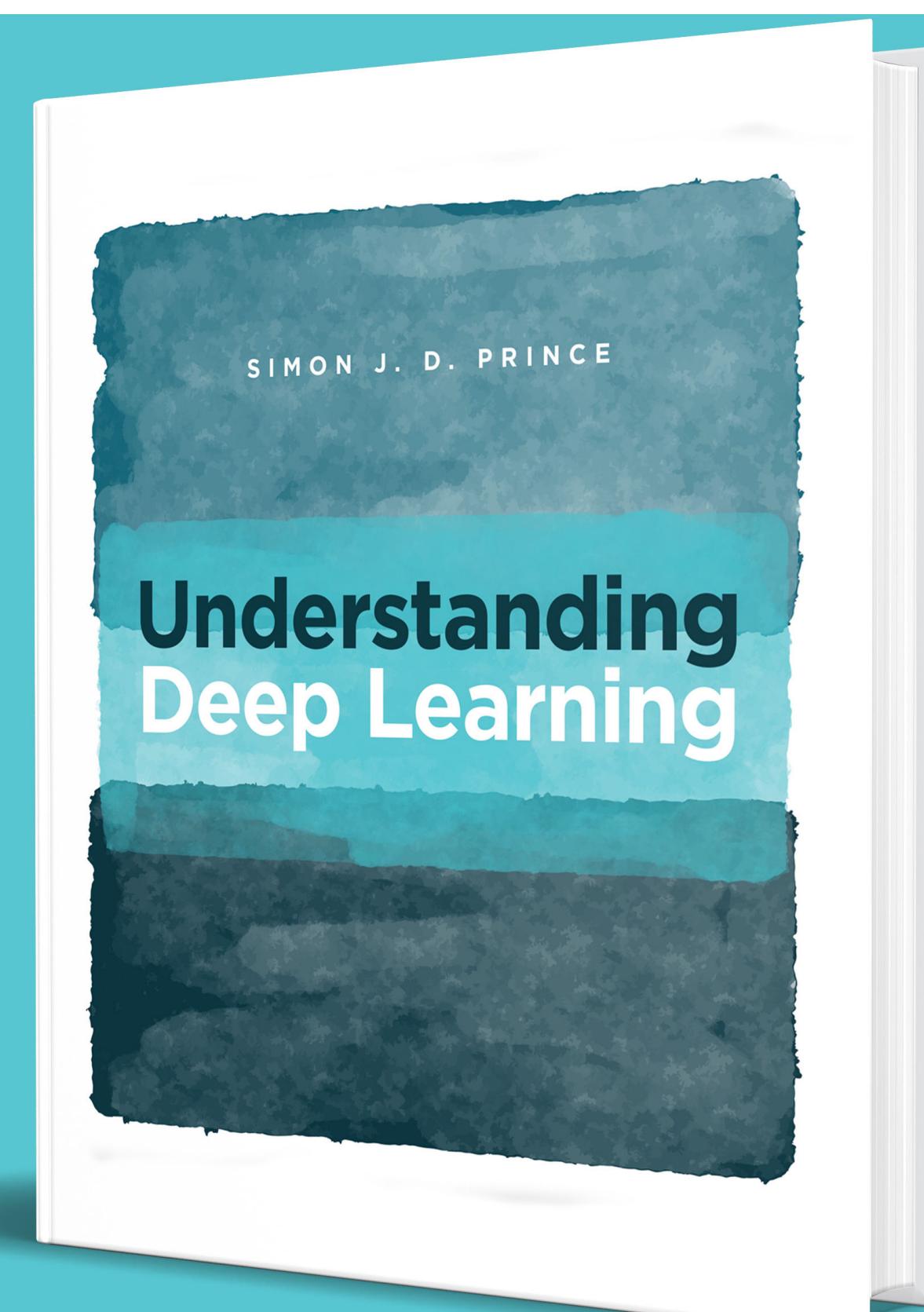




Lecture 12: normalizing Flows



Program for Today

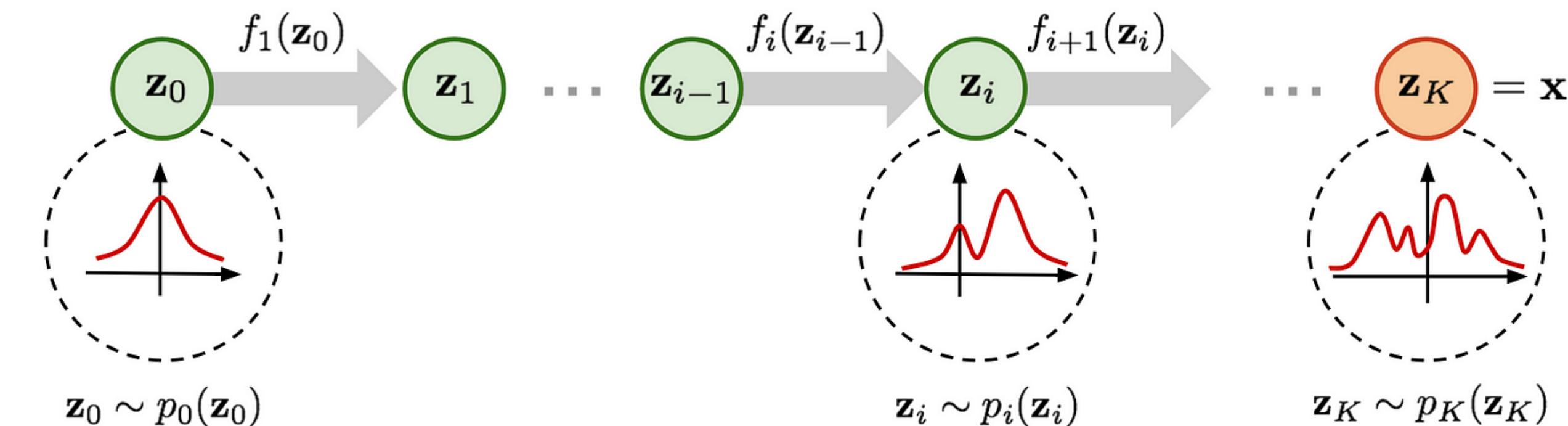


A	B	C
Date	Topic	Tutorial
Sep 16	Intro & class description	
Sep 23	probability Theory + Basic of machine learning + Dense NN	<i>Basic jupyter + DNN on mnist (give jet dnn as homework)</i>
Sep 30	Statistics + Convolutional NN	<i>Convolutional NNs with MNIST</i>
Oct 7	Training in practice: regularization, optimization, etc	<i>Free Practice</i>
Oct 14	Unsupervised learning and anomaly detection	<i>Autoencoders with MNIST</i>
Oct 21	Graph NNs	<i>Tutorial on Graph NNs</i>
Oct 28	Transformers	
Nov 4	Network compression (pruning, quantization, Knowledge Distillation)	
Nov 11		<i>Tutorial on hls4ml</i>
Nov 18	Generative models: GANs, VAEs, etc	<i>GAN and VAE</i>
Nov 25		<i>Reinforcement Learning</i>
Dec 2	Normalizing flows	<i>Normalizing flows</i>
Dec 9		<i>Quantum Machine Learning</i>
Dec 16		<i>Exams To Be Confirmed</i>



Going Beyond GANs (and VAEs)

- GANs are trained to create new realistic examples learning from a given dataset
- But they gain no knowledge on the distribution over the data sample
- Because of that, they might generate images of one kind more than another, w/o respecting the relative population of the reference dataset
- We saw that VAEs are trained to “force” this knowledge, imposing a certain prior in the latent space through the KL regularization. But they have other issues (e.g., blurry images)
- **Normalizing Flows** address the sampling problem:
 - They learn to transform an easy-to-control distribution (e.g., Gaussian) into an arbitrary one (i.e., they fit the reference data distribution with a neural network)
 - **GENERATORS:** One can then sample a generic function, sampling the easy one and applying the morphing network
 - **ANOMALY DETECTION:** One can associate a likelihood p-value to a given event



Invertible Change of Variable

- We change variable from z to x , transporting an infinitesimal of probability density using a bijector (invertible and differentiable function)

$$x = f(z)$$

$$\int \Pi(z) dz = \int \Pi(z) \left| \frac{dz}{dx} \right| dx = \int \Pi(f^{-1}(x)) \left| \frac{df^{-1}(x)}{dx} \right| dx = \int P(x) dx$$

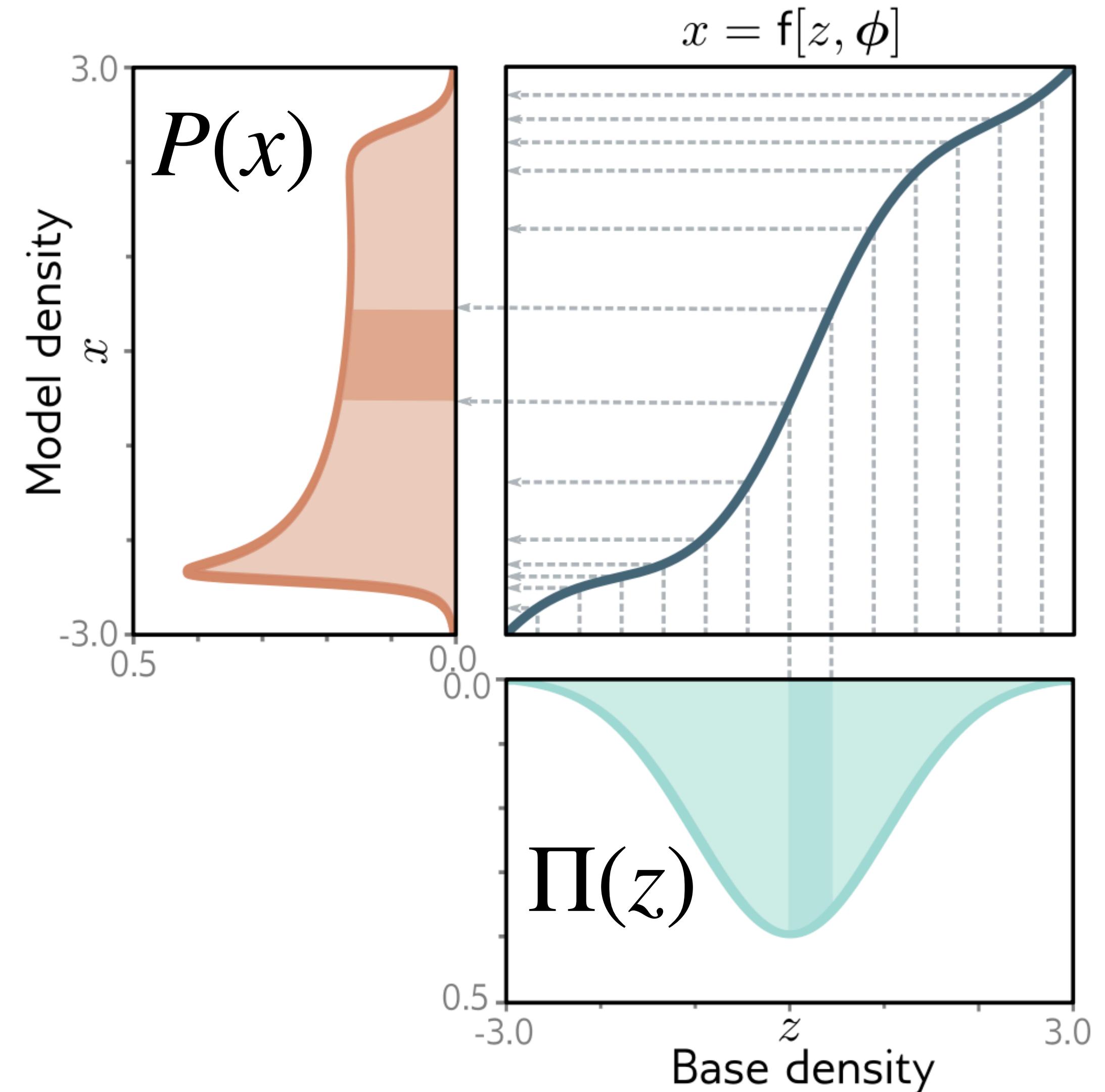
- But this implies

$$P(x) = \Pi(f^{-1}(x)) \left| \frac{df^{-1}(x)}{dx} \right| = \Pi(z) \left| \frac{df(z)}{dz} \right|^{-1}$$

where the last equality follows from the inverse function theorem

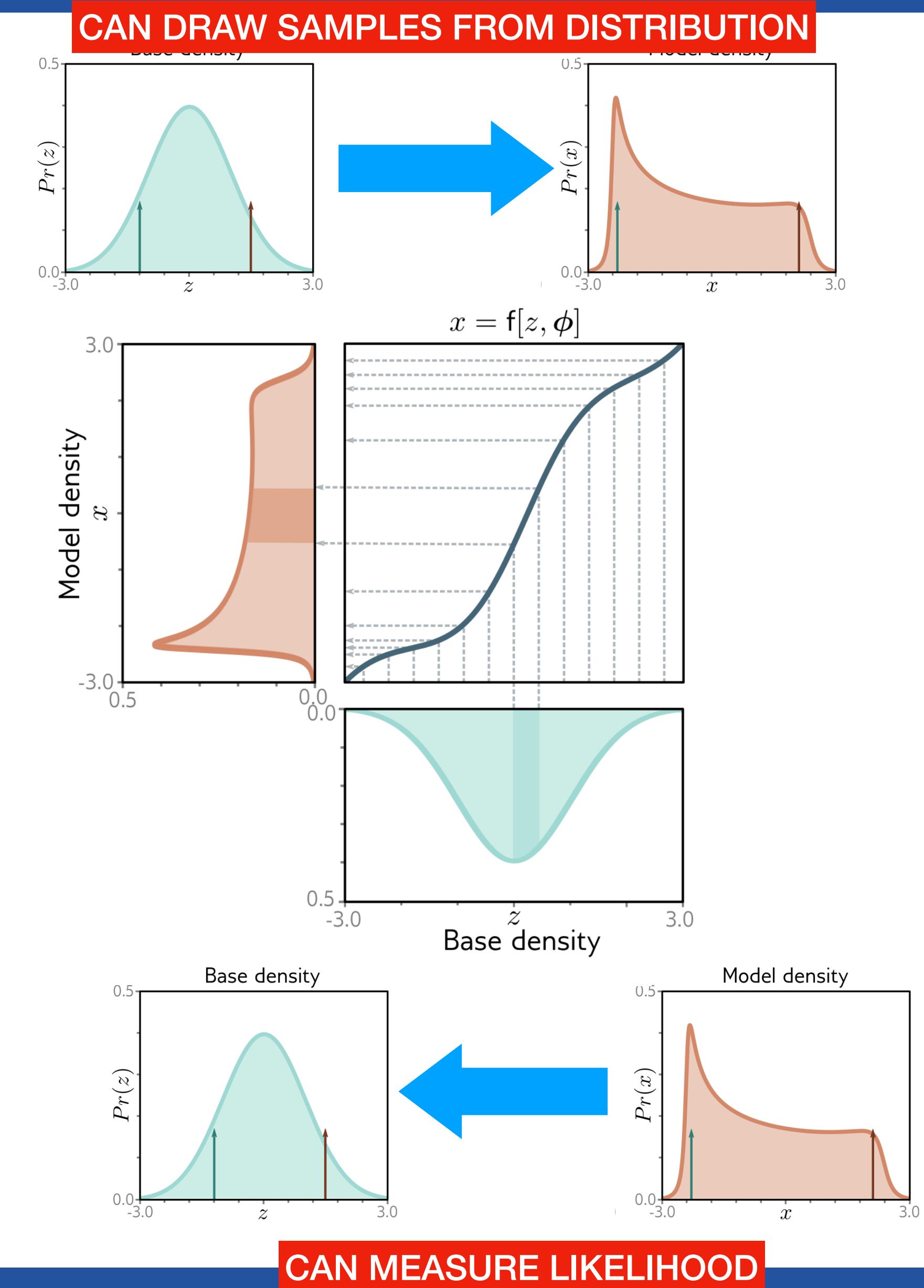
- For normalizing flows $x = f(z|\phi)$, where ϕ are the parameters of the neural network f .

This implies that $P(x|\phi) = \Pi(z) \left| \frac{\partial f(z|\phi)}{\partial z} \right|^{-1}$



1D normalizing Flow

- Start with a Gaussian $G(z)$
 - Apply a transformation $x = f(z)$ so that x is distributed according to the desired function
- $$P(x|\phi) = \left| \frac{\partial f(z, \phi)}{\partial z} \right|^{-1} \Pi(z)$$
- Notice that depending on the Jacobian being $>$ or < 1 , probability is compressed or stretched
 - The transformation is reversible \rightarrow one can do the opposite and map the reference distribution back to the latent one
 - This is useful for anomaly detection: can associate a probability (technically, a p-value) to a point in x without knowing $P(x)$





Training a Normalizing Flow

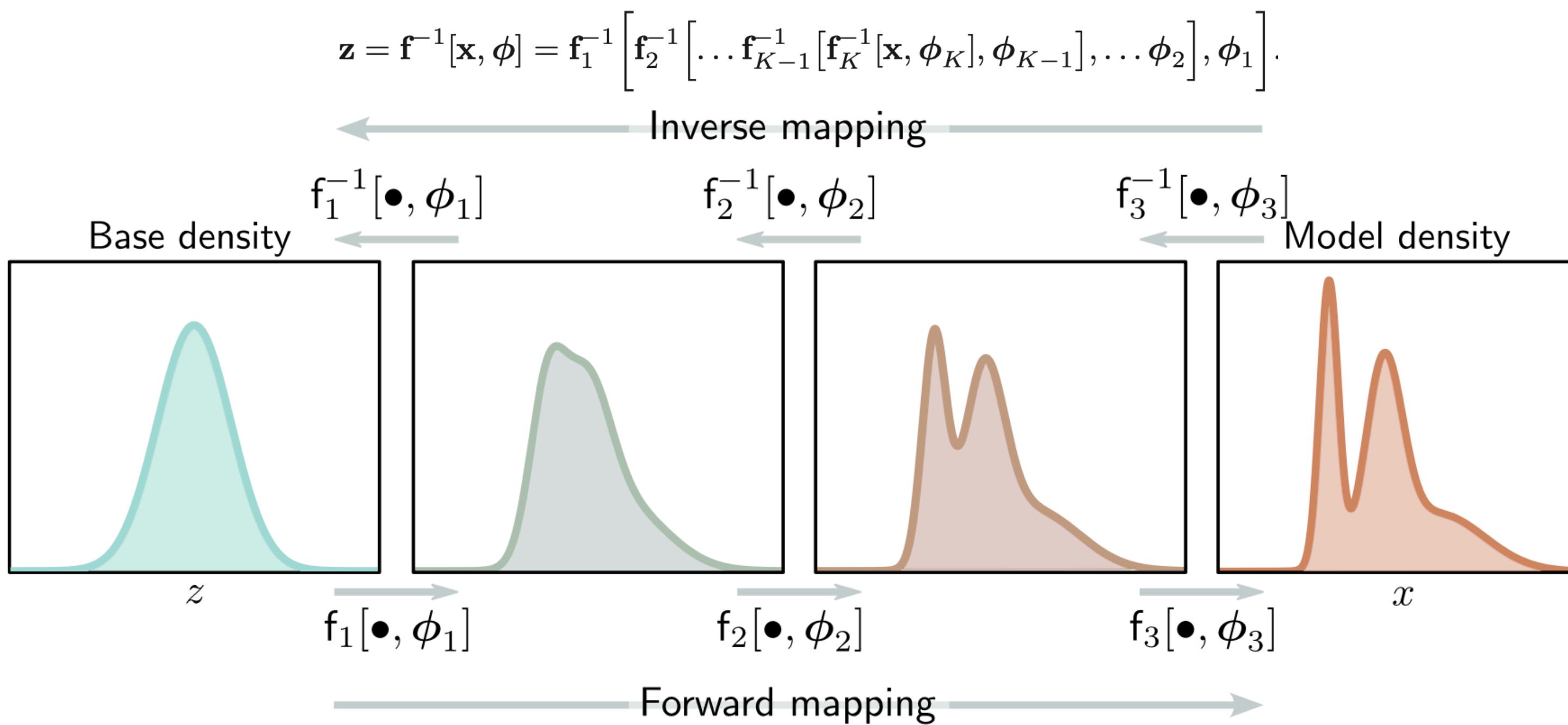
- A Normalizing Flow is trained looking for the parameters of the network that maximize the likelihood of the training data distribution $P(x)$
- It is just a likelihood fit, but done through NN

$$\hat{\phi} = \operatorname{argmax}_{\phi} \left[\prod_i P(x_i | \phi) \right] = \operatorname{argmin}_{\phi} \left[-\log \prod_i P(x_i | \phi) \right] = \operatorname{argmin}_{\phi} \left[\log \prod_i \left| \frac{\partial f(z_i, \phi)}{\partial z_i} \right| - \log \prod_i \Pi(z_i) \right] =$$
$$\operatorname{argmin}_{\phi} \left[\sum_i \log \left| \frac{\partial f(z_i, \phi)}{\partial z_i} \right| - \sum_i \log \Pi(z_i) \right]$$

- One can easily generalize this to the case of N-dim, replacing the abs value of the partial derivative with the determinant of the Jakobian

A Deep Normalizing Flow

- A deep NF is a chain of NFs, gradually morphing the latent function $\Pi(z)$ into the target function $P(x)$



$$\mathbf{x} = \mathbf{f}[\mathbf{z}, \phi] = \mathbf{f}_K\left[\mathbf{f}_{K-1}\left[\dots \mathbf{f}_2\left[\mathbf{f}_1[\mathbf{z}, \phi_1], \phi_2\right], \dots \phi_{K-1}\right], \phi_K\right].$$

$$\frac{\partial \mathbf{f}[\mathbf{z}, \phi]}{\partial \mathbf{z}} = \frac{\partial \mathbf{f}_K[\mathbf{f}_{K-1}, \phi_K]}{\partial \mathbf{f}_{K-1}} \cdot \frac{\partial \mathbf{f}_{K-1}[\mathbf{f}_{K-2}, \phi_{K-1}]}{\partial \mathbf{f}_{K-2}} \cdots \frac{\partial \mathbf{f}_2[\mathbf{f}_1, \phi_2]}{\partial \mathbf{f}_1} \cdot \frac{\partial \mathbf{f}_1[\mathbf{z}, \phi_1]}{\partial \mathbf{z}},$$

TRAINING

$$\begin{aligned}\hat{\phi} &= \underset{\phi}{\operatorname{argmax}} \left[\prod_{i=1}^I Pr(\mathbf{z}_i) \cdot \left| \frac{\partial \mathbf{f}[\mathbf{z}_i, \phi]}{\partial \mathbf{z}_i} \right|^{-1} \right] \\ &= \underset{\phi}{\operatorname{argmin}} \left[\sum_{i=1}^I \log \left| \frac{\partial \mathbf{f}[\mathbf{z}_i, \phi]}{\partial \mathbf{z}_i} \right| - \log [Pr(\mathbf{z}_i)] \right]\end{aligned}$$

The NF Architecture

- So far we did not specify what f looks like
- The choice of f , i.e., of the network architecture, should be made so that certain properties of f are obtained
- The set of f functions chained in the deep NF should be expressive enough to be capable of learning many distributions
- Each layer should be invertible (1-to-1 mappings)
- The inverse should be computable efficiently
- The determinant of the Jacobian should be computed efficiently

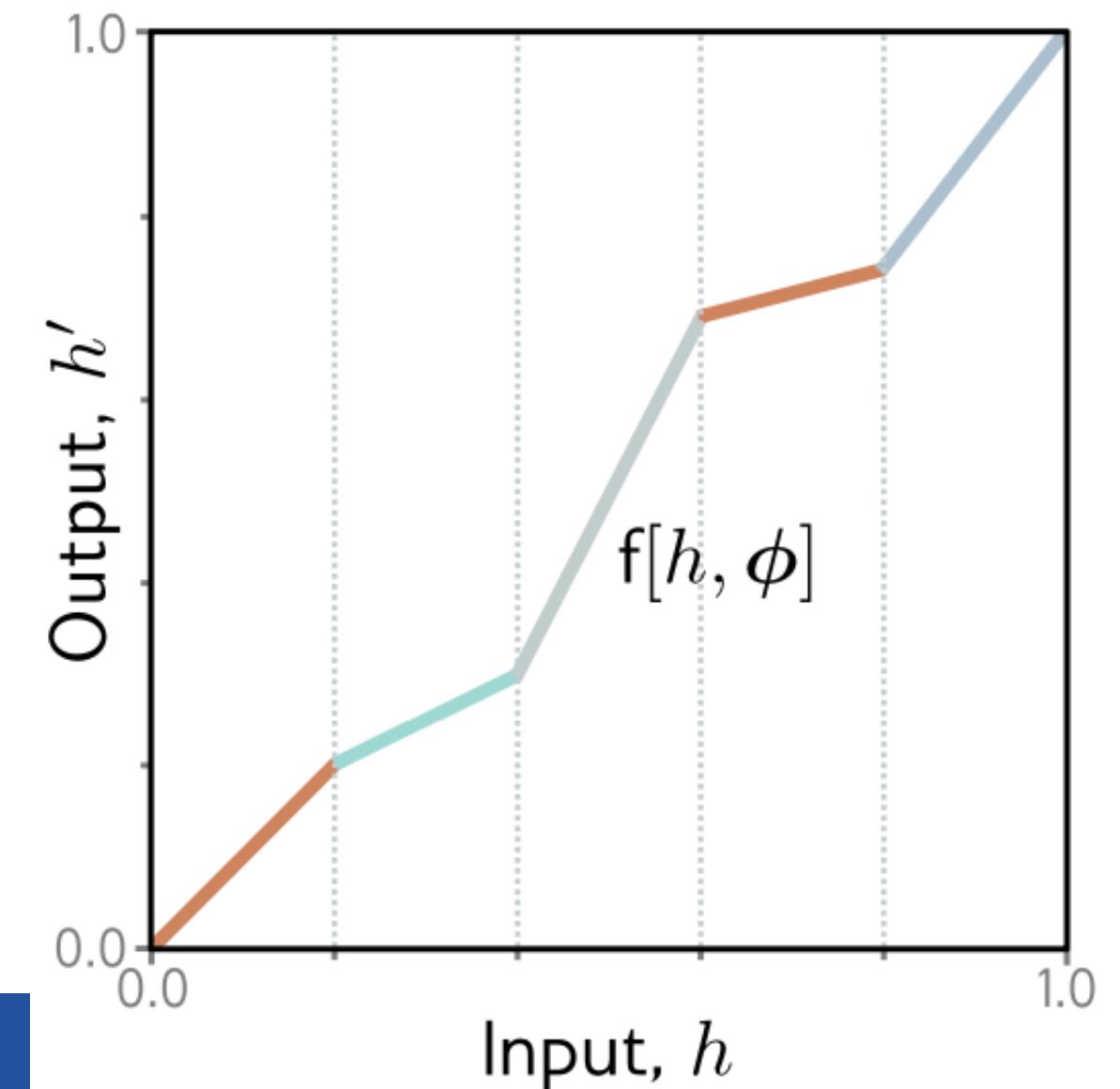
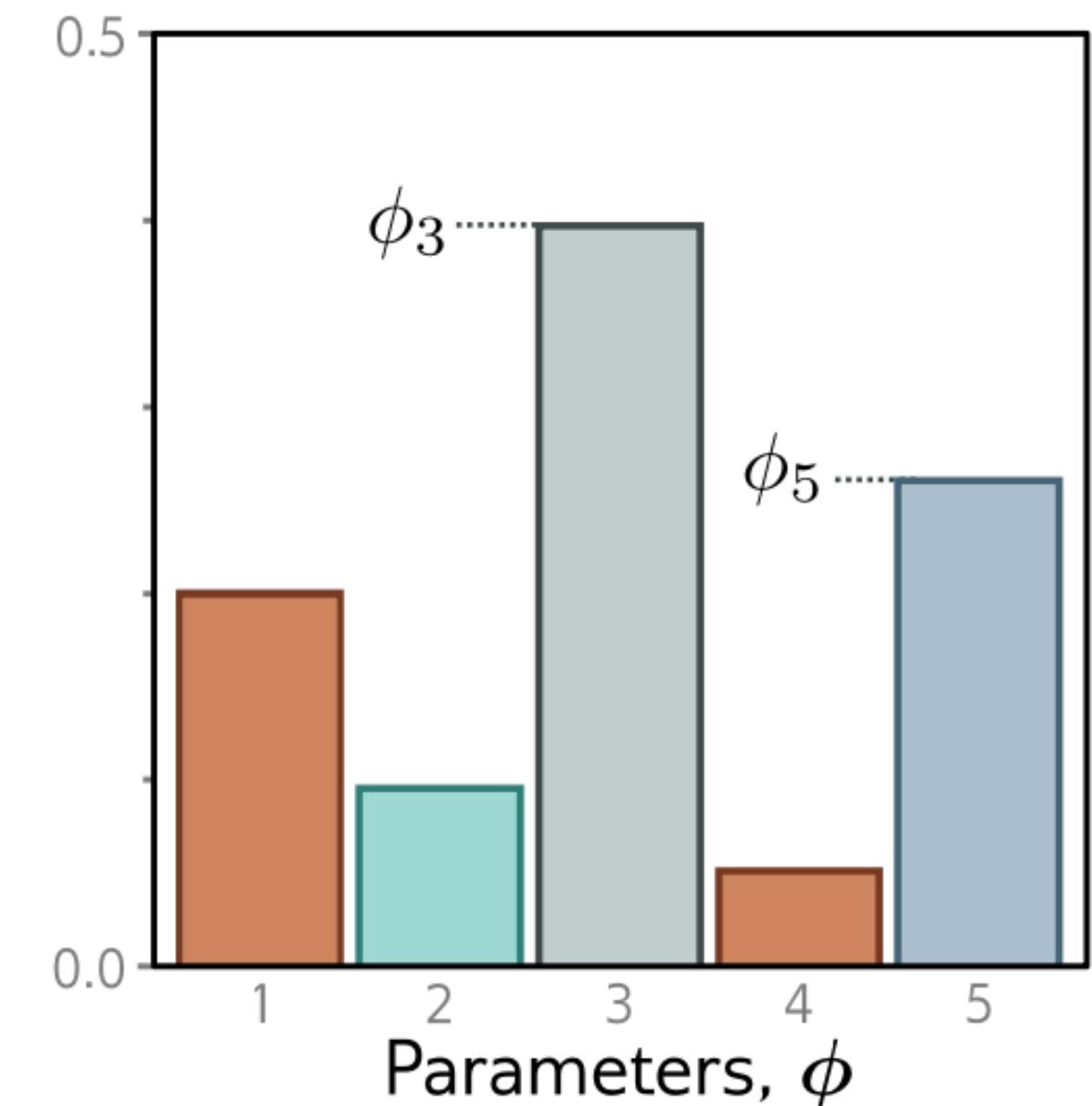
Linear Flows

- As the name suggests, a linear flow is a transformation $f[z] = \Omega z + \beta$, which is invertible if Ω is invertible
- Linear flows are expensive at large dimensions D
- Computing the determinant and the inverse are $\mathcal{O}[D^3]$ slow, less if matrices are special (triangular, diagonal, etc.)
- Linear flows are not very expressive: they can't learn non-linear transformations and a Gaussian stays a Gaussian

Elementwise flows

- Elementwise flows are the simplest examples of non-linear flows
- They correspond to apply the same non linear functions to each element of the input
- The Jacobian is diagonal, since there is no cross-talk between the elements in the transformation. So the determinant calculation is easy
- It adds non-linearity, but it cannot alter correlations between input quantities
- The function could be any invertible non-linear function, as simple as LeakyReLu (no trainable parameters) or a generalization of that, with trainable slope (piecewise linear function)

$$f[h, \phi] = \left(\sum_{k=1}^{b-1} \phi_k \right) + (hK - b)\phi_b$$



Coupling flows

- Couplings flows operate as follows

- Split input h in two vectors h_1 and h_2

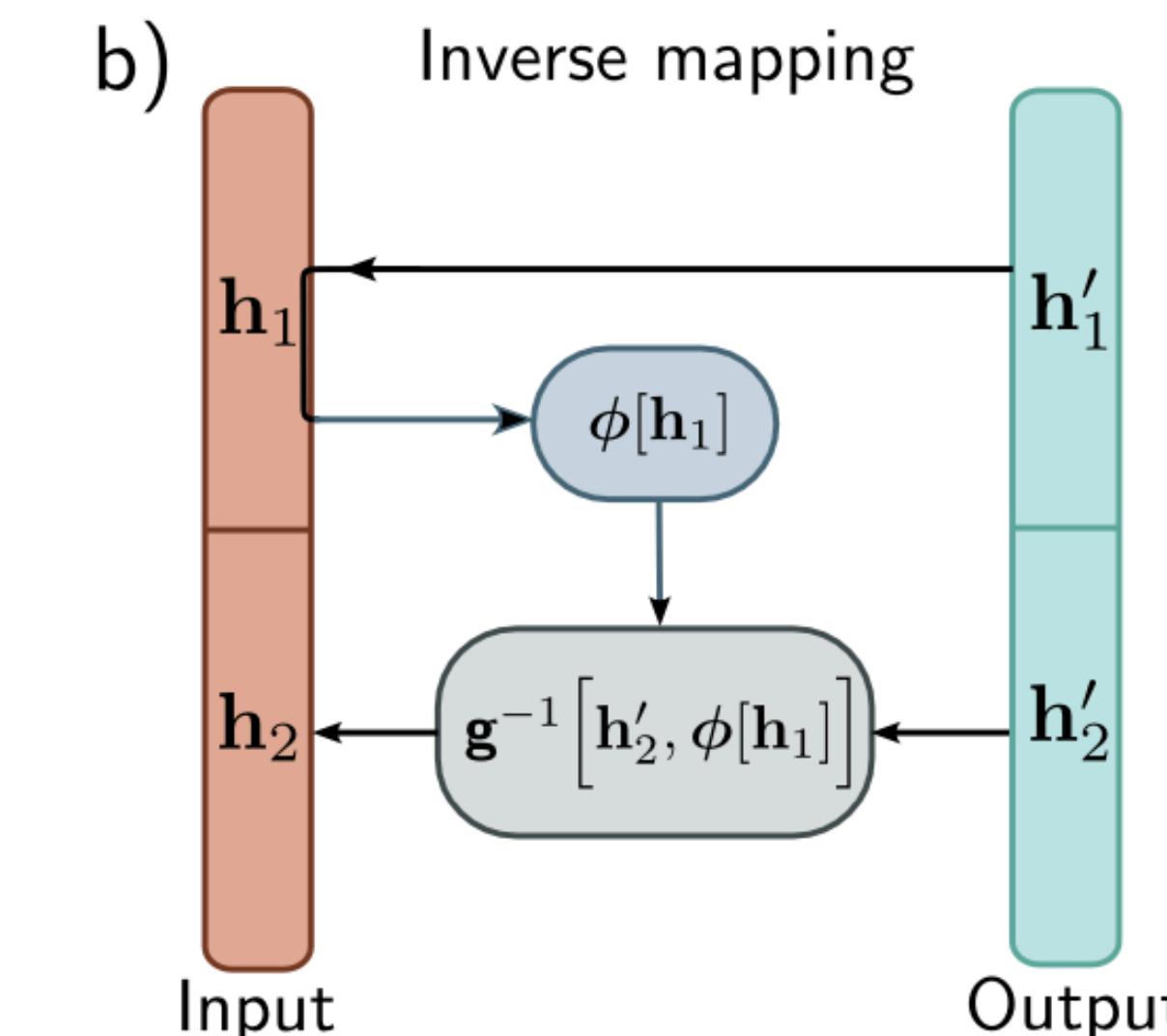
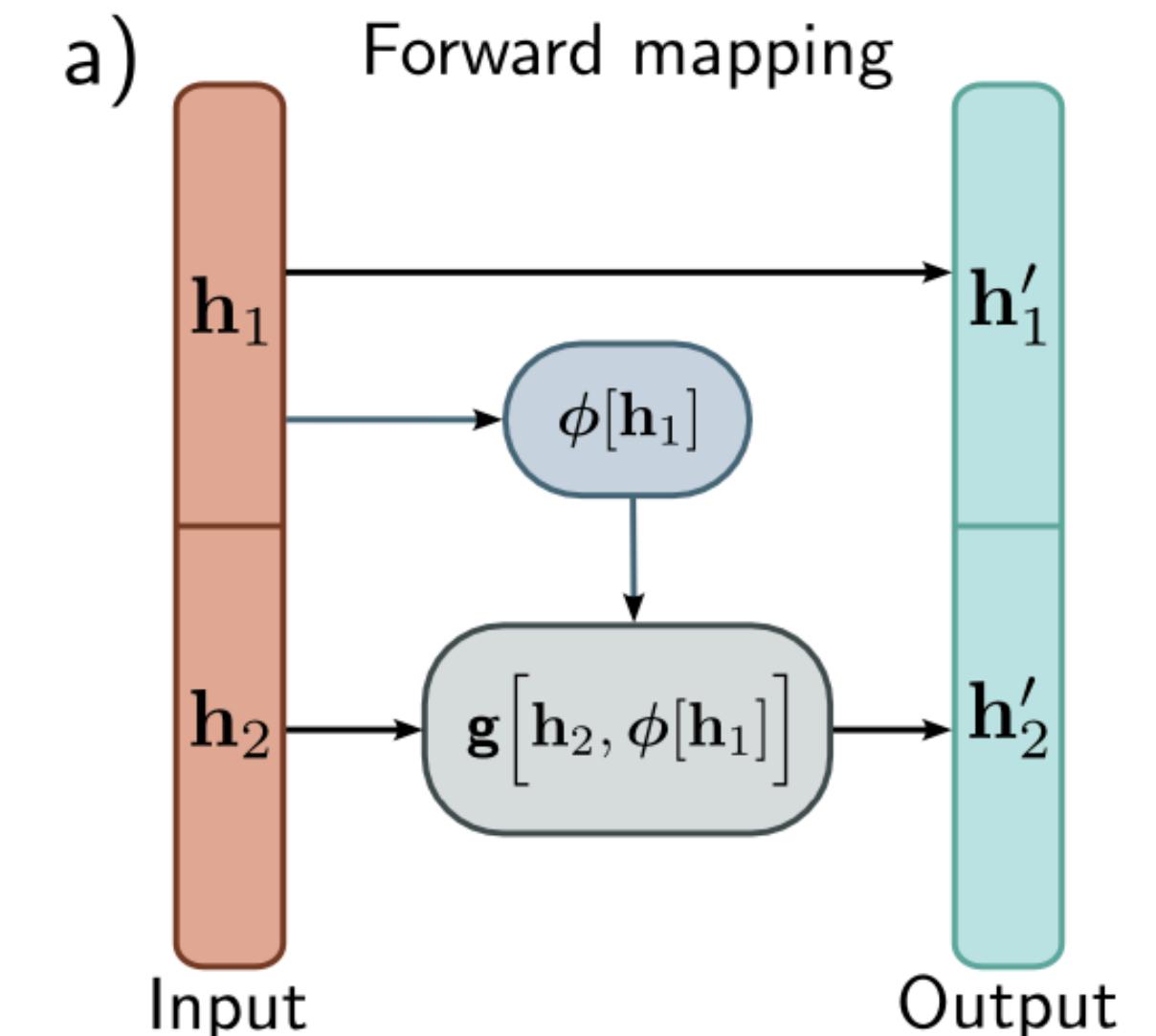
- h_1 is passed to output as it is

- h_1 is used to compute some vector of parameters $\phi[h_1]$

- $\phi[h_1]$ are used as the parameters of an elemetwise (or any other invertible) flow that acts on h_2

- Jacobian is triangular (h'_2 depends on h_1 , but h'_1 does not depend on h_2)

- One usually shuffles the elements in h before splitting it into h_1 and h_2 at every layer, so that going deep in the chain every input lands in one of the h_2 (so that everything is transformed)



$$h'_1 = h_1$$

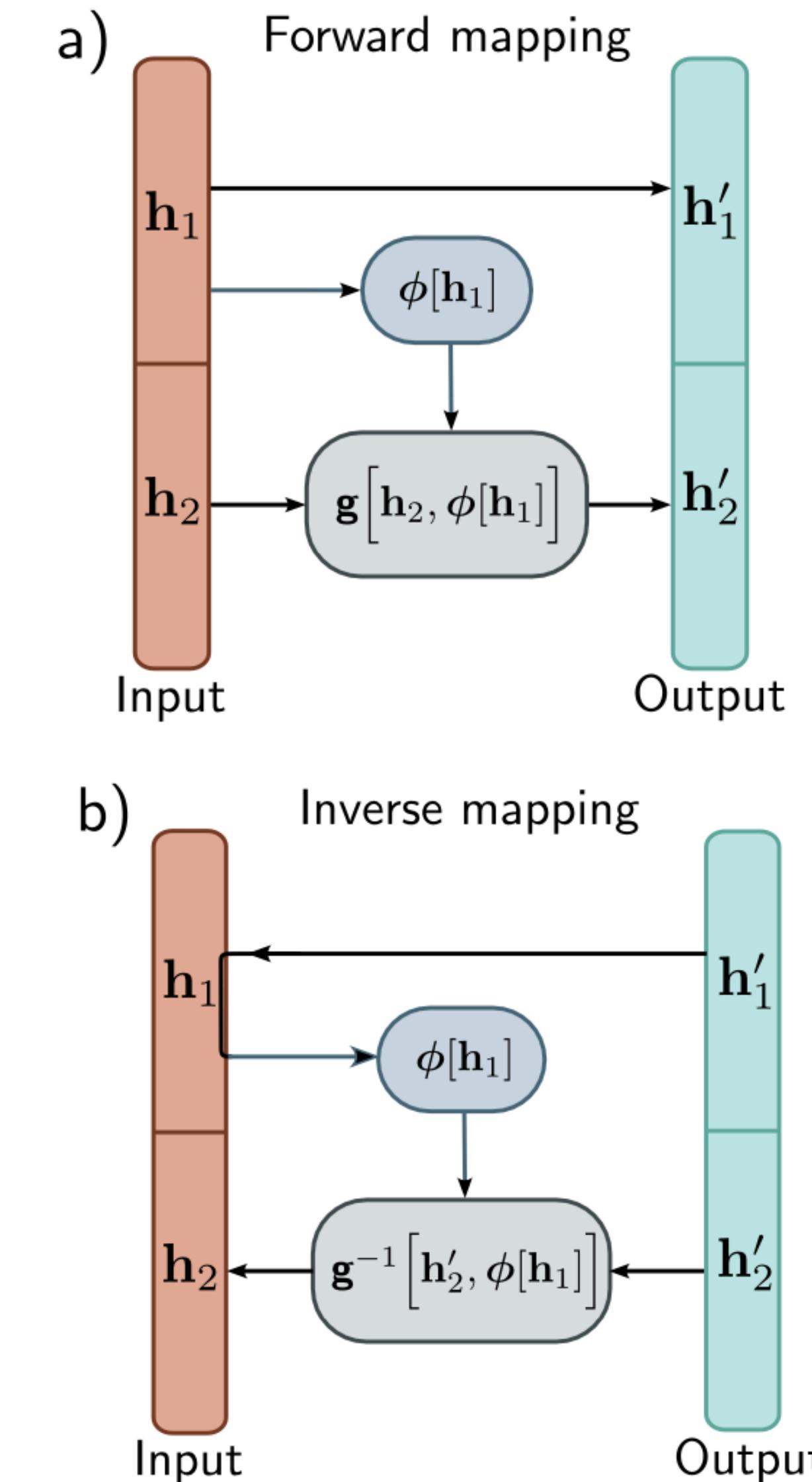
$$h'_2 = g[h_2, \phi[h_1]]$$

$$h_1 = h'_1$$

$$h_2 = g^{-1}[h'_2, \phi[h_1]]$$

- Real-valued Non-Volume Preserving (or RealNVP) are special kinds of flows in which
- a mask is used to shuffle the inputs
- The affine transformation g is then obtained
- Computing the output of two dense networks
 - $s(h_1)$ is the scale to apply on h_2
 - $t(h_1)$ is the shift to apply to h_2
- Computing h'_2 as

$$h'_2 = h_2 \cdot e^{s(h_1)} + t(h_1)$$



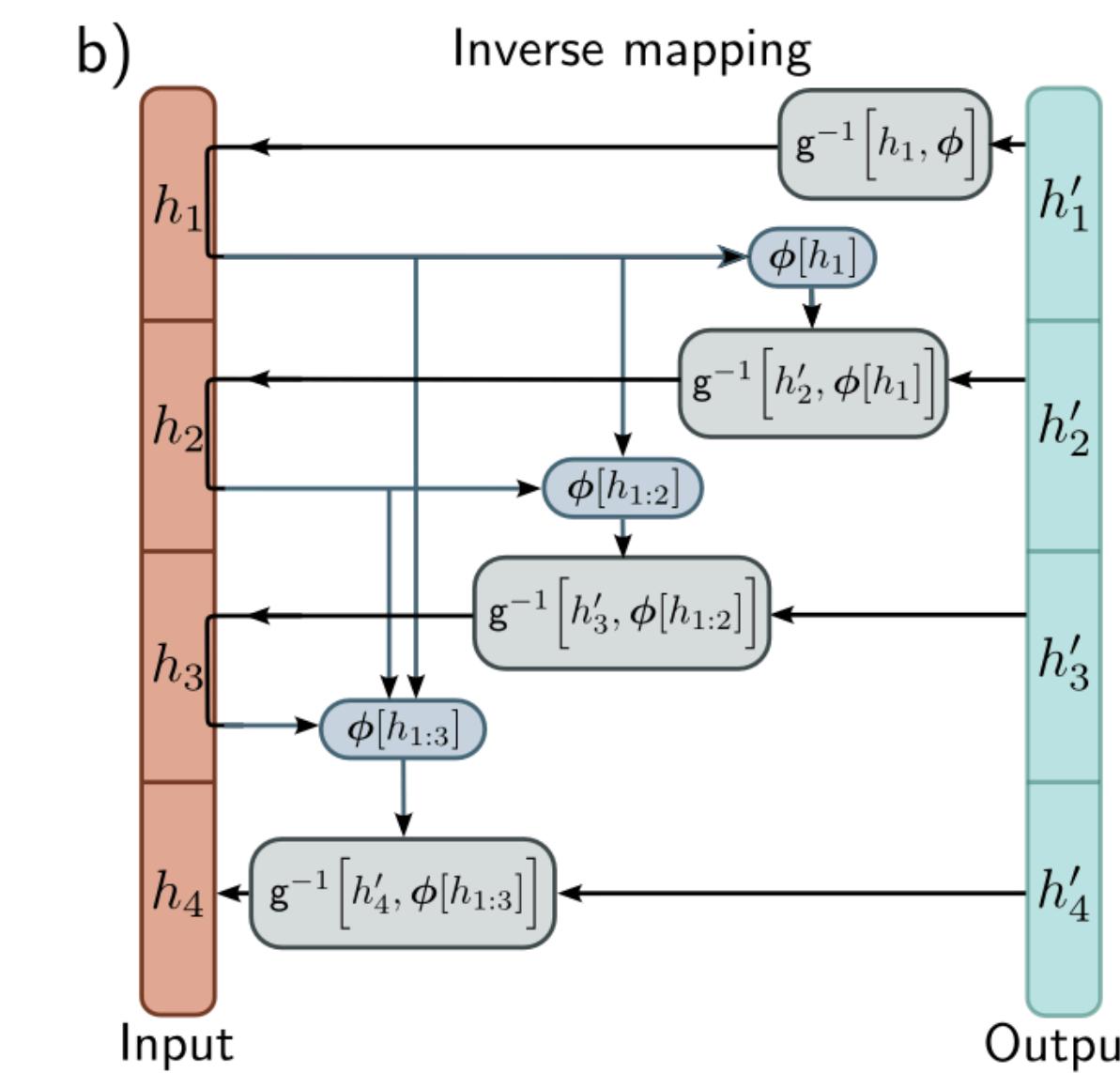
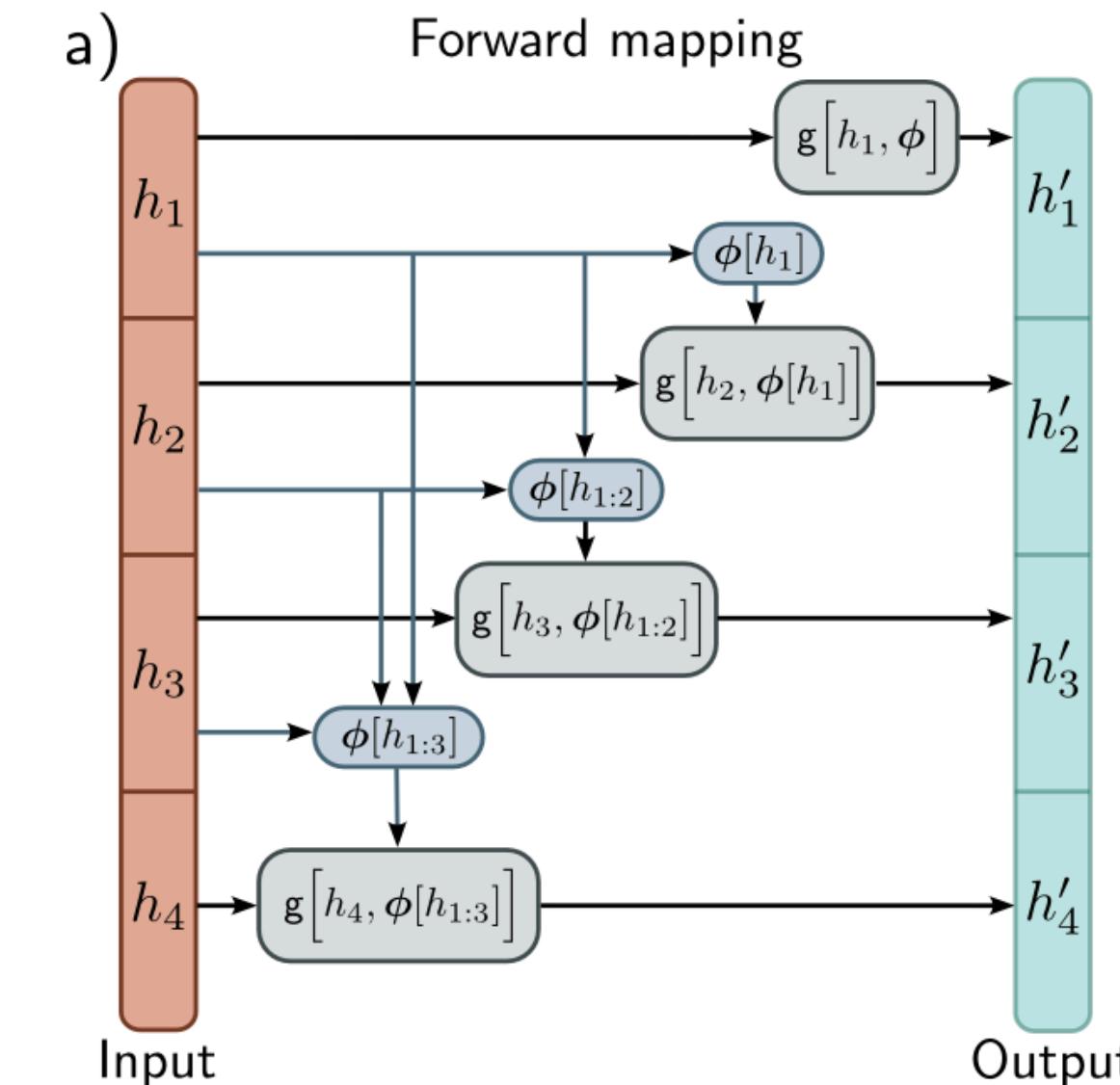
$$\begin{aligned}
 h'_1 &= h_1 \\
 h'_2 &= g[h_2, \phi[h_1]] \\
 \\
 h_1 &= h'_1 \\
 h_2 &= g^{-1}[h'_2, \phi[h_1]]
 \end{aligned}$$

Autoregressive flows

- Generalize coupling flows to enhance complexity and expressivity
- Each input is a separate block
- Each block is transformed according to a transformer function $g[h_i, \phi]$ (not related to transformer networks!!!), depending on conditioners $\phi = [\phi_{1:i}]$

$$h'_d = g[h_d, \phi[\mathbf{h}_{1:d-1}]]$$

- Notice: since the computation of h_d depends on $\mathbf{h}_{1:d-1} = [h_1, h_2, \dots, h_{d-1}]$, it cannot proceed in parallel (which slows down the inference)
- For similar reasons, the inversion has to be done sequentially



$$\begin{aligned} h'_1 &= g[h_1, \phi] \\ h'_2 &= g[h_2, \phi[h_1]] \\ h'_3 &= g[h_3, \phi[h_{1:2}]] \\ h'_4 &= g[h_4, \phi[h_{1:3}]] \end{aligned}$$

$$\begin{aligned} h_1 &= g^{-1}[h'_1, \phi] \\ h_2 &= g^{-1}[h'_2, \phi[h_1]] \\ h_3 &= g^{-1}[h'_3, \phi[h_{1:2}]] \\ h_4 &= g^{-1}[h'_4, \phi[h_{1:3}]] \end{aligned}$$

Residual flows

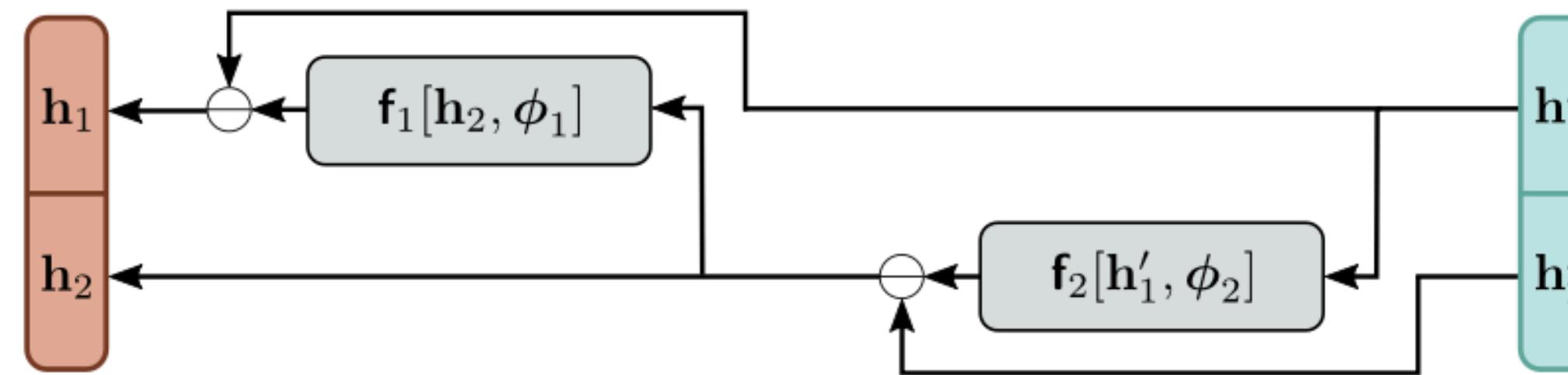
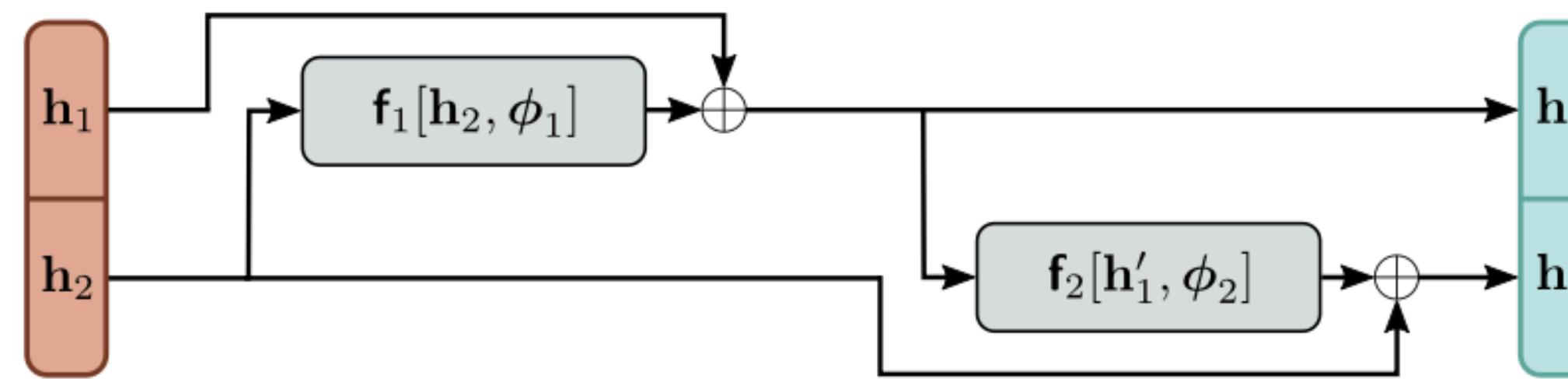
- Architecture inspired from ResNet CNNs

- as for couplings flows, split h in $[h_1, h_2]$

- Both are transformed (sequentially)

- h_2 determines $h'_1 = g[h_2, \phi_1] + h_1$ (adding h_1 through a skip connection is what comes from the ResNet inspiration)

- h'_1 goes to outputs and it also determines $h'_2 = g[h_1, \phi_2] + h_2$



$$h'_1 = h_1 + f_1[h_2, \phi_1]$$

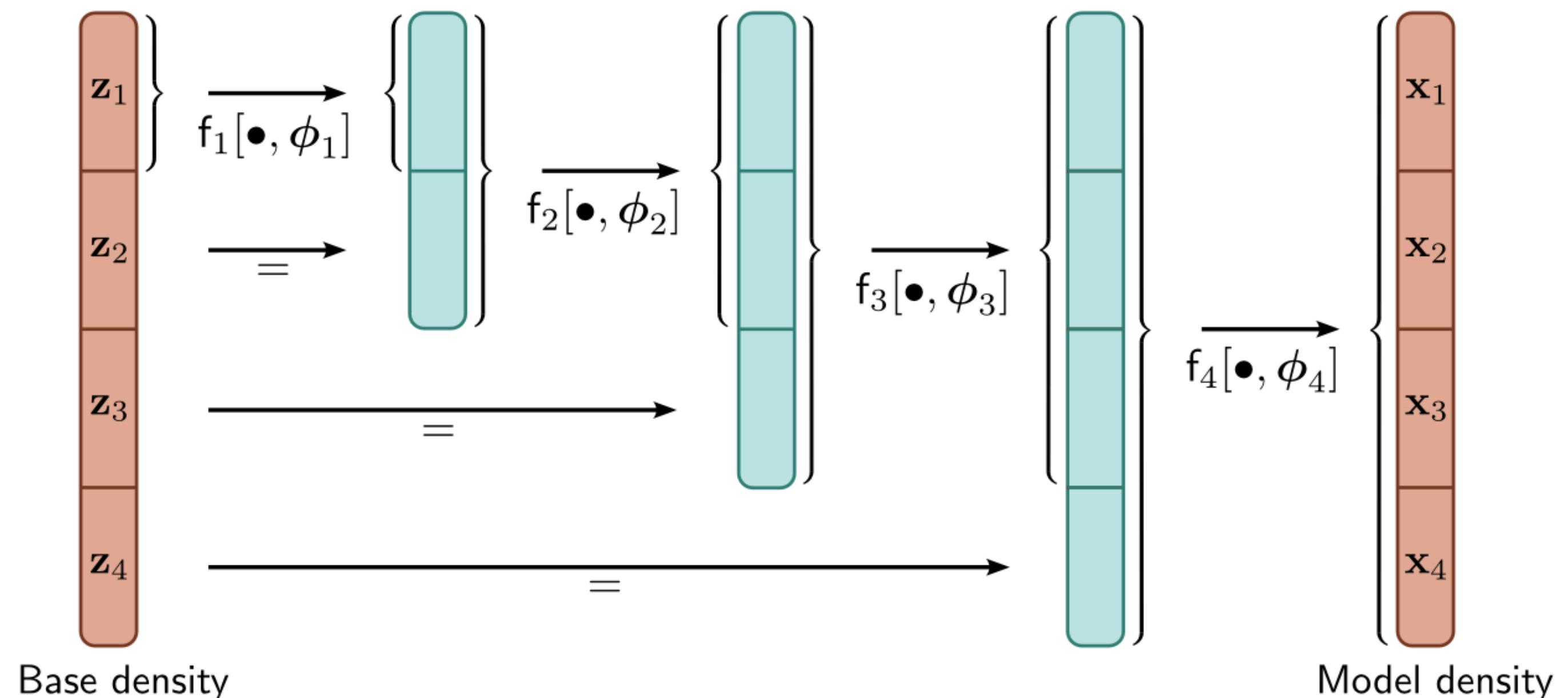
$$h'_2 = h_2 + f_2[h'_1, \phi_2]$$

$$h_2 = h'_2 - f_2[h'_1, \phi_2]$$

$$h_1 = h'_1 - f_1[h_2, \phi_1]$$

multi-scale flows

- To be invertible, flows need to have outputs of the same dimension as inputs
- On the other hand, some part of the computation might only depend on a subset (in which case, carrying around the entire input when only a fraction is needed would be inefficient)
- In multi-scale flows, this is achieved factorizing the computation and introducing additional elements progressively



Applications

- *Modeling Densities*: unlike other generative models, one can learn the likelihood of the sample

- Why GANs and VAEs can't?

- *Synthesis*: can generate images

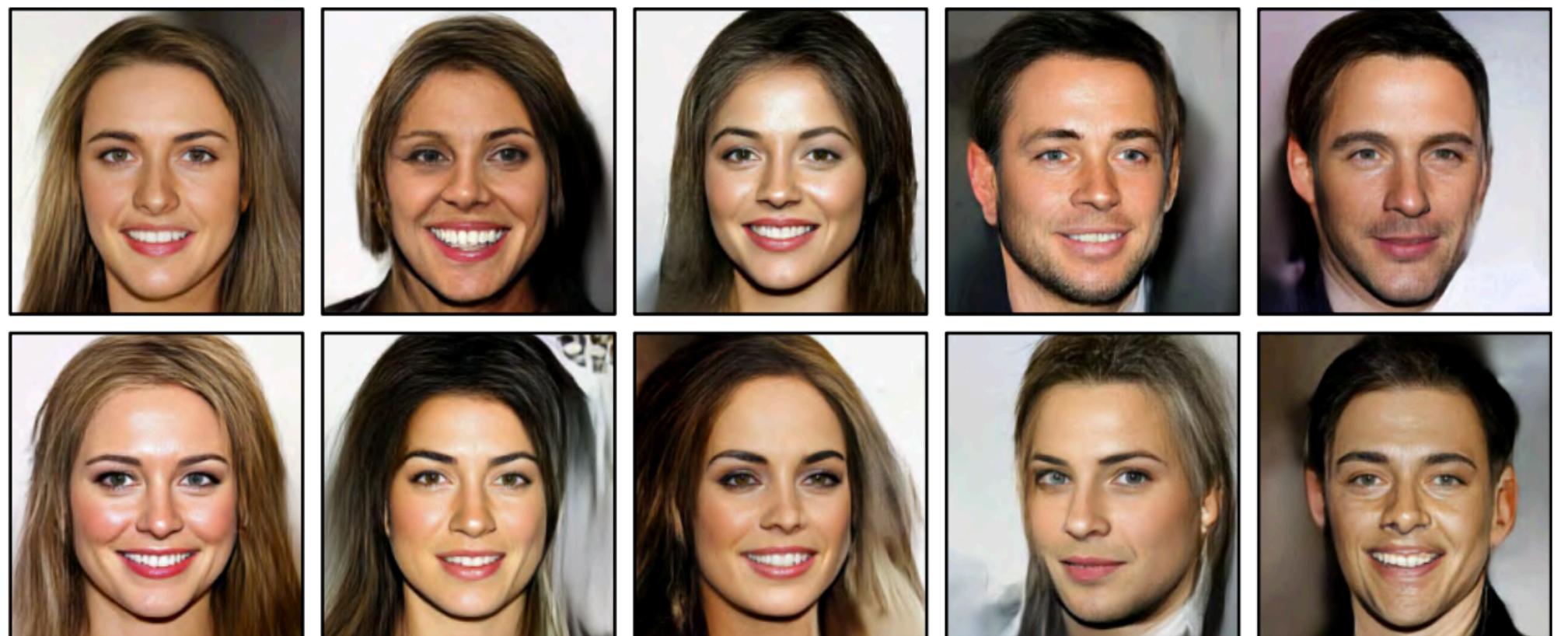
- *GLOW is an example of NormalizingFlow generating faces.*

- It starts from images (256x256x3)

- It uses coupling flows, in which the transformation imposed on h_2 depends on the position and is learned as a 2D CNN applied on h_1

- *Approximating other density models*: they can learn to approximate other functions that might be more complicated to sample from. This will be more clear when we will discuss knowledge distillation

GLOW can generate faces



One can interpolate between faces, following a linear trajectory in the latent space and applying GLOW to each point

