

# Normalizing Flows

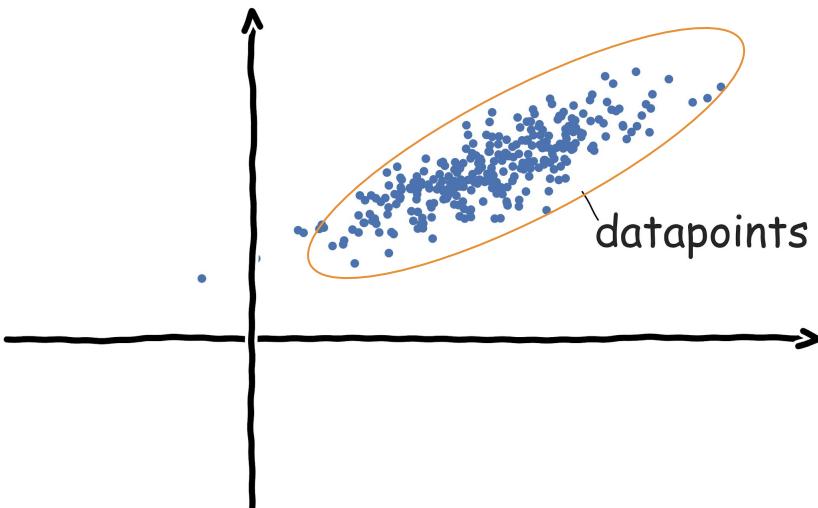
Arsenii Ashukha

PhD Candidate at Samsung AI Center Moscow



# Probabilistic View on Generative models

**Generative models** – aim to learn a data distribution.



**Data** points are vectors  $x \in \mathbb{R}^2$

**Goal** is to learn a model so that:

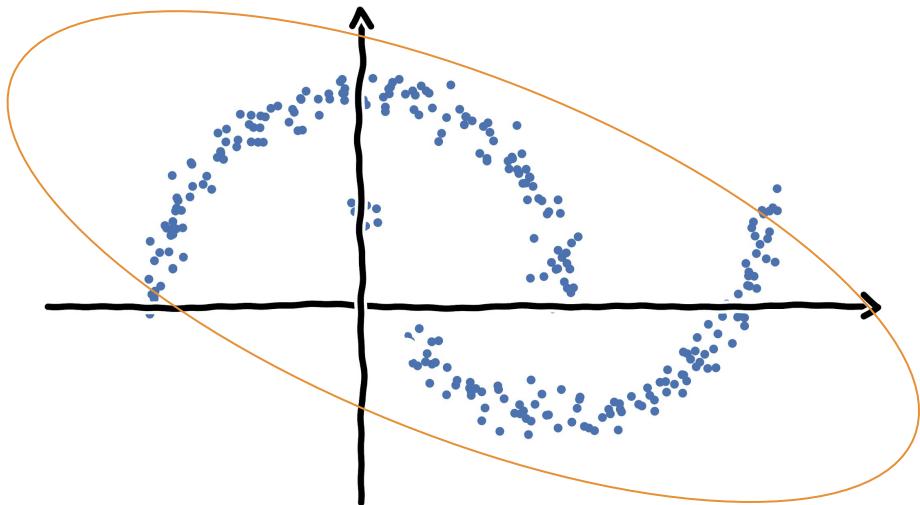
$$p_{model}(x) \approx p_{data}(x)$$

**Reasonable model for the case is  
2d-Gaussian distribution :**

$$p_{model} = \mathcal{N}(x | \mu, \sigma^2)$$

# Probabilistic View on Generative models

**Generative models** – aim to learn a data distribution.



Data points are vectors  $x \in \mathbb{R}^2$

**Goal** is to learn a model so that:

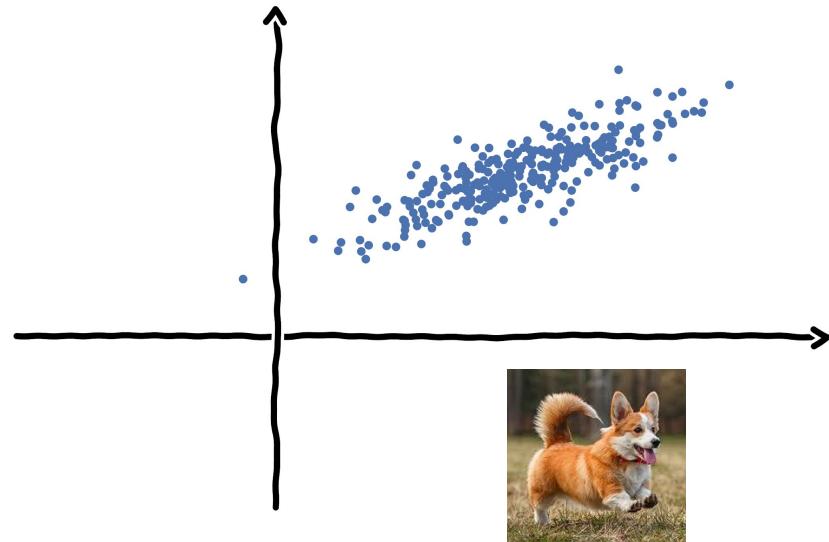
$$p_{model}(x) \approx p_{data}(x)$$

**Reasonable model for the case is  
2d-Gaussian distribution :**

$$p_{model} = \mathcal{N}(x | \mu, \sigma^2)$$

The model mismatches the data :(

# Ok, but what do we need this for?



1. Generation of new objects:
  - a. data augmentation
  - b. generation of private data for learning
  - c. specific problems: translation / text2speech, superresolution, ...
2. Probability estimate (out-of-distribution detection )
3. Representation learning (semi-supervised learning)
4. ....

# Gradient-based training of probabilistic GMs

**Get samples from target data distribution (dataset):**

$$X = (x_1, \dots, x_N), \quad x_i \sim p_{\text{data}}(x)$$

**Define a parametric model:**

$$p_{\theta}(x)$$

- assume that a model is suitable for gradient optimization by  $\theta$
- a model gives a p.d.f. value for any point  $x$  (sum-ups to 1, nonnegative)
- desirable to have a tractable sampling procedure

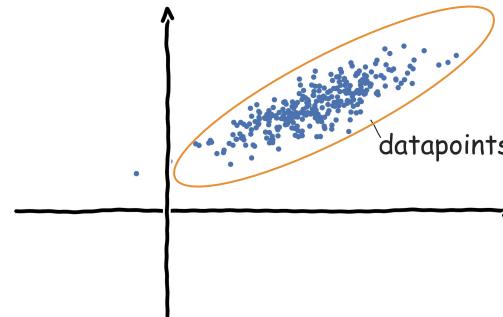
**Train a model by Maximum likelihood:**

$$\log p_{\theta}(X) = \sum_i \log p_{\theta}(x_i) \rightarrow \max_{\theta}$$

can be solved by SGA or Adam-like methods

# Gradient-based training of probabilistic GMs

Get samples from target  
data distribution (dataset):



Define a parametric model:

$$p_{\theta}(x) = \mathcal{N}(x | \mu, \Sigma) \quad \theta = (\mu, \Sigma)$$

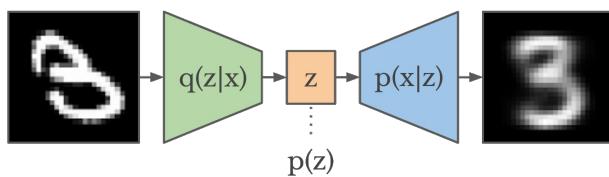
Train a model by Maximum  
likelihood:

$$\sum_i \log \mathcal{N}(x_i | \mu, \Sigma) \rightarrow \max_{\mu, \Sigma}$$
$$\begin{aligned} \mu &= [2 \ 2] \\ \Sigma &= \\ &\quad [[0.4 \ 0.2] \\ &\quad [0.2 \ 0.2]] \end{aligned}$$

# Gradient based training of probabilistic GMs

How to choose a model for more complex data?

Variational autoencoders  
(semi-implicit models)



- But can we define an expressive and tractable pdf based on DNNs?

- Yes!

**$p_\theta(x) = \text{Normalizing Flow} :)$**

$$\log p_\theta(x) = \log \mathbb{E}_{p(z)} p_\theta(x|z)$$

$$\log p_\theta(X) \rightarrow \max_{\theta}$$

$$\log p_\theta(X) \geq \mathcal{L}(X, \theta) \rightarrow \max_{\theta}$$

- Expressive parametric model
- Tractable computation/optimization of an exact  $p_\theta(x)$
- Fast and tractable sampling

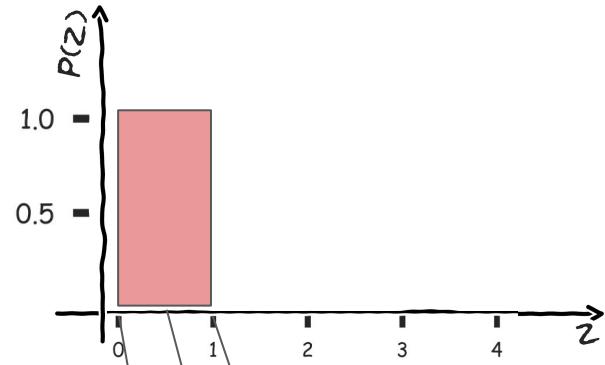
# Change of variables formula

1d-case:

Given:  $p(z) \quad z = f(x)$       What is  $p(x)$  ?

$$p(x)|dx| = p(z)|dz| \quad p(x) = p(z) \left| \frac{dz}{dx} \right|$$

$$p(x) = p(f(x)) \left| \frac{df(x)}{dx} \right|$$

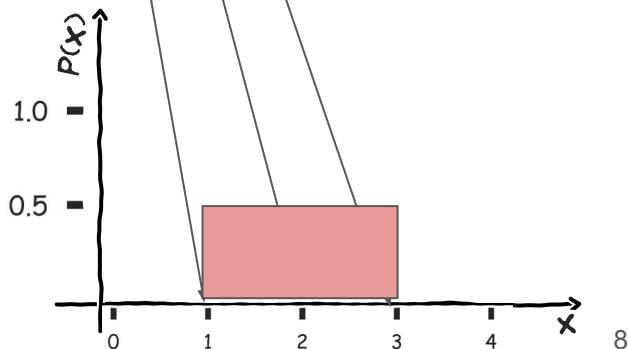


$$x = 2z + 1 \quad \frac{df(x)}{dx} = 0.5$$
$$z = (x - 1)/2$$

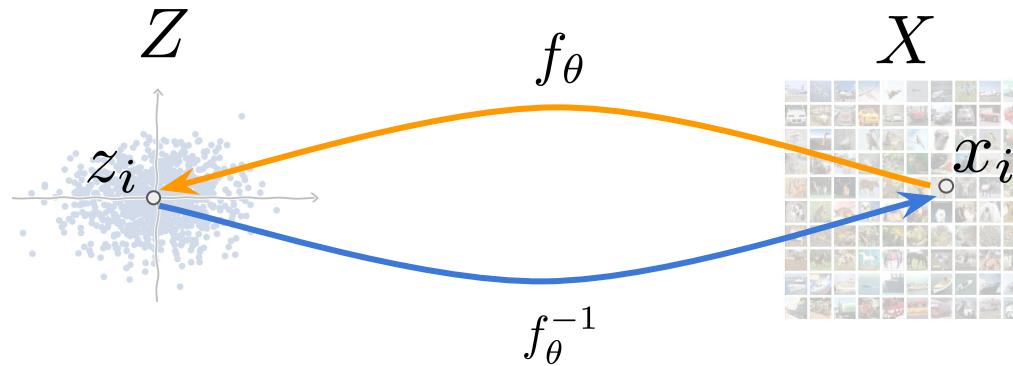
Nd-case:

$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x^T} \right|$$

$$\left( \frac{\partial f(x)}{\partial x^T} \right)_{ij} = \frac{\partial f(x)_i}{\partial x_j^T}$$



# Normalizing flows concept



$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x^T} \right|$$

# Coupling layer

Forward:

$$\begin{aligned}y_{1:d} &= x_{1:d} \\y_{d+1:D} &= x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d})\end{aligned}$$

Inverse:

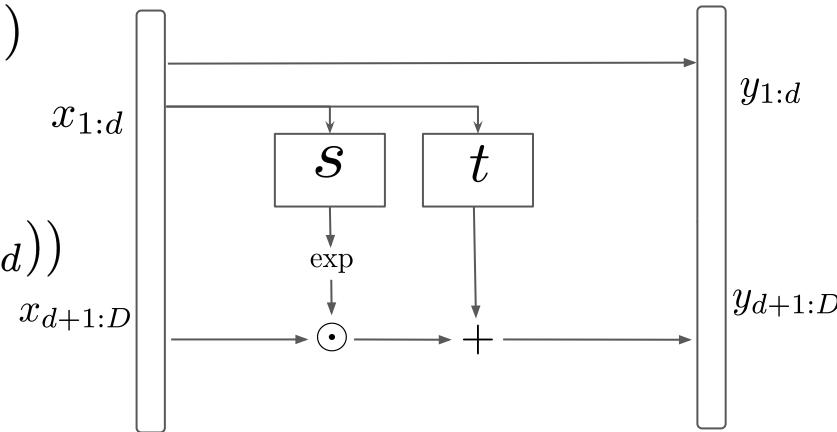
$$\begin{aligned}x_{1:d} &= y_{1:d} \\x_{d+1:D} &= (y_{d+1:D} - t(y_{1:d})) \odot \exp(-s(y_{1:d}))\end{aligned}$$

Jacobian:

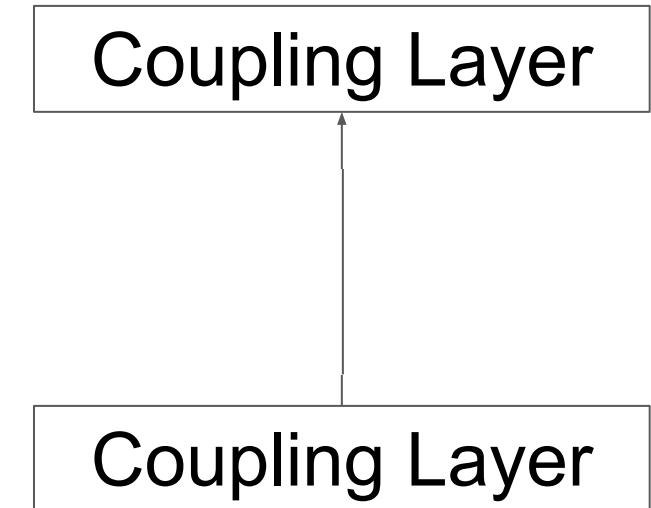
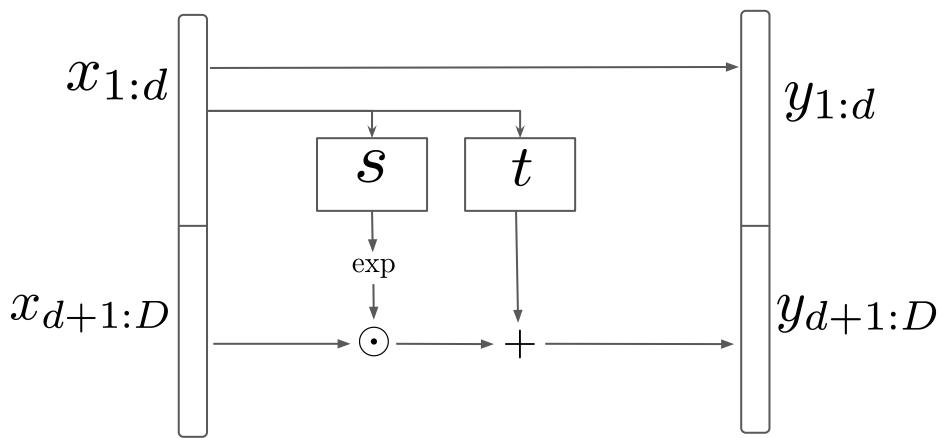
$$\frac{\partial y}{\partial x^T} = \left[ \quad \right]$$

$\det J$ :

$$\det \left| \frac{\partial y}{\partial x^T} \right| = \exp \sum_j s(x_{1:d})_j$$



# Combining coupling layers



# The model

$$f_\theta(x) : \mathbb{R}^D \rightarrow \mathbb{R}^D \quad f_\theta(x) = f^N \circ \cdots \circ f^1(x)$$

$$\log p_\theta(x) = \log p(f_\theta(x)) + \log \left| \det \frac{\partial f_\theta(x)}{\partial x^T} \right| = \log p(f(x)) + \sum_{i=1}^N \log \left| \det \frac{\partial f^i}{\partial x^{i-1}} \right|$$

$$\log p_\theta(X) = \sum_i \log p_\theta(x_i) \rightarrow \max_\theta$$

# Simple data

## Inference

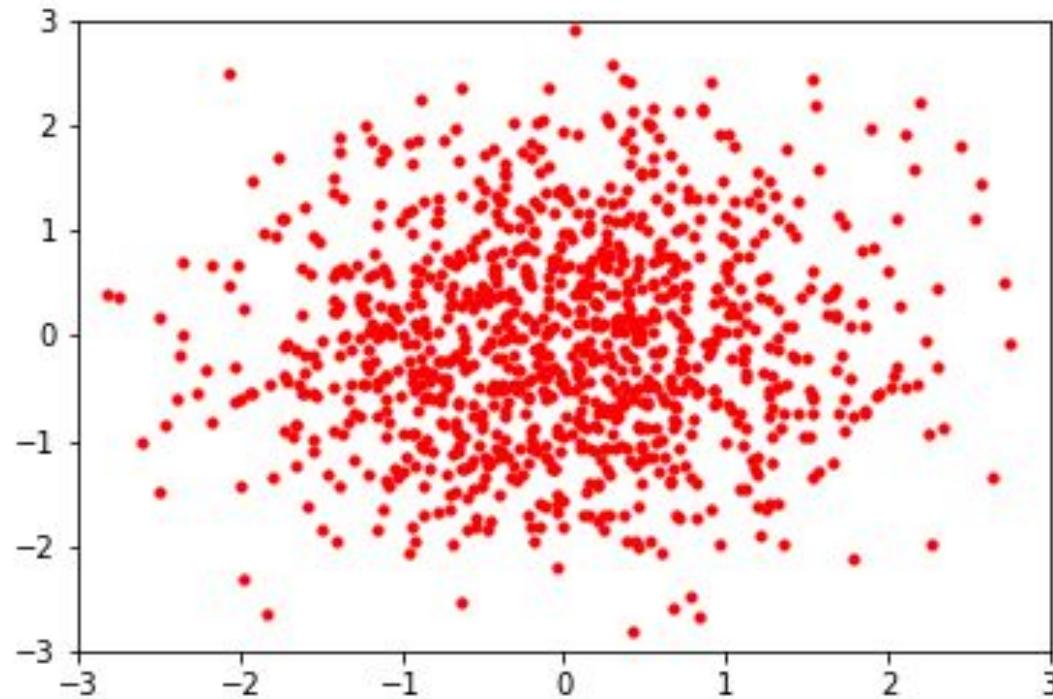
$$x \sim \hat{p}_X$$

$$z = f(x)$$

?

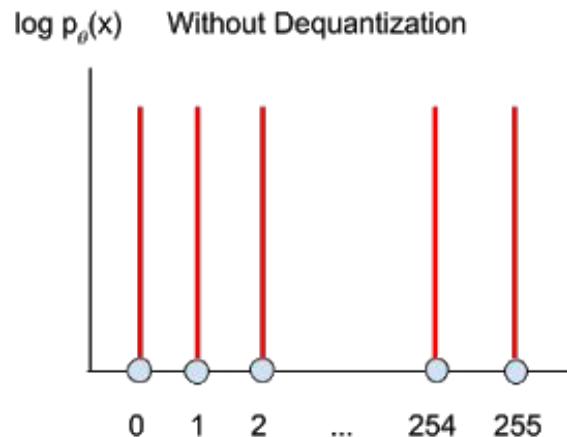
# Simple data

$$f^{-1}(z) = f_1^{-1} \circ \cdots \circ f_N^{-1}(z)$$



# Image generation: Data Dequantization

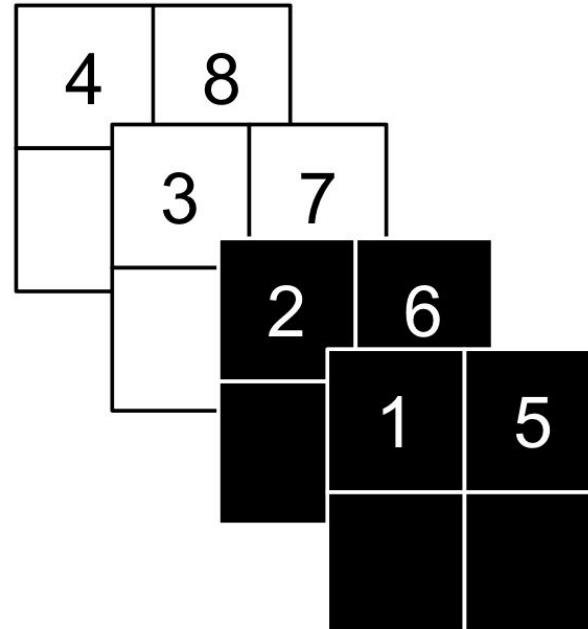
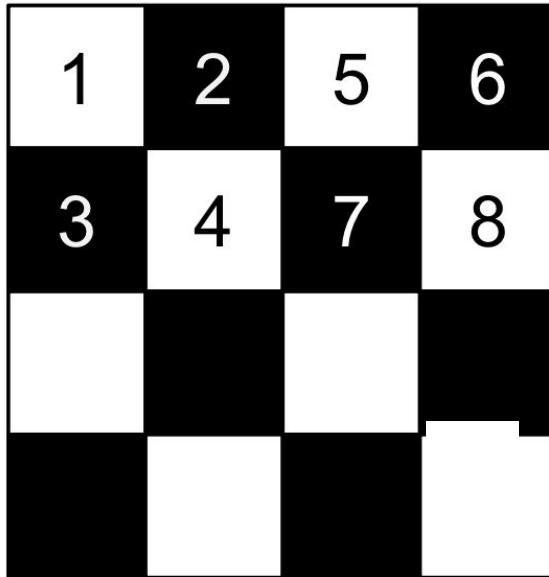
$$x_i \in \{0, \dots, 255\}$$



# Image generation: Masking

$$y_{1:d} = x_{1:d}$$

$$y_{d+1:D} = x_{d+1:D} \odot \exp(s(x_{1:d})) + t(x_{1:d})$$



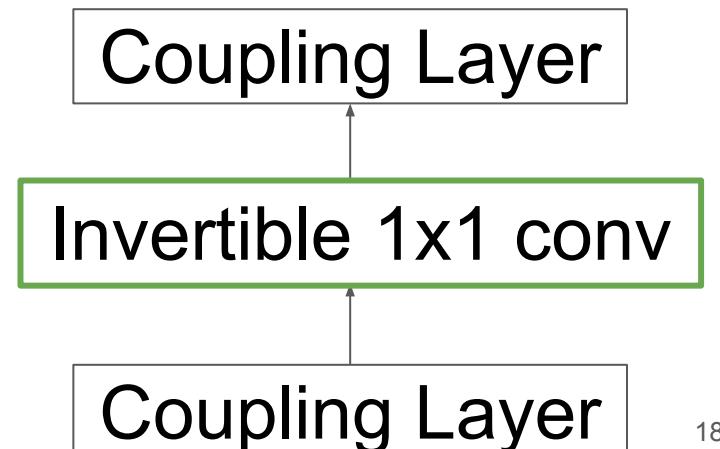
# Image results: Real NVP



# Image results: Glow



Figure 5: Linear interpolation in latent space between real images



# Measuring Quality

**Continuous data:** average test log-likelihood

$$\frac{1}{|X_{\text{test}}|} \sum_{x \in X_{\text{test}}} \log p_{\theta}(x)$$

**De-quantized data:** an estimate of probability (so-called bits-per-dimension)

$$\text{Bits-per-dim}(x) = -\frac{1}{D} \log \hat{P}(x)$$

# Measuring Quality: Results

| Dataset                      | PixelRNN [46] | Real NVP    | Conv DRAW [22] | IAF-VAE [34] |
|------------------------------|---------------|-------------|----------------|--------------|
| <b>CIFAR-10</b>              | 3.00          | 3.49        | < 3.59         | < 3.28       |
| <b>Imagenet (32 × 32)</b>    | 3.86 (3.83)   | 4.28 (4.26) | < 4.40 (4.35)  |              |
| <b>Imagenet (64 × 64)</b>    | 3.63 (3.57)   | 3.98 (3.75) | < 4.10 (4.04)  |              |
| <b>LSUN (bedroom)</b>        |               | 2.72 (2.70) |                |              |
| <b>LSUN (tower)</b>          |               | 2.81 (2.78) |                |              |
| <b>LSUN (church outdoor)</b> |               | 3.08 (2.94) |                |              |
| <b>CelebA</b>                |               | 3.02 (2.97) |                |              |

| Model   | CIFAR-10    | ImageNet 32x32 | ImageNet 64x64 | LSUN (bedroom) | LSUN (tower) | LSUN (church outdoor) |
|---------|-------------|----------------|----------------|----------------|--------------|-----------------------|
| RealNVP | 3.49        | 4.28           | 3.98           | 2.72           | 2.81         | 3.08                  |
| Glow    | <b>3.35</b> | <b>4.09</b>    | <b>3.81</b>    | <b>2.38</b>    | <b>2.46</b>  | <b>2.67</b>           |

# Conditional Normalizing Flows

$$D = \{(x, y)\}_{i=0, \dots, N} \quad p(x, y) = p(x|y)p(y)$$

$$\log p(x | y) = \log p(f(x) | y) + \log \left| \det \frac{\partial f_y(x)}{\partial x^T} \right|$$

$$s(x_{1:d}), t(x_{1:d}) \rightarrow s(x_{1:d}, y), t(x_{1:d}, y)$$

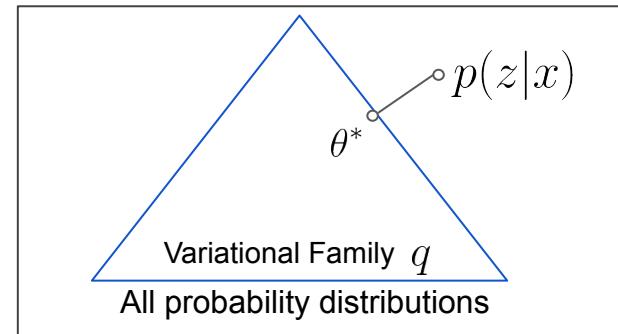
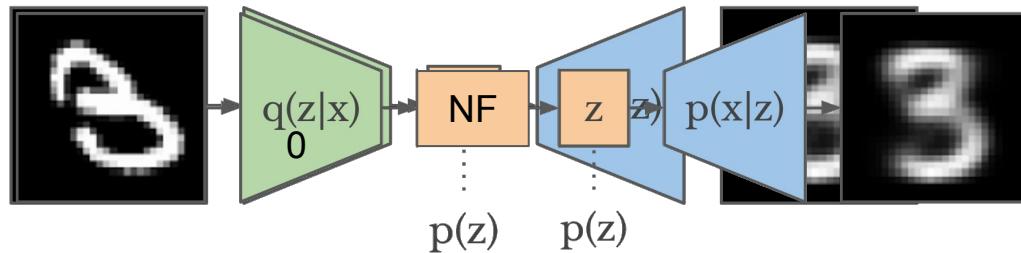


(b) Class conditional  $32 \times 32$  ImageNet samples

# Variational Inference with Normalizing Flows

$$\log p(X) = \mathcal{L}(X, \theta) + KL(q_\theta(z|x) || p(z|x))$$

$$N(z | \mu(x), \text{diag}(\sigma(x)))$$



---

## Algorithm 1 Variational Inf. with Normalizing Flows

Parameters:  $\phi$  variational,  $\theta$  generative

**while** not converged **do**

$\mathbf{x} \leftarrow \{\text{Get mini-batch}\}$

$\mathbf{z}_0 \sim q_0(\bullet|\mathbf{x})$

$\mathbf{z}_K \leftarrow f_K \circ f_{K-1} \circ \dots \circ f_1(\mathbf{z}_0)$

$\mathcal{F}(\mathbf{x}) \approx \mathcal{F}(\mathbf{x}, \mathbf{z}_K)$

$\Delta\theta \propto -\nabla_\theta \mathcal{F}(\mathbf{x})$

$\Delta\phi \propto -\nabla_\phi \mathcal{F}(\mathbf{x})$

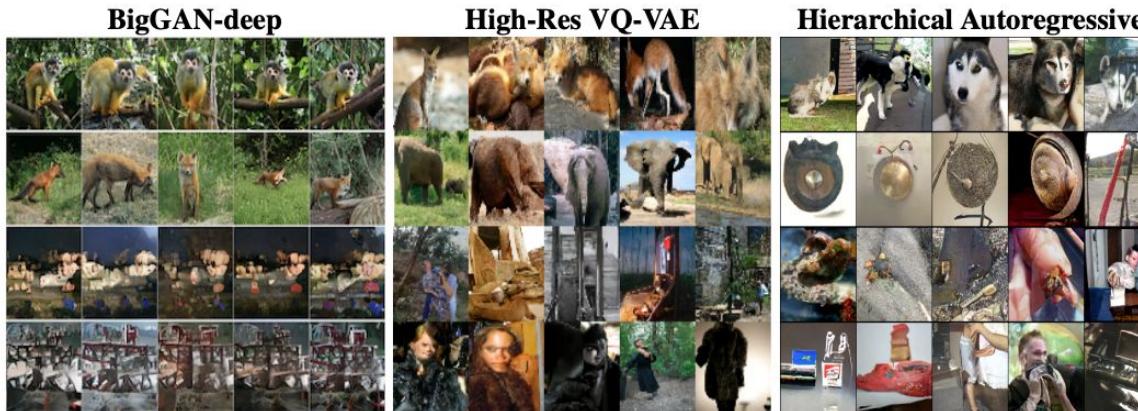
**end while**

---

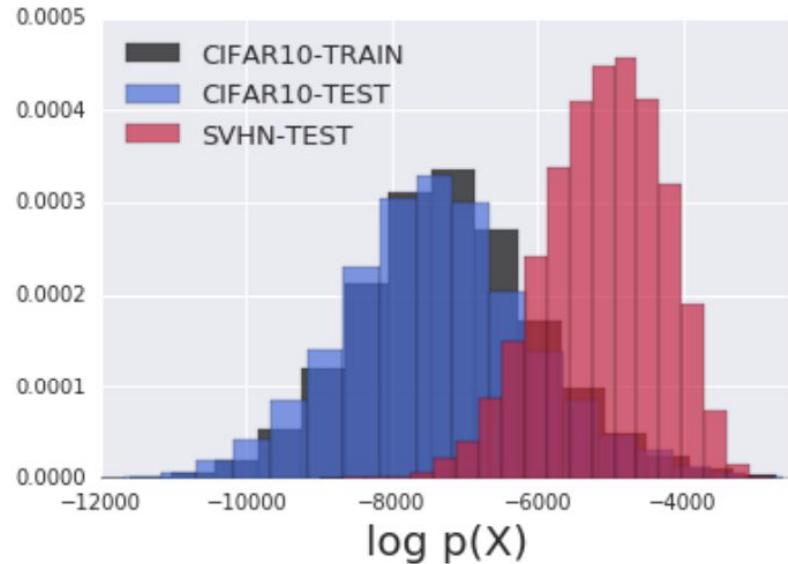
# Classification Accuracy Score for Conditional Generative Models

Table 1: CAS for different models at  $128 \times 128$  and  $256 \times 256$  resolutions. BigGAN-deep samples taken from best truncation parameter of 1.5.

| Training Set                                       | Resolution | Top-5 Accuracy | Top-1 Accuracy | Inception Score   | FID-50K |
|--|------------|----------------|----------------|-------------------|---------|
| Real<br>BigGAN-deep<br>Hierarchical Autoregressive | 128×128    | 88.79%         | 68.82%         | $165.38 \pm 2.84$ | 1.61    |
|  | 128×128    | 64.44%         | 40.64%         | $71.31 \pm 1.57$  | 4.22    |
|  | 128×128    | <b>77.33%</b>  | <b>54.05%</b>  | $17.02 \pm 0.79$  | 46.05   |
| Real<br>BigGAN-deep<br>High-Res VQ-VAE             | 256×256    | 91.47%         | 73.09%         | $331.83 \pm 5.00$ | 2.47    |
|  | 256×256    | 65.92%         | 42.65%         | $109.39 \pm 1.56$ | 11.78   |
|  | 256×256    | <b>77.59%</b>  | <b>54.83%</b>  | $43.44 \pm 0.87$  | 38.05   |



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



(b) Train on CIFAR-10, Test on SVHN

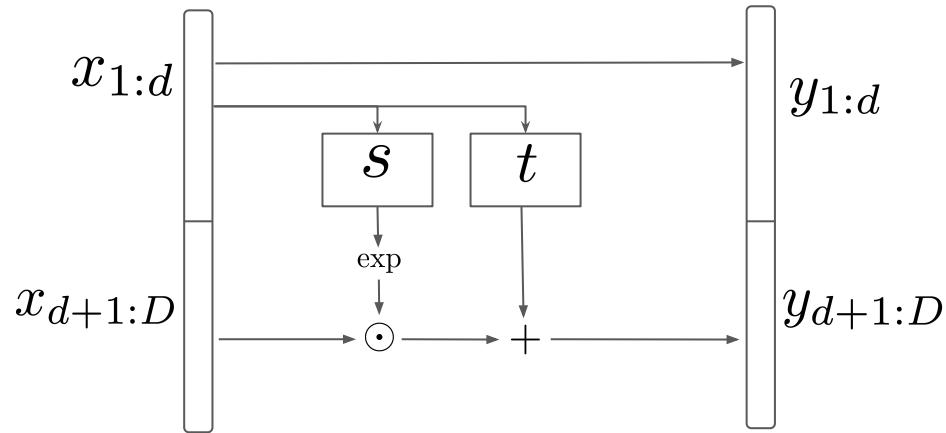
| Data Set                        | Avg. Bits Per Dimension |
|---------------------------------|-------------------------|
| <i>Glow Trained on CIFAR-10</i> |                         |
| CIFAR10-Train                   | 3.386                   |
| CIFAR10-Test                    | 3.464           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

# NFs

$$f_\theta(x) = f^N \circ \cdots \circ f^1(x)$$

$$p(x) = p(f(x)) \left| \det \frac{\partial f(x)}{\partial x^T} \right|$$

$$\log p_\theta(X) = \sum_i \log p_\theta(x_i) \rightarrow \max_\theta$$

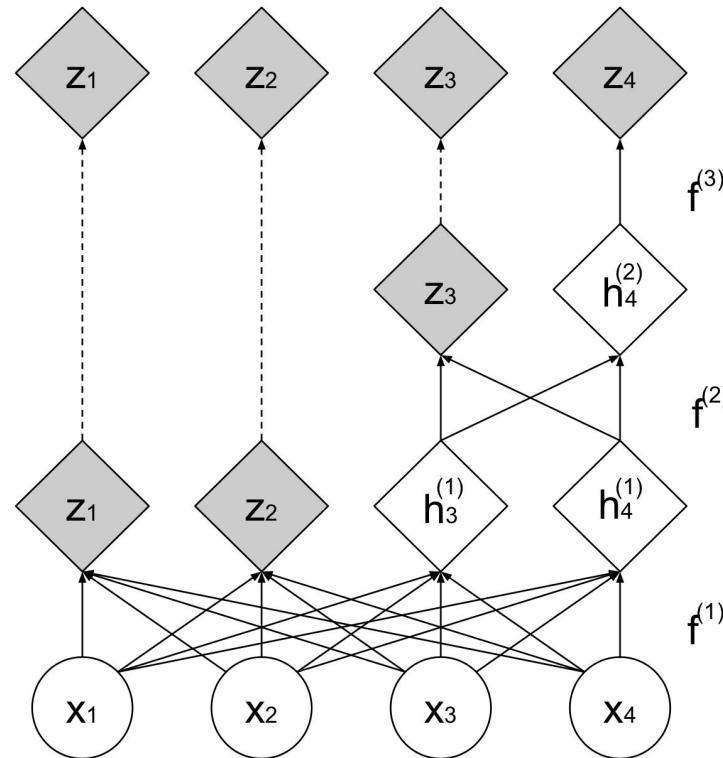


- Expressive parametric model
- Tractable computation/optimization of an exact  $p_\theta(x)$
- Fast and tractable sampling

# Reading List

- Tutorial on normalizing flows, [slideslive.com/38917907/tutorial-on-normalizing-flows](https://slideslive.com/38917907/tutorial-on-normalizing-flows)
- Tips for Training Likelihood Models, [blog.evjang.com/2019/07/likelihood-model-tips.html](https://blog.evjang.com/2019/07/likelihood-model-tips.html)
- FFJORD tutorial, [https://youtu.be/\\_ALdCSSVYkw](https://youtu.be/_ALdCSSVYkw)
- Must read papers:
  - Variational Inference with Normalizing Flows, <https://arxiv.org/abs/1505.05770>
  - Density estimation using Real NVP, <https://arxiv.org/abs/1605.08803>
  - Glow: Generative Flow with Invertible  $1 \times 1$  Convolutions <https://arxiv.org/abs/1807.03039>
  - Sylvester Normalizing Flows for Variational Inference, <https://arxiv.org/abs/1803.05649>
  - FFJORD: Free-form Continuous Dynamics for Scalable Reversible Generative Models, <https://arxiv.org/abs/1810.01367>
  - Do Deep Generative Models Know What They Don't Know?, <https://arxiv.org/abs/1810.09136>
  - Classification Accuracy Score, <https://arxiv.org/abs/1905.10887>

# Image generation: Multi-scale architecture

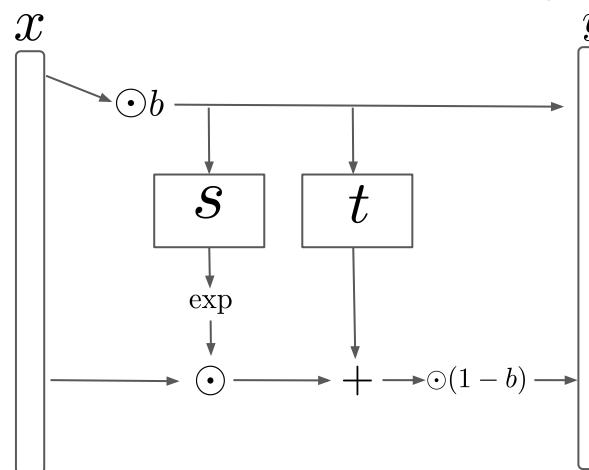


# Masked Coupling Layers (MCL)

## Coupling Layer

$$y = b \odot x + (1 - b) \odot \left( x \odot \exp(s(b \odot x)) + t(b \odot x) \right)$$

## Permutation



- $b$  is a layer specific binary mask
- $b_i = 1$ : copy
- $b_i = 0$ : transform