

# Denoising Lévy Probabilistic Models - DLPM

## Heavy-Tailed Diffusion Models

**Umut Simsekli**

joint work with Dario Shariatian, Alain Durmus

March 19, 2026

Generative AI for Extreme Events Workshop

# Introduction on Diffusion Models

## DDPM – Overview

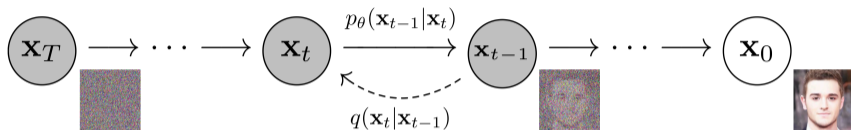


Figure 1: Forward/backward structure, discrete time [HJA20]

### Setup, discrete time:

- **Forward process**  $\{\mathbf{X}_t\}_{t=0}^T$  is a Markov chain with Gaussian transition kernels  $p_{t+1|t}(\cdot|\cdot)$ , such that

$$\mathbf{X}_0 \sim p_0 \text{ (the data)}, \quad \mathbf{X}_T \sim p_T \approx \mathcal{N}(0, \mathbf{I}_d) \text{ (the noise)} \quad (1)$$

- **Generative process**  $\{\bar{\mathbf{X}}_t^\theta\}_{t=0}^T$  will be a Markov chain running in reverse time ; we want  $p_t$  and  $p_t^\theta$  to match
- **Training loss** Fit the joint distributions with an **ELBO loss**, like in VAEs.

## DDPM – Forward Process

- **Forward process (Markov chain):**

$$X_{t+1} = \sqrt{\alpha_t}X_t + \sqrt{1 - \alpha_t}\epsilon_t, \quad (2)$$

where  $\{\alpha_t\}_{t=0}^{T-1}$  is a noise schedule,  $0 < \alpha_t < 1$ .

- **Closed form for  $X_t | X_0$ , by stability of the Gaussian distribution:**

$$X_t | X_0 \stackrel{d}{=} \sqrt{\bar{\alpha}_t}X_0 + \sqrt{(1 - \bar{\alpha}_t)I_d}\bar{\epsilon}_t, \quad (3)$$

with  $\bar{\alpha}_t = \prod_{s=1}^t \alpha_s$ , chosen such that  $X_T$  is approximately distributed as  $\mathcal{N}(0, I_d)$ .

## DDPM – Backward Process

- **Reformulating the forward process** Let us examine its joint distribution

$$\begin{aligned}
 p(x_0, \dots, x_T) &= p_0(x_0) \cdot \prod_{t=1}^T p_{t|t-1}(x_t|x_{t-1}) \\
 &= p_0(x_0) \cdot p_{1|0}(x_1|x_0) \cdot \prod_{t=2}^T p_{t|t-1}(x_t|x_{t-1}, x_0) \\
 &= p_0(x_0) \cdot p_{1|0}(x_1|x_0) \cdot \prod_{t=2}^T \frac{p_{t-1|t,0}(x_{t-1}|x_t, x_0) p_{t|0}(x_t|x_0)}{p_{t-1|0}(x_{t-1}|x_0)} \quad , \text{ by Bayes rule} \\
 &= \underbrace{p_0(x_0)}_{\text{data}} \cdot \underbrace{p_{T|0}(x_T|x_0)}_{\text{noise}} \cdot \prod_{t=2}^T \underbrace{p_{t-1|t,0}(x_{t-1}|x_t, x_0)}_{\text{Gaussian transitions}}
 \end{aligned}$$

- **Gaussian transitions**  $p_{t-1|t,0}(\cdot|x_t, x_0)$  is the density of  $\mathcal{N}(\tilde{m}_t(x_t, x_0), \tilde{\Sigma}_t)$ .

## DDPM – Generative Process

- Backward process

$$p(x_0, \dots, x_t) = \underbrace{p_0(x_0)}_{\text{data}} \cdot \underbrace{p_{T|0}(x_T|x_0)}_{\text{noise}} \cdot \prod_{t=2}^T \underbrace{p_{t-1|t,0}(x_{t-1}|x_t, x_0)}_{\text{Gaussian transitions}}, \quad (4)$$

where  $p_{t-1|t,0}(\cdot|x_t, x_0) = \mathcal{N}(\cdot; \tilde{m}_t(x_t, x_0), \tilde{\Sigma}_t)$ .

- Generative process** This suggests using the following structure for the generative model

$$p^\theta(x_0, \dots, x_t) = \underbrace{p_T^\theta(x_T)}_{\text{noise}} \cdot \prod_{t=1}^T \underbrace{p_{t-1|t}^\theta(x_{t-1}|x_t)}_{\text{Gaussian transitions}}, \quad (5)$$

with  $p_{t-1|t}^\theta(\cdot|x_t) = \mathcal{N}(\cdot; \hat{m}_t^\theta(x_t), \tilde{\Sigma}_t)$ .

## DDPM – Training Objective

- **Variational bound (ELBO)** We want to fit  $p^\theta$  to  $p$ :

$$\begin{aligned}
 \log p_\theta(x_0) &= \log \left( \int p^\theta(X_{0:T}) dX_{1:T} \right) \\
 &\geq \log \left( \mathbb{E}_{p(X_{1:T}|x_0)} \frac{p^\theta(X_{0:T})}{p(X_{1:T}|X_0)} \right) \\
 &\geq \mathbb{E}_{p(X_{1:T})} \log \left( \frac{p^\theta(X_{0:T})}{p(X_{1:T}|X_0)} \right) \quad \text{By Jensen's ineq.} \\
 &= -\mathcal{L}_{\text{ELBO}}(\theta)
 \end{aligned}$$

Rearranging terms, we obtain

$$\mathcal{L}_{\text{ELBO}}(\theta) = \mathbb{E} \left[ \underbrace{\text{KL}(p_{T|0}(\cdot|x_0) \parallel p_T^\theta(\cdot))}_{L_T} + \sum_{t=2}^T \underbrace{\text{KL}(p_{t-1|t,0}(\cdot|x_t, x_0) \parallel p_{t-1|t}^\theta(\cdot|x_t))}_{L_{t-1}} - \underbrace{\log p_{0|1}^\theta(x_0|x_1)}_{L_0} \right].$$

The terms  $L_T, L_0$  are typically neglected.

# DDPM – Training Objective

- **ELBO loss**

$$\mathcal{L}(\theta) = \mathbb{E} \left[ \frac{\|\tilde{\mathbf{m}}_t(X_t, X_0) - \hat{\mathbf{m}}_t^\theta(X_t)\|^2}{2\tilde{\Sigma}_t} \right], \quad (6)$$

- **Denoiser reparameterization** Instead of predicting  $\tilde{\mathbf{m}}_t(x_t, x_0)$  from  $x_t$ , we can predict  $x_0$  or the noise  $\bar{\epsilon}$ . Typically, instead of the true ELBO:

$$\mathcal{L}_{\text{simple}}(\theta) = \mathbb{E}_{X_0, t, \bar{\epsilon}_t, X_t = \bar{\alpha}_t X_0 + \sqrt{1 - \bar{\alpha}_t} \bar{\epsilon}_t} [\|\bar{\epsilon}_t - \hat{\epsilon}_t^\theta(X_t)\|^2]. \quad (7)$$

Thus, the model learns  $\mathbb{E}[\bar{\epsilon}_t | X_t]$ , or  $\mathbb{E}[X_0 | X_t]$ ...

# Advantages

- High quality samples
  
  
  
  
  
  
  
  
  
  
- Stable/easy training (e.g., contrary to GANs)
  
  
  
  
  
  
  
  
  
  
- Equivalence between multiple approaches (continuous time with SDEs, flow matching etc.)





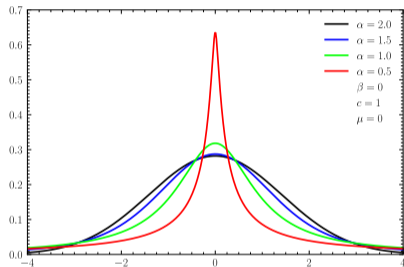




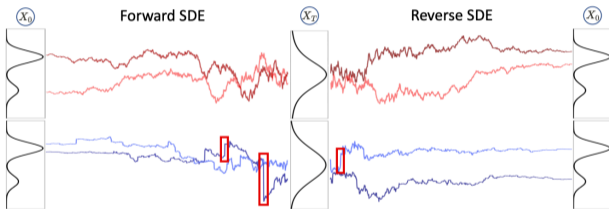




# $\alpha$ -stable Lévy distributions



(a) Symmetric  $\alpha$ -Stable distribution, varying  $\alpha$  [Wik24]

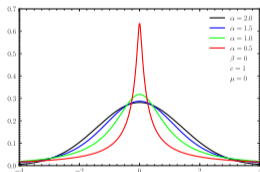


(b) Lévy Process vs Brownian Motion ( $\alpha = 2$ ) [Yoo+23]

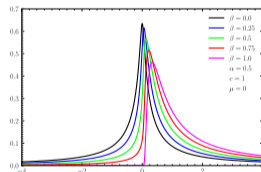
## Definition and properties

The  $\alpha$ -stable distributions  $\mathcal{S}_{\alpha,\beta}(\mu, \sigma)$  are characterized by four parameters  $(\alpha, \beta, \mu, \sigma)$  :

- $\alpha \in (0, 2)$ , the tail heaviness parameter
- $\beta \in (-1, 1)$ , the skewness parameter
- $\mu$ , the location parameter
- $\sigma$ , the scale parameter



(a)  $\beta = 0, \mu = 0, \sigma = 1$ ,  
varying  $\alpha$  [Wik24]



(b)  $\alpha = 0.5, \mu = 0, \sigma = 1$ ,  
varying  $\beta$  [Wik24]



# Stability

- This family of distributions is **stable** by addition, i.e.,

$$X_{S_{\alpha,\beta_0}(\mu_0,\sigma_0)} + X_{S_{\alpha,\beta_1}(\mu_1,\sigma_1)} \sim X_{S_{\alpha,\beta}(\mu,\sigma)}$$

where

$$\sigma^\alpha = \sigma_0^\alpha + \sigma_1^\alpha, \quad \beta = \frac{\beta_0 \sigma_0^\alpha + \beta_1 \sigma_1^\alpha}{\sigma^\alpha}, \quad \mu = \mu_0 + \mu_1$$

- **Gaussian case** ( $\alpha = 2, \beta = 0$ ):

$$\sigma^2 = \sigma_0^2 + \sigma_1^2, \quad \mu = \mu_0 + \mu_1$$

# Gaussian Trick

## Gaussian Trick

Let  $A \sim \mathcal{S}_{\alpha/2,1}(0, c_A)$ , and  $G \sim \mathcal{N}(0, 1)$ , where  $c_A := \cos^{2/\alpha}(\pi\alpha/4)$ . Then

$$A^{1/2}G \sim \mathcal{S}_\alpha(0, 1). \quad (8)$$

- **Isotropic noise.**  $A^{1/2} \cdot G$ .
- **Non-isotropic (independent) noise.** With  $A = \{A_i\}_{i=1}^d$  i.i.d.:  $A^{1/2} \odot G$ .

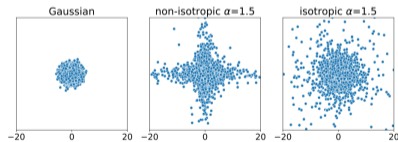


Figure 6: Different multidimensional heavy-tailed noise distributions, Gaussian vs  $\alpha = 1.5$  [Yoo+23]



# Forward Process - first approach

- **Forward process (Markov chain)** Consider  $\{X_t\}_{t=0}^T$  defined by:

$$X_0 \sim p_0, \quad X_t = \gamma_t X_{t-1} + \sigma_t \epsilon_t^{(\alpha)}, \quad (9)$$

where  $\epsilon_t^{(\alpha)} \sim \mathcal{S}_\alpha^i(0, I_d)$  i.i.d., and  $\{(\gamma_t, \sigma_t)\}_{t=1}^T$  is the noising schedule.

- **Closed form for  $X_t|X_0$**

$$X_t \stackrel{d}{=} \gamma_{1 \rightarrow t} X_0 + \sigma_{1 \rightarrow t} \epsilon_t^{(\alpha)}, \quad (10)$$

where  $\epsilon_t^{(\alpha)} \sim \mathcal{S}_\alpha^i(0, I_d)$ .

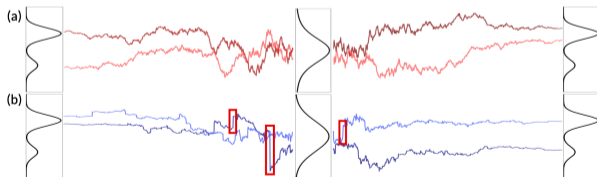


Figure 7: (a) Gaussian transitions (b) Heavy-tailed transitions





# Backward Process – Data Augmentation Approach

- **Conditioning on  $\{A_t\}_{t=1}^T$**  The joint distribution admits the decomposition

$$\begin{aligned}
 p(x_0, \dots, x_T, \mathbf{a}_{1:T}) &= p_0(x_0) \cdot \prod_{t=1}^T p_{t|t-1}(x_t|x_{t-1}, \mathbf{a}_{1:T}) \psi_{(\alpha)}^{\otimes T}(\mathbf{a}_{1:T}) \\
 &= \underbrace{p_0(x_0)}_{\text{data}} \cdot \underbrace{p_{T|0}(x_T|x_0, \mathbf{a}_{1:T})}_{\text{noise}} \cdot \prod_{t=2}^T \underbrace{p_{t-1|t,0}(x_{t-1}|x_t, x_0, \mathbf{a}_{1:T})}_{\text{Gaussian transitions}} \psi_{(\alpha)}^{\otimes T}(\mathbf{a}_{1:T}),
 \end{aligned}$$

where  $\psi_{(\alpha)}$  is the density of  $\mathcal{S}_{\alpha/2,1}(0, cA)$ .

- **Gaussian transitions**  $p_{t-1|t,0,\mathbf{a}_{1:T}}(\cdot|x_t, x_0, \mathbf{a}_{1:T})$  is the density of  $\mathcal{N}(\tilde{\mathbf{m}}_t(x_t, x_0, \mathbf{a}_{1:t}), \tilde{\Sigma}_t(\mathbf{a}_{1:t}))$ .



# Loss function - alpha-stable case

## Reminder: ELBO loss, Gaussian case

$$\mathcal{L}(\theta) = \mathbb{E} \left[ \frac{\|\tilde{\mathbf{m}}_t(X_t, X_0) - \hat{\mathbf{m}}_t^\theta(X_t)\|^2}{2\tilde{\Sigma}_t} \right], \quad (18)$$

- **A naive solution:** by Jensen's inequality:

$$\text{KL}(p_0 \| p_0^\theta) \leq \mathbb{E} (\text{KL} [p_0(\cdot) \| p_{0|a}^\theta(\cdot | A_{1:T})]) . \quad (19)$$

- As we see in (18), this expression would involve taking expectation of  $A_t$

- However,  $A_t$  is distributed as  $\mathcal{S}_{\alpha/2,1}(0, c_A)$ , and does not admit a first order moment.



## Further Design Choices

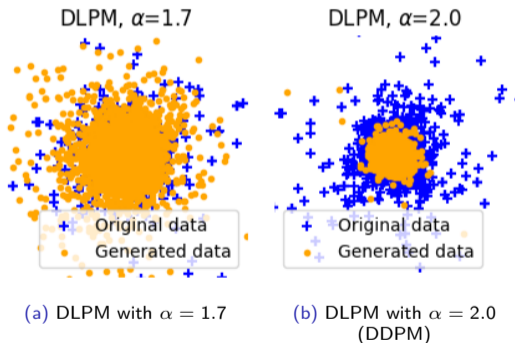


# Experiments



## 2D data - covering the dataset and capturing heavy-tails

- **Challenge:** cover the dataset and correctly capture the tails.
- **Dataset** 20000 samples of  $S_{\alpha}^i(0, 0.05 \cdot I_2)$ , with  $\alpha = 1.7$ .



- The lighter tailed process fails to capture the distribution's tail.

## 2D data - covering the dataset and capturing heavy-tails

- Drawing inspiration from [AGG22], we define the MSLE:

$$\text{MSLE}(\xi) = \int_{\xi}^1 \left( \log \hat{F}^{-1}(p) - \log \hat{F}^{\theta^{-1}}(p) \right)^2 dp, \quad (24)$$

where  $\hat{F}, \hat{F}^{\theta}$  denote respectively the cdf of the true data and the generated data.

Method	$\alpha = 1.5$	$\alpha = 1.6$	$\alpha = 1.7$	$\alpha = 1.8$	$\alpha = 1.9$	$\alpha = 2.0$
DLPM	<b>0.160</b> $\pm$ 0.128	<b>0.081</b> $\pm$ 0.078	<b>0.071</b> $\pm$ 0.028	<b>0.099</b> $\pm$ 0.044	<b>0.132</b> $\pm$ 0.101	0.798 $\pm$ 0.601
DDPM	-	-	-	-	-	0.528 $\pm$ 0.400
LIM	0.743 $\pm$ 0.290	0.497 $\pm$ 0.311	0.267 $\pm$ 0.077	0.653 $\pm$ 0.413	2.444 $\pm$ 1.067	1.239 $\pm$ 0.240
	<i>1.0e-08</i>	<i>8.6e-06</i>	<i>1.3e-10</i>	<i>8.8e-06</i>	<i>7.9e-09</i>	<i>5.0e-3</i>

**Table 1:**  $\text{MSLE}_{\xi=0.95} \downarrow$  averaged over 20 runs. Figures below scores corresponds to  $p$ -values from Welch's  $t$ -test (assuming unequal variances), comparing the mean of DLPM with the given method.

## 2D data - managing class imbalance

- **Challenge:** correctly approximate the mixture weights.  $F_1^{\text{pr}}$  score of precision and recall metrics
- **Dataset** Mixture of nine Gaussian distributions arranged in a grid

$$\sum_{i=1}^9 w_i \mathcal{N}(\mu_i, 0.05^2 \cdot I_2). \quad (25)$$

Mixture weights range from .01 to .3:  $\{.01, .02, .02, .05, .05, .1, .1, .15, .2, .3\}$ .

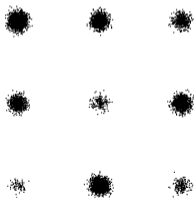


Figure 9: Gaussian grid

Method	$\alpha = 1.5$	$\alpha = 1.6$	$\alpha = 1.7$	$\alpha = 1.8$	$\alpha = 1.9$	$\alpha = 2.0$
DLPM	$0.933 \pm 0.018$	$0.923 \pm 0.005$	$0.933 \pm 0.028$	$0.923 \pm 0.024$	$0.907 \pm 0.034$	$0.862 \pm 0.028$
DLPM <sub>s</sub>	<b><math>0.944 \pm 0.013</math></b>	<b><math>0.943 \pm 0.021</math></b>	<b><math>0.943 \pm 0.010</math></b>	<b><math>0.941 \pm 0.014</math></b>	<b><math>0.928 \pm 0.016</math></b>	-
	<i>9.0e-3</i>	<i>1.6e-05</i>	<i>7.4e-2</i>	<i>9.0e-4</i>	<i>3.9e-3</i>	
LIM	$0.842 \pm 0.039$	$0.850 \pm 0.046$	$0.868 \pm 0.034$	$0.874 \pm 0.030$	$0.884 \pm 0.017$	$0.874 \pm 0.027$
	<i>1.7e-14</i>	<i>1.3e-09</i>	<i>5.7e-11</i>	<i>3.9e-09</i>	<i>1.9e-3</i>	<i>9.6e-2</i>
DDPM	-	-	-	-	-	$0.867 \pm 0.029$
						<i>5.0e-1</i>

Table 2:  $F_1^{\text{pr}}$   $\uparrow$  score, averaged over 30 runs. Figures below scores corresponds to  $p$ -values from Welch's  $t$ -test (assuming unequal variances), comparing the mean of DLPM with the given method.

## 2D data - faster convergence

- **Challenge:** get to the data distribution with the smallest  $T$  possible
- DLIM vs LIM-ODE with varying total diffusion steps  $T$ , on the Gaussian grid.

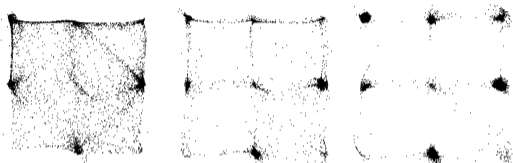


Figure 10: DLIM with  $T = 5, 10, 25$  diffusion steps on the Gaussian grid

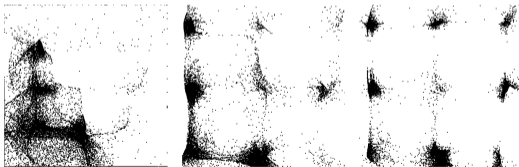


Figure 11: LIM-ODE with  $T = 5, 10, 25$  diffusion steps on the Gaussian grid

# Image data - LIM vs DLPM

- **Dataset** MNIST and CIFAR10\_LT.
- Convergence speed for the different methods, varying total number of diffusion steps  $T$ .

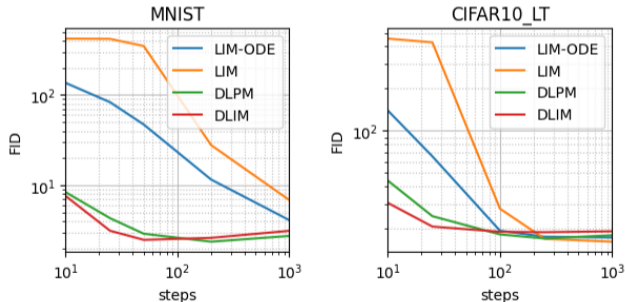


Figure 12: FID $\downarrow$  with varying step size,  $\alpha = 1.7$

## Image data - LIM vs DLPM vs DDPM

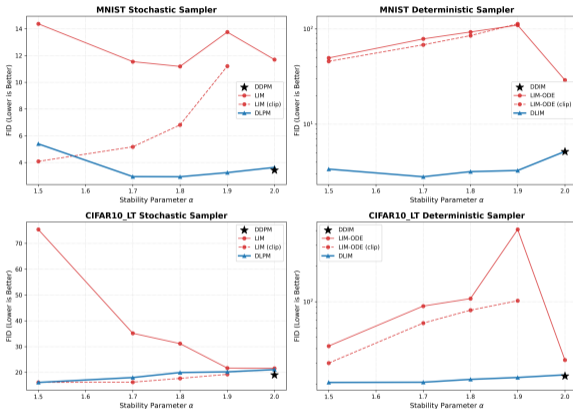


Figure 13: FID $\downarrow$  on MNIST, CIFAR10\_LT for different methods with the stochastic sampler (1000 steps) and the deterministic sampler (25 steps)

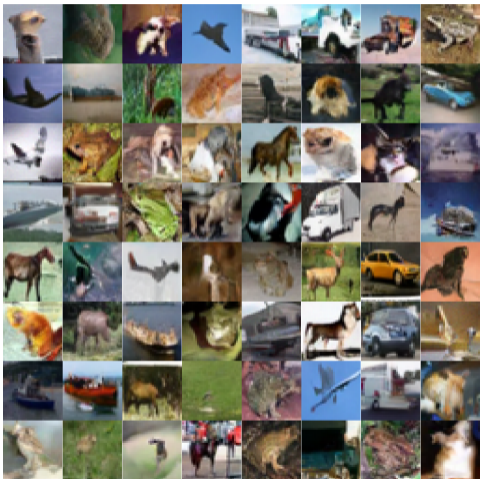
DLPM beats LIM, and smaller  $\alpha$  induce better performance



## Reference II

- [Yoo+23] Eun Bi Yoon et al. “Score-based Generative Models with Lévy Processes”. In: *Advances in Neural Information Processing Systems*. Ed. by A. Oh et al. Vol. 36. Curran Associates, Inc., 2023, pp. 40694–40707. URL: [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/8011b23e1dc3f57e1b6211ccad498919-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/8011b23e1dc3f57e1b6211ccad498919-Paper-Conference.pdf).
- [Wik24] Wikipedia contributors. *Stable distribution* — *Wikipedia, The Free Encyclopedia*. [Online; accessed 2-July-2024]. 2024. URL: [https://en.wikipedia.org/w/index.php?title=Stable\\_distribution&oldid=1227672574](https://en.wikipedia.org/w/index.php?title=Stable_distribution&oldid=1227672574).

# Some images - DLPM



(a) CIFAR10,  $T = 4000$



(b) MNIST,  $T = 1000$

## Some images - DLIM

(a) CIFAR10,  $T = 200$ (b) MNIST,  $T = 50$

# Loss function - design choices **D1**

- **D1 (Fixed variance)** We set  $\hat{\Sigma}_t^\theta = \tilde{\Sigma}_t$ .

## Loss function - design choice **D2**

- **D2 (Denoiser Reparameterization)** We predict the *injected noise*  $\epsilon_t(y_t, y_0)$  rather than  $\tilde{m}_{t-1}(y_t, y_0, a_{1:t})$ . Since

$$\tilde{m}_{t-1}(Y_t, Y_0, A_{1:t}) = \frac{1}{\gamma_t} (Y_t - \sigma_{1 \rightarrow t} \Gamma_t(A_{1:t}) \epsilon_t(Y_t, Y_0)) , \quad (26)$$

we re-parameterize  $\hat{m}_{t-1}^\theta$  as

$$\hat{m}_{t-1}^\theta(Y_t, A_{1:t}) = \frac{1}{\gamma_t} (Y_t - \sigma_{1 \rightarrow t} \Gamma_t(A_{1:t}) \hat{\epsilon}_t^\theta(Y_t)) . \quad (27)$$

with  $\hat{\epsilon}_t^\theta$  the output of the model.

- The model  $\hat{\epsilon}_t^\theta$  does not take any heavy-tailed  $A_{1:t}$  as input.
- Assuming **D1**, the loss  $\mathcal{L}^L$  becomes

$$\mathcal{L}^L(\theta) = \mathbb{E} \left[ \lambda_{t, A_{1:t}}^2 \|\hat{\epsilon}_t^\theta(Y_t) - \epsilon_t(Y_t, Y_0)\|^2 \right] , \quad (28)$$

$$\lambda_{t, a_{1:t}} = \frac{\Gamma_t(a_{1:t}) \sigma_{1 \rightarrow t}}{2\gamma_t \tilde{\Sigma}_{t-1}} , \quad \epsilon_t(Y_t, Y_0) = \frac{(Y_t - \gamma_{1 \rightarrow t} Y_0)}{\sigma_{1 \rightarrow t}} . \quad (29)$$

# Loss function - design choice **D3**

- **D3 (Simple loss)** With design choices **D1**, **D2**, the loss  $\mathcal{L}^L$  is

$$\mathcal{L}^L(\theta) = \mathbb{E} \left[ \lambda_{t, A_{1:t}}^2 \|\hat{\epsilon}_t^\theta(Y_t) - \epsilon_t(Y_t, Y_0)\|^2 \right]. \quad (30)$$

We choose to set  $\lambda_{t, a_{1:t}} = 1$ , which improves performance, and draws similarities to the continuous  $\alpha$ -stable score-based perspective.

We obtain a simplified denoising objective function

$$\mathcal{L}^{\text{Simple}}(\theta) = \mathbb{E} \left[ \mathbb{E} \left( \|\hat{\epsilon}_t^\theta(Y_t) - \epsilon_t(Y_t, Y_0)\|^2 \mid A_{1:t} \right)^{1/2} \right]. \quad (31)$$

## Bonus - faster sampling

Assume design choices **D1**, **D2**, **D3** are satisfied. Then one can obtain the following simplified denoising objective function:

$$\mathcal{L}_{t-1}^{\text{SimpleLess}}(\theta) = \mathbb{E} \left[ \mathbb{E} \left( \left\| \hat{\epsilon}_t^\theta(Y_t^{\text{Less}}) - \epsilon_t(Y_t^{\text{Less}}, Y_0^{\text{Less}}) \right\|^2 \mid \bar{A}_t \right) \right]^{1/2}, \quad t \in \{2, \dots, T\} \quad (32)$$

where

$$Y_t^{\text{Less}} = \gamma_{1 \rightarrow t} Y_0^{\text{Less}} + \sigma_{1 \rightarrow t} \bar{A}_t^{1/2} G_t, \quad \epsilon_t(Y_t^{\text{Less}}, Y_0^{\text{Less}}) = \frac{Y_t^{\text{Less}} - \gamma_{1 \rightarrow t} Y_0^{\text{Less}}}{\sigma_{1 \rightarrow t}}. \quad (33)$$

with  $G_t \sim \mathcal{N}(0, I_d)$ ,  $\bar{A}_t \sim \mathcal{S}_{\alpha/2, 1}(0, c_A)$ .

- Idea: marginalization/sufficient statistic, as

$$Y_t \stackrel{d}{=} \gamma_{1 \rightarrow t} Y_0 + \Sigma_{1 \rightarrow t}(A_{1:t})^{1/2} \bar{G}_t \stackrel{d}{=} \gamma_{1 \rightarrow t} Y_0 + \sigma_{1 \rightarrow t} \epsilon_t^{(\alpha)} \stackrel{d}{=} \gamma_{1 \rightarrow t} Y_0 + \sigma_{1 \rightarrow t} \bar{A}_t^{1/2} G_t \quad (34)$$

- Cheaper than sampling a list  $A_{1:t}$  for each datapoint.
- The final denoising loss is similar to LIM (continuous  $\alpha$ -stable case), but guaranteed to be finite.

## LIM - forward

- **Forward process** The forward process  $\{X_t\}_{0 \leq t \leq T}$ , with  $X_0 \sim p_0$ , is obtained with

$$dX_t = \gamma(t, X_{t-})dt + \sigma(t)dL_t^\alpha, \quad (35)$$

where  $X_{t-}$  denotes the left limit of  $X$  at time  $t$ . **LIM only defines scale-preserving schedule:**

$$\gamma(t, x) = -\frac{\beta_t}{\alpha}x, \quad \sigma(t) = \beta_t^{1/\alpha}. \quad (36)$$

- **Closed-form expression of  $X_t|X_0$**

$$X_t \stackrel{d}{=} \gamma_{1 \rightarrow t} X_0 + \sigma_{1 \rightarrow t} \bar{\epsilon}, \quad (37)$$

where  $\bar{\epsilon}_t \sim \mathcal{S}_\alpha^i(0, I_d)$ . The values of  $\gamma_{1 \rightarrow t}$  and  $\sigma_{1 \rightarrow t}$  match with the DLPD definition on integer timesteps.

## LIM - backward

- **Backward process** The following backward process  $\bar{X}_t$  is obtained:

$$d\bar{X}_t = (-\gamma(t, \bar{X}_{t+}) + \alpha\sigma^\alpha(t, \bar{X}_{t+})s_t(\bar{X}_{t+})) dt + \sigma(t)d\bar{L}_t^\alpha + d\bar{Z}_t \quad (38)$$

where

- $\bar{Z}_t$  is the backward version of a Levy-type stochastic integral  $Z_t$  s.t  $\mathbb{E}[Z_t] = 0$  with finite variation
- $s_t$  is the fractional score function:

$$s_t(x) = \frac{\Delta^{\frac{\alpha-2}{2}} \nabla p_t(x)}{p_t(x)}, \quad (39)$$

where  $\Delta^{\eta/2}$  is the fractional Laplacian of order  $\eta/2$ , defined with Fourier transform  $\mathcal{F}$ :

$$\Delta^{\eta/2} f(x) = \mathcal{F}^{-1}\{\|u\|^\eta \mathcal{F}\{f\}(u)\}. \quad (40)$$

# LIM - training

- The true score  $s_t(x_t|x_0)$  can be expressed as

$$s_t(x_t|x_0) = -\frac{1}{\alpha\sigma_{1\rightarrow t}^{\alpha-1}(t)}\epsilon_t(x_t, x_0), \quad (41)$$

where  $\epsilon_t(x_t, x_0) = \frac{x_t - \gamma_{1\rightarrow t}x_0}{\sigma_{1\rightarrow t}}$ , thus we re-parametrize

$$s_\theta(x_t, t) = -\frac{1}{\alpha\sigma_{1\rightarrow t}^{\alpha-1}(t)}\hat{\epsilon}_t^\theta(x_t, x_0), \quad (42)$$

so that we rather work with  $\hat{\epsilon}_t^\theta$ .

- Training loss obtained using denoising score matching technique:

$$L : \theta \mapsto \mathbb{E}\|s_\theta(X_t, t) - s_t(X_t)\|^2, \quad L' : \theta \mapsto \mathbb{E}\|s_\theta(X_t, t) - s_t(X_t|X_0)\|^2, \quad (43)$$

are equivalent objective functions, with  $s_\theta$  the score approximation given by the model.

## LIM vs DLPM - forward/backward

With  $\{G'_t\}_{t=T}^1$  i.i.d.  $\mathcal{N}(0, I_d)$ ,  $\{\epsilon'_t\}_{t=T}^1$  i.i.d.  $\mathcal{S}_\alpha^i(0, I_d)$ , and  $\hat{\epsilon}_t^\theta$  the model at time  $t$ :

	Stochastic	Deterministic
Continuous (LIM)	$\frac{\bar{X}_t^\theta}{\gamma_t} - \frac{\alpha(1/\gamma_t - 1)}{\sigma_{1 \rightarrow t}^{\alpha-1}} \hat{\epsilon}_t^\theta + \left(\frac{1}{\gamma_t^\alpha} - 1\right)^{1/\alpha} \epsilon'_t$	$\frac{\bar{X}_t^\theta}{\gamma_t} - \left(\frac{\sigma_{1 \rightarrow t}^{1-\alpha}}{\gamma_t} - \sigma_{1 \rightarrow t}^{1-\alpha}\right) \hat{\epsilon}_t^\theta$
Denoising (DLPM)	$\frac{\bar{Y}_t^\theta}{\gamma_t} - \Gamma_t \sigma_{1 \rightarrow t} \hat{\epsilon}_t^\theta + \Gamma_t \Sigma_{1 \rightarrow t-1} G'_t$	$\frac{\bar{Y}_t^\theta}{\gamma_t} - \left(\frac{\sigma_{1 \rightarrow t}}{\gamma_t} - \sigma_{1 \rightarrow t-1}\right) \hat{\epsilon}_t^\theta$

- **Stochastic sampling** Different sampling procedures. Moreover:
  - When  $\alpha = 2$ ,  $0 \leq \Gamma_t \leq 1$  becomes deterministic, and one recovers DDPM formulas
  - $\Gamma_t$  brings additional stochasticity
  - $\Gamma_t$  scales (i) the noise added at time  $t-1$  (ii) the output of the noise model.
- **Deterministic sampling** Different sampling procedures.

# LIM vs DLPM - training

- Alike the Gaussian case ( $\alpha = 2$ ), the score  $s_t(x_t|x_0)$  is a linear expression of the noise term:

$$s_t(x_t|x_0) = -\frac{1}{\alpha\sigma_{1\rightarrow t}^{\alpha-1}(t)}\epsilon_t(x_t, x_0), \quad (44)$$

leading to a similar denoising loss:

$$\mathcal{L}_{t-1} : \theta \mapsto \mathbb{E} \left( \|\hat{\epsilon}_t^\theta(X_t) - \epsilon_t(X_t, X_0)\|_p^\eta \right). \quad (45)$$

- DLPM: use  $p = 2$  and  $\eta = 1$ .
- LIM (theory): use  $p = 2$  and  $\eta = 2$ , for denoising score matching loss equivalence. But  $\epsilon_t(X_t, X_0)$  is heavy-tailed: no variance!
- LIM (experiments): use  $p = 1$  and  $\eta = 1$ . Indicates potential shortcoming of the theoretical approach.

## Denoising Lévy Implicit Models: Deterministic Generation



# Deterministic Generation – Gaussian case (DDIM)

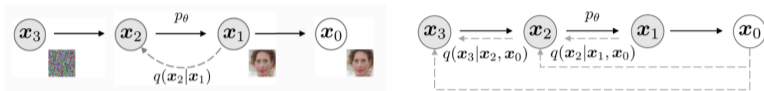


Figure 16: Non-Markovian forward process [SME20]

**Distribution of  $Z_t|Z_0$**  Same as DLPM. Informal proof:

$$\begin{aligned}
 Z_{t-1} &= \gamma_{1 \rightarrow t-1} Z_0 + (\sigma_{1 \rightarrow t-1}^\alpha - \rho_t^\alpha)^{1/\alpha} \cdot \underbrace{\frac{Z_t - \gamma_{1 \rightarrow t} Z_0}{\sigma_{1 \rightarrow t}}}_{\equiv \text{injected noise term } \epsilon_t(Z_t, Z_0)} + \underbrace{\rho_t A_t^{1/2} G_t}_{\text{stochasticity}} \\
 &\stackrel{d}{=} \gamma_{1 \rightarrow t-1} Z_0 + (\sigma_{1 \rightarrow t-1}^\alpha - \rho_t^\alpha)^{1/\alpha} \cdot \underbrace{\epsilon_t^0}_{\equiv \text{injected noise term } \epsilon_t(Z_t, Z_0)} + \underbrace{\rho_t \epsilon_t^1}_{\text{stochasticity}} \\
 &\stackrel{d}{=} \gamma_{1 \rightarrow t-1} Z_0 + \sigma_{1 \rightarrow t-1} \epsilon_t \quad (\text{Stability of } \alpha\text{-stable})
 \end{aligned}$$

# DLIM - Denoising Lévy Implicit Models

- **Recovers DLPM loss** Models  $\hat{\epsilon}_t^\theta(Z_t)$  trained for DLPM can be reused
- **Possibly better loss** Since  $p_{t-1|t,0}$  is  $\alpha$ -stable, we can bypass data-augmentation if closed-form KL exists between  $\mathcal{S}(\mu_1, \sigma_1)$  and  $\mathcal{S}(\mu_2, \sigma_2)$ . It is the case for Cauchy ( $\alpha = 1$ ):

$$\text{KL} [\text{Cauchy}(\mu_1, \sigma) \mid \text{Cauchy}(\mu_2, \sigma)] = \log \left( 1 + \frac{(\mu_1 - \mu_2)^2}{4} \right). \quad (49)$$

- **Deterministic generation** Deterministic sampling process with  $\rho_t = 0$ .