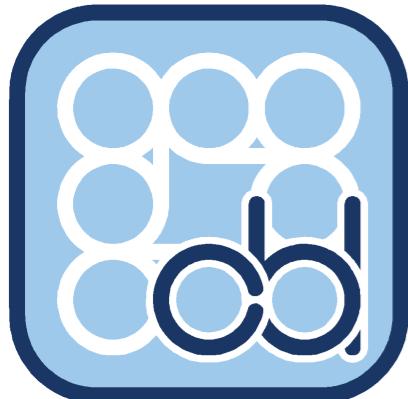
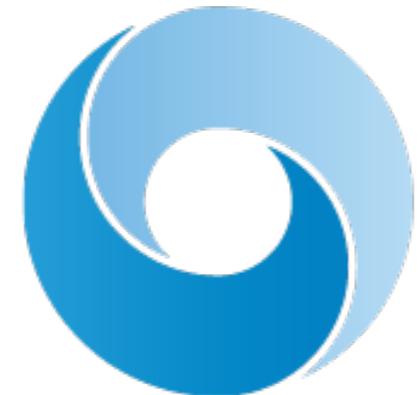

Evaluating Deep Generative Models on Out-of-Distribution Inputs

Eric Nalisnick

OxCML Seminar
31.5.19



Computational and
Biological Learning
University of Cambridge



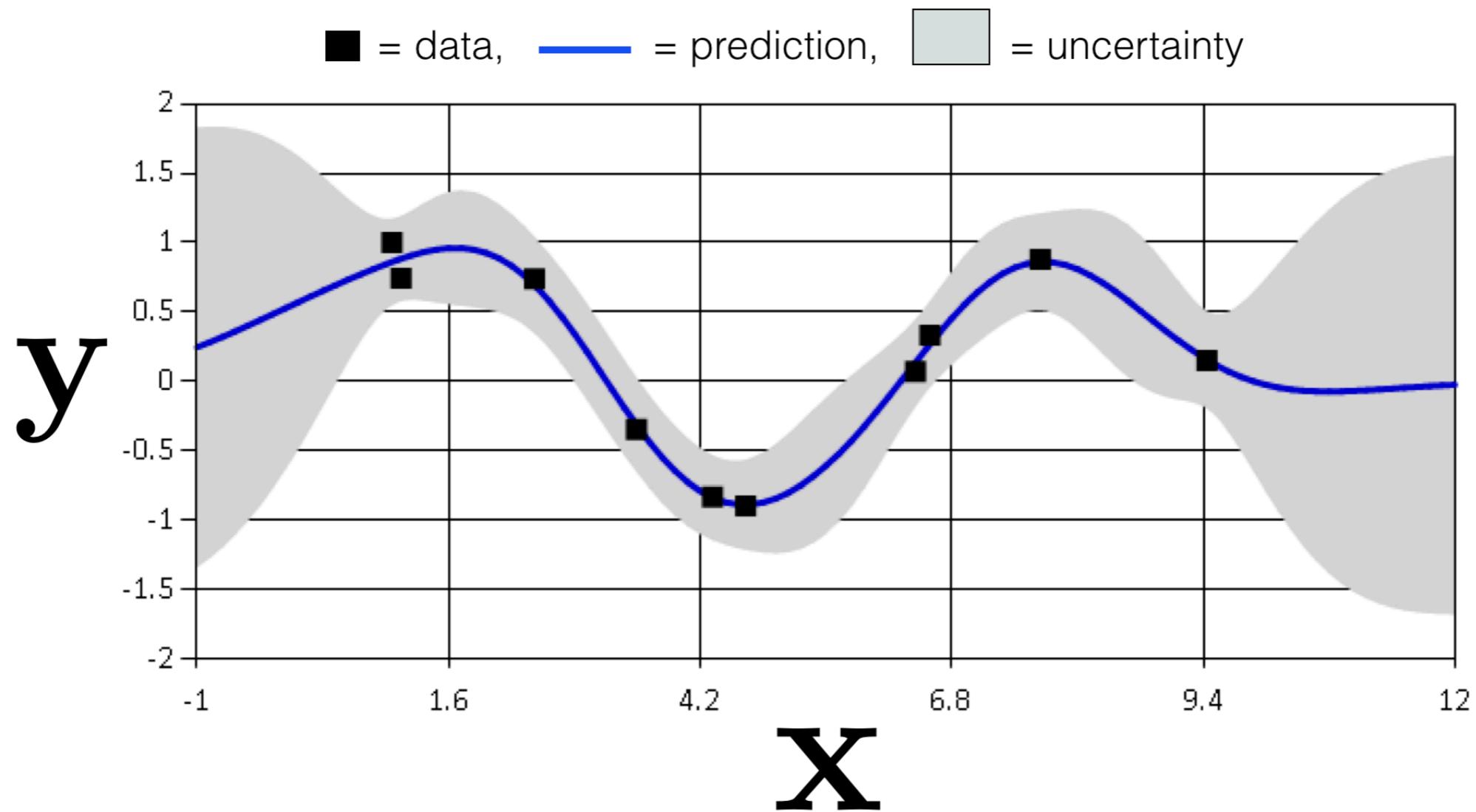
DeepMind

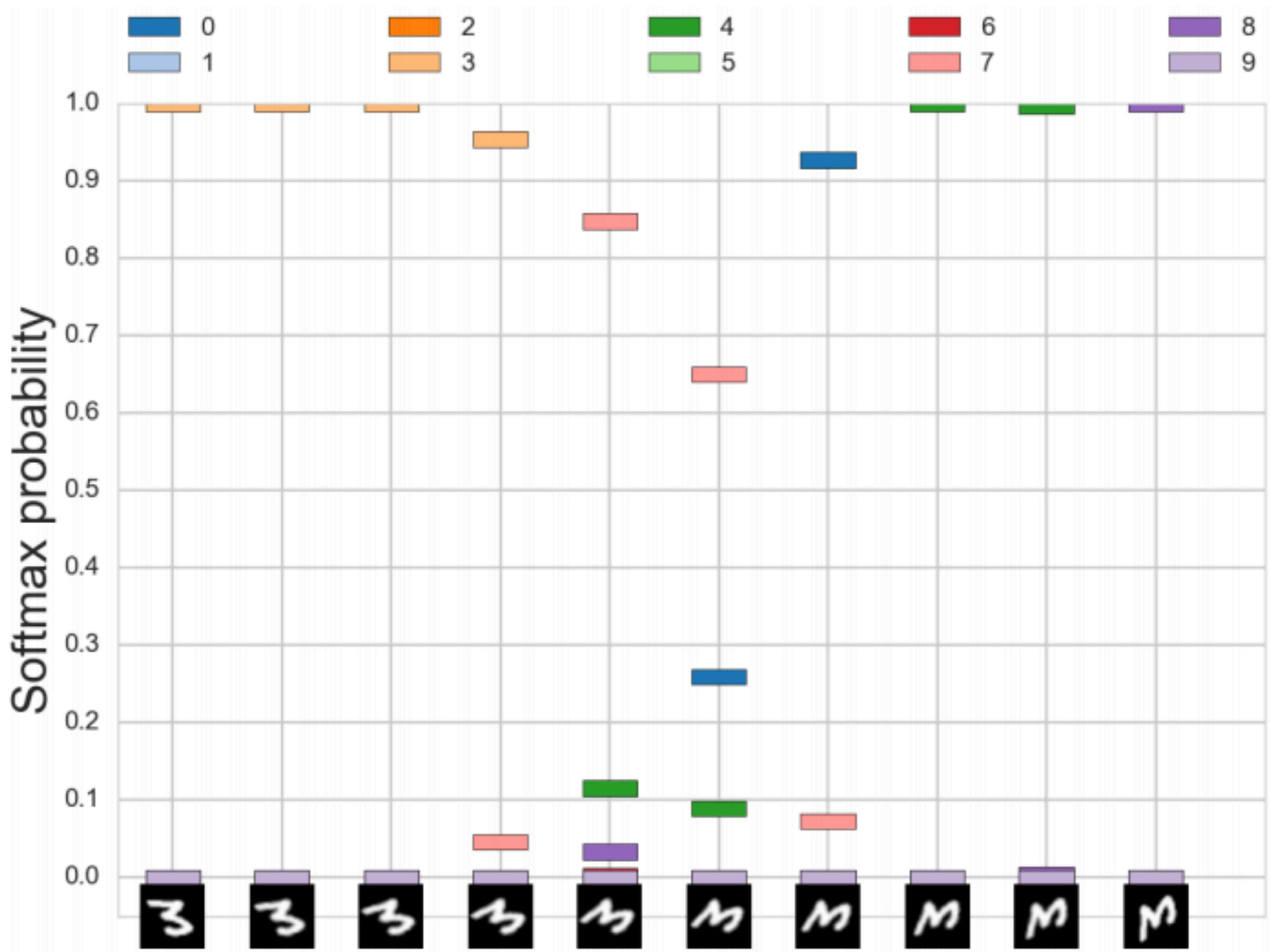
PART #1

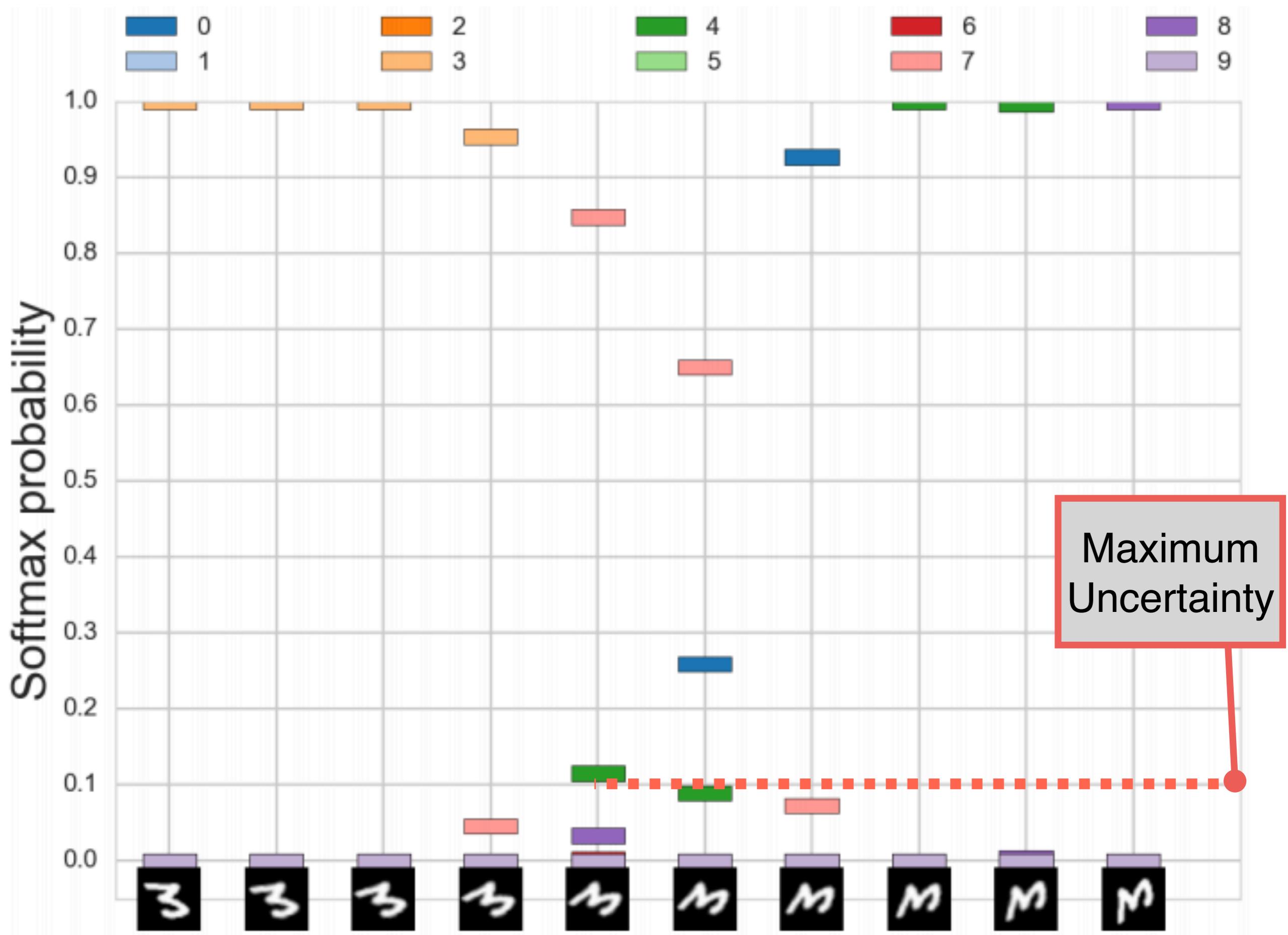
Validating Predictive Models

$$p(\mathbf{y} | \mathbf{X}; \theta)$$

Labels Features Parameters



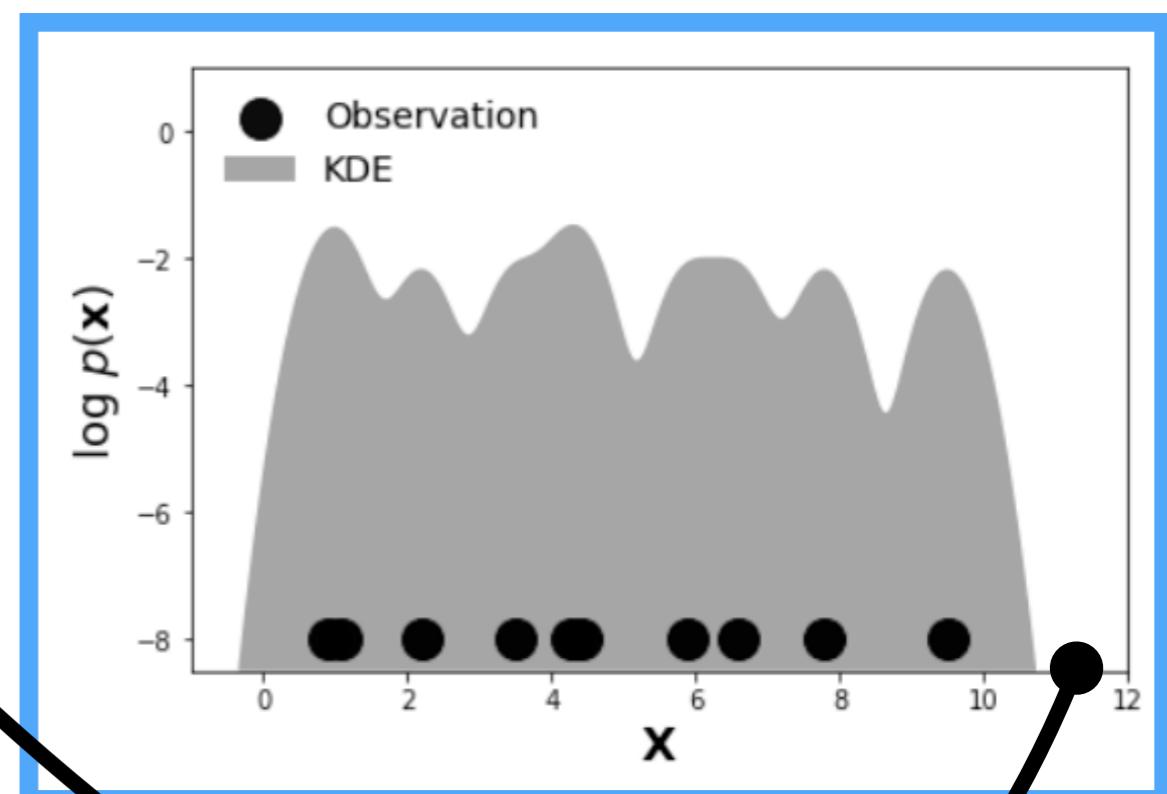
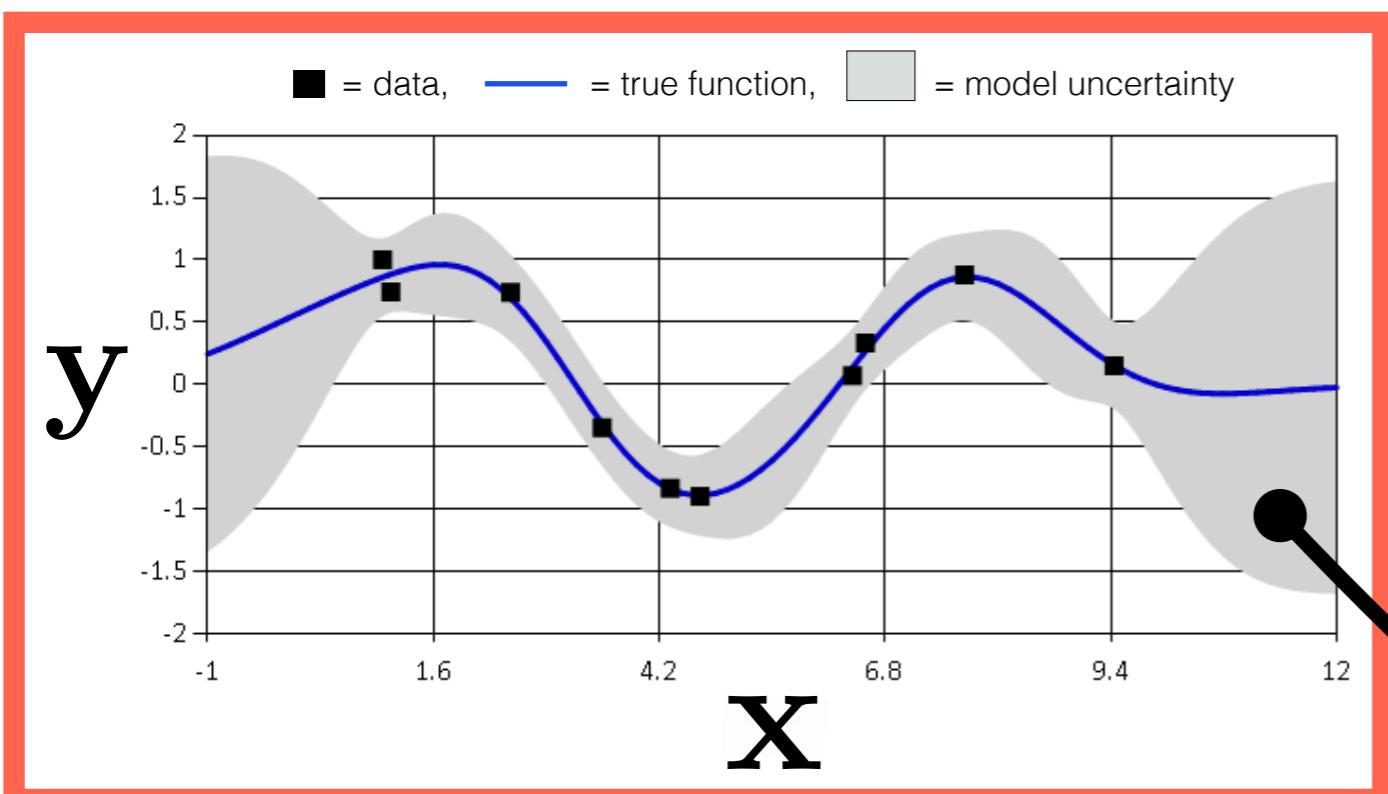




$$p(\mathbf{y}, \mathbf{X}) = p(\mathbf{y}|\mathbf{X}; \boldsymbol{\theta}) \frac{p(\mathbf{X}; \boldsymbol{\phi})}{p(\mathbf{X}; \boldsymbol{\phi})}$$

Predictive
Model

Generative
Model

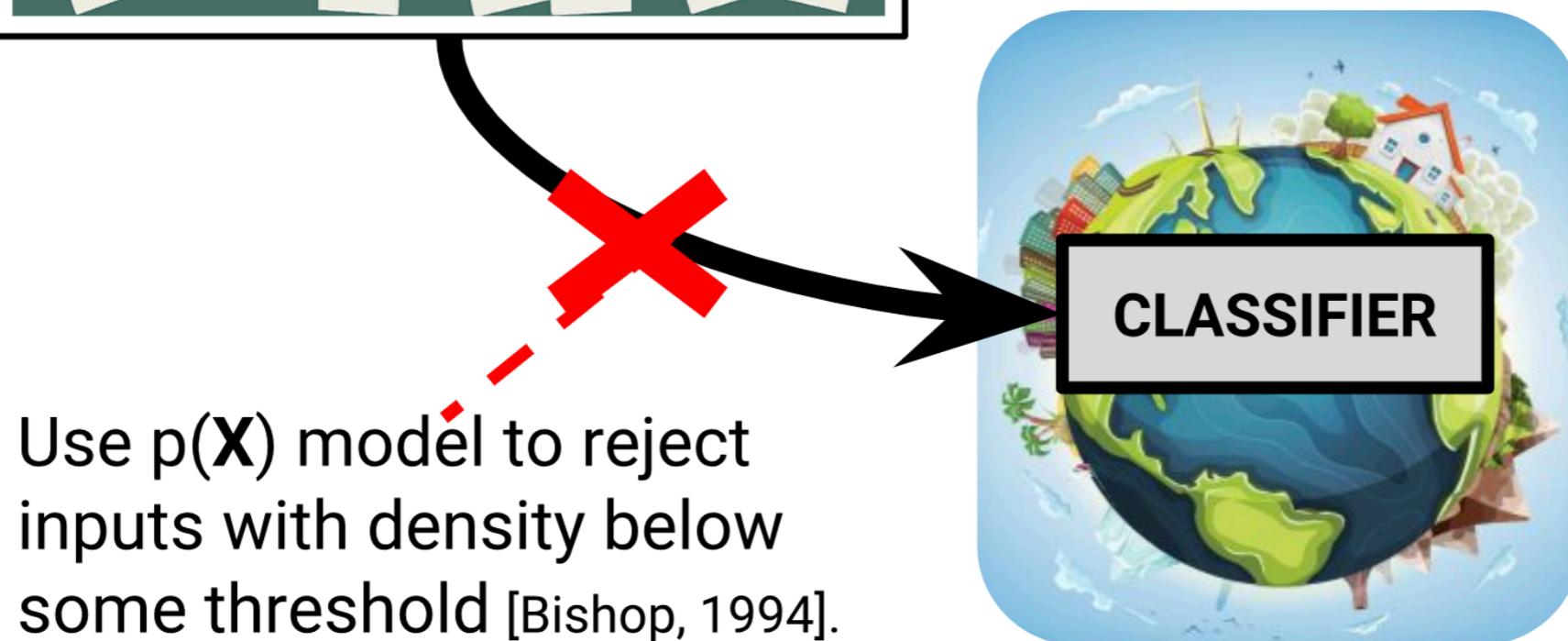


Threshold-Based Rejection

Inputs Unlike Training Data



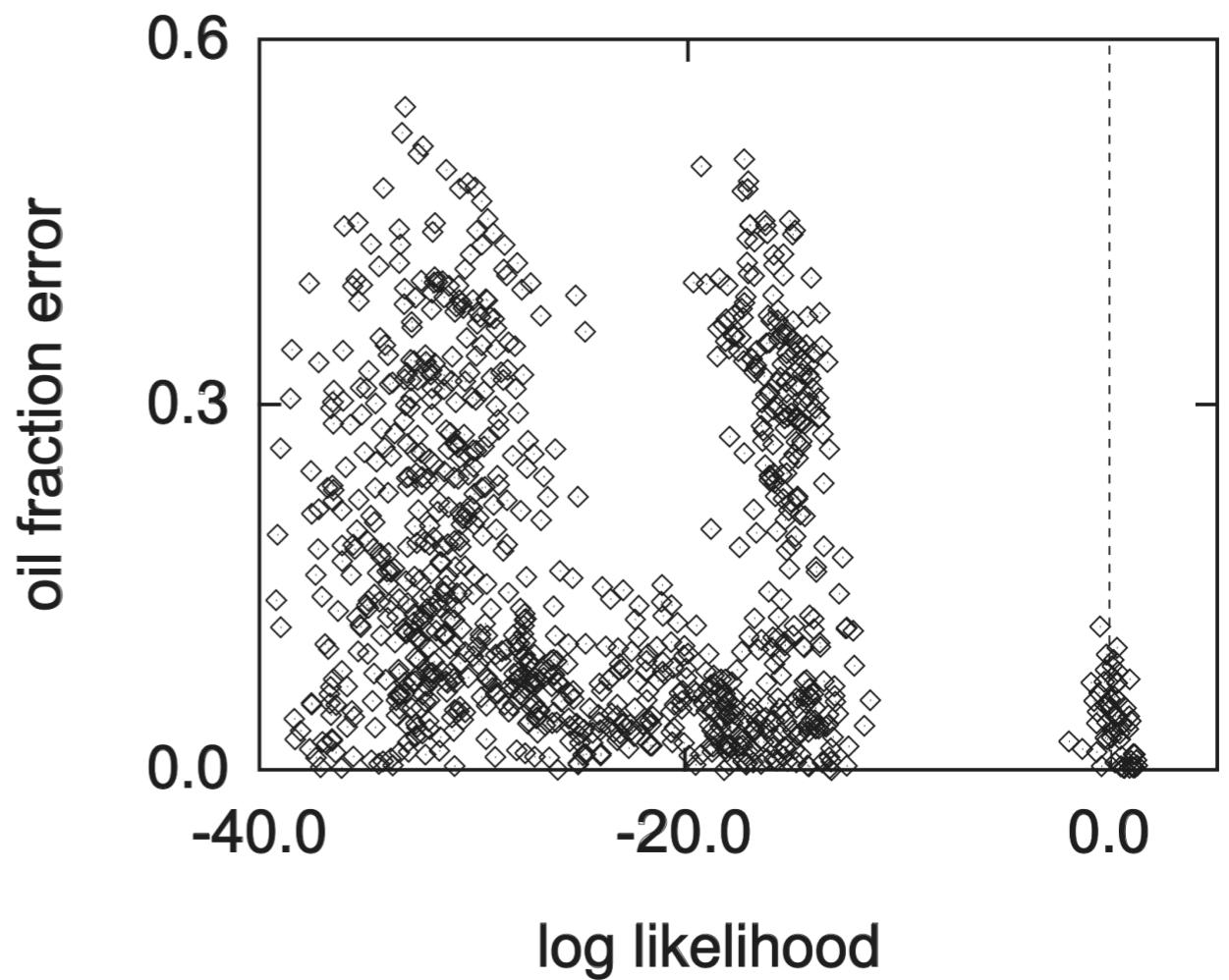
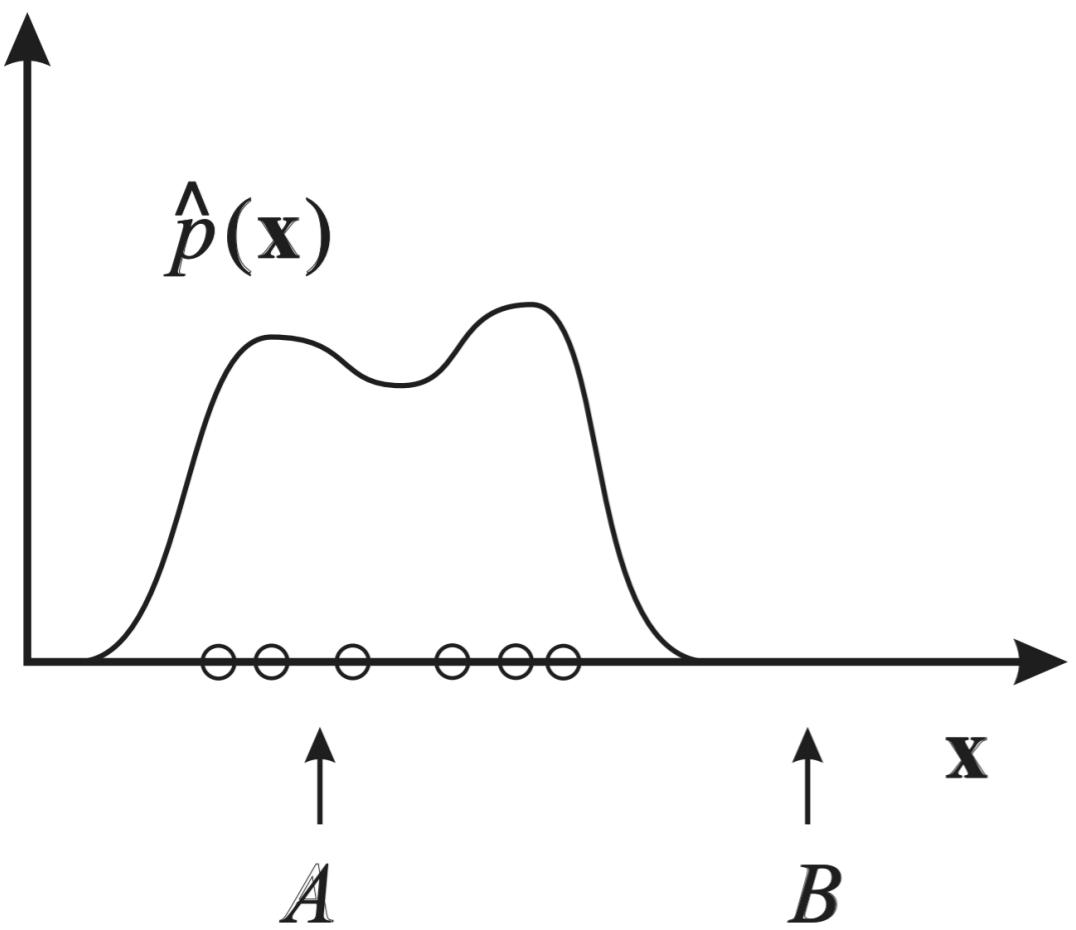
if $p(x^*; \phi) < \tau$,
then reject x^*





Novelty Detection and Neural Network Validation

Chris M. Bishop (May 1994)



ZOUBIN: [The Bishop (1994) procedure] should be built into the software.

MODERATOR: Isn't that hard?

ZOUBIN: If you stick a picture of a chicken into an MNIST classifier, it should tell you it's neither a seven nor a one.

[AUDIENCE LAUGHS]

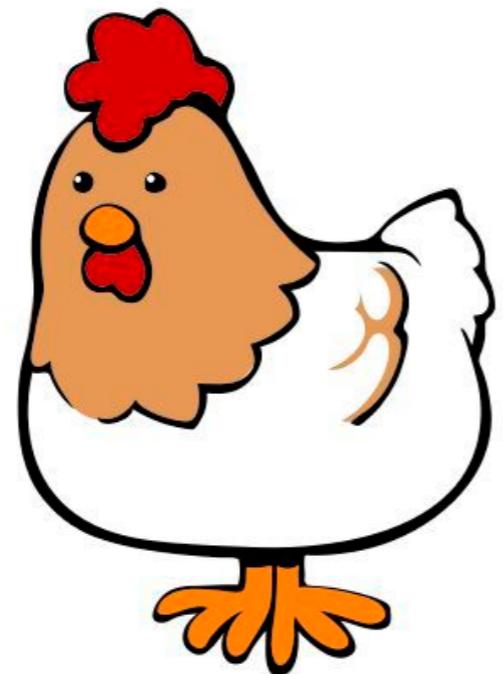


Panel Discussion

Advances in Approximate Bayesian Inference, Dec 2017

PART #2

Chicken or Seven?



Deep Generative Models

■ Autoregressive Models (e.g. WaveNet)

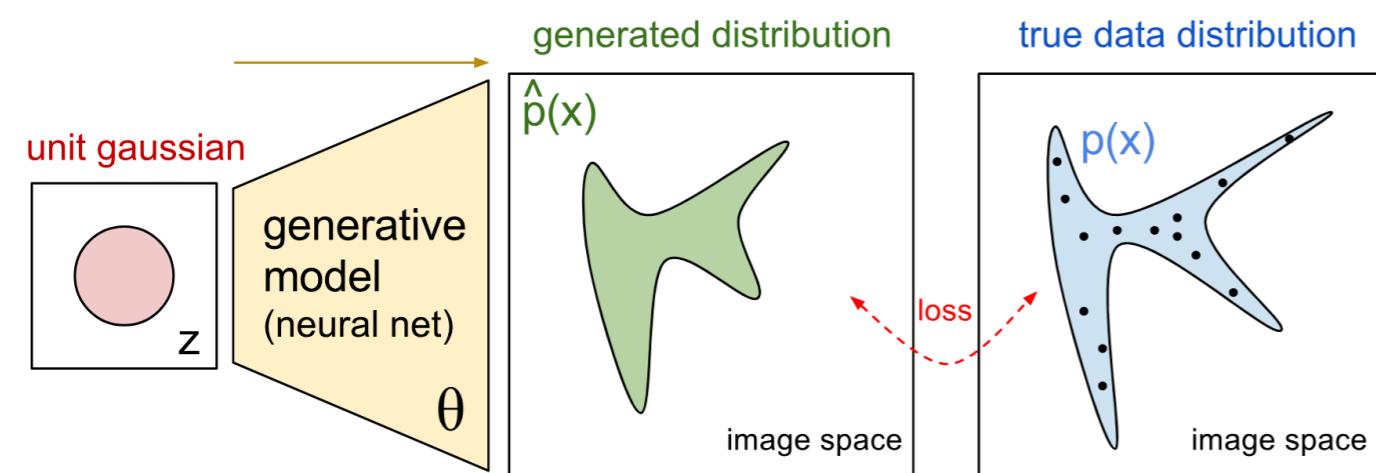
[Larochelle & Murray, AISTATS 2011; van den Oord, ICLR 2016; van den Oord, NeurIPS 2016]

■ Variational Autoencoders (VAEs)

[Kingma & Welling, ICLR 2014; Rezende et al., ICML 2014]

■ Normalizing Flows

[Tabak & Turner, 2013; Rezende & Mohamed, ICML 2015; Dinh et al., ICLR 2017]

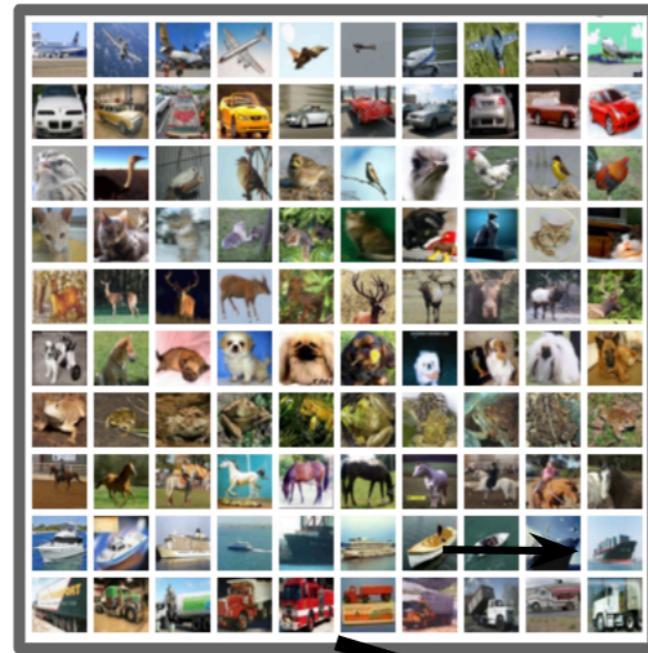


<https://blog.openai.com/generative-models/>



Chicken or Seven?

Training: CIFAR-10



Testing: SVHN



$$p(\mathbf{x}_{\text{CIFAR-10}}) > p(\mathbf{x}_{\text{SVHN}})$$

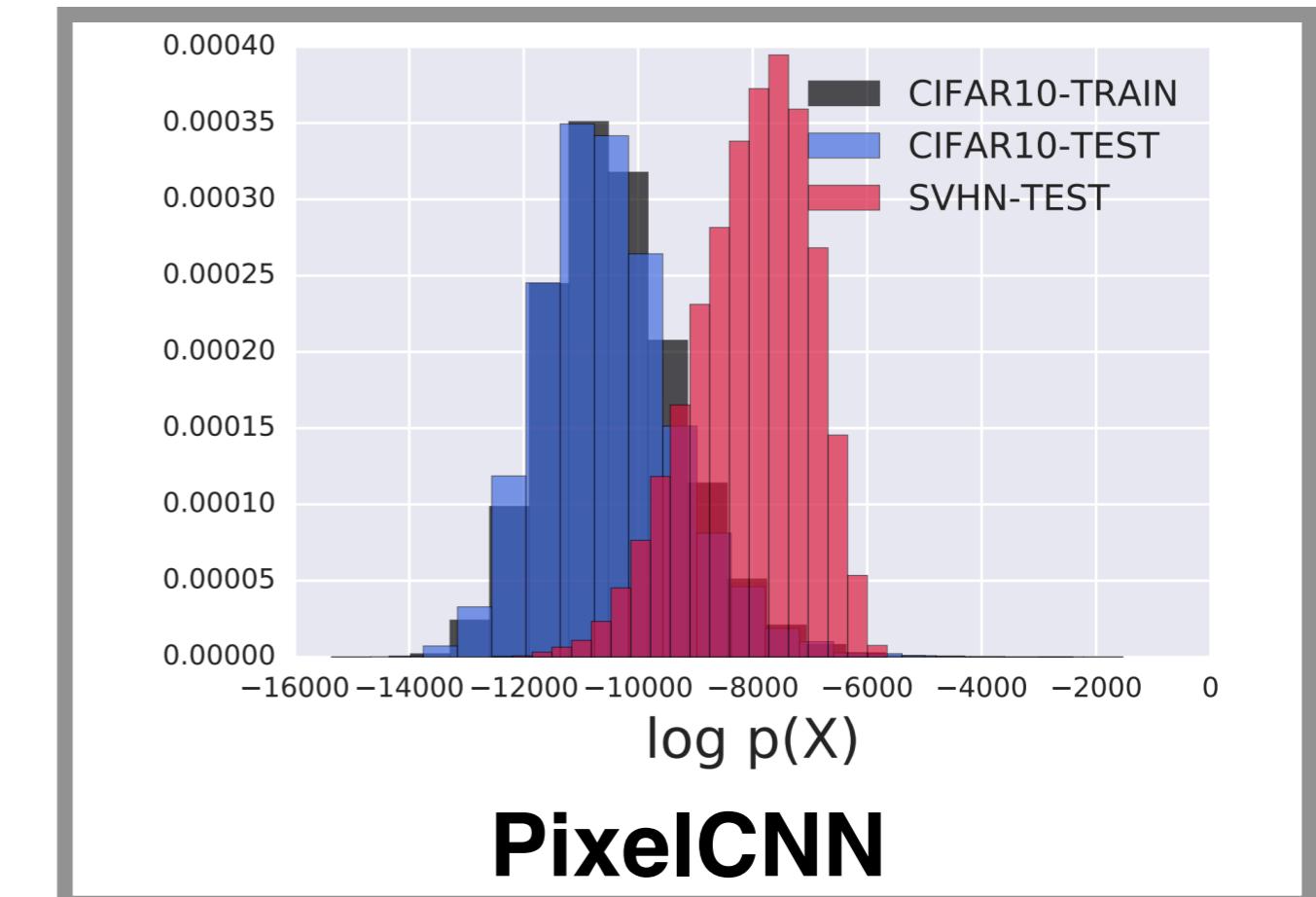
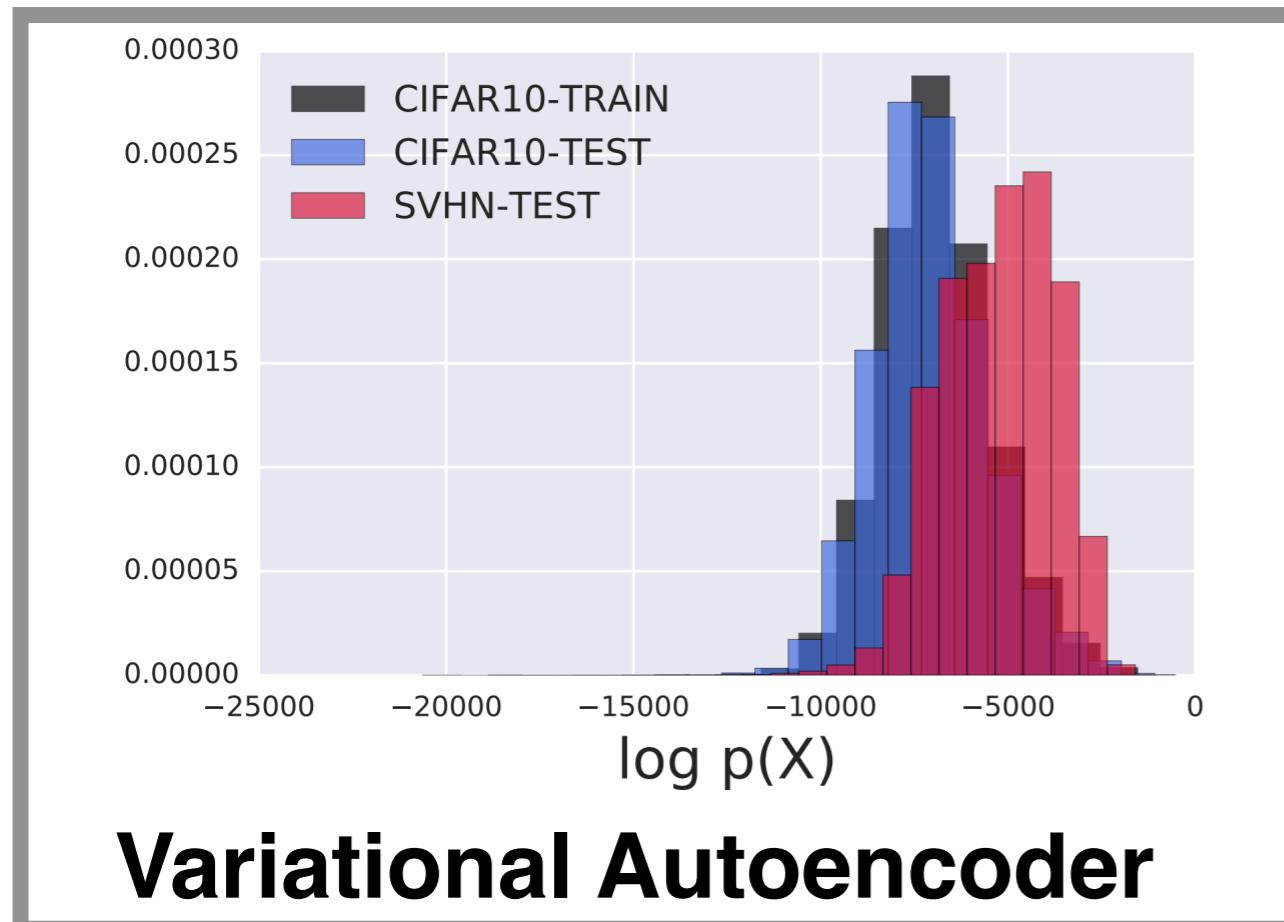
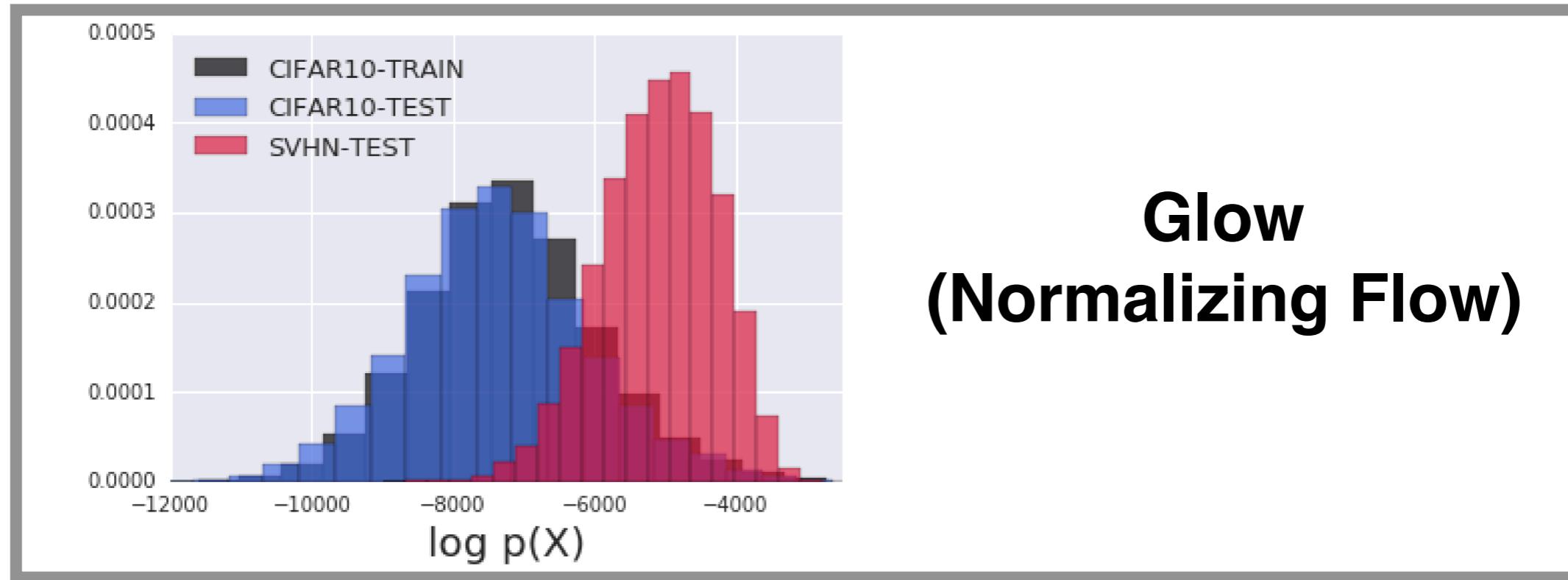
A large red question mark is placed between the two probability terms.

CIFAR-10 vs SVHN

First let's try a normalizing flow, Glow [Kingma & Dhariwal, 2018]:

Data Set	Avg. Bits Per Dimension	(lower is better)
<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	

CIFAR-10 vs SVHN

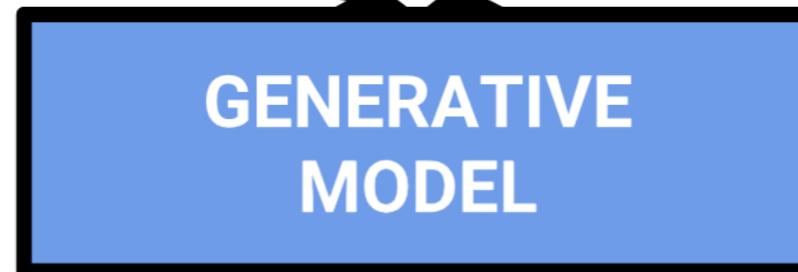


Shirt or Seven?

Training: *FashionMNIST* Testing: *MNIST*



3	4	2	1	9	5	6	2	1	8
8	9	1	2	5	0	0	6	6	4
6	7	0	1	6	3	6	3	7	0
3	7	7	9	4	6	6	1	8	2
2	9	3	4	3	9	8	7	2	5
1	5	9	8	3	6	5	7	2	3
9	3	1	9	1	5	8	0	8	4
5	6	2	6	8	5	8	8	9	9
3	7	7	0	9	4	8	5	4	3
7	9	6	4	7	0	6	9	2	3



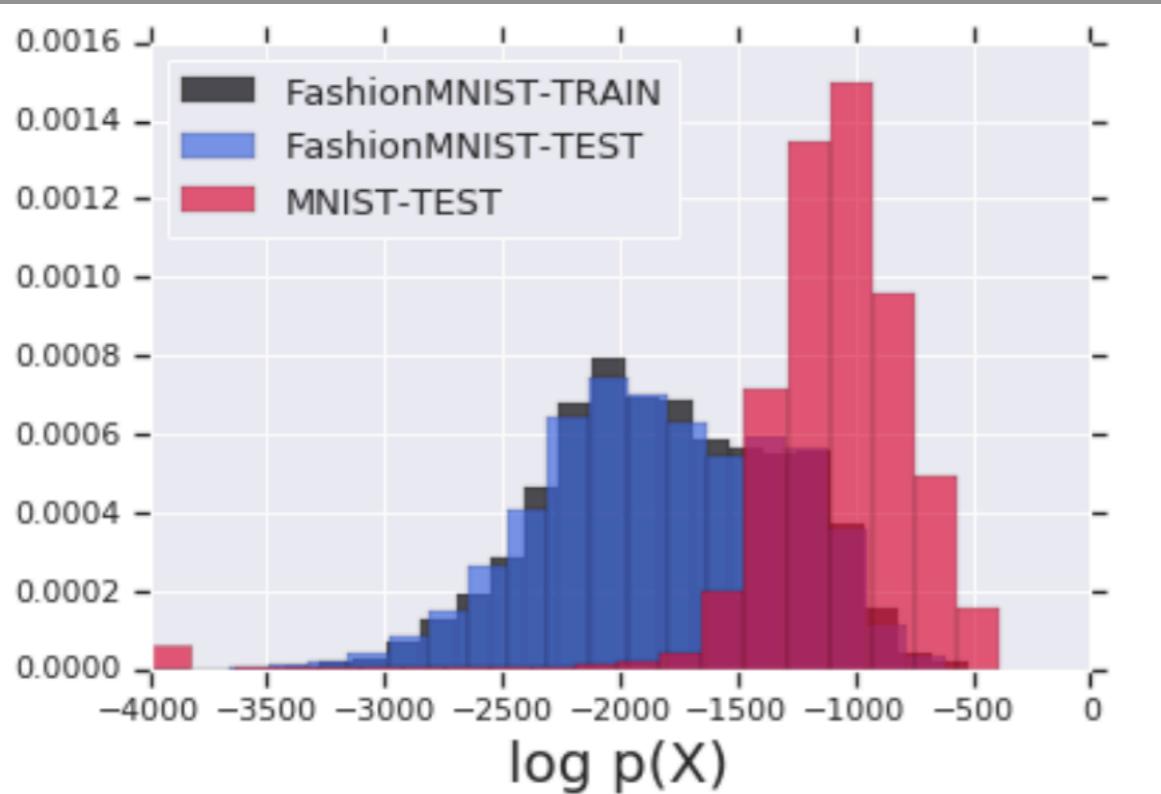
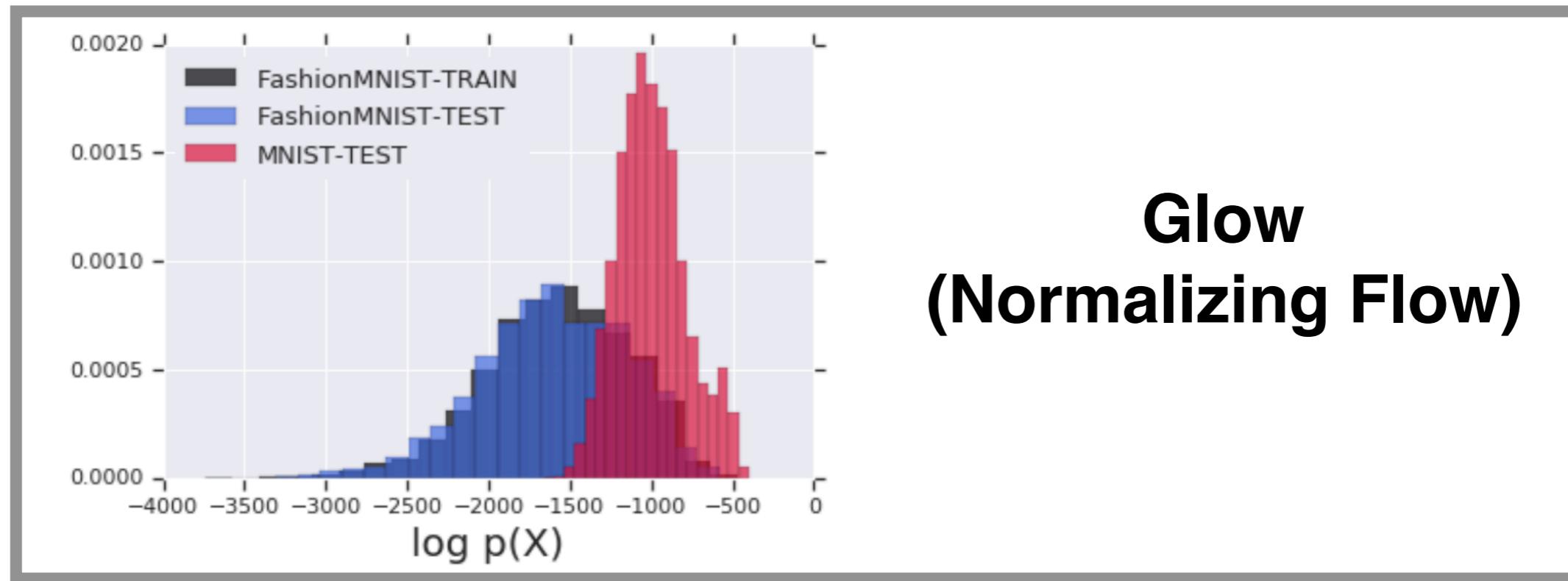
$$p(\mathbf{x}_{\text{FashionMNIST}}) > p(\mathbf{x}_{\text{MNIST}})$$

FashionMNIST vs MINST

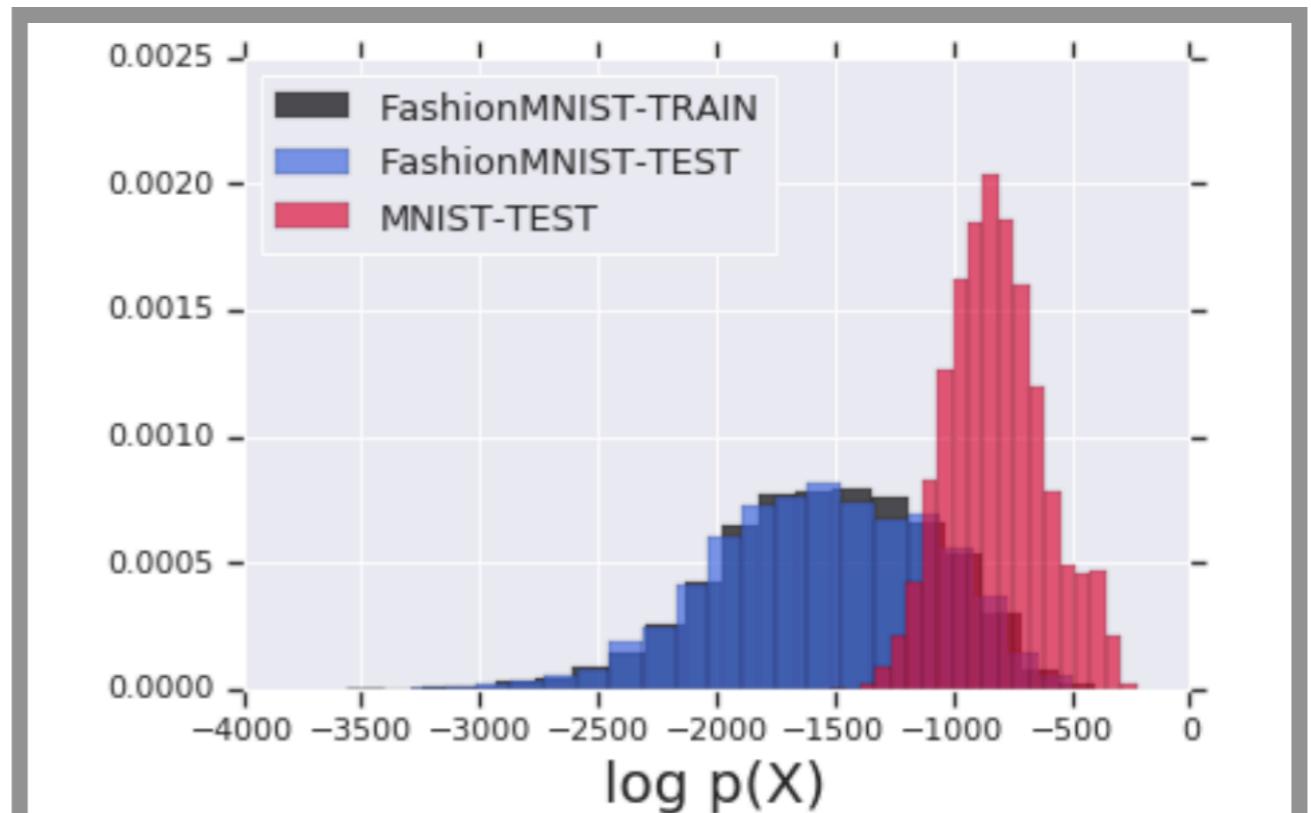
Again, let's first try Glow [Kingma & Dhariwal, 2018]:

Data Set	Avg. Bits Per Dimension	(lower is better)
<i>Glow Trained on FashionMNIST</i>		
FashionMNIST-Train	2.902	
FashionMNIST-Test	2.958	
MNIST-Test	1.833	
<i>Glow Trained on MNIST</i>		
MNIST-Test	1.262	

FashionMNIST vs MNIST



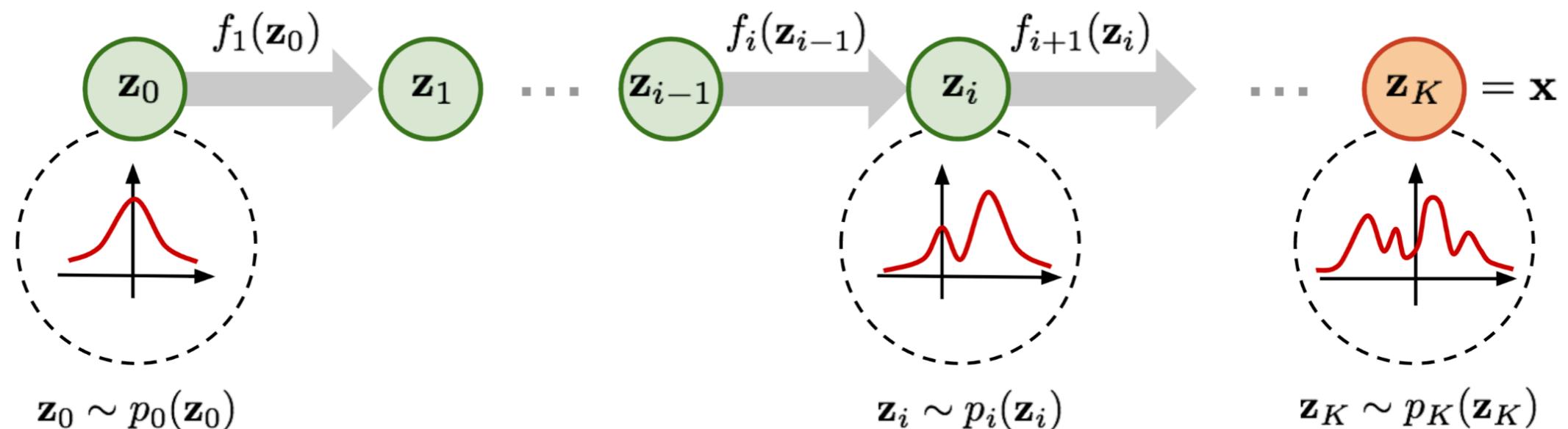
Variational Autoencoder



PixelCNN

PART #3

Digging Deeper into Flow-Based Models



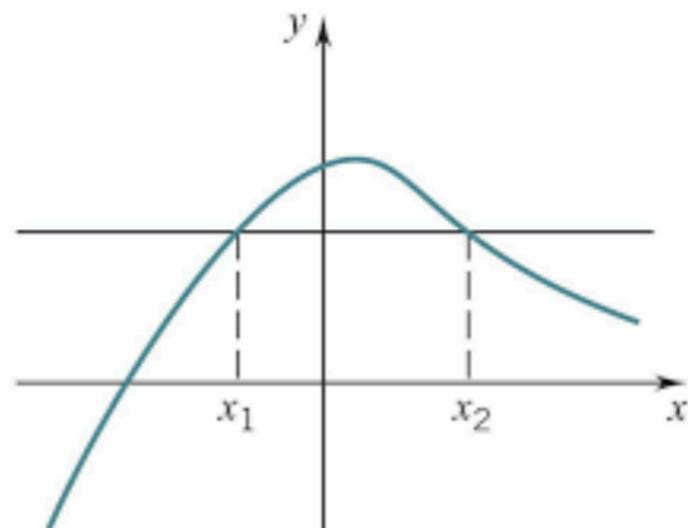
Normalizing Flows: Background

Change of Variables Formula ($X \rightarrow Z$):

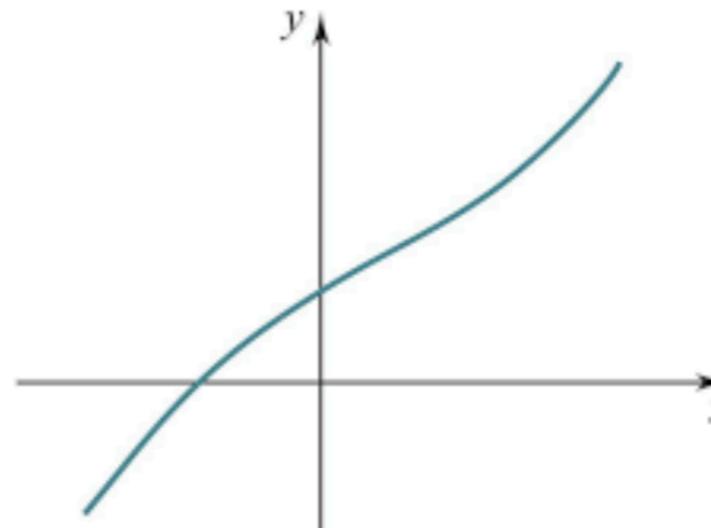
$$p_z(f(X)) \left| \frac{df(X)}{dX} \right| = p(X)$$

So what's the catch?

$f(x)$ must be a *bijection*

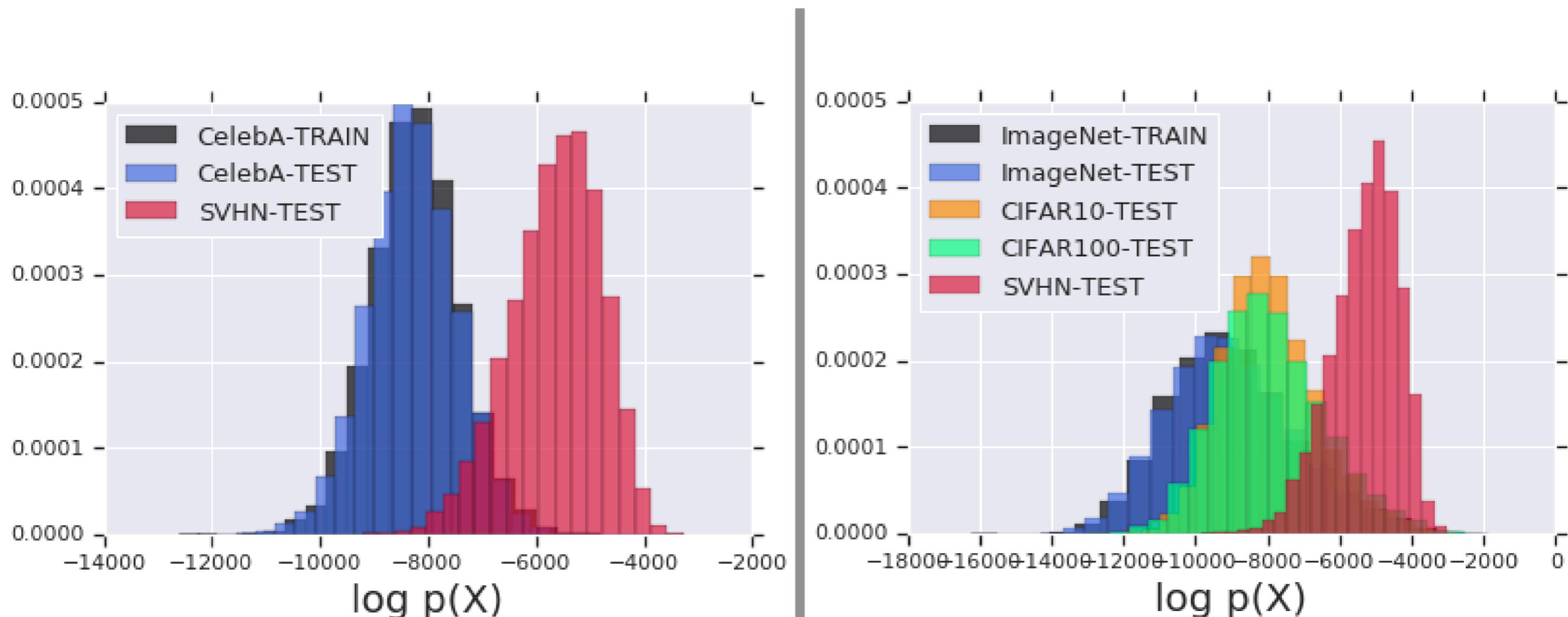


f is not one-to-one: $f(x_1) = f(x_2)$



f is one-to-one:

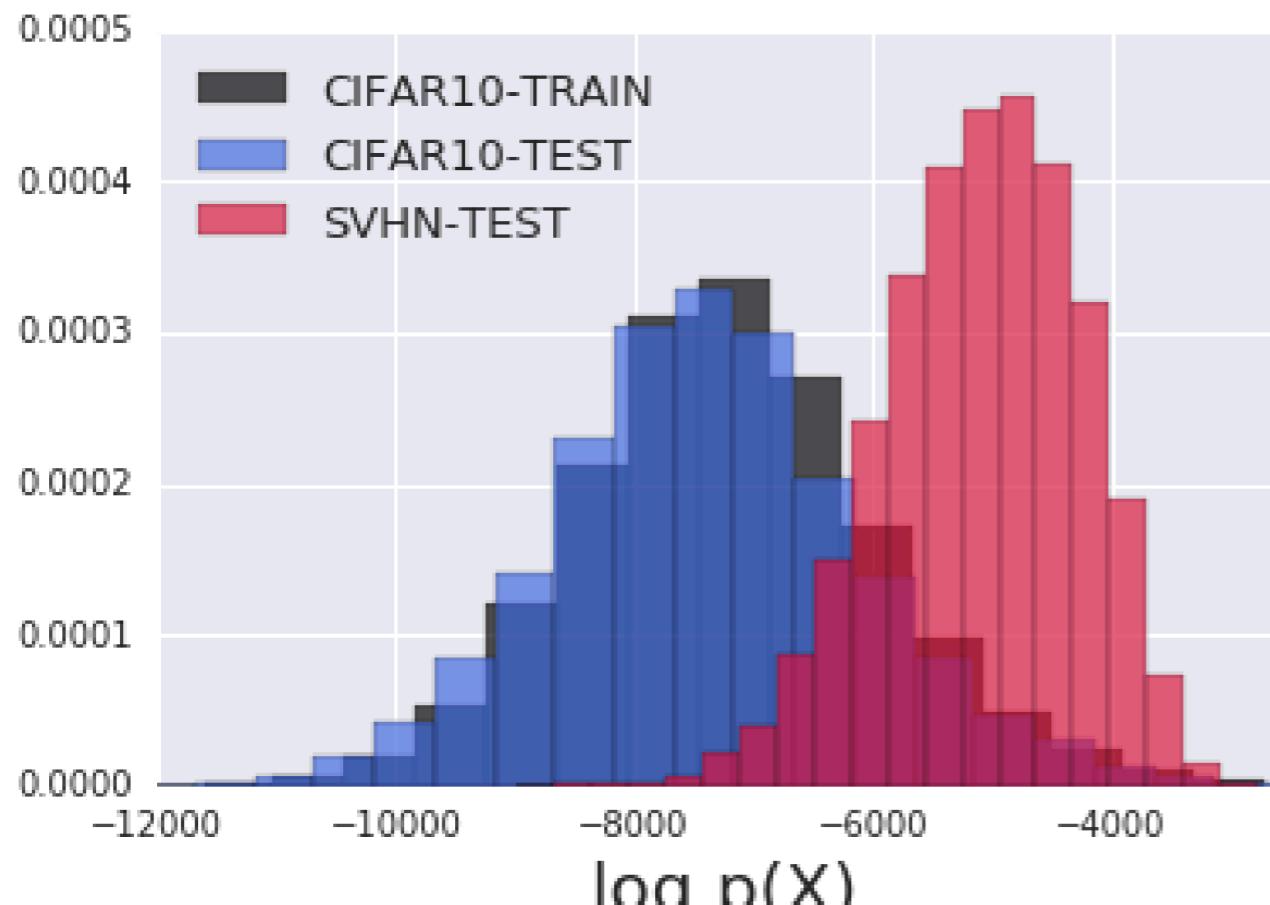
Additional Data Sets



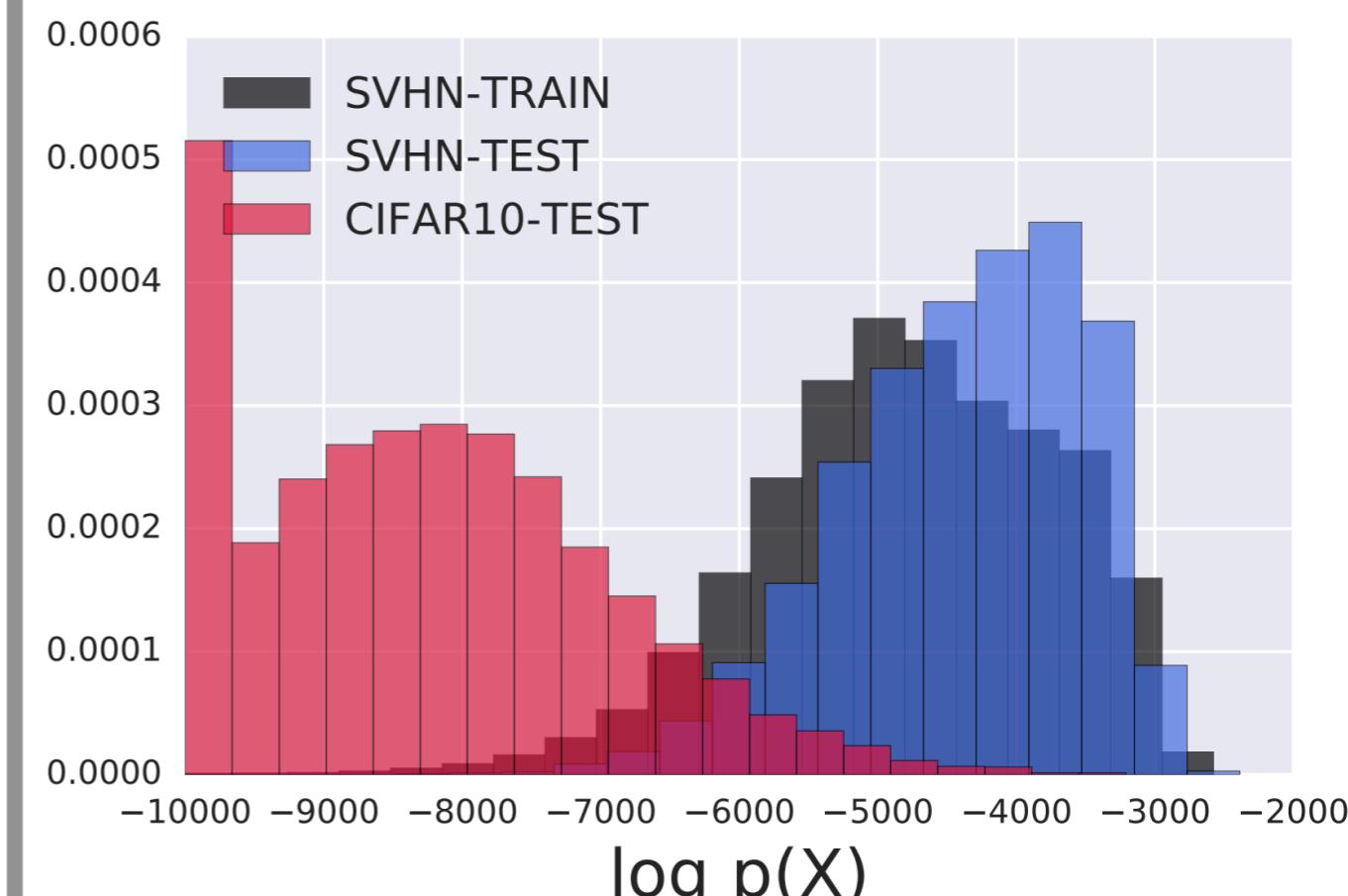
CelebA vs SVHN

**ImageNet vs CIFAR-10
vs SVHN**

The Phenomenon is Not Symmetric



CIFAR-10 vs SVHN

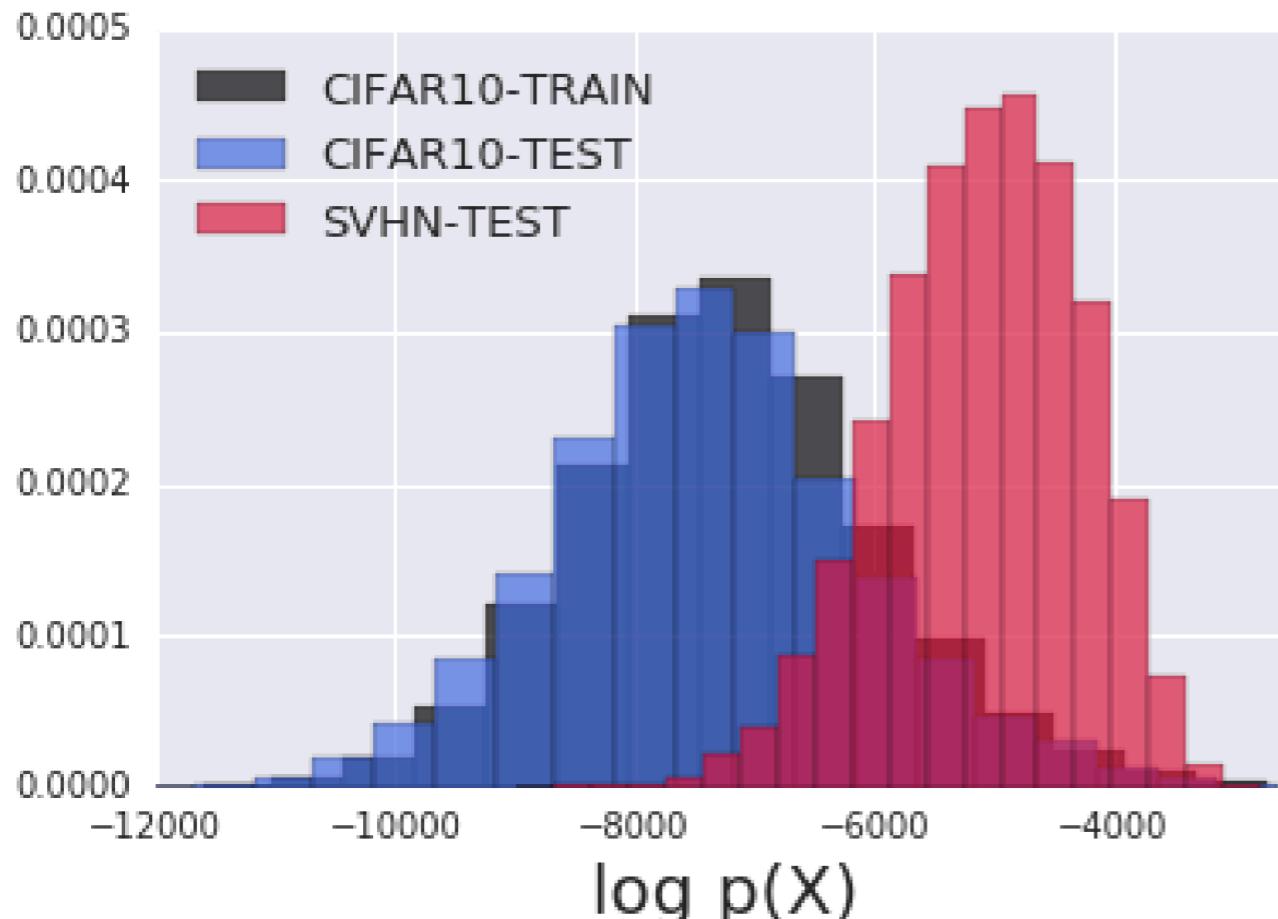


SVHN vs CIFAR-10

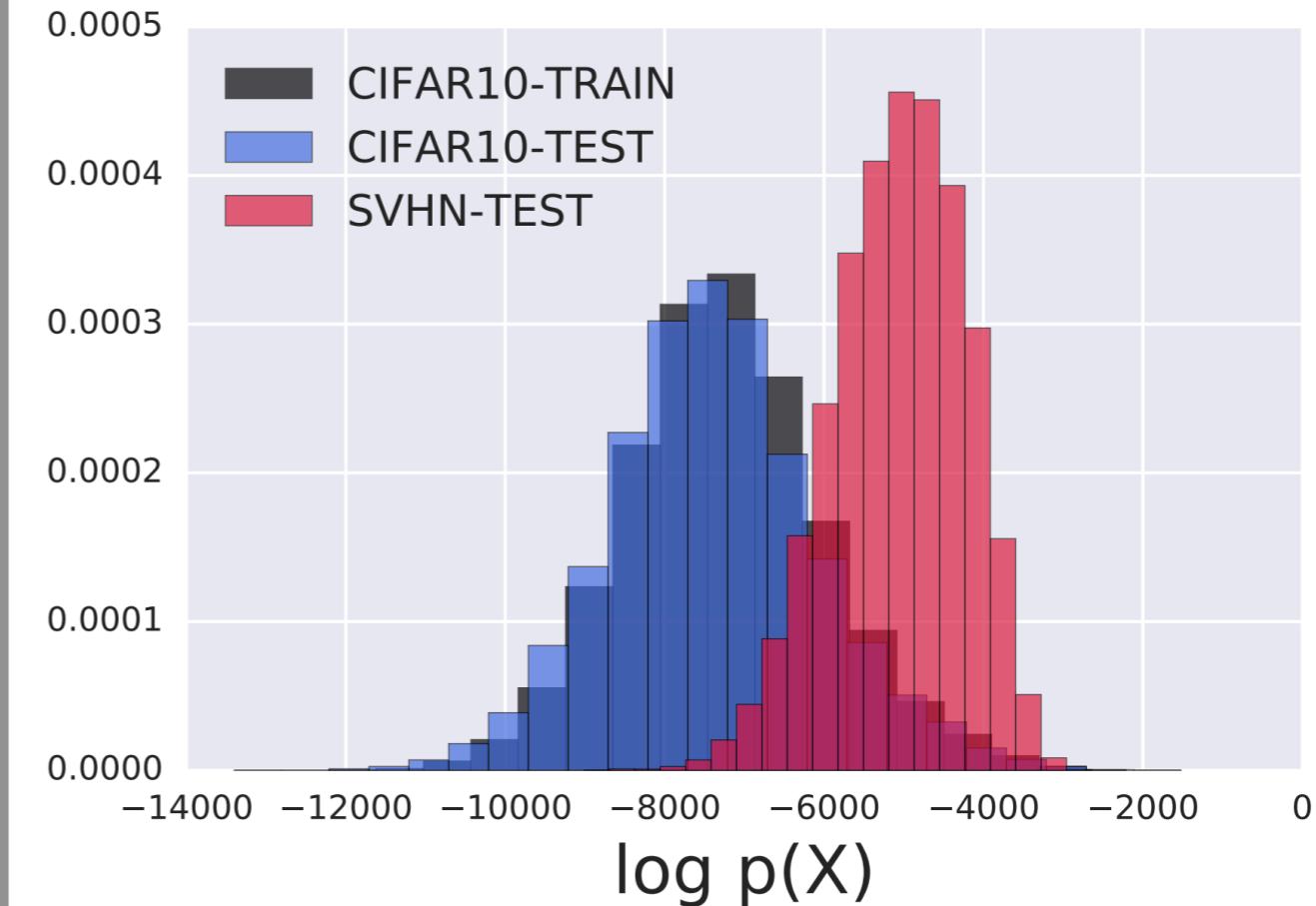
Uniform (Constant) Inputs

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

Ensembling Makes No Difference

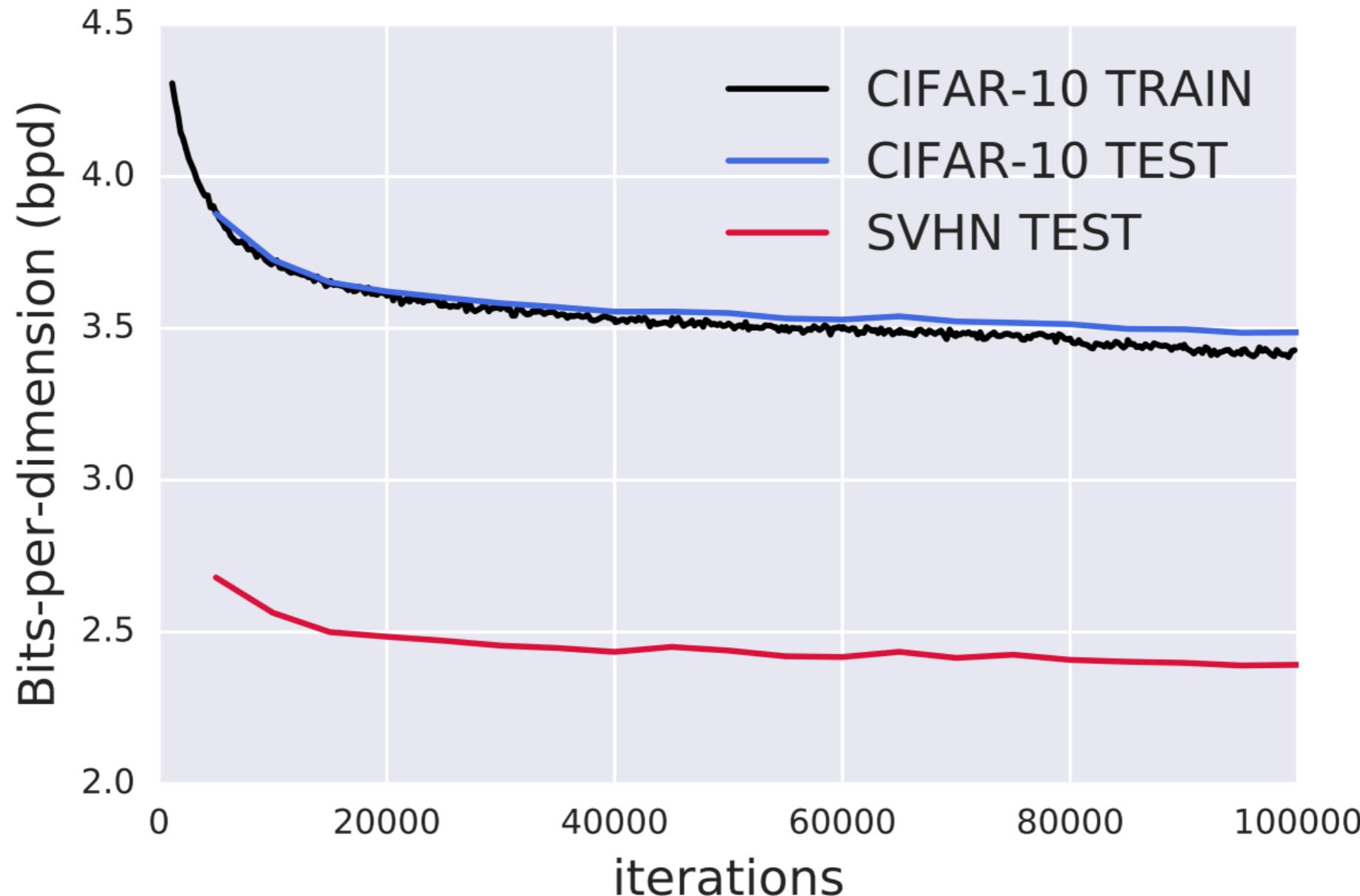


CIFAR-10 vs SVHN
1 Glow



CIFAR-10 vs SVHN
Ensemble of 10 Glows

Early-Stopping Makes No Difference



During Optimization

Switching to an Analytical Approach

Mathematical characterization:

$$0 < \underbrace{\mathbb{E}_q[\log p(\mathbf{x}; \boldsymbol{\theta})]}_{\text{Non-Training Distribution}} - \underbrace{\mathbb{E}_{p^*}[\log p(\mathbf{x}; \boldsymbol{\theta})]}_{\text{Training Distribution}}$$
$$\approx \frac{1}{2} \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p_z(f(\mathbf{x}_0; \boldsymbol{\phi})) + \nabla_{\mathbf{x}_0}^2 \log \left| \frac{\partial f_{\boldsymbol{\phi}}}{\partial \mathbf{x}_0} \right| \right] (\underbrace{\Sigma_q - \overline{\Sigma}_{p^*}}_{\text{Second Moment of Non-Training Distribution}}) \right\}$$

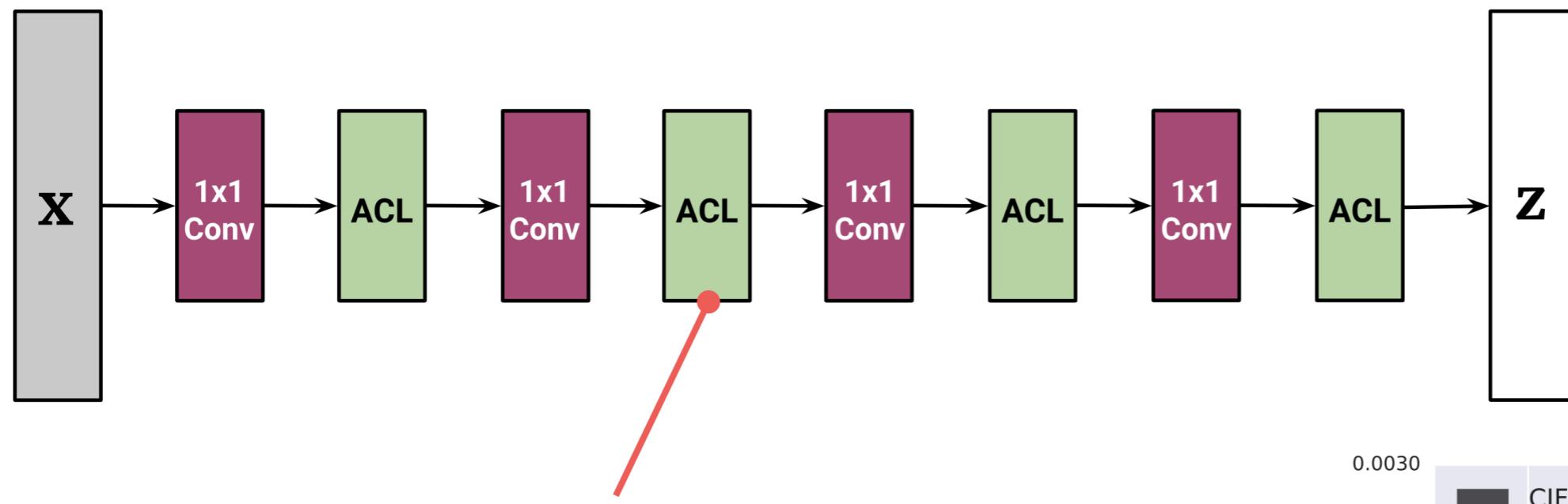
Change-of-Variable Terms

Second
Moment of
Non-Training
Distribution

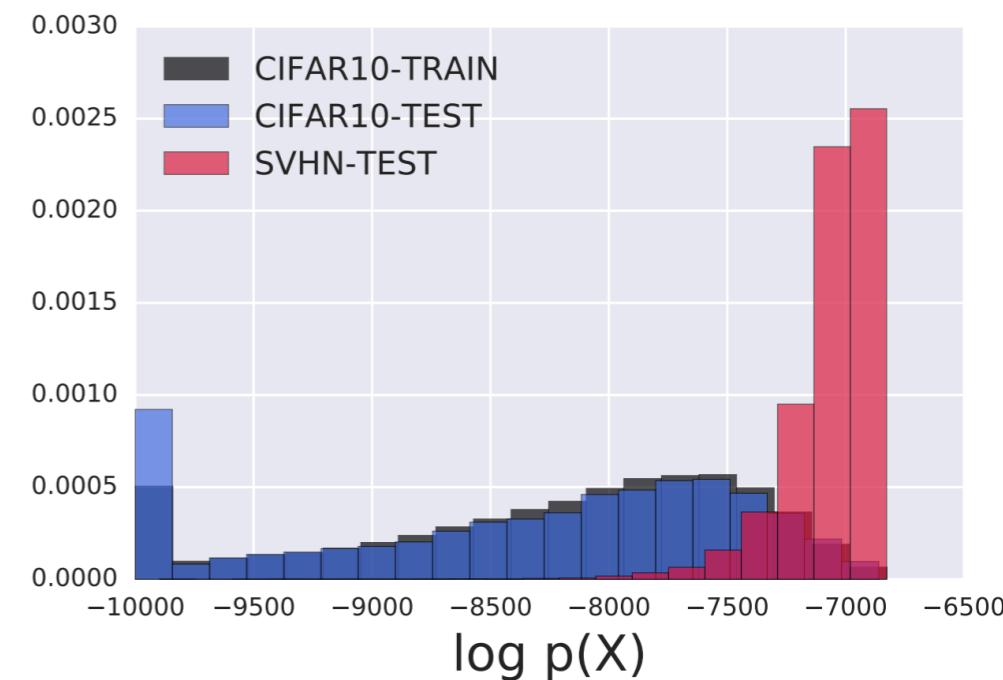
Analysis via Constant-Volume Flows

We define a sub-class we term *constant-volume* (w.r.t. input) flows.

Glow Architecture:



Use only *translation* operations.



CIFAR-10 vs SVHN

Analysis via Constant-Volume Flows

Mathematical characterization:

$$0 < \underbrace{\mathbb{E}_q[\log p(\mathbf{x}; \theta)] - \mathbb{E}_{p^*}[\log p(\mathbf{x}; \theta)]}_{\text{Non-Training Distribution}} + \underbrace{\frac{1}{2} \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p_z(f(\mathbf{x}_0; \phi)) + \nabla_{\mathbf{x}_0}^2 \log \left| \frac{\partial f_\phi}{\partial \mathbf{x}_0} \right| \right] (\Sigma_q - \Sigma_{p^*}) \right\}}_{\begin{array}{l} \text{Training Distribution} \\ \text{Second Moment of Training Distribution} \\ \text{Second Moment of Non-Training Distribution} \end{array}}$$

The diagram illustrates the mathematical characterization of the analysis via constant-volume flows. It shows the difference between the expected log probability under the non-training distribution (\mathbb{E}_q) and the training distribution (\mathbb{E}_{p^*}). This difference is decomposed into two main components: "Change-of-Variable Terms" (indicated by a grey double-headed arrow) and the "Second Moment of Training Distribution" minus the "Second Moment of Non-Training Distribution". A red cross is drawn over the term involving the Jacobian determinant to emphasize its removal or inaccuracy.

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \boldsymbol{\theta}) \right] (\boldsymbol{\Sigma}_q - \boldsymbol{\Sigma}_{p^*}) \right\}$$

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\begin{aligned} & \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \boldsymbol{\theta}) \right] (\boldsymbol{\Sigma}_q - \boldsymbol{\Sigma}_{p^*}) \right\} \\ &= \frac{\partial^2}{\partial z^2} \log p(z; \boldsymbol{\psi}) \sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2 \sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2) \end{aligned}$$

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \boldsymbol{\theta}) \right] (\boldsymbol{\Sigma}_q - \boldsymbol{\Sigma}_{p^*}) \right\}$$

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

**Base
Distribution**

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\begin{aligned} & \text{Tr} \left\{ \left[\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \boldsymbol{\theta}) \right] (\boldsymbol{\Sigma}_q - \boldsymbol{\Sigma}_{p^*}) \right\} \\ &= \frac{\partial^2}{\partial z^2} \log p(\mathbf{z}; \boldsymbol{\psi}) \sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2 \sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2) \end{aligned}$$

Base Distribution **1x1 Conv. Params**

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\text{Tr} \left\{ [\nabla_{\mathbf{x}_0}^2 \log p(\mathbf{x}_0; \theta)] (\Sigma_q - \Sigma_{p^*}) \right\}$$

$$= \frac{\partial^2}{\partial z^2} \log p(z; \psi) \sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2 \sum_{h,w} (\underbrace{\sigma_{q,h,w,c}^2}_{\text{Second Moment of Non-Training Distribution}} - \underbrace{\sigma_{p^*,h,w,c}^2}_{\text{Second Moment of Training Distribution}})$$

Base Distribution **1x1 Conv. Params**

Analysis via Constant-Volume Flows

Plugging in the CV-Glow transform:

$$\text{Tr} \left\{ [\nabla_{x_0}^2 \log p(x_0; \theta)] (\Sigma_q - \Sigma_{p^*}) \right\}$$

$$= \frac{\partial^2}{\partial z^2} \log p(z; \psi)$$

Base
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \sum_{j=1}^C u_{k,c,j} \right)^2$$

1x1 Conv.
Params

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Second
Moment of
Training
Distribution

Second Moment of
Non-Training
Distribution

Negative for all
log-concave
densities
(e.g. Gaussian)

Non-negative
due to square

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

Training
Distribution

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

Training
Distribution

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

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training
Distribution

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Second Moment of
Non-Training
Distribution

Analysis via Constant-Volume Flows

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

**Non-Training
Distribution**

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

**Training
Distribution**

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

**Second
Moment of
Training
Distribution**

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

**Second Moment of
Non-Training
Distribution**

CIFAR-10 vs SVHN

Asymmetry

Uniform Inputs

Ensembling

Early Stopping

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \underline{\sigma_{p^*,h,w,c}^2})$$

Second Moment of
Non-Training
Distribution

- CIFAR-10 vs SVHN (plugging in empirical moments)

- Asymmetry

- Uniform Inputs

- Ensembling

- Early Stopping

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training
Distribution

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Second Moment of
Non-Training
Distribution

CIFAR-10 vs SVHN (plugging in empirical moments)

Asymmetry (due to sub. being non-commutative)

Uniform Inputs

Ensembling

Early Stopping

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training
Distribution

$$\sum_{h,w} (\cancel{\sigma_{q,n,w,c}^2} - \sigma_{p^*,h,w,c}^2)$$

Second Moment of
Non-Training
Distribution

CIFAR-10 vs SVHN (plugging in empirical moments)

Asymmetry (due to sub. being non-commutative)

Uniform Inputs

Ensembling

Early Stopping

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training

$$\sum_{h,w} (\sigma_{q,n,w,c}^2 - \sigma_{p^*,n,w,c}^2)$$

Second Moment of
Non-Training
Distribution

CIFAR-10 vs SVHN (plugging in empirical moments)

Asymmetry (due to sub. being non-commutative)

Uniform Inputs (non-training 2nd moment is zero)

Ensembling

Early Stopping

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training
Distribution

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Second Moment of
Non-Training
Distribution

CIFAR-10 vs SVHN (plugging in empirical moments)

Asymmetry (due to sub. being non-commutative)

Uniform Inputs (non-training 2nd moment is zero)

Ensembling }
 Early Stopping } (sign doesn't depend on model param. values)

Analysis via Constant-Volume Flows

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

Non-Training
Distribution

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

Training
Distribution

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

Second
Moment of
Training
Distribution

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

Second Moment of
Non-Training
Distribution

Hypothesis: If the second-order statistics do indeed dominate, we should be able to control the likelihoods by **graying** the images...

Analysis via Constant-Volume Flows

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

**Non-Training
Distribution**

$$\approx \frac{\partial^2}{\partial z^2} \left(\underline{\log p(z; \psi)} \right)$$

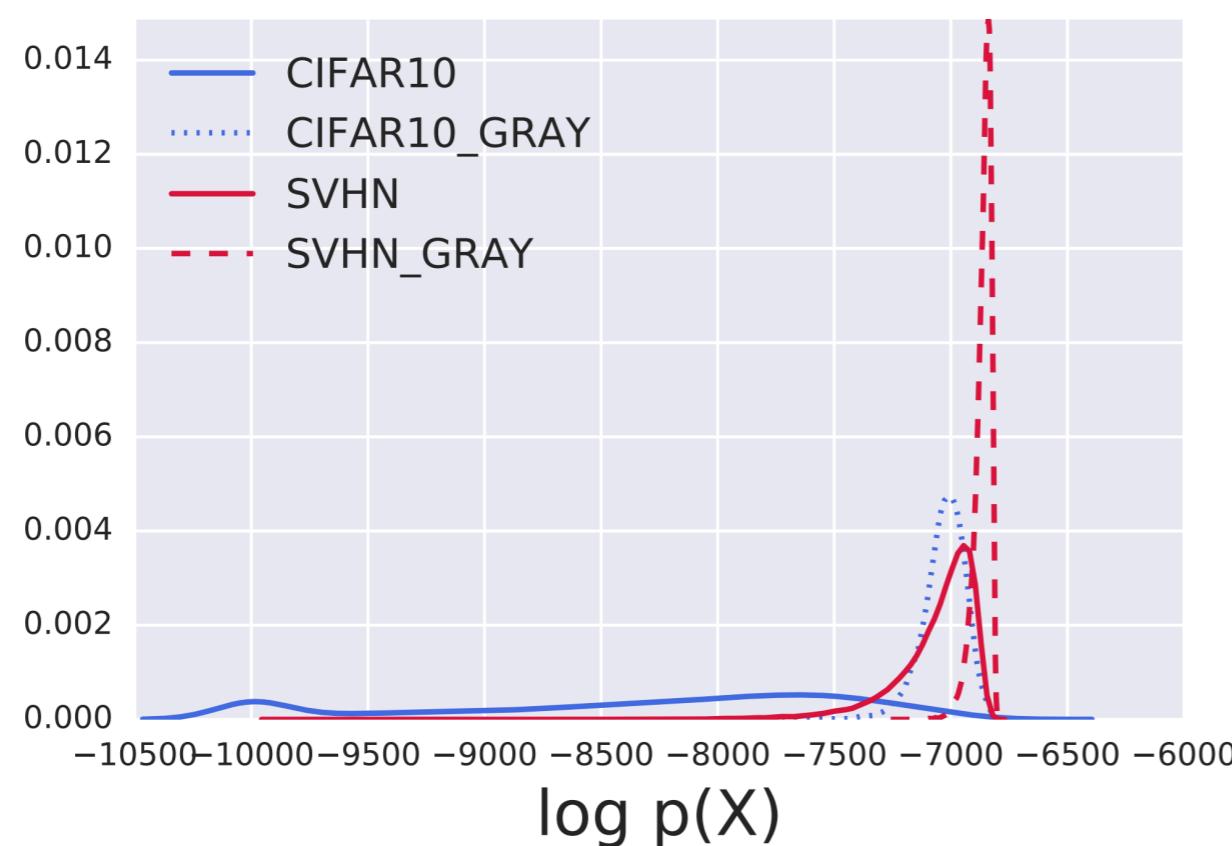
**Training
Distribution**

$$\sum_{c=1}^C \left(\prod_{k=1}^K \left(\sum_{j=1}^C u_{k,c,j} \right) \right)$$

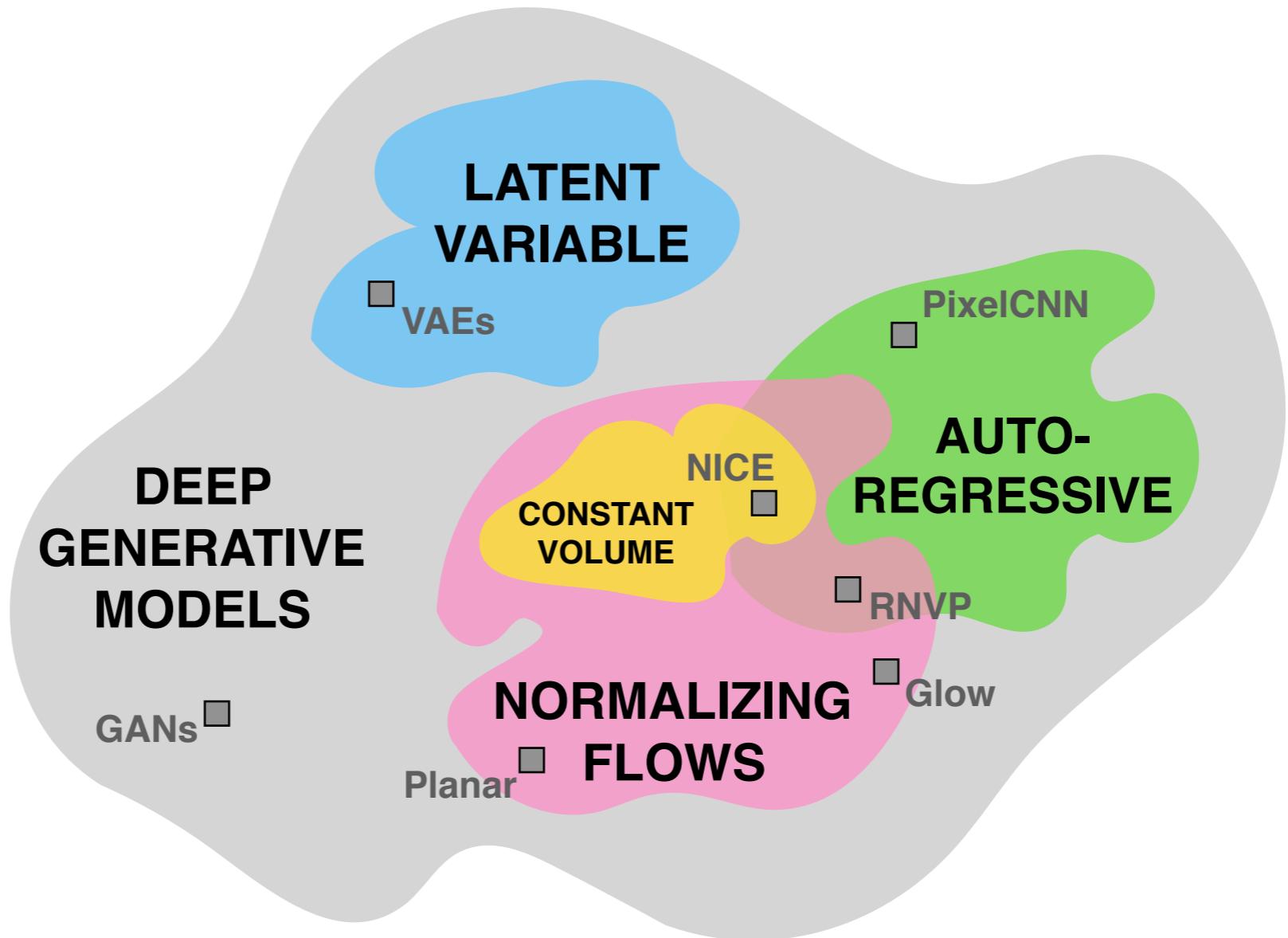
**Second
Moment of
Training
Distribution**

$$\sum_{h,w} (\sigma_{q,h,w,c}^2 - \sigma_{p^*,h,w,c}^2)$$

**Second Moment of
Non-Training
Distribution**

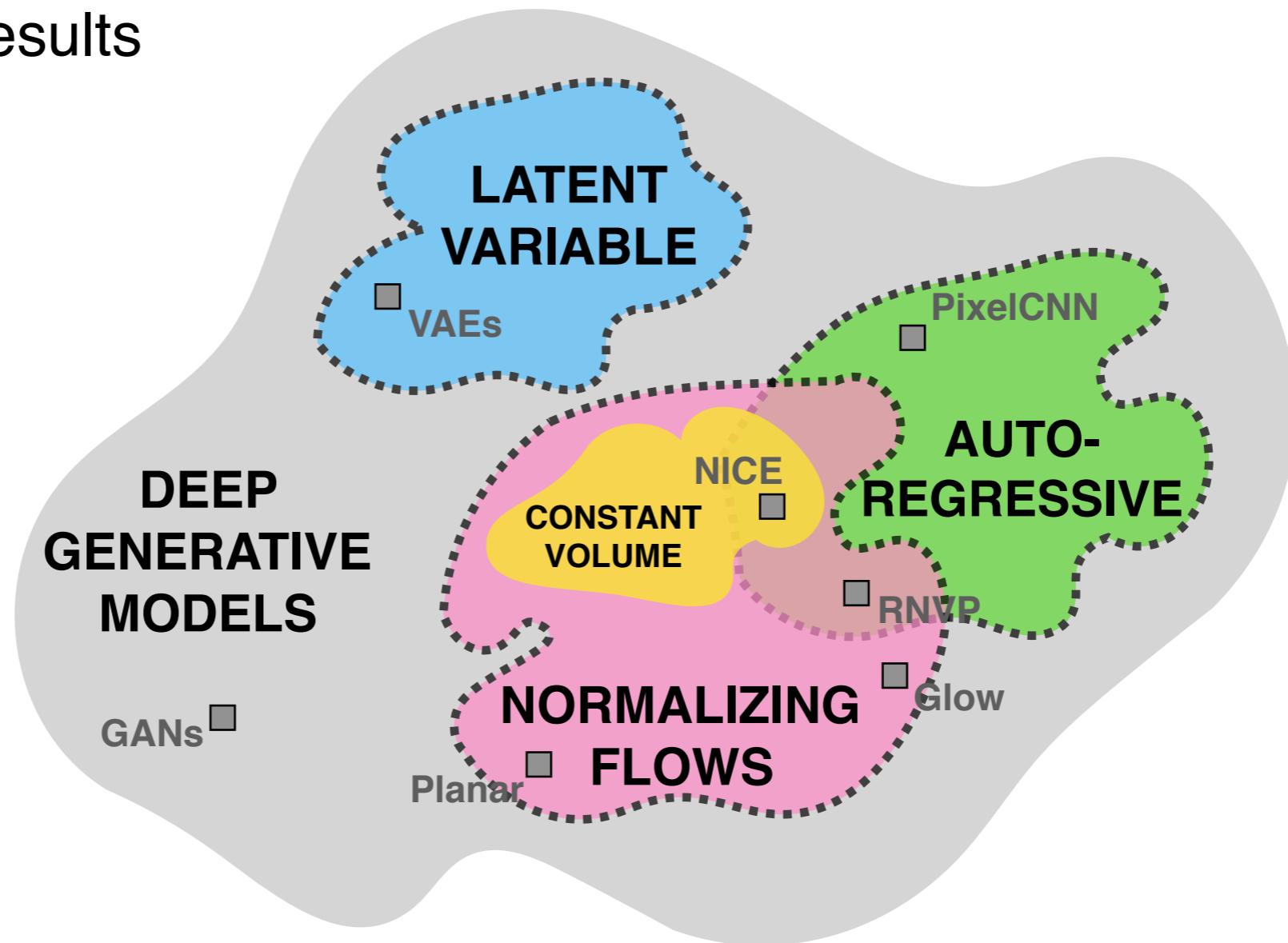


Summary of Results



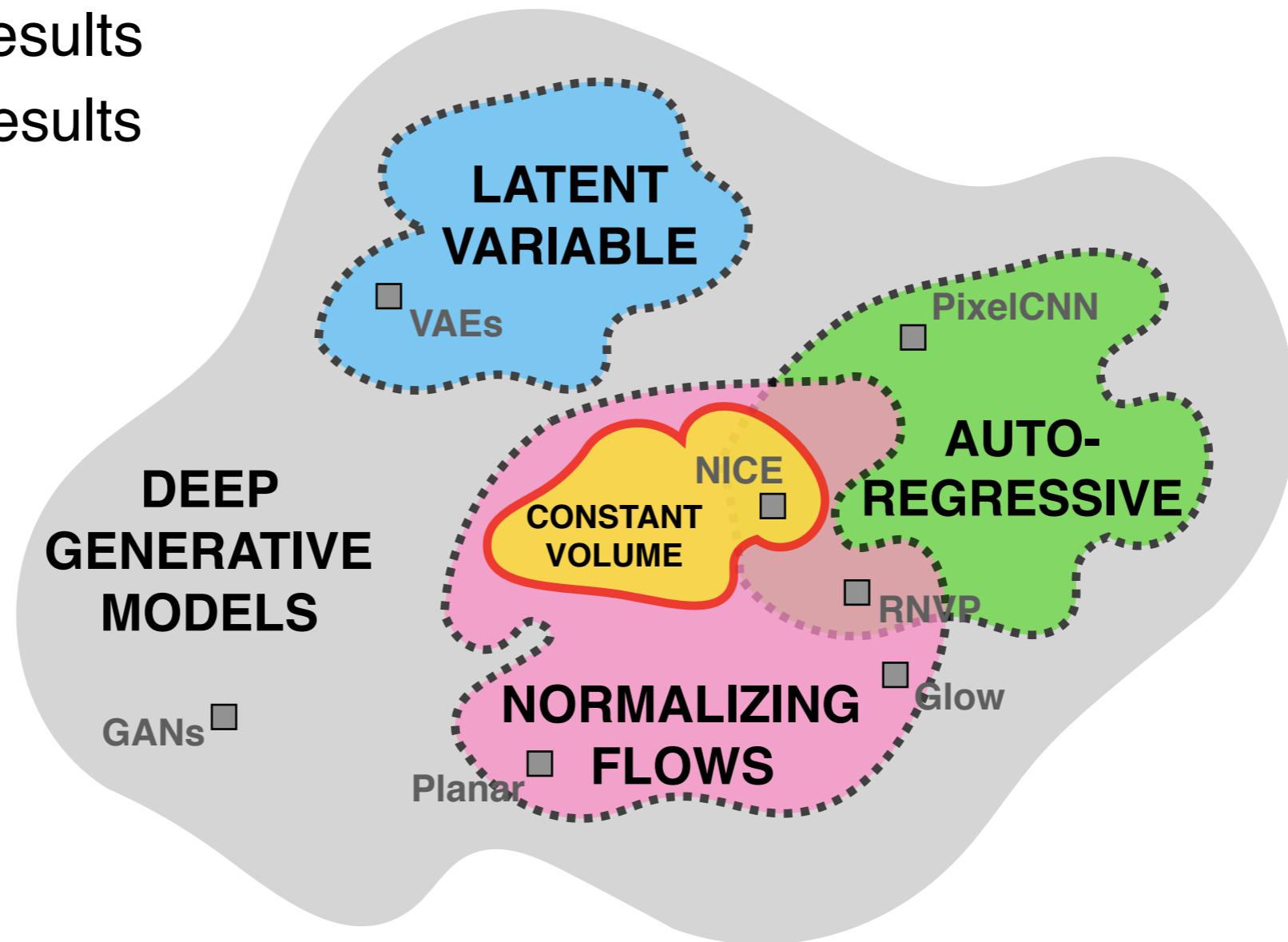
Summary of Results

---- Empirical Results



Summary of Results

- Empirical Results
- Analytical Results



PART #4

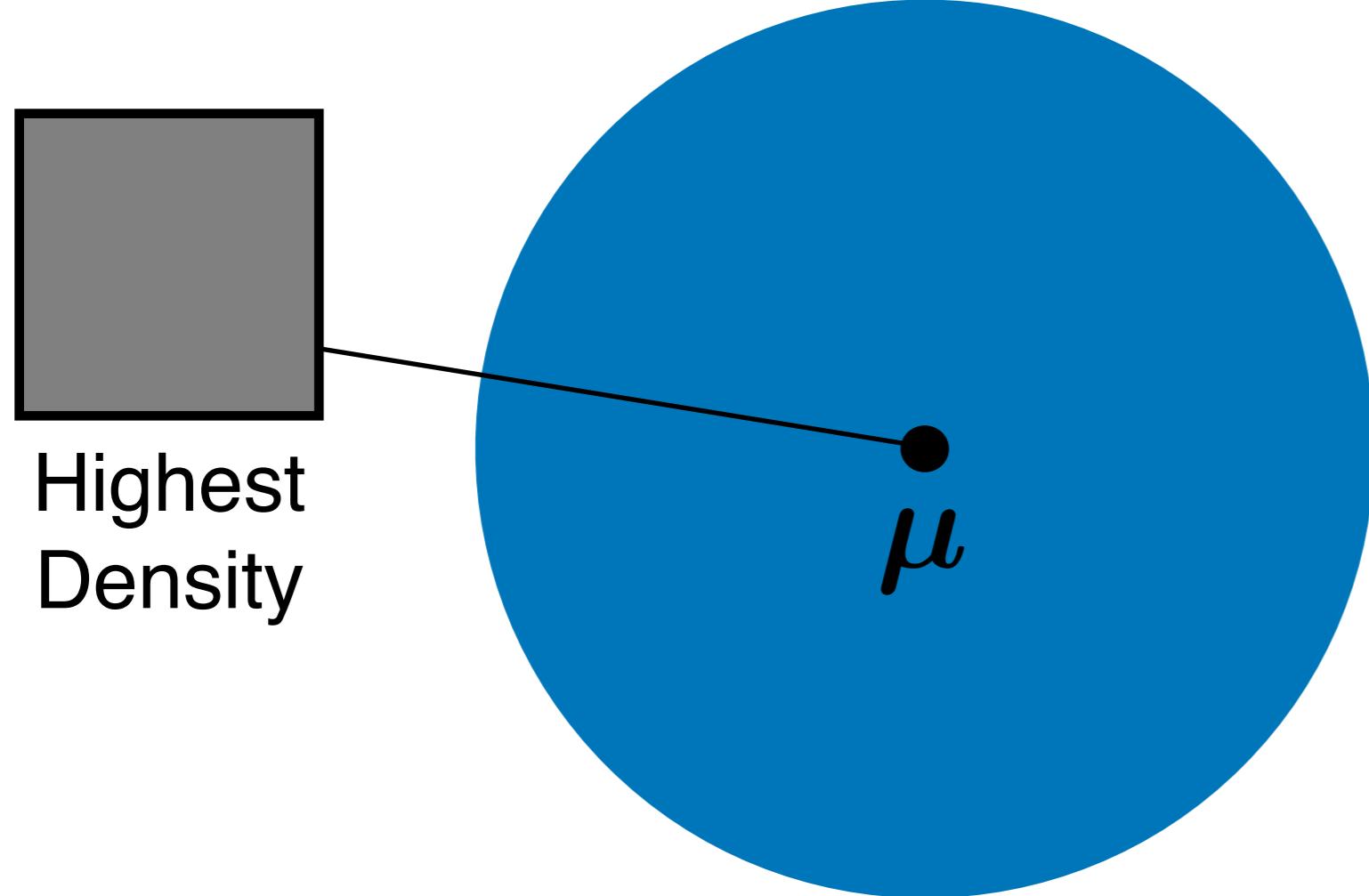
Concentration and Typicality

Question: In the CIFAR vs SVHN case, why don't we ever see samples from SVHN?

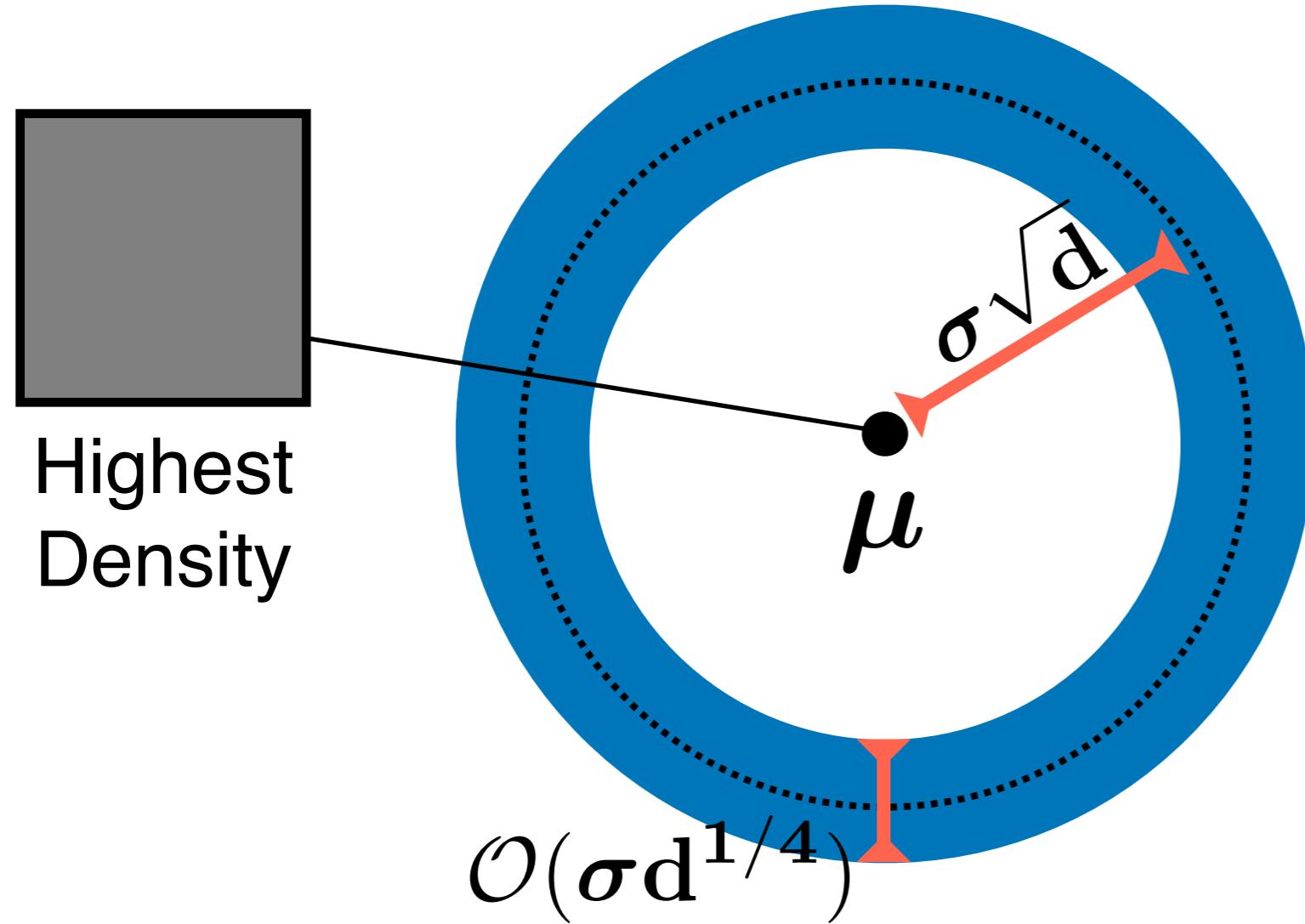


**Samples from Glow when
trained on CIFAR-10**

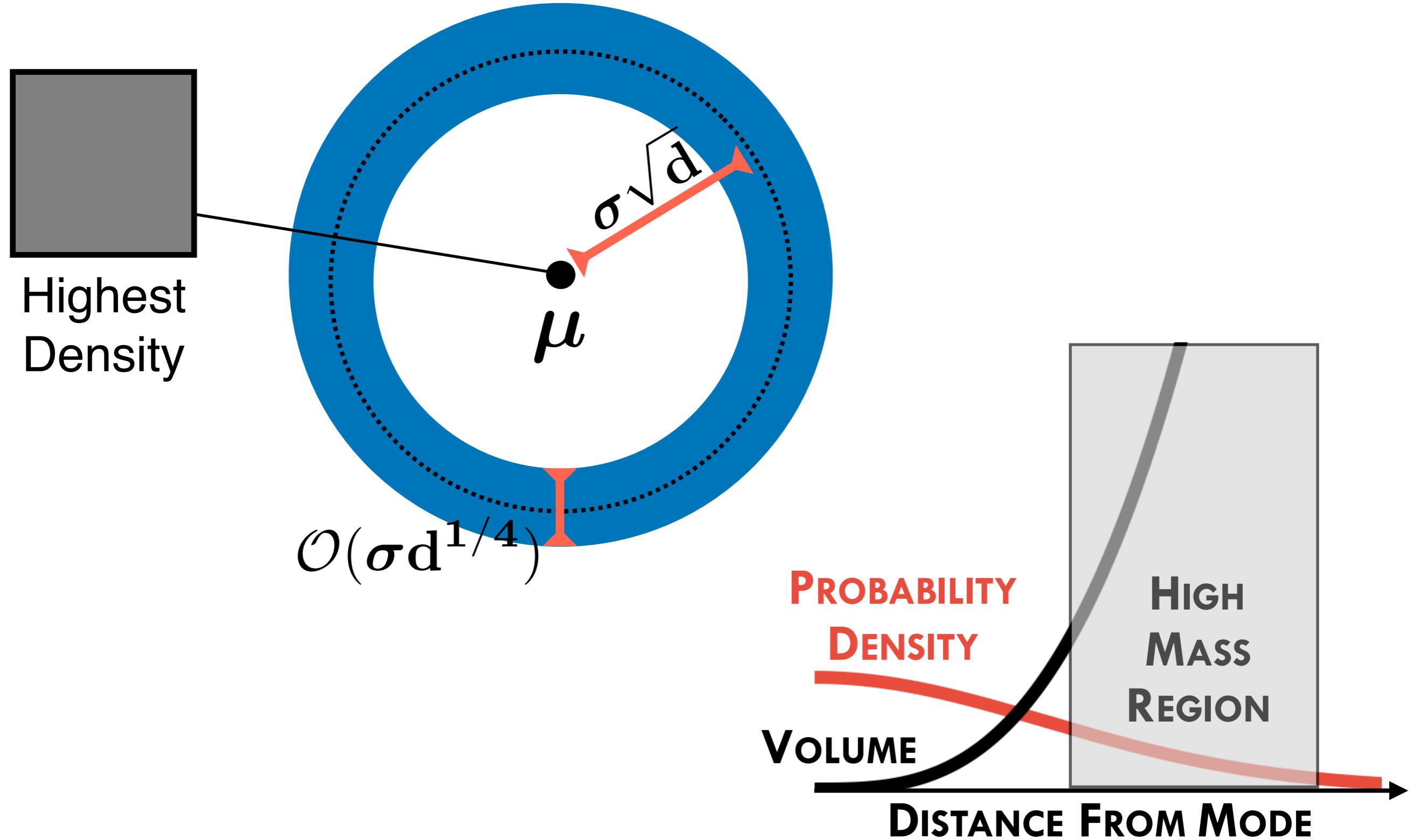
Back to basics: consider a high-dimensional Gaussian centered on the all-gray image...



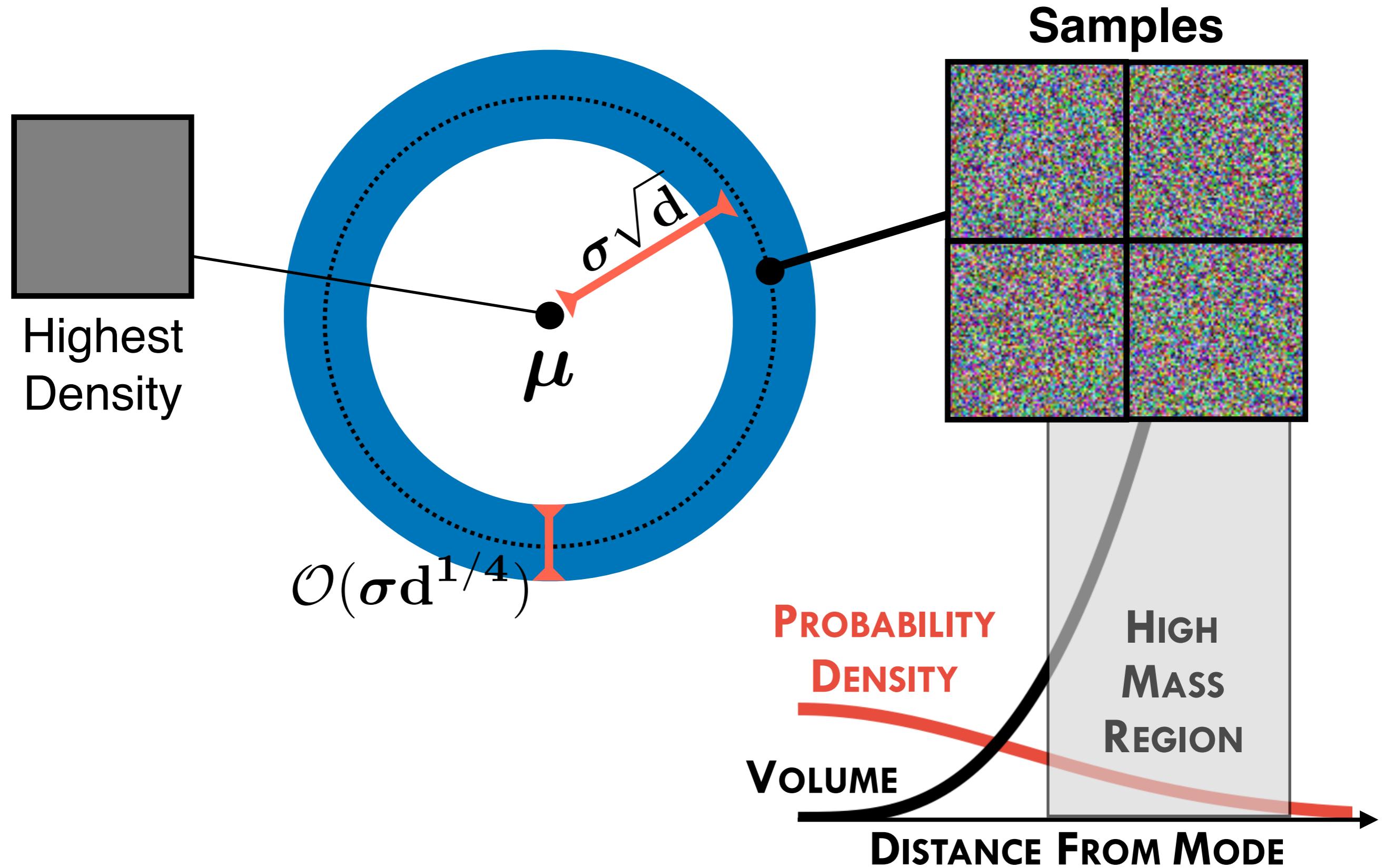
Back to basics: consider a high-dimensional Gaussian centered on the all-gray image...



Back to basics: consider a high-dimensional Gaussian centered on the all-gray image...



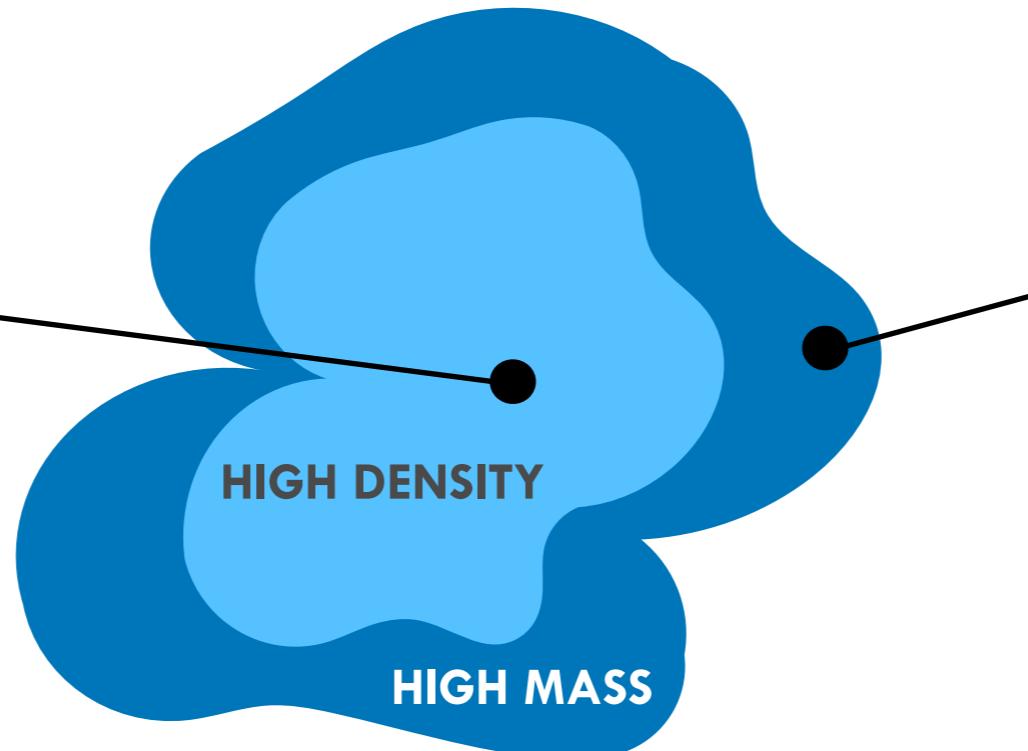
Back to basics: consider a high-dimensional Gaussian centered on the all-gray image...



Conjecture: We are probably just seeing the same concentration phenomenon in deep generative models...



High Density



High Probability
(Samples)

FashionMNIST train vs **MNIST** test

Thus we should be checking for out-of-distribution inputs **not via density but via mass...**

Thus we should be checking for out-of-distribution inputs **not via density but via mass...**

But how do we do this?...

Thus we should be checking for out-of-distribution inputs **not via density but via mass...**

But how do we do this?...

Idea #1: Integrate the density function.

No surprise that the canonical goodness-of-fit tests use the model CDF (e.g. Kolmogorov–Smirnov, Anderson–Darling, etc.).

Thus we should be checking for out-of-distribution inputs **not via density but via mass...**

But how do we do this?...

Idea #1: Integrate the density function.

No surprise that the canonical out-of-distribution or-fit tests use the model CDF (e.g. Kolmogorov–Smirnov, Anderson–Darling, etc.).

Thus we should be checking for out-of-distribution inputs **not via density but via mass...**

But how do we do this?...

Idea #1: Integrate the density function.

No surprise that the canonical out-of-distribution or-fit tests use the model CDF (e.g. Kolmogorov–Smirnov, Anderson–Darling, etc.).

Idea #2: Check for *typicality*.

Typical Sets

Definition 2.1. ϵ -Typical Set [11] *For a distribution $p(\mathbf{x})$ with support $\mathbf{x} \in \mathcal{X}$, the ϵ -typical set $\mathcal{A}_\epsilon^N[p(\mathbf{x})] \in \mathcal{X}^N$ is comprised of all N -length sequences that satisfy*

$$\mathbb{H}[p(\mathbf{x})] - \epsilon \leq \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n) \leq \mathbb{H}[p(\mathbf{x})] + \epsilon$$

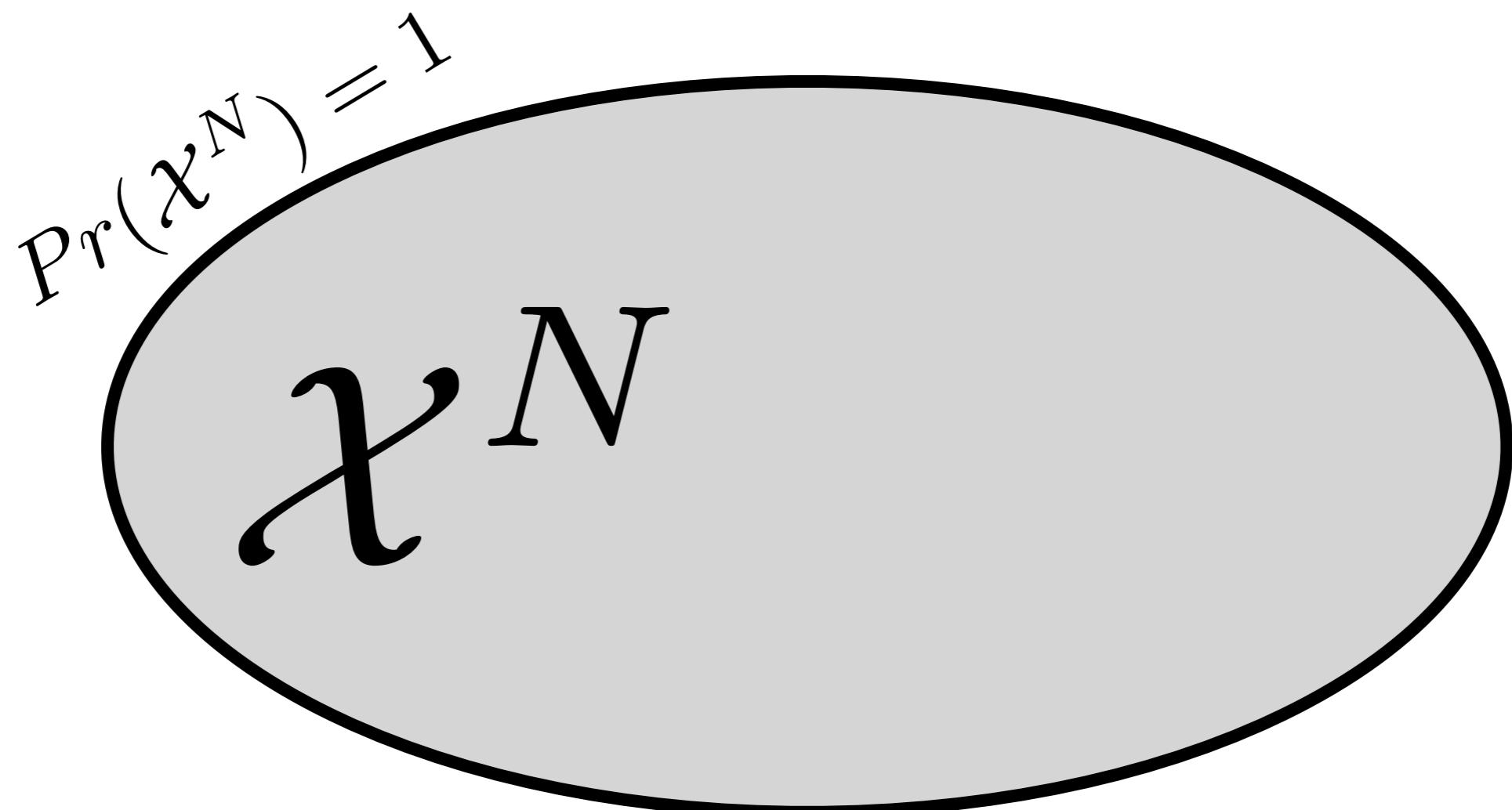
where $\mathbb{H}[p(\mathbf{x})] = \int_{\mathcal{X}} p(\mathbf{x})[-\log p(\mathbf{x})]d\mathbf{x}$ and $\epsilon \in \mathbb{R}^+$ is a small constant.

Typical Sets

Definition 2.1. ϵ -Typical Set [11] For a distribution $p(\mathbf{x})$ with support $\mathbf{x} \in \mathcal{X}$, the ϵ -typical set $\mathcal{A}_\epsilon^N[p(\mathbf{x})] \in \mathcal{X}^N$ is comprised of all N -length sequences that satisfy

$$\mathbb{H}[p(\mathbf{x})] - \epsilon \leq \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n) \leq \mathbb{H}[p(\mathbf{x})] + \epsilon$$

where $\mathbb{H}[p(\mathbf{x})] = \int_{\mathcal{X}} p(\mathbf{x})[-\log p(\mathbf{x})]d\mathbf{x}$ and $\epsilon \in \mathbb{R}^+$ is a small constant.

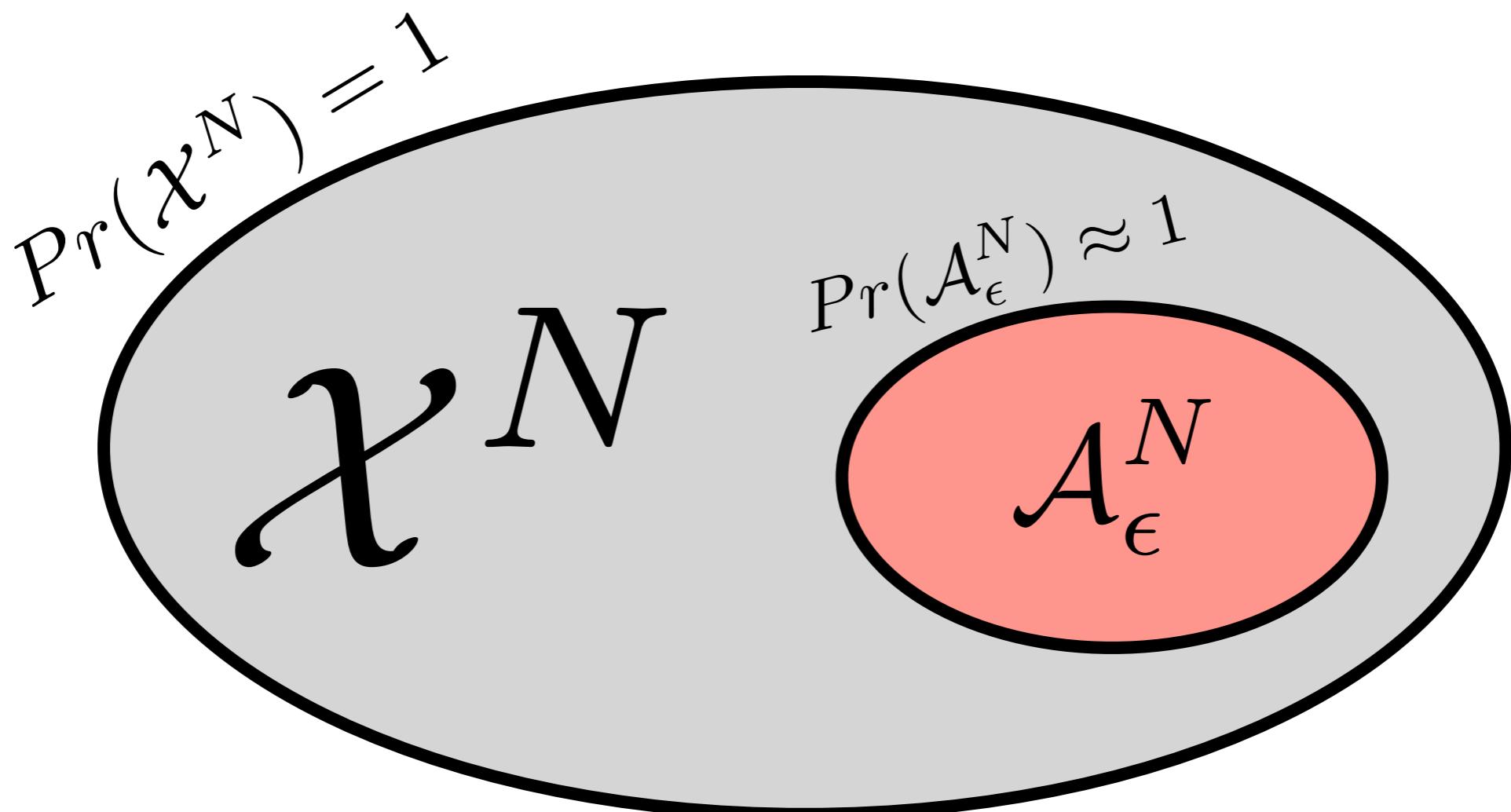


Typical Sets

Definition 2.1. ϵ -Typical Set [11] For a distribution $p(\mathbf{x})$ with support $\mathbf{x} \in \mathcal{X}$, the ϵ -typical set $\mathcal{A}_\epsilon^N[p(\mathbf{x})] \in \mathcal{X}^N$ is comprised of all N -length sequences that satisfy

$$\mathbb{H}[p(\mathbf{x})] - \epsilon \leq \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n) \leq \mathbb{H}[p(\mathbf{x})] + \epsilon$$

where $\mathbb{H}[p(\mathbf{x})] = \int_{\mathcal{X}} p(\mathbf{x})[-\log p(\mathbf{x})]d\mathbf{x}$ and $\epsilon \in \mathbb{R}^+$ is a small constant.



Using Typicality to Assess Model Fit

For an M -sized test batch

$$\widetilde{\mathbf{X}} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_M\}$$

if $\widetilde{\mathbf{X}} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$ then $\widetilde{\mathbf{X}} \sim p(\mathbf{x}; \boldsymbol{\theta})$,

otherwise $\widetilde{\mathbf{X}} \not\sim p(\mathbf{x}; \boldsymbol{\theta})$

Using Typicality to Assess Model Fit

For an M -sized test batch

$$\tilde{\mathbf{X}} = \{\tilde{\mathbf{x}}_1, \dots, \tilde{\mathbf{x}}_M\}$$

if $\tilde{\mathbf{X}} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$ then $\tilde{\mathbf{X}} \sim p(\mathbf{x}; \boldsymbol{\theta})$,

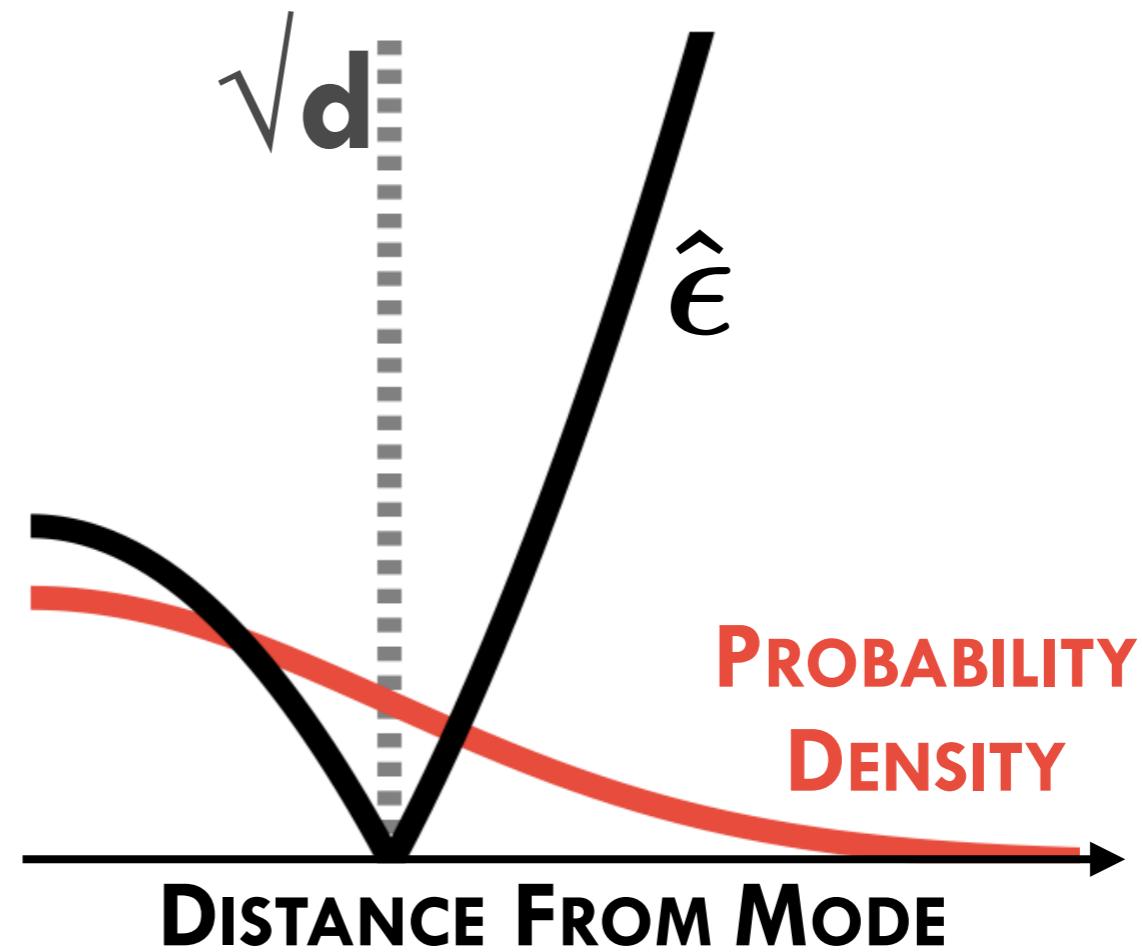
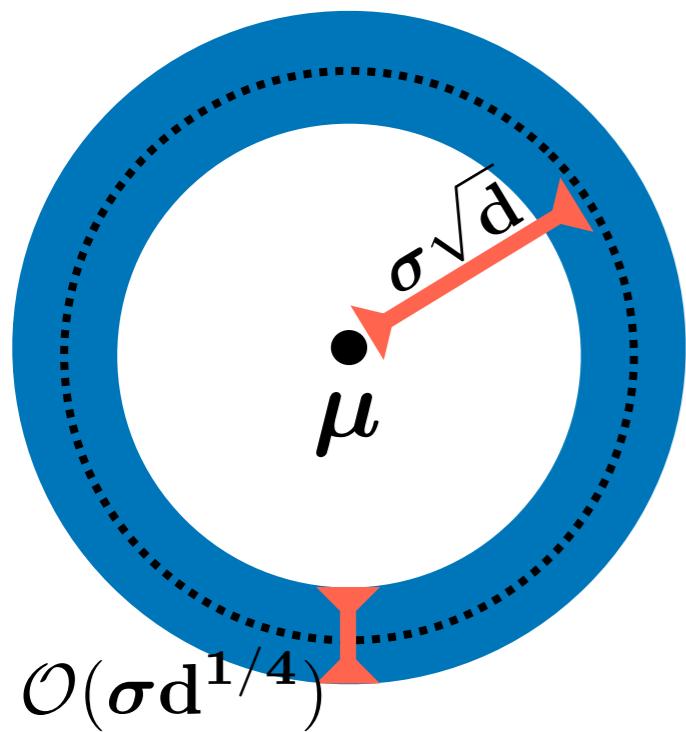
otherwise $\tilde{\mathbf{X}} \not\sim p(\mathbf{x}; \boldsymbol{\theta})$

Plugging in the entropy bound...

if $\left| \frac{-1}{M} \sum_{m=1}^M \log p(\tilde{\mathbf{x}}_m; \boldsymbol{\theta}) - \mathbb{H}[p(\mathbf{x}; \boldsymbol{\theta})] \right| \leq \epsilon$ then $\tilde{\mathbf{X}} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$

Gaussian Simulation

$$\left| \frac{-1}{M} \sum_{m=1}^M \log p(\tilde{\mathbf{x}}_m; \boldsymbol{\theta}) - \mathbb{H}[p(\mathbf{x}; \boldsymbol{\theta})] \right| \leq \epsilon$$



Implementation via Bootstrap CI

How should we set ϵ in practice?

Implementation via Bootstrap CI

How should we set ϵ in practice?

We construct the test statistic via a bootstrap confidence interval and validation data...

Implementation via Bootstrap CI

How should we set ϵ in practice?

We construct the test statistic via a bootstrap confidence interval and validation data...

$$H_0 : \widetilde{X} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$$

Implementation via Bootstrap CI

How should we set ϵ in practice?

We construct the test statistic via a bootstrap confidence interval and validation data...

$$H_0 : \widetilde{\bar{X}} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$$

$$\left| \frac{-1}{M} \sum_{m=1}^M \log p(\mathbf{x}'_{k,m}; \boldsymbol{\theta}) - \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \right| = \hat{\epsilon}_k$$

Implementation via Bootstrap CI

How should we set ϵ in practice?

We construct the test statistic via a bootstrap confidence interval and validation data...

$$H_0 : \widetilde{\bar{X}} \in \mathcal{A}_\epsilon^M[p(\mathbf{x}; \boldsymbol{\theta})]$$

$$\left| \frac{-1}{M} \sum_{m=1}^M \log p(\mathbf{x}'_{k,m}; \boldsymbol{\theta}) - \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta}) \right| = \hat{\epsilon}_k$$

$$F(\epsilon) = \frac{1}{K} \sum_{k=1}^K \delta[\hat{\epsilon}_k]$$

Implementation via Bootstrap CI

Algorithm 1 A Bootstrap Test for Typicality

Input: Training data \mathbf{X} , validation data \mathbf{X}' , trained model $p(\mathbf{x}; \boldsymbol{\theta})$, number of bootstrap samples K , significance level α , M -sized batch of possibly OOD inputs $\tilde{\mathbf{X}}$.

Offline prior to deployment

1. **Compute** $\hat{\mathbb{H}}^N[p(\mathbf{x}; \boldsymbol{\theta})] = \frac{-1}{N} \sum_{n=1}^N \log p(\mathbf{x}_n; \boldsymbol{\theta})$.
2. **Sample** K M -sized data sets from \mathbf{X}' using bootstrap resampling.
3. **For all** $k \in [1, K]$:
 Compute $\hat{\epsilon}_k = \left| \frac{-1}{M} \sum_{m=1}^M \log p(\mathbf{x}'_{k,m}; \boldsymbol{\theta}) - \hat{\mathbb{H}}^N[p(\mathbf{x}; \boldsymbol{\theta})] \right|$ *(Equation 7)*
4. **Set** $\epsilon_\alpha^M = \text{quantile}(F(\epsilon), \alpha)$ *(e.g. $\alpha = .99$)*

Online during deployment

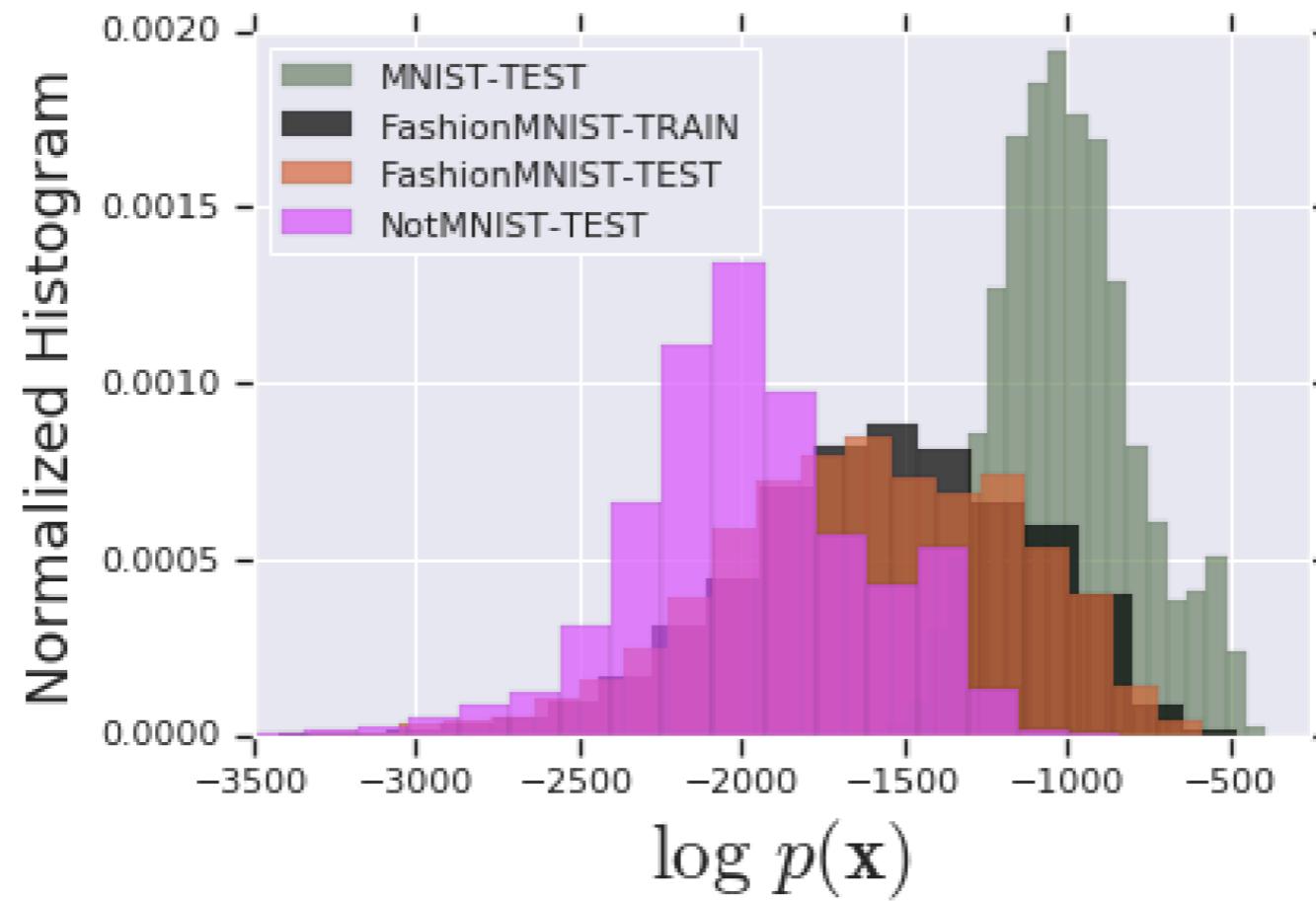
If $\left| \frac{-1}{M} \sum_{m=1}^M \log p(\tilde{\mathbf{x}}_m) - \hat{\mathbb{H}}^N[p(\mathbf{x}; \boldsymbol{\theta})] \right| > \epsilon_\alpha^M$:

Return $\tilde{\mathbf{X}}$ is out-of-distribution

Else:

Return $\tilde{\mathbf{X}}$ is in-distribution

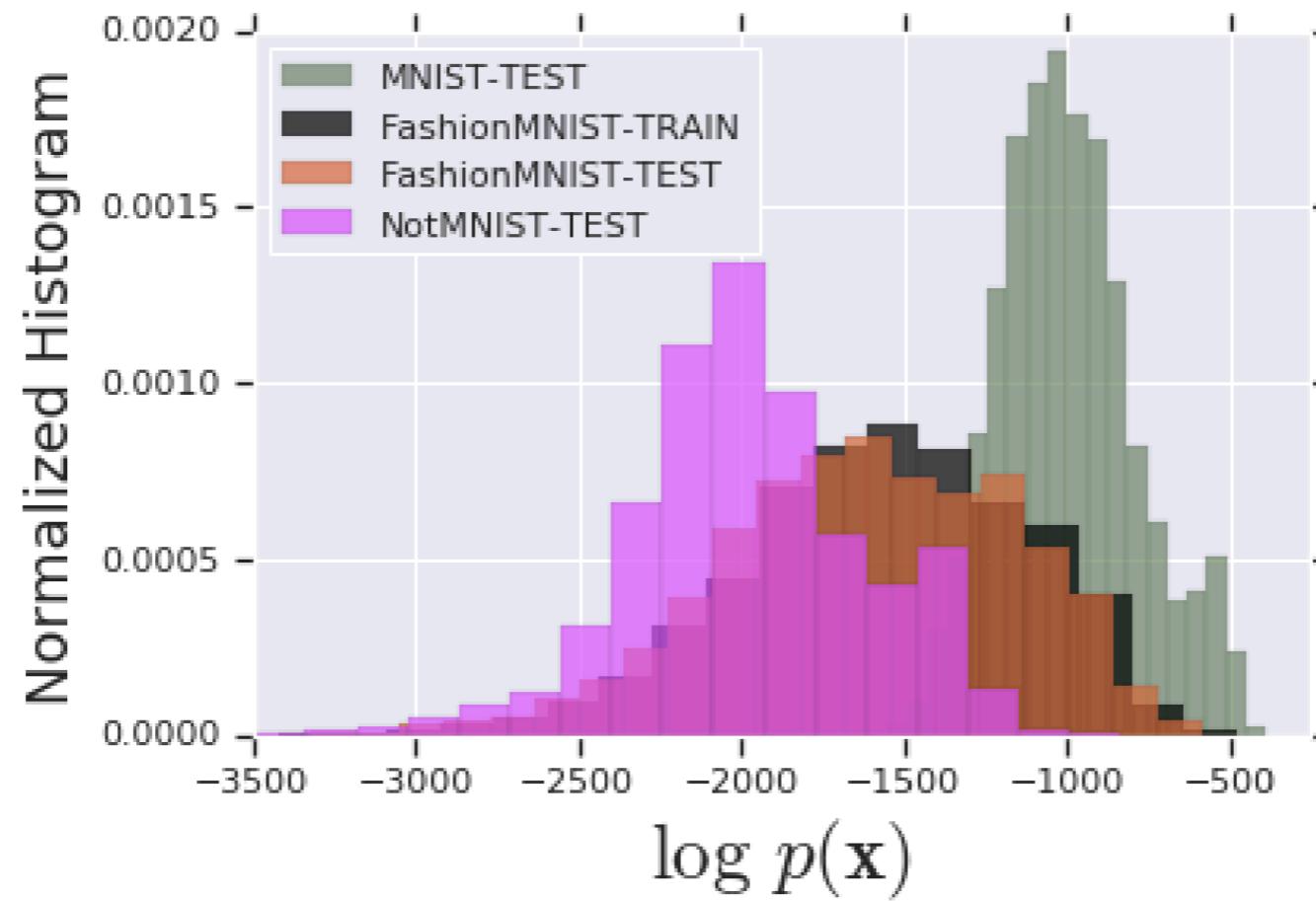
Results: Glow Trained on FashionMNIST



METHOD	$M = 2$		$M = 10$		$M = 25$			
	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST
<i>Glow Trained on FashionMNIST</i>								
Typicality Test								
KS-Test								
Max Mean Dis.								
Kern. Stein Dis.								

Fraction of M-sized batches classified as OOD ($\alpha = 99\%$)

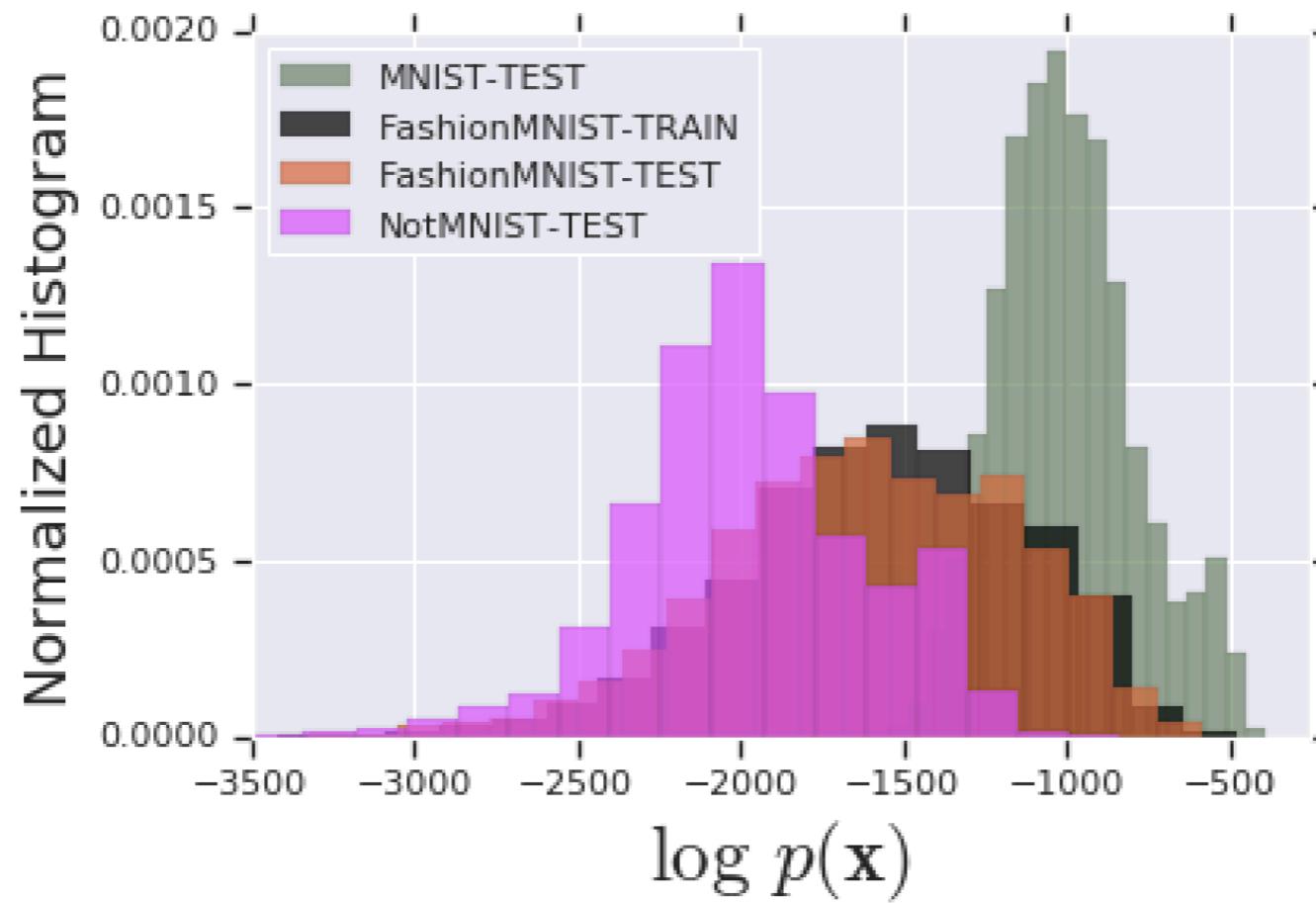
Results: Glow Trained on FashionMNIST



METHOD	$M = 2$		$M = 10$		$M = 25$				
	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST	NOTMNIST
<i>Glow Trained on FashionMNIST</i>									
Typicality Test							$0.01 \pm .00$	$1.00 \pm .00$	$1.00 \pm .00$
KS-Test							$0.00 \pm .00$	$1.00 \pm .00$	$0.98 \pm .01$
Max Mean Dis.							$0.04 \pm .04$	$1.00 \pm .00$	$1.00 \pm .00$
Kern. Stein Dis.							$0.02 \pm .03$	$0.76 \pm .21$	$0.00 \pm .00$

Fraction of M-sized batches classified as OOD ($\alpha = 99\%$)

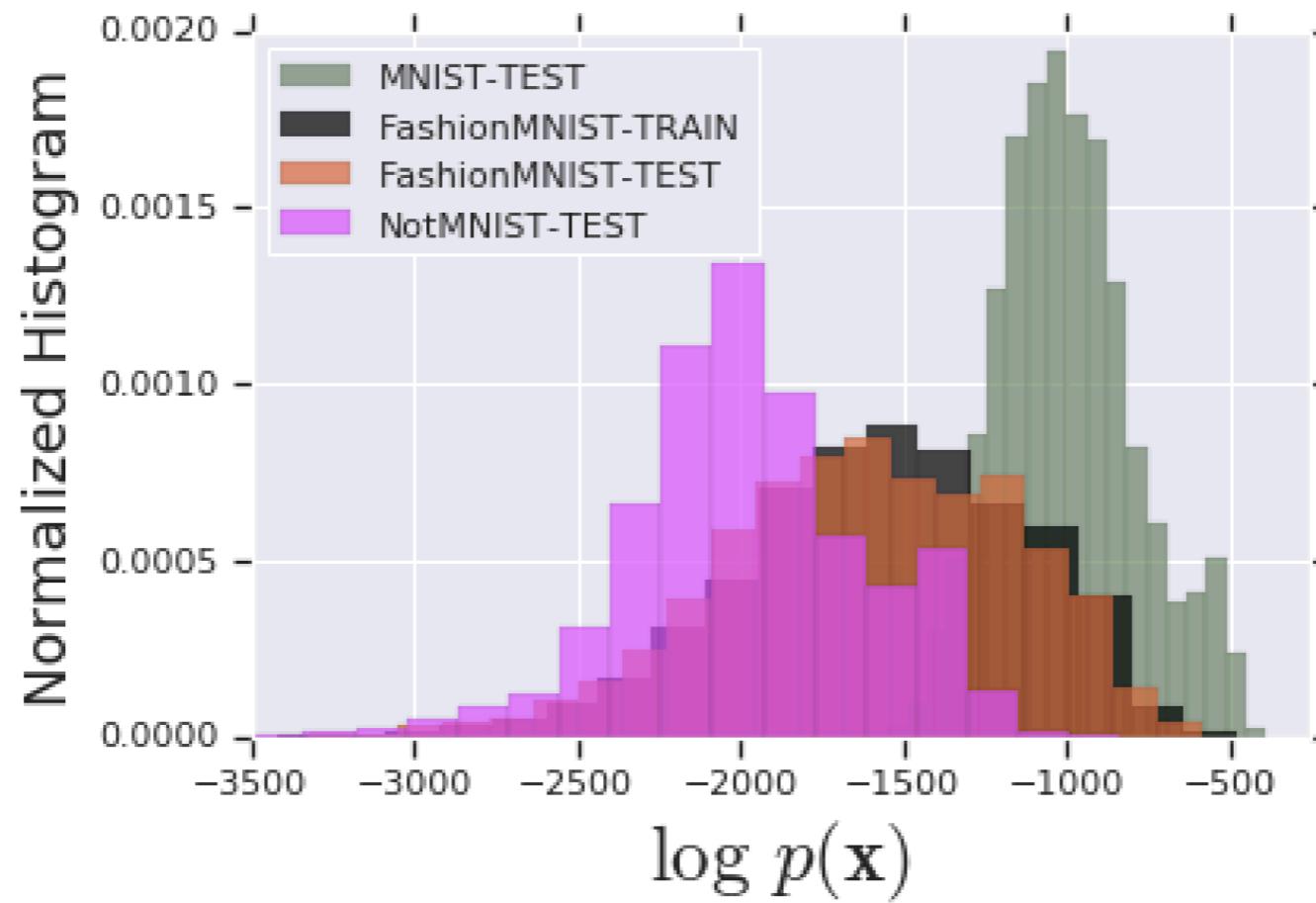
Results: Glow Trained on FashionMNIST



METHOD	IN-DIST.	$M = 2$		$M = 10$		$M = 25$		
		MNIST	NOTMNIST	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST
<i>Glow Trained on FashionMNIST</i>								
Typicality Test				$0.02 \pm .02$	$1.00 \pm .00$	$0.69 \pm .11$	$0.01 \pm .00$	$1.00 \pm .00$
KS-Test				$0.01 \pm .00$	$1.00 \pm .00$	$0.61 \pm .01$	$0.00 \pm .00$	$1.00 \pm .00$
Max Mean Dis.				$0.02 \pm .02$	$0.63 \pm .12$	$0.37 \pm .24$	$0.04 \pm .04$	$1.00 \pm .00$
Kern. Stein Dis.				$0.01 \pm .01$	$0.21 \pm .11$	$0.01 \pm .00$	$0.02 \pm .03$	$0.76 \pm .21$
								$0.00 \pm .00$

Fraction of M-sized batches classified as OOD ($\alpha = 99\%$)

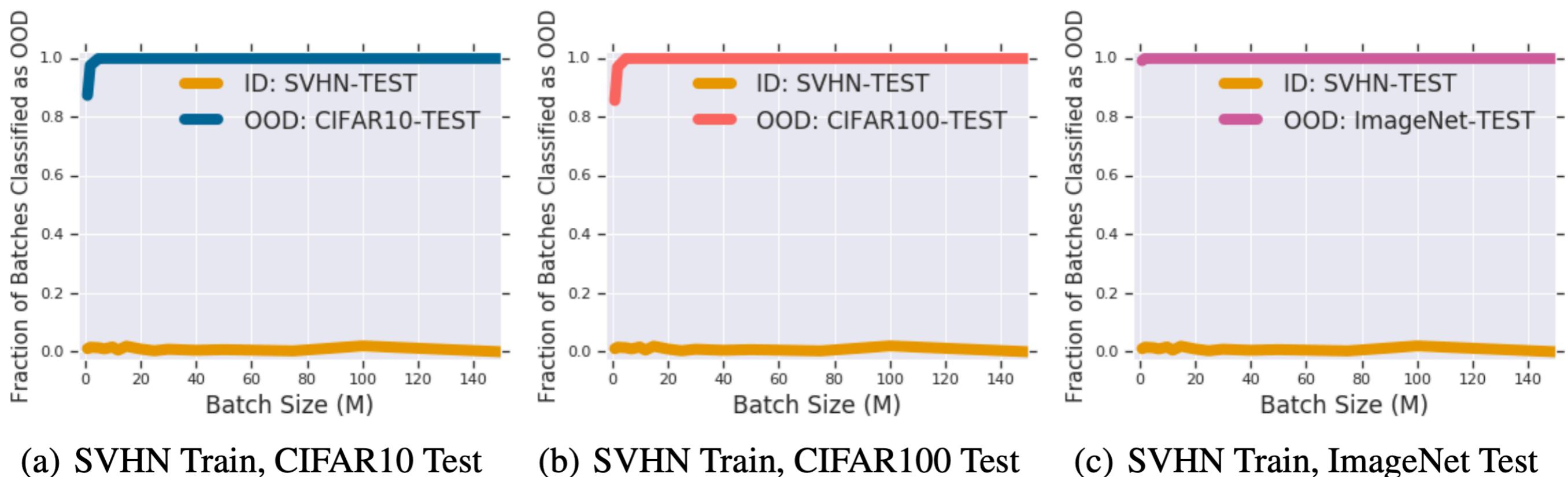
Results: Glow Trained on FashionMNIST



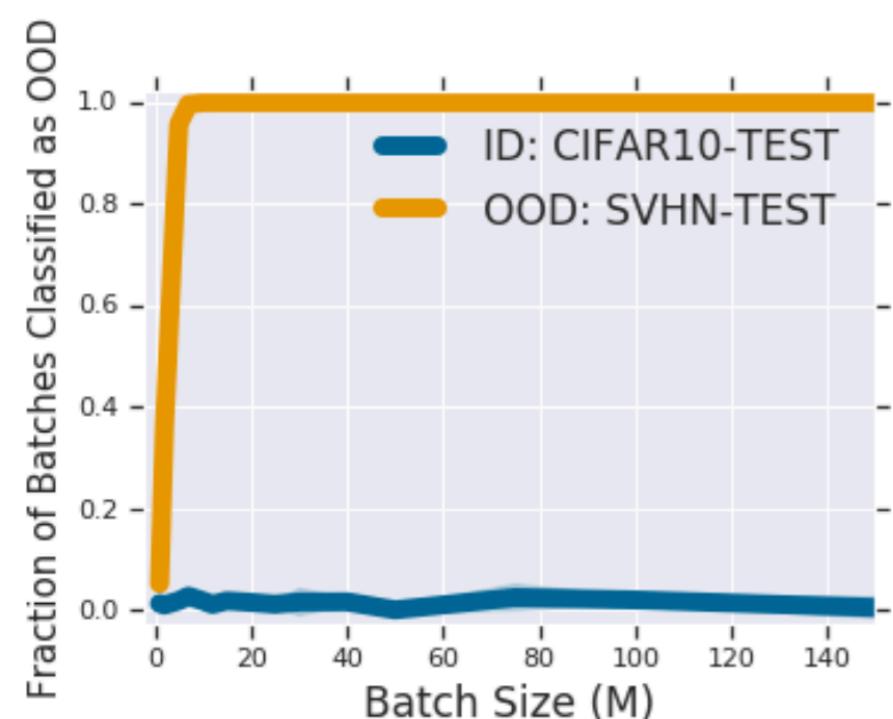
METHOD	IN-DIST.	$M = 2$		$M = 10$		$M = 25$		
		MNIST	NOTMNIST	IN-DIST.	MNIST	NOTMNIST	IN-DIST.	MNIST
<i>Glow Trained on FashionMNIST</i>								
Typicality Test	$0.02 \pm .01$	$0.14 \pm .10$	$0.08 \pm .04$	$0.02 \pm .02$	$1.00 \pm .00$	$0.69 \pm .11$	$0.01 \pm .00$	$1.00 \pm .00$
KS-Test	$0.00 \pm .00$	$0.00 \pm .00$	$0.00 \pm .00$	$0.01 \pm .00$	$1.00 \pm .00$	$0.61 \pm .01$	$0.00 \pm .00$	$1.00 \pm .00$
Max Mean Dis.	$0.05 \pm .02$	$0.17 \pm .06$	$0.04 \pm .03$	$0.02 \pm .02$	$0.63 \pm .12$	$0.37 \pm .24$	$0.04 \pm .04$	$1.00 \pm .00$
Kern. Stein Dis.	$0.05 \pm .05$	$0.16 \pm .14$	$0.01 \pm .01$	$0.01 \pm .01$	$0.21 \pm .11$	$0.01 \pm .00$	$0.02 \pm .03$	$0.76 \pm .21$
								$0.00 \pm .00$

Fraction of M-sized batches classified as OOD ($\alpha = 99\%$)

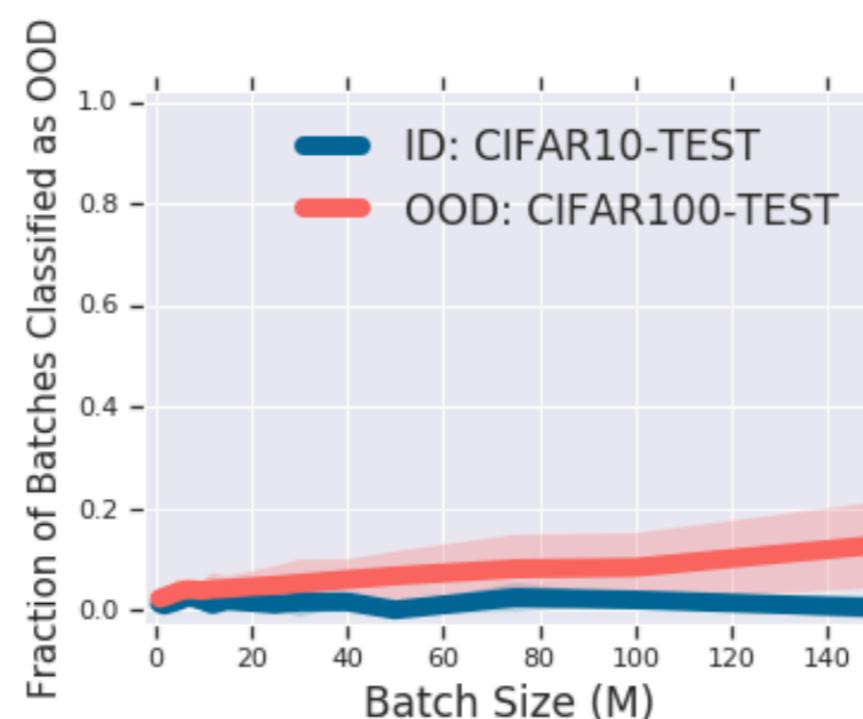
Results: Glow Trained on SVHN



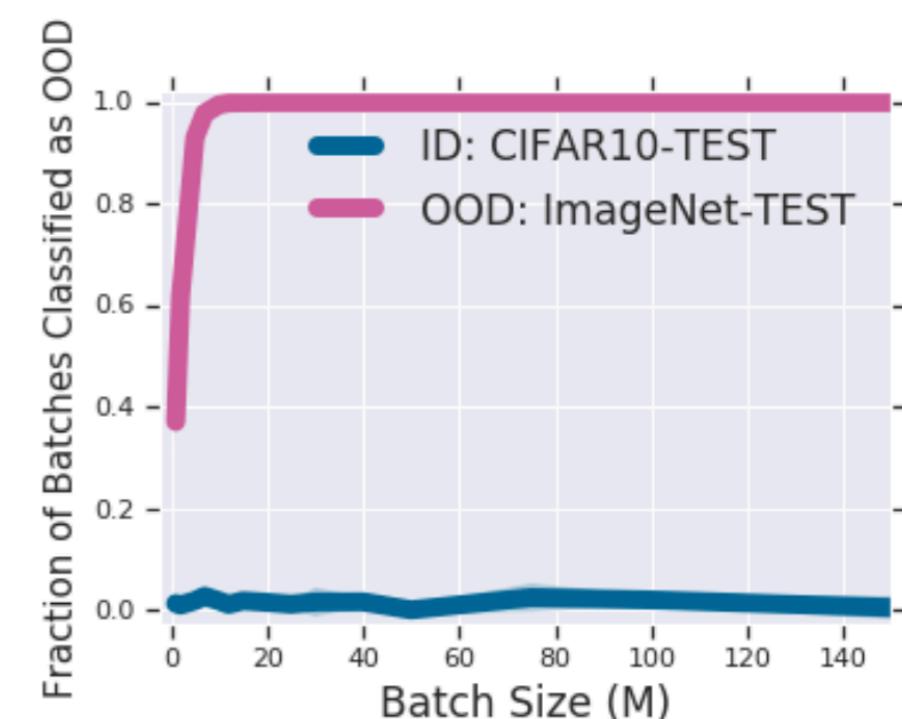
Results: Glow Trained on CIFAR-10



(d) CIFAR10 Train, SVHN Test

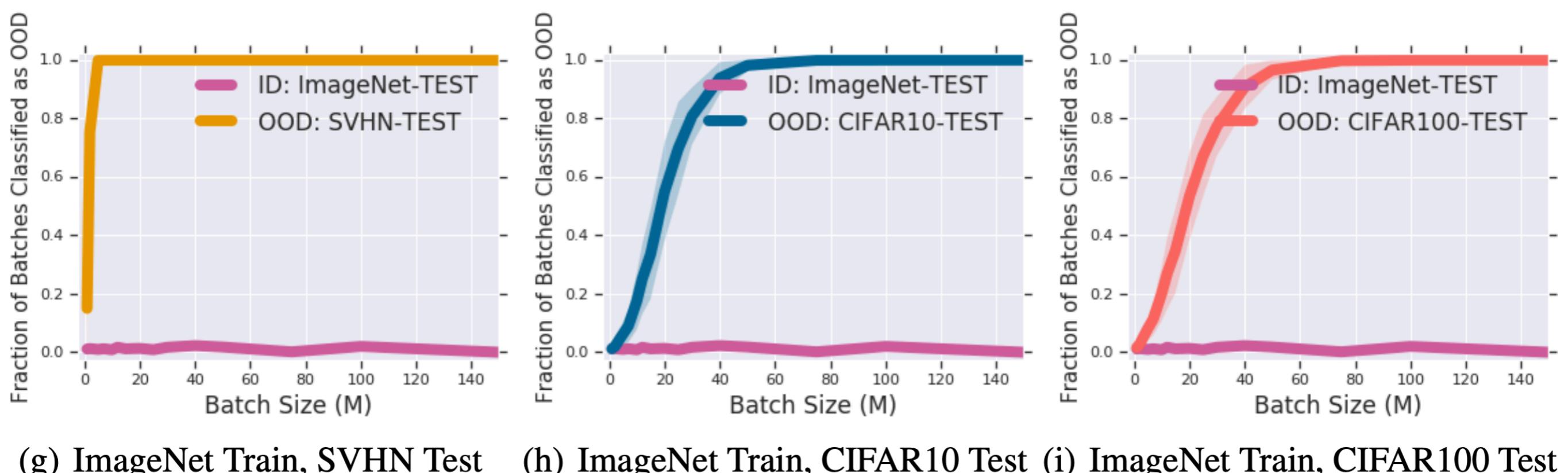


(e) CIFAR10 Train, CIFAR100 Test



(f) CIFAR10 Train, ImageNet Test

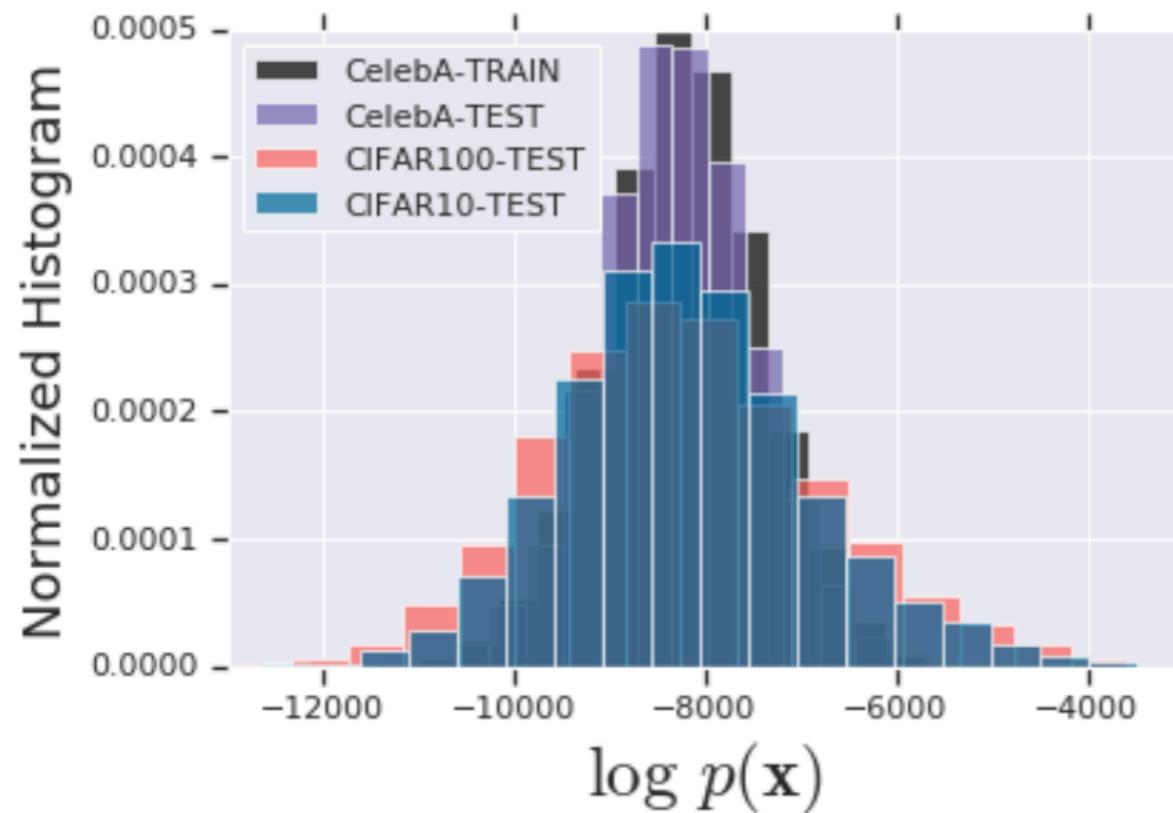
Results: Glow Trained on ImageNet



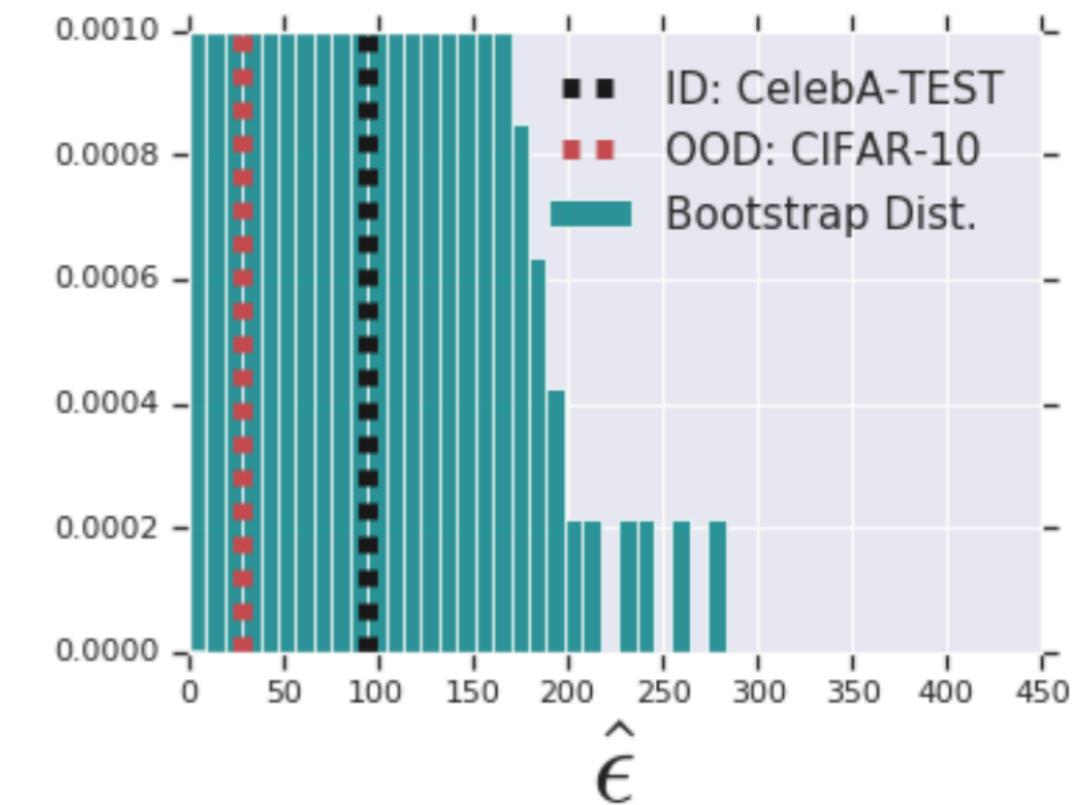
(g) ImageNet Train, SVHN Test

(h) ImageNet Train, CIFAR10 Test (i) ImageNet Train, CIFAR100 Test

Limitations: Only as Good as Likelihood Dist.



(b) Glow: CelebA vs CIFARs



(c) Bootstrap Dist. ($M = 200$)

PART #5

Conclusions

Conclusions

- 1 Failure of Likelihood-Based Anomaly Detection:** Using a likelihood / density threshold is a bad idea for OOD / anomaly detection, especially in high dimensions.
- 2 Connecting Likelihoods to Typical Sets:** Using the empirical dist. of likelihoods can still work surprisingly well, if used properly (via connection to typicality).
- 3 Moving Beyond Generation:** If we really want to leverage these new generative models for principled probabilistic inference, we need to devise goodness-of-fit tests and other critiques that can appropriately scale.

DO DEEP GENERATIVE MODELS KNOW WHAT THEY DON'T KNOW?

Eric Nalisnick^{*†}, Akihiro Matsukawa, Yee Whye Teh, Dilan Gorur, Balaji Lakshminarayanan*
DeepMind

<https://openreview.net/forum?id=H1xwNhCcYm>

Detecting Out-of-Distribution Inputs to Deep Generative Models Using a Test for Typicality

Eric Nalisnick*

DeepMind

enalisnick@google.com

Akihiro Matsukawa

DeepMind

amatsukawa@google.com

Yee Whye Teh

DeepMind

ywteh@google.com

Balaji Lakshminarayanan*

DeepMind

balajiln@google.com

Coming soon to ArXiv...

Thank you. Questions?

In collaboration with...



Aki Matsukawa



Balaji
Lakshminarayanan



Dilan Gorur



Yee Whye Teh