



On the design space of one step diffusion via shortcutting flow paths

Reporter: Haitao Lin

Supervisor: Stan Z. Li & Tailin Wu

School of Engineering, Westlake University
Computer Science and Technology, Zhejiang University



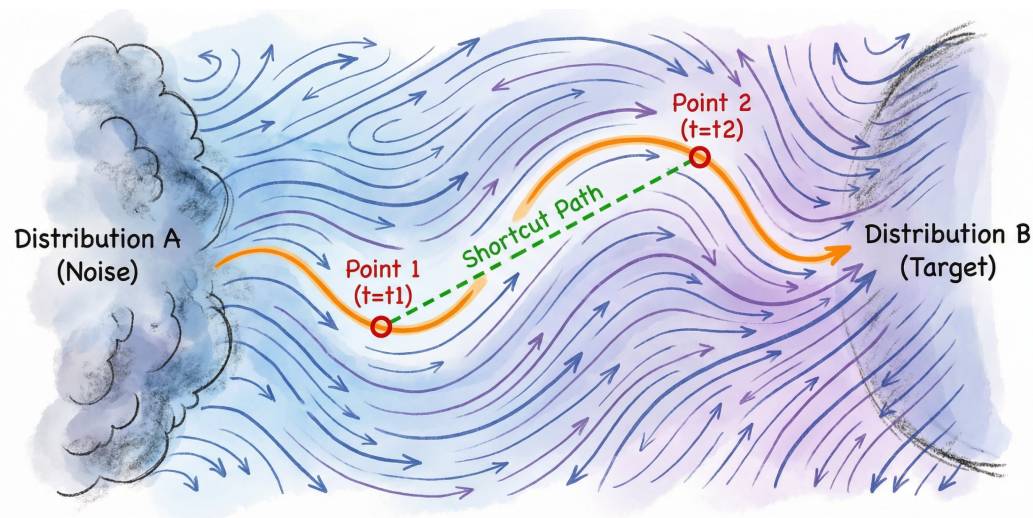
Denosed Diffusion

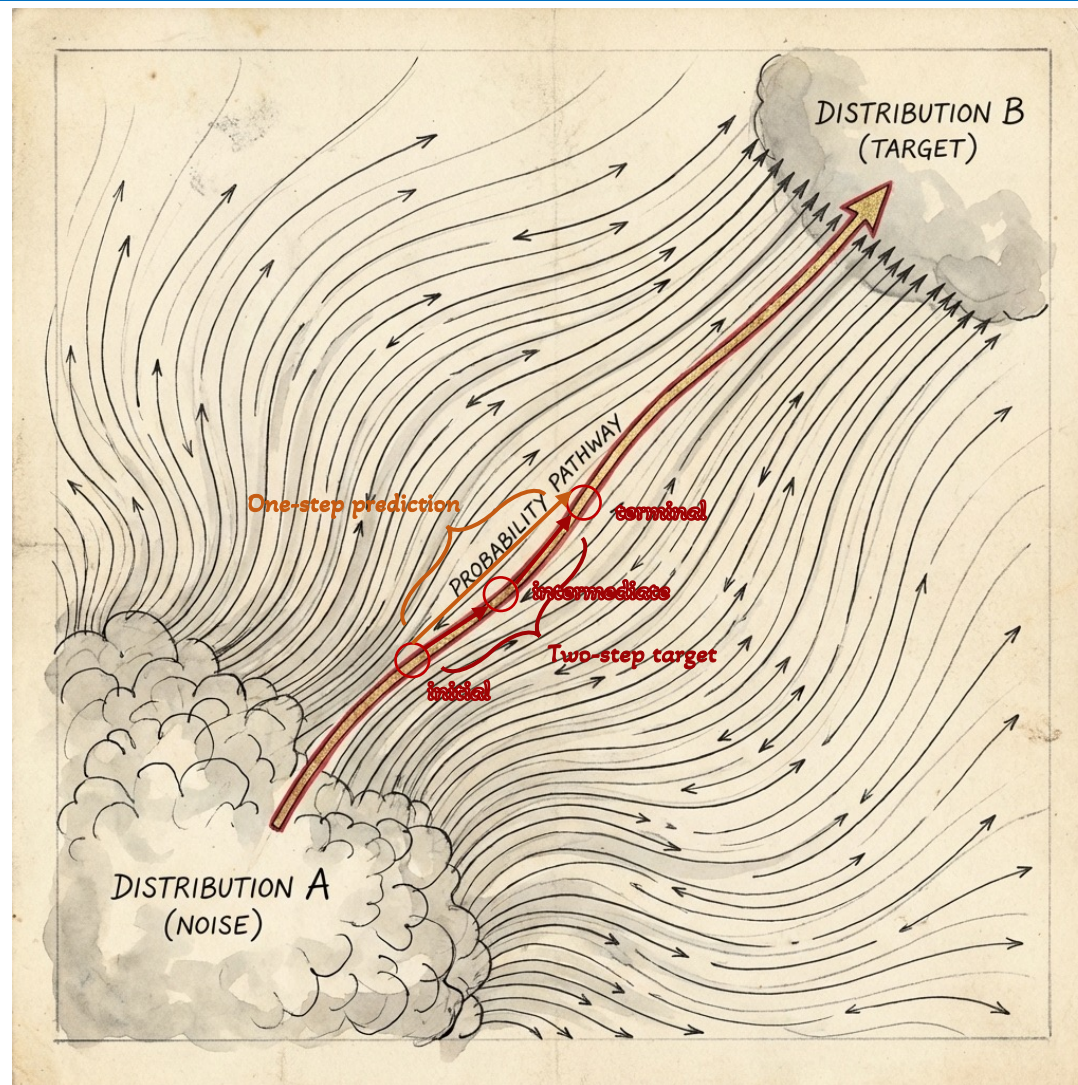


Few-step Diffusion



One-step Diffusion





Preliminary

1. Diffusion Models

$$\mathbf{x}_t = \alpha_t \mathbf{x}_0 + \sigma_t \boldsymbol{\varepsilon} \quad \mathbf{x}_0 \sim p_0, \boldsymbol{\varepsilon} \sim p_1 \quad \text{(Diffusion process)}$$

$$\mathbf{v}_t(\mathbf{x}) = \dot{\alpha}_t \mathbb{E}(\mathbf{x}_0 | \mathbf{x}_t = \mathbf{x}) + \dot{\sigma}_t \mathbb{E}(\boldsymbol{\varepsilon} | \mathbf{x}_t = \mathbf{x}). \quad \text{(Marginal Velocity)}$$

$$\dot{\mathbf{x}}_t = \mathbf{v}_t(\mathbf{x}_t). \quad \text{(PF-ODE)}$$

$$\mathbf{v}_{t|0} = \mathbf{v}_t(\mathbf{x}_t | \mathbf{x}_0) = \dot{\alpha}_t \mathbf{x}_0 + \dot{\sigma}_t \boldsymbol{\varepsilon} \quad \text{(Conditional Velocity)}$$

$\sigma = \sigma_{\text{data}}, \alpha_t = \cos(\frac{\pi}{2}t), \text{ and } \sigma_t = \sin(\frac{\pi}{2}t)$ (Coefficient in cosine paths)

$\sigma = 1, \alpha_t = 1 - t \text{ and } \sigma_t = t$ (Coefficient in linear paths)

Preliminary

2. Flow Map Solvers

$$\mathbf{x}_r = X_{t,r}(\mathbf{x}_t) = \mathbf{x}_t + \int_t^r \mathbf{v}_\tau(\mathbf{x}_\tau) d\tau,$$

(Flow Map from t to r)

$$X_{t,r}(\mathbf{x}_t) = \mathbf{x}_t + (r - t) \cdot \mathbf{u}_{t,r}(\mathbf{x}_t)$$

$$\text{where } \mathbf{u}_{t,r}(\mathbf{x}_t) = \frac{1}{r - t} \int_t^r \mathbf{v}_\tau(\mathbf{x}_\tau) d\tau,$$

(Average Velocity Solver)

$$X_{t,r}(\mathbf{x}_t) \approx \text{DDIM}(\mathbf{x}_t, \mathbf{v}_t, t, r) = \bar{\alpha}_{t,r} \mathbf{x}_t + \bar{\beta}_{t,r} \mathbf{v}_t,$$

$$\bar{\alpha}(t, r) = \cos\left(\frac{\pi}{2}(r - t)\right), \quad \bar{\beta}(t, r) = \frac{2}{\pi} \sin\left(\frac{\pi}{2}(r - t)\right)$$

$$\bar{\alpha}(t, r) = 1, \quad \bar{\beta}(t, r) = r - t$$

(One-step DDIM Solver)

(Coefficient in cosine paths)

(Coefficient in linear paths)

General Designing Framework

1. Consistency Property & Learning Objective & Time Sampler

$$X_{s,r}(X_{t,s}(\mathbf{x}_t)) = X_{t,r}(\mathbf{x}_t).$$

(Consistency Property)

$$\arg \min_{\theta} \mathbb{E}_{r,s,t \sim p(\tau), \mathbf{x}_t \sim p_t} \left[\underbrace{w(r, s, t) \cdot d \left(\overbrace{X_{t,r}^{\theta}(\mathbf{x}_t)}^{\text{one-step prediction}}, \overbrace{\text{sg}(\hat{X}_{s,r} \circ \hat{X}_{t,s}(\mathbf{x}_t))}^{\text{two-step target}} \right)}_{l(\mathbf{x}_t, r, s, t; \theta)} \right] \quad (\text{Learning Objective})$$

where w is the weight term, \hat{X} and X^{θ} are flow maps obtained with the conditional velocity or the neural network F^{θ} , $d(\cdot, \cdot)$ is a loss metric function, such as the squared l_2 -distance, and $\text{sg}(\cdot)$ is the stop gradient operator in backpropagation.

We call $\hat{X}_{s,r} \circ \hat{X}_{t,s}(\mathbf{x}_t)$ two-step flow map targets and $X_{t,r}^{\theta}(\mathbf{x}_t)$ one-step flow map predictions, and write the inner loss term of expectation as $l(\mathbf{x}_t, r, s, t; \theta)$.

$\{r, s, t\}$ discrete time points \Rightarrow DTSC

$s \rightarrow t \Rightarrow$ CTSC

(Time Sampler)

Examples: DTSC & CTSC

1. Consistency Training Model (DTSC)

$$\begin{aligned}
 \mathbf{x}_0 &\sim p_0, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, 1) \\
 \mathbf{x}_t &= \cos\left(\frac{\pi}{2}t\right)\mathbf{x}_0 + \sin\left(\frac{\pi}{2}t\right)\boldsymbol{\varepsilon} \\
 \mathbf{v}_{t|0} &= -\frac{\pi}{2}\sin\left(\frac{\pi}{2}t\right)\mathbf{x}_0 + \frac{\pi}{2}\cos\left(\frac{\pi}{2}t\right)\boldsymbol{\varepsilon} \\
 \hat{\mathbf{x}}_s &= \cos\left(\frac{\pi}{2}s\right)\mathbf{x}_0 + \sin\left(\frac{\pi}{2}s\right)\boldsymbol{\varepsilon},
 \end{aligned}
 \tag{Conditional Flow Path}$$

$$\begin{aligned}
 \hat{X}_{t,s}(\mathbf{x}_t) &= \hat{\mathbf{x}}_s = \text{DDIM}(\mathbf{x}_t, \mathbf{v}_{t|0}, t, s) \\
 \hat{X}_{s,r}(\hat{\mathbf{x}}_s) &= \hat{\mathbf{x}}_r = \text{DDIM}(\hat{\mathbf{x}}_s, F^\theta(\hat{\mathbf{x}}_s), s, r) \\
 X_{t,r}^\theta(\mathbf{x}_t) &= \mathbf{x}_r^\theta = \text{DDIM}(\mathbf{x}_t, F^\theta(\mathbf{x}_t), t, r)
 \end{aligned}
 \tag{Flow Map Construction}$$

$$\begin{aligned}
 l_{\text{ct}}(\mathbf{x}_t, r, s, t; \theta) &= \text{LPIPS}(f^\theta(\mathbf{x}_t), \text{sg}(f^\theta(\hat{\mathbf{x}}_s))), \\
 &= \text{LPIPS}(X_{t,r}^\theta(\mathbf{x}_t), \text{sg}(\hat{X}_{s,r}(\hat{\mathbf{x}}_s))) \\
 &= \text{LPIPS}\left(\text{DDIM}(\mathbf{x}_t, \mathbf{v}_t^\theta(\mathbf{x}_t), t, r), \text{sg}(\text{DDIM}(\hat{\mathbf{x}}_s, \mathbf{v}_s^\theta(\hat{\mathbf{x}}_s), s, r))\right)
 \end{aligned}
 \tag{Learning Objective}$$

		CT
Diffusion basis[†]		
Flow path		Cosine
Network	Architecture	U-Net
F^θ	Output	\mathbf{v}^θ
Flow map construction		
		(note: *)
		$t = \frac{\pi}{2} \arctan([\sigma_{\max}^{1/\rho} + \frac{\tau}{K}(\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho})]^\rho)$
		$s = \frac{\pi}{2} \arctan([\sigma_{\max}^{1/\rho} + \frac{\tau+1}{K}(\sigma_{\min}^{1/\rho} - \sigma_{\max}^{1/\rho})]^\rho)$
Time sampler		$r = 0$, where $\tau \sim \mathcal{U}\{0, \dots, K-1\}$
Two-step target	1st-step ($\hat{\mathbf{x}}_s$)	$\text{DDIM}(\mathbf{x}_t, \mathbf{v}_{t 0}, t, s)$
	2nd-step ($\hat{\mathbf{x}}_r$)	$\text{DDIM}(\hat{\mathbf{x}}_s, \mathbf{v}_s^\theta, s, r)$
One-step prediction (\mathbf{x}_r^θ)		$\text{DDIM}(\mathbf{x}_t, \mathbf{v}_t^\theta, t, r)$
Training		
Loss metric d		LPIPS
sg EMA decay		✓

Examples: DTSC & CTSC

2. Shortcut Diffusion (DTSC)

$$\begin{aligned} \mathbf{x}_0 &\sim p_0, \quad \boldsymbol{\varepsilon} \sim \mathcal{N}(0, 1) \\ \mathbf{x}_t &= (1-t)\mathbf{x}_0 + t\boldsymbol{\varepsilon} && \text{(Conditional Flow Path)} \\ \mathbf{v}_{t|0} &= -\mathbf{x}_0 + \boldsymbol{\varepsilon} \end{aligned}$$

$$\begin{aligned} \hat{X}_{t,s}(\mathbf{x}_t) &= \hat{\mathbf{x}}_s = \mathbf{x}_t - h\mathbf{u}_{t,s}^\theta(\mathbf{x}_t) \\ \hat{X}_{s,r}(\mathbf{x}_s) &= \hat{\mathbf{x}}_r = \mathbf{x}_s - h\mathbf{u}_{s,r}^\theta(\hat{\mathbf{x}}_s) && \text{(Flow Map Construction)} \\ X_{t,r}^\theta(\mathbf{x}_t) &= \hat{\mathbf{x}}_t = \mathbf{x}_t - 2h\mathbf{u}_{t,r}^\theta(\mathbf{x}_t) \end{aligned}$$

$$\begin{aligned} l_{\text{scd}}(\mathbf{x}_t, r, s, t; \theta) &= \frac{1}{4h^2} \cdot \left\| \mathbf{x}_t - 2h\mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - \text{sg}(\mathbf{x}_t - h\mathbf{u}_{t,s}^\theta(\mathbf{x}_t) - h\mathbf{u}_{s,r}^\theta(\mathbf{x}_t - h\mathbf{u}_{t,s}^\theta(\mathbf{x}_t))) \right\|_2^2 \\ &= \frac{1}{4h^2} \cdot \left\| \mathbf{x}_t - 2h\mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - \text{sg}(\mathbf{x}_t - h\mathbf{u}_{t,s}^\theta(\mathbf{x}_t) - h\mathbf{u}_{s,r}^\theta(\hat{\mathbf{x}}_s)) \right\|_2^2 \\ &= \left\| \mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - \frac{1}{2} \text{sg}(\mathbf{u}_{t,s}^\theta(\mathbf{x}_t) + \mathbf{u}_{s,r}^\theta(\hat{\mathbf{x}}_s)) \right\|_2^2, && \text{(Learning Objective)} \end{aligned}$$

			SCD
Diffusion basis[†]			
Flow path			Linear
Network	Architecture		DiT
F^θ	Output		\mathbf{u}^θ
Flow map construction			
			$t = \tau$
			$s = \tau - h$
			$r = \tau - 2h$
			and with p_{teq} ,
			$r = s = t$,
			where $\tau, h \sim$
			Uniform $\log_2(\tau, h)$
Time sampler			
Two-step target	1st-step ($\hat{\mathbf{x}}_s$)		$\mathbf{x}_t - h\mathbf{u}_{t,s}^\theta(\mathbf{x}_t)$
	2nd-step ($\hat{\mathbf{x}}_r$)		$\hat{\mathbf{x}}_s - h\mathbf{u}_{s,r}^\theta(\hat{\mathbf{x}}_s)$
One-step prediction (\mathbf{x}_r^θ)			$\mathbf{x}_t - 2h\mathbf{u}_{t,r}^\theta(\mathbf{x}_t)$
Training			
Loss metric d			Squared l_2 -distance
sg EMA decay			\times

Examples: DTSC & CTSC

3. MeanFlow (CTSC)

Consistency property

$$\begin{aligned} & (\mathbf{x}_t + (r-t)\mathbf{u}_{t,r}^\theta(\mathbf{x}_t)) - (\mathbf{x}_t + (s-t)\mathbf{v}_t + (r-s)\mathbf{u}_{s,r}^\theta(\mathbf{x}_s)) \\ & = (r-t)\mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - (s-t)\mathbf{v}_t - (r-s)\mathbf{u}_{s,r}^\theta(\mathbf{x}_s). \end{aligned}$$

$s \rightarrow t$ and normalized by dt

$$\begin{aligned} & ((r-t)\mathbf{u}_{t,r}^\theta(\mathbf{x}_t) + dt \cdot \mathbf{v}_t - (r-t+dt)\mathbf{u}_{t-dt,r}^\theta(\mathbf{x}_{t-dt}))/dt \\ & = dt \cdot \left(\mathbf{v}_t + \frac{d[(r-t)\mathbf{u}_{t,r}^\theta(\mathbf{x}_t)]}{dt} \right) / dt \\ & = \mathbf{v}_t - \mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - (t-r)\frac{d}{dt}\mathbf{u}_{t,r}^\theta(\mathbf{x}_t). \end{aligned}$$

Replace \mathbf{v}_t with $\mathbf{v}_{t|0}$

$$l(\mathbf{x}_t, r, t-dt, t; \theta) = w \cdot \left\| \mathbf{u}_{t,r}^\theta(\mathbf{x}_t) - \text{sg} \left(\mathbf{v}_{t|0} + (r-t)\frac{d\mathbf{u}_{t,r}^\theta(\mathbf{x}_t)}{dt} \right) \right\|^2$$

(Learning Objective)

		MeanFlow
Diffusion basis[†]		
Flow path		Linear
Network	Architecture	DiT
F^θ	Output	\mathbf{u}^θ
Flow map construction		
		$r, t = \{\text{sigmoid}(\tau_1), \text{sigmoid}(\tau_2)\}$
		$s.t. r \leq t,$
		$s = t - dt,$ and with
		$p_{\text{req}}, r = s = t,$ where
		$\tau_1, \tau_2 \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$
Time sampler		
		(note: ‡)
Two-step target	1st-step ($\hat{\mathbf{x}}_s$)	DDIM($\mathbf{x}_t, \mathbf{v}_{t 0}, t, s$)
	2nd-step ($\hat{\mathbf{x}}_r$)	$\hat{\mathbf{x}}_s + (r-s)\mathbf{u}_{s,r}^\theta(\mathbf{x}_s)$
One-step prediction (\mathbf{x}_r^θ)		$\mathbf{x}_t + (r-t)\mathbf{u}_{t,r}^\theta(\mathbf{x}_t)$
Training		
Loss metric d		Squared l_2 -distance
sg EMA decay		\times

Examples: DTSC & CTSC

4. sConsistency Training Model (CTSC)

Gradient approximation

$$\nabla_{\theta} l(\mathbf{x}_t, r, s, t; \theta) \approx \nabla_{\theta} \|\mathbf{v}_t^{\theta}(\mathbf{x}_t) - \text{sg}(\mathbf{v}_t^{\theta}(\mathbf{x}_t) + w(t) \frac{d}{dt} X_{t,r}^{\theta}(\mathbf{x}_t))\|_2^2.$$

Replace with DDIM solver

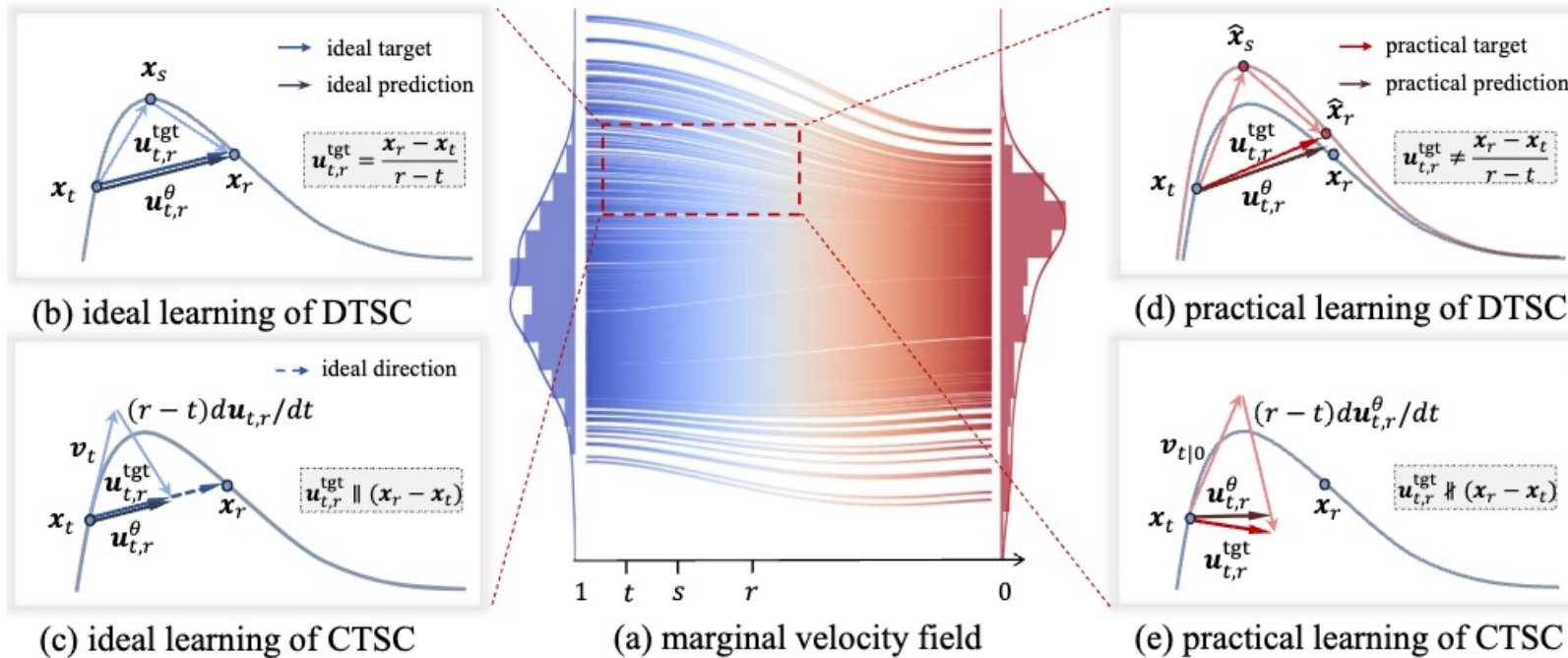
$$l_{\text{sct}}(\mathbf{x}_t, r, t - dt, t; \theta) = \left\| \mathbf{v}_t^{\theta}(\mathbf{x}_t) - \text{sg} \left(\mathbf{v}_t^{\theta}(\mathbf{x}_t) + w(t) \frac{d\text{DDIM}(\mathbf{x}_t, \mathbf{v}_t^{\theta}(\mathbf{x}_t), t, r)}{dt} \right) \right\|_2^2,$$

(Learning Objective)

sCT with linear paths is of the same form as MeanFlow

		sCT
Diffusion basis[†]		
Flow path		Cosine
Network	Architecture	U-Net
F^{θ}	Output	\mathbf{v}^{θ}
Flow map construction		
		$t = \frac{2}{\pi} \arctan(\exp(\tau))$
		$s = t - dt$
		$r = 0$, where
		$\tau \sim \mathcal{N}(P_{\text{mean}}, P_{\text{std}}^2)$
Time sampler		
		(note: ‡)
Two-step target	1st-step ($\hat{\mathbf{x}}_s$)	DDIM($\mathbf{x}_t, \mathbf{v}_{t 0}, t, s$)
	2nd-step ($\hat{\mathbf{x}}_r$)	DDIM($\hat{\mathbf{x}}_s, \mathbf{v}_s^{\theta}, s, r$)
One-step prediction (\mathbf{x}_r^{θ})		DDIM($\mathbf{x}_t, \mathbf{v}_t^{\theta}, t, r$)
Training		
Loss metric d		Squared l_2 -distance
sg EMA decay		X

Discussion



Error bound of DTSC&CTSC

Under the assumption:

- (i) one-sided Lipschitz condition of marginal velocity
- (ii) twice continuous differentiability with bounded second derivatives of $X_{\tau_1, \tau_2}^{\theta}$

$$W_2^2(p_0, p_0^{\theta}) \leq C_1 \mathcal{L}_{dtsc}(\theta) + C_2(t - s);$$

$$W_2^2(p_0, p_0^{\theta}) \leq C_3 \mathcal{L}_{ctsc}(\theta),$$

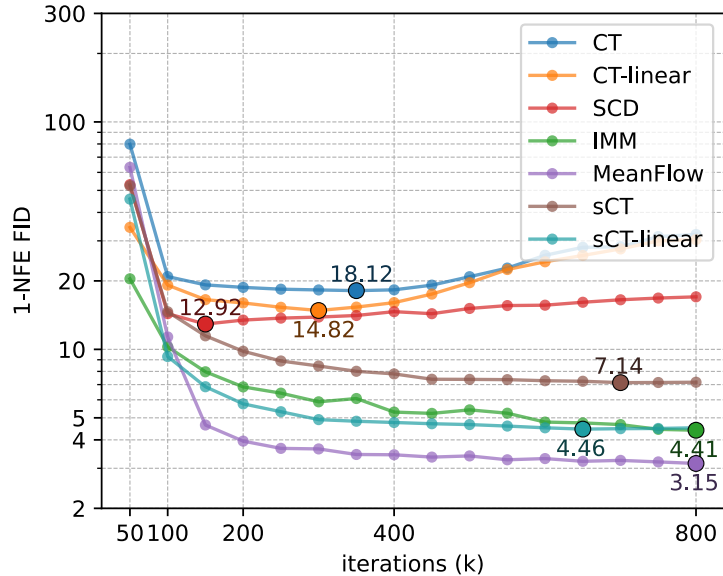
$W_2(\cdot, \cdot)$ is the Wasserstein-2 distance.

Why sharing a common design frame?

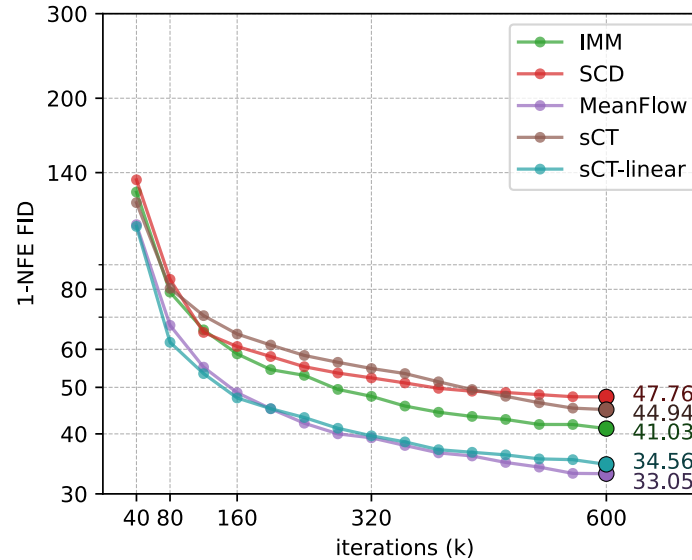
What challenges in constructing flow map targets?

Why distillation from pretrained velocity fields performs better?

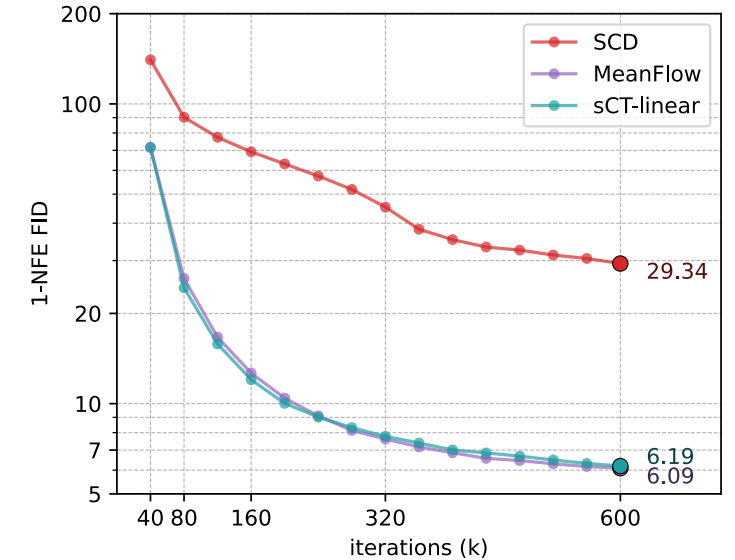
Elucidating the Design Space



Uncond. CIFAR



Cond. ImageNet



CFG. ImageNet

Following a linear or cosine path?

Shortcutting flow paths discretely or continuously?

Fixing the terminal time or not?

1. Linear paths in the setting of shortcut models are optimal under Fisher information metrics in the setting of shortcut models
2. The variance of conditional velocity $\sigma_{\mathbf{v}_t|0}^2 := \text{Var}(\mathbf{v}_t(\mathbf{x}_t|\mathbf{x}_0))$ is the key to training stability and inference fidelity

Improvements to Training

Table 2: Evaluation of training improvements under one-step generation with SiT-B/2 as F^θ .

Training configuration	FID50k
MeanFlow under CFG. (Baseline)	6.09
+A1 Plug-in velocity ($p_{\text{plugin}} = 1.0$)	6.01
+A2 Plug-in velocity ($p_{\text{plugin}} = 0.5$)	5.98
+B1 Plug-in velocity ($p_{\text{plugin}} = 1.0$) & class-consistent batching	6.08
+B2 Plug-in velocity ($p_{\text{plugin}} = 0.5$) & class-consistent batching	5.96
+C Gradual time sampler	5.99
+D sCM training techniques	5.95
ESC (Baseline + B2 + C + D)	5.77

Algorithm 1 Calculation of Plug-in Velocity

```
# x: training batch (B,D)
# t: sampled time
e = randn_like(x)
xt = (1- t) * x + t * e
x_ex, xt_ex = x[:,None,:], xt[None,:,:]
eps = (xt_ex - (1- t) * x_ex) / t

logp_fn = Normal(0, 1).log_prob
logp = sum(logp_fn(eps), dim=2)
weight = softmax(logp, dim=0)

v_cnd = eps - x_ex
v_plugin = matmul(weight.T, v_cnd)
```

Use plug-in velocity instead

Ideal velocity

Assume the data distribution is the empirical distribution, as

$$p_0(\mathbf{y}) = \frac{1}{N} \sum_{i=1}^N \mathbb{1}_{\mathbf{y}_i}(\mathbf{y})$$

The marginal velocity reads

$$\mathbf{v}_t^*(\mathbf{x}_t | \{\mathbf{y}^{(i)}\}_{i=1}^N) = \sum_i^N \frac{\mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{y}^{(i)}, \sigma_t^2 \mathbf{I})}{\sum_j^N \mathcal{N}(\mathbf{x}_t; \alpha_t \mathbf{y}^{(j)}, \sigma_t^2 \mathbf{I})} (\dot{\alpha}_t \mathbf{y}^{(i)} + \frac{\dot{\sigma}_t}{\sigma_t} (\mathbf{x}_t - \alpha_t \mathbf{y}^{(i)})),$$

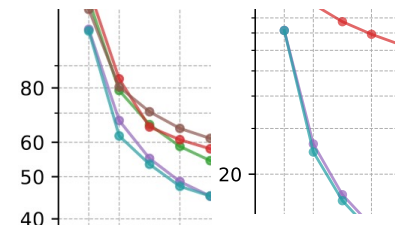
Plug-in velocity instead of conditional one:

- Reducing the level of variance $\sigma_{v_t|0}$ to $O(1/B)$, at the minor cost of increased bias
- Employing class-consistent mini-batching in guidance training to avoiding dilution the class specific signal

Gradual time sampler from sCT to MeanFlow:

- Fast convergence to local minima around the global one

Adoption of training techniques from sCT



Scaling up

Table 3: Evaluation on ImageNet-256×256. Values with underline denote the best except shortcut models, values in **bold** is the best shortcut diffusion model under one-step generation.

Family	Method	Param.	NFE	FID50k
GAN	BigGAN (Brock et al., 2019)	112M	1	6.95
	GigaGAN (Kang et al., 2023)	569M	1	3.45
	StyleGAN-XL (Karras et al., 2019)	166M	1	2.30
AR/Mask	AR w/ VQGAN (Esser et al., 2021)	227M	1024	26.52
	MaskGIT (Chang et al., 2022)	227M	8	6.18
	VAR-d30 (Tian et al., 2024)	2B	10×2	1.92
	MAR-H (Li et al., 2024)	943M	256×2	1.55
Diff/Flow	ADM (Karras et al., 2024)	554M	250×2	10.94
	LDM-4-G (Rombach et al., 2021)	400M	250×2	3.60
	SimDiff (Hoogeboom et al., 2023)	2B	512×2	2.77
	DiT-XL/2 (Peebles & Xie, 2022)	675M	250×2	2.27
	SiT-XL/2 (Ma et al., 2024)	675M	250×2	2.06
	SiT-XL/2+REPA (Yu et al., 2025)	675M	250×2	<u>1.42</u>
Shortcut	iCT (Song & Dhariwal, 2023)	675M	1	34.24
	SCD (Frans et al., 2025)	675M	1	10.60
	IMM (Zhou et al., 2025)	675M	1×2	7.77
	MeanFlow (Geng et al., 2025a)	676M	1	3.43
			2	2.93
	ESC (w/o-class-consist.)	676M	1	2.92
	ESC (w/-class-consist.)	676M	1	2.85
ESC+ (w/-class-consist.)	676M	1	2.53	

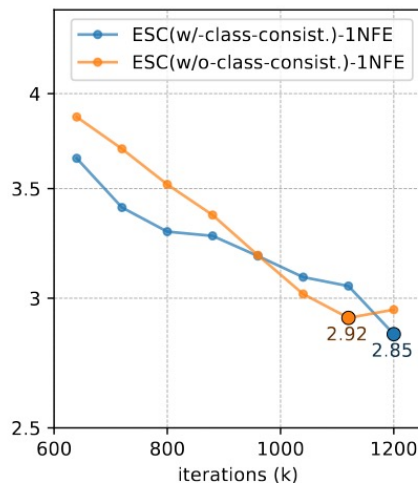
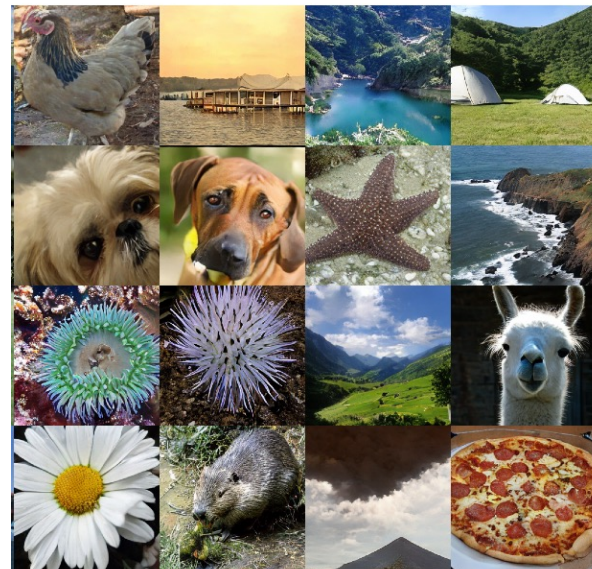


Figure 3: Convergence of FID50k.

Table 4: Uncond. CIFAR-10.

method	NFE	FID
iCT	1	2.83
ECT	1	3.60
sCT	1	2.97
IMM	1	3.20
MeanFlow	1	2.92
ESC	1	2.83

- State-of-the-art Performance
- Minimal time cost
- Faster convergence



ESC with SiT-XL/2 (2.85 FID)



DP Technology
深势科技

Thanks!

Center of AI for Research and Innovation, School of Engineering, Westlake University;

AI for Scientific Simulation and Discovery Lab, School of Engineering, Westlake University;

Academy of Mathematics and Systems Science (AMSS) of the Chinese Academy of Sciences;

DP Technology

Paper:



*Haitao Lin**, *Peiyan Hu**, *Minsi Ren*, *Zhifeng Gao*, *Zhi-Ming Ma*, *Guolin Ke*, *Tailin Wu*, *Stan Z. Li*, ON THE DESIGN OF ONE-STEP DIFFUSION VIA SHORTCUTTING FLOW PATHS, url: <https://openreview.net/pdf?id=k6q8rRYVQR>

Code:



<https://github.com/EDAPINENUT/ExplicitShortCut>