

# VAEs for Anomalous Jet Tagging

Taoli Cheng  
Mila, University of Montreal



In collaboration with

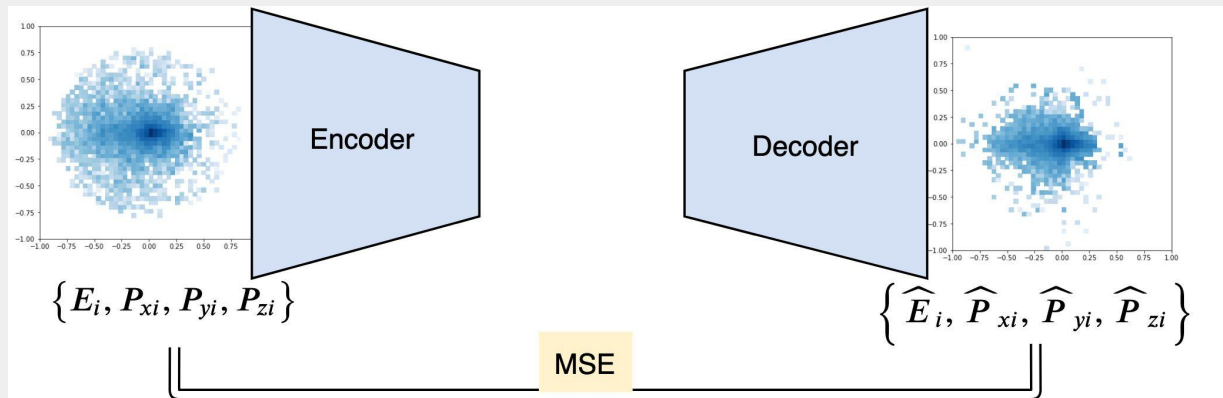
Jean-Francois Arguin, Julien Leissner-Martin, Jacinthe Pilette (Universite de Montreal (CA))  
Tobias Golling, Takuya Nobe, Johnny Raine (Universite de Geneve (CH))  
Amir Farbin, Debottam Bakshi Gupta (University of Texas at Arlington (US))

Jan. 16, 2020  
ML4Jets, NYU

# Review: Generative Model for Anomaly Detection

- Explicit (e.g. VAE) or implicit (e.g. GAN) estimation of  $\log p(x)$
- Building Blocks
  - Model
  - Embedding Architecture
  - Anomaly Metric

## Autoencoder for Anomalous Jet Tagging

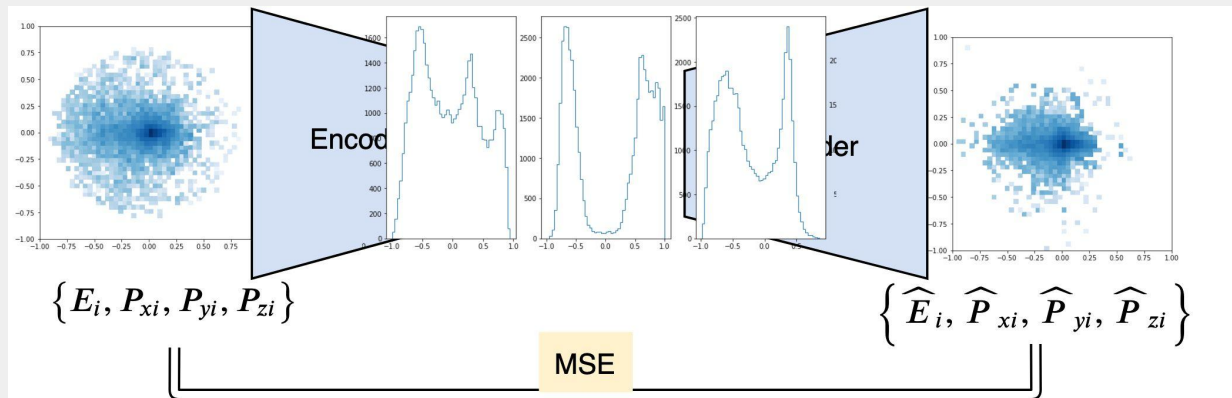


- T. HeimeI , G. Kasieczka , T. Plehn , and J. M Thompson. QCD or What? arXiv:1808.08979.
- M. Farina, Y. Nakai, D. Shih. Searching for New Physics with Deep Autoencoders. arXiv: 1808.08992
- Tuhin S. Roya and Aravind H. Vijayb. A robust anomaly finder based on autoencoder. arXiv: 1903.02032

# Review: Generative Model for Anomaly Detection

- Explicit (e.g. VAE) or implicit (e.g. GAN) estimation of  $\log p(x)$
- Building Blocks
  - Model
  - Embedding Architecture
  - Anomaly Metric

## Autoencoder for Anomalous Jet Tagging



- T. Heimpl, G. Kasieczka, T. Plehn, and J. M Thompson. QCD or What? arXiv:1808.08979.
- M. Farina, Y. Nakai, D. Shih. Searching for New Physics with Deep Autoencoders. arXiv: 1808.08992
- Tuhin S. Roy and Aravind H. Vijayb. A robust anomaly finder based on autoencoder. arXiv: 1903.02032

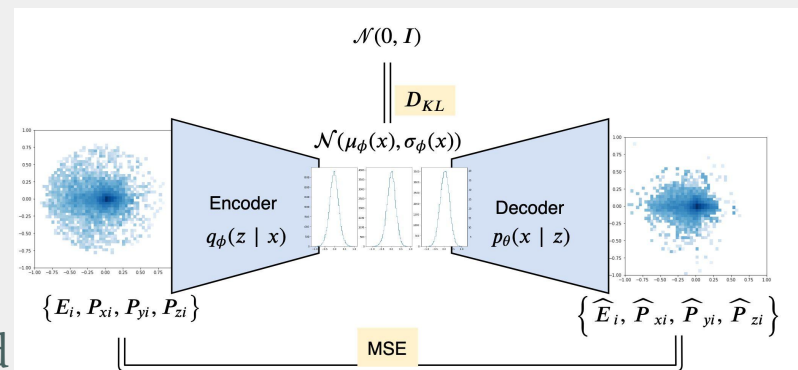
# From Autoencoder to Variational Autoencoder

- VAE:
  - enforce a prior distribution in the Latent space through a  $D_{KL}$  (Kullback-Leibler Divergence) term
$$D(q(x)||p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$

(regularization term to autoencoder)

- Likelihood estimation:  
 $\text{logp}(x) > -L_{\text{VAE}}$  : Evidence Lower Bound

$$q_\phi(z|x) = \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x)) \quad p(z) = \mathcal{N}(0, I)$$

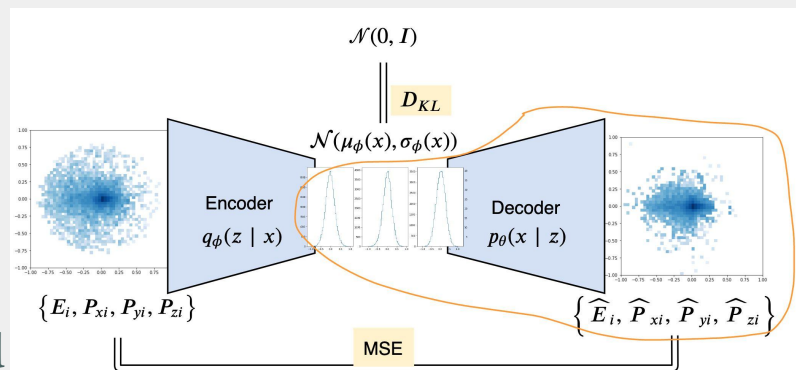


$$L = \frac{1}{4n} \sum_i \|\hat{x}_i - x_i\|_2^2 + \beta D_{KL}(q(z|x)||p(z))$$

# From Autoencoder to Variational Autoencoder

- VAE:
    - enforce a prior distribution in the Latent space through a  $D_{KL}$  (Kullback-Leibler Divergence) term
- $$D(q(x)||p(x)) = \sum_x q(x) \log \frac{q(x)}{p(x)}$$
- (regularization term to autoencoder)
- Likelihood estimation:  
 $\log p(x) > -L_{VAE}$  : Evidence Lower Bound

$$q_\phi(z|x) = \mathcal{N}(\mu_\phi(x), \Sigma_\phi(x)) \quad p(z) = \mathcal{N}(0, I)$$



$$L = \frac{1}{4n} \sum_i \|\hat{x}_i - x_i\|_2^2 + \beta D_{KL}(q(z|x)||p(z))$$

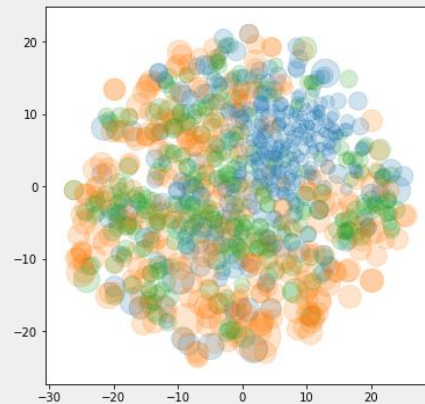
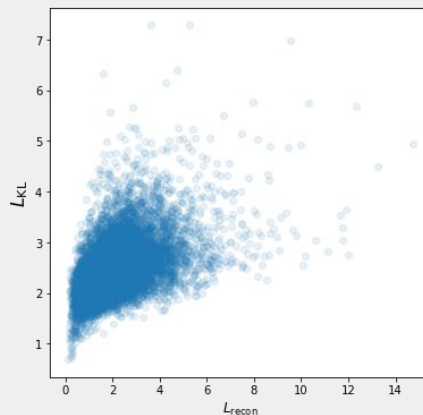
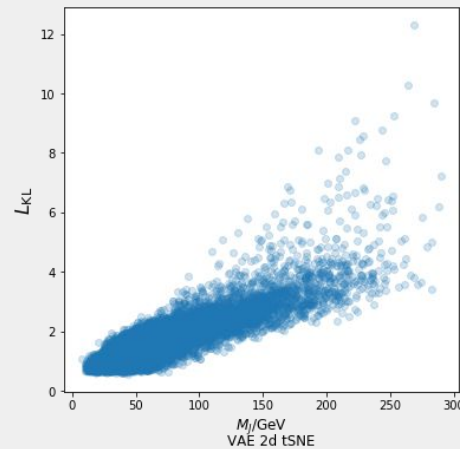
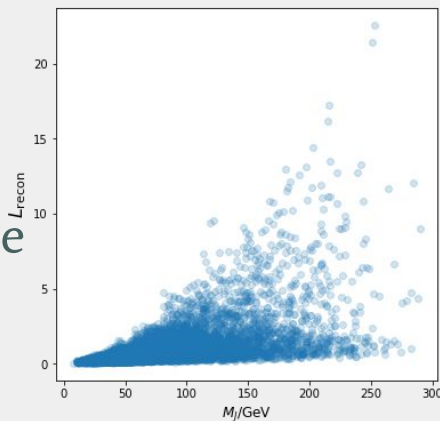
- Generative model:  
Sample from gaussians  $\rightarrow$  generate new jets
- Anomaly Detection: distance in input space (Mean Squared Error (MSE)); latent space (KL divergence)

# Settings

- Simple FCN/LSTM architecture
- Taking the first 20 pt-ordered jet constituents (zero-padded)
  - Inputs: four vectors (E, P<sub>x</sub>, P<sub>y</sub>, P<sub>z</sub>) of jet constituents (particle flow objects)
  - Preprocessing: Boost to jet rest frame, Centering, Rotating → Principal axis alignment
- Train on 600,000 QCD jets (of which 20% serve as validation set)
  - QCD dijet production:  $pp \rightarrow jj$
  - ATLAS fatjet trigger:  $R = 1.0$  antikt jets,  $p_T > 450$  GeV
  - No trimming applied
- VAE
  - $d_{\text{hidden}} = 10$
  - $\beta = 0.1, 0.5, 1, 5$

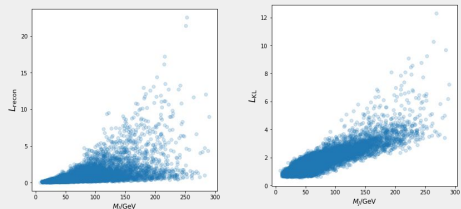
# KL Regularization Strength -- beta = 0.1

- MSE and KL both has jet mass correlation
- Very strong mass correlation in latent space

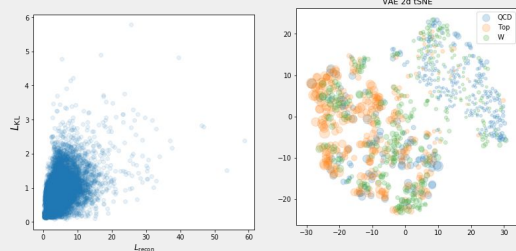
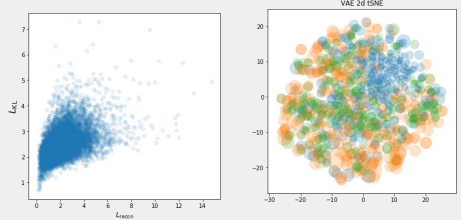
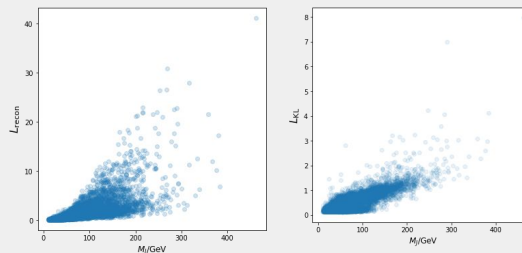


# KL Regularization Strength

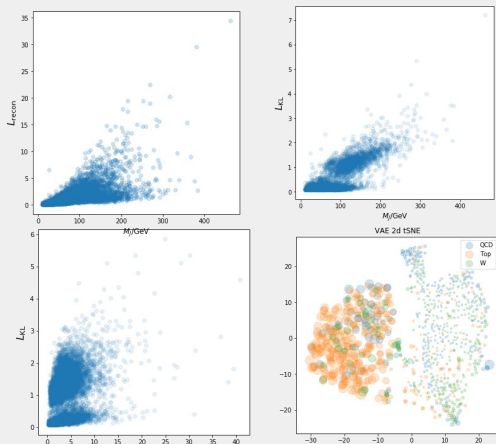
beta=0.1



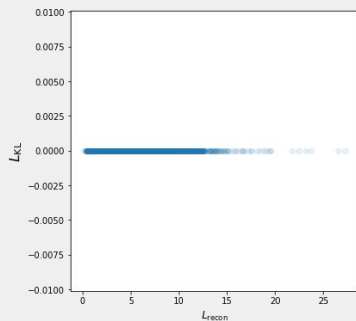
beta=0.5



beta=1



beta>~1



- As beta increases, reconstruction performance decreases. Latents develop different modes (mass modes).
- KL vanishing for  $\beta > \sim 1$



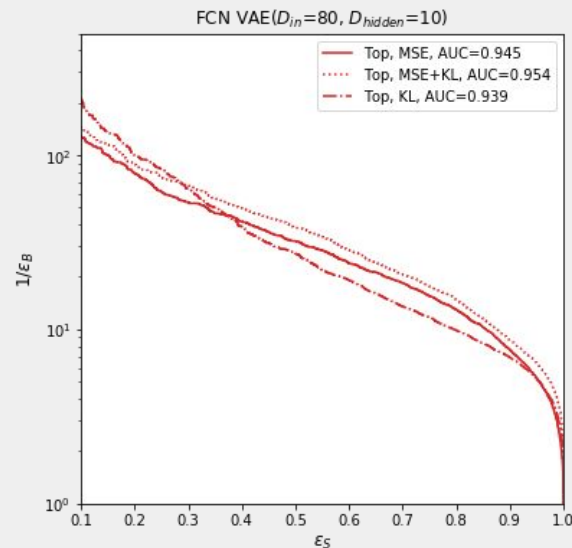
# Jet Tagging Performance

- Datasets:
  - Background: QCD,  $R=0.1$ ,  $p_T > 450$  GeV
  - Testing on  $p_T$  : [550, 650] GeV
    - 2-prong:  $W$  ( $W' \rightarrow W Z$ )
      - $M$ : 59 GeV, 80 GeV, 120 GeV
    - 3-prong: Top ( $Z' \rightarrow t t\bar{}$ )
      - $M$ : 80 GeV, 174 GeV
    - $H(->hh->4j)$ ,  $M_H=174$  GeV
      - $M_h = 20$  GeV, 80 GeV

- Anomaly Metric

- Examine:

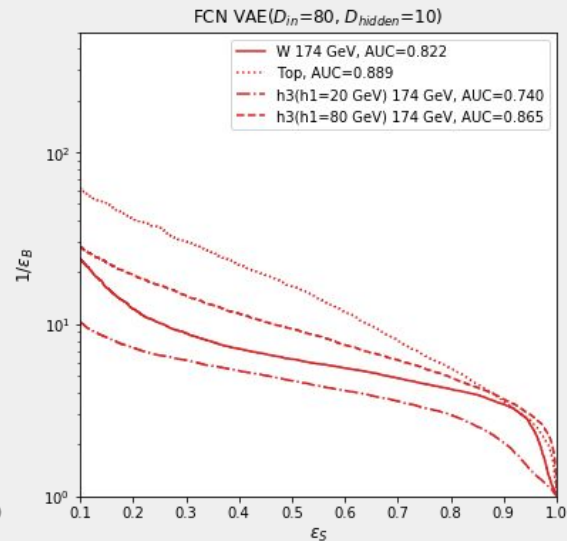
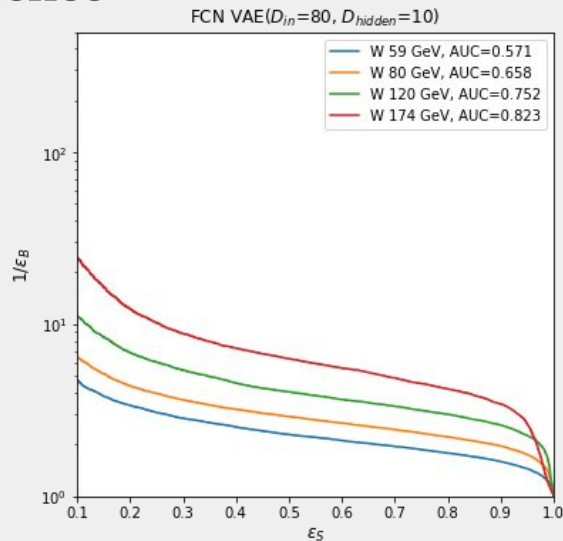
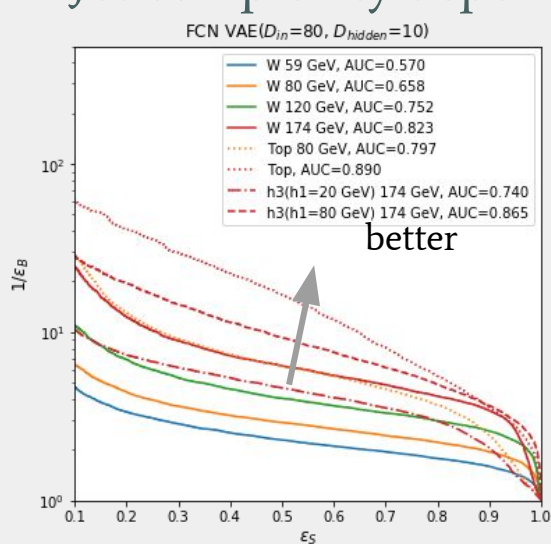
- Jet mass effects
- $p_T$  effects (training on full  $p_T$ , test on fixed  $p_T$ )
- Jet type (focus on prong-ness)



Test results for Top Tagging Reference Data  
(\* T. Heime1 , G. Kasieczka , T. Plehn , and J. M. Thompson. QCD or What? arXiv:1808.08979.)

# Jet Tagging Performance

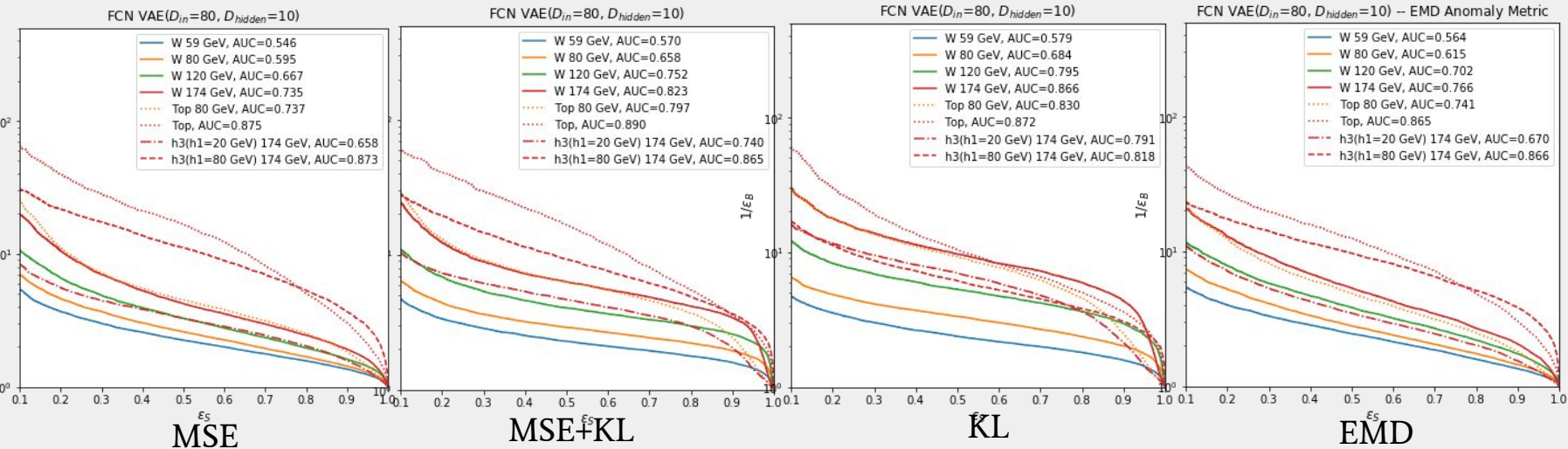
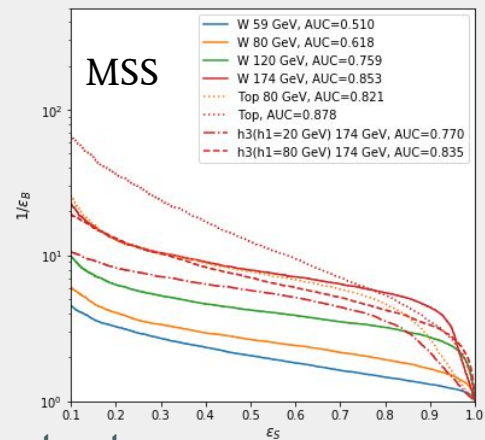
- Jet mass dependence: discriminative power decreases when jet mass decreases (mass correlation) ---> works well for Top, but low significance for W jets
- Jet complexity dependence



# Jet Tagging Performance

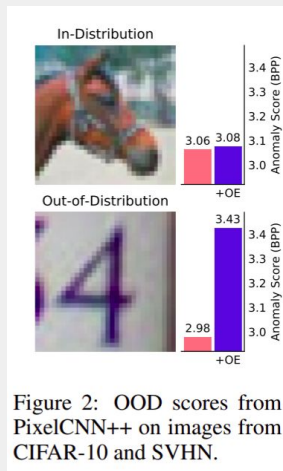
- Anomaly Metric

- Reconstruction error: MSE
- Negative Log Likelihood: MSE+KL
- Latent space: KL divergence
- EMD(Energy Mover Distance) between inputs and outputs
- MSS:  $l_2$  norm the input feature vector ( $\chi^2$ )

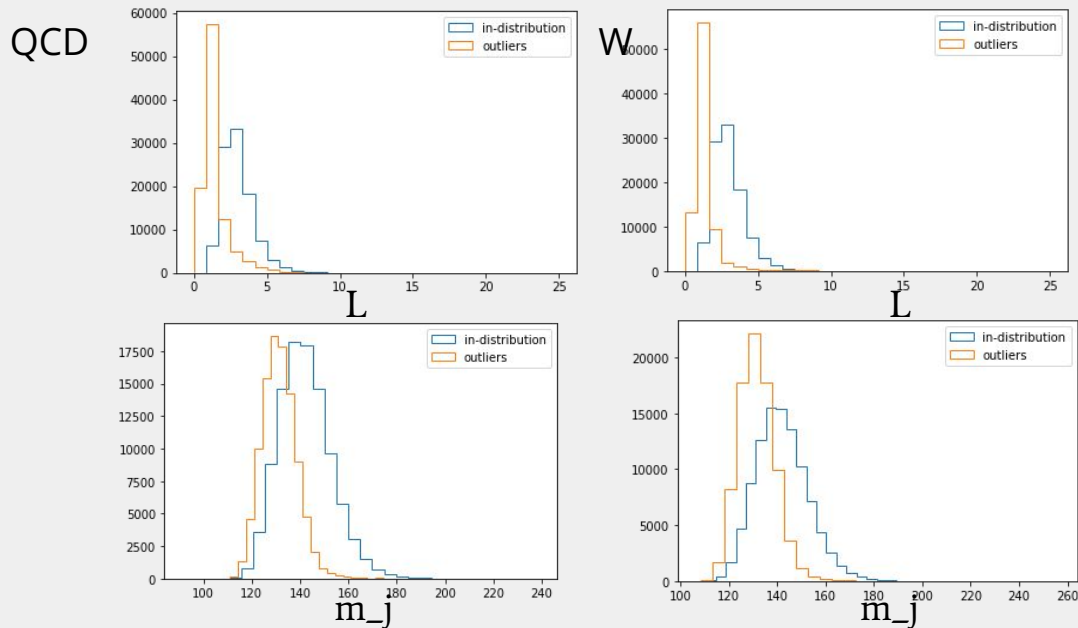


# Anomaly Detection can Fail

- Outliers can be assigned higher probability sometimes, this happens in a general scope of anomaly detection using generative models
- Quick example: MSE based anomaly metric has intrinsic mass dependence  $\rightarrow$  naive VAE assigns higher probability to lower mass jets



- D. Hendrycks, M. Mazeika, T. Dietterich.  
Deep Anomaly Detection with Outlier Exposure.  
arXiv: 1812.04606



# Outlier Exposure (OE)

- Is the VAE learning useful enough representations?
  - Restricted by the format of loss function
  - Need extra information to guide directions for better anomaly detection
- Semi-supervised Learning: encourage specific directions in the loss landscape
  - Relative weight lambda controls the OE strength  $L = L_{VAE} - \lambda L_{OE}$
  - Restricting reconstruction error strength between Out-of-distribution (OoD) and In-distribution (InD) samples

$$L_{OE} = \text{sigmoid}(MSE(OoD) - MSE(InD))$$

- Restricting KL divergence in latent space

$$L_{OE} = \min\{0, KL(OoD) - KL(InD) - \text{margin}\}$$

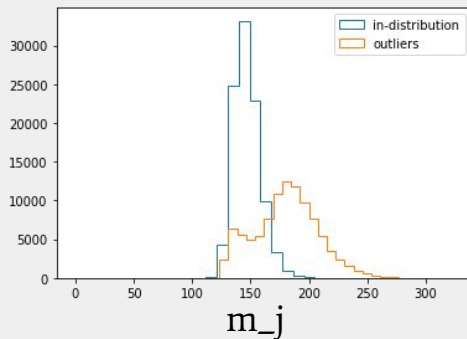
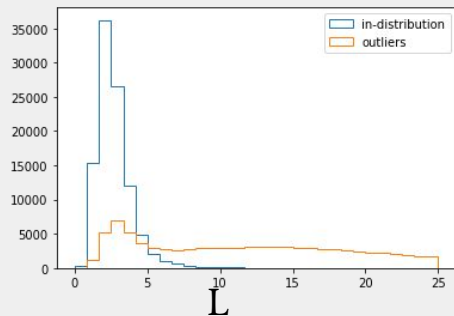
- Training Scenarios:
  - Fine-tuning using outlier exposure
  - Train with outlier exposure from scratch (results shown for this scenario)

# Quick Test -- Top Tagging Reference Data -- OE(QCD) Training on Top

- Quick test on training VAE with top jets, using QCD as outlier exposure samples → test on QCD and W jets

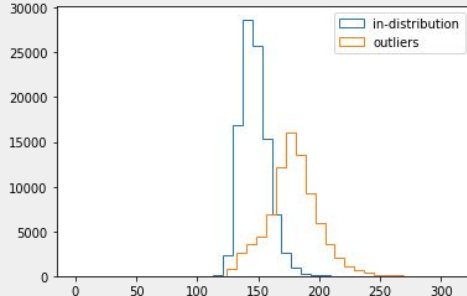
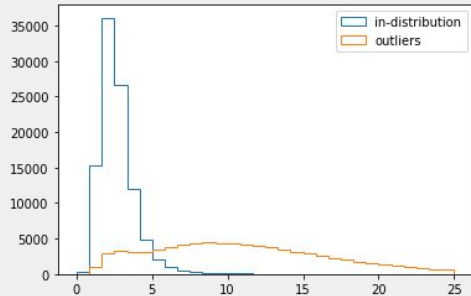
AUC = 0.920

QCD



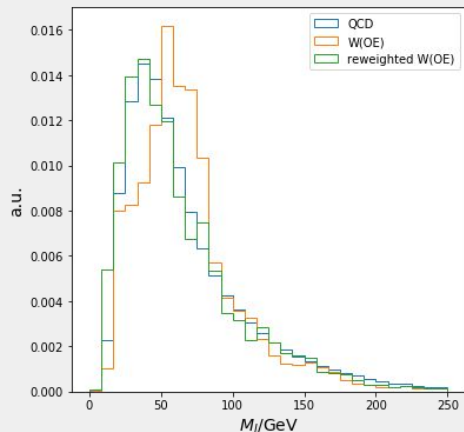
AUC = 0.939

W



# Outlier Exposure -- Results

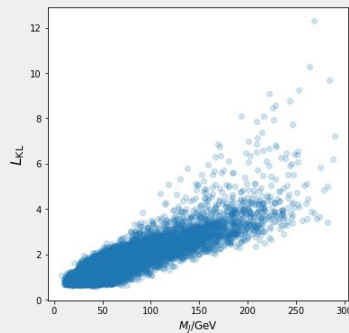
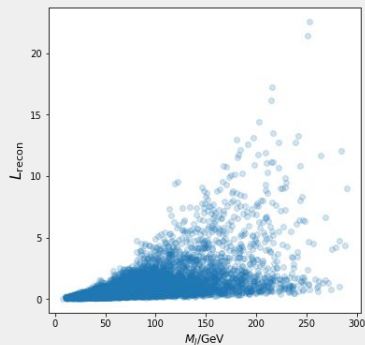
- Outlier samples:  $W$  (mass rescaled) jets with mass distribution reweighted to match QCD jets ( $\rightarrow$  mass decorrelation)
- Annealing training of OE weight lambda (cyclically annealing, from 0 - 2)
- Results:
  - Test on different jet type and jet mass
  - Mass decorrelation



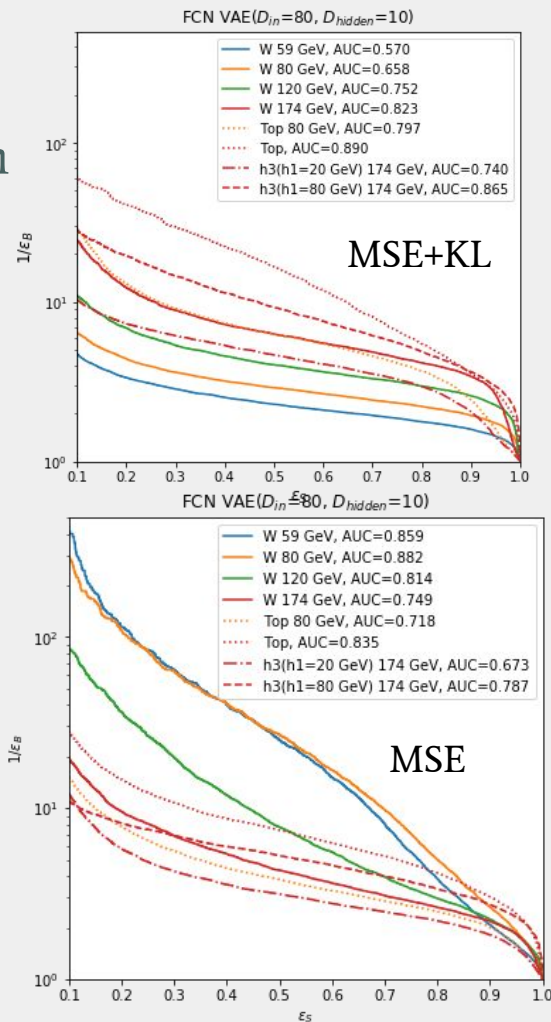
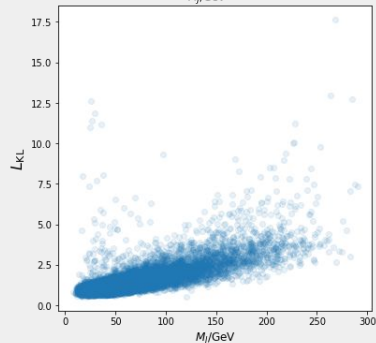
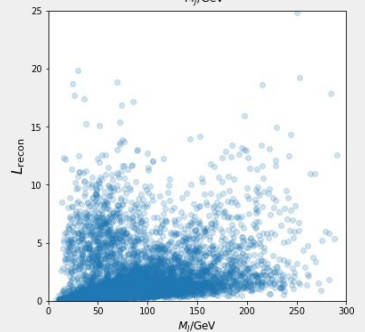
# Outlier Exposure -- Results -- MSE-OE

- Train using MSE outlier exposure loss term
- Anomaly metric: MSE

before



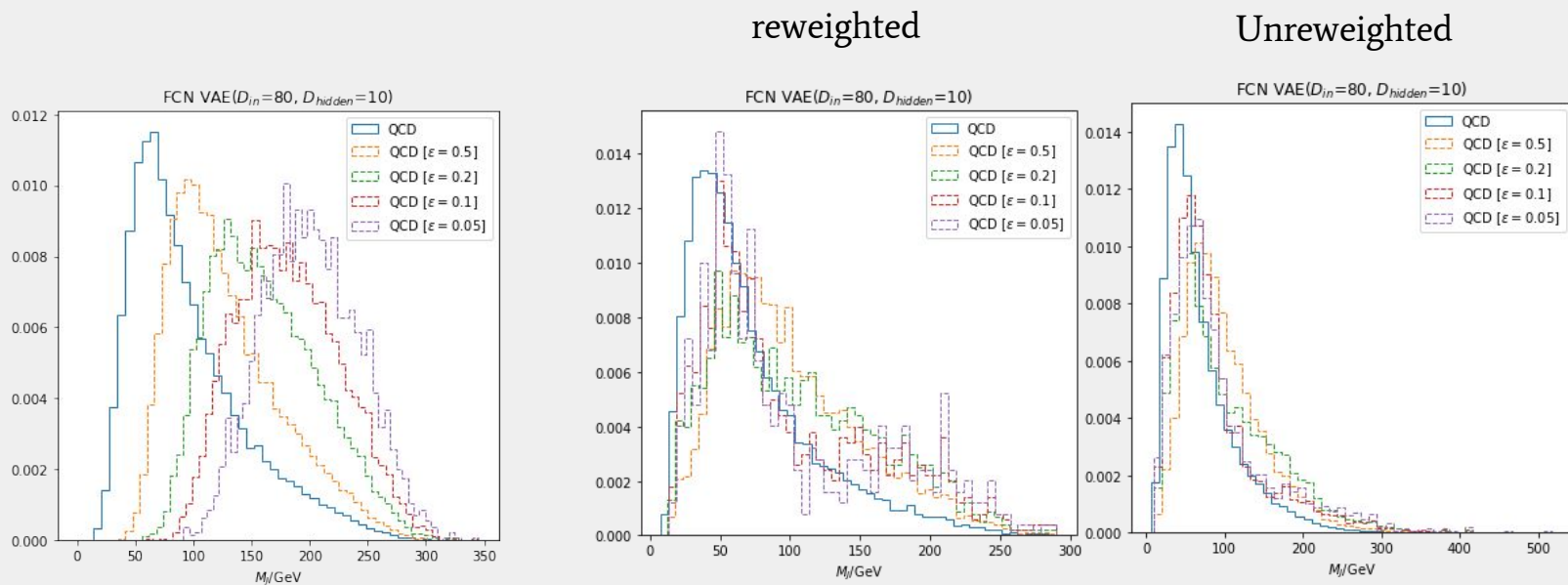
after





# Outlier Exposure -- Results -- MSE-OE

- Train using MSE outlier exposure loss term
- Anomaly metric: MSE+KL (since KL and MSE are correlated; MSE+KL better in MSE-OE case)

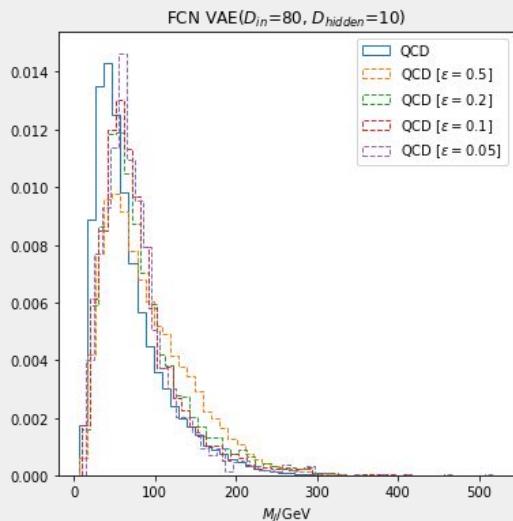


before

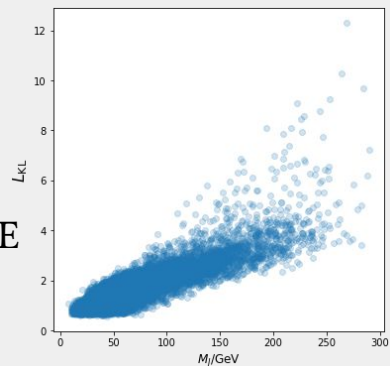
after

# Outlier Exposure -- Results -- KL-OE

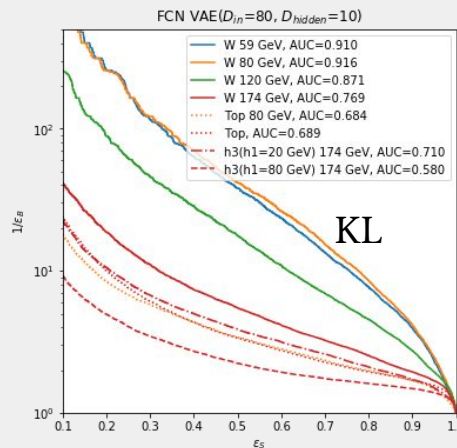
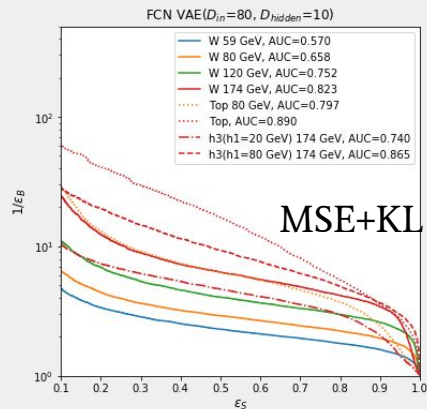
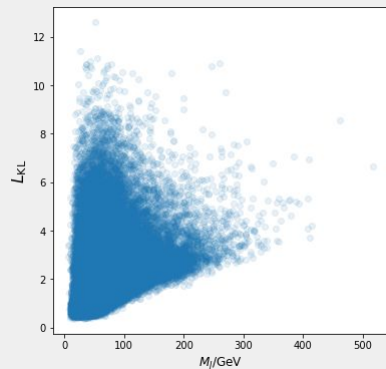
- Training using latent KL OE loss
- Reconstruction performance unchanged
- Using latent KL divergence directly as anomaly metric  $\rightarrow$  very good mass decorrelation



Without OE



With OE



\*results shown here are trained using OE unweighted samples

# Summary

- Explored Generative Model (VAE) for anomaly detection
- KL Regularization
- Anomalous Jet Tagging
  - Different jet masses
  - Different jet types
  - Anomaly metrics
- Outlier exposure to increase sensitivity to out-of-distribution samples
  - Especially in latent space
- Mass correlation affected by outlier exposure ← mass sculpting

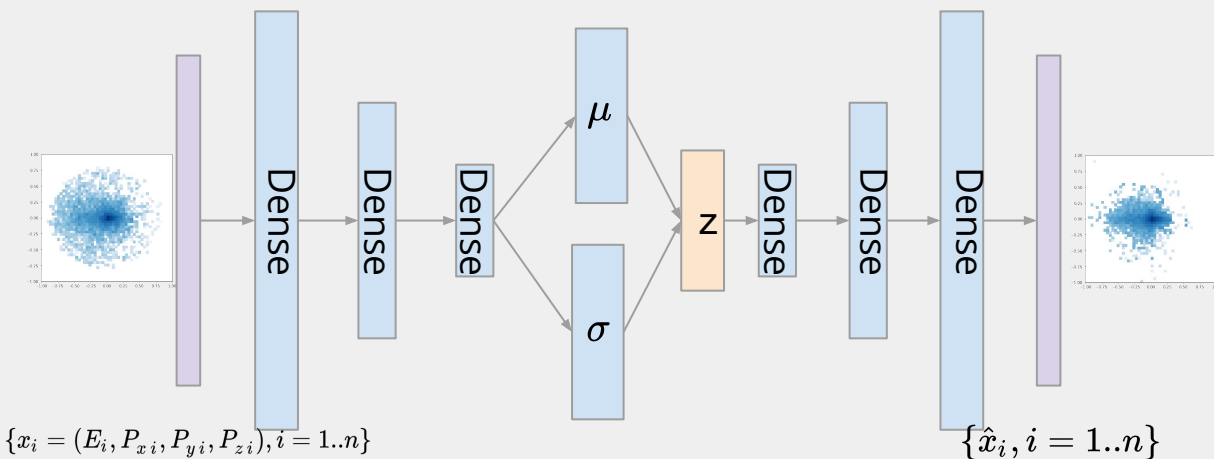
# Outlook

- Architecture and input representation
- Reconstruction loss

# Backup

# VAE Architecture -- FCN

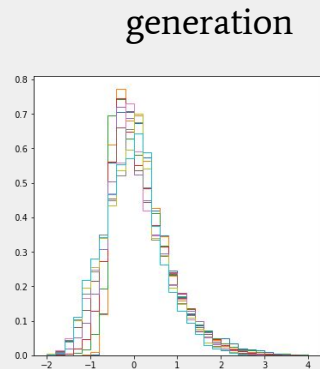
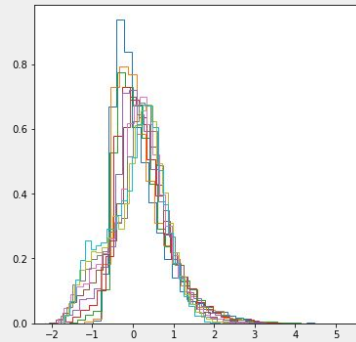
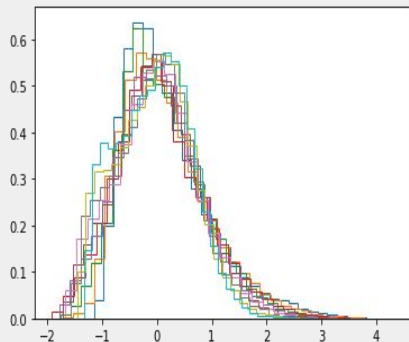
- Latent dimension = 10 (best: 6 - 20)



$$L = \frac{1}{4n} \sum_i \|\hat{x}_i - x_i\|_2^2 + \beta D_{KL}(q(z|x) \| p(z))$$

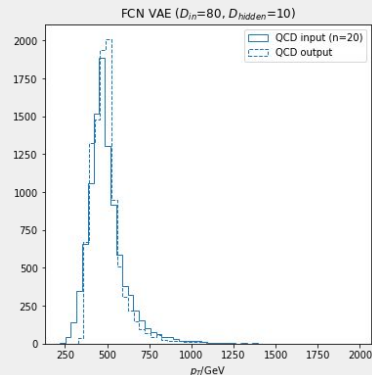
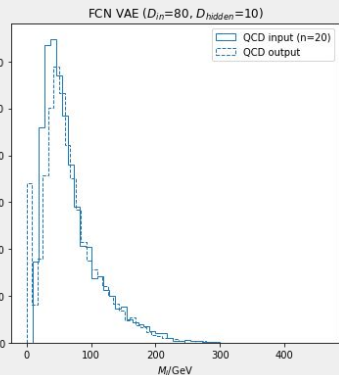
# Reconstruction Performance

- Input features

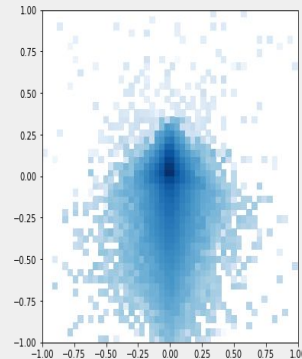


- High Level features

beta=0.1

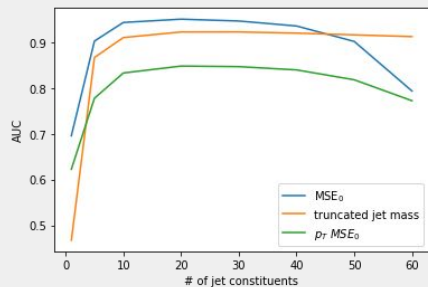


Sampling from latent priors  $\rightarrow$  decoder

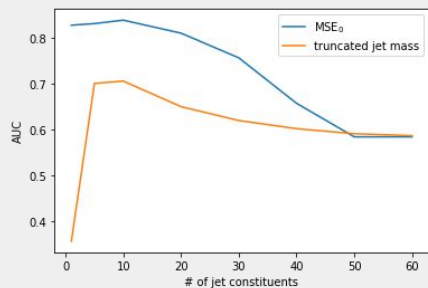


# Investigation on MSE Anomaly Metric

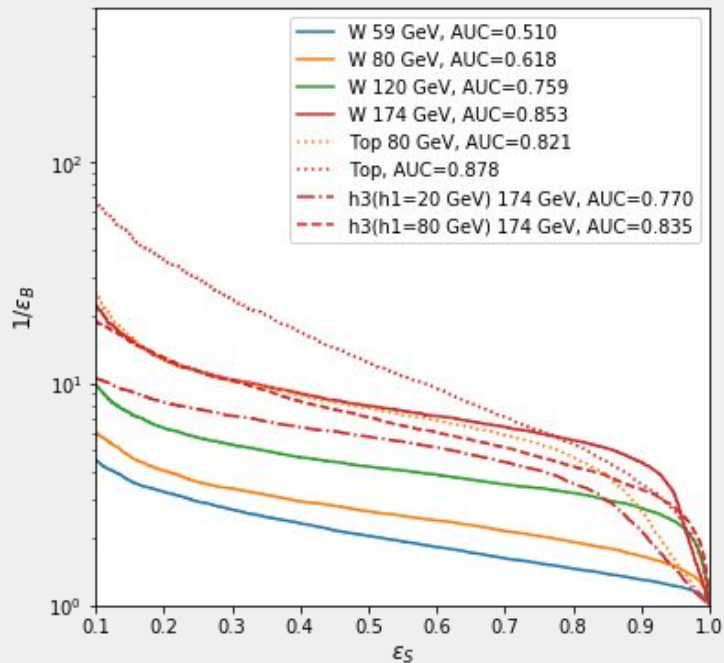
- Pure non-ML metric:  $L = \sum_i \|x\|^2 \sim \chi^2$      $M_J = \sum_i z_i \theta_i^2$
- MSE-based anomaly metric doesn't require perfect reconstruction



Top Tagging Reference

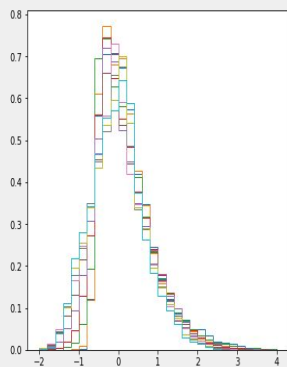


W test jets

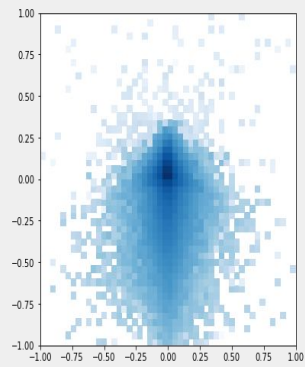


# Generator

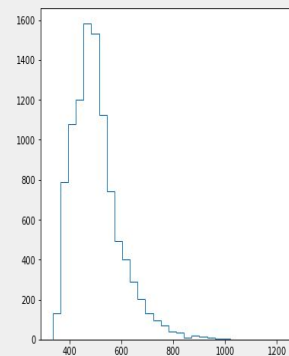
- Sample from prior latent distribution  $\mathcal{N}(0, I)$
- Specific dimensions more correlated with mass



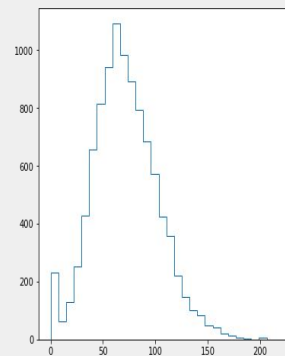
Input space



Jet Images



Jet pt



Jet mass