# EXPLAINABLE AI FOR ML JET TAGGERS

Garvita Agarwal, Lauren Hay, Ia Iashvili, Benjamin Mannix, *Christine McLean*, Margaret Morris, Salvatore Rappoccio, Ulrich Schubert
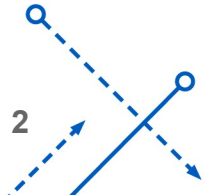
ML4Jets2021

July 8, 2021

**UB** University at Buffalo The State University of New York

https://doi.org/10.1007/978-3-030-28954-6_10
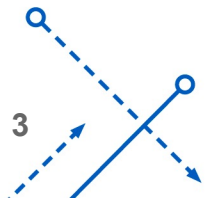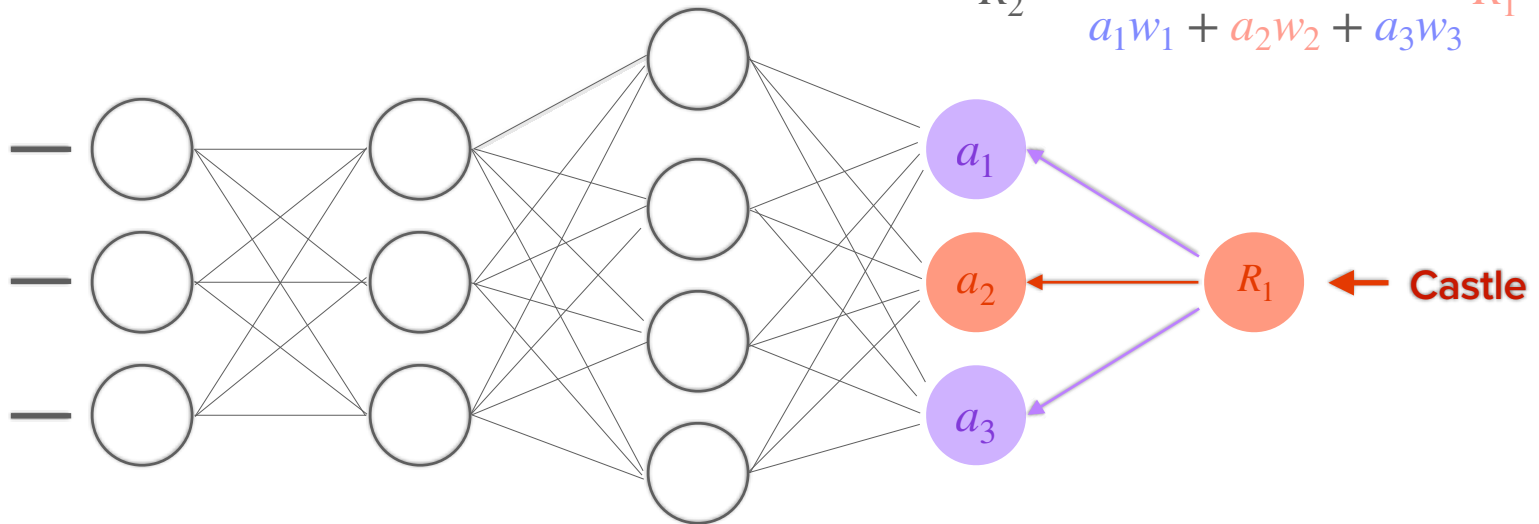
ML Network Black Box

Classification: Castle

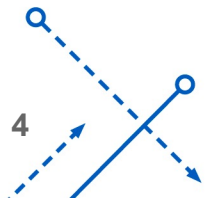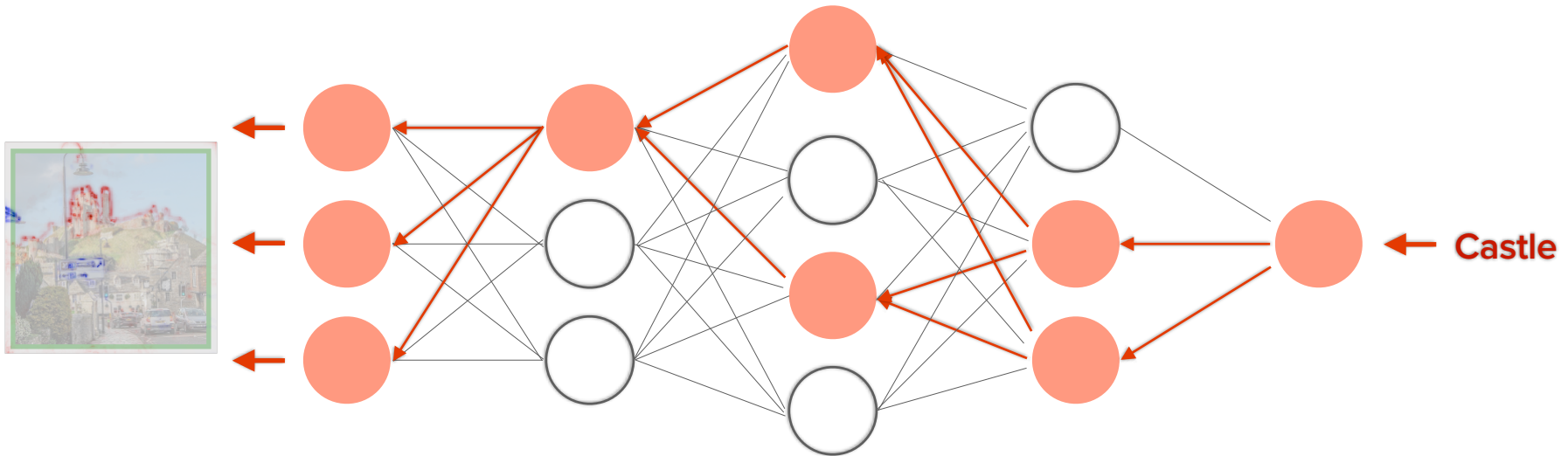*How do we understand the network's decision-making process?*

**LRP (layer-wise relevance propagation)** propagates a prediction backwards through the network, assigning a relevance to each input

$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

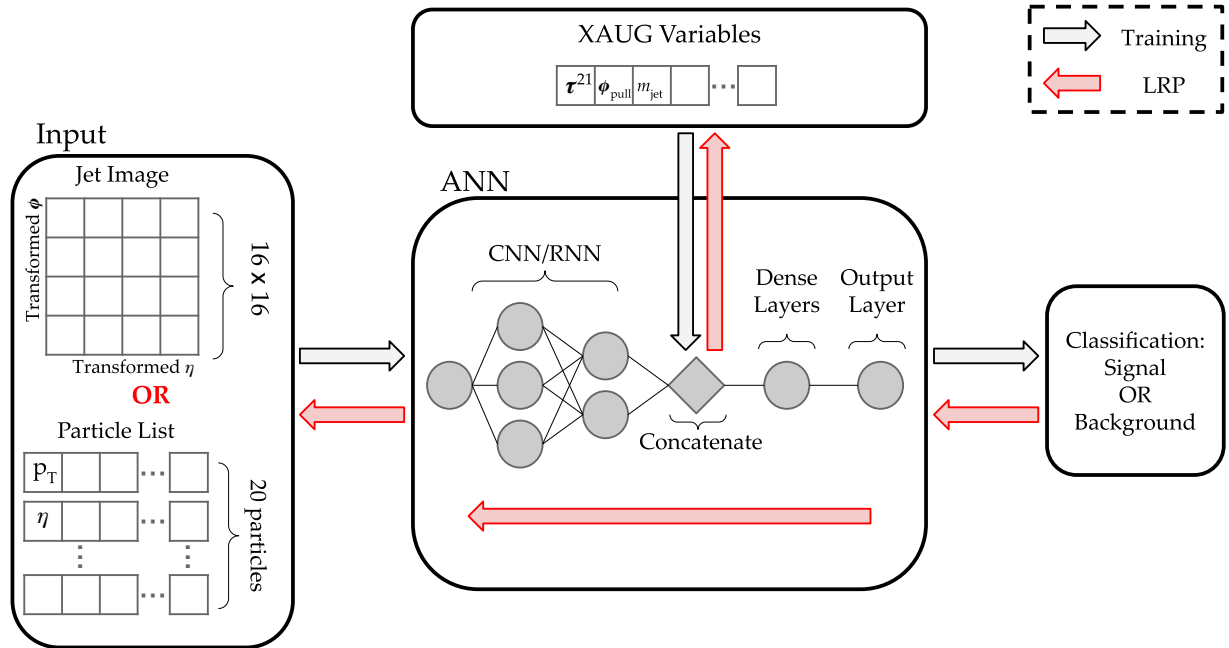$$R_2 = \frac{a_2 w_2}{a_1 w_1 + a_2 w_2 + a_3 w_3} R_1$$



Castle

# ML explainability with LRP

- **Relevance is conserved -** the prediction is not changed

- **LRP attributes the entirety of the network's decision to the inputs**
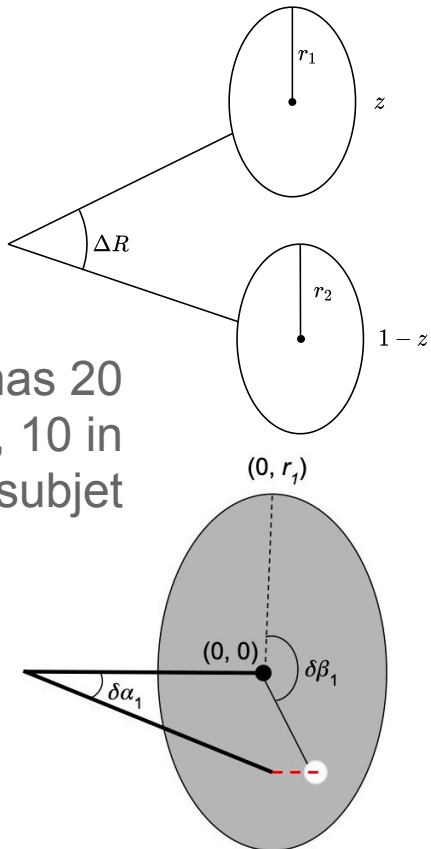  - Visualized as a heat map, in the case of images

# ML explainability with XAUG Variables

- **Goal:** explain decisions of ML jet classifiers using expert augmented (XAUG) variables

- **Method:** Input XAUGs into jet tagger, analyze network decision with LRP, and compare to network without XAUGs
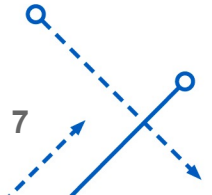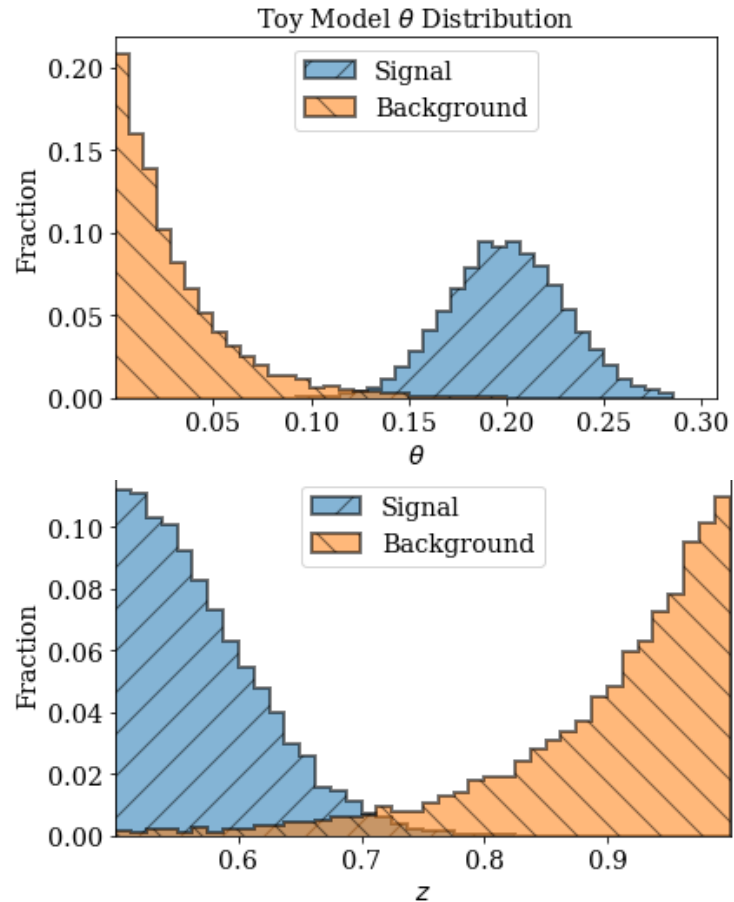
# TOY MODEL

- **Toy events simulated to mimic particle-level events**
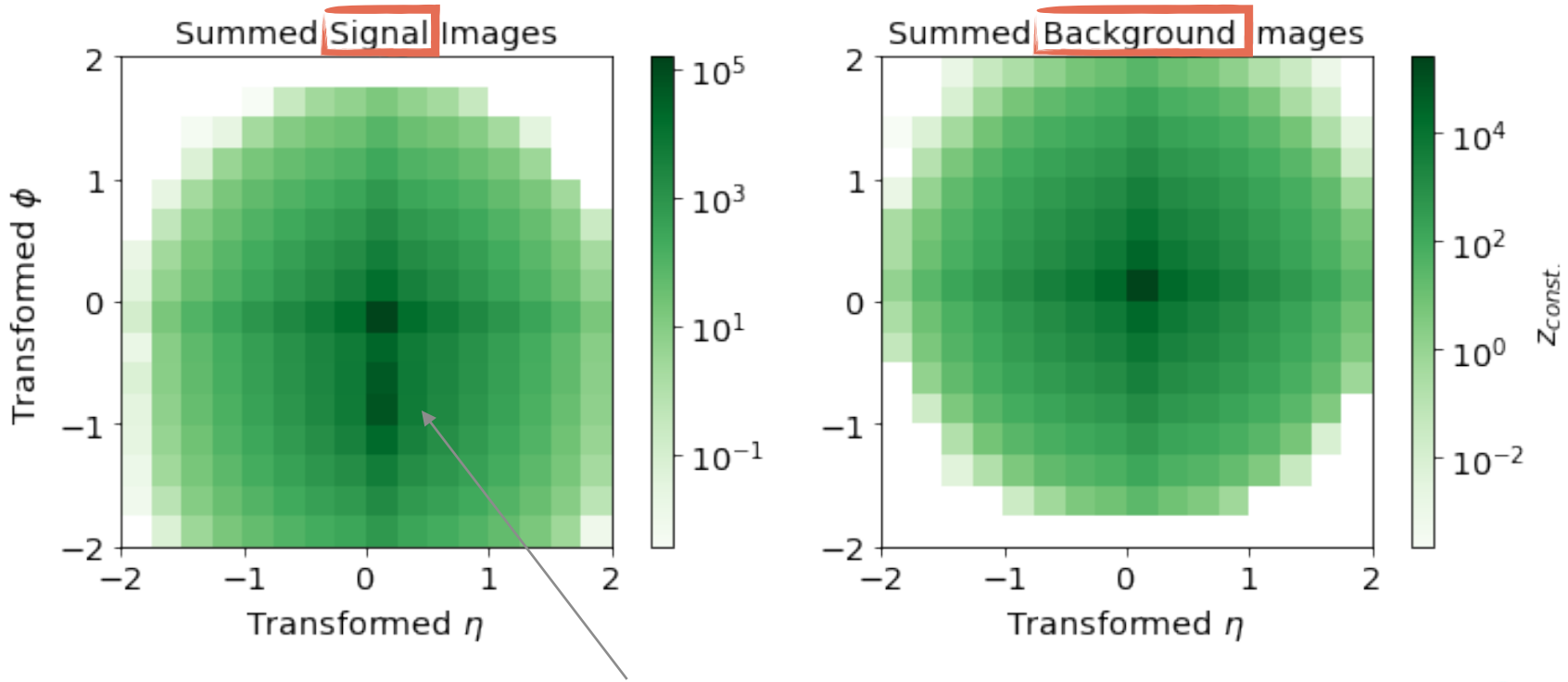- **Goal: capture all event information with a few variables**
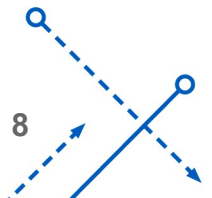
Each "jet" has 20 particles, 10 in each subjet

Toy Model $\theta$ Distribution

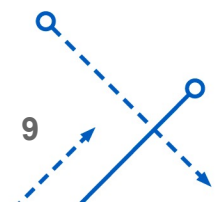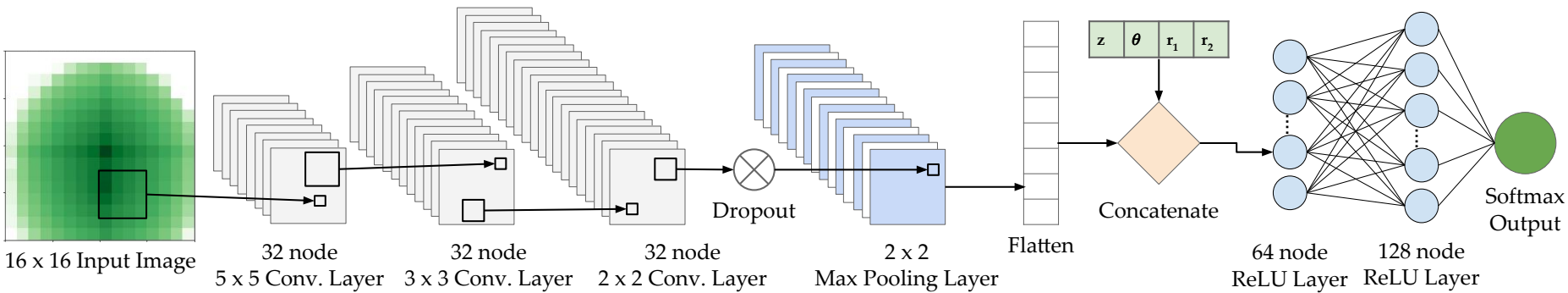- ## **Image pre-processing**
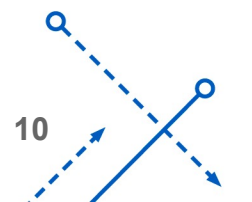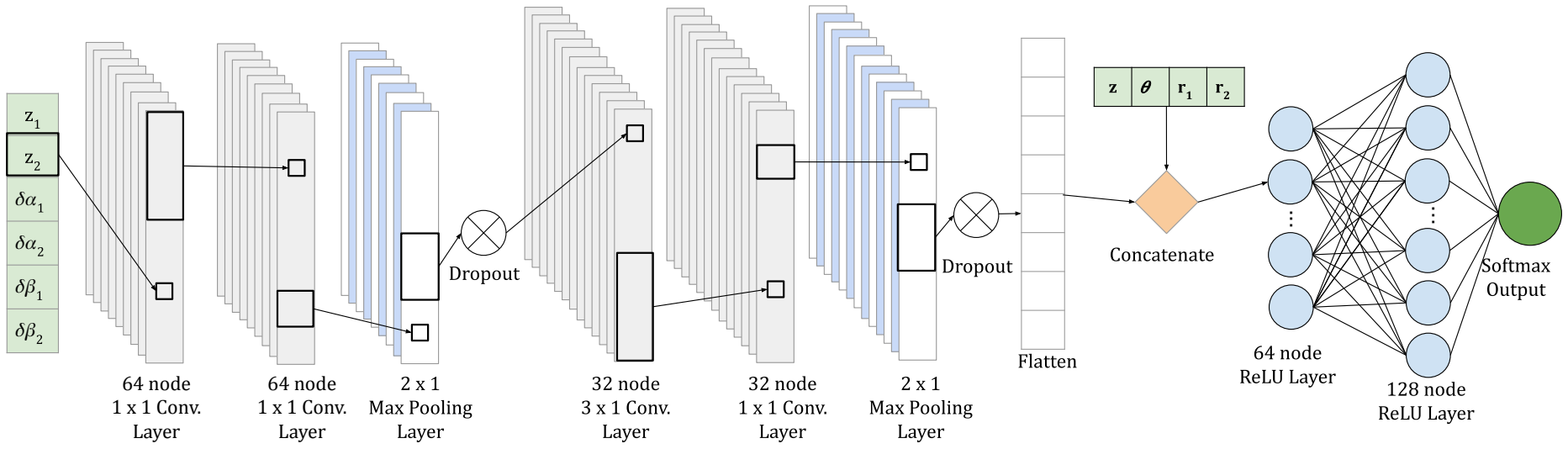  - Leading-$p_T$ subjet at (0,0), sub-leading at (0,-1)
  - Parity flip



Much more pronounced second subjet

# Architecture based on ImageTop network



16 x 16 Input Image    32 node 5 x 5 Conv. Layer    32 node 3 x 3 Conv. Layer    32 node 2 x 2 Conv. Layer    Dropout    2 x 2 Max Pooling Layer    Flatten    Concatenate    64 node ReLU Layer    128 node ReLU Layer    Softmax Output
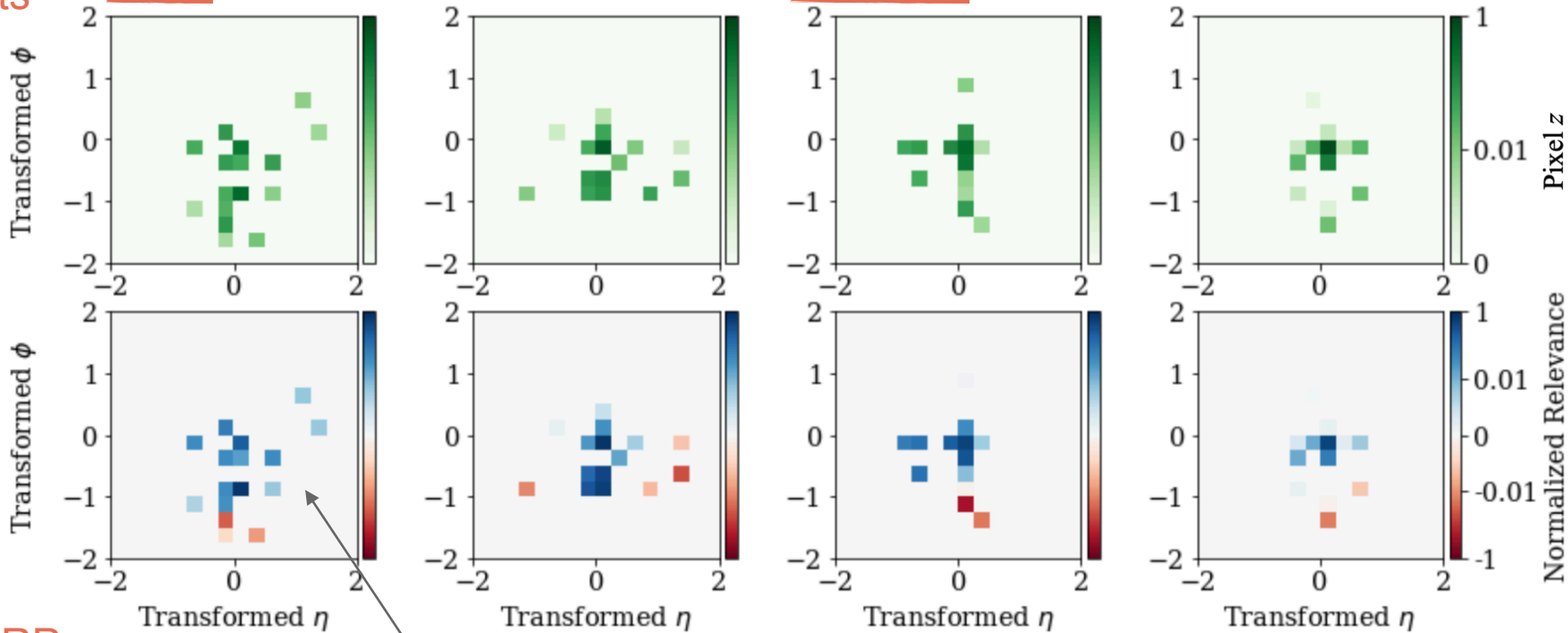
$z$   $\theta$   $r_1$   $r_2$
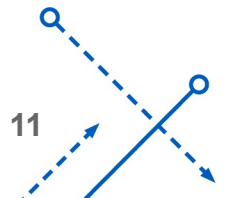
## Architecture based on DeepAK8 jet classifier

Inputs

Signal Images and their relevance heatmaps
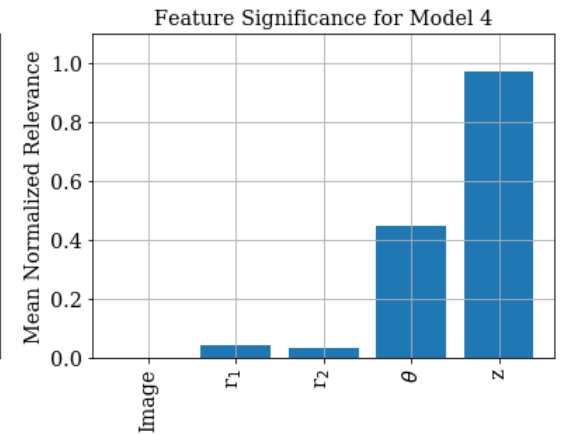
Background Images and their relevance heatmaps

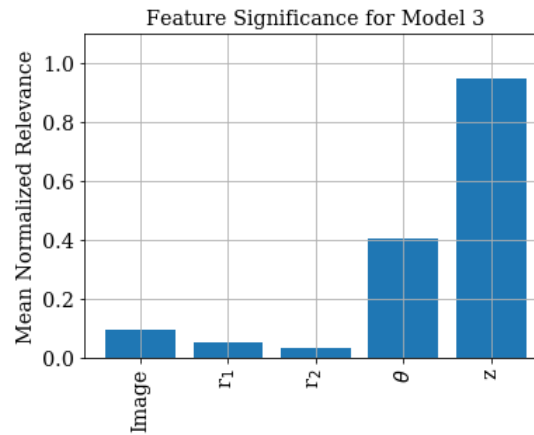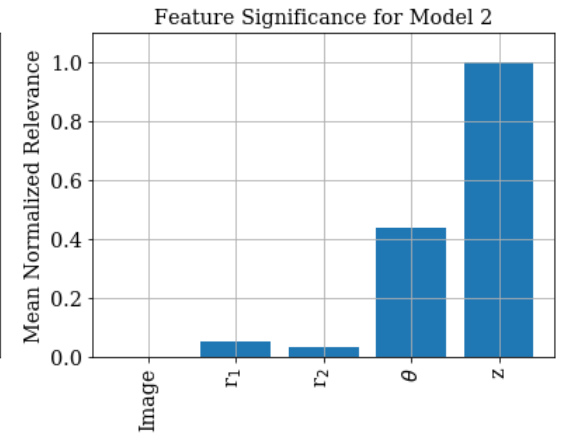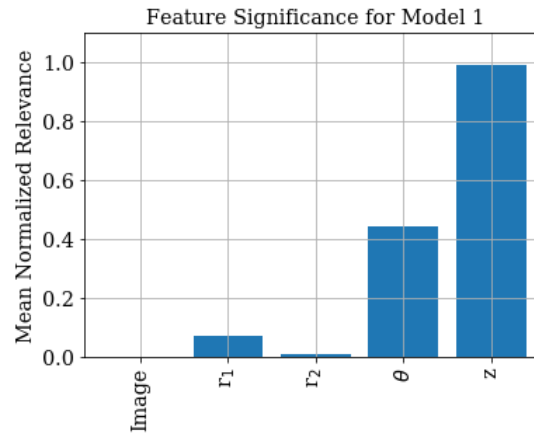LRP Results

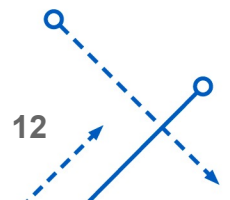Signal: more relevance along $\phi$ axis
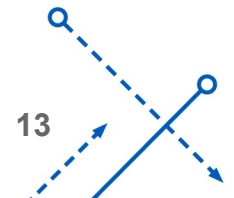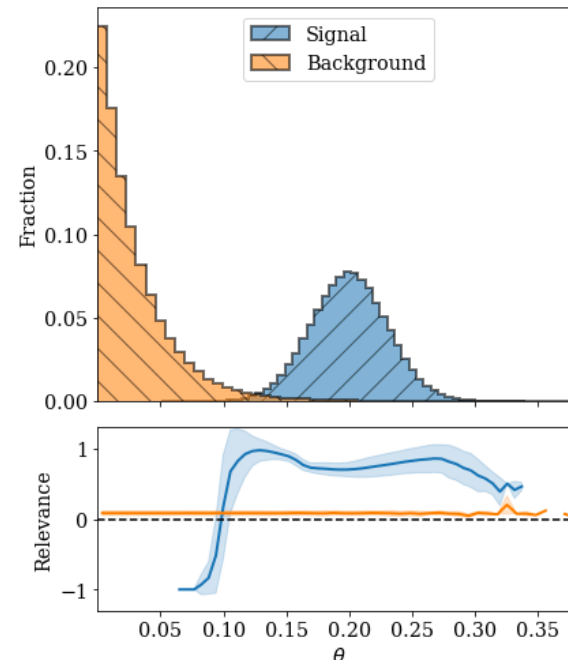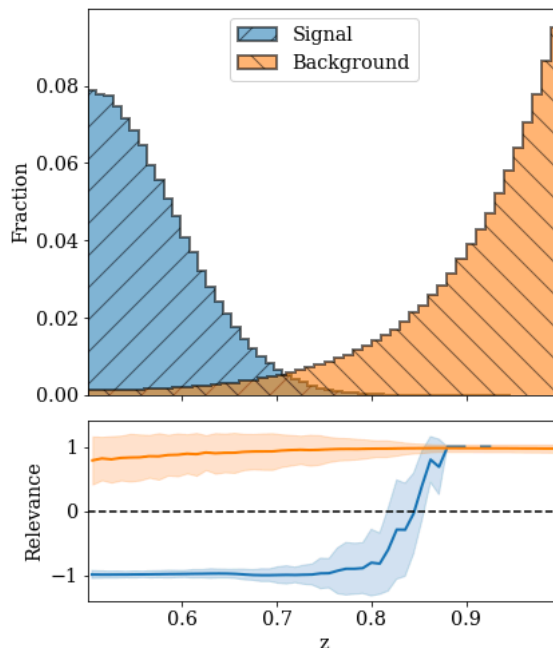
- ## **Mean normalized relevance**

  - **For each event:** find feature with max absolute LRP score, divide all scores by this max value
    - **For each image:** sum absolute value of normalized pixels to get a single image LRP score

  - **For each feature:** average normalized relevance scores across all events



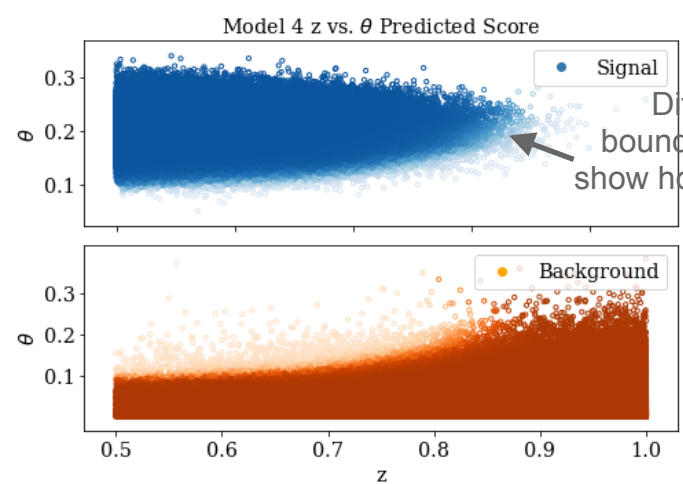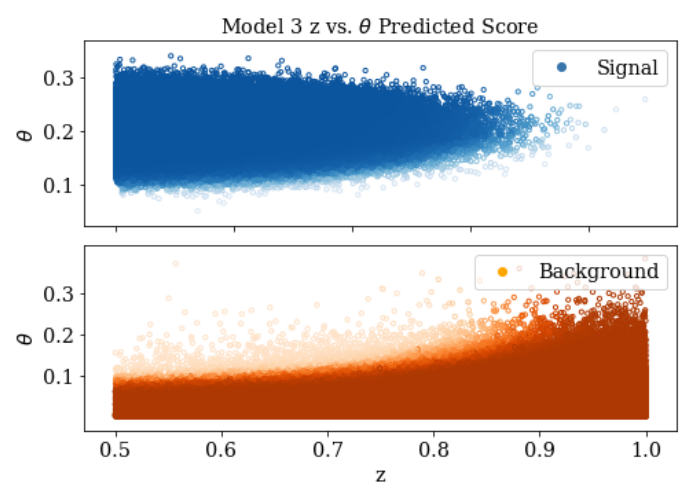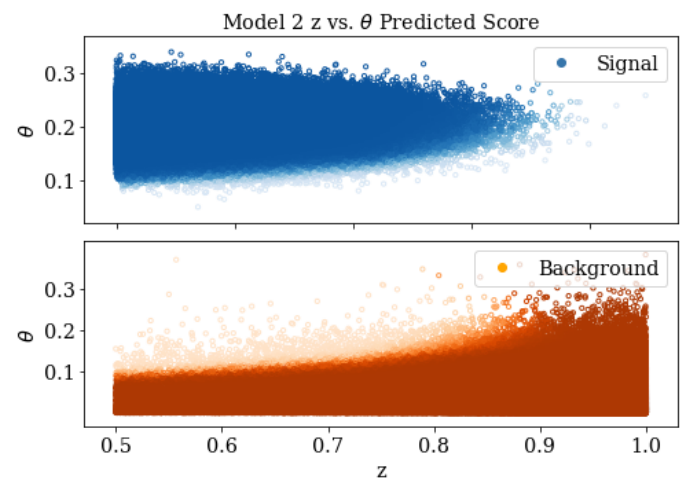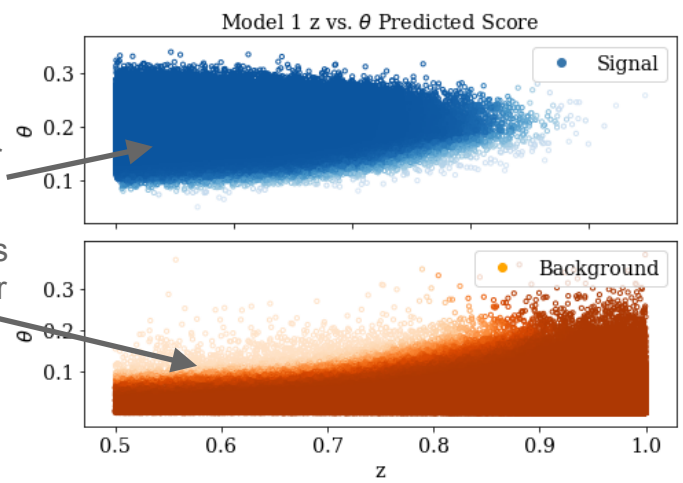Some variation between trainings

- **Profile plots:** relevance vs corresponding input variable
- For some profiles relevance appears to reflect input distribution, but other don't - **networks' decision boundaries live in a higher dimensional space**

# Toy 2DCNN Results



Darker markers corresponds to higher relevance scores.

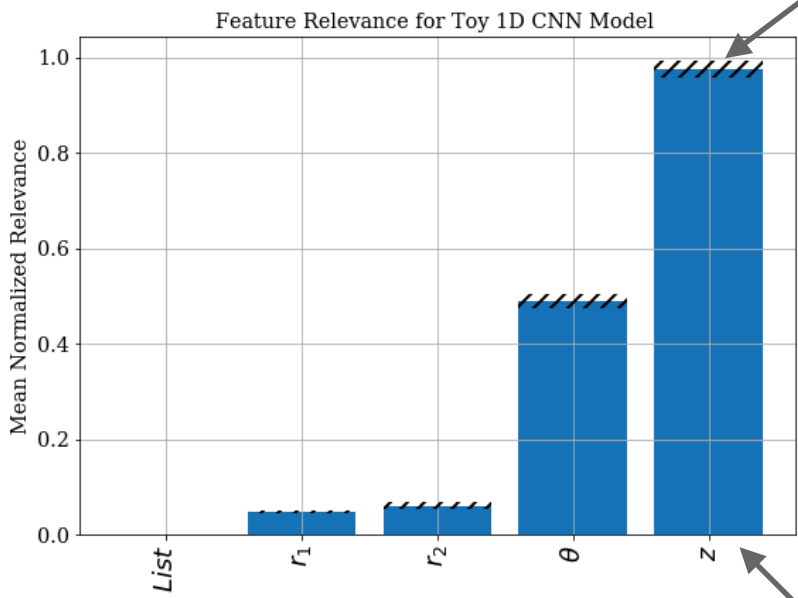Sharp gradient shows decision boundary for these variables

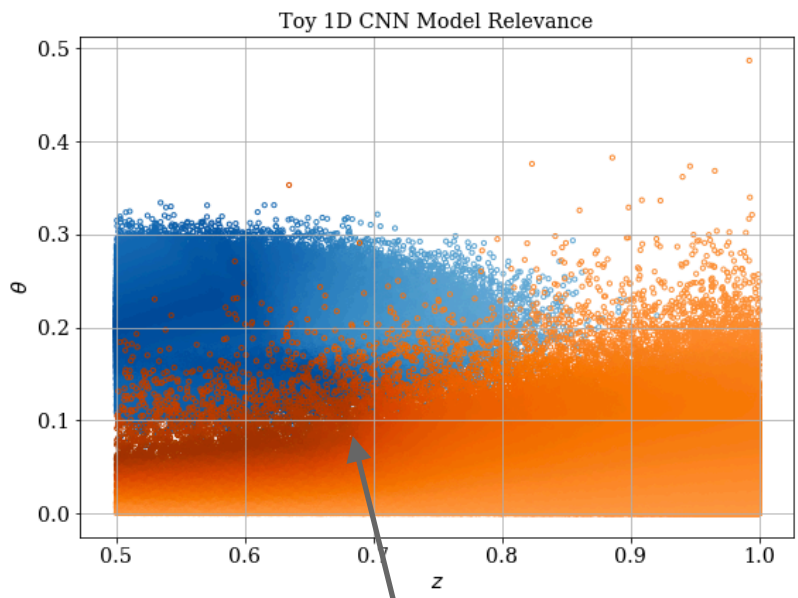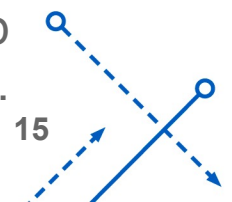Differences in boundary shapes show how trainings vary

Error bars show standard deviation of relevance after multiple trainings.

Feature Relevance for Toy 1D CNN Model

Mean Normalized Relevance

List $r_1$ $r_2$ $\theta$ $z$

Toy 1D CNN Model Relevance

$\theta$

$z$

Most relevant features are same as 2DCNN.

More robust "substructure" within relevance of the top two variables.
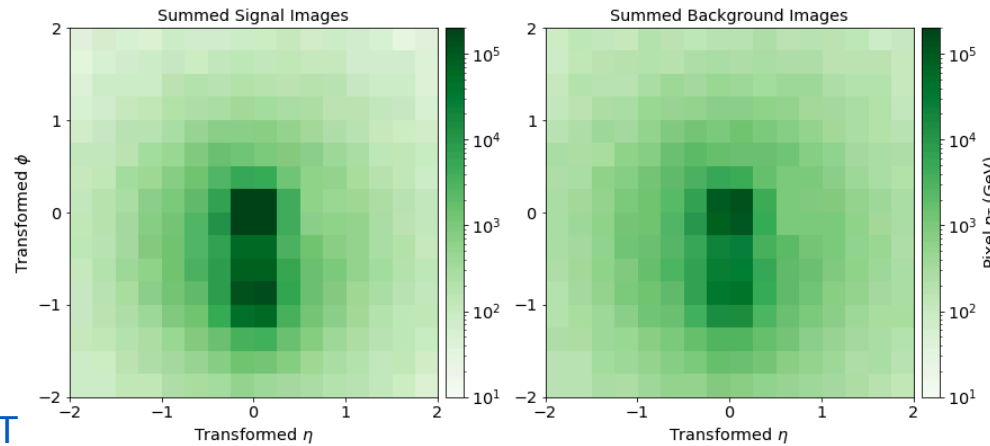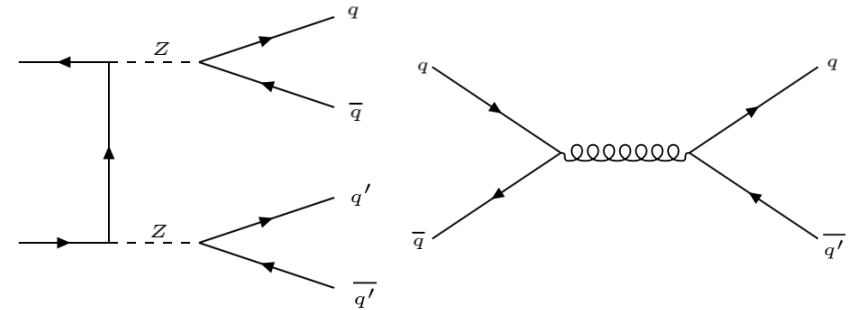
# PYTHIA MODEL

# Pythia Model

- ## Simulated with Pythia8
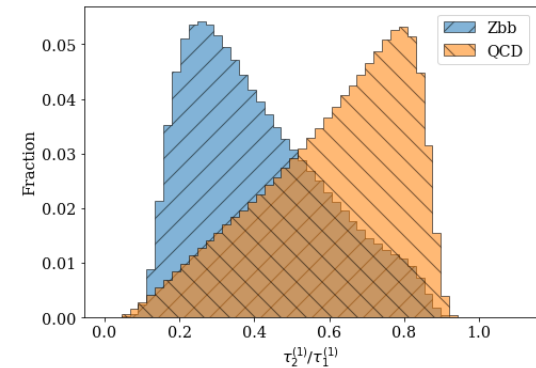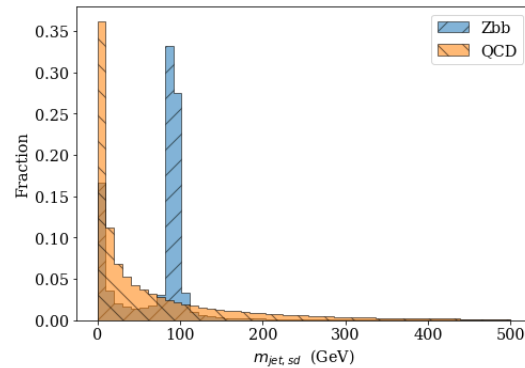  - Signal: SM $ZZ, Z \rightarrow b\bar{b}$
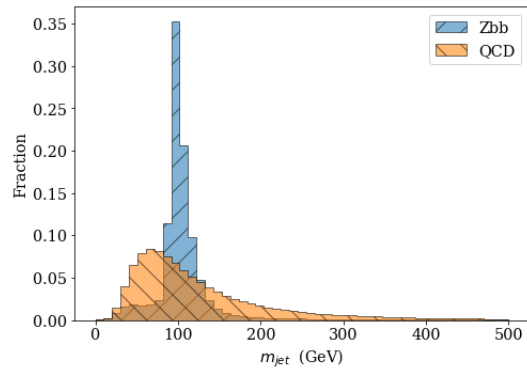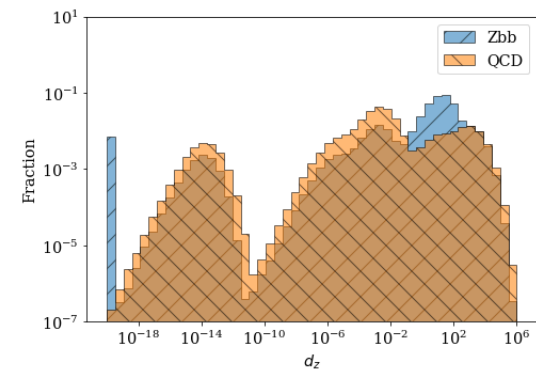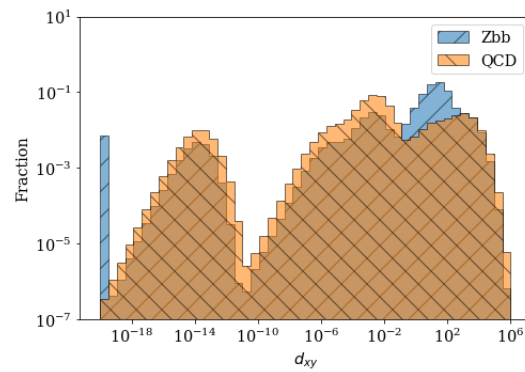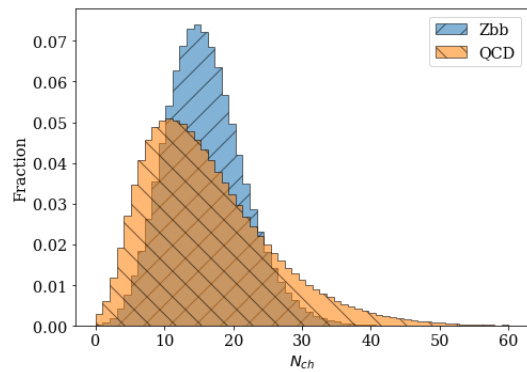  - QCD

- ## Jets
  - Consider leading AK8 jet
  - $p_T$ > 200 GeV
  - mMDT: $z_{cut} = 0.1, \beta = 0$

- ## Preprocessing
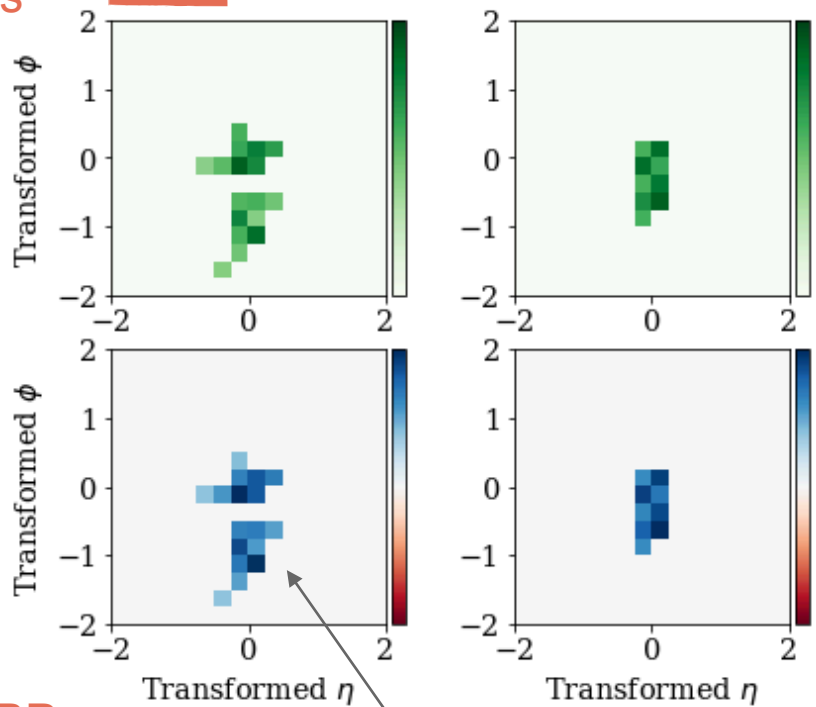  - Normalize inputs wrt to jet $p_T$
  - Same as toy model



Summed Signal Images

Summed Background Images

## Use same network structure as toy model; replace particle-level inputs

# Pythia LRP Heatmaps
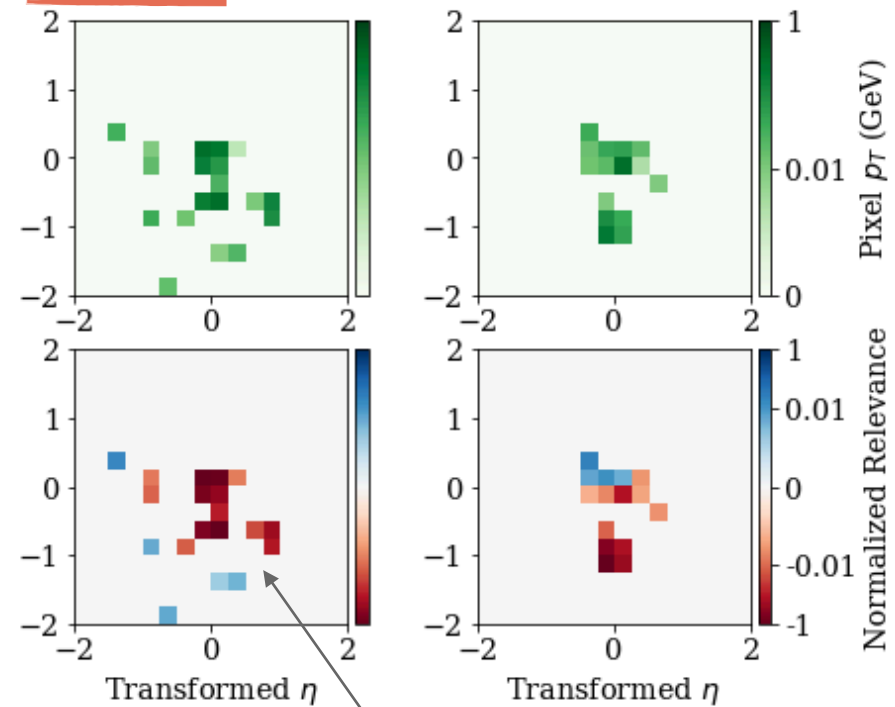
Inputs

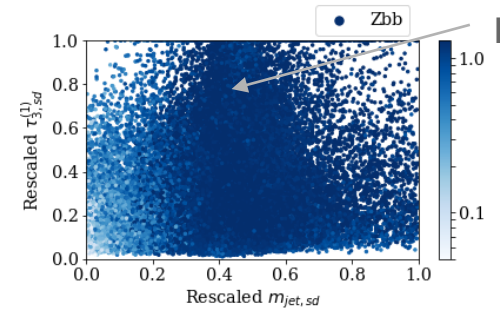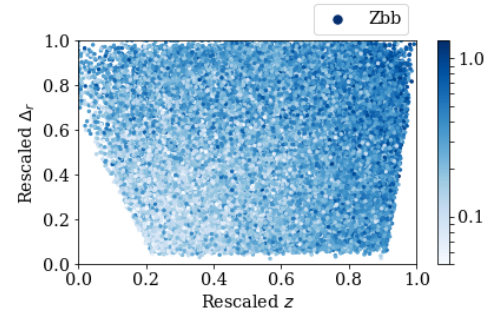Signal Images and their relevance heatmaps

Background Images and their relevance heatmaps
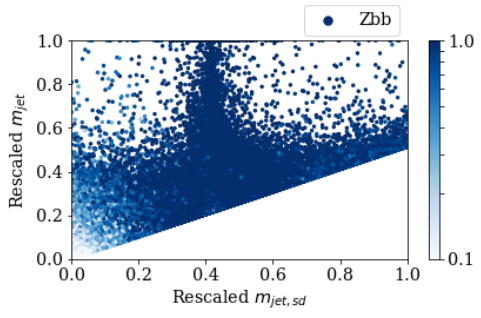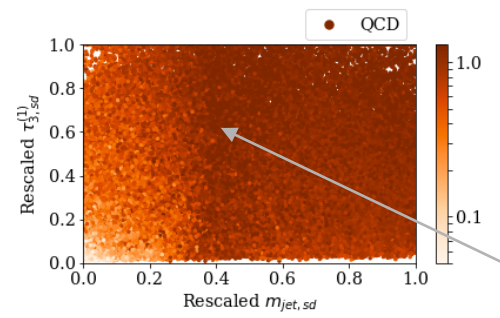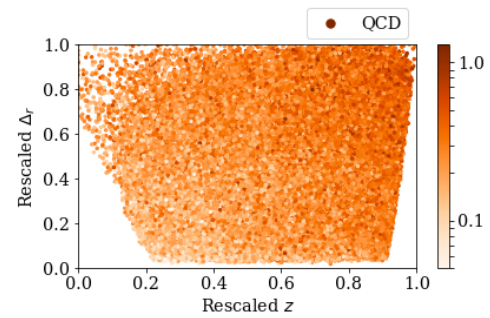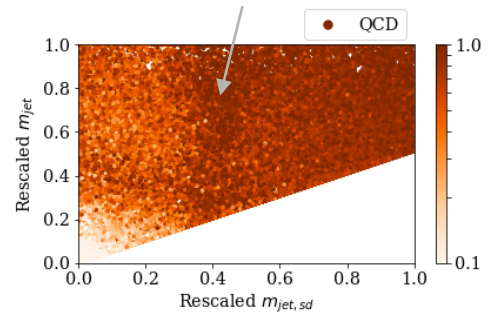
LRP
Results

**Signal:** mostly positive relevance, primarily along $\phi$ axis

Background: mostly negative relevance, more diffuse

ML4Jets2021
July 8th, 2021
Christine McLean

**Darker markers:** higher absolute relevance score



**Decision boundaries:** not as clear as for toy model

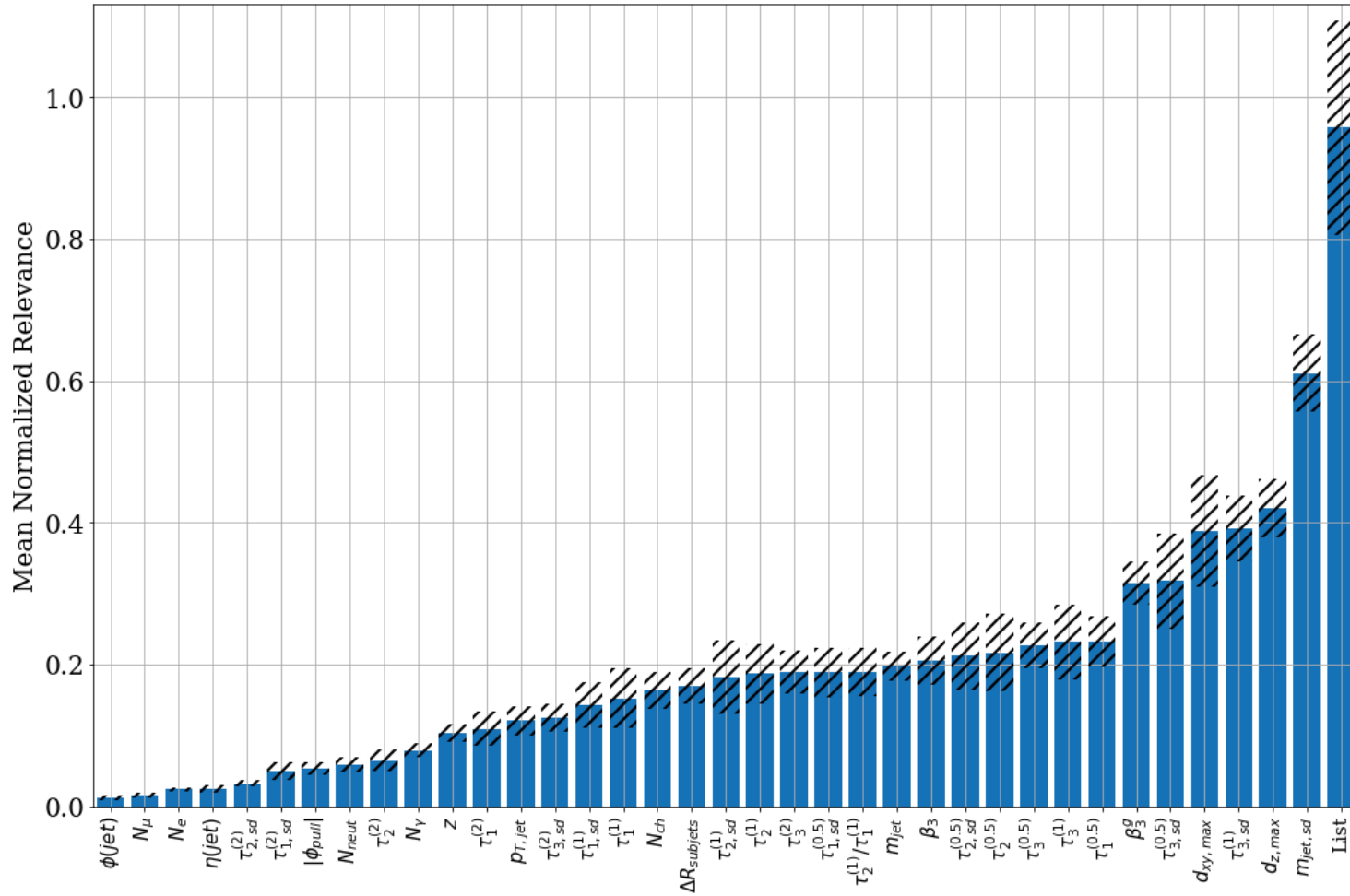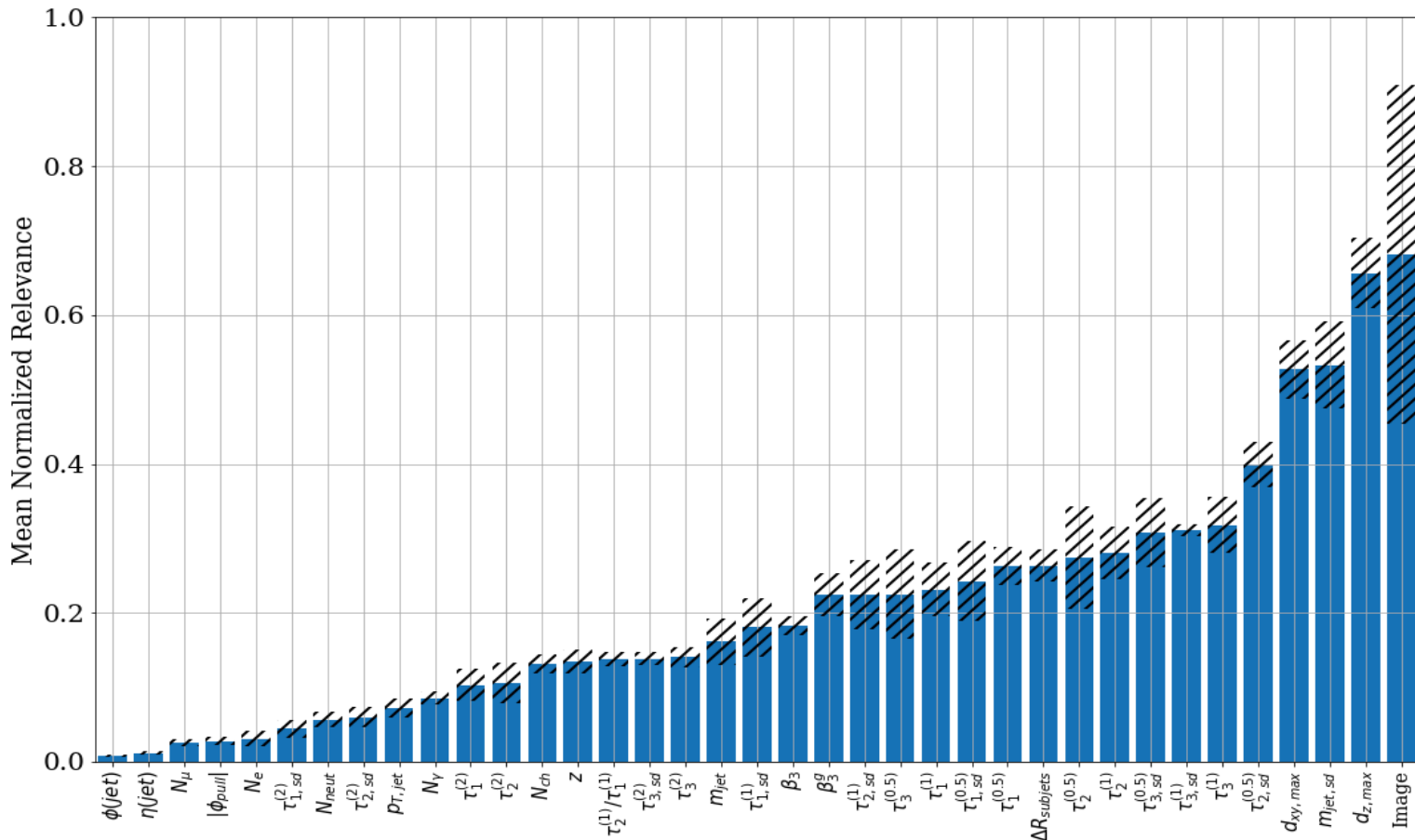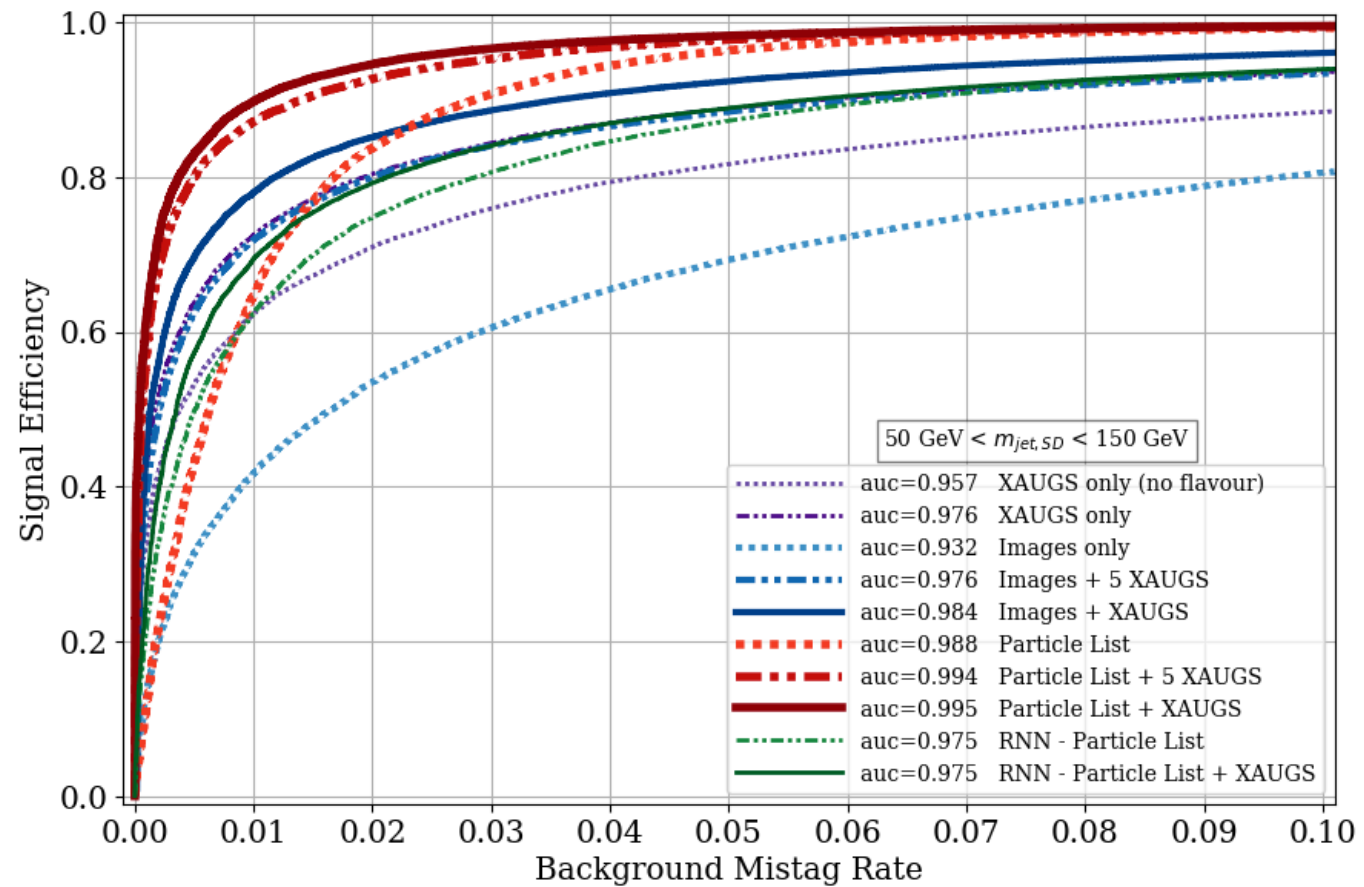## **Particle list:** highest relevance for all models

**Image and $d_{z,max}$:** highest relevance, depending on the model

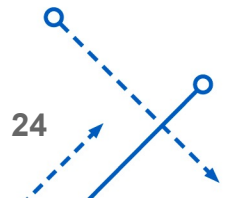# Pythia Results: Model Comparisons

- **Images only does worst**
- **Particle List + XAUGs does best**
  - Particle list + 5 XAUGS comparable

# Conclusions

- **Introduced novel method for ML tagger explainability:** LRP + expert augmented variables
  - Help explain network decisions, and relevant subspaces

- **XAUGs**
  - Can boost classification performance
  - Can entirely capture relevant information of lower-level networks

- **XAUGs + LRP**
  - Can be used to reduce list of network inputs
  - Can be used to quantify numerical uncertainty in DNN training

# ADDITIONAL MATERIAL

# LRP Propagation Rules

- **LRP-z**
  - Redistributes the relevance in proportion to the contributions to the neuron activation.
  - Gradient X Input → Noisy

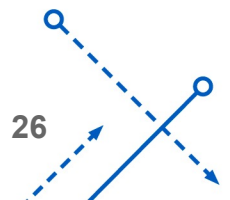$$R_j = \sum_k \frac{a_j w_{jk}}{\sum_{0,j} a_j w_{jk}} R_k$$

- **LRP-$\epsilon$**
  - $\epsilon$ absorbs some relevance for weak and/or contradictory contributions.
  - For large $\epsilon$ only salient explanation factors survive the absorption → Less Noisy
  - **Used in our networks' dense layers**

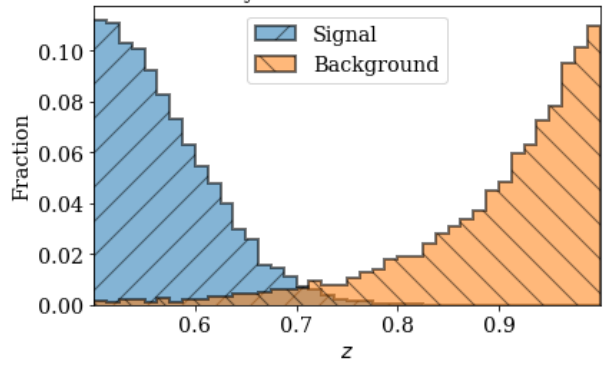$$R_j = \sum_k \frac{a_j w_{jk}}{\epsilon + \sum_{0,j} a_j w_{jk}} R_k$$

- **LRP-$\alpha_1 \beta_0$**
  - Limiting effect on how large positive and negative relevance can grow → Stable Explanations
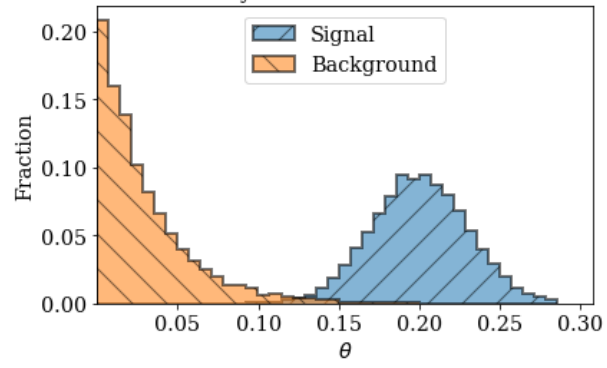  - **Used in our networks' convolution layers**

$$R_j = \sum_k \left( \alpha \frac{(a_j w_{jk})^+}{\sum_{0,j} (a_j w_{jk})^+} - \beta \frac{(a_j w_{jk})^-}{\sum_{0,j} (a_j w_{jk})^-} \right) R_k$$
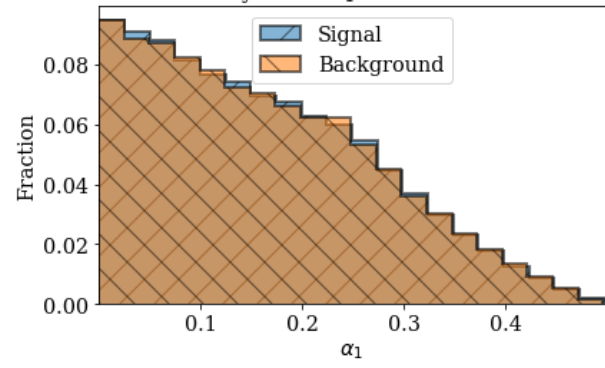
# Toy Model Inputs

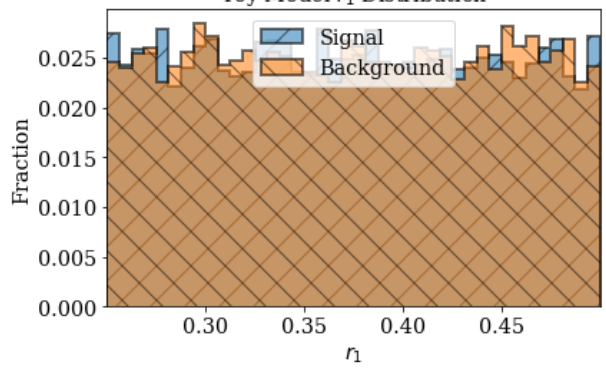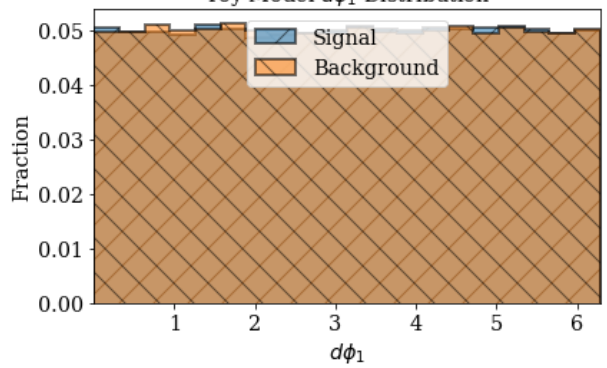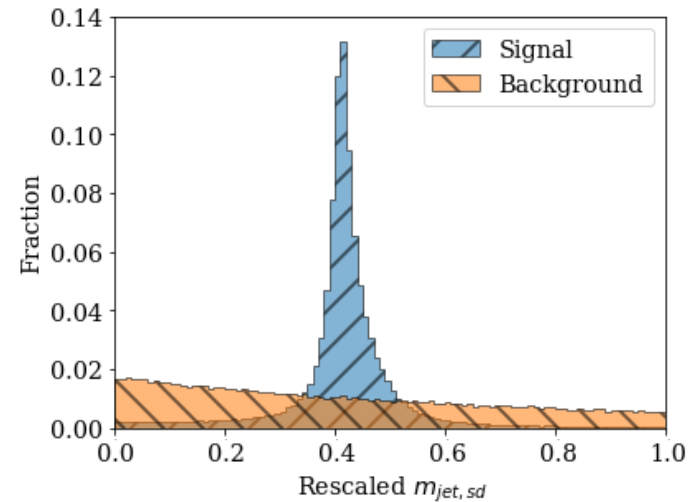| Variable |
|---|
| $log(p_T)$ |
| $log(p_T/p_{T_{jet}})$ |
| $log(E)$ |
| $\lvert \eta \rvert$ |
| $\Delta\phi(jet)$ |
| $\Delta\eta(jet)$ |
| $\Delta R(jet)$ |
| $\Delta R(subjet1)$ |
| $\Delta R(subjet2)$ |
| Charge $q$ |
| isMuon |
| isElectron |
| isPhoton |
| isChargedHadron |
| isNeutralHadron |
| $d_{xy}$ |
| $d_z$ |

# Pythia Model Preprocessing

## 1. Cut on softdrop mass: keep jets with m$_{SD}$ 50-150 GeV
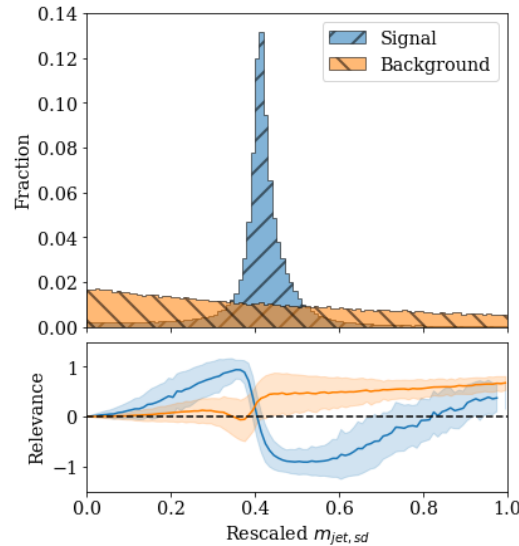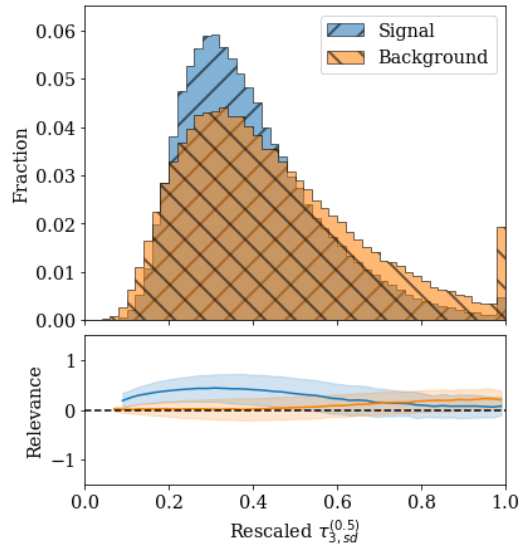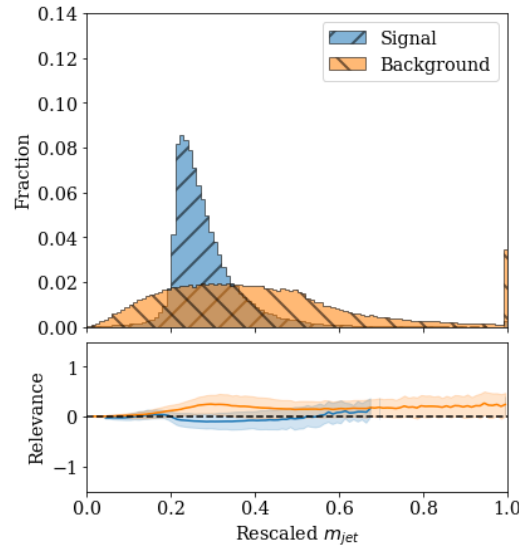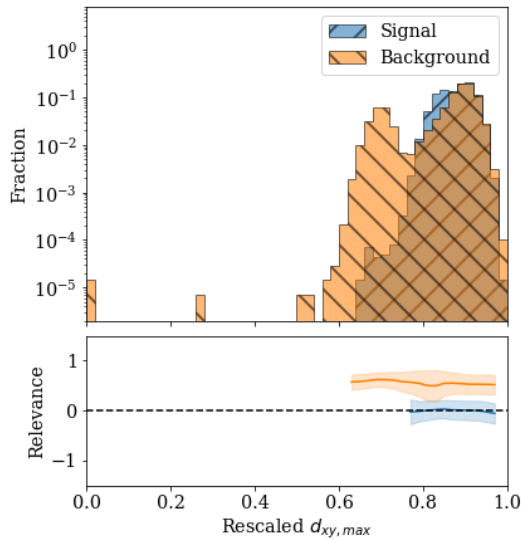## 2. Numerical rescaling
   1. Rebin outliers to *mean + 3(std)* and *mean - 3(std)*
   2. Input distributions are then rescaled from 0 to 1:   $\dfrac{x - x_{min}}{x_{min} - x_{max}}$



Mass cut + rescaling

Profiles don't show clear decision boundary - need higher dimensional plots