



Tractable density models- Normalizing Flows

Generative and Deep Learning (GDL)
Davide Bacciu (davide.bacciu@unipi.it)



UNIVERSITÀ DI PISA

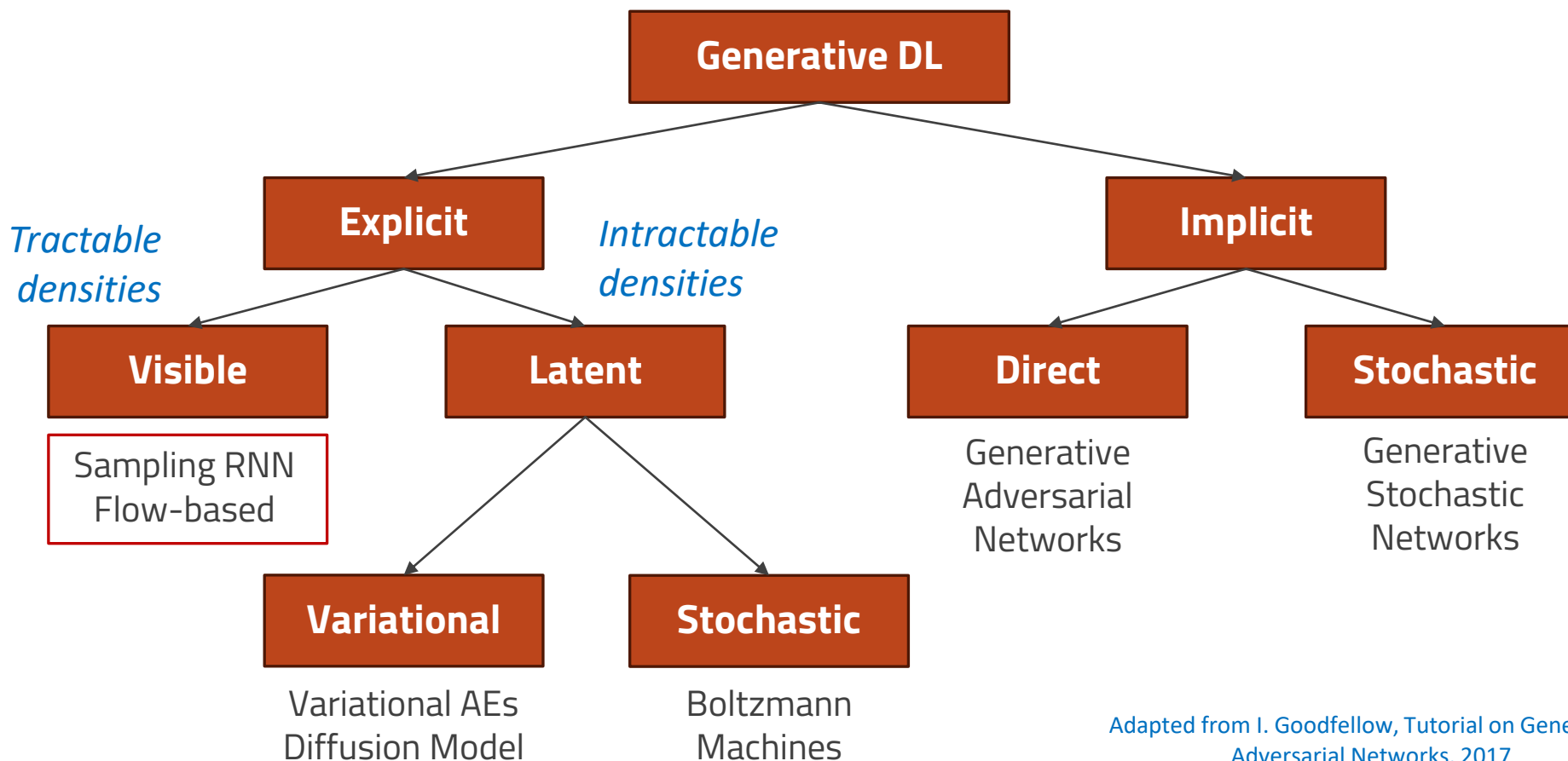


Lecture(s) outline

Generative models with explicit likelihood

- ◇ Autoregressive learning with fully visible information
- ◇ Explicitly likelihood by change of variables
- ◇ Normalizing flows
 - ◇ Coupling flows
 - ◇ Masking & squeezing
 - ◇ Invertible convolutions
 - ◇ Autoregressive flows
- ◇ Residual and continuous normalizing flows

A Taxonomy



Adapted from I. Goodfellow, Tutorial on Generative Adversarial Networks, 2017

Explicit likelihood models

Learning with Fully Visible Information

If all information is fully visible the joint distribution can be computed from the **chain rule factorization**

Bayesian Networks $\rightarrow P(\mathbf{x}) = \prod_i^N P(x_i | x_1, \dots, x_{i-1})$



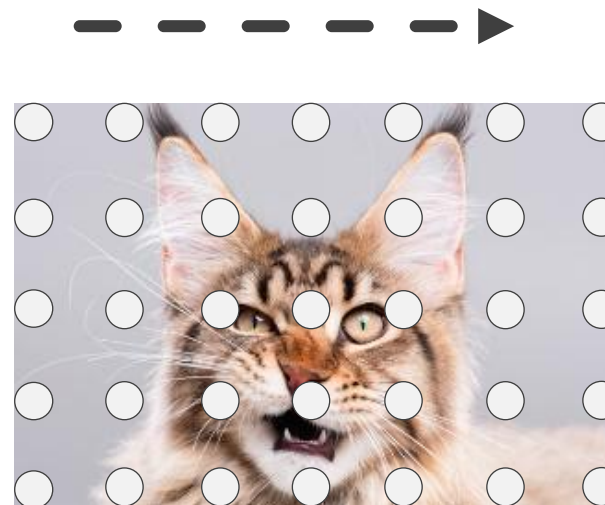
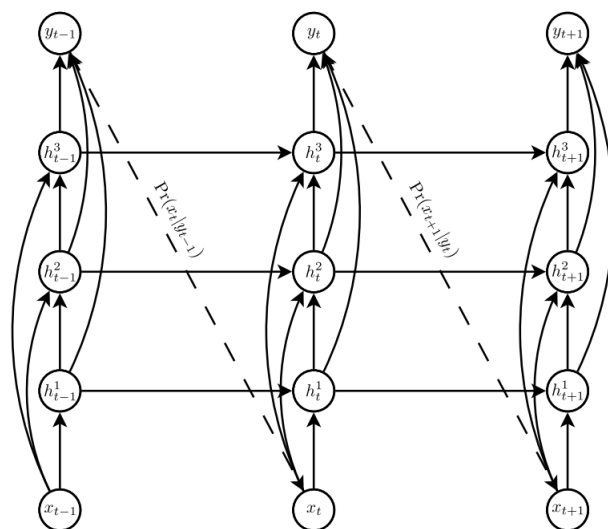
Probability of a pixel having a certain intensity value, given the known intensity of its predecessor

Need to be able to define a sensible ordering for the chain rule

Conditional distribution difficult to compute

Approximating the Conditional Probability

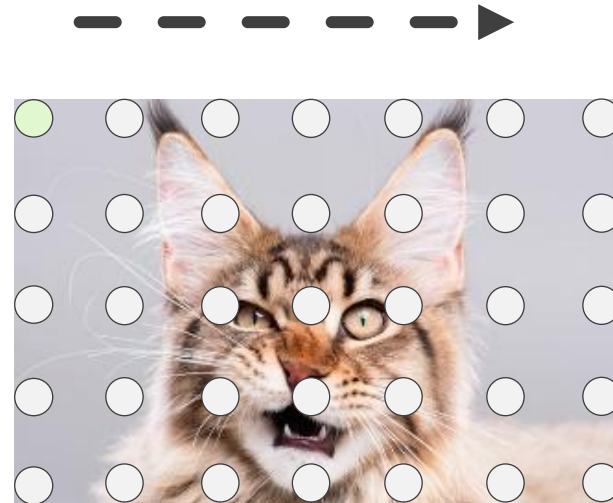
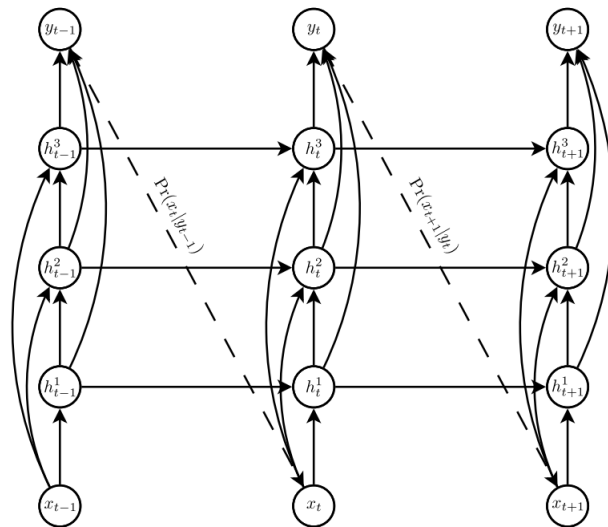
If all information is fully visible the joint distribution can be computed from the **chain rule factorization**



Scan the image according to a schedule and **encode the dependency** from previous pixels in the **states of an RNN**

Approximating the Conditional Probability

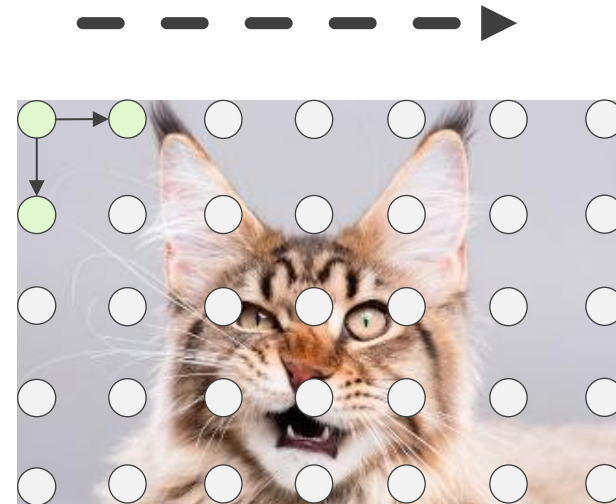
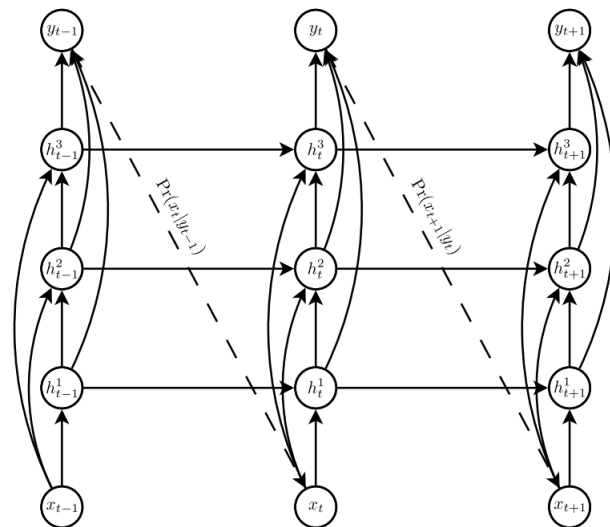
If all information is fully visible the joint distribution can be computed from the **chain rule factorization**



Scan the image according to a schedule and **encode the dependency** from previous pixels in the **states of an RNN**

Approximating the Conditional Probability

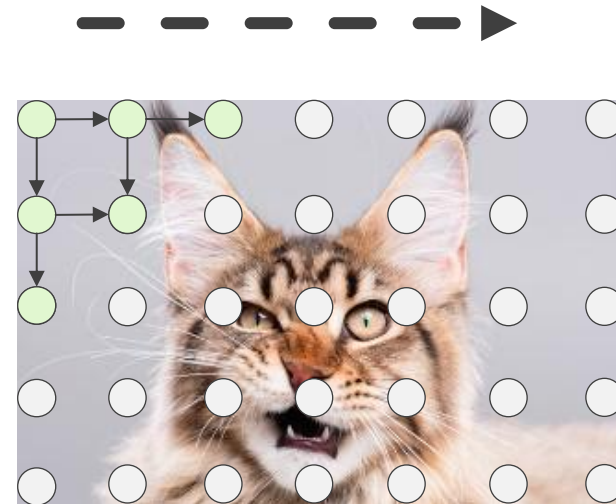
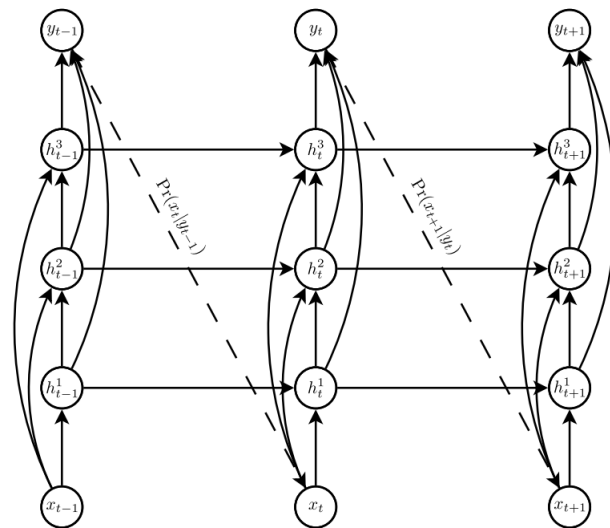
If all information is fully visible the joint distribution can be computed from the **chain rule factorization**



Scan the image according to a schedule and **encode the dependency** from previous pixels in the **states of an RNN**

Approximating the Conditional Probability

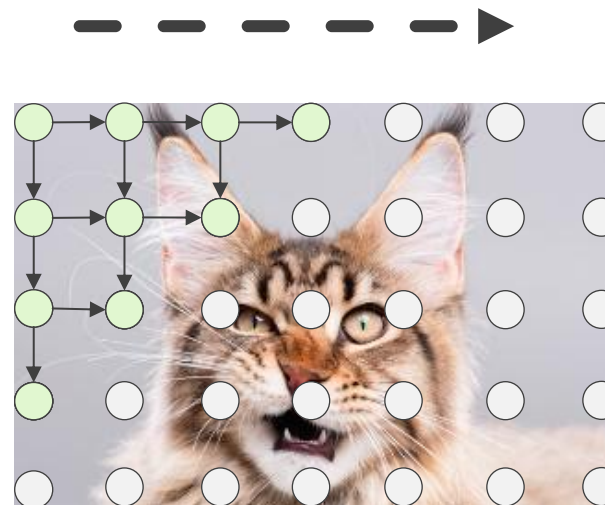
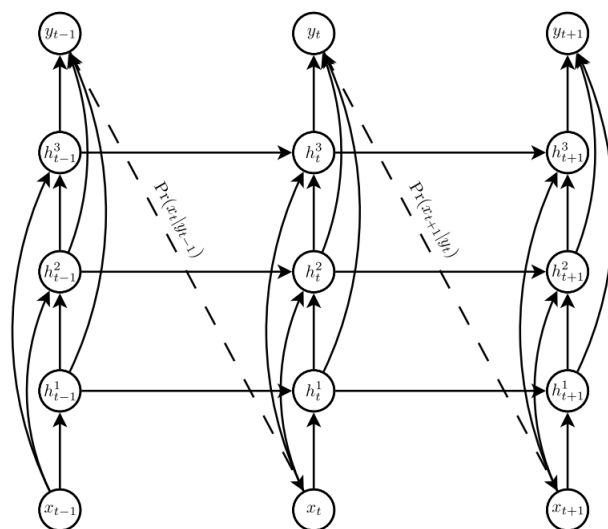
If all information is fully visible the joint distribution can be computed from the **chain rule factorization**



Scan the image according to a schedule and **encode the dependency** from previous pixels in the **states of an RNN**

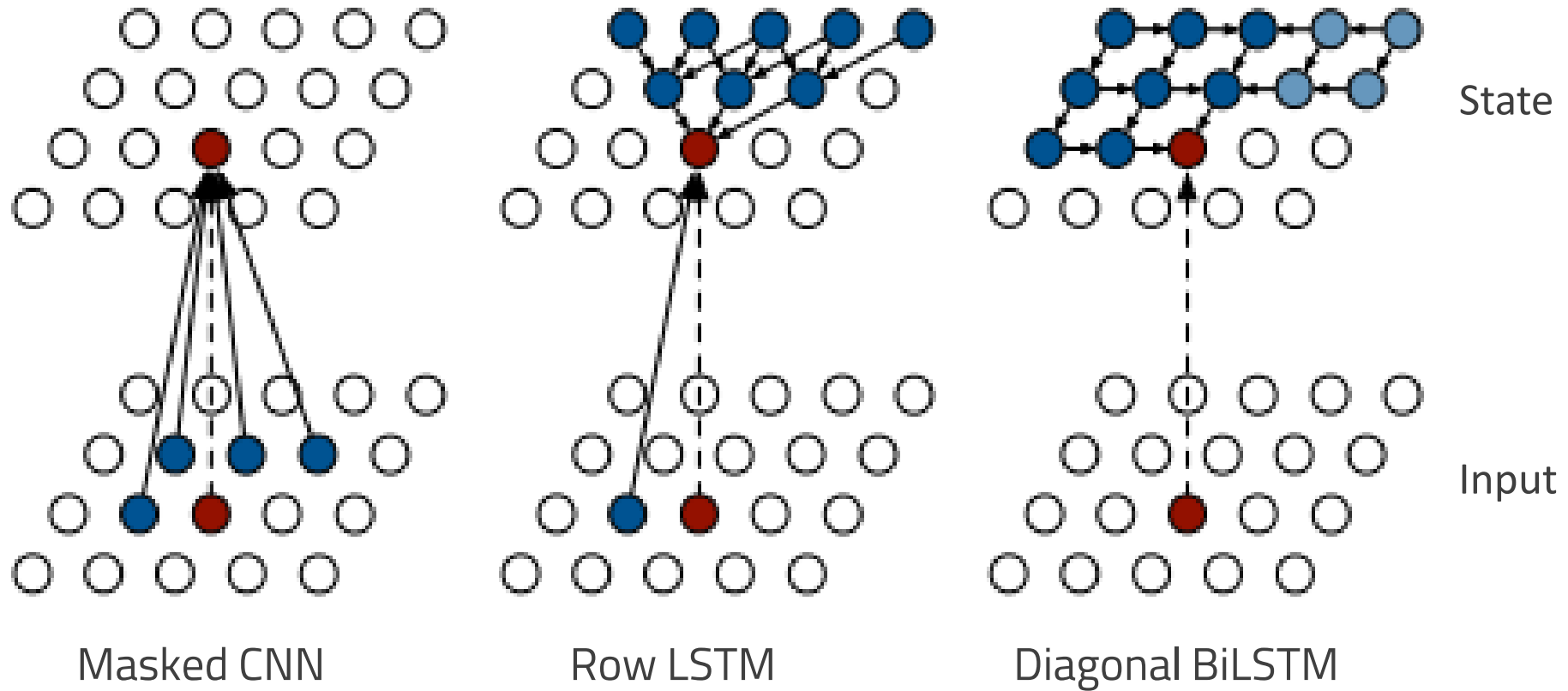
Approximating the Conditional Probability

If all information is fully visible the joint distribution can be computed from the **chain rule factorization**



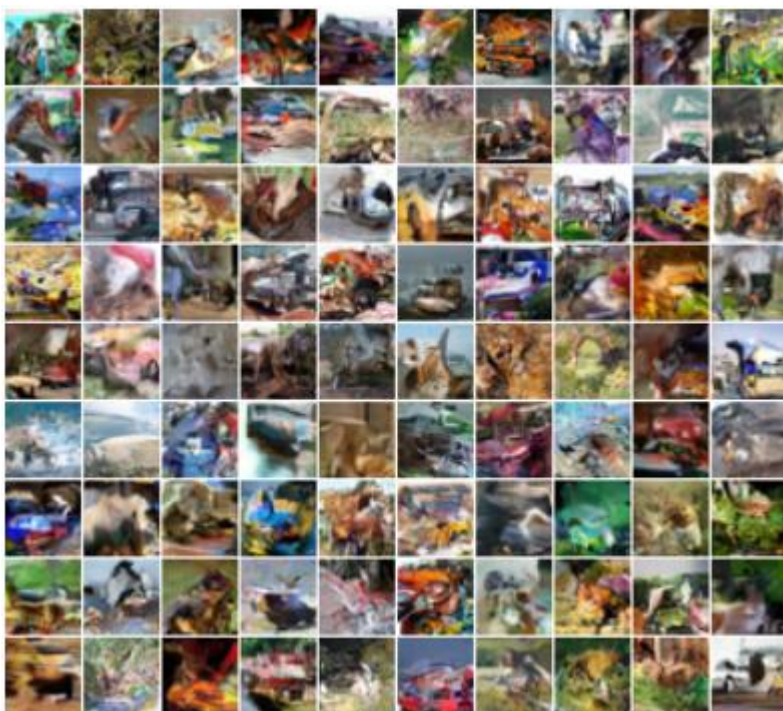
Scan the image according to a schedule and **encode the dependency** from previous pixels in the **states of an RNN**

Generating Images Pixel by Pixel



A. van der Oord et al., Pixel Recurrent Neural Networks, 2016

Generating Images Pixel by Pixel - Results



32x32 CIFAR-10



32x32 ImageNet

A. van der Oord et al., Pixel Recurrent Neural Networks, 2016

Summary on autoregressive models

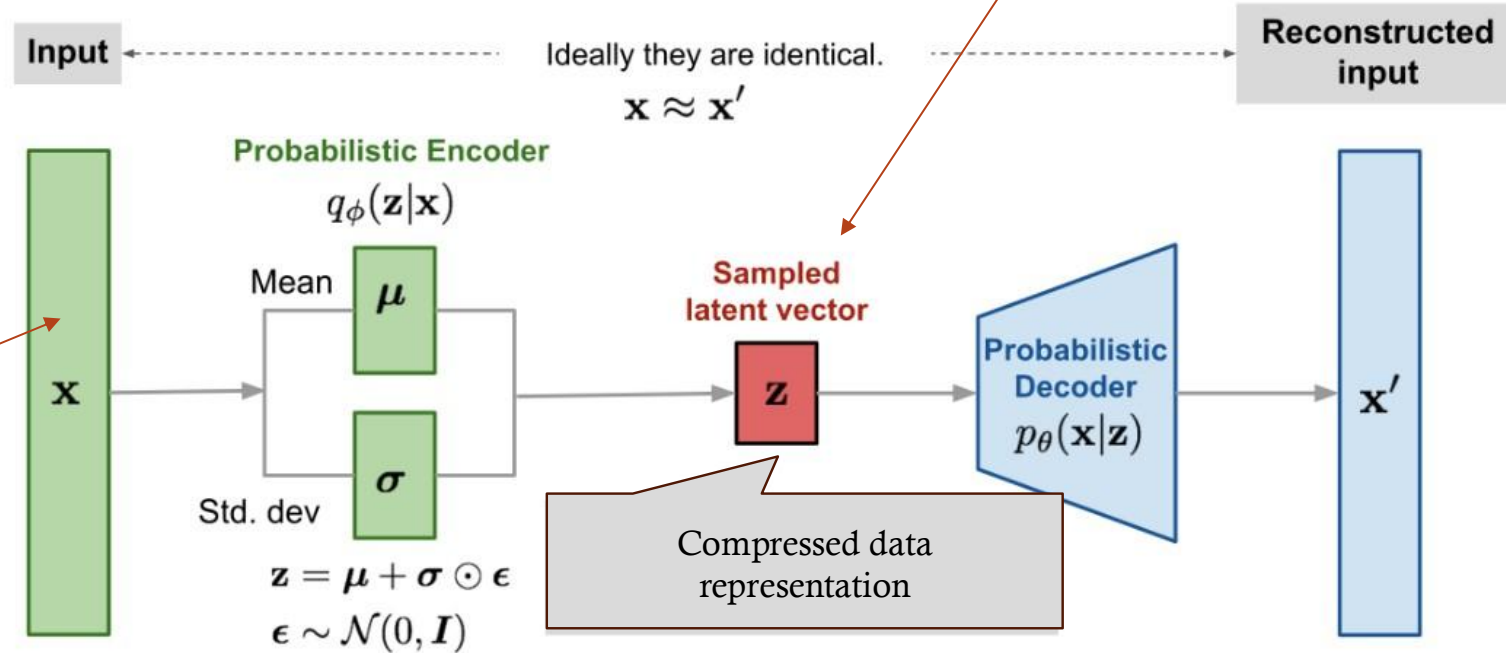
- ◇ Autoregressive models are easy to understand and design
 - ◇ Their left-to-right factorization is a **natural fit for modeling text**
 - ◇ They are at the core of modern LLMs (using transformers in place of RNNs/CNNs)
- ◇ **Computing the likelihood** is easy in general and **maximizing it is easy** for text (made of discrete tokens)
- ◇ **Variable-length generations** can be sampled one token at a time
- ◇ They are **very slow to sample** (sequential)
- ◇ No natural way to get meaningful **latent representations**
- ◇ Sub-optimal sample quality beyond text

Intuition and Motivation

Remember me? The VAE

Requires to capture all relevant variations in a low-dimensional embedding \mathbf{z}

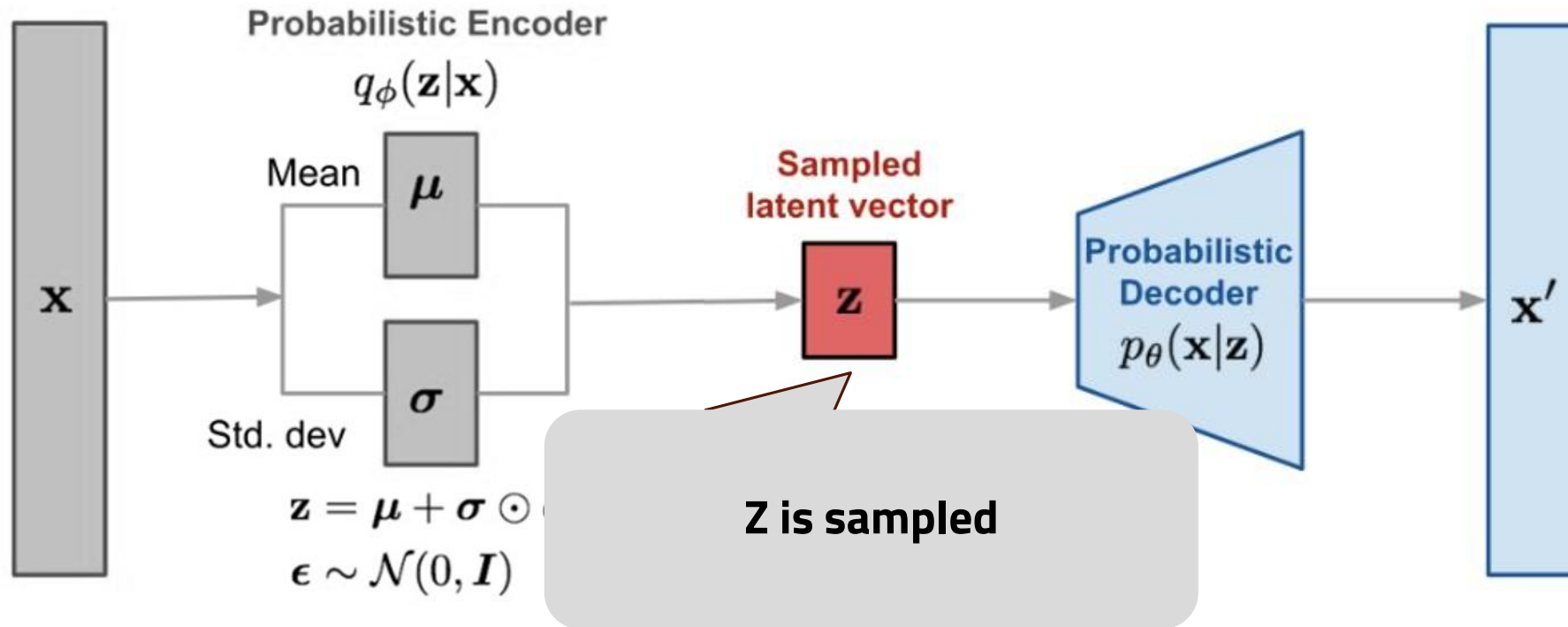
Needs separate encoder/decoder networks



$$\mathbb{E}_{\mathbf{x} \sim D} [\mathbb{E}_{\mathbf{z} \sim Q} [\log P(\mathbf{x}|\mathbf{z})] - KL(Q(\mathbf{z}|\mathbf{x}, \phi) || P(\mathbf{z}))]$$

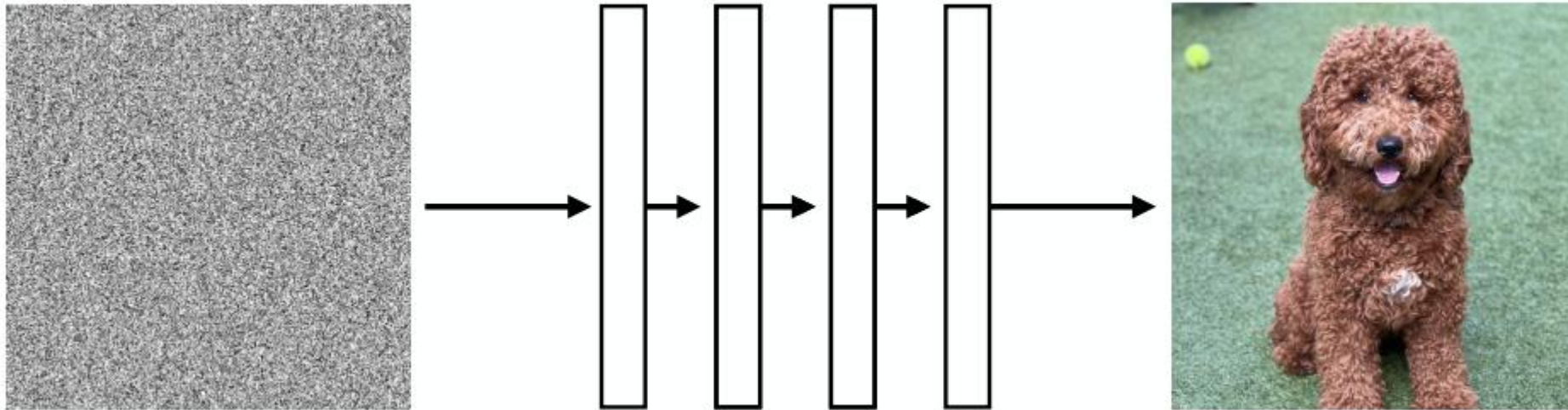
Cannot directly estimate $P_\theta(\mathbf{x})$

Can we get rid of the encoder?



Invertible Transformations

The **generative** direction



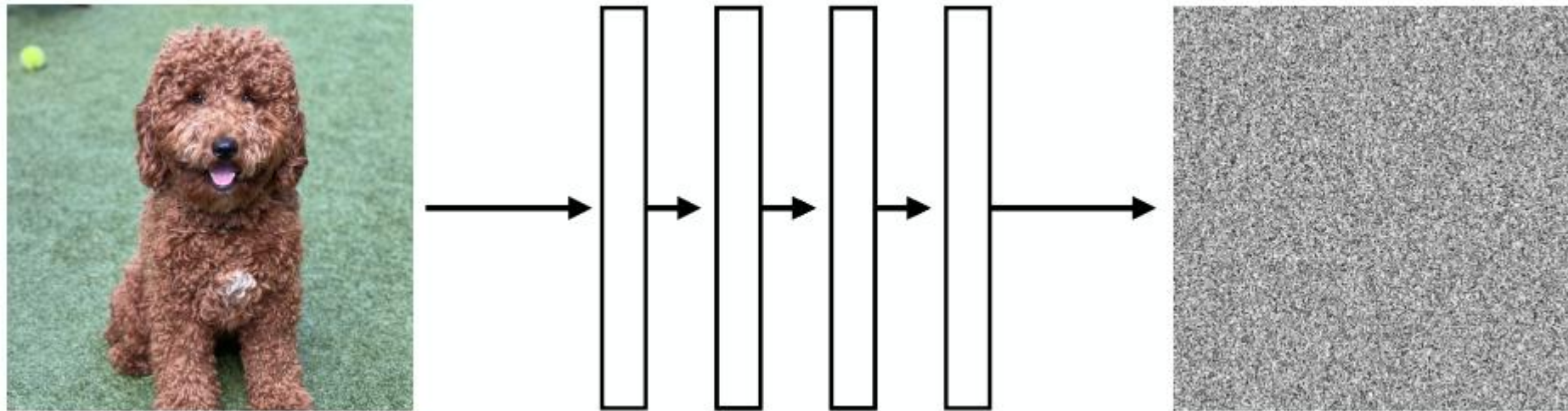
\mathbf{z} sampled from a simple distribution

$$\mathbf{x} = g(\mathbf{z})$$

Reconstructed as a sample from the data distribution \mathbf{x}

Invertible Transformations

The **normalizing** direction



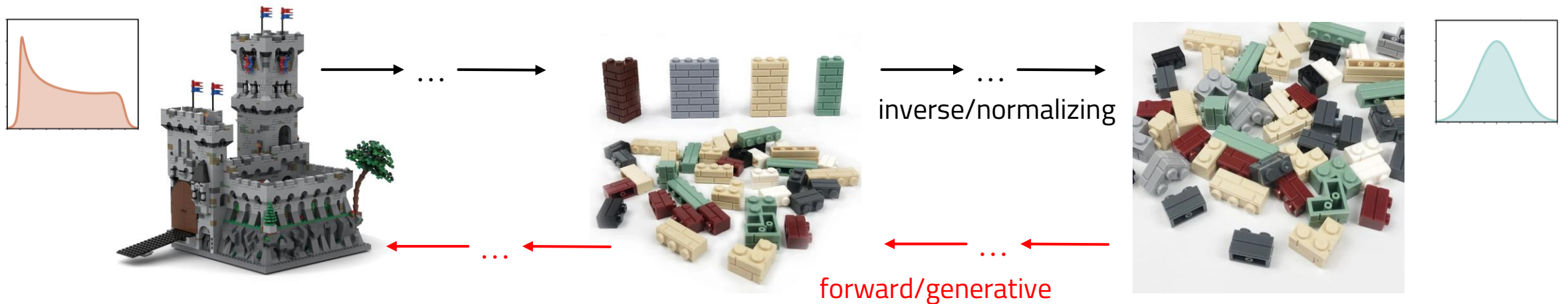
\mathbf{x} data sample

$$\begin{aligned}\mathbf{z} &= f(\mathbf{x}) \\ &= g^{-1}(\mathbf{x})\end{aligned}$$

Transformed in
noise \mathbf{z}

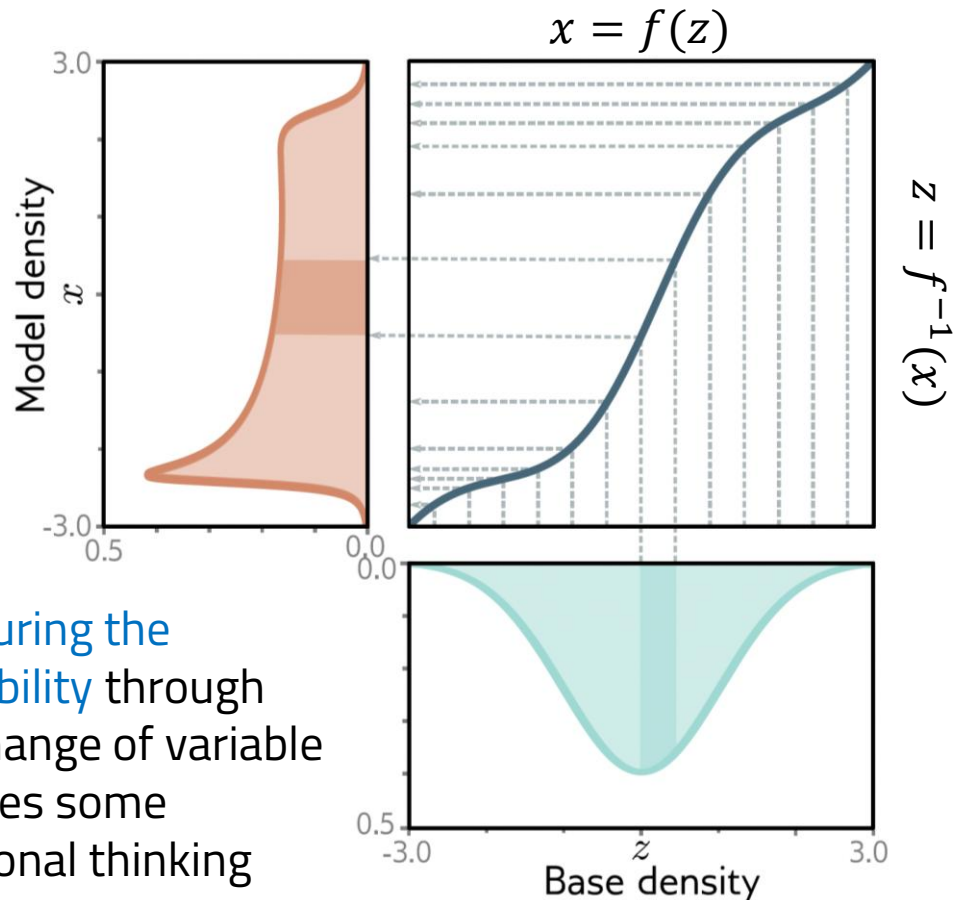
Normalizing Flow (NF) – The Intuition

- ◇ Learn a probabilistic model by **transforming a simple distribution into the complex data generating distribution** using a deep network
- ◇ Easy to sample and evaluate the probability
- ◇ Requires a specialized architecture where each layer must be invertible



Change of variable

Probabilistic Change of Variable



Measuring the probability through the change of variable requires some additional thinking

- ◆ Take a tractable base distribution $P(z)$ over latent variable z and a model density $P(x)$ over data x
- ◆ Apply a change of variable function (possibly learned with parameters θ)

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

- ◆ In addition, we are going to require that f is invertible

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

Linear 1D Change of Variable

NF define complex densities by transforming a base one by invertible mappings (bijections)

- ◇ Simplest case in 1D is a **univariate Gaussian base density**

$$z \sim \mathcal{N}(0,1)$$

- ◇ Simplest change of variable (forward) by **linear transformation**

$$x = f(z; \mu, \sigma) = \mu + z\sigma$$

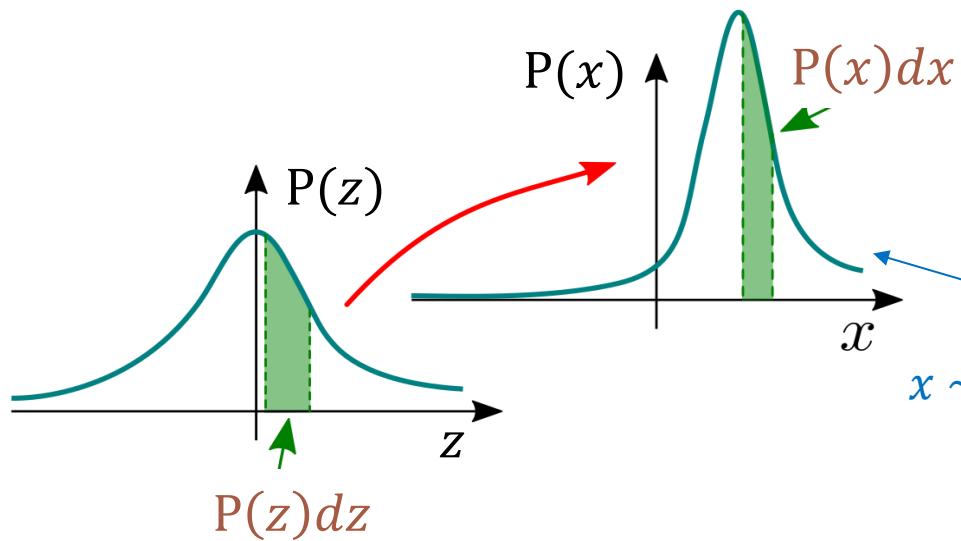
- ◇ **Inverse** then (under $\sigma \neq 0$)

$$z = f^{-1}(x; \mu, \sigma) = (x - \mu)/\sigma$$

- ◇ With $P(z)$ known we want to **find $P(x)$**

Linear 1D – Mass conservation

The volume may change but the density must be preserved



The necessary condition for this is

$$P(z)dz = P(x)dx$$

The probability of data x under the transformed distribution is

$$x \sim \mathcal{N}(\mu, \sigma^2)$$

$$P(x) = P(z) \left| \frac{dx}{dz} \right|^{-1} = \frac{P(z)}{\sigma} \frac{x - \mu}{\sigma}$$

Change of volume \rightarrow $\left| \frac{dx}{dz} \right|^{-1}$

$x = \mu + \sigma z$

Linear 1D – Iterated forward pass

$$P(x) = P(z) \left| \frac{dx}{dz} \right|^{-1} \quad \text{Forward transformation equation}$$

- ◇ Sample x through 2 mappings (transformations)

$$z_0 \sim P(z) \quad z_1 = f_1(z_0) \quad x = f_2(z_1)$$

- ◇ Density obtained by composing forward transformations

$$P(x) = P(z_0) \left| \frac{dz^1}{dz^0} \right|^{-1} \left| \frac{dx}{dz^1} \right|^{-1}$$

Linear 1D – Inverse Flow

- ◇ We may be interested in **estimating the density** of a given input sample x
- ◇ Requires building the **inverse flow** ($g = f^{-1}$)

$$z_1 = f_2^{-1}(x) = g_2(x) \quad z_0 = f_1^{-1}(z_1) = g_1(z_1)$$

- ◇ And computing the density accordingly

$$P(x) = \left| \frac{dz^1}{dx} \right| \left| \frac{dz^0}{dz^1} \right| P(z_0)$$

Multidimensional flow

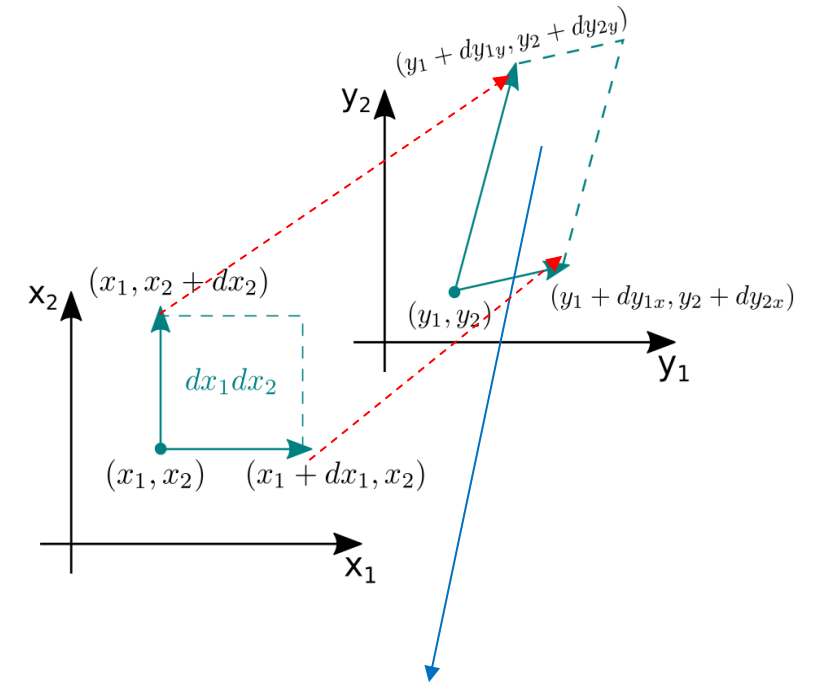
- ◇ Extend the approach to multi-dimensional case
 - ◇ \mathbf{x}, \mathbf{z} vectorial RVs with density $P(\mathbf{z})$ and $P(\mathbf{x})$
 - ◇ Flow $f(\mathbf{z})$ invertible and differentiable (closed under composition)
- ◇ Transformation $\mathbf{x} = f(\mathbf{z})$ leads to the probability change

$$P(\mathbf{x}) = P(\mathbf{z}) \left| \det \frac{\partial f(\mathbf{z})}{\partial \mathbf{z}} \right|^{-1}$$

Determinant

Jacobian

Provides information on the **rate of change of the volume** affected by the f transformation



Area of this field can be computed with vector calculus and turns out to be the Jacobian determinant

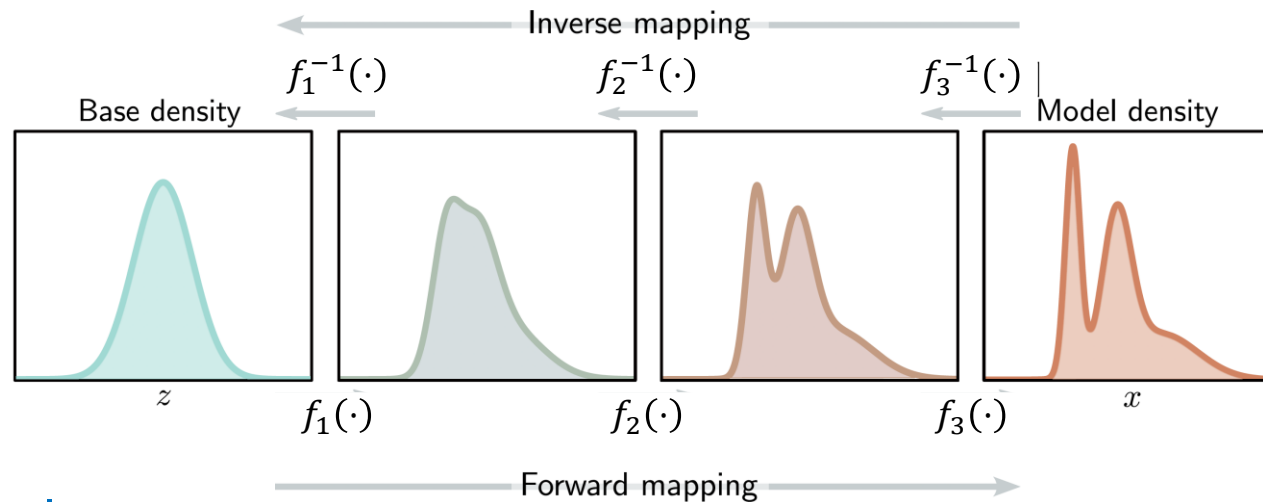
General Multistep Case

Density of the input

Used to “know” the likelihood
(e.g. learning, anomaly
detection)

$$\mathbf{z}_0 = f_1^{-1}(\mathbf{z}_1)$$

$$P(\mathbf{x}) = P(\mathbf{z}_0) \prod_{i=1}^N \left| \det \frac{\partial f_i^{-1}(\mathbf{z}_i)}{\partial \mathbf{z}_i} \right| = P(\mathbf{z}_0) \prod_{i=1}^N |\det J_{\mathbf{z}_i}(f_i^{-1})| \quad \begin{array}{l} \mathbf{z}_N = \mathbf{x} \\ \mathbf{z}_0 = \mathbf{z} \end{array}$$



Density of the sample

Used for sampling
 $\mathbf{z}_0 \sim P(\mathbf{z}_0)$

$$P(\mathbf{x}) = P(\mathbf{z}_0) \prod_{i=1}^N \left| \det \frac{\partial f_i(\mathbf{z}_{i-1})}{\partial \mathbf{z}_{i-1}} \right|^{-1} = P(\mathbf{z}_0) \prod_{i=1}^N |\det J_{\mathbf{z}_{i-1}}(f_i)|^{-1} \quad \begin{array}{l} \mathbf{z}_N = \mathbf{x} \\ \mathbf{z}_0 = \mathbf{z} \end{array}$$

Some considerations & desiderata

- ◇ Can use **log densities** for stability and learning

$$\log P(\mathbf{x}) = \log P(\mathbf{z}_0) + \sum_{i=1}^N \log |\det J_{\mathbf{z}_{i-1}}(f_i^{-1})| = \log P(\mathbf{z}_0) - \sum_{i=1}^N \log |\det J_{\mathbf{z}_i}(f_i)|$$

- ◇ Can optimize the parameters θ of the $f_i(\cdot; \theta)$ **by gradient based optimization of the log-likelihood** above
 - ◇ f_i needs to be **invertible and differentiable** (and remain so throughout learning)
 - ◇ f_i composition needs to be **expressive** enough to map Normal into arbitrary distributions
 - ◇ Need to **compute determinant easily** (e.g. Jacobian diagonal or triangular matrix)
 - ◇ Computation of f_i needs to be efficient for sampling
 - ◇ Computation of f_i^{-1} needs to be efficient for learning
 - ◇ Computation of f_i needs to be stable numerically

Normalizing Flow Layers

Invertible neural layers - Linear

- ◇ **Affine** (i.e. linear) flows can be invertible

$$\mathbf{f}(\mathbf{z}) = \mathbf{b} + \mathbf{W}\mathbf{z}$$

- ◇ **W** needs to be full rank \Rightarrow expensive to compute determinant in the general case
- ◇ Interesting special cases
 - ◇ **W diagonal** (poorly expressive)
 - ◇ **W triangular** (quadratic cost to invert; still poorly expressive)
 - ◇ **W permutation** (volume preserving ($\det = 1$); often used jointly with other flow layers)

$$\begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 1 & 0 \end{bmatrix}$$

Invertible neural layers – Pointwise nonlinear

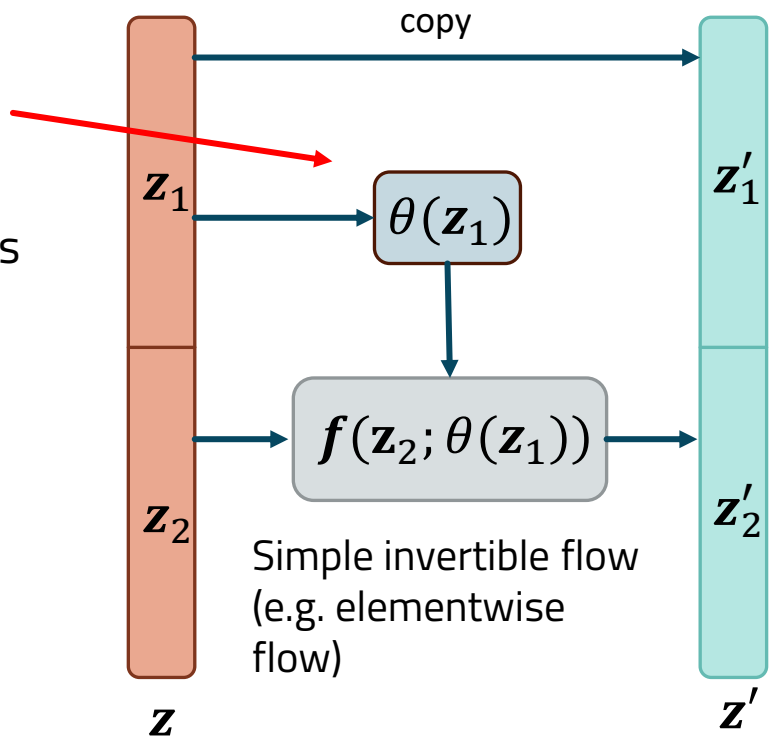
- ◆ **Pointwise nonlinear** where f are piecewise linear or smooth splines (nonlinear and easy to compute)

$$\mathbf{f}(\mathbf{z}) = [f(z^{(1)}, \theta), f(z^{(2)}, \theta), \dots, f(z^{(D)}, \theta)]$$

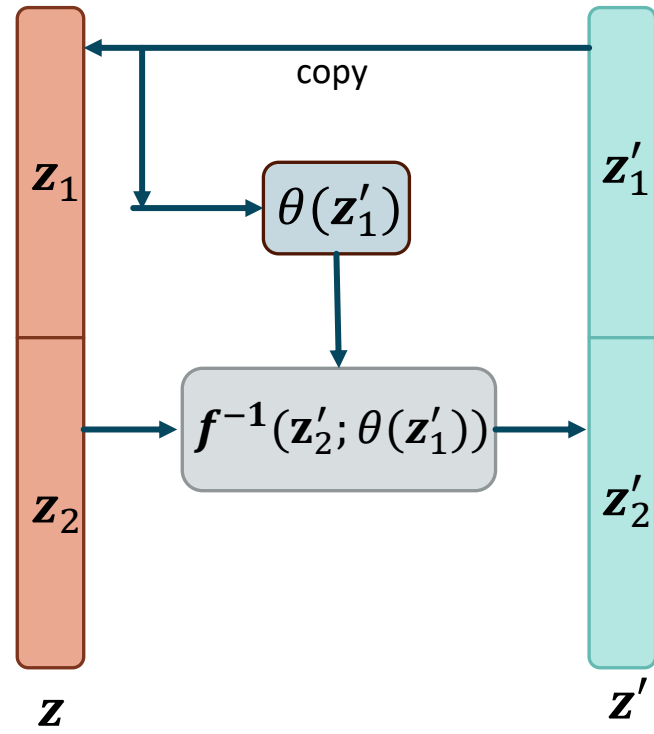
- ◆ If f is invertible, so it is \mathbf{f}
- ◆ Diagonal Jacobian
- ◆ Pointwise does not allow **capturing correlations** between dimensions

Coupling Flows

Neural network that generates parameters θ for the invertible function

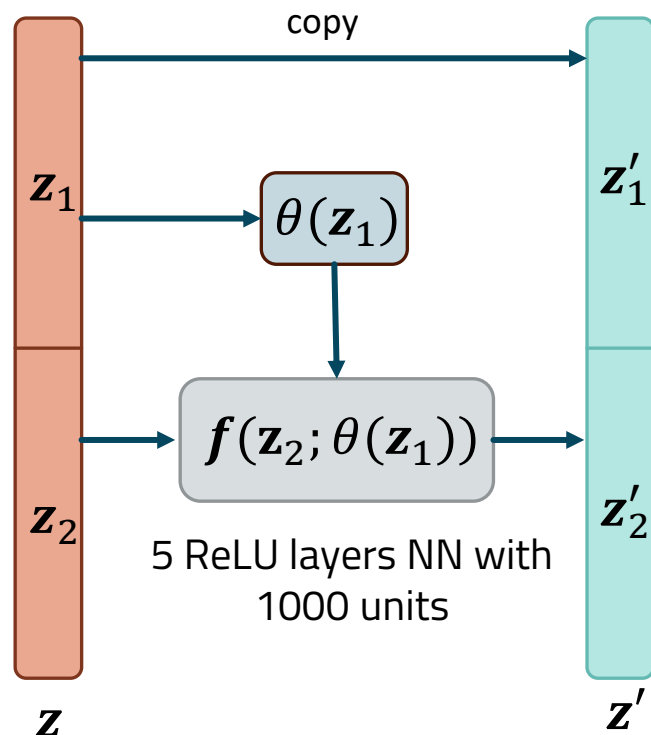


Forward



Inverse

Non-linear Independent Components Estimation (NICE)



$$\mathbf{z}'_2 = \mathbf{z}_2 + \theta(\mathbf{z}_1) \xrightarrow{\text{Inverse}} \mathbf{z}_2 = \mathbf{z}'_2 - \theta(\mathbf{z}'_1)$$

↓ Jacobian

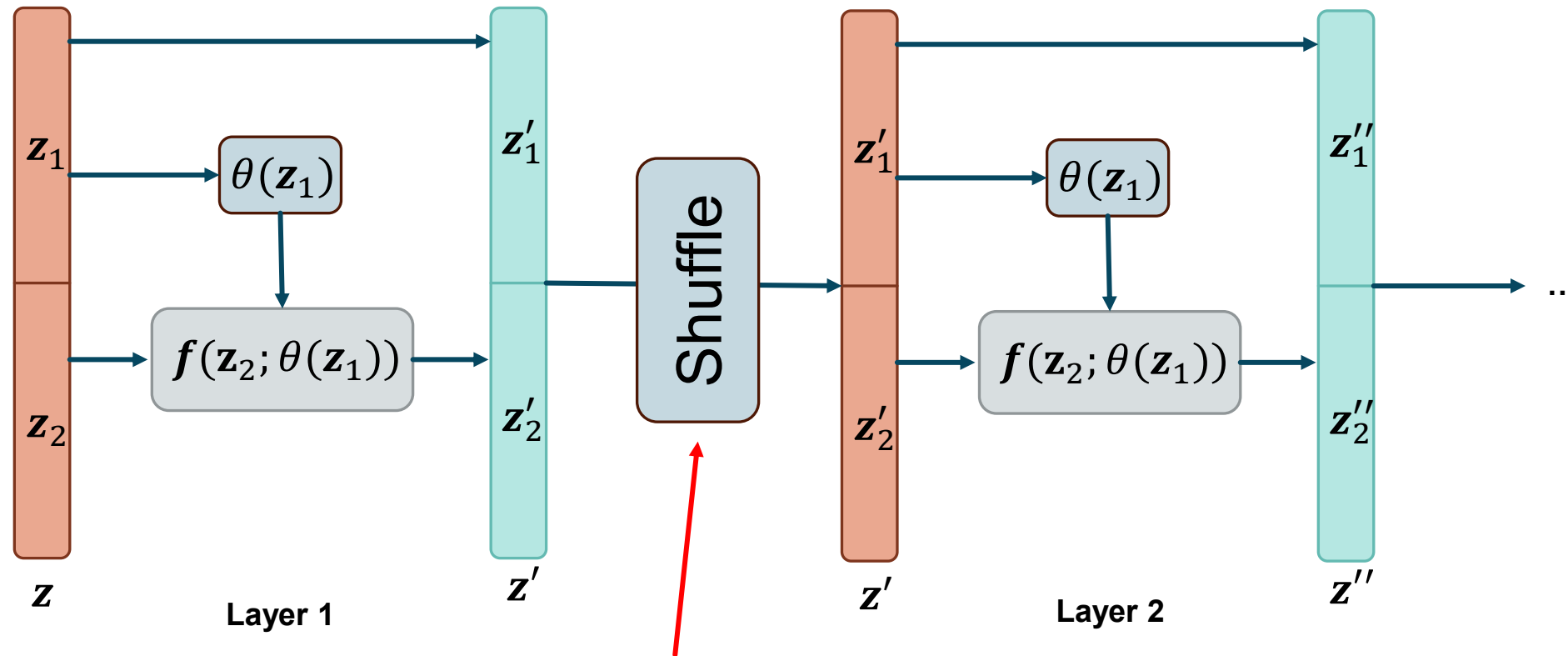
$$\begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \frac{\partial \theta(\mathbf{z}_1)}{\partial \mathbf{z}_1} & \mathbf{I} \end{pmatrix}$$

First example of coupling layer

- Simple affine transformation

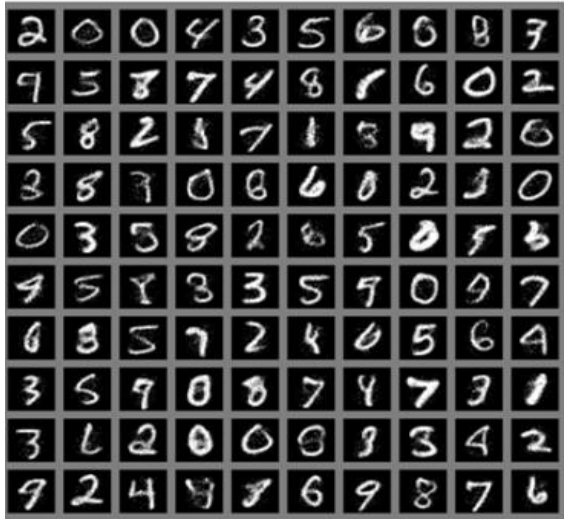
L Dinh et al, Non-linear Independent Components Estimation (NICE), ICLR-WS 2014

NICE – Stacked Coupling Flows



Random shuffle allows more general transformations than between only elements in 1st and 2nd half

(Not so) NICE Results



(a) Model trained on MNIST



(b) Model trained on TFD



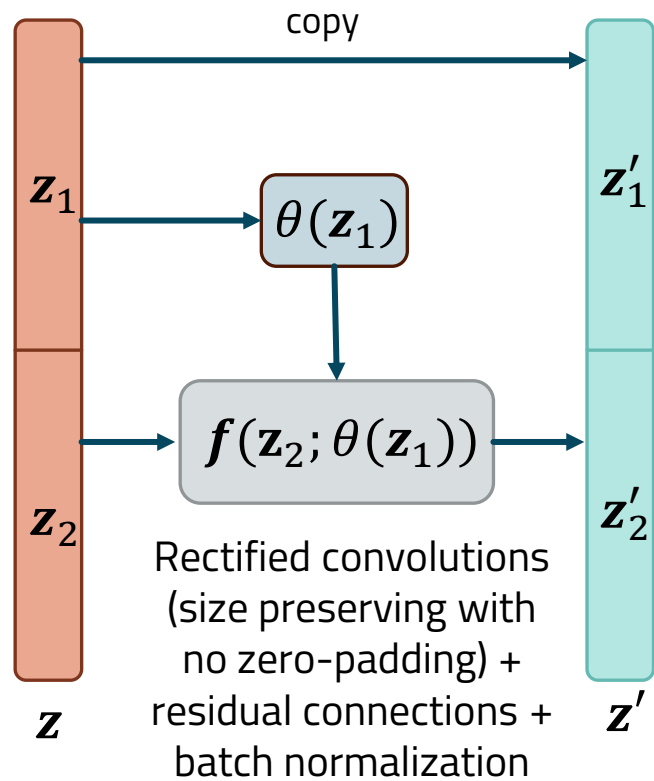
(c) Model trained on SVHN



(d) Model trained on CIFAR-10

L Dinh et al, Non-linear Independent Components Estimation (NICE), ICLR-WS 2014

RealNVP – Multiscale Nonlinear Flow



$$\mathbf{z}'_2 = \overbrace{\exp(\boldsymbol{\theta}_A(\mathbf{z}_1)) \odot \mathbf{z}_2}^{\text{Scale}} + \overbrace{\boldsymbol{\theta}_B(\mathbf{z}_1)}^{\text{Shift}} \quad (\text{Forward})$$

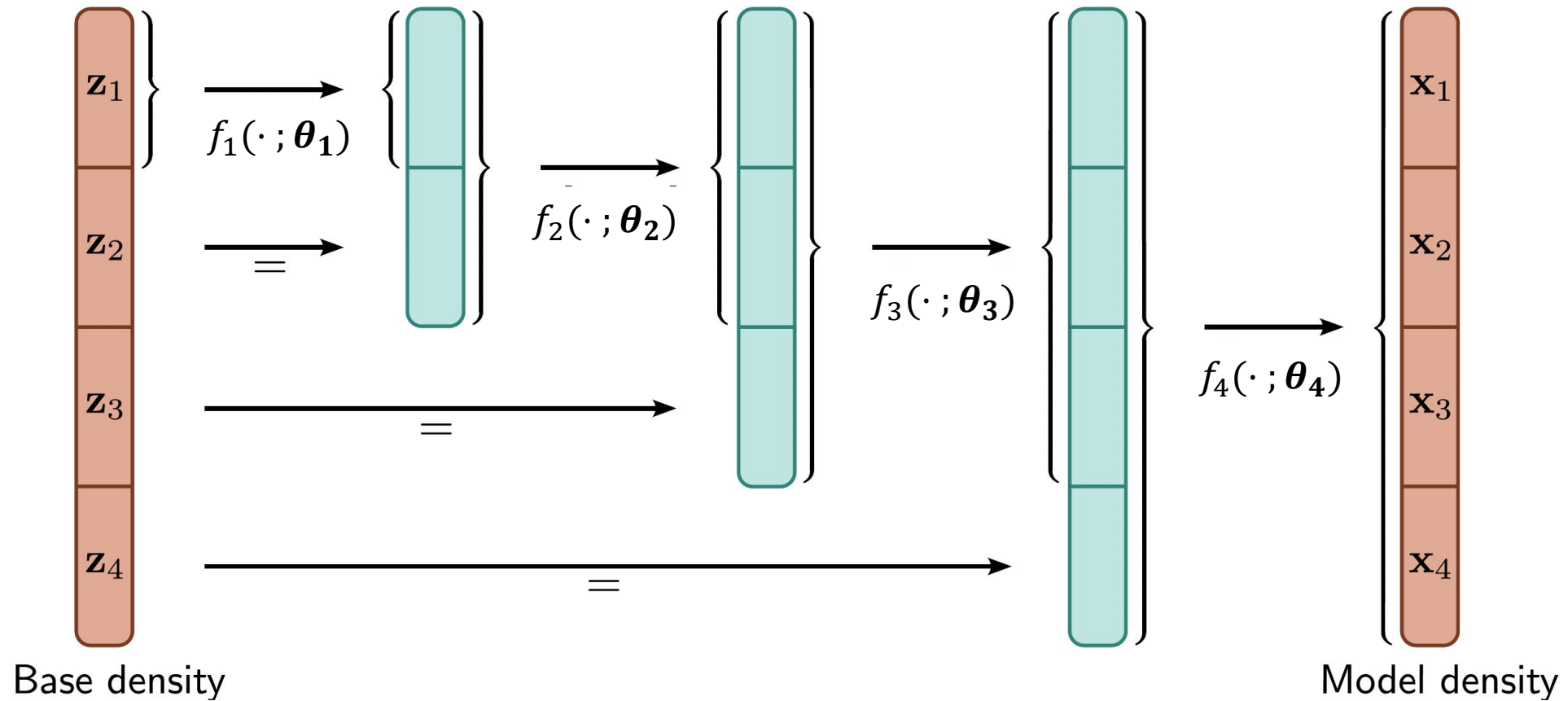
$$\mathbf{z}_2 = \frac{\mathbf{z}'_2 - \boldsymbol{\theta}_B(\mathbf{z}'_1)}{\exp(\boldsymbol{\theta}_A(\mathbf{z}'_1))} \quad (\text{Inverse})$$

Jacobian

$$\begin{pmatrix} \mathbf{I} & \mathbf{0} \\ \frac{\partial \boldsymbol{\theta}_b(\mathbf{z}_1)}{\partial \mathbf{z}_1} & \text{diag exp}(\boldsymbol{\theta}_A(\mathbf{z}_1)) \end{pmatrix}$$

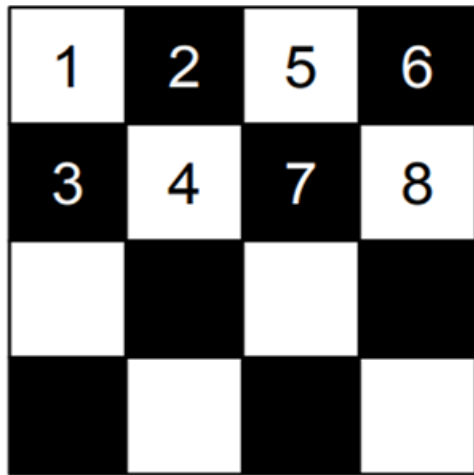
L Dinh et al, Density Estimation using real NVP, ICLR 2017

Multiscale Flows



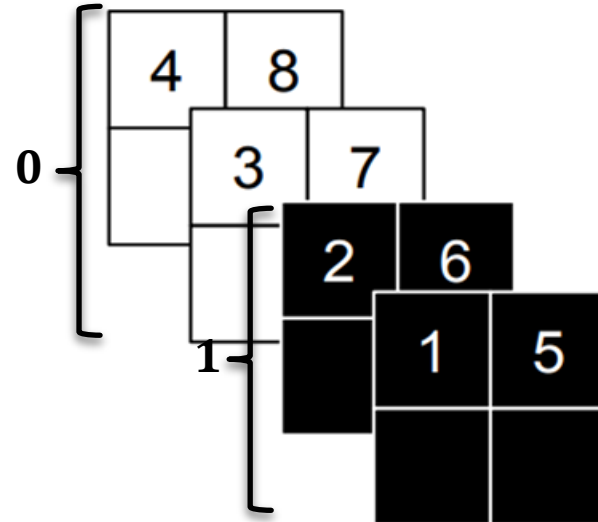
RealNVP – Masking and Squeezing

$$\mathbf{z}' = \mathbf{b} \odot \mathbf{z} + (1 - \mathbf{b}) \odot \{\exp(\boldsymbol{\theta}_A(\mathbf{b} \odot \mathbf{z})) \odot \mathbf{z} + \boldsymbol{\theta}_B(\mathbf{b} \odot \mathbf{z})\}$$



Partitioning

with checkerboard pattern



Squeezing

followed by channel-wise masking

$$S \times S \times C \xrightarrow{\text{Squeezing}} \frac{S}{2} \times \frac{S}{2} \times 4C$$

Multiscale flow implemented by alternating **binary masking** ($b_i \in \{0,1\}$) and **squeezing**

- Pixel masking before squeeze
- Channel masking after squeeze

L Dinh et al, Density Estimation using real NVP, ICLR 2017

RealNVP Results

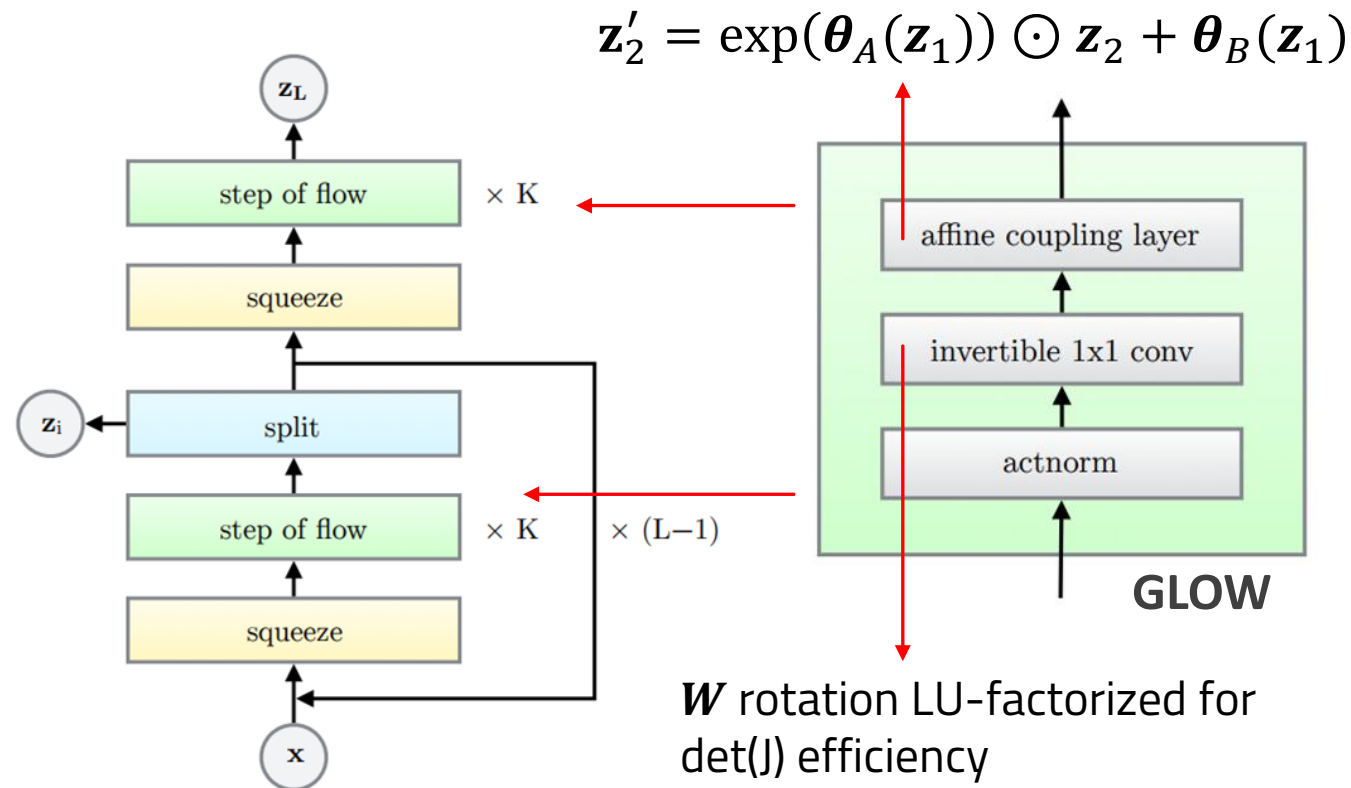
L Dinh et al, Density Estimation using real NVP, ICLR 2017



dataset

sampled

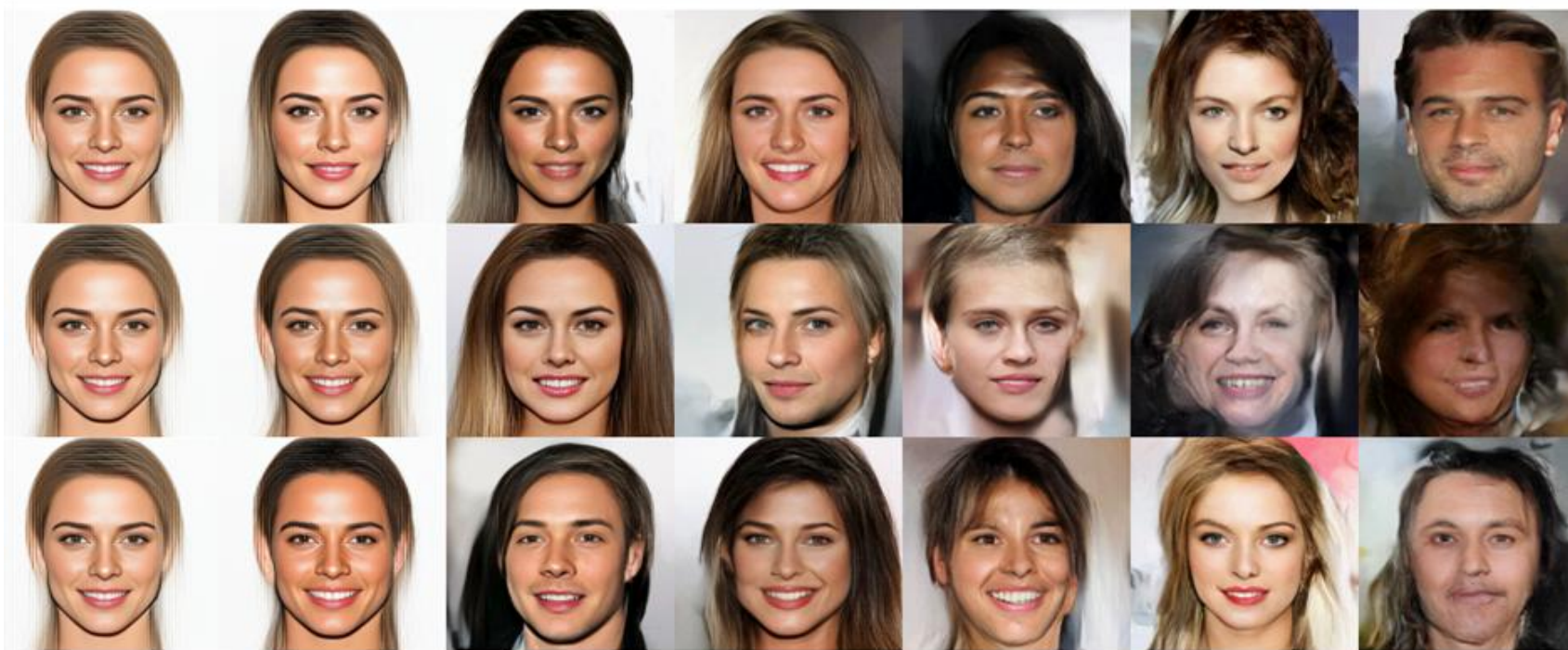
GLOW – Multiscale Coupling Flow with Invertible 1x1 Convolutions



- ◇ Start with RGB tensor
- ◇ **Split** channels in 2 halves
- ◇ Run 1x1 **convolutions** parameterized with an LU decomposition (channel mixing/permutation)
- ◇ **Affine** transform each spatial position in second half
- ◇ **Multiscale** & periodic squeeze

Kingma & Dhariwal, P, Glow: Generative flow with invertible 1x1 convolutions, NeurIPS 2018

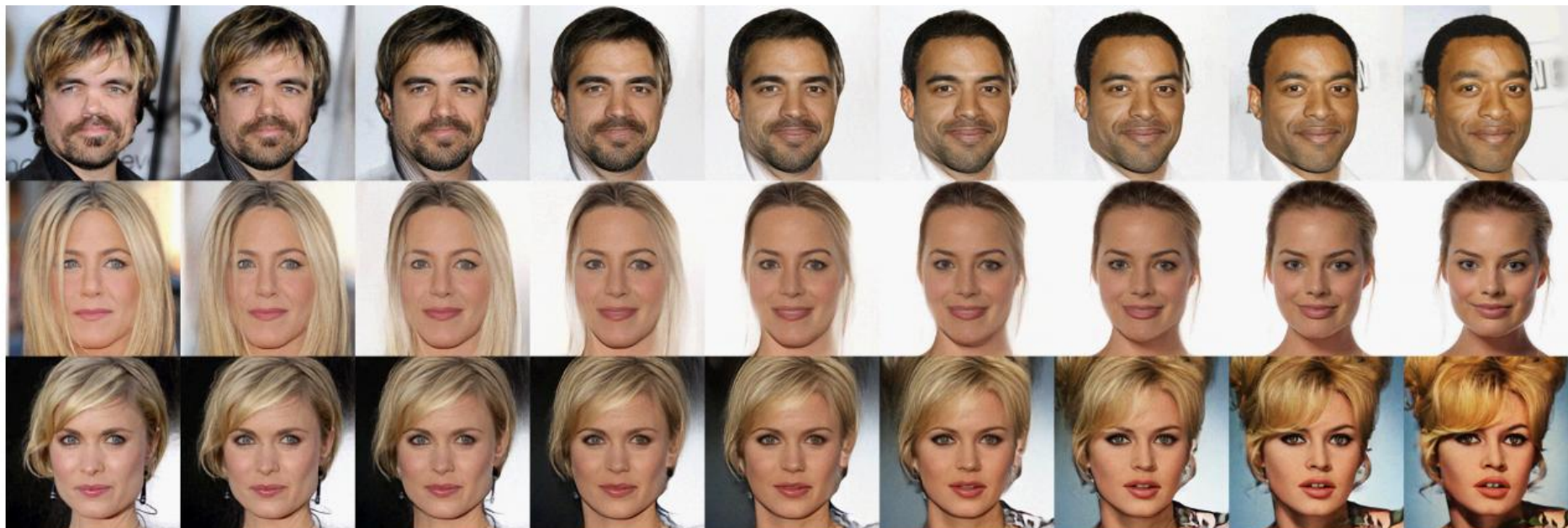
GLOW Results - Sampling



Increasing temperature

Kingma & Dhariwal, P, Glow: Generative flow with invertible 1x1 convolutions, NeurIPS 2018

GLOW Results - Interpolation



Kingma & Dhariwal, P, Glow: Generative flow with invertible 1x1 convolutions, NeurIPS 2018

GLOW Results - Manipulation



(a) Smiling

(b) Pale Skin



(c) Blond Hair

(d) Narrow Eyes

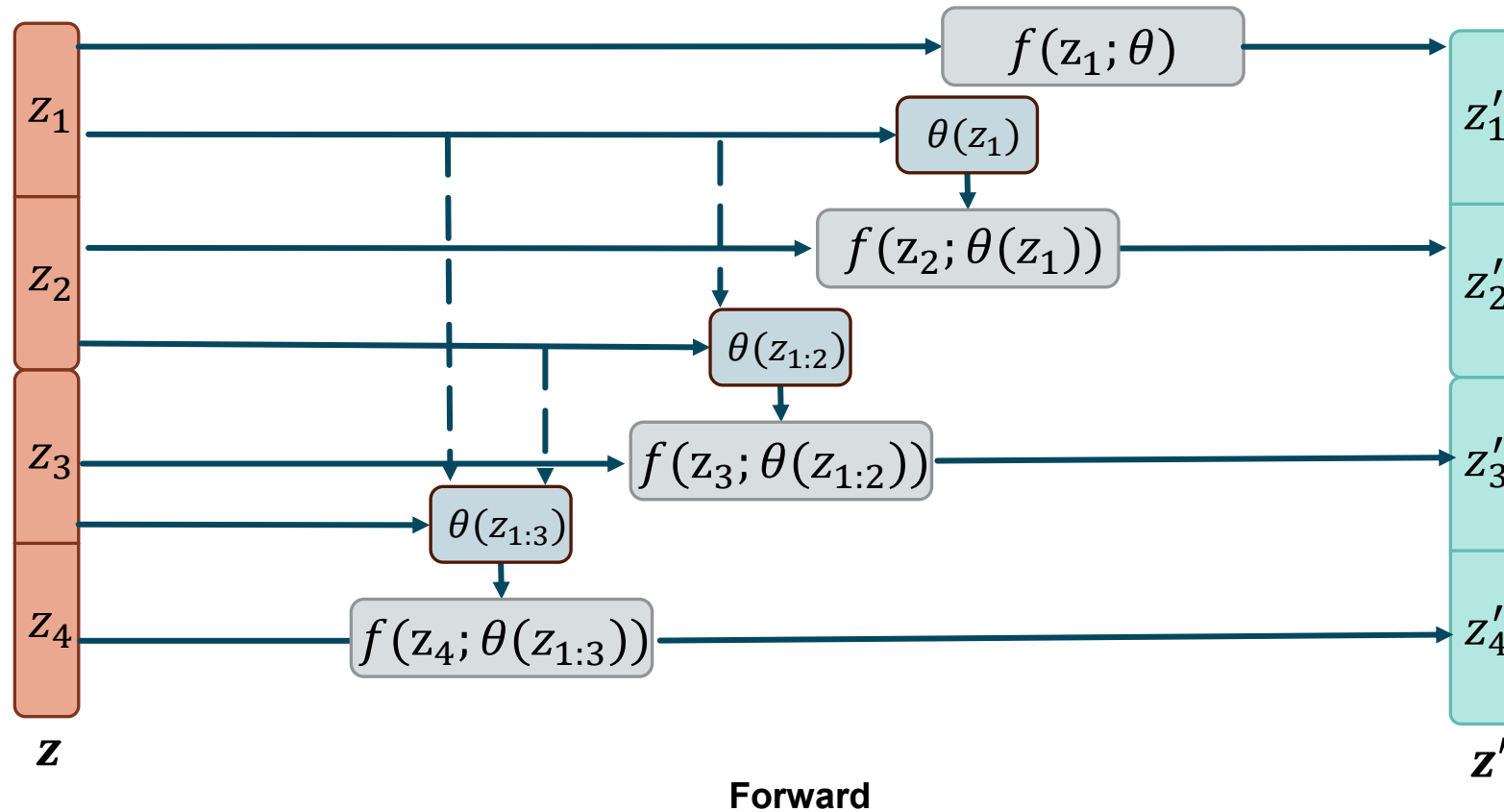


(e) Young

(f) Male

Kingma & Dhariwal,
P, Glow:
Generative flow
with invertible 1x1
convolutions,
NeurIPS 2018

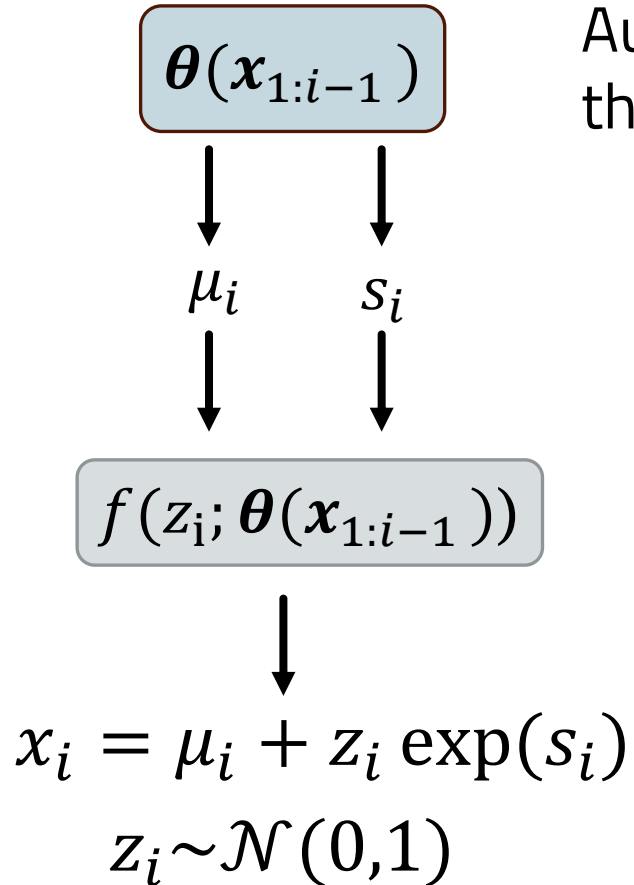
Autoregressive Flows



Generalization of coupling flows that treats each input dimension as a separate block

Forward and inverse directions have different costs (parallel/sequential)

Masked Autoregressive Flow



Autoregressive model as a transformation f from the space of random vectors \mathbf{z} to space of data \mathbf{x}

$$P(x_i | \mathbf{x}_{1:i-1}) = \mathcal{N}(\mu_i, (\exp(s_i))^2)$$

f is easily invertible, Jacobian is triangular and easily computable determinant

$$z_i = f^{-1}(x_i) = (x_i - \mu_i) \exp(-s_i)$$

$$\left| \det \left(\frac{\partial f^{-1}}{\partial \mathbf{z}} \right) \right| = \exp - \sum_i s_i$$

Residual flows

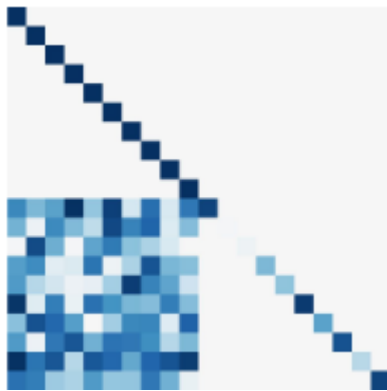
Designing scalable flows

So far, to keep Normalizing Flows scalable we have restricted our layers to have “well” behaved Jacobians

$$\log P(\mathbf{x}) = \log P(f(\mathbf{x})) + \log|\det J_x(f)|$$

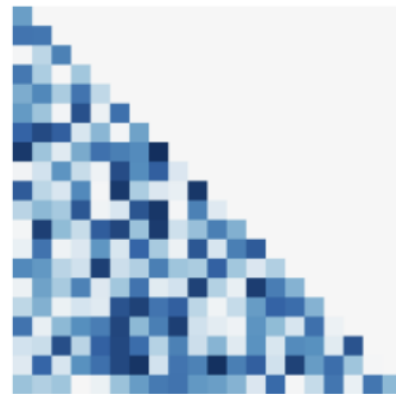
Coupling layers

NICE
Real NVP
Glow

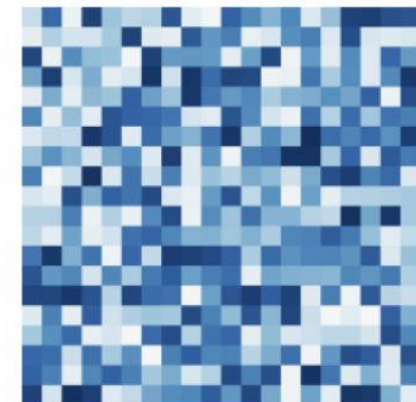


Autoregressive

Masked
Inverse
...



Can we get a more expressive (full rank) Jacobian and still remain scalable?



Residual flows

Use a **residual-based** neural layer

$$\mathbf{z}' = \mathbf{z} + \theta(\mathbf{z})$$

Intuition

- ◇ **Small updates** allow for easier invertibility
- ◇ Composing many small updates (i.e. using a deep network) yields expressive transformations

Formally

- ◇ If θ is Lipschitz with constant $L < 1$

$$\|\theta(\mathbf{z}) - \theta(\mathbf{z}')\| \leq L \|\mathbf{z} - \mathbf{z}'\|$$

- ◇ Then by Banach fixed-point $\mathbf{z}' = \mathbf{z} + \theta(\mathbf{z})$ is invertible

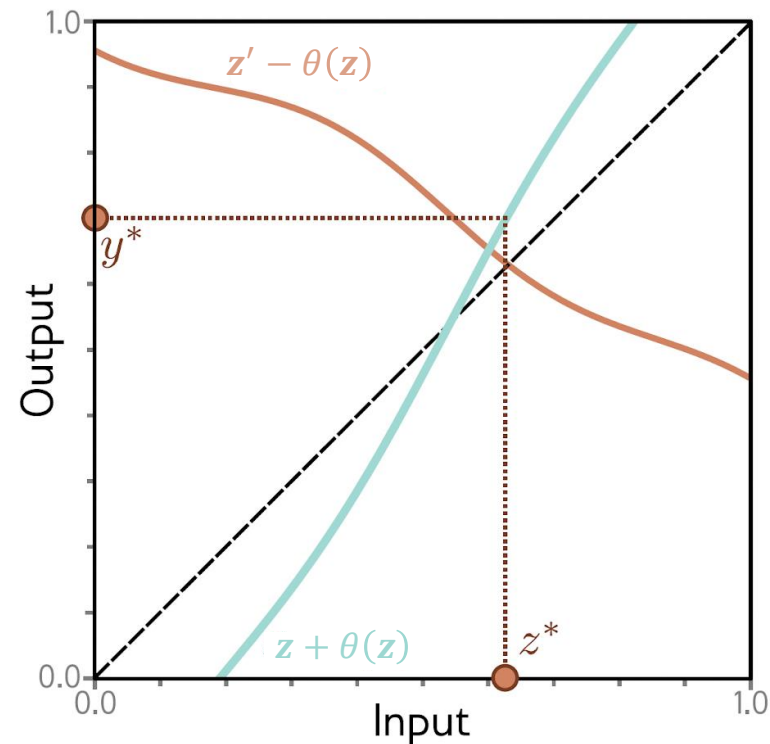
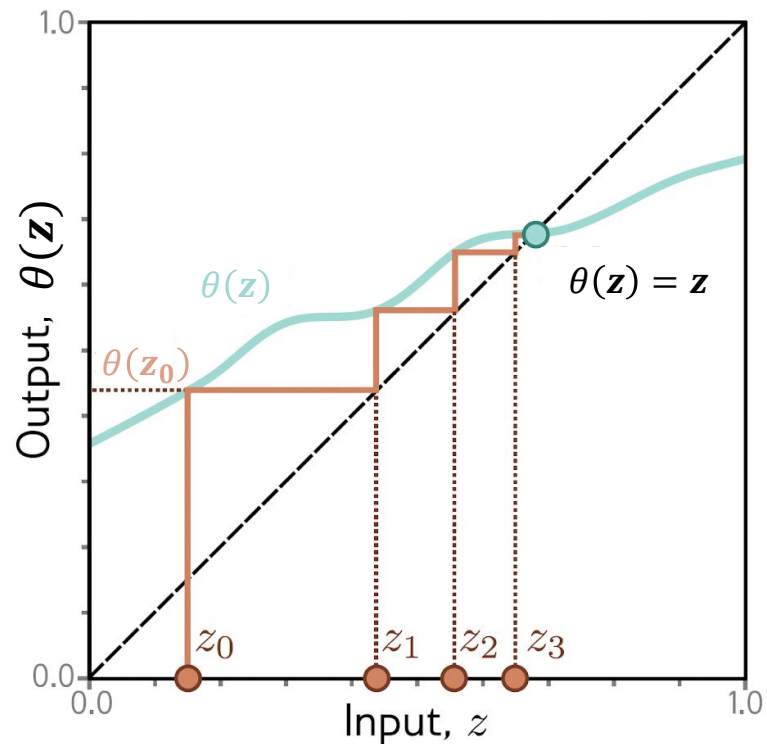
In practice

The inverse can be obtained via **fixed-point iterative method**

$$\mathbf{z}^{(i)} = \mathbf{z}' - \theta(\mathbf{z}^{(i-1)})$$

until it converges to the \mathbf{z}^* that is the starting point for \mathbf{z}'

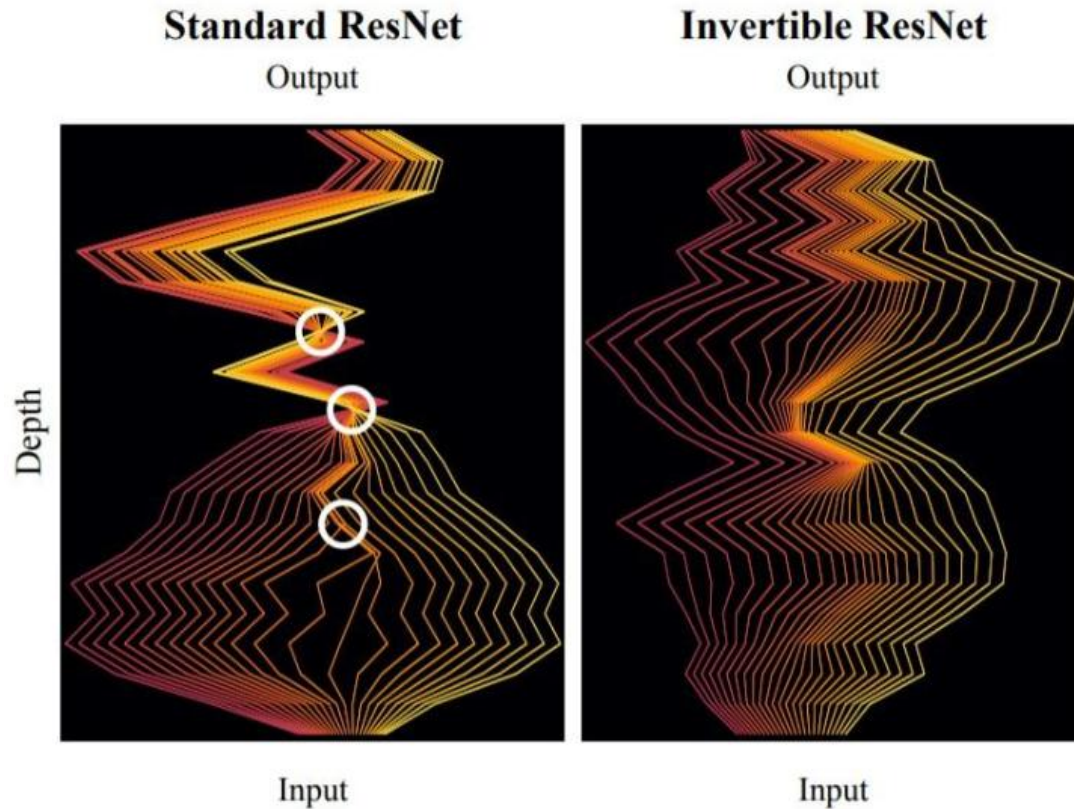
Contraction Mappings



When a contractive mapping is iterated (i.e., the output is repeatedly passed back in as an input), the result converges to a fixed point where $\theta(z) = z$

This can be done by starting with any point z_0 and iterating $z^{(i)} = z' - \theta(z^{(i-1)})$, which has a fixed point at $z + \theta(z) = z'$

Invertible ResNets



Obtaining invertible residual layers entails

- ◇ Choosing Lipschitz-constrained activation functions (i.e. $\sigma'(\cdot) < 1$)
- ◇ Bounding the spectral radius of weight matrices

Invertible ResNets and their Jacobian

- ◆ For an invertible residual change of variable

$$\mathbf{z}' = f(\mathbf{z}) = \mathbf{z} + \theta(\mathbf{z})$$

- ◆ The target **log-likelihood rewrites** as

$$\log P(\mathbf{z}) = \log P(\mathbf{z}') + \log |\det J_{\mathbf{z}}(f)| = \log P(\mathbf{z}') + \underbrace{\left(\sum_{k=1}^{\infty} \frac{(-1)^{k+1}}{k} \text{tr}(\mathbf{J}_{\mathbf{z}}(\boldsymbol{\theta}))^k \right)}$$

Can be approximated efficiently (linearly) using the Hutchinson trace estimator $\text{tr}(\mathbf{A}) = \mathbb{E}[\boldsymbol{\epsilon}^T \mathbf{A} \boldsymbol{\epsilon}]$ with $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$

One final twist

Let's reformulate the invertible residual adding a scaling $\delta > 0$

$$\mathbf{z}' = f(\mathbf{z}) = \mathbf{z} + \delta\theta(\mathbf{z})$$

Rearranging terms

$$\theta(\mathbf{z}) = \frac{\mathbf{z}' - \mathbf{z}}{\delta}$$

which looks a lot like a [derivative](#).

We can make the thing more explicit by defining the [infinitesimal update](#) (for $\Delta t = \delta \rightarrow 0$)

$$\mathbf{z}_{t+\delta} = \mathbf{z}_t + \delta \frac{\partial \mathbf{z}_t}{\partial t} \text{ s.t. } \underbrace{\frac{\partial \mathbf{z}_t}{\partial t}}_{\text{Instantaneous change}} = \theta_t(\mathbf{z}_t)$$

Instantaneous change

Continuous Normalizing Flows

Now suppose we are **cascading multiple residual layers** $f_i(\cdot)$, yielding the composition

$$f_T \circ f_{T-1} \circ \dots \circ f_1$$

Using the **instantaneous update** defined in previous slide this can be represented by an **Ordinary Differential Equation (ODE)**

A neural ODE \longrightarrow
$$\begin{cases} \mathbf{z}_0 \sim N(0,1) & \text{Initial conditions} \\ \frac{\partial \mathbf{z}_t}{\partial t} = \theta_t(\mathbf{z}_t) & \text{Instantaneous update (vector field)} \end{cases}$$

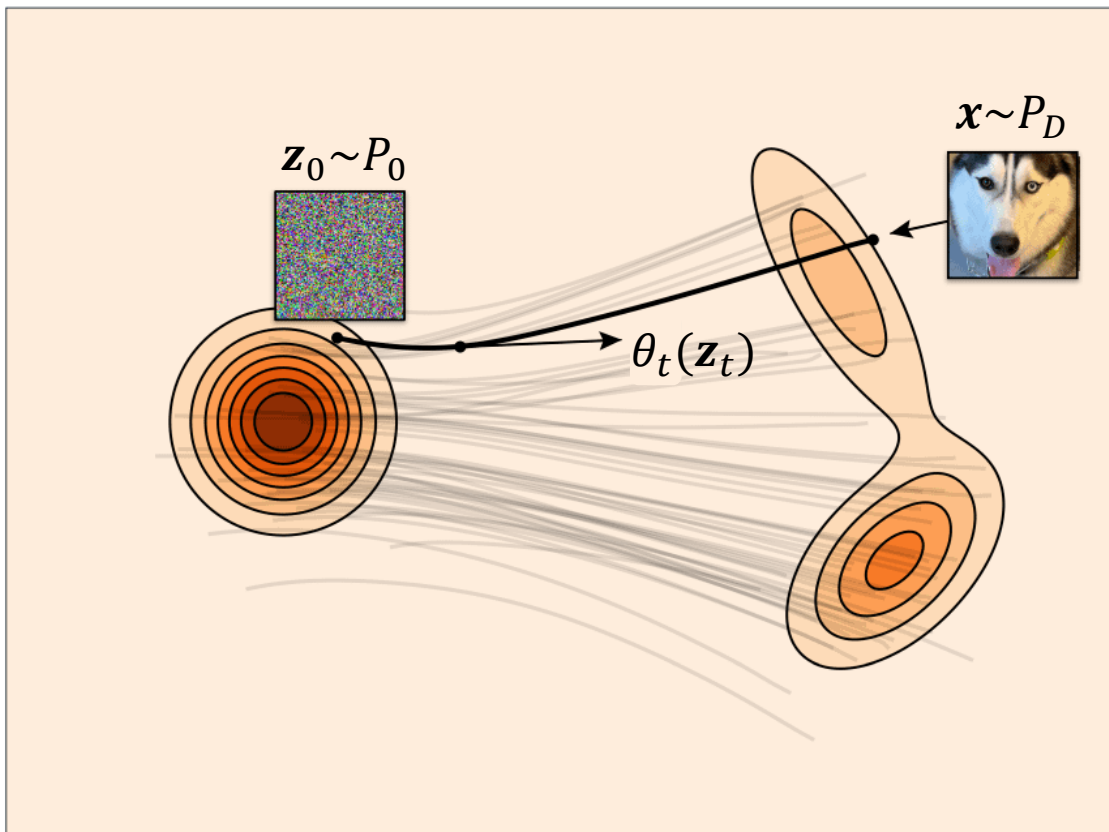
Plugging the continuous formulation above in the change of variable setting of Normalizing Flows gives us the loglikelihood for **instantaneous change of variables**

$$\log P(\mathbf{x}) = \log P(\mathbf{z}_0) - \int_0^T \text{Tr} \left(\mathbf{J}_{\mathbf{z}_t}(\theta_t) \right) dt$$

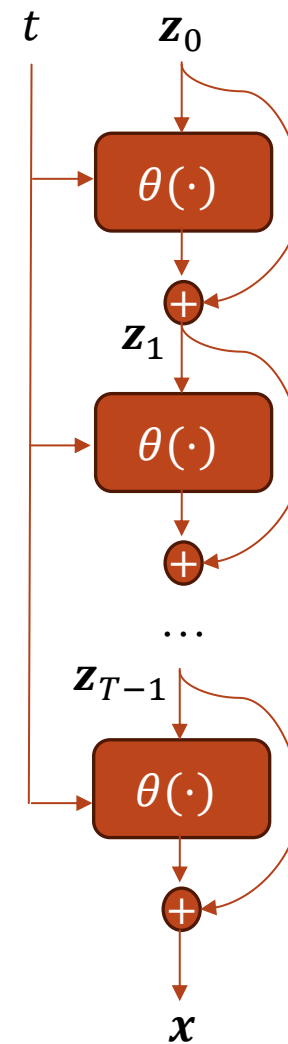
No Jacobian determinant!
Change in log probability is related to the trace of the derivative of the forward propagation

Which can be **solved using numerical integration**

Morphing distributions by continuous change of variable



Continuous change of variable trajectories



Discretized
ODE as a
ResNet

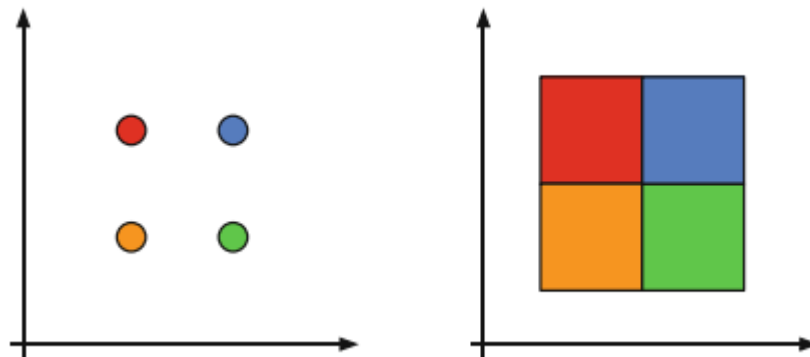
Wrap-up

Implementations & Libraries

- ◇ Normalizing flows are natively supported by **Tensorflow** (through the TF Probability module)
 - ◇ `tf.probability.distribution` (for base distributions)
 - ◇ `tf.probability.bijector` (for predefined layers, e.g. masked autoregressive)
 - ◇ `Chain()` (to chain bijectors and compose complex modules)
 - ◇ You can of course define you own bijectors according to a template
- ◇ [Normflows](#) - PyTorch package for Normalizing Flows
- ◇ [Flowtorch](#) – PyRo based Pytorch library for Normalizing Flows

A pragmatic note - Dequantization

- ◇ Flow-based models assume that \mathbf{x} is a vector of real-valued random variables
- ◇ In practice, \mathbf{x} is often many objects discrete valued (e.g. pixel colors have values in $\{0, \dots, 255\}$)
- ◇ [Uniform dequantization](#) - Add uniform noise $\mathbf{u} \in [-0.5, 0.5]^d$ to original data $\mathbf{x}' = \mathbf{x} + \mathbf{u}$
- ◇ Dequantization is supported by TF/Torch APIs



Take Home Messages

- ◇ Normalizing flows as an effective and tractable way to generate new samples (**efficient**) and to evaluate the likelihood of samples (**not so efficient**)
- ◇ **Universality property** - The flow can learn any target density to any required precision given sufficient capacity and data
 - ◇ Flow can be used to generate samples that **approximate a density** easy to evaluate but difficult to sample
- ◇ Normalizing flow design needs to take care of
 - ◇ Keeping flow invertible and efficient
 - ◇ Making the determinant of the Jacobian easy to compute
- ◇ Normalizing flows can be made **continuous** using a **neural ODE** scheme
 - ◇ Can be implemented using ResNets
 - ◇ Expressive and powerful: can parameterize a large class of probability distributions
 - ◇ Can be **extremely slow due to ODE numerical integration** **We may be doing something about this in a couple of lectures**

Next Lecture

No lecture on 05/05/2026

Diffusion models

- ◇ Generative learning as a noising-denoising process
- ◇ Guided diffusion
- ◇ Conditional diffusion
- ◇ Latent space diffusion