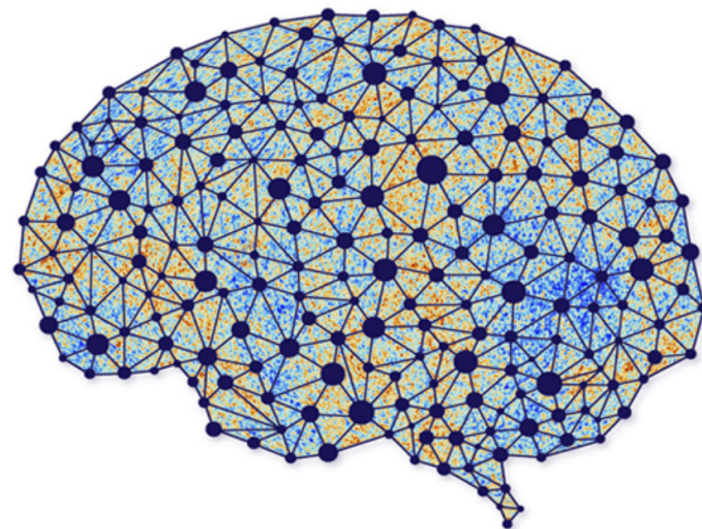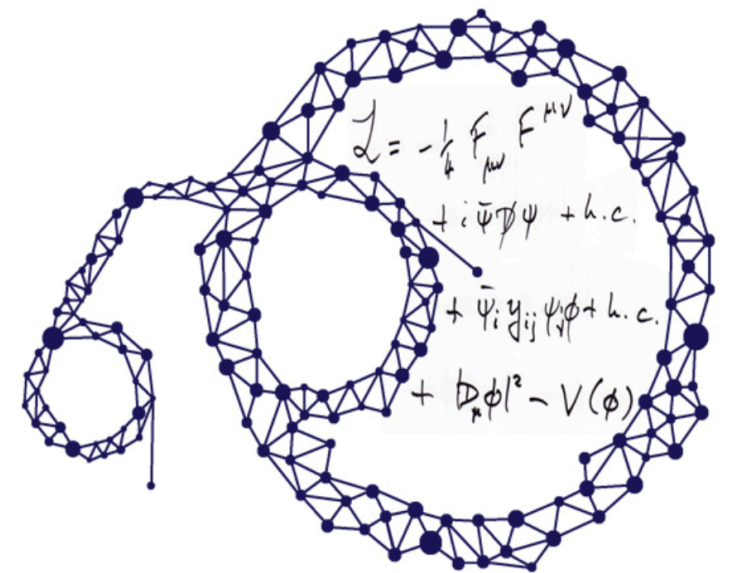# Physics 361 - Machine Learning in Physics

# Lecture 23 – More Generative Models

**April 17th 2025**
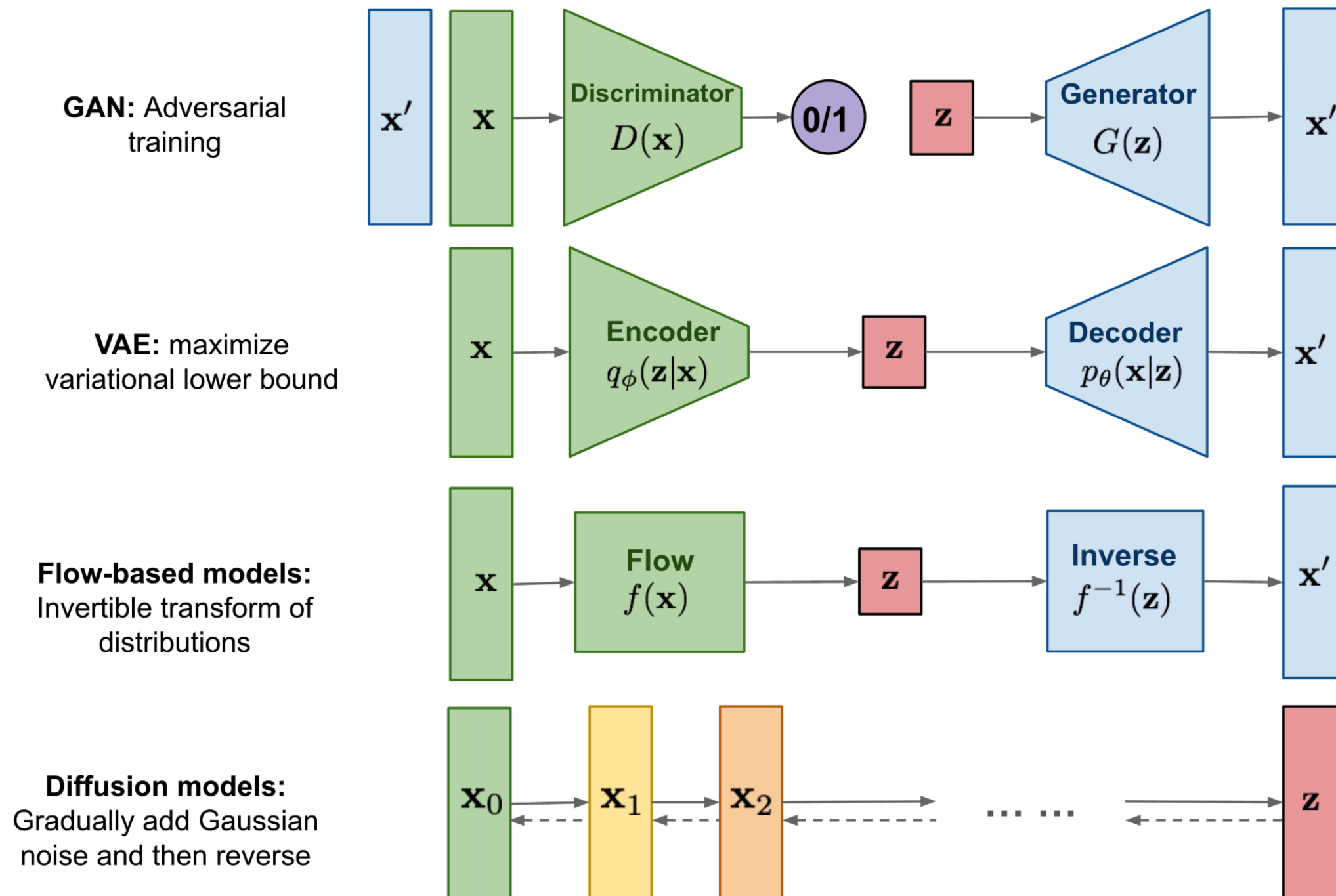


AI
∩
Universe

**Moritz Münchmeyer**

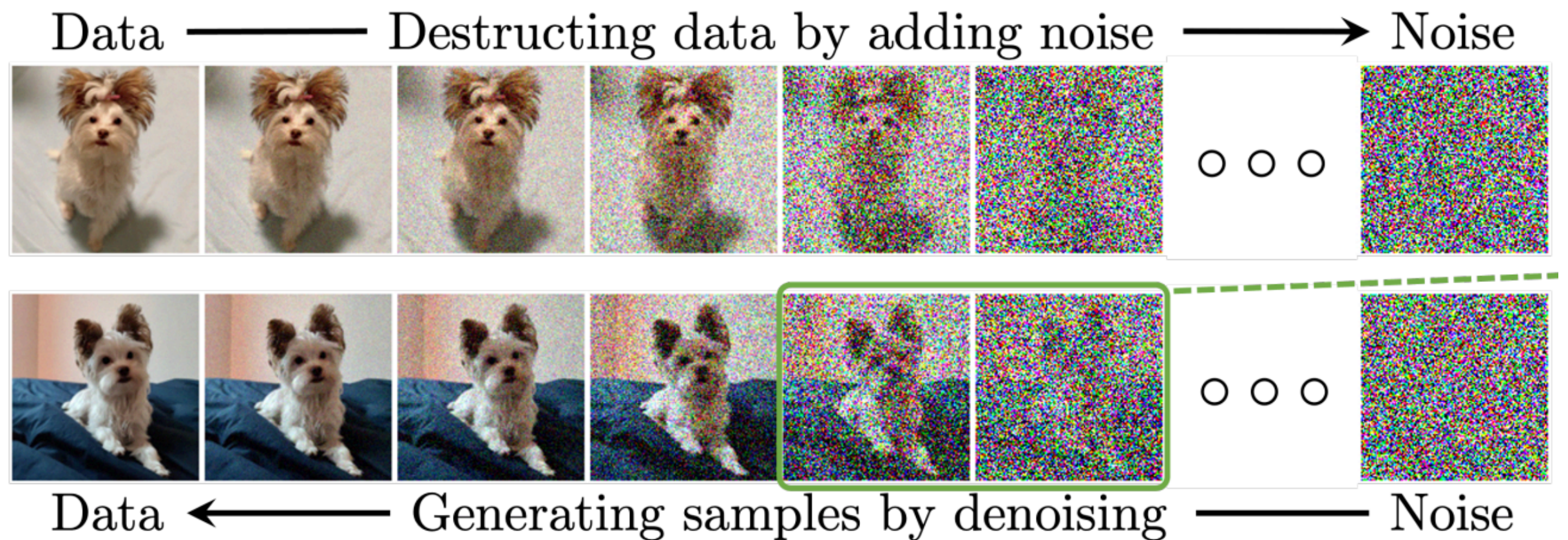# Common Generative Models



**GAN:** Adversarial training

$\mathbf{x}'$   $\mathbf{x}$ → Discriminator $D(\mathbf{x})$ → 0/1   $\mathbf{z}$ → Generator $G(\mathbf{z})$ → $\mathbf{x}'$

**VAE:** maximize variational lower bound

$\mathbf{x}$ → Encoder $q_\phi(\mathbf{z}|\mathbf{x})$ → $\mathbf{z}$ → Decoder $p_\theta(\mathbf{x}|\mathbf{z})$ → $\mathbf{x}'$

**Flow-based models:** Invertible transform of distributions

$\mathbf{x}$ → Flow $f(\mathbf{x})$ → $\mathbf{z}$ → Inverse $f^{-1}(\mathbf{z})$ → $\mathbf{x}'$

**Diffusion models:** Gradually add Gaussian noise and then reverse

$\mathbf{x}_0$ ⇄ $\mathbf{x}_1$ ⇄ $\mathbf{x}_2$ ⇄ · · · · · ⇄ $\mathbf{z}$

A fifth class are **autoregressive models based on transformers,** which we discussed before.

Figure credit: https://lilianweng.github.io/posts/2021-07-11-diffusion-models/

# Recall: DDPM

- We will focus on the most popular version of diffusion models, "denoising diffusion probabilistic models" (DDPM).

- Pictorially the process works as follows



Data ——————— Destructing data by adding noise ——————→ Noise

Data ←——————— Generating samples by denoising ——————— Noise

- The noise to data (denoising) process is **learned by a neural network (e.g. a U-Net)**, which is applied many times (roughly 100 to 1000 times).

https://arxiv.org/pdf/2209.00796.pdf

# Short introduction to GANs

Plots and discussion from Bishop book

**One night in 2014, Ian Goodfellow went drinking to celebrate with a fellow** doctoral student who had just graduated. At Les 3 Brasseurs (The Three Brewers), a favorite Montreal watering hole, some friends asked for his help with a thorny project they were working on: a computer that could create photos by itself.

Researchers were already using neural networks, algorithms loosely modeled on the web of neurons in the human brain, as "generative" models to create plausible new data of their own. But the results were often not very good: images of a computer-generated face tended to be blurry or have errors like missing ears. The plan Goodfellow's friends were proposing was to use a complex statistical analysis of the elements that make up a photograph to help machines come up with images by themselves. This would have required a massive amount of number-crunching, and Goodfellow told them it simply wasn't going to work.

But as he pondered the problem over his beer, he hit on an idea. What if you pitted two neural networks against each other? His friends were skeptical, so once he got home, where his girlfriend was already fast asleep, he decided to give it a try. Goodfellow coded into the early hours and then tested his software. It worked the first time.

What he invented that night is now called a GAN, or "generative adversarial network." The technique has sparked huge excitement in the field of machine learning and turned its creator into an AI celebrity.

# Are GANs still useful?

- GANs are still actively used although they have lost a lot of ground to diffusion models.

- An advantage over diffusion: GANs generate samples in a single forward pass, whereas diffusion models typically require hundreds of denoising steps (though that's improving with techniques like DDIM, LCM, etc.).

- The adversarial training objective is very interesting and I think it is still worth to spend a few slides on them.

- Ideas often become important again in machine learning (and physics).

# Adversarial training

Latent space distribution $\qquad P(\vec{z}) = N(z \mid 0, I)$

Gaussian

Non-Linear transformation "Generator"

$$\vec{x} = g(\vec{z}, \vec{w})$$

weights to be learned from $\{x_n\}$

NOT a bijection $\neq$ normalizing flows

$\Rightarrow$ defines $P(\vec{x})$

How to learn $g(\vec{z}, \vec{w})$ ?

# Idea

real images



synthetic images

Discriminator

$d(\mathbf{x}, \phi)$

$t$

$\mathbf{z} \longrightarrow$ Generator

$\mathbf{g}(\mathbf{z}, \mathbf{w})$

Schematic illustration of a GAN in which a discriminator neural network $d(\mathbf{x}, \phi)$ is trained to distinguish between real samples from the training set, in this case images of kittens, and synthetic samples produced by the generator network $\mathbf{g}(\mathbf{z}, \mathbf{w})$. The generator aims to maximize the error of the discriminator network by producing realistic images, whereas the discriminator network tries to minimize the same error by becoming better at distinguishing real from synthetic examples.

# Loss function

Define binary target variable

$$t = 1, \quad \text{real data,}$$
$$t = 0, \quad \text{synthetic data}$$

Descriminator network outputs 1d probability (via sigmoid)

$$P(t = 1) = d(\mathbf{x}, \boldsymbol{\phi})$$

Loss is the standard cross-entropy with 2 classes:

$$E(\mathbf{w}, \boldsymbol{\phi}) = -\frac{1}{N} \sum_{n=1}^{N} \{t_n \ln d_n + (1 - t_n) \ln(1 - d_n)\}$$

output of descriminator network

# Loss function

Training set contains real data and synthetic data (from the generator).

$\Longrightarrow$

$$E_{\text{GAN}}(\mathbf{w}, \boldsymbol{\phi}) = -\frac{1}{N_{\text{real}}} \sum_{n \in \text{real}} \ln d(\mathbf{x}_n, \boldsymbol{\phi})$$

$$-\frac{1}{N_{\text{synth}}} \sum_{n \in \text{synth}} \ln(1 - d(\mathbf{g}(\mathbf{z}_n, \mathbf{w}), \boldsymbol{\phi}))$$

↑ random sample from $p(z) = \mathcal{N}(z)$

We train generator $g$ and discriminator $d$ together, end to end.

$$\Delta \boldsymbol{\phi} = -\lambda \nabla_{\boldsymbol{\phi}} E_n(\mathbf{w}, \boldsymbol{\phi})$$ minimize wrt. $\phi$

$$\Delta \mathbf{w} = \lambda \nabla_{\mathbf{w}} E_n(\mathbf{w}, \boldsymbol{\phi})$$ maximize wrt. $w$

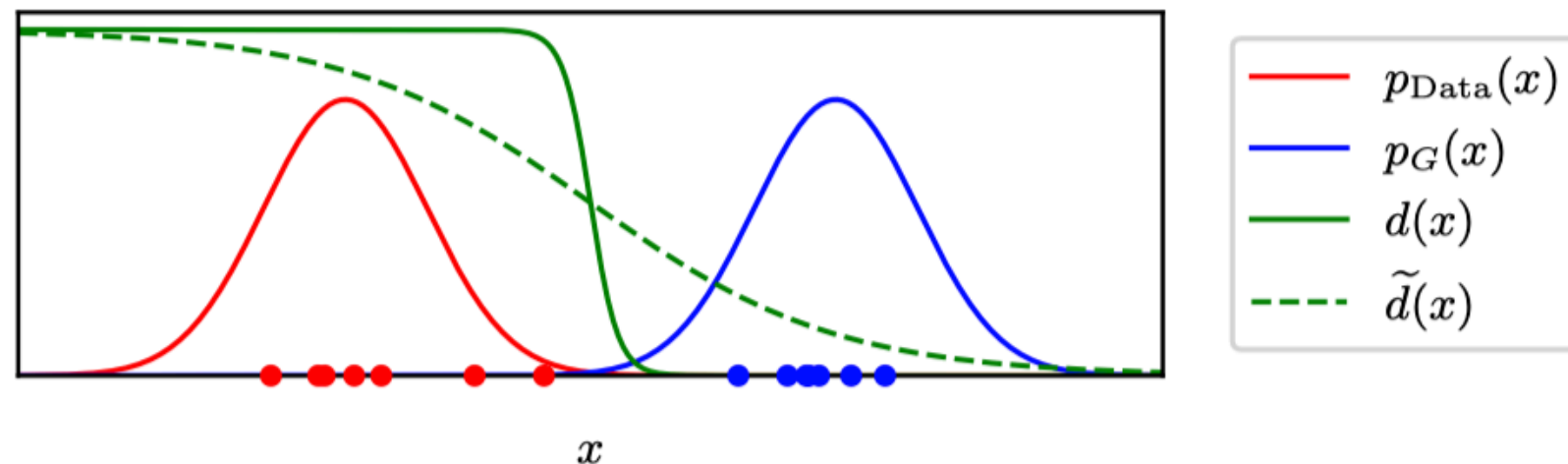Generator is trained to INCREASE the error rate.

# Loss function

- Training alternates between updating parameters of $g$ and of $d$ (one step each per minibatch).

- If generator were perfect we would have
$$d = 0.\bar{5}$$

- Discriminator can be discarded after training.

- As usual, one can make conditional version, by adding a class $c$,
$$p(\vec{x} \mid \vec{c}) \quad \text{from training data} \quad \{x_n, c_n\}$$

# GAN training

- GANs are not easy to train successfully due to the adversarial learning.

- There is no clear metric of progress because the objective can go up as well as down during training.

- Problem of mode collapse:

  - the generator network weights adapt during training such that all latent-variable samples z are mapped to a subset of possible valid outputs (e.g. only images of the number 3).

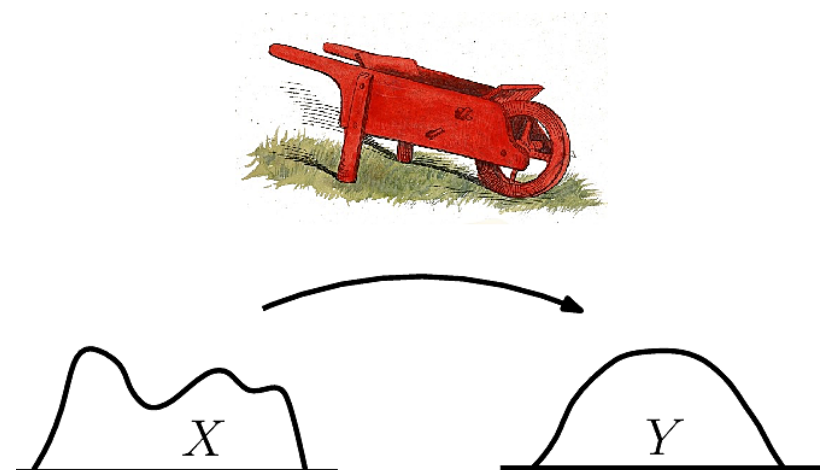  - The discriminator then assigns the value 0.5 to these instances, and training ceases.

# GAN training problems



Conceptual illustration of why it can be difficult to train GANs, showing a simple one-dimensional data space $x$ with the fixed, but unknown, data distribution $p_{\text{Data}}(x)$ and the initial generative distribution $p_G(x)$. The optimal discriminator function $d(x)$ has virtually zero gradient in the vicinity of either the training or synthetic data points, making learning very slow. A smoothed version $\widetilde{d}(x)$ of the discriminator function can lead to faster learning.

- Because the data and generative distributions are so different, the optimal discriminator function d(x) is easy to learn and has a very steep fall-off with virtually zero gradient in the vicinity of either the real or synthetic samples.

- This can be addressed e.g. by using a smoothed version d~(x) of the discriminator function.

- Numerous other modifications to the GAN error function and training procedure have been proposed to improve training.

# WGAN (briefly)

- A more direct way to ensure that the generator distribution $p_G(x)$ moves towards the data distribution $p_{data}(x)$ is to modify the error criterion to reflect how far apart the two distributions are in data space.

- This can be measured using the **Wasserstein distance**, also known as the **earth mover's distance.**

  - Imagine the distribution $p_G(x)$ as a pile of earth that is transported in small increments to construct the distribution $p_{data}(x)$. The Wasserstein metric is the total amount of earth moved multiplied by the mean distance moved.

  - In practice, this cannot be implemented directly, and it is approximated by using a discriminator network that has real-valued outputs.

- This gives rise to the **Wasserstein GAN or WGAN**

- **WGAN have more stable training, meaningful loss curves, better convergence behavior.**



https://sbl.inria.fr/doc/Earth_mover_distance-user-manual.html

# GAN vs WGAN

**Ordinary GAN (Goodfellow et al., 2014)**

- Uses the **Jensen-Shannon (JS) divergence** to measure the distance between the real data distribution $P_r$ and the generated data distribution $P_g$.

- Objective:

$$\min_G \max_D \mathbb{E}_{x \sim P_r}[\log D(x)] + \mathbb{E}_{z \sim P_z}[\log(1 - D(G(z)))]$$

- This setup can suffer from **vanishing gradients** and **mode collapse**, especially when $P_r$ and $P_g$ have little overlap.

**Wasserstein GAN (Arjovsky et al., 2017)**

- Replaces JS divergence with the **Wasserstein-1 distance (a.k.a. Earth Mover's Distance)**, which provides a better behaved and meaningful gradient signal, even when $P_r$ and $P_g$ are disjoint.

- Objective:

$$\min_G \max_{D \in \mathcal{D}} \mathbb{E}_{x \sim P_r}[D(x)] - \mathbb{E}_{z \sim P_z}[D(G(z))]$$

where $\mathcal{D}$ is the set of **1-Lipschitz functions**, enforced via:

    $\}$ *smoothness condition*

- **Weight clipping** in the original WGAN

- **Gradient penalty** in WGAN-GP (improved version)

Think of it as the **minimal "cost" of transporting mass** from the generated distribution to match the real one. Even if the two distributions don't overlap at all, it still gives a finite and informative gradient.

# GANs for images

- GANs for images use versions of CNNs for both the generator and the discriminator.

- **Discriminator uses a standard CNN classifier**

- **Generator maps lower dimensional latent space to higher dimensional image**, often using transposed convolutions.
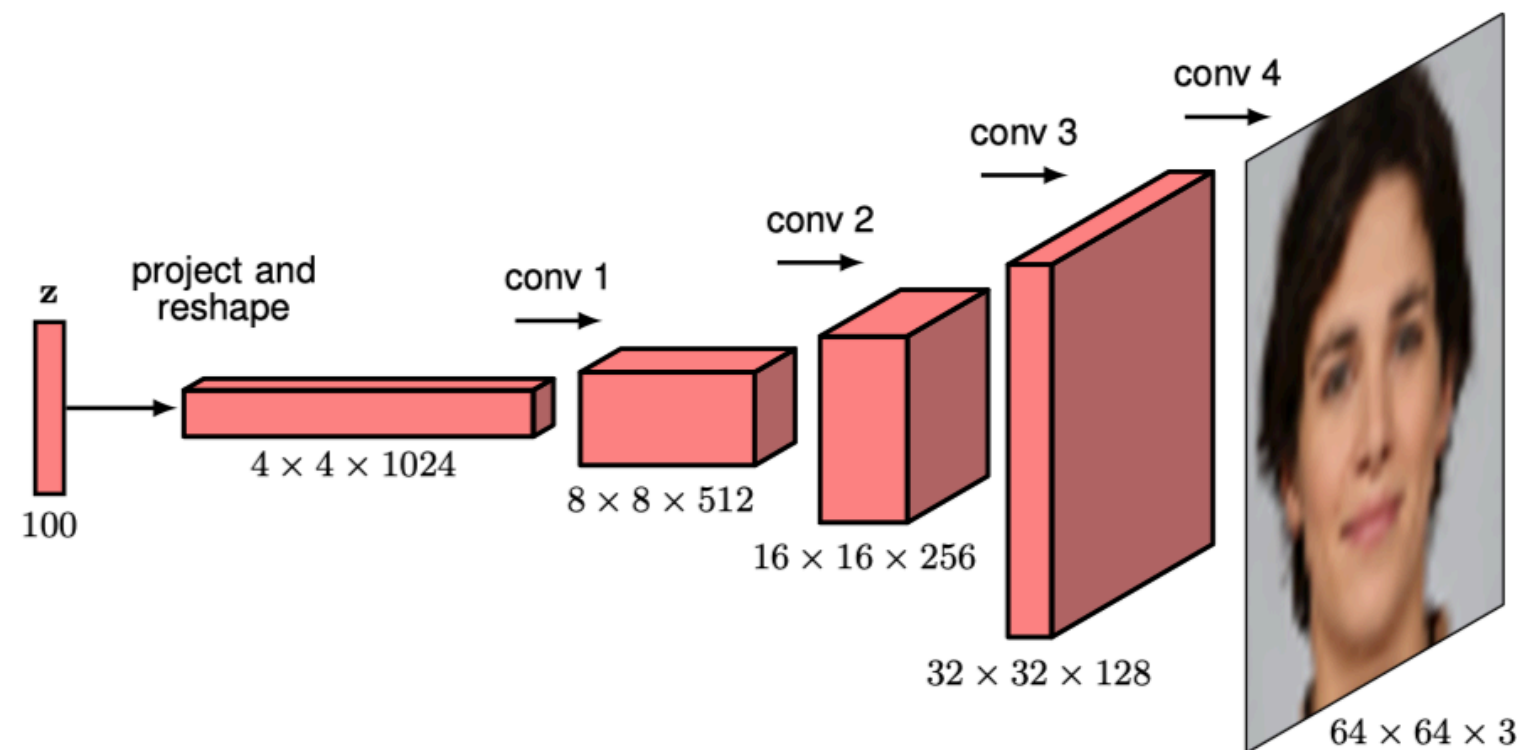


**Figure 17.4** Example architecture of a deep convolutional GAN showing the use of transpose convolutions to expand the dimensionality in successive blocks of the network.
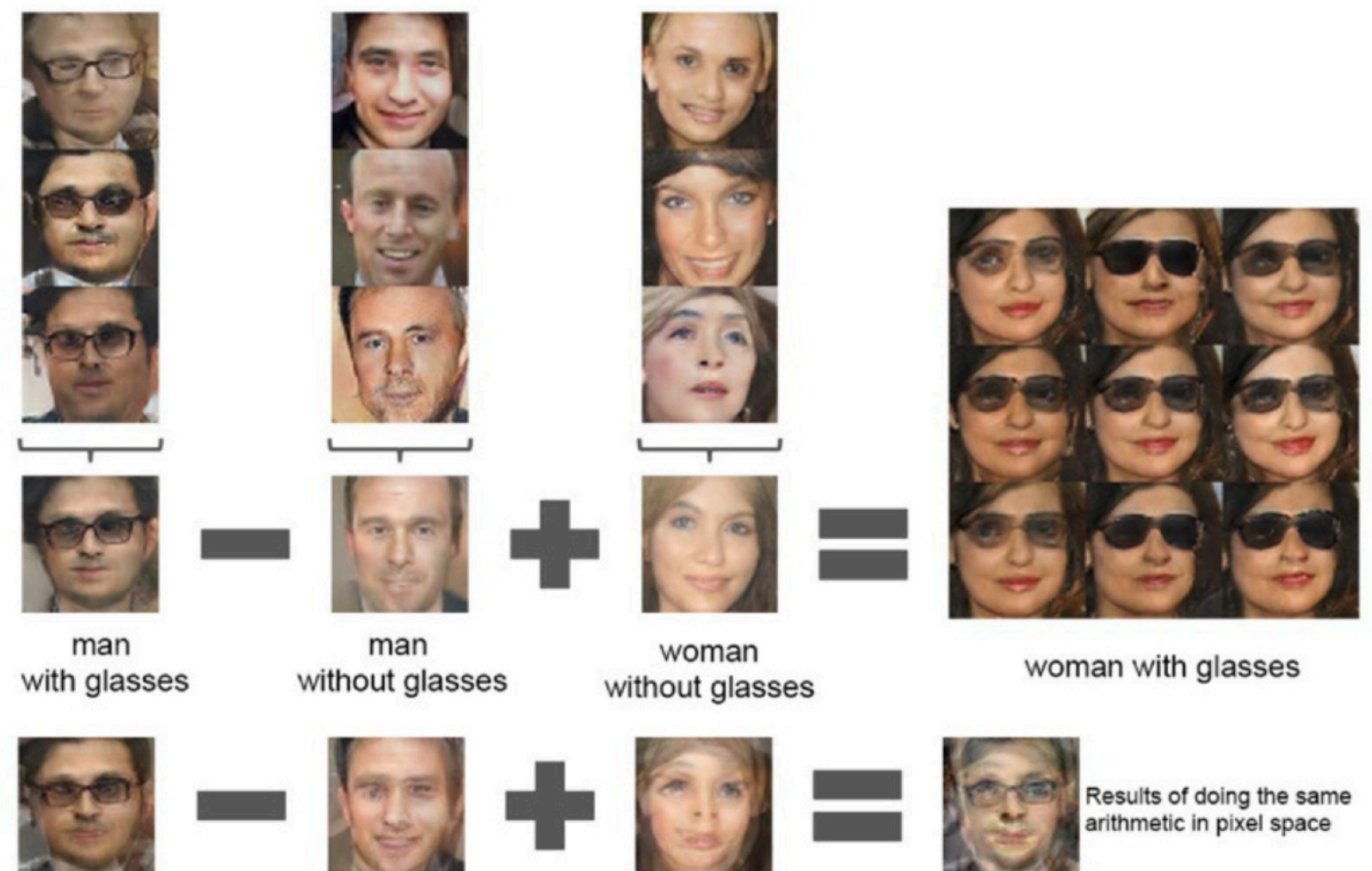
# GANs learn a meaningful latent space



**Continuous deformations in latent space**

# GANs learn a meaningful latent space



**Figure 17.10** An example of vector arithmetic in the latent space of a trained GAN. In each of the three columns, the latent space vectors that generated these images are averaged and then vector arithmetic is applied to the resulting mean vectors to create a new vector corresponding to the central image in the $3 \times 3$ array on the right. Adding noise to this vector generates another eight sample images. The four images on the bottom row show that the same arithmetic applied directly in data space simply results in a blurred image due to misalignment. [From Radford, Metz, and Chintala (2015) with permission.]

man with glasses

man without glasses

woman without glasses

woman with glasses

Results of doing the same arithmetic in pixel space

**Arithmetic in latent space**
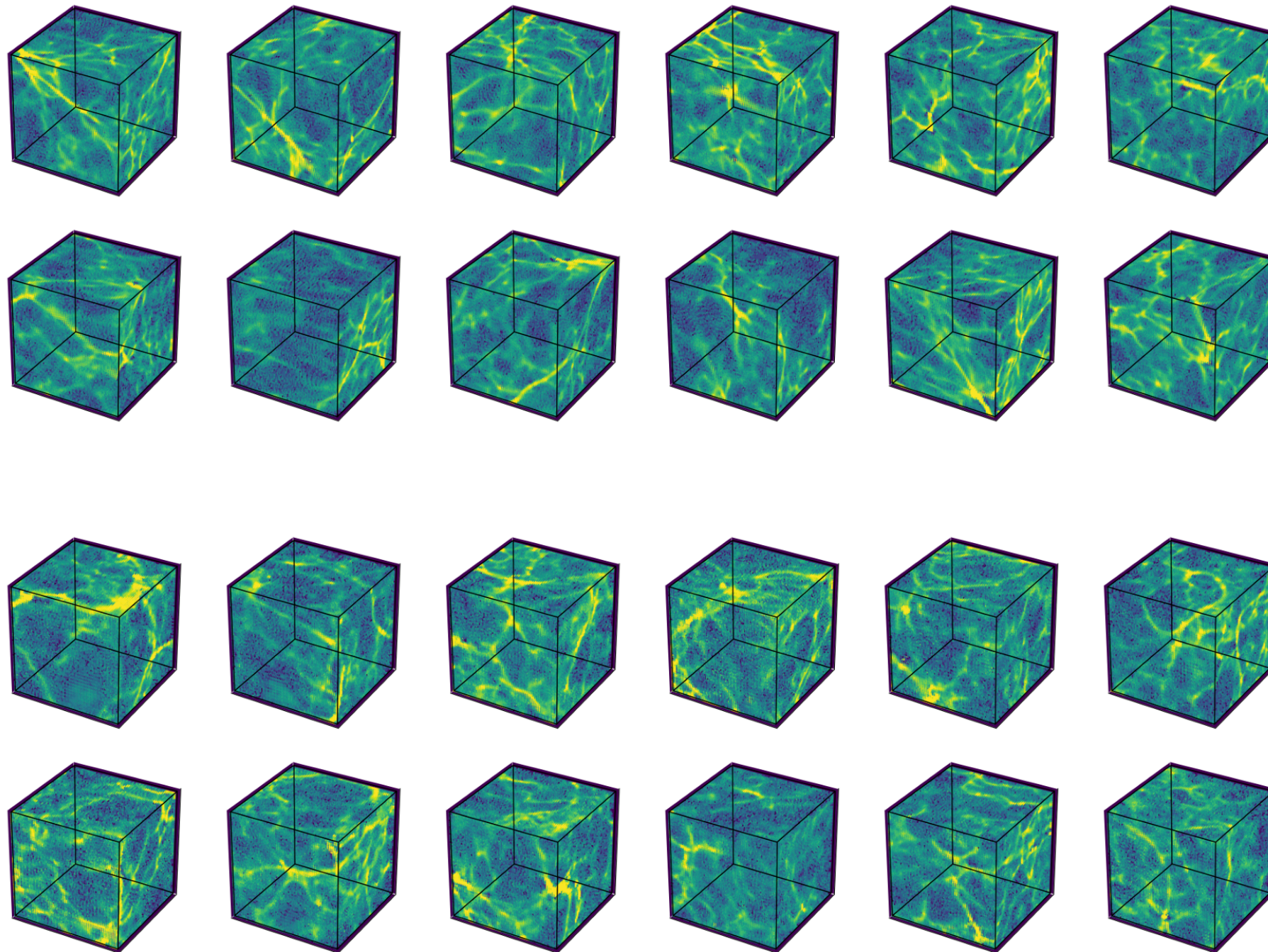
# WGAN GAN in cosmology

**Figure 1.** 3D HI distributions from IllustrisTNG (top rows) and WGAN (bottom rows). The network successfully produces spatial distributions with all the elements of the HI web: filaments, voids and dense regions.
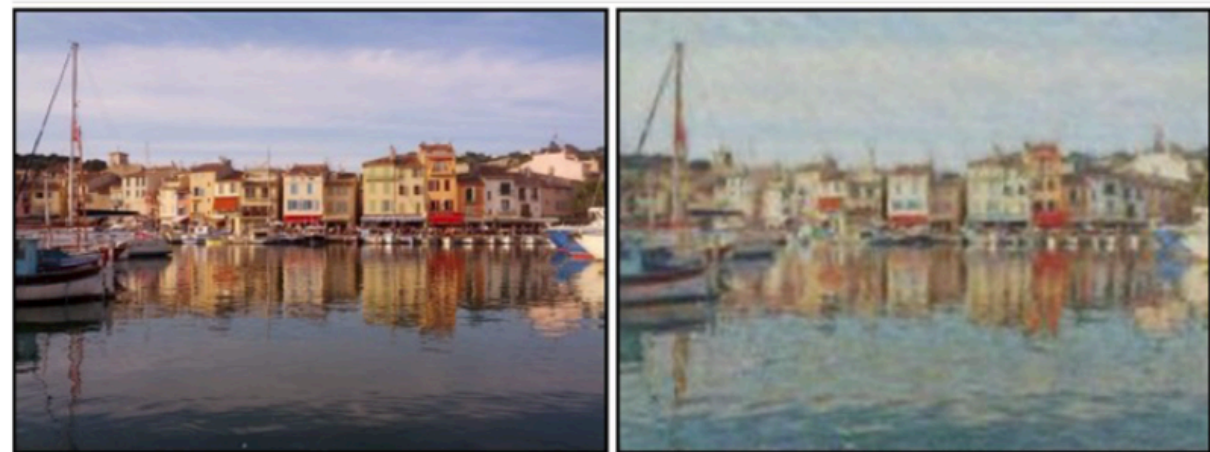
# Cycle GANs

# CycleGANs for image-to-image translation

**Figure 17.6**  Examples of image translation using a CycleGAN showing the synthesis of a photographic-style image from a Monet painting (top row) and the synthesis of an image in the style of a Monet painting from a photograph (bottom row). [From Zhu *et al.* (2017) with permission.]



Monet → photograph

photograph → Monet

- **Style transfer (image-to-image) is useful in physics.**

- **CycleGANs do not require paired training data. This is a key strength.**

# CycleGANs for image-to-image translation

- The aim is to learn two bijective (one-to-one) mappings, one that goes from the **domain X of photographs** to the **domain Y of Monet paintings** and one in the reverse direction.

- To achieve this, CycleGAN makes use o**f two conditional generators, gX and gY, and two discriminators, dX and dY.**

- We need to ensure that **when a photograph is translated into a painting and then back into a photograph it should be close to the original photograph**, thereby ensuring that the generated painting retains sufficient information about the photograph to allow the photograph to be reconstructed.
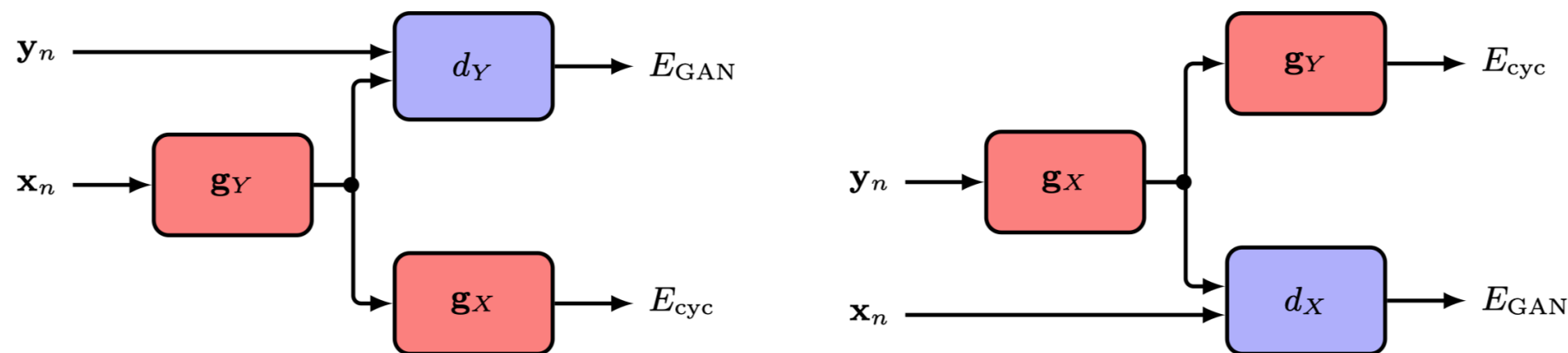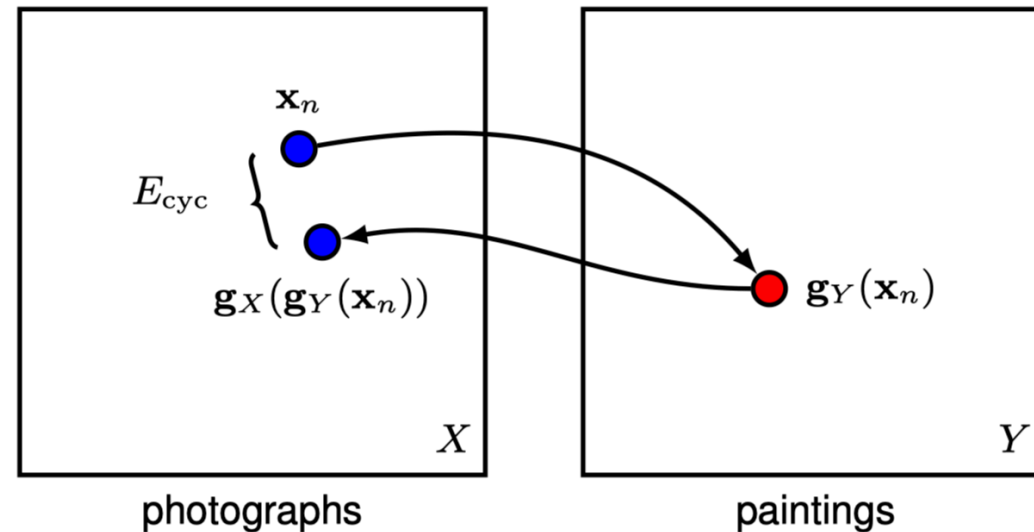


**Figure 17.8**   Flow of information through a CycleGAN. The total error for the data points $\mathbf{x}_n$ and $\mathbf{y}_n$ is the sum of the four component errors.

$x_n$ **is used as a conditioning input** to $g_Y$. There is **no base distribution** $z$ sampled in this architecture.

The transformation is **fully deterministic and conditional**.

# CycleGANs for image-to-image translation

**Figure 17.7** Diagram showing how the cycle consistency error is calculated for an example photograph $\mathbf{x}_n$. The photograph is first mapped into the painting domain using the generator $\mathbf{g}_Y$, and the resulting vector is then mapped back into the photograph domain using the generator $\mathbf{g}_X$. The discrepancy between the resulting photograph and the original $\mathbf{x}_n$ defines a contribution to the cycle consistency error. An analogous process is used to calculate the contribution to the cycle consistency error from a painting $\mathbf{y}_n$ by mapping it to a photograph using $\mathbf{g}_X$ and then back to a painting using $\mathbf{g}_Y$.
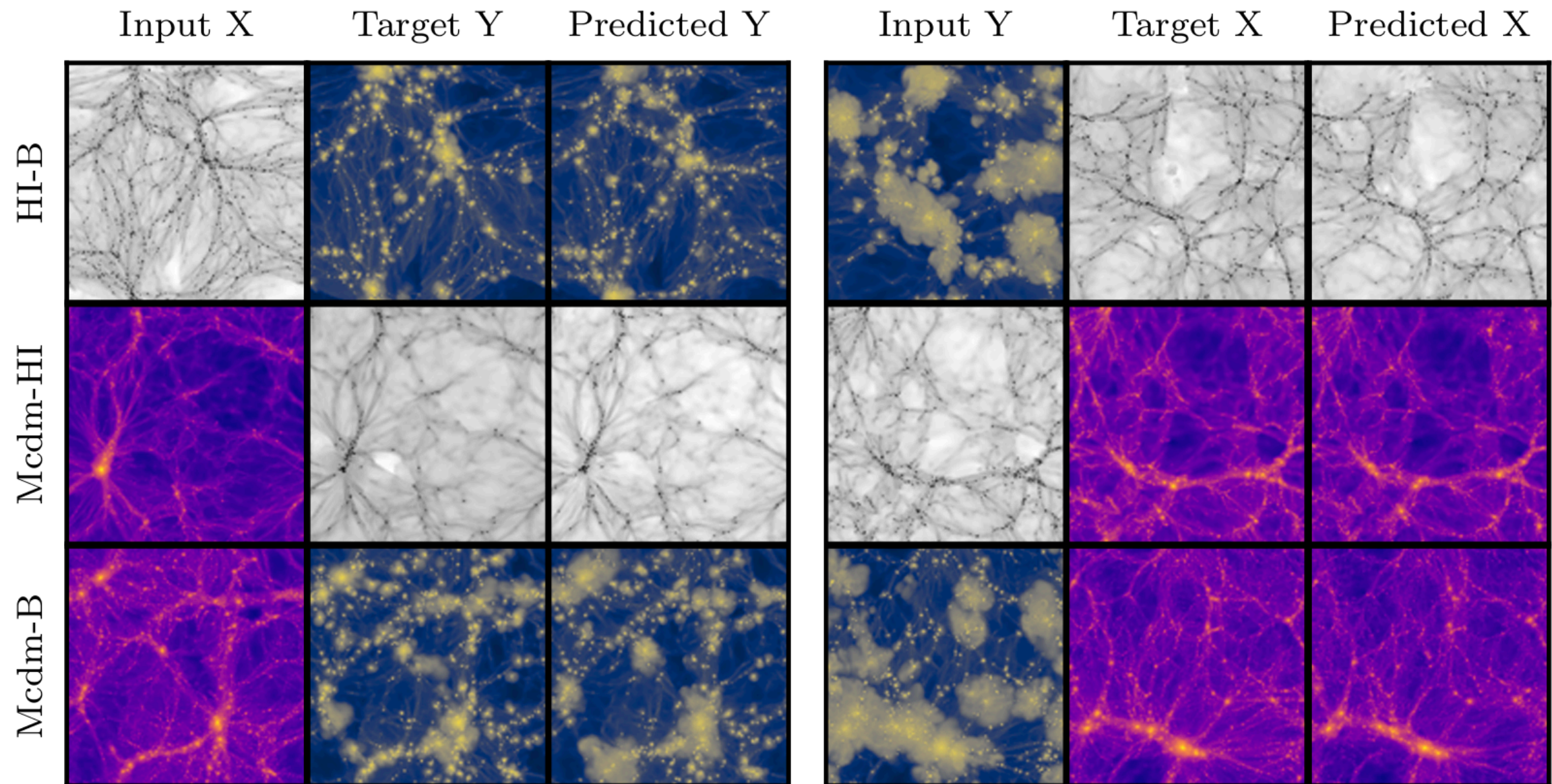


photographs   paintings

$$E_{\mathrm{cyc}}(\mathbf{w}_X, \mathbf{w}_Y) = \frac{1}{N_X} \sum_{n \in X} \|\mathbf{g}_X(\mathbf{g}_Y(\mathbf{x}_n)) - \mathbf{x}_n\|_1$$

$$+ \frac{1}{N_Y} \sum_{n \in Y} \|\mathbf{g}_Y(\mathbf{g}_X(\mathbf{y}_n)) - \mathbf{y}_n\|_1$$

total loss: $\mathcal{L} = E_{\mathrm{GAN}}(\mathbf{w}_X, \phi_X) + E_{\mathrm{GAN}}(\mathbf{w}_Y, \phi_Y) + \eta E_{\mathrm{cyc}}(\mathbf{w}_X, \mathbf{w}_Y)$

# CycleGAN in cosmology

**https://arxiv.org/abs/2303.07473**   **Invertible mapping between fields in CAMELS**

# Recall: Normalizing Flows

# Recall Normalizing flows

Method to learn the PDF of data $p(x)$.

Start with a base distribution:

$$p_\mathbf{z}(\mathbf{z})$$

Transform using an invertible function

$$\mathbf{x} = \mathbf{f}(\mathbf{z}, \mathbf{w})$$
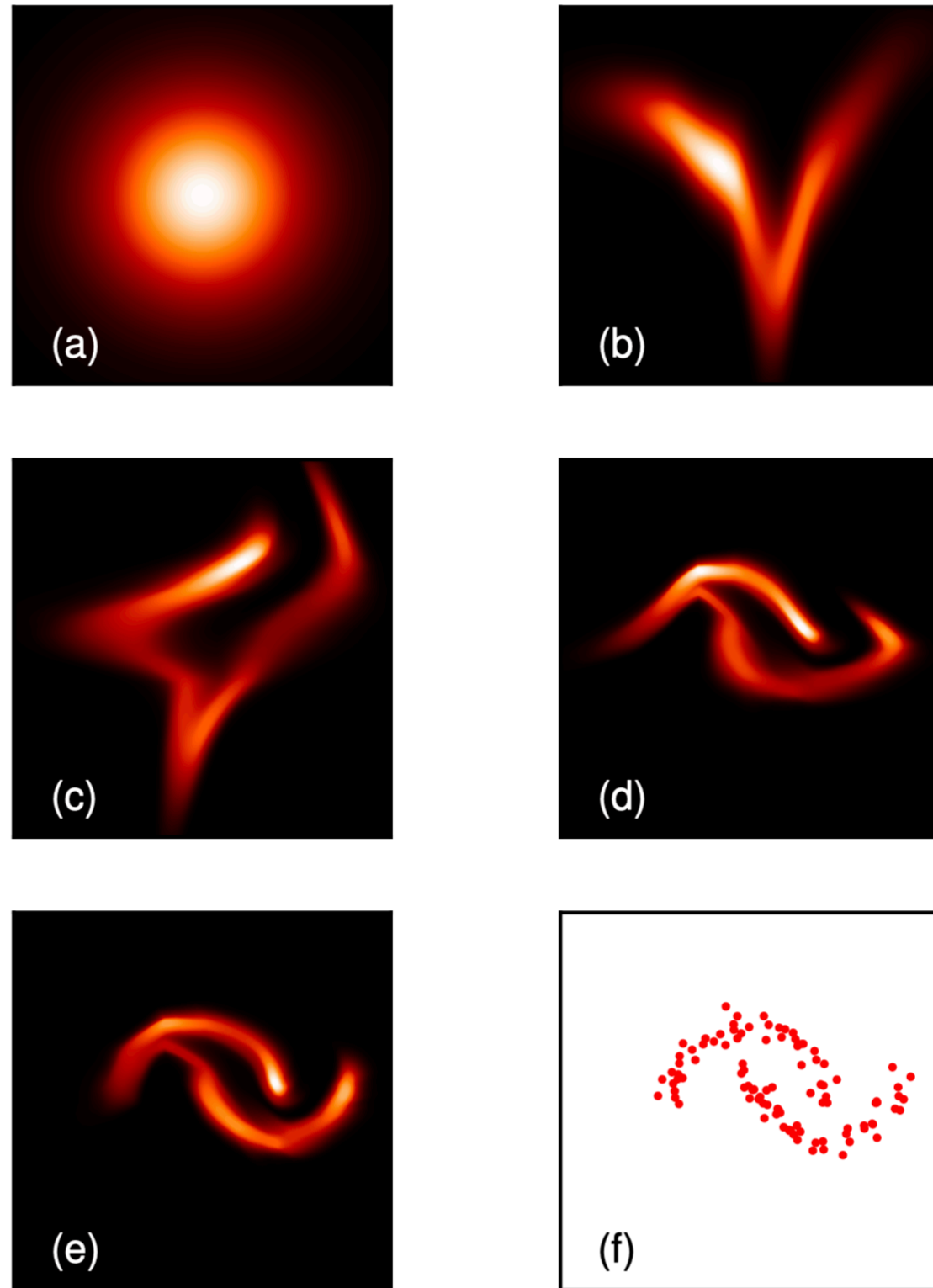
The density transforms as

$$p_\mathbf{x}(\mathbf{x}|\mathbf{w}) = p_\mathbf{z}(\mathbf{g}(\mathbf{x}, \mathbf{w})) \left| \det \mathbf{J}(\mathbf{x}) \right|$$

with Jacobian $\quad J_{ij}(\mathbf{x}) = \dfrac{\partial g_i(\mathbf{x}, \mathbf{w})}{\partial x_j}$

usually we compose $f$ from simple stacked functions.

# Example in 2d

**Figure 18.3** Illustration of the real NVP normalizing flow model applied to the two-moons data set showing (a) the Gaussian base distribution, (b) the distribution after a transformation of the vertical axis only, (c) the distribution after a subsequent transformation of the horizontal axis, (d) the distribution after a second transformation of the vertical axis, (e) the distribution after a second transformation of the horizontal axis, and (f) the data set on which the model was trained.

# Recall auto-regressive flows
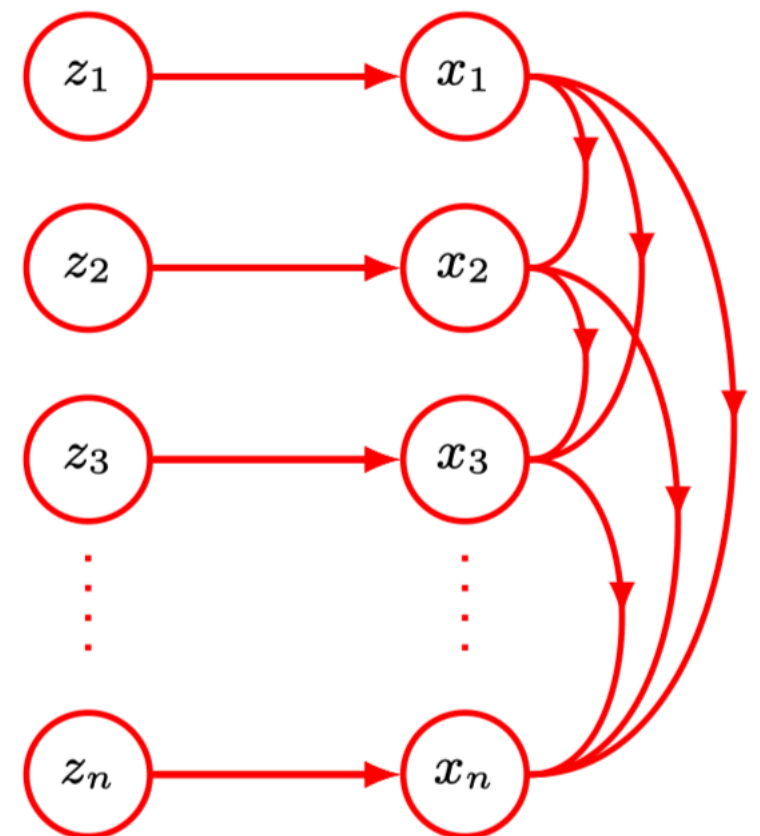
Any PDF can be written as follows

$$p(x_1, \ldots, x_D) = \prod_{i=1}^{D} p(x_i | \mathbf{x}_{1:i-1})$$

Masked auto regressive flows construct a flow based on transformations

$$x_i = h(z_i, \mathbf{g}_i(\mathbf{x}_{1:i-1}, \mathbf{w}_i))$$

conditioner

"transformer" or "coupling function"

# Neural ODE, Neural ODE flows

# Neural ODE

- We have seen that neural networks are especially useful when they comprise many layers of processing, and so we can ask **what happens if we explore the limit of an infinitely large number of layers.**

- Consider a residual network where each layer of processing generates an output given by the input vector with the addition of some parameterized nonlinear function of that input vector:

$$\mathbf{z}^{(t+1)} = \mathbf{z}^{(t)} + \mathbf{f}(\mathbf{z}^{(t)}, \mathbf{w})$$

- $t = 1, \ldots, T$ labels the layers in the network.

- Note that we have used the same function at each layer, with a shared parameter vector w.

- In the infinite limit of the number of layers we get an **ordinary (single parameter t) differential equation**:

$$\frac{d\mathbf{z}(\mathbf{t})}{dt} = \mathbf{f}(\mathbf{z}(t), \mathbf{w})$$

# Neural ODE

The solution to ODE $\dfrac{d\mathbf{z(t)}}{dt} = \mathbf{f}(\mathbf{z}(t), \mathbf{w})$

for initial condition / input $\mathbf{z}(0)$

is simply $\mathbf{z}(T) = \displaystyle\int_0^T \mathbf{f}(\mathbf{z}(t), \mathbf{w})\, dt$

Can be evaluated using s**tandard numerical integration packages.**

The simplest method for solving differential equations is Euler's forward integration method, which corresponds to the expression:

$$\frac{d\mathbf{z(t)}}{dt} = \mathbf{f}(\mathbf{z}(t), \mathbf{w})$$

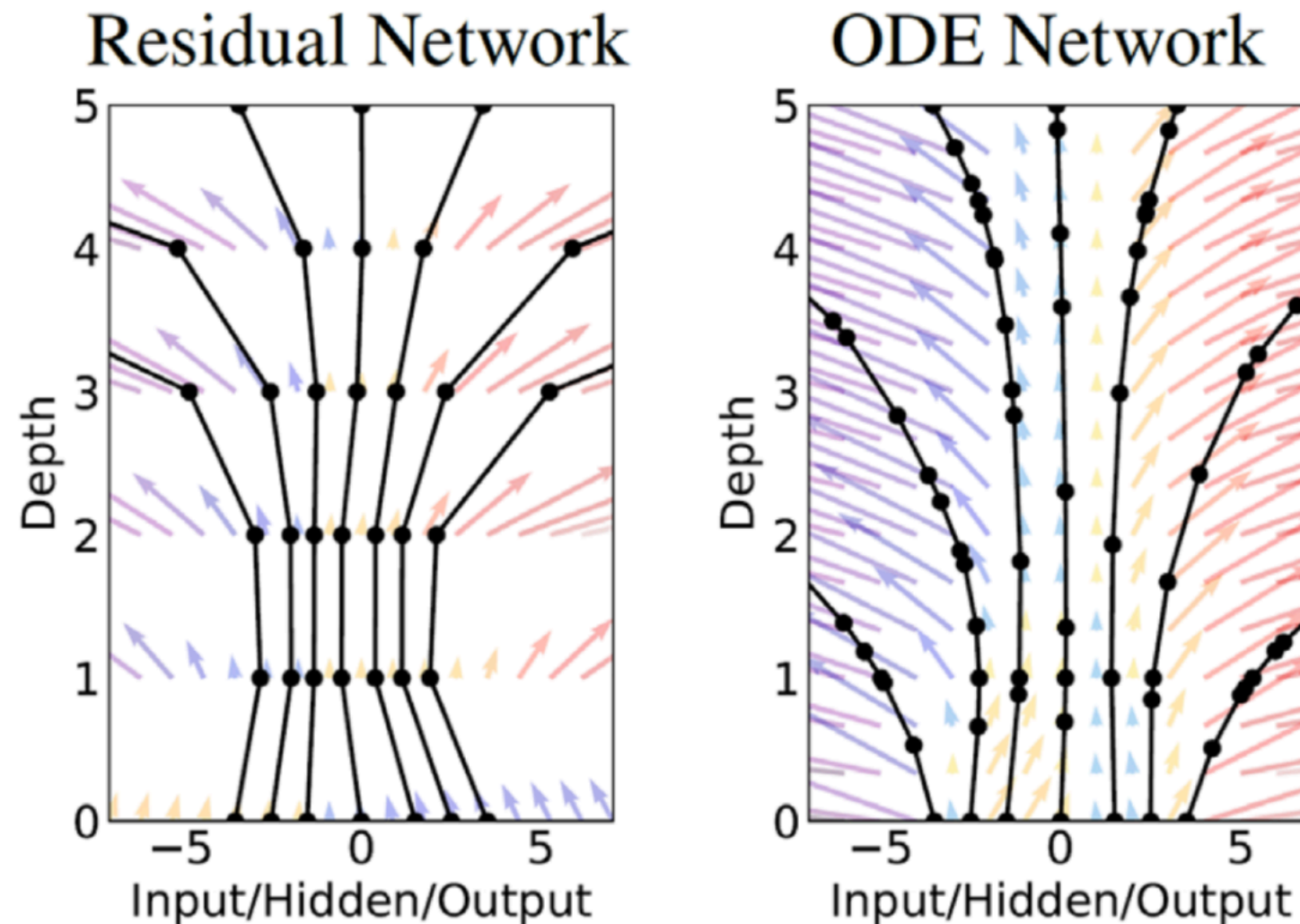But there are better ODE solvers / numerical integrators.

**Figure 18.5** Comparison of a conventional layered network with a neural differential equation. The diagram on the left corresponds to a residual network with five layers and shows trajectories for several starting values of a single scalar input. The diagram on the right shows the result of numerical integration of a continuous neural ODE, again for several starting values of the scalar input, in which we see that the function is not evaluated at uniformly-spaced time intervals, but instead the evaluation points are chosen adaptively by the numerical solver and depend on the choice of input value. [From Chen *et al.* (2018) with permission.]

# Backprop. in Neural ODE

- We have input $\mathbf{z}(0)$ and want to learn output $\mathbf{z}(T)$ which minimizes loss $\mathcal{L}$.
- Simple auto-diff through ODE solver is too costly.
- Instead: "adjoint sensitivity method"; a continuous analogue to backprop.

- Define the adjoint:

$$\mathbf{a}(t) = \frac{\mathrm{d}L}{\mathrm{d}\mathbf{z}(t)}$$

which satisfies
$$\frac{\mathrm{d}\mathbf{a}(t)}{\mathrm{d}t} = -\mathbf{a}(t)^{\mathrm{T}}\nabla_{\mathbf{z}}f(\mathbf{z}(t), \mathbf{w})$$

This leads to the following integral for the loss derivative:

$$\nabla_{\mathbf{w}}L = -\int_{0}^{T} \mathbf{a}(t)^{\mathrm{T}}\nabla_{\mathbf{w}}f(\mathbf{z}(t), \mathbf{w})\,\mathrm{d}t$$

# What are neural ODE good for?

- Advantages:

  - Neural ODEs can naturally handle continuous-time data in which observations occur at arbitrary times.

    - Physics of course has a lot of such data

  - Adaptive computation: A high level of accuracy in the solver can be used during training, with a lower accuracy, and hence fewer function evaluations, during inference in applications for which compute resources are limited.

  - Efficient memory: One benefit of neural ODEs trained using the adjoint method, compared to conventional layered networks, is that there is no need to store the intermediate results of the forward propagation, and hence the memory cost is constant.

# Neural ODE flows

- We can make **use of a neural ordinary differential equation to define an alternative approach to the construction of tractable normalizing flow models**.

- A neural ODE defines a highly flexible transformation from an input vector z(0) to an output vector z(T) in terms of a differential equation of the form

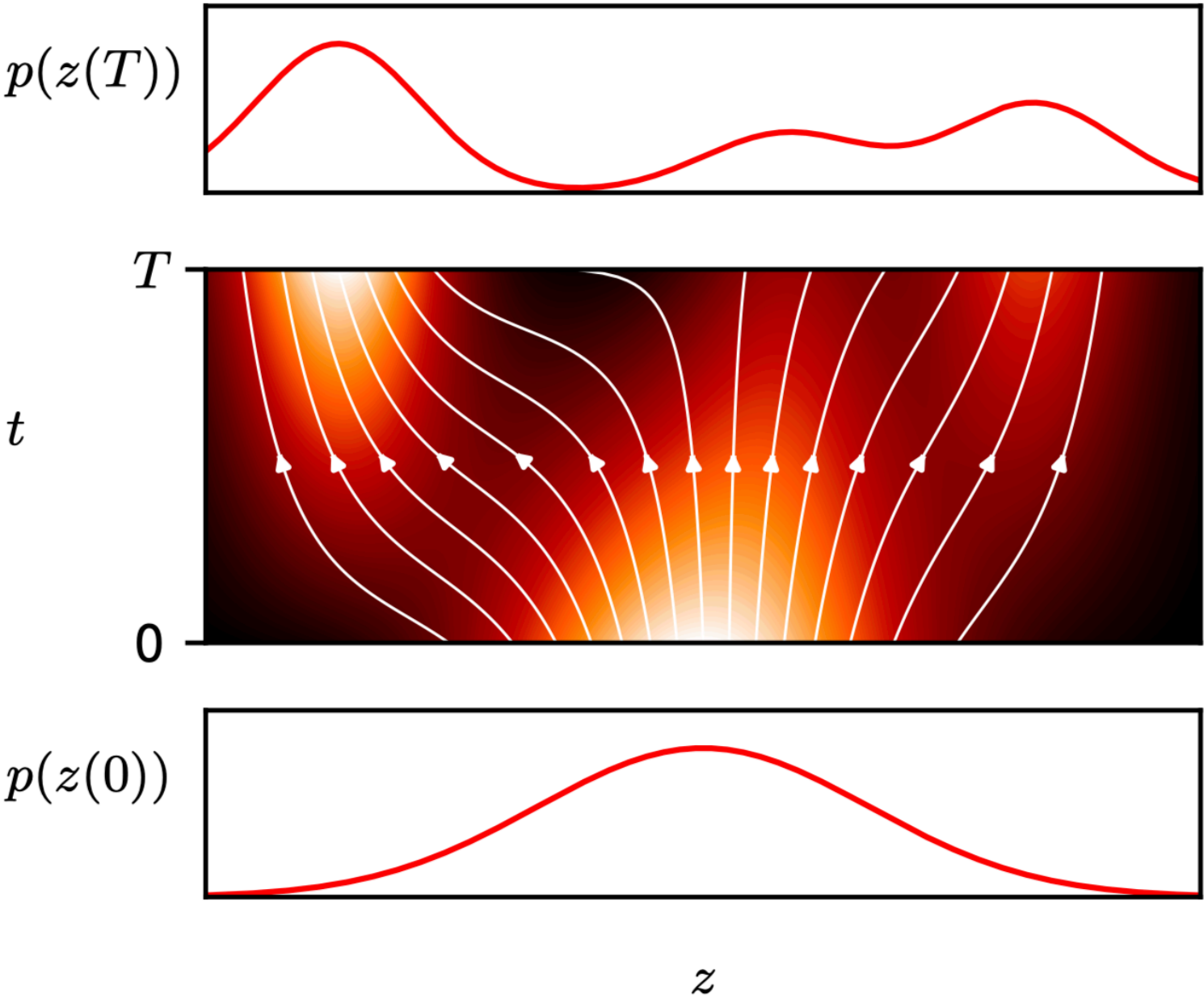$$\frac{d\mathbf{z}(t)}{dt} = \mathbf{f}(\mathbf{z}(t), \mathbf{w})$$

- Under this transformation densities transform as

$$\frac{d \ln p(\mathbf{z}(t))}{dt} = -\text{Tr}\left(\frac{\partial \mathbf{f}}{\partial \mathbf{z}(t)}\right)$$

  Which we can integrate.

- The resulting framework is known as a **continuous normalizing flows.**

**Figure 18.6** Illustration of a continuous normalizing flow showing a simple Gaussian distribution at $t = 0$ that is continuously transformed into a multimodal distribution at $t = T$. The flow lines show how points along the $z$-axis evolve as a function of $t$. Where the flow lines spread apart the density is reduced, and where they move together the density is increased.

# Example: Neural ODE in cosmology

- Neural ODE are attractive because of their elegance, but they are not necessarily the best performing tool.

- As an example, here is a paper that uses them in cosmology:

## Hybrid Physical–Neural ODEs for Fast N-body Simulations

Denise Lanzieri, François Lanusse, Jean-Luc Starck

We present a new scheme to compensate for the small-scales approximations resulting from Particle-Mesh (PM) schemes for cosmological N-body simulations. This kind of simulations are fast and low computational cost realizations of the large scale structures, but lack resolution on small scales. To improve their accuracy, we introduce an additional effective force within the differential equations of the simulation, parameterized by a Fourier-space Neural Network acting on the PM-estimated gravitational potential. We compare the results for the matter power spectrum obtained to the ones obtained by the PGD scheme (Potential gradient descent scheme). We notice a similar improvement in term of power spectrum, but we find that our approach outperforms PGD for the cross-correlation coefficients, and is more robust to changes in simulation settings (different resolutions, different cosmologies).

# Symmetric flows for molecules (graphs)

Computer Science > Machine Learning

[Submitted on 19 May 2021 (v1), last revised 14 Jan 2022 (this version, v4)]

## E(n) Equivariant Normalizing Flows

Victor Garcia Satorras, Emiel Hoogeboom, Fabian B. Fuchs, Ingmar Posner, Max Welling

This paper introduces a generative model equivariant to Euclidean symmetries: E(n) Equivariant Normalizing Flows (E-NFs). To construct E-NFs, we take the discriminative E(n) graph neural networks and integrate them as a differential equation to obtain an invertible equivariant function: a continuous-time normalizing flow. We demonstrate that E-NFs considerably outperform baselines and existing methods from the literature on particle systems such as DW4 and LJ13, and on molecules from QM9 in terms of log-likelihood. To the best of our knowledge, this is the first flow that jointly generates molecule features and positions in 3D.
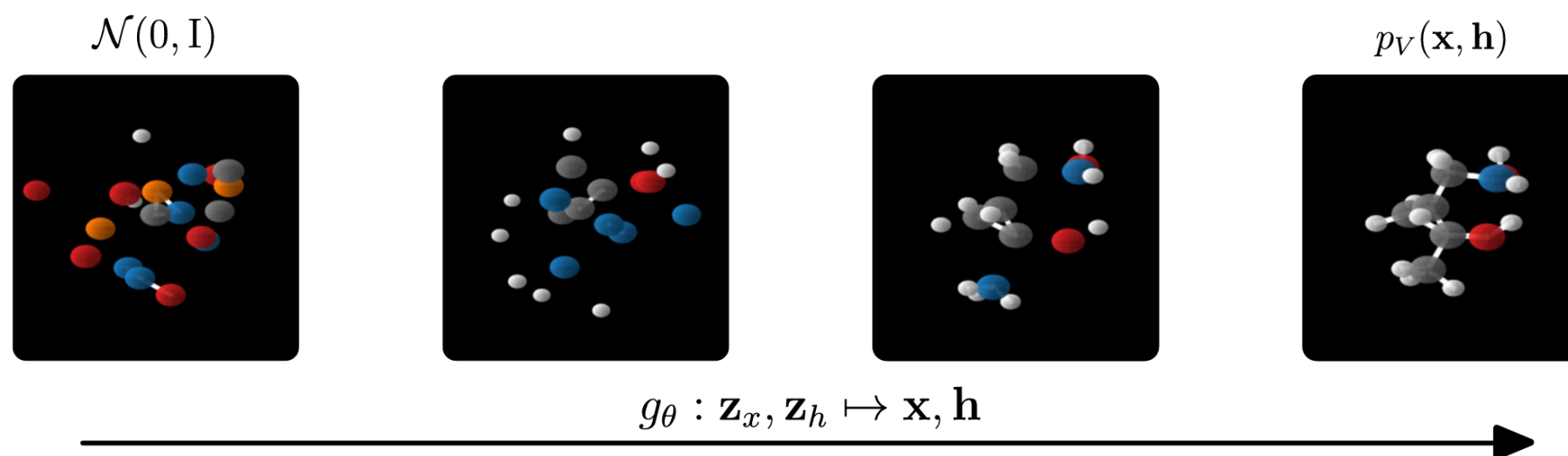
$\mathcal{N}(0, \mathbf{I})$ $p_V(\mathbf{x}, \mathbf{h})$

$g_\theta : \mathbf{z}_x, \mathbf{z}_h \mapsto \mathbf{x}, \mathbf{h}$

Figure 1: Overview of our method in the sampling direction. An equivariant invertible function $g_\theta$ has learned to map samples from a Gaussian distribution to molecules in 3D, described by $\mathbf{x}, \mathbf{h}$.

# Flow Matching

- **Flow matching** is a method for training generative models using **vector fields** that transform a simple base distribution (like a Gaussian) into a complex data distribution — without solving the ODE during training. It's closely related to neural ODE flows but avoids some of the expensive computations.

- https://diffusionflow.github.io/

  - Flow matching and diffusion models are two popular frameworks in generative modeling. Despite seeming similar, there is some confusion in the community about their exact connection. In this post, we aim to clear up this confusion and show that **diffusion models and Gaussian flow matching are the same, although different model specifications can lead to different network outputs and sampling schedules**. This is great news, it means you can use the two frameworks interchangeably.

# Side note: Other physical processes for sampling?

# Other physical processes as generative models?

- Diffusion models can be seen example of how physics can be used to do machine learning.

- **Since diffusion is a physical process, one may ask if there are other physical processes which can be use as generative models**. This idea was developed here:

  - https://arxiv.org/abs/2209.11178 Poisson Flow Generative Models

  - https://arxiv.org/abs/2302.04265 PFGM++: Unlocking the Potential of Physics-Inspired Generative Models

  - https://arxiv.org/abs/2304.02637 GenPhys: From Physical Processes to Generative Models

- While so far these results have not been very important in practice, let's have a quick look in the papers because they are a great example of combining physics and ML.
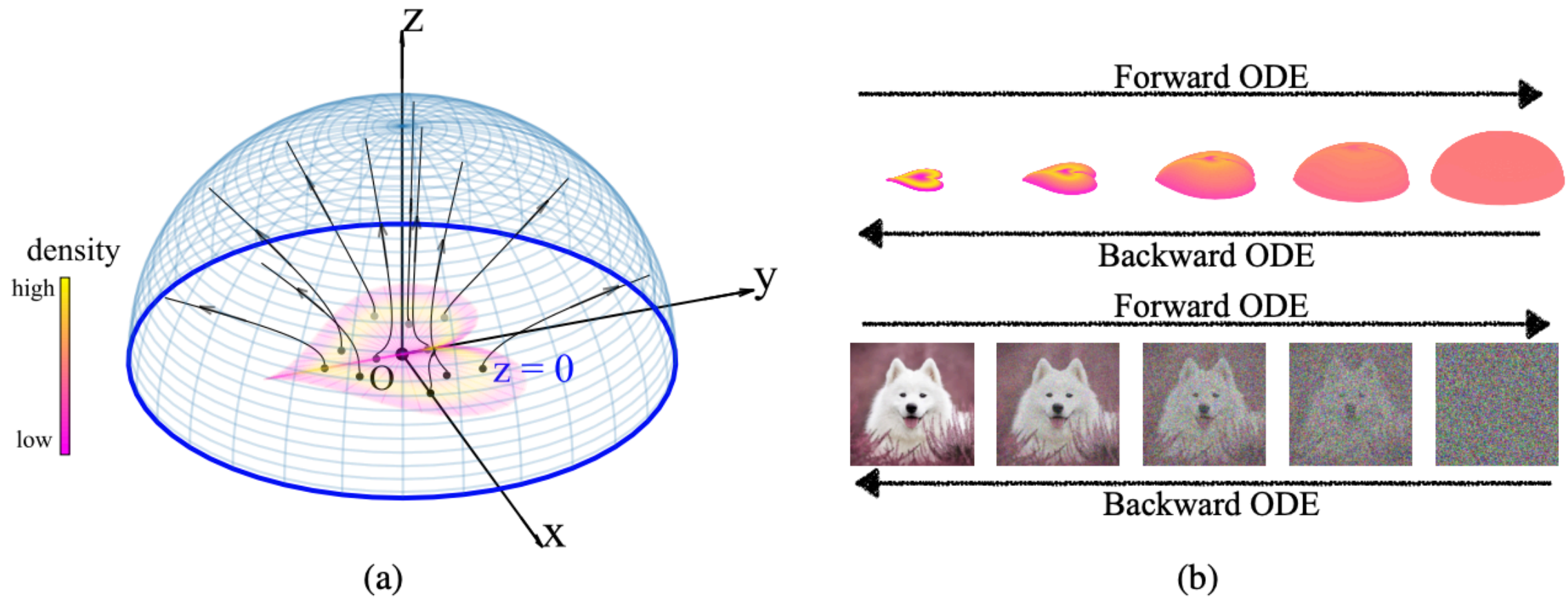
Poisson flow from https://arxiv.org/pdf/2209.11178.pdf



Figure 1: **(a)** 3D Poisson field trajectories for a heart-shaped distribution **(b)** The evolvements of a distribution (**top**) or an (augmented) sample (**bottom**) by the forward/backward ODEs pertained to the Poisson field.

# Course logistics

- **Reading for this lecture:**
  - This lecture was based in part on the book by Bishop linked on the website.