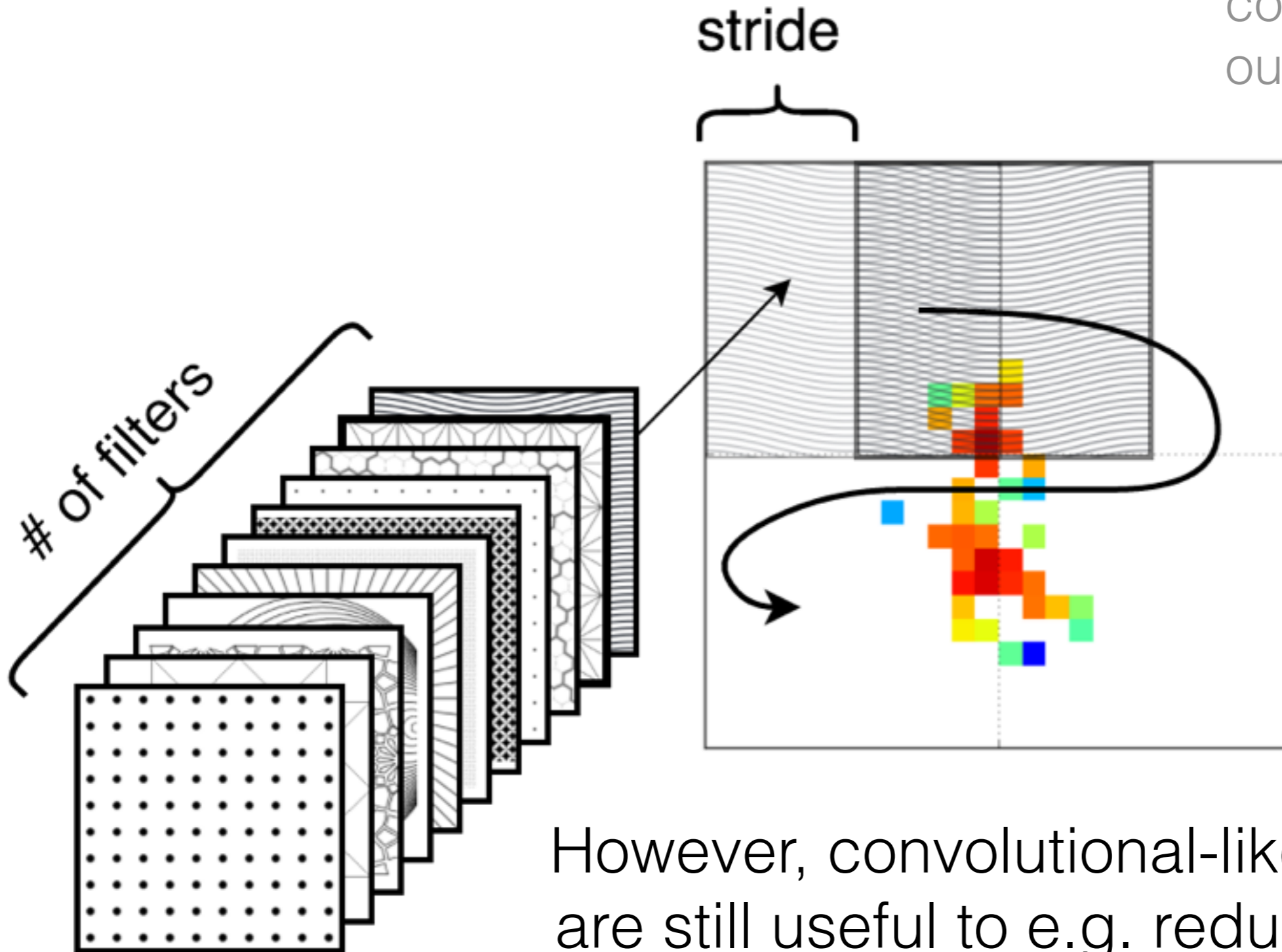
Pre-processing & spacetime symmetries



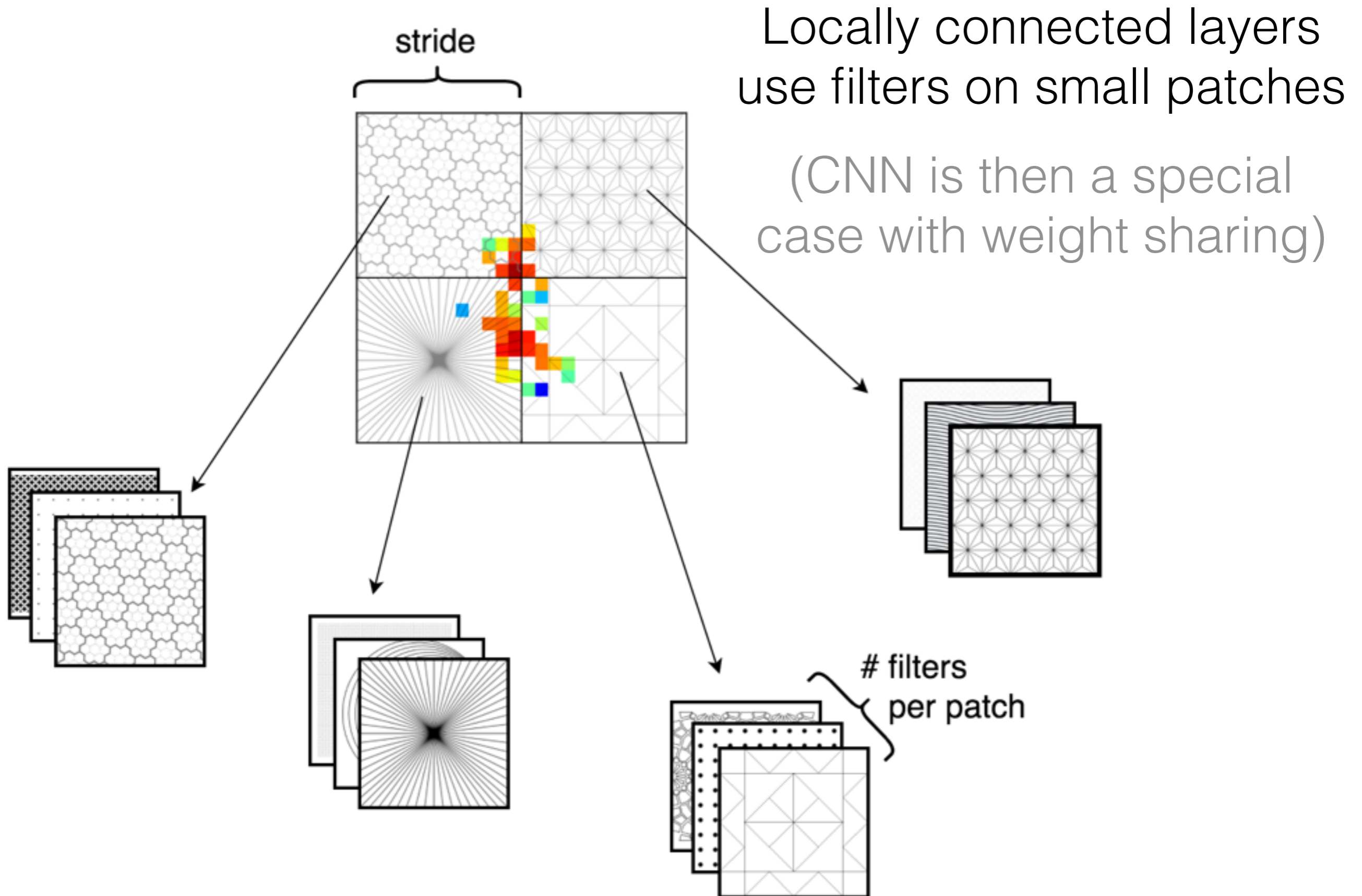
**If we rotate the jet
by $\pi/2$, then the new
jet mass is ~ 2.4**

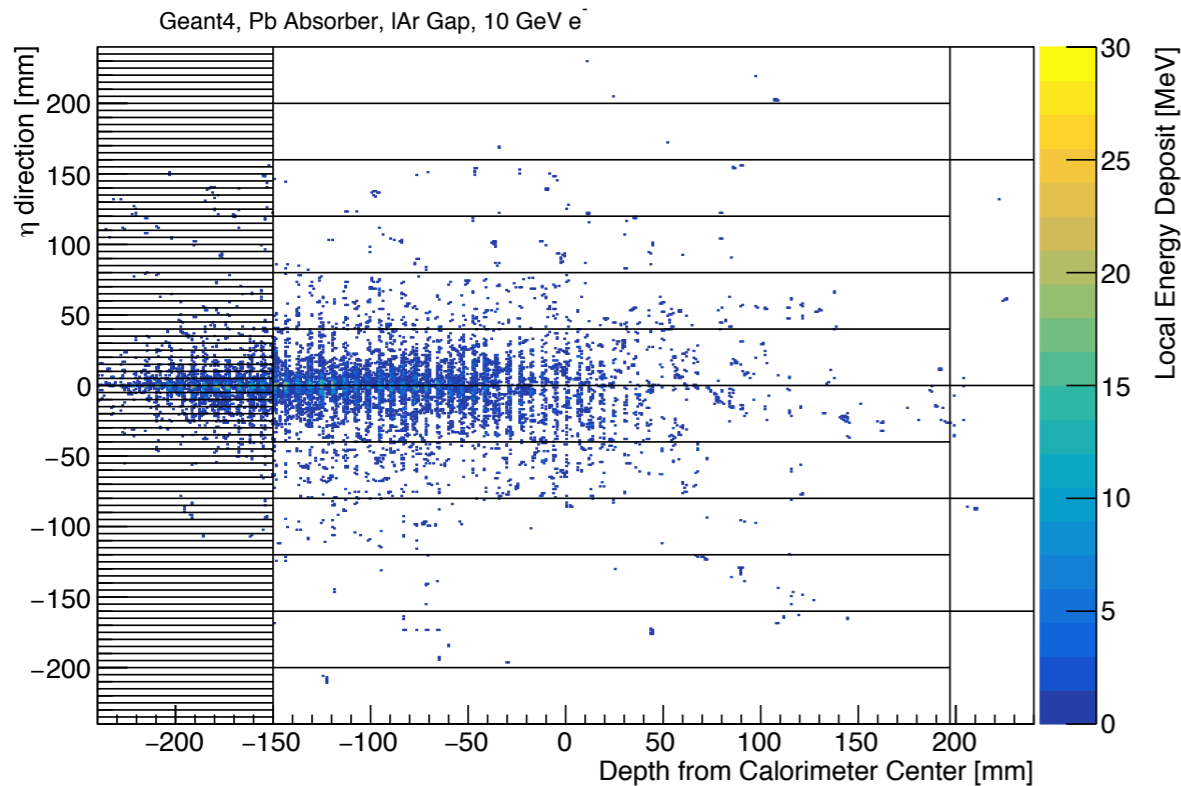
Due to the structure of the problem, we do not have translation invariance.

Classification studies found fully connected networks outperformed CNNs



However, convolutional-like architectures are still useful to e.g. reduce parameters

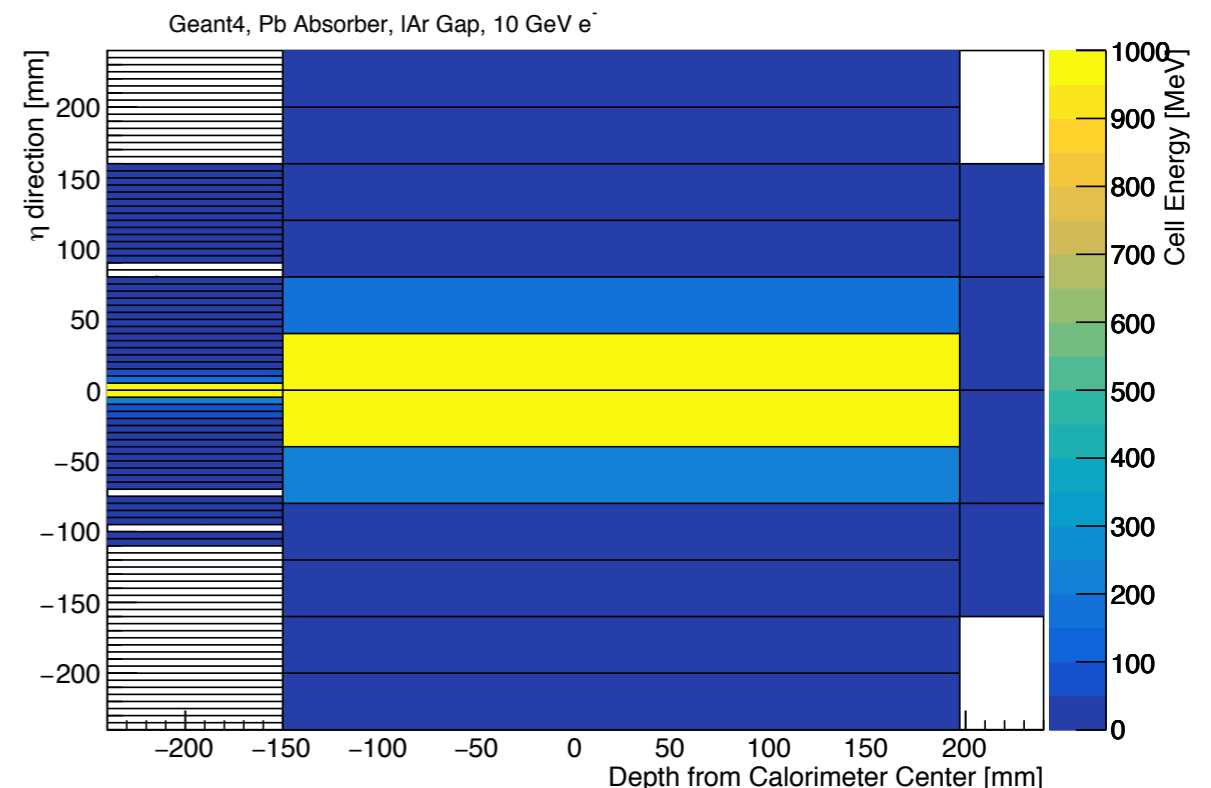




We take as our model a 3-layer LAr calorimeter, inspired by the ATLAS barrel EM calorimeter

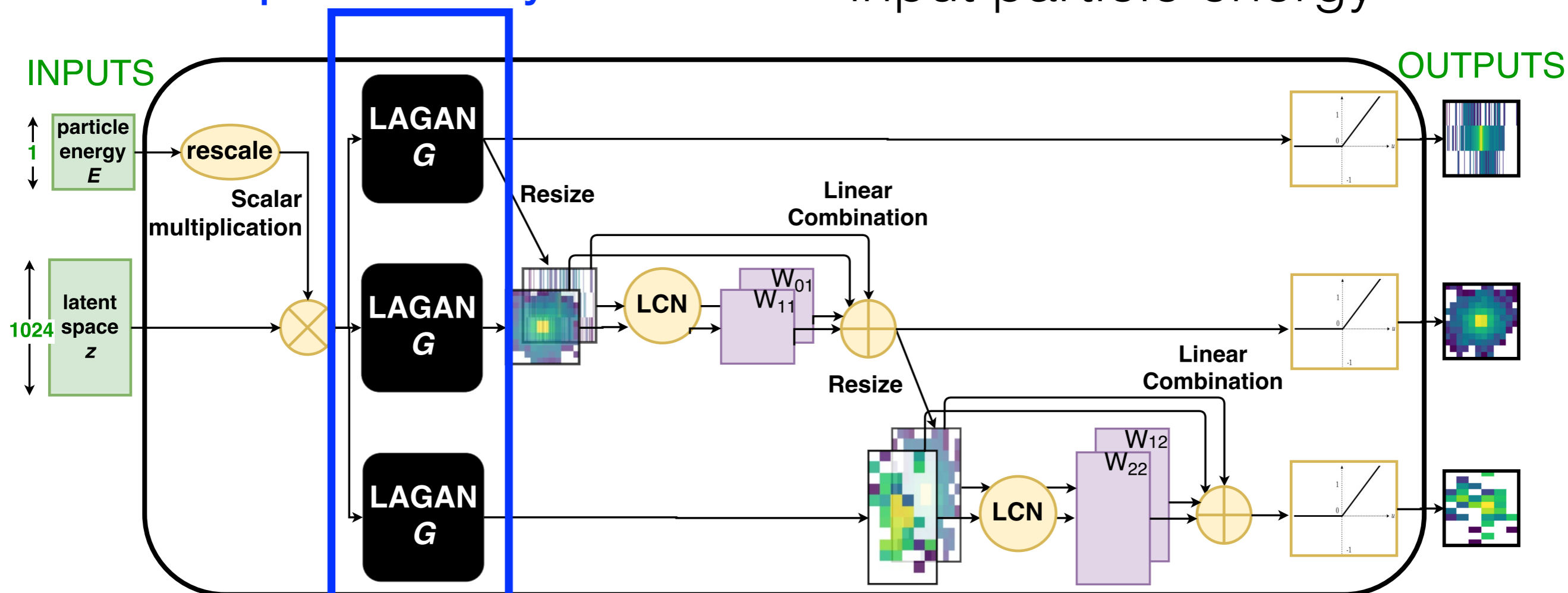
A single event may have $O(10^3)$ of particles showering in the calorimeter - too cumbersome to do all at once (now)

We exploit factorization of energy depositions



One 'jet image'
per calo layer

One network per particle type;
input particle energy



use layer i as
input to layer $i+1$

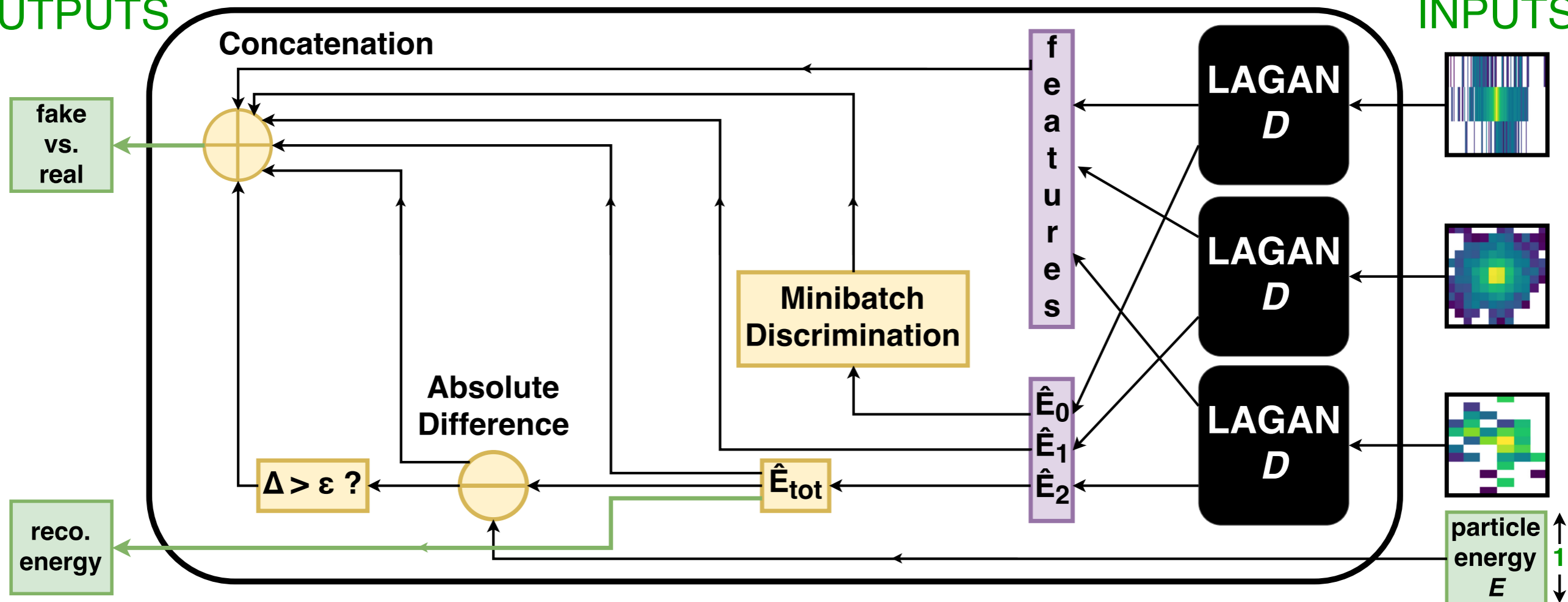
ReLU to
encourage
sparsity

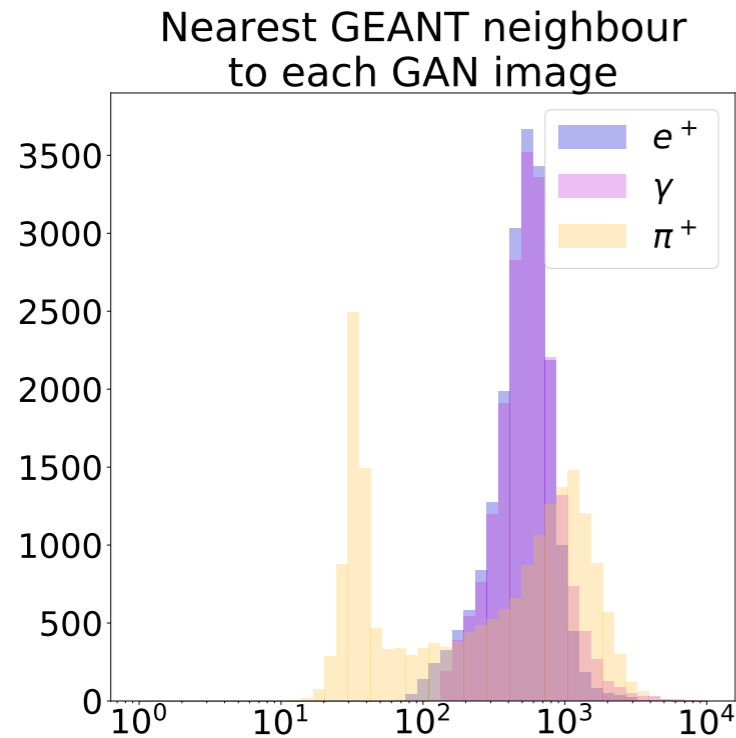
help avoid
'mode collapse'



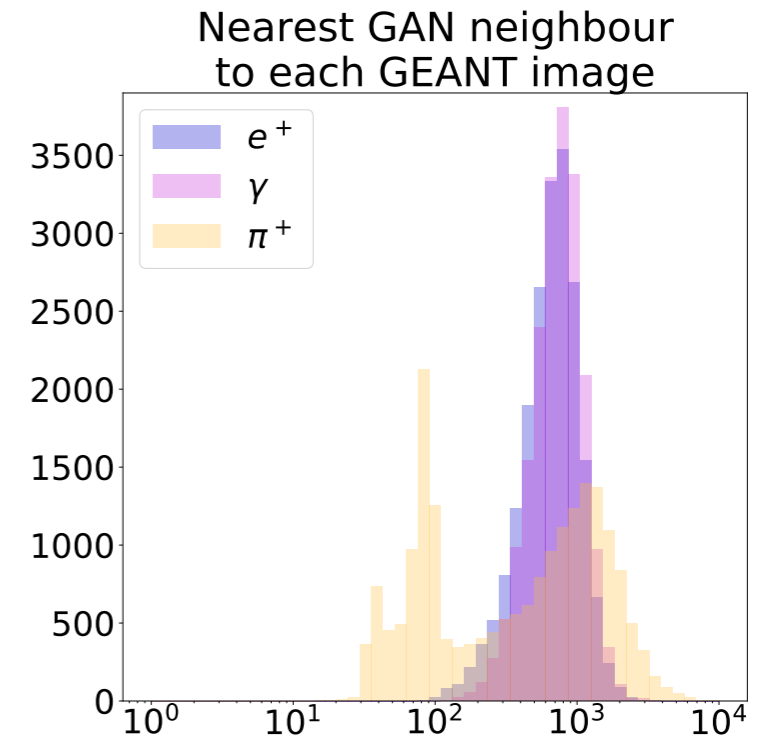
OUTPUTS

INPUTS

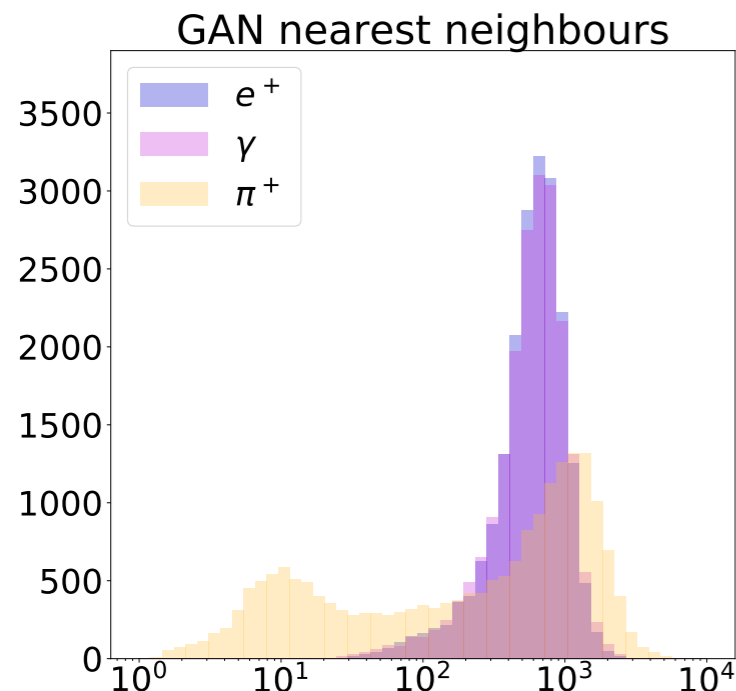




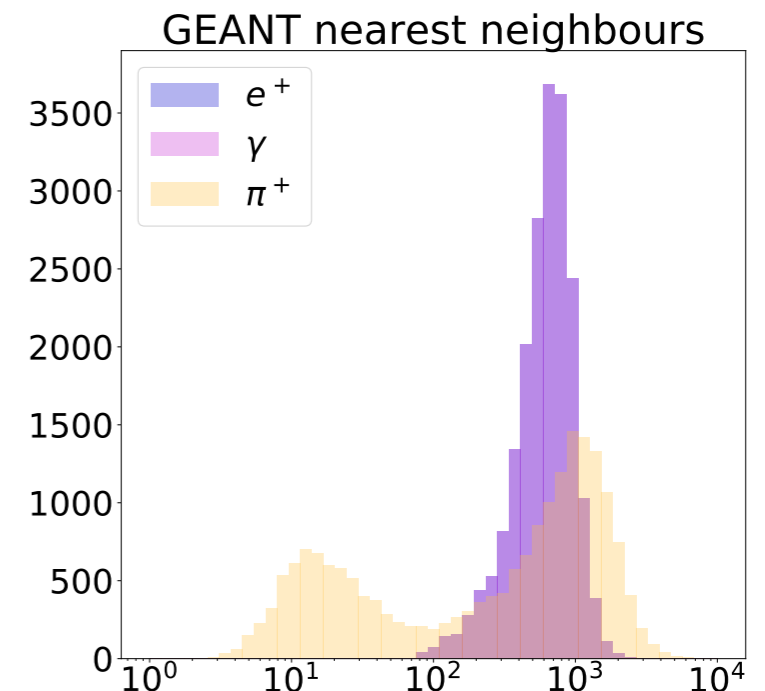
not
memorizing



A key challenge in training GANs is the diversity of generated images. This does not seem to be a problem for CaloGAN.

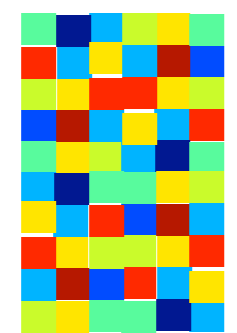


no mode
collapse

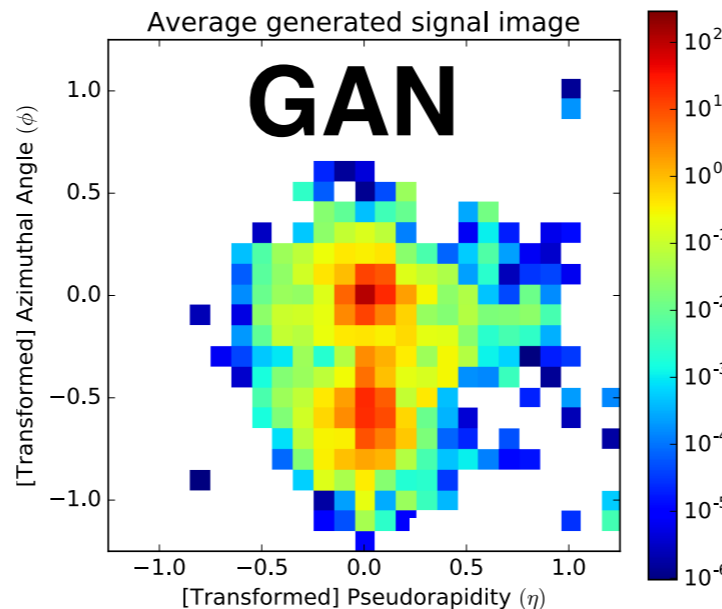


*M. Paganini, L. de Oliveira, and **BPN** 1705.05927, 1705.02355*

Generative Adversarial Networks (GAN):
*A two-network game where one **maps noise to images**
and one **classifies images as fake or real.***

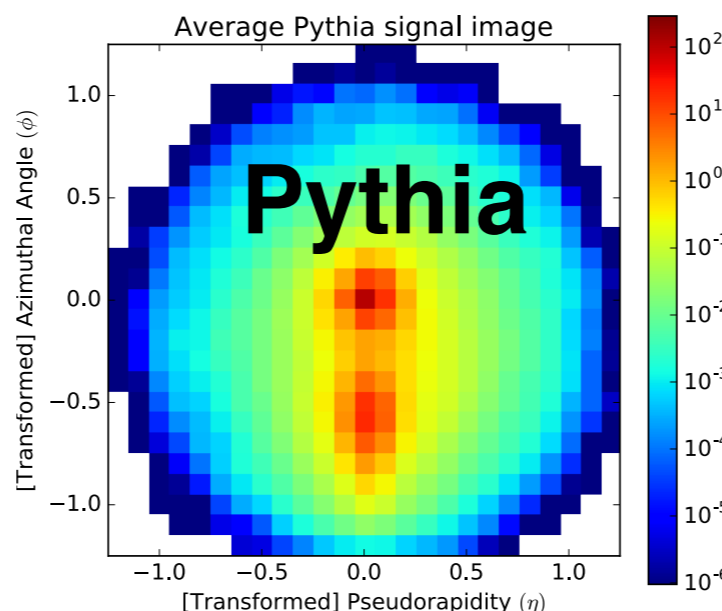


noise



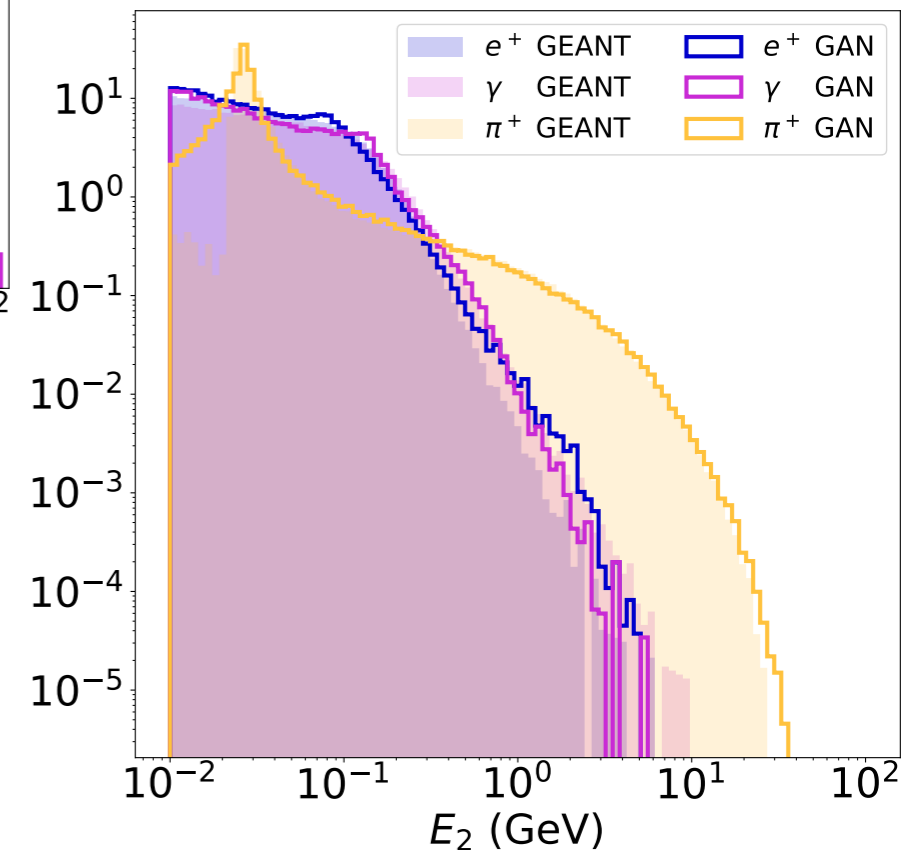
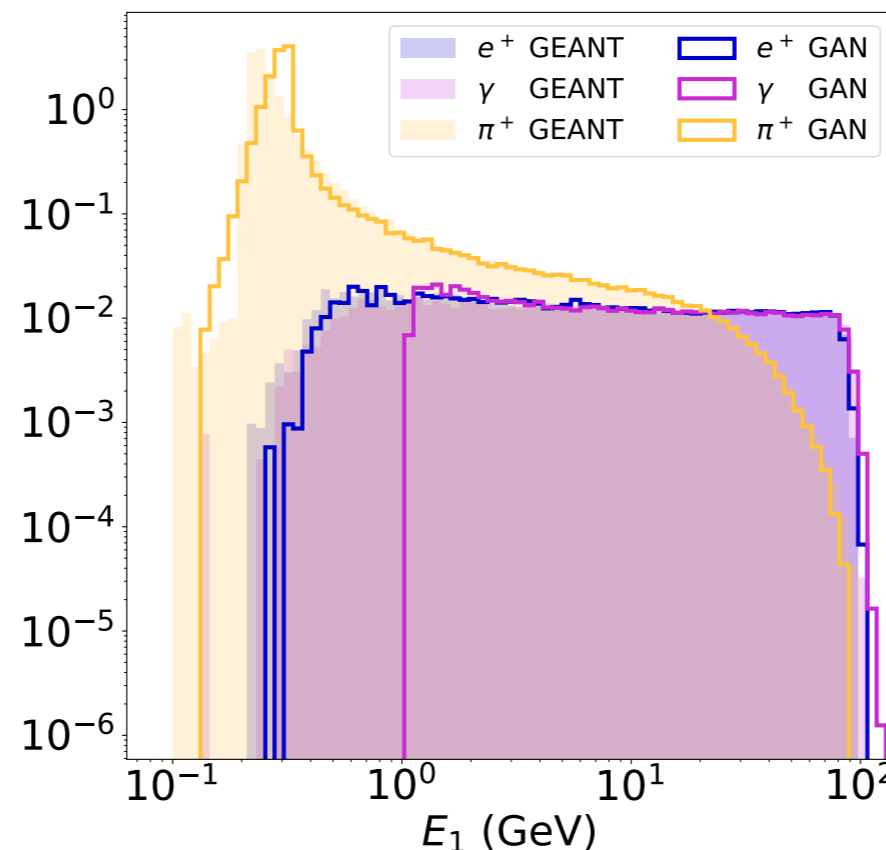
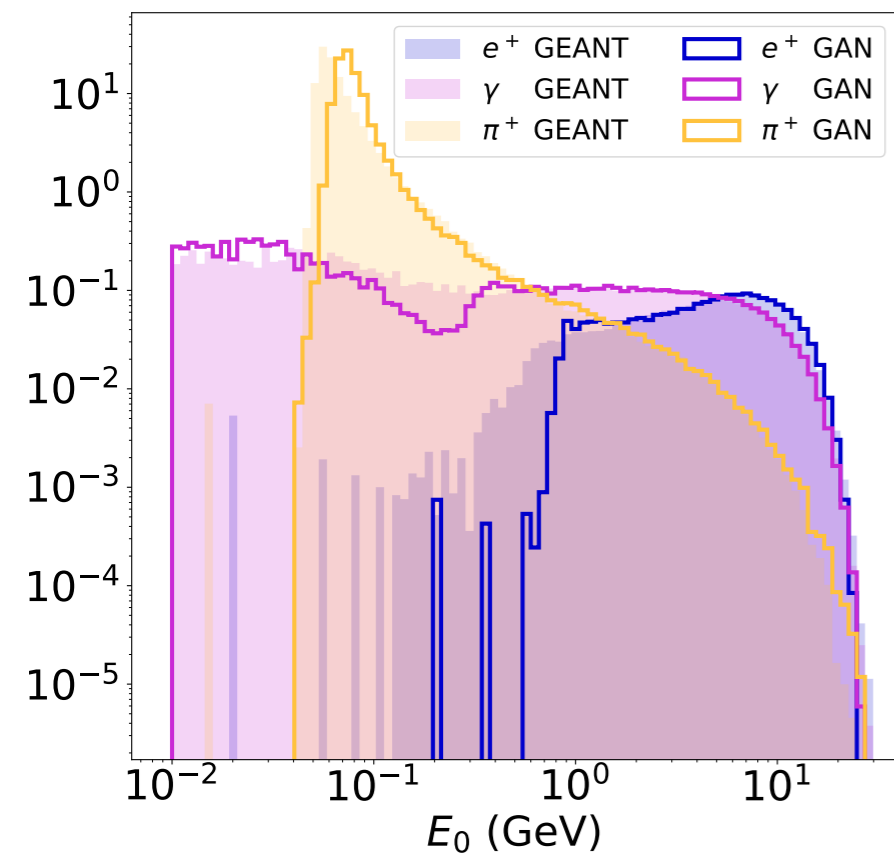
{real, fake}

When **D** is maximally confused, **G** will be a good generator

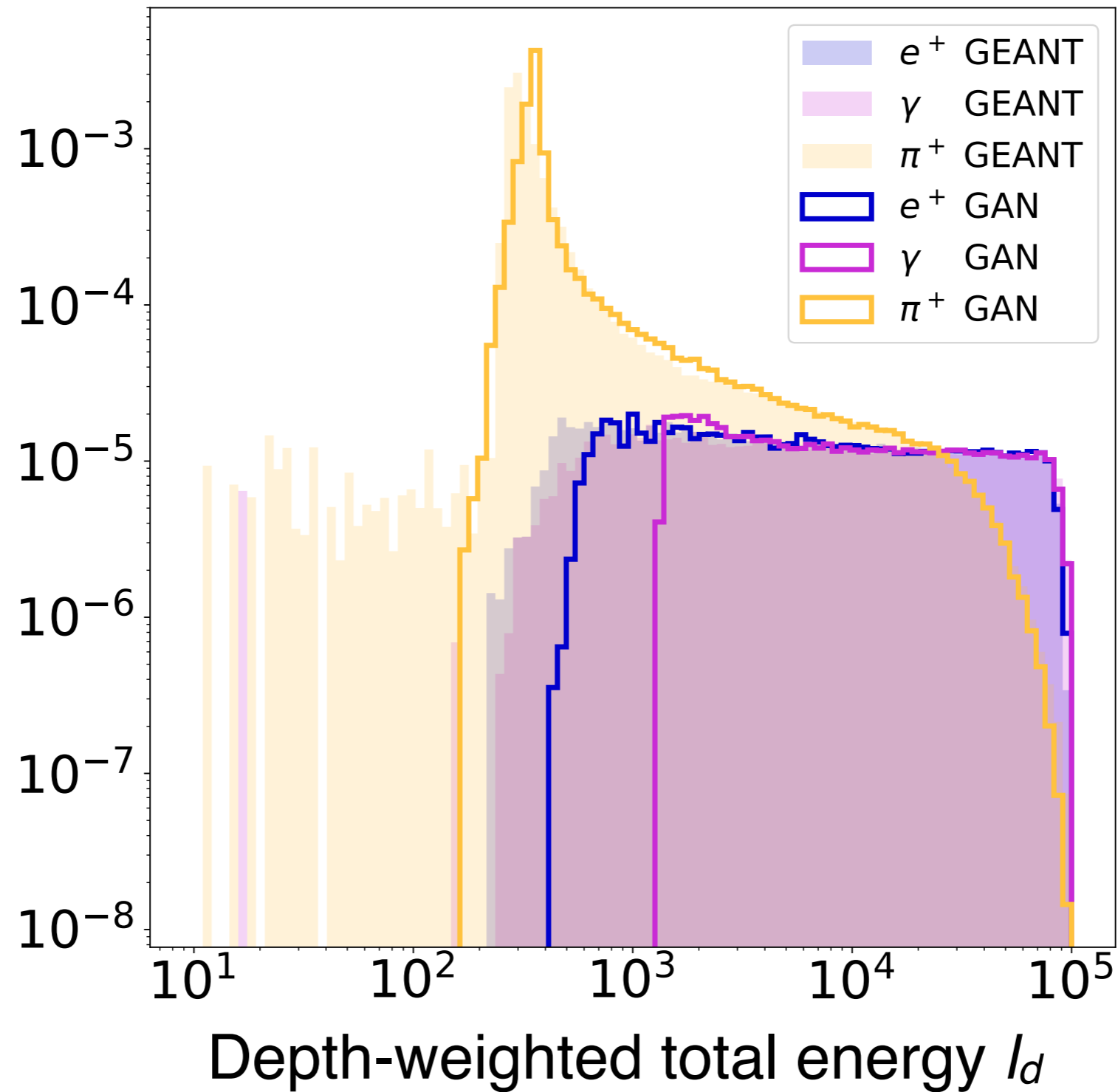


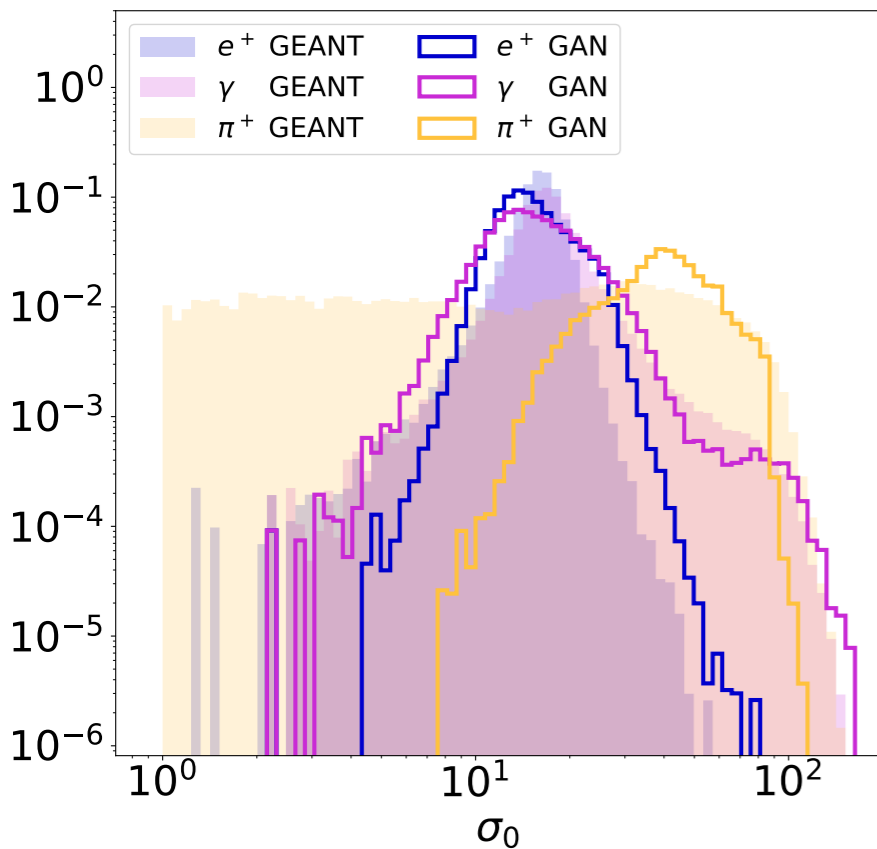
Physics-based simulator

Pions deposit much less energy in the first layers; leave the calorimeter with significant energy

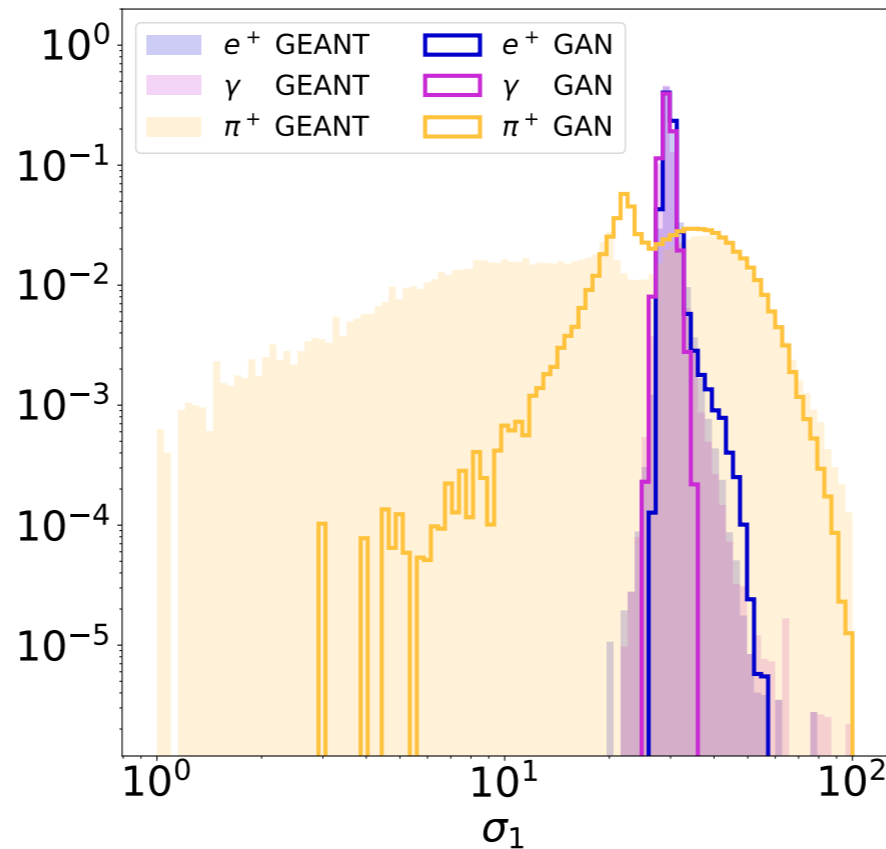


N.B. can always add these (and others) explicitly to the training

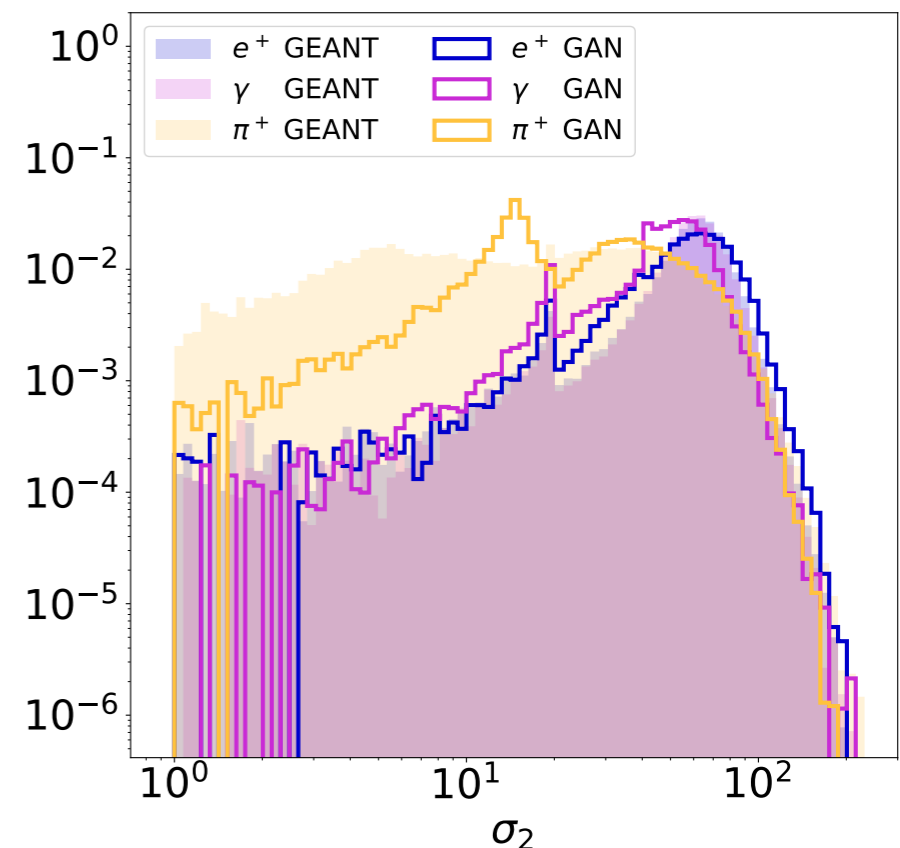


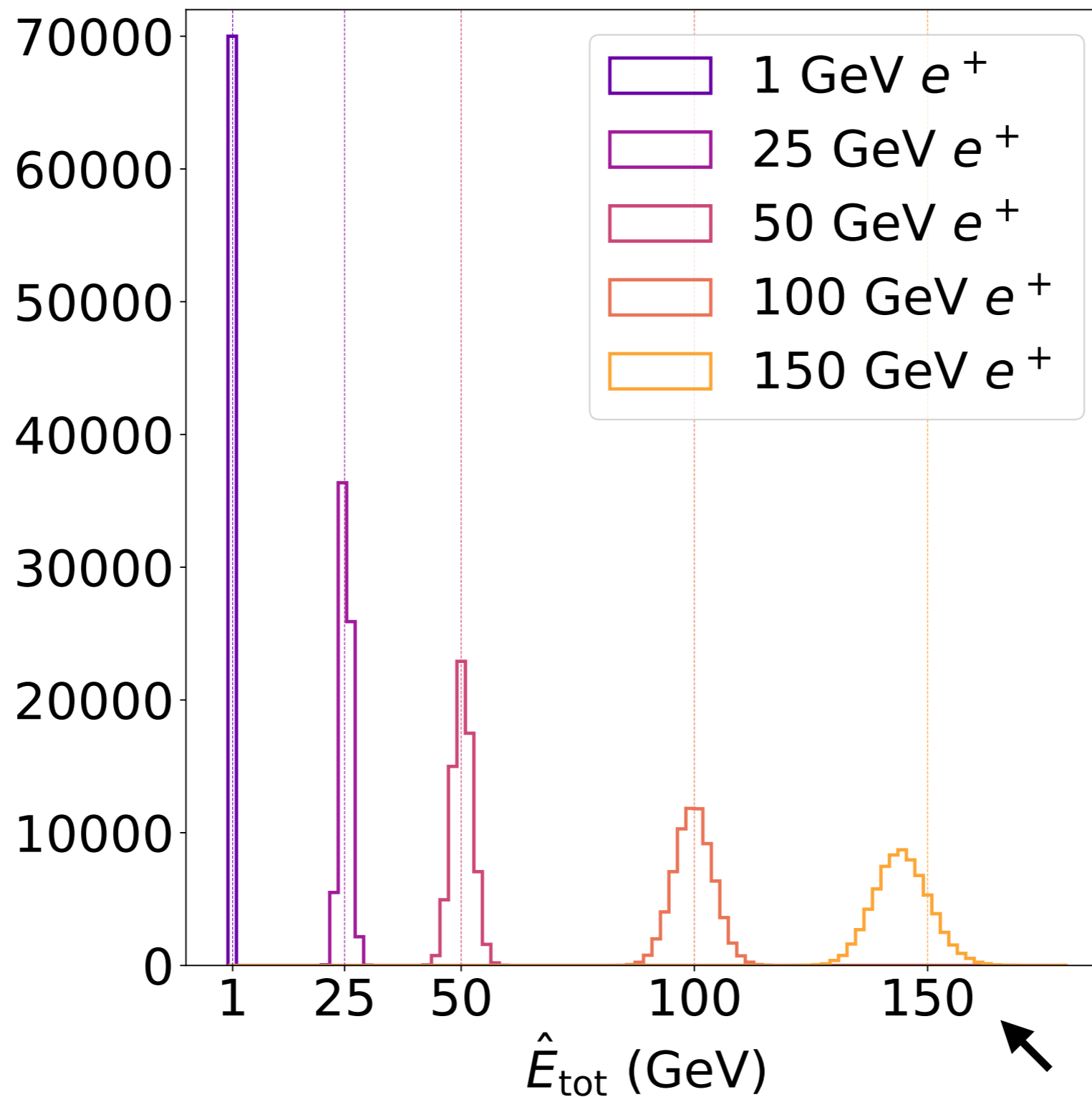


The much larger variation in the pion showers is a challenge for the network.



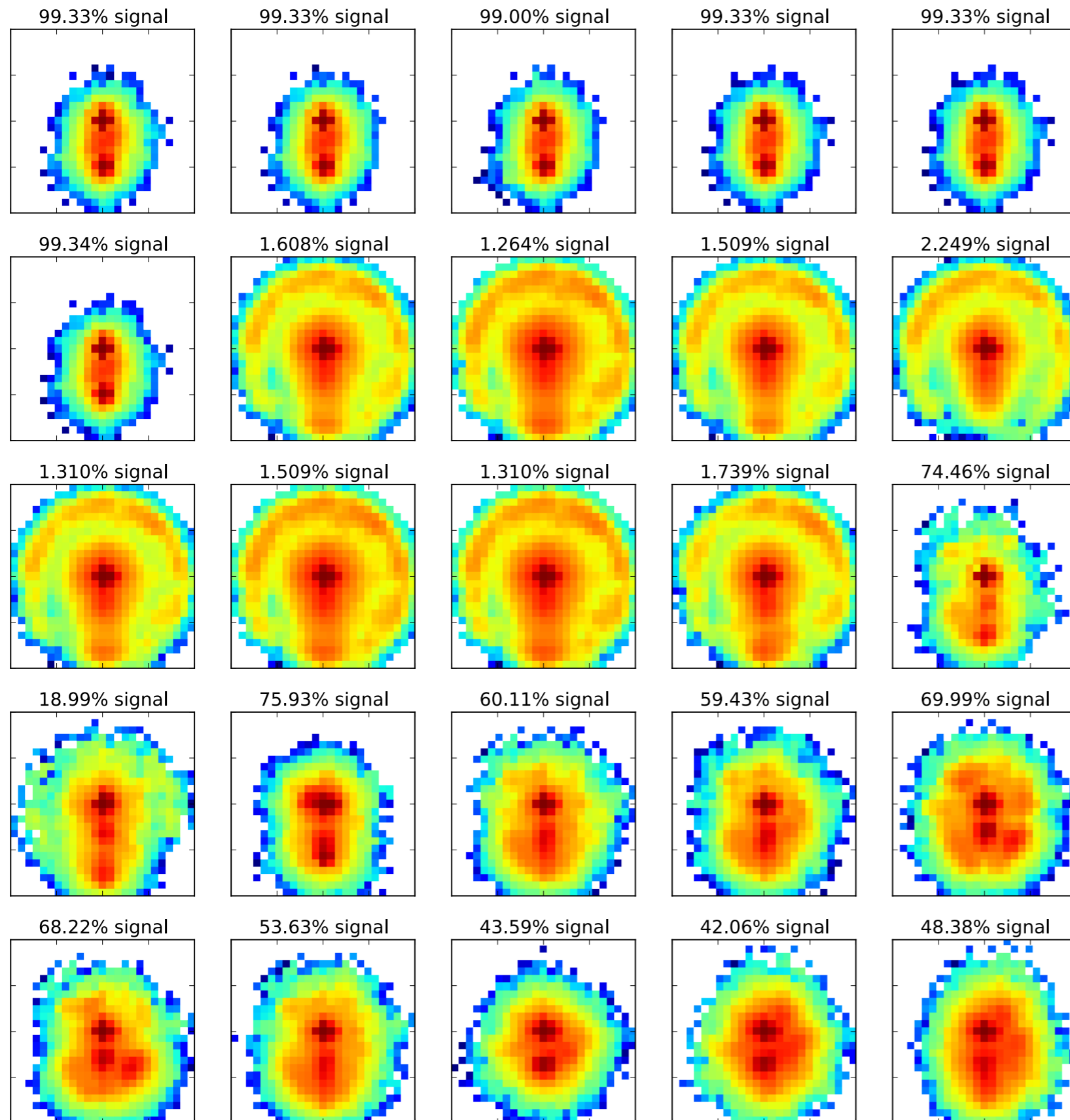
These moments and others are useful for classification; we have also tested this as a metric (NN on 3D images)





← Beyond our training sample!

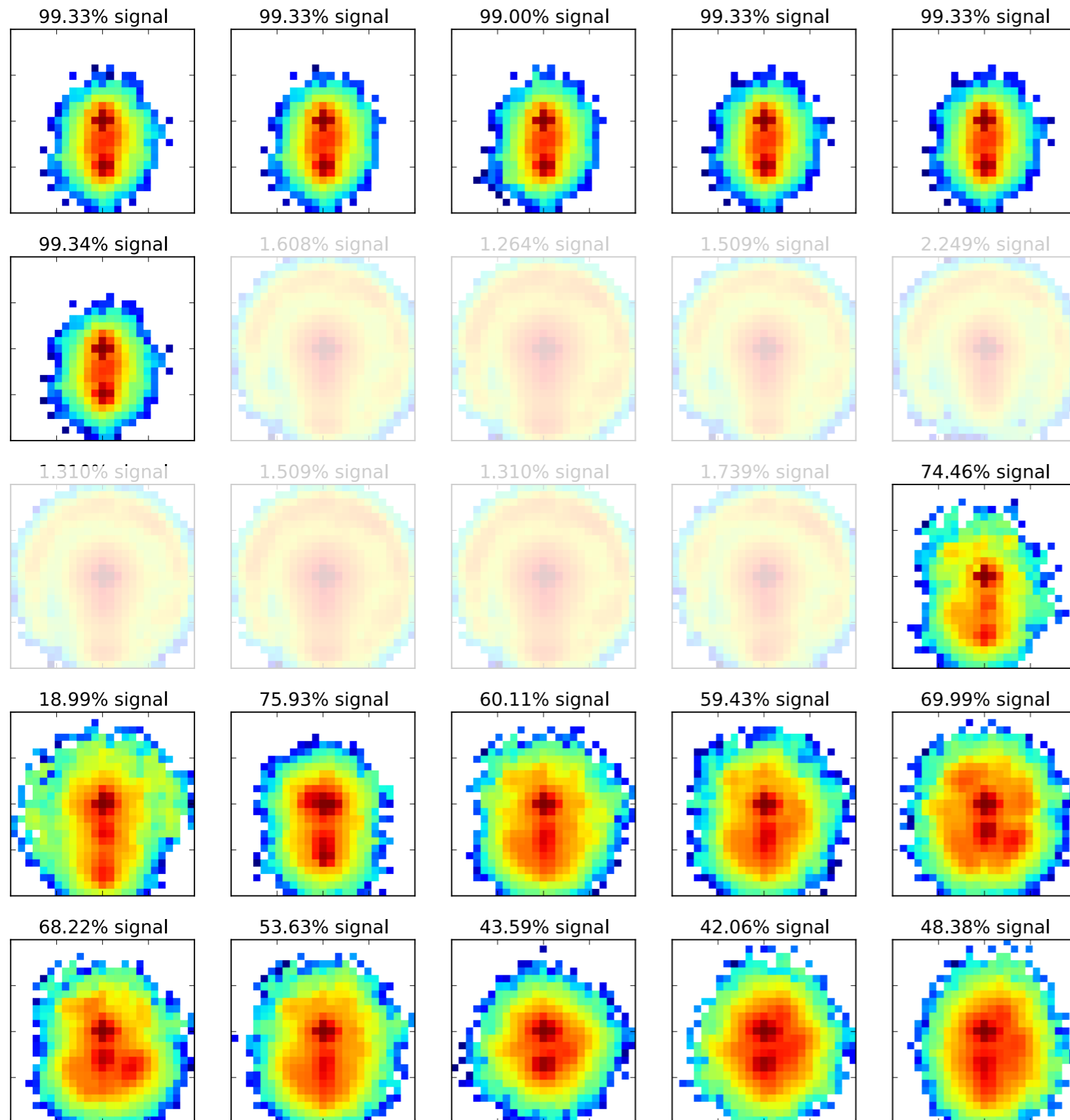
Most activating images



Take a node in the NN and ask which input images activate it the most

Some nodes learn about subjects and some learn about peripheral radiation

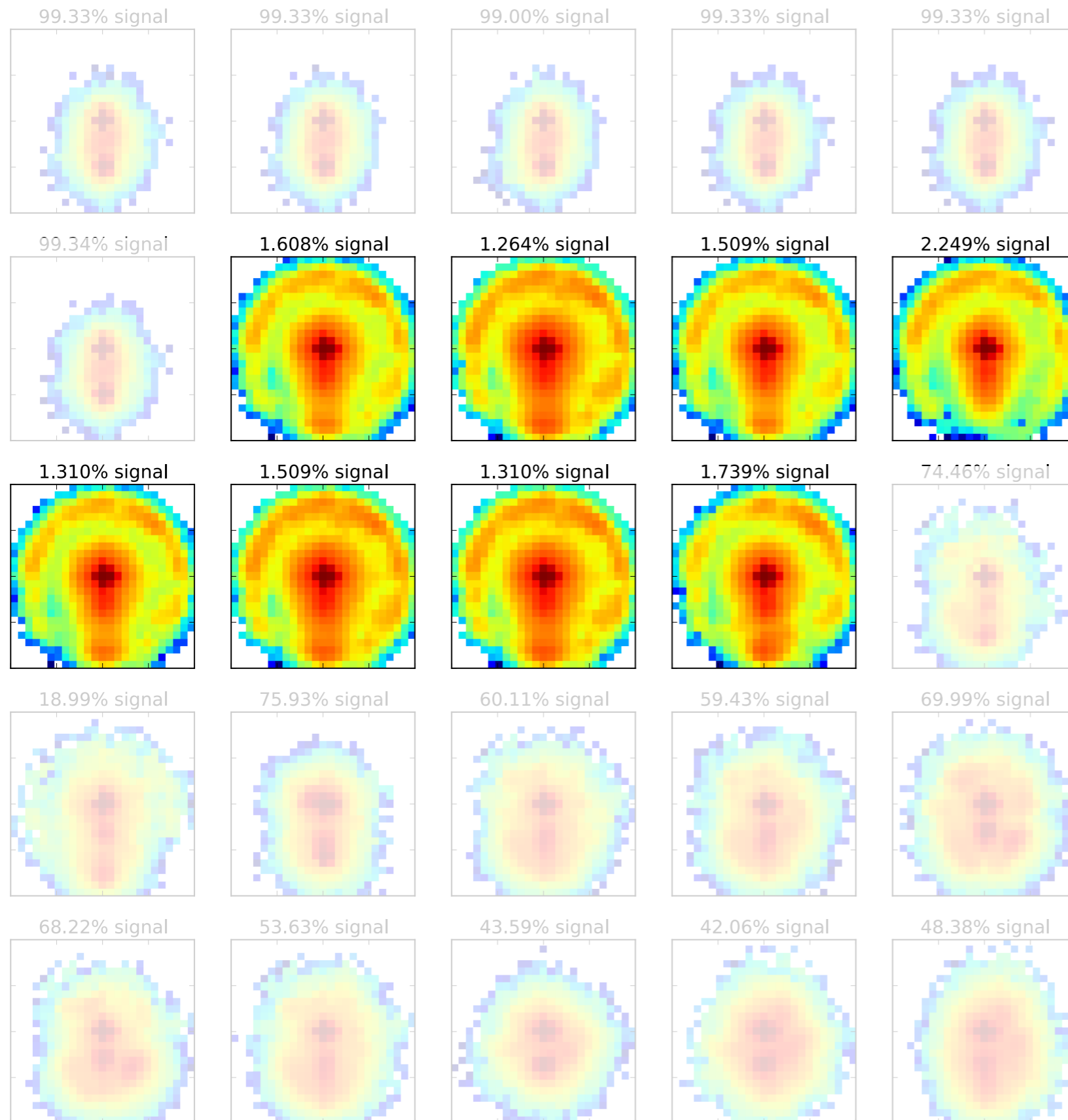
Most activating images



Take a node in the NN and ask which input images activate it the most

Some nodes learn about subjects and some learn about peripheral radiation

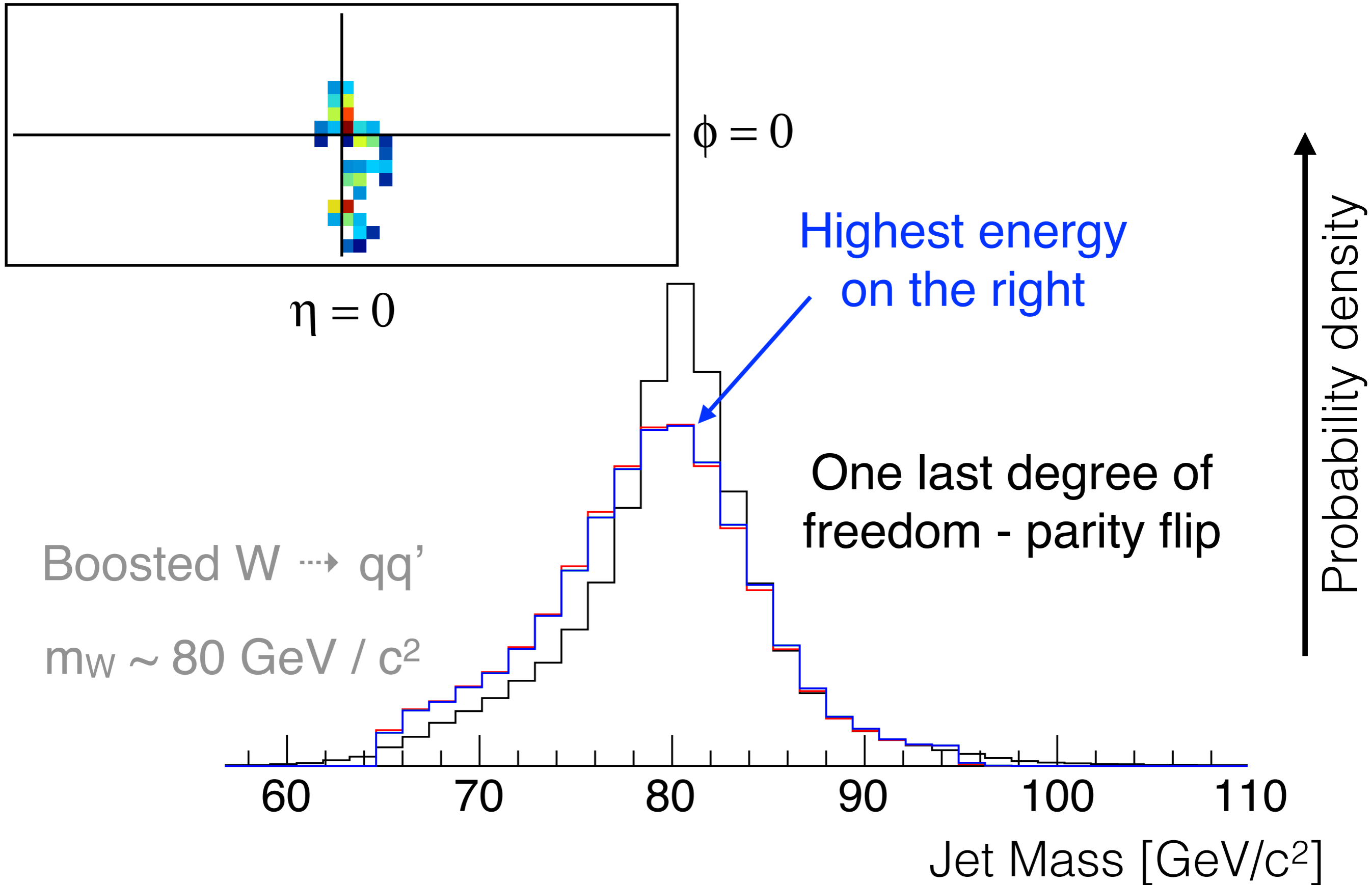
Most activating images



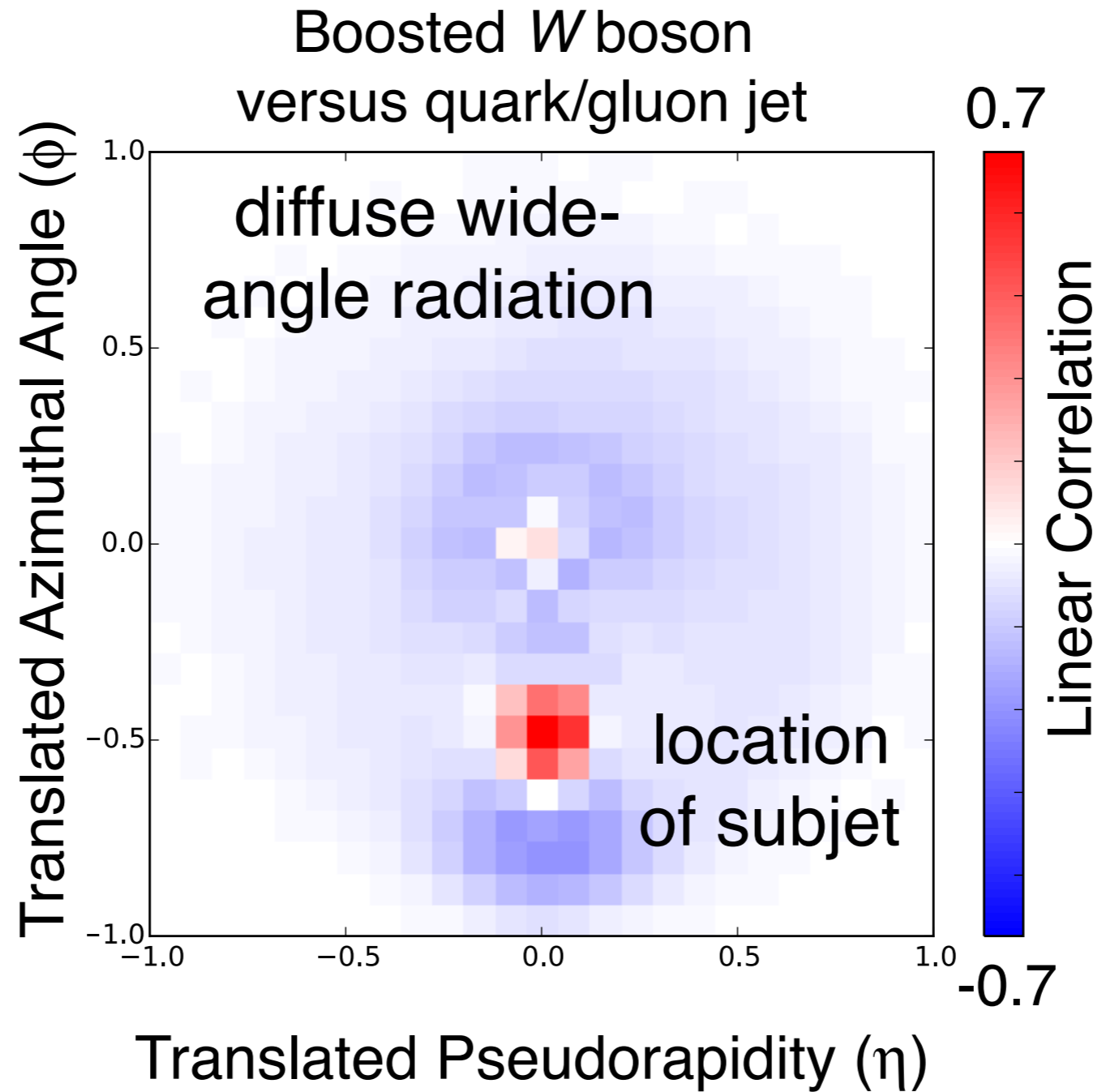
Take a node in the NN and ask which input images activate it the most

Some nodes learn about subjects and some learn about peripheral radiation

Pre-processing & spacetime symmetries



Correlation between input and output



Red = network is more activated (more signal-like)