

# Enhancing Accuracy and Efficiency in Diffusion Models

Sony AI @ Tokyo  
Chieh-Hsin (Jesse) LAI

Series Talks in 2024

# My Collaborators



**Chieh-Hsin (Jesse) Lai**

- Researcher/  
DGM Tech Lead @ Sony AI  
- Ph.D. Math, U. of MN



**Dongjun Kim**

- Postdoc @ Stanford  
- Intern @ Sony AI  
- Ph.D. IE, KAIST



**Stefano Ermon**

- Prof. @ Stanford

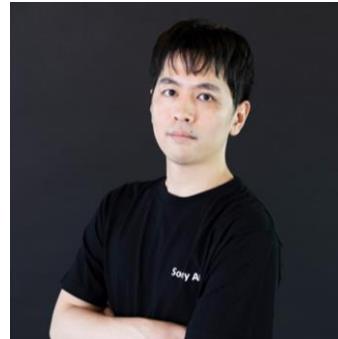
**WeiHsiang Liao**



**Murata Naoki**



**Yuhta Takida**



**Toshimitsu Uesaka**

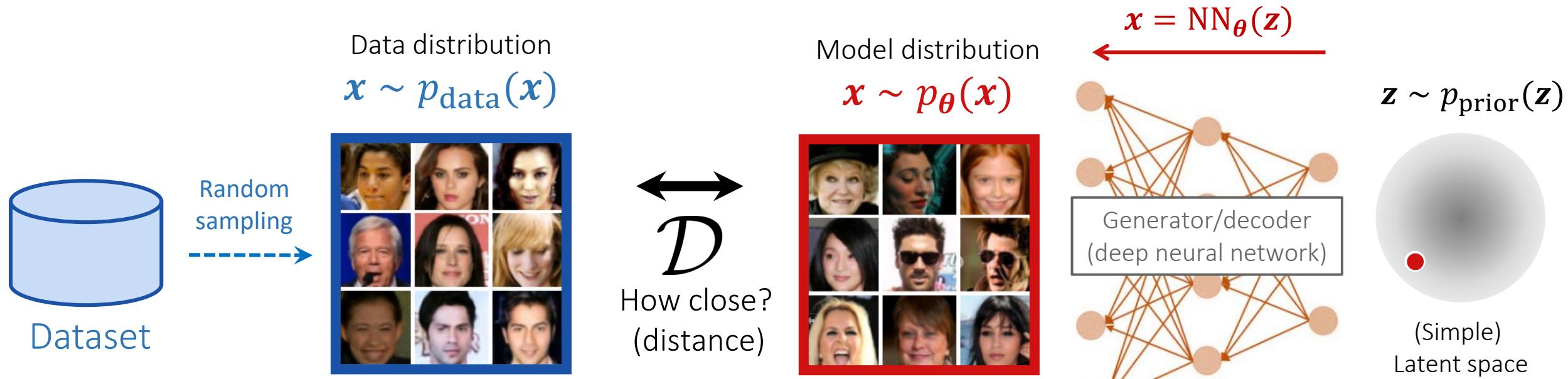


**Yuki Mitsufuji**



# Deep Generative Model

Goal: use NN to learn target distribution explicitly/implicitly



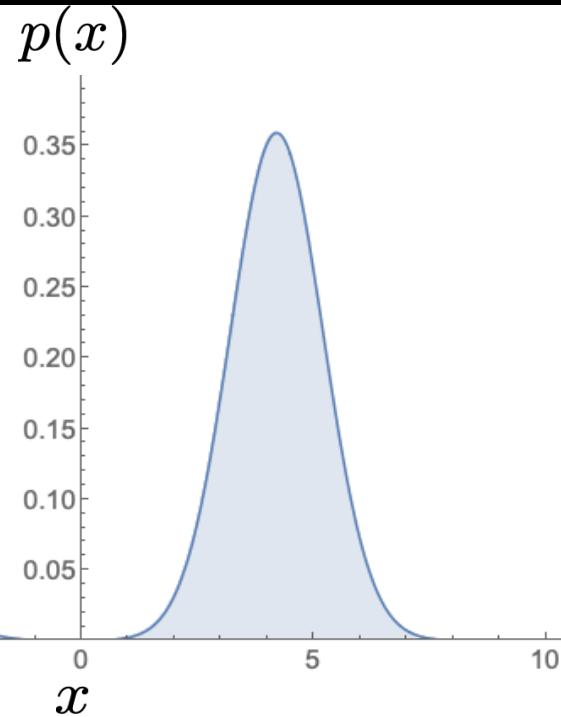
$$\min_{\theta} \mathcal{D}(p_{\theta}(x), p_{\text{data}}(x))$$

# WHAT IS DIFFUSION MODELS?



# Using “Score” (Slope) for Generation

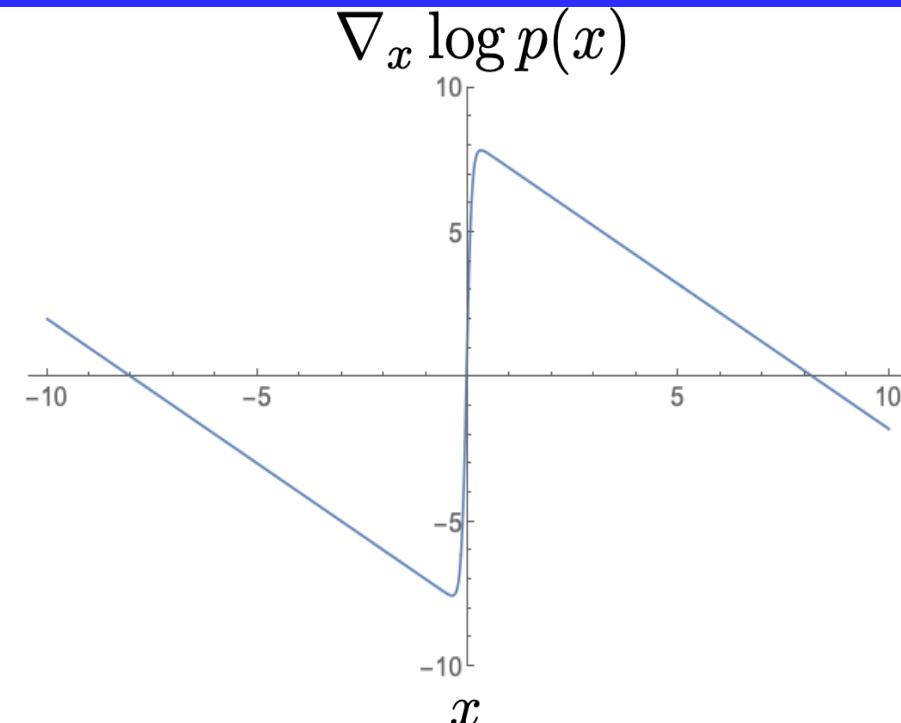
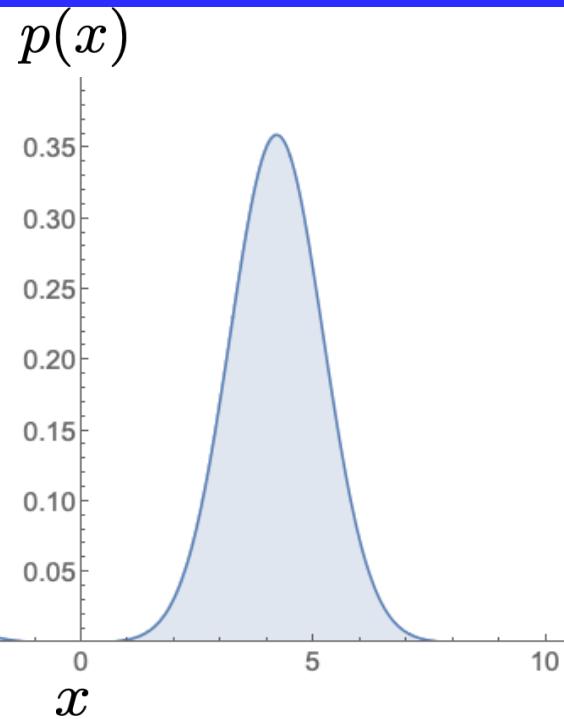
$p_{\text{data}}(x)$



# Using “Score” (Slope) for Generation

Moving toward more likely direction...

$$\nabla_x \log p_{\text{data}}(x)$$



# Using “Score” (Slope) for Generation

Moving toward more likely direction...

$$\nabla_x \log p_{\text{data}}(x)$$

- Phase I. Score Matching
- Phase II. Denoising Score Matching
- Phase III. Noised Conditioned Denoising Score Matching
- Phase IV. Score-Based Continuous Time Diffusion Model

# Phase I (2005). Score Matching

Score function:

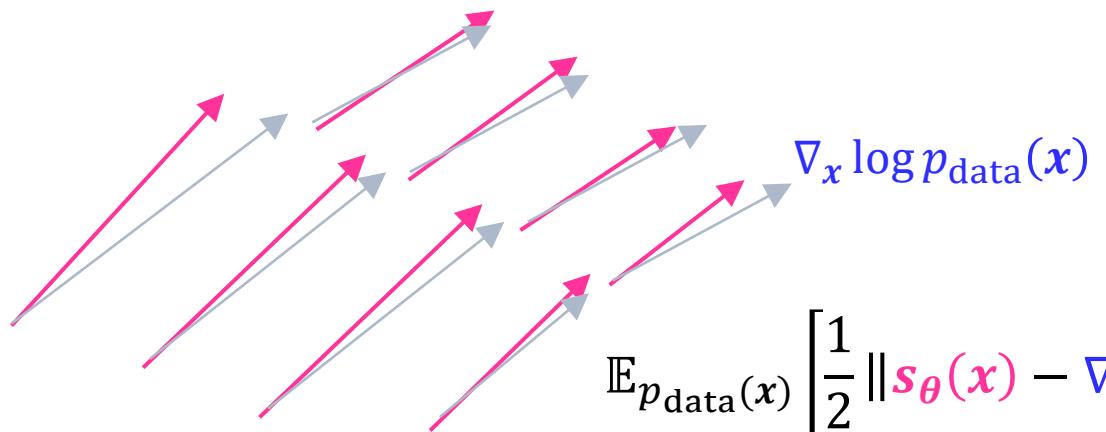
$$\nabla_x \log p_{\text{data}}(x)$$

[Generation] Langevin sampling:

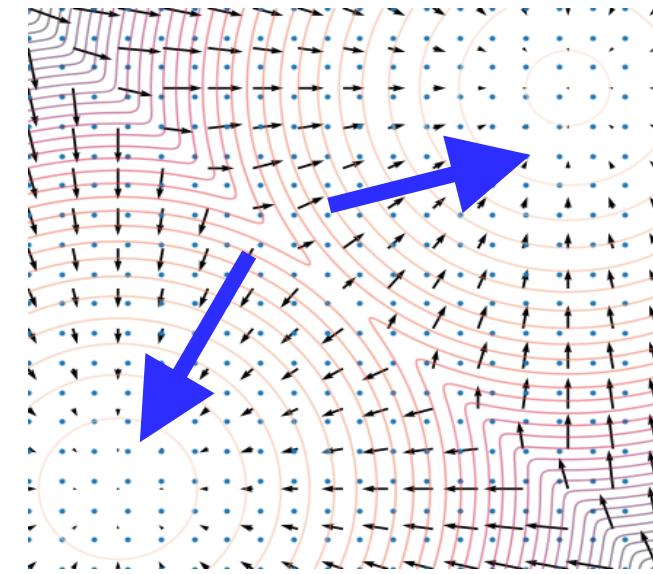
$$x_{t+1} \leftarrow x_t + \eta \nabla_x \log p_{\text{data}}(x) + \sqrt{2} \epsilon_t, \epsilon_t \sim \mathcal{N}(\mathbf{0}, I)$$

[Training] Score Matching:

$$\mathcal{s}_\theta(x)$$



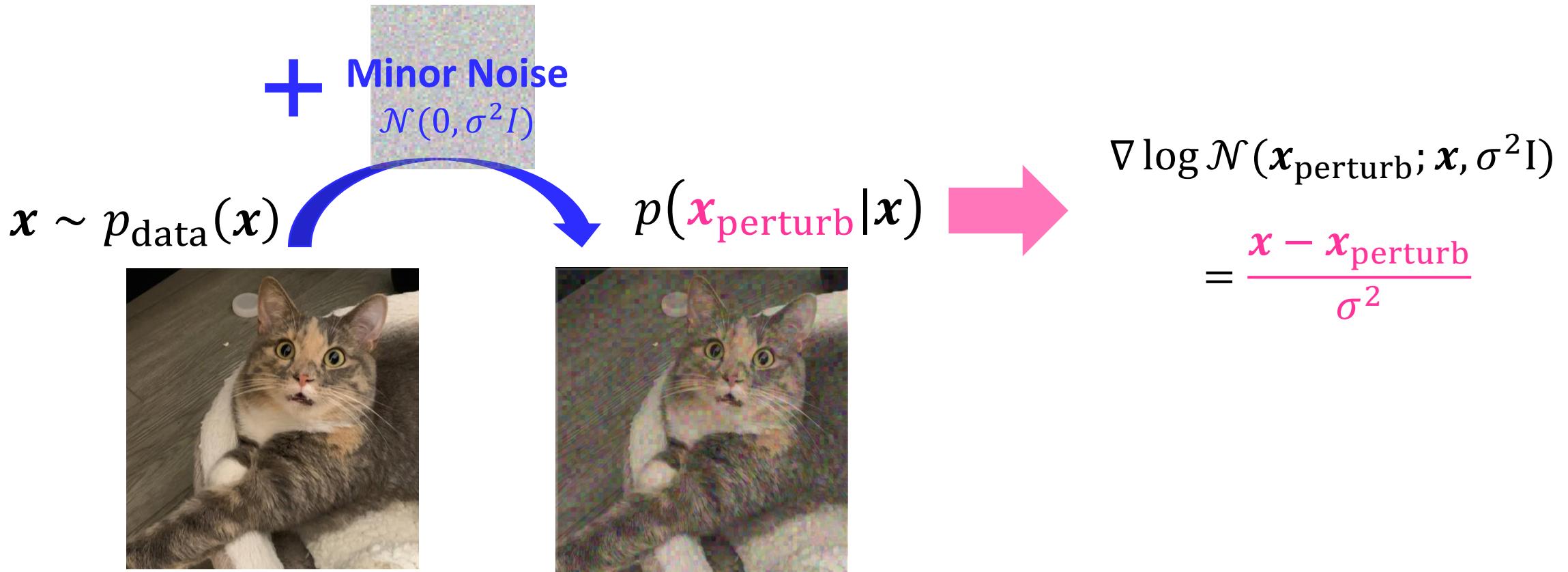
$$\mathbb{E}_{p_{\text{data}}(x)} \left[ \frac{1}{2} \|s_\theta(x) - \nabla_x \log p_{\text{data}}(x)\|_2^2 \right] = \mathbb{E}_{p_{\text{data}}(x)} \left[ \frac{1}{2} \|s_\theta(x)\|_2^2 + \text{tr}(\nabla_x s_\theta(x)) \right] + C$$



# Phase II (2011). Denoising Score Matching

[Training] Denoising Score Matching

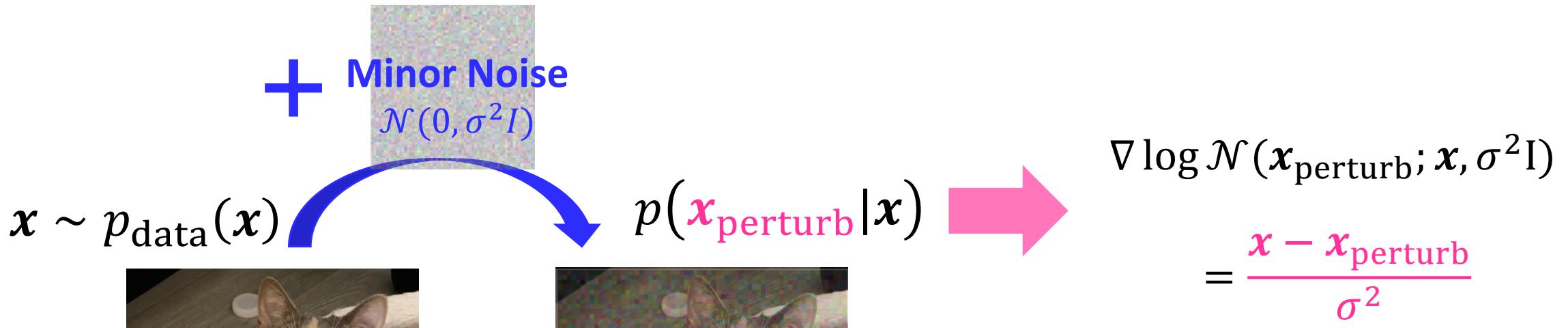
$$\mathbb{E}_{p(x_{\text{perturb}}|x)} \left[ \frac{1}{2} \left\| s_{\theta}(x) - \left( \frac{x - x_{\text{perturb}}}{\sigma^2} \right) \right\|_2^2 \right]$$



# Phase II (2011). Denoising Score Matching

[Training] Denoising Score Matching

$$\mathbb{E}_{p(x_{\text{perturb}}|x)} \left[ \frac{1}{2} \left\| s_{\theta}(x) - \left( \frac{x - x_{\text{perturb}}}{\sigma^2} \right) \right\|_2^2 \right]$$



[Generation] Langevin sampling

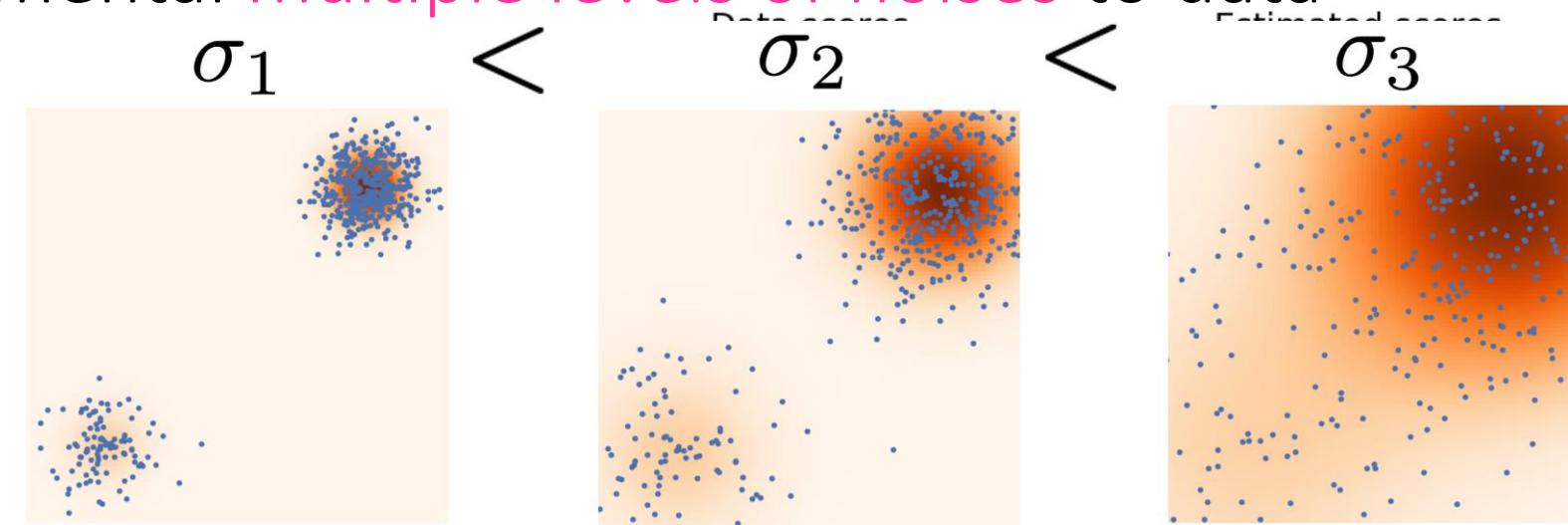
$$x_{t+1} \leftarrow x_t + \eta \nabla \log p_{\text{perturb}}(x_{\text{perturb}}|x) + \sqrt{2} \epsilon_t, \epsilon_t \sim \mathcal{N}(\mathbf{0}, I)$$

$$\approx s_{\theta}(x)$$

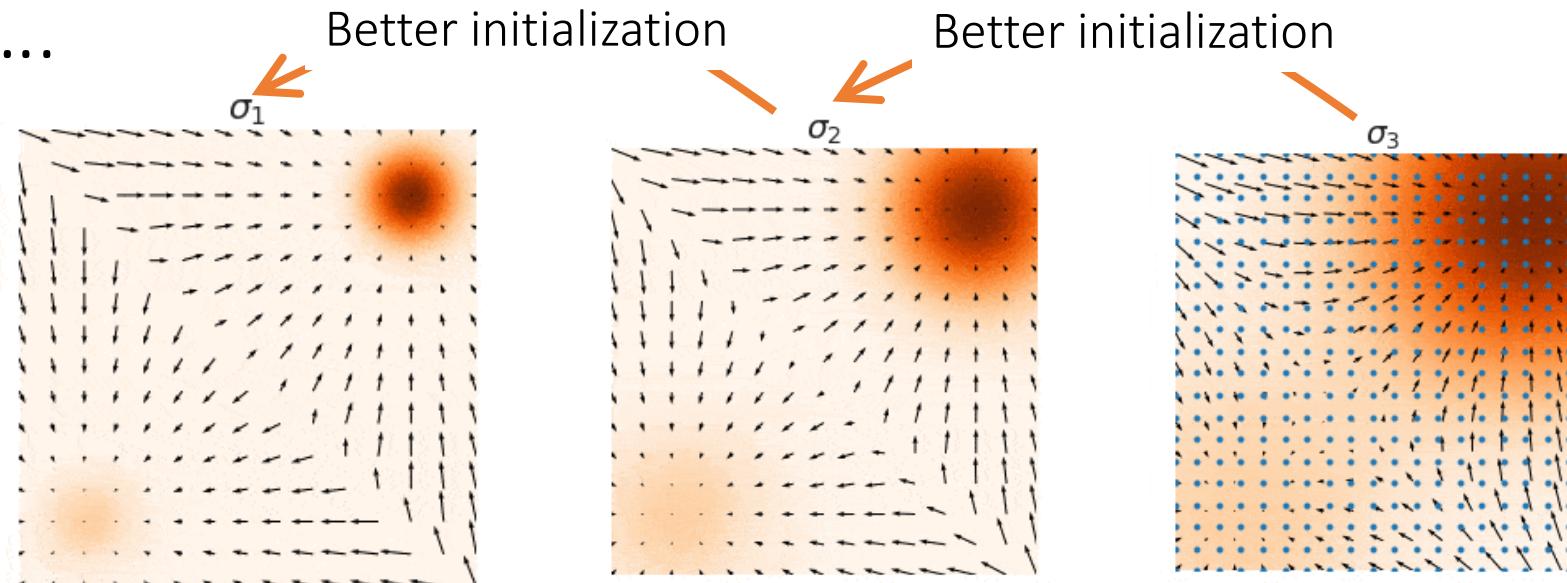


# Phase III (2019). Noised Conditioned Denoising Score Matching

Injecting incremental **multiple levels of noises** to data



In generation....



# Phase III (2019). Noised Conditioned Denoising Score Matching

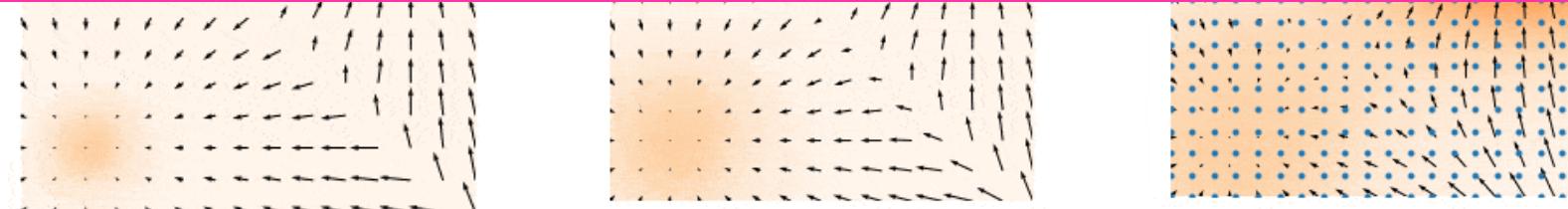
Injecting incremental multiple levels of noises to data

$$\sigma_1 < \sigma_2 < \sigma_3$$

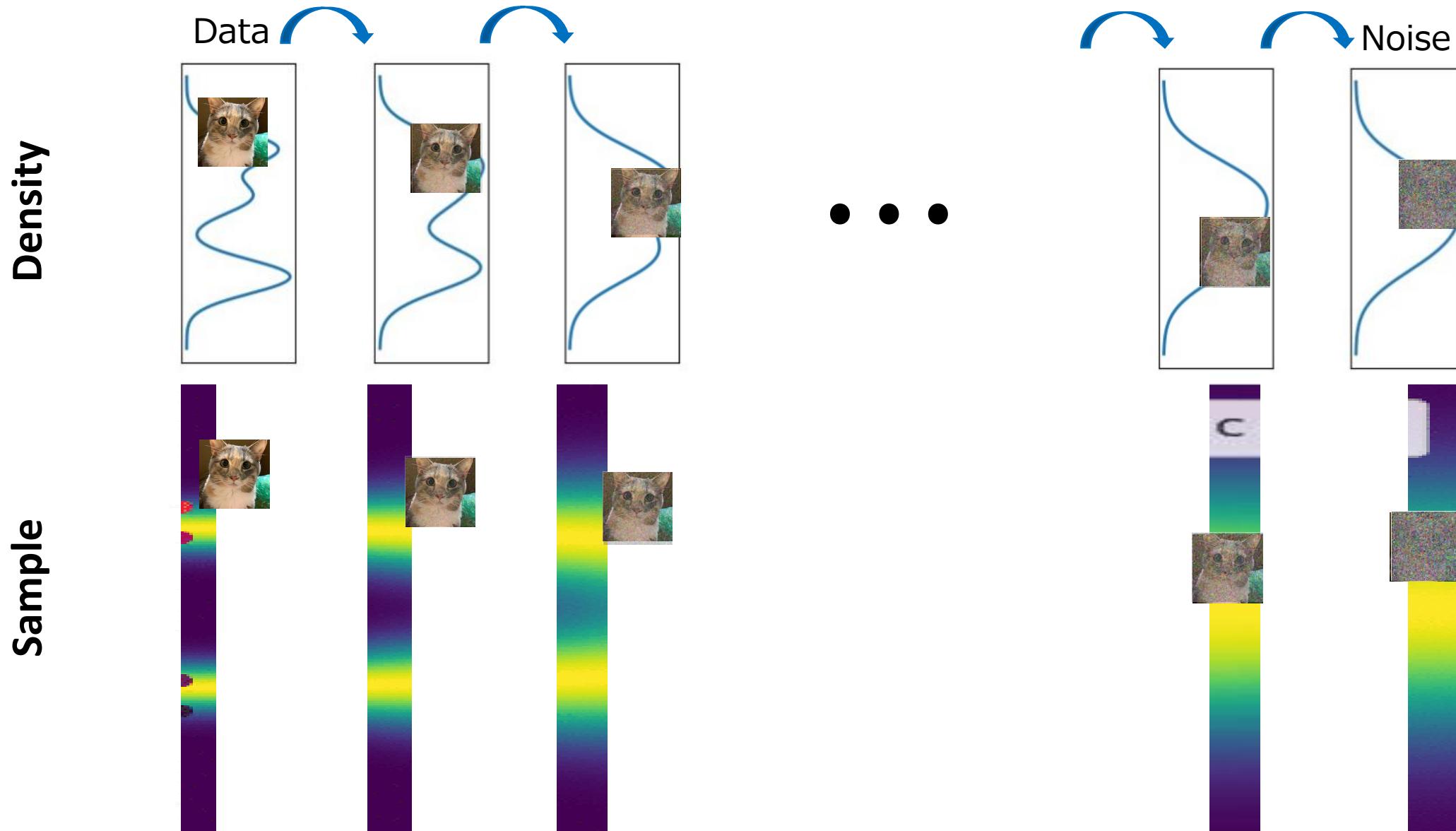


$$\sum_{i=1}^3 \omega_i \mathbb{E}_{p(x_{\text{perturb}, \sigma_i} | x)} \left[ \frac{1}{2} \left\| s_{\theta}(x, \sigma_i) - \left( \frac{x - x_{\text{perturb}, \sigma_i}}{\sigma_i^2} \right) \right\|_2^2 \right]$$

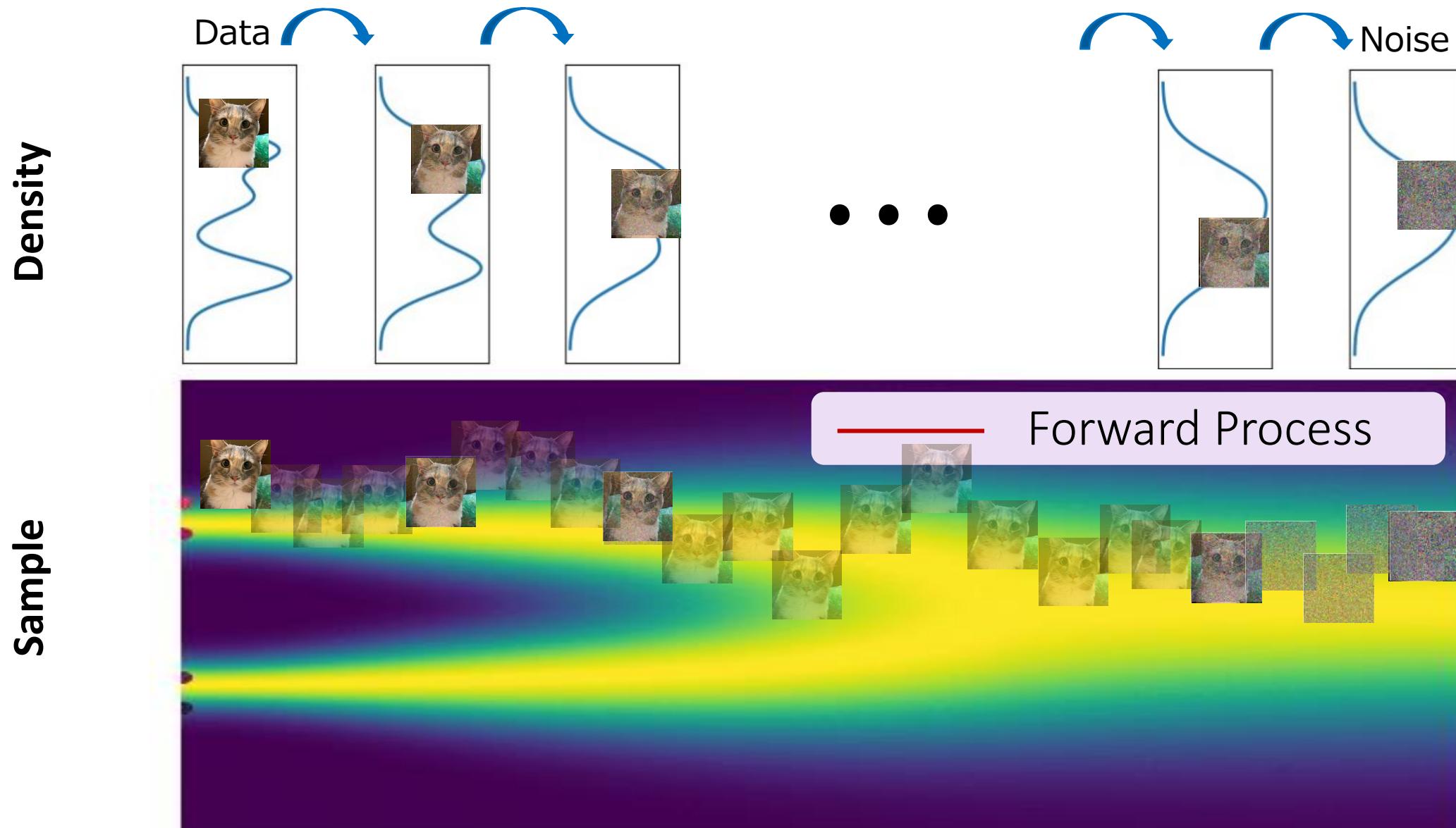
The vibe of Diffusion Model appeared...



# Phase IV (2020). Score-Based Continuous Time Diffusion Model

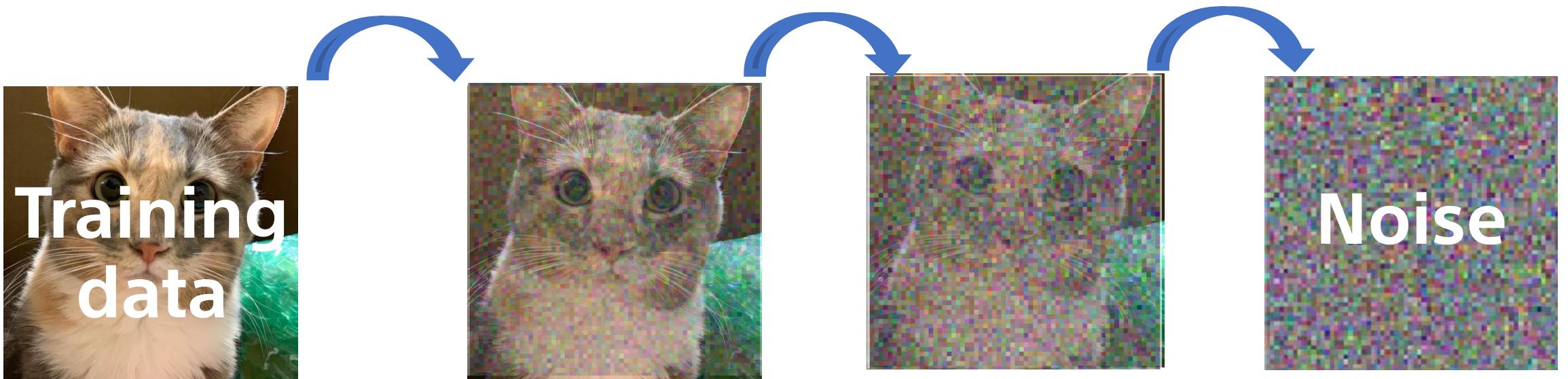


# Phase IV (2020). Score-Based Continuous Time Diffusion Model



# Diffusion Model at High Level

Forward Process: Training data → white noise (simple space)



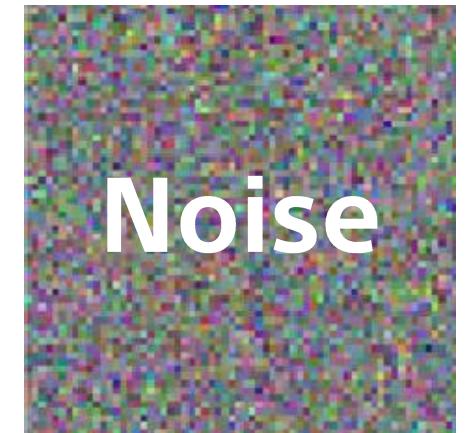
# Diffusion Model at High Level

Backward Process: Learn NN to gradually denoise

Generated new image

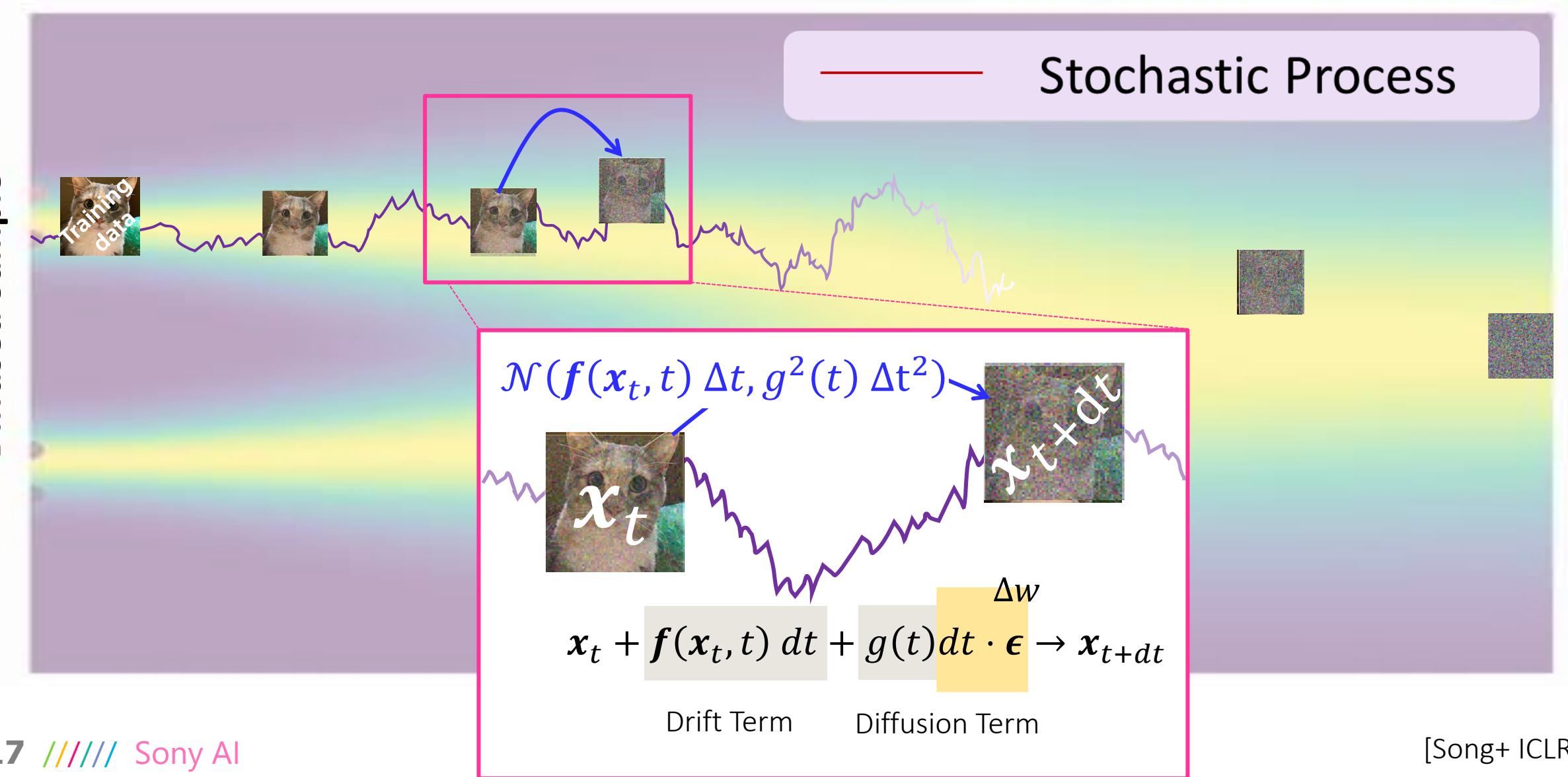


Gradually transforming noise to a new image



# Diffusion Model's Forward Process

Users predefined noise schedulers:  $f(x_t, t)$  &  $g(t)$



# Diffusion Model's Forward Process

Users predefined noise schedulers:  $f(x_t, t)$  &  $g(t)$

Stochastic Process

Getting finer timesteps ( $\Delta t \rightarrow 0$ ) leads to

$$dx_t = f(x_t, t)dt + g(t)d\omega_t,$$

with  $d\omega_t \sim \mathcal{N}(0, dt)$



$$x_t + f(x_t, t) dt + g(t)dt \cdot \epsilon \rightarrow x_{t+dt}$$

Drift Term

Diffusion Term

# Diffusion Model's Forward Process

Users predefined noise schedulers:  $f(x_t, t)$  &  $g(t)$

Stochastic Process

For example,

$$f(x_t, t) = 0 \text{ & } g(t) = \sqrt{2t}$$


$$x_t + f(x_t, t) dt + g(t)dt \cdot \epsilon \rightarrow x_{t+dt}$$

Drift Term      Diffusion Term

$\Delta w$

# Diffusion Model's Forward Process

Users predefined noise schedulers:  $f(x_t, t)$  &  $g(t)$

Stochastic Process

Getting finer timesteps ( $\Delta t \rightarrow 0$ ) leads to

$$dx_t = \sqrt{2t} dw_t,$$

with  $dw_t \sim \mathcal{N}(0, dt)$

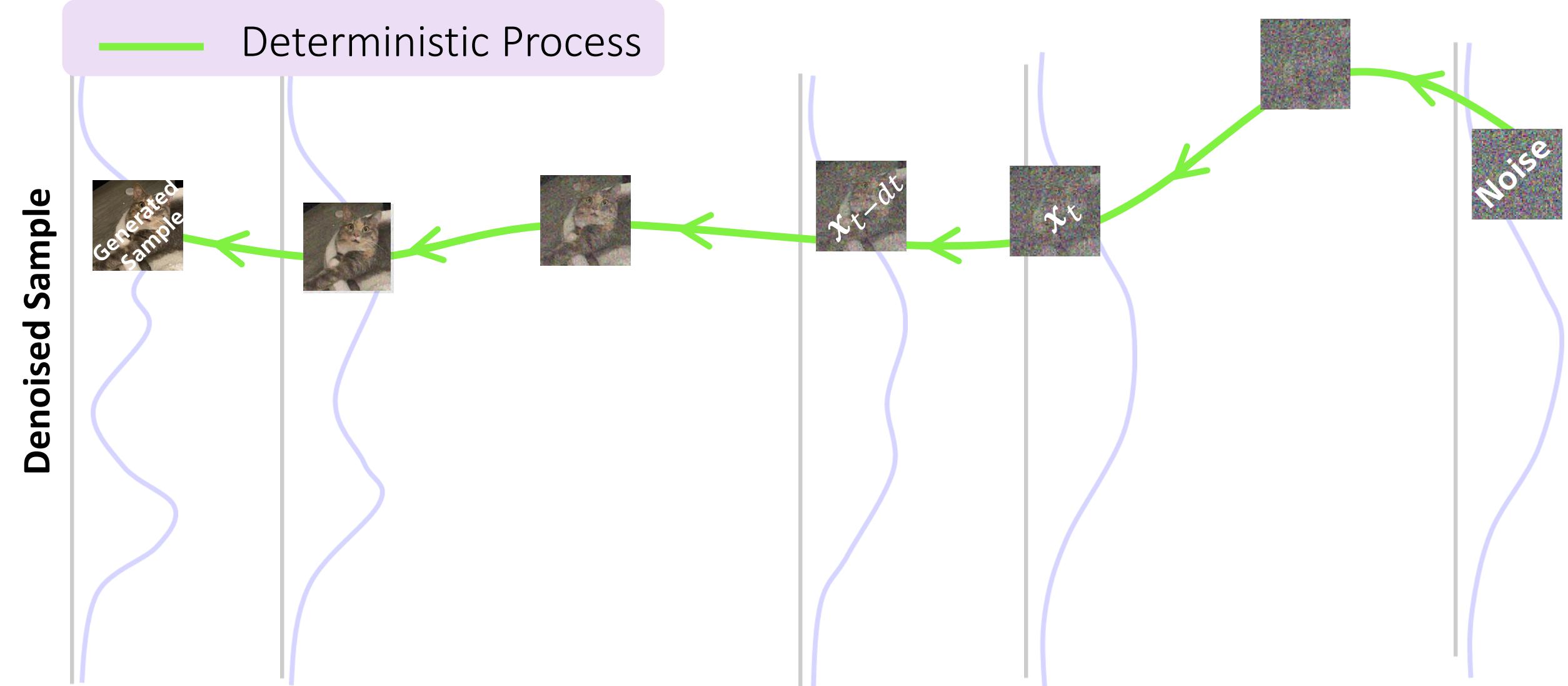


$$x_t + f(x_t, t) dt + g(t)dt \cdot \epsilon \rightarrow x_{t+dt}$$

Drift Term

Diffusion Term

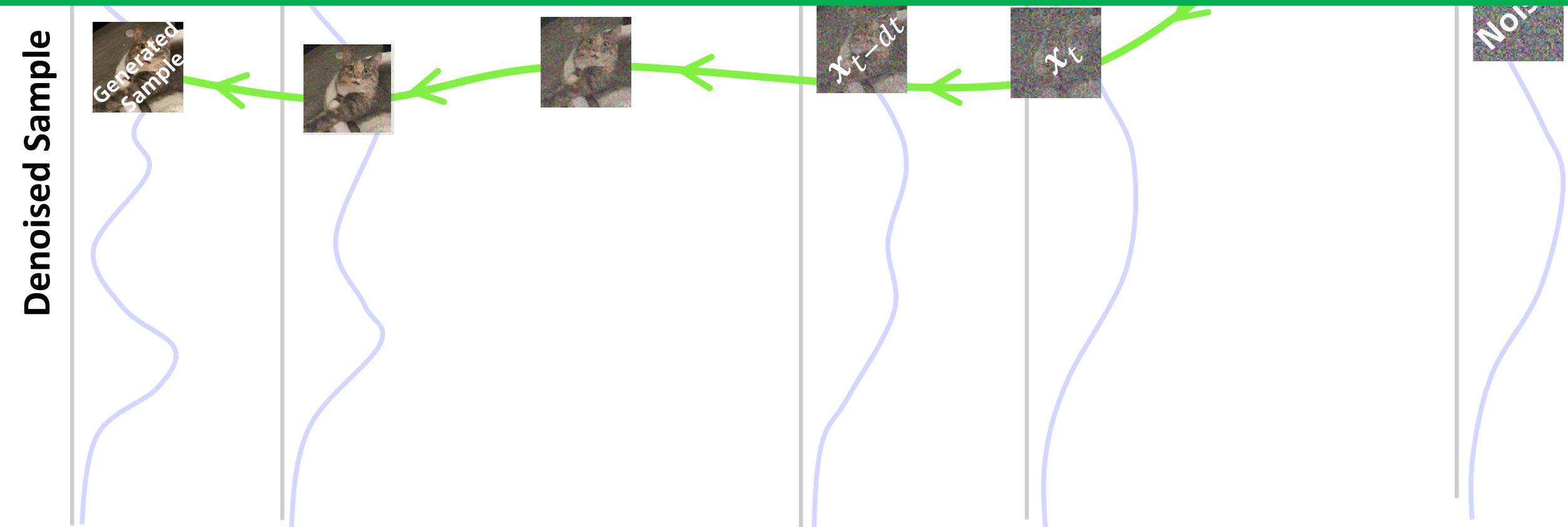
# Diffusion Model's Backward Deterministic Process (Generation)



# Diffusion Model's Backward Deterministic Process (Generation)

Getting finer timesteps ( $dt \rightarrow 0$ ) leads to

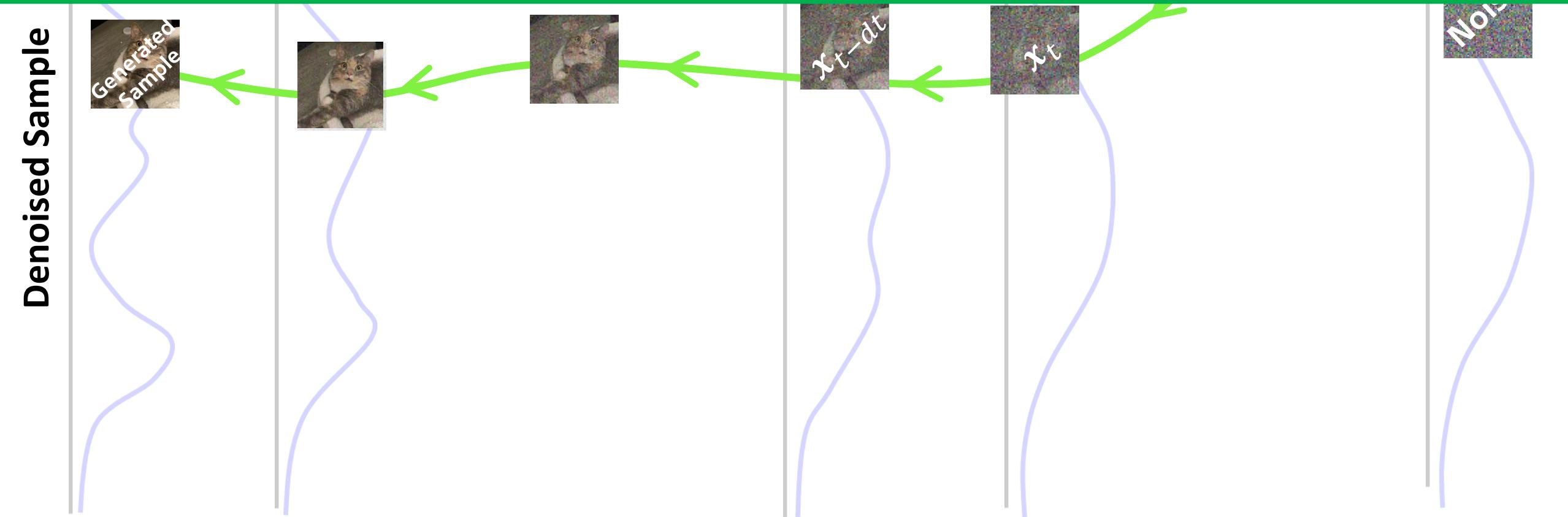
$$\frac{dx_t}{dt} = f(t) - \frac{1}{2} g^2(t) \nabla \log p_t(x_t), x_T \sim p_{\text{prior}}$$



# Diffusion Model's Backward Deterministic Process (Generation)

In our example,

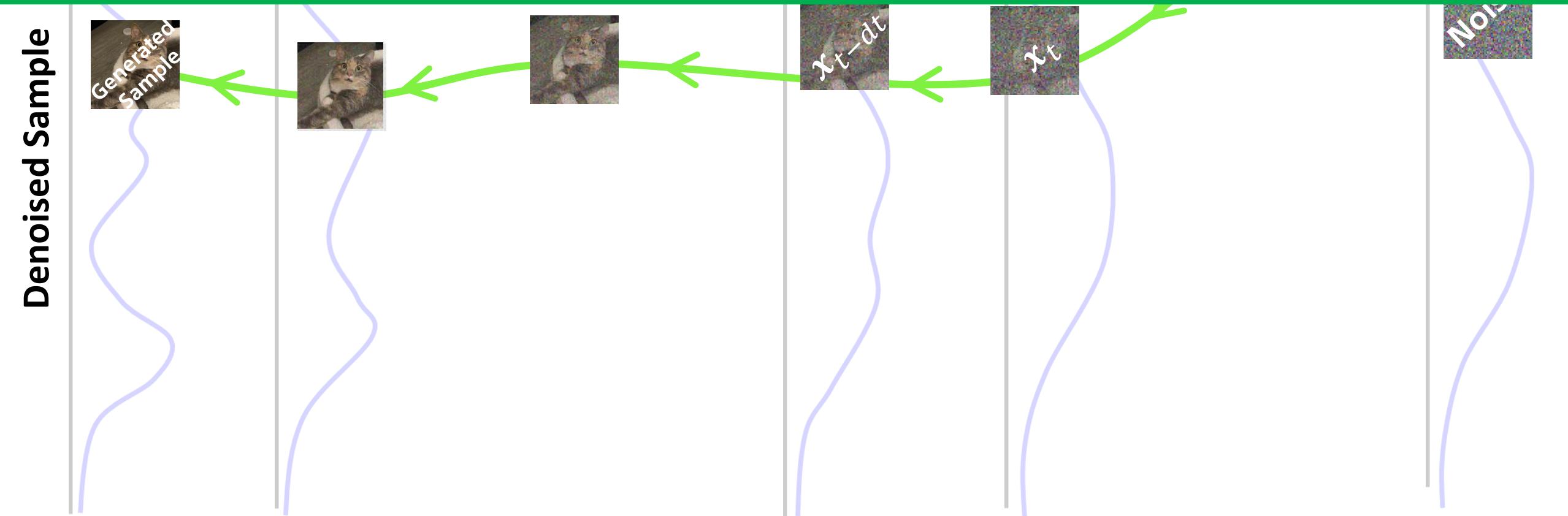
$$\frac{d\mathbf{x}_t}{dt} = -t \nabla \log p_t(\mathbf{x}_t), \mathbf{x}_T \sim p_{\text{prior}}$$



# Diffusion Model's Backward Deterministic Process (Generation)

For example,

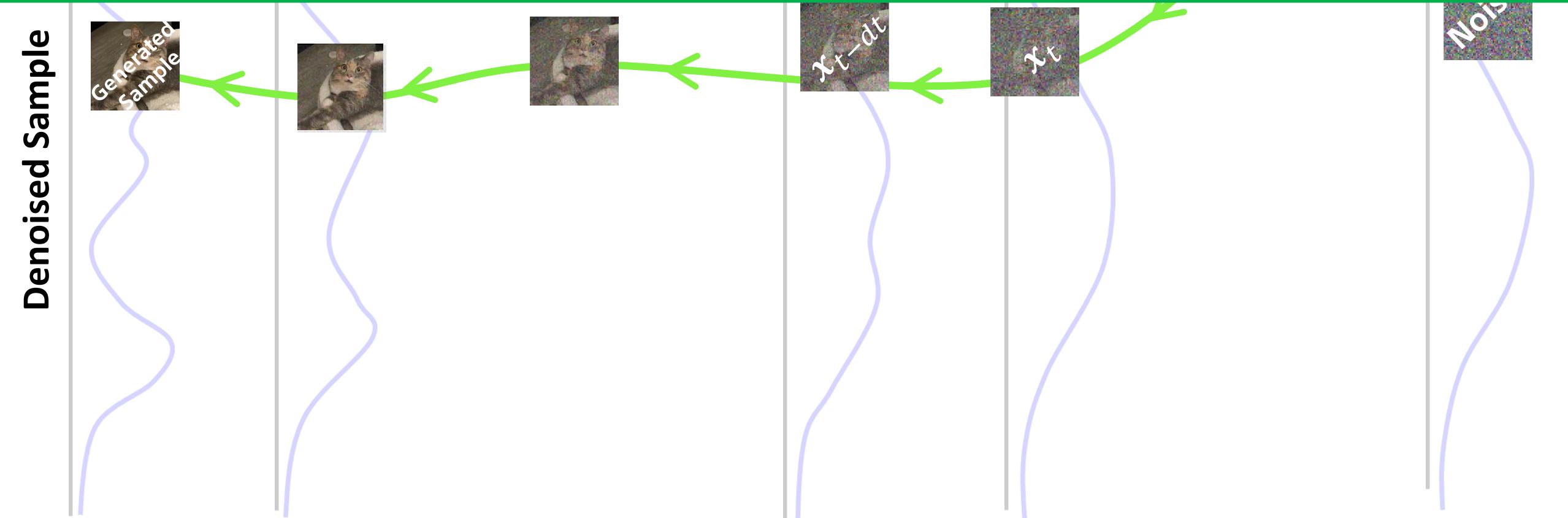
$$f(x_t, t) = 0 \text{ & } g(t) = \sqrt{2t}$$



# Diffusion Model's Backward Deterministic Process (Generation)

In our example,

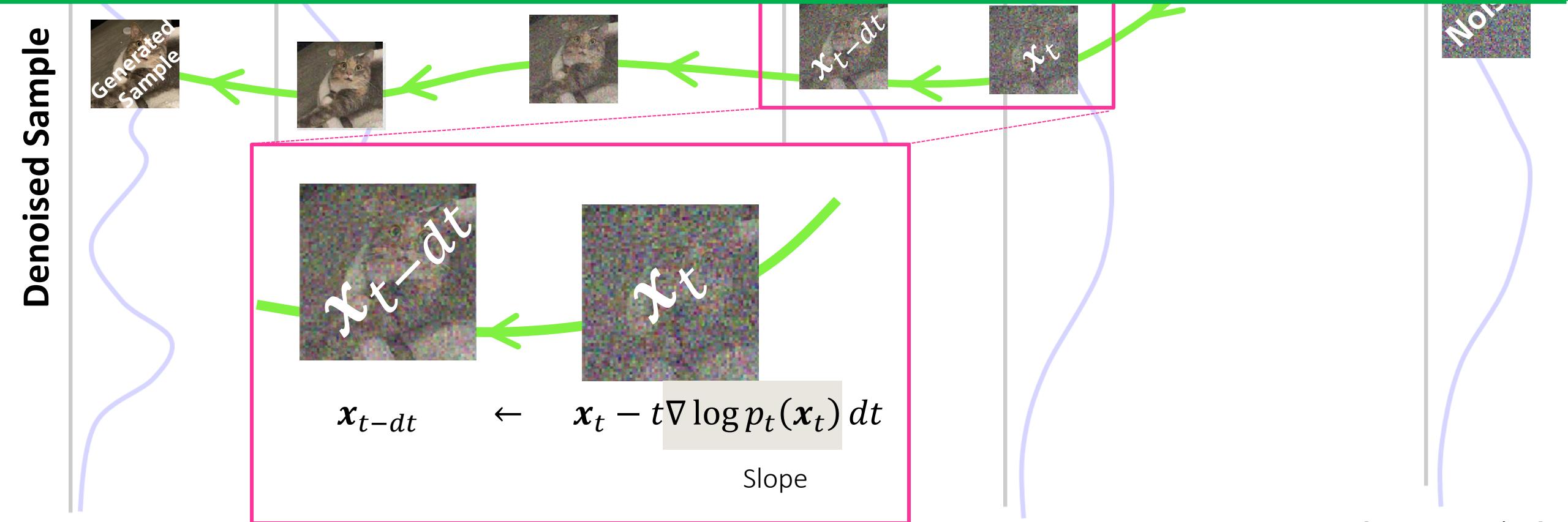
$$\frac{d\mathbf{x}_t}{dt} = -t \nabla \log p_t(\mathbf{x}_t), \mathbf{x}_T \sim p_{\text{prior}}$$



## Diffusion Model's Backward Deterministic Process (Generation)

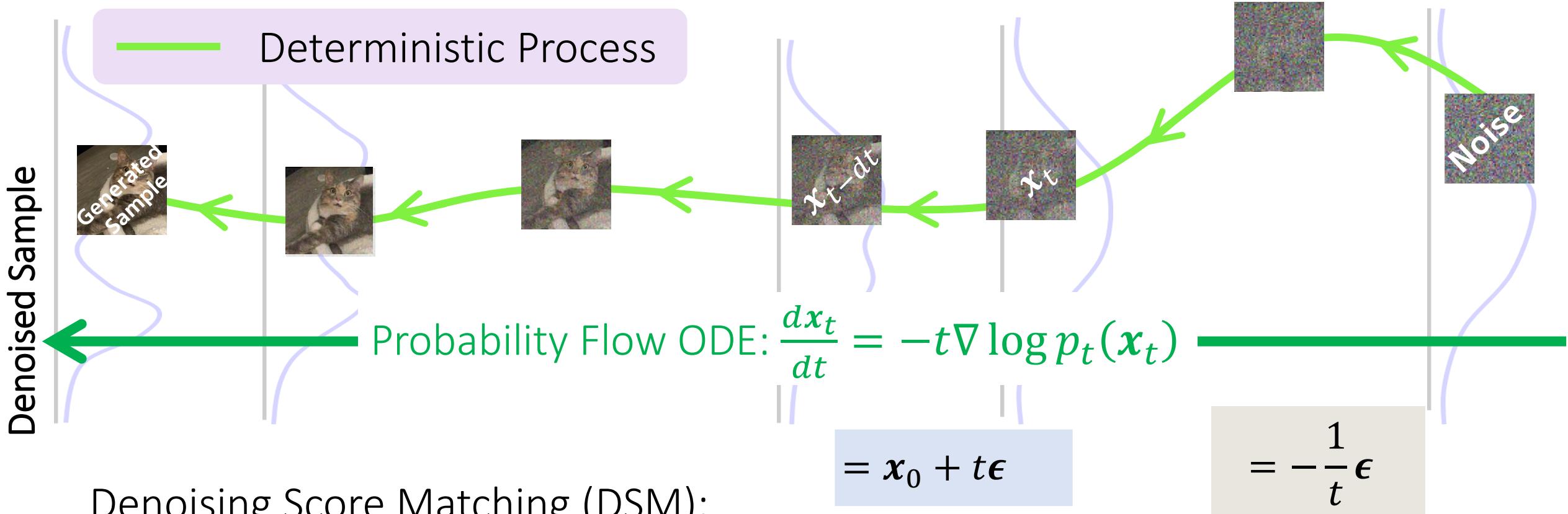
In our example,

$$\frac{dx_t}{dt} = -t \nabla \log p_t(x_t), \quad x_T \sim p_{\text{prior}}$$



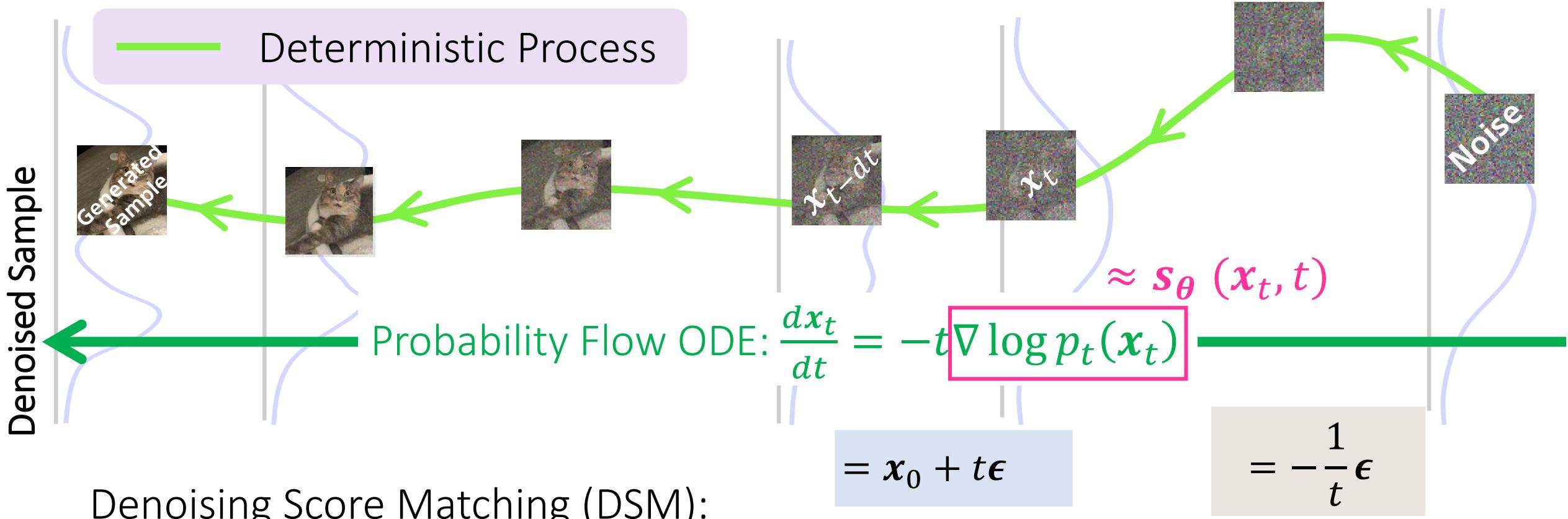
# Score-Based Diffusion Model's Training and Sampling

Users predefined noise schedulers:  $x_t = a_t x_0 + b_t \epsilon$  with  $a_t = 1$  &  $b_t = t$



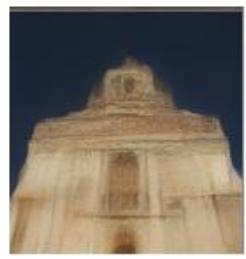
# Score-Based Diffusion Model's Training and Sampling

Users predefined noise schedulers:  $x_t = a_t x_0 + b_t \epsilon$  with  $a_t = 1$  &  $b_t = t$



$$\mathbb{E}_{x \sim p_{\text{data}}(x)} \mathbb{E}_{t \in U[0,1]} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,I)} \left[ \lambda_t \|s_\theta(x_t, t) - \nabla \log p_t(x_t | x_0)\|^2 \right]$$

# What is Happening in Generation?



High-frequency

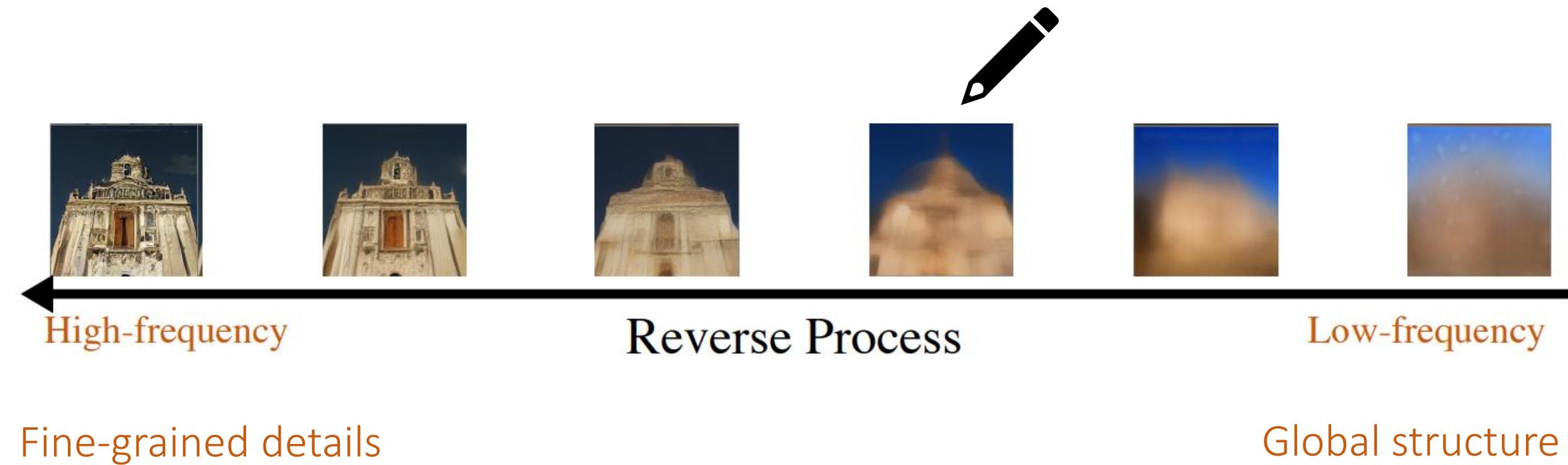
Reverse Process

Low-frequency

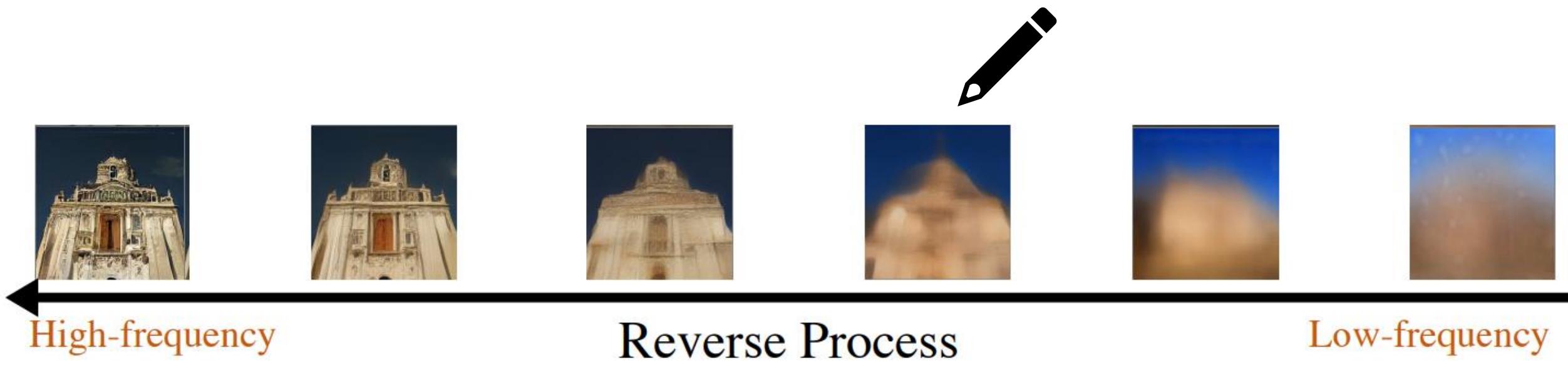
Fine-grained details

Global structure

# What is Happening in Generation?



# What is Happening in Generation?

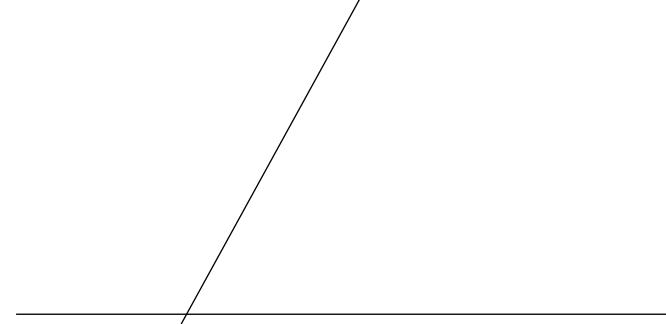


Fine-grained details

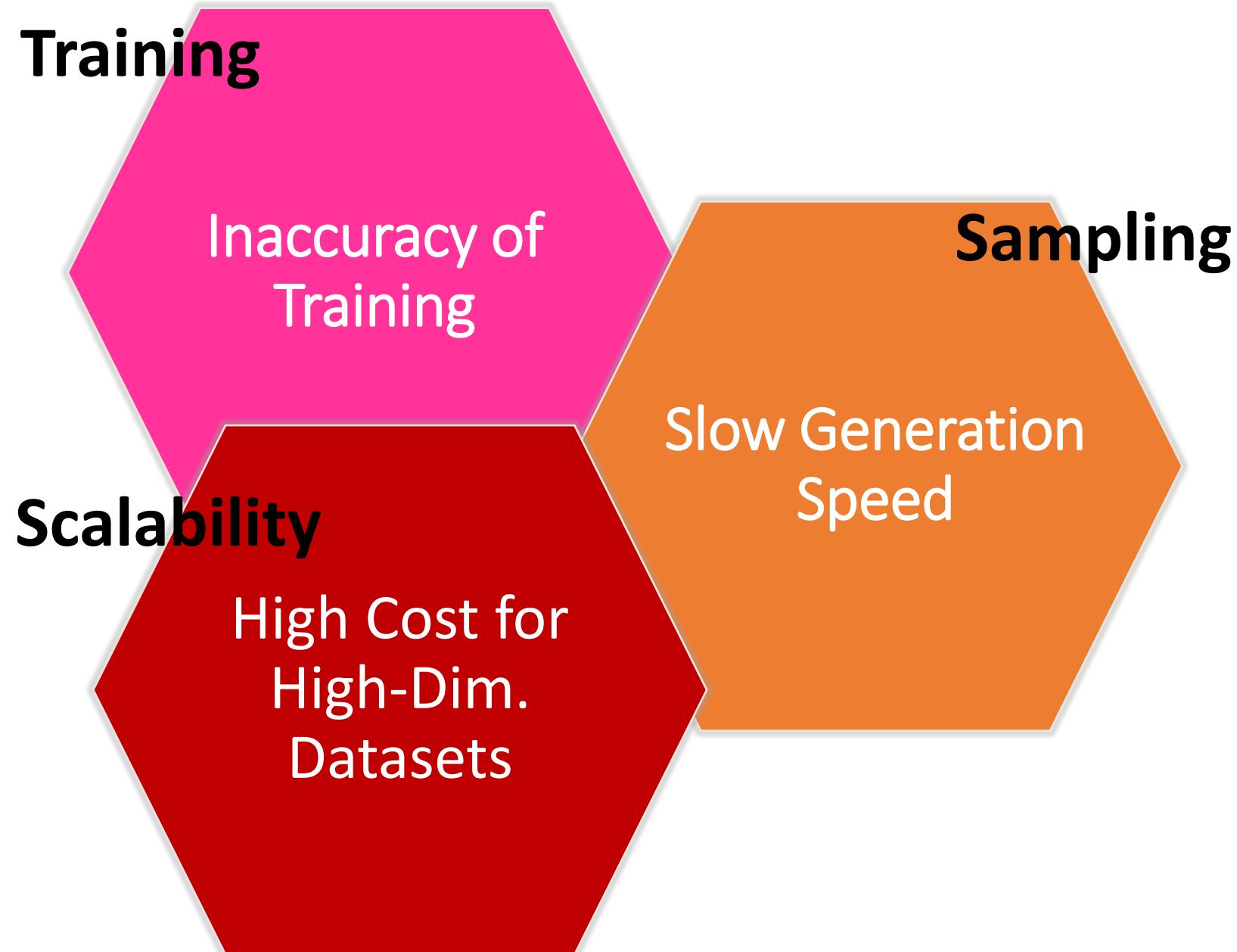
Global structure

Gradual denoising process, on the other hand,  
offers flexibility and control!

# Problems of Diffusion Models



# Problems of Diffusion Models



# Making Diffusion Models More Powerful

# Training

# [ICML'23 Lai+] FP-Diffusion



# Accurate Training

# Scalability

# [NeurIPS'24 Kim&Lai+] PaGoDA



# Cheap Training and Sampling in High-Dim.

# Sampling

# [ICLR'24 Kim&Lai+] CTM



# Making Diffusion Model More Powerful

# Training

# [ICML'23 Lai+]

# FP-Diffusion



# accuracy of Training

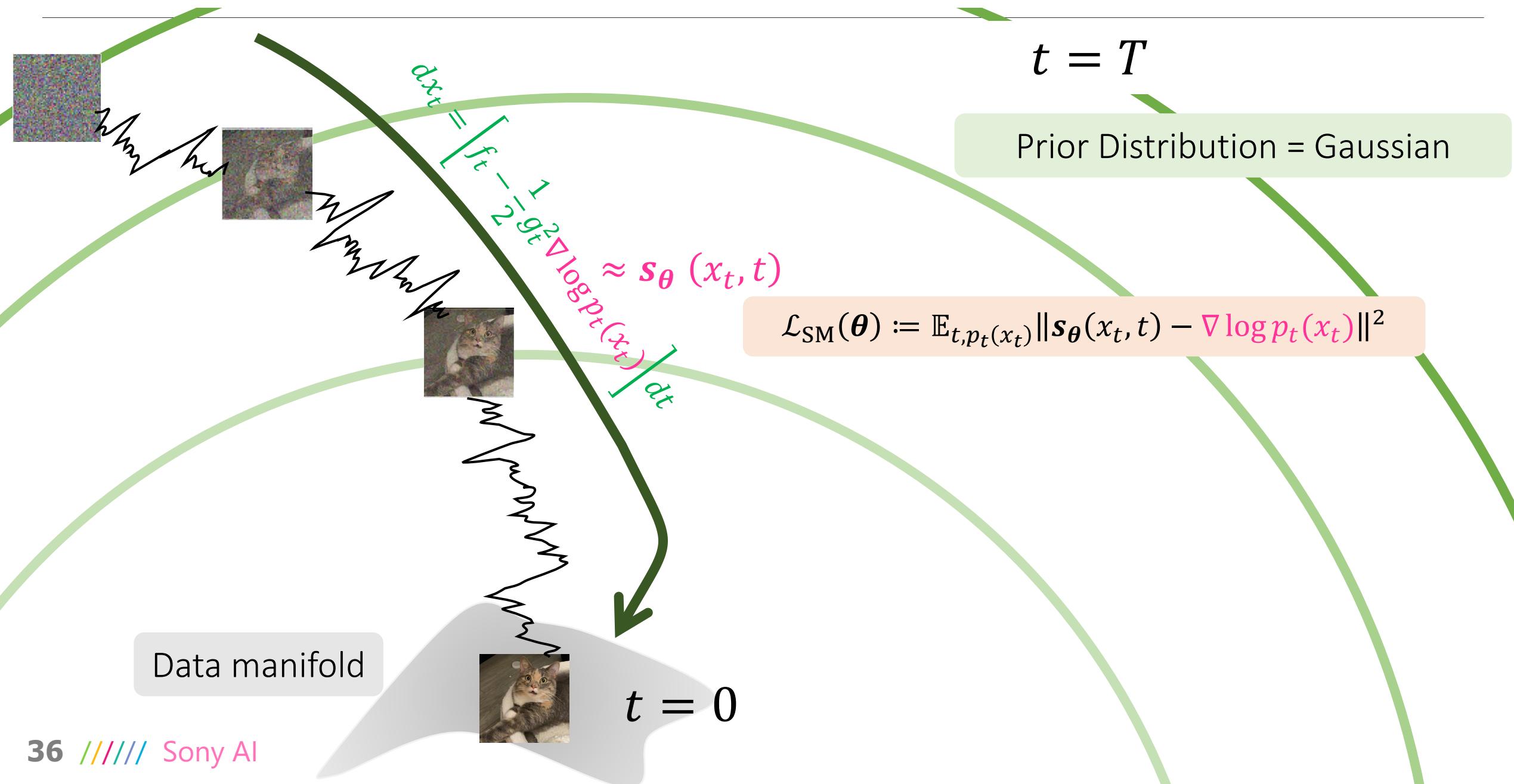
# Scalability

# High Cost for High-Dim. Datasets

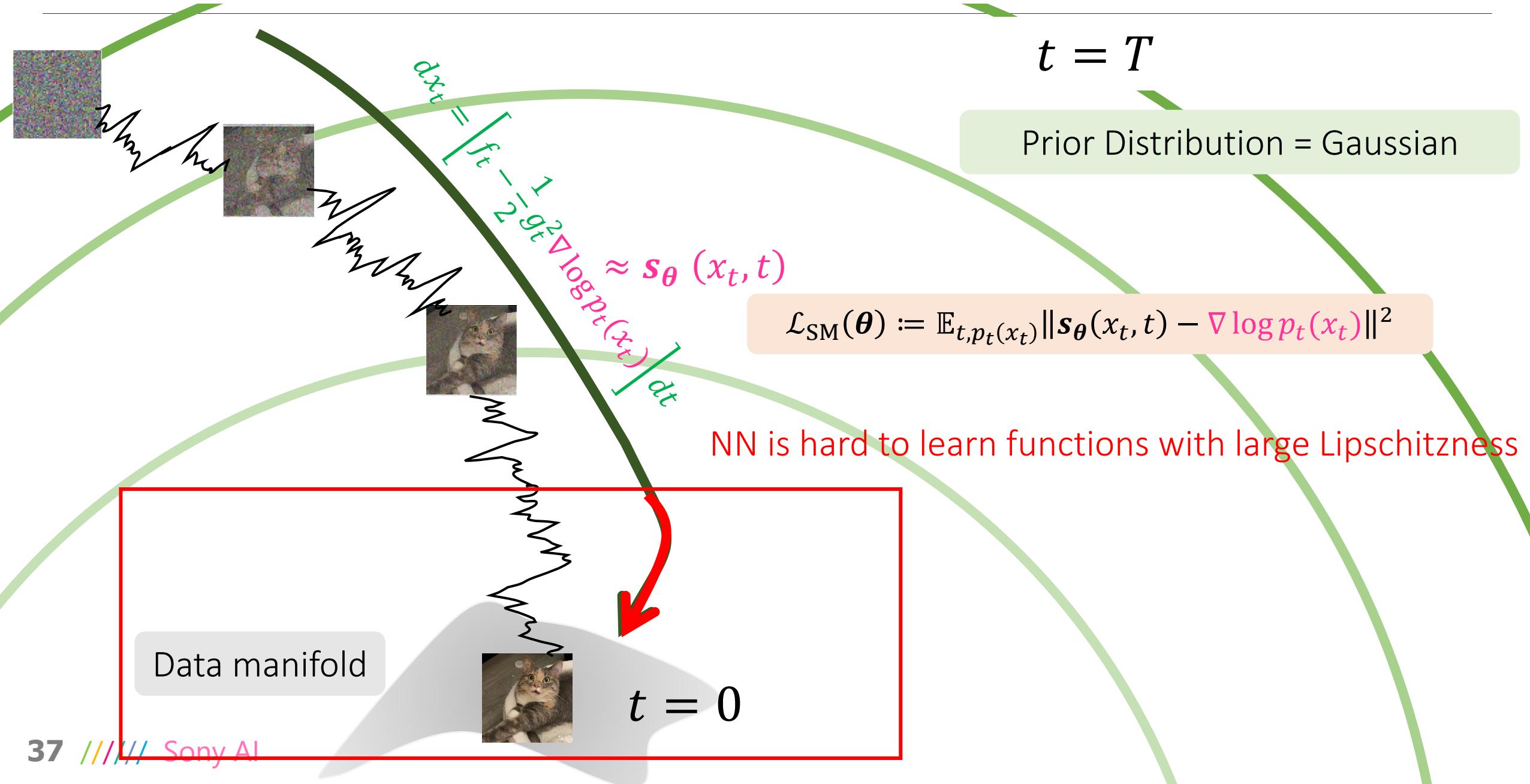
# Sampling

# Slow Generation Speed

# Motivation of FP-Diffusion

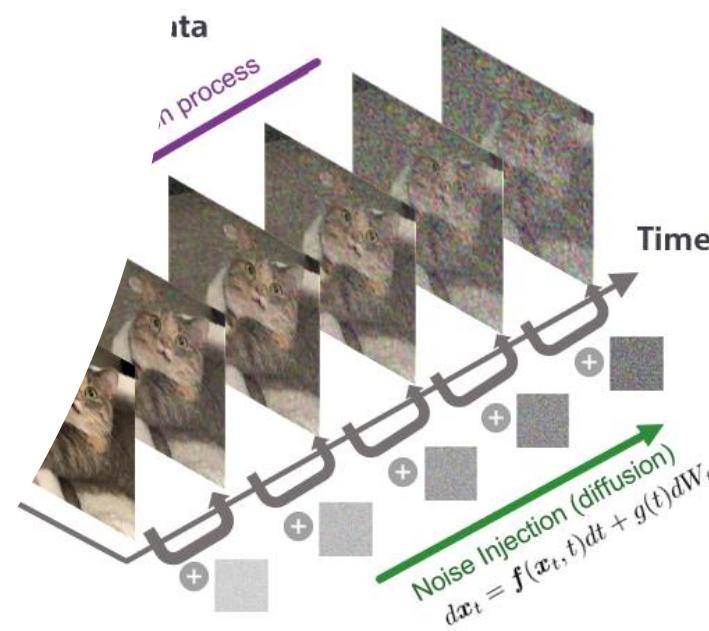


# Motivation of FP-Diffusion



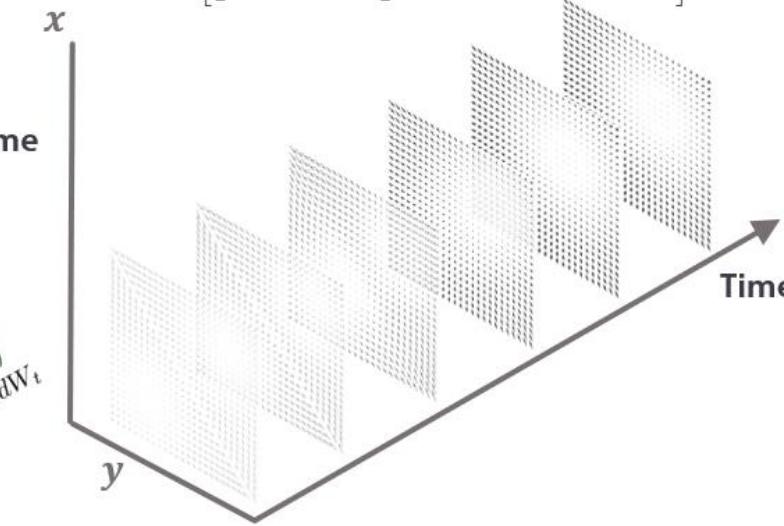
# FP-Diffusion

## FP-Diffusion

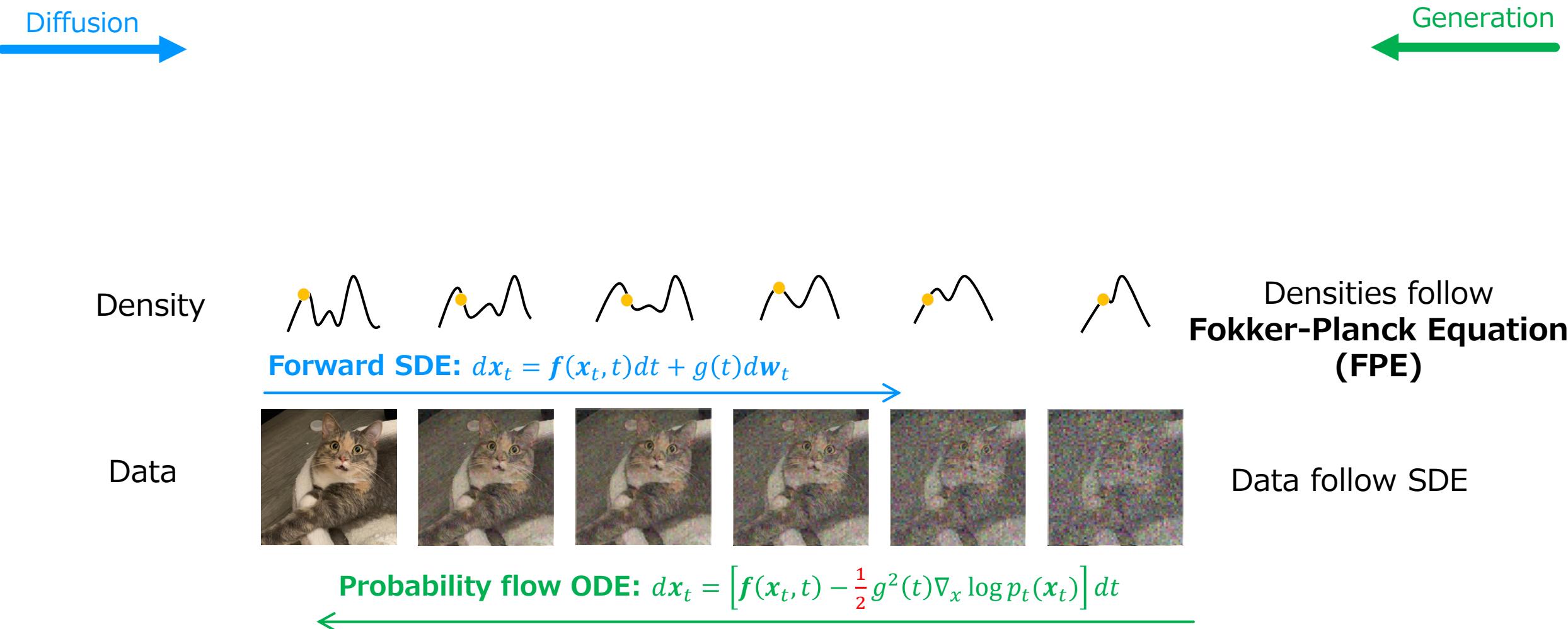


Evolution of scores follow  
**Score Fokker-Planck Equation (score FPE)**

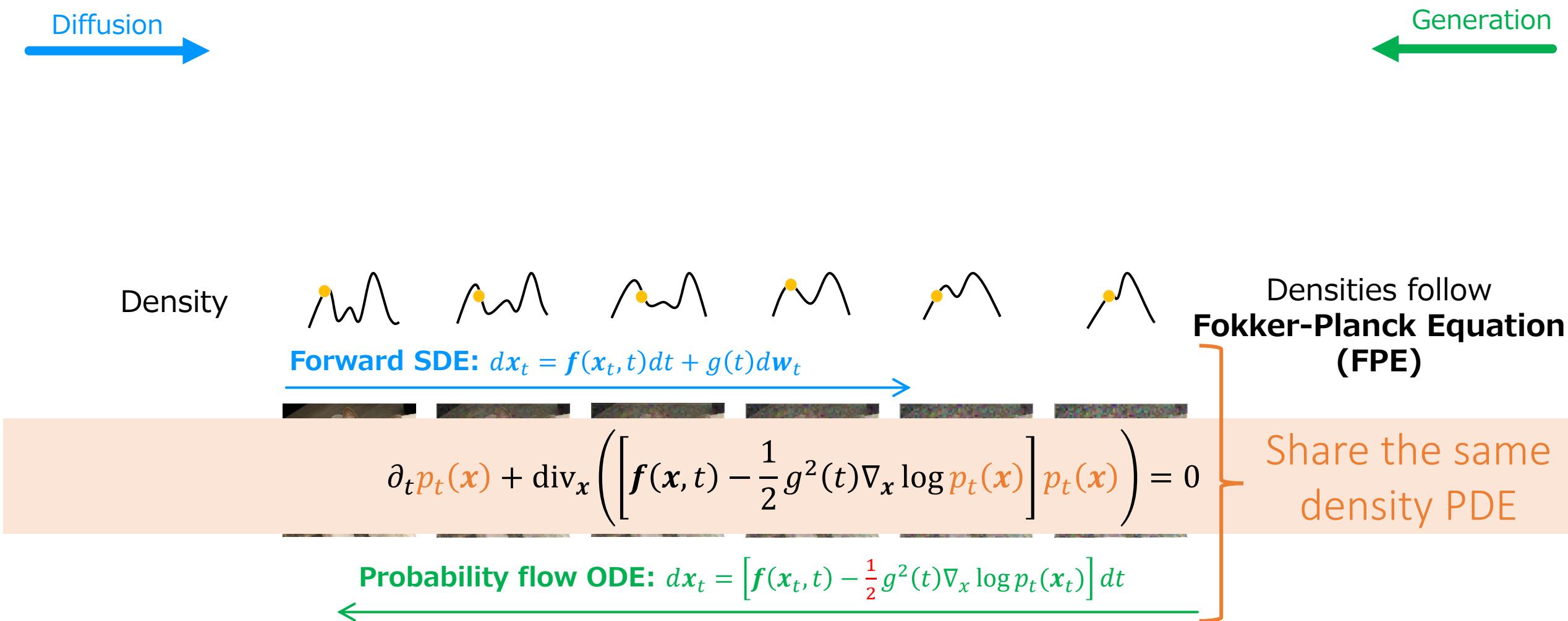
$$\partial_t \mathbf{s} - \nabla \left[ \frac{1}{2} g^2(t) \operatorname{div}(\mathbf{s}) + \frac{1}{2} g^2(t) \|\mathbf{s}\|^2 - \langle \mathbf{f}, \mathbf{s} \rangle - \operatorname{div}(\mathbf{f}) \right] = 0$$



# New perspective of score evolution: Score FPE



# New perspective of score evolution: Score FPE



# New perspective of score evolution: Score FPE

## Diffusion



# Generation



$$\text{Score: } s(x, t) := \nabla_x \log p_t(x)$$

**Forward SDE:**  $dx_t = f(x_t, t)dt + g(t)d\omega_t$

$$\partial_t p_t(x) + \operatorname{div}_x \left( \left[ f(x, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x) \right] p_t(x) \right) = 0$$

(FPE)

Share the same  
density PDE

**Probability flow ODE:**  $dx_t = \left[ f(x_t, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x_t) \right] dt$

# New perspective of score evolution: Score FPE

[ICML'23  
Lai+]

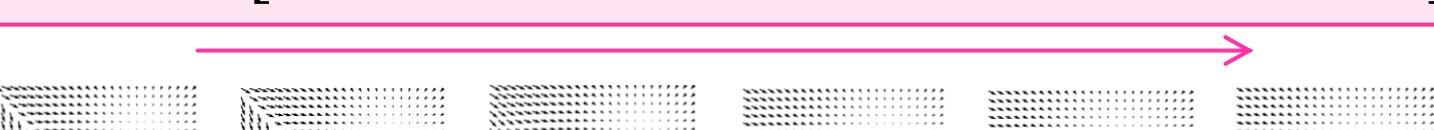
$$\partial_t \mathbf{s} - \nabla_{\mathbf{x}} \left[ \frac{1}{2} g^2(t) \operatorname{div}_{\mathbf{x}}(\mathbf{s}) + \frac{1}{2} g^2(t) \|\mathbf{s}\|_2^2 - \langle \mathbf{f}, \mathbf{s} \rangle - \operatorname{div}(\mathbf{f}) \right] = 0$$

Diffusion →

← Generation

Score

Scores follow  
Score FPE



$$\text{Score: } s(x, t) := \nabla_x \log p_t(x)$$

**Forward SDE:**  $dx_t = f(x_t, t)dt + g(t)d\omega_t$

$$\partial_t p_t(x) + \operatorname{div}_x \left( \left[ f(x, t) - \frac{1}{2} g^2(t) \nabla_x \log p_t(x) \right] p_t(x) \right) = 0$$

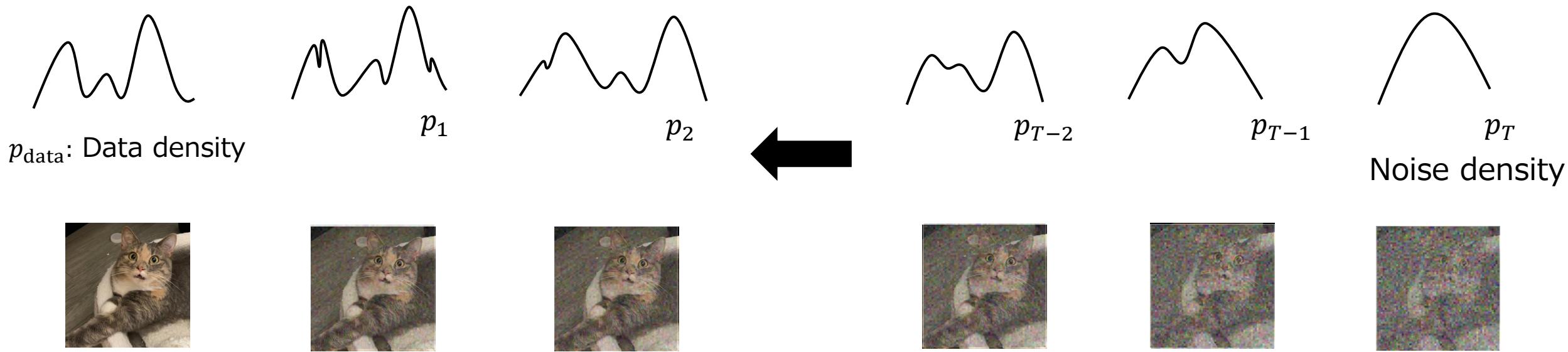
(FPE)

Share the same  
density PDE

**Probability flow ODE:**  $dx_t = \left[ f(x_t, t) - \frac{1}{2}g^2(t)\nabla_x \log p_t(x_t) \right] dt$

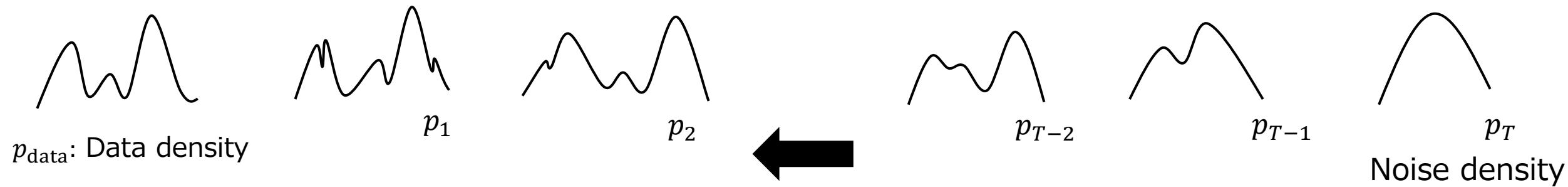
# Score FPE-error $\downarrow \Rightarrow$ density estimation quality $\uparrow$

---



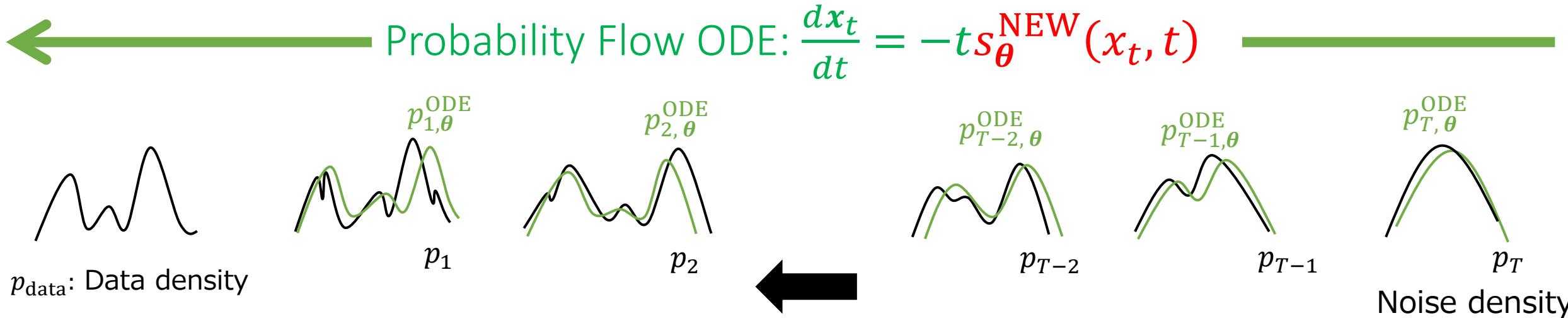
Score FPE-error  $\downarrow \Rightarrow$  density estimation quality  $\uparrow$

---



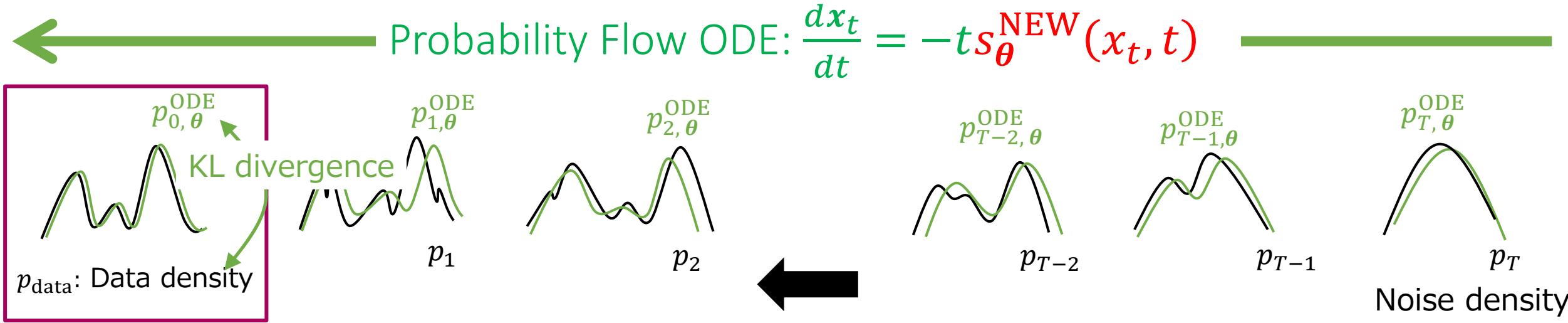
$$\mathcal{J}_{\text{ScoreFPE}}(\theta) := \left\| \partial_t \mathbf{s}_\theta^{\text{NEW}} - t \nabla_{\mathbf{x}} \left[ \text{div}_{\mathbf{x}}(\mathbf{s}_\theta^{\text{NEW}}) + \|\mathbf{s}_\theta^{\text{NEW}}\|_2^2 \right] \right\|_{\ell_p}^m$$

Score FPE-error  $\downarrow \Rightarrow$  density estimation quality  $\uparrow$



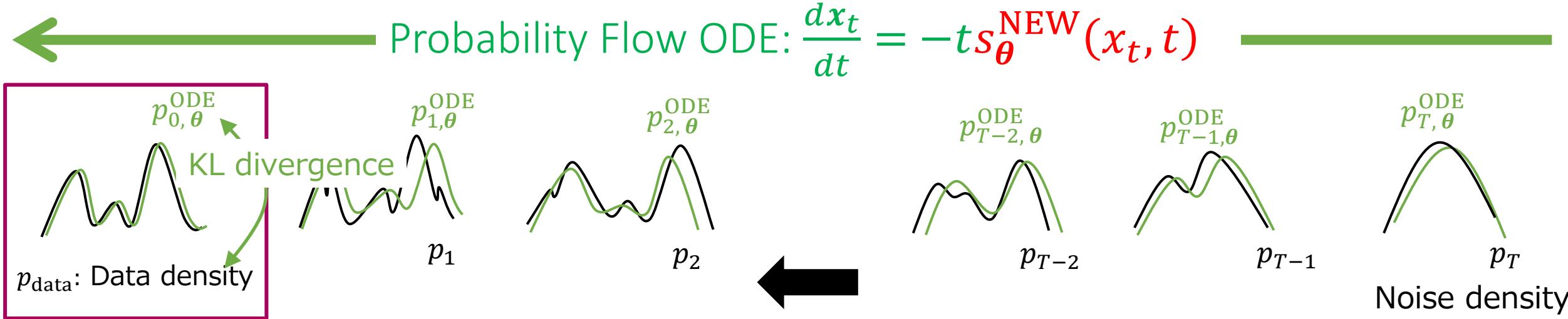
$$\mathcal{J}_{\text{ScoreFPE}}(\theta) := \left\| \partial_t s_{\theta}^{\text{NEW}} - t \nabla_x \left[ \text{div}_x(s_{\theta}^{\text{NEW}}) + \|s_{\theta}^{\text{NEW}}\|_2^2 \right] \right\|_{\ell_p}^m$$

Score FPE-error  $\downarrow \Rightarrow$  density estimation quality  $\uparrow$



$$\mathcal{J}_{\text{ScoreFPE}}(\theta) := \left\| \partial_t s_{\theta}^{\text{NEW}} - t \nabla_x \left[ \text{div}_x(s_{\theta}^{\text{NEW}}) + \|s_{\theta}^{\text{NEW}}\|_2^2 \right] \right\|_{\ell_p}^m$$

Score FPE-error  $\downarrow \Rightarrow$  density estimation quality  $\uparrow$



$$\mathcal{J}_{\text{ScoreFPE}}(\theta) := \left\| \partial_t s_{\theta}^{\text{NEW}} - t \nabla_x \left[ \text{div}_x(s_{\theta}^{\text{NEW}}) + \|s_{\theta}^{\text{NEW}}\|_2^2 \right] \right\|_{\ell_p}^m$$

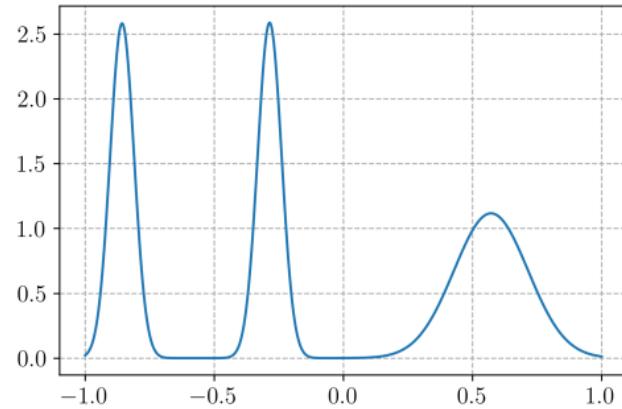
IICML'23 Lai+

### Theorems.

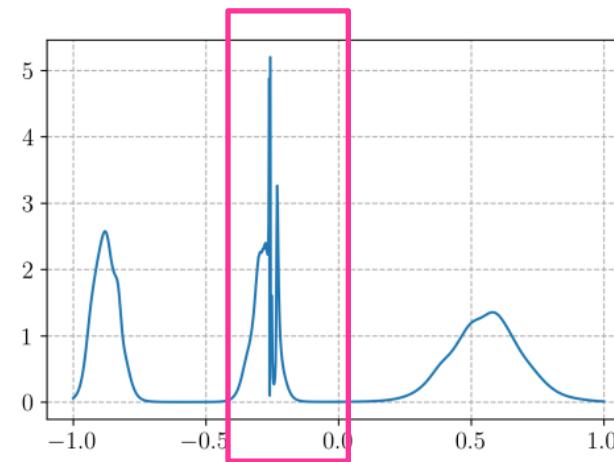
$$D_{\text{KL}}(p_{\text{data}} \parallel p_{0,\theta}^{\text{ODE}})$$

$$\lesssim D_{\text{KL}}(p_T \parallel p_{T,\theta}^{\text{ODE}}) + \mathcal{J}_{\text{DSM}}(\theta) + \mathcal{J}_{\text{DSM}}^{\frac{1}{2}}(\theta) \left( \mathcal{J}_{\text{ScoreFPE}}(\theta) + \mathcal{J}_{\text{ScoreFPE}}^{\frac{1}{2}}(\theta) + C \right)^{1/2}$$

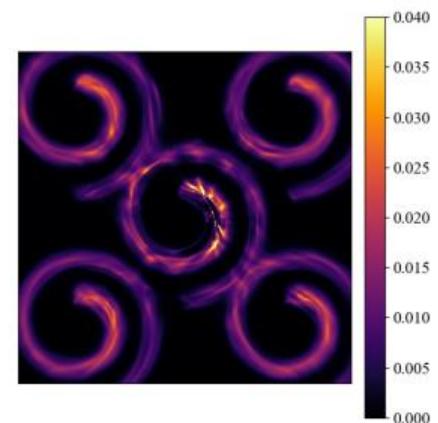
# Results on toy datasets



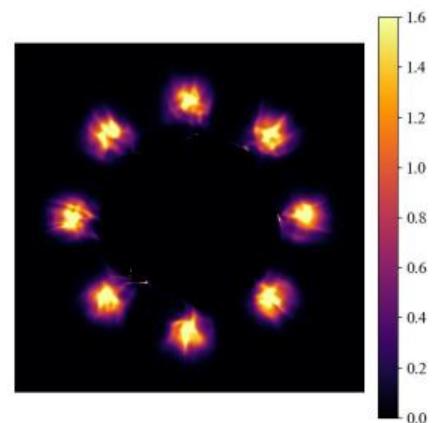
### (a) Data density



### (b) Vanilla DSM ( $\alpha = \beta = 0.0$ )

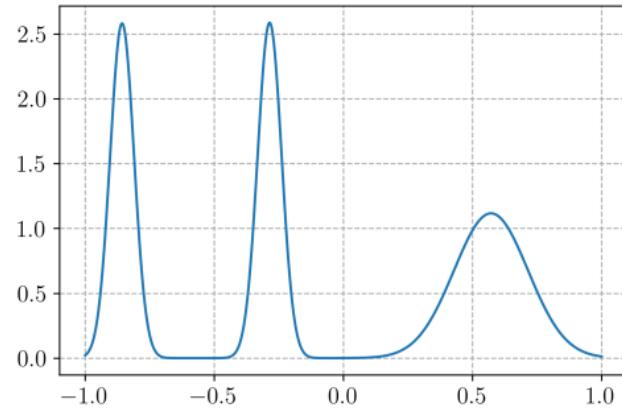


(a) Vanilla DSM  
 $(\alpha = \beta = 0.0)$

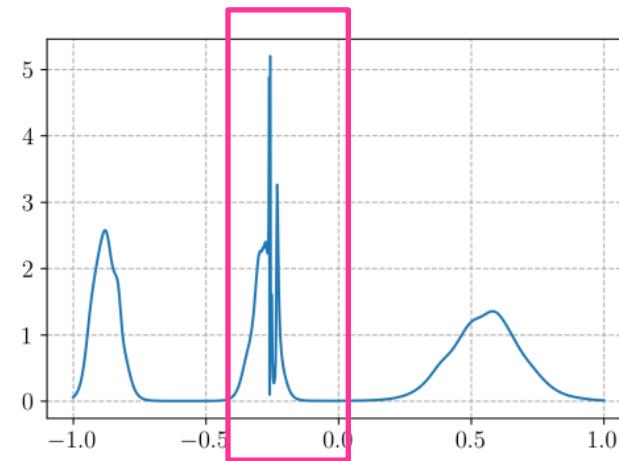


(e) Vanilla DSM  
 $(\alpha = \beta = 0.0)$

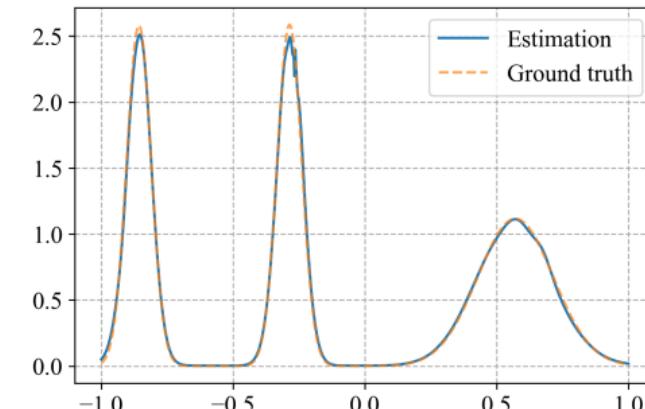
# Results on toy datasets



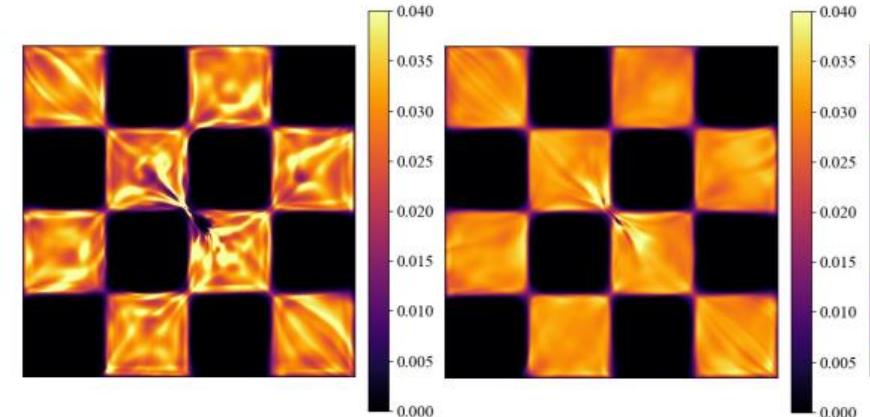
### (a) Data density



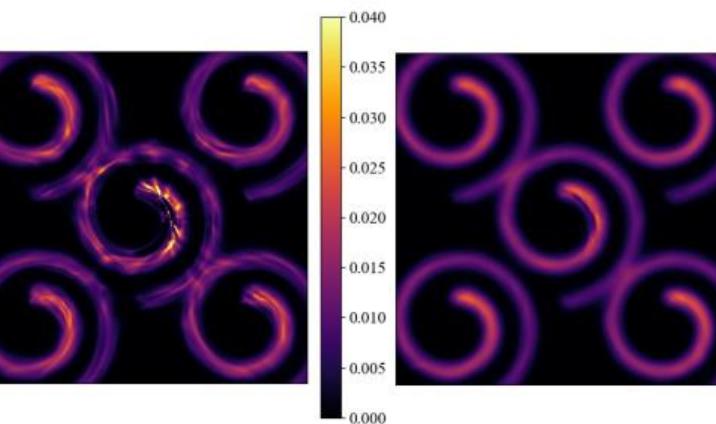
(b) Vanilla DSM ( $\alpha = \beta = 0.0$ )



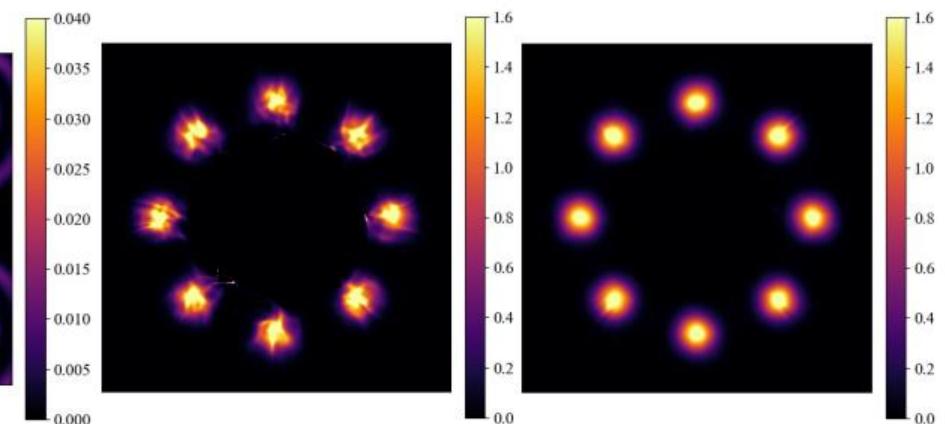
(c) FP-Diffusion (with  $\alpha = 0.0015$ )



(a) Vanilla DSM  
 $(\alpha = \beta = 0.0)$



(c) Vanilla DSM  
 $(\alpha = \beta = 0.0)$



(e) Vanilla DSM  
 $(\alpha = \beta = 0.0)$

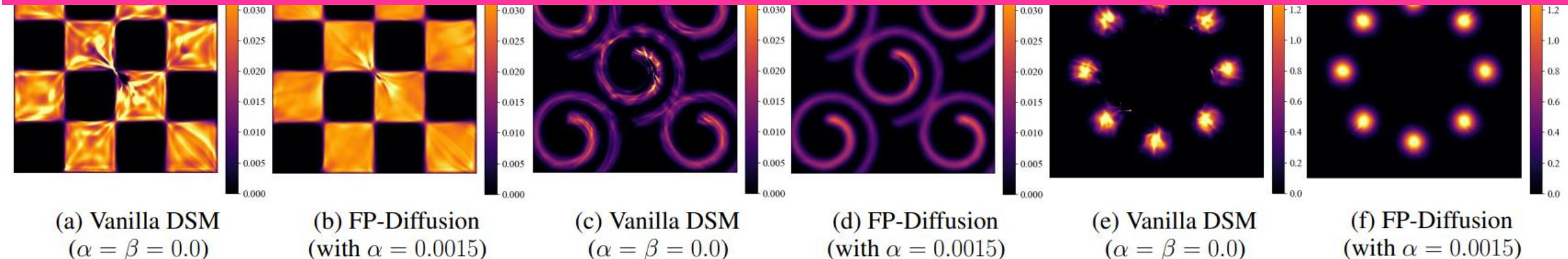
# Results on toy datasets



Imposing Score FPE + Vanilla DSM → Improve density learning



# Score FPE bridges PDE to practical Diffusion Model



# More Than This...

- Relating to conservative field (or energy-based method)
  - Using PINNs for solving scores...?
  - ...

Let's catch up offline if you're interested in this!

# Making Diffusion Model More Powerful

# Training

# Inaccuracy of Training

# Scalability

# High Cost for High-Dim. Datasets

# Sampling

# [ICLR'24 Kim&Lai+] CTM

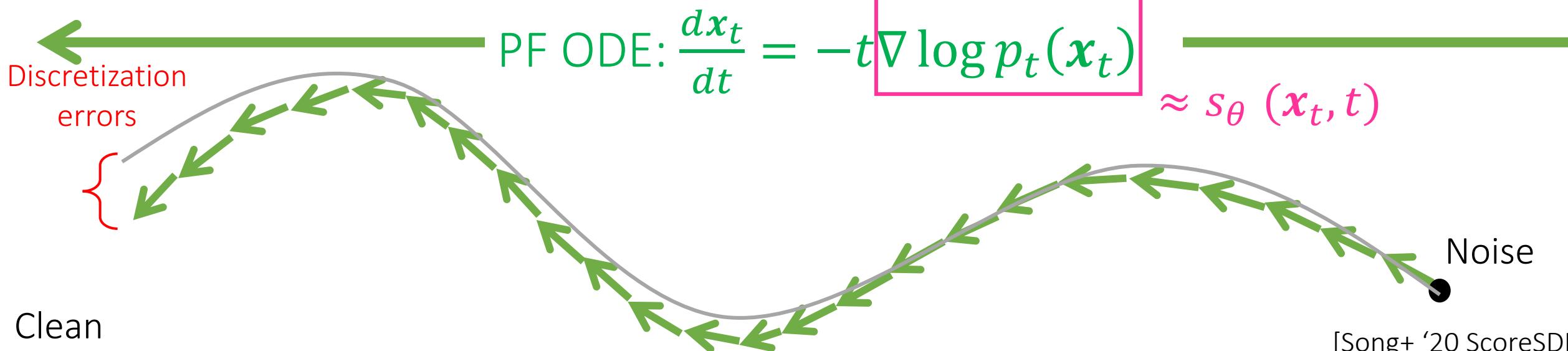


# Motivation of CTM

Generated new image



Gradually transforming noise to a new image



# Motivation of CTM

Generated new image

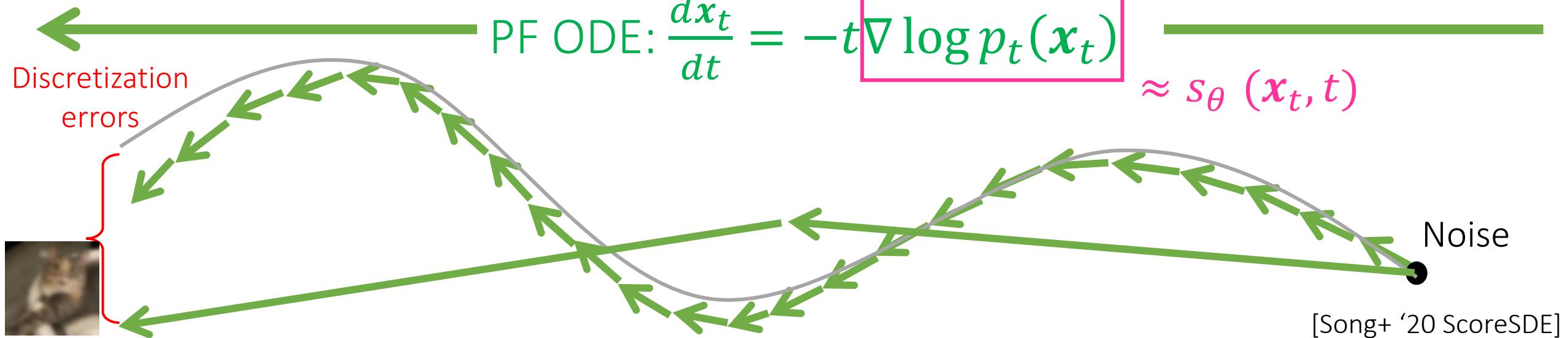


Gradually transforming noise to a new image



$$\text{PF ODE: } \frac{dx_t}{dt} = -t \nabla \log p_t(x_t)$$

$$\approx s_\theta(x_t, t)$$

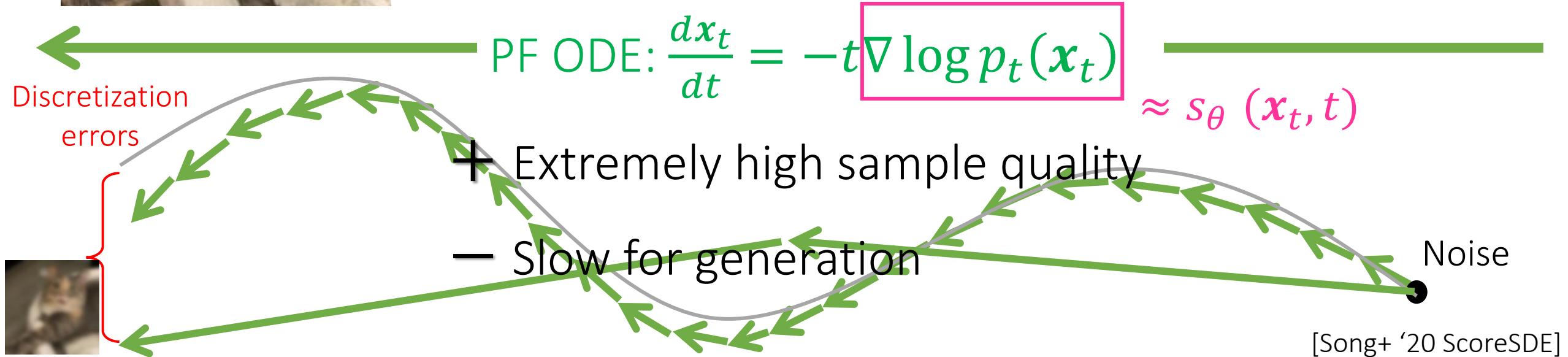


# Motivation of CTM

Generated new image



Gradually transforming noise to a new image



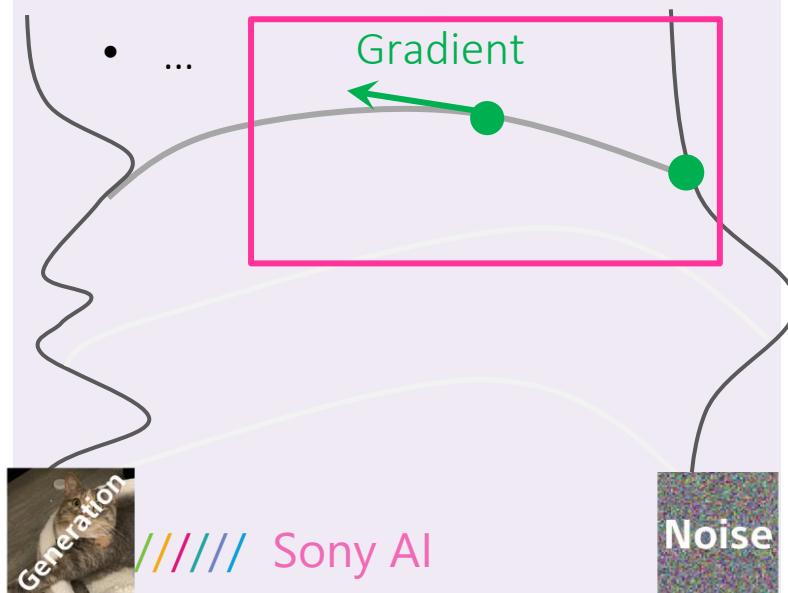
# Fast Sampling in Diffusion Model

PF-ODE (2021)

# Training-free

# Sophisticated solvers (2021-2022+)

- [ICLR'21] DDIM
  - [NeurIPS'22] DPM; DPM++
  - [ICLR'23] DEIS
  - ...



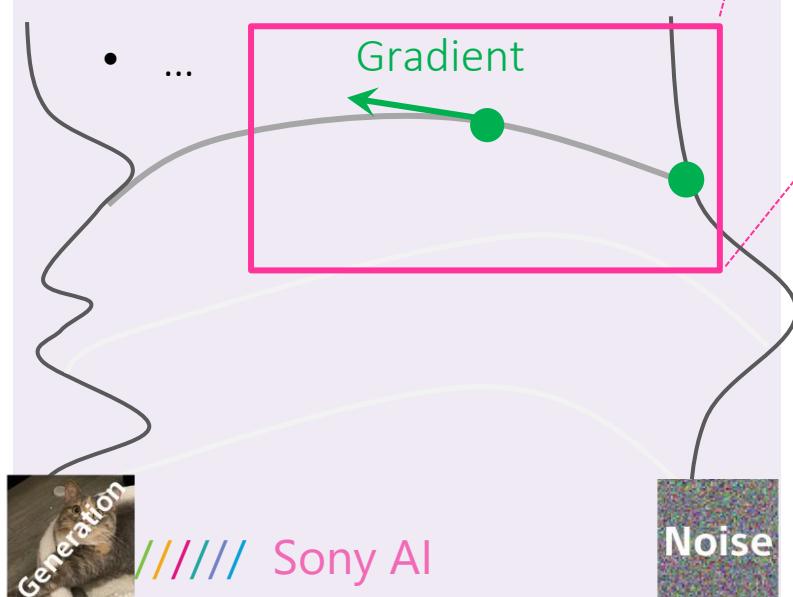
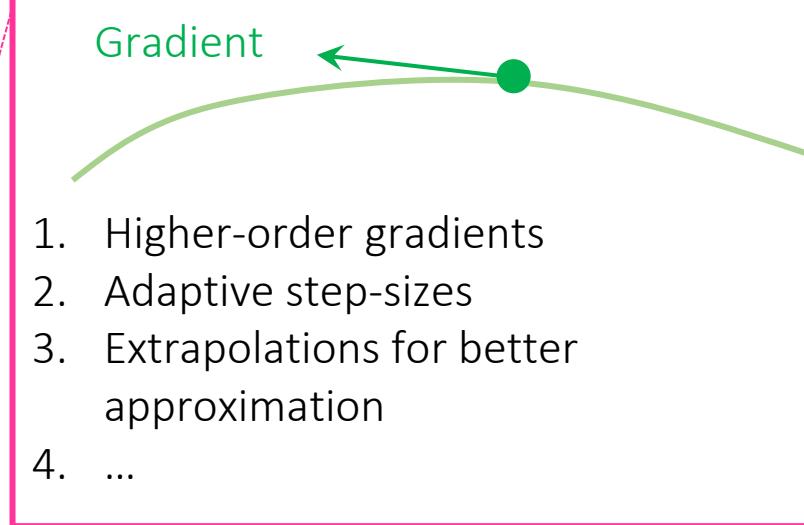
# Fast Sampling in Diffusion Model

PF-ODE (2021)

Training-free

Sophisticated solvers  
(2021-2022+)

- [ICLR'21] DDIM
- [NeurIPS'22] DPM; DPM++
- [ICLR'23] DEIS
- ...



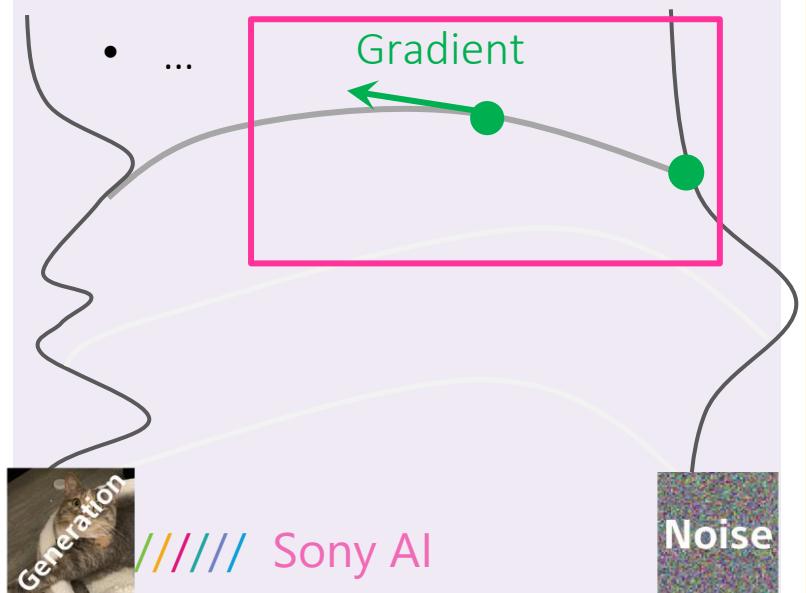
# Fast Sampling in Diffusion Model

PF-ODE (2021)

## Training-free

Sophisticated solvers  
(2021-2022+)

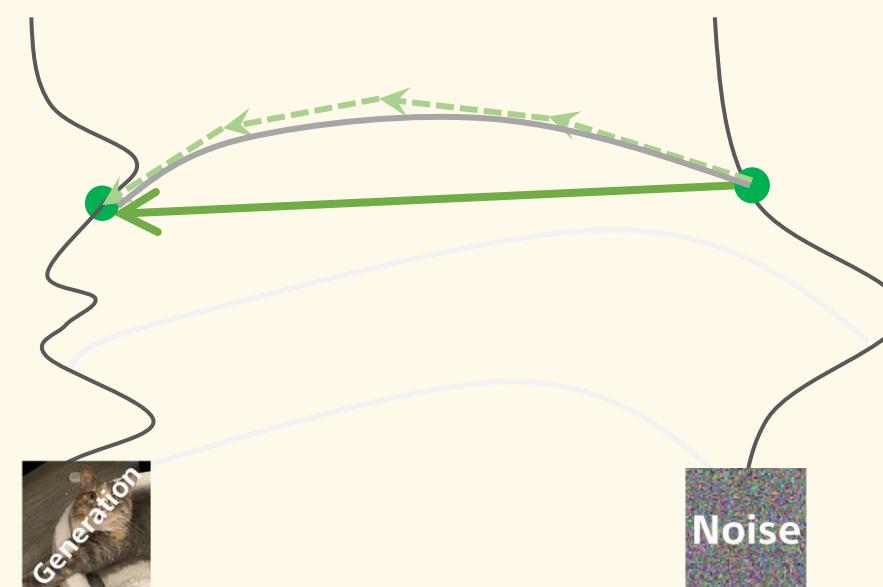
- [ICLR'21] DDIM
- [NeurIPS'22] DPM; DPM++
- [ICLR'23] DEIS
- ...



## Training-based

Looking for more optimal  
trajectories for sampling (2023+)

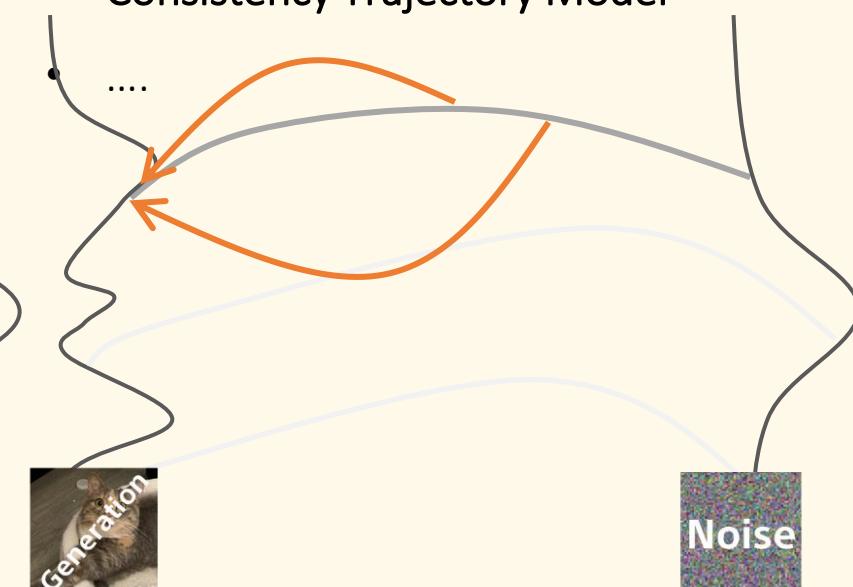
- [ICLR'23] Rectified Flow
- [ICLR'23] Flow Matching



Distillation from DM (2023-)

- [ICML'23 Song+] Consistency Model
- [ICLR'24 Kim&Lai+]

## Consistency Trajectory Model



# CONSISTENCY TRAJECTORY MODEL



istency Trajectory Models: Learning  
ability Flow ODE Trajectory of Diffusion

Dongjun Kim<sup>\*,▲1</sup>, Chieh-Hsin Lai<sup>\*,1</sup>,  
Ji-Hsiang Liao<sup>1</sup>, Naoki Murata<sup>1</sup>, Yuhta Takida<sup>1</sup>, Toshimitsu Uesaka<sup>1</sup>, Yutong He<sup>▲1,3</sup>,  
Yuki Mitsufuji<sup>1,2</sup>, Stefano Ermon<sup>4</sup>,  
<sup>1</sup>Sony AI, <sup>2</sup>Sony Group Corporation, <sup>3</sup>Carnegie Mellon University, <sup>4</sup>Stanford University  
ICLR 2024

\*Equal Contribution ([✉ Dongjun Kim](#); [✉ Chieh-Hsin Lai](#))

▲Internship at Sony AI

# Diffusion → Distillation

PF-ODE (2021)

$$dx_u = -u \nabla_x \log p_u^{\text{oracle}}(x_u)$$

# Diffusion Model

- Training:

$$\nabla_x \log p_u^{\text{oracle}}(x_u) \approx s_\phi(x_u, u)$$

- Sampling:

$$\begin{aligned} x_0^{\text{oracle}} &= x_T - \int_T^0 u \nabla_x \log p_u^{\text{oracle}}(x_u) du \\ &\approx x_T - \int_T^0 u s_\phi(x_u, u) du \end{aligned}$$

# Distillation-Based Model

- **Training:**

$$F^{\text{oracle}}(\boldsymbol{x}_t, t) = \boldsymbol{x}_t - \int_t^0 u \nabla_{\boldsymbol{x}} \log p_u^{\text{oracle}}(\boldsymbol{x}_u) \, du$$

# Diffusion → Distillation

PF-ODE (2021)

$$dx_u = -u \nabla_x \log p_u^{\text{oracle}}(x_u)$$

# Diffusion Model

- **Training:**

$$\nabla_x \log p_u^{\text{oracle}}(x_u) \approx s_\phi(x_u, u)$$

- Sampling:

$$x_0^{\text{oracle}} = x_T - \int^0 u \nabla_x \log p_u^{\text{oracle}}(x_u) du$$

$$\approx x_T - \int_T^0 us_\phi(x_u, u) du$$

# Distillation-Based Model

- **Training:**

$$F^{\text{oracle}}(x_t, t) \approx x_t - \int_t^0 u s_{\phi}^{\text{PreTrained}}(x_u, u) du$$

$$\approx f_{\theta}(x_t, t)$$

- Sampling: few steps!

$$f_{\theta}(x_T, T)$$

This direction is drastically blooming since 2023

PF-UDE (2021)

$$dx_u = -u \nabla_x \log p_u^{\text{oracle}}(x_u)$$

# Diffusion Model

- Training:

$$\nabla_x \log p_u^{\text{oracle}}(x_u) \approx s_\phi(x_u, u)$$

- Sampling:

$$x_0^{\text{oracle}} = x_T - \int_0^T u \nabla_x \log p_u^{\text{oracle}}(x_u) du$$

$$\approx x_T - \int_T^0 u s_\phi(x_u, u) du$$

# Distillation-Based Model

- Training:

$$F^{\text{oracle}}(\mathbf{x}_t, t) \approx \boxed{\mathbf{x}_t - \int_t^0 u \ s_{\phi}^{\text{PreTrained}}(\mathbf{x}_u, u) \ du}$$

$$\approx f_{\theta}(x_t, t)$$

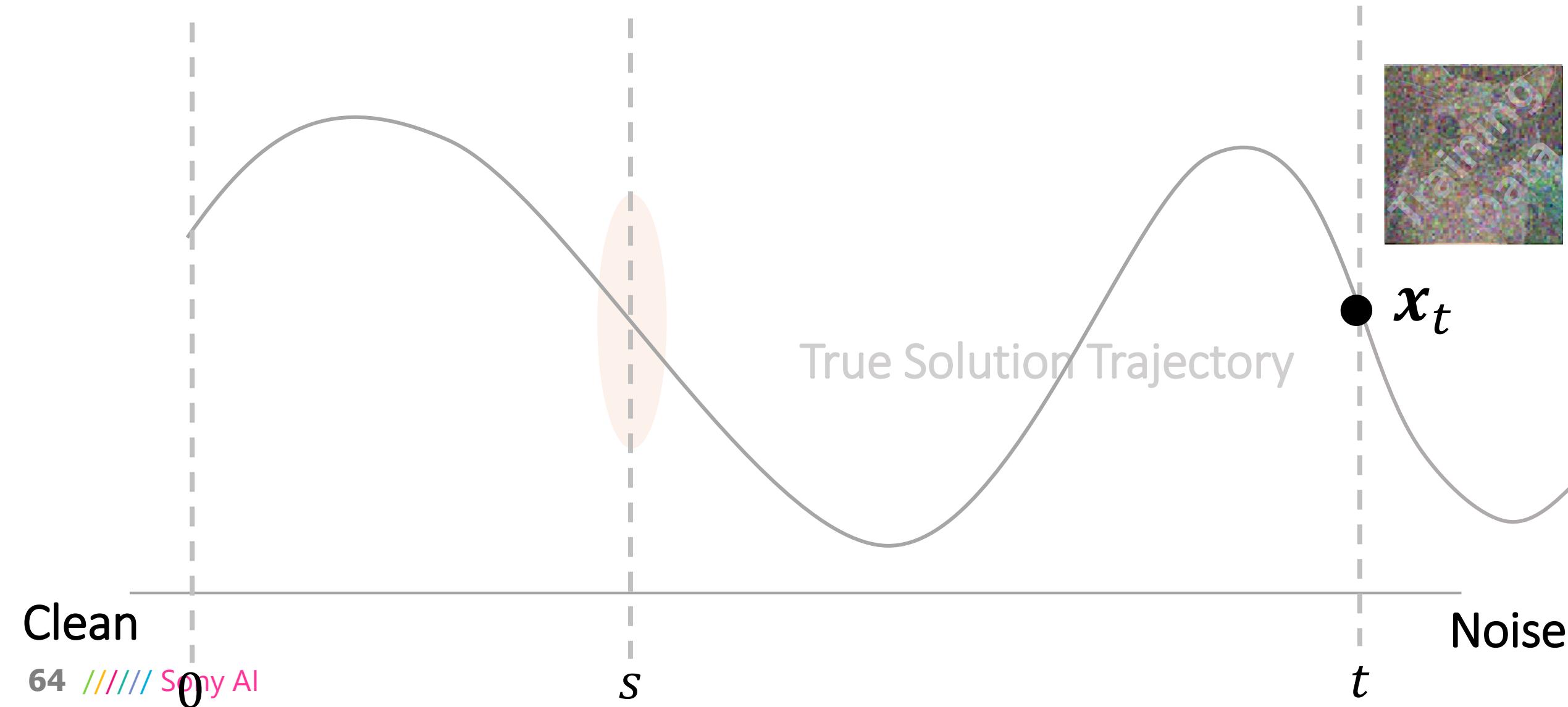
- Sampling: few steps!

$$f_{\theta}(x_T, T)$$

# CTM's Training and Mechanism

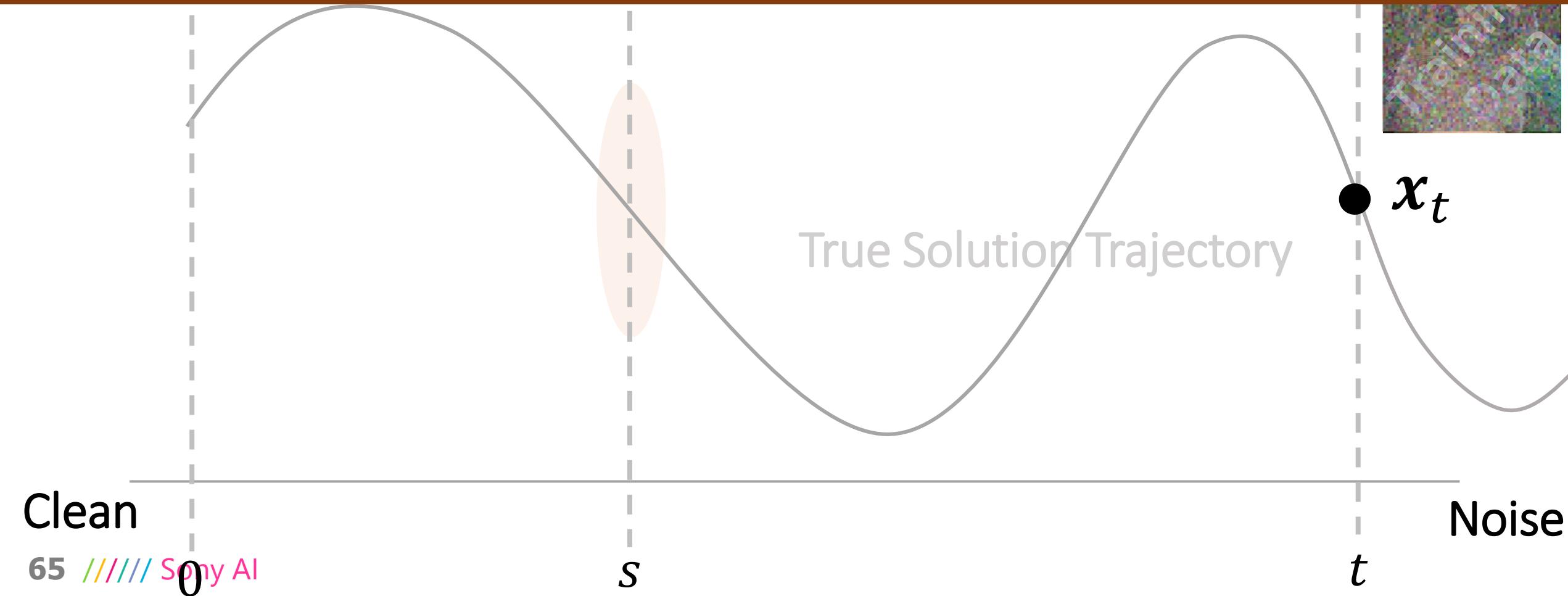
---

# Consistency Trajectory Model (CTM)'s Training



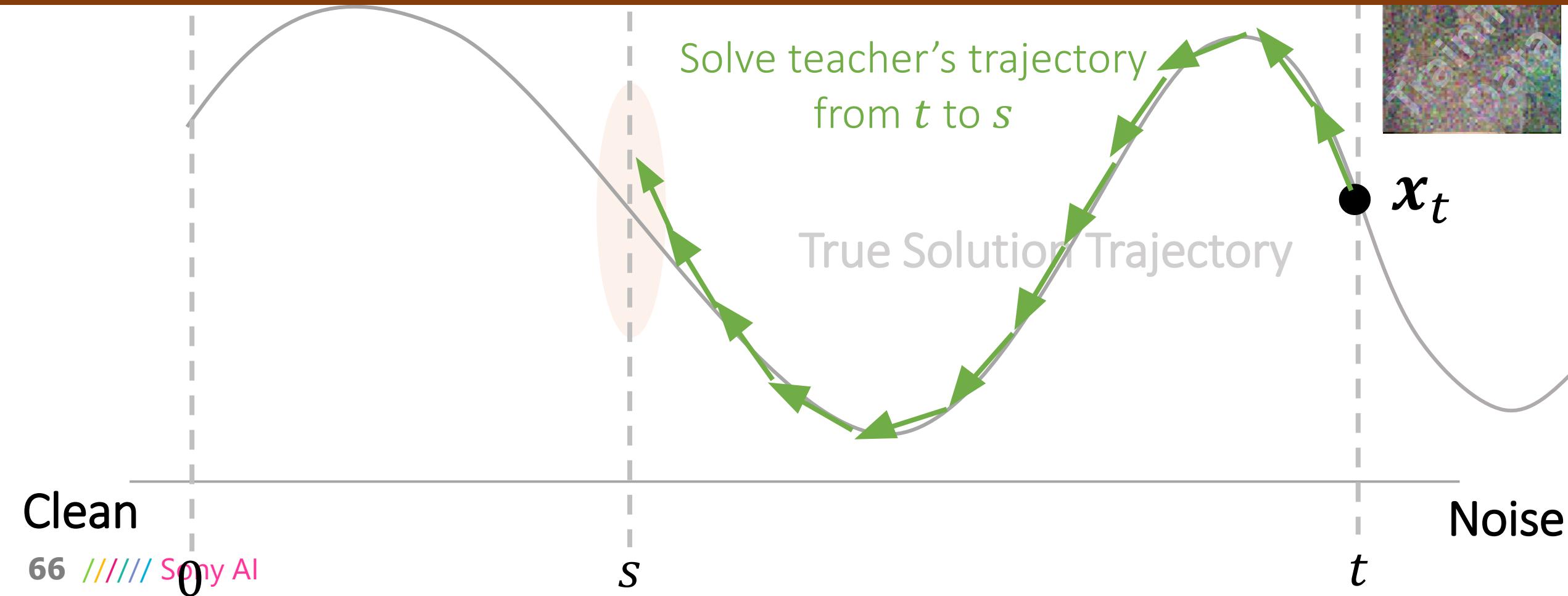
# CTM learns any-to-any jump on the trajectory

$$x_t - \int_t^S \tau \nabla \log p_\tau(x_\tau) d\tau$$



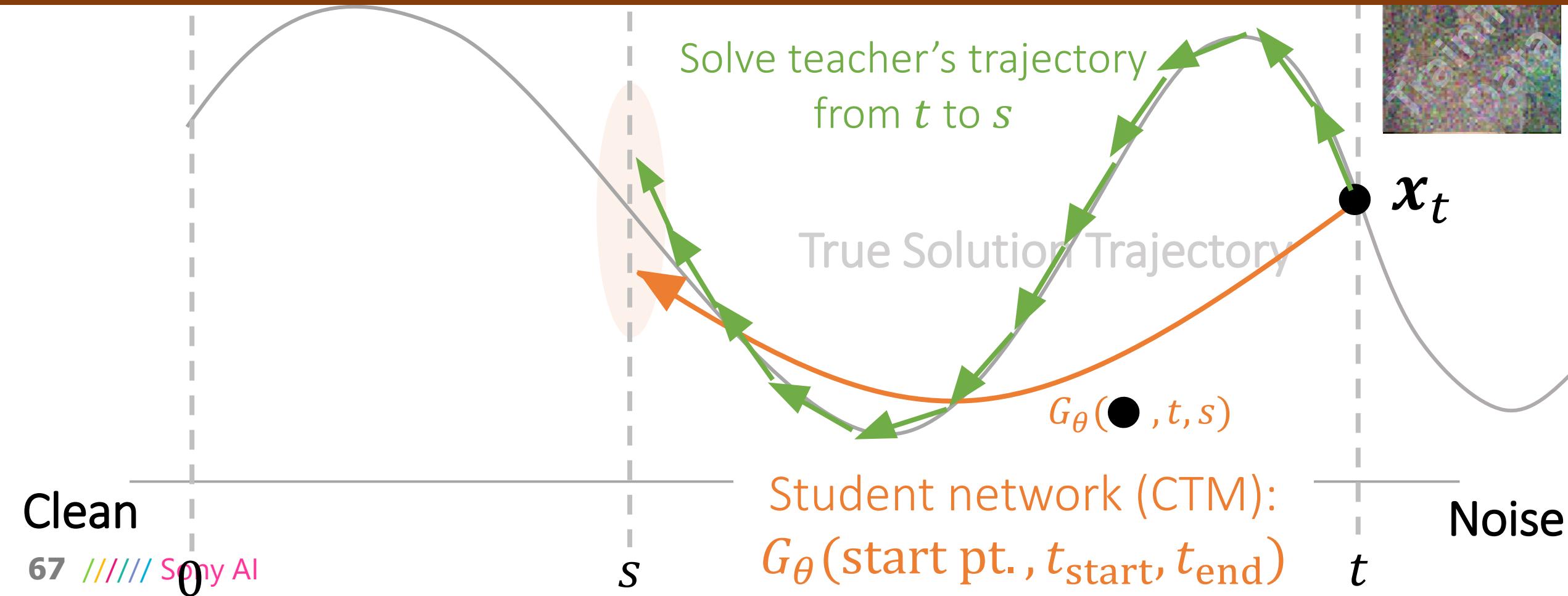
# CTM learns any-to-any jump on the trajectory

$$\boldsymbol{x}_t - \int_t^s \tau \mathbf{s}_{\text{teacher}}(\boldsymbol{x}_\tau, \tau) d\tau$$



# CTM learns any-to-any jump on the trajectory

$$G_\theta(x_t, t, s) \approx x_t - \int_t^s \tau s_{\text{teacher}}(x_\tau, \tau) d\tau$$



# Consistency Trajectory Model (CTM)'s Training

Student network (CTM):  $G_\theta(\text{start pt.}, t_{\text{start}}, t_{\text{end}})$

## CTM's Soft Matching:

## Teacher diffusion model with multi-step ODE solver



# True Solution Trajectory

$G_\theta(\bullet, t, s)$

$G_{\text{sg}(\theta)}(\bullet, u,$

1

1

1

## Noise

2

5

# Clean

# Consistency Trajectory Model (CTM)'s Training

Student network (CTM):  $G_\theta(\text{start pt.}, t_{\text{start}}, t_{\text{end}})$

# CTM's Soft Matching: Teacher diffusion model with multi-step ODE solve



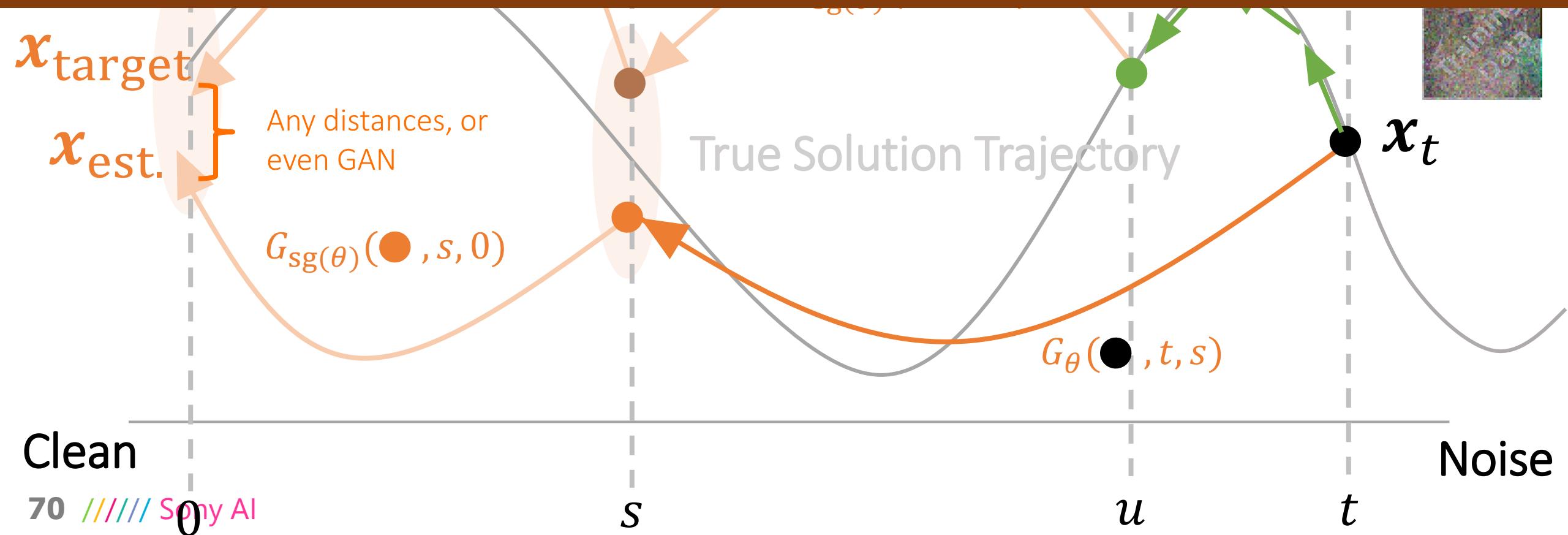
# Consistency Trajectory Model (CTM)'s Training

Student network (CTM):  $G_\theta(\text{start pt.}, t_{\text{start}}, t_{\text{end}})$

CTM's Soft Matching:

True Solution Trajectory

CTM learns any-to-any jump on the trajectory



# Consistency Trajectory Model (CTM)'s Training

Student network (CTM):  $G_\theta(\text{start pt.}, t_{\text{start}}, t_{\text{end}})$

# CTM's Soft Matching

## Technical Difficulties

# CTM learns any-to-any jump on the trajectory

$x_{\text{target}}$

$x_{\text{est.}}$



Any distances, or even GAN

# Any distances, or even GAN

# True Solution Trajectory

1

$x_t$

$$\mathcal{L}_{\text{CTM}} + \mathcal{L}_{\text{DSM}}$$

$G_\theta(\bullet, t, s)$

# Clean

# Consistency Trajectory Model (CTM)'s Training

Student network (CTM):  $G_\theta(\text{start pt.}, t_{\text{start}}, t_{\text{end}})$

CTM's Soft Matching:

True Solution Trajectory

CTM learns any-to-any jump on the trajectory



$$d(x_{\text{target}}, x_{\text{est.}}) + d(G_\theta(x_t, t, s \approx t), x_{\text{data}})$$

Clean

72 // Sony AI

$s$

$u$

$t$

Noise

# CTM's flexibility in sampling

---

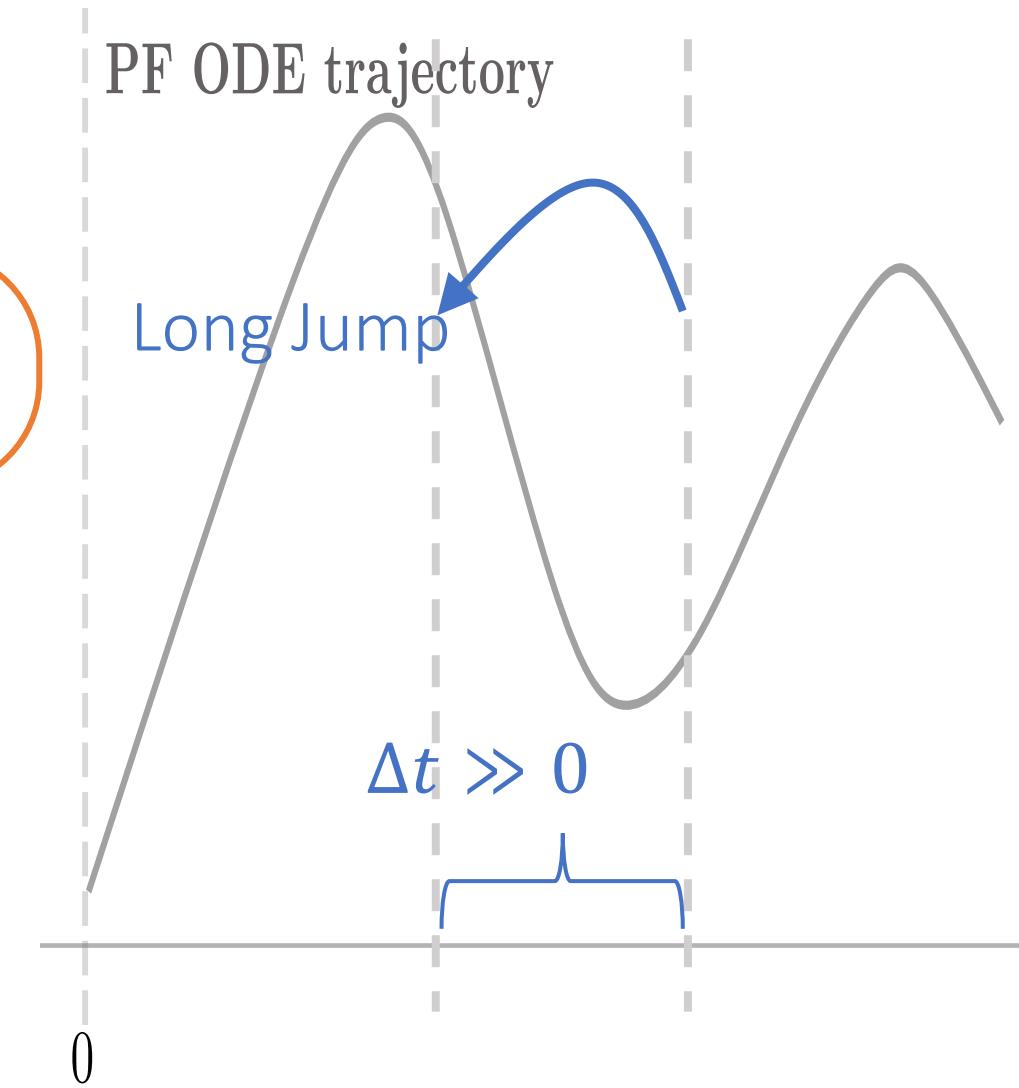
# CTM's Sampling Flexibility

CTM enables both long “jumps” along the solution trajectory and score evaluation!

$$G_\theta(x_t, t, t - \Delta t)$$

# Long jump ( $\Delta t \gg 0$ )

CTM enables long jump along the trajectory  
→ Allow our new sampling method



# CTM's Sampling Flexibility

CTM enables both long “jumps” along the solution trajectory and score evaluation!

$$G_{\theta}(x_t, t, t - \Delta t)$$

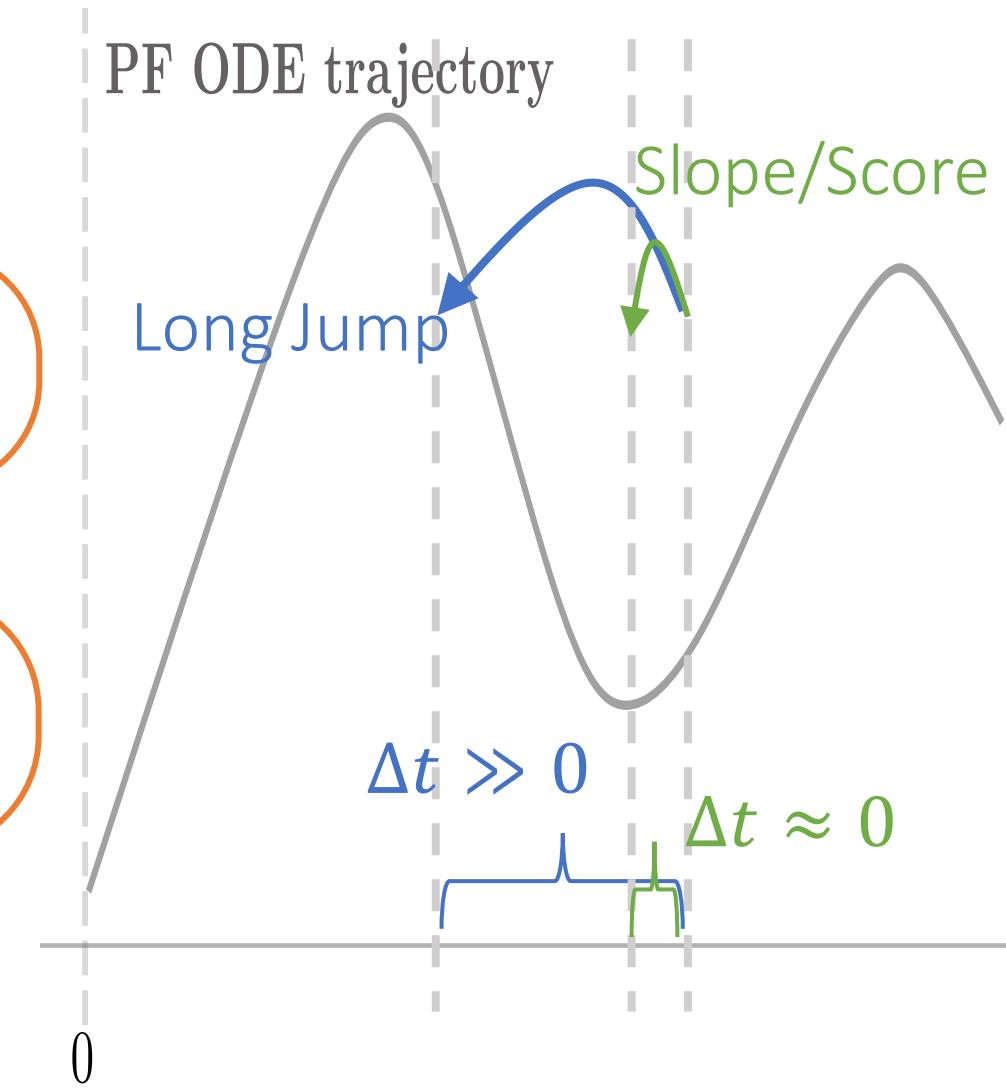
# Long jump ( $\Delta t \gg 0$ )

CTM enables long jump along the trajectory  
→ Allow our new sampling method

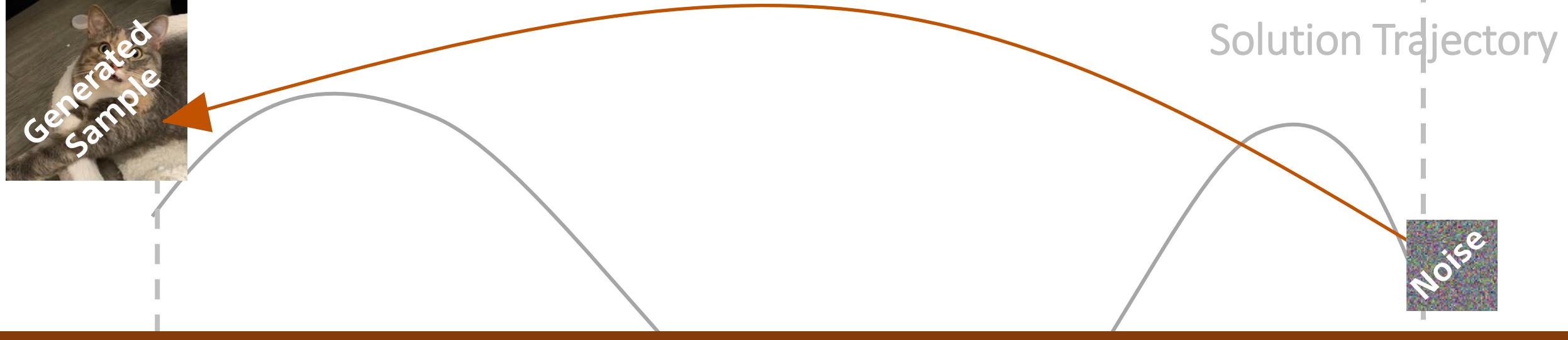
Slope (infinitesimal jump  $\Delta t \approx 0$ )

CTM enables score evaluations

→ Allow score-based sampling & likelihood computation



# One-step Generation with CTM



$G_{\theta}(\text{Noise}, t_{\text{start}} = T, t_{\text{end}} = 0)$

Clean

76 // Sony AI

Noise

T

# SOTA (beats DM) with One-Step



Teacher: Sampling steps = 79

Takes minutes

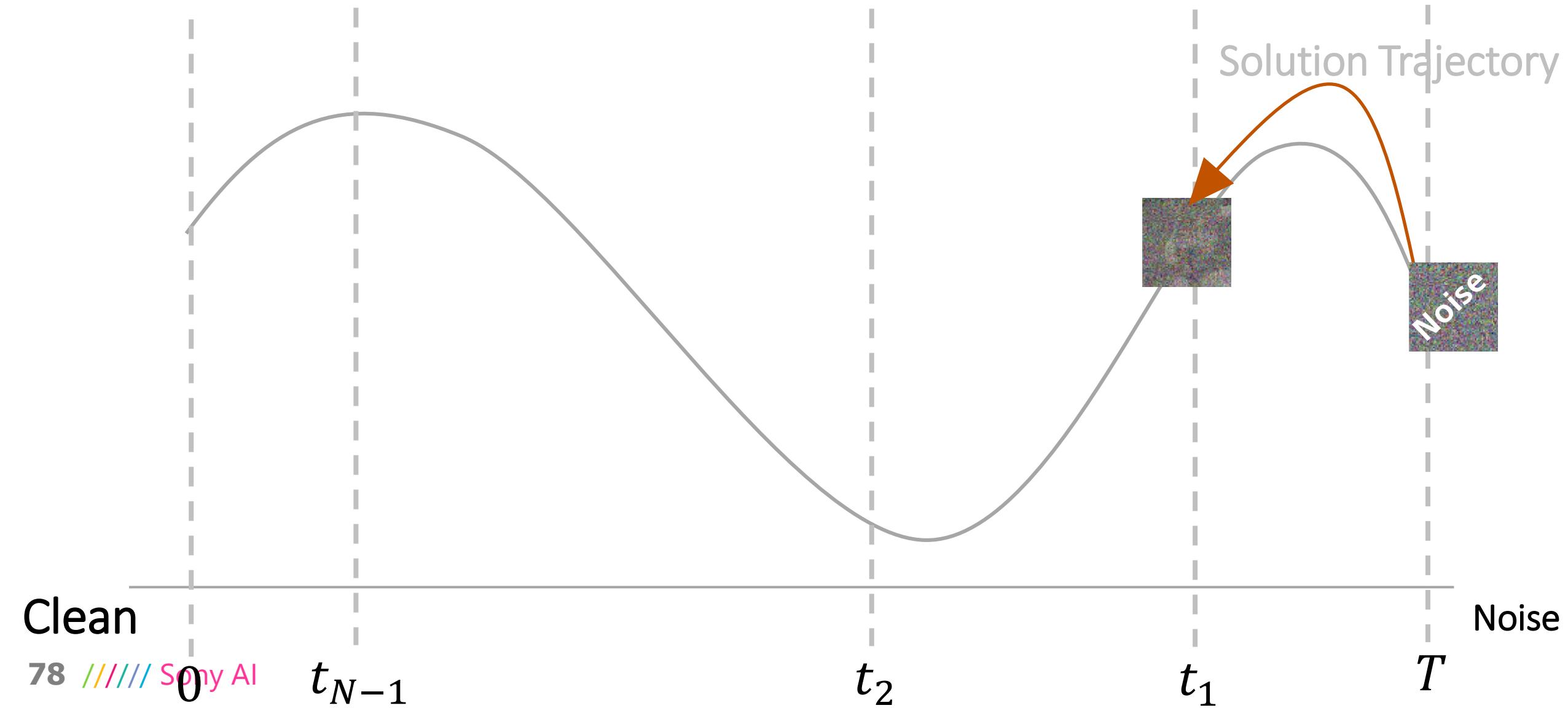


Student (CTM): Sampling step = 1

Within a second

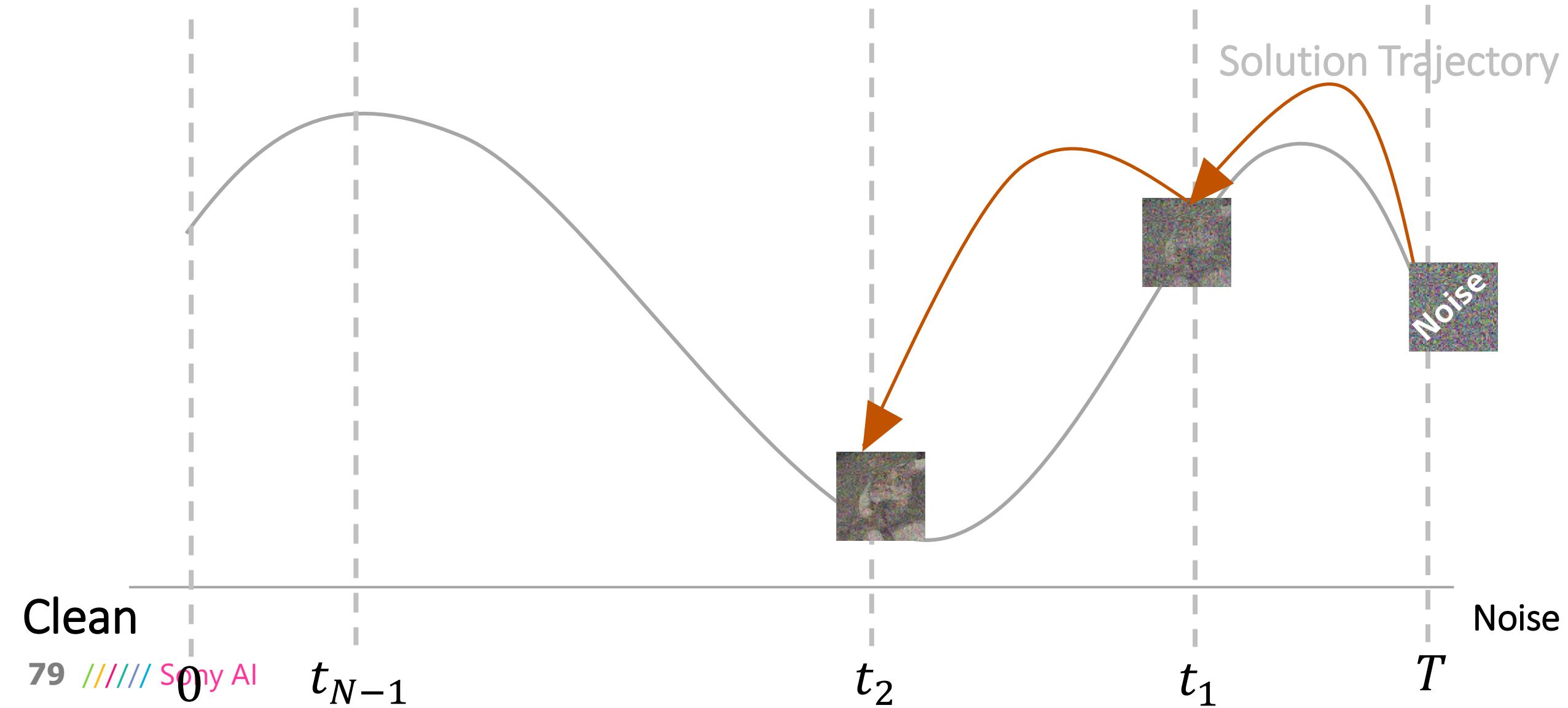
# Multistep Generation with CTM: $\gamma$ -sampling

## $\gamma = 0$ : Fully Deterministic Sampling



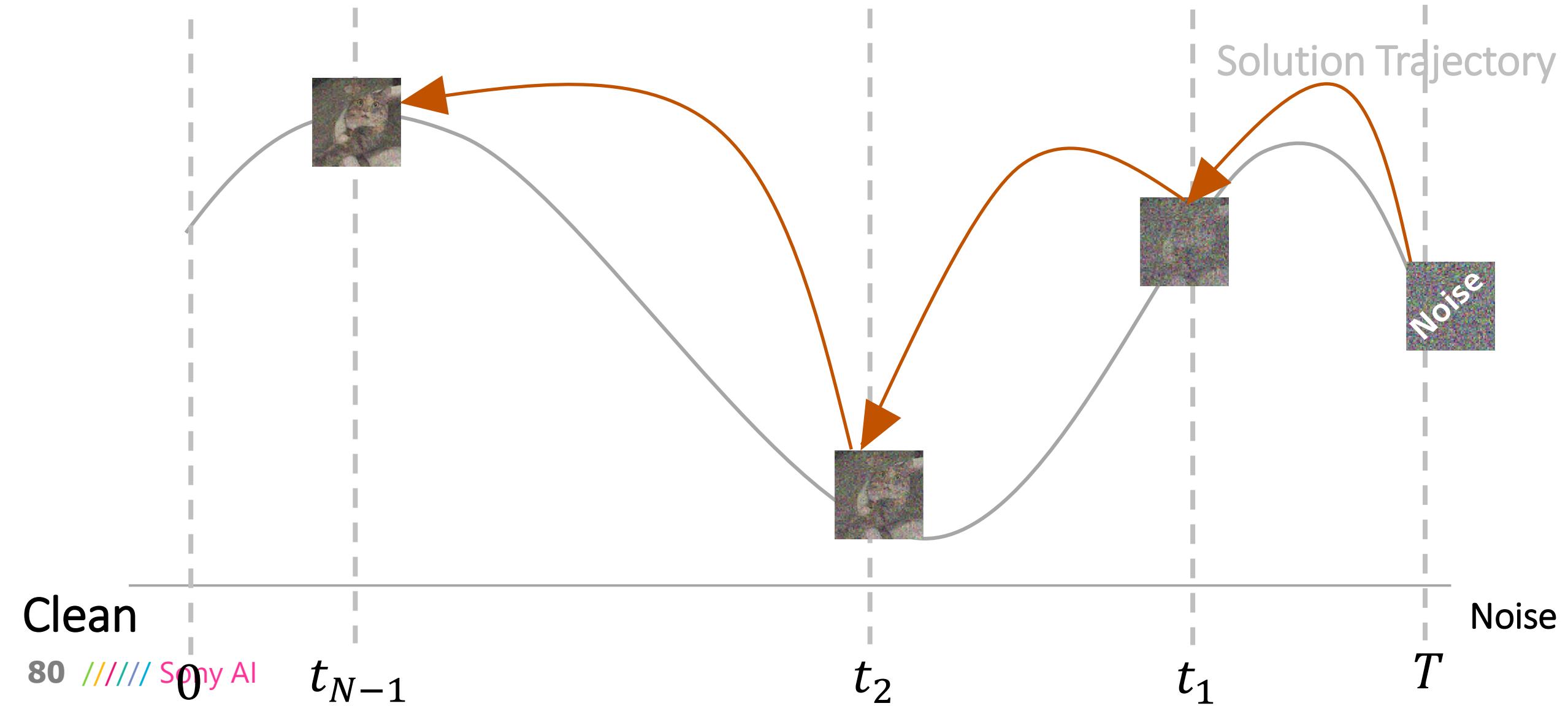
# Multistep Generation with CTM: $\gamma$ -sampling

## $\gamma = 0$ : Fully Deterministic Sampling



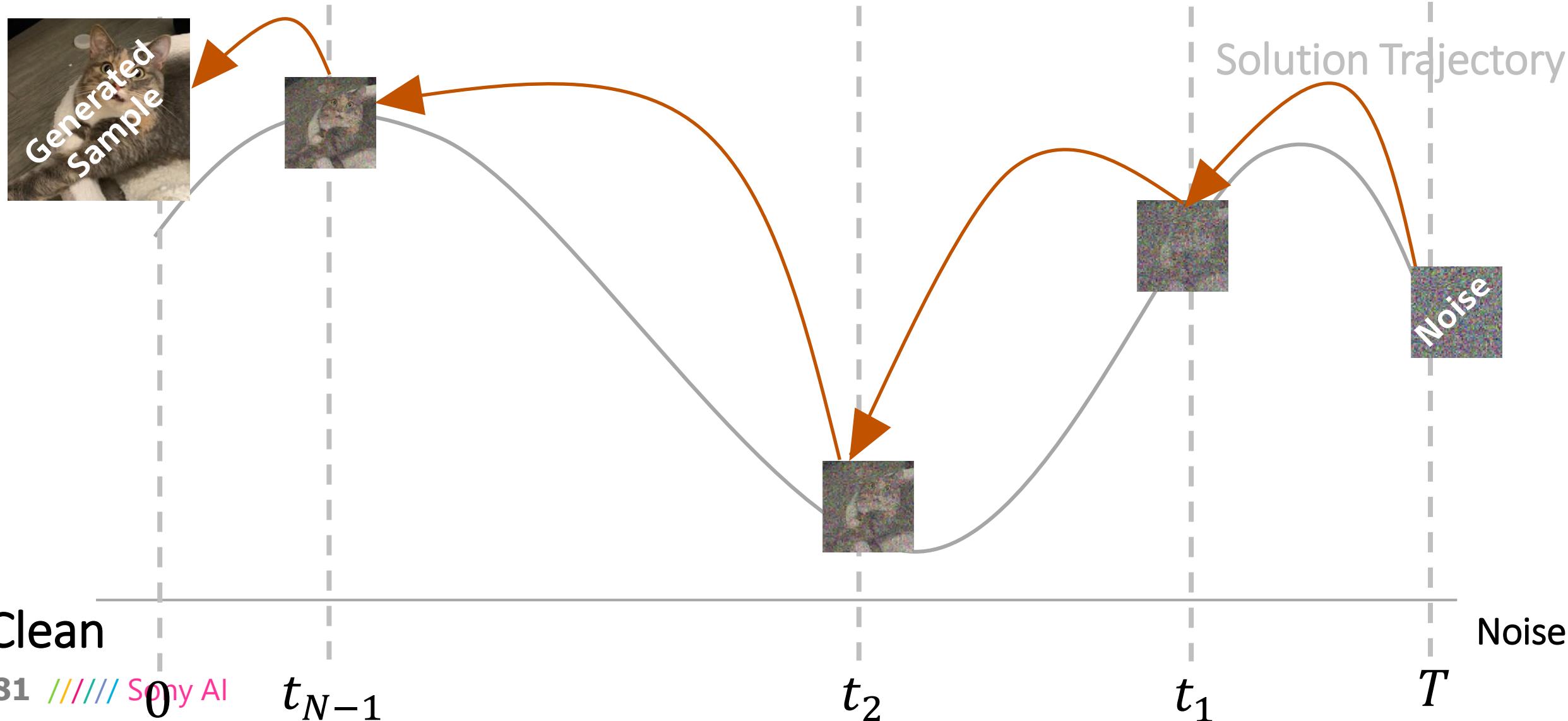
# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 0$ : Fully Deterministic Sampling



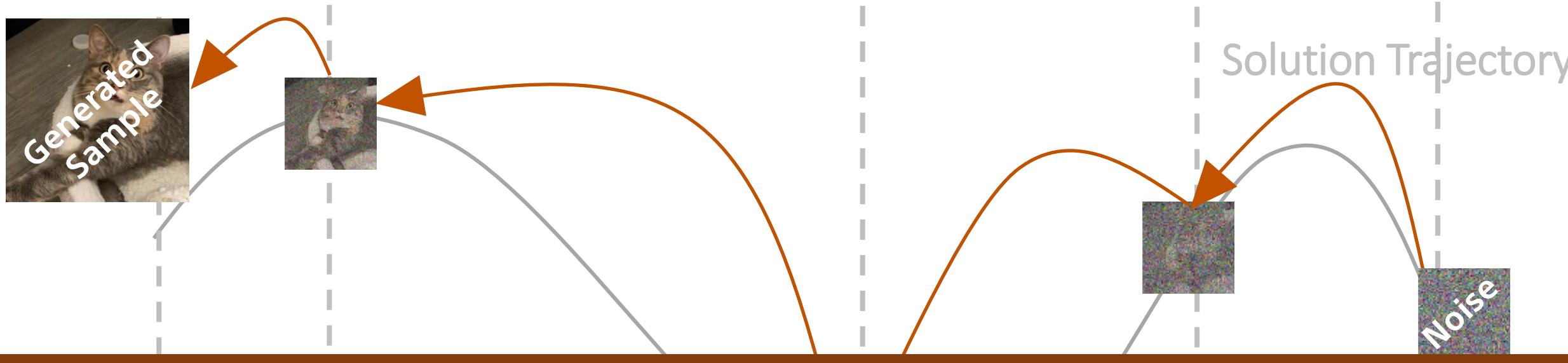
# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 0$ : Fully Deterministic Sampling



# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 0$ : Fully Deterministic Sampling



$$\mathcal{O}(T)$$

Clean

82 // Sony AI

$t_{N-1}$

$t_2$

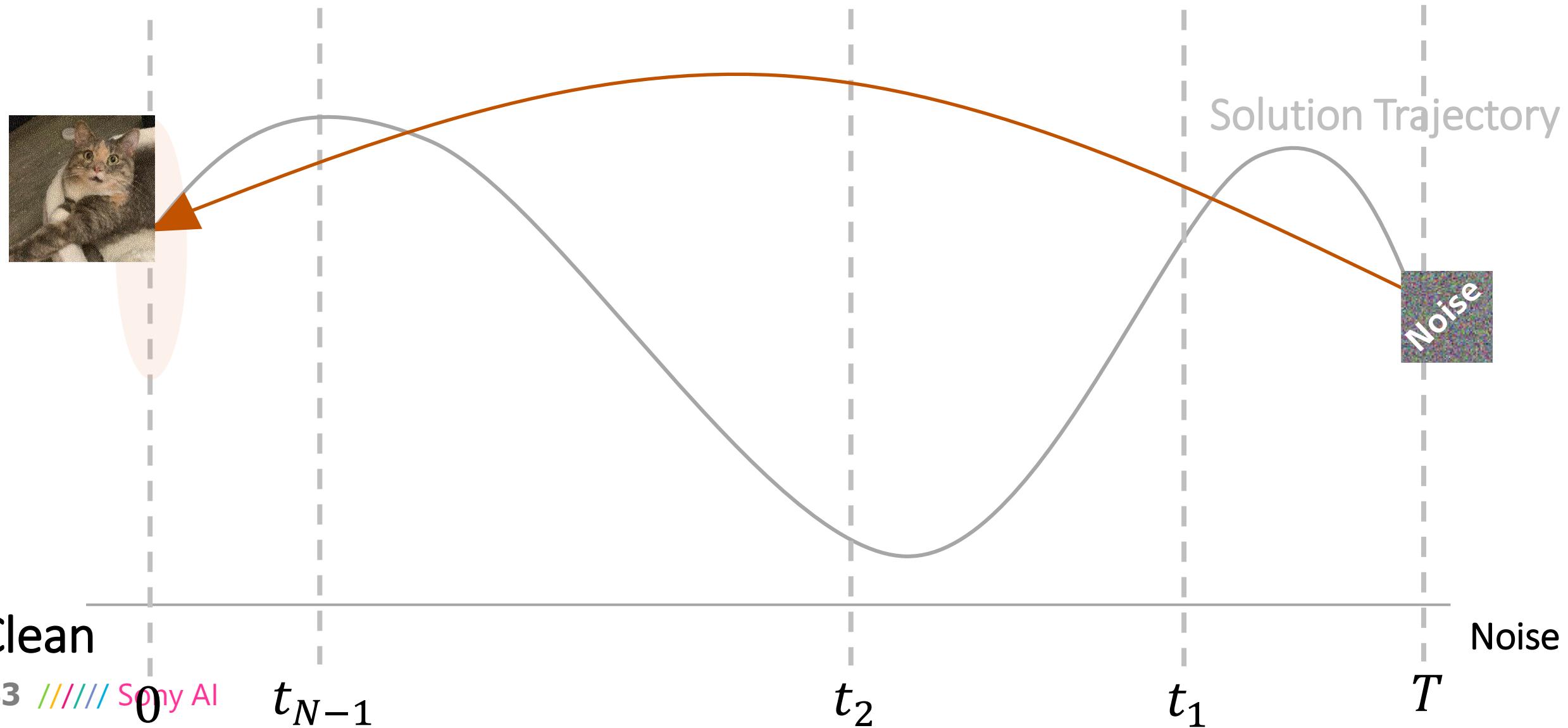
$t_1$

$T$

Noise

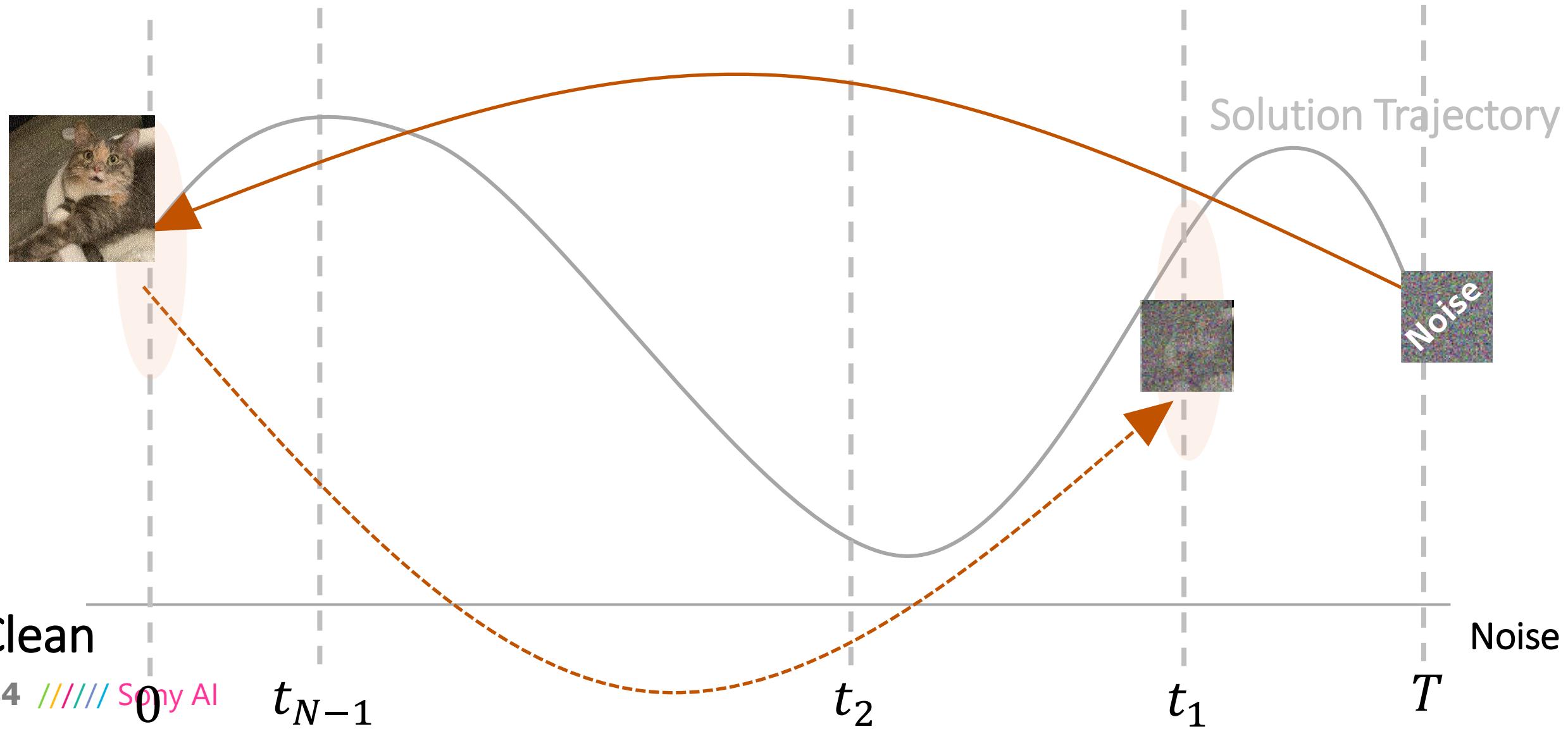
# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling



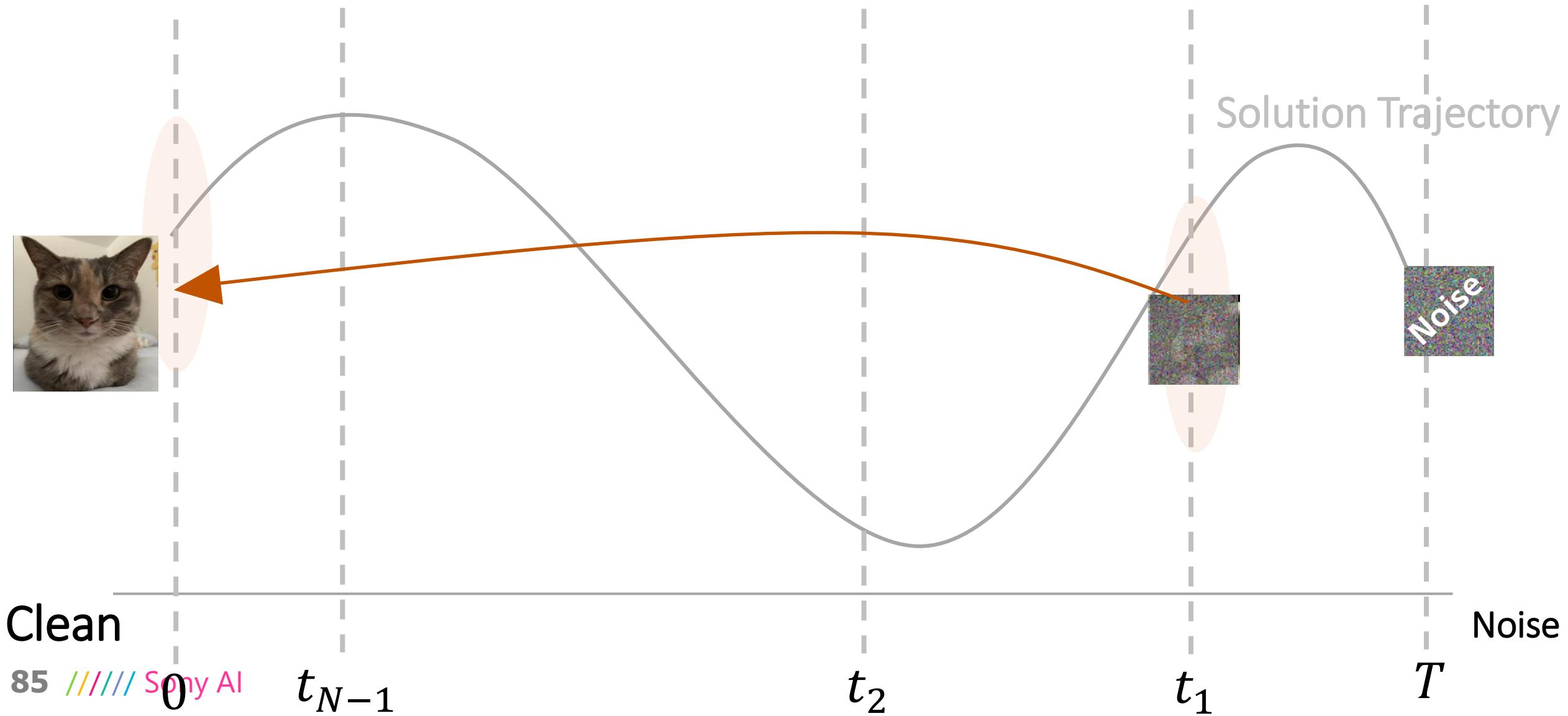
# Multistep Generation with CTM: $\gamma$ -sampling

## $\gamma = 1$ : Fully Stochastic Sampling



# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling



Clean

85 // Sony AI

$t_{N-1}$

$t_2$

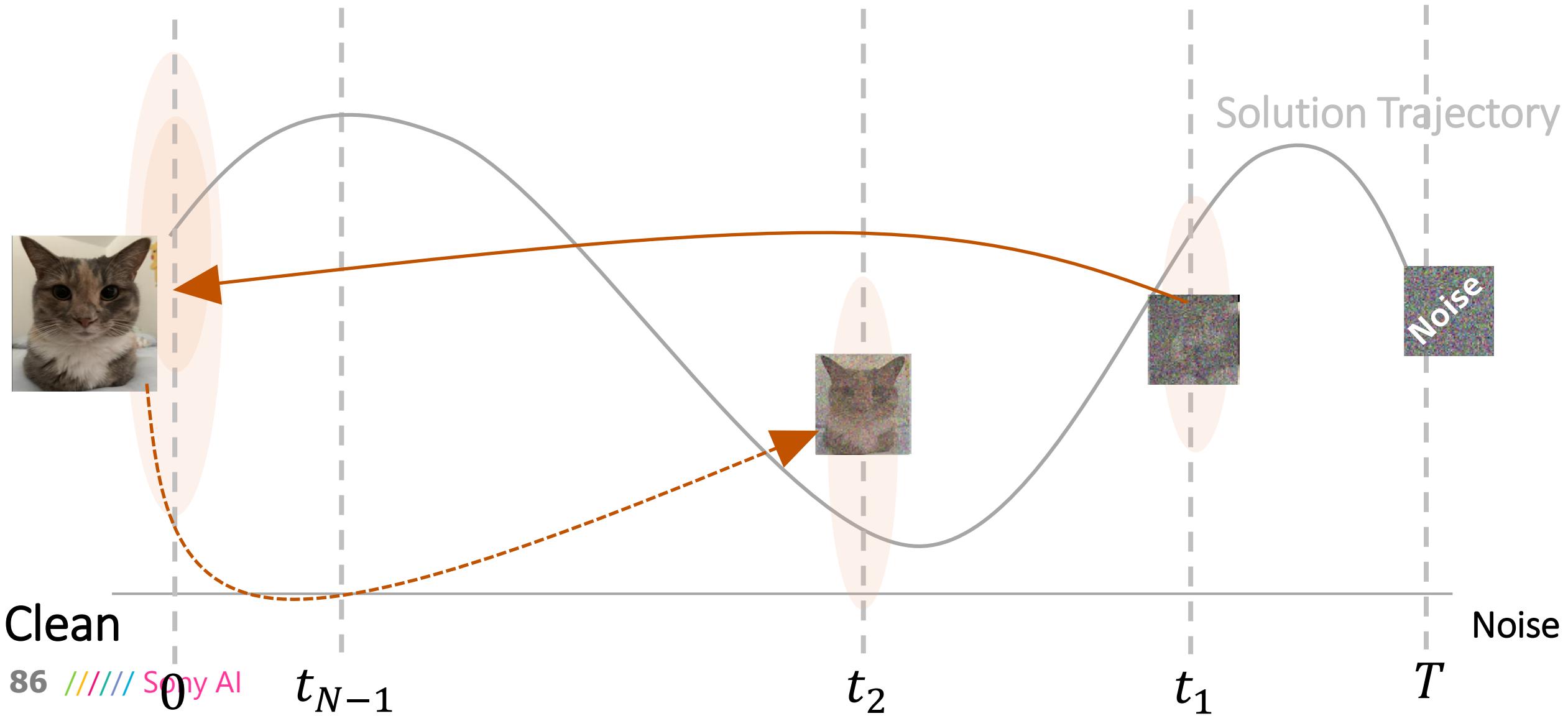
$t_1$

Noise

$T$

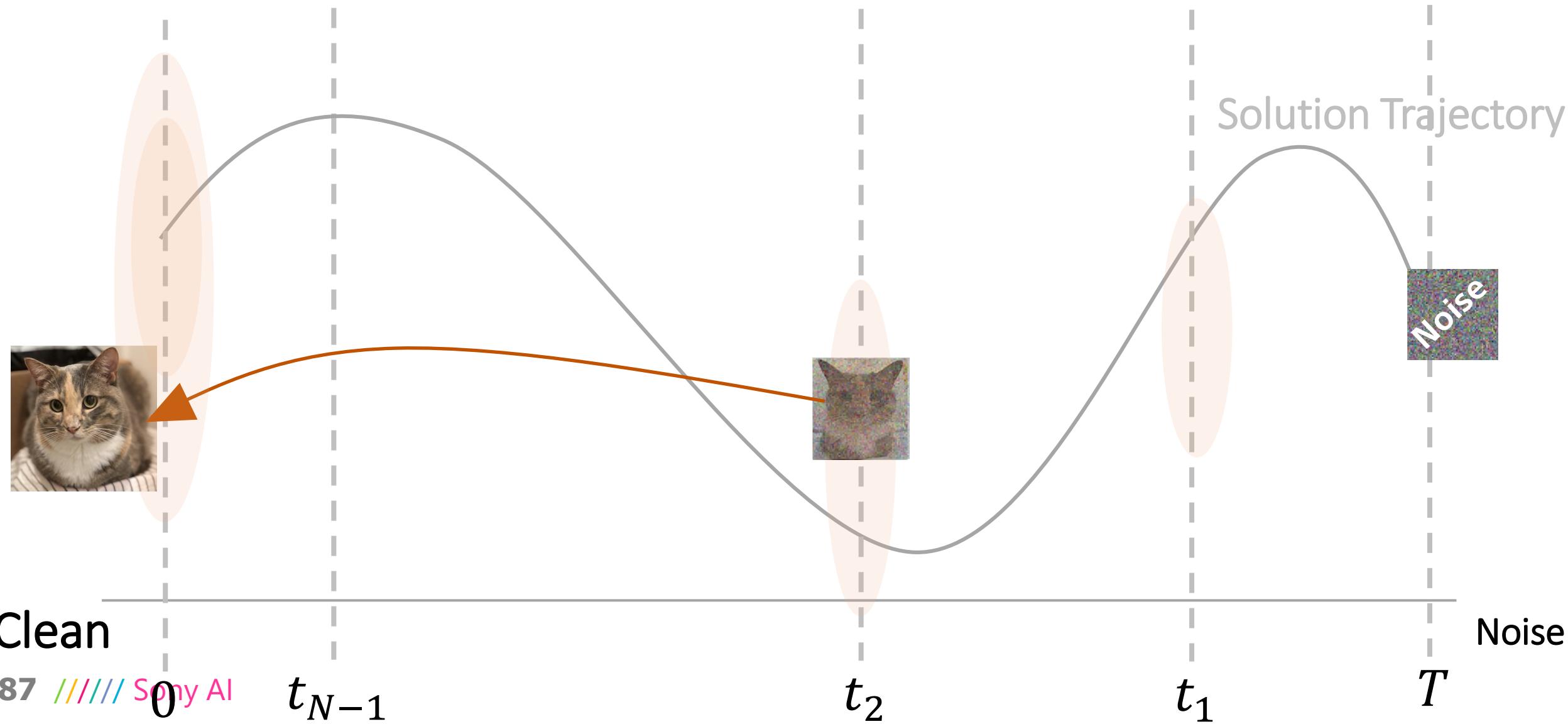
# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling



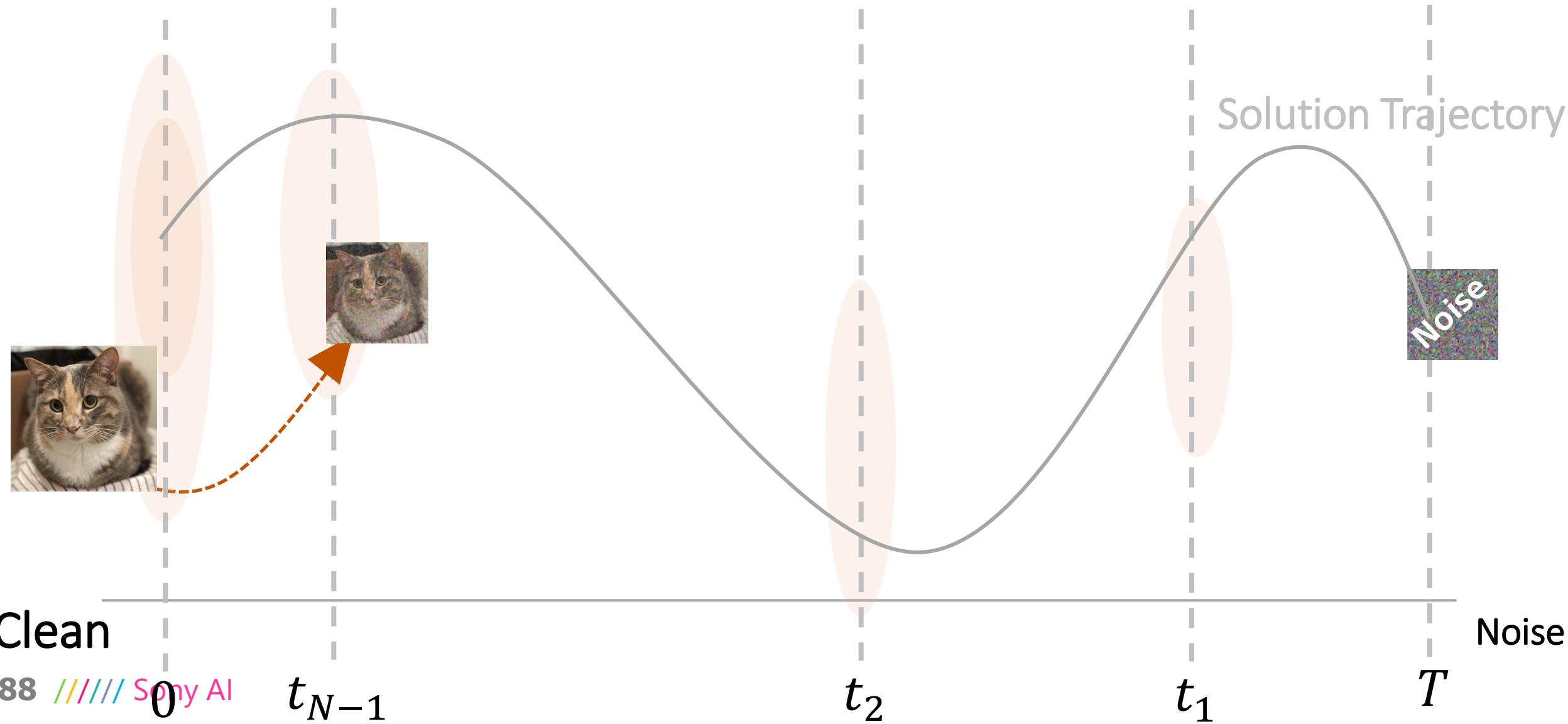
# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling

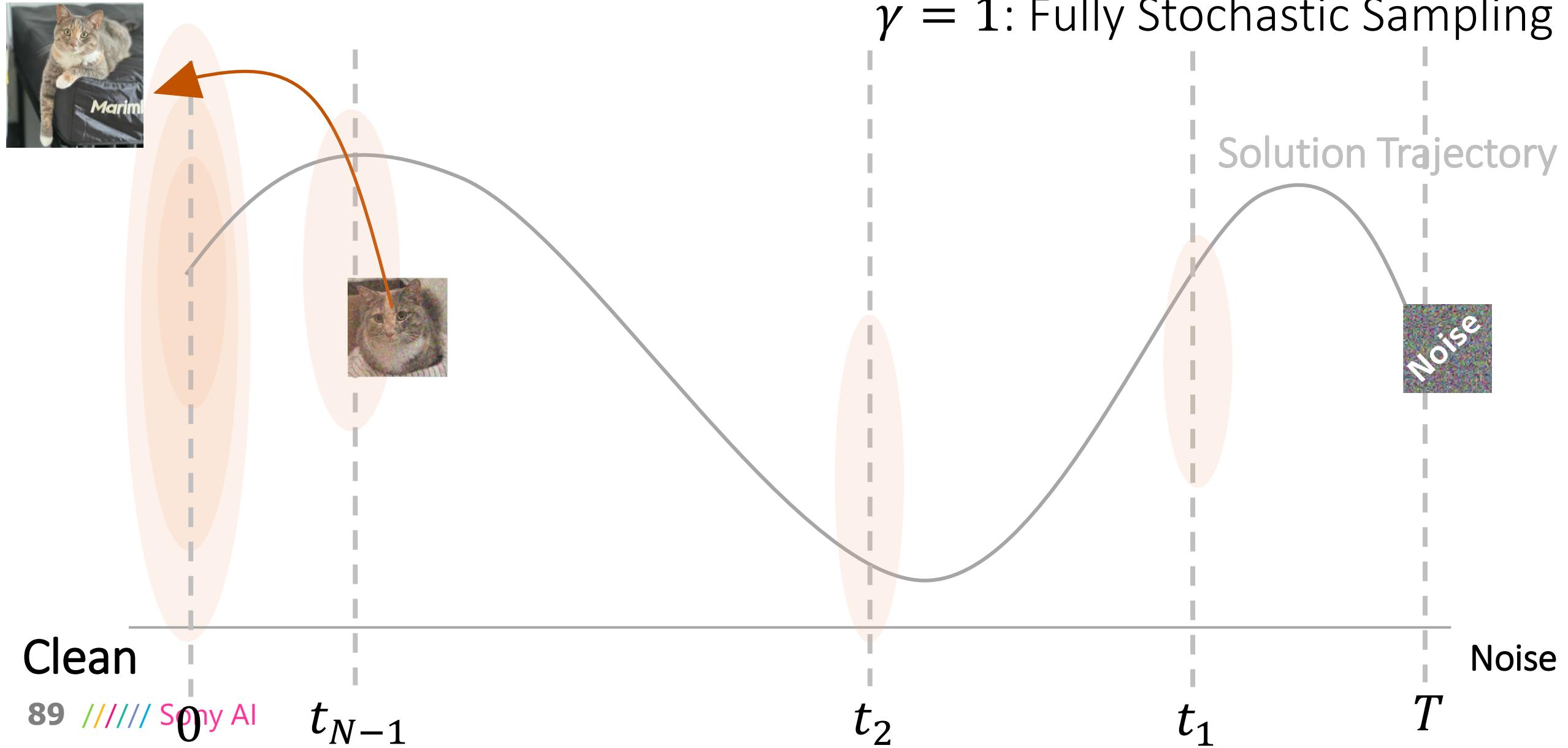


# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling

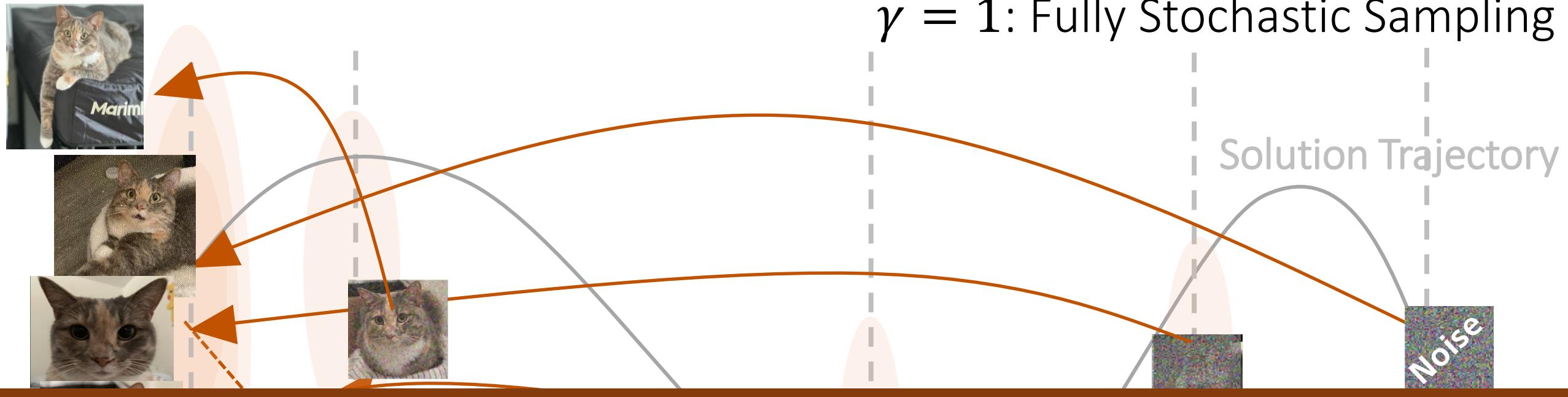


# Multistep Generation with CTM: $\gamma$ -sampling

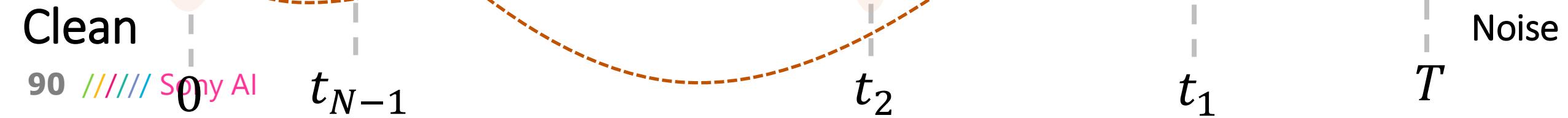


# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma = 1$ : Fully Stochastic Sampling

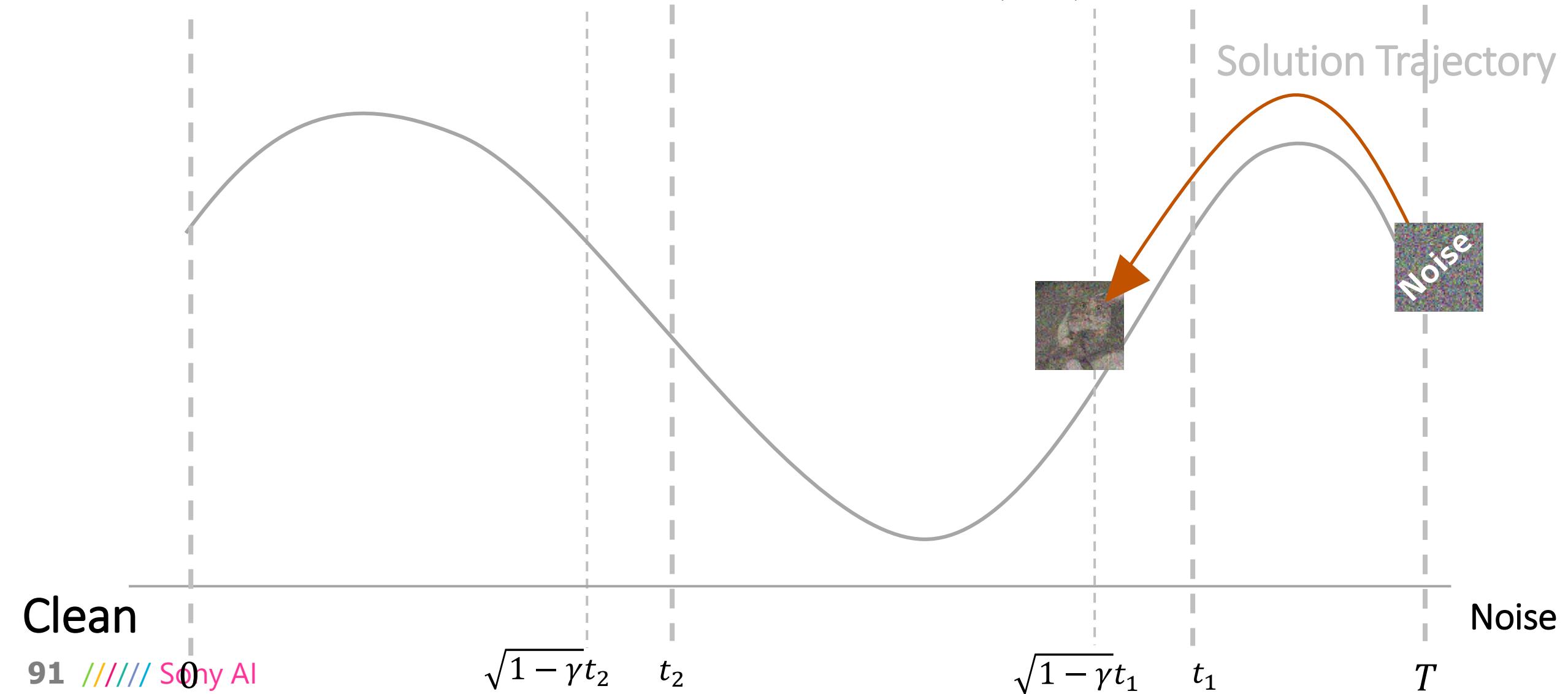


$$\mathcal{O}(T + t_1 + \cdots + t_{N-1})$$



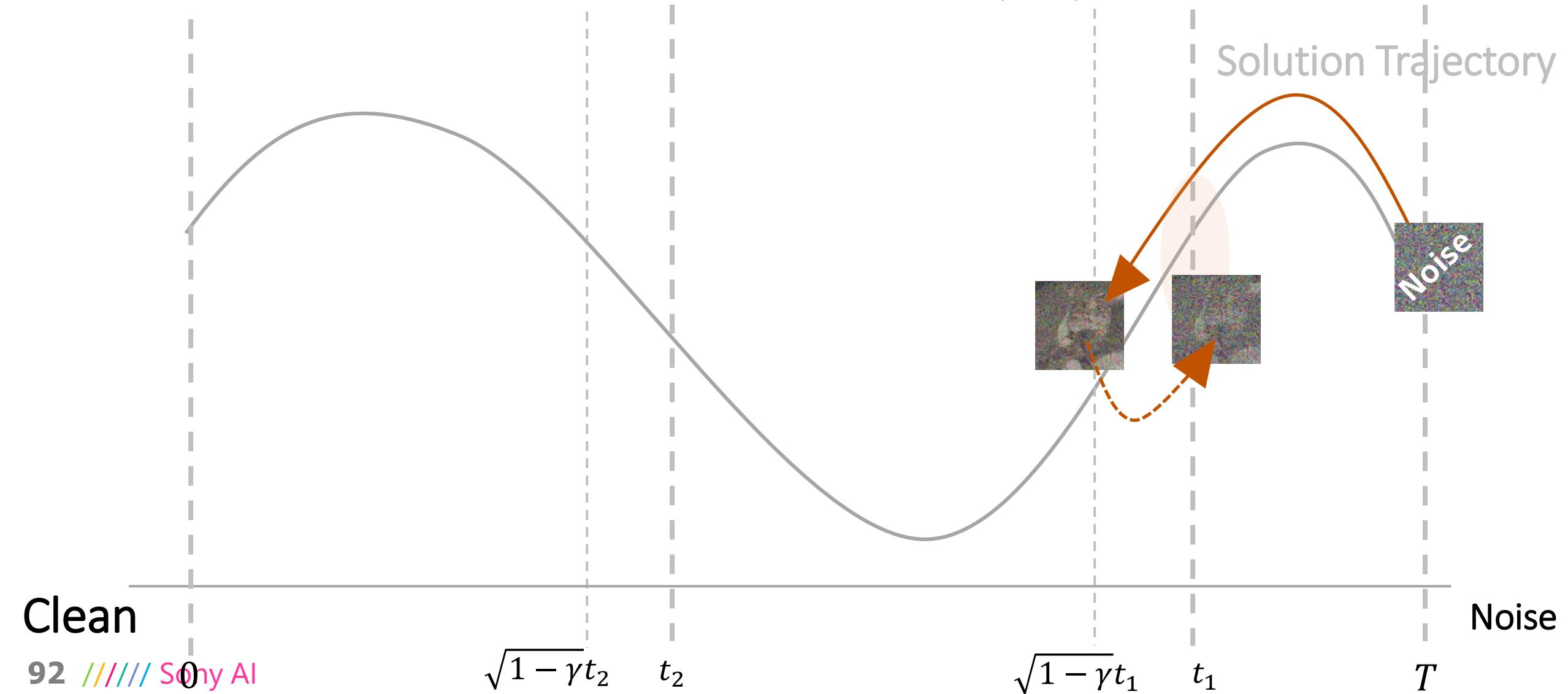
# Multistep Generation with CTM: $\gamma$ -sampling

# $\gamma \in (0,1)$ : Stochastic Sampling



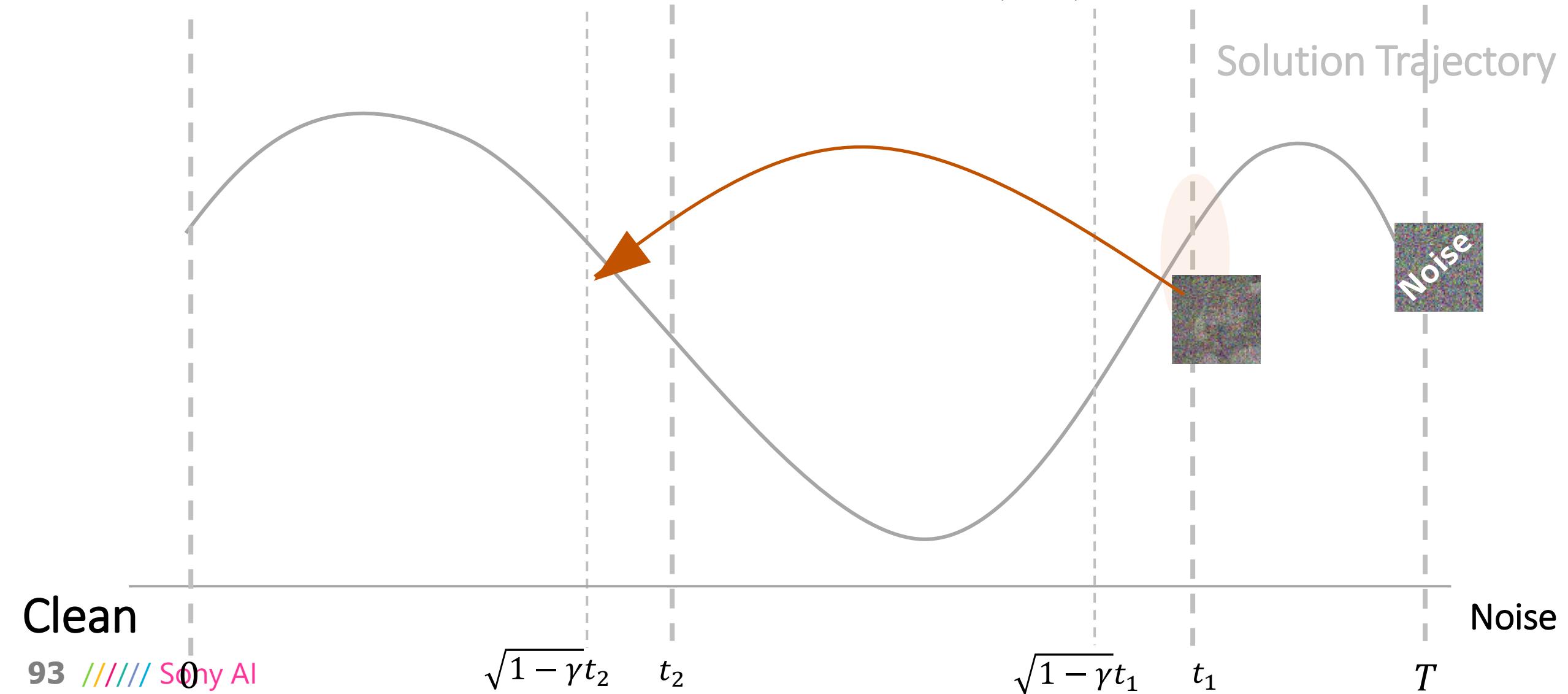
# Multistep Generation with CTM: $\gamma$ -sampling

# $\gamma \in (0,1)$ : Stochastic Sampling



# Multistep Generation with CTM: $\gamma$ -sampling

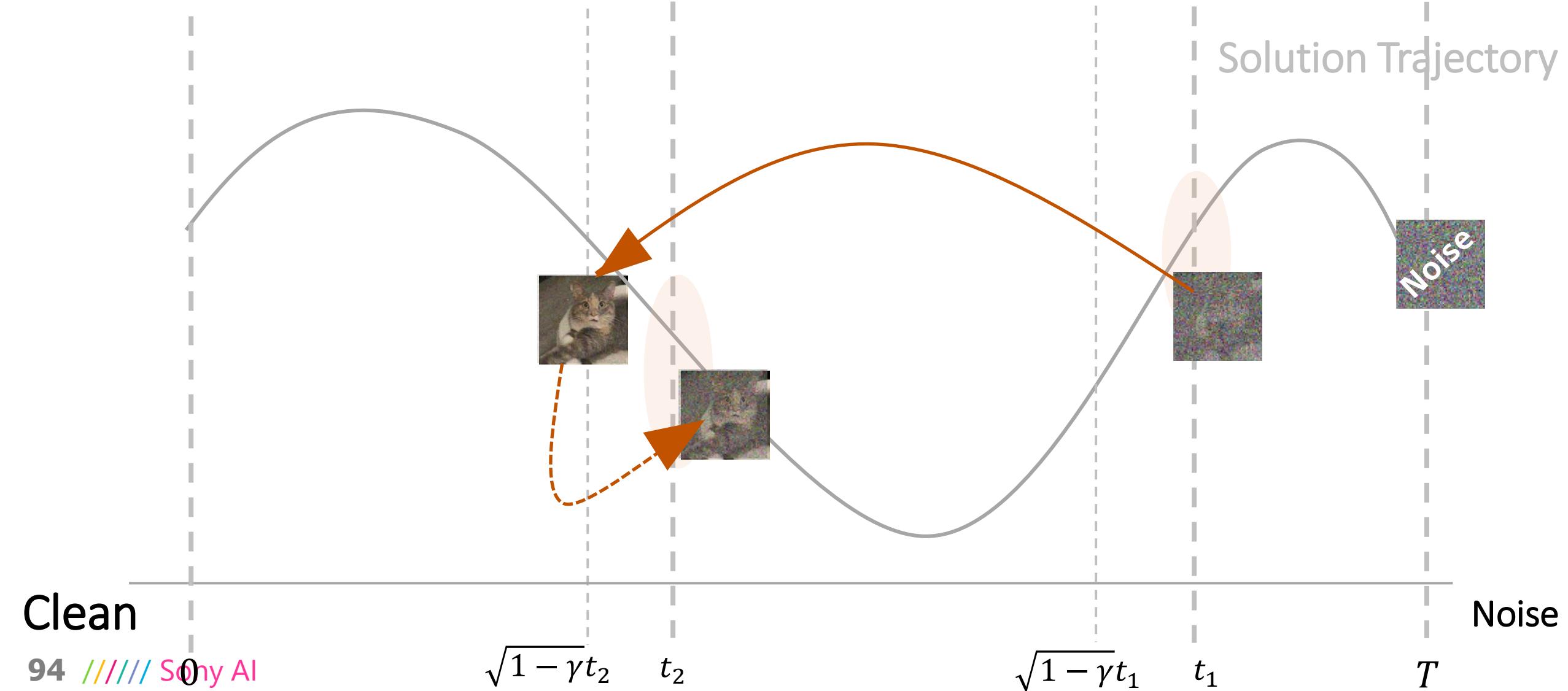
$\gamma \in (0,1)$ : Stochastic Sampling



# Multistep Generation with CTM: $\gamma$ -sampling

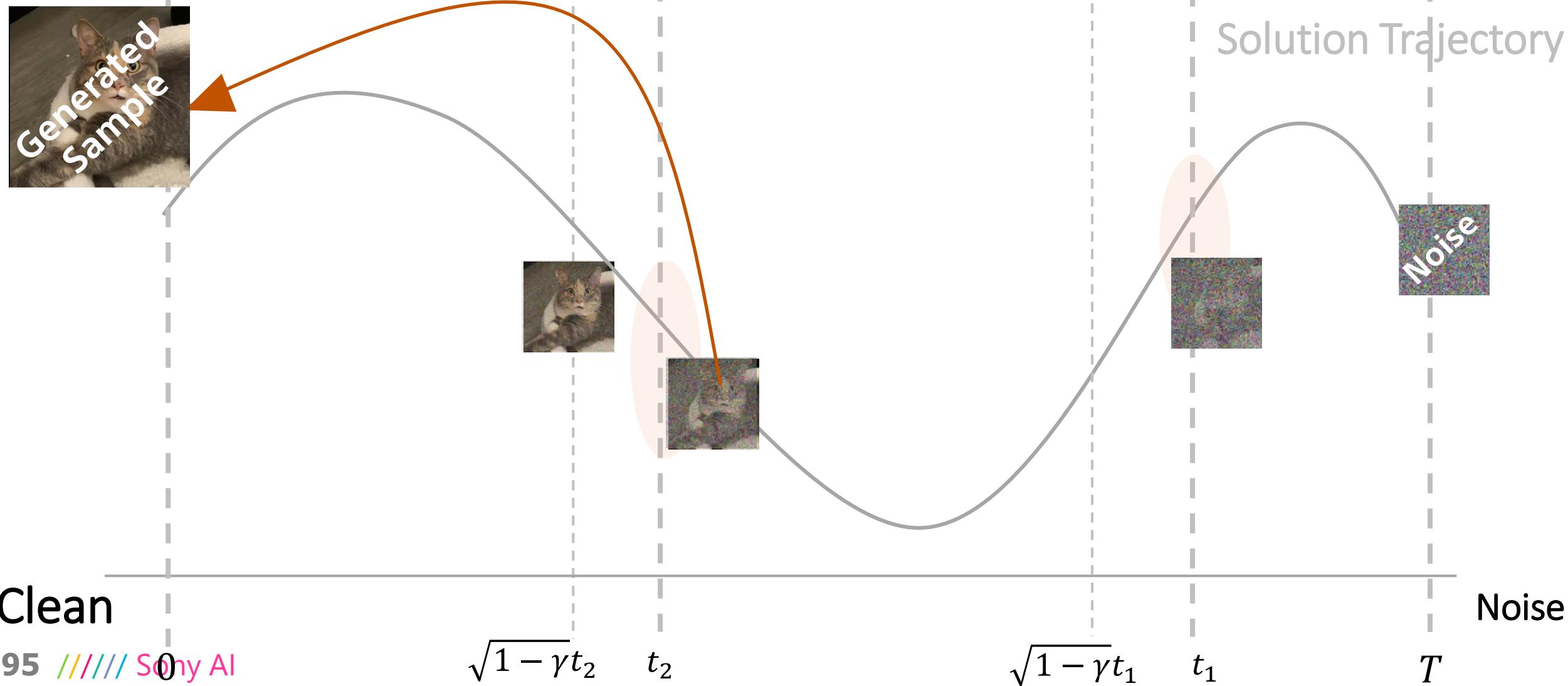
$\gamma \in (0,1)$ : Stochastic Sampling

Solution Trajectory



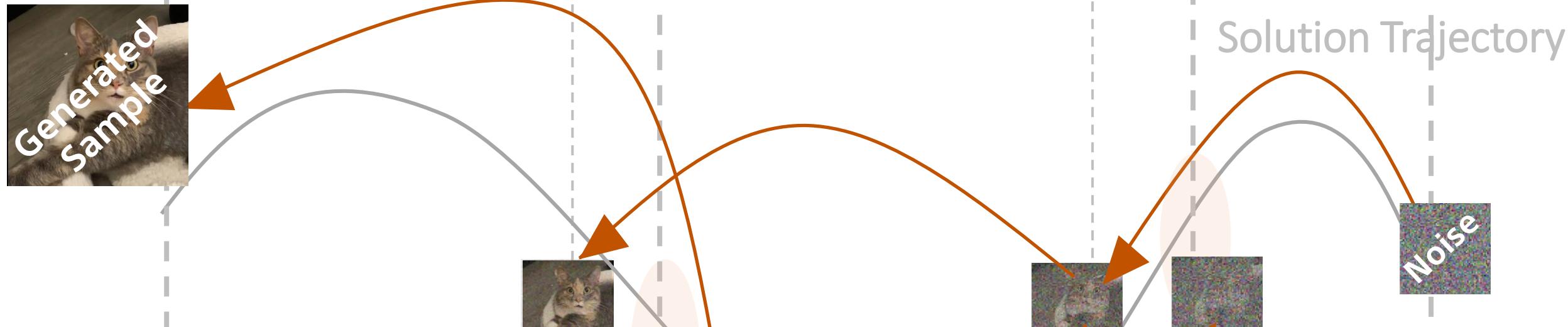
# Multistep Generation with CTM: $\gamma$ -sampling

# $\gamma \in (0,1)$ : Stochastic Sampling



# Multistep Generation with CTM: $\gamma$ -sampling

$\gamma \in (0,1)$ : Stochastic Sampling



$$\mathcal{O} \left( \sum_{n=0}^{N-1} \sqrt{t_n - \sqrt{1 - \gamma^2} t_{n+1}} \right)$$

# CTM is a General Recipe to Many Modalities

# SoundCTM: Uniting Score-based and Consistency Models for Text-to-Sound Generation

**Koichi Saito**  
Sony AI  
NY, USA  
[koichi.saito@sony.com](mailto:koichi.saito@sony.com)

Dongjun Kim  
Stanford University  
CA USA

**Takashi Shibuya** Sony AI  
Tokyo, Japan

**Zhi Zhong**  
Sony Group Corporation  
Tokyo, Japan

**Yuhta Takida**  
Sony AI  
Tokyo, Japan

**Yuki Mitsufuji**  
Sony AI, Sony Group Corporation  
NY, USA

## DITTO-2: Distilled Diffusion Inference-Time T-Optimization for Music Generation

Zachary Novack<sup>1</sup>

Julian McAuley<sup>1</sup>

Taylor Berg-Kirkpatrick<sup>1</sup>

Nicholas Bryan<sup>2</sup>

<sup>1</sup> University of California – San Diego

<sup>2</sup> Adobe Research

[znovack@ucsd.edu](mailto:znovack@ucsd.edu), [njb@ieee.org](mailto:njb@ieee.org)

# SoundCTM's Generation with $\gamma$ -sampling

---

## SoundCTM: Uniting Score-based and Consistency Models for Text-to-Sound Generation

Koichi Saito<sup>1</sup>, Dongjun Kim<sup>2</sup>, Takashi Shibuya<sup>1</sup>, Chieh-Hsin Lai<sup>1</sup>,  
Zhi Zhong<sup>3</sup>, Yuhta Takida<sup>1</sup>, Yuki Mitsufuji<sup>1,3</sup>,

<sup>1</sup>Sony AI, <sup>2</sup>Stanford University, <sup>3</sup>Sony Group Corporation



Paper PDF



arXiv



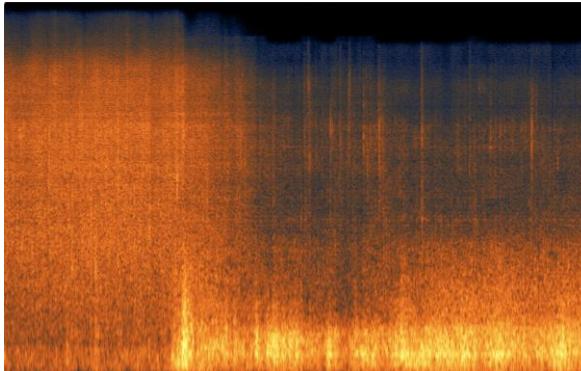
Code

Hugging Face

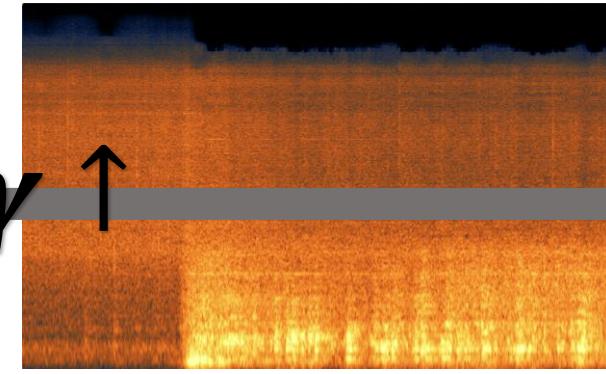
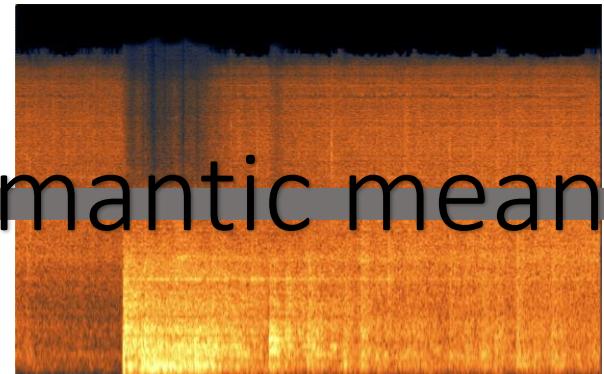
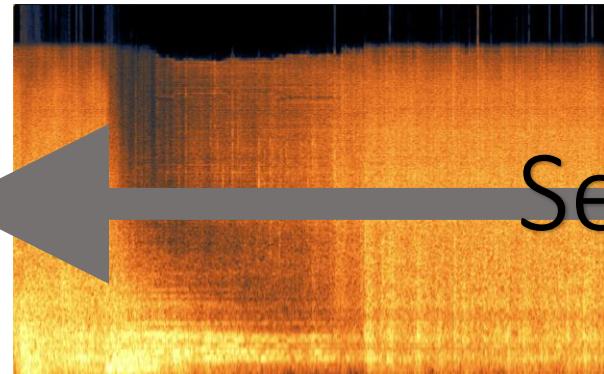
# SoundCTM's Generation with $\gamma$ -sampling

**User specified Text prompt:** Thunderclaps, and hard rain falls and splashes on surfaces.

One step



16 steps



Semantic meaning varies as  $\gamma \uparrow$



$\gamma=1$



$\gamma=0.5$



$\gamma=0.25$



$\gamma=0$

# Making Diffusion Model More Powerful

# Training

# Inaccuracy of Training

# Sampling

# Slow Generation Speed

# Scalability

# [NeurIPS'24 Kim&Lai+] High Cost for PaGoDA



## High-Dim. Datasets



# Motivation of PaGoDA

---

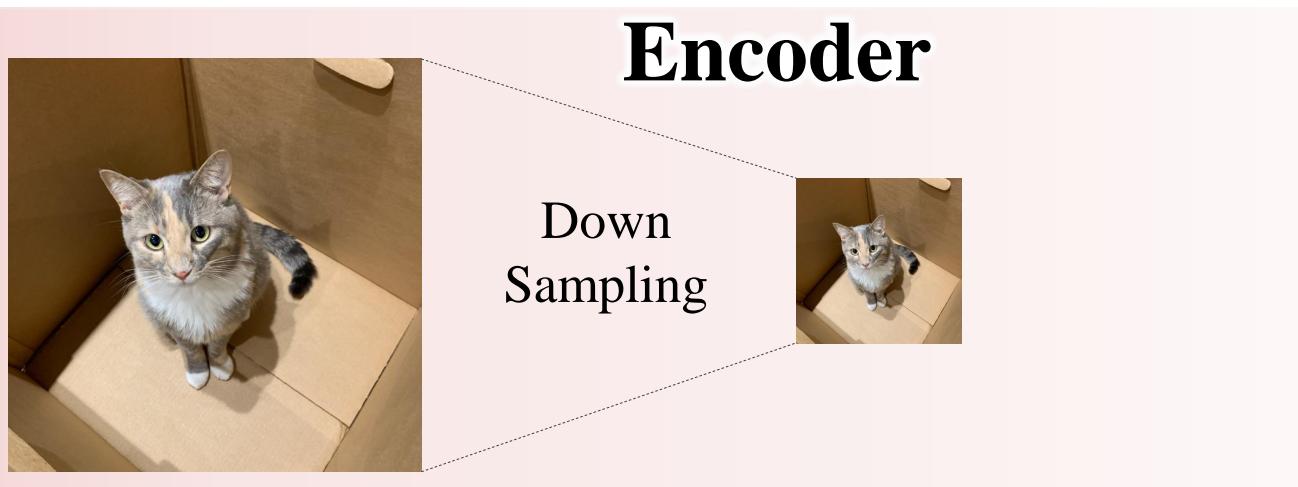
- DMs need to stay in the same dimension
- When image resolution  $\uparrow$ , training and sampling will be costly...
- Can we distill low-resolution DM for a **high-resolution 1-step** generation?
  - Training and Sampling Efficiency

# **P**rogressive **G**rowing of **D**iffusion **A**utoencoder

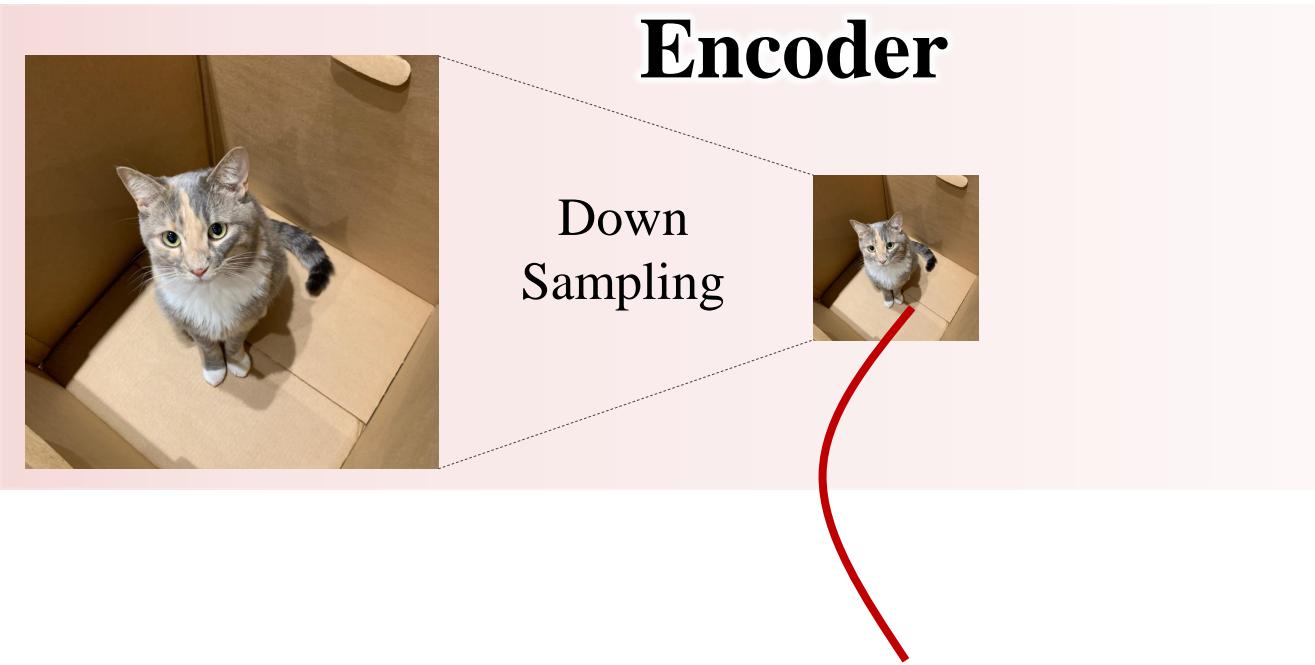
# PaGoDA



# Stage 1. Pretrain DM on low-dim data

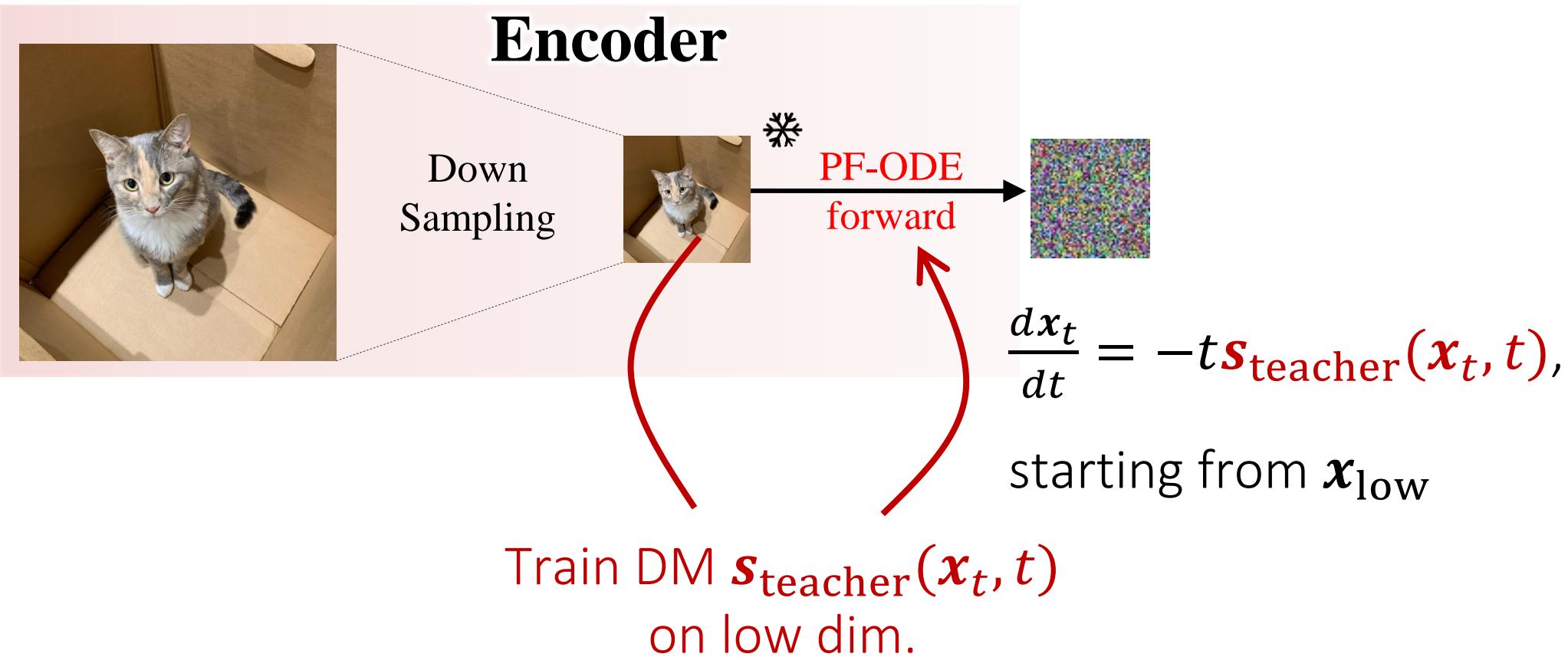


# Stage 1. Pretrain DM on low-dim data

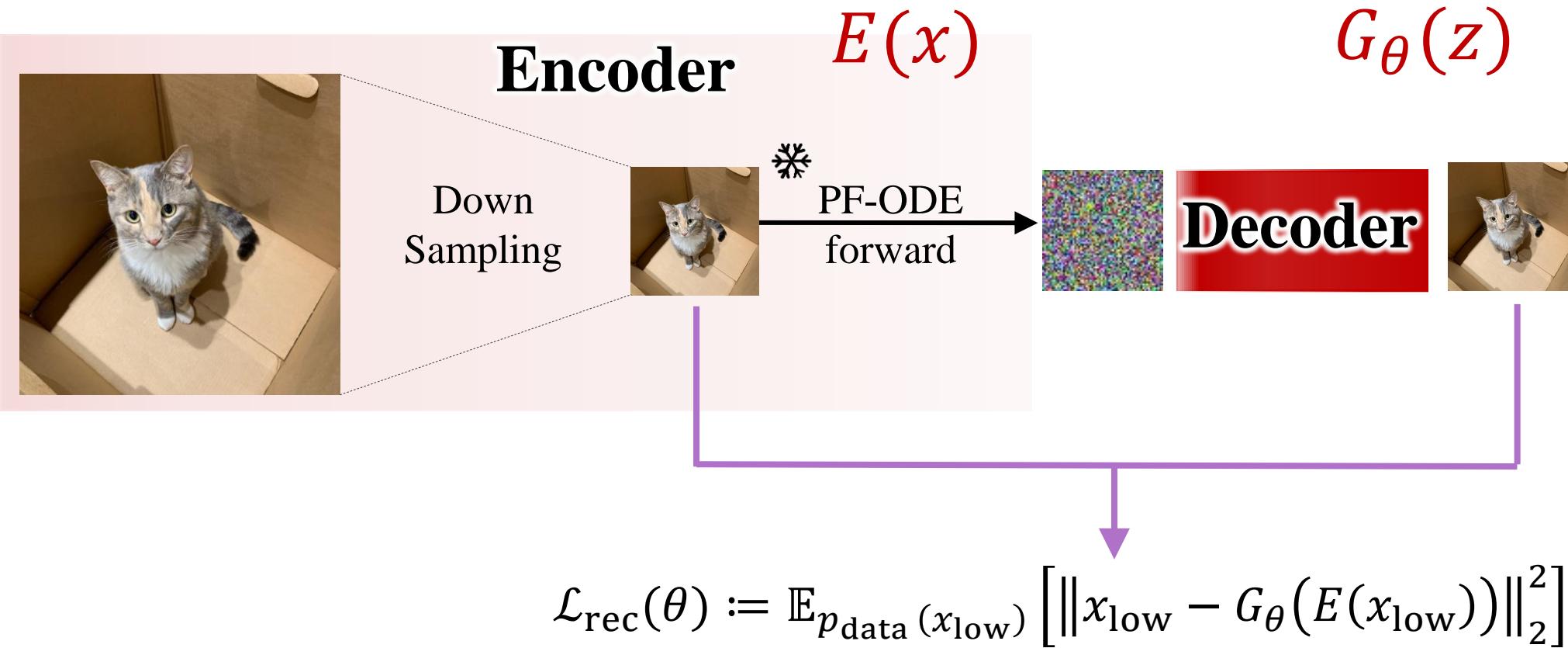


Train DM  $s_{\text{teacher}}(x_t, t)$   
on low dim.

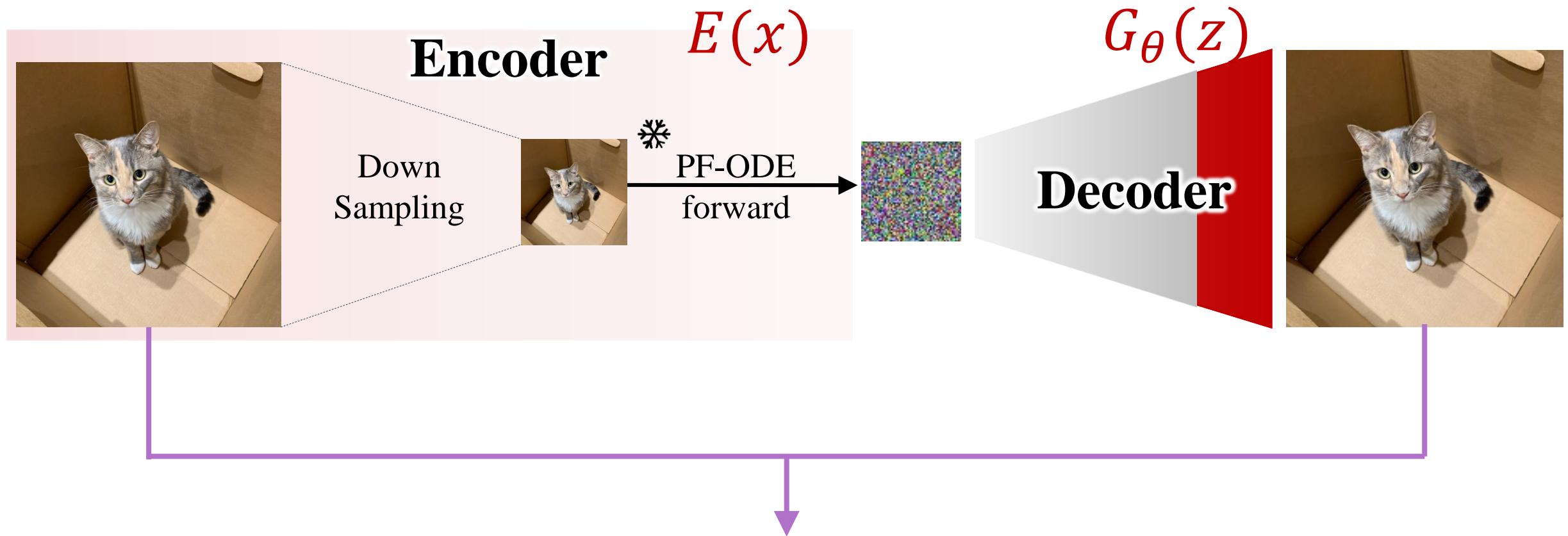
# Stage 1. Pretrain DM on low-dim data



## Stage 2. Distill DM with *Forward* PF-ODE

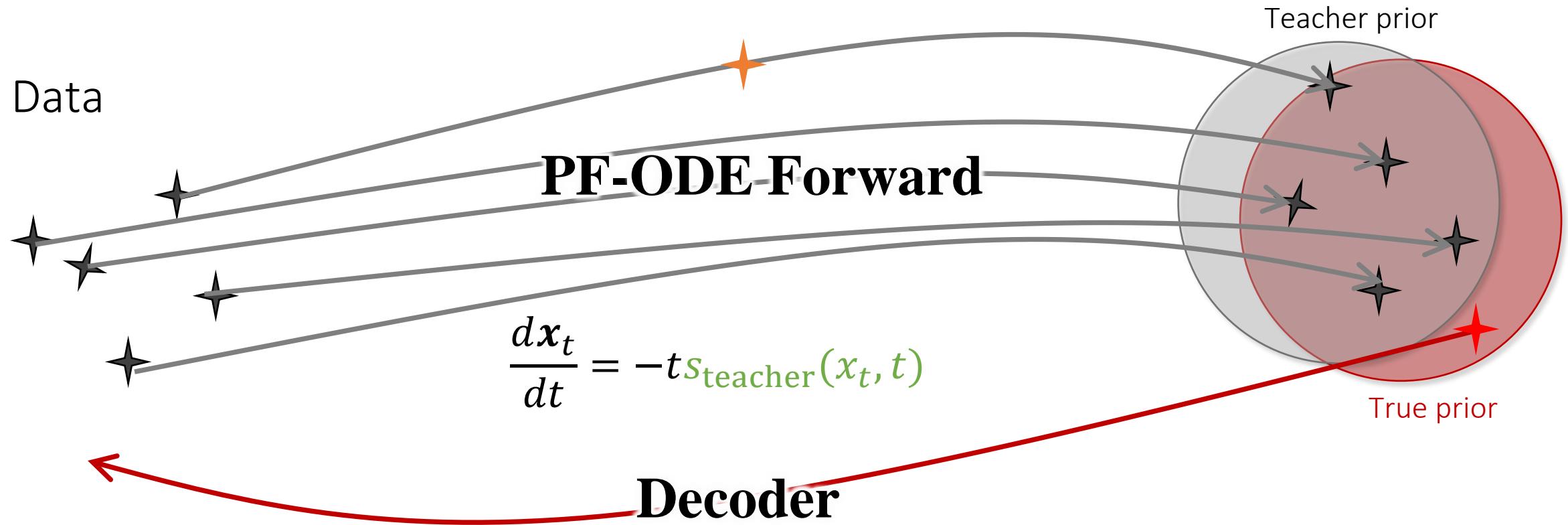


## Stage 3. Progressively Grows Generator's Output Dim.



$$\mathcal{L}_{\text{rec}}(\theta) \coloneqq \mathbb{E}_{p_{\text{data}}(x_{\text{high}})} \left[ \left\| x_{\text{high}} - G_\theta(E(x_{\text{high}})) \right\|_2^2 \right]$$

# Issues of Solely Learning with Rec.

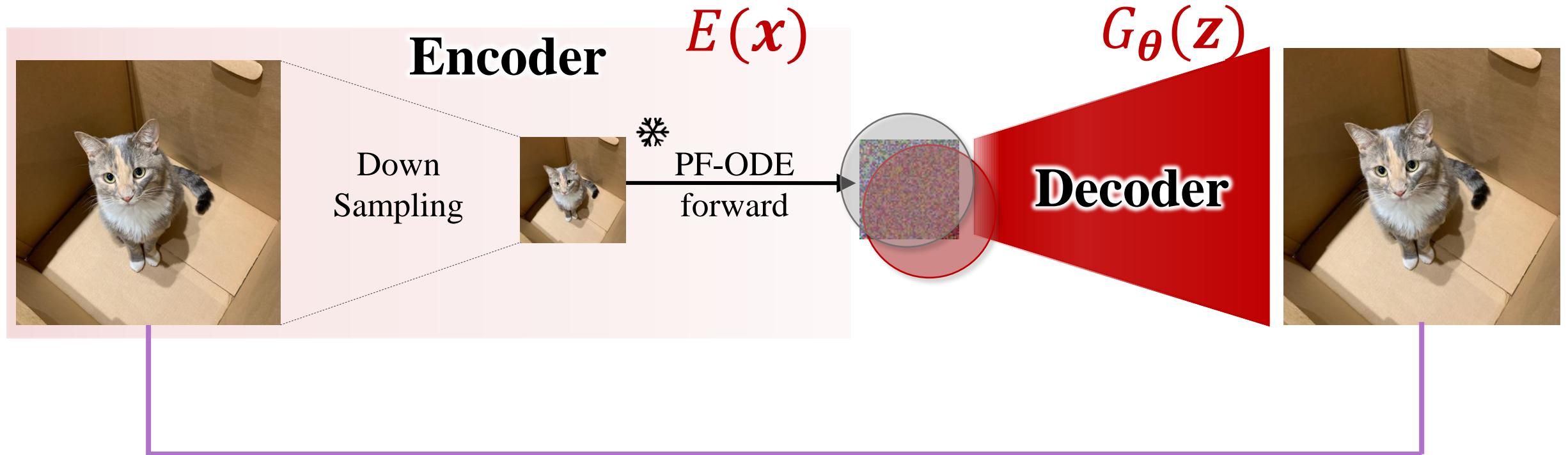


Theorem 1 (Informal) [Density can be improved with enhanced teacher]

$$W_q(p_{\text{PaGoDA\_noGAN}, \theta}, p_{\text{data}}) \lesssim \mathcal{L}_{\text{rec}}(\theta) + \epsilon_{\text{DSM}}, q = 1, 2$$

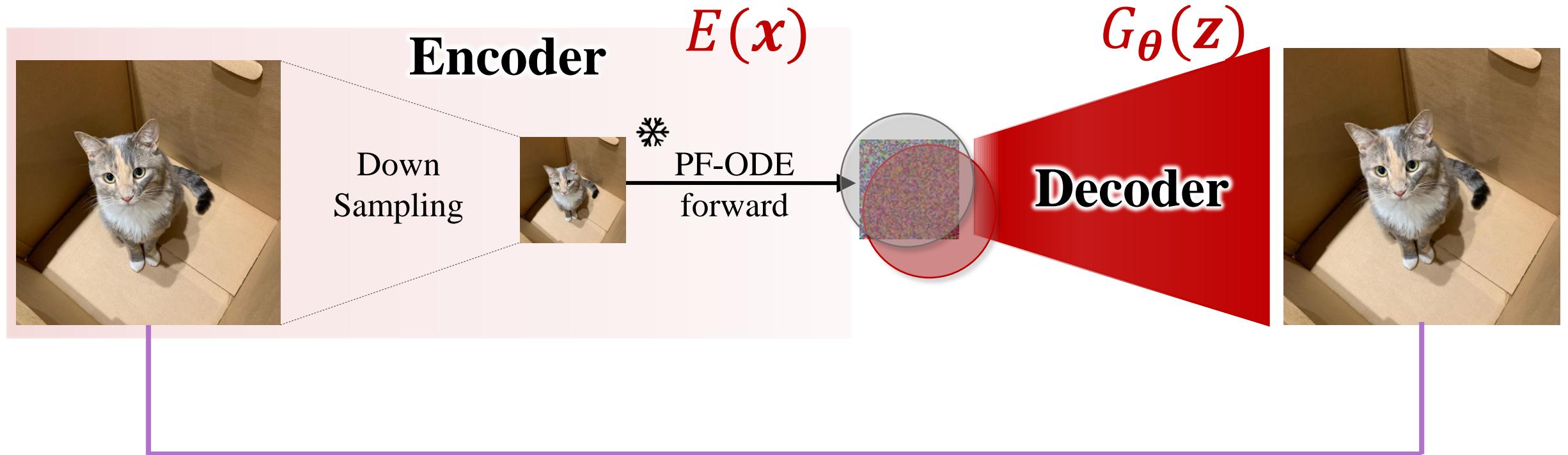
$$\mathcal{L}_{\text{rec}}(\theta) \coloneqq \mathbb{E}_{p_{\text{data}}(x)} \left[ \|x - G_\theta(E(x))\|_2^2 \right]$$

# Overall PaGoDA's Objective



$$\mathcal{L}_{\text{rec}}(G_{\theta}) \coloneqq \mathbb{E}_{p_{\text{data}}(x_{\text{high}})} \left[ \left\| x_{\text{high}} - G_{\theta} \left( E(x_{\text{high}}) \right) \right\|_2^2 \right]$$

# Overall PaGoDA's Objective

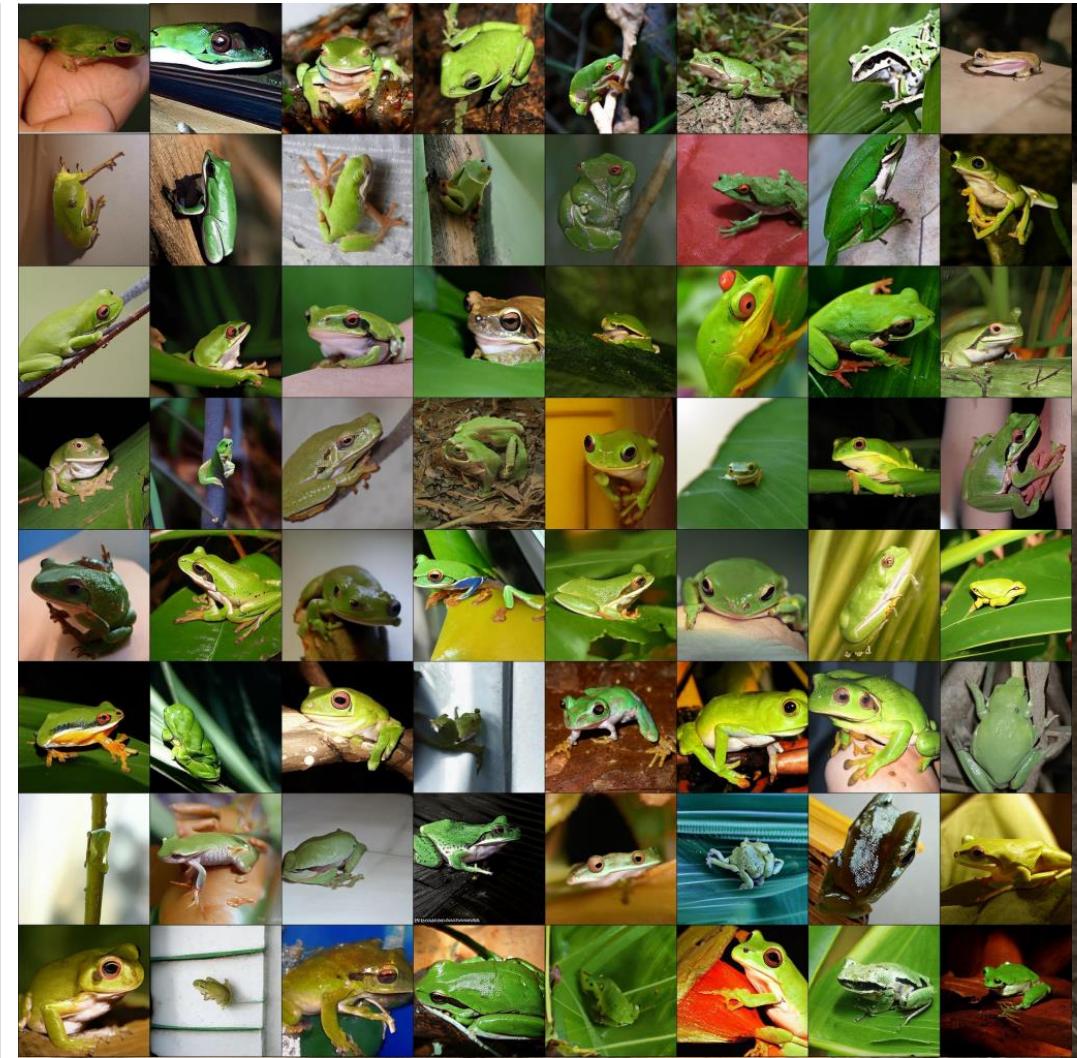
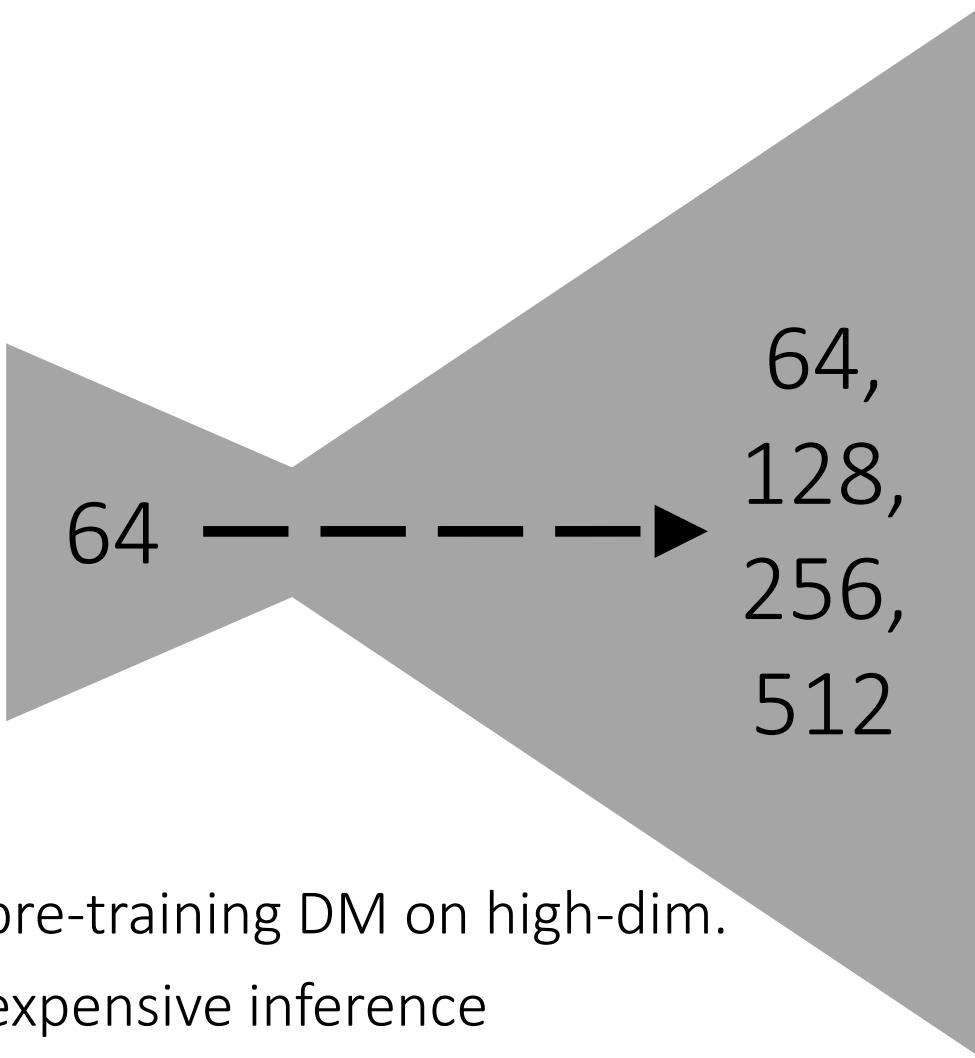


$$\mathcal{L}_{\text{rec}}(G_{\theta}) \coloneqq \mathbb{E}_{p_{\text{data}}(x_{\text{high}})} \left[ \left\| x_{\text{high}} - G_{\theta} \left( E(x_{\text{high}}) \right) \right\|_2^2 \right]$$

$$\mathcal{L}_{\text{GAN}}(G_{\theta}, D_{\psi}) \coloneqq \mathbb{E}_{p_{\text{data}}(x_{\text{high}})}[\log D_{\psi}(x)] + \mathbb{E}_{p_{\text{prior}}(\mathbf{z}_{\text{low}})}\left[\log\left(1 - D_{\psi}(G_{\theta}(\mathbf{z}_{\text{low}}))\right)\right]$$

$$\mathcal{L}_{\text{PaGoDA}}(G_{\theta}, D_{\psi}) := \mathcal{L}_{\text{rec}}(G_{\theta}) + \lambda \mathcal{L}_{\text{GAN}}(G_{\theta}, D_{\psi})$$

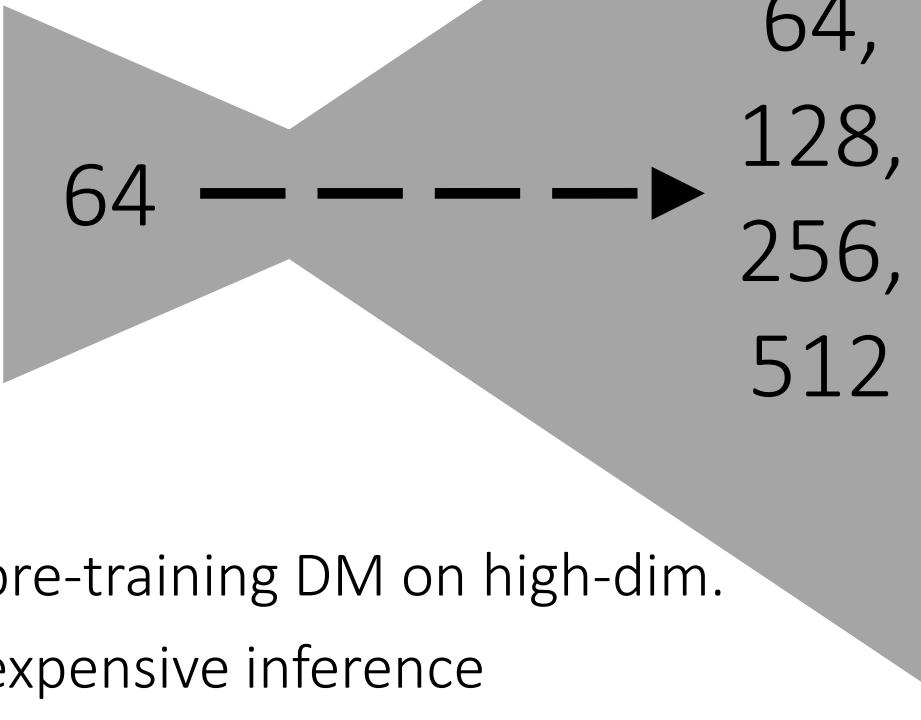
# SOTA 1-step “Res.-Agnostic” Generation



- ✓ Don't need pre-training DM on high-dim.
- ✓ Don't need expensive inference

# SOTA 1-step “Res.-Agnostic” Generation

Text-to-Image generation



- ✓ Don't need pre-training DM on high-dim.
- ✓ Don't need expensive inference



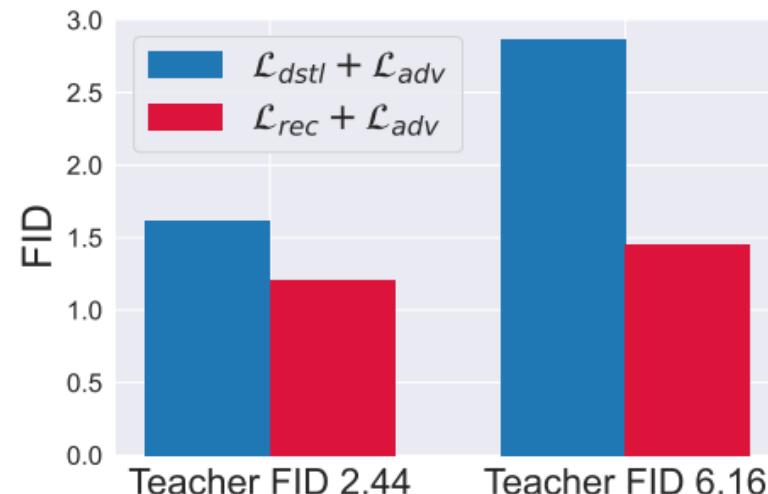
# Insight I. PaGoDA can learn the true $p_{\text{data}}$

$$\mathcal{L}_{\text{PaGoDA}}(G_{\theta}, D_{\psi}) := \mathcal{L}_{\text{rec}}(G_{\theta}) + \lambda \mathcal{L}_{\text{GAN}}(G_{\theta}, D_{\psi})$$

**Theorem 2 (Informal) [PaGoDA learns the true data]** Assume an optimal discriminator  $\psi^*$  in the above loss. Then with some conditions, the following holds:

$$p_{\text{PaGoDA}, \theta^*} = p_{\text{data}},$$

where  $p_{\text{PaGoDA}, \theta^*} := G_{\theta^*} \# p_{\text{prior}}$ . Not the case for conventional distillation-based method.



# Insight II. PaGoDA is not GAN

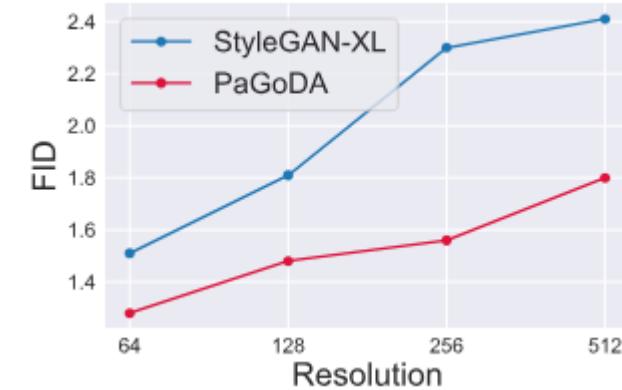
$$\mathcal{L}_{\text{PaGoDA}}(G_{\theta}, D_{\psi}) := \mathcal{L}_{\text{rec}}(G_{\theta}) + \lambda \mathcal{L}_{\text{GAN}}(G_{\theta}, D_{\psi})$$



(e) Upsampling of PaGoDA



(f) Upsampling of StyleGAN-XL



Theorem 3 (Informal) [PaGoDA is STABLE!] Let  $E$  be fixed ODE-based encoder. Assume at optimal generator  $\theta^*$  that

$$x = G_{\theta^*}(E(x)) \text{ on } x \in \text{supp}(p_{\text{data}})$$

- I. Then under some conditions, PaGoDA's alternative gradient descent of  $(\theta, \psi)$  is stable at  $(\theta^*, \psi^*)$  in Lyapunov sense.
- II. Potentially not the case for GAN.

# Making Diffusion Models More Powerful

# Training

# [ICML'23 Lai+] FP-Diffusion



# Accurate Training

# Scalability

# [NeurIPS'24 Kim&Lai+] PaGoDA



# Cheap Training and Sampling in High-Dim.

# Sampling

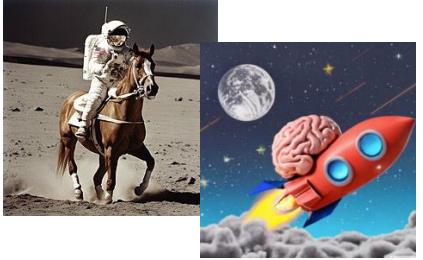
# [ICLR'24 Kim&Lai+] CTM



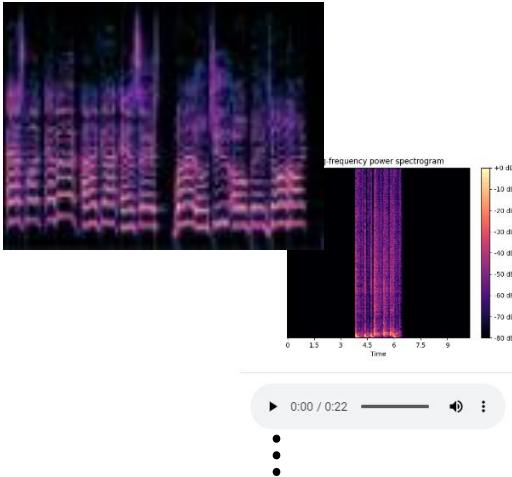
# More than that...

## Pre-trained Diffusion Models

### Image Diffusion Model



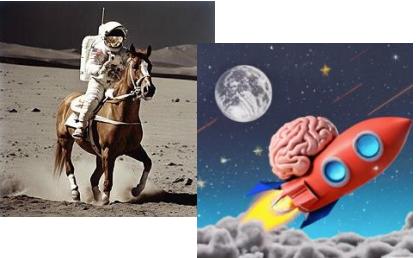
### Music/audio Diffusion Model



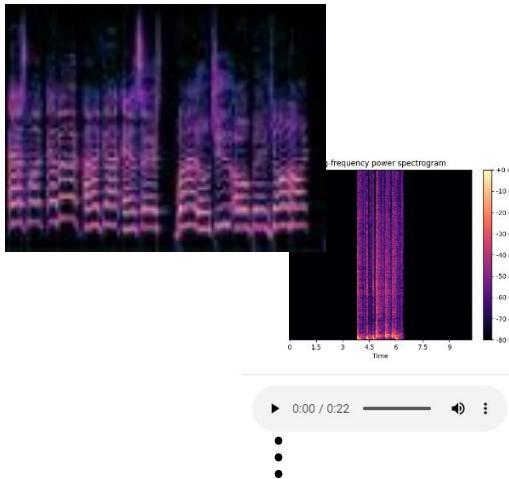
# More than that...

## Pre-trained Diffusion Models

### Image Diffusion Model



### Music/audio Diffusion Model



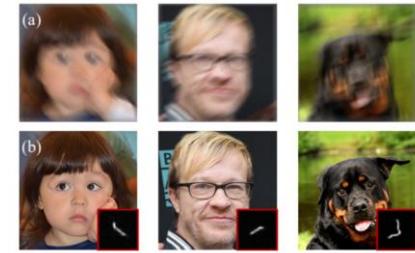
Pre-trained Diffusion  
for various applications

Without  
any extra  
training!!

ICML'23  
GibbsDDRM

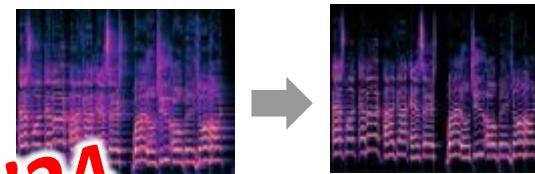
## Various applications

### Inverse Problems/Restoration



#### Image Restoration

- GibbsDDRM [ICML'23]

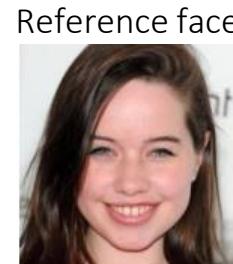


#### Audio Restoration

- DiffDereverb [ICASSP'23]
- Diffiner in [Interspeech'23]
- VRDMG in [ICASSP'24]

### Guided (Controllable) Generation

Prompt  
*"a headshot of a girl with blond hair and space background"*



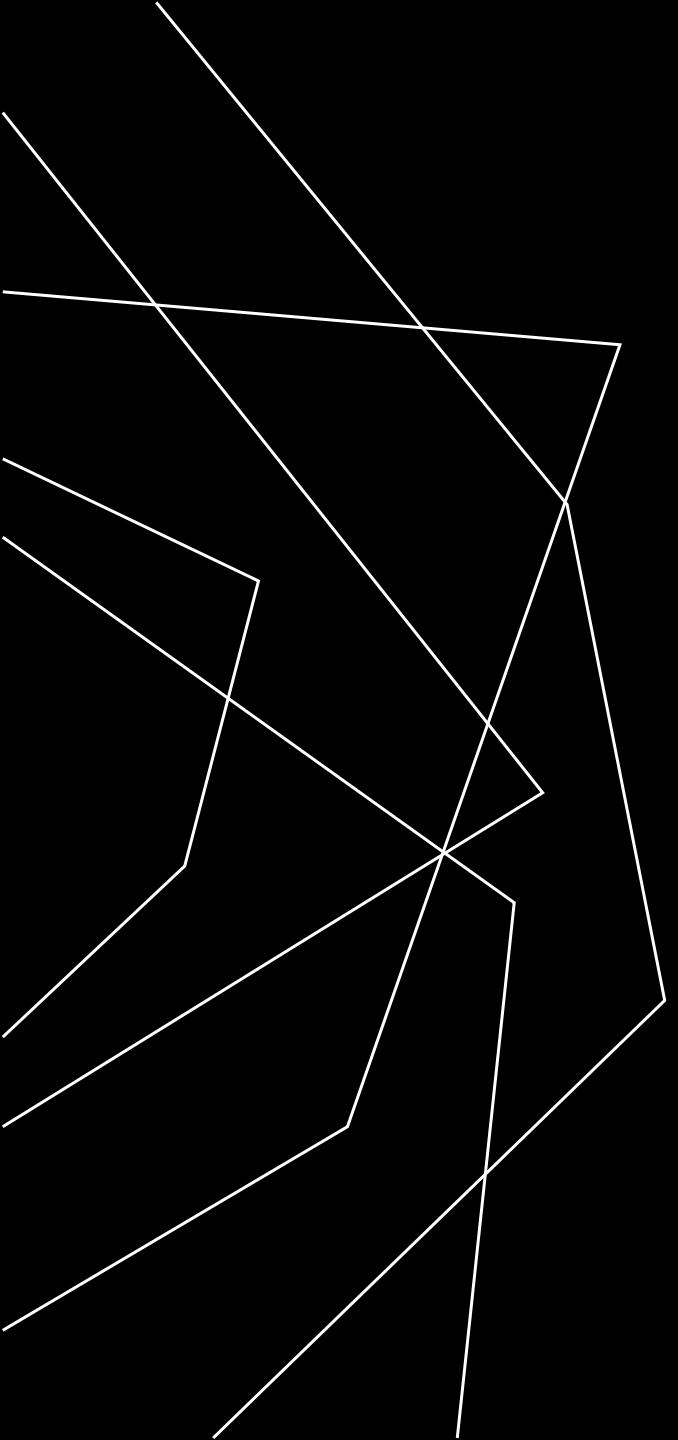
ICLR'24  
MPGD

# Missing Collaborator...

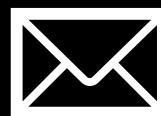


# MawMaw (貓貓)

...literally means Cat Cat



# THANK YOU!



Chieh-Hsin.Lai@sony.com



@JCJesseLai

