

# Training flow and Diffusion Models

Reminder - Marginal vector field:

$$u_t^{\text{target}}(x) = \int u_t^{\text{target}}(x|z) \frac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)} dz$$

*Conditional vector field*



*High computation cost!!*

**Marginalization trick:** The marginal vector field defines an ODE that "follows" the marginal probability path:

$$X_0 \sim p_{\text{init}}, \quad \frac{d}{dt} X_t = u_t^{\text{target}}(X_t) dt \quad \Rightarrow X_t \sim p_t$$

This means that the final point fulfills:  $X_1 \sim p_{\text{data}}$

3.1 Flow Matching.

$$u_t^\theta$$

Goal:  $u_t^\theta \simeq u_t^{\text{target}}$  (marginal VF)  $\Rightarrow$  Learn marginal VF.

Flow matching loss:  $L_{fm}(\theta) = \mathbb{E} \left[ \|u_t^\theta(x) - u_t^{\text{target}}(x)\|_2^2 \right]$

$t \sim \text{Unif}$ :  $t$  is uniform in  $[0, T]$ .

$x \sim p_{\text{data}}$ : draw samples from data

$x \sim p_t(\cdot|z)$ : draw samples from  $p(\cdot|z)$  for fixed  $z$ .

We want: minimizer of  $L_{fm}$ .  
 But This loss is not tractable!  
*Because  $u_t^{\text{target}}(x)$  is an integral.  
 Hard to compute.*

Conditional Flow matching loss:

$$L_{cfm}(\theta) = \mathbb{E} \left[ \|u_t^\theta(x) - u_t^{\text{target}}(x|z)\|_2^2 \right]$$

① Tractable.

② But minimize CFM loss  $t \sim \text{Unif}$ :  $t$  is uniform in  $[0, T]$ .

Will help minimize FM loss?  $z \sim p_{\text{data}}$ : draw samples from data  
 $x \sim p_t(\cdot|z)$ : draw samples from  $p(\cdot|z)$  for fixed  $z$ .

Main Thm: Minimize  $L_{\text{cfm}}(\theta) \iff \text{Minimize } L_{\text{fm}}(\theta)$ .

Conditional Prob. Path, Vector Field, and Score		
Notation	Key property	Gaussian example
Conditional Probability Path	$p_t(\cdot z)$	Interpolates $p_{\text{init}}$ and a data point $z$
Conditional Vector Field	$u_t^{\text{target}}(x z)$	ODE follows conditional path
Conditional Score Function	$\nabla \log p_t(x z)$	Gradient of log-likelihood

Can be derived directly if we set  $a_t, b_t$  in advance.

Thm 4.1  $L_{\text{fm}}(\theta) = L_{\text{cfm}}(\theta) + C$ , where  $C > 0$  is independent of  $\theta$ .

① For minimizer  $\theta^*$  of  $L_{\text{cfm}}$ :  $u_t^{\theta^*} = u_t^{\text{target}}$

②  $\nabla_{\theta} L_{\text{cfm}}(\theta) = \nabla_{\theta} L_{\text{fm}}(\theta)$   
 $\Rightarrow$  SGD process will be exactly the same.

### Algorithm 3 Flow Matching Training Procedure (General)

**Require:** A dataset of samples  $z \sim p_{\text{data}}$ , neural network  $u_t^{\theta}$

- 1: **for** each mini-batch of data **do** *SetD.*
- 2: Sample a data example  $z$  from the dataset.
- 3: Sample a random time  $t \sim \text{Unif}_{[0,1]}$ . Another idea: not uniform time sampling.
- 4: Sample  $x \sim p_t(\cdot|z)$
- 5: Compute loss

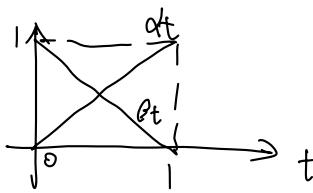
$$\mathcal{L}(\theta) = \|u_t^{\theta}(x) - u_t^{\text{target}}(x|z)\|^2$$

- 6: Update the model parameters  $\theta$  via gradient descent on  $\mathcal{L}(\theta)$
- 7: **end for**

nn. Architecture: Takes input  $(x, t) \rightarrow$  output  $u_t^{\theta}(x)$ .  
 $\theta$   
(param by  $\theta$ )

L<sub>cmf</sub> for Gaussian cond path

$$P_t(\cdot | z) = \mathcal{N}(\alpha_t z, \beta_t^2 \text{Id})$$



$$\hat{u}_t^{\text{target}}(x|z) = \left( \dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z + \frac{\dot{\beta}_t}{\beta_t} x$$

$$\epsilon \sim \mathcal{N}(0, \text{Id}) \Rightarrow \alpha_t z + \beta_t \epsilon \stackrel{def}{=} x \sim f_t(\cdot | z).$$

Then  $L_{\text{cmf}}(0) = \mathbb{E}_{\substack{t \sim \text{unif} \\ z \sim \text{pdata} \\ x \sim N(\alpha_t z, \beta_t^2 \text{Id})}} \left[ \| \hat{u}_t^0(x) - \left( \dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t \right) z - \frac{\dot{\beta}_t}{\beta_t} x \|_2^2 \right]$

Replace  $\underset{\parallel}{x} = \mathbb{E}_{\substack{t \sim \text{unif} \\ z \sim \text{data} \\ \epsilon \sim \mathcal{N}(0, \text{Id})}} \left[ \| \hat{u}_t^t(\alpha_t z + \beta_t \epsilon) - (\alpha_t z + \beta_t \epsilon) \|_2^2 \right]$

A easy choice for  $\alpha_t, \beta_t$ :  $\begin{cases} \alpha_t = t & \dot{\alpha}_t = 1 \\ \beta_t = 1-t & \dot{\beta}_t = -1 \end{cases}$

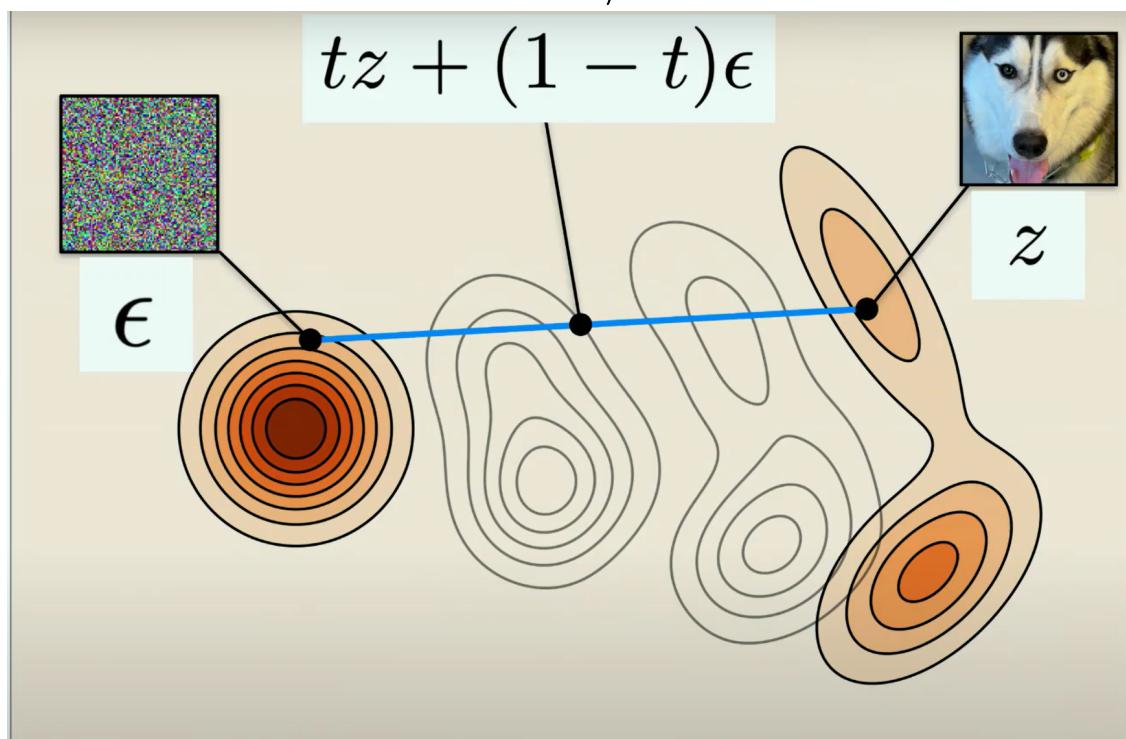
Cond OT path.

(Optimal Transport Path)

$$\epsilon \rightarrow z$$

A straight line  
between noise & data

$$tz + (1-t)\epsilon$$



#### Algorithm 4 Flow Matching Training for CondOT path

**Require:** A dataset of samples  $z \sim p_{\text{data}}$ , neural network  $u_t^\theta$

- 1: **for** each mini-batch of data **do**
- 2:   Sample a data example  $z$  from the dataset.
- 3:   Sample a random time  $t \sim \text{Unif}_{[0,1]}$ .
- 4:   Sample noise  $\epsilon \sim \mathcal{N}(0, I_d)$  ✓. *Easy to op.*
- 5:   Set  $x = tz + (1-t)\epsilon$
- 6:   Compute loss

$$\mathcal{L}(\theta) = \|u_t^\theta(x) - (z - \epsilon)\|^2$$

- 7:   Update the model parameters  $\theta$  via gradient descent on  $\mathcal{L}(\theta)$ .
- 8: **end for**

Stable diffusion: still use this training algo.

proof of Thm 4.1.

$$L_{\text{fm}}(\theta) = \mathbb{E}_{t \geq 0, x} [\|u_t^\theta(x) - u_t^{\text{target}}(x)\|_2^2] = \mathbb{E} [\|u_t^\theta(x)\|_2^2 + \|u_t^{\text{target}}(x)\|_2^2 - 2 u_t^\theta(x)^T u_t^{\text{target}}(x)]$$

$$L_{\text{cfm}}(\theta) = \mathbb{E} [\|u_t^\theta(x)\|_2^2 + \|u_t^{\text{target}}(x|z)\|_2^2 - 2 u_t^\theta(x)^T u_t^{\text{target}}(x|z)].$$

Denote  $C = \mathbb{E} [\|u_t^{\text{target}}(x)\|_2^2 - \|u_t^{\text{target}}(x|z)\|_2^2]$ . ... A constant.

$$L_{\text{fm}}(\theta) - L_{\text{cfm}}(\theta) = 2 \mathbb{E} [u_t^\theta(x)^T u_t^{\text{target}}(x|z) - u_t^\theta(x)^T \underline{u_t^{\text{target}}(x)}] + C$$

$$u_t^{\text{target}}(x) = \int u_t^{\text{target}}(x|z) \frac{p_t(x|z) p_{\text{data}}(z)}{p_t(x)} dz$$

$$\mathbb{E} [u_t^\theta(x)^T u_t^{\text{target}}(x)] = \mathbb{E} [u_t^\theta(x)^T \int_z u_t^{\text{target}}(x|z) \frac{p_t(x|z) p_{\text{data}}(z)}{p_t(x)} dz]$$

$$= \int_x \int_z \int_z p_t(x) \underbrace{u_t^\theta(x)^T u_t^{\text{target}}(x|z)}_{\frac{p_t(x|z) p_{\text{data}}(z)}{p_t(x)}} dz dx dt$$

$$= \mathbb{E} [u_t^\theta(x)^T u_t^{\text{target}}(x|z)].$$

## 4.2 Score Matching.

Reminder - Marginal score function:



$$\nabla \log p_t(x) = \int \nabla \log p_t(x|z) \frac{p_t(x|z)p_{\text{data}}(z)}{p_t(x)} dz$$

*Conditional score function*

**SDE extension trick:** The marginal score function allows to extend the ODE to a SDE with arbitrary diffusion coefficient:

$$X_0 \sim p_{\text{init}}, \quad dX_t = \left[ u_t^{\text{target}}(X_t) + \frac{\sigma_t^2}{2} \nabla \log p_t(X_t) \right] dt + \sigma_t dW_t$$

Already solved by FM  $\Rightarrow X_t \sim p_t$  Train these 2.

$$\begin{aligned} \nabla \log p_t(x) &= \frac{\nabla p_t(x)}{p_t(x)} = \frac{\nabla \int p_t(x|z) p_{\text{data}}(z) dz}{p_t(x)} \\ &= \underbrace{\int \nabla \log p_t(x|z) \frac{p_t(x|z) p_{\text{data}}(z)}{p_t(x)} dz}_{\text{Also hard to directly compute } \nabla \log p_t(x).} \end{aligned}$$

Score Network:  $s_t^\theta$  ( $\theta$  parameters) goal:  $s_t^\theta \approx \nabla \log p_t$ .

$$\begin{cases} \text{Score Matching loss: } L_{\text{sm}}(\theta) = \mathbb{E}_{t,z,x} [\| s_t^\theta(x) - \nabla \log p_t(x) \|_2^2] \\ \text{Denoising Matching loss: } L_{\text{dsm}}(\theta) = \mathbb{E}_{t,z,x} [\| s_t^\theta(x) - \nabla \log p_t(x|z) \|_2^2] \end{cases}$$

Similar to FM,  
SM loss: ✗ tractable  
DSM loss: ✓ tractable.

Thm 4.2.  $L_{\text{sm}}(\theta) = L_{\text{dsm}}(\theta) + C$ . for  $C \ll \epsilon$  of  $\theta$ .

① for  $\theta^*$  of  $L_{\text{dsm}}$ .  $s_t^{\theta^*} = \nabla \log p_t$ .

$$\Leftrightarrow \nabla_{\theta} L_{sm}(\theta) = \nabla_{\theta} L_{dm}(\theta).$$

The pf is exactly the same as DM part.

### Algorithm 6 Score Matching Training Procedure (General)

**Require:** A dataset of samples  $z \sim p_{\text{data}}$ , score network  $s_t^\theta$

- 1: **for** each mini-batch of data **do**
- 2:    Sample a data example  $z$  from the dataset.
- 3:    Sample a random time  $t \sim \text{Unif}_{[0,1]}$ .
- 4:    Sample  $x \sim p_t(\cdot|z)$
- 5:    Compute loss

$$\mathcal{L}(\theta) = \|s_t^\theta(x) - \nabla \log p_t(x|z)\|^2$$

- 6:    Update the model parameters  $\theta$  via gradient descent on  $\mathcal{L}(\theta)$
- 7: **end for**

Eg. For Gaussian prob path,

$$\nabla \log p_t(x|z) = -\frac{x - \alpha_t z}{\beta_t^2}$$

$$L_{sm}(\theta) = \mathbb{E}_{\substack{t \sim \text{Unif} \\ z \sim \text{data} \\ x \sim \mathcal{N}(\alpha_t z, \beta_t^2 I_d)}} \left[ \|s_t^\theta(x) + \frac{x - \alpha_t z}{\beta_t^2}\|_2^2 \right]$$

$$x = \alpha_t z + \beta_t \Sigma.$$

$$\Sigma \sim \mathcal{N}(0, I_d)$$

$$= \mathbb{E}_{\substack{t \sim \text{Unif} \\ z \sim \text{data} \\ \Sigma \sim \mathcal{N}(0, I_d)}} \left[ \|s_t^\theta(\alpha_t z + \beta_t \Sigma) + \frac{\Sigma}{\beta_t}\|_2^2 \right]$$

### Algorithm 5 Score Matching Training Procedure for Gaussian probability

**Require:** A dataset of samples  $z \sim p_{\text{data}}$ , score network  $s_t^\theta$  or noise predictor  $e_t^\theta$

**Require:** Schedulers  $\alpha_t, \beta_t$  with  $\alpha_0 = \beta_1 = 0, \alpha_1 = \beta_0 = 1$

- 1: **for** each mini-batch of data **do**
- 2:    Sample a data example  $z$  from the dataset.
- 3:    Sample a random time  $t \sim \text{Unif}_{[0,1]}$ .
- 4:    Sample noise  $\epsilon \sim \mathcal{N}(0, I_d)$
- 5:    Set  $x_t = \alpha_t z + \beta_t \epsilon$
- 6:    Compute loss

$$\mathcal{L}(\theta) = \|s_t^\theta(x_t) + \frac{\epsilon}{\beta_t}\|^2$$

Numerically unstable for low beta

- 7:    Update the model parameters  $\theta$  via gradient descent on  $\mathcal{L}(\theta)$ .
- 8: **end for**

For Gaussian prob path with  $\beta_t = 1 - t_j$

$$f_t \rightarrow 0 \text{ as } t \rightarrow 1. \quad \frac{\epsilon}{f_t} \rightarrow \infty ?$$

To deal with this problem, DDPM was proposed. ① Prop  $\hat{p}_t$  ②  $L_t(x) = -\hat{p}_t \hat{s}_t^0(x)$ .

Let  $E_t^{\theta}(x) = -f_t S_t^{\theta}(x)$ , a noisy predictor.  $S_t^{\theta}$  is another nn.

$$E^t : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}$$

$$\mathcal{L}_{DDPM}(\theta) = \mathbb{E}_{\substack{t \sim \text{Unif} \\ z \sim \text{Pdata.} \\ \varepsilon \sim \mathcal{N}(0, \text{Id})}} \left[ \left\| \mathbb{E}_{\theta}^{\hat{p}_t} (\phi_t z + \beta_t \varepsilon) - \varepsilon \right\|_2^2 \right].$$

---

**Algorithm 4** Score Matching Training Procedure for Gaussian probability path

**Require:** A dataset of samples  $z \sim p_{\text{data}}$ , score network  $s_t^\theta$  or noise predictor  $\epsilon_t^\theta$

- 1:** **for** each mini-batch of data **do**

  - 2:** Sample a data example  $z$  from the dataset.
  - 3:** Sample a random time  $t \sim \text{Unif}_{[0,1]}$ .
  - 4:** Sample noise  $\epsilon \sim \mathcal{N}(0, I_d)$
  - 5:** Set  $x_t = \alpha_t z + \beta_t \epsilon$
  - 6:** Compute loss

(General case:  $x_t \sim p_t(\cdot|z)$ )

$$\mathcal{L}(\theta) = \|s_t^\theta(x_t) + \frac{\epsilon}{\beta_t}\|^2 \quad (\text{General case: } = \|s_t^\theta(x_t) - \nabla \log p_t(x_t|z)\|^2)$$

Alternatively:  $\mathcal{L}(\theta) = \|\epsilon_t^\theta(x_t) - \epsilon\|^2$  *when  $t \rightarrow 1$ .*

- 7:   Update the model parameters  $\theta$  via gradient descent on  $\mathcal{L}(\theta)$ .  
8: **end for**

# Stochastic sampling with diffusion models



$$X_0 \sim p_{\text{init}}, \quad dX_t = \left[ u_t^{\text{target}}(X_t) + \frac{\sigma_t^2}{2} \nabla \log p_t(X_t) \right] dt + \sigma_t dW_t$$

## *Plugin neural networks*

$$X_0 \sim p_{\text{init}}, \quad dX_t = \left[ u_t^\theta(X_t) + \frac{\sigma_t^2}{2} s_t^\theta(X_t) \right] dt + \sigma_t dW_t$$

*After training*  $\Rightarrow X_t \sim p_t$



## Denoising diffusion models

= Diffusion models with a Gaussian probability path

(our standard example)  $\mathcal{N}(\alpha_t z, \beta_t^2 I_d)$

### Terminology (by many people)

**Denoising diffusion model = diffusion model**

i.e. many people drop the word “denoising” and just say “diffusion model”.

Special property for DDM:

$$\text{Obs: } \begin{cases} u_t^{\text{target}}(x|z) = (\dot{\alpha}_t - \frac{\dot{\beta}_t}{\beta_t} \alpha_t) z + \frac{\dot{\beta}_t}{\beta_t} x \\ \nabla \log p_t(x|z) = -\frac{x - \alpha_t z}{\beta_t^2} = -\frac{\dot{\alpha}_t}{\beta_t^2} z - \frac{1}{\beta_t^2} x \end{cases} \quad \begin{array}{l} \text{Both comb of } z \text{ & } x. \\ \Rightarrow \text{Convert them to each other.} \end{array}$$

Easy to show:

$$u_t^{\text{target}}(x|z) = \left( \beta_t^2 \frac{\dot{\alpha}_t}{\alpha_t} - \dot{\beta}_t \beta_t \right) \nabla \log p_t(x|z) + \frac{\dot{\alpha}_t}{\alpha_t} x.$$

$$u_t^{\text{target}}(x) = \left( \beta_t^2 \frac{\dot{\alpha}_t}{\alpha_t} - \dot{\beta}_t \beta_t \right) \nabla \log p_t(x) + \frac{\dot{\alpha}_t}{\alpha_t} x.$$

### Conversion formula for DDMs



One can also convert the score network into vector field network post training:

$$u_t^\theta = \left( \beta_t^2 \frac{\dot{\alpha}_t}{\alpha_t} - \dot{\beta}_t \beta_t \right) s_t^\theta(x) + \frac{\dot{\alpha}_t}{\alpha_t} x$$

In denoising diffusion models, we don't have to train the vector field and score network separately

→ they can be converted into one another post-training

The first generation of diffusion models only did score matching!

We only need to train one of  $u_t^\theta$  &  $s_t^\theta(x)$ .

Next steps: nn. Architectures.

conditioned on a prompt?

Brief literature overview (see lecture notes for details)

- There is a whole zoo of different flow and diffusion model (Don't be confused!)
- First diffusion models: **Discrete time** (no ODE/SDEs)
- First SDE-based model: Forward-backward process using **time-reversal** of SDEs → construction of a probability path
- **Inverted time convention**: Data distribution at t=0

Here: Flow matching and stochastic interpolants:

- Arguably most simple flow and diffusion algorithms
- Allows you to restrict yourself to flows
- Allows you go from arbitrary  $p_{\text{init}}$  to arbitrary  $p_{\text{data}}$

**Note: Our method allows to convert arbitrary distributions into arbitrary distributions!**

We not always start from Gaussian. The distribution can be arbitrary.

### Bridging arbitrary distributions - Example

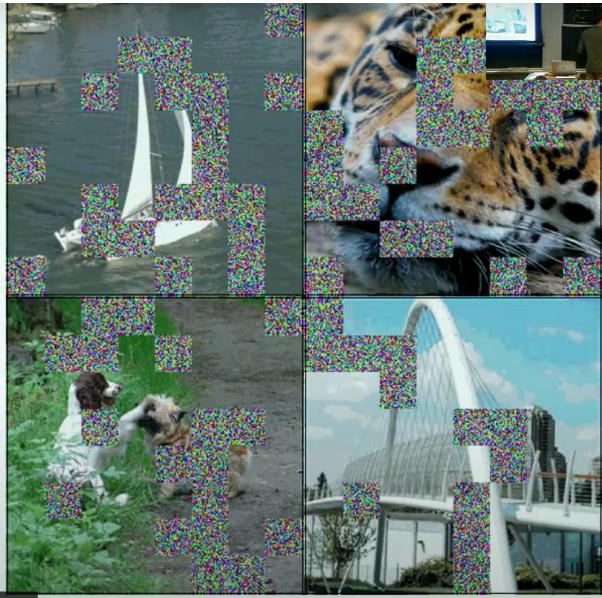
Videos without audio → videos with audio

Low resolution images → high resolution images

Unperturbed cells → perturbed cells

etc.

noise that's  
for what we write



These are not pure Gaussian noise at start.