

EECE 571F: Advanced Topics in Deep Learning

Lecture 8: Flow Matching

Qi Yan, Yuanpei Gao, Renjie Liao

University of British Columbia

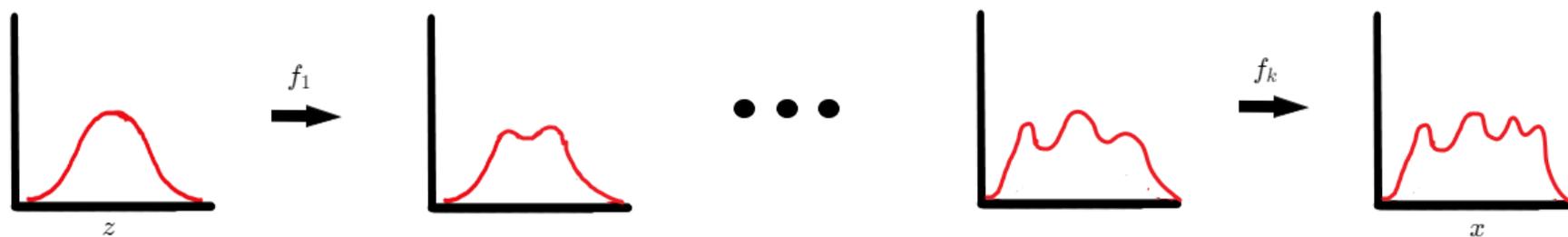
Winter, Term 2, 2025

Outline

- **Normalizing Flows and Continuous Normalizing Flows**
 - The Continuity Equation
- The Fokker Plank Equation
- Flow matching
- Variants:
 - Batch Optimal Transport Flow Matching

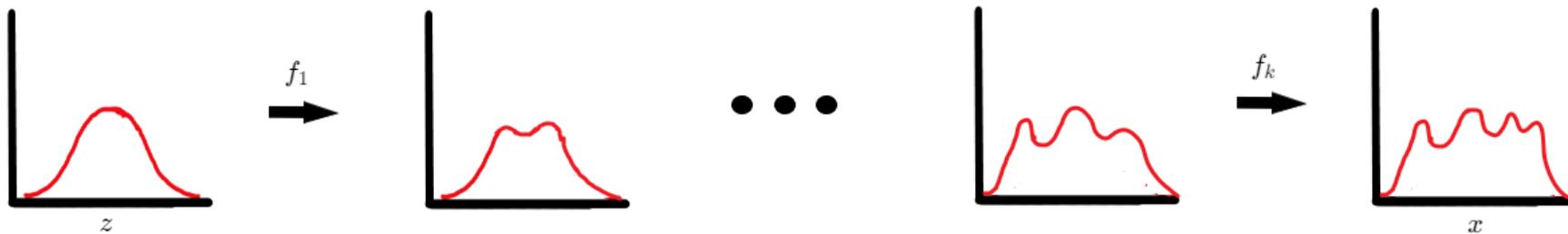
Normalizing Flows

- Our goal with this setup is to learn the transformation from $p(\cdot)$ to the complex data distribution $q(\cdot)$.



Normalizing Flows

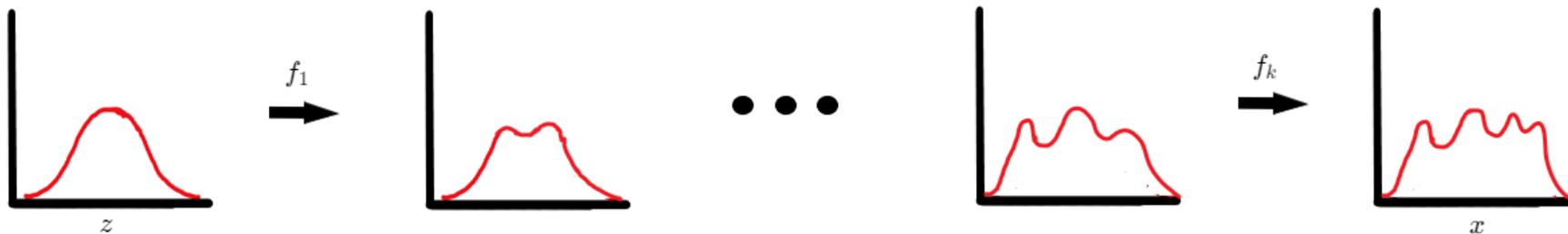
- Our goal with this setup is to learn the transformation from $p(\cdot)$ to the complex data distribution $q(\cdot)$.
- We can do this by learning the invertible transformation f_θ using neural networks.



Normalizing Flows

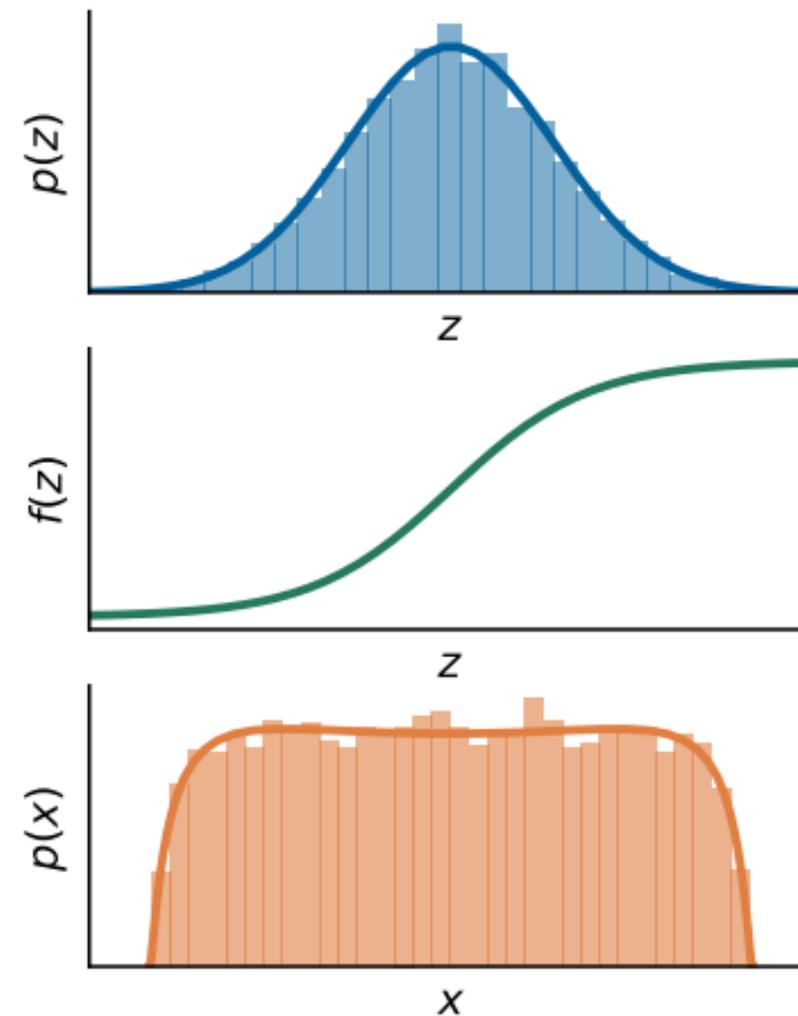
- Our goal with this setup is to learn the transformation from $p(\cdot)$ to the complex data distribution $q(\cdot)$.
- We can do this by learning the invertible transformation f_θ using neural networks.
- f_θ can contain multiple transformations. Each transformation transforms an input distribution into a slightly more complex distribution.

$$f_\theta = f_k \circ f_{k-1} \cdots f_2 \circ f_1.$$



Normalizing Flows

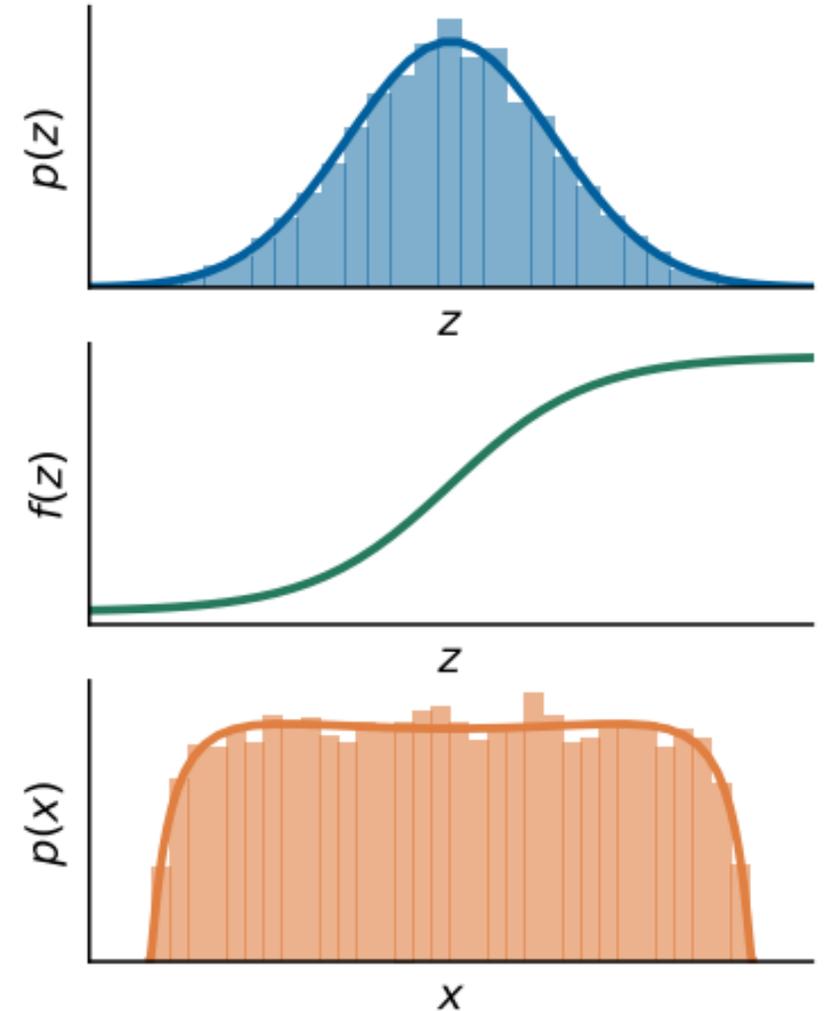
- Starting with known distribution $z \sim p_z(\cdot)$



Normalizing Flows

- Starting with known distribution $z \sim p_z(\cdot)$
- Let f_θ be an invertible and differentiable function, apply the transformation to z :

$$p_\theta(\mathbf{x}) = p_z(f_\theta^{-1}(\mathbf{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right|$$



Normalizing Flows – Multivariate Change of Variable

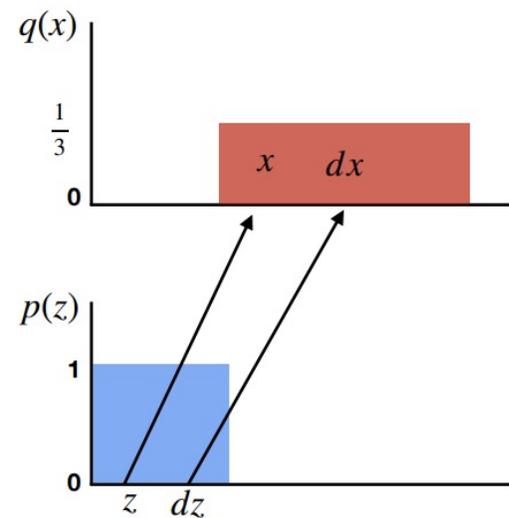
- Consider we have $\boldsymbol{x} = f(\boldsymbol{z})$, when transforming coordinates from \boldsymbol{z} -space to \boldsymbol{x} -space, we are interested in understanding how infinitesimal regions around a point in the original space change under the transformation.
- The function $f(\boldsymbol{z})$ can be approximated using first-order Taylor expansion:

$$\boldsymbol{x} \simeq f(\boldsymbol{z}_0) + J(\boldsymbol{z}_0)(\boldsymbol{z} - \boldsymbol{z}_0)$$

Normalizing Flows – Multivariate Change of Variable

- Based on the probability density preservation under transformation, we can have:

$$p_{\mathbf{x}}(\mathbf{x})d\mathbf{x} = p_{\mathbf{z}}(\mathbf{z})d\mathbf{z}$$

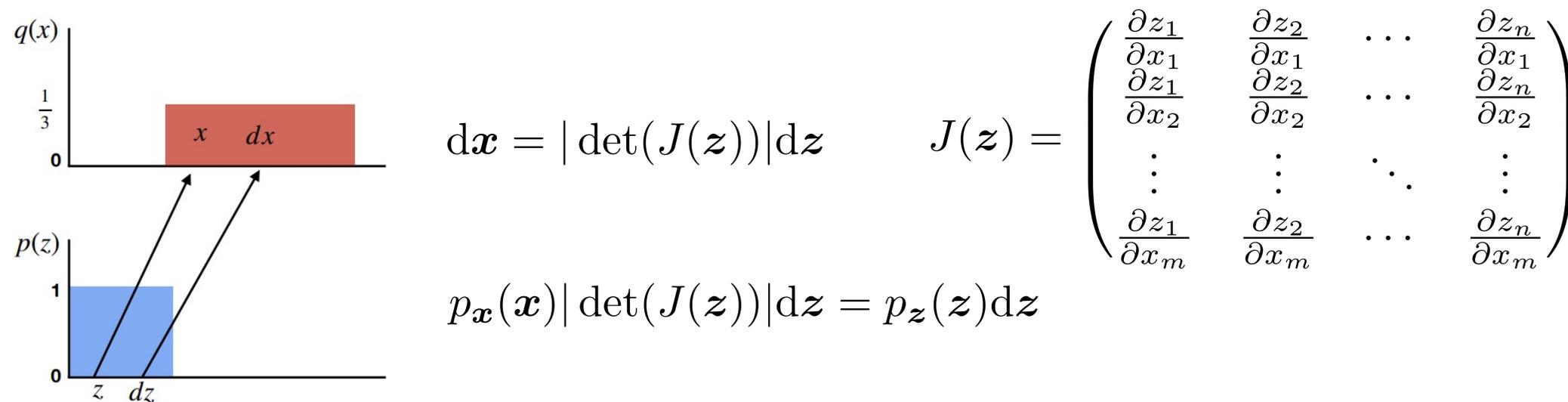


Normalizing Flows – Multivariate Change of Variable

- Based on the probability density preservation under transformation, we can have:

$$p_{\mathbf{x}}(\mathbf{x})d\mathbf{x} = p_{\mathbf{z}}(\mathbf{z})d\mathbf{z}$$

- The infinitesimal volume transform is (only the linear term matters):



Normalizing Flows – Multivariate Change of Variable

- Rearrange the equations and we will have:

$$p_{\mathbf{x}}(\mathbf{x}) = p_{\mathbf{z}}(\mathbf{z}) |\det(J(\mathbf{z}))|^{-1} = p_{\mathbf{z}}(f^{-1}(\mathbf{x})) |\det(J(f^{-1}(\mathbf{x})))|^{-1}$$

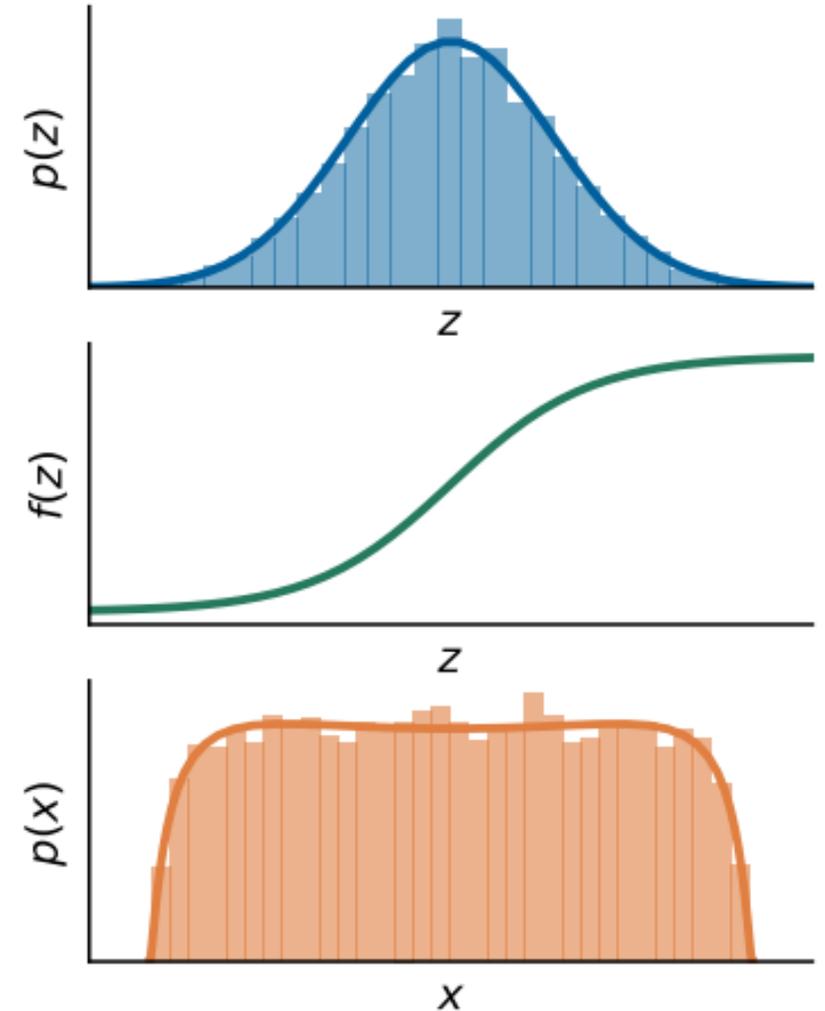
Normalizing Flows

- Starting with known distribution $z \sim p_z(\cdot)$
- Let f_θ be an invertible and differentiable function, apply the transformation to z :

$$p_\theta(\mathbf{x}) = p_z(f_\theta^{-1}(\mathbf{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right|$$

- Maximize likelihood of data:

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log p_\theta(x_i)$$



Normalizing Flows

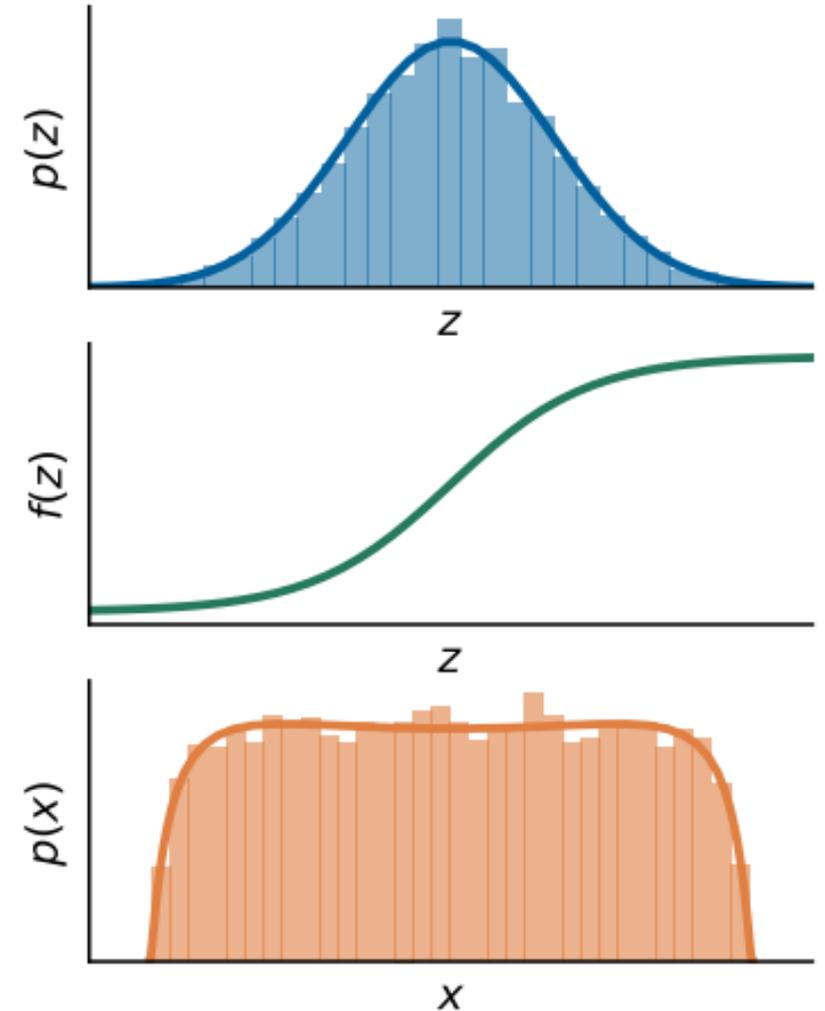
- Starting with known distribution $z \sim p_z(\cdot)$
- Let f_θ be an invertible and differentiable function, apply the transformation to z :

$$p_\theta(\mathbf{x}) = p_z(f_\theta^{-1}(\mathbf{x})) \cdot \left| \det \frac{\partial f_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}} \right|$$

- Maximize likelihood of data:

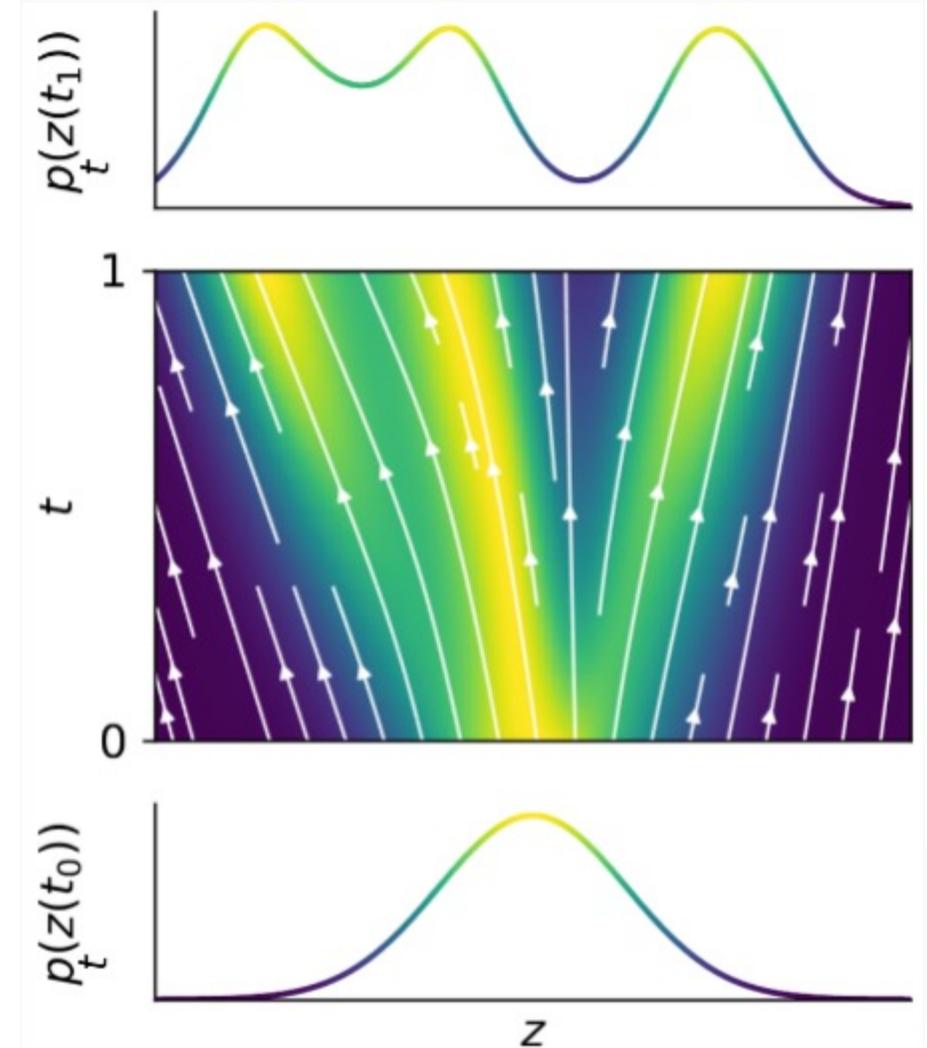
$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log p_\theta(x_i)$$

- $\frac{\partial f_\theta^{-1}(\mathbf{x})}{\partial \mathbf{x}}$ is the Jacobian of the transformation $f_\theta^{-1}(\mathbf{x})$



Continuous Normalizing Flows

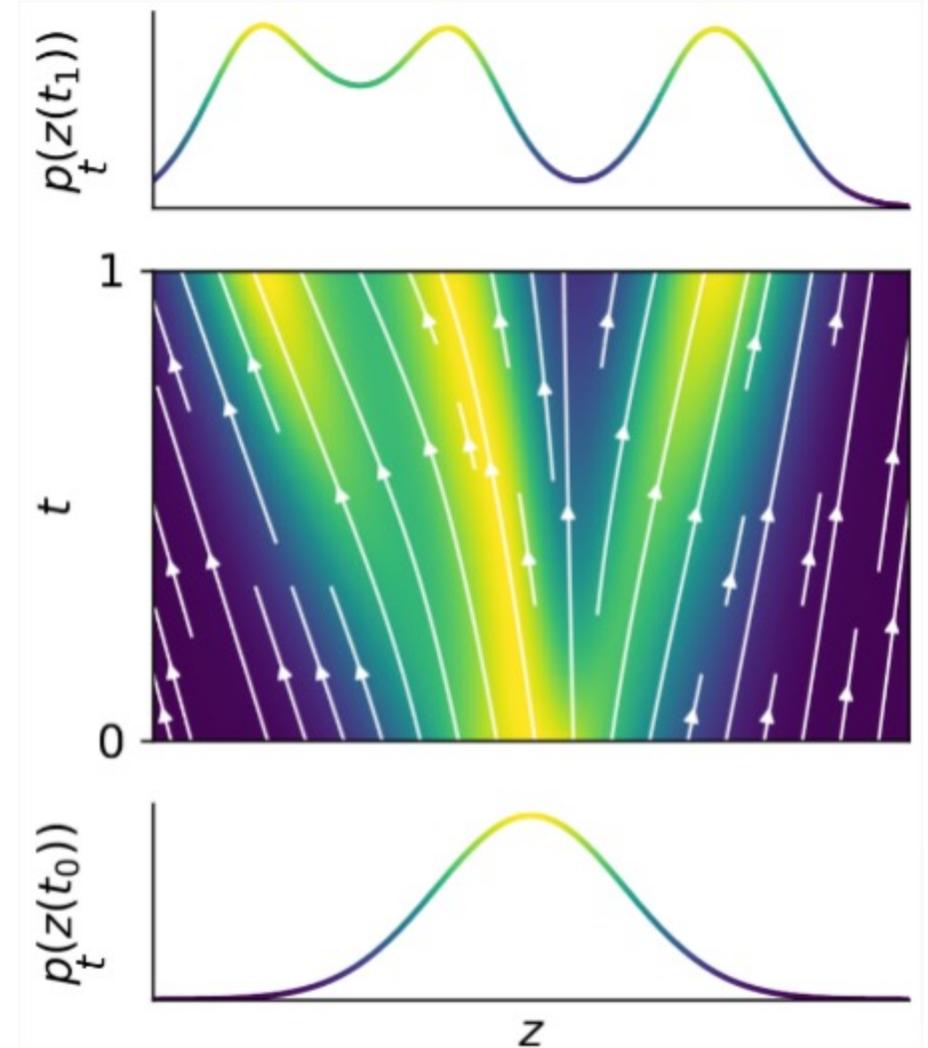
- Continuously normalizing flows are a **generalization of normalizing flows** where the transformations are parameterized by continuous dynamics governed by an ordinary differential equation (ODE).



Continuous Normalizing Flows

- Define the transformation as an ODE

$$\mathbf{x} = \mathbf{z}(t_1) = \int_{t_0}^{t_1} \mathbf{v}_\theta(\mathbf{z}(t), t) dt$$

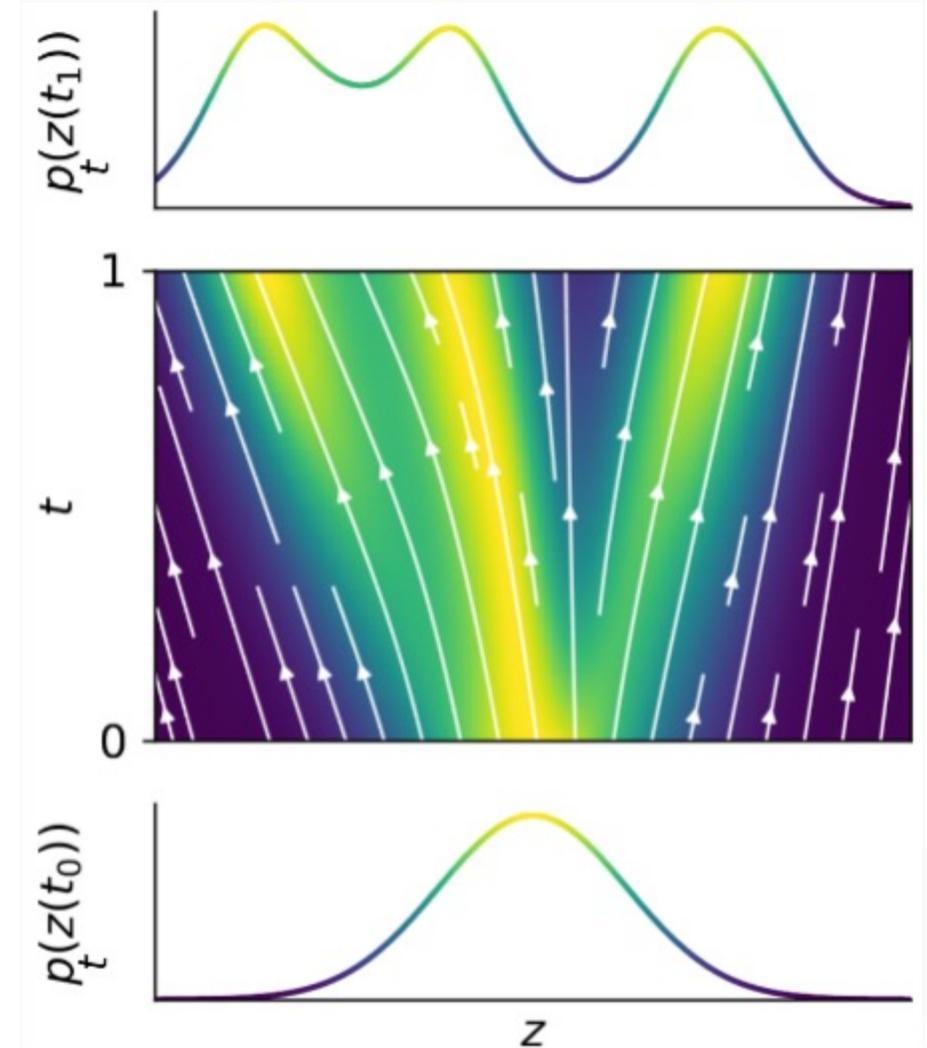


Continuous Normalizing Flows

- Define the transformation as an ODE

$$\boldsymbol{x} = \boldsymbol{z}(t_1) = \int_{t_0}^{t_1} \boldsymbol{v}_\theta(\boldsymbol{z}(t), t) dt$$

- Here $\boldsymbol{v}_\theta(\boldsymbol{z}(t), t)$ represents the velocity field of the latent variable \boldsymbol{z} as it evolves under a continuous transformation

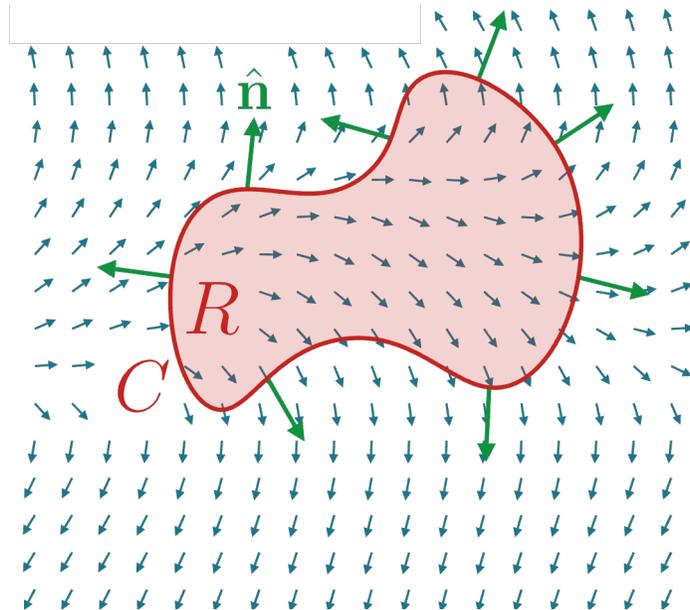


Outline

- Normalizing Flows and Continuous Normalizing Flows
 - **The Continuity Equation**
- The Fokker Plank Equation
- Flow matching
- Variants:
 - Batch Optimal Transport Flow Matching

Continuous Normalizing Flows

- Gauss's Divergence Theorem: the flux of a vector field through a closed surface equals the volume integral of its divergence over the enclosed region.

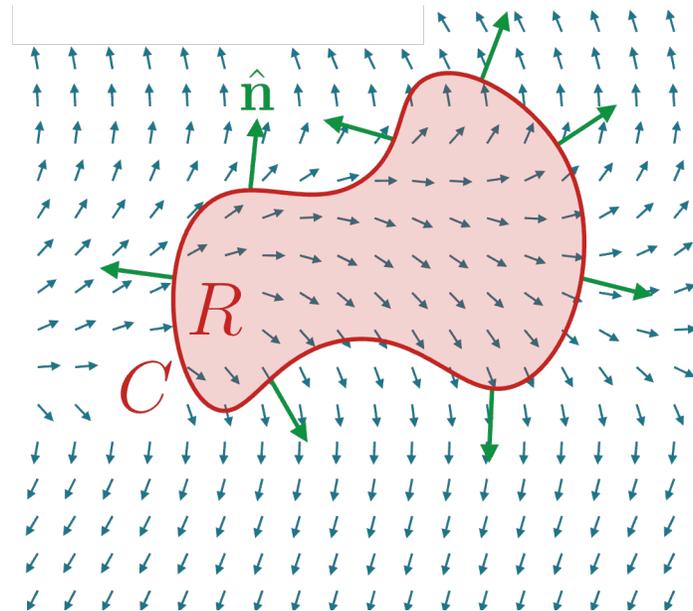


$$\oint (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) \cdot \mathbf{n} dC = \iint_R \nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) dR$$

Flux integral through the boundary C
Divergence integral over the region R

Continuous Normalizing Flows

- Gauss's Divergence Theorem: the flux of a vector field through a closed surface equals the volume integral of its divergence over the enclosed region.



$$\oint (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) \cdot \mathbf{n} dC = \iint_R \nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) dR$$

Flux integral through the boundary C

Divergence integral over the region R

$p(\mathbf{z}(t))$ is the density at position $\mathbf{z}(t)$

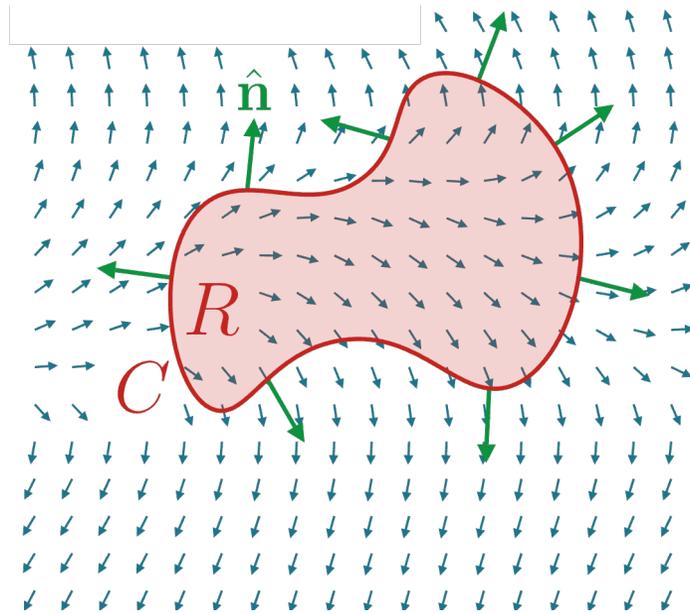
$\mathbf{v}_\theta(\mathbf{z}(t), t)$ describes the relevant flow

$p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)$ describes how much density flows per unit time in a unit area.

Physical analogy: Think of the flow of fluid mass!

Continuous Normalizing Flows

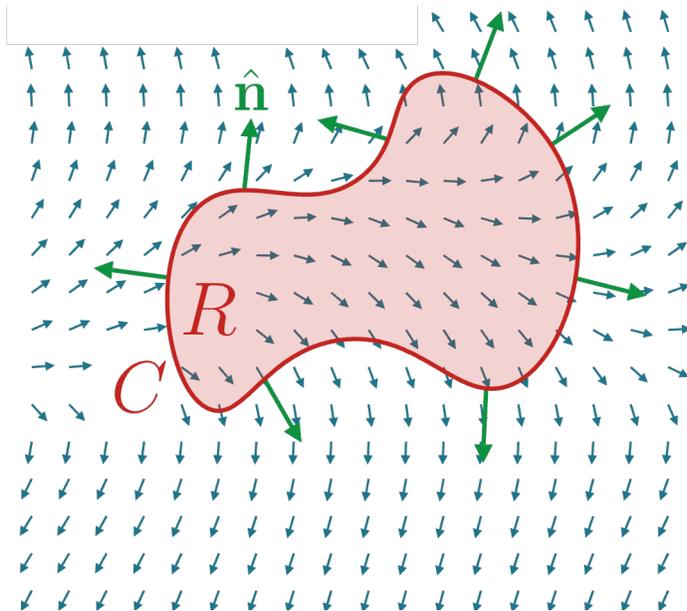
- Consider the law of conservation (the continuity equation):



$$\underbrace{\iint_R \frac{\partial p(\mathbf{z}(t))}{\partial t} dR}_{\text{Flux in over the region } R} + \underbrace{\oint (p(\mathbf{z}(t)) \mathbf{v}_\theta(\mathbf{z}(t), t)) \cdot \mathbf{n} dC}_{\text{Flux out through the boundary } C} = 0$$

Continuous Normalizing Flows

- Consider the law of conservation (the continuity equation):



$$\iint_R \frac{\partial p(\mathbf{z}(t))}{\partial t} dR + \oint (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) \cdot \mathbf{n} dC = 0$$

Apply Gauss's Divergence Theorem

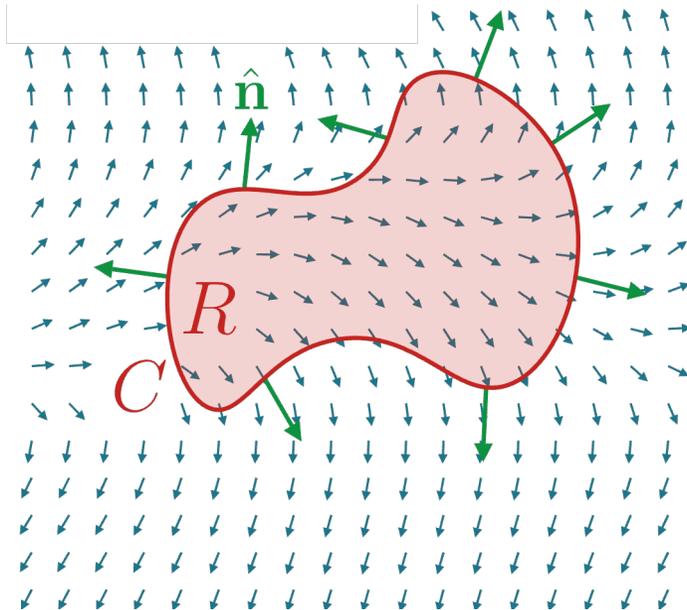
$$\iint_R \frac{\partial p(\mathbf{z}(t))}{\partial t} dR + \iint_R \nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) dR = 0$$

Flux in over the region R

Divergence integral over the region R

Continuous Normalizing Flows

- Consider the law of conservation (the continuity equation):



$$\iint_R \frac{\partial p(\mathbf{z}(t))}{\partial t} dR + \oint (p(\mathbf{z}(t)) \mathbf{v}_\theta(\mathbf{z}(t), t)) \cdot \mathbf{n} dC = 0$$

$$\iint_R \frac{\partial p(\mathbf{z}(t))}{\partial t} dR + \iint_R \nabla \cdot (p(\mathbf{z}(t)) \mathbf{v}_\theta(\mathbf{z}(t), t)) dR = 0$$

$$\frac{\partial p(\mathbf{z}(t))}{\partial t} + \nabla \cdot (p(\mathbf{z}(t)) \mathbf{v}_\theta(\mathbf{z}(t), t)) = 0$$

Continuity equation (differential form)

This is due to the fact that the conservation law holds for all kinds of regions, densities, and velocity fields!

Continuous Normalizing Flows

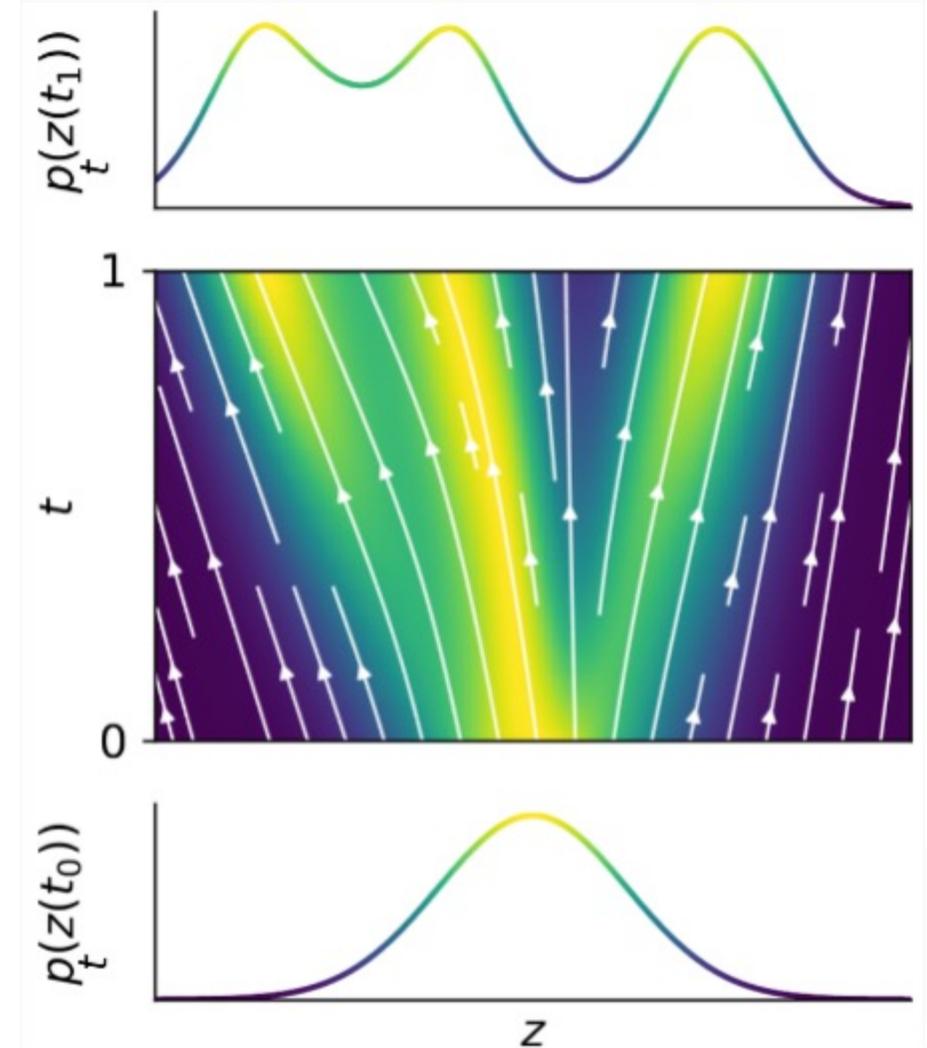
- The continuity equation:

$$\frac{\partial p(z(t))}{\partial t} + \nabla \cdot (p(z(t))\mathbf{v}_\theta(z(t), t)) = 0$$

Flux in

Flux out

- The continuity equation is a **principle of conservation** in fluid dynamics and other physical systems. It states that the change in density over time is balanced by the flux of density due to the velocity field.



Continuous Normalizing Flows

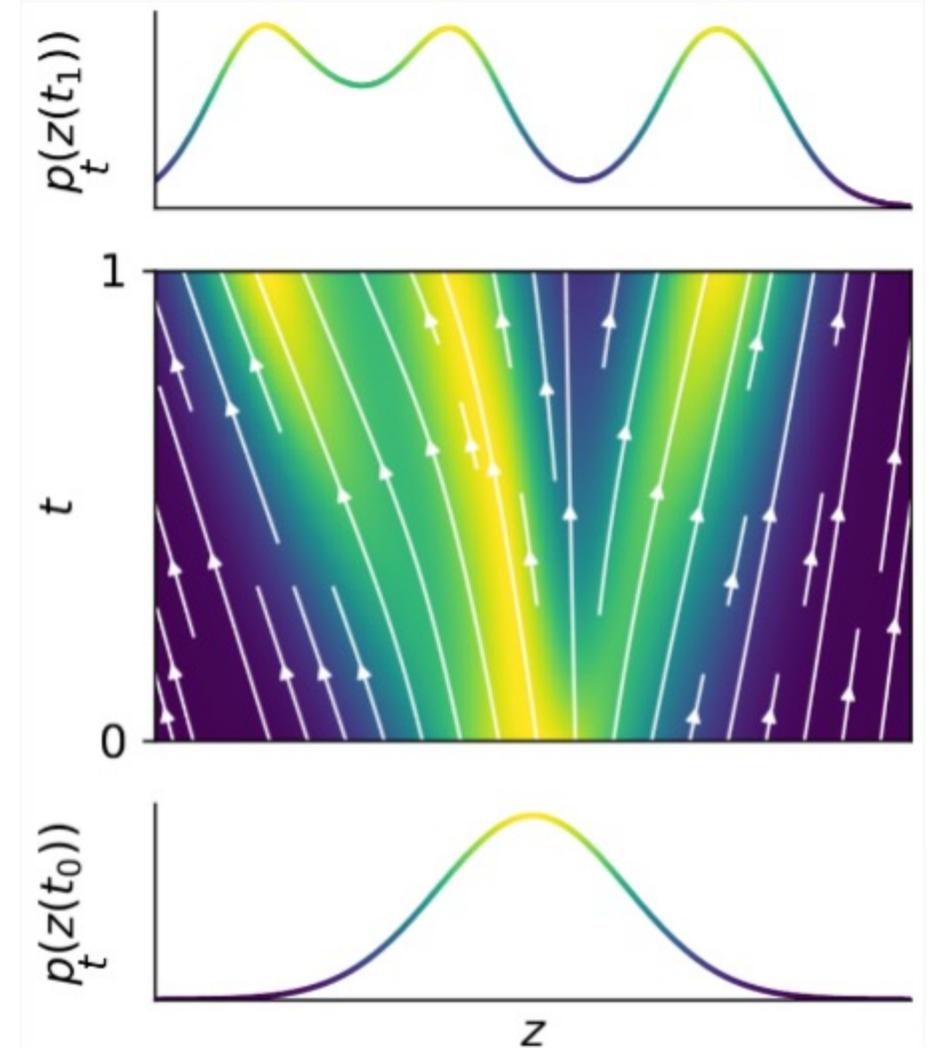
- The continuity equation:

$$\frac{\partial p(\mathbf{z}(t))}{\partial t} + \nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t)) = 0$$

Flux in

Flux out

- The divergence symbol $\nabla \cdot$ measures the "net flow" of a vector field out of a point in space.



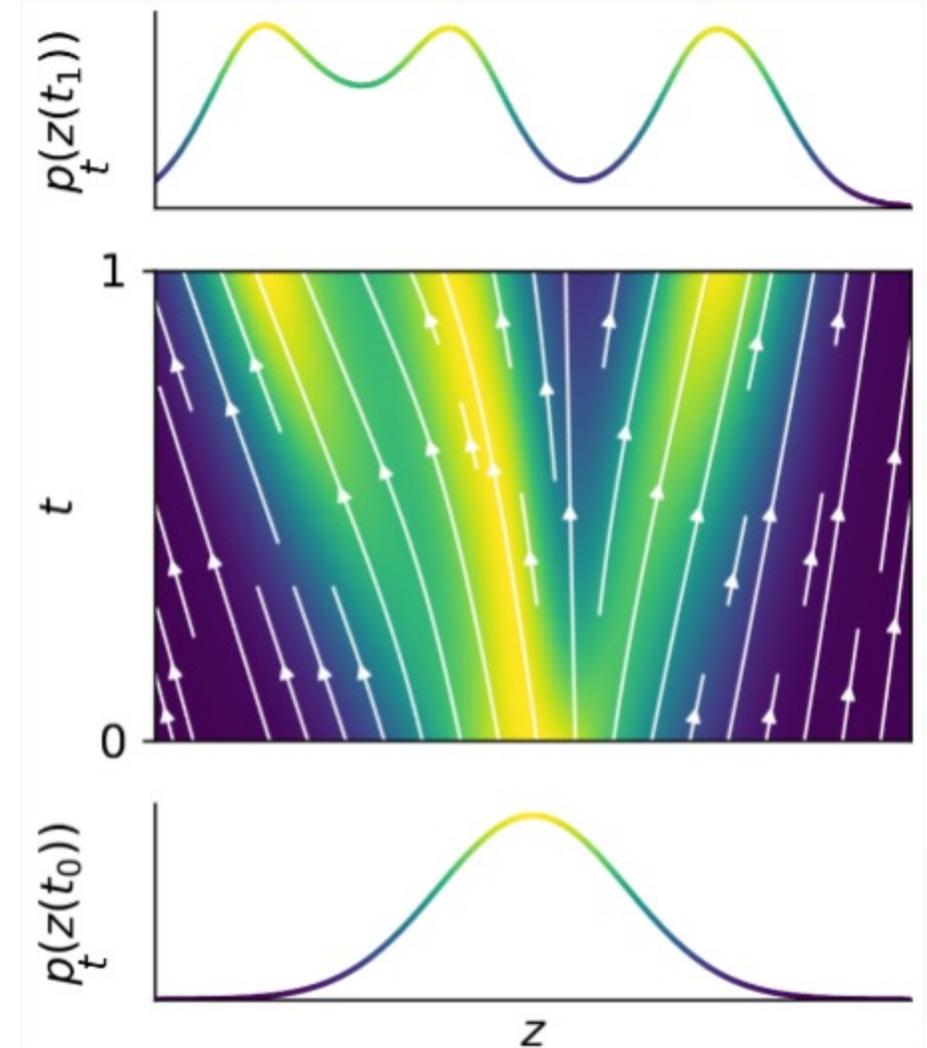
Continuous Normalizing Flows

- The continuity equation:

$$\underbrace{\frac{\partial p(\mathbf{z}(t))}{\partial t}}_{\text{Flux in}} + \underbrace{\nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t))}_{\text{Flux out}} = 0$$

- The divergence symbol $\nabla \cdot$ measures the "net flow" of a vector field out of a point in space.
- For $\mathbf{v}(\mathbf{z}) = [v_1(\mathbf{z}), v_2(\mathbf{z}), \dots, v_n(\mathbf{z})]$

$$\nabla \cdot \mathbf{z} = \frac{\partial v_1}{\partial z_1} + \frac{\partial v_2}{\partial z_2} + \dots + \frac{\partial v_n}{\partial z_n}$$



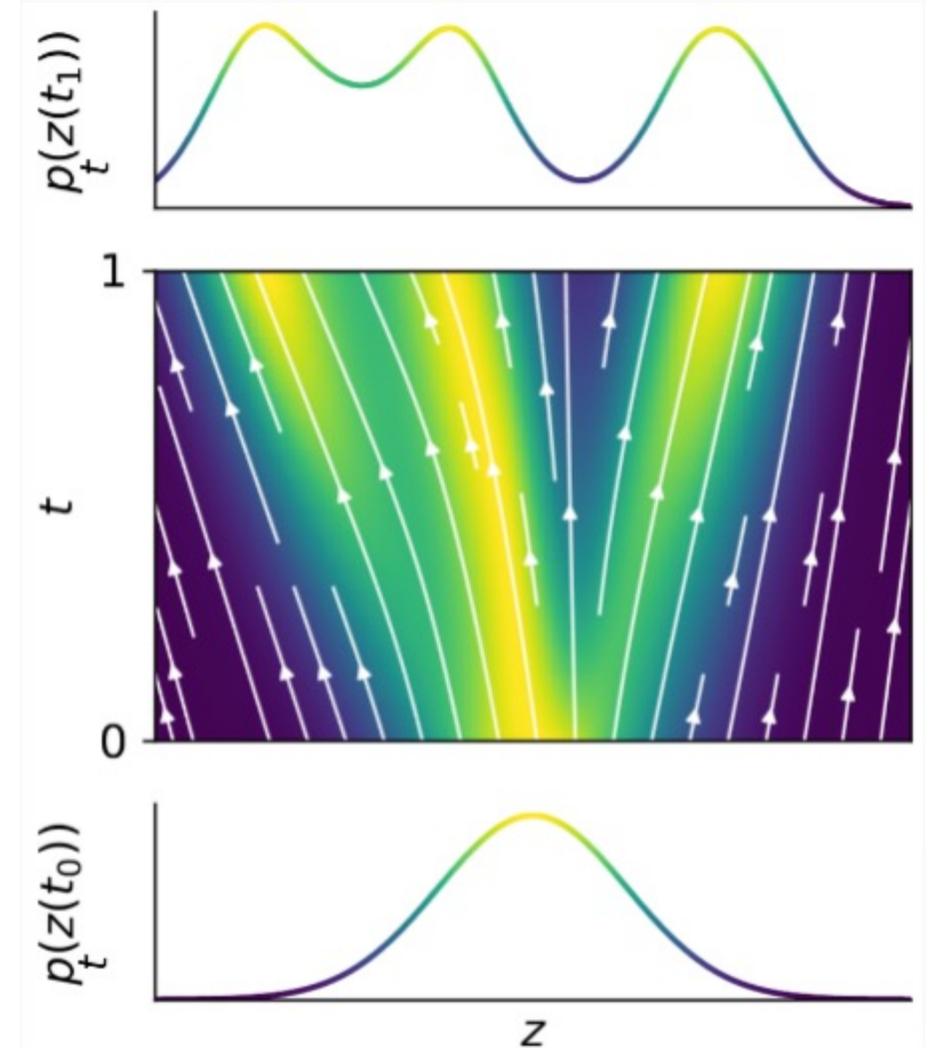
Continuous Normalizing Flows

- The continuity equation:

$$\underbrace{\frac{\partial p(\mathbf{z}(t))}{\partial t}}_{\text{Flux in}} + \underbrace{\nabla \cdot (p(\mathbf{z}(t))\mathbf{v}_\theta(\mathbf{z}(t), t))}_{\text{Flux out}} = 0$$

- One property of the divergence operator is the product rule:

$$\nabla \cdot (p_t(\mathbf{x})u_t(\mathbf{x})) = p_t(\mathbf{x})\nabla \cdot u_t(\mathbf{x}) + u_t(\mathbf{x})^T \nabla_{\mathbf{x}} p_t(\mathbf{x})$$



Continuous Normalizing Flows - Instantaneous change of density

- In CNFs, we transform a simple distribution to a more complex target distribution, and the challenge is understanding how the probability density changes during the transformation. And this change is governed by **the instantaneous change of density**.
- Here we use $\phi_t(\mathbf{x})$ to denote the flow trajectory.
- Consider the total derivative of $\log p_t(\phi_t(\mathbf{x}))$

$$\begin{aligned}\frac{d \log p_t(\phi_t(\mathbf{x}))}{dt} &= \frac{\partial \log p_t(\phi_t(\mathbf{x}))}{\partial t} \cdot \frac{\partial t}{\partial t} + \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x})) \cdot \frac{d\phi_t(\mathbf{x})}{dt} \\ &= \frac{\partial \log p_t(\phi_t(\mathbf{x}))}{\partial t} + \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x})) \cdot \frac{d\phi_t(\mathbf{x})}{dt} \\ &= \frac{\partial \log p_t(\phi_t(\mathbf{x}))}{\partial t} + \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x})) \cdot u_t(\phi_t(\mathbf{x}))\end{aligned}$$

Continuous Normalizing Flows - Instantaneous change of density

- The continuity equation with the product rule of divergence:

$$\frac{\partial}{\partial t} p_t(\phi_t(\mathbf{x})) + p_t(\phi_t(\mathbf{x})) \nabla \cdot u_t(\phi_t(\mathbf{x})) + u_t(\phi_t(\mathbf{x}))^T \nabla_{\mathbf{x}} p_t(\phi_t(\mathbf{x})) = 0$$

$$\frac{1}{p_t(\phi_t(\mathbf{x}))} \left(\frac{\partial}{\partial t} p_t(\phi_t(\mathbf{x})) + p_t(\phi_t(\mathbf{x})) \nabla \cdot u_t(\phi_t(\mathbf{x})) + u_t(\phi_t(\mathbf{x}))^T \nabla_{\mathbf{x}} p_t(\phi_t(\mathbf{x})) \right) = 0$$

$$\frac{\partial}{\partial t} \log p_t(\phi_t(\mathbf{x})) = -\nabla \cdot u_t(\phi_t(\mathbf{x})) - u_t(\phi_t(\mathbf{x}))^T \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x}))$$

Continuous Normalizing Flows - Instantaneous change of density

- Consider the total derivative of $\log p_t(\phi_t(\mathbf{x}))$

$$\frac{d \log p_t(\phi_t(\mathbf{x}))}{dt} = \frac{\partial \log p_t(\phi_t(\mathbf{x}))}{\partial t} + \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x})) \cdot u_t(\phi_t(\mathbf{x}))$$

- The continuity equation with the product rule of divergence:

$$\frac{\partial}{\partial t} \log p_t(\phi_t(\mathbf{x})) = -\nabla \cdot u_t(\phi_t(\mathbf{x})) - u_t(\phi_t(\mathbf{x})) \cdot \nabla_{\mathbf{x}} \log p_t(\phi_t(\mathbf{x}))$$

- Now, replace the first term with continuity equation, we will have:

$$\log p_1(\phi_1(\mathbf{x})) = \log p_0(\phi_0(\mathbf{x})) - \int_0^1 \nabla \cdot u_t(\phi_t(\mathbf{x})) dt$$

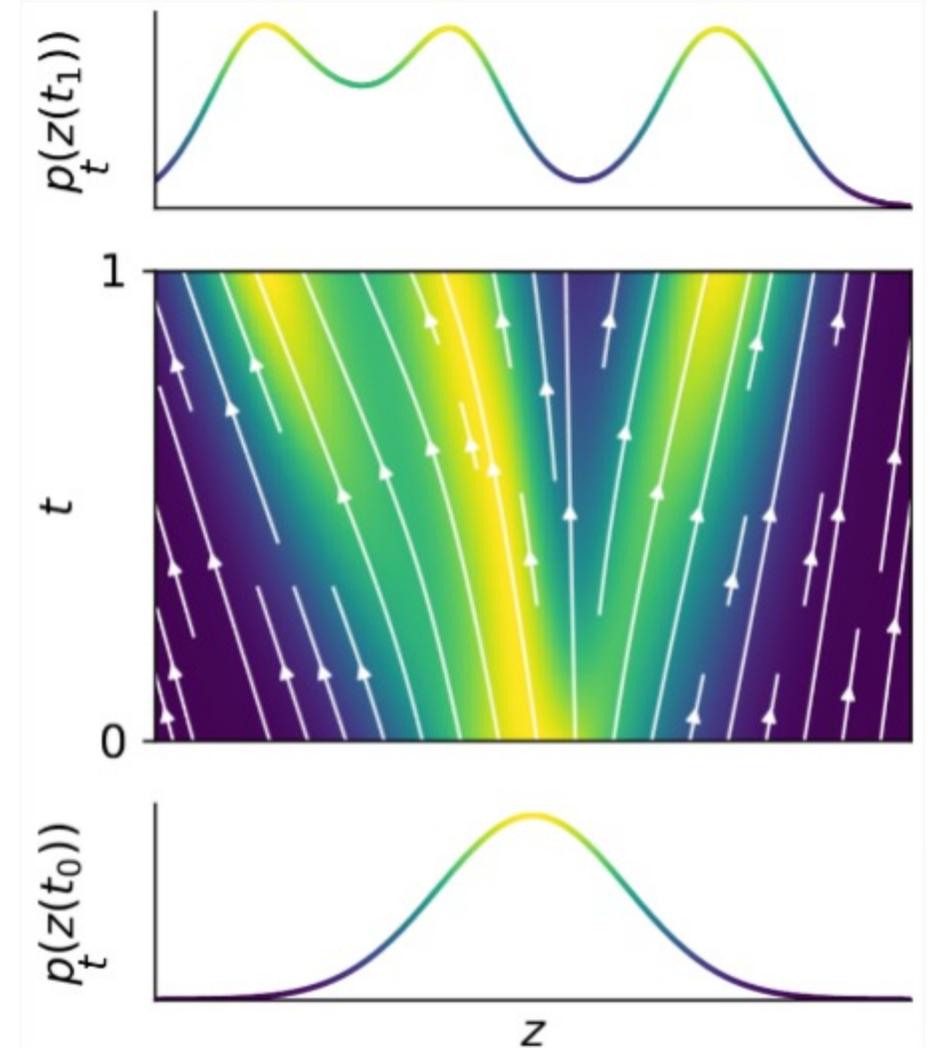
Continuous Normalizing Flows

- Define the transformation as an ODE

$$\mathbf{x} = \mathbf{z}(t_1) = \int_{t_0}^{t_1} \mathbf{v}_\theta(\mathbf{z}(t), t) dt$$

- Instantaneous change of density

$$\frac{\partial \log p_t(\mathbf{z}(t))}{\partial t} = -\nabla \cdot \mathbf{v}_\theta(\mathbf{z}(t), t)$$



Continuous Normalizing Flows

- Define the transformation as an ODE

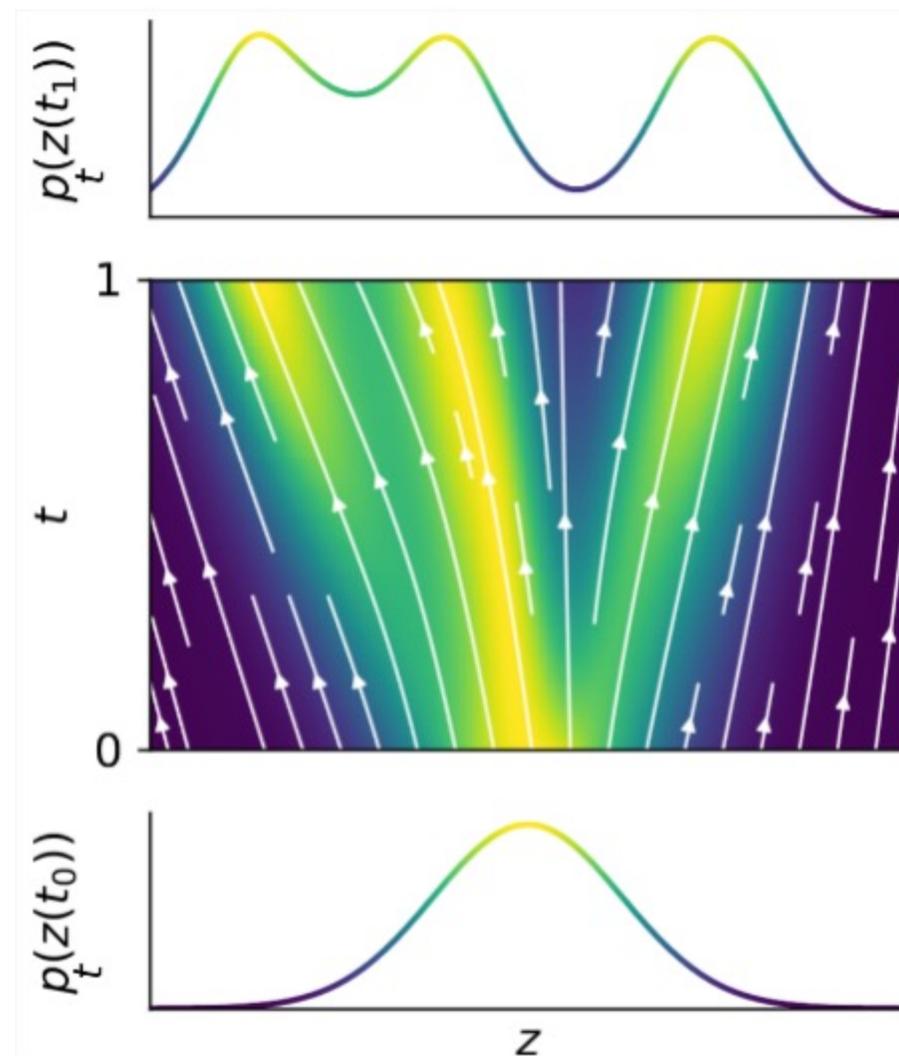
$$\mathbf{x} = \mathbf{z}(t_1) = \int_{t_0}^{t_1} \mathbf{v}_\theta(\mathbf{z}(t), t) dt$$

- Instantaneous change of density

$$\frac{\partial \log p_t(\mathbf{z}(t))}{\partial t} = -\nabla \cdot \mathbf{v}_\theta(\mathbf{z}(t), t)$$

- Solve the ODE for $\log p_{t_1}(\mathbf{z}(t_1))$

$$\log p_{t_0}(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \nabla \cdot \mathbf{v}_\theta(\mathbf{z}(t), t) dt$$



Training of the Neural ODEs

- The ODEs parameterized by neural networks are called Neural ODEs.
- We still adopt maximum likelihood training objective.

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^N \log p_{\theta}(\mathbf{x}_i)$$

Training of the Neural ODEs

- Training requires simulation (solving ODE) to obtain exact likelihood

$$\begin{aligned}\boxed{\log p_{\theta}(\mathbf{x})} &= \log p_{t_0}(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt \\ &= \log p_{t_0}(\mathbf{z}(t_0)) + \int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt\end{aligned}$$

Training of the Neural ODEs

- Training requires simulation (solving ODE) to obtain exact likelihood

$$\begin{aligned}\log p_{\theta}(\mathbf{x}) &= \log p_{t_0}(\mathbf{z}(t_0)) - \int_{t_0}^{t_1} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt \\ &= \log p_{t_0}(\mathbf{z}(t_0)) + \int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt\end{aligned}$$

Both need to be numerically solved through ODEs

Training of the Neural ODEs

- Training requires simulation (solving ODE) to obtain exact likelihood

$$\log p_{\theta}(\mathbf{x}) = \log p_{t_0}(\mathbf{z}(t_0)) + \int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt$$

$$\underbrace{\frac{d}{d\tilde{t}}}_{\text{inversed } t_1 \rightarrow t_0} \left[\int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt \right] = \begin{bmatrix} -\mathbf{v}_{\theta}(\mathbf{z}(t), t) \\ \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) \end{bmatrix}$$

Training of the Neural ODEs

- Training requires simulation (solving ODE) to obtain exact likelihood

$$\log p_{\theta}(\mathbf{x}) = \log p_{t_0}(\mathbf{z}(t_0)) + \int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt$$

$$\underbrace{\frac{d}{d\tilde{t}}}_{\text{inversed } t_1 \rightarrow t_0} \left[\int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt \right] = \begin{bmatrix} -\mathbf{v}_{\theta}(\mathbf{z}(t), t) \\ \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) \end{bmatrix}$$

Trace of Jacobian.

We can use Hutchinson's trace estimator.

Training of the Neural ODEs

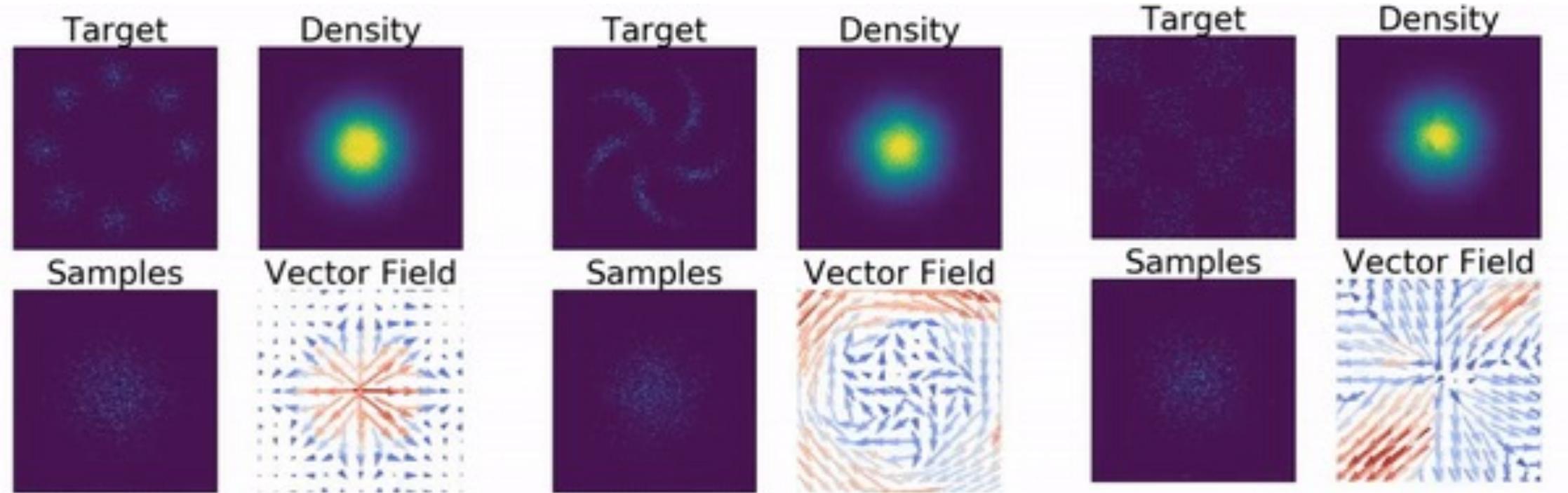
- Training requires simulation (solving ODE) to obtain exact likelihood

$$\log p_{\theta}(\mathbf{x}) = \log p_{t_0}(\mathbf{z}(t_0)) + \int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt$$

$$\frac{d}{d\tilde{t}} \left[\int_{t_1}^{t_0} \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) dt \right] = \begin{bmatrix} -\mathbf{v}_{\theta}(\mathbf{z}(t), t) \\ \nabla \cdot \mathbf{v}_{\theta}(\mathbf{z}(t), t) \end{bmatrix}$$

- Solving ODEs numerically at each training iteration is slow!
- Gradient computation for backpropagation requires careful handling (adjoint method).

Continuous Normalizing Flows



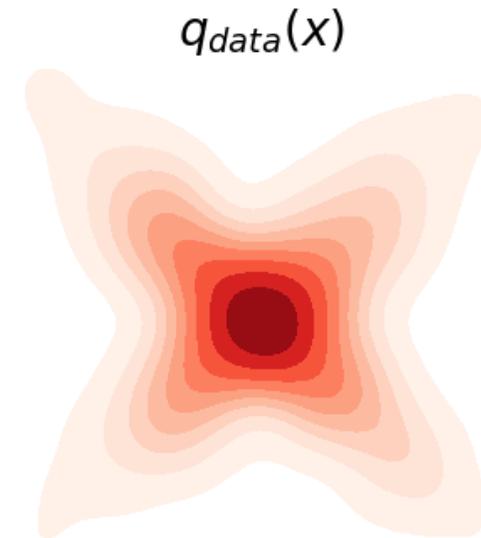
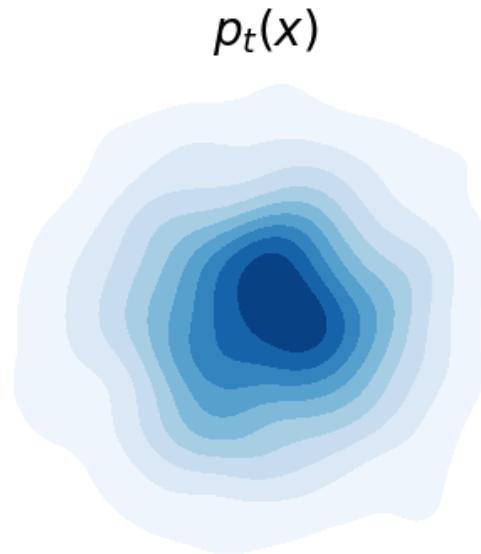
Outline

- Normalizing Flows and Continuous Normalizing Flows
 - The Continuity Equation
- **The Fokker Plank Equation**
- Flow matching
- Variants:
 - Batch Optimal Transport Flow Matching

The Fokker Plank Equation

- What happens to the continuity equation if there is stochastic noise?

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\mathbf{W}_t$$



The Fokker Plank Equation

- What happens to the continuity equation if there is stochastic noise?

$$dx = \underbrace{u(x, t)dt}_{\substack{\text{Drift term} \\ \text{(deterministic)}}} + \underbrace{\sigma(x, t)d\mathbf{W}_t}_{\substack{\text{Diffusion term} \\ \text{(stochastic)}}}$$

Wiener process

- The ODE now becomes a *stochastic differential equation* (SDEs).

The Fokker Plank Equation

- What defines the Wiener process (aka Brownian motion)?

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\mathbf{W}_t$$

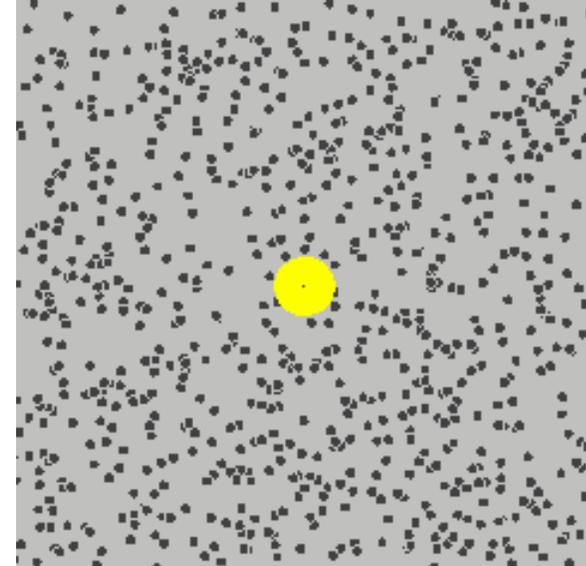
- Its increments are independent Gaussians.

$$\forall t, u > 0, s < t$$

$$(\mathbf{W}_{t+u} - \mathbf{W}_t) \sim \mathcal{N}(\mathbf{0}, u\mathbf{I})$$

$$(\mathbf{W}_{t+u} - \mathbf{W}_t) \perp \mathbf{W}_s$$

$$\mathbf{W}_0 = \mathbf{0}$$



The Fokker Plank Equation

- What defines the Wiener process (aka Brownian motion)?

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\mathbf{W}_t$$

- Increment in infinitesimal time interval is Gaussian.

$$d\mathbf{W}_t \sim \mathcal{N}(0, dt\mathbf{I})$$

$$\mathbf{W}_{t+\Delta t} - \mathbf{W}_t \approx \sqrt{\Delta t}\boldsymbol{\epsilon}$$

$$\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{1})$$

The Fokker Plank Equation

- How does $p(\mathbf{x}, t)$ change w.r.t. time if \mathbf{x} is governed by the SDE?

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\mathbf{W}_t$$

- This is given by the famous Fokker-Plank equation:

$$\frac{\partial}{\partial t}p(\mathbf{x}, t) = -\nabla \cdot [\mathbf{u}(\mathbf{x}, t)p(\mathbf{x}, t)] + \nabla^2 \cdot \left[\frac{\sigma^2(\mathbf{x}, t)}{2} p(\mathbf{x}, t) \right]$$

- Also known as the Kolmogorov forward equation.
- The initial distribution at $t = 0$ must be known.

The Fokker Plank Equation

- SDEs:

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt + \sigma(\mathbf{x}, t)d\mathbf{W}_t$$

$$\frac{\partial}{\partial t}p(\mathbf{x}, t) = -\nabla \cdot [\mathbf{u}(\mathbf{x}, t)p(\mathbf{x}, t)] + \nabla^2 \cdot \left[\frac{\sigma^2(\mathbf{x}, t)}{2} p(\mathbf{x}, t) \right]$$

- ODEs:

$$d\mathbf{x} = \mathbf{u}(\mathbf{x}, t)dt$$

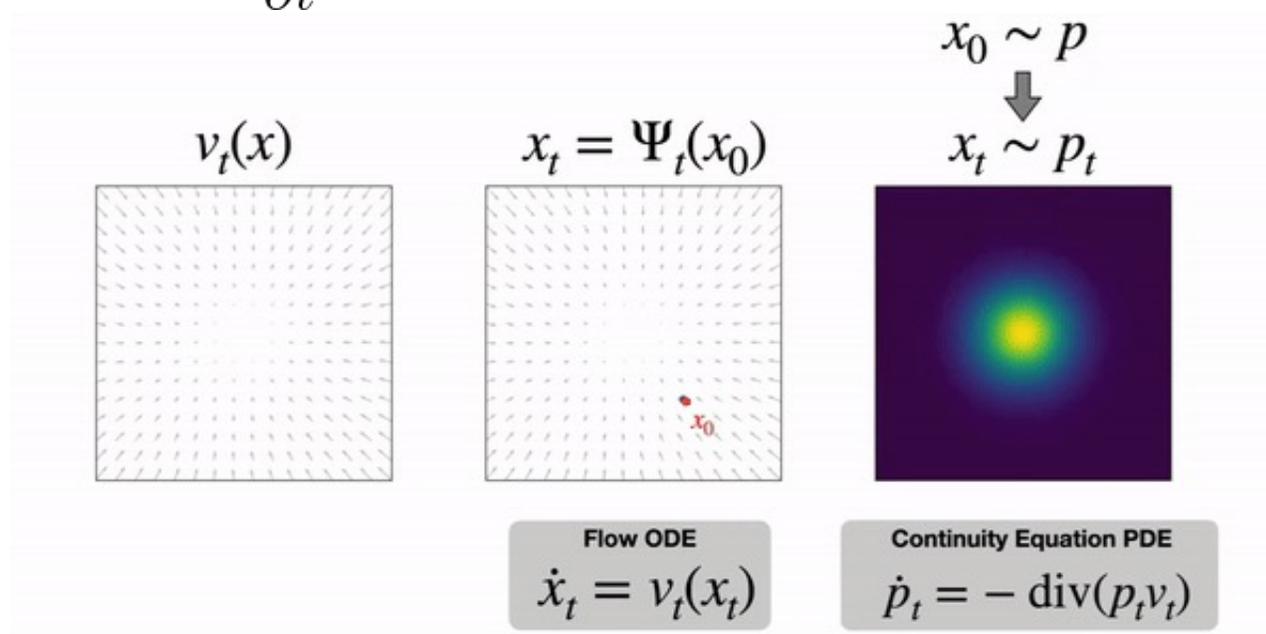
$$\frac{\partial}{\partial t}p(\mathbf{x}, t) = -\nabla \cdot [\mathbf{u}(\mathbf{x}, t)p(\mathbf{x}, t)]$$

Outline

- Normalizing Flows and Continuous Normalizing Flows
 - The Continuity Equation
- The Fokker Plank Equation
- **Flow matching**
- Variants:
 - Batch Optimal Transport Flow Matching

Flow Matching Model

- Motivation: Recall the continuity equation, which shows how the probability and the velocity field is coupled, since probability density is conserved as it flows through the space.

$$\frac{\partial p(z(t))}{\partial t} + \nabla \cdot (p(z(t))v_\theta(z(t), t)) = 0$$


$x_0 \sim p$
 \downarrow
 $x_t \sim p_t$

$v_t(x)$

$x_t = \Psi_t(x_0)$

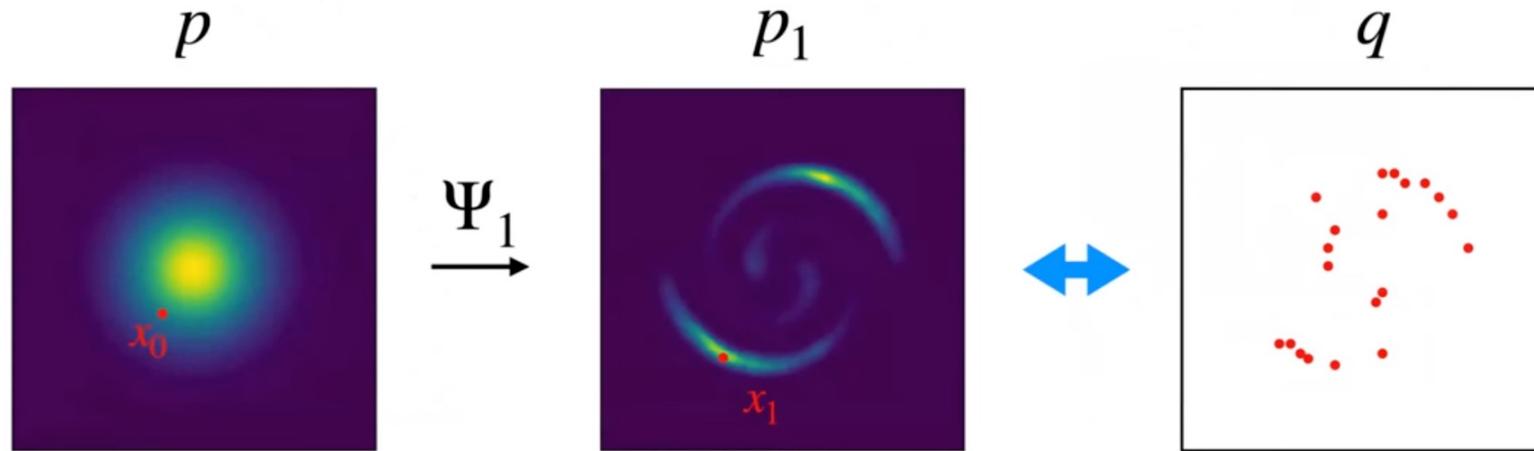
Flow ODE
 $\dot{x}_t = v_t(x_t)$

Continuity Equation PDE
 $\dot{p}_t = -\text{div}(p_t v_t)$

Flow Matching Model

- We want to have a loss for this generative model that is differentiable and tractable. Here we can try to minimize the distance between the target distribution p_1 and the data distribution q by minimizing the KL Divergence:

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



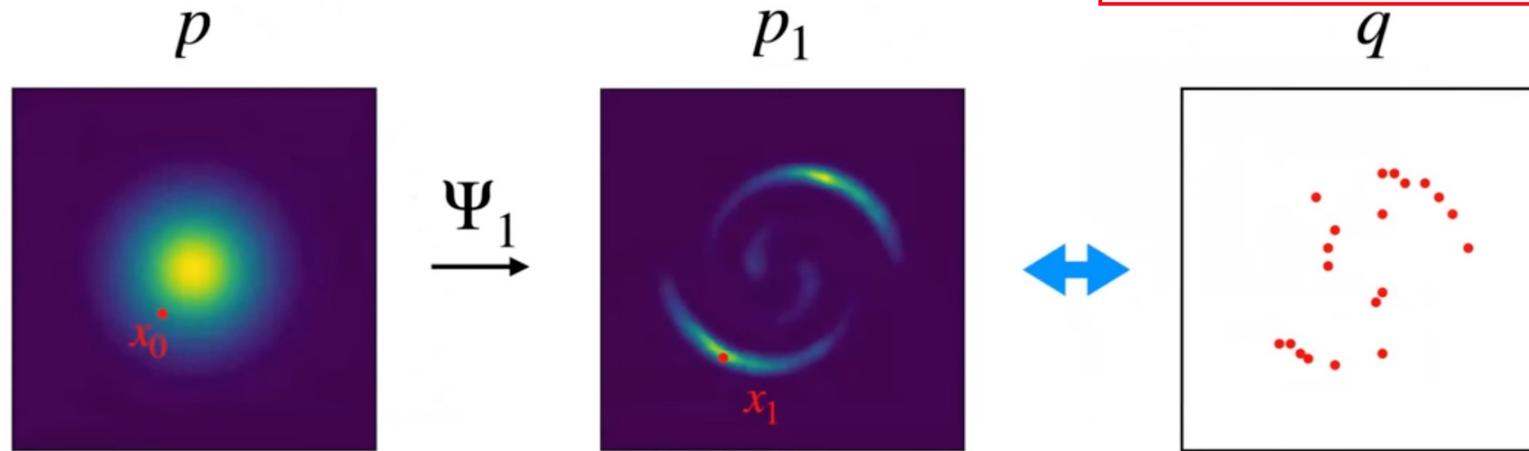
Flow Matching Model

- We want to have a loss for this generative model that is differentiable and tractable. Here we can try to minimize the distance between the target distribution p_1 and the data distribution q by minimizing the KL Divergence:

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

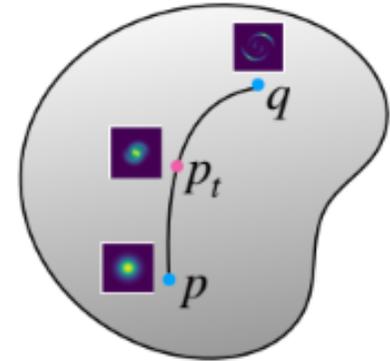
$$\frac{\partial \log p_t(\mathbf{z}(t))}{\partial t} = -\nabla \cdot \mathbf{v}_\theta(\mathbf{z}(t), t)$$

Need simulation if use instantaneous change of variable



Flow Matching Model

- So, revisit the continuity equation, we can observe that based on a known vector field, we could know how the density evolves.



Flow Matching Model

- Therefore, instead of directly optimizing the probability density path, we can optimize the vector field.

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t, p_t(\mathbf{x})} \|\mathbf{v}_t(\mathbf{x}) - \mathbf{u}_t(\mathbf{x})\|^2,$$

Flow Matching Model

- Therefore, instead of directly optimizing the probability density path, we can optimize the vector field.

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t, p_t(\mathbf{x})} \|\mathbf{v}_t(\mathbf{x}) - \mathbf{u}_t(\mathbf{x})\|^2,$$

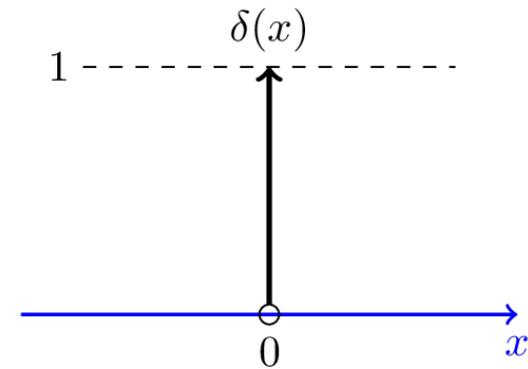
- However, we still cannot compute this loss because we don't know p_t or u_t .

Flow Matching Model

- We need to find a tractable loss.
- Assume we have samples from data distribution $q(\mathbf{x}_1)$, construct conditional probability paths $p_t(\mathbf{x}|\mathbf{x}_1)$, and marginalize the probability over data distribution:

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x}|\mathbf{x}_1)q(\mathbf{x}_1)d\mathbf{x}_1$$

$$\begin{cases} p_0(\mathbf{x}|\mathbf{x}_1) = p(\mathbf{x}_0) & t = 0 \\ p_1(\mathbf{x}|\mathbf{x}_1) \simeq \delta(\mathbf{x}_1) & t = 1 \end{cases}$$

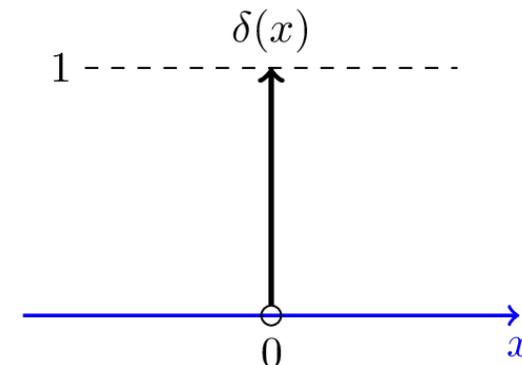


Flow Matching Model

- We need to find a tractable loss.
- Assume we have samples from data distribution $q(\mathbf{x}_1)$, construct conditional probability paths $p_t(\mathbf{x}|\mathbf{x}_1)$, and marginalize the probability over data distribution:

$$p_t(\mathbf{x}) = \int p_t(\mathbf{x}|\mathbf{x}_1)q(\mathbf{x}_1)d\mathbf{x}_1$$

$$\begin{cases} p_0(\mathbf{x}|\mathbf{x}_1) = p(\mathbf{x}_0) & t = 0 \\ p_1(\mathbf{x}|\mathbf{x}_1) \simeq \delta(\mathbf{x}_1) & t = 1 \end{cases}$$



- The conditional vector field is $u_t(\mathbf{x}|\mathbf{x}_1)$ and the marginal vector field is:

$$u_t(\mathbf{x}) = \int u_t(\mathbf{x}|\mathbf{x}_1) \frac{p_t(\mathbf{x}|\mathbf{x}_1)q(\mathbf{x}_1)}{p_t(\mathbf{x})} d\mathbf{x}_1$$

Flow Matching Model

Theorem 1. *Given vector fields $u_t(x|x_1)$ that generate conditional probability paths $p_t(x|x_1)$, for any distribution $q(x_1)$, the marginal vector field u_t in equation 8 generates the marginal probability path p_t in equation 6, i.e., u_t and p_t satisfy the continuity equation (equation 26).*

$$p_t(x) = \int p_t(x|x_1)q(x_1)dx_1, \quad (6)$$

$$u_t(x) = \int u_t(x|x_1) \frac{p_t(x|x_1)q(x_1)}{p_t(x)} dx_1, \quad (8)$$

Flow Matching Model

- The conditional flow matching loss:

$$\mathcal{L}_{CFM}(\theta) = \mathbb{E}_{t, q(\mathbf{x}_1), p_t(\mathbf{x}|\mathbf{x}_1)} \|v_t(\mathbf{x}) - u_t(\mathbf{x}|\mathbf{x}_1)\|^2$$

- Performing regression on conditional velocities has the same gradient as the flow matching loss.

$$\mathcal{L}_{FM}(\theta) = \mathbb{E}_{t, p_t(\mathbf{x})} \|v_t(\mathbf{x}) - u_t(\mathbf{x})\|^2,$$

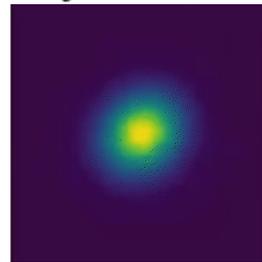
Theorem 2. *Assuming that $p_t(x) > 0$ for all $x \in \mathbb{R}^d$ and $t \in [0, 1]$, then, up to a constant independent of θ , \mathcal{L}_{CFM} and \mathcal{L}_{FM} are equal. Hence, $\nabla_{\theta} \mathcal{L}_{FM}(\theta) = \nabla_{\theta} \mathcal{L}_{CFM}(\theta)$.*

Flow Matching Model

Supervision p_t, u_t

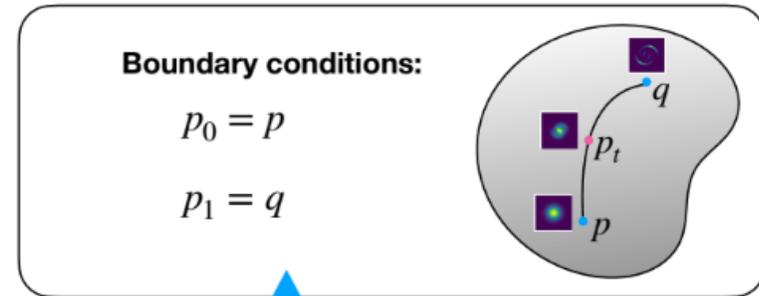
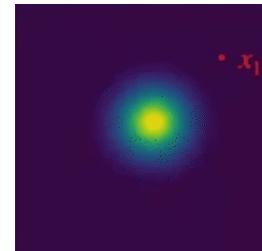
Marginal path

$$p_t(x) = \int p_t(x|x_1)q(x_1)dx_1$$



Conditional path

$$p_t(x|x_1)$$



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

$$p_1(\cdot | x_1) = \delta_{x_1}$$

Flow Matching Model

- The conditional flow based on the conditional vector field is:

$$\frac{d}{dt} \Psi_t(\mathbf{x}) = u_t(\Psi_t(\mathbf{x}) | \mathbf{x}_1)$$

Flow Matching Model

- The conditional flow based on the conditional vector field is:

$$\frac{d}{dt}\Psi_t(\mathbf{x}) = u_t(\Psi_t(\mathbf{x})|\mathbf{x}_1)$$

- Simply, let $p_t(\mathbf{x}|\mathbf{x}_1) = N(\mathbf{x}|\mu_t(\mathbf{x}), \sigma_t(\mathbf{x})^2 I)$ and $\Psi_t(\mathbf{x}) = \sigma_t(\mathbf{x}_1)\mathbf{x} + u_t(\mathbf{x}_1)$, where $\sigma_0(\mathbf{x}_1) = 1, \sigma_1(\mathbf{x}_1) = \sigma_{min}, \mu_0(\mathbf{x}_1) = 0, \mu_1(\mathbf{x}_1) = \mathbf{x}_1$. Here we introduce σ_{min} to mimic $\delta(\mathbf{x}_1)$ without losing smoothness when t is close to 1.

$$\Psi_t(\mathbf{x}) = \sigma_t(\mathbf{x}_1)\mathbf{x} + u_t(\mathbf{x}_1)$$

Flow Matching Model

- Specifically, the mean and standard deviation change linearly with time:

$$\mu_t(\mathbf{x}) = t\mathbf{x}_1, \quad \text{and} \quad \sigma_t(\mathbf{x}) = 1 - (1 - \sigma_{min})t$$

Flow Matching Model

- Specifically, the mean and standard deviation change linearly with time:

$$\mu_t(\mathbf{x}) = t\mathbf{x}_1, \quad \text{and} \quad \sigma_t(\mathbf{x}) = 1 - (1 - \sigma_{min})t$$

- This gives a straight path:

$$u_t(\mathbf{x}|\mathbf{x}_1) = \frac{\mathbf{x}_1 - (1 - \sigma_{min})\mathbf{x}}{1 - (1 - \sigma_{min})t}$$

Flow Matching Model

- The conditional flow with optimal transport is:

$$\Psi_t(\mathbf{x}) = [1 - (1 - \sigma_{min})t]\mathbf{x} + t\mathbf{x}_1$$

Flow Matching Model

- The conditional flow with optimal transport is:

$$\Psi_t(\mathbf{x}) = [1 - (1 - \sigma_{min})t]\mathbf{x} + t\mathbf{x}_1$$

- The reparametrized conditional flow matching loss is:

$$\mathbb{E}_{t,q(\mathbf{x}_1),p(\mathbf{x}_0)} \|\mathbf{v}_t(\Psi_t(\mathbf{x}_0|\mathbf{x}_1)) - (\mathbf{x}_1 - (1 - \sigma_{min})\mathbf{x}_0)\|^2.$$

Flow Matching Model

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

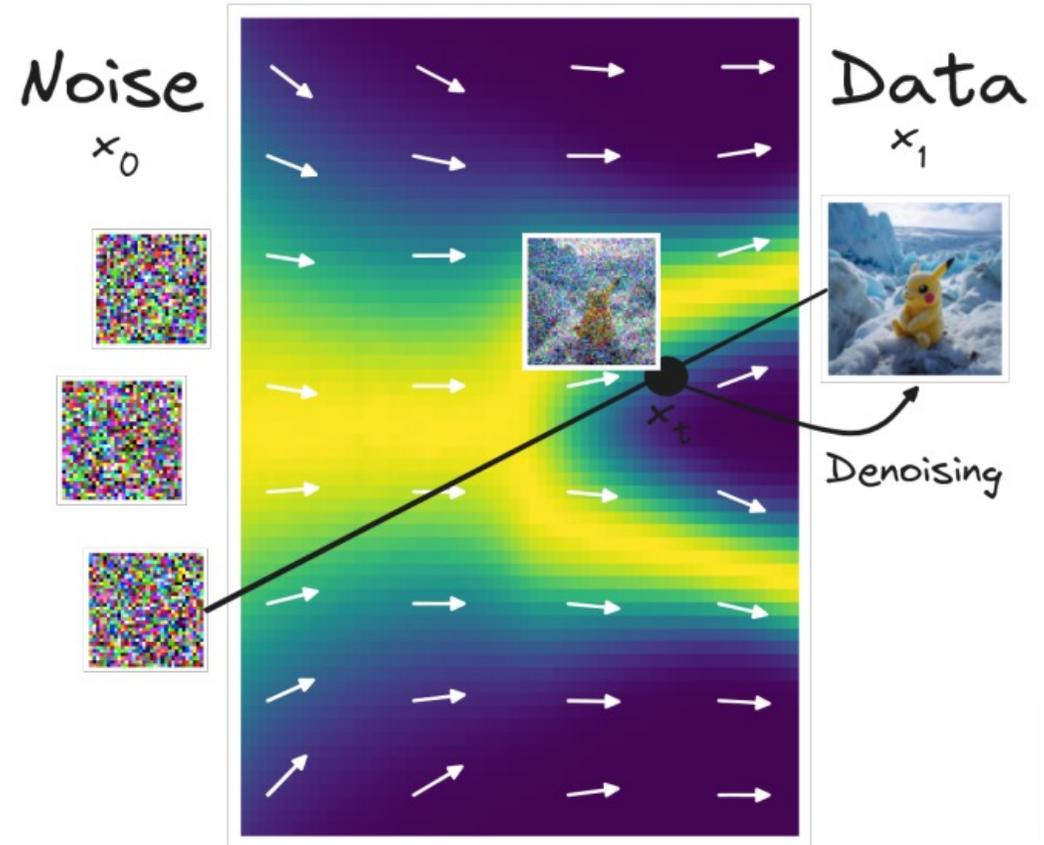
Algorithm 2: Flow Matching sampling.

Input: trained model v^θ

$x_0 \sim p(x_0)$ // sample noise

Numerically solve ODE $\dot{x}_t = v_t^\theta(x_t)$

Output: x_1



Flow Matching Model vs. Diffusion Model

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

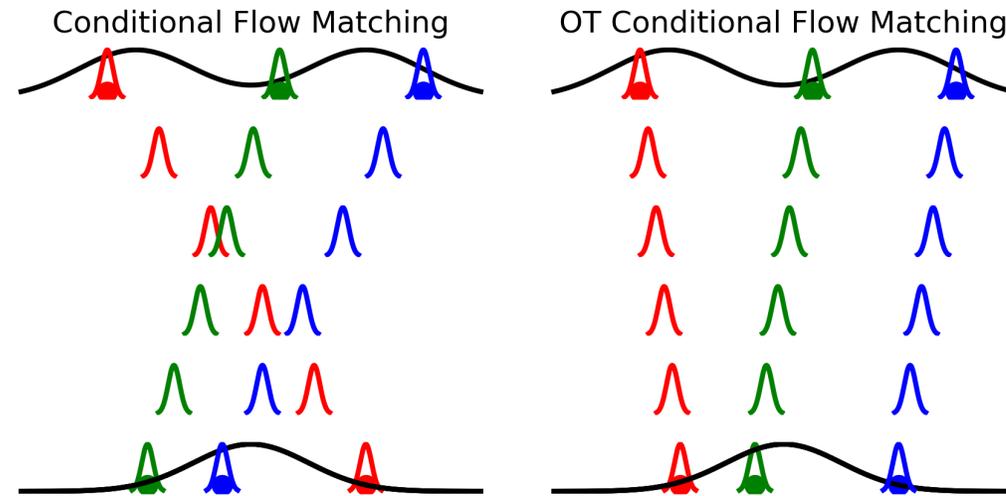
$p(x_0)$ is Gaussian

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

Outline

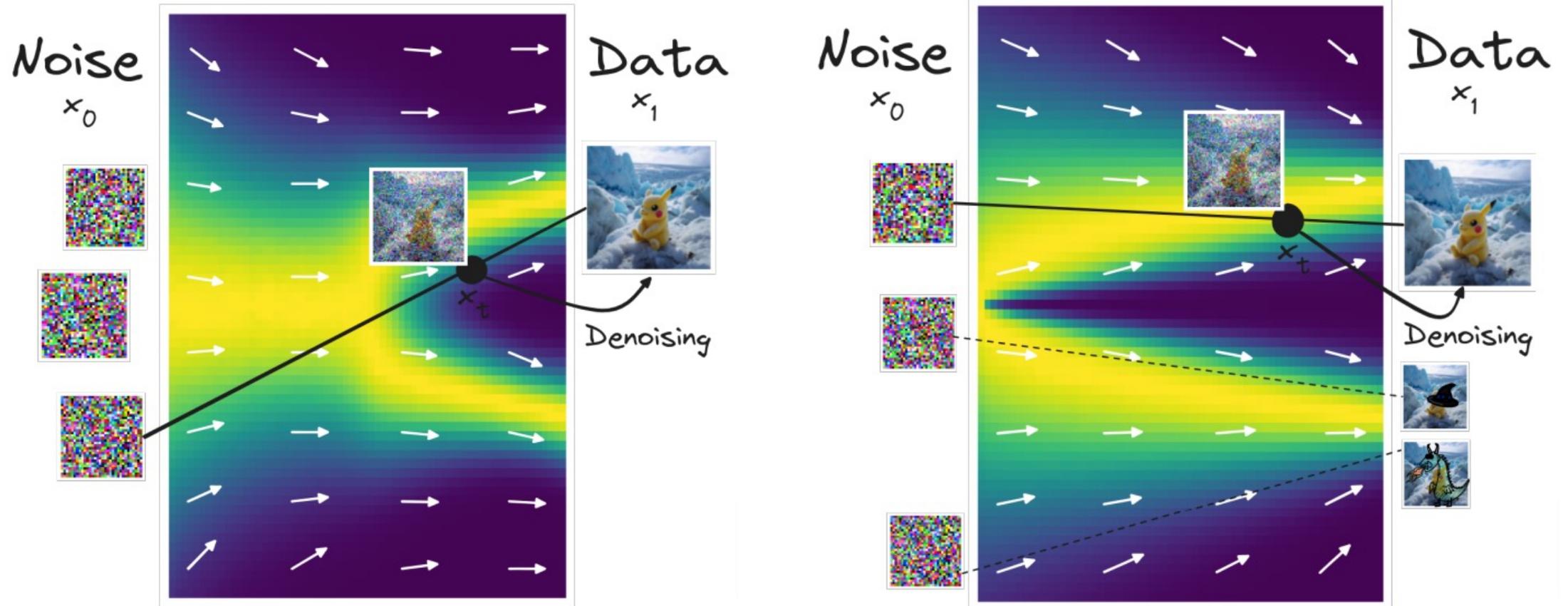
- Normalizing Flows and Continuous Normalizing Flows
 - The Continuity Equation
- The Fokker Plank Equation
- Flow matching
- **Variants:**
 - **Batch Optimal Transport Flow Matching**

Mini Batch OT Flow Matching Model



- Paths of various flow matching model design
 - Vanilla Conditional Flow Matching: Each conditional path is straight, but some paths intersect.
 - OT Conditional Flow Matching: Within the mini-batch, all paths are assigned as non-intersecting straight lines.

Mini Batch OT Flow Matching Model



Questions?