

---

# Improving Variational Autoencoders for New Physics Detection at the LHC with Normalizing Flows

Pratik Jawahar<sup>1,\*</sup>, Thea Aarrestad<sup>1</sup>, Nadezda Chernyavskaya<sup>1</sup>,  
Maurizio Pierini<sup>1</sup>, Kinga A. Wozniak<sup>1,2</sup>, Jennifer Ngadiuba<sup>3,4</sup>, Javier Duarte<sup>5</sup>,  
Steven Tsan<sup>5</sup>

<sup>1</sup>European Center for Nuclear Research (CERN), CH 1211, Geneva 23, Switzerland

<sup>2</sup>University of Vienna, 1010 Wien, Austria

<sup>3</sup>Fermi National Accelerator Laboratory (FNAL), Batavia, IL 60510, USA

<sup>4</sup>California Institute of Technology, Pasadena, CA 91125, USA

<sup>5</sup>University of California San Diego, La Jolla, CA 92093, USA

Correspondence\*:  
Pratik Jawahar  
pjawahar@wpi.edu

## ABSTRACT

We investigate how to improve new physics detection strategies exploiting variational autoencoders and normalizing flows for anomaly detection at the Large Hadron Collider. As a working example, we consider the DarkMachines challenge dataset. We show how different design choices (e.g., event representations, anomaly score definitions, network architectures) affect the result on specific benchmark new physics models. Once a baseline is established, we discuss how to improve the anomaly detection accuracy by exploiting normalizing flow layers in the latent space of the variational autoencoder.

## 1 INTRODUCTION

Most searches for new physics at the CERN Large Hadron Collider (LHC) target specific experimental signatures. The underlying assumption of a specific new physics model could enter at various stages in the search design, e.g., when reducing the data rate from 40 M to 1,000 collision events per second in real time (Aad et al., 2020; Sirunyan et al., 2020; Trocino, 2014), when designing the event selection, or when running the final hypothesis testing. When searching for pre-established and theoretically well-motivated particles (e.g., the Higgs boson), this strategy is extremely successful because the underlying assumption can be exploited to maximize the search sensitivity. On the other hand, the lack of a predefined target might turn this strength into a limitation.

To compensate for this potential problem, *model independent* searches are also carried out (Aaboud et al., 2019; Aaltonen et al., 2009; Aaron et al., 2009; CMS-PAS-EXO-14-016, 2017; D0 Collaboration, 2012) at hadron colliders. These searches consist in an extensive set of comparisons between the data distribution and the expectation derived from Monte Carlo simulation. Many comparisons are carried out in parallel for multiple physics-motivated features while applying different event selections. However, when searching for new physics among many channels, the “global” significance of observing a particular discrepancy must

take into account the probability of observing such a discrepancy anywhere. This so called look-elsewhere effect can be quantified in terms of a trial factor (Gross and Vitells, 2010). While the large trial factor typically reduces the statistical power of this strategy in terms of significance, model independent searches are valuable tools to identify possible regions of interest and provide data-driven motivations for traditional, more targeted searches to be performed on future data.

Recently, the use of machine learning techniques has been advocated as a mean to reduce the model dependence (Weisser and Williams, 2016; Cerri et al., 2019; D’Agnolo and Wulzer, 2019; De Simone and Jacques, 2019; Farina et al., 2020; Collins et al., 2018; Blance et al., 2019; Hajer et al., 2020; Heimel et al., 2019; Collins et al., 2019; D’Agnolo et al., 2021; Nachman and Shih, 2020; Andreassen et al., 2020; Amram and Suarez, 2021; Dillon et al., 2020; Cheng et al., 2020; Khosa and Sanz, 2020; Nachman, 2020; Park et al., 2020; Bortolato et al., 2021; Collins et al., 2021; Finke et al., 2021; Gonski et al., 2021; Hallin et al., 2021; Ostdiek, 2021). In this context, the particle-physics community engaged in two data challenges: the LHC Olympics 2020 (Kasieczka et al., 2021) and the DarkMachines challenge (Aarrestad et al., 2021), where different approaches were explored to attempt to detect an unknown signal of new physics hidden in simulated data.

As part of our contribution to the DarkMachines challenge, we investigated the use of a particle-based variational autoencoder (VAE) (Kingma and Welling, 2014; Rezende et al., 2014) and the possibility of enhancing its anomaly detection capability by using normalizing flows (NFs) (Papamakarios et al., 2021) in the latent space to optimize the choice of the latent-space prior. In this paper, we document those studies and expand that effort, investigating the impact of specific architecture choices (event representation, network architecture, usage of expert features, and definition of the anomaly score). This study is an update of our contribution to the DarkMachine challenge (Aarrestad et al., 2021), which benefits from the lessons learned by the DarkMachines challenge. Taking inspiration from solutions presented by other groups in the challenge (e.g., Refs. (Caron et al., 2021; Ostdiek, 2021)), we evaluate the impact of some of their findings on our specific setup. In some cases (but not always), these solutions translate in an improved performance, quantified using the same metrics presented in Aarrestad et al. (2021). In this way, we establish an improved baseline model, on top of which we evaluate the impact of the normalizing flow layers in the latent space.

## 2 DATA SAMPLES AND EVENT REPRESENTATION

This study is based on the datasets released on the Zenodo platform (DarkMachines Community, 2020) in relation to the Dark Machines Anomaly Score Challenge (Aarrestad et al., 2021). They consist of a set of processes predicted in the standard model (SM) of particle physics, mixed according to their production cross section in proton-proton collisions at 13 TeV center-of-mass energy, and a set of benchmark signal samples. The datasets contains labels, identifying the process that generated each event. Labels are ignored during training and used to evaluate performance metrics.

For each sample, four datasets are provided, corresponding to four different event selections (called *channels* (Aarrestad et al., 2021)):

- Channel 1:  $H_T \geq 600$  GeV,  $p_T^{\text{miss}} \geq 200$  GeV, and  $p_T^{\text{miss}}/H_T \geq 0.2$ .
- Channel 2a:  $p_T^{\text{miss}} \geq 50$  GeV and at least three light leptons (muons or electrons) with  $p_T > 15$  GeV.
- Channel 2b:  $p_T^{\text{miss}} \geq 50$  GeV,  $H_T \geq 50$  GeV and at least two light leptons (muons or electrons) with  $p_T > 15$  GeV.
- Channel 3:  $H_T \geq 600$  GeV,  $p_T^{\text{miss}} \geq 100$  GeV.

where  $p_T$  is the magnitude of a particle’s transverse momentum,  $H_T$  is the scalar sum of the jet  $p_T$  in the event, and  $\vec{p}_T^{\text{miss}}$  is the vector equal and opposite to the vector sum of the transverse momenta of the reconstructed particles in the event, while  $p_T^{\text{miss}}$  is its magnitude<sup>1</sup>. More details are provided in Aarrestad et al. (2021).

**Table 1.** Summary of the available dataset size.

Dataset	Channel 1	Channel 2a	Channel 2b	Channel 3
Training	193,800	13,425	238,450	7,100,934
Validation	10,200	707	12,550	373,733
Bkg. Test	10,000	5,868	89,000	1,025,333
Sig. Test	38,666	5,868	89,676	1,023,320

The input consists of the momenta of all the reconstructed physics objects in the event (jets, b jets, electrons  $e$ , muons  $\mu$ , and photons), ordered by decreasing  $p_T$ . Each list of objects is zero-padded to force each event into a fixed-length matrix with the same order: up to 15 jets, and up to 4 each of b jets,  $\mu^\pm$ ,  $e^\pm$ , and photons. We pre-process the input by applying the `scikit-learn` (Pedregosa et al., 2011) standard scaling and arranging the list of objects into a matrix of 39 particles times four momentum features  $(E, p_T, \eta, \phi)$ , where  $E$  is the particle energy. For  $e$ ,  $\mu$ , and photons, the energy is computed assuming zero mass. For jets, the measured jet mass is used. The input matrix is interpreted as an image or an unordered point cloud, depending on the underlying VAE architecture.

The training and validation dataset consists of background events from the SM mixture. The available dataset size is detailed in Table 1 for each of the channels. The background test samples are combined with the benchmark signal samples listed in Table 2 to form the labeled test dataset on which performance is evaluated.

### 3 TRAINING SETUP AND EVALUATION METRICS

Variational Autoencoders (Kingma and Welling, 2014; Rezende et al., 2014; Kingma and Welling, 2019) are a class of likelihood-based generative models that maximize the likelihood of the training data according to the generative model  $\prod_{x \in \text{data}} p_\theta(x)$  for the set of observed variables  $x$  in the training data. To achieve this in a tractable way, the generative model is augmented by the introduction of a set of latent variables  $z$ , such that the the marginal distribution over the observed variables  $p_\theta(x)$ , is given by:  $p_\theta(x) = \int p_\theta(x|z)q_\theta(z)dz$ . In this way,  $q_\theta(z)$  can be a relatively simple distribution, such as a Gaussian, while maintaining high expressivity for the marginal distribution  $p_\theta(x)$  as an infinite mixture of simple distributions controlled by  $z$ . Besides being used as generative models, VAEs have been shown to be effective as anomaly detection algorithms (An and Cho, 2015).

In this work, the VAE models are trained on the training and validation datasets, minimizing the loss function:

$$L_{\text{total}} = \beta D_{\text{KL}} + (1 - \beta)L_C, \tag{1}$$

<sup>1</sup> We use a Cartesian coordinate system with the  $z$  axis oriented along the beam axis, the  $x$  axis on the horizontal plane, and the  $y$  axis oriented upward. The  $x$  and  $y$  axes define the transverse plane, while the  $z$  axis identifies the longitudinal direction. The azimuth angle  $\phi$  is computed with respect to the  $x$  axis. The polar angle  $\theta$  is used to compute the pseudorapidity  $\eta = -\log(\tan(\theta/2))$ . The transverse momentum ( $p_T$ ) is the projection of the particle momentum on the  $(x, y)$  plane. We fix units such that  $c = \hbar = 1$ .

**Table 2.** BSM processes contributing to the signal dataset in each channel. The process code, adopted in this study, is taken from Aarrestad et al. (2021).

BSM process	Code	Ch.1	Ch.2a	Ch.2b	Ch.3
$Z' + \text{jet}$	monojet_Zp2000.0_DM_50.0	×	×		×
$Z' + W/Z$	monoV_Zp2000.0_DM_50.0				×
$Z' + t$	monotop_200_A	×			×
$Z'$ in LFV $U(1)_{L_\mu-L_\tau}$	pp23mt_50		×	×	
	pp24mt_50		×	×	
$\cancel{R}$ -SUSY $\tilde{t}\tilde{t}$	stlp_st1000	×		×	×
$\cancel{R}$ -SUSY $\tilde{q}\tilde{q}$	sqsq1_sq1400_neut800	×			×
SUSY $\tilde{g}\tilde{g}$	glgl1400_neutralino1100	×	×	×	×
	glgl1600_neutralino800	×	×	×	×
SUSY $\tilde{t}\tilde{t}$	stop2b1000_neutralino300	×			×
SUSY $\tilde{q}\tilde{q}$	sqsq_sq1800_neut800	×			×
SUSY $\tilde{\chi}^\pm\tilde{\chi}^0$	chaneut_cha200_neut50		×	×	
	chaneut_cha250_neut150		×	×	
SUSY $\tilde{\chi}^\pm\tilde{\chi}^\pm$	chacha_cha300_neut140			×	
	chacha_cha400_neut60			×	
	chacha_cha600_neut200			×	

where  $L_C$  is a reconstruction loss, which is chosen to be an  $L_1$ -type permutation-invariant Chamfer loss (Barrow et al., 1977):

$$L_C = \sum_{\vec{x} \in S_{\text{input}}} \min_{\vec{y} \in S_{\text{output}}} |\vec{x} - \vec{y}| + \sum_{\vec{y} \in S_{\text{output}}} \min_{\vec{x} \in S_{\text{input}}} |\vec{x} - \vec{y}|, \quad (2)$$

similar to the  $L_2$ -type Chamfer distance used in Refs. (Fan et al., 2017; Zhang et al., 2020). In Eq. (2),  $D_{\text{KL}}$  is the Kullback–Liebler divergence term usually employed to force the data distribution in the latent space to a multidimensional Gaussian with unitary covariance matrix (Rezende and Mohamed, 2015), and  $\beta$  is a parameter that controls the relative importance of the two terms (Higgins et al., 2017).

All of our models are optimized using the Adam minimizer (Kingma and Ba, 2015). A learning rate of  $10^{-4}$  is applied along with a brute force early stopping strategy used on an ad-hoc basis. A batch size of 32 is chosen to train models. All models are implemented with the PyTorch (Paszke et al., 2019) deep learning framework and are hosted on GitHub (Jawahar and Pierini, 2021).

We train and test all our models on the WPI Turing Research Cluster<sup>2</sup>, using 8 CPU nodes and 1 GPU node (NVIDIA Tesla V100 or Tesla P100).

At inference time,  $L_C$  is used as an anomaly detection score, to quantify the distance between the input and the output. By applying a lower-bound threshold on  $L_C$ , we identify every event with an  $L_C$  value larger than the threshold as an anomaly. By comparing this prediction to the ground truth, we can assess the performance of the VAE on specific signal benchmark models.

To evaluate model performance we follow the same strategy and code used in Aarrestad et al. (2021) to enable comparison with other models tested on this dataset. As explained in Aarrestad et al. (2021), we extract four main performance parameters from the receiver operating characteristic (ROC) curves based on the chosen anomaly metric for each model, namely the area under the curve (AUC) and true positive

<sup>2</sup> <https://arc.wpi.edu/computing/hpc-clusters/>

rate (also known as the signal efficiency  $\epsilon_S$ ) at three different, fixed values of the false positive rate (also known as background efficiency  $\epsilon_B$ ). We then combine these scores from all models on all available signal regions across all channels of the dataset to form box-and-whisker plots, using 6 different combination and comparison strategies namely, the highest mean score method, highest median score method, average rank method, top scorer method, top-5 scorer method, and highest minimum scorer method. A box is drawn spanning the inner half (50% quantile centered at the median) of the data as shown in Fig. 1. A line through the box marks the median. Whiskers extend from the box to either the maximum and minimum unless these are further away from the edge of the box than 1.5 box lengths. The outlier points are shown as circles.

For Fig. 1 and the other figures, the representative ranking as denoted by the legend corresponds to the performance based on the highest mean score method unless mentioned otherwise. However, to choose the best model for each experiment described in this paper, we consider all six comparison methods to arrive at a consensus. The code to perform these comparisons and to generate the corresponding plots is available in Aarrestad et al. (2021).

## 4 BASELINE VAE MODEL

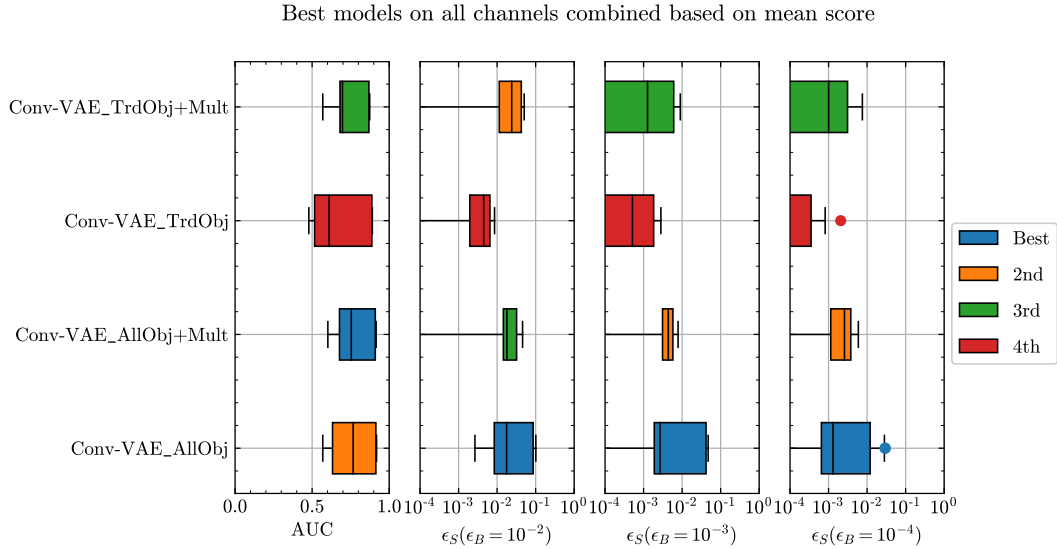
The main goal of this study is to evaluate the impact of normalizing flow layers in the latent space on the anomaly detection capability of a reference VAE model. This and the following sections describe how this reference model is built, starting from the VAE based on convolutional layers (Conv-VAE) presented in Aarrestad et al. (2021) and modifying its architecture based on some of the lessons learned during the DarkMachine challenge.

The encoder of the initial Conv-VAE consists of three convolutional layers, with 32, 16, and 8 kernels of size (3, 4), (5, 1), and (7, 1), respectively. For all layers, the stride is set to 1 and zero padding to “same”. The output of the convolutional layers is flattened and passed to 2 fully-connected neural network (FCN) layers that output the mean and variance for the latent space. The cardinality of the latent space is fixed to 15. The decoder mirrors the encoder architecture, returning an output of the same size as the input.

In order to define the reference model, the architecture of the starting model is modified in different ways, each time evaluating the impact of a given choice on the test dataset. Several possibilities are considered: how to embed the event in the two-dimensional (2D) array (see section 4.1); how to interpret the array, e.g., as an image or a graph (see section 4.2); whether to extend the event representation beyond the particle momenta, adding domain-specific high level features as an additional input (see section 4.3); and which anomaly score to use (see section 4.4). We study various options for each of these points, following this order. Doing so, we establish a candidate model, which replaces the initial model. We evaluate on this new model the benefit of using normalizing flow layers in the latent space (see section 5) to improve the anomaly detection accuracy.

### 4.1 Data representation

By their nature, events consist of a variable number of objects. To some extent, this conflicts with most neural network architectures, which assume a fixed-size input. As a baseline, we adopt the simplest solution, i.e., to zero-pad all events to standardized event sizes for all available samples. To get a better idea of how padding affects results, we study performance across alternative input encodings. We consider two main types of encodings, listed as AllObj and TrdObj in Fig. 1. The former involves considering the entire event which implies allowing for a large enough padding such that every object per event is taken into consideration across the entire dataset. The latter involves cutting down the padding and the input sequence by considering only up to four leading jets and three objects each of the other types per event.



**Figure 1.** Anomaly detection performance for the Conv-VAE with different inputs given (see text for more details): all physics objects in the event (AllObj); truncated input object list (TrdObj); all objects and array of object multiplicity (AllObj+Mult); truncated input object list and array of object multiplicity (TrdObj+Mult).

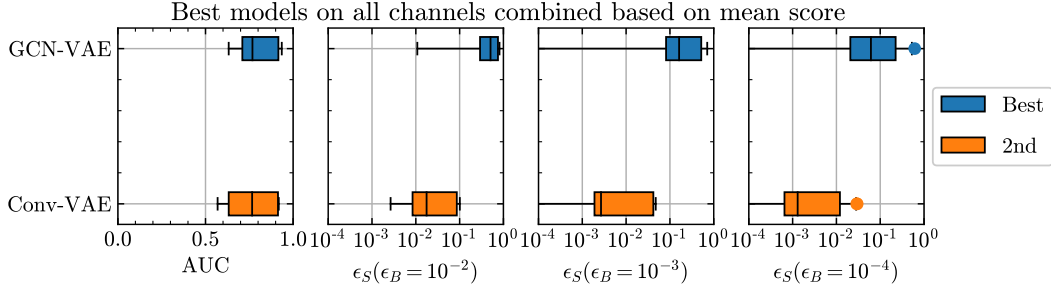
When using the truncated sequence, the model loses information regarding the number of objects of each type per event, which is implicitly learned when the whole sequence is considered. To compensate for this loss, one can explicitly add this information passing a second input to the model, consisting of a vector containing the multiplicities of each object type. This input is concatenated to the flattened output received from the convolutional layers in the encoder before passing them to the fully connected layers. For the sake of comparison, we also do the same for the AllObj case (labeled as “+Mult” in Fig. 1).

The results in Fig. 1 show that the truncated sequence does worse than the full sequence. We also see little improvement in performance with the addition of multiplicity information per event in both the AUC as well as performance at lower background efficiencies. As a result, we keep the input encoding that considers the complete sequence per event.

## 4.2 VAE architecture

The convolutional architecture used for the baseline VAE is not the only option to handle the input considered in this study. The ensemble of reconstructed particles in an event can be represented as a point cloud. Doing so, we can process it with a graph neural network. The main advantage of this choice stands with the permutation invariance of the graph processing, which pairs that of the loss in Eq. 2 and complies with the unordered nature of the input list of particles. Graph-based architectures have already been shown to perform better with sparse, non-Euclidean data representations in general (Bronstein et al., 2017; Zhou et al., 2020) and in particle physics in particular (Shlomi et al., 2020; Duarte and Vlimant, 2020).

To this end, we consider a GCN-VAE model composed of multilayer graph convolutional network layers (GCNs) (Kipf and Welling, 2017) and FCN layers in both the encoder and the decoder. As for the VAE, the input graphs are built from the input list described in section 2, each particle representing one vertex of the graph in the space identified by five particle features:  $E$ ,  $p_T$ ,  $\eta$ ,  $\phi$ , and object type. The object type is a label-encoded integer that signifies the object type. The input is structured as a fully connected, undirected



**Figure 2.** Comparison of the GCN-VAE and Conv-VAE performances, in terms of the benchmark figures of merit adopted in the paper.

graph which is passed to the GCN layers of the encoder, defined as (Kipf and Welling, 2017):

$$H_{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H_{(l)} W_{(l)}), \quad (3)$$

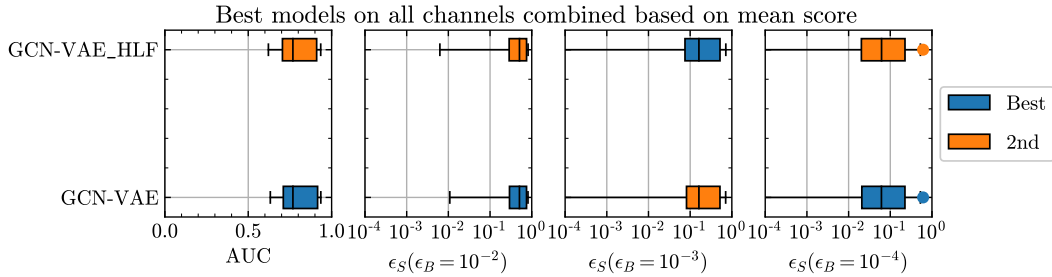
where  $H_{(l)}$  is the input to the  $(l + 1)$ th GCN layer with  $H_{(0)} = X$  where  $X$  represents the node feature matrix.  $H_{(l+1)}$  is the layer output,  $\tilde{A} = A + I$ , where  $A$  is the adjacency of the graph, with  $I$  being the identity matrix which implies added self connections for each node.  $\tilde{D}_{ii} = \sum_j A_{ij}$  is defined for the normalized adjacency based message passing regime,  $W_{(l)}$  is the layer weights matrix and  $\sigma(\bullet)$  is a suitable nonlinear activation function. The output of the last GCN layer is flattened and passed to an FCN layer which populates the latent space. The encoder has 3 GCN layers that scale the 5 node features to 32, 16, and 2 respectively, followed by a single FCN layer which generates a 15-dimensional latent space. The decoder has a symmetrically inverted structure with the sampled point being upscaled through an FCN layer first and the resulting output is reshaped and passed to GCN layers that reconstruct the node features.

Considering all comparison metrics along with the representative results shown in Fig. 2, graph architectures exhibit a definitive improvement in performance compared to the Conv-VAE. The improvement is seen not only in the AUC metric, but more significantly in the  $\epsilon_S$  at low  $\epsilon_B$ . Because of this, the GCN-VAE is used as the reference architecture in the rest of this section and in section 5.

### 4.3 Physics-motivated high-level features

We also experiment with adding physics-motivated high-level features, as explicit inputs to the model, similar to what was done with object multiplicities in section 4.1. Doing so, we intend to check if domain knowledge helps in improving anomaly detection capability. We pass event information such as the missing transverse momentum in the event ( $p_T^{\text{miss}}$ ), the scalar sum of the jet  $p_T$  ( $H_T$ ) and  $m_{\text{Eff}} = H_T + p_T^{\text{miss}}$  to the model, by concatenating these with the output of the convolutional layers of the encoder. The concatenated output is then passed to the fully connected layers in the encoder to form the latent space. After the point sampled from the latent space passes through the fully connected layers of the decoder, the reconstructed  $p_T^{\text{miss}}$ ,  $H_T$  and  $m_{\text{Eff}}$  are extracted and the rest of the layer output is re-shaped and further passed to the subsequent layers of the decoder.

To include the reconstruction of these features in the loss, we add to Eq. (1) a mean-squared error (MSE) term, computed from the reconstructed and input high-level features and weighted by a coefficient. This coefficient is treated as a hyperparameter that is scanned until the best performance is found.



**Figure 3.** Comparison of the GCN-VAE performance with and without high-level features added as a separate input.

Figure 3 shows that adding high-level features brings no definitive improvement in performance, thereby leading us to conclude that the baseline model with marginally lower number of trainable parameters is a good choice.

#### 4.4 Anomaly scores

While so far the Chamfer loss has been used as the anomaly score, this is not the only possibility. We consider two alternative metrics: the  $D_{KL}$  term in Eq. (1) and (Aarrestad et al., 2021):

$$R_z = \sum_i \left( \frac{\mu_i}{\sigma_i} \right)^2 \tag{4}$$

where  $\mu$  and  $\sigma$  are the mean and RMS returned by the encoder and the index  $i$  runs across the latent-space dimensions.

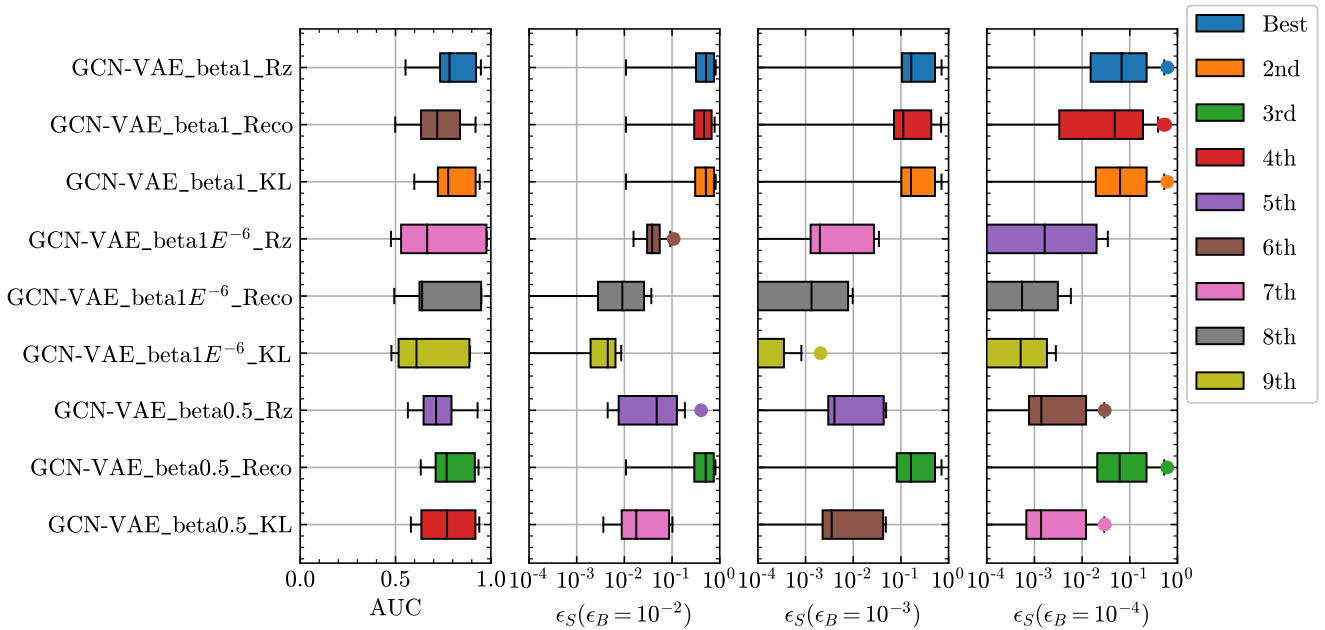
The use of different anomaly scores requires a tuning of the  $\beta$  hyperparameter. Since  $\beta$  determines the relative importance of the  $D_{KL}$  and Chamfer loss terms in the loss, the use of one or the other as anomaly score is certainly related to the choice of the optimal  $\beta$  value. Similarly, the use of  $R_z$  (i.e., anomaly detection in the latent space) might not be optimal when using a  $\beta$  value that was tuned to emphasize the reconstruction accuracy (i.e., the minimization of the Chamfer term in the loss). On the other hand, the study in Aarrestad et al. (2021) shows that an excessive tuning of the hyperparameters affects generalization of performance negatively beyond the available dataset.

In order to address this point, we compare three weights for the  $\beta$  term. The first case ( $\beta = 1$ ) corresponds to training the VAE without the contribution of the reconstruction loss. In the second case ( $\beta = 0.5$ ) the two contributions are equally weighted. The final case ( $\beta = 10^{-6}$ ) corresponds to suppressing the  $D_{KL}$  term to a negligible level.

Figure 4 shows that all three anomaly scores underperform in the  $\beta = 10^{-6}$  case. The best performing models overall are the  $\beta = 1$  and  $\beta = 0.5$  cases. Comparing across the three different anomaly scores, we see that the  $\beta = 1$  model that uses  $D_{KL}$  and  $R_z$  metrics, as well as the  $\beta = 0.5$  model that uses the reconstruction metric perform the best. All three cases also show very similar performance across all comparison metrics as well as methods, implying that either model-anomaly score combination is equally suitable. We also find that the  $\beta = 1$   $D_{KL}$  score and the  $\beta = 0.5$  reconstruction score show a similar correlation pattern on signal and background. As a result, we expect that only a limited improvement would be obtained by combining the two, which spares us the cost of introducing a new hyperparameter (the



Best models on all channels combined based on mean score



**Figure 4.** Comparison of anomaly detection performance from different anomaly score definitions, applied to the GCN-VAE.

relative weight of the two terms) whose optimal value would be signal-specific, as in the case of Caron et al. (2021).

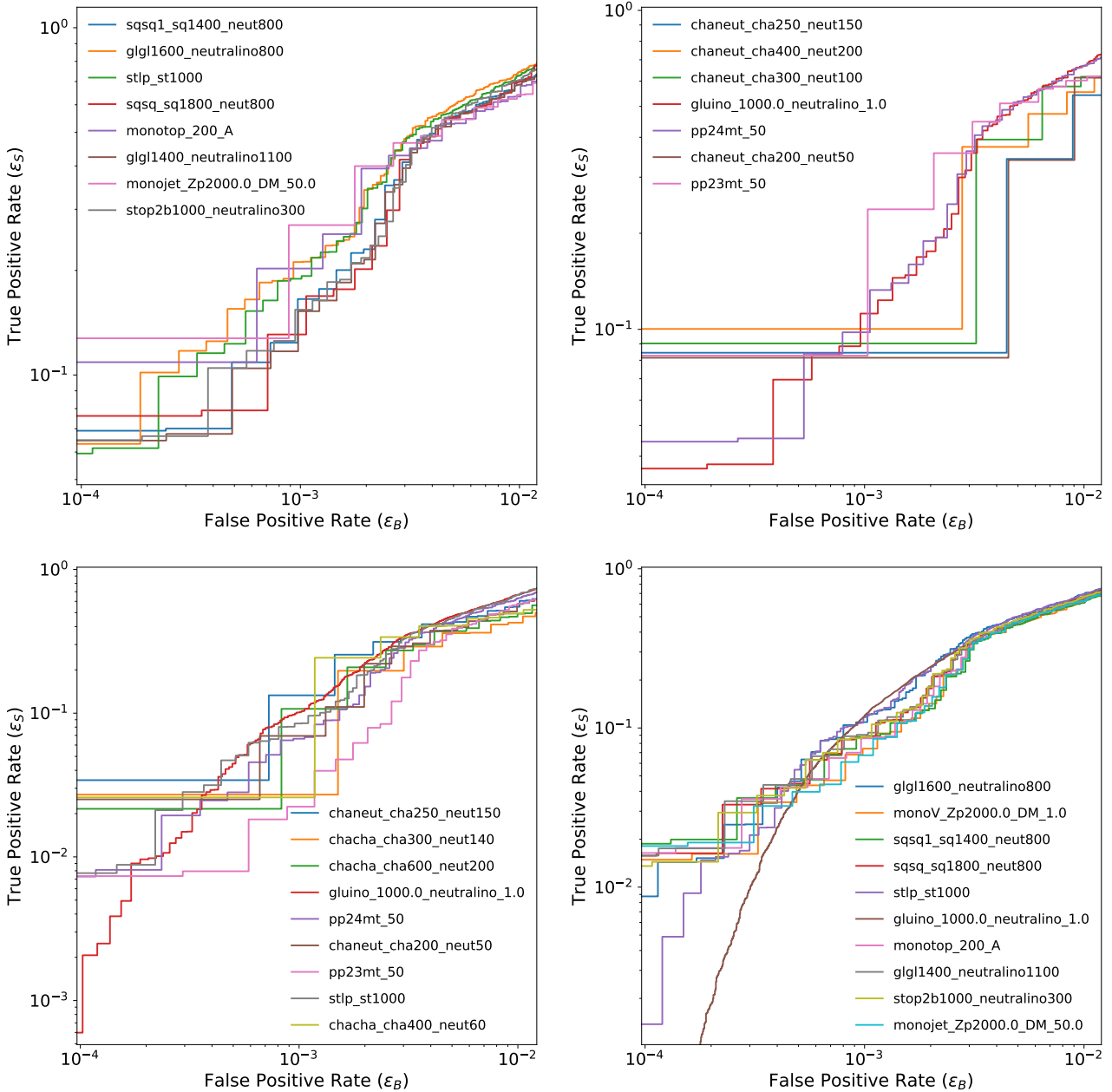
### 4.5 Baseline discrimination

As a result of the tests presented so far, the baseline VAE model is established as a GCN-VAE taking as input the whole set of reconstructed physics object but no domain-specific high level features. The Chamfer loss function is used as the anomaly score. The GCN-VAE is trained and tested only with data available within a given channel and the dataset sizes per channel are described in Table 1. Figure 5 shows the ROC curves for the baseline VAE model on benchmark signals in the four channels. It is evident that we suffer from a shortage of events for some signal models at very low  $\epsilon_B$ . We still show ROC curves down to  $\epsilon_B = 10^{-4}$  to allow one to compare our results to those in Arrestad et al. (2021), where this range was chosen. We see an overall improvement in  $\epsilon_S$  at very low  $\epsilon_B$  for the GCN-VAE compared to our Conv-VAE submission in Arrestad et al. (2021).

## 5 NORMALIZING FLOWS

With the GCN-VAE serving as the baseline, we investigate how the use of NFs (Tabak and Vanden-Eijnden, 2010; Tabak and Turner, 2013) impacts the anomaly-detection performance. Normalizing flow layers are inserted between the Gaussian sampling and the decoder. They provide additional complexity to learn better posterior distributions (Rezende and Mohamed, 2015) by morphing the multivariate prior of the latent space to a more suitable, learned function.

In other words, we use the NF layers to handle the fact that a VAE converging to a good output-to-input matching does not necessarily correspond to a configuration with a Gaussian prior in the latent space,  $p(z) = \prod G(z)$ . To reach this configuration (e.g., when training a VAE as a generative model), one typically uses a  $\beta$ -VAE with an increased weighting of the  $D_{KL}$  regularizer. This typically results in a



**Figure 5.** ROC curves for the baseline GCN-VAE model in channel 1 (top left), channel 2a (top right), channel 2b (bottom left), and channel 3 (bottom right), computed from the  $\epsilon_S$  and  $\epsilon_B$  values obtained on the background sample and the benchmark signal samples. Most of the ROC curves are not smooth, due to the small dataset size for some of the channels.

degradation of the output-to-input matching. With NFs, we learn a generic prior  $p(z)$  as  $f(G(z))$ , where  $f$  is the transformation function learned by the NF layers. This is different from the way NFs are traditionally used in VAE training, i.e., to improve the convergence of  $f(z)$  to  $G(z)$  with a stronger evidence lower bound (ELBO) condition. Because of this, we do not modify the  $D_{KL}$  term in the loss, as done in Rezende and Mohamed (2015). The results obtained following this more traditional training procedure are described in the supplementary material. Doing so, we observe worse  $\epsilon_S$  for the same  $\epsilon_B$ . This is expected because the ELBO improvement with NFs was introduced in Tomczak and Welling (2017) as a way to improve the VAE generative properties, and it does not imply a better anomaly detection capability.

A NF can be generalized as any invertible, diffeomorphic transformation that can be applied to a given distribution to produce tractable distributions (Papamakarios et al., 2021; Kobyzev et al., 2020). In order to be compatible with variational inference, it is desirable for the transformations to have an efficient mechanism for computing the determinant of the Jacobian, while being invertible (Rezende and Mohamed, 2015). The NFs are trained sequentially, together with the baseline VAE model.

We utilize four major families of flow models:

- **Planar flows** are invertible transformations whose Jacobian determinant can be computed rather efficiently, making them suitable candidates for variational inference (Rezende and Mohamed, 2015). PF transformations are defined as:

$$\mathbf{z}' = \mathbf{z} + \mathbf{u}h(\mathbf{w}^T \mathbf{z} + b) \quad , \quad (5)$$

where  $\mathbf{u}, \mathbf{w} \in \mathbb{R}^D, b \in \mathbb{R}$  and  $h$  is a suitable smooth activation function.

- **Sylvester normalizing flows** (SNFs) (Berg et al., 2018) build on the planar flow formulation and extend it to be analogous to a multilayer perceptron with one hidden layer of  $M$  units and a residual connection as:

$$\mathbf{z}' = \mathbf{z} + \mathbf{A}h(\mathbf{B}\mathbf{z} + b) \quad , \quad (6)$$

where  $\mathbf{A} \in \mathbb{R}^{D \times M}, \mathbf{B} \in \mathbb{R}^{M \times D}, b \in \mathbb{R}^M$  and  $M \leq D$ . Computing the Jacobian determinant for such a formulation is made more efficient by utilizing the Sylvester determinant identity (Berg et al., 2018). Depending on the way  $A$  and  $B$  are parametrized, we get different types of SNFs. In this paper we use orthogonal, Householder, and triangular SNFs, as described in Berg et al. (2018).

- **Inverse autoregressive flows** (IAFs) (Kingma et al., 2016) are computation-efficient normalizing flows based on autoregressive models. Autoregressive transformations are invertible, making them suitable candidates for our case. However, computing the transformation requires multiple sequential steps (Berg et al., 2018). The inverse transformation however, leads to certain simplifications as described in Berg et al. (2018), allowing more efficient parallel computing, thereby making it a more desirable transformation for our case. We use the IAFs formulated as:

$$z_i^t = \mu_i^t(z_{1:i-1}^{t-1}) + \sigma_i^t(z_{1:i-1}^{t-1}) \cdot z_i^{t-1} \quad , \quad i = 1, 2, \dots, D \quad . \quad (7)$$

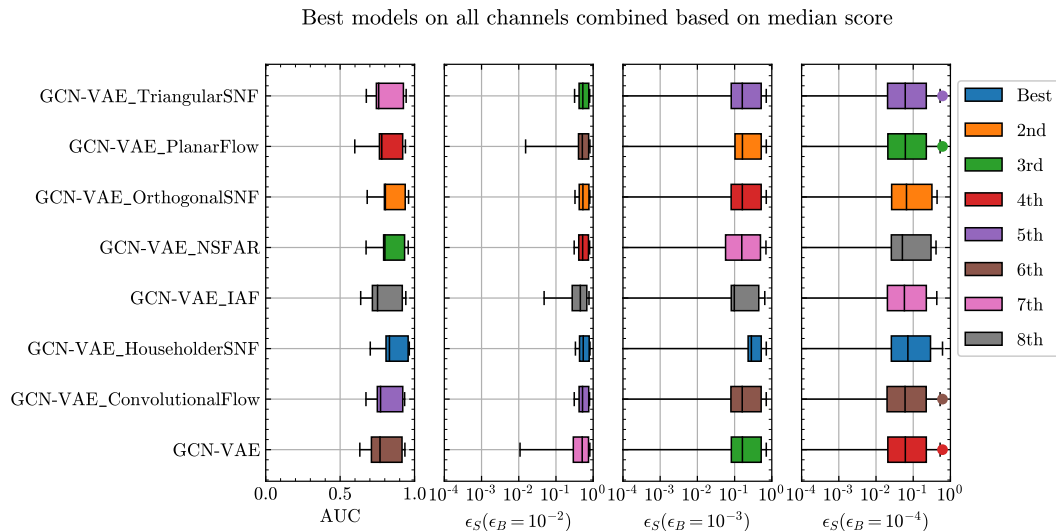
Such a formulation allows one to stack multiple transformations to achieve more flexibility in producing target distributions.

- **Convolutional normalizing flows** (ConvolutionalFlows) (Zheng et al., 2018) are an extension of single-hidden-unit planar flows (Kingma et al., 2016) to the case of multiple hidden units, further enhanced by replacing the fully connected network operation with a one-dimensional (1D) convolution, to achieve bijectivity. They are defined by the following transformation:

$$\mathbf{z}' = \mathbf{z} + \mathbf{u} \odot h(\text{conv}(\mathbf{z}, \mathbf{w})) \quad , \quad (8)$$

where  $w \in \mathbb{R}^k$  is the parameter of the 1D convolution filter with  $k$ -sized kernel,  $h$  is a monotonic nonlinear activation function and  $\odot$  denotes pointwise multiplication.

- **Autoregressive neural spline flows** (NSFARs) (Durkan et al., 2019) are similar to IAFs, where affine transforms are replaced by monotonic rational-quadratic spline transforms as described in Durkan et al.



**Figure 6.** Comparison of anomaly detection performance for GCN-VAE models with different normalizing flow architectures in the latent space

(2019). They resemble a traditional feed-forward neural network architecture, alternating between linear transformations and elementwise non-linearities, while retaining an exact, analytic inverse.

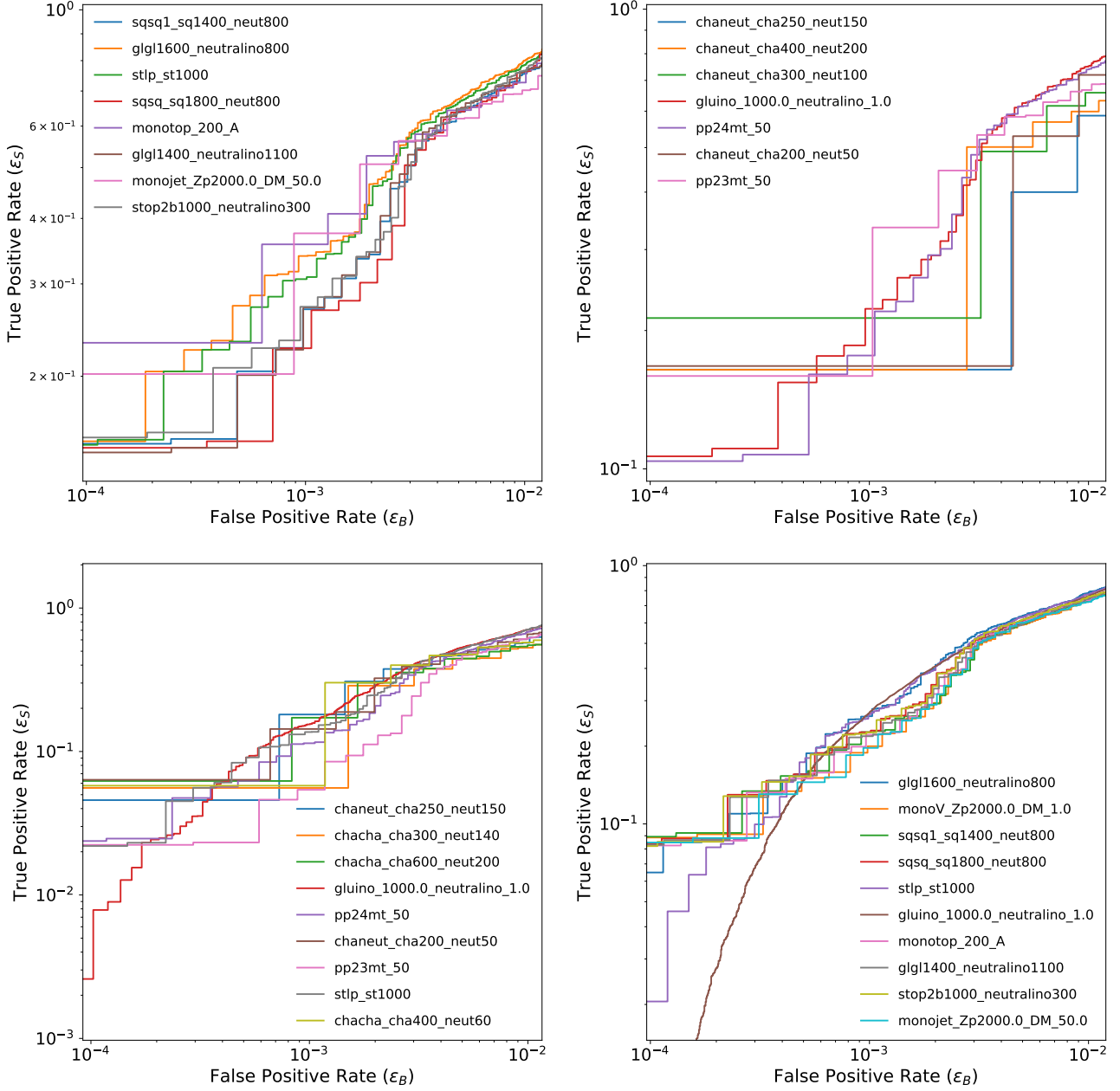
The hyperparameters for each normalizing flow architecture are chosen arbitrarily to avoid overtuning on the available dataset as learned from Aarrestad et al. (2021). The planar flow model consists of a stack of six flows, each made of three dense layers with 90 neurons each. SNFs are defined by stacking six flows with 8 orthogonal, householder and triangular transformations for each of the respective types of SNF. IAFs are constructed with four masked autoencoder for distribution estimation (MADE) (Germain et al., 2015) layers as described in Kingma et al. (2016), each containing 330 neurons. ConvolutionalFlows include four flow layers with kernel size  $k = 7$  and applying kernel dilation as described in Zheng et al. (2018). NSFARs are defined by stacking four flow layers each with  $K = 64$  bins and eight hidden features.

Figure 6 shows the results of all GCN-VAE models combined with all the different types of flows as described in section 5. Based on results from all data channels combined through all six strategies mentioned in section 3, and considering variance across trainings from different random seeds (see supplementary material), it is evident that using normalizing flows improves not only the AUC metric but also the signal efficiencies at low background efficiencies. We find that the Householder variant of SNFs produces the best improvement with respect to the baseline GCN-VAE model. The exercise was also repeated with a Conv-VAE model and similar trends were observed. There, the normalizing flows showed a larger improvement from the baseline Conv-VAE than for the GCN-VAE model but the overall results are less accurate than that of GCN-VAE with normalizing flows.

Figure 7 shows the ROC curves for the best presented model, GCN-VAE\_HouseholderSNF across all available signal samples in all data channels. For some of the samples, the small dataset size translates in a discontinuous curve and larger uncertainties.

## 6 CONCLUSIONS

We constructed a graph-based anomaly detection model to identify new physics events in the DarkMachines challenge dataset. Inspired by the outcome of this challenge, specific model design choices (data representation, use of physics-motivated high-level features, and anomaly score definition) were further



**Figure 7.** ROC curves of GCN-VAE\_HouseholderSNF for all signals in each of channel 1 (top left), channel 2a (top right), channel 2b (bottom left), and channel 3 (bottom right).

optimized in order to maximize anomaly detection performance. As the case for many other deep learning applications to particle-physics data, we observed that the graph architecture better captures the point-cloud nature of this data, resulting in an enhanced performance.

In this baseline, we investigate the impact of using a stack of normalizing flows in the latent space of the variational autoencoder (VAE), between the Gaussian sampling and the decoding, in order to improve the accuracy of the prior learning process, by morphing the Gaussian prior to a more suitable function, learned during the training.

Testing the trained model on a set of benchmark signal samples, we observe an overall improvement when normalizing flows are used, with the Householder variant of the Sylvester normalizing flow model giving the best results. With that, we reach a median anomaly identification probability of 72% (34%) for an  $\epsilon_B$  of 1% (0.1%) across all signal samples over all available channels. The median anomaly identification probability increases to 95% (96%) for an  $\epsilon_B$  of 30% (60%).

This work presents an improvement over our Conv-VAE model, submitted to the DarkMachines challenge (Arrestad et al., 2021).

## FUNDING

P. J., T. A., M. P., and K. A. W. are supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation program (Grant Agreement No. 772369). J. D. is supported by the U.S. Department of Energy (DOE), Office of Science, Office of High Energy Physics Early Career Research program under Award No. DE-SC0021187. S. T. was supported by the University of California San Diego Triton Research and Experiential Learning Scholars (TRELS) program. J. N. is supported by Fermi Research Alliance, LLC under Contract No. DE-AC02-07CH11359 with the U.S. Department of Energy, Office of Science, Office of High Energy Physics.

## REFERENCES

- Aaboud, M. et al. (2019). A strategy for a general search for new phenomena using data-derived signal regions and its application within the ATLAS experiment. *Eur. Phys. J. C* 79, 120. doi:10.1140/epjc/s10052-019-6540-y
- Aad, G. et al. (2020). Operation of the ATLAS trigger system in Run 2. *JINST* 15, P10004. doi:10.1088/1748-0221/15/10/P10004
- Aaltonen, T. et al. (2009). Global Search for New Physics with 2.0 fb<sup>-1</sup> at CDF. *Phys. Rev. D* 79, 011101. doi:10.1103/PhysRevD.79.011101
- Aaron, F. D. et al. (2009). A General Search for New Phenomena at HERA. *Phys. Lett. B* 674, 257–268. doi:10.1016/j.physletb.2009.03.034
- Arrestad, T. et al. (2021). The Dark Machines Anomaly Score Challenge: Benchmark Data and Model Independent Event Classification for the Large Hadron Collider. arXiv:2105.14027. Submitted to *SciPost Phys*.
- Amram, O. and Suarez, C. M. (2021). Tag N’ Train: a technique to train improved classifiers on unlabeled data. *JHEP* 01, 153. doi:10.1007/JHEP01(2021)153
- An, J. and Cho, S. (2015). Variational autoencoder based anomaly detection using reconstruction probability. *Special Lecture on IE* 2, 1
- Andreassen, A., Nachman, B., and Shih, D. (2020). Simulation assisted likelihood-free anomaly detection. *Physical Review D* 101. doi:10.1103/physrevd.101.095004
- Barrow, H. G., Tenenbaum, J. M., Bolles, R. C., and Wolf, H. C. (1977). Parametric correspondence and Chamfer matching: Two new techniques for image matching. In *Proceedings of the 5th International Joint Conference on Artificial Intelligence (IJCAI)* (San Francisco, CA, USA: Morgan Kaufmann Publishers Inc.), vol. 2, 659
- Berg, R. v. d., Hasenclever, L., Tomczak, J. M., and Welling, M. (2018). Sylvester normalizing flows for variational inference. In *Conference on Uncertainty in Artificial Intelligence (UAI) 2018*
- Blance, A., Spannowsky, M., and Waite, P. (2019). Adversarially-trained autoencoders for robust unsupervised new physics searches. *JHEP* 10, 047. doi:10.1007/JHEP10(2019)047
- Bortolato, B., Dillon, B. M., Kamenik, J. F., and Smolkovič, A. (2021). Bump Hunting in Latent Space

- Bronstein, M. M., Bruna, J., LeCun, Y., Szlam, A., and Vandergheynst, P. (2017). Geometric deep learning: Going beyond Euclidean data. *IEEE Signal Processing Magazine* 34, 18. doi:[10.1109/MSP.2017.2693418](https://doi.org/10.1109/MSP.2017.2693418)
- Caron, S., Hendriks, L., and Verheyen, R. (2021). Rare and Different: Anomaly Scores from a combination of likelihood and out-of-distribution models to detect new physics at the LHC
- Cerri, O., Nguyen, T. Q., Pierini, M., Spiropulu, M., and Vlimant, J.-R. (2019). Variational Autoencoders for New Physics Mining at the Large Hadron Collider. *JHEP* 05, 036. doi:[10.1007/JHEP05\(2019\)036](https://doi.org/10.1007/JHEP05(2019)036)
- Cheng, T., Arguin, J.-F., Leissner-Martin, J., Pilette, J., and Golling, T. (2020). Variational Autoencoders for Anomalous Jet Tagging
- [Dataset] CMS-PAS-EXO-14-016 (2017). MUSiC, a Model Unspecific Search for New Physics, in pp Collisions at  $\sqrt{s} = 8$  TeV. [CMS-PAS-EXO-14-016](https://doi.org/10.1007/JHEP05(2019)036)
- Collins, J. H., Howe, K., and Nachman, B. (2018). Anomaly Detection for Resonant New Physics with Machine Learning. *Phys. Rev. Lett.* 121, 241803. doi:[10.1103/PhysRevLett.121.241803](https://doi.org/10.1103/PhysRevLett.121.241803)
- Collins, J. H., Howe, K., and Nachman, B. (2019). Extending the search for new resonances with machine learning. *Phys. Rev. D* 99, 014038. doi:[10.1103/PhysRevD.99.014038](https://doi.org/10.1103/PhysRevD.99.014038)
- Collins, J. H., Martín-Ramiro, P., Nachman, B., and Shih, D. (2021). Comparing weak- and unsupervised methods for resonant anomaly detection. *Eur. Phys. J. C* 81, 617. doi:[10.1140/epjc/s10052-021-09389-x](https://doi.org/10.1140/epjc/s10052-021-09389-x)
- D0 Collaboration (2012). Model independent search for new phenomena in  $p\bar{p}$  collisions at  $\sqrt{s} = 1.96$  TeV. *Phys. Rev. D* 85, 092015. doi:[10.1103/PhysRevD.85.092015](https://doi.org/10.1103/PhysRevD.85.092015)
- D’Agnolo, R. T., Grosso, G., Pierini, M., Wulzer, A., and Zanetti, M. (2021). Learning multivariate new physics. *Eur. Phys. J. C* 81, 89. doi:[10.1140/epjc/s10052-021-08853-y](https://doi.org/10.1140/epjc/s10052-021-08853-y)
- D’Agnolo, R. T. and Wulzer, A. (2019). Learning New Physics from a Machine. *Phys. Rev. D* 99, 015014. doi:[10.1103/PhysRevD.99.015014](https://doi.org/10.1103/PhysRevD.99.015014)
- [Dataset] DarkMachines Community (2020). Unsupervised-hackathon. doi:[10.5281/zenodo.3961917](https://doi.org/10.5281/zenodo.3961917)
- De Simone, A. and Jacques, T. (2019). Guiding New Physics Searches with Unsupervised Learning. *Eur. Phys. J. C* 79, 289. doi:[10.1140/epjc/s10052-019-6787-3](https://doi.org/10.1140/epjc/s10052-019-6787-3)
- Dillon, B. M., Faroughy, D. A., Kamenik, J. F., and Szewc, M. (2020). Learning the latent structure of collider events. *JHEP* 10, 206. doi:[10.1007/JHEP10\(2020\)206](https://doi.org/10.1007/JHEP10(2020)206)
- Duarte, J. and Vlimant, J.-R. (2020). Graph neural networks for particle tracking and reconstruction. In *Artificial Intelligence for High Energy Physics*, eds. P. Calafiura, D. Rousseau, and K. Terao (World Scientific Publishing). doi:[10.1142/12200](https://doi.org/10.1142/12200). Submitted to *Int. J. Mod. Phys. A*
- Durkan, C., Bekasov, A., Murray, I., and Papamakarios, G. (2019). Neural spline flows. *Advances in Neural Information Processing Systems* 32, 7511–7522
- Fan, H., Su, H., and Guibas, L. J. (2017). A point set generation network for 3D object reconstruction from a single image. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2463. doi:[10.1109/CVPR.2017.264](https://doi.org/10.1109/CVPR.2017.264)
- Farina, M., Nakai, Y., and Shih, D. (2020). Searching for New Physics with Deep Autoencoders. *Phys. Rev. D* 101, 075021. doi:[10.1103/PhysRevD.101.075021](https://doi.org/10.1103/PhysRevD.101.075021)
- Finke, T., Krämer, M., Morandini, A., Mück, A., and Oleksiyuk, I. (2021). Autoencoders for unsupervised anomaly detection in high energy physics. *JHEP* 06, 161. doi:[10.1007/JHEP06\(2021\)161](https://doi.org/10.1007/JHEP06(2021)161)
- Germain, M., Gregor, K., Murray, I., and Larochelle, H. (2015). MADE: Masked autoencoder for distribution estimation. In *Proceedings of the 32nd International Conference on Machine Learning*, eds. F. Bach and D. Blei (Lille, France: PMLR), vol. 37 of *Proceedings of Machine Learning Research*, 881
- Gonski, J., Lai, J., Nachman, B., and Ochoa, I. (2021). High-dimensional Anomaly Detection with Radiative Return in  $e^+e^-$  Collisions

- Gross, E. and Vitells, O. (2010). Trial factors for the look elsewhere effect in high energy physics. *Eur. Phys. J. C* 70, 525. doi:10.1140/epjc/s10052-010-1470-8
- Hajer, J., Li, Y.-Y., Liu, T., and Wang, H. (2020). Novelty Detection Meets Collider Physics. *Phys. Rev. D* 101, 076015. doi:10.1103/PhysRevD.101.076015
- Hallin, A., Isaacson, J., Kasieczka, G., Krause, C., Nachman, B., Quadfasel, T., et al. (2021). Classifying Anomalies THrough Outer Density Estimation (CATHODE)
- Heimel, T., Kasieczka, G., Plehn, T., and Thompson, J. M. (2019). QCD or What? *SciPost Phys.* 6, 030. doi:10.21468/SciPostPhys.6.3.030
- Higgins, I., Matthey, L., Pal, A., Burgess, C. P., Glorot, X., Botvinick, M., et al. (2017). beta-VAE: Learning basic visual concepts with a constrained variational framework. In *5th International Conference on Learning Representations*
- Jawahar, P. and Pierini, M. (2021). mpp-hep/DarkFlow repository. <https://github.com/mpp-hep/DarkFlow>
- Kasieczka, G. et al. (2021). The LHC Olympics 2020: A Community Challenge for Anomaly Detection in High Energy Physics. *Rep. Prog. Phys.* 84, 124201. doi:10.1088/1361-6633/ac36b9
- Khosa, C. K. and Sanz, V. (2020). Anomaly Awareness. [arXiv:2007.14462](https://arxiv.org/abs/2007.14462)
- Kingma, D. P. and Ba, J. (2015). Adam: A method for stochastic optimization. In *3rd International Conference for Learning Representations*
- Kingma, D. P., Salimans, T., Jozefowicz, R., Chen, X., Sutskever, I., and Welling, M. (2016). Improving variational inference with inverse autoregressive flow. In *Advances in Neural Information Processing Systems*, eds. D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett (Curran Associates, Inc.), vol. 29
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational Bayes. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*
- Kingma, D. P. and Welling, M. (2019). An introduction to variational autoencoders. *Found. Trends Mach. Learn.* 12, 307. doi:10.1561/22000000056
- Kipf, T. N. and Welling, M. (2017). Semi-supervised classification with graph convolutional networks. In *5th International Conference on Learning Representations*
- Kobyzev, I., Prince, S., and Brubaker, M. (2020). Normalizing flows: An introduction and review of current methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*
- Nachman, B. (2020). Anomaly Detection for Physics Analysis and Less than Supervised Learning. [arXiv:2010.14554](https://arxiv.org/abs/2010.14554)
- Nachman, B. and Shih, D. (2020). Anomaly Detection with Density Estimation. *Phys. Rev. D* 101, 075042. doi:10.1103/PhysRevD.101.075042
- Ostdiek, B. (2021). Deep Set Auto Encoders for Anomaly Detection in Particle Physics
- Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., and Lakshminarayanan, B. (2021). Normalizing flows for probabilistic modeling and inference. *J. Mach. Learn. Res.* 22, 1
- Park, S. E., Rankin, D., Udrescu, S.-M., Yunus, M., and Harris, P. (2020). Quasi Anomalous Knowledge: Searching for new physics with embedded knowledge. *JHEP* 21, 030. doi:10.1007/JHEP06(2021)030
- Paszke, A. et al. (2019). PyTorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems*, eds. H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Curran Associates, Inc.), vol. 32
- Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., et al. (2011). Scikit-learn: Machine learning in Python. *J. Mach. Learn. Res.* 12, 2825



- Rezende, D. and Mohamed, S. (2015). Variational inference with normalizing flows. In *Proceedings of the 32nd International Conference on Machine Learning*, eds. F. Bach and D. Blei (Lille, France: PMLR), vol. 37, 1530
- Rezende, D. J., Mohamed, S., and Wierstra, D. (2014). Stochastic backpropagation and approximate inference in deep generative models. In *Proceedings of the 31st International Conference on Machine Learning*. vol. 32 of *Proceedings of Machine Learning Research*, 1278
- Shlomi, J., Battaglia, P., and Vlimant, J.-R. (2020). Graph neural networks in particle physics. *Mach. Learn.: Sci. Tech.* 2, 021001. doi:10.1088/2632-2153/abbf9a
- Sirunyan, A. M. et al. (2020). Performance of the CMS Level-1 trigger in proton-proton collisions at  $\sqrt{s} = 13$  TeV. *JINST* 15, P10017. doi:10.1088/1748-0221/15/10/P10017
- Tabak, E. G. and Turner, C. V. (2013). A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics* 66, 145
- Tabak, E. G. and Vanden-Eijnden, E. (2010). Density estimation by dual ascent of the log-likelihood. *Communications in Mathematical Sciences* 8, 217
- Tomczak, J. M. and Welling, M. (2017). Improving variational auto-encoders using convex combination linear inverse autoregressive flow. In *Benelearn 2017*
- Trocino, D. (2014). The CMS High Level Trigger. *J. Phys. Conf. Ser.* 513, 012036. doi:10.1088/1742-6596/513/1/012036
- Weisser, C. and Williams, M. (2016). Machine learning and multivariate goodness of fit
- Zhang, Y., Hare, J., and Prügel-Bennett, A. (2020). FSPool: Learning set representations with featurewise sort pooling. In *8th International Conference on Learning Representations*
- Zheng, G., Yang, Y., and Carbonell, J. (2018). Convolutional normalizing flows. In *ICML 2018 Workshop on Theoretical Foundations and Applications of Deep Generative Models*
- Zhou, J., Cui, G., Hu, S., Zhang, Z., Yang, C., Liu, Z., et al. (2020). Graph neural networks: A review of methods and applications. *AI Open* 1, 57. doi:10.1016/j.aiopen.2021.01.001

**Conflict of Interest:** The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

# Supplementary Material

## 1 MODEL VARIANCE

In general, deep learning models are bound to have a variance in the results arising from different trainings with different random seeds for the stochastic gradient descent optimization method. To ensure that does not affect our inferences based on the final anomaly identification results, we performed five separate trainings with different random seeds. Table S1 summarizes the median anomaly identification performance across all channels with the added variance across these five separate trainings, for the Conv-VAE model, the baseline GCN-VAE, and the GCN-VAE with the addition of the various normalizing flow models. Compared to the observed uncertainties, the improvements discussed in the paper are statistically significant.

**Table S1.** Anomaly detection performance across all channels, combined based on median scores along with the variance over multiple trainings.

Model	AUC	$\epsilon_S(\epsilon_B = 10^{-2})$	$\epsilon_S(\epsilon_B = 10^{-3})$	$\epsilon_S(\epsilon_B = 10^{-4})$
Conv-VAE	75.6% $\pm$ 0.5%	1.7% $\pm$ 0.6%	0.26% $\pm$ 0.04%	0.13% $\pm$ 0.07%
GCN-VAE	76.8% $\pm$ 0.2%	55.7% $\pm$ 1.1%	16.2% $\pm$ 0.6%	7.1% $\pm$ 0.1%
<b>GCN-VAE_HouseholderSNF</b>	<b>86.1% <math>\pm</math> 0.4%</b>	<b>69.6% <math>\pm</math> 1.9%</b>	<b>34.2% <math>\pm</math> 0.9%</b>	<b>8.2% <math>\pm</math> 0.5%</b>
GCN-VAE_OrthogonalSNF	82.4% $\pm$ 0.6%	65.0% $\pm$ 1.4%	16.1% $\pm$ 1.4%	7.6% $\pm$ 0.1%
GCN-VAE_NSFAR	82.1% $\pm$ 0.4%	64.1% $\pm$ 1.8%	15.6% $\pm$ 1.0%	5.1% $\pm$ 0.3%
GCN-VAE_PlanarFlow	80.0% $\pm$ 0.7%	61.1% $\pm$ 1.4%	16.2% $\pm$ 1.2%	7.1% $\pm$ 0.1%
GCN-VAE_ConvolutionalFlow	79.2% $\pm$ 0.2%	62.7% $\pm$ 1.4%	16.0% $\pm$ 1.4%	6.0% $\pm$ 0.2%
GCN-VAE_TriangularSNF	75.9% $\pm$ 0.5%	64.7% $\pm$ 1.5%	16.1% $\pm$ 1.3%	6.1% $\pm$ 0.6%
GCN-VAE_IAF	75.1% $\pm$ 0.8%	59.2% $\pm$ 2.7%	9.8% $\pm$ 0.8%	5.8% $\pm$ 0.7%

## 2 CHOICE OF LOSS FUNCTION

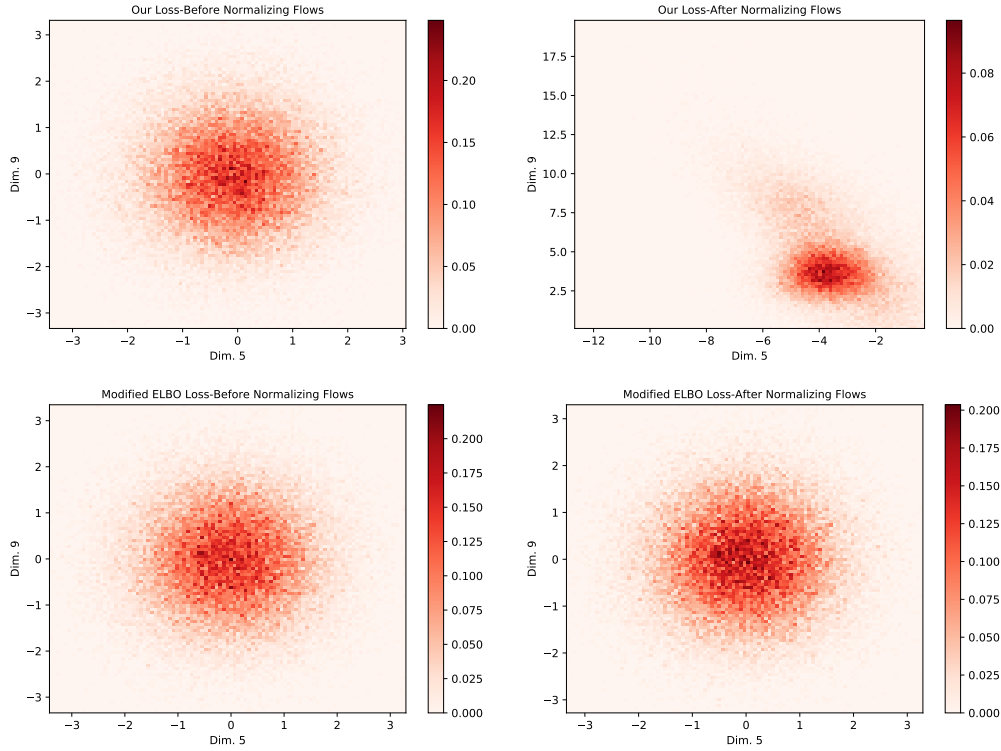
The optimization function used for normalizing flow models in variational inference is commonly formulated as the free energy bound (Rezende and Mohamed, 2015):

$$\mathcal{F}(x) = \mathbb{E}_{q_0(z_0)}[\ln q_0(z_0)] - \mathbb{E}_{q_0(z_0)}[\log q(x, z_K)] - \mathbb{E}_{q_0(z_0)} \left[ \sum_{k=1}^K \log |\det(J)| \right] \quad (9)$$

where  $\log |\det(J)|$  stand for the log-det-Jacobian term (Rezende and Mohamed, 2015) for the corresponding flow model. In our setup, this term would be added to the task-specific reconstruction loss. For our GCN-VAE\_HouseholderSNF model, we compare using this loss function with the added Chamfer loss described in Eq. (2), against using the same loss function for our baseline VAE model as described in Eq. (1). All other model parameters are kept identical between the two trainings. Table S2 shows the results of this comparison with “Modified ELBO” signifying Eq. 9 with the added Chamfer term, and “Our Loss” signifying Eq. (1). It is evident that not modifying the loss for the VAE after adding the flow layers results in significantly better anomaly identification performance and as a result we utilize this strategy to train all other flow models presented in this study.

**Table S2.** Anomaly detection performance across all channels, combined based on median scores.

Model	AUC	$\epsilon_S(\epsilon_B = 10^{-2})$	$\epsilon_S(\epsilon_B = 10^{-3})$	$\epsilon_S(\epsilon_B = 10^{-4})$
Modified ELBO	78.3%	57.1%	16.1%	6.6%
<b>Our Loss</b>	<b>86.4%</b>	<b>71.6%</b>	<b>34.0%</b>	<b>8.2%</b>

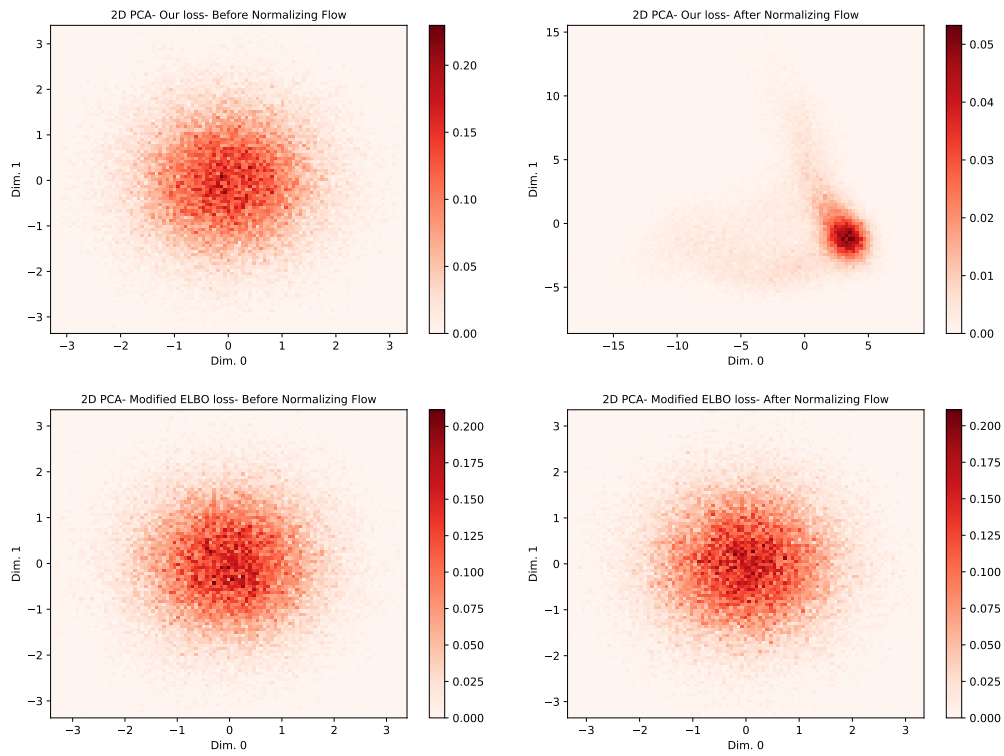


**Figure S1.** Latent space visualization by making histograms across arbitrarily chosen dimensions 5 and 9, before (left) and after normalizing flow transformations (right) with our loss function (top) and the modified ELBO loss (bottom).

### 3 IMPACT OF NORMALIZING FLOWS ON THE LATENT SPACE PRIOR

To understand exactly how non-Gaussian the latent distributions become after passing through the normalizing flow layers, we attempt to visualize our 15 dimensional latent space via two methods. First, we create 2D histograms from multiple randomly chosen pairs of dimensions. Fig. S1 shows one such distribution between dimensions 5 and 9. We also make an approximate visualisation by first performing a principal component analysis (PCA) to express the latent space in 2 dimensions, and then plotting the resulting 2D histogram as shown in Fig. S2.

We also show a comparison of the latent space obtained from the trainings with the two different loss functions, as described in the previous section of supplementary material. We see that our loss function results in a more complex, non-Gaussian distribution compared to the modified ELBO loss, and this is desirable to improve anomaly detection performance using VAEs. It is important to note that using our loss function may not necessarily correspond to better reconstruction of the input for the trained class (background samples) or better generation, but it rather contributes to a larger separation between the trained class and the non-trained class (signal samples) by making it harder for the decoder to reconstruct signal samples. As a result the anomaly identification performance increases regardless of whether the reconstruction of the background samples improves or not.



**Figure S2.** Latent space visualization after 2D PCA, before (left) and after normalizing flow transformations (right) with our loss function (top) and the modified ELBO loss (bottom).