

Why Normalizing Flows Fail to Detect Out-of-Distribution Data

Polina Kirichenko*, Pavel Izmailov*, Andrew Gordon Wilson
New York University

Abstract

Detecting out-of-distribution (OOD) data is crucial for robust machine learning systems. Normalizing flows are flexible deep generative models that often surprisingly fail to distinguish between in- and out-of-distribution data: a flow trained on pictures of clothing assigns higher likelihood to handwritten digits. We investigate why normalizing flows perform poorly for OOD detection. We demonstrate that flows learn local pixel correlations and generic image-to-latent-space transformations which are not specific to the target image dataset. We show that by modifying the architecture of flow coupling layers we can bias the flow towards learning the semantic structure of the target data, improving OOD detection. Our investigation reveals that properties that enable flows to generate high-fidelity images can have a detrimental effect on OOD detection.

1 Introduction

Normalizing flows [39, 9, 10] seem to be ideal candidates for out-of-distribution detection, since they are simple generative models that provide an exact likelihood. However, Nalisnick et al. [27] revealed the puzzling result that flows often assign higher likelihood to out-of-distribution data than the data used for maximum likelihood training. In Figure 1(a), we show the log-likelihood histogram for a RealNVP flow model [10] trained on the ImageNet dataset [35] subsampled to 64×64 resolution. The flow assigns higher likelihood to both the CelebA dataset of celebrity photos, and the SVHN dataset of images of house numbers, compared to the target ImageNet dataset.

While there has been empirical progress in improving OOD detection with flows [27, 7, 28, 36, 37, 45], the fundamental reasons for why flows fail at OOD detection in the first place are not fully understood. In this paper, we show how the *inductive biases* [26, 44] of flow models — implicit assumptions in the architectures and training procedures — can hinder OOD detection.

In particular, our contributions are the following:

- We show that flows learn latent representations for images largely based on local pixel correlations, rather than semantic content, making it difficult to detect data with anomalous semantics.
- We identify mechanisms through which normalizing flows can simultaneously increase likelihood for all structured images. For example, in Figure 1(b, c), we show that the coupling layers of RealNVP transform the in-distribution ImageNet in the same way as the OOD CelebA.

*Equal contribution.

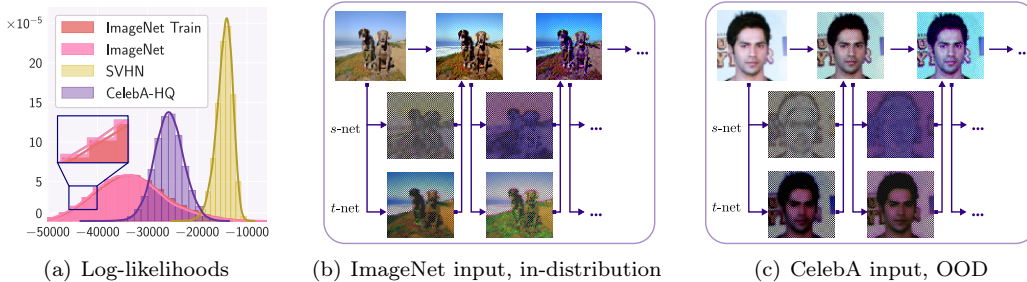


Figure 1: **RealNVP flow on in- and out-of-distribution images.** (a): A histogram of log-likelihoods that a RealNVP flow trained on ImageNet assigns to ImageNet, SVHN and CelebA. The flow assigns higher likelihood to out-of-distribution data. (b, c): A visualization of the intermediate layers of a RealNVP model on an (b) in-distribution image and (c) OOD image. The first row shows the coupling layer activations, the second and third rows show the scale s and shift t parameters predicted by a neural network applied to the corresponding coupling layer input. Both on in-distribution and out-of-distribution images, s and t accurately approximate the structure of the input, even though the model has not observed inputs (images) similar to the OOD image during training. *Flows learn generic image-to-latent-space transformations that leverage local pixel correlations and graphical details rather than the semantic content needed for OOD detection.*

- We show that by changing the architectural details of the coupling layers, we can encourage flows to learn transformations specific to the target data, improving OOD detection.
- We show that OOD detection is improved when flows are trained on high-level features which contain semantic information extracted from image datasets.

We also provide code at https://github.com/PolinaKirichenko/flows_ood.

2 Background

We briefly introduce normalizing flows based on coupling layers. For a more detailed introduction, see Papamakarios et al. [31] and Kobyzev et al. [23].

Normalizing flows Normalizing flows [39] are a flexible class of deep generative models that model a target distribution $p^*(x)$ as an invertible transformation f of a base distribution $p_Z(z)$ in the latent space. Using the change of variables formula, the likelihoods for an input x and a dataset \mathcal{D} are

$$p_X(x) = p_Z(f^{-1}(x)) \left| \det \frac{\partial f^{-1}}{\partial x} \right|, \quad p(\mathcal{D}) = \prod_{x \in \mathcal{D}} p_X(x). \quad (1)$$

The latent space distribution $p_Z(z)$ is commonly chosen to be a standard Gaussian. Flows are typically trained by maximizing the log-likelihood (1) of the training data with respect to the parameters of the invertible transformation f .

Coupling layers We focus on normalizing flows based on *affine coupling layers*. In these flows, the transformation performed by each layer is given by

$$f_{\text{aff}}^{-1}(x_{\text{id}}, x_{\text{change}}) = (y_{\text{id}}, y_{\text{change}}), \quad \begin{cases} y_{\text{id}} = x_{\text{id}} \\ y_{\text{change}} = (x_{\text{change}} + t(x_{\text{id}})) \odot \exp(s(x_{\text{id}})) \end{cases} \quad (2)$$

where x_{id} and x_{change} are disjoint parts of the input x , y_{id} and y_{change} are disjoint parts of the output y , and the scale and shift parameters $s(\cdot)$ and $t(\cdot)$ are usually implemented by

a neural network (which we will call the *st-network*). The split of the input into x_{id} and x_{change} is defined by a *mask*: a coupling layer transforms the masked part $x_{\text{change}} = \text{mask}(x)$ of the input based on the remaining part x_{id} . The transformation (2) is invertible and allows for efficient Jacobian computation in (1):

$$\log \left| \det \frac{\partial f_{\text{aff}}^{-1}}{\partial x} \right| = \sum_{i=1}^{\dim(x_{\text{change}})} s(x_{\text{id}})_i. \quad (3)$$

Flows with coupling layers Coupling layers can be stacked together into flexible normalizing flows: $f = f^K \circ f^{K-1} \circ \dots \circ f^1$. Examples of flows with coupling layers include NICE [9], RealNVP [10], Glow [22], and many others [e.g., 4, 5, 11, 15, 16, 21, 25, 32].

Out-of-distribution detection using likelihood Flows can be used for out-of-distribution detection based on the likelihood they assign to the inputs. One approach is to choose a likelihood threshold ϵ on a validation dataset, e.g. to satisfy a desired false positive rate, and during test time identify inputs which have likelihood lower than ϵ as OOD. Qualitatively, we can estimate the performance of the flows for OOD detection by plotting a histogram of the log-likelihoods such as Figure 1(a): the likelihoods for in-distribution data should generally be higher compared to OOD. Alternatively, we can treat OOD detection as a binary classification problem using likelihood scores, and compute accuracy with a fixed likelihood threshold ϵ , or AUROC (area under the receiver operating characteristic curve).

3 Related Work

Recent works have shown that normalizing flows, among other deep generative models, can assign higher likelihood to out-of-distribution data [27, 7]. The work on OOD detection with deep generative models falls into two distinct categories. In group anomaly detection (GAD), the task is to label a batch of $n > 1$ datapoints as in- or out-of-distribution. Point anomaly detection (PAD) involves the more challenging task of labelling single points as out-of-distribution.

Group anomaly detection Nalisnick et al. [28] introduce the typicality test which distinguishes between a high density set and a typical set of a distribution induced by a model. However, the typicality test cannot detect OOD data if the flow assigns it with a similar likelihood distribution to that of in-distribution data. Song et al. [37] showed that out-of-distribution datasets have lower likelihoods when batch normalization statistics are computed from a current batch instead of accumulated over the train set, and proposed a test based on this observation. Zhang et al. [45] introduce a GAD algorithm based on measuring correlations of flow’s latent representations corresponding to the input batch. The main limitation of GAD methods is that for most practical applications the assumption that the data comes in batches of inputs that are all in-distribution or all OOD is not realistic.

Point anomaly detection Choi et al. [7] proposed to estimate the Watanabe-Akaike Information Criterion using an ensemble of generative models, showing accurate OOD detection on some of the challenging dataset pairs. Ren et al. [33] explain the poor OOD detection performance of deep generative models by the fact that the likelihood is dominated by background statistics. They propose a test based on the ratio of the likelihoods for the image and background likelihood estimated using a separate *background model*. Serrà et al. [36] show that normalizing flows assign higher likelihoods to simpler datasets and propose to normalize the flow’s likelihood by an image complexity score.

In this work we argue that it is the inductive biases of the model that determine its OOD performance. While most work treats flows as black-box density estimators, we conduct a careful study of the latent representations and image-to-latent-space transformations learned by the flows. Throughout the paper, we connect our findings with prior work and provide new insights.

4 Why flows fail to detect OOD data

Normalizing flows consistently fail at out-of-distribution detection when applied to common benchmark datasets (see Appendix D). In this paper, we discuss the reasons behind this surprising phenomenon. We summarize our thesis as follows:

The maximum likelihood objective has a limited influence on OOD detection, relative to the *inductive biases* of the flow, captured by the modelling assumptions of the architecture.

Why should flows be able to detect OOD inputs? Flows are trained to maximize the likelihood of the training data. Likelihood is a probability density function $p(\mathcal{D})$ defined on the image space and hence has to be normalized. Thus, likelihood cannot be simultaneously increased for all the inputs (images). In fact, the optimal maximizer of (1) would only assign positive density to the datapoints in the training set, and, in particular, would not even generalize to the test set of the same dataset. In practice, flows do not seem to overfit, assigning similar likelihood distributions to train and test (see e.g. Figure 1(a)). Thus, despite their flexibility, flows are not maximizing the likelihood (1) to values close to the global optimum.

What is OOD data? There are infinitely many distributions that give rise to any value of the likelihood objective in (1) except the global optimum. Indeed, any non-optimal solution assigns probability mass outside of the training data distribution; we can arbitrarily re-assign this probability mass to get a new solution with the same value of the objective (see Appendix A for a detailed discussion). Therefore the inductive biases of a model determines which specific solution is found through training. In particular, the inductive biases will affect what data is assigned high likelihood (in-distribution) and what data is not (OOD).

What inductive biases are needed for OOD detection? The datasets in computer vision are typically defined by the semantic content of the images. For example, the CelebA dataset consists of images of faces, and SVHN contains images of house numbers. In order to detect OOD data, the inductive biases of the model have to be aligned with learning the semantic structure of the data, i.e. what objects are represented in the data.

What are the inductive biases of normalizing flows? In the remainder of the paper, we explore the inductive biases of normalizing flows. We argue that flows are biased towards learning *graphical* properties of the data such as local pixel correlations (e.g. nearby pixels usually have similar colors) rather than semantic properties of the data (e.g. what objects are shown in the image).

Flows have capacity to distinguish datasets In Appendix B, we show that if we explicitly train flows to distinguish between a pair of datasets, they can assign large likelihood to one dataset and low likelihood to the other. However, when trained with the standard maximum likelihood objective, flows do not learn to make this distinction. The inductive biases of the flows prefer solutions that assign high likelihood to most structured datasets simultaneously.

5 Flow latent spaces

Normalizing flows learn highly non-linear image-to-latent-space mappings often using hundreds of millions of parameters. One could imagine that the learned latent representations have a complex structure, encoding high-level semantic information about the inputs. In this section, we visualize the learned latent representations on both in-distribution and

out-of-distribution data and demonstrate that they encode simple graphical structure rather than semantic information.

Observation: There exists a correspondence between the coordinates in an image and in its learned representation. We can recognize edges of the inputs in their latent representations.

Significance for OOD detection: In order to detect OOD images, a model has to assign likelihood based on the semantic content of the image (see Sec. 4). Flows do not represent images based on their semantic contents, but rather directly encode their visual appearance.

In the first four columns of Figure 2, we show latent representations¹ of a RealNVP model trained on FashionMNIST for an in-distribution FashionMNIST image and an out-of-distribution MNIST digit. The first column shows the original image x , and the second column shows the corresponding latent z . The latent representations appear noisy both for in- and out-of-distribution samples, but the edges of the MNIST digit can be recognized in the latent. In the third column of Figure 2, we show latent representations averaged over $K = 40$ samples of dequantization noise² $\epsilon_k: \frac{1}{K} \sum_{k=1}^K f^{-1}(x + \epsilon_k)$. In the averaged representation, we can clearly see the edges from the original image. Finally, in the fourth column of Figure 2, we visualize the latent representations (for a single sample of dequantization noise) from a flow when batch normalization layers are in train mode [18]. In train mode, batch normalization layers use the activation statistics of the current batch, and in evaluation mode they use the statistics accumulated over the train set. While for in-distribution data there is no structure visible in the latent representation, the out-of-distribution latent clearly preserves the shape of the 7-digit from the input image. In the remaining panels of Figure 2, we show an analogous visualization for a RealNVP trained on CelebA using an SVHN image as OOD. In the third panel of this group, we visualize the blue channel of the latent representations. Again, the OOD input can be recognized in the latent representation; some of the edges from the in-distribution CelebA image can also be seen in the corresponding latent variable. Additional visualizations (e.g. for Glow) are in Appendix F.

Insights into prior work The group anomaly detection algorithm proposed in Zhang et al. [45] uses correlations of the latent representations as an OOD score. Song et al. [37] showed that normalizing flows with batch normalization layers in train mode assign much lower likelihood to out-of-distribution images than they do in evaluation mode, while for in-distribution data the difference is not significant. Our visualizations explain the presence of correlations in the latent space and shed light into the difference between the behaviour of the flows in train and test mode.

6 Transformations learned by coupling layers

To better understand the inductive biases of coupling-layer based flows, we study the transformations learned by individual coupling layers.

What are coupling layers trained to do? Each coupling layer updates the masked part x_{change} of the input x to be $x_{\text{change}} \leftarrow (x_{\text{change}} + t(x_{\text{id}})) \cdot \exp(s(x_{\text{id}}))$, where x_{id} is the non-masked part of x , and s and t are the outputs of the st -network given x_{id} (see Section 2). The flow is encouraged to predict high values for s since for a given coupling layer the Jacobian term in the likelihood of Eq. (1) is given by $\sum_j s(x_{\text{id}})_j$ (see Section 4). Intuitively,

¹ For the details of the visualization procedure and the training setup please see Appendices E and C.

² When training flow models on images or other discrete data, we use dequantization to avoid pathological solutions [43, 41]: we add uniform noise $\epsilon \sim U[0, 1]$ to each pixel $x_i \in \{0, 1, \dots, 255\}$. Every time we pass an image through the flow $f(\cdot)$, the resulting latent representation z will be different.

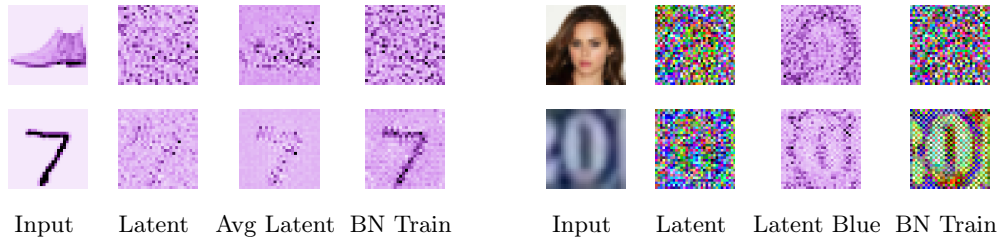


Figure 2: **Latent spaces.** Visualization of latent representations for the RealNVP model on in-distribution and out-of-distribution inputs. Panels 1-4: original images, latent representations, latent representation averaged over 40 samples of dequantization noise, and latent representations for batch normalization in train mode for a flow trained on FashionMNIST and using MNIST for OOD data. Panels 5-8: same as 1-4 but for a model trained on CelebA with SVHN for OOD, except in panel 7 we show the blue channel of the latent representation from panel 6 instead of an averaged latent representation. For both dataset pairs, we can recognize the shape of the input image in the latent representations. The flow represents images based on their graphical appearance rather than semantic content.

to afford large values for scale s without making the latent representations large in norm and hence decreasing the density term $p_Z(z)$ in (1), the shift $-t$ has to be an accurate approximation of the masked input x_{change} . For example, in Figure 1(b, c) the $-t$ outputs of the first coupling layers are a very close estimate of the input to the coupling layer. The likelihood for a given image will be high whenever the coupling layers can accurately predict masked pixels. To the best of our knowledge, this intuition has not been discussed in any previous work.

Observation: We describe two mechanisms through which coupling layers learn to predict the masked pixels: (1) leveraging local color correlations and (2) using information about the masked pixels encoded by the previous coupling layer (coupling layer co-adaptation).

Significance for OOD detection: These mechanisms allow the flows to predict the masked pixels equally accurately on in- and out-of-distribution datasets. As a result, flows assign high likelihood to OOD data.

6.1 Leveraging local pixel correlations

In Figure 3(a, b), we visualize intermediate coupling layer activations of a small RealNVP model with 2 coupling layers and checkerboard masks trained on FashionMNIST. For the masked inputs, the outputs of the st -network are shown in black. Even though the flow was trained on FashionMNIST and has never seen an MNIST digit, the st -networks can easily predict masked from observed pixels on both FashionMNIST and MNIST. Figure 1 shows the same behaviour in the first coupling layers of RealNVP trained on ImageNet.

With the checkerboard mask, the st -networks predict the masked pixels from neighbouring pixels (see Appendix G for a discussion of different masks). Natural images have local structure and correlations: with a high probability, a particular pixel value will be similar to its neighbouring pixels. The checkerboard mask creates an inductive bias for the flow to pick up on these local correlations. In Figure 3, we can see that the outputs of the s -network are especially large for the background pixels and large patches of the same color (larger

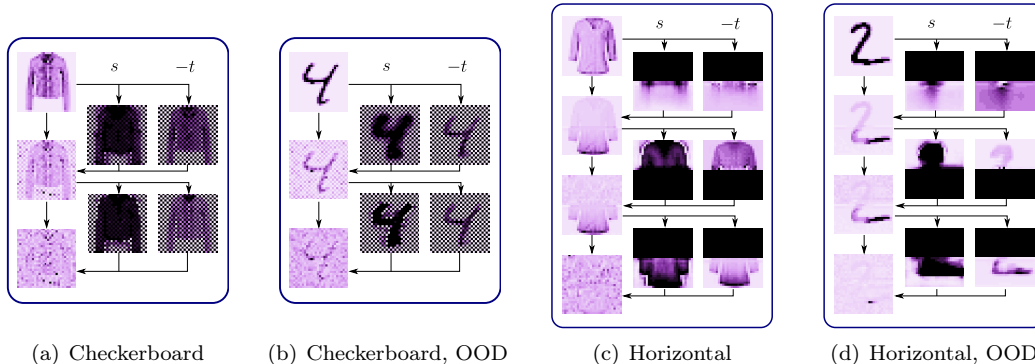


Figure 3: **Coupling layers.** Visualization of RealNVP’s intermediate coupling layer activations, as well as scales s and shifts t predicted by each coupling layer on in-distribution (panels a, c) and out-of-distribution inputs (panels b, d). RealNVP was trained on FashionMNIST. (a), (b): RealNVP with a standard checkerboard masks. The st -networks are able to predict the masked pixels well both on in-distribution and OOD inputs from neighbouring pixels. (c), (d): RealNVP with a horizontal mask. Despite being trained on FashionMNIST, the st -networks are able to correctly predict the bottom half of MNIST digits in the second coupling layer due to coupling layer co-adaptation.

values are shown with lighter color), where the flow simply predicts for example that a pixel surrounded by black pixels would itself be black.

In addition to the checkerboard mask, RealNVP and Glow also use channel-wise masks. These masks are applied after a squeeze layer, which puts different subsampled versions of the image in different channels. As a result, the st -network is again trained to predict pixel values from neighbouring pixels. We provide additional visualizations for RealNVP and Glow in Appendix H.

6.2 Coupling layer co-adaptation

To better understand the transformations learned by the coupling layers, we replaced the standard masks in RealNVP with a sequence of *horizontal masks* that cover one half of the image (either top or bottom). For example, the first coupling layer of the flow shown in panels (c, d) of Figure 3 transforms the bottom half of the image based on the top half, the second layer transforms the top half based on the bottom half, and so on. In Figure 3(c, d) we visualize the coupling layers for a 3-layer RealNVP with horizontal masks on in-distribution (FashionMNIST) and OOD (MNIST) data.

In the first coupling layer, the shift output $-t$ of the st -network predicts the bottom half of the image poorly and the layer does not seem to transform the input significantly. In the second and third layer, $-t$ presents an almost ideal reconstruction of the masked part of the image on both the in-distribution and, surprisingly, the OOD input. It is not possible for the st -network that was only trained on FashionMNIST to predict the top half of an MNIST digit based on the other half. The resolution is that the first layer encodes information about the top half into the bottom half of the image; the second layer then decodes this information to accurately predict the top half. Similarly, the third layer leverages information about the bottom half of the image encoded by the second layer. We refer to this phenomenon as *coupling layer co-adaptation*. Additional visualizations are in Appendix H.

Horizontal masks allow us to conveniently visualize the coupling layer co-adaptation, but we hypothesize that the same mechanism applies to standard checkerboard and channel-wise masks in combination with local color correlations.

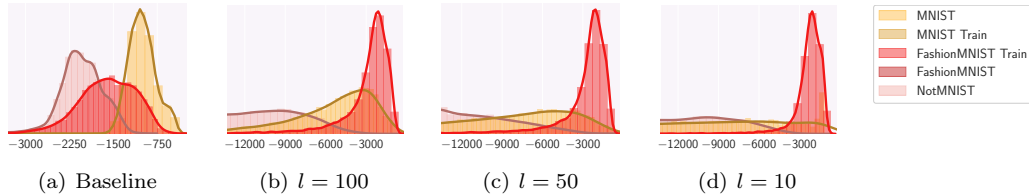


Figure 4: **Effect of st -networks capacity.** Histograms of log-likelihoods of in- and out-of-distribution data for RealNVP trained on FashionMNIST, varying the dimension l of the bottleneck in the st -networks. Flows with lower l work better for OOD detection: the baseline assigns higher likelihood to the out-of-distribution MNIST images, while the flows with $l = 50$ and $l = 10$ assign much higher likelihood to in-distribution FashionMNIST data. With $l = 100$ the flow assigns higher likelihood to in-distribution data, but the overlap of the likelihood distribution with OOD MNIST is higher than for $l = 50$ and $l = 10$.

Insights into prior work Prior work showed that the likelihood score is heavily affected by the input complexity [36] and background statistics [33]; however, prior work does not explain *why* flows exhibit such behavior. Simpler images (e.g. SVHN compared to CIFAR-10) and background often contain large patches of the same color, which makes it easy to predict masked pixels from their neighbours and to encode and decode the information via coupling layer co-adaptation.

7 Changing biases in flows for better OOD detection

Our observations in Sections 5 and 6 suggest that normalizing flows are biased towards learning transformations that increase likelihood simultaneously for all structured images. We discuss two simple ways of changing the inductive biases for better OOD detection.

By changing the masking strategy or the architecture of st -networks in flows we can improve OOD detection based on likelihood.

Changing masking strategy We consider two three types of masks. We introduced the horizontal mask in Section 6.2: in each coupling layer the flow updates the bottom half of the image based on the top half or vice versa. With a horizontal mask, flows cannot simply use the information from neighbouring pixels when predicting a given pixel, but they exhibit coupling layer co-adaptation (see Section 6.2). To combat coupling layer co-adaptation, we additionally introduce the *cycle-mask*, a masking strategy where the information about a part of the image has to travel through three coupling layers before it can be used to update the same part of the image (details in Appendix I.1). To compare the performance of the checkerboard mask, horizontal mask and cycle-mask, we construct flows of exactly the same size and architecture (RealNVP with 8 coupling layers and no squeeze layers) with each of these masks, trained on CelebA and FashionMNIST. We present the results in the Appendix I.1. As expected, for the checkerboard mask, the flow assigns higher likelihood to the simpler OOD datasets (SVHN for CelebA and MNIST for FashionMNIST). With the horizontal mask, the OOD data still has higher likelihood on average, but the relative ranking of the in-distribution data is improved. Finally, for the cycle-mask, on FashionMNIST the likelihood is higher compared to MNIST on average. On CelebA the likelihood is similar but slightly lower compared to SVHN.

st -networks with bottleneck Another way to force the flow to learn global structure rather than local pixel correlations and to prevent coupling layer co-adaptation is to restrict

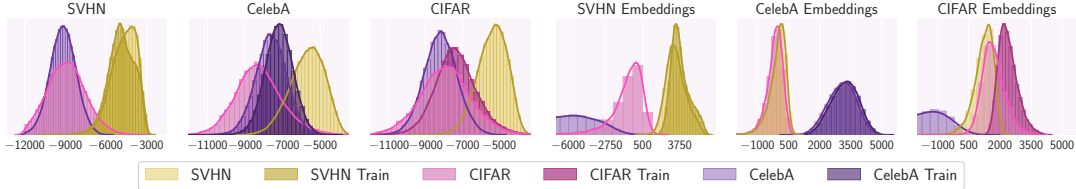


Figure 5: **Image embeddings.** Log-likelihood histograms for RealNVP trained on raw pixel data (first three panels) and embeddings extracted for the same image datasets using EfficientNet trained on ImageNet. On raw pixels, the flow assigns the highest likelihood to SVHN regardless of its training dataset. On image embeddings, flows always assign higher likelihood to in-distribution data. When trained on features capturing the semantic content of the input, flows can detect OOD.

the capacity of the *st*-networks. To do so, we introduce a *bottleneck* to the *st*-networks: a pair of fully-connected layers projecting to a space of dimension l and back to the original input dimension. We insert these layers after the middle layer of the *st*-network. If the latent dimension l is small, the *st*-network cannot simply reproduce its input as its output, and thus cannot exploit the local pixel correlations discussed in Section 6. Passing information through multiple layers with a low-dimensional bottleneck also reduces the effect of coupling layer co-adaptation. We train a RealNVP flow varying the latent dimension l on CelebA and on FashionMNIST. We present the results in Figure 4 and Appendix I. On FashionMNIST, introducing the bottleneck forces the flow to assign lower likelihood to OOD data (Figure 4). Furthermore, as we decrease l , the likelihood of the OOD data decreases but FashionMNIST likelihood stays the same. On CelebA the relative ranking of likelihood for in-distribution data is similarly improved when we decrease the dimension l of the bottleneck, but SVHN is still assigned slightly higher likelihood than CelebA. See Appendix I for detailed results.

While the proposed modifications do not completely resolve the issue of OOD data having higher likelihood, the experiments support our observations in Section 6: preventing the flows from leveraging local color correlations and coupling layer co-adaptation, we improve the relative likelihood ranking for in-distribution data.

8 Out-of-distribution detection using image embeddings

In Section 4 we argued that in order to detect OOD data the model has to assign likelihood based on high-level semantic features of the data, which the flows fail to do when trained on images. In this section, we test out-of-distribution detection using image representations from a deep neural network.

Normalizing flows can detect OOD images when trained on high-level semantic representations instead of raw pixels.

We extract embeddings for CIFAR-10, CelebA and SVHN using an EfficientNet [40] pretrained on ImageNet [35] which yields 1792-dimensional features³. We train RealNVP on each of the representation datasets, considering the other two datasets as OOD. We present the likelihood histograms for all datasets in Figure 5(b). Additionally, we report AUROC scores in Appendix Table 2. For the models trained on SVHN and CelebA, both OOD datasets have lower likelihood and the AUROC scores are close to 100%. For the model trained on CIFAR-10, CelebA has lower likelihood. Moreover, the likelihood distribution on SVHN,

³ The original images are 3072-dimensional, so the dimension of the embeddings is only two times smaller. Thus, the inability to detect OOD images *cannot* be explained just by the high dimensionality of the data.

while significantly overlapping with CIFAR-10, still has a lower average: the AUROC score between CIFAR-10 and SVHN is 73%. Flows are much better at OOD detection on image embeddings than on the original image datasets. For example, a flow trained on CelebA images assigns higher likelihood to SVHN, while a flow trained on CelebA embeddings assigns low likelihood to SVHN embeddings (see Appendix D for likelihood distribution and AUROC scores on image data).

Non-image data In Appendix K we evaluate flows on tabular UCI datasets, where the features are relatively high-level compared to images. On these datasets, normalizing flows assign higher likelihood to in-distribution data.

9 Conclusion

Many of the puzzling phenomena in deep learning can be boiled down to a matter of *inductive biases*. Neural networks in many cases have the flexibility to overfit datasets, but they do not because the biases of the architecture and training procedures can guide us towards reasonable solutions. In performing OOD detection, the biases of normalizing flows can be more of a curse than a blessing. Indeed, we have shown that flows tend to learn representations that achieve high likelihood through generic graphical features and local pixel correlations, rather than discovering semantic structure that would be specific to the training distribution.

To provide insights into prior results [e.g., 27, 7, 28, 37, 45, 36], part of our discussion has focused on an in-depth exploration of the popular class of normalizing flows based on affine coupling layers. We hypothesize that many of our conclusions about coupling layers extend at a high level to other types of normalizing flows [e.g., 3, 6, 12, 20, 13, 30, 38, 17, 8]. A full study of these other types of flows is a promising direction for future work.

Acknowledgements

PK, PI, and AGW are supported by an Amazon Research Award, Amazon Machine Learning Research Award, Facebook Research, NSF I-DISRE 193471, NIH R01 DA048764-01A1, NSF IIS-1910266, and NSF 1922658 NRT-HDR: FUTURE Foundations, Translation, and Responsibility for Data Science. We thank Marc Finzi, Greg Benton, Wesley Maddox, and Alex Wang for helpful discussions.

References

- [1] Andrei Atanov, Alexandra Volokhova, Arsenii Ashukha, Ivan Sosnovik, and Dmitry Vetrov. Semi-conditional normalizing flows for semi-supervised learning. *arXiv preprint arXiv:1905.00505*, 2019.
- [2] Pierre Baldi, Kyle Cranmer, Taylor Faucett, Peter Sadowski, and Daniel Whiteson. Parameterized machine learning for high-energy physics. *arXiv preprint arXiv:1601.07913*, 2016.
- [3] Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Invertible residual networks. *arXiv preprint arXiv:1811.00995*, 2018.
- [4] Apratim Bhattacharyya, Shweta Mahajan, Mario Fritz, Bernt Schiele, and Stefan Roth. Normalizing flows with multi-scale autoregressive priors. In *Proceedings of the Conference on Computer Vision and Pattern Recognition*, pages 8415–8424, 2020.

- [5] Jianfei Chen, Cheng Lu, Biqi Chenli, Jun Zhu, and Tian Tian. Vflow: More expressive generative flows with variational data augmentation. *arXiv preprint arXiv:2002.09741*, 2020.
- [6] Ricky TQ Chen, Jens Behrmann, David Duvenaud, and Jörn-Henrik Jacobsen. Residual flows for invertible generative modeling. *arXiv preprint arXiv:1906.02735*, 2019.
- [7] Hyunsun Choi, Eric Jang, and Alexander A Alemi. Waic, but why? generative ensembles for robust anomaly detection. *arXiv preprint arXiv:1810.01392*, 2018.
- [8] Nicola De Cao, Ivan Titov, and Wilker Aziz. Block neural autoregressive flow. *arXiv preprint arXiv:1904.04676*, 2019.
- [9] Laurent Dinh, David Krueger, and Yoshua Bengio. NICE: Non-linear independent components estimation. *arXiv preprint arXiv:1410.8516*, 2014.
- [10] Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio. Density estimation using Real NVP. *arXiv preprint arXiv:1605.08803*, 2016.
- [11] Conor Durkan, Artur Bekasov, Iain Murray, and George Papamakarios. Neural spline flows. In *Advances in Neural Information Processing Systems*, pages 7509–7520, 2019.
- [12] Marc Finzi, Pavel Izmailov, Wesley Maddox, Polina Kirichenko, and Andrew Gordon Wilson. Invertible convolutional networks. In *Workshop on Invertible Neural Nets and Normalizing Flows, International Conference on Machine Learning*, 2019.
- [13] Will Grathwohl, Ricky TQ Chen, Jesse Betterncourt, Ilya Sutskever, and David Duvenaud. Ffjord: Free-form continuous dynamics for scalable reversible generative models. *arXiv preprint arXiv:1810.01367*, 2018.
- [14] Dan Hendrycks, Mantas Mazeika, and Thomas Dietterich. Deep anomaly detection with outlier exposure. *arXiv preprint arXiv:1812.04606*, 2018.
- [15] Jonathan Ho, Xi Chen, Aravind Srinivas, Yan Duan, and Pieter Abbeel. Flow++: Improving flow-based generative models with variational dequantization and architecture design. *arXiv preprint arXiv:1902.00275*, 2019.
- [16] Emiel Hooeboom, Rianne van den Berg, and Max Welling. Emerging convolutions for generative normalizing flows. *arXiv preprint arXiv:1901.11137*, 2019.
- [17] Chin-Wei Huang, David Krueger, Alexandre Lacoste, and Aaron Courville. Neural autoregressive flows. *arXiv preprint arXiv:1804.00779*, 2018.
- [18] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [19] Pavel Izmailov, Polina Kirichenko, Marc Finzi, and Andrew Gordon Wilson. Semi-supervised learning with normalizing flows. In *International Conference on Machine Learning*, 2020.
- [20] Mahdi Karami, Dale Schuurmans, Jascha Sohl-Dickstein, Laurent Dinh, and Daniel Duckworth. Invertible convolutional flow. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages 5635–5645. Curran Associates, Inc., 2019. URL <http://papers.nips.cc/paper/8801-invertible-convolutional-flow.pdf>.
- [21] Sungwon Kim, Sang-gil Lee, Jongyoon Song, Jaehyeon Kim, and Sungroh Yoon. Flowavenet: A generative flow for raw audio. *arXiv preprint arXiv:1811.02155*, 2018.
- [22] Durk P Kingma and Prafulla Dhariwal. Glow: Generative flow with invertible 1×1 convolutions. In *Advances in Neural Information Processing Systems*, pages 10215–10224, 2018.

- [23] Ivan Kobyzev, Simon Prince, and Marcus A Brubaker. Normalizing flows: Introduction and ideas. *arXiv preprint arXiv:1908.09257*, 2019.
- [24] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [25] Xuezhe Ma, Xiang Kong, Shanghang Zhang, and Eduard Hovy. Macow: Masked convolutional generative flow. In *Advances in Neural Information Processing Systems*, pages 5891–5900, 2019.
- [26] Tom M Mitchell. *The need for biases in learning generalizations*. Department of Computer Science, Laboratory for Computer Science Research . . . , 1980.
- [27] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, and Balaji Lakshminarayanan. Do deep generative models know what they don’t know? *arXiv preprint arXiv:1810.09136*, 2018.
- [28] Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, and Balaji Lakshminarayanan. Detecting out-of-distribution inputs to deep generative models using a test for typicality. *arXiv preprint arXiv:1906.02994*, 2019.
- [29] Aaron van den Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel recurrent neural networks. *arXiv preprint arXiv:1601.06759*, 2016.
- [30] George Papamakarios, Theo Pavlakou, and Iain Murray. Masked autoregressive flow for density estimation. In *Advances in Neural Information Processing Systems*, pages 2338–2347, 2017.
- [31] George Papamakarios, Eric Nalisnick, Danilo Jimenez Rezende, Shakir Mohamed, and Balaji Lakshminarayanan. Normalizing flows for probabilistic modeling and inference. *arXiv preprint arXiv:1912.02762*, 2019.
- [32] Ryan Prenger, Rafael Valle, and Bryan Catanzaro. Waveglow: A flow-based generative network for speech synthesis. In *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3617–3621. IEEE, 2019.
- [33] Jie Ren, Peter J Liu, Emily Fertig, Jasper Snoek, Ryan Poplin, Mark Depristo, Joshua Dillon, and Balaji Lakshminarayanan. Likelihood ratios for out-of-distribution detection. In *Advances in Neural Information Processing Systems*, pages 14680–14691, 2019.
- [34] Byron P Roe, Hai-Jun Yang, Ji Zhu, Yong Liu, Ion Stancu, and Gordon McGregor. Boosted decision trees as an alternative to artificial neural networks for particle identification. *Nuclear Instruments and Methods in Physics Research Section A: Accelerators, Spectrometers, Detectors and Associated Equipment*, 543(2-3):577–584, 2005.
- [35] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- [36] Joan Serra, David Álvarez, Vicenç Gómez, Olga Slizovskaia, José F Núñez, and Jordi Luque. Input complexity and out-of-distribution detection with likelihood-based generative models. *arXiv preprint arXiv:1909.11480*, 2019.
- [37] Jiaming Song, Yang Song, and Stefano Ermon. Unsupervised out-of-distribution detection with batch normalization. *arXiv preprint arXiv:1910.09115*, 2019.
- [38] Yang Song, Chenlin Meng, and Stefano Ermon. Mintnet: Building invertible neural networks with masked convolutions. In *Advances in Neural Information Processing Systems*, pages 11002–11012, 2019.

- [39] Esteban G Tabak and Cristina V Turner. A family of nonparametric density estimation algorithms. *Communications on Pure and Applied Mathematics*, 66(2):145–164, 2013.
- [40] Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
- [41] Lucas Theis, Aäron van den Oord, and Matthias Bethge. A note on the evaluation of generative models. *arXiv preprint arXiv:1511.01844*, 2015.
- [42] Antonio Torralba, Rob Fergus, and William T Freeman. 80 million tiny images: A large data set for nonparametric object and scene recognition. *IEEE transactions on pattern analysis and machine intelligence*, 30(11):1958–1970, 2008.
- [43] Benigno Uria, Iain Murray, and Hugo Larochelle. Rnade: The real-valued neural autoregressive density-estimator. In *Advances in Neural Information Processing Systems*, pages 2175–2183, 2013.
- [44] Andrew Gordon Wilson and Pavel Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. *arXiv preprint arXiv:2002.08791*, 2020.
- [45] Yufeng Zhang, Wanwei Liu, Zhenbang Chen, Ji Wang, Zhiming Liu, Kenli Li, Hongmei Wei, and Zuoning Chen. Out-of-distribution detection with distance guarantee in deep generative models. *arXiv preprint arXiv:2002.03328*, 2020.

Appendix outline

This appendix is organized as follows.

- In Section A, we provide additional discussion and a formal statement of the argument presented in Section 4.
- In Section B, we show that normalizing flows can be trained to assign high likelihood to the target data and low likelihood to a given OOD dataset.
- In Section C, we provide the hyperparameters that we used for the experiments in this paper.
- In Section D, we report the log-likelihood histograms and OOD detection AUROC scores for the baseline RealNVP and Glow models on various datasets.
- In Section E, we explain the visualization procedure that we use to visualize the latent representations and coupling layers of normalizing flows.
- In Section F, we provide additional latent representation visualizations.
- In Section G, we explain the different masking strategies for coupling layers of normalizing flows.
- In Section H, we provide additional coupling layer visualizations.
- In Section I, we provide additional details on the experiments of Section 7.
- In Section J, we provide samples from baseline models on various datasets. We also discuss an experiment on resampling parts of the latent variables corresponding to different images with normalizing flows.
- In Section K, we provide additional details and results for the experiments on image embeddings and tabular data from Section 8.

A Maximum likelihood objective is agnostic to what data is OOD

In Section 4 we argued that the maximum likelihood objective by itself does not define out-of-distribution detection performance of a normalizing flow. Instead, it is the inductive biases of the flow that define what data will be assigned with high or low likelihood. We illustrate this point in Figure 6.

The yellow and red shaded regions illustrate the high-probability regions of two distributions defined on the image space \mathcal{X} . The distribution in yellow assigns high likelihood to the train (CelebA) images corrupted by a small level of noise, or brightness adjustments. This distribution represents how a human could describe the target dataset. The red distribution on the other hand assigns high likelihood to all structured images including those from ImageNet and SVHN, but does not support noisy train images. The red distribution represents a distribution learned by normalizing flow.

For simplicity, we could think that the distributions are uniform on the highlighted sets, and the sets have the same volume. Then, both distributions assign equally high likelihood to the training data, but the split of the data into in-distribution and OOD is different. As

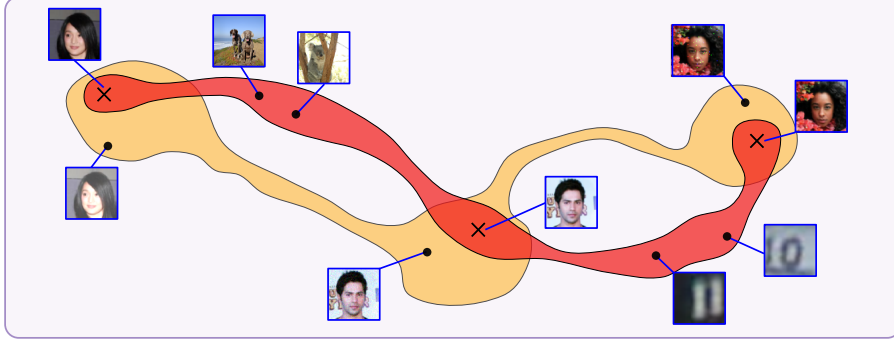


Figure 6: **Inductive biases define what data is OOD.** A conceptual visualization of two distributions in the image space (shown in yellow and red), training CelebA data is shown with crosses, and other images are shown with circles. The distribution shown in yellow could represent inductive biases of a human: it assigns high likelihood to all images of human faces, regardless of small levels of noise, and small brightness changes. The second distribution, shown in red, could represent a normalizing flow: it assigns high likelihood to all smooth structured images, including images from SVHN and ImageNet. Both distributions assign the same likelihood to the training set, but their high-probability sets are different.

both distributions provide the same density to the target data, the value of the maximum likelihood objective in Equation (1) would be the same for the corresponding models.

More generally, for any distribution that only assigns finite density to the train set, we can construct another distribution that assigns the same density to the train data, but also high density to a given set of (OOD) datapoints. In particular, the new distribution will achieve the same value of the maximum likelihood objective in Equation (1). We formalize our reasoning in the following simple proposition.

Proposition 1. *Let $p(\cdot)$ be a probability density on the space \mathcal{X} , and let $\mathcal{D} = \{x_i\}_{i=1}^N$ be the training dataset, where $x_i \in \mathcal{X}$ for $i = 1, \dots, N$. Assume for simplicity that p is upper bounded: for any x $p(x) \leq u$. Let \mathcal{D}_{OOD} be an arbitrary finite set of points. Then, for any $c \geq 0$ there exists a distribution with density $p'(\cdot)$ such that $p'(x) = p(x)$ for all $x \in \mathcal{D}$, and $p'(x') \geq c$ for all $x' \in \mathcal{D}_{OOD}$.*

Proof. Consider the set $\mathcal{S}(r) = \cup_{x_i \in \mathcal{D}} B(x_i, r)$, where $B(x, r)$ is a ball of radius r centered at x . The probability mass of this set $P(\mathcal{S}(r)) = \int_{x \in \mathcal{S}(r)} p(x) dx$. As $r \rightarrow 0$, the volume $V(\mathcal{S}(r))$ of the set $\mathcal{S}(r)$ goes to zero. We have

$$P(\mathcal{S}(r)) = \int_{x \in \mathcal{S}(r)} p(x) dx \leq V(\mathcal{S}(r)) \cdot u \xrightarrow{r \rightarrow 0} 0. \quad (4)$$

Hence, there exists r_0 such that $P(\mathcal{S}(r_0)) \leq \frac{1}{2}$.

Now, define the neighborhood of the set \mathcal{D}_{OOD} as

$$\mathcal{S}_{OOD} = \cup_{x' \in \mathcal{D}_{OOD}} B(x', \hat{r}), \quad (5)$$

where \hat{r} is selected so that the total volume of set \mathcal{S}_{OOD} is $1/2c$. Then, we can define a new density p' by redistributing the mass in $p(\cdot)$ from outside the set $\mathcal{S}(r_0)$ to the neighborhood \mathcal{S}_{OOD} as follows:

$$p'(x) = \begin{cases} p(x), & \text{if } x \in \mathcal{S}(r_0), \\ 2c \cdot (1 - P(\mathcal{S}(r_0))), & \text{if } x \in \mathcal{S}_{OOD}, \\ 0, & \text{otherwise.} \end{cases} \quad (6)$$

The density $p'(\cdot)$ integrates to one, coincides with p on the training data, and assigns density of at least c to points in \mathcal{D}_{OOD} . \square

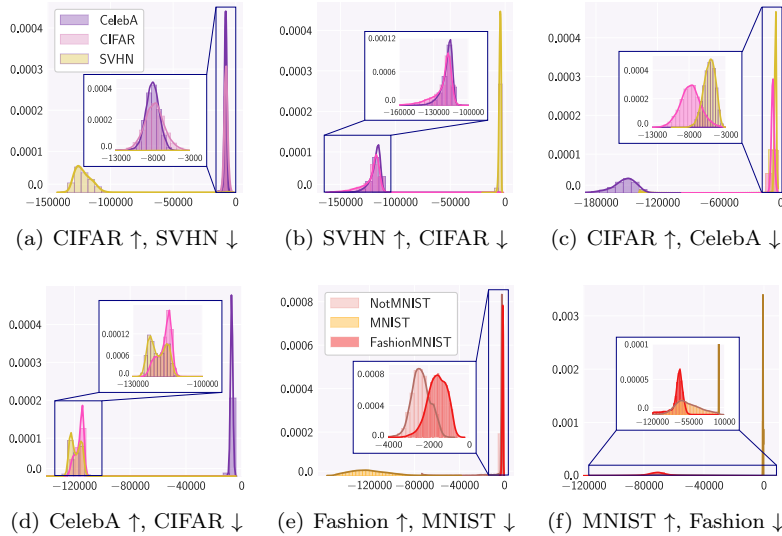


Figure 7: **Negative training.** The histograms of log-likelihood for RealNVP when in training likelihood is maximized on one dataset and minimized on another dataset: (a) maximized on CIFAR, minimized on SVHN; (b) maximized on SVHN, minimized on CIFAR; (c) maximized on CIFAR, minimized on CelebA; (d) maximized on CelebA, minimized on CIFAR. (e) maximized on FashionMNIST, minimized on MNIST; (f) maximized on MNIST, minimized on FashionMNIST;

B Flows have capacity to distinguish datasets

Normalizing flows are unable to detect OOD image data when trained to maximize likelihood on the train set. It is natural to ask whether these models are at all capable of distinguishing different image datasets. In this section we demonstrate the following:

Observation: Flows can assign high likelihood to the train data and low likelihood to a given OOD dataset if they are explicitly trained to do so.

Relevance to OOD detection: While flows have sufficient capacity to distinguish different data, they are biased towards learning solutions that assign high likelihood to all structured data and consequently fail to detect OOD inputs.

We introduce an objective that encouraged the flow to maximize likelihood on the target dataset and to minimize likelihood on a specific OOD dataset. The objective we used is

$$\frac{1}{N_{\mathcal{D}}} \sum_{x \in \mathcal{D}} \log p(x) - \frac{1}{N_{\text{OOD}}} \sum_{x \in \mathcal{D}_{\text{OOD}}} \log p(x) \cdot I[\log p(x) > c], \quad (7)$$

where $I[\cdot]$ is an indicator function and the constant c allows us to encourage the flow to only push the likelihood of OOD data to a threshold rather than decreasing it to $-\infty$; $N_{\mathcal{D}}$ is the number of train datapoints and $N_{\text{OOD}} = \sum_{x \in \mathcal{D}_{\text{OOD}}} I[\log p(x) > c]$ is the number of OOD datapoints that have likelihood above the threshold c .

We trained a RealNVP flow with the objective (7) using different pairs of target and OOD datasets: CIFAR-10 vs CelebA, CIFAR-10 vs SVHN and FashionMNIST vs MNIST. We present the results in Figure 7. In each case, the flow is able to push the likelihood of the OOD dataset to very low values, and simultaneously maximize the likelihood on the target dataset creating a clear separation between the two.

Hyper-parameters For the flow architecture and training used the same hyper-parameters as we did for the baselines, described in Appendix C. For CelebA, CIFAR and SVHN models we set $c = -100000$, and for MNIST, FashionMNIST and NotMNIST we set $c = -30000$.

Connection with prior work Flows can be used as classifiers separating different classes of the same dataset [28, 19, 1], which further highlights the fact that flows can distinguish images based on their contents when trained to do so. A similar experiment for the PixelCNN model [29] was presented in Hendrycks et al. [14]. The authors maximized the likelihood of CIFAR-10 and minimized the likelihood of the TinyImages dataset [42]. In their experiments, this procedure consistently led to CIFAR-10 having higher likelihood than any of the other benchmark datasets. In Figures 7, for each experiment in addition to the two datasets that were used in training we show the log-likelihood distribution on another OOD dataset. For example, when we train the flow to separate CIFAR-10 from CelebA (panels c, d), the flow successfully does so but assigns SVHN with likelihood similar to that of CIFAR. When we train the flow to separate CIFAR-10 from SVHN (panels c, d), the flow successfully does so but assigns CelebA with likelihood similar to that of CIFAR. Similar observations can be made for MNIST, FashionMNIST and notMNIST. At least for normalizing flows, minimizing the likelihood on a single OOD dataset does not lead to all the other OOD datasets achieving low-likelihood.

C Details of the experiments

RealNVP For all RealNVP models, we generally follow the architecture design of Dinh et al. [10]. We use multi-scale architecture where after a block of coupling layers half of the variables are factored out and copied forward directly to the latent representation. Each scale consists of 3 coupling layers with checkerboard mask, followed by a squeeze operation and 3 coupling layers with channel-wise mask (see Figure 9). For the *st*-network we use deep convolutional residual networks with additional skip connections following Dinh et al. [10]. In all experiments, we use Adam optimizer. On grayscale images (MNIST, FashionMNIST), we used 2 scales in RealNVP, 6 blocks in residual *st*-network, learning rate 5×10^{-5} , batch size 32 and trained model for 80 epochs. On CIFAR-10, CelebA and SVHN, we used 3 scales, 8 blocks in *st*-network, learning rate 10^{-4} , batch size 32, weight decay 5×10^{-5} and trained the model for 100 epochs. On ImageNet, we used 5 scales, 2 blocks in *st*-network, learning rate 10^{-3} , batch size 64, weight decay 5×10^{-5} and trained the model for 42 epochs. On CelebA 64×64 , we used 4 scales, 4 blocks in *st*-network, learning rate 10^{-4} , batch size 64, weight decay 5×10^{-5} and trained the model for 100 epochs.

Glow We follow the training details of Nalisnick et al. [27] for multi-scale Glow models. Each scale consists of a sequence of actnorm, invertible 1×1 convolution and coupling layers [22]. The squeeze operation is applied before each scale, and half of the variables are factored out after each scale. In all experiments, we use RMSprop optimizer. On grayscale images (MNIST, FashionMNIST), we used 2 scales with 16 coupling layers, a 3-layer Highway network with 200 hidden units for *st*-network, learning rate 5×10^{-5} , batch size 32 and trained model for 80 epochs. On color images (CIFAR-10, CelebA, SVHN), we used 3 scales with 8 coupling layers, a 3-layer Highway network with 400 hidden units for *st*-network, learning rate 5×10^{-5} , batch size 32 and trained model for 80 epochs.

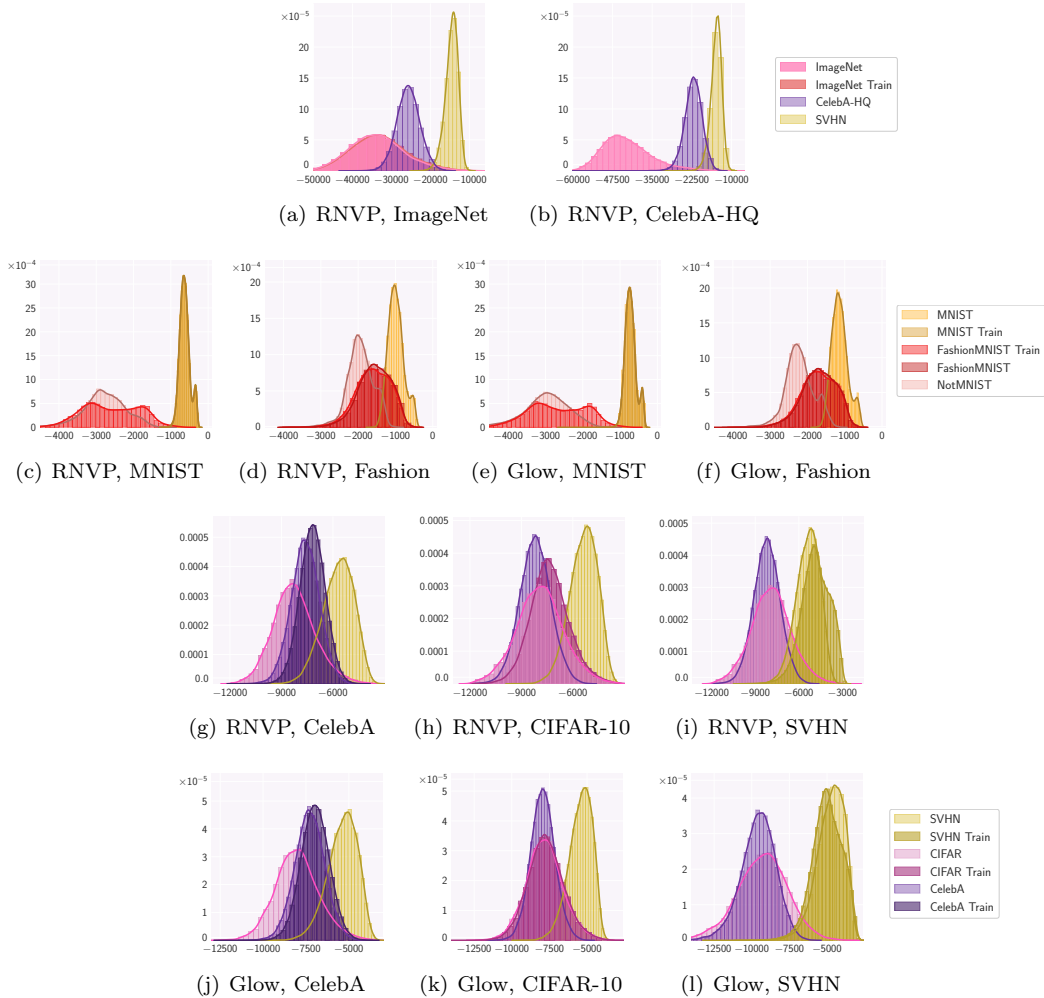


Figure 8: **Baseline log-likelihoods.** The histograms of log-likelihood for RealNVP and Glow models trained on various datasets. Both flows consistently assign similar or higher likelihood to OOD data compared to the target dataset. The likelihood distribution for train and test sets of the target data is typically very similar.

D Baseline models likelihood distributions and AUROC scores

In Figure 8, we plot the histograms of the log likelihoods on in-distribution and out-of-distribution datasets RealNVP and Glow models. In Table 1 we report AUROC scores for OOD detection with these models. As reported in prior work, Glow and RealNVP consistently fail at OOD detection.

E Visualization implementation

Normalizing flows such as RealNVP and Glow consist of a sequence of coupling layers which change the content of the input and squeeze layers (see Figure 9) which reshape it. Due to the presence of squeeze layers, the latent representations of the flow have a different

Model	Train data	OOD data			OOD data			NotMNIST
		CelebA	CIFAR-10	Data	SVHN	MNIST	Fashion	
RealNVP	CelebA	–	67.7	6.3	MNIST	–	99.99	99.99
	CIFAR-10	56.0	–	6.0	Fashion	10.8	–	72.1
	SVHN	99.0	98.4	–				
Glow	CelebA	–	69.1	6.4	MNIST	–	99.96	100.0
	CIFAR-10	52.9	–	5.5	Fashion	13.3	–	80.2
	SVHN	99.9	99.1	–				

Table 1: **Baseline AUROC.** AUROC scores on OOD detection for RealNVP and Glow models trained on various image data. Flows consistently assign higher likelihoods to OOD dataset except when trained on MNIST and SVHN. The AUROC scores for RealNVP and Glow are close.

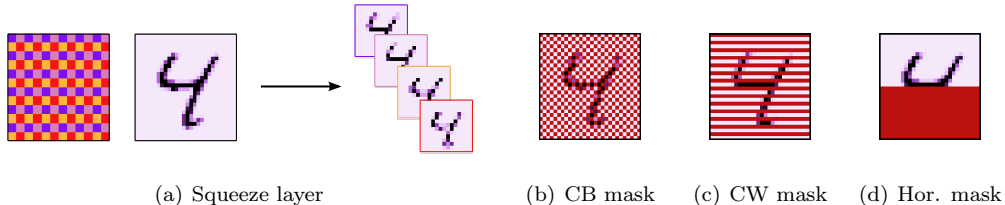
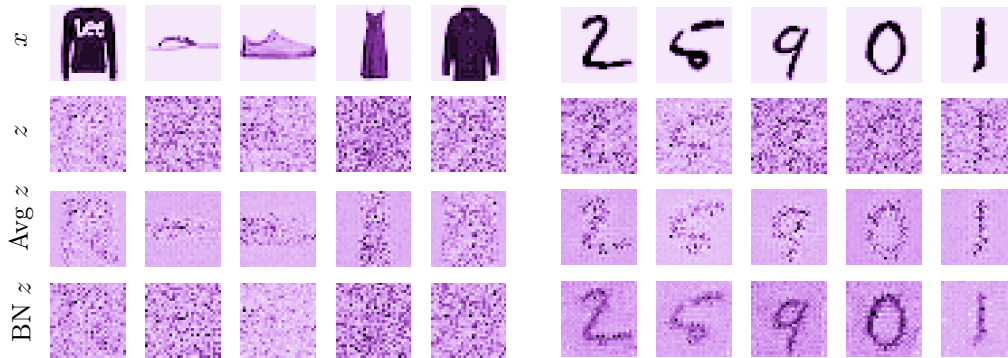


Figure 9: **Squeeze layers and masks.** (a): A squeeze layer squeezes an image of size $c \times h \times w$ into $4c \times h/2 \times w/2$. The first panel shows the mask, where each color corresponds to a channel added by the squeeze layer (for visual clarity we show the mask for a 12×12 image). The second panel shows a $1 \times 28 \times 28$ MNIST digit, and the last panel shows the 4 channels produced by the squeeze layer. The colors of the boundaries of the channel visualizations correspond to the colors of the pixels in the mask. Each channel produced by the squeeze layer is a subsampled version of the input image. (b)-(d): Checkerboard, channel-wise and horizontal masks applied to the same input image. Masked regions are shown in red. Channel-wise mask is obtained by applying a squeeze layer and masking two of the channels (e.g. the last two); here we show the masked pixels in the un-squeezed image. Masks are typically alternated: in the subsequent layers the masked and observed positions are swapped.

shape compared to the input. In order to visualize latent representations, we revert all squeezing operations of the flow and visualize $\text{unsqueeze}(z)$. Similarly, for visualization of coupling layer activations and scale and shift parameters predicted by st -network, we revert all squeezing operations and join all factored out tensors in the case of multi-scale architecture (i.e., we feed the corresponding tensor through inverse sub-flow without applying coupling layers or invertible convolutions).

F Additional latent representation visualizations

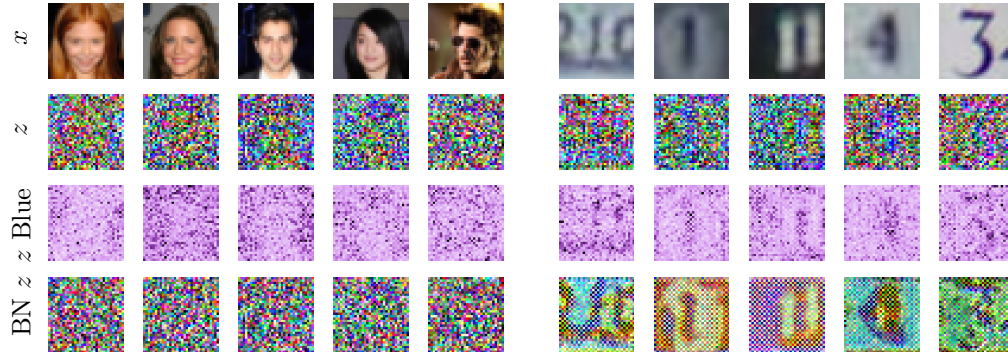
in Figure 10, we plot additional latent representations for RealNVP and Glow trained on FashionMNIST with MNIST as OOD dataset, RealNVP trained on CelebA with SVHN as OOD. The results agree with Section 5: we can recognize edges from the original inputs in their latent representations.



(a) RealNVP trained on FashionMNIST

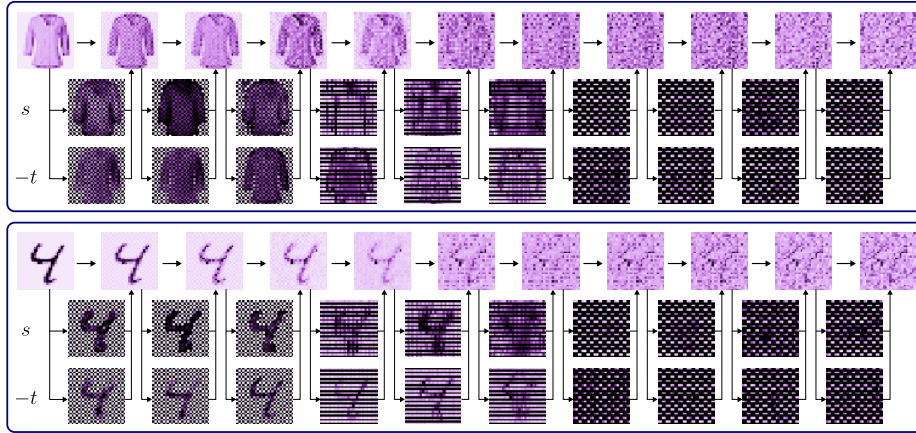


(b) Glow trained on FashionMNIST

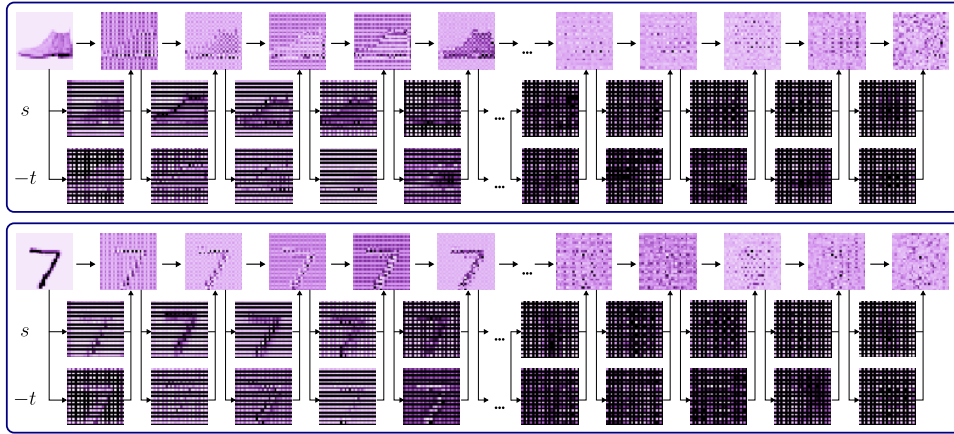


(c) RealNVP trained on CelebA

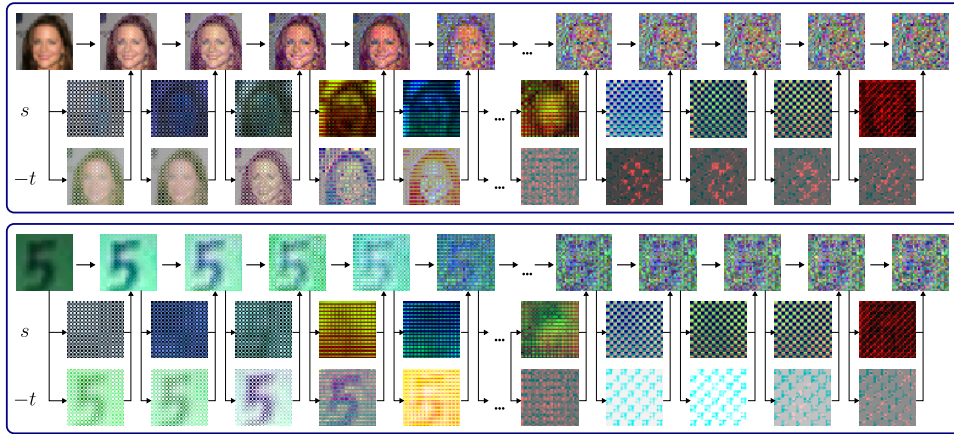
Figure 10: **Latent spaces.** Visualization of latent representations for RealNVP and Glow models on in-distribution and out-of-distribution inputs. **Rows 1-3 in (a) and (b):** original images, latent representations, latent representation averaged over 40 samples of dequantization noise for RealNVP and Glow model trained on FashionMNIST and using MNIST for OOD data. **Row 4 in (a):** latent representations for batch normalization in train mode. **Rows 1-4 in (c):** original images, latent representations, the blue channel of the latent representation, and the latent representations for batch normalization in train mode for a RealNVP model trained on CelebA and using SVHN as OOD data. For both dataset pairs, we can recognize the shape of the input image in the latent representations. The flow represents images based on their graphical appearance rather than semantic content.



(a) RealNVP trained on FashionMNIST



(b) Glow trained on FashionMNIST



(c) RealNVP trained on CelebA

Figure 11: **Coupling layer visualizations.** Visualization of intermediate coupling layer activations and st -network predictions for (a): RealNVP trained on FashionMNIST; (b): Glow trained on FashionMNIST; (c): RealNVP trained on CelebA. The top half of each subfigure shows the visualizations for an in-distribution image (FashionMNIST or CelebA) while the bottom half shows the visualizations for an OOD image (MNIST or SVHN). For all models, the shape of the input both for in- and out-of-distribution image is clearly visible in s and t predictions of the coupling layers.

G Masking strategies

In Figure 9, we visualize checkerboard, channel-wise masks and horizontal masks on a single-channel image. The checkerboard and channel-wise masks are commonly used in RealNVP, Glow and other coupling layer-based flows for image data. We use the horizontal mask to better understand the transformations learned by the coupling layers in Section 6.

H Additional coupling layer visualizations

In Figure 11, we plot additional visualizations of coupling layer activations and scale s and shift t parameters predicted by st -networks. In Figure 12 we visualize the coupling layer activations for the small flow with horizontal mask from Section 6.2 on several additional OOD inputs. These visualizations provide additional empirical support for Section 6.

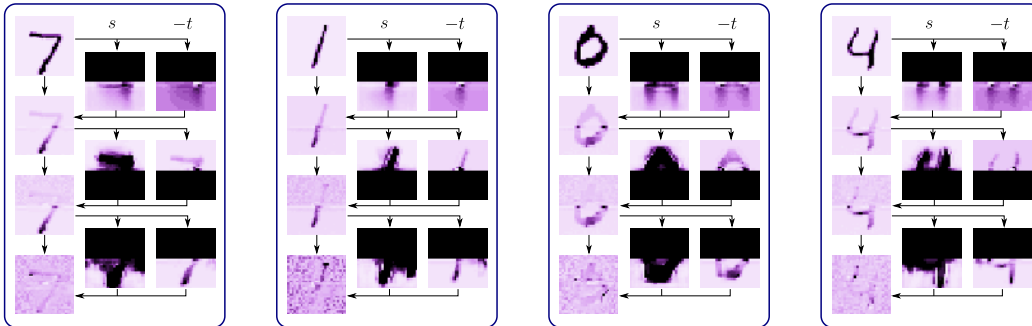


Figure 12: **Coupling layer co-adaptation.** Visualization of intermediate coupling layer activations, as well as scales s and shifts t predicted by each coupling layer of a RealNVP model with a horizontal mask on out-of-distribution MNIST inputs. Although RealNVP was trained on FashionMNIST, the st -networks are able to correctly predict the bottom half of MNIST digits in the second coupling layer due to coupling layer co-adaptation.

I Changing biases in flow models for better OOD detection

I.1 Cycle-mask

In Section 6 we identified two mechanisms through which normalizing flows learn to predict masked pixels from observed pixels on OOD data: leveraging local color correlations and coupling layer co-adaptation. We reduce the applicability of these mechanisms with *cycle-mask*: a new masking strategy for the coupling layers illustrated in Figure 13.

With cycle-mask, the coupling layers do not have access to neighbouring pixels when predicting the masked pixels, similarly to the horizontal mask. Furthermore, cycle mask reduce the effect of coupling layer co-adaptation: the information about a part of the image has to travel through 4 coupling layers before it can be used to update the same part of the image.

Changing masking strategy In Figure 14 we show the log-likelihood histograms and samples for a RealNVP of a fixed size with checkerboard, horizontal and cycle-mask.

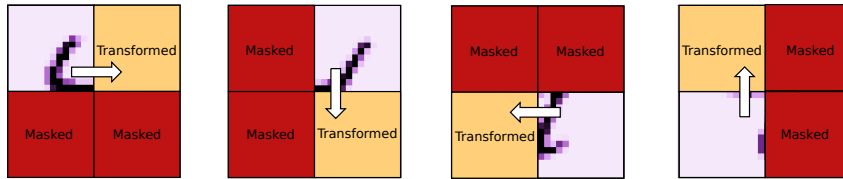


Figure 13: **Cycle-mask**. A new sequence of masks for coupling layers in RealNVP that we evaluate in Section 7. We separate the input image space of size $c \times h \times w$ into four quadrants of size $c \times h/2 \times w/2$ each. Each coupling layer transforms one quadrant based on the previous quadrant. Cycle-mask prevents co-adaptation between subsequent coupling layers discussed in Section 6: the information from a quadrant has to propagate through four coupling layers before reaching the same quadrant.

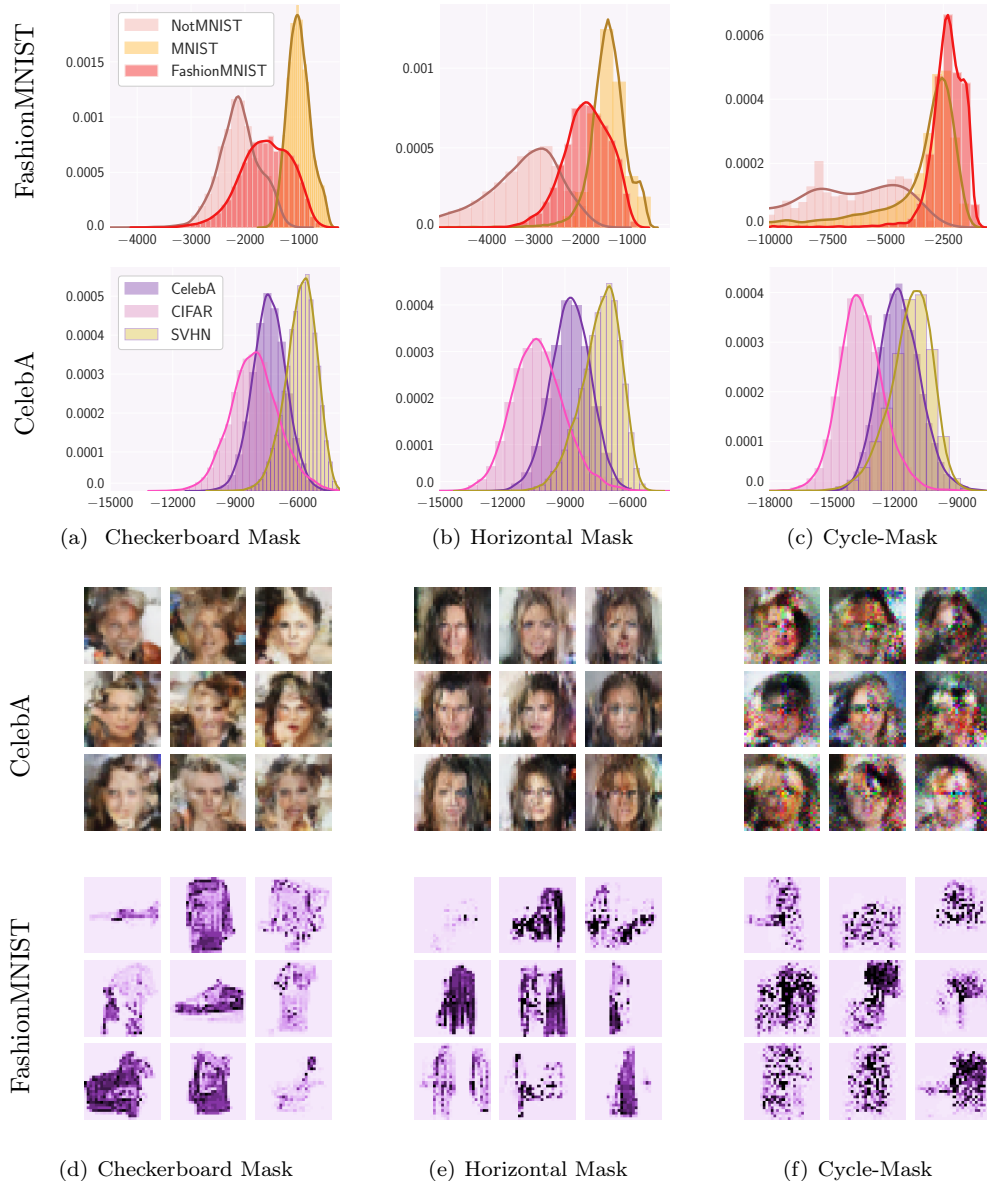


Figure 14: **Effect of masking strategy** The first two rows show log likelihood distribution for RealNVP models trained on FashionMNIST and CelebA with (a) checkerboard mask; (b) horizontal mask; and (c) cycle-mask. The third and the fourth rows show samples produced by the corresponding models.

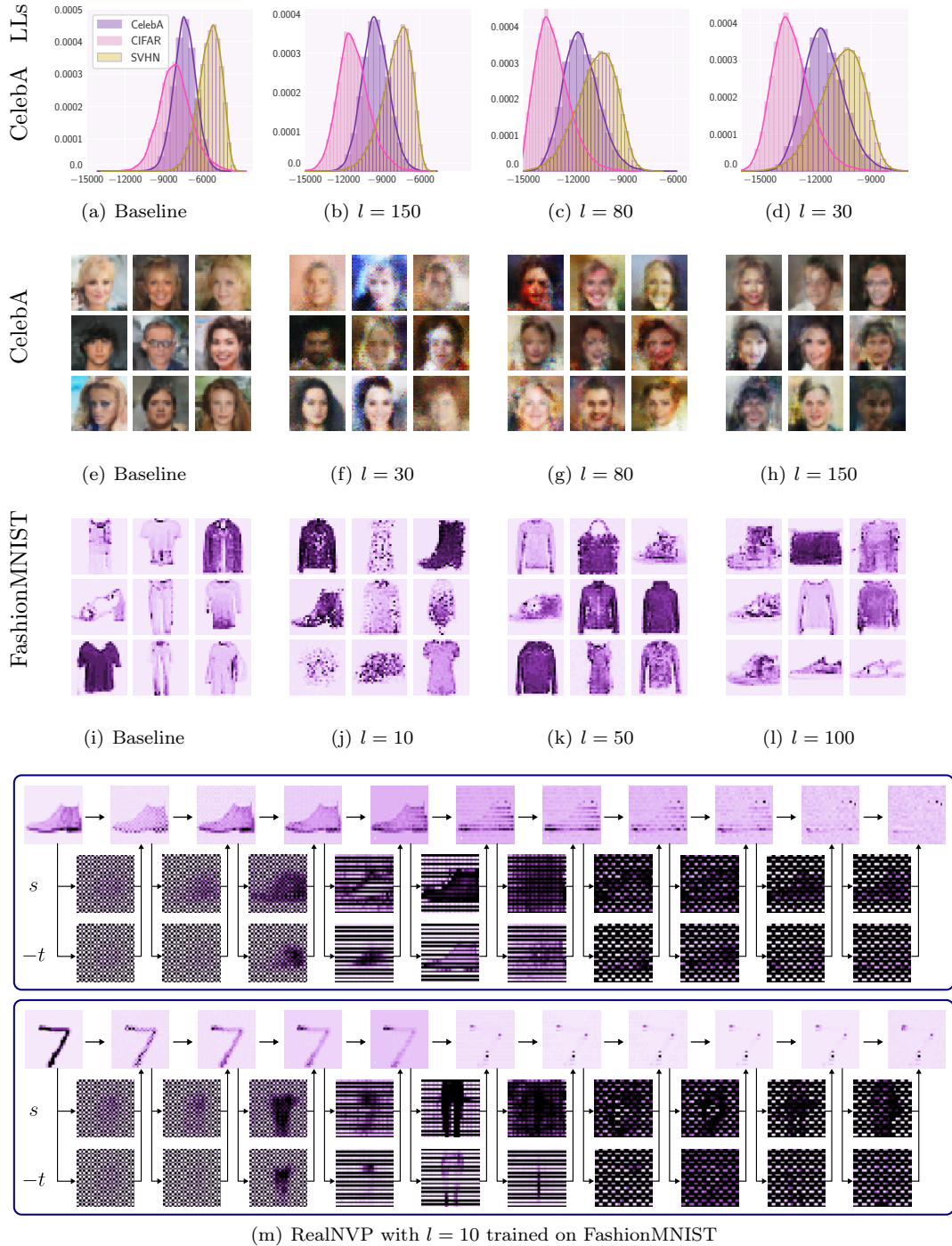


Figure 15: **Effect of st -network capacity.** The first row shows the histogram of log likelihoods for a RealNVP model trained on CelebA dataset: (a) for a baseline model, and (b)-(d) for models with different bottleneck dimensions l in st -network. The second and third rows show samples from RealNVP model trained on CelebA and FashionMNIST respectively: (e) and (i) for baseline models, and (f)-(h) and (j)-(l) for models with different bottleneck dimensions l . In (m), we show the visualization of the coupling layer activations and st -network predictions for a RealNVP model trained on FashionMNIST with a bottleneck of dimension $l = 10$. The top half shows the visualizations for an in-distribution FashionMNIST image while the bottom half shows the visualizations for an OOD MNIST image. st -network with restricted capacity cannot accurately predict masked pixels of the OOD image in the intermediate coupling layers. Moreover, in the middle coupling layers for the MNIST input the activations resemble FashionMNIST images in s and t predictions.

Changing the architecture of *st*-networks In Figure 15, we show likelihood distributions, samples and coupling layer visualization for RealNVP model with *st*-network with a bottleneck trained on FashionMNIST and CelebA datasets. The considered bottleneck dimensions for FashionMNIST are $\{10, 50, 100\}$, and for CelebA the dimensions are $\{30, 80, 150\}$. In the baseline RealNVP model, we use a standard deep convolutional residual network without additional skip connections from the intermediate layers to the output which were used in Dinh et al. [10].

J Samples

In Figure 17, we show samples for RealNVP and Glow models trained on CelebA, CIFAR-10, SVHN, FashionMNIST and MNIST, and a RealNVP model trained on ImageNet 64×64 and CelebA 64×64 .

J.1 Latent variable resampling

To further understand the structure of the latent representations learned by the flow, we study the effect of resampling part of the latent representations corresponding to images from different datasets from the base Gaussian distribution. In Figure 16, using a RealNVP model trained on CelebA we compute the latent representations corresponding to input images from CelebA, SVHN, and CIFAR-10 datasets, and randomly re-sample the subset of latent variables corresponding to a 10×10 square in the center of the image (to find the corresponding latent variables we apply the squeeze layers from the flow to the 32×32 mask). We then invert the flow and compute the reconstructed images from the altered latent representations.

Both for in-distribution and out-of-distribution data, the model almost ideally preserves the part of the image other than the center, confirming the alignment between the latent space and the original input space discussed in Section 5. The model adds a face to the resampled part of the image, preserving the consistency with the background to some extent.

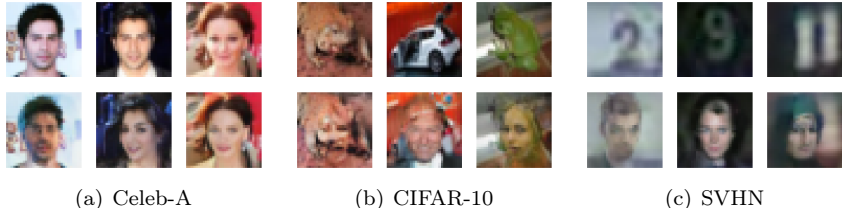


Figure 16: **Latent variable resampling.** Original images (**top row**) and reconstructions with the latent variables corresponding to a 10×10 square in the center of the image randomly re-sampled for a RealNVP model trained on Celeb-A (**bottom row**). The model adds faces (as it was trained Celeb-A) to the part of the image that is being re-sampled.

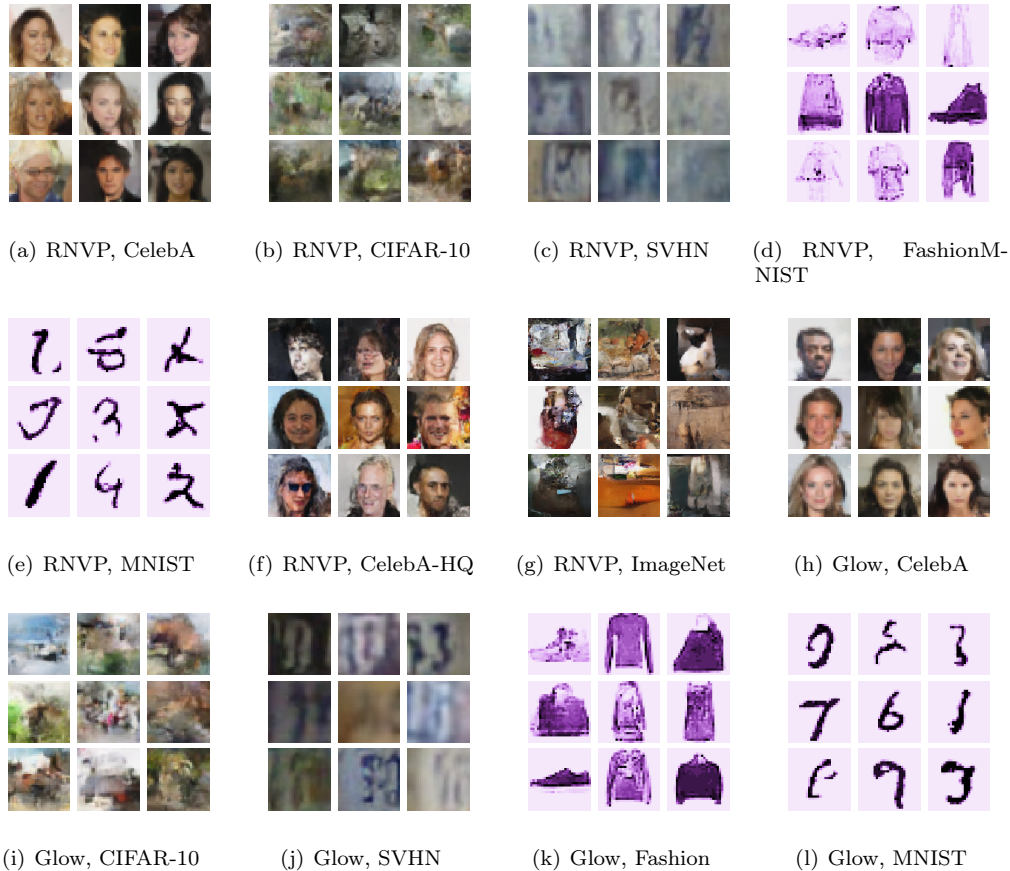


Figure 17: **Baseline Samples.** Samples from baseline RealNVP and Glow models. For ImageNet and CelebA-HQ we used datasets with (64×64) definition.

K Out-of-distribution detection on tabular data

(a) Image embeddings				(b) Tabular data		
Train data	OOD data			Train class (OOD class)	Dataset	
	CelebA	CIFAR-10	SVHN		HEPMASS	MINIBOONE
CelebA	–	99.99	99.99	Background (Signal)	83.78	72.71
CIFAR-10	99.99	–	73.31	Signal (Background)	70.73	87.56
SVHN	100.0	99.98	–			

Table 2: **Image embedding and UCI AUROC.** (a): AUROC scores on OOD detection for RealNVP model trained on image embeddings extracted from EfficientNet. The model is trained on one of the embedding datasets while the remaining two are considered OOD. The models consistently assign higher likelihood to in-distribution data, and in particular AUROC scores are significantly better compared to flows trained on the original images (see Table 1). (b): AUROC scores on OOD detection for RealNVP trained on one class of Hepmass and Miniboone datasets while the other class is treated as OOD data.

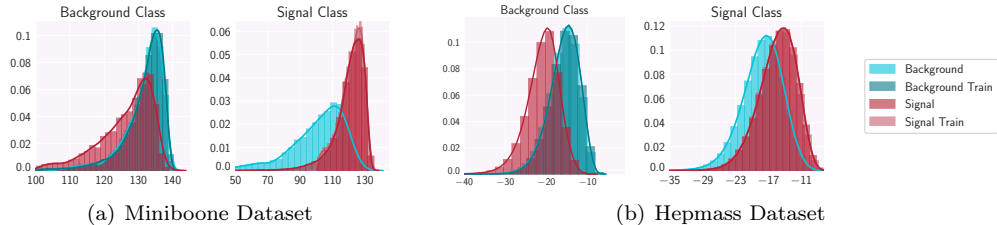


Figure 18: **UCI datasets.** The histograms of log-likelihood for RealNVP on Hepmass and Miniboone tabular datasets when trained on one class and the other class is viewed as OOD. The train and test likelihood distributions are almost identical when trained on either class, and the OOD class receives lower likelihoods on average. There is however a significant overlap between the likelihoods for in- and out-of-distribution data.

K.1 Model

We use RealNVP with 8 coupling layers, fully-connected *st*-network and masks which split input vector by half in an alternating manner. For UCI experiments, we use 1 hidden layer and 256 hidden units in *st*-networks, learning rate 10^{-4} , batch size 32 and train the model for 100 epochs. For image embeddings experiments, we use 3 hidden layer and 512 hidden units in *st*-networks, learning rate 10^{-3} , batch size 1024 and train the model for 120 epochs. For all experiments, we use the AdamW optimizer [24] and weight decay 10^{-3} .

K.2 EfficientNet embeddings

We train RealNVP model on image embeddings for CIFAR-10, CelebA and SVHN extracted from EfficientNet train on ImageNet, and report AUROC scores in Table 2(a).

K.3 UCI datasets

In this experiment, we use 2 UCI classification datasets which were used for unsupervised modeling in prior works on normalizing flows [30, 11, 13]: HEPMASS [2] and MINIBOONE [34]. HEPMASS and MINIBOONE are both binary classification datasets originating from physics, and the two classes represent *background* and *signal*. We follow data preprocessing steps of Papamakarios et al. [30]. We filter features which have too many reoccurring values, after that the dimensionality of the data is 15 for HEPMASS and 50 for MINIBOONE. For HEPMASS, we use the “1000” dataset which contains subset of particle signal with mass 1000. For MINIBOONE data, for each class we take a random split of 10% for a test set.

To test OOD detection performance, for each dataset we train a model on one class while treating the second class as OOD data. We plot the resulting train, test and OOD likelihood distributions for each dataset in Figure 18. We also report AUROC scores for each setup in Table 2(b). While test and OOD likelihoods overlap, the in-distribution class has higher average likelihood in all cases, and AUROC values are ranging between 70% and 87% which is a significantly better result compared to the results for image benchmarks reported in Nalisnick et al. [27].