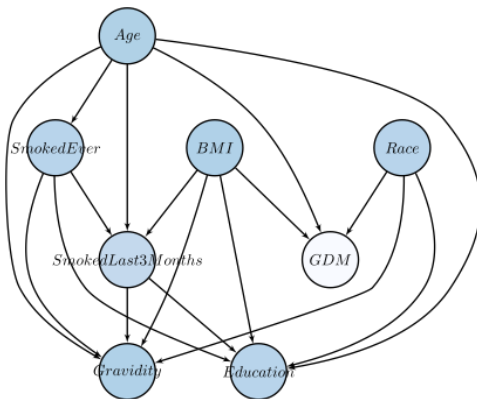# Normalizing flows

**Materials adapted from:**

- Normalizing Flows and Invertible Neural Networks, ECCV 2020 Tutorial

- Nordic Probabilistic AI School (ProbAI) 2022.
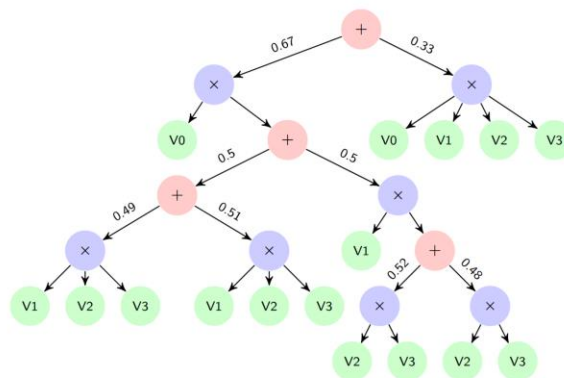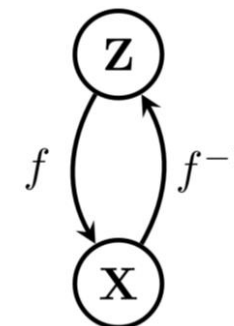
# Joint distributions



**PGMs**

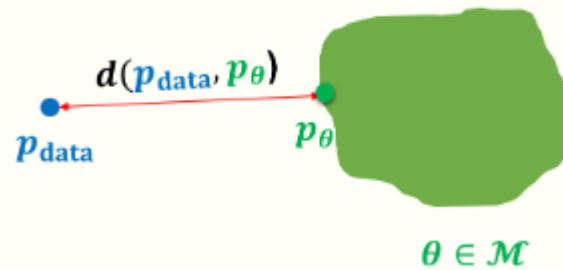Explainable

**PCs**

Tractable

**Normalizing flows**

Expressive

# Learning Probability Distributions

$x^{(j)} \sim p_{\text{data}}$
$j = 1, 2, \dots, |\mathcal{D}|$

$d(p_{\text{data}}, p_\theta)$

$p_{\text{data}}$

$p_\theta$

$\theta \in \mathcal{M}$

- Typical image resolution – 700 x 1400
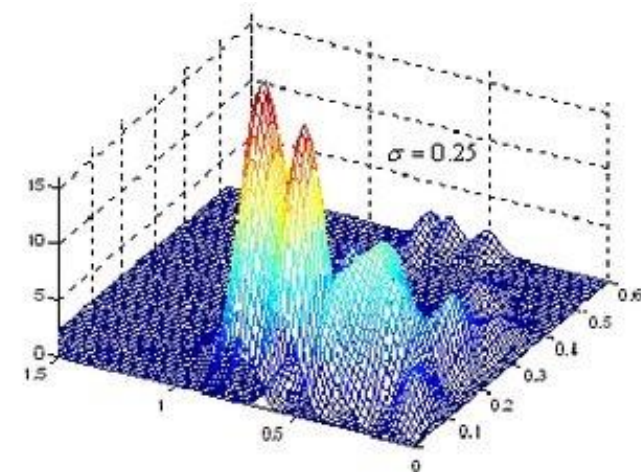- Each pixel has 3 channels  - RGB
- Each channel can take values 0-255

**How many possible images?**

$256^{700 \times 1400 \times 3} \approx 10^{800000}$

What if the data was continuous ?
- Infinite no. of values
- Complex shapes for probability densities

**We have only a limited amount of data to fit the model**



$\sigma = 0.25$

3

# Give up ?

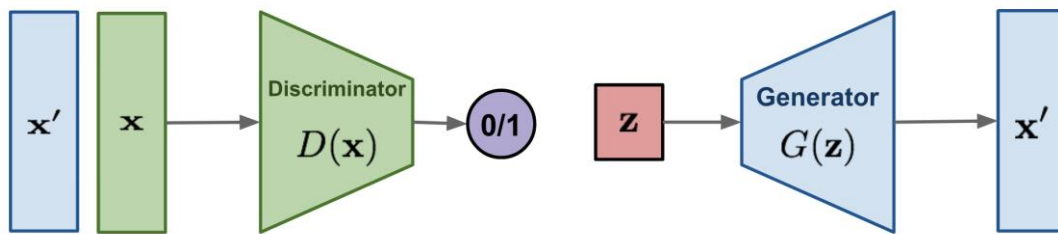Not yet. Look at these images .. Can you say which of them are real ?



All of them are generated by a model. **None of these people exist in real life !**
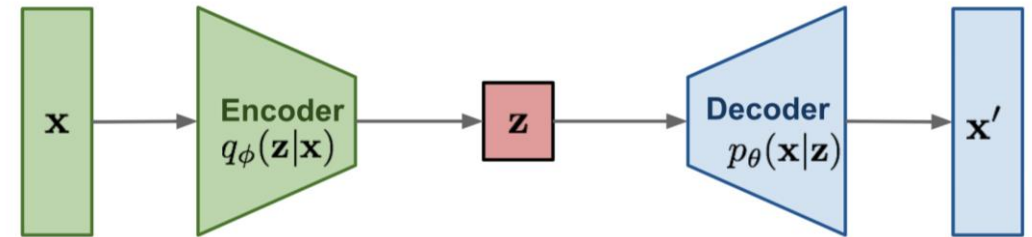
# Deep Generative Models

Using deep neural networks to learn probability distributions

- Deep latent variable models
  - Assume a simple distribution over low dimensional latent factors of variation
  - Map to data space using neural networks



**Generative Adversarial Networks (GANs)**
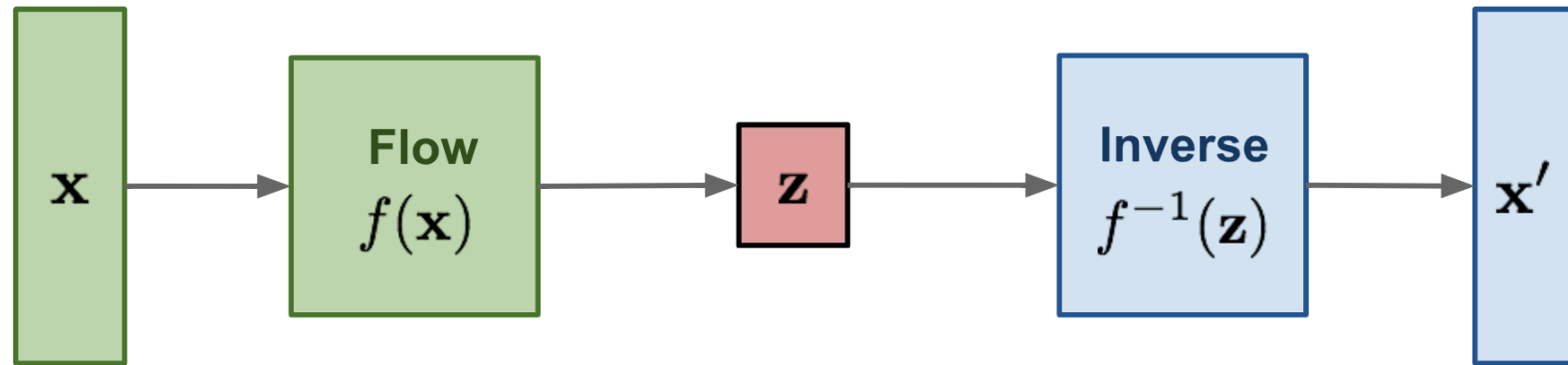         **Variational Auto Encoders (VAEs)**

Expressive models, can generate high quality samples

But **cannot** perform exact inference over the modeled probability distribution

# Normalizing flows

A new class of deep generative model

  • Utilizes invertible transformations



Why are they exciting?

  • Expressive models utilizing flexible neural networks
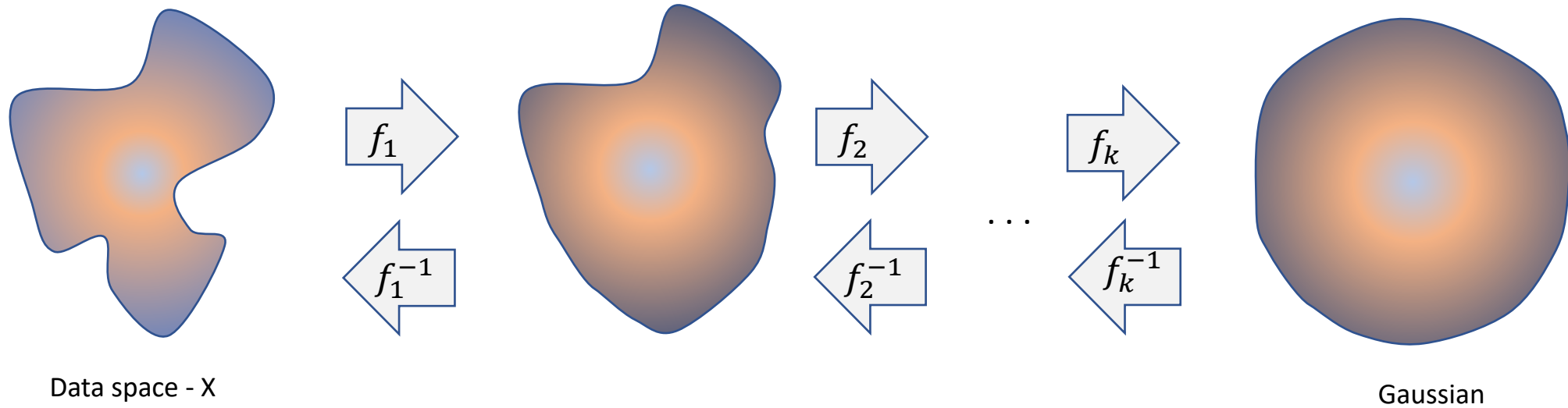  • Can perform exact likelihood evaluation

# Some cool results using normalizing flows



Almost as good as GANs / VAEs but better tractability

# Normalizing Flows:  Overview



Data space - X

Gaussian

Model complex probability distributions by
- flowing a simple distribution through a sequence of invertible transformations

Components
- Simple base distribution – typically a gaussian
- Invertible Transformation –  parameterized using neural networks

# Normalizing Flows: Formulation

**Change of variables:** Linking the probability densities in the two spaces

$$p_X(x) = \boxed{p_Z(f(x))}\boxed{|\det(J_f(x))|} \quad \text{, where} \quad J_f(x) = \frac{\partial f(x)}{\partial x}$$

Density in the gaussian space
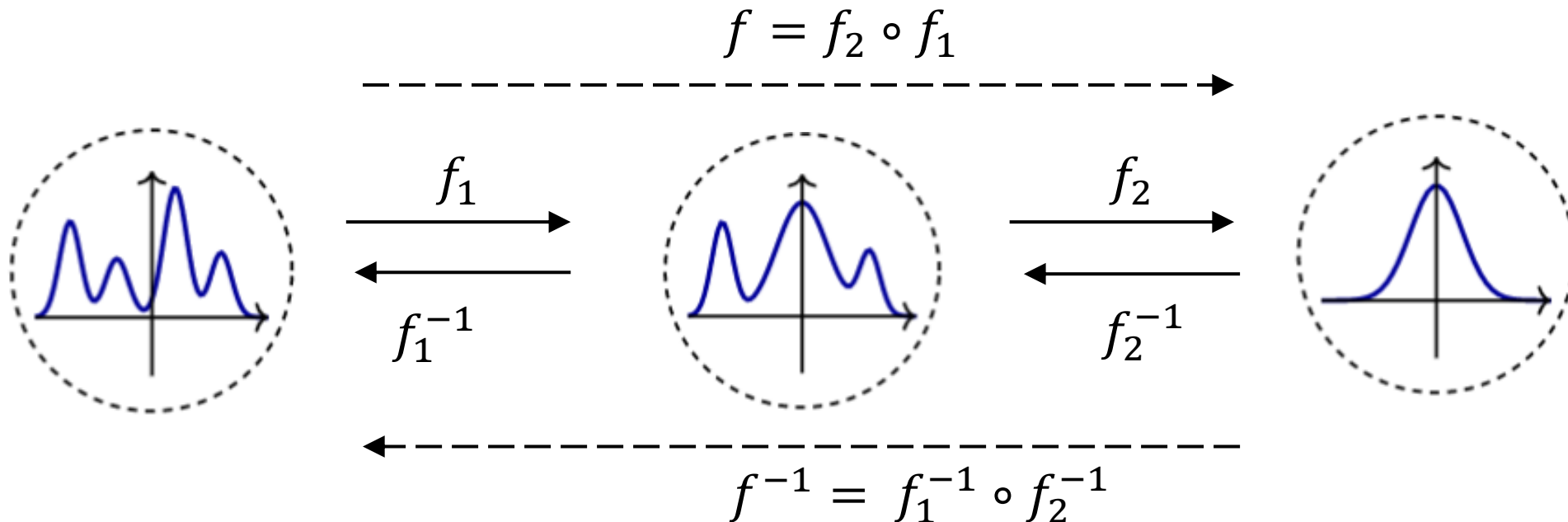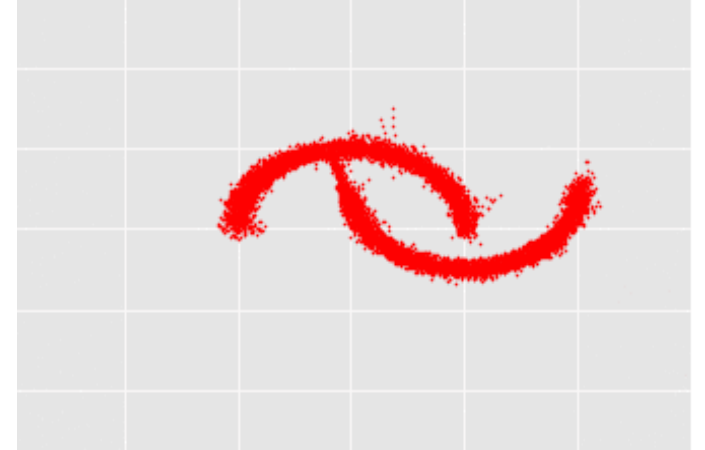
Volume correction term

**Properties**

A flow $f$ is a parametric function which is

- Invertible

- Differentiable

- Has efficiently computable Jacobian determinant

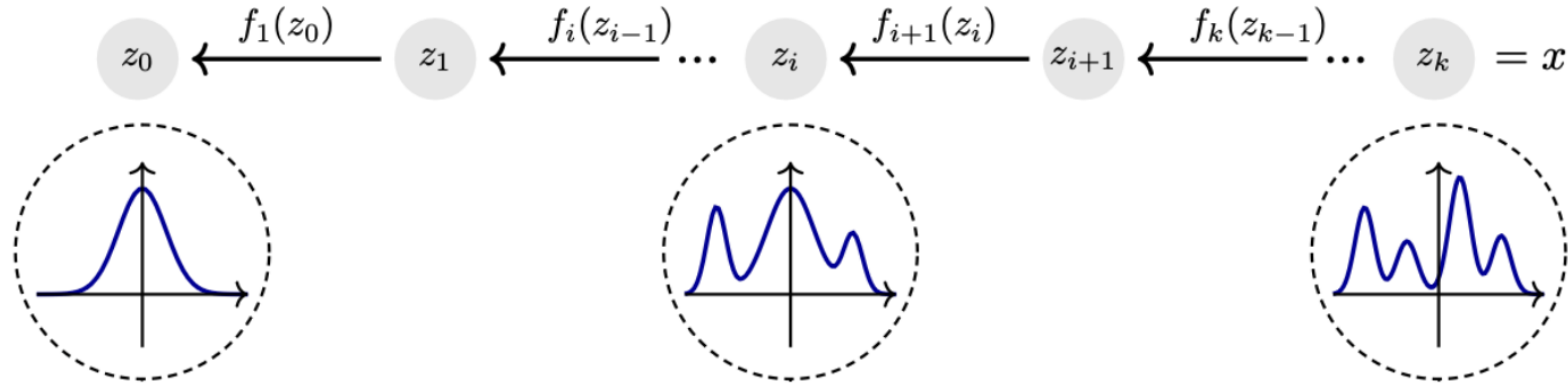Designing such functions is the core research problem in the field of flows

# Why are they called normalizing flows ?

- Composing transformations
  - Invertible differentiable functions are closed under composition

- Can slowly '*flow*' a complex distribution step by step into a normal

$$f = f_2 \circ f_1$$

$$f_1 \qquad f_2$$

$$f_1^{-1} \qquad f_2^{-1}$$

$$f^{-1} = f_1^{-1} \circ f_2^{-1}$$

# How do you train them?



$$z_0 \xleftarrow{\quad f_1(z_0) \quad} z_1 \xleftarrow{\quad f_i(z_{i-1}) \quad} \dots z_i \xleftarrow{\quad f_{i+1}(z_i) \quad} z_{i+1} \xleftarrow{\quad f_k(z_{k-1}) \quad} \dots z_k = x$$

We can compute likelihood exactly

$$p_X(x) = p_Z(f(x)) \left| \det (J_f(x)) \right| \quad = p_Z(f_1 \circ f_2 \circ \dots f_k(x)) \left| \det (J_{f_1 \circ f_2 \circ \dots f_k}(x)) \right|$$

$$= p_Z(f_1 \circ f_2 \circ \dots f_k(x)) \prod_i \left| \det (J_{f_i}) \right|$$

Maximize the log-likelihood of the data points w.r.t to the parameters of the flow

$$\arg\max_\theta \sum_m \log p_X(x_m) = \arg\max_\theta \sum_m \log p_Z(f_\theta(x_m)) + \sum_m \sum_i \log \left| \det (J_{f_{\theta_i}}) \right|$$

$f$ is differentiable, solve using **gradient descent**

# Tractable for what ?



$$z_0 \xrightarrow{f_1^{-1}(z_0)} z_1 \cdots \xrightarrow{f_i^{-1}(z_{i-1})} z_i \xrightarrow{f_{i+1}^{-1}(z_i)} z_{i+1} \cdots \xrightarrow{f_k^{-1}(z_{k-1})} z_k = x$$

$$z_0 \sim p_0(z_0) \qquad z_i \sim p_i(z_i) \qquad z_k \sim p_k(z_k)$$

- **Sampling**
  - Sample from gaussian and apply inverse transform
- **Density estimation**
  - Use change of variables formula to compute $p_X(x)$ exactly.

| ✓ Evidential | ✗ Marginal | ✗ Conditional | ✗ MAP |
|---|---|---|---|

# Parameterizing flow transformations

## Linear Flows
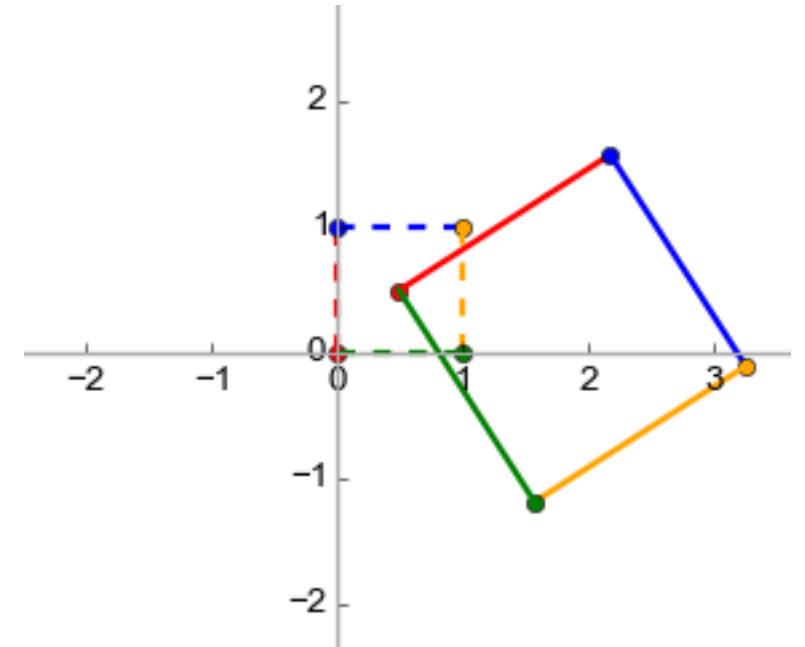
We can define a flow using a linear transformation if the matrix is invertible

$$f(\boldsymbol{x}) = A\boldsymbol{x} + b$$

- Inverse: $f^{-1}(\boldsymbol{z}) = A^{-1}(\boldsymbol{z} - b)$

- Jacobian is same as $A$. Thus, $\det J_f = \det A$

### Caveats

- Will have to ensure invertibility during learning

- Affine gaussians are gaussians! Not quite expressive.

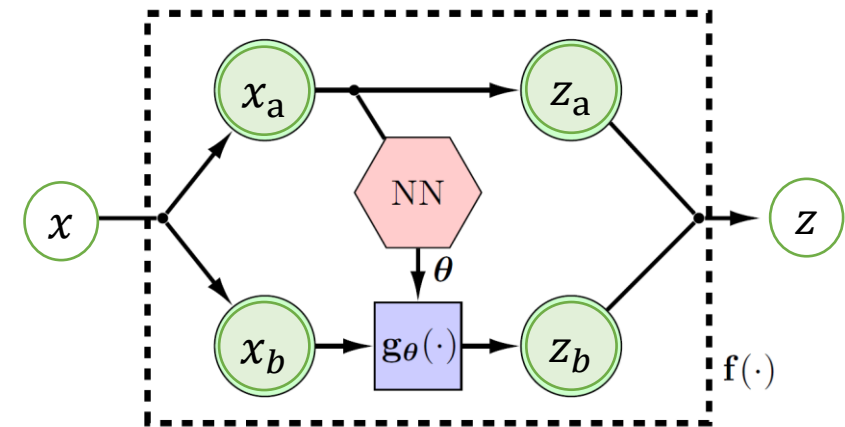- Determinant computation is expensive $O(d^3)$

# Coupling Flows

Partition the input into two disjoint sets $x = (x_A, x_B)$

Define $f(x) = (x_A, g(x_B | \theta(x_A))$

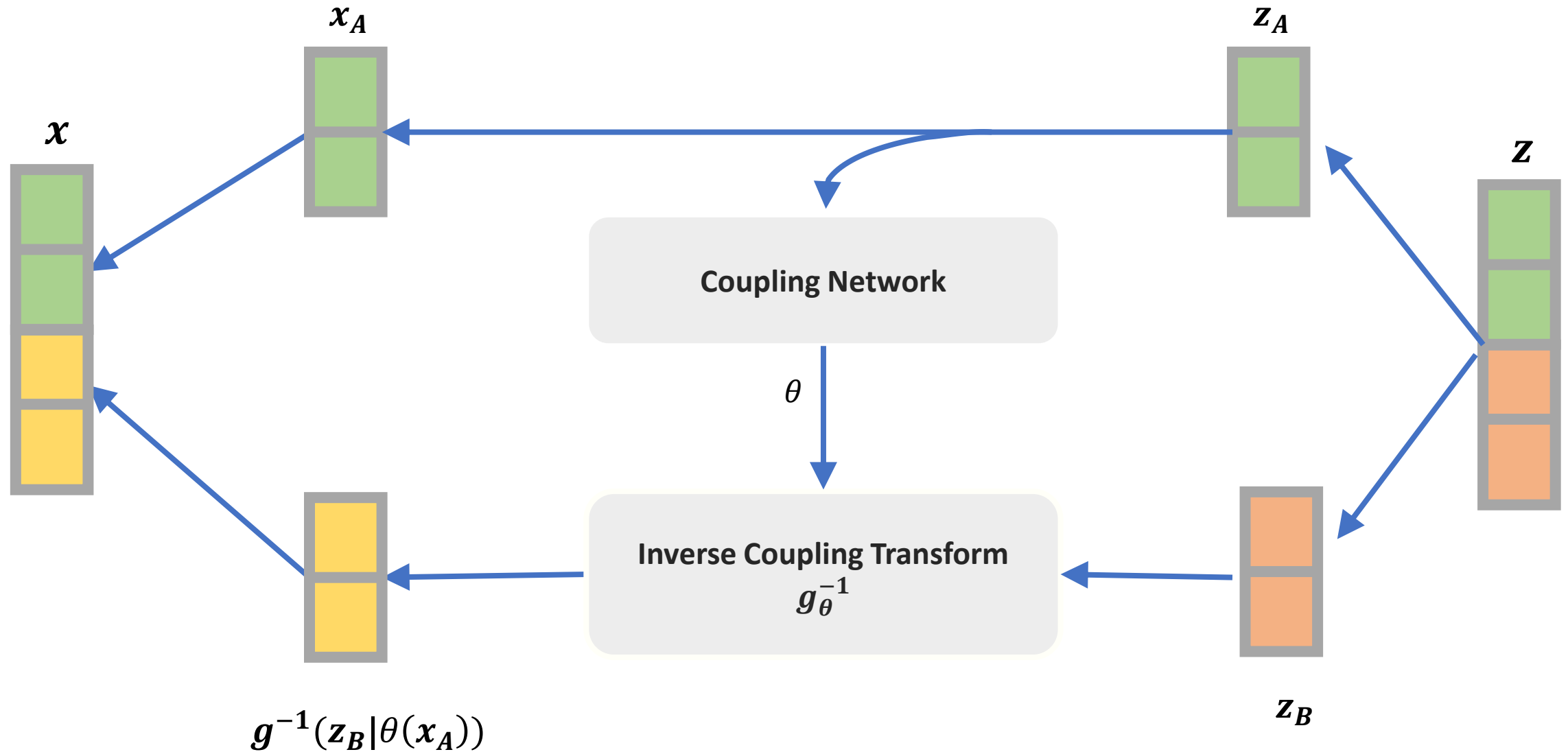$g$ is an invertible transformation whose parameters depend on $x_A$



- Will use this as a case study demonstrate the key properties
  - invertibility
  - efficiently Jacobian determinant
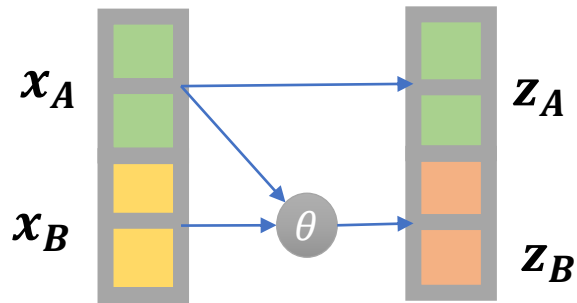- Also implement in hands on part

14

# Coupling Flow: The forward transformation



$x$

$x_A$

$x_A$

$z$

**Coupling Network**

$\theta$

$x_B$

**Coupling Transform** $g_{\theta}$

$g(x_B | \theta(x_A))$

# Coupling Flow: The inverse transformation



$x_A$

$z_A$

$x$

$z$

**Coupling Network**

$\theta$

**Inverse Coupling Transform**
$g_\theta^{-1}$

$z_B$

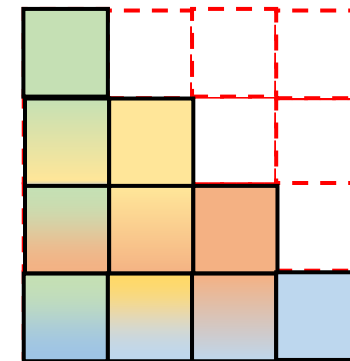$g^{-1}(z_B | \theta(x_A))$

# Coupling Flow: Jacobian

$$J_f(\boldsymbol{x}) = \begin{bmatrix} \dfrac{\partial z_A}{\partial x_A} & \dfrac{\partial z_A}{\partial x_B} \\[2em] \dfrac{\partial z_B}{\partial x_A} & \dfrac{\partial z_B}{\partial x_B} \end{bmatrix} = \begin{bmatrix} \mathbb{I} & 0 \\[2em] \dfrac{\partial g(x_B|\theta(x_A))}{\partial x_A} & J_g \end{bmatrix}$$

$x_A$
$x_B$
$z_A$
$z_B$

We don't need this for determinant!

$\theta$ can be an arbitrarily complex (non-invertible) model like a neural network

Block diagonal jacobian, determinant computation is efficient
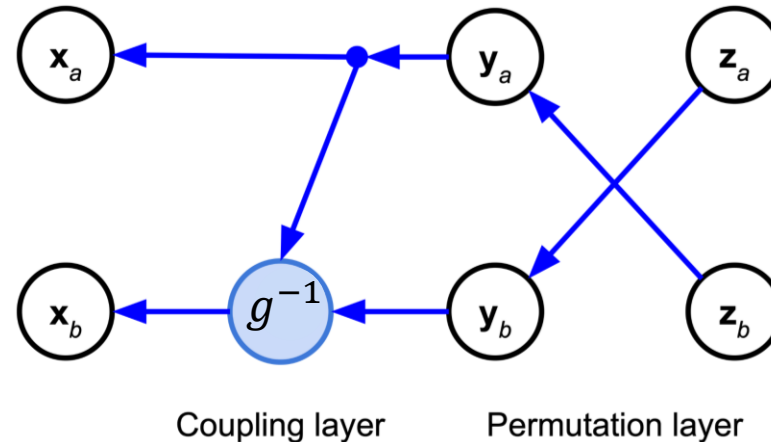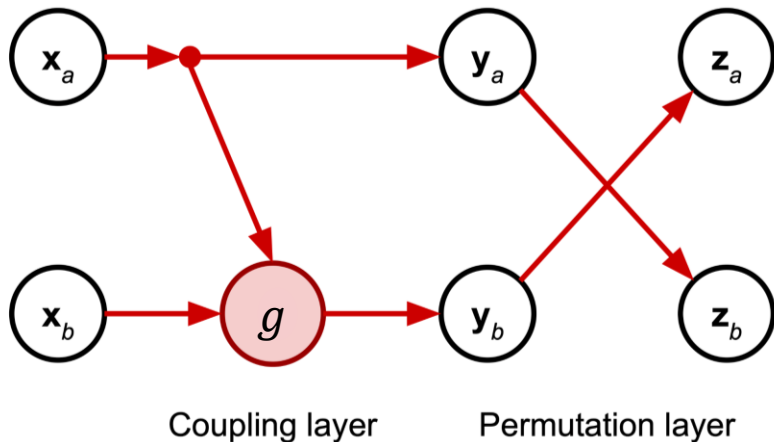
$O(d)$

# Coupling Transforms

- Permutation is an invertible linear operation
  - Has unit determinant

$$\underbrace{\begin{bmatrix} 0 & 0 & 0 & 1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{bmatrix}}_{\mathbf{P^T}\ \text{matrix}} \cdot \underbrace{\begin{bmatrix} 4 & 6 & 8 & 1 \\ 5 & 3 & 1 & 0 \\ 7 & 21 & 0 & 9 \\ 3 & -4 & 1 & 3 \end{bmatrix}}_{\mathbf{A}\ \text{matrix}} = \underbrace{\begin{bmatrix} 3 & -4 & 1 & 3 \\ 4 & 6 & 8 & 1 \\ 5 & 3 & 1 & 0 \\ 7 & 21 & 0 & 9 \end{bmatrix}}_{\mathbf{P^T \cdot A}\ \text{matrix}}$$



Coupling layer    Permutation layer            Coupling layer    Permutation layer

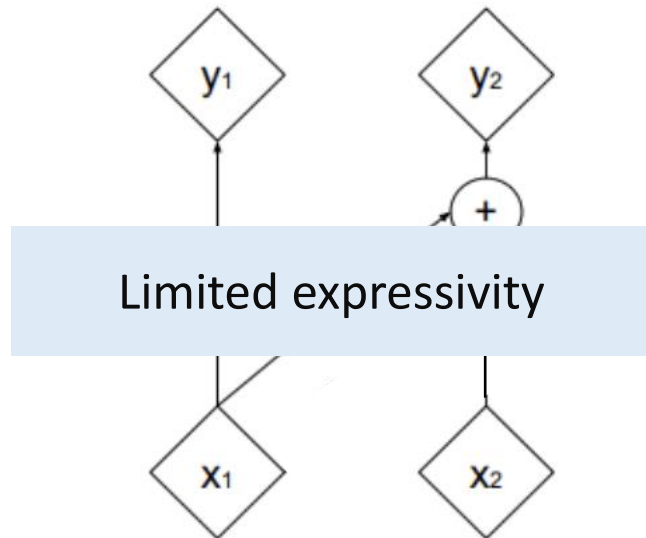# Coupling Transforms

Defining Coupling transforms $g$ ..

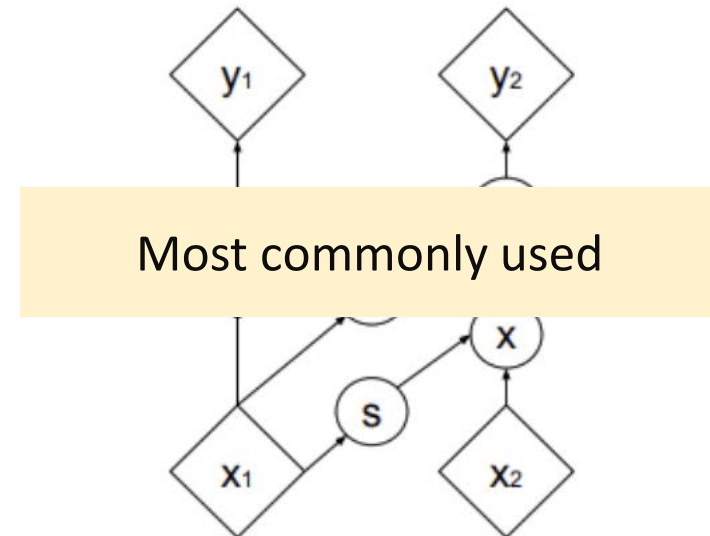**Simple** invertible transformations suffice

- Additive [NICE-Dinh et al., 2014]
  - $g(x|t) \quad = x + t$
  - $g^{-1}(z|t) = z - t$

- Affine [RealNVP-Dinh et al., 2017]
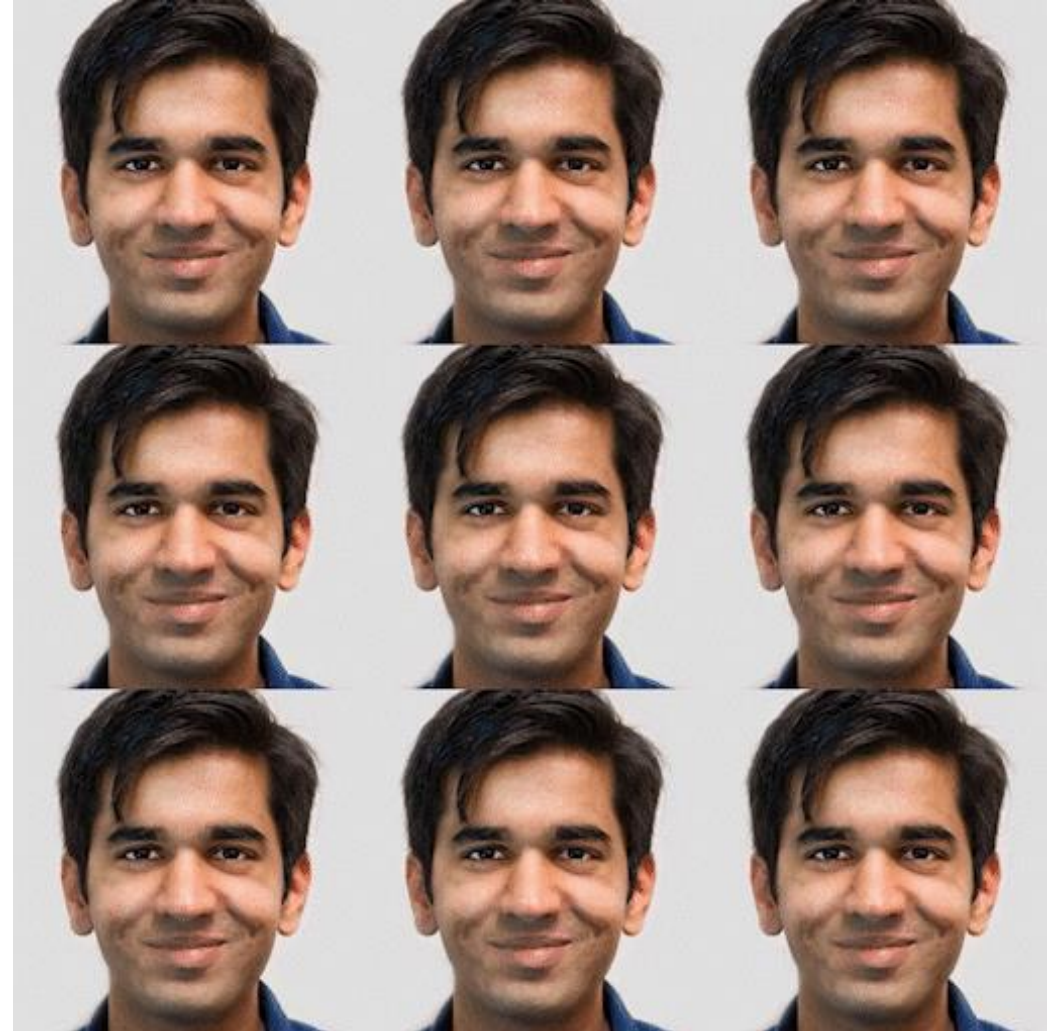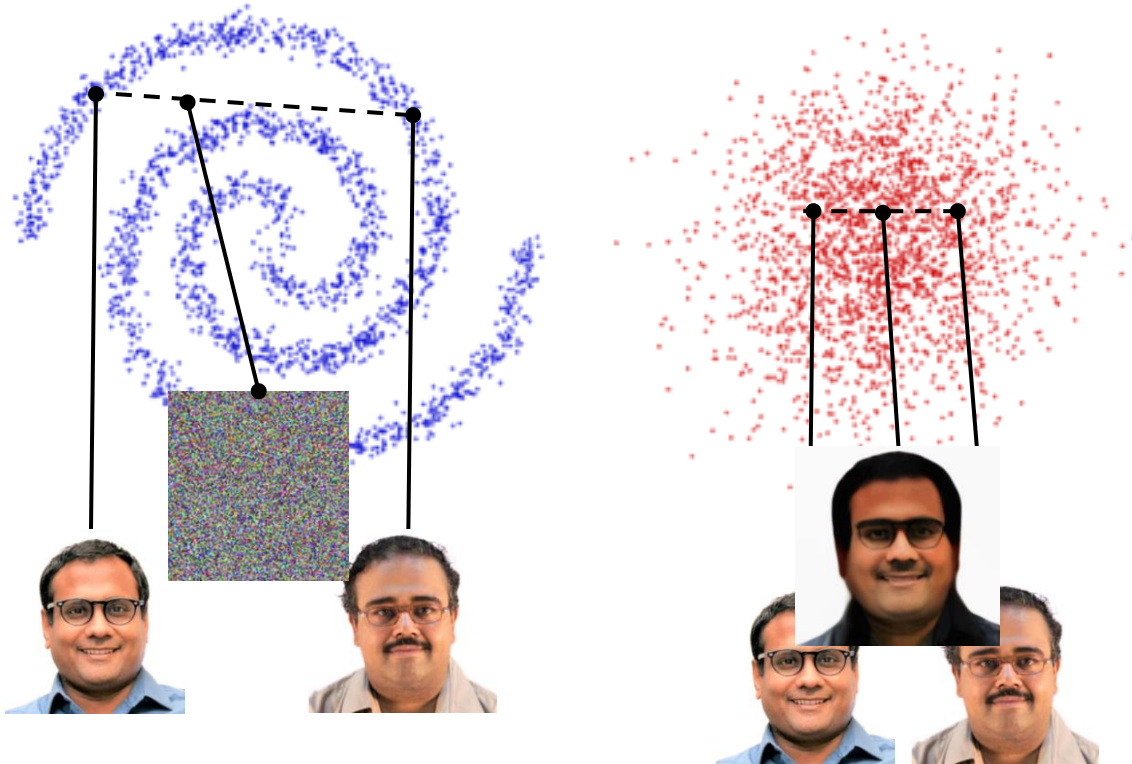  - $g(x|s,t) \quad = \exp(s) \odot x + t$
  - $g^{-1}(z|s,t) = \exp(s)^{-1} \odot (z - t)$



Limited expressivity

Most commonly used

Can have more complex coupling transform formulations .. MLP, splines, etc. [Kingma et al., 2016, Huang et al., 2018; de Cao et al., 2019; Durkan et al., 2019]
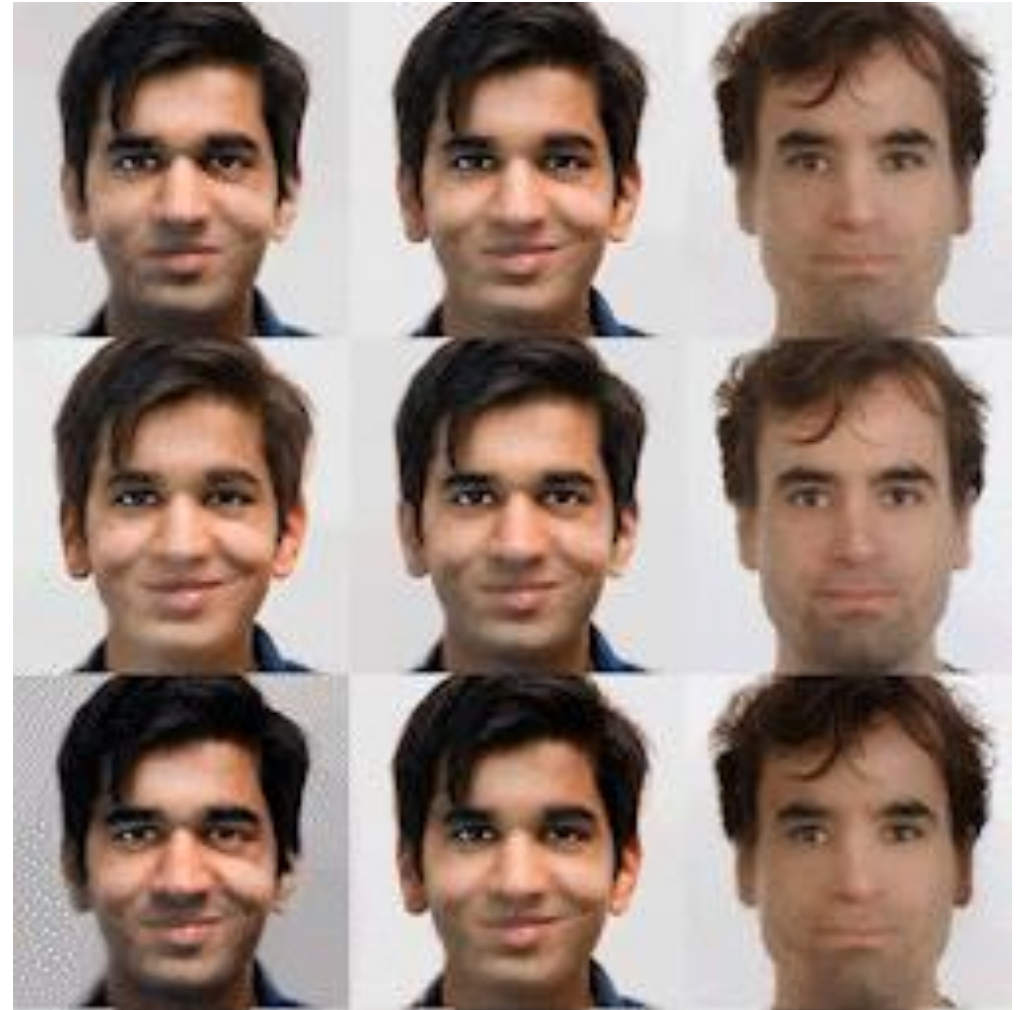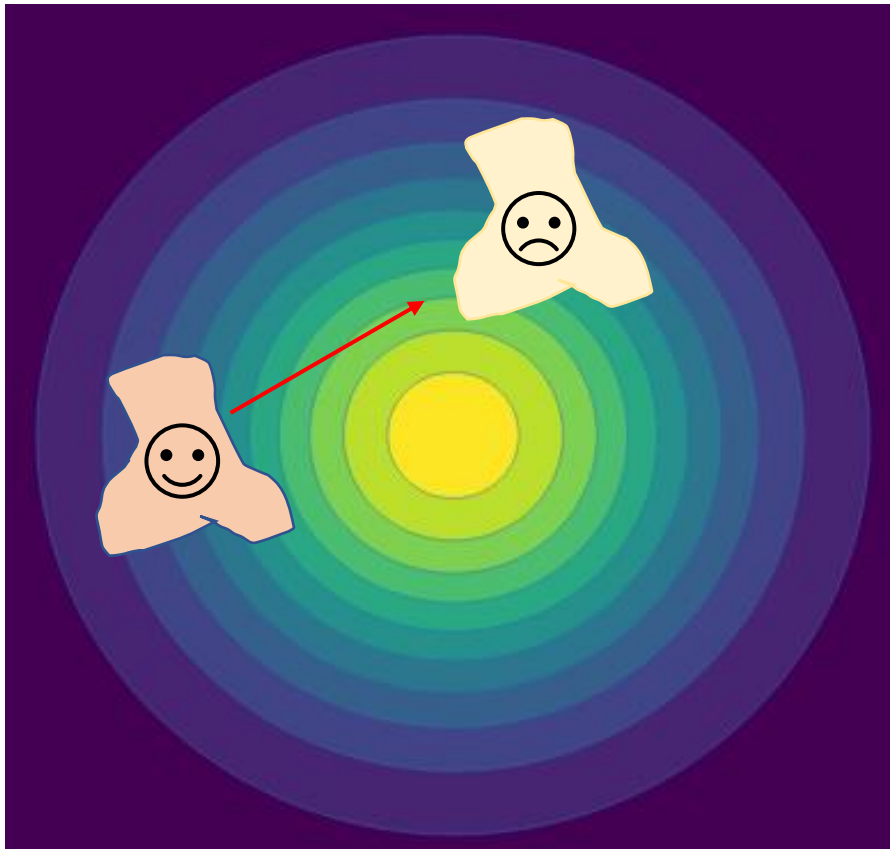
# Utilizing the latent space of the flow

Can Interpolate between points in the latent space

# Utilizing the latent space of the flow

Manipulate semantic features by moving in certain directions of the latent space

# Normalizing Flows
# Hands-On <span style="color:red">Demo</span>
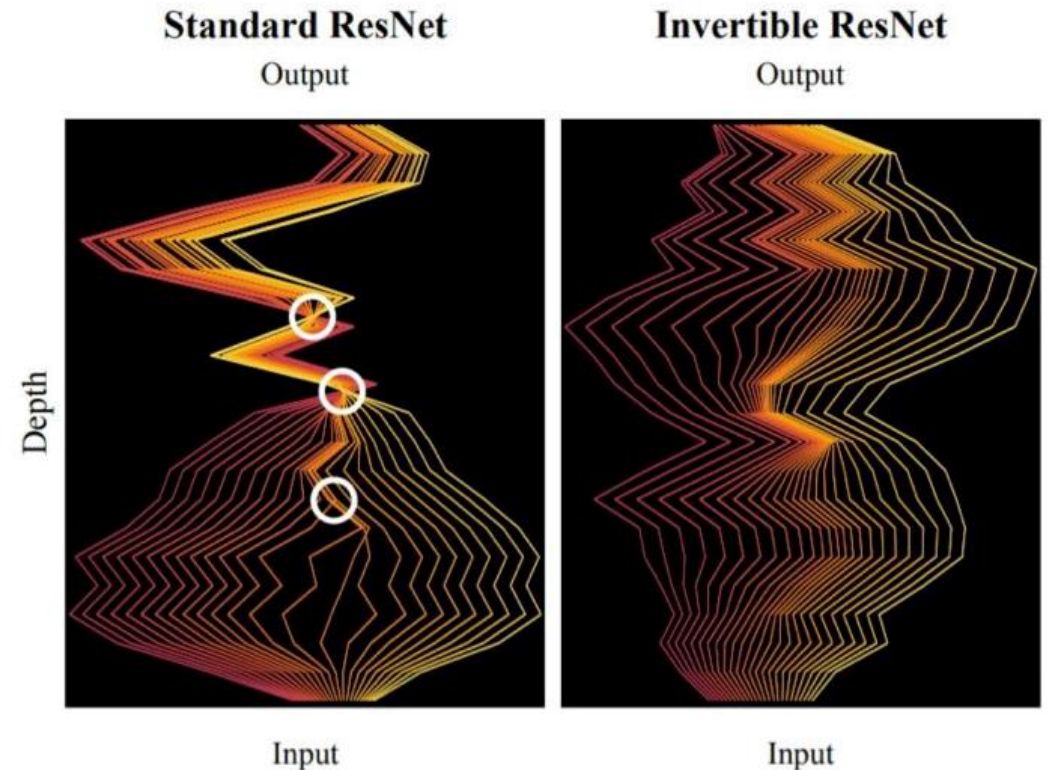
https://bit.ly/tpm-day2-flows

# Conclusion

- Other approaches
  - Discrete time invertible transforms
    - Contractive Flows [Behrmann et al. 2019, Perugachi-Diaz et. Al 2020, Chen et al.]
  - Continuous time invertible transforms
    - Neural ODE flows [Grathwohl et al. 2018]

- Limitations and takeaways
  - Tractability compared to PCs
  - Dimensionality preserving
  - Topological restrictions of diffeomorphisms



**Standard ResNet**
Output

**Invertible ResNet**
Output

Depth

Input

Input

# Some Resources on Flows

**Review Paper**

Papamakarios, G., Nalisnick, E., Rezende, D. J., Mohamed, S., & Lakshminarayanan, B. (2021). **Normalizing flows for probabilistic modeling and inference**. *The Journal of Machine Learning Research*,

**Tutorial**

Normalizing Flows and Invertible Neural Networks in Computer Vision, CVPR 2021 by Marcus A. Brubaker, Ullrich Köthe

**Github Repo**

janosh / **awesome-normalizing-flows** (Public)

# Thank You!