

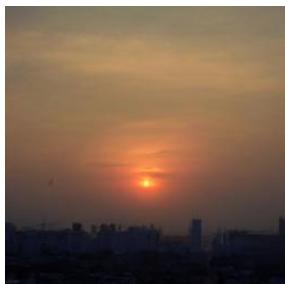
# Reward fine-tuning for flow and diffusion models

Carles Domingo-Enrich,  
Nov. 19 2024

Main reference: *Adjoint Matching: Fine-tuning Flow and Diffusion Generative Models with Memoryless Stochastic Optimal Control*. C. Domingo-Enrich, M. Drozdzal, B. Karrer, R. T. Q. Chen, ICLR 2025. <https://arxiv.org/abs/2409.08861>

# Open questions

- What is the best way to perform RLHF on diffusion and flow models?
- How do we solve stochastic optimal control problems in high dimensions?



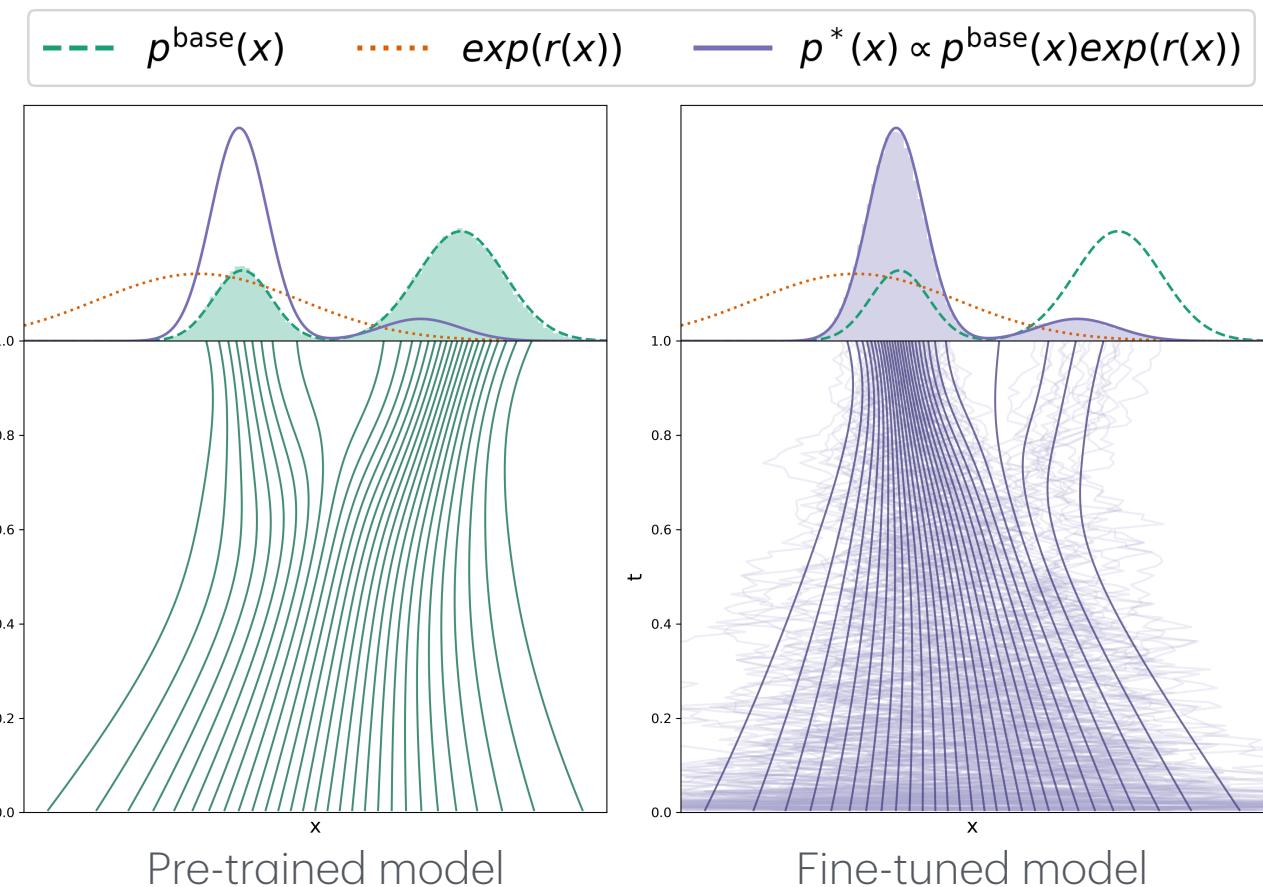
*"Beautiful colorful sunset midst of building in Bangkok Thailand"*



*"Beautiful grandma and granddaughter are mixing salad and smiling while cooking in kitchen"*

# Fine-tuning setup

- Pre-trained diffusion or flow matching model that generates distribution  $p^{base}$
- Reward model  $r(x)$ , trained using human preferences or encoding conditional sampling:  
$$r(x) = \log p(o | x)$$
- Goal: modify pre-trained diffusion model such that it generates  $p^*(x) \propto p(x)\exp(r(x))$



# Flow Matching and diffusion models: notation

- $\bar{X}_0 \sim p_0 = N(0, I)$  and  $\bar{X}_1 \sim p_{\text{data}}$
- Reference flow  $\bar{\mathbf{X}} = (\bar{X}_t)_{t \in [0,1]}$ , where

$$\bar{X}_t = \beta_t \bar{X}_0 + \alpha_t \bar{X}_1.$$

- Generative flow  $\mathbf{X} = (X_t)_{t \in [0,1]}$ , where

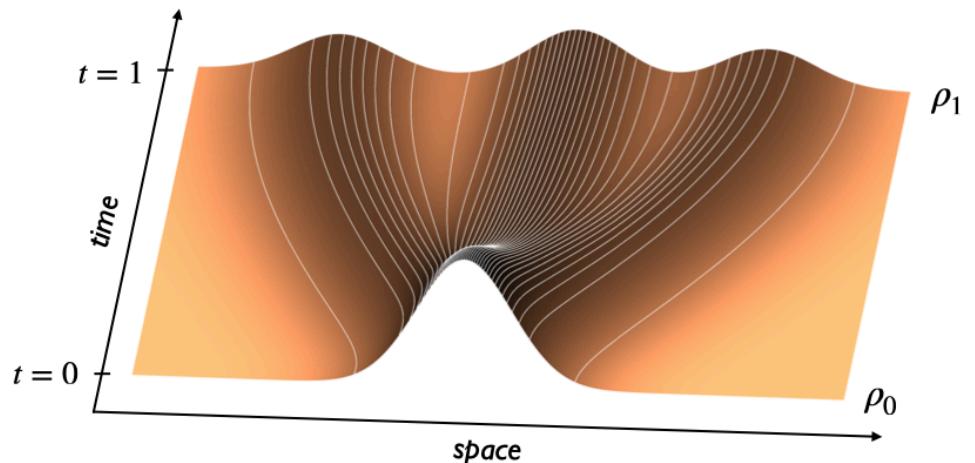
$$\frac{dX_t}{dt} = v(X_t, t), \quad X_0 \sim N(0, I)$$

- $\bar{X}_t$  and  $X_t$  are equally distributed, for all  $t \in [0,1]$

- How do we train the vector field  $v$ ?

By matching the derivative of the reference flow (*Flow Matching*):

$$v(\bar{X}_t, t) = \operatorname{argmin}_{\hat{v}} \mathbb{E} \left\| \hat{v}(\bar{X}_t, t) - \frac{d\bar{X}_t}{dt} \right\|^2$$



Source: Albergo & Vanden-Eijnden, ICLR 2023

# The Flow Matching Generative SDE

- Let  $p_t$  be the distribution of the reference flow  $\bar{X}_t = \beta_t \bar{X}_0 + \alpha_t \bar{X}_1$ . Score function  $\mathfrak{s}(x, t) := \nabla \log p_t(x)$ .
- FM vector field  $v(x, t)$  in terms of  $\mathfrak{s}(x, t)$ :  $v(x, t) = \frac{\dot{\alpha}_t}{\alpha_t} x + \beta_t \left( \frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right) \mathfrak{s}(x, t)$ .
- Flow Matching Generative SDE:  $dX_t = (v(X_t, t) + \sigma(t)^2 \mathfrak{s}(X_t, t)) dt + \sigma(t) dB_t$ , or equivalently,

$$dX_t = \left( v(X_t, t) + \frac{\sigma(t)^2}{2\beta_t \left( \frac{\dot{\alpha}_t}{\alpha_t} - \dot{\beta}_t \right)} \left( v(X_t, t) - \frac{\dot{\alpha}_t}{\alpha_t} X_t \right) \right) dt + \sigma(t) dB_t, \quad X_0 \sim N(0, I)$$

where the noise schedule  $\sigma(t)$  is arbitrary.  $X_t$  still has the same distribution as  $\bar{X}_{t'}$  for all  $t \in [0, 1]!$

- The FM Generative SDE reads:  $dX_t = b(X_t, t) dt + \sigma(t) dB_t$ ,  $X_0 \sim \mathcal{N}(0, I)$ ,  
 where  $b(x, t) = \kappa_t x + \left( \frac{\sigma(t)^2}{2} + \eta_t \right) \mathfrak{s}(x, t)$ ,  $\kappa_t = \frac{\dot{\alpha}_t}{\alpha_t}$ ,  $\eta_t = \beta_t \left( \frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right)$

# Stochastic optimal control

Definition of the problem

$$\begin{aligned} \min_u \quad & \overbrace{\mathbb{E} \left[ \int_0^1 \left( \underbrace{\frac{1}{2} \|u(X_t^u, t)\|^2}_{\text{control cost}} + \underbrace{f(X_t^u, t)}_{\text{state cost}} \right) dt + \underbrace{g(X_1^u)}_{\text{terminal cost}} \right]}^{\text{control objective}}, \\ s.t. \quad & dX_t^u = \underbrace{(b(X_t^u, t) + \sigma(t)u(X_t^u, t))}_{\substack{\text{base} \\ \text{drift}}} dt + \underbrace{\sigma(t)}_{\text{noise schedule}} dB_t, \quad X_0^u \sim p_0. \end{aligned}$$

Value function:

$$V(x, t) := \min_{u \in \mathcal{U}} \mathbb{E}_{X^u} \left[ \int_t^1 \left( \frac{1}{2} \|u(X_s^u, s)\|^2 + f(X_s^u, s) \right) ds + g(X_1^u) \mid X_t^u = x \right] = -\log \mathbb{E}_X \left[ \exp \left( - \int_t^1 f(X_s, s) ds - g(X_1) \right) \mid X_t = x \right].$$

# SOC = KL-regularized RL

$\mathbb{P}^u(\cdot | X_0), \mathbb{P}(\cdot | X_0)$ : conditional probability measures of the controlled and uncontrolled processes.

From the Girsanov theorem:  $D_{\text{KL}}(\mathbb{P}^u(\cdot | X_0) \| \mathbb{P}(\cdot | X_0)) = \mathbb{E}_{X^u} \left[ \int_0^1 \frac{1}{2} \|u(X_t^u, t)\|^2 dt \right].$

$$\implies \max_{u \in \mathcal{U}} \mathbb{E}_{X_0 \sim p_0} \left[ \underbrace{\mathbb{E}_{X^u \sim \mathbb{P}^u(\cdot | X_0)} \left[ \int_0^1 -f(X_t^u, t) dt - g(X_1^u) \right] - D_{\text{KL}}(\mathbb{P}^u(\cdot | X_0) \| \mathbb{P}(\cdot | X_0))}_{\text{Solution : } \mathbb{P}^{u^*}(X | X_0) \propto \mathbb{P}(X | X_0) \exp \left( - \int_0^1 f(X_t, t) dt - g(X_1) \right)}. \right]$$

# The tilted path measure

$$\max_{u \in \mathcal{U}} \mathbb{E}_{X_0 \sim p_0} \left[ \underbrace{\mathbb{E}_{X^u \sim \mathbb{P}^u(\cdot | X_0)} \left[ \int_0^1 -f(X_t^u, t) dt - g(X_1^u) \right] - D_{\text{KL}}(\mathbb{P}^u(\cdot | X_0) \| \mathbb{P}(\cdot | X_0))}_{\text{Solution : } \mathbb{P}^{u^*}(X | X_0) \propto \mathbb{P}(X | X_0) \exp \left( - \int_0^1 f(X_t, t) dt - g(X_1) \right)}. \right],$$

Normalization constant of  $\mathbb{P}(X | X_0) \exp \left( - \int_0^1 f(X_t, t) dt - g(X_1) \right)$ :

$$\mathbb{E}_{X \sim \mathbb{P}(\cdot | X_0)} \left[ \exp \left( - \int_0^1 f(X_t, t) dt - g(X_1) \right) \right] = \exp(-V(X_0, 0)),$$

where  $V(x, t) = -\log \mathbb{E}_{X \sim \mathbb{P}} \left[ \exp \left( - \int_t^1 f(X_s, s) ds - g(X_1) \right) \mid X_t = x \right]$  is the value function.

$$\implies \mathbb{P}^{u^*}(X) = \mathbb{P}(X) \exp \left( - \int_0^1 f(X_t, t) dt - g(X_1) + V(X_0, 0) \right).$$

# Memoryless generative processes

Setting  $f(x, t) = 0$  and  $g(x) = -r(x)$ , we obtain  $\mathbb{P}^{u^*}(X) = \mathbb{P}(X)\exp(r(X_1) + V(X_0, 0))$ .

$$\implies \mathbb{P}^{u^*}(X_0, X_1) = \mathbb{P}(X_0, X_1)\exp(r(X_1) + V(X_0, 0)).$$

*Definition:* A generative process is **memoryless** if  $X_0$  and  $X_1$  are independent, i.e.  $\mathbb{P}(X_0, X_1) = \mathbb{P}(X_0)\mathbb{P}(X_1)$ .

$$\text{Then, } \mathbb{P}^{u^*}(X_1) = \int \mathbb{P}^{u^*}(X_0, X_1) dX_0 = \int \mathbb{P}(X_0)\mathbb{P}(X_1)\exp(r(X_1) + V(X_0, 0)) dX_0 \propto \mathbb{P}(X_1)\exp(r(X_1)).$$

Recall:  $dX_t = b(X_t, t) dt + \sigma(t) dB_t$ ,  $X_0 \sim \mathcal{N}(0, I)$ ,

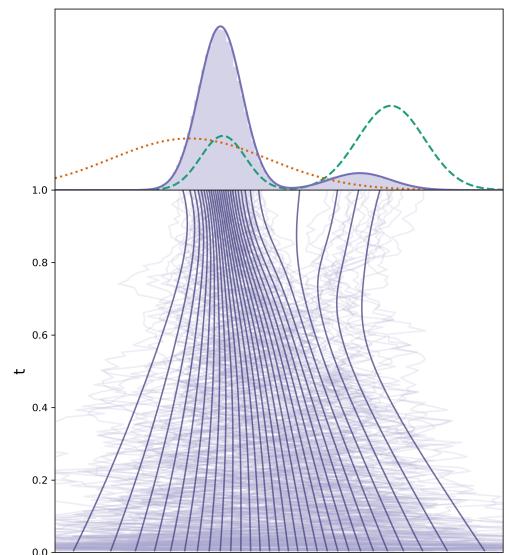
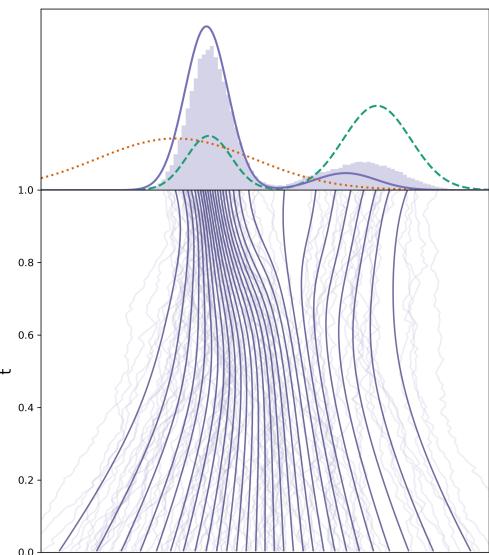
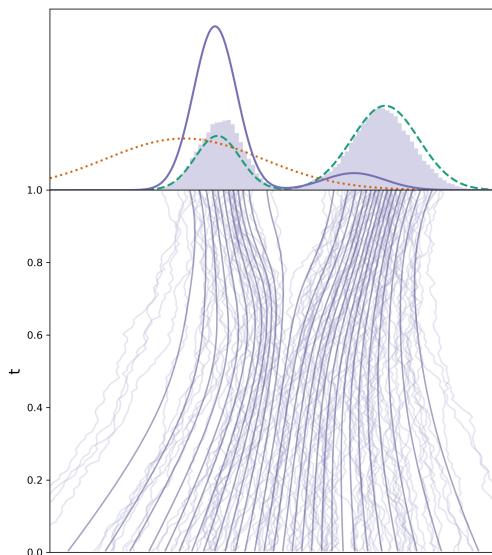
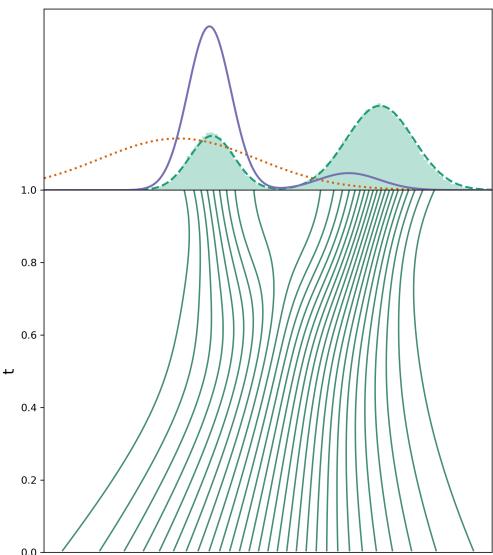
$$\text{where } b(x, t) = \kappa_t x + \left(\frac{\sigma(t)^2}{2} + \eta_t\right) \mathfrak{s}(x, t), \quad \kappa_t = \frac{\dot{\alpha}_t}{\alpha_t}, \quad \eta_t = \beta_t \left( \frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t \right)$$

*Proposition:*  $X$  is memoryless iff

$$\sigma(t)^2 = 2\eta_t + \chi(t), \text{ where } \chi : [0, 1] \rightarrow \mathbb{R} \text{ is s.t. } \forall t \in (0, 1], \quad \lim_{t' \rightarrow 0^+} \alpha_{t'} \exp\left(-\int_{t'}^t \frac{\chi(s)}{2\beta_s^2} ds\right) = 0.$$

# Fine-tuning with different noise schedules

—  $p^{\text{base}}(x)$     ...  $\exp(r(x))$     —  $p^*(x) \propto p^{\text{base}}(x)\exp(r(x))$



# Fine-tuning recipe

to sample the fine-tuned model with arbitrary noise schedules

Recall:  $dX_t = b(X_t, t) dt + \sigma(t) dB_t, \quad X_0 \sim \mathcal{N}(0, I),$

where  $b(x, t) = \kappa_t x + \left(\frac{\sigma(t)^2}{2} + \eta_t\right) \mathfrak{s}(x, t), \quad \kappa_t = \frac{\dot{\alpha}_t}{\alpha_t}, \quad \eta_t = \beta_t \left(\frac{\dot{\alpha}_t}{\alpha_t} \beta_t - \dot{\beta}_t\right)$

*Theorem:* In order to allow the use of arbitrary noise schedules and still generate samples according to the tilted distribution, the fine-tuning problem must be done with the memoryless noise schedule

$$\sigma(t) = \sqrt{2\eta_t}.$$

# Implicit parameterization of the control

$$\text{Recall: } b(x, t) = v(X_t, t) + \frac{\sigma(t)^2}{2\beta_t(\frac{\dot{\alpha}_t}{\alpha_t}\beta_t - \dot{\beta}_t)} \left( v(X_t, t) - \frac{\dot{\alpha}_t}{\alpha_t} X_t \right)$$

$$\sigma(t) = \sqrt{2\eta_t} \implies b(x, t) = 2v(x, t) - \frac{\dot{\alpha}_t}{\alpha_t} x$$

$$\left. \begin{array}{l} b(x, t) = 2v^{\text{base}}(x, t) - \frac{\dot{\alpha}_t}{\alpha_t} x \\ b(x, t) + \sigma(t)u(x, t) = 2v^{\text{finetune}}(x, t) - \frac{\dot{\alpha}_t}{\alpha_t} x \end{array} \right\} \implies u(x, t) = \sqrt{\frac{2}{\beta_t(\frac{\dot{\alpha}_t}{\alpha_t}\beta_t - \dot{\beta}_t)}} (v^{\text{finetune}}(x, t) - v^{\text{base}}(x, t))$$

# Methods for SOC: the adjoint method

Discrete adjoint method:  $\mathcal{L}(u; X^u) := \int_0^1 \left( \frac{1}{2} \|u(X_t^u, t)\|^2 + f(X_t^u, t) \right) dt + g(X_1^u).$

Continuous Adjoint method / Basic Adjoint Matching:

$$\mathcal{L}_{\text{Basic-Adj-Match}}(u; X^u) := \frac{1}{2} \int_0^1 \|u(X_t, t) + \sigma(t)^\top a(t; X^{\bar{u}}, \bar{u})\|^2 dt, \quad \bar{u} = \text{stopgrad}(u),$$

where the adjoint state satisfies the adjoint ODE:

$$\frac{d}{dt} a(t; X^u, u) = - \left[ \left( \nabla_x (b(X_t^u, t) + \sigma(t)u(X_t^u, t)) \right)^\top a(t; X^u, u) + \nabla_x \left( f(X_t^u, t) + \frac{1}{2} \|u(X_t^u, t)\|^2 \right) \right],$$

$$a(1; X^u, u) = \nabla g(X_1^u).$$

# Basic Adjoint Matching

$$\mathcal{L}_{\text{Basic-Adj-Match}}(u; X^u) := \frac{1}{2} \int_0^1 \|u(X_t, t) + \sigma(t)^\top a(t; X^{\bar{u}}, \bar{u})\|^2 dt, \quad \bar{u} = \text{stopgrad}(u)$$

Interpretation:  $a(t; X^u, u) := \nabla_{X_t^u} \left( \int_t^1 \left( \frac{1}{2} \|u(X_{t'}, t')\|^2 + f(X_{t'}, t') \right) dt' + g(X_1^u) \right),$

And the optimal control satisfies:

$$u^*(x, t) = -\sigma(t)^\top \nabla V(x, t) = -\sigma(t)^\top \nabla_x \mathbb{E} \left[ \int_t^1 \left( \frac{1}{2} \|u^*(X_s^{u^*}, s)\|^2 + f(X_s^{u^*}, s) \right) ds + g(X_1^{u^*}) \mid X_t^{u^*} = x \right].$$

Guarantee: The only critical point of  $\mathbb{E}[\mathcal{L}_{\text{Basic-Adj-Match}}]$  is the optimal control  $u^*$ .

Proof idea:  $u^*$  is the only solution of the fixed point equation above.

# Adjoint Matching

Adjoint ODE

$$\begin{aligned}\frac{d}{dt}a(t; X^u, u) &= - \left[ \left( \nabla_{X_t}(b(X_t, t) + \sigma(t)u(X_t, t)) \right)^\top a(t; X^u, u) + \nabla_{X_t} \left( f(X_t, t) + \frac{1}{2} \|u(X_t, t)\|^2 \right) \right], \\ a(1; X^u, u) &= \nabla_x g(X_1^u).\end{aligned}$$

Lean Adjoint ODE

$$\begin{aligned}\frac{d}{dt}\tilde{a}(t; X^u) &= - (\nabla_x b(X_t^u, t)^\top \tilde{a}(t; X^u) + \nabla_x f(X_t^u, t)), \\ \tilde{a}(1; X^u) &= \nabla_x g(X_1^u).\end{aligned}$$

Adjoint Matching loss:

$$\mathcal{L}_{\text{Adj-Match}}(u; X^u) := \frac{1}{2} \int_0^1 \|u(X_t, t) + \sigma(t)^\top \tilde{a}(t; X^u)\|^2 dt.$$

Key observation:

$$\mathbb{E} \left[ \nabla_{X_t} (\sigma(t)u^*(X_t^{u^*}, t))^\top a(t; X^{u^*}, u^*) + \nabla_{X_t} \frac{1}{2} \|u(X_t^{u^*}, t)\|^2 \right] = 0$$

The only critical point of  $\mathbb{E}[\mathcal{L}_{\text{Adj-Match}}]$  is the optimal control  $u^*$

# Recap

- New algorithm: **Adjoint Matching**
- Sample trajectories  $\mathbf{X}$  according to the SDE  $dX_t^u = (b(X_t^u, t) + \sigma(t)u(X_t^u, t)) dt + \sigma(t) dB_t$
- For each trajectory, compute:

$$\mathcal{L}_{\text{Adj-Match}}(u_\theta; \mathbf{X}) := \frac{1}{2} \int_0^1 \|u_\theta(X_t, t) + \sigma(t)\tilde{a}(t; \mathbf{X})\|^2 dt,$$

where  $\tilde{a}$  is the solution to the *Lean Adjoint ODE*:

$$\frac{d}{dt}\tilde{a}(t; \mathbf{X}) = -\nabla_x b(X_t, t)^\top \tilde{a}(t; \mathbf{X}), \quad \tilde{a}(1; \mathbf{X}) = -\nabla_x r(X_1).$$

- Compute  $\nabla_\theta \mathcal{L}_{\text{Adj-Match}}(u_\theta; \mathbf{X})$  and update  $\theta$  using e.g. Adam.

# Comparison to non-SOC method

$\lambda = 1000$



$\lambda = 2500$



$\lambda = 12500$



Adjoint Matching

*"Handsome Smiling man  
in blue jacket portrait"*



1000 it.



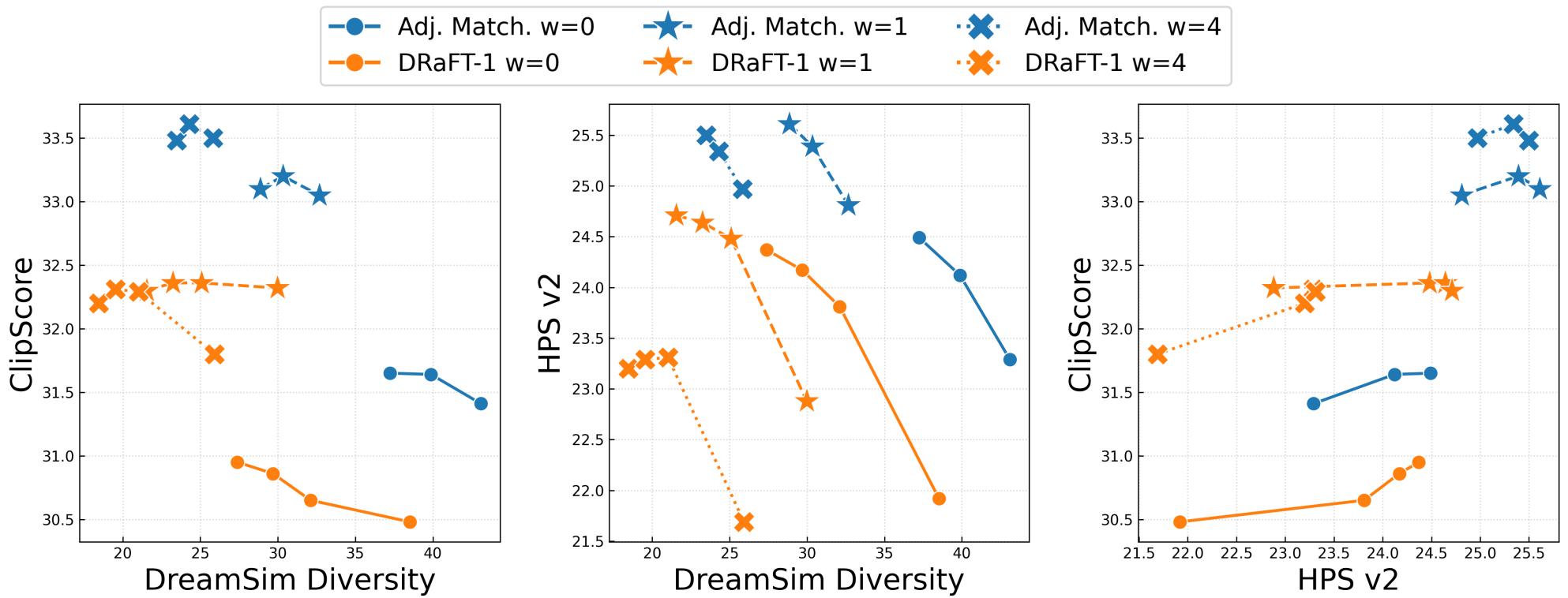
2000 it.



4000 it.

DRaFT-1 (Clark. et al., 2024)

# Comparison to non-SOC method



Tradeoffs between different aspects of generative models: text-to-image consistency (ClipScore), sample diversity for each prompt (DreamSim Diversity), and generalization to unseen human preferences (HPS v2).

Thank you!

# Backup: Stochastic optimal control

Important objects

Cost functional:

$$J(u; x, t) := \mathbb{E}_{X^u} \left[ \int_t^1 \left( \frac{1}{2} \|u(X_s^u, s)\|^2 + f(X_s^u, s) \right) ds + g(X_1^u) \mid X_t^u = x \right].$$

Value function:  $V(x, t) := \min_{u \in \mathcal{U}} J(u; x, t) = J(u^*; x, t)$ ,

Path integral representation of the value function:

$$V(x, t) = -\log \mathbb{E}_X \left[ \exp \left( - \int_t^1 f(X_s, s) ds - g(X_1) \right) \mid X_t = x \right].$$

Optimal control:  $u^*(x, t) = -\sigma(t)^\top \nabla_x V(x, t) = -\sigma(t)^\top \nabla_x J(u^*, x, t)$ .

Reminder: the SOC problem

$$\begin{aligned} \min_u \quad & \mathbb{E} \left[ \int_0^1 \left( \frac{1}{2} \|u(X_t^u, t)\|^2 + f(X_t^u, t) \right) dt + g(X_1^u) \right], \\ \text{s.t.} \quad & dX_t^u = (b(X_t^u, t) + \sigma(t)u(X_t^u, t)) dt + \sigma(t) dB_t, \\ & X_0^u \sim p_0. \end{aligned}$$