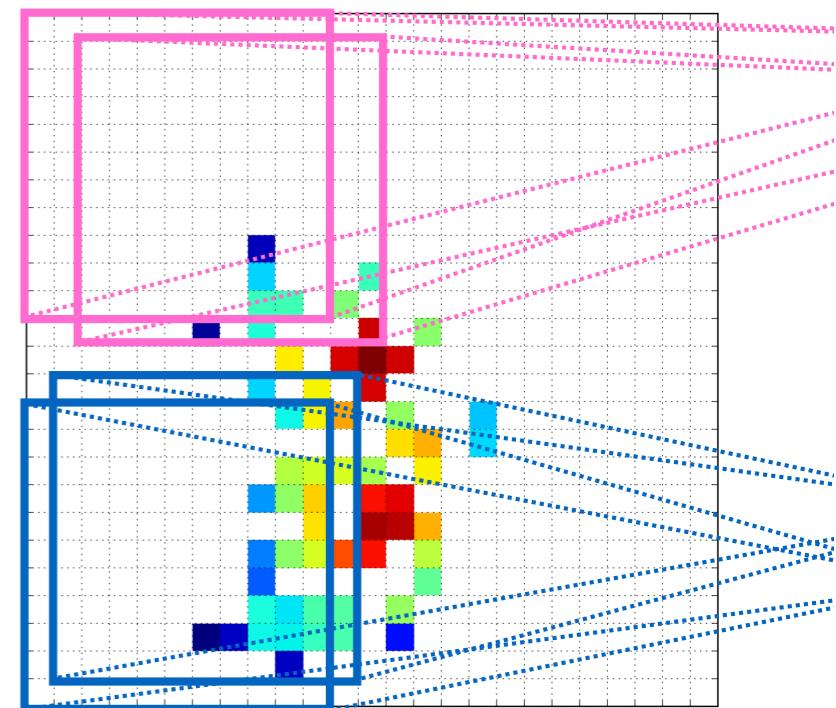


Machine Learning for Jets

2019 - Fermilab edition

Benjamin Nachman

Lawrence Berkeley National Laboratory



January 23, 2019

What are jets?

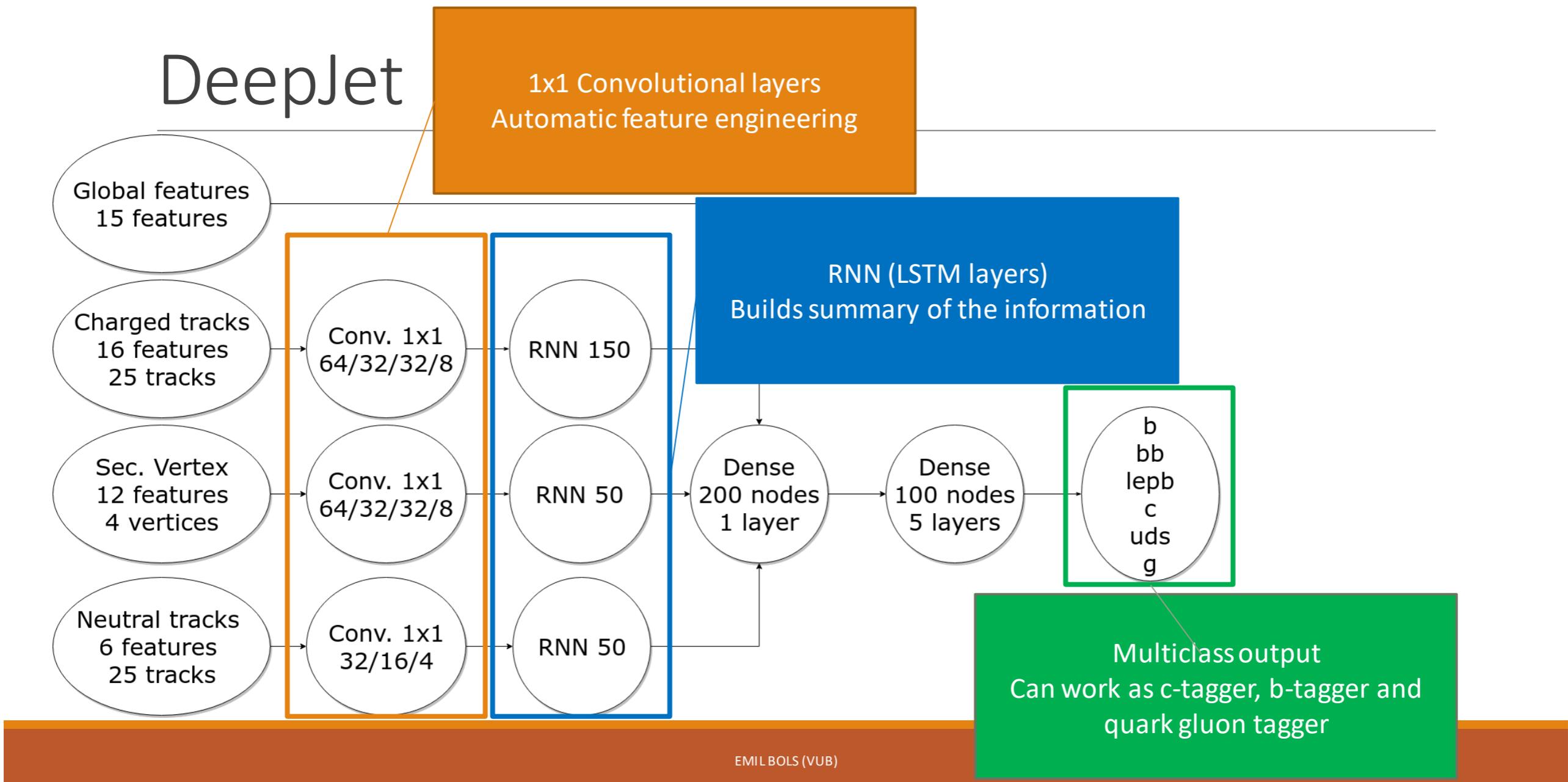


ML4Jets

3



Flavor tagging in CMS

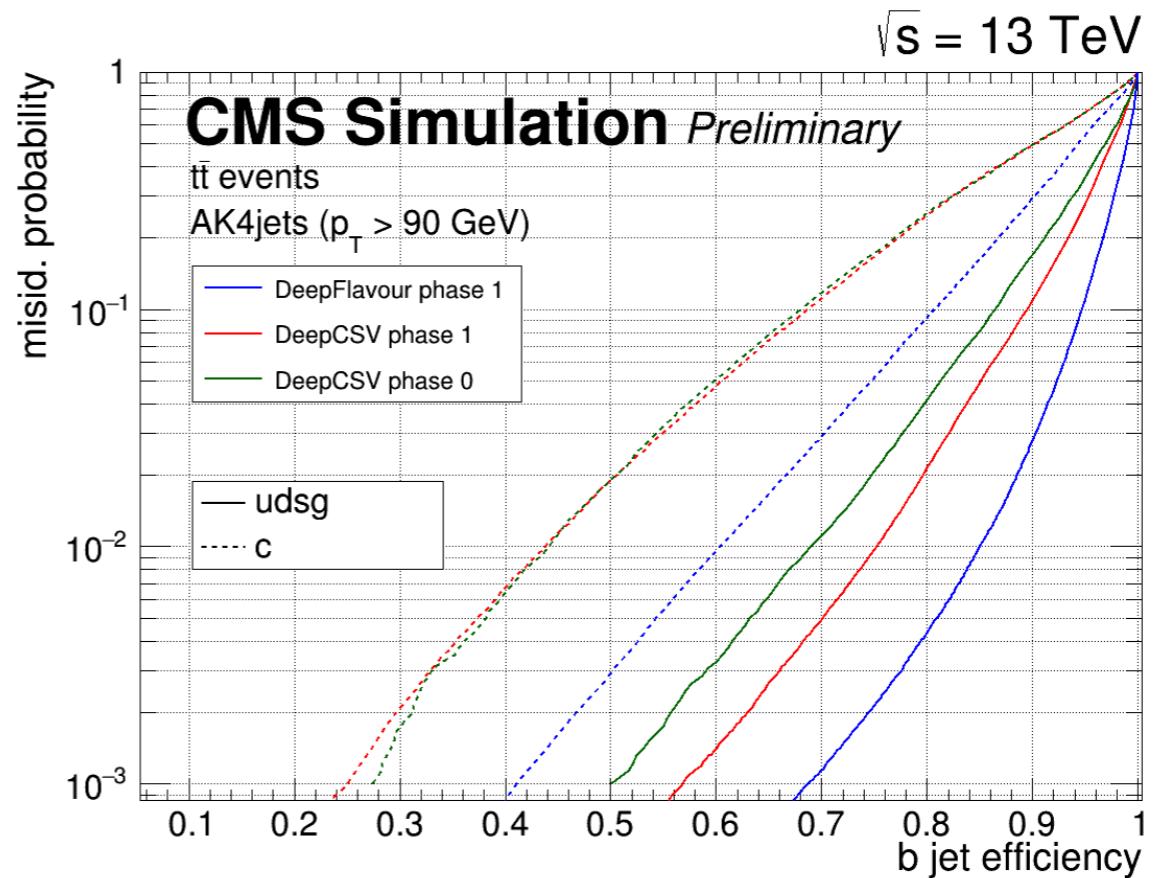


Flavor tagging in CMS

DeepJet

CMS DP-2018/033

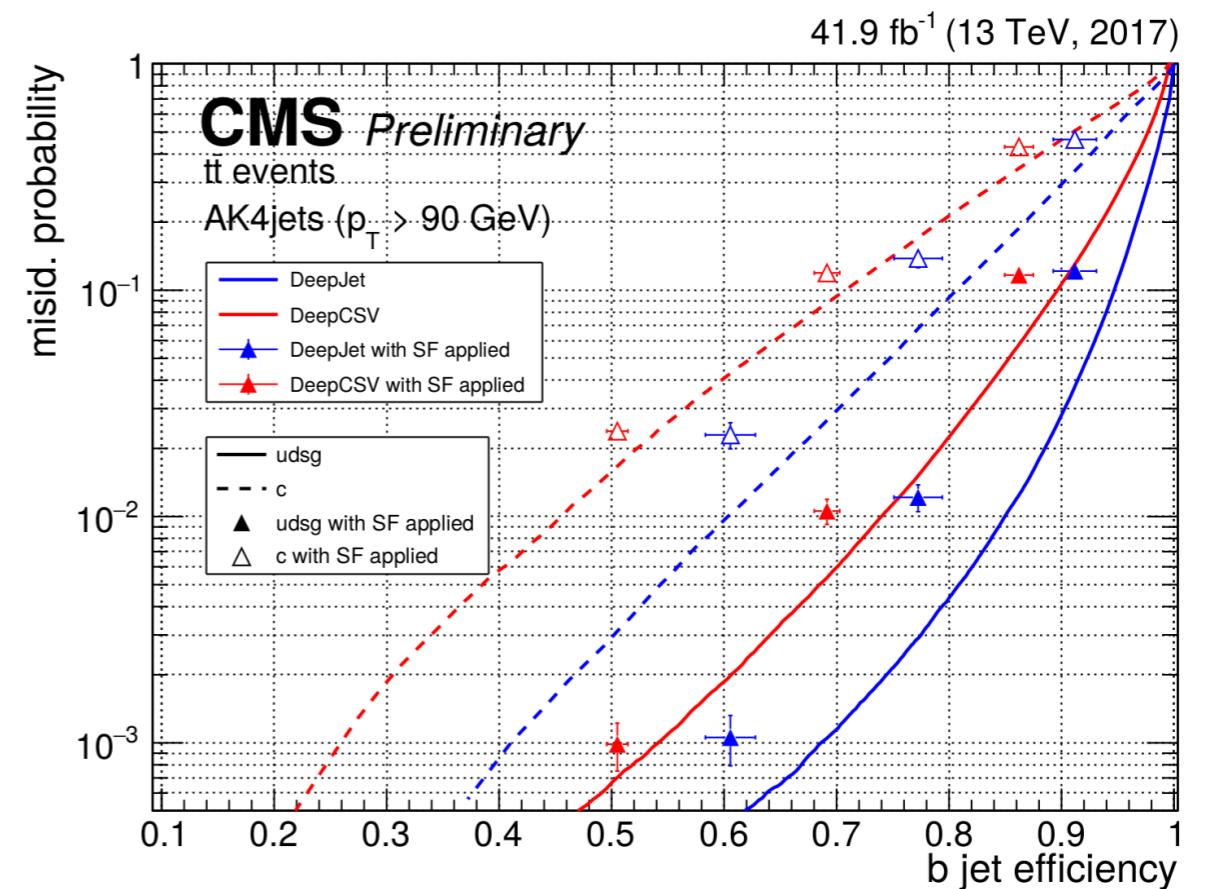
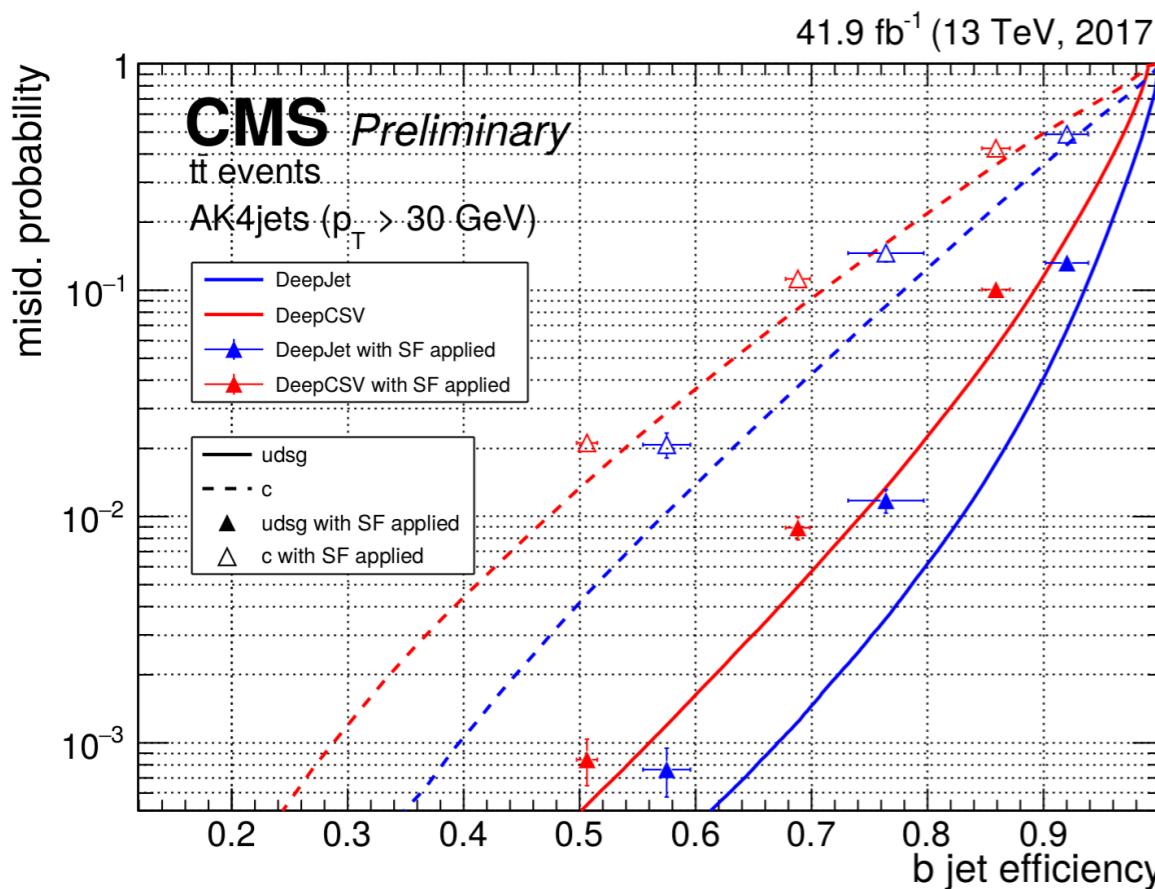
- Large improvement in performance in simulation
- However simulation is never perfect
- Will this gain translate into data?



Flavor tagging in CMS

DeepJet

CMS DP-2018/058



Deep Sets

Deep Sets for Particle Jets

[PTK, Metodiev, Thaler, [1810.05165](#)]

Particle Flow Network (PFN)

$$\text{PFN}(\{p_1^\mu, \dots, p_M^\mu\}) = F \left(\sum_{i=1}^M \Phi(p_i^\mu) \right)$$

Fully general latent space

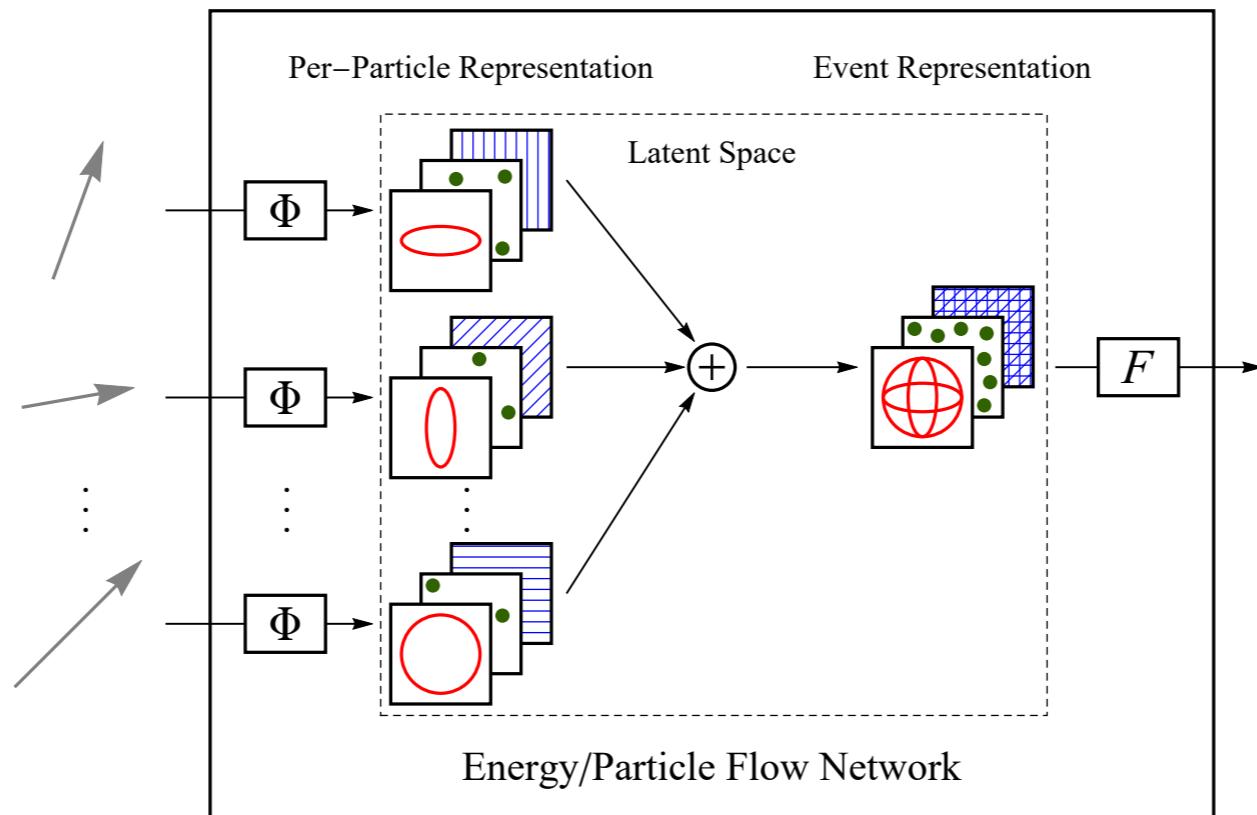
Energy Flow Network (EFN)

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

IRC-safe latent space

Particles

Observable



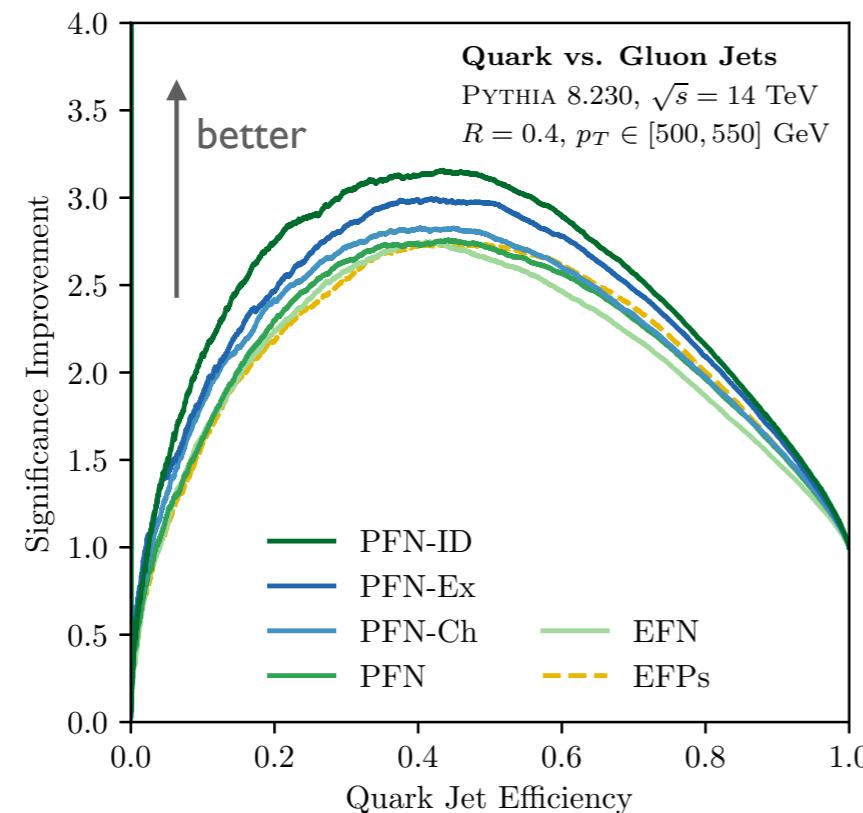
Deep Sets

Classification Performance

PFN-ID: Full particle flavor info
 $(\gamma, \pi^\pm, K^\pm, K_L, p, \bar{p}, n, \bar{n}, e^\pm, \mu^\pm)$

PFN-Ex: Experimentally accessible info
 $(\gamma, h^{\pm,0}, e^\pm, \mu^\pm)$

PFN-Ch: Particle charge info
 $(+, 0, -)$

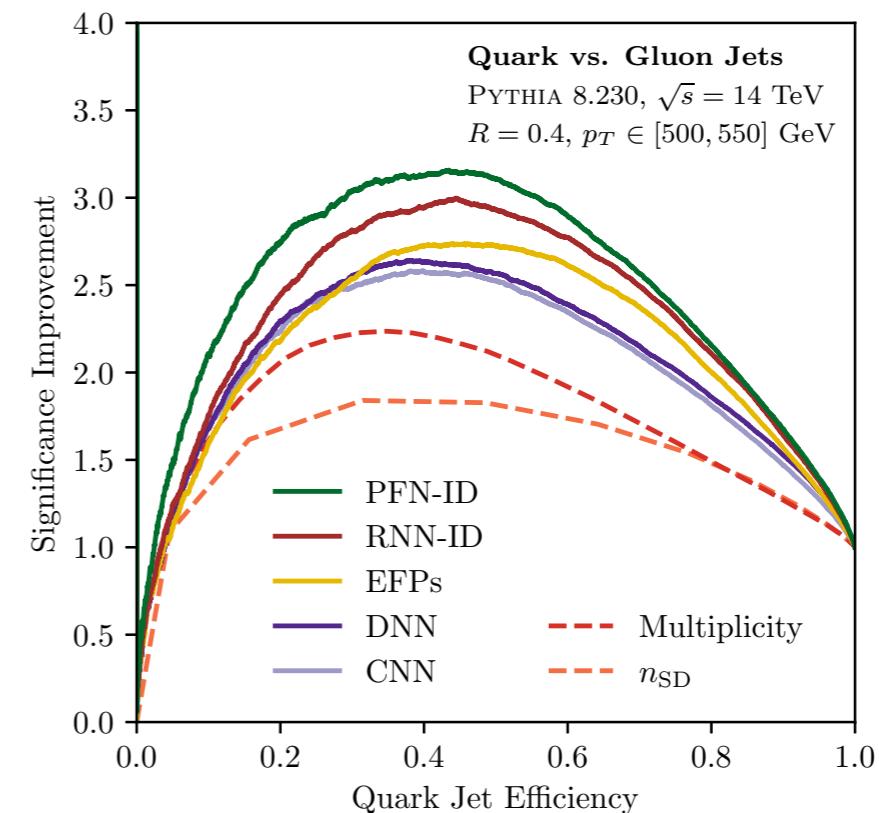


Latent space dimension $\ell = 256$

EFPs are comparable to EFN

PFN: No particle type info, arbitrary energy dependence

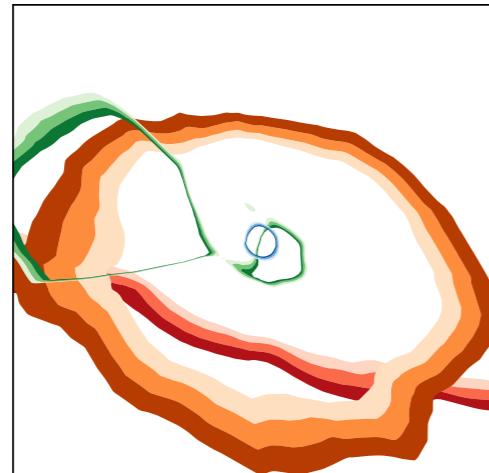
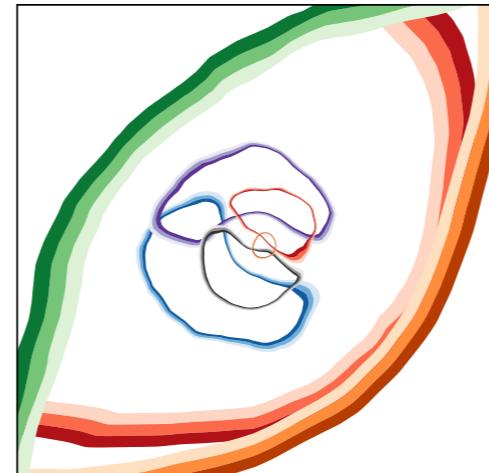
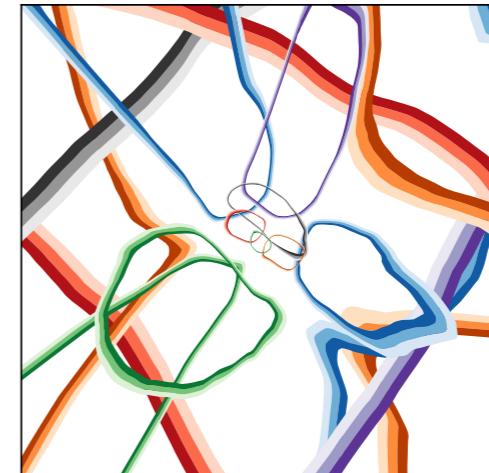
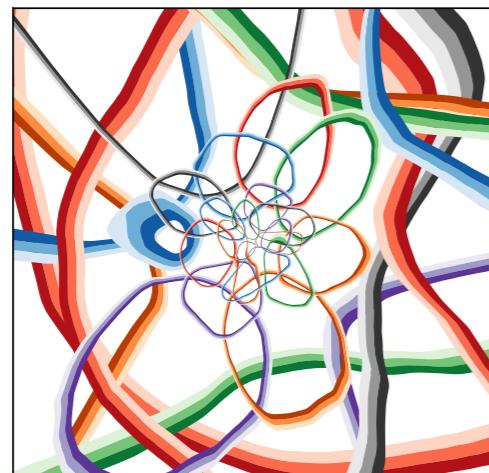
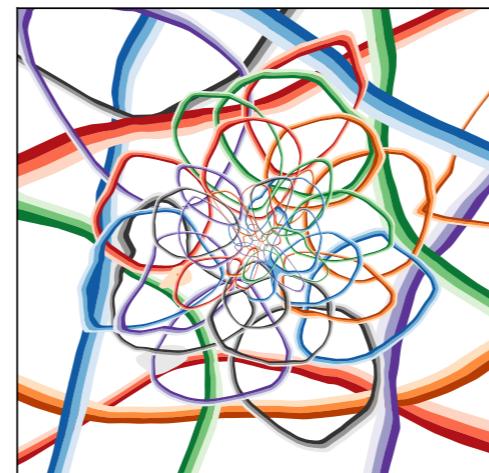
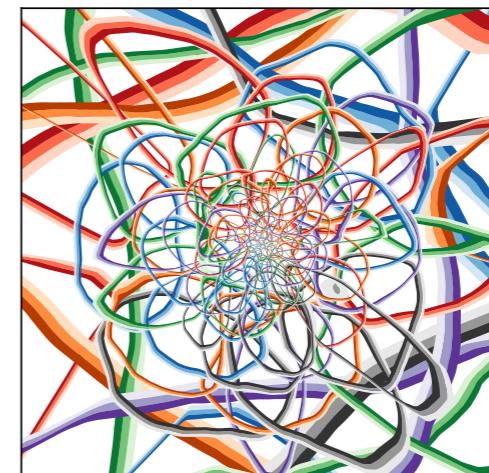
EFN: **IRC**-safe latent space



PFN-ID slightly better than RNN-ID

Deep Sets

Visualizing Q/G EFN Filters

 $\ell = 4$  $\ell = 8$  $\ell = 16$  $\ell = 32$  $\ell = 64$  $\ell = 128$

Deep Sets

Measuring Q/G EFN Filters

Power-law dependence between filter size and distance from center is observed

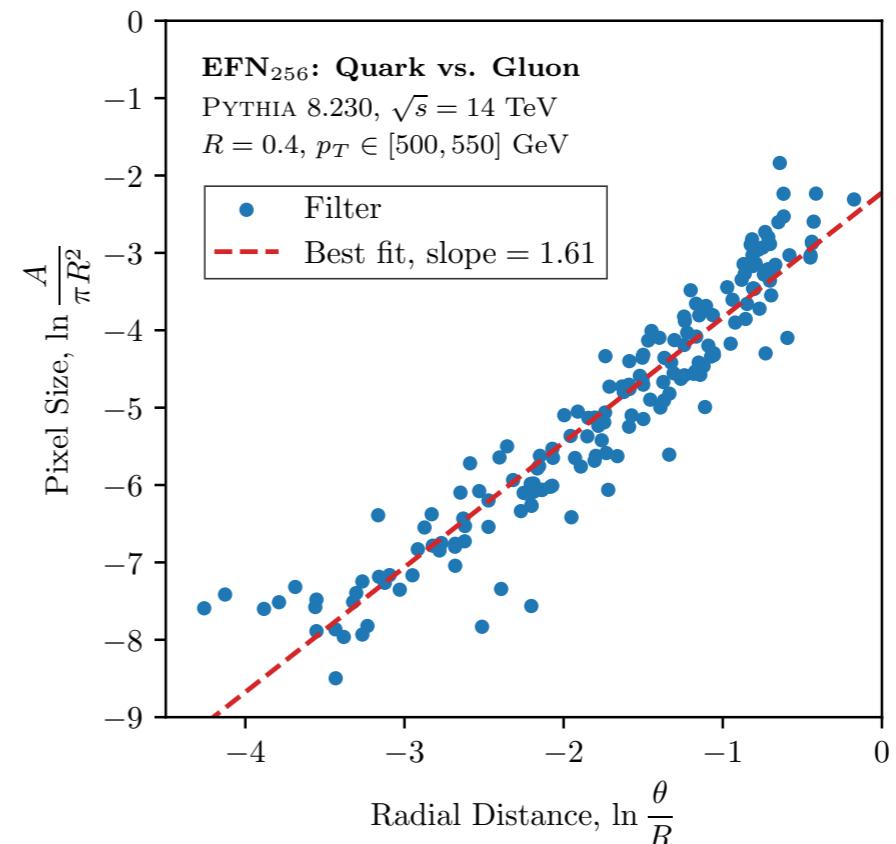
Slope of 2 is predicted at leading log

$$\left| d \ln \frac{\theta}{R} d\varphi \right| = \theta^2 |dy d\phi|$$

Emission plane area element

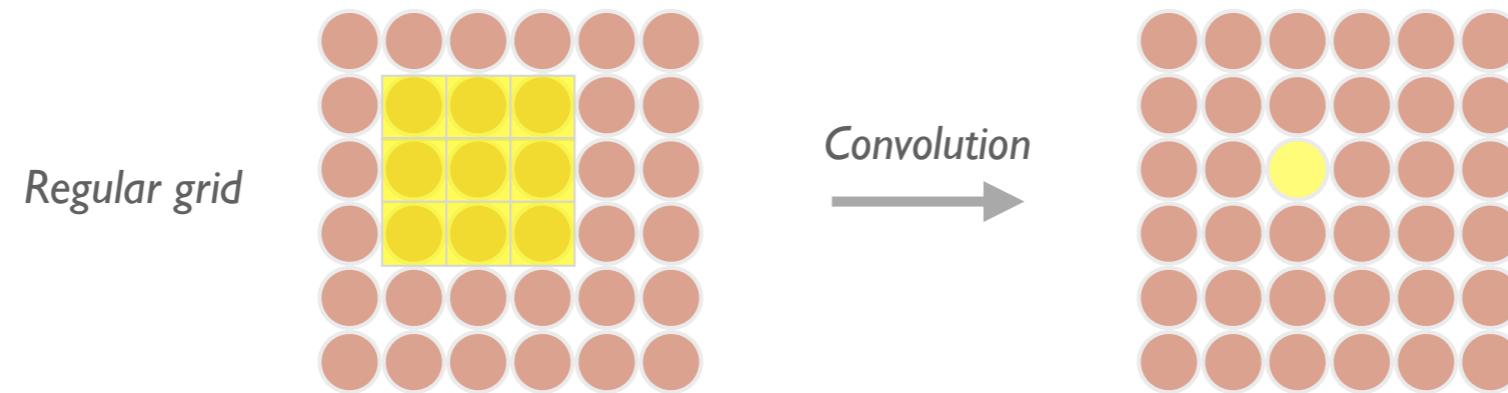
Area element in rap-phi plane

Non-perturbative physics, axis recoil, higher order effects cause deviations from slope of 2

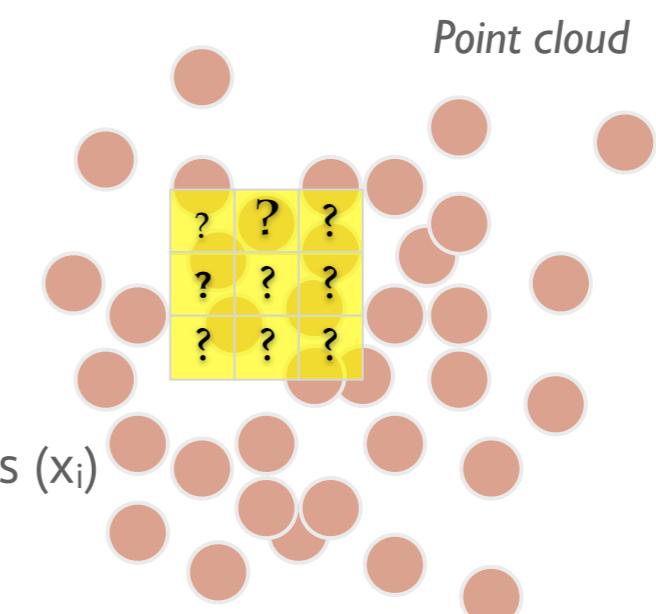


Another particle cloud approach

CONVOLUTION ON REGULAR GRIDS



- Conventional convolution only operates on regular grids and cannot be applied on point clouds
 - point clouds are **irregular**
 - how to define a “local” patch to convolve?
 - point clouds are **unordered**
 - conventional convolution operation ($\sum_i K_i x_i$) is not invariant under permutation of the points (x_i)



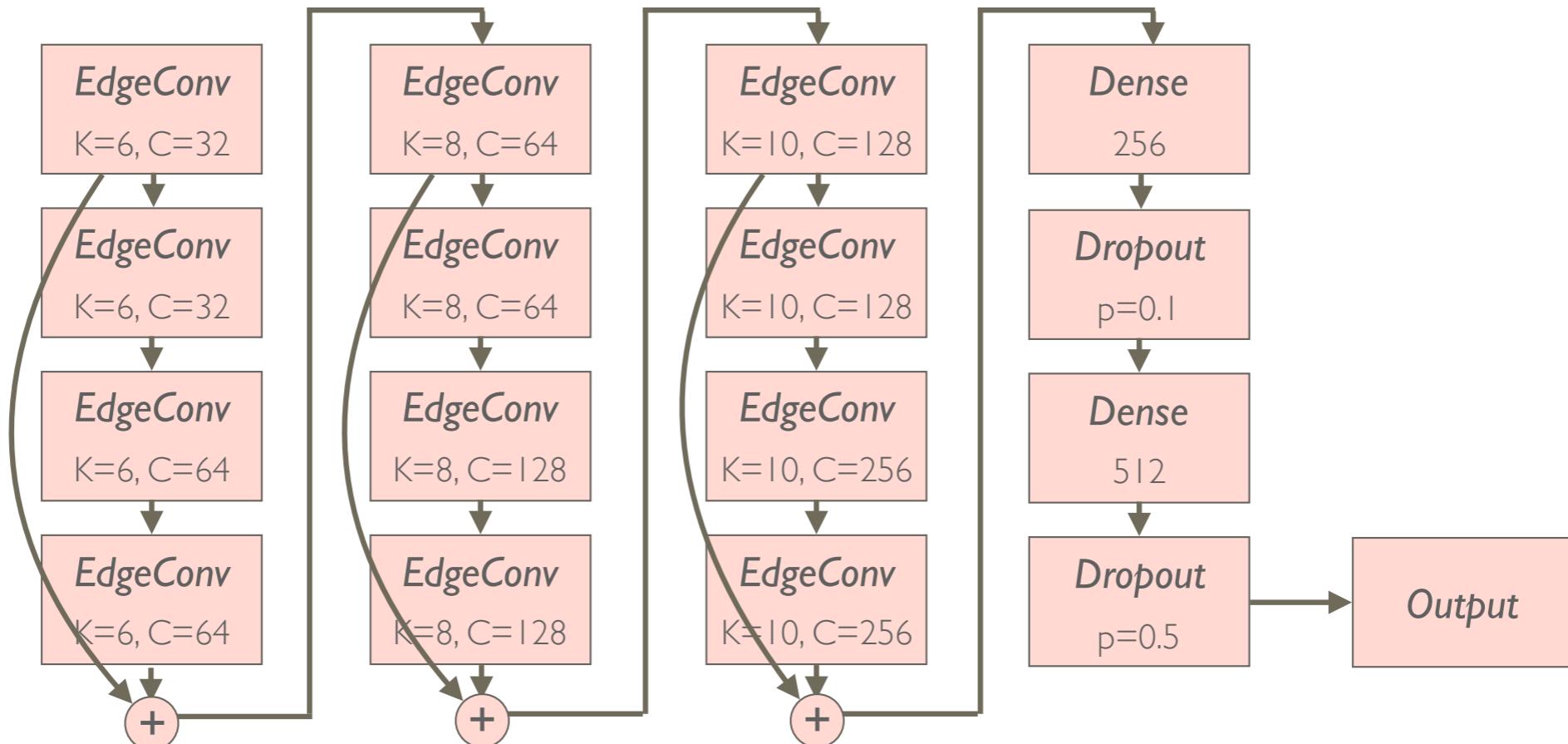
Another particle cloud approach

NETWORK ARCHITECTURE

- Implemented with

mxnet

Jet as a Particle Cloud - November 15, 2018 - Huilin Qu (UCSB)

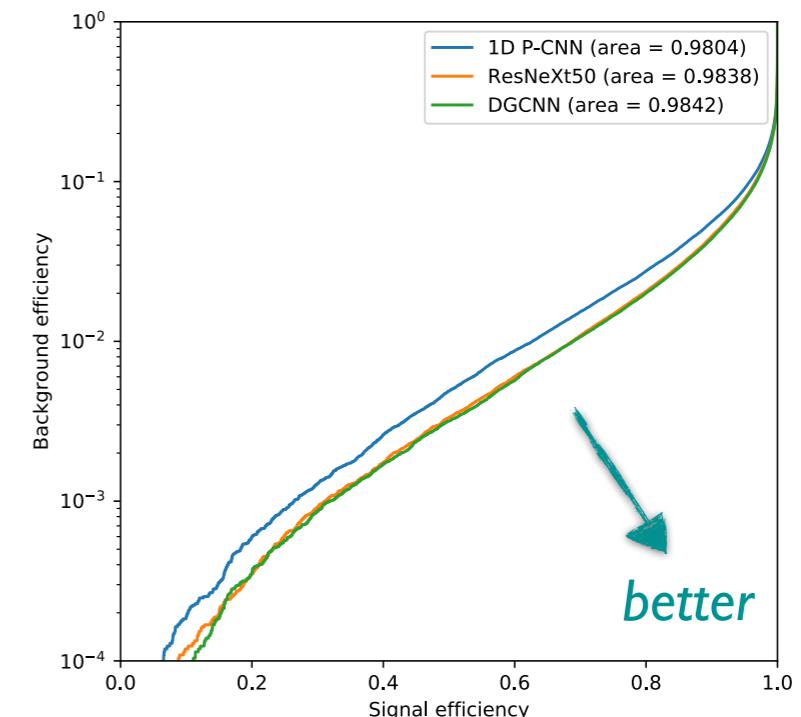


- 3 stages: k-nearest neighbors updated at the beginning of each stage
- batch normalization (BN) used after each EdgeConv operation
- residual connection (RC) [1512.03385, 1603.05027] added between EdgeConv layers
- BN and RC helped greatly for stabilizing the training and also improving the performance

Another particle cloud approach

PERFORMANCE: TOP TAGGING

- Top tagging:
 - only particle 4-momentum is available
 - train/val/test 1.2M/400k/400k
 - results of more algorithms available at [link](#)
- We managed to push the boundary a bit further
 - >20% lower background at signal efficiency of 30%
 - *Is the gain real? Or is it just learning more details of the parton shower model?*
 - but personally I would not really consider the jet tagging problem as “solved”
 - especially facing realistic experimental challenges like pileup, detector effects, and additional information (e.g., tracking, timing, etc.)

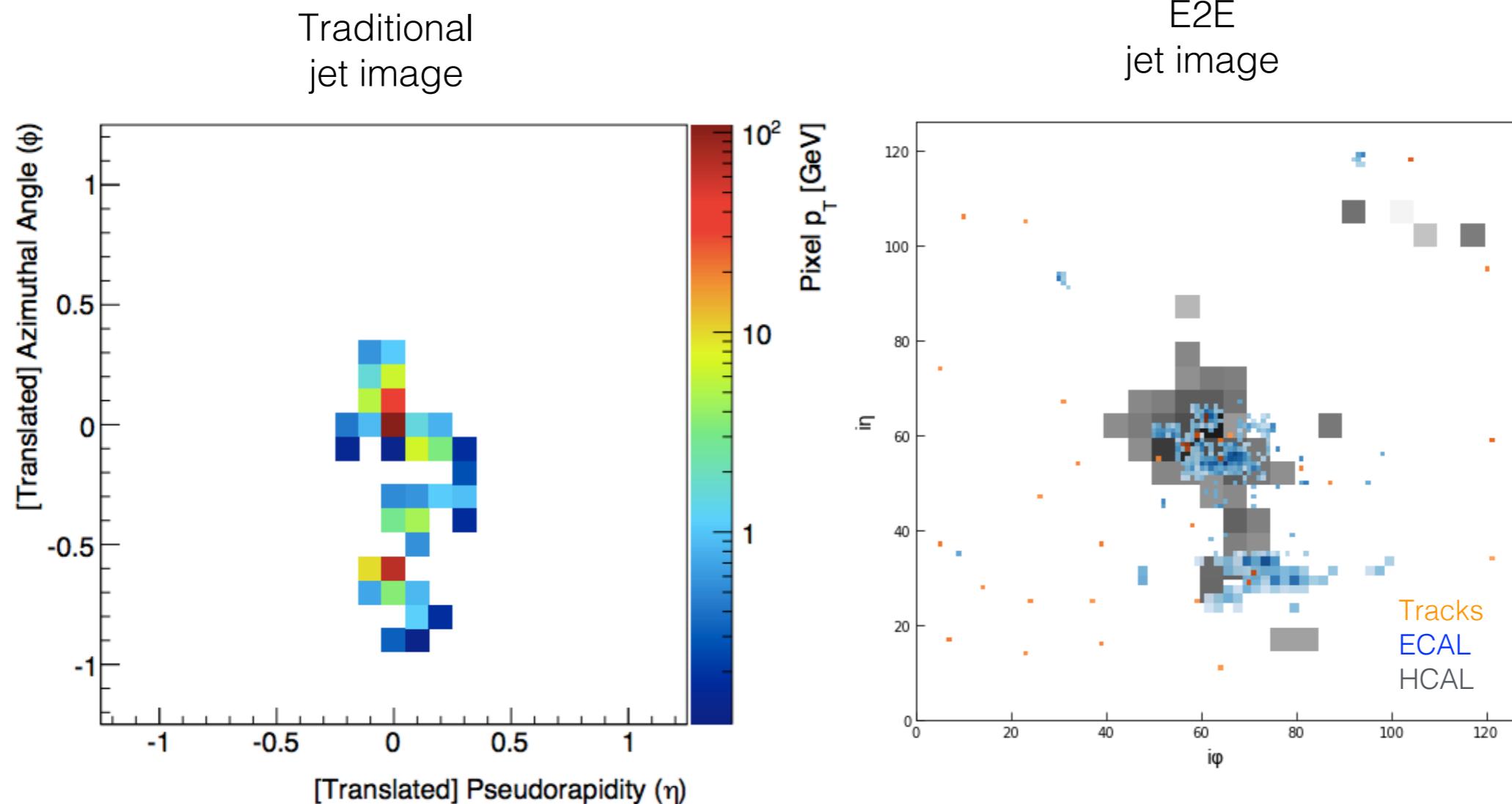


Performance on top-tagging dataset

Algorithm	Accuracy	ROC AUC	$1/\epsilon_{\text{bkg}} @ \epsilon_{\text{sig}}=30\%$
1D P-CNN	0.930	0.9804	780
2D CNN [ResNeXt50]	0.936	0.9838	1086
DGCNN	0.937	0.9842	1160
<i>PFN-r.r. [arXiv:1810.05165]</i>	0.932	0.9819 ± 0.0001	888 ± 17

Raw image processing

Jet ID | Traditional vs E2E Image



Note: Not the same jet.

Prior-independent regression

Introduction: Numerical Inversion

SLAC

- Calibration with *numerical inversion*:
 - Train on jets with known energies x , take output value with maximum likelihood: $f(x) = \langle p_T^{reco} | p_T^{true} = x \rangle$
 - Invert function for unknown observed values: $p_T^{reco} \rightarrow f^{-1}(p_T^{reco})$
-
- ```

 graph LR
 x1[x1] --> Instrument[Instrument]
 x2[x2] --> Instrument
 x3[x3] --> Instrument
 Instrument --> y1["y1 ± σ1"]
 Instrument --> y2["y2 ± σ2"]
 Instrument --> y3["y3 ± σ3"]
 y1 --> x1
 y2 --> x2
 y3 --> x3

```
- *Independent* of the distribution of true energies used in training
    - Processes in LHC have different true energy distributions!
    - Guarantees closure\* when conditioning on  $p_T^{true} = x$
  - *This is* what we do in ATLAS right now for MCJES and GSC
  - \*Formal details: [arXiv:1609.05195](https://arxiv.org/abs/1609.05195)

# Prior-independent regression

## Generalized Numerical Inversion

SLAC

- Generalized numerical inversion:

- $L(x; \theta) = \langle p_T^{\text{reco}} | p_T^{\text{true}} = x; \theta \rangle$
- Calibration:  $p_T^{\text{reco}} \rightarrow C(p_T^{\text{reco}}, \theta) = L^{-1}(p_T^{\text{reco}}, \theta)$



1. Learn a neural network approximation to the function  $L(x, \theta) = \langle p_T^{\text{reco}} | p_T^{\text{true}} = x, \theta \rangle$ . Note that  $L(x, \theta) : \mathbb{R}^{n+1} \rightarrow \mathbb{R}$ .
2. Learn a neural network  $C(L(x, \theta), \theta)$  that tries to predict  $x$  given  $\theta$  and  $L(x, \theta)$ . This is an approximation to the family of functions  $f_\theta^{-1}(x)$ . Note that learning the inverse this way is technically simple since  $L$  is single-valued.
3. Calibrate with  $p_T^{\text{reco}} \mapsto C(p_T^{\text{reco}}, \theta)$ . The calibration non-closure is given the deviation of  $C(p_T^{\text{reco}}, \theta|x)$  from  $x$ .

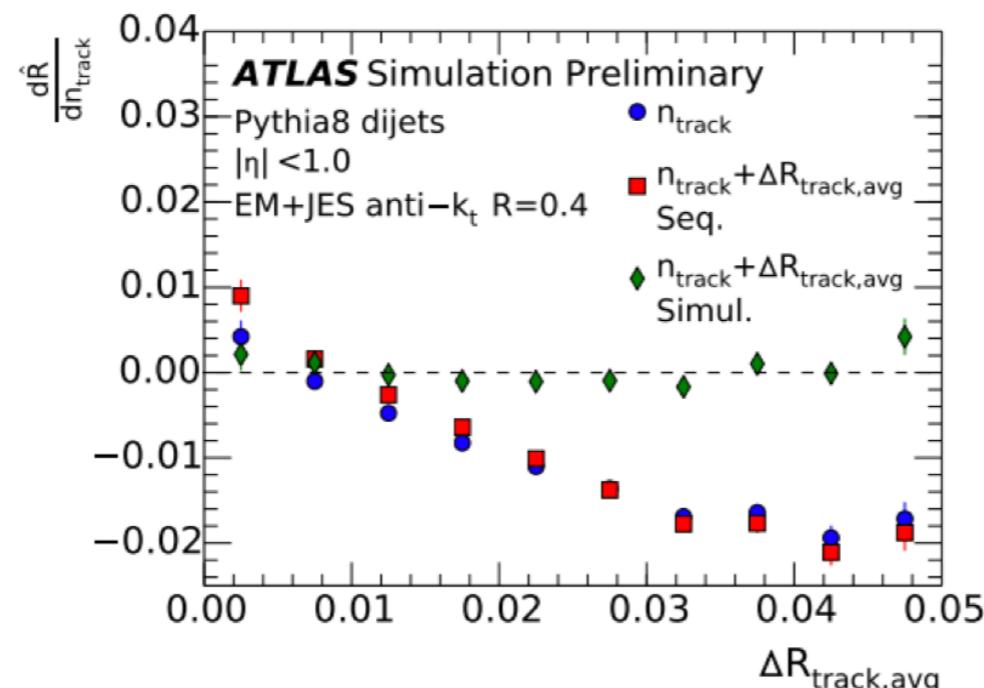
- Advantages:

- Still independent of underlying  $p_T$  distribution
- Easy to add in other features, unbinned
- Can take into account correlations between features

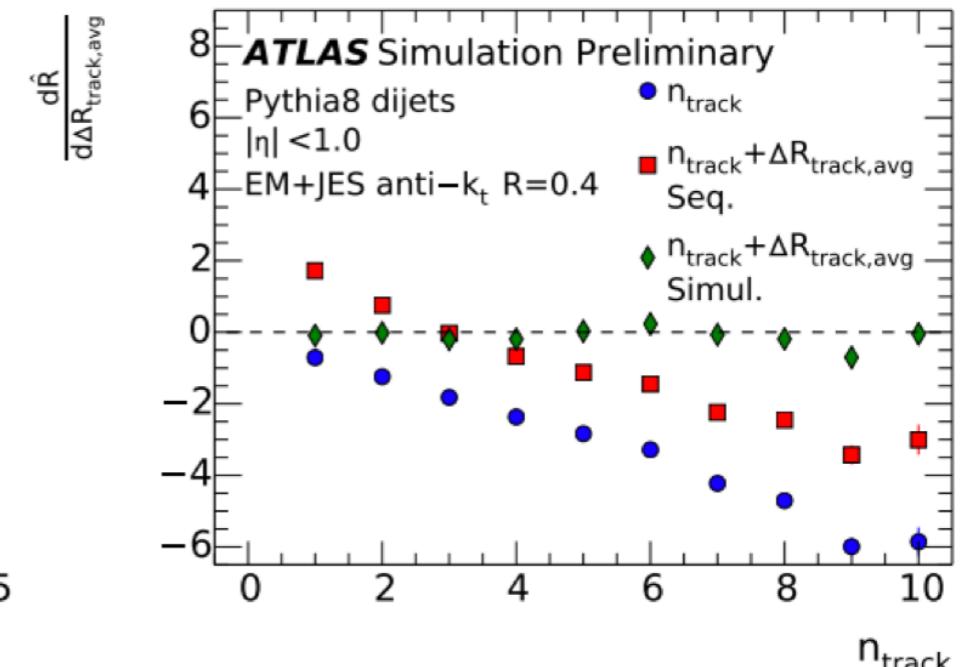
# Prior-independent regression

## Generalized Numerical Inversion – Final Results

- Residual dependence
  - Simultaneous is flat at 0 – demonstration of method



(a)

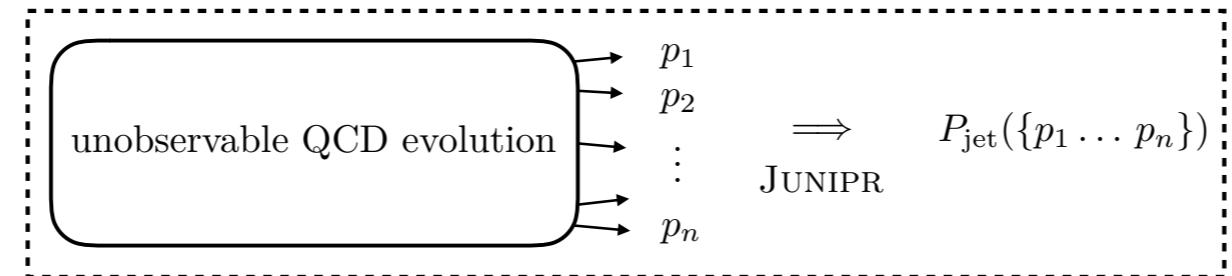


(b)

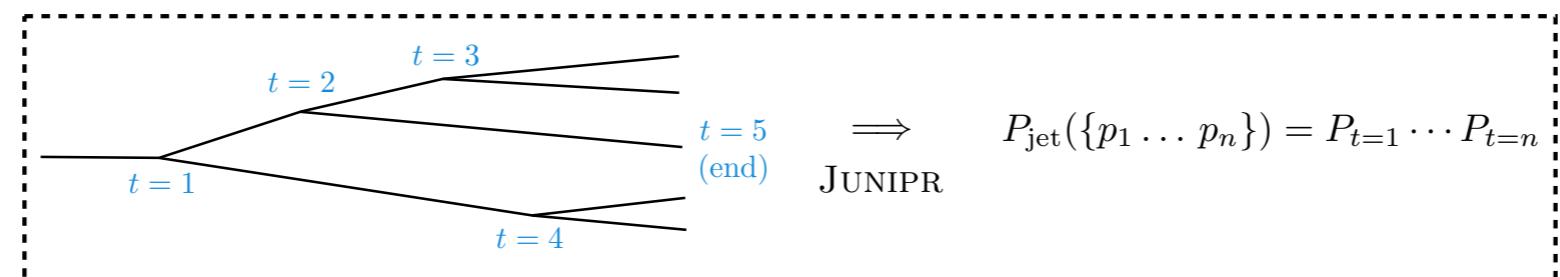
# Physics-inspired generative models

## Summary So Far

JUNIPR computes the probability of a jet...



...as a product over time steps in its clustering tree...



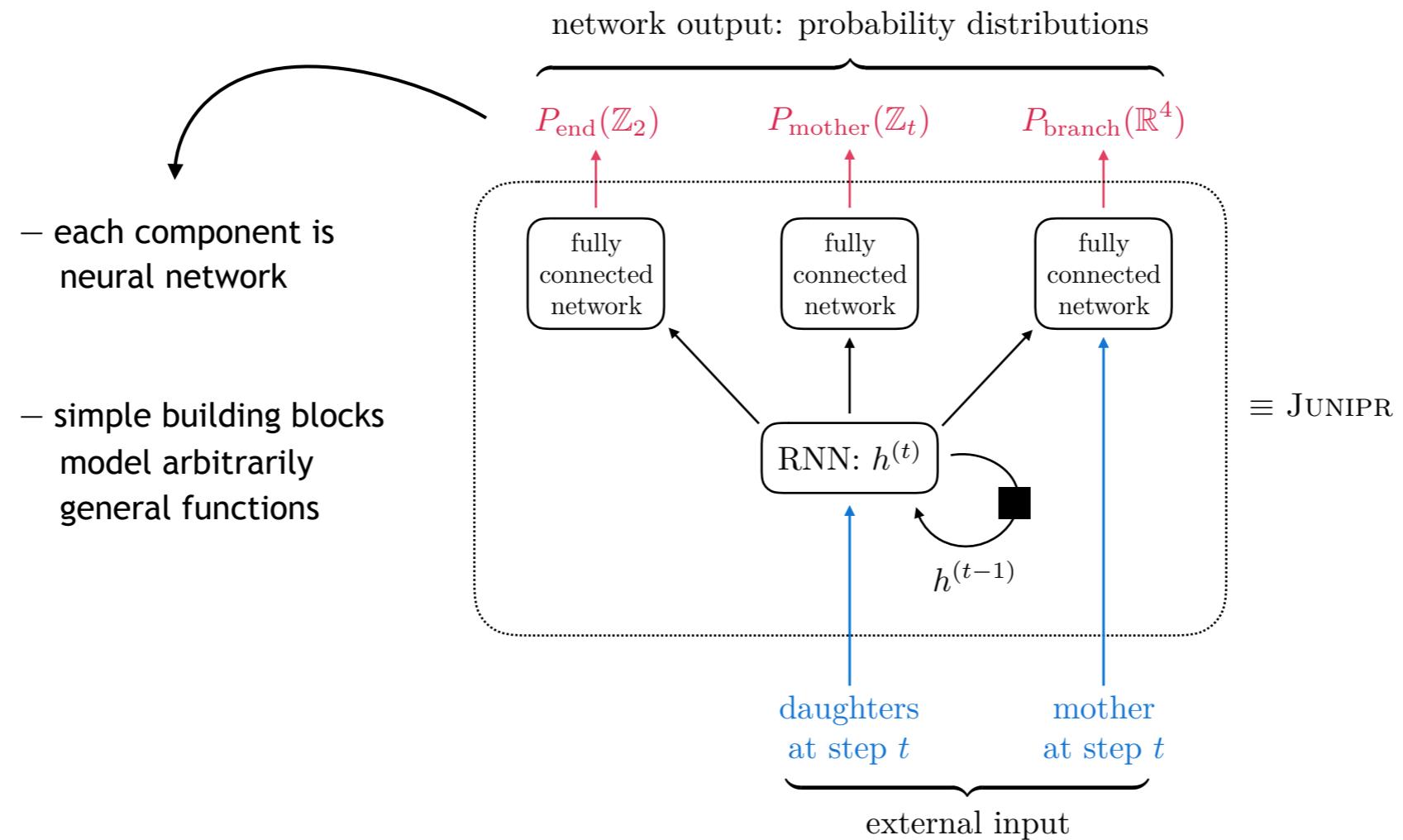
...where each time step is decomposed into 3 parts:

$$P_t = P_{\text{end}} \cdot P_{\text{mother}} \cdot P_{\text{branch}}$$

# Physics-inspired generative models

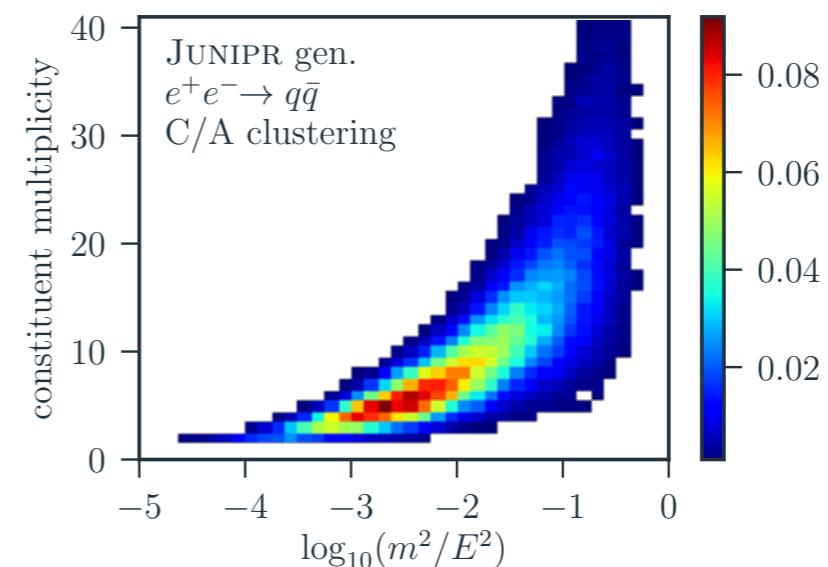
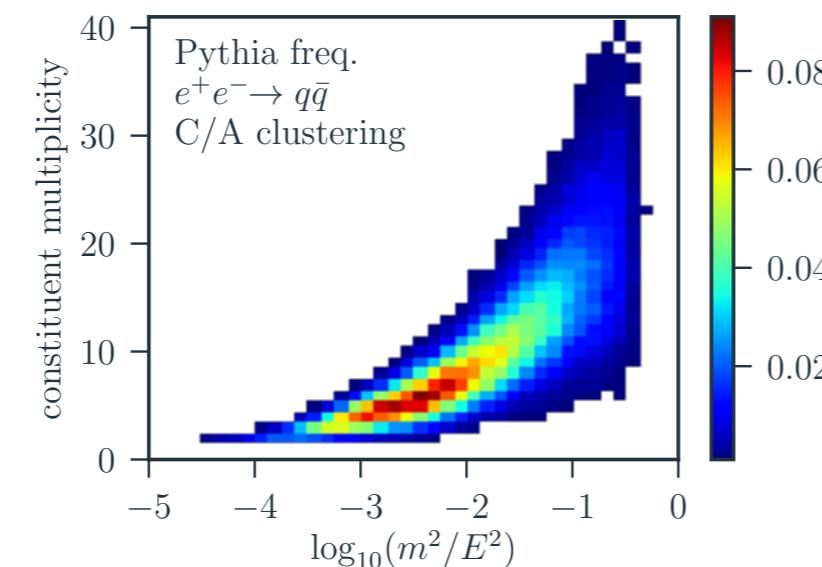
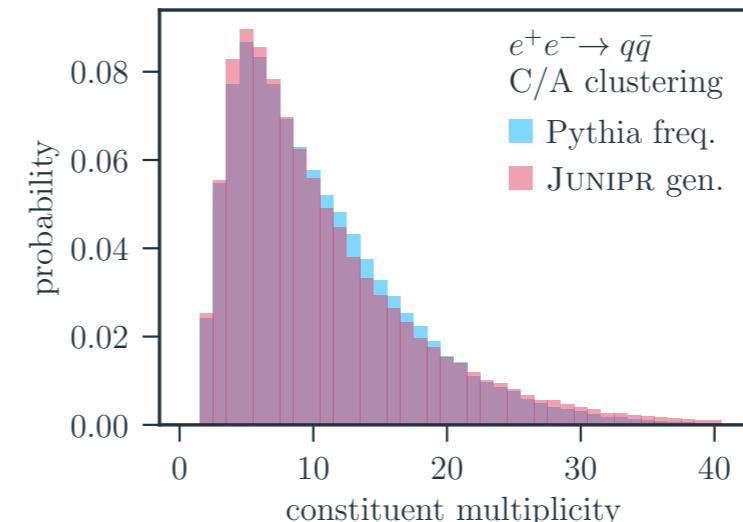
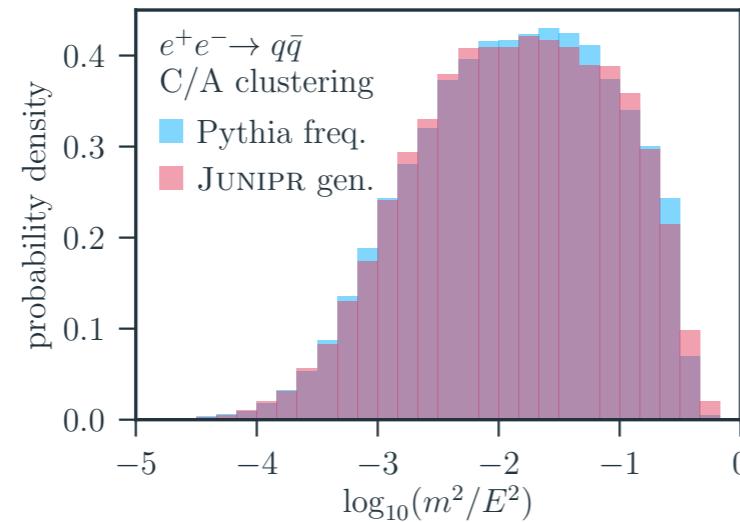
## Implementation with RNN

- STEP 2) feed  $h^{(t)}$  into neural networks computing probability distributions



# Physics-inspired generative models

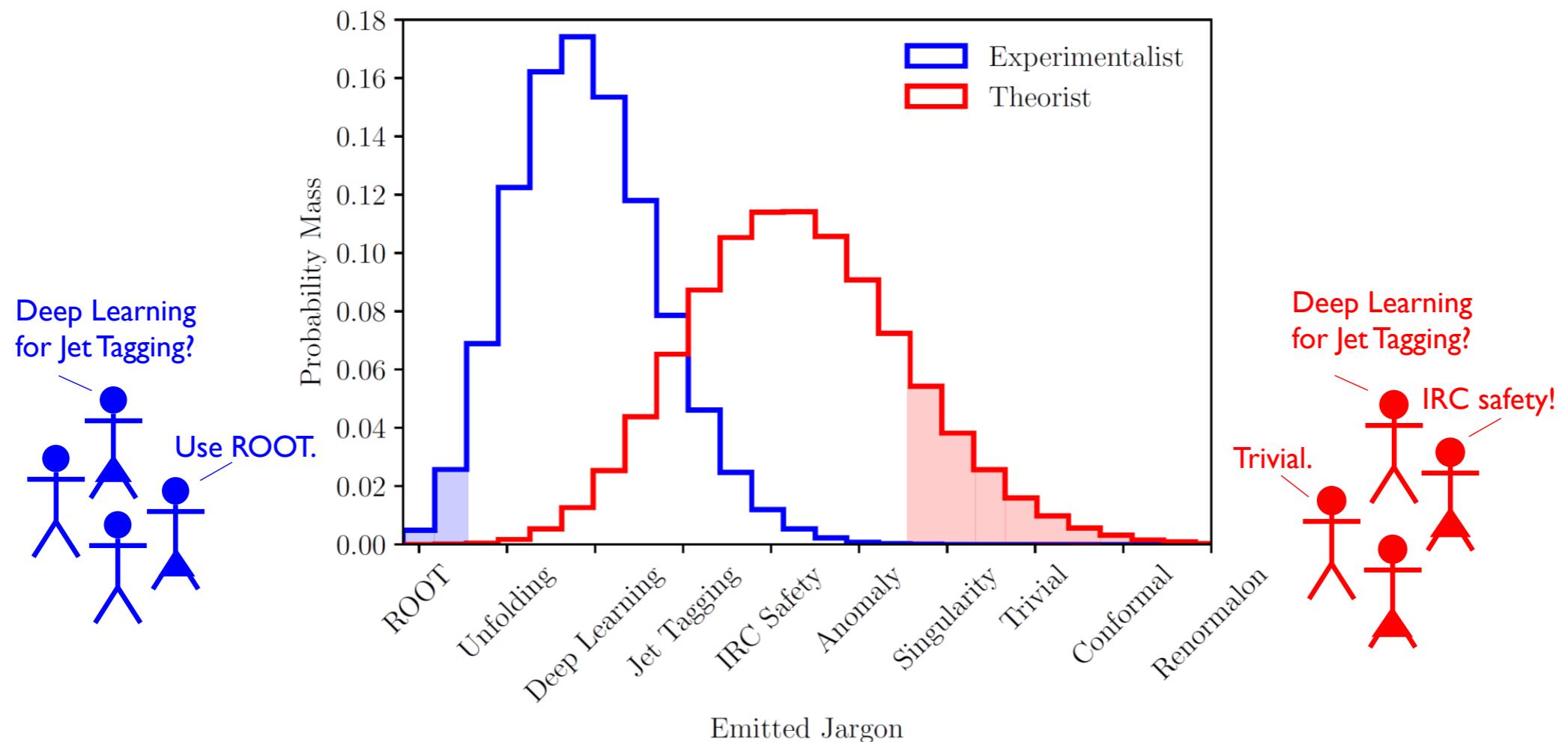
## (2) Generation



# Weak supervision turned on its head

## An Example

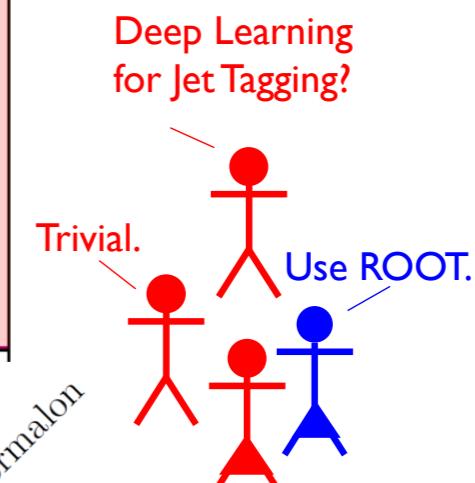
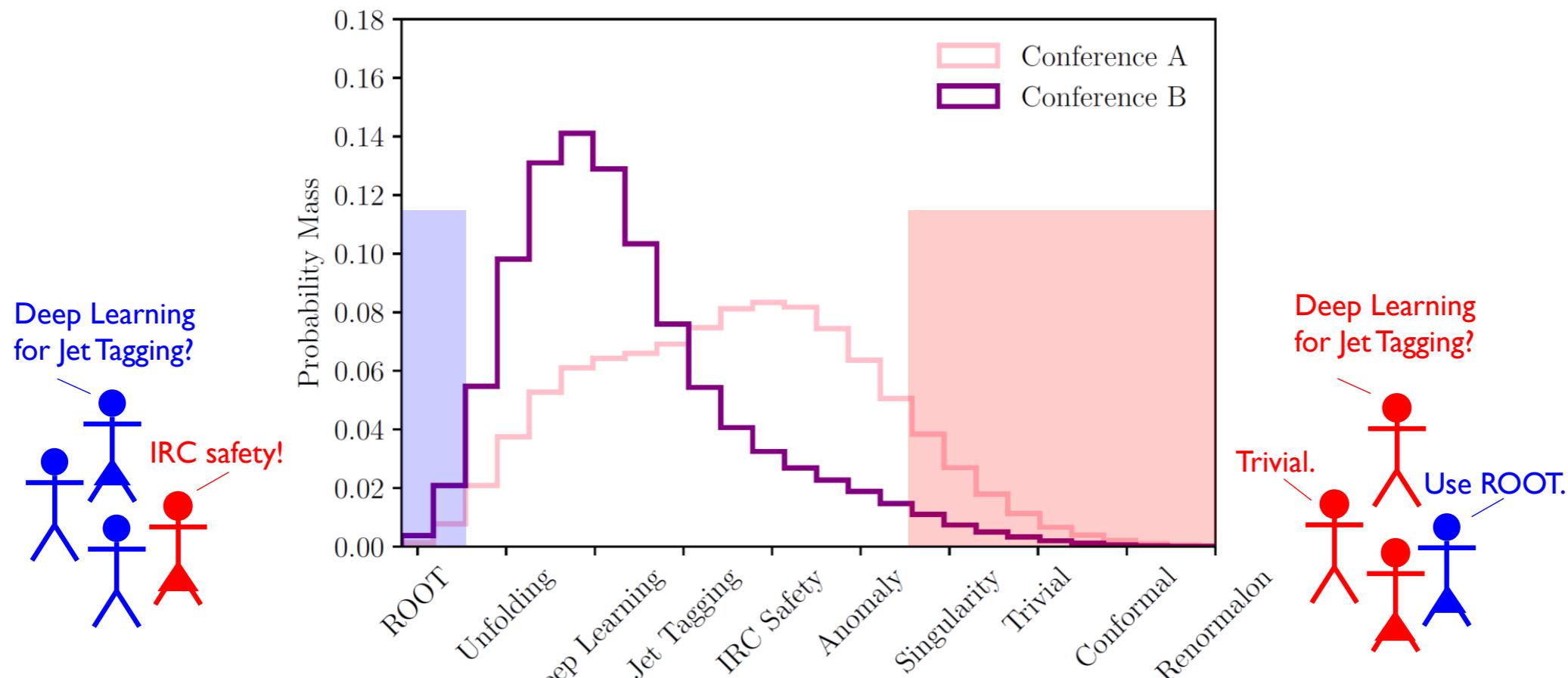
Let's model physicists as random jargon emitters.



# Weak supervision turned on its head

## An Example

Listen to the jargon emitted from two different conferences.



$$\frac{N^{\text{Conf. } A} \text{"ROOT"}}{N^{\text{Conf. } B} \text{"ROOT}} = \frac{f_{\text{Expt.}}^{\text{Conf. } A}}{f_{\text{Expt.}}^{\text{Conf. } B}}$$

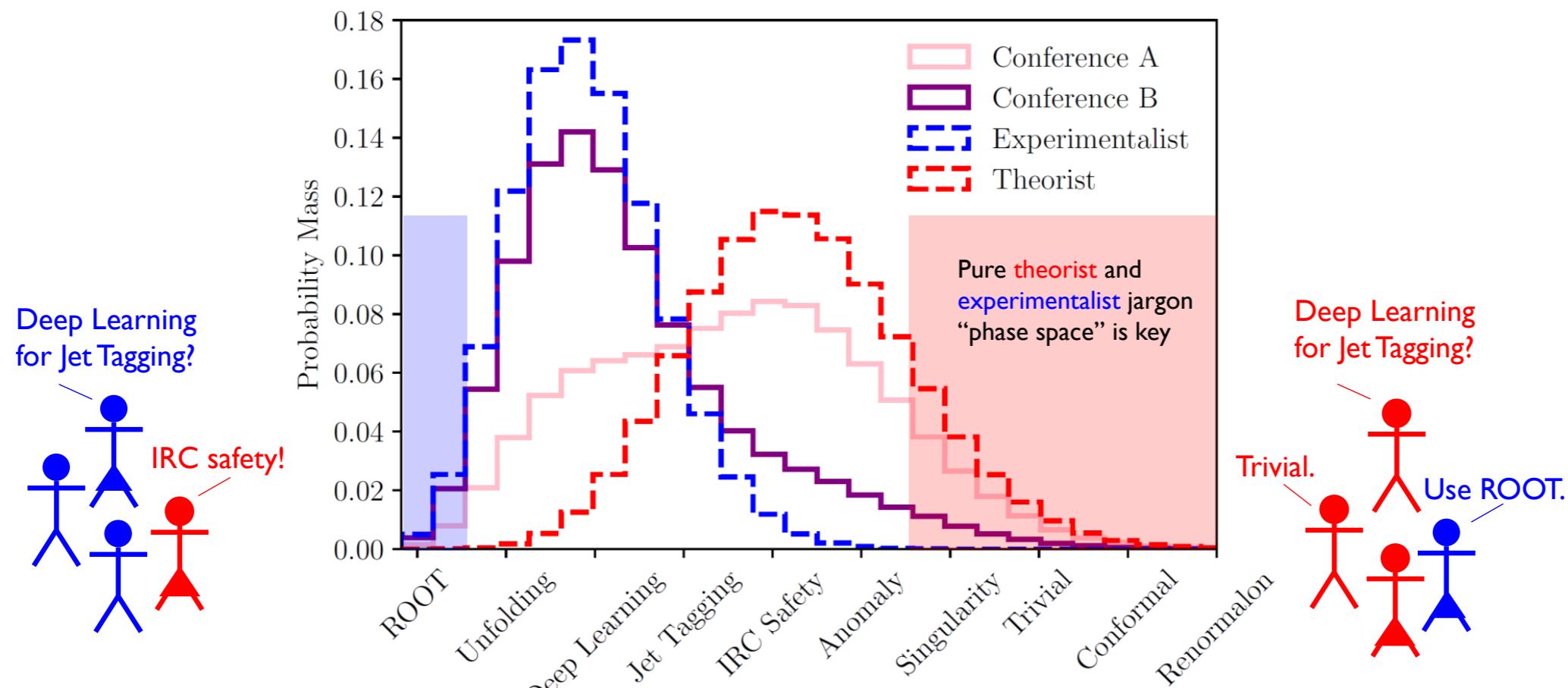
Emitted Jargon

$$\frac{N^{\text{Conf. } A} \text{"Trivial"}}{N^{\text{Conf. } B} \text{"Trivial}} = \frac{1 - f_{\text{Expt.}}^{\text{Conf. } A}}{1 - f_{\text{Expt.}}^{\text{Conf. } B}}$$

# Weak supervision turned on its head

## An Example

Disentangle **theorist** and **experimentalist** vocabularies from the jargon at conferences.



$$\frac{N^{\text{Conf. } A} \text{"ROOT"} }{N^{\text{Conf. } B} \text{"ROOT"} } = \frac{f_{\text{Expt.}}^{\text{Conf. } A}}{f_{\text{Expt.}}^{\text{Conf. } B}}$$

Emitted Jargon

$$\frac{N^{\text{Conf. } A} \text{"Trivial"} }{N^{\text{Conf. } B} \text{"Trivial"} } = \frac{1 - f_{\text{Expt.}}^{\text{Conf. } A}}{1 - f_{\text{Expt.}}^{\text{Conf. } B}}$$

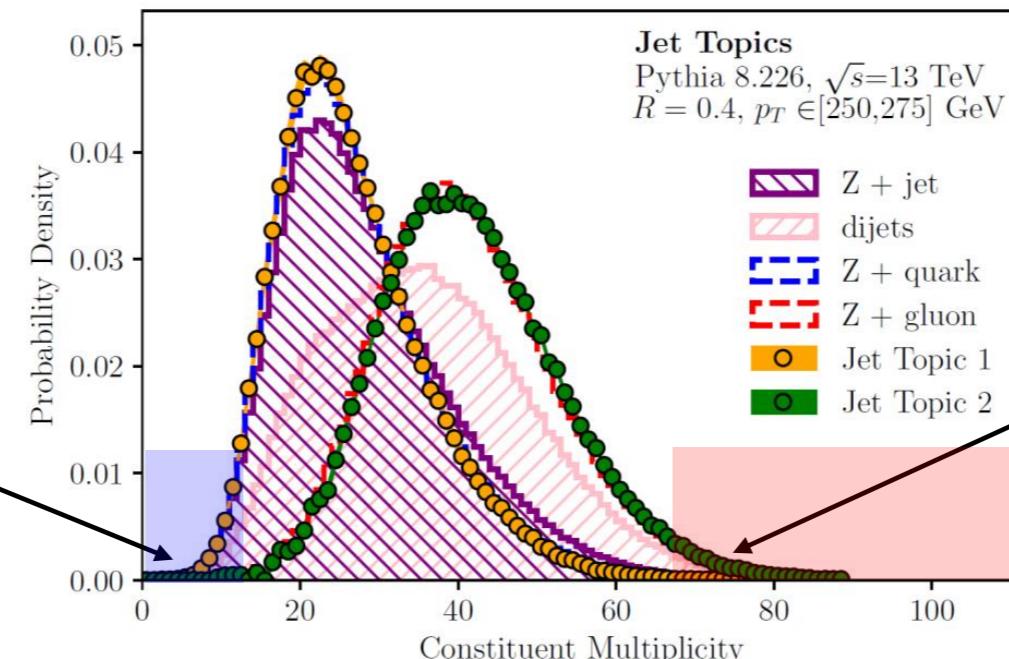
# Weak supervision turned on its head

## Demixing the mixtures

$$p_A(x) = f_A^q p_{\text{quark}}(x) + (1 - f_A^q) p_{\text{gluon}}(x)$$

$$p_B(x) = f_B^q p_{\text{quark}}(x) + (1 - f_B^q) p_{\text{gluon}}(x)$$

$$\begin{aligned} \kappa_{BA} &\equiv \min_x \frac{p_B(x)}{p_A(x)} \\ &= \frac{f_B^q}{f_A^q} \end{aligned}$$



$$\begin{aligned} \kappa_{AB} &\equiv \min_x \frac{p_A(x)}{p_B(x)} \\ &= \frac{1 - f_A^q}{1 - f_B^q} \end{aligned}$$

With **reducibility factors**  $\kappa_{AB}$  and  $\kappa_{BA}$ , solve for the quark and gluon fractions and distributions:

$$f_A^q = \frac{1 - \kappa_{AB}}{1 - \kappa_{AB} \kappa_{BA}}$$

$$f_B^q = \frac{\kappa_{BA}(1 - \kappa_{AB})}{1 - \kappa_{AB} \kappa_{BA}}$$

$$p_{\text{quark}}(x) = \frac{p_A(x) - \kappa_{AB} p_B(x)}{1 - \kappa_{AB}}$$

$$p_{\text{gluon}}(x) = \frac{p_B(x) - \kappa_{BA} p_A(x)}{1 - \kappa_{BA}}$$

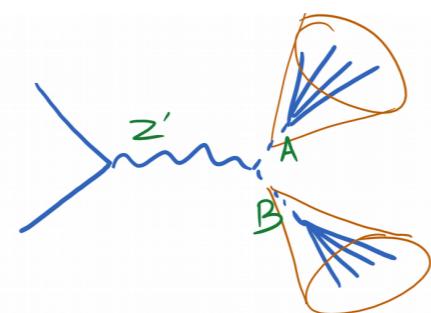
# Weak supervision for new physics?

Jack Collins

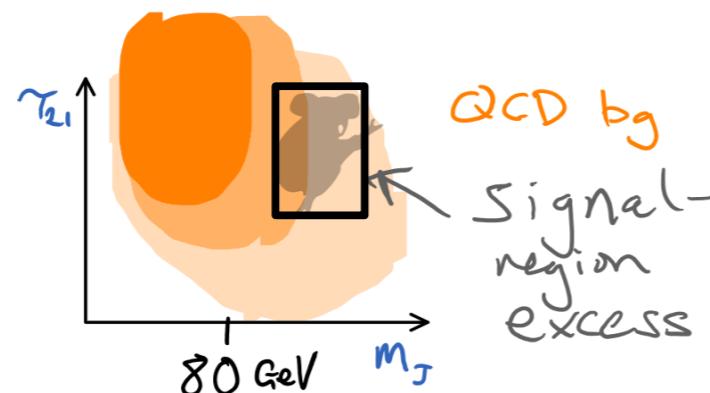
## Case Study: Dijet Resonances

### CWoLa Hunting Proposal

- 1) Propose generic signal resonance hypothesis



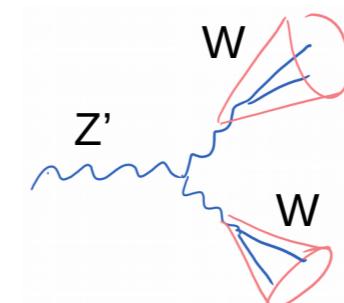
- 2) Let Machine look for signal region overdensity in some substructure phase space region



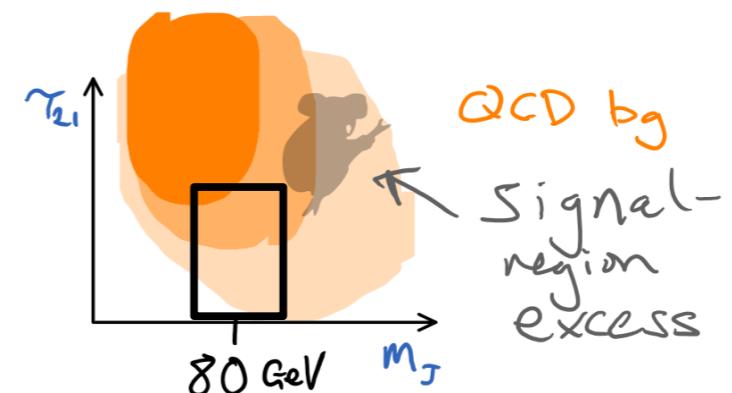
- 3) Perform bump-hunt in distribution of selected events

### How do existing di-SM searches use substructure information?

- 1) Propose specific signal resonance hypothesis



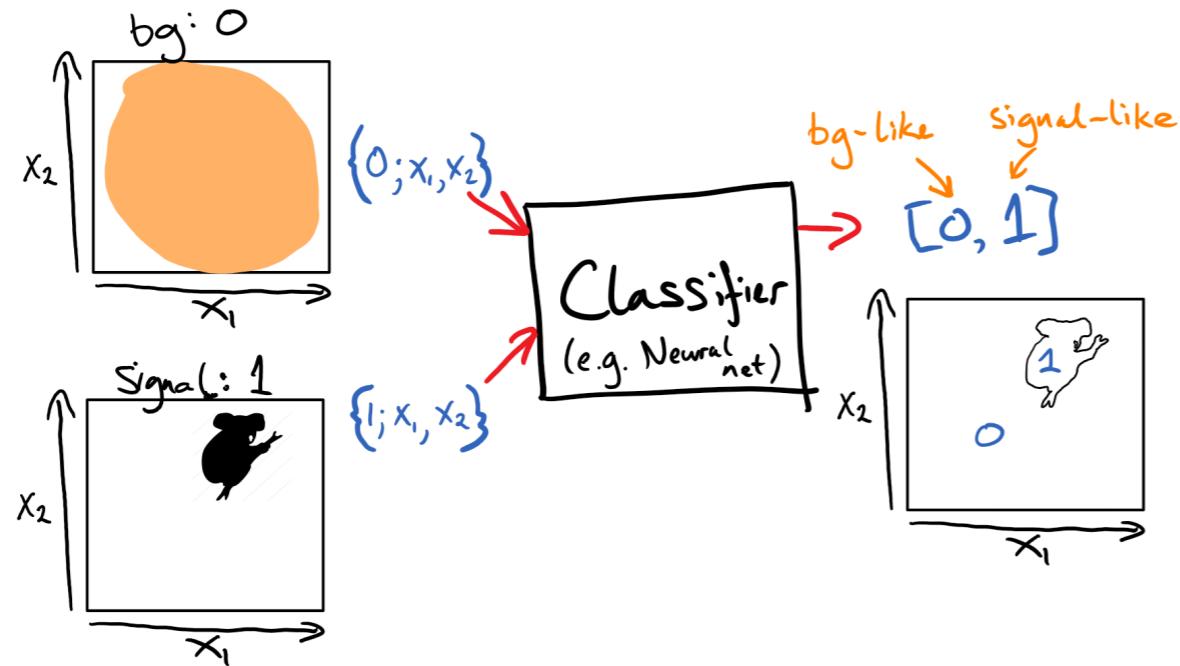
- 2) Apply pre-defined model-specific substructure cut using a couple of expert variables or supervised NN



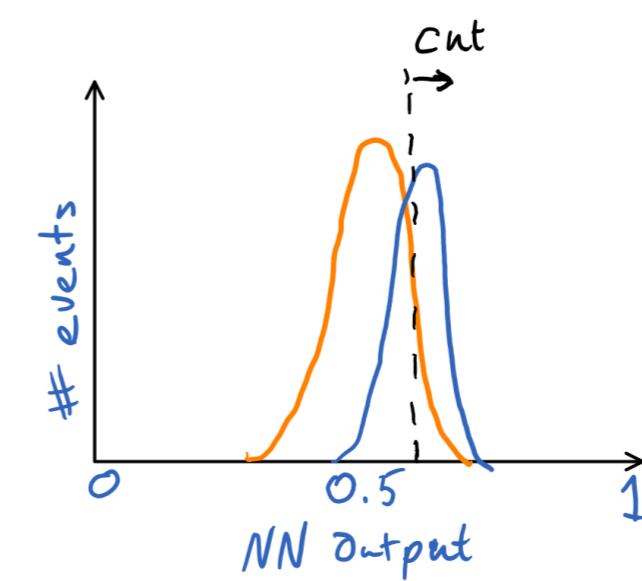
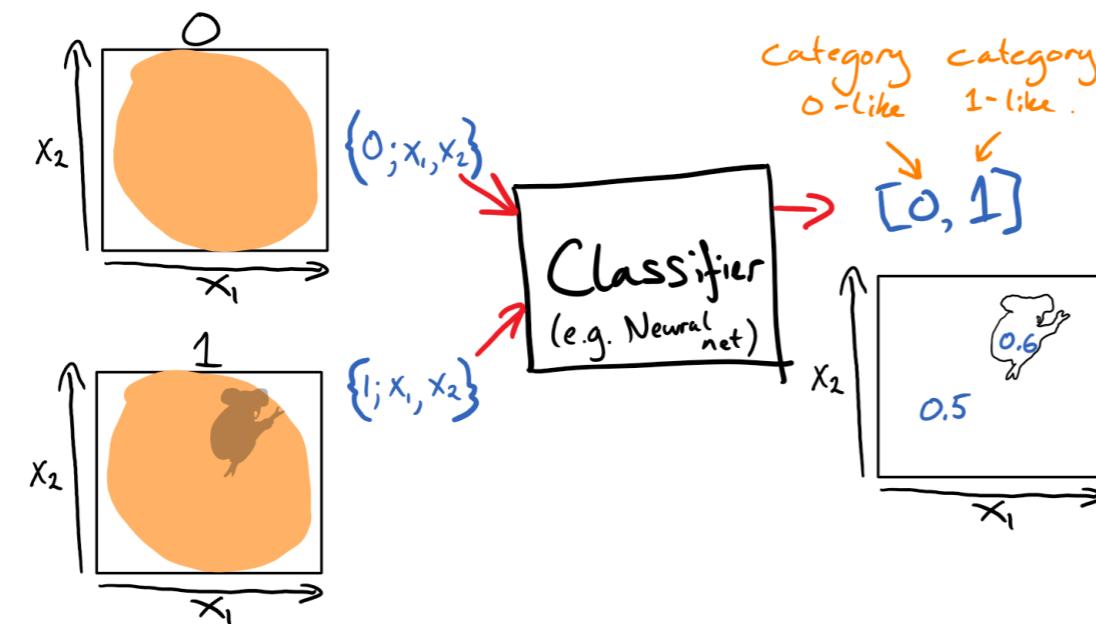
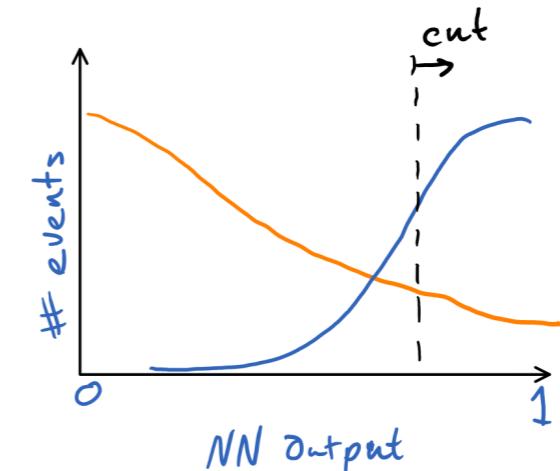
- 3) Perform bump-hunt in distribution of selected events

# Weak supervision for new physics?

## Fully Supervised learning vs Weak Supervision

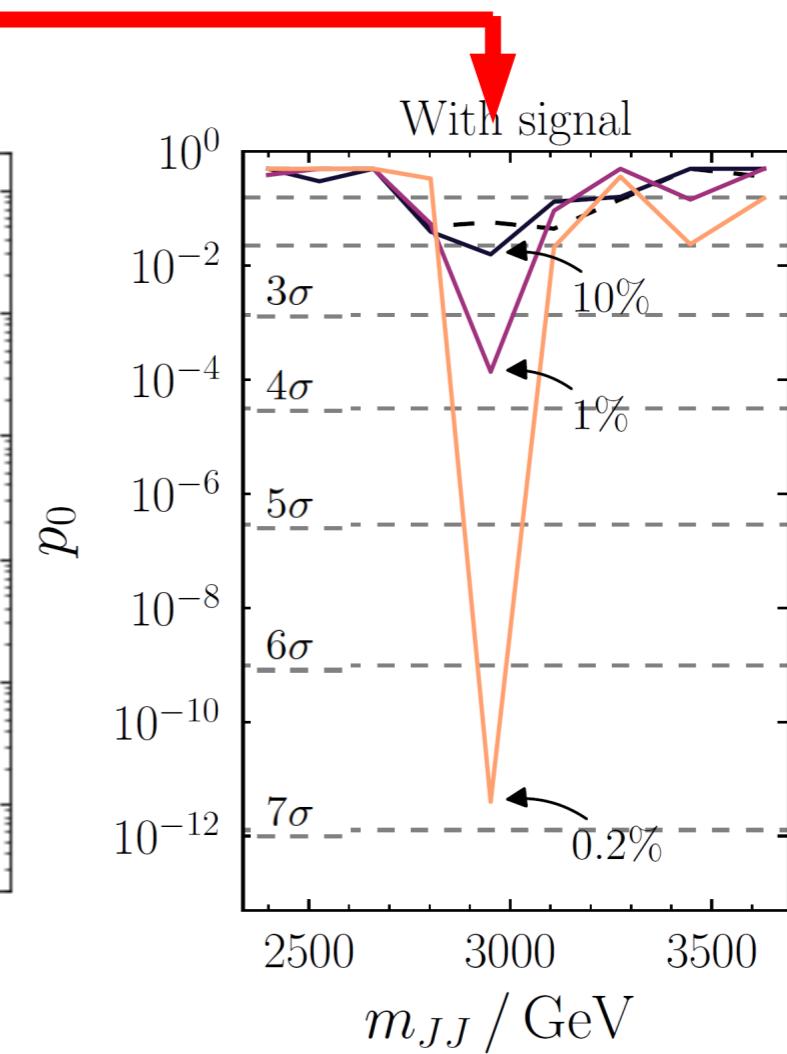
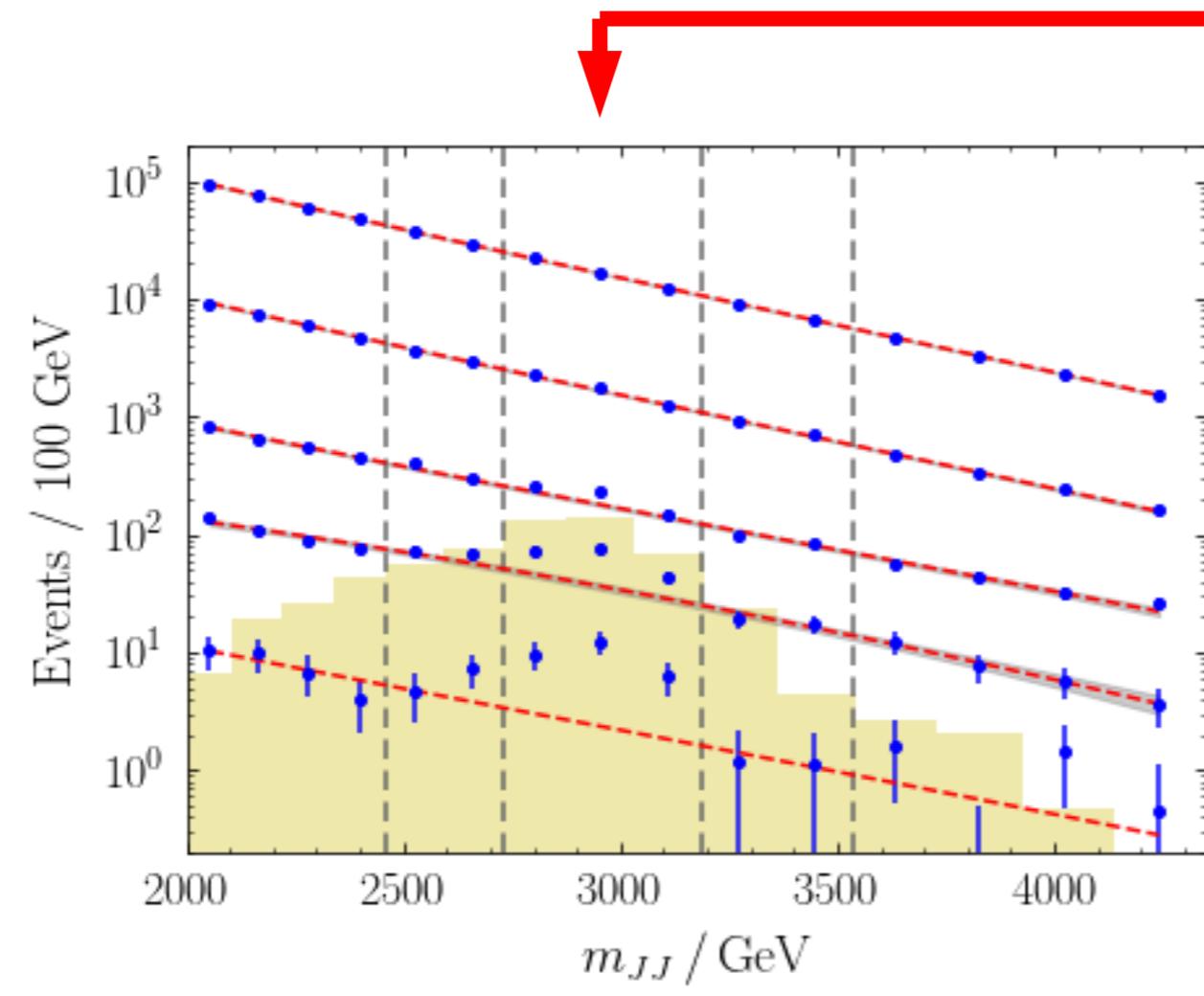


- True signal  
- True BG



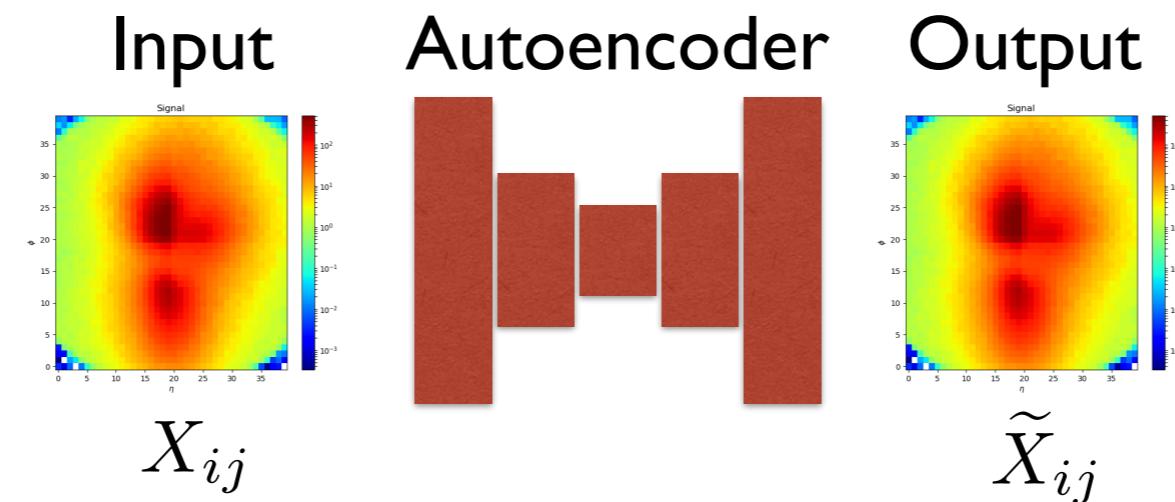
# Weak supervision for new physics?

## Mass Scan



# Autoencoders!

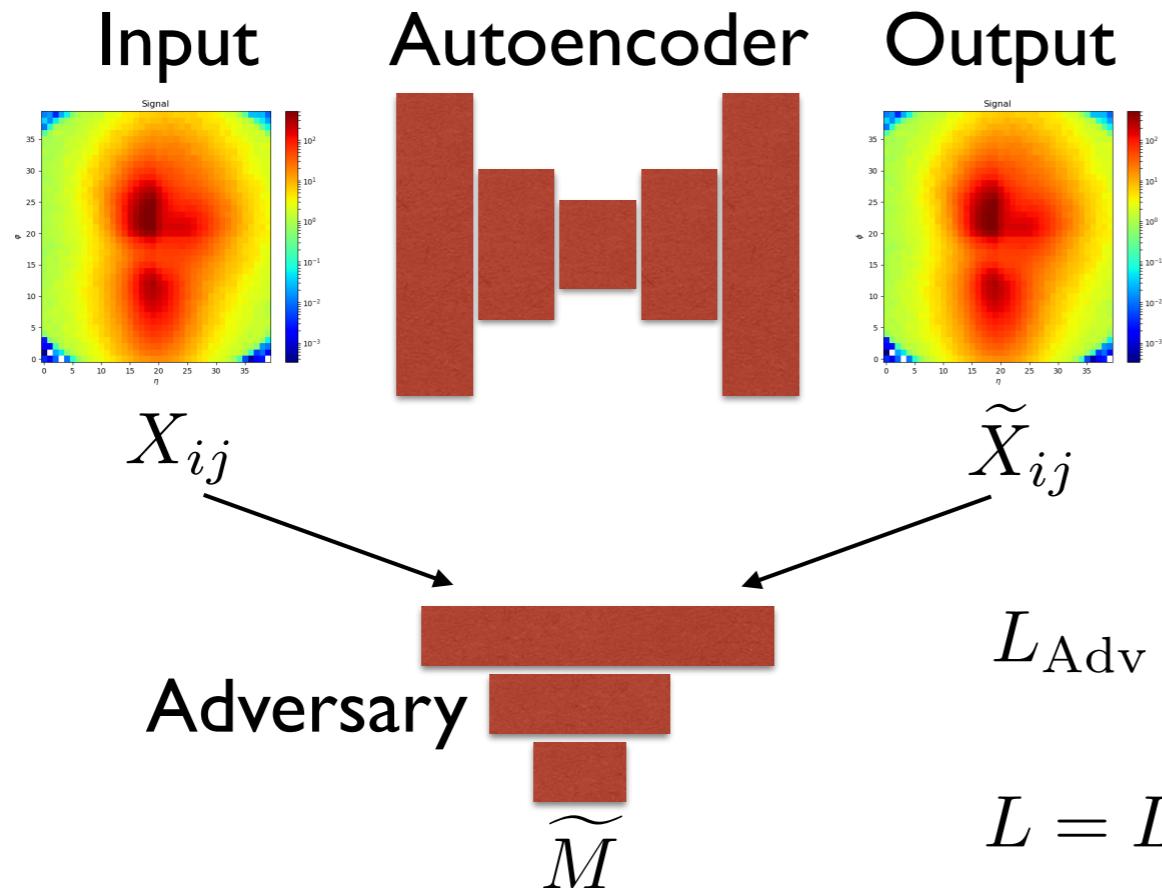
## Autoencoder for Physics



- Train on pure QCD light quark/gluon jets
  - Use top tagging reference sample  
<https://goo.gl/XGYju3>
  - Train **only on QCD** events
- New physics identified as anomaly
  - Tail of the loss function

# Autoencoders!

## Combined Setup



$$L_{\text{Auto}} = \sum_{\text{Pixels } ij} (X_{ij} - \tilde{X}_{ij})$$

$$L_{\text{Adv}} = \text{CCE} \left( M, \tilde{M}(X_{ij} - \tilde{X}_{ij}) \right)$$

$$L = L_{\text{Auto}} - \lambda L_{\text{Adv}}$$

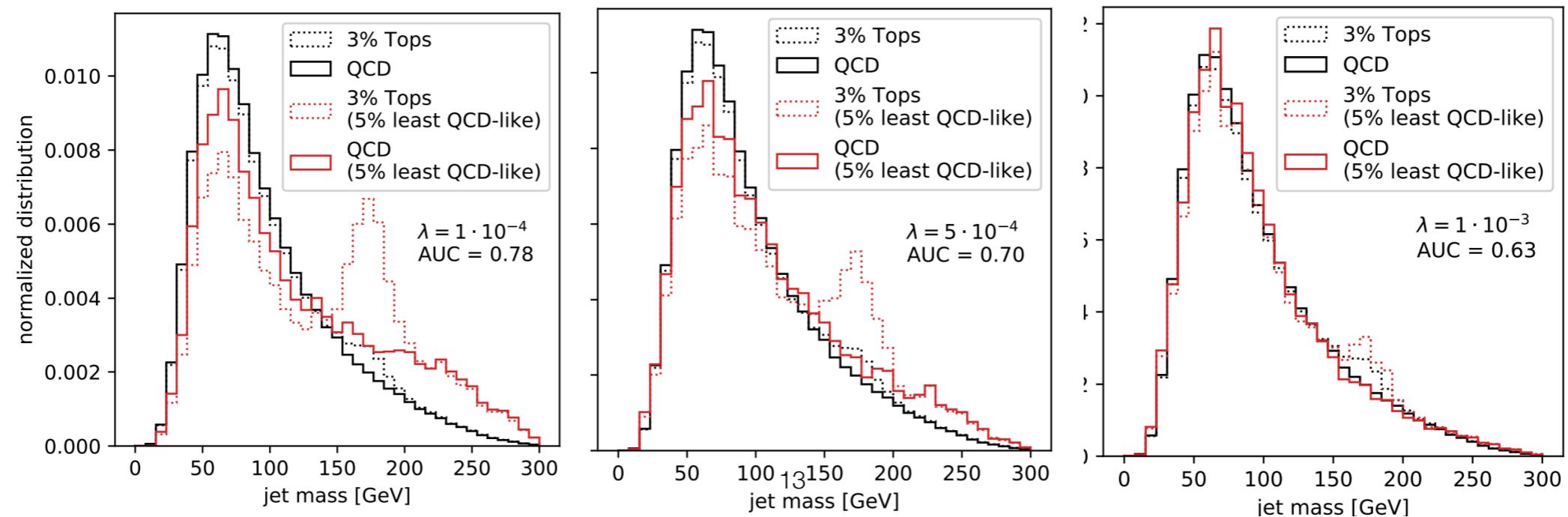
# Autoencoders!

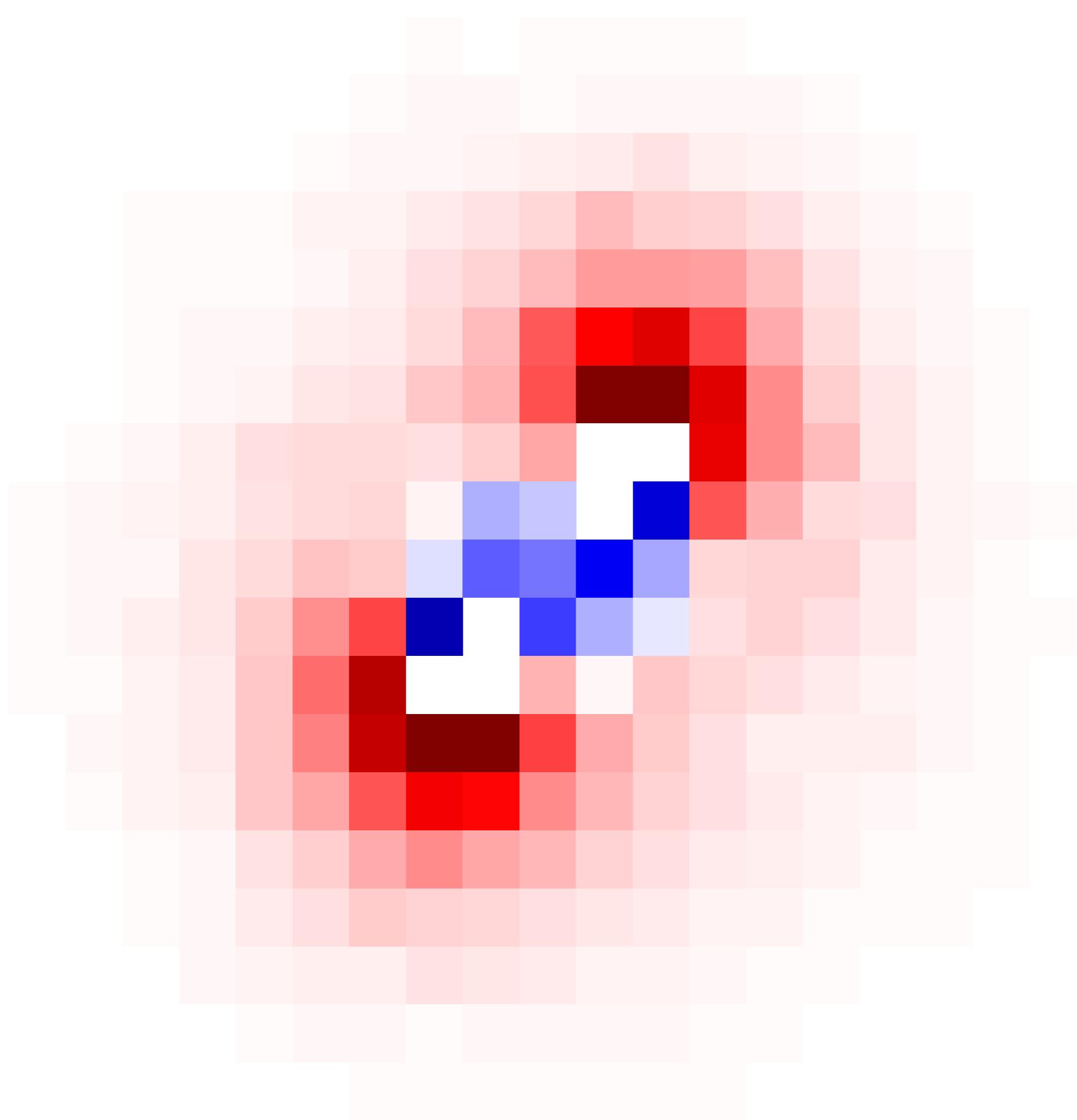
## Mass Sculpting

- Counteract with adversary:

$$L = L_{\text{Auto}} - \lambda L_{\text{Adv}}$$

- Tune mass dependency with Lagrange multiplier
  - Defines control regions in data





Fin.