

# Stochastic Optimal Control Matching

Carles Domingo-Enrich, Jiequn Han, Brandon Amos, Joan Bruna, Ricky T. Q.  
Chen

Thomas Mousseau

August 28, 2025

# Overview

---

1. Setup and Preliminaries
2. Stochastic Optimal Control Matching
3. Experiments and results
4. Conclusion

# Generative Models

## The Core Challenge: Unnormalized Densities

Generative models must sample from complex distributions  $p_{\text{data}}(x) = \frac{1}{Z} \tilde{p}_{\text{data}}(x)$  where the normalization constant  $Z = \int \tilde{p}_{\text{data}}(x) dx$  is intractable to compute. This intractability arises from the curse of dimensionality when integrating over high-dimensional spaces.

**2020**      **DDPM:** Denoising Diffusion Probabilistic Models interpret generation as reversing a discrete noise-adding process, learning to denoise at each step. They produced high-quality samples but required thousands of slow sampling steps.

**2023**      **Flow Matching:** Flow matching views generation as learning a deterministic ODE vector field that directly transports a simple distribution (e.g., Gaussian) to data. This removed stochasticity and significantly improved efficiency compared to diffusion/score methods.

# SOC as the Foundation of Generative Models

---

## SOC Connection

Transform tractable distributions (Gaussian) to complex target distributions through **optimal control policies**.

This bridges the gap between:

- Simple sampling (easy)
- Complex data distributions (hard)

From Langevin dynamics to modern diffusion models, all major breakthroughs in generative modeling can be understood through the lens of stochastic optimal control theory.

# What is a Stochastic Control Problem?

## Control-Affine Stochastic Differential Equation

The general form of a controlled stochastic process:

$$dX_t^u = (b(X_t^u, t) + \sigma(t)u(X_t^u, t))dt + \sqrt{\lambda}\sigma(t)dB_t \quad (1)$$

**State Process:**  $X_t^u \in \mathbb{R}^d$  (system state under control  $u$  at time  $t$ )

**Drift Term:**  $b(X_t^u, t) \in \mathbb{R}^d$  (natural evolution of the system)

**Control Term:**  $u(X_t^u, t) \in \mathbb{R}^d$  (how control influences the system)

**Diffusion Coefficient:**  $\sigma(t) \in \mathbb{R}^{d \times d}$  (volatility matrix)

**Noise Process:**  $B_t \in \mathbb{R}^d$  (Brownian motion, external randomness)

**Noise Intensity:**  $\lambda > 0$  (controls the strength of stochastic perturbations)

# The Optimal Control Objective

## Cost Function to Minimize

Find the optimal control policy  $u^* \in \mathcal{U}$  that minimizes:

$$\min_{u \in \mathcal{U}} \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u(X_t^u, t)\|^2 + f(X_t^u, t) \right) dt + g(X_T^u) \right] \quad (2)$$

**Control Effort:**  $\frac{1}{2} \|u(X_t^u, t)\|^2$  (penalizes large control actions)

**Running Cost:**  $f(X_t^u, t)$  (ongoing cost during the process evolution)

**Terminal Cost:**  $g(X_T^u)$  (final cost based on end state at time  $T$ )

**Control Space:**  $\mathcal{U}$  (set of admissible control policies)

# SOC in Practice

Aspect	Steering Actuator	Diffusion Models	Flow Models
State $X_t$	Vehicle position, velocity	Noisy data sample	Clean data sample
Control $u_t$	Steering angle, throttle	Denoising direction	Flow velocity field
Dynamics Source	Newtonian mechanics	Forward noising process	No natural drift (learnable)
Drift $b(X_t, t)$	Vehicle dynamics	Predetermined schedule	$b = 0$ (control learns drift)
Control Goal	Reach target safely	Reverse noise process	Transport distributions
Noise $\sqrt{\lambda}\sigma dB_t$	Road disturbances	Brownian motion	Optional stochasticity

# Applications of SOC

## Key ML Applications of SOC

- **Reward fine-tuning of diffusion and flow models:** Optimizing generation quality using reward signals
- **Conditional sampling on diffusion and flow models:** Steering generation towards specific conditions or constraints
- **Sampling from unnormalized densities:** Efficiently drawing samples from complex, intractable distributions
- **Importance sampling of rare events in SDEs:** Computing probabilities of low-probability but critical events

## Key Insight

The prevalence of SOC formulations in modern ML motivates the need for more efficient and stable solving methods.



# Solve SOC Problems in Low Dimension

## Classical Approach: Hamilton-Jacobi-Bellman PDE

In small state spaces, discretize the state space and solve the HJB PDE directly. After applying dynamic programming and Bellman's optimality principle we get:

$$\frac{\partial V}{\partial t} + \min_u \left[ \frac{1}{2} \|u\|^2 + f(x, t) + \nabla V \cdot (b + \sigma u) + \frac{1}{2} \lambda \text{tr}(\sigma^T \nabla^2 V \sigma) \right] = 0 \quad (3)$$

## What You Get

- **Value function:**  $V(x, t)$
- **Optimal control:**  
 $u^*(x, t) = -\sigma(t)^T \nabla V(x, t)$
- **Global optimum** guaranteed
- **Theoretical guarantees**  
(convergence, stability)

## The Curse of Dimensionality

- **Grid size:**  $\mathcal{O}(N^d)$  where  $d$  is dimension
- **Memory:** Exponential growth with  $d$
- **Computation:** Infeasible for  $d > 3$

# Solve SOC Problems in High Dimension

## Modern Approach: Adjoint Methods

Cannot discretize high-dimensional spaces, so use **adjoint methods** to optimize control directly:

$$\min_u \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u(X_t^u, t)\|^2 + f(X_t^u, t) \right) dt + g(X_T^u) \right] \quad (4)$$

## What You Get

- **Scalable:** Works in high dimensions
- **Neural networks:** Can parameterize complex controls
- **Gradient-based:** Standard optimization techniques

## Major Problem

- **Non-convex landscape:** Full of local minima
- **Unstable training:** Difficult optimization
- **No guarantees:** May not find global optimum

# Similar Trend in Generative Modeling

## The Same Pattern: Non-convex $\rightarrow$ Least-Squares

Generative modeling experienced the same transition from non-convex to convex optimization

### Continuous Normalizing Flows

**Method:** Learn invertible transformations

$$\frac{dx}{dt} = f_{\theta}(x, t) \quad (5)$$

**Problem:**

- Use **adjoint methods** for gradients
- **Non-convex** optimization landscape
- Difficult training, unstable dynamics

### Denoising Diffusion Models

**Method:** Learn to reverse noise process

$$\mathbb{E}[\|\epsilon_{\theta}(x_t, t) - \epsilon\|^2] \quad (6)$$

**Advantage:**

- Use **least-squares loss**
- **Convex** functional landscape
- Stable, reliable optimization

# Optimization Landscapes

---

Task	Non-convex	Least Squares
Generative Modeling	Maximum Likelihood CNFs	Diffusion models and Flow Matching
Stochastic Optimal Control	Adjoint Methods	<b>Stochastic Optimal Control Matching</b>

# From Classical SOC to SOCM

## Dynamics (Same for Both)

$$dX_t^v = (b(X_t^v, t) + \sigma(t)v(X_t^v, t))dt + \sqrt{\lambda}\sigma(t)dB_t, \text{ with } X_0^v \sim p_0 \quad (7)$$

## Original SOC Cost Function

$$\min_{u \in \mathcal{U}} \mathcal{L}_{SOC} := \mathbb{E} \left[ \int_0^T \left( \frac{1}{2} \|u(X_t^u, t)\|^2 + f(X_t^u, t) \right) dt + g(X_T^u) \right] \quad (8)$$

## SOCM Cost Function

$$\min_{u, M} \mathcal{L}_{SOCM}(u, M) := \mathbb{E} \left[ \frac{1}{T} \int_0^T \|u(X_t^v, t) - w(t, v, X^v, B, M_t)\|^2 dt \times \alpha(v, X^v, B) \right] \quad (9)$$

**Key Point:** Same underlying stochastic dynamics, but SOCM uses a [least-squares matching approach](#) instead of direct minimization of the original cost.

# SOCM Parameters Definition

$$\min_{u, M} \mathcal{L}_{SOCM}(u, M) := \mathbb{E} \left[ \frac{1}{T} \int_0^T \|u(X_t^v, t) - w(t, v, X^v, B, M_t)\|^2 dt \times \alpha(v, X^v, B) \right]$$

Where:

$u : \mathbb{R}^d \times [0, 1] \rightarrow \mathbb{R}^d$  is the **control** (policy being learned)

$v : \mathbb{R}^d$  is a **fixed arbitrary control**,  $X^v$  is the solution of the SDE with control  $v$

$M : [0, 1]^2 \rightarrow \mathbb{R}^{d \times d}$  is the **reparameterization matrix**

$w$  : is the **matching vector field** (target for  $u$  to match)

$\alpha$  is the **importance weight** (for measure correction)

Key Dependencies

$w$  and  $\alpha$  depend on  $f, g, \lambda, \sigma$

# Presenting Matching Vector Field

## Reparameterization Function $w(t, v, X^\nu, B, M_t)$

The **matching vector field** computed via path-wise reparameterization:

$$\begin{aligned} w(t, v, X^\nu, B, M_t) = & \sigma(t)^\top \left( - \int_t^T M_t(s) \nabla_x f(X_s^\nu, s) ds - M_t(T) \nabla g(X_T^\nu) \right. \\ & + \int_t^T (M_t(s) \nabla_x b(X_s^\nu, s) - \partial_s M_t(s)) (\sigma^{-1}(s))^\top v(X_s^\nu, s) ds \quad (10) \\ & \left. + \sqrt{\lambda} \int_t^T (M_t(s) \nabla_x b(X_s^\nu, s) - \partial_s M_t(s)) (\sigma^{-1}(s))^\top dB_s \right) \end{aligned}$$

Where  $M_t(s)$  is the **reparameterization matrix** (learned jointly with  $u$ )

# Matching Vector Field (1/3) - Path Integral I

## Step 1 — Path-integral form of $u^*$ (Kappen 2005)

### What to Show

We want to eliminate the value function  $V$  from the optimal control formula and express  $u^*$  directly in terms of path integrals.

The value function is defined as the infimum over all controls:

$$V(x, t) = \inf_u J(u; x, t) = \inf_u \mathbb{E} \left[ \int_t^T \left( \frac{1}{2} \|u(X_s^u, s)\|^2 + f(X_s^u, s) \right) ds + g(X_T^u) \mid X_t^u = x \right]$$



# Matching Vector Field (1/3) - Path Integral II

---

We introduce a small step  $h$  so we can expand the cost-to-go function with dynamic programming recursion. Then using Itô's lemma (the stochastic version of the chain rule) this expansion produces the gradient and Hessian terms in the HJB. Shrinking  $h$  to zero converts the global Bellman principle into a local PDE.

$$V(x, t) = \min_u \mathbb{E} \left[ \int_t^{t+h} \left( \frac{1}{2} \|u_s\|^2 + f(X_s, s) \right) ds + V(X_{t+h}, t+h) \mid X_t = x \right]$$

## Based on Itô's lemma

Expand  $V(X_{t+h}, t+h)$  using Itô's lemma:

$$dV = \partial_t V ds + \nabla V^\top dX_s + \frac{1}{2} \text{tr}[\sigma \sigma^\top \nabla^2 V] \lambda ds.$$

## Matching Vector Field (1/3) - Path Integral III

---

Substituting  $dX_s = b ds + \sigma u ds + \sqrt{\lambda} \sigma dB_s$ , and taking expectation (the Brownian term vanishes):

$$\mathbb{E}[dV] = \left[ \partial_t V + \nabla V \cdot (b + \sigma u) + \frac{\lambda}{2} \text{tr}(\sigma^\top \nabla^2 V \sigma) \right] ds.$$

The expected increment equation plugged into the dynamic programming principle gave you the HJB PDE

$$\frac{\partial V}{\partial t} + f(x, t) + \nabla V \cdot (b + \sigma u) + \frac{\lambda}{2} \text{tr}(\sigma^\top \nabla^2 V \sigma) + \frac{1}{2} \|u\|^2 = 0$$

Minimizing this gives us the **optimal control**:

$$u^*(x, t) = -\sigma^\top(t) \nabla V(x, t)$$

# Matching Vector Field (1/3) - Path Integral IV

---

And the reduced HJB using  $u^*$  is:

$$\frac{\partial V}{\partial t} + f + b \cdot \nabla V + \frac{\lambda}{2} \text{tr}(\sigma^\top \nabla^2 V \sigma) - \frac{1}{2} \|\sigma^\top \nabla V\|^2 = 0, \quad V(x, T) = g(x)$$

**The HJB is nonlinear because of the quadratic gradient term. To linearize it, we do the Cole–Hopf transformation:**

$$V(x, t) = -\lambda \log \phi(x, t).$$

# Matching Vector Field (1/3) - Path Integral V

---

Plugging into the reduced HJB cancels the  $-\frac{1}{2}\|\sigma^\top \nabla V\|^2$  term and gives a linear PDE for  $\phi$ :

$$\frac{\partial \phi}{\partial t} + L\phi + \frac{1}{\lambda}f(x, t)\phi = 0, \quad \phi(x, T) = e^{-g(x)/\lambda}.$$

Here  $L$  is the infinitesimal generator of the uncontrolled diffusion:

$$L\phi = b \cdot \nabla \phi + \frac{\lambda}{2} \text{tr}(\sigma \sigma^\top \nabla^2 \phi).$$

**The Feynman–Kac theorem says: the solution to this linear PDE equals an expectation over uncontrolled trajectories:**

$$\phi(x, t) = \mathbb{E}_0 \left[ \exp \left( -\frac{1}{\lambda} \int_t^T f(X_s, s) ds - \frac{1}{\lambda} g(X_T) \right) \middle| X_t = x \right].$$

# Matching Vector Field (1/3) - Path Integral VI

---

Thus

$$V(x, t) = -\lambda \log \phi(x, t).$$

**Remember**

$$u^*(x, t) = -\sigma^\top \nabla V(x, t).$$

Substitute  $V = -\lambda \log \phi$ :

$$u^*(x, t) = \lambda \sigma^\top(t) \nabla_x \log \phi(x, t).$$

# Matching Vector Field (1/3) - Path Integral VII

Define the "partition function" (unnormalized density):

$$Z(x, t) = \mathbb{E} \left[ \exp \left( -\frac{1}{\lambda} \int_t^T f(X_s, s) ds - \frac{1}{\lambda} g(X_T) \right) \middle| X_t = x \right].$$

Then the final path–integral representation is:

$$u^*(x, t) = \lambda \sigma^\top(t) \nabla_x \log Z(x, t).$$

This is exactly Lemma 1 (Kappen 2005): the optimal control is the score of a

## Path-integral Optimal Control

$$u^*(x, t) = \lambda \sigma^\top(t) \nabla_x \log \mathbb{E} \left[ \exp \left( -\frac{1}{\lambda} \int_t^T f(X_s, s) ds - \frac{1}{\lambda} g(X_T) \right) \middle| X_t = x \right]$$

# Matching Vector Field (1/3) - Path Integral VIII

---

## Why This Matters

**Key Achievement:** We eliminate  $V$  to avoid solving a nonlinear PDE and reframe the control as a **score function**  $\nabla_x \log \mathbb{E}[\cdot | X_t = x]$ .

We avoid solving a nonlinear PDE and instead target a score that admits a pathwise estimator (next slide). This sets up a score we can approximate with a **low-variance pathwise estimator** in Step 2, which will become the matching vector field  $w$ . This is the foundation for SOCM's **least-squares matching** approach.

# Matching Vector Field (2/3) — Optimal Control I

## Step 2 — Rendering $u^*$ Trackable using Pathwise reparameterization (under the uncontrolled law)

### Reasons Why $u^*$ is not Trackable

Je dois expliquer pourquoi on ne peut pas prendre  $u$  optimal directement.

### Ideal Loss Function

$$\tilde{\mathcal{L}}(u) = \mathbb{E} \left[ \frac{1}{T} \int_0^T \|u(X_t, t) - u^*(X_t, t)\|^2 dt \exp \left( -\lambda^{-1} \int_0^T f(X_t, t) dt - \lambda^{-1} g(X_T) \right) \right] \quad (11)$$



## Matching Vector Field (2/3) — Optimal Control II

Let's expand the squared term  $\|u(X_t, t) - u^*(X_t, t)\|^2$ :

$$\begin{aligned} \tilde{\mathcal{L}}(u) = \mathbb{E} & \left[ \frac{1}{T} \int_0^T \left( \|u(X_t, t)\|^2 - 2\langle u(X_t, t), u^*(X_t, t) \rangle + \|u^*(X_t, t)\|^2 \right) dt \right. \\ & \left. \times \exp \left( -\lambda^{-1} \int_0^T f(X_t, t) dt - \lambda^{-1} g(X_T) \right) \right] \end{aligned} \quad (12)$$

### The Fundamental Challenge

**We cannot compute  $u^*(X_t, t)$  directly!** This is exactly why we need SOCM's matching approach - to find a computable target  $w$  that serves as a proxy for  $u^*$ .

# Matching Vector Field (2/3) — Optimal Control III

## Why Focus on the Cross Term?

We concentrate on transforming the cross term  $\langle u(X_t, t), u^*(X_t, t) \rangle$  because:

**First term**  $\|u(X_t, t)\|^2$ : Straightforward to compute - depends only on our learned control  $u$

**Third term**  $\|u^*(X_t, t)\|^2$ : Constant with respect to  $u$  - doesn't affect the optimization

**Cross term:** The **only term** that requires the path-integral formulation since it involves both  $u$  and  $u^*$

Using the path-integral representation from Step 1, we can transform the problematic cross term:

## Matching Vector Field (2/3) — Optimal Control IV

### Reminder: Path-integral Optimal Control

$$u^*(x, t) = \lambda \sigma^\top(t) \nabla_x \log \mathbb{E} \left[ \exp \left( -\frac{1}{\lambda} \int_t^T f(X_s, s) ds - \frac{1}{\lambda} g(X_T) \right) \middle| X_t = x \right]$$

$$\begin{aligned} & \mathbb{E} \left[ \frac{1}{T} \int_0^T \langle u(X_t, t), u^*(X_t, t) \rangle dt \exp \left( -\lambda^{-1} \int_0^T f(X_t, t) dt - \lambda^{-1} g(X_T) \right) \right] \\ &= \lambda \mathbb{E} \left[ \frac{1}{T} \int_0^T \left\langle u(X_t, t), \sigma(t)^\top \nabla_x \mathbb{E} \left[ \exp \left( -\lambda^{-1} \int_t^T f(X_s, s) ds - \lambda^{-1} g(X_T) \right) \middle| X_t = x \right] \right\rangle \right. \\ & \quad \left. \times \exp \left( -\lambda^{-1} \int_0^t f(X_s, s) ds \right) dt \right] \end{aligned} \tag{13}$$

# Matching Vector Field (2/3) — Optimal Control V

## Next Challenge

The transformed cross term still contains an intractable expectation. This is where path-wise reparameterization becomes essential to make this computable since Monte-Carlo “score-function” (REINFORCE) estimators for this expectation have prohibitively high variance.

## The Challenge: Intractable Conditional Expectation

From Step 2, we need to make this computable for training:

$$\nabla_x \mathbb{E} \left[ \exp \left( -\lambda^{-1} \int_t^T f(X_s, s) ds - \lambda^{-1} g(X_T) \right) \middle| X_t = x \right]$$

**Standard approach:** Monte-Carlo score estimators  $\rightarrow$  **prohibitively high variance**

# Matching Vector Field (2/3) — Optimal Control (continued) I

This transformation comes from a general theorem for Fréchet-differentiable functionals  $F(X)$ :

## Proposition 4: Path-wise Reparameterization

$$\begin{aligned} \nabla_x \mathbb{E}[\exp(-F(X)) | X_0 = x] = & \mathbb{E} \left[ (-\nabla_z F(X + \psi)) \Big|_{z=0} \right. \\ & + \frac{1}{\lambda} \int_0^T (\nabla_z \psi \nabla_x b - \nabla_z \partial_s \psi) (\sigma^{-1})^T dB_s \Big) \quad (14) \\ & \left. \times e^{-F(X)} \Big|_{X_0 = x} \right] \end{aligned}$$

# Matching Vector Field (2/3) — Optimal Control (continued) II

---

where  $\psi(z, s)$  is a "path shift" function with  $\psi(z, 0) = z$  and  $\psi(0, s) = 0$ .

## SOCM's Specific Choice

SOCM chooses a **linear path shift**:  $\psi_t(z, s) = M_t(s)^T z$  with  $M_t(t) = I$

This gives:  $\nabla_z \psi_t(0, s) = M_t(s)$  and  $\nabla_z \partial_s \psi_t(0, s) = \partial_s M_t(s)$

The reparameterization matrix  $M_t(s)$  becomes a learnable parameter for variance reduction!

# Matching Vector Field (2/3) — Optimal Control (continued) III

## Path-wise Reparameterization Solution

Transform the intractable gradient into a computable expectation:

$$\begin{aligned} & \nabla_x \mathbb{E} \left[ \exp \left( -\lambda^{-1} \int_t^T f(X_s, s) ds - \lambda^{-1} g(X_T) \right) \middle| X_t = x \right] \\ &= \mathbb{E} \left[ \left( -\lambda^{-1} \int_t^T M_t(s) \nabla_x f(X_s, s) ds - \lambda^{-1} M_t(T) \nabla g(X_T) \right. \right. \\ & \quad \left. \left. + \lambda^{-1/2} \int_t^T (M_t(s) \nabla_x b(X_s, s) - \partial_s M_t(s)) (\sigma^{-1})^T(s) dB_s \right) \right. \\ & \quad \left. \times \exp \left( -\lambda^{-1} \int_t^T f(X_s, s) ds - \lambda^{-1} g(X_T) \right) \middle| X_t = x \right] \end{aligned} \tag{15}$$

# Matching Vector Field (3/3) - Girsanov Theorem I

---

## Step 3 — Girsanov Theorem



# Matching Vector Field (3/3) - Girsanov Theorem II

## SOCM Loss in Construction

From Steps 1-2, we have the theoretical loss function:

$$\begin{aligned}\tilde{\mathcal{L}}(u) = \mathbb{E} \left[ \frac{1}{T} \int_0^T \left\| u(X_t, t) + \sigma(t)^T \left( \int_t^T M_t(s) \nabla_x f(X_s, s) ds + M_t(T) \nabla g(X_T) \right. \right. \right. \\ \left. \left. \left. - \lambda^{1/2} \int_t^T (M_t(s) \nabla_x b(X_s, s) - \partial_s M_t(s)) (\sigma^{-1})^T(s) dB_s \right) \right. \right. \\ \left. \left. \times \exp \left( -\lambda^{-1} \int_0^T f(X_t, t) dt - \lambda^{-1} g(X_T) \right) + K \right] \quad (16)\end{aligned}$$

**Problem:** This expectation is under the **uncontrolled** measure, but we need to sample using an **arbitrary control**  $v$  for practical training.

# Matching Vector Field (3/3) - Girsanov Theorem (continued) I

---

## Girsanov Measure Change

To transform from uncontrolled to controlled sampling, Girsanov theorem:

**Changes the drift:**  $dX_t = bdt + \sqrt{\lambda}\sigma dB_t \rightarrow dX_t^\nu = (b + \sigma\nu)dt + \sqrt{\lambda}\sigma dB_t$

**Introduces importance weight:**

$$\alpha(\nu, X^\nu, B) = \exp\left(-\frac{1}{\sqrt{\lambda}} \int_0^T \langle \nu, dB_t \rangle - \frac{1}{2\lambda} \int_0^T \|\nu\|^2 dt\right)$$

**Modifies the exponential weight:** The original  $\exp(-\lambda^{-1} \int f - \lambda^{-1} g)$  becomes part of the full importance weight.

# Matching Vector Field (3/3) - Girsanov Theorem (continued) II

## Final SOCM Loss

After Girsanov transformation, we get the practical SOCM loss:

$$\mathcal{L}_{SOCM}(u, M) = \mathbb{E} \left[ \frac{1}{T} \int_0^T \|u(X_t^v, t) - w(t, v, X^v, B, M_t)\|^2 dt \times \alpha(v, X^v, B) \right] \quad (17)$$

where:

- $w$  is the matching vector field (from the reparameterization terms)
- $\alpha$  is the complete importance weight including costs and control effort
- Trajectories  $X^v$  are generated using the practical control  $v$

# Matching Vector Field (3/3) - Girsanov Theorem (continued) III

---

## Key Achievement

**Girsanov bridges theory and practice:** We can train on controlled trajectories  $X^v$  while the importance weight  $\alpha$  ensures we're still optimizing the original SOC objective. The extra drift shift from  $v$  generates both the  $v$ -dependent terms in  $w$  and the corrective importance weight  $\alpha$ .

# SOCM Algorithm

---

---

**Algorithm 2** Stochastic Optimal Control Matching (SOCM)

---

**Input:** State cost  $f(x, t)$ , terminal cost  $g(x)$ , diffusion coeff.  $\sigma(t)$ , base drift  $b(x, t)$ , noise level  $\lambda$ , number of iterations  $N$ , batch size  $m$ , number of time steps  $K$ , initial control parameters  $\theta_0$ , initial matrix parameters  $\omega_0$ , loss  $\mathcal{L}_{\text{SOCM}}$  in (125)

```
1 for  $n \in \{0, \dots, N-1\}$  do
2   Simulate  $m$  trajectories of the process  $X^v$  controlled by  $v = u_{\theta_n}$ , e.g., using Euler-Maruyama updates
3   Detach the  $m$  trajectories from the computational graph, so that gradients do not backpropagate
4   Using the  $m$  trajectories, compute an  $m$ -sample Monte-Carlo approximation  $\hat{\mathcal{L}}_{\text{SOCM}}(u_{\theta_n}, M_{\omega_n})$  of the loss
       $\mathcal{L}_{\text{SOCM}}(u_{\theta_n}, M_{\omega_n})$  in (125)
5   Compute the gradients  $\nabla_{(\theta, \omega)} \hat{\mathcal{L}}_{\text{SOCM}}(u_{\theta_n}, M_{\omega_n})$  of  $\hat{\mathcal{L}}_{\text{SOCM}}(u_{\theta_n}, M_{\omega_n})$  at  $(\theta_n, \omega_n)$ 
6   Obtain  $\theta_{n+1}, \omega_{n+1}$  with via an Adam update on  $\theta_n, \omega_n$ , resp.
7 end
```

**Output:** Learned control  $u_{\theta_N}$

---

Figure: Stochastic Optimal Control Matching (SOCM) Algorithm

# Ornstein-Uhlenbeck Process Benchmarks I

## What's an Ornstein-Uhlenbeck (OU) Process?

A **mean-reverting** stochastic differential equation with linear drift and Gaussian noise:

$$dX_t = \theta(\mu - X_t)dt + \sigma dB_t \quad (1D \text{ case})$$

For multi-dimensional benchmarks:  $dX_t = Axdt + \sigma dB_t$  with linear drift  $b(x, t) = Ax$

## Double-Well (Challenging Nonlinear Case)

**Nonlinear drift from quartic potential:**  $\Psi(x) = \sum_i \kappa_i (x_i^2 - 1)^2$

**Multimodal terminal cost:**  $g(x)$  has  $2^d$  modes (exponentially many!)

**Highly challenging:** Tests algorithm stability in complex, multimodal landscapes

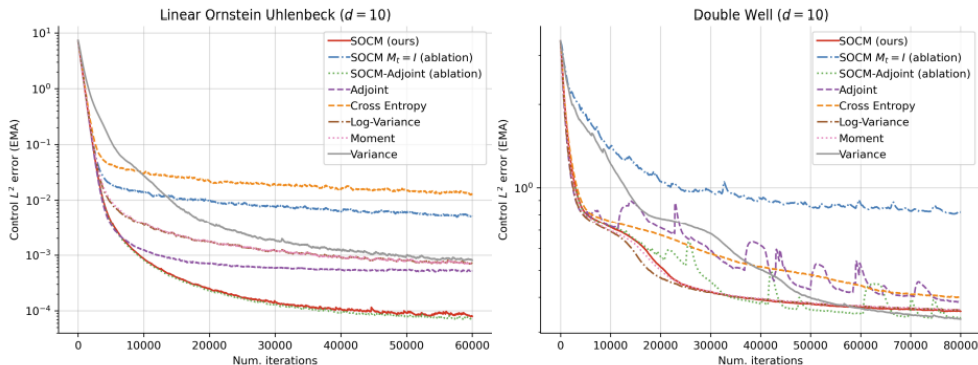
# Ornstein-Uhlenbeck Process Benchmarks II

---

## Why OU-type Problems are Standard SOC Testbeds

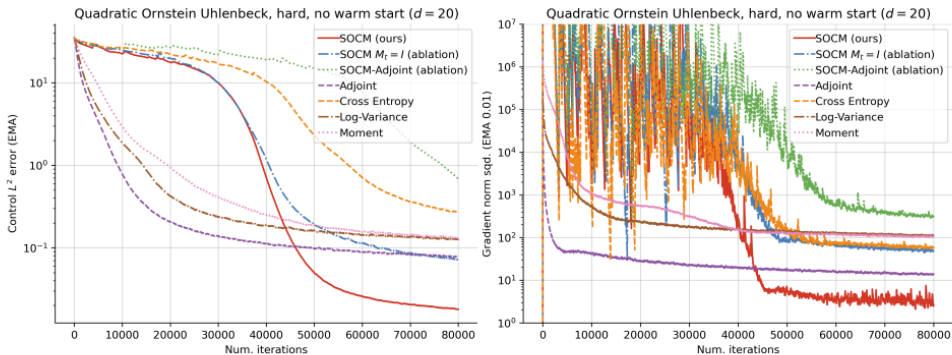
- **Ground truth available:** Analytic  $u^*$  for LQR/linear cases  $\rightarrow$  direct  $L_2$  control error measurement
- **Controlled difficulty:** Can dial complexity via matrix scaling ( $A$ ,  $P$ ,  $Q$  parameters)
- **Stress test variance:** Reveals importance weight and gradient signal-to-noise issues
- **Multimodal landscapes:** Double-Well tests stability in challenging optimization scenarios

# Experimental Results (1/2)





## Experimental Results (2/2)



**Figure 3** Plots of the  $L^2$  error incurred by the learned control (*top*), and the norm squared of the gradient with respect to the parameters  $\theta$  of the control (*bottom*), for the QUADRATIC ORNSTEIN UHLENBECK (HARD) setting and for each IDO loss. All the algorithms use a warm-started control (see [Appendix D](#)).

# Stochastic Optimal Control Matching

## Experiments and results

### Experimental Results (2/2)

At the end of training, SOCM obtains the lowest L2 error, improving over all existing methods by a factor of around ten. The two SOCM ablations come in second and third by a substantial difference, which underlines the importance of the path-wise reparameterization trick.

JE DOIS COMPRENDRE CE QUE EST UN ORNSTEIN UHLENBECK PROCESS

## Experimental Results (2/2)

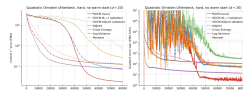


Figure 3: Plots of the  $L^2$  error (top), and the norm squared of the gradient with respect to the parameters  $\theta$  of the control (bottom), for the QG solver (Quasi-Newton), the SOCM (Stochastic Optimal Control Matching) solver and for each DSO line. All the algorithms use a reparameterized control (see Appendix D).

# Conclusion

---

## Personal thoughts and conclusion

SOCM is a framework grounded in theory that has shown interesting empirical results in specific benchmarks, but to truly judge its potential further application of this framework needs to be explored.

This said the authors seems to have already identified a few difficulties with this approach being the high variance related to the factor  $\alpha(v, Xv, B)$  due to a mismatch between the probability measures induced by the learned control and the optimal control. This is a problem that is also encountered in out-of-distribution generalization for reinforcement learning.

Additionally, I would love to come back to this paper in the next year after having acquired a proper foundation in stochastic calculus, partial differential equations, and statistical mechanics.

# Stochastic Optimal Control Matching

└ Conclusion

└ Conclusion

## Conclusion

### Personal thoughts and conclusion

SOCM is a framework grounded in theory that has shown interesting empirical results in specific benchmarks, but to truly judge its potential further application of this framework needs to be explored.

This said the authors seems to have already identified a few difficulties with this approach being the high variance related to the factor  $\alpha(v, Xv, B)$  due to a mismatch between the probability measures induced by the learned control and the optimal control. This is a problem that is also encountered in out-of-distribution generalization for reinforcement learning.

Additionally, I would love to come back to this paper in the next year after having acquired a proper foundation in stochastic calculus, partial differential equations, and statistical mechanics.

The main roadblock when we try to apply SOCM to more challenging problems is that the variance of the factor  $\alpha(v, Xv, B)$  explodes when  $f$  and/or  $g$  are large, or when the dimension  $d$  is high. The control L2 error for the SOCM and cross-entropy losses remains high and fluctuates heavily due to the large variance of  $\alpha$ . The large variance of  $\alpha$  is due to the mismatch between the probability measures induced by the learned control and the optimal control. Similar problems are encountered in out-of-distribution generalization for reinforcement learning, and some approaches may be carried over from that area (Munos et al., 2016).