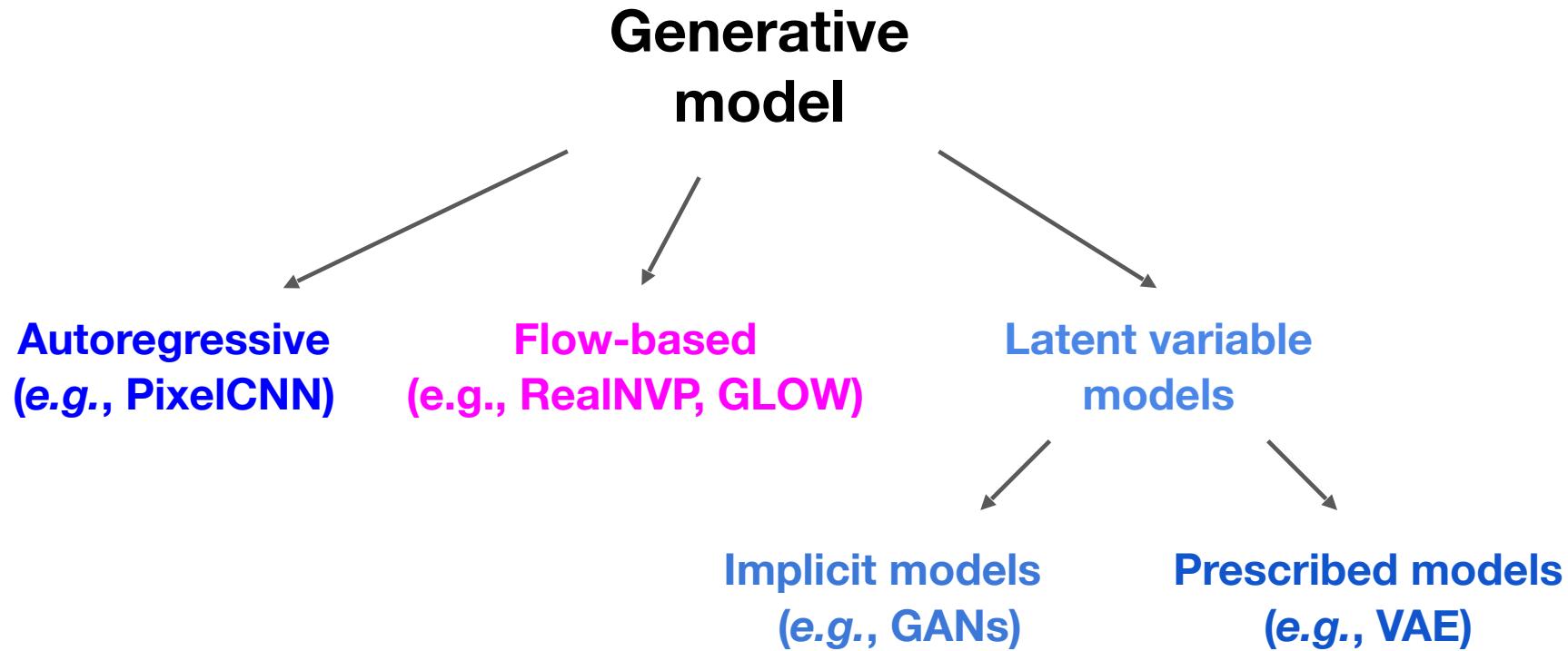


Deep generative modeling: ARMs and Normalizing Flows

Jakub M. Tomczak
Deep Learning

TYPES OF GENERATIVE MODELS



GENERATIVE MODELS

	Training	Likelihood	Sampling	Compression
Autoregressive models (e.g., PixelCNN)	Stable	Exact	Slow	No
Flow-based models (e.g., RealNVP)	Stable	Exact	Fast/Slow	No
Implicit models (e.g., GANs)	Unstable	No	Fast	No
Prescribed models (e.g., VAEs)	Stable	Approximate	Fast	Yes

GENERATIVE MODELS

	Training	Likelihood	Sampling	Compression
Autoregressive models (e.g., PixelCNN)	Stable	Exact	Slow	No
Flow-based models (e.g., RealNVP)	Stable	Exact	Fast/Slow	No
Implicit models (e.g., GANs)	Unstable	No	Fast	No
Prescribed models (e.g., VAEs)	Stable	Approximate	Fast	Yes

ARMS: AUTOREGRESSIVE MODELS

REPRESENTING A JOINT DISTRIBUTION

There are two main rules in the probability theory:

- Sum rule: $p(x) = \sum_y p(x, y)$
- Product rule: $p(x, y) = p(y | x) p(x)$

REPRESENTING A JOINT DISTRIBUTION

There are two main rules in the probability theory:

- Sum rule: $p(x) = \sum_y p(x, y)$
- Product rule: $p(x, y) = p(y | x) p(x)$

Before, we used these two rules for latent-variable models:

$$\begin{aligned} p(\mathbf{x}) &= \int p(\mathbf{x}, \mathbf{z}) \, d\mathbf{z} \\ &= \int p(\mathbf{x} | \mathbf{z}) p(\mathbf{z}) \, d\mathbf{z} \end{aligned}$$

REPRESENTING A JOINT DISTRIBUTION

There are two main rules in the probability theory:

- Sum rule: $p(x) = \sum_y p(x, y)$
- Product rule: $p(x, y) = p(y | x) p(x)$

Now, we will use the product rule to express the distribution of $\mathbf{x} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = p(x_1) \sum_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$

REPRESENTING A JOINT DISTRIBUTION

We use the product rule to express the distribution of $\mathbf{x} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = p(x_1) \sum_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$.

REPRESENTING A JOINT DISTRIBUTION

We use the product rule to express the distribution of $\mathbf{x} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = p(x_1) \sum_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$.

The order of variables isn't important.

REPRESENTING A JOINT DISTRIBUTION

We use the product rule to express the distribution of $\mathbf{x} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = p(x_1) \sum_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$.

The order of variables isn't important.

However, modeling all conditionals separately is **infeasible**...

REPRESENTING A JOINT DISTRIBUTION

We use the product rule to express the distribution of $\mathbf{x} \in \mathbb{R}^D$:

$$p(\mathbf{x}) = p(x_1) \sum_{d=2}^D p(x_d | \mathbf{x}_{<d})$$

where $\mathbf{x}_{<d} = [x_1, x_2, \dots, x_{d-1}]^\top$.

The order of variables isn't important.

However, modeling all conditionals separately is **infeasible**...

Can we do better that?

REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

For instance, for two last variables:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \sum_{d=3}^D p(x_d|x_{d-2}, x_{d-1})$$

Now, we can model $p(x_d|x_{d-2}, x_{d-1})$ by a **single model**.

→ E.g., we can take a **neural network**.

REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

For instance, for two last variables:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \sum_{d=3}^D p(x_d|x_{d-2}, x_{d-1})$$

REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

For instance, for two last variables:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \sum_{d=3}^D p(x_d|x_{d-2}, x_{d-1})$$

Now, we can model $p(x_d|x_{d-2}, x_{d-1})$ by a **single model**.

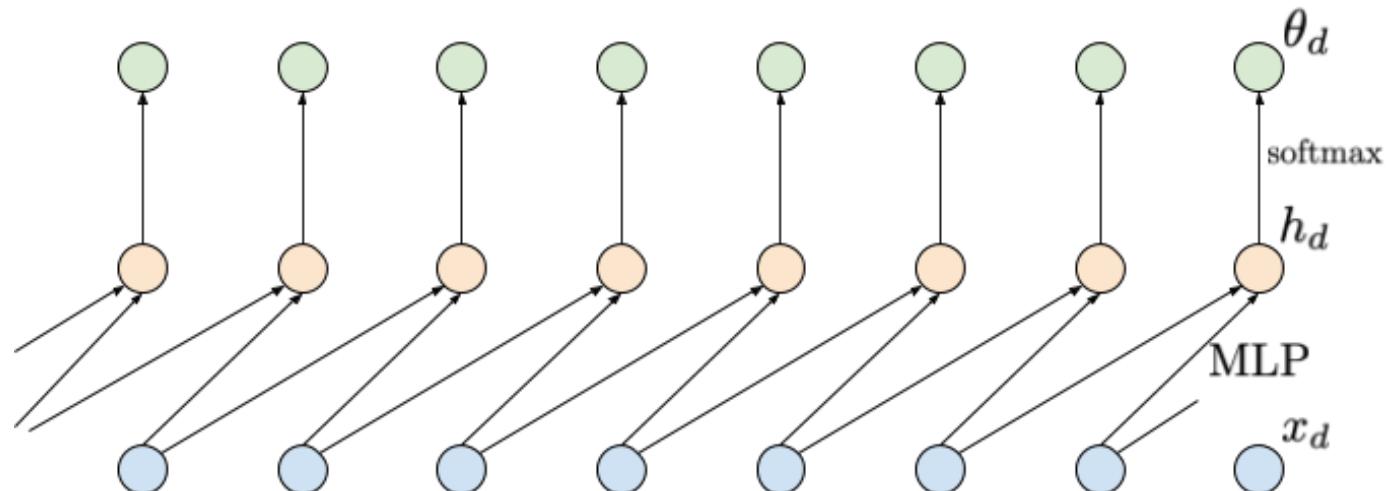
→ E.g., we can take a **neural network**.

REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

For instance, for two last variables:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \sum_{d=3}^D p(x_d|x_{d-2}, x_{d-1})$$



REPRESENTING A JOINT DISTRIBUTION

We can assume a **finite dependency**.

For instance, for two last variables:

$$p(\mathbf{x}) = p(x_1)p(x_2|x_1) \sum_{d=3}^D p(x_d|x_{d-2}, x_{d-1})$$

Now, we can model $p(x_d|x_{d-2}, x_{d-1})$ by a **single model**.

→ E.g., we can take a **neural network**.

However, it is still pretty **limiting**, because we need to decide on the length of the dependency.

AUTOREGRESSIVE MODELS (ARM)

Instead, we can use RNNs to model the conditionals:

$$p(x_d | \mathbf{x}_{<d}) = p(x_d | RNN(x_{d-1}, h_{d-1}))$$

where $h_d = RNN(x_{d-1}, h_{d-1})$.

AUTOREGRESSIVE MODELS (ARM)

Instead, we can use RNNs to model the conditionals:

$$p(x_d | \mathbf{x}_{)} = p(x_d | RNN(x_{d-1}, h_{d-1}))$$

where $h_d = RNN(x_{d-1}, h_{d-1})$.

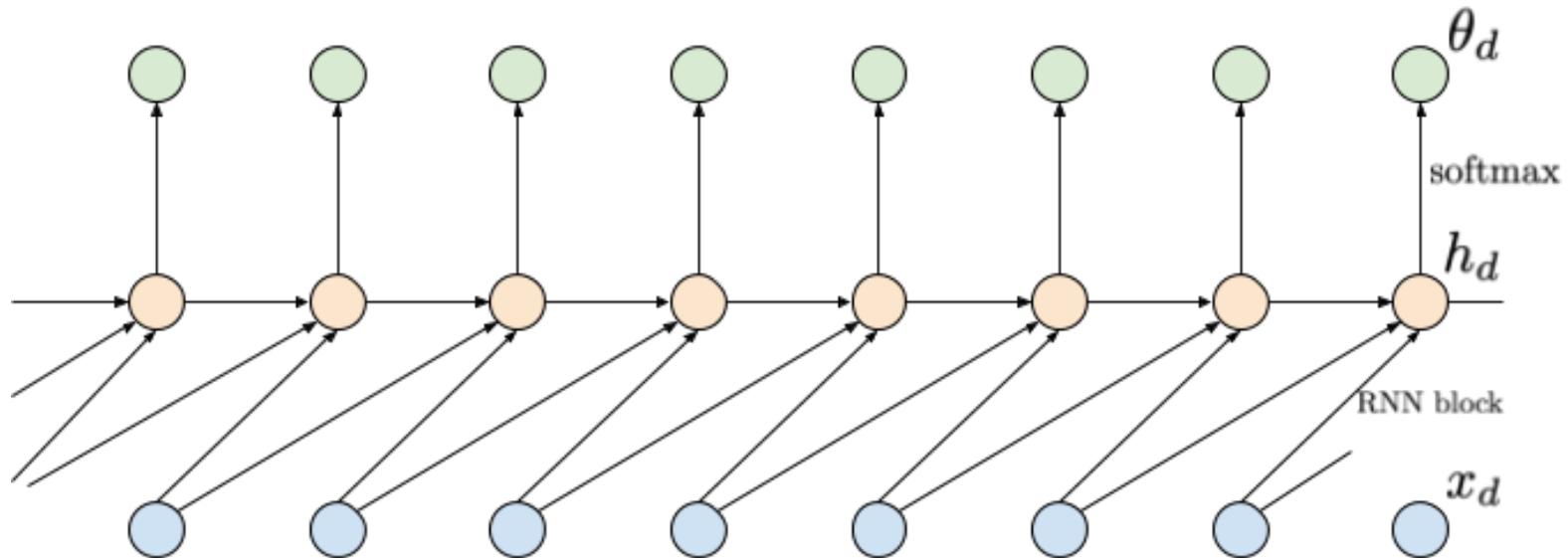
Advantages:

- We don't need to define dependencies.
- A **single parameterization**.

AUTOREGRESSIVE MODELS (ARM)

Instead, we can use RNNs to model the conditionals:

$$p(x_d | \mathbf{x}_{$$



AUTOREGRESSIVE MODELS (ARM)

Instead, we can use RNNs to model the conditionals:

$$p(x_d | \mathbf{x}_{<d}) = p(x_d | RNN(x_{d-1}, h_{d-1}))$$

where $h_d = RNN(x_{d-1}, h_{d-1})$.

Advantages:

- We don't need to define dependencies.
- A **single parameterization**.

RNN are slow, because they're sequential.

AUTOREGRESSIVE MODELS (ARM)

Instead, we can use RNNs to model the conditionals:

$$p(x_d | \mathbf{x}_{)} = p(x_d | RNN(x_{d-1}, h_{d-1}))$$

where $h_d = RNN(x_{d-1}, h_{d-1})$.

Advantages:

- We don't need to define dependencies.
- A **single parameterization**.

RNN are slow, because they're sequential. Can we do better?

USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

Let us consider a sequence $\mathbf{x} = [x_1, x_2, \dots, x_D]^\top$.

We assume all observed data are D -dimensional.

USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

Let us consider a sequence $\mathbf{x} = [x_1, x_2, \dots, x_D]^T$.

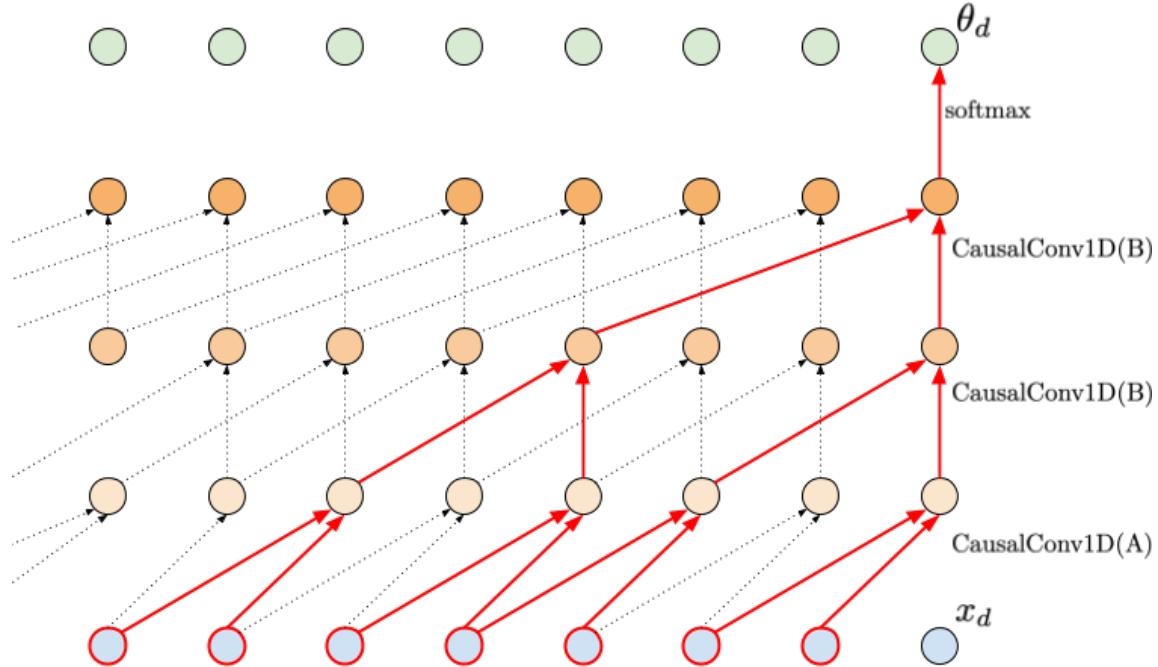
We assume all observed data are D -dimensional.

We can use **1D convolutional layers** to process all signals at once.

Moreover, we can use **dilation** to learn **long-range dependencies**.

USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

Let us consider a sequence $\mathbf{x} = [x_1, x_2, \dots, x_D]^T$.

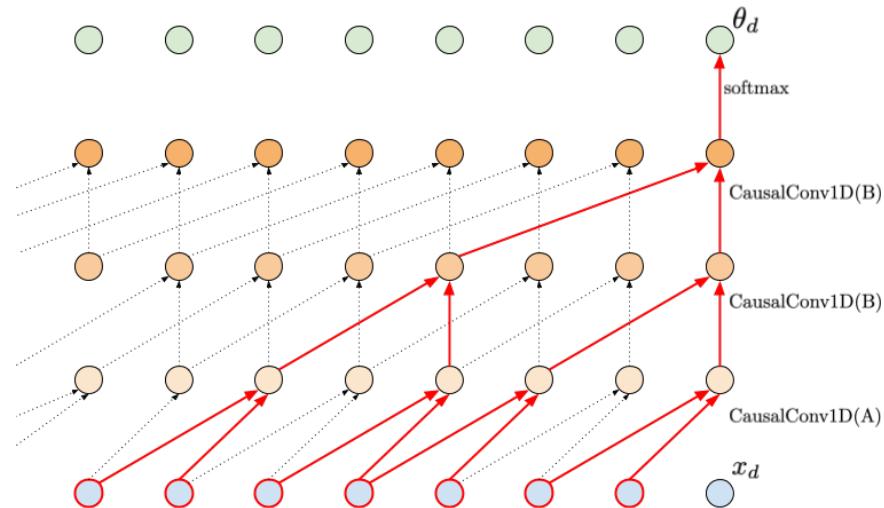


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)

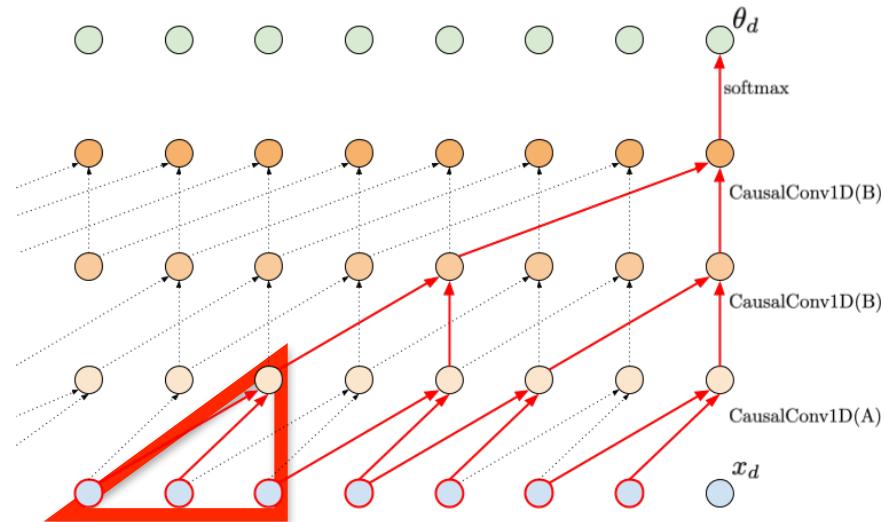


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)

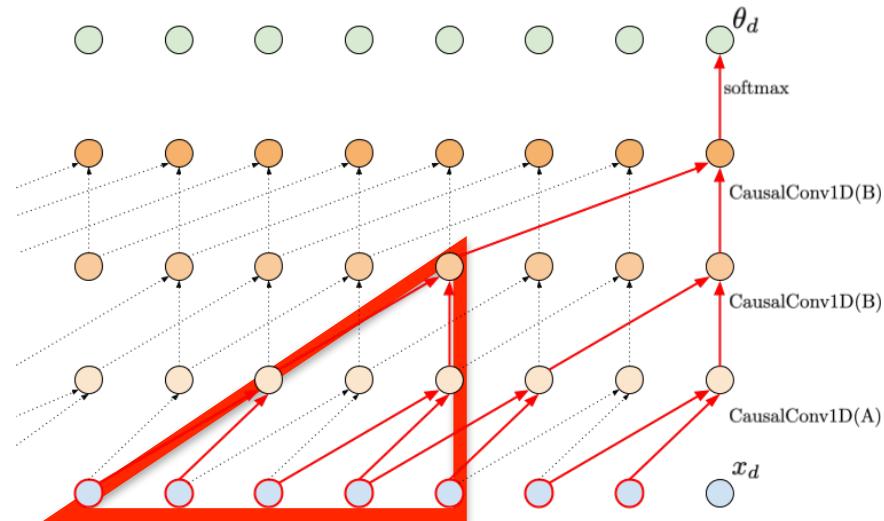


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)

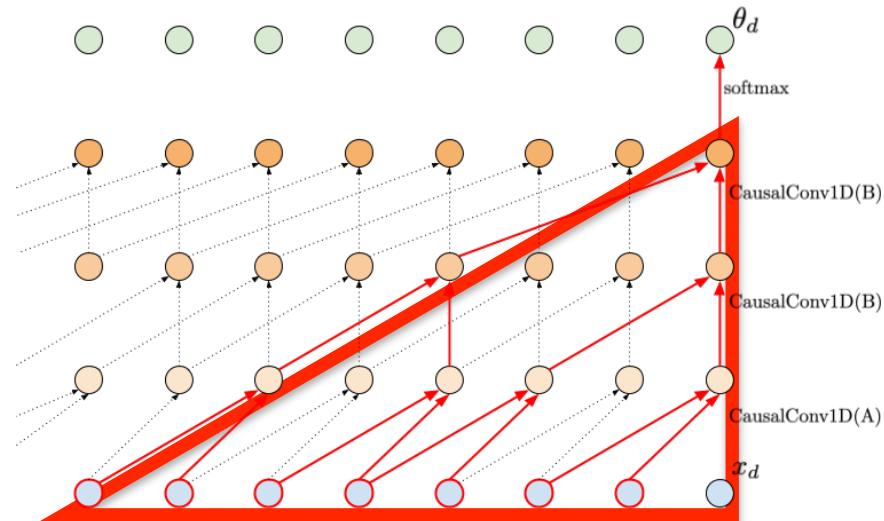


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)

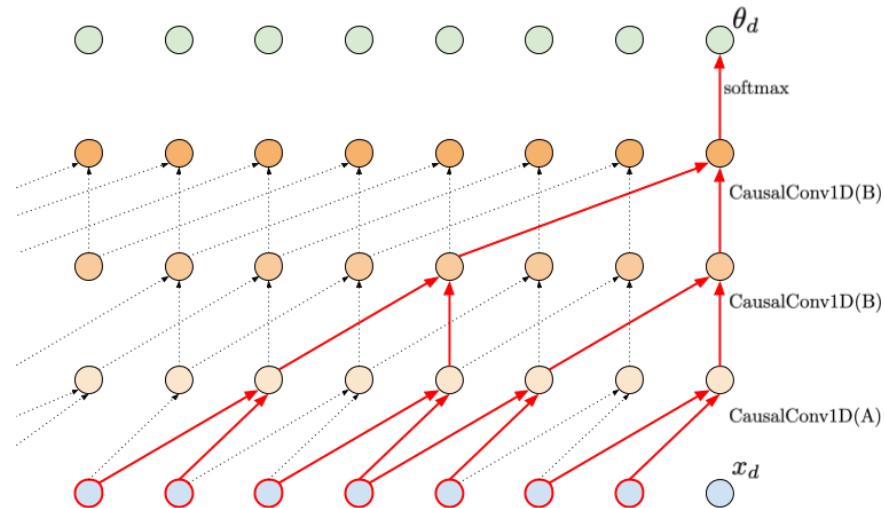


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)
- Very **efficient** using current
DL frameworks.

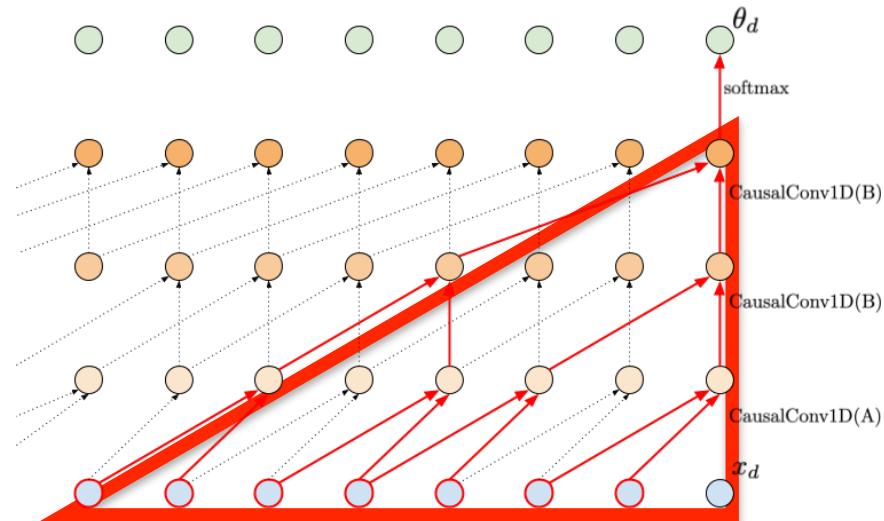


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

Notice:

- **Causal convolution**
(i.e., looking only to the “past”)
- Very **efficient** using current
DL frameworks.
- For deep neural networks, NNs learn long-range dependencies.

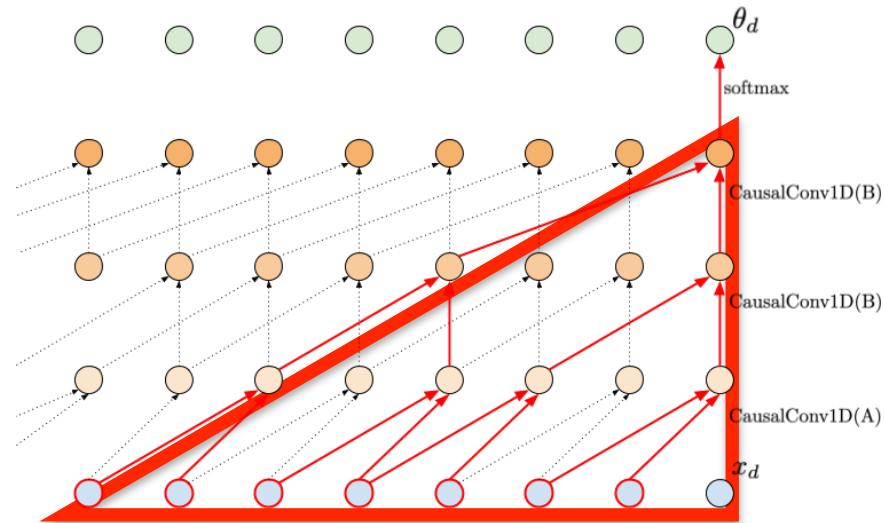


USING CONVOLUTIONAL NEURAL NETWORKS FOR MODELING SEQUENCES

We can use **1D convolutional layers** to process all signals at once.

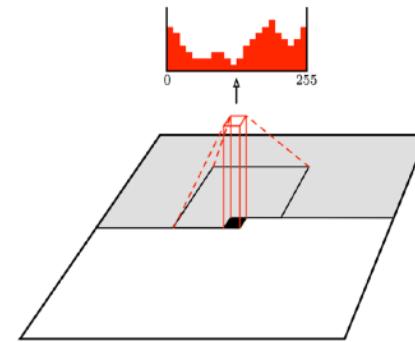
Notice:

- **Causal convolution**
(i.e., looking only to the “past”)
- Very **efficient** using current
DL frameworks.
- For deep neural networks, NNs learn long-range dependencies.



We can utilize the very same idea for images, but using **2D convolutions**.

We need to remember about **causal convolutions!**

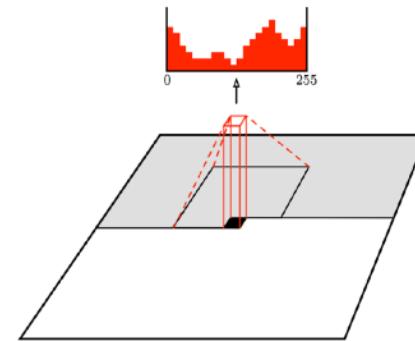


1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

We can utilize the very same idea for images, but using **2D convolutions**.

We need to remember about **causal convolutions!**

Moreover, we should think of 2 or even **3 dimensions** ($C \times H \times W$).



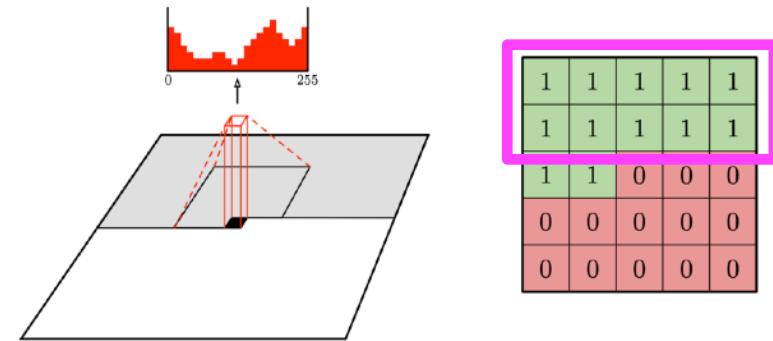
1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

We can accomplish that by composing **two Conv2D layers**.

We can utilize the very same idea for images, but using **2D convolutions**.

We need to remember about **causal convolutions!**

Moreover, we should think of 2 or even **3 dimensions** ($C \times H \times W$).



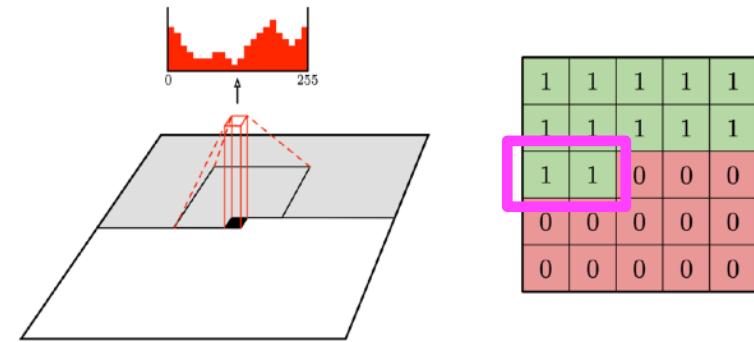
We can accomplish that by composing **two Conv2D layers**.

The first Conv2D layer covers the upper part.

We can utilize the very same idea for images, but using **2D convolutions**.

We need to remember about **causal convolutions!**

Moreover, we should think of 2 or even **3 dimensions** ($C \times H \times W$).



We can accomplish that by composing **two Conv2D layers**.

The first Conv2D layer covers the upper part.

The second Conv2D layer covers the left part.

We can utilize the very same idea for images, but using **2D convolutions**.

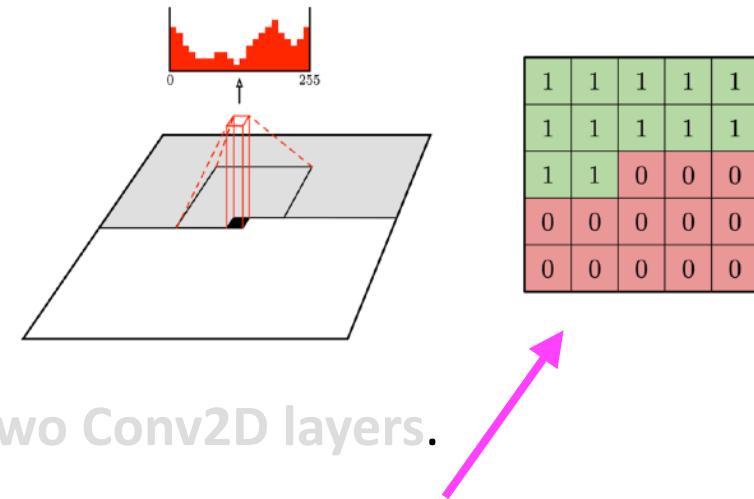
We need to remember about **causal convolutions!**

Moreover, we should think of 2 or even **3 dimensions** ($C \times H \times W$).

We can accomplish that by composing two Conv2D layers.

The first Conv2D layer covers the upper part.

The second Conv2D layer covers the left part.



using masking for kernel weights.

Originally, PixelCNN used the **softmax non-linearity** at the end to output integers between 0 and 255 (i.e., pixel values).

Currently, a **mixture of discretized logistic distributions** is used:

$$P(x \mid \pi, \mu, s) = \sum_{i=1}^K \pi_i \left[\sigma\left(\frac{(x + 0.5 - \mu_i)}{s_i}\right) - \sigma\left(\frac{(x - 0.5 - \mu_i)}{s_i}\right) \right]$$

³⁹ Salimans, Tim, et al. "PixelCNN++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications." *arXiv* (2017).

Originally, PixelCNN used the **softmax non-linearity** at the end to output integers between 0 and 255 (i.e., pixel values).

Currently, a **mixture of discretized logistic distributions** is used:

$$P(x \mid \pi, \mu, s) = \sum_{i=1}^K \pi_i \left[\sigma\left(\frac{(x + 0.5 - \mu_i)}{s_i}\right) - \sigma\left(\frac{(x - 0.5 - \mu_i)}{s_i}\right) \right]$$

sigmoid function

⁴⁰ Salimans, Tim, et al. "PixelCNN++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications." *arXiv* (2017).

Originally, PixelCNN used the **softmax non-linearity** at the end to output integers between 0 and 255 (i.e., pixel values).

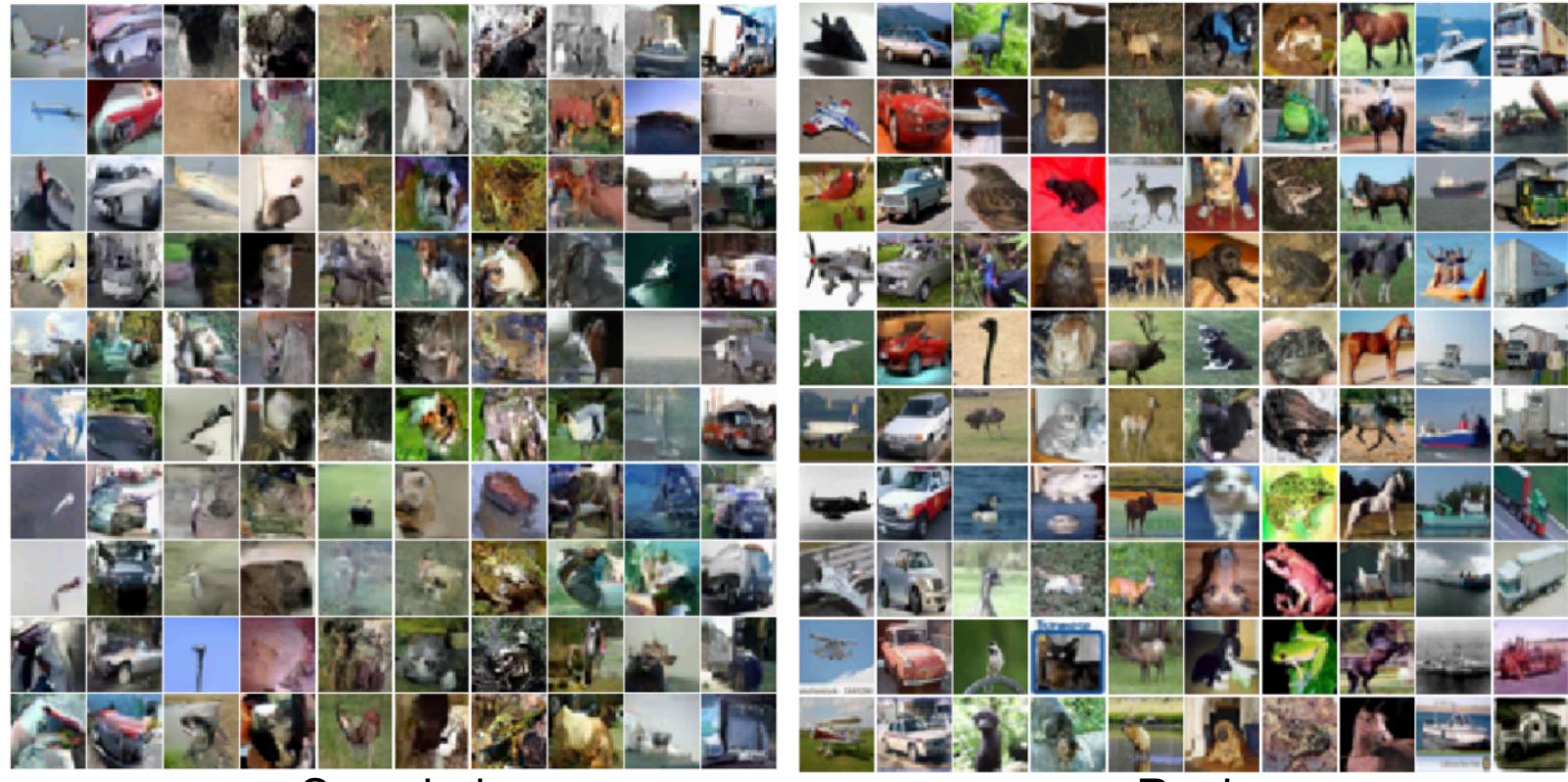
Currently, a **mixture of discretized logistic distributions** is used:

$$P(x \mid \pi, \mu, s) = \sum_{i=1}^K \pi_i \left[\sigma\left(\frac{(x + 0.5 - \mu_i)}{s_i}\right) - \sigma\left(\frac{(x - 0.5 - \mu_i)}{s_i}\right) \right]$$

Learnable as a parameter

⁴¹ Salimans, Tim, et al. "PixelCNN++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications." *arXiv* (2017).

PIXELCNN



Sampled

Real

⁴² Salimans, Tim, et al. "PixelCNN++: Improving the pixelcnn with discretized logistic mixture likelihood and other modifications." *arXiv* (2017).

Advantages

- ✓ Exact likelihood.
- ✓ Stable training.

Disadvantages

- Very slow sampling.
- No compression.
- Sometimes, low visual quality.

FLOW-BASED MODELS

CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

Then, we take the following transformation: $u = 0.75 \cdot v + 1$.

Q: What is the pdf for u ?

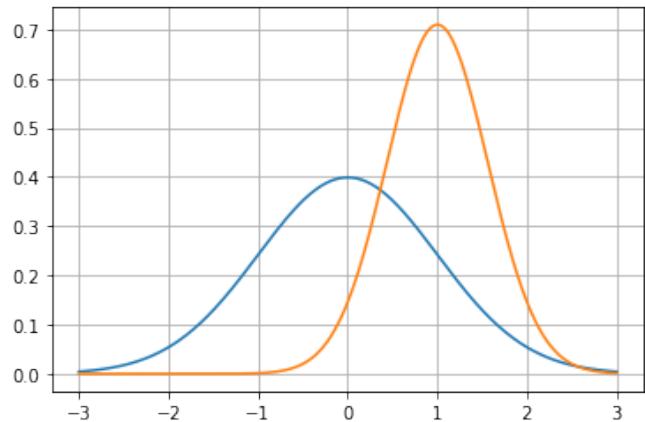
CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

Then, we take the following transformation: $u = 0.75 \cdot v + 1$.

Q: What is the pdf for u ?

A: The pdf is $p(u) = \mathcal{N}(u | 1, 0.75^2)$.



CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

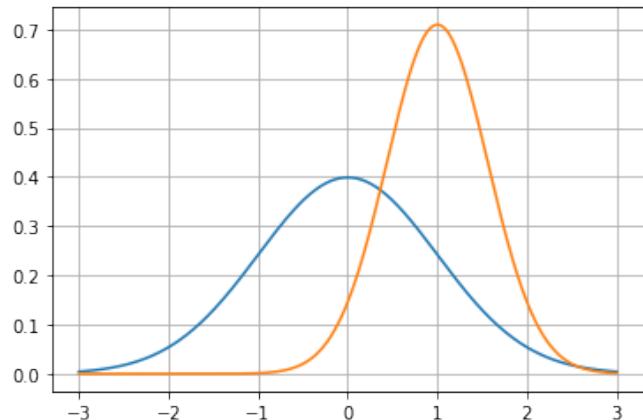
Then, we take the following transformation: $u = 0.75 \cdot v + 1$.

Q: What is the pdf for u ?

A: The pdf is $p(u) = \mathcal{N}(u | 1, 0.75^2)$.

In general, we have:

$$p(u) = p(f^{-1}(u)) \left| \frac{\partial f^{-1}(u)}{\partial u} \right|$$



CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

Then, we take the following transformation: $u = 0.75 \cdot v + 1$.

Q: What is the pdf for u ?

A: The pdf is $p(u) = \mathcal{N}(u | 1, 0.75^2)$.

In general, we have:

$$p(u) = p(f^{-1}(u)) \left| \frac{\partial f^{-1}(u)}{\partial u} \right|$$

$$f^{-1}(u) = \frac{u - 1}{0.75}$$

$$\left| \frac{\partial f^{-1}(u)}{\partial u} \right| = \frac{4}{3}$$

CHANGE OF VARIABLES

Let us consider a random variable $v \in \mathbb{R}$ with $p(v) = \mathcal{N}(v | 0, 1)$.

Then, we take the following transformation: $u = 0.75 \cdot v + 1$.

Q: What is the pdf for u ?

A: The pdf is $p(u) = \mathcal{N}(u | 1, 0.75^2)$.

In general, we have:

$$p(u) = p(f^{-1}(u)) \left| \frac{\partial f^{-1}(u)}{\partial u} \right|$$
$$p(u) = p\left(\frac{u-1}{0.75}\right) \frac{4}{3} = \frac{1}{\sqrt{2\pi} 0.75^2} \exp\left\{-\frac{(u-1)^2}{0.75^2}\right\}$$

Multidimensional case:

$$p(\mathbf{u}) = p(f^{-1}(\mathbf{u})) \left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right|$$

where:

$$\left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right| = \left| \det \mathbf{J}_{f^{-1}(\mathbf{u})} \right|$$

CHANGE OF VARIABLES

Multidimensional case:

$$p(\mathbf{u}) = p(f^{-1}(\mathbf{u})) \left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right|$$

where:

$$\left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right| = \left| \det \mathbf{J}_{f^{-1}(\mathbf{u})} \right|$$

Jacobian

$$\mathbf{J}_{f^{-1}} = \begin{bmatrix} \frac{\partial f_1^{-1}}{\partial u_1} & \cdots & \frac{\partial f_1^{-1}}{\partial u_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D^{-1}}{\partial u_1} & \cdots & \frac{\partial f_D^{-1}}{\partial u_D} \end{bmatrix}$$

CHANGE OF VARIABLES

Multidimensional case:

$$p(\mathbf{u}) = p(f^{-1}(\mathbf{u})) \left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right|$$

where:

$$\left| \frac{\partial f^{-1}(\mathbf{u})}{\partial \mathbf{u}} \right| = \left| \det \mathbf{J}_{f^{-1}(\mathbf{u})} \right|$$

Jacobian

$$\mathbf{J}_{f^{-1}} = \begin{bmatrix} \frac{\partial f_1^{-1}}{\partial u_1} & \cdots & \frac{\partial f_1^{-1}}{\partial u_D} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_D^{-1}}{\partial u_1} & \cdots & \frac{\partial f_D^{-1}}{\partial u_D} \end{bmatrix}$$

How can we utilize this idea?

APPLYING CHANGE OF VARIABLES AND INVERTIBLE TRANSFORMATIONS

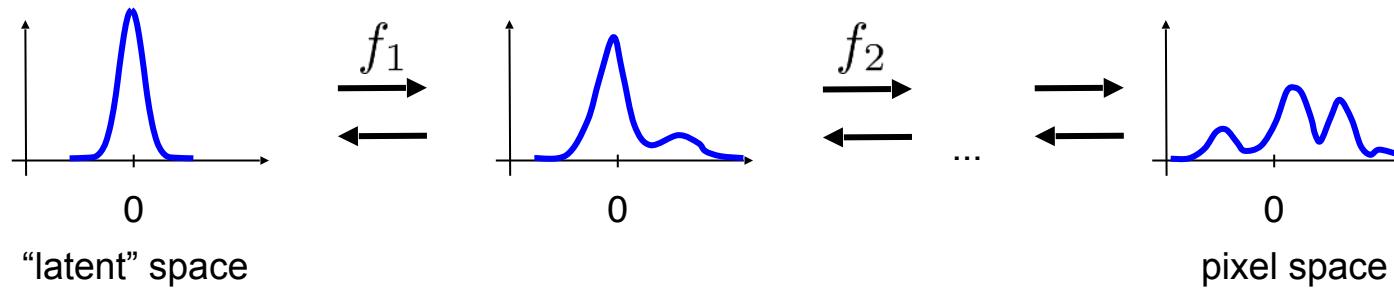
Let us consider a sequence of invertible transformations $f_k : \mathbb{R}^D \rightarrow \mathbb{R}^D$.

We can start with a *simple* distribution, e.g., $\pi(\mathbf{z}) = \mathcal{N}(\mathbf{z} | \mathbf{0}, \mathbf{I})$.

APPLYING CHANGE OF VARIABLES AND INVERTIBLE TRANSFORMATIONS

Let us consider a sequence of invertible transformations $f_k : \mathbb{R}^D \rightarrow \mathbb{R}^D$.

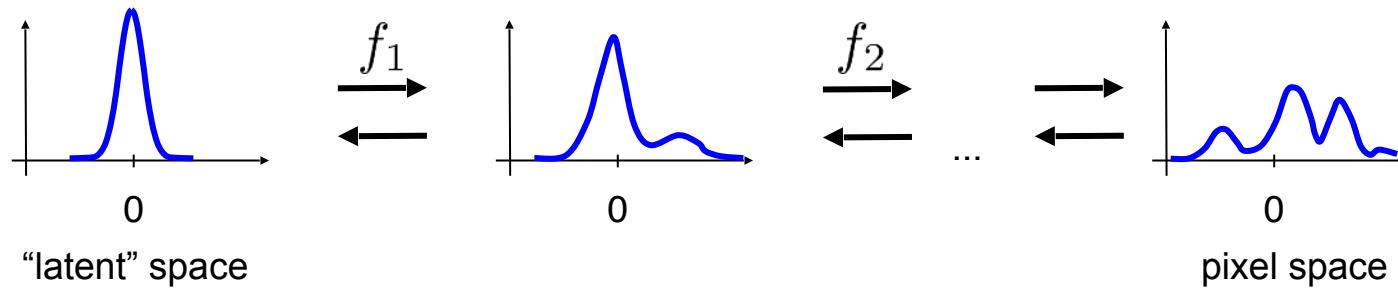
We can start with a *simple* distribution, e.g., $\pi(\mathbf{z}) = \mathcal{N}(\mathbf{z} | 0, \mathbf{I})$.



APPLYING CHANGE OF VARIABLES AND INVERTIBLE TRANSFORMATIONS

Let us consider a sequence of invertible transformations $f_k : \mathbb{R}^D \rightarrow \mathbb{R}^D$.

We can start with a *simple* distribution, e.g., $\pi(\mathbf{z}) = \mathcal{N}(\mathbf{z} | 0, \mathbf{I})$.



This results in: $p(\mathbf{x}) = \pi(\mathbf{z}_0) \prod_{i=1}^K \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}$

2D EXAMPLE

Inference

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

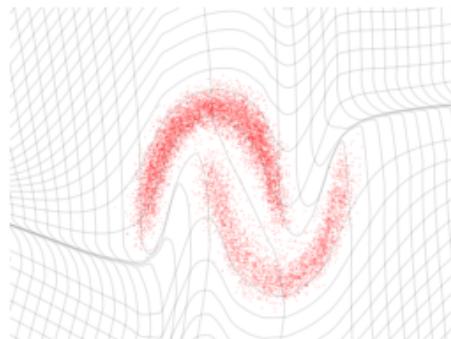
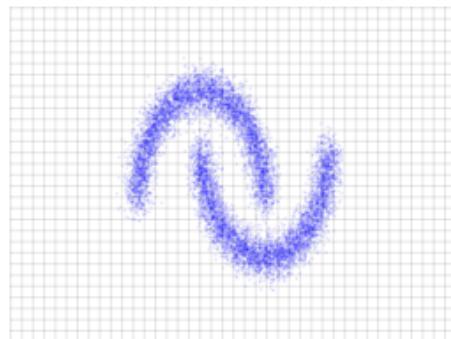
$$z = f(x)$$

Generation

$$z \sim p_Z$$

$$x = f^{-1}(z)$$

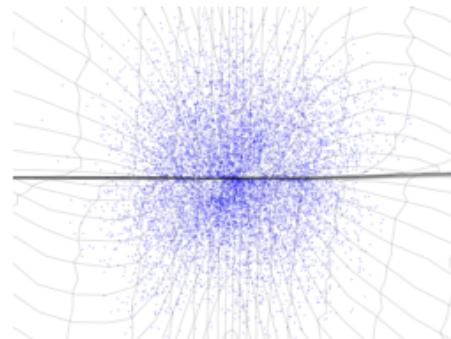
Data space \mathcal{X}



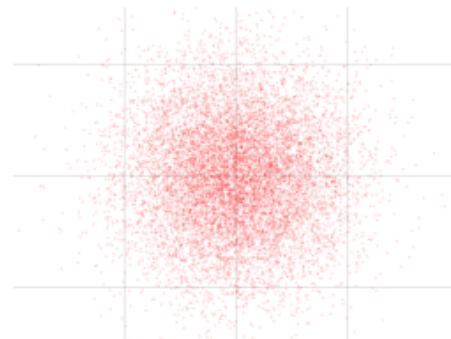
$$\{f_k\}$$

\Rightarrow

Latent space \mathcal{Z}



$$\{f_k^{-1}\}$$



DENSITY MODELING WITH INVERTIBLE NEURAL NETWORKS

The density model: $p(\mathbf{x}) = \pi(\mathbf{z}_0) \prod_{i=1}^K \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}$

DENSITY MODELING WITH INVERTIBLE NEURAL NETWORKS

The density model: $p(\mathbf{x}) = \pi(\mathbf{z}_0) \prod_{i=1}^K \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}$

In order to obtain flexible transformations f_k , we use **neural networks**.

However, we need to ensure that they are **invertible**.

The density model: $p(\mathbf{x}) = \pi(\mathbf{z}_0) \prod_{i=1}^K \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}$

In order to obtain flexible transformations f_k , we use **neural networks**.

However, we need to ensure that they are **invertible**.

Moreover, we cannot use any invertible neural network, because we need to remember about the **Jacobian**!

DENSITY MODELING WITH INVERTIBLE NEURAL NETWORKS

$$\text{The density model: } p(\mathbf{x}) = \pi(\mathbf{z}_0) \prod_{i=1}^K \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1}$$

In order to obtain flexible transformations f_k , we use **neural networks**.

However, we need to ensure that they are **invertible**.

Moreover, we cannot use any invertible neural network, because we need to remember about the **Jacobian**!

Calculating Jacobian is the main challenge in flow-based models.

Design the invertible transformations as follows:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp\left(s(\mathbf{x}_{1:d})\right) + t(\mathbf{x}_{1:d})$$

where: $s(\cdot)$ and $t(\cdot)$ are **arbitrary** neural networks.

Design the invertible transformations as follows:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp\left(s(\mathbf{x}_{1:d})\right) + t(\mathbf{x}_{1:d})$$

where: $s(\cdot)$ and $t(\cdot)$ are **arbitrary** neural networks.

This is **invertible by design**, because:

$$\mathbf{x}_{d+1:D} = \left(\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d}) \right) \odot \exp\left(-s(\mathbf{y}_{1:d})\right)$$

$$\mathbf{x}_{1:d} = \mathbf{y}_{1:d}$$

Design the invertible transformations as follows:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp\left(s(\mathbf{x}_{1:d})\right) + t(\mathbf{x}_{1:d})$$

Known as
affine coupling layer

where: $s(\cdot)$ and $t(\cdot)$ are **arbitrary** neural networks.

This is **invertible by design**, because:

$$\mathbf{x}_{d+1:D} = \left(\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d}) \right) \odot \exp\left(-s(\mathbf{y}_{1:d})\right)$$

$$\mathbf{x}_{1:d} = \mathbf{y}_{1:d}$$

The invertible transformation:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

What about the Jacobian?

The invertible transformation:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

What about the Jacobian?

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

The invertible transformation:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

What about the Jacobian?

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

$$\text{and } \det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp(s(\mathbf{x}_{1:d}))_j = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_{1:d})_j\right)$$

The invertible transformation:

$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

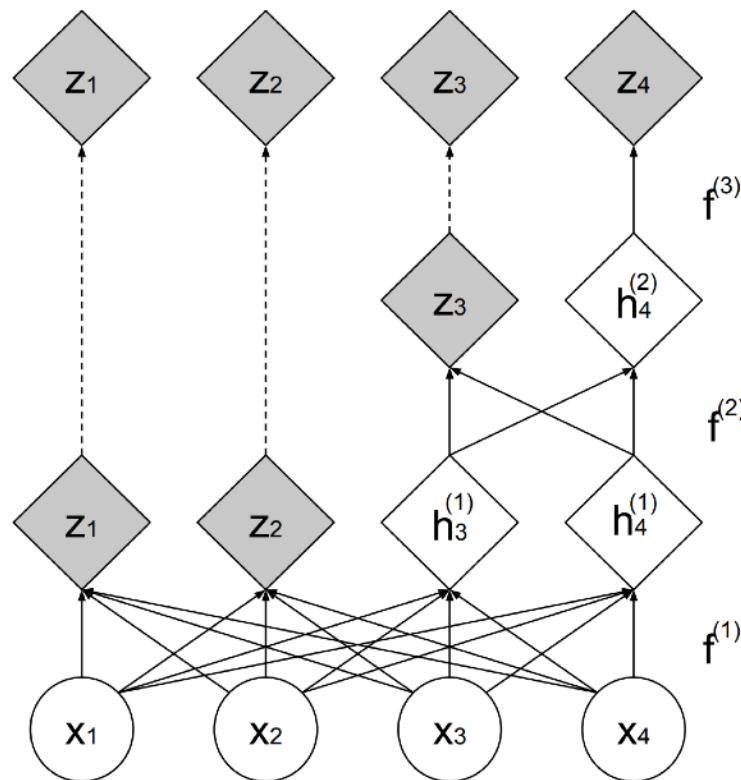
What about the Jacobian?

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

$$\text{and } \det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp(s(\mathbf{x}_{1:d}))_j = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_{1:d})_j\right)$$

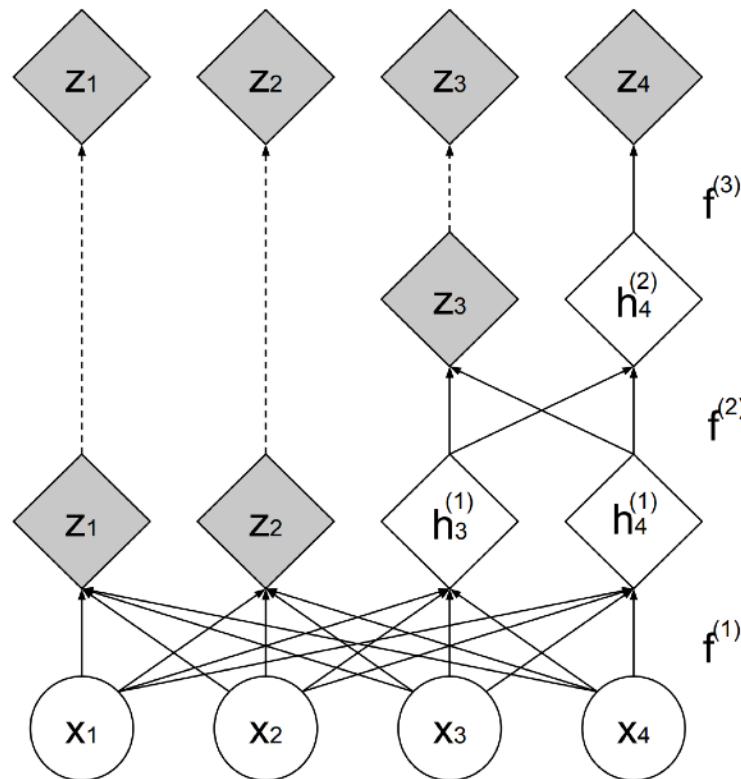
Easy to calculate!

We can also introduce the idea of autoregressive modeling here as well:



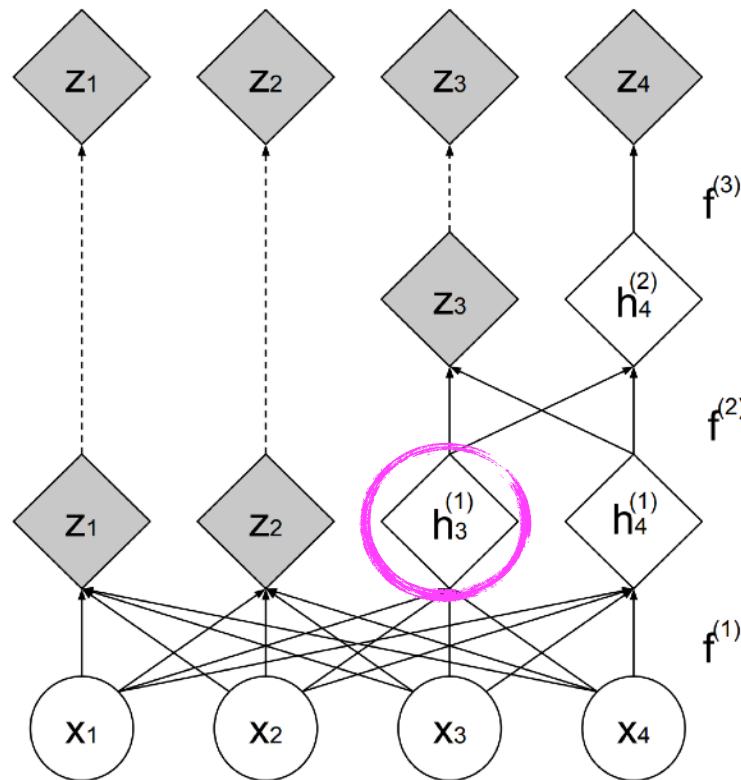
We can also introduce the idea of autoregressive modeling here as well:

$$p(z_1)p(z_2)p(z_3 | z_1, z_2)p(z_4 | z_1, z_2, z_3)$$



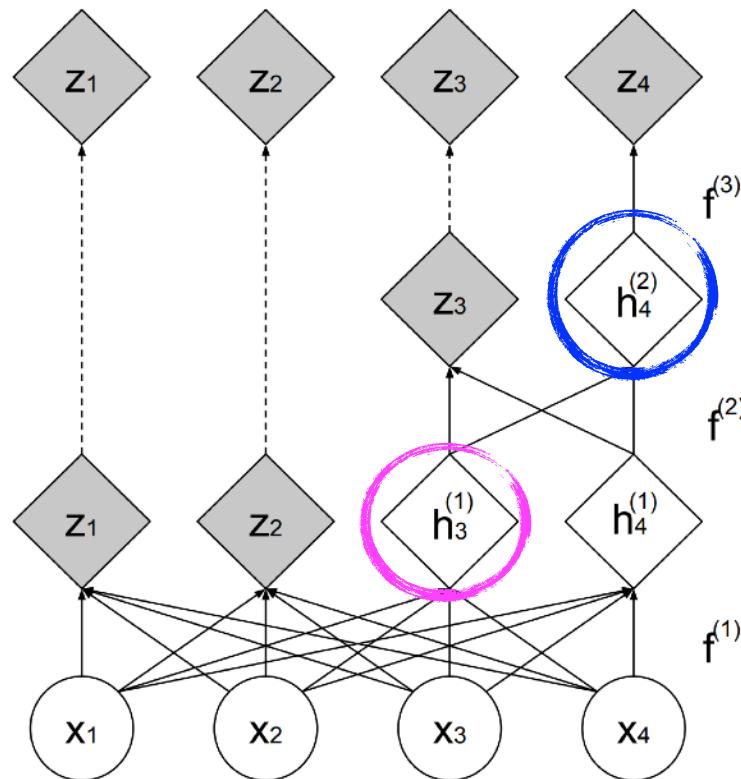
We can also introduce the idea of autoregressive modeling here as well:

$$p(z_1)p(z_2)p(z_3 | z_1, z_2)p(z_4 | z_1, z_2, z_3)$$



We can also introduce the idea of autoregressive modeling here as well:

$$p(z_1)p(z_2)p(z_3 | z_1, z_2)p(z_4 | z_1, z_2, z_3)$$



Moreover, we can use additional transformations:

1. **Permutations** of variables (this is invertible).

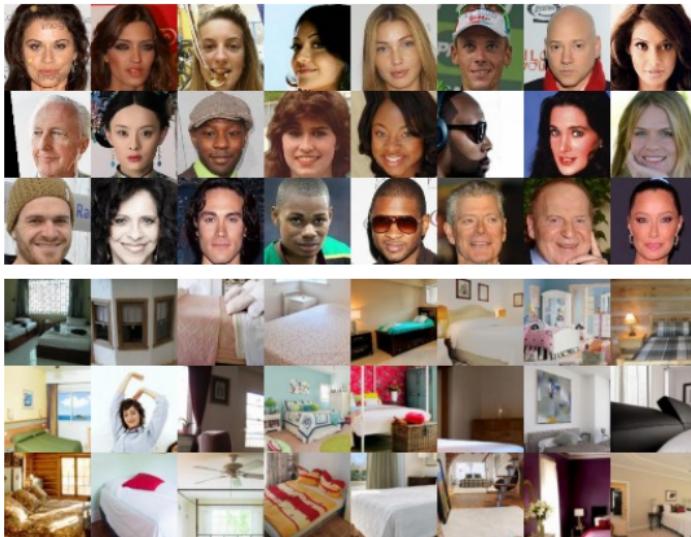
→ this helps to “mix” variables.

2. Divide variables using a **checkerboard pattern**.

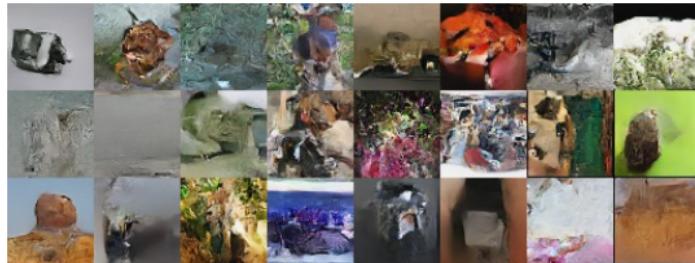
→ this helps to learn higher-order dependencies.

3. Use **squeezing**: reshape input tensor

→ reshaping can help to “mix” variables.



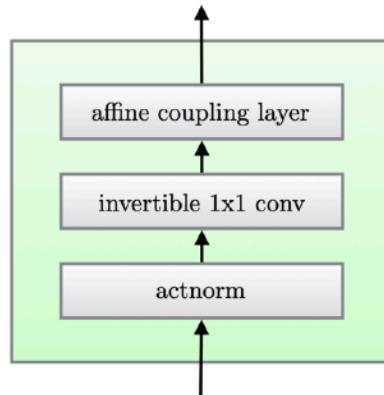
REALNVP



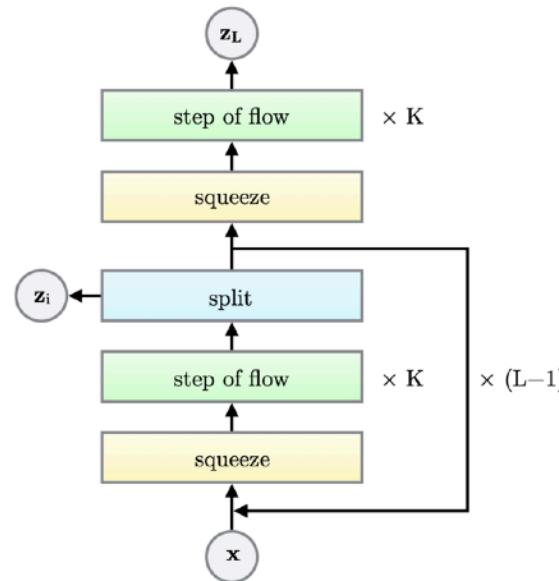
GLOW: REALNVP WITH 1X1 CONVOLUTIONS

A model contains ~1000 convolutions.

A new component: 1x1 convolution instead of a permutation matrix.



(a) One step of our flow.



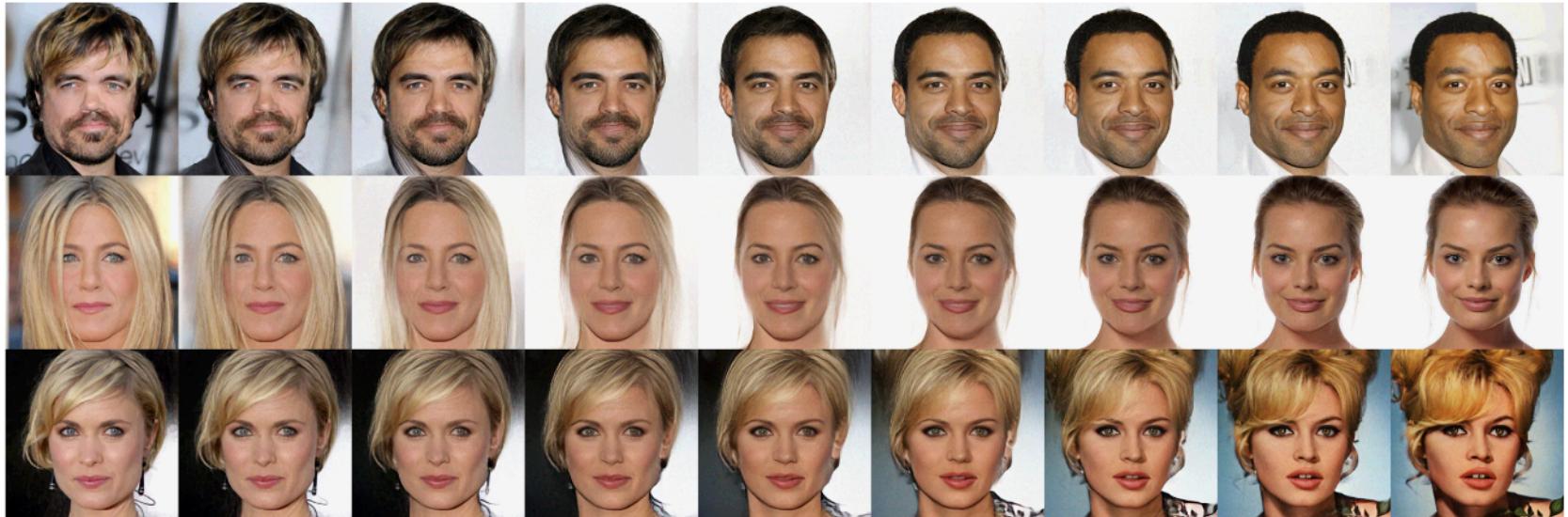
(b) Multi-scale architecture (Dinh et al., 2016).

GLOW: SAMPLES



CelebAHQ

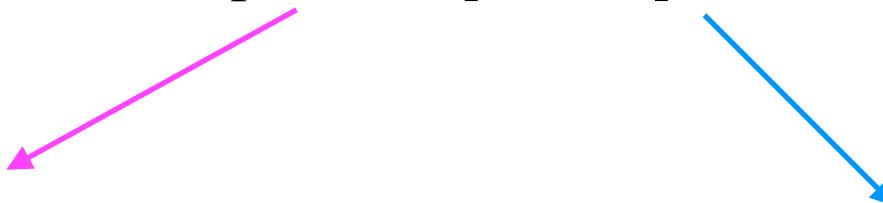
GLOW: LATENT INTERPOLATION



CelebAHQ

VAES WITH NORMALIZING FLOWS

$$q(\mathbf{z} \mid \mathbf{x}) \propto p(\mathbf{x} \mid \mathbf{z}) p(\mathbf{z})$$



Variational inference with normalizing flows

Rezende & Mohamed. "Variational inference with normalizing flows."

v.d. Berg, Hasenclever, Tomczak, Welling, "Sylvester normalizing flows for variational inference"

Kingma, Salimans, Jozefowicz, Chen, Sutskever, Welling "Improved variational inference with inverse autoregressive flow"

Tomczak, Welling, "Improving variational auto-encoders using householder flow"

Flow-based priors

Chen, Kingma, Salimans, Duan, Dhariwal, Schulman, Abbeel, "Variational lossy autoencoder"

Gatopoulos, Tomczak. "Self-Supervised Variational Auto-Encoders."

Thank you!