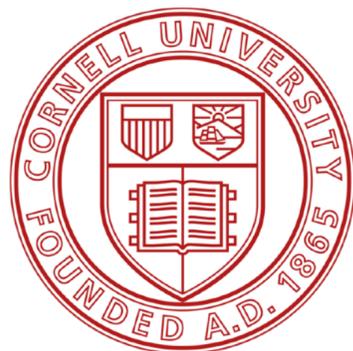


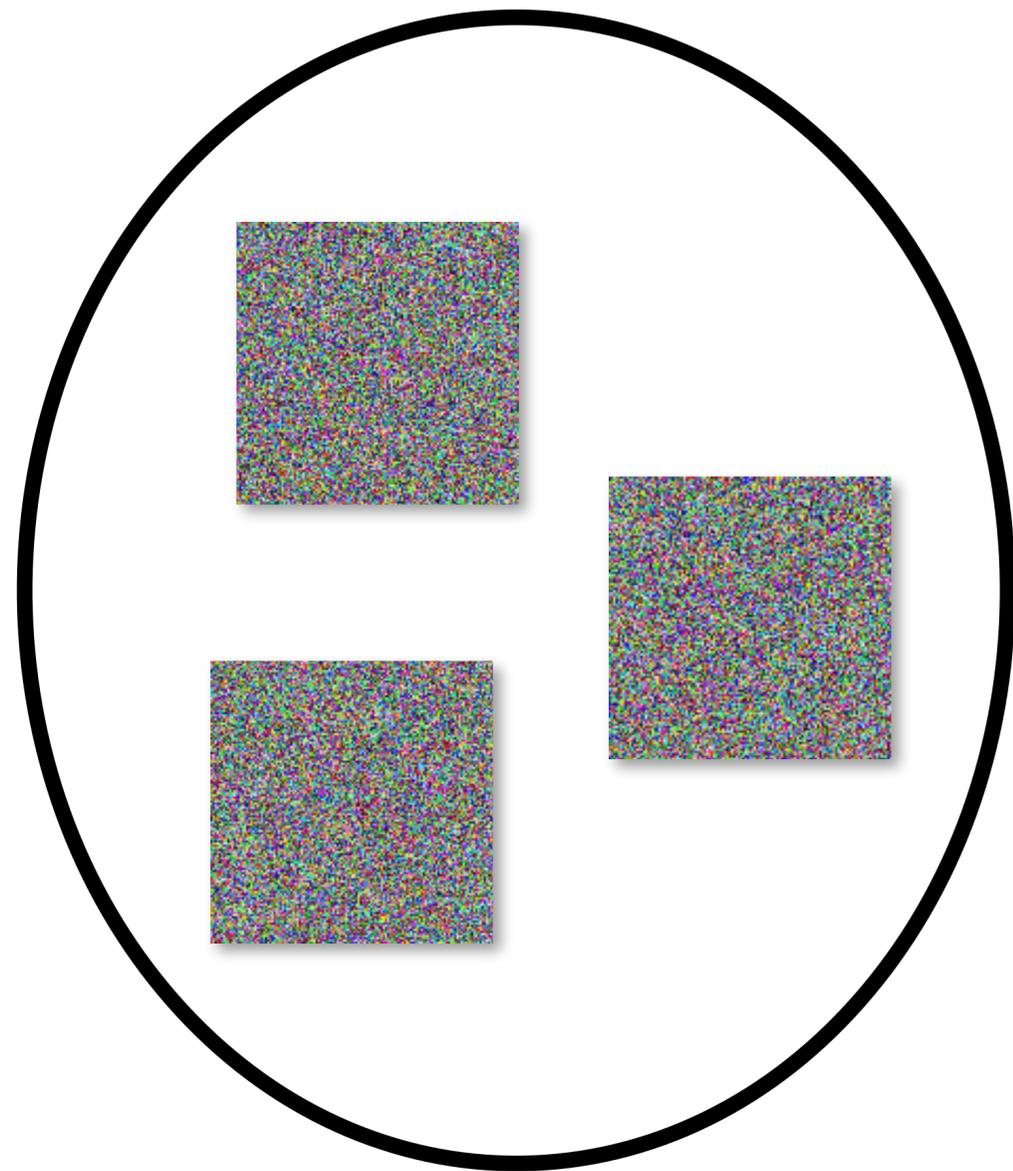
Lecture 16: Flow matching

CS 5788: Introduction to Generative Models



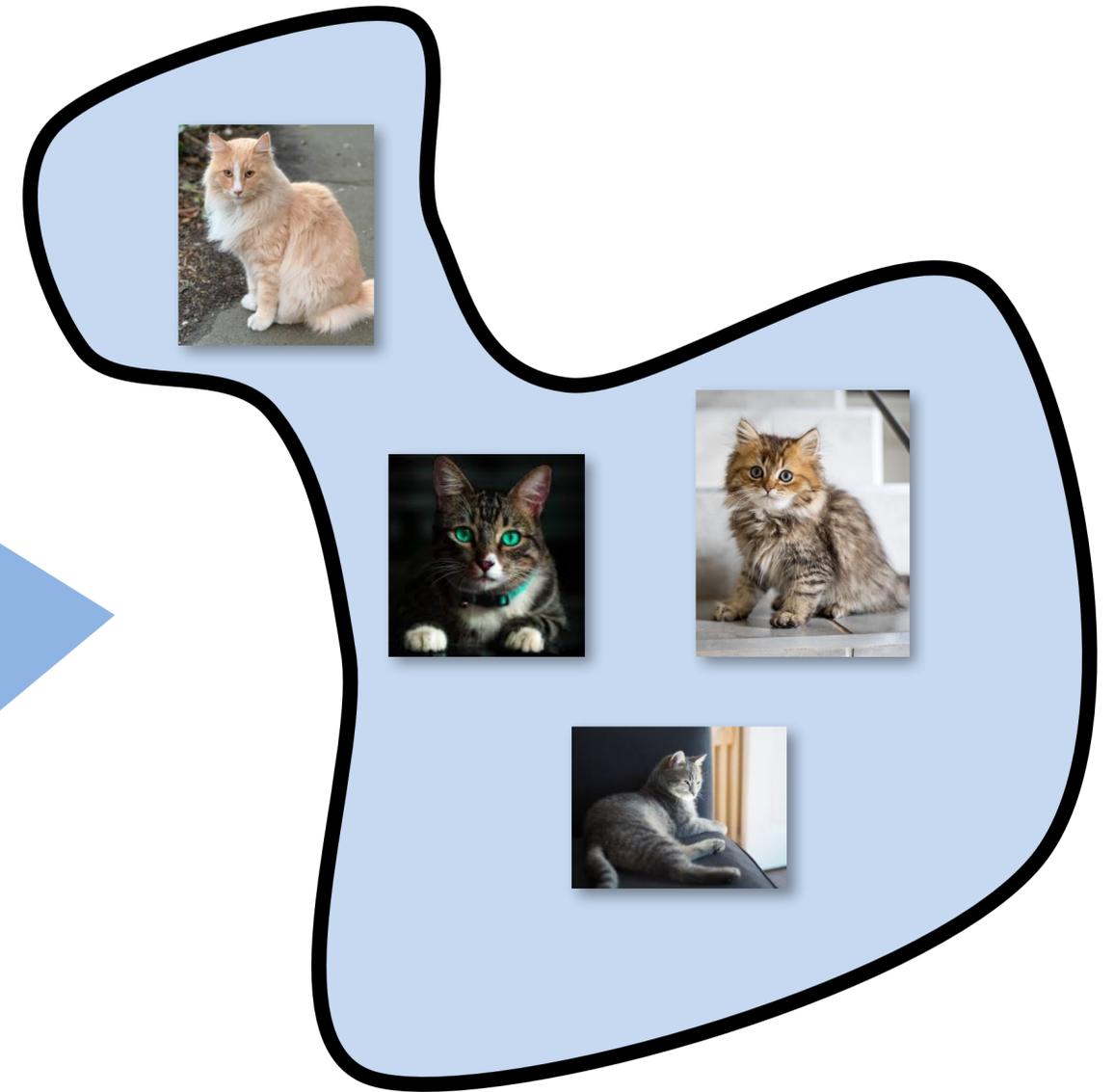
Most slides from Angjoo Kanazawa (also includes slides from Yaron Lipman and Steve Seitz)

Recall: diffusion models

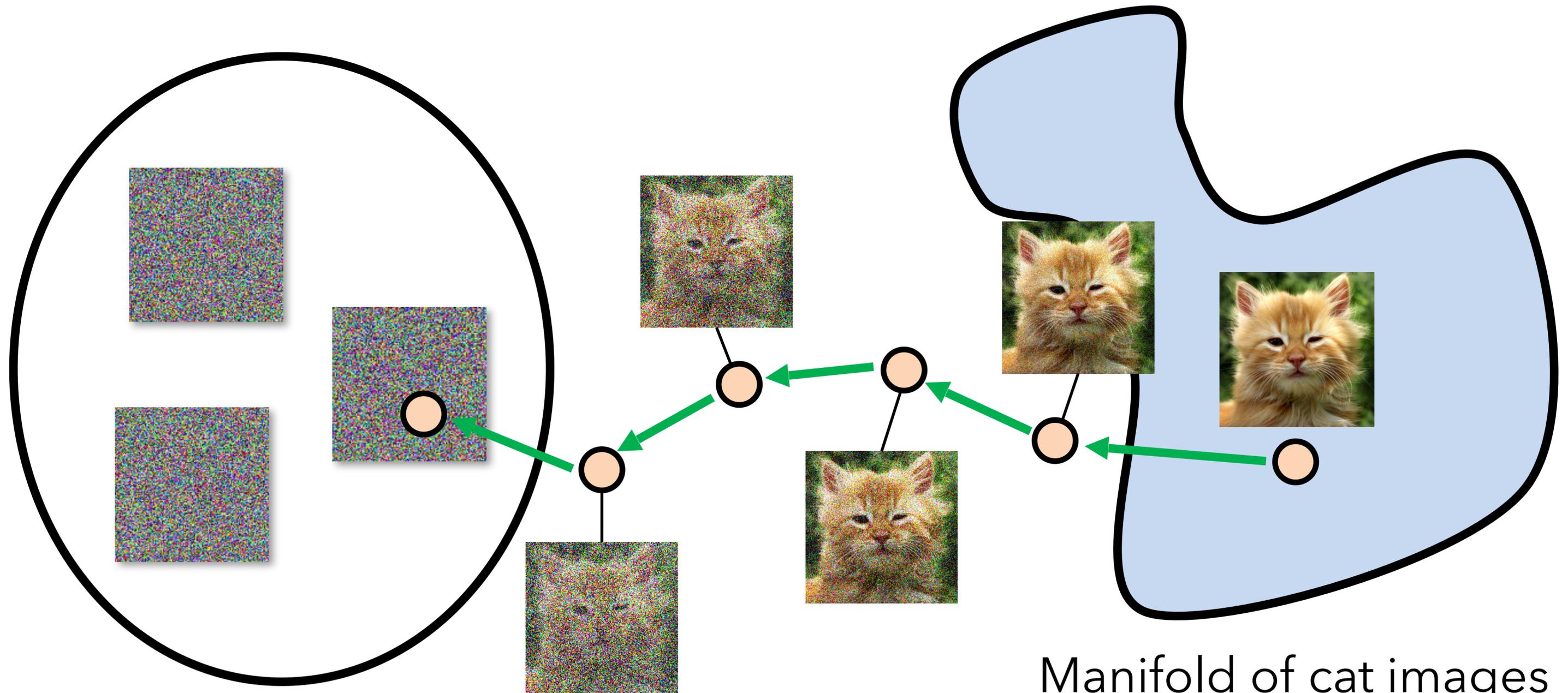


Random images

Diffusion

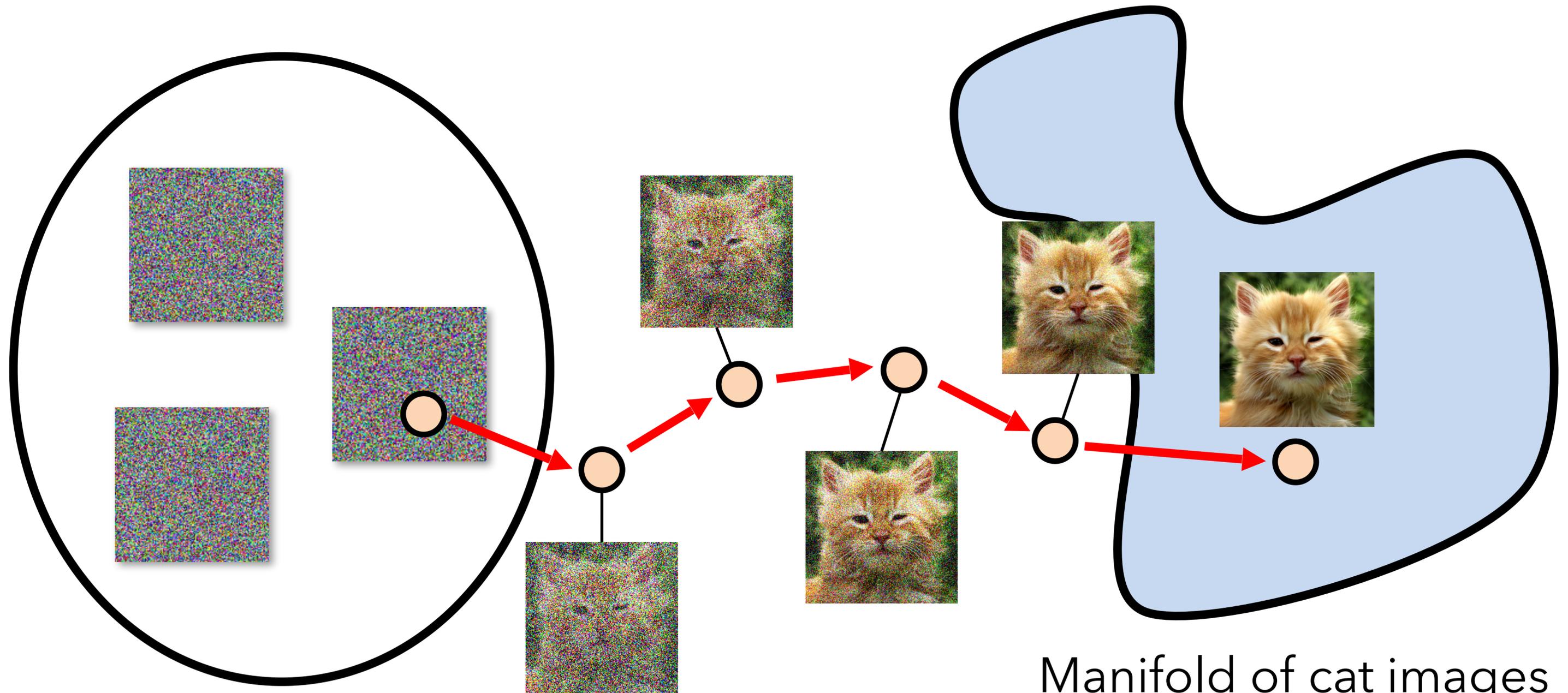
A blue arrow pointing from the random images to the manifold of cat images, indicating the diffusion process.

Manifold of cat images



Random images

Manifold of cat images



Random images

Manifold of cat images

Physical motivation for diffusion

Heat Diffusion



Physical motivation for diffusion

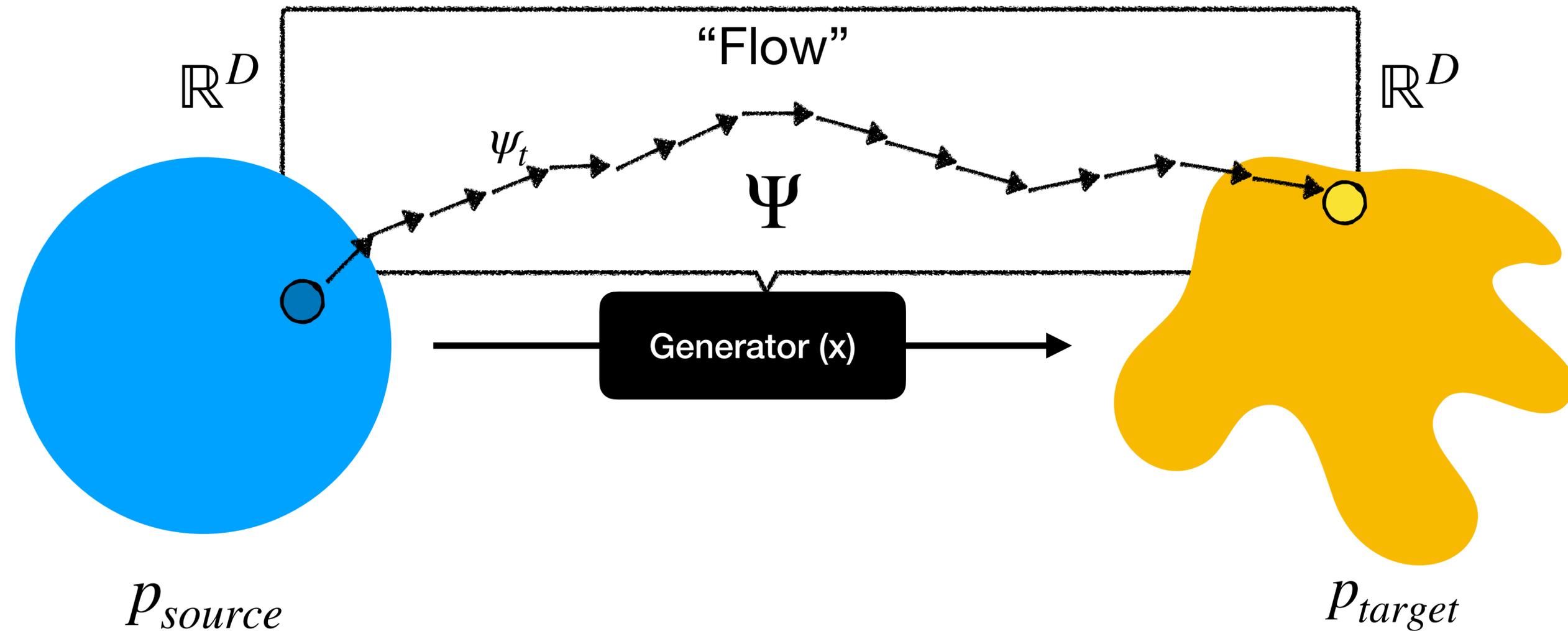
Reversing the process



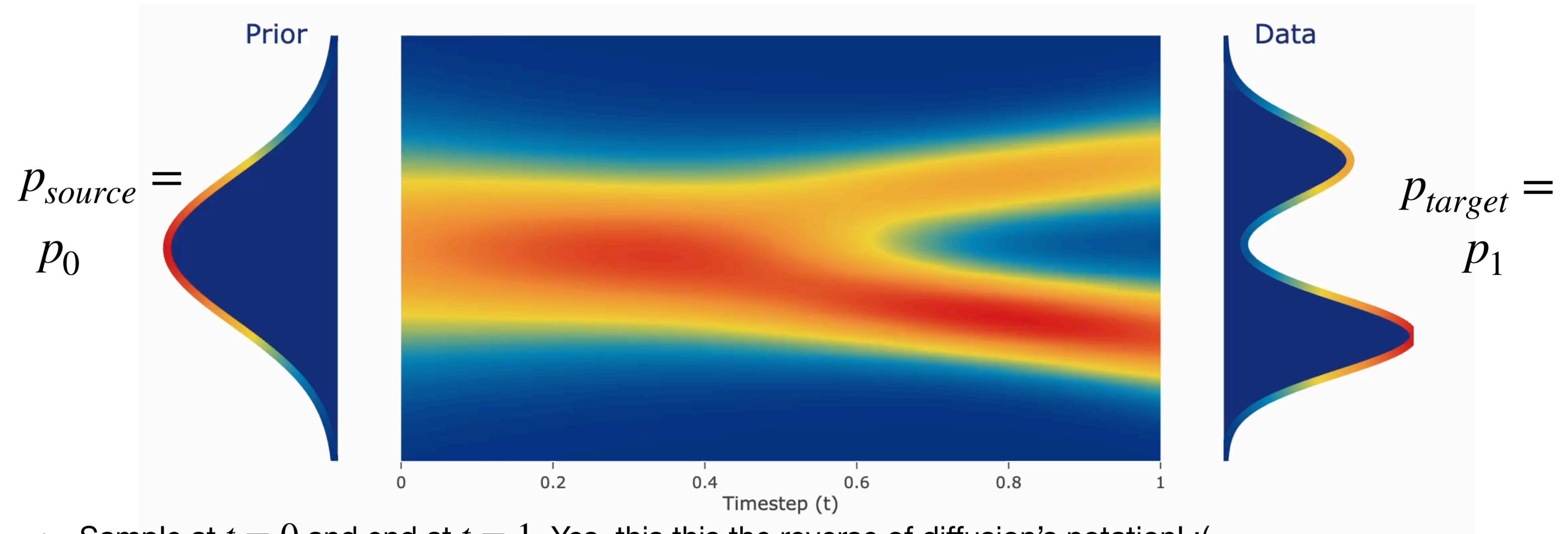
Flow matching

- Another perspective on generative modeling
- Widely used in state-of-art image and video generators
- Training and sampling process is quite similar to diffusion models.
 - In fact, many popular formulations are equivalent [Gao et al., “Diffusion Meets Flow Matching”, 2024].
- This formulation makes a number of things simpler, though!

Flow matching models

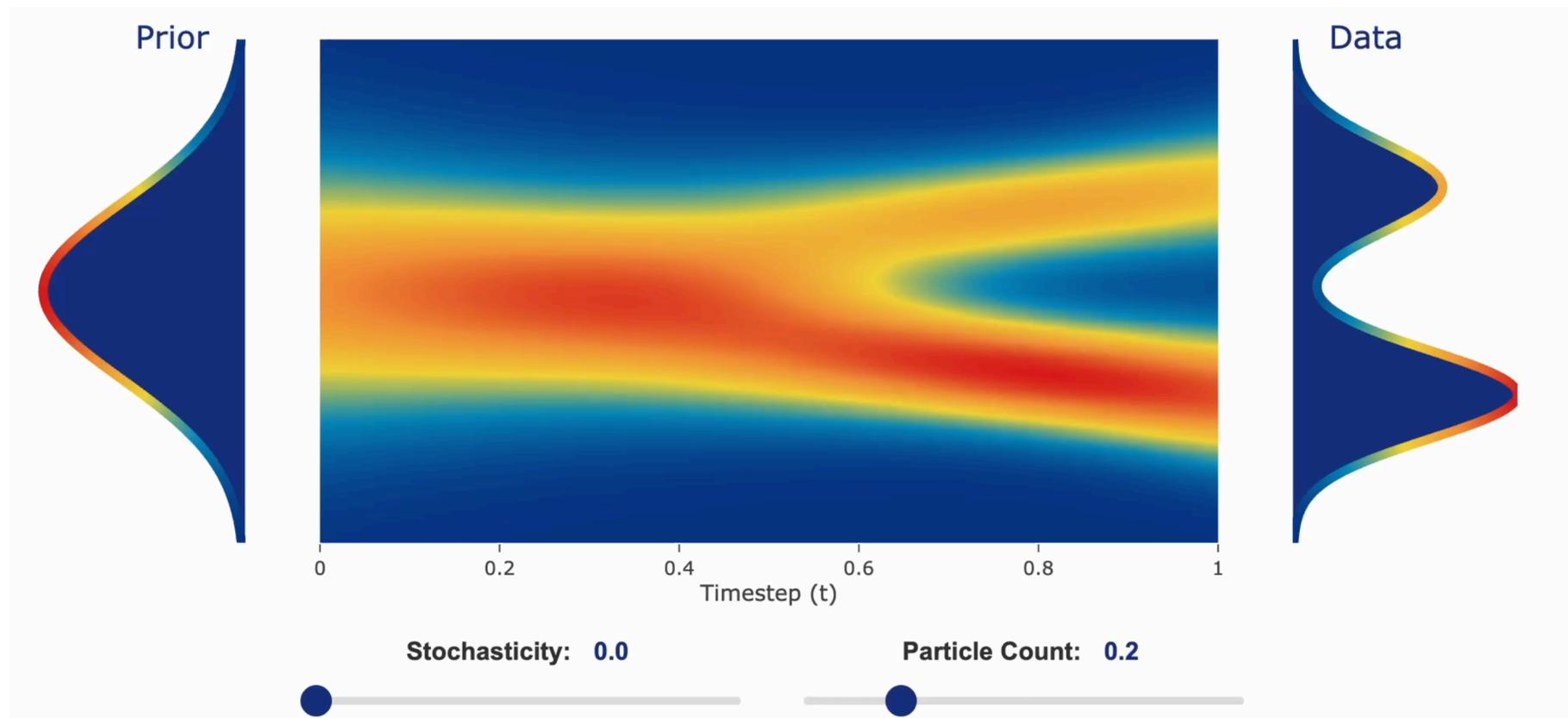


What is a Flow?



- Sample at $t = 0$ and end at $t = 1$. Yes, this is the reverse of diffusion's notation! :(
- It is a **velocity field**.
- It's like a river with some currents, every point defines how fast you move (velocity)
- You ride this river to go from one distribution to next

Riding the river = integration

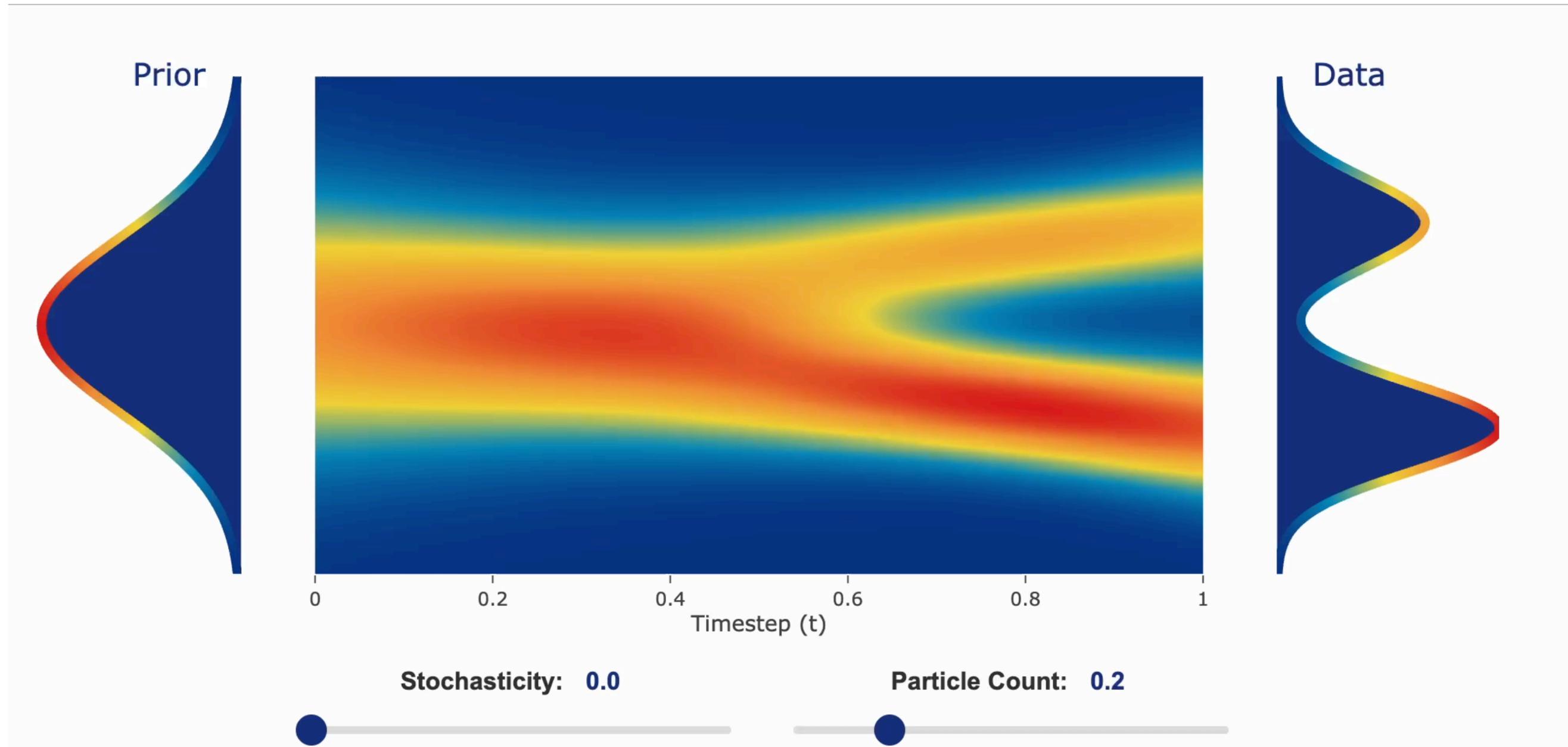


Simplest *Euler Integration*:

$$x_{t+\Delta t} = x_t + v_{\theta}(x_t, t)\Delta t$$

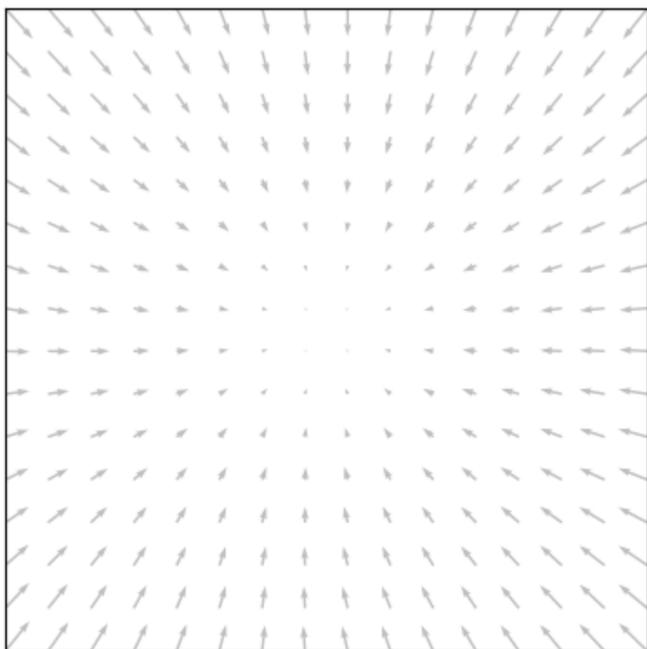
- Riding this river means you add little bits of velocity defined at each location
- This is called *integration*: solving the Ordinary Differential Equation (ODE) with initial state x_0 , through some differential parametrized by a network: $\frac{dx}{dt} = v_{\theta}(x, t)$
- You can add stochasticity when riding it, then it becomes a Stochastic Differential Equation (SDE)

How can we learn this flow?

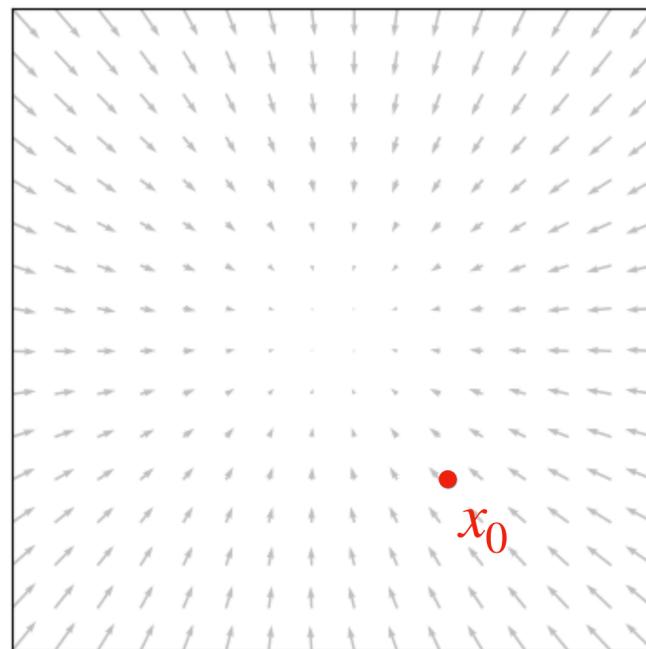


In 2D

$$v_t(x)$$



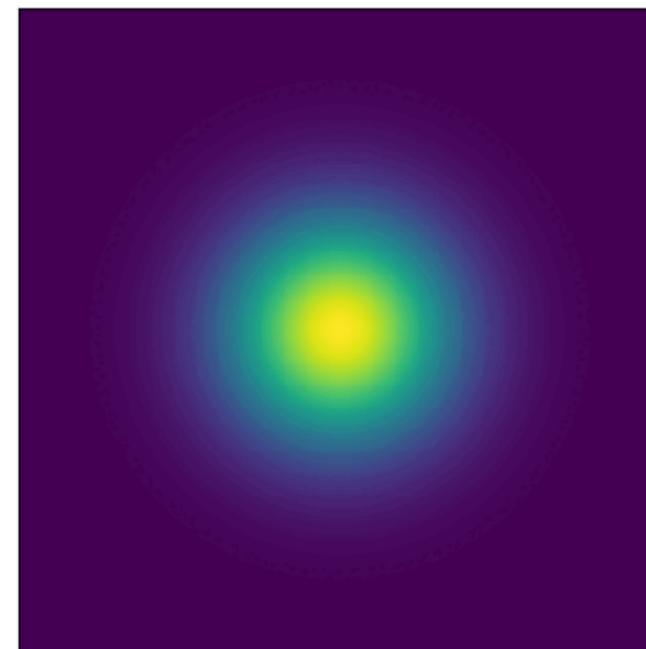
$$x_t = \Psi_t(x_0)$$



$$x_0 \sim p$$



$$x_t \sim p_t$$

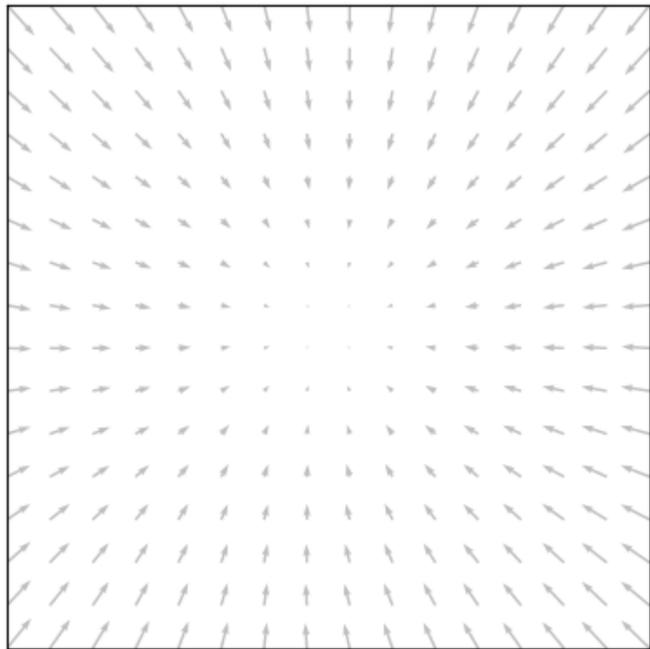


Flow ODE

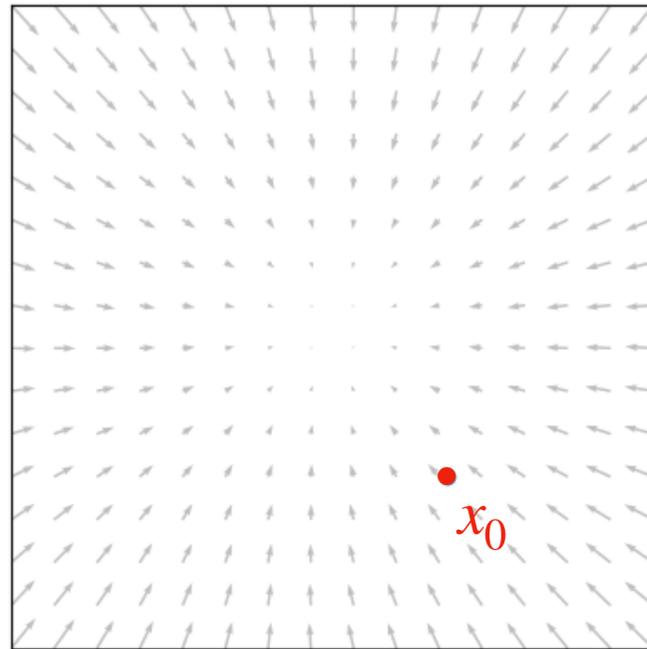
$$\dot{x}_t = v_t(x_t)$$

How can we train this?

$$v_t(x)$$



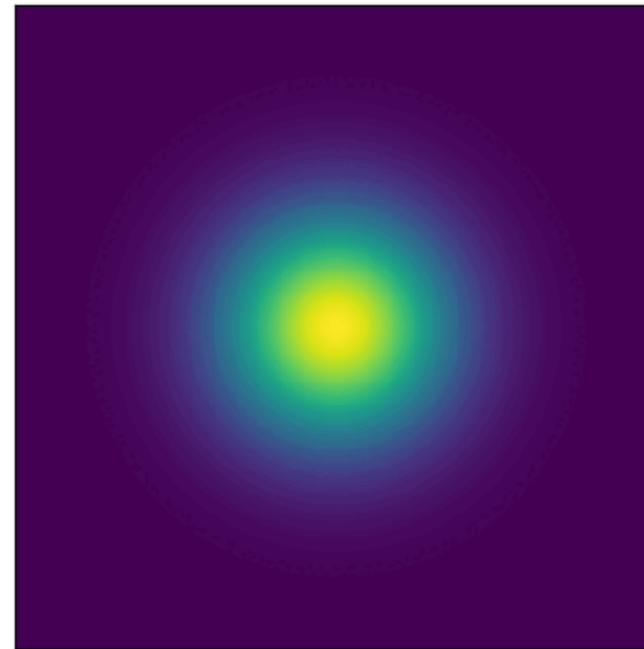
$$x_t = \Psi_t(x_0)$$



$$x_0 \sim p$$

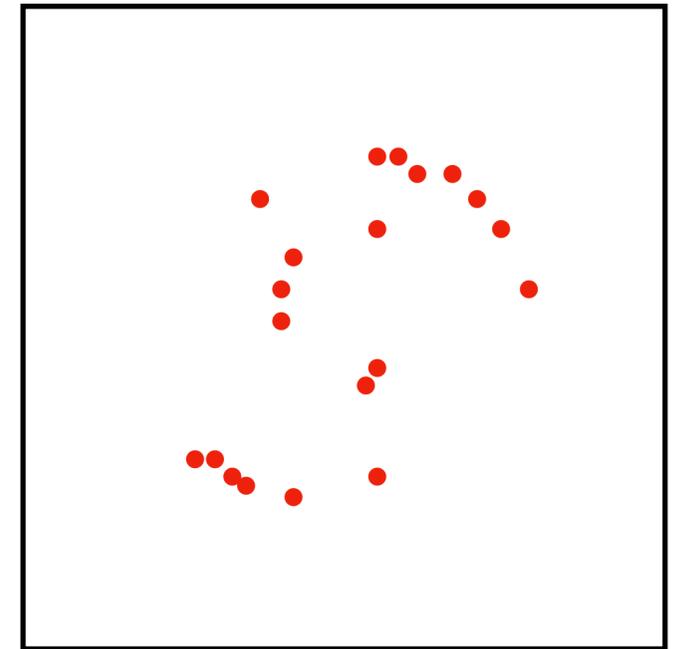
↓

$$x_t \sim p_t$$



Samples of p_1

$$q(x)$$

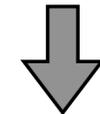


Flow ODE

$$\dot{x}_t = v_t(x_t)$$

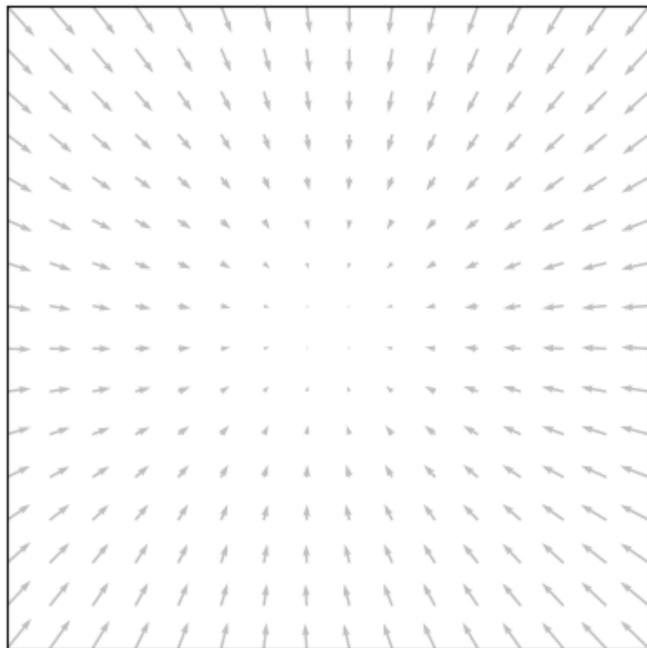
Caveat: Continuity equation

$$x_0 \sim p$$

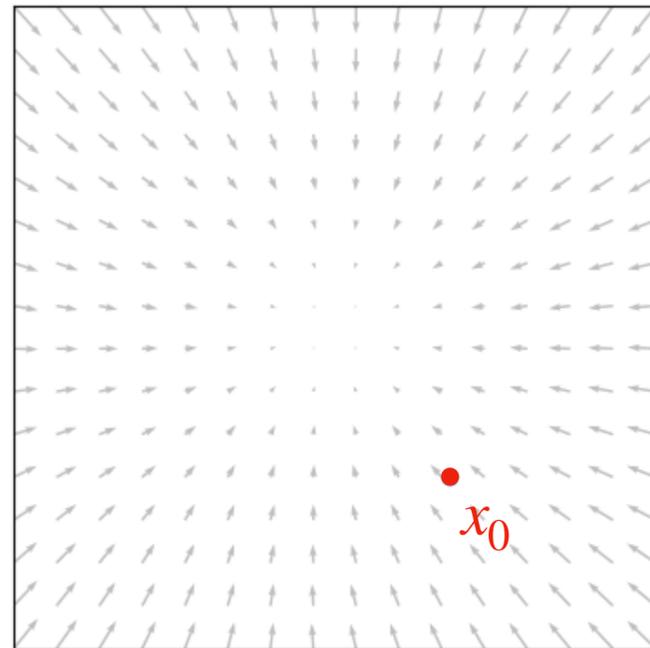


$$x \sim p$$

$$v_t(x)$$



$$x_t = \Psi_t(x_0)$$



Continuity Equation PDE (fixed x)

$$\frac{d}{dt} p_t(x) = -\operatorname{div}(p_t v_t)$$

where $\operatorname{div} v(x) = \sum_{i=1}^d \frac{\partial}{\partial x_i} v_i(x)$

Flow ODE

$$\dot{x}_t = v_t(x_t)$$

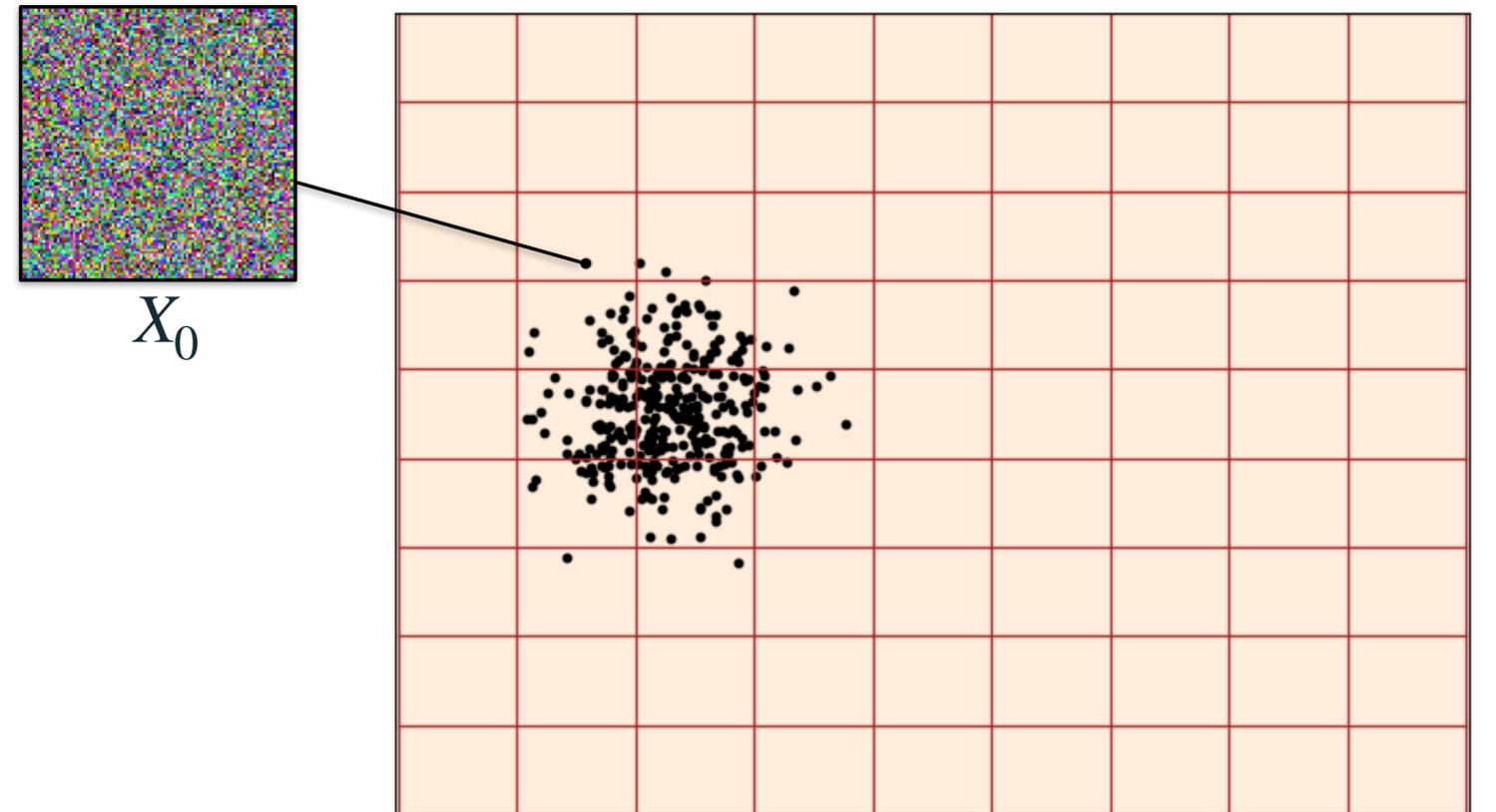
- For probability path to be valid, it needs to conserve probability mass.
- In the river, analogy you cannot add or remove water
- It has to come from somewhere and go somewhere

Treating the flow as a warping function

$$X_t = \psi_t(X_0), \quad t \in [0,1]$$

Warping

Source $X_0 \sim p$



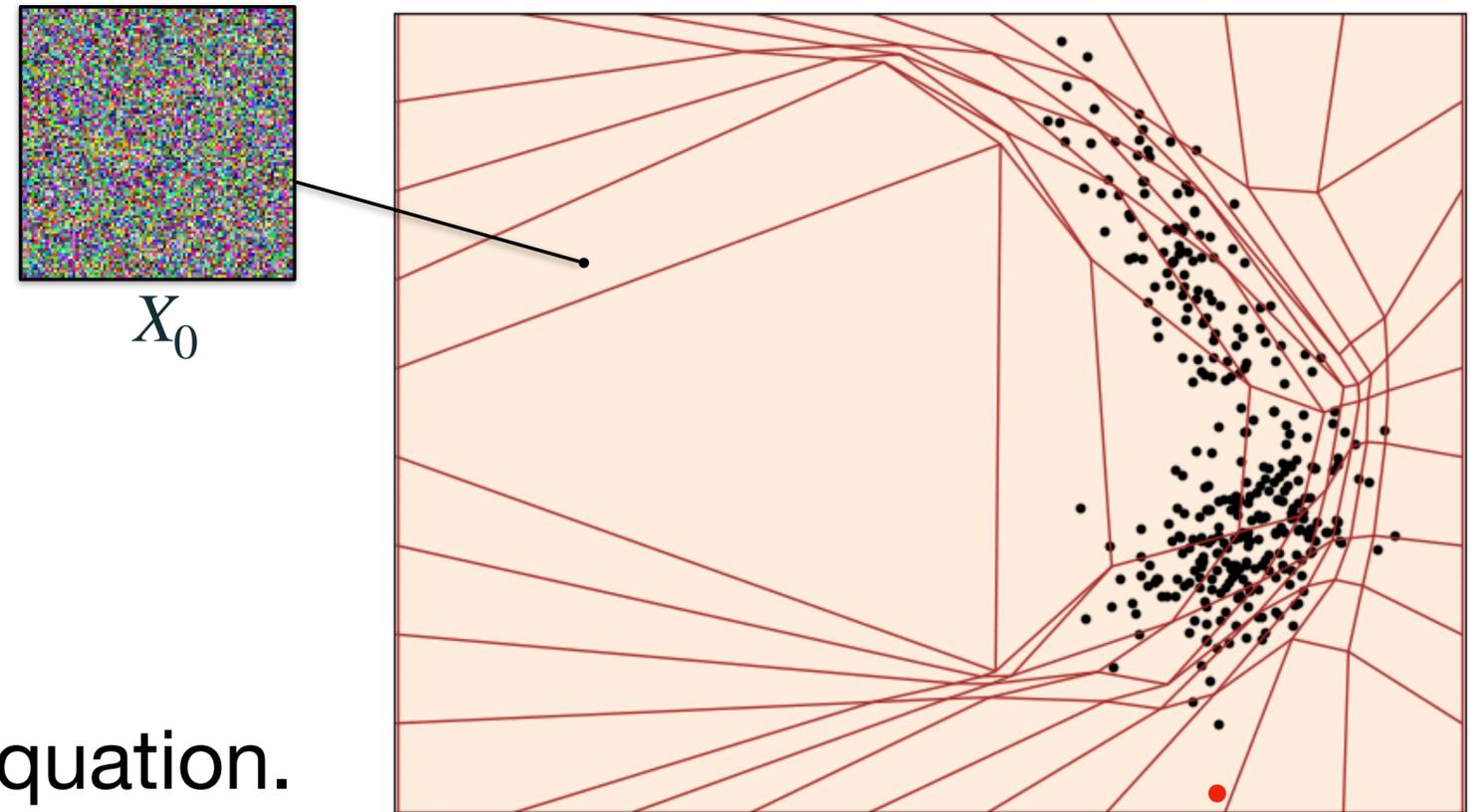
Early work trained flow by maximum likelihood

$$D_{\text{KL}}(q \parallel p_1) = - \mathbb{E}_{x \sim q} \log p_1(x) + c$$

$$X_t = \psi_t(X_0), \quad t \in [0, 1]$$

Warping

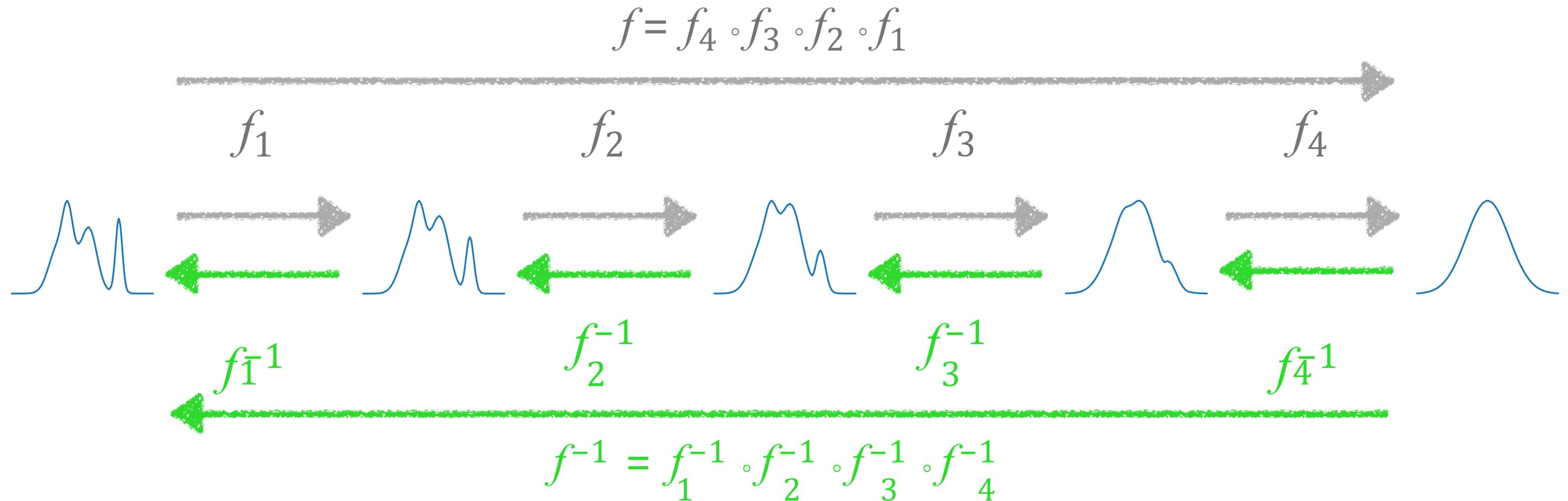
Source $X_0 \sim p$



- Chaining ψ_t needs to satisfy the continuity equation.
- This requires ODE integration *during training* with invertible neural networks.
- This is a *Continuous* Normalizing Flow (CNF) model!

$\log p_1(x) = ?$

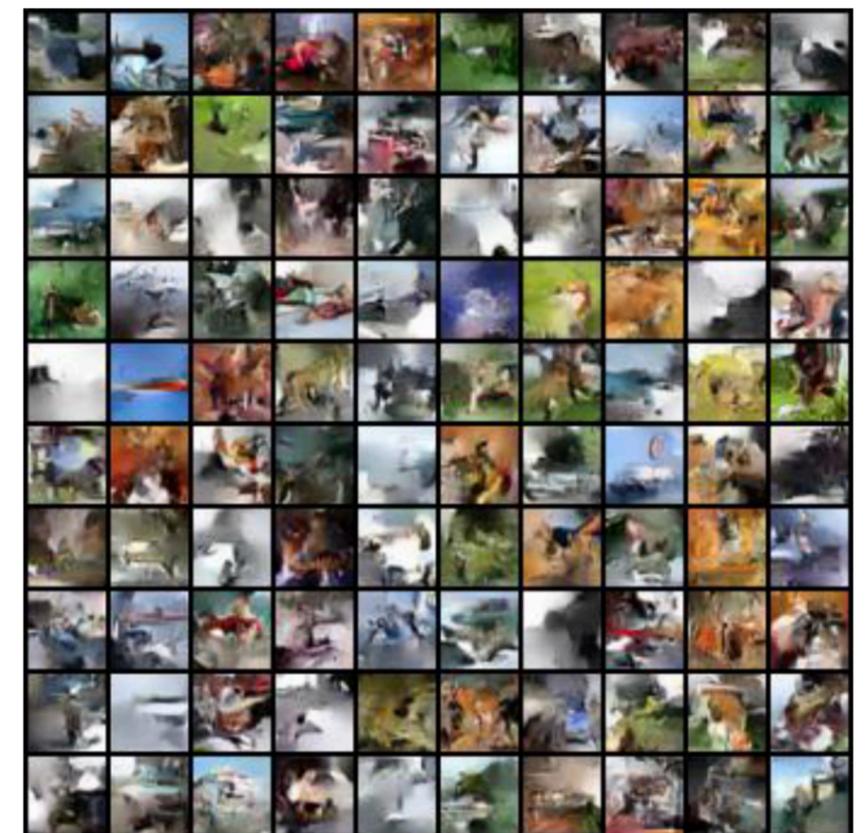
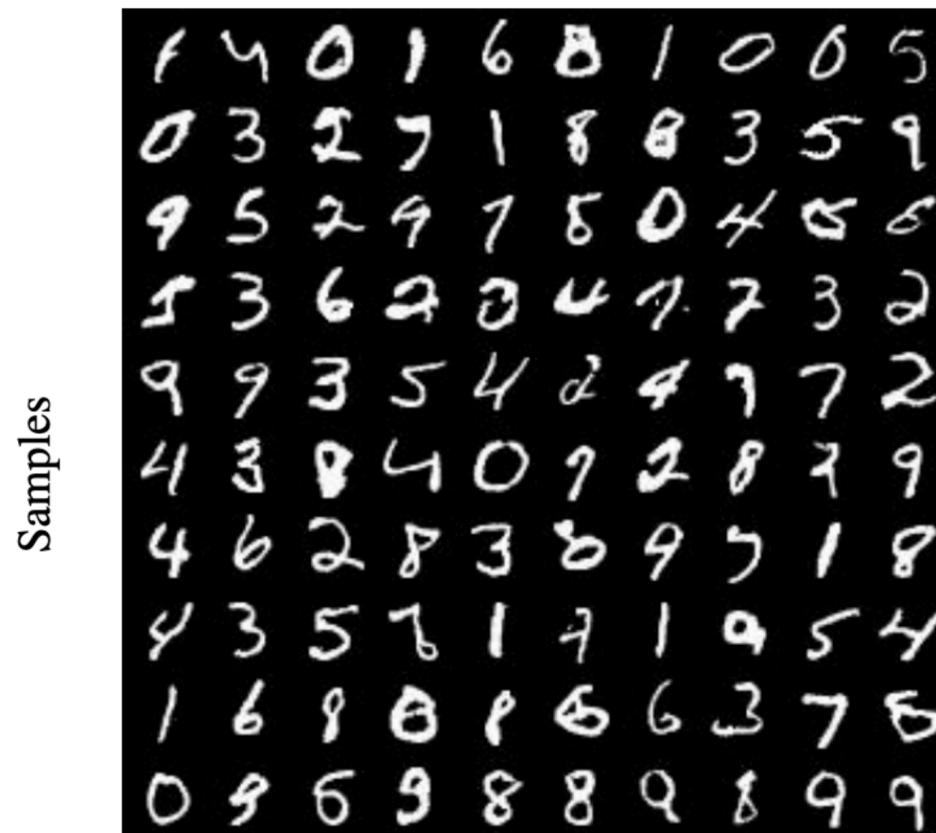
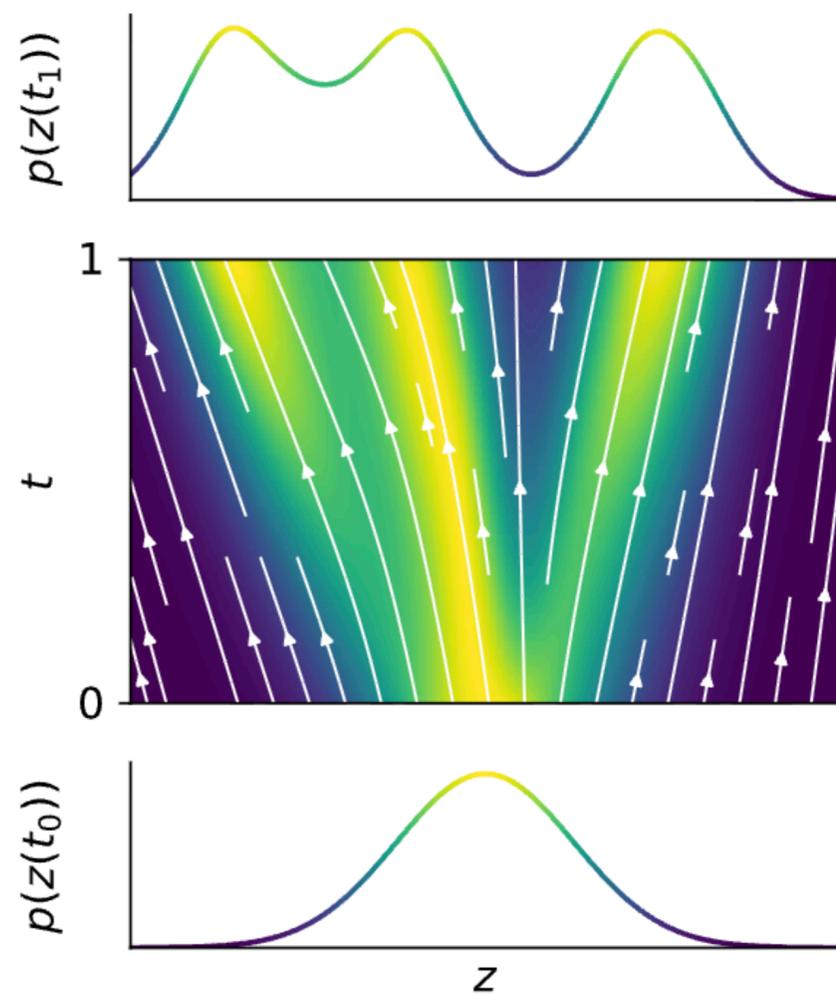
Recall: normalizing flow (with discrete layers)



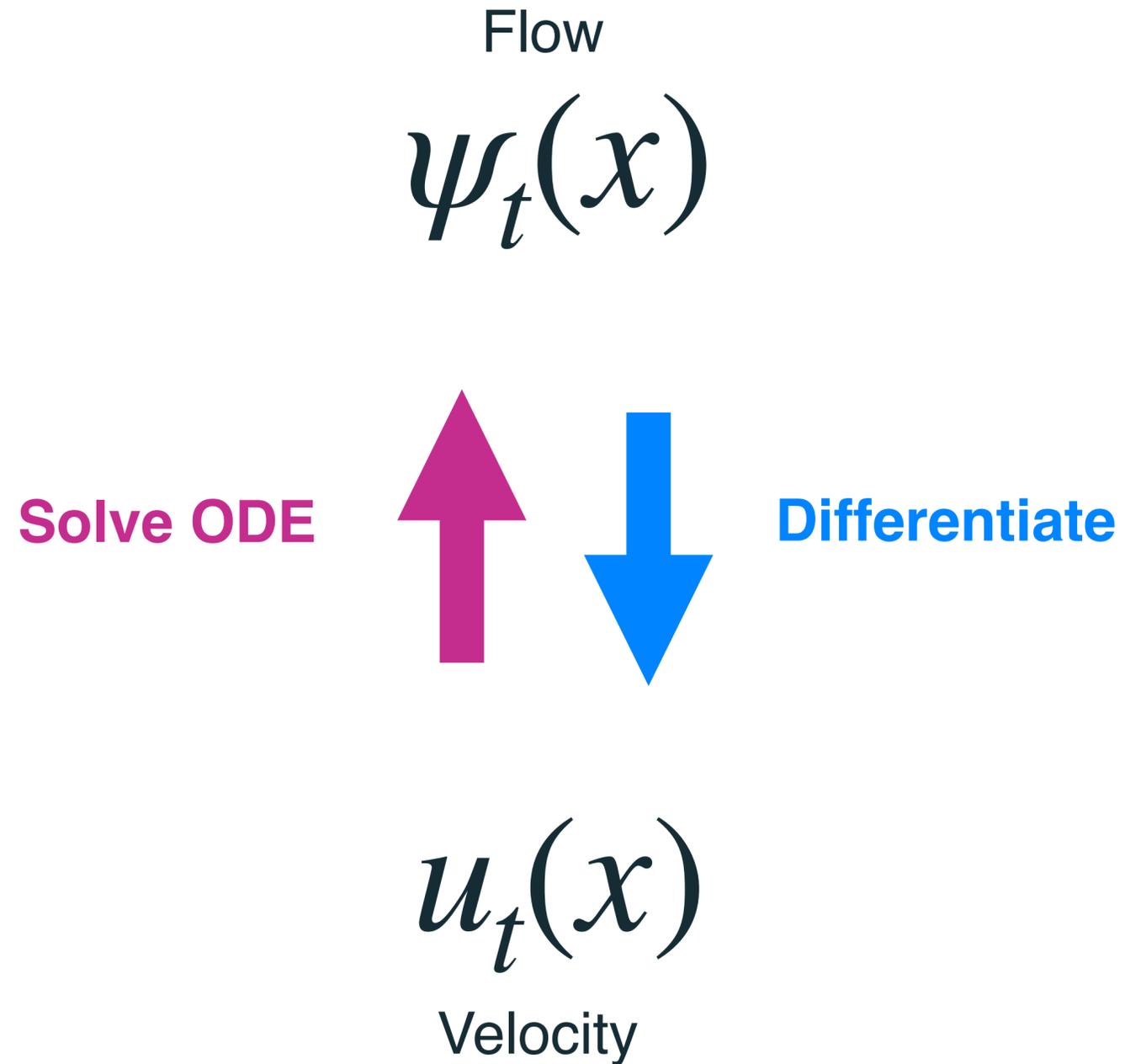
Early flow learning methods

- Treat like a Continuous Normalizing Flow (CNF) model (e.g., [Chen et al., “Neural ODE”, 2018])
- Tries to directly deal with this continuity equation constraint
- Very slow to train (need to integrate while training)
- Requires other constraints like invertibility of the ψ_t network.
- Due to these limitations, they are hard to train and require special architectures.

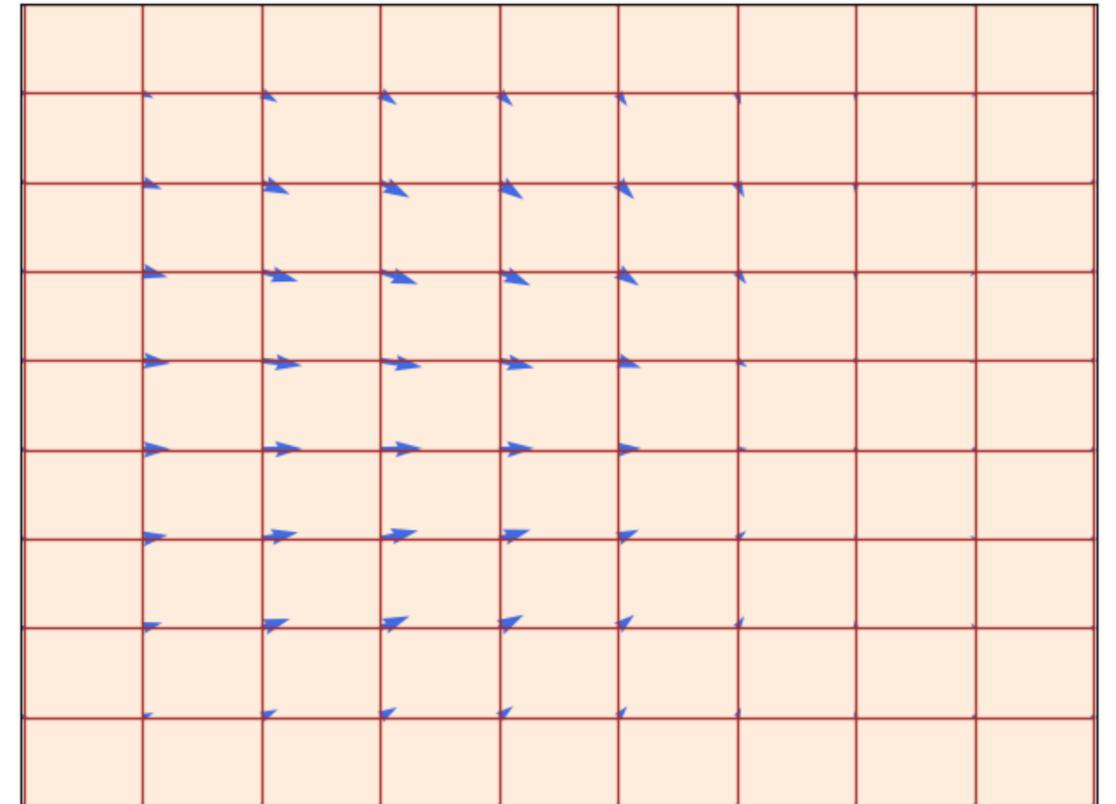
Continuous normalizing flow example



Instead, model flow with velocity



$$\frac{d}{dt}\psi_t(x) = u_t(\psi_t(x))$$

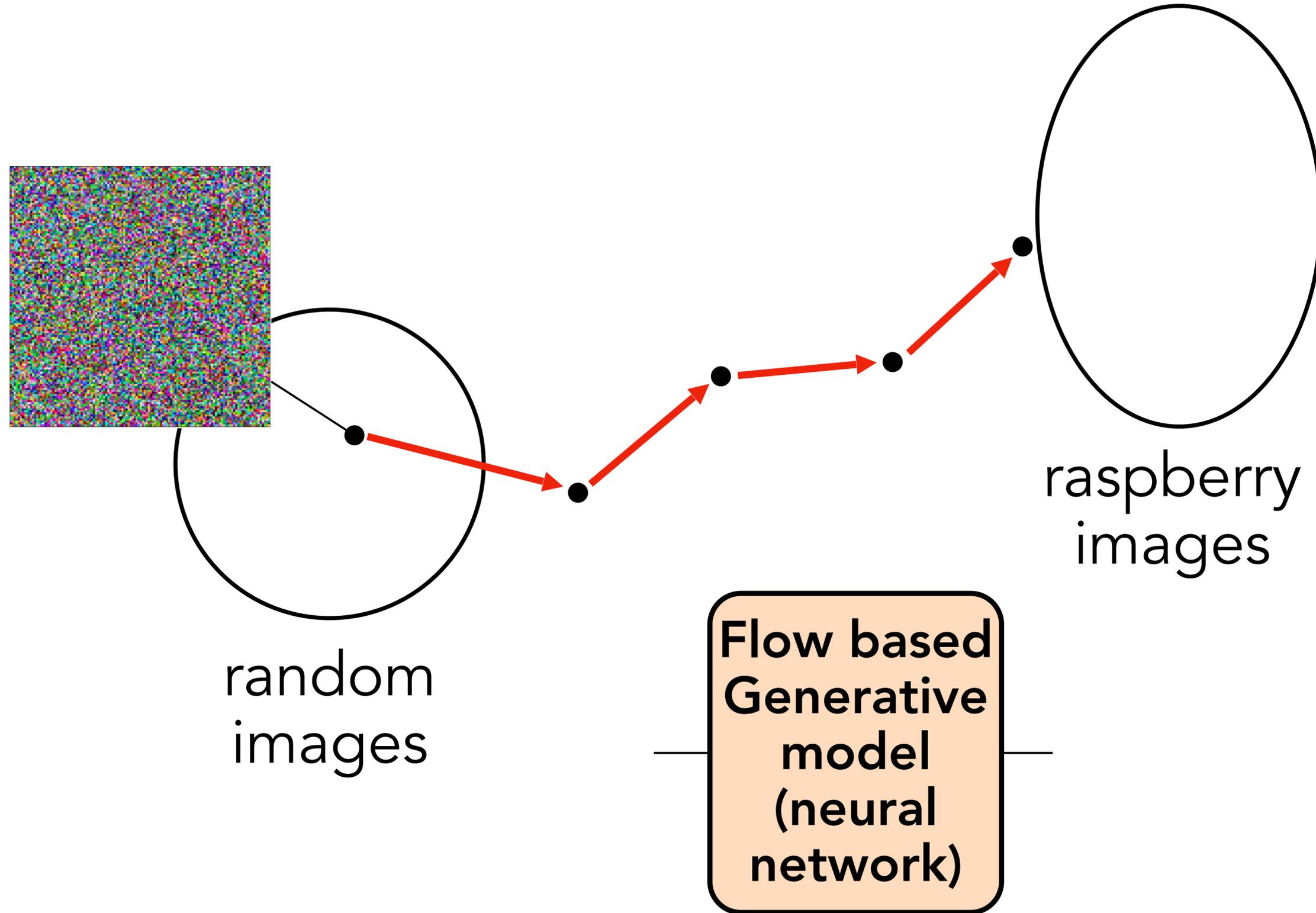


Flow Matching

- Directly learn the velocity field [Lipman et al. 2022]
- Don't have to worry about the continuity equation because velocity fields can't add / subtract mass. You only re-distribute
- Continuity equation satisfied by construction.
- Can learn velocity fields for each data sample (conditional velocity field), and everything works fine.

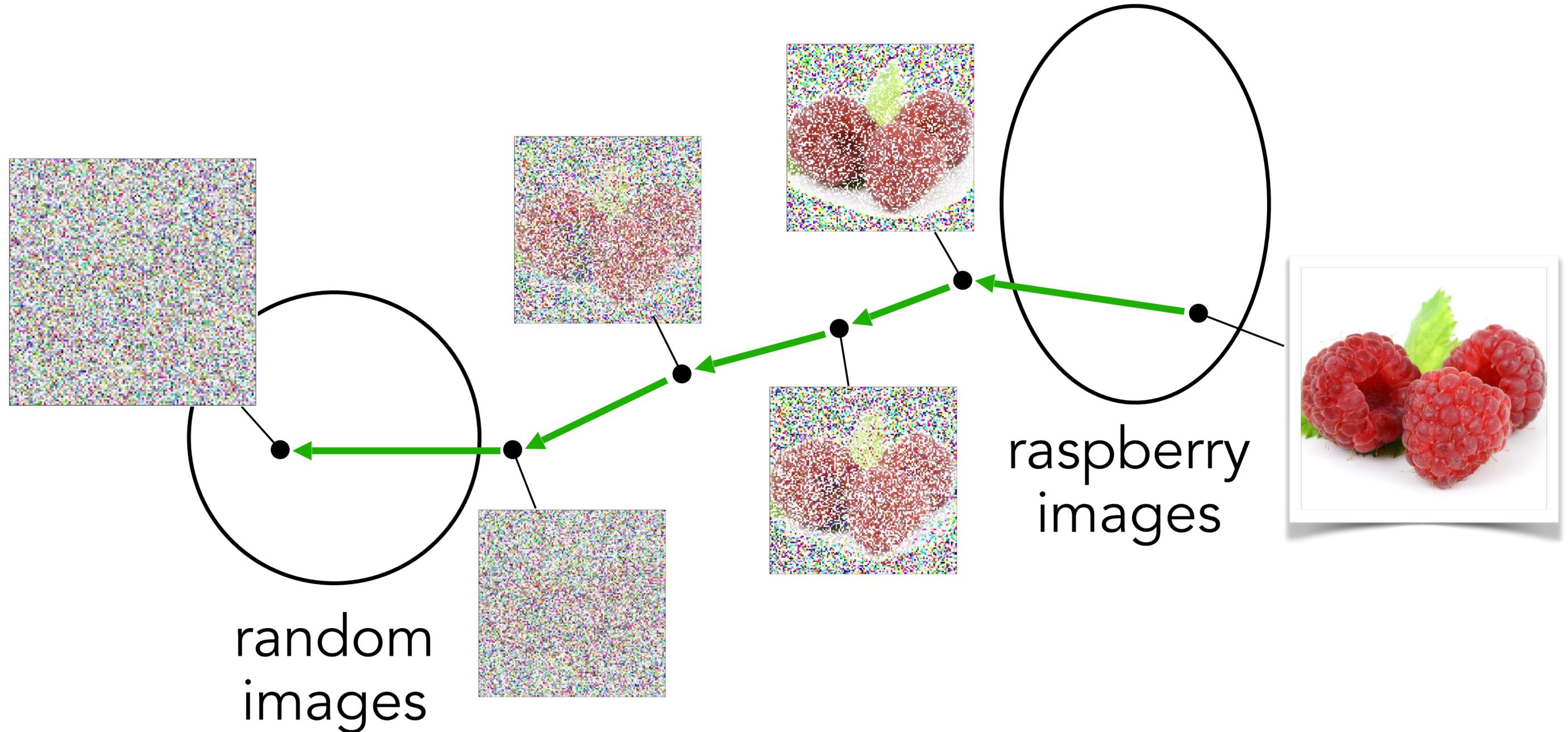
How do you train the flow?

Just like diffusion:



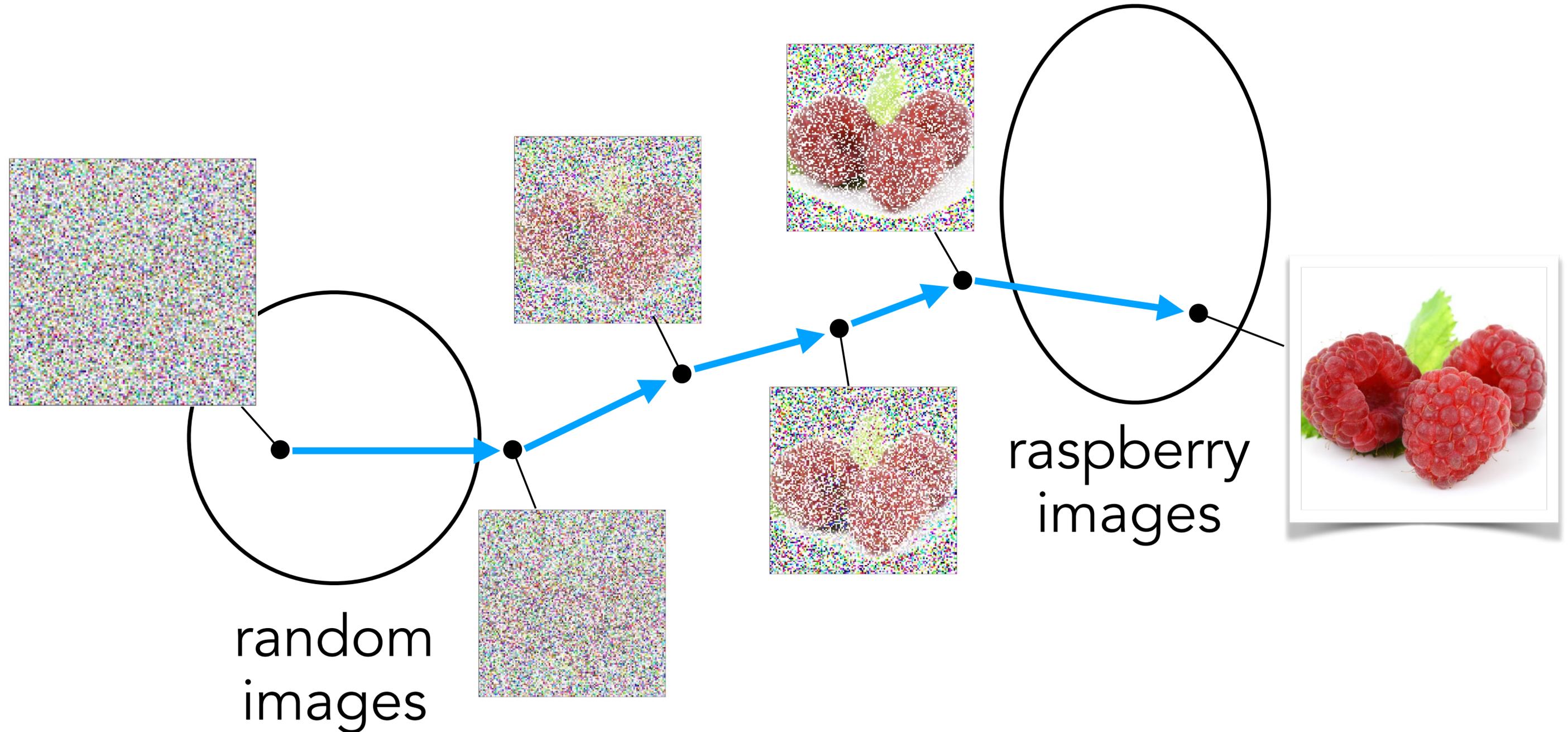
Training

1. Take real data and corrupt it

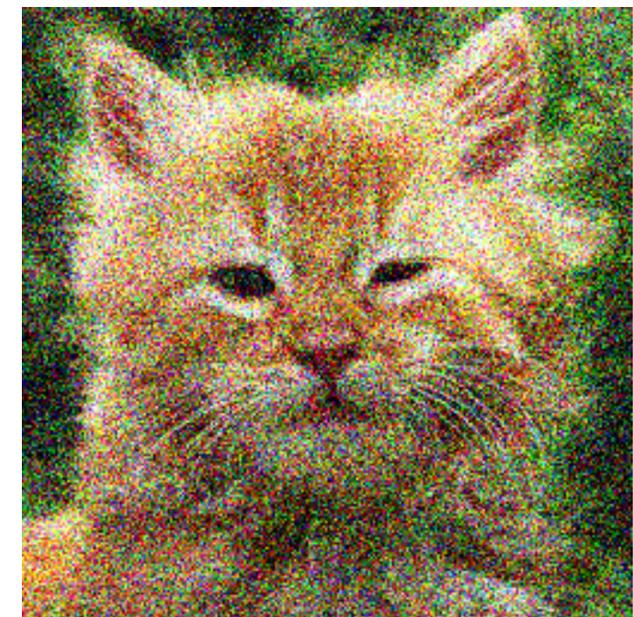
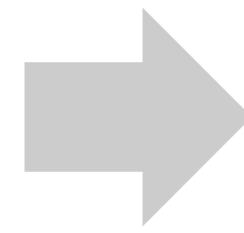
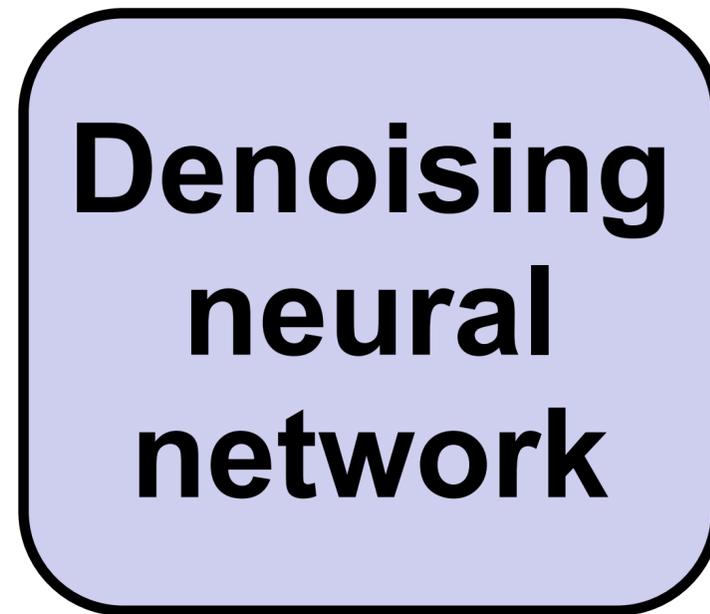
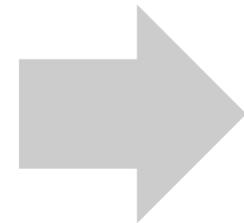
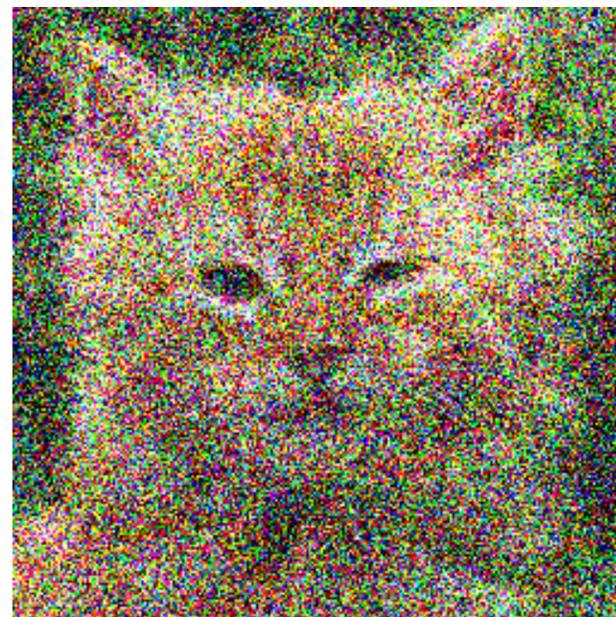


Training

1. Take real data and corrupt it
2. Learn to undo the process!



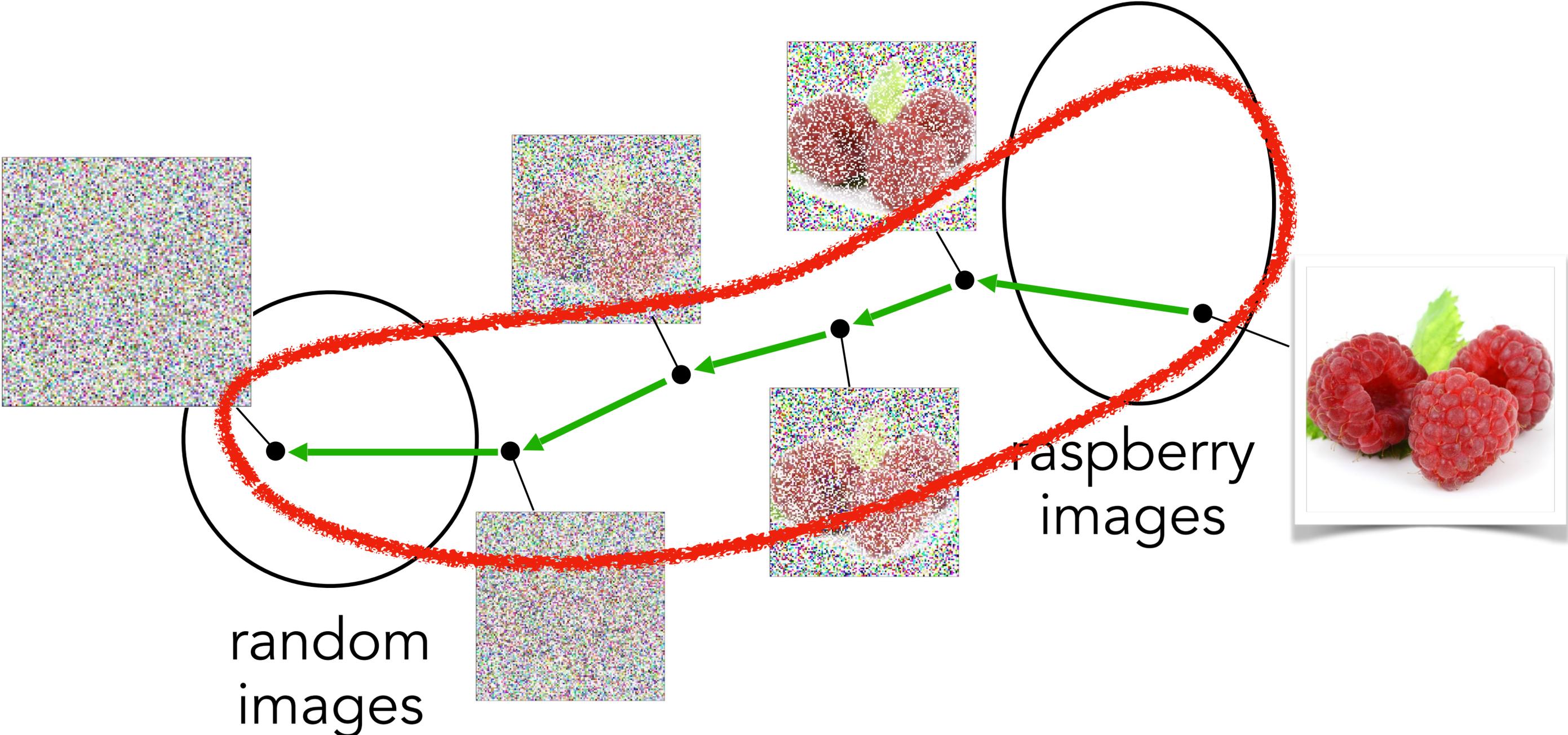
Denoising with a neural network



Just like in diffusion, can be a U-Net or
transformer (DiT)

How do we pick the intermediate path?

How to generate the path?



What is the path?

- How to add noise? What kind of noise? What schedule do we use?
- This is a bit complicated in diffusion. Noise process needed to have a specific form to make the math work.
- Can we make it simpler?
- ▶ Flow matching says that you can (essentially) add noise however you like!

How do we construct x_t ?

TLDR: Sample noise, add it, then reconstruct the data

- Flow matching is quite general! For now, just consider time-dependent weighted combinations.
- Flow matching says we can choose any weighted combination, as long we start from a sample in the source (e.g., Gaussian) and end with a sample in the target distribution (image).

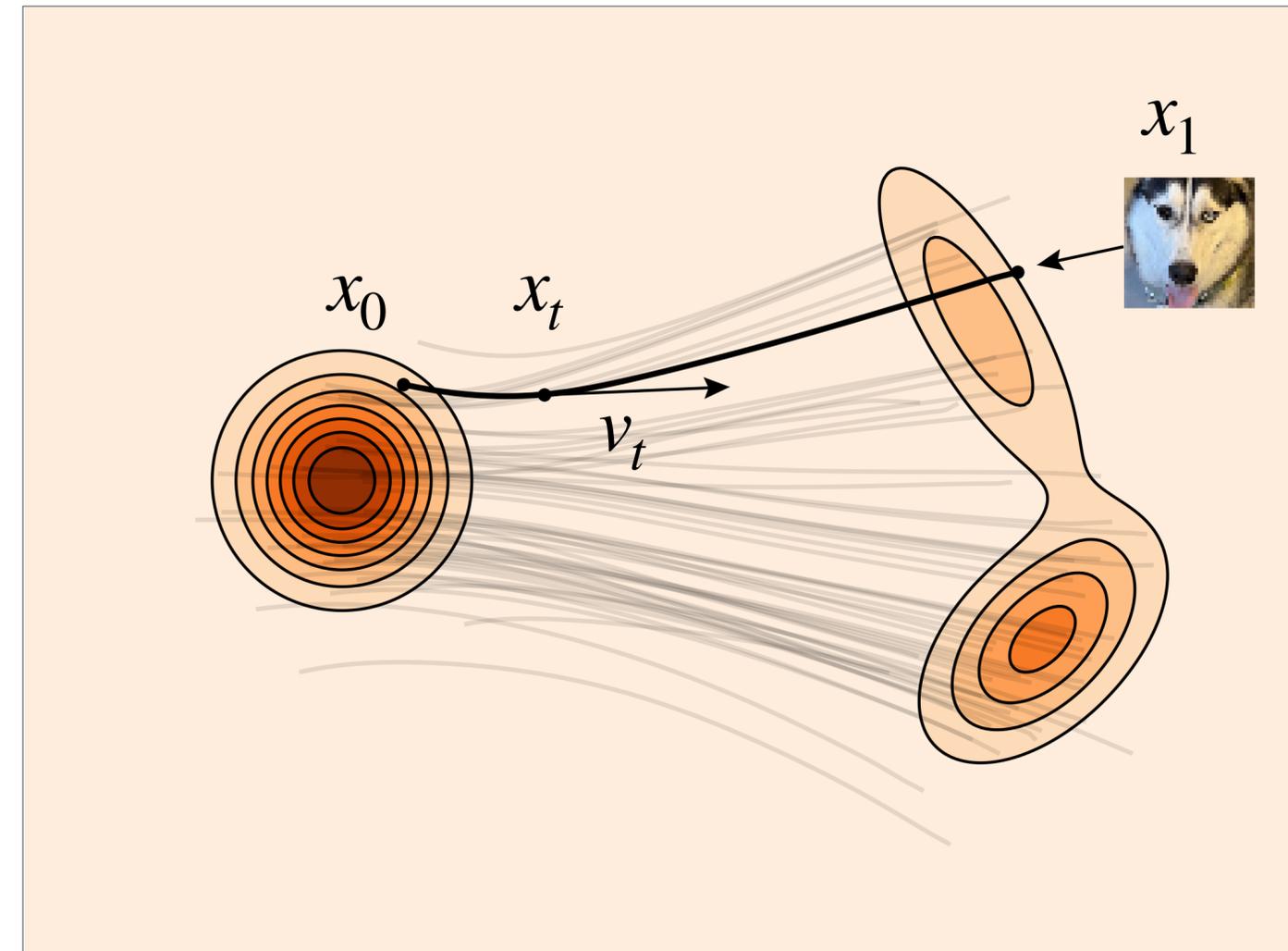
$$x_t = \alpha_t x_0 + \sigma_t x_1$$

$$x_0 \sim p_0(x)$$

$$x_1 \sim p_1(x)$$

Flow training

- For each image x_1
 - Sample some noise x_0
 - Combine them however you'd like to get x_t
 - Now learn to predict the velocity at x_t
 - What is the velocity? It depends on how you got x_t .

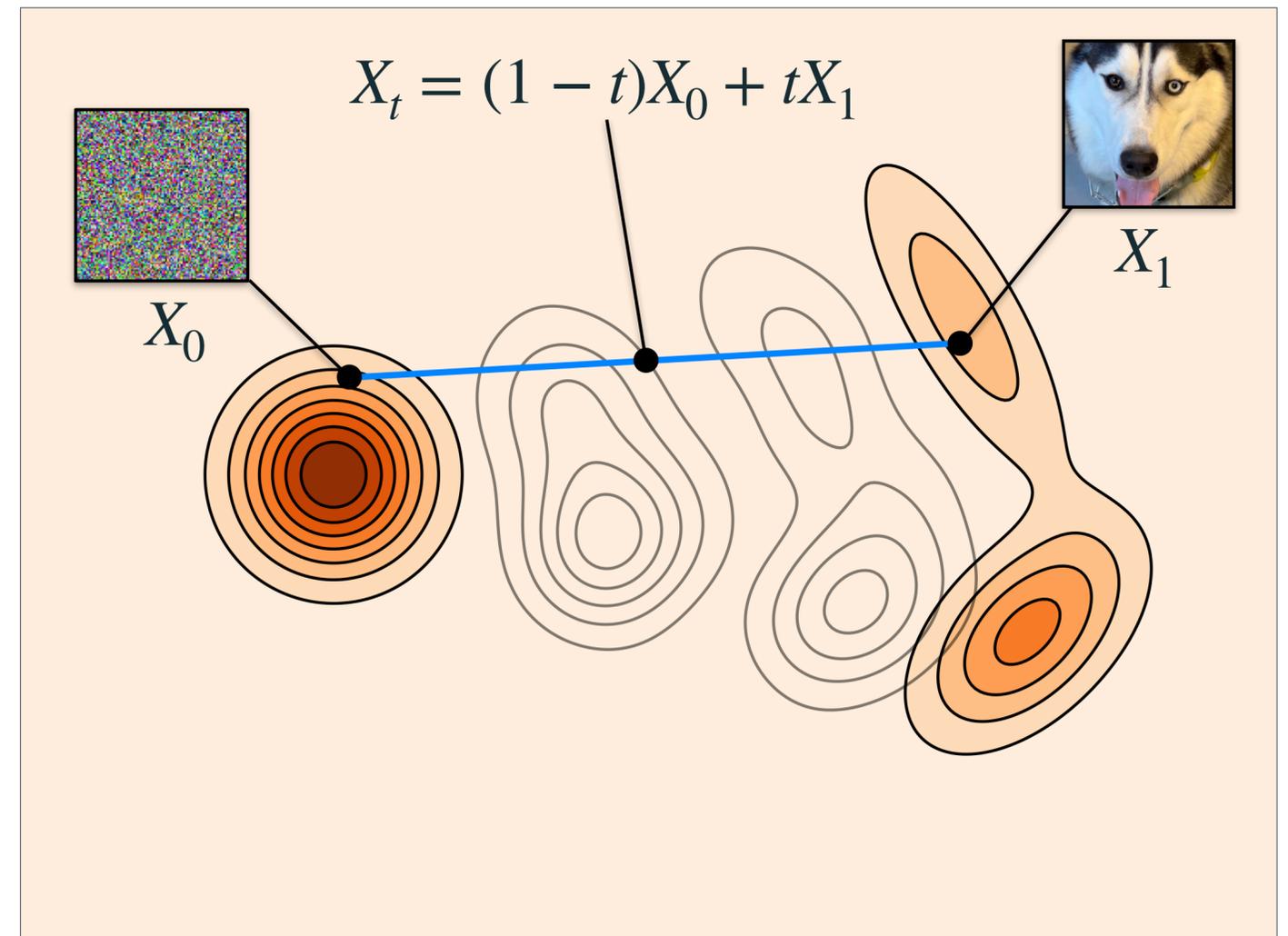


A Very Simple Way

Linear interpolation!

$$x_t = \underbrace{\alpha_t}_{(1-t)} x_0 + \underbrace{\sigma_t}_{t} x_1$$

$$x_t = (1-t)x_0 + tx_1$$



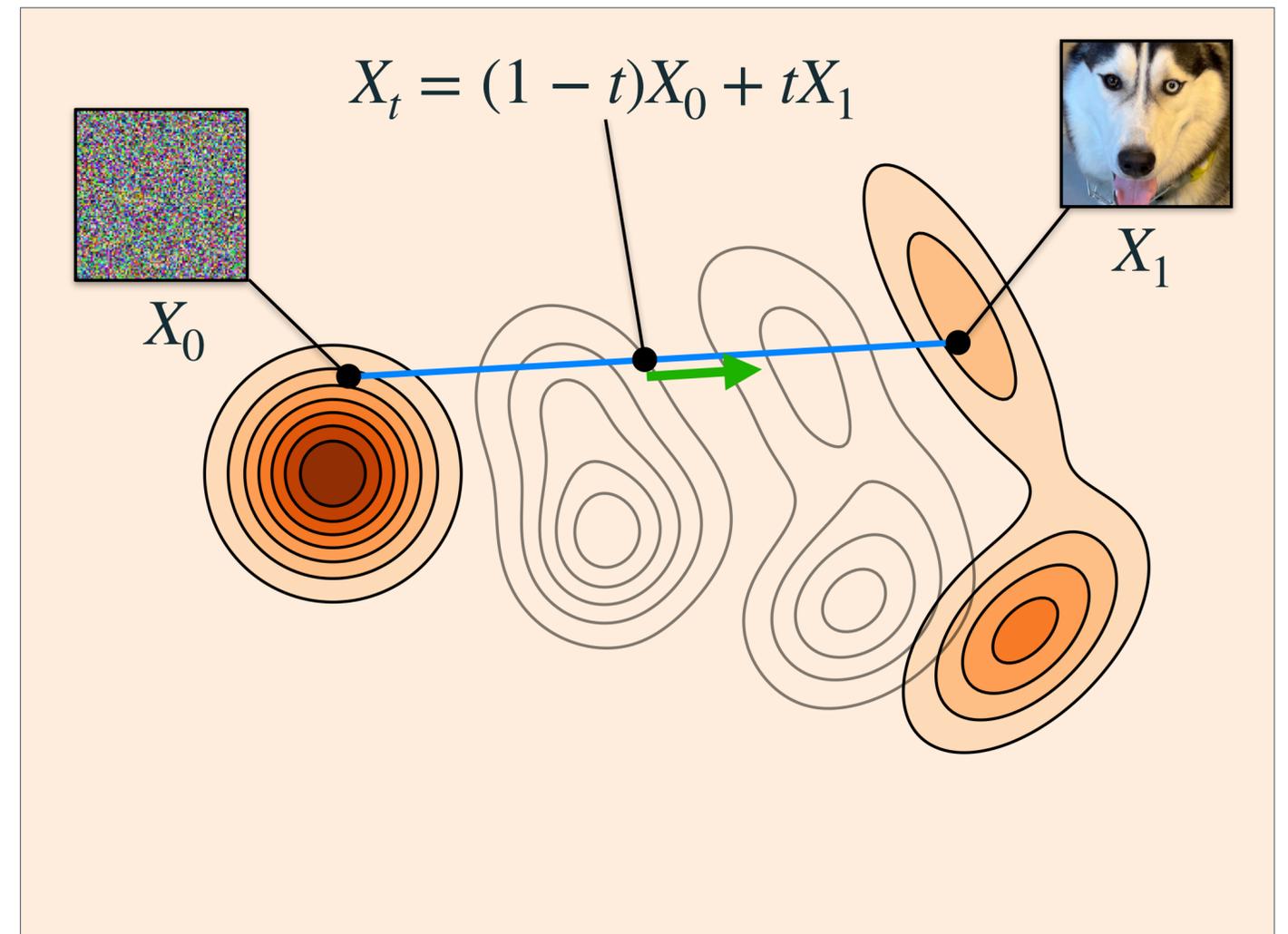
What is the velocity supervision?

$$x_t = \alpha_t x_0 + \sigma_t x_1$$

$$x_t = (1 - t)x_0 + tx_1$$

$$\begin{aligned} \frac{dx_t}{dt} &= -x_0 + x_1 \\ &= x_1 - x_0 \end{aligned}$$

$$\mathbb{E}_{t, X_0, X_1} \left\| u_t^\theta(x_t) - (X_1 - X_0) \right\|^2$$



*Conditioned on a single sample

Inside a Training Loop

Flow Matching

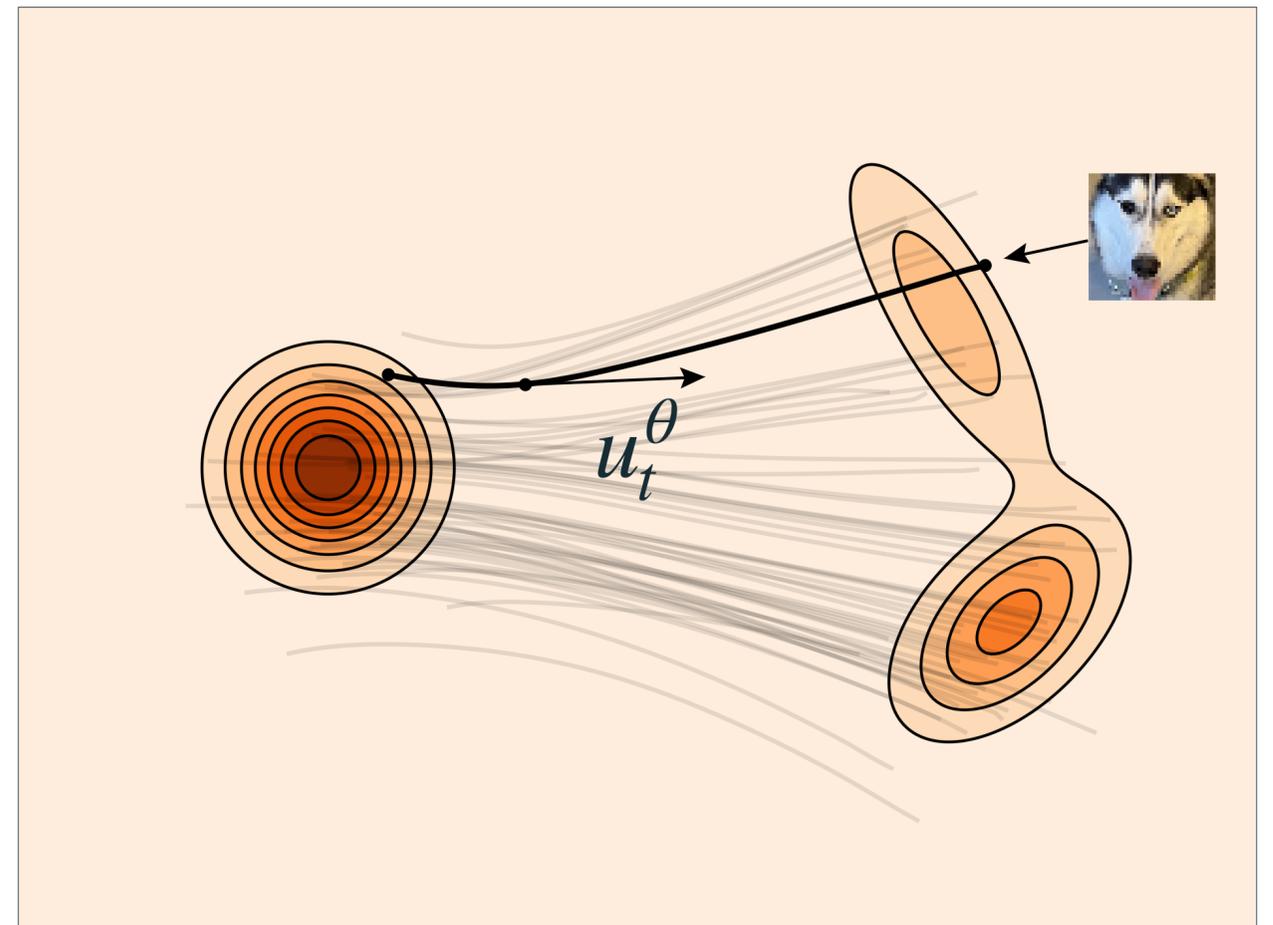
```
x = next(dataset)
t = torch.rand(1) # Sample timestep (0,1)
noise = torch.randn_like(x) # Sample noise
x_t = (1-t) * noise + (t) * x # Get noisy x_t

flow_pred = model(x_t, t) # Predict noise in x_t
flow_gt = x - noise # ground truth flow (w/ linear sched)
loss = F.mse_loss(flow_pred, flow_gt) # Update model
loss.backward()
optimizer.step()
```

Test-time sampling

- Just take a small step in the velocity
- Use any ODE Solver, i.e. integration you like, e.g., Euler integration:

$$x_{t+\Delta t} = x_t + \Delta t \cdot \left. \frac{dx}{dt} \right|_{x_t, t}$$



Sample
from $X_0 \sim p$

Inside Sampling Loop

```
velocity = model(x_t, t) # Predict noise in x_t  
x_t = x_t + dt * velocity # Step in velocity
```

Training: Model parameterization

- You can make your network output undo the noise in many different ways, predicting x_1 , v , noise, or flow

$$v_t = \alpha_t x_1 - \sigma_t x_0 \quad \begin{aligned} u_t &= x_1 - x_0 = \epsilon - x_0 \\ u_t &= x_t - x_0 \end{aligned}$$

- These are all equivalent because of the linear relationship with x_t . You can derive all of these as long as you know one of them

$$x_t = \alpha_t x_0 + \sigma_t x_1$$

- For example:

$$\mathbb{E}[(\hat{\mathbf{x}}_0 - \mathbf{x}_0)^2] = \mathbb{E} \left[\left(\frac{\mathbf{x}_t - \sigma(t)\hat{\epsilon}}{\alpha(t)} - \frac{\mathbf{x}_t - \sigma(t)\epsilon}{\alpha(t)} \right)^2 \right] = \mathbb{E} \left[\frac{\sigma(t)^2}{\alpha(t)^2} (\hat{\epsilon} - \epsilon)^2 \right].$$

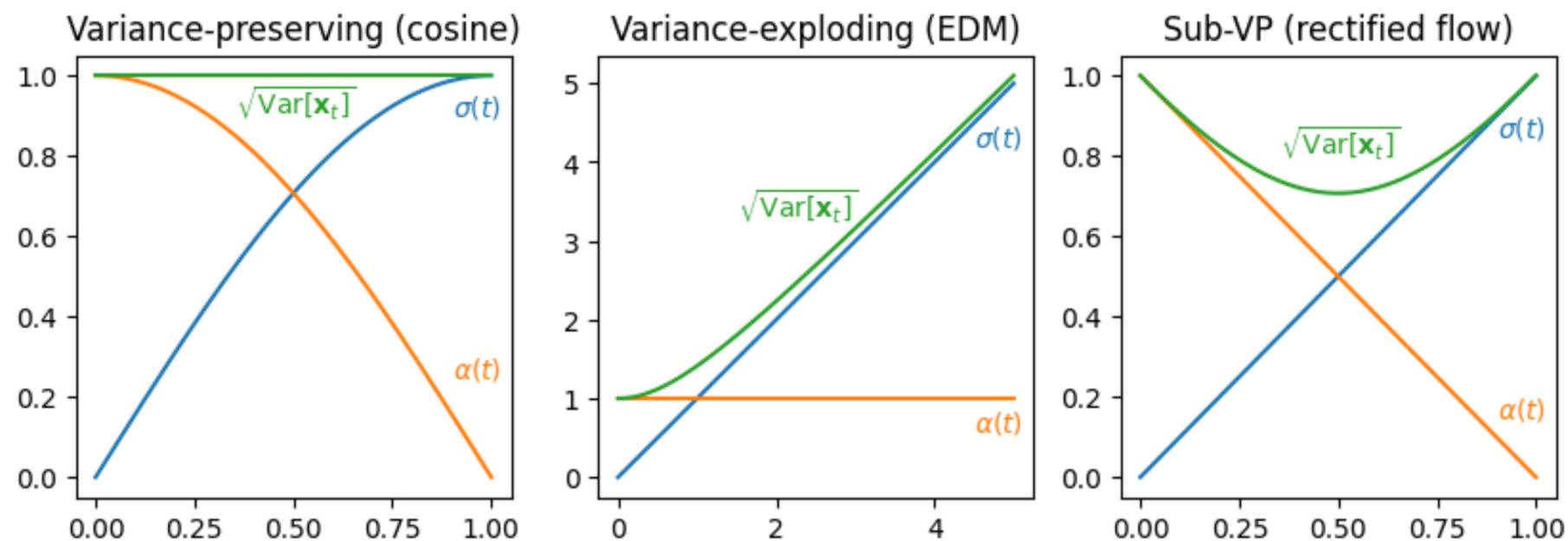
Weighting schemes

Different choices for: $\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \mathbf{x}_1$

- “Variance preserving”: DDPM (the main diffusion formulation you’ve seen already):

- i.e., $\mathbf{x}_t = \sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$

- Linear interpolation (Flow Matching variant that we’ve discussed so far, which is also known as Rectified Flow)



Training: Flow Matching vs. Diffusion

Algorithm 1: Flow Matching training.

Input : dataset q , noise p

Initialize v^θ

while *not converged* **do**

$t \sim \mathcal{U}([0, 1])$ ▷ sample time

$x_1 \sim q(x_1)$ ▷ sample data

$x_0 \sim p(x_0)$ ▷ sample noise

$x_t = \Psi_t(x_0|x_1)$ ▷ conditional flow

Gradient step with $\nabla_\theta \|v_t^\theta(x_t) - \dot{x}_t\|^2$

Output: v^θ

$p_t(x_t|x_1)$ general

$p(x_0)$ is general

Algorithm 2: Diffusion training.

Input : dataset q , noise p

Initialize s^θ

while *not converged* **do**

$t \sim \mathcal{U}([0, 1])$ ▷ sample time

$x_1 \sim q(x_1)$ ▷ sample data

$x_t = p_t(x_t|x_1)$ ▷ sample conditional prob

Gradient step with

$\nabla_\theta \|s_t^\theta(x_t) - \nabla_{x_t} \log p_t(x_t|x_1)\|^2$

Output: v^θ

$p_t(x_t|x_1)$ closed-form from of SDE $dx_t = f_t dt + g_t dw$

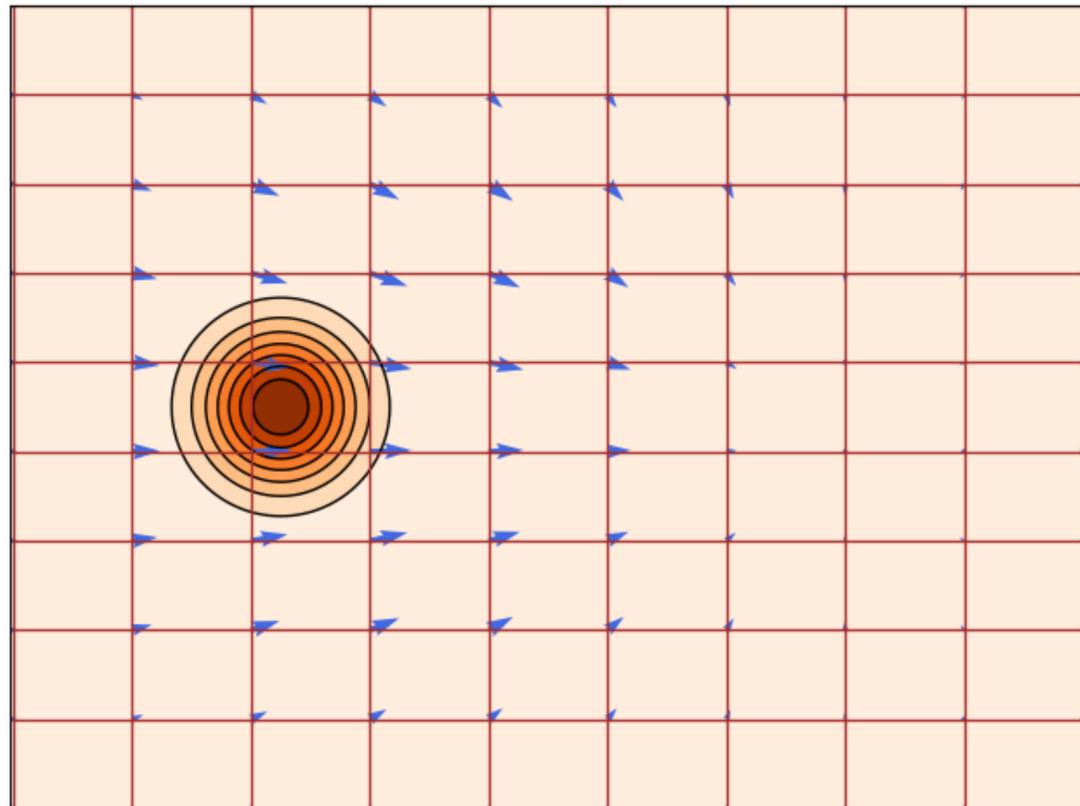
- **Variance Exploding:** $p_t(x|x_1) = \mathcal{N}(x|x_1, \sigma_{1-t}^2 I)$
- **Variance Preserving:** $p_t(x|x_1) = \mathcal{N}(x|\alpha_{1-t}x_1, (1 - \alpha_{1-t}^2)I)$
 $\alpha_t = e^{-\frac{1}{2}T(t)}$

$p(x_0)$ is Gaussian

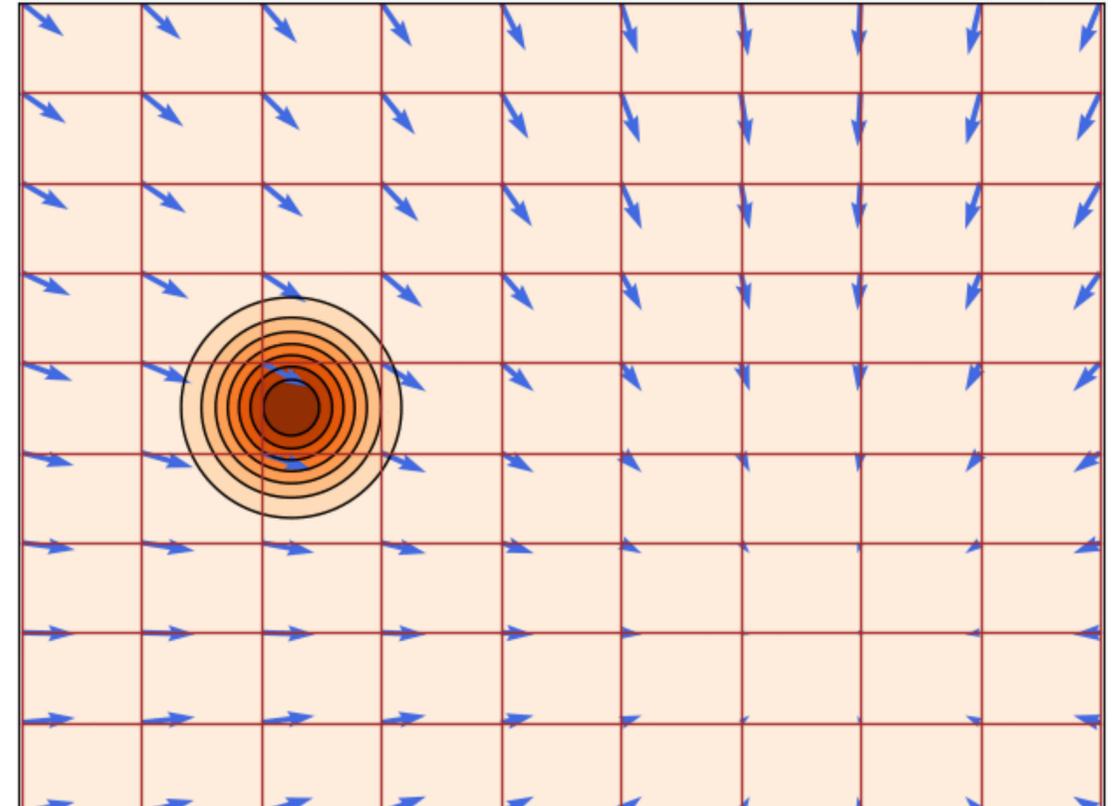
$p_0(\cdot|x_1) \approx p$

Why does this work?

- What we want is the flow (velocity field) that takes samples from p_0 to p_1 when integrated (called Marginal Flow)
- But what we did was to train flow for each sample (Conditional Flow)



Marginal Flow (what we want)



Conditional Flow (what we trained)

It turns out that the gradient is the same!

- **Flow Matching loss:**

$$\mathcal{L}_{\text{FM}}(\theta) = \mathbb{E}_{t, X_t} \left\| u_t^\theta(X_t) - u_t(X_t) \right\|^2$$

where $u_t(X_t)$ is the target (marginal) velocity field

- **Conditional Flow Matching loss:**

$$\mathcal{L}_{\text{CFM}}(\theta) = \mathbb{E}_{t, X_1, X_t} \left\| u_t^\theta(X_t) - u_t(X_t | X_1) \right\|^2$$

$u_t(X_t | X_1) = X_1 - X_0$

We can't do this because we don't know what ground truth flow is!

Theorem: Gradient of the losses are equivalent:

$$\nabla_{\theta} \mathcal{L}_{\text{FM}}(\theta) = \nabla_{\theta} \mathcal{L}_{\text{CFM}}(\theta)$$

Why does this work?

- High level: You can average conditional flow (marginalization)

$$u_t(x) = \mathbb{E} [u_t(X_t | X_1) | X_t = x]$$

Marginal flow (what we want) Conditional flow (easy to obtain from each sample)

$$= \int u_t(x | x_1) p(x_1 | X_t = x) dx_1$$

$$u_t(x) = \int u_t(x | x_1) \frac{p_t(x | x_1) q(x_1)}{p_t(x)} dx_1$$

Bayes rule with notation:

$$p(x_1) = q(x_1)$$

It's all just weighted average

$$u_t(x) = \int u_t(x | x_1) \frac{p_t(x | x_1) q(x_1)}{p_t(x)} dx_1$$

$$u_t(x_t) = \sum \underbrace{u_t(x_t | x_1)}_{\text{Path from } x_1 \text{ to } x_t} \underbrace{\frac{p_t(x_t | x_1)}{p_t(x)}}_{\text{Path weight (how relevant } x_1 \text{ is to } x_t)} q(x_1)$$

Weight of sample $x_0 = 1$

Marginal flow

Just a weighted average of the flow to each data sample!

Widely used in state-of-the-art generators

Scaling Rectified Flow Transformers for High-Resolution Image Synthesis

Patrick Esser* Sumith Kulal Andreas Blattmann Rahim Entezari Jonas Müller Harry Saini Yam Levi
Dominik Lorenz Axel Sauer Frederic Boesel Dustin Podell Tim Dockhorn Zion English
Kyle Lacey Alex Goodwin Yannik Marek Robin Rombach*
Stability AI



***FLUX.1 Kontext*: Flow Matching for In-Context Image Generation and Editing in Latent Space**

Black Forest Labs*

Stephen Batifol Andreas Blattmann Frederic Boesel Saksham Consul Cyril Diagne
Tim Dockhorn Jack English Zion English Patrick Esser Sumith Kulal
Kyle Lacey Yam Levi Cheng Li Dominik Lorenz Jonas Müller
Dustin Podell Robin Rombach Harry Saini Axel Sauer Luke Smith

Abstract

We present evaluation results for *FLUX.1 Kontext*, a generative flow matching model that unifies image generation and editing. The model generates novel output views by incorporating semantic context from text and image inputs. Using a simple

WAN: OPEN AND ADVANCED LARGE-SCALE VIDEO GENERATIVE MODELS

Wan Team, Alibaba Group

ABSTRACT

This report presents Wan, a comprehensive and open suite of video foundation models designed to push the boundaries of video generation. Built upon the mainstream diffusion transformer paradigm, Wan achieves significant advancements in generative capabilities through a series of innovations, including our novel spatio-temporal variational autoencoder (VAE), scalable pre-training strategies, large-scale data curation, and automated evaluation metrics. These contributions collectively enhance the model's performance and versatility. Specifically, Wan is characterized by four key features: *Leading Performance*: The 14B model of Wan, trained on a vast dataset comprising billions of images and videos, demonstrates the scaling laws of video generation with respect to both data and model size. It consistently outperforms the existing open-source models as well as state-of-

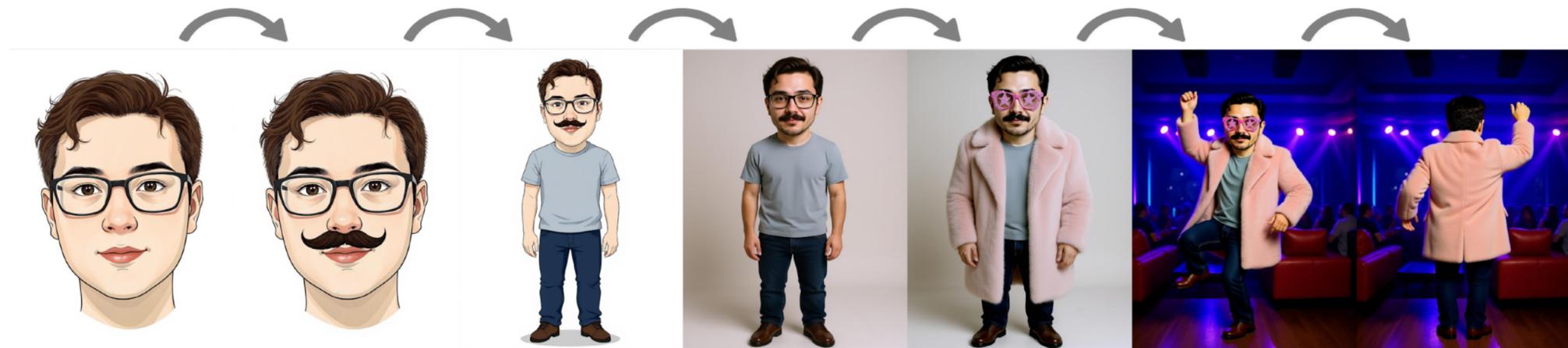
FRIEREN: Efficient Video-to-Audio Generation Network with Rectified Flow Matching

Yongqi Wang[†] Wenxiang Guo[‡] Rongjie Huang, Jiawei Huang, Zehan Wang,
Fuming You, Ruiqi Li, Zhou Zhao*
Zhejiang University
cyanbox@zju.edu.cn

Abstract

Video-to-audio (V2A) generation aims to synthesize content-matching audio from silent video, and it remains challenging to build V2A models with high generation quality, efficiency, and visual-audio temporal synchrony. We propose FRIEREN, a V2A model based on rectified flow matching. FRIEREN regresses the conditional transport vector field from noise to spectrogram latent with straight paths and

Flow matching in state-of-art generators



Diffusion vs. Flow matching

- They both end up learning flow, but diffusion poses the problem as learning a specific noising process and learning to denoise it.
- Diffusion: spread out like heat diffusion, learn how to undo it
- The flow matching perspective is very general: just directly learn a velocity field from one distribution to another.

Next class: image manipulation with diffusion models