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



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





UTORYO Universal Biology Institute



# Reference and collaborators

#### Main topic (Diffusion model)

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



Kotaro Ikeda (UTokyo)

#### 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)

Tomoya Uda (UTokyo)



Daisuke Okanohara (Preferred Networks Inc.)



# 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. Okanohara 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

Dream studio by stability.ai https://beta.dreamstudio.ai/generate



# Generative model

### Stable diffusion (2022)

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

Prompt

Frontiers in Non-equilibrium Physics 2024

Dream studio by stability.ai https://beta.dreamstudio.ai/generate







# Diffusion model - Original paper (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.



J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, PMLR, pp. 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.

### Forward diffusion process [learning]

Reverse diffusion process [data generation]











# Essential idea

J. Sohl-Dickstein, E. Weiss, N. Maheswaranathan, and S. Ganguli, PMLR, pp. 2256–2265 (2015). Training data q  $P_0(\mathbf{x}_0) = q(\mathbf{x}_0)$  $P_{\tau}(\boldsymbol{x}_N)$ Forward diffusion process [learning] Estimating the reverse process  $\hat{T}_{i}^{\dagger} = T_{i}^{\dagger}$ Generated data *p* Reverse diffusion process [data generation]  $p(\mathbf{x}_0)(\simeq q(\mathbf{x}_0))$  $P_0^{\dagger}(\mathbf{x}_N)(\simeq P_{\tau}(\mathbf{x}_N))$ 

$$\mathscr{P}^{\mathrm{F}}(\{\boldsymbol{x}\}) = q(\boldsymbol{x}_0) \prod_{i} T_i(\boldsymbol{x}_{i+1} \mid \mathbf{x}_i)$$

$$\mathscr{P}^{\mathrm{E}}(\{\boldsymbol{x}\}) = P_0^{\dagger}(\boldsymbol{x}_{\tau}) \prod_i T_i^{\dagger}(\boldsymbol{x}_i | \boldsymbol{x}_i)$$



## 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)

 $\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\nu_t(\mathbf{x}) P_t(\mathbf{x})) \qquad \nu_t(\mathbf{x}) = F_t(\mathbf{x}) - T_t \nabla \ln P_t(\mathbf{x})$ 

-Data generation by the reverse stochastic differential equation

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

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

# Estimating $\hat{\nu}_t(\mathbf{x}) = F_t(\mathbf{x}) - T_t s_t(\mathbf{x})$ via the score function $s_t = \nabla \ln P_t$

- $(\tilde{t} = \tau t : \text{Reversed time})$
- -Data generation by the ordinary differential equation (probability flow ODE)







## 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)

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\nu_t(\mathbf{x}) P_t(\mathbf{x}))$$
 Estimating

-Data generation by the ordinary differential equation

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

ng the velocity field  $\hat{\nu}_t(x) = \nu_t(x)$ 

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







## Examples: Forward diffusion process for accurate data generation

Linear force  $F_t(x) = A_t x + b_t$ 

Gaussian transition probability

 $P_t^{c}(\boldsymbol{x} | \boldsymbol{y}) = \mathcal{N}(\boldsymbol{x} | \boldsymbol{\mu}_t(\boldsymbol{y}), \boldsymbol{\Sigma}_t) \qquad \boldsymbol{\mu}_0(\boldsymbol{y}) = \boldsymbol{y} \qquad \boldsymbol{\Sigma}_0 = \boldsymbol{y}$ 

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

$$\mu_t(\mathbf{y}) = m_t \mathbf{y} \qquad \Sigma_t = \sigma_t^2 \mathbf{I} \qquad m_t = \cos\left(\frac{\pi t}{2\tau}\right) \qquad m_t^2 + \sigma_t^2 = 1$$

$$\boldsymbol{\mu}_t(\mathbf{y}) = m_t \mathbf{y} \qquad \Sigma_t = \sigma_t^2 \mathbf{I} \qquad m_t = 1 - \frac{\iota}{\tau}$$

### 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)

$$\sigma_t = \frac{t}{\tau} \qquad t \in [0, \tau]$$



(Figure from) K. Ikeda, T. Uda, D. Okanohara and SI, arXiv:2407.04495.





# 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

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



## Optimal transport: p-Wasserstein distance

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

$$P, Q) = \left( \inf_{\pi \in \Pi(P,Q)} \int d\mathbf{x} \int d\mathbf{y} \pi(\mathbf{x}, \mathbf{y}) \|\mathbf{x} - \mathbf{y}\|^p \right)^{\frac{1}{p}}$$

$$Q(x) = \left\{ \left. \pi(x, y) \right| \int dy \pi(x, y) = P(x), \int dx \pi(x, y) = Q(y), \pi(x, y) \ge 0 \right\}$$

 $(1) \mathcal{W}_p(P,Q) \ge 0 \quad (2) \mathcal{W}_p(P,Q) = 0 \Leftrightarrow P = Q \quad (3) \mathcal{W}_p(P,Q) = \mathcal{W}_p(Q,P)$ 



## Optimal transport: Expressions based on dual problems

# 1-Wasserstein distance(Kantorovich-Rubinstein duality)

### 2-Wasserstein distance (Benamou-Brenier formula)

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

$$\mathscr{W}_{1}(P,Q) = \sup_{f \in \operatorname{Lip}^{1}} [\langle f \rangle_{P} - \langle f \rangle_{Q}]$$

$$\langle f \rangle_P = \int d\mathbf{x} f(\mathbf{x}) P(\mathbf{x}) \qquad \text{Lip}^1 = \{f(\mathbf{x}) | \| \nabla f(\mathbf{x}) \|^2 \le$$

$$\mathcal{W}_{2}(P,Q) = \sqrt{\inf_{\{u_{t},Q_{t}\}_{0 \le t \le \tau}} \tau \int_{0}^{\tau} dt \int dx \|u_{t}(x)\|^{2} Q_{t}(x)}$$

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

J-D. Benamou & Y. Brenier. Numerische Mathematik 84, 375-393 (2000).



## Stochastic thermodynamics for the diffusion systems

#### Fokker-Planck equation

$$\partial_t P_t(\mathbf{x}) = -\nabla \cdot (\nu_t(\mathbf{x})P_t(\mathbf{x}))$$

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

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

### 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 Speed in the space of the 2-Wasserstein distance (Excess entropy production rate\*)

$$\dot{S}_{t}^{\text{tot}} \geq \frac{[\nu_{2}(t)]^{2}}{T_{t}} = \frac{1}{T_{t}} \int d\mathbf{x} \| \boldsymbol{\nu}_{t}^{\text{ex}}(\mathbf{x}) \|^{2} P_{t}(\mathbf{x})$$

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

 $v_t^{\text{ex}}(x) = \nabla \phi_t(x)$  :conservative (gradient flow)

 $\nabla \cdot (\nu_t^{hk}(x)P_t(x)) = 0 \quad \text{inon-conservative}$  $\nu_t^{hk}(x) = \nu_t(x) - \nu_t^{ex}(x) \quad \text{(cyclic)}$ 

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

cf.) Benamou-Brenier formula

) ent flow)  $P_t(x) = v_t + v_t^{ex} + v_t^{hk}$ 

(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  $T_t = T = const$ .

Thermodynamic speed limit

 $S_{\tau}^{\text{tot}} \geq \frac{\left[\int_{0}^{\tau}\right]}{}$ 

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

Animum entropy production: Geodesic + Conservative  

$$S_{\tau}^{\text{tot}} = \frac{[\mathscr{W}_{2}(P_{0}, P_{\tau})]^{2}}{\tau T}$$

$$\dot{S}_{t}^{\text{tot}} = \frac{[v_{2}(t)]^{2}}{T} \quad \text{:Conservative} \ (\nu_{t}(\boldsymbol{x}) = \nabla \phi_{t}(\boldsymbol{x}) \text{ or } F_{t}(\boldsymbol{x}) = -\nabla U_{t}(\boldsymbol{x}))$$

$$\mathcal{W}_{2}(P_{0}, P_{\tau})$$

$$\mathcal{W}_{2}(t) = \frac{\mathscr{W}_{2}(P_{0}, P_{\tau})}{\tau} = \text{const.} \quad \text{:Geodesic (optimal transport)}$$

$$P_{0} \quad \text{Geodesic :} v_{2}(t) = \text{const.}$$

Minimum entropy production: Geodesic + Conservative  

$$S_{\tau}^{\text{tot}} = \frac{[\mathscr{W}_{2}(P_{0}, P_{\tau})]^{2}}{\tau T}$$

$$\dot{S}_{t}^{\text{tot}} = \frac{[v_{2}(t)]^{2}}{T} \quad \text{:Conservative} \ (\nu_{t}(\boldsymbol{x}) = \nabla \phi_{t}(\boldsymbol{x}) \text{ or } \boldsymbol{F}_{t}(\boldsymbol{x}) = -\nabla U_{t}(\boldsymbol{x}))$$

$$v_{2}(t) = \frac{\mathscr{W}_{2}(P_{0}, P_{\tau})}{\tau} = \text{const.} \quad \text{:Geodesic (optimal transport)}$$

$$P_{0} \quad \begin{array}{c} \mathcal{W}_{2}(P_{0}, P_{\tau}) \\ \mathcal{W}_{2}(t) = \mathcal{W}_{2}(t) = \text{const.} \end{array}$$

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





## Thermodynamic uncertainty relation for the excess entropy production rate

Thermodynamic uncertainty relation

$$\dot{S}_{t}^{\text{tot}} \geq \frac{[v_{2}(t)]^{2}}{T_{t}} \geq \frac{|\partial_{t}\langle r \rangle_{P_{t}}|^{2}}{T_{t}\langle ||\nabla r||^{2}\rangle_{P_{t}}}$$
cf.) Cramér-Rao bound: SI and A. Dechant, *Physical Review X*, 10, 02  
 $r(x)$ : time-independent observable  
 $v_{2}(t) \geq v_{r}(t)$ 
(Normalized) speed of observable  $r(x)$   $v_{r}(t) = \frac{|\partial_{t}\langle r \rangle_{P_{t}}|}{\sqrt{\langle ||\nabla r||^{2}\rangle_{P_{t}}}}$ 

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

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).

21056 (2020).

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

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







# Analogous to thermodynamic speed limit and thermodynamic uncertainty relation

#### **Question:**

Stochastic thermodynamics

**Optimal transport** = Minimum entropy production

# Trade-offs $\dot{S}_{t}^{\text{tot}} \ge \frac{[v_{2}(t)]^{2}}{T_{\star}} \quad S_{\tau}^{\text{tot}} \ge \frac{[\int_{0}^{\tau} dt v_{2}(t)]^{2}}{\tau T} \quad v_{2}(t) \ge v_{r}(t)$

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

**Diffusion models** 

Analogy

(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. Okanohara and SI, arXiv:2407.04495.





### 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).





Forward process/Reverse process

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

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

### Initial perturbation

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

#### **Response function**

 $D_0$ 

#### **Estimation error**

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

*D*<sub>0</sub> Perturbation

 $\frac{\Delta \mathscr{W}_1^2}{D_0}$ 

# Perturbation and response

Estimated process

[Probability flow ODE/Flow-based generative modeling]

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



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



Conservative case  $(\nu_t(\mathbf{x}) = \nabla \phi_t(\mathbf{x}) \text{ or } F_t(\mathbf{x}) = -\nabla U_t(\mathbf{x}))$ 



The robustness of data generation is generally limited by the diffusion speed  $v_2(t)$ 

(or the entropy production rate  $\dot{S}_{t}^{\text{tot}}$ ) in the forward process.





### Main results: Speed-accuracy trade-off for the diffusion models (Instantaneous)

Instantaneous speed-accuracy trade-off

 $|\partial_t \mathcal{W}_1(\tilde{P}^{\dagger}_{\tau-t})|$ 

Conservative case ( $\nu_t(x) = \nabla \phi_t(x)$  or  $F_t$ 

 $v_{\rm los}$ 

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

$$\frac{P_{\tau-t}^{\dagger}}{T_{t}} \Big|^{2} \leq T_{t} \dot{S}_{t}^{\text{tot}}$$

$$f(\mathbf{x}) = -\nabla U_t(\mathbf{x}))$$

$$v_{\text{loss}}(t) \le v_2(t) \qquad \qquad v_{\text{loss}}(t) = \frac{\left|\partial_t \mathcal{W}_1(\tilde{P}_{\tau-t}^{\dagger}, P_{\tau-t}^{\dagger})\right|}{\sqrt{D_0}}$$

# Sketch of proof: Instantaneous trade-off

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

 $f \in \operatorname{Lip}^{1} \qquad |\partial_{t}(\langle f \rangle_{P^{\dagger}_{\tau-t}} - \langle f \rangle_{\tilde{P}^{\dagger}_{\tau-t}})|^{2} = \left( \int_{\mathbb{T}} |\partial_{t}(\langle f \rangle_{P^{\dagger}_{\tau-t}} - \langle f \rangle_{\tilde{P}^{\dagger}_{\tau-t}})|^{2} \right)$ = ( | Cauchy-Schwartz inequality  $\leq$ + 1-Lipshitz ( $\|\nabla f(\mathbf{x})\| \le 1$ )

+ Kantrovich-Rubinstein duality

Instantaneous speed-accuracy trade

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

Continuity equation

$$= \left( \int d\mathbf{x} f(\mathbf{x}) \partial_t [P_{\tilde{t}}^{\dagger}(\mathbf{x}) - \tilde{P}_{\tilde{t}}^{\dagger}(\mathbf{x})] \right)^2$$

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

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

$$= \int d_t \mathcal{W}_1(P_{\tau-t}^{\dagger}, \tilde{P}_{\tau-t}^{\dagger}) |^2 \leq |\partial_t (\langle f \rangle_{P_{\tau-t}^{\dagger}} - \langle f \rangle_{\tilde{P}_{\tau-t}^{\dagger}})|^2$$

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







# Sketch of proof: speed-accuracy trade-off

Instantaneous speed-accuracy trade-off

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

Cauchy-Schwartz inequality  $\geq$   $(\Delta)$ 

### Speed-accuracy trade-off

$$\frac{\left|\partial_{t}\mathcal{W}_{1}(\tilde{P}_{\tau-t}^{\dagger}, P_{\tau-t}^{\dagger})\right|^{2}}{D_{0}} \leq T_{t}\dot{S}_{t}^{\text{tot}}$$

$$\frac{(\Delta \mathcal{W}_1)^2}{\tau D_0}$$

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

## "Optimal" forward process for accurate data generation

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

$$\int_{0}^{\tau} dt [v_{2}(t)]^{2} \ge \frac{\mathcal{W}_{2}(P_{0}, P_{\tau})^{2}}{\tau}$$

cf.) Minimum entropy production

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

:Geodesic (optimal transport)

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

Minimizing the upper bound

$$\int_{0}^{\tau} dt [v_2(t)]^2 = \frac{\mathscr{W}_2(P_0, P_{\tau})^2}{\tau} \qquad \text{Minimum value}$$

## "Suboptimal" forward process for accurate data generation

Theorem If the number of data  $N_{\rm D}$  is small enough compared to the dimension of the data  $n_{\rm d}$  $(N_{\rm D}/\sqrt{n_{\rm d}} \rightarrow 0),$  $\int_{0}^{1} dt [v_2(t)]^2 \simeq n_d$ 

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

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

$$\int_0^{\tau} dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2]$$

$$P_{t}(\boldsymbol{x}) = \int d\boldsymbol{y} P_{t}^{c}(\boldsymbol{x} | \boldsymbol{y}) P_{0}(\boldsymbol{y})$$
$$P_{t}^{c}(\boldsymbol{x} | \boldsymbol{y}) = \mathcal{N}(\boldsymbol{x} | \boldsymbol{m}_{t} \boldsymbol{y}, \sigma_{t}^{2})$$

Minimizing the approximate upper bound (suboptimal)





## "Suboptimal" forward process: Conditional optimal transport schedule

$$n_{\mathrm{d}} \int_{0}^{\tau} dt [(\partial_{t} \sigma_{t})^{2} + (\partial_{t} m_{t})^{2}] \geq n_{\mathrm{d}} \frac{(\sigma_{0} - \sigma_{\tau})^{2} + (m_{0} - m_{\tau})^{2}}{\tau}$$

$$m_t = 1 - \frac{t}{\tau}$$
  $\sigma_t = \frac{t}{\tau}$ 

:Conditional optimal transport schedule

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

$$n_{\rm d} \int_0^{\tau} dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] = n_{\rm d} \frac{(\sigma_0 - \sigma_\tau)^2 + (m_0 - m_0)^2}{\tau}$$
  
Minimum valu







## "Suboptimal" forward process: Cosine schedule

Constraint: 
$$m_t^2 + \sigma_t^2 = 1 \implies (m_t, \sigma_t) = (constraint)$$

$$n_{\rm d} \int_0^{\tau} dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] \ge n_{\rm d} \frac{(\theta_0 - \theta_0)^2}{\tau}$$

$$m_t = \cos\left(\frac{\pi}{2}\frac{t}{\tau}\right) \quad \sigma_t = \sin\left(\frac{\pi}{2}\frac{t}{\tau}\right)$$

:Cosine schedule

driven by the cosine schedule.

 $\cos\theta_t, \sin\theta_t$ 

 $(\theta_{\tau})^2$ 

$$n_{\rm d} \int_0^\tau dt [(\partial_t \sigma_t)^2 + (\partial_t m_t)^2] = n_{\rm d} \frac{(\theta_0 - \theta_\tau)^2}{\tau}$$

Minimum value

The "suboptimal" forward process under the constraint is a dynamics



# 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





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



 $P_t(x)$ :Forward process  $\tilde{P}_{-}^{\dagger}(x)$ :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<sup>32/34</sup> (Instantaneous)



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

$$(t)]^2 \le [v_2(t)]^2$$

# Speed-accuracy trade-off for the diffusion models



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.

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

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 $(\Delta \mathcal{W})$
	Cosine	$9.1884 \times 10^{-2}$
$(t)]^2$	Cond-OT	$8.7810 \times 10^{-2}$
	OT	$8.5375 \times 10^{-2}$



# Summary

- discuss the accurate data generation in the diffusion models.
- or the entropy production rate.
- methods (i.e., the cosine and the conditional optimal transport schedule.)

Take-home message: Stochastic thermodynamics (based on optimal transport) is useful for generative Al.

• We used the technique of stochastic thermodynamics and optimal transport to

• 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

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







