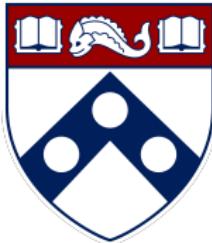


Diffusion model theory: low-dimensional adaptation, and acceleration



Yuxin Chen

Wharton Statistics & Data Science

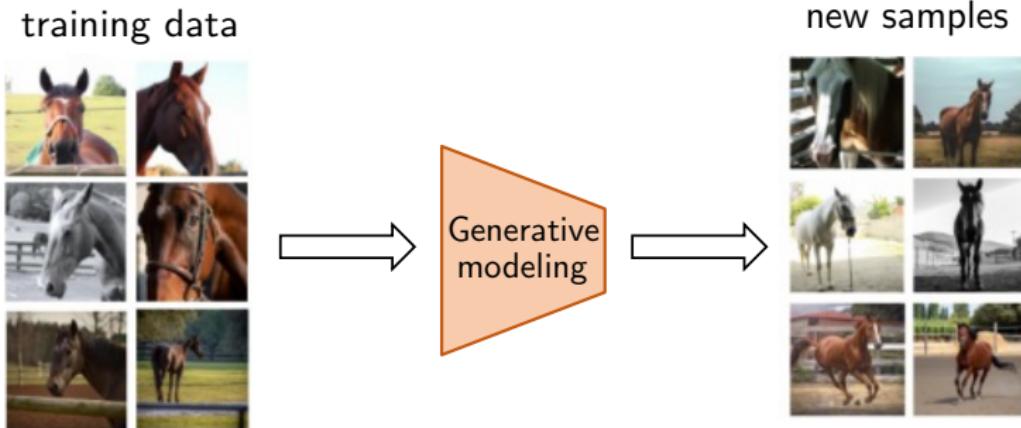
Generative modeling

training data



- Given training data $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a general distribution}} (1 \leq i \leq N)$ in \mathbb{R}^d

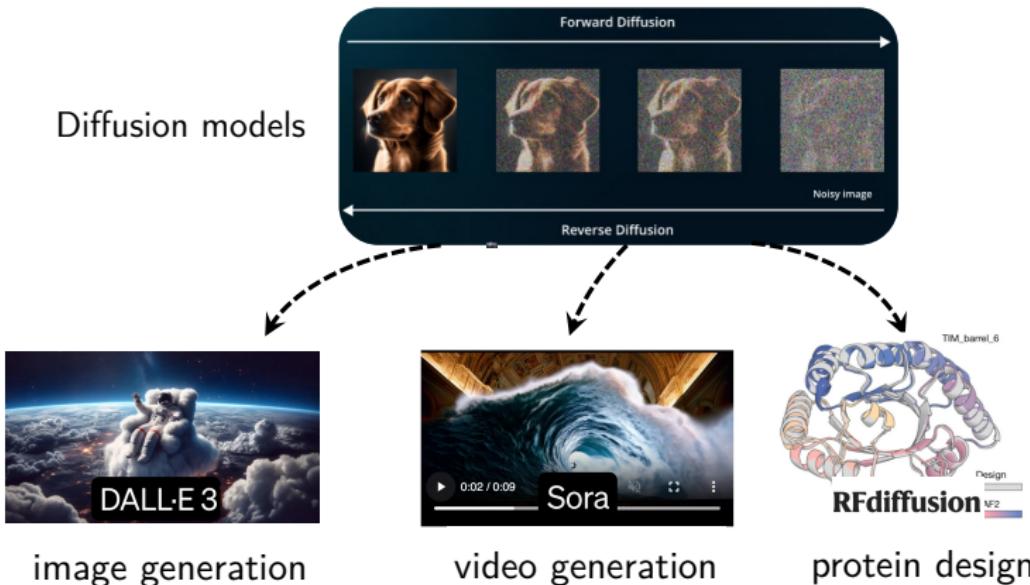
Generative modeling

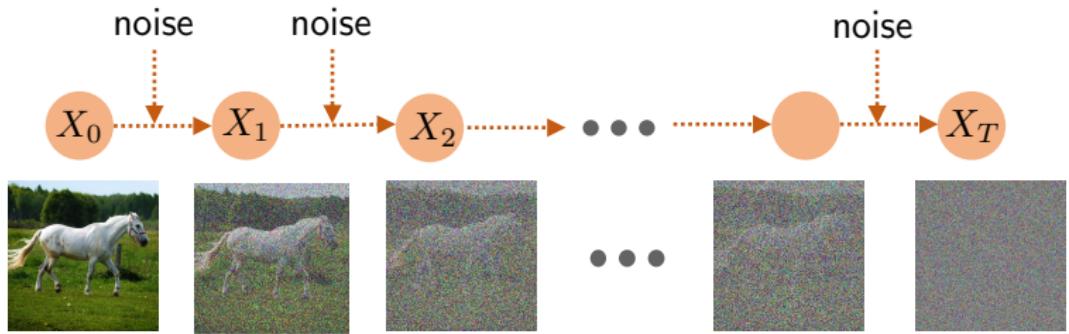


- Given training data $\underbrace{X^{\text{train},i} \sim p_{\text{data}}}_{\text{from a general distribution}} \quad (1 \leq i \leq N)$ in \mathbb{R}^d
- Generate **new** samples $Y \sim p_{\text{data}}$

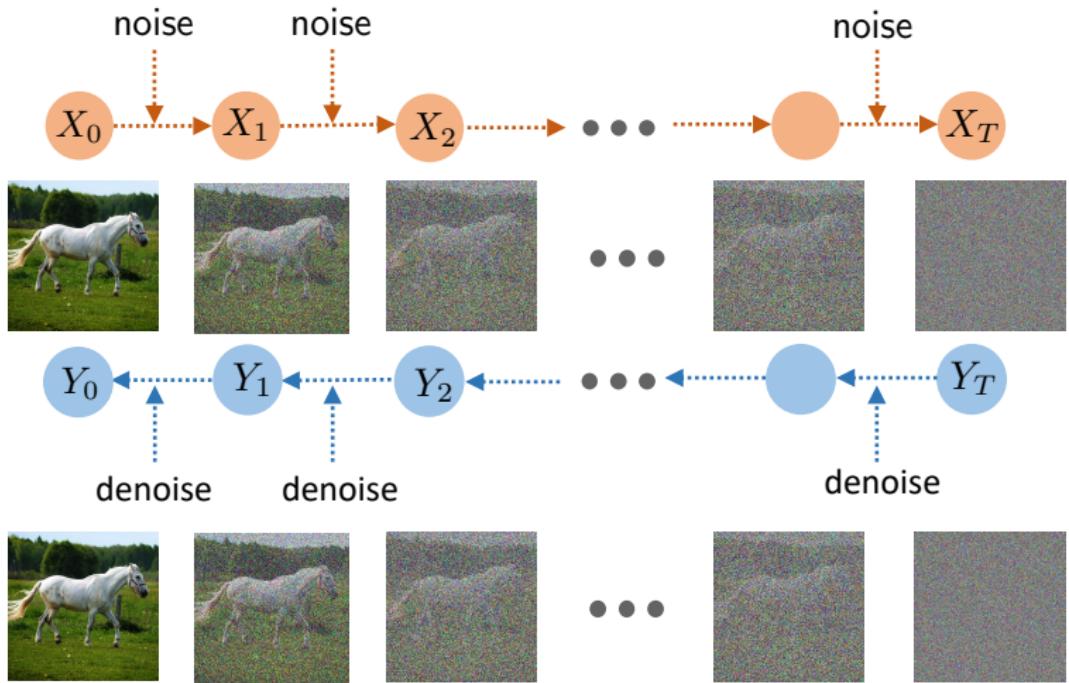
Inspired by nonequilibrium thermodynamics

— *Sohl-Dickstein, Weiss, Maheswaranathan, Ganguli '15*

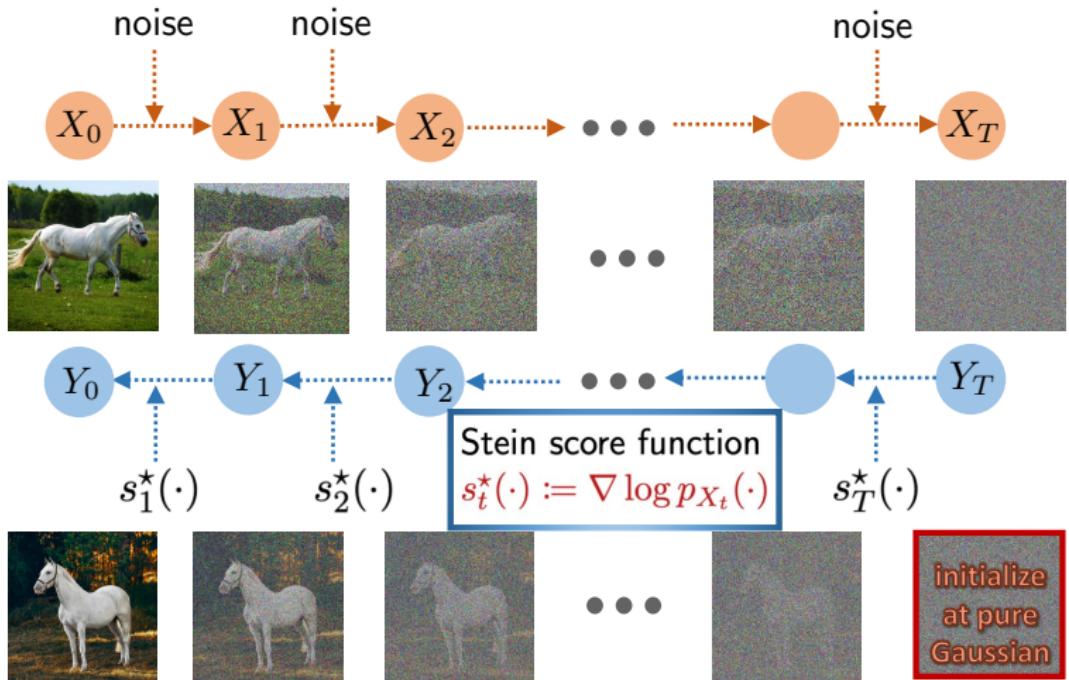




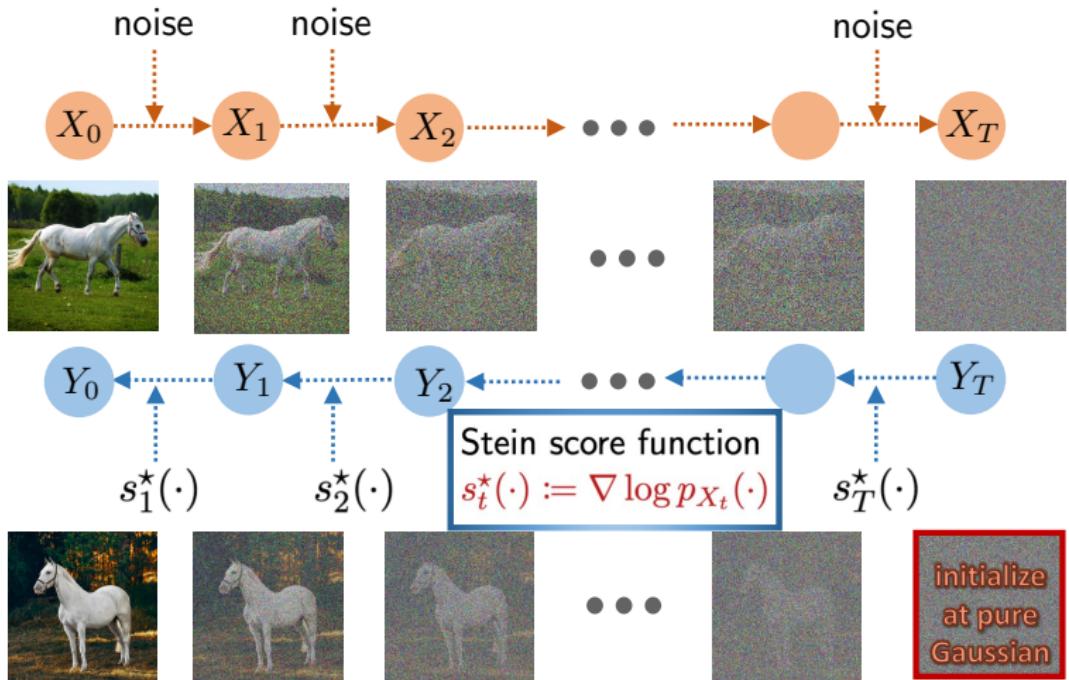
- **forward process:** diffuse data into noise



- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions



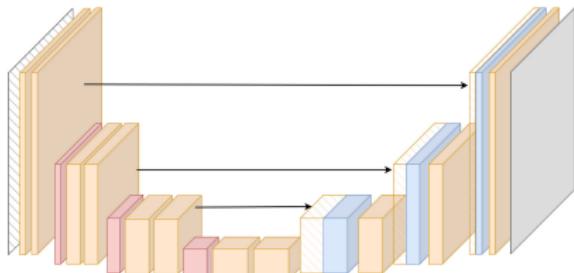
- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions



- **forward process:** diffuse data into noise
- **reverse process:** convert pure noise into data-like distributions

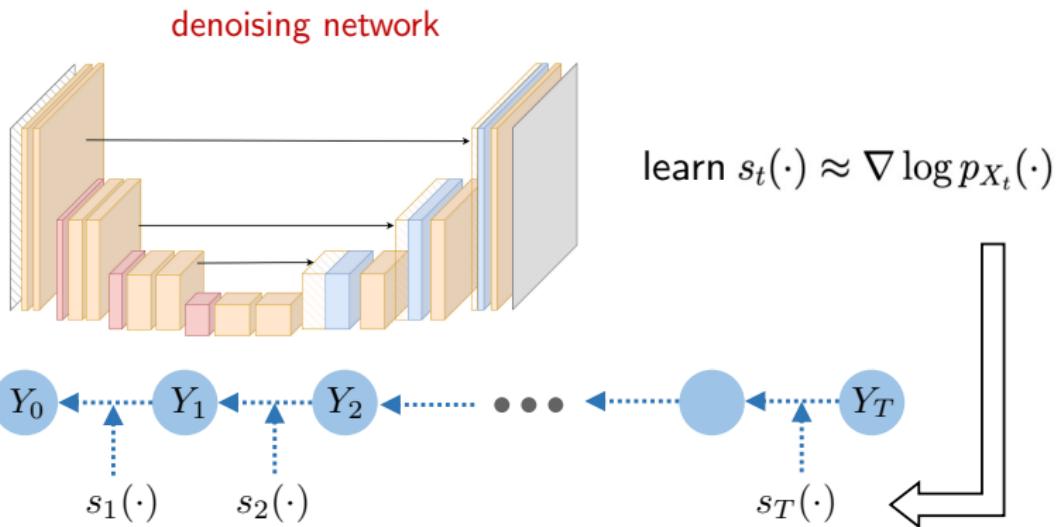
Goal: $Y_t \xrightarrow{d} X_t, \quad t = T, \dots, 1$

denoising network



learn $s_t(\cdot) \approx \nabla \log p_{X_t}(\cdot)$

1. **score learning/matching:** learn estimates $s_t(\cdot)$ for $\nabla \log p_{X_t}(\cdot)$



- 1. score learning/matching:** learn estimates $s_t(\cdot)$ for $\nabla \log p_{X_t}(\cdot)$
- 2. data generation:** sampling w/ the aid of score estimates $\{s_t(\cdot)\}$

Towards mathematical theory for diffusion models

- hard to develop full-fledged **end-to-end** theory

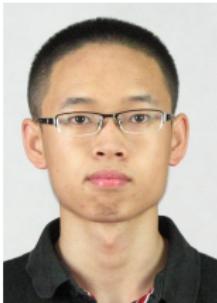
Towards mathematical theory for diffusion models

- hard to develop full-fledged **end-to-end** theory
 - “divide-and-conquer”: score learning $\leftarrow \cancel{X} \rightarrow$ sampling
-  decouple

Agenda:

1. non-asymptotic convergence theory
2. adaptation to (unknown) low dimensionality
3. acceleration via higher-order approximation

Part 1: nonasymptotic convergence theory



Gen Li
CUHK



Yuting Wei
UPenn



Yuejie Chi
Yale

Two mainstream approaches

Denoising Diffusion Probabilistic Models

Jonathan Ho

UC Berkeley

jonathanho@berkeley.edu

Ajay Jain

UC Berkeley

ajayj@berkeley.edu

Pieter Abbeel

UC Berkeley

pabbeel@cs.berkeley.edu

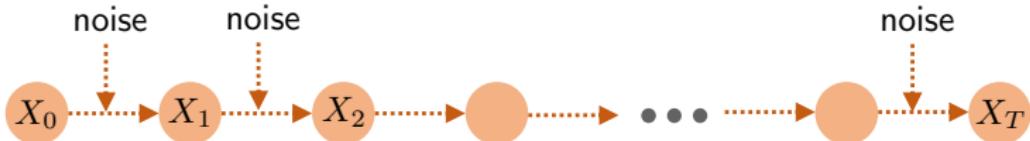
DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon

Stanford University

{tsong,chenlin,ermon}@cs.stanford.edu

DDPM vs. DDIM

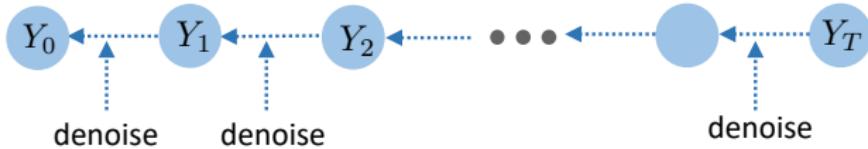


forward process: $X_0 \sim p_{\text{data}}$ (target distribution)

$$X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I_d) \quad t \geq 1$$

- $\beta_t := 1 - \alpha_t$ controls variance of injected noise

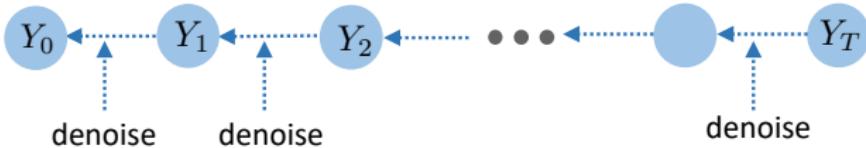
DDPM vs. DDIM



— Ho, Jain, Abbeel '20

1. A stochastic sampler: denoising diffusion probabilistic models
DDPM

DDPM vs. DDIM



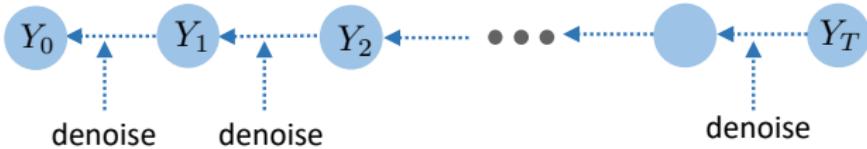
— Ho, Jain, Abbeel '20

1. A stochastic sampler: denoising diffusion probabilistic models
DDPM

$$Y_T \sim \mathcal{N}(0, I_d)$$

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\underbrace{Y_t + \eta_t^{\text{ddpm}} \textcolor{red}{s}_t(Y_t)}_{\text{deterministic}} + \underbrace{\sigma_t^{\text{ddpm}} \mathcal{N}(0, I_d)}_{\text{stochastic}} \right), \quad t = T, \dots, 1$$

DDPM vs. DDIM

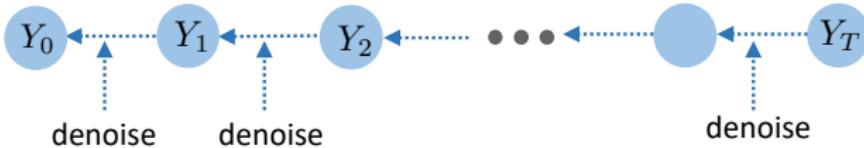


— Song, Meng, Ermon '20

— Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20

2. A deterministic sampler: **denoising diffusion implicit models**
DDIM; or probability flow ODE

DDPM vs. DDIM



— Song, Meng, Ermon '20

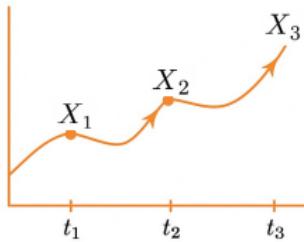
— Song, Sohl-Dickstein, Kingma, Kumar, Ermon, Poole '20

2. A deterministic sampler: denoising diffusion implicit models
DDIM; or probability flow ODE

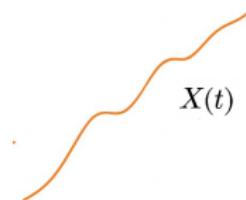
$$Y_T \sim \mathcal{N}(0, I_d)$$

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(\underbrace{Y_t + \eta_t^{\text{ddim}} s_t(Y_t)}_{\text{deterministic}} \right), \quad t = T, \dots, 1$$

Interpretations from lens of SDE/ODE

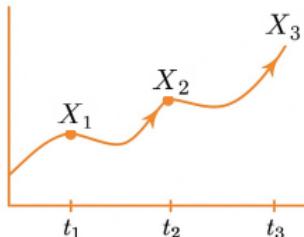


discrete-time

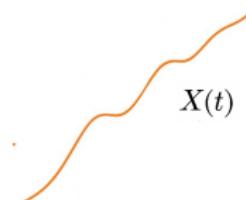


continuous-time

Interpretations from lens of SDE/ODE



discrete-time

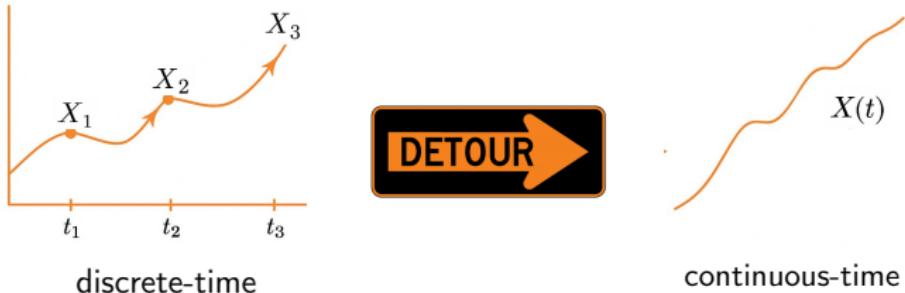


continuous-time

forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\implies dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

Interpretations from lens of SDE/ODE



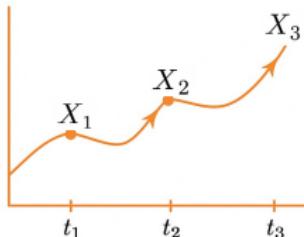
forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

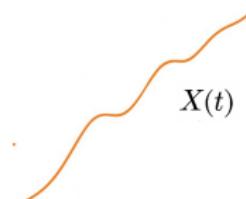
- \exists reverse-time SDE w/ same path distribution (Anderson '82)

$$dY_t = (Y_t + 2s_{T-t}^{\star}(Y_t)) \beta(T-t) dt + \sqrt{2\beta(T-t)} dW_t$$

Interpretations from lens of SDE/ODE



discrete-time



continuous-time

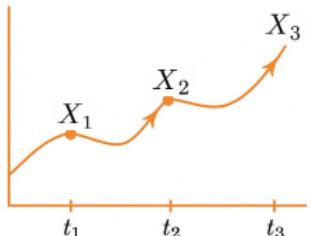
forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

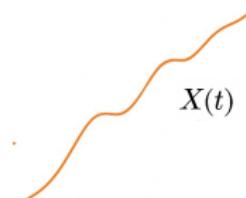
- \exists reverse-time SDE w/ same path distribution (Anderson '82)

time discretization
 \longrightarrow DDPM

Interpretations from lens of SDE/ODE



discrete-time



continuous-time

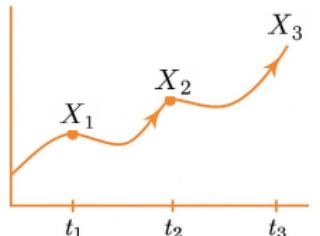
forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

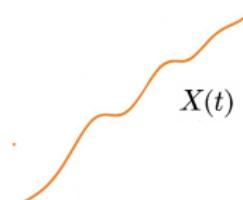
- \exists reverse-time ODE w/ same *marginal dist* (Song et al. '20)

$$dY_t = (Y_t + s_{T-t}^*(Y_t)) \beta(T - t) dt$$

Interpretations from lens of SDE/ODE



discrete-time



continuous-time

forward process

$$X_t = \sqrt{1 - \beta_t} X_{t-1} + \sqrt{\beta_t} \mathcal{N}(0, I_d)$$
$$\Rightarrow dX_t = -\frac{1}{2} \beta(t) X_t dt + \sqrt{\beta(t)} dW_t \quad (\text{SDE})$$

- \exists reverse-time ODE w/ same *marginal dist* (Song et al. '20)

time discretization $\xrightarrow{\hspace{1cm}}$ DDIM

Key takeaway: in continuous-time limits, sampling is feasible once exact score functions are available

- *almost no restriction on target data distributions*

Key takeaway: in continuous-time limits, sampling is feasible once exact score functions are available

- *almost no restriction on target data distributions*

Questions:

- what happens in discrete time? — effect of discretization error
- what if we only have imperfect scores? — effect of score error

A small sample of convergence theory

- Lee, Lu, Tan '22
- Chen, Chewi, Li, Li, Salim, Zhang '22
- Chen, Lee, Lu '22
- Lee, Lu, Tan '23
- Chen, Daras, Dimakis '23
- Chen, Chewi, Lee, Li, Lu, Salim '23
- Benton, De Bortoli, Doucet, Deligiannidis '23
- Li, Wei, Chen, Chi '23
- Benton, Deligiannidis, Doucet '23
- Cheng, Lu, Tan, Xie '23
- Tang '23
- Li, Wei, Chi, Chen '24
- Li, Yan '24a, '24b
- Azangulov, Deligiannidis, Rousseau '24
- Potapchik, Azangulov, Deligiannidis '24
- Huang, Wei, Chen '24
- Gao, Zhu '24
- Huang, Huang, Lin '24
- Li, Jiao '24
- Li, Di, Gu '24
- Liang, Ju, Liang, Shroff '24
- Tang, Zhao '24
- Liang, Huang, Chen '25
- Li, Cai, Wei '25
- Tang, Yan '25
- Yu, Yu '25
- Gentiloni-Silveri, Ocello '25
- ...

Assumptions: target data distribution p_{data}

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

Assumptions: target data distribution p_{data}

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large

Assumptions: target data distribution p_{data}

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large
- very general: *no need of log-concavity, smoothness, etc*

Assumptions: target data distribution p_{data}

$$\mathbb{P}(\|X_0\|_2 \leq T^{c_R}) = 1 \text{ for arbitrarily large const } c_R > 0$$

- support size can be very large
- very general: *no need of log-concavity, smoothness, etc*
- can also be replaced by $\mathbb{E}[\|X_0\|_2] \leq T^{c_M}$ for large const c_M

Assumptions: score estimates $\{s_t(\cdot)\}$

- ℓ_2 score estimation error: $s_t^*(X) := \nabla \log p_{X_t}(X)$,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[\|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

Assumptions: score estimates $\{s_t(\cdot)\}$

- ℓ_2 score estimation error: $s_t^*(X) := \nabla \log p_{X_t}(X)$,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[\|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption

Assumptions: score estimates $\{s_t(\cdot)\}$

- ℓ_2 score estimation error: $s_t^*(X) := \nabla \log p_{X_t}(X)$,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[\|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption
- suffices for DDPM **but not DDIM** (counterexample in [Li et al. '24](#))

Assumptions: score estimates $\{s_t(\cdot)\}$

- ℓ_2 score estimation error: $s_t^*(X) := \nabla \log p_{X_t}(X)$,

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[\|s_t(X) - s_t^*(X)\|_2^2 \right] \leq \varepsilon_{\text{score}}^2$$

- much weaker than *pointwise* score error assumption
- suffices for DDPM **but not DDIM** (counterexample in [Li et al. '24](#))
- Jacobian estimation error (for DDIM only):

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{X \sim p_{X_t}} \left[\left\| \frac{\partial s_t}{\partial x}(X) - \frac{\partial s_t^*}{\partial x}(X) \right\| \right] \leq \varepsilon_{\text{Jacobi}}$$

Assumptions: learning rates

$$X_0 \sim p_{\text{data}}, \quad X_t = \sqrt{\alpha_t} X_{t-1} + \sqrt{1 - \alpha_t} \mathcal{N}(0, I_d)$$

- **learning rates:** for some consts $c_0, c_1 > 0$,

$$1 - \alpha_1 = \frac{1}{T^{c_0}}$$

$$1 - \alpha_t = \underbrace{\frac{c_1 \log T}{T} \min \left\{ \left(1 - \alpha_1\right) \left(1 + \frac{c_1 \log T}{T}\right)^t, 1 \right\}}_{\text{2 phases: exp growth} \rightarrow \text{flat}}$$

A glimpse of convergence theory (up to log factor)

DDPM: $\text{KL}(p_{X_1} \parallel p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

A glimpse of convergence theory (up to log factor)

DDPM: $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

- $\underbrace{\text{iteration complexity: } d/\varepsilon^2}_{\text{to yield } \text{KL} \leq \varepsilon^2}$ (assuming accurate scores)
- $\underbrace{\text{iteration complexity: } d/\varepsilon}_{\text{to yield } \text{TV} \leq \varepsilon}$ (assuming accurate scores)

A glimpse of convergence theory (up to log factor)

DDPM: $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

- $\underbrace{\text{iteration complexity: } d/\varepsilon^2}_{\text{to yield } \text{KL} \leq \varepsilon^2}$ (assuming accurate scores)
- $\underbrace{\text{iteration complexity: } d/\varepsilon}_{\text{to yield } \text{TV} \leq \varepsilon}$ (assuming accurate scores)
 - Pinsker inequality ($\text{TV} \leq \sqrt{\frac{1}{2}\text{KL}}$) is loose when bounding TV

A glimpse of convergence theory (up to log factor)

DDPM: $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

- $\underbrace{\text{iteration complexity: } d/\varepsilon^2}_{\text{to yield } \text{KL} \leq \varepsilon^2}$ (assuming accurate scores)
- $\underbrace{\text{iteration complexity: } d/\varepsilon}_{\text{to yield } \text{TV} \leq \varepsilon}$ (assuming accurate scores)
- **stability:** degrades gracefully as $\varepsilon_{\text{score}} \uparrow$

A glimpse of convergence theory (up to log factor)

DDPM: $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

- $\underbrace{\text{iteration complexity: } d/\varepsilon^2}_{\text{to yield } \text{KL} \leq \varepsilon^2}$ (assuming accurate scores)
- $\underbrace{\text{iteration complexity: } d/\varepsilon}_{\text{to yield } \text{TV} \leq \varepsilon}$ (assuming accurate scores)
- **stability:** degrades gracefully as $\varepsilon_{\text{score}} \uparrow$ and $\varepsilon_{\text{Jacobi}} \uparrow$
- **general data distribution:** no need of smoothness, log-concavity

A glimpse of convergence theory (up to log factor)

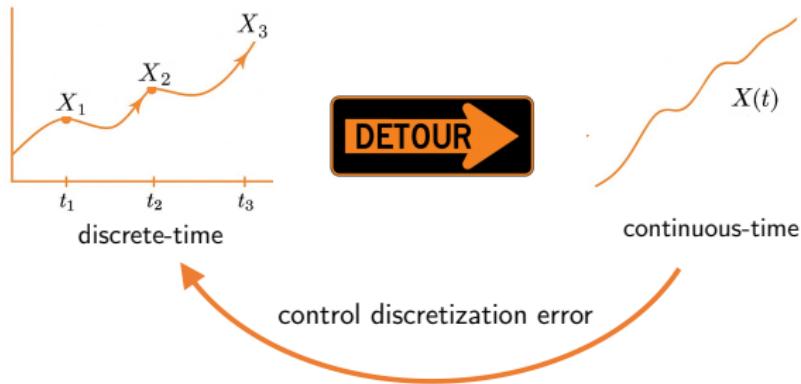
DDPM: $\text{KL}(p_{X_1} \| p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}^2$ (Benton et al. '23)
 $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \varepsilon_{\text{score}}$ (Li, Yan '24)

DDIM: $\text{TV}(p_{X_1}, p_{Y_1}) \lesssim d/T + \sqrt{d} \varepsilon_{\text{score}} + d\varepsilon_{\text{Jacobi}}$ (Li et al. '24)

- **iteration complexity:** d/ε^2 (assuming accurate scores)
to yield $\text{KL} \leq \varepsilon^2$
- **iteration complexity:** d/ε (assuming accurate scores)
to yield $\text{TV} \leq \varepsilon$
- **stability:** degrades gracefully as $\varepsilon_{\text{score}} \uparrow$ and $\varepsilon_{\text{Jacobi}} \uparrow$
- **general data distribution:** no need of smoothness, log-concavity

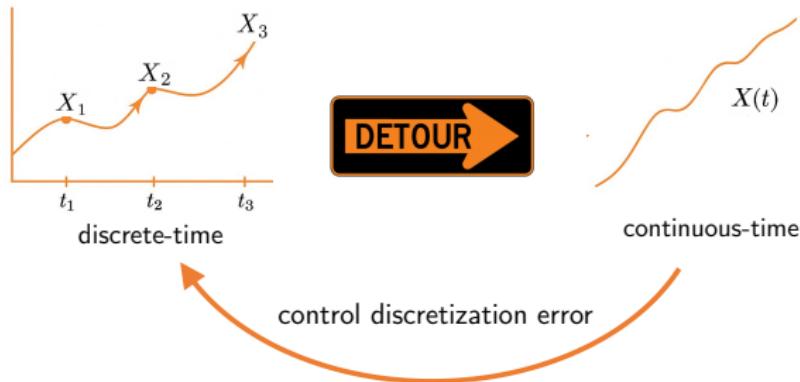
Analysis strategy # 1 (for DDPM)

- Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24
- Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24



Analysis strategy # 1 (for DDPM)

- Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24
- Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24

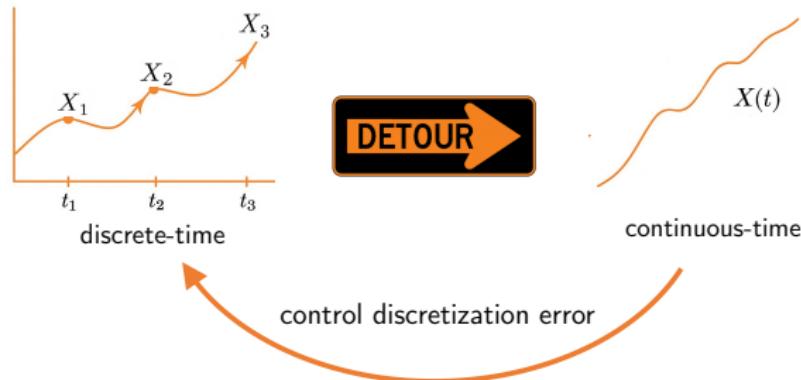


Analogy: (stochastic) gradient descent vs. gradient flow, TD learning via ODE

yields state-of-the-art **KL-based theory** for DDPM!

Analysis strategy # 1 (for DDPM)

- Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24
- Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24



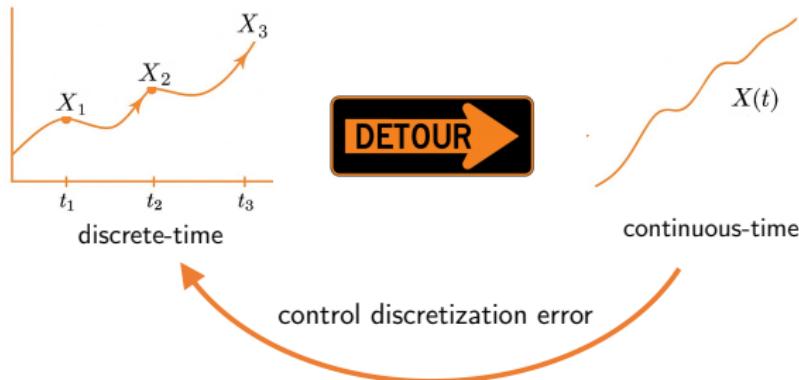
2 key steps:

- apply change of measure (e.g. Girsanov thm) to show

$$\text{KL}(P^{\text{true}} \| P^{\text{ddpm}}) \leq \int w(t) \underbrace{\mathbb{E} \left[\left\| \text{drift}^{\text{true}}(t) - \text{drift}^{\text{ddpm}}(t) \right\| \right]^2}_{\text{score error} + \text{discretization error}} dt + \text{small-term}$$

Analysis strategy # 1 (for DDPM)

- Chen, Chewi, Li, Li, Salim, Zhang '22, Chen, Lee, Lu '22, Tang, Zhao '24
- Benton, De Bortoli, Doucet, Deligiannidis '23, Huang, Wei, Chen '24



2 key steps:

- leverage stochastic localization to characterize

$$\text{discretization error} \xleftrightarrow{\text{link}} \mathbb{E}[\text{Cov}(X_0 | X_t)]$$

Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

yields state-of-the-art **TV-based theory** for DDIM & DDPM!

Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0$$

Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

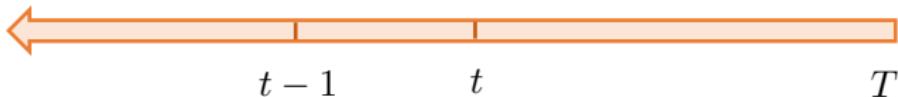
$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \quad \iff \quad \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$

Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$

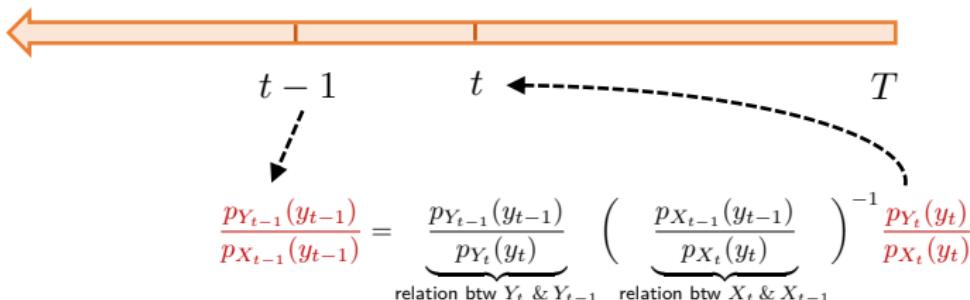


Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$

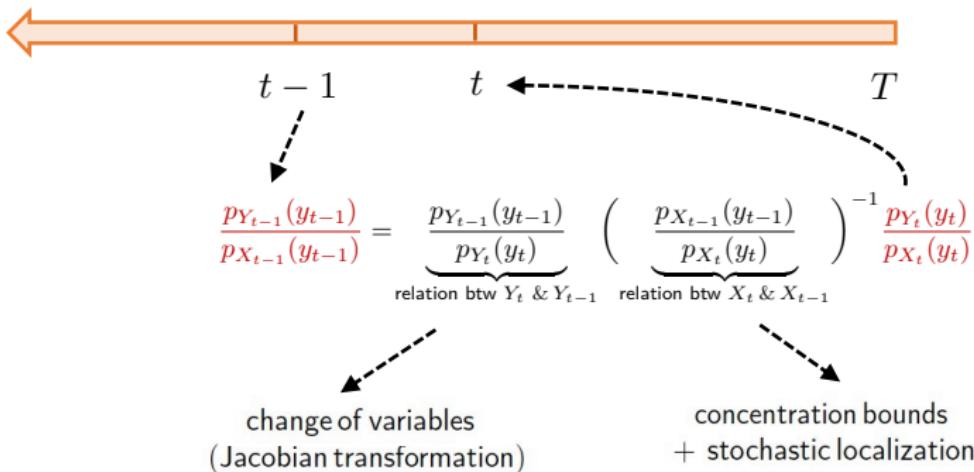


Analysis strategy # 2 (for DDIM & DDPM)

- Li, Wei, Chen, Chi '24, Li, Wei, Chi, Chen '24
- Li, Yan '24, Liang, Huang, Chen '25

Tackle discrete-time process directly & track changes of TV distance

$$\text{TV}(p_{X_t}, p_{Y_t}) \approx 0 \iff \frac{p_{Y_t}(y_t)}{p_{X_t}(y_t)} \approx 1 \quad \forall y_t \in \mathcal{E}_t \text{ (some "typical" set)}$$



Part 2: adaptation to (unknown) low dimensionality



Jiadong Liang
UPenn



Zhihan Huang
UPenn



Yuting Wei
UPenn

Recap: theory for mainstream diffusion models

Denoising Diffusion Probabilistic Models

Jonathan Ho
UC Berkeley

jonathanho@berkeley.edu

Ajay Jain
UC Berkeley

ajayj@berkeley.edu

Pieter Abbeel
UC Berkeley

pabbeel@cs.berkeley.edu

DENOISING DIFFUSION IMPLICIT MODELS

Jiaming Song, Chenlin Meng & Stefano Ermon
Stanford University
`{tsong,chenlin,ermon}@cs.stanford.edu`

Theorem 1 (Li, Wei, Chi, Chen '24, Li, Yan '24)

With perfect scores, both DDIM & DDPM yield $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in
 $\tilde{O}(\textcolor{red}{d}/\varepsilon)$ iterations

- d : ambient dimension

d/ε **iterations are too slow . . .**



ImageNet: $d = 150,528$ pixels per image

d/ε **iterations are too slow . . .**



ImageNet: $d = 150,528$ pixels per image (so $\frac{d}{\varepsilon} > 10^6$ for moderate ε)

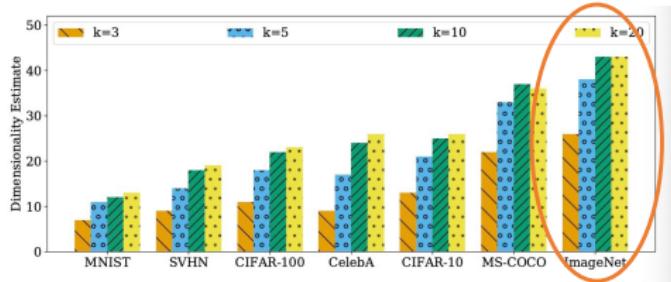
d/ε **iterations are too slow** . . .



ImageNet: $d = 150,528$ pixels per image (so $\frac{d}{\varepsilon} > 10^6$ for moderate ε)

In practice, DDIM/DDPM yield good samples in hundreds (or tens) of iterations . . .

d/ε iterations are too slow ...



ImageNet: $d = 150,528$ pixels per image (so $\frac{d}{\varepsilon} > 10^6$ for moderate ε)
 $k = 43$ intrinsic dimension (Pope et al. '21)

In practice, DDIM/DDPM yield good samples in hundreds (or tens) of iterations ...

Can diffusion models adapt to intrinsic low dimensionality?

Intrinsic dimension

The target distribution p_{data} is said to have **intrinsic dimension k** if

$$\log \underbrace{N^{\text{cover}}(\text{support}(p_{\text{data}}), \|\cdot\|_2, \varepsilon_0)}_{\text{covering number of support of } p_{\text{data}}} \lesssim k \log \frac{1}{\varepsilon_0}$$

Intrinsic dimension

The target distribution p_{data} is said to have **intrinsic dimension k** if

$$\log \underbrace{N^{\text{cover}}(\text{support}(p_{\text{data}}), \|\cdot\|_2, \varepsilon_0)}_{\text{covering number of support of } p_{\text{data}}} \lesssim k \log \frac{1}{\varepsilon_0}$$

- k -dimensional linear subspace
- low-dimensional manifold
- doubling dimension, Minkowski dimension
- ...

— see *Huang, Wei, Chen '24*

Assumptions

- **minimal data assumptions:**

$$\mathbb{P}(\|X_0\|_2 \leq \underbrace{T^{c_R}}_{\text{polynomially large diameter}}) = 1$$

for arbitrarily large constant $c_R > 0$

- **perfect score estimates:** $s_t(\cdot) = \nabla \log p_{X_t}(\cdot)$
→ not needed; only to simplify presentation

Convergence theory in total variation

Theorem 2 (Liang, Huang, Chen '24)

Both DDPM & DDIM (their original form) yield $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(k/\varepsilon) \text{ iterations}$$

— concurrent work [Li, Yan '25](#), [Tang, Yan '25](#)

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \eta_t^{\text{ddim}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \eta_t^{\text{ddpm}} s_t(Y_t) + \sigma_t^{\text{ddpm}} \mathcal{N}(0, I_d) \right)$$

Convergence theory in total variation

Theorem 2 (Liang, Huang, Chen '24)

Both DDPM & DDIM (their original form) yield $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(k/\varepsilon) \text{ iterations}$$

— concurrent work [Li, Yan '25](#), [Tang, Yan '25](#)

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \frac{1 - \alpha_t}{1 + \sqrt{\frac{\alpha_t - \bar{\alpha}_t}{1 - \bar{\alpha}_t}}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + (1 - \alpha_t) s_t(Y_t) + \sqrt{\frac{(1 - \alpha_t)(\alpha_t - \bar{\alpha}_t)}{1 - \bar{\alpha}_t}} \mathcal{N}(0, I_d) \right)$$

$$\text{where } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

Convergence theory in total variation

Theorem 2 (Liang, Huang, Chen '24)

Both DDPM & DDIM (their original form) yield $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(k/\varepsilon) \text{ iterations}$$

— concurrent work [Li, Yan '25](#), [Tang, Yan '25](#)

$$\text{DDIM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \frac{1 - \alpha_t}{1 + \sqrt{\frac{\alpha_t - \bar{\alpha}_t}{1 - \bar{\alpha}_t}}} s_t(Y_t) \right)$$

$$\text{DDPM: } Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} \left(Y_t + (1 - \alpha_t) s_t(Y_t) + \sqrt{\frac{(1 - \alpha_t)(\alpha_t - \bar{\alpha}_t)}{1 - \bar{\alpha}_t}} \mathcal{N}(0, I_d) \right)$$

$$\text{where } \bar{\alpha}_t = \prod_{i=1}^t \alpha_i$$

- originally derived to optimize variational lower bounds!

Convergence theory in KL divergence

Theorem 3 (Huang, Wei, Chen '24)

DDPM sampler (its original form) yields $\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(k/\varepsilon) \text{ iterations}$$

- prior work *Li and Yan '24, Azangulov et al. '24*
- concurrent work *Potapchik et al.'24*

Convergence theory in KL divergence

Theorem 3 (Huang, Wei, Chen '24)

DDPM sampler (its original form) yields $\text{KL}(p_{X_1} \parallel p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(k/\varepsilon) \text{ iterations}$$

- prior work *Li and Yan '24, Azangulov et al. '24*
- concurrent work *Potapchik et al.'24*

- optimal scaling in k

Interpretation from lens of SDE/ODE

reverse-time SDE (same distribution as X_t):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2}dB_t$$

Interpretation from lens of SDE/ODE

reverse-time SDE (same distribution as X_t):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2}dB_t$$

Tweedie's formula \Downarrow

$$dY_t = (\underbrace{c_{1,t}Y_t}_{\text{linear drift}} + \underbrace{c_{2,t}\mathbb{E}[X_0 \mid X_{T-t} = Y_t]}_{\text{cond. mean of } X_0})dt + \sqrt{2}dB_t$$

- **key enabler:** $\mathbb{E}[X_0 \mid X_t]$ is “projection” onto low-dimensional structure

Interpretation from lens of SDE/ODE

reverse-time SDE (same distribution as X_t):

$$dY_t = (Y_t + 2s_{T-t}^*(Y_t))dt + \sqrt{2}dB_t$$

Tweedie's formula



$$dY_t = (\underbrace{c_{1,t}Y_t}_{\text{linear drift}} + c_{2,t}\underbrace{\mathbb{E}[X_0 | X_{T-t} = Y_t]}_{\text{cond. mean of } X_0})dt + \sqrt{2}dB_t$$

time-discretize



DDIM or DDPM

- **key enabler:** $\mathbb{E}[X_0 | X_t]$ is “projection” onto low-dimensional structure
- **discretization scheme matters:** time-discretize carefully to retain low-dimensional adaptation

Crucial choices of coefficients: a lower bound

DDPM & DDIM updates take the form

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} (Y_t + \eta_t s_t(Y_t) + \sigma_t \mathcal{N}(0, I_d))$$

Crucial choices of coefficients: a lower bound

DDPM & DDIM updates take the form

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} (Y_t + \eta_t s_t(Y_t) + \sigma_t \mathcal{N}(0, I_d))$$

Theorem 4 (Liang, Huang, Chen '24)

Even when starting from $X_t = Y_t$, one can have

$$\text{TV}(X_{t-1}, Y_{t-1}) \gtrsim \sqrt{d} \cdot \left| \frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \left(1 - \frac{\eta_t}{1 - \bar{\alpha}_t} \right)^2 + \frac{\sigma_t^2}{\alpha_t - \bar{\alpha}_t} - 1 \right|$$

Crucial choices of coefficients: a lower bound

DDPM & DDIM updates take the form

$$Y_{t-1} = \frac{1}{\sqrt{\alpha_t}} (Y_t + \eta_t s_t(Y_t) + \sigma_t \mathcal{N}(0, I_d))$$

Theorem 4 (Liang, Huang, Chen '24)

Even when starting from $X_t = Y_t$, one can have

$$\text{TV}(X_{t-1}, Y_{t-1}) \gtrsim \sqrt{d} \cdot \left| \frac{1 - \bar{\alpha}_t}{\alpha_t - \bar{\alpha}_t} \left(1 - \frac{\eta_t}{1 - \bar{\alpha}_t} \right)^2 + \frac{\sigma_t^2}{\alpha_t - \bar{\alpha}_t} - 1 \right|$$

To avoid scaling in d , one needs red term ≈ 0

- both DDIM and DDPM satisfy this!

Part 3: acceleration via higher-order approximation



Gen Li*
CUHK



Yuchen Zhou*
UIUC



Yu Huang*
UPenn



Timofey Efimov
CMU



Yuting Wei
UPenn



Yuejie Chi
Yale

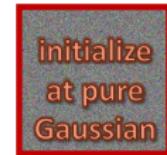
DDPM and DDIM are still slow ...

Low sampling speed!

100s-1000s steps



• • •



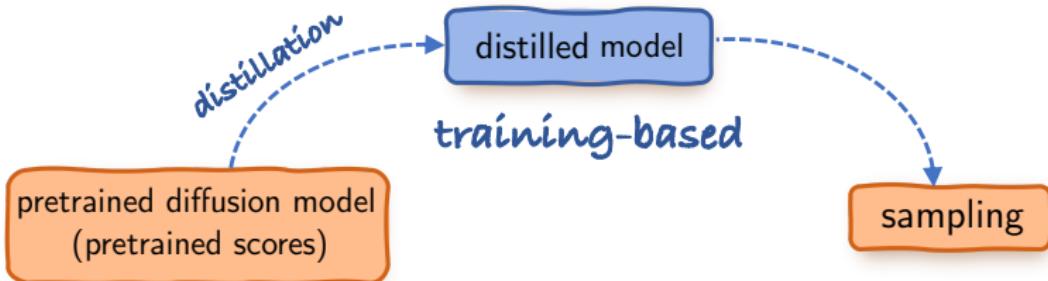
— Song, Meng, Ermon '20

DDPM and DDIM are still slow ...



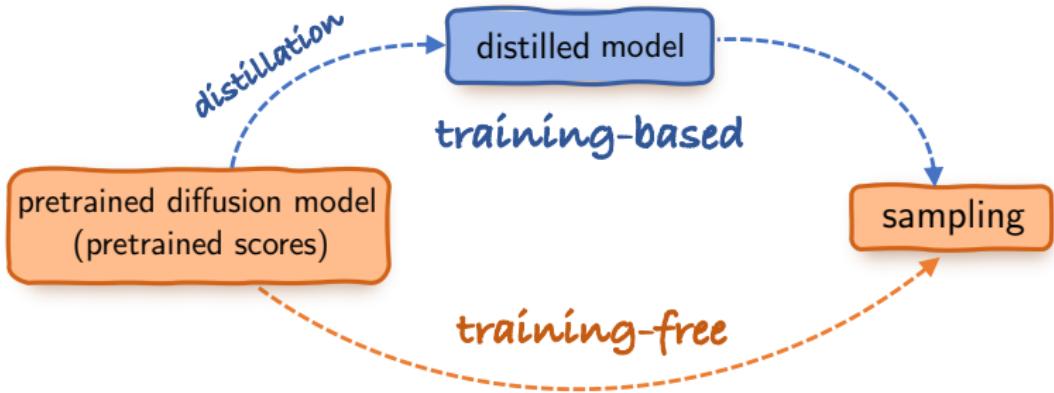
50K 32×32 images: DDPM (20h) vs. single-step GANs (< 1min)

— Song, Meng, Ermon '20



- **training-based:** distill pre-trained diffusion model into another model that can be executed rapidly
 - e.g., progressive distillation, consistency model

requires additional training



- **training-free:** directly invoke pre-trained score estimates for sampling w/o additional training
 - e.g., DPM-Solver/++ (Lu et al. '22), UniPC (Zhao et al. '23), ...

*Can we design a **training-free** sampler that is
provably faster than DDIM/DDPM?*

Discretization \longleftrightarrow approximation

A starting point: equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_\gamma^*(Y_\gamma^{\text{ode}}) d\gamma}_{\text{represent } Y_t}$$

where $s_\gamma^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

Discretization \longleftrightarrow approximation

A starting point: equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_\gamma^*(Y_\gamma^{\text{ode}}) d\gamma}_{\text{infinitely many score evaluations!}}$$

where $s_\gamma^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

Discretization \longleftrightarrow approximation

A starting point: equiv solution to probability flow ODE

$$\underbrace{Y_{\bar{\alpha}_{t-1}}^{\text{ode}}}_{\text{represent } Y_{t-1}} = \frac{1}{\sqrt{\alpha_t}} \underbrace{Y_{\bar{\alpha}_t}^{\text{ode}}}_{\text{represent } Y_t} + \underbrace{\int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} s_\gamma^*(Y_\gamma^{\text{ode}}) d\gamma}_{\text{infinitely many score evaluations!}}$$

where $s_\gamma^*(x) := \nabla \log p_{\sqrt{\gamma}X_0 + \sqrt{1-\gamma}\mathcal{N}(0, I_d)}(x)$

- can we approximate the integral by a few score evals?

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$$\implies Y_{t-1} \approx \frac{1}{\sqrt{\alpha_t}} \left(Y_t + \frac{1 - \alpha_t}{2} s_t(Y_t) \right) \quad \text{original DDIM}$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

refined approximation?

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

refined approximation?

$$s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_t(Y_t) + \frac{\gamma - \bar{\alpha}_t}{\bar{\alpha}_t - \bar{\alpha}_{t+1}} (s_t(Y_t) - s_{t+1}(Y_{t+1}))$$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

2nd order approx: (Li, Huang, Efimov, Wei, Chi, Chen '24)

$$\sqrt{\alpha_t} Y_{t-1} \approx Y_t + \frac{1 - \alpha_t}{2} s_t(Y_t) + \frac{(1 - \alpha_t)^2}{4(1 - \alpha_{t+1})} (s_t(Y_t) - \sqrt{\alpha_{t+1}} s_{t+1}(Y_{t+1}))$$

— similar in spirit to DPM-Solver-2 (Lu et al '22)

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

2nd order approx: (Li, Huang, Efimov, Wei, Chi, Chen '24)

$\frac{\text{poly}(d)}{\sqrt{\varepsilon}}$ iterations; 2 score evals per iteration

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

even higher-order approximation?

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

even higher-order approximation? for order K :

$$\frac{1}{\gamma^{3/2}} s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx \sum_{0 \leq i < K} \psi_i(\gamma) \frac{s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})}{(\gamma_{t,i})^{3/2}}$$

- K anchor points: $\gamma_{t,0}, \dots, \gamma_{t,K-1}$

$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approximate by?}} \, d\gamma$$

1st order approx: $s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx s_{\bar{\alpha}_t}^\star(Y_{\bar{\alpha}_t}^{\text{ode}}) \approx s_t(Y_t)$

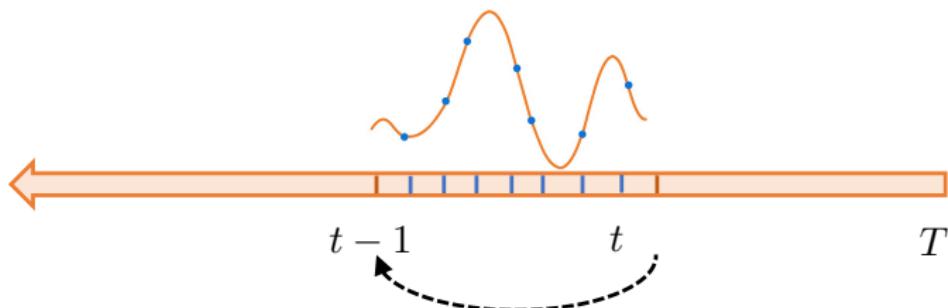
$\frac{d}{\varepsilon}$ iterations; 1 score eval per iteration (DDIM)

even higher-order approximation? for order K :

$$\frac{1}{\gamma^{3/2}} s_\gamma^\star(Y_\gamma^{\text{ode}}) \approx \sum_{0 \leq i < K} \psi_i(\gamma) \frac{s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})}{(\gamma_{t,i})^{3/2}}$$

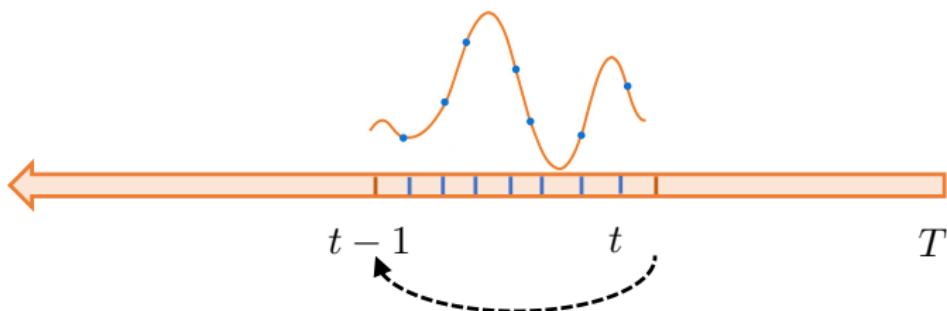
- K anchor points: $\gamma_{t,0}, \dots, \gamma_{t,K-1}$
- Lagrange basis polynomial: $\psi_i(\gamma) := \frac{\prod_{i': i' \neq i} (\gamma - \gamma_{t,i'})}{\prod_{i': i' \neq i} (\gamma_{t,i} - \gamma_{t,i'})}$

Proposed K -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} d\gamma$$

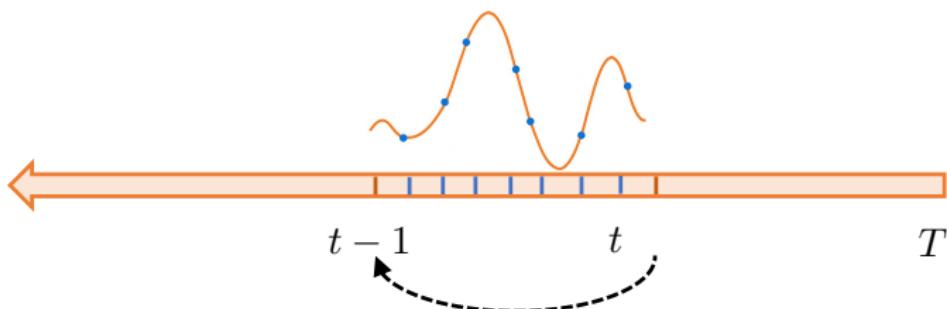
Proposed K -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} d\gamma$$

- successively, alternately refine $Y_{\gamma_{t,i}}^{\text{ode}}$ and $s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})$

Proposed K -th order sampler (Li et al. '25)



$$Y_{\bar{\alpha}_{t-1}}^{\text{ode}} = \frac{1}{\sqrt{\alpha_t}} Y_{\bar{\alpha}_t}^{\text{ode}} + \int_{\bar{\alpha}_t}^{\bar{\alpha}_{t-1}} \frac{1}{\sqrt{\gamma^3}} \underbrace{s_\gamma^\star(Y_\gamma^{\text{ode}})}_{\text{approx by deg-}(K-1) \text{ Lagrange polynomials}} \, d\gamma$$

- successively, alternately refine $Y_{\gamma_{t,i}}^{\text{ode}}$ and $s_{\gamma_{t,i}}(Y_{\gamma_{t,i}}^{\text{ode}})$

K score evals per iteration; $\tilde{O}(1)$ rounds of refinements

Convergence theory for our accelerated sampler

Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any $K = O(1)$. With perfect scores, our accelerated deterministic sampler yields $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(d^{1+2/K}/\varepsilon^{1/K}) \text{ iterations}$$

Convergence theory for our accelerated sampler

Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any $K = O(1)$. With perfect scores, our accelerated deterministic sampler yields $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(d^{1+2/K}/\varepsilon^{1/K}) \text{ iterations}$$

- **# score function evaluations:** $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$

Convergence theory for our accelerated sampler

Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any $K = O(1)$. With perfect scores, our accelerated deterministic sampler yields $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(d^{1+2/K}/\varepsilon^{1/K}) \text{ iterations}$$

- **# score function evaluations:** $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$
- outperforms vanilla DDIM (d/ε)
 - substantially improved ε -dependency
 - almost no loss in d -dependency;

Convergence theory for our accelerated sampler

Theorem 5 (Li, Zhou, Wei, Chen '25)

Consider any $K = O(1)$. With perfect scores, our accelerated deterministic sampler yields $\text{TV}(p_{X_1}, p_{Y_1}) \leq \varepsilon$ in

$$\tilde{O}(d^{1+2/K}/\varepsilon^{1/K}) \text{ iterations}$$

- **# score function evaluations:** $\frac{d^{1+o(1)}}{\varepsilon^{1/K}}$
- outperforms vanilla DDIM (d/ε)
 - substantially improved ε -dependency
 - almost no loss in d -dependency;
- **minimal assumptions** on data distributions
 - see also Huang et al. '24, '25 (Runge-Kutta; stronger assumptions)

Summary

- nonasymptotic convergence theory for diffusion models
- adaptation to unknown low dimensionality
- provable training-free acceleration

Summary

- nonasymptotic convergence theory for diffusion models
- adaptation to unknown low dimensionality
- provable training-free acceleration

Future directions:

- end-to-end theory that accounts for score learning + sampling?
- adaptive improvement under stylized statistical models
- design of high-order stochastic samplers
- parallelization
- discrete-valued/structured problems

Papers

- "A sharp convergence theory for the probability flow ODEs of diffusion models,"
G. Li, Y. Wei, Y. Chi, Y. Chen, arXiv:2408.02320, 2024
- "Towards non-asymptotic convergence for diffusion-based generative models,"
G. Li, Y. Wei, Y. Chen, Y. Chi, arXiv:2306.09251, ICLR 2024
- "Low-dimensional adaptation of diffusion models: convergence in total variation,"
J. Liang, Z. Huang, Y. Chen, arXiv:2501.12982, COLT 2025
- "Denoising diffusion probabilistic models are optimally adaptive to unknown low dimensionality," Z. Huang, Y. Wei, Y. Chen, arXiv:2410.18784, 2024
- "Accelerating convergence of score-based diffusion models, provably," G. Li*,
Y. Huang*, T. Efimov, Y. Wei, Y. Chi, Y. Chen, arXiv:2403.03852, ICML 2024
- "Stochastic Runge-Kutta methods: Provable acceleration of diffusion models,"
Y. Wu, Y. Chen, Y. Wei, arXiv:2410.04760, 2024
- "Faster diffusion models via higher-order approximation," G. Li*, Y. Zhou*,
Y. Wei, Y. Chen, arXiv:2506.24042, 2025

(* =equal contributions)