

1 We thank all reviewers for their useful comments. First we address the common concern regarding the **importance and**
2 **impact of our work** in the following three arguments:

3 1. Although classifier based OOD detectors often outperform OOD detectors based on generative models, we should
4 note that many machine learning systems deployed in practice are trained on unlabeled data, and generative models are
5 widely used in tasks like language modeling and image editing. Therefore, OOD detection with unsupervised generative
6 models is an important task.

7 2. Among likelihood based generative models, VAEs have the promising property of extracting representations from
8 data, and they are particularly popular in applications (for example, VAE is a component of most image-to-image
9 translation frameworks). Recent work obtains impressive OOD detection on flow/auto-regressive models, but **our**
10 **motivation is that, if VAE is being used in a task, it would be beneficial to enable the VAE itself to detect OOD inputs,**
11 **rather than training a separate generative model merely for OOD detection.** However, as we show, recent SOTA OOD
12 scores for generative models cannot be applied to VAEs. Our work makes an important contribution by introducing an
13 easy yet effective OOD score that works well on VAEs.

14 3. It is most natural to apply LR to VAEs due to the bottleneck structure serving as a regularization for optimizable
15 model configuration, however, later we also explore the possibility of applying LR to flow/pixelCNN. **In Table 1, we**
16 **report the AUCROC on some tasks obtained from LR on Glow and PixelCNN.** We only optimize the last coupling block
17 and the last masked conv layer for Glow and PixelCNN respectively to avoid overfitting on a single input. Our results
18 are competitive to those of previous SOTA methods on generative models. **We therefore highlight LR as a general**
19 **approach for OOD detection of all likelihood based generative models, with unique effectiveness on VAEs.**

	Glow	PixelCNN
MNIST	0.995	0.987
KMNIST	0.990	0.973
NotMNIST	0.991	0.995
Noise	0.989	0.964
Constant	1	1

(a) Models trained on Fashion MNIST

	Glow	PixelCNN
SVHN	0.894	0.889
CelebA	0.735	0.774
LSUN	0.643	0.686
Noise	1	0.975
Constant	0.989	1

(b) Models trained on CIFAR-10

Table 1: AUCROC obtained from Likelihood Regret on Glow and PixelCNN.

20 We now address detailed concerns of each reviewer.

21 **Reviewer 1:** We apologize for the confusion caused by the notation. Indeed, the decoder should be denoted with θ .
22 Our VAE is based on the DCGAN structure, which is fully convolutional and the last FC layer is replaced by conv
23 layer that reduces the spatial dimension to 1×1 (which returns the vectors for posterior mean/variance). We tried
24 different network structures, including those with FC layers, and have found little difference in terms of OOD detection.
25 The precise definition of OOD samples is interesting and debatable, as we only have an in-distribution *dataset* rather
26 than a true distribution. Strictly speaking, both blue circle and cars are OOD, but the former is much harder to detect,
27 especially for generative models.

28 **Reviewer 2:** Please refer to items 1-3 for concerns regarding why we focus on VAE’s OOD detection. In short, our
29 motivation is to let VAEs detect OOD when it is used, and our extended experiments on flow/pixelCNN indicate that LR
30 is general across generative models. As for why optimizing the encoder works better, for in-dist samples, optimizing
31 the encoder is more constrained than directly optimizing z , which prevents the latent variables from moving too much.
32 We will provided a detailed explanation in the revision. We will carefully fix writing issues as pointed out.

33 **Reviewer 3:** Please refer to items 1-3 for concerns regarding why we focus on VAE’s OOD detection. We found that
34 β -VAEs behave similarly to plain VAEs in terms of OOD detection. We will include more details on applying our
35 method to other VAE variants. We will also include the MNIST results in the appendix as suggested. In particular,
36 we obtain AUCROC 0.963 on your mentioned MNIST vs. EMNIST (VAE’s likelihood gives 0.781). CIFAR-10
37 vs. CIFAR-100 is indeed a very hard task for generative models. We tested LR on this task and obtained AUCROC
38 0.582 (compared to 0.489 using likelihood), while likelihood ratio gives 0.564 and IC gives 0.535. This result is also
39 comparable to IC on pixelCNN. We thank the reviewer for pointing out writing issues and typos in the code. These will
40 be fixed, suggested additions will be included, as well as the correct reference for the definition of OOD.

41 **Reviewer 4:** Please refer to items 1-3 for concerns regarding why we focus on VAE’s OOD detection and applying LR
42 to other generative models. We only tested on image data because most recent OOD papers, both classifier based and
43 generative model based, focus on image data. The Likelihood Ratio paper by Ren et al. introduces a gene sequence
44 OOD dataset, and we tested our Likelihood Regret on a VAE trained on this dataset. Our LR obtains 0.745 AUCROC
45 (the VAE’s likelihood baseline gives 0.538). This result is competitive to Ren et al. (which uses a auto-regressive model
46 whose likelihood gives 0.626 and their proposed likelihood ratio improves to 0.732).