

# Online Control-Informed Learning

Anonymous authors

Paper under double-blind review

## Abstract

This paper proposes an Online Control-Informed Learning (OCIL) framework, which synthesizes the well-established control theories to solve a broad class of learning and control tasks in real time. This novel integration effectively handles practical issues in machine learning such as noisy measurement data, online learning, and data efficiency. By considering any robot as a tunable optimal control system, we propose an online parameter estimator based on extended Kalman filter (EKF) to incrementally tune the system in real time, enabling it to complete designated learning or control tasks. The proposed method also improves robustness in learning by effectively managing noise in the data. Theoretical analysis is provided to demonstrate the convergence and regret of OCIL. Three learning modes of OCIL, i.e. Online Imitation Learning, Online System Identification, and Policy Tuning On-the-fly, are investigated via experiments, which validate their effectiveness.

## 1 Introduction

Informed Machine Learning (IML) (Von Rueden et al., 2021) represents an emerging approach integrating prior knowledge into the machine learning (ML) process. While classic classification tasks in unsupervised, semi-supervised, or supervised ML primarily focus on extracting patterns from labeled or unlabeled data (LeCun et al., 2015), IML leverages prior knowledge such as physical laws, expert knowledge, or existing models to uncover underlying connections within data (Karniadakis et al., 2021). This integration enables models to produce more reliable predictions with enhanced interpretability, particularly in scenarios where data is scarce, noisy, or complex. This approach is especially advantageous in the domains where theoretical understanding is well-established and thus can guide ML. One notable example of IML is physics-informed machine learning (Karniadakis et al., 2021), which proves particularly valuable in solving partial differential equations for computational fluid dynamics.

Control-informed learning (CIL) is a subset of IML tailored for system control, autonomy, and robotics. This approach merges standard control theory, especially the established optimal control principles, with ML techniques to enhance the functionality of autonomous systems. This integration is beneficial for robotic applications, where optimal control systems typically govern robots. By incorporating prior knowledge about the system, such as dynamic models, control laws, and optimization strategies, into the learning process, CIL reduces the required training data, accelerates deployment, and enhances safety and efficiency by ensuring adherence to established control theories (Jin et al., 2020; 2021b).

Online control-informed learning (OCIL) represents a subset of CIL that continually updates its model as new data streams. This is beneficial for robotic applications, where the environment can change unpredictably. One advantage of online methods is their ability to learn and adapt on the fly. Unlike offline methods, which require a static dataset for training, online learning algorithms adjust their parameters or models in response to new information without training from scratch. This continual learning capability allows autonomous systems to improve their policies in real time. Additionally, OCIL can handle dynamic environments by continually updating its parameters. This feature is critical in applications such as autonomous driving, where adapting to new scenarios quickly is crucial for safety and performance. This concept intersects with many ML methods such as transfer learning (Pan & Yang, 2009), continual learning (Aljundi et al., 2019), and learning on-the-fly (Ornik et al., 2019).

## 1.1 Related Work

The proposed OCIL framework includes three modes for three sets of classic problems in both ML and system control communities. Hence, the following literature review presents some representative methods from both perspectives, for all three modes.

**Online System Identification** To identify a nonlinear system with possibly noisy measurement in an online fashion, Markov-decision-process-based methods are widely used, such as linear regression (Haruno et al., 2001), observation-transition modeling (Finn et al., 2016), latent-space modeling (Watter et al., 2015), (deep) neural networks (Fragkiadaki et al., 2015), Gaussian processes (Deisenroth & Rasmussen, 2011), and transition graphs (Zhang et al., 2018). Despite their widespread use, these methods often must balance data efficiency against prediction accuracy. To improve both metrics, physics-informed learning approaches Raissi et al. (2019); Saemundsson et al. (2020); Lutter et al. (2019); Zhong et al. (2019) incorporate physical laws into learning models. Recently, there has been an emerging trend to interpret (deep) neural networks through the lens of dynamical systems, leading to the development of several new algorithms (Chen et al., 2018; Han & E, 2016; Li et al., 2018; Li & Hao, 2018; Han et al., 2019; Zhang et al., 2019; Benning et al., 2019; Liu & Markowich, 2020). Koopman operator theory offers a method to lift states to an infinite-dimensional linear observable space (Mauroy et al., 2020; Williams et al., 2015). The recent deep Koopman representation (DKR) utilizes neural networks (NNs) to represent the observables that are difficult to formulate by hand without expertise in a particular dynamical system (Liang et al., 2023; Hao et al., 2023). This paper proposes an online learning framework to estimate system dynamics in real time as noisy measurement data continually comes in. We aim to incorporate inductive knowledge from optimal control theory to enhance both run-time computational efficiency, robustness against measurement noise, and prediction accuracy.

**Online Imitation Learning** Online imitation and objective learning are typically referred to as inverse reinforcement learning (IRL) in the ML community and online inverse optimal control (IOC) in the system control community. IRL aims to deduce a control objective function with observed optimal demonstrations. The objective function is generally represented as a weighted sum of features (Abbeel & Ng, 2004; Ziebart et al., 2008; Ratliff et al., 2006). Approaches to find these unknown weights include feature matching (Abbeel & Ng, 2004), maximum entropy (Ziebart et al., 2008), and maximum margin (Ratliff et al., 2006). These IRL methods update the learned objective estimate within a defined feature space, exploiting the linearity of feature weights. As for learning nonlinear parameter mapping of objective functions, prior and system-dependent knowledge is required to further extend the methods above. On the other hand, with system dynamics, IOC aims for efficient learning approaches (Keshavarz et al., 2011; Mombaur et al., 2010; Liang et al., 2022; 2023; Jin et al., 2019; 2021a; Jin & Mou, 2021). For instance, some methods (Keshavarz et al., 2011; Liang et al., 2022; 2023; Jin et al., 2019; 2021a; Jin & Mou, 2021) directly calculate unknown weights by minimizing the violation of optimality conditions by the observed demonstration data, which avoid solving optimal control problems repetitively. This paper presents an online framework to facilitate learning objective functions online as noisy demonstration data continually streams in. By synthesizing the online state estimation techniques in control theory, the proposed method can recover the object function in real time while ensuring robustness against noise.

**Tuning Policy On-the-fly** Tuning policy on the fly is typically referred to as transfer learning and tuning optimal control (OC) systems in ML and control communities, respectively. Transfer learning exploits the generalization of existing knowledge such that it can be transferred across different domains (Taylor & Stone, 2009). Recently, transfer learning has been investigated in the reinforcement learning (RL) community because the knowledge gained in one task may improve learning performance in a related but different task. Such a concept of transfer learning may speed up the learning process in RL (Taylor & Stone, 2009). If the knowledge transfers from a domain of experts to a target domain, imitation learning (Ho & Ermon, 2016; Wu et al., 2019; Oh et al., 2018), behavior cloning (Torabi et al., 2018; Sasaki & Yamashina, 2021), or learning from demonstrations (Schaal, 1996; Kim et al., 2013; Jin et al., 2022a) can also be viewed as transfer learning. Knowledge from a source domain can also be represented by an action probability distribution of an expert policy (Czarnecki et al., 2019) or even directional corrections from expert (Jin et al., 2022b). In the control community, tuning OC systems initially refers to neighboring extremal optimal control (NEOC) (Bryson, 1975; Ghaemi et al., 2009). Recently, Jin et al. (2020) proposes a framework to tune an OC system based on differentiating Pontryagin’s Maximum Principle. Tuning an OC system by either NEOC or differentiating

Pontryagin’s Principle requires the knowledge of an entire optimal trajectory. However, for online tasks, the knowledge of the entire optimal trajectory is unavailable at run-time until the end. The proposed online learning framework in this paper includes a specific mode for model-based control tasks, where some information is unavailable until the current time.

## 1.2 Contributions

This paper introduces an online learning framework called Online Control-Informed Learning (OCIL). This framework is designed to be data efficient for various learning and control tasks while providing robustness against noisy data. In this paper, we consider a robot as an OC system, which is parameterized by tunable parameters within different components of the system, including dynamics, policy, and objective function. By tuning the OC system in real time, the proposed OCIL tackles three learning tasks in robotics, namely Online Imitation Learning, Online System Identification, and Policy Tuning On-the-fly. The proposed OCIL consists of two main components, both of which are inspired by control theory. Specifically, the framework first proposes an online parameter estimator based on the classic online state estimation techniques in control theory. The estimator continually updates the parameter estimates in real time as new data becomes available, aiming to minimize a cumulative loss defined for a specific task. To do so, the gradient information for the loss with respect to the tunable parameter is required. Therefore, OCIL employs a gradient generator (GG) based on Pontryagin Differential Programming in optimal control theory to calculate the exact gradient.

**Notations.**  $\|\cdot\|$  denotes the Euclidean norm. Given a matrix  $A \in \mathbb{R}^{n \times m}$ , let  $A'$  denotes its transpose. For positive integers  $n$  and  $m$ , let  $\mathbf{I}_n$  be the  $n \times n$  identity matrix;  $\mathbf{0}_n \in \mathbb{R}^n$  denotes a vector with all value 0;  $\mathbf{0}_{n \times m}$  denotes a  $n \times m$  matrix with all value 0. Let  $\text{col}\{\mathbf{v}_1, \dots, \mathbf{v}_a\}$  denote a column stack of elements  $\mathbf{v}_1, \dots, \mathbf{v}_a$ , which may be scalars, vectors or matrices, i.e.  $\text{col}\{\mathbf{v}_1, \dots, \mathbf{v}_a\} \triangleq [\mathbf{v}'_1 \dots \mathbf{v}'_a]$ .

## 2 Problem Formulation

Consider the following class of optimal control (OC) systems  $\Sigma(\boldsymbol{\theta})$ :

$$\begin{aligned} \text{system dynamics:} \quad & \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\theta}), \quad \text{with } \mathbf{x}_0 \text{ given,} \\ \text{objective to be minimized:} \quad & J(\boldsymbol{\theta}) = \sum_{t=0}^{T-1} c(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\theta}) + h(\mathbf{x}_T, \boldsymbol{\theta}), \end{aligned} \quad (1)$$

where  $t = 0, 1, 2, \dots, T$  is the time index with  $T$  being the final time;  $\mathbf{x}_t \in \mathbb{R}^n$  and  $\mathbf{u}_t \in \mathbb{R}^m$  denote the system state and control input, respectively;  $\boldsymbol{\theta} \in \mathbb{R}^p$  denotes a tunable parameter;  $\mathbf{f} : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \rightarrow \mathbb{R}^n$  denotes a twice-differentiable time-invariant system dynamics;  $c : \mathbb{R}^n \times \mathbb{R}^m \times \mathbb{R}^p \mapsto \mathbb{R}$  and  $h : \mathbb{R}^n \times \mathbb{R}^p \mapsto \mathbb{R}$  denote running cost the final cost, respectively, both of which are assumed to be twice-differentiable. For a choice of  $\boldsymbol{\theta}$ , the trajectory of optimal control system (1) can be determined by solving the optimal control problem:

$$\begin{aligned} \{\mathbf{x}_{0:T}(\boldsymbol{\theta}), \mathbf{u}_{0:T-1}(\boldsymbol{\theta})\} \in \arg \min_{\substack{\mathbf{x}_{0:T} \\ \mathbf{u}_{0:T-1}}} J(\boldsymbol{\theta}) \\ \text{subject to} \quad & \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\theta}), \quad \forall t \text{ with given } \mathbf{x}_0. \end{aligned} \quad (2)$$

where  $\mathbf{x}_{0:T}(\boldsymbol{\theta}) \triangleq \text{col}\{\mathbf{x}_0(\boldsymbol{\theta}), \dots, \mathbf{x}_T(\boldsymbol{\theta})\}$  and  $\mathbf{u}_{0:T-1}(\boldsymbol{\theta}) \triangleq \text{col}\{\mathbf{u}_0(\boldsymbol{\theta}), \dots, \mathbf{u}_{T-1}(\boldsymbol{\theta})\}$  denote the states and inputs trajectory given parameter  $\boldsymbol{\theta}$ , respectively. For notation simplicity, we define the trajectory of the optimal control system as  $\boldsymbol{\xi}(\boldsymbol{\theta}) \triangleq \text{col}\{\mathbf{x}_{0:T}(\boldsymbol{\theta}), \mathbf{u}_{0:T-1}(\boldsymbol{\theta})\} \in \mathbb{R}^{(T+1)n+Tm}$ . Since our goal is to develop an online update rule, we introduce the iteration index  $k$ , where  $\boldsymbol{\theta}_k \in \mathbb{R}^p$  represents the tunable parameter at iteration  $k$ . However, since the parameter is updated at each time  $t$  in an online algorithm,  $k$  is equivalent to  $t$ , leading to  $\boldsymbol{\theta}_k \equiv \boldsymbol{\theta}_t \in \mathbb{R}^p$ . At each time  $t$ , given a slice of trajectory  $\boldsymbol{\xi}(\boldsymbol{\theta}_t)$  at time  $t$ , denoted as  $\boldsymbol{\xi}_t(\boldsymbol{\theta}_t) \triangleq \text{col}\{\mathbf{x}_t(\boldsymbol{\theta}_t), \mathbf{u}_t(\boldsymbol{\theta}_t)\} \in \mathbb{R}^{n+m}$ . There is also information  $\mathbf{O}_t \in \mathbb{R}^r$ , which is used to evaluate  $\boldsymbol{\xi}_t(\boldsymbol{\theta}_t)$  and remains inaccessible until time  $t$ . The information  $\mathbf{O}_t$  could represent the desired state, control input, or reference measured output. Examples of  $\mathbf{O}_t$  will be provided at the end of this section. We assume that the information  $\mathbf{O}_t$  is subject to Gaussian measurement noise, denoted as  $\mathbf{v}_t \sim \mathcal{N}(\mathbf{0}_r, \mathbf{R}_t)$ , where  $\mathbf{R}_t \in \mathbb{R}^{r \times r}$  represents the covariance matrix of the noise. In other words, at each time step  $t$ , the information received can be noisy and expressed as  $\mathbf{O}_t = \mathbf{O}_t^* + \mathbf{v}_t$ , where  $\mathbf{O}_t^*$  represents the true but unknown information.

The *problem of interest* is to develop an online method to update  $\theta$  at every time  $t$ , such that its trajectory  $\xi(\theta)$  from (2) minimizes a task-specific cumulative loss  $L(\xi(\theta))$ :

$$\min_{\theta} L(\xi(\theta)) \quad \text{subject to } \xi(\theta) \text{ is in 2.} \quad (3)$$

To define  $L(\xi(\theta))$ , we consider a predefined differentiable stage loss function  $l(\xi_t(\theta), \mathbf{O}_t) : \mathbb{R}^{n+m} \times \mathbb{R}^r \rightarrow \mathbb{R}^r$ . Then, the performance of the entire trajectory can be evaluated by the cumulative loss:

$$L(\xi(\theta)) = \sum_{t=0}^T \|l(\xi_t(\theta), \mathbf{O}_t)\|^2. \quad (4)$$

To achieve a specific learning or control task, one needs to design a specific stage loss function  $l(\xi_t(\theta), \mathbf{O}_t)$  and select the most suitable information  $\mathbf{O}_t$ . Below, we will present several examples to illustrate how this can be done. It is worth noting that in different applications, adjustments to the configuration of system  $\Sigma(\theta)$  are required according to the task.

**Online SysID:** For a SysID problem, the goal is to identify the dynamics model of a physical system from the state-input trajectory  $\xi^o = \{\mathbf{x}_{0:T}^o, \mathbf{u}_{0:T-1}\}$ , where the superscript  $o$  denotes the observed trajectory. The trajectory is often generated by persistent excitation of the system without considering any control objectives (Keesman, 2011). Therefore, we can set  $J(\theta) = 0$ :

$$\Sigma(\theta) : \begin{array}{ll} \text{dynamics:} & \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \theta), \quad \text{with } \mathbf{x}_0 \text{ given,} \\ \text{objective:} & J(\theta) = 0. \end{array} \quad (5)$$

To identify the model dynamics, namely finding the  $\theta$  in the dynamics  $\mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \theta)$ , one could design the stage loss function to represent the discrepancy between the observed trajectory and the trajectory produced by  $\theta$ , i.e.  $l(\xi_t(\theta), \xi_t^o) = \xi_t^o - \xi_t(\theta)$ , where  $\xi_t^o$  is a slice of  $\xi^o$  at time  $t$ . In the SysID mode, the information  $\mathbf{O}_t$  received at time  $t$  is a slice of the trajectory of a physical system  $\xi_t^o$ .

**Online Imitation Learning:** The objective function and the model dynamics are parameterized by an unknown  $\theta$ . The OC system follows (1). Suppose one can observe the measurement of the optimal demonstration  $\mathbf{y}_t^*$  at each time  $t$ , where  $\mathbf{y}_t = \mathbf{g}(\mathbf{x}_t(\theta), \mathbf{u}_t(\theta))$  is a known differentiable mapping from the state and input to a measured output. Then, the stage loss function can be designed as  $l(\xi_t(\theta), \mathbf{y}_t^*) = \mathbf{y}_t^* - \mathbf{g}(\mathbf{x}_t(\theta), \mathbf{u}_t(\theta))$ . In this case, the information  $\mathbf{O}_t$  received at time  $t$  is the measurement of the optimal demonstration  $\mathbf{y}_t^*$ . The optimal demonstration can vary between being continuous or sparse, depending on practical application scenarios.

**Tuning Policy On-the-fly:** For a specific system, we want to obtain a feedback controller such that the trajectory minimizes certain task loss. We consider a feedback controller which is parameterized by  $\theta$ , i.e.  $\mathbf{u}_t = \mu(\mathbf{x}_t, \theta)$ . Then the OC system is written as follows:

$$\Sigma(\theta) : \begin{array}{ll} \text{dynamics:} & \mathbf{x}_{t+1} = \mathbf{f}(\mathbf{x}_t, \mu(\mathbf{x}_t, \theta)), \quad \text{with } \mathbf{x}_0 \text{ given,} \\ \text{objective:} & J = \sum_{t=0}^{T-1} c(\mathbf{x}_t, \mu(\mathbf{x}_t, \theta)) + h(\mathbf{x}_T). \end{array} \quad (6)$$

Then we can introduce any loss function to be minimized, which represents a specific tuning objective or specification, such as stabilizing the system on a particular point or tracking a desired trajectory.

### 3 Main Results

The proposed OCIL consists of two main components, both of which are inspired by control theory. Specifically, OCIL first proposes an online parameter estimator based on the extended Kalman filter (EKF). Going forward, we will show the challenge of obtaining the Kalman gain. To tackle this challenge, the gradient information for the loss with respect to the tunable parameter is required. Therefore, OCIL employs a gradient generator (GG) based on Pontryagin Differential Programming to calculate the exact gradient. Then the proposed OCIL framework will be introduced and supported with theoretical analysis.

### 3.1 Online Parameter Estimator

To minimize the cumulative task loss  $L(\boldsymbol{\xi}(\boldsymbol{\theta}))$  with information  $\mathbf{O}_t$ , which is unavailable until time  $t$ , the optimization problem that needs to be solved in real time is:

$$\min_{\boldsymbol{\theta}} \sum_{t=0}^T \|\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)\|^2 \quad \text{subject to } \boldsymbol{\xi}(\boldsymbol{\theta}) \text{ is in } \mathcal{Z}. \quad (7)$$

The optimization problem (7) is essentially a least squares problem, although under constraints. One of the most famous methods to solve the least squares problems incrementally is the EKF (Bertsekas, 1996). The EKF was proposed to incrementally estimate the state of a system using measured output available at each time step. In our problem setting, instead of estimating the state of a system, our goal is to estimate the parameter  $\boldsymbol{\theta}$  by utilizing the information  $\mathbf{O}_t$  that is available at each time  $t$ . Therefore, by considering the tunable parameter  $\boldsymbol{\theta}$  as the state to be estimated, its corresponding system can be written as:

$$\text{dynamics: } \boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \mathbf{w}_t, \quad \text{measurement: } \mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}_t), \mathbf{O}_t) + \mathbf{v}_t \equiv \mathbf{0}_s \quad (8)$$

where  $\mathbf{w}_t \sim \mathcal{N}(\mathbf{0}_p, \mathbf{Q}_t)$  is the process noise which is assumed to be multivariate Gaussian noise.  $\mathbf{Q}_t \in \mathbb{R}^{p \times p}$  is the covariance matrices of the process noise. In our experiments, we omit process noise since we are estimating a static parameter rather than the states of a dynamic system. Process noise usually accounts for dynamic uncertainty in such systems, which does not apply to our scenario. Nevertheless, we include the process noise here to support the convergence analysis. The measurement equation in (8) indicates that the stage loss with optimal  $\boldsymbol{\theta}$  is assumed to be zero. The  $\boldsymbol{\theta}$  estimate via EKF can be done as follows:

$$\hat{\boldsymbol{\theta}}_t^- := \hat{\boldsymbol{\theta}}_{t-1}, \quad \mathbf{P}_t^- := \mathbf{P}_{t-1} + \mathbf{Q}_{t-1} \quad (9a)$$

$$\mathbf{K}_t := \mathbf{P}_t^- \mathbf{L}_t' (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}_t' + \mathbf{R}_t)^{-1}, \quad \mathbf{P}_t := (\mathbf{I}_p - \mathbf{K}_t \mathbf{L}_t) \mathbf{P}_t^-, \quad \hat{\boldsymbol{\theta}}_t := \hat{\boldsymbol{\theta}}_t^- + \mathbf{K}_t (\mathbf{0}_s - \mathbf{l}(*)), \quad (9b)$$

where  $\mathbf{l}(\ast) \triangleq \mathbf{l}(\boldsymbol{\xi}_t(\hat{\boldsymbol{\theta}}_t^-), \mathbf{O}_t)$ ;  $\mathbf{L}_t \triangleq \left. \frac{d\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)}{d\boldsymbol{\theta}} \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_t^-} \in \mathbb{R}^{r \times p}$ ; (9a) predicts the dynamics; (9b) updates the parameter estimate. Here,  $\hat{\boldsymbol{\theta}} \in \mathbb{R}^p$  is the estimation of parameter  $\boldsymbol{\theta}$ , the superscript  $-$  means the term is not yet updated by measurement residual;  $\mathbf{P}_t \in \mathbb{R}^{p \times p}$  is a positive-definite matrix that denotes the covariance of the estimate;  $\mathbf{K}_t \in \mathbb{R}^{p \times r}$  denotes the Kalman gain. Throughout the estimation process, all of the terms are known except  $\mathbf{L}_t$ . It is challenging to obtain this term as the stage loss  $\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)$  is not an explicit function of  $\boldsymbol{\theta}$ . In the next subsection, we will present a *gradient generator* which computes the exact value for  $\mathbf{L}_t$ .

### 3.2 Gradient Generator

In this section, for brevity, the stage loss  $\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)$  is written as  $\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}))$  since the gradient is not related to  $\mathbf{O}_t$ . Additionally, for the remainder of the paper, to maintain brevity, the notation  $\left. \cdot \right|_{\boldsymbol{\theta}=\hat{\boldsymbol{\theta}}_t^-}$  will be omitted for all the partial derivatives that involve  $\boldsymbol{\theta}$ . To obtain the gradient  $\frac{d\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}))}{d\boldsymbol{\theta}}$ , one can employ the chain rule by definition,

$$\frac{d\mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}))}{d\boldsymbol{\theta}} = \frac{\partial \mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}))}{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})} \frac{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \quad (10)$$

where  $\frac{\partial \mathbf{l}(\boldsymbol{\xi}_t(\boldsymbol{\theta}))}{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}$  is known since the stage loss is pre-designed. The challenge that remains is to find the partial derivative  $\frac{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ , i.e. an analytical relation between trajectory  $\boldsymbol{\xi}_t$  and the tunable parameter  $\boldsymbol{\theta}$ . To tackle this challenge, the gradient generator in Jin et al. (2020) is used to obtain the exact value of  $\frac{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ .

Given the OC system (2), one can obtain the Hamiltonian equation

$$H_t = c(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\theta}) + \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t, \boldsymbol{\theta})' \boldsymbol{\lambda}_{t+1} \quad (11)$$

for all  $t = 0, \dots, T-1$ , where  $\boldsymbol{\lambda}_t \in \mathbb{R}^n$  denotes the Lagrangian multiplier associated with the equality constraint of model dynamics. With the definition of  $\boldsymbol{\xi}(\boldsymbol{\theta})$ , one has  $\frac{\partial \boldsymbol{\xi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \text{col}\left\{ \frac{\partial \mathbf{x}_{0:T}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}, \frac{\partial \mathbf{u}_{0:T-1}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \right\}$ . By

defining

$$X_t \triangleq \frac{\partial \mathbf{x}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \in \mathbb{R}^{n \times p}, \quad U_t \triangleq \frac{\partial \mathbf{u}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} \in \mathbb{R}^{m \times p}, \quad (12)$$

one can utilize the following lemma from Jin et al. (2020) to obtain the partial derivatives  $\frac{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$ :

**Lemma 1.** *Jin et al. (2020)* By defining the Jacobian and Hessian matrices related to  $\boldsymbol{\xi}(\boldsymbol{\theta})$ :

$$\begin{aligned} \mathbf{F}_t &= \frac{\partial \mathbf{f}}{\partial \mathbf{x}_t}, \quad \mathbf{G}_t = \frac{\partial \mathbf{f}}{\partial \mathbf{u}_t}, \quad \mathbf{E}_t = \frac{\partial \mathbf{f}}{\partial \boldsymbol{\theta}}, \quad \mathbf{H}_t^{xx} = \frac{\partial^2 H_t}{\partial \mathbf{x}_t \partial \mathbf{x}_t}, \quad \mathbf{H}_t^{xu} = \frac{\partial^2 H_t}{\partial \mathbf{x}_t \partial \mathbf{u}_t} = (\mathbf{H}_t^{ux})', \\ \mathbf{H}_t^{uu} &= \frac{\partial^2 H_t}{\partial \mathbf{u}_t \partial \mathbf{u}_t}, \quad \mathbf{H}_t^{x\theta} = \frac{\partial^2 H_t}{\partial \mathbf{x}_t \partial \boldsymbol{\theta}}, \quad \mathbf{H}_t^{u\theta} = \frac{\partial^2 H_t}{\partial \mathbf{u}_t \partial \boldsymbol{\theta}}, \quad \mathbf{H}_T^{xx} = \frac{\partial^2 h}{\partial \mathbf{x}_T \partial \mathbf{x}_T}, \quad \mathbf{H}_T^{x\theta} = \frac{\partial^2 h}{\partial \mathbf{x}_T \partial \boldsymbol{\theta}}, \end{aligned} \quad (13)$$

if  $\mathbf{H}_t^{uu}$  is invertible for all  $t = 0, \dots, T-1$ , define the following recursions

$$\begin{aligned} \mathbf{V}_t &= \mathbf{C}_t + \mathbf{A}_t'(\mathbf{I} + \mathbf{V}_{t+1}\mathbf{B}_t)^{-1}\mathbf{V}_{t+1}\mathbf{A}_t, \\ \mathbf{W}_t &= \mathbf{A}_t'(\mathbf{I} + \mathbf{V}_{t+1}\mathbf{B}_t)^{-1}(\mathbf{W}_{t+1} + \mathbf{V}_{t+1}\mathbf{M}_t) + \mathbf{N}_t, \end{aligned} \quad (14)$$

with  $\mathbf{V}_T = \mathbf{H}_T^{xx}$  and  $\mathbf{W}_T = \mathbf{H}_T^{x\theta}$ . Here,  $\mathbf{A}_t = \mathbf{F}_t - \mathbf{G}_t(\mathbf{H}_t^{uu})^{-1}\mathbf{H}_t^{ux}$ ,  $\mathbf{B}_t = \mathbf{G}_t(\mathbf{H}_t^{uu})^{-1}\mathbf{G}_t'$ ,  $\mathbf{M}_t = \mathbf{E}_t - \mathbf{G}_t(\mathbf{H}_t^{uu})'\mathbf{H}_t^{u\theta}$ ,  $\mathbf{C}_t = \mathbf{H}_t^{xx} - \mathbf{H}_t^{xu}(\mathbf{H}_t^{uu})^{-1}\mathbf{H}_t^{ux}$ ,  $\mathbf{N}_t = \mathbf{H}_t^{x\theta} - \mathbf{H}_t^{xu}(\mathbf{H}_t^{uu})'\mathbf{H}_t^{u\theta}$  are all known given (13).

Then, the partial derivative  $\frac{\partial \boldsymbol{\xi}(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$  can be obtained by recursively solving the following equations from  $t = 0$  to  $T-1$  with  $\mathbf{X}_0(\boldsymbol{\theta}) = \mathbf{0}$ :

$$\begin{aligned} \mathbf{U}_t &= -(\mathbf{H}_t^{uu})^{-1}(\mathbf{H}_t^{ux}\mathbf{X}_t + \mathbf{H}_t^{u\theta} + \mathbf{G}_t'(\mathbf{I} + \mathbf{V}_{t+1}\mathbf{B}_t)^{-1}(\mathbf{V}_{t+1}\mathbf{A}_t\mathbf{X}_t + \mathbf{V}_{t+1}\mathbf{M}_t + \mathbf{W}_{t+1})), \\ \mathbf{X}_{t+1} &= \mathbf{F}_t\mathbf{X}_t + \mathbf{G}_t\mathbf{U}_t + \mathbf{E}_t. \end{aligned} \quad (15)$$

The terms in (13) are based on the trajectory  $\boldsymbol{\xi}(\boldsymbol{\theta})$  and the associated Lagrangian multiplier  $\boldsymbol{\lambda}_{0:T-1}$ . According to the discrete-time Pontryagin Maximum Principle (Jin et al., 2020), the trajectory of the Lagrangian multiplier can be obtained by

$$\boldsymbol{\lambda}_T = \frac{\partial h}{\partial \mathbf{x}_T}, \quad \boldsymbol{\lambda}_t \triangleq \frac{\partial H_t}{\partial \mathbf{x}_t} = \frac{\partial c}{\partial \mathbf{x}_t} + \frac{\partial h}{\partial \mathbf{x}_t} \boldsymbol{\lambda}_{t+1}, \quad \text{for } t = T-1, \dots, 1. \quad (16)$$

### 3.3 OCIL Framework

With the online parameter estimator and the gradient generator, we propose the Online Control-Informed Learning framework in Fig. 1. The framework is summarized in Algorithm 2.

As shown in Fig. 1, at each time step, the predefined OC system  $\Sigma(\boldsymbol{\theta})$  generates a system trajectory  $\boldsymbol{\xi}(\boldsymbol{\theta})$  by performing optimal control with given  $\mathbf{x}_0$  and  $\boldsymbol{\theta}$ . The trajectory  $\boldsymbol{\xi}(\boldsymbol{\theta})$  is then fed into the stage loss function  $l(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)$  and the gradient generator. Along with the information  $\mathbf{O}_t$  obtained at time  $t$ , the stage loss function generates  $\frac{\partial l(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{O}_t)}{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}$ , while  $\frac{\partial \boldsymbol{\xi}_t(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$  is generated by the gradient generator in Algorithm 1. The chain rule is then performed to obtain the Jacobian matrix  $\mathbf{L}_t$ , which is then passed into the online parameter estimator for the estimation of  $\boldsymbol{\theta}$ .

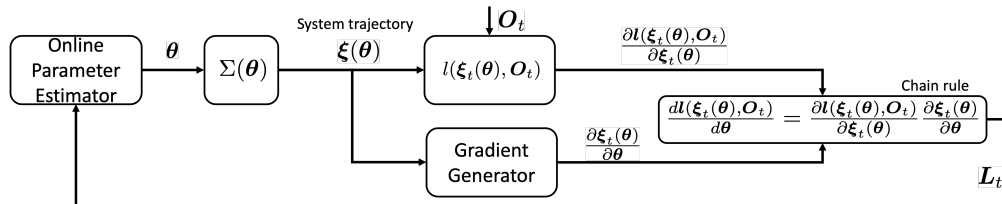


Figure 1: Framework of Online Control-Informed Learning.

**Algorithm 1:** Gradient Generator (GG)**Input:** Trajectory  $\xi(\theta)$  from  $\Sigma(\theta)$ 

- 1 Compute the coefficient matrices in (13) ;
- 2 Set  $V_T = H_T^{xx}$  and  $W_T = H_T^{x\theta}$ ;
- 3 **for**  $t \leftarrow T$  to 0 by  $\Delta t$  **do**
- 4   | Update  $V_t$  and  $W_t$  using (14)
- 5 Set  $X_0(\theta) = \mathbf{0}$ ;
- 6 **for**  $t \leftarrow 0$  to  $T$  by  $\Delta t$  **do**
- 7   | Update  $X_t(\theta)$  and  $U_t(\theta)$  using (15)

**Output:**  $\frac{\partial \xi(\theta)}{\partial \theta} = \{X_{0:T}(\theta), U_{0:T-1}(\theta)\}$ **Algorithm 2:** Online Control-Informed Learning**System and Stage Loss:**  $\Sigma(\theta)$  and  $l(\xi_t(\theta), O_t)$ **Initialize:**  $\theta_0, P_0$ 

- 1 **for**  $t = t_0, t_1, \dots$  **do**
- 2   | Obtain new information  $O_t$ ;
- 3   | Solve  $\xi(\theta_t)$  from current OC system  $\Sigma(\theta_t)$ ;
- 4   | Obtain  $\frac{\partial \xi_t(\theta_t)}{\partial \theta_t}$  with GG in Algorithm 1;
- 5   | Obtain  $\frac{\partial l(\xi_t(\theta_t), O_t)}{\partial \xi_t(\theta_t)}$  from  $l(\xi_t(\theta_t), O_t)$ ;
- 6   | Obtain  $L_t$  via the chain rule (10);
- 7   | Update  $\theta_t$  using the estimator (9);

### 3.4 Convergence Analysis

This subsection presents the convergence analysis of the online parameter estimator. The analysis employs a candidate Lyapunov function and introduces how the covariance matrices  $Q_t$  and  $R_t$  affect the convergence of the cumulative loss  $L(\xi(\theta))$ . In this section, for brevity, the stage loss  $l(\xi_t(\theta), O_t)$  is written as  $l(\theta)$ . Suppose for a specific task, there is an unknown true parameter  $\theta_t$ , where  $\theta_t = \theta_{t-1}$  for all  $t$ , and  $L(\xi(\theta_t)) = 0$ . Then, we define the estimation error as  $\tilde{\theta}_t = \theta_t - \hat{\theta}_t$ . Furthermore, we define

$$\begin{aligned} \text{Measurement error: } e_t &= l(\theta_t) - l(\hat{\theta}_t^-) \\ \text{Prediction error: } \tilde{\theta}_t^- &= \theta_t - \hat{\theta}_t^- \end{aligned} \quad (17)$$

To perform the convergence analysis, a candidate Lyapunov function is employed:

$$V_t = \tilde{\theta}_t' P_t^{-1} \tilde{\theta}_t. \quad (18)$$

The goal here is to determine conditions for which the candidate Lyapunov function  $\{V_t\}_{t=1,2,\dots}$  is a decreasing sequence, i.e.  $V_{t+1} - V_t \leq 0, \forall t$ . For rigorous analysis of the candidate Lyapunov function, as proposed in Boutayeb et al. (1997), unknown diagonal matrices  $\alpha_t \in \mathbb{R}^{r \times r}$  and  $\beta_t \in \mathbb{R}^{p \times p}$  are introduced to model the measurement and prediction error defined in (17):

$$\alpha_t e_t = L_t \tilde{\theta}_t^-, \quad \tilde{\theta}_t^- = \beta_t \tilde{\theta}_{t-1}. \quad (19)$$

To ensure convergence of the proposed estimator, the following assumptions need to be made.

**Assumption 1.** The derivative  $L_t = \frac{dl(\theta)}{d\theta}$  is of full rank for every  $\theta$ .

**Remark 1.** The discrete-time dynamical system (8) satisfies the observability rank condition, i.e., for every  $\theta$ ,  $\text{rank}(\text{col}\{\frac{dl(\theta)}{d\theta}, \frac{dl(\theta)}{d\theta} \mathbf{I}_p, \dots, \frac{dl(\theta)}{d\theta} \mathbf{I}_p^{p-1}\}) = p$  (Song & Grizzle, 1992). That means if Assumption 1 is satisfied for every  $\theta$ , the system (8) is observable for every  $\theta$ . The observability condition assures that  $P_t$  is a bounded matrix from above and below (Song & Grizzle, 1992; Boutayeb & Aubry, 1999).

As common in the EKF analysis, we adopt the following assumption:

**Assumption 2.**  $L_t$  is a uniformly bounded matrix.

We have the following lemma to show how the covariance matrices  $Q_t$  and  $R_t$  affect the convergence of the tunable parameter. The proof can be found in Appendix A.

**Lemma 2.** Let Assumptions 1 and 2 hold. If the matrices  $R_t$  and  $Q_t$  satisfy the following inequalities:

$$(\alpha_t - \mathbf{I}_s)^2 \leq R_t (L_t P_t^- L_t' + R_t)^{-1}, \quad (20)$$

$$\beta_t' (P_t + Q_t)^{-1} \beta_t - P_t^{-1} \leq 0, \quad (21)$$

Then the proposed estimator (9), when used as an observer for the system (8), ensures local asymptotic convergence, i.e.  $\lim_{t \rightarrow \infty} \tilde{\theta}_t = \mathbf{0}$ .

**Remark 2.** Lemma 2 provides sufficient conditions for the convergence of  $\hat{\boldsymbol{\theta}}$ . As the diagonal matrices  $\boldsymbol{\alpha}_t$  and  $\boldsymbol{\beta}_t$  are unknown, one can design the matrices  $\mathbf{Q}_t$  and  $\mathbf{R}_t$  to satisfy inequalities (20) and (21). For example, one can set the matrix  $\mathbf{Q}_t$  to be sufficiently large so that (21) is satisfied, which means the parameter estimator can tolerate arbitrary large initial state estimation error, i.e. large  $\boldsymbol{\beta}_t$ . It is worth to note that as long as (20) and (21) are satisfied,  $\hat{\boldsymbol{\theta}}_t$  converges to  $\boldsymbol{\theta}_t$  and consequently  $\boldsymbol{\alpha}_t$  and  $\boldsymbol{\beta}_t$  become identity matrix. In the case when there is no noise, i.e.  $\mathbf{Q}_t = \mathbf{0}_{p \times p}$  and  $\mathbf{R}_t = \mathbf{0}_{s \times s}$ ,  $\boldsymbol{\alpha}_t$  and  $\boldsymbol{\beta}_t$  can only be identity matrices to satisfy the inequalities (20) and (21), indicating the convergence of  $\hat{\boldsymbol{\theta}}_t$  to  $\boldsymbol{\theta}_t$ .

**Remark 3.** Equation (20) and (21) indicate one of the limitations of the estimator, which is the selection of initial guess. If the initial guess of  $\boldsymbol{\theta}$  results in  $\boldsymbol{\alpha}_0$  and  $\boldsymbol{\beta}_0$  that does not satisfy (20) and (21), the value of the Lyapunov function (18) becomes larger, which leads to even larger  $\boldsymbol{\alpha}_t$  and  $\boldsymbol{\beta}_t$ , causing the estimation to fail.

We have the following main theorem show how the covariance matrix  $\mathbf{R}_t$  and  $\mathbf{Q}_t$  affect the convergence of cumulative loss  $L(\boldsymbol{\xi}(\boldsymbol{\theta}_t))$  by utilizing the inequalities introduced in Lemma 2. The proof can be found in Appendix B.

**Theorem 1.** Let Assumptions 1 and 2 hold. If  $\mathbf{R}_t$  and  $\mathbf{Q}_t$  satisfy the following inequalities:

$$(\boldsymbol{\alpha}_t - \mathbf{I}_s)^2 \leq \mathbf{R}_t(\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}_t' + \mathbf{R}_t)^{-1}, \quad \boldsymbol{\beta}_t'(\mathbf{P}_t + \mathbf{Q}_t)^{-1} \boldsymbol{\beta}_t - \mathbf{P}_t^{-1} \leq 0, \quad (22)$$

then finding  $\boldsymbol{\theta}$  with the proposed estimator (9) employing the gradient generator in (14)-(15) ensures local asymptotic convergence of the cumulative loss  $L$  in (3) to 0, i.e.  $\lim_{t \rightarrow \infty} L(\boldsymbol{\xi}(\boldsymbol{\theta}_t)) = 0$ .

To evaluate the effectiveness of online algorithms, the regret analysis is often adopted (Li et al., 2021; Zinkevich, 2003; Hazan et al., 2007). In our case, the regret at time  $T$  is defined as:

$$\text{Regret}_T = \sum_{t=1}^T (\|\mathbf{l}(\hat{\boldsymbol{\theta}}_t)\|^2 - \|\mathbf{l}(\boldsymbol{\theta}_t)\|^2), \quad (23)$$

which represents the accumulative performance discrepancy between  $\hat{\boldsymbol{\theta}}_t$  and  $\boldsymbol{\theta}_t$ . Assume  $\|\mathbf{l}(\boldsymbol{\theta})\|^2$  is convex with respect to  $\boldsymbol{\theta}$ . As Lemma 2 states, the sufficient conditions (20) and (21) are satisfied for all  $t$ , then  $\hat{\boldsymbol{\theta}}_t$  will converge to  $\boldsymbol{\theta}_t$  asymptotically. Furthermore, according to (19),  $\mathbf{e}_t$  will converge to zero, which indicates that  $\|\mathbf{l}(\hat{\boldsymbol{\theta}}_t)\|^2$  will converge to  $\|\mathbf{l}(\boldsymbol{\theta}_t)\|^2$  asymptotically. Therefore, we can say that the regret in (23) will converge to a bounded value as  $t$  goes to infinity. Please refer to C in Appendix for details.

## 4 Applications to Different Online Learning Modes and Experiments

This section demonstrates the capability of the proposed OCIL framework with its three modes by three applications, Online Imitation Learning, Online System Identification, and Learning Policy on-the-fly. This section includes a performance comparison with some state-of-the-art frameworks for three environments that are summarized in Table 1. For every environment, the unknown parameter is static, which means  $\mathbf{Q}_t = \mathbf{0}_{p \times p}$ . The covariance matrix of measurement noise is set to be  $\mathbf{R}_t = 10^{-8}$  for OCIL, and  $\mathbf{R}_t = 0$  for other methods as they fail to complete the tasks with measurement noise. To highlight the flexibility of OCIL, each experiment includes two phases: 1) online phase, where OCIL keeps learning the unknown parameter while new data comes in before the final time  $T$ ; 2) offline phase, where OCIL keeps learning the parameter given the learned parameter at time  $T$  and the entire trajectory obtained from time  $t = 0$  to time  $T$ . For each environment and task, a terminal time  $T \in \mathbb{Z}$  is defined to represent a desired time duration where the system shall finish the task.

To unify the data visualization of both online and offline phases, the horizontal axis represents the number of data points, where a vertical red line corresponds to the final time  $T$ , i.e. the end of the online phase. The number of data points reflects the number of iterations multiplied by the total number of time steps for each iteration. The solid blue curves indicate the online portion of OCIL, whereas the dashed blue curves indicate the offline portion. For every environment and every method, 5 trials are performed given random initial conditions due to the high computational cost for other methods. The computational performance and analysis for OCIL are shown in Section 5 of the Appendix.

**Online Imitation Learning.** The control objective is parameterized as a weighted distance to the goal. Set the stage loss of imitation learning  $l(\boldsymbol{\xi}_t(\boldsymbol{\theta}), \mathbf{y}_t^*) = \mathbf{y}_t^* - \mathbf{g}(\mathbf{x}_t(\boldsymbol{\theta}), \mathbf{u}_t(\boldsymbol{\theta}))$ . Four existing methods are used for



Table 1: Experiment Environments

Systems	Dynamics parameter $\theta_{dyn}$	Objective parameter $\theta_{obj}$
Cartpole	cart mass, pole mass and length	$c(\mathbf{x}, \mathbf{u}) = \theta_{obj} \ \mathbf{x} - \mathbf{x}_g\ ^2 + \ \mathbf{u}\ ^2$ $h(\mathbf{x}) = \theta_{obj} \ \mathbf{x} - \mathbf{x}_g\ ^2$
6-DoF Quadrotor	mass, wing length, inertia matrix	
6-DoF Rocket	mass, rocket length, inertia matrix	

comparisons: (i) inverse KKT (Englert et al., 2017) (ii) neural policy cloning (Bojarski et al., 2016) and (iii) PDP (Jin et al., 2020). These methods don't handle measurement noise well because of their limitations, so we performed the experiments without including measurement noise for these methods. Fig. 2a-2c summarize the comparison result, where OCIL converges faster and obtains lower loss than the other offline methods, in both online and offline phases. The initial loss for each method is different because the learning representation (parameterization) is different. Thus, it is hard to guarantee that an initial neural network has the same loss as another initial parameter vector. Nevertheless, the initial representation of each method is adjusted such that OCIL does not take advantage of good initialization. Fig. 2a-2c validate the effectiveness of OCIL's both online and offline performance, even with measurement noise.

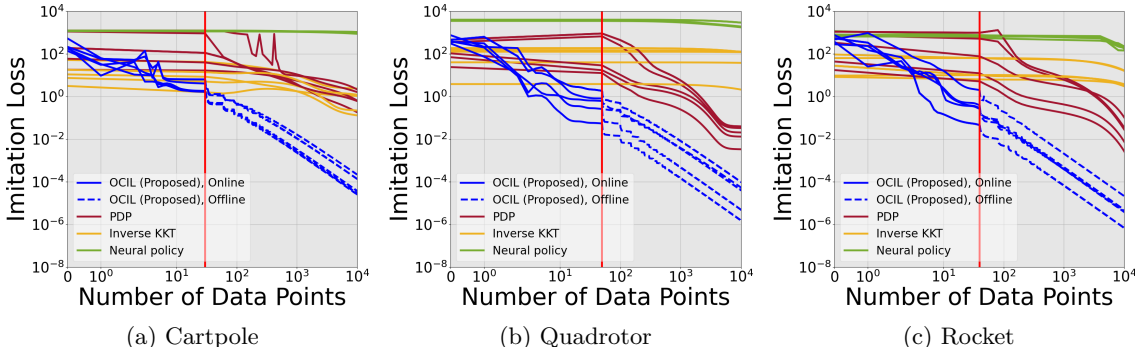


Figure 2: Imitation loss v.s. number of data points

**Online System Identification** The stage loss is set to be  $l(\xi_t(\theta), \xi_t^o) = \xi_t^o - \xi_t(\theta)$ . Three other methods are used for comparison: (i) Pytorch Adam solver, (ii) DMDC and (iii) PDP. No measurement noise are injected into observed data for existing methods due to their inherent limitations. Fig. 3a-3c summarize the result, where OCIL outperforms PDP for faster convergence and lower loss, in both online and offline phases. Different than Online Imitation Learning, OCIL does not decrease its SysID loss significantly at first because the number of data points is not sufficient for online learning. Once the number of data points becomes sufficient, the SysID loss starts decreasing significantly. This phenomenon can also be observed in the other methods, but their critical number of data points is significantly larger than OCIL's. In Fig. 3d-3f, OCIL and other methods are applied to learn the neural dynamics using the same observed trajectory. It can be seen that OCIL outperforms other methods for lower loss. Fig. 4 demonstrates the capability of OCIL dealing with neural dynamics that has different size of NN.

**Policy Tuning On-the-fly.** The parameterized OC system in 6 is used here, where the policy is in a state-feedback form and parameterized by the tunable parameter  $\theta$ . The stage loss is set to be  $l(\xi_t(\theta), \xi_t^*) = \xi_t^* - \xi_t(\theta)$ , where  $\xi_t^*$  is the trajectory that needs to be tracked. Other methods are used for comparison (i) iLQR (ii) GPS, and (iii) PDP. No measurement noise is included for existing methods due to their limitations. Fig. 5a-5c summarize the result, where the loss and its variation of OCIL converge very quickly. The buffers in Fig. 5a-5f indicate 3 times of standard deviation. Fig. 5d-5f presents the online phase of OCIL given 1000 random trials, which further validates the effectiveness and robustness of OCIL given measurement noise.

In general, OCIL from all figures does not have a smooth loss trajectory as the other offline methods. This is because at the online phase, an optimal gain matrix  $\mathbf{K}_t$  from (9a) is computed to update  $\theta$ , whereas the other methods either use a constant or iteration-dependent step size. The optimal gain is conceptually similar

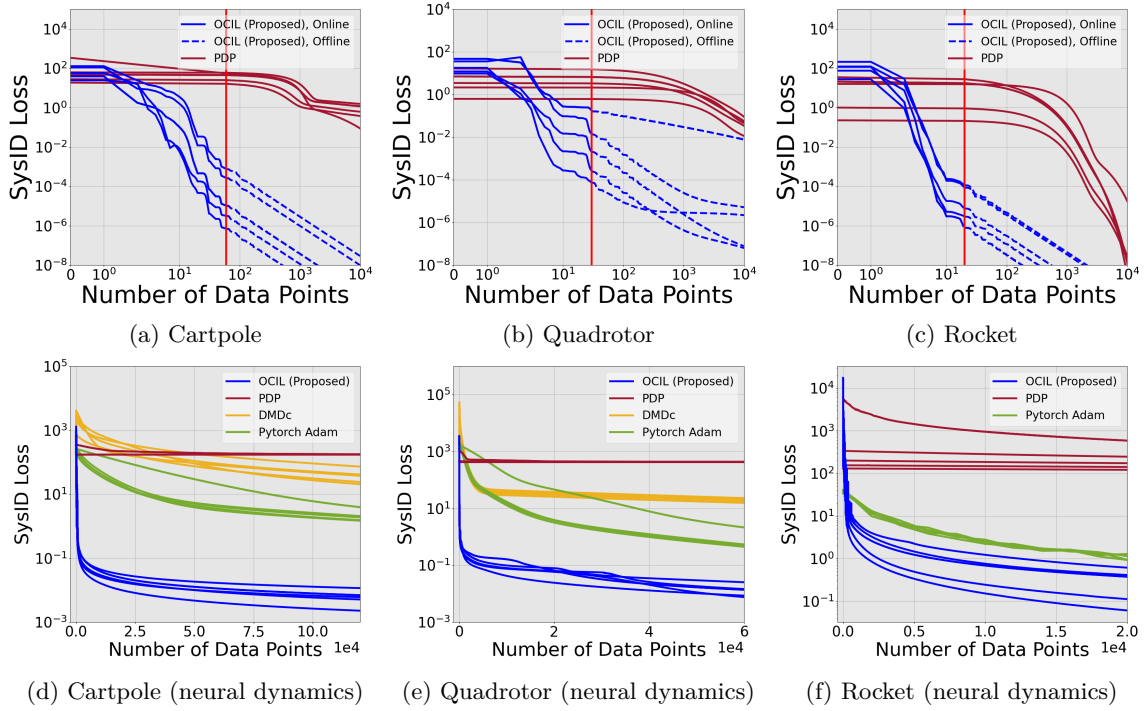


Figure 3: SysID loss v.s. number of data points

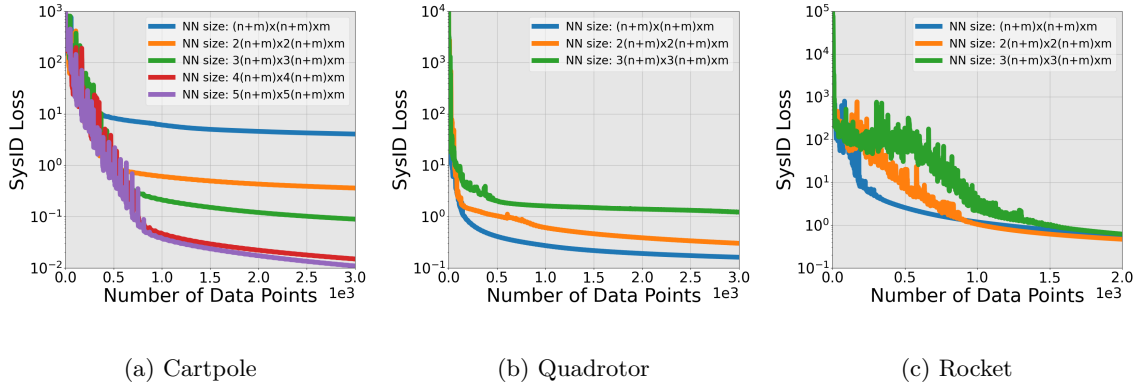


Figure 4: SysID Loss v.s. number of data points with different size of neural dynamics

to searching an optimal step size in the line-search optimization algorithms. Thus, it is observed that the loss variation, as represented by blue buffers, is relatively high initially but starts decreasing significantly as new data comes in because  $\mathbf{K}_t$  is continually updated. In contrast, the loss variation barely changes for the other offline methods after some data points.

## 5 Real-time Computational Performance

The experiments with OCIL were performed on a desktop with one Intel Core i7-8700k CPU with 8GB RAM. No GPU was used. The experiments with other methods were performed on a desktop with one AMD Ryzen 9 5900X CPU, one Nvidia Geforce RTX 4070ti, and 32 GB RAM. A more powerful PC was selected for the other methods because of their high computational cost. As noted at the beginning of Section 4, only 5 trials were conducted due to the computational expense.

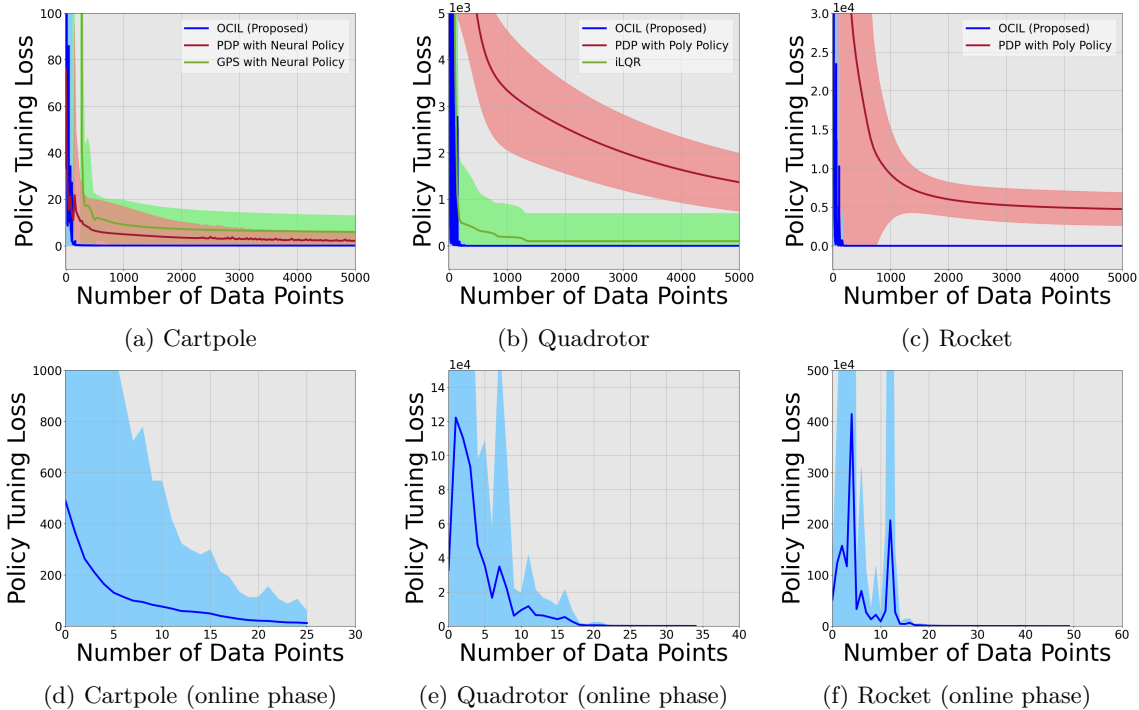


Figure 5: Policy Tuning Loss v.s. number of data points. Buffers represent loss variation under  $3\sigma$  with random initial conditions.

To demonstrate that the computational performance of OCIL is enough to be used in an online fashion, we recorded the computational time for OCIL in different modes for 100 trials. Table 3 summarizes the OCIL’s computational performance for the system identification task for three environments, where OCIL Time indicates the computational time of running OCIL at each time  $t$ , i.e. the iteration within the for-loop of Algorithm 2; GG Time indicates the computational time of running gradient generator (GG) at each time  $t$ , i.e. Algorithm 1; Estimator Time indicates the computational time of updating  $\hat{\theta}$ , i.e. Line 7 of Algorithm 2;  $\Delta$  indicates the time step of each environment, i.e., the time duration between two consecutive data measurements or the maximum allowed time duration of online algorithms to perform computation; Percentage indicates the percentage of the average OCIL time with respect to  $\Delta$ . The header of Table 4 and Table 2 are the same. Roughly speaking, OCIL time = GG time + Estimator time + Optimal Control computation time.

Table 3 illustrates that OCIL can estimate the dynamical system with neural network representation in real time, within the system frequency of getting new data. Table 4 illustrates that OCIL can tune the neural policy in real time. As indicated in Line 2 of Algorithm 2, the most computationally heavy part is solving optimal control trajectory in real time, instead of GG and the parameter estimator. As demonstrated at the beginning of this section, OCIL does not require huge computational resources, such as GPU. Therefore, this paper believes that OCIL has the capability to run in an online fashion.

Table 2: Computational Performance for Online Imitation Learning

Env.	OCIL Time [ms]	GG Time [ms]	Estimator Time [ms]	$\Delta$ [ms]	Percentage
Cartpole	$62.10 \pm 6.63$	$7.47 \pm 0.25$	$0.031 \pm 0.0023$	100	62.10 %
Quadrotor	$81.70 \pm 2.51$	$21.72 \pm 0.84$	$0.058 \pm 0.039$	100	81.70 %
Rocket	$72.25 \pm 13.91$	$19.77 \pm 6.26$	$0.060 \pm 0.012$	100	72.25%

Table 3: Computational Performance for SysID with Neural System

Env.	OCIL Time [ms]	GG Time [ms]	Estimator Time [ms]	$\Delta$ [ms]	Percentage
Cartpole	$17.18 \pm 5.15$	$6.05 \pm 2.59$	$1.93 \pm 0.85$	50	34.36 %
Quadrotor	$35.53 \pm 8.98$	$12.18 \pm 4.77$	$16.18 \pm 6.11$	100	35.53 %
Rocket	$29.54 \pm 8.29$	$11.25 \pm 4.67$	$12.89 \pm 5.29$	200	14.77 %

Table 4: Computational Performance for Policy Tuning with Neural Policy

Env.	OCIL Time [ms]	GG Time [ms]	Estimator Time [ms]	$\Delta$ [ms]	Percentage
Cartpole	$16.41 \pm 5.35$	$7.72 \pm 3.07$	$3.96 \pm 1.91$	50	32.82 %
Quadrotor	$62.67 \pm 9.51$	$30.86 \pm 2.25$	$22.47 \pm 8.33$	100	62.67 %
Rocket	$59.02 \pm 7.25$	$33.99 \pm 2.98$	$12.94 \pm 5.62$	100	59.02 %

## 6 Limitations.

This section discusses the major limitations of the proposed framework from three perspectives.

**Local convergence:** Since OCIL is based on first-order gradient, it can only achieve local minima for general non-convex optimization problems in (3). Furthermore, the general problem in (3) belongs to a bi-level optimization framework. Under certain assumptions such as convexity and smoothness on models (e.g., dynamics model, policy, loss function, and control objective function), global convergence of the bi-level optimization can be established. However, this paper believes that such conditions are too restrictive in the context of dynamical control systems. Therefore, this paper believes that the local convergence analysis based on general nonlinear optimization is enough.

**Parameterization matters for global convergence:** When performing experiments, we find that how models are parameterized matters for good convergence performance. For example, in online SysID mode, we observe that using a neural network dynamics (in Fig. 3d-3f) is more likely to get trapped in local minima than using the true dynamics with unknown parameters (in Fig. 3a-3c). In general, more complex parameterization will bring extreme non-convexity to the optimization problem, making the algorithm more easily trapped in local minima. Determining the parameterization of an object to be learned requires prior or expert knowledge, which is common in ML.

**Initialization matters:** As OCIL borrows how optimal gain updates from EKF, they share the same drawback that convergence depends on the selection of initialization. As shown in Remark 3, a bad initial guess might cause the estimator to diverge according to Lemma 2. Therefore, if a relatively good initial guess is hard to retrieve, one might need to use other methods to cold start OCIL.

## 7 Conclusions

This paper proposes Online Control-Informed Learning (OCIL), an online learning method tailored for diverse learning tasks. By considering an optimal control system tunable through a tunable parameter, OCIL effectively addresses tasks such as online imitation learning, online system identification, and tuning policy on-the-fly. By designing a stage loss specific to each task and treating the tunable parameter as a state of a new system, we employ the online parameter estimator to estimate the parameter in real-time and minimize loss at each time step. Theoretical analysis establishes the convergence conditions for OCIL, while experiments on various environments, tasks, and existing methods are done to validate its data efficiency, versatility, and robustness against measurement noise.

## References

- Pieter Abbeel and Andrew Y Ng. Apprenticeship learning via inverse reinforcement learning. In *International Conference on Machine Learning*, pp. 1–8, 2004.
- Rahaf Aljundi, Klaas Kelchtermans, and Tinne Tuytelaars. Task-free continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 11254–11263, 2019.
- Martin Benning, Elena Celledoni, Matthias J Ehrhardt, Brynjulf Owren, and Carola-Bibiane Schönlieb. Deep learning as optimal control problems: models and numerical methods. *arXiv preprint arXiv:1904.05657*, 2019.
- Dimitri P. Bertsekas. Incremental least squares methods and the extended kalman filter. *SIAM Journal on Optimization*, 6(3):807–822, 1996. doi: 10.1137/S1052623494268522.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseem Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- M. Boutayeb and D. Aubry. A strong tracking extended kalman observer for nonlinear discrete-time systems. *IEEE Transactions on Automatic Control*, 44(8):1550–1556, 1999. doi: 10.1109/9.780419.
- M. Boutayeb, H. Rafaralahy, and M. Darouach. Convergence analysis of the extended kalman filter used as an observer for nonlinear deterministic discrete-time systems. *IEEE Transactions on Automatic Control*, 42(4):581–586, 1997. doi: 10.1109/9.566674.
- Arthur Earl Bryson. *Applied optimal control: optimization, estimation and control*. CRC Press, 1975.
- Tian Qi Chen, Yulia Rubanova, Jesse Bettencourt, and David K Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems*, pp. 6571–6583, 2018.
- Wojciech M Czarnecki, Razvan Pascanu, Simon Osindero, Siddhant Jayakumar, Grzegorz Swirszcz, and Max Jaderberg. Distilling policy distillation. In *The 22nd International Conference on Artificial Intelligence and Statistics*, pp. 1331–1340. PMLR, 2019.
- Marc Deisenroth and Carl E Rasmussen. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on Machine Learning (ICML-11)*, pp. 465–472, 2011.
- Peter Englert, Ngo Anh Vien, and Marc Toussaint. Inverse kkt: Learning cost functions of manipulation tasks from demonstrations. *The International Journal of Robotics Research*, 36(13-14):1474–1488, 2017.
- Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. In *Advances in Neural Information Processing Systems*, pp. 64–72, 2016.
- Katerina Fragkiadaki, Sergey Levine, Panna Felsen, and Jitendra Malik. Recurrent network models for human dynamics. In *IEEE International Conference on Computer Vision*, pp. 4346–4354, 2015.
- Reza Ghaemi, Jing Sun, and Ilya V Kolmanovsky. Neighboring extremal solution for nonlinear discrete-time optimal control problems with state inequality constraints. *IEEE Transactions on Automatic Control*, 54(11):2674–2679, 2009.
- Jiequn Han and Weinan E. Deep learning approximation for stochastic control problems. *Deep Reinforcement Learning Workshop, Advances in Neural Information Processing Systems*, 2016.
- Jiequn Han, Qianxiao Li, et al. A mean-field optimal control formulation of deep learning. *Research in the Mathematical Sciences*, 6(1):1–41, 2019.
- Wenjia Hao, Paulo C Heredia, Bowen Huang, Zehui Lu, Zihao Liang, and Shaoshuai Mou. Policy learning based on deep koopman representation. *arXiv preprint arXiv:2305.15188*, 2023.

- Masahiko Haruno, Daniel M Wolpert, and Mitsuo Kawato. Mosaic model for sensorimotor learning and control. *Neural Computation*, 13(10):2201–2220, 2001.
- Elad Hazan, Amit Agarwal, and Satyen Kale. Logarithmic regret algorithms for online convex optimization. *Machine Learning*, 69(2):169–192, 2007.
- Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning. *Advances in Neural Information Processing Systems*, 29, 2016.
- Wanxin Jin and Shaoshuai Mou. Distributed inverse optimal control. *Automatica*, 129:109658, 2021. ISSN 0005-1098. doi: <https://doi.org/10.1016/j.automatica.2021.109658>. URL <https://www.sciencedirect.com/science/article/pii/S0005109821001783>.
- Wanxin Jin, Dana Kulić, Jonathan Feng-Shun Lin, Shaoshuai Mou, and Sandra Hirche. Inverse optimal control for multiphase cost functions. *IEEE Transactions on Robotics*, 35(6):1387–1398, 2019.
- Wanxin Jin, Zhaoran Wang, Zhuoran Yang, and Shaoshuai Mou. Pontryagin differentiable programming: An end-to-end learning and control framework. *Advances in Neural Information Processing Systems*, 33:7979–7992, 2020.
- Wanxin Jin, Dana Kulić, Shaoshuai Mou, and Sandra Hirche. Inverse optimal control from incomplete trajectory observations. *The International Journal of Robotics Research*, 40(6-7):848–865, 2021a.
- Wanxin Jin, Shaoshuai Mou, and George J Pappas. Safe pontryagin differentiable programming. *Advances in Neural Information Processing Systems*, 34:16034–16050, 2021b.
- Wanxin Jin, Todd D Murphey, Dana Kulić, Neta Ezer, and Shaoshuai Mou. Learning from sparse demonstrations. *IEEE Transactions on Robotics*, 39(1):645–664, 2022a.
- Wanxin Jin, Todd D Murphey, Zehui Lu, and Shaoshuai Mou. Learning from human directional corrections. *IEEE Transactions on Robotics*, 39(1):625–644, 2022b.
- George Em Karniadakis, Ioannis G Kevrekidis, Lu Lu, Paris Perdikaris, Sifan Wang, and Liu Yang. Physics-informed machine learning. *Nature Reviews Physics*, 3(6):422–440, 2021.
- Karel J Keesman. *System identification: an introduction*. Springer Science & Business Media, 2011.
- Arezou Keshavarz, Yang Wang, and Stephen Boyd. Imputing a convex objective function. In *IEEE International Symposium on Intelligent Control*, pp. 613–619, 2011.
- Beomjoon Kim, Amir-massoud Farahmand, Joelle Pineau, and Doina Precup. Learning from limited demonstrations. *Advances in Neural Information Processing Systems*, 26, 2013.
- Jack B Kuipers. *Quaternions and rotation sequences*, volume 66. Princeton University Press, 1999.
- Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. Deep learning. *Nature*, 521(7553):436–444, 2015.
- Sergey Levine and Pieter Abbeel. Learning neural network policies with guided policy search under unknown dynamics. In *Advances in Neural Information Processing Systems*, pp. 1071–1079, 2014.
- Qianxiao Li and Shuji Hao. An optimal control approach to deep learning and applications to discrete-weight neural networks. *arXiv preprint arXiv:1803.01299*, 2018.
- Qianxiao Li, Long Chen, Cheng Tai, and E Weinan. Maximum principle based algorithms for deep learning. *Journal of Machine Learning Research*, 18(165):1–29, 2018.
- Weiwei Li and Emanuel Todorov. Iterative linear quadratic regulator design for nonlinear biological movement systems. In *International Conference on Informatics in Control, Automation and Robotics*, pp. 222–229, 2004.

- Yingying Li, Guannan Qu, and Na Li. Online optimization with predictions and switching costs: Fast algorithms and the fundamental limit. *IEEE Transactions on Automatic Control*, 66(10):4761–4768, 2021. doi: 10.1109/TAC.2020.3040249.
- Zihao Liang, Wanxin Jin, and Shaoshuai Mou. An iterative method for inverse optimal control. In *2022 13th Asian Control Conference (ASCC)*, pp. 959–964, 2022. doi: 10.23919/ASCC56756.2022.9828009.
- Zihao Liang, Wenjian Hao, and Shaoshuai Mou. A data-driven approach for inverse optimal control. In *2023 62nd IEEE Conference on Decision and Control (CDC)*, pp. 3632–3637, 2023. doi: 10.1109/CDC49753.2023.10383220.
- Hailiang Liu and Peter Markowich. Selection dynamics for deep neural networks. *Journal of Differential Equations*, 269(12):11540–11574, 2020.
- Michael Lutter, Christian Ritter, and Jan Peters. Deep lagrangian networks: Using physics as model prior for deep learning. *arXiv preprint arXiv:1907.04490*, 2019.
- Alexandre Mauroy, Y Susuki, and Igor Mezić. *Koopman operator in systems and control*. Springer, 2020.
- Katja Mombaur, Anh Truong, and Jean-Paul Laumond. From human to humanoid locomotion—an inverse optimal control approach. *Autonomous Robots*, 28(3):369–383, 2010.
- Junhyuk Oh, Yijie Guo, Satinder Singh, and Honglak Lee. Self-imitation learning. In *International Conference on Machine Learning*, pp. 3878–3887. PMLR, 2018.
- Melkior Ornik, Steven Carr, Arie Israel, and Ufuk Topcu. Control-oriented learning on the fly. *IEEE Transactions on Automatic Control*, 65(11):4800–4807, 2019.
- Sinno Jialin Pan and Qiang Yang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 22(10):1345–1359, 2009.
- Joshua L Proctor, Steven L Brunton, and J Nathan Kutz. Dynamic mode decomposition with control. *SIAM Journal on Applied Dynamical Systems*, 15(1):142–161, 2016.
- Maziar Raissi, Paris Perdikaris, and George E Karniadakis. Physics-informed neural networks: A deep learning framework for solving forward and inverse problems involving nonlinear partial differential equations. *Journal of Computational Physics*, 378:686–707, 2019.
- Nathan D Ratliff, J Andrew Bagnell, and Martin A Zinkevich. Maximum margin planning. In *International Conference on Machine Learning*, pp. 729–736, 2006.
- Steindor Saemundsson, Alexander Terenin, Katja Hofmann, and Marc Deisenroth. Variational integrator networks for physically structured embeddings. In *International Conference on Artificial Intelligence and Statistics*, pp. 3078–3087, 2020.
- Fumihiro Sasaki and Ryota Yamashina. Behavioral cloning from noisy demonstrations. In *International Conference on Learning Representations*, 2021.
- Stefan Schaal. Learning from demonstration. *Advances in Neural Information Processing Systems*, 9, 1996.
- Yongkyu Song and Jessy W. Grizzle. The extended kalman filter as a local asymptotic observer for nonlinear discrete-time systems. In *1992 American Control Conference*, pp. 3365–3369, 1992. doi: 10.23919/ACC.1992.4792775.
- Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(7), 2009.
- Faraz Torabi, Garrett Warnell, and Peter Stone. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*, 2018.

- Laura Von Rueden, Sebastian Mayer, Katharina Beckh, Bogdan Georgiev, Sven Giesselbach, Raoul Heese, Birgit Kirsch, Julius Pfrommer, Annika Pick, Rajkumar Ramamurthy, et al. Informed machine learning—a taxonomy and survey of integrating prior knowledge into learning systems. *IEEE Transactions on Knowledge and Data Engineering*, 35(1):614–633, 2021.
- Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in Neural Information Processing Systems*, pp. 2746–2754, 2015.
- Matthew O Williams, Ioannis G Kevrekidis, and Clarence W Rowley. A data-driven approximation of the koopman operator: Extending dynamic mode decomposition. *Journal of Nonlinear Science*, 25(6): 1307–1346, 2015.
- Yueh-Hua Wu, Nontawat Charoenphakdee, Han Bao, Voot Tangkaratt, and Masashi Sugiyama. Imitation learning from imperfect demonstration. In *International Conference on Machine Learning*, pp. 6818–6827. PMLR, 2019.
- Amy Zhang, Sainbayar Sukhbaatar, Adam Lerer, Arthur Szlam, and Rob Fergus. Composable planning with attributes. In *International Conference on Machine Learning*, pp. 5842–5851, 2018.
- Dinghuai Zhang, Tianyuan Zhang, Yiping Lu, Zhanxing Zhu, and Bin Dong. You only propagate once: Painless adversarial training using maximal principle. *arXiv preprint arXiv:1905.00877*, 2019.
- Yaofeng Desmond Zhong, Biswadip Dey, and Amit Chakraborty. Symplectic ode-net: Learning hamiltonian dynamics with control. *arXiv preprint arXiv:1909.12077*, 2019.
- Brian D Ziebart, Andrew Maas, J Andrew Bagnell, and Anind K Dey. Maximum entropy inverse reinforcement learning. In *AAAI Conference on Artificial Intelligence*, pp. 1433–1438, 2008.
- Martin Zinkevich. Online convex programming and generalized infinitesimal gradient ascent. In *Proceedings of the 20th International Conference on Machine Learning (ICML-03)*, pp. 928–936, 2003.



## Appendix

### A Proof of Lemma 2

Since the matrices  $\mathbf{P}_t$  and  $\mathbf{L}_t$  are bounded according to Assumption 1 and 2, from (9b), one will have:

$$\mathbf{K}_t = \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \quad (24)$$

$$= \mathbf{P}_t^- \mathbf{L}'_t (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1}. \quad (25)$$

Then, by taking the inverse of (24) and (25), one will get:

$$\mathbf{P}_t^{-1} = (\mathbf{P}_t^-)^{-1} + \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{L}_t. \quad (26)$$

Substituting (24) into (9b) and subtracting both sides from  $\boldsymbol{\theta}_t$ , one will have:

$$\tilde{\boldsymbol{\theta}}_t = \tilde{\boldsymbol{\theta}}_t^- - \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t. \quad (27)$$

Then, plug (27) into the Lyapunov function (18):

$$V_t = \tilde{\boldsymbol{\theta}}_t' \mathbf{P}_t^{-1} \tilde{\boldsymbol{\theta}}_t \quad (28)$$

$$= (\tilde{\boldsymbol{\theta}}_t^- - \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t)' \mathbf{P}_t^{-1} (\tilde{\boldsymbol{\theta}}_t^- - \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t) \quad (29)$$

$$= (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{P}_t^{-1} \tilde{\boldsymbol{\theta}}_t^- - (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t - \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- + \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t \quad (30)$$

Next, we plug (26) into (30):

$$V_t = (\tilde{\boldsymbol{\theta}}_t^-)' ((\mathbf{P}_t^-)^{-1} + \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{L}_t) \tilde{\boldsymbol{\theta}}_t^- - (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t - \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- + \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t \quad (31)$$

$$= V_t^- + (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- - (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t - \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- + \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t, \quad (32)$$

where

$$V_t^- = (\tilde{\boldsymbol{\theta}}_t^-)' (\mathbf{P}_t^-)^{-1} \tilde{\boldsymbol{\theta}}_t^- \quad (33)$$

$$= (\tilde{\boldsymbol{\theta}}_{t-1})' \boldsymbol{\beta}'_t (\mathbf{P}_{t-1} + \mathbf{Q}_{t-1})^{-1} \boldsymbol{\beta}_t \tilde{\boldsymbol{\theta}}_{t-1}. \quad (34)$$

Using (19), (32) becomes:

$$V_t = V_t^- + (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- - (\tilde{\boldsymbol{\theta}}_t^-)' \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t - \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \tilde{\boldsymbol{\theta}}_t^- + \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t \quad (35)$$

$$= V_t^- + \mathbf{e}_t' \boldsymbol{\alpha}_t \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t \mathbf{e}_t - \mathbf{e}_t' \boldsymbol{\alpha}_t \mathbf{R}_t^{-1} \mathbf{e}_t - \mathbf{e}_t' \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t \mathbf{e}_t + \mathbf{e}_t' \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \mathbf{e}_t \quad (36)$$

$$= V_t^- + \mathbf{e}_t' (\boldsymbol{\alpha}_t \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t - \boldsymbol{\alpha}_t \mathbf{R}_t^{-1} - \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t + \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1}) \mathbf{e}_t. \quad (37)$$

To ensure that the Lyapunov function  $\{V_t\}_{t=1,2,\dots}$  is a decreasing sequence,  $V_t - V_{t-1} \leq 0$ .

$$V_t - V_{t-1} \quad (38)$$

$$= \mathbf{e}_t' (\boldsymbol{\alpha}_t \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t - \boldsymbol{\alpha}_t \mathbf{R}_t^{-1} - \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t + \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1}) \mathbf{e}_t \quad (39)$$

$$+ (\tilde{\boldsymbol{\theta}}_{t-1})' (\boldsymbol{\beta}'_t (\mathbf{P}_{t-1} + \mathbf{Q}_{t-1})^{-1} \boldsymbol{\beta}_t - \mathbf{P}_{t-1}^{-1}) \tilde{\boldsymbol{\theta}}_{t-1} \leq 0. \quad (40)$$

Therefore, to ensure the Lyapunov function is a decreasing sequence,

$$\boldsymbol{\alpha}_t \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t - \boldsymbol{\alpha}_t \mathbf{R}_t^{-1} - \mathbf{R}_t^{-1} \boldsymbol{\alpha}_t + \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \leq 0, \quad (41)$$

and

$$\boldsymbol{\beta}'_t (\mathbf{P}_{t-1} + \mathbf{Q}_{t-1})^{-1} \boldsymbol{\beta}_t - \mathbf{P}_{t-1}^{-1} \leq 0. \quad (42)$$

With some manipulations:

$$(\boldsymbol{\alpha}_t - \mathbf{I}_s) \mathbf{R}_t^{-1} (\boldsymbol{\alpha}_t - \mathbf{I}_s) - \mathbf{R}_t^{-1} + \mathbf{R}_t^{-1} \mathbf{L}_t \mathbf{P}_t \mathbf{L}'_t \mathbf{R}_t^{-1} \leq 0, \quad (43)$$

$$(\boldsymbol{\alpha}_t - \mathbf{I}_s) \mathbf{R}_t^{-1} (\boldsymbol{\alpha}_t - \mathbf{I}_s) - \mathbf{R}_t^{-1} (\mathbf{I}_s - \mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1}) \leq 0. \quad (44)$$

By letting  $\mathbf{I}_s = (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)(\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1}$ , we have

$$(\boldsymbol{\alpha}_t - \mathbf{I}_s) \mathbf{R}_t^{-1} (\boldsymbol{\alpha}_t - \mathbf{I}_s) - (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1} \leq 0. \quad (45)$$

Since  $\boldsymbol{\alpha}_t$  and  $\mathbf{R}_t$  are diagonal matrices, we will have

$$\mathbf{R}_t^{-1} (\boldsymbol{\alpha}_t - \mathbf{I}_s)^2 - (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1} \leq 0, \quad (46)$$

which at the end yields:

$$(\boldsymbol{\alpha}_t - \mathbf{I}_s)^2 \leq \mathbf{R}_t (\mathbf{L}_t \mathbf{P}_t^- \mathbf{L}'_t + \mathbf{R}_t)^{-1}, \quad (47)$$

therefore the proof is completed.

## B Proof of Theorem 1

This proof is straightforward once Lemma 2 is provided. Consider the assumptions 1 and 2 are met, according to Lemma 2, with the exact gradient generated by the gradient generator in (14)-(15),  $\lim_{t \rightarrow \infty} \hat{\boldsymbol{\theta}}_t = \boldsymbol{\theta}$ . As the estimated  $\hat{\boldsymbol{\theta}}$  converges to the true  $\boldsymbol{\theta}$ , where the true parameter gives zero cumulative loss, the cumulative loss  $L(\boldsymbol{\xi}(\hat{\boldsymbol{\theta}}))$  goes to 0.

## C Scratch for Regret Analysis

In the regret analysis for online optimization algorithms Zinkevich (2003); Hazan et al. (2007), the regret rate is often determined by the step size used for their gradient descent, which is typically a constant or an iteration-dependent value. Such step size can make the sequence of regret over time a Cauchy sequence. Consequently, one can induce that the accumulative regret is bounded.

As for the proposed OCIL framework, as mentioned at the end of Section 4, the optimal gain matrix  $\mathbf{K}_t$  is updated every time  $t$  as new data comes in. Then given the proposed estimate update (9b),  $\mathbf{K}_t$  updates the parameter estimate  $\hat{\boldsymbol{\theta}}_t$  by minimizing the loss function. Therefore, the combination  $\mathbf{K}_t(\mathbf{0} - \mathbf{l}(\boldsymbol{\xi}_t(\hat{\boldsymbol{\theta}}_t^-)))$  is conceptually similar to the product of the gradient and optimal step size, where the optimal step is size calculated by line search algorithms in the context of optimization. This statement is also verified by our experiments in Section 4.

To sum up, if 20 and 21 are satisfied for all  $t$ ,  $\mathbf{l}(\hat{\boldsymbol{\theta}}_t^-)$  diminishes as time increases, which indicates that the accumulative regret will converge to a bounded value.

## D Experiment Details

### D.1 System/Environment Setups

**Cartpole.** We consider the following continuous dynamics of the cartpole

$$\begin{bmatrix} \dot{p} \\ \ddot{p} \\ \dot{\theta} \\ \ddot{\theta} \end{bmatrix} = \begin{bmatrix} \dot{p} \\ (F + \frac{m_p l \dot{\theta}^2 \sin(\theta)}{m_t}) - \frac{m_p l \ddot{\theta} \cos(\theta)}{m_t} \\ \dot{\theta} \\ \frac{g \sin(\theta) - \cos(\theta)(F + \frac{m_p l \dot{\theta}^2 \sin(\theta)}{m_t})}{l(\frac{4}{3} - \frac{m_p \cos(\theta)^2}{m_t})} \end{bmatrix}, \quad (48)$$

where  $p \in \mathbb{R}$  is the horizontal displacement of the cart;  $\theta \in \mathbb{R}$  is the pole angle;  $F \in \mathbb{R}$  denotes the horizontal force applied to the cart which is between  $-1$  and  $+1$ ;  $l \in \mathbb{R}$  is the length of the pole;  $m_p, m_t \in \mathbb{R}$  are the masses of the pole and total cartpole, respectively. By defining the states and control inputs of the cartpole

$$\mathbf{x} \triangleq [p \quad \dot{p} \quad \theta \quad \dot{\theta}]' \quad \text{and} \quad \mathbf{u} \triangleq F \quad (49)$$

respectively.

**Quadrotor UAV.** We consider a quadrotor UAV with the following dynamics

$$\begin{aligned}
\dot{\mathbf{p}}_I &= \mathbf{v}_I, \\
m\dot{\mathbf{v}}_I &= m\mathbf{g}_I + \mathbf{F}_I, \\
\dot{\mathbf{q}}_{B/I} &= \frac{1}{2}\boldsymbol{\Omega}(\boldsymbol{\omega}_B)\mathbf{q}_{B/I}, \\
J_B\dot{\boldsymbol{\omega}}_B &= \mathbf{M}_B - \boldsymbol{\omega} \times J_B\boldsymbol{\omega}_B.
\end{aligned} \tag{50}$$

Here, the subscription  $B$  and  $I$  denote a quantity is expressed in the body frame and inertial (world) frame, respectively;  $m$  and  $J_B \in \mathbb{R}^{3 \times 3}$  are the mass and moment of inertia with respect to body frame of the UAV, respectively.  $g$  is the gravitational constant ( $g = 10\text{kg/m}^2$ ),  $\mathbf{g}_I = [0, 0, g]'$ .  $\mathbf{p} \in \mathbb{R}^3$  and  $\mathbf{v} \in \mathbb{R}^3$  are the position and velocity vector of the UAV;  $\boldsymbol{\omega}_B \in \mathbb{R}^3$  is the angular velocity vector of the UAV;  $\mathbf{q}_{B/I} \in \mathbb{R}^4$  is the unit quaternion Kuipers (1999) that describes the attitude of UAV with respect to the inertial frame;  $\boldsymbol{\Omega}(\boldsymbol{\omega}_B)$  is defined as:

$$\boldsymbol{\Omega}(\boldsymbol{\omega}_B) = \begin{bmatrix} 0 & -\omega_x & -\omega_y & -\omega_z \\ \omega_x & 0 & \omega_z & -\omega_y \\ \omega_y & -\omega_z & 0 & \omega_x \\ \omega_z & \omega_y & -\omega_x & 0 \end{bmatrix}, \tag{51}$$

$\mathbf{M}_B \in \mathbb{R}^3$  is the torque applied to the UAV;  $\mathbf{F}_I \in \mathbb{R}^3$  is the force vector applied to the UAV center of mass. The total force magnitude  $f = \|\mathbf{F}_I\| \in \mathbb{R}$  (along  $z$ -axis of the body frame) and torque  $\mathbf{M}_B = [M_x, M_y, M_z]'$  are generated by thrust from four rotating propellers  $[T_1, T_2, T_3, T_4]'$ , their relationship can be expressed as:

$$\begin{bmatrix} f \\ M_x \\ M_y \\ M_z \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 0 & -l_w/2 & 0 & l_w/2 \\ -l_w/2 & 0 & l_w/2 & 0 \\ c & -c & c & -c \end{bmatrix} \begin{bmatrix} T_1 \\ T_2 \\ T_3 \\ T_4 \end{bmatrix}, \tag{52}$$

where  $l_w$  is the wing length of the UAV and  $c$  is a fixed constant. The state and input vectors of the UAV are defined as:

$$\begin{aligned}
\mathbf{x} &\triangleq [\mathbf{p}' \quad \mathbf{v}' \quad \mathbf{q}' \quad \boldsymbol{\omega}']' \in \mathbb{R}^{13}, \\
\mathbf{u} &\triangleq [T_1 \quad T_2 \quad T_3 \quad T_4]' \in \mathbb{R}^4.
\end{aligned} \tag{53}$$

**Rocket.** The rocket is treated as a rigid body subject to constant gravitational acceleration,  $g_I \in \mathbb{R}^3$ , and neglect aerodynamic forces. The vehicle is assumed to actuate a single gimbaled rocket engine to generate a thrust vector within a feasible range of magnitudes and gimbal angles. We assume that at the landing phase, the depletion of fuel is insignificant. Therefore, we omit the dynamics of rocket mass. The rocket has the following dynamics:

$$\begin{aligned}
\dot{\mathbf{p}}_I &= \mathbf{v}_I, \\
\dot{\mathbf{v}}_I &= \frac{1}{m}C_{I/B}\mathbf{T}_B + \mathbf{g}_I, \\
\dot{\mathbf{q}}_{B/I} &= \frac{1}{2}\boldsymbol{\Omega}(\boldsymbol{\omega}_B)\mathbf{q}_{B/I}, \\
J_B\dot{\boldsymbol{\omega}}_B &= \mathbf{M}_B - [\boldsymbol{\omega}_B \times]J_B\boldsymbol{\omega}_B.
\end{aligned} \tag{54}$$

Here, the subscription  $B$  and  $I$  denote a quantity is expressed in the body frame and inertial (world) frame, respectively;  $m$  and  $J_B \in \mathbb{R}^{3 \times 3}$  are the mass and moment of inertia with respect to body frame of the rocket, respectively.  $\mathbf{p} \in \mathbb{R}^3$  and  $\mathbf{v} \in \mathbb{R}^3$  are the position and velocity vector of the rocket;  $\boldsymbol{\omega}_B \in \mathbb{R}^3$  is the angular velocity vector of the rocket;  $\mathbf{q}_{B/I} = [q_0, q_1, q_2, q_3]$  is the unit quaternion that describes the attitude of rocket with respect to the inertial frame;  $\mathbf{T}_B \in \mathbb{R}^3$  is the commanded thrust vector;  $\mathbf{M}_B \in \mathbb{R}^3$  is the torque applied to the rocket;  $C_{B/I}$  is the direction cosine matrix that encodes the attitude transformation from body frame to inertia frame and related to  $\mathbf{q}_{B/I}$  by the following relationship:

$$C_{\mathcal{B}/\mathcal{I}} = \begin{bmatrix} 1 - 2(q_2^2 + q_3^2) & 2(q_1q_2 + q_0q_3) & 2(q_1q_3 - q_0q_2) \\ 2(q_1q_2 - q_0q_3) & 1 - 2(q_1^2 + q_3^2) & 2(q_2q_3 + q_0q_1) \\ 2(q_1q_3 + q_0q_2) & 2(q_2q_3 - q_0q_1) & 1 - 2(q_1^2 + q_2^2) \end{bmatrix},$$

The inverse transformation is denoted as  $C_{\mathcal{I}/\mathcal{B}} = C_{\mathcal{B}/\mathcal{I}}^T$ ;

The skew-symmetric matrices  $[\boldsymbol{\omega}_{\mathcal{B}} \times]$  and  $\Omega(\boldsymbol{\omega}_{\mathcal{B}})$  are defined as follow:

$$[\boldsymbol{\omega}_{\mathcal{B}} \times] \triangleq \begin{bmatrix} 0 & -\omega_z & \omega_y \\ \omega_z & 0 & -\omega_x \\ -\omega_y & \omega_x & 0 \end{bmatrix}, \quad \Omega(\boldsymbol{\omega}_{\mathcal{B}}) \triangleq \begin{bmatrix} 0 & -\omega_x & -\omega_y & -\omega_z \\ \omega_x & 0 & \omega_z & -\omega_y \\ \omega_y & -\omega_z & 0 & \omega_x \\ \omega_z & \omega_y & -\omega_x & 0 \end{bmatrix},$$

The state and input vectors of the rocket are defined as:

$$\begin{aligned} \mathbf{x} &= [\mathbf{p}'_{\mathcal{I}} \quad \mathbf{v}'_{\mathcal{I}} \quad \mathbf{q}'_{\mathcal{B}/\mathcal{I}} \quad \boldsymbol{\omega}'_{\mathcal{B}}]' \in \mathbb{R}^{13}, \\ \mathbf{u} &= \mathbf{T}_{\mathcal{B}} = [T_x \quad T_y \quad T_z]' \in \mathbb{R}^3, \end{aligned} \tag{55}$$

**Discretization** Discretization is done by the following discrete-time form

$$\mathbf{x}_{t+1} \approx \mathbf{x}_t + \Delta \cdot \mathbf{g}(\mathbf{x}_t, \mathbf{u}_t) \triangleq \mathbf{f}(\mathbf{x}_t, \mathbf{u}_t), \tag{56}$$

where  $\Delta$  is the discretization interval.

## D.2 Online Imitation Learning

**Data acquisition.** The dataset of expert demonstrations is generated by solving an optimal control system with the true dynamics and control objective parameter  $\boldsymbol{\theta}^* = \{\boldsymbol{\theta}_{dyn}, \boldsymbol{\theta}_{obj}\}$  given. We generate a number of five trajectories, where different trajectories have different initial conditions  $\mathbf{x}_0$  and time horizons  $T$ .

**PDP.** We employed the PDP in Jin et al. (2020) to solve this problem. The learning rate is set to be  $\eta = 10^{-4}$ . Five trials were run given random initial  $\boldsymbol{\theta}_0$ .

**Inverse KKT method.** We choose the inverse KKT method Englert et al. (2017) for comparison because it is suitable for learning objective functions for high-dimensional continuous-space systems. We adapt the inverse KKT method, and define the KKT loss as the norm-2 violation of the KKT condition by the demonstration data:

$$\min_{\boldsymbol{\theta}, \boldsymbol{\lambda}_{1:T}} \left( \left\| \frac{\partial L}{\partial \mathbf{x}_{0:T}}(\mathbf{x}_{0:T}^*, \mathbf{u}_{0:T-1}^*) \right\|^2 + \left\| \frac{\partial L}{\partial \mathbf{u}_{0:T-1}}(\mathbf{x}_{0:T}^*, \mathbf{u}_{0:T-1}^*) \right\|^2 \right).$$

**Neural policy cloning.** For the neural policy cloning, we directly learn a neural network policy  $\mathbf{u} = \boldsymbol{\mu}(\mathbf{x}, \boldsymbol{\theta})$  from the dataset using supervised learning, that is

$$\min_{\boldsymbol{\theta}} \sum_{t=0}^{T-1} \|\mathbf{u}_t^* - \boldsymbol{\mu}(\mathbf{x}_t^*, \boldsymbol{\theta})\|^2 \tag{57}$$

## D.3 Online System Identification

**Data acquisition.** In the system identification experiment, we collect a total number of five trajectories from systems with dynamics known, wherein different trajectories  $\boldsymbol{\xi}^o = \{\mathbf{x}_{0:T}^o, \mathbf{u}_{0:T-1}\}$  have different initial conditions  $\mathbf{x}_0$  and horizons  $T$  ( $T$  ranges from 10 to 20 depending on different environment and task), with random inputs  $\mathbf{u}_{0:T-1}$  drawn from uniform distribution.

**PDP.** We employed the PDP in Jin et al. (2020) to solve this problem. The learning rate is set to be  $\eta = 10^{-4}$ . Five trials were run given random initial  $\theta_0$ . For the neural objective function case, the learning rate is set to be  $\eta = 10^{-5}$ .

**Pytorch Adam** We use Pytorch Adam to learn a neural network  $f(\mathbf{x}, \mathbf{u}, \theta)$  to represent the system dynamics, where the input of the network is state and control vectors, and output is the state of next step. We train the neural network by minimizing the following residual

$$\min_{\theta} \sum_{t=0}^{T-1} \|\mathbf{x}_{t+1}^o - f(\mathbf{x}_t^o, \mathbf{u}_t, \theta)\|^2. \quad (58)$$

**DMDc.** The DMDc method Proctor et al. (2016) is a method that based on Koopman theory to represent a nonlinear dynamics with linear dynamics of observables. Observables are basis functions of states. The observable space has much higher dimension compared to state space. The estimation of the dynamics is achieved by the following optimization:

$$\min_{\mathbf{A}, \mathbf{B}} \sum_{t=0}^{T-1} \|\mathbf{x}_{t+1}^o - \mathbf{A}\mathbf{x}_t^o - \mathbf{B}\mathbf{u}_t\|^2. \quad (59)$$

**Learning a neural dynamics.** We consider the dynamics of each system (cartpole, quadrotor and rocket) are represented by a fully-connected feed-forward neural network  $\hat{f}(\mathbf{x}_t, \mathbf{u}_t, \theta)$ . The neural network has a layer structure of  $(n+m)$ - $2(n+m)$ - $n$  with *tanh* activation functions, i.e., there is an input layer with  $(n+m)$  neurons equal to the dimension of state, one hidden layer with  $2(n+m)$  neurons and one output layer with  $n$  neurons. The  $\xi^o = \{\mathbf{x}_{0:T}^o, \mathbf{u}_{0:T-1}\}$  obtained previously are used in stage loss. We conducted five trials for each method with different initial  $\theta$ .

#### D.4 Policy Tuning On-the-fly

**Neural State Feedback Policy.** In this application, we learn a the parameter of a neural state feedback policy by minimizing given control objective functions. Specifically, we use a fully-connected feed-forward neural network which has a layer structure of  $3n$ - $3n$ - $m$  with *tanh* activation functions, i.e., there is an input layer with  $3n$  neurons equal to the dimension of state, one hidden layer with  $3n$  neurons and one output layer with  $m$  neurons. The policy parameter  $\theta$  is the neural network parameter. For comparison, we apply the guided policy search (GPS) method Levine & Abbeel (2014) and iLQR Li & Todorov (2004) to solve the same problem.

**PDP.** We employed the PDP in Jin et al. (2020) to solve this problem. The learning rate is set to be  $\eta = 10^{-4}$  or  $= 10^{-6}$ . Five trials were ran given random initial  $\theta_0$ . For the neural objective function case, the learning rate is set to be  $\eta = 10^{-5}$ .