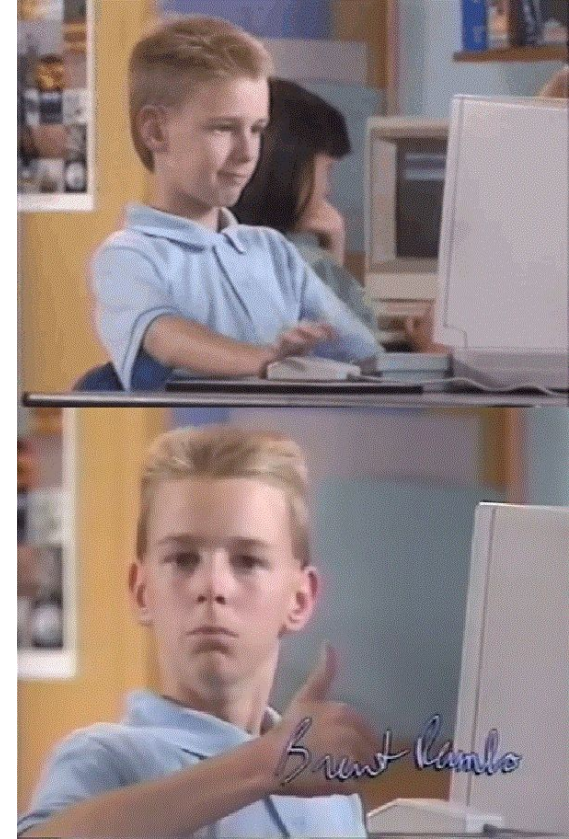# Jet Transformer - Update

Reality:

Expectation:

A Network for quark/gluon
Jet Tagging

# Towards Explainable AI

- Froze network parameters

Consider a jet $\boldsymbol{J}$, i.e. list of $n$ constituents and their $m$ features, and a NN:

$$\text{NN} : \mathbb{R}^{n \times m} \to \mathbb{R}^2$$

$$\begin{bmatrix} \Delta\eta_1 & \Delta\phi_1 & \sigma_1 & \sigma_1 & m_1^2 & p_{\text{T}\,1} & q_1 & \text{PID}_1 \\ \Delta\eta_2 & \Delta\phi_2 & \sigma_2 & \sigma_2 & m_2^2 & p_{\text{T}\,2} & q_2 & \text{PID}_2 \\ \Delta\eta_3 & \Delta\phi_3 & \sigma_3 & \sigma_3 & m_3^2 & p_{\text{T}\,3} & q_3 & \text{PID}_3 \\ & & & \vdots & & & & \end{bmatrix} \mapsto \begin{bmatrix} \text{class}_1 & \text{class}_2 \end{bmatrix}$$

Now, calculate:

$$\frac{\partial \text{class}_1}{\partial \boldsymbol{J}} = \begin{bmatrix} \frac{\partial \text{class}_1}{\partial \Delta\eta_1} & \frac{\partial \text{class}_1}{\partial \Delta\phi_1} & \cdots \\ \frac{\partial \text{class}_1}{\partial \Delta\eta_2} & \frac{\partial \text{class}_1}{\partial \Delta\phi_2} & \cdots \\ \vdots & \vdots & \ddots \end{bmatrix}$$

Then, look jets which are predicted to be class 2, but their true label is class 1. We now want to increase the classification output score for class 1, so we calculate

$$\boldsymbol{J}' = \boldsymbol{J} + \text{scale} \cdot \frac{\partial \text{class}_1}{\partial \boldsymbol{J}}$$

which has a higher class 1 output score.

* sorry for the colour mismatch

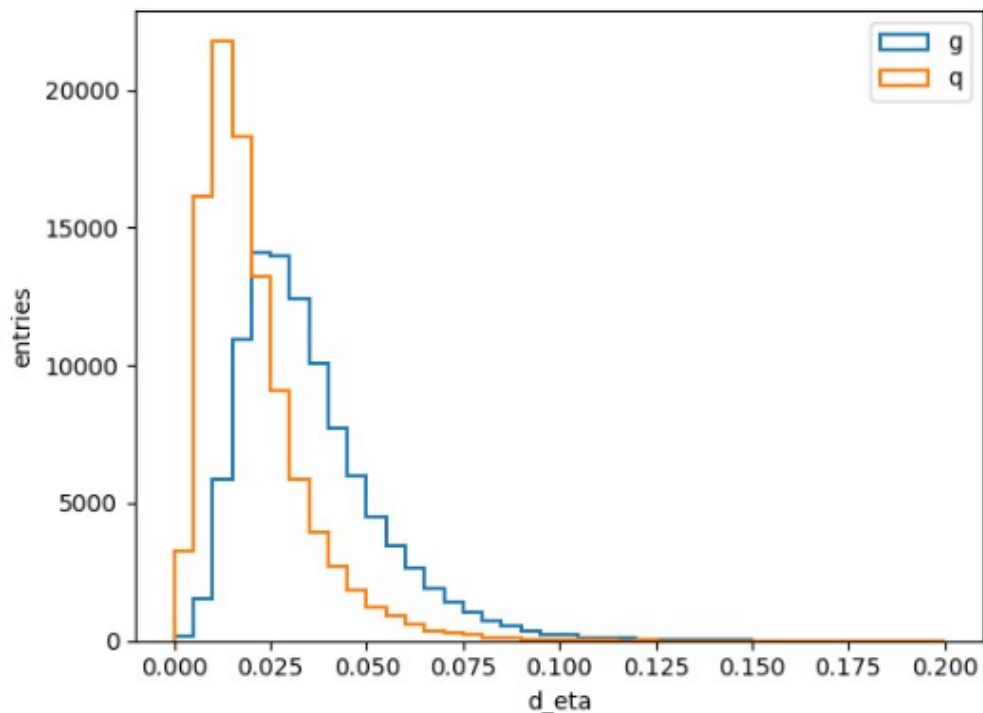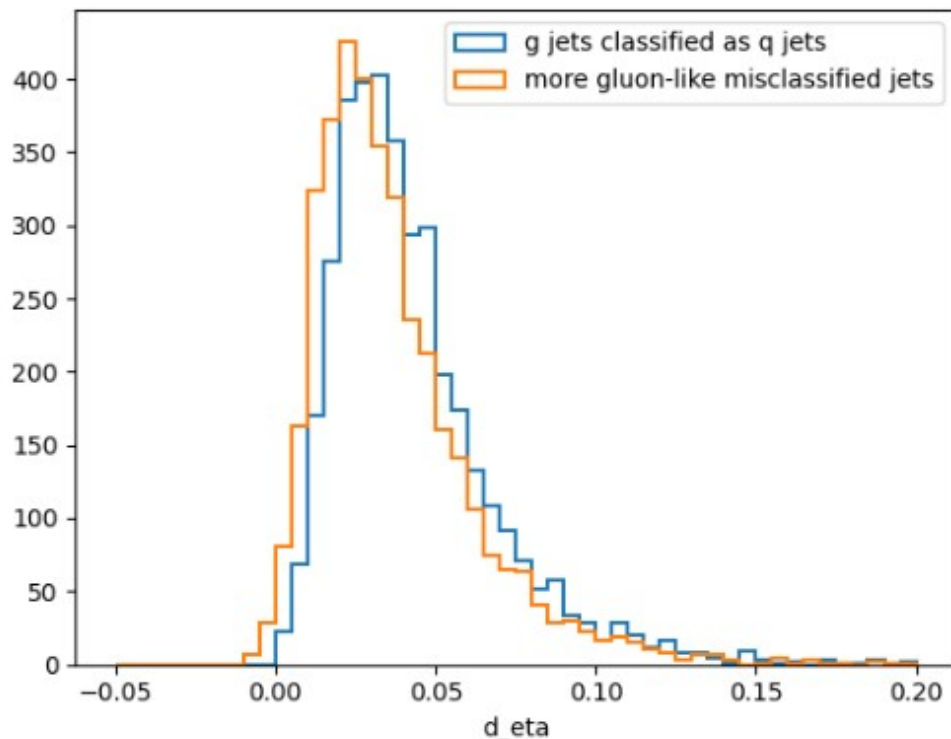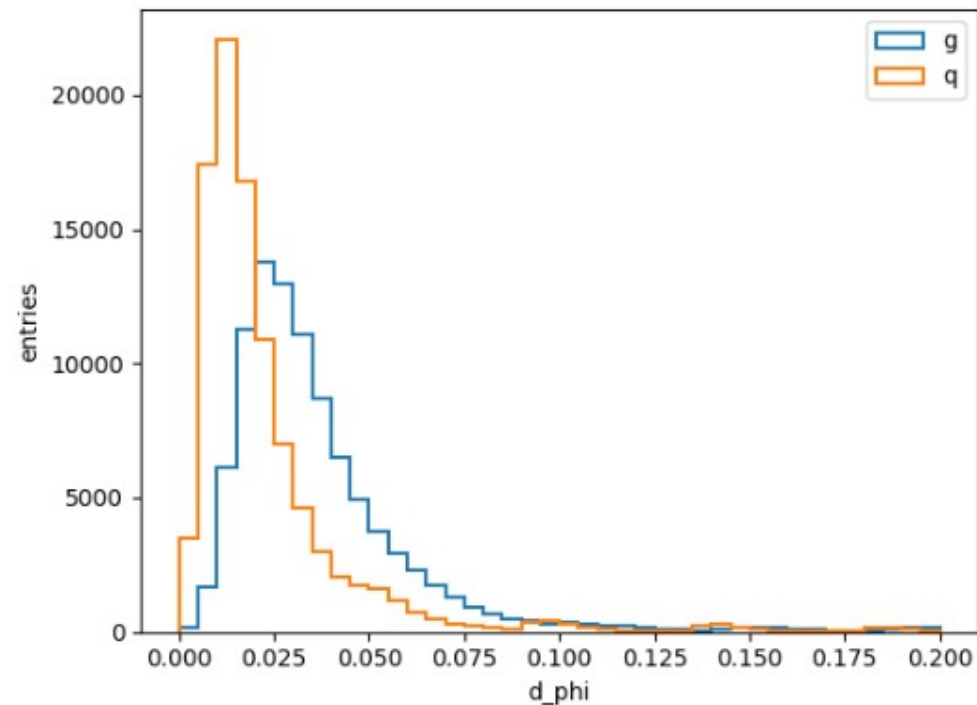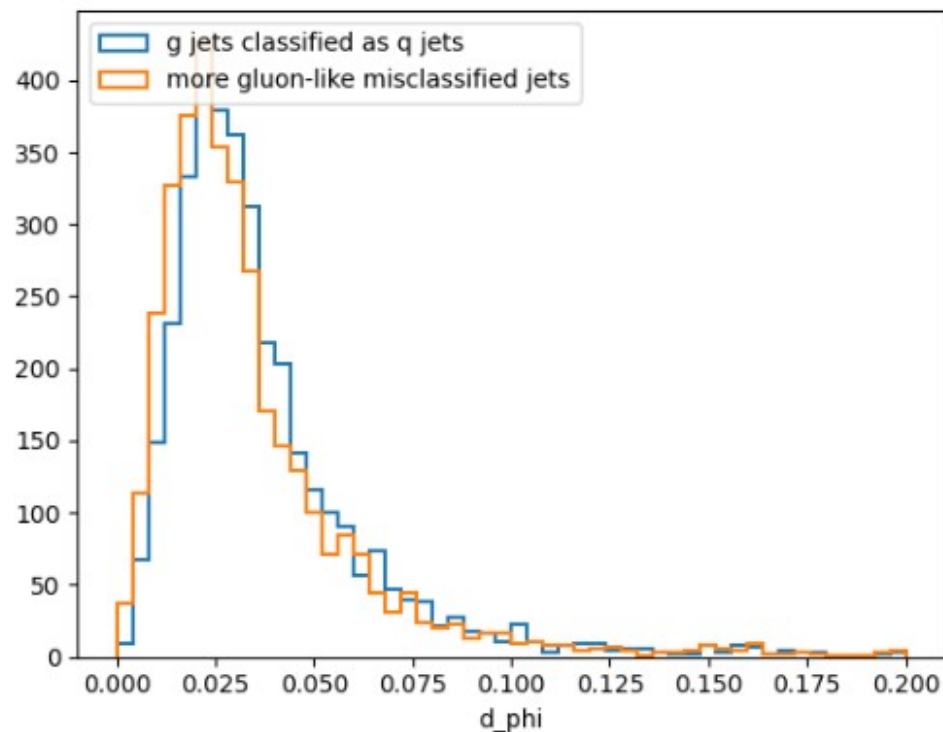** scale=1 here

*** NN trained for 3 epochs

# Does the network pick up Characteristics of q/g Jets?

**DATA:**

**NN:**

* sorry for the colour mismatch
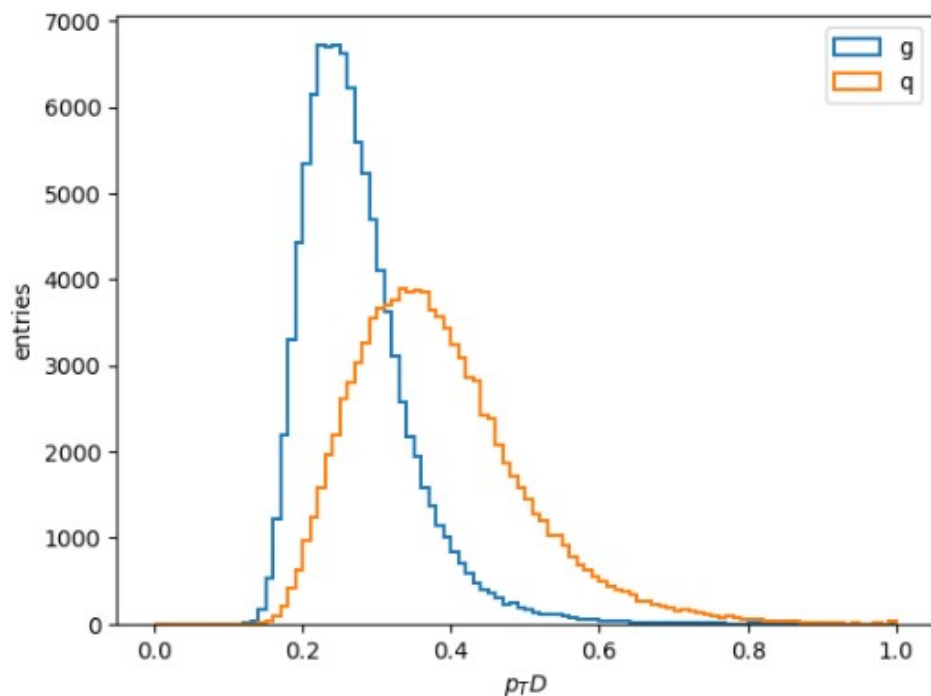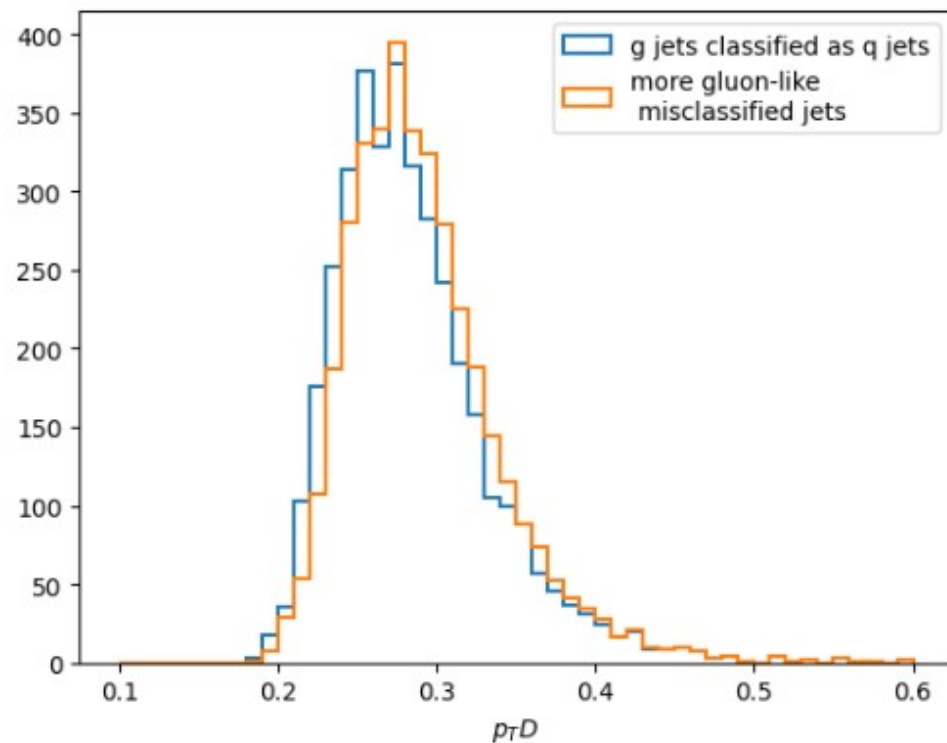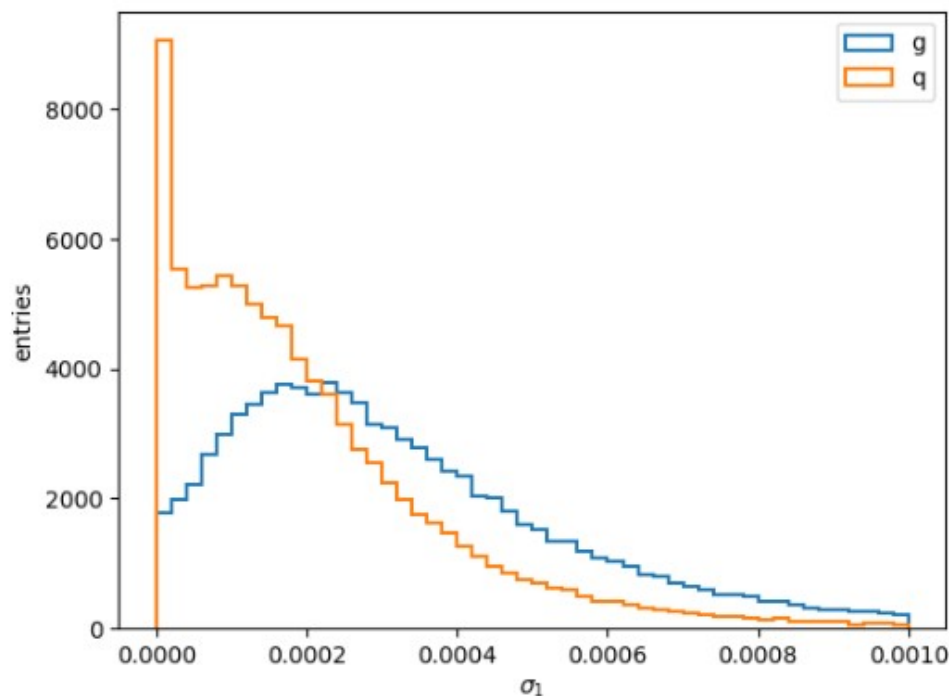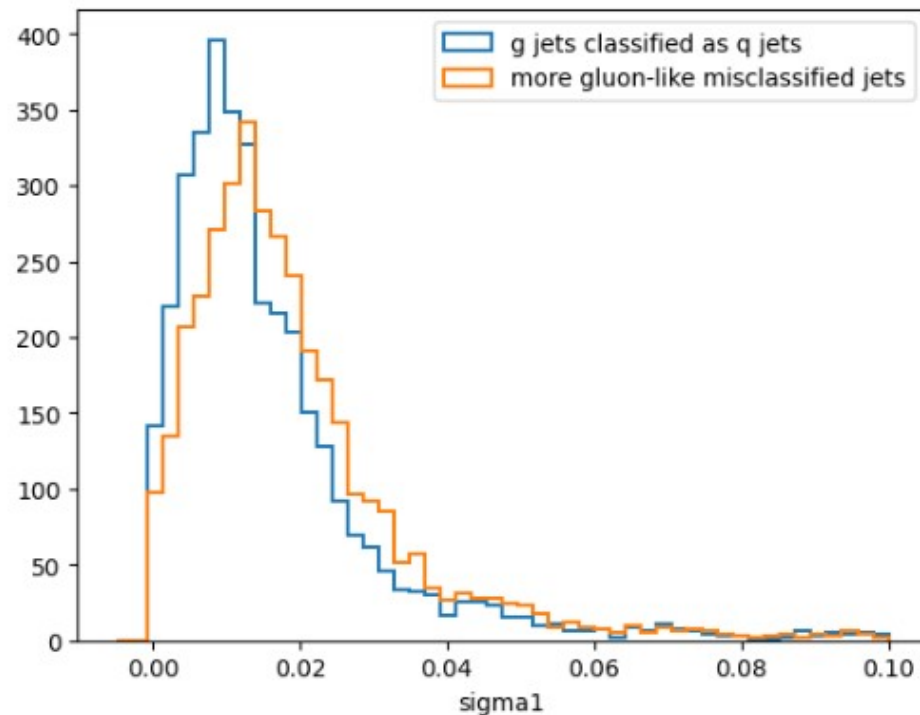
** scale=1 here
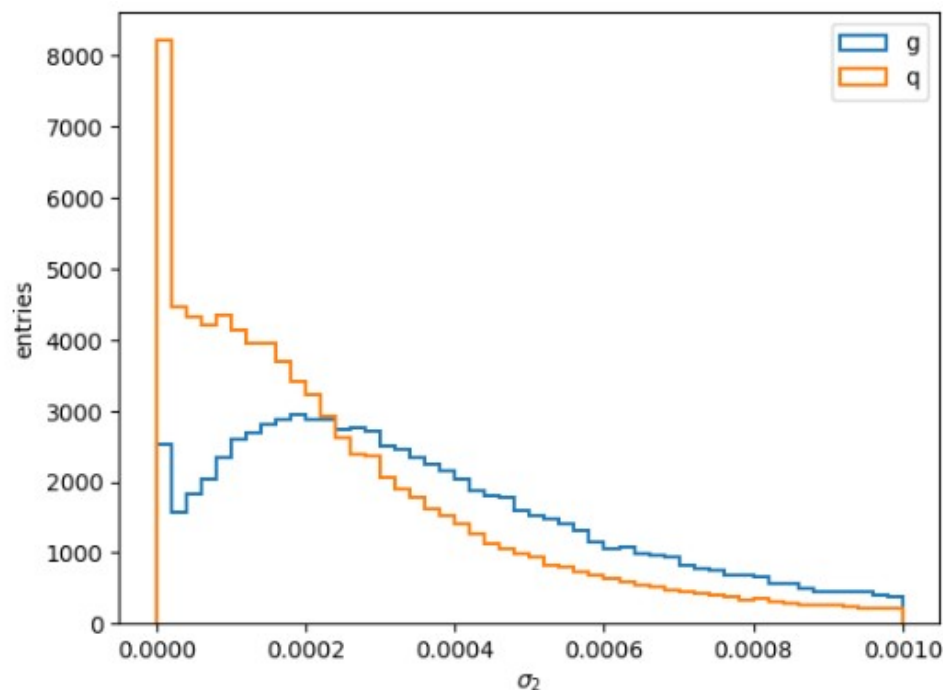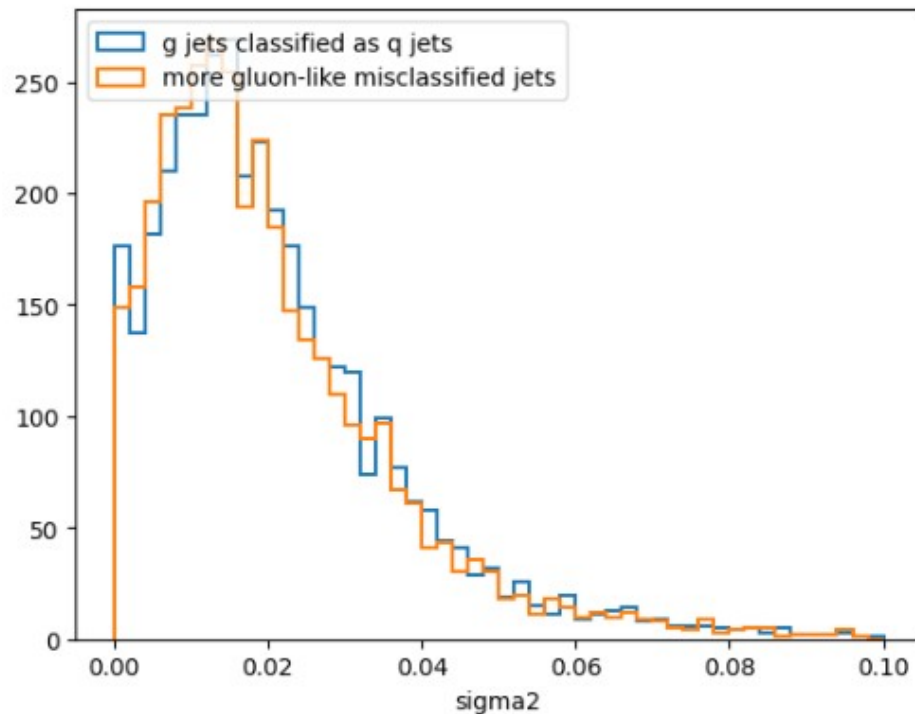
*** NN trained for 3 epochs

# Does the network pick up Characteristics of q/g Jets?

**DATA:**

**NN:**

* sorry for the colour mismatch

** scale=1 here
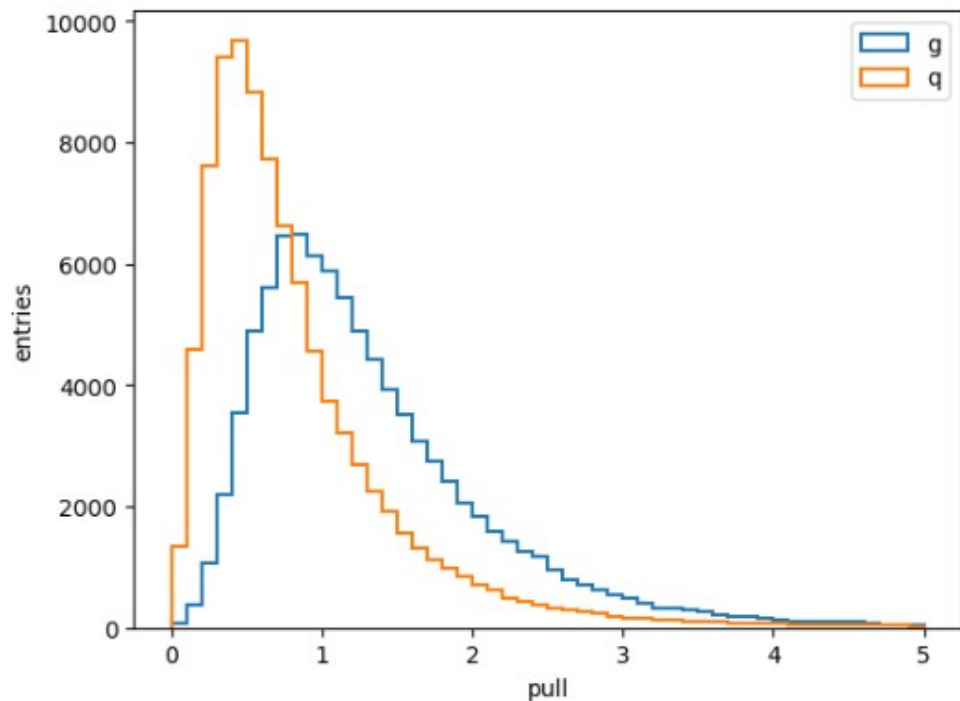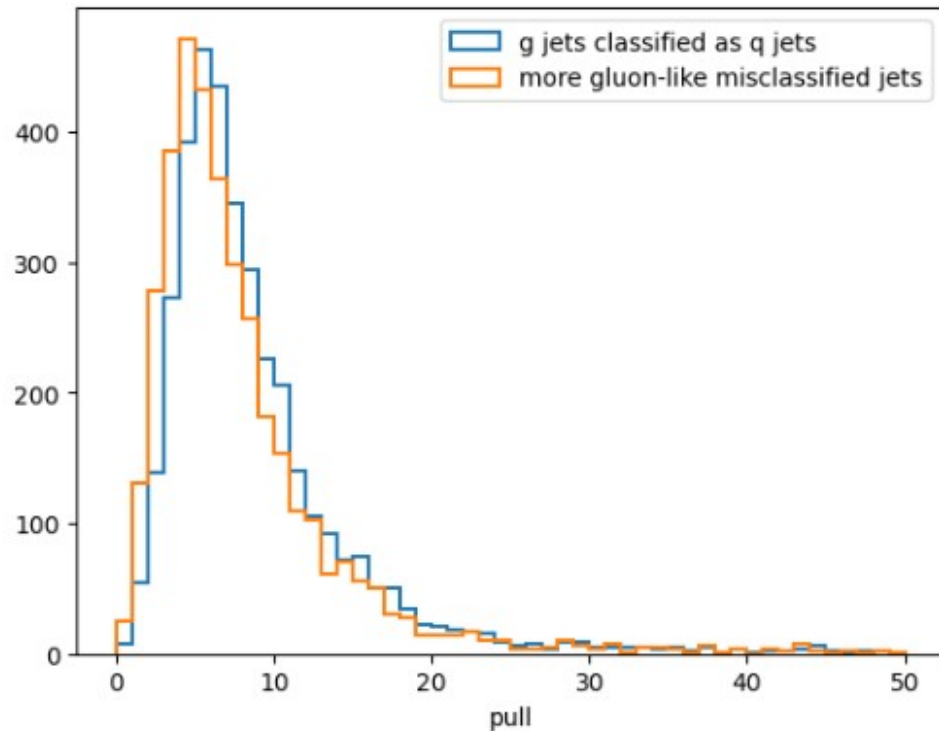
*** NN trained for 3 epochs

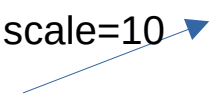# Does the network pick up Characteristics of q/g Jets?

**DATA:**

**NN:**

# To Do:



scale=10
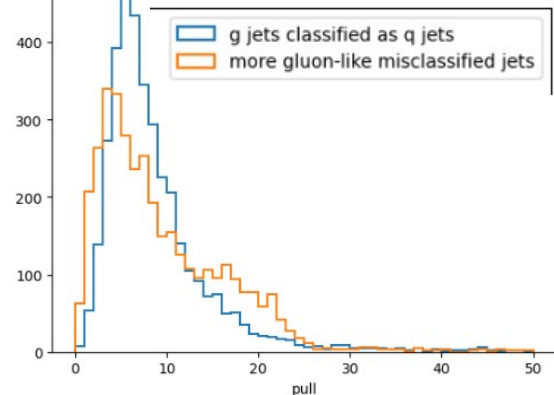
- What is the optimal value for 'scale'

- How does the classification output score vary with the 'scale'

- What feature does the network 'want to change the most'?

- Apply this to TOF Transformer

  → plot vector field of gradients for showers!

# Bonus



Pearson product-moment correlation coefficients