

SPLAM: Accelerating Image Generation with **Sub-Path Linear Approximation Model**

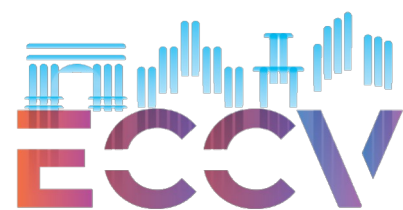
Chen Xu^{1,2*}, Tianhui Song^{1,2*}, Weixin Feng², Xubin Li²,

Tiezheng Ge², Bo Zheng², Limin Wang^{1,3}

¹Nanjing University ²Alibaba Group ³Shanghai AI Lab

* Equal contribution

Oral Presentation



Preliminaries : Effectiveness of Consistency models

- Training for one ideal denoiser \mathbf{D}_θ for \mathbf{x}_t :

$$\text{Minimize } \mathcal{L} = \mathbb{E}[|\mathbf{D}_\theta(\mathbf{x}_t, t) - \alpha(t)\mathbf{x}_0|]$$

- The approximation in Consistency Models :

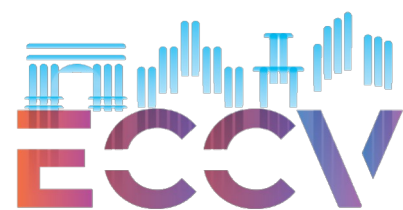
$$\text{One well-learned generator } \mathbf{f}_\theta: \mathbf{f}_\theta(\mathbf{x}_{t-1}, t-1) \approx \mathbf{x}_0$$

- and the approximated learning objective ($\mathbf{D}_\theta = \alpha(t)\mathbf{f}_\theta$):

$$\text{Minimize } \mathcal{L} = \mathbb{E}[|D_\theta(\mathbf{x}_t, t) - \alpha(t)\mathbf{f}_\theta(\mathbf{x}_{t-1}, t-1)|]$$

- also as

$$\mathcal{L}_{\text{Approx}}(\theta) = \mathbb{E}[|\mathbf{x}_t - \frac{\alpha(t)}{\alpha(t-1)}\mathbf{x}_{t-1} + \frac{\alpha(t)}{\alpha(t-1)}\sigma(t-1)\epsilon_\theta(\mathbf{x}_{t-1}, t-1) - \sigma(t)\epsilon_\theta(\mathbf{x}_t, t)|]$$



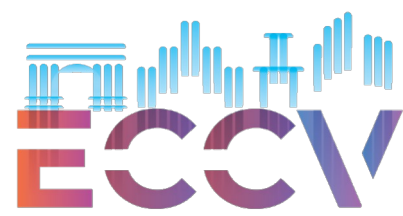
Preliminaries : Effectiveness of Consistency models

- The approximated learning objective for Consistency Models:

$$\mathcal{L}_{Approx}(\theta) = \mathbb{E}[|\mathbf{x}_t - \frac{\alpha(t)}{\alpha(t-1)}\mathbf{x}_{t-1} + \frac{\alpha(t)}{\alpha(t-1)}\sigma(t-1)\epsilon_{\theta}(\mathbf{x}_{t-1}, t-1) - \sigma(t)\epsilon_{\theta}(\mathbf{x}_t, t)|]$$

- The error estimation with the accumulative approximation $\mathbf{f}_{\theta}(\mathbf{x}_{t-1}, t-1) \approx \mathbf{x}_0$

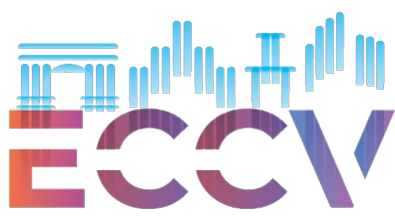
$$|\mathbf{f}_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_0| \leq \sum_{t' \in [1, 2, \dots, t]} |\mathbf{f}_{\theta}(\mathbf{x}_{t'}, t') - \mathbf{f}_{\theta}(\mathbf{x}_{t'-1}, t'-1)|$$



Motivation : Optimization on the upper bound

$$|f_{\theta}(x_t, t) - x_0| \leq \sum_{t' \in [1, 2, \dots, t]} |f_{\theta}(x_{t'}, t') - f_{\theta}(x_{t'-1}, t' - 1)|$$

- Converge slowly when T is large



Motivation : Optimization on the upper bound

$$|f_{\theta}(x_t, t) - x_0| \leq \sum_{t' \in [1, 2, \dots, t]} |f_{\theta}(x_{t'}, t') - f_{\theta}(x_{t'-1}, t' - 1)|$$

↓

$$|f_{\theta}(x_t, t) - x_0| \leq \sum_{t' \in [k, 2k, \dots, t]} |f_{\theta}(x_{t'}, t') - f_{\theta}(x_{t'-k}, t' - k)|$$



Motivation : Optimization on the upper bound

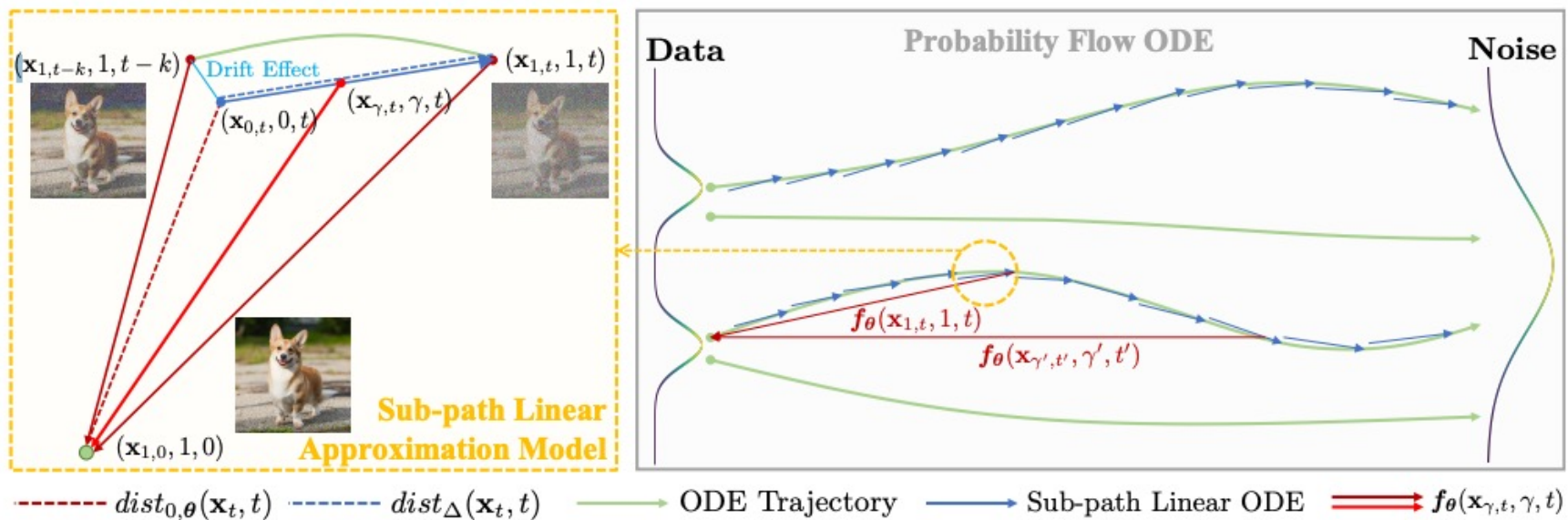
- The changed learning objective

$$\mathcal{L}_{Approx}(\boldsymbol{\theta}) = \mathbb{E} \left[\underbrace{\left| \mathbf{x}_t - \frac{\alpha(t)}{\alpha(t-k)} \mathbf{x}_{t-1} \right|}_{dist_{\Delta}} + \underbrace{\frac{\alpha(t)}{\alpha(t-k)} \sigma(t-k) \epsilon_{\boldsymbol{\theta}}(\mathbf{x}_{t-k}, t-k) - \sigma(t) \epsilon_{\boldsymbol{\theta}}(\mathbf{x}_t, t)}_{dist_{0, \boldsymbol{\theta}}} \right]$$

↑
to be learned

Our SPLAM

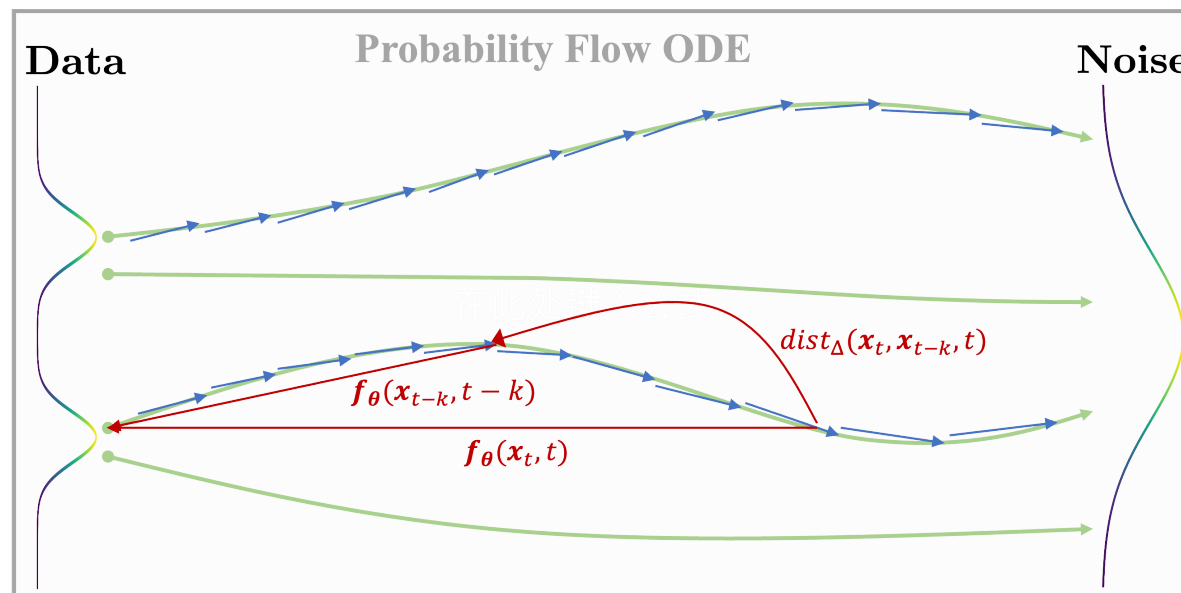
- Treat the Sampled PF-ODE trajectory as a series of connected sub-paths
- Build a better estimation for each sub-path from its start point \mathbf{x}_t to the end \mathbf{x}_{t-k} especially for optimizing $dist_{\Delta}$



Methods

- PF-ODE trajectory as a series of connected sub-paths
 - error estimation for two sample points can be defined as

$$\mathcal{L}_{Sub-p}(\theta, k) = \mathbb{E}[|dist_{\Delta}(\mathbf{x}_t, \mathbf{x}_{t-k}, t) + dist_{0, \theta}(\mathbf{x}_{t-k}, t - k, t) - \sigma(t)\epsilon_{\theta}(\mathbf{x}_t, t)|]$$

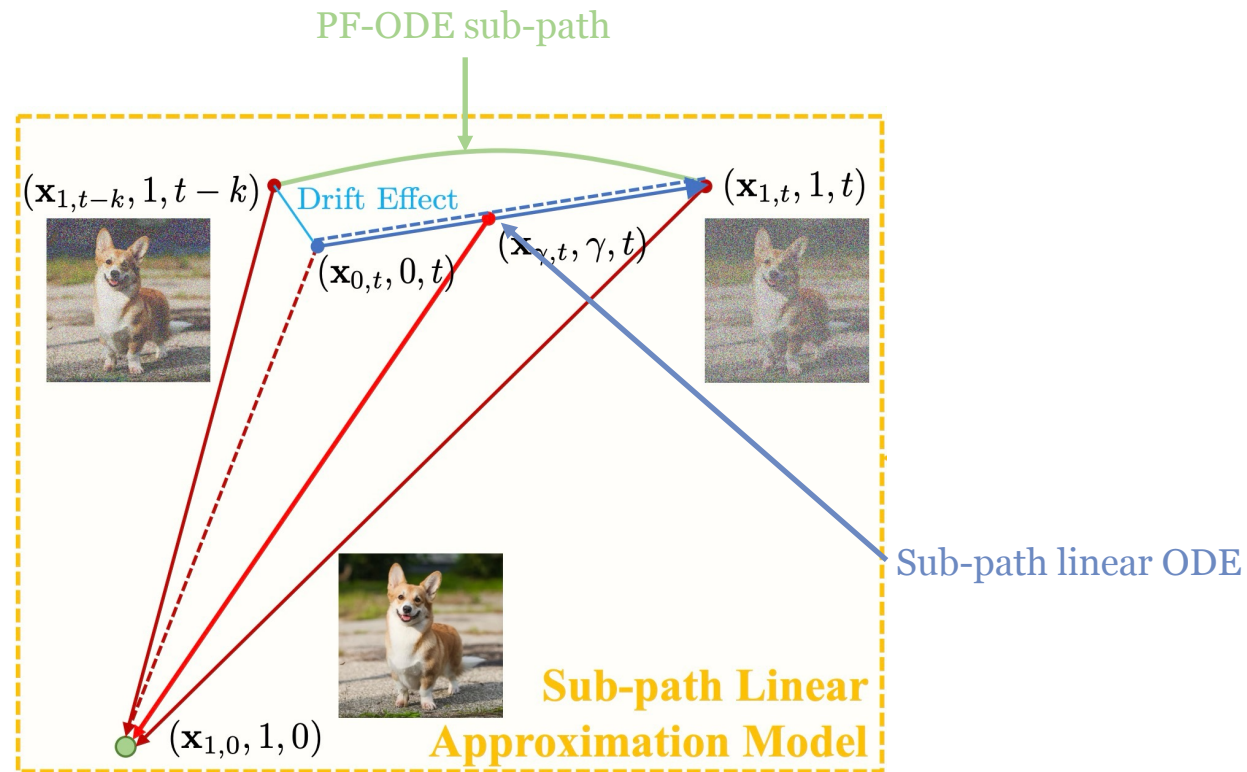


Methods

- We introduce Sub-path Linear ODE $\{x_{\gamma,t}\}_{\gamma \in [0,1]}$

$$\begin{aligned} \mathbf{x}_{\gamma,t} &= \frac{\alpha(t)}{\alpha(t-k)} \mathbf{x}_{t-k} + \gamma * \text{dist}_{\Delta}(\mathbf{x}_t, \mathbf{x}_{t-k}, t) \\ &= (1 - \gamma) \frac{\alpha(t)}{\alpha(t-k)} \mathbf{x}_{t-k} + \gamma \mathbf{x}_t \end{aligned}$$

- from a sub-path start \mathbf{x}_t to a drifted end $\frac{\alpha(t)}{\alpha(t-k)} \mathbf{x}_{t-k}$

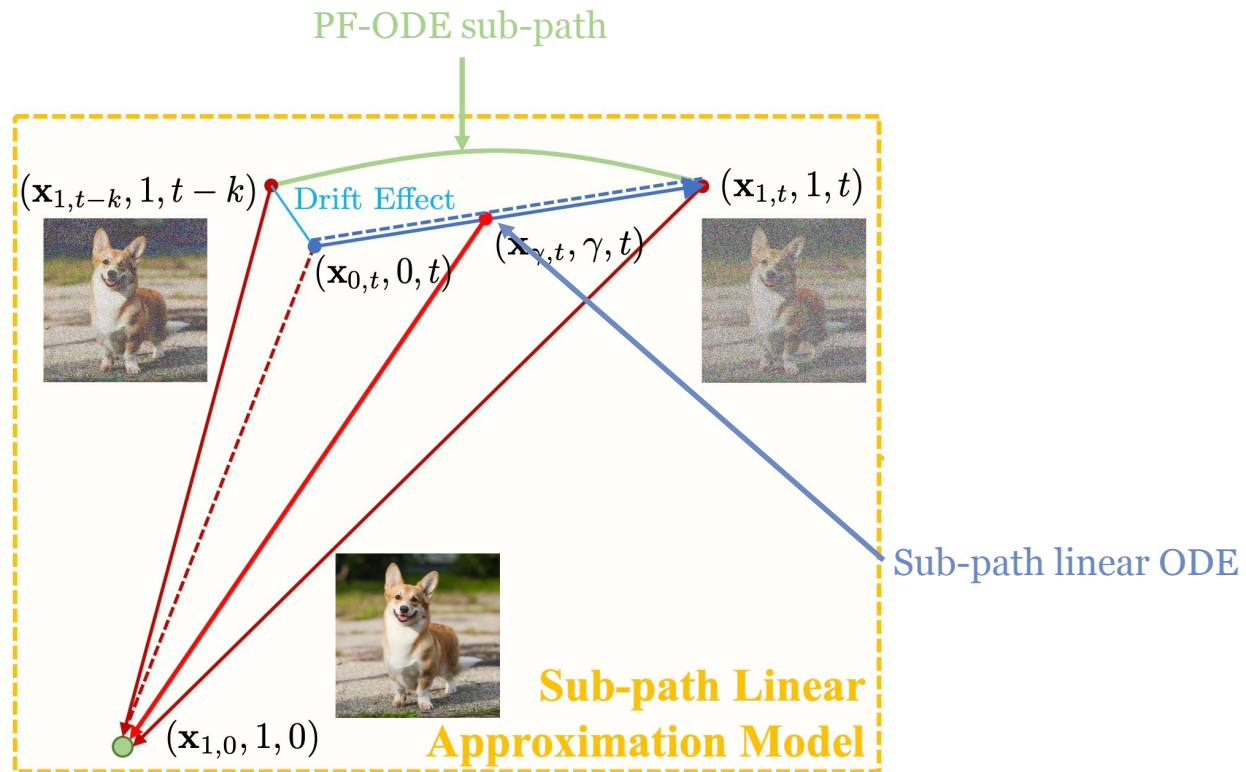


$$x_{\gamma,t} = (1 - \gamma) \frac{\alpha(t)}{\alpha(t-k)} x_{t-k} + \gamma x_t$$

Methods

- Important property for our SL-ODE $\{x_{\gamma,t}\}$

$$dx_{\gamma,t} = \gamma * dist_{\Delta}(x_t, x_{t-k}, t) d\gamma$$



$$x_{\gamma,t} = (1 - \gamma) \frac{\alpha(t)}{\alpha(t-k)} x_{t-k} + \gamma x_t$$

Methods

- Important property for our SL-ODE $\{x_{\gamma,t}\}$

$$dx_{\gamma,t} = \gamma * dist_{\Delta}(x_t, x_{t-k}, t) d\gamma$$



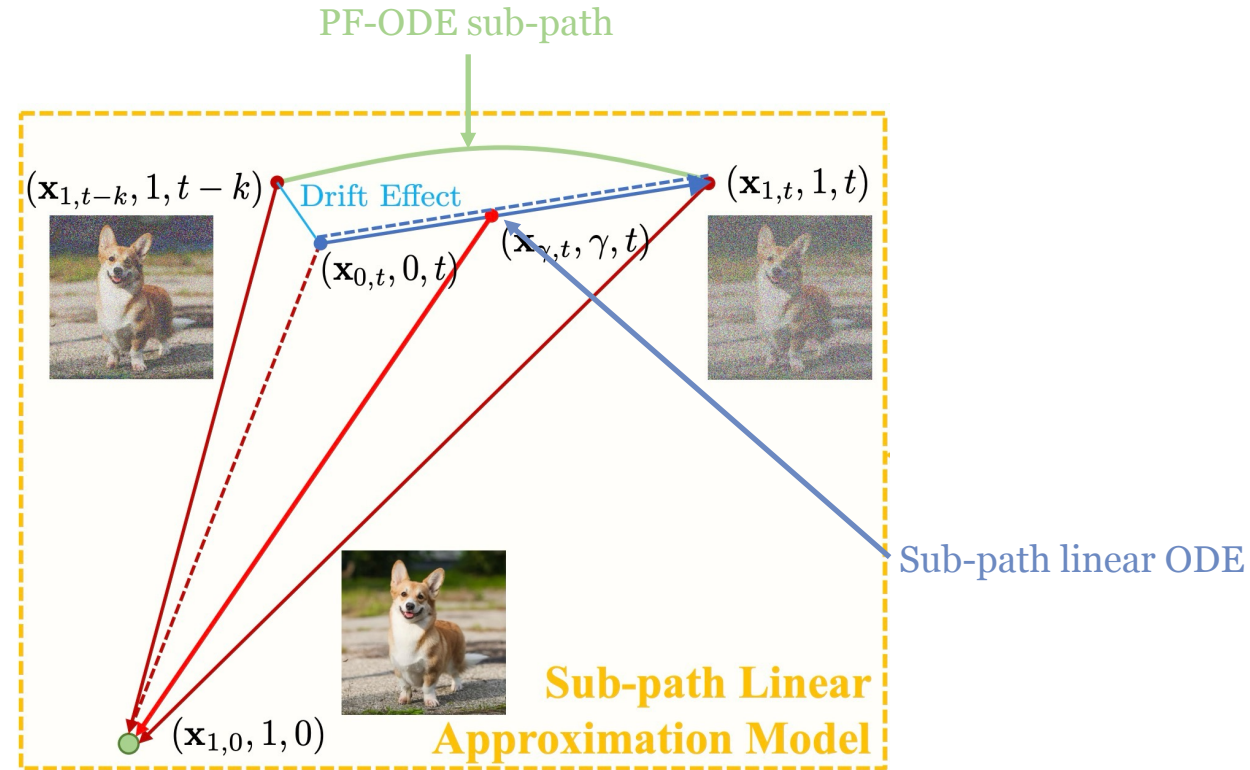
- The error estimation when $x_{\gamma,t}$ is involved:

$$f_{\theta}(x_{\gamma,t}, \gamma, t) = \frac{x_{\gamma,t} - \sigma(\gamma, t)\epsilon_{\theta}(x_{\gamma,t}, \gamma, t)}{\alpha(t)}$$

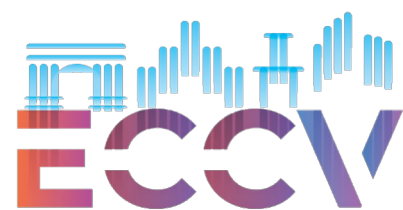
- Using the approximation strategy in CM:

$$\mathcal{L}_{SL-ODE} = |f_{\theta}(x_{\gamma,t}, \gamma, t) - f_{\theta}(x_{1,t-k}, 1, t-k)|$$

$$\text{and } x_{1,t-k} \equiv x_{t-k}$$



$$x_{\gamma,t} = (1 - \gamma) \frac{\alpha(t)}{\alpha(t-k)} x_{t-k} + \gamma x_t$$



Methods

- Important property for our SL-ODE $\{x_{\gamma,t}\}$

$$\mathcal{L}_{SL-ODE} = |\mathbf{f}_{\theta}(\mathbf{x}_{\gamma,t}, \gamma, t) - \mathbf{f}_{\theta}(\mathbf{x}_{1,t-k}, 1, t - k)|$$

- Final learning objective for SPLAM:

$$\mathcal{L}_{SPLAM}(\theta, k) = \mathbb{E}[\gamma * dist_{\Delta}(\mathbf{x}_t, \mathbf{x}_{t-k}, t) + dist_{0,\theta}(\mathbf{x}_{t-k}, t - k, t) - \sigma(\gamma, t)\epsilon_{\theta}(\mathbf{x}_{\gamma,t}, \gamma, t)]]$$

- The original learning objective

$$\mathcal{L}_{Sub-p}(\theta, k) = \mathbb{E}[|dist_{\Delta}(\mathbf{x}_t, \mathbf{x}_{t-k}, t) + dist_{0,\theta}(\mathbf{x}_{t-k}, t - k, t) - \sigma(t)\epsilon_{\theta}(\mathbf{x}_t, t) |]$$

Methods : Distillation from Latent Diffusion Models

Algorithm 1 Sub-Path Linear Approximation Distillation (SPLAD)

Input: dataset \mathcal{D} , initial model parameter θ , learning rate η , EMA decay rate μ , ODE solver $\Phi(\cdot, \cdot; \phi)$, distance estimation $|\cdot|$, a fixed guidance scale w , step size k , VAE encoder $\mathcal{E}(\cdot)$, noise schedule $\alpha(t), \sigma(t)$

$\theta^- \leftarrow \theta$

repeat

sample $(x, c) \sim \mathcal{D}, t \sim \mathcal{U}[k, T]$ and $\gamma \sim \mathcal{U}[0, 1]$ ← Uniform sampling for γ

convert x into latent space: $z = \mathcal{E}(x)$

sample $\mathbf{z}_t \sim \mathcal{N}(\alpha(t)z, \sigma(z)^2 I)$

$\hat{\mathbf{z}}_{t\phi,0}^\Phi \leftarrow \mathbf{z}_t, i \leftarrow 0$

repeat ← Multiple Estimation for ODE solvers

$\hat{\mathbf{z}}_{t\phi,i+1}^\Phi \leftarrow \hat{\mathbf{z}}_{t\phi,i}^\Phi + w\Phi(\hat{\mathbf{z}}_{t\phi,i}^\Phi, t_{\phi,i}, t_{\phi,i+1}, c; \phi) + (1-w)\Phi(\hat{\mathbf{z}}_{t\phi,i}^\Phi, t_{\phi,i}, t_{\phi,i+1}, \emptyset; \phi)$

$i \leftarrow i + 1$

until $k = i * k_\phi$

$\mathbf{z}_{\gamma,t} \leftarrow (1 - \gamma) * \frac{\alpha(t)}{\alpha(t-k)} \hat{\mathbf{z}}_{i-k}^\Phi + \gamma * \mathbf{z}_t$ ▷ Sample a point on the SL-ODE.

$\mathcal{L}(\theta, \theta^-; \phi) \leftarrow |(\mathbf{F}_\theta(\mathbf{z}_{\gamma,t}, c, \gamma, t) - \mathbf{F}_{\theta^-}(\hat{\mathbf{z}}_{1,t-k}^\Phi, c, 1, t - k))|$

$\theta \leftarrow \theta - \eta \nabla_\theta \mathcal{L}(\theta, \theta^-; \phi)$

$\theta^- \leftarrow \text{stopgrad}(\mu\theta^- + (1 - \mu)\theta)$

until convergence

Experiment Results: better FIDs with faster convergence

Table 2: Quantitative results for SDv1.5. Baseline numbers are cited from [50] and [49]. All the results of LCM are our reproduction whose performance is aligned as stated in the paper. † Results are evaluated by us using the released models.

(a) Results on MSCOCO2014-30k, $w = 3$.

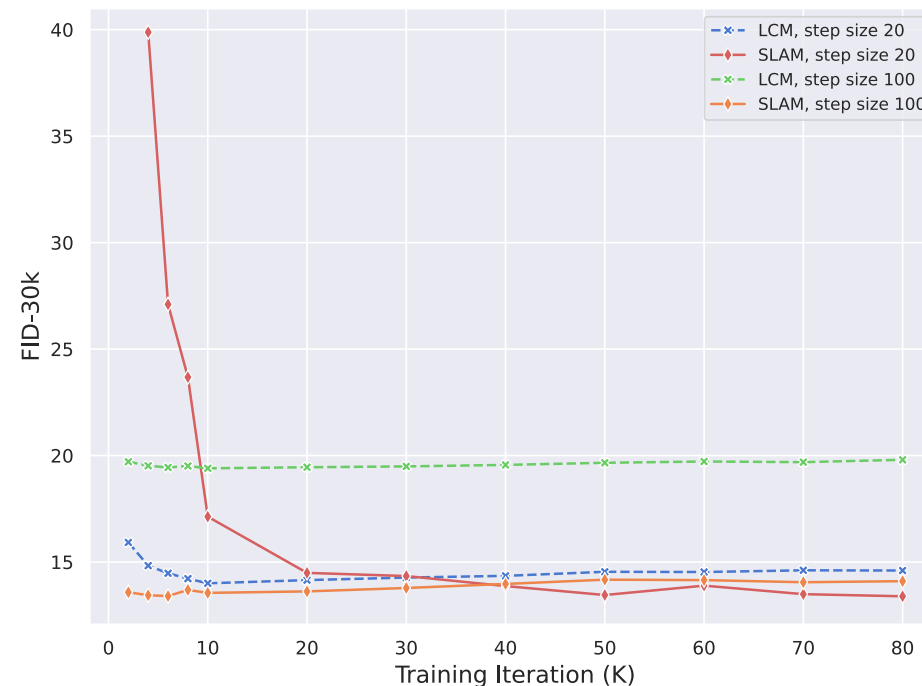
Family	Methods	Latency(↓)	FID(↓)
Unaccelerated	DALL-E [33]	-	27.5
	DALL-E2 [32]	-	10.39
	Parti-750M [51]	-	10.71
	Parti-3B [51]	6.4s	8.10
	Parti-20B [51]	-	7.23
	Make-A-Scene [5]	25.0s	11.84
	Muse-3B [4]	1.3	7.88
	GLIDE [29]	15.0s	12.24
	LDM [34]	3.7s	12.63
	Imagen [35]	9.1s	7.27
eDiff-I [1]	32.0s	6.95	
GANs	LAFITE [54]	0.02s	26.94
	StyleGAN-T [38]	0.10s	13.90
	GigaGAN [13]	0.13s	9.09
Accelerated Diffusion	DPM++ (4step) [24]	0.26s	22.36
	UniPC (4step) [52]	0.26s	19.57
	LCM-LoRA (4step) [27]	0.19s	23.62
	InstaFlow-0.9B [21]	0.09s	13.10
	InstaFlow-1.7B [21]	0.12s	11.83
	UFOGen [49]	0.09s	12.78
	DMD [50]	0.09s	11.49
	LCM (2step) [26]	0.12s	14.29
	SPLAM (2step)	0.12s	12.31
	LCM (4step) [26]	0.19s	10.68
SPLAM (4step)	0.19s	10.06	
Teacher	SDv1.5 [34]†	2.59s	8.03

(b) Results on MSCOCO2017-5k, $w = 3$.

Methods	#Step	Latency(↓)	FID(↓)
DPM Solver++ [24]†	4	0.21s	35.0
	8	0.34s	21.0
	1	0.09s	37.2
Progressive Distillation [36]	2	0.13s	26.0
	4	0.21s	26.4
CFG-Aware Distillation [16]	8	0.34s	24.2
InstaFlow-0.9B [21]	1	0.09s	23.4
InstaFlow-1.7B [21]	1	0.12s	22.4
UFOGen [49]	1	0.09s	22.5
LCM [26]	2	0.12s	25.22
	4	0.19s	21.41
SPLAM	2	0.12s	23.07
	4	0.19s	20.77

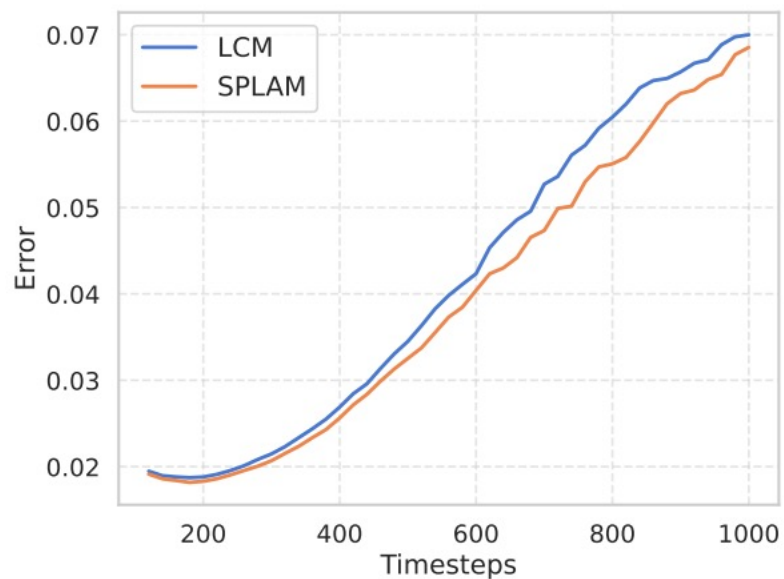
(c) Results on MSCOCO2014-30k, $w = 8$.

Family	Methods	Latency(↓)	FID(↓)
Accelerated Diffusion	DPM++ (4step)	0.26s	22.44
	UniPC (4step) [52]	0.26s	23.30
	LCM-LoRA (4step) [27]	0.19s	23.62
	DMD [50]	0.09s	14.93
Accelerated Diffusion	LCM (2step) [26] [26]	0.12s	15.56
	SPLAM (2step)	0.12s	14.50
	LCM (4step) [26] [26]	0.19s	14.53
	SPLAM (4step)	0.19s	13.39
Teacher	SDv1.5 [34]†	2.59s	13.05





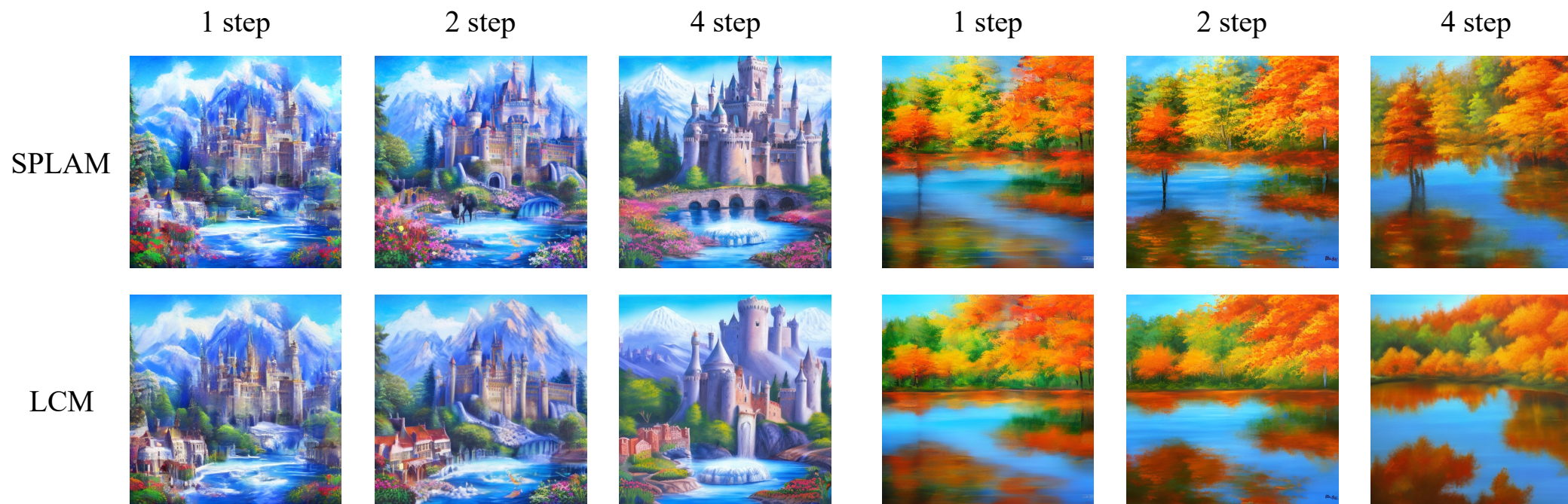
Experiment Results: Directly estimate the errors



$$|x_0 - f_\theta(x_t, t)| \leq \sum_{t' \in [k, 2k, \dots, t]} |f_\theta(x_{t'}, t') - f_\theta(x_{t'-k}, t'-k)|$$



Experiment Results: Generation of images



Experiment Results: Generation of images

"A purple vase with flowers on the table"



"On a fresh summer morning, dew-laden grass glistening in the light, a fawn grazes quietly among the mist"



"A high-speed chase through a cyberpunk metropolis, with advanced motorcycles and holographic billboards"



"A steampunk pirate ship sailing through the clouds with mechanical parrots perched on the masts"





<https://hypnosxc.github.io>
chenxu24568@gmail.com

Thanks



<https://mcg.nju.edu.cn/index.html>
<https://github.com/MCG-NJU>



Paper, code, and models are available