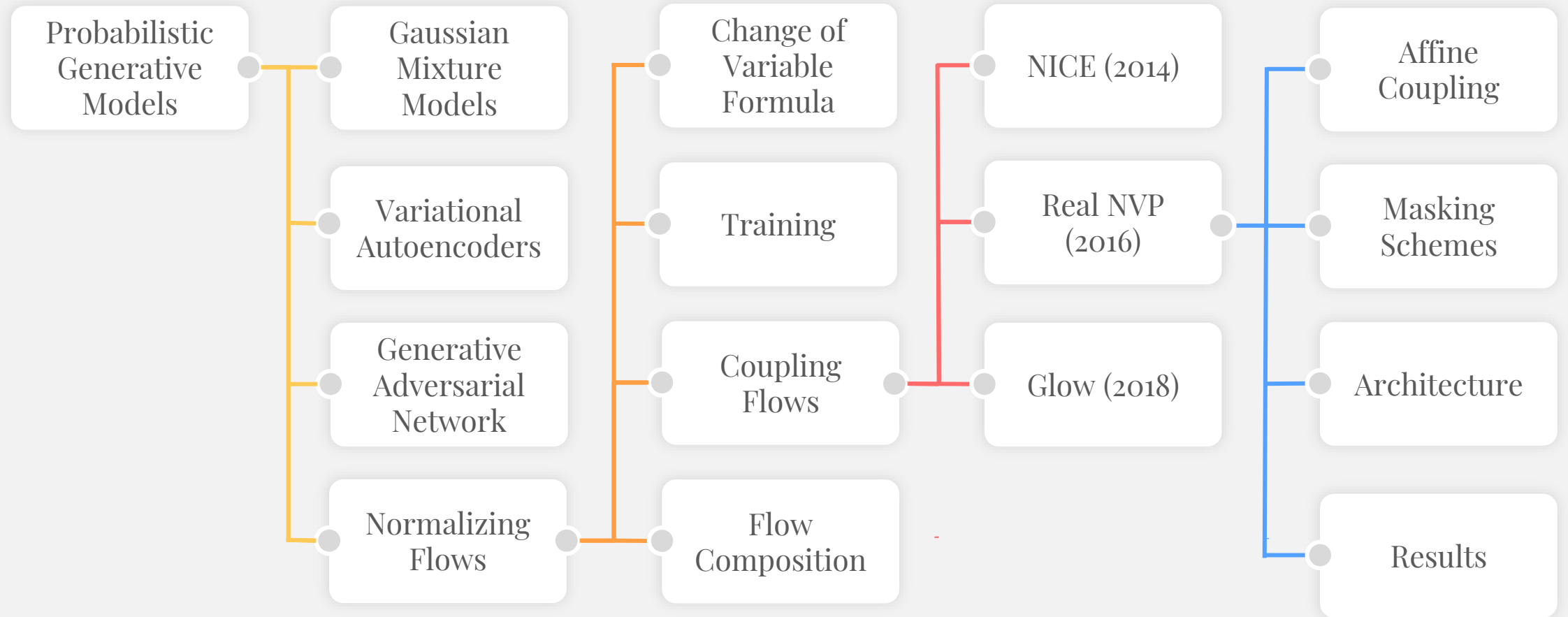


Density Estimation Using Real NVP

Laurent Dinh, et al. · 2016

By Flavio Schneider · March 2021 · **ETH** zürich

Map



Probabilistic Generative Model

Definition • 2 Models • Adding Latent z • Uses

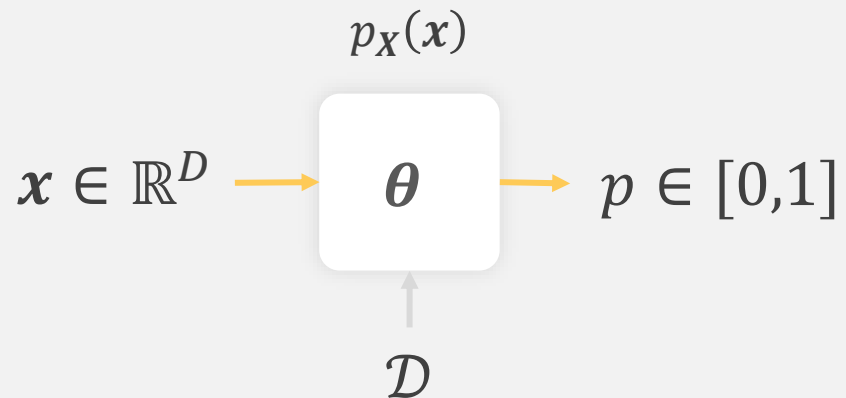
Def 1. A **probabilistic generative model (PGM)** is a model that attempts to estimate the probability distribution $p_X(\mathbf{x})$ parametrized by $\boldsymbol{\theta}$ over a random variable \mathbf{X} given a dataset $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with $\mathbf{x}_i \in \mathbb{R}^D$.

Probabilistic Generative Model

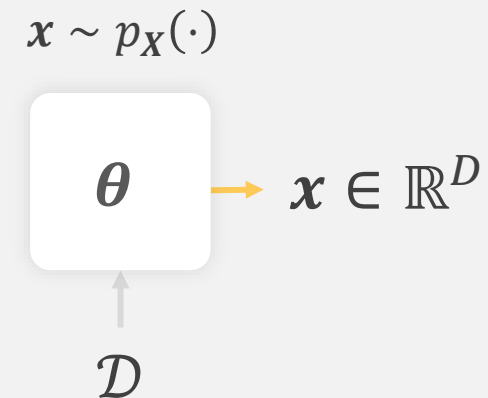
Definition • 2 Models • Adding Latent z • Uses

Def 1. A **probabilistic generative model (PGM)** is a model that attempts to estimate the probability distribution $p_X(\mathbf{x})$ parametrized by θ over a random variable \mathbf{X} given a dataset $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ with $\mathbf{x}_i \in \mathbb{R}^D$.

Density Estimation



Sampling

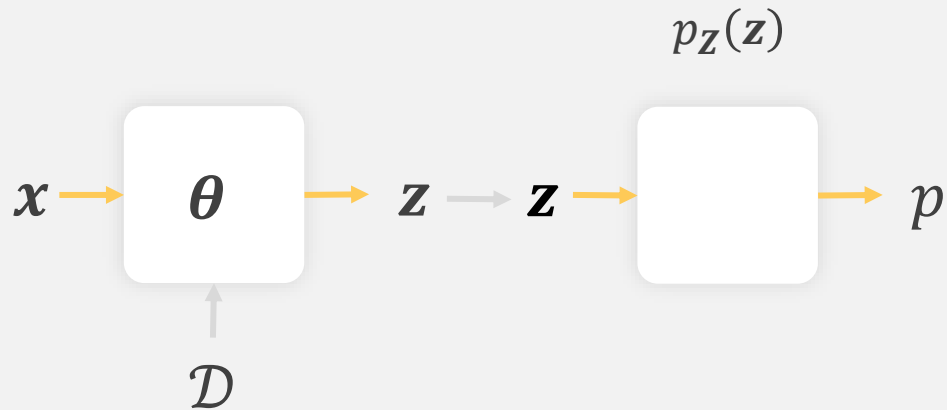


Probabilistic Generative Model

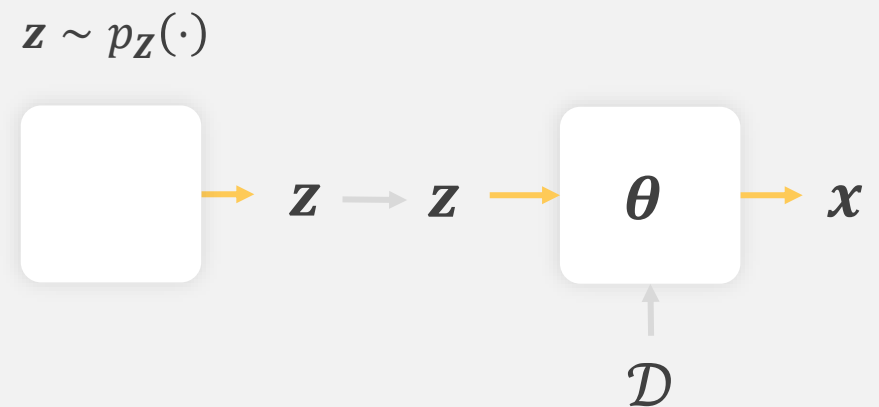
Definition • 2 Models • Adding Latent \mathbf{z} • Uses

Add an additional variable $\mathbf{z} \in \mathbb{R}^{D'}$ with some known simple distribution $p_{\mathbf{z}}(\mathbf{z})$ (e.g. $p_{\mathbf{z}}(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, \mathbf{I})$) which is easy to evaluate and sample from.

Density Estimation



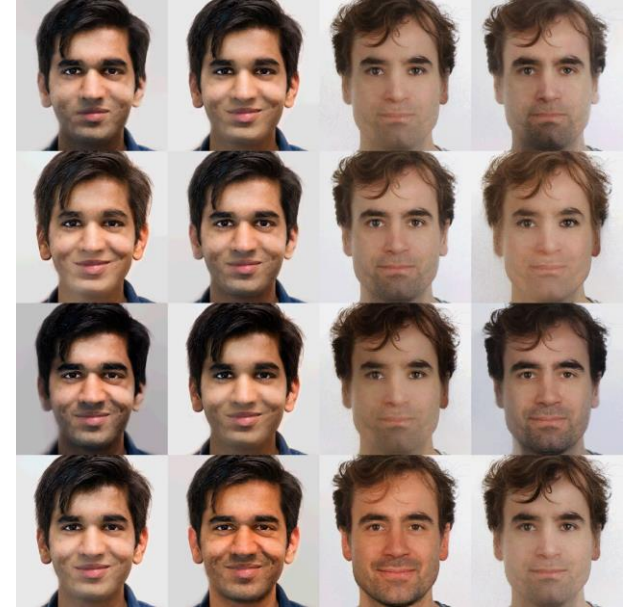
Sampling



Probabilistic Generative Model

Definition • 2 Models • Adding Latent z • Uses

- Generate data (e.g., images) like the samples in \mathcal{D}
- Interpolate between samples in a meaningful way
- Unsupervised clustering
- Denoising
- Data reconstruction
- Build a prior $p_X(\mathbf{x})$ to some other model (Bayesian setting)



(Gaussian) Mixture Model

- Used mostly for unsupervised clustering.
- Parameters $\theta = \{\pi, \mu_{z_1, \dots, z_k}, \Sigma_{z_1, \dots, z_k}\}$ estimated with EM algorithm.
- + Density estimation is easy: $p_X(x) = \sum_z p_{X|Z}(x|z)p_Z(z)$
- + Sampling is easy: $z \sim p_Z(\cdot)$ then $x \sim p_{X|Z}(\cdot | z)$
- EM requires initial guess for θ .
- Scales poorly with large D .
- Number of clusters k is fixed.

$k \in \mathbb{N}$: number of clusters

$x \in \mathbb{R}^D$: data point

$z \in \{0,1\}^k$: cluster membership

$$p_{Z|X}(z|x) = \frac{p_{X|Z}(x|z)p_Z(z)}{p_X(x)}$$

- $p_{X|Z}(x|z) = \mathcal{N}(x|\mu_z, \Sigma_z)$
- $p_Z(z) = \text{Cat}(z|\pi_1, \dots, \pi_k)$

Variational Autoencoder

- Used for compression, data generation/reconstruction/denoising.
- Parameters θ estimated by optimizing ELBO.

± Density estimation is approximate: $p_X(\mathbf{x}) \approx \frac{p_{X|Z}(\mathbf{x}|\mathbf{z})p_Z(\mathbf{z})}{q_{Z|X}(\mathbf{z}|\mathbf{x})}$

+ Sampling is easy: $\mathbf{z} \sim p_Z(\cdot)$ then $\mathbf{x} \sim p_{X|Z}(\cdot|\mathbf{z})$

$\mathbf{x} \in \mathbb{R}^D$: data point

$\mathbf{z} \in \mathbb{R}^{D'}$: latent data point

$$p_{Z|X}(\mathbf{z}|\mathbf{x}) = \frac{p_{X|Z}(\mathbf{x}|\mathbf{z})p_Z(\mathbf{z})}{p_X(\mathbf{x})}$$

- $p_Z(\mathbf{z}) = \mathcal{N}(\mathbf{z}|\mathbf{0}, I)$

$$q_{Z|X}(\mathbf{z}|\mathbf{x}) \approx p_{Z|X}(\mathbf{z}|\mathbf{x})$$

Generative Adversarial Networks

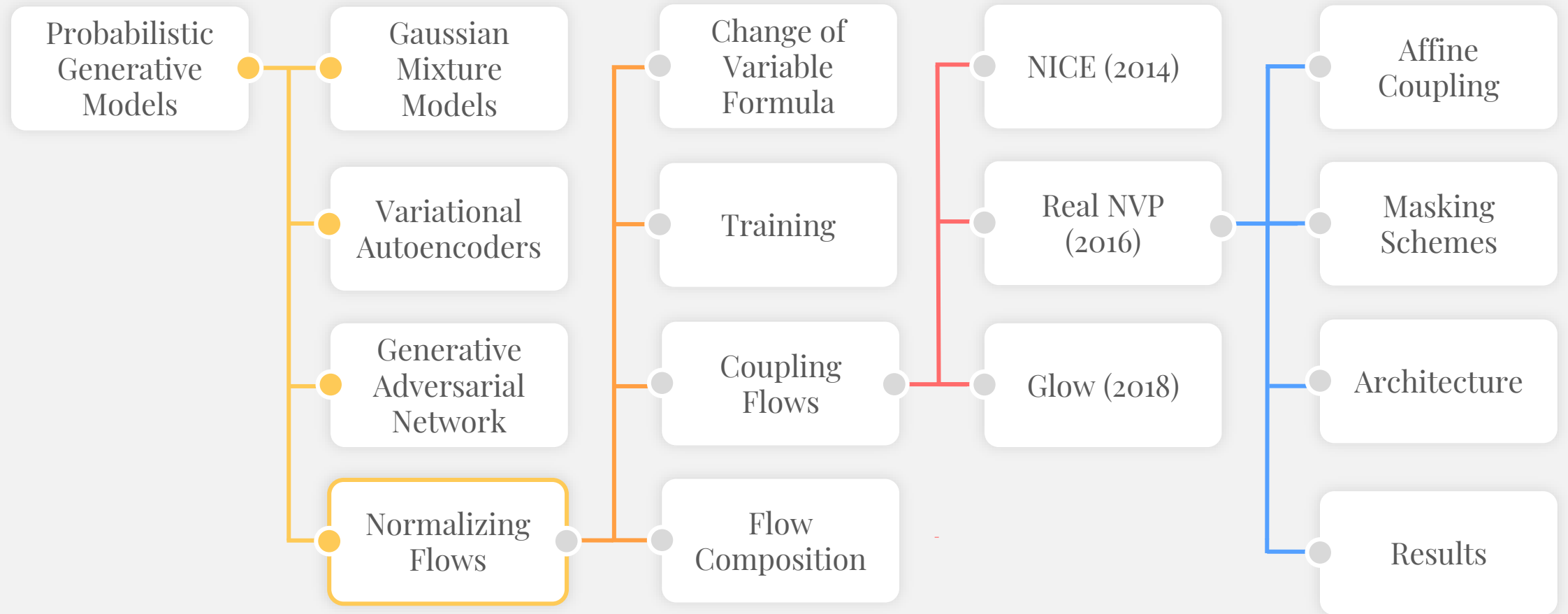
- Used mainly for data generation.
- Parameters θ estimated by playing the minimax game.
- + Sampling is easy (built for this): $\mathbf{z} \sim p_{\mathbf{z}}(\cdot)$ then $\mathbf{x} = \mathbf{G}(\mathbf{z}|\theta_{\mathbf{G}})$
- Density estimation $p_{\mathbf{x}}(\mathbf{x})$ is not possible.

$\mathbf{x} \in \mathbb{R}^D$: data point

$\mathbf{z} \in \mathbb{R}^{D'}$: latent data point

$$\min_G \max_D V(D, G)$$

Map



Normalizing Flow

Introduction • Change of Variable • Diagram • Example

- A PGM to directly approximate $p_{\mathbf{x}}(\mathbf{x})$.
- Turn a simple distribution $p_{\mathbf{z}}(\mathbf{z})$ into a complex one in an *invertible manner*.
- + Easier to train than GANs and VAEs.
- + Easy to scale.
- + We can directly compute and optimize the likelihood.
- + Density estimation $p_{\mathbf{x}}(\mathbf{x})$ is easy.
- + Sampling $\mathbf{x} \sim p_{\mathbf{x}}(\cdot)$ is easy.

Normalizing Flow

Introduction • Change of Variable • Diagram • Example

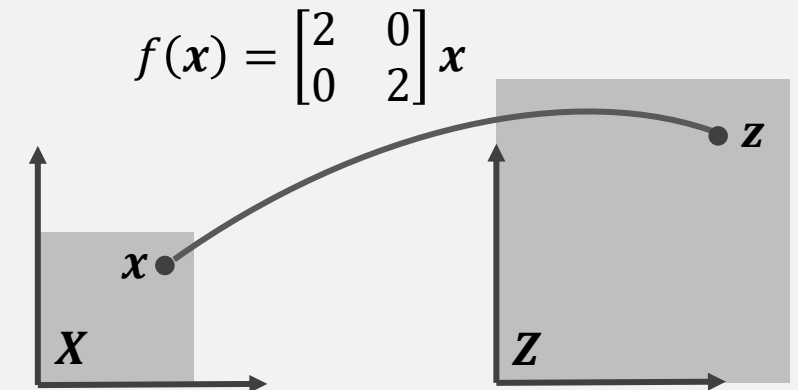
Lemma 1. (Change of Variable) Let \mathbf{X} and \mathbf{Z} be random variables related by an *invertible* and *differentiable* mapping $f: \mathbb{R}^D \rightarrow \mathbb{R}^D$ such that $\mathbf{Z} = f(\mathbf{X})$, then the following equality holds:

$$p_{\mathbf{X}}(\mathbf{x}) = p_{\mathbf{Z}}(f(\mathbf{x})) |\det(Df(\mathbf{x}))|$$

$\mathbf{x} \in \mathbb{R}^D$: data point

$\mathbf{z} \in \mathbb{R}^D$: “latent” variable

$$Df(\mathbf{x}) = \begin{bmatrix} \partial_{x_1} f_1(\mathbf{x}) & \cdots & \partial_{x_D} f_1(\mathbf{x}) \\ \vdots & \ddots & \vdots \\ \partial_{x_1} f_D(\mathbf{x}) & \cdots & \partial_{x_D} f_D(\mathbf{x}) \end{bmatrix}$$



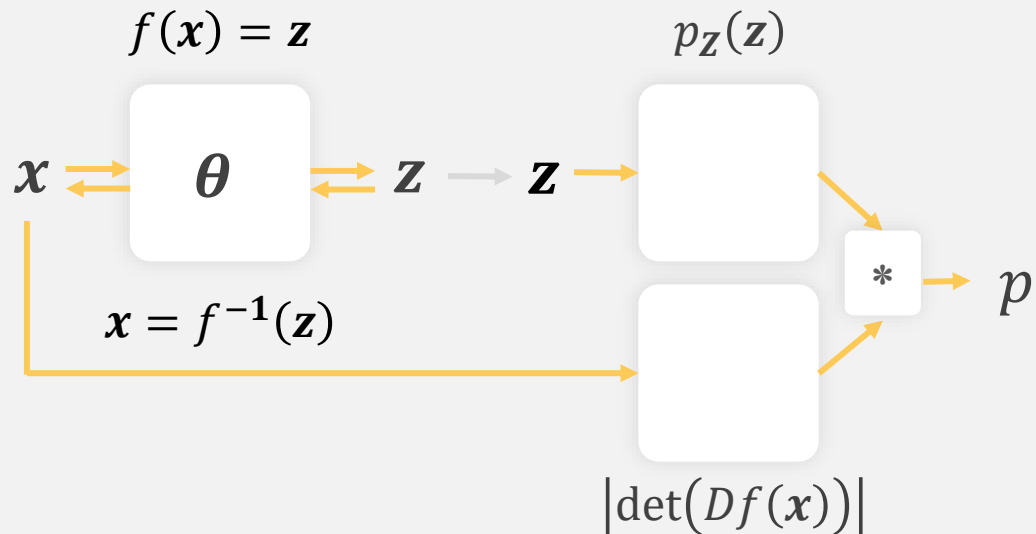
$$p_{\mathbf{X}}(\mathbf{x}) = p_{\mathbf{Z}}\left(\begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix} \mathbf{x}\right) \cdot 4$$

Normalizing Flow

Introduction · Change of Variable · **Diagram** · Example

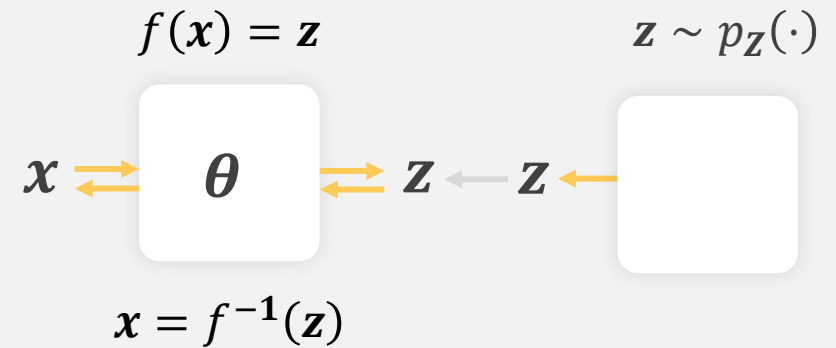
Density Estimation

$$p_X(\mathbf{x})$$



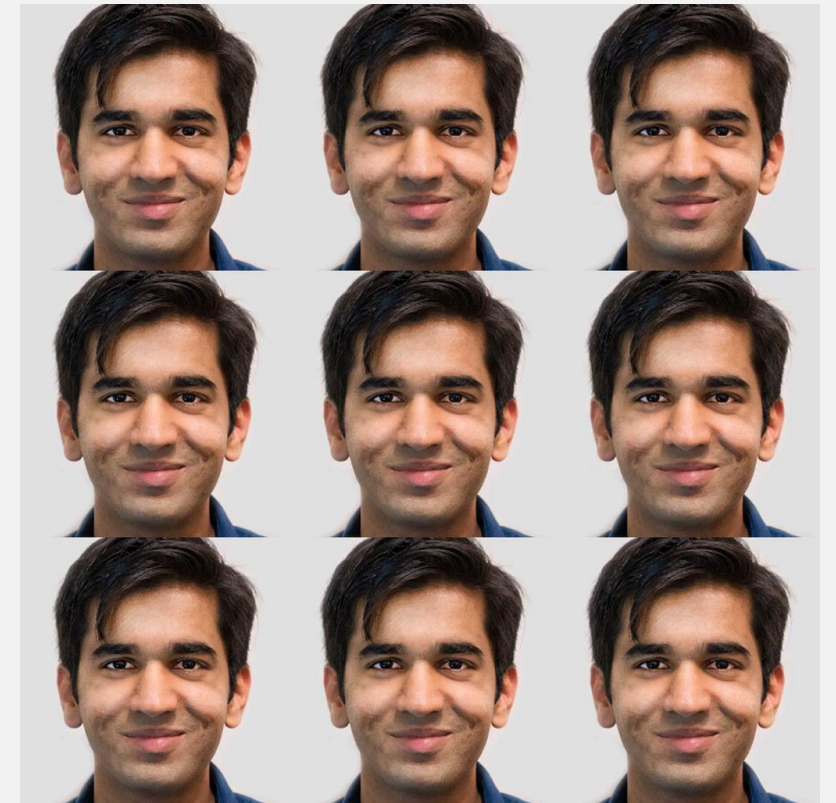
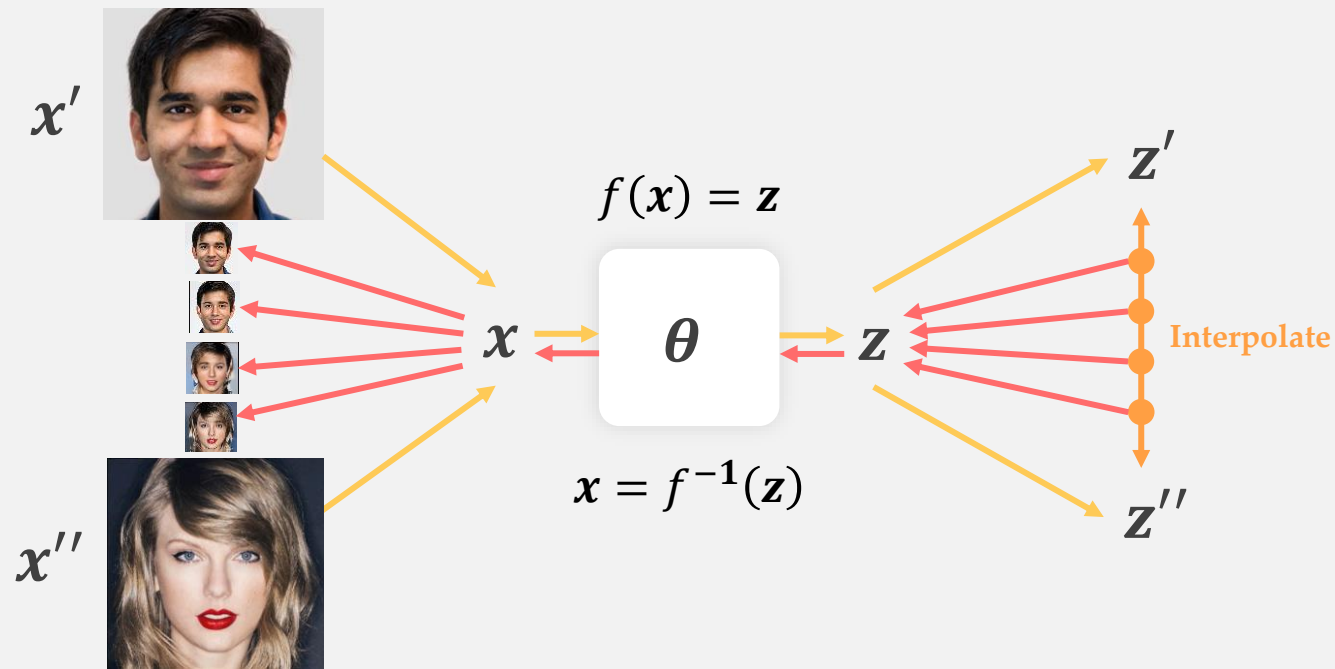
Sampling

$$\mathbf{x} \sim p_X(\cdot)$$



Normalizing Flow

Introduction • Change of Variable • Diagram • Examples



"Glow: Generative flow with invertible 1x1 convolutions." (2018).

Training

Maximum Likelihood • Problems

Given a dataset $\mathcal{D} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ where $\mathbf{x}_1, \dots, \mathbf{x}_N \sim_{i.i.d.} p_{\mathbf{X}}(\cdot)$ we want to find the best parameters $\boldsymbol{\theta}$ of $f(\mathbf{x}|\boldsymbol{\theta})$ that maximize the likelihood of \mathcal{D} .

$$\begin{aligned} L(\mathcal{D}|\boldsymbol{\theta}) &= L(\mathbf{x}_1, \dots, \mathbf{x}_N|\boldsymbol{\theta}) \\ &= \prod_{i=1}^N p_{\mathbf{X}}(\mathbf{x}_i|\boldsymbol{\theta}) \\ &= \prod_{i=1}^N p_{\mathbf{Z}}(f(\mathbf{x}_i|\boldsymbol{\theta})) |\det(Df(\mathbf{x}_i|\boldsymbol{\theta}))| \end{aligned}$$

$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} L(\mathcal{D}|\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \log L(\mathcal{D}|\boldsymbol{\theta}) \\ &= \arg \max_{\boldsymbol{\theta}} \log \prod_{i=1}^N p_{\mathbf{Z}}(f(\mathbf{x}_i|\boldsymbol{\theta})) |\det(Df(\mathbf{x}_i|\boldsymbol{\theta}))| \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(f(\mathbf{x}_i|\boldsymbol{\theta})) + \log |\det(Df(\mathbf{x}_i|\boldsymbol{\theta}))| \end{aligned}$$

Training

Maximum Likelihood • Problems

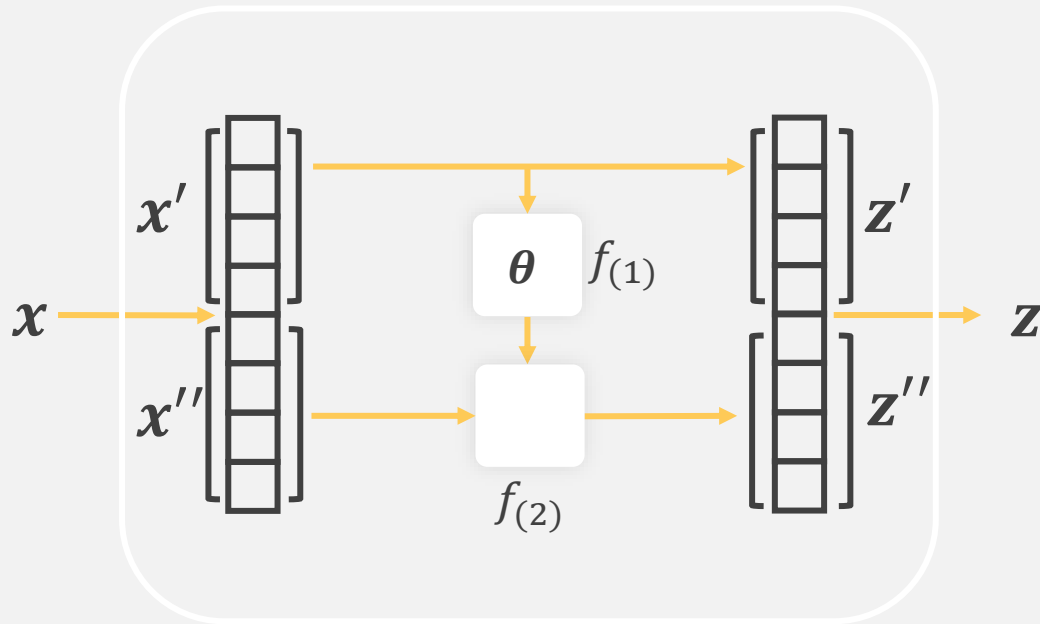
- Problem 1: How do we enforce the fact that f must be invertible (standard NNs aren't) ?
- Problem 2: How do we compute the determinant which is usually computationally prohibitive for large D ?

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(\textcolor{red}{f}(\mathbf{x}_i|\boldsymbol{\theta})) + \log |\textcolor{red}{det}(Df(\mathbf{x}_i|\boldsymbol{\theta}))|$$

Coupling Flows

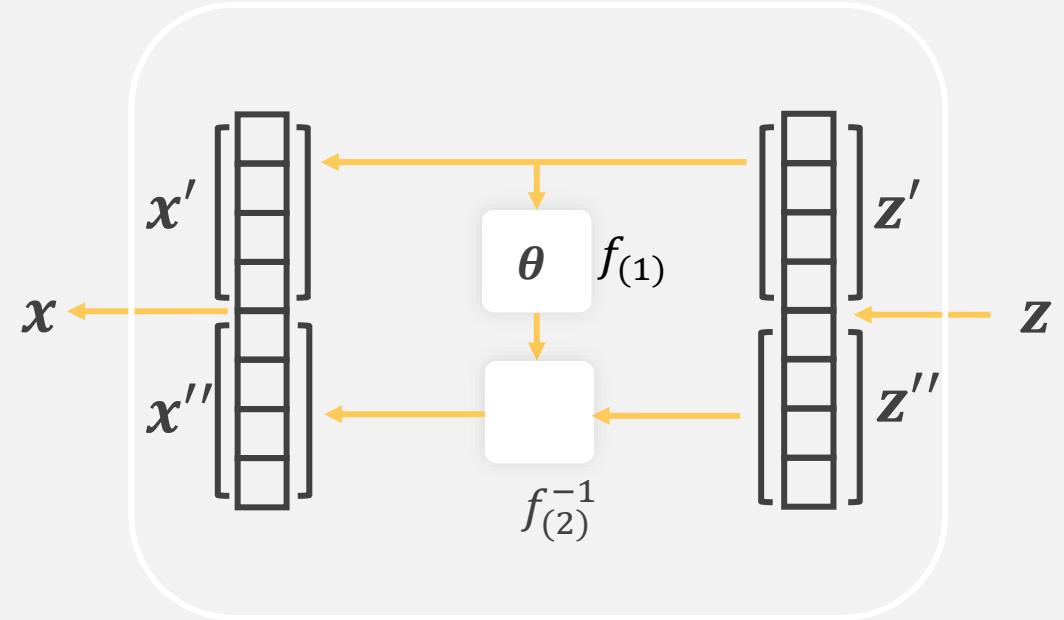
Solution To Problem 1 • Solution To Problem 2

$$f(x | \theta) = z$$



$$f(x | \theta) = \left(x', f_{(2)} \left(x'', f_{(1)}(x' | \theta) \right) \right)$$

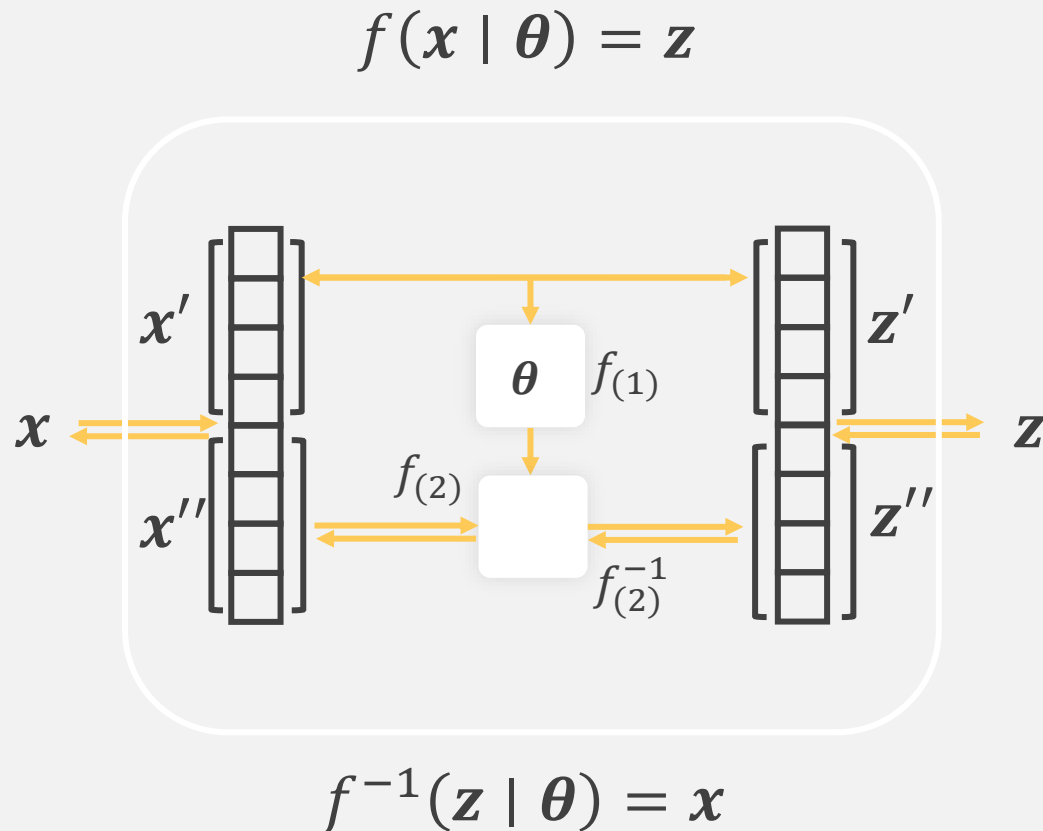
$$f^{-1}(z | \theta) = x$$



$$f^{-1}(z | \theta) = \left(z', f_{(2)}^{-1} \left(z'', f_{(1)}(z' | \theta) \right) \right)$$

Coupling Flows

Solution To Problem 1 • Solution To Problem 2



$$f_{(1)}(\cdot \mid \boldsymbol{\theta}) = \mathbf{NN}(\cdot \mid \boldsymbol{\theta}) = \mathbf{t}$$

Assuming...

$$f_{(2)}(\mathbf{x}'', \mathbf{t}) = \mathbf{x}'' + \mathbf{t}$$

$$f_{(2)}^{-1}(\mathbf{z}'', \mathbf{t}) = \mathbf{z}'' - \mathbf{t}$$

We found a way to compute f in an invertible manner!

“NICE: NON-LINEAR INDEPENDENT COMPONENTS ESTIMATION (2015).

Coupling Flows

Solution To Problem 1 • Solution To Problem 2

- ~~Problem 1: How do we enforce the fact that f must be invertible (standard NNs aren't)?~~
- Problem 2: How do we compute the determinant which is usually computationally prohibitive for large D ?

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(\textcolor{teal}{f}(\mathbf{x}_i|\boldsymbol{\theta})) + \log |\textcolor{red}{\det}(Df(\mathbf{x}_i|\boldsymbol{\theta}))|$$

Coupling Flows

Solution To Problem 1 • Solution To Problem 2

Lemma 2. (Determinant LT Matrix) Let $\mathbf{A} \in \mathbb{R}^{D \times D}$ be a lower triangular matrix, then:

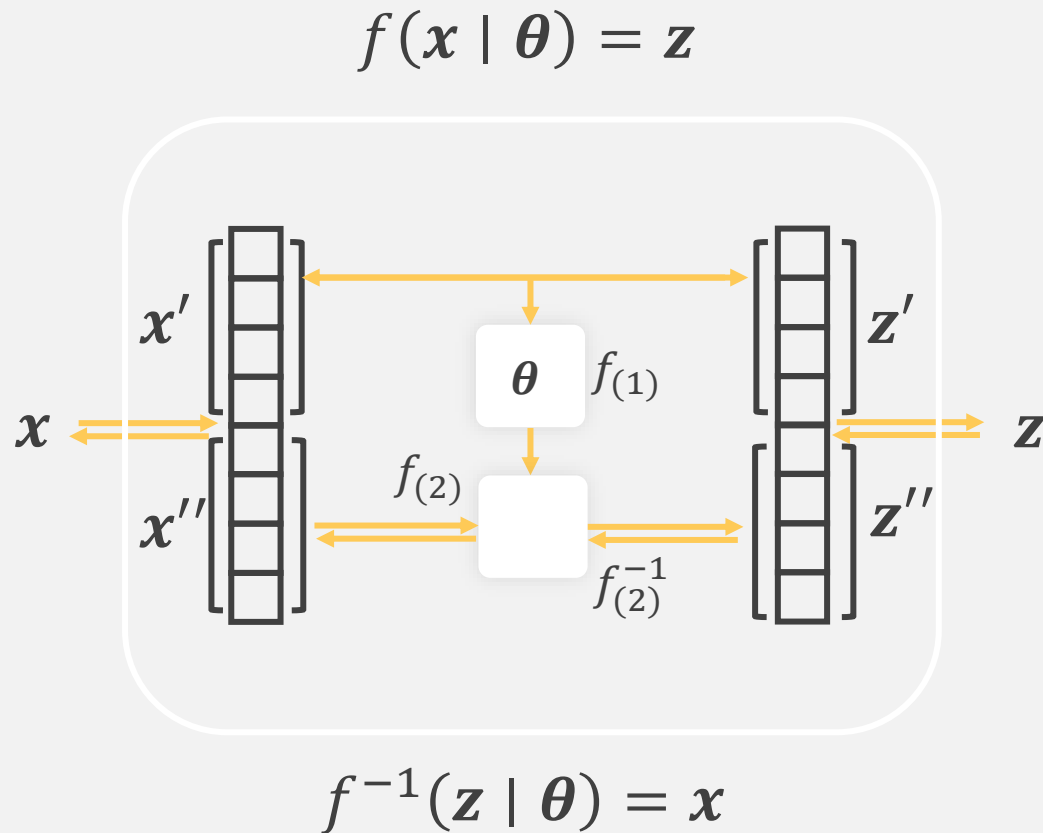
$$\det \mathbf{A} = \prod_{i=1}^D a_{ii}$$

$$\text{E.g.: } \det \begin{bmatrix} a_{11} & 0 & 0 \\ * & a_{22} & 0 \\ * & * & a_{33} \end{bmatrix} = a_{11} a_{22} a_{33}$$

Coupling Flows

$$|\det(Df(x_i|\theta))|$$

Solution To Problem 1 • Solution To Problem 2



Derivatives

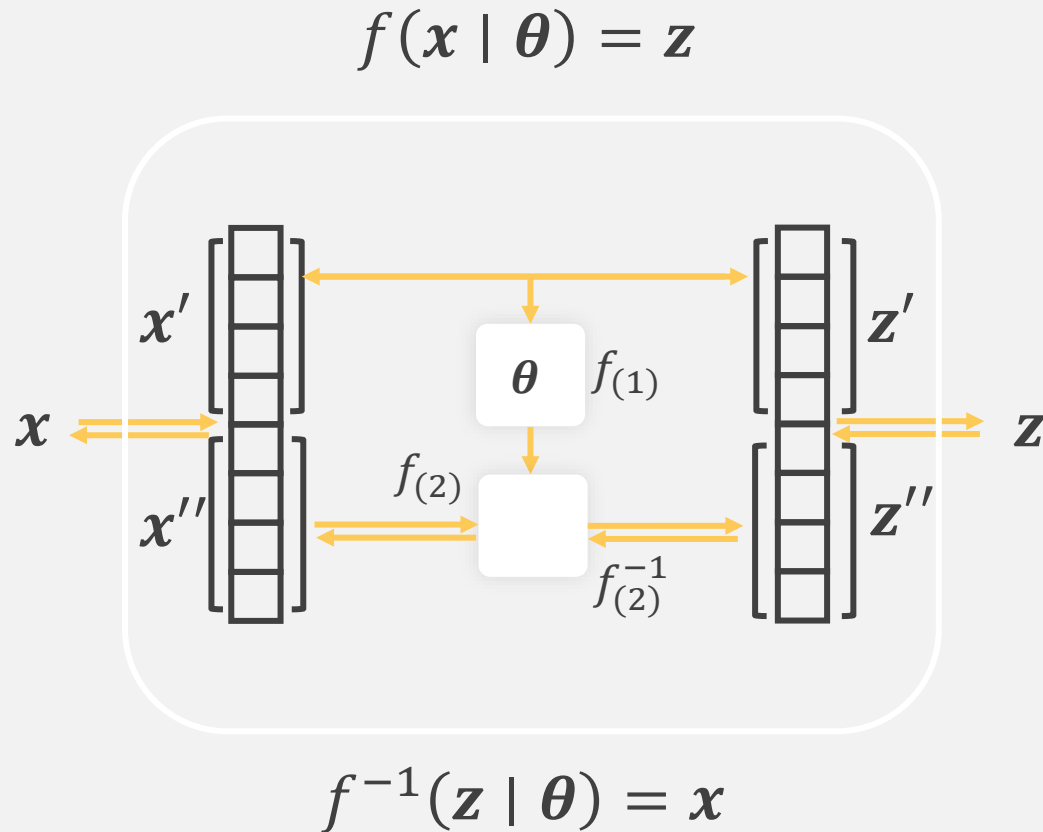
- How much does z' change if we change x'' ?

$$Df(x) = \begin{bmatrix} * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix}$$

Coupling Flows

$$|\det(Df(\mathbf{x}_i | \boldsymbol{\theta}))|$$

Solution To Problem 1 • Solution To Problem 2



Derivatives

- How much does \mathbf{z}' change if we change \mathbf{x}'' ?
- How much does \mathbf{z}' change if we change \mathbf{x}' ?

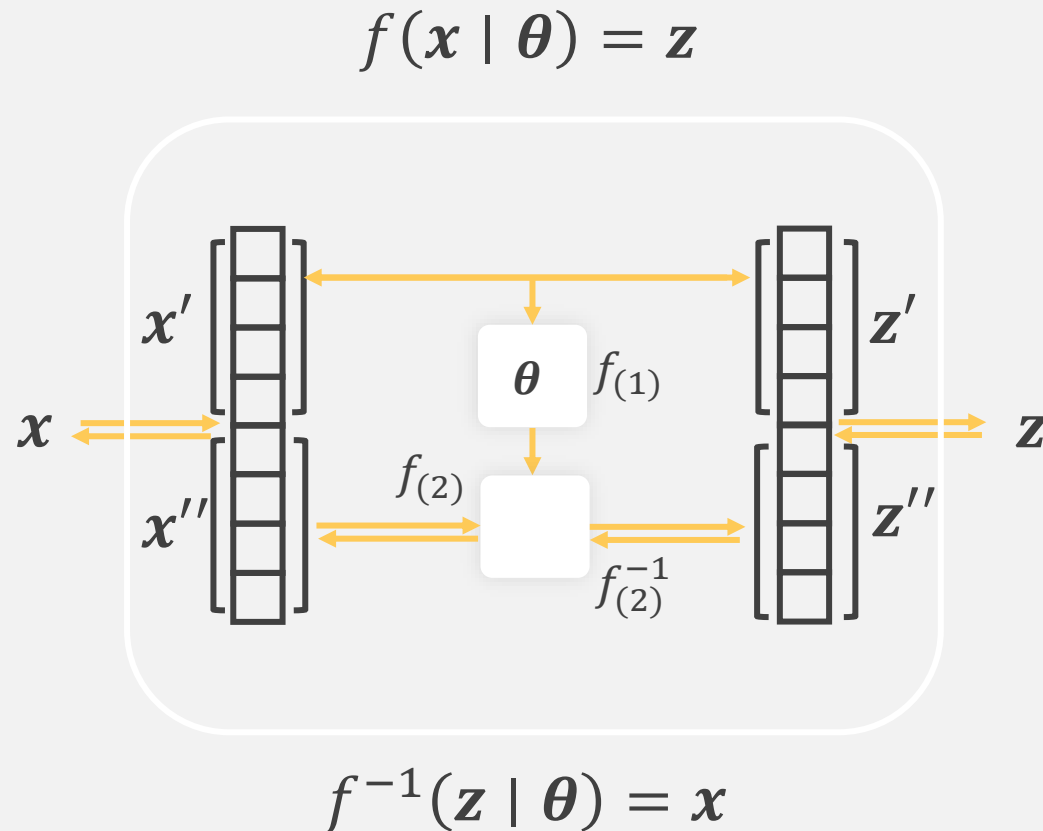
$$Df(\mathbf{x}) = \begin{bmatrix} * & * & * & 0 & 0 & 0 \\ * & * & * & 0 & 0 & 0 \\ * & * & * & 0 & 0 & 0 \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix} \begin{matrix} \mathbf{z}' \\ \mathbf{z}'' \end{matrix}$$

$\underbrace{\hspace{1.5cm}}_{\mathbf{x}'} \quad \underbrace{\hspace{1.5cm}}_{\mathbf{x}''}$

Coupling Flows

$$|\det(Df(\mathbf{x}_i|\boldsymbol{\theta}))|$$

Solution To Problem 1 • Solution To Problem 2



Derivatives

- How much does \mathbf{z}' change if we change \mathbf{x}'' ?
- How much does \mathbf{z}' change if we change \mathbf{x}' ?
- How much does \mathbf{z}'' change if we change \mathbf{x}'' ?

Assuming $f_{(2)}(\mathbf{x}'', \mathbf{t}) = \mathbf{x}'' + \mathbf{t}$

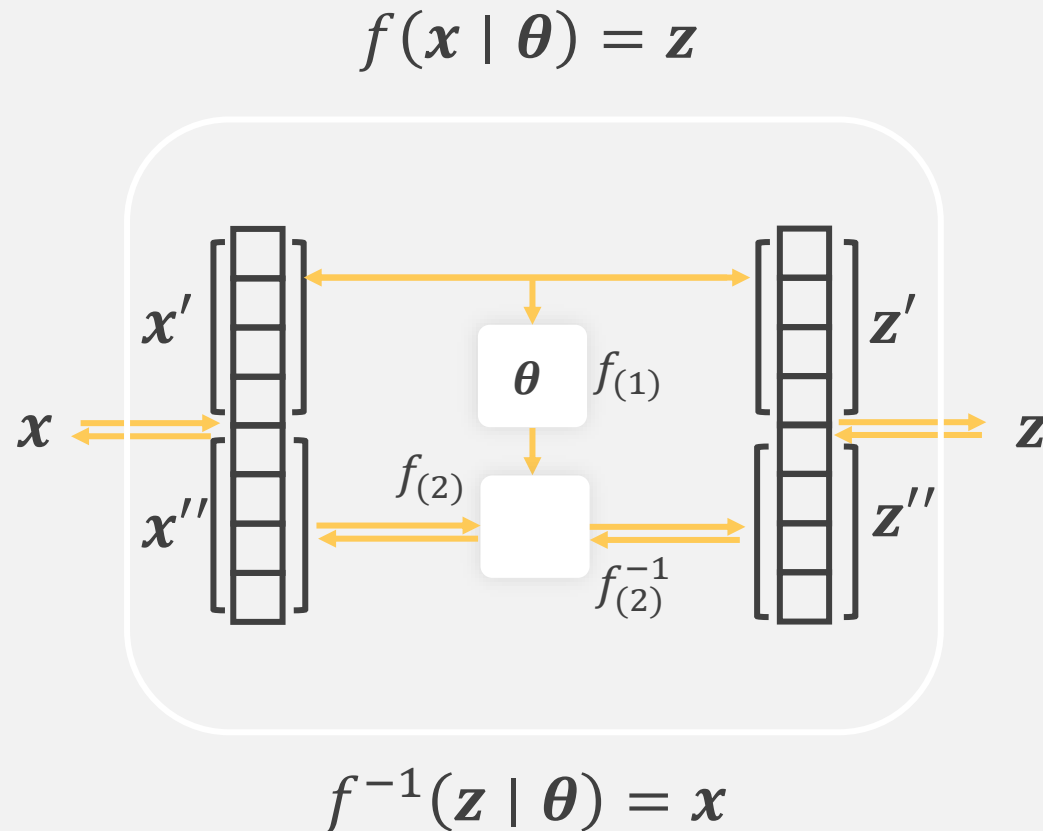
$$Df(\mathbf{x}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ * & * & * & * & * & * \\ * & * & * & * & * & * \\ * & * & * & * & * & * \end{bmatrix} \begin{matrix} \mathbf{z}' \\ \mathbf{z}'' \end{matrix}$$

$\underbrace{\quad}_{\mathbf{x}'} \quad \underbrace{\quad}_{\mathbf{x}''}$

Coupling Flows

$$|\det(Df(\mathbf{x}_i|\boldsymbol{\theta}))|$$

Solution To Problem 1 • Solution To Problem 2



Derivatives

- How much does \mathbf{z}' change if we change \mathbf{x}'' ?
- How much does \mathbf{z}' change if we change \mathbf{x}' ?
- How much does \mathbf{z}'' change if we change \mathbf{x}'' ?

Assuming $f_{(2)}(\mathbf{x}'', \mathbf{t}) = \mathbf{x}'' + \mathbf{t}$

$$Df(\mathbf{x}) = \begin{bmatrix} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ \hline * & * & * & 1 & 0 & 0 \\ * & * & * & 0 & 1 & 0 \\ * & * & * & 0 & 0 & 1 \end{bmatrix} \begin{matrix} \mathbf{z}' \\ \mathbf{z}'' \end{matrix}$$

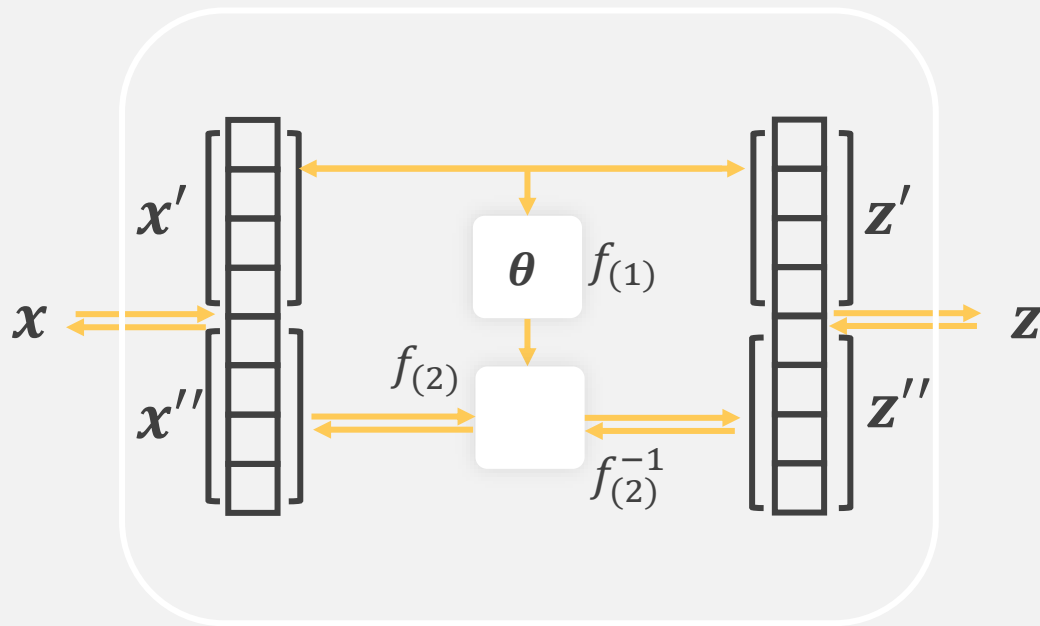
$\underbrace{\begin{matrix} * & * & * \end{matrix}}_{\mathbf{x}'} \quad \underbrace{\begin{matrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{matrix}}_{\mathbf{x}''}$

Coupling Flows

$$|\det(Df(x_i|\theta))|$$

Solution To Problem 1 • Solution To Problem 2

$$f(x | \theta) = z$$



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

Assuming $f_{(2)}(x'', t) = x'' + t$

$$Df(x) = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ \hline * & * & * & 1 & 0 & 0 \\ * & * & * & 0 & 1 & 0 \\ * & * & * & 0 & 0 & 1 \end{array} \right] \begin{array}{l} z' \\ z'' \end{array}$$

$x' \quad x''$

The Jacobian is lower triangular; hence we can easily compute the determinant!

Coupling Flows

Solution To Problem 1 • Solution To Problem 2

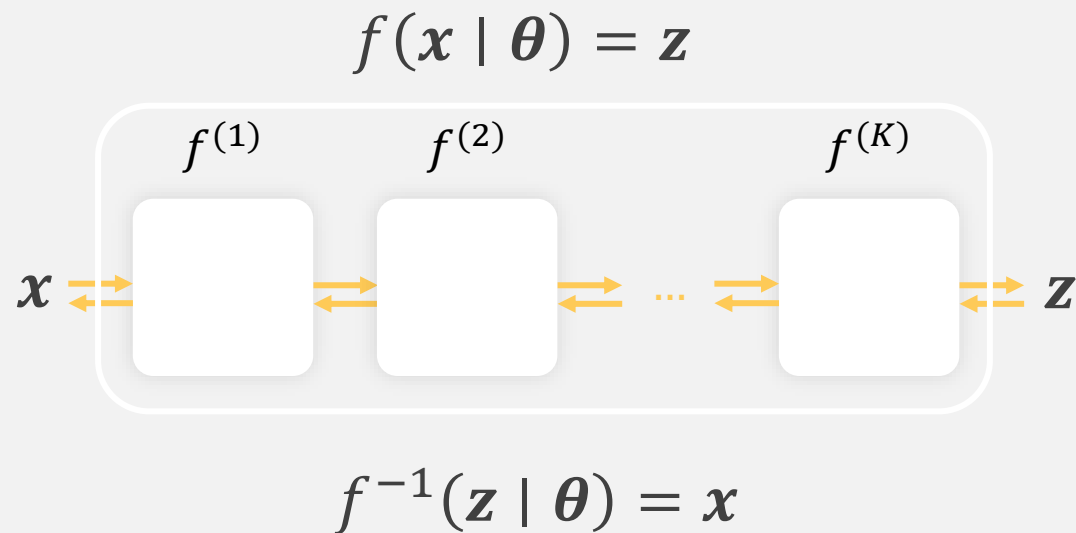
- ~~Problem 1: How do we enforce the fact that f must be invertible (standard NNs aren't)?~~
- ~~Problem 2: How do we compute the determinant which is usually computationally prohibitive for large D ?~~

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(\textcolor{teal}{f}(\mathbf{x}_i|\boldsymbol{\theta})) + \log|\textcolor{teal}{det}(D\textcolor{teal}{f}(\mathbf{x}_i|\boldsymbol{\theta}))|$$

Flow Composition

Diagram • Determinant Jacobian Lemma

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(f(\mathbf{x}_i | \boldsymbol{\theta})) + \log |\det(Df(\mathbf{x}_i | \boldsymbol{\theta}))|$$



Flow Composition

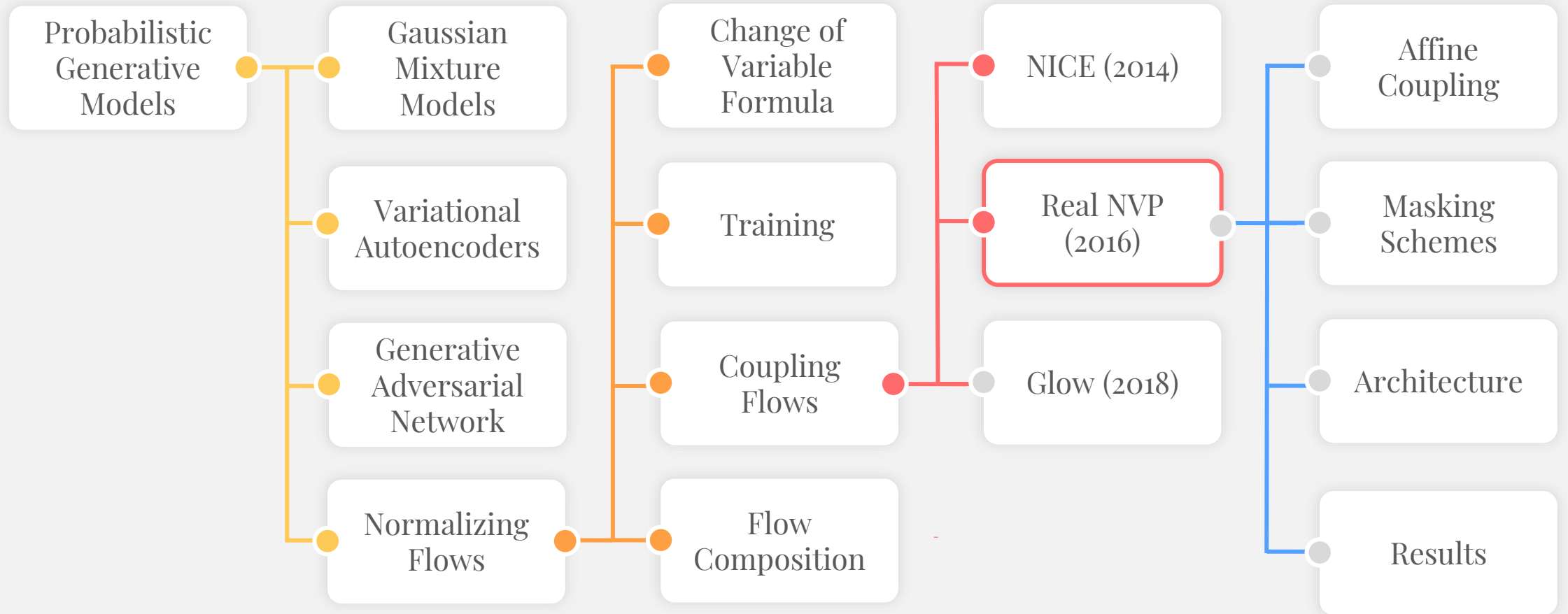
Diagram • Determinant Jacobian Lemma

Lemma 3. (Determinant Jacobian Composition) Let $f = f^{(1)} \circ f^{(2)} \circ \dots \circ f^{(K)}$ with $f^{(j)} : \mathbb{R}^D \rightarrow \mathbb{R}^D$, then:

$$\det Df(\mathbf{x}) = \det \prod_{j=1}^K Df^{(j)}(\mathbf{x}) = \prod_{j=1}^K \det Df^{(j)}(\mathbf{x})$$

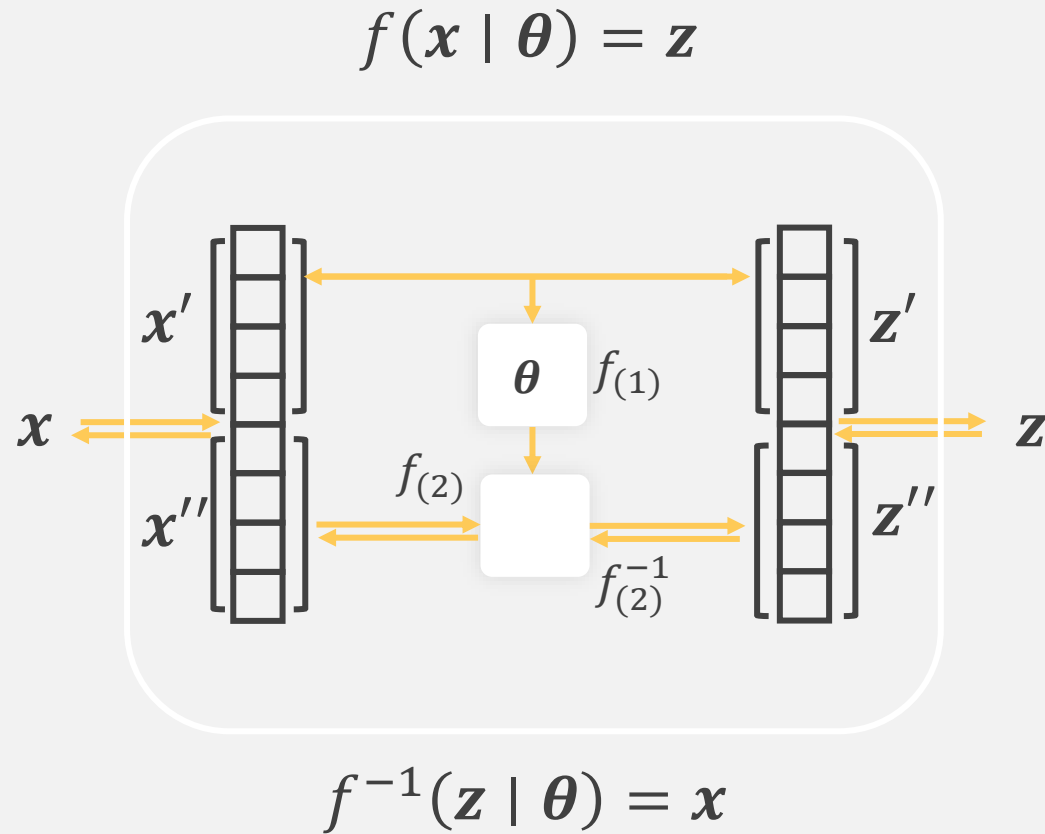
$$\begin{aligned} \hat{\boldsymbol{\theta}} &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(f(\mathbf{x}_i | \boldsymbol{\theta})) + \log |\det(Df(\mathbf{x}_i | \boldsymbol{\theta}))| \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(f(\mathbf{x}_i | \boldsymbol{\theta})) + \log \left| \prod_{j=1}^K \det Df^{(j)}(\mathbf{x}) \right| \\ &= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p_{\mathbf{Z}}(f(\mathbf{x}_i | \boldsymbol{\theta})) + \sum_{j=1}^K \log |\det Df^{(j)}(\mathbf{x})| \end{aligned}$$

Map



RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results



$$f_{(1)}(\cdot \mid \boldsymbol{\theta}) = \text{NN}(\cdot \mid \boldsymbol{\theta}) = (\mathbf{t}, \mathbf{s})$$

Assuming...

$$f_{(2)}(\mathbf{x}'', (\mathbf{t}, \mathbf{s})) = \mathbf{x}'' \odot \exp(\mathbf{s}) + \mathbf{t}$$

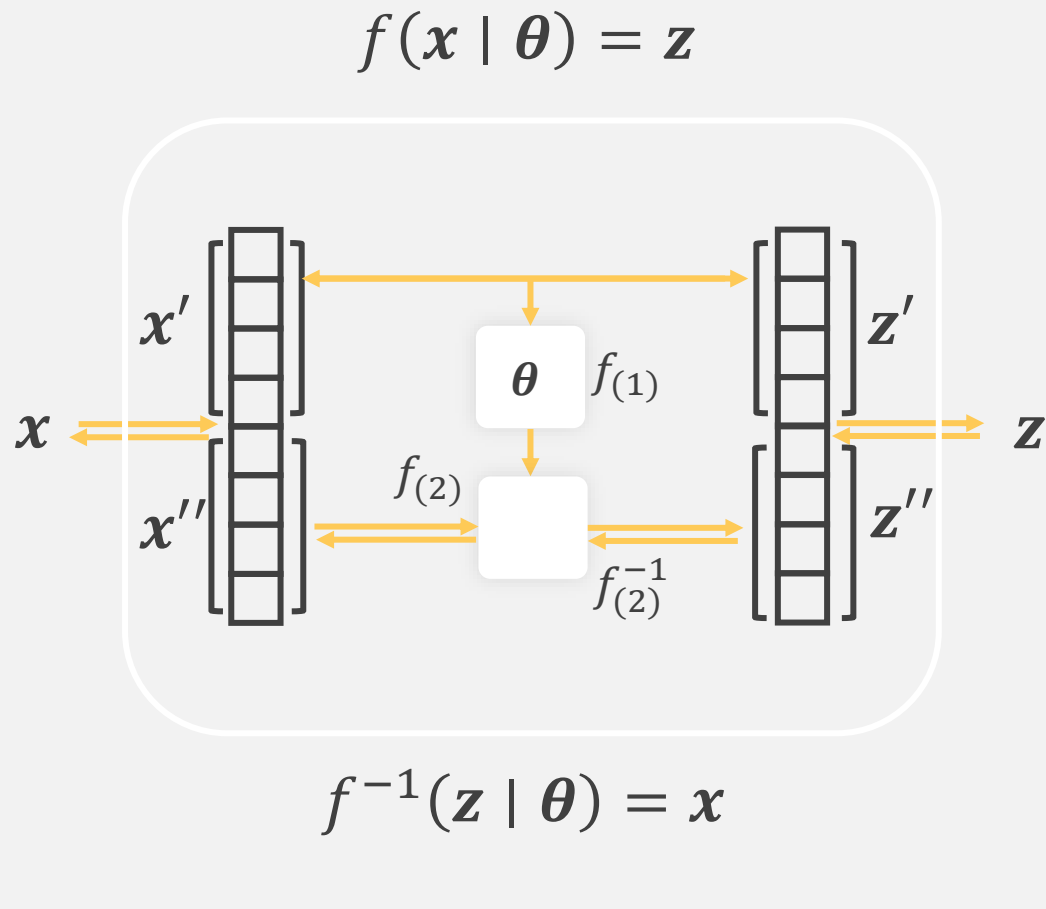
$$f_{(2)}^{-1}(\mathbf{z}'', (\mathbf{t}, \mathbf{s})) = (\mathbf{z}'' - \mathbf{t}) \odot \exp(-\mathbf{s})$$

$$Df(\mathbf{x}) = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ \hline * & * & * & \exp(s_1) & 0 & 0 \\ * & * & * & 0 & \ddots & 0 \\ * & * & * & 0 & 0 & \exp(s_D) \end{array} \right] \begin{array}{l} \mathbf{z}' \\ \mathbf{z}'' \end{array}$$

$\underbrace{\quad}_{\mathbf{x}'} \quad \underbrace{\quad}_{\mathbf{x}''}$

RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results



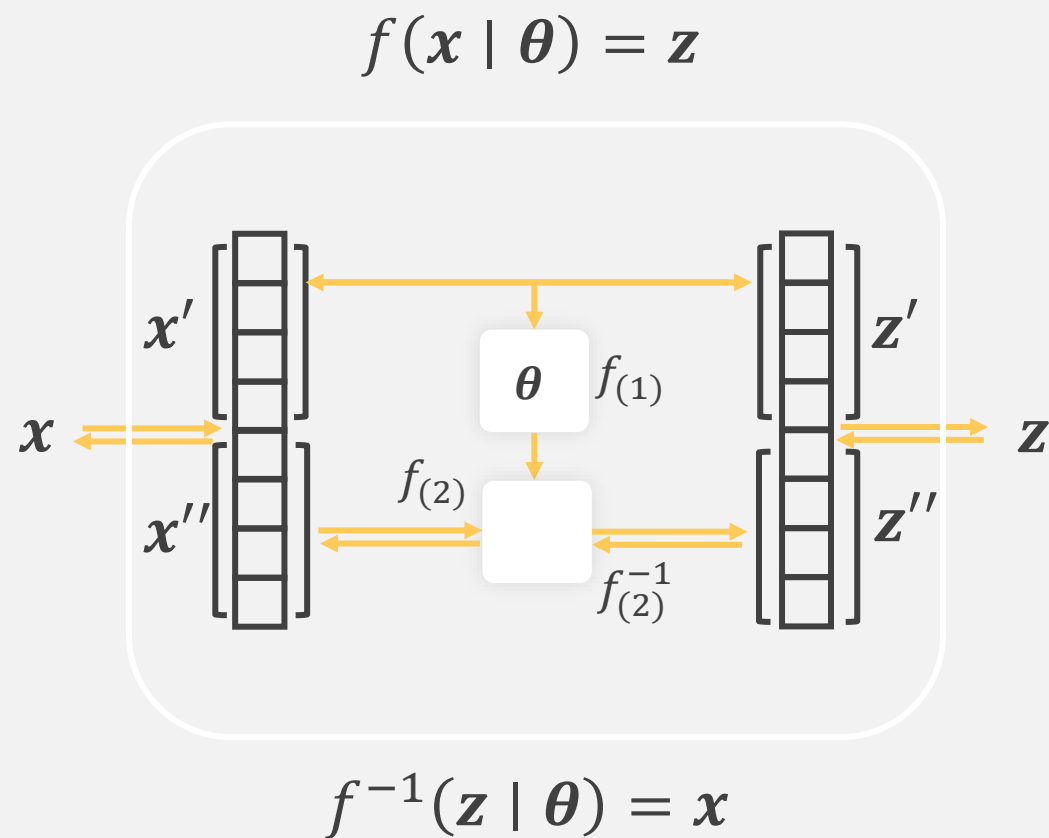
$$\log|\det Df(\mathbf{x})| = \log \left| \prod_{i=1}^D \exp s_i \right| = \log \left| \exp \sum_{i=1}^D s_i \right| = \sum_{i=1}^D s_i$$

$$Df(\mathbf{x}) = \left[\begin{array}{ccc|ccc} 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 \\ \hline * & * & * & \exp(s_1) & 0 & 0 \\ * & * & * & 0 & \ddots & 0 \\ * & * & * & 0 & 0 & \exp(s_D) \end{array} \right] \begin{array}{l} \mathbf{z}' \\ \mathbf{z}'' \end{array}$$

$\underbrace{\quad}_{\mathbf{x}'} \quad \underbrace{\quad}_{\mathbf{x}''}$

RealNVP (2017)

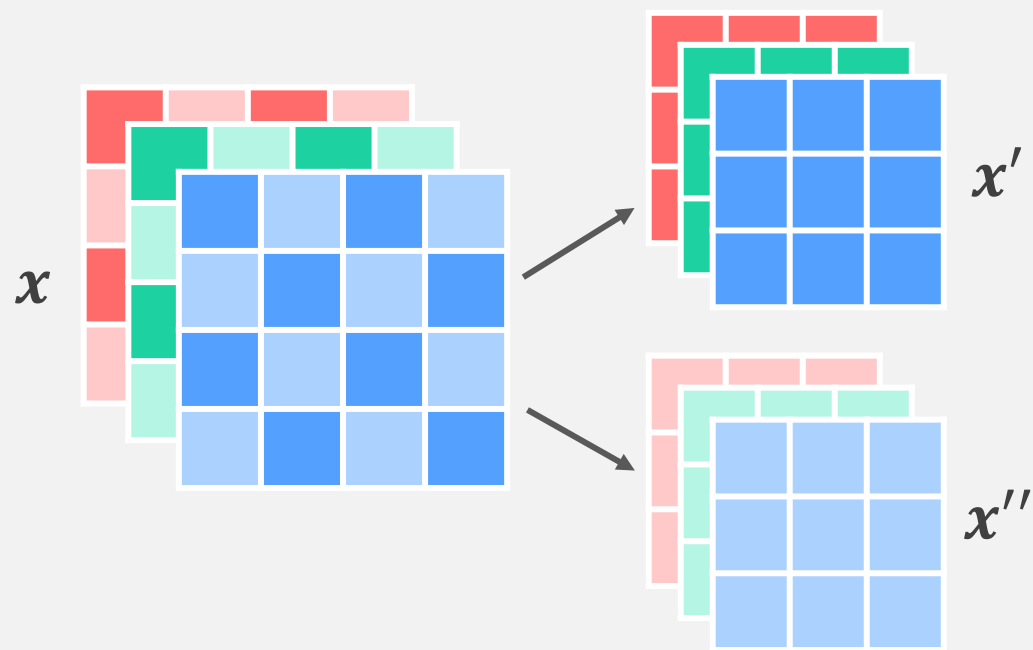
Affine Coupling Transform • Masking Schemes • Architecture • Results



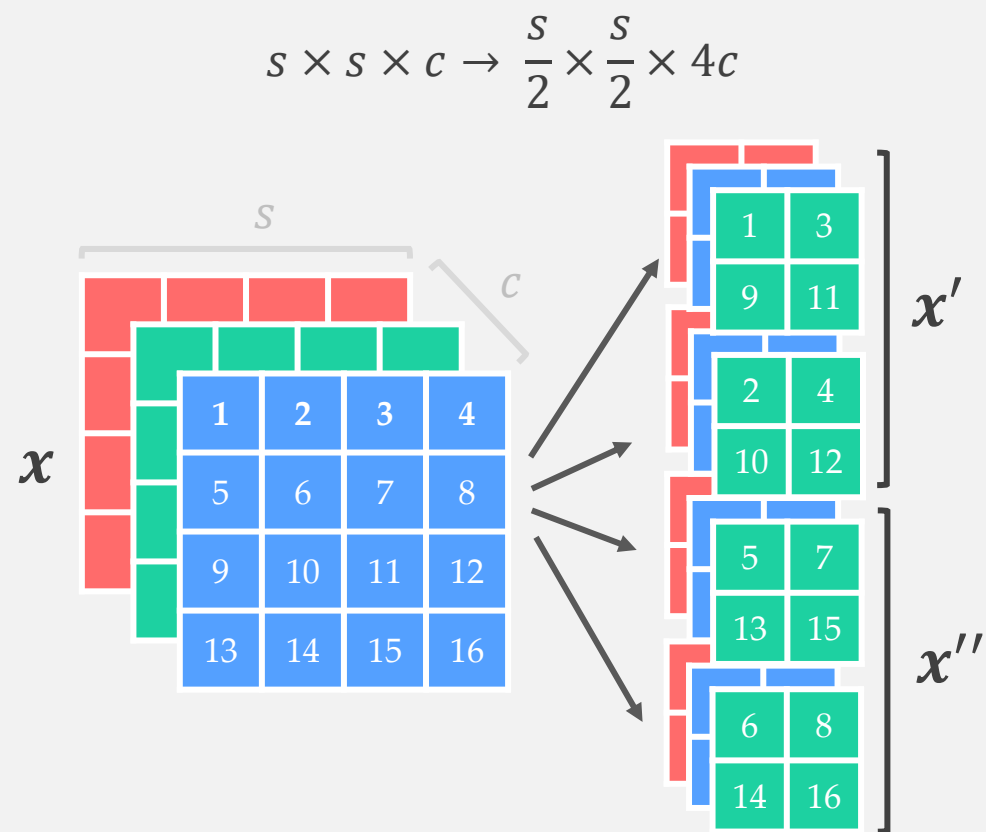
RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results

Checkboard Masking (M1)



Channel-Wise Masking (M2)

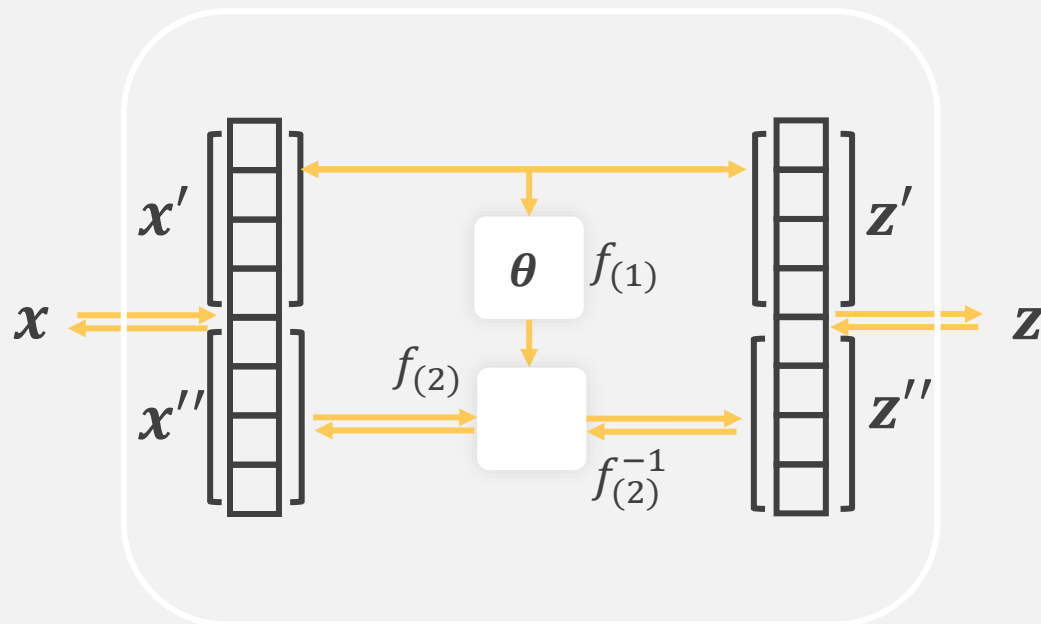


RealNVP (2017)

Affine Coupling Transform • Masking Schemes • **Architecture** • Results

“Affine Coupling Layer”

f_{ACL}



$$f_{(1)}(\cdot, \theta) = \text{ResNet}(\cdot, n_{\text{blocks}}, \theta)$$

$$f_{(2)}(x'', (t, s)) = x'' \odot \exp(s) + t$$

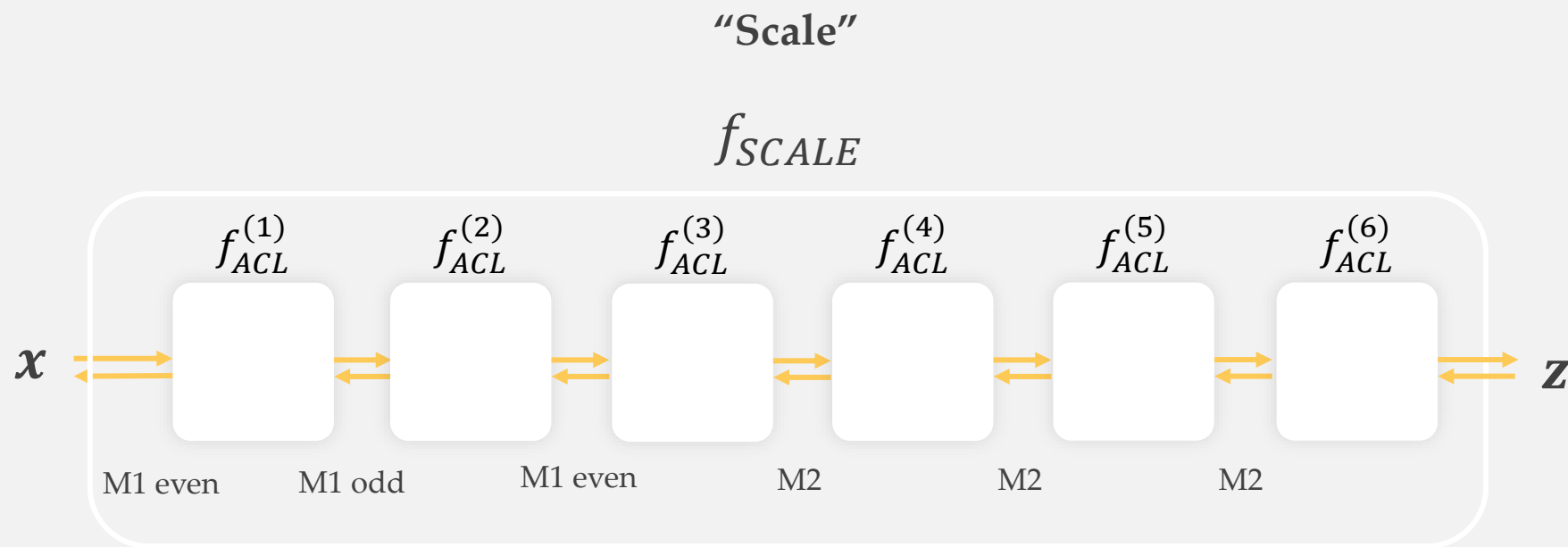
$$f_{(2)}^{-1}(z'', (t, s)) = (z'' - t) \odot \exp(-s)$$

RealNVP uses:

- $n_{\text{blocks}} = 4$ for 32x32 images.
- $n_{\text{blocks}} = 2$ for 64x64 images.

RealNVP (2017)

Affine Coupling Transform • Masking Schemes • **Architecture** • Results

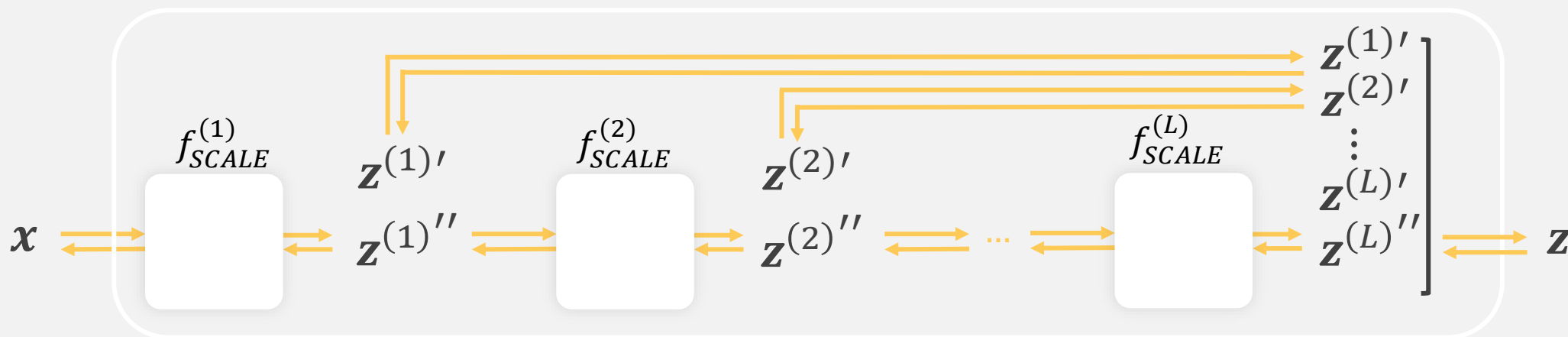


RealNVP (2017)

Affine Coupling Transform • Masking Schemes • **Architecture** • Results

Final Multi-Scale Architecture

$$f(x | \theta) = z$$

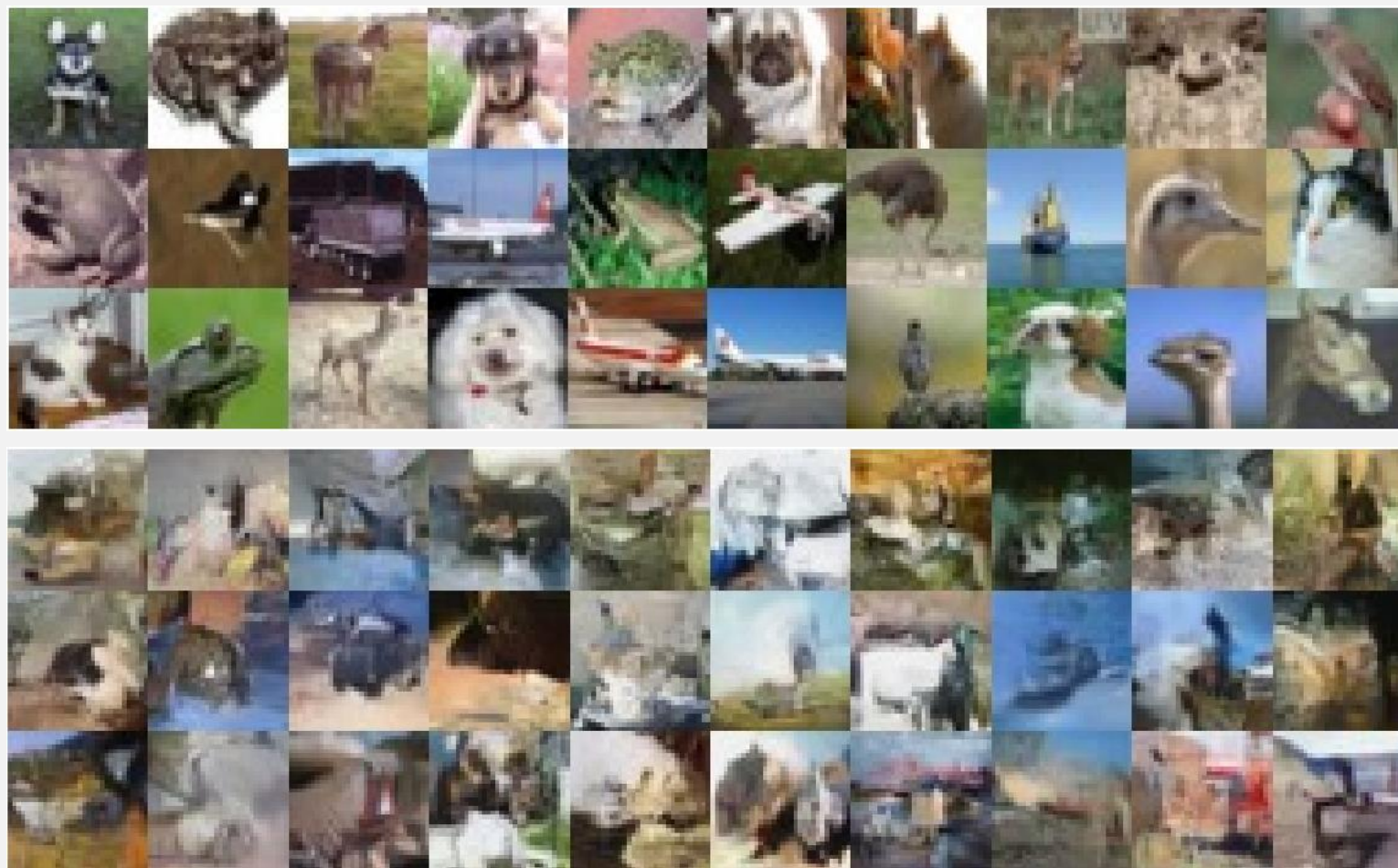


$$f^{-1}(z | \theta) = x$$

RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results

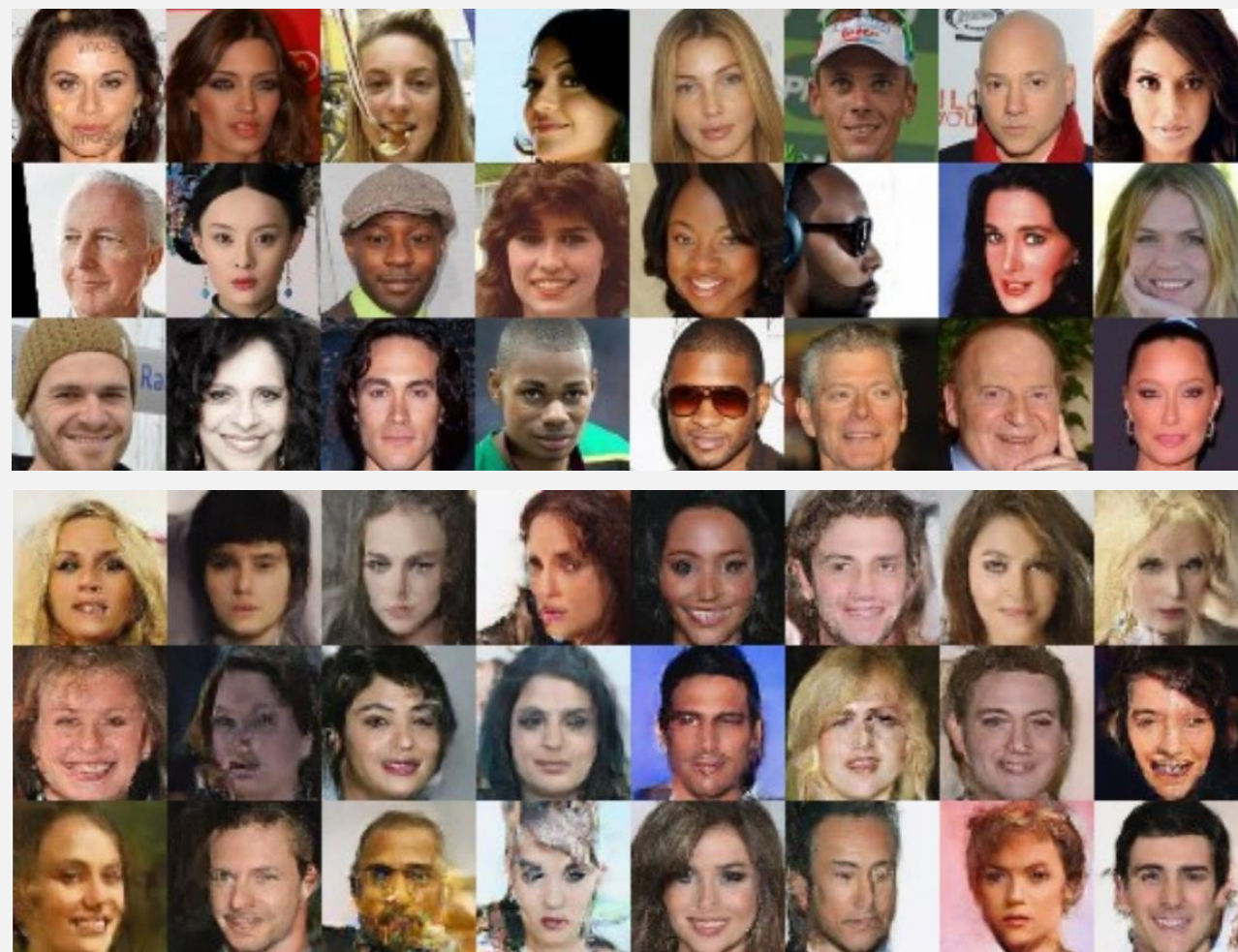
Sampling: CIFAR-10



RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results

Sampling: CelebA



RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results

Interpolating: CelebA



RealNVP (2017)

Affine Coupling Transform • Masking Schemes • Architecture • Results

Interpolating: LSUN (Tower)



Glow (2018)

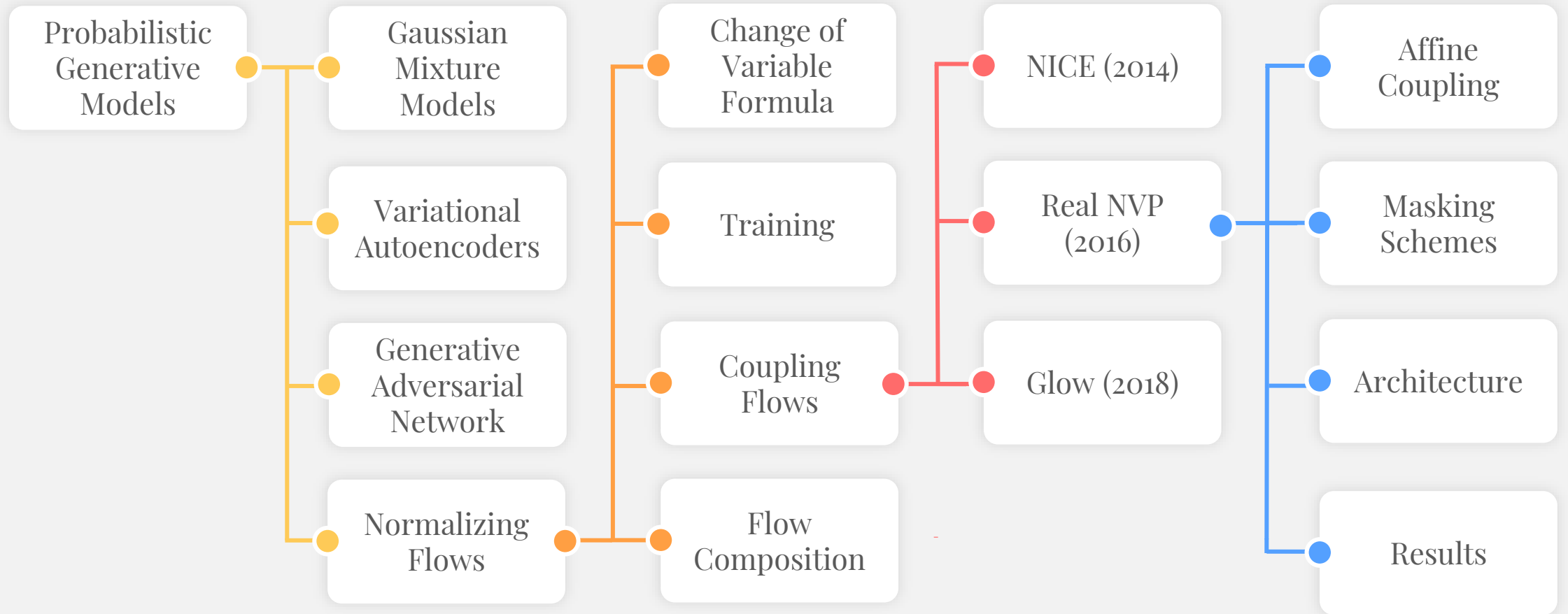
Sampling: CelebA



My Take

- + Even if the results are not so good, RealNVP laid the groundwork for papers like GLOW.
- + Invertible NN are a very interesting idea.
- The latent space is hard to interpret and “customize”.
- Some info about RealNVP is only found in the code.

Map



Thank you for your attention.

References

- RealNVP Paper: <https://arxiv.org/abs/1605.08803>
- NICE Paper: <https://arxiv.org/abs/1410.8516>
- Glow Paper: <https://arxiv.org/abs/1807.03039>
- Code Implementations RealNVP: <https://paperswithcode.com/paper/density-estimation-using-real-nvp>