

# Stochastic Optimal Control Matching

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

Thomas Mousseau

August 26, 2025

# Overview

---

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

# Evolution of Generative Models

---

- 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.
- 2021**     **Score-based Models:** Score-based generative models extended diffusion to continuous-time SDEs, learning the score function ( $\nabla_x \log p_t(x)$ ) to reverse a stochastic diffusion process. This unified diffusion with stochastic control, allowed probability flow ODEs, and sped up sampling.
- 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

## 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.

## SOC Connection

### Key Insight:

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

This bridges the gap between:

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

## Modern Implementations

### Diffusion Models:

$$u_t = -\frac{1}{2} \nabla_x \log p_t(x) \text{ (denoising)}$$

### Score-based Models:

$$u_t = \nabla_x \log p_t(x) \text{ (score function)}$$

### Flow Matching:

$$u_t = \frac{x_1 - x_0}{T - t} \text{ (deterministic flow)}$$

All learn **optimal control policies** to transport distributions!

# Stochastic Optimal Control Matching

## └ Setup and Preliminaries

## └ SOC as the Foundation of Generative Models

**Unnormalized Densities:** The fundamental challenge in generative modeling is sampling from distributions  $p(x) = \frac{1}{Z} e^{-E(x)}$  where  $Z$  is unknown. SOC provides the mathematical framework to construct sampling procedures.

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

### SOC as the Foundation of Generative Models

#### The Core Challenge: Unnormalized Densities

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

#### SOC Connection

##### Key Insight:

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

This bridges the gap between:

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

#### Modern Implementations

##### Diffusion Models:

$u_t = -\frac{1}{\epsilon} \nabla_x \log p_t(x)$  (denoising)

##### Score-based Models:

$u_t = \nabla_x \log p_t(x)$  (score function)

##### Flow Matching:

$u_t = \eta \frac{dx}{dt}$  (deterministic flow)

All learn **optimal control policies** to transport distributions!

# 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:**  $\sigma(t)u(X_t^u, t) \in \mathbb{R}^d$  (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)

# Stochastic Optimal Control Matching

## └ Setup and Preliminaries

## └ What is a Stochastic Control Problem?

### 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:**  $\sigma(t)u(X_t^u, t) \in \mathbb{R}^d$  (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)

**Control-Affine Structure:** The "control-affine" property means the control  $u$  enters linearly (affinely) in the drift term. This is the most general practical form for controlled SDEs, encompassing most applications in finance, robotics, and machine learning.

Je dois aussi expliquer l'importance que le steering input ainsi que le noise sont tous les 2 multipliés par sigma(t) ce qui signifie que ...

# 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)



# Stochastic Optimal Control Matching

## └ Setup and Preliminaries

## └ The Optimal Control Objective

### The Optimal Control Objective

#### Cost Function to Minimize

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

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

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

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

**Terminal Cost:**  $g(X_T^*)$  (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 (vehicle dynamics)	Forward noising process	No natural drift (learnable)
Drift $b(X_t, t)$	Kinematic equations	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

# Stochastic Optimal Control Matching

## └ Setup and Preliminaries

## └ SOC in Practice

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 (vehicle dynamics)	Forward noising process	No natural drift (learnable)
Drift $b(X_t, t)$	Kinematic equations	Predetermined schedule	$b = 0$ (control learns drift)
Control Goal	Reach target safely	Reverse noise process	Transport distributions
Noise $\sqrt{\gamma} \sigma dB_t$	Road disturbances	Brownian motion	Optional stochasticity

**Steering Actuator:** The dynamics come from well-established physics - Newton's laws, kinematics, and vehicle dynamics. The control optimizes safety and efficiency.

**Diffusion Models:** The drift is determined by the forward noising process (e.g.,  $\beta(t)$ ), and control learns to reverse this predetermined corruption.

**Flow Models:** No natural drift exists - the control  $u_t$  directly becomes the drift term, learning the entire velocity field that transports distributions.

# 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.

# Stochastic Optimal Control Matching

- └ Stochastic Optimal Control Matching
- └ Applications of SOC

## 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.

Je dois etre capable de vraiment bien expliquer les applications du SOC dans le ML moderne pour pouvoir demontrer l'importance d'avoir pris ce papier

# Solve SOC Problems in Low Dimension

## Classical Approach: Hamilton-Jacobi-Bellman PDE

In small state spaces (1D, 2D, sometimes 3D), discretize the state space and solve the HJB PDE directly:

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

## 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  $d > 3$
- **Impossible** to discretize for  $d \gg 1$

# Stochastic Optimal Control Matching

## └ Stochastic Optimal Control Matching

## └ Solve SOC Problems in Low Dimension

**Works perfectly in low dimensions, but state space grows exponentially with dimension.**

**Classical HJB:** In low dimensions, you can discretize the entire state space on a grid and solve the HJB PDE using finite difference methods. This gives you the exact solution but becomes computationally impossible as dimension increases.

**Curse of Dimensionality:** For a  $d$ -dimensional problem with  $N$  grid points per dimension, you need  $N^d$  total grid points. Even modest problems ( $d = 10$ ,  $N = 100$ ) require  $100^{10} = 10^{20}$  grid points.

### Solve SOC Problems in Low Dimension

#### Classical Approach: Hamilton-Jacobi-Bellman PDE

In small state spaces (1D, 2D, sometimes 3D), discretize the state space and solve the HJB PDE directly:

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

#### What You Get

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

#### The Curse of Dimensionality

- **Grid size:**  $O(N^d)$  where  $d$  is dimension
- **Memory:** Exponential growth with  $d$
- **Computation:** Infeasible for  $d > 3$
- **Impossible** to discretize for  $d \gg 1$

# 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
- **Sensitive initialization:** Results vary significantly



# Stochastic Optimal Control Matching

## └ Stochastic Optimal Control Matching

## └ Solve SOC Problems in High Dimension

### Solve SOC Problems in High Dimension

#### Modern Approach: Adjoint Methods

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

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

#### 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
- **Sensitive initialization:** Results vary

**Highly non-convex functional landscape makes optimization extremely challenging.**

Traditional SOC methods struggle with local minima and unstable training dynamics.

**Adjoint Methods:** These methods solve the forward SDE and then use the adjoint equation to compute gradients efficiently. This allows gradient-based optimization of the control parameters without discretizing the state space.

**Non-convex Optimization:** The major challenge is that the resulting optimization landscape is highly non-convex, leading to difficult optimization with many local minima and unstable training dynamics.

# 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

# Stochastic Optimal Control Matching

## └ Stochastic Optimal Control Matching

## └ Similar Trend in Generative Modeling

### 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

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

Advantage:

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

**Insight:** Apply the same principle to SOC! Replace non-convex adjoint methods with **least-squares matching**, creating a **convex optimization landscape** for stochastic optimal control.

**Historical Parallel:** Continuous Normalizing Flows (CNFs) suffered from the same optimization challenges as traditional SOC - they used adjoint methods and had non-convex landscapes. DDPMs revolutionized generative modeling by reformulating the problem as least-squares regression.

**SOCM's Contribution:** SOCM brings the same insight to stochastic optimal control, replacing direct optimization with least-squares matching to achieve stable, convex optimization.

# SOCM in Context: 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>

# SOCM Framework: 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)

$v$  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  $f, g, \lambda, \sigma$

# Stochastic Optimal Control Matching

## └ Stochastic Optimal Control Matching

### └ SOCM Parameters Definition

#### SOCM Parameters Definition

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

Where:

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

$v$  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  $f, g, \lambda, \sigma$

**Control  $u$ :** This is what we're trying to learn - the optimal policy that minimizes our cost function.

**Arbitrary Control  $v$ :** Used to generate sample trajectories for training. The choice of  $v$  affects the efficiency but not the correctness of SOCM.

**Reparameterization Matrix  $M$ :** Jointly optimized with  $u$  to reduce variance in the gradient estimation through path-wise reparameterization.

# The SOCM Framework (2/3)

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

The **target matching function** 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  $u$ )



# The SOCM Framework (2/3)

## Importance Weight $\alpha(v, X^v, B)$

The **importance sampling weight** for measure correction:

$$\alpha(v, X^v, B) = \exp \left( -\frac{1}{\lambda} \int_0^T f(X_t^v, t) dt - \frac{1}{\lambda} g(X_T^v) \right. \\ \left. - \frac{1}{\sqrt{\lambda}} \int_0^T \langle v(X_t^v, t), dB_t \rangle - \frac{1}{2\lambda} \int_0^T \|v(X_t^v, t)\|^2 dt \right) \quad (11)$$

Incorporates **running costs**, **terminal costs**, and **control effort**

# Stochastic Optimal Control Matching

- Stochastic Optimal Control Matching
- The SOCM Framework (2/3)

Importance Weight  $\alpha(v, X^*, B)$ The **importance sampling weight** for measure correction:

$$\alpha(v, X^*, B) = \exp \left( -\frac{1}{\lambda} \int_0^T \ell(X_t^*, t) dt - \frac{1}{\lambda} g(X_T^*) - \frac{1}{\sqrt{\lambda}} \int_0^T (v(X_t^*, t), dB_t) - \frac{1}{2\lambda} \int_0^T \|v(X_t^*, t)\|^2 dt \right) \quad (11)$$

Incorporates **running costs**, **terminal costs**, and **control effort**

**Reparameterization Matrix  $M_t$ :** This matrix enables path-wise reparameterization, a technique to reduce variance in gradient estimation by reparameterizing the stochastic process. It's optimized jointly with the control  $u$  to minimize the overall SOCM loss.

**Importance Weight  $\alpha$ :** This exponential weight corrects for the mismatch between the learned control measure and the optimal control measure. However, it can have high variance when costs are large or in high dimensions, which is the main limitation of SOCM.

# 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)

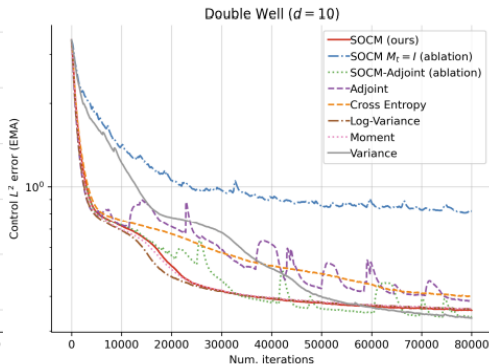
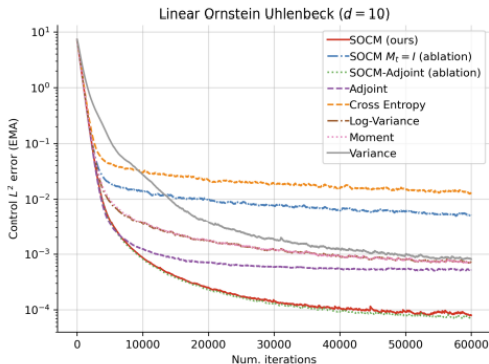
```
1 for  $n \in \{0, \dots, N-1\}$  do
2   Simulate  $m$  trajectories of the process  $X^v$  controlled by  $v = u_{\theta_n}$ , e.g., using Euler-Maruyama updates
3   Detach the  $m$  trajectories from the computational graph, so that gradients do not backpropagate
4   Using the  $m$  trajectories, 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  $u_{\theta_N}$

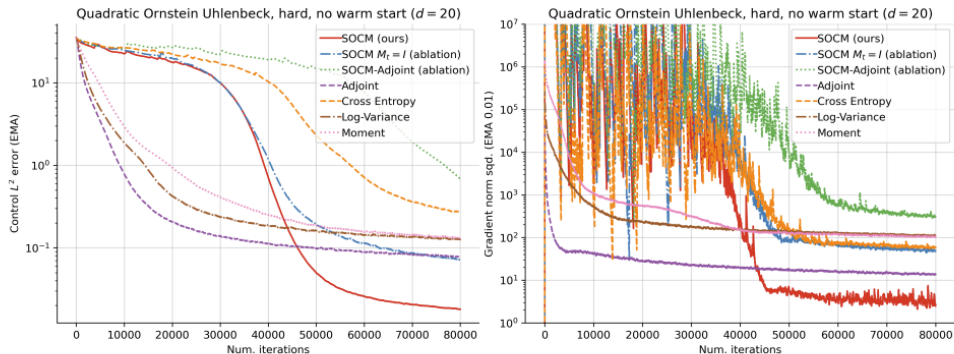
---

Figure: Stochastic Optimal Control Matching (SOCM) Algorithm

# 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)

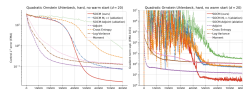


Figure 3: Plots of the  $L^2$  error (top), and the norm squared of the gradient with respect to the parameters  $\theta$  of the control (bottom), for the QG solver (Quadratic Ornstein-Uhlenbeck, hard) using and for each 1000 time. All the algorithms use a reparameterized control (see Appendix D).

# Conclusion

---

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).



# References

---