

# InstaFlow:

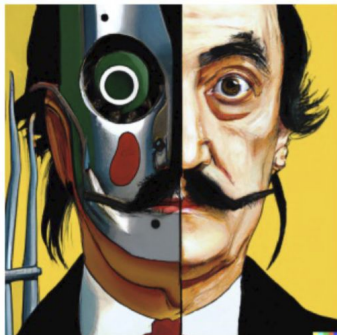
## One Step is Enough for High-Quality Diffusion-Based Text-to-Image Generation

Authors:

Xingchao Liu, Xiwen Zhang, Jianzhu Ma, Jian Peng, Qiang Liu

Presenter:

James Soole



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it

## Text-to-Image DALL-E 2 (OpenAI)



## Inpainting RePaint

- Denoising
- Super Resolution
- ....



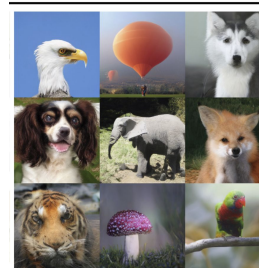
Input images



## Stylized Images

LoRA on SDXL (UIUC, Google)

ImageNet



## Unconditional Generation Latent Diffusion



## Customizing Images DreamBooth (Google)

# Background: Diffusion Summary

## Diffusion

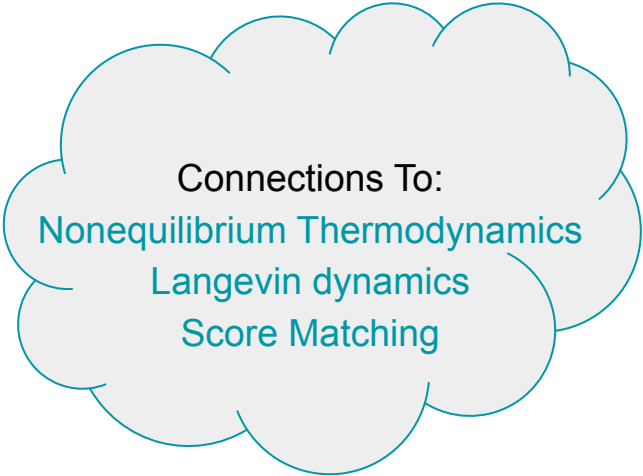
- [Diffusion Probabilistic Models](#) (2015, Sohl-Dickstein et al.)
- [Denoising Diffusion Probabilistic Models](#) (2019, Ho et al.)

## Latent Diffusion

- [Latent Diffusion](#) (Rombach et al., 2021)
  - → Stable Diffusion

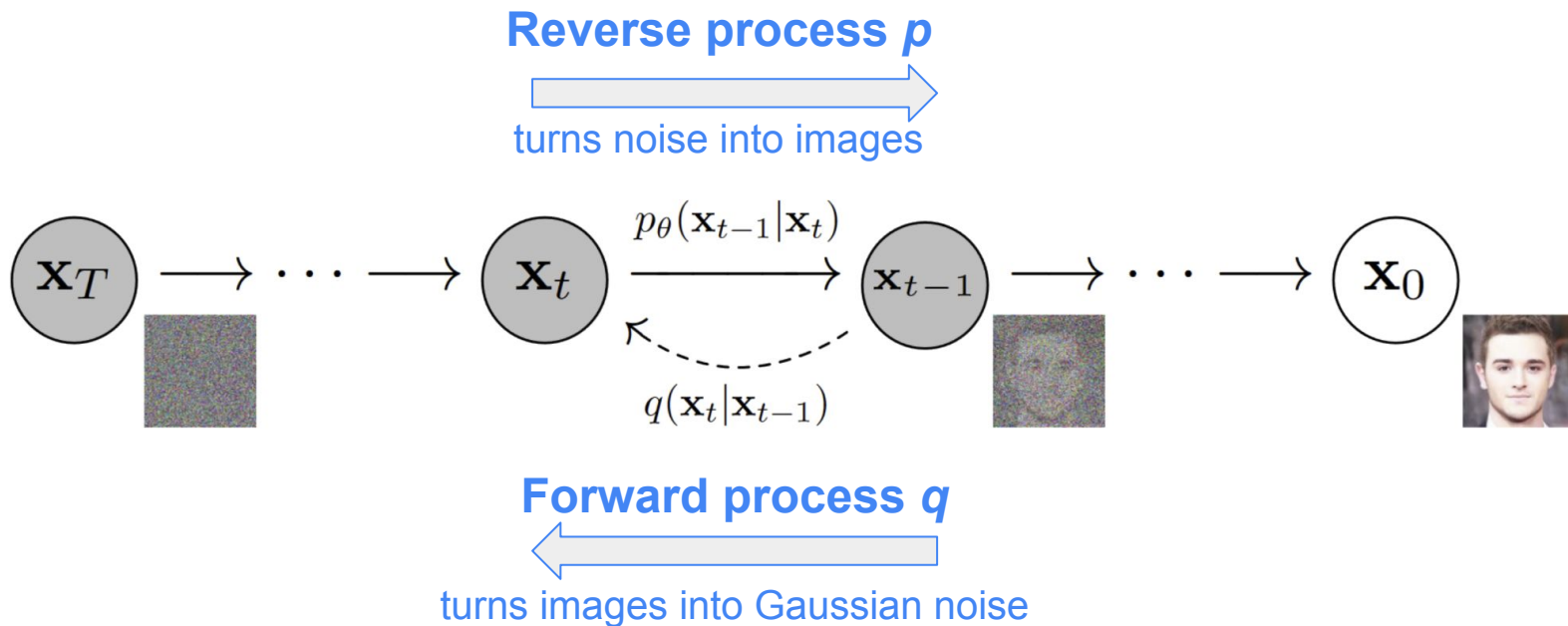
## Even More Diffusion

- [InstaFlow](#) (this)
- and many more ...



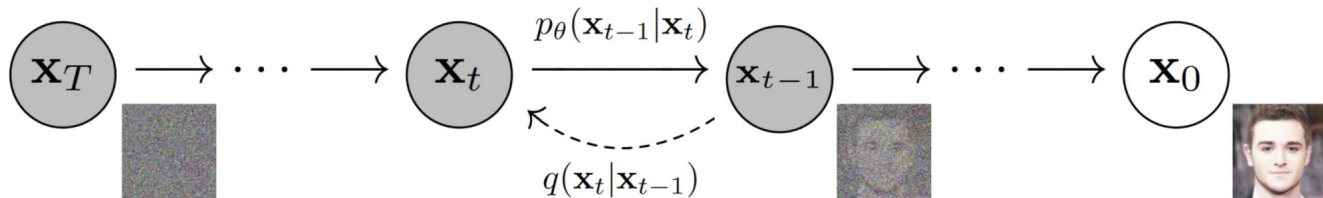
Connections To:  
Nonequilibrium Thermodynamics  
Langevin dynamics  
Score Matching

# Background: Diffusion



- Provided transitions  $t$  are small enough,  $p_\theta(\mathbf{x}_{t-1}|\mathbf{x}_t)$  is gaussian
- we can train a neural network to estimate  $\mathbf{x}_{t-1}$  from  $\mathbf{x}_t$

# Background: Diffusion



## Forward Process ( $\leftarrow$ )

- Gradually add gaussian noise
  - $\beta$  varies by timestep (from  $10^{-4}$  to 0.02)
  - until the image approaches gaussian noise
- Can sample in closed form

$$q(\mathbf{x}_t | \mathbf{x}_{t-1}) = \mathcal{N}(\mathbf{x}_t; \sqrt{1 - \beta_t}\mathbf{x}_{t-1}, \beta_t\mathbf{I})$$

constant

## Reverse Process ( $\rightarrow$ )

- Gradually predict the added noise
  - our trained model predicts noise **mean** and **variance**
- Reconstruct image by subtracting predicted noise

$$p_\theta(\mathbf{x}_{t-1} | \mathbf{x}_t) = \mathcal{N}(\mathbf{x}_{t-1}; \boldsymbol{\mu}_\theta(\mathbf{x}_t, t), \boldsymbol{\Sigma}_\theta(\mathbf{x}_t, t))$$

predicted

# Background: Diffusion

---

## Algorithm 1 Training

---

- 1: **repeat**
- 2:  $\mathbf{x}_0 \sim q(\mathbf{x}_0)$
- 3:  $t \sim \text{Uniform}(\{1, \dots, T\})$
- 4:  $\boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 5: Take gradient descent step on  
$$\nabla_{\theta} \left\| \boldsymbol{\epsilon} - \boldsymbol{\epsilon}_{\theta}(\sqrt{\bar{\alpha}_t} \mathbf{x}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, t) \right\|^2$$
- 6: **until** converged

add some amount of noise

predict noise

---

## Algorithm 2 Sampling

---

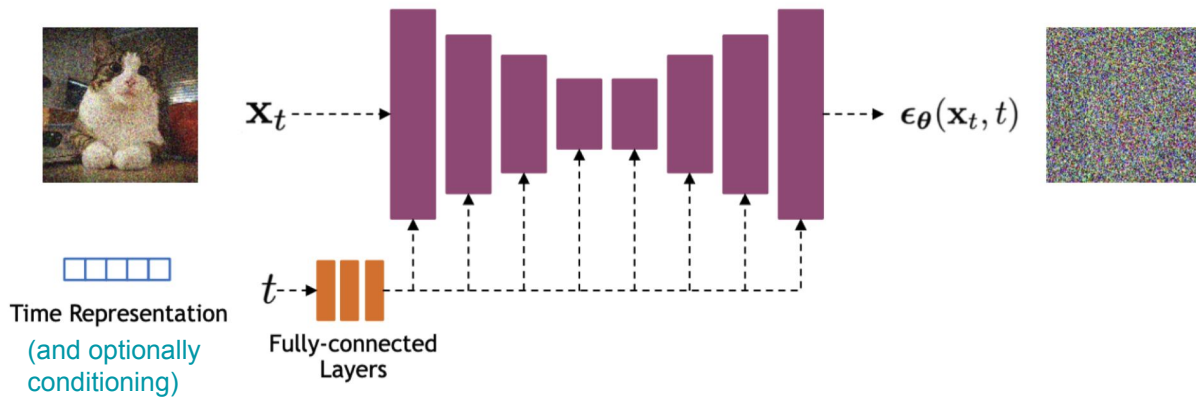
- 1:  $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 2: **for**  $t = T, \dots, 1$  **do**
- 3:  $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$  if  $t > 1$ , else  $\mathbf{z} = \mathbf{0}$
- 4:  $\mathbf{x}_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left( \mathbf{x}_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \boldsymbol{\epsilon}_{\theta}(\mathbf{x}_t, t) \right) + \sigma_t \mathbf{z}$
- 5: **end for**
- 6: **return**  $\mathbf{x}_0$

start with noise

subtract noise

predict noise

# Background: Diffusion

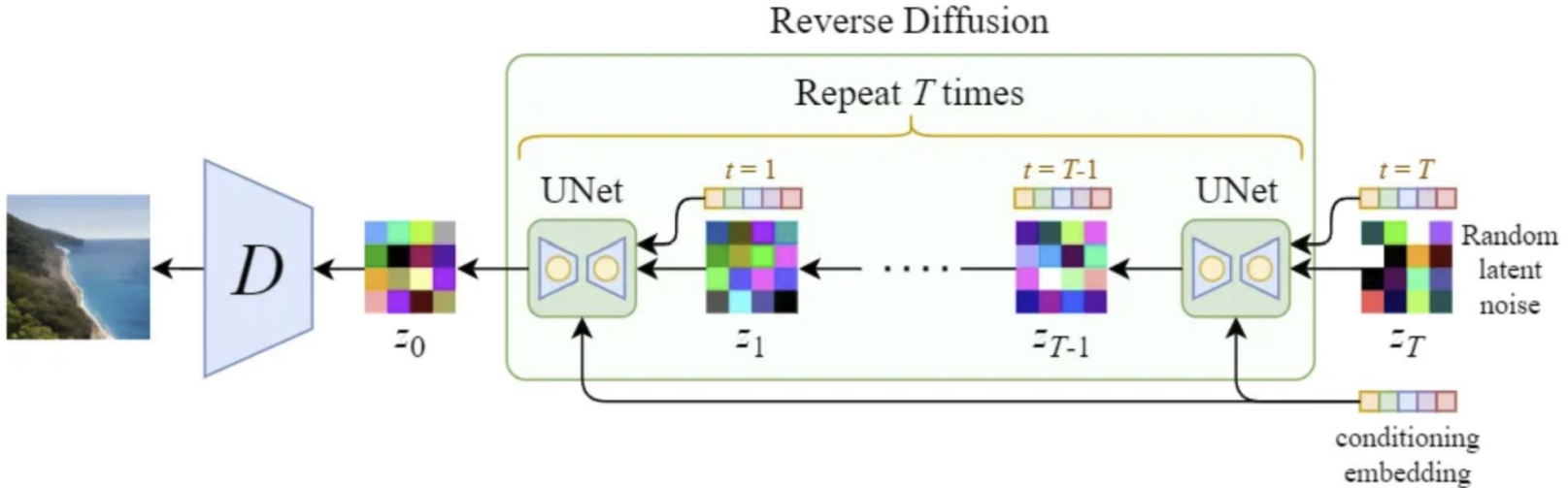
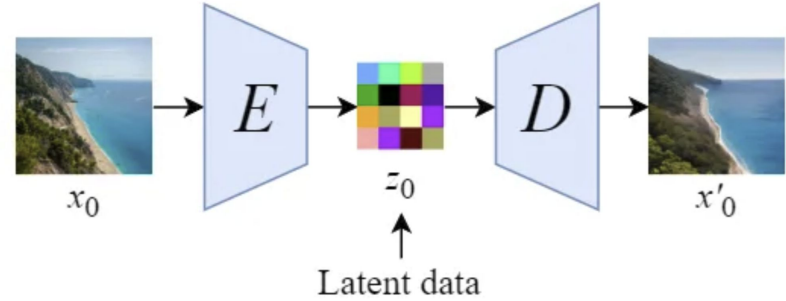


## Hindrances:

- Iterative
- Image-sized inputs + latents

# Background: Latent Diffusion

- Encode image into smaller latent space
- Run diffusion in latent space

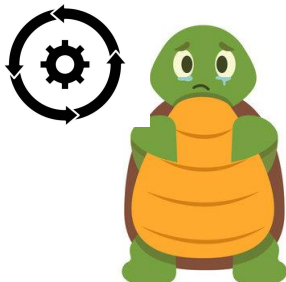




# Motivation

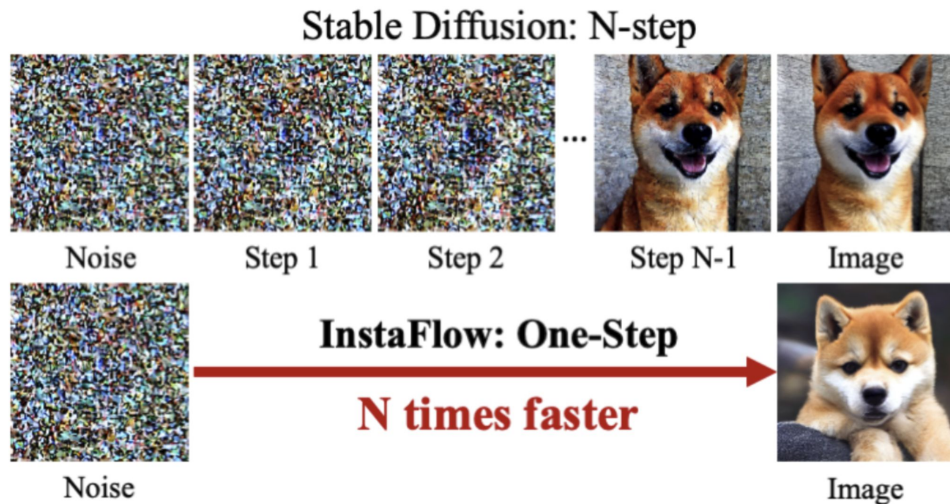
Sampling is still **iterative**

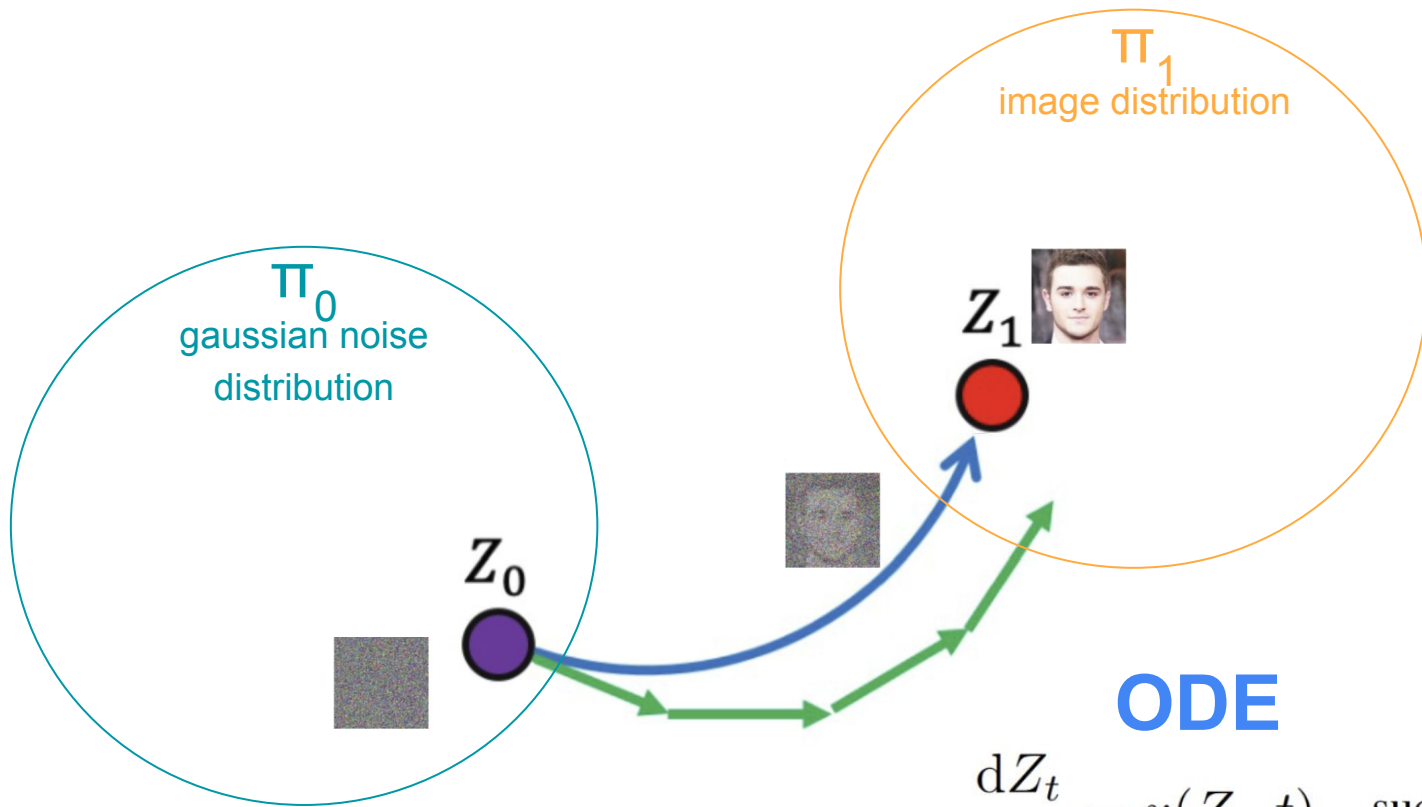
That makes it **slow**



**InstaFlow** claims

- “an ultra-fast one-step model”
- “with SD-level image quality”





$\pi_1$   
image distribution

$\pi_0$   
gaussian noise  
distribution

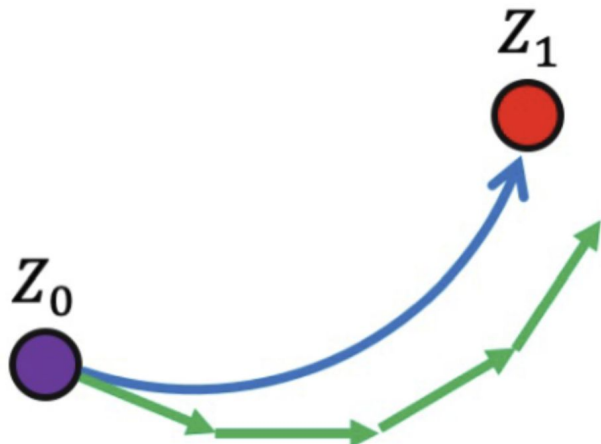
$Z_1$

$Z_0$

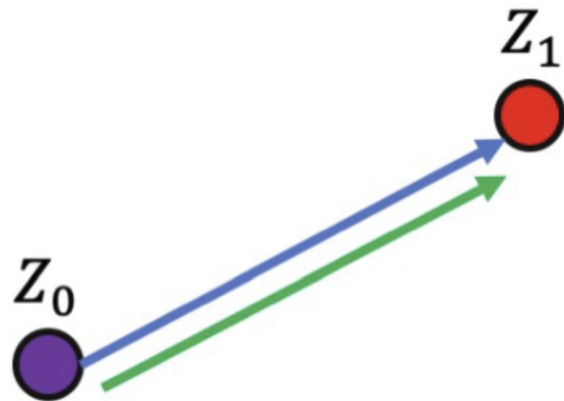
**ODE**

$$\frac{dZ_t}{dt} = v(Z_t, t), \quad \text{such that } Z_1 \sim \pi_1 \\ Z_0 \sim \pi_0$$

**v is the ODE's velocity field**



$$Z_{t+\frac{1}{N}} = Z_t + \frac{1}{N}v(Z_t, t),$$



Perfect Simulation in  
just one step

Rectified Flow: an **ODE** Framework

- Learns to transfer  $\pi_0$  to  $\pi_1$
- approximated by numerical solvers (e.g. Euler Method)
- choice of  $N$  yields a cost-accuracy trade-off
- **want a straight flow**

## Neural Net

$$\min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \left\| \frac{d}{dt} X_t - v(X_t, t) \right\|^2 dt \right]$$

Can choose different **interpolations** for  $X_t$



$$X_t = \phi(X_0, X_1, t)$$

$$X_t = tX_1 + (1 - t)X_0$$

$$\frac{d}{dt} X_t = X_1 - X_0$$

can use any  
time-differentiable  
interpolation

DDIM uses this one

$$= \min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \left\| (X_1 - X_0) - v(X_t, t) \right\|^2 dt \right]$$


# Rectified Flow


Now we have  $v$ ,  
approximated by a neural net

We can now solve this ODE  
starting from  $Z_0 \sim \pi_0$  to transfer  $\pi_0$  to  $\pi_1$

$$\frac{dZ_t}{dt} = v(Z_t, t),$$

$$(Z_0, Z_1) = \text{Rectify}((X_0, X_1))$$

  
the rectified coupling

  
can be an arbitrary coupling of  $\pi_0, \pi_1$

why is this any better?

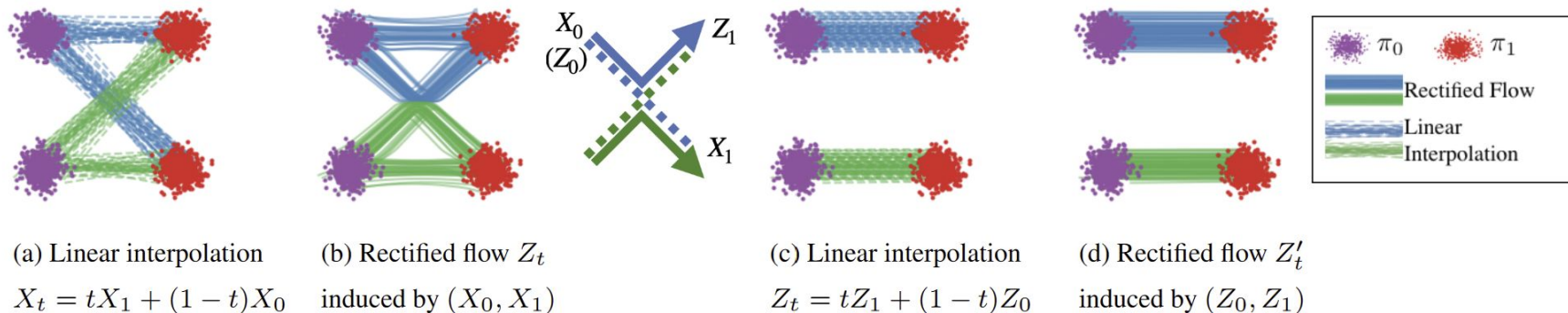


Figure 2: (a) Linear interpolation of data input  $(X_0, X_1) \sim \pi_0 \times \pi_1$ . (b) The rectified flow  $Z_t$  induced by  $(X_0, X_1)$ ; the trajectories are “rewired” at the intersection points to avoid the crossing. (c) The linear interpolation of the end points  $(Z_0, Z_1)$  of flow  $Z_t$ . (d) The rectified flow induced from  $(Z_0, Z_1)$ , which follows straight paths.

- Paths for a well-defined ODE cannot cross each other
  - That would mean the ODE solution is not unique
- **Rectified Flow rewires individual trajectories to avoid crossing**
  - Still keeps the same density map

**Flows avoid crossing** A key to understanding the method is the non-crossing property of flows: the different paths following a well defined ODE  $dZ_t = v(Z_t, t)dt$ , whose solution exists and is unique, *cannot cross each other* at any time  $t \in [0, 1)$ . Specifically, there exists no location  $z \in \mathbb{R}^d$  and time  $t \in [0, 1)$ , such that two paths go across  $z$  at time  $t$  along different directions, because otherwise the solution of the ODE would be non-unique. On the other hand, the paths of the interpolation process  $X_t$  may intersect with each other (Figure 2a), which makes it non-causal. Hence, as shown in Figure 2b, the rectified flow *rewires* the individual trajectories passing through the intersection points to avoid crossing, while tracing out the same density map as the linear interpolation paths due to the optimization of (1). We can view the linear interpolation  $X_t$  as building roads (or tunnels) to connect  $\pi_0$  and  $\pi_1$ , and the rectified flow as traffics of particles passing through the roads in a myopic, memoryless, non-crossing way, which allows them to ignore the global path information of how  $X_0$  and  $X_1$  are paired, and rebuild a more deterministic pairing of  $(Z_0, Z_1)$ .

---

**Algorithm 1** Rectified Flow: Main Algorithm

---

**Procedure:**  $\mathbf{Z} = \text{RectFlow}((X_0, X_1))$ :

*Inputs:* Draws from a coupling  $(X_0, X_1)$  of  $\pi_0$  and  $\pi_1$ ; velocity model  $v_\theta: \mathbb{R}^d \rightarrow \mathbb{R}^d$  with parameter  $\theta$ .

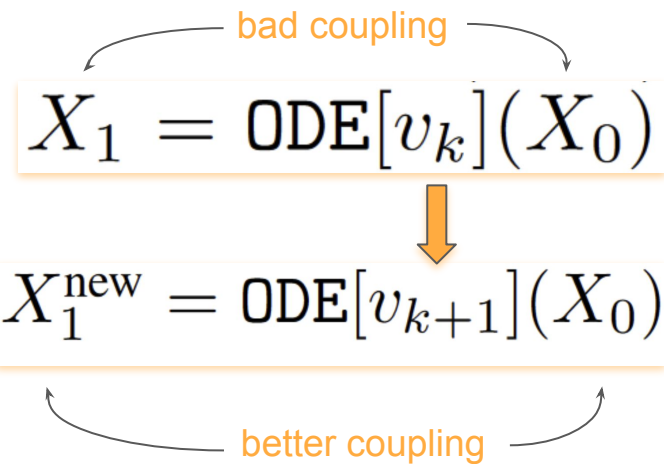
*Training:*  $\hat{\theta} = \arg \min_{\theta} \mathbb{E} \left[ \|X_1 - X_0 - v(tX_1 + (1-t)X_0, t)\|^2 \right]$ , with  $t \sim \text{Uniform}([0, 1])$ .

*Sampling:* Draw  $(Z_0, Z_1)$  following  $dZ_t = v_{\hat{\theta}}(Z_t, t)dt$  starting from  $Z_0 \sim \pi_0$  (or backwardly  $Z_1 \sim \pi_1$ ).

*Return:*  $\mathbf{Z} = \{Z_t: t \in [0, 1]\}$ .

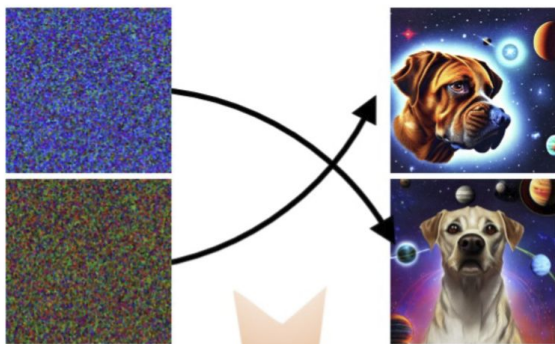
**Reflow** (optional):  $\mathbf{Z}^{k+1} = \text{RectFlow}((Z_0^k, Z_1^k))$ , starting from  $(Z_0^0, Z_1^0) = (X_0, X_1)$ .

- 1) Draw  $(X_0, X_1)$
- 2) Train  $v$
- 3) Follow  $v$  to get new mapping  $(Z_0, Z_1)$
- 4) Do it again.
  - Prev  $(Z_0, Z_1)$  is now the input  $(X_0, X_1)$
  - “k-flow”



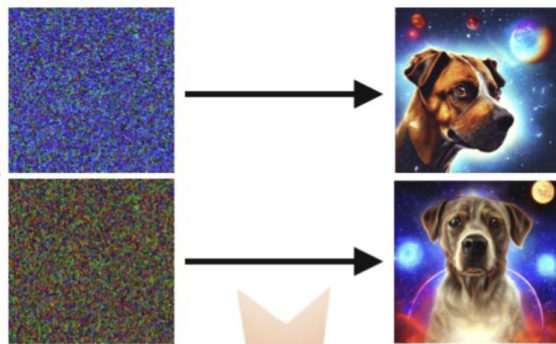


**Original Model**  
Curved Trajectories; **Bad** Coupling

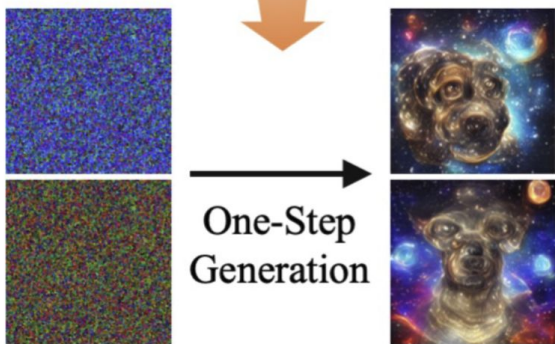


**Reflow**

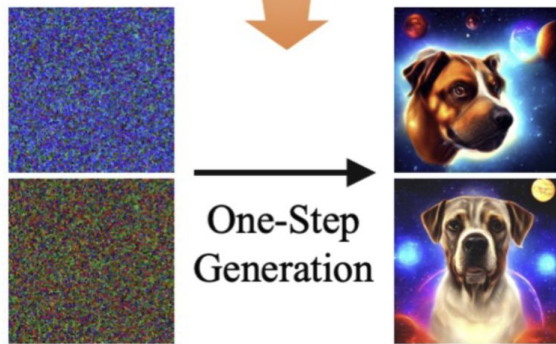
**Reflowed Model**  
Straighter Trajectories; **Good** Coupling



Distill



Distill



 **Failure: blurry, low-quality**

 **Success: clear, high-quality**

## Text-Conditioned

$$\min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \| (X_1 - X_0) - v(X_t, t) \|^2 dt \right]$$

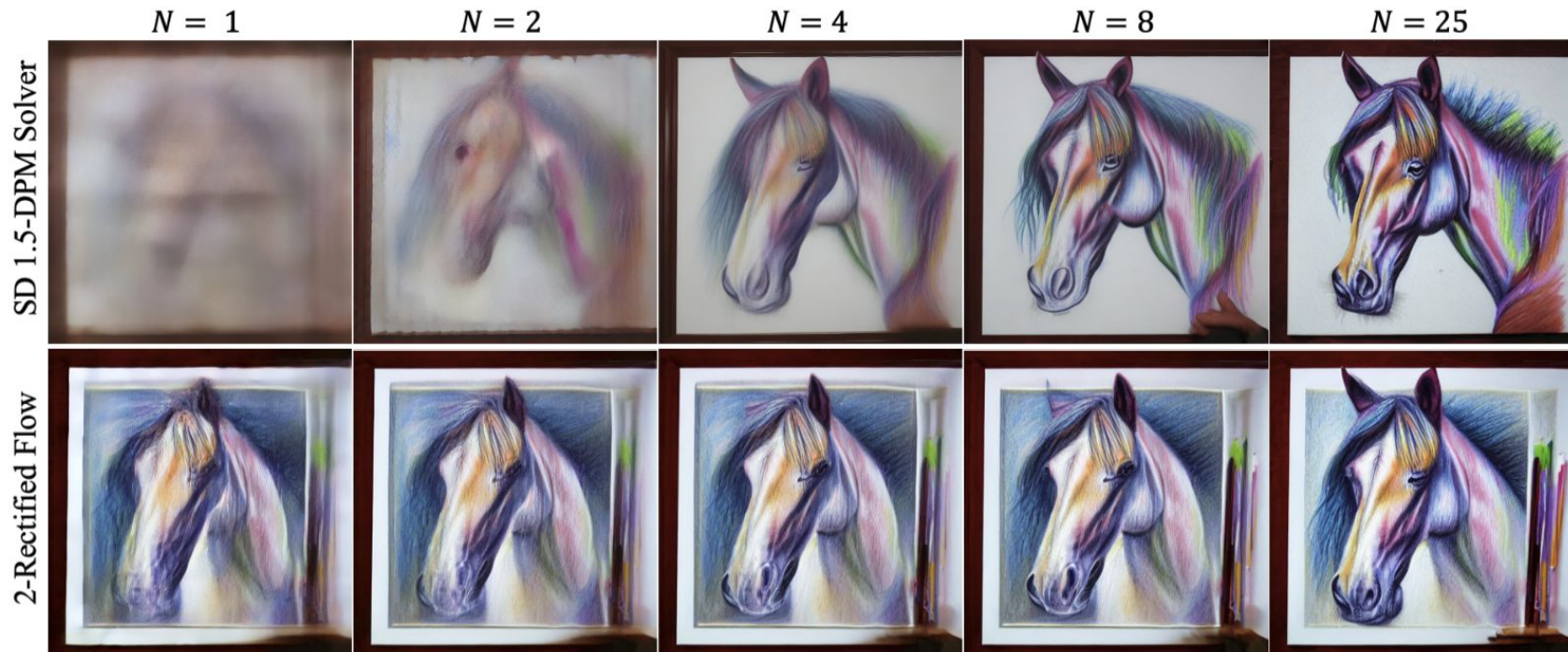
$$v_{k+1} = \arg \min_v \mathbb{E}_{X_0 \sim \pi_0, \mathcal{T} \sim D_{\mathcal{T}}} \left[ \int_0^1 \| (X_1 - X_0) - v(X_t, t | \mathcal{T}) \|^2 dt \right]$$

$D_{\mathcal{T}}$  is a dataset of text prompts

$$X_1 = \text{ODE}[v_k](X_0 | \mathcal{T}) = X_0 + \int_0^1 v_k(X_t, t | \mathcal{T}) dt$$

# Evaluation





*'Masterpiece color pencil drawing of a horse; bright vivid color'*

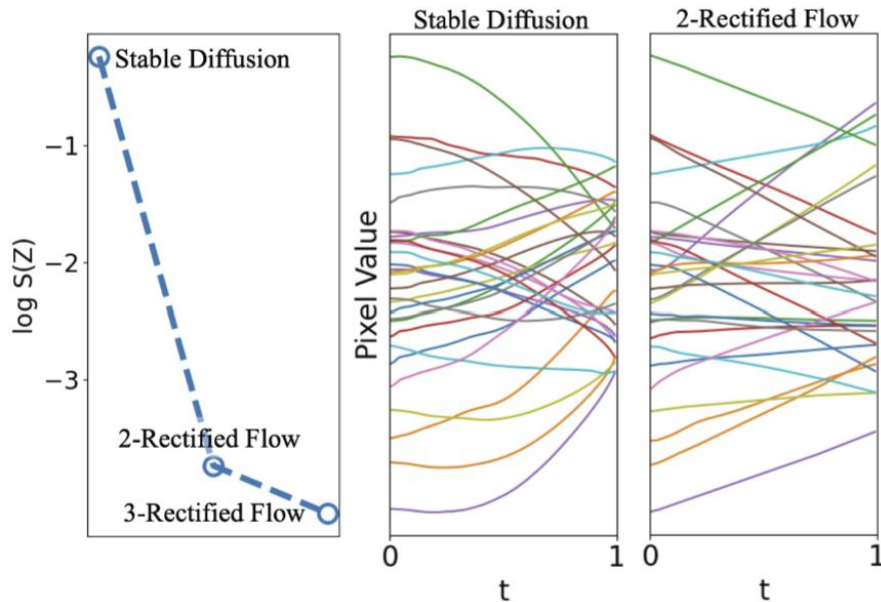
# Evaluation: Flow

Flow 'straightness'

$$S(Z) = \int_{t=0}^1 \mathbb{E} [\| (Z_1 - Z_0) - v(Z_t, t) \|^2] dt$$

How pixel values change over time

Straighter trajectories

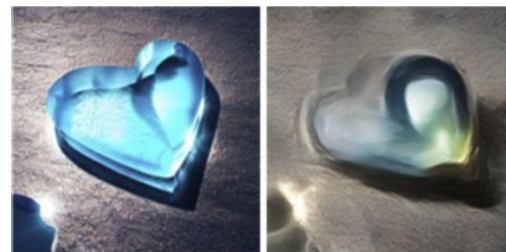
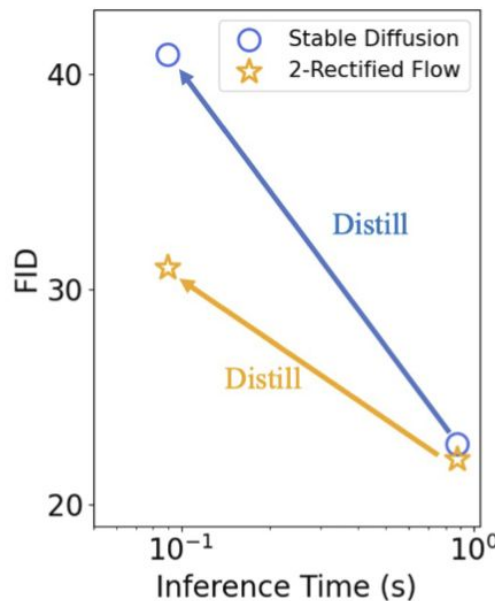


# Evaluation: Inference Time

SD cannot effectively distill

k-Rectified can

Distilled version of k-Flow has a smaller gap with teacher



Stable Diffusion

+Distill



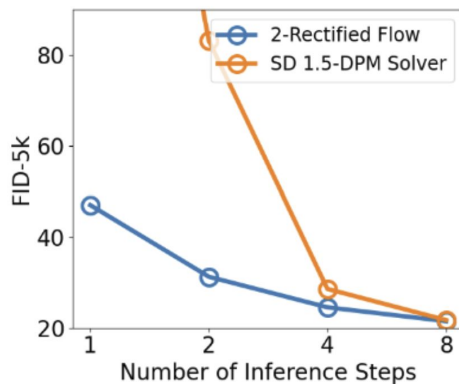
2-Rectified Flow

+Distill

# Evaluation: FID + CLIP

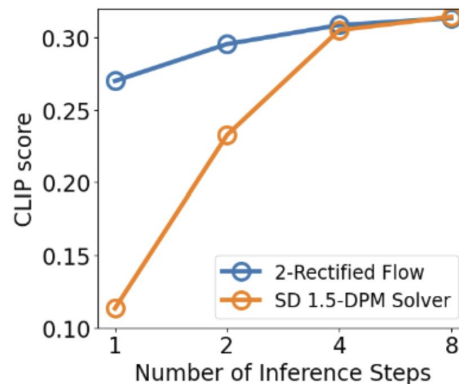
Method	Inf. Time	FID-5k	CLIP
SD 1.4 (25 step)[70]	0.88s	22.8	<b>0.315</b>
<b>(Pre) 2-RF (25 step)</b>	0.88s	<b>22.1</b>	0.313
PD (1 step)[58]	0.09s	37.2	0.275
SD 1.4+Distill	0.09s	40.9	0.255
<b>(Pre) 2-RF (1 step)</b>	0.09s	68.3	0.252
<b>(Pre) 2-RF+Distill</b>	0.09s	<b>31.0</b>	<b>0.285</b>

(a) MS COCO 2017



Method	Inf. Time	FID-30k
SD* [70]	2.9s	<b>9.62</b>
<b>(Pre) 2-RF (25 step)</b>	0.88s	13.4
SD 1.4+Distill	0.09s	34.6
<b>(Pre) 2-RF+Distill</b>	0.09s	<b>20.0</b>

(b) MS COCO 2014



# Thoughts

- It's all based on Rectified Flow
- New text-conditioning
- boasts strong improvement (FID + Results)
- failures on complicated prompts
- usefulness of FID ?
- better reflow?
- 1-step ?

InstaFlow-0.9B

SD 1.5-DPM Solver



*A dog and a cat shake hands with each other*

DPM Solver

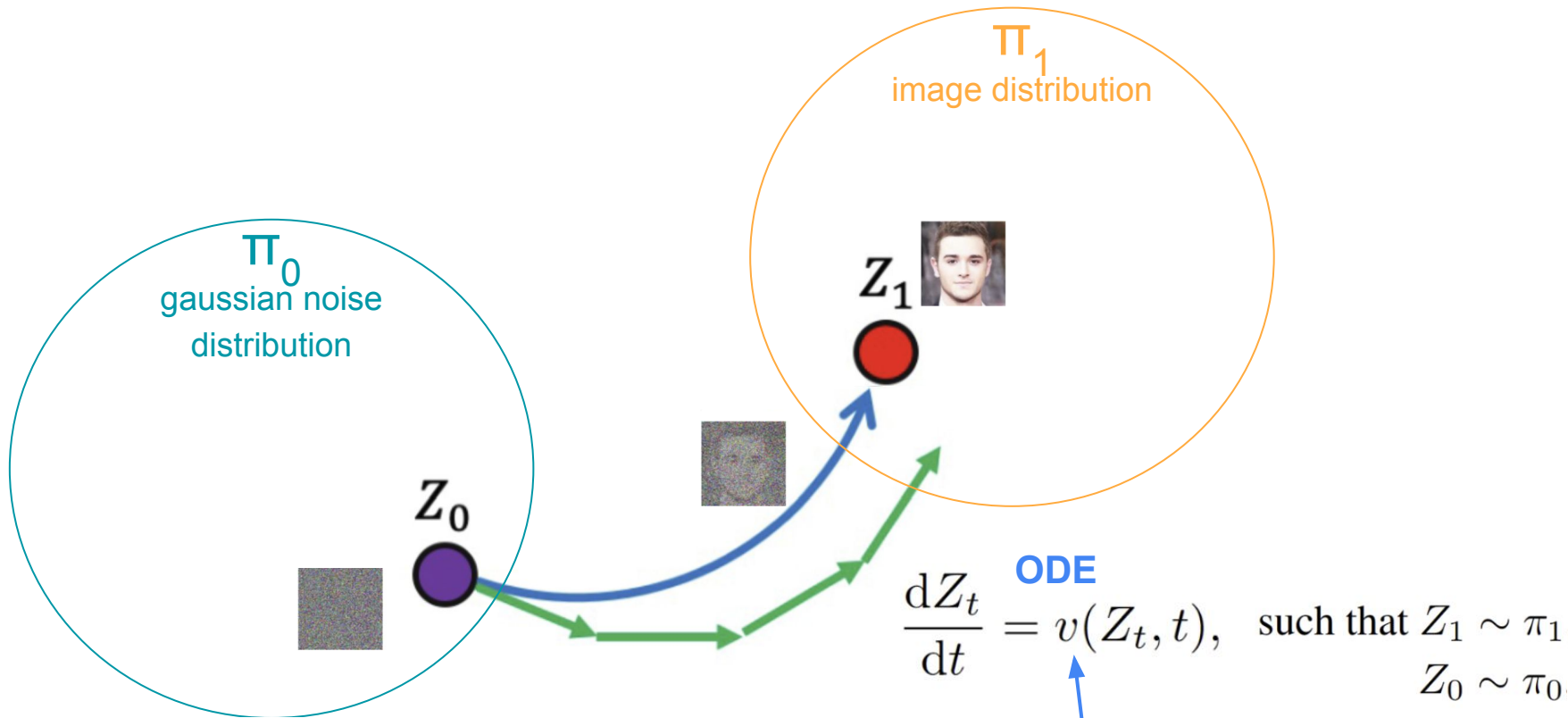


2-Rectified Flow







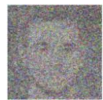


$$\min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \left\| \frac{d}{dt} X_t - v(X_t, t) \right\|^2 dt \right]$$

$$\uparrow$$

$$\phi(X_0, X_1, t),$$

Can choose different **interpolations** for  $X_t$



$$X_t = \phi(X_0, X_1, t)$$

can use any  
time-differentiable  
interpolation

$$X_t = tX_1 + (1 - t)X_0$$

$$\frac{d}{dt}X_t = X_1 - X_0$$

DDIM uses this

**Neural Net**

$$\begin{aligned} & \min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \left\| \frac{d}{dt}X_t - v(X_t, t) \right\|^2 dt \right] \\ &= \min_v \mathbb{E}_{(X_0, X_1) \sim \gamma} \left[ \int_0^1 \left\| (X_1 - X_0) - v(X_t, t) \right\|^2 dt \right] \end{aligned}$$