

Kullback-Leibler Divergence-Based Out-of-Distribution Detection with Flow-Based Generative Models

Yufeng Zhang, Jialu Pan, Wanwei Liu, Zhenbang Chen, Kenli Li, Ji Wang, Zhiming Liu, Hongmei Wei

Abstract—Recent research has revealed that deep generative models including flow-based models and Variational Autoencoders may assign higher likelihoods to out-of-distribution (OOD) data than in-distribution (ID) data. However, we cannot sample OOD data from the model. This counterintuitive phenomenon has not been satisfactorily explained and brings obstacles to OOD detection with flow-based models. In this paper, we prove theorems to investigate the Kullback-Leibler divergence in flow-based model and give two explanations for the above phenomenon. Based on our theoretical analysis, we propose a new method KLODS to leverage KL divergence and local pixel dependence of representations to perform anomaly detection. Experimental results on prevalent benchmarks demonstrate the effectiveness and robustness of our method. For group anomaly detection, our method achieves 98.1% AUROC on average with a small batch size of 5. On the contrary, the baseline typicality test-based method only achieves 64.6% AUROC on average due to its failure on challenging problems. Our method also outperforms the state-of-the-art method by 9.1% AUROC. For point-wise anomaly detection, our method achieves 90.7% AUROC on average and outperforms the baseline by 5.2% AUROC. Besides, our method has the least notable failures and is the most robust one.

Index Terms—Out-of-distribution detection, deep learning, flow-based model, Kullback-Leibler divergence, Gaussian distribution.

1 INTRODUCTION

ANOMALY detection is the process of “finding patterns in data that do not conform to expected behavior” [1], [2]. Under an unsupervised learning setting, the model is trained on a set of unlabeled data $\{x_1, \dots, x_n\}$ which are drawn independently from an unknown distribution p^* . Group anomaly detection (GAD) [3] aims to determine whether a given group of test inputs $\{x_1, \dots, \tilde{x}_m\} (m > 1)$ is sampled from p^* . Typical applications of GAD include discovering high-energy particle physics, [4], anomalous galaxy clusters in astronomy [5], [6], unusual vorticity in fluid dynamics [7], and stealthy attacks [3], [8]. Point-wise anomaly detection (PAD) [1], [9] aims to determine whether an individual input is sampled from p^* . PAD is applied in many areas including detecting intrusion [1], fraud [10], malware [11], and medical anomalies [1]. It is worth noting that GAD cannot be implemented by PAD because the individual members of the input group may not be anomalies [2], [3], [12]. In literature, the term *anomaly* is also referred to as outlier, peculiarity, out-of-distribution (OOD) data, etc. In the following, we mainly use terms *OOD data* and *anomaly* as in

most related works.

Counterintuitive Phenomenon. This paper focuses on unsupervised OOD detection using explicit deep generative models (DGM) including flow-based models and Variational Autoencoders (VAE). Recent research shows that explicit deep generative models including flow-based models [13], [14], VAE [15], and auto-regressive models [16], [17] are not capable of distinguishing OOD data from in-distribution (ID) data (training data) according to the model likelihood (*i.e.*, Type II errors) [18], [19], [20], [21], [22], [23]. For example, as shown in Figures K.1(a) and K.1(b) in the supplementary material, Glow [13] assigns higher likelihoods for SVHN (MNIST) when trained on CIFAR-10 (FashionMNIST). Figure K.2 in the supplementary material shows similar results in recent proposed residual flows [24]. However, as pointed out by Nalisnick *et al.* [22] *we cannot sample OOD data from the model*. We can also observe a similar phenomenon in class conditional Glow (GlowGMM), which contains a Gaussian mixture model on the top layer with one Gaussian distribution for each class [13], [25], [26]. For example, GlowGMM does not achieve the same performance as prevalent discriminative models such as ResNet [27] on FashionMNIST. We observe that the centroids of different components are close to each other (see Figure K.3 in the supplementary material). One component may assign higher likelihoods for other classes (see Table J.11 in the supplementary material). However, *we always sample images of the correct class from the corresponding component*.

Nalisnick *et al.* explain the above phenomenon by the discrepancy of the typical set and high probability density regions of the model distribution [22]. They propose using typicality test to detect OOD data. However, their explanation and method fail on problems where the likelihoods of

- Yufeng Zhang, Jialu Pan, and Kenli Li are with the College of Computer Science and Electronic Engineering, Hunan University, Changsha, China. E-mail: yufengzhang@hnu.edu.cn, jialupan@hnu.edu.cn, lkl@hnu.edu.cn
- Wanwei Liu, Zhenbang Chen, and Ji Wang are with the College of Computer, National University of Defense Technology, Changsha, China. E-mail: wwliu@nudt.edu.cn, zbchen@nudt.edu.cn, wj@nudt.edu.cn
- Zhiming Liu is with the Centre for Research and Innovation in Software Engineering, Southwest University, Chongqing, China. E-mail: zhimingliu88@swu.edu.cn
- Hongmei Wei is with the National Research Center of Parallel Computer Engineering and Technology, China. E-mail: wei_hongmei@163.com
- Kenli Li and Ji Wang are the corresponding authors.

Manuscript received xx xx, 2022; revised xx xx, 2022.

ID and OOD data coincide (*e.g.*, CIFAR-10 vs CIFAR-100, CelebA vs CIFARs). In this paper, we manipulate the model likelihoods such that ID and OOD data have coinciding likelihoods (see Subsection 3.1). Such manipulation could make all existing likelihood-based OOD detection methods [22], [28], [29], [30] fail. Some researchers investigate the behaviors of flow-based models in OOD detection. Kirichenko *et al.* reveal that flow-based model learns local pixel correlations and generic image-to-latent-space transformations [23]. Such learned knowledge may also exist in OOD dataset. Zhang *et al.* state that the estimation error of the flow-based model is the reason for the failure of anomaly detection [31].

Research Questions. Currently, the above counterintuitive phenomenon has not been explained satisfactorily. In this paper, we rethink the existing conclusions relating to OOD detection using flow-based model. We focus on the following two research questions:

- **Q1: Explanation**¹. Why can we not sample OOD data from flow-based model? We need a unified answer to this question whenever OOD data have lower, higher, or coinciding likelihoods.
- **Q2: OOD detection.** How to detect OOD data using flow-based model and VAE without supervision?

We start our research from the sampling process. Flow-based model constructs diffeomorphism $z = f(x)$ from visible (data) space to latent space. The model maps each input data point x to a unique representation z in latent space. We can sample noise ε from prior (usually standard Gaussian distribution) and generate new data $f^{-1}(\varepsilon)$. So we should ask *why we cannot sample the representations of OOD data from prior*. In this paper, we explain why we cannot sample OOD data. We abandon the model likelihood and leverage Kullback-Leibler (KL) divergence and local pixel dependence of representations for OOD detection.

Contributions. The contributions of this paper are:

- 1) We prove several theorems to investigate the KL divergence in flow-based model. We answer why we cannot sample OOD data from two perspectives. The first answer reveals the large KL divergence between the distribution of representations of OOD data and the prior. The second answer states that the representations of OOD data locate in specific directions.
- 2) We propose a unified OOD detection method in three steps based on our analysis. Firstly, we propose leveraging the KL divergence between the distribution of representations and prior for GAD. We also propose using fitted Gaussian to estimate the (lower bound of) KL divergence. Secondly, we decompose the KL divergence and leverage the last-scale KL divergence for OOD detection. Finally, we leverage the local pixel dependence of representations to improve our method further and support PAD.
- 3) We conduct experiments to demonstrate the effectiveness and robustness of our method.

The remaining part of this paper is organized as follows. Section 2 discusses the related work. Section 3 discusses problem settings. Section 4 presents our theoretical analysis to answer Q1. Section 5 elaborates on the details of our

OOD detection method. Section 6 presents experimental results. Finally, Section 7 concludes. More details of the methods, experimental results, discussion, and related work are presented in the supplementary material.

2 RELATED WORK

We discuss the most related work here. More discussion is presented in Section I in the supplementary material.

GAD and PAD. In [3], Toth *et al.* give a survey on GAD methods and plenty of real-world GAD applications. In [9], Chalapathy *et al.* survey a wide range of deep learning-based GAD and PAD methods. In [2], Pang *et al.* review the deep learning-based anomaly detection methods. It is worth noting that in GAD an individual data point in the input group can be normal [2], [3], [12]. So GAD and PAD have different contexts. According to the availability of supervision information, OOD detection can be classified into supervised, semi-supervised, and unsupervised settings. In this paper, we focus on unsupervised OOD detection using flow-based model, so we mainly compare with methods in the same category.

OOD Detection Using Flow-Based Model. Generally, it seems straightforward to use model likelihood $p(x)$ (if any) of a generative model to detect OOD data [3], [32]. However, these methods fail when OOD data have higher or similar likelihoods. Choi *et al.* propose using the Watanabe-Akaike Information Criterion (WAIC) to detect OOD data [20]. WAIC penalizes points that are sensitive to the particular choice of posterior model parameters. However, Nalisnick *et al.* [22] could not reproduce the results of WAIC. Choi *et al.* also propose using typicality test in the latent space to detect OOD data. Our results reported in Subsection 3.1 demonstrate that typicality test in the latent space can be attacked. Sabeti *et al.* propose detecting anomalies based on typicality [33], but their method is not suitable for deep generative models. Nalisnick *et al.* propose using typicality test on model distribution (Ty-test) for GAD [22]. Jiang *et al.* propose GOD2KS which combines random projection and two-sample KS test to perform GAD based on flow-based model [34]. Ren *et al.* propose to use likelihood ratios for OOD detection [35]. Serrà *et al.* propose using likelihood compensated by input complexity for OOD detection [28]. In [29], Schirrmeister *et al.* find the likelihood contributed by the last scale of Glow (L_{last}) is a better criterion than $\log p(x)$ for PAD. We find L_{last} should not be explained as likelihood consistently for OOD data. See Section I in the supplementary material for more discussion. In [30], Morningstar *et al.* train density estimator (DoSE) and one-class SVM on the statistics of deep generative models to detect OOD data. Before this writing, GOD2KS [34] and DoSE [30] are the SOTA GAD and PAD methods applicable to flow-based models under unsupervised setting, respectively. We will show that many baseline methods could degenerate into being not better than random guessing under data manipulation. These results demonstrate the difficulty of OOD detection using flow-based model.

3 PROBLEM SETTINGS

This paper mainly focuses on flow-based generative model, which constructs diffeomorphism $z = f(x)$ from visible space \mathcal{X} to latent space \mathcal{Z} [13], [14], [36], [37]. Our work also

1. We focus on the reason behind Q1 rather than aiming to sample OOD data in this paper.

involves Variational Autoencoder (VAE) [15]. Please refer to Section A in the supplementary material for background. In this section, we first discuss how to manipulate the model likelihoods. Then we note the target problems of this work.

3.1 Manipulating Likelihoods

In [22], Nalisnick *et al.* conjecture that the counterintuitive phenomena in Q1 stem from the distinction of high probability density regions and the typical set of the model distribution [20], [22]. For example, Figure K.4 in the supplementary material shows the typical set of d -dimensional standard Gaussian distribution, which is an annulus with a radius of \sqrt{d} [38]. When sampling from the Gaussian distribution, it is highly likely to get points in the typical set rather than the highest density region (*i.e.* the center) or the lowest density region far from the mean. Based on this explanation, Nalisnick *et al.* propose using typicality test (Ty-test in short) to detect OOD data [22]. However, their explanation and method do not apply to problems where OOD data reside in the typical set of model distribution (*i.e.*, OOD data has coinciding likelihoods with ID data). Researchers have also proposed other likelihood-related OOD detection methods, including input complexity compensated likelihood [28], likelihood contributed by the last scale [29], and DoSE [30]. In the following, we show how to manipulate OOD data to make the likelihood of ID and OOD dataset coincide. Such manipulation could make all existing likelihood-based methods fail.

M1: Manipulating $p(z)$ by Rescaling z to Typical Set of Prior. We train Glow with 768-dimensional standard Gaussian prior on FashionMNIST. Figure K.1(c) in the supplementary material shows the histogram of log-likelihood of representations under prior (*i.e.*, $\log p(z)$)². Note that $\log p(z)$ of FashionMNIST is around $-768 \times (0.5 \times \ln 2\pi e) \approx -1089.74$, which is the log-probability of typical set of the prior [39]. Here it seems that we can detect OOD data by $p(z)$ or typicality test in the latent space [20]. However, as shown in Figure K.4 in the supplementary material, we can decode each OOD data point x as $z = f(x)$ and rescale z to the typical set by setting $z' = \sqrt{d} \times z/|z|$ ($d = 768$). Then we decode z' to generate image $x' = f^{-1}(z')$. We find that x' corresponds to the similar image with x . Figure K.5 in the supplementary material shows some examples of x' . These results demonstrate that flow-based model cannot expel representations of OOD data from the typical set of the prior. Note that, Glow model uses multi-scale architecture and has three stages of representations with different scales. In our experiments, rescaling the last scale yields similar results as rescaling all scales simultaneously (see Figure K.5 in the supplementary material). To the best of our knowledge, we are the first to discover that the latents rescaled to the typical set of prior still can be mapped back to legal images. In this paper, we will see that, such manipulation can make multiple existing OOD detection methods fail.

M2: Manipulating $p(x)$ by Adjusting Contrast. Nalisnick *et al.* find that the likelihoods can be manipulated by adjusting the variance of inputs [18]. As shown in Figure K.1(d) in the supplementary material, SVHN with

increased contrast by a factor of 2.0 has coinciding likelihood distribution with CIFAR-10 on Glow trained on CIFAR-10. So it is impossible to detect OOD data by $p(x)$ or typicality test on the model distribution (see Figure K.1(b) in the supplementary material too). In our experiments, we can manipulate the likelihoods of OOD dataset in this way for almost all problems (see Figure K.6~K.10 in the supplementary material). We will see that (in Section 6) multiple existing OOD detection methods could degenerate into being not better than random guessing under data manipulation. Similarly, in VAE, we can also manipulate the likelihoods by adjusting the contrast of input images.

Summary. We can manipulate both $p(x)$ and $p(z)$ of OOD data without knowing the model parameters. In this paper, we abandon the model likelihood and propose an OOD detection method that is robust to data manipulations.

3.2 Problems

We use ID vs OOD to represent an OOD detection problem and use “ID (OOD) representations” to denote the representations of ID (OOD) data. According to the statistics of OOD dataset, we group OOD detection problems into two categories:

- **Category I: smaller/similar variance, higher/similar likelihoods.** OOD dataset has smaller or similar variance with ID dataset and tends to have higher or similar likelihoods;
- **Category II: larger variance, lower likelihoods.** OOD dataset has larger variance than ID data and tends to have lower likelihoods.

As shown in Subsection 3.1, we can use data manipulation **M2** (adjusting contrast) to convert one problem from one category to another.

4 EXPLAINING WHY CANNOT SAMPLE OOD DATA

In this section, we explain why we cannot sample OOD data from two perspectives. Based on these analyses, we will derive our OOD detection method in Section 5.

Figure 1 shows the overview of our analysis of the KL divergence in flow-based model for a certain case (discussed in Subsection 4.1.2). The top half of Figure E.1 in the supplementary material also summarizes our discussion in this section. Please refer to Figure 1 and Figure E.1 in the supplementary material when reading this section.

4.1 Explanation 1: Divergence Perspective

Our analysis involves the following distributions: the distributions of ID data (p_X) and OOD data (q_X), the distributions of ID representations (p_Z) and OOD representations (q_Z), the prior p_Z^r , and the model induced distribution p_X^r such that $Z_r \sim p_Z^r$ and $X_r = f^{-1}(Z_r) \sim p_X^r$. Table C.1 in the supplementary material summarizes the notations involved in our analysis and how they influence each other. In this subsection, we first discuss the general case. Then we conduct further analysis for *Category I* problems (smaller/similar variance, higher likelihoods).

4.1.1 General case

We can analyze the KL divergence in flow-based model in the following steps.

2. In official Glow model, $\log p(z)$ is implemented as the log-likelihood of the representation of the last scale of Glow under prior.

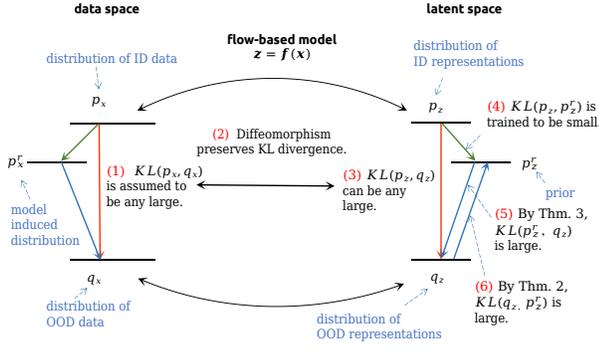


Fig. 1: The key steps of our analysis for Gaussian case (Subsection 4.1.2). Arrows represent KL divergences.

- (1) We treat ID and OOD datasets as samples from different unknown distributions. Therefore, it is reasonable to consider the following assumption.

Assumption 1 *The KL divergence between the distributions of ID and OOD datasets is large.*

So we can assume both $KL(p_X||q_X)$ and $KL(q_X||p_X)$ can be any large.

- (2) According to the following Theorem 1, we know diffeomorphism preserves KL divergence.

Theorem 1 (See [40]) *Given a diffeomorphism $z = f(x)$, let $X_1 \sim p_X$, $X_2 \sim q_X$, $Z_1 = f(X_1) \sim p_Z$ and $Z_2 = f(X_2) \sim q_Z$. Let D_ϕ^h be a (h, ϕ) -divergence measure,*

$$D_\phi^h(p_X, q_X) = D_\phi^h(p_Z, q_Z)$$

Proof *KL divergence is a member of the (h, ϕ) -divergence family (See Section B in the supplementary material). The proof of Theorem 1 relies on diffeomorphisms. See [40] for proof.*

Thus, we can know $KL(p_X||q_X) = KL(p_Z||q_Z)$ is large.

- (3) We can suppose the model is expressible enough and trained by maximum likelihood estimation. This is equal to minimizing forward KL divergence $KL(p_X||p_X^r)$ [37]. By Theorem 1, we also have $KL(p_X||p_X^r) = KL(p_Z||p_Z^r)$. Thus, $KL(p_Z||p_Z^r)$ is small.
- (4) KL divergence is not symmetric and does not satisfy the triangle inequality (*i.e.*, not a proper statistical distance)³. Otherwise, we would know that the reverse KL divergence $KL(p_Z^r||p_Z)$ is small and that $KL(q_Z||p_Z^r)$ is large by triangle inequality. Researchers have investigated other statistical divergences in different contexts [41], [42], [43]. However, flow-based model is usually trained by minimizing KL divergence. In order to explain the phenomenon of flow-based model, we should conduct further analysis on KL divergence. In this paper, we seek stronger conclusions for a special case.

We perform generalized Shapiro-Wilk test for multivariate normality [44] on representations. As shown in Table C.3 in the supplementary material, ID representations always have high p -values. This indicates that ID representations always manifest strong normality. Therefore,

3. For example, we can construct two distributions p and q such that $KL(p||q)$ is any small but $KL(q||p)$ is any large.

we can use a Gaussian distribution \mathcal{N}_p to approximate p_Z and have $KL(p_Z||p_Z^r) \approx KL(\mathcal{N}_p||p_Z^r)$. Now we can apply the following Theorem 2 which reveals the approximate symmetry of small KL divergence between Gaussian distributions.

Theorem 2 (Approximate symmetry of small KL divergence between Gaussian distributions) *For any n -dimensional Gaussian distributions $\mathcal{N}(\mu_1, \Sigma_1)$ and $\mathcal{N}(\mu_2, \Sigma_2)$, if $KL(\mathcal{N}(\mu_1, \Sigma_1)||\mathcal{N}(\mu_2, \Sigma_2)) \leq \varepsilon$ ($\varepsilon \geq 0$),*

$$KL(\mathcal{N}(\mu_2, \Sigma_2)||\mathcal{N}(\mu_1, \Sigma_1)) \leq \varepsilon + 2\varepsilon^{1.5} + O(\varepsilon^2) \quad (1)$$

Proof *The proof is too long. See our manuscript [45] for details. Importantly, the supremum is independent of the dimension n . So Theorem 2 is applicable to high-dimensional problems (e.g., flow-based model).*

By Theorem 2, we can know the reverse KL divergence $KL(p_Z^r||\mathcal{N}_p) \approx KL(p_Z^r||p_Z)$ must be small too. Thus, we can consider the following assumption.

Assumption 2 *The distribution of ID representations and the prior are close enough.*

- (5) Now that the forward and reverse KL divergence between p_Z and prior p_Z^r are both small, we can consider a stronger assumption $p_Z \approx p_Z^r$. Thus, we have $KL(q_Z||p_Z^r) \approx KL(q_Z||p_Z)$. In step (1), we have known $KL(q_X||p_X) = KL(q_Z||p_Z)$ is large, so $KL(q_Z||p_Z^r)$ is large too.

4.1.2 The Gaussian case

In the above Step (4), we use a strong assumption $p_Z \approx p_Z^r$. In fact, for *Category I* problems (smaller/similar variance, higher/similar likelihoods), we do not need such assumption. The results of normality test on OOD representations demonstrate OOD representations in all *Category I* problems except for SVHN vs Constant have p -values greater than 0.05 (see Table C.3 in supplementary material). It seems that OOD datasets “sitting inside” the training data are also “Gaussianized” along with the training data. As far as we know, we are the first to observe this phenomenon.

Based on this observation, we can conduct more analysis using the following Theorem 3, which reveals that KL divergence between Gaussian distributions follows a relaxed triangle inequality.

Theorem 3 (Relaxed triangle inequality) *For any three n -dimensional Gaussian distributions $\mathcal{N}(\mu_i, \Sigma_i)$ ($i \in \{1, 2, 3\}$) such that $KL(\mathcal{N}(\mu_1, \Sigma_1)||\mathcal{N}(\mu_2, \Sigma_2)) \leq \varepsilon_1$ and $KL(\mathcal{N}(\mu_2, \Sigma_2)||\mathcal{N}(\mu_3, \Sigma_3)) \leq \varepsilon_2$ for small $\varepsilon_1, \varepsilon_2 \geq 0$,*

$$KL(\mathcal{N}(\mu_1, \Sigma_1)||\mathcal{N}(\mu_3, \Sigma_3)) < 3\varepsilon_1 + 3\varepsilon_2 + 2\sqrt{\varepsilon_1\varepsilon_2} + o(\varepsilon_1) + o(\varepsilon_2) \quad (2)$$

Proof *The proof is complex and too long. See our work [45] for details. The bound is small for small $\varepsilon_1, \varepsilon_2$ and is 0 when $\varepsilon_1 = \varepsilon_2 = 0$. Similarly, the bound is independent of the dimension n and applicable to high-dimensional problems.*

As shown in Figure 1 and Figure E.1 in the supplementary material, when q_Z is Gaussian-like, we can use a Gaussian distribution \mathcal{N}_q to approximate q_Z and have

$KL(q_Z||p_Z^r) \approx KL(\mathcal{N}_q||p_Z^r)$, $KL(p_Z||q_Z) \approx KL(p_Z||\mathcal{N}_q)$. Now that $KL(p_Z||q_Z)$ is large and $KL(p_Z||p_Z^r)$ is small. According to the relaxed triangle inequality in Theorem 3, $KL(p_Z^r||\mathcal{N}_q)$ must not be small. Furthermore, we can apply Theorem 2 on $KL(p_Z^r||\mathcal{N}_q)$ and know that $KL(\mathcal{N}_q||p_Z^r)$ is large. Finally, we know $KL(q_Z||p_Z^r)$ is large too.

4.1.3 Summary

Overall, we can explain why we cannot sample OOD data from the divergence perspective.

Answer 1 to Q1: The KL divergence between the distribution of OOD representations and prior is large regardless of when the likelihoods of OOD data are higher, lower, or coinciding with that of ID data. So it is hard to sample OOD-like data from the model.

4.2 Explanation 2: Geometric Perspective

We can obtain another explanation from a geometric perspective based on the analysis in the last subsection. The first step is to use the following Theorem 4 to decompose forward KL divergence. Besides, we will use Theorem 4 to derive OOD detection method in Section 5.

Theorem 4 Let $X \sim p_X^*$ be an n -dimensional random vector, $X_i \sim p_{X_i}^*$ be the i -th dimensional element of X . Then

$$KL(p_X^*||\mathcal{N}(0, I_n)) \tag{3}$$

$$= \underbrace{KL(p_X^*||\prod_{i=1}^n p_{X_i}^*(x))}_{I_d[p_X^*]} + \underbrace{\sum_{i=1}^n KL(p_{X_i}^*||\mathcal{N}(0, 1))}_{D_d[p_X^*]=\sum_{i=1}^n D_d^i[p_{X_i}^*]} \tag{4}$$

Proof We can decompose KL divergence as in [46]. See Section D.1 in the supplementary material for proof.

Theorem 4 decomposes forward KL divergence into two non-negative parts: I_d is total correlation (generalized mutual information) measuring the mutual dependence between dimensions [47]; D_d is dimension-wise KL divergence between the marginal distribution of each dimension and prior. We use $[p_X^*]$ to denote one term is computed from p_X^* .

Theorem 4 can help us further investigate the forward KL divergence. For ID data, we have known that $KL(p_Z||p_Z^r)$ is small. Applying Theorem 4 to $KL(p_Z||p_Z^r)$, we can know the total correlation $I_d[p_Z]$ must be small. This indicates that ID data tends to have independent representations. On the contrary, for OOD data, a large $KL(q_Z||p_Z^r)$ allows a large total correlation $I_d[q_Z]$. Although it is hard to estimate total correlation [47], we can use an alternative dependence measure, *i.e.*, the most commonly used correlation coefficient, to investigate the linear dependency. We train Glow on FashionMNIST and test on MNIST/notMNIST. Figure K.11 in the supplementary material shows the histogram of the non-diagonal elements in the correlation matrix of representations. We can see that OOD representations are more correlated. In fact, this happens for all the problems in our experiments. See Figure K.12 to K.17 in the supplementary material for more details.

From a geometric perspective, *a high correlation between dimensions indicates the representations of OOD dataset locate in*

specific directions [48] (see Figure K.18 in the supplementary material for a 3-d example). It is hard to obtain data on specific directions in high dimensional space when sampling from standard Gaussian distribution.

Sampling OOD Data. To verify the above conclusion further, we have tried to restore the information of OOD dataset from the covariance of OOD representations. Ordinarily, after training a flow-based model f , we sample noise $\varepsilon \sim \mathcal{N}(0, I)$ and feed back to the model, we can generate new image $f^{-1}(\varepsilon)$ seeming like training data. Now we feed the model with an OOD dataset and fit a Gaussian distributions $\mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$ from OOD representations, where $\tilde{\mu}$ and $\tilde{\Sigma}$ are the sample mean and covariance of OOD representations, respectively. Then we sample noise $\varepsilon' \sim \mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$ and generate new image $f^{-1}(\varepsilon')$. We find that these generated images are meaningful OOD data. For example, we train Glow on CIFAR-10 and perform the above OOD sampling using notMNIST as OOD dataset (gray-scale images are preprocessed for consistency, see Subsection 6.1). As shown in Figure K.19 in the supplementary material, we can generate images similar to notMNIST, although the images are blurred. In this way, using a single Glow model trained on one training dataset, we can generate images like multiple OOD datasets, including MNIST, notMNIST, SVHN, CelebA, *etc.*, as long as we replace prior with the fitted Gaussian from the representations of the corresponding dataset (See Figure K.20~ K.22 in the supplementary material for details). *These results demonstrate that OOD representations reside in specific directions that can be partially characterized by the mean and covariance of OOD representations.* Such a similar phenomenon is also reported in [49], where Gambardella *et al.* only use the mean of OOD representations. Their manuscript [49] is released contemporaneously with the first edition of this paper.

Furthermore, we scale the norm of OOD representations with different factors. The decoded images also vary from ID data to OOD data gradually. See Figure K.23 in the supplementary material for details. Overall, this leads to the second answer to Q1.

Answer 2 to Q1: OOD representations locate in specific directions with specific norms. The mean and covariance of OOD representations partially characterizes such specific directions. In high dimensional space, it is hard to sample data in specific directions from standard Gaussian distribution (prior) regardless of whether these data reside in the typical set or not.

Note. In the proposed question Q1 “why we cannot sample OOD data from the model”, we mean we cannot generate OOD data when sampling noise ε from prior. In this section, we sample OOD data from flow-based model with fitted Gaussian distribution from OOD representations. This does not contradict the proposed question Q1 because we need the mean and covariance of OOD representations in advance. More research on sampling OOD data is beyond the scope of this paper. We will explore this direction in the future.

5 ANOMALY DETECTION METHOD

In this section, we elaborate on our OOD detection method in three steps in three subsections, respectively. In Subsection 5.1, we propose leveraging KL divergence for OOD detection. In Subsection 5.2, we reduce the computation cost. Finally, in Subsection 5.3, we present a unified OOD method supporting PAD and GAD with small batch sizes. Please refer to Figure 2 and Figure E.1 in the supplementary material for an overview when reading this section.

5.1 Step 1: Leveraging KL divergence

Answer 1 in Subsection 4.1.3 reminds us to detect OOD data by estimating $KL(p||p_Z^r)$, where p is the distribution of representations of inputs. However, when only samples are available, divergence estimation is provable hard, and the estimation error decays slowly in high dimension space [50], [51], [52]. This brings difficulty in applying existing divergence estimation [52], [53], [54], [55], [56] to high dimensional problems with small sample size. Luckily, as shown in Table C.3 in the supplementary material, we observe that both ID data and OOD data of *Category I* problems (smaller/similar variance, higher/similar likelihood) follow a Gaussian-like distribution. This provides us with a facility to estimate the KL divergence for GAD.

5.1.1 Flow-based Model

ID Data. As discussed in Section 4, we can use a Gaussian distribution \mathcal{N}_p to approximate p_Z . Here we use sample expectation $\tilde{\mu}$ and covariance $\tilde{\Sigma}$ of representations to estimate the parameters of \mathcal{N}_p ⁴. Experiments also show that we can generate high-quality images by sampling from \mathcal{N}_p rather than the prior (standard Gaussian distribution). Now we can calculate the KL divergence between two Gaussian distributions $\mathcal{N}(\tilde{\mu}, \tilde{\Sigma})$ and $\mathcal{N}(\mu, \Sigma)$ analytically by

$$KL(\mathcal{N}(\tilde{\mu}, \tilde{\Sigma})||\mathcal{N}(\mu, \Sigma)) \quad (5)$$

$$= \frac{1}{2} \left\{ \log \frac{|\Sigma|}{|\tilde{\Sigma}|} + \text{Tr}(\Sigma^{-1}\tilde{\Sigma}) + (\mu - \tilde{\mu})^T \Sigma^{-1}(\mu - \tilde{\mu}) - n \right\}$$

When the prior ($\mathcal{N}(\mu, \Sigma)$) is standard Gaussian distribution $\mathcal{N}(0, I)$, Equation (5) equals to

$$\frac{1}{2} \left\{ -\log |\tilde{\Sigma}| + \text{Tr}(\tilde{\Sigma}) + \tilde{\mu}^T \tilde{\mu} - n \right\} \quad (6)$$

where generalized variance $|\tilde{\Sigma}|$ and total variation $\text{Tr}(\tilde{\Sigma})$ both measure the dispersion of representations. $KL(\mathcal{N}_p||p_Z^r)$ can be calculated in $O(n^3)$ where n is the dimension.

OOD Data in *Category I* Problems. As discussed in Subsection 4, OOD representations of *Category I* problems (smaller/similar variance, higher/similar likelihood) tend to follow a Gaussian-like distribution. Similar to ID data, we can use fitted Gaussian distribution \mathcal{N}_q to approximate q_Z and estimate $KL(q_Z||p_Z^r)$.

OOD Data in *Category II* Problems. Our normality test results (see Table C.3 in the supplementary material) show that OOD representations in *Category II* problems (larger variance, lower likelihood) do not follow a Gaussian-like distribution. However, we find that Equation (6) performs even better on *Category II* problems. The rationality of using

Equation (6) for *Category II* problems can be explained both intuitively and theoretically.

Intuitively, the first two items of Equation (6) compensate each other. For *Category I* problems (smaller/similar variance, higher/similar likelihood), OOD representations are less dispersed than ID representations and have a larger $-\log |\tilde{\Sigma}|$. For *Category II* problems, OOD representations tend to be more dispersed and have a larger $\text{Tr}(\tilde{\Sigma})$. Besides, we find OOD representations tend to have a larger $\tilde{\mu}^T \tilde{\mu}$ than ID representations. Thus, Equation (6) always produces a larger result for OOD than ID data. Note that the term $\tilde{\mu}^T \tilde{\mu}$ alone cannot achieve high performance in GAD. It can also be manipulated by moving the center of dataset (i.e., adding a vector to the input dataset). We can treat Equation (6) as a more comprehensive statistic than that used in t-test, Maximum Mean Discrepancy, etc.

Theoretically, the following Theorem 5 can explain the rationality of using Equation 6 in *Category II* problems.

Theorem 5 (see [43]) *Let $\mathcal{N}_1(\mu_1, \Sigma_1)$ and $\mathcal{N}_1(\mu_2, \Sigma_2)$ be two n -dimensional Gaussian distributions. Assume that $Z \sim P_Z(z)$ is an arbitrary n -dimensional continuous random variable with mean vector μ_1 and covariance matrix Σ_1 , then*

$$KL(\mathcal{N}_1(\mu_1, \Sigma_1)||\mathcal{N}_1(\mu_2, \Sigma_2)) \leq KL(P_Z(z)||\mathcal{N}_1(\mu_2, \Sigma_2))$$

According to Theorem 5, when we use fitted Gaussian \mathcal{N}_q from OOD representations, $KL(\mathcal{N}_q||p_Z^r)$ is a lower bound of $KL(q_Z||p_Z^r)$. If the lower bound is large, $KL(q_Z||p_Z^r)$ must be large.

Summary. Equation (6) is a unified conservative criterion for GAD due to the following reasons.

- 1) For ID data, Equation (6) approximates $KL(p_Z||p_Z^r)$ and should be small;
- 2) For OOD data whose representations follow a Gaussian-like distribution, Equation (6) approximates $KL(q_Z||p_Z^r)$ and should be large;
- 3) For OOD data whose representations do not follow a Gaussian-like distribution, Equation (6) computes the lower bound of $KL(q_Z||p_Z^r)$. If the lower bound is large, then $KL(q_Z||p_Z^r)$ must be large.

Note that Equation (6) also applies to Gaussian prior with diagonal covariance $\text{diag}(\sigma^2)$ and mean μ . In such a case, we only need to normalize the data by a linear operation $Z' = (Z - \mu)/\sigma$ while keeping $KL(p_Z||\mathcal{N}(\mu, \text{diag}(\sigma^2))) = KL(p_{Z'}||\mathcal{N}(0, I))$ (by Theorem 1). This equals to using Equation (5) directly. We also note that we are not pursuing precise divergence estimation or parameter estimation that are proven to be hard with very small batch sizes in high-dimensional problems.

5.1.2 VAE

It is well-known that VAE and its variations learn independent representations [58], [59], [60], [61], [62]. In VAE, the probabilistic encoder $q_\phi(z|x)$ is often chosen as Gaussian form $\mathcal{N}(\mu(x), \text{diag}(\sigma(x)^2))$, where $z \sim q_\phi(z|x)$ is used as sampled representation, $\mu(x)$ is used as mean representation. The KL term in variational evidence lower bound objective (ELBO, see Equation (16) in the supplementary material) can be rewritten as $E_{p(x)}[KL(q_\phi(z|x)||p(z))] = I(x; z) + KL(q(z)||p(z))$, where $p(z)$ is the prior, $q(z)$ the aggregated posterior, and $I(x; z)$ the mutual information between x and z [63]. Here the term $KL(q(z)||p(z))$ pulls p_Z

4. This is equal to using maximum likelihood estimation [57].

to the Gaussian prior and encourages independent sampled representations. We also investigate the representations in VAE. The results show that:

- 1) ID representations in VAE do not always have p -value greater than 0.05 in Shapiro-Wilk (normality) test;
- 2) the representations of all OOD datasets do not have p -value greater than 0.05 in normality test;
- 3) the representations of OOD datasets are more correlated (see Figure K.25~K.27 in the supplementary material).

Furthermore, there is no theoretical guarantee that $KL(q_Z||p_Z^r)$ is large enough because Theorem 1 does not apply to non-diffeomorphisms. Nevertheless, we find that Equation (6) also works for GAD with VAE.

5.2 Step 2: Leveraging Last-Scale KL Divergence

Although we can use Equation (6) as a preliminary criterion for GAD, it is expensive to compute the sample covariance of representations in $O(n^3)$ when the dimension n reaches several thousand in flow-based model. We propose to use the last scale of representations instead.

Glow model uses multi-scale architecture and has three stages of representations [64]. At the end of the first two stages, outputs are split into two parts \mathbf{h}_i and \mathbf{z}_i ($i = 1, 2$), where \mathbf{h}_i is processed by the next stage. The output of the final stage (*i.e.*, \mathbf{z}_3) contains a quarter of the whole dimensions. Among the three scales, the last scale is the most special one. Interpolating between two representations of the last scale can generate gradually varying images between two real-world images. Schirrmeister *et al.* have shown that Glow network scales manifest a hierarchy of features [29]. Earlier scales learn low-level features that may be generic in different datasets. The last scale learns high-level features that are more specific to the training dataset. The results in [29] also demonstrate that the likelihood contributed by the last scale is a better metric than the whole likelihood for OOD detection. Other work such as [65] also demonstrates the effectiveness of the higher scale. Therefore, the last scale of OOD representations should differ more from ID representations than earlier stages. More precisely, let $q_{Z_1} \sim q_{Z_3}$ be the marginal distribution of the three scales of OOD representations, respectively. We should observe $KL(q_{Z_3}||\mathcal{N}(0, I)) > KL(q_{Z_1}||\mathcal{N}(0, I))$ ($i \in \{1, 2\}$).

Theoretically, similar to Theorem 4, we can decompose the whole KL divergence into local divergence inside each scale and total correlation between different scales as follows.

$$\begin{aligned}
 KL(p_Z(\mathbf{z})||\mathcal{N}) &= \underbrace{KL(p_Z(\mathbf{z})||p_{Z_{1,2}}(\mathbf{z}_1\mathbf{z}_2)p_{Z_3}(\mathbf{z}))}_{\text{total correlation between scales}} \\
 &+ \underbrace{KL(p_{Z_{1,2}}(\mathbf{z}_1\mathbf{z}_2)||\mathcal{N})}_{\text{KL divergence from first two scales}} + \underbrace{KL(p_{Z_3}(\mathbf{z})||\mathcal{N})}_{\text{last-scale KL divergence}} \quad (7)
 \end{aligned}$$

where $\mathbf{z} = \mathbf{z}_1\mathbf{z}_2\mathbf{z}_3$, $\mathbf{z}_1\mathbf{z}_2 \sim p_{Z_{1,2}}$, $\mathbf{z}_3 \sim p_{Z_3}$ and \mathcal{N} is standard Gaussian distribution. Figure 2 shows the decomposition. We call the last item of Equation (7) as *last-scale KL divergence*. The rationality of using last-scale KL divergence as the criterion for OOD detection is based on the following inequality.

$$KL(q_{Z_3}||\mathcal{N}) > KL(p_{Z_3}||\mathcal{N}) \quad (8)$$

where q_{Z_3} and p_{Z_3} are the marginal distributions of the last scale of OOD and ID representations, respectively. Since

the last scale contains fewer dimensions, we can efficiently calculate the last-scale KL divergence. For the non-Gaussian case, we can still rely on Theorem 5 to compute the lower bound.

5.3 Step 3: Leveraging Group-Wise KL divergence in the Last Scale

Up to now, we are still facing two issues. First, when batch size is small (e.g., <5), the performance of last-scale KL divergence is unsatisfactory. Second, the last-scale KL divergence does not support PAD. In this subsection, we address these two issues. The key idea is splitting representation into groups.

The factorizability of standard Gaussian distribution allows us to investigate representations in groups. Intuitively, if $\mathbf{z} \sim \mathcal{N}(0, I)$, then each dimension group of \mathbf{z} follows $\mathcal{N}(0, I)$; Otherwise, it is unlikely that each part of \mathbf{z} follows $\mathcal{N}(0, I)$. Thus, we can split one single \mathbf{z} into multiple subvectors and investigate these subvectors separately. This also generates multiple samples from one data point artificially. Formally, we split random vector Z into k l -dimensional ($k = n/l$) subvectors $\bar{Z}_1, \dots, \bar{Z}_k$. We note the marginal distribution of \bar{Z}_i as $p_{\bar{Z}_i}$ ($1 \leq i \leq k$). Then we can use the following Theorem 6 to further decompose the last-scale KL divergence.

Theorem 6 Let $X \sim p_X^*$ be an n -dimensional random vector. Note $X = \bar{X}_1 \dots \bar{X}_k$ where $\bar{X}_i \sim p_{\bar{X}_i}^*$ be the i -th l -dimensional ($k = n/l$) subvector of X , $\bar{X}_{ij} \sim p_{\bar{X}_{ij}}^*$ is the j -th element of \bar{X}_i . Then,

$$\begin{aligned}
 &KL(p_X^*(\mathbf{x})||\mathcal{N}(0, I_n)) \\
 &= \underbrace{KL(p_X^*(\mathbf{x})||\prod_{i=1}^k p_{\bar{X}_i}^*(\mathbf{x}))}_{I_g[p_X^*]} + \underbrace{\sum_{i=1}^k KL(p_{\bar{X}_i}^*(\mathbf{x})||\mathcal{N}(0, I_l))}_{D_g[p_X^*] = \sum_{i=1}^k D_g^i[p_{\bar{X}_i}^*]} \quad (9) \\
 &= \underbrace{KL(p_X^*(\mathbf{x})||\prod_{i=1}^k p_{\bar{X}_i}^*(\mathbf{x}))}_{I_g[p_X^*]} + \underbrace{\sum_{i=1}^k KL(p_{\bar{X}_i}^*(\mathbf{x})||\prod_{j=1}^l p_{\bar{X}_{ij}}^*(\mathbf{x}))}_{I_l[p_X^*] = \sum_{i=1}^k I_d^i[p_{\bar{X}_i}^*]} \\
 &\quad + \underbrace{\sum_{i=1}^n KL(p_{X_i}^*(\mathbf{z})||\mathcal{N}(0, 1))}_{D_d[p_X^*]} \quad (10)
 \end{aligned}$$

Proof The proof of Theorem 6 is similar to Theorem 4. See Subsection D.2 in the supplementary material for details. \square

In Equation (9), I_g is the generalized mutual information between dimension groups [47]. D_g is *group-wise KL divergence*. Furthermore, in Equation (10) D_g is decomposed as $I_l + D_d$, where I_l is the generalized mutual information inside each group, D_d is dimension-wise KL divergence that also occurs in Equation (3). Combining Equation (3) and 10, we have $I_d = I_g + I_l$ and $D_g = I_l + D_d$. Equation (9) distributes more divergence into the second term than Equation (3). In principle, there are multiple strategies to split Z into k subvectors $\bar{Z}_1, \dots, \bar{Z}_k$. The splitting strategy affects how the whole KL divergence is distributed into I_g and D_g in Equation (9). When $k = n$, Equation (9) is equal to Equation (3).

As shown in Figure 2, we can apply Theorem 6 on p_{Z_3} and q_{Z_3} and get

$$KL(p_{Z_3}||p_Z^r) = I_g[p_{Z_3}] + D_g[p_{Z_3}] = I_g[p_{Z_3}] + \sum_{i=1}^k D_g^i[p_{\bar{Z}_i}]$$

$$KL(q_{Z_3}||p_Z^r) = I_g[q_{Z_3}] + D_g[q_{Z_3}] = I_g[q_{Z_3}] + \sum_{i=1}^k D_g^i[q_{\bar{Z}_i}]$$

where $p_{\bar{Z}_i}$, $q_{\bar{Z}_i}$ are the marginal distributions of subvectors of the last scale of ID and OOD representations, respectively. Combining Equation (8), we can know

$$I_g[q_{Z_3}] + D_g[q_{Z_3}] > I_g[p_{Z_3}] + D_g[p_{Z_3}] \quad (11)$$

Final Criterion. Based on the analysis up to now, we can obtain a final criterion for both GAD and PAD. Figure 2 shows our analysis in this Section. For ID data, $KL(p_Z||\mathcal{N})$ is trained to be small (see Subsection 4.1.1). According to Equation (7), the last-scale KL divergence $KL(p_{Z_3}||\mathcal{N}) = I_g[p_{Z_3}] + D_g[p_{Z_3}]$ must be smaller. We can assume the mutual information between groups $I_g[p_{Z_3}]$ is sufficiently small, *i.e.*, $I_g[p_{Z_3}] < \varepsilon$. To make Equation (11) hold, it suffices that the group-wise KL divergence part satisfies $D_g[q_{Z_3}] > D_g[p_{Z_3}] + \varepsilon$. If we choose an appropriate splitting strategy and distribute more divergence to group-wise KL divergence part ($D_g[q_{Z_3}]$) in Equation (9), it is highly likely that we can make

$$D_g[q_{Z_3}] > D_g[p_{Z_3}] \quad (12)$$

Then we can use group-wise KL divergence of the last scale D_g as the criterion to detect OOD data.

The remaining problems are: (1) how to choose a strategy to split Z into k subvectors so that more divergence is distributed into D_g and (2) how to leverage group-wise KL divergence for OOD detection.

5.3.1 Splitting Strategy: Leveraging Local Pixel Dependence

From Equation (9) and (10), we can know a good splitting strategy should retain enough intragroup dependence in $I_l[q_{Z_3}]$ to make group-wise KL divergence part satisfy $D_g[q_{Z_3}] > D_g[p_{Z_3}] + \varepsilon$.

Take the Glow model for example, the last scale has a shape of $(H \times W \times C)^5$ where H, W, C are the height, width, and the channels, respectively. We can split the last scale into multiple groups. The most natural choices are as follows.

- 1) **horizontal**: treat dimensions in the same pixel position in different channels as one group and split z as $H \times W$ C -dimensional vectors;
- 2) **vertical**: treat dimensions in one channel as one group and split z as C $(H \times W)$ -dimensional vectors.

Here horizontal strategy retains inter-channel dependence into group-wise KL divergence part (*i.e.*, D_g). Vertical strategy retains pixel dependence into D_g .

Figure K.28 in the supplementary material shows the idea behind this subsection. Precisely, we split a single representation z into k subvectors z_1, \dots, z_k and treat z_i as a sample of random vector $\bar{Z}_i \sim p_{\bar{Z}_i}$. Then we can treat z_1, \dots, z_k as k samples of one random vector \bar{Z}_m which follows a mixture of distributions $p_{\bar{Z}_m} = (1/k) \sum_{i=1}^k p_{\bar{Z}_i}$. If the r -th element $\bar{Z}_{i,r}$ and s -th element $\bar{Z}_{i,s}$ are strongly correlated for all $1 \leq i \leq k$, we can say that $\bar{Z}_{m,r}$ and $\bar{Z}_{m,s}$ are also strongly correlated. More generally, if $\bar{Z}_1, \dots, \bar{Z}_k$ have a similar dependence structure, \bar{Z}_m would also have a similar dependence structure. Based on this intuition, we conduct experiments and find that OOD representations manifest local pixel dependence. For example, we test ImageNet32 on Glow trained on SVHN. For each OOD dataset, we visualize

5. The shape of the last scale of the representation in Glow is $4 \times 4 \times 48$.

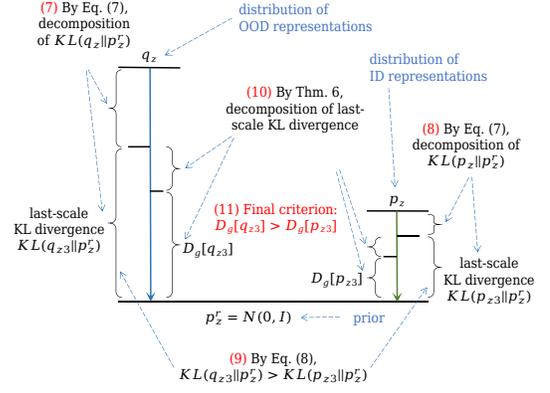


Fig. 2: Decomposition of KL divergence for OOD detection. Steps (1) ~ (6) are shown in Figure 1.

the correlation between pixels. We find that in almost all channels each pixel always has stronger correlation with its neighbors. For example, Figure K.29 in the supplementary material shows the correlation between each pixel with its neighbors in a randomly selected channel. Therefore, we can say that $\bar{Z}_1, \dots, \bar{Z}_k$ tend to have a similar dependence structure. This means that the vertical strategy tends to retain more divergence to D_g . On the contrary, we cannot observe a similar dependence structure between channels when using the horizontal strategy. Thus, the vertical strategy can leverage pixel dependence of representations and is more suitable for OOD detection. Besides, we have also tried other splitting strategies. Evaluation results show that the vertical strategy is the best one.

5.3.2 How to leverage Group-wise KL Divergence in the Last Scale

We want to leverage group-wise KL divergence D_g for OOD detection. For ID data, we treat each representation as k data points sampled from a mixture of distributions $p_{\bar{Z}_m} = (1/k) \sum_{i=1}^k p_{\bar{Z}_i}$ where $p_{\bar{Z}_i}$ ($1 \leq i \leq k$) is very close to $\mathcal{N}(0, I)$. Thus, we can use a single Gaussian distribution $\mathcal{N}_{\bar{Z}_s}$ to approximate each $p_{\bar{Z}_i}$. Therefore, $D_g[p_{Z_3}]$ can be approximated as

$$D_g[p_{Z_3}] = \sum_{i=0}^k KL(p_{\bar{Z}_i}(z)||\mathcal{N}(0, I)) \approx k \times KL(\mathcal{N}_{\bar{Z}_s}||\mathcal{N}(0, I)) \quad (13)$$

Now we can plug Equation (6) in Equation (13), except that each representation z is treated as k samples of $p_{\bar{Z}_m}$.

For OOD data, we cannot use a single Gaussian distribution to approximate $q_{\bar{Z}_m} = (1/k) \sum_{i=1}^k q_{\bar{Z}_i}$ when $q_{\bar{Z}_i}$ are far from each other or q_Z is not Gaussian-like. Nevertheless, we can still use fitted Gaussian and Equation (6) to compute the lower bound according to Theorem 5.

Summary. Overall, we get the following answer to Q2.

Answer to Q2: We use group-wise KL divergence in last scale (*i.e.*, D_g in Equation 9) as a unified criterion for both GAD and PAD with flow-based generative models.

5.4 Algorithm

Algorithm 1 shows the details of our OOD detection method. The inputs are a set of data points $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_m\}$ ($m \geq$

Algorithm 1 KL divergence-based Out-of-Distribution Detection with Split representations (KLODS)

```

1: Input:  $f(x)$ : a well-trained flow-based model or the encoder of VAE using Gaussian prior  $\mathcal{N}(\mu, \text{diag}(\sigma))$ ;  $\mathbf{X} = \{x_1, \dots, x_m\}$  ( $m \geq 1$ ): a batch of inputs;  $(H, W, C)$ : the height, width, and channels of last scale of representations.  $t$ : threshold
2:  $\bar{\mathbf{Z}} = \emptyset$ 
3: for  $i = 1$  to  $m$  do
4:    $z_i = f(x_i)$ 
5:    $\bar{z}_i = (z_i - \mu) / \sigma$ 
6:   if  $z_i$  consists of multiple stages then
7:      $\bar{z}_i =$  last scale of  $z_i$ 
8:   end if
9:   split  $\bar{z}_i$  as  $C$  ( $H \times W$ )-dimensional subvectors  $\bar{z}_{i,1} \dots \bar{z}_{i,C}$ 
10:   $\bar{\mathbf{Z}} = \bar{\mathbf{Z}} \cup \{\bar{z}_{i,1}, \dots, \bar{z}_{i,C}\}$ 
11: end for
12: calculate sample covariance  $\tilde{\Sigma}$  and sample mean  $\tilde{\mu}$  of  $\bar{\mathbf{Z}}$ 
13:  $score = (1/2)\{-\log|\tilde{\Sigma}| + \text{Tr}(\tilde{\Sigma}) + \tilde{\mu}^\top \tilde{\mu} - n\}$ 
14: if  $score > t$  then
15:   return “ $\mathbf{X}$  is OOD data”
16: else
17:   return “ $\mathbf{X}$  is ID data”
18: end if

```

1), where each x_i is an individual input. Here we support both GAD and PAD (in the case of $m = 1$). For each input x_i , we first compute the representation $z_i = f(x_i)$ (line 4). Here f represents flow-based model or the encoder part of VAE. Then we normalize representations $\{z_1, \dots, z_m\}$ as $\bar{z}_i = (z_i - \mu) / \sigma$, where μ and $\text{diag}(\sigma^2)$ are the mean and covariance matrix of Gaussian prior, respectively (line 5). If z_i has multiple stages, we choose only the last-scale representation to leverage last-scale KL divergence (line 6, Section 5.2). Then we split \bar{z}_i as C ($H \times W$)-dimensional subvectors to leverage local pixel dependence as discussed in Subsection 5.3.1. We collect the subvectors from each x_i and treat $\bar{\mathbf{Z}}$ as $m \times C$ data points sampled from a mixture of distributions (line 10). Then we calculate the sample covariance $\tilde{\Sigma}$ and sample mean $\tilde{\mu}$ of $\bar{\mathbf{Z}}$ (line 12). Finally, we use the following *anomaly score*

$$score = (1/2)\{-\log|\tilde{\Sigma}| + \text{Tr}(\tilde{\Sigma}) + \tilde{\mu}^\top \tilde{\mu} - n\} \quad (14)$$

as the criterion (line 13). For ID data, $score$ is the estimated group-wise KL divergence in last scale (i.e., $D_g[p_{Z_3}]$) except for neglecting the constant k in Equation (13). For OOD data, $score$ is the lower bound of $D_g[q_{Z_3}]$. The larger $score$ is, the more like OOD the input. If $score$ is greater than a threshold t , the input is determined as OOD data (line 15). Otherwise, the input is determined as ID data (line 17).

We name our method as **KLODS** for *KL divergence-based Out-of-Distribution Detection with Split representations*. We also call our method without split representations as **KLOD**.

5.5 Summary

In Figure 1, we have illustrated our analysis steps in explaining why we cannot sample OOD data. In Figure 2, we summarize how to leverage KL divergence for OOD detection in Section 5. To help readers have a bird’s eye view

of our whole work, we summarize all critical steps in a big flowchart in Figure E.1 in the supplementary material.

6 EXPERIMENTS

We conduct experiments to evaluate the effectiveness and robustness of our OOD detection method.

6.1 Experimental Setting

Benchmarks. We evaluate our method with prevalent benchmarks in deep anomaly detection research [18], [19], [22], [66], [67], [68], including Constant, Uniform, MNIST [69], FashionMNIST [70], notMNIST [71], KMNIST [72], Omniglot [73], CIFAR-10/100 [74], SVHN [75], CelebA [76], TinyImageNet [77], ImageNet32 [78], and LSUN [79]. We use different dataset compositions falling into *Category I* (smaller/similar variance, higher/similar likelihoods, e.g., CIFAR-10 vs SVHN) and *Category II* (larger variance, lower likelihoods, e.g., SVHN vs CIFAR-10) problems. All datasets are resized to $32 \times 32 \times 3$ for consistency. We use $S\text{-}C(k)$ ($k \geq 0$) to denote dataset S with adjusted contrast by a factor k . More details of the benchmarks are described in Section F in the supplementary material.

Baselines. We choose the following recently published OOD detection methods as baselines.

GAD:

- 1) *t*-test: two-sample students’ *t*-test for a difference in means in the empirical likelihoods.
- 2) **Kolmogorov-Smirnov test (KS-test)**: two-sample KS-test to the likelihood empirical distribution functions.
- 3) **Maximum Mean Discrepancy (MMD)** [80]: two-sample MMD test.
- 4) **Kernelized Stein Discrepancy (KSD)** [81]: test for Goodness of Fit to the generative model.
- 5) **Annulus Method** [20]: Typicality test in latent space. Inputs whose latents are far from the annulus with radius \sqrt{n} are classified as OOD data.
- 6) **Ty-test** [22]: typicality test in model distribution.
- 7) **GOD2KS** [34]: combining random projection and two-sample KS test.

Among the above GAD methods, Annulus Method, Ty-test, and GOD2KS are the best ones. We reimplement Annulus Method and Ty-test to produce more results.

PAD:

- 1) **S** [28]: input complexity compensated likelihood.
- 2) L_{last} [29]: likelihood contributed by the last-scale representation of Glow.
- 3) **DoSE** [30]: density estimators on the statistics of models to detect OOD data.
- 4) **ODIN** [82]: Liang *et al.* introduce ODIN method for OOD detection.
- 5) **Joint confidence loss** [83]: Lee *et al.* introduce joint confidence loss for OOD detection.
- 6) **Joint confidence loss+ODIN** [83]: combination of Joint confidence loss and ODIN (better than each method alone).

For a more comprehensive evaluation, we reimplement the first three PAD baselines applicable to flow-based model. DoSE is the SOTA PAD method applicable to flow-based model. The rest baselines apply to classification networks rather than flow-based models. See Section H in the supplementary material for more discussion about baselines.

Models. We use the official Glow model [64] and the model released by the authors of Ty-test (DeepMind [84]). See Section G in the supplementary material for details.

Metrics. We use the same metrics as baseline methods in their original publications. These metrics include false positive rate (FPR), true positive rate (TPR), threshold-independent metrics area under the receiver operating characteristic curve (AUROC) and area under the precision-recall curve (AUPR) [85], and threshold-dependent Equal Error Rate (EER). We treat OOD data as positive data. For GAD, each dataset is shuffled and then divided into groups of size m . We run each method for 5 times and show “mean \pm standard deviation” for each GAD problem.

6.2 Experimental Results

6.2.1 Group Anomaly Detection

Main Results on Unconditional Glow.

FashionMNIST vs Others. Table J.1 in the supplementary material shows the GAD results of Glow trained on FashionMNIST. The ID column reflects false positive rate (ideally should be 0). The MNIST and notMNIST columns reflect true positive rate (ideally should be 1). The authors of baselines apply bootstrap procedure on validation data to establish thresholds. See Section J.1 in the supplementary materials for the details on how we establish thresholds. We can see that all methods cannot achieve satisfactory results with small batch size $m = 2$. Our method achieves the highest true positive rate with the lowest false positive rate for larger batch sizes (*i.e.*, 10, 25).

The first two subfigures of Figure 3 show the comparison of KLODS, our reimplementation of Ty-test, and Annulus Method on FashionMNIST vs Others. The corresponding numerical results of Figure 3 are shown in Table J.2, J.3, and J.4 in the supplementary material. In our reimplementation, Annulus Method achieves much better results than that reported in [22] (and Table J.1 in the supplementary material). Nevertheless, our method outperforms all baselines significantly.

SVHN/CIFAR-10/CelebA vs Others. Figure 3 also shows the GAD results on Glow trained on SVHN, CIFAR-10, and CelebA. Our method is the best one. We adjust the contrast of OOD dataset to make the likelihood distributions of ID and OOD data coincide. For these kinds of problems, the performance of Annulus Method and Ty-test degenerate severely. Our method is more robust against data manipulation.

CelebA vs CIFAR-10/100 are challenging for Ty-test as reported by [22]. Our method can achieve 100% AUROC with batch size 10. In our experiments, it is hard to make the likelihood distributions of CelebA train and test split fit very well on the official Glow model⁶. This affects the performance of Ty-test. Please see Section J.1 in the supplementary material for more discussion.

CIFAR-10 vs CIFAR-100 is one of the most challenging problems. Annulus Method and Ty-test achieve 47.2% and 72.4% AUROCs with batch size $m = 200$, respectively. KLOD and KLODS achieve around 70% AUROC when the batch size $m = 200$. We think this is due to unsuccessful model and the similarity between ID and OOD datasets. Please see Section H in the supplementary material for more discussion on CIFAR-10 vs CIFAR-100.

6. We stop training after 2,000 epochs.

Smaller Batch Sizes. KLODS outperforms Ty-test when batch size is smaller (*i.e.*, $2 \sim 4$). See Table J.5 in the supplementary material for details.

Comparison with GOD2KS. Table J.6 in the supplementary material compares our method and GOD2KS [34] on Glow. We use the same problems reported in [34]. Our method outperforms GOD2KS.

Robustness. The results presented above have demonstrated the robustness of our method against data manipulation method M2 (adjusting contrast). Experimental results show that KLODS achieves the same performance under M1 (rescaling representations), except that a slightly larger batch size (+5) is needed for CIFAR-10-related problems. As shown in Figure 3, Table J.2, J.3, and J.4 in the supplementary material, Annulus Method and Ty-test is affected by data manipulation M2 (adjusting contrast). Besides, *Annulus Method achieves exactly 0 AUROCs for all problems under data manipulation M1 (rescaling representations)*. This is because all OOD representations are rescaled to the annulus of typical set of prior and hence definitely closer to the typical set annulus than ID representations (see Section 3.1). The results are omitted for brevity. Currently, the performance of GOD2KS under data manipulations is still not clear.

Summary. For GAD, our method achieves 98.1% AUROC, 98.2% AUPR, and 4.6% EER on average with batch size 5 and outperforms Ty-test by 33.5%, 29.2%, 29.3% on average in AUROC, AUPR, and EER, respectively. Our method also outperforms GOD2KS by 9.1%, 12.1% on average in AUROC and AUPR with batch size 5, respectively. Our method is robust against data manipulations, while the baseline methods Ty-test and Annulus Method can be attacked in almost all cases.

More Results.

Mixture of OOD Datasets. We also use the mixture of two datasets among SVHN, CelebA, and CIFAR-10 as one OOD dataset. We can treat samples from multiple distributions as from a mixture of distributions. We randomly choose 5,000 samples from each dataset and get 10,000 samples as one OOD dataset. Table J.7 in the supplementary material shows the results of KLODS. Our method outperforms Ty-test by 38.9% AUROC on average with batch size 5.

Ablation Study. We compare the following four methods to evaluate how the techniques proposed in Section 5 affect the performance.

- 1) **Ty-test:** the baseline.
- 2) **KLOD-all:** GAD with all scales of representation, without splitting dimensions.
- 3) **KLOD:** GAD using the last-scale representation, without splitting dimensions.
- 4) **KLODS:** GAD using the last-scale representation with splitting dimensions.

Table J.8 in the supplementary material shows the results. Neglecting CIFAR10 vs ImageNet32-C(0.3), the order of the methods by performance is KLODS > KLOD > KLOD-all > Ty-test. The only exception is CIFAR10 vs ImageNet32-C(0.3), where KLOD outperforms KLODS. The low contrast leads to weak local pixel dependence and affects our splitting strategy. Overall, we can see that both using the last scale and splitting dimensions into groups can improve the

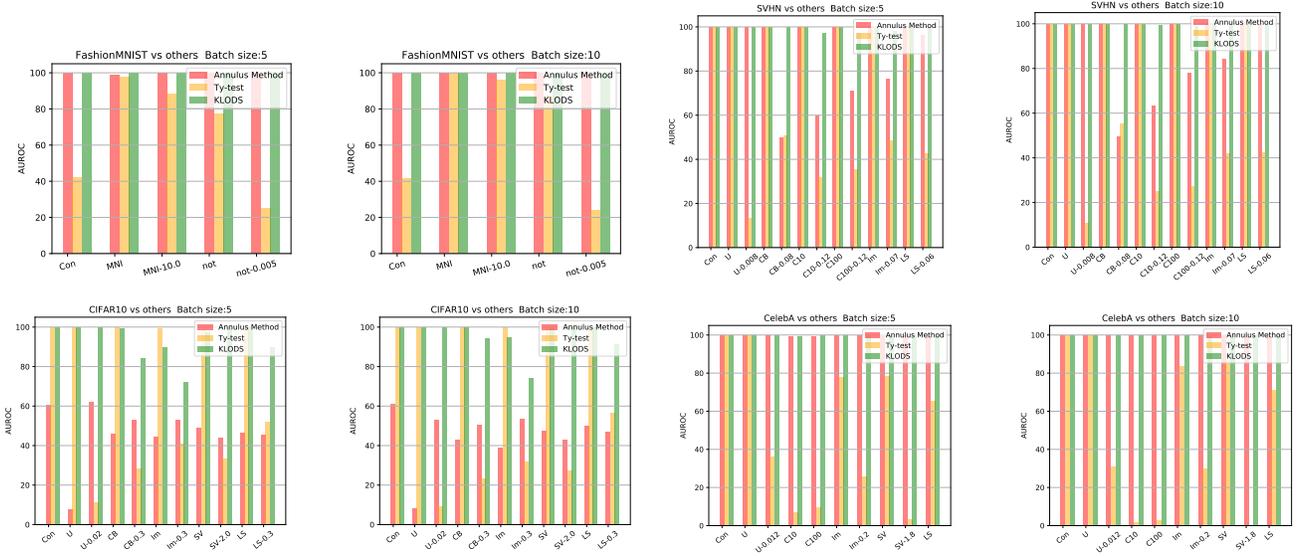


Fig. 3: GAD results (AUROC) on Glow with batch sizes 5 and 10. The X-axis are labeled with OOD datasets including Con: Constant, MNI: MNIST, not: notMNIST, U: Uniform, CB: CelebA, C10/100: CIFAR-10/100, Im: ImageNet, SV: SVHN, LS: LSUN. The corresponding numerical results are shown in Table J.2, J.3, and J.4 in the supplementary material.

performance of GAD. Note that splitting dimensions also makes PAD feasible. Besides, when the batch size is smaller (e.g., 5), KLODS outperforms KLOD more obviously. More results are not shown for brevity.

One-vs-Rest. We conduct one-vs-rest evaluation on MNIST. For each class from 0 to 9, we use images in that class as ID data and the rest classes as OOD data. We train one Glow model under each setting for 120 epochs. As shown in Table J.9 in the supplementary material, our method achieves 85.4% AUROC and 85.8% AUPR on average with batch size 5, outperforming the baseline by 8% and 5%, respectively.

GAD on GlowGMM. We train GlowGMM on FashionMNIST. We treat each class as ID data and the rest as OOD data. KLODS can achieve 100% AUROC on average when batch size is 25. On the contrary, Ty-test is worse than random guessing in most cases. See Figure J.1 and Table J.10 in the supplementary material for results. Experimental results also demonstrate that each component may assign higher likelihoods to other classes (See Table J.11 in the supplementary material).

Generating OOD Images Using GlowGMM. In GlowGMM, we can generate *high-quality* OOD images. See Section J in the supplementary material for more discussion.

GAD on VAE. We train convolutional VAE with 8-/16-/32-dimensional latent space on FashionMNIST, SVHN, and CIFAR-10, respectively. The latent space is not large enough, so we did not split representations and only used KLOD in experiments. The results are shown in Figure J.3 and Table J.12 in the supplementary material. KLOD achieves 99.9% AUROC on average when $m = 25$ for most problems. CIFAR-10 vs CIFAR-100 is also the most challenging problem on VAE. KLOD needs a batch size of 150 to achieve 98%+ AUROC (See Table J.13 in the supplementary material). Nevertheless, KLOD still outperforms Ty-test. Again, Ty-test can be attacked by data manipulations M2 (adjusting contrast). Finally, as pointed out by [86], for vanilla VAE the

reconstruction probability is not a reliable criterion for OOD detection (See Table J.14 in the supplementary material).

6.2.2 Point-wise Anomaly Detection

The PAD results of KLODS, S , L_{last} , and DoSE are shown in Table 1.

SVHN vs Others. The problems above the dash line in Table 1 fall in *Category II* (larger variance, lower likelihoods). KLODS can achieve 98.8%+ AUROC and outperforms the baselines. In [28], although the authors state that their method S can detect OOD data with more complexity than ID data (roughly *Category II*), they did not evaluate their method thoroughly on *Category II* problems. We find S does not perform well on these problems.

The problems for SVHN vs others below the dash line in Table 1 fall in *Category I* (smaller/similar variance, higher likelihoods). For these problems, L_{last} and DoSE degenerate into being not better than random guessing. KLODS is comparable with S and outperforms L_{last} and DoSE significantly. The reason is all the distributions of $\log p(\mathbf{x})$, $\log p(\mathbf{z})$, and $\log p(\mathbf{x})$ contributed by the last scale overlap with those of ID data. Figure K.7 and K.8 in the supplementary material shows the histograms of these three statistics. These issues make DoSE fail because DoSE relies on the effectiveness of its based statistics.

CelebA vs Others. The performance of S degenerates severely on these problems. Our method is slightly affected because the likelihoods of the train and test split of CelebA do not fit very well (see Figure K.10 in the supplementary material).

CIFAR-10 vs Others. As discussed in the last subsection, Glow model fails to generate high-quality CIFAR-10-like images. Our method is affected on CIFAR-10 vs others. As discussed before, we argue that it is hard to require an “unsuccessful” model can detect OOD data.

TABLE 1: PAD results (AUROC in percentage) on Glow. Notable failures (below 60%) are underlined.

ID	OOD	S	L_{last}	DoSE	KLODS
SVHN	Uniform	100.0	100.0	100.0	100.0
	ImageNet32	78.7	99.8	99.9	99.9
	CelebA	83.1	100.0	100.0	100.0
	CIFAR-10	43.8	97.7	96.2	98.9
	CIFAR-100	<u>44.9</u>	97.3	96.5	98.8
	LSUN	91.8	100.0	91.6	100.0
	Uniform-C(0.008)	<u>97.9</u>	0.0	96.8	98.6
	CelebA-C(0.08)	81.4	<u>41.6</u>	<u>48.0</u>	<u>82.2</u>
	CIFAR-10-C(0.12)	75.3	47.7	<u>50.5</u>	72.5
	CIFAR-100-C(0.12)	75.2	48.6	<u>54.5</u>	75.3
	ImageNet32-C(0.07)	99.6	<u>42.2</u>	<u>55.7</u>	84.0
	LSUN-C(0.06)	81.3	3.0	69.5	91.6
	notMNIST	100.0	98.7	99.9	99.6
	Constant	100.0	0.4	99.9	99.8
CelebA	Constant	98.0	99.8	99.9	100.0
	Uniform	91.0	100	100.0	100.0
	Uniform-C(0.012)	97.2	98.1	90.9	99.5
	ImageNet	16.5	99.7	99.8	100.0
	ImageNet-C(0.2)	88.5	97.9	91	93.3
	CIFAR-10	<u>55.0</u>	90.4	94.9	69.0
	CIFAR-100	<u>53.2</u>	90.6	95.6	72.3
	SVHN	83.9	99.3	99.7	94.7
CIFAR-10	SVHN-C(1.8)	90.5	99.9	85.2	98.9
	LSUN	65.4	99.6	84.9	99.2
	Constant	100	1.4	99.8	98.9
	Uniform	100	100	100	100
	Uniform-C(0.2)	98.8	1.9	64.7	99.7
	CelebA	86.3	96.6	99.5	85.2
	CelebA-C(0.3)	95.0	7.8	<u>46.5</u>	64.9
	SVHN	95.0	92.9	95.5	82.6
	SVHN-C(2.0)	94.0	98.9	93.7	95
	TinyImageNet	71.6	90.7	76.7	83.9
	CIFAR-100	73.6	60.0	<u>57.1</u>	<u>54.1</u>
	LSUN	91.1	82.8	98.0	98.9
	LSUN-C(0.3)	96.4	94.8	61.2	83.3
	average	82.7	73.7	85.5	90.7
#notable failures	5	5	6	1	

Finally, our method has only one notable failure (*i.e.*, below 60% AUROC). S , L_{last} , and DoSE have 5, 5, and 6 notable failures in total, respectively.

Other comparisons.

We compare KLODS with Joint confidence loss, ODIN, and Joint confidence loss+ODIN. These three baseline methods do not apply to flow-based model. The results are shown in Table J.15 in the supplementary material, where we use the same datasets reported in [83]. Our method is the best one. See Section H in the supplementary material for more discussion on our method and results.

Summary. For PAD, our method achieves 90.7% AUROC on average and outperforms the SOTA baseline DoSE by 5.2% in AUROC. Our method also has the least notable failures.

7 CONCLUSION

In this paper, we prove theorems to investigate KL divergences in flow-based models. We observe the normality of ID and OOD representations in flow-based model for a wide range of problems. Based on our analysis, we explain why we cannot sample OOD data from flow-based model from two perspectives. We propose leveraging KL divergence for OOD detection. We further decompose the KL divergence to leverage the last-scale KL divergence of Glow model. Furthermore, we split representations into groups to leverage group-wise KL divergence as the final OOD

detection criterion. Experimental results have demonstrated the effectiveness and robustness of our method.

REFERENCES

- [1] V. Chandola, A. Banerjee, and V. Kumar, "Anomaly detection: A survey," *ACM Comput. Surv.*, vol. 41, no. 3, Jul. 2009.
- [2] G. Pang, C. Shen, L. Cao, and A. V. D. Hengel, "Deep learning for anomaly detection: A review," *ACM Comput. Surv.*, vol. 54, no. 2, Mar. 2021.
- [3] E. Toth and S. Chawla, "Group deviation detection methods: A survey," *ACM Comput. Surv.*, vol. 51, no. 4, Jul. 2018.
- [4] K. Muandet and B. Schölkopf, "One-class support measure machines for group anomaly detection," in *Proceedings of the Twenty-Ninth Conference on Uncertainty in Artificial Intelligence*, ser. UAI'13. Arlington, Virginia, USA: AUAI Press, 2013, p. 449–458.
- [5] G. Jorge, C. Stephane, and H. R., "Support measure data description for group anomaly detection," *ODDx3 Workshop on Outlier Definition, Detection, and Description at ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD 2015)*, 2015.
- [6] L. Xiong, B. Póczos, J. Schneider, A. Connolly, and J. VanderPlas, "Hierarchical Probabilistic Models for Group Anomaly Detection." *Journal of Machine Learning Research - Proceedings Track*, vol. 15, pp. 789–797, 2011.
- [7] L. Xiong, B. Póczos, and J. Schneider, "Group anomaly detection using flexible genre models," in *Proceedings of the 24th International Conference on Neural Information Processing Systems*, ser. NIPS'11. Red Hook, NY, USA: Curran Associates Inc., 2011, p. 1071–1079.
- [8] A. Kuppa, S. Grzonkowski, M. R. Asghar, and N. Le-Khac, "Finding rats in cats: Detecting stealthy attacks using group anomaly detection," in *2019 18th IEEE International Conference On Trust, Security And Privacy In Computing And Communications/13th IEEE International Conference On Big Data Science And Engineering (TrustCom/BigDataSE)*, 2019, pp. 442–449.
- [9] R. Chalapathy and S. Chawla, "Deep learning for anomaly detection: A survey," 2019.
- [10] S. P. Mishra and P. Kumari, "Analysis of Techniques for Credit Card Fraud Detection: A Data Mining Perspective," in *New Paradigm in Decision Science and Management*, S. Patnaik, A. W. H. Ip, M. Tavana, and V. Jain, Eds. Singapore: Springer Singapore, 2020, pp. 89–98.
- [11] Y. Ye, T. Li, D. Adjeroh, and S. S. Iyengar, "A survey on malware detection using data mining techniques," *ACM Comput. Surv.*, vol. 50, no. 3, Jun. 2017.
- [12] L. Xiong, B. Póczos, and J. Schneider, "Group anomaly detection using flexible genre models," in *Proceedings of the 24th International Conference on Neural Information Processing Systems*, ser. NIPS'11. Red Hook, NY, USA: Curran Associates Inc., 2011, p. 1071–1079.
- [13] D. P. Kingma and P. Dhariwal, "Glow: Generative flow with invertible 1x1 convolutions," in *Advances in Neural Information Processing Systems*, 2018, pp. 10 215–10 224.
- [14] L. Dinh, J. Sohl-Dickstein, and S. Bengio, "Density estimation using Real NVP," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [15] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2014.
- [16] A. Van den Oord, N. Kalchbrenner, L. Espeholt, O. Vinyals, A. Graves *et al.*, "Conditional image generation with pixelcnn decoders," in *Advances in neural information processing systems*, 2016, pp. 4790–4798.
- [17] T. Salimans, A. Karpathy, X. Chen, and D. P. Kingma, "Pixel-CNN++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications," *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [18] E. Nalisnick, A. Matsukawa, Y. W. Teh, D. Gorur, and B. Lakshminarayanan, "Do deep generative models know what they don't know?" *ICLR*, 2019.
- [19] A. Shafaei, M. Schmidt, and J. J. Little, "Does your model know the digit 6 is not a cat? a less biased evaluation of "outlier" detectors," *arXiv preprint arXiv:1809.04729*, 2018.
- [20] H. Choi and E. Jang, "WAIC, but why?: Generative ensembles for robust anomaly detection," *arXiv preprint arXiv:1810.01392*, 2018.
- [21] V. Škvára, T. Pevný, and V. Šmídl, "Are generative deep models for novelty detection truly better?" *KDD Workshop on Outlier Detection De-Constructed (ODD v5.0)*, 2018.

- [22] E. Nalisnick, A. Matsukawa, Y. W. Teh, and B. Lakshminarayanan, "Detecting out-of-distribution inputs to deep generative models using typicality," *4th workshop on Bayesian Deep Learning (NeurIPS 2019)*, 2019.
- [23] P. Kirichenko, P. Izmailov, and A. G. Wilson, "Why normalizing flows fail to detect out-of-distribution data," *ICML workshop on Invertible Neural Networks and Normalizing Flows, 2020 (NeurIPS 2020)*, 2020.
- [24] R. T. Q. Chen, J. Behrmann, D. Duvenaud, and J. Jacobsen, "Residual flows for invertible generative modeling," in *Advances in Neural Information Processing Systems*, 2019.
- [25] E. Fetaya, J. Jacobsen, and R. S. Zemel, "Conditional generative models are not robust," *CoRR*, vol. abs/1906.01171, 2019.
- [26] P. Izmailov, P. Kirichenko, M. Finzi, and A. G. Wilson, "Semi-supervised learning with normalizing flows," 2019.
- [27] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. Los Alamitos, CA, USA: IEEE Computer Society, jun 2016, pp. 770–778.
- [28] J. Serrà, D. Álvarez, V. Gómez, O. Slizovskaia, J. F. Núñez, and J. Luque, "Input complexity and out-of-distribution detection with likelihood-based generative models," in *International Conference on Learning Representations*, 2020.
- [29] R. T. Schirmeister, Y. Zhou, T. Ball, and D. Zhang, "Understanding Anomaly Detection with DeepInvertible Networks through Hierarchies of Distributions and Features," in *Advances in Neural Information Processing Systems 33*. Curran Associates, Inc., 2020.
- [30] W. Morningstar, C. Ham, A. Gallagher, B. Lakshminarayanan, A. Alemi, and J. Dillon, "Density of States Estimation for Out of Distribution Detection," in *Proceedings of The 24th International Conference on Artificial Intelligence and Statistics*, vol. 130. PMLR, 2021, pp. 3232–3240.
- [31] L. H. Zhang, M. Goldstein, and R. Ranganath, "Understanding failures in out-of-distribution detection with deep generative models," *CoRR*, vol. abs/2107.06908, 2021.
- [32] M. A. F. Pimentel, D. A. Clifton, C. Lei, and L. Tarassenko, "A review of novelty detection," *Signal Processing*, vol. 99, no. 6, pp. 215–249, 2014.
- [33] E. Sabeti and A. Hostmadsen, "Data discovery and anomaly detection using atypicality for real-valued data," *Entropy*, vol. 21, no. 3, p. 219, 2019.
- [34] D. Jiang, S. Sun, and Y. Yu, "Revisiting flow generative models for out-of-distribution detection," in *International Conference on Learning Representations*, 2022.
- [35] J. Ren, P. J. Liu, E. Fertig, J. Snoek, R. Poplin, M. A. DePristo, J. V. Dillon, and B. Lakshminarayanan, "Likelihood ratios for out-of-distribution detection," 2019.
- [36] L. Dinh, D. Krueger, and Y. Bengio, "NICE: Non-linear independent components estimation," *arXiv preprint arXiv:1410.8516*, 2014.
- [37] G. Papamakarios, E. Nalisnick, D. J. Rezende, S. Mohamed, and B. Lakshminarayanan, "Normalizing flows for probabilistic modeling and inference," 2019.
- [38] R. Vershynin, *High-dimensional probability: An introduction with applications in data science*. Cambridge University Press, 2018, vol. 47.
- [39] T. M. Cover and J. A. Thomas, *Elements of information theory*. John Wiley & Sons, 2012.
- [40] F. Nielsen, "An elementary introduction to information geometry," *arXiv preprint arXiv:1808.08271*, 2018.
- [41] S. Nowozin, B. Cseke, and R. Tomioka, "*f*-GAN: Training generative neural samplers using variational divergence minimization," in *Advances in Neural Information Processing Systems*, D. Lee, M. Sugiyama, U. Luxburg, I. Guyon, and R. Garnett, Eds., vol. 29. Curran Associates, Inc., 2016.
- [42] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of Wasserstein GANs," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, ser. NIPS'17. Red Hook, NY, USA: Curran Associates Inc., 2017, p. 5769–5779.
- [43] L. Pardo, *Statistical inference based on divergence measures*. CRC press, 2018.
- [44] N. Mohd Razali and B. Yap, "Power comparisons of shapiro-wilk, kolmogorov-smirnov, lilliefors and anderson-darling tests," *J. Stat. Model. Analytics*, vol. 2, 01 2011.
- [45] Y. Zhang, W. Liu, Z. Chen, K. Li, and J. Wang, "On the properties of Kullback-Leibler divergence between multivariate Gaussian distributions," *arXiv preprint arXiv:2102.05485*, 2021.
- [46] T. Q. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud, "Isolating sources of disentanglement in variational autoencoders," in *Advances in Neural Information Processing Systems*, 2018, pp. 2610–2620.
- [47] M. Giraudo, L. Sacerdote, and R. Sirovich, "Non-parametric estimation of mutual information through the entropy of the linkage," *Entropy*, vol. 15, no. 12, p. 5154–5177, Nov 2013.
- [48] J. L. Rodgers and W. A. Nicewander, "Thirteen ways to look at the correlation coefficient," *The American Statistician*, vol. 42, no. 1, pp. 59–66, 1988.
- [49] A. Gambardella, A. G. Baydin, and P. H. S. Torr, "Transflow learning: Repurposing flow models without retraining," 2019.
- [50] H. Hoihtink, I. Klugkist, L. D. Broemeling, R. Jensen, Q. Shen, S. Mukherjee, R. A. Bailey, J. L. Rosenberger, J. D. Leeuw, E. Meijer, B. G. Leroux, A. B. Tsybakov, W. Wedelmeyer, P. C. Consul, F. Famoye, and D. Richards, "Introduction to nonparametric estimation." 2009.
- [51] X. Nguyen, M. J. Wainwright, and M. I. Jordan, "Estimating divergence functionals and the likelihood ratio by penalized convex risk minimization," in *In Advances in Neural Information Processing Systems (NIPS)*, 2007.
- [52] P. K. Rubenstein, O. Bousquet, J. Djolonga, C. Riquelme, and I. O. Tolstikhin, "Practical and consistent estimation of *f*-divergences." *Annual Conference on Neural Information Processing Systems*, vol. abs/1905.11112, pp. 4072–4082, 2019.
- [53] Qing Wang, S. R. Kulkarni, and S. Verdu, "Divergence estimation of continuous distributions based on data-dependent partitions," *IEEE Transactions on Information Theory*, vol. 51, no. 9, pp. 3064–3074, 2005.
- [54] Q. Wang, S. R. Kulkarni, and S. Verdu, "Divergence estimation for multidimensional densities via *k*-nearest-neighbor distances," *IEEE Transactions on Information Theory*, vol. 55, no. 5, pp. 2392–2405, 2009.
- [55] X. Nguyen, M. J. Wainwright, and M. I. Jordan, "Estimating divergence functionals and the likelihood ratio by convex risk minimization," *IEEE Trans. Inf. Theor.*, vol. 56, no. 11, p. 5847–5861, Nov. 2010.
- [56] K. R. Moon and A. O. Hero, "Ensemble estimation of multivariate *f*-divergence," in *2014 IEEE International Symposium on Information Theory*, 2014, pp. 356–360.
- [57] C. M. Bishop, *Pattern Recognition and Machine Learning (Information Science and Statistics)*. Berlin, Heidelberg: Springer-Verlag, 2006.
- [58] C. P. Burgess, I. Higgins, A. Pal, L. Matthey, N. Watters, G. Desjardins, and A. Lerchner, "Understanding disentangling in β -VAE," in *Workshop on Learning Disentangled Representations at the 31st Conference on Neural Information Processing Systems*, 2018.
- [59] H. Kim and A. Mnih, "Disentangling by factorising," in *Proceedings of the 35th International Conference on Machine Learning*, vol. 80. PMLR, 10–15 Jul 2018, pp. 2649–2658.
- [60] T. Q. Chen, X. Li, R. B. Grosse, and D. K. Duvenaud, "Isolating sources of disentanglement in variational autoencoders," in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Inc., 2018, pp. 2610–2620.
- [61] A. Kumar, P. Sattigeri, and A. Balakrishnan, "Variational inference of disentangled latent concepts from unlabeled observations," in *International Conference on Learning Representations*, 2017.
- [62] F. Locatello, S. Bauer, M. Lucic, G. Rätsch, S. Gelly, B. Schölkopf, and O. Bachem, "Challenging common assumptions in the unsupervised learning of disentangled representations," in *Proceedings of the 36th International Conference on Machine Learning*, 2019.
- [63] M. D. Hoffman and M. J. Johnson, "ELBO surgery: yet another way to carve up the variational evidence lower bound," in *Workshop in Advances in Approximate Bayesian Inference, NIPS*, vol. 1, 2016.
- [64] OpenAI, "Glow," <https://github.com/openai/glow>, 2018.
- [65] J. D. D. Havtorn, J. Frellsen, S. Hauberg, and L. Maaløe, "Hierarchical vaes know what they don't know," in *Proceedings of the 38th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, M. Meila and T. Zhang, Eds., vol. 139. PMLR, 18–24 Jul 2021, pp. 4117–4128.
- [66] K. Lee, K. Lee, H. Lee, and J. Shin, "A simple unified framework for detecting out-of-distribution samples and adversarial attacks," in *Advances in Neural Information Processing Systems*, 2018, pp. 7167–7177.
- [67] D. Hendrycks and K. Gimpel, "A baseline for detecting misclassified and out-of-distribution examples in neural networks," in

- Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.
- [68] D. Hendrycks, M. Mazeika, and T. G. Dietterich, "Deep anomaly detection with outlier exposure," *International Conference on Learning Representations (ICLR)*, 2019.
- [69] Y. LeCun, L. Bottou, Y. Bengio, P. Haffner *et al.*, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [70] H. Xiao, K. Rasul, and R. Vollgraf, "Fashion-MNIST: a novel image dataset for benchmarking machine learning algorithms," 2017.
- [71] Y. Bulatov, "notMNIST," <http://yaruslavvb.blogspot.com/2011/09/notmnist-dataset.html>, 2011, accessed October 4, 2019.
- [72] T. Clanuwat, M. Bober-Irizar, A. Kitamoto, A. Lamb, K. Yamamoto, and D. Ha, "Deep learning for classical japanese literature," *CoRR*, vol. abs/1812.01718, 2018.
- [73] B. M. Lake, R. Salakhutdinov, and J. B. Tenenbaum, "Human-level concept learning through probabilistic program induction," *Science*, vol. 350, no. 6266, pp. 1332–1338, 2015.
- [74] A. Krizhevsky, G. Hinton *et al.*, "Learning multiple layers of features from tiny images," Citeseer, Tech. Rep., 2009.
- [75] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading digits in natural images with unsupervised feature learning," 2011.
- [76] Z. Liu, P. Luo, X. Wang, and X. Tang, "Deep learning face attributes in the wild," in *Proceedings of International Conference on Computer Vision (ICCV)*, December 2015.
- [77] Stanford, <https://tiny-imagenet.herokuapp.com/>.
- [78] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.
- [79] F. Yu, A. Seff, Y. Zhang, S. Song, T. Funkhouser, and J. Xiao, "Lsun: Construction of a large-scale image dataset using deep learning with humans in the loop," 2015.
- [80] A. Gretton, K. M. Borgwardt, M. J. Rasch, B. Schölkopf, and A. Smola, "A kernel two-sample test," vol. 13, no. null, p. 723–773, mar 2012.
- [81] Q. Liu, J. D. Lee, and M. Jordan, "A kernelized stein discrepancy for goodness-of-fit tests," in *Proceedings of the 33rd International Conference on Machine Learning - Volume 48*, ser. ICML'16. JMLR.org, 2016, p. 276–284.
- [82] S. Liang, Y. Li, and R. Srikant, "Enhancing the reliability of out-of-distribution image detection in neural networks," 2017.
- [83] K. Lee, H. Lee, K. Lee, and J. Shin, "Training confidence-calibrated classifiers for detecting out-of-distribution samples," in *International Conference on Learning Representations*, 2018.
- [84] DeepMind, <https://github.com/y0ast/Glow-PyTorch>.
- [85] M. Buckland and F. Gey, "The relationship between recall and precision," *Journal of the American society for information science*, vol. 45, no. 1, pp. 12–19, 1994.
- [86] J. An and S. Cho, "Variational autoencoder based anomaly detection using reconstruction probability," *Special Lecture on IE*, vol. 2, no. 1, 2015.
- [87] S. M. Ali and S. D. Silvey, "A general class of coefficients of divergence of one distribution from another," *Journal of the Royal Statistical Society: Series B (Methodological)*, vol. 28, no. 1, pp. 131–142, 1966.
- [88] J. Sneyers and P. Wuille, "FLIF: Free lossless image format based on maniac compression," in *2016 IEEE International Conference on Image Processing (ICIP)*, 2016, pp. 66–70.
- [89] N. . D. Challenge, <https://www.aicrowd.com/challenges/neurips-2019-disentanglement-challenge>.
- [90] I. Higgins, D. Amos, D. Pfau, S. Racaniere, L. Matthey, D. Rezende, and A. Lerchner, "Towards a definition of disentangled representations," 2018.
- [91] C. Eastwood and C. K. I. Williams, "A framework for the quantitative evaluation of disentangled representations," in *International Conference on Learning Representations*, 2018.
- [92] I. Higgins, L. Matthey, A. Pal, C. Burgess, X. Glorot, M. Botvinick, S. Mohamed, and A. Lerchner, "beta-vae: Learning basic visual concepts with a constrained variational framework." *ICLR*, vol. 2, no. 5, p. 6, 2017.
- [93] G. Osada, T. Tsubasa, B. Ahsan, and T. Nishide, "Out-of-distribution detection with reconstruction error and typicality-based penalty," 2022.
- [94] J. Bian, X. Hui, S. Sun, X. Zhao, and M. Tan, "A novel and efficient cvae-gan-based approach with informative manifold for semi-supervised anomaly detection," *IEEE Access*, vol. 7, pp. 88 903–88 916, 2019.
- [95] A. Grover, M. Dhar, and S. Ermon, "Flow-GAN: Combining maximum likelihood and adversarial learning in generative models," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence (AAAI-18), the 30th innovative Applications of Artificial Intelligence (IAAI-18), and the 8th AAAI Symposium on Educational Advances in Artificial Intelligence (EAAI-18)*, New Orleans, Louisiana, USA, February 2-7, 2018, S. A. McIlraith and K. Q. Weinberger, Eds. AAAI Press, 2018, pp. 3069–3076.
- [96] Z. Zhao, K. G. Mehrotra, and C. K. Mohan, "Ensemble Algorithms for Unsupervised Anomaly Detection," in *Current Approaches in Applied Artificial Intelligence*, M. Ali, Y. S. Kwon, C.-H. Lee, J. Kim, and Y. Kim, Eds. Springer International Publishing, 2015, pp. 514–525.
- [97] L. E. J. Brouwer, "Beweis der invarianz des n-dimensionalen gebiets," *Mathematische Annalen*, vol. 71, no. 3, pp. 305–313, 1911.
- [98] G. Papamakarios, T. Pavlakou, and I. Murray, "Masked Autoregressive Flow for Density Estimation," in *Advances in Neural Information Processing Systems 30*. Curran Associates, Inc., 2017, pp. 2338–2347.
- [99] Z. Liu, Z. Cen, V. Isenbaev, W. Liu, Z. S. Wu, B. Li, and D. Zhao, "Constrained variational policy optimization for safe reinforcement learning," in *The 5th Multidisciplinary Conference on Reinforcement Learning and Decision Making (RLDM)*, 2022.
- [100] S. Stergiopoulos, *Advanced Signal Processing Handbook*. CRC Press, 2001.
- [101] Y. Zhang, W. Liu, Z. Chen, J. Wang, K. Li, H. Wei, and Z. Chen, "Out-of-distribution detection with distance guarantee in deep generative models," *arXiv preprint arXiv:2002.03328v1*, 2020.
- [102] A. Sinha, K. Ayush, J. Song, B. Uzkent, H. Jin, and S. Ermon, "Negative data augmentation," *International Conference on Learning Representations (ICLR)*, 2021.
- [103] N. Dionelis, M. Yaghoobi, and S. A. Tsafaris, "Boundary of distribution support generator (BDSG): Sample generation on the boundary," in *2020 IEEE International Conference on Image Processing (ICIP)*. IEEE, oct 2020. [Online]. Available: <https://doi.org/10.1109%2Ficp40778.2020.9191341>
- [104] A. Atanov, A. Volokhova, A. Ashukha, I. Sosnovik, and D. Vetrov, "Semi-conditional normalizing flows for semi-supervised learning," 2019.
- [105] "Residual flows," <https://github.com/rtqichen/residual-flows>.
- [106] "Residual flows model checkpoints," <https://github.com/rtqichen/residual-flows/releases/tag/v1.0.0>.

APPENDIX A BACKGROUND

Flow-based generative model constructs diffeomorphism f from visible space \mathcal{X} to latent space \mathcal{Z} [13], [14], [36], [37]. The model uses a series of diffeomorphisms implemented by multilayered neural networks

$$\mathbf{x} \xleftarrow{f_1} \mathbf{h}_1 \xleftarrow{f_2} \mathbf{h}_2 \dots \xleftarrow{f_n} \mathbf{z}$$

like flow. The whole bijective transformation $f(\mathbf{x}) = f_n \circ f_{n-1} \dots \circ f_1(\mathbf{x})$ can be seen as encoder, and the inverse function $f^{-1}(\mathbf{z})$ is used as decoder. According to the change of variable rule, the probability density function of the model can be formulated as

$$\begin{aligned} \log p_{\mathcal{X}}(\mathbf{x}) &= \log p_{\mathcal{Z}}(f(\mathbf{x})) + \log \left| \det \frac{\partial \mathbf{z}}{\partial \mathbf{x}^T} \right| \\ &= \log p_{\mathcal{Z}}(f(\mathbf{x})) + \sum_{i=1}^n \log \left| \det \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}^T} \right| \end{aligned} \quad (15)$$

where $\mathbf{x} = \mathbf{h}_0$, $\mathbf{z} = \mathbf{h}_n$, $\frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}^T}$ is the Jacobian of f_i , \det is the determinant.

Here prior $p_{\mathcal{Z}}(\mathbf{z})$ is chosen as tractable density function. For example, the most popular prior is standard Gaussian distribution $\mathcal{N}(0, I)$, which makes $\log p_{\mathcal{Z}}(\mathbf{z}) = -(1/2) \times \sum_i z_i^2 + C$ (C is a constant). After training, one can sample noise ε from prior and generate new samples $f^{-1}(\varepsilon)$.

In this paper, we replace prior $\mathcal{N}(0, I)$ with fitted Gaussian distributions $\mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ from OOD representations, where $\tilde{\boldsymbol{\mu}}$ and $\tilde{\boldsymbol{\Sigma}}$ are the sample mean and covariance, respectively. Then we sample noise $\varepsilon' \sim \mathcal{N}(\tilde{\boldsymbol{\mu}}, \tilde{\boldsymbol{\Sigma}})$ and generate OOD samples $f^{-1}(\varepsilon')$.

Variational Autoencoder (VAE) is directed graphical model approximating the data distribution $p(\mathbf{x})$ with encoder-decoder architecture [15]. The probabilistic encoder $q_\phi(\mathbf{z}|\mathbf{x})$ approximates the unknown intractable posterior $p(\mathbf{z}|\mathbf{x})$. The probabilistic decoder $p_\theta(\mathbf{x}|\mathbf{z})$ approximates $p(\mathbf{x}|\mathbf{z})$. In VAE, the variational lower bound of the marginal likelihood of data points (ELBO)

$$\mathcal{L}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N E_{\mathbf{z} \sim q_\phi} [\log p_\theta(\mathbf{x}^i|\mathbf{z})] - KL(q_\phi(\mathbf{z}|\mathbf{x}^i)||p(\mathbf{z})) \quad (16)$$

can be optimized using stochastic gradient descent. After training, one can sample \mathbf{z} from prior $p(\mathbf{z})$ and use the decoder $p_\theta(\mathbf{x}|\mathbf{z})$ to generate new samples.

APPENDIX B DEFINITIONS

Definition 1 (ϕ -divergence) *The ϕ -divergence between two densities $p(\mathbf{x})$ and $q(\mathbf{x})$ is defined by*

$$D_\phi(p, q) = \int \phi(p(\mathbf{x})/q(\mathbf{x}))q(\mathbf{x})d\mathbf{x}, \quad (17)$$

where ϕ is a convex function on $[0, \infty)$ such that $\phi(1) = 0$. When $q(\mathbf{x}) = 0$, $0\phi(0/0) = 0$ and $0\phi(p/0) = \lim_{t \rightarrow \infty} \phi(t)/t$ [87].

ϕ -divergence family is used widely in machine learning fields. As shown in Table B.1, many commonly used measures, including the KL divergence, Jensen-Shannon divergence, and squared Hellinger distance, belong to the ϕ -divergence family. Many ϕ -divergences are not proper distance metrics and do not satisfy the triangle inequality.

TABLE B.1: Examples of ϕ -divergence family

$\phi(x)$	Divergence
$x \log x - x + 1$	Kullback-Leibler
$-\log x + x - 1$	Minimum Discrimination Information
$(x-1) \log x$	J -Divergence
$\frac{1}{2} 1-x $	Total Variation Distance
$(1-\sqrt{x})^2$	Squared Hellinger distance
$x \log \frac{2x}{x+1} + \log \frac{2}{x+1}$	Jensen-Shannon divergence

APPENDIX C TABLES

Table C.1 summarizes the notations of distributions and KL divergences involved in our analysis.

Table C.2 shows some approximate values of the supremum of KL divergence in Theorem 2.

TABLE C.1: Distributions and KL divergences in our analysis.

Notations	Explanations
$\mathbf{z} = f(\mathbf{x})$	the flow-based model function
p_X	the distribution of ID data
q_X	the distribution of OOD data
p_Z	the distribution of ID representations
q_Z	the distribution of OOD representations
p_Z^r	the prior of the flow-based model
p_X^r	the model induced distribution such that $Z_r \sim p_Z^r$ and $X_r = f^{-1}(Z_r) \sim p_X^r$
$KL(p_X q_X)$	the KL divergence between p_X and q_X , assumed to be any large (by Assumption 1).
$KL(p_Z q_Z)$	the KL divergence between p_Z and q_Z , equals to $KL(p_X q_X)$ (by Theorem 1).
$KL(p_Z p_Z^r)$	the KL divergence between p_Z and prior, trained to be small (by Assumption 2).
$KL(p_Z^r q_Z)$	the KL divergence between prior and q_Z , influenced by $KL(p_Z q_Z)$ and $KL(p_Z p_Z^r)$. By relaxed triangle inequality (Theorem 3), a small $KL(p_Z q_Z)$ and a large $KL(p_Z p_Z^r)$ imply $KL(p_Z^r q_Z)$ is large.
$KL(q_Z p_Z^r)$	the KL divergence between prior and q_Z , influenced by $KL(p_Z^r q_Z)$. By the approximate symmetry property (Theorem 2), a large $KL(p_Z^r q_Z)$ must lead to a large $KL(q_Z p_Z^r)$.

TABLE C.2: Some approximate values of the supremum of KL divergence

ε	0.001	0.005	0.01	0.05	0.1	0.5
sup	0.001	0.006	0.011	0.069	0.016	1.732

APPENDIX D PROOFS

D.1 Proof of Theorem 4

Proof

$$\begin{aligned} & KL(p_X^*(\mathbf{x})||\mathcal{N}(0, I_n)) \\ &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\mathcal{N}(0, I_n)} \right) \right] \\ &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\prod_{i=1}^n p_{X_i}^*(x)} \frac{\prod_{i=1}^n p_{X_i}^*(x)}{\mathcal{N}(0, I_n)} \right) \right] \\ &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\prod_{i=1}^n p_{X_i}^*(x)} \right) \right] + \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{\prod_{i=1}^n p_{X_i}^*(x)}{\prod_{i=1}^n \mathcal{N}(0, 1)} \right) \right] \\ &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\prod_{i=1}^n p_{X_i}^*(x)} \right) \right] + \mathbb{E}_{p_X^*(\mathbf{x})} \left[\sum_{i=1}^n \log \left(\frac{p_{X_i}^*(x)}{\mathcal{N}(0, 1)} \right) \right] \\ &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\prod_{i=1}^n p_{X_i}^*(x)} \right) \right] + \sum_{i=1}^n \mathbb{E}_{p_{X_i}^*(x)} \left[\log \left(\frac{p_{X_i}^*(x)}{\mathcal{N}(0, 1)} \right) \right] \\ &= KL(p_X^*(\mathbf{x})||\prod_{i=1}^n p_{X_i}^*(x)) + \sum_{i=1}^n KL(p_{X_i}^*(x)||\mathcal{N}(0, 1)) \end{aligned}$$

□

TABLE C.3: Results of Generalized Shapiro-Wilk test for multivariate normality on the representations of datasets under Glow. See Section 6.1 for the explanation of dataset names. For each dataset, we randomly select 2000 inputs for normality test. The larger W and p are, the more Gaussian-like the distribution is. When $p \geq 0.05$, there is no evidence to reject the normality hypothesis. In our experiments, ID representations under all models manifest strong normality. For *Category I* problems, all OOD representations except for SVHN vs Constant manifest normality. Some OOD representation (e.g., Uniform, ImageNet32) even has a higher p -value than ID data (CelebA).

ID	Input(ID/OOD)	Category	W	p -value	
Fashion.	Fashion.	-	0.9996	0.9479	
	Constant	I	0.9992	0.5872	
	Constant-C(0.1)	I	0.9995	0.9212	
	MNIST	I	0.9985	0.0733	
	MNIST-C(10.0)	I	0.9991	0.4114	
	notMNIST	I	0.9989	0.2337	
	notMNIST-C(0.005)	I	0.9993	0.6411	
SVHN	SVHN	-	0.9993	0.6227	
	Constant	I	0.9911	$9.6e-10$	
	Constant-C(0.1)	I	0.9992	0.5442	
	Uniform	II	0.9992	0.5273	
	Uniform-C(0.008)	II	0.9993	0.6203	
	CelebA	II	0.9336	$\ddagger 2.2e-16$	
	CelebA-C(0.08)	I	0.9993	0.6503	
	CIFAR-10	II	0.99429	$5.7e-07$	
	CIFAR-10-C(0.12)	I	0.9995	0.8838	
	CIFAR-100	II	0.9528	$\ddagger 2.2e-16$	
	CIFAR-100-C(0.12)	I	0.9985	0.0760	
	ImageNet32	II	0.8618	$\ddagger 2.2e-16$	
	ImageNet32-C(0.07)	I	0.9670	$\ddagger 2.2e-16$	
	CIFAR-10	-	0.9995	0.9064	
CIFAR-10	Constant	I	0.9992	0.5512	
	Constant-C(0.1)	I	0.9991	0.4725	
	Uniform	I	0.70958	$\ddagger 2.2e-16$	
	Uniform-C(0.02)	II	0.99931	0.6964	
	CIFAR-100	I	0.9994	0.8426	
	CelebA	I	0.9987	0.1390	
	CelebA-C(0.3)	I	0.9994	0.7960	
	ImageNet32	I	0.9977	0.0048	
	TinyImageNet	I	0.9995	0.3092	
	SVHN	I	0.9989	0.2532	
	SVHN-C(2.0)	I	0.9989	0.2547	
	CelebA	-	0.9992	0.6064	
	CelebA	Constant	I	0.9989	0.2605
		Constant-C(0.1)	I	0.9984	0.7184
Uniform		II	0.9993	0.6922	
Uniform-C(0.012)		II	0.9992	0.5815	
CIFAR-10		I	0.9992	0.5953	
CIFAR-100		I	0.9990	0.3313	
ImageNet32		I	0.9993	0.6410	
ImageNet32-C(0.2)		I	0.9990	0.3676	
SVHN		I	0.9991	0.4351	
SVHN-C(1.8)		I	0.9990	0.3600	

Then we apply Theorem 4 on each $D_g^i[p_{\bar{X}_i}^*]$ and have

$$\begin{aligned}
 & KL(p_{\bar{X}_i}^*(\mathbf{x}) || \mathcal{N}(0, I_l)) \\
 &= KL(p_{\bar{X}_i}^*(\mathbf{x}) || \prod_{j=1}^l p_{\bar{X}_{ij}}^*(x)) + \sum_{j=1}^l KL(p_{\bar{X}_{ij}}^*(x) || \mathcal{N}(0, 1))
 \end{aligned} \tag{18}$$

Finally, combining Equation (9) and 18 we can obtain Equation (10). \square

APPENDIX E SUMMARY OF OUR WORK

The flowchart in Figure E.1 can help readers to have an overview of our work.

D.2 Proof of Theorem 6

Proof We can use the similar deduction in Theorem 4 and get Equation (9).

$$\begin{aligned}
 & KL(p_X^*(\mathbf{x}) || \mathcal{N}(0, I_n)) \\
 &= \mathbb{E}_{p_X^*(\mathbf{x})} \left[\log \left(\frac{p_X^*(\mathbf{x})}{\prod_{i=1}^k p_{\bar{X}_i}^*(\mathbf{x})} \frac{\prod_{i=1}^k p_{\bar{X}_i}^*(\mathbf{x})}{\mathcal{N}(0, I_n)} \right) \right] \\
 &= I_g[p_X^*] + D_g[p_X^*]
 \end{aligned}$$

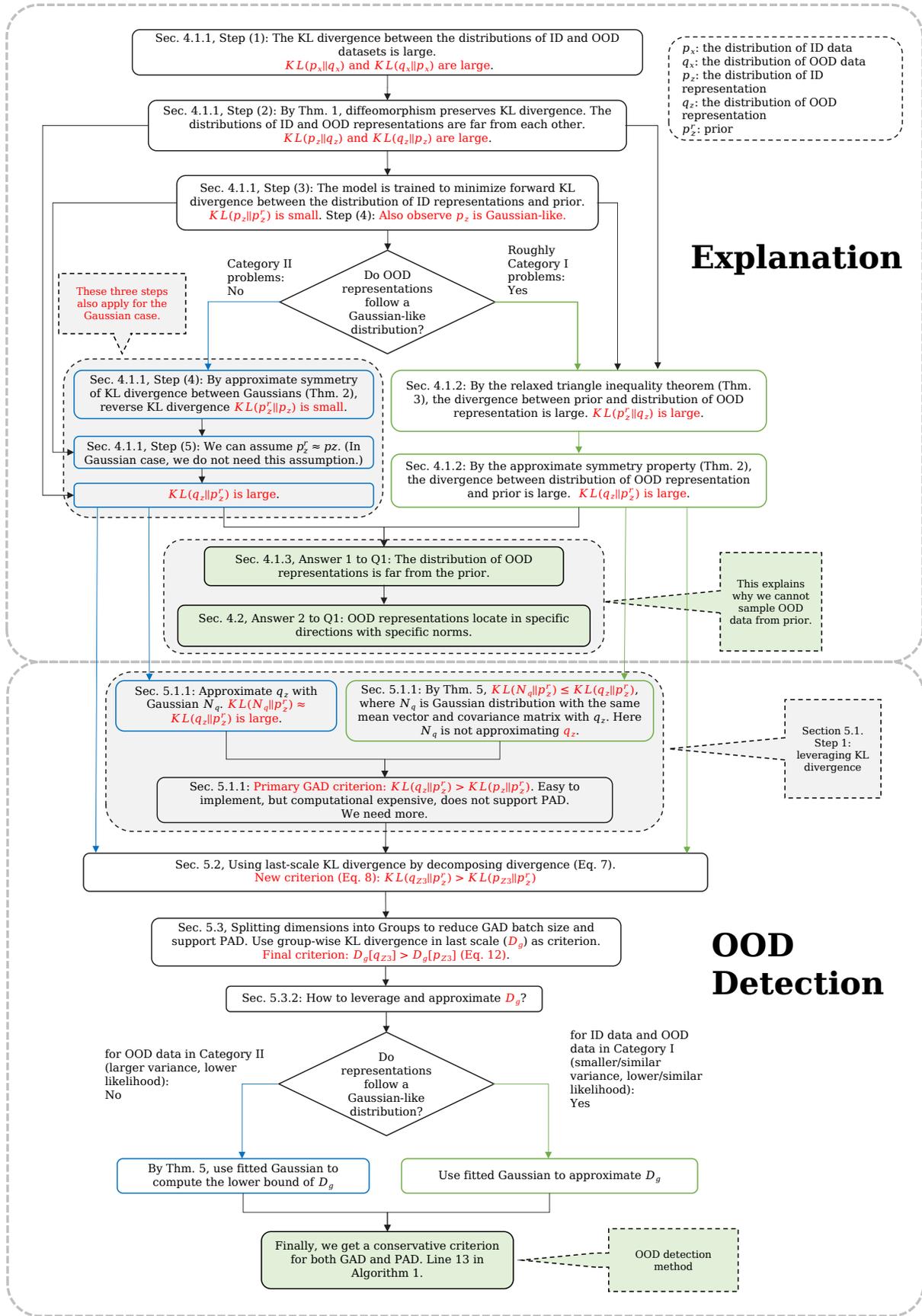


Fig. E.1: Overview of our method. The top half explains why we cannot sample OOD data from the flow-based model. The bottom half is for OOD detection method. Green lines are for the Gaussian case. Blue lines are for the non-Gaussian case. Please also refer to Figure 1 and 2 in the main text.

APPENDIX F BENCHMARKS

Here we briefly introduce the benchmarks used in our experiments:

- 1) Constant: The Constant dataset consists of images with all pixels equal to the same constant $C \sim U\{0, 255\}$.
- 2) Uniform: The Uniform dataset consists of images with each pixel sampled independently from $U\{0, 255\}$.
- 3) MNIST [69]: MNIST is a dataset of handwritten digits.
- 4) FashionMNIST [70]: FashionMNIST is a dataset of images of clothes and shoes.
- 5) notMNIST [71]: notMNIST is a dataset of fonts and extracting glyphs similar to MNIST.
- 6) KMNIST [72]: KMNIST is a dataset of Japanese characters.
- 7) Omniglot [73]: Omniglot is a dataset of handwritten digits of a set of alphabets.
- 8) CIFAR-10/100 [74]: CIFAR-10/100 are datasets of natural images including animals and vehicles.
- 9) SVHN [75]: SVHN is a dataset of street view housing numbers.
- 10) CelebA [76]: CelebA is a dataset of face images of celebrities.
- 11) TinyImageNet [77]: TinyImageNet is a subset of ImageNet.
- 12) ImageNet32 [78]: Imagenet32 is a dataset of small images called the down-sampled version of Imagenet.
- 13) LSUN [79]: LSUN is a dataset of scene categories including bedrooms, classroom, etc.

All datasets are resized to $32 \times 32 \times 3$ for consistency. The size of each test dataset is fixed to 10,000 for comparison. For grayscale datasets of size $28 \times 28 \times 1$, we replicate channels and pad zeros around images. We use the same method to process LSUN as the baseline method GOD2KS [34]. See Figure K.30 in the supplementary material for example images of different datasets.

APPENDIX G MODEL DETAILS

We use the released model (checkpoints) by the author of the baseline to conduct experiments if possible. Otherwise, we train the model by ourselves.

The authors of GAD baseline method (*i.e.*, Ty-test) reimplement Glow with PyTorch and release only one model checkpoint trained on CIFAR-10 [22]. We use their model for CIFAR-10 vs others. For other problems, we train the official Glow model [64] by ourselves. For VAE, we train convolutional VAE and use sampled representation for all problems. Among baseline methods, only the authors of L_{last} released their model checkpoints. We use their checkpoints to produce results on problems not evaluated in the original paper.

The Glow model consists of three stages, each containing 32 coupling layers with width 512. After each stage, the latent variables are split into two parts. One half is treated as the final representations and another half is processed by the next stage. We use additive coupling layers for grayscale datasets and CelebA and use affine coupling layers for SVHN and CIFAR-10. We find no difference between

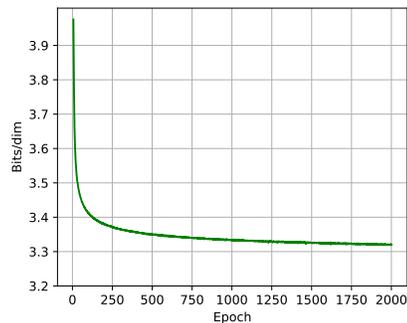


Fig. G.1: The training curve of Glow on CelebA32.

these two coupling layers for OOD detection. All priors are standard Gaussian distribution except for CIFAR-10, which has learned mean and diagonal covariance. All models are trained using Adamax optimization method with a batch size of 64. The learning rate is increased from 0 up to 0.001 in the first 10 epochs and keeps invariable in the remaining epochs. Flow-based models are resource consuming. We train Glow on FashionMNIST/SVHN/CelebA32 for 130/390/2000 epochs, respectively. The training curve of Glow on CelebA32 is shown in Figure G.1. We have also conducted experiments using the checkpoint released by OpenAI [64] for CIFAR-10 vs others. The results are similar.

For VAE, we use convolutional architecture in the encoder and decoder. The encoder consists of three $4 \times 4 \times 64$ convolution layers. On top of convolutional layers, two dense layer heads output the mean $\mu(\mathbf{x})$ and the standard variance $\sigma(\mathbf{x})$ respectively. The decoder has the mirrored architecture as the encoder. All activations are LeakyReLU with $\alpha = 0.3$. For FashionMNIST, SVHN, and CIFAR-10, we use 8-, 16- and 32-dimensional latent space, respectively. Models are trained using Adam without dropout. The learning rate is 5×10^{-4} with no decay.

Details of Baseline Methods.

- 1) \mathcal{S} . The authors of \mathcal{S} modified the official Glow model by using zero padding and removing ActNorm layer [28]. They did not explain the reason for such modification. In principle, such modification to models should not affect the baseline method. Since the authors did not release their source code and model checkpoint, we reimplement \mathcal{S} method using the official Glow model [84] for those problems not evaluated in [28]. We also use FLIF [88] as the compressor, which is the best compressor in [28]. We find the performance of \mathcal{S} degenerates on the official Glow model.
- 2) L_{last} . We use the model checkpoints released by the authors of L_{last} for results not reported in the original publication [29].
- 3) DoSE. The authors of DoSE did not release their source code and model. We reimplement DoSE_{svm} on the official Glow model for problems not evaluated in [30]. We use the same parameters as the original paper. The performance of DoSE on the official Glow model trained on SVHN is slightly better than the results in [30]. We also find the performance of DoSE degenerates severely on CelebA32 vs CIFAR-10/100 compared with

the results reported in their original publication (below 75% AUROC). We did not present these results in the main text.

APPENDIX H DISCUSSION

Normality of Representations. The normality of ID and OOD representation facilitates our theoretical analysis and OOD detection algorithm on flow-based model. In our experiments, we find that the normality of OOD representation is a widely existing phenomenon under flow-based models. We are investigating the underlying reason. Most importantly, our method performs even better on *Category II* problems, although the criterion computes the lower bound of the KL divergence.

Both Flow-based model and VAE are trained to minimize KL divergence between p_Z and prior. However, it seems that the normality of ID/OOD representations is a characteristic of flow-based model.

In principle, we can construct latents following any distribution and decode these latents back to data space to construct an OOD dataset. Note that such manipulation does not necessarily make our OOD detection method fail, although it can violate the normality hypothesis. Besides, such manipulation is much more difficult to conduct than the data manipulations presented in this paper because we need the model parameters.

PAD Results. From the results of three Glow models trained on SVHN/CelebA/CIFAR-10 (Table 1 in the main text), we can see that the more high-quality generated images are, the better performance our method can achieve. This is consistent with our theoretical analysis. Our method has the a solid theoretical foundation. We believe that our method can achieve better performance with the increasing success of flow-based model in the future.

Limitation. Our method requires the model to capture the distribution of training data. For example, Glow trained on CIFAR-10 does not generate meaningful images. Thus, the KL divergence between the distributions of ID representation and prior is not small enough. So our theoretical analysis does not apply to this problem well. This is why our method does not achieve high AUROC on CIFAR-10 vs CIFAR-100. We think it is hard to achieve high AUROC with a model which does not succeed in generating meaningful images (see Figure K.31 in the supplementary material). In such a case, the model generates “OOD data” that differ from the training set. Besides, the similarity between CIFAR-10 and CIFAR-100 also brings obstacles to OOD methods.

There are two possible solutions for the most challenging problem CIFAR-10 vs CIFAR-100. The first one is to improve the model. Modeling data is a long-standing goal of unsupervised learning [57]. Up to now, it is still hard to generate high-quality CIFAR-10-like images using unconditional flow-based models. The second possible solution is to use a more sensitive criterion to estimate KL divergence or dependence. We leave this direction as future work.

Ty-test applies to flow-based model, VAE, and autoregressive model. Our method applies to models which learn independent or disentangled representations [59], [60], [61], [89], [90], [91], [92], not including auto-regressive model.

Other Baselines.

In this paper, we mainly choose recently proposed methods applicable to flow-based model as baselines. We did not choose WAIC [20] whose results could not be reproduced by Nalisnick *et al* [22]. We did not choose the likelihood ratios method [35] as the baseline either for several reasons. First, in [28], Serrà *et al.* interpret their method \mathcal{S} as a likelihood-ratio test statistic and achieve better performance than likelihood ratios. Second, the authors of the likelihood ratios method [35] did not report results on flow-based models. So we choose method \mathcal{S} rather than likelihood ratios as baseline. Finally, we choose SOTA method DoSE as baseline, which is better than $\log p(\mathbf{x})$, $\log p(\mathbf{z})$, WAIC, and likelihood ratio as reported in [30].

In [34], the authors propose GOD2KS mainly for GAD. In the appendix of [34], the authors also use data augmentation to support PAD using GOD2KS. However, they only report a few PAD results based on RealNVP. The PAD results of GOD2KS on CIFAR-10 vs SVHN/CelebA/LSUN with RealNVP are 85%/57%/46% AUROCs, respectively. Our method achieves 82.6%/85.2%/99.2% AUROCs on the same three problems with Glow. Due to this situation, we did not use GOD2KS as PAD baseline.

Just after we receive the first round of review comments, Osada *et al.* propose PRE method [93], which uses reconstruction error and typicality-based penalty to perform point-wise anomaly detection with flow-based model. PRE uses the original Equation (4) in [93] as anomaly score. The larger the score, the more likely the input is OOD. We did not choose PRE as baseline for two reasons. First, the authors of PRE do not use as many dataset compositions in evaluation as ours. They also use ID datasets of different sizes. It is hard to compare two methods in this situation. Second, similar to Annulus Method, PRE can also be attacked by data manipulation **M1** (rescaling representations). For each OOD dataset $S = \{\mathbf{x}\}$, we can use data manipulation **M1** (see Subsection 3.1) to construct an OOD dataset $S' = \{\mathbf{x}' = f^{-1}(\sqrt{d} \frac{f(\mathbf{x})}{|f(\mathbf{x})|}) | \mathbf{x} \in S\}$ (see Figure K.4 and Figure K.5). For input $\mathbf{x}' \in S'$, the representation $\mathbf{z}' = f(\mathbf{x}')$ locates in the typical set annulus of prior precisely and Equation (4) in [93] equals 0. This can make PRE method achieve near 0 AUROC.

Other Comparisons. Explicit generative models, including autoregressive models, flow-based models, and VAEs, can provide users with likelihoods (or lower bound). An ideal explicit generative model should: 1) generate new data from the training data distribution and 2) provide likelihood indicating the confidence of whether the data belongs to the training data distribution. Implicit generative models (*i.e.*, GAN) do not produce likelihood, so these two kinds of models are under different settings. Commonly, explicit generative models are compared together in anomaly detection publications. All the baseline methods applying to flow-based model are compared with explicit generative models in evaluation. (*e.g.*, [20], [23], [28], [29], [30], [35]).

We notice that the existing hybrid model [94] achieves better performance on leave-one-out setting on MNIST⁷. It is unfair to compare a flow-based method with a hybrid

7. The author of [94] did not report experimental results on cross-dataset problems.

model combining explicit and implicit generative models. For example, existing work has shown that the combination of GAN and flow-based model can improve the quality of the generated images of flow-based model. For example, in Flow-GAN [95], flow-based model is used to avoid mode collapse. Adversarial training is used to improve the image quality of flow-based model. The method is to sample noise z from the typical set of prior and use a discriminator to distinguish $f^{-1}(z)$ and training data. In our data manipulation M1 (rescaling representations to the typical set), we find that Glow trained by maximum likelihood estimation cannot expel OOD representation from the typical set of prior (see Subsection 3.1). In Flow-GAN, adversarial training tends to compel Glow (f^{-1}) to map latents in the typical set of prior to in-distribution images. Importantly, our theoretical analysis also applies to flow-GAN whose loss function includes the basic loss function of Glow. We did not conduct experiments on Flow-GAN [95] because Flow-GAN uses old flow-based models NICE [36] and RealNVP [14]. Besides, adversarial training would affect the divergences in flow-based model. We will explore the properties of hybrid models which combine the SOTA flow-based model and GAN in the future.

We conduct comprehensive experiments to evaluate our method on different dataset compositions falling into *Category I* and *II*. Flow-based models have different behaviors for these two categories of problems. Commonly, one OOD detection method’s performance may vary on different problems. Several existing OOD detection methods have been evaluated with very few datasets (*e.g.*, only CIFAR-10) as the training dataset. We did not compare our method with such methods for two reasons. First, there is no result reported on more problems. Second, our method requires the model to succeed in modeling the training dataset. Unluckily, flow-based model is not as successful as other datasets on CIFAR-10. This affects our method. We think it is not comprehensive if using only CIFAR-10 as training dataset in evaluation.

Researchers also propose ensembling algorithms for anomaly detection [96], which is orthogonal to our method. We plan to explore this direction in the future.

Models. We did not conduct more experiments on flow-based models with various architectures. In principle, a more expressive model can make the forward KL divergence smaller, and our method can benefit more.

Our theoretical analysis relies on the assumption that the KL divergence between the distributions of ID and OOD representations is large (see Section 4.1.1 in the main text). So our analysis does not apply to VAE and autoregressive models directly. According to the Brouwer Invariance of Domain Theorem [97], R^n cannot be homeomorphic to R^m if $n \neq m$. The Brouwer Invariance of Domain Theorem also implies that there is no dead neuron in flow-based model. Otherwise, we can construct diffeomorphism from high to low dimensional space. For VAE, a high-dimensional latent space may contain nearly dead neurons. This may reduce the performance of our method. We did not conduct experiments on other VAE variations, *e.g.*, β -VAE [92], FactorVAE [59], β -TCVAE [60], and DIP-VAE [61]. These variations add more regularization strength on disentanglement and have more independent representations than vanilla VAE [62]. We will conduct experiments on larger VAE models and variations

in the future.

Finally, we use the model checkpoints released by baselines as long as possible. These released models should be tuned elaborately for their methods, so our method benefits less from fine-tuning.

APPENDIX I MORE RELATED WORK

OOD Detection. In [29], Schirrmeister *et al.* find the likelihood contributed by the last scale of Glow (L_{last}) is better than $\log p(x)$ for PAD. Their method relies on the decomposition of the likelihood (original Equation (3) in [29]). Such decomposition requires that the split two parts at each stage of Glow are independent. This may not hold for OOD data due to covariate shift. Experimental results show that the last-scale log-likelihood of OOD data may be larger than 0 (See Figure K.8 in the supplementary material). So the criterion used by L_{last} should not be explained as likelihood for OOD data. Finally, as shown in Table 1, L_{last} is also affected by data manipulation.

Theoretical Analysis. Previous works [37], [98] analyze the training objective of flow-based model in KL divergence form. We apply the property of diffeomorphism to investigate the divergences between distributions in flow-based models in the setting of OOD detection. We also propose new theorems on the properties of KL divergence between Gaussian distributions for further analysis. Currently, there is no similar work on the properties of KL divergence between Gaussian distributions. Theorems 2, and 3 can be used as basic theorems in machine/deep learning and information theory. For example, after we post the manuscript containing the proofs of Theorems 2 and 3 on Arxiv [45], the relaxed triangle inequality (Theorem 3) has been used in constrained variational policy optimization for safe reinforcement learning [99].

Other Approximator and Divergence Estimation. In principle, GMM can approximate a target density better than a single Gaussian distribution [100]. We tried to use GMM to approximate the distribution of representations and use Monte Carlo sampling to estimate the KL divergence between GMMs. The results show GMM is worse than using Gaussian distribution for OOD detection. The reasons are twofold: a) it is inappropriate to use GMM for modeling ID representations that follow a Gaussian-like distribution. b) the batch size is too small (usually $5 \sim 10$) to estimate the parameters of GMM.

We also tried the SOTA ϕ -divergence estimation method applicable for VAE, *i.e.* RAM-MC [52]. Results show that RAM-MC can also be affected by data manipulation M2 (adjusting contrast, see Section 3.1)⁸.

OOD Sampling. In this paper, we sample OOD data using the fitted Gaussian from OOD representations to verify that OOD representations reside in specific directions. Such directions can be partially captured by the sample mean and sample covariance. In [49], the authors sample noise $\varepsilon \sim \mathcal{N}(\tilde{\mu}, I)$ and generate OOD data $f^{-1}(\varepsilon)$ using flow-based model, where $\tilde{\mu}$ is the sample mean of OOD representa-

⁸. This does not prove that RAM-MC is not applicable to general-purpose divergence estimation.

tions. They did not use the sample covariance of OOD representations. Their manuscript is released contemporaneously with the first edition of this paper [101]. Some researchers have explored other OOD sampling methods. Sinha *et al.* use various operations including Jigsaw, Stitching on normal images to generate negative data [102]. These negative data can be used to help GAN to improve generation quality and OOD detection ability. However, their performance of OOD detection on CIFAR-10 vs others is worse than our method. Dionelis *et al.* propose to generate samples [103] on the boundary of the support of data distributions which is learned by flow-based model. Their method does not modify the parameters of flow-based model and has the same OOD detection ability as the original flow-based model. In this paper, we sample OOD data in order to verify that OOD representations reside in specific directions that can be partially captured by the sample mean and covariance of representations. We will explore more work on OOD sampling in the future.

Local Pixel Dependence. In [23], Kirichenko *et al.* reshape the representations of flow-based models to the original input shape ($32 \times 32 \times 3$) and analyze the induction biases of flow-based model. Their work reveals the reshaped representation manifests local pixel dependence. Our work show that the representations with a raw shape ($4 \times 4 \times 48$) also manifest local pixel dependence.

OOD Detection With Auxiliary Data. OOD detection can be improved with the help of an auxiliary outlier dataset. In [29], Schirrmeyer *et al.* improve likelihood-ratio-based method by the help of a huge outlier dataset (80 Million Tiny ImageNet). This is not unsupervised learning due to the exposure to outliers in training as like [68]. Besides, the huge outlier dataset includes almost all the image classes in the testing phase. We did not compare with such methods due to different problem settings.

Classification of Problems. We classify OOD problems into *Category I* and *II* according to the variance and likelihoods of datasets. This criterion is roughly similar to the complexity used in [28]. See Figure K.32 in the supplementary material for details.

APPENDIX J

MORE EXPERIMENTAL RESULTS

J.1 GAD Results on Glow

FashionMNIST vs Others. Table J.1 shows the GAD results of Glow trained on FashionMNIST. The results of baselines are referenced from [22], in which the authors use bootstrap method to establish thresholds. We use a low false positive rate to establish a threshold t (see Algorithm 1) and then compute the corresponding true positive rate with t . Take FashionMNIST vs MNIST with batch size $m = 2$ as example. We find the threshold t corresponding to false positive rate 0.01 which is the second lowest one among all the baselines (column 2 in Table J.1). Then we use t to compute the corresponding true positive rate 0.43 ± 0.02 . For larger batch sizes ($m = 10, 25$), we set the thresholds t , achieving a most rigorous 0 false positive rate.

SVHN/CIFAR-10/CelebA vs Others. Table J.2, J.3 and J.4 shows numerical GAD results corresponding to Figure 3

in the main text. Table J.5 shows GAD results with smaller batch sizes 2 and 4.

CelebA vs CIFAR-10/100 are challenging for Ty-test. In principle, if the train and test split of ID dataset have coinciding likelihoods, the worst AUROC of Ty-test should be around 50%. But Ty-test only achieves 1.7% and 2.9% AUROCs on these two problems. The reasons are two-fold. First, CIFAR-10/100 have coinciding likelihoods with CelebA. Please see Figure K.10 for details. Second, we find it is hard to make the likelihood distributions of CelebA train and test split fit very well on the official Glow model even within 2,000 epochs (see Figure G.1 for training curve). The likelihoods of CIFAR-10/100 are closer to CelebA train set than the CelebA test set. This misleads Ty-test to make wrong decisions (below 10% AUROC). This also makes Ty-test perform worse when batch size is larger because a larger batch size eliminates randomness (see Table J.2). On the contrary, our method is not affected by such possible underfitting or overfitting.

Comparison with GOD2KS. Table J.6 shows the comparison of our method with GOD2KS. Our method is better.

Mixture of OOD Datasets. Table J.7 shows the results of KLODS when OOD dataset is a mixture of two of the three datasets: SVHN, CelebA, and CIFAR-10.

Ablation Study. Table J.8 shows the results of ablation study. Except for CIFAR-10 vs ImageNet32-C(0.3), the order of performance is KLODS > KLOD > KLOD-all > Ty-test. The only exception is CIFAR-10 vs ImageNet32-C(0.3). Note that KLOD only applies to GAD.

One-vs-Rest. Table J.9 shows GAD Results (AUROC and AUPR in percentage) of Glow on One-vs-Rest on MNIST.

GlowGMM. Figure J.1 and Table J.10 shows the results of GAD on FashionMNIST. Our method achieves 100% AUROC on average when the batch size is 25. The baseline only reaches 22.7% AUROC. Recent works have improved the accuracy of conditional Glow on classification problems [26], [104]. However, as long as GlowGMM does not achieve 100% classification accuracy, the question proposed in the introduction remains.

Table J.11 shows the results of using $p(\mathbf{z})$ for one-vs-rest classification on FashionMNIST with GlowGMM. $p(\mathbf{z})$ is not a good criterion for OOD detection. For example, the AUROC for class 8 vs rest is only 55.5%.

Generating OOD data Using GlowGMM. In Section 5.1, we sample blurred OOD data with fitted Gaussian distribution. In GlowGMM, we can generate *high-quality* OOD images with fitted Gaussian distribution from OOD representations. Figure J.2(a) shows the generated images using noise sampled from each Gaussian component $\mathcal{N}_i(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2))$ ($1 \leq i \leq 10$) as prior. The i -th column corresponds to the i -th Gaussian \mathcal{N}_i . Figure J.2(b) shows the generated images using the similar operation in Section 5.1.1. For each i , we compute the representations $\{\mathbf{z}\}$ of the $((i+1)\%10)$ -th class and normalize them under $\mathcal{N}_i(\boldsymbol{\mu}_i, \text{diag}(\boldsymbol{\sigma}_i^2))$ as $\mathbf{z}' = (\mathbf{z} - \boldsymbol{\mu}_i)/\boldsymbol{\sigma}_i$. We use the normalized representation $\{\mathbf{z}'\}$ to fit a Gaussian distribution $\tilde{\mathcal{N}}_i(\tilde{\boldsymbol{\mu}}_i, \tilde{\boldsymbol{\Sigma}}_i)$. Then We sample $\varepsilon_{i'} \sim \tilde{\mathcal{N}}_i(\tilde{\boldsymbol{\mu}}_i, \tilde{\boldsymbol{\Sigma}}_i)$ and unnormalize $\varepsilon_{i'}$ back using parameters of the i -th component as $\varepsilon_{i'} \cdot \boldsymbol{\sigma}_i + \boldsymbol{\mu}_i$. Finally, we compute $f^{-1}(\varepsilon_{i'} \cdot \boldsymbol{\sigma}_i + \boldsymbol{\mu}_i)$ to generate new images. As shown in Figure J.2(b), we can generate almost high quality images of the $((i+1)\%10)$ -th class from the fitted

TABLE J.1: GAD Results on Glow trained on FashionMNIST. The ID column reflects FPR (ideally should be 0) and the MNIST and notMNIST columns are TPR (ideally should be 1). The results of baselines are referenced from [22]. Notable failures (under 0.5 TPR) are underlined.

Method	$m=2$			$m=10$			$m=25$		
	ID	MNIST	notMNIST	ID	MNIST	notMNIST	ID	MNIST	notMNIST
Ty-test	0.02±0.01	0.14±0.10	0.08±0.04	0.02±0.02	1.00±0.00	0.69±0.11	0.01±0.00	1.00±0.00	1.00±0.00
t -test	0.01±0.00	0.08±0.00	0.06±0.00	0.01±0.00	1.00±0.00	0.67±0.01	0.01±0.00	1.00±0.00	0.99±0.00
KS-Test	0.00±0.00	0.00±0.00	0.00±0.00	0.01±0.00	1.00±0.00	0.61±0.01	0.00±0.00	1.00±0.00	0.98±0.01
Max Mean Dis.	0.05±0.02	0.17±0.06	0.04±0.03	0.02±0.02	0.63±0.12	0.37±0.24	0.04±0.04	1.00±0.00	1.00±0.00
Kern. Stein Dis.	0.05±0.05	0.16±0.14	0.01±0.01	0.01±0.01	0.21±0.11	0.01±0.00	0.02±0.03	0.76±0.21	0.00±0.00
Annulus Method	0.01±0.01	0.00±0.00	0.96±0.03	0.02±0.00	0.00±0.00	1.00±0.00	0.03±0.03	0.00±0.00	1.00±0.00
KLODS	0.01±0.00	0.43±0.02	0.95±0.00	0.00±0.00	1.00±0.00	1.00±0.00	0.00±0.00	1.00±0.00	1.00±0.00

TABLE J.2: GAD Results (AUROC and AUPR in percentage) of KLODS and Ty-test on Glow with batch sizes 5 and 10. The higher the better. The performance of Ty-test on CelebA vs CIFAR-10/100 decreases when the batch size is larger. See our explanation in Section 6.2.1 in the main text.

ID↓	OOD↓	Batch size→		$m=5$				$m=10$			
		Method→		KLODS		Ty-test		KLODS		Ty-test	
		Metric→		AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
Fash.	Constant	100.0±0.0	100.0±0.0	42.1±0.3	42.1±0.2	100.0±0.0	100.0±0.0	41.7±0.5	41.9±0.2		
	MNIST	99.8±0.0	99.8±0.0	97.6±0.1	95.8±0.5	100.0±0.0	100.0±0.0	99.7±0.1	99.6±0.1		
	MNIST-C(10.0)	100.0±0.0	100.0±0.0	88.2±0.3	81.8±0.2	100.0±0.0	100.0±0.0	95.8±0.5	93.5±1.2		
	notMNIST	100.0±0.0	100.0±0.0	77.5±0.3	74.6±0.4	100.0±0.0	100.0±0.0	87.1±0.2	85.4±0.4		
	notMNIST-C(0.005)	100.0±0.0	100.0±0.0	25.0±0.6	35.8±0.2	100.0±0.0	100.0±0.0	23.8±0.4	35.5±0.1		
SVHN	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.008)	100.0±0.0	100.0±0.0	13.5±0.5	33.0±0.1	100.0±0.0	100.0±0.0	11.1±0.5	32.6±0.1		
	CelebA	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CelebA-C(0.08)	99.7±0.0	99.7±0.0	50.7±0.7	47.0±0.3	100.0±0.0	100.0±0.0	55.2±0.4	49.1±0.3		
	CIFAR-10	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CIFAR-10-C(0.12)	97.0±0.2	97.4±0.2	31.6±0.5	37.9±0.2	99.3±0.1	99.4±0.1	25.0±0.3	35.6±0.1		
	CIFAR-100	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CIFAR-100-C(0.12)	96.9±0.1	97.3±0.1	35.3±0.5	39.4±0.2	98.9±0.3	99.0±0.3	27.2±0.8	36.3±0.2		
	ImageNet32	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	ImageNet32-C(0.07)	99.8±0.0	99.8±0.0	45.5±0.9	46.0±0.5	100.0±0.0	100.0±0.0	42.1±0.7	44.1±0.5		
LSUN	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0			
LSUN-C(0.06)	99.9±0.0	99.0±0.0	42.8±0.6	42.5±0.3	100.0±0.0	100.0±0.0	42.3±0.5	42.2±0.1			
CIFAR-10	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.02)	100.0±0.0	100.0±0.0	11.5±0.0	32.9±0.0	100.0±0.0	100.0±0.0	9.3±0.0	32.5±0.0		
	CelebA	99.2±0.1	99.4±0.1	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CelebA-C(0.3)	84.3±0.3	84.4±0.4	28.4±0.5	36.7±0.2	94.5±0.3	94.7±0.3	23.5±0.5	35.2±0.1		
	ImageNet32	90.0±0.2	92.1±0.1	99.2±0.1	99.3±0.1	95.0±0.4	96.2±0.2	100.0±0.0	100.0±0.0		
	ImageNet32-C(0.3)	72.0±0.3	72.6±0.4	40.9±0.4	43.2±0.2	74.3±0.6	74.8±0.8	32.0±0.7	38.5±0.3		
	SVHN	97.6±0.2	97.8±0.2	98.6±0.1	98.4±0.1	99.8±0.0	99.8±0.0	99.9±0.1	99.9±0.1		
	SVHN-C(2.0)	100.0±0.0	100.0±0.0	33.5±0.4	61.0±0.2	100.0±0.0	100.0±0.0	27.2±0.5	58.2±0.1		
	LSUN	100.0±0.0	100.0±0.0	99.9±0.0	99.9±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
LSUN-C(0.3)	90.0±0.2	90.8±0.2	52.2±0.8	48.8±0.4	91.2±0.2	92.0±0.2	56.6±0.4	51.3±0.3			
CelebA	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.012)	100.0±0.0	100.0±0.0	36.2±0.7	39.6±0.2	100.0±0.0	100.0±0.0	30.9±0.7	37.9±0.2		
	CIFAR-10	99.6±0.0	99.6±0.0	7.2±0.2	31.4±0.0	100.0±0.0	100.0±0.0	1.7±0.1	30.8±0.0		
	CIFAR-100	99.8±0.0	99.8±0.0	9.5±0.3	31.8±0.1	100.0±0.0	100.0±0.0	2.9±0.2	30.9±0.0		
	ImageNet32	100.0±0.0	100.0±0.0	78.1±0.4	85.6±0.3	100.0±0.0	100.0±0.0	83.9±0.4	89.6±0.2		
	ImageNet32-C(0.2)	100.0±0.0	100.0±0.0	26.0±0.3	36.4±0.2	100.0±0.0	100.0±0.0	18.2±0.3	33.8±0.0		
	SVHN	100.0±0.0	100.0±0.0	78.7±0.3	73.3±0.9	100.0±0.0	100.0±0.0	86.6±0.8	83.3±1.4		
	SVHN-C(1.8)	100.0±0.0	100.0±0.0	3.5±0.2	31.0±0.0	100.0±0.0	100.0±0.0	0.5±0.1	30.7±0.0		
	LSUN	100.0±0.0	100.0±0.0	65.4±0.3	64.3±0.2	100.0±0.0	100.0±0.0	71.2±0.3	70.0±0.5		
average		98.1	98.2	64.6	69.0	98.8	98.9	64.0	68.4		

Gaussian. These results verify that OOD representations reside in specific directions that can be characterized by the mean and covariance matrix of OOD representations. We did not conduct more experiments on OOD sampling because it is beyond the scope of this paper.

J.2 GAD Results on VAE

Figure J.3 and Table J.12 show GAD results on convolutional VAE on problems FashionMNIST/SVHN/CIFAR10 vs others.

Table J.13 shows the GAD results on CIFAR10 vs CIFAR100/ImageNet32.

Table J.14 shows the results of using reconstruction probability $E_{z \sim q_\phi}[\log p_\theta(x|z)]$ for OOD detection in VAE.

J.3 PAD results on Glow

Table J.15 shows PAD results (AUROC in percentage) of Joint confidence loss method, ODIN, Joint confidence loss method+ODIN, DoSE and KLODS.

TABLE J.3: GAD Results (AUROC and AUPR in percentage) of KLODS and Annulus Method (R) (reimplementation) on Glow with batch sizes 5 and 10. The higher the better. *Our reimplementation of Annulus Method achieves much better results than that reported in [22] (referenced by Table J.1).*

ID↓	OOD↓	Batch size→		$m=5$				$m=10$			
		Method→	KLODS		Annulus Method (R)		KLODS		Annulus Method (R)		
		Metric→	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	
Fash.	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	MNIST	99.8±0.0	99.8±0.0	98.6±0.0	98.7±0.0	100.0±0.0	100.0±0.0	99.9±0.0	99.9±0.0	100.0±0.0	100.0±0.0
	MNIST-C(10.0)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	notMNIST	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	notMNIST-C(0.005)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
SVHN	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	Uniform-C(0.008)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CelebA	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CelebA-C(0.08)	99.7±0.0	99.7±0.0	50.0±0.2	49.3±0.1	100.0±0.0	100.0±0.0	49.6±0.2	49.2±0.1	100.0±0.0	100.0±0.0
	CIFAR-10	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR-10-C(0.12)	97.0±0.2	97.4±0.2	59.6±0.0	58.2±0.0	99.3±0.1	99.4±0.1	63.5±0.1	62.0±0.5	100.0±0.0	100.0±0.0
	CIFAR-100	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR-100-C(0.12)	96.9±0.1	97.3±0.1	70.9±0.2	70.4±0.3	98.9±0.3	99.0±0.3	78.2±0.3	77.9±0.4	100.0±0.0	100.0±0.0
	ImageNet32	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32-C(0.07)	99.8±0.0	99.8±0.0	76.3±0.2	75.7±0.3	100.0±0.0	100.0±0.0	84.2±0.5	83.8±0.5	100.0±0.0	100.0±0.0
	LSUN	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	LSUN-C(0.06)	99.9±0.0	99.0±0.0	96.4±0.1	96.1±0.1	100.0±0.0	100.0±0.0	99.4±0.1	99.4±0.1	100.0±0.0	100.0±0.0
CIFAR-10	Constant	100.0±0.0	100.0±0.0	62.0±1.9	67.2±2.5	100.0±0.0	100.0±0.0	66.7±4.5	75.6±2.8	100.0±0.0	100.0±0.0
	Uniform	100.0±0.0	100.0±0.0	6.3±1.0	31.8±0.4	100.0±0.0	100.0±0.0	7.3±1.4	34.1±1.6	100.0±0.0	100.0±0.0
	Uniform-C(0.02)	100.0±0.0	100.0±0.0	54.4±2.1	63.3±2.6	100.0±0.0	100.0±0.0	61.3±2.9	71.8±2.7	100.0±0.0	100.0±0.0
	CelebA	99.2±0.1	99.4±0.1	36.8±1.9	48.3±2.2	100.0±0.0	100.0±0.0	45.4±4.9	60.6±4.1	100.0±0.0	100.0±0.0
	CelebA-C(0.3)	84.3±0.3	84.4±0.4	55.9±2.8	63.0±2.7	94.5±0.3	94.7±0.3	48.4±2.8	62.1±2.8	100.0±0.0	100.0±0.0
	ImageNet32	90.0±0.2	92.1±0.1	45.8±1.9	55.2±1.4	95.0±0.4	96.2±0.2	43.8±2.5	57.7±2.1	100.0±0.0	100.0±0.0
	ImageNet32-C(0.3)	72.0±0.3	72.6±0.4	47.2±2.6	56.1±2.5	74.3±0.6	74.8±0.8	51.3±3.5	64.5±2.7	100.0±0.0	100.0±0.0
	SVHN	97.6±0.2	97.8±0.2	47.3±1.6	56.9±2.4	99.8±0.0	99.8±0.0	49.4±2.4	63.0±2.4	100.0±0.0	100.0±0.0
	SVHN-C(2.0)	100.0±0.0	100.0±0.0	45.0±1.5	56.0±1.9	100.0±0.0	100.0±0.0	39.4±0.9	54.9±1.5	100.0±0.0	100.0±0.0
	LSUN	100.0±0.0	100.0±0.0	29.1±1.2	44.6±2.2	100.0±0.0	100.0±0.0	29.9±4.2	46.5±4.2	100.0±0.0	100.0±0.0
LSUN-C(0.3)	90.0±0.2	90.8±0.2	49.4±2.1	58.1±2.3	91.2±0.2	92.0±0.2	43.6±1.1	58.1±1.0	100.0±0.0	100.0±0.0	
CelebA	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	Uniform-C(0.012)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR-10	99.6±0.0	99.6±0.0	99.5±0.0	99.6±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR-100	99.8±0.0	99.8±0.0	99.6±0.0	99.6±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32-C(0.2)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	SVHN	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	SVHN-C(1.8)	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	LSUN	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
average		98.1	98.2	80.3	83.3	98.8	98.9	81.1	85.2		



Fig. J.1: GAD results(AUROC) on GlowGMM trained on FashionMNIST. Numerical results are shown in Table J.10 in the supplementary material.

TABLE J.4: GAD Results (EER in percentage) of KLODS, Ty-test and Annulus Method (Annulus.) on Glow with batch sizes 5 and 10. The lower the better.

ID↓	Method→	Batch size→		$m=5$			$m=10$		
		OOD↓	Method→	KLODS	Ty-test	Annulus.	KLODS	Ty-test	Annulus.
Fash.	Constant			0.0±0.0	54.2±0.4	0.0±0.0	0.0±0.0	53.9±0.8	0.0±0.0
	MNIST			0.9±0.2	5.6±0.3	6.1±0.4	0.0±0.0	1.2±0.2	1.6±0.3
	MNIST-C(10.0)			0.0±0.0	14.7±0.5	0.6±0.1	0.0±0.0	7.3±0.2	0.0±0.0
	notMNIST			0.0±0.0	29.5±0.3	0.0±0.0	0.0±0.0	20.8±0.3	0.0±0.0
	notMNIST-C(0.005)			0.0±0.0	44.5±0.4	0.0±0.0	0.0±0.0	39.4±0.5	0.0±0.0
SVHN	Constant			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	Uniform			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	Uniform-C(0.008)			0.0±0.0	79.1±0.6	0.0±0.0	0.0±0.0	81.7±0.4	0.0±0.0
	CelebA			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	CelebA-C(0.08)			8.0±0.3	43.2±0.3	50.3±0.3	2.3±0.2	40.2±0.8	50.2±0.3
	CIFAR-10			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	CIFAR-10-C(0.12)			24.0±0.5	64.1±0.3	43.2±0.4	18.7±0.5	70.9±0.3	40.6±0.6
	CIFAR-100			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	CIFAR-100-C(0.12)			23.3±0.4	61.0±0.3	34.6±0.3	19.4±0.5	68.2±0.2	29.2±0.7
	ImageNet32			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	ImageNet32-C(0.07)			9.3±0.2	53.7±0.6	30.8±0.3	4.0±0.3	56.4±0.6	23.8±0.4
	LSUN			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
LSUN-C(0.06)			0.8±0.1	54.1±0.4	9.6±0.3	0.0±0.0	42.3±0.5	3.7±0.5	
CIFAR-10	Constant			0.0±0.0	0.0±0.0	41.4±2.0	0.0±0.0	0.0±0.0	37.8±4.7
	Uniform			0.0±0.0	0.0±0.0	87.6±1.2	0.0±0.0	0.0±0.0	89.2±3.3
	Uniform-C(0.02)			0.0±0.0	85.9±0.5	48.0±2.3	0.0±0.0	87.5±1.0	41.3±3.4
	CelebA			3.9±0.2	0.9±0.1	59.2±1.8	0.4±0.1	0.0±0.0	55.1±4.1
	CelebA-C(0.3)			23.3±0.7	65.5±0.6	46.0±2.7	13.4±0.5	69.5±0.8	50.4±3.0
	ImageNet32			18.9±0.6	3.8±0.3	53.5±0.9	12.3±0.6	0.8±0.1	55.4±3.2
	ImageNet32-C(0.3)			35.3±0.6	57.1±0.5	51.4±2.8	32.6±1.3	64.0±0.7	50.9±2.7
	SVHN			8.0±0.2	5.6±0.2	52.2±1.8	2.5±0.4	1.1±0.2	52.5±2.4
	SVHN-C(2.0)			0.2±0.1	61.3±0.4	53.7±1.7	0.0±0.0	67.0±0.6	59.0±1.5
	LSUN			0.0±0.0	1.2±0.1	65.9±0.8	0.0±0.0	0.0±0.0	65.4±4.7
LSUN-C(0.3)			18.2±0.3	47.7±0.6	51.1±1.0	17.2±0.4	44.7±0.9	54.1±2.1	
CelebA	Constant			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	Uniform			0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0	0.0±0.0
	Uniform-C(0.012)			0.0±0.0	71.0±0.5	0.1±0.0	0.0±0.0	78.3±0.6	0.0±0.0
	CIFAR-10			3.2±0.3	85.1±0.6	3.4±0.3	0.1±0.1	93.1±0.6	0.4±0.2
	CIFAR-100			2.8±0.2	82.9±0.1	3.0±0.2	0.0±0.0	91.6±0.5	0.3±0.1
	ImageNet32			0.0±0.0	26.4±0.2	0.0±0.0	0.0±0.0	20.7±0.8	0.0±0.0
	ImageNet32-C(0.2)			0.0±0.0	67.7±0.6	0.0±0.0	0.0±0.0	73.7±0.5	0.0±0.0
	SVHN			0.0±0.0	28.5±0.2	0.0±0.0	0.0±0.0	22.0±0.4	0.0±0.0
	SVHN-C(1.8)			0.0±0.0	90.0±0.4	0.0±0.0	0.0±0.0	97.0±0.3	0.0±0.0
	LSUN			0.0±0.0	38.8±0.5	0.0±0.0	0.0±0.0	34.6±0.3	0.0±0.0
average			4.6	33.9	33.9	3.2	34.0	19.5	

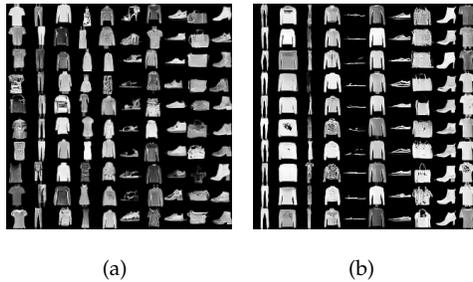


Fig. J.2: GlowGMM with 10 components trained on FashionMNIST. (a) sampling from $\mathcal{N}_i(\mu_i, \text{diag}(\sigma_i^2))$. The i -th column corresponds to Gaussian distribution \mathcal{N}_i . (b) For \mathcal{N}_i , we fit another Gaussian distribution $\tilde{\mathcal{N}}_i(\tilde{\mu}_i, \tilde{\Sigma}_i)$ using the normalized representations (by parameters of \mathcal{N}_i) of inputs of the $((i + 1)\%10)$ -th class. The i -th column shows images generated from $\tilde{\mathcal{N}}_i$.

TABLE J.5: GAD Results (AUROC and AUPR in percentage) of KLODS and Ty-test on Glow with batch sizes 2 and 4. The higher the better.

ID↓	OOD↓	Batch size →		$m=2$				$m=4$			
		Method →		KLODS		Ty-test		KLODS		Ty-test	
		Metric →		AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
Fash.	Constant	100.0±0.0	100.0±0.0	40.6±0.4	41.3±0.2	100.0±0.0	100.0±0.0	41.0±0.5	41.6±0.2		
	MNIST	91.6±0.1	91.8±0.2	88.7±0.2	81.1±0.4	99.4±0.0	99.4±0.0	96.1±0.1	93.2±0.1		
	MNIST-C(10.0)	97.2±0.1	97.3±0.1	74.0±0.3	65.2±0.2	100.0±0.0	100.0±0.0	85.4±0.2	77.4±0.5		
	notMNIST	99.2±0.0	99.4±0.0	64.0±0.3	61.8±0.3	100.0±0.0	100.0±0.0	74.3±0.4	71.2±0.3		
	notMNIST-C(0.005)	100.0±0.0	100.0±0.0	23.2±0.2	35.3±0.0	100.0±0.0	100.0±0.0	24.8±0.3	35.7±0.1		
SVHN	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.008)	99.9±0.0	99.8±0.0	14.6±0.5	33.2±0.1	100.0±0.0	100.0±0.0	14.5±0.4	33.2±0.1		
	CelebA	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CelebA-C(0.08)	90.1±0.2	88.6±0.4	49.0±0.3	46.6±0.2	96.2±0.2	95.1±0.3	55.8±0.4	50.4±0.3		
	CIFAR-10	100.0±0.0	100.0±0.0	99.9±0.0	99.9±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CIFAR-10-C(0.12)	77.2±0.2	76.1±0.2	36.6±0.2	40.0±0.1	81.9±0.3	80.6±0.2	32.9±0.8	38.3±0.3		
	CIFAR-100	99.9±0.0	99.9±0.0	99.9±0.0	99.9±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	CIFAR-100-C(0.12)	79.8±0.3	79.3±0.3	40.6±0.4	42.1±0.2	83.5±0.2	82.7±0.1	36.5±0.7	39.9±0.3		
	ImageNet32	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	ImageNet32-C(0.07)	97.8±0.1	98.1±0.1	48.4±0.3	48.2±0.1	100.0±0.0	100.0±0.0	42.5±0.3	44.1±0.1		
	LSUN	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
LSUN-C(0.06)	98.1±0.0	97.7±0.1	41.3±0.4	42.0±0.2	99.8±0.0	99.7±0.0	42.4±0.4	42.4±0.1			
CIFAR-10	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.02)	100.0±0.0	100.0±0.0	9.9±0.0	32.5±0.0	100.0±0.0	100.0±0.0	11.2±0.0	32.8±0.0		
	CelebA	93.3±0.1	94.6±0.1	98.0±0.1	98.1±0.0	98.4±0.1	98.7±0.1	99.9±0.0	99.9±0.0		
	CelebA-C(0.3)	72.4±0.3	71.6±0.3	32.4±0.3	38.1±0.1	81.3±0.3	81.2±0.3	29.7±0.4	37.1±0.1		
	ImageNet32	82.2±0.2	85.2±0.1	93.1±0.2	94.6±0.2	87.8±0.2	90.2±0.2	98.3±0.2	98.7±0.1		
	ImageNet32-C(0.3)	68.2±0.1	69.1±0.3	47.8±0.3	48.0±0.2	70.2±0.3	71.0±0.2	42.6±0.9	44.0±0.6		
	SVHN	90.0±0.1	90.7±0.2	91.2±0.1	88.1±0.3	96.2±0.1	96.5±0.1	97.6±0.1	96.8±0.2		
	SVHN-C(2.0)	99.1±0.1	99.2±0.0	39.2±0.1	64.0±0.1	100.0±0.0	100.0±0.0	35.2±0.5	61.9±0.2		
	LSUN	100.0±0.0	100.0±0.0	97.4±0.0	97.8±0.0	100.0±0.0	100.0±0.0	99.8±0.0	99.8±0.0		
LSUN-C(0.3)	86.8±0.1	87.6±0.2	46.6±0.4	45.5±0.2	89.2±0.3	90.0±0.2	50.4±0.3	47.7±0.2			
CelebA	Constant	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0		
	Uniform-C(0.12)	99.2±0.0	99.3±0.0	35.1±0.2	39.3±0.1	100.0±0.0	100.0±0.0	29.4±0.4	37.5±0.1		
	CIFAR-10	86.3±0.2	86.4±0.2	21.4±0.2	34.6±0.1	98.5±0.1	98.6±0.1	10.2±0.3	31.9±0.1		
	CIFAR-100	89.6±0.2	90.0±0.2	25.0±0.2	35.9±0.0	99.2±0.1	99.3±0.1	12.9±0.1	32.5±0.0		
	ImageNet32	100.0±0.0	100.0±0.0	76.4±0.3	83.5±0.1	100.0±0.0	100.0±0.0	76.9±0.4	84.5±0.2		
	ImageNet32-C(0.2)	99.2±0.0	99.3±0.0	34.6±0.2	40.3±0.1	100.0±0.0	100.0±0.0	28.4±0.2	37.4±0.0		
	SVHN	99.9±0.0	99.9±0.0	69.3±0.1	62.5±0.1	100.0±0.0	100.0±0.0	76.0±0.2	70.5±0.4		
	SVHN-C(1.8)	100.0±0.0	100.0±0.0	14.5±0.2	32.9±0.1	100.0±0.0	100.0±0.0	5.6±0.2	31.3±0.0		
	LSUN	100.0±0.0	100.0±0.0	60.2±0.1	58.8±0.3	100.0±0.0	100.0±0.0	63.8±0.1	62.5±0.3		
average		94.8	94.9	64.4	65.2	97.0	97.0	64.7	68.6		

TABLE J.6: GAD Results (AUROC and AUPR in percentage) of KLODS and GOD2KS on Glow with batch sizes 5 and 10. We run our method for 5 times. The higher the better. We reference the results on all problems of GOD2KS reported in [34], where the authors did not report average results of multiple runs.

ID↓	OOD↓	Batch size →		$m=5$				$m=10$			
		Method →		KLODS		GOD2KS		KLODS		GOD2KS	
		Metric →		AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
FashionMNIST	MNIST	99.8±0.0	99.8±0.0	98	98	100.0±0.0	100.0±0.0	100	100		
	KMNIST	99.9±0.0	99.9±0.0	97	96	100.0±0.0	100.0±0.0	100	100		
	Omniglot	100.0±0.0	100.0±0.0	100	100	100.0±0.0	100.0±0.0	100	100		
SVHN	CelebA	100.0±0.0	100.0±0.0	100	99	100.0±0.0	100.0±0.0	100	100		
	CIFAR-10	100.0±0.0	100.0±0.0	92	84	100.0±0.0	100.0±0.0	99	98		
	CIFAR-100	100.0±0.0	100.0±0.0	93	86	100.0±0.0	100.0±0.0	99	98		
	LSUN	100.0±0.0	100.0±0.0	99	98	100.0±0.0	100.0±0.0	100.0	100.0		
CIFAR-10	CelebA	99.2±0.1	99.4±0.1	86	92	100.0±0.0	100.0±0.0	96	98		
	SVHN	97.6±0.2	97.8±0.2	96	98	99.8±0.0	99.8±0.0	100	100		
	LSUN	100.0±0.0	100.0±0.0	60	58	100.0±0.0	100.0±0.0	58	56		
CelebA	CIFAR-10	99.6±0.0	99.6±0.0	84	73	100.0±0.0	100.0±0.0	94	91		
	CIFAR-100	99.8±0.0	99.8±0.0	82	71	100.0±0.0	100.0±0.0	94	90		
	SVHN	100.0±0.0	100.0±0.0	97	98	100.0±0.0	100.0±0.0	100	100		
	LSUN	100.0±0.0	100.0±0.0	85	75	100.0±0.0	100.0±0.0	96	92		
average		99.7	99.7	90.6	87.6	100.0	100.0	95.4	94.5		

TABLE J.7: GAD Results (AUROC and AUPR in percentage) of KLODS and Ty-test on Glow with batch size 5 and 10. For SVHN, CIFAR-10, and CelebA, we choose one dataset as ID data and the mixture of the other two datasets as OOD data.

ID↓	OOD ↓	Batch size→		m=5				m=10			
		Method→		KLODS		Ty-test		KLODS		Ty-test	
		Metric→		AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
SVHN	CelebA+CIFAR-10	100.0±0.0									
CIFAR-10	SVHN+CelebA	98.2±0.2	98.5±0.2	60.8±0.3	64.4±0.5	99.8±0.0	99.9±0.0	52.8±1.0	58.1±1.1		
CelebA	SVHN+CIFAR-10	100.0±0.0	100.0±0.0	20.7±0.2	34.9±0.1	100.0±0.0	100.0±0.0	11.8±0.5	32.3±0.1		
average		99.4	99.5	60.5	66.4	99.9	100.0	54.9	63.5		

TABLE J.8: GAD Results (AUROC in percentage) of ablation study. We compare four methods Ty-test, KLOD-all, KLOD, and KLODS. KLODS is the best one.

ID	OOD↓	Batch size Method	m=10				m=25			
			Ty-test	KLOD-all	KLOD	KLODS	Ty-test	KLOD-all	KLOD	KLODS
Fash.M	Constant		41.6±0.41	100.0±0.0	100.0±0.0	100.0±0.0	39.1±0.8	100.0±0.0	100.0±0.0	100.0±0.0
	MNIST		99.2±0.1	100.0±0.0						
	MNIST-C(10.0)		84.9±0.3	100.0±0.0	100.0±0.0	100.0±0.0	94.7±0.3	100.0±0.0	100.0±0.0	100.0±0.0
	notMNIST		92.7±0.5	100.0±0.0	100.0±0.0	100.0±0.0	98.9±0.2	100.0±0.0	100.0±0.0	100.0±0.0
	notMNIST-C(0.005)		7.0±0.6	100.0±0.0	100.0±0.0	100.0±0.0	2.7±0.2	100.0±0.0	100.0±0.0	100.0±0.0
SVHN	Constant		100.0±0.0							
	Uniform		100.0±0.0							
	Uniform-C(0.008)		11.8±0.3	100.0±0.0	100.0±0.0	100.0±0.0	5.2±0.6	100.0±0.0	100.0±0.0	100.0±0.0
	CelebA		100.0±0.0							
	CelebA-C(0.08)		54.7±0.5	100.0±0.0	100.0±0.0	100.0±0.0	58.2±0.3	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR10		100.0±0.0							
	CIFAR10-C(0.12)		54.7±0.5	86.2±0.3	100.0±0.0	100.0±0.0	12.6±0.9	98.3±0.5	99.1±0.3	100.0±0.0
	CIFAR100		100.0±0.0							
	CIFAR100-C(0.12)		26.9±1.3	86.0±0.8	95.5±0.4	100.0±0.0	12.0±1.1	96.2±1.0	97.2±0.2	100.0±0.0
	ImageNet32		100.0±0.0							
ImageNet32-C(0.07)		42.6±0.4	100.0±0.0	100.0±0.0	100.0±0.0	35.7±0.3	100.0±0.0	100.0±0.0	100.0±0.0	
CIFAR10	Constant		100.0±0.0							
	Uniform		100.0±0.0							
	Uniform-C(0.02)		10.7±0.1	100.0±0.0	100.0±0.0	100.0±0.0	5.1±1.0	100.0±0.0	100.0±0.0	100.0±0.0
	CelebA		100.0±0.0							
	CelebA-C(0.3)		23.4±5.3	100.0±0.0	100.0±0.0	100.0±0.0	12.6±0.7	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32		100.0±0.0	99.7±0.1	99.3±0.0	94.7±0.1	100.0±0.0	95.3±0.7	99.0±0.3	98.9±0.4
	ImageNet32-C(0.3)		31.7±0.7	98.4±0.2	94.8±0.3	73.1±1.5	15.0±1.0	97.5±0.3	96.7±0.5	76.9±1.3
	SVHN		99.9±0.0	96.7±0.2	99.1±0.0	99.8±0.1	100.0±0.0	87.6±0.5	99.6±0.1	100.0±0.0
	SVHN-C(2.0)		26.7±0.6	100.0±0.0	100.0±0.0	100.0±0.0	58.2±0.2	100.0±0.0	100.0±0.0	100.0±0.0
	CIFAR10		1.0±0.1	94.3±0.8	99.8±0.0	100.0±0.0	0.0±0.0	95.6±0.5	100.0±0.0	100.0±0.0
CelebA	CIFAR100		2.0±0.2	94.7±0.4	99.8±0.0	100.0±0.0	0.0±0.0	95.2±0.4	100.0±0.0	100.0±0.0
	ImageNet32		87.9±0.3	100.0±0.0	100.0±0.0	100.0±0.0	96.7±0.4	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32-C(0.2)		18.2±0.3	100.0±0.0	100.0±0.0	100.0±0.0	7.8±0.3	100.0±0.0	100.0±0.0	100.0±0.0
	SVHN		91.5±0.6	100.0±0.0	100.0±0.0	100.0±0.0	98.6±0.2	100.0±0.0	100.0±0.0	100.0±0.0
	SVHN-C(1.8)		1.4±0.2	100.0±0.0	100.0±0.0	100.0±0.0	0.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	average except CIFAR10 vs ImageNet32-C(0.3)		65.49	98.76	99.79	99.82	64.43	98.94	99.83	99.96
average		64.40	98.75	99.63	98.95	62.84	98.89	99.73	99.22	

TABLE J.9: GAD Results (AUROC and AUPR in percentage) of Glow on One-vs-Rest on MNIST. The higher is the better.

Batch size Method	m=5				m=10			
	KLODS		Ty-test		KLODS		Ty-test	
	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
class 0 vs rest	92.9±1.2	93.4±1.1	86.5±1.5	89.9±1.1	99.1±0.6	99.2±0.5	95.7±1.6	96.8±1.0
class 1 vs rest	96.6±0.5	96.9±0.6	100.0±0.0	100.0±0.0	99.8±0.1	99.8±0.1	100.0±0.0	100.0±0.0
class 2 vs rest	85.0±0.7	85.8±0.9	69.8±1.5	74.8±1.8	95.0±1.2	95.2±1.2	76.7±1.4	81.7±1.4
class 3 vs rest	65.5±2.6	64.9±3.2	65.7±1.5	70.0±1.3	76.2±1.5	77.0±1.1	68.1±2.1	72.7±1.7
class 4 vs rest	94.8±0.6	95.3±0.5	77.4±0.8	82.2±1.4	99.8±0.1	99.8±0.1	81.6±2.1	86.0±1.8
class 5 vs rest	66.3±2.3	66.1±2.6	61.6±2.2	64.6±2.2	78.3±2.4	78.6±3.0	62.6±2.7	64.6±2.8
class 6 vs rest	88.0±1.0	88.4±0.9	64.3±1.7	66.9±1.1	98.0±0.7	98.2±0.7	58.6±1.2	63.0±1.2
class 7 vs rest	94.0±0.5	94.3±0.6	91.0±0.3	93.2±0.1	99.2±0.3	99.3±0.3	96.3±1.1	97.4±0.6
class 8 vs rest	76.0±1.4	76.8±1.6	81.2±1.6	85.7±1.3	86.2±1.4	86.7±1.8	92.1±0.8	93.7±1.0
class 9 vs rest	95.0±0.7	95.6±0.8	76.3±1.5	80.2±1.8	99.7±0.2	99.7±0.2	80.1±2.6	84.2±1.9
average	85.4	85.8	77.4	80.8	93.1	93.4	81.2	84.0

TABLE J.10: GAD results (AUROC and AUPR in percentage) on GlowGMM trained on FashionMNIST.

Batch size		$m=25$		
Method	KLODS		Ty-test	
Metrics	AUROC	AUPR	AUROC	AUPR
class 0 vs rest	100.0±0.0	100.0±0.0	5.4±1.6	31.2±0.3
class 1 vs rest	100.0±0.0	100.0±0.0	15.7±2.4	33.4±4.9
class 2 vs rest	100.0±0.0	100.0±0.0	0.5±0.5	30.7±0.0
class 3 vs rest	99.9±0.1	99.9±0.1	89.6±2.5	91.3±2.3
class 4 vs rest	100.0±0.0	100.0±0.0	0.7±0.6	30.7±0.0
class 5 vs rest	100.0±0.0	100.0±0.0	64.2±1.4	66.4±2.9
class 6 vs rest	99.9±0.1	99.9±0.1	0.0±0.0	30.7±0.0
class 7 vs rest	100.0±0.0	100.0±0.0	31.4±2.8	46.6±3.3
class 8 vs rest	100.0±0.0	100.0±0.0	0.4±0.5	30.7±0.0
class 9 vs rest	100.0±0.0	100.0±0.0	69.0±3.6	76.0±1.7
average	100	100	27.7	46.8

TABLE J.11: GlowGMM trained on FashionMNIST. Use $p(z)$ as criterion for 1 vs rest classification.

Method	$p(z)$	
Metrics	AUROC	AUPR
class 0 vs rest	72.7±1.6	72.0±1.4
class 1 vs rest	85.1±0.6	86.2±0.6
class 2 vs rest	74.8±4.5	76.9±4.0
class 3 vs rest	68.9±4.7	71.2±4.5
class 4 vs rest	77.1±2.1	78.4±3.2
class 5 vs rest	71.7±1.4	71.9±1.2
class 6 vs rest	73.5±7.8	73.7±8.6
class 7 vs rest	86.9±0.4	88.6±0.4
class 8 vs rest	55.5±0.9	53.8±0.5
class 9 vs rest	86.6±0.3	87.1±0.3
average	75.3	76.0

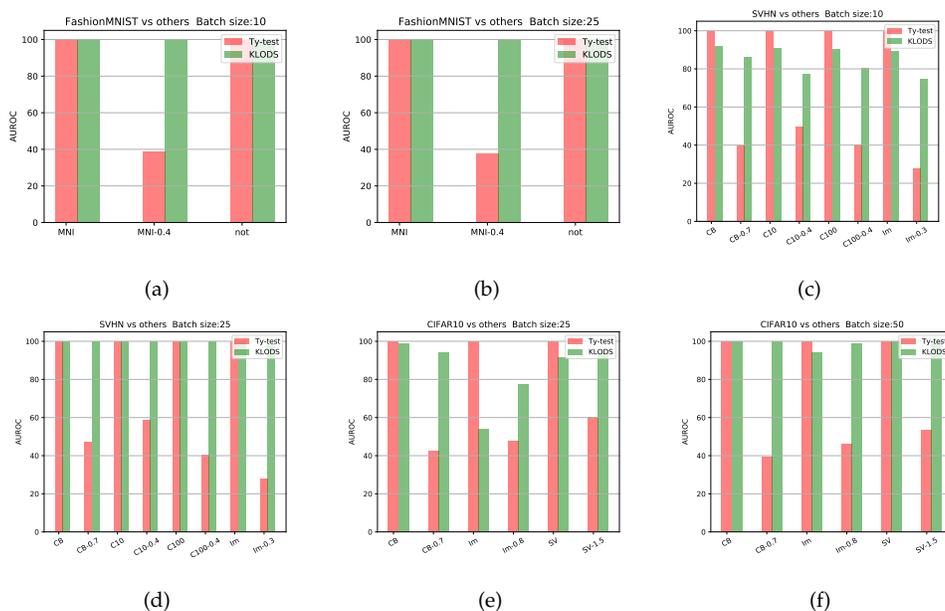


Fig. J.3: GAD results (AUROC) on convolutional VAE. with batch sizes 5 and 10. The X-axis labeled with OOD datasets with abbreviated names. MNI: MNIST, not: notMNIST, CB: CelebA, C10/100: CIFAR-10/100, Im: ImageNet, SV: SVHN. The number k after the dataset name indicates the dataset with adjusted contrast with a factor k . For example, CB-0.7 means CelebA-C(0.7). Numerical results are shown in Table J.12 in the supplementary material.

TABLE J.12: GAD results (AUROC and AUPR in percentage) of KLOD on VAE.

ID↓	OOD↓	Batch size→		$m=10$				$m=25$			
		Method→		KLOD		Ty-test		KLOD		Ty-test	
		Metric→		AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR
Fash.	MNIST	99.7±0.1	99.5±0.2	100.0±0.0							
	MNIST-C(0.4)	99.8±0.0	99.8±0.0	39.1±0.7	40.5±0.3	100.0±0.0	100.0±0.0	37.6±1.9	39.8±0.7	100.0±0.0	100.0±0.0
	notMNIST	100.0±0.0									
SVHN	CelebA	92.2±0.6	82.3±1.1	100.0±0.0							
	CelebA-C(0.7)	86.2±0.9	76.5±1.5	39.9±1.2	41.2±0.5	100.0±0.0	100.0±0.0	47.4±1.5	44.3±0.7	100.0±0.0	100.0±0.0
	CIFAR-10	90.9±1.3	81.3±2.3	100.0±0.0							
	CIFAR-10-C(0.4)	77.6±8.8	69.9±1.3	49.8±0.6	45.8±0.3	100.0±0.0	100.0±0.0	58.8±0.9	50.2±0.4	100.0±0.0	100.0±0.0
	CIFAR-100	90.4±0.4	80.3±0.6	100.0±0.0							
	CIFAR-100-C(0.4)	80.5±1.0	73.2±1.8	40.3±0.8	40.7±1.3	100.0±0.0	100.0±0.0	40.5±0.4	41.3±0.2	100.0±0.0	100.0±0.0
	ImageNet32	89.3±8.6	80.1±1.5	100.0±0.0							
	ImageNet32-C(0.3)	74.6±0.6	67.8±0.7	27.9±1.0	36.5±0.3	100.0±0.0	100.0±0.0	27.9±1.0	36.5±0.3	100.0±0.0	100.0±0.0
average		89.2	82.8	72.5	73.2	99.9	99.9	73.8	73.8	100.0±0.0	100.0±0.0
ID ↓	OOD ↓	Batch size		$m=25$				$m=50$			
CIFAR-10	CelebA	99.1±0.4	99.1±0.4	100.0±0.0							
	CelebA-C(0.7)	94.2±0.6	93.8±0.8	42.3±1.1	42.8±0.6	100.0±0.0	100.0±0.0	39.3±2.0	41.1±1.0	100.0±0.0	100.0±0.0
	ImageNet32	54.0±1.9	53.4±0.7	99.8±0.1	99.8±0.1	94.0±0.6	94.0±0.5	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	ImageNet32-C(0.8)	77.4±1.4	77.3±1.8	47.8±1.5	48.0±1.5	98.8±0.5	98.9±0.4	46.4±1.7	46.8±1.2	100.0±0.0	100.0±0.0
	SVHN	91.8±1.5	91.1±2.3	99.8±0.0	99.8±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0
	SVHN-C(1.5)	94.2±1.5	91.1±2.3	60.0±1.7	61.4±1.7	100.0±0.0	100.0±0.0	53.6±2.7	55.7±1.6	100.0±0.0	100.0±0.0
	average		85.1	84.3	75.0	75.3	98.8	98.8	73.2	73.9	100.0±0.0

TABLE J.13: GAD results (AUROC and AUPR) of KLOD without split representations on VAE trained on CIFAR10 and tested on CIFAR100. Each row is for one batch size.

Problem	CIFAR10 vs CIFAR100				CIFAR10 vs ImageNet32			
	KLOD		Ty-test		KLOD		Ty-test	
	Metric	AUROC	AUPR	AUROC	AUPR	AUROC	AUPR	AUROC
$m=50$	72.9±0.7	73.7±2.1	73.8±0.5	74.3±1.8	94.0±0.6	94.0±0.5	100.0±0.0	100.0±0.0
$m=100$	90.9±1.0	91.3±1.3	82.6±0.5	83.5±1.1	99.9±0.2	99.9±0.2	100.0±0.0	100.0±0.0
$m=150$	98.0±0.4	98.1±0.5	88.4±1.3	88.6±2.3	100.0±0.0	100.0±0.0	100.0±0.0	100.0±0.0

TABLE J.14: VAE trained on CIFAR10. Use reconstruction probability for OOD data detection.

Method	reconstruction probability	
Metrics	AUROC	AUPR
SVHN	17.6±0.0	34.3±0.0
CelebA	83.1±0.0	82.5±0.0
ImageNet32	72.4±0.2	75.0±0.1
CIFAR100	52.3±0.0	53.6±0.0
average	56.4	61.4

TABLE J.15: PAD results (AUROC in percentage) of Joint confidence loss method, ODIN, Joint confidence loss method+ODIN, DoSE, and KLODS. We use all the problems evaluated in the original publication of Joint confidence loss method [83]. Our method outperforms all baselines.

* The authors of Joint confidence loss method [83] did not report AUROC result of Joint confidence loss method+ODIN on SVHN vs TinyImageNet. Since joint confidence loss+ODIN is reported to be better than Joint confidence loss method, so we just use 100% AUROC.

ID	OOD	confidence loss	joint confidence loss	ODIN	joint confidence loss+ODIN	DoSE	KLODS
SVHN	CIFAR-10	83	97.6	94	99	96.2	98.9
	TinyImageNet	98	99.5	95	100*	100	99.8
	LSUN	98.5	99.8	94	100	91.6	100
CIFAR-10	SVHN	46.5	67.5	85	85	95.5	82.6
	TinyImageNet	67	72.5	76	86	76.7	83.9
	LSUN	62.5	76	78	87.5	98	98.9
average		75.9	85.5	87	92.9	93	94.0

APPENDIX K

FIGURES

This section contains figures referred to in the other parts.

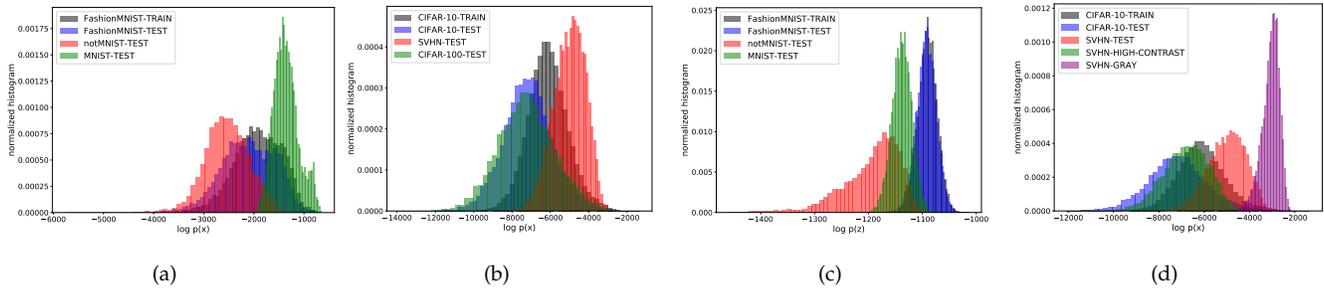


Fig. K.1: Distributions of likelihoods of ID dataset (train and test) and OOD dataset. (a) and (b) show the normalized histogram of $\log p(z)$ and $\log p(x)$ on Glow trained on FashionMNIST, respectively. (c) shows the normalized histogram of $\log p(x)$ on Glow trained on CIFAR-10. (d) shows that $\log p(x)$ of OOD data can be manipulated by adjusting the contrast of images. SVHN-HIGH-CONTRAST and SVHN-GRAY are SVHN with adjusted contrast by a factor of 2.0 and 0.3, respectively.

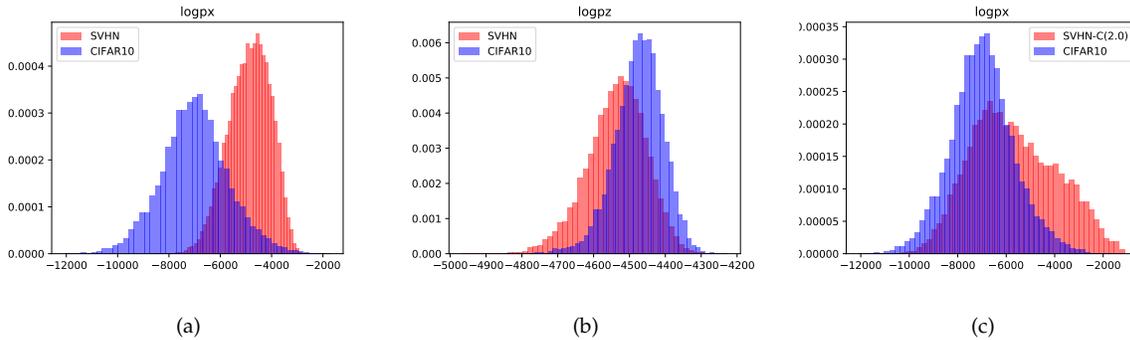


Fig. K.2: Residual flow trained on CIFAR-10 assigns (a) higher $\log p(x)$ for SVHN; (b) similar $\log p(z)$ for SVHN; and (c) coinciding $\log p(x)$ for SVHN with increased contrast with a factor of 2. We use the official implementation at [105] and the model checkpoint released at [106].

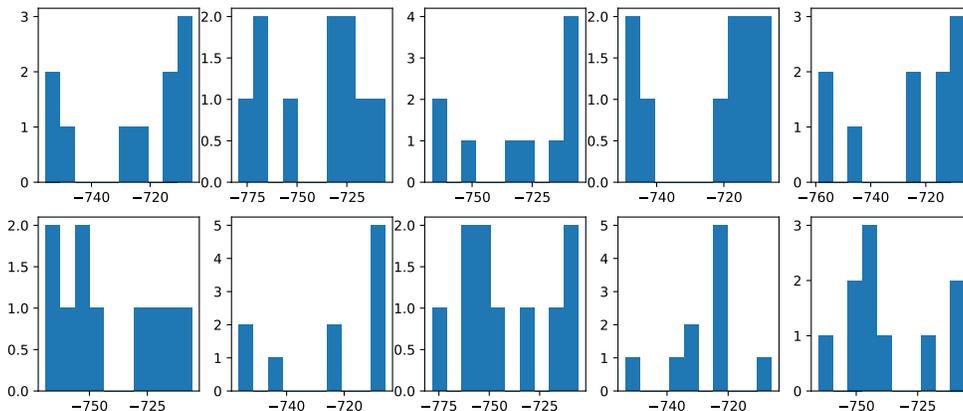


Fig. K.3: Train GlowGMM on FashionMNIST. The i -th subfigure shows the histogram of log-probabilities of 10 centroids under the i -th Gaussian component. All log-probabilities are close to $768 \times \log(1/\sqrt{2\pi}) \approx -705.74$, which is the log-probability of the center of 768-dimensional standard Gaussian distribution. These results indicate that these centroids are close to each other.

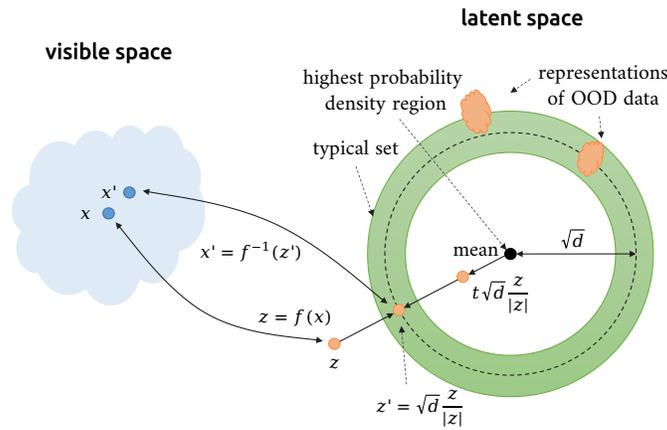


Fig. K.4: Rescaling z to the typical set of prior. For each input x , we can compute $z = f(x)$, and rescale z to the typical set annulus of Gaussian prior as $z' = \sqrt{d} \frac{z}{|z|}$. Then we map z' back to visible space as $x' = f^{-1}(z')$. We observe that x' is similar to x . See Figure K.5 for examples.

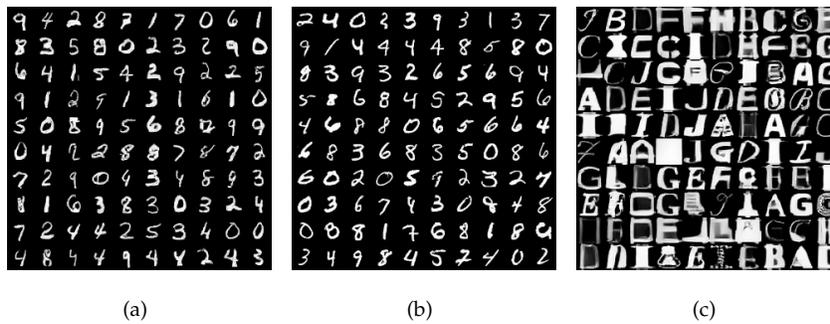


Fig. K.5: Train Glow on FashionMNIST and test on (a), (b) MNIST and (c) notMNIST. We scale the representations of OOD dataset to the typical set of prior Gaussian distribution. The scaled latent vectors still correspond to nearly the same images. (a) and (c): rescale only the last scale of OOD representation to the typical set of prior. The first and second scales are kept as standard Gaussian noise. (b) rescale the OOD representations at all scales to the typical set of prior.

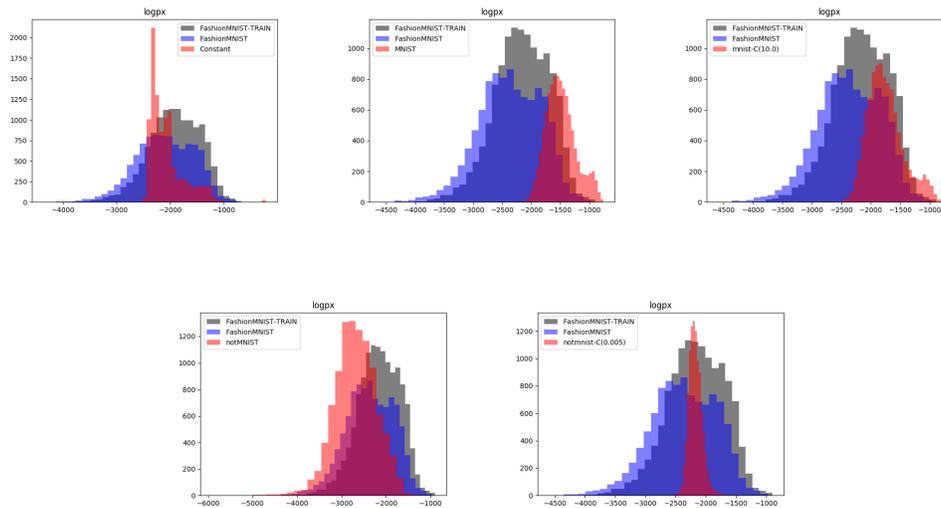


Fig. K.6: Glow trained on FashionMNIST. Histogram of $\log p(x)$. We can manipulate the likelihood distribution of OOD dataset by adjusting the contrast. “-C(k)” means the dataset with adjusted contrast by a factor of k .

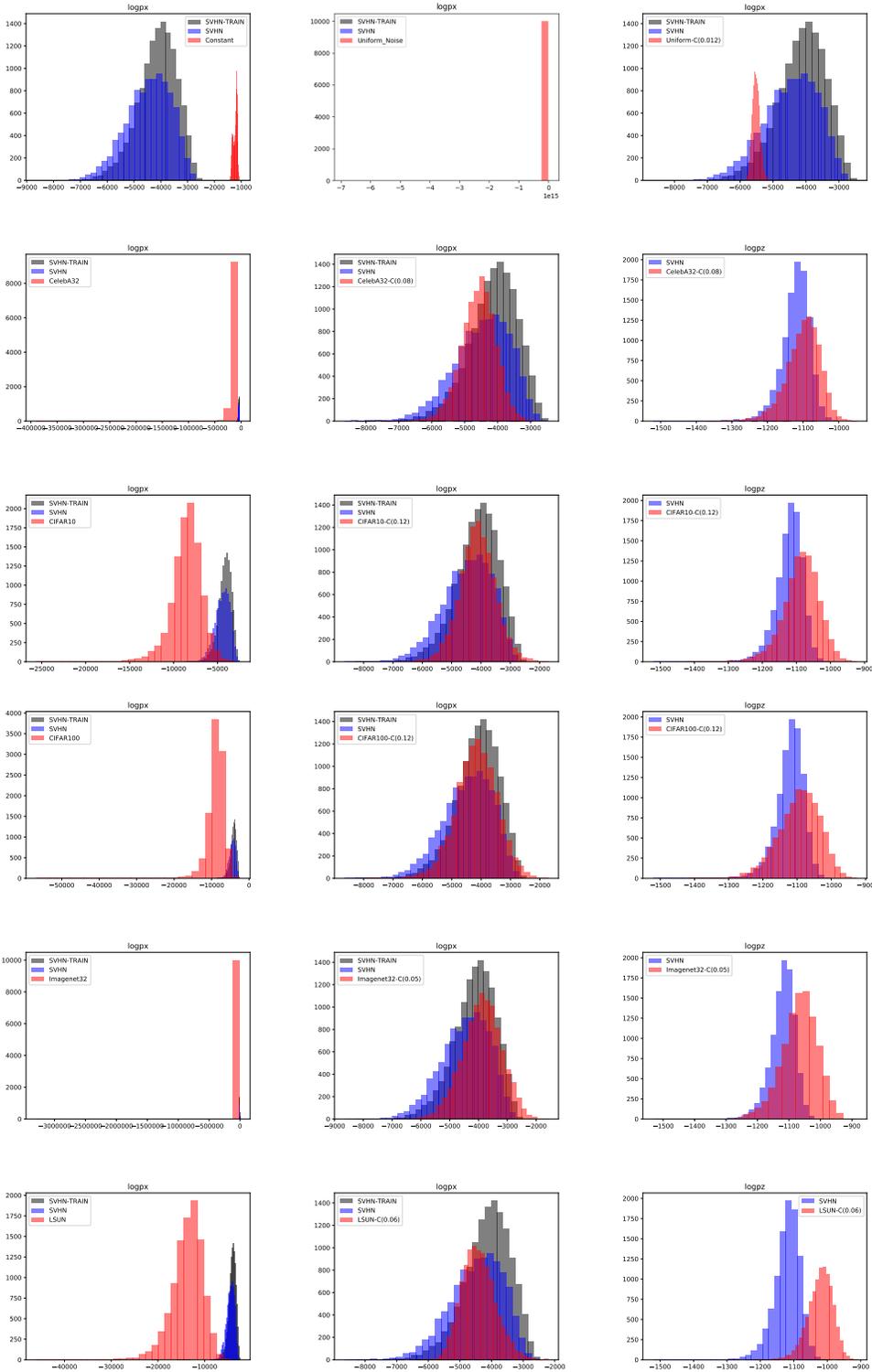


Fig. K.7: Official Glow trained on SVHN. Histogram of $\log p(\mathbf{x})$, $\log p(\mathbf{z})$, and $\log p(\mathbf{x})$ contributed by the last scale. We can manipulate the likelihood distribution of OOD dataset by adjusting the contrast. “-C(k)” means the dataset with adjusted contrast by a factor of k . Note that the distribution of $\log p(\mathbf{x})$ of the last scale and $\log p(\mathbf{z})$ have a similar shape. This is because the log-determinant of the last scale is similar for every data point in the same dataset. We do not observe this phenomenon in Glow trained on CIFAR-10.

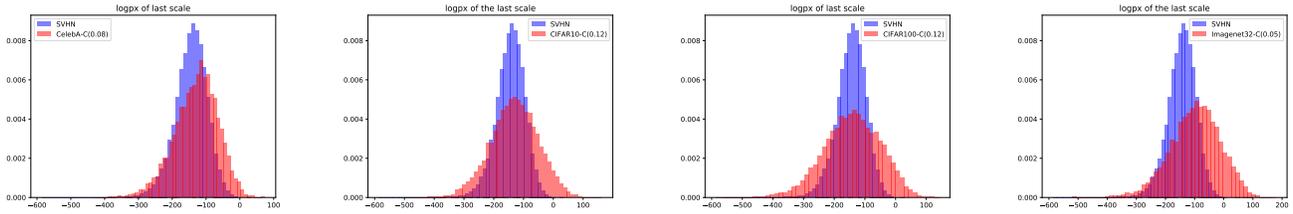


Fig. K.8: To reproduce the results of baseline method L_{last} precisely, we use the implementation of Glow model and checkpoint trained on SVHN released by the authors of L_{last} . The $\log p(x)$ of the last scale of OOD and ID data also coincide under data manipulation. We can also see that the $\log p(x)$ of the last scale even becomes positive for OOD data. The authors also said that the metric used by L_{last} could not be explained as log-likelihood for OOD data.

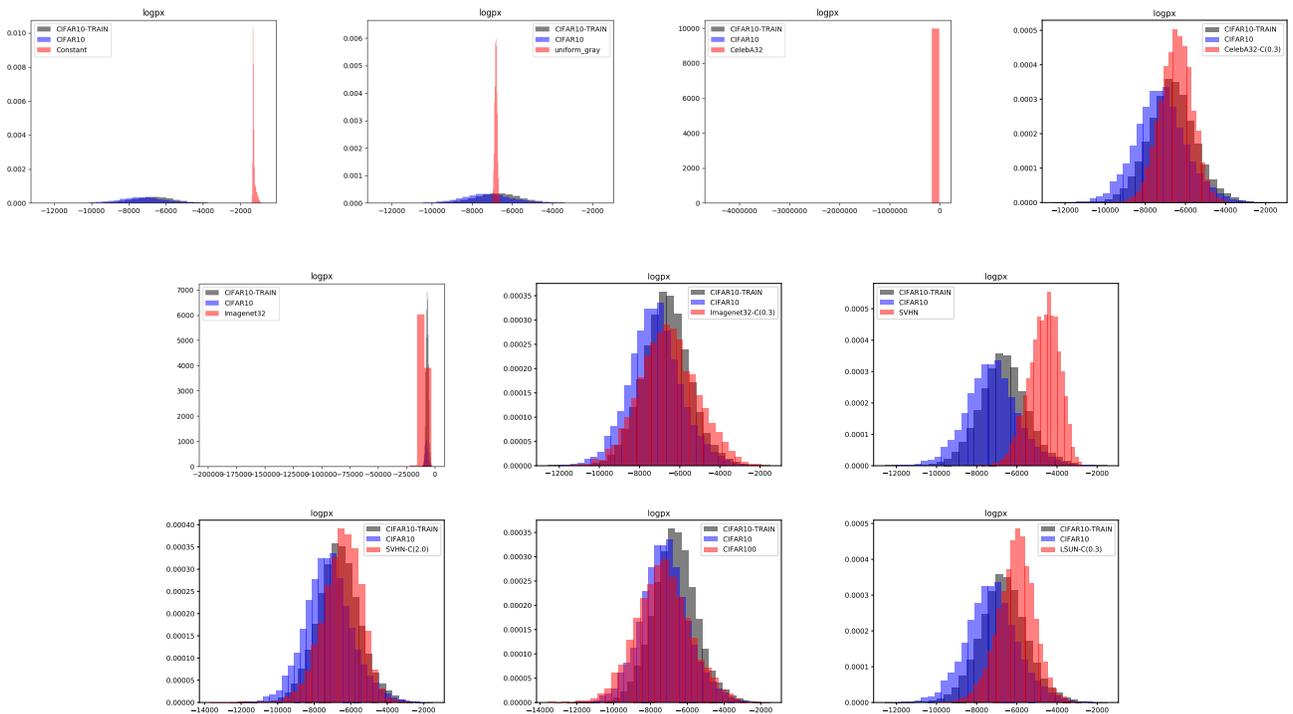


Fig. K.9: Glow trained on CIFAR10. Histogram of $\log p(x)$. We can manipulate the likelihood distribution of OOD dataset by adjusting the contrast. “-C(k)” means the dataset with adjusted contrast by a factor of k . The ranges of $\log p(x)$ of CelebA and LSUN are too large to break the scale of the figure. For CIFAR10 vs Uniform, $\log p(x)$ of Uniform are too small.

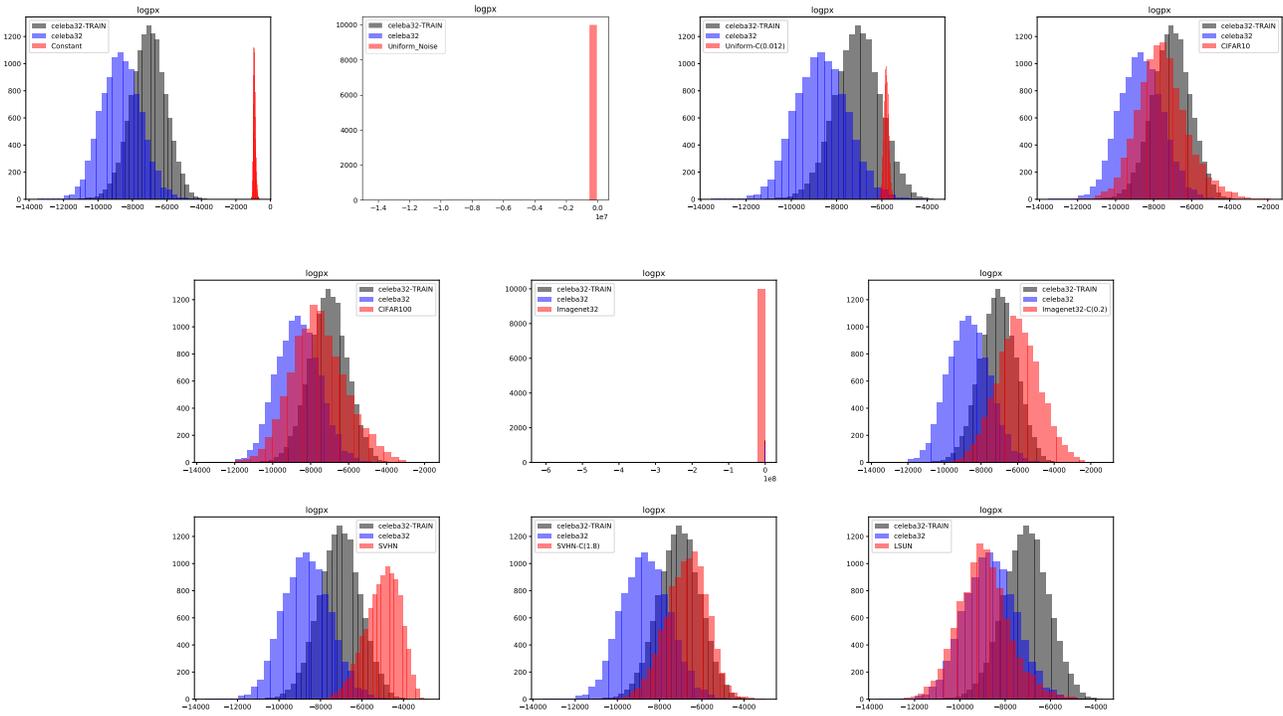


Fig. K.10: Glow trained on CelebA. Histogram of $\log p(\mathbf{x})$. We can manipulate the likelihood distribution of OOD dataset by adjusting the contrast. “-C(k)” means the dataset with adjusted contrast by a factor of k . It is hard to make the likelihoods of train and test split of CelebA fit well on the official Glow model.

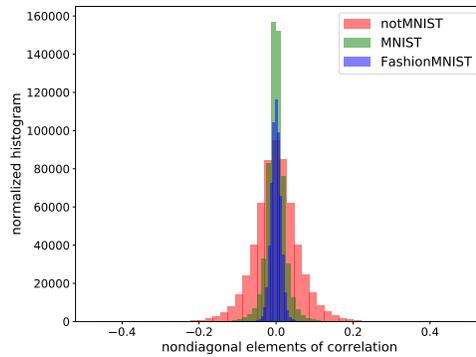


Fig. K.11: Glow trained on FashionMNIST and tested on MNIST/notMNIST. Histogram of non-diagonal elements in the correlation coefficient of representations.

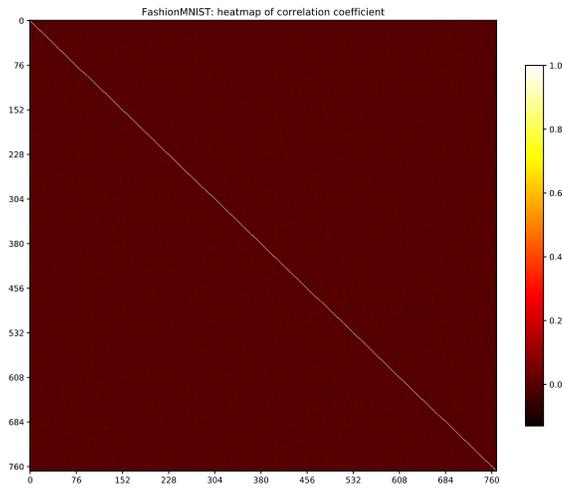


Fig. K.12: Glow trained on FashionMNIST. Heatmap of correlation of FashionMNIST representations.

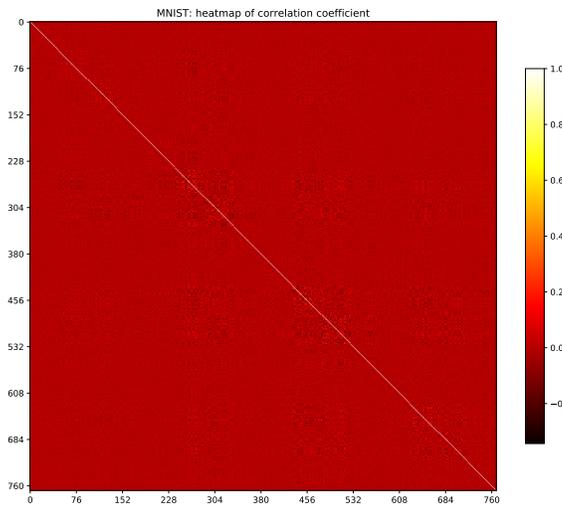


Fig. K.13: Glow trained on FashionMNIST. Heatmap of correlation of MNIST representations.

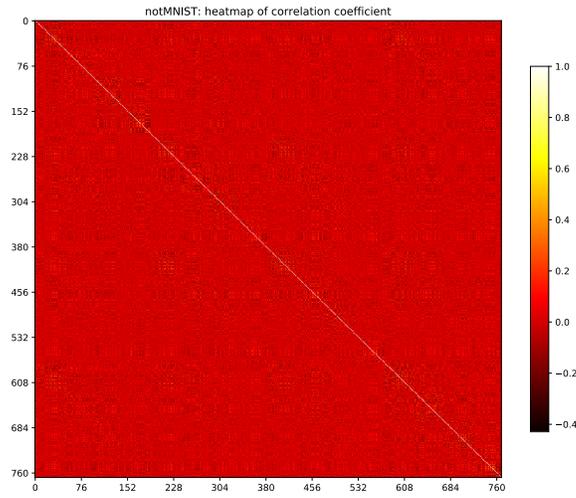


Fig. K.14: Glow trained on FashionMNIST. Heatmap of correlation of notMNIST representations.

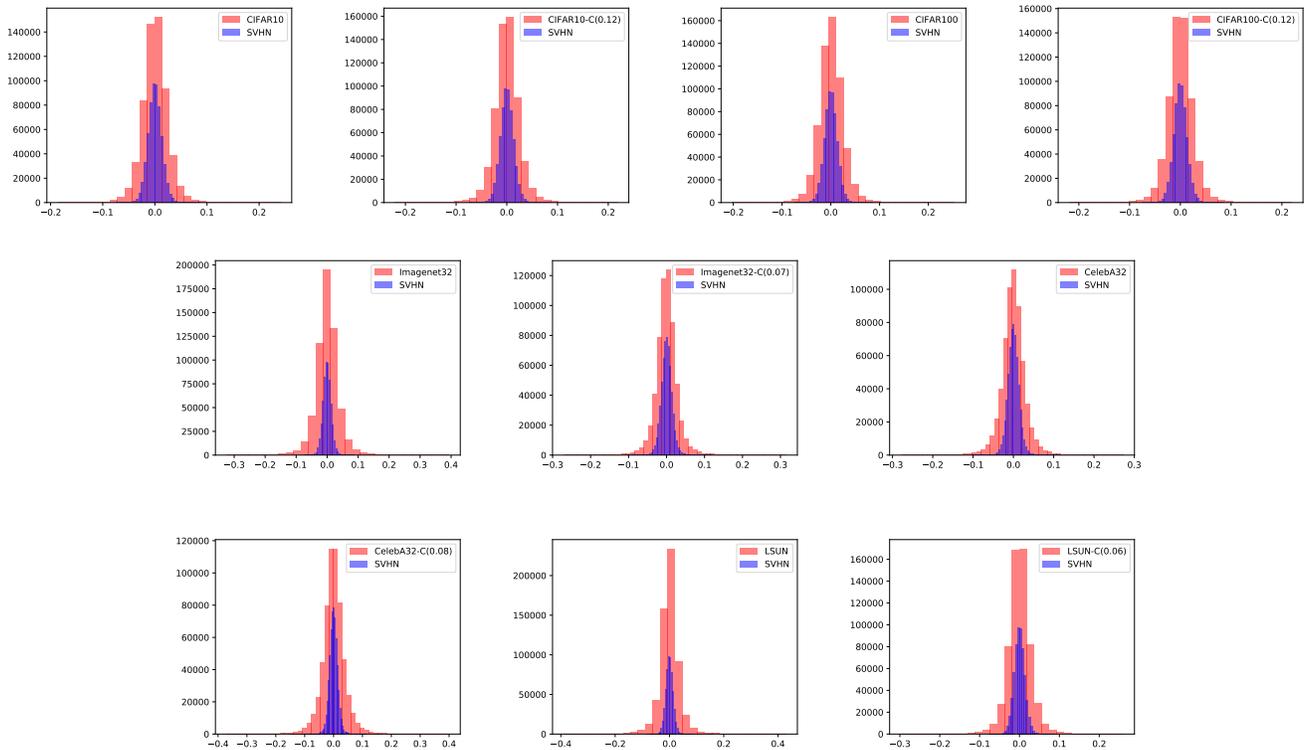


Fig. K.15: Glow trained on SVHN. Histogram of non-diagonal elements of correlation of representations.

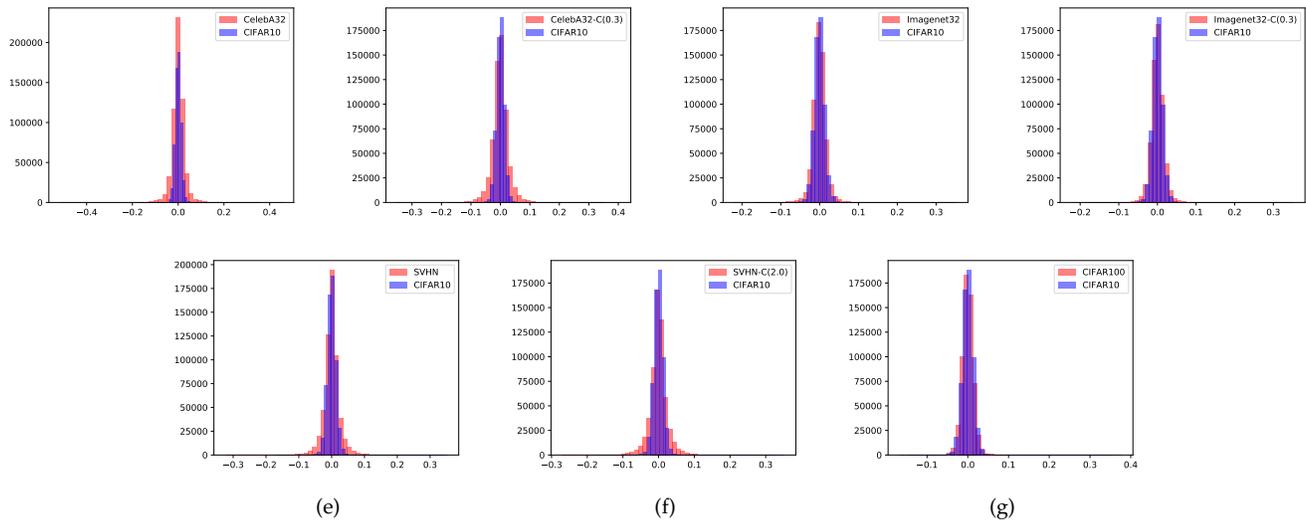


Fig. K.16: Glow trained on CIFAR10. Histogram of non-diagonal elements of correlation of representations.

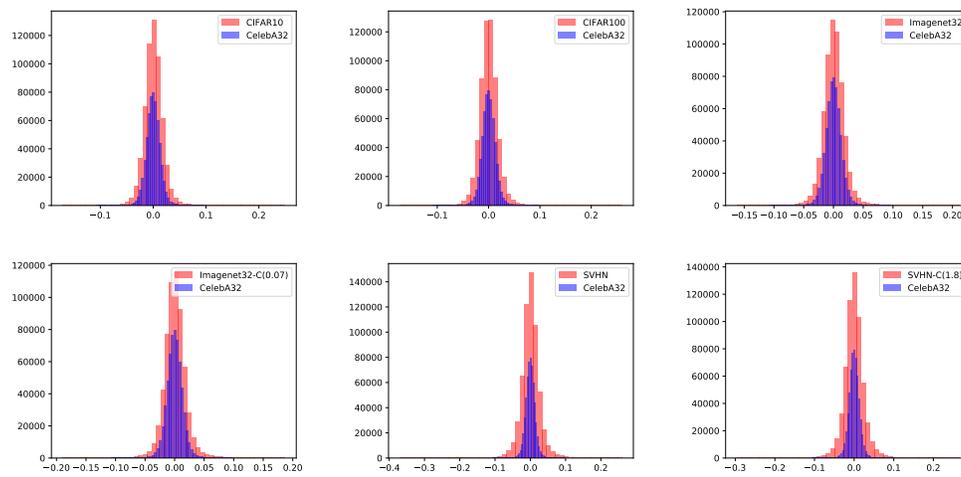


Fig. K.17: Glow trained on CelebA. Histogram of non-diagonal elements of correlation of representations.

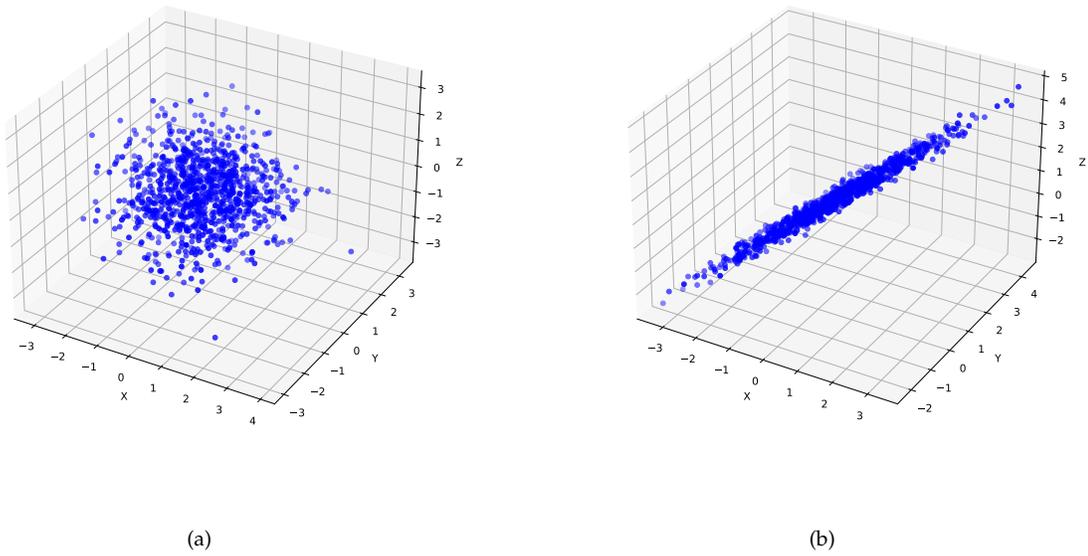


Fig. K.18: Samples from 3-d Gaussian distribution $\mathcal{N}(\mu, \Sigma)$. The mean μ and covariance matrix Σ determines where the data locate in. (a) $\mu = (0, 0, 0)$, $\Sigma = ((1, 0, 0)^\top, (0, 1, 0)^\top, (0, 0, 1)^\top)$. (b) $\mu = (0, 1, 1)$, $\Sigma = ((1, 0.98, 0.98)^\top, (0.98, 1, 0.98)^\top, (0.98, 0.98, 1)^\top)$.



Fig. K.19: Glow trained on CIFAR-10. Generated images from prior (up), fitted Gaussian distribution from the representations of OOD dataset notMNIST (down).

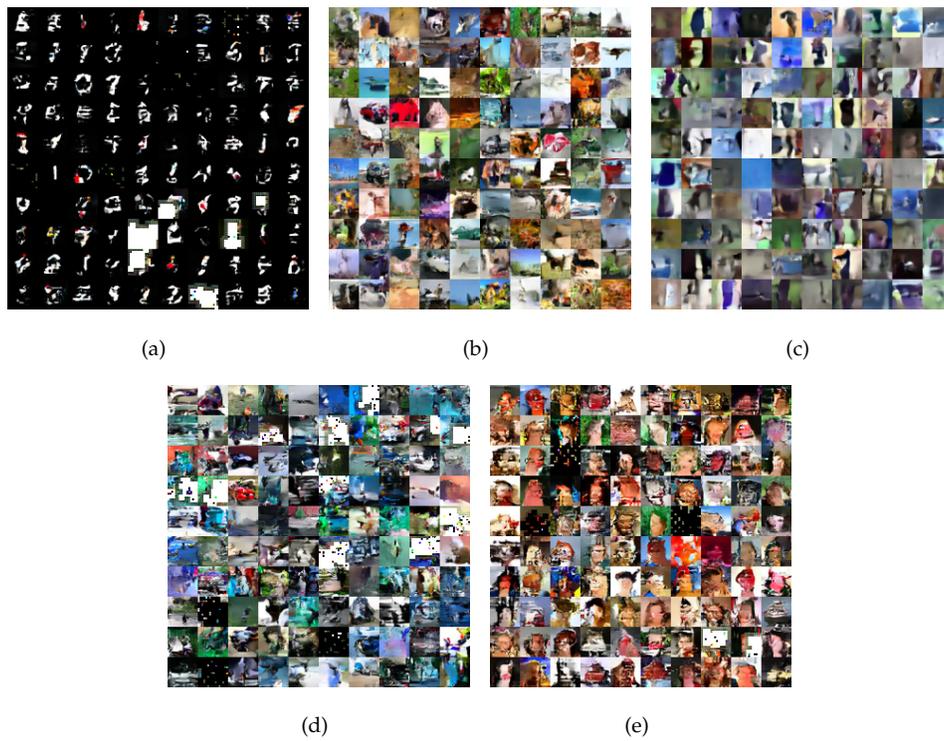


Fig. K.20: Glow trained on CIFAR10. Generated images according to the fitted Gaussian distribution from representations of (a) MNIST; (b) CIFAR100; (c) SVHN; (d) ImageNet32; (e) CelebA. We replicate MNIST into three channels and pad zeros for consistency. These results demonstrate that the covariance of representations contains important information of an OOD dataset.

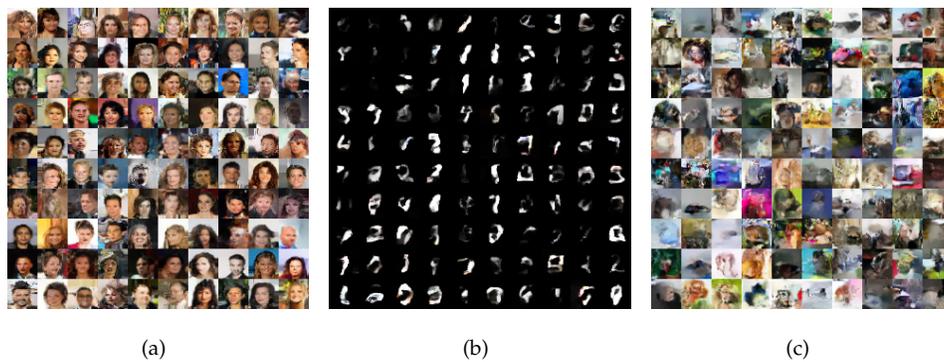


Fig. K.21: Glow trained on CelebA 32×32 , sampling according to (a) standard Gaussian distribution; (b) fitted Gaussian distribution from MNIST representations; (c) fitted Gaussian distribution from CIFAR10 representations.



Fig. K.22: Glow trained on FashionMNIST. Sampling according to prior (up), fitted Gaussian distribution from representations of MNSIT (middle) and notMNIST (down).

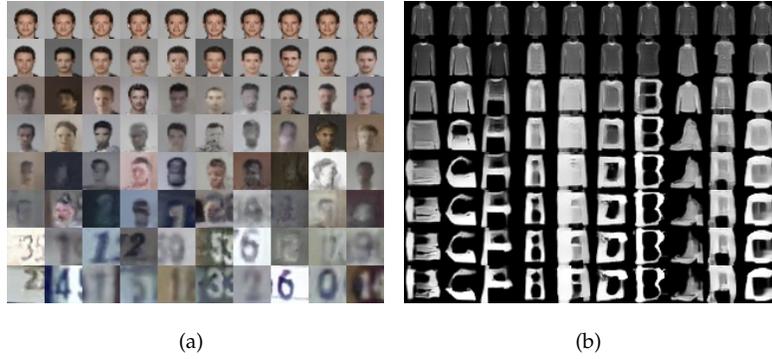


Fig. K.23: (a) Train Glow on CelebA and sample from the fitted Gaussian distribution of SVHN. (b) Train on FashionMNIST and sample from the fitted Gaussian distribution of notMNIST. From top to down, the sampled noises from Gaussian distribution are scaled by temperature 0, 0.25, 0.5, 0.6, 0.7, 0.8, 0.9, 1.0, respectively.

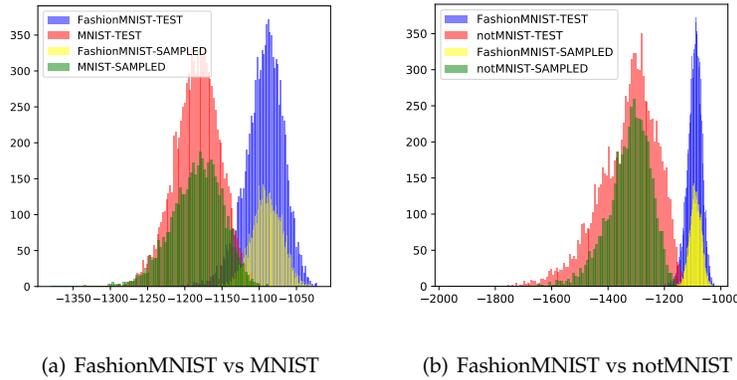


Fig. K.24: Glow trained on FashionMNIST. Histogram of $\log p(z)$ of (a) FashionMNIST vs MNIST, (b) FashionMNIST vs notMNIST under Glow. The green part corresponds to the $\log p(z)$ of noises sampled from the fitted Gaussian distribution of OOD datasets.

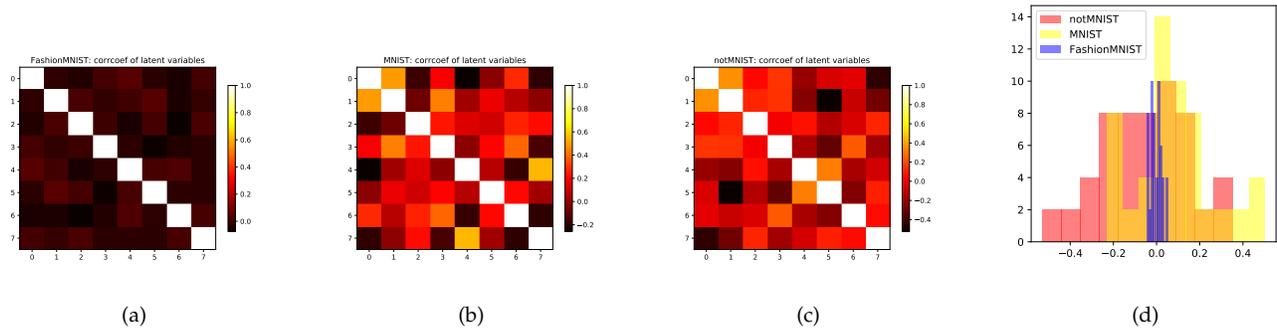


Fig. K.25: VAE trained on FashionMNIST. Heatmap of correlation of (a) FashionMNIST (b) MNIST (c) notMNIST representations. (d) Histogram of non-diagonal elements of correlation of sampled representations.

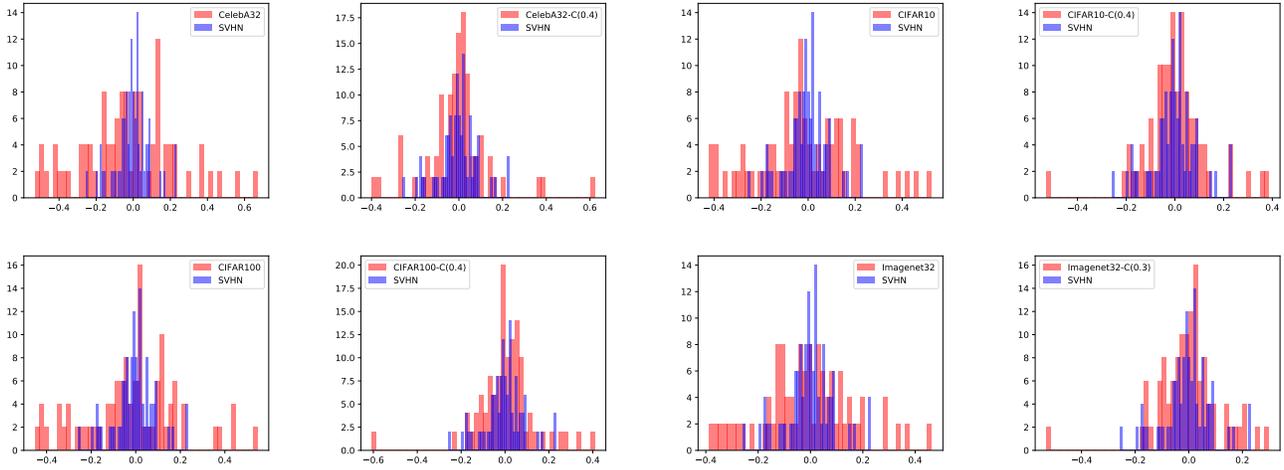


Fig. K.26: VAE trained on SVHN. Histogram of non-diagonal elements of correlation of sampled representations.

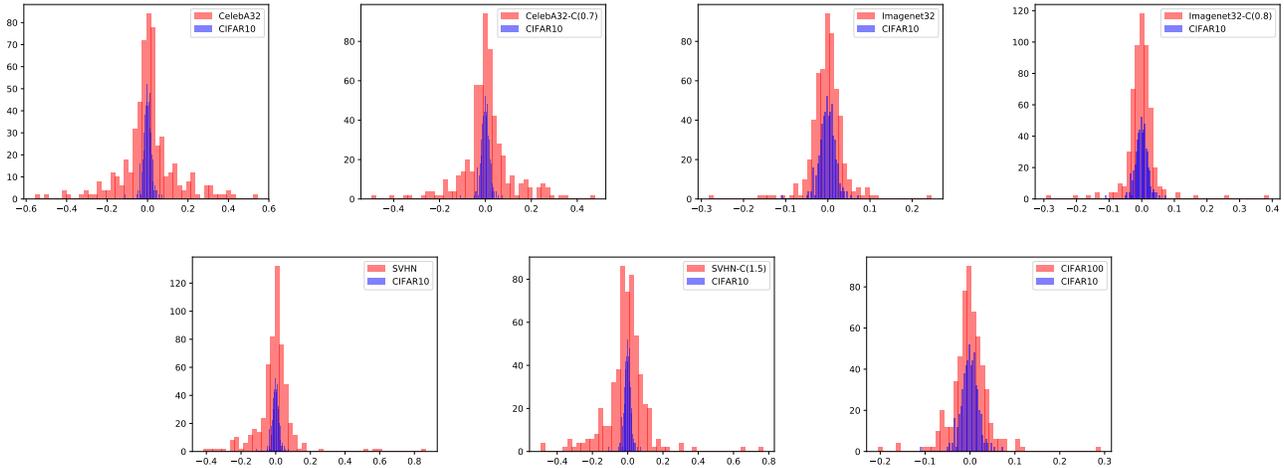


Fig. K.27: VAE trained on CIFAR10. Histogram of non-diagonal elements of correlation of sampled representations.

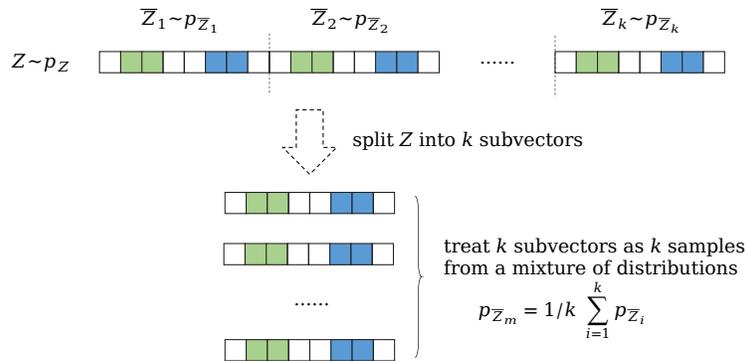


Fig. K.28: Split a random vector $Z \sim p_Z$ into k subvectors $\bar{Z}_i \sim p_{\bar{Z}_i}$ ($1 \leq i \leq k$). We treat k subvectors as k samples from a mixture of distributions $p_{\bar{Z}_i} = 1/k \sum_{i=1}^k p_{\bar{Z}_i}$. In the figure, we use the same color to indicate neighboring pixels that are strongly correlated. For example, if the second element $\bar{Z}_{i,2}$ and the third element $\bar{Z}_{i,3}$ are strongly correlated for all $1 \leq i \leq k$, we can say that $\bar{Z}_{m,2}$ and $\bar{Z}_{m,3}$ are also strongly correlated. This is why we can leverage local pixel dependence in our method.

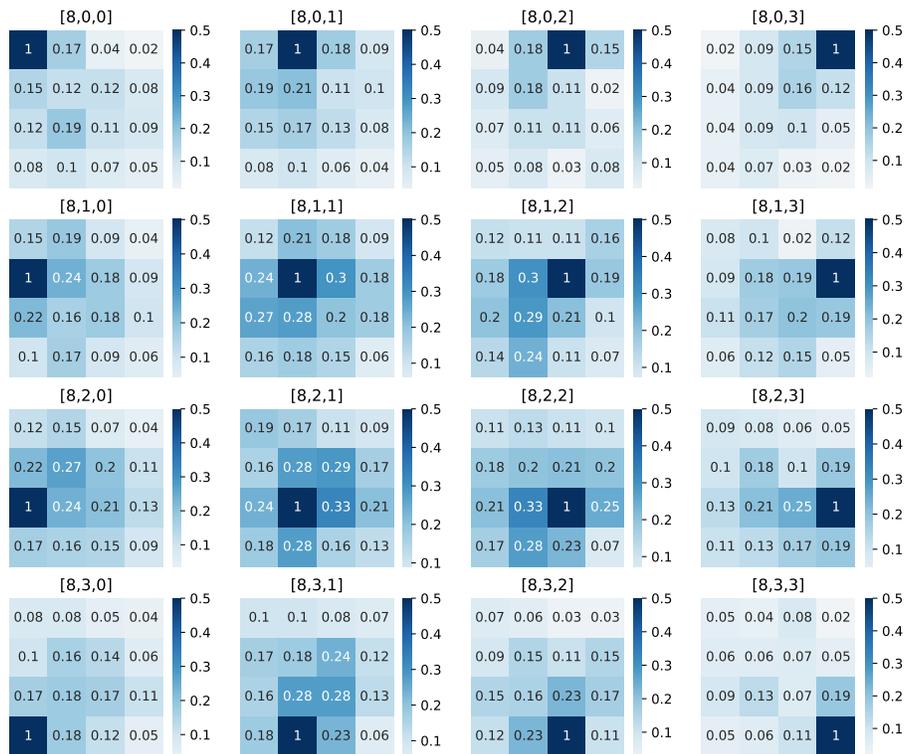


Fig. K.29: Train Glow on SVHN and test on ImageNet32. We randomly select the 8-th channel. The subfigure at i -th row and j -th column shows the correlation between the pixel at position (i, j) and all other pixels. Adjacent pixels tend to have stronger correlation.



(a) SVHN



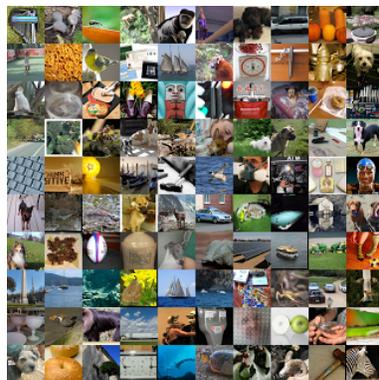
(b) SVHN with increased contrast by a factor of 2, have lower likelihood



(c) CelebA32



(d) CelebA32 with decreased contrast by a factor of 0.3, have higher likelihood



(e) ImageNet32



(f) ImageNet32 with decreased contrast by a factor of 0.3, have higher likelihood

Fig. K.30: Examples of datasets and their mutations. Under Glow trained on CIFAR10, these mutated datasets have a similar likelihood distribution with CIFAR10 test split.

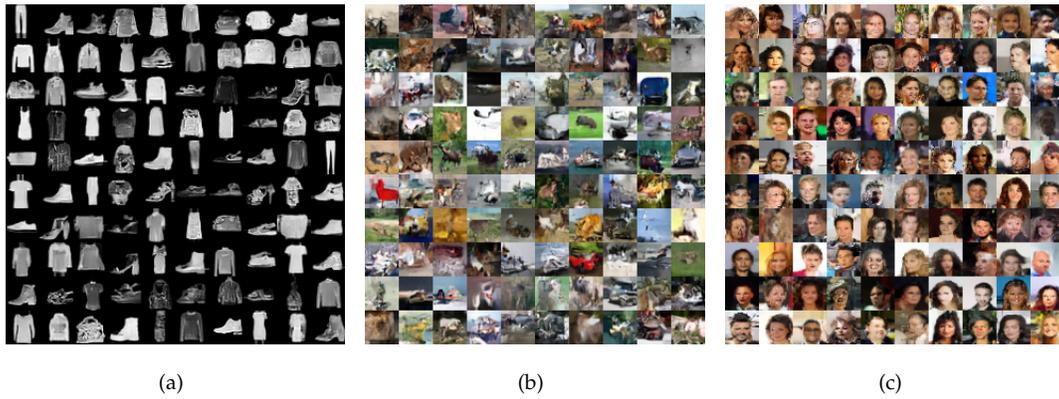


Fig. K.31: Generated images from Glow trained on (a)FashionMNIST; (b)CIFAR-10; (c)CelebA32.

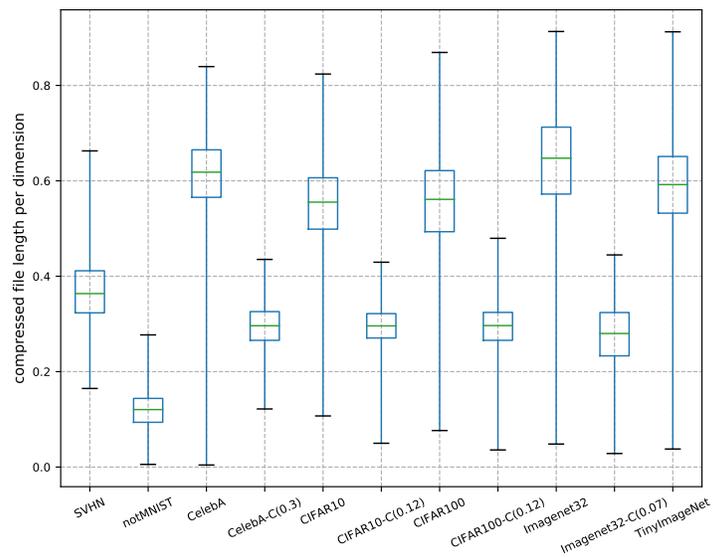


Fig. K.32: The distributions of complexity estimated by the lengths of compressed files of datasets. We use FLIF as compressor and compute lengths in bits per dimension. Datasets with decreased contrast has lower complexity.