

# Initial Value Problem Enhanced Sampling for Closed-Loop Optimal Control Design with Deep Neural Networks

Jiequn Han

Flatiron Institute, Simons Foundation  
collaboration with Xuanxi Zhang, Jihao Long, Wei Hu, Weinan E

4th AFOSR Monterey Workshop on  
Computational Issues in Nonlinear Control  
May 2023

# Open-Loop and Closed-Loop Optimal Control

Consider the following optimal control problem:

$$\begin{aligned} \min_{\mathbf{x}(t), \mathbf{u}(t)} \int_0^T L(t, \mathbf{x}(t), \mathbf{u}(t)) dt + M(\mathbf{x}(T)) \\ \text{s.t. } \dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t)), \mathbf{x}(0) = \mathbf{x}, \end{aligned}$$

$\mathbf{x}(t)$  = state,  $\mathbf{u}(t)$  = control.

- Open-loop optimal control: find the optimal path  $(\mathbf{x}^*(t), \mathbf{u}^*(t))$  for a **specific initial point**.
- Closed-loop optimal control: find the optimal policy function  $\mathbf{u}^*(t, \mathbf{x})$ , applicable for **a set of initial points**  $\mathbf{x}(0) \in \mathcal{X}$ . More powerful, but more difficult to solve.
- Traditional methods by solving the associated Hamilton-Jacobi-Bellman equation suffers from **the curse of dimensionality**.

# Neural Network-Based Closed-Loop Control

Neural networks have demonstrated astonishing capability in dealing with high-dimensional functions.

How to use neural networks to design closed-loop optimal control for **high-dimensional problems**?

# Neural Network-Based Closed-Loop Control

Neural networks have demonstrated astonishing capability in dealing with high-dimensional functions.

How to use neural networks to design closed-loop optimal control for **high-dimensional problems**?

Approach 1: direct policy search (Han and E (2016), Böttcher, Antulov-Fantulin and Asikis (2022))

$$\begin{aligned} \min_{\theta} \mathbb{E} \int_0^T L(t, \mathbf{x}(t), \mathbf{u}^{\text{NN}}(t, \mathbf{x}(t); \theta)) dt + M(\mathbf{x}(T)) \\ \text{s.t. } \dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}^{\text{NN}}(t, \mathbf{x}(t); \theta)), \mathbf{x}(0) \sim \mu_0, \end{aligned}$$

# Neural Network-Based Closed-Loop Control

Neural networks have demonstrated astonishing capability in dealing with high-dimensional functions.

How to use neural networks to design closed-loop optimal control for **high-dimensional problems**?

Approach 1: direct policy search (Han and E (2016), Böttcher, Antulov-Fantulin and Asikis (2022))

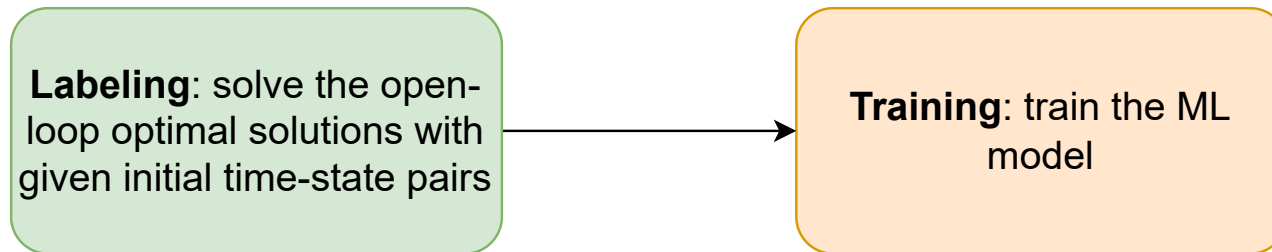
$$\begin{aligned} \min_{\theta} \mathbb{E} \int_0^T L(t, \mathbf{x}(t), \mathbf{u}^{\text{NN}}(t, \mathbf{x}(t); \theta)) dt + M(\mathbf{x}(T)) \\ \text{s.t. } \dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}^{\text{NN}}(t, \mathbf{x}(t); \theta)), \mathbf{x}(0) \sim \mu_0, \end{aligned}$$

Approach 2: supervised learning (Nakamura-Zimmerer, Gong, and Kang (2021))

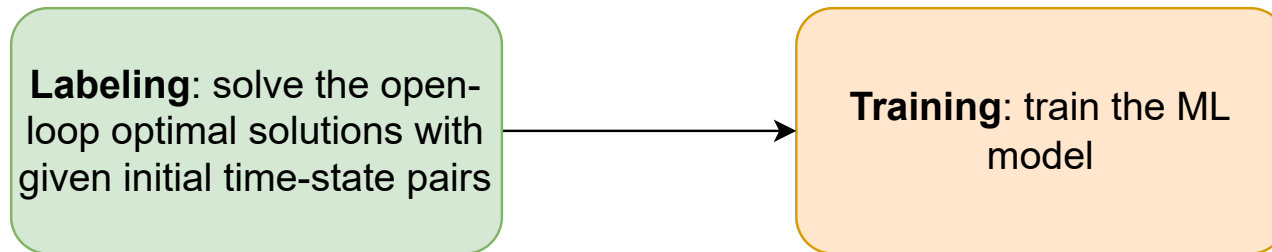
- Sample initial points and solve the corresponding open-loop optimal control problems.
- Choose the time-state-control tuples along every paths to collect the training data:  
 $\mathcal{D} = \{(t^i, \mathbf{x}^i), \mathbf{u}^i\}_{1 \leq i \leq M}$ .
- Train a neural network to approximate the closed-loop optimal control  $\mathbf{u}^{\text{NN}}$  :

$$\min_{\theta} \sum_{i=1}^M \|\mathbf{u}^i - \mathbf{u}^{\text{NN}}(t^i, \mathbf{x}^i; \theta)\|^2.$$

# Supervised-Learning-Based Approach

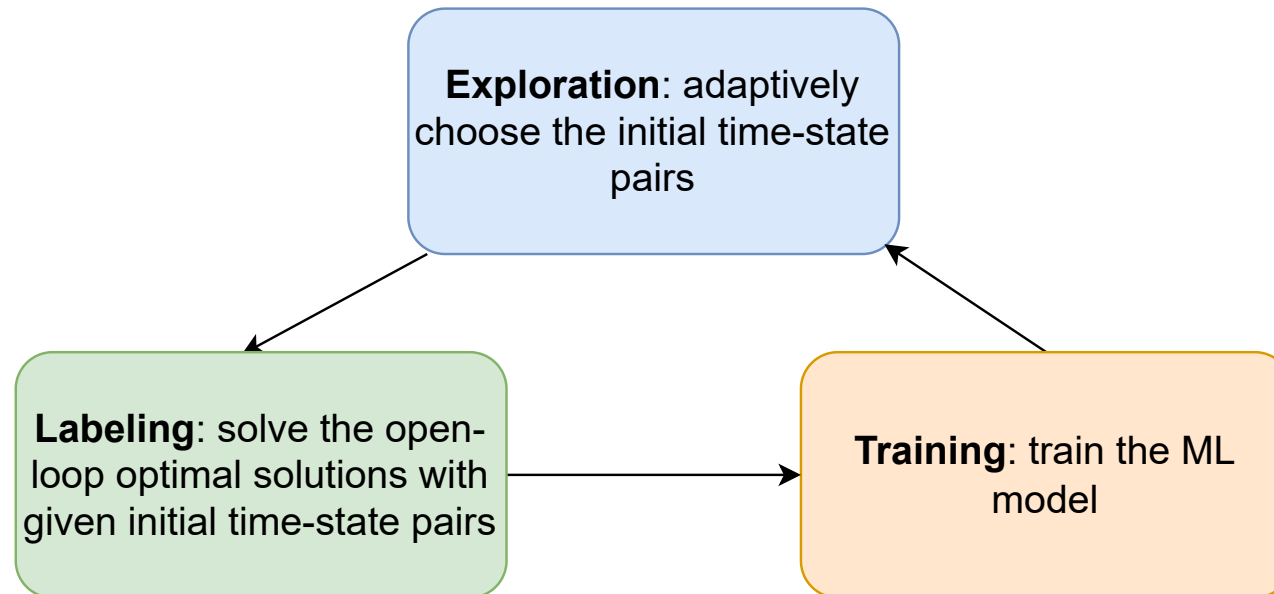


# Supervised-Learning-Based Approach



Unlike a lot of classical ML tasks, we can choose data to label by our own choice.

# Supervised-Learning-Based Approach



Unlike a lot of classical ML tasks, we can choose data to label by our own choice.

This work: focus on the exploration part to sample data that are more valuable to training



# Adaptive Sampling Methods

- It is observed that the NN controller trained by the SL-based approach does not perform well **even when both the training error and testing error are fairly small.**
- Adaptive sampling methods aim to find the difficult points the NN controller suffers and add these points into the training data to improve the performance.
- Most existing approaches focus on choosing the **initial points** such as choosing the initial points with large gradients or bad performance.

# Adaptive Sampling Methods

- It is observed that the NN controller trained by the SL-based approach does not perform well **even when both the training error and testing error are fairly small.**
- Adaptive sampling methods aim to find the difficult points the NN controller suffers and add these points into the training data to improve the performance.
- Most existing approaches focus on choosing the **initial points** such as choosing the initial points with large gradients or bad performance.
- However, these methods can not mitigate the **distribution mismatch phenomenon** brought by **controlled dynamics.**

# Distribution Mismatch Phenomenon

Define  $\mu_{\mathbf{u}}(t)$  the distribution of  $\mathbf{x}(t)$ :

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t, \mathbf{x}(t))), \quad \mathbf{x}(0) \sim \mu_0,$$

- $\mu_{\mathbf{u}^*}(t)$ : the distribution of the state at time  $t$  in the training data.
- $\mu_{\mathbf{u}^{\text{NN}}}(t)$ : the distribution of the input state of  $\mathbf{u}^{\text{NN}}$  at time  $t$  when applying the learned NN controller in the dynamics.

# Distribution Mismatch Phenomenon

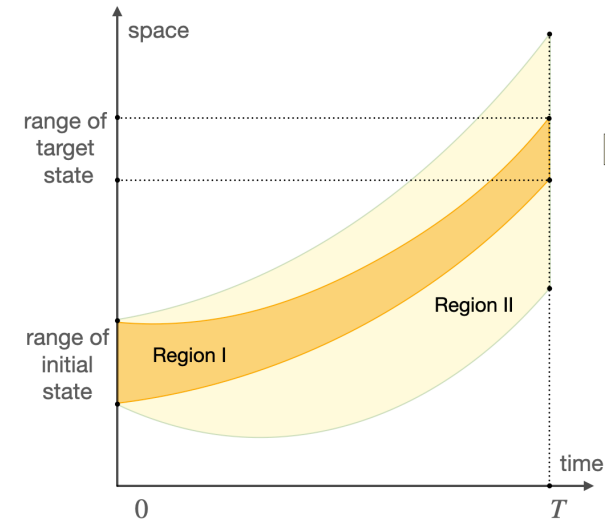
Define  $\mu_{\mathbf{u}}(t)$  the distribution of  $\mathbf{x}(t)$ :

$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t, \mathbf{x}(t))), \quad \mathbf{x}(0) \sim \mu_0,$$

- $\mu_{\mathbf{u}^*}(t)$ : the distribution of the state at time  $t$  in the training data.
- $\mu_{\mathbf{u}^{\text{NN}}}(t)$ : the distribution of the input state of  $\mathbf{u}^{\text{NN}}$  at time  $t$  when applying the learned NN controller in the dynamics.

Due to the learning error,  $\mu_{\mathbf{u}^*}(t) \neq \mu_{\mathbf{u}^{\text{NN}}}(t)$ , and **its discrepancy increases over time due to compounding error.**

When  $t$  is large, the training data fails to represent the states encountered when keeping applying  $\mathbf{u}^{\text{NN}}$ , and the error between  $\mathbf{u}^*$  and  $\mathbf{u}^{\text{NN}}$  dramatically increases.



Region I: optimal paths

Region II: Paths controlled by the NN controller trained on data in Region I

# Distribution Mismatch Phenomenon

Define  $\mu_{\mathbf{u}}(t)$  the distribution of  $\mathbf{x}(t)$ :

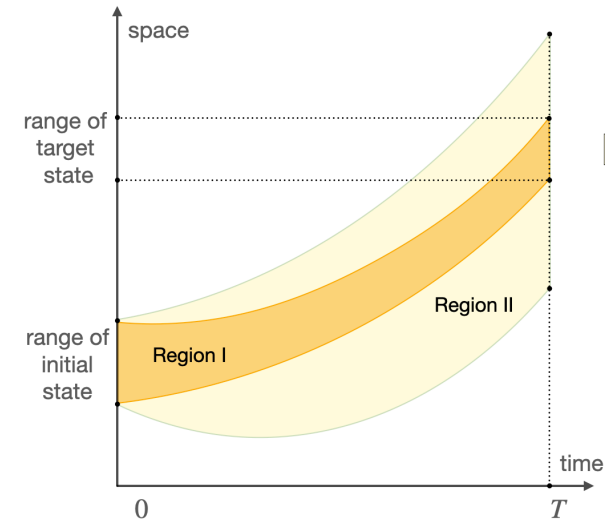
$$\dot{\mathbf{x}}(t) = \mathbf{f}(t, \mathbf{x}(t), \mathbf{u}(t, \mathbf{x}(t))), \quad \mathbf{x}(0) \sim \mu_0,$$

- $\mu_{\mathbf{u}^*}(t)$ : the distribution of the state at time  $t$  in the training data.
- $\mu_{\mathbf{u}^{\text{NN}}}(t)$ : the distribution of the input state of  $\mathbf{u}^{\text{NN}}$  at time  $t$  when applying the learned NN controller in the dynamics.

Due to the learning error,  $\mu_{\mathbf{u}^*}(t) \neq \mu_{\mathbf{u}^{\text{NN}}}(t)$ , and **its discrepancy increases over time due to compounding error.**

When  $t$  is large, the training data fails to represent the states encountered when keeping applying  $\mathbf{u}^{\text{NN}}$ , and the error between  $\mathbf{u}^*$  and  $\mathbf{u}^{\text{NN}}$  dramatically increases.

**distribution mismatch phenomenon** is common when involving machine learning and dynamical systems, such as reinforcement learning and imitation learning.



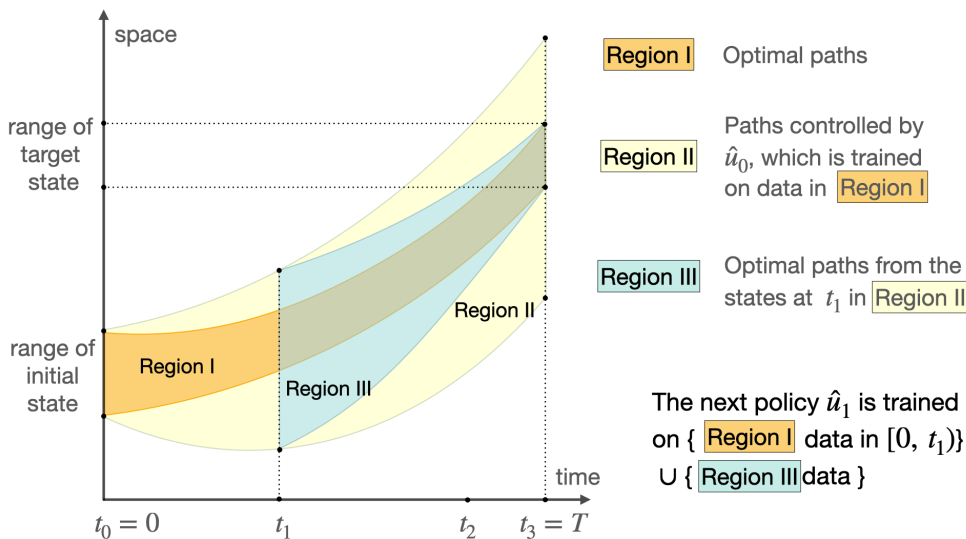
Region I: optimal paths

Region II: Paths controlled by the NN controller trained on data in Region I

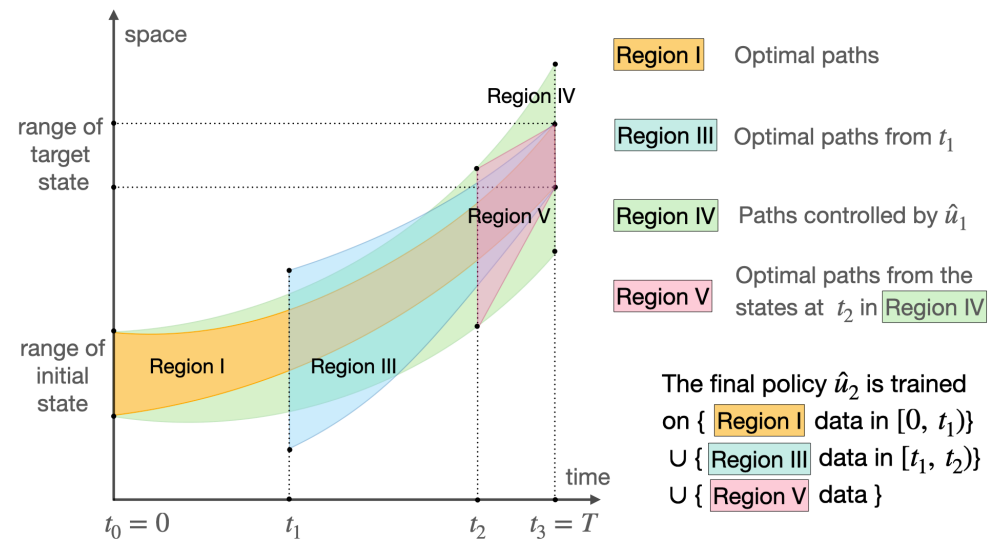
# IVP Enhanced Sampling Method

We propose **initial value problem (IVP) enhanced sampling method** to mitigate the distribution mismatch phenomenon.

**Key idea:** improve the quality of the NN controller iteratively by enlarging the training dataset with the states seen by the NN controller at previous times.

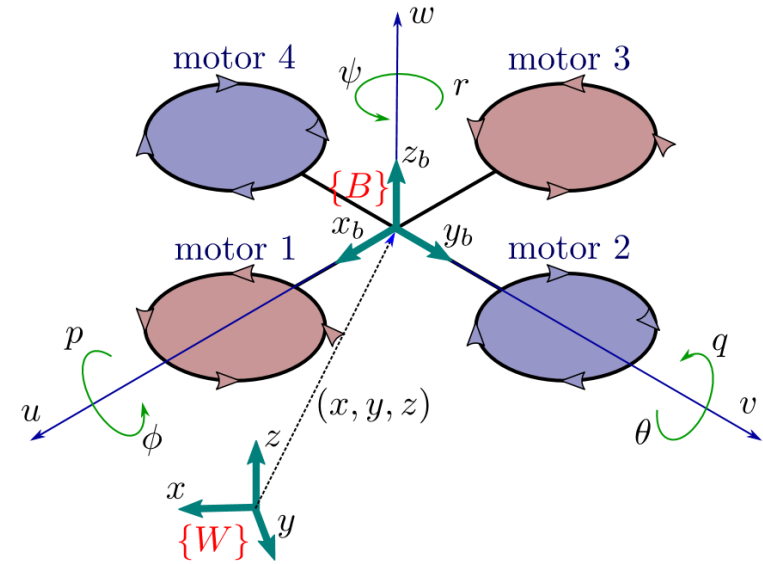


**Figure:** after first training



**Figure:** after second training

# The Optimal Landing Problem of Quadrotor

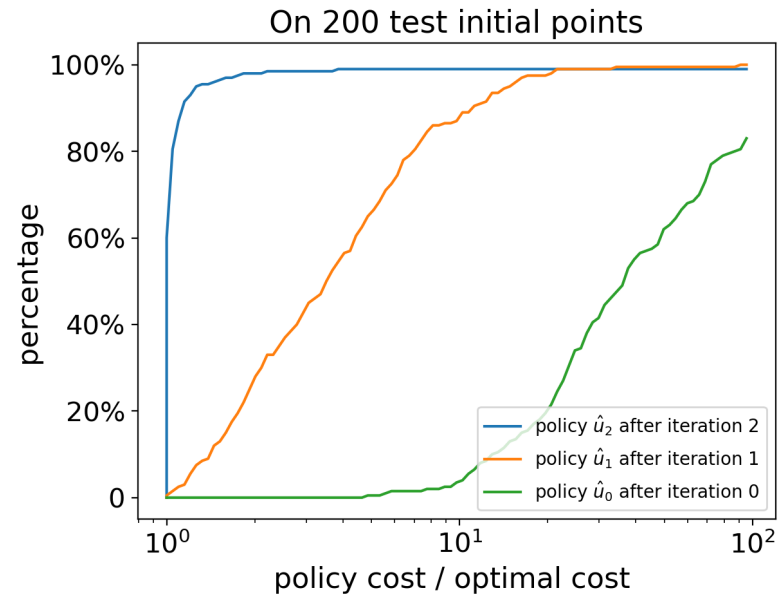


Full dynamic model of a quadrotor with 12-dimensional state and 4-dimensional control.

The goal is to find the optimal landing paths with the minimum control effort from 0 to  $T = 16$ .

The IVP enhanced sampling method is implemented on the time grid  $[0, 10, 14, 16]$ .

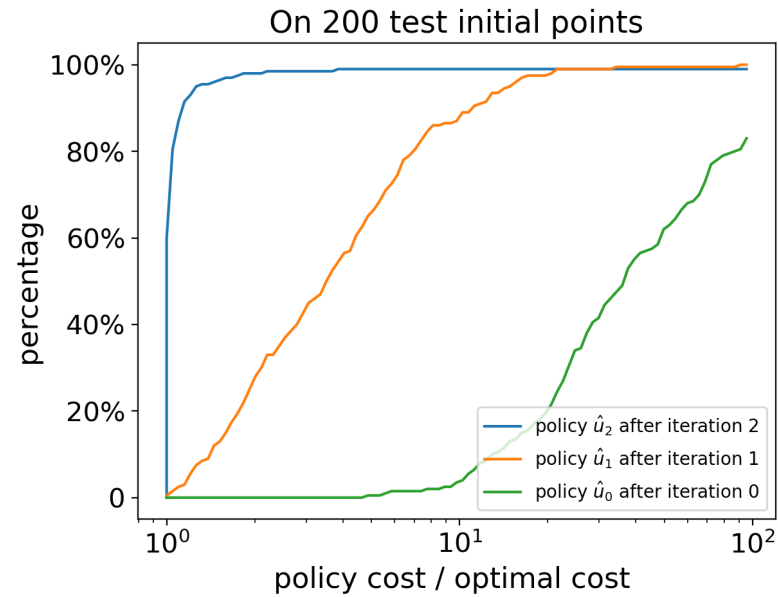
# Results on the Optimal Landing Problem



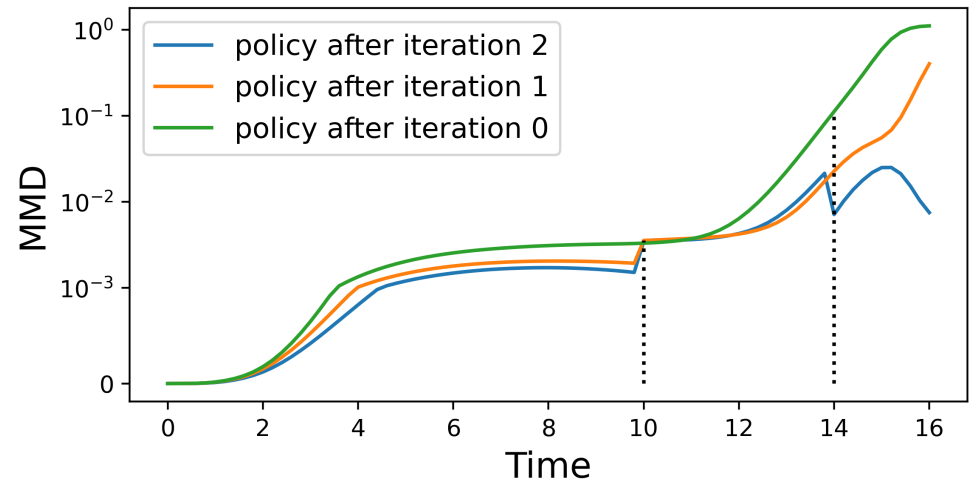
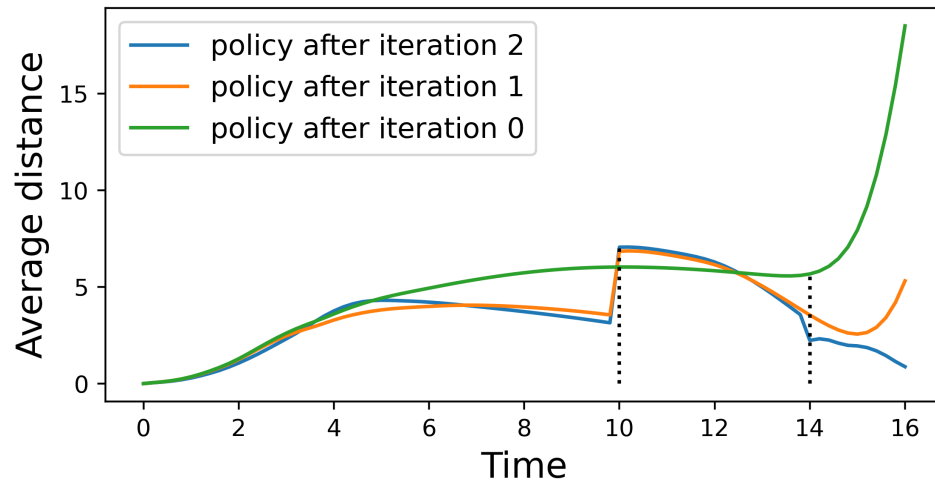
**Figure:** Cumulative distribution on 200 test initial points



# Results on the Optimal Landing Problem

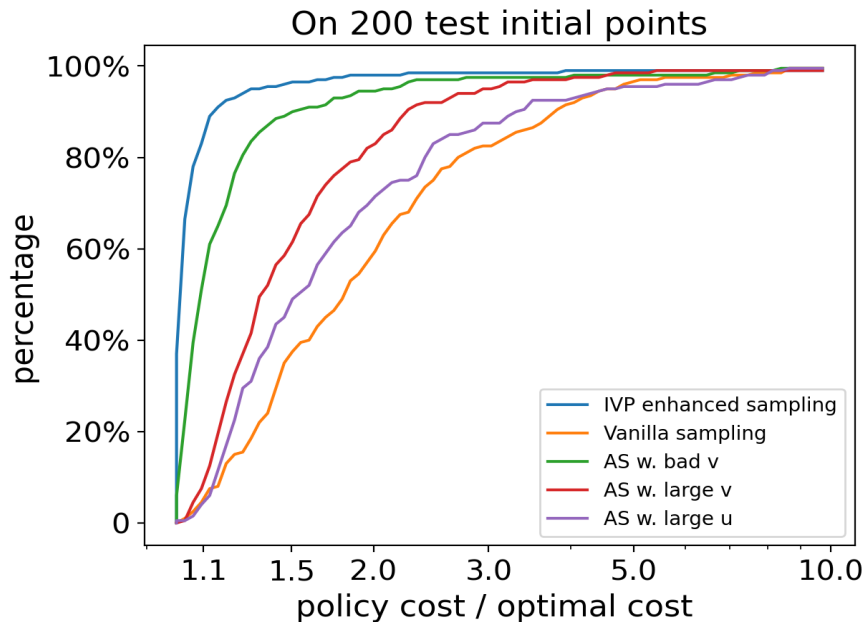


**Figure:** Cumulative distribution on 200 test initial points



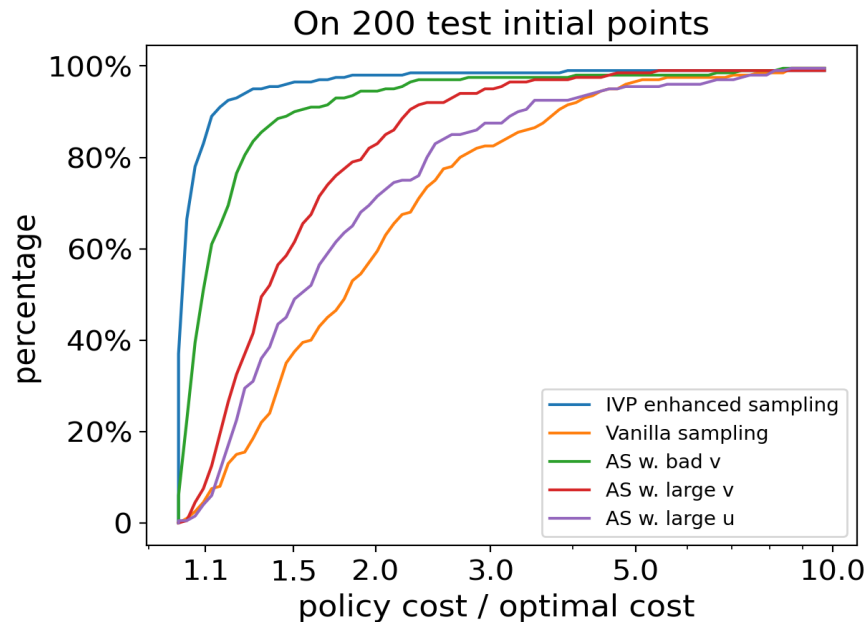
**Figure:** Discrepancy between the training data and the data reached by controllers at every time

# Comparison with Other Methods

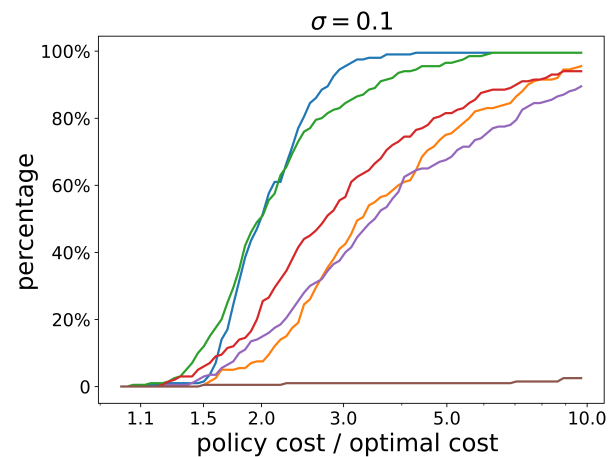
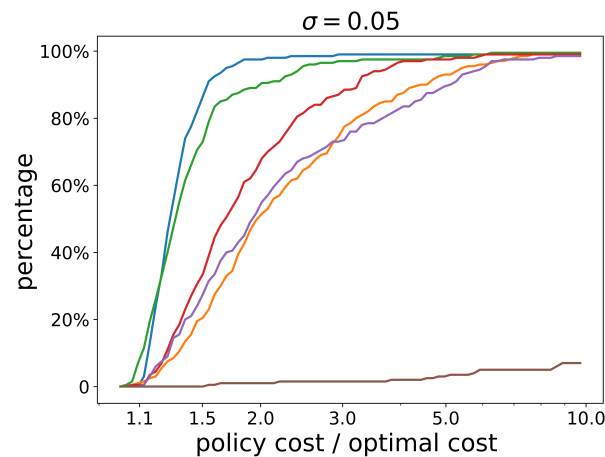
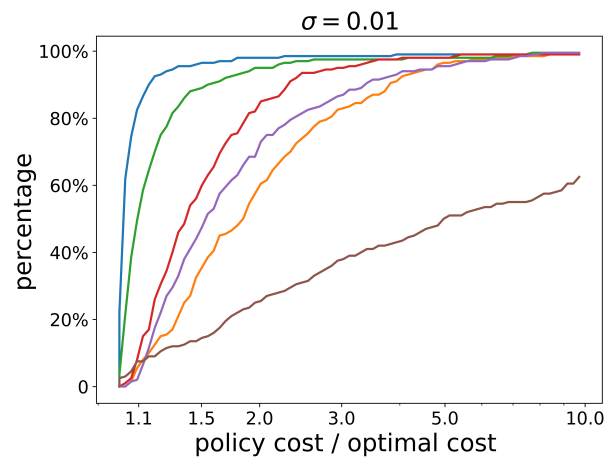


- Vanilla sampling: no adaptive
- AS w. bad v: choose initial points on which the NN controller performs bad
- AS w. large v: choose initial points on which the IVPs induced by the NN controller have large cost
- AS w. large u: choose initial points whose NN-predicted controls are large

# Comparison with Other Methods



- Vanilla sampling: no adaptive
- AS w. bad v: choose initial points on which the NN controller performs bad
- AS w. large v: choose initial points on which the IVPs induced by the NN controller have large cost
- AS w. large u: choose initial points whose NN-predicted controls are large



— IVP enhanced sampling — Vanilla sampling — AS w. bad v — AS w. large v — AS w. large u — Open-loop optimal controller

**Figure:** Cumulative distribution function of cost ratio between NN controlled value and optimal value under disturbance.

# Optimal Reaching Problem of a 7-DoF Manipulator

The reaching problem on a 7-DoF torque-controlled manipulator, the KUKA LWR iiwa R820 14 with 14-dimensional state and 7-dimensional control.

# Theoretical Analysis on an LQR Example

Consider the 1-d linear quadratic regulator (LQR) problem:

$$\begin{aligned} & \min_{x(t), u(t)} \frac{1}{T} \int_{t_0}^T |u(t)|^2 dt + |x(T)|^2 \\ \text{s.t.} \quad & \begin{cases} \dot{x}(t) = u(t), t \in [t_0, T] \\ x(t_0) = x_0 \end{cases} \end{aligned}$$

The optimal controls are

$$\begin{cases} u^*(t; t_0, x_0) = -\frac{T}{T(T-t_0)+1}x_0, & \text{(open-loop optimal control)} \\ u^*(t, x) = -\frac{T}{T(T-t)+1}x. & \text{(closed-loop optimal control)} \end{cases}$$

# Theoretical Analysis on an LQR Example

Consider the 1-d linear quadratic regulator (LQR) problem:

$$\begin{aligned} \min_{x(t), u(t)} \quad & \frac{1}{T} \int_{t_0}^T |u(t)|^2 dt + |x(T)|^2 \\ \text{s.t.} \quad & \begin{cases} \dot{x}(t) = u(t), t \in [t_0, T] \\ x(t_0) = x_0 \end{cases} \end{aligned}$$

The optimal controls are

$$\begin{cases} u^*(t; t_0, x_0) = -\frac{T}{T(T-t_0)+1}x_0, & \text{(open-loop optimal control)} \\ u^*(t, x) = -\frac{T}{T(T-t)+1}x. & \text{(closed-loop optimal control)} \end{cases}$$

Model 1:  $u_\theta(t, x) = -\frac{T}{T(T-t)+1}x + b(t)$ , where  $\theta = \{\theta_t\}_{0 \leq t \leq T} = \{b(t)\}_{0 \leq t \leq T}$ .

# Theoretical Analysis on an LQR Example

Consider the 1-d linear quadratic regulator (LQR) problem:

$$\begin{aligned} & \min_{x(t), u(t)} \frac{1}{T} \int_{t_0}^T |u(t)|^2 dt + |x(T)|^2 \\ \text{s.t.} \quad & \begin{cases} \dot{x}(t) = u(t), t \in [t_0, T] \\ x(t_0) = x_0 \end{cases} \end{aligned}$$

The optimal controls are

$$\begin{cases} u^*(t; t_0, x_0) = -\frac{T}{T(T-t_0)+1}x_0, & \text{(open-loop optimal control)} \\ u^*(t, x) = -\frac{T}{T(T-t)+1}x. & \text{(closed-loop optimal control)} \end{cases}$$

$$\text{Model 1: } u_\theta(t, x) = -\frac{T}{T(T-t)+1}x + b(t), \text{ where } \theta = \{\theta_t\}_{0 \leq t \leq T} = \{b(t)\}_{0 \leq t \leq T}.$$

Assume the open-loop optimal control solver gives the data with noise  $Z$

$$\begin{cases} \hat{u}(t; t_0, x_0) = -\frac{T}{T(T-t_0)+1}x_0 + \epsilon Z, \\ \hat{x}(t; t_0, x_0) = x_0 + \int_{t_0}^t \hat{u}(t; t_0, x_0) dt = \frac{T(T-t)+1}{T(T-t_0)+1}x_0 + (t - t_0)\epsilon Z. \end{cases}$$

# Theoretical Superiority of the IVP Enhanced Sampling

## Theorem

With Model 1, we do vanilla sampling with  $NT$  samples and IVP enhanced sampling with  $N$  initial points on temporal grids  $0 < 1 < 2 < \dots < T$ . Let  $u_o$ ,  $u_v$  and  $u_a$  be the optimal controller, the controller learned by the vanilla method, and the controller learned by the IVP enhanced sampling method, respectively. Define

$$\dot{x}_s(t) = u_s(t) = u_s(t, x_s(t)), \quad x_s(0) = x_{init}, \quad 0 \leq t \leq T, \quad s \in \{o, v, a\}.$$

1. If  $x_{init}$  is a standard normal random variable. Let  $\{\hat{x}_v^j(t)\}_{j=1}^{NT}$  and  $\{\hat{x}_a^j(t)\}_{j=1}^N$  be the state variables in the training data of the vanilla method and the last iteration of the IVP enhanced sampling method. Then,  $\hat{x}_v^j(t)$ ,  $\hat{x}_a^j(t)$ ,  $x_v(t)$  and  $x_a(t)$  are mean-zero normal random variables and

$$|\mathbb{E}|\hat{x}_v^j(t)|^2 - \mathbb{E}|x_v(t)|^2| = \left(1 - \frac{1}{NT}\right)\epsilon^2 t^2, \quad |\mathbb{E}|\hat{x}_a^j(t)|^2 - \mathbb{E}|x_a(t)|^2| \leq \epsilon^2.$$

2. If  $x_{init}$  is a fixed initial point, define the total cost

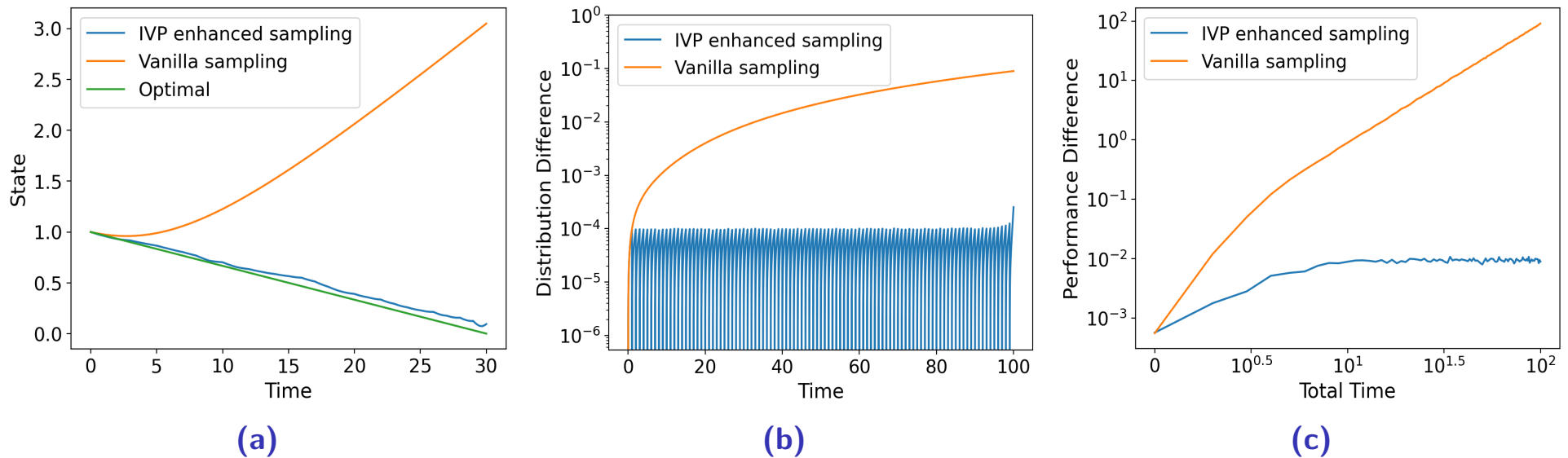
$$J_s = \frac{1}{T} \int_0^T |u_s(t)|^2 dt + |x_s(T)|^2, \quad s \in \{o, v, a\}.$$

$$\text{Then, } \mathbb{E}J_v - J_o = \frac{T^2 + 1}{NT}\epsilon^2, \quad \mathbb{E}J_a - J_o \leq \frac{3\epsilon^2}{N}.$$



# Superiority of the IVP Enhanced Sampling

Model 2:  $u_\theta(t, x) = a(t)x + b(t)$ , where  $\theta = \{\theta_t\}_{0 \leq t \leq T} = \{(a(t), b(t))\}_{0 \leq t \leq T}$ ,



**Figure:** Numerical results on learning Model 2. (a) The optimal path and the paths generated by the vanilla sampling method and the IVP enhanced sampling method. (b) Differences of the second order moments (in the logarithm scale) between the distributions of the training data and the data reached by the controllers at different times. (c) Performance differences (in the logarithm scale) of the vanilla sampling method and the IVP enhanced sampling method for different total times (in the logarithm scale).

# Conclusions

- Traditional methods for the closed-loop optimal control suffer from the curse of dimensionality while deep learning is promising in high dimensional closed-loop optimal control problems.
- Distribution mismatch phenomenon is an essential challenge in the supervised-learning-based approach for optimal control
- IVP enhanced sampling method can mitigate the distribution mismatch phenomenon and significantly improve the performance of the NN controller.

Thank you for your attention!