

Do Deep Generative Models Know What They Don't Know?



Eric Nalisnick, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, Balaji Lakshminarayanan

1. INTRODUCTION

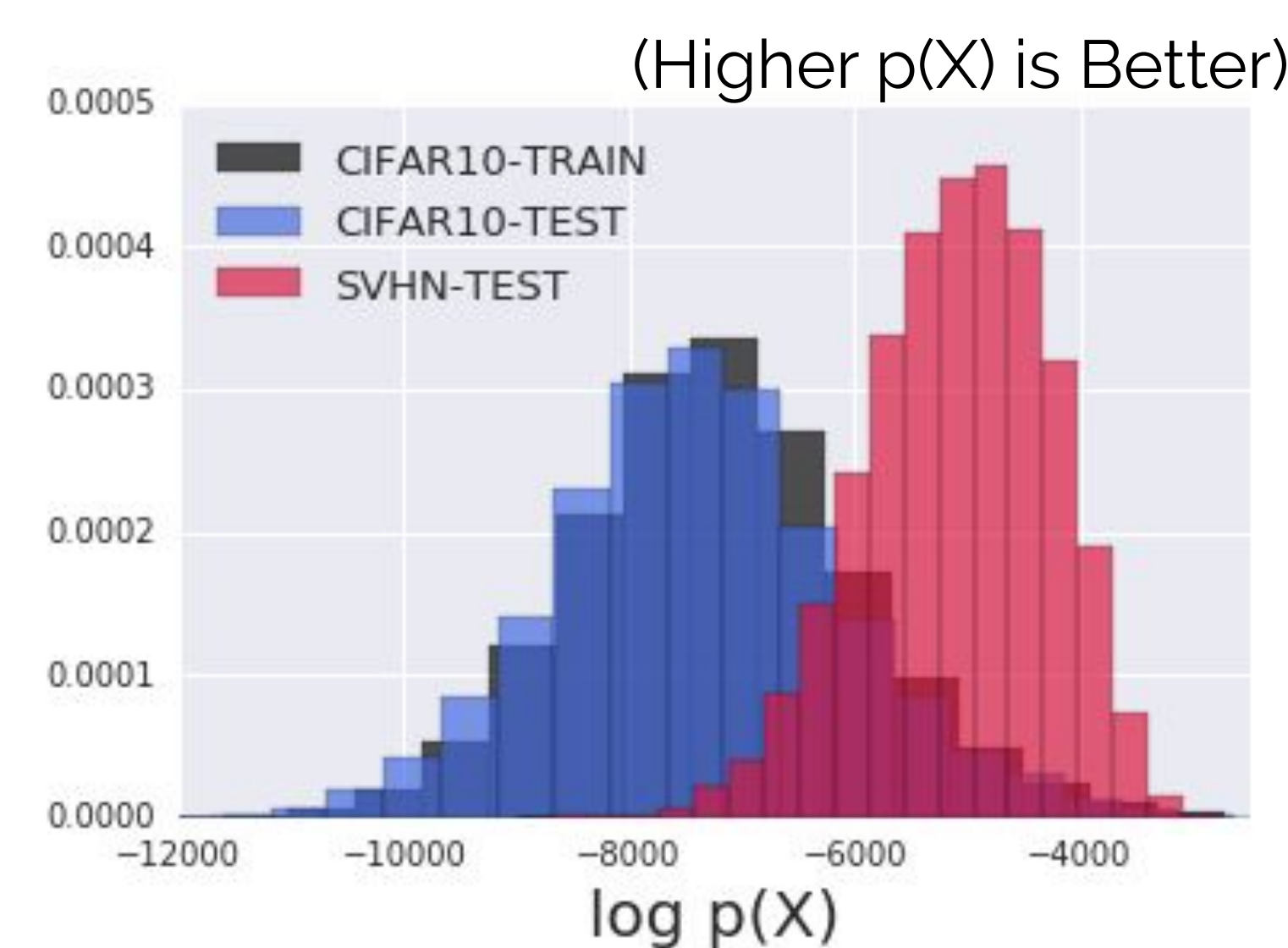
- Discriminative models are susceptible to overconfidence on **out-of-distribution (OOD) inputs**. Generative models are widely believed to be more robust to such inputs as they also model $p(\mathbf{x})$ [Bishop, 1994].
- We challenge this assumption, showing that **deep generative models can assign higher density estimates to an OOD dataset than to the training data!**
- This phenomenon has implications not just for anomaly detection but also for detecting covariate shift, open-set classification, active learning, semi-supervised learning, etc

2. MOTIVATING OBSERVATION: CIFAR-10 vs SVHN



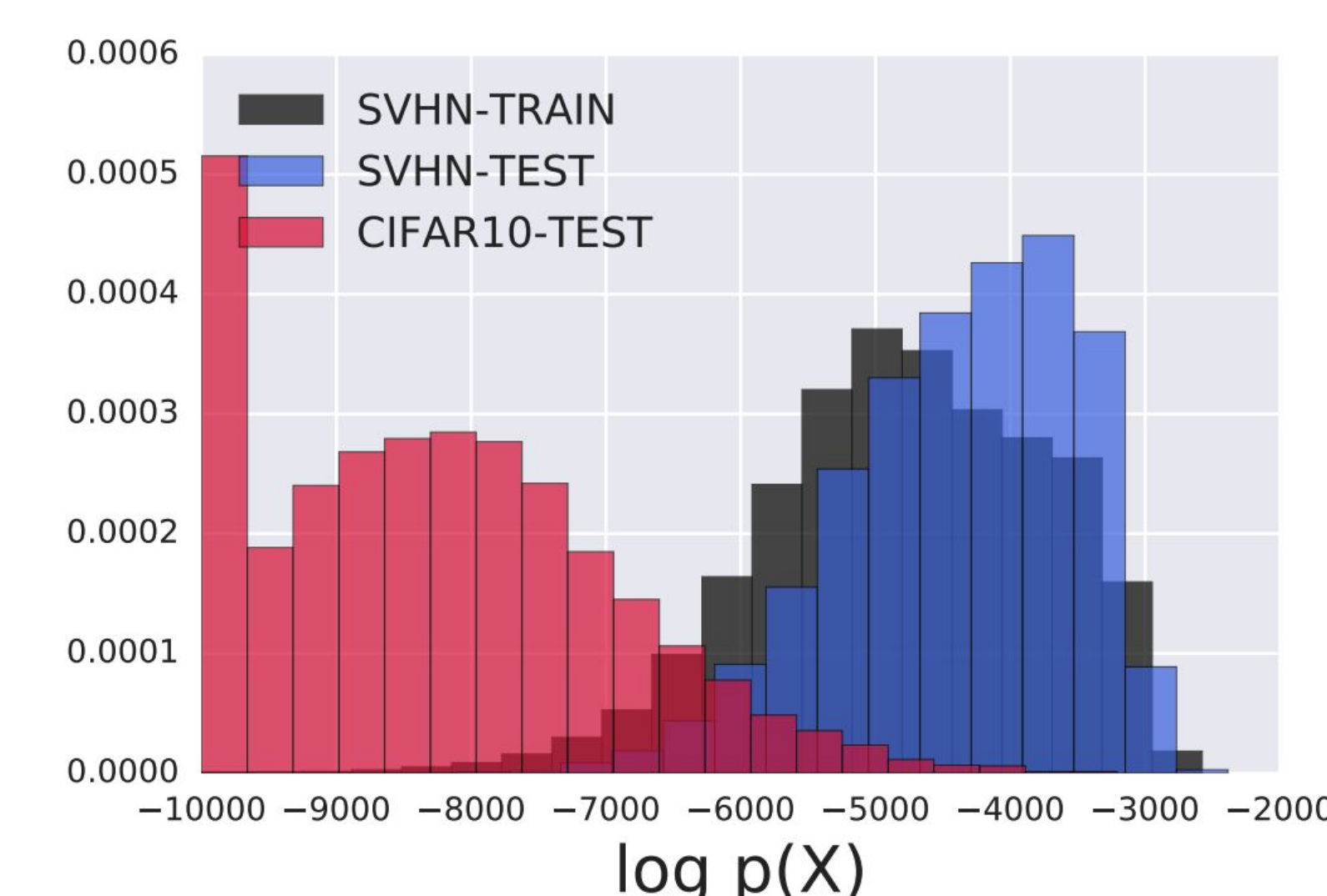
We trained **Glow** [Kingma & Dhariwal, 2018] on CIFAR-10 and evaluated on the model on SVHN. We find that **Glow assigns a higher likelihood to SVHN than to CIFAR-10** (both train/test splits).

Data Set	Avg. Bits Per Dimension
<i>Glow Trained on CIFAR-10</i>	
CIFAR10-Train	3.386
CIFAR10-Test	3.464
SVHN-Test	2.389
<i>Glow Trained on SVHN</i>	
SVHN-Test	2.057



Histogram of Glow Log-Likelihoods

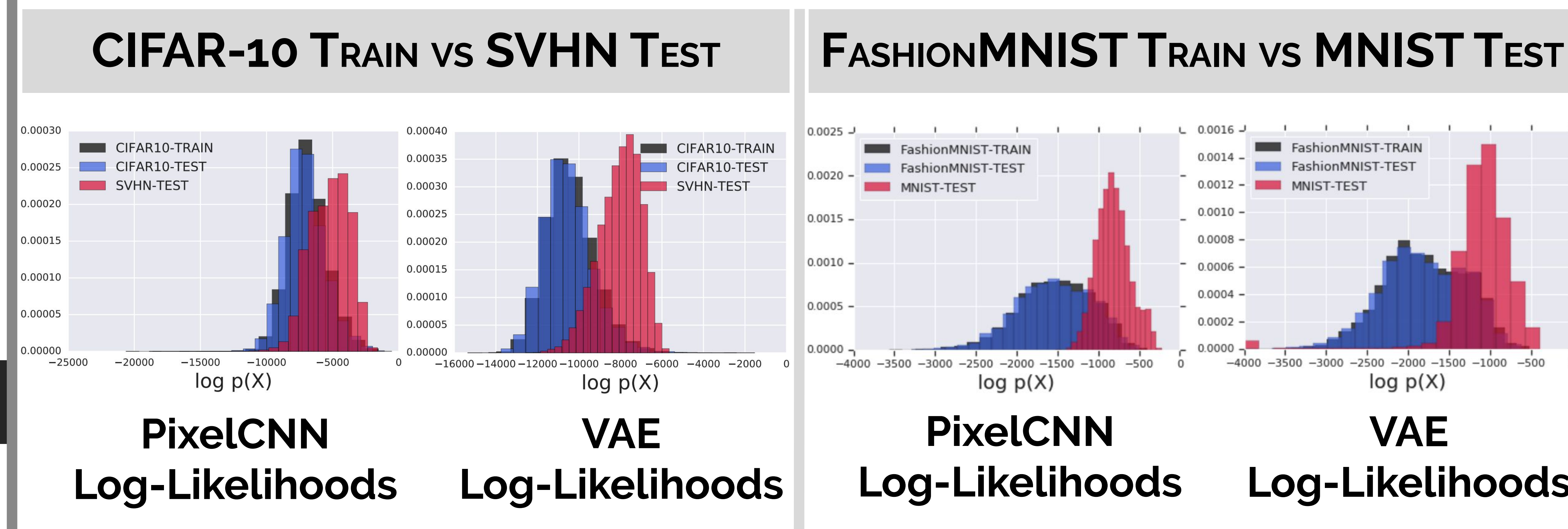
The **phenomenon is asymmetric w.r.t. datasets**: Training on SVHN and evaluating on CIFAR-10 results in the expected ordering (SVHN is assigned higher likelihood).



Glow Log-Likelihoods: SVHN train, CIFAR-10 test

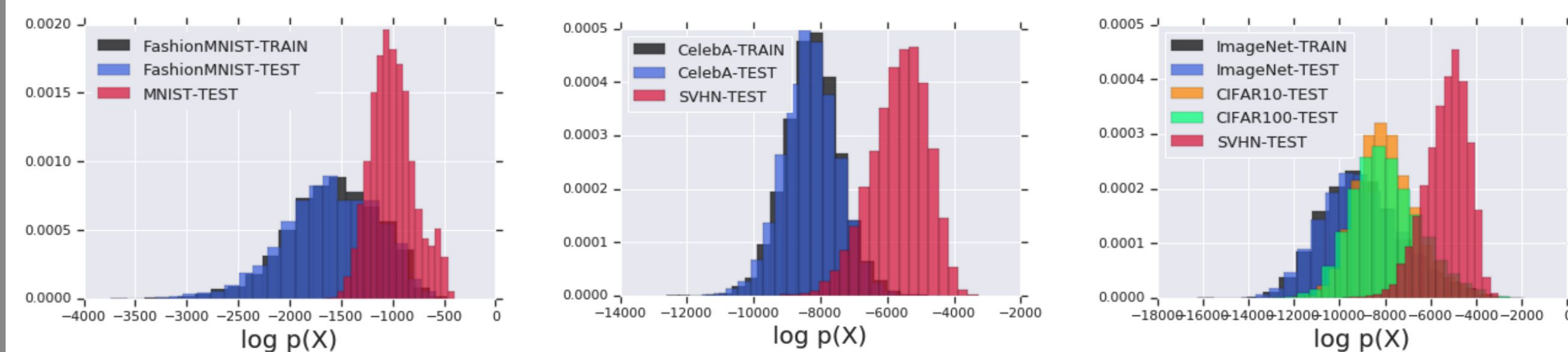
3. TESTING OTHER DEEP GENERATIVE MODEL CLASSES

This **phenomenon is also observed in two other classes of deep generative models**: auto-regressive (PixelCNN) and latent variable models (Variational Auto-Encoders).



4. TESTING GLOW ON OTHER DATA SETS

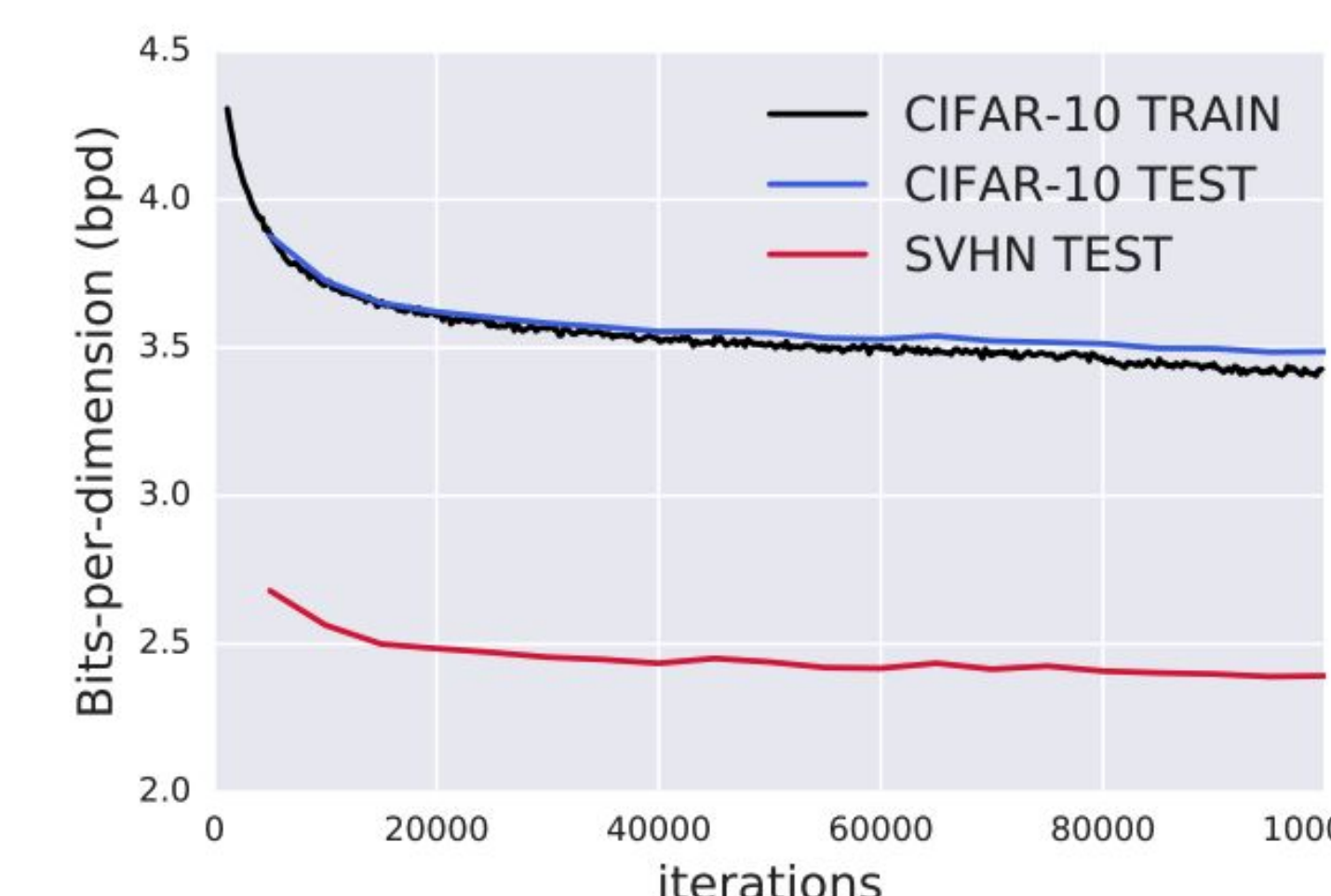
We find further evidence of the phenomenon in five other data set pairs:



FashionMNIST vs MNIST CelebA vs SVHN ImageNet vs SVHN / CIFAR

- We also observe that **constant inputs have the highest log-likelihood of any (tested) input**.
- Furthermore, we find that **SVHN has higher likelihood over the entire duration of training**.
- Ensembling generative models does not help.

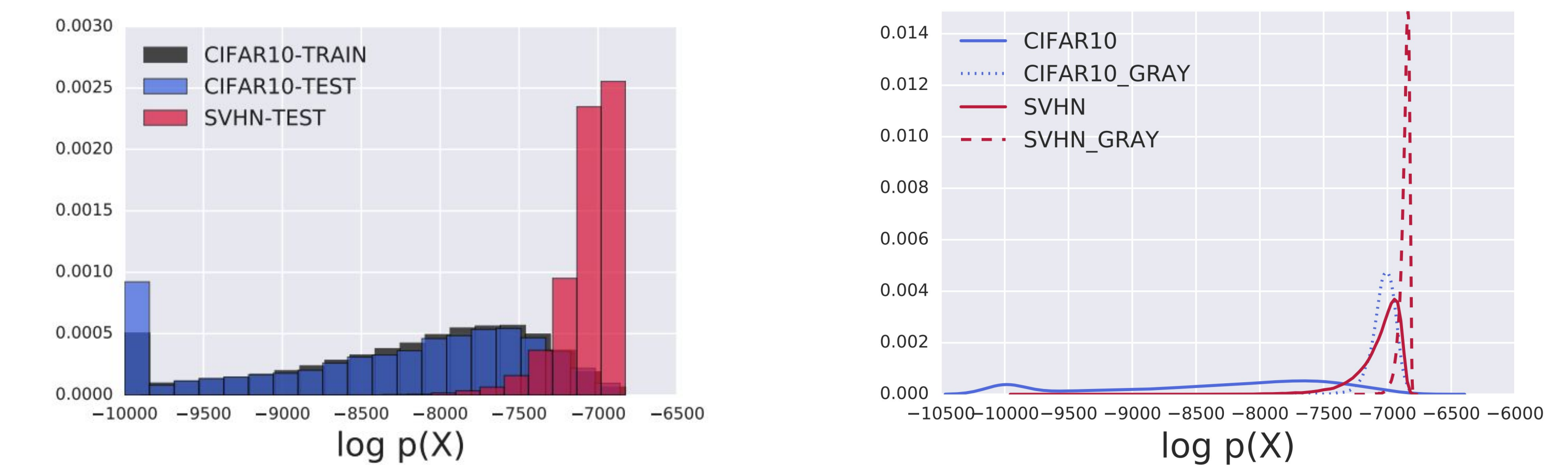
Data Set	Avg. Bits Per Dimension	Data Set	Avg. Bits Per Dimension
<i>Glow Trained on FashionMNIST</i>		<i>Glow Trained on CIFAR-10</i>	
Random	8.686	Random	15.773
Constant (0)	0.339	Constant (128)	0.589



BPD vs Training Iteration

5. DIGGING DEEPER INTO GLOW

To make theoretical analysis more tractable, we **restrict Glow to have constant volume (CV) transformations (w.r.t. input)**. We see similar CIFAR-vs-SVHN results for this model.



CV-Glow: CIFAR-10 vs SVHN

Graying Images Increases Likelihoods

For CV-Glow, we can approximate the difference in likelihoods between the training and OOD data as follows:

$$0 < \mathbb{E}_q[\log p(\mathbf{x}; \theta)] - \mathbb{E}_{p^*}[\log p(\mathbf{x}; \theta)]$$

$$\approx \frac{\partial^2}{\partial z^2} \log p(z; \psi) \sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2 \sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Non-Training Distribution Training Distribution Second Moment of Non-Training Distribution Second Moment of Training Distribution

< 0 for all log-concave densities (e.g. Gaussian) Non-negative due to square

This expression helps explain several observations:

- Asymmetry**: difference between 2nd moments does not commute.
- Constant / grayscale inputs**: equivalent to non-training moment being zero. **Graying images increases likelihoods**.
- Early stopping / ensembling would not help**: expression holds true for all values of CV-Glow's parameters.

6. SUMMARY

Density estimates from (current) deep generative models are **not always able to detect out-of-distribution inputs**.

Paper: <https://arxiv.org/abs/1810.09136>

