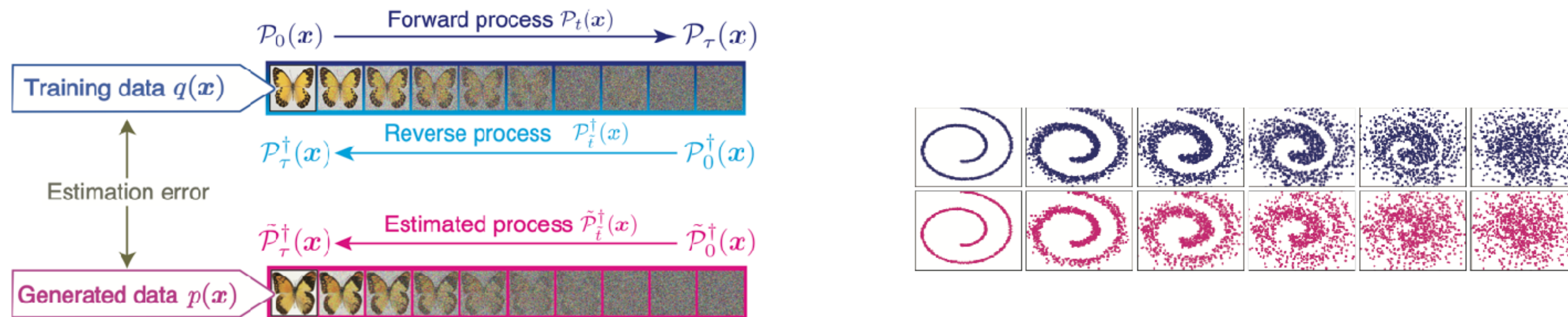


Optimal transport and thermodynamics for the learning: Application to the diffusion model



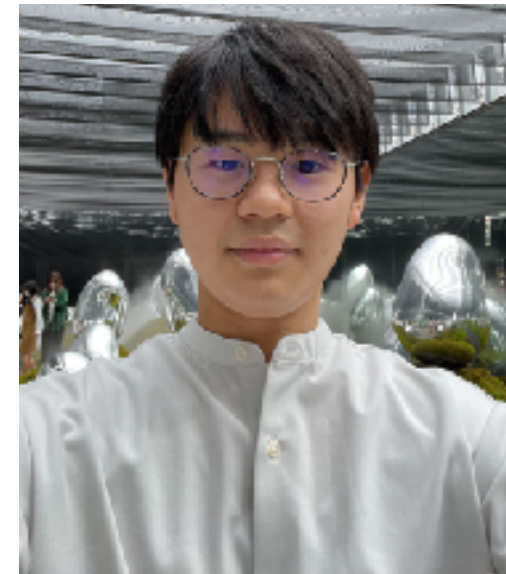
Sosuke Ito

Frontiers in Non-equilibrium Physics 2024, YITP, Jul. 18th, 2024

Reference and collaborators

Main topic (Diffusion model)

K. Ikeda, T. Uda, D. Okanohara and SI, arXiv:2407.04495.



Kotaro Ikeda (UTokyo)



Tomoya Uda (UTokyo)



Daisuke Okanohara (Preferred Networks Inc.)

Related topic (Thermodynamics and optimal transport)

SI, Information geometry, *Information Geometry* 7.Suppl 1, 441-483 (2024).

M. Nakazato and SI. *Phys. Rev. Res.* 3, 043093 (2021).

A. Dechant, S-I Sasa and SI. *Phys. Rev. Res.* 4, L012034 (2022).

A. Dechant, S-I Sasa and SI, *Phys. Rev. E.* 106, 024125 (2022).

K. Yoshimura, A. Kolchinsky, A. Dechant and SI. *Phys. Rev. Res.* 5, 013017 (2023).

Y. Fujimoto and SI, *Phys. Rev. Res.* 6, 013023 (2024).

K. Yoshimura and SI, *Phys. Rev. Res.* 6, L022057 (2024).

A. Kolchinsky, A. Dechant, K. Yoshimura and SI, arXiv:2206.14599.

R. Nagayama, K. Yoshimura, A. Kolchinsky and SI. arXiv: 2311.16569.

D. Sekizawa, SI, M. Oizumi, arXiv:2312.03489.

Collaborators:

Lab members (+alumni): Muka Nakazato, Kohei Yoshimura, Yuma Fujimoto, Artemy Kolchinsky, Ryan Nagayama
Andreas Dechant (KyotoU), Shin-ichi Sasa (KyotoU), Daiki Sekizawa (UTokyo), Masafumi Oizumi (UTokyo)

Outline

- Introduction: Generative models and diffusion models
- Stochastic thermodynamics based on optimal transport
- Main results: Speed-accuracy trade-off for the diffusion models

K. Ikeda, T. Uda, D. Okanojara and SI, arXiv:2407.04495.

Generative model

Stable diffusion (2022)

- Generative artificial intelligence
- Text-to-image model
- The diffusion models

✓ Prompt

Frontiers in Non-equilibrium Physics 2024

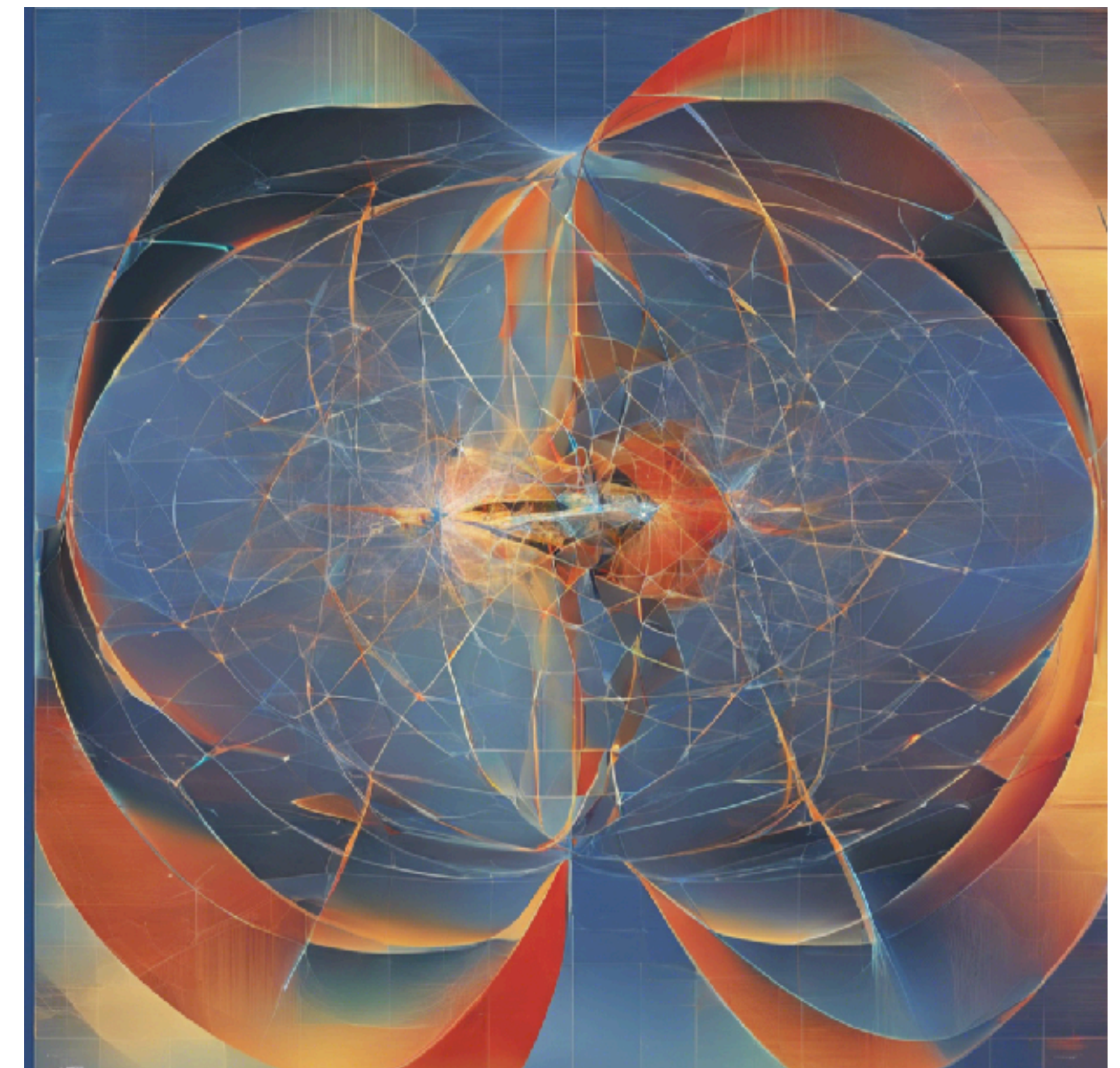
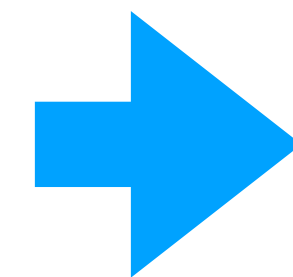
Generative model

Stable diffusion (2022)

- Generative artificial intelligence
- Text-to-image model
- The diffusion models

✓ Prompt

Frontiers in Non-equilibrium Physics 2024



Diffusion model - Original paper (2015)

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, PMLR, pp. 2256–2265 (2015).

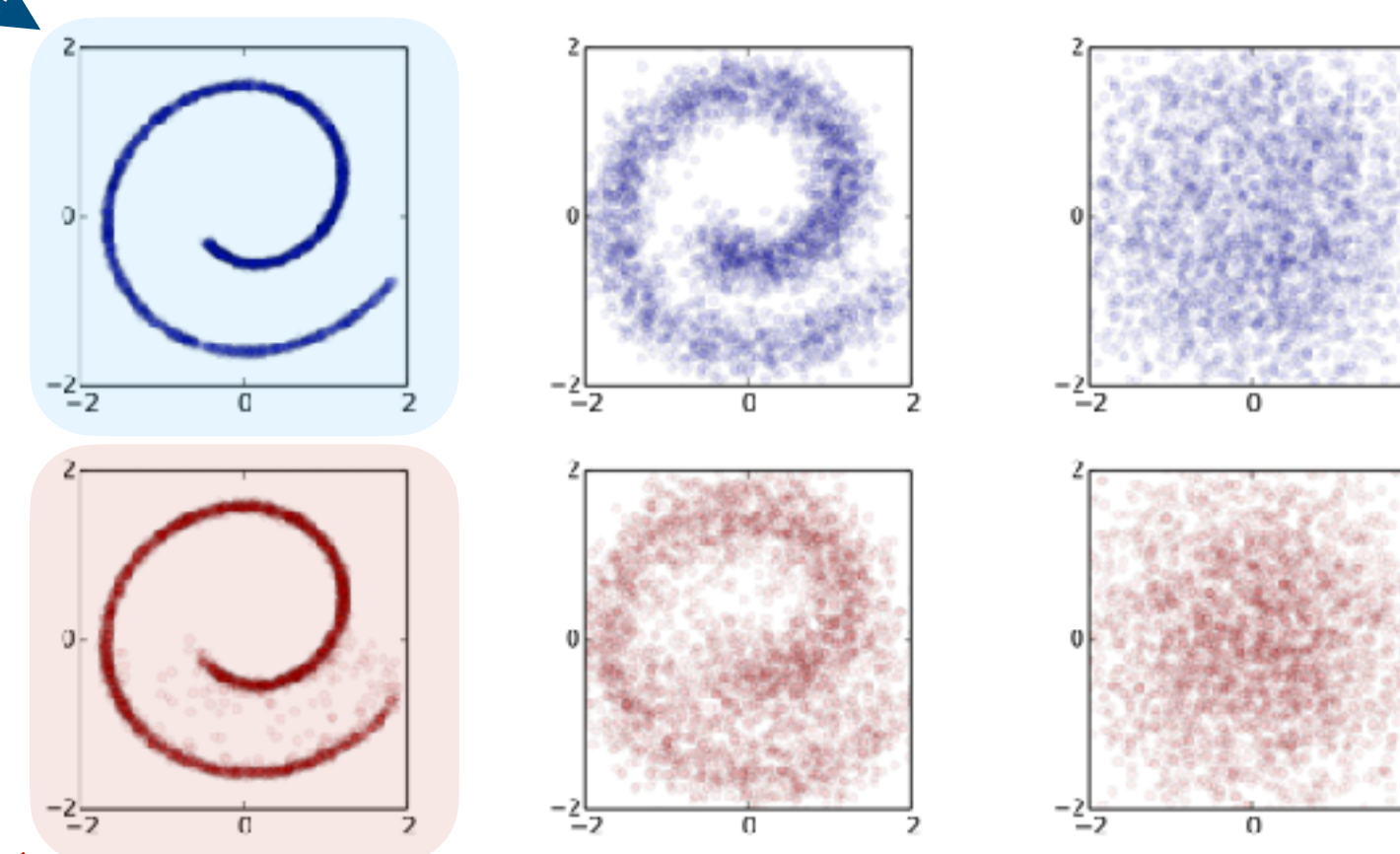
Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, Surya Ganguli Proceedings of the 32nd International Conference on Machine Learning, PMLR 37:2256-2265, 2015.

Abstract

A central problem in machine learning involves modeling complex data-sets using highly flexible families of probability distributions in which learning, sampling, inference, and evaluation are still analytically or computationally tractable. Here, we develop an approach that simultaneously achieves both flexibility and tractability. The essential idea, inspired by non-equilibrium statistical physics, is to systematically and slowly destroy structure in a data distribution through an iterative forward diffusion process. We then learn a reverse diffusion process that restores structure in data, yielding a highly flexible and tractable generative model of the data. This approach allows us to rapidly learn, sample from, and evaluate probabilities in deep generative models with thousands of layers or time steps, as well as to compute conditional and posterior probabilities under the learned model. We additionally release an open source reference implementation of the algorithm.

Training data



Forward diffusion process
[learning]

Reverse diffusion process
[data generation]

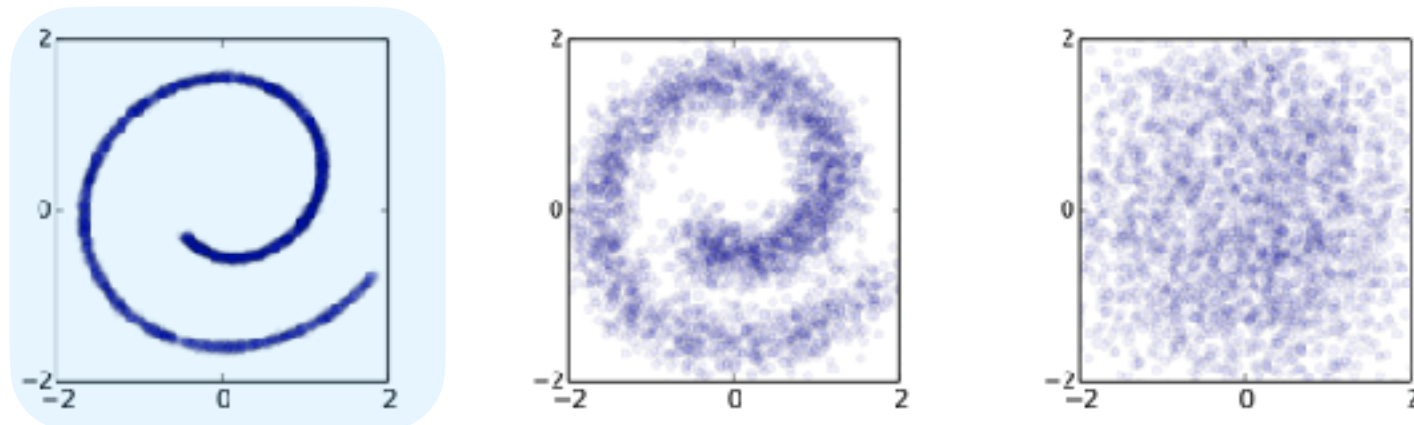
Generated data

Essential idea

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, PMLR, pp. 2256–2265 (2015).

Training data q

$$P_0(\mathbf{x}_0) = q(\mathbf{x}_0) \longrightarrow P_\tau(\mathbf{x}_N)$$



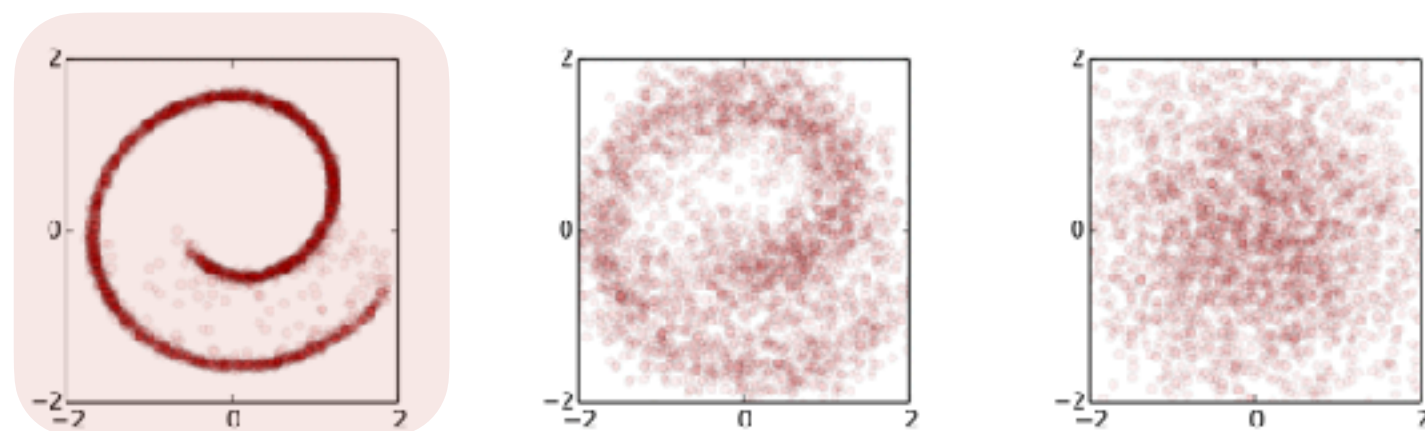
Forward diffusion process [learning]

$$\mathcal{P}^F(\{\mathbf{x}\}) = q(\mathbf{x}_0) \prod_i T_i(\mathbf{x}_{i+1} | \mathbf{x}_i)$$

Estimating the reverse process $\hat{T}_i^\dagger = T_i^\dagger$

$$\mathcal{P}^E(\{\mathbf{x}\}) = P_0^\dagger(\mathbf{x}_\tau) \prod_i T_i^\dagger(\mathbf{x}_i | \mathbf{x}_{i+1})$$

Generated data p



Reverse diffusion process [data generation]

$$p(\mathbf{x}_0) (\simeq q(\mathbf{x}_0)) \longleftarrow P_0^\dagger(\mathbf{x}_N) (\simeq P_\tau(\mathbf{x}_N))$$

Variants of the diffusion models

- Score-based generative model

Song, Y., Sohl-Dickstein, J., Kingma, D. P., Kumar, A., Ermon, S., & Poole, B. In International Conference on Learning Representations. (2021).

Score-based generative model

Fokker-Planck equation (Forward diffusion process)

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) P_t(\mathbf{x})) \quad \boldsymbol{\nu}_t(\mathbf{x}) = F_t(\mathbf{x}) - T_t \nabla \ln P_t(\mathbf{x})$$

Estimating $\hat{\boldsymbol{\nu}}_t(\mathbf{x}) = F_t(\mathbf{x}) - T_t s_t(\mathbf{x})$

via the score function $s_t = \nabla \ln P_t$

-Data generation by the reverse stochastic differential equation

$$\dot{\mathbf{x}}_{\tilde{t}} = F_{\tau-\tilde{t}}(\mathbf{x}_{\tilde{t}}) - 2\hat{\boldsymbol{\nu}}_{\tau-\tilde{t}}(\mathbf{x}_{\tilde{t}}) + \sqrt{2T_{\tau-\tilde{t}}}\boldsymbol{\xi}_{\tau-\tilde{t}}$$

($\tilde{t} = \tau - t$:Reversed time)

-Data generation by the ordinary differential equation (probability flow ODE)

$$\dot{\mathbf{x}}_{\tilde{t}} = -\hat{\boldsymbol{\nu}}_{\tau-\tilde{t}}(\mathbf{x}_{\tilde{t}})$$

Variants of the diffusion models

- Flow-based generative model

Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, International Conference on Learning Representations (2022)

Flow-based generative model

Continuity equation (Forward process)

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) P_t(\mathbf{x}))$$

Estimating the velocity field $\hat{\boldsymbol{\nu}}_t(\mathbf{x}) = \boldsymbol{\nu}_t(\mathbf{x})$

-Data generation by the ordinary differential equation

$$\dot{\mathbf{x}}_{\tilde{t}} = -\hat{\boldsymbol{\nu}}_{\tau-\tilde{t}}(\mathbf{x}_{\tilde{t}})$$

($\tilde{t} = \tau - t$:Reversed time)

Examples: Forward diffusion process for accurate data generation

Linear force $F_t(\mathbf{x}) = A_t \mathbf{x} + \mathbf{b}_t$

Gaussian transition probability

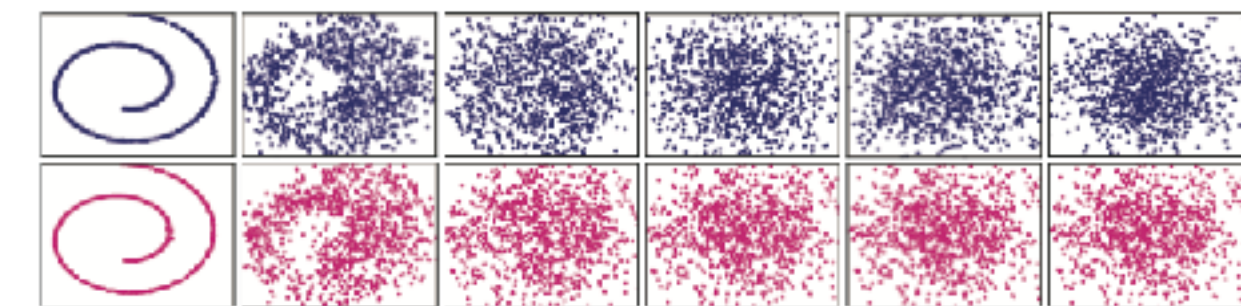
$$P_t^c(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | \boldsymbol{\mu}_t(\mathbf{y}), \boldsymbol{\Sigma}_t) \quad \boldsymbol{\mu}_0(\mathbf{y}) = \delta(\mathbf{x} - \mathbf{y}) \quad \boldsymbol{\Sigma}_0 = \mathbf{O} \quad \longrightarrow \quad P_t(\mathbf{x}) = \int d\mathbf{y} P_t^c(\mathbf{x} | \mathbf{y}) P_0(\mathbf{y})$$

Cosine schedule A. Q. Nichol, & P. Dhariwal, In *International conference on machine learning* (pp. 8162-8171). PMLR (2021)

$$\boldsymbol{\mu}_t(\mathbf{y}) = m_t \mathbf{y} \quad \boldsymbol{\Sigma}_t = \sigma_t^2 \mathbf{I}$$

$$m_t^2 + \sigma_t^2 = 1$$

$$m_t = \cos\left(\frac{\pi t}{2\tau}\right) \quad \sigma_t = \sin\left(\frac{\pi t}{2\tau}\right) \quad t \in [0, \tau]$$

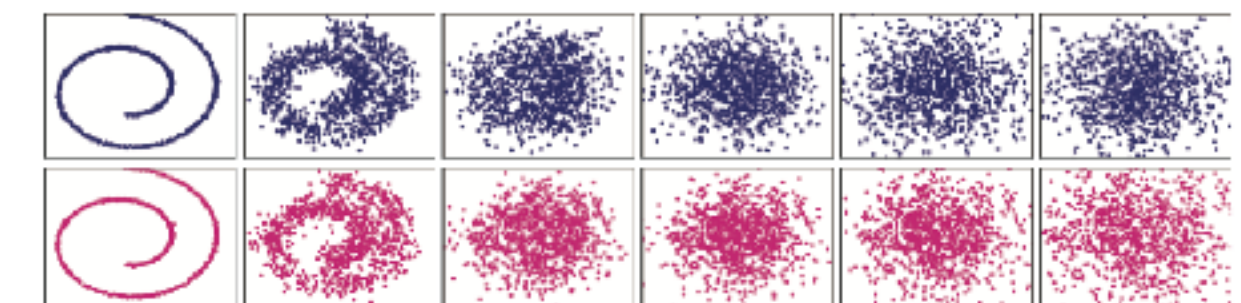


Conditional optimal transport schedule (Approximate optimal transport)

Y. Lipman, R. T. Chen, H. Ben-Hamu, M. Nickel, and M. Le, *International Conference on Learning Representations* (2022)

$$\boldsymbol{\mu}_t(\mathbf{y}) = m_t \mathbf{y} \quad \boldsymbol{\Sigma}_t = \sigma_t^2 \mathbf{I}$$

$$m_t = 1 - \frac{t}{\tau} \quad \sigma_t = \frac{t}{\tau} \quad t \in [0, \tau]$$



Motivation

The diffusion models are inspired by nonequilibrium thermodynamics.

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, PMLR, pp. 2256–2265 (2015).

Question:

Is stochastic thermodynamics still useful for understanding the current technique (e.g., optimal transport) in the diffusion models?

Our results:

In terms of stochastic thermodynamics based on optimal transport, the accuracy of data generation in the diffusion models can be discussed thermodynamically.

Outline

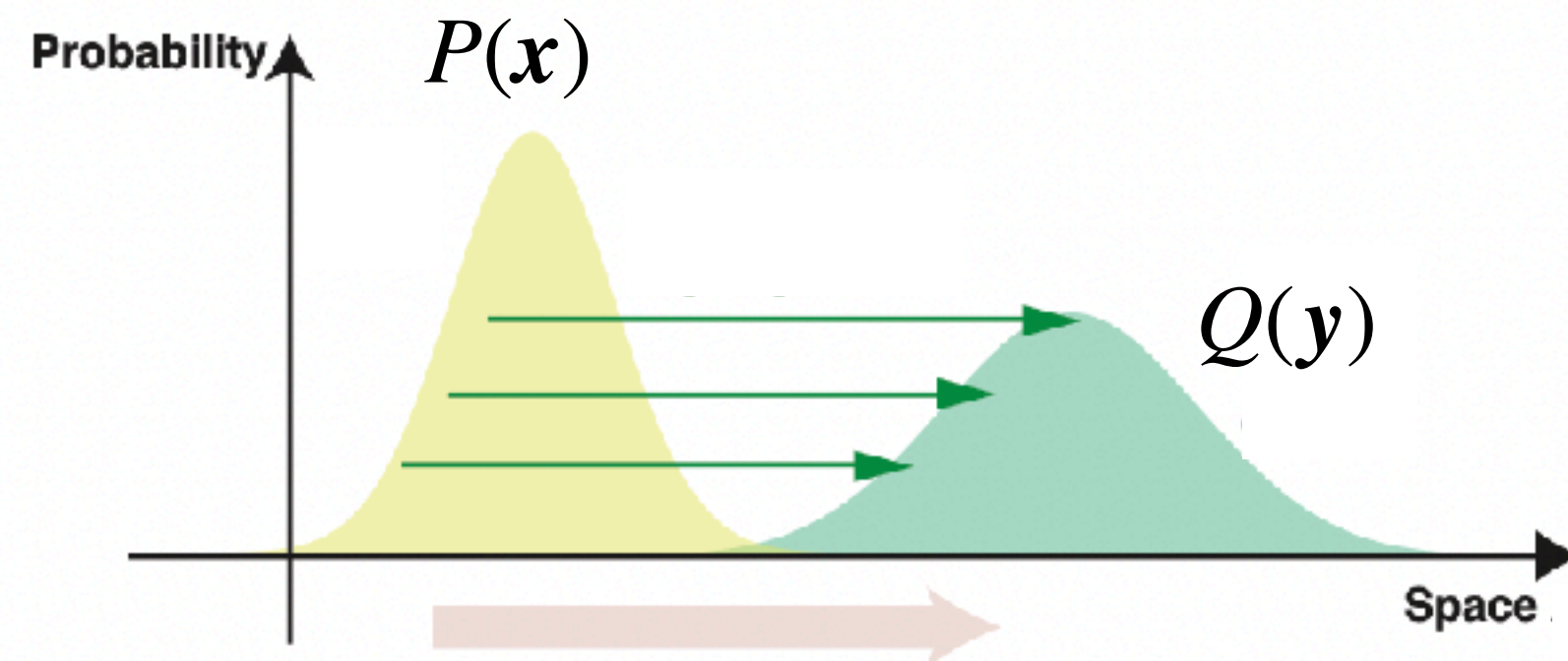
- Introduction: Generative models and diffusion models
- Stochastic thermodynamics based on optimal transport
- Main results: Speed-accuracy trade-off for the diffusion models

Optimal transport: p-Wasserstein distance

Textbook: Villani, C. (2009). *Optimal transport: old and new* (Vol. 338, p. 23). Berlin: springer.

p-Wasserstein distance

$$\mathcal{W}_p(P, Q) = \left(\inf_{\pi \in \Pi(P, Q)} \int dx \int dy \pi(x, y) \|x - y\|^p \right)^{\frac{1}{p}}$$



$$\Pi(P, Q) = \left\{ \pi(x, y) \mid \int dy \pi(x, y) = P(x), \int dx \pi(x, y) = Q(y), \pi(x, y) \geq 0 \right\}$$

Metric: ① $\mathcal{W}_p(P, Q) \geq 0$ ② $\mathcal{W}_p(P, Q) = 0 \Leftrightarrow P = Q$ ③ $\mathcal{W}_p(P, Q) = \mathcal{W}_p(Q, P)$

④ $\mathcal{W}_p(P, R) + \mathcal{W}_p(R, Q) \geq \mathcal{W}_p(P, Q)$

Inequality: $p \geq q \geq 1 \Rightarrow \mathcal{W}_p(P, Q) \geq \mathcal{W}_q(P, Q)$

Optimal transport:

Expressions based on dual problems

Textbook: Villani, C. (2009). *Optimal transport: old and new* (Vol. 338, p. 23). Berlin: springer.

1-Wasserstein distance

(Kantorovich-Rubinstein duality)

$$\mathcal{W}_1(P, Q) = \sup_{f \in \text{Lip}^1} [\langle f \rangle_P - \langle f \rangle_Q]$$

$$\langle f \rangle_P = \int dx f(\mathbf{x}) P(\mathbf{x}) \quad \text{Lip}^1 = \{f(\mathbf{x}) \mid \|\nabla f(\mathbf{x})\|^2 \leq 1\}$$

2-Wasserstein distance

(Benamou-Brenier formula)

$$\mathcal{W}_2(P, Q) = \sqrt{\inf_{\{u_t, Q_t\}_{0 \leq t \leq \tau}} \tau \int_0^\tau dt \int dx \|u_t(\mathbf{x})\|^2 Q_t(\mathbf{x})}$$

$$\partial_t Q_t(\mathbf{x}) = -\nabla \cdot (u_t(\mathbf{x}) Q_t(\mathbf{x})) \quad Q_0(\mathbf{x}) = P(\mathbf{x}) \quad Q_\tau(\mathbf{x}) = Q(\mathbf{x})$$

Stochastic thermodynamics for the diffusion systems

Review: U. Seifert, Reports on progress in physics, 75, 126001 (2012).

Fokker-Planck equation

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) P_t(\mathbf{x}))$$

$$\boldsymbol{\nu}_t(\mathbf{x}) = F_t(\mathbf{x}) - T_t \nabla \ln P_t(\mathbf{x})$$

The entropy production rate

$$\dot{S}_t^{\text{tot}} = \frac{1}{T_t} \int d\mathbf{x} \|\boldsymbol{\nu}_t(\mathbf{x})\|^2 P_t(\mathbf{x})$$

The entropy production

$$S_\tau^{\text{tot}} = \int_0^\tau dt \dot{S}_t^{\text{tot}}$$

Lower bound on the entropy production rate

Lower bound on the entropy production rate
(Excess entropy production rate*)

$$\dot{S}_t^{\text{tot}} \geq \frac{[v_2(t)]^2}{T_t} = \frac{1}{T_t} \int d\mathbf{x} \|\boldsymbol{\nu}_t^{\text{ex}}(\mathbf{x})\|^2 P_t(\mathbf{x})$$

Speed in the space of the 2-Wasserstein distance

$$v_2(t) = \lim_{\Delta t \rightarrow +0} \frac{\mathcal{W}_2(P_t, P_{t+\Delta t})}{\Delta t} = \sqrt{\int d\mathbf{x} \|\boldsymbol{\nu}_t^{\text{ex}}(\mathbf{x})\|^2 P_t(\mathbf{x})}$$

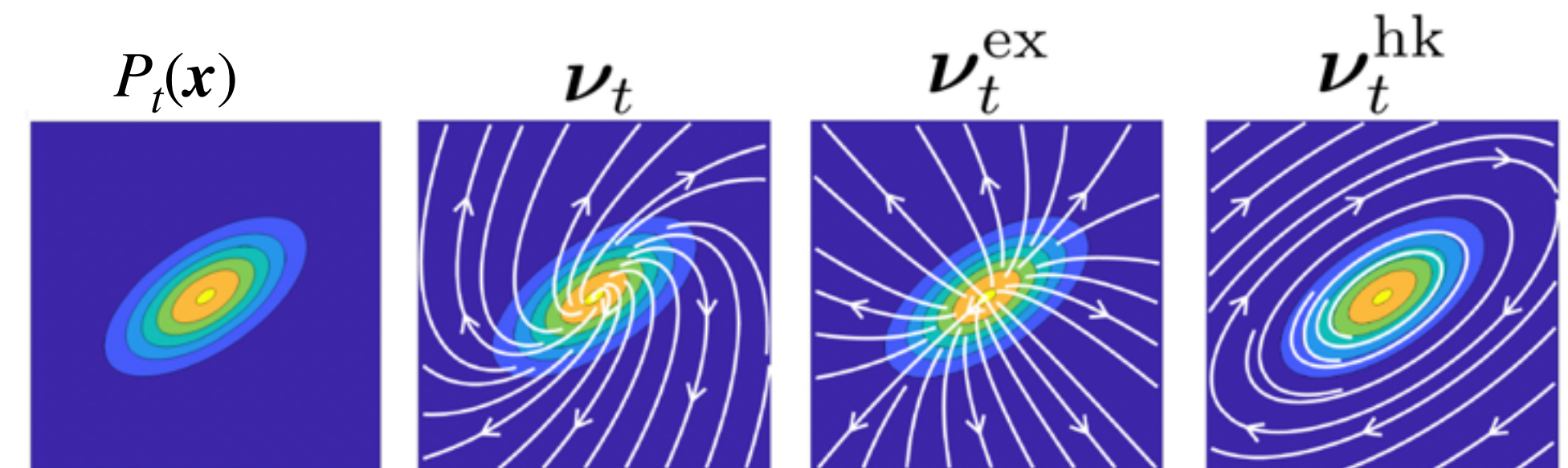
cf.) Benamou-Brenier formula

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) P_t(\mathbf{x})) = -\nabla \cdot (\boldsymbol{\nu}_t^{\text{ex}}(\mathbf{x}) P_t(\mathbf{x}))$$

$$\boldsymbol{\nu}_t^{\text{ex}}(\mathbf{x}) = \nabla \phi_t(\mathbf{x}) \quad \text{:conservative (gradient flow)}$$

$$\nabla \cdot (\boldsymbol{\nu}_t^{\text{hk}}(\mathbf{x}) P_t(\mathbf{x})) = 0 \quad \text{:non-conservative}$$

$$\boldsymbol{\nu}_t^{\text{hk}}(\mathbf{x}) = \boldsymbol{\nu}_t(\mathbf{x}) - \boldsymbol{\nu}_t^{\text{ex}}(\mathbf{x}) \quad \text{(cyclic)}$$



(Figure from) D. Sekizawa, SI and M. Oizumi, arXiv:2312.03489.

M. Nakazato and SI. Phys. Rev. Res. 3, 043093 (2021).

A. Dechant, S-I Sasa and SI. Phys. Rev. Res. 4, L012034 (2022).

* Maes, C., & Netočný, K. *Journal of Statistical Physics*, 154, 188-203 (2014).

Minimum entropy production and geodesic (optimal transport)

Time-independent temperature $T_t = T = \text{const.}$

Thermodynamic speed limit

$$S_\tau^{\text{tot}} \geq \frac{[\int_0^\tau dt v_2(t)]^2}{\tau T} \geq \frac{[\mathcal{W}_2(P_0, P_\tau)]^2}{\tau T}$$

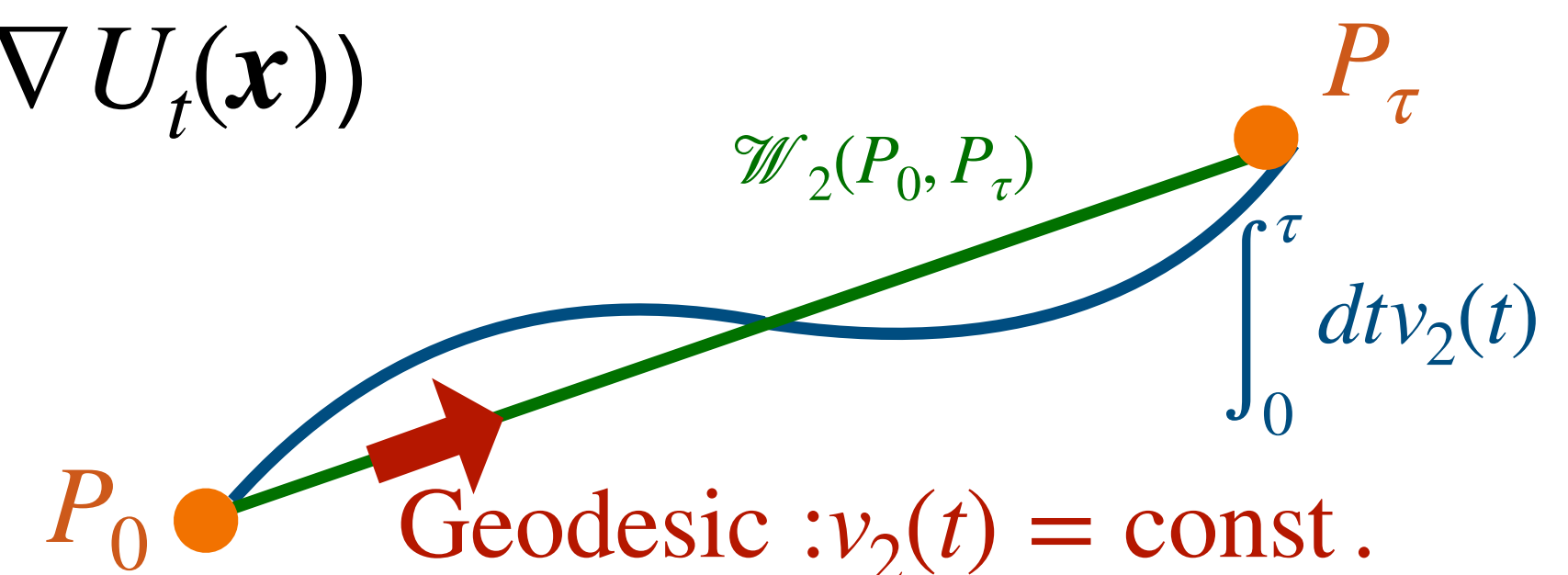
E. Aurell, K. Gawędzki, C. Mejía-Monasterio, R. Mohayaei, & P. Muratore-Ginanneschi, *Journal of statistical physics*, 147, 487-505 (2012).
M. Nakazato and SI. Phys. Rev. Res. 3, 043093 (2021).

Minimum entropy production: Geodesic + Conservative

$$S_\tau^{\text{tot}} = \frac{[\mathcal{W}_2(P_0, P_\tau)]^2}{\tau T}$$

$$\dot{S}_t^{\text{tot}} = \frac{[v_2(t)]^2}{T} \quad \text{:Conservative } (\nu_t(\mathbf{x}) = \nabla \phi_t(\mathbf{x}) \text{ or } F_t(\mathbf{x}) = -\nabla U_t(\mathbf{x}))$$

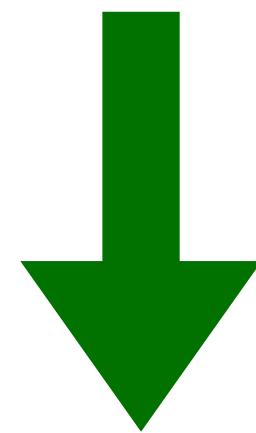
$$v_2(t) = \frac{\mathcal{W}_2(P_0, P_\tau)}{\tau} = \text{const.} \quad \text{:Geodesic (optimal transport)}$$



Thermodynamic uncertainty relation for the excess entropy production rate

Thermodynamic uncertainty relation

$$\dot{S}_t^{\text{tot}} \geq \frac{[v_2(t)]^2}{T_t} \geq \frac{|\partial_t \langle r \rangle_{P_t}|^2}{T_t \langle \|\nabla r\|^2 \rangle_{P_t}}$$



$$v_2(t) \geq v_r(t)$$

(Normalized) speed of observable $r(\mathbf{x})$

$$v_r(t) = \frac{|\partial_t \langle r \rangle_{P_t}|}{\sqrt{\langle \|\nabla r\|^2 \rangle_{P_t}}}$$

Speed in the space of the 2-Wasserstein distance is the upper bound on the speed of any observable.

cf.) $\mathcal{W}_2(P, Q) \geq \mathcal{W}_1(P, Q)$, $r(\mathbf{x}) \in \text{Lip}^1$

R. Nagayama, K. Yoshimura, A. Kolchinsky and SI. arXiv: 2311.16569.

A. Dechant, S-I Sasa and SI. Phys. Rev. Res. 4, L012034 (2022).

A. Dechant, S-I Sasa and SI, Phys. Rev. E. 106, 024125 (2022).

cf.) Cramér–Rao bound: SI and A. Dechant, *Physical Review X*, 10, 021056 (2020).

$r(\mathbf{x})$: time-independent observable

Analogous to thermodynamic speed limit and thermodynamic uncertainty relation

Question:

Is stochastic thermodynamics still useful for understanding the current technique (e.g., optimal transport) in the diffusion models?

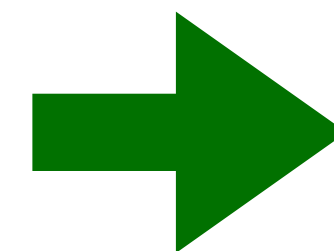
Stochastic thermodynamics

Optimal transport
= Minimum entropy production

Trade-offs

$$\dot{S}_t^{\text{tot}} \geq \frac{[v_2(t)]^2}{T_t} \quad S_\tau^{\text{tot}} \geq \frac{[\int_0^\tau dt v_2(t)]^2}{\tau T} \quad v_2(t) \geq v_r(t)$$

Analogy



Diffusion models

(Approximate) optimal transport
= Accurate data generation
(empirical finding)

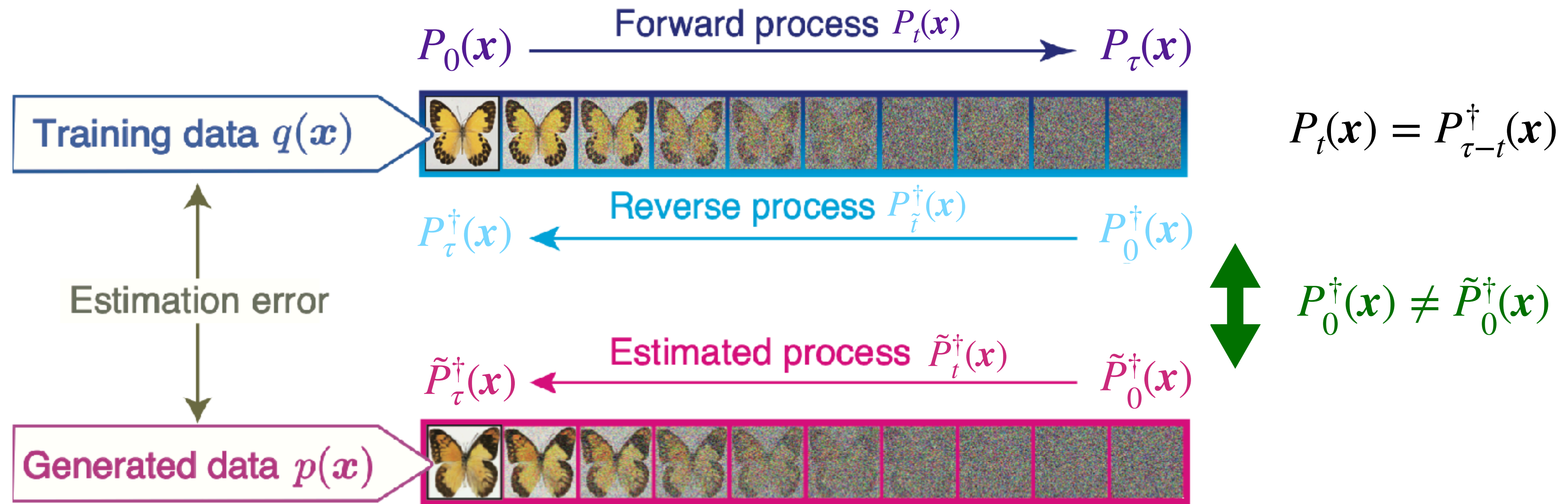
Trade-offs: Our results

Outline

- Introduction: Generative models and diffusion models
- Stochastic thermodynamics based on optimal transport
- Main results: Speed-accuracy trade-off for the diffusion models

K. Ikeda, T. Uda, D. Okanojara and SI, arXiv:2407.04495.

Estimation error in the diffusion models



Estimation error (measured by the 1-Wasserstein distance)

$$\mathcal{W}_1(p, q)$$

e.g.,) K. Oko, S. Akiyama & T. Suzuki, In *International Conference on Machine Learning* (pp. 26517-26582). PMLR (2023).

Perturbation and response

Forward process/Reverse process

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) P_t(\mathbf{x})) \quad P_t(\mathbf{x}) = P_{\tau-t}^\dagger(\mathbf{x})$$

$$\partial_{\tilde{t}} P_{\tilde{t}}^\dagger(\mathbf{x}) = \nabla \cdot (\boldsymbol{\nu}_{\tau-\tilde{t}}(\mathbf{x}) P_{\tilde{t}}^\dagger(\mathbf{x})) \quad \tilde{t} = \tau - t$$

Estimated process

[Probability flow ODE/Flow-based generative modeling]

$$\partial_{\tilde{t}} \tilde{P}_{\tilde{t}}^\dagger(\mathbf{x}) = \nabla \cdot (\boldsymbol{\nu}_{\tau-\tilde{t}}(\mathbf{x}) \tilde{P}_{\tilde{t}}^\dagger(\mathbf{x}))$$

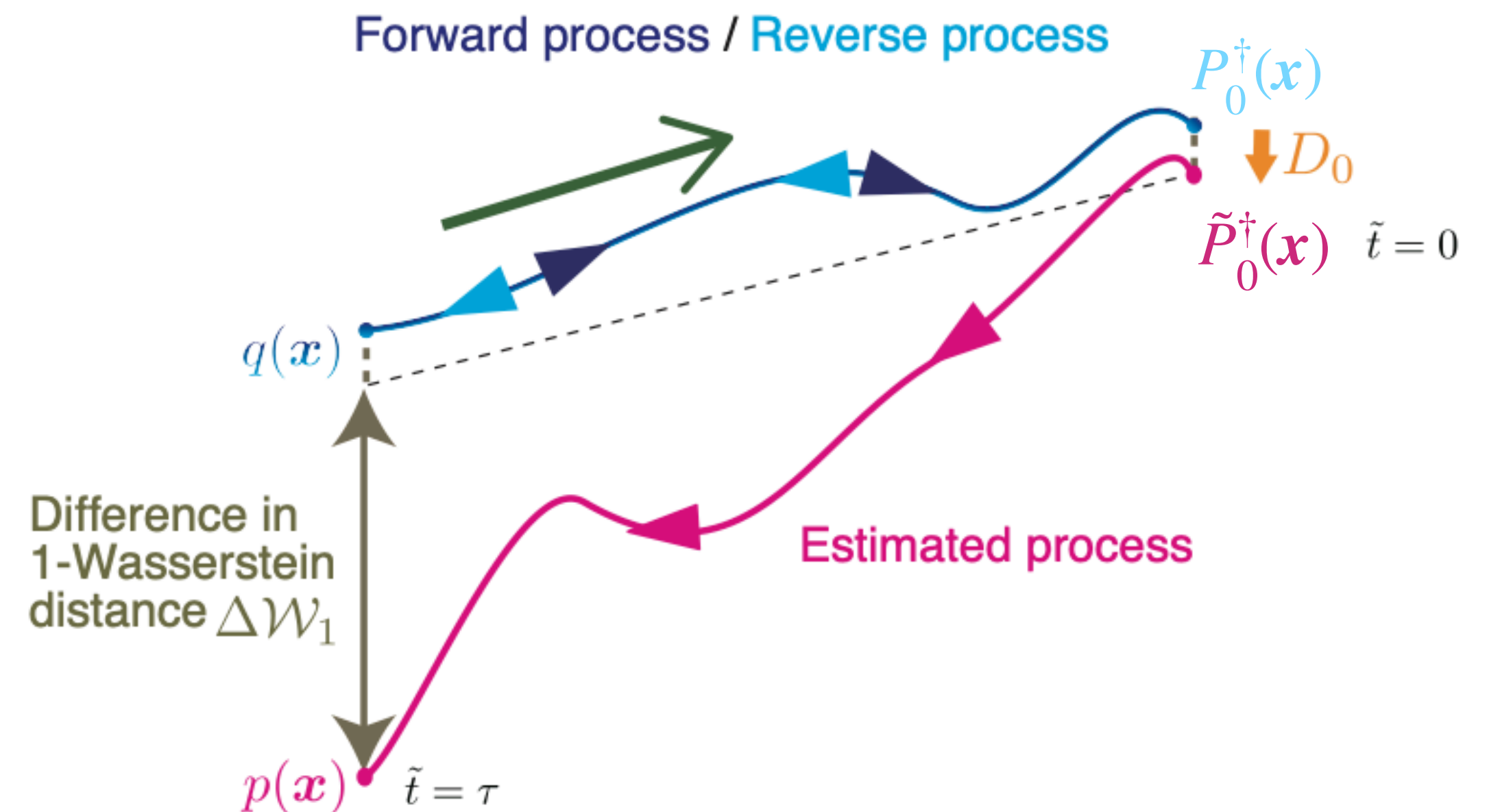
Initial perturbation

$$D_0 = \int dx \frac{(P_0^\dagger(\mathbf{x}) - \tilde{P}_0^\dagger(\mathbf{x}))^2}{P_0^\dagger(\mathbf{x})} : \chi^2\text{-divergence}$$

Response function

$$\frac{\Delta \mathcal{W}_1^2}{D_0} = \frac{[\mathcal{W}_1(p, q) - \mathcal{W}_1(P_0^\dagger, \tilde{P}_0^\dagger)]^2}{D_0} \text{ Perturbation}$$

$\frac{\Delta \mathcal{W}_1^2}{D_0}$ is **small**. \Rightarrow Data generation is **robust** to the initial perturbation.



Main results:

Speed-accuracy trade-off for the diffusion models

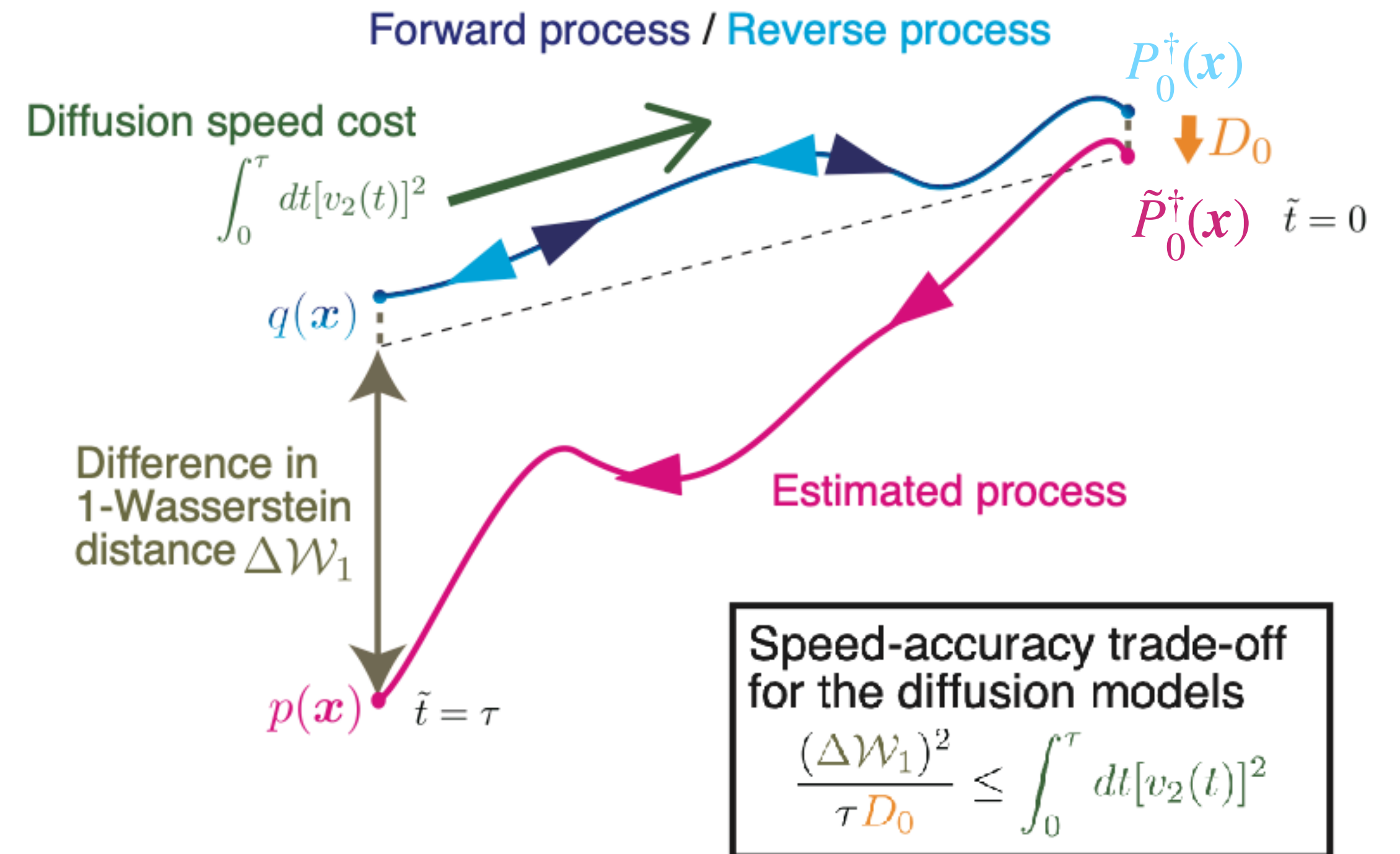
Speed-accuracy trade-off

$$\frac{\Delta \mathcal{W}_1^2}{\tau D_0} \leq \int_0^\tau dt T_t \dot{S}_t^{\text{tot}}$$

Conservative case

$$(\boldsymbol{v}_t(\boldsymbol{x}) = \nabla \phi_t(\boldsymbol{x}) \text{ or } \boldsymbol{F}_t(\boldsymbol{x}) = -\nabla U_t(\boldsymbol{x}))$$

$$\frac{\Delta \mathcal{W}_1^2}{\tau D_0} \leq \int_0^\tau dt [v_2(t)]^2$$



The **robustness** of data generation is generally limited by **the diffusion speed** $v_2(t)$ (or the entropy production rate \dot{S}_t^{tot}) in the forward process.

Main results:

Speed-accuracy trade-off for the diffusion models (Instantaneous)

Instantaneous speed-accuracy trade-off

$$\frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^\dagger, P_{\tau-t}^\dagger)|^2}{D_0} \leq T_t \dot{S}_t^{\text{tot}}$$

Conservative case ($\boldsymbol{\nu}_t(\mathbf{x}) = \nabla \phi_t(\mathbf{x})$ or $\mathbf{F}_t(\mathbf{x}) = -\nabla U_t(\mathbf{x})$)

$$v_{\text{loss}}(t) \leq v_2(t)$$

$$v_{\text{loss}}(t) = \frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^\dagger, P_{\tau-t}^\dagger)|}{\sqrt{D_0}}$$

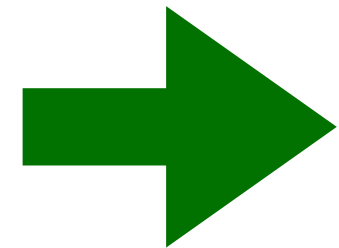
cf.) Thermodynamic uncertainty relation $v_r(t) \leq v_2(t)$

Sketch of proof: Instantaneous trade-off

$$\partial_{\tilde{t}} P_{\tilde{t}}^{\dagger}(\mathbf{x}) = \nabla \cdot (\boldsymbol{\nu}_{\tau-\tilde{t}}(\mathbf{x}) P_{\tilde{t}}^{\dagger}(\mathbf{x}))$$

$$\partial_{\tilde{t}} \tilde{P}_{\tilde{t}}^{\dagger}(\mathbf{x}) = \nabla \cdot (\boldsymbol{\nu}_{\tau-\tilde{t}}(\mathbf{x}) \tilde{P}_{\tilde{t}}^{\dagger}(\mathbf{x}))$$

$$\tilde{t} = \tau - t$$



$$\partial_t [P_{\tau-t}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tau-t}^{\dagger}(\mathbf{x})] = -\nabla \cdot (\boldsymbol{\nu}_t(\mathbf{x}) [P_{\tau-t}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tau-t}^{\dagger}(\mathbf{x})])$$

Continuity equation

$$f \in \text{Lip}^1 \quad |\partial_t (\langle f \rangle_{P_{\tau-t}^{\dagger}} - \langle f \rangle_{\tilde{P}_{\tau-t}^{\dagger}})|^2 = \left(\int d\mathbf{x} f(\mathbf{x}) \partial_t [P_{\tau-t}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tau-t}^{\dagger}(\mathbf{x})] \right)^2$$

$$= \left(\int d\mathbf{x} \nabla f(\mathbf{x}) \cdot \boldsymbol{\nu}_t(\mathbf{x}) [P_{\tau-t}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tau-t}^{\dagger}(\mathbf{x})] \right)^2$$

Cauchy-Schwartz inequality
+ 1-Lipshitz ($\|\nabla f(\mathbf{x})\| \leq 1$)

$$\leq \left(\int d\mathbf{x} \|\boldsymbol{\nu}_t(\mathbf{x})\|^2 P_t(\mathbf{x}) \right) \left(\int d\mathbf{x} \frac{[P_{\tau-t}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tau-t}^{\dagger}(\mathbf{x})]^2}{P_{\tau-t}^{\dagger}(\mathbf{x})} \right)$$

$T_t \dot{S}_t^{\text{tot}}$

D_0 (Time-independent)

+ Kantorovich-Rubinstein duality

$$\exists f \in \text{Lip}^1 \quad |\partial_t \mathcal{W}_1(P_{\tau-t}^{\dagger}, \tilde{P}_{\tau-t}^{\dagger})|^2 \leq |\partial_t (\langle f \rangle_{P_{\tau-t}^{\dagger}} - \langle f \rangle_{\tilde{P}_{\tau-t}^{\dagger}})|^2$$

Instantaneous speed-accuracy trade-off

$$\frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^{\dagger}, P_{\tau-t}^{\dagger})|^2}{D_0} \leq T_t \dot{S}_t^{\text{tot}}$$

Sketch of proof: speed-accuracy trade-off

Instantaneous speed-accuracy trade-off

$$\frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^\dagger, P_{\tau-t}^\dagger)|^2}{D_0} \leq T_t \dot{S}_t^{\text{tot}}$$

$$\int_0^\tau dt T_t \dot{S}_t^{\text{tot}} \geq \int_0^\tau dt \frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^\dagger, P_{\tau-t}^\dagger)|^2}{D_0}$$

Cauchy-Schwartz inequality $\geq \frac{(\Delta \mathcal{W}_1)^2}{\tau D_0}$

Speed-accuracy trade-off

$$\frac{\Delta \mathcal{W}_1^2}{\tau D_0} \leq \int_0^\tau dt T_t \dot{S}_t^{\text{tot}}$$

“Optimal” forward process for accurate data generation

$$\frac{\Delta \mathcal{W}_1^2}{\tau D_0} \leq \int_0^\tau dt [v_2(t)]^2$$

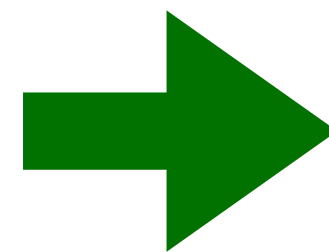
Minimizing the upper bound

$$\int_0^\tau dt [v_2(t)]^2 \geq \frac{\mathcal{W}_2(P_0, P_\tau)^2}{\tau}$$

cf.) Minimum entropy production

$$v_2(t) = \frac{\mathcal{W}_2(P_0, P_\tau)}{\tau} = \text{const.}$$

:Geodesic (optimal transport)



$$\int_0^\tau dt [v_2(t)]^2 = \frac{\mathcal{W}_2(P_0, P_\tau)^2}{\tau}$$

Minimum value

The "optimal" forward process is a dynamics driven by optimal transport (i.e., geodesic in the space of the 2-Wasserstein distance).

“Suboptimal” forward process for accurate data generation

Theorem

N. Shaul, R. T. Chen, M. Nickel, M. Le, and Y. Lipman, in International Conference on Machine Learning, PMLR, pp. 30883–30907 (2023)

If the number of data N_D is small enough compared to the dimension of the data n_d
($N_D/\sqrt{n_d} \rightarrow 0$),

$$\int_0^\tau dt [v_2(t)]^2 \simeq n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2]$$

$$P_t(\mathbf{x}) = \int d\mathbf{y} P_t^c(\mathbf{x} | \mathbf{y}) P_0(\mathbf{y})$$

$$P_t^c(\mathbf{x} | \mathbf{y}) = \mathcal{N}(\mathbf{x} | m_t \mathbf{y}, \sigma_t^2 \mathbf{I})$$

$$\frac{\Delta \mathcal{W}_1^2}{\tau D_0} \leq \int_0^\tau dt [v_2(t)]^2 \simeq n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2]$$

Minimizing the approximate upper bound
(suboptimal)

“Suboptimal” forward process: Conditional optimal transport schedule

$$n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] \geq n_d \frac{(\sigma_0 - \sigma_\tau)^2 + (m_0 - m_\tau)^2}{\tau}$$

$$m_t = 1 - \frac{t}{\tau} \quad \sigma_t = \frac{t}{\tau} \quad \longrightarrow \quad n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] = n_d \frac{(\sigma_0 - \sigma_\tau)^2 + (m_0 - m_\tau)^2}{\tau}$$

:Conditional optimal transport schedule

Minimum value

The “suboptimal” forward process is a dynamics driven by the conditional optimal transport schedule.

“Suboptimal” forward process: Cosine schedule

Constraint: $m_t^2 + \sigma_t^2 = 1 \Rightarrow (m_t, \sigma_t) = (\cos \theta_t, \sin \theta_t)$

$$n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] \geq n_d \frac{(\theta_0 - \theta_\tau)^2}{\tau}$$

$$m_t = \cos\left(\frac{\pi t}{2\tau}\right) \quad \sigma_t = \sin\left(\frac{\pi t}{2\tau}\right) \quad \longrightarrow \quad n_d \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] = n_d \frac{(\theta_0 - \theta_\tau)^2}{\tau}$$

:Cosine schedule

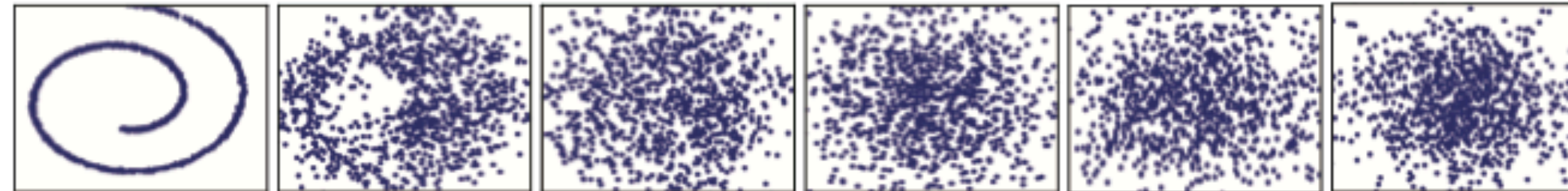
Minimum value

The “suboptimal” forward process under the constraint is a dynamics driven by the cosine schedule.

Examples of optimal and suboptimal dynamics for the diffusion models: Swiss roll

Cosine schedule

Forward
process

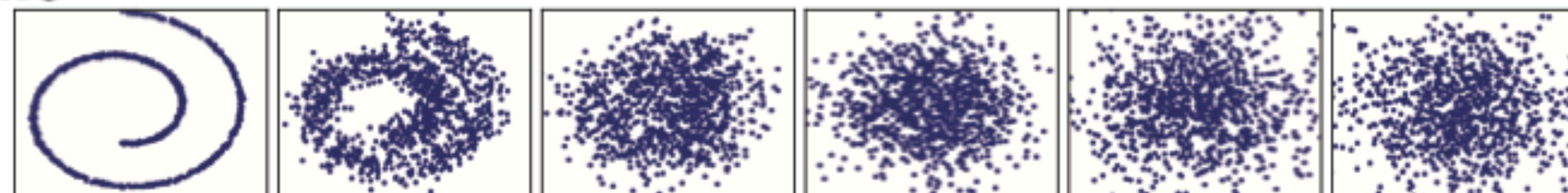


Estimated
process



Cond-OT schedule

Forward
process

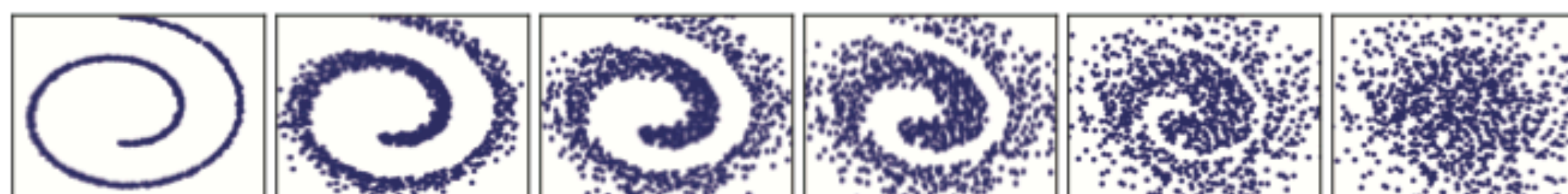


Estimated
process

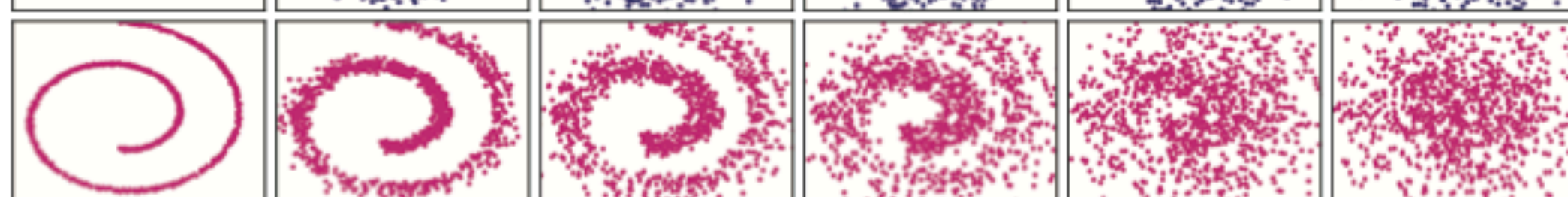


Optimal transport

Forward
process

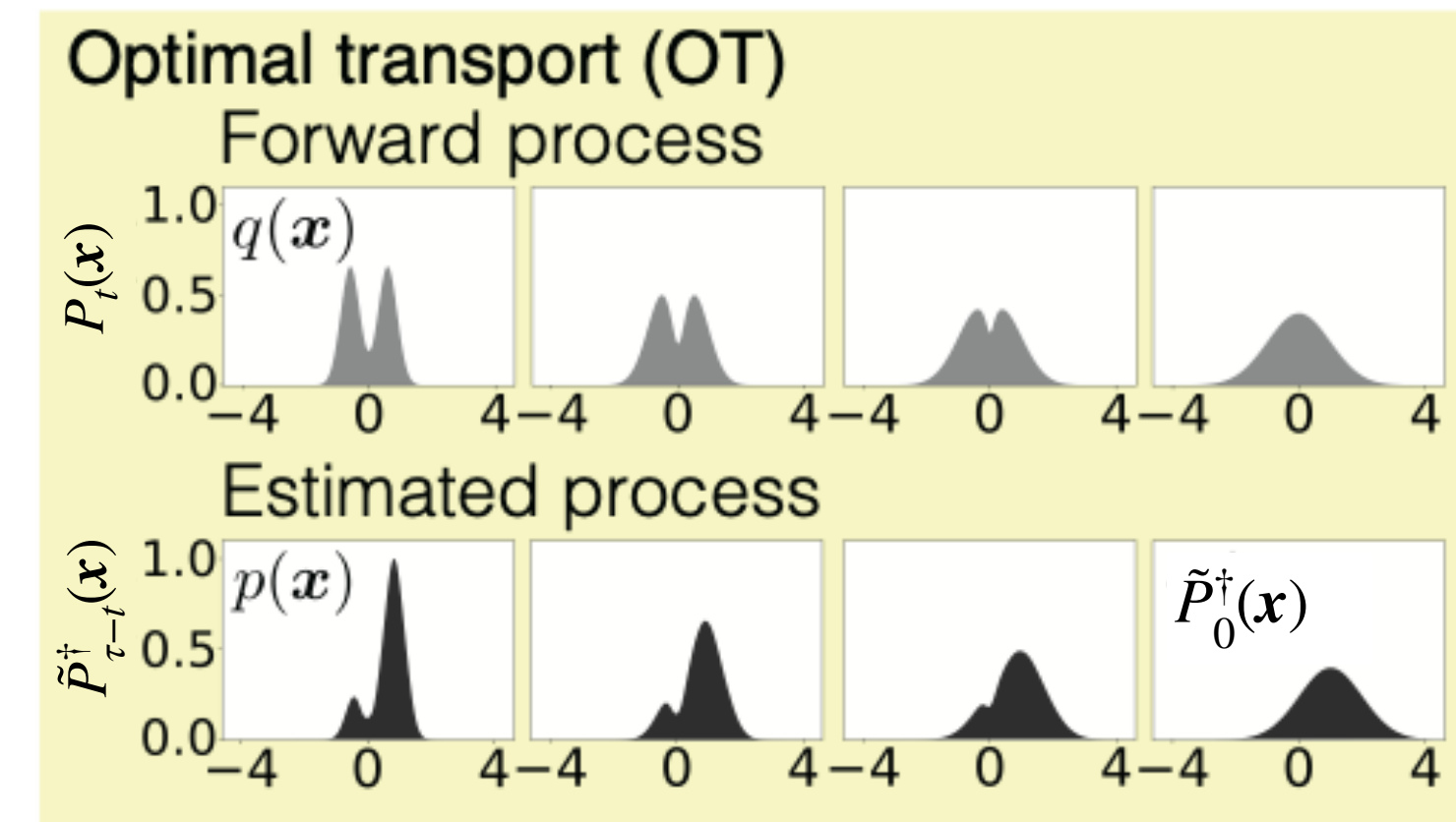
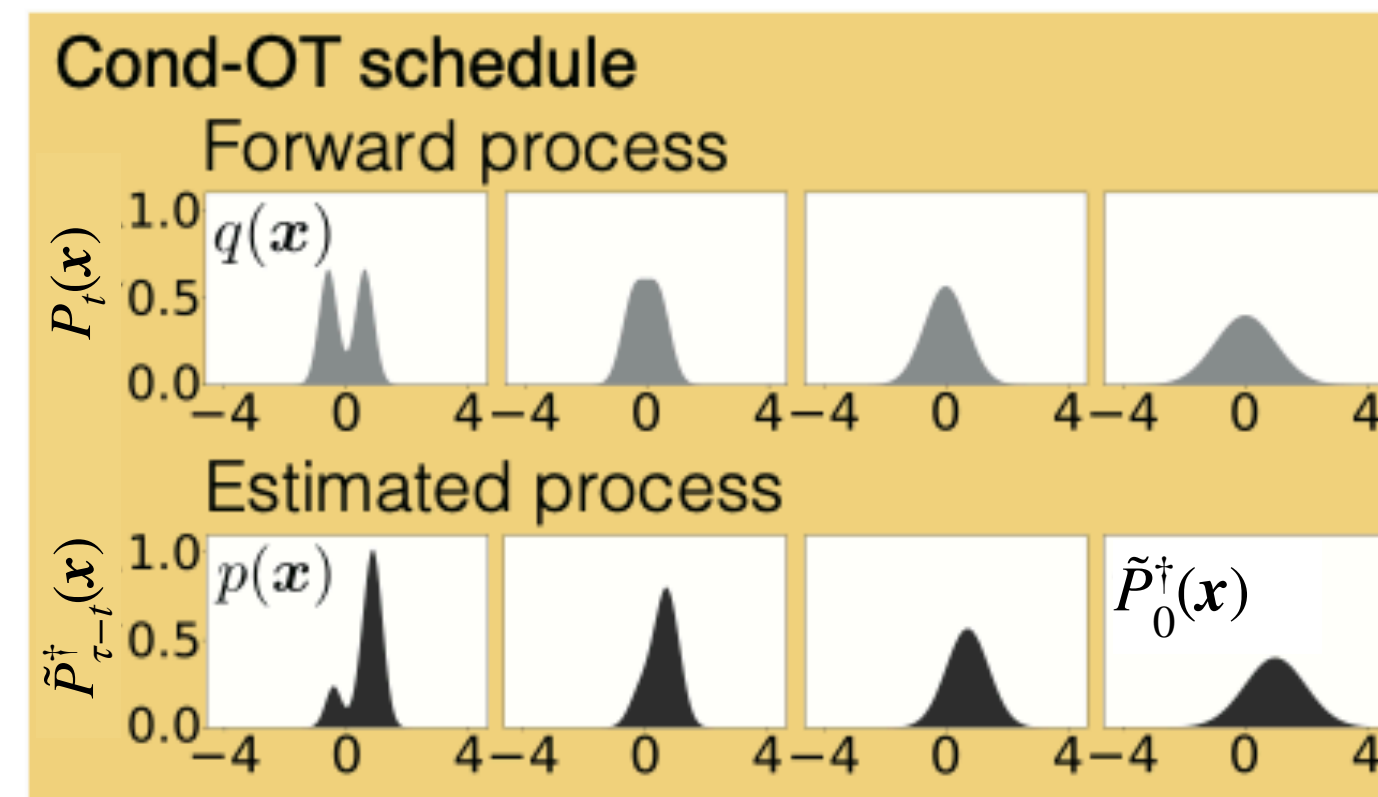
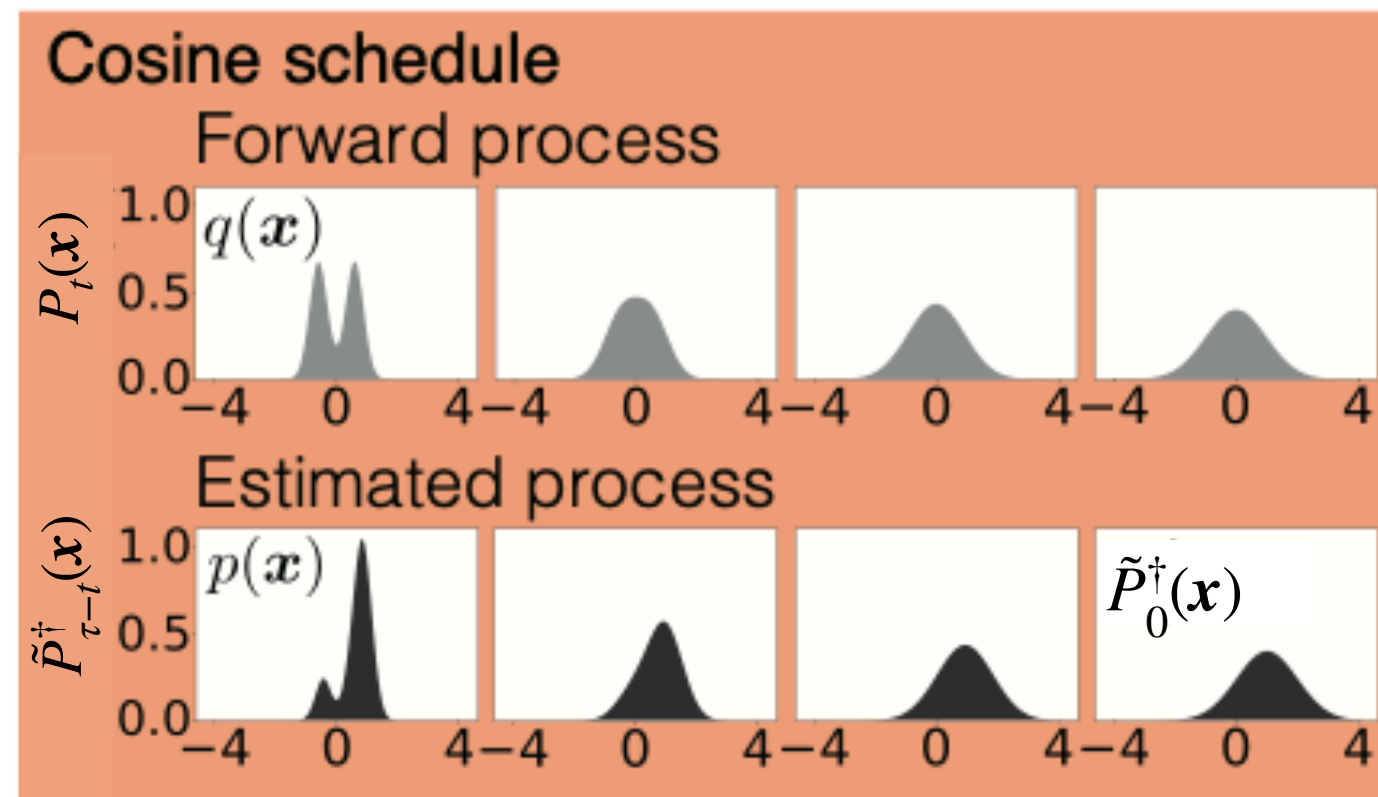


Estimated
process



← Most robust

Examples of optimal and suboptimal dynamics for the diffusion models: Gaussian mixture



0 1/3 2/3 1 t/τ

0 1/3 2/3 1 t/τ

0 1/3 2/3 1 t/τ

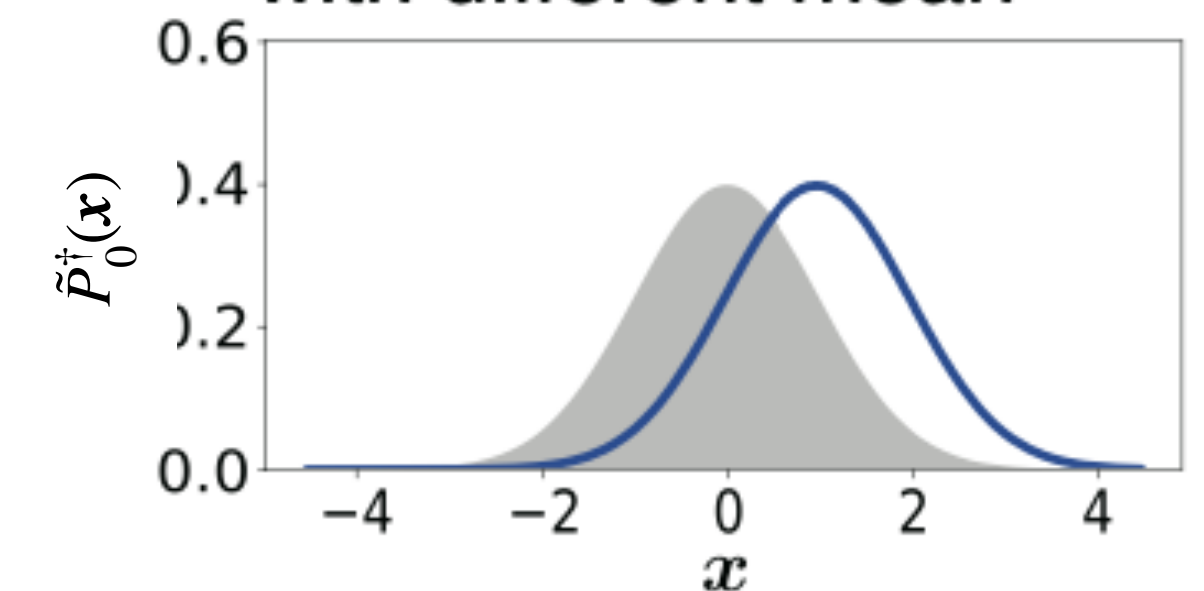
$P_t(\mathbf{x})$: Forward process

$\tilde{P}_{\tau-t}^\dagger(\mathbf{x})$: Estimated process

The data structure (the two peaks) is well recovered even during the dynamics of the estimated process in the case of optimal transport.

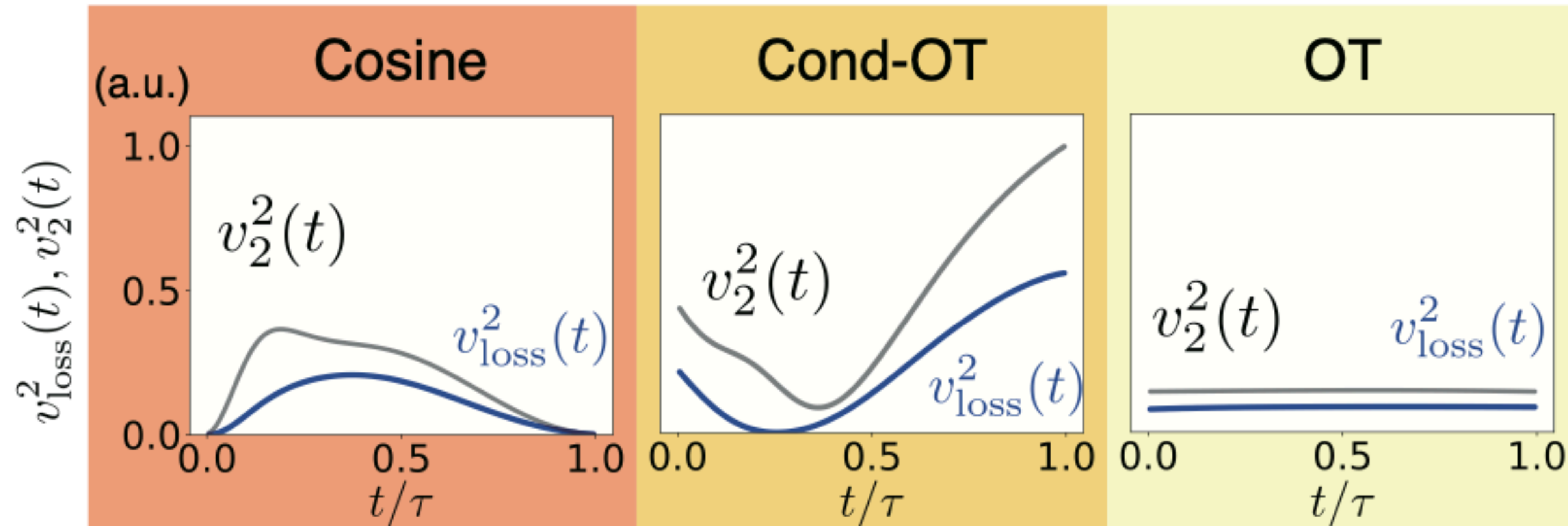
Initial perturbation

Gaussian distribution with different mean



Speed-accuracy trade-off for the diffusion models ^{32/34}

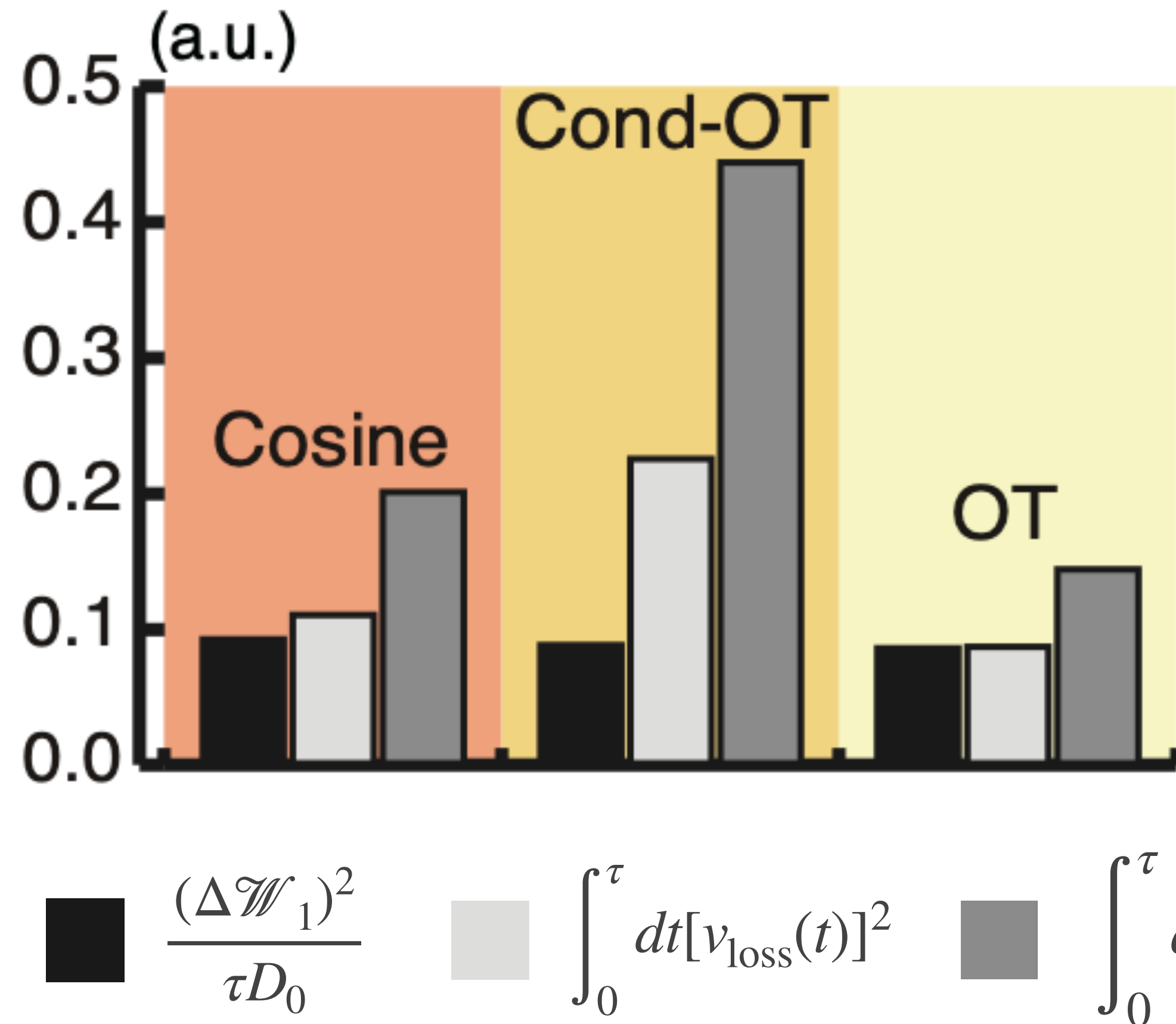
(Instantaneous)



$$\frac{|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^\dagger, P_{\tau-t}^\dagger)|^2}{D_0} = [v_{\text{loss}}(t)]^2 \leq [v_2(t)]^2$$

In the case of optimal transport, the data structure is not well rapidly changed during the dynamics of the estimated process.

Speed-accuracy trade-off for the diffusion models



$$\frac{(\Delta \mathcal{W}_1)^2}{\tau D_0} \leq \int_0^\tau dt [v_{\text{loss}}(t)]^2 \leq \int_0^\tau dt [v_2(t)]^2$$

The bounds are tighter in the case of the optimal transport compared to the cosine and conditional optimal transport schedules.

The value of $(\Delta \mathcal{W}_1)^2 / (\tau D_0)$ for the optimal transport is the smallest for any schedules.

Noise schedules	Values of $(\Delta \mathcal{W}_1)^2 / (\tau D_0)$
Cosine	9.1884×10^{-2}
Cond-OT	8.7810×10^{-2}
OT	8.5375×10^{-2}

Interestingly, the cosine and conditional optimal transport schedules work well in the data generation for this simple case because the response function $(\Delta \mathcal{W}_1)^2 / (\tau D_0)$ is small enough.

Summary

- We used the technique of stochastic thermodynamics and optimal transport to discuss the accurate data generation in the diffusion models.
- We derived the trade-off relationship between the robust data generation to the initial perturbation and the diffusion speed cost given by the 2-Wasserstein distance or the entropy production rate.
- We discuss the optimality and suboptimality of the forward diffusion process in terms of the trade-off, and we found the theoretical validity of the well-used methods (i.e., the cosine and the conditional optimal transport schedule.)

For more information and examples, see K. Ikeda, T. Uda, D. Okanojara and SI, arXiv:2407.04495.

Take-home message:

Stochastic thermodynamics (based on optimal transport) is useful for generative AI.