



FROM WHITE NOISE TO SYMPHONY : DIFFUSION MODELS FOR MUSIC AND SOUND

@ ISMIR 2024

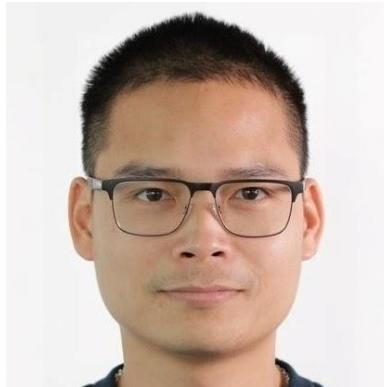
Sony AI



MEET OUR TEAM



Chieh-Hsin (Jesse) Lai



Bac Nguyen



Koichi Saito

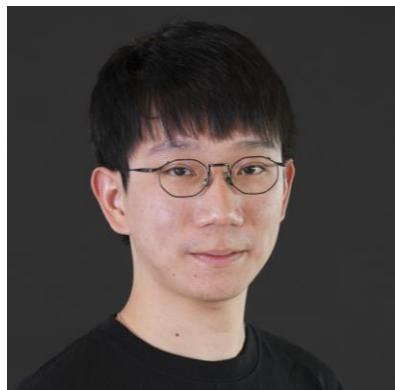


Yuki Mitsufuji

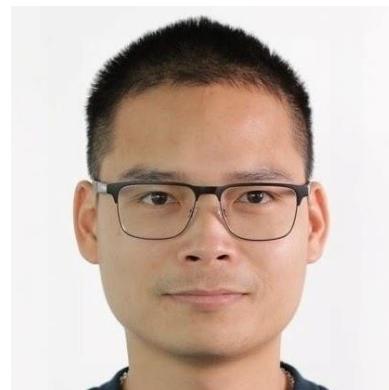


Stefano Ermon

MEET OUR SPEAKERS



Chieh-Hsin (Jesse) Lai



Bac Nguyen



Koichi Saito



Yuki Mitsufuji

AGENDA

Opening (9:00-9:10)

Session 1. (9:10-10:00) Foundation of Score-Based Diffusion Model

Session 2. (10:00-10:40) Practical Implementation and Application of Diffusion Model

Break (10:40-11:10)

Session 3. (11:10-12:00) Demo on Controllable Music Generation

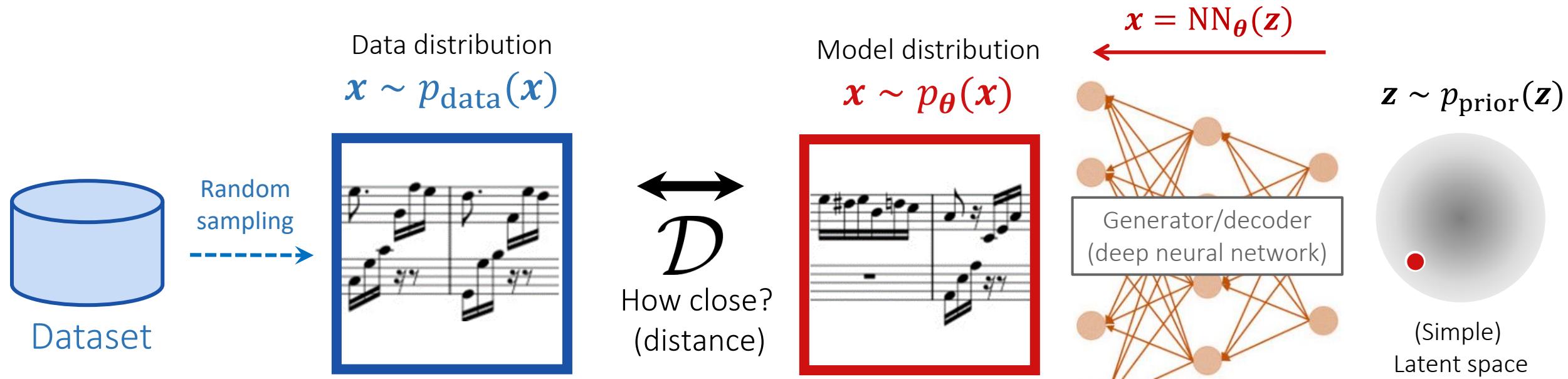
Session 4. (12:00-12:30) Future Trend

Disclaimer

For illustration, we occasionally use images.
But you can think of them as music data!

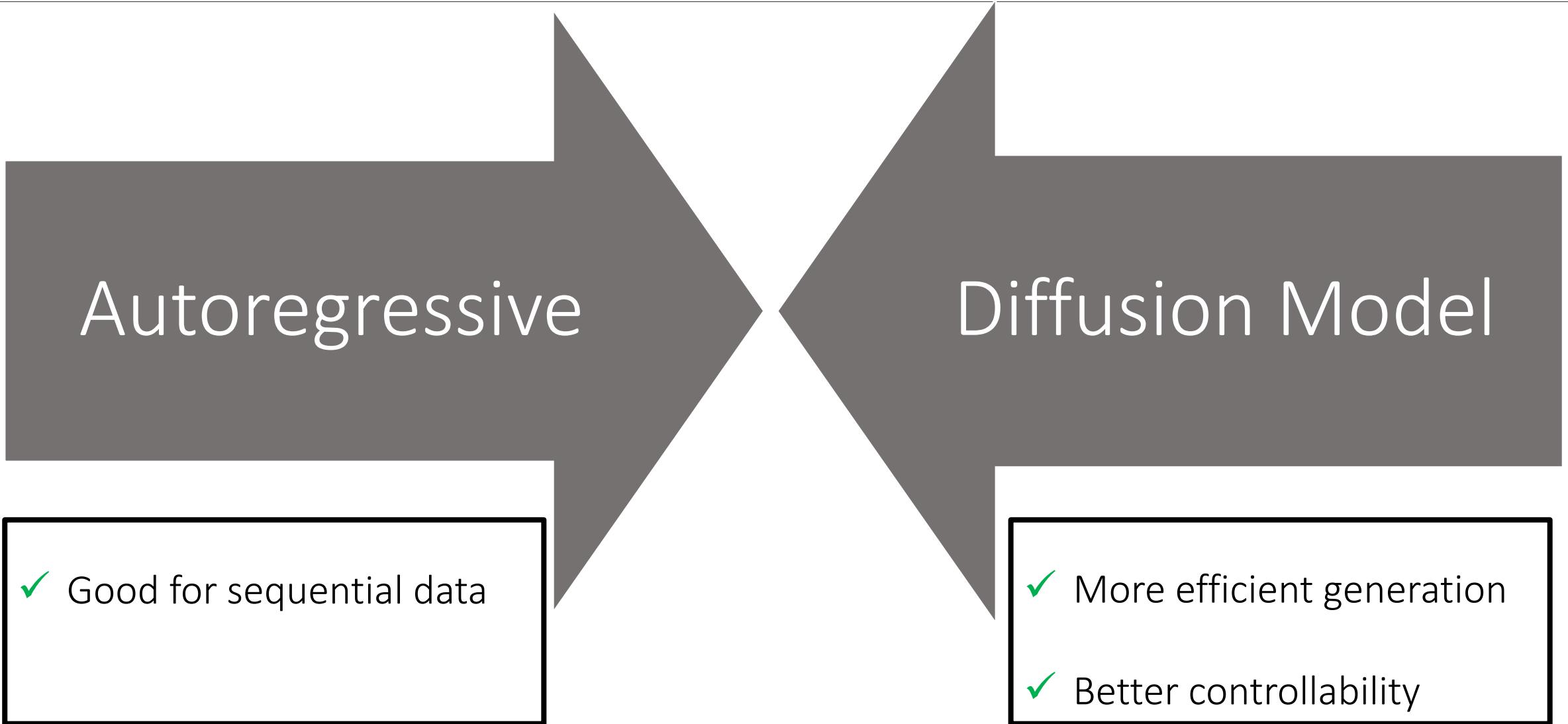
Deep Generative Model

Goal: use NN to learn target distribution explicitly/implicitly

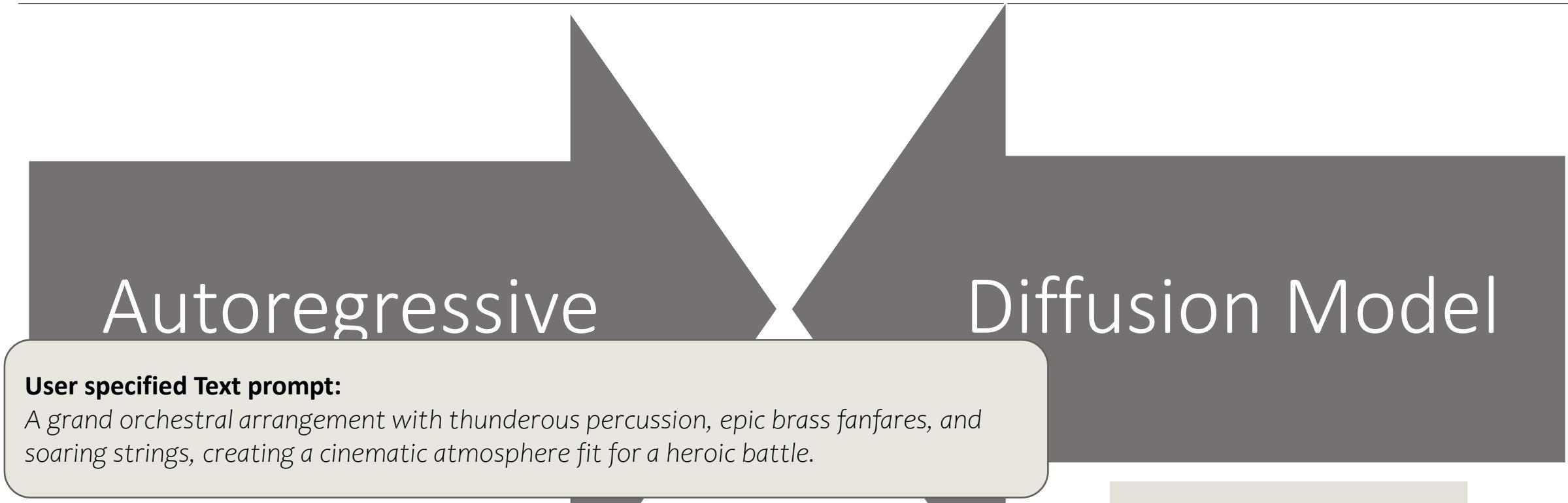


$$\min_{\theta} \mathcal{D}(p_{\theta}(x), p_{\text{data}}(x))$$

DGM in Music Generation



DGM in Music Generation



[Google Deepmind '16] WaveNet

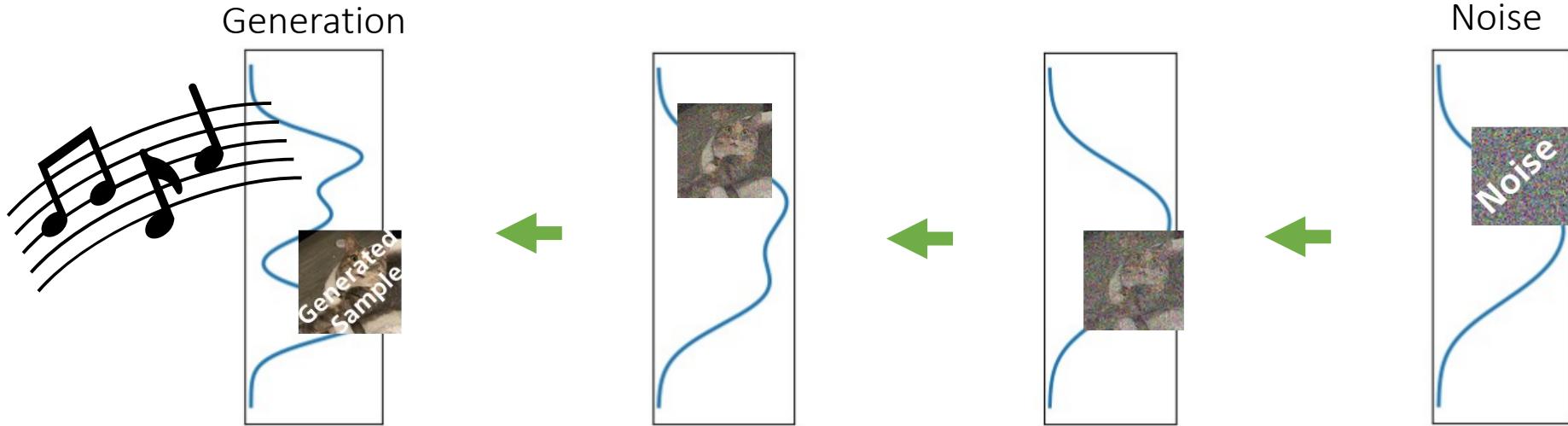


[Meta '23] MusicGen

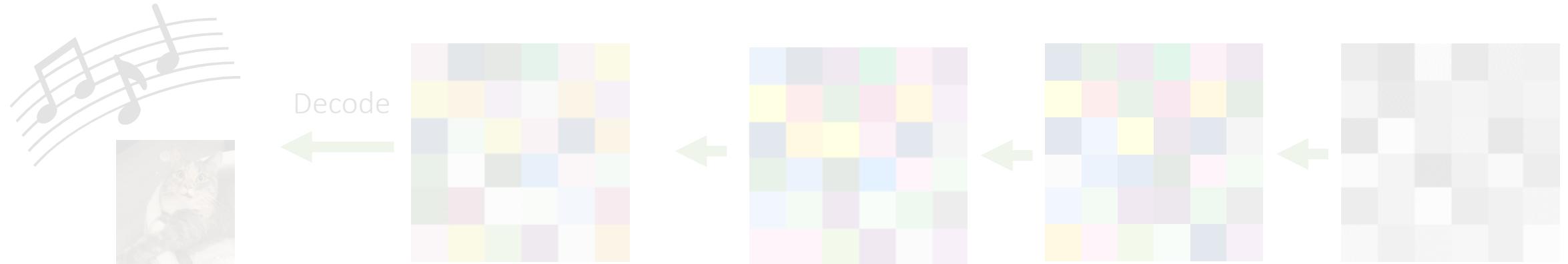


Regime of Continuous/Discrete Diffusion Model

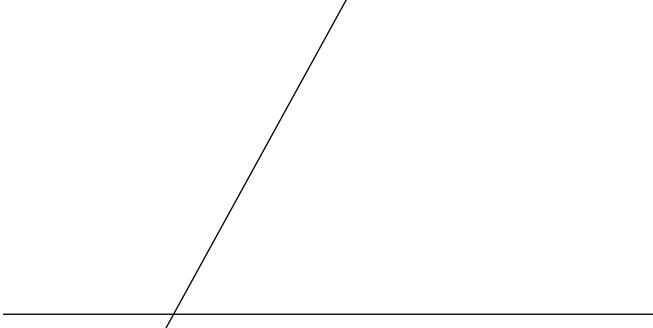
Continuous Data Diffusion Model:



Discrete Data Diffusion Model:

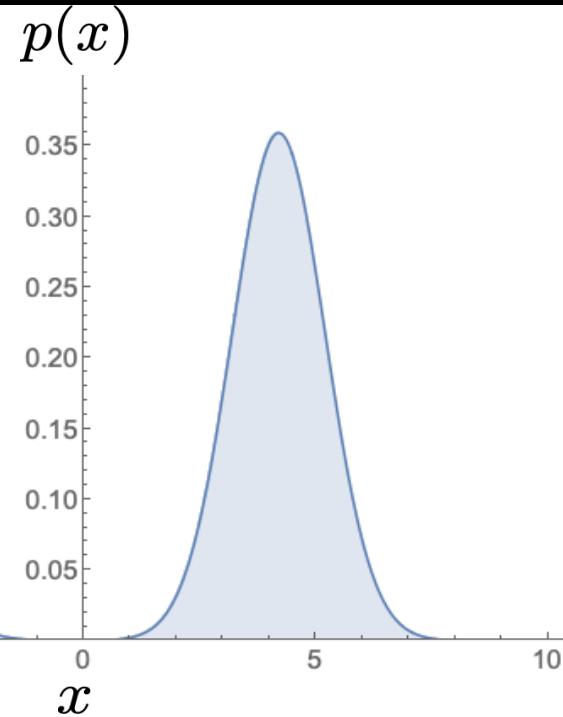


DEVELOPMENT OF SCORE-BASED DIFFUSION MODEL



Using “Score” (Slope) for Generation

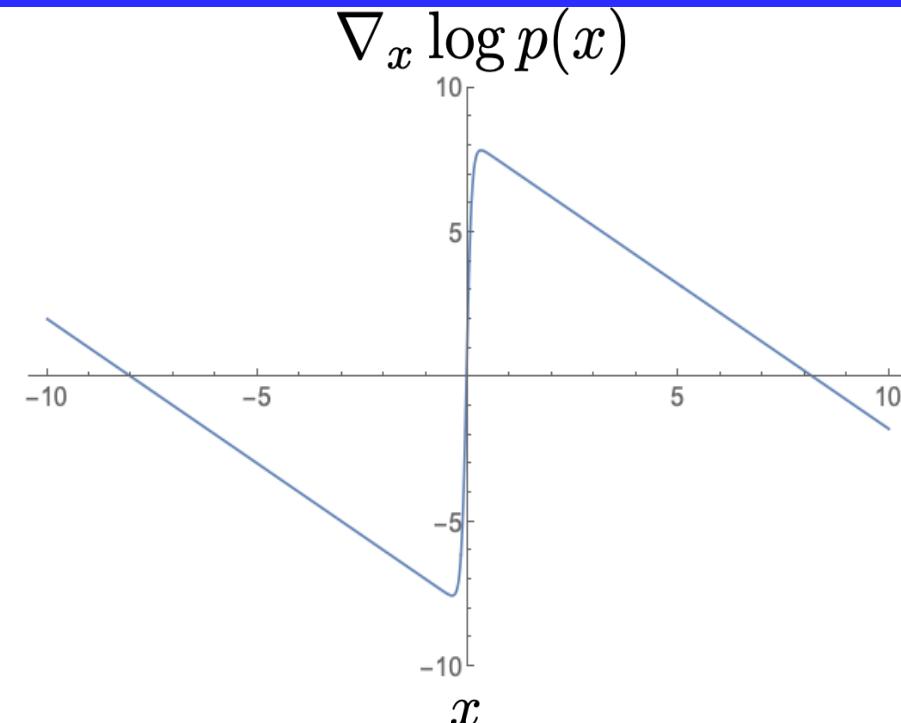
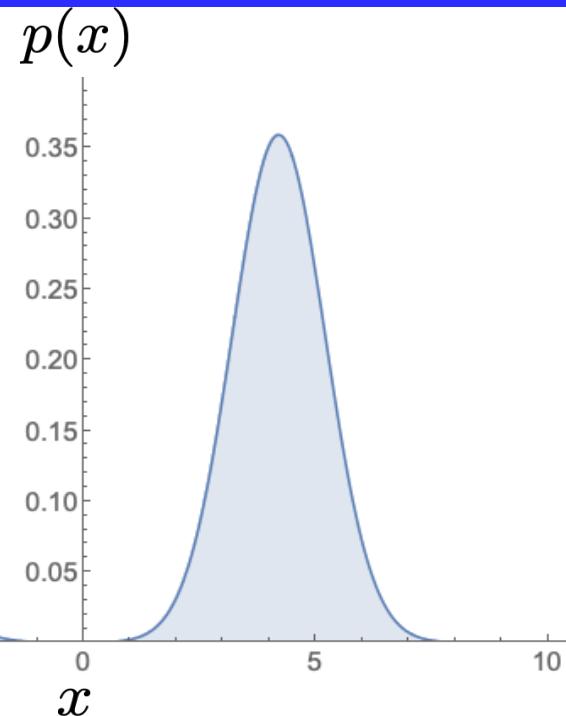
$p_{\text{data}}(x)$



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:

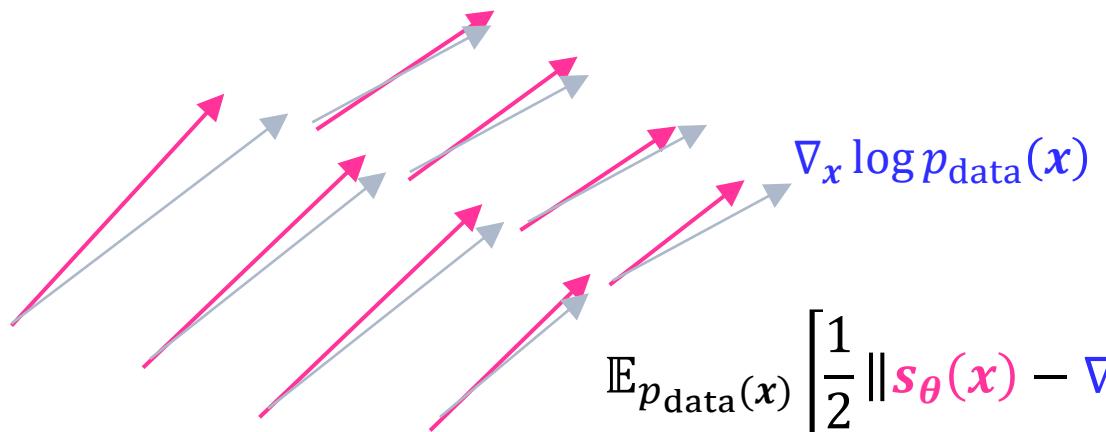
$$\nabla_x \log p_{\text{data}}(\mathbf{x})$$

[Generation] Langevin sampling:

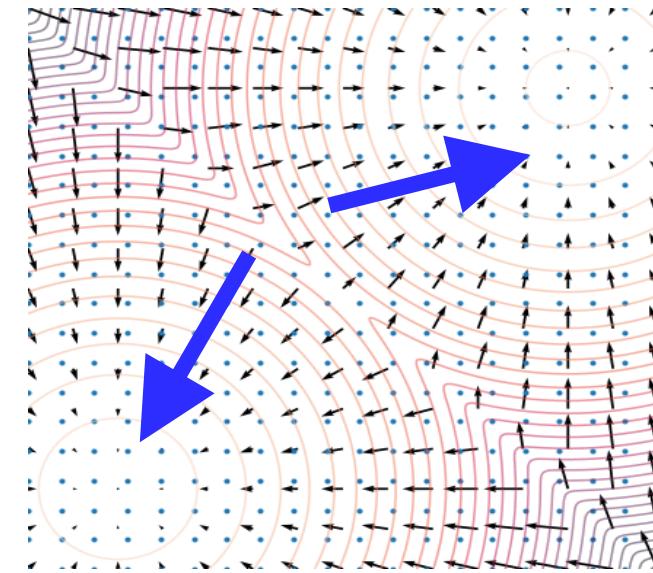
$$\mathbf{x}_{t+1} \leftarrow \mathbf{x}_t + \eta \nabla_x \log p_{\text{data}}(\mathbf{x}) + \sqrt{2} \boldsymbol{\epsilon}_t, \boldsymbol{\epsilon}_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$$

[Training] Score Matching:

$$\mathfrak{s}_{\theta}(\mathbf{x})$$



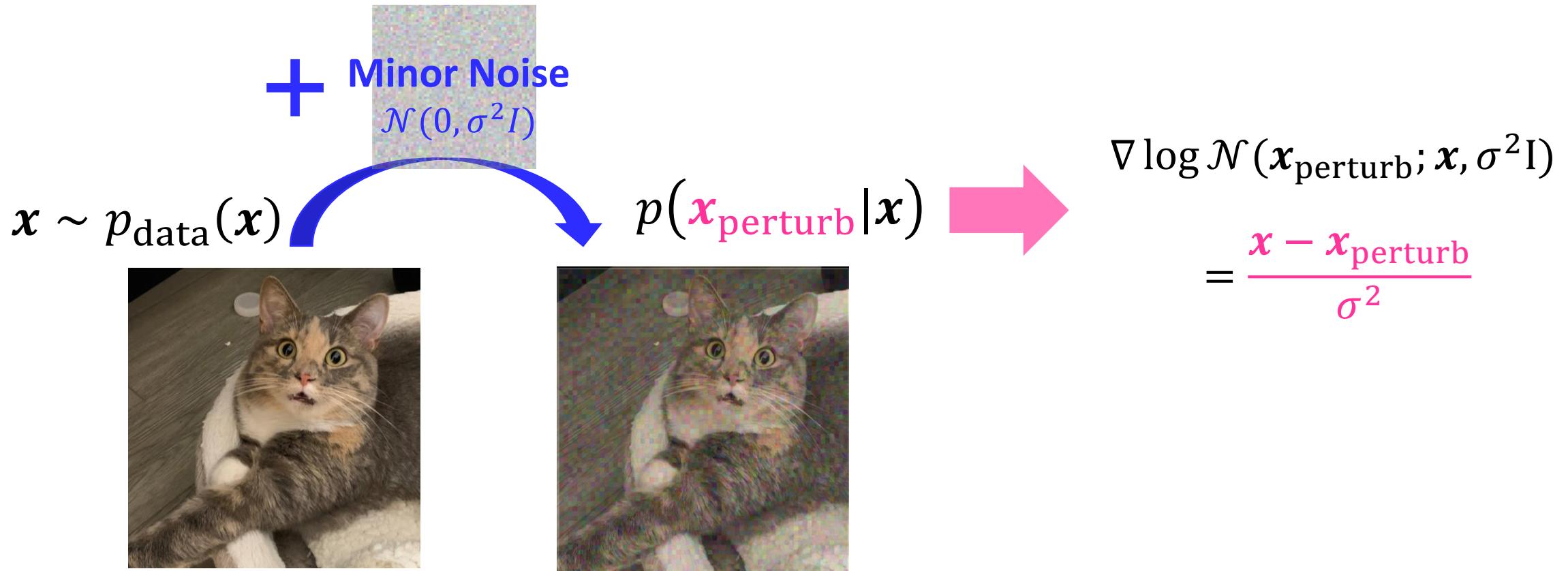
$$\mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\frac{1}{2} \| \mathfrak{s}_{\theta}(\mathbf{x}) - \nabla_x \log p_{\text{data}}(\mathbf{x}) \|_2^2 \right] = \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \left[\frac{1}{2} \| \mathfrak{s}_{\theta}(\mathbf{x}) \|_2^2 + \text{tr}(\nabla_{\mathbf{x}} \mathfrak{s}_{\theta}(\mathbf{x}))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



Phase II (2011). Denoising Score Matching

[Training] Denoising Score Matching

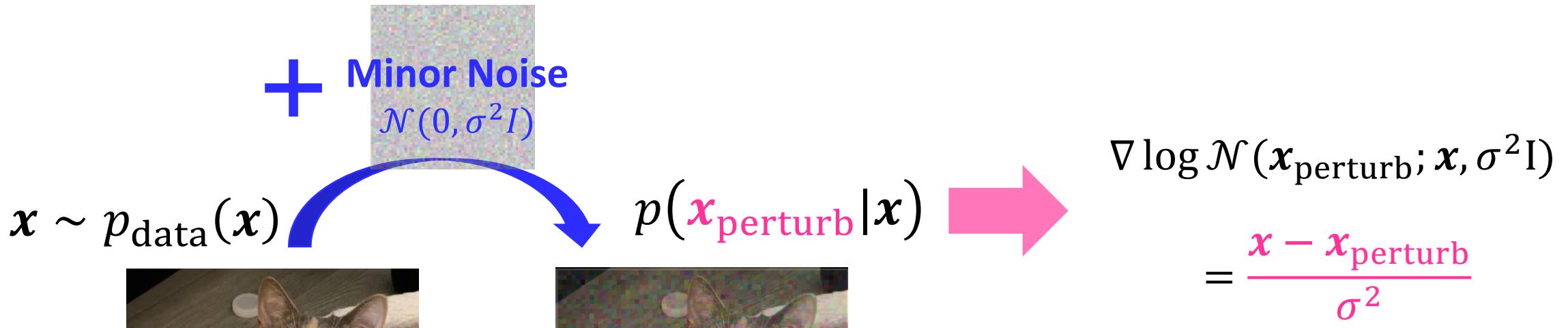
$$\mathbb{E}_{p(x_{\text{perturb}}|x)} \left[\frac{1}{2} \left\| s_{\theta}(x) - \left(\frac{x - x_{\text{perturb}}}{\sigma^2} \right) \right\|_2^2 \right]$$



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

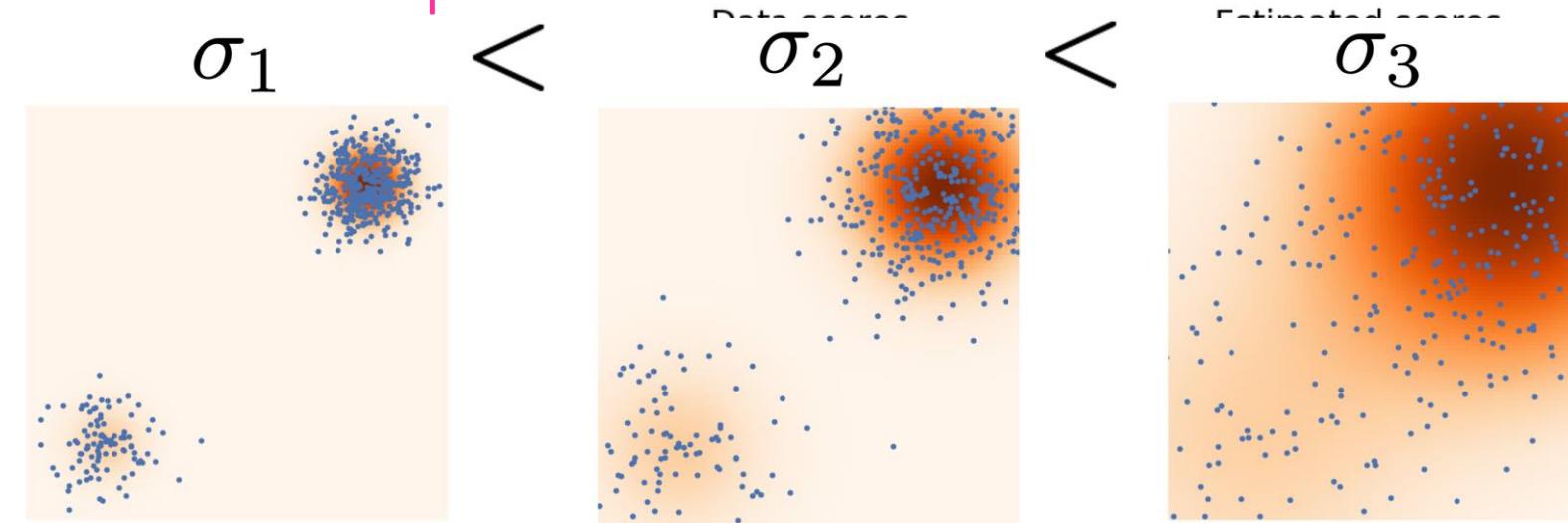
$$x_{t+1} \leftarrow x_t + \eta \nabla \log p_{\text{perturb}}(x_{\text{perturb}}|x) + \sqrt{2} \epsilon_t, \epsilon_t \sim \mathcal{N}(\mathbf{0}, I)$$

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

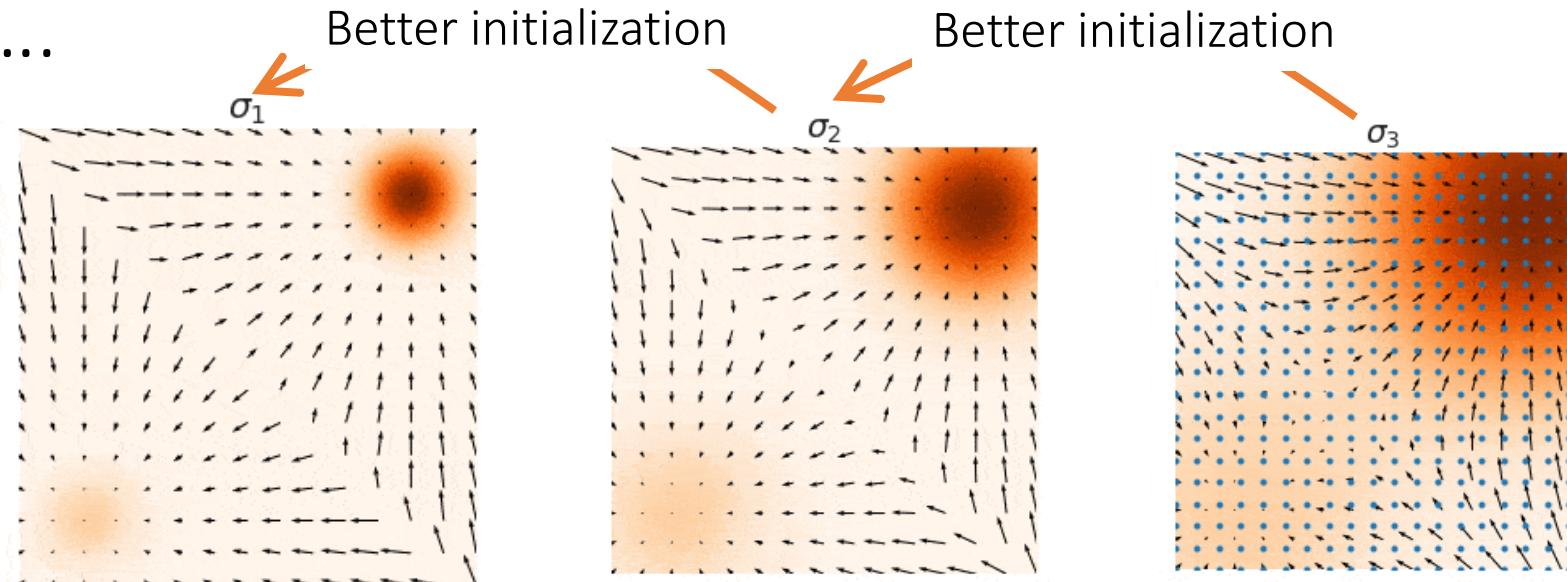


Phase III (2019). Noised Conditioned Denoising Score Matching

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



In generation....



Phase III (2019). Noised Conditioned Denoising Score Matching

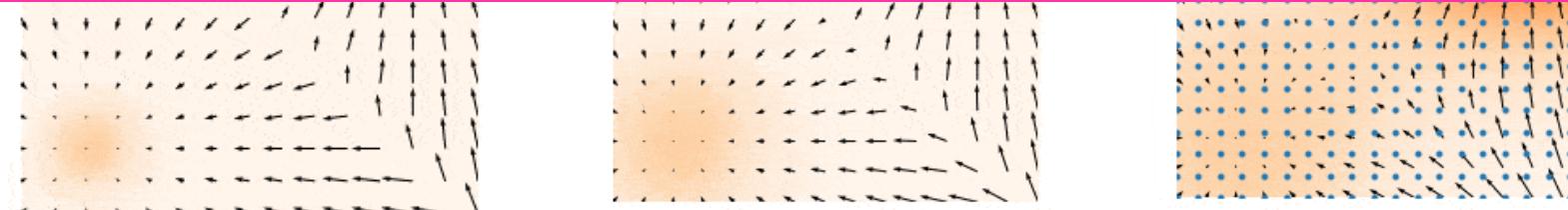
Injecting incremental multiple levels of noises to data

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



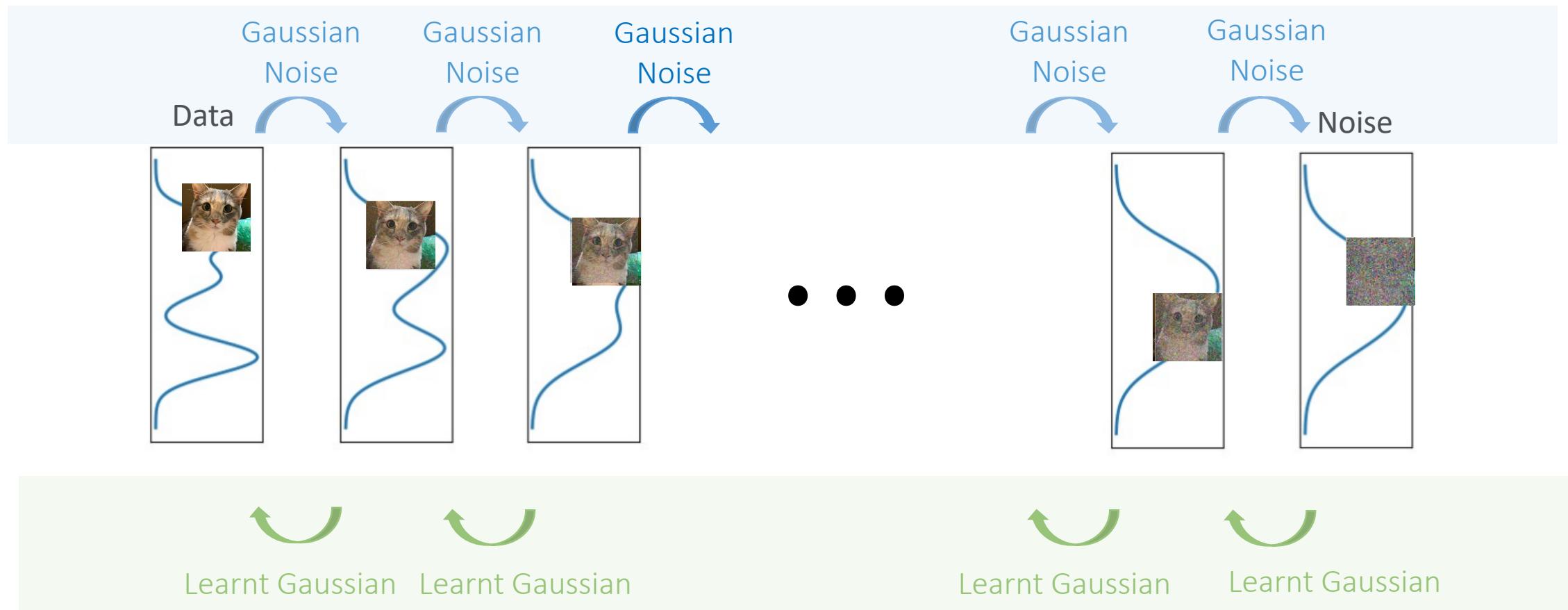
$$\sum_{i=1}^3 \omega_i \mathbb{E}_{p(x_{\text{perturb}, \sigma_i} | x)} \left[\frac{1}{2} \left\| s_{\theta}(x, \sigma_i) - \left(\frac{x - x_{\text{perturb}, \sigma_i}}{\sigma_i^2} \right) \right\|_2^2 \right]$$

The vibe of Diffusion Model appeared...



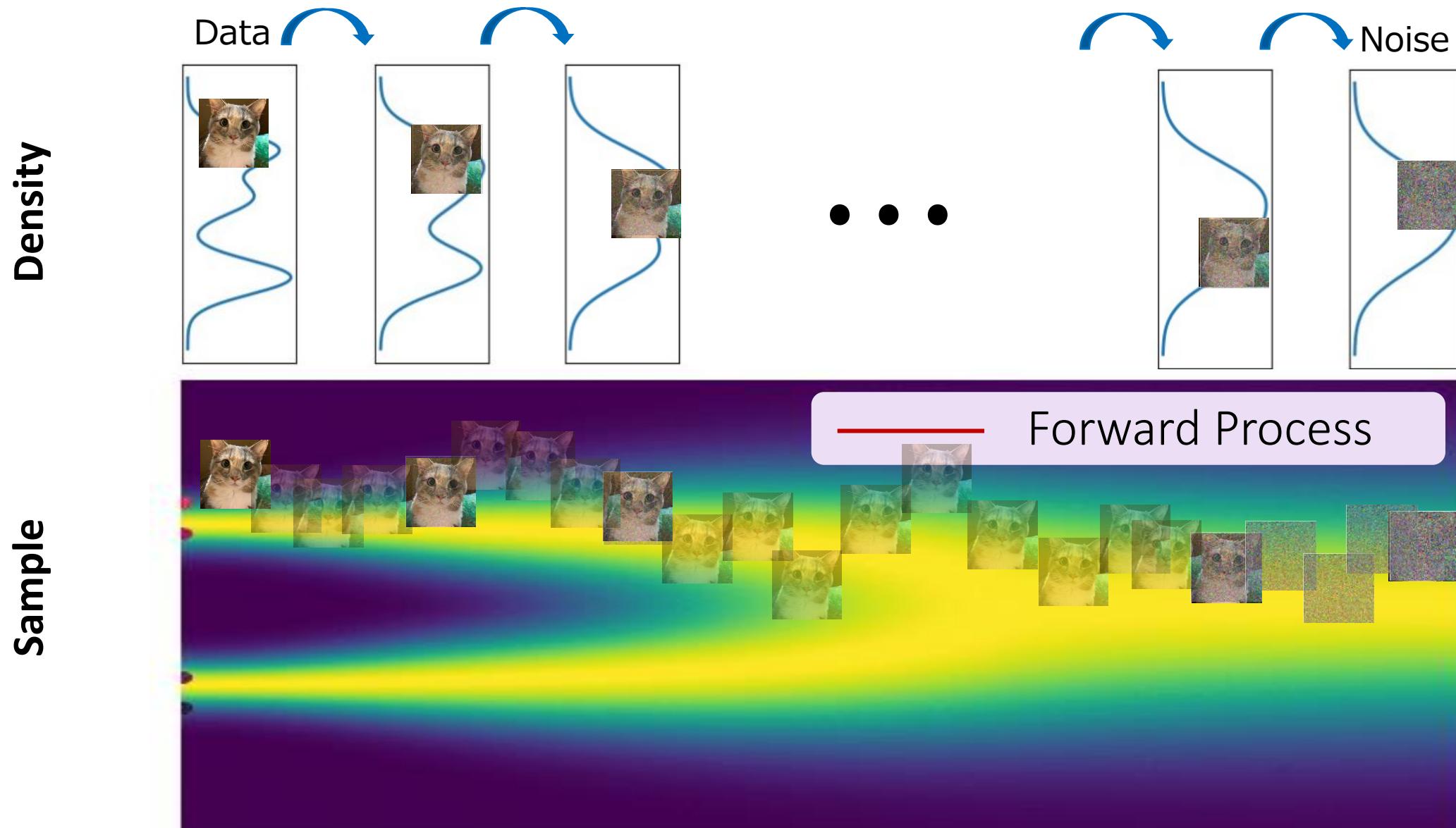
Concurrently, Denoising Diffusion Probabilistic Model (DDPM)

- Different forward noise schedulers



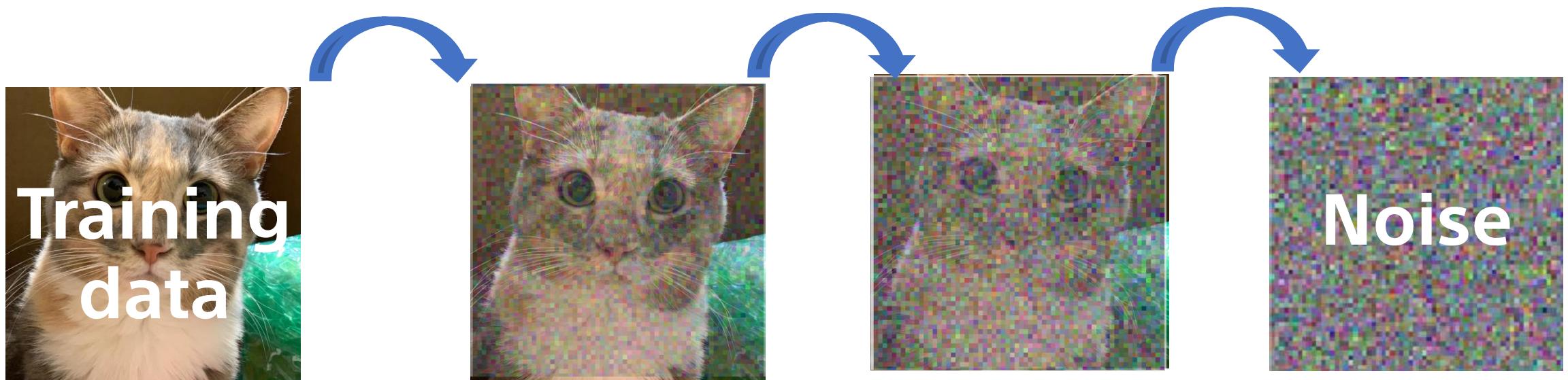
- Bayesian viewpoint for reverse process – learn reversal Gaussian

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



Diffusion Model at High Level

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



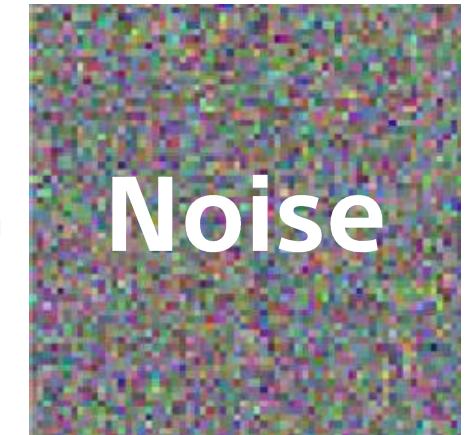
Diffusion Model at High Level

Backward Process: Learn NN to gradually denoise

Generated new image



Gradually transforming noise to a new image

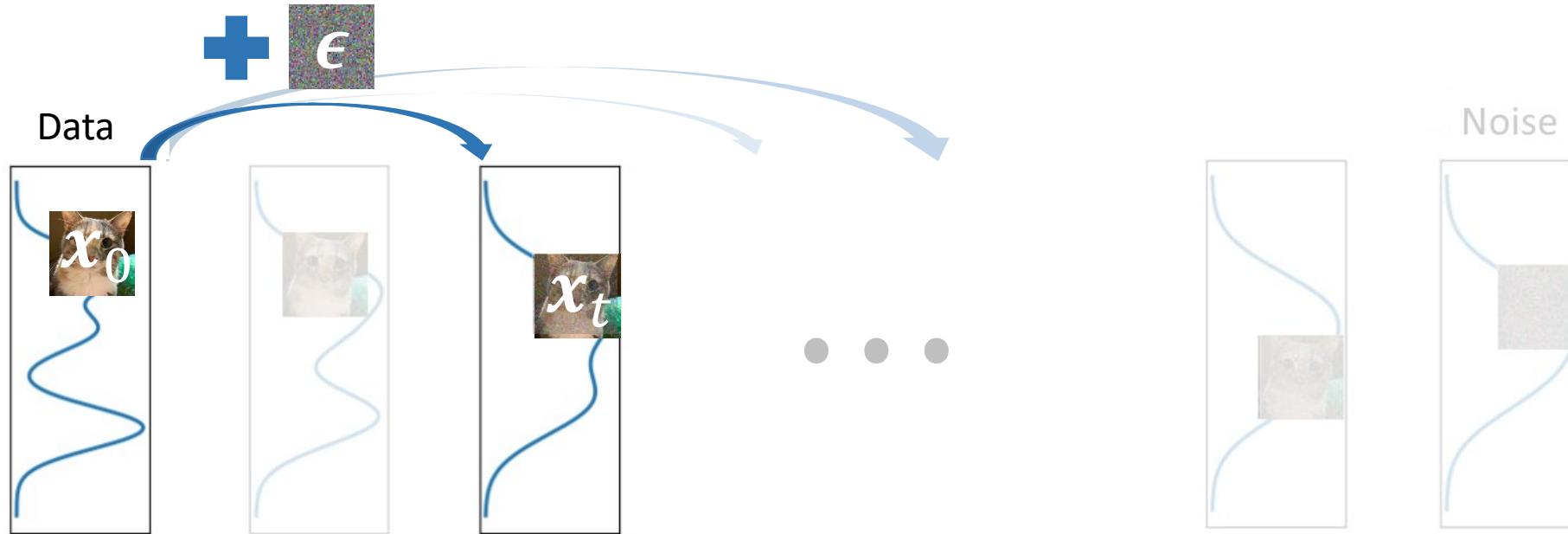


Noise

Diffusion Model's Forward Process

Forward Process:

- Data \rightarrow Noise (Common simple space)
- With user specified noise schedulers a_t and b_t

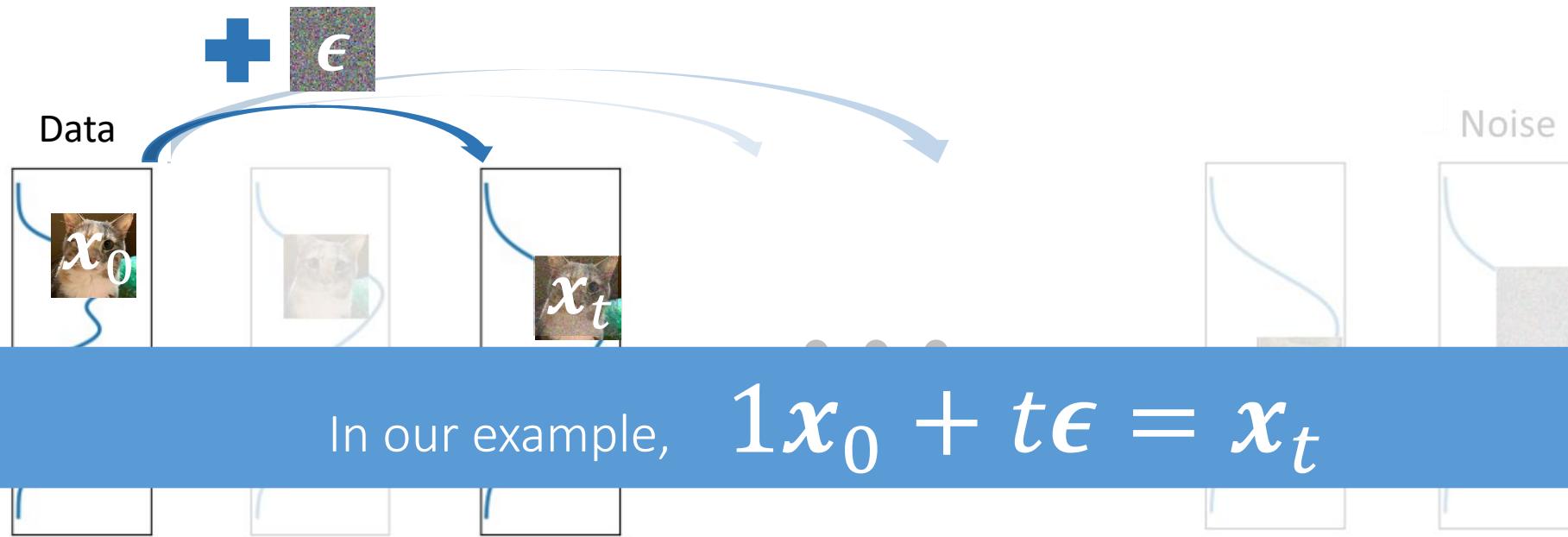


$$a_t x_0 + b_t \epsilon = x_t$$

Diffusion Model's Forward Process

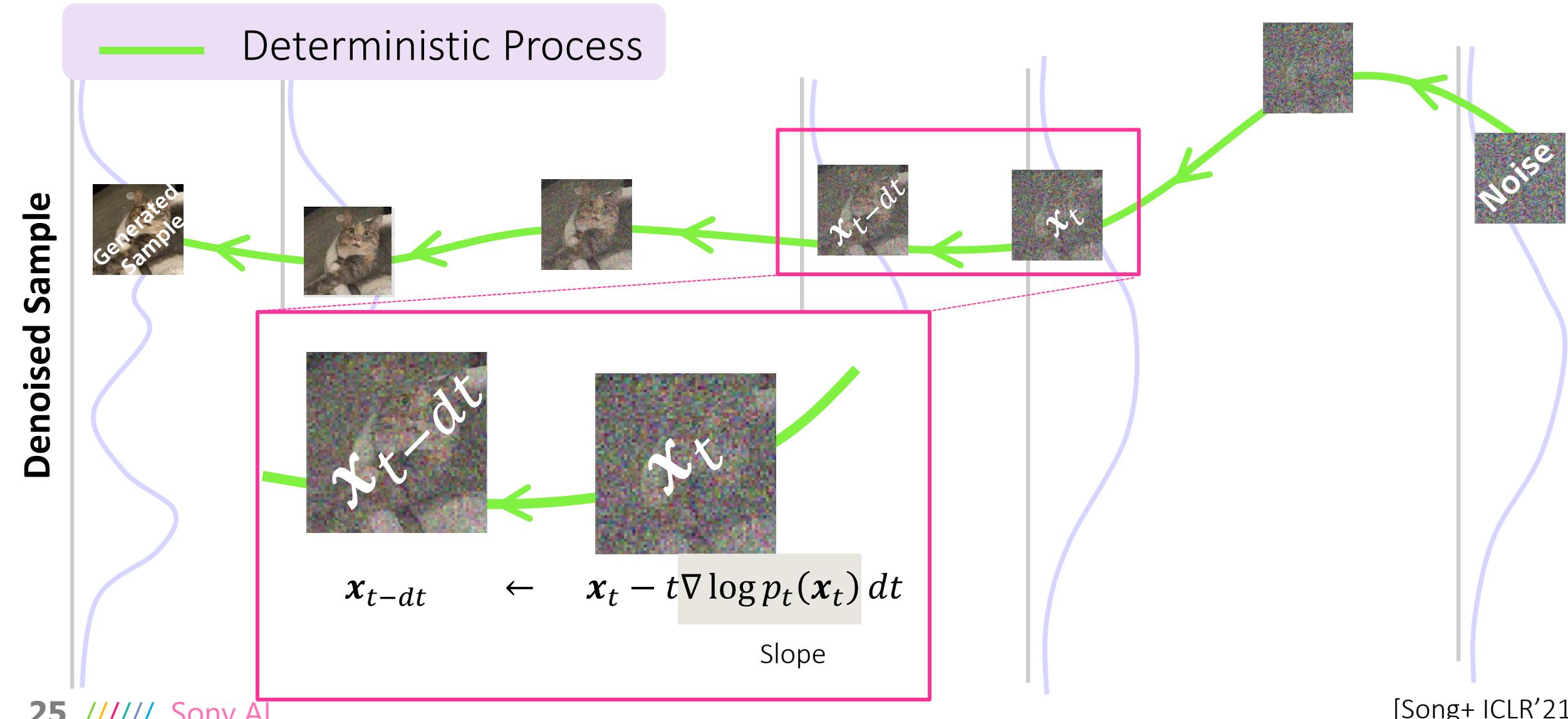
Forward Process:

- Data \rightarrow Noise (Common simple space)
- With user specified noise schedulers a_t and b_t

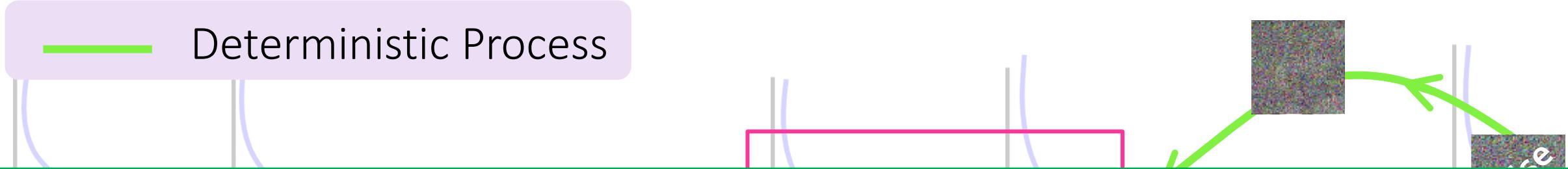


$$a_t x_0 + b_t \epsilon = x_t$$

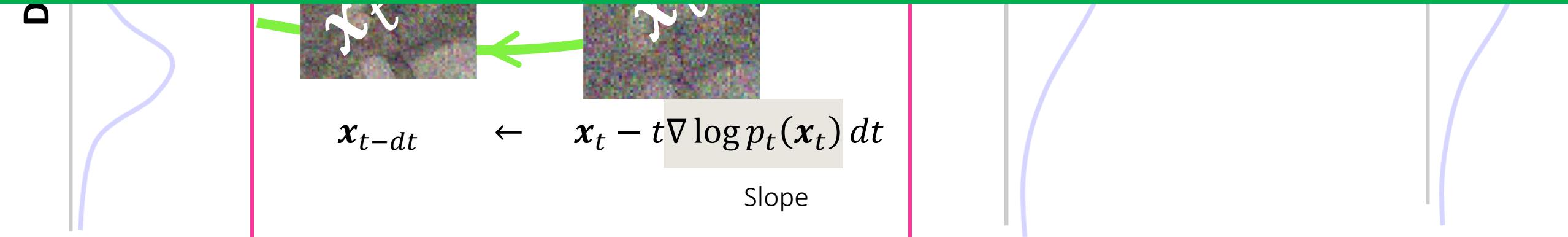
Diffusion Model's Backward Deterministic Process (Generation)



Diffusion Model's Backward Deterministic Process (Generation)



$$x_{t-dt} - x_t = -t \nabla \log p_t(x_t) dt$$



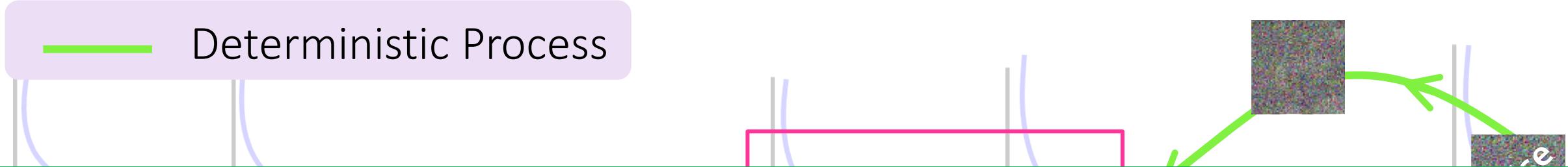
Diffusion Model's Backward Deterministic Process (Generation)



$$\frac{x_{t-dt} - x_t}{dt} = -t \nabla \log p_t(x_t)$$



Diffusion Model's Backward Deterministic Process (Generation)



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

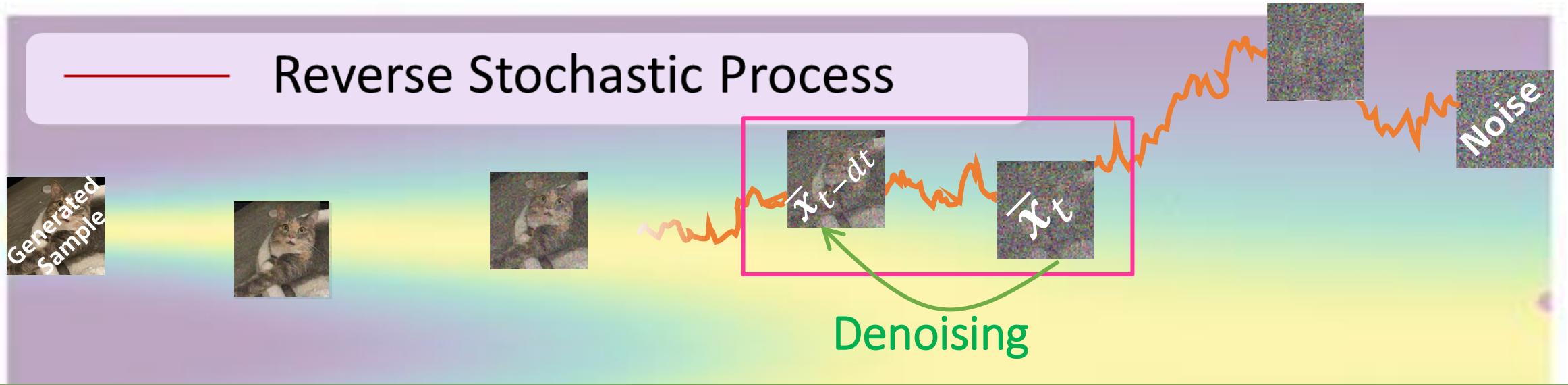
$$\frac{dx_t}{dt} = -t \nabla \log p_t(x_t)$$

$$x_{t-dt} \leftarrow x_t - t \nabla \log p_t(x_t) dt$$

Slope

Diffusion Model's Backward Stochastic Process (Generation)

Users predefined noise schedulers: $f(x_t, t) = 0$ & $g(t) = \sqrt{2t}$

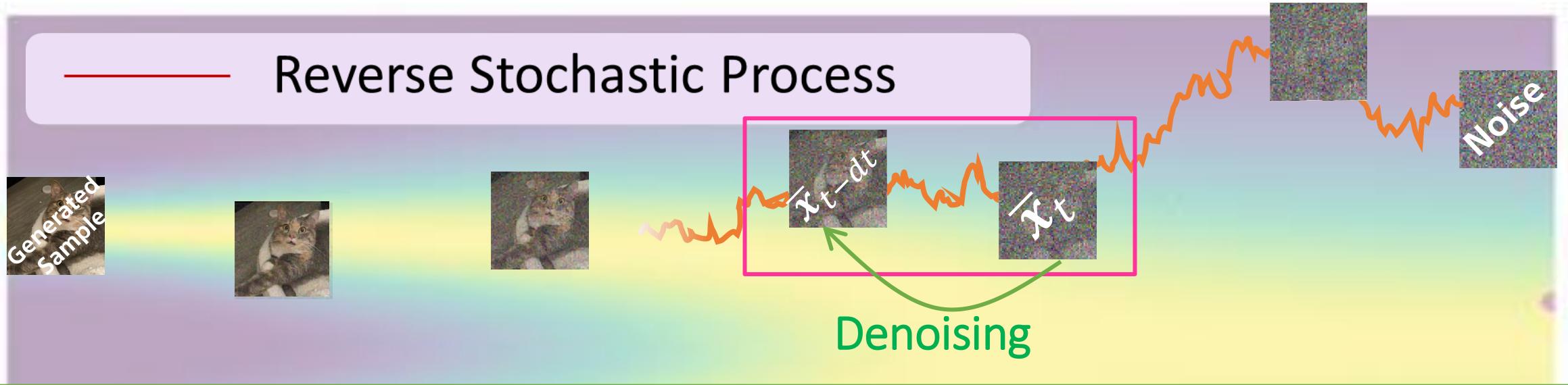


Getting finer timesteps ($\Delta t \rightarrow 0$) leads to

$$d\bar{x}_t = [f(\bar{x}_t, t) - g^2(t)\nabla \log p_t(\bar{x}_t)]d\bar{t} + g(t)d\bar{w}_t$$

Diffusion Model's Backward Stochastic Process (Generation)

Users predefined noise schedulers: $f(x_t, t) = 0$ & $g(t) = \sqrt{2t}$

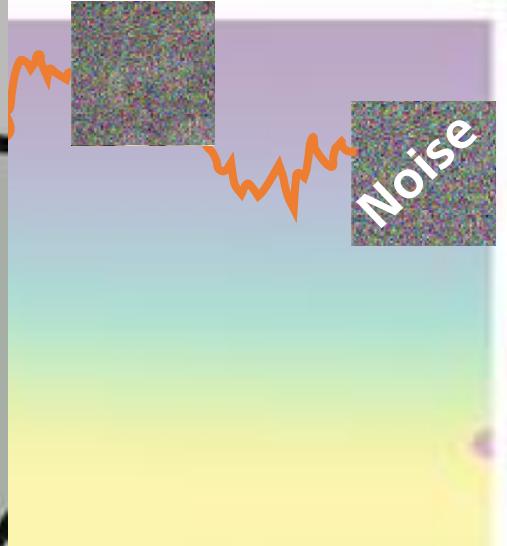
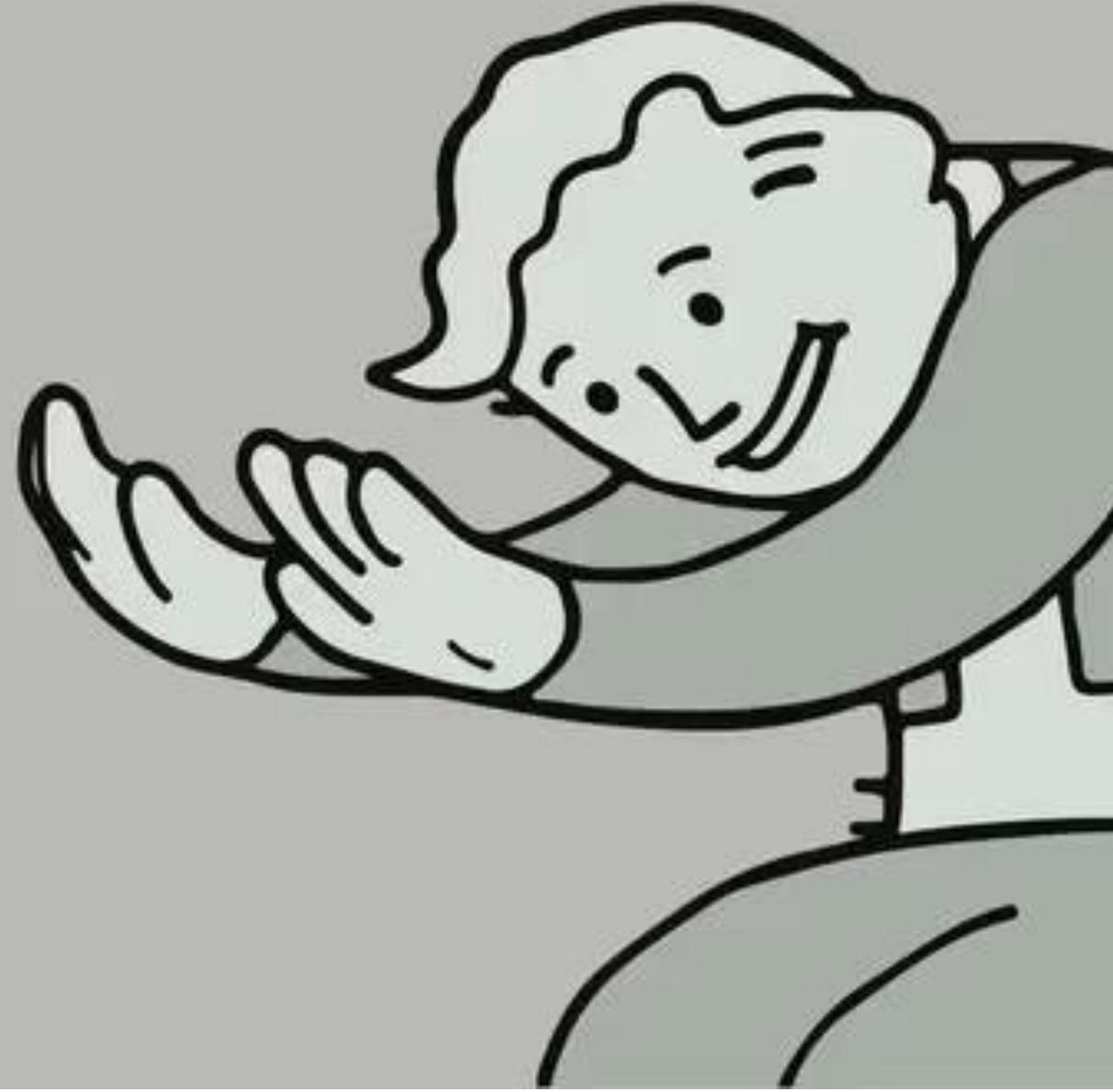
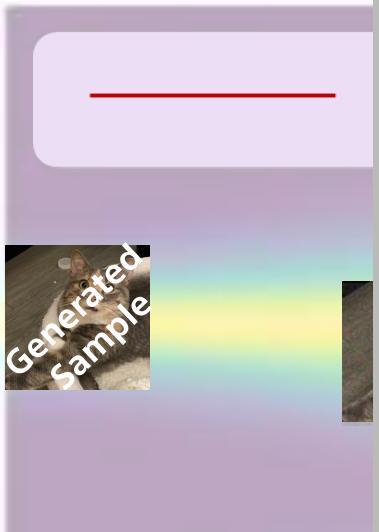


In our example,

$$d\bar{x}_t = [-2t\nabla \log p_t(\bar{x}_t)]d\bar{t} + \sqrt{2t}d\bar{w}_t$$

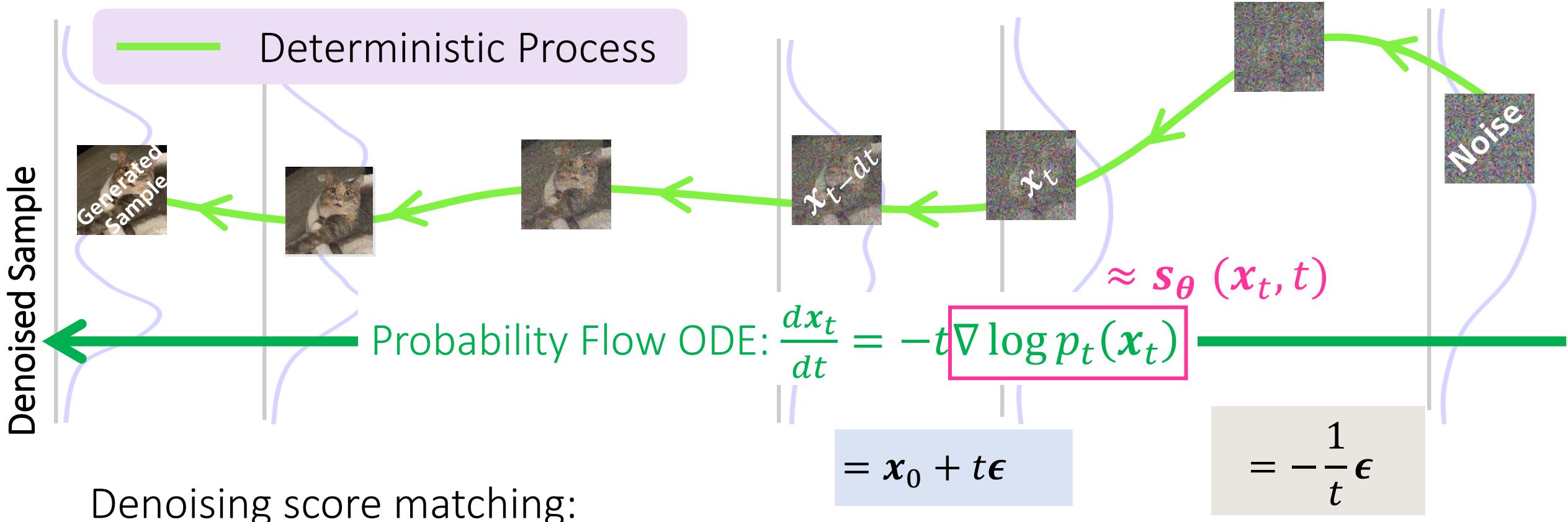
Diffusion Model

Users predefined



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 $a_t = 1$ & $b_t = t$

Denoising score matching:

$$= -\frac{1}{t} \epsilon$$

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

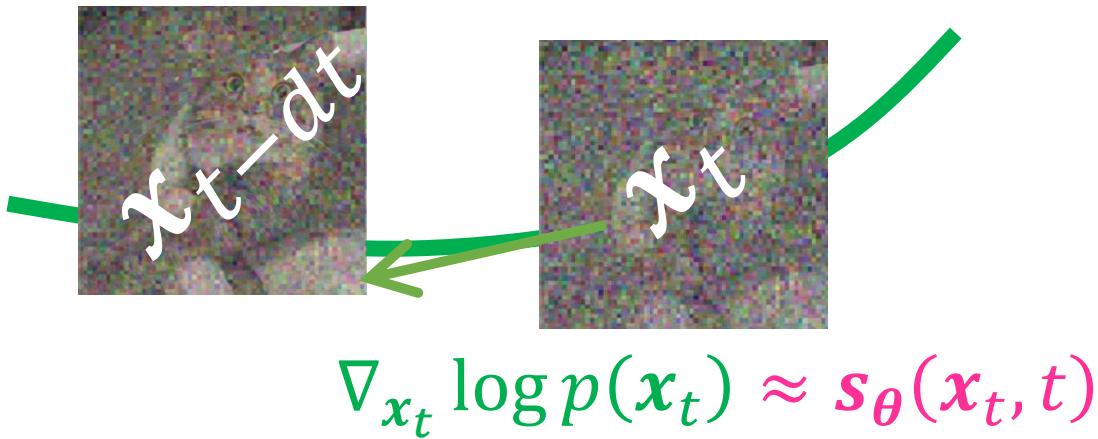
Training DM

Initialize: $s_{\theta}(x_t, t)$

- Get data x_0
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_{\theta} \| s_{\theta}(x_t, t) - \left(-\frac{1}{t} \epsilon \right) \|^2$$

Score-Based Diffusion Model's Training and Sampling



Training DM

Initialize: $s_\theta(x_t, t)$

- Get data \mathbf{x}_0
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- Perturb $\mathbf{x}_t = \mathbf{x}_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_\theta \|s_\theta(x_t, t) - \left(-\frac{1}{t}\epsilon\right)\|^2$$

For instance, Euler method

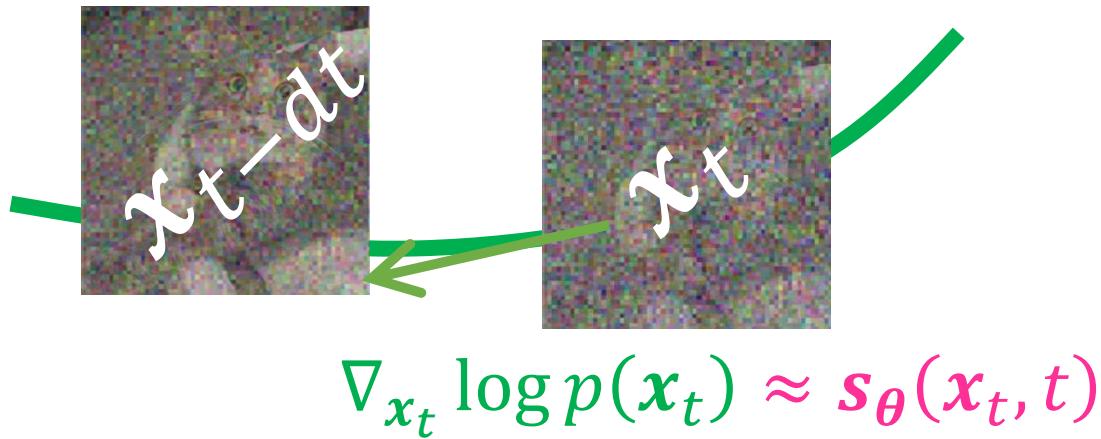
Starting from $\mathbf{x}_T \sim \mathcal{N}(\mathbf{0}, T^2 \mathbf{I})$, any ODE solver can be used to solve

$$\frac{d\mathbf{x}_t}{dt} = -t s_\theta(\mathbf{x}_t, t),$$

by calling a library, such as

```
scipy.integrate.odeint
```

Score-Based Diffusion Model's Training and Sampling



Training DM

Initialize: $s_\theta(x_t, t)$

- Get data x_0
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_\theta \|s_\theta(x_t, t) - \left(-\frac{1}{t}\epsilon\right)\|^2$$

Sampling with DM

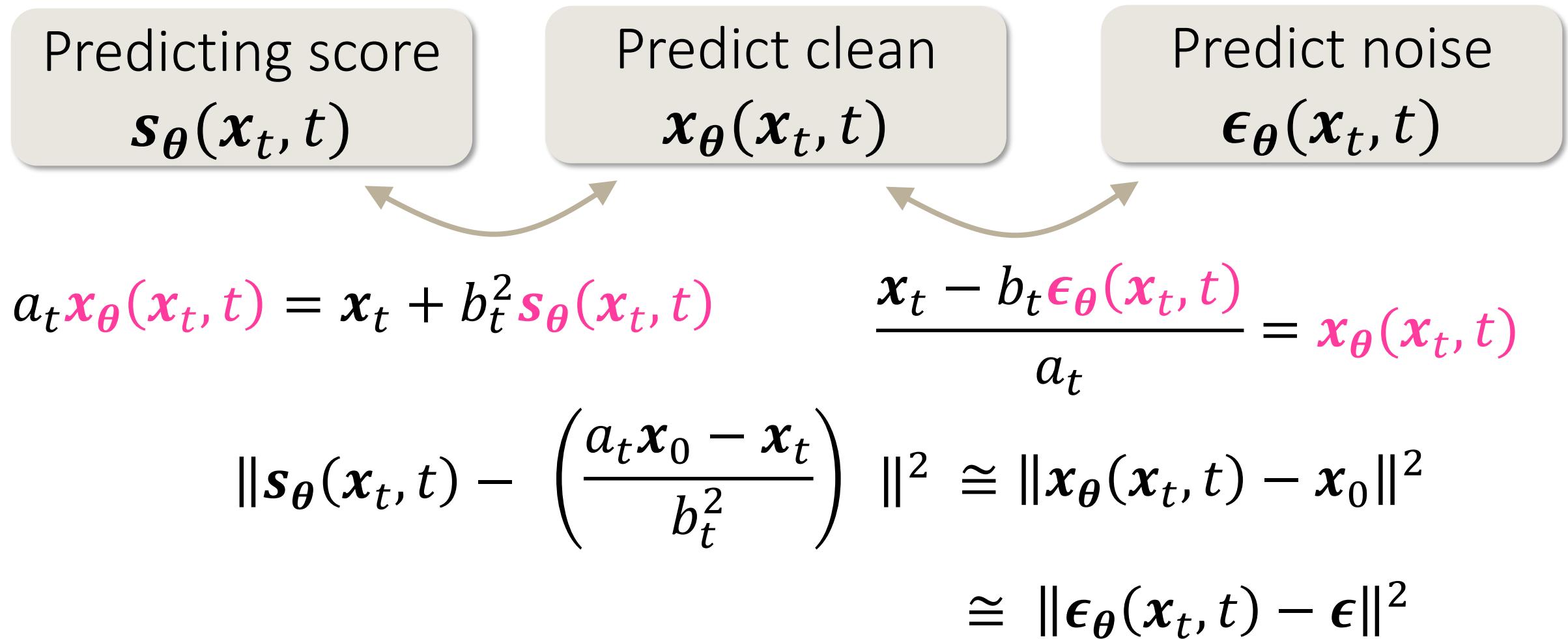
Request: $T = t_{N-1} > \dots > t_1 > t_0 = 0$

Initialize: $x_T \sim \mathcal{N}(\mathbf{0}, T^2 I)$

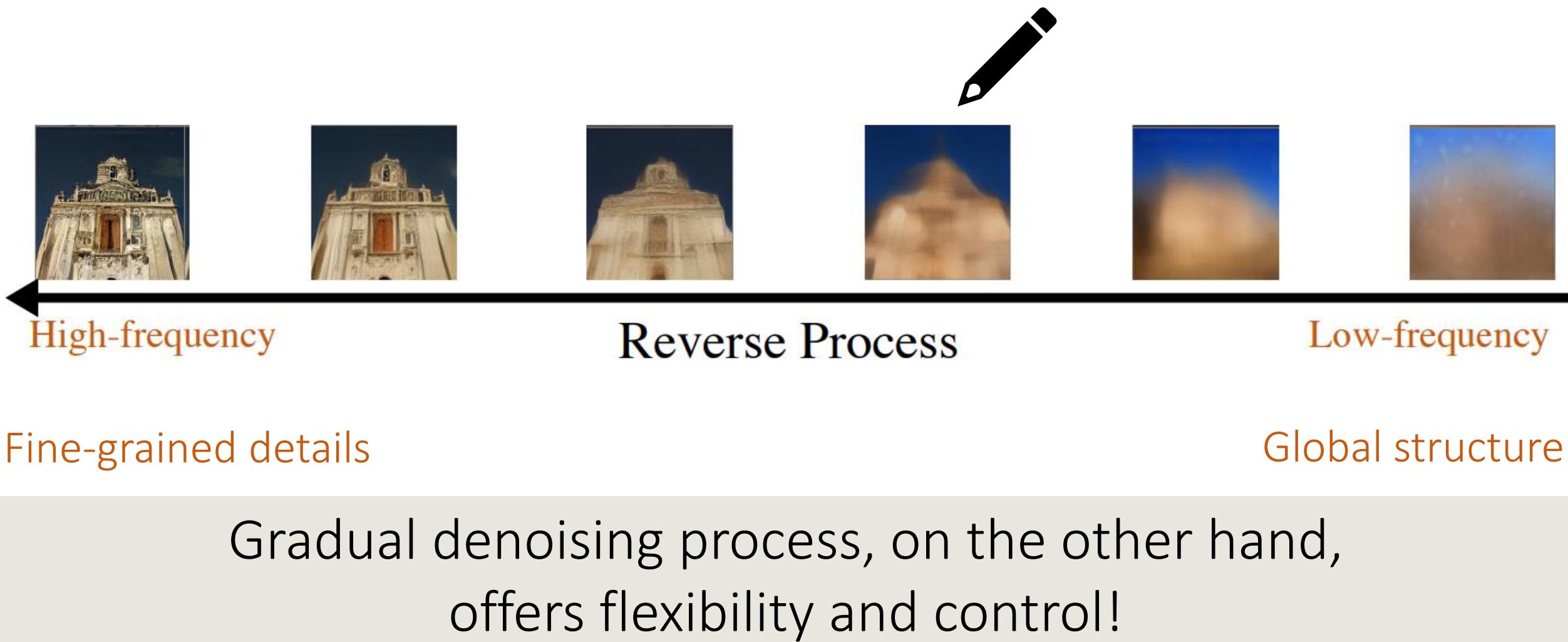
- For $i = N-1, \dots, 1, 0$
 - $dt = t_{i-1} - t_i$
 - $x_{t_{i-1}} \leftarrow x_{t_i} - t_i s_\theta(x_{t_i}, t_i) dt$

Different (but Equivalent) Parametrizations

Given x_t at time t

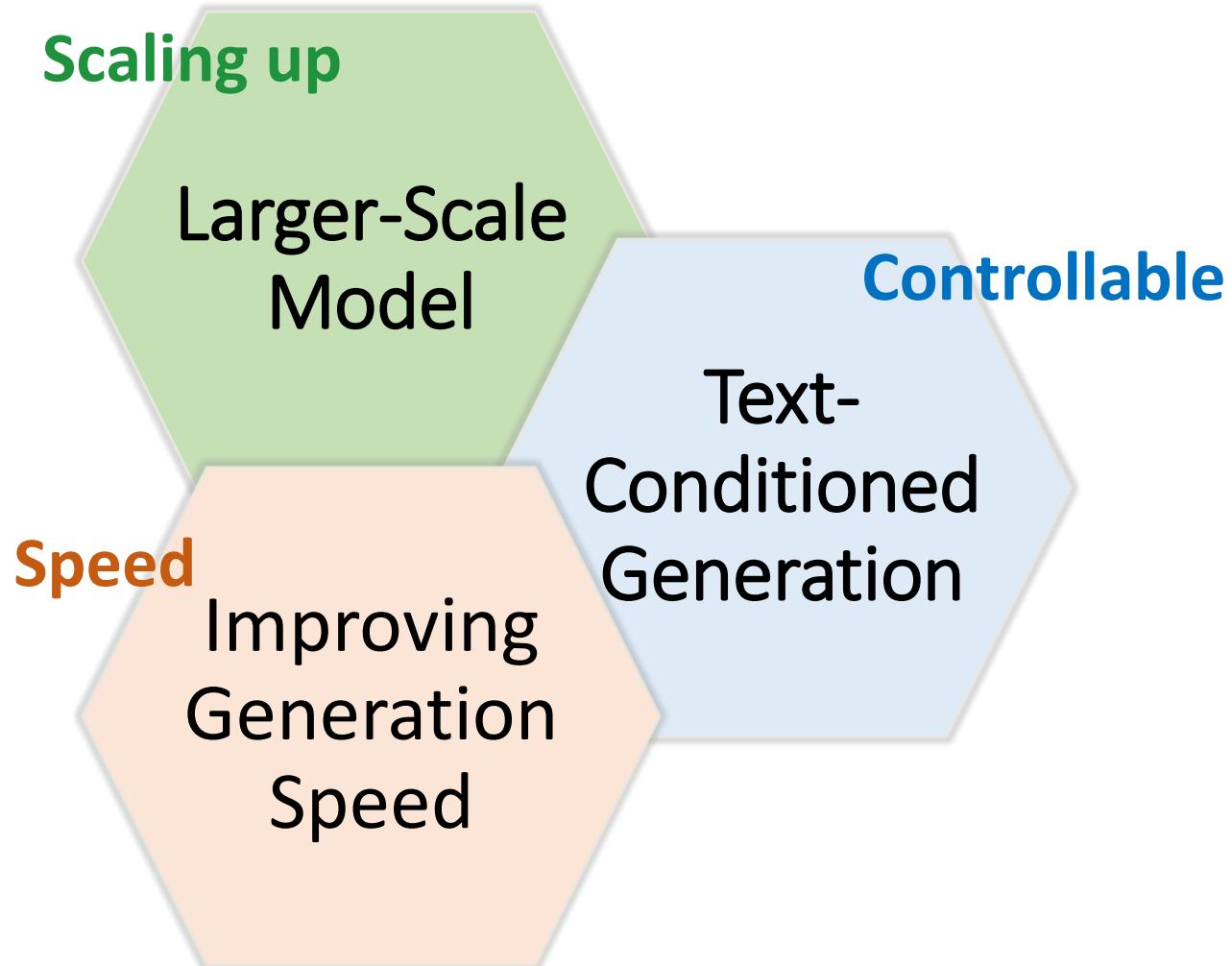


What is Happening in Generation?

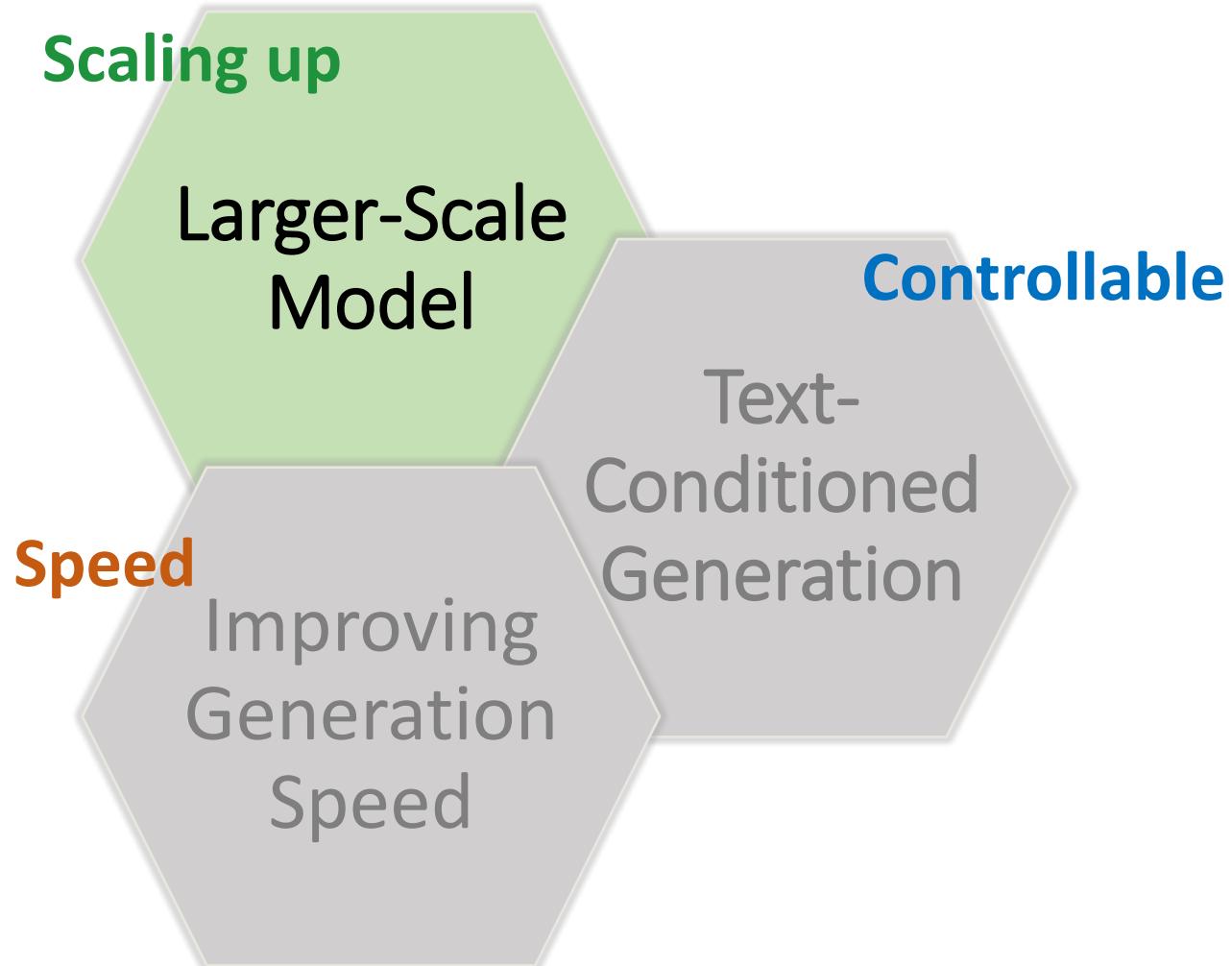


Making Diffusion Model Powerful

Making Diffusion Model More Powerful



Making Diffusion Model More Powerful



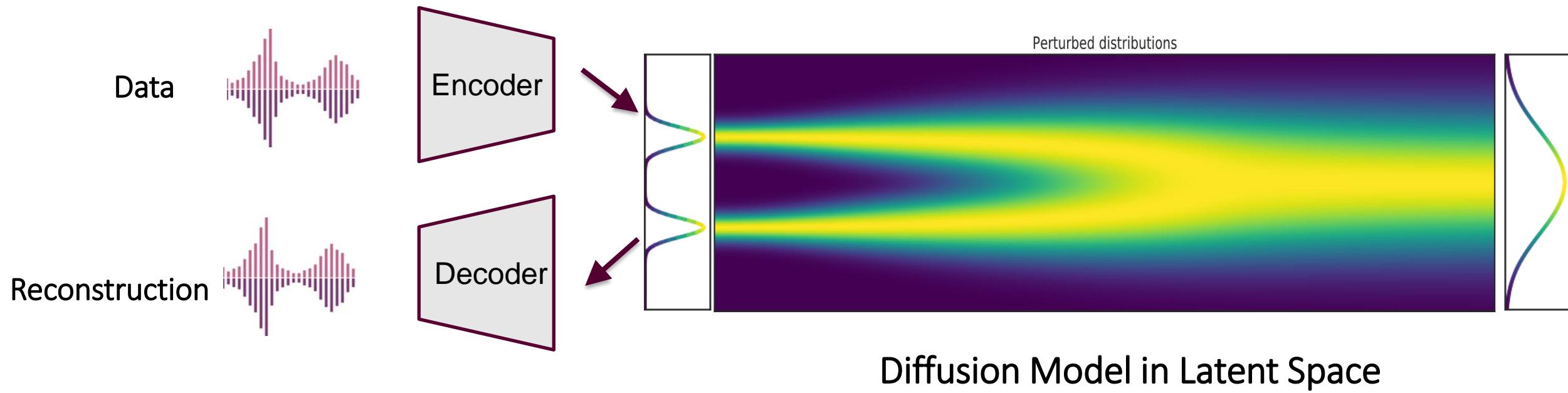
Latent Diffusion Model

Challenge: Music generation requires modeling long temporal sequences.

A solution:

Compress music data via autoencoder →

Examples: MusicLDM, Riffusion, Stable Audio etc..



Diffusion Model in Latent Space

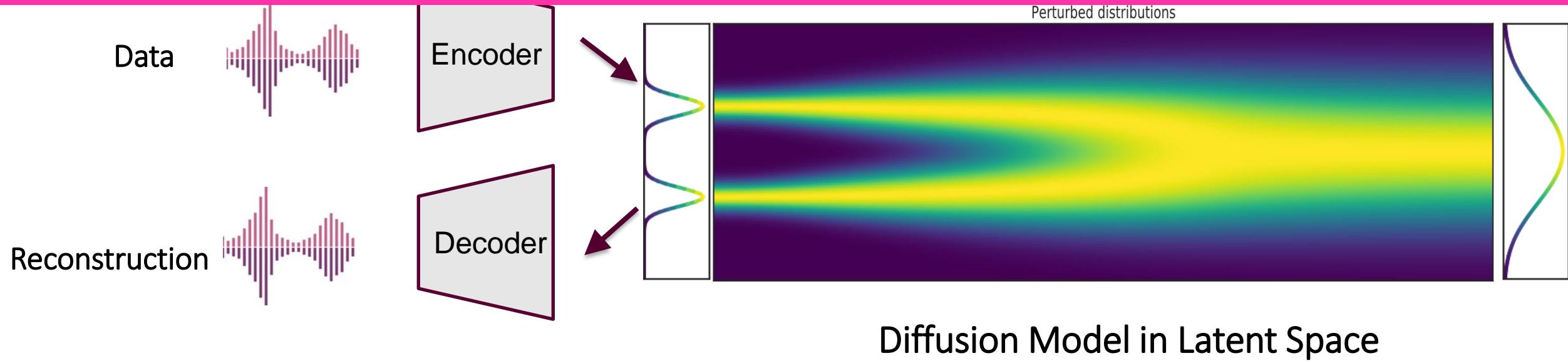
Latent Diffusion Model

Challenge: Music generation requires modeling long temporal sequences.

A solution:

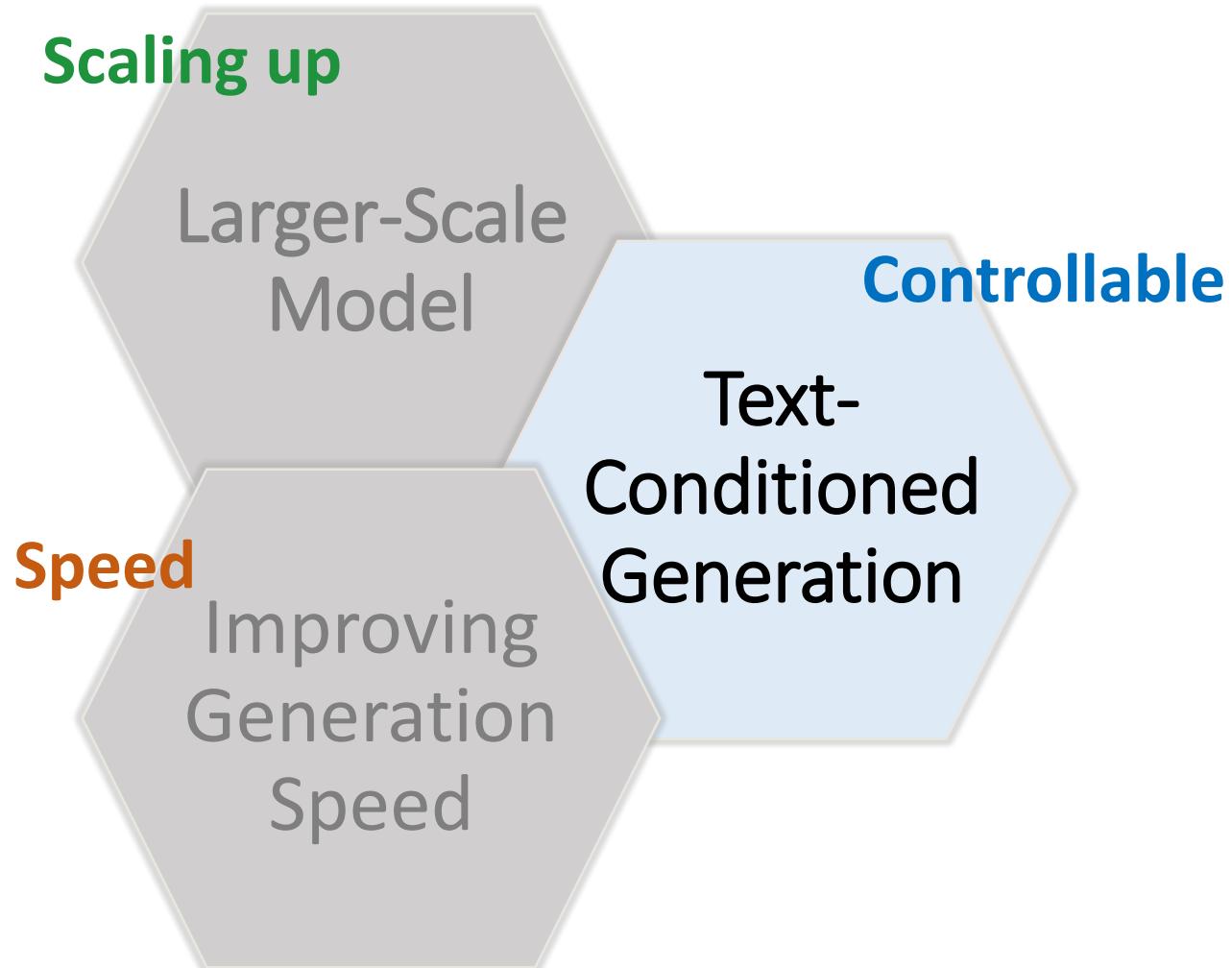
Compress music data via autoencoder →

More Introduction in Session 2...



Diffusion Model in Latent Space

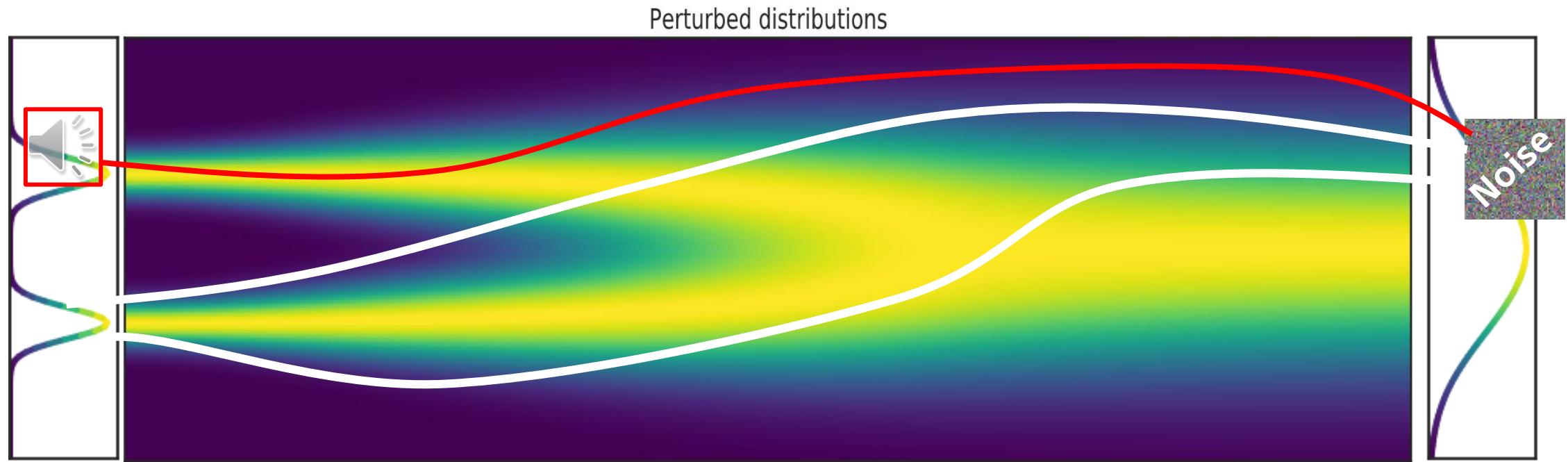
Making Diffusion Model More Powerful



Text-Conditional Generation for Controllability

Goal: Make diffusion model generates music in response to user specified text prompts

Frogs croaking and a humming with insects vocalizing.

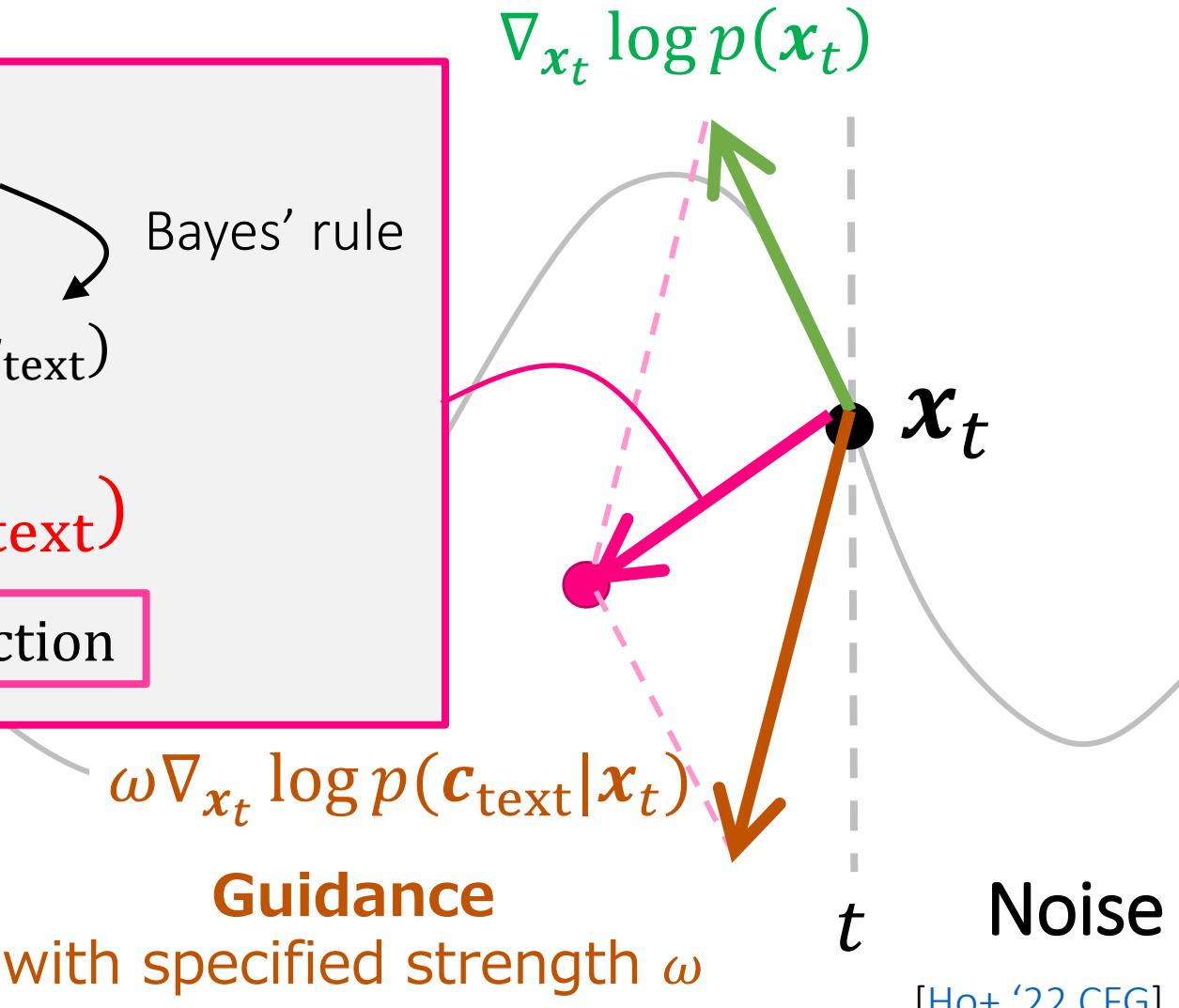


Classifier-Free Guidance (CFG) for Conditional Generation

Given a text condition: \mathbf{c}_{text} .

$$\begin{aligned} & \nabla_{x_t} \log p(x_t) + \omega \nabla_{x_t} \log p(\mathbf{c}_{\text{text}} | x_t) \\ &= (1 - \omega) \nabla_{x_t} \log p(x_t) + \omega \nabla_{x_t} \log p(x_t | \mathbf{c}_{\text{text}}) \\ &\quad \Downarrow \qquad \Downarrow \\ & \mathbf{s}_\theta(x_t, t, \text{null}) \quad \mathbf{s}_\theta(x_t, t, \mathbf{c}_{\text{text}}) \\ &= (1 - \omega) \text{Uncon. Direction} + \omega \text{Con. Direction} \end{aligned}$$

Unconditional Direction



Training and Generation with CFG

Joint training of (un-) conditional DM

Require: κ ratio of Null condition

Initialize: $s_\theta(x_t, t, c_{\text{text}})$

- Get paired data (x, c_{text})
- $c_{\text{text}} \leftarrow \text{Null}$ with ratio κ
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_\theta \|s_\theta(x_t, t, c_{\text{text}}) - \left(-\frac{1}{t}\epsilon\right)\|^2$$

Training DM

Initialize: $s_\theta(x_t, t)$

- Get data x_0
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_\theta \|s_\theta(x_t, t) - \left(-\frac{1}{t}\epsilon\right)\|^2$$

Training and Generation with CFG

Joint training of (un-) conditional DM

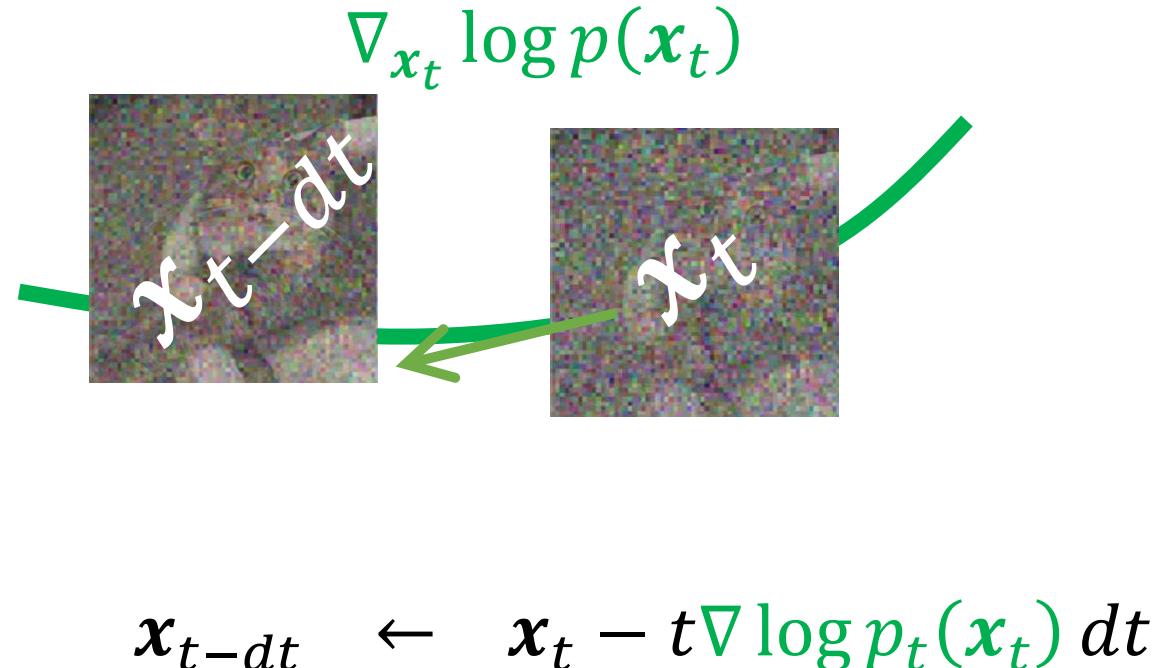
Require: κ ratio of Null condition

Initialize: $s_\theta(x_t, t, c_{\text{text}})$

- Get paired data (x, c_{text})
- $c_{\text{text}} \leftarrow \text{Null}$ with ratio κ
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

$$\nabla_\theta \| s_\theta(x_t, t, c_{\text{text}}) - \left(-\frac{1}{t}\epsilon \right) \|^2$$

Sampling with CFG



Training and Generation with CFG

Joint training of (un-) conditional DM

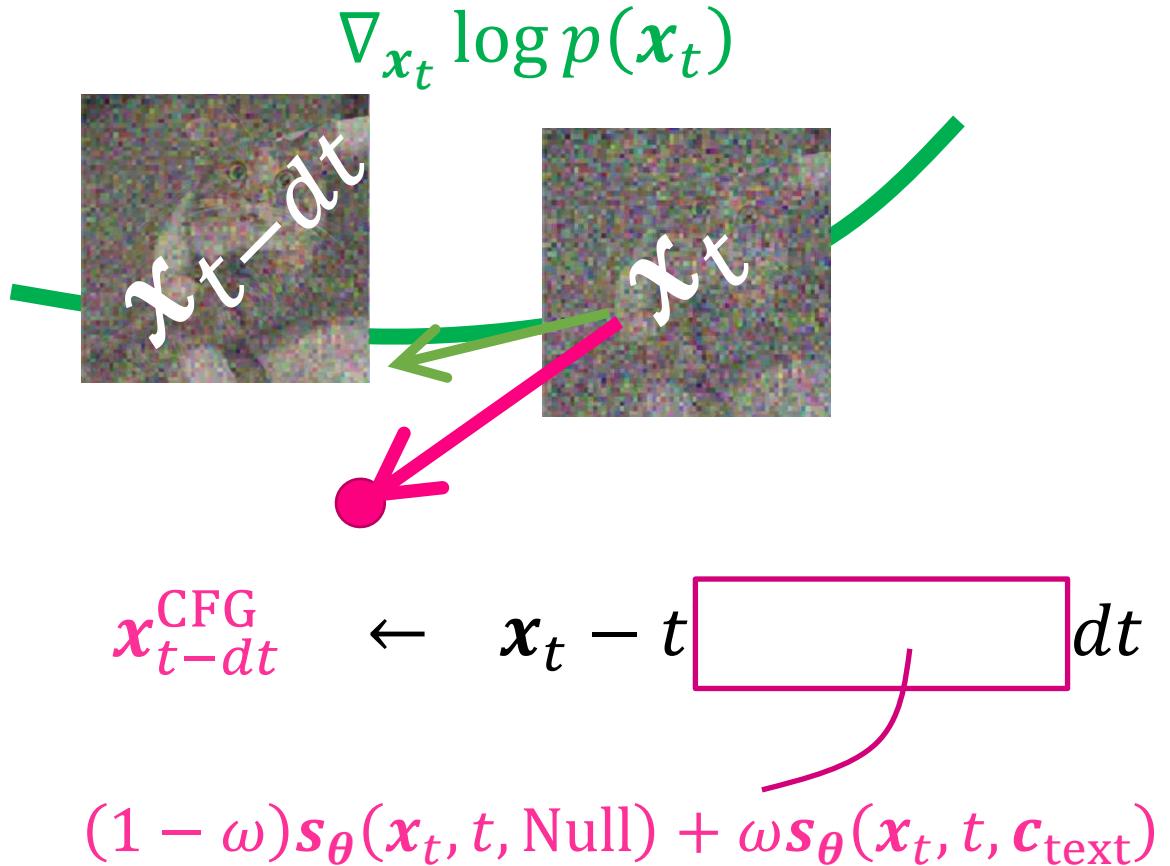
Require: κ ratio of Null condition

Initialize: $s_\theta(x_t, t, c_{\text{text}})$

- Get paired data (x, c_{text})
- $c_{\text{text}} \leftarrow \text{Null}$ with ratio κ
- Sample time $t \sim \mathcal{U}(0,1)$ and noise $\epsilon \sim \mathcal{N}(\mathbf{0}, I)$
- Perturb $x_t = x_0 + t\epsilon$
- Training with Denoising Score Matching:

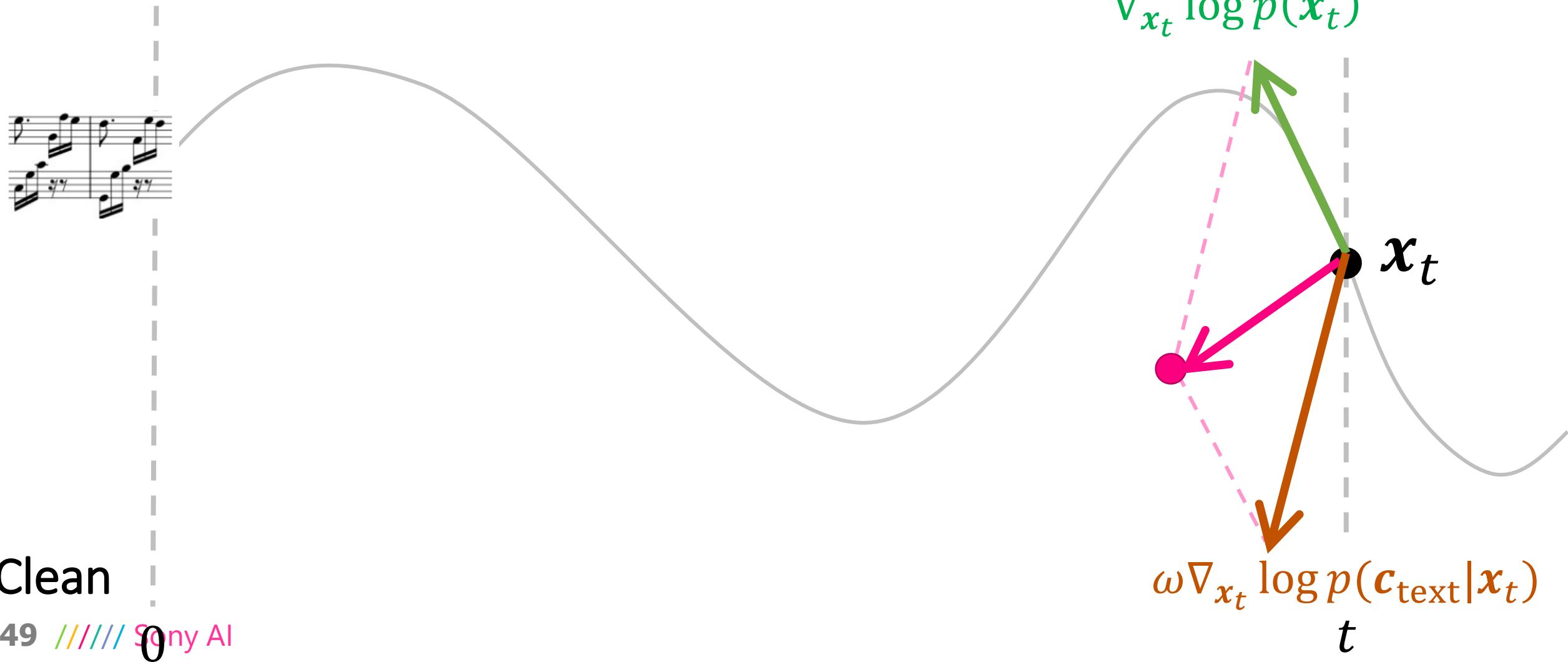
$$\nabla_\theta \| s_\theta(x_t, t, c_{\text{text}}) - \left(-\frac{1}{t} \epsilon \right) \|^2$$

Sampling with CFG



Training and Generation with CFG

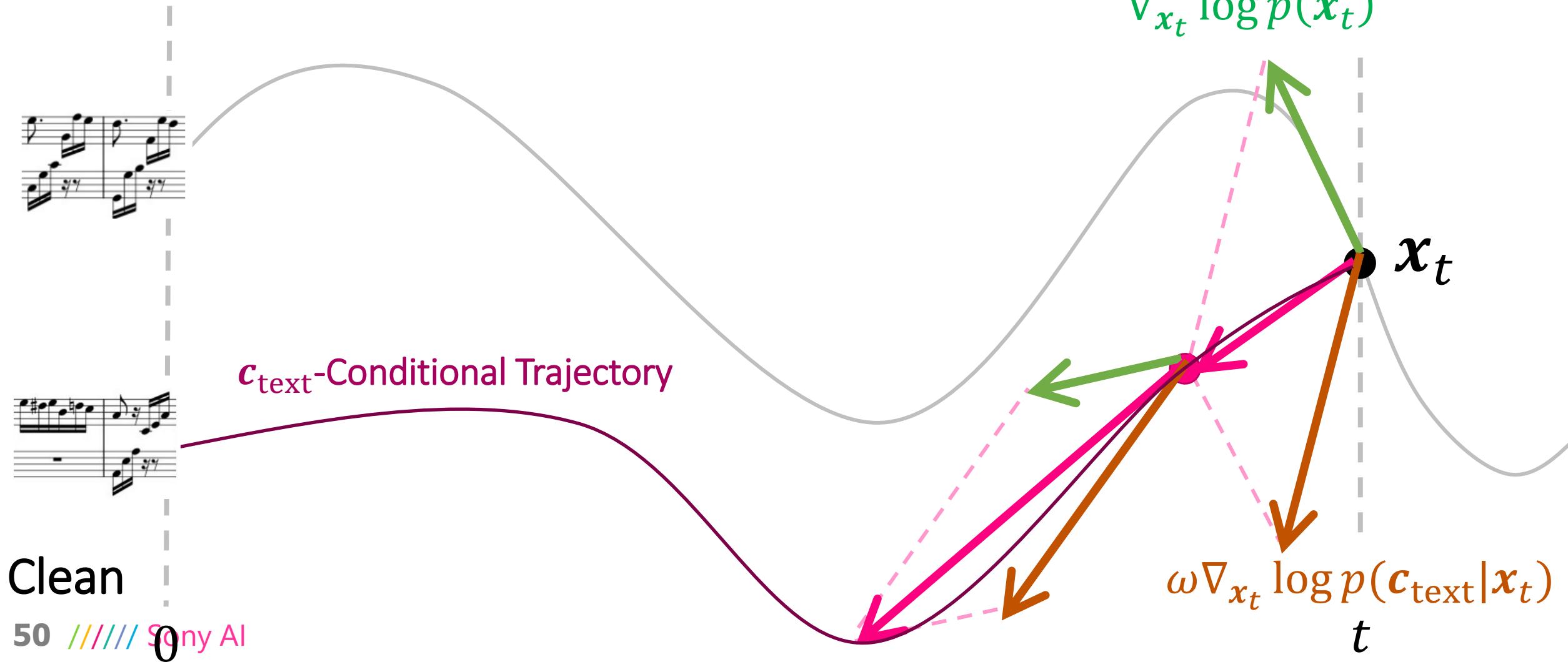
Given a text condition: c_{text} .



Training and Generation with CFG

Given a text condition: c_{text} .

Unconditional Direction



Training and Generation with CFG

Given a text condition: c_{text} .

Unconditional Direction

$$\nabla_{x_t} \log p(x_t)$$



More Introduction in Session 2 & 3...

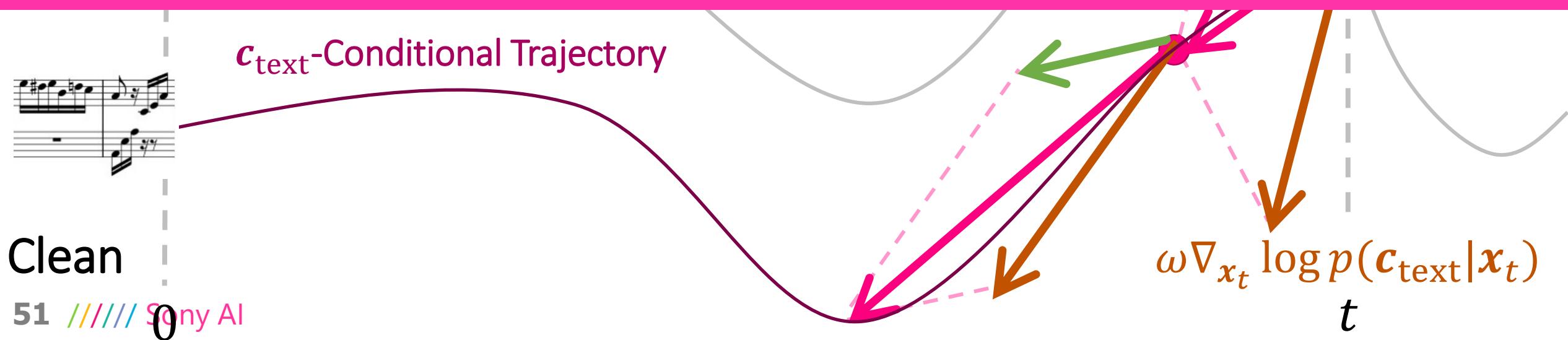
c_{text} -Conditional Trajectory

Clean

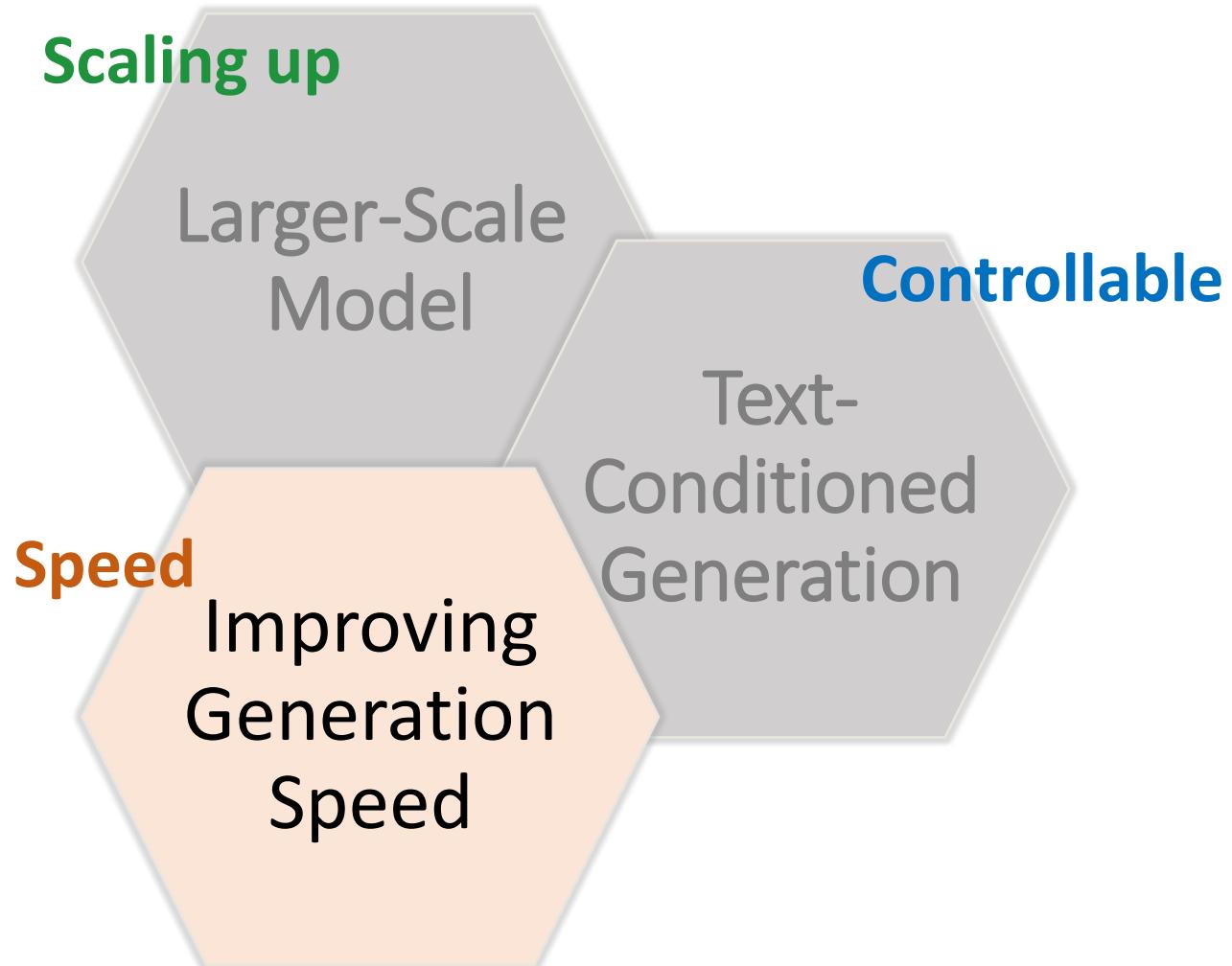
51 ////////////////////////////////////////////////////////////////// Sony AI 0

$$\omega \nabla_{x_t} \log p(c_{\text{text}} | x_t)$$

t

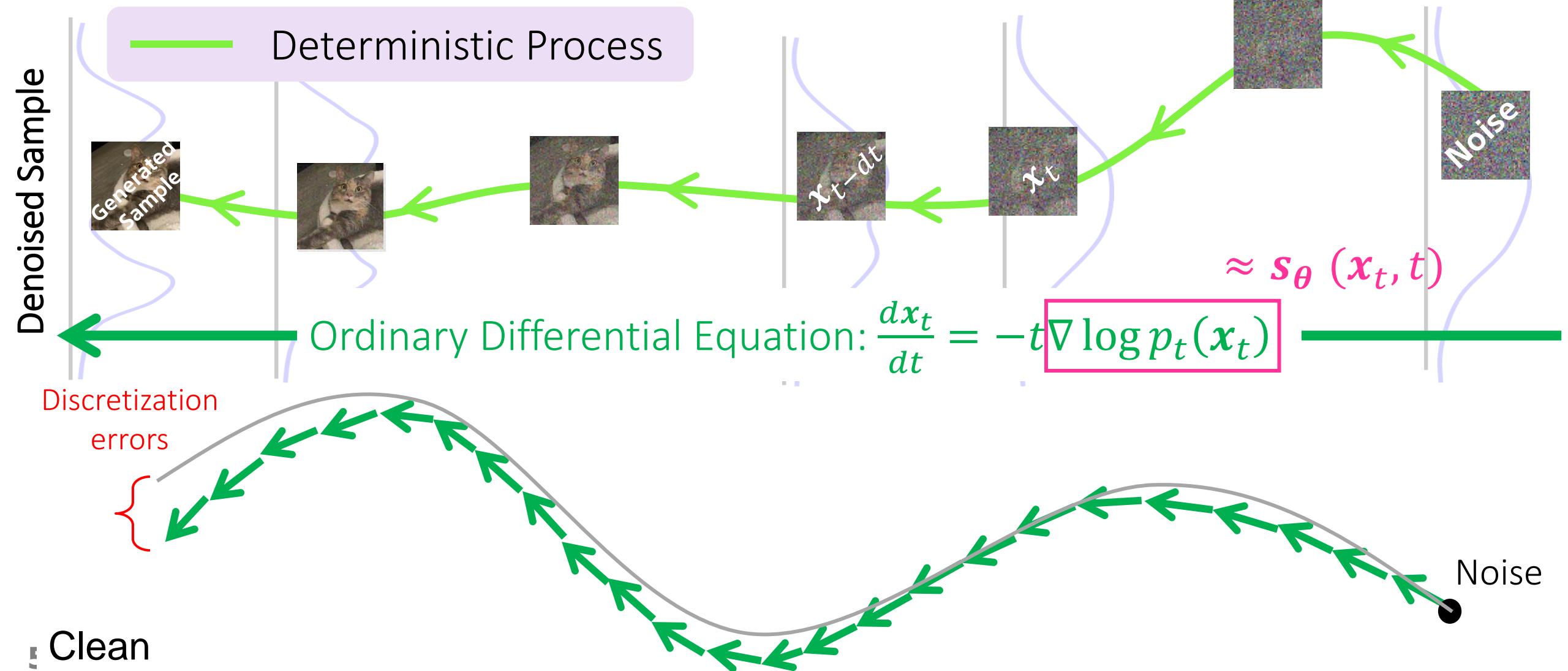


Making Diffusion Model More Powerful



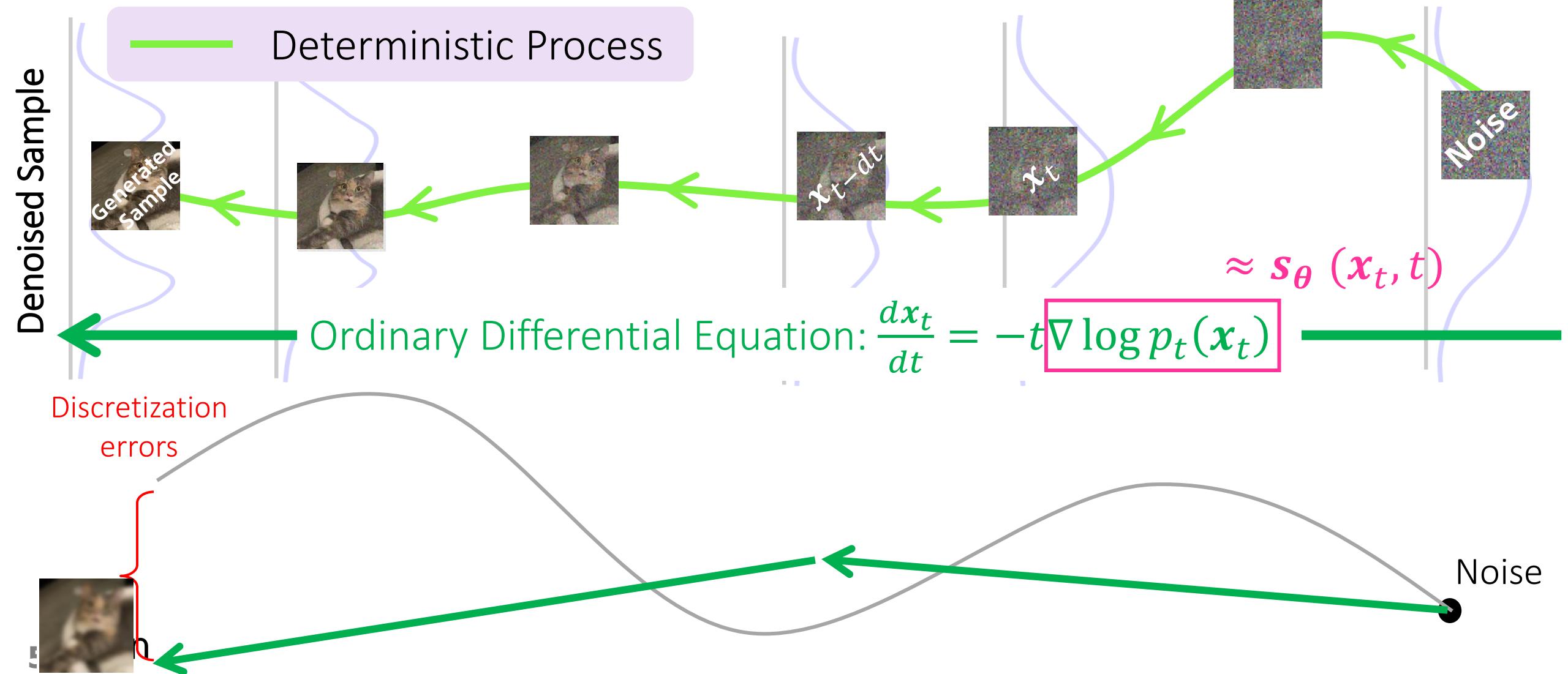
Score-Based Diffusion Model's Training and Sampling

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



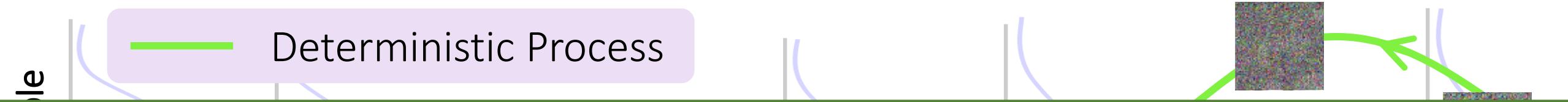
Score-Based Diffusion Model's Training and Sampling

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

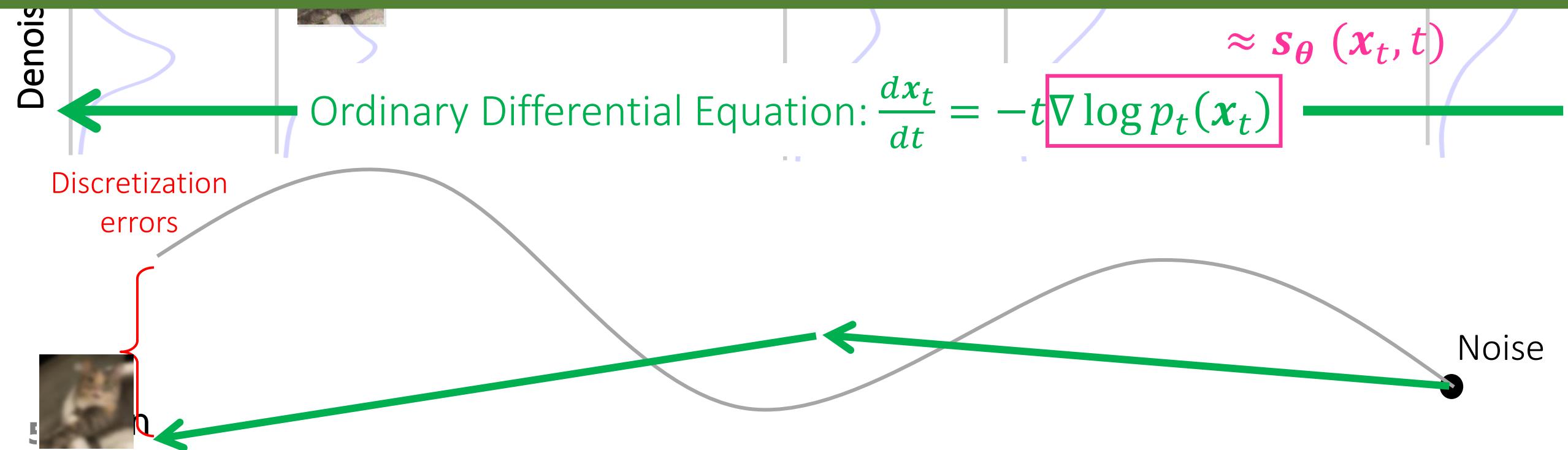


Score-Based Diffusion Model's Training and Sampling

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



Trade-off Between Generation Speed and Quality



Fast Sampling in Diffusion Model

PF-ODE (2022)

Training-free

Sophisticated solvers
(2021-2022+)

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

Gradient

1. Higher-order gradients
2. Adaptive step-sizes
3. Extrapolations for better approximation
4. ...

Gradient



//// Sony AI

Noise

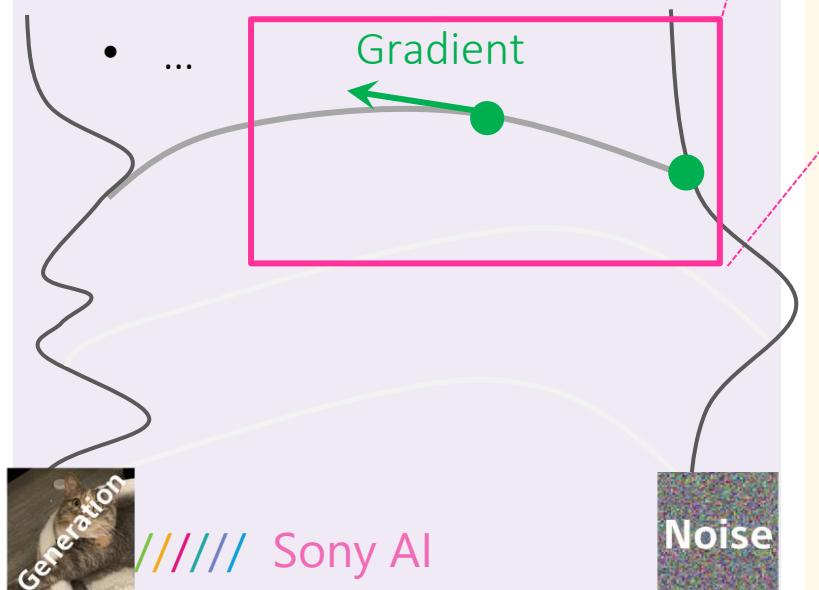
Fast Sampling in Diffusion Model

PF-ODE (2022)

Training-free

Sophisticated solvers
(2021-2022+)

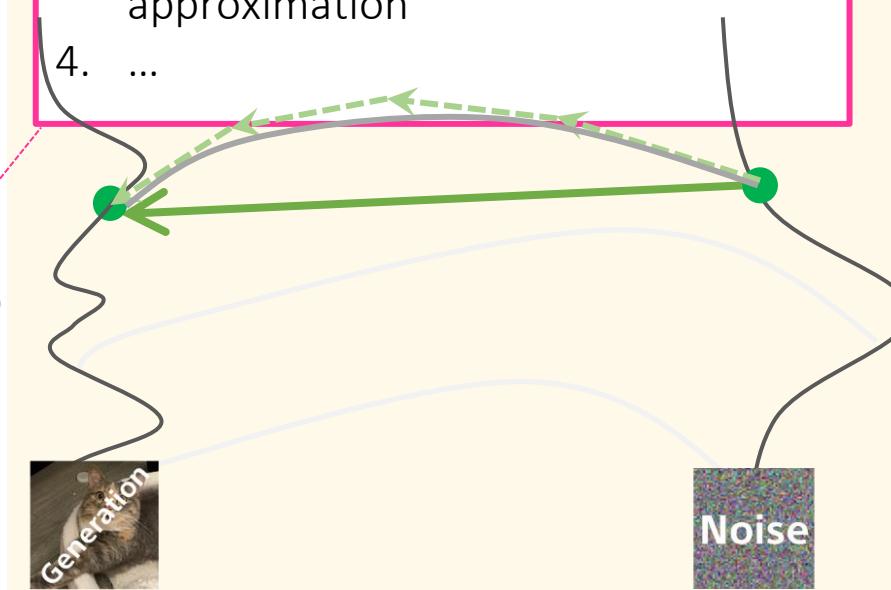
- [ICLR'21] [DDIM](#)
- [NeurIPS'22] [DPM](#); DPM++
- [ICLR'23] [DEIS](#)
- ...



Training-based

Gradient

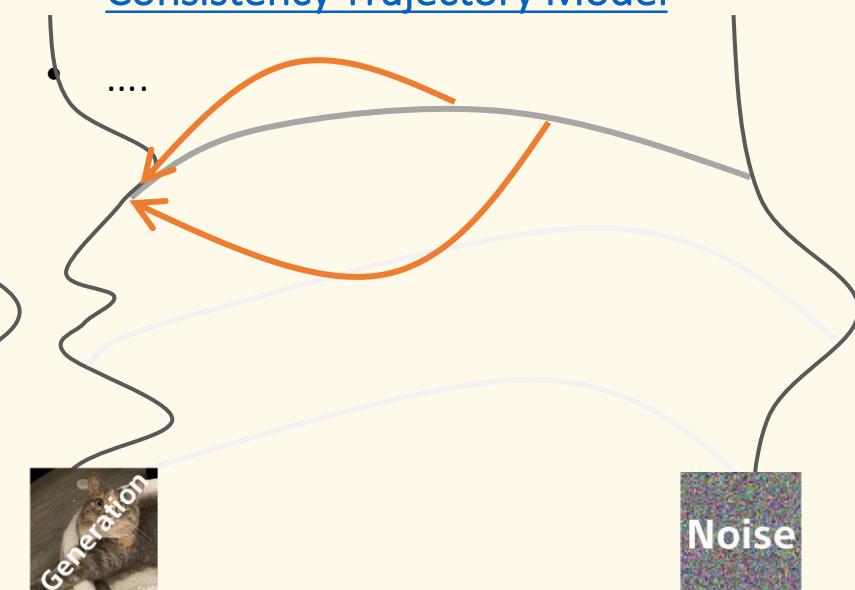
1. Higher-order gradients
2. Adaptive step-sizes
3. Extrapolations for better approximation
4. ...



Distillation from DM (2023-)

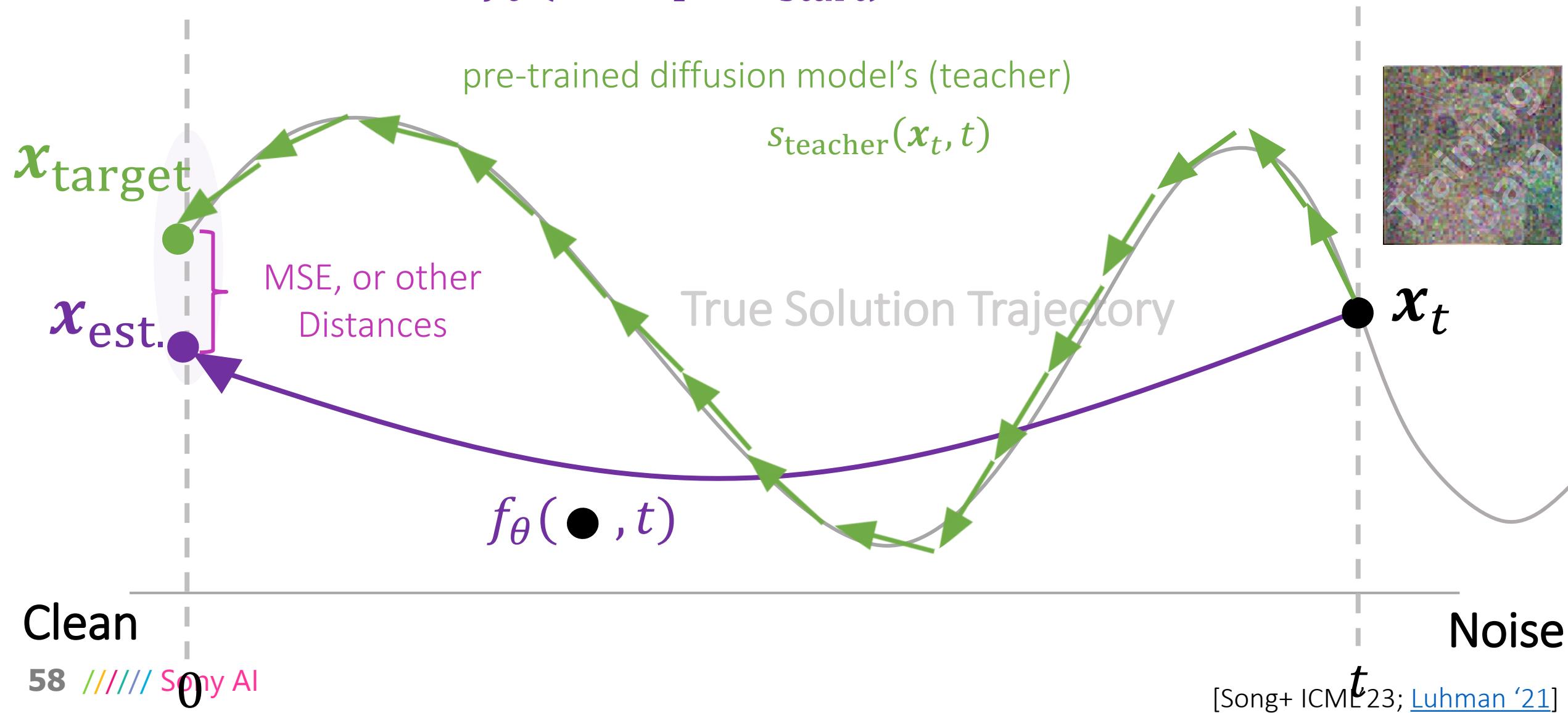
- [ICML'23 Song+] [Consistency Model](#)
- [ICLR'24 Kim&Lai+]

[Consistency Trajectory Model](#)



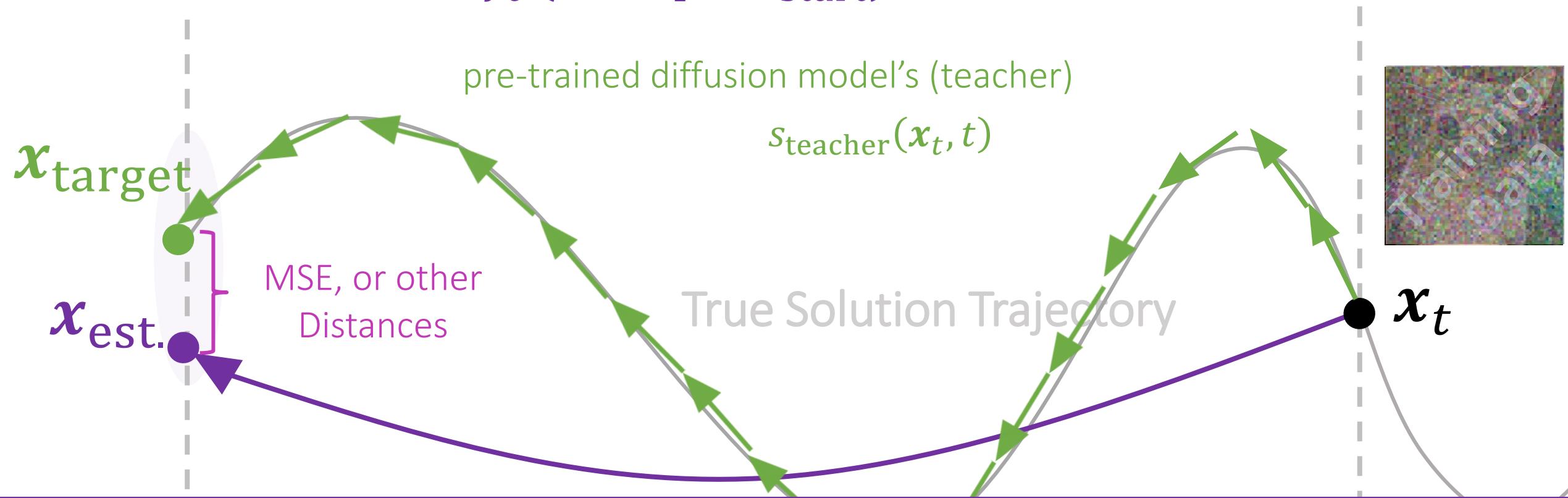
Distillation-Based Method's Training

Student network $f_\theta(\text{start pt.}, t_{\text{start}})$



Distillation-Based Method's Training

Student network $f_\theta(\text{start pt.}, t_{\text{start}})$



It learns any-to-0 jump on the trajectory

Distillation-Based Method's Training

Student network $f_\theta(\text{start pt.}, t_{\text{start}})$

Solving ODE in Training is Expensive!

\hat{x}_{target}

Any-to-0 jump → Less flexible generation

It learns any-to-0 jump on the trajectory

Clean

Noise

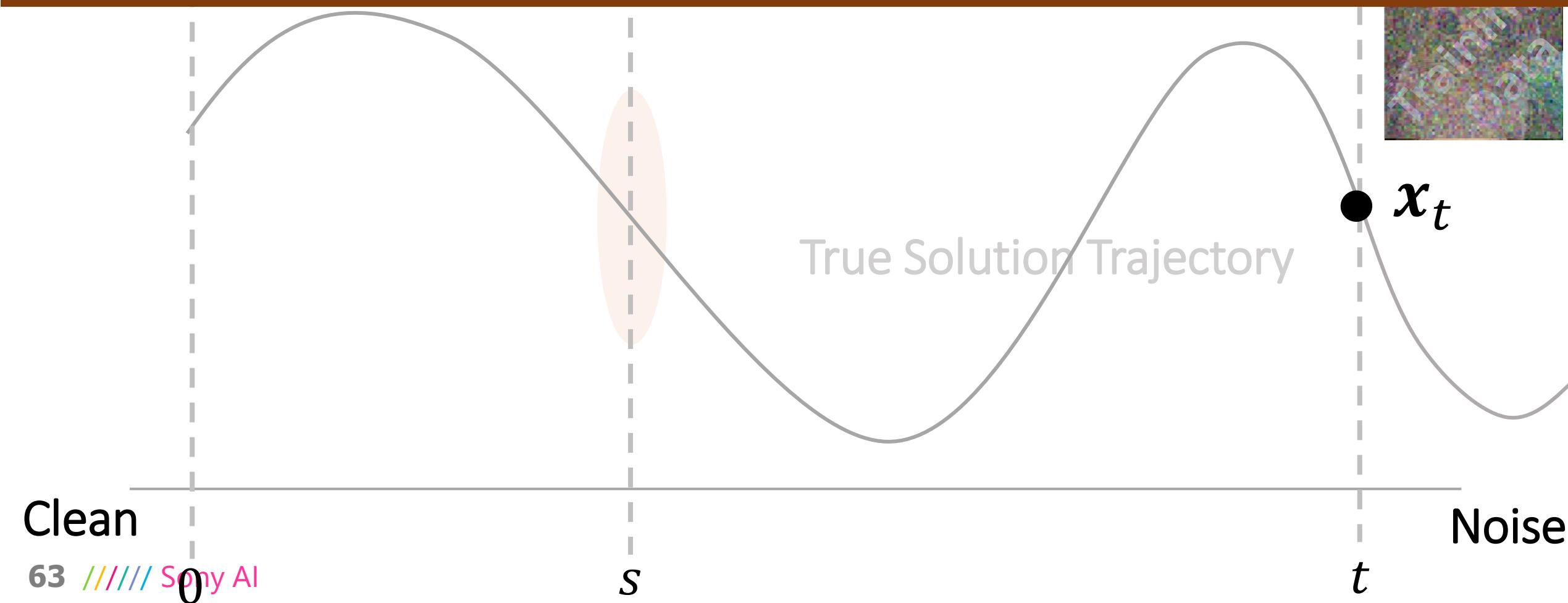
CONSISTENCY TRAJECTORY MODEL (CTM)

- ✓ Efficient Training and Sampling
- ✓ Unifying Diffusion Model and Distillation-Based Method for *Controllable Generation*

CTM's Training and Mechanism

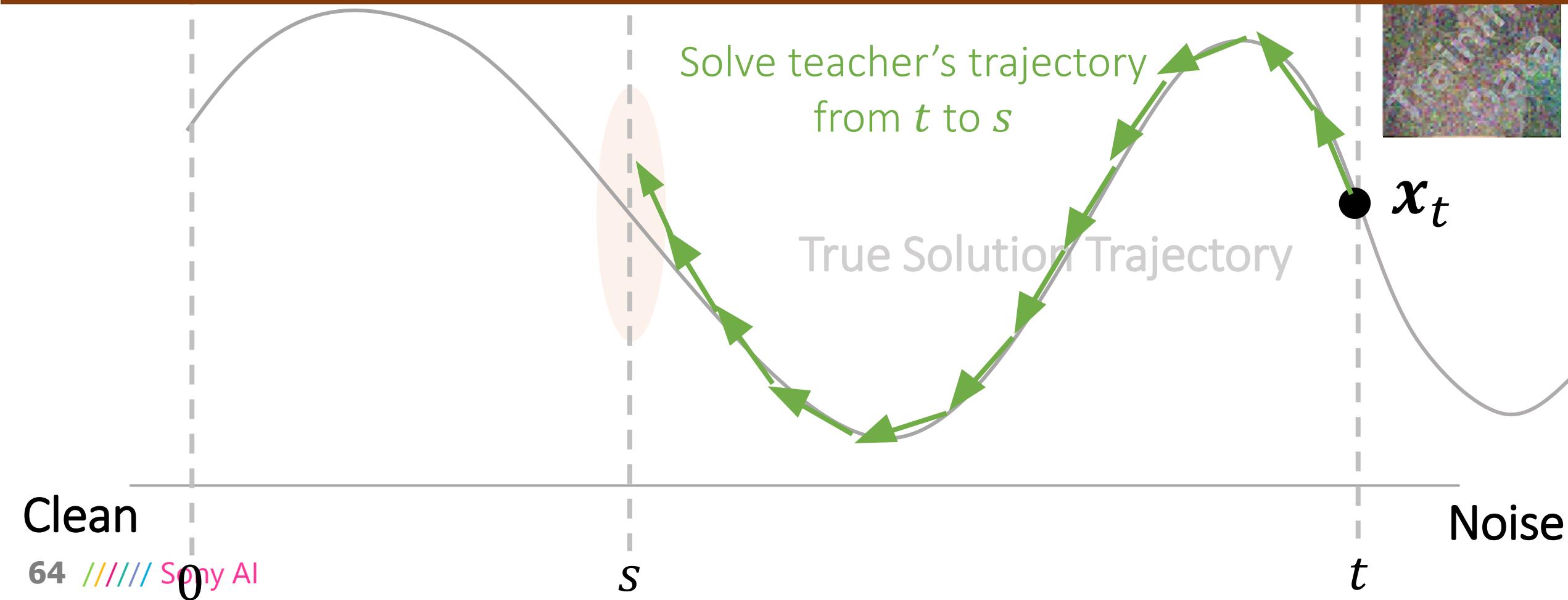
Consistency Trajectory Model (CTM)'s Training

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



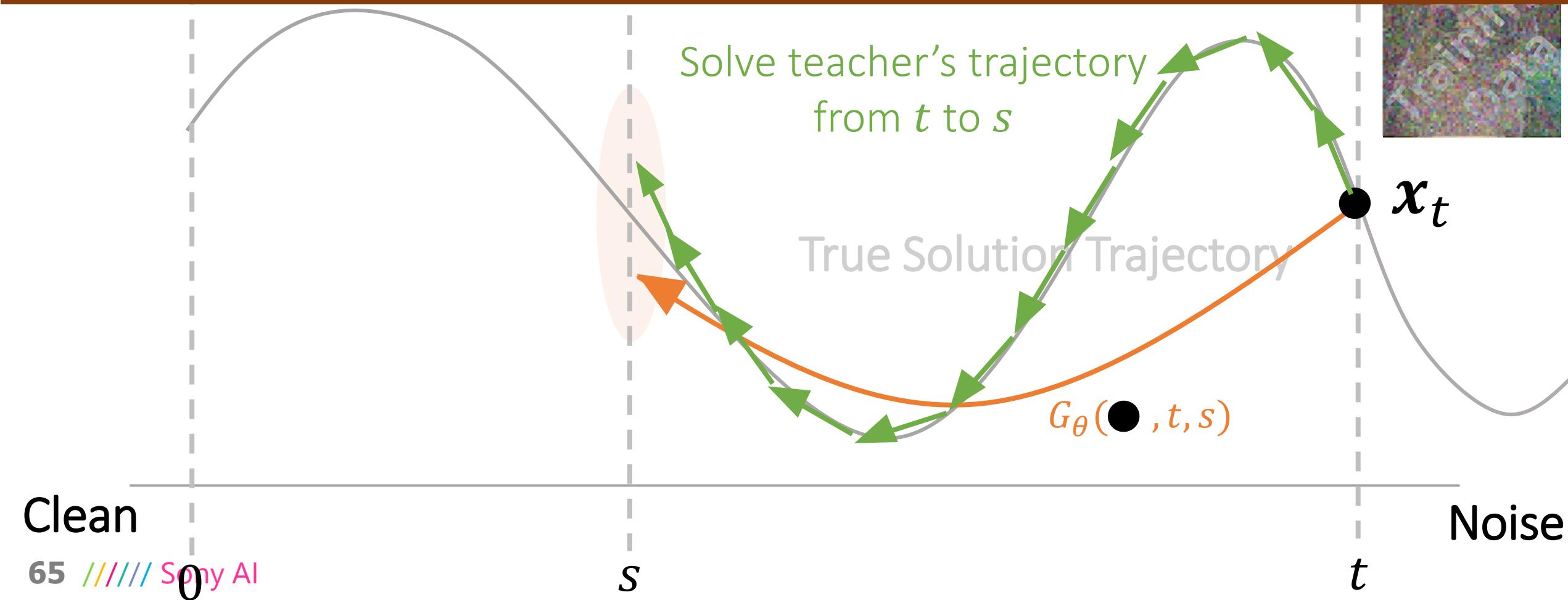
Consistency Trajectory Model (CTM)'s Training

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



Consistency Trajectory Model (CTM)'s Training

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



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

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

2

Clean

Noise

S

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



x_{target}

$x_{\text{est.}}$

Any distances, or even GAN

$G_{sg(\theta)}(\cdot, s, 0)$

True Solution Trajectory

Clean

Noise

s

u

t

x_t

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

True Solution

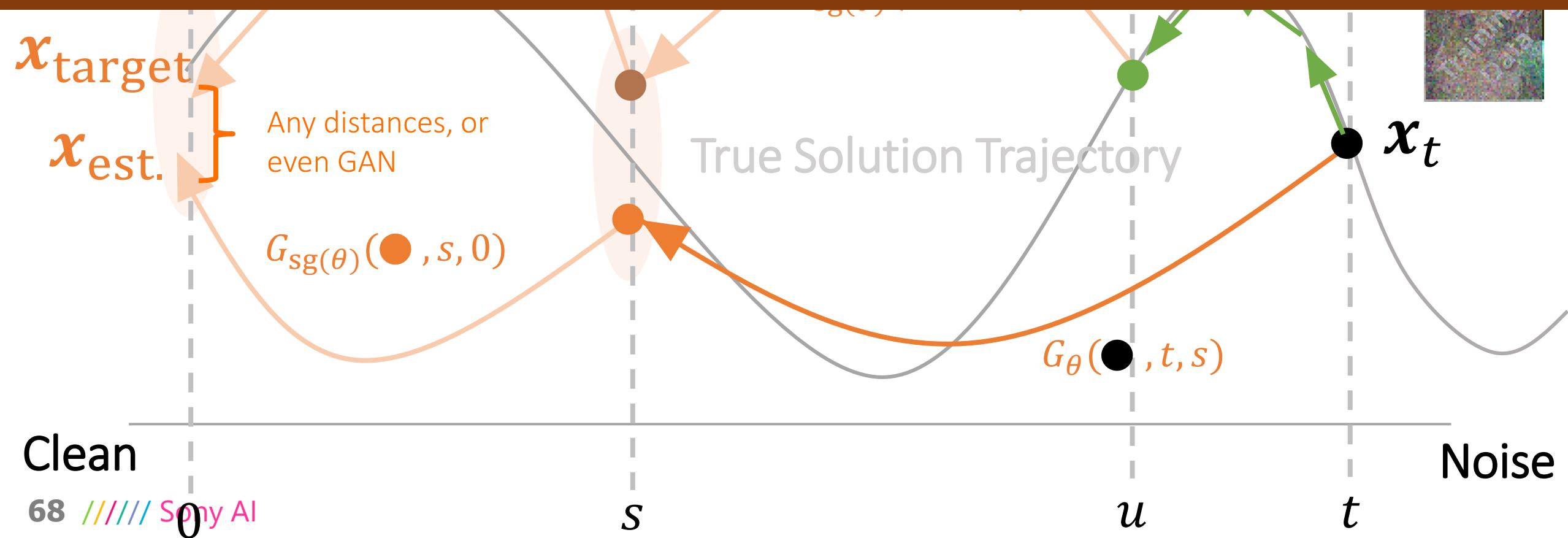
Consistency Trajectory Model (CTM)'s Training

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

CTM's Soft Matching

Technical difficulties avoided

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:

Textbook Diffusion Model

CTM learns any-to-any jump on the trajectory

x_{target}

$x_{est.}$

Any distances, or even GAN

Any distances, or even GAN

True Solution Trajectory

100

x_t

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

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

Clean

Noise

S

w

t

Consistency Trajectory Model (CTM)'s Training

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

CTM's Soft Matching

To conclude, I'd like to thank you all.

CTM learns any-to-any jump on the trajectory

The diagram shows a trajectory starting from a x_{target} point (orange circle) and ending at an x_t point (black circle). The trajectory is labeled "True Solution Trajectory". A green arrow points from x_t towards x_{target} . A bracket on the left indicates the distance between x_{target} and $x_{\text{est.}}$, with the text "Any distances, or even GAN" below it. An orange arrow points from x_{target} towards $x_{\text{est.}}$. A green arrow points from $x_{\text{est.}}$ towards x_t . A small image of a noisy image is shown next to x_t .

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

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

Clean

Noise

CTM's Flexible Sampling

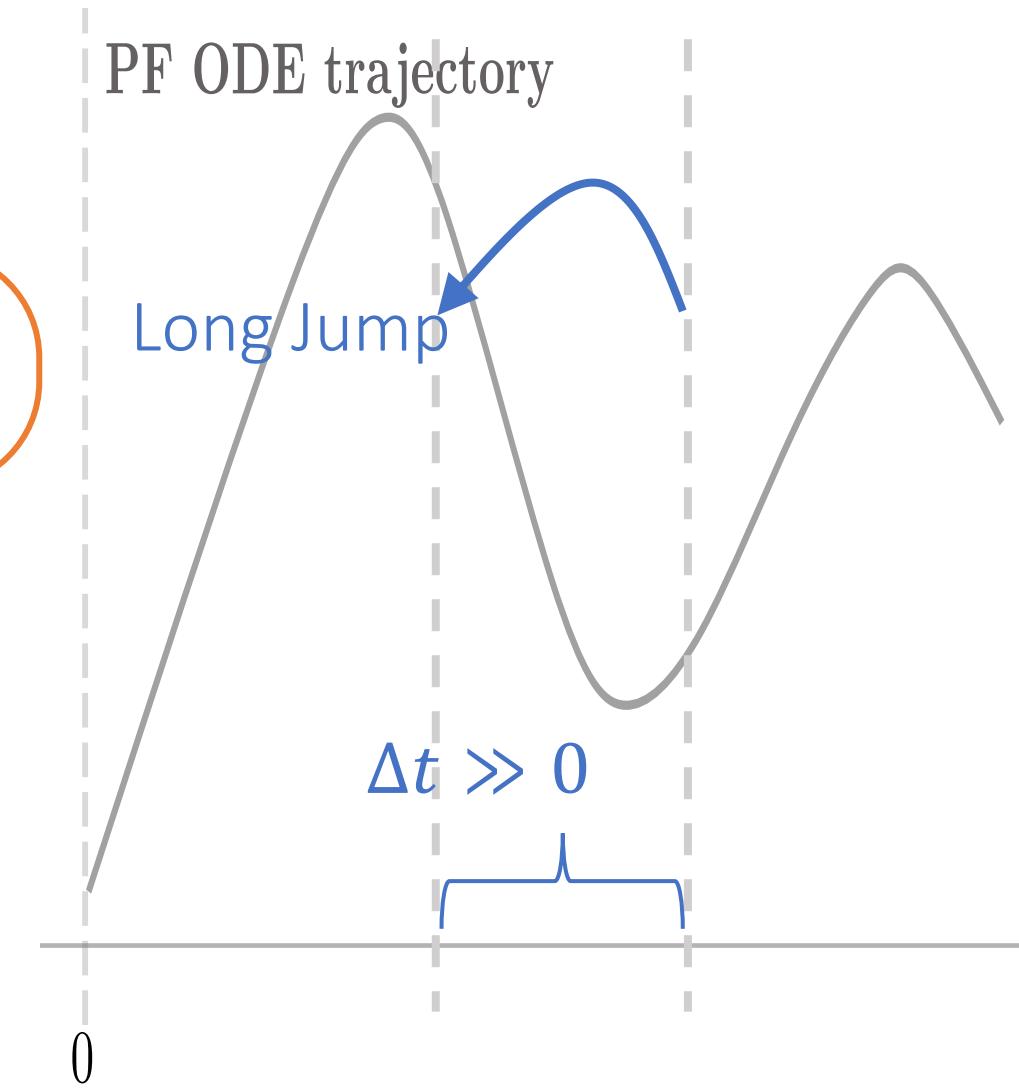
CTM's Sampling Flexibility

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!

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

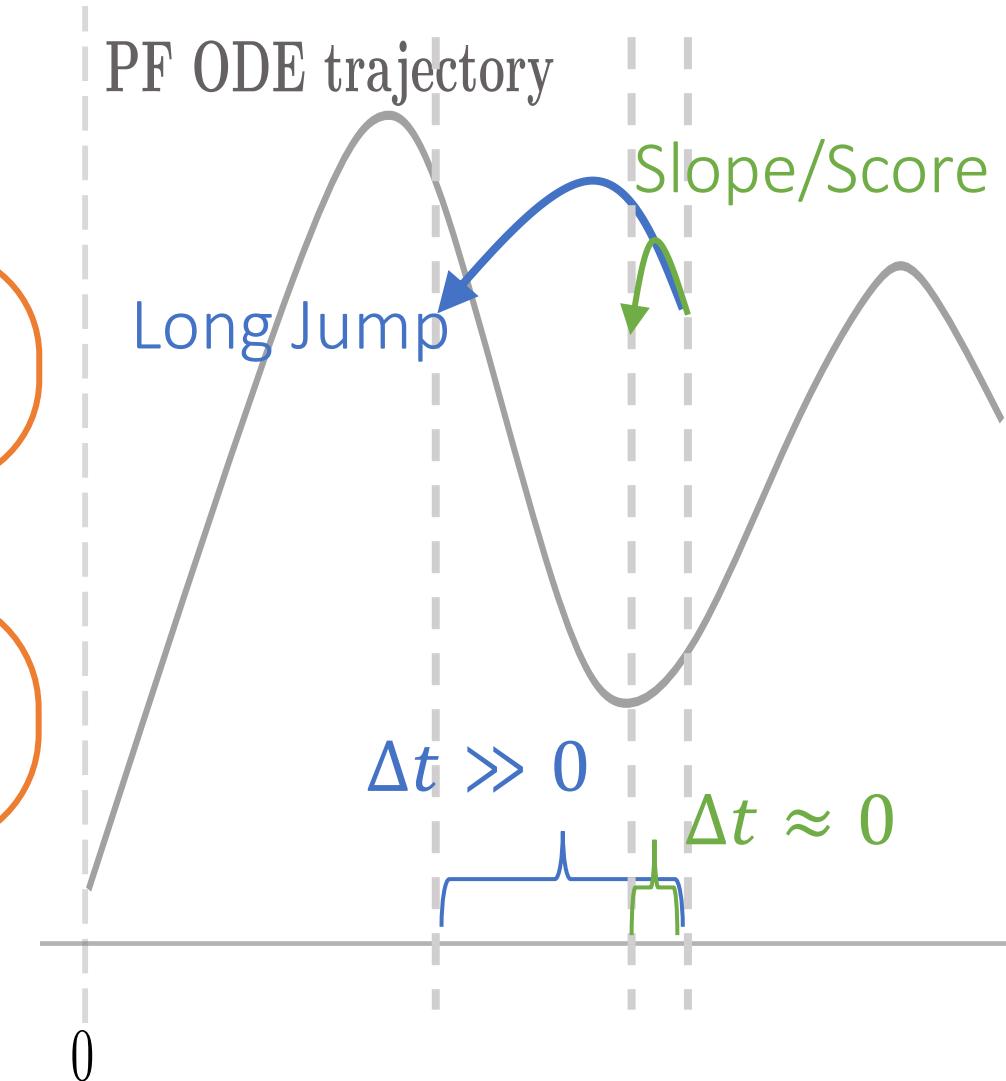
Long jump ($\Delta t \gg 0$)

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

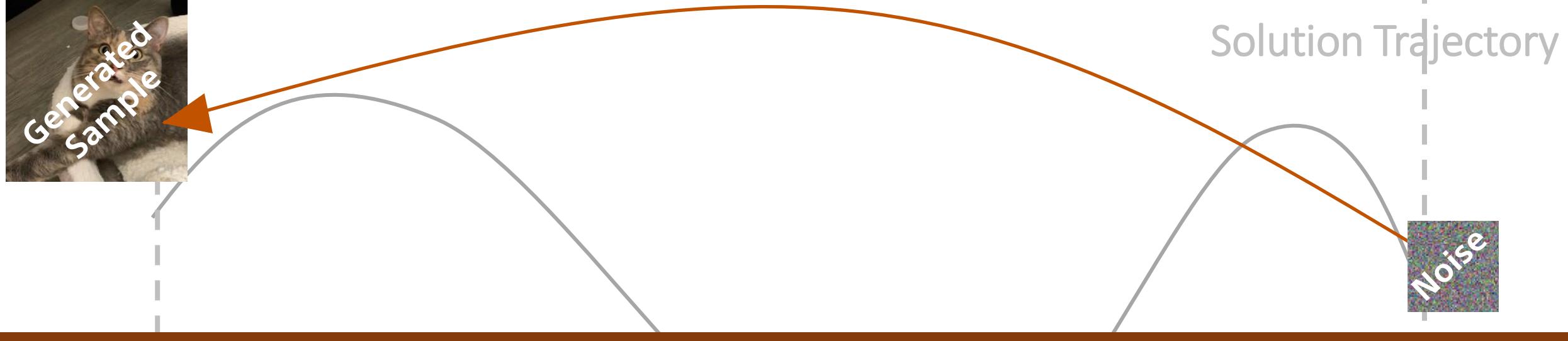
Slope (infinitesimal jump $\Delta t \approx 0$)

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

74 ////////////////////////////////////////////////////////////////// Sony AI

Noise

T

SOTA (beats DM) with One-Step



Teacher: Sampling steps = 79

Takes minutes

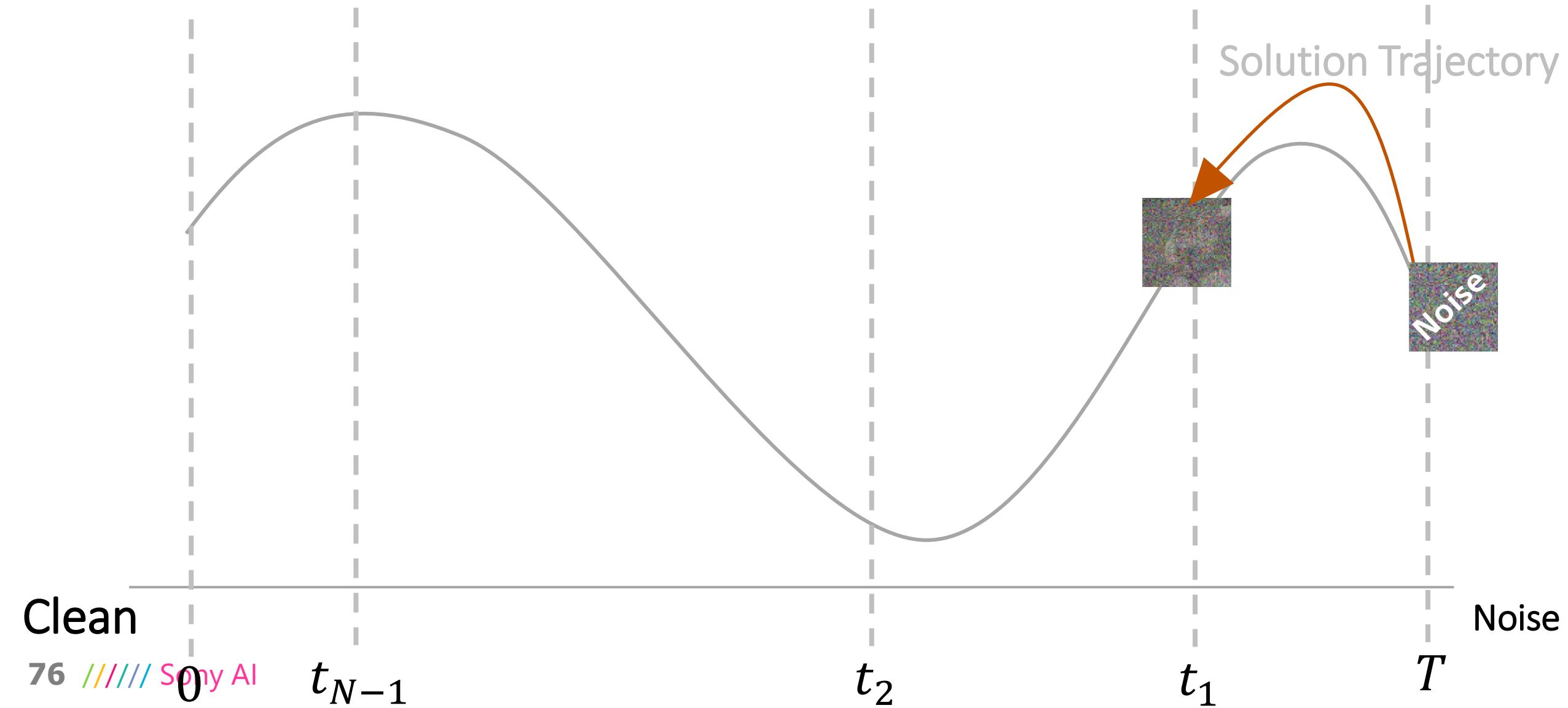


Student (CTM): Sampling step = 1

Within a second

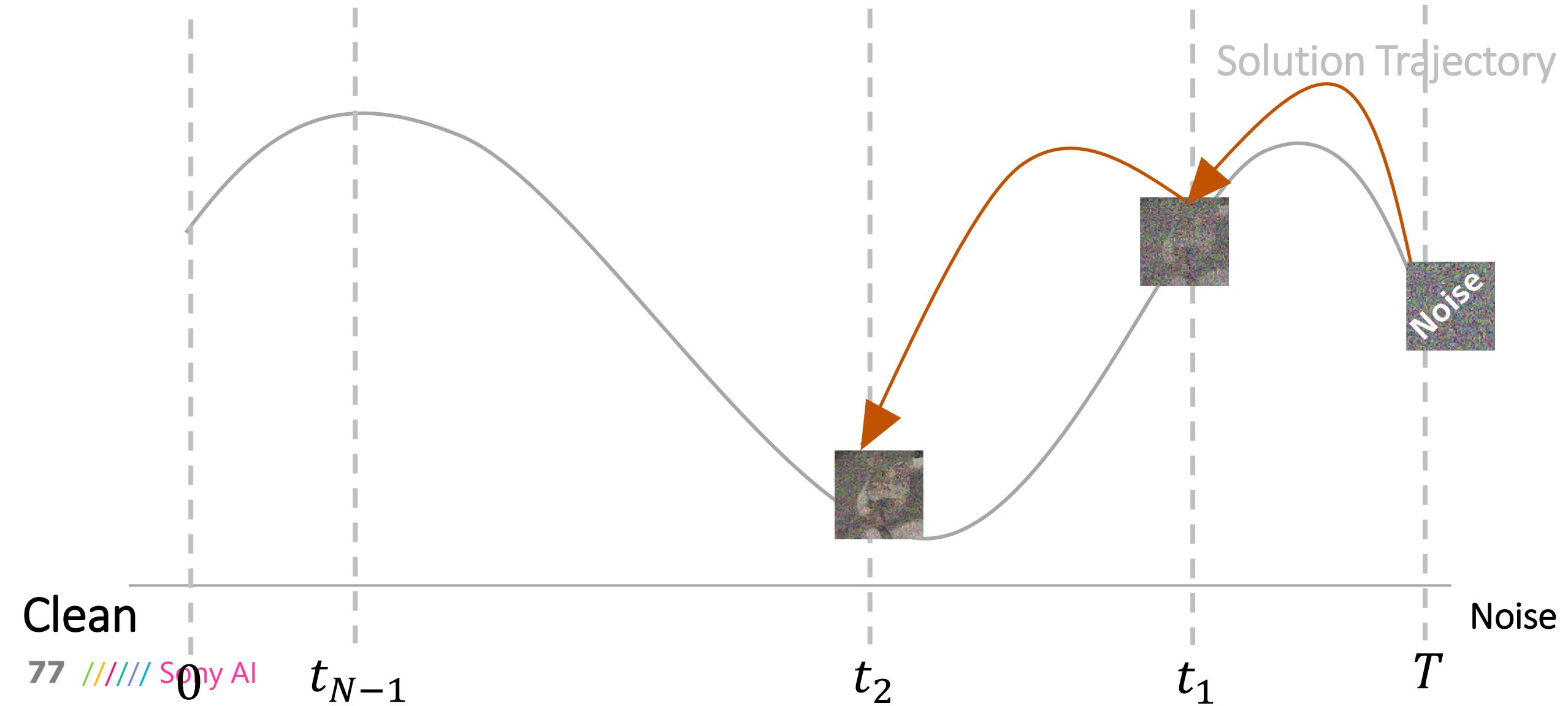
Multistep Generation with CTM: γ -sampling

$\gamma = 0$: Fully Deterministic Sampling



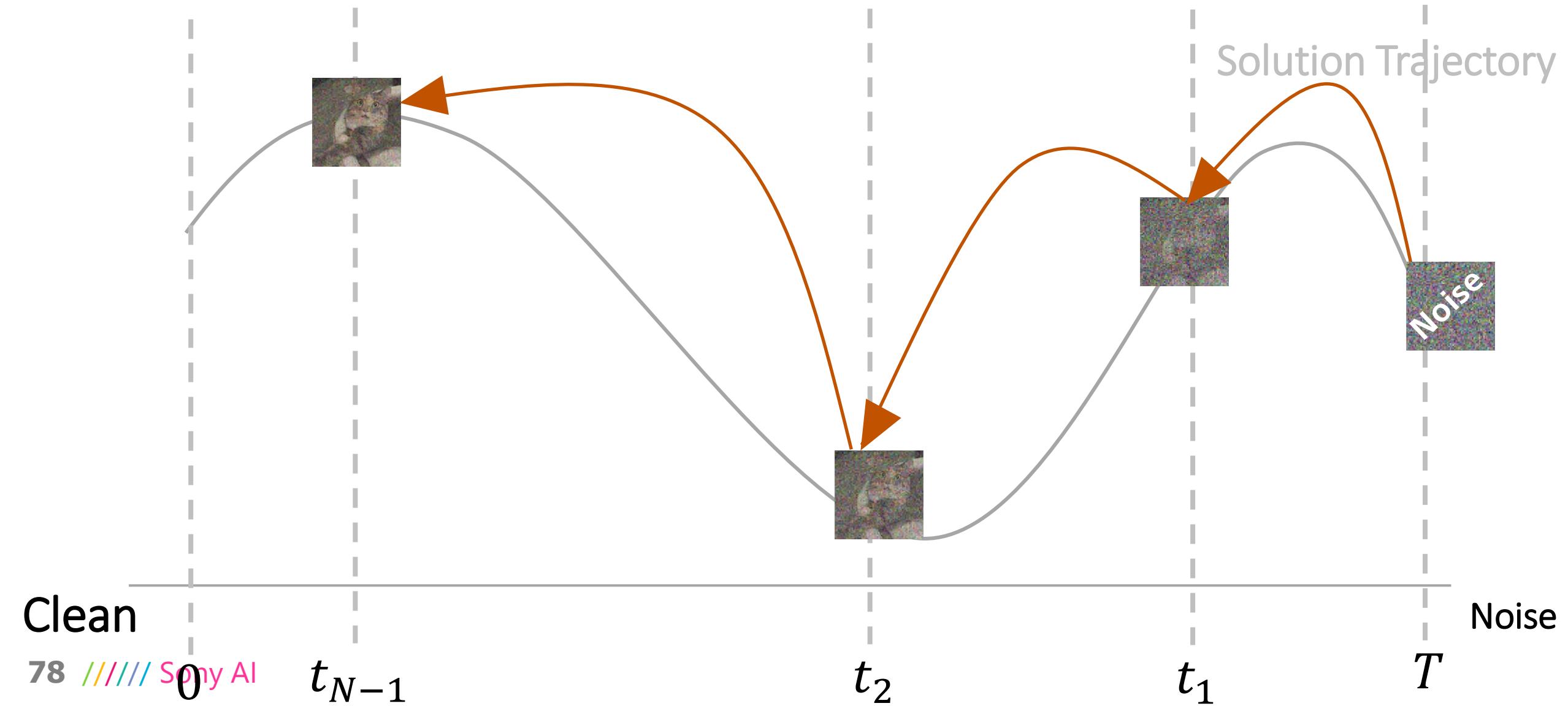
Multistep Generation with CTM: γ -sampling

$\gamma = 0$: Fully Deterministic Sampling



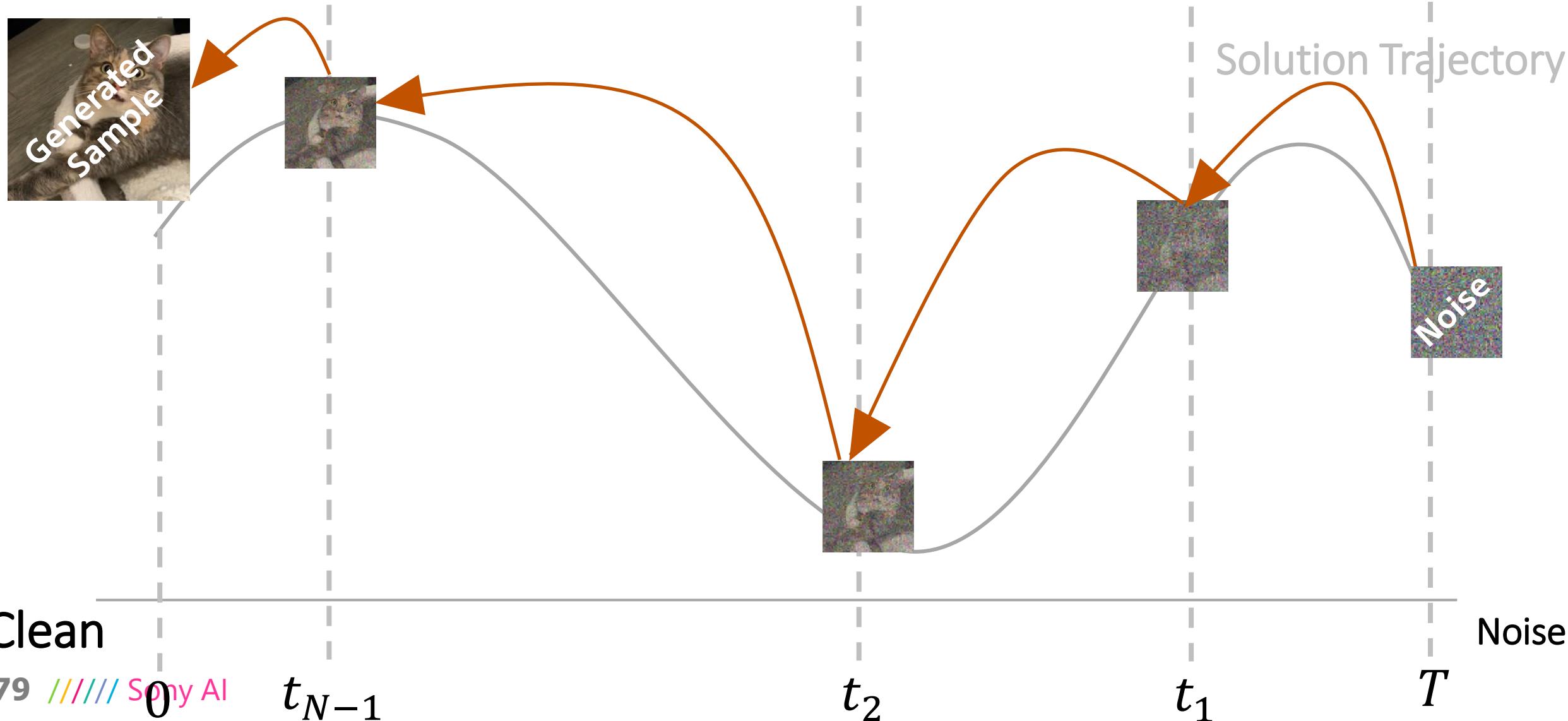
Multistep Generation with CTM: γ -sampling

$\gamma = 0$: Fully Deterministic Sampling



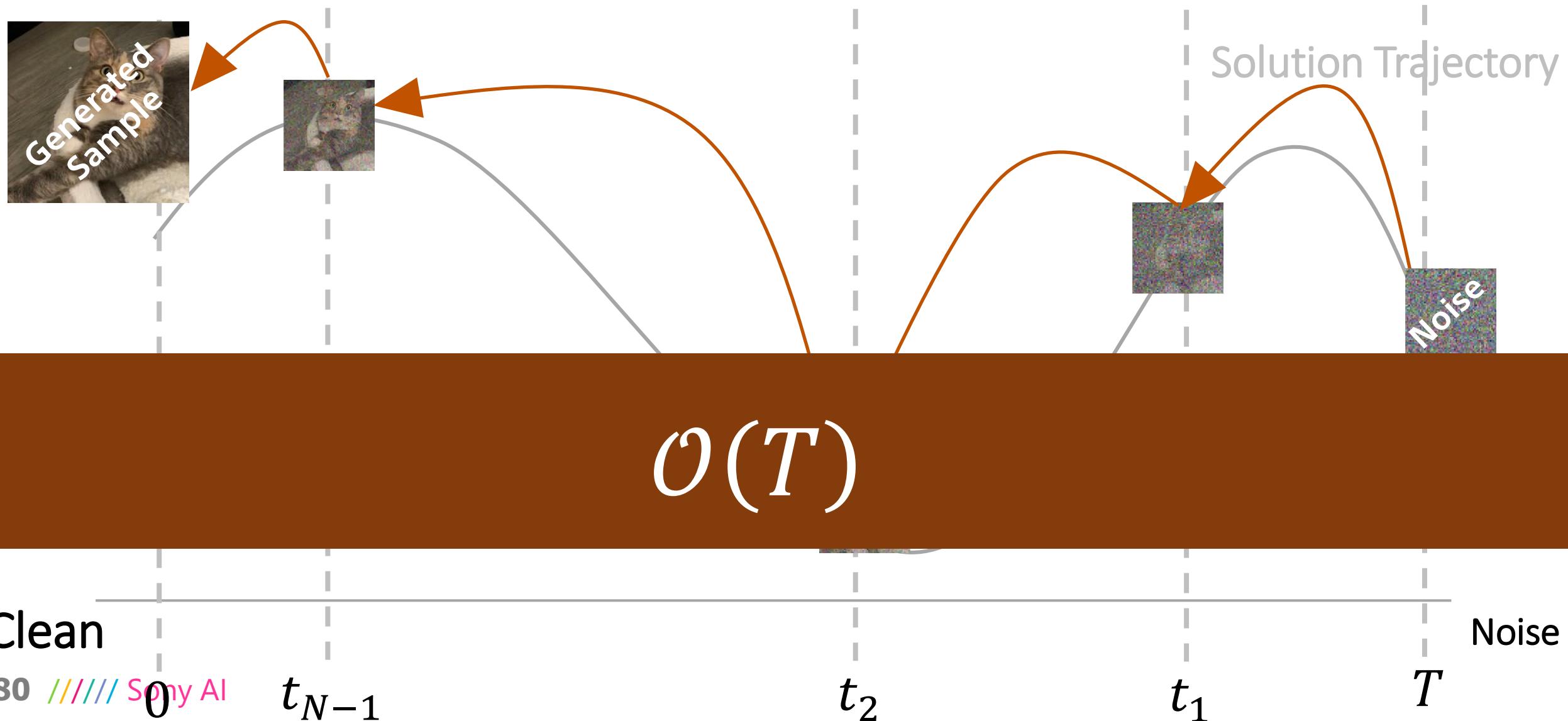
Multistep Generation with CTM: γ -sampling

$\gamma = 0$: Fully Deterministic Sampling



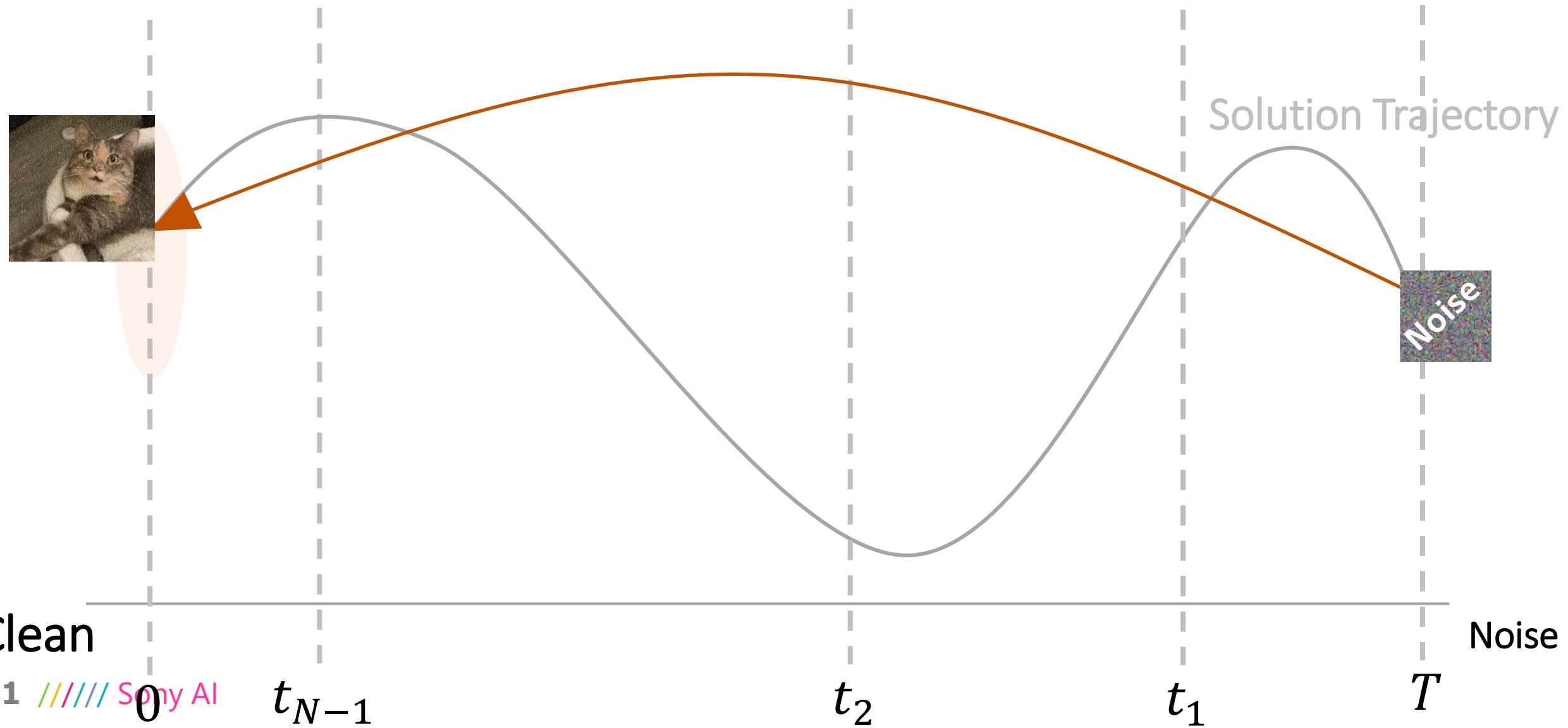
Multistep Generation with CTM: γ -sampling

$\gamma = 0$: Fully Deterministic Sampling



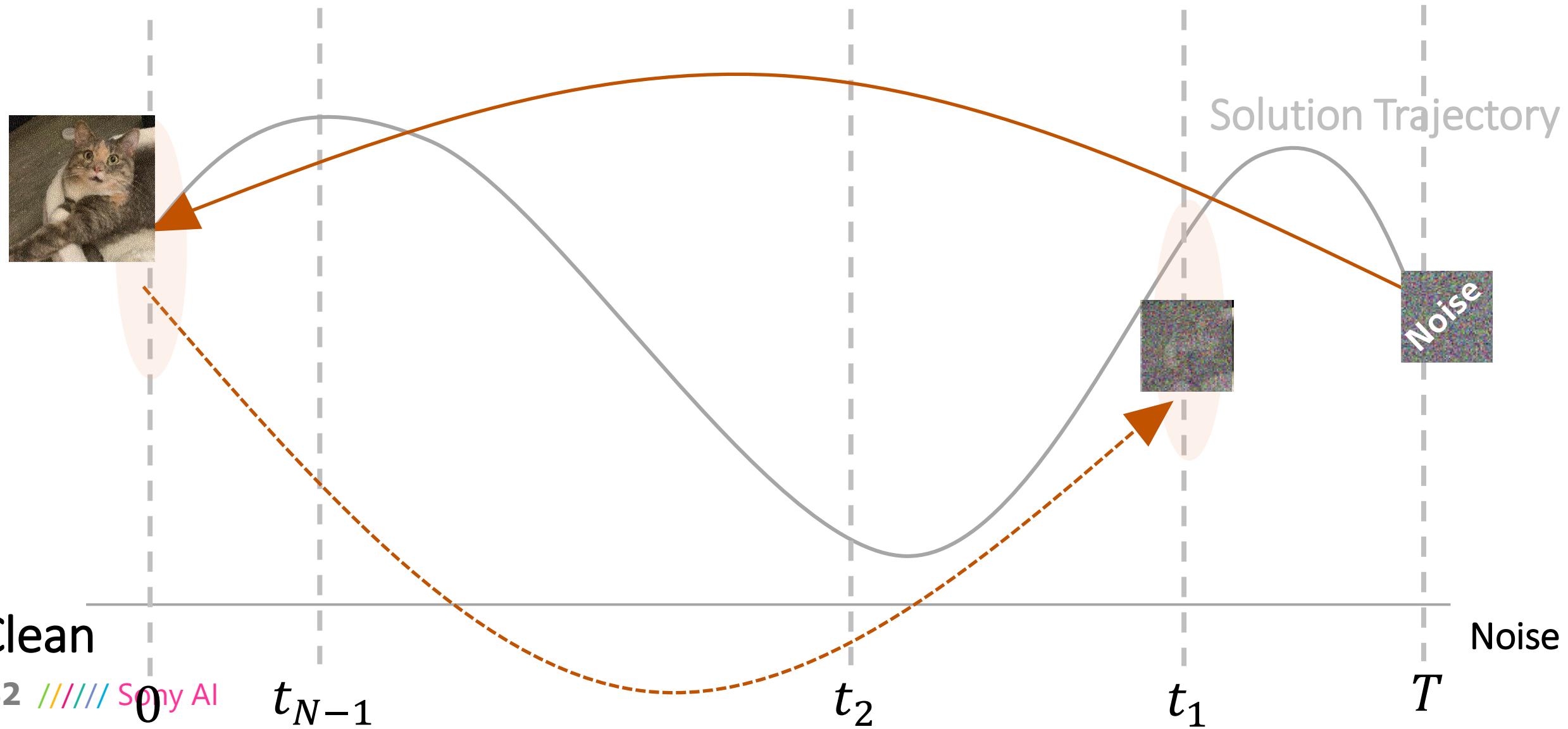
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



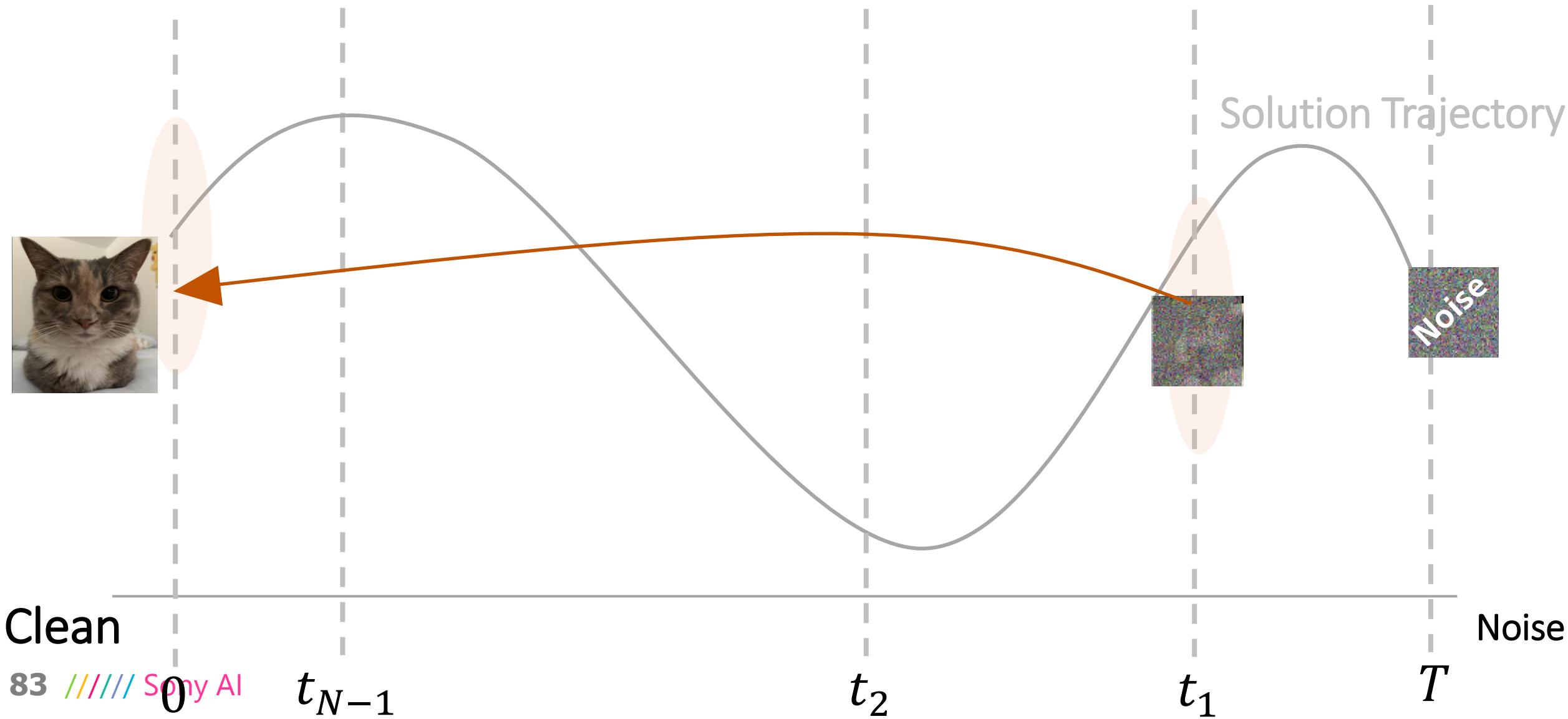
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



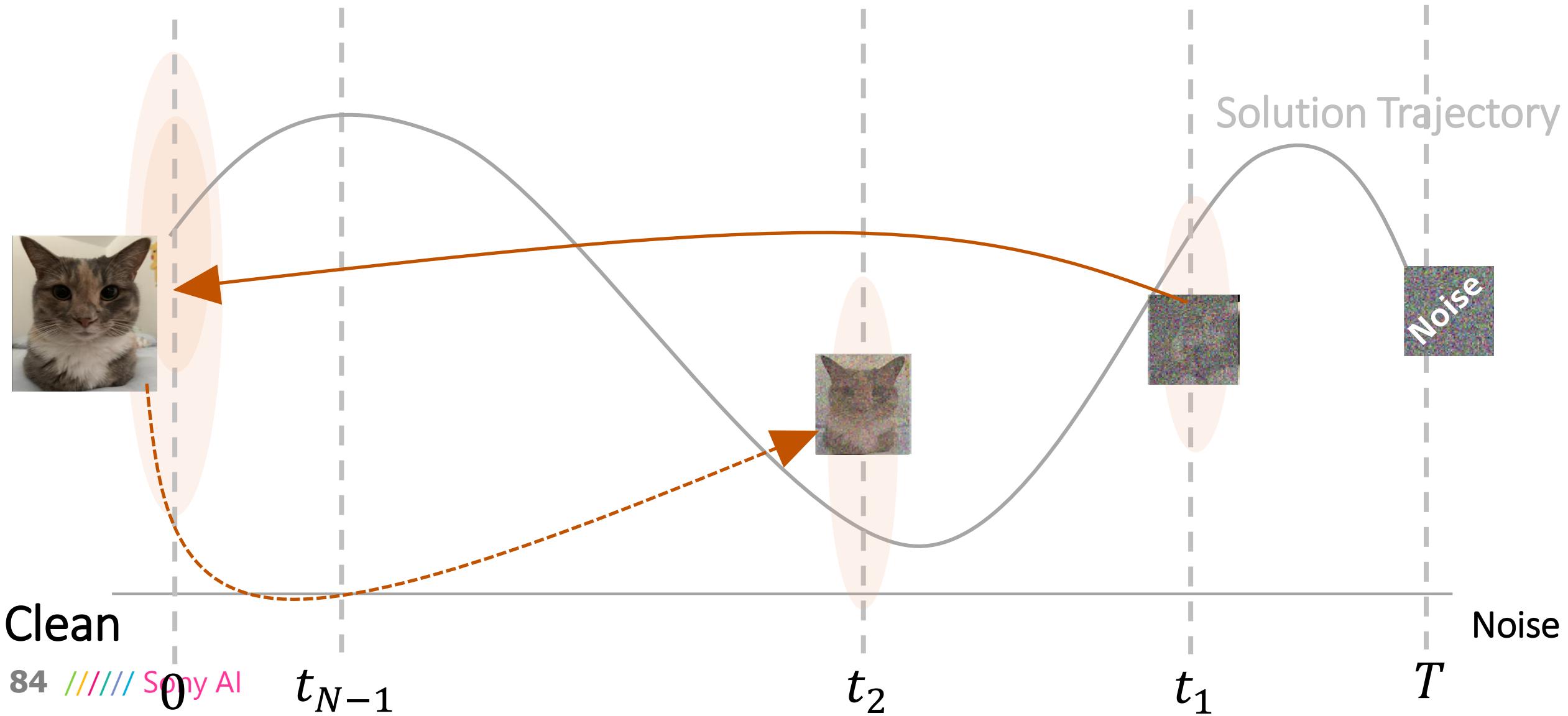
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



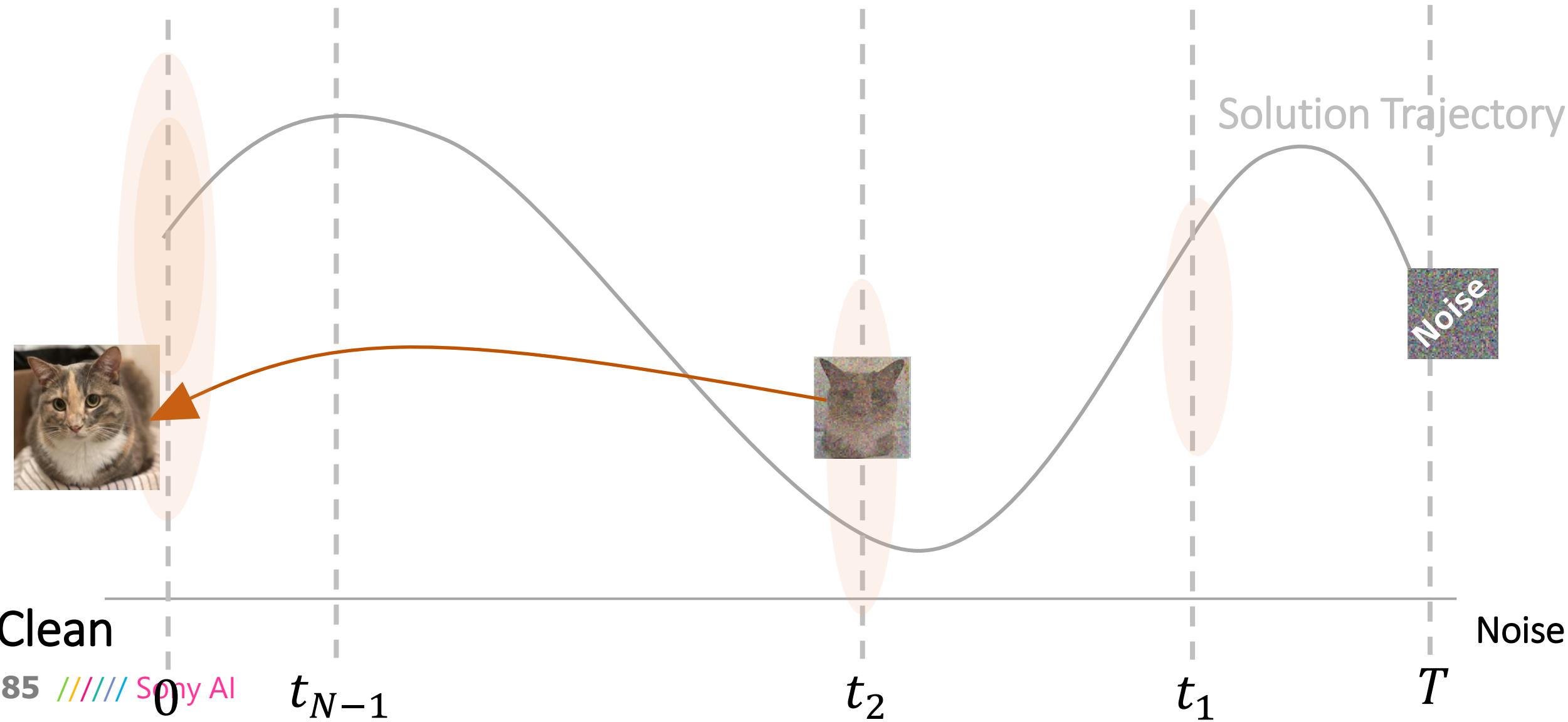
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



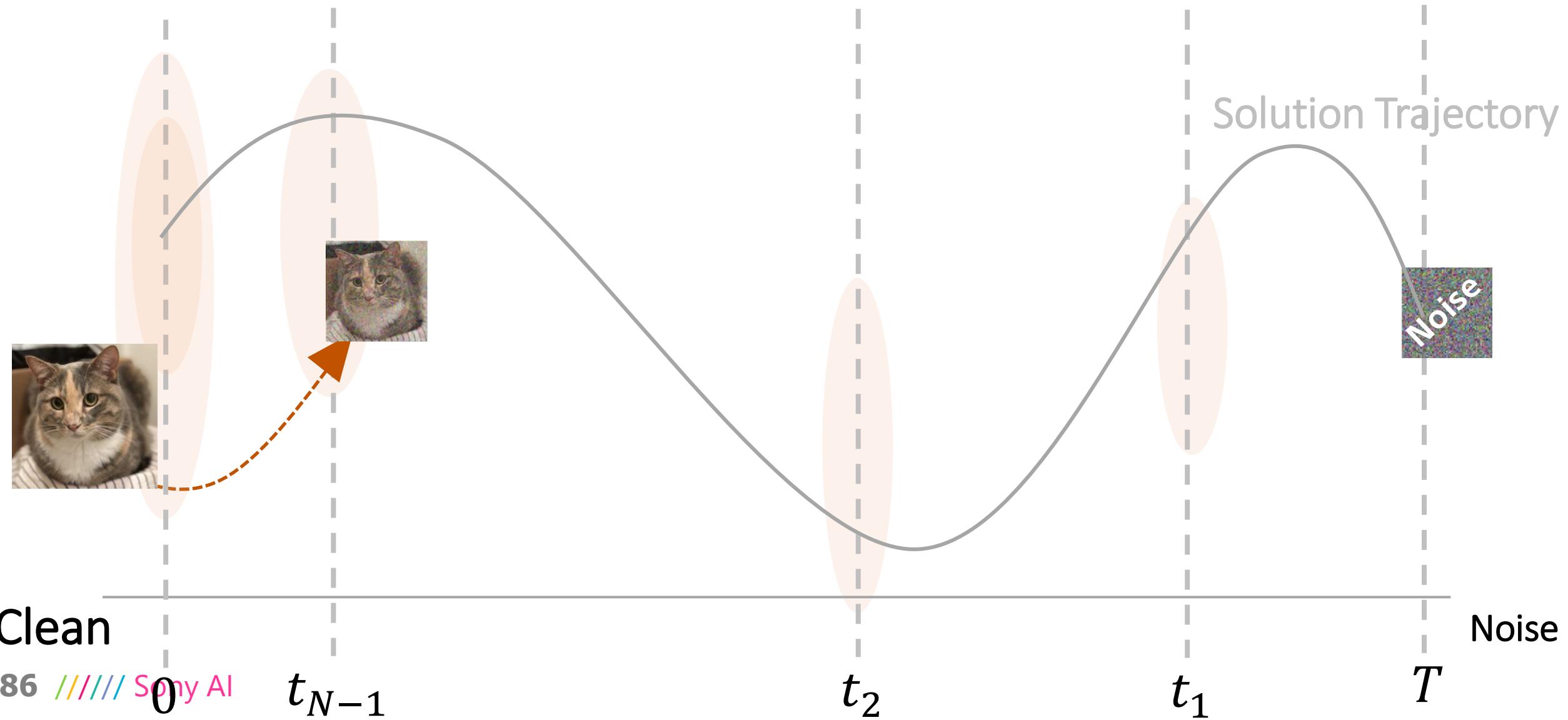
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



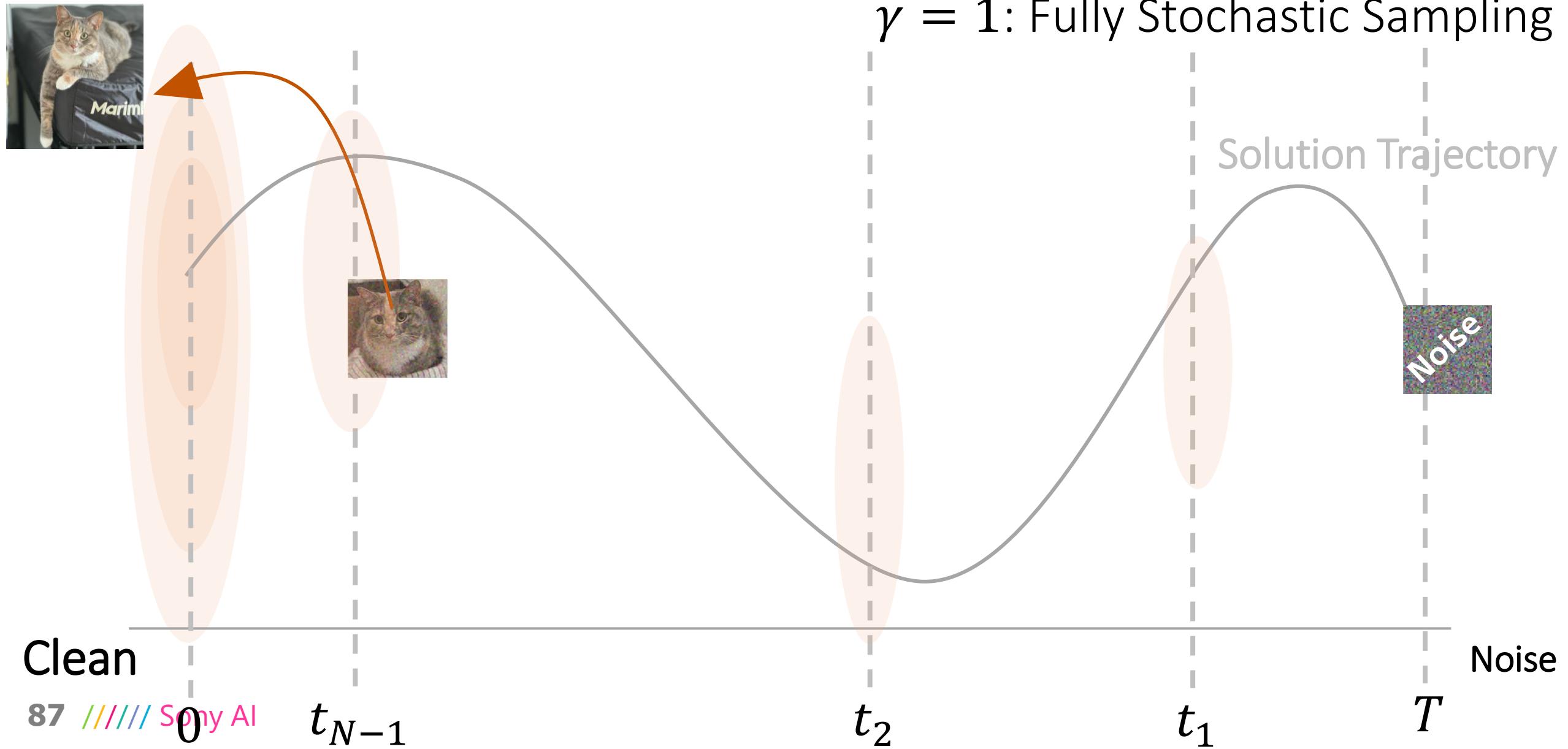
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



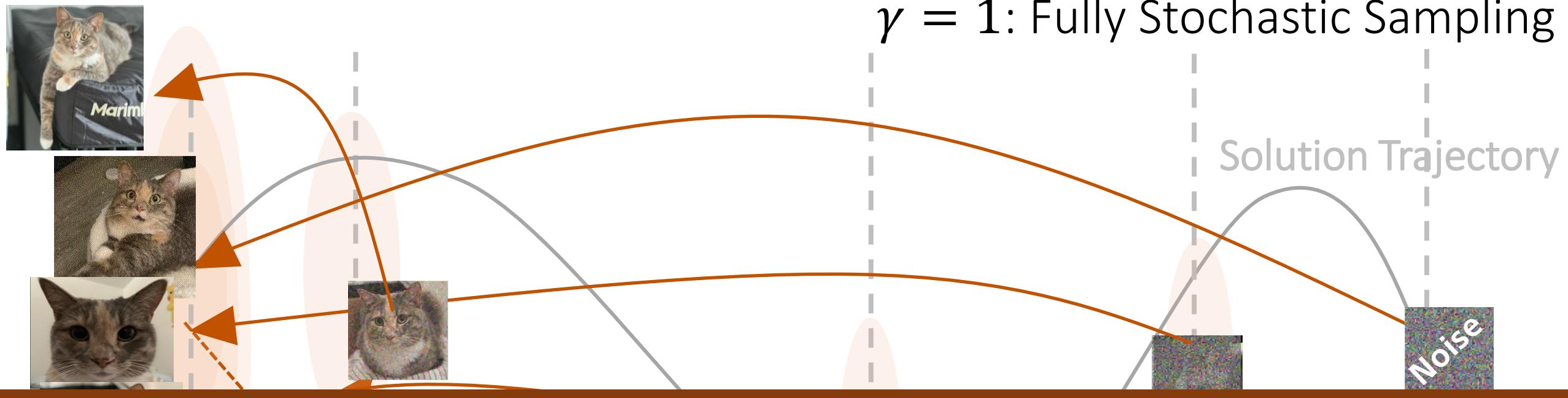
Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling

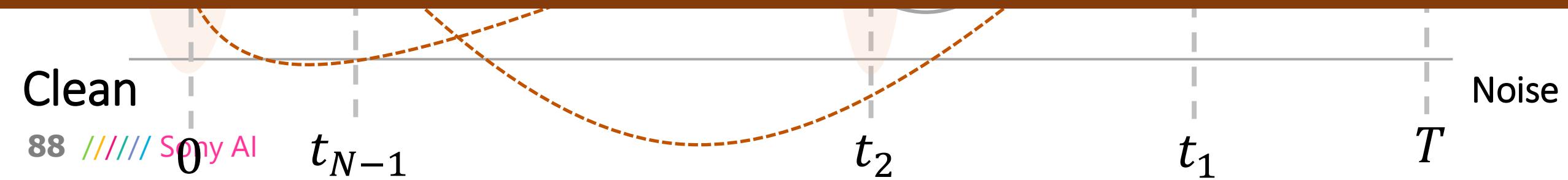


Multistep Generation with CTM: γ -sampling

$\gamma = 1$: Fully Stochastic Sampling



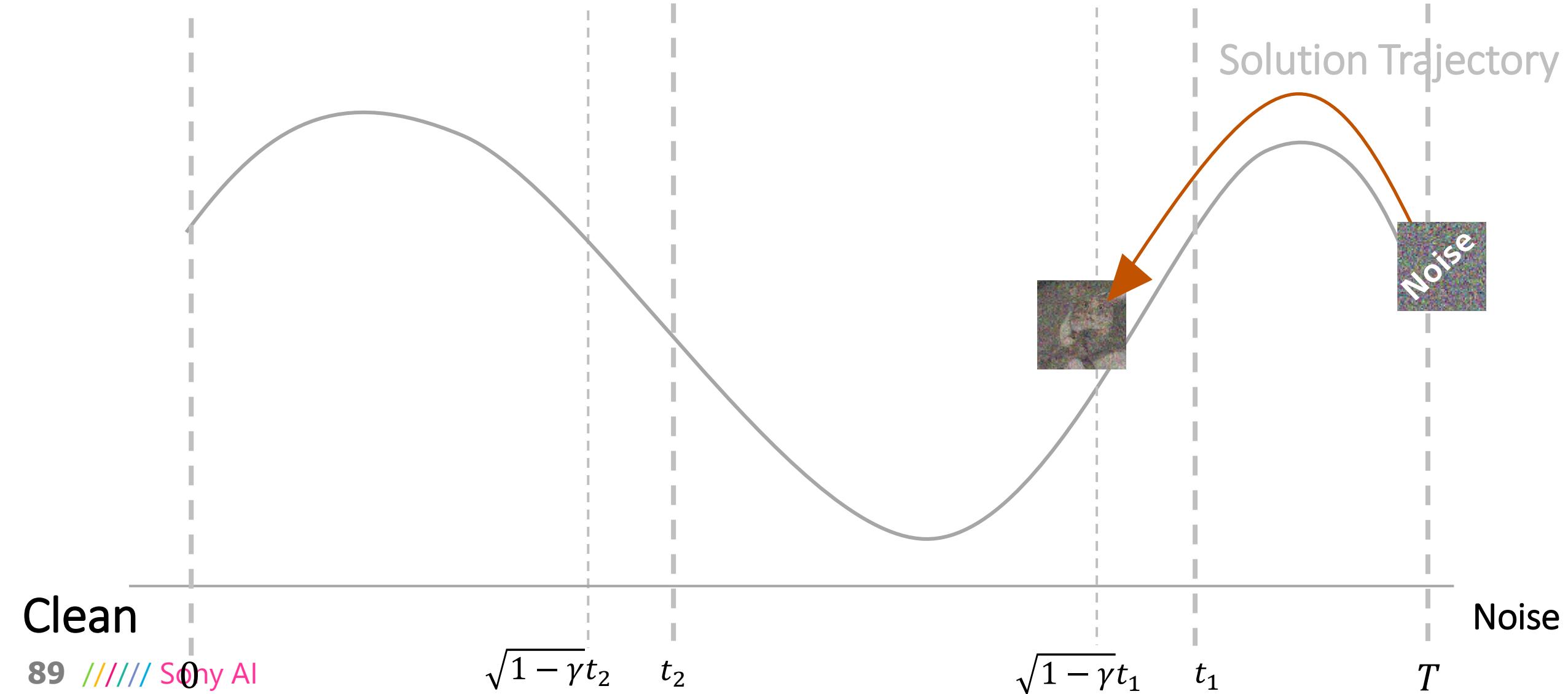
$$\mathcal{O}(T + t_1 + \cdots + t_{N-1})$$



Multistep Generation with CTM: γ -sampling

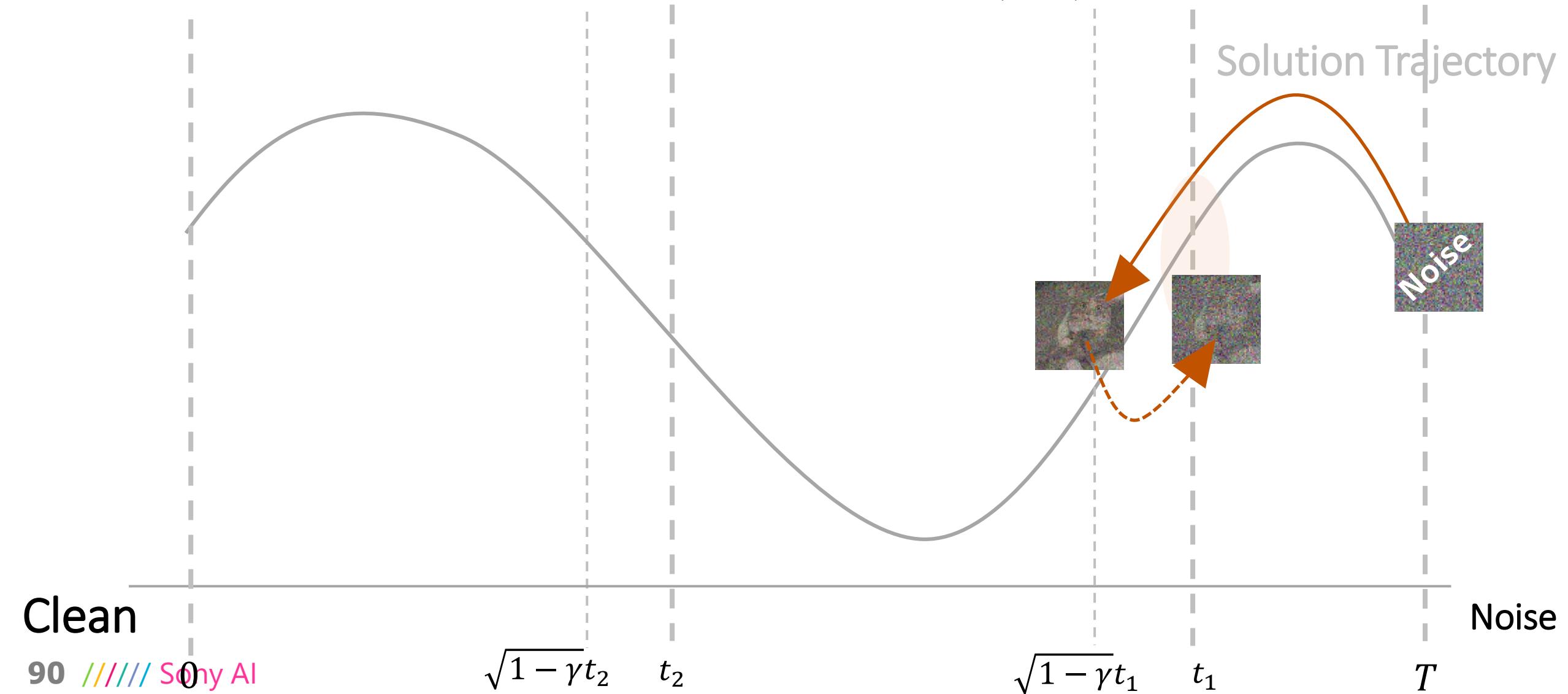
$\gamma \in (0,1)$: Stochastic Sampling

Solution Trajectory



Multistep Generation with CTM: γ -sampling

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



Clean

90 /////////////////////////////////////////////////////////////////// 0ony AI

$\sqrt{1 - \gamma}t_2$

t_2

$\sqrt{1 - \gamma}t_1$

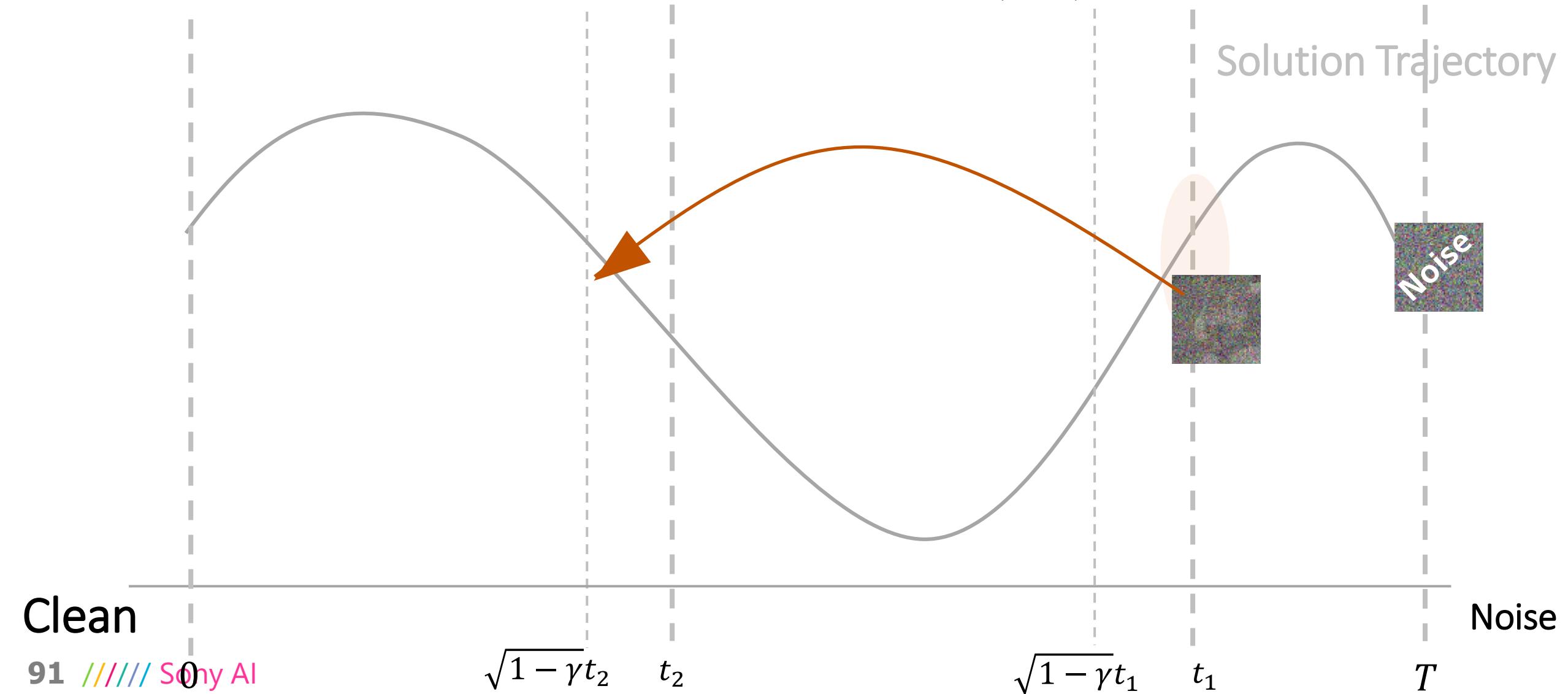
t_1

T

Noise

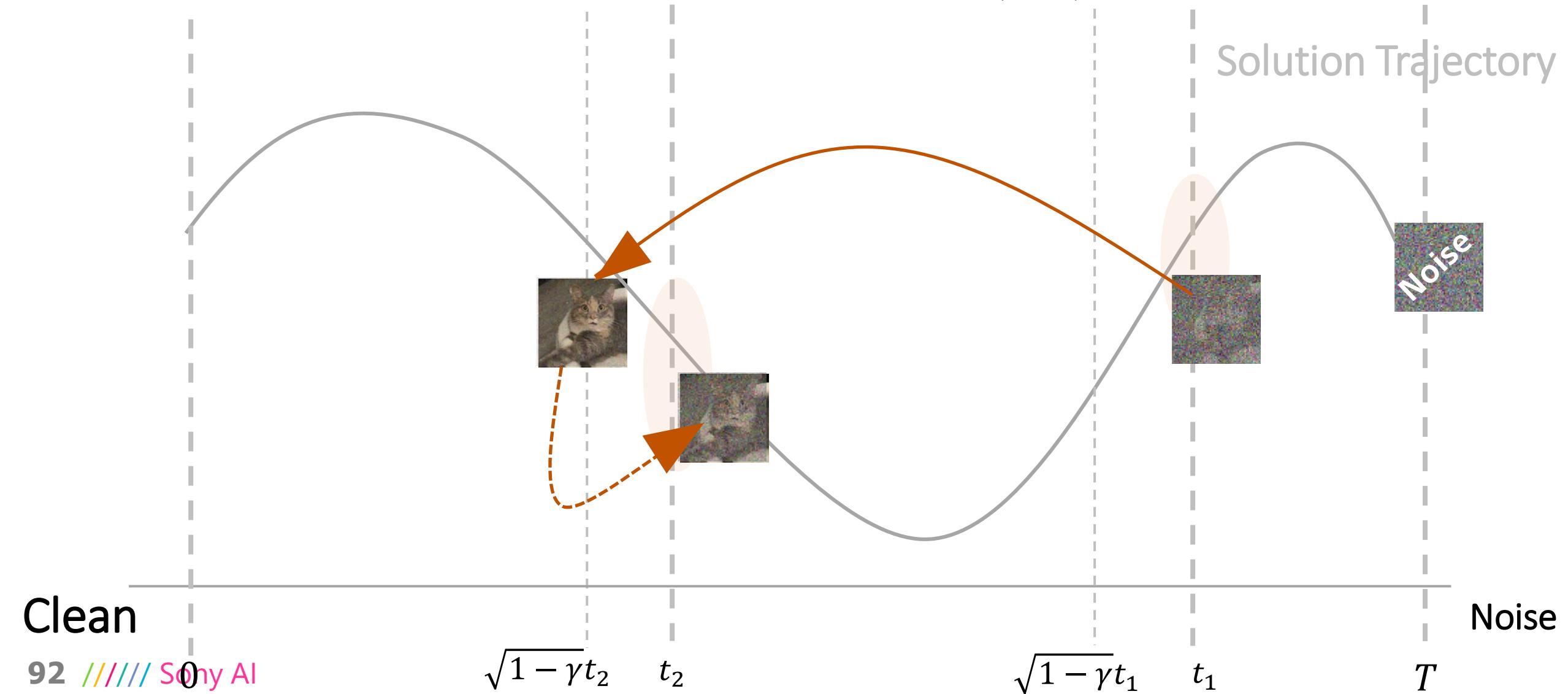
Multistep Generation with CTM: γ -sampling

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



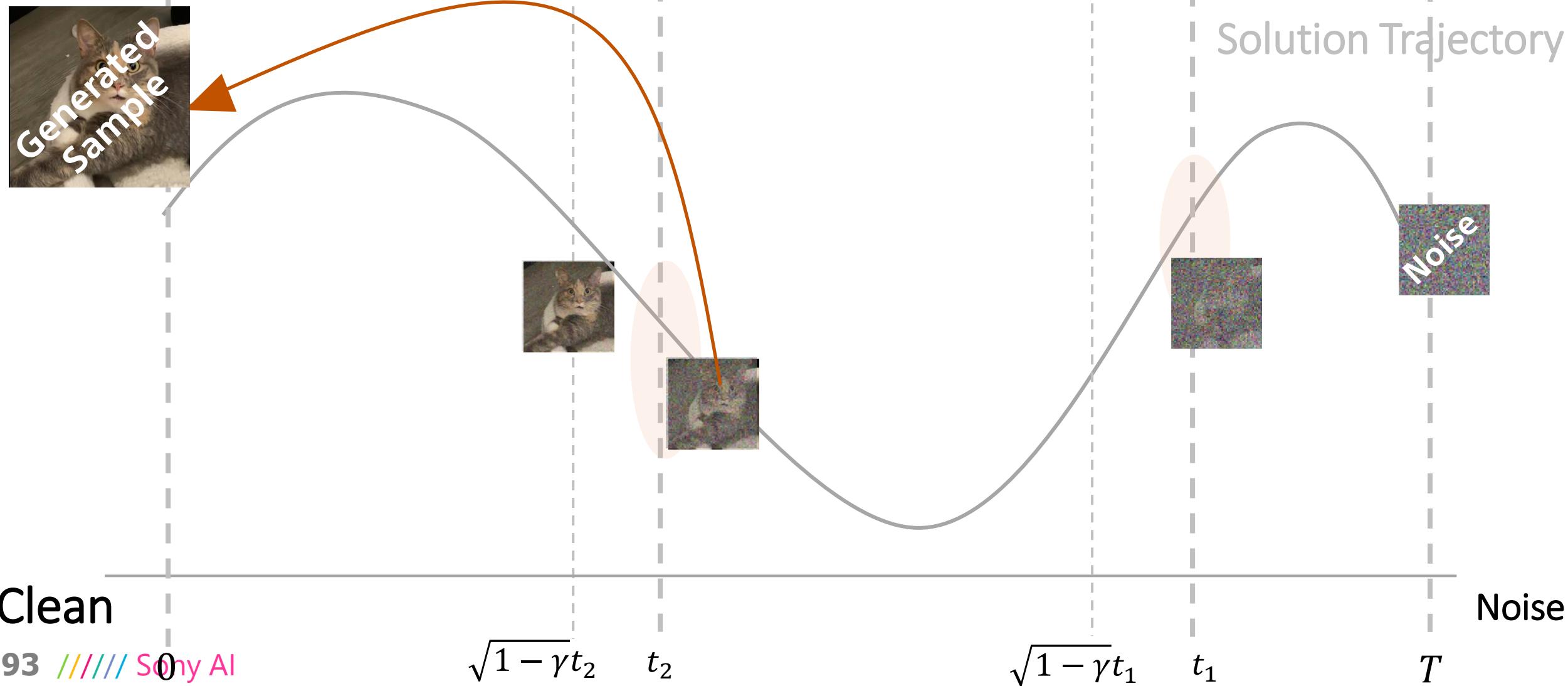
Multistep Generation with CTM: γ -sampling

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



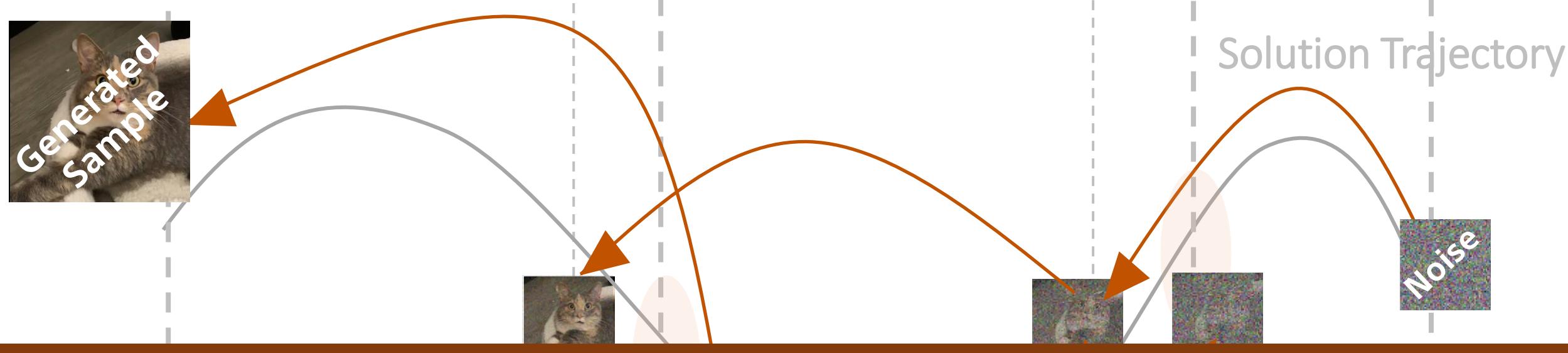
Multistep Generation with CTM: γ -sampling

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



Multistep Generation with CTM: γ -sampling

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



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

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

Dongjun Kim
Stanford University
CA USA

Takashi Shibuya **Chieh-Hsin Lai**
Sony AI Sony AI
Tokyo, Japan 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 Novack¹

Julian McAuley¹

Taylor Berg-Kirkpatrick¹

Nicholas Bryan²

¹ University of California – San Diego

² Adobe Research

znovack@ucsd.edu, njb@ieee.org

SoundCTM's Generation with γ -sampling

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

Koichi Saito¹, Dongjun Kim², Takashi Shibuya¹, Chieh-Hsin Lai¹,
Zhi Zhong³, Yuhta Takida¹, Yuki Mitsufuji^{1,3},

¹Sony AI, ²Stanford University, ³Sony Group Corporation



Paper PDF



arXiv



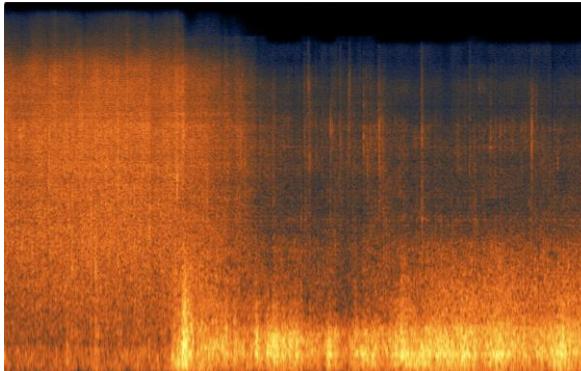
Code

Hugging Face

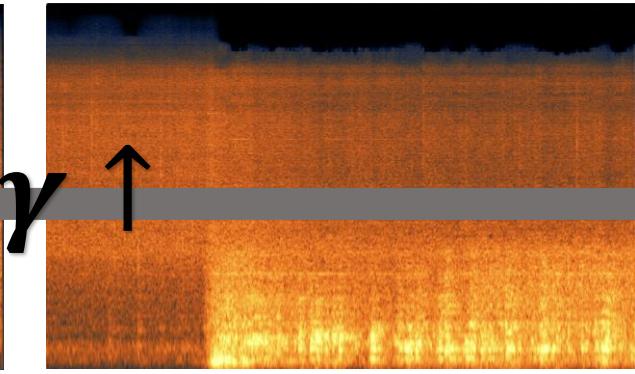
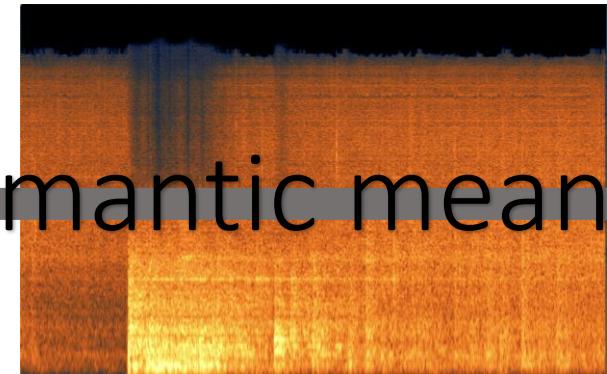
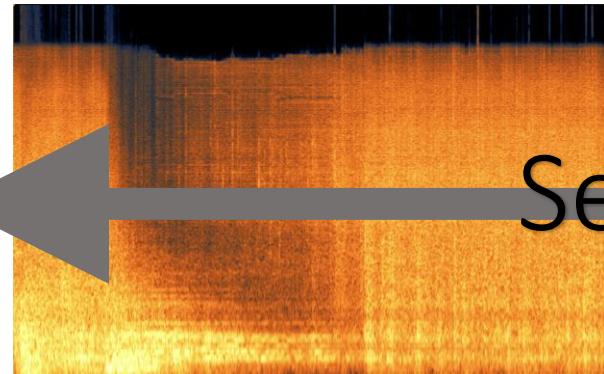
SoundCTM's Generation with γ -sampling

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

One step



16 steps



Semantic meaning varies as $\gamma \uparrow$



$\gamma=1$



$\gamma=0.5$

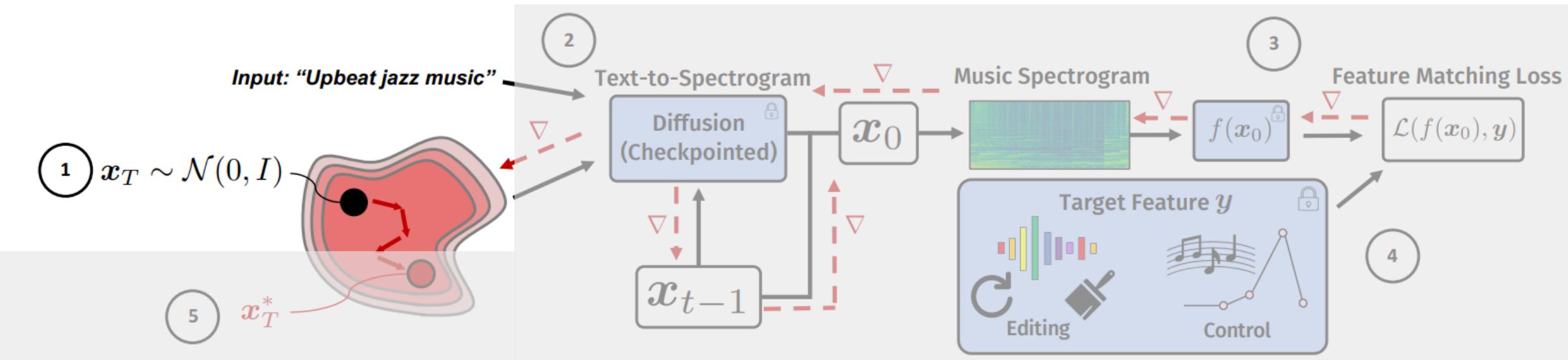


$\gamma=0.25$



$\gamma=0$

DITTO-2: CTM-Based for Controllable Text-to-Music



DITTO:
Diffusion Inference-Time T-Optimization
for Music Generation

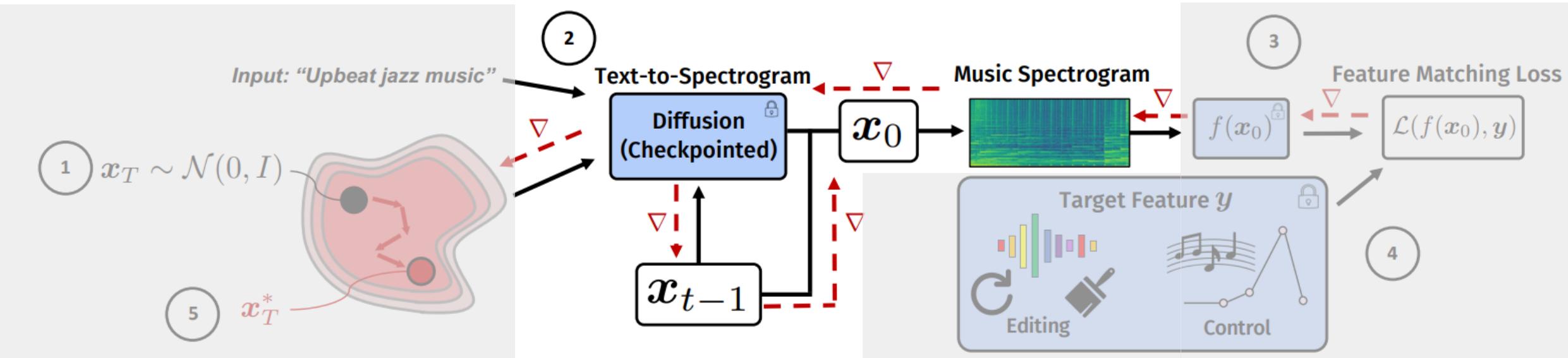
Zachary Novack^{1,2*} Julian McAuley¹ Taylor Berg-Kirkpatrick¹ Nicholas J. Bryan²

¹University of California, San Diego

²Adobe Research

*Work done during an internship at Adobe Research

DITTO-2: CTM-Based for Controllable Text-to-Music



DITTO:
Diffusion Inference-Time T-Optimization
for Music Generation

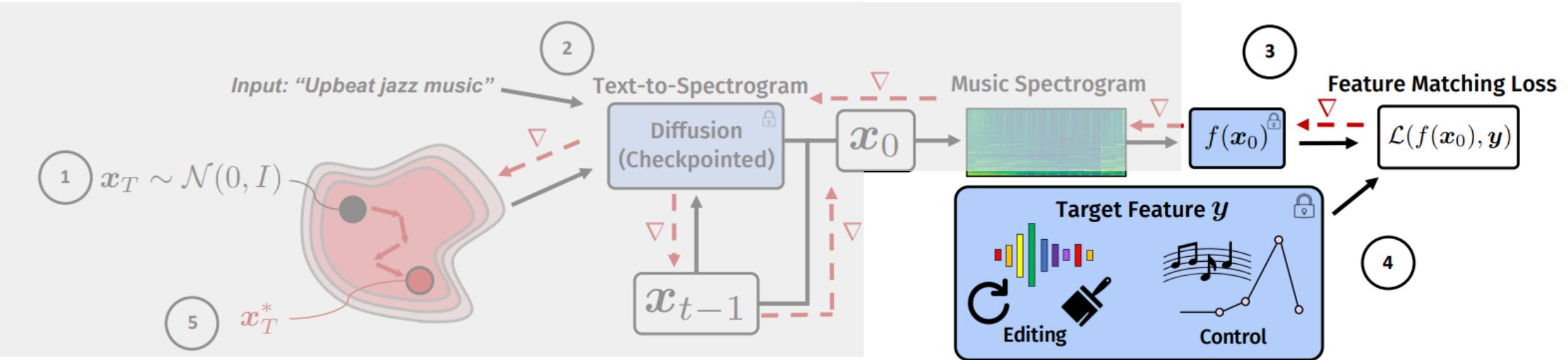
Zachary Novack^{1,2*} Julian McAuley¹ Taylor Berg-Kirkpatrick¹ Nicholas J. Bryan²

¹University of California, San Diego

²Adobe Research

*Work done during an internship at Adobe Research

DITTO-2: CTM-Based for Controllable Text-to-Music



DITTO:
Diffusion Inference-Time T-Optimization
for Music Generation

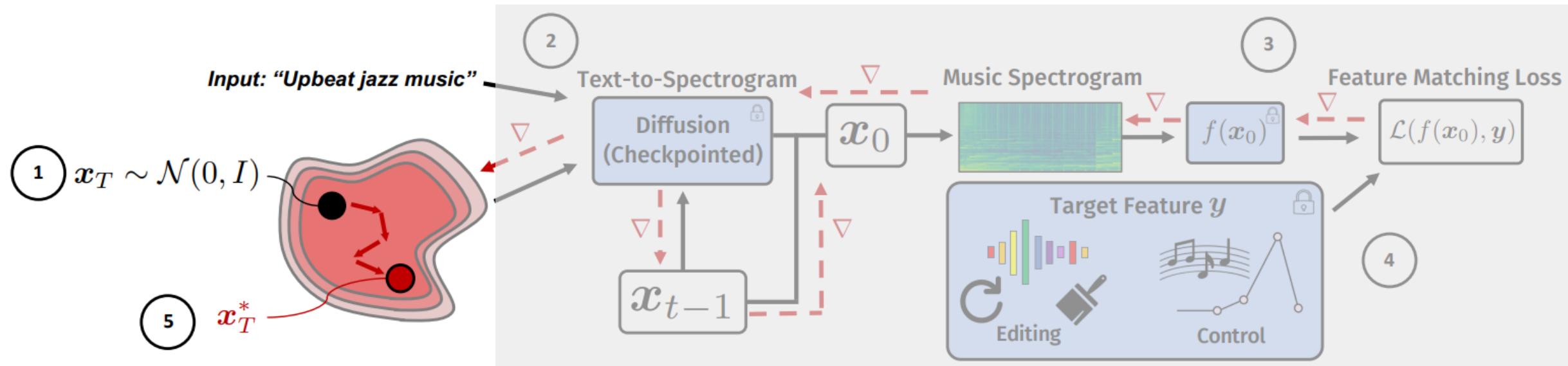
Zachary Novack^{1,2*} Julian McAuley¹ Taylor Berg-Kirkpatrick¹ Nicholas J. Bryan²

¹University of California, San Diego

²Adobe Research

*Work done during an internship at Adobe Research

DITTO-2: CTM-Based for Controllable Text-to-Music



DITTO:

**Diffusion Inference-Time T-Optimization
for Music Generation**

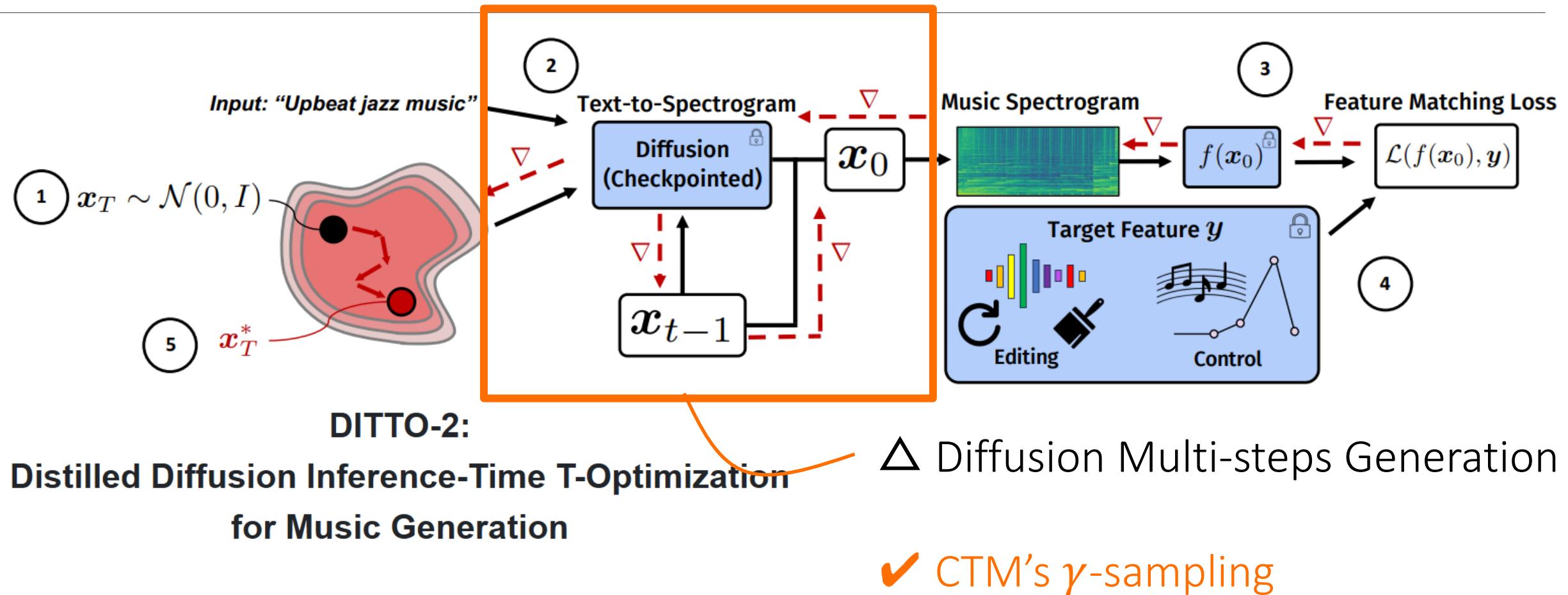
Zachary Novack^{1,2*} Julian McAuley¹ Taylor Berg-Kirkpatrick¹ Nicholas J. Bryan²

¹University of California, San Diego

²Adobe Research

*Work done during an internship at Adobe Research

DITTO-2: CTM-Based for Controllable Text-to-Music



Zachary Novack¹ Julian McAuley¹ Taylor Berg-Kirkpatrick¹ Nicholas J. Bryan²

¹University of California, San Diego

²Adobe Research

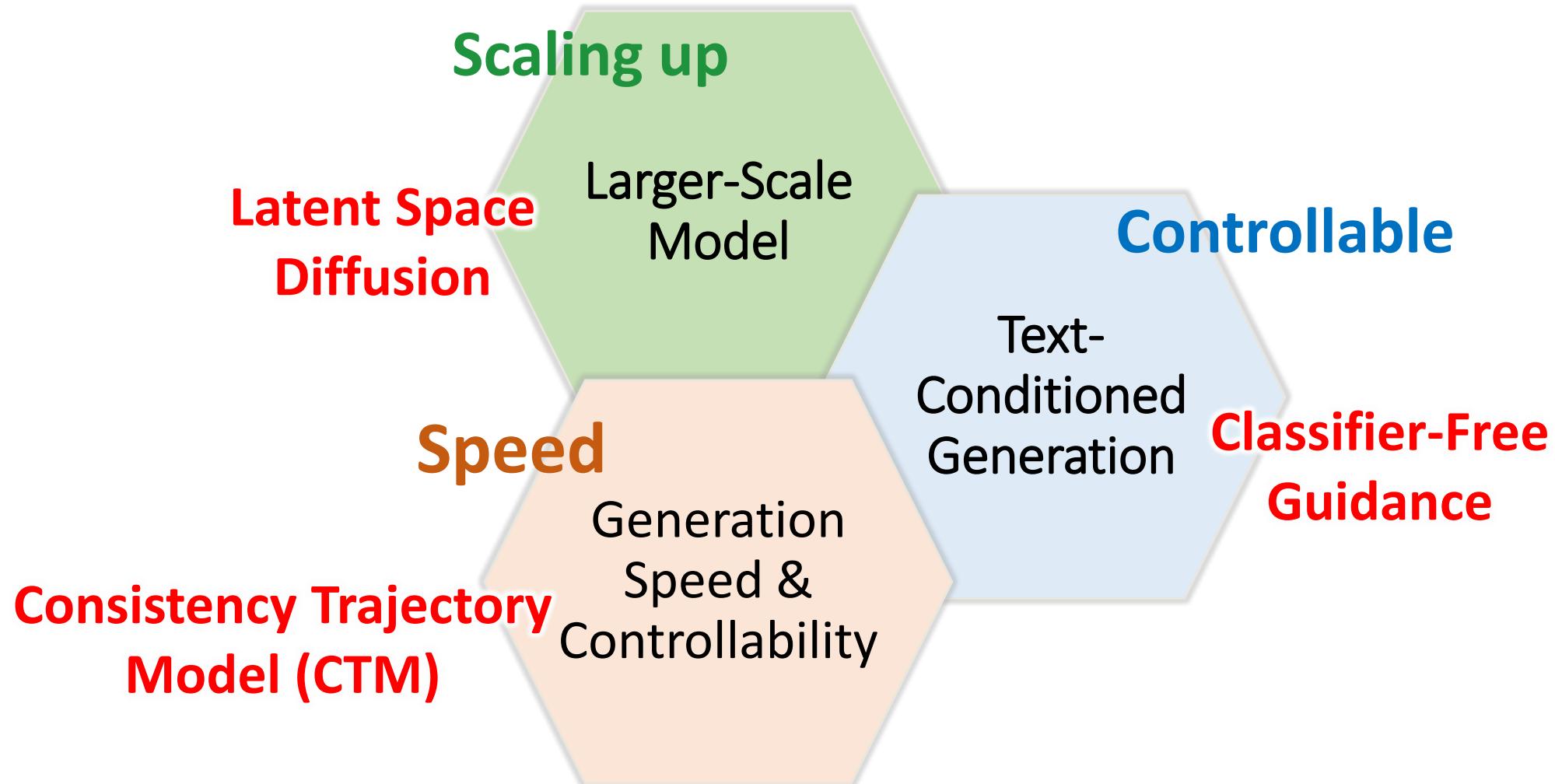


Paper



HF paper

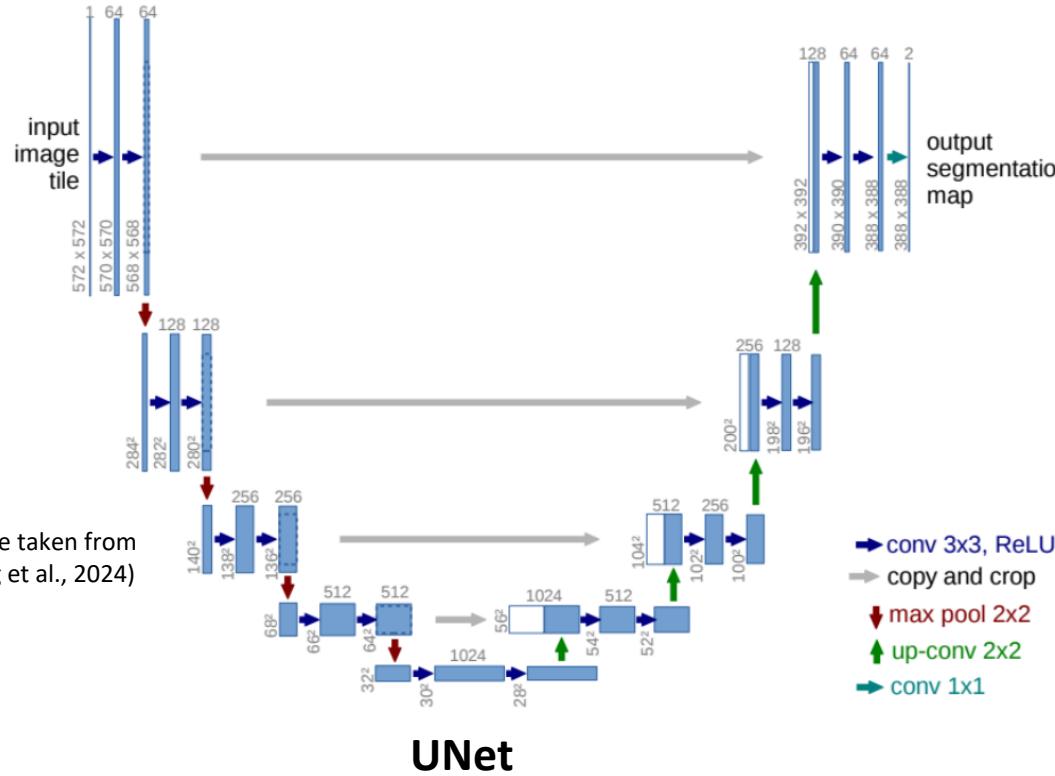
Making Diffusion Model More Powerful



PRACTICAL IMPLEMENTATIONS

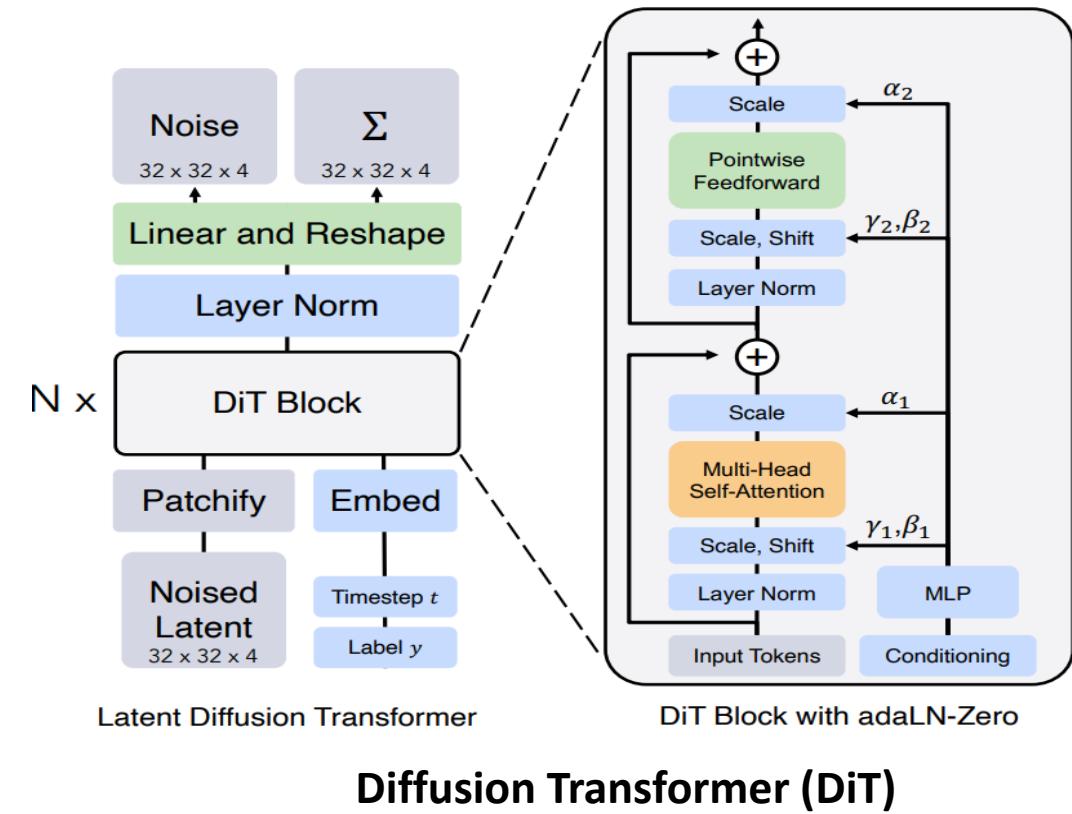
Network Architectures

(Figure taken from
Zhang et al., 2024)



UNet

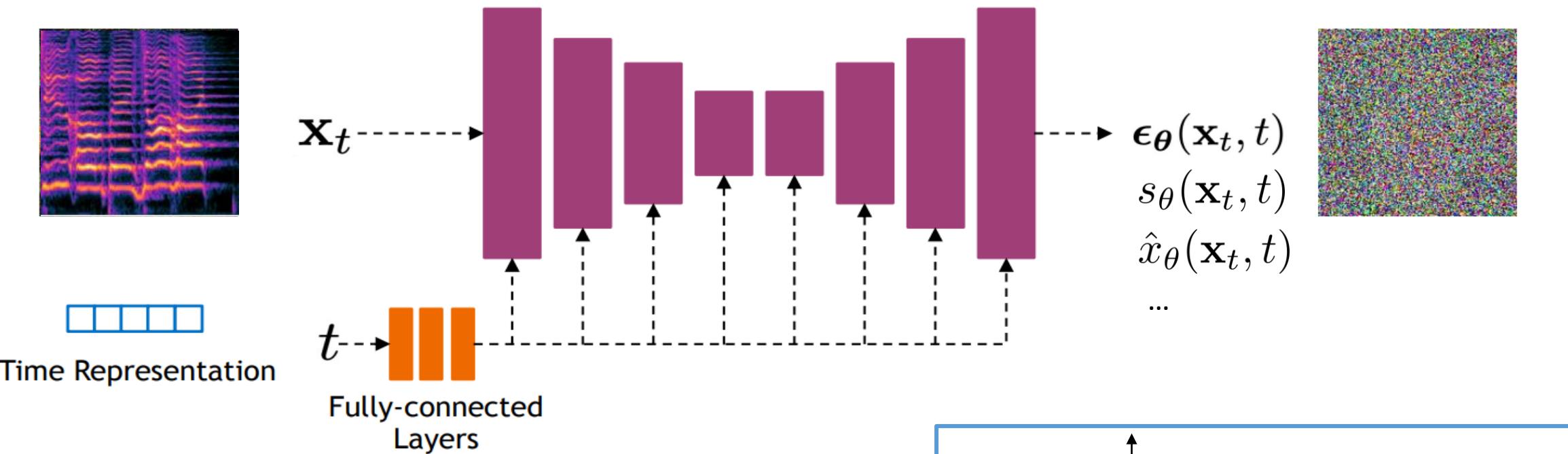
- Downsampling and Upsampling stacks
- ResNet blocks and self-attention layers



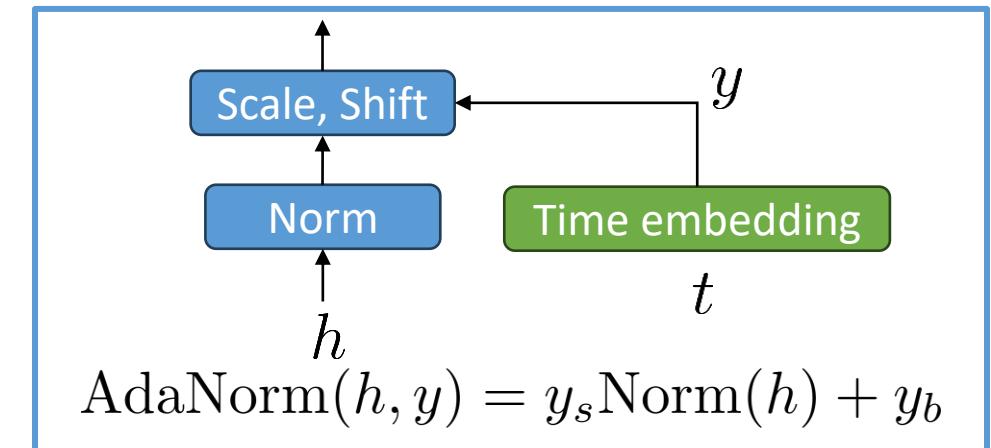
Diffusion Transformer (DiT)

- Transformer architecture
- AdaLN-Zero blocks

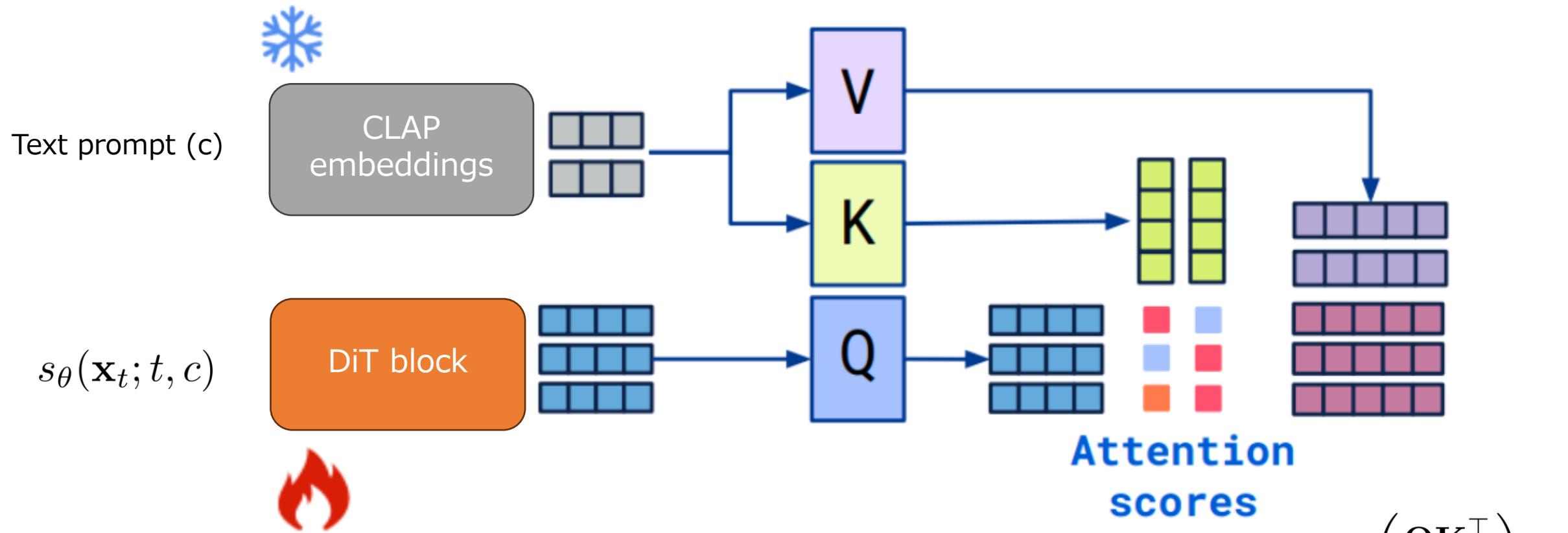
Timestep Embedding Conditioning



- Time embeddings are sinusoidal positional embeddings or random Fourier features
- Time embeddings can be concatenated, added or through adaptative normalization layers



Cross-Attention Conditioning



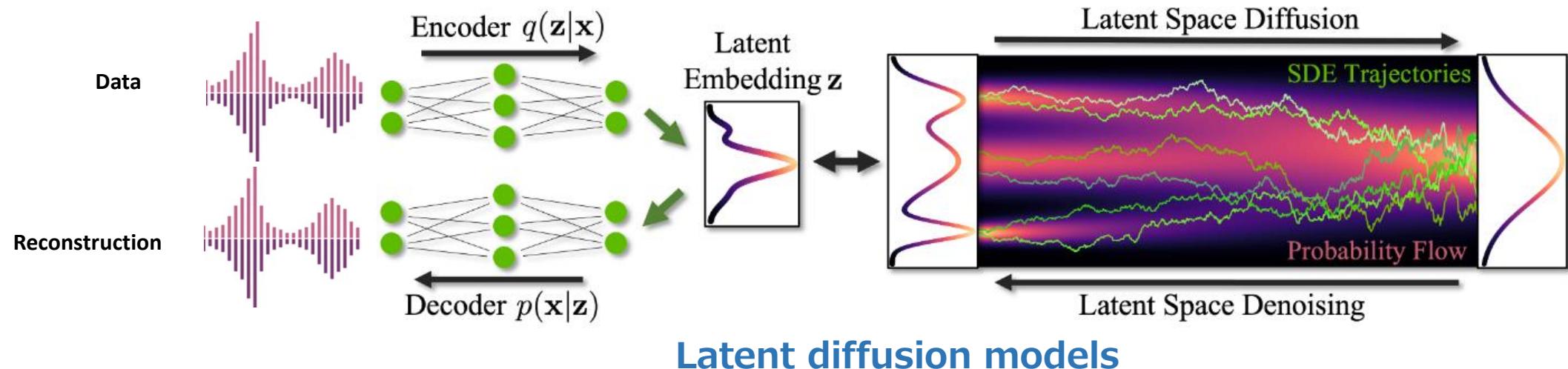
- Text prompts are embedded with CLAP, T5, ...
 - Using cross attention to condition the text prompt

$$\text{Cross-Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}$$

Latent Diffusion Models

- Music generation requires modelling long-range sequences
 - Using the full frequency spectrum (48 kHz vs 16 kHz)
- Music contains harmonies and melodies from different instruments, create complex structures
 - Human listeners are highly sensitive to disharmony

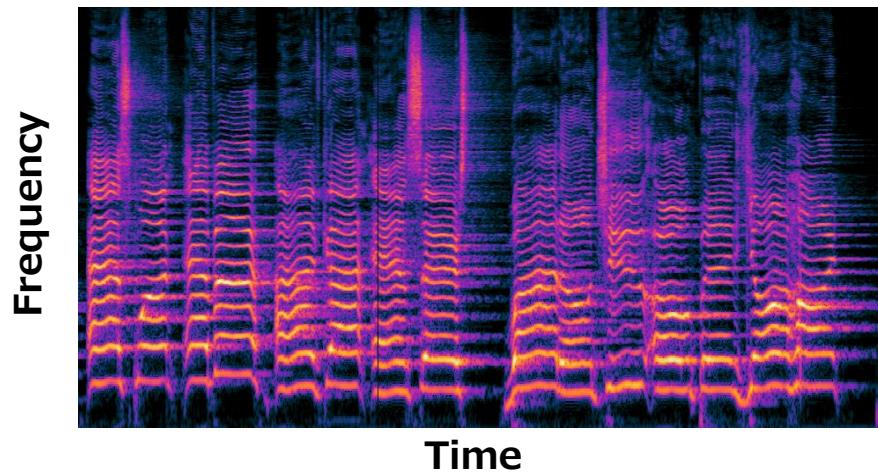
Compress the audio with a pre-trained auto-encoder



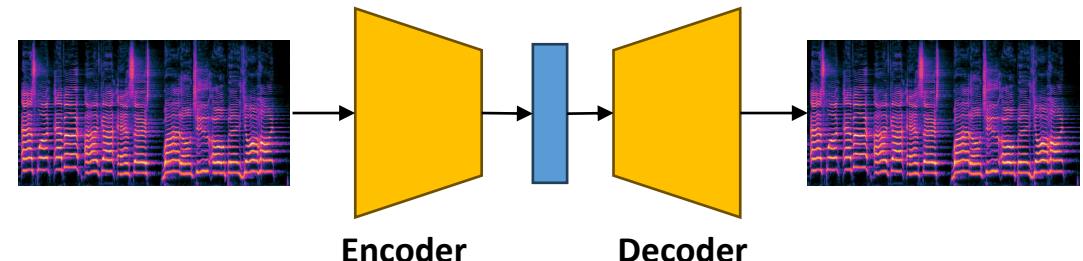
Encoder-Decoder

➤ **Audio spectrograms**

- Can be seen as an image, but requires special consideration: temporal structure, data augmentation, frequencies are non-locally distributed, …)

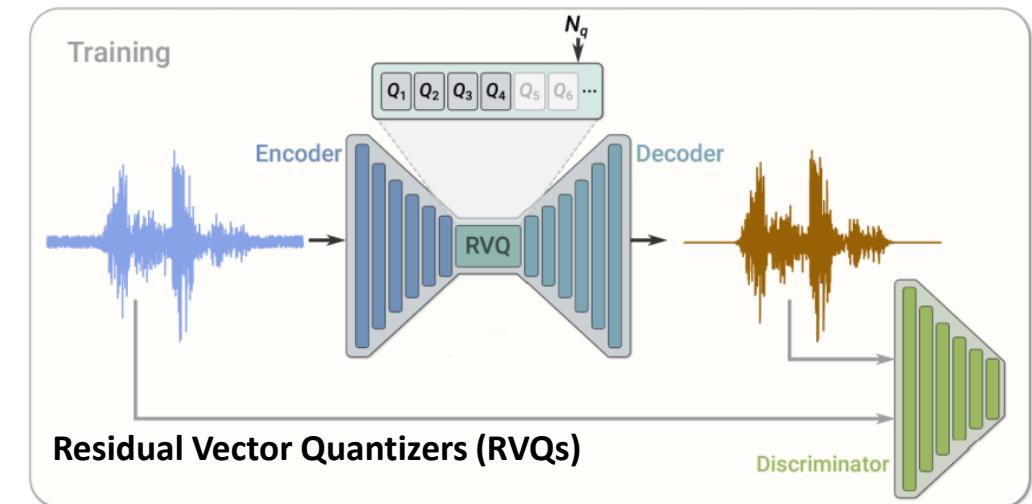


➤ **Variational Auto-Encoders (VAE)**



➤ Neural Audio Codecs and RVQs

- SoundStream, EncoDec, DAC, ...



90x compression

[Kumar et al., NeurIPS 2023] High-Fidelity Audio Compression with Improved RVQGAN
[Zeghidour et al., Trans. Audio, 2021] SoundStream: An End-to-End Neural Audio Codec
[Défossez et al., TMLR 2024] High Fidelity Neural Audio Compression

APPLICATIONS

- Controllable generation
- Inverse problem

CONTROLLABLE GENERATION

- **Text-to-Music Generation**
- **Finetuning**
- **Inference-Time Optimization**
- **Inference-Time Guidance**

Text-to-Music Generation (I)

➤ Problem

- A generative model that generates music in response to text prompts

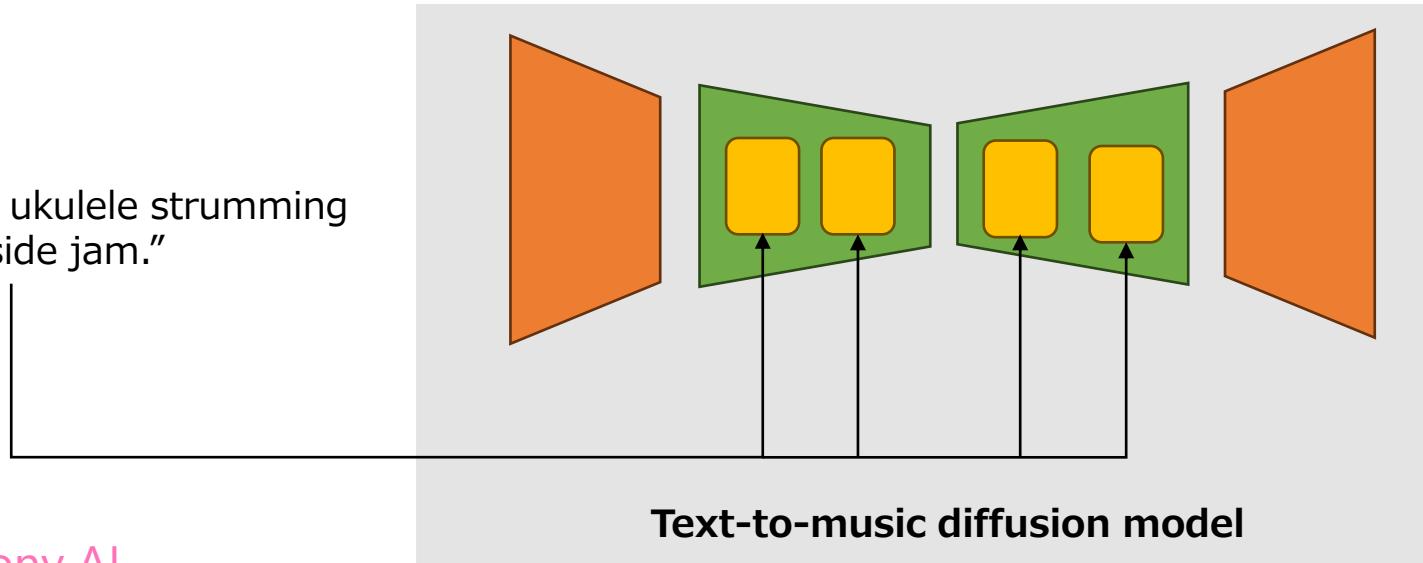
➤ Applications

- Music composition and soundtrack creation
 - Assistance for musicians
 - Music editing and personalized

➤ Challenges

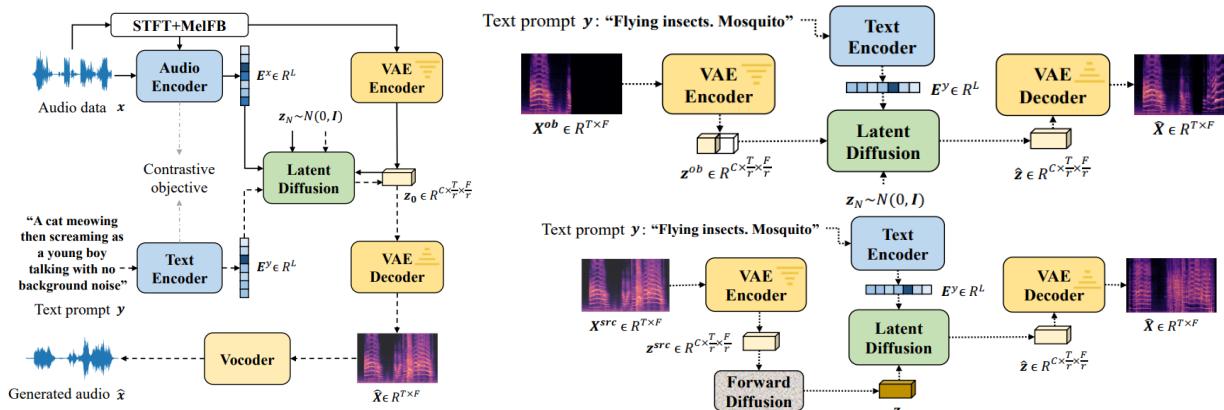
- Modelling long sequence
 - Sound quality
 - A large amount of supervised data

“A cheerful ukulele strumming
in a beachside jam.”

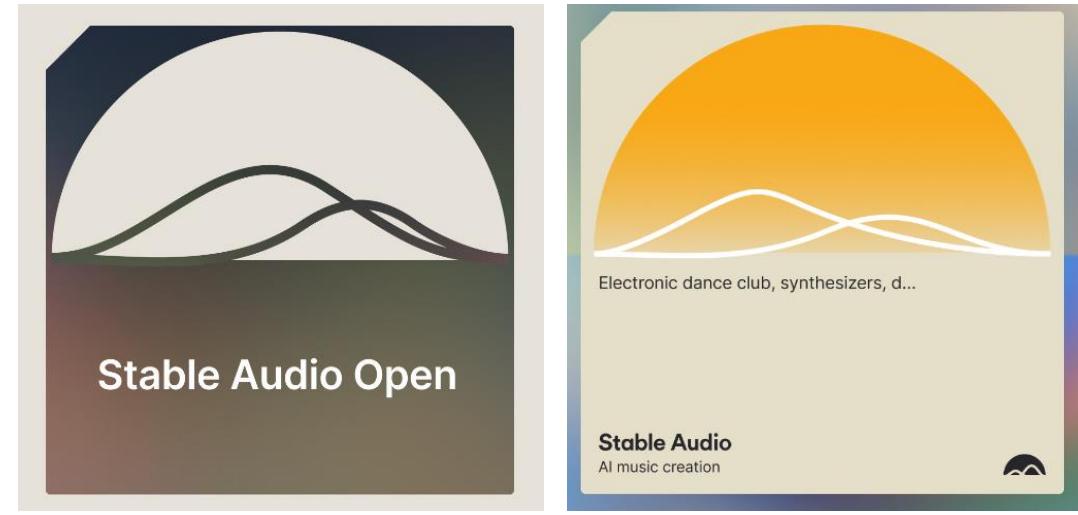


Text-to-Music Generation (II)

AudioLDM: Text-to-Audio Generation with Latent Diffusion Models



<https://github.com/haoheliu/AudioLDM>

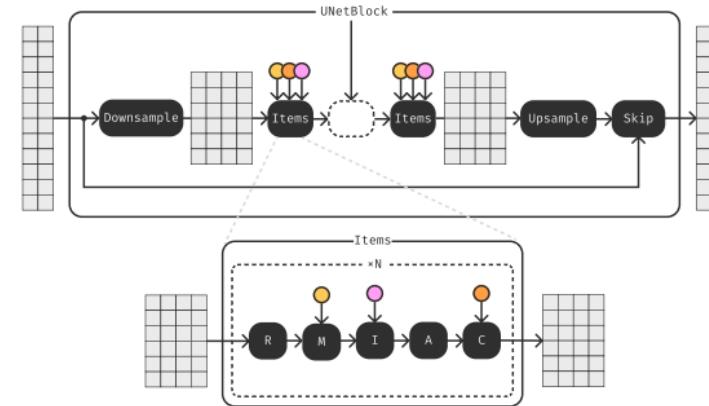


<https://github.com/Stability-AI/stable-audio-tools>



113 ////////////////////////////////////////////////////////////////// Sony AI

Moûsai: Efficient Text-to-Music Diffusion Models



<https://github.com/archinetai/audio-diffusion-pytorch>

Text-to-Music Generation (III)

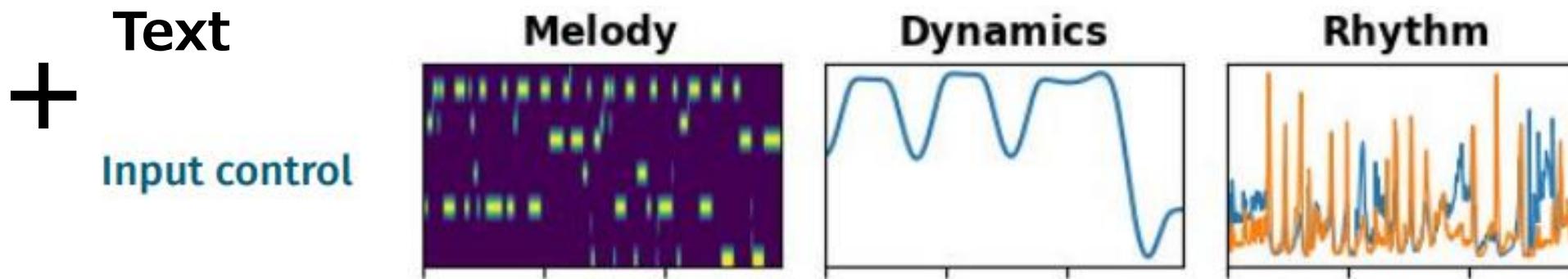
Text input: A traditional Irish fiddle playing a lively reel.

Up Next: The sound of a light saber

Text-to-Music Generation (IV)

➤ Text-to-Music generation models can produce high-fidelity music audio, but they **lack fine-grained controllability**

- Text descriptions are insufficient for detailed controls
- Suitable for global attributes like genre, mood, tempo, ⋯
- Not suitable for fine-grained controls over time-varying attributes



How to enable control with additional information alongside the text description?

Controllable Generation (I)

- Learning conditional controls in an end-to-end way is challenging
 - Limited training data for a specific condition
 - Training cost is high

long audio files. Finally, we have in total 3,302,553 ten-seconds audio samples for model training. It should be noted that even if some datasets, e.g., AudioCaps and BBC SFX, have text captions for the audio, we do not utilize them during the training of LDMs. We only use the audio samples for training. We resample all the datasets into 16kHz sampling rate and mono format, and all samples are padded to 10 seconds.

AudioLDM

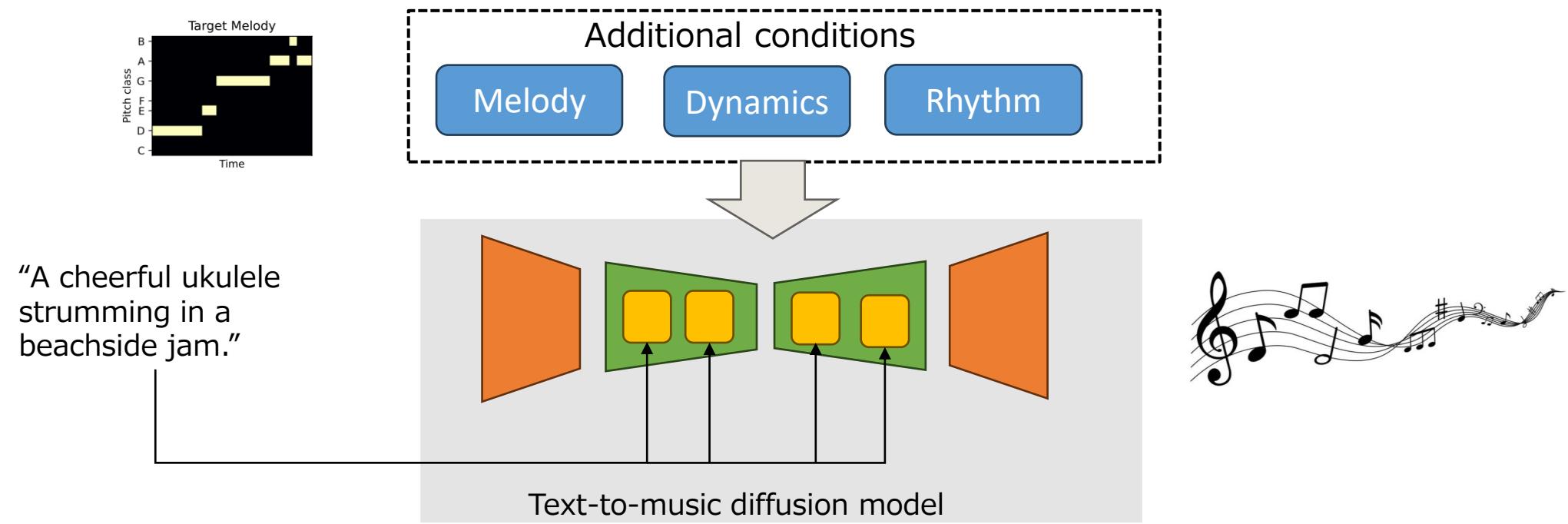
Stability AI used a cluster of 4,000 Nvidia A100 GPUs running in AWS to train Stable Diffusion over the course of a month. CompVis, the machine vision and learning research group at Ludwig Maximilian University of Munich, oversaw the training, while Stability AI donated the compute power.

Stable Diffusion was trained on LAION-5B

- **Goal:** Leveraging the pretrained diffusion models to model **new types of conditions** while ensuring high-fidelity outputs

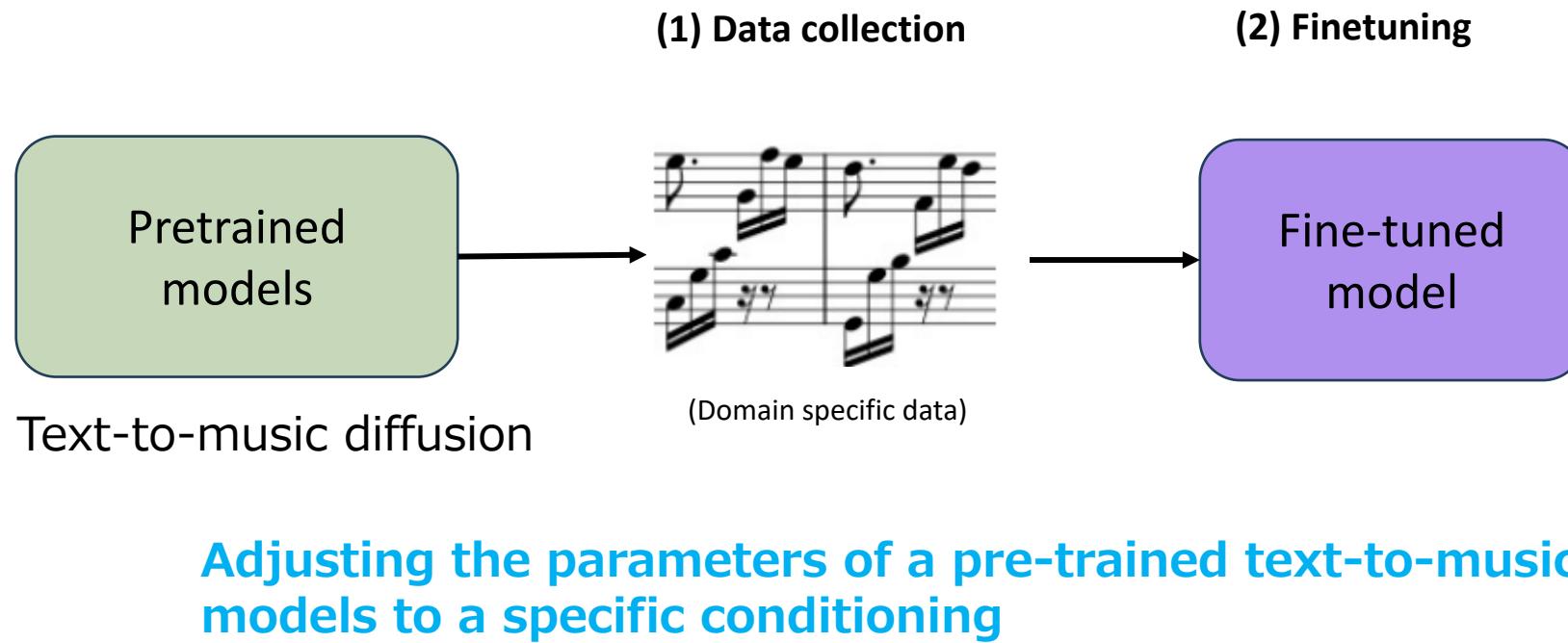
$$p(\mathbf{x} | \mathbf{c}_{\text{text}}, \mathbf{y})$$

Controllable Generation (II)



Training-based	Training-free	
Finetuning	Inference-time optimization	Inference-time guidance
<ul style="list-style-type: none">• MusicControlNet• ...	<ul style="list-style-type: none">• DITTO• DITTO-2• ...	<ul style="list-style-type: none">• MusicMagus• CQT-Diff• ...

Finetuning (I)



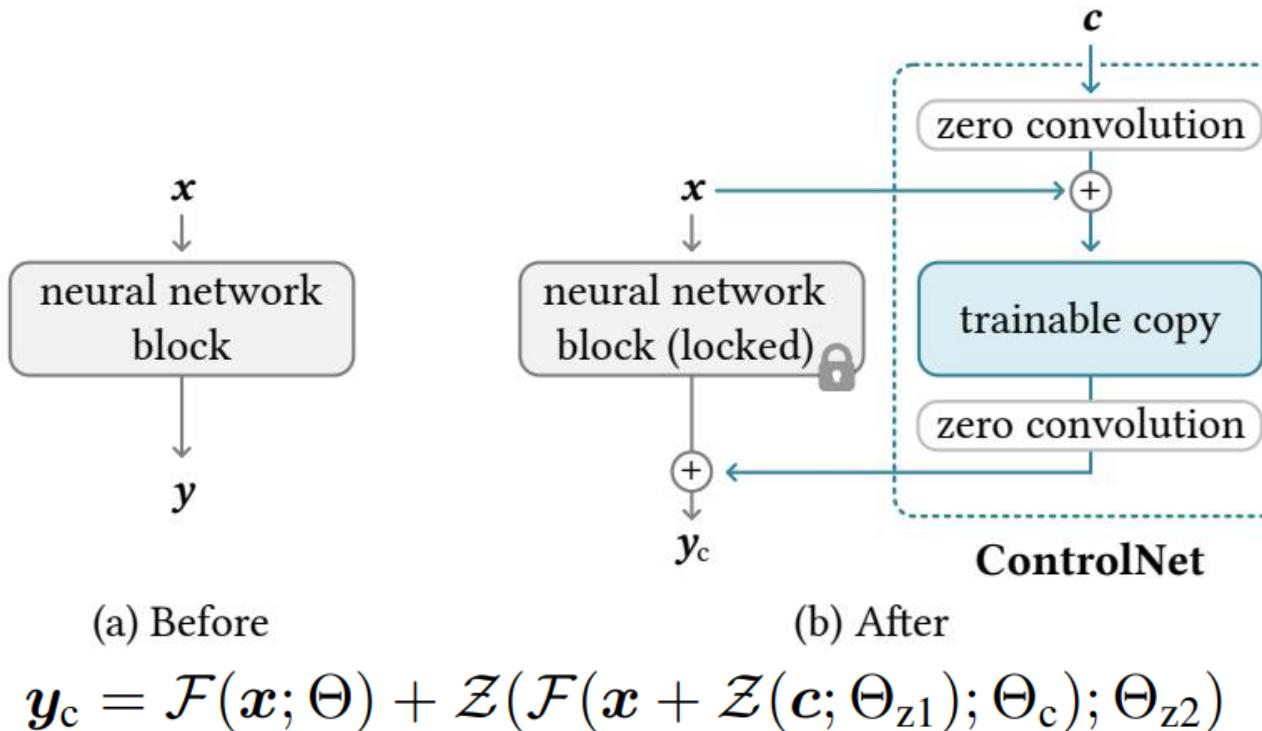
➤ Challenges

- Limitation of finetuning data causes **overfitting** and **catastrophic forgetting**

Finetuning (II)

➤ ControlNet

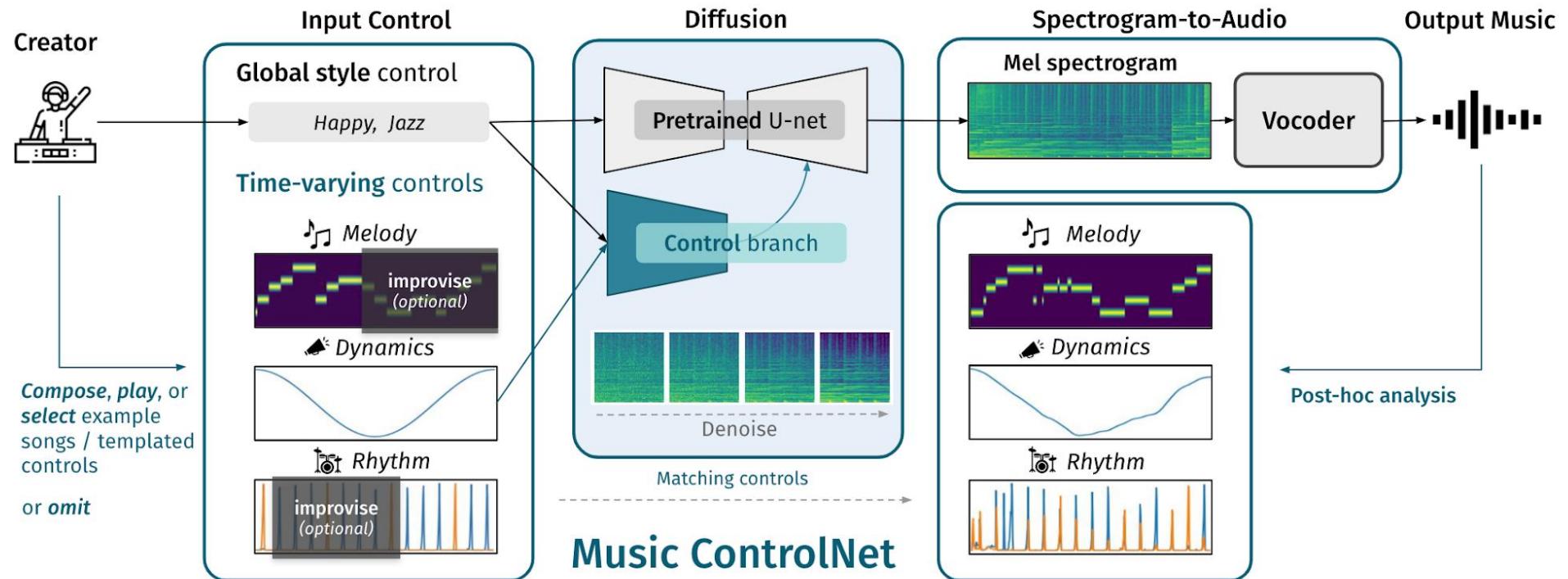
- Adding conditioning controls to large, pretrained diffusion models
- Reusing the deep and robust encoding layers
- “Zero convolutions” ensure that no harmful noise could affect the finetuning



Finetuning (III)

➤ **MusicControlNet**: Enables fine-grained control generation over

- Positions of beats in time
- Dynamics
- Melody
- Rhythm



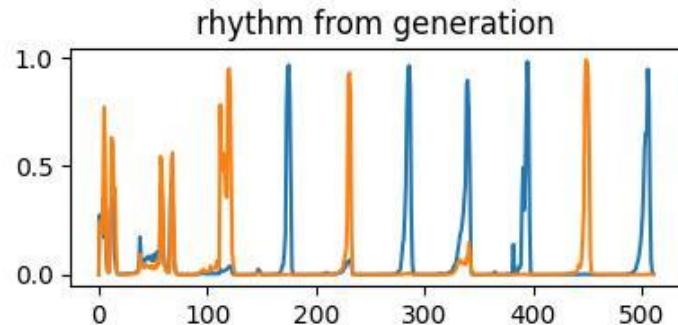
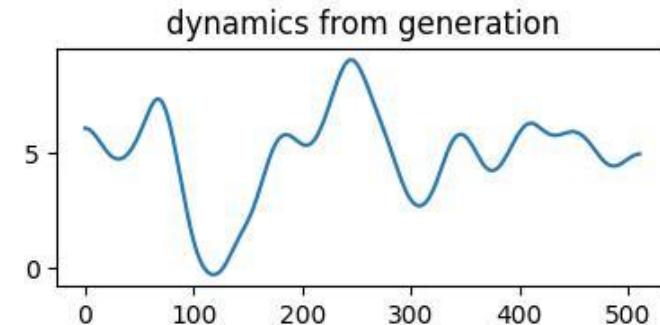
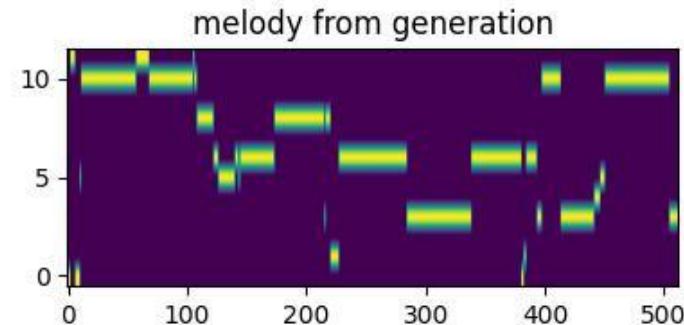
Finetuning (IV)

- **MusicControlNet**: Enables fine-grained control generation over

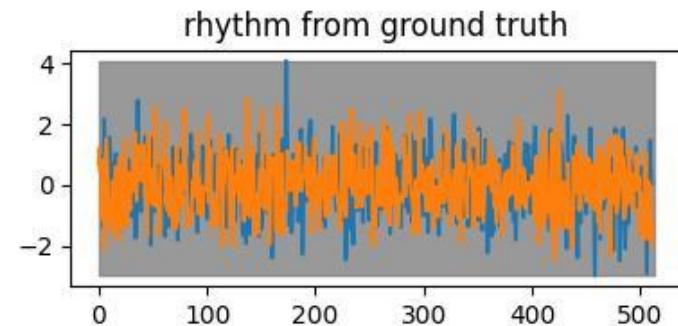
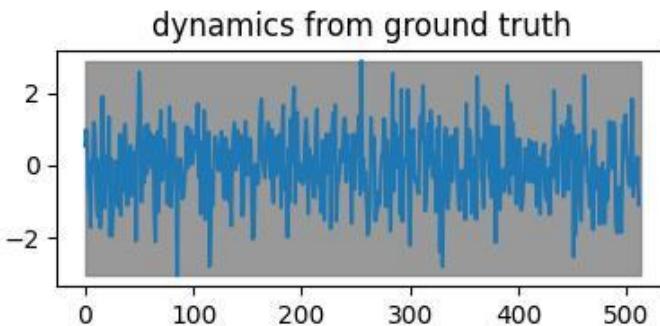
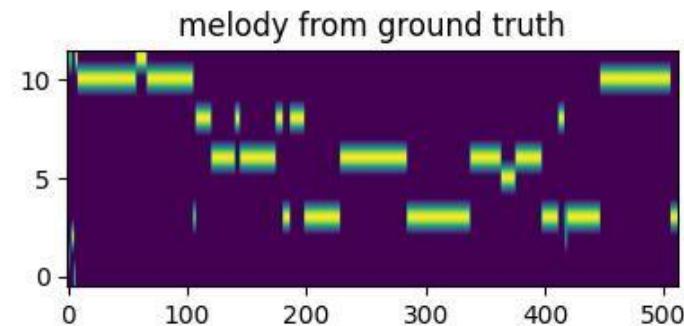
Text prompt: happy, jazz



Generated



Reference



Melody Control

Inference-Time Optimization (I)

➤ Formulate the control task as an optimization problem

- Searching for the noise such that the model output follows the target control \mathcal{A}

$$\arg \min_{\mathbf{x}_T \in \mathcal{N}(0, \mathbf{I})} \mathcal{L}(\mathcal{A}(\mathbf{x}_0), \mathbf{y})$$
$$\mathbf{x}_0 = \text{Sampler}_T(\epsilon_\theta, \mathbf{x}_T, \mathbf{c}_{\text{text}})$$

\mathbf{x}_T is optimized rather than model weights

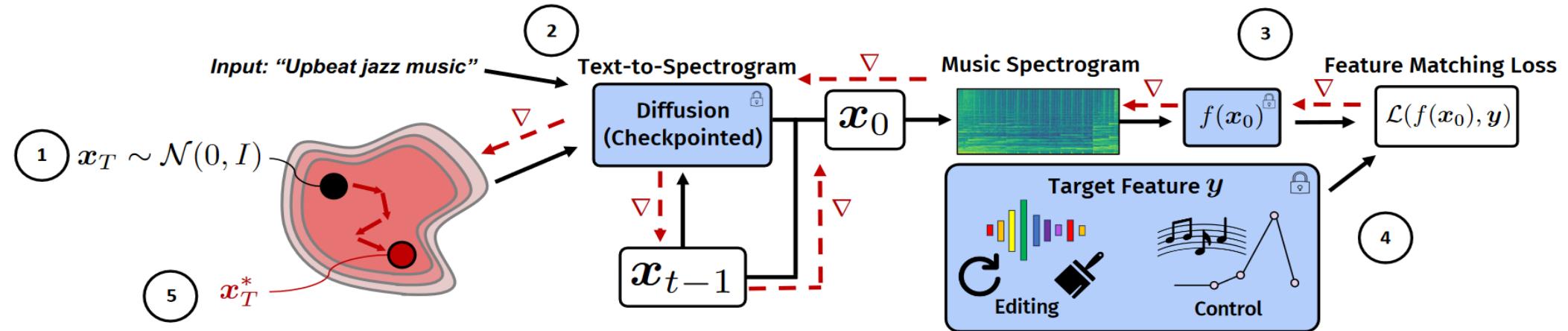
➤ Challenges

- Diffusion sampling process consists of multiple steps
- Solving the optimization problem requires huge memory

Inference-Time Optimization (II)

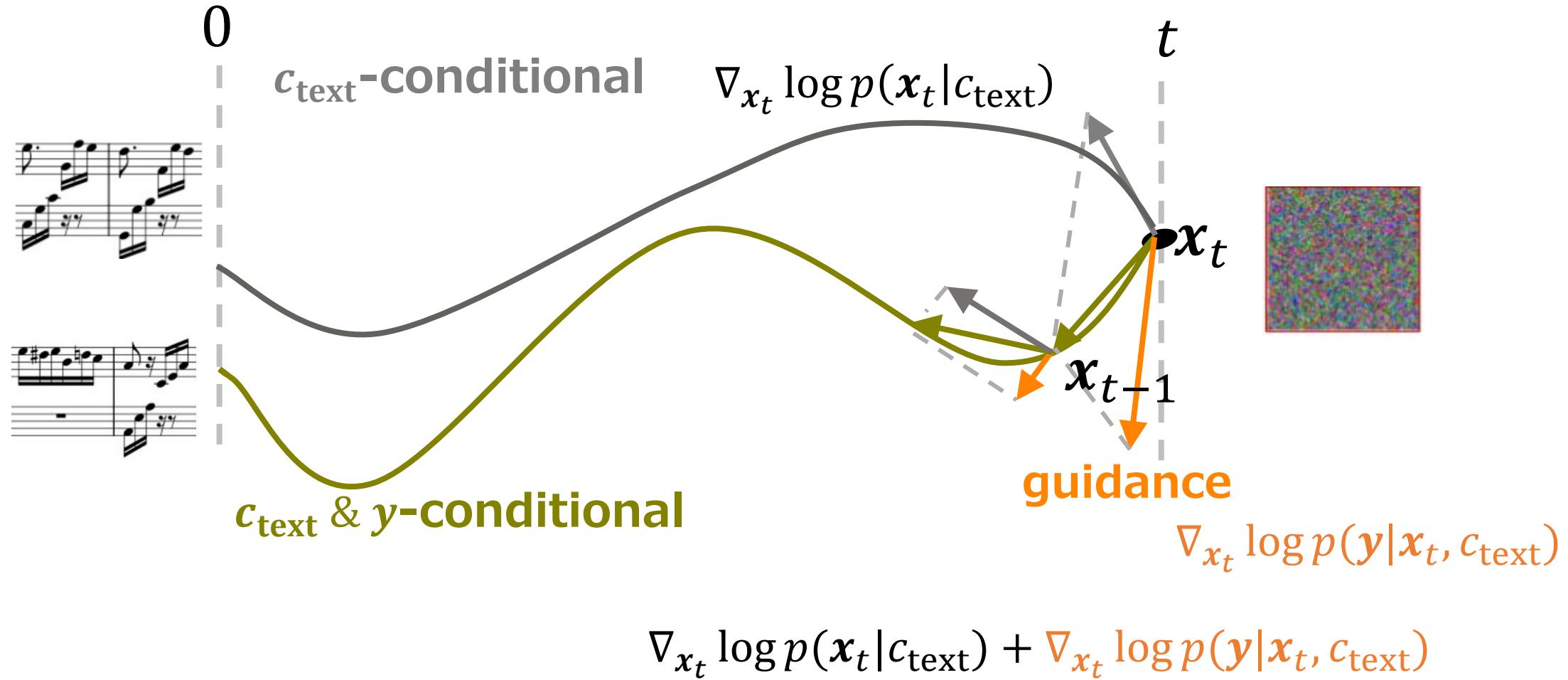
➤ DITTO

- Gradient Checkpointing to circumvent the large memory usage
- Store the intermediate noisy diffusion tensors



- Provide controllability over intensity, melody, musical structure, and looping

Inference-Time Guidance (I)



➤ Challenges

- Require an approximation of the guidance during sampling

Inference-Time Guidance (II)

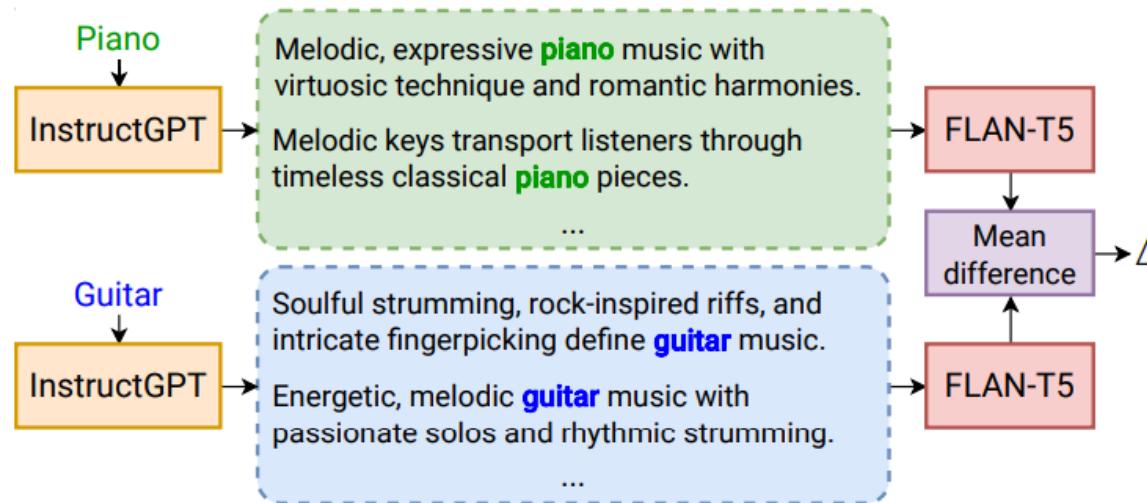
➤ MusicMagus

- Focuses on modifications such as adding effects or changing instruments
 - E.g., `add reverb to this soundtrack`, `transfer the timbre of the specific notes`, etc.
- It can perform zero-shot editing without additional training pairs



Inference-Time Guidance (II)

- Finding editing direction

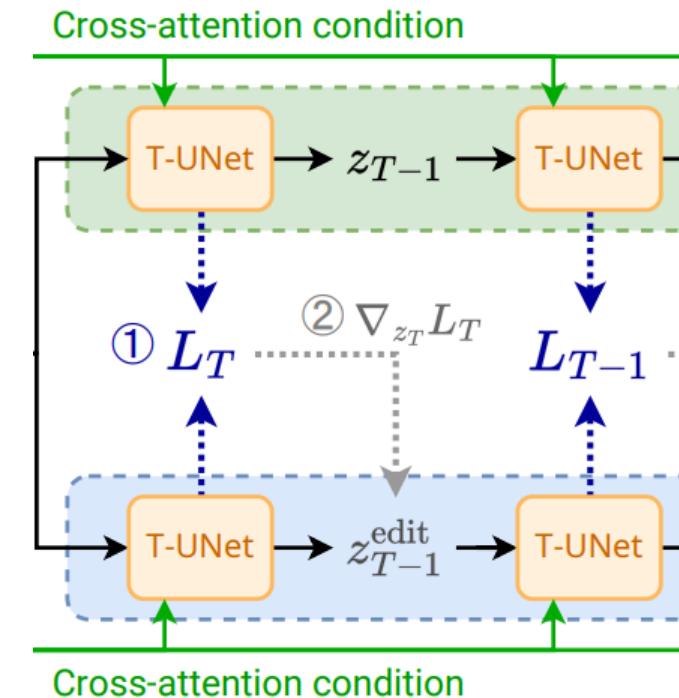


$$E^{\text{edit}} = \{E_{\text{T5}} + \Delta, E'_{\text{GPT}}\}$$

Example: Piano to violin



- Adding constraints over cross-attention



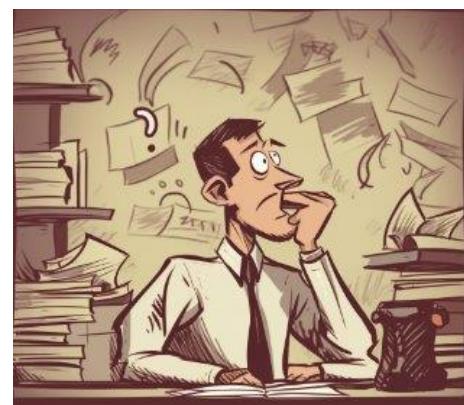
$$L_t = \|M_t^{\text{edit}} - M_t^{\text{origin}}\|_2$$

$$\epsilon_{\theta}^{\text{edit}} = \epsilon_{\theta}(z_t - \alpha \nabla_{z_t} L_t, E^{\text{edit}}, t)$$

Takeaways

Training-based	Training-free	
Finetuning	Inference-time optimization	Inference-time guidance
✓ Arbitrary controls	✓ Zero training	✓ Zero training
✓ Fast inference	✓ Applicable to many tasks	✓ Applicable to many tasks
- Large-scale training data	- Slow inference	- Limited fine-grained controllability
- Only applicable to specific problems	- Differentiable control	

- Which method should I use to enable controllable generation for my problem?
 - Do I have **enough budgets** (data + compute + skills)? → Finetuning
 - I'm working on small data and want **quick results** → Inference-time guidance
 - I'm willing to pay some cost (slow at inference) because I care about **quality**
 - Inference-time optimization



INVERSE PROBLEM

- **General Inverse Problem**
- **Linear Inverse Problem**
- **Linear Blind Inverse Problem**

Motivation

- ## ➤ Finding **X** given **y**

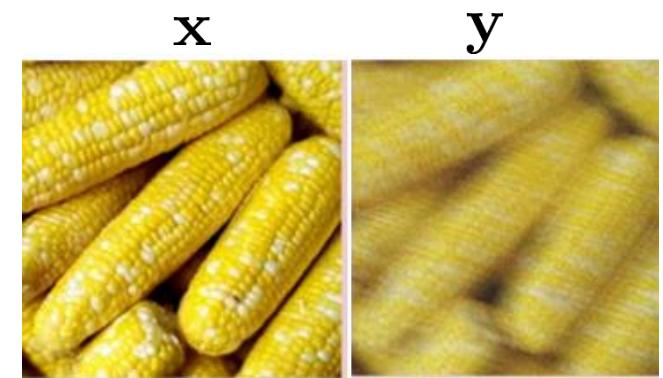
$$p(\mathbf{x}|\mathbf{y})$$

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{z}$$

Applications



y



Phase retrieval

deblurring

Ill-posed: There are infinitely many solutions ...

$\mathbf{z} \sim \mathcal{N}(0, \sigma_{\mathbf{v}}^2 \mathbf{I})$ Gaussian noise with known variance

\mathcal{A} Measurement operator

- How can we use diffusion models for solving general inverse problems?

- Training is expensive
 - It is data-specific and not applicable to many inverse problems

Challenge

- ## ➤ How to sample with diffusion model?

$$\mathbf{x}_{t-dt} \leftarrow \mathbf{x}_t - t \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) dt$$

Using Bayesian theorem

$$p(\mathbf{x}_t | \mathbf{y}) = \frac{p(\mathbf{y} | \mathbf{x}_t) p(\mathbf{x}_t)}{p(\mathbf{y})}$$

$$\log p(\mathbf{x}_t|\mathbf{y}) = \log p(\mathbf{y}|\mathbf{x}_t) + \log p(\mathbf{x}_t) - \cancel{\log p(\mathbf{y})}$$

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t | \mathbf{y}) = \nabla_{\mathbf{x}_t} \log p(\mathbf{y} | \mathbf{x}_t) + \nabla_{\mathbf{x}_t} \log p(\mathbf{x}_t)$$

Don't know

Diffusion model

General Inverse Problem (I)

$$\mathbf{y} = \mathcal{A}(\mathbf{x}) + \mathbf{z}$$

Given:

$$\mathcal{A}: \mathbb{R}^n \rightarrow \mathbb{R}^m$$

A known measurement operator

$$\mathbf{y} \in \mathbb{R}^m$$

Measurement

$$\mathbf{z} \sim \mathcal{N}(0, \sigma_y^2 \mathbf{I})$$

Gaussian noise with known variance

Find:

$$\mathbf{x} \in \mathbb{R}^n$$

The true signa

How to estimate?

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t)$$

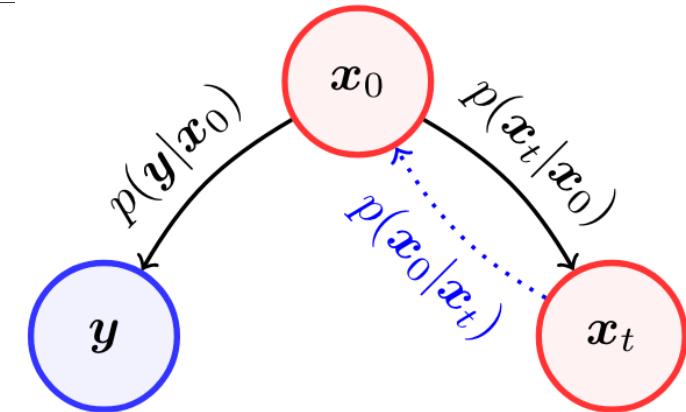
General Inverse Problem (II)

- Diffusion Posterior Sampling (DPS)

$$p(\mathbf{y}|\mathbf{x}_t) = \int p(\mathbf{y}|\mathbf{x}_0, \mathbf{x}_t)p(\mathbf{x}_0|\mathbf{x}_t)d\mathbf{x}_0$$

$$= \int p(\mathbf{y}|\mathbf{x}_0)p(\mathbf{x}_0|\mathbf{x}_t)d\mathbf{x}_0$$

$$= \mathbb{E}_{\mathbf{x}_0 \sim p(\mathbf{x}_0|\mathbf{x}_t)}[p(\mathbf{y}|\mathbf{x}_0)] \simeq p(\mathbf{y}|\hat{\mathbf{x}}_0) \text{ where } \hat{\mathbf{x}}_0 := \mathbb{E}[\mathbf{x}_0|\mathbf{x}_t]$$



$\log p(\mathbf{y}|\mathbf{x}_0)$ is known

$\log p(\mathbf{y}|\mathbf{x}_t)$ is unknown

$$\nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\mathbf{x}_t) \simeq \nabla_{\mathbf{x}_t} \log p(\mathbf{y}|\hat{\mathbf{x}}_0)$$

$$-\frac{1}{\sigma_y^2} \nabla_{\mathbf{x}_t} \|\mathbf{y} - \mathcal{A}(\hat{\mathbf{x}}_0(\mathbf{x}_t))\|_2^2$$

Linear Inverse Problem (I)

$$\mathbf{y} = \mathbf{Hx} + \mathbf{z}$$

Given:

$$\mathbf{H} \in \mathbb{R}^{m \times n}$$

$$\mathbf{y} \in \mathbb{R}^m$$

$$\mathbf{z} \sim \mathcal{N}(0, \sigma_{\mathbf{y}}^2 \mathbf{I})$$

A known linear degradation matrix

Measurement

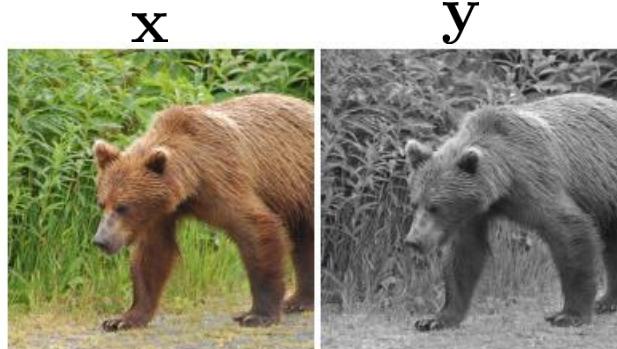
Gaussian noise with known variance

Find:

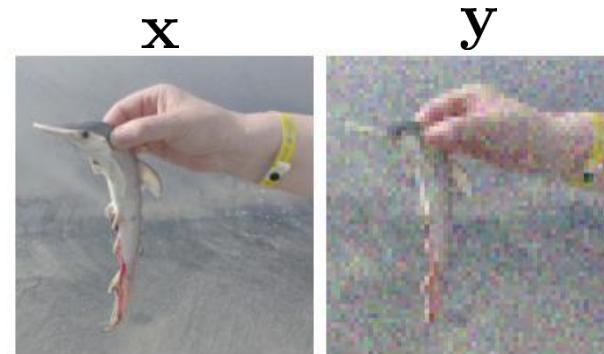
$$\mathbf{x} \in \mathbb{R}^n$$

The true signal

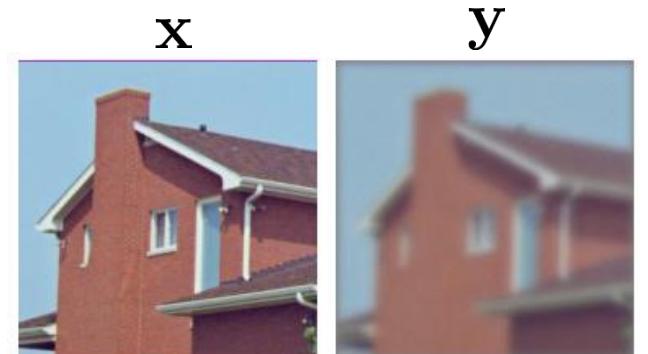
Ill-posed: There are infinitely many solutions ...



colorization



super-resolution

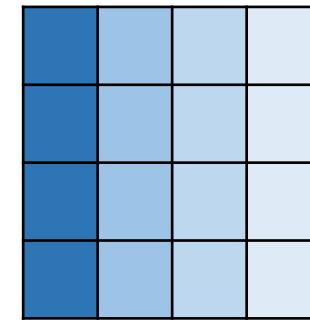
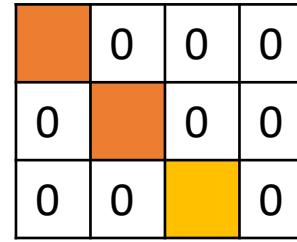
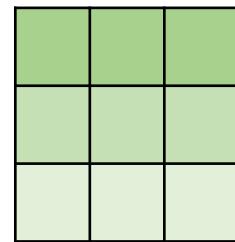
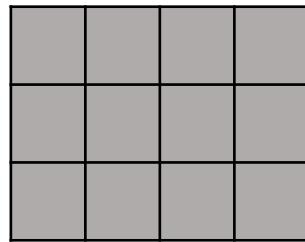


deblurring

Linear Inverse Problem (II)

- Denoising Diffusion Restoration Models (DDRM)

\mathbf{H} can be diagonalized



$$\mathbf{H} = \mathbf{U} \times \Sigma \times \mathbf{V}^T$$

$$\mathbf{U}^\top \mathbf{y} = \Sigma(\mathbf{V}^\top \mathbf{x}) + \mathbf{U}^\top \mathbf{z}$$

- Linear inverse problem can be transformed into “denoising and inpainting” in spectral space

Linear Blind Inverse Problem (I)

Given:

$$\mathbf{y} = \mathbf{H}_\varphi(\mathbf{x}) + \mathbf{z}$$

$$\begin{aligned}\mathbf{H}_\varphi &: \mathbb{R}^n \rightarrow \mathbb{R}^m \\ \mathbf{y} &\in \mathbb{R}^m \\ \mathbf{z} &\sim \mathcal{N}(0, \sigma_{\mathbf{y}}^2 \mathbf{I})\end{aligned}$$

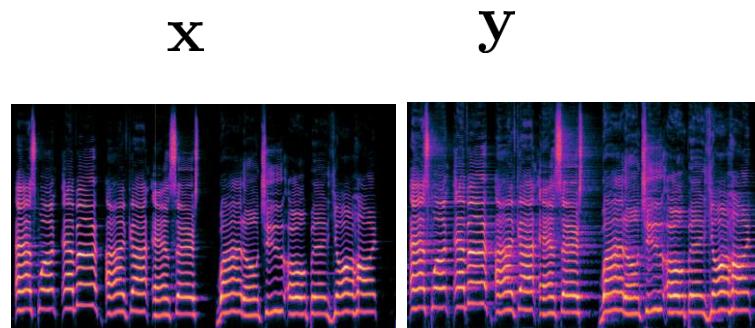
An unknown linear degradation matrix
Measurement
Gaussian noise with known variance

Find:

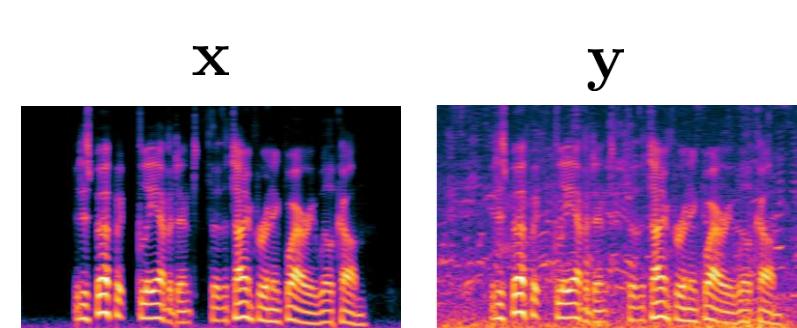
$$\mathbf{x} \in \mathbb{R}^n$$

The true signal

Ill-posed: There are infinitely many solutions ...



Vocal dereverberation



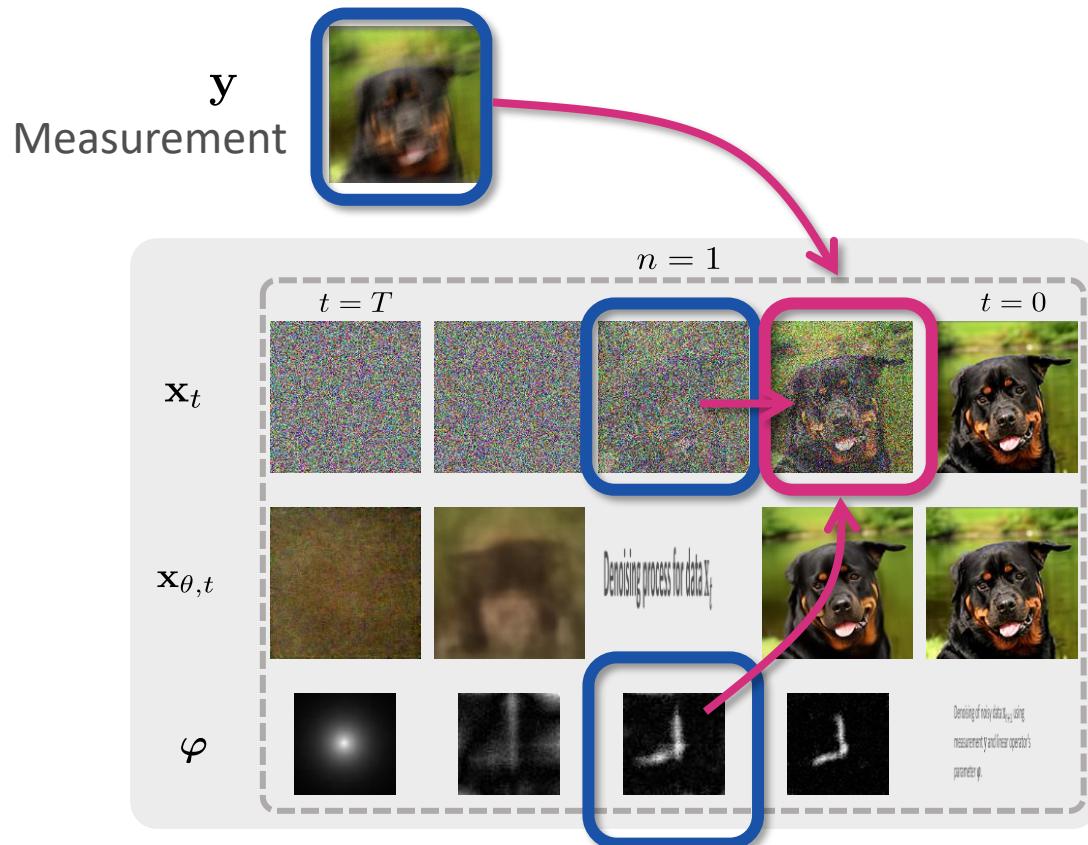
Speech enhancement



deblurring

Linear Blind Inverse Problem (II)

- Simultaneous estimation of data (sampling) and linear operator parameters



Denoising process for data \mathbf{x}_t

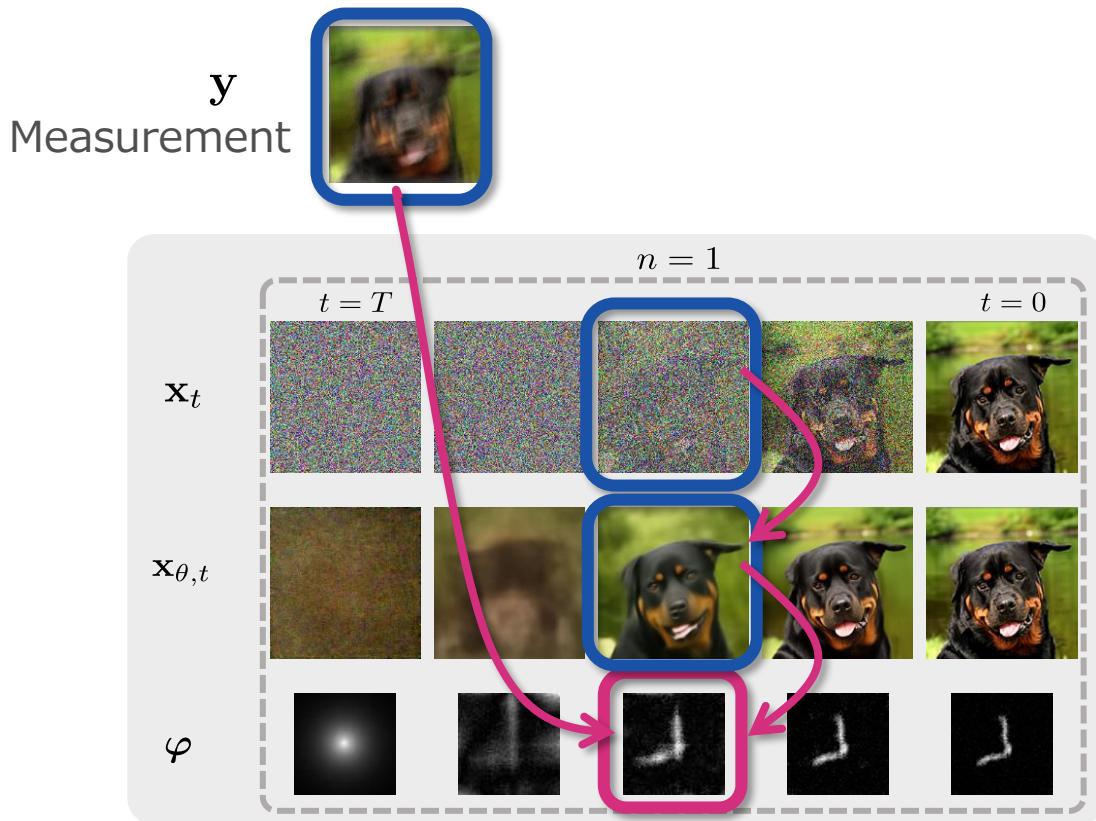
Denoising of noisy data \mathbf{x}_{t+1} using measurement \mathbf{y} and linear operator's parameter φ .

$$\mathbf{x}_t \sim p_{\theta}(\mathbf{x}_t | \mathbf{x}_{t+1}, \varphi, \mathbf{y})$$



Linear Blind Inverse Problem (III)

- Simultaneous estimation of data (sampling) and linear operator parameters



Sampling of linear operator's parameter φ

Parameter φ is estimated using measurement \mathbf{y} and clean data estimate $\mathbf{x}_{\theta,t}$

$$\nabla_{\varphi} \log p(\varphi | \mathbf{x}_{t:T}, \mathbf{y}) \simeq -\frac{1}{2\sigma_{\mathbf{y}}^2} \nabla_{\varphi} \|\mathbf{y} - \mathbf{H}_{\varphi} \mathbf{x}_{\theta,t}\|_2^2 + \nabla_{\varphi} \log p(\varphi)$$

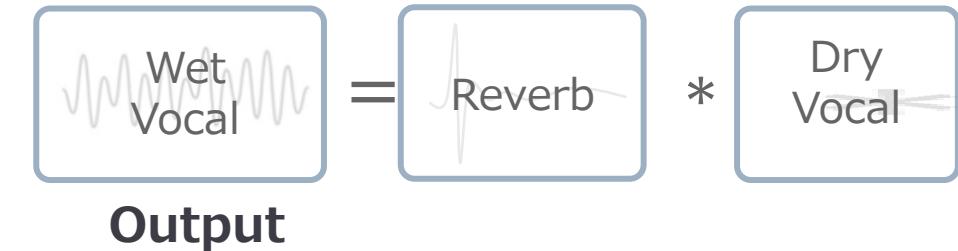
Measurement \mathbf{y}

Estimation of clean data $\mathbf{x}_{\theta,t}$

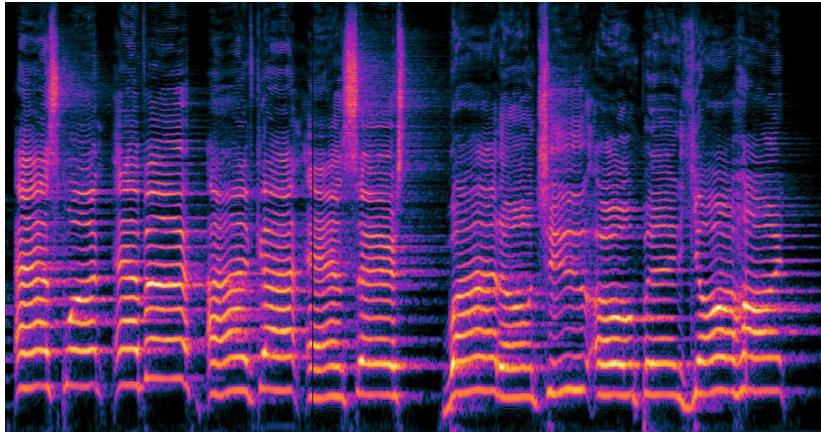
Prior on φ

Application: Vocal dereverberation

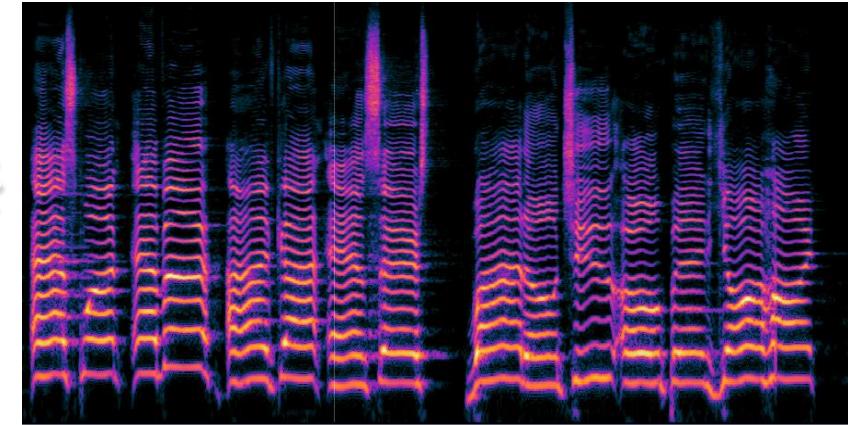
- Task: Restore the **dry vocal** from the **wet (reverberant) vocal**
- Pre-trained model: **Diffusion model trained on dry vocal data**
- Application: **Remix** of recorded contents



Input



(Reverberant vocal)



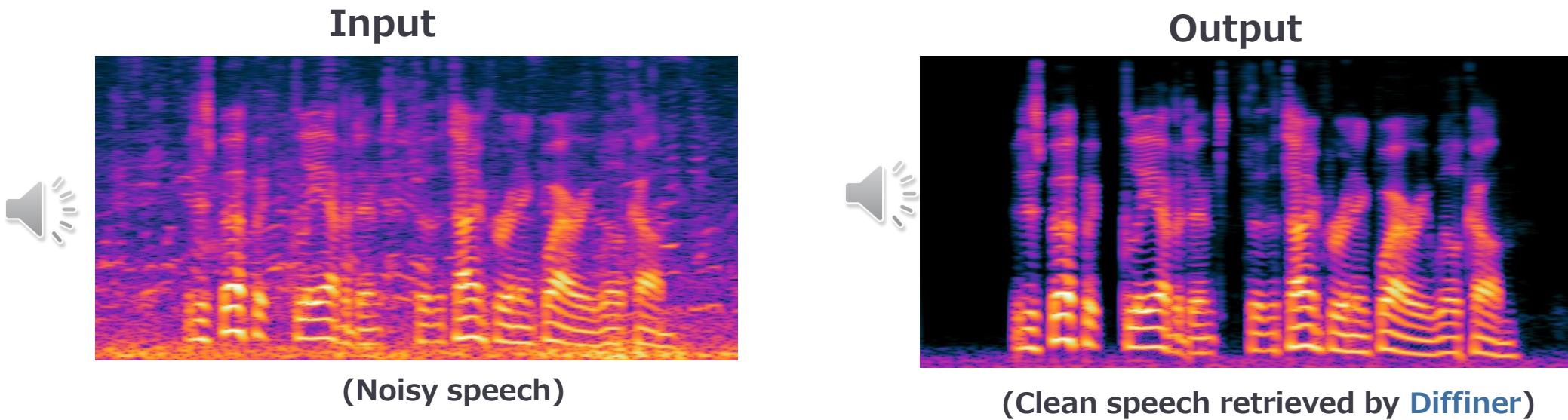
(Clean vocal retrieved by **DiffDereverb**)



Audio
samples

Application: Speech Enhancement

- Task: Extract **clean speech** from the **noisy speech**
- Pre-trained model: Diffusion model trained on clean speech data
- Application: Tele-communication, video post-production

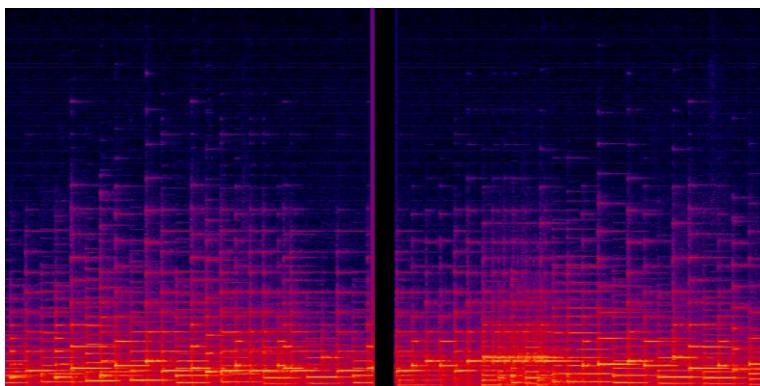


Application: Audio Inpainting

- Task: Inpainting audio from the **masked** input
- Pre-trained model: Diffusion model trained on clean audio data
- Application: Audio editing

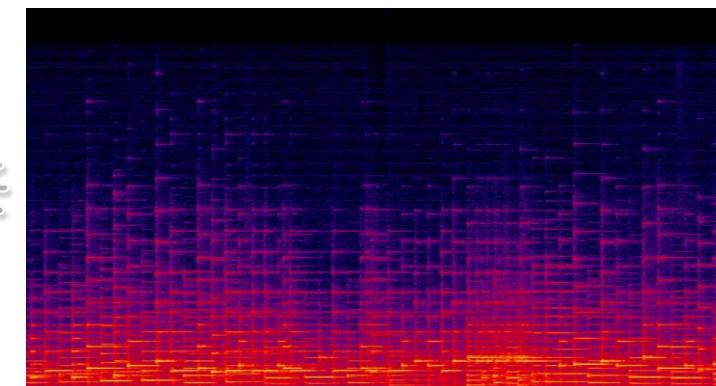
$$\mathcal{A}(\mathbf{x}) = (1 - \mathbb{1}_{[t_{\text{start}}, t_{\text{end}}]})\mathbf{x}$$

Input



Masked audio

Output



generation



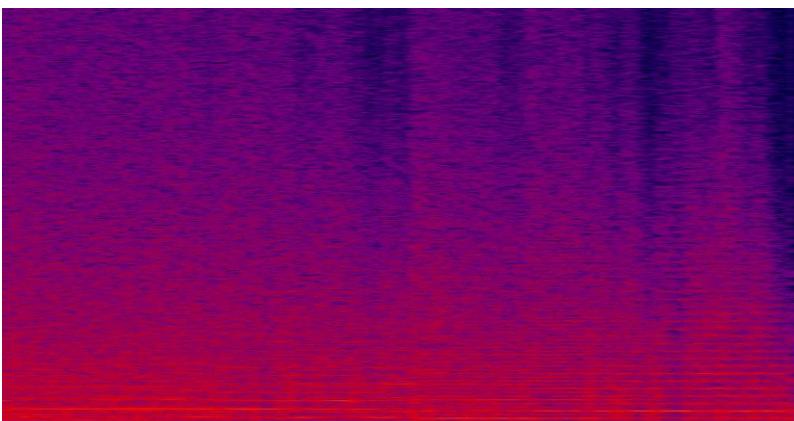
Application: Audio Declipping

- Task: Cancelling the distortion in audio signal that has been hard-clipped
 - Pre-trained model: Diffusion model trained on clean audio data
 - Application: Audio editing

Original audio

$$\mathcal{A}(\mathbf{x}) = \frac{|\mathbf{x} + c| - |\mathbf{x} - c|}{2}$$

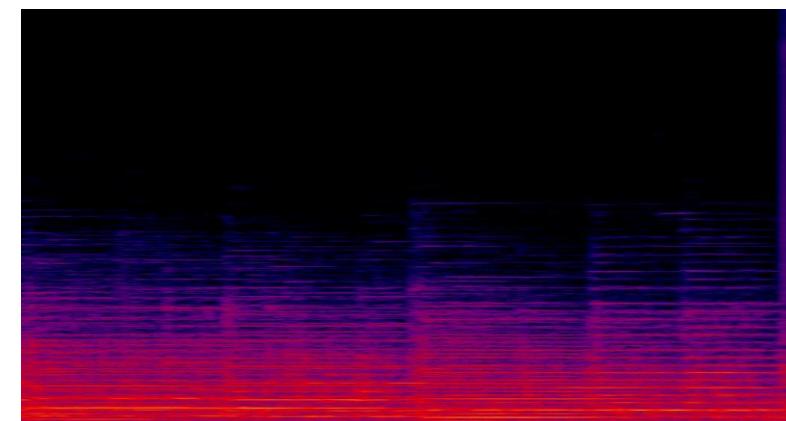
Input



Clipped audio



Output



generation

[Hernandez-Olivan, et al., ICASSP 2024] Vrdmg: Vocal restoration via diffusion posterior sampling with multiple guidance

[Moliner et al., ICASSP 2023] Solving audio inverse problems with a diffusion model



NEW TREND

New Trend

1. Audio-Visual Generation / Sounding Video Generation
 2. Stem-based Music Generation
 3. Memorization

Audio-Visual Generation / Sounding Video Generation

2023

2024

2025



DALL-E 3 · An expressive oil painting of a chocolate chip cookie being dipped in a glass of milk, depicted as an...

Video Generation

Sora



Sounding Video Generation ?



Unified vs Composable Modeling

1) Unified Modeling

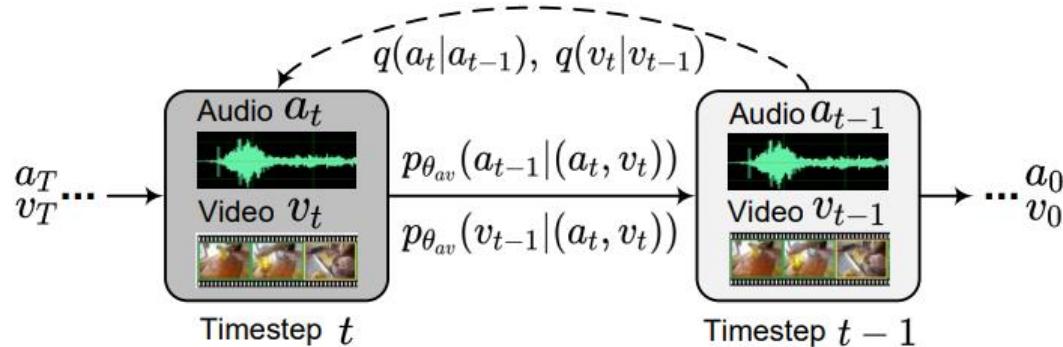
(e.g., MM-Diffusion [[Ruan+ CVPR2023](#)])

[Pros]

Multimodal alignment can be explicitly learned
-> **better in alignment**

[Cons]

A lot of pair data are needed



2) Composable Modeling

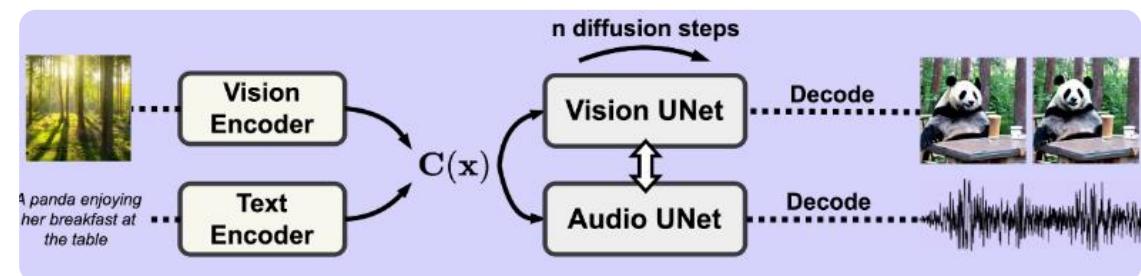
(e.g., CoDi [[Tang+ NeurIPS2023](#)])

[Pros]

Off-the-shelf single-modal pretrained models can be incorporated in the framework -> **better in quality**

[Cons]

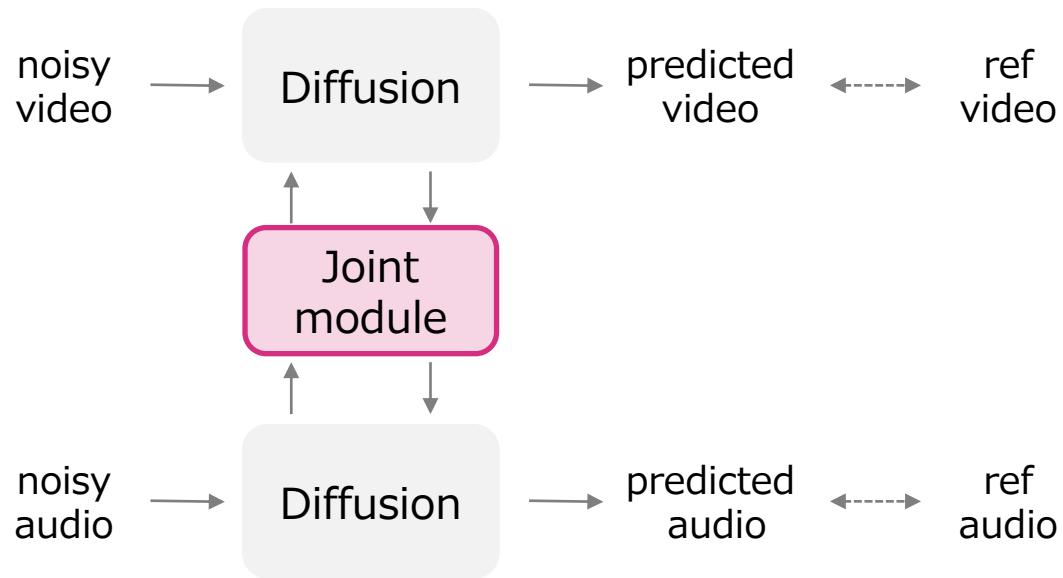
Bridging modalities need to be learned w/ pair data



Composable Modeling

Finetuning

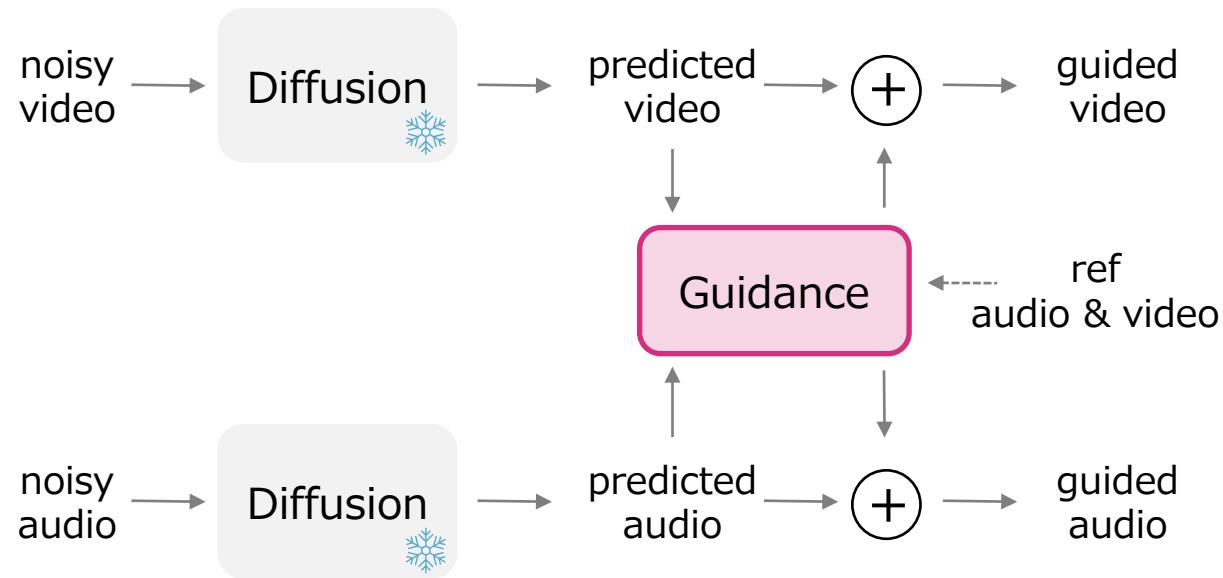
[[Tang+ NeurIPS2023](#)][[Ishii+ ECCV24 AVGenL](#)]



- Finetune pretrained model to predict joint distribution $p(x, y)$
 - Depends on model architectures

Guidance

[Xing+ CVPR2024][Hayakawa+ 2024]



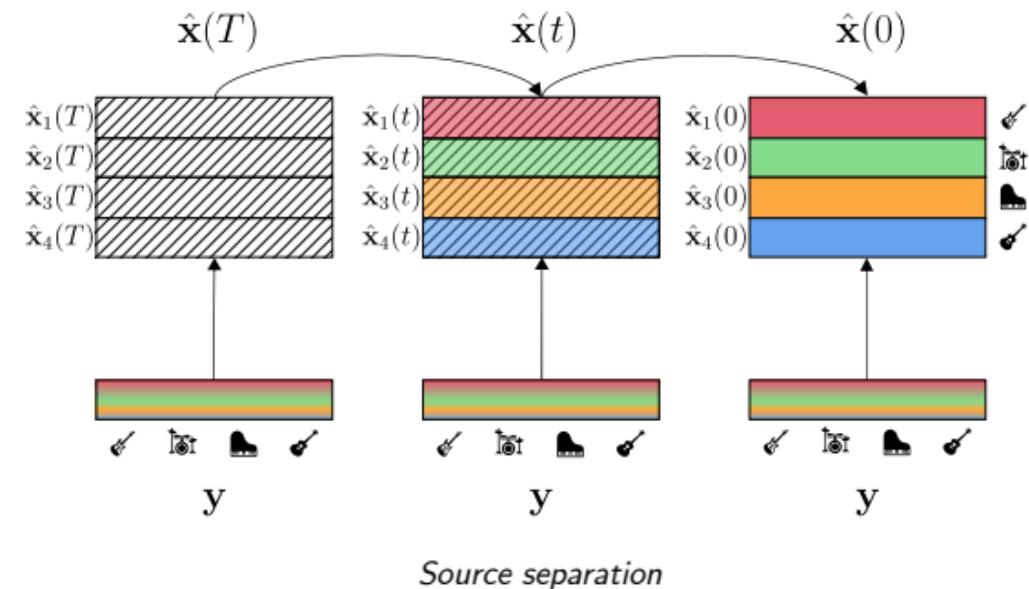
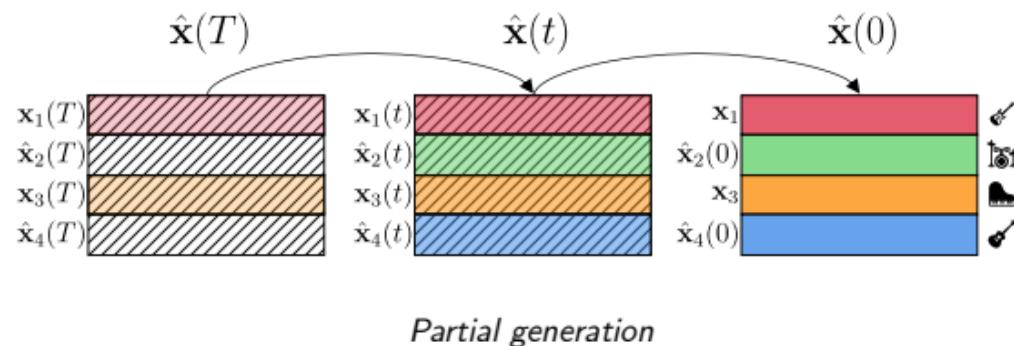
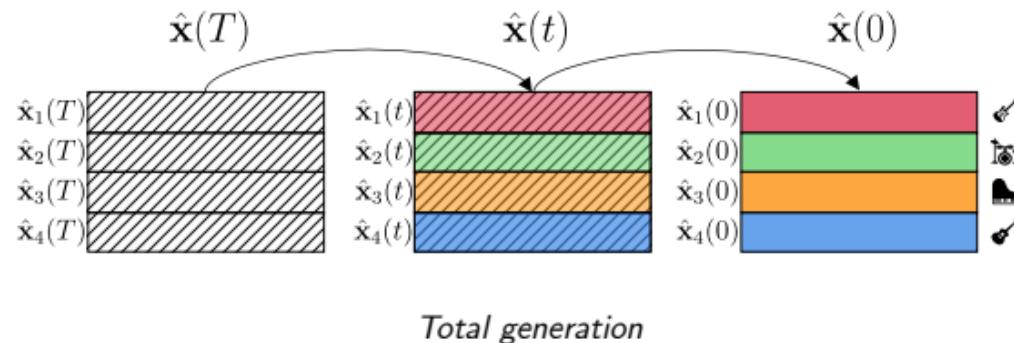
- Refine prediction from pretrained model towards joint distribution $p(x, y)$
 - Model architecture agnostic

Academic Challenge on Sounding Video Generation [[URL](#)]

The image shows a screenshot of the AIcrowd platform. At the top left is the AIcrowd logo with a red devil icon. A search bar is positioned next to it. On the right side, there are links for 'Challenges' and three vertical dots for more options. Below the header, a red banner indicates a 'Warm-up Round: 32 days left'. The main title of the challenge is 'SONY Sounding Video Generation Challenge' in large yellow text. To the right of the title is a large, stylized graphic featuring a blue play button inside a triangle, set against a background of colorful, glowing sound waves in shades of blue, purple, and red. Below the title, the challenge's purpose is described as 'Generate Synchronised & Contextually Accurate Videos' and a '\$ 35,000 USD Cash Prize' is mentioned. At the bottom, there are social sharing icons and a summary of the challenge's metrics: 1710 submissions, 30 participants, 0 teams, 6 rockets (likely referring to AI models), 1 like, and a 'Share' button.

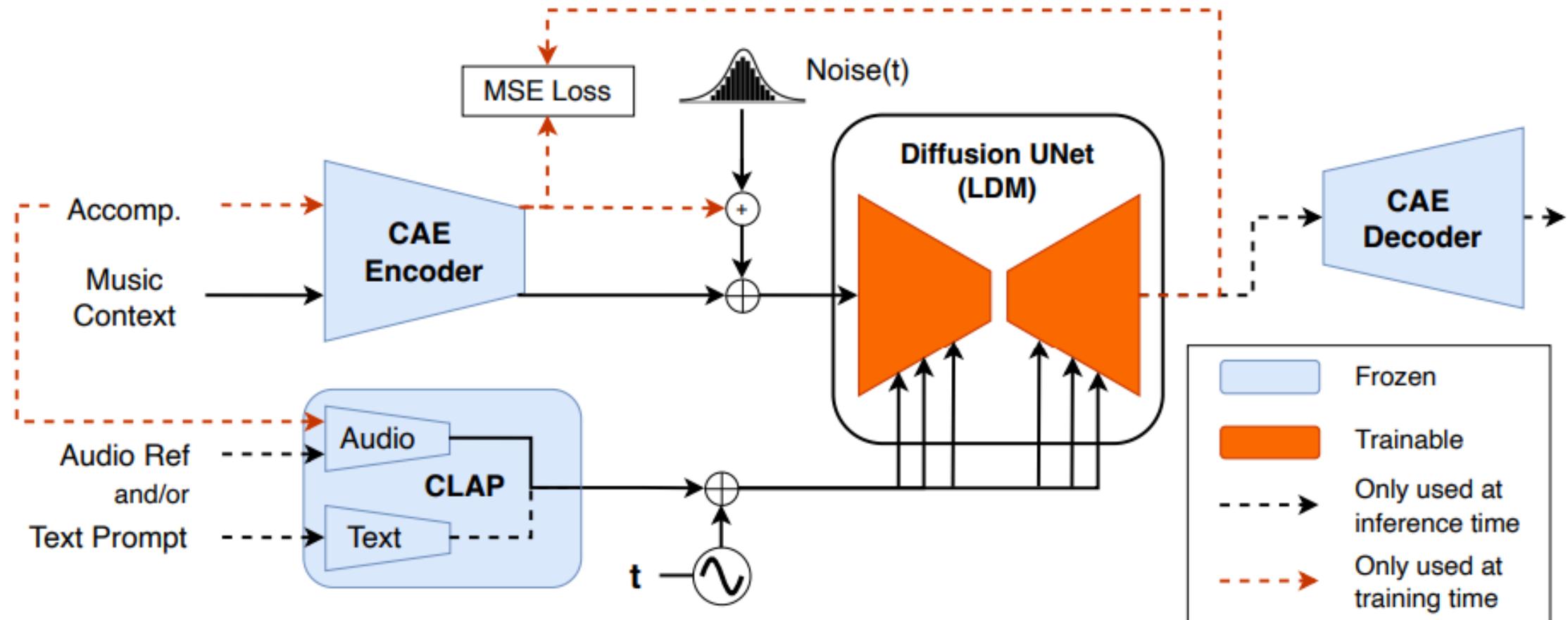
Multi-Source Diffusion [[Mariani+ ICLR2024](#)]

The model can achieve total generation, partial generation, and source separation



Single Instrument Accompaniment Generation [[Nistal+ ISMIR2024](#)]

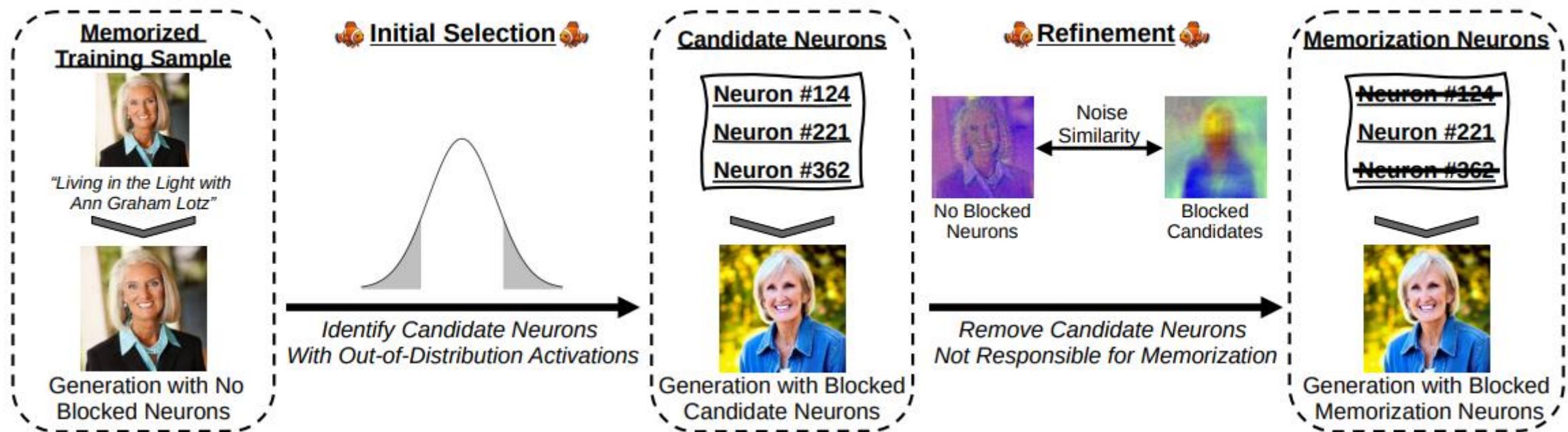
Diff-A-Riff: LDM designed to generate high-quality instrumental accompaniments
adaptable to any musical context



Memorization of Diffusion Models [[Hintersdorf+ NeurIPS2024](#)]

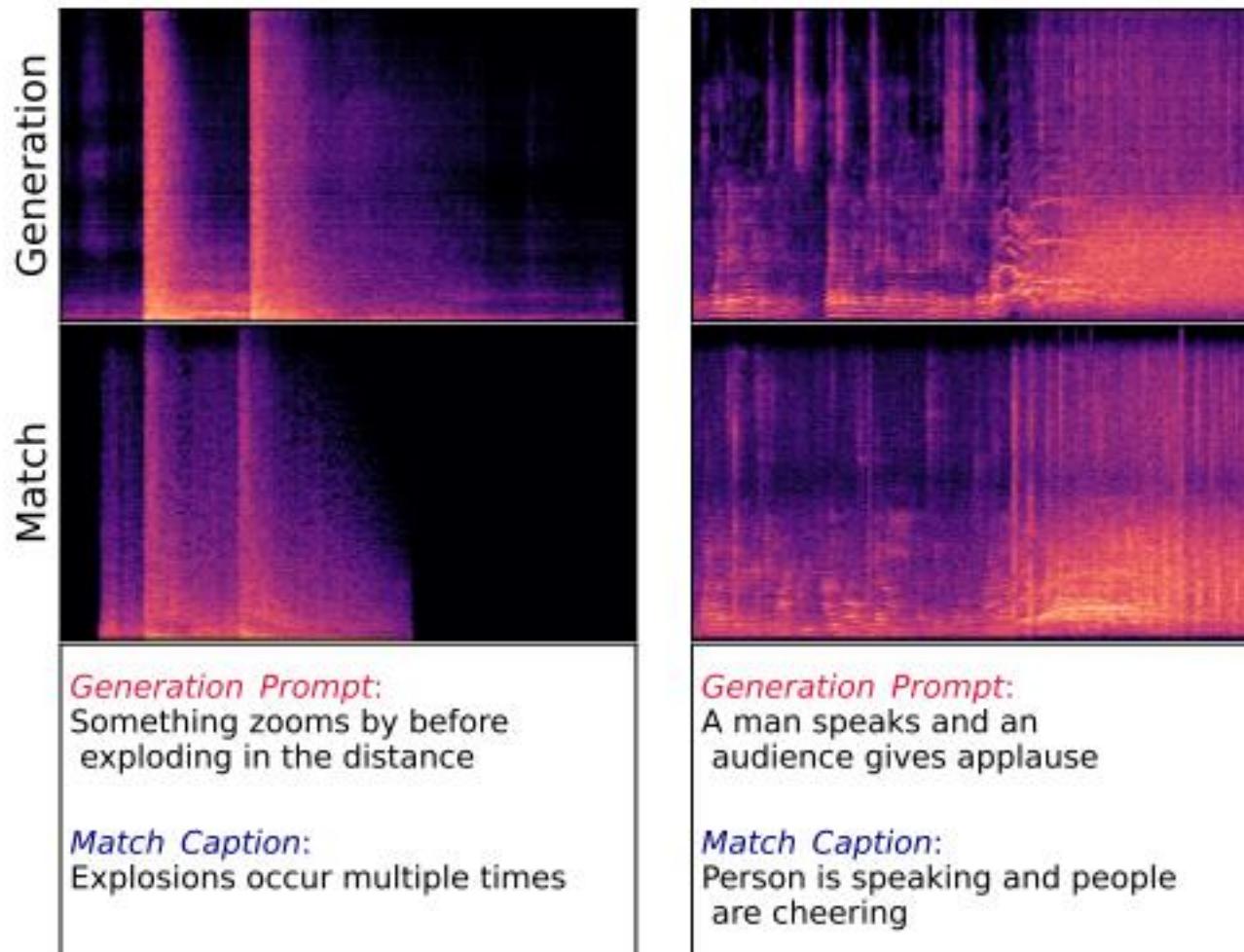
Single neurons are responsible for memorizing training samples

- ✓ Verbatim memorization: replicating the training image exactly
 - ✓ Template memorization: reproducing the general composition of the training image while having some non-semantic variations at fixed image positions



Memorization Problem in Audio Diffusion [[Bralios+ ICASSP2024](#)]

Samples generated using TANGO and their top matches in the AudioCaps



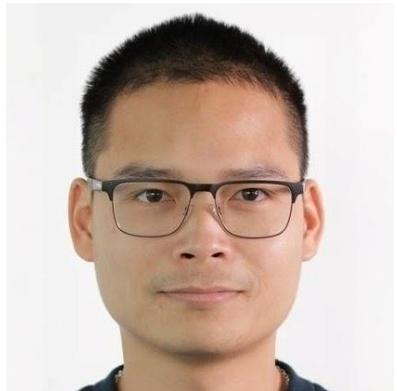
More work from
CreativeAI @Sony



MEET OUR SPEAKERS



Chieh-Hsin (Jesse) Lai



Bac Nguyen



Koichi Saito



Yuki Mitsufuji



@JCJesseLai



@nguyencongbacbk



@Koichi_Saito



@mittu1204



THANK YOU

