

# Lecture 5: Normalizing Flows

## Deep Generative Models

Sajjad Amini

Department of Electrical Engineering  
Sharif University of Technology

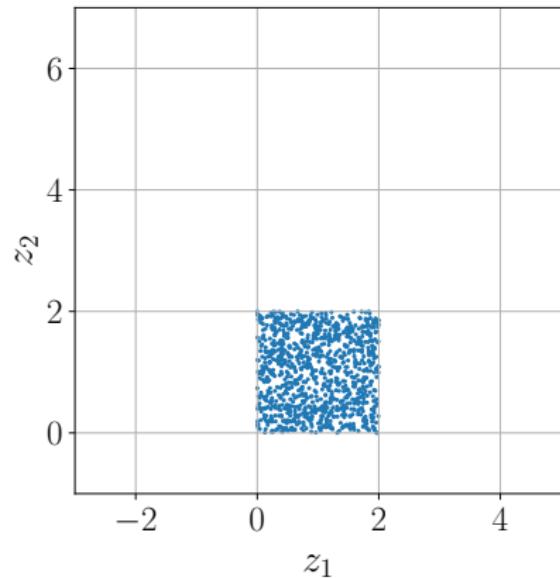
# Contents

- 1 Intuition
- 2 Change of Variable Theorem
- 3 Normalizing Flow Models
- 4 Learning
- 5 Nonlinear Independent Component Estimation
- 6 Real-valued Non-volume Preserving Transformations
- 7 Masked Autoregressive Flows
- 8 Inverse Autoregressive Flows
- 9 Conclusion

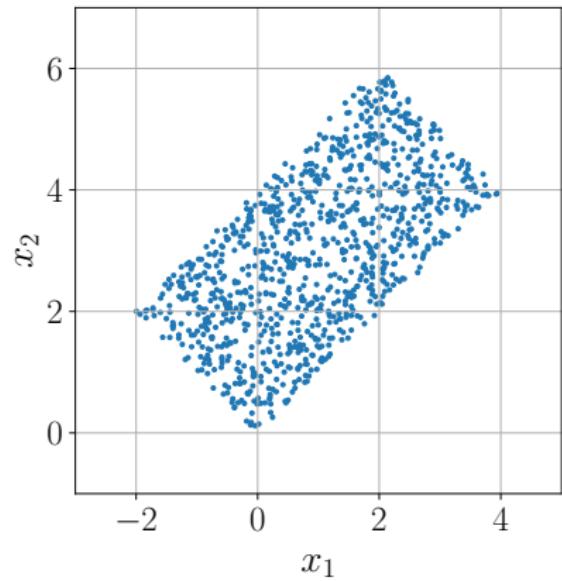
# Section 1

## Intuition

# Intuition



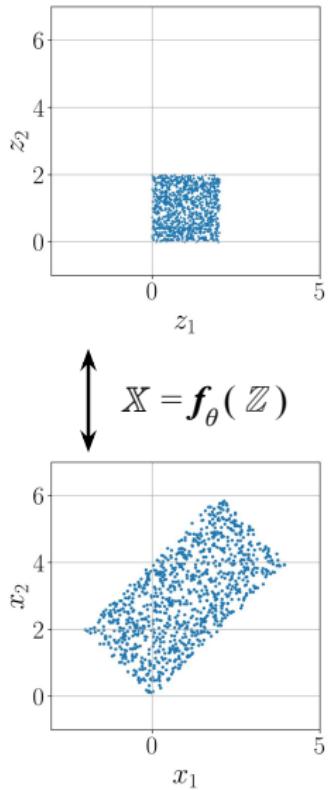
(a) Base random variable  $p(\mathbb{Z})$



(b) Data Samples  $p_{\text{data}}(\mathbb{X})$

**Figure:** The scenario: We have access to random variable  $\mathbb{Z}$  with  $p(\mathbb{Z})$  distribution while our target distribution is  $p_{\text{data}}(\mathbb{X})$

# General Idea



## General Idea

To achieve the target distribution using the base one, we follow these steps:

- Design a parametric transformation  $f_\theta(\cdot)$
  - Calculate the  $p_\theta(x)$
  - Maximize the dataset likelihood to find the transformation parameters  $\theta$
- ☞ But, how can we find  $p_\theta(x)$ ?

## Change of Variable

Assume  $\mathbf{z} \sim p(\mathbb{Z})$  and  $\mathbf{x} = f(\mathbf{z})$ . Given that the transformation  $f(\mathbf{z})$  is invertible, you can calculate the probability distribution  $p(\mathbb{X})$ .

## Section 2

### Change of Variable Theorem

# Transformation [1]

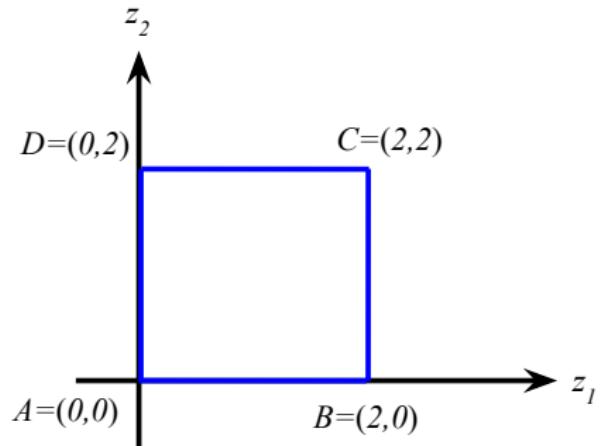


Figure: Base random variable

# Transformation [1]

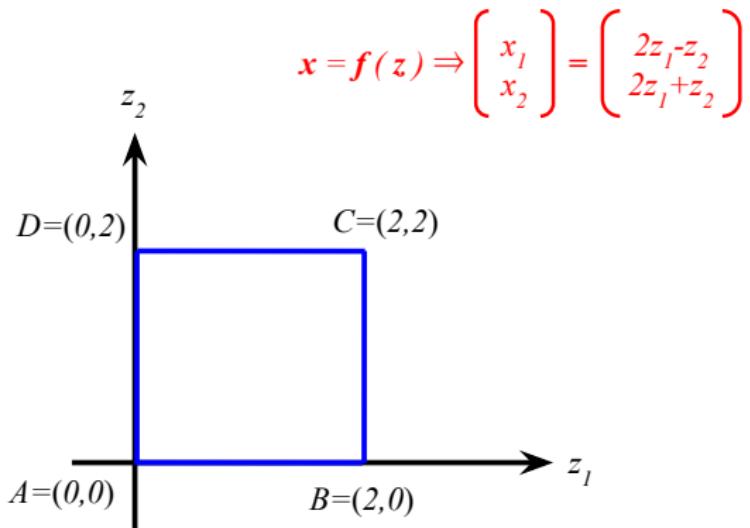


Figure: Transformation

# Transformation [1]

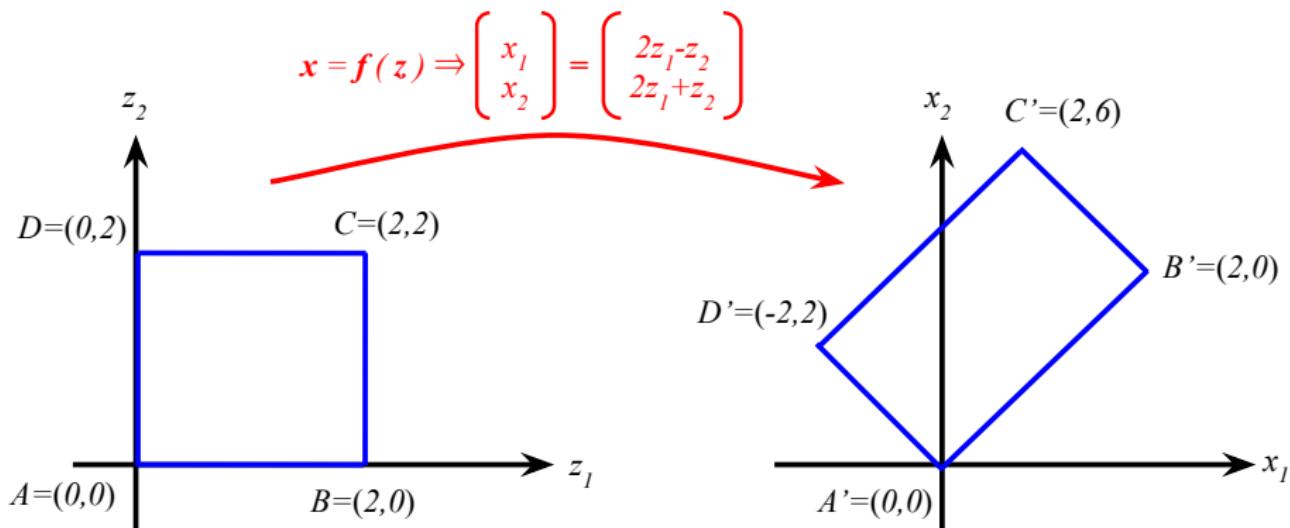
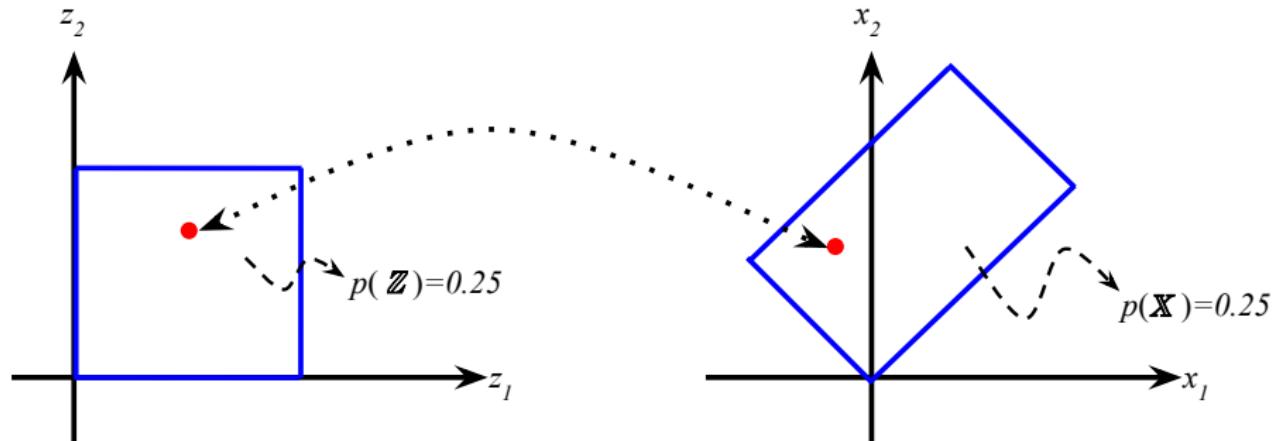


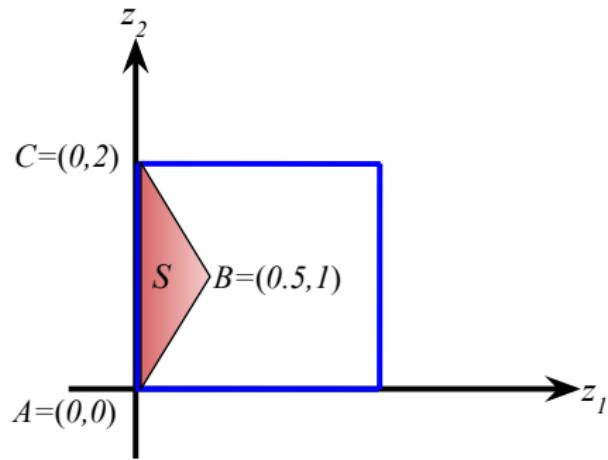
Figure: Transformed random variable

# Correspondence

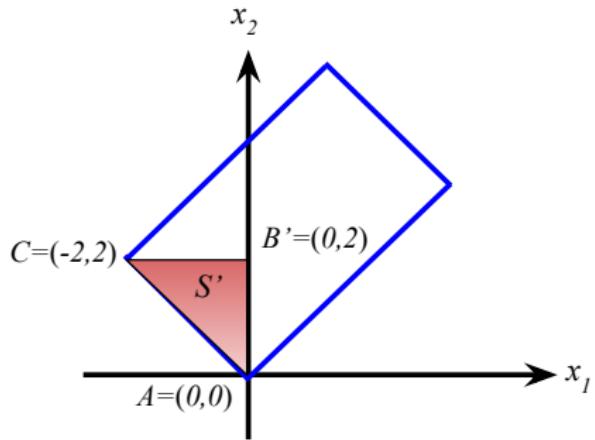


**Figure:** The point-to-point correspondence and probability density functions (we assume  $p(\mathbb{Z}) = 0.25$  and conclude that  $p(\mathbb{X}) = 0.25$  because of the point-to-point correspondence)

# Equivalent Events



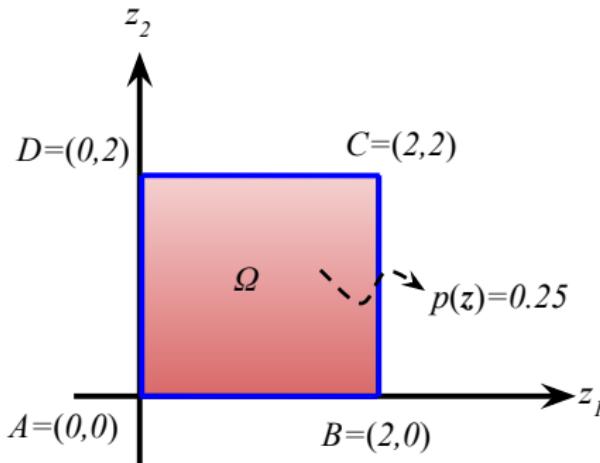
$$(a) p(S) = \text{Area}_S \times p(\mathbb{Z}) = \frac{1}{2} \times \frac{1}{4} = \frac{1}{8}$$



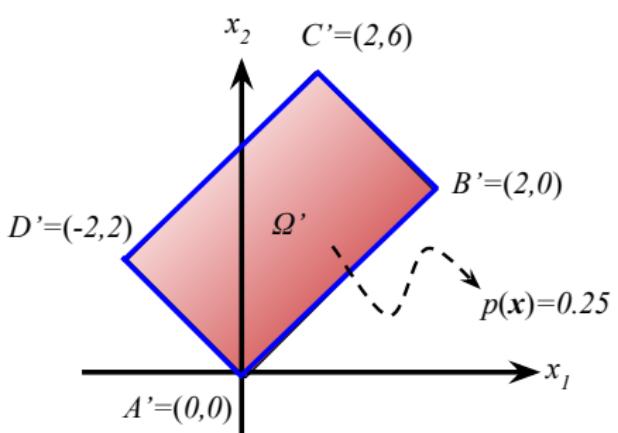
$$(b) p(S') = \text{Area}_{S'} \times p(\mathbb{X}) = 2 \times \frac{1}{4} = \frac{1}{2}$$

**Figure:** Probability calculations for equivalent events (while we expect  $p(S) = p(S')$  because they are equivalent events, we see a contradiction)

# Whole Sample Space Probability



$$(a) \quad p(\Omega) = \text{Area}_{\Omega} \times p(\mathbb{Z}) = 4 \times \frac{1}{4} = 1$$



$$(b) \quad p(\Omega') = \text{Area}_{\Omega'} \times p(\mathbb{X}) = 16 \times \frac{1}{4} = 4$$

**Figure:** Probability calculations for the whole sample space (while the Sample space must have probability 1, but for  $p(\mathbb{X})$  this value is greater than one which shows a clear contradiction with probability axioms)

# Change of Variable

## Theorem

Assume  $\mathbb{Z} \in \mathbb{R}^D$  is a random variable with base distribution  $p_z(\mathbb{Z})$ . Also assumes:

- $\mathbb{Z}$  passes through an invertible function  $\mathbf{f} : \mathbb{R}^D \rightarrow \mathbb{R}^D$  to generate random variable  $\mathbb{X}$  with a more complex distribution.
- We denote the inverse function  $\mathbf{f}^{-1}$  by  $\mathbf{g}$  such that: 
$$\begin{cases} \mathbf{x} = \mathbf{f}(\mathbf{z}) \\ \mathbf{z} = \mathbf{g}(\mathbf{x}) \end{cases}$$

then:

$$\begin{aligned} p_x(\mathbf{x}) &= p_z(\mathbf{g}(\mathbf{x})) \left| \det \mathbf{J}_g(\mathbf{x}) \right| \\ &= p_z(\mathbf{z}) \left| \det \mathbf{J}_f(\mathbf{z}) \right|^{-1} \end{aligned}$$

# General Idea

## General Idea

Normalizing flow models approximate the data distribution by:

- Selecting a simple base distribution  $p(\mathbb{Z})$  (e.g. standard Gaussian distribution)
- Selecting a flexible and *invertible* parameterized transformation  $\mathbf{f}_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$
- Assuming the observed random vector  $\mathbb{X} = \mathbf{f}_\theta(\mathbb{Z})$
- Calculating the exact distribution over  $\mathbb{X}$  using change of variable

# General Idea

$$\mathbf{x} = \mathbf{f}_\theta(\mathbf{z}) \Rightarrow \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} \theta_1 z_1 + \theta_2 z_2 \\ \theta_3 z_1 + \theta_4 z_2 \end{pmatrix}$$

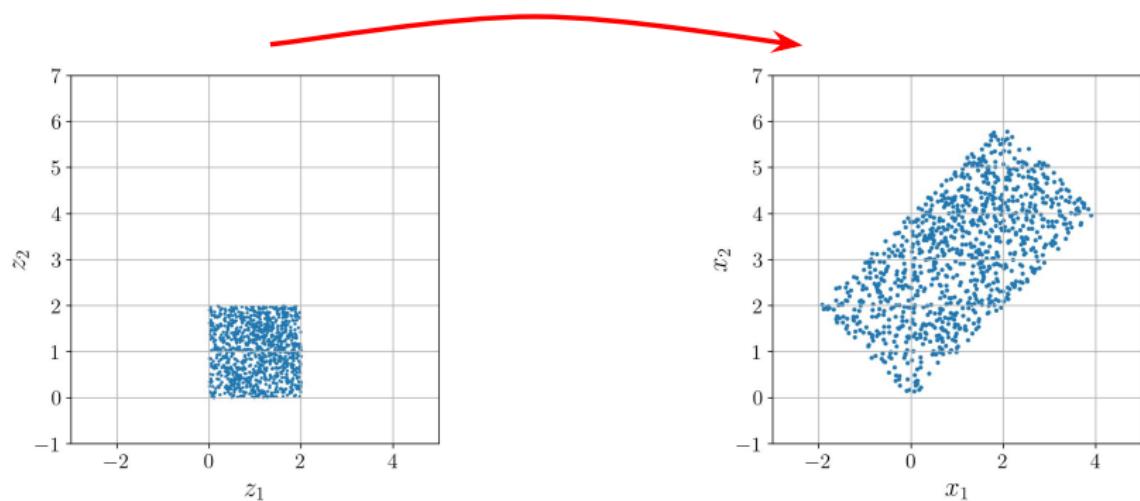


Figure: Sample of complex distribution using normalizing flow

# Likelihood Calculation

## Likelihood Calculation

Assume:

$$p(\mathbb{Z}) = \frac{1}{4}, \begin{cases} 0 \leq z_1 \leq 2 \\ 0 \leq z_2 \leq 2 \end{cases} \quad \text{and} \quad \mathbf{f}_\theta(\mathbf{z}) = \begin{bmatrix} \theta_1 z_1 + \theta_2 z_2 \\ \theta_3 z_1 + \theta_4 z_2 \end{bmatrix} = \begin{bmatrix} \theta_1 & \theta_2 \\ \theta_3 & \theta_4 \end{bmatrix} \underbrace{\begin{bmatrix} z_1 \\ z_2 \end{bmatrix}}_{\mathbf{z}}$$

Then:

$$p_\theta(\mathbf{x}) = \frac{p(\mathbf{z})}{|\det \mathbf{J}_f(\mathbf{u})|} = \frac{\frac{1}{4}}{\left| \det \begin{bmatrix} \theta_1 & \theta_2 \\ \theta_3 & \theta_4 \end{bmatrix} \right|} = \frac{\frac{1}{4}}{|\theta_1 \theta_4 - \theta_2 \theta_3|}$$

Note that the range for  $x_1$  and  $x_2$  is defined based on the inverse transformation and range for  $z_1$  and  $z_2$ .

# Transformation [1]

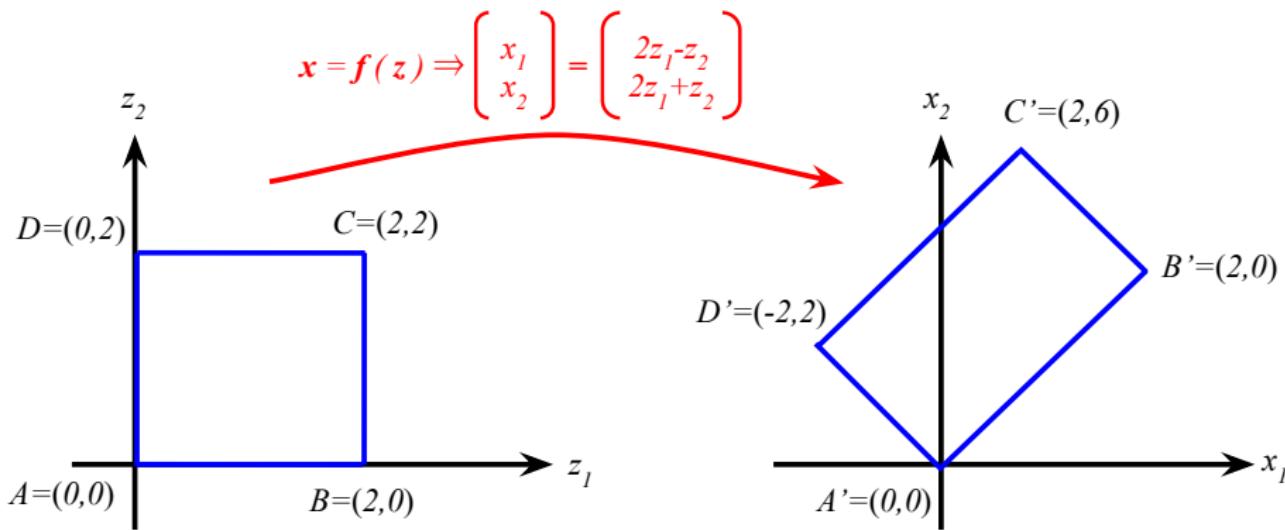
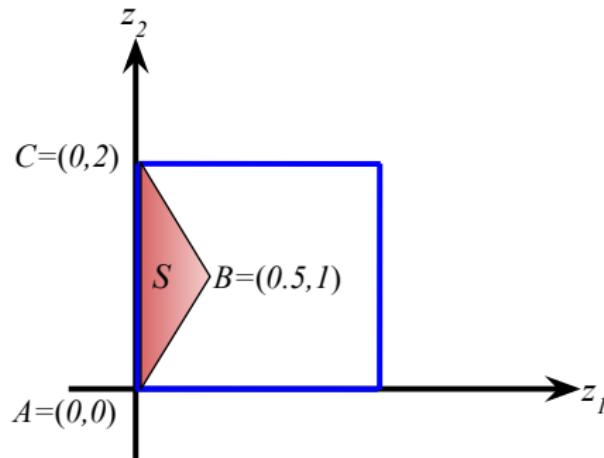


Figure: Sample transformation

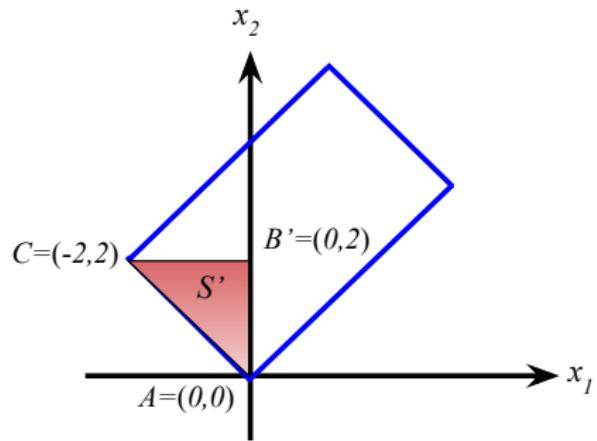
## Probability Density Estimation

In the above example  $\begin{bmatrix} \theta_1 & \theta_2 \\ \theta_3 & \theta_4 \end{bmatrix} = \begin{bmatrix} 2 & -1 \\ 2 & 1 \end{bmatrix}$ , thus  $p_{\theta}(\mathbb{X}) = \frac{1}{4} = \frac{1}{16}$

# Equivalent Events (Revisited)



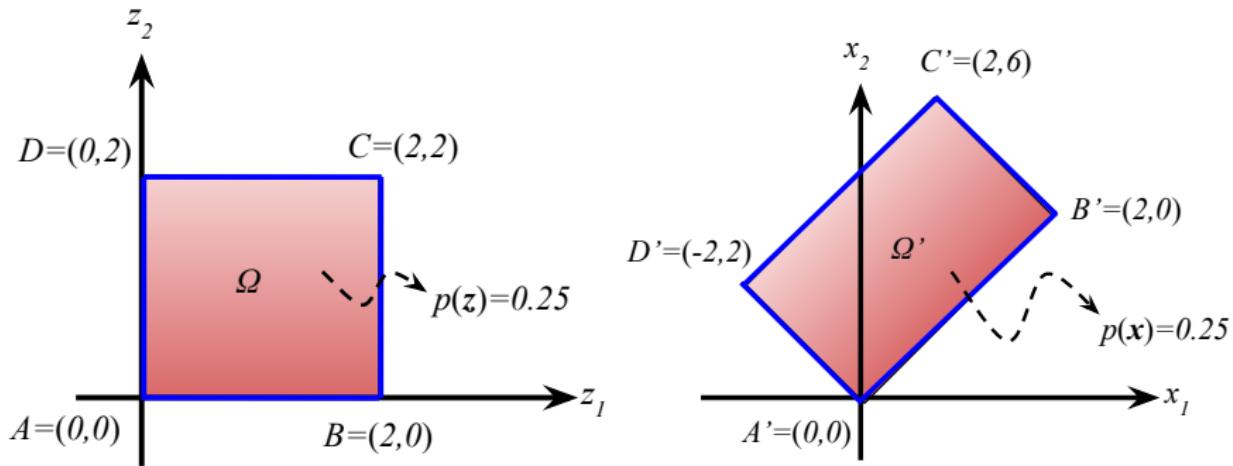
$$(a) p(S) = \text{Area}_S \times p(\mathbb{Z}) = \frac{1}{2} \times \frac{1}{4} = \frac{1}{8}$$



$$(b) p(S') = \text{Area}_{S'} \times p(\mathbb{X}) = 2 \times \frac{1}{16} = \frac{1}{8}$$

**Figure:** Probability calculations for equivalent events (Using change of variable technique, we see the probabilities are equal)

# Whole Sample Space Probability (Revisited)



$$(a) p(\Omega) = \text{Area}_{\Omega} \times p(\mathbb{Z}) = 4 \times \frac{1}{4} = 1$$

$$(b) p(\Omega') = \text{Area}_{\Omega'} \times p(\mathbb{X}) = 16 \times \frac{1}{16} = 1$$

**Figure:** Probability calculations for the whole sample space (Using change of variable technique, we see the probability for the sample space in  $\mathbb{X}$  equals to one as expected)

## Section 3

### Normalizing Flow Models

# Normalizing Flow

## Definition

In normalizing flow, we have the following building blocks:

- Assume  $\mathbf{z} \sim p(\mathbb{Z})$ , where  $p(\mathbb{Z})$  is an *easy to sample* base distribution.
- We define a flexible parameterized transformation  $\mathbf{f}_\theta$  from base distribution  $p(\mathbb{Z})$  to estimate data distribution by  $p_\theta(\mathbb{X})$
- We optimize the transformation parameters by:

$$\boldsymbol{\theta}^* = \operatorname{argmin}_{\boldsymbol{\theta}} \text{KL} \left( p_{\text{data}}(\mathbb{X}) \| p_{\boldsymbol{\theta}}(\mathbb{X}) \right) \equiv \operatorname{argmax}_{\boldsymbol{\theta}} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}(\mathbb{X})} \left[ p_{\boldsymbol{\theta}}(\mathbf{x}) \right]$$

# Normalizing Flow vs VAE

## Comparison to VAE

- In contrast to VAE, here we can easily calculate the exact data likelihood as:

$$\begin{aligned} p_{\theta}(\mathbf{x}_i) &= p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \left| \det \mathbf{J}_{\mathbf{g}}(\mathbf{x}_i) \right| \\ &= \frac{p(\mathbf{z}_i)}{\left| \det \mathbf{J}_{\mathbf{f}}(\mathbf{z}_i) \right|}, \text{ where } \mathbf{z}_i = \mathbf{g}(\mathbf{x}_i) \end{aligned}$$

- Due to the invertibility of normalizing flow transformation, the latent vector  $\mathbb{Z}$  and observed vector  $\mathbb{X}$  are of the same dimension, thus:
  - There is no abstraction in the latent space.
  - Finding meaningful direction in the latent space is not possible.

# Single Invertible Transformation

## Forward Transformation

Assume  $\mathbf{f}_\theta : \mathbb{R}^D \rightarrow \mathbb{R}^D$  is an invertible transformation defined as:

$$\mathbf{x} = \mathbf{f}_\theta(\mathbf{z}) \Rightarrow \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{bmatrix} = \mathbf{f}_\theta \left( \begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_D \end{bmatrix} \right) \Rightarrow \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{bmatrix} = \begin{bmatrix} u_1(\mathbf{z}) \\ u_2(\mathbf{z}) \\ \dots \\ u_D(\mathbf{z}) \end{bmatrix}$$

then the Jacobian matrix is:

$$\mathbf{J}_{\mathbf{f}} = \begin{bmatrix} \frac{\partial u_1}{\partial z_1} & \dots & \frac{\partial u_1}{\partial z_D} \\ \vdots & \vdots & \vdots \\ \frac{\partial u_D}{\partial z_1} & \dots & \frac{\partial u_D}{\partial z_D} \end{bmatrix}$$

and we can calculate the data likelihood as:

$$p_\theta(\mathbf{x}) = p(\mathbf{g}(\mathbf{x})) \left| \det \mathbf{J}_f(\mathbf{z}) \right|^{-1}$$

## Recap Previous Section

### Inverse Transformation

Now assume the inverse transformation  $\mathbf{g}_\theta = \mathbf{f}_\theta^{-1}$  defined as:

$$\mathbf{z} = \mathbf{g}_\theta(\mathbf{x}) \Rightarrow \begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_D \end{bmatrix} = \mathbf{g}_\theta \left( \begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_D \end{bmatrix} \right) \Rightarrow \begin{bmatrix} z_1 \\ z_2 \\ \dots \\ z_D \end{bmatrix} = \begin{bmatrix} v_1(\mathbf{x}) \\ v_2(\mathbf{x}) \\ \dots \\ v_D(\mathbf{x}) \end{bmatrix}$$

then the Jacobian matrix is:

$$\mathbf{J}_{\mathbf{g}} = \begin{bmatrix} \frac{\partial v_1}{\partial x_1} & \dots & \frac{\partial v_1}{\partial x_D} \\ \vdots & \vdots & \vdots \\ \frac{\partial v_D}{\partial x_1} & \dots & \frac{\partial v_D}{\partial x_D} \end{bmatrix}$$

and we can calculate the data likelihood as:

$$p_\theta(\mathbf{x}) = p(\mathbf{g}(\mathbf{x})) \left| \det \mathbf{J}_{\mathbf{g}}(\mathbf{x}) \right|$$

# Composite Transformations

## Forward Transformation

To have more flexibility in Normalizing flow, the transformation is realized using  $L$  layers of invertible transformations  $\mathbf{f}_i, i = 1, \dots, L$  as:

$$\mathbf{f}_\theta = \mathbf{f}_{L,\theta_L} \circ \mathbf{f}_{L-1,\theta_{L-1}} \circ \dots \circ \mathbf{f}_{2,\theta_2} \circ \mathbf{f}_{1,\theta_1}$$

where  $\mathbf{f}_{i,\theta_i} : \mathbb{R}^D \rightarrow \mathbb{R}^D$  for  $1 \leq i \leq L$  and  $\theta_i$  represents the parameters for the  $i$ -th transformation.

- We omit  $\theta_i$  subscript for simplicity.

# Composite Transformations

## Inverse Transformation

Assume the composite forward transformation in Slide 23 as :

$$\mathbf{f}_\theta = \mathbf{f}_L \circ \mathbf{f}_{L-1} \circ \dots \circ \mathbf{f}_2 \circ \mathbf{f}_1$$

We know each transformation  $\mathbf{f}_i$  is invertible, denoted by  $\mathbf{g}_i$ , then the inverse transformation can also be represented in a compositional way as:

$$\mathbf{g}_\theta = \mathbf{g}_1 \circ \mathbf{g}_2 \circ \dots \circ \mathbf{g}_{L-1} \circ \mathbf{g}_L$$

- ☞ Note that the sequence is reversed in the inverse transformation.

# Intermediate Random Variables

## Intermediate Random Variables

Assume  $\mathbf{x} = \mathbf{f}(\mathbf{z})$ , then we name the intermediate variables as:

$$\mathbf{x} = \underbrace{\mathbf{f}_L \circ \mathbf{f}_{L-1} \circ \dots \circ \mathbf{f}_2}_{\mathbf{z}(2)} \circ \underbrace{\mathbf{f}_1(\mathbf{z})}_{\mathbf{z}(1)}$$

$\overbrace{\hspace{10em}}$   
 $\mathbf{z}(L-1)$

If we define  $\mathbf{z}(L) \triangleq \mathbf{x}$  and  $\mathbf{z}(0) \triangleq \mathbf{z}$ , then we have:

$$\mathbf{z}(l) = \mathbf{f}_l(\mathbf{z}(l-1)), \text{ for } l = 1, \dots, L$$

Equivalently using the inverse transformation we have:

$$\mathbf{z}(l-1) = \mathbf{g}_l(\mathbf{z}(l)), \text{ for } i = 1 \dots, L$$

# Jacobian of Composite

## Forward Transformation

As we see before, we have:

$$\mathbf{f}_\theta = \mathbf{f}_L \circ \mathbf{f}_{L-1} \circ \dots \circ \mathbf{f}_2 \circ \mathbf{f}_1$$

Using the chain rule, we can easily show that:

$$\mathbf{J}_f(\mathbf{z}) = \mathbf{J}_{f_L}(\mathbf{z}(L-1)) \times \mathbf{J}_{f_{L-1}}(\mathbf{z}(L-2)) \times \dots \times \mathbf{J}_{f_2}(\mathbf{z}(1)) \times \mathbf{J}_{f_1}(\mathbf{z})$$

Having access to the Jacobian of forward transformation, you can calculate the likelihood as:

$$p_\theta(\mathbf{x}) = p(\mathbf{g}_\theta(\mathbf{x})) \left| \det \mathbf{J}_f(\mathbf{z}) \right|^{-1}$$

# Jacobian of Composite

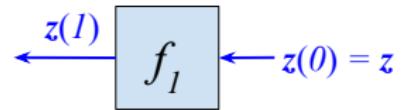


Figure: Base random variable

# Jacobian of Composite

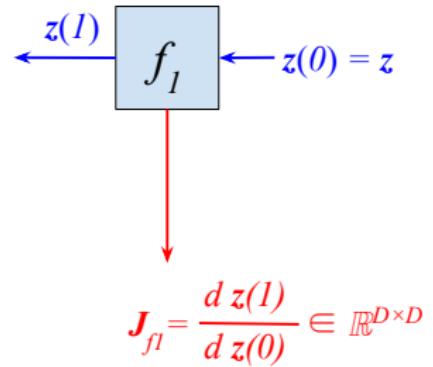


Figure: Transformation

# Jacobian of Composite

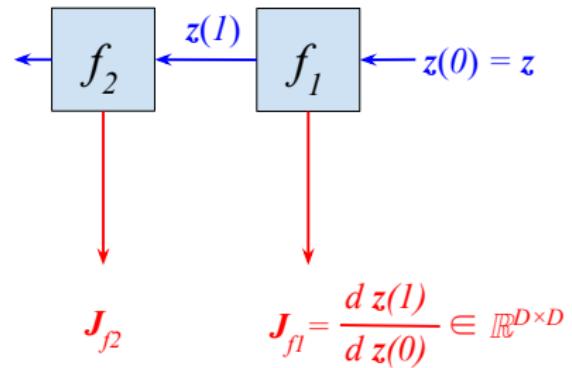


Figure: Transformed random variable

# Jacobian of Composite

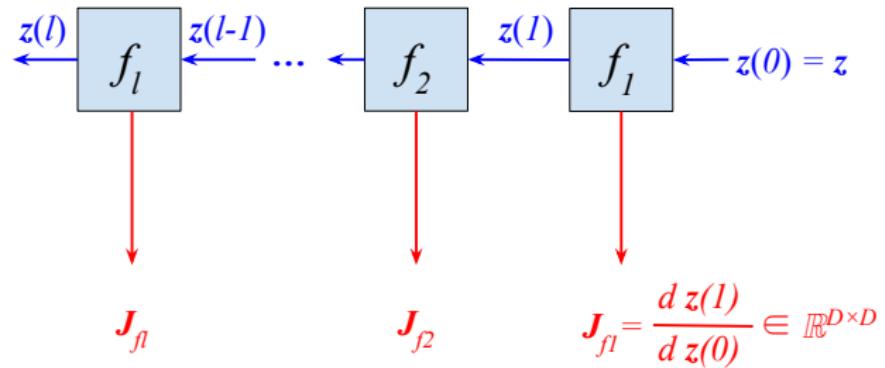


Figure: Transformed random variable

# Jacobian of Composite

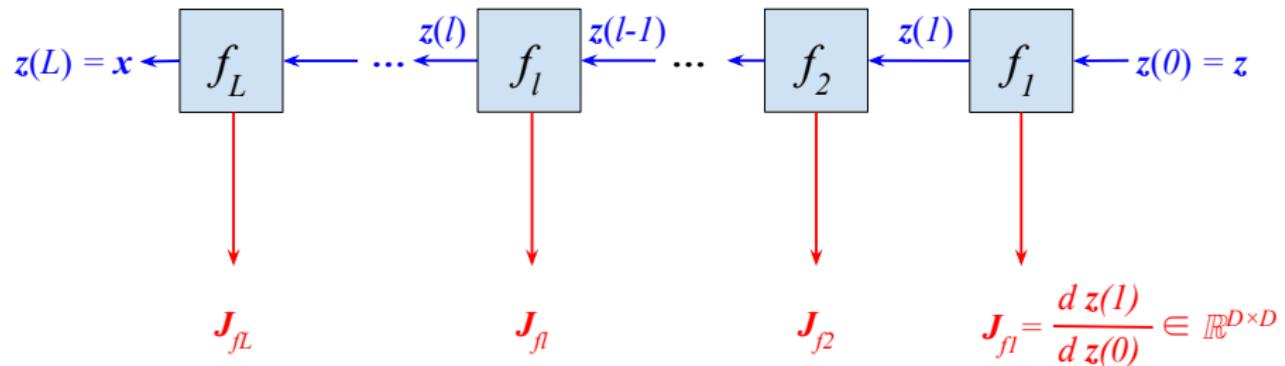


Figure: Transformed random variable

# Jacobian of Composite

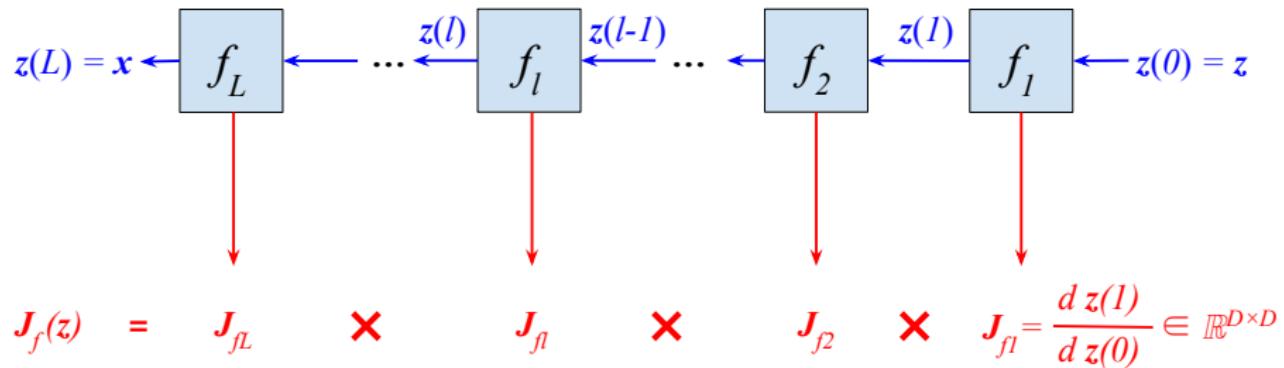


Figure: Transformed random variable

# Jacobian of Composite

## Inverse Transformation

Similarly to the forward composite transformation, we have:

$$\mathbf{g}_\theta = \mathbf{g}_1 \circ \mathbf{g}_2 \circ \dots \circ \mathbf{g}_{L-1} \circ \mathbf{g}_L$$

Using the chain rule, we can easily show that:

$$\mathbf{J}_g(\mathbf{x}) = \mathbf{J}_{g_1}(z(1)) \times \mathbf{J}_{g_2}(z(2)) \times \dots \times \mathbf{J}_{g_{L-1}}(z(L-1)) \times \mathbf{J}_{g_L}(\mathbf{x})$$

Having access to the Jacobian of inverse transformation, you can calculate the likelihood as:

$$p_\theta(\mathbf{x}) = p(\mathbf{g}(\mathbf{x})) \left| \det \mathbf{J}_g(\mathbf{x}) \right|$$

# Jacobian of Composite

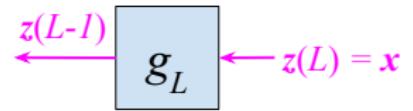


Figure: Base random variable

# Jacobian of Composite

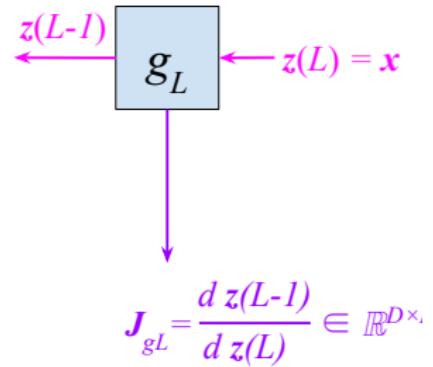


Figure: Transformation

# Jacobian of Composite

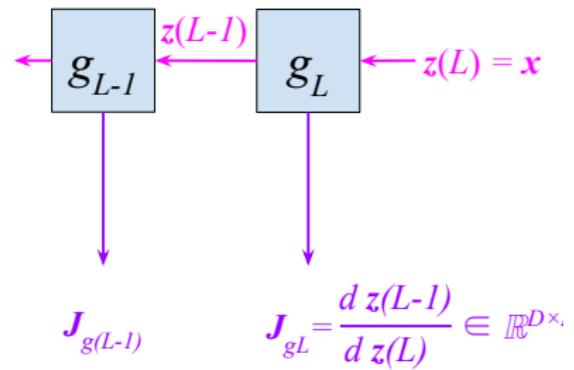


Figure: Transformed random variable

# Jacobian of Composite

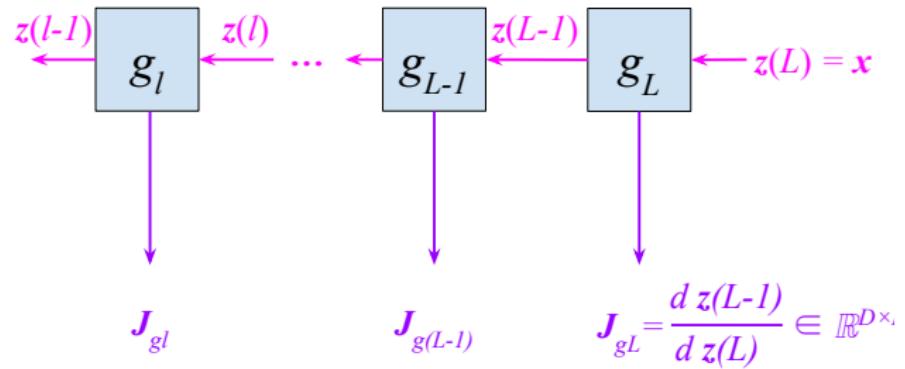


Figure: Transformed random variable

# Jacobian of Composite

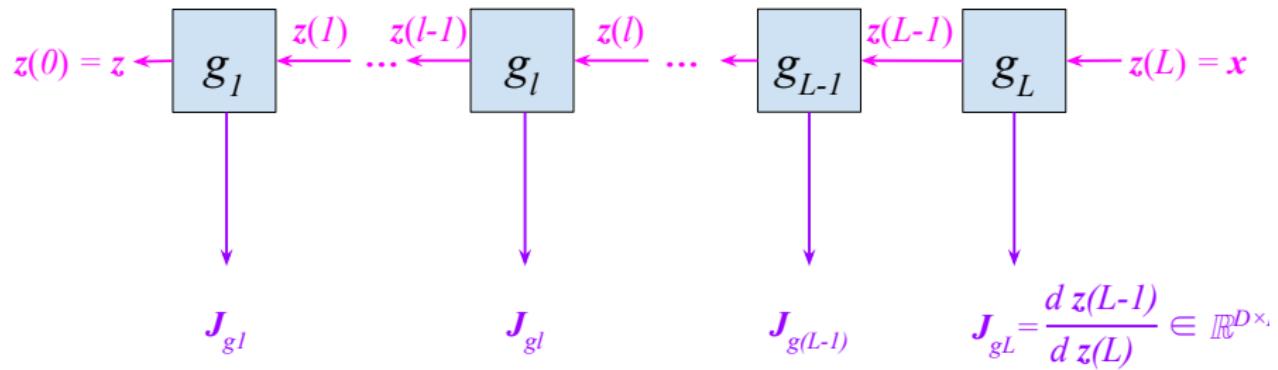


Figure: Transformed random variable

# Jacobian of Composite

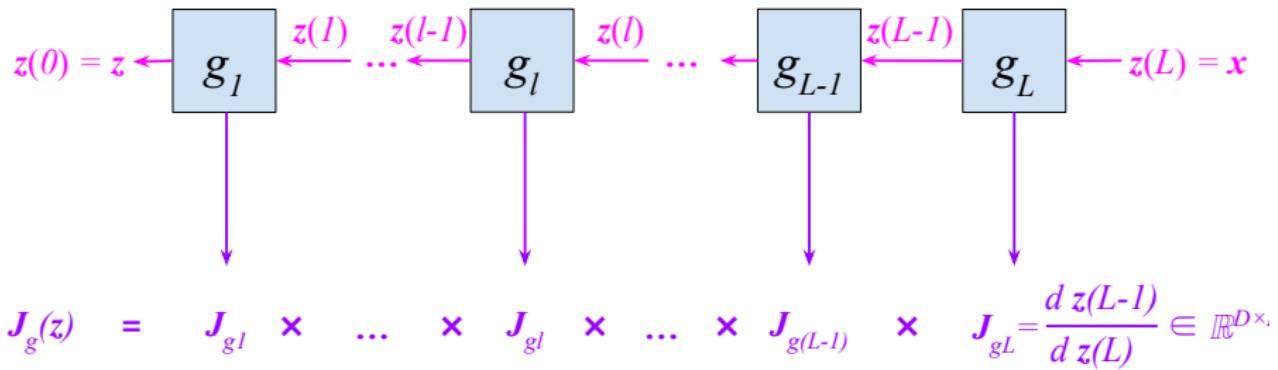


Figure: Transformed random variable

# Why Normalizing Flow?

## Why Flow?

The random variable  $\mathbb{Z}$  with easy-to-sample distribution  $p(\mathbb{Z})$  **flow** through  $L$  invertible transformations  $\mathbf{f}_l, l = 1 \dots L$  to generate the random variable  $\mathbb{X}$  with parametric distribution  $p_\theta(\mathbb{X})$ .

## Why Normalizing?

As we see before, the distribution over  $\mathbb{X}$  can be calculated as:

$$p_\theta(\mathbf{x}_i) = \frac{p(\mathbf{g}_\theta(\mathbf{x}_i))}{\left| \det \mathbf{J}_{\mathbf{f}}(\mathbf{z}_i) \right|}$$

in other word the term  $\left| \det \mathbf{J}_{\mathbf{f}}(\mathbf{z}_i) \right|$  Normalizes  $p(\mathbf{g}_\theta(\mathbf{x}_i))$  to be a valid probability distribution function.

# Generative Modeling Tasks: Generation

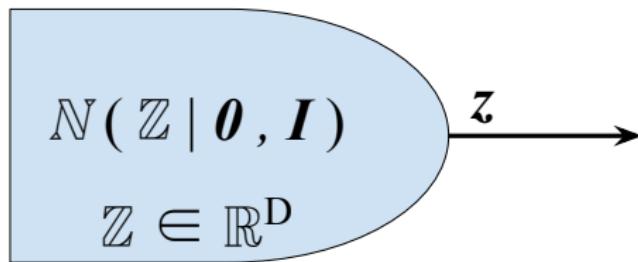


Figure: Sampling base distribution

# Generative Modeling Tasks: Generation

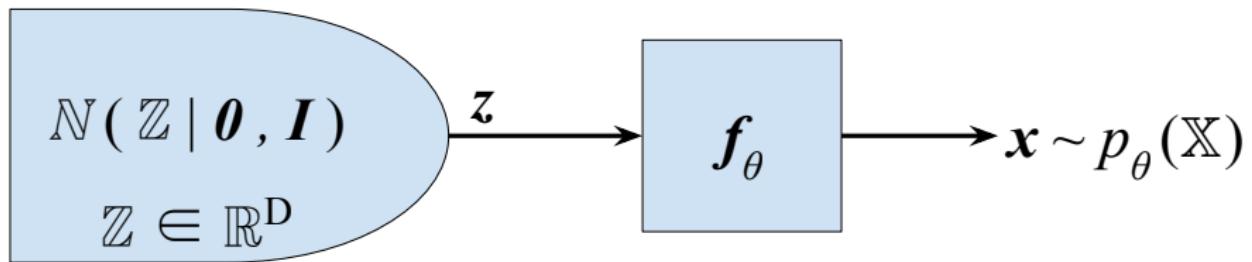


Figure: Passing generated sample through forward transformation

# Generative Modeling Tasks: Density Estimation

$\boldsymbol{x}$

Figure: Input observed vectore  $\boldsymbol{x}$

# Generative Modeling Tasks: Density Estimation



Figure: Corresponding latent vector  $z$

# Generative Modeling Tasks: Density Estimation

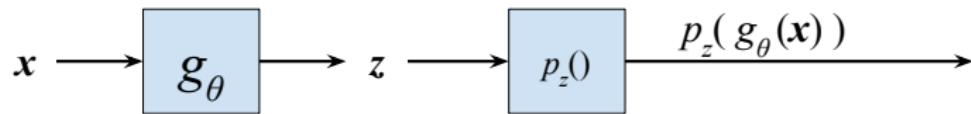


Figure: Corresponding latent vector  $z$

# Generative Modeling Tasks: Density Estimation

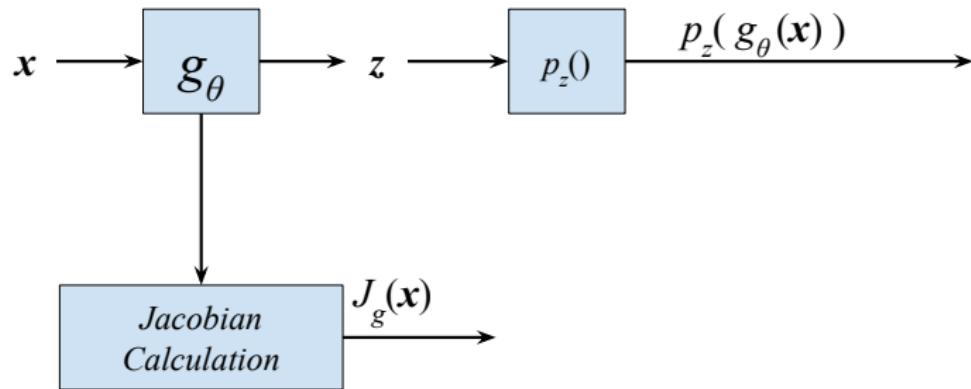


Figure: Jacobian calculations

# Generative Modeling Tasks: Density Estimation

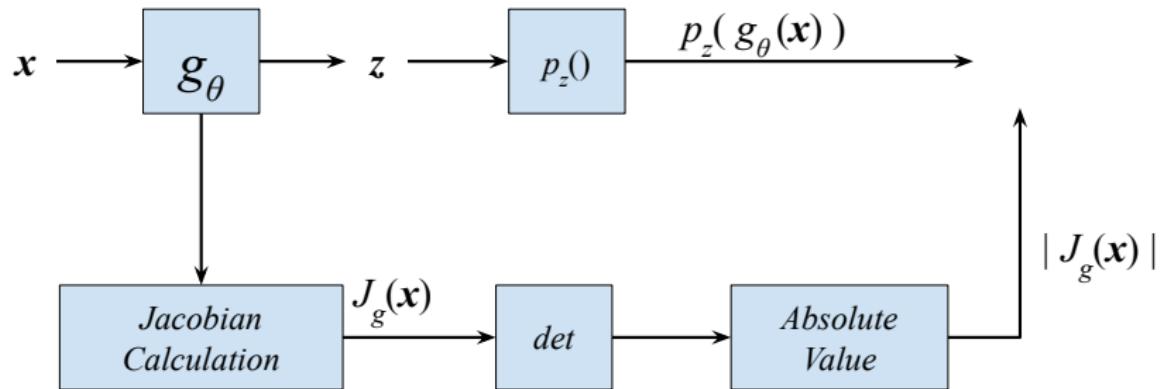


Figure: Calculating the normalization value

# Generative Modeling Tasks: Density Estimation

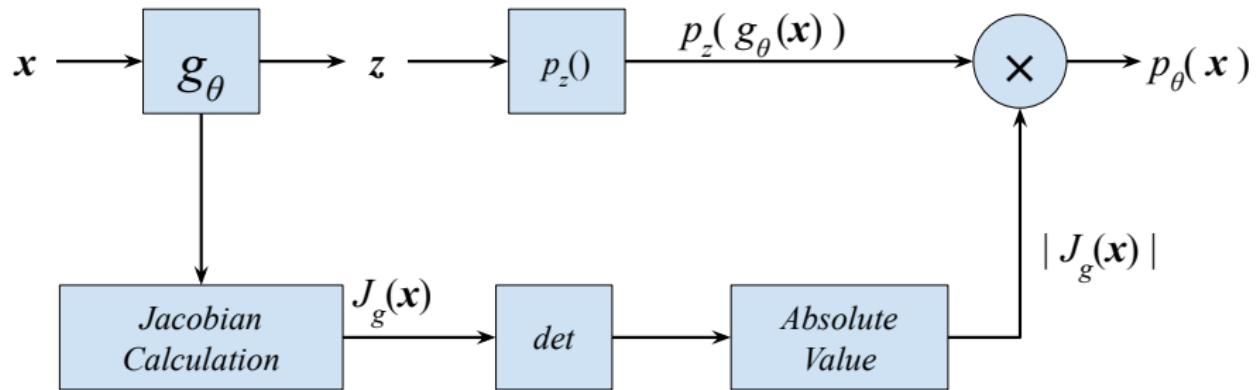
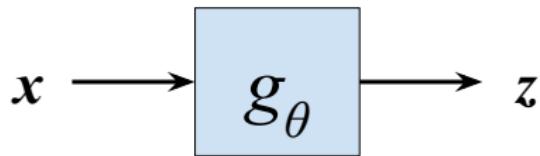


Figure: Density estimation

# Generative Modeling Tasks: Latent Vector



**Figure:** Sample of complex distribution using normalizing flow

## Latent Representation

The latent variable representation corresponding to  $x$  is easily found by inverse transformation  $g_\theta(\cdot)$  but you can't find semantically meaningful latent directions (in contrast to VAE).

# Normalizing Flow Limitations

## Continuous Random Variable

The change of variable technique is only applicable to continuous random variables, thus  $\mathbb{X}$  and  $\mathbb{Z}$  must be continuous random variables.

## Base Distribution

The base distribution should have the following properties:

- *Easy-to-sample*
- Straightforward to compute the likelihood  $p(\mathbf{z})$

## Equal Dimension

Invertibility of transformation necessitates  $\mathbb{X}$  and  $\mathbb{Z}$  to be of the same dimension.  
Thus:

- $\mathbf{J}_g(\mathbf{x}) \in \mathbb{R}^{D \times D}$
- The determinant calculations is of  $\mathcal{O}(D^3)$
- ☞ The transformation  $\mathbf{f}_\theta$  should be designed carefully to reduce the calculations

## Section 4

Learning

## Maximum Likelihood Estimation

To minimize the KL divergence between  $p_{\text{data}}(\mathbb{X})$  and  $p_{\theta}(\mathbb{X})$ , we should maximize the training data likelihood with respect to model parameters  $\theta$  as:

$$\theta^* = \operatorname{argmax}_{\theta} \mathbb{E}_{\mathbf{x} \sim p_{\text{data}}} [\log p_{\theta}(\mathbf{x})] \stackrel{(a)}{\approx} \operatorname{argmax}_{\theta} \frac{1}{N} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

where (a) follows from Monte Carlo estimation.

# Data Likelihood

## Using Inverse Transformation

To train the Normalizing Flow model, we need to calculate the data likelihood as:

$$\log p_{\theta}(\mathbf{x}_i) = \log \left( p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \left| \det \mathbf{J}_{\mathbf{g}}(\mathbf{x}_i) \right| \right)$$

$$= \log \left( p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \left| \det \prod_{l=1}^L \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right| \right) \quad \# \text{Using chain rule}$$

$$= \log \left( p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \left| \prod_{l=1}^L \det \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right| \right) \quad \# \text{Determinant properties}$$

# Data Likelihood

## Using Inverse Transformation (cont.)

$$\begin{aligned}\log p_{\theta}(\mathbf{x}_i) &= \log \left( p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \left| \prod_{l=1}^L \det \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right| \right) \\ &= \log \left( p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) \prod_{l=1}^L \left| \det \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right| \right) \quad \# \text{Absolute value properties} \\ &= \log p\left(\mathbf{g}_{\theta}(\mathbf{x}_i)\right) + \sum_{l=1}^L \log \left| \det \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right| \quad \# \text{Logarithm properties}\end{aligned}$$

# Data Likelihood

## Using Inverse Transformation (Cont.)

Thus to calculate the data likelihood, we can use the following:

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \log p(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{x}_i)) + \sum_{l=1}^L \log \left| \det \mathbf{J}_{g_l}(\mathbf{z}_i(l)) \right|, \quad \mathbf{z}_i(l) = \mathbf{f}_l(\mathbf{z}_i(l-1))$$

## Using Forward Transformation

Similar to the inverse transformation, we can calculate the likelihood using forward transformation which leads to the following problem:

$$\boldsymbol{\theta}^* = \operatorname{argmax}_{\boldsymbol{\theta}} \frac{1}{N} \sum_{i=1}^N \log p(\mathbf{g}_{\boldsymbol{\theta}}(\mathbf{x}_i)) - \sum_{l=1}^L \log \left| \det \mathbf{J}_{f_l}(\mathbf{z}_i(l-1)) \right|$$

## Section 5

Nonlinear Independent Component Estimation

# Additive Coupling Layer [2]

## Transformation

- $f_\theta : \mathbf{z} \rightarrow \mathbf{x}$ : Assume  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^D$ , then the forward transformation is defined as:
  - Select a random integer  $d$  such that  $1 \leq d < D$
  - $\mathbf{x}_{1:d} = \mathbf{z}_{1:d}$
  - $\mathbf{x}_{d+1:D} = \mathbf{z}_{d+1:D} + \mathbf{m}_\theta(\mathbf{z}_{1:d})$where  $\mathbf{m}_\theta(\cdot)$  can be any neural network architecture.
- $g_\theta : \mathbf{x} \rightarrow \mathbf{z}$ : Using the forward transformation, one can easily show that:
  - $\mathbf{z}_{1:d} = \mathbf{x}_{1:d}$
  - $\mathbf{z}_{d+1:D} = \mathbf{x}_{d+1:D} - \mathbf{m}_\theta(\mathbf{z}_{1:d}) = \mathbf{x}_{d+1:D} - \mathbf{m}_\theta(\mathbf{x}_{1:d})$

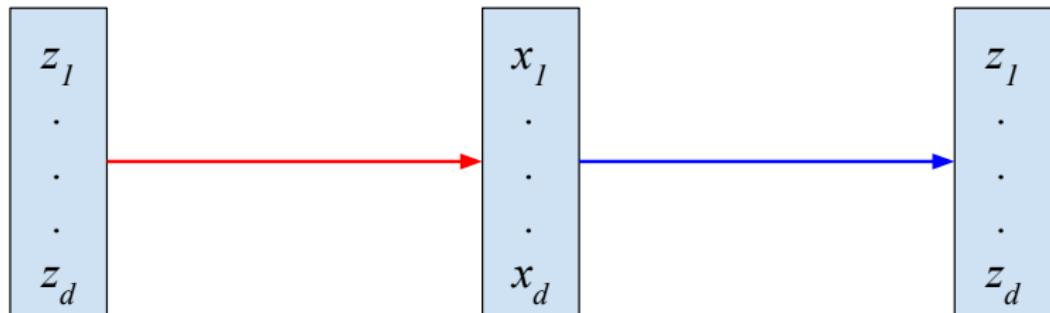
## Additive Coupling Layer [2]



$$f_{\theta}$$

Figure: Forward transformation of upper part

## Additive Coupling Layer [2]



$$f_{\theta}$$

$$g_{\theta}$$

Figure: Inverse transformation of upper part

## Additive Coupling Layer [2]

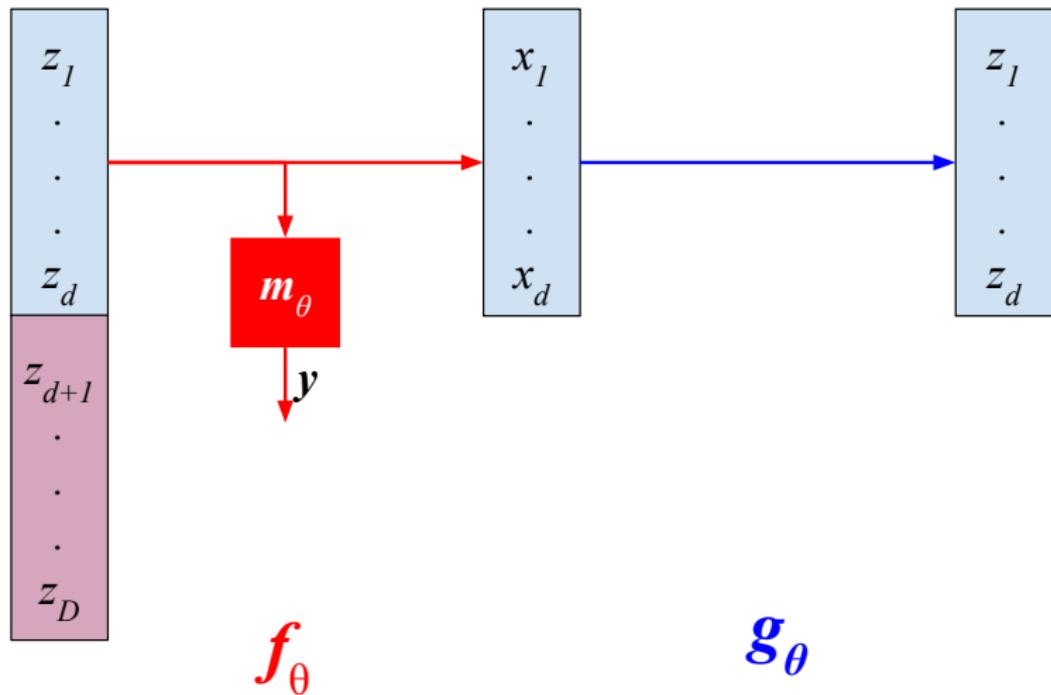


Figure: Forward transformation of lower part

## Additive Coupling Layer [2]

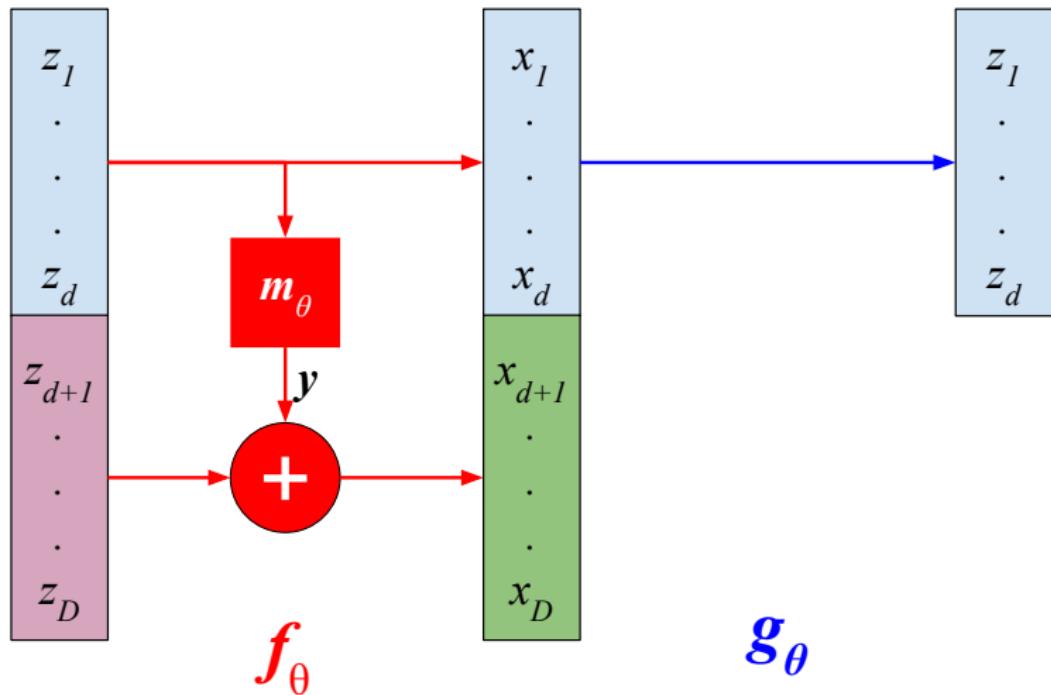


Figure: Forward transformation of lower part

# Additive Coupling Layer [2]

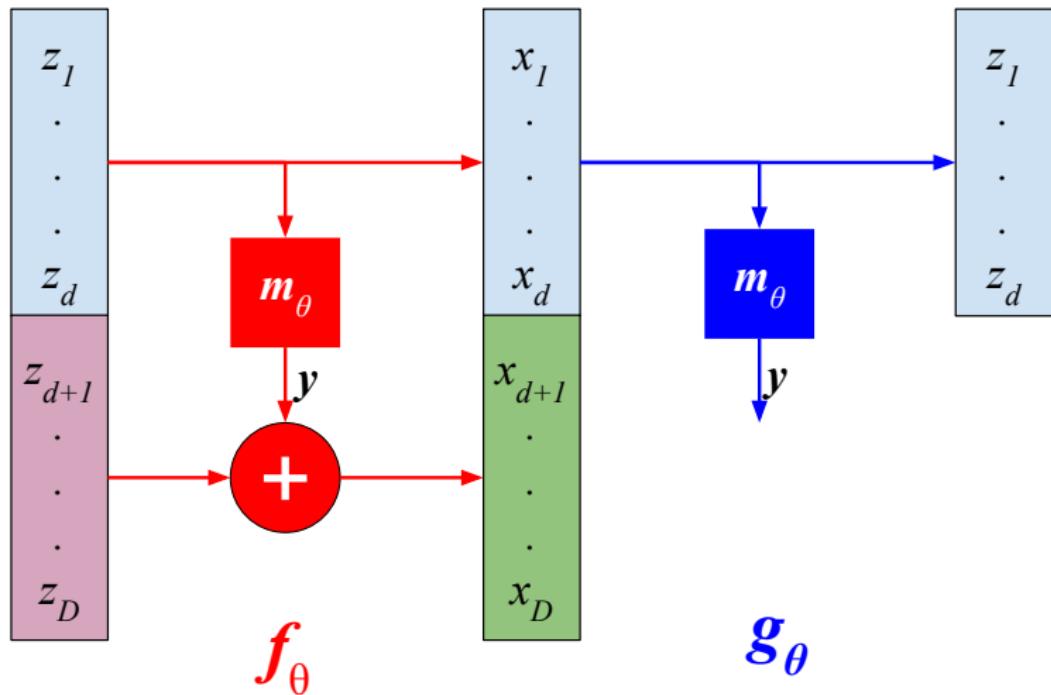


Figure: Inverse transformation of lower part

# Additive Coupling Layer [2]

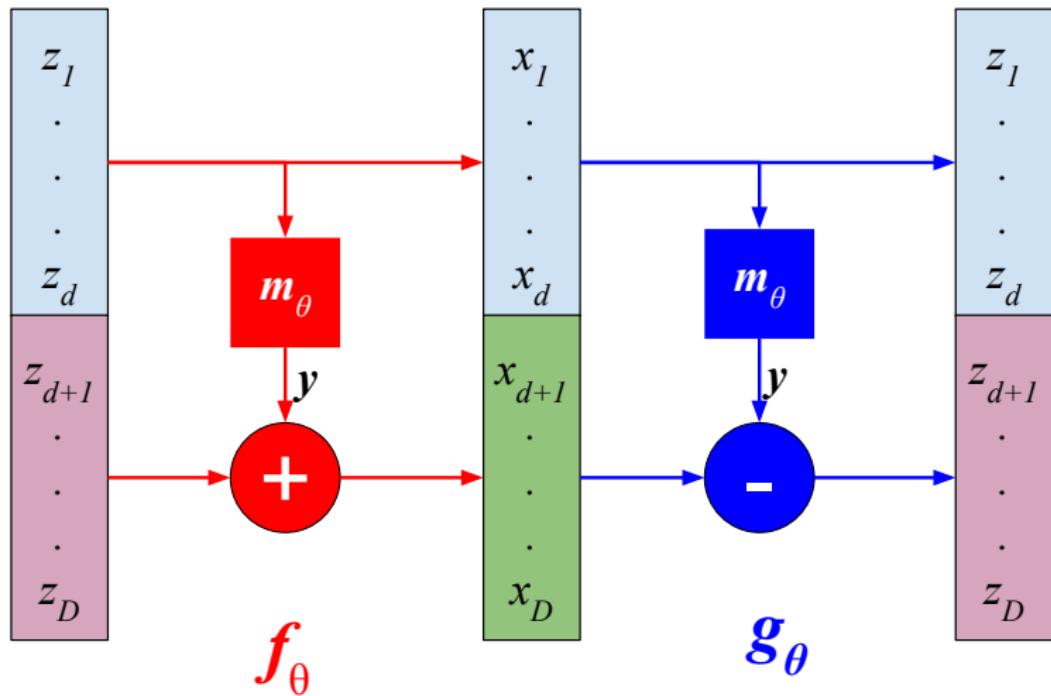


Figure: Inverse transformation of lower part

# Additive Coupling Layer [2]

## Jacobian

The Jacobian for the forward transformation of the additive coupling layer is:

$$\begin{aligned}\mathbf{J}_f &= \det \begin{bmatrix} \left[ \frac{\partial \mathbf{x}_{1:d}}{\partial \mathbf{z}_{1:d}} \right]_{d \times d} & \left[ \frac{\partial \mathbf{x}_{1:d}}{\partial \mathbf{z}_{d+1:D}} \right]_{d \times (n-d)} \\ \left[ \frac{\partial \mathbf{x}_{d+1:D}}{\partial \mathbf{z}_{1:d}} \right]_{(n-d) \times d} & \left[ \frac{\partial \mathbf{x}_{d+1:D}}{\partial \mathbf{z}_{d+1:D}} \right]_{(n-d) \times (n-d)} \end{bmatrix} \\ &= \det \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \frac{\partial \mathbf{m}_\theta(\mathbf{z}_{1:d})}{\partial \mathbf{z}_{1:d}} & \mathbf{I} \end{bmatrix} \\ &\stackrel{(a)}{=} 1\end{aligned}$$

where (a) follows the fact that for triangular matrices, the determinant is the product of diagonal elements.

- ☞ This value for the determinant can greatly reduce the computations to calculate the log-likelihood.
- ☞ This transformation is volume preserving.

# Rescaling Layer [2]

## Transformation

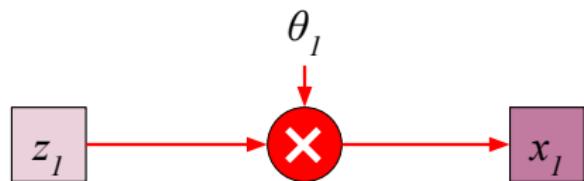
- $f_\theta : \mathbf{z} \rightarrow \mathbf{x}$ : Assume  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^D$ , then the forward transformation is defined as:

$$x_i = \theta_i z_i , \quad i = 1, \dots, D$$

- $g_\theta \mathbf{x} \rightarrow \mathbf{z}$ : Using the forward transformation, one can easily show that:

$$z_i = \frac{x_i}{\theta_i} , \quad i = 1, \dots, D$$

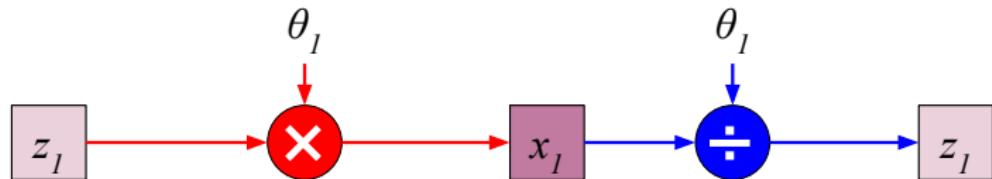
## Rescaling Layer [2]



$$f_{\theta}$$

Figure: Forward transformation of first element

## Rescaling Layer [2]

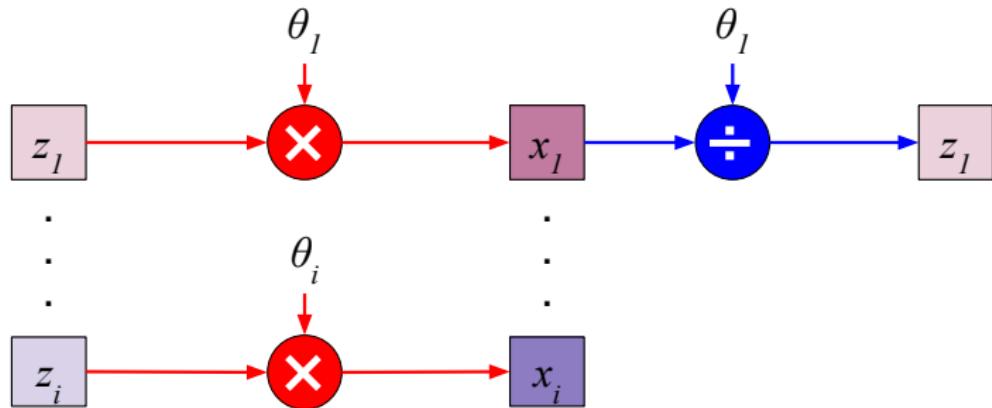


$$f_{\theta}$$

$$g_{\theta}$$

Figure: Inverse transformation of first element

## Rescaling Layer [2]

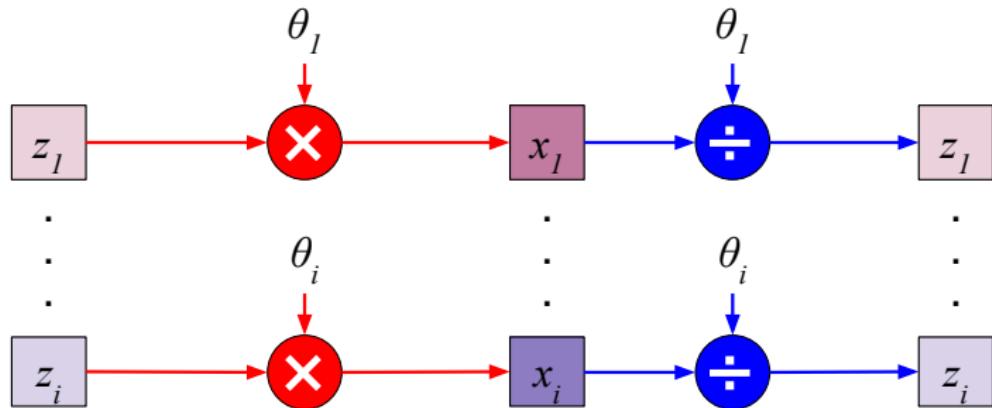


$$f_{\theta}$$

$$g_{\theta}$$

Figure: Forward transformation of  $i$ -th element

## Rescaling Layer [2]



$$f_{\theta}$$

$$g_{\theta}$$

Figure: Inverse transformation of  $i$ -th element

## Rescaling Layer [2]

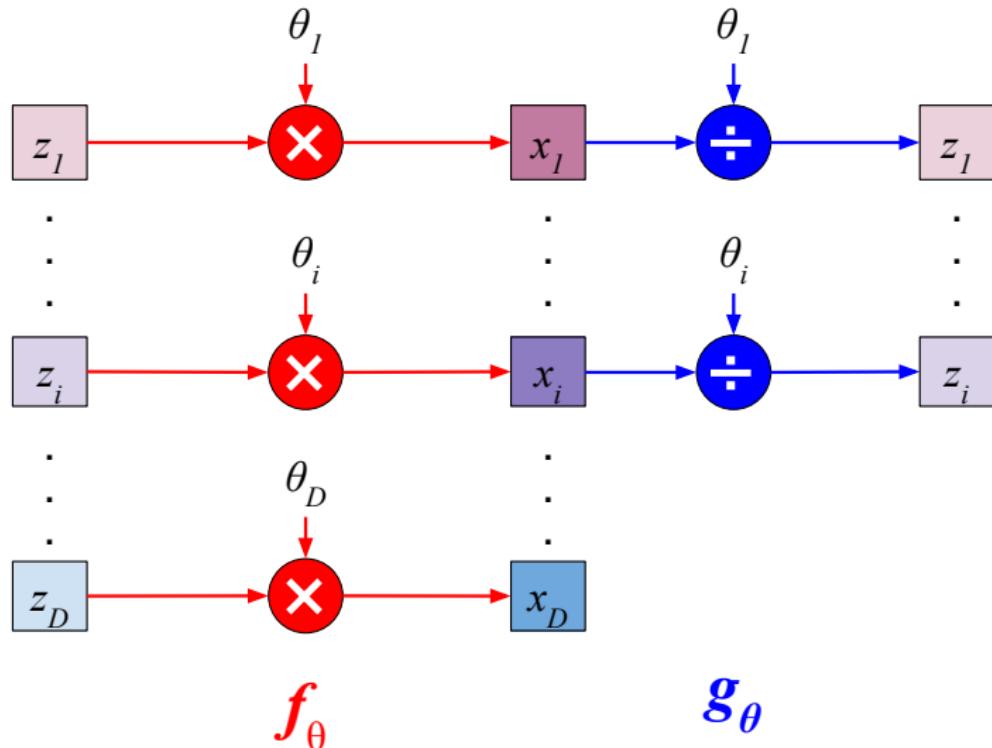


Figure: Forward transformation of last element

## Rescaling Layer [2]

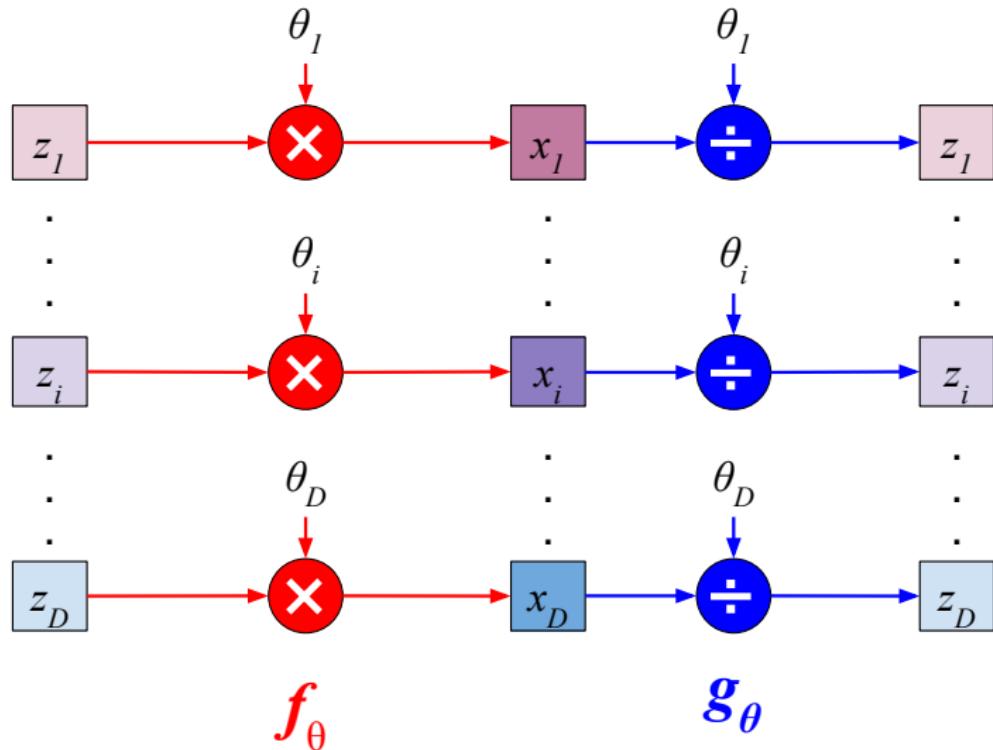


Figure: Inverse transformation of last element

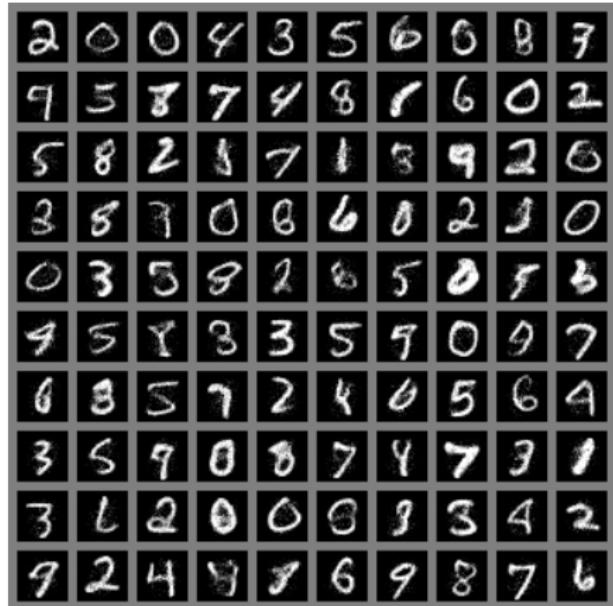
# Rescaling Layer [2]

## Jacobian

The Jacobian for the forward transformation of the rescaling layer is:

$$\begin{aligned}\mathbf{J}_f &= \det \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \frac{\partial x_1}{\partial z_2} & \cdots & \frac{\partial x_1}{\partial z_D} \\ \frac{\partial x_2}{\partial z_1} & \frac{\partial x_2}{\partial z_2} & \cdots & \frac{\partial x_2}{\partial z_D} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial x_D}{\partial z_1} & \frac{\partial x_D}{\partial z_2} & \cdots & \frac{\partial x_D}{\partial z_D} \end{bmatrix} \\ &= \det \begin{bmatrix} \theta_1 & 0 & \cdots & 0 \\ 0 & \theta_2 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & \theta_D \end{bmatrix} \\ &= \prod_{i=1}^D \theta_i\end{aligned}$$

# NICE Performance [2]



(a) MNIST dataset



(b) TFD dataset

**Figure:** NICE performance. The architecture comprises four additive coupling layers and one rescaling layer

# NICE Performance [2]



(a) SVHN dataset



(b) CIFAR10 dataset

**Figure:** NICE performance. The architecture comprises four additive coupling layers and one rescaling layer

# Image Inpainting Using NICE [2]

## Image Inpainting

Assume:

- $p_\theta(\mathbf{x})$  to be a trained model over the MNIST dataset using NICE.
  - You are given a sample  $\mathbf{x}$  where the  $\mathbf{x}_o$  is the observed part (clean) and  $\mathbf{x}_h$  is hidden part (corrupted).
- ☞ Find an estimate to the corrupted part?

## Maximum a Posteriori Estimation

To find MAP approximation to the corrupted part, we have:

$$\begin{aligned}\mathbf{x}_{h, \text{map}} &= \underset{\mathbf{x}_h}{\operatorname{argmax}} p_\theta(\mathbf{x}_h | \mathbf{x}_o) \stackrel{(a)}{\equiv} \underset{\mathbf{x}_h}{\operatorname{argmax}} \frac{p_\theta(\mathbf{x}_h, \mathbf{x}_o)}{p_\theta(\mathbf{x}_o)} \\ &\equiv \underset{\mathbf{x}_h}{\operatorname{argmax}} p_\theta(\mathbf{x}_h, \mathbf{x}_o)\end{aligned}$$

# Image Inpainting Using NICE [2]

## Maximum a Posteriori Estimation (cont.)

Thus we should solve the following problem:

$$\mathbf{x}_{h, \text{map}} = \operatorname{argmax}_{\mathbf{x}_h} p_{\theta}(\mathbf{x}_h, \mathbf{x}_o)$$

Note:

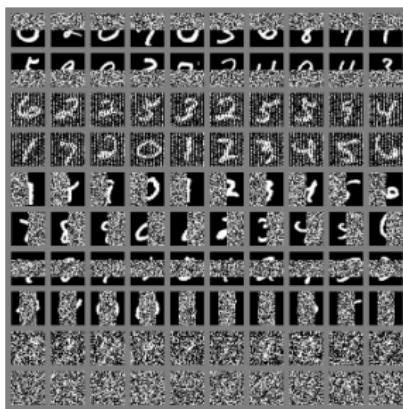
- The model parameters are fixed (trained).
- We should maximize  $p_{\theta}(\mathbf{x}_h, \mathbf{x}_o)$  w.r.t.  $\mathbf{x}_h$  using gradient descent:

$$\mathbf{x}_{h,k+1} = \mathbf{x}_{h,k} + \frac{\partial p_{\theta}(\mathbf{x}_h, \mathbf{x}_o)}{\partial \mathbf{x}_h}$$

# NICE Performance [2]



(a) Original sample



(b) corrupted Sample



(c) MAP estimation

**Figure:** NICE inpainting. The corruption in the middle image are: top rows, bottom rows, odd pixels, even pixels, left side, right side, middle vertically, middle horizontally, 75% random and 90% random.

## Section 6

Real-valued Non-volume Preserving Transformations

# Coupling Layer [3]

## Transformation

- $f_\theta : \mathbf{z} \rightarrow \mathbf{x}$ : Assume  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^D$ , then the forward transformation is defined as:
  - Select a random integer  $d$  such that  $1 \leq d < D$
  - $\mathbf{x}_{1:d} = \mathbf{z}_{1:d}$
  - $\mathbf{x}_{d+1:D} = \mathbf{z}_{d+1:D} \odot \exp(\mathbf{s}_\theta(\mathbf{z}_{1:d})) + \mathbf{t}_\theta(\mathbf{z}_{1:d})$where  $\mathbf{s}_\theta(\cdot)$  and  $\mathbf{t}_\theta(\cdot)$  can be any neural network architecture.
- $g_\theta : \mathbf{x} \rightarrow \mathbf{z}$ : Using the forward transformation, one can easily show that:
  - $\mathbf{z}_{1:d} = \mathbf{x}_{1:d}$
  -

$$\begin{aligned}\mathbf{z}_{d+1:D} &= (\mathbf{x}_{d+1:D} - \mathbf{t}_\theta(\mathbf{z}_{1:d})) \odot \exp(-\mathbf{s}_\theta(\mathbf{z}_{1:d})) \\ &= (\mathbf{x}_{d+1:D} - \mathbf{t}_\theta(\mathbf{x}_{1:d})) \odot \exp(-\mathbf{s}_\theta(\mathbf{x}_{1:d}))\end{aligned}$$

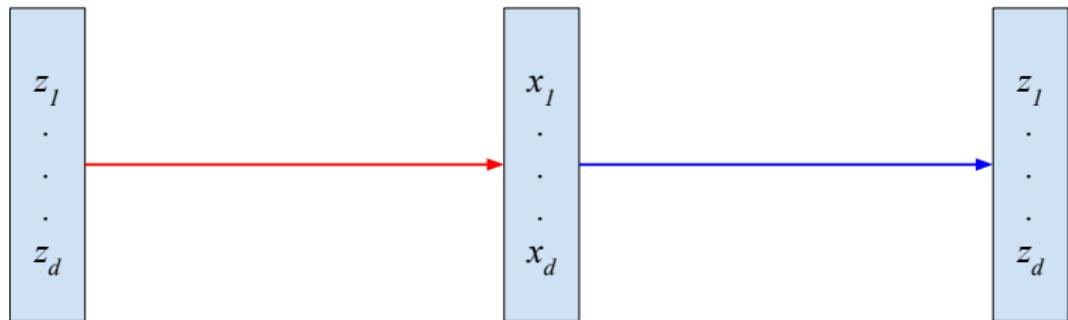
## Coupling Layer [3]



$$\mathbf{f}_\theta$$

Figure: Forward transformation of upper part

## Coupling Layer [3]



$$f_\theta$$

$$g_\theta$$

Figure: Inverse transformation of upper part

# Coupling Layer [3]

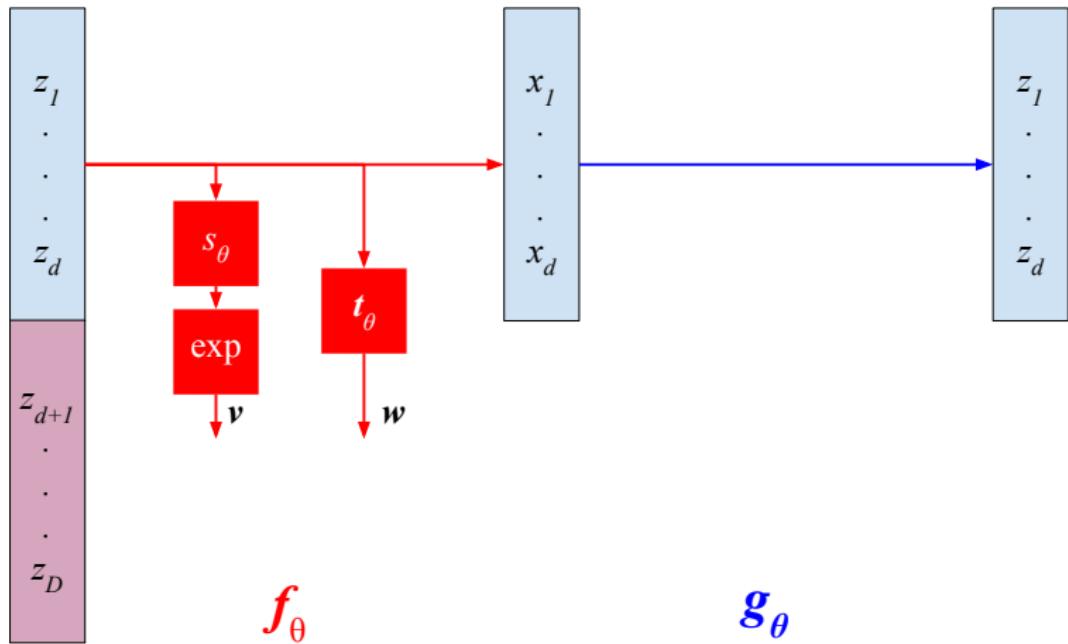


Figure: Forward transformation of lower part

# Coupling Layer [3]

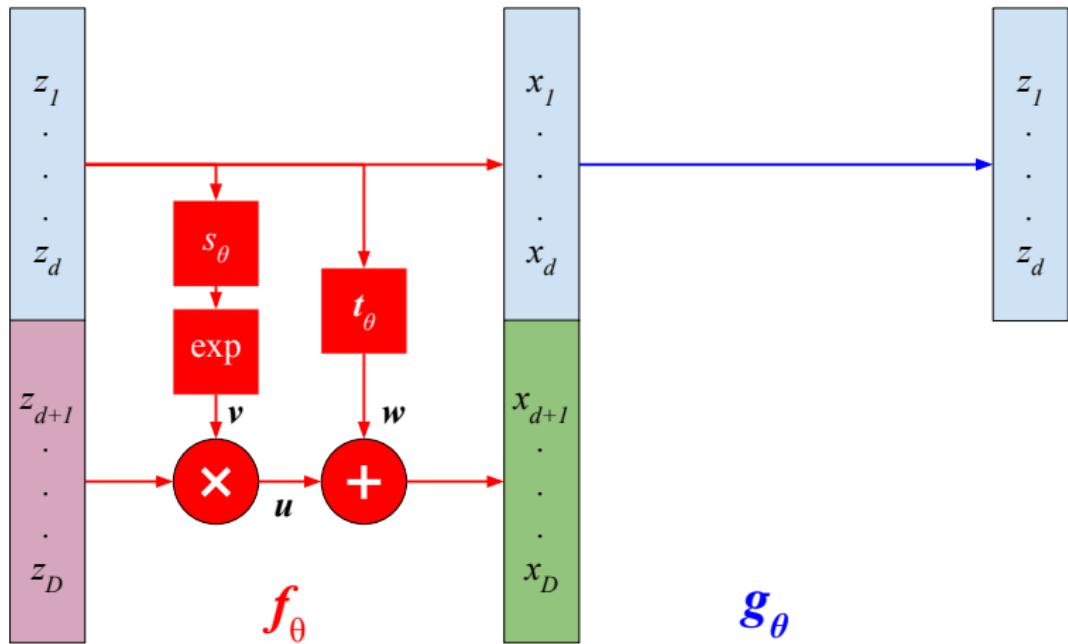


Figure: Forward transformation of lower part

# Coupling Layer [3]

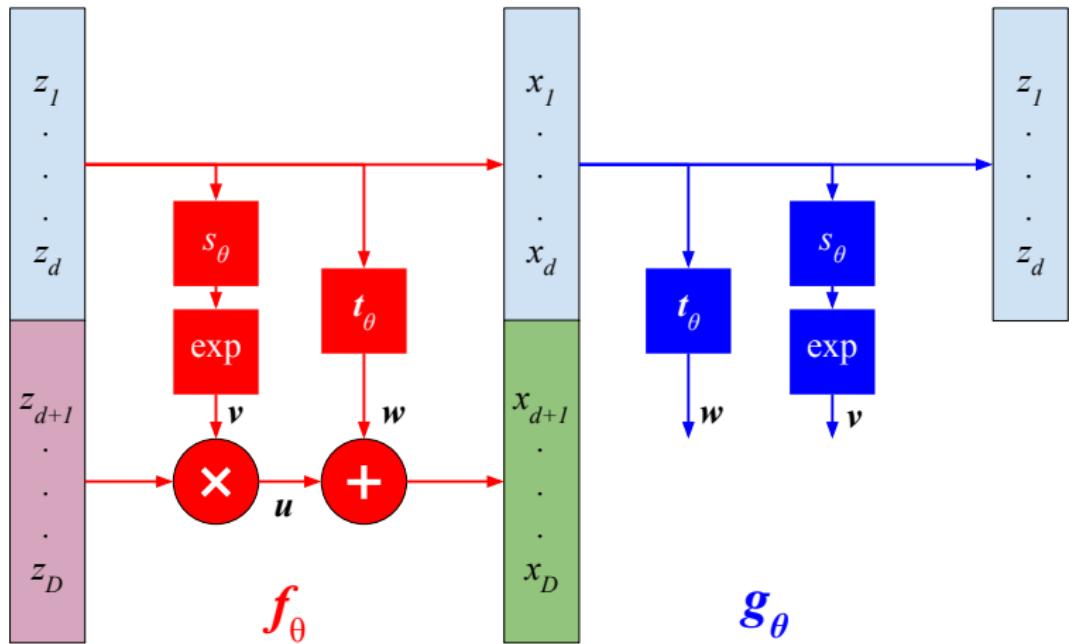


Figure: Inverse transformation of lower part

# Coupling Layer [3]

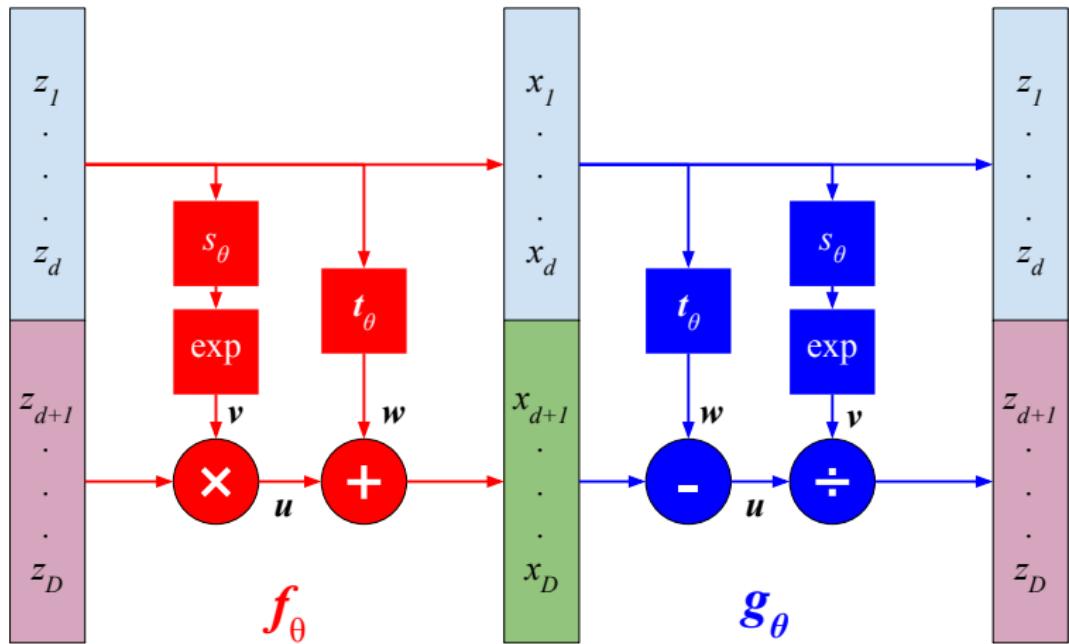


Figure: Inverse transformation of lower part

# Coupling Layer [3]

## Jacobian

The Jacobian for the forward transformation of the coupling layer is:

$$\begin{aligned}\mathbf{J}_f &= \det \begin{bmatrix} \left[ \frac{\partial \mathbf{x}_{1:d}}{\partial \mathbf{z}_{1:d}} \right]_{d \times d} & \left[ \frac{\partial \mathbf{x}_{1:d}}{\partial \mathbf{z}_{d+1:D}} \right]_{d \times (n-d)} \\ \left[ \frac{\partial \mathbf{x}_{d+1:D}}{\partial \mathbf{z}_{1:d}} \right]_{(n-d) \times d} & \left[ \frac{\partial \mathbf{x}_{d+1:D}}{\partial \mathbf{z}_{d+1:D}} \right]_{(n-d) \times (n-d)} \end{bmatrix} \\ &= \det \begin{bmatrix} \mathbf{I} & \mathbf{0} \\ \frac{\partial \mathbf{m}_\theta(\mathbf{z}_{1:d})}{\partial \mathbf{z}_{1:d}} & \begin{bmatrix} \exp(\mathbf{s}_{\theta,1}(\mathbf{x}_{1:d})) & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \exp(\mathbf{s}_{\theta,D-d}(\mathbf{x}_{1:d})) \end{bmatrix} \end{bmatrix} \\ &= \prod_{i=1}^{D-d} \exp(\mathbf{s}_{\theta,i}(\mathbf{x}_{1:d})) = \exp \left( \sum_{i=1}^{D-d} \mathbf{s}_{\theta,i}(\mathbf{x}_{1:d}) \right)\end{aligned}$$

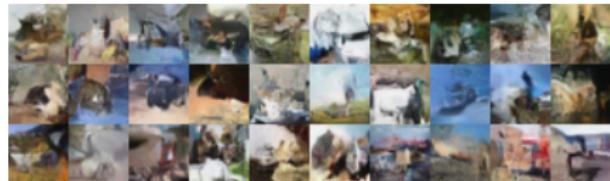
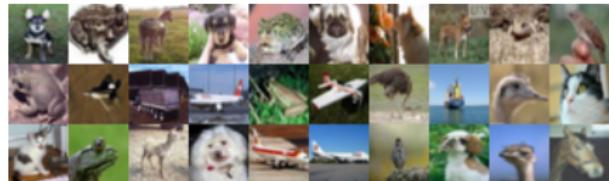
- ☞ The likelihood becomes more complicated in comparison to NICE.
- ☞ The transformation is not volume preserving anymore which justifies Real-NVP.

# Real NVP Performance [3]

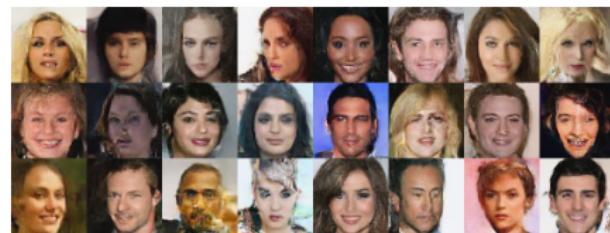
## Other Note on Real NVP

- Combining several coupling layers
- Masked convolution for partitioning
- Multiscale architectures
- Batch-normalization

# Real NVP Generation [3]



(a) Cifar10 dataset



(b) CelebA dataset

Figure: Left column: Dataset examples, right column: generated samples

# Real NVP Generation [3]



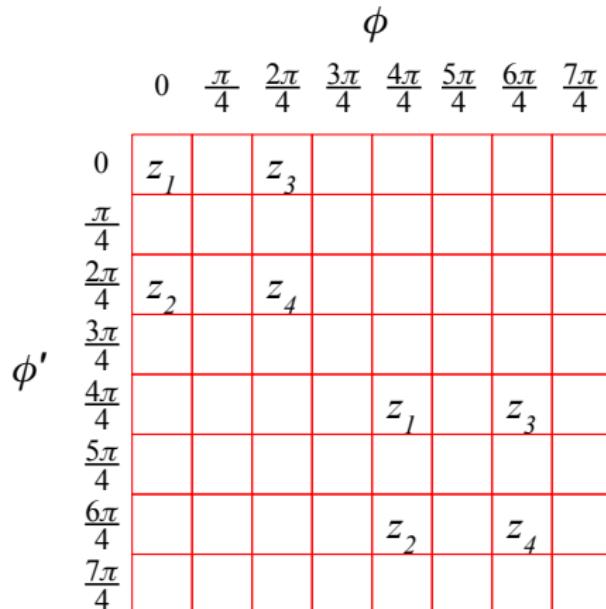
(a) ImageNet  $32 \times 32$  dataset



(b) Imagenet  $64 \times 64$  dataset

Figure: Left column: Dataset examples, right column: generated samples

# Real NVP Latent Variable Interpolation [3]



## Interpolation Formula

$$\mathbf{z} = \cos \phi (\cos \phi' \mathbf{z}_1 + \sin \phi' \mathbf{z}_2) + \sin \phi (\cos \phi' \mathbf{z}_3 + \sin \phi' \mathbf{z}_4)$$

# Real NVP Latent Variable Interpolation [3]

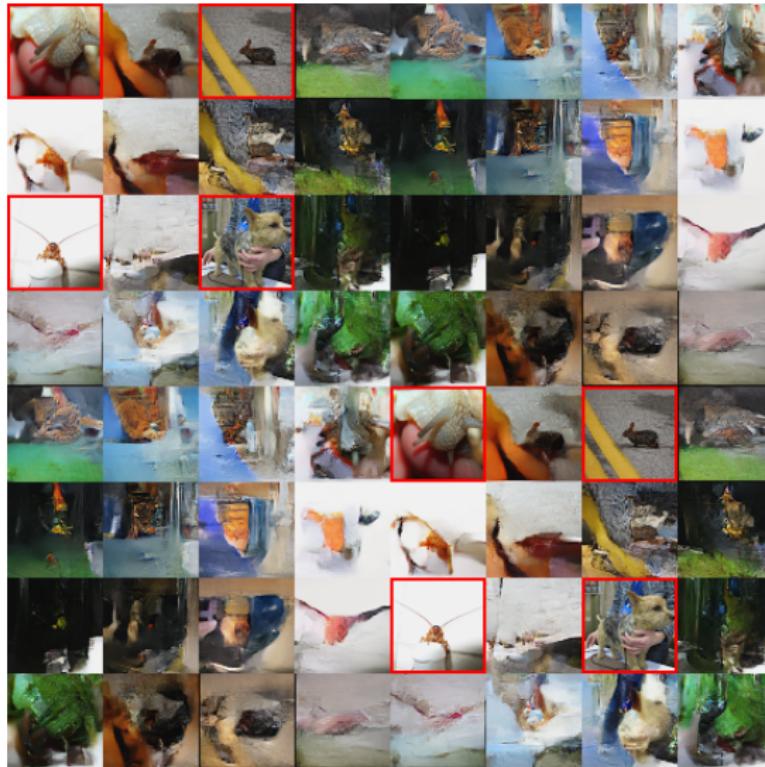


Figure: ImageNet  $64 \times 64$  dataset

# Real NVP Latent Variable Interpolation [3]



Figure: CelebA dataset

# Real NVP Latent Variable Interpolation [3]



Figure: LSUN (bedroom category) dataset

# Real NVP Latent Variable Interpolation [3]



Figure: LSUN (church outdoor category) dataset

# Real NVP Latent Variable Interpolation [3]



Figure: LSUN (tower category) dataset

## Section 7

### Masked Autoregressive Flows

# Masked Autoregressive Flows [4]

## Transformation

- $f_\theta : \mathbf{z} \rightarrow \mathbf{x}$ : Assume  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^D$ , then the forward transformation is defined as:
  - $x_i = z_i \times \exp(\alpha_i(\mathbf{x}_{<i})) + \mu_i(\mathbf{x}_{<i})$
  - $i = 1$ :  $x_1 = z_1 \times \exp(\alpha_1(\mathbf{x}_{<1})) + \mu_1(\mathbf{x}_{<1}) = z_1 \times \exp(\alpha) + \mu$
  - $i = 2$ :  $x_2 = z_2 \times \exp(\alpha_2(\mathbf{x}_{<2})) + \mu_2(\mathbf{x}_{<2}) = z_2 \times \exp(\alpha_2(x_1)) + \mu_2(x_1)$
  - $\vdots$

where  $\alpha_i(\cdot)$  and  $\mu_i(\cdot)$  can be any neural network architecture with single output.

- $g_\theta : \mathbf{x} \rightarrow \mathbf{z}$ : Using the forward transformation, one can easily show that:

$$z_i = \frac{x_i - \mu_i(\mathbf{x}_{<i})}{\exp(\alpha_i(\mathbf{x}_{<i}))}$$

# Masked Autoregressive Flows [4]

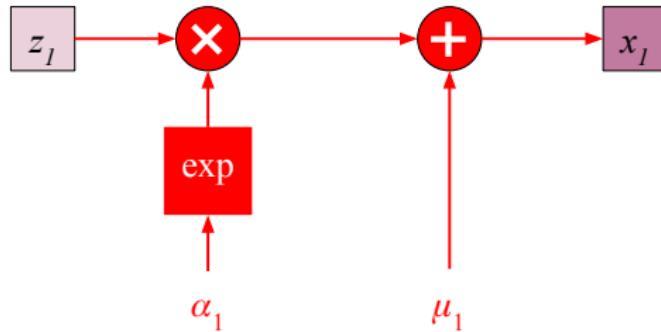
$$z_l$$

$$\alpha_1 \quad \mu_1$$

$$f_\theta$$

Figure: Forward transformation of first element

# Masked Autoregressive Flows [4]



$$f_{\theta}$$

Figure: Forward transformation of first element

# Masked Autoregressive Flows [4]

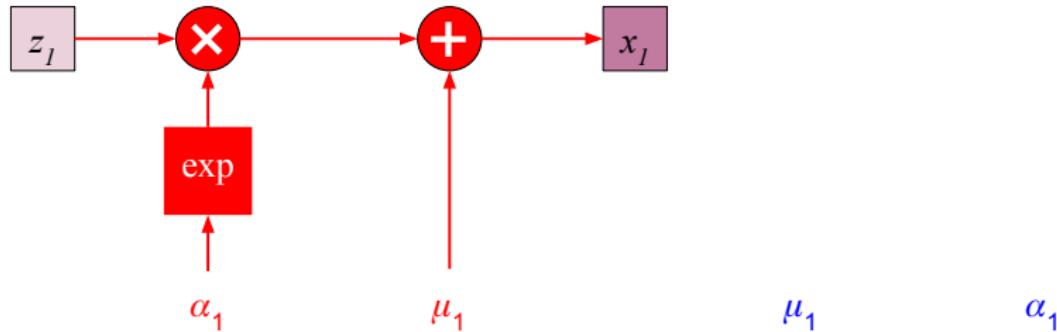
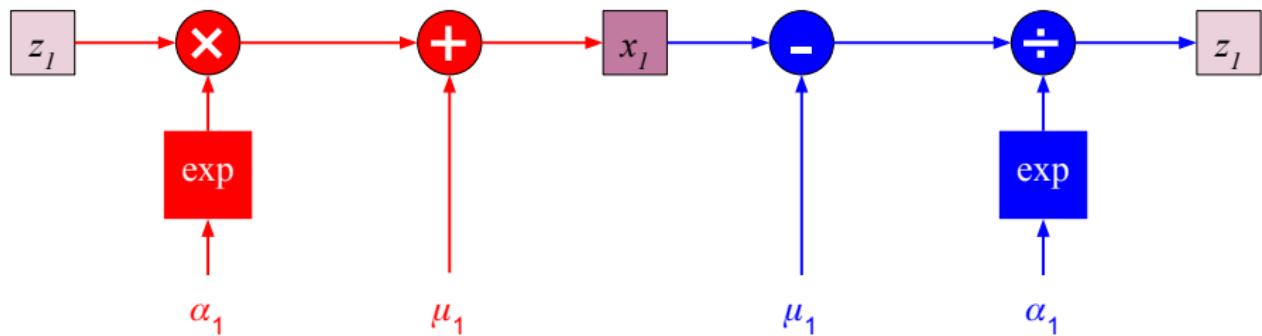


Figure: Inverse transformation of first element

# Masked Autoregressive Flows [4]



$$f_{\theta}$$

$$g_{\theta}$$

Figure: Inverse transformation of first element

# Masked Autoregressive Flows [4]

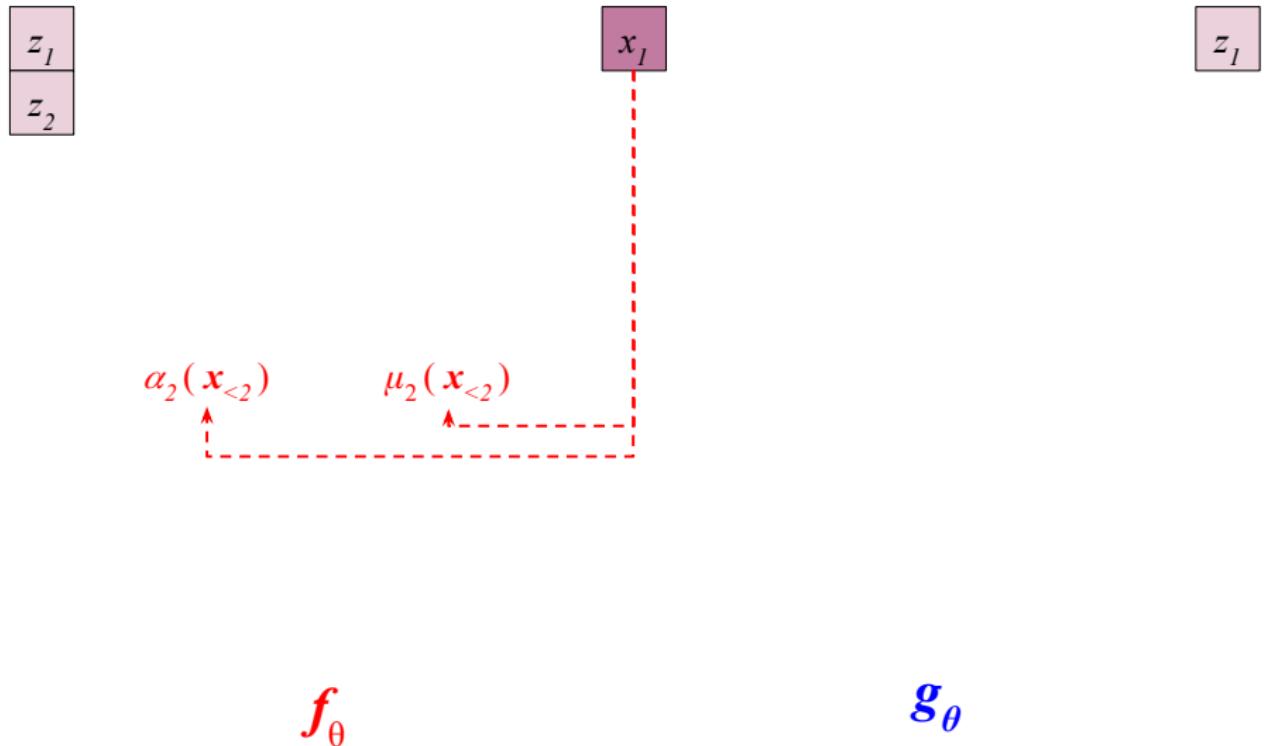
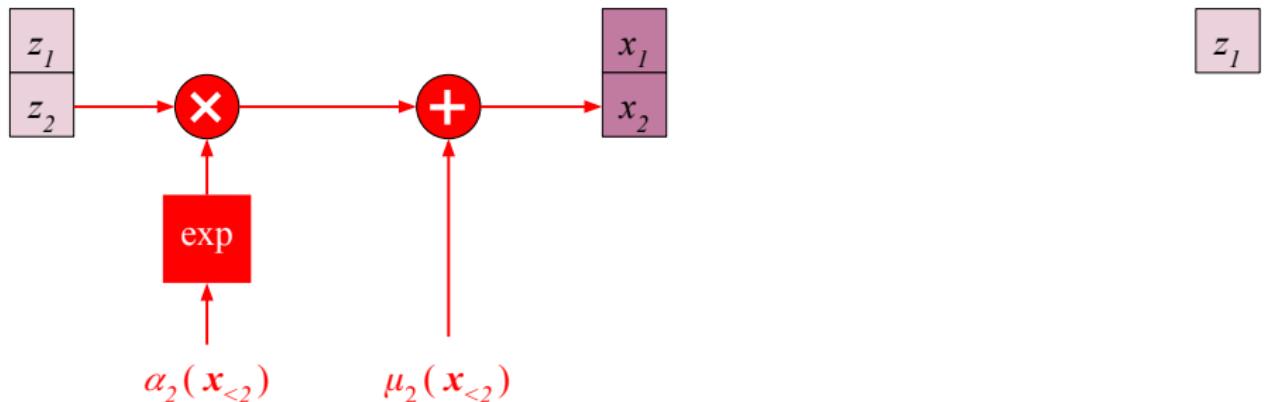


Figure: Forward transformation of second element

# Masked Autoregressive Flows [4]

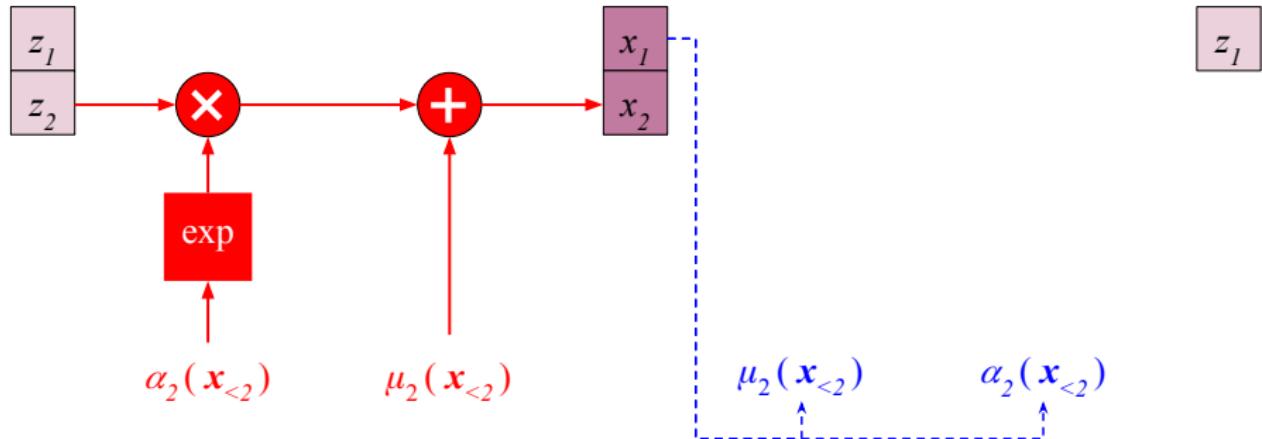


$$f_\theta$$

$$g_\theta$$

Figure: Forward transformation of second element

# Masked Autoregressive Flows [4]



$$f_\theta$$

$$g_\theta$$

Figure: Inverse transformation of second element

# Masked Autoregressive Flows [4]

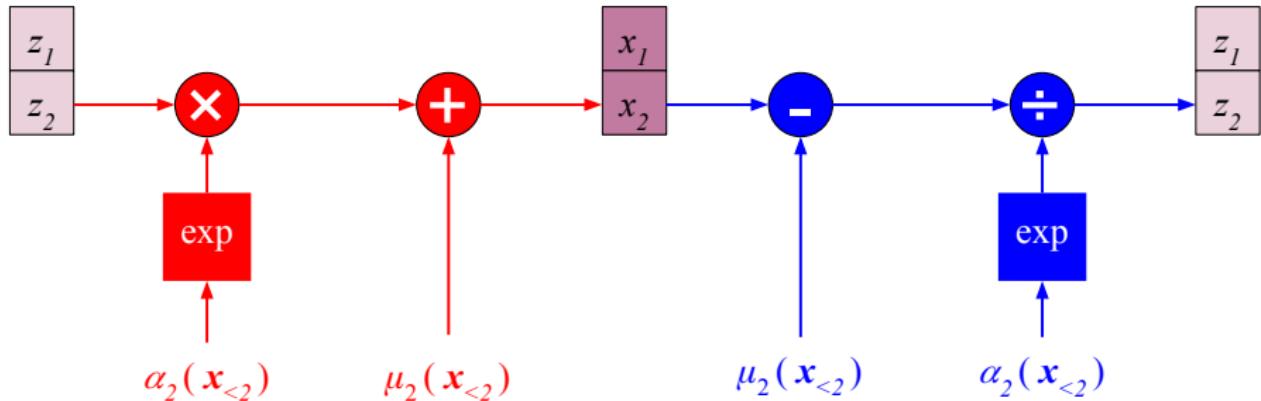


Figure: Inverse transformation of second element

# Masked Autoregressive Flows [4]

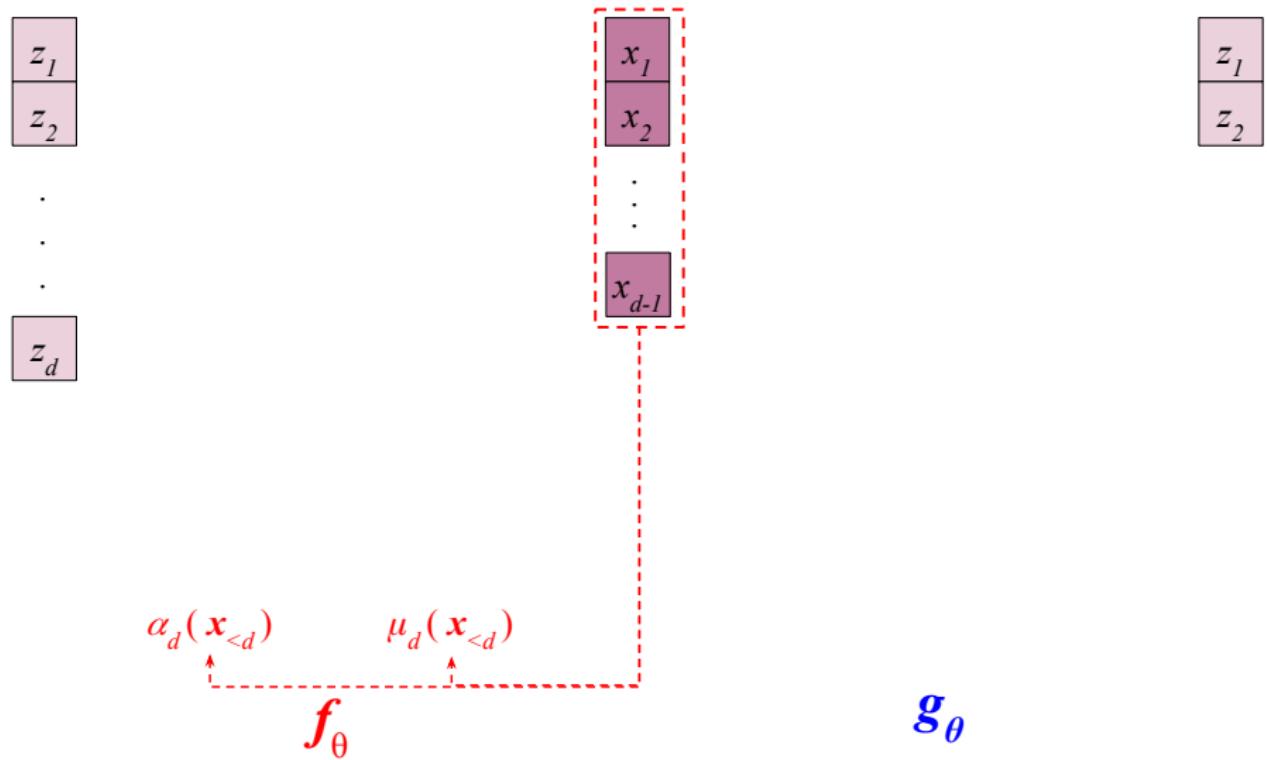


Figure: Forward transformation of  $d$ -th part

# Masked Autoregressive Flows [4]

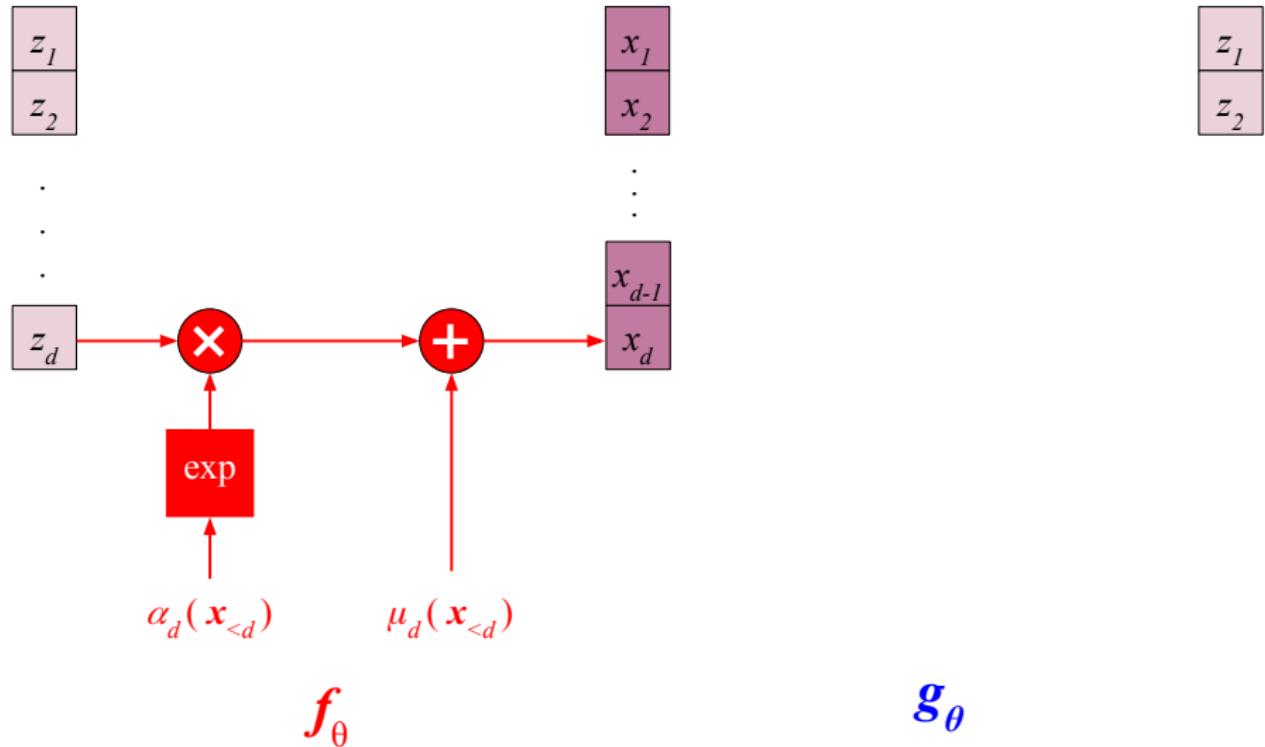


Figure: Forward transformation of  $d$ -th part

# Masked Autoregressive Flows [4]

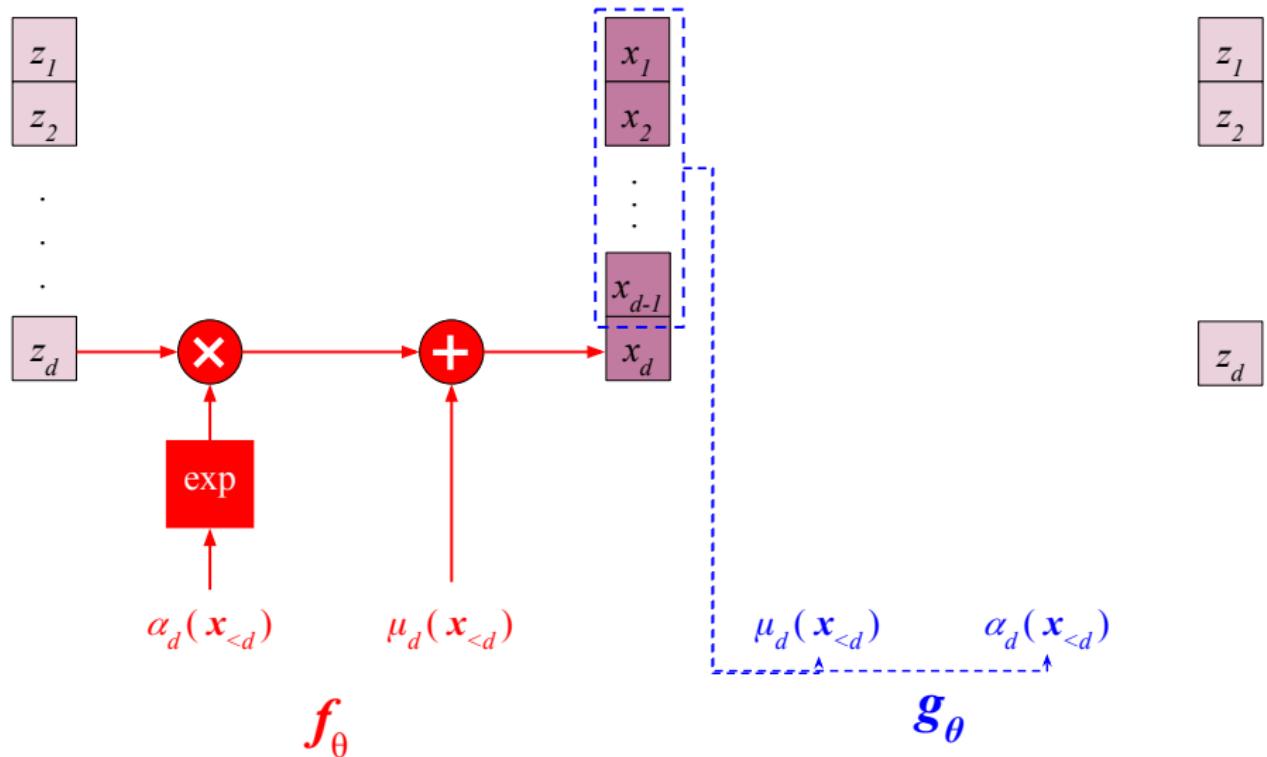


Figure: Inverse transformation of  $d$ -th part

# Masked Autoregressive Flows [4]

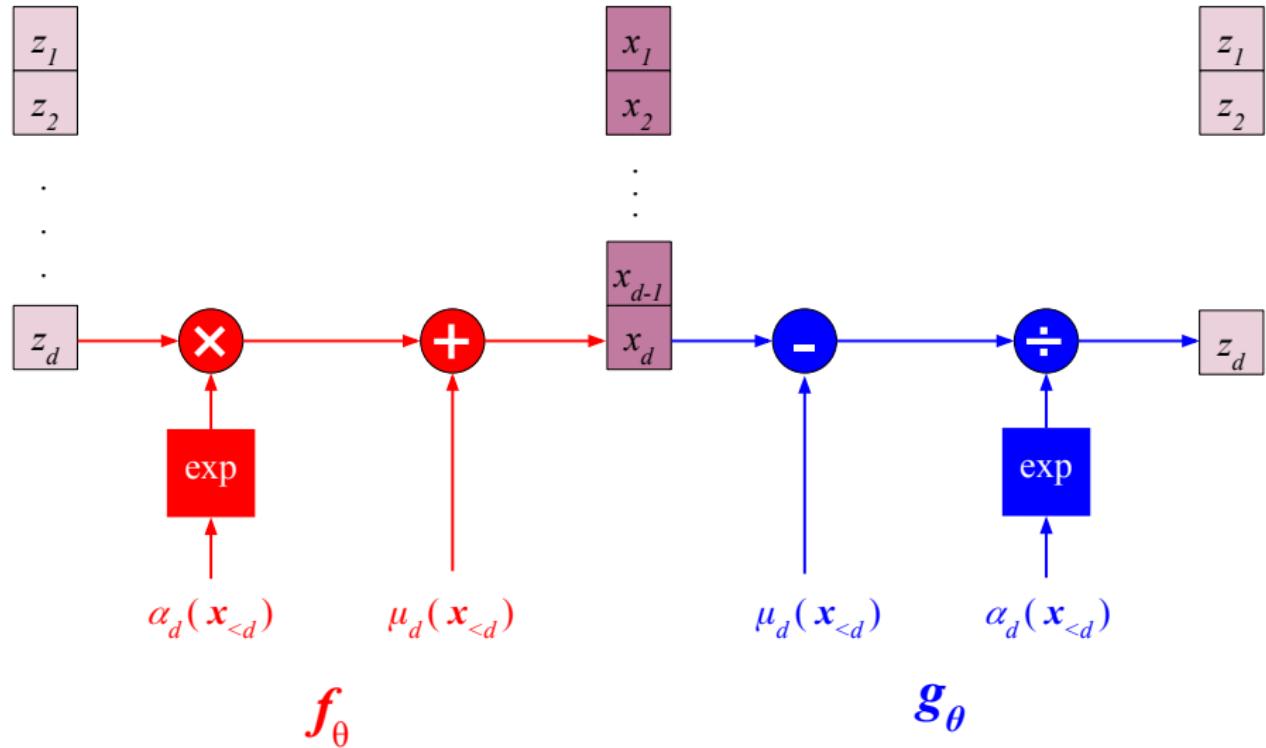


Figure: Inverse transformation of  $d$ -th part

# Masked Autoregressive Flow [4]

## Jacobian

The Jacobian for the forward transformation of MAF is:

$$\begin{aligned}\log |\det \mathbf{J}_f| &= \log \left| \det \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \cdots & \frac{\partial x_1}{\partial z_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_D}{\partial z_1} & \cdots & \frac{\partial x_D}{\partial z_D} \end{bmatrix} \right| \\ &= \log \left| \det \begin{bmatrix} \exp(\alpha_1(\mathbf{x}_{<1})) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial x_D}{\partial z_1} & \cdots & \exp(\alpha_D(\mathbf{x}_{$$

- ☞ The transformation is not volume preserving similar to Real-NVP.

# Masked Autoregressive Flow [4]

## Sampling

Sampling from an MAF can be done through a sequential procedure:

- Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- $x_1 = z_1 \exp(\alpha_1) + \mu_1$ , Calculate  $\begin{cases} \mu_2 \triangleq \mu_2(\mathbf{x}_{<2}) \\ \alpha_2 \triangleq \alpha_2(\mathbf{x}_{<2}) \end{cases}$
- $x_2 = z_2 \exp(\alpha_2) + \mu_2$ , Calculate  $\begin{cases} \mu_3 \triangleq \mu_3(\mathbf{x}_{<3}) \\ \alpha_3 \triangleq \alpha_3(\mathbf{x}_{<3}) \end{cases}$
- $\vdots$
- $x_D = z_D \exp(\alpha_D) + \mu_D$

## Sequential Sampling

The sampling is sequential in MAF as we expect from the autoregressive nature.

# Masked Autoregressive Flow [4]

## Likelihood Calculation

Probability calculations can be done in parallel as:

- Given  $\mathbf{x}$  vector, calculate the following vectors (elements can be calculated in parallel):

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1(\mathbf{x}_{<1}) \\ \mu_2(\mathbf{x}_{<2}) \\ \vdots \\ \mu_D(\mathbf{x}_{<D}) \end{bmatrix}, \quad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1(\mathbf{x}_{<1}) \\ \alpha_2(\mathbf{x}_{<2}) \\ \vdots \\ \alpha_D(\mathbf{x}_{<D}) \end{bmatrix}$$

- Calculate  $\mathbf{z} = \frac{\mathbf{x} - \boldsymbol{\mu}}{\exp(\boldsymbol{\alpha})}$  (all the operations are assumed element-wise)
- Calculate probability using  $p_{\theta}(\mathbf{x}) = p(\mathbf{z}) \left| \prod_{d=1}^D \exp(\alpha_d(\mathbf{x}_{<d})) \right|^{-1}$

## Parallel Likelihood

The likelihood calculation and thus learning which needs transformation Jacobian is in parallel for MAF.

# Conditional Modeling [4]

## Conditional Modeling

Assume we are interested in modeling the conditional distribution  $p_{\theta}(\mathbb{X}|\mathbf{y})$ . Then we can use our previous discussion about unconditional distribution while pay attention to:

- $\mathbf{y}$  is not a random variable (such as  $\mathbb{X}$  or latent random variable  $\mathbb{Z}$ )
- $\mathbf{y}$  is given. So we must design invertible transformation assuming you have access to  $\mathbf{y}$  in both forward and inverse transformations.

# Conditional Modeling [4]

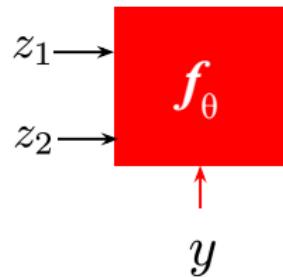


Figure: Forward transformation

## Conditional Modeling [4]

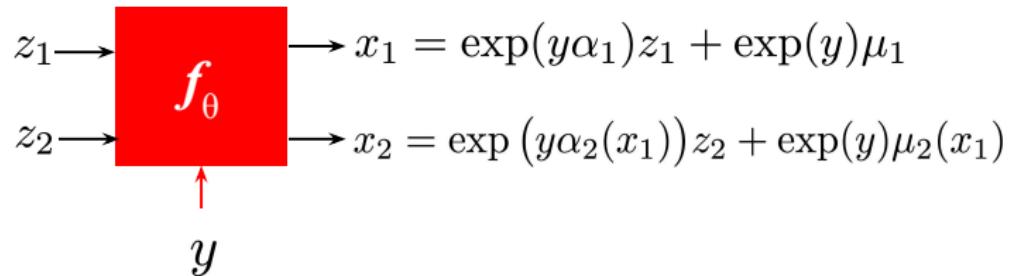


Figure: Forward transformation

## Conditional Modeling [4]

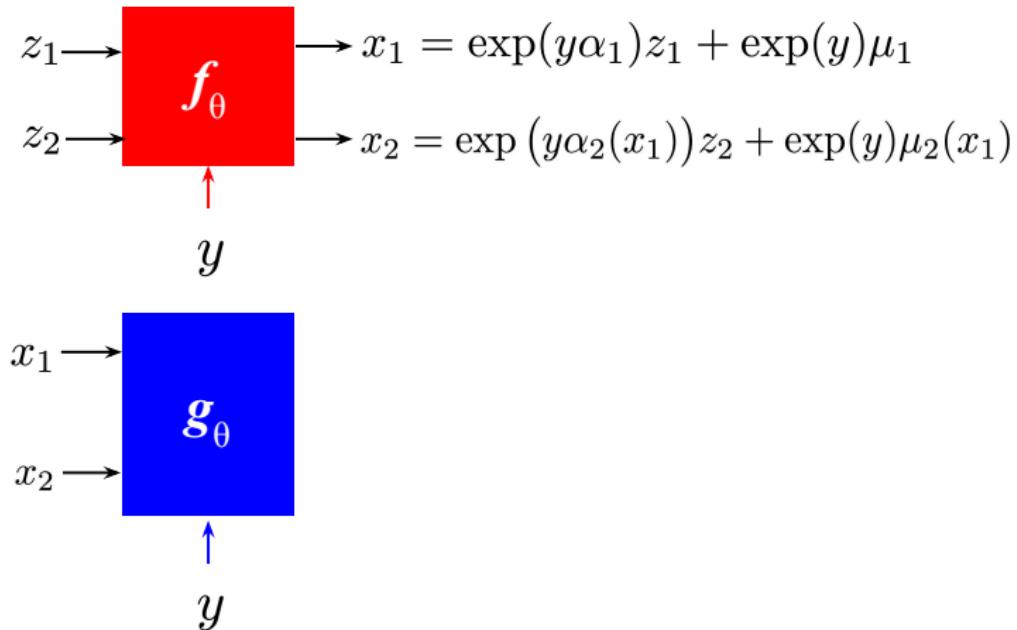


Figure: Inverse transformation

## Conditional Modeling [4]

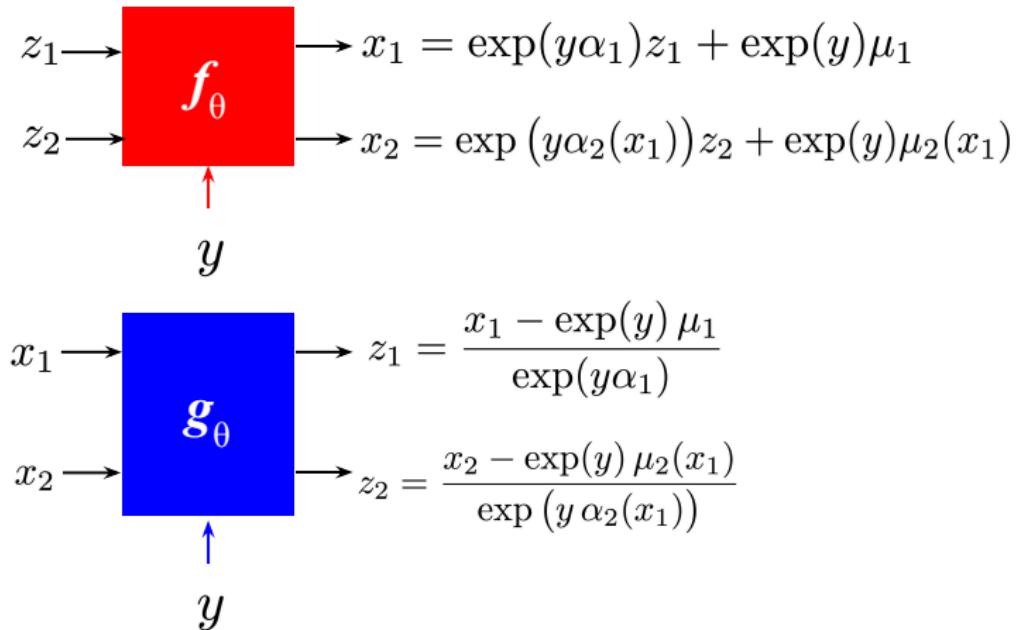


Figure: Inverse transformation

# Class Conditional Generation (source: [4])



(a) Dataset samples



(b) Class-conditional samples

## Section 8

### Inverse Autoregressive Flows

# Inverse Autoregressive Flows [5]

## Forward Transformation

$f_\theta : \mathbf{z} \rightarrow \mathbf{x}$ : Assume  $\mathbf{z}, \mathbf{x} \in \mathbb{R}^D$ , then the forward transformation is defined as:

$$x_d = z_d \times \exp(\alpha_d(\mathbf{z}_{<d})) + \mu_d(\mathbf{z}_{<d})$$

$$\Rightarrow \begin{cases} d = 1 : & x_1 = z_1 \times \exp(\alpha_1(\mathbf{z}_{<1})) + \mu_1(\mathbf{z}_{<1}) = z_1 \times \exp(\alpha_1) + \mu_1 \\ d = 2 : & x_2 = z_2 \times \exp(\alpha_2(\mathbf{z}_{<2})) + \mu_2(\mathbf{z}_{<2}) = z_2 \times \exp(\alpha_2(z_1)) + \mu_2(z_1) \\ \vdots \end{cases}$$

where  $\alpha_d(\cdot)$  and  $\mu_d(\cdot)$  can be any neural network architecture with single output.

# Inverse Autoregressive Flows [5]

## Inverse Transformation

$\mathbf{g}_\theta : \mathbf{x} \rightarrow \mathbf{z}$ : Using the forward transformation, one can easily show that:

$$z_d = \frac{x_d - \mu_d(\mathbf{z}_{<d})}{\exp(\alpha_d(\mathbf{z}_{<d}))}$$

$$\Rightarrow \begin{cases} d = 1 : & z_1 = \frac{x_1 - \mu_1(\mathbf{z}_{<1})}{\exp(\alpha_1(\mathbf{z}_{<d}))} = \frac{x_1 - \mu_1}{\exp(\alpha_1)} \\ d = 2 : & z_2 = \frac{x_2 - \mu_2(\mathbf{z}_{<2})}{\exp(\alpha_2(\mathbf{z}_{<2}))} = \frac{x_1 - \mu_2(z_1)}{\exp(\alpha_2(z_1))} \\ \vdots \end{cases}$$

# Inverse Autoregressive Flows [5]

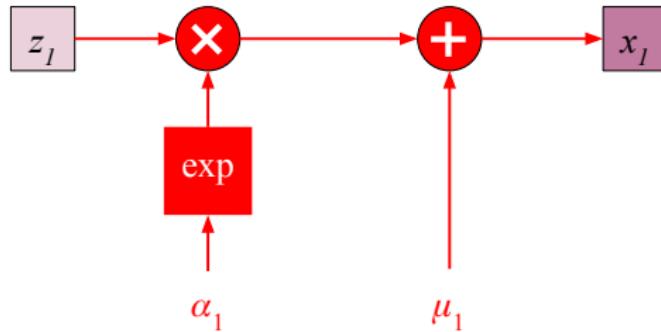
$$z_l$$

$$\alpha_1 \quad \mu_1$$

$$f_\theta$$

Figure: Forward transformation of first element

# Inverse Autoregressive Flows [5]



$f_\theta$

Figure: Forward transformation of first element

# Inverse Autoregressive Flows [5]

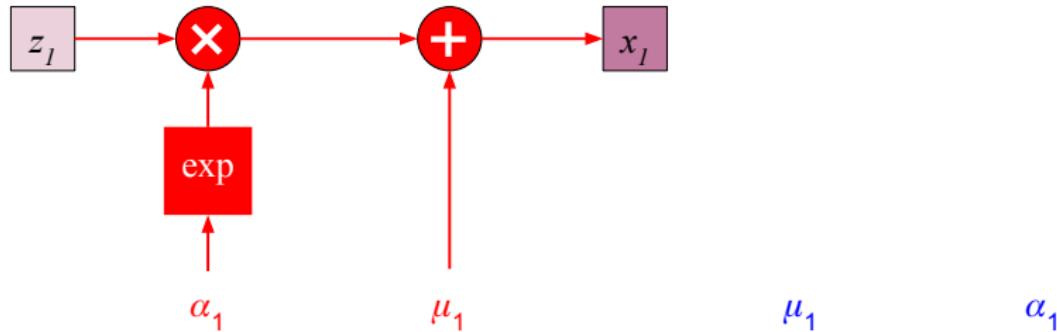
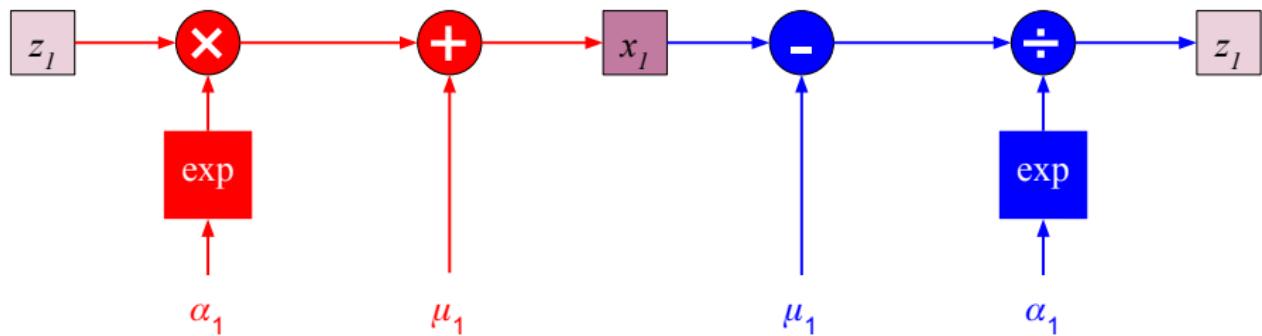


Figure: Inverse transformation of first element

# Inverse Autoregressive Flows [5]



$$f_{\theta}$$

$$g_{\theta}$$

Figure: Inverse transformation of first element

# Inverse Autoregressive Flows [5]

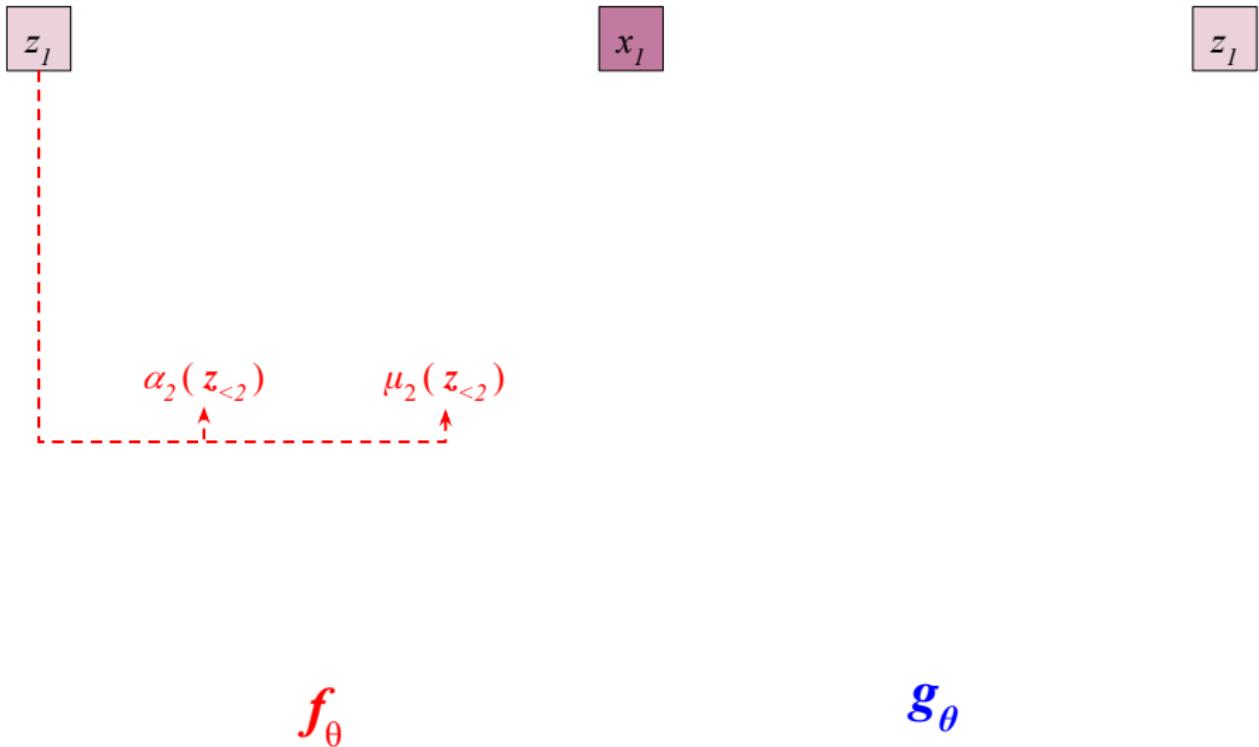
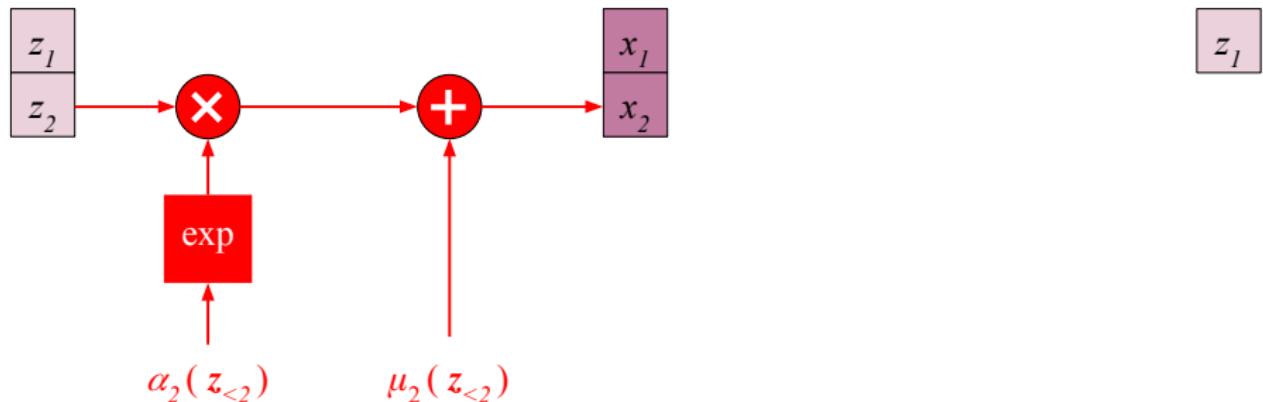


Figure: Forward transformation of second element

# Inverse Autoregressive Flows [5]

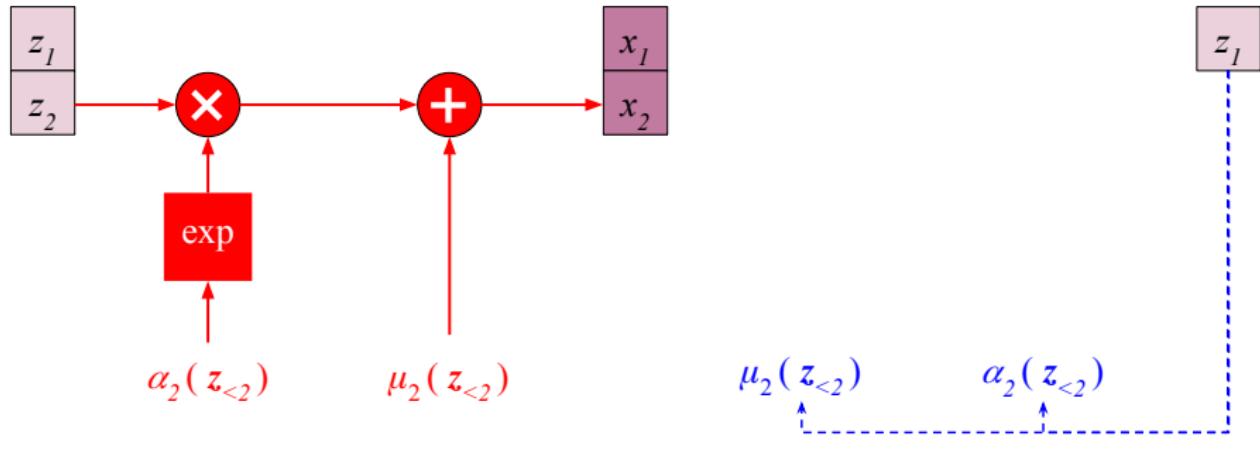


$$f_\theta$$

$$g_\theta$$

Figure: Forward transformation of second element

# Inverse Autoregressive Flows [5]



$$f_\theta$$

$$g_\theta$$

Figure: Inverse transformation of second element

# Inverse Autoregressive Flows [5]

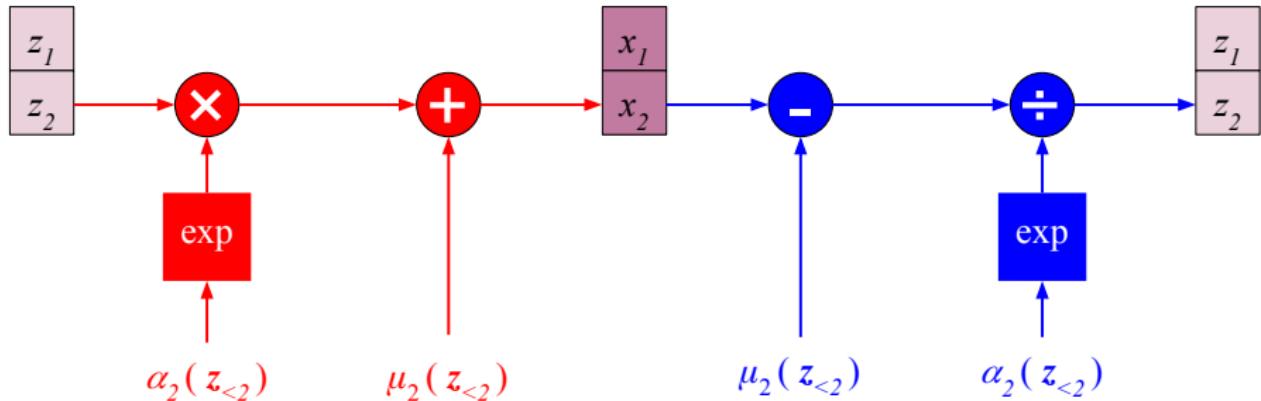


Figure: Inverse transformation of second element

# Inverse Autoregressive Flows [5]

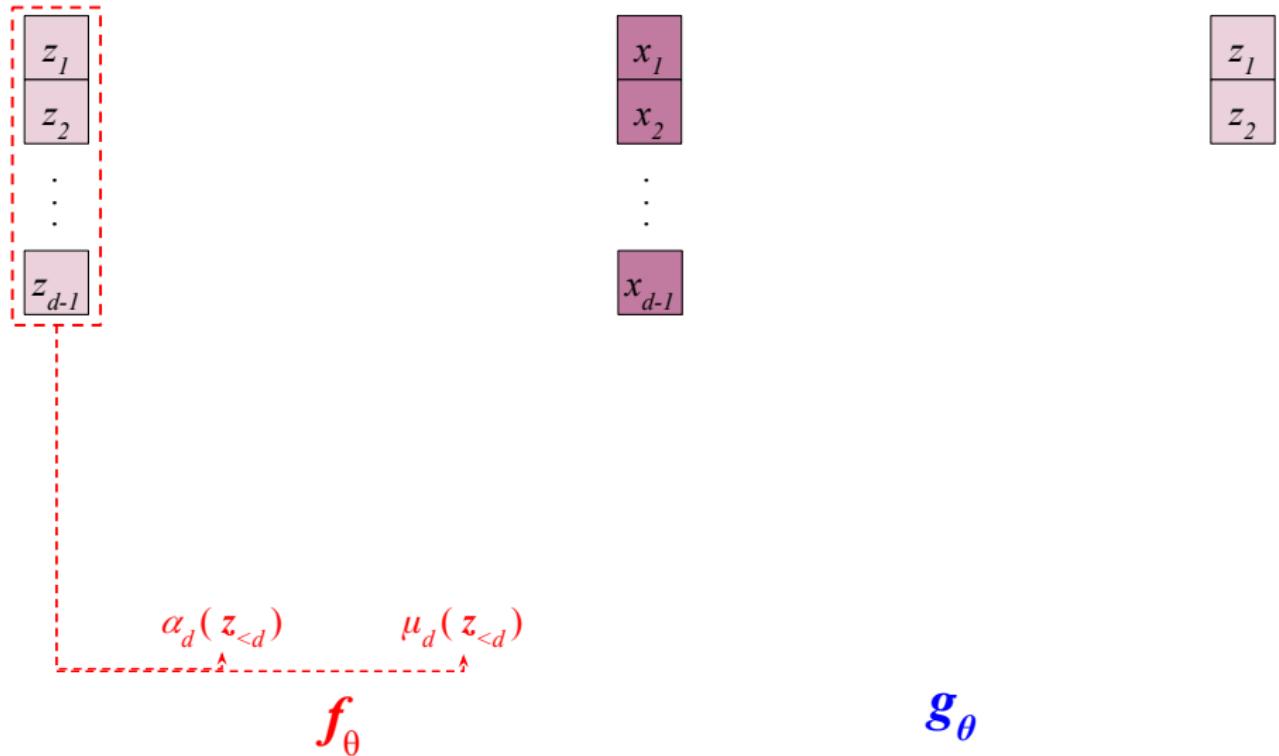


Figure: Forward transformation of  $d$ -th part

# Inverse Autoregressive Flows [5]

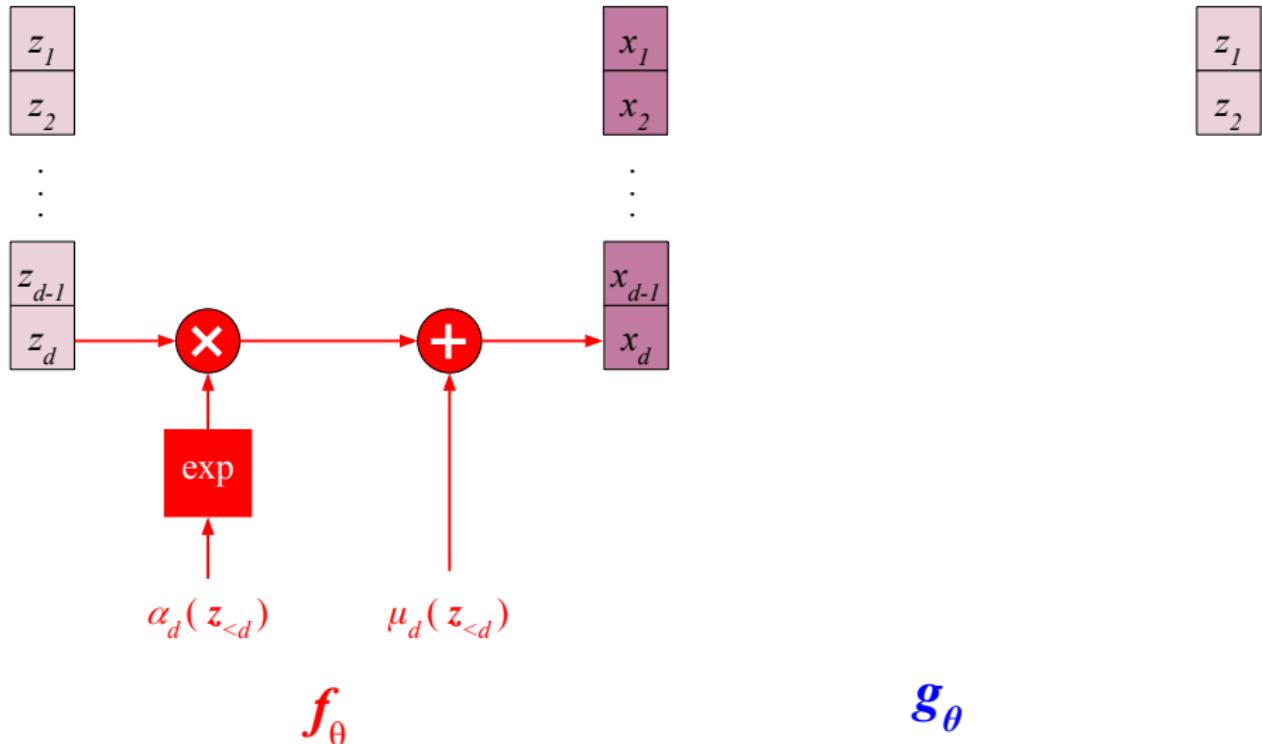


Figure: Forward transformation of  $d$ -th part

# Inverse Autoregressive Flows [5]

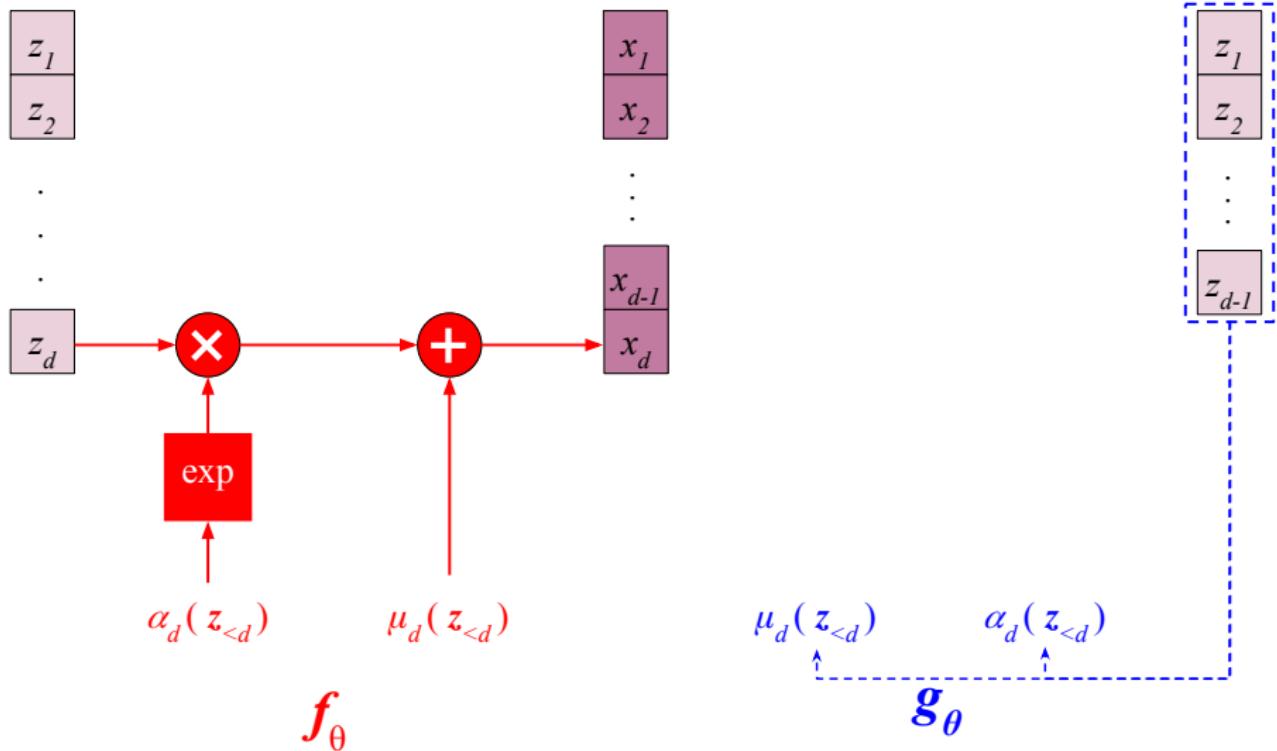


Figure: Inverse transformation of  $d$ -th part

# Inverse Autoregressive Flows [5]

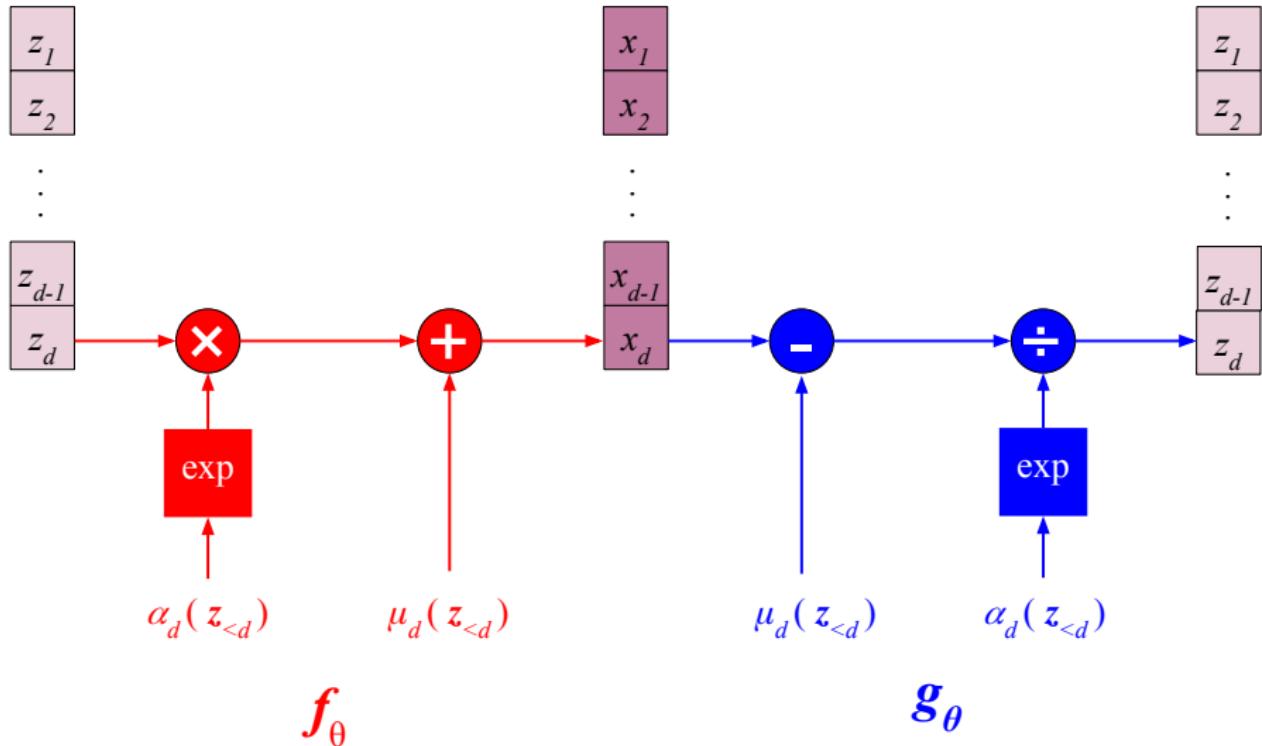


Figure: Inverse transformation of  $d$ -th part

# Inverse Autoregressive Flow [5]

## Jacobian

The Jacobian for the forward transformation of IAF is:

$$\begin{aligned}\log |\det \mathbf{J}_f| &= \log \left| \det \begin{bmatrix} \frac{\partial x_1}{\partial z_1} & \cdots & \frac{\partial x_1}{\partial z_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial x_D}{\partial z_1} & \cdots & \frac{\partial x_D}{\partial z_D} \end{bmatrix} \right| \\ &= \log \left| \det \begin{bmatrix} \exp(\alpha_1(\mathbf{z}_{<1})) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ \frac{\partial x_D}{\partial z_1} & \cdots & \exp(\alpha_D(\mathbf{z}_{$$

- ☞ The transformation is not volume preserving similar to Real-NVP and MAF.

# Inverse Autoregressive Flow [5]

## Sampling

Sampling from an IAF can be done in parallel as:

- Sample  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Calculate the following vectors (elements can be calculated in parallel)

$$\boldsymbol{\mu} = \begin{bmatrix} \mu_1(\mathbf{z}_{<1}) \\ \mu_2(\mathbf{z}_{<2}) \\ \vdots \\ \mu_D(\mathbf{z}_{<D}) \end{bmatrix}, \quad \boldsymbol{\alpha} = \begin{bmatrix} \alpha_1(\mathbf{z}_{<1}) \\ \alpha_2(\mathbf{z}_{<2}) \\ \vdots \\ \alpha_D(\mathbf{z}_{<D}) \end{bmatrix}$$

- Generate  $\mathbf{x}$  as:

$$\mathbf{x} = \mathbf{z} \odot \exp(\boldsymbol{\alpha}) + \boldsymbol{\mu}$$

## Parallel Sampling

The sampling can be done in parallel which can be done fast.

# Inverse Autoregressive Flow [5]

## Likelihood Calculation

Probability calculation for a given  $\mathbf{x}$  can be done sequentially as:

- $z_1 = \frac{x_1 - \mu_1}{\exp(\alpha_1)}$ , Calculate  $\begin{cases} \mu_2 \triangleq \mu_2(\mathbf{z}_{<2}) \\ \alpha_2 \triangleq \alpha_2(\mathbf{z}_{<2}) \end{cases}$
- $z_2 = \frac{x_2 - \mu_2}{\exp(\alpha_2)}$ , Calculate  $\begin{cases} \mu_3 \triangleq \mu_3(\mathbf{z}_{<3}) \\ \alpha_3 \triangleq \alpha_3(\mathbf{z}_{<3}) \end{cases}$
- $\vdots$
- $z_D = \frac{x_D - \mu_D}{\exp(\alpha_D)}$

Calculate probability using  $p_{\theta}(\mathbf{x}) = p(\mathbf{z}) \left| \prod_{d=1}^D \exp(\alpha_d(\mathbf{z}_{<d})) \right|^{-1}$

## Sequential Likelihood Calculation

The likelihood calculation and thus learning which needs transformation Jacobian is sequential.

## Section 9

### Conclusion

# Conclusion

## Normalizing Flow Models

- Pros
  - Exact likelihood calculation
  - Invertible latent transformation
  - Latent representation extrapolation
- Cons
  - Same latent and visible space dimension
  - Limitation to invertible transformation
  - Scalability issues due to latent dimension

# List of Abbreviations

Complete	Abbreviation
Inverse autoregressive Flow	IAF
Masked Autoregressive Flow	MAF
Nonlinear Independent Component Estimation	NICE
Real-valued Non-volume Preserving Transformations	Real-NVP
Variational AutoEncoder	VAE

# References I

-  [Christine Breiner](#),  
“Change of variables | mit 18.02sc multivariable calculus, fall 2010,” .
-  [Laurent Dinh, David Krueger, and Yoshua Bengio](#),  
“Nice: Non-linear independent components estimation,”  
*arXiv preprint arXiv:1410.8516*, 2014.
-  [Laurent Dinh, Jascha Sohl-Dickstein, and Samy Bengio](#),  
“Density estimation using real nvp,”  
*arXiv preprint arXiv:1605.08803*, 2016.
-  [George Papamakarios, Theo Pavlakou, and Iain Murray](#),  
“Masked autoregressive flow for density estimation,”  
*Advances in neural information processing systems*, vol. 30, 2017.
-  [Durk P Kingma, Tim Salimans, Rafal Jozefowicz, Xi Chen, Ilya Sutskever, and Max Welling](#),  
“Improved variational inference with inverse autoregressive flow,”  
*Advances in neural information processing systems*, vol. 29, 2016.